



YÜZ FEN EDEBİYAT FAKULTESİ İSTATİSTİK BÖLÜMÜ

RUSYA UKRAYNA SAVAŞININ DOLAR KURUNA ETKİSİNİN MODELLENMESİ

İlknur Ayaç ÇİÇEK -19023045

Danışman: Doç. Dr. Serpil KILIÇ DEPREN

ÖZET

Ekonomi ve politika sürekli olarak karşılıklı bir etkileşim içerisindedir. Özellikle ekonomik değişkenler demokratikleşme arasında ilişki hakkında kesin bir yargıya varılmaması, politik istikrarsızlık ile ekonomik değişkenler arasında ortaya çıkan ilişkilerin incelenmesinin gerekliliğini ortaya koymaktadır. Politik istikrarsızlık birçok sebebe bağlıdır. Sebeplerden bazıları terörizm, anayasa değişiklikleri, seçim sıklığı, komşu ülkeler arasında yaşanan olumlu ve olumsuz olaylar ve savaşlardır. Bu yüzden yaşanabilecek bu durumlarda dünyanın geçmişte yaşanmış olan olaylardan yola çıkarak ekonomik olarak ne kayıplara yaşayacağını belirlemesi ve tahmin etmesi beklenmektedir. Bu duruma örnek teşkil edebilecek en büyük olaylardan bir tanesi 11 Eylül Olayı ve sonucunda oluşan Amerika - Afganistan savaşlarıdır. Bu araştırma önerisinde özellikle ülkemizde dolar kurundaki oynaklığı arttıran ve güncel bir konu olan Rusya Ukrayna savaşının etkisinin analiz edilmesi ve modellenmesi amaçlanmıştır. Dolar kuruna etki eden BIST 100 endeksi, VIX endeksi, CDX endeksi gibi makroekonomik göstergeler bağımsız değişkenler olarak alınmıştır. Bu değişkenlere ek olarak, Rusya - Ukrayna savaşının etkisini incelemek amacıyla farklı yerli ve yabancı ajanslara ait gazete metinlerinden faydalanılarak sayısal bir bağımsız değişken oluşturularak modele dahil edilmiştir.

METİN MADENCİLİĞİ

Metin madenciliği (text mining), metin verilerini analiz ederek anlamlı bilgiler çıkarmayı amaçlayan bir veri madenciliği dalıdır. Metin madenciliği, doğal dil işleme, istatistiksel analiz, makine öğrenimi ve bilgisayar dilbilimi gibi disiplinlerin tekniklerini kullanarak metin verilerinin yapılandırılması, keşfedilmesi, modellenmesi ve yorumlanması sağlar. Metin madenciliği, büyük miktarda metin verisine sahip belgelerin otomatik olarak analiz edilmesini ve anlamlı bilgilerin çıkarılmasını mümkün kılar.

METİN MADENCİLİĞİ YÖNTEMLERİ

1- DUYGU ANALİZİ (SENTIMENT ANALYSIS)

Duygu analizi (sentiment analysis), metin veya belgelere uygulanan bir doğal dil işleme tekniğidir. Amacı, bir metnin veya belgenin duygusal tonunu anlamak ve sınıflandırmaktır. Duygu analizi, metinlerdeki duygusal ifadeleri tespit ederek pozitif, negatif veya nötr gibi duygusal etiketler atar. Bu sayede, sosyal medya mesajları, müşteri yorumları, ürün incelemeleri, haber başlıkları ve diğer metin verileri gibi büyük miktardaki metinleri otomatik olarak analiz etmek ve anlamak mümkün hale gelir.

SENTIMENT ANALYSIS



2- ÖZETLEME (SUMMARIZATION)

Text Summarization olarak da adlandırılan bu teknik, uzun metinlerin ana hatlarından anlamlı, kısa ve tutarlı bir özet oluşturulmasına yardımcı olur.

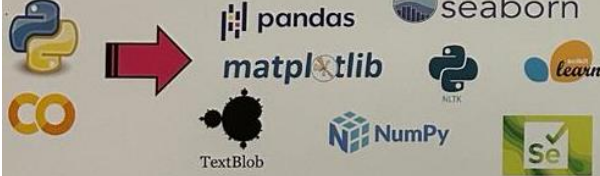
3- METİN SINIFLANDIRMA (TEXT CLASSIFICATION)

Metin sınıflandırma, metin verilerini farklı kategorilere veya sınıflara ayırmak için kullanılan bir makine öğrenimi tekniğidir. Temel olarak, belirli bir metin belgesini analiz ederek, içeriğine dayanarak belgenin hangi kategoriye veya sınıfa ait olduğunu belirlemeye çalışır. Metin sınıflandırma, metin verilerinin otomatik olarak etiketlenmesini ve organize edilmesini sağlar.



VADER ANALİZİ

Doğal dil işleme alanında sıklıkla kullanılan bir duygusal analiz aracıdır. VADER, metin verilerini pozitif, negatif veya nötr duygusal değerlere göre puanlar. Ek olarak, isteklerine göre duygusal bileşiklik (emotional compound) adında bir değer daha bulunur. Bu değer, metnin genel duygusal yoğunluğunu temsil eder. Duygusal bileşiklik, pozitif ve negatif duygusal ifadelerin yoğunluğunu dikkate alarak hesaplanır.



Proje python dili ve kütüphaneleri kullanılmıştır. Projede özellikle tool olarak Google colab kullanılmıştır. Yukarıda kullanılmış olan dil ve kütüphaneler görsel olarak verilmiştir. Selenium veri toplamak için pandas, numpy kütüphaneleri veri manipülasyonunda kullanılmıştır. Matplotlib ve seaborn görselleştirme için NLTK ve TextBlob kütüphanesi veri ön işlemede kullanılmıştır. En son olarak model kurmak için scikit learn kütüphanesi kullanılmıştır.

VERİNİN ELDE EDİLMESİ

Yerli ve yabancı gazeteler olmak üzere veriler sekiz ayrı kaynaktan toplanmıştır. Bu haber ajanslarından 24.02.2022-24.02.2023 tarihleri arasındaki Rusya Ukrayna savaşı adına yer alan haber metinleri ve bu haber metinlerinin tarih, link, açıklama ve başlık bilgileri selenium, newspaper3k kütüphaneleri sayesinde toplanıp .xlsx formatında kayıt edilmiştir.



VERİ ÖN İŞLEME

Veriyi kullanılabilir hale getirmek için aşağıdaki adımları sırasıyla gerçekleştirecek fonksiyon oluşturulmuştur.
Genel olarak bulunan bir yazım hatasını düzeltme
Kelimeler arası birden fazla boşluk olan yerlerin tek boşluk ile değiştirilmesi
Noktalama işaretlerinin boşluğa çevirme
Sayıları boşluğa çevirme

Öncüde ön işlemeden geçen haber metinleri vader analizi uygulanarak 87 farklı değişkene çevrilmiştir.

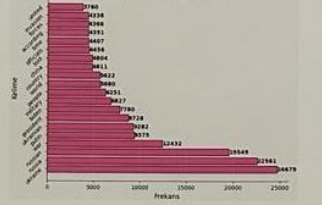
KAYNAKÇA

- Alpay, Ö. (2017). Ekonomi Haberlerinin BIST 100 Endeksinin Veri Madenciliği ile İncelenmesi Üzerine Bir Çalışma, Yüksek Lisans Tezi, Elazığ Üniversitesi, Sosyal Bilimler Enstitüsü, Elazığ.
- Kanago, B., & McCormick, K. (2013). The Dollar-Pound Exchange Rate During the First Nine Months of World War II. Atlantic Economic Journal, 385-395.

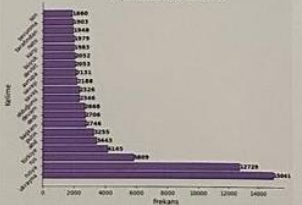
VERİLERİN GÖRSELLEŞTİRİLMESİ

Veriler seaborn ve matplotlib kütüphaneleri yardımıyla yabancı ve Türkçe kaynaklar arasındaki farkı iyi anlayabilmek, en çok tekrar eden ve metinleri yapısal olarak en çok hangi kelimelerden oluştuğunu tespit edebilmek için metinler görselleştirilmiştir. İlk olarak Türkçe ve yabancı haber metinleri olarak ayrılandırılmış son word cloud kullanılarak değişim büyüklüklerine göre görselleştirilmiştir. Farklı olarak Türkçe kaynaklarda «abd» kelimesi yer alırken yabancı kaynaklarda «china» kelimesinin sıklığı göze çarpmaktadır.

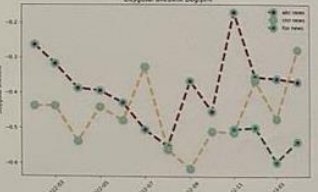
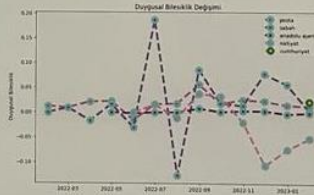
En Çok Tekrar Eden 20 Kelime



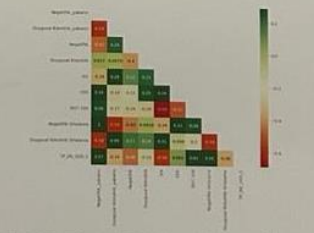
En Çok Tekrar Eden 20 Kelime



İkinci adımda ise en çok tekrar eden ilk 20 kelime frekans bazında sütun grafiği kullanılarak görselleştirilmiştir. Hem yabancı hem de Türkçe kaynaklarda kullanılan kelimelerin büyük oranda benzediğini görmekteyiz.



Üçüncü adım olarak ay bazında veri tabanımızda tuttuğumuz haber metinlerinin aylık duygusal bileşik değerlerinin çizgi grafiğini görmekteyiz.



Dördüncü adım olarak dolar kuruyla hem aylık hem de günlük olarak bağımsız değişkenler arasındaki korelasyonları heatmap grafikleriyle görselleştirilmiştir.

MODELLEME

Regresyon tabanlı makine öğrenimi algoritmaları denenmiştir. Yanda tablo şeklinde en iyi sonuç veren algoritmadan en kötü sonuç veren algoritmaya göre sıralanmıştır. Karar verme değeri olarak RMSE kullanılmıştır.

MODEL	RMSE
RandomForestRegressor	0.225
GradientBoostingRegressor	0.235
XGBRegressor	0.251
LGBMRegressor	0.258
DecisionTreeRegressor	0.263
KNeighborsRegressor	0.275
SVR	0.372
MLPRegressor	1.036

SONUÇ

Bu çalışmada çeşitli yerli ve yabancı haber ajanslarının Rusya-Ukrayna Savaşı adı altında web sitelerinde yer verdikleri haber metinleri ele alınmıştır ve bu haber metinlerine literatür dikkate alınarak çeşitli ön işleme adımları uygulanmıştır. Duygu analizi yöntemlerinden biri olan VADER yöntemi kullanılmış olup CIX endeksi, VIX endeksi ve BIST 100 endeksi ile birlikte on bir adet bağımsız değişkenle dolar kurunu tahmin edilmiştir. Haber ajanslarından alınan bu metinler kelime bulutu ve bar grafikleriyle görselleştirilmeleri yapılmış olup yabancı kaynaklı haber ajanslarının yerli haber ajanslarına göre tutumlarının daha olumsuz olduğu saptandı. En çok kullanılan kelimelere bakıldığında yerli ve yabancı haber ajansları arasında çok büyük bir fark saptanmamıştır. Haber ajanslarının metinlerinden yola çıkılarak oluşturulan sayısal değişkenlerin korelasyon analizi sonuçları göz önünde bulundurulduğunda aylık bazda bakıldığında yabancı haber ajanslarının dolar kuruna anlamlı ve göze alınabilecek sonuçlar ortaya çıkaran yerli haber ajanslarıyla oluşturulan sayısal değişkenlerin modelde etkisinin olduğu oranda ve şekilde gerçekleştirilemeyeceği ön görülmüştür. Diğer bağımsız değişkenler göz önüne alındığında modelde etkisinin en kuvvetli olduğu görülen değişken BIST 100 endeksi olmuştur. Regresyon tabanlı makine öğrenmesi algoritmalarıyla oluşturulan modelde RMSE değerleri dikkate alındığında veriye en uygun modelin Random Forest



İSTATİSTİK BÖLÜMÜ

MAKİNE ÖĞRENMESİ İLE OECD ÜLKELERİNİN ÇEVRESEL PERFORMANSININ TAHMİN EDİLMESİ

Beste ÖZÜLKÜ

Doç.Dr. Serpil KILIÇ DEPREN

ÖZET

Bu çalışma, çevresel performans üzerinde farklı bağımsız faktörlerin etkisini derinlemesine analiz etmeyi amaçlamaktadır. Çalışma, çevresel performans endeksinde odaklanarak, 181 farklı ülkeyi kapsayan 509 gözlemler içeren bir nihaî veri kümesi kullanılarak veri kümesi, yedi farklı kaynaktan derlenen verilerin birleştirilmesiyle oluşturulmuş ve altı farklı makine öğrenimi algoritmasıyla modellenmiştir. En önemlisi, Super Learner algoritması da birden fazla makine öğrenimi algoritmasını birleştiren ensemble modeli oluşturmak için kullanılmış ve tüm bireysel makine öğrenimi algoritmalarının önüne geçmiştir.

ÇEVRESEL PERFORMANS

Çevresel performans, sürdürülebilir kalkınmanın kritik bir boyutudur ve bir ülkenin bu hedeflere ne kadar yakın olduğunu ölçer. Çevresel performans göstergeleri, çevresel politika sonuçlarına odaklanarak hükümetlerin kapsamlı kirillik kontrolü ve doğal kaynak yönetimi hedeflerine ilerlemesini ölçmelerine yardımcı olur. Hedefler, mevcut uluslararası anlaşmalara, kirillik insanlar ve ekosistemler üzerindeki zararlı etkilerine dair bilimsel kanıtlara ve ekonomik olarak sürdürülebilir çevre koruma stratejilerine dayanmaktadır. Hedefler, temiz hava, içme suyu ve biyolojik çeşitliliğin korunması gibi temel çevresel konuları kapsar. Her hedef, teorik mantık, politika uyumu, ölçülebilirlik ve mevcut veri kapsamı gibi katı kriterlere göre seçilen birkaç performans göstergesiyle ilişkilidir.

Çevresel Performans Endeksi (Environmental Performance Index: EPI), öncelikle kirillik, doğal kaynak yönetimi, çevresel sağlık, ekosistem kalitesi ve iklim değişikliği gibi çevresel konuları değerlendirmek için kullanılan bir dizi ölçüme dayanmaktadır. Bu göstergeler arasında su ve hava kalitesi, atık yönetimi, biyolojik çeşitlilik ve orman koruması gibi unsurlar yer alır.

Endeks, çevresel performansın yanı sıra sosyal ve ekonomik konuları da dikkate alır. Bu, çevresel sürdürülebilirliği diğer unsurlarla birlikte kapsamlı bir şekilde incelemeyi sağlar. Örneğin, çevresel performans sosyal eşitsizlik, yoksulluk ve ekonomik kalkınma gibi unsurlardan etkilenebilir ve değerlendirilmeye dahil edilmeleri önemlidir. Bu çalışma, makine öğrenimi (ML) algoritmalarını kullanarak 2018, 2020 ve 2022 yıllarında 181 farklı ülkede çevresel performansı etkileyen faktörlerin etkilerini incelemeyi amaçlamaktadır. Özellikle, çevresel yönetim faaliyetlerini teşvik etmek için en iyi algoritmayı ve politikayı belirleme sürecinde EPI ile ilişkiler arasındaki bağlantıları incelemeye katkıda bulunmaktadır. Bu şekilde, çevresel performansı iyileştirme çabaları, ülkelerin çevresel düzenlemelere uyum sağlamasını teşvik edecektir.

Bu çalışma aşağıdaki araştırma sorularını ele almaktadır:

- Çevresel performans üzerinde hangi faktörler önemli bir etkiye sahiptir? Etki yönü nedir?
- Önemli değişkenlerdeki kritik eşik değerleri nelerdir?
- Seçilen değişkenlerle çevresel performans modelenebilir mi? Hangi ML algoritmaları diğerlerinden daha güçlü tahminler üretir?
- Çevresel performans göstergesi en iyi hangi algoritma ile modelenebilir? Eğitim ve Test Verilerindeki model performansını tatmin edici midir?
- Seçilen modelden öne çıkan öncelikli hareket alanları nelerdir? Hangi politika önerileri sunulmalıdır?

METODOLOJİ

İncelenen değişkenlerin çevresel performans üzerindeki etkisini incelemek için Makine Öğrenimi (ML) algoritmalarını kullanılmaktadır. Son zamanlarda birçok farklı disiplinde sınıflandırma ve tahmin yapmak tercih edilmektedir. ML, verilere dayalı olarak kendini öğrenebilen matematiksel modellerin geliştirilmesine izin veren bir yapay zeka uygulamasıdır. Bu, karmaşık görünen bir veri yığınındaki problemleri tespit etmenizi, gelecekteki senaryoları öngörmenizi ve verideki desenleri keşfetmenizi sağlar.

MAKİNE ÖĞRENİMİ (ML) ALGORİTMALARINI

Sınıflandırma ve Regresyon Ağaçları (CART)

Rastgele Orman (RF)

Aşırı Gradyan Artırma Makineleri (XGB)

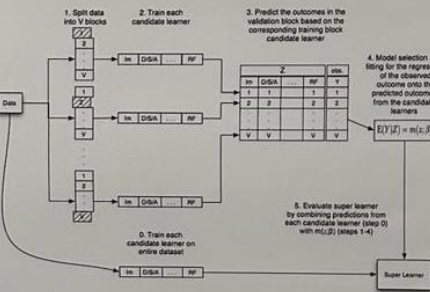
Gradyan Artırma Makineleri (GBM)

Destek Vektör Makineleri (SVM)

k-En Yakın Komşular (k-NN)

SÜPER ÖĞRENİCİ ALGORİTMASI

Süper öğrenme, çeşitli öğrenme algoritmalarının optimal bir kombinasyonunu bulan bir ensemble yöntemiyle bu sorunu ele alır. Bu, tahminsel modelleme problemleri için araştırılabileceğiniz tüm modelleri ve yapılandırmaları birleştirerek, araştırmanız herhangi bir tek modelden daha iyi veya ona eşit bir tahmin yapmak için kullanılan bir tekniktir.

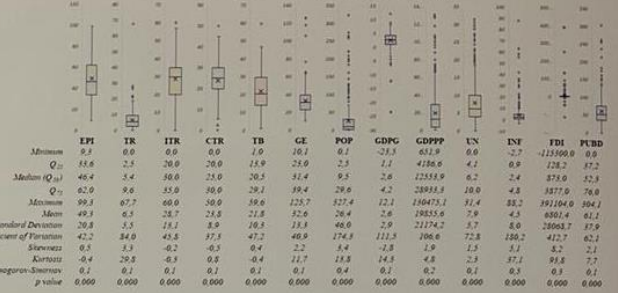


Yukarıda görüldüğü gibi, süper öğrenici, bireysel modellerin k-katlımlı veri bölünmesinde eğitildiği bir varyasyon olan istifleme veya k-katlımlı çapraz doğrulamayı içerir. Ardından, her modelin çıktısı olan kat dışı tahminlerden de anılan bir nihaî meta modeli eğitilir.

KAYNAKÇA

- Van Der Laan, M. J., Polley, E. C., & Hubbard, A. (2007). Super Learner. *Statistical Applications in Genetics and Molecular Biology*, 6(1). <https://doi.org/10.2202/1544-6115.1309>
- Young, S. L., Abdou, T., & Bener, A. (2018a). Deep Super Learner: A Deep Ensemble for Classification Problems. In *Lecture Notes in Computer Science* (pp. 84-95). Springer Science+Business Media. https://doi.org/10.1007/978-3-319-89656-4_7
- Rumao, P. (2022, March 30). Super Learner versus Deep Neural Network - Towards Data Science. Medium. <https://towardsdatascience.com/super-learner-versus-deep-neural-network-706472c377>

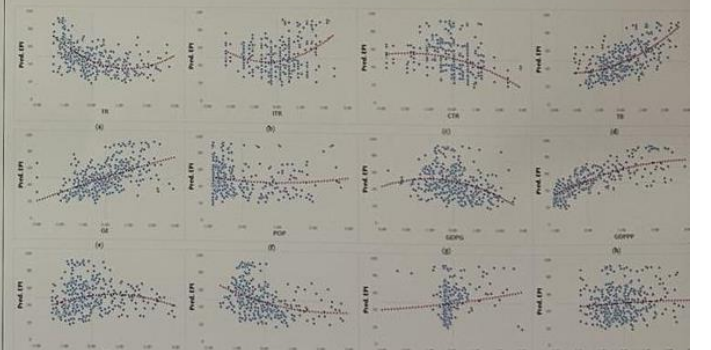
İSTATİSTİKLER VE BOX-JENKINS GRAFİKLERİ



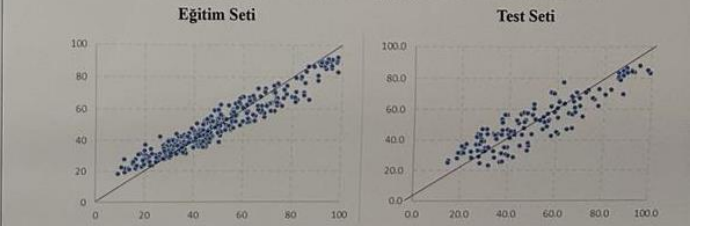
MAKİNE ÖĞRENME ALGORİTMALARININ BULGULARI SÜPER ÖĞRENİCİ ALGORİTMASININ PERFORMANS METRİ

ML Algoritmaları	Eğitim Seti (70% - 373 gözlem)			Test Seti (30% - 136 gözlem)			Eğitim Seti Performansı (70% - 373 gözlem)			Test Seti Performansı (30% - 136 gözlem)			
	R ²	RMSE	MAE	R ²	RMSE	MAE	R ²	RMSE	MAE	R ²	RMSE	MAE	
CART	0.460	13.1076	10.4831	0.480	20.9591	8.0257	0.2174	0.400					
GBM	0.458	13.1714	10.8988	0.480	20.9591	8.0478	0.2174	0.400					
XGB	0.458	13.3359	10.6013	0.479	20.8146	8.0240	0.2174	0.400					
k-NN	0.447	14.15098	12.3839	0.472	20.9999	8.1136	0.2174	0.400	0.402	6.3948	5.1531	0.479	9.6707
CART	0.443	15.4409	12.4509	0.452	14.5279	12.1075	0.2212	0.400					
SVM	0.418	18.9072	14.5938	0.467	17.8502	14.3446	0.2576	0.400					

DEĞİŞKENLER ARASINDAKİ İKİLİ İLİŞKİ VE ÖNEMLİ EŞİK DEĞERLERİ



SÜPER ÖĞRENİCİ ALGORİTMASININ GERÇEK VE TAHMİN EDİLEN DEĞERLERİ



Sonuçlar

Analiz sonuçlarına göre, kişi başına düşen gayri safi yurtiçi hasıla (GSYİH), tarife oranları, vergi yükü, devlet harcamaları ve enflasyon çevresel performansı geliştirmek için politika yapıcıların odaklanması gereken beş temel faktör olarak belirlenmiştir. Ayrıca, nüfus, gayri safi yurtiçi hasıla büyüme oranı, gelir vergisi oranı, kamu borcu, doğrudan yabancı yatırım ve kurumsal vergisi oranı çevresel performans üzerinde önemli bir etkiye olduğu belirlenmiştir. Bunun yanı sıra, ilgili faktörlerin çevresel performans üzerindeki etkisini gösteren kritik eşikler belirlenmiştir. Sonuç olarak, bu makale, politika yapıcılara çevresel performansı geliştirebilecek faktörler konusunda değerli öngörüler sunmaktadır. Yenilenebilir enerji kaynaklarına yatırım yaparak, yenilenebilir ve yeşil enerji kaynaklarına yönelik yabancı yatırımları teşvik ederek, ekonomik büyüme ve GSYİH'yi kontrol altına alan politika yapıcılar, çevresel performansı önemli ölçüde ilerletilebilirler. Ayrıca, politika yapıcılar, vergi yükü, gelir vergisi ve kü vergisinin çevresel performans üzerindeki etkisini de dikkate almalı ve gerekli reformları yapmalıdır.



CUSTOMER CHURN ANALYSIS

EZGİ GÖKBUA 21023604

Thesis Advisor: Prof. Dr. Filiz KARAMAN

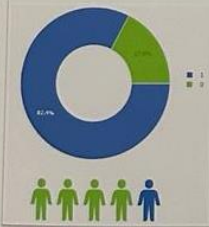
Customer churn analysis is a crucial topic for companies to succeed in today's world. The aim of this thesis is to predict customers with a high likelihood of churn by reviewing their purchasing behaviors. Anonymous data from customers who renewed and did not renew their comprehensive insurance policies in a private insurance company in Turkey was used. The study employed data balancing methods such as oversampling and undersampling. Machine learning classification methods were applied on the balanced dataset, with Random Forest, XGBoost, and LightGBM algorithms yielding the highest success rates. The data set was processed using the Python programming language.

Customer Churn Analysis in the Insurance Sector

Churn analysis involves various methods and techniques to identify the reasons behind customer attrition and make strategic decisions. Typically, data such as customer behavior, demographic information, past purchasing habits, and customer satisfaction feedback are utilized in this analysis. The insights derived can help companies offer tailored solutions to customers, increase customer loyalty, and facilitate the growth of their customer base. In this study, customer data from a specific insurance company regarding the renewal and cancellation of comprehensive insurance policies has been used. In addition to demographic information such as age and gender, the data also includes variables related to the insured vehicle, such as its age, brand, and type.



Churn Customer Percentage



Customer churn analysis is closely related to the growth strategies and revenue goals of companies. Determining why customers are leaving and identifying the factors that contribute to customer churn provides valuable insights for companies to improve customer relationships and take measures to minimize customer attrition.

During the observation, it was found that 61,014 customers have renewed their policies, while 13,001 customers have chosen not to renew their policies and ceased to be customers of the insurance company.

Before applying classification methods for customer churn prediction on the dataset, the issue of imbalanced dataset, which is an important point in classification problems, has been addressed. Two separate sampling techniques, namely Under Sampling and Over Sampling have been thoroughly examined.

Techniques to Handle Imbalanced Data

Undersampling is a sampling technique used to address the problem of imbalanced datasets in classification problems. In this method, the number of instances belonging to the majority class is reduced to achieve a more balanced dataset with the minority class.

- > Random Under Sampler
- > Near Miss
- > Cluster Centroids
- > Tomek Links
- > ENN

Oversampling aims to artificially increase the number of samples in classes with a small number of instances. This way, it enhances the importance and impact of classes with a small number of instances, enabling the model to learn them better. This can potentially improve the classification performance of the model and reduce the occurrence of erroneous results due to class imbalance.

- > Random Over Sampler
- > SMOTE
- > ADASYN



Machine Learning Classification Methods

XGBoost

SMOTE



ADASYN



Cluster Centroids



XGBoost is a tree-based learning method, which means it utilizes an ensemble of decision trees. Each tree divides the dataset using a set of decision rules to predict a target variable. The decision rule is created by comparing the features (or attributes) in the dataset with a specific threshold value. Each split further divides the dataset into more homogeneous subsets.

LightGBM Classifier

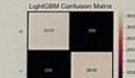
SMOTE



ADASYN



Cluster Centroids



LGBM (LightGBM) is a machine learning algorithm that is lightweight and fast, and it is a gradient boosting method. It divides the dataset into rectangular regions and updates the sample weights using histograms for each region. This allows for high performance on large datasets and helps reduce overfitting.

Random Forest

SMOTE



ADASYN



Cluster Centroids



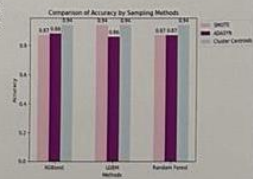
Random Forest is a machine learning algorithm that consists of an ensemble of decision trees. Each decision tree processes the dataset by randomly sampling and splitting the data using random features. As a result, this method reduces high variance, improves generalization on the data, and makes stable predictions.

Conclusion

"The Lazy Classifier" method was used to evaluate all classification techniques. This evaluation was performed on the datasets obtained using balancing techniques. The evaluation identified three machine learning classification techniques, namely XGBoost, LGBM, and Random Forest, that demonstrated high performance across all balancing techniques. These methods were found to have high accuracy and F-values, indicating their effectiveness in predicting customer churn. When examining the obtained results, methods like XGBoost, LGBM, and Random Forest demonstrate themselves as powerful tools that can be used to predict churn customers in the insurance company's dataset.

Method	SMOTE	ADASYN	Cluster Centroids
XGBoost	0.87	0.88	0.94
LGBM	0.94	0.86	0.94
Random Forest	0.87	0.87	0.94

Accuracy values of Classification Methods



Recommendations

There are several strategies that can be implemented for customers who have the potential to churn:

- Customer Relationship Management (CRM)
- Personalized Marketing
- Improved Customer Experience
- Loyalty Programs
- Win-back Campaigns
- Data Analytics and Predictive Models

REFERENCES

[1] Akyüz, H. E., & TAŞCI, T. (2021). SİGORTACILIK SEKTÖRÜNDE MAKİNE ÖĞRENMESİ İLE MÜŞTERİ KAYBI ANALİZİ. Tasarım Mimarlık ve Mühendislik Dergisi, 2(1), 66-79.
 [2] AYDIN, M. A. (2021). Müşteri Kaybı Tahmininde Sınıf Dengesizliği Problemi. Politeknik Dergisi, 1-1.
 [3] KAYNAR, O., TUNA, M. F., GÖRMEZ, Y., & DEVECİ, M. A. (2017). Makine öğrenmesi yöntemleriyle müşteri kaybı analizi. Cumhuriyet Üniversitesi İktisadi ve İdari Bilimler Dergisi, 18(1), 1-14.



DEPARTMENT OF STATISTICS

Convolutional Neural Networks for Generating, Modeling and Testing

Statistical Data from Traffic Images

Fikret Efe DIRİ 18023058

Advisor: Prof. Dr. Filiz KARAMAN

Abstract

As the human population continues to grow, various costs arise, one of which is traffic congestion. The blockage or complete halt of traffic flow leads to negative consequences such as time lost, increased carbon emissions, and higher fuel consumption. Efforts to minimize these adversities are made by entities such as highway administrations, municipalities, and transportation engineers, who rely on accurate traffic flow data as one of the most important sources of information. While traffic flow data can be obtained through multiple methods, cameras emerge as one of the best solutions due to their affordability and flexibility. This study aims to generate and analyze traffic flow data derived from camera recordings for the purpose of mitigating these negative impacts.

What is Traffic Flow Data?

Traffic flow data is a type of data that involves counting the number of vehicles in traffic based on certain breakdowns, as well as including vehicle speeds. These breakdowns can include vehicle types such as cars, trucks, motorcycles, etc., and information about the routes of the vehicles on the road. Below is an example of a result obtained from a study where this data was generated along a road with two-way traffic.

Cls	Direction	Kmph
2	Left	140,10
2	Right	134,44
2	Left	113,54
7	Right	88,93
2	Right	134,76

- **Cls** : It includes the types of vehicles, specifically 2 cars and 7 trucks.
- **Direction** : It includes the route information of the vehicles. 'Left' indicates the direction of travel, while 'Right' indicates the opposite direction based on the location where the recording is taken.
- **Kmph**: It includes the speeds of the vehicles.

The Algorithm and Statistical Tests

This study utilized the YoloV5_StrongSort_OsNet Algorithm, which is based on an artificial neural network. The algorithm combines the one-stage object detection algorithm YoloV5 with the construction of an object tracking system. Similar to object tracking, the object detection and classification components of the algorithm have been pre-trained to create an enhanced algorithm. YoloV5 is a convolutional neural network (CNN) based algorithm capable of object detection and processing live images at a high speed.

Although the algorithm demonstrates success in object recognition, prediction, and detection, several modifications need to be made to adapt it for traffic flow data. The study makes use of statistical tests such as Shapiro-Wilk, Kruskal-Wallis, and Mann-Whitney to analyze the data.

Customization of the Algorithm

1. Entrance and Exit Detection:

Firstly, it is necessary to digitally detect the entrance and exit pixels of vehicles on the road from the recorded traffic image. This process is carried out by analyzing a single frame of the image, as shown in the photograph.



During the execution of the algorithm, apart from the entrance and exit detection, the video recording is played live on the screen. In order to display the calculations performed on this window, counters will be added for the vehicles entering the frame. The positions of these counters and texts have been calculated.

2. Object Count:

For this task, if-else blocks have been created as shown below. These blocks trigger the code responsible for storing various raw data when the center pixel coordinates of the detected and tracked vehicle, obtained from the algorithm, fall within the predetermined coordinates. This allows for the collection of raw data related to the vehicle.

```
def count_obj(box_w, h, id, cls):
    global ccl, ctl, ccr, ctr

    if id not in control:
        if int(box[1] + int(box[3] - box[1]) / 2) > (h - 370):
            if int(box[0] + (box[2] - box[0]) / 2) > 50 and int(box[0] + (box[2] - box[0]) / 2) < (615):

                control.append(id)
                dctSt['TimeSt'].append(time.time())
                dctSt['ID'].append(id)
                dctSt['Cls'].append(cls)
                dctSt['Direction'].append('Left')

                if cls == 2:
                    ccl += 1
                elif cls == 7:
                    ctl += 1
```

3. Processing of the Data

The collected raw data includes the route and class of the vehicle as desired. However, additional calculations are required to determine the speed of the vehicle. After the vehicles enter the designated entrance and exit areas, the computer retrieves the timestamp information, and the speed is calculated in pixel-seconds based on the time difference. Since the length of the specified region is not known, the speed data is estimated hypothetically.

Results

Vehicle	Direction	Real	Algorithm Count	Miss Rate
Car	Left	217	209	3,69%
	Right	202	188	6,93%
Truck	Left	36	31	13,89%
	Right	37	33	10,81%



The output during the execution of the algorithm appears as shown in the following figure. The vehicles on the road consist only of cars and trucks. Vehicles classified as vans and trucks are considered under the category of 'truck' at a single level. The difference between the algorithmic observations and manual road counting is evident in the images above.

```
# Perform Mann-Whitney U test
p_value = stats.mannwhitneyu(right_data, left_data, alternative='two-sided')
alpha = 0.05 # Significance level

if p_value < alpha:
    if right_data.median() > left_data.median():
        print("The 'Right' direction is significantly faster than the 'Left' direction.")
    else:
        print("The 'Left' direction is significantly faster than the 'Right' direction.")
else:
    print("There is no significant difference in the median 'kmph' values between the 'Right' and 'Left' directions.")
The 'Left' direction is significantly faster than the 'Right' direction.
```

Based on the conducted analysis, it can be concluded that there exists a statistically significant variation in vehicle speeds across different classes. Moreover, a noticeable discrepancy in speeds is observed between the right and left sides of the road. This study successfully facilitated the real-time generation of object detection data and speed measurements. Furthermore, an exemplified application of this data is presented, demonstrating its potential utilization.

```
# Create an empty dictionary to store the arrays for each group
group_arrays = {}

# Iterate over the groups and create arrays for each group
for cls, group in grouped_df:
    group_arrays[cls] = group['kmph'].values

# Perform Kruskal-Wallis U test
p_value = stats.kruskal(*group_arrays.values())
alpha = 0.05 # Significance level

if p_value < alpha:
    print("There is a significant difference in the means of the 'kmph' variable among the 'cls' groups.")
    print()

# Calculate the median for each group
medians = {cls: np.median(values) for cls, values in group_arrays.items()}

# Find the fastest cls based on the highest median
fastest_cls = max(medians, key=medians.get)

print(f"The 'cls' with the highest median 'kmph' value is {fastest_cls}." )
print(f"The median 'kmph' value for {fastest_cls} is {medians[fastest_cls]}." )
else:
    print("There is no significant difference in the means of the 'kmph' variable among the 'cls' groups.")
There is a significant difference in the means of the 'kmph' variable among the 'cls' groups.
The 'cls' with the highest median 'kmph' value is 2.
The median 'kmph' value for 2 is 116.10800000000004.
```

REFERENCES

[1] YoloV5_StrongSort_OsNet https://github.com/hndh2006/YoloV5_DeepSort_Pytorch
 [2] Yardim, M. S., & Akylidiz, G., (2005). Akıllı Ulaştırma Sistemleri ve Türkiye'deki Uygulamaları . 6. Ulaştırma Kongresi (pp.405-414). Istanbul, Turkey
 [3] <https://towardsdatascience.com/artificial-neural-networks-for-total-beginners-d8cd07abaae4>



FAKULTESİ

İSTATİSTİK BÖLÜMÜ

Uykusuzluğa Neden Olabilecek Faktörler Ve Hayatımıza Olası Etkileri

Tuğçe Gül YILMAZ 18023030

Danışman: Doç. Dr. Atif Ahmet EVREN

ÖZET

Araştırmada hayatımızda önemli bir faktör olan uyku ele alınmıştır. Bu ana başlık altında asıl incelenen ise uykusuzluğa neden olabilecek faktörler ve hayatımıza olası etkileri konusudur. Rastgele seçilen 256 kişilik örneklem üzerinde uygulanmıştır. Tezin amacı uykusuzluğa neden olabilecek faktörler ve uykusuzluğun hayatımıza olası etkilerini, ankette sorulan 27 altı değişken ile istatistiksel yöntemlerle araştırmak ve çıkan sonuçlar doğrultusunda ele alınan konu hakkında çıkarımlarda bulunmaktır. SPSS programı kullanılarak anket sonucundan elde edilen veriler incelenmiş ve analizler yapılmıştır. Genel olarak uykusuzluğu etkileyen faktörlerin birden fazla olduğu ve bunların özellikle; cinsiyet, medeni durum ve çalışma durumu değişkenlerine göre farklılıkları görülmüştür.

TEZİN AMACI

Amaç uykusuzluğa neden olabilecek faktörler ve uykusuzluğun hayatımıza olası etkilerini, ankette sorulan 27 altı değişken ile birlikte istatistiksel yöntemlerle araştırmak ve çıkan sonuçlar doğrultusunda ele alınan konu hakkında çıkarımlarda bulunmaktır.

TEZİN YÖNTEMİ

SPSS programı kullanılarak anket sonucuna çıkan veriler incelenmiş ve analizler yapılmıştır. İlk olarak tanımlayıcı istatistikler ele alınmıştır. Burada değişkenlerin frekansları ve dağılımları incelenmiştir. İkinci olarak normallik testleri yapılarak değişkenlerin normal dağılıp dağılmadığı incelenmiştir. Üçüncü olarak hipotezler kurulup değişkenlerin bir önceki bölümde elde ettiğimiz normallik sonuçlarına göre uygun olan testler yapılarak hipotezler yorumlanmıştır. İlk olarak, iki değişken arasında ortalamamız anlamlı olup olmadığını test eden bağımsız örneklem t-testi ve mann-whitney u testi uygulanmıştır. İkinci olarak, iki değişken arasında anlamlı bir ilişki olup olmadığını test eden ki-kare testi ve anlamlı bir ilişki olup olmadığına beraber bu ilişkinin yönünü de test eden korelasyon testi uygulanmıştır. Sonuçları aksi belirtilmediği sürece 0,95 güven aralığında $p \leq 0,05$ değeri anlamlı olarak kabul edilerek değerlendirilmiştir.

ÖN BİLGİLER

Sonuçları aksi belirtilmediği sürece 0,95 güven aralığında $p \leq 0,05$ değeri anlamlı olarak kabul edilerek değerlendirilmiştir.

Ki-kare testinde tablonun altında verilen bilgilerde 5'in altında frekans değeri %20'nin üzerinde bulunmadığına Pearson Ki-Kare değerini, %20'nin üzerinde bulunuyorsa yine tablonun altında verilen "beklenen sayı..." kısmını ele alarak 5'in altında olması durumunda Fisher's Exact test satırı, 5 ile 25 arasında ise Continuity Correction satırı ve eğer 25 üzerinde ise yine Pearson Ki-Kare satırı ele alınarak sonuçta varılmıştır. Monte Carlo simülasyonu ile elde edilen Exact testi p-değerinin (10.000 örneklem ve %99 güven aralığıyla) Exact seçeneğiyle elde edilen Exact testi p-değerine virgülden sonraki üç sıfır kadar aynı olduğu kabul edilmiştir [4].

Normallik testinde verilerin normal dağılıp dağılmadıklarına yapılan test sonucunda ortaya çıkan basıklık ve çarpıklık değerlerine bakılarak karar verilmiştir. Değerler "+1,0,-1,0" arasında ise normal dağılıklarına karar verilmiştir [5]

DEMOGRAFİK BULGULAR

Ankete Katılanların Cinsiyetleri

Ankete katılan 256 kişinin cinsiyet dağılımı Tablo 2.1.1 de verilmiştir ve Şekil 2.1.1 de gösterilmiştir. Buna göre anket katılımcılarının %68'inde kadınlar yer alırken, %33'ünde erkekler yer almaktadır.

Cinsiyetiniz nedir?			
	Frequency	Percent	Cumulative Percent
Valid			
Erkek	82	32,0	32,0
Kadın	174	68,0	100,0
Total	256	100,0	100,0

Ankete Katılanların Yaşları

Ankete katılan 256 kişinin yaş dağılımı Tablo 2.1.2 de verilmiştir ve Şekil 2.1.2 de gösterilmiştir. Buna göre anket katılımcılarının %47,3'ü 18-25 yaş arası olarak en fazla katılım sağlayan yaş aralığı olurken, %2,3'le 10-17 yaşları arasındaki katılımcılar ankete en az katılım sağlamıştır.

Yaş aralığınız nedir?			
	Frequency	Percent	Cumulative Percent
Valid			
10-17	6	2,3	2,3
18-25	121	47,3	49,6
26-35	72	28,1	77,7
36-50	39	15,2	93,0
50 ve üstü	18	7,0	100,0
Total	256	100,0	100,0

BETİMSSEL BULGULAR

Son Bir Ay İçinde Uykusuzluğun Rahatsız Etme Sıklığı

Ankete katılan katılımcıların %30,5'le çoğunluk olan kısmı bir ay içinde uykusuzluğun rahatsız etme sıklığının orta düzeyde olduğunu belirtmiştir.

1'den 5'e kadar bir ölçekte, 5 en yüksek olmak üzere, son bir ay içinde uykusuzluk sizi ne sıklıkta rahatsız etti?			
	Frequency	Percent	Cumulative Percent
Valid			
5	43	16,8	16,8
4	49	19,1	35,9
3	78	30,5	66,4
2	48	18,8	85,2
1	38	14,8	100,0
Total	256	100,0	100,0

Bir Hafta İçinde Kötü Uyuyan Gece Sayısı

Ankete katılan katılımcıların çoğunluk olan %22,3'ü en çok bir hafta içinde 3 ile 2 gün kötü uyuduklarını belirtmişlerdir. Bununla birlikte bir haftanın her günü kötü uyuyan katılımcı sayısı ile hiçbir günü kötü uyumayan katılımcı sayısının %6,6'lık oranlarla eşit sayıda oldukları görülmüştür.

Bir hafta içinde kaç gece kötü uyuduklarınıza değeri nedir?			
	Frequency	Percent	Cumulative Percent
Valid			
0	17	6,6	6,6
1	24	9,4	16,0
2	40	15,6	31,6
3	57	22,3	53,9
4	32	12,5	66,4
5	26	10,1	76,5
6	17	6,6	83,1
Total	256	100,0	100,0

HİPOTEZLER VE TESTLER

MANN-WHITNEY U TESTİ

Cinsiyet * Gecenin Bir Yarısı Uyuma Nedeni

Hipotez;

H0: Cinsiyet değişkeni ile gecenin bir yarısı uyuma nedeni değişkeninin ortalamaları arasında anlamlı bir fark yoktur.

H1: Cinsiyet değişkeni ile gecenin bir yarısı uyuma nedeni değişkeninin ortalamaları arasında anlamlı bir fark vardır.

Tablo 3.2'de gösterilen test sonucundaki anlamlılık değerine bakıldığında anlamlılık değerinin 0,05'ten küçük bir değer olduğu görülmüştür (0,042<0,05). Sonuç olarak H0 hipotezi reddedilerek, cinsiyet değişkeni ile uyku problemi için yardım alma değişkeninin ortalamaları arasında anlamlı bir farkın olduğu görülmüştür.

Kİ-KARE TESTİ

Son bir ay içinde uykusuzluk sizi ne sıklıkta rahatsız etti? * Medeni durumunuz?

Hipotez;

H0: Son bir ay içinde uykusuzluğun rahatsız etme sıklığı medeni durumdan bağımsızdır.

H1: Son bir ay içinde uykusuzluğun rahatsız etme sıklığı medeni duruma bağlıdır.

Tablo 3.3.3'te görüldüğü gibi anlamlılık değeri 0,05 anlamlılık değerinden küçük olduğu için H0 reddedilmiştir (0,005<0,05). Medeni durum ile son bir ay içinde uykusuzluğun rahatsız etme sıklığı arasında bir ilişki bulunmaktadır.

KORELASYON TESTİ

H0: Bir hafta içinde kötü geçen gece sayısı ile 24 saat içinde uyunan saat arasında ilişki yoktur.

H1: Bir hafta içinde kötü geçen gece sayısı ile 24 saat içinde uyunan saat arasında ilişki vardır.

Tablo 3.4.1'e bakıldığında 0,01 anlamlılık seviyesine göre ilişkinin anlamlı olduğu görülmüştür (0,001<0,01) ve böylelikle H0 hipotezi reddedilmiştir. Bir hafta içinde kötü geçen gece sayısı ile 24 saat içinde uyunan saat arasında istatistiksel anlamda ilişki vardır. Pearson korelasyon katsayısına bakıldığında negatif yönlü bir ilişkinin olduğu görülmüştür (-0,215).

Test Statistic ^a		Chi-Square Tests		Likelihood Ratio		Fisher's Exact Test		Linear-by-Linear Association		N of Valid Cases	
	df	Value	Asymp. Sig.	Value	df	Asymp. Sig.	Value	df	Asymp. Sig.		
Chi-Square	1	10,200	,002	10,200	1	,002				256	
Continuity Correction ^b	1	9,880	,002							256	
Fisher's Exact Test				10,200	1	,002				256	
Linear-by-Linear Association	1	10,200	,002							256	

SONUÇ

Araştırmada hayatımızda önemli bir faktör olan uyku ele alınmıştır. Bu ana başlık altında asıl incelenen ise uykusuzluğa neden olabilecek faktörler ve hayatımıza olası etkileri konusudur. SPSS programı kullanılarak anket sonucuna çıkan veriler incelenmiş ve analizler yapılmıştır. İlk olarak tanımlayıcı istatistikler ele alınmıştır. Burada değişkenlerin frekansları ve dağılımları incelenmiştir. İkinci olarak normallik testleri yapılarak değişkenlerin normal dağılıp dağılmadığı incelenmiştir. Üçüncü olarak hipotezler kurulup değişkenlerin bir önceki bölümde elde ettiğimiz normallik sonuçlarına göre uygun olan testler yapılarak hipotezler yorumlanmıştır. İlk olarak, iki değişken arasında ortalamamız anlamlı olup olmadığını test eden bağımsız örneklem t-testi ve mann-whitney u testi uygulanmıştır. İkinci olarak, iki değişken arasında anlamlı bir ilişki olup olmadığını test eden ki-kare testi ve anlamlı bir ilişki olup olmadığına beraber bu ilişkinin yönünü de test eden korelasyon testi uygulanmıştır. Yapılan testler sonucunda göze çarpan sonuçlar olmuştur. Genel olarak uykusuzluğu etkileyen faktörlerin birden fazla olduğu ve bunların özellikle; cinsiyet, medeni durum ve çalışma durumu değişkenlerine göre farklılıkları görülmüştür.

KAYNAKÇA

- [1] Şenol V, Soyuer F, v.d. (2012), "Adolesanlarda Uyku Kalitesi ve Etkileyen Faktörler", Kocatepe Tıp Dergisi, Cilt 13, No:2, s.93-102
- [4] Mehta, C. R., & Patel, N. R. (2011). IBM SPSS exact tests. Armonk, NY: IBM Corporation.
- [5] Hair, J. F., Black, W. C., Babin, B. J., Anderson, R. E., & Tatham, R. L. (2013). Multivariate Data Analysis: Pearson Education Limited.



İSTATİSTİK BÖLÜMÜ

SEÇİLİ MAKROEKONOMİK FAKTÖRLERİN S&P 500 ENDEKSİNE ETKİSİ

İREM TUĞRUL 18023010

Danışman: Doç. Dr. Serpil KILIÇ DEPREN

ÖZET

S&P 500 endeksi Amerika'da yer alan 500 büyük firmanın hisselerini taşıyan bir endeks olarak belirtilebilir. Bu nedenle borsa üzerinden yatırım amaçlı güntümüzde tasarruf sahipleri tarafından oldukça önemli bir yere sahip olduğunu söyleyebiliriz. Bu çalışmada seçili makroekonomik faktörlerden Amerika'da faiz oranları, para arzı, petrol fiyatı ve döviz kuru faktörlerinin S&P 500 endeksi üzerindeki etkisini araştırmak olacaktır. Çalışmada, zaman serileri analizi kullanılmıştır. Serilerin durağanlığı ADF, KPSS ve PP birim kök testleri ile araştırılmıştır. Bu bağlamda S&P 500 endeksine olan etkinin Johansen eşbütünlüme yaklaşımıyla uzun dönem tahminleri yapılmıştır. Değişkenler arasındaki ilişkileri Granger nedensellik testini kullanarak kısa dönemli ilişkinin varlığı incelenmiştir. Johansen eşbütünlüme testiyle değişkenler arasında uzun dönemli bir ilişkinin olduğu sonucuna ulaşılmıştır. Granger nedensellik testi sonucunda da para arzı değişiminin S&P 500 endeksi değişimini üzerinde anlamlı bir nedensel etkisi olduğu sonucuna varılmıştır.

S&P 500 Kavramı

S&P 500, Amerika Birleşik Devletleri borsasında listelenen en büyük 500 şirketin hisse senedi performansını izleyen bir borsa endeksidir. En çok takip edilen hisse senedi endekslerinden biridir. Endeks 500 ileri gelen şirket kapsar ve mevcut piyasa değerinin yaklaşık % 80'ini oluşturur. S&P 500 endeksi, halka açık ağırlıklı bir endekstir.

Endeks Nedir?

Borsada işlem gören hisse senetlerinin fiyat ve getirilerinin performanslarını ölçmek için kullanılan göstergedir. Borsa endeksleri hissecedarları fiyat hareketlerinden yola çıkarak borsanın genel trendlerini belirlemesine kullanılır.

S&P 500 Endeksi Nedir?

S&P 500 endeksi, Amerika Birleşik Devletleri borsasında işlem gören en büyük 500 şirketin performansını ölçen bir borsa endeksidir. Bu endeks dünyadaki en önemli hisse senedi endekslerinden biridir. ABD ekonomisinin genel sağlığına dair de önemli bilgiler verir.

Makroekonomik Göstergeler

Bir ekonominin genel sağlığı, büyümesi ve performansını ölçmek için kullanılan göstergelerdir. Hisse senetleri ve diğer finansal varlıklar üzerinde önemli etkilere neden olabilir. Çalışmada para arzı, bir ekonomideki para biriminin miktarını ve likiditesini belirler. Amerika Birleşik Devletleri'nde para arzı, ABD Merkez Bankası (Federal Reserve) tarafından yönetilen para politikası aracılığıyla belirlenir. ABD'de petrol fiyatları, genellikle West Texas Intermediate (WTI) ve Brent petrolü fiyatları olarak ölçülür. Petrol fiyatları zaman içinde değişebilir ve gelecekte ne olacağı belirsizdir. Döviz kuru, Amerika Birleşik Devletleri'nde döviz kuru, Amerikan dolarının başka bir para birimi karşındaki değerini ifade eder. Faiz oranı, ABD 10 yıllık tahvil faizi denildiğinde ABD Hazine'nin çıkardığı 10 yıllık tahvilin piyasa faiz oranı anlaşılır.

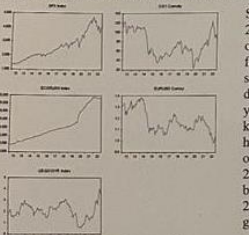
ARAŞTIRMANIN VERİ SETİ

Çalışmada 31/01/2012-30/11/2022 tarihleri arasındaki aylık veriler kullanılarak seçili makroekonomik göstergelerin S&P 500 endeksine etkisi incelenmiştir. Makroekonomik değişkenler olarak faiz oranı, para arzı, petrol fiyatı ve döviz kuru seçilmiştir. Çalışmanın bağımlı değişkeni S&P 500 endeksi, bağımsız değişkenleri ise faiz oranı, para arzı, petrol fiyatı ve döviz kuru'dur. Bu çalışmada 131 aylık veri kullanılmıştır. Çalışmanın verisi Bloomberg'ten alınmıştır. S&P 500 endeksi, faiz oranı, para arzı, petrol fiyatı ve döviz kuru'nun kapama fiyatları baz alınmıştır.

TANIMLAYICI İSTATİSTİKLER

	S&P İNDEKSİ	CO1	COM2	ECORUSN	EURUSD	USGG10YR
Mean	2613.978	74.78289	14402.28	1.179205	2.081238	
Median	2423.416	67.20020	13475.10	1.150600	2.134800	
Minimum	4786.180	122.8800	21945.10	1.388700	4.047000	
Maximum	120.320	22.40200	9877.500	0.902200	0.520200	
Std. Dev.	916.4229	26.07435	3704.138	0.207889	0.870738	
Skewness	0.86476	0.28650	0.788910	0.490777	-0.050613	
Kurtosis	2.49969	1.82678	2.32979	2.32649	3.180252	

S&P 500 endeksi 2012-2015 döneminde minimum 1310.330, maksimum 4766.180 değerini almıştır. S&P 500 endeksi ortalaması bahsedilen dönemde 2613.979, standart sapması ise 916.4029'dur.



S&P 500 endeksi için 2019-2020 yılları arasında düşüş yaşandığı 2020'den 2021'in ortalarına kadar yükseliş geçip 2021'in ortalarında zirve yapıp sonrasında hafif düşüş geçmiştir. Petrol fiyatı (Co1 Comdy), 2012'den 2014'ün ortalarına kadar belirli bir seviyede dalgalanma gösterirken 2014'ün ortalarında ani bir düşüş yaşamıştır. 2020 yılında da ani bir düşüş yaşayıp sonra yükseliş geçmiştir. Para arzı (Ecorusn Index), 2012'den 2020'ye kadar sürekli olarak artmıştır. 2020'den 2022'ye kadar yükselişi hız kazanmıştır. Döviz Kuru (EurUSD Curncy), 2012'den 2014'ün ortalarına kadar belirli seviyelerde dalgalanmıştır. 2014'ten 2015'e kadar düşüş yaşamıştır. 2021'den 2022'ye kadar da belirli bir seviyede düşmüştür. Faiz Oranı (Usgg10yr Index), 2018'den 2020 ortalarına kadar azalmıştır. 2020 ortalarından 2022'ye kadar genel olarak yükseliş geçtiğini söyleyebiliriz.

BİRİM KÖK TESTLERİ

Null Hypothesis: DGPX_INDEKSI has a unit root	Exogenous: Constant	Lag Length: 12 (Automatic - based on SIC, maxlags=12)
1-Statistic	Prob. >	
Augmented Dickey-Fuller test statistic	(-13.20604)	0.0000
T-test critical values	1% level	-3.483232
	5% level	-2.862600
	10% level	-2.577888

*MacKinnon (1996) one-sided p-values.

ÖKÜL REGRESYON MODELİ

$$SPX = \beta_0 + \beta_1 CO1 + \beta_2 ECORUSN + \beta_3 EURUSD + \beta_4 USGG10YR + \epsilon_t$$

Modeldeki Prob (F-Statistic) değeri, 0.05 ten küçük olduğu için model anlamlıdır. Bağımlı değişken olan S&P 500 endeksinin %21.44'lük kısmı modele dahil edilen değişkenler tarafından açıklanmıştır. Geri kalan kısım ise hata terimi aracılığıyla modele dahil edilmeyen değişkenlere aittir.

Variable	Coefficient	Std. Error	t-Statistic	Prob.
FARMSCO1_COMDY	7.98889	1.74284	4.58071	0.0000
FARMSCO2_COMDY	1.47896	0.63572	2.32614	0.0219
FARMSEURUSD_CUR	888.1928	421.8882	2.10634	0.0372
FARMSUSGG10YR_INDEX	0.14036	0.02243	6.251461	0.0000
C	11.92128	11.18870	1.06580	0.2746

VARSAYIMLAR

Modelde sağlanması gereken varsayımları sıralayalım.

Variable	Coefficient	Uncentered	Centered
FARMSCO1_COMDY	7.98889	1.28706	1.28704
FARMSCO2_COMDY	1.50284	0.71024	0.71024
FARMSEURUSD_CUR	1.77814	1.11830	1.11830
FARMSUSGG10YR_INDEX	0.14166	1.18277	1.18277
C	134.402	1.28847	NA

Çoklu doğrusal bağlantı varsayımıyla modelin incelenmesi yanda gösterilmiştir. VIF değeri 1 ile 5 arasında olması gerekmektedir. Gözlemlü üzere tüm değerler 1 ile 5 arasında değer almıştır. Bu nedenle çoklu doğrusal bağlantı sorunu yoktur. Bu şart sağlanmıştır.



Normallik testine göre probability değeri 0.05'ten büyükse kalıntılar normal dağılımı sonucuna ulaşılır. Bu testin sonucu probability değeri 0.6611 > 0.05 olduğu için kalıntılar normal dağılımıdır.

Burada ters kural mantığı geçerlidir. Probability değeri 0.0574 > 0.05'ten büyük olduğu için otokorelasyon problemi yoktur.

Burada da ters kural mantığı vardır. Bu hipotezi, probability değeri 0.0604 > 0.05'ten olduğu için kabul edilir. Bu demek oluyor ki, değişen varyans sorunu yoktur. Bir diğer ifadeyle sabit varyans vardır. Böylelikle bu varsayım da sağlanmış olmaktadır.

GRANGER NEDENSİLLİK TESTİ

Null Hypothesis	Obs	F-Statistic	Prob.
FARMSCO1_COMDY does not Granger Cause FARMSUSGG10YR_INDEX	128	0.2638	0.7792
FARMSCO2_COMDY does not Granger Cause FARMSUSGG10YR_INDEX	128	2.3078	0.1028
FARMSEURUSD_CUR does not Granger Cause FARMSUSGG10YR_INDEX	128	8.51378	0.0001
FARMSUSGG10YR_INDEX does not Granger Cause FARMSEURUSD_CUR	128	2.9965	0.1027
FARMSEURUSD_CUR does not Granger Cause FARMSCO1_COMDY	128	0.47320	0.6287
FARMSEURUSD_CUR does not Granger Cause FARMSCO2_COMDY	128	0.438495	0.50513
FARMSEURUSD_CUR does not Granger Cause FARMSUSGG10YR_INDEX	128	7.48118	0.0008
FARMSUSGG10YR_INDEX does not Granger Cause FARMSEURUSD_CUR	128	0.58312	0.5825
FARMSUSGG10YR_INDEX does not Granger Cause FARMSCO1_COMDY	128	3.47423	0.0320
FARMSUSGG10YR_INDEX does not Granger Cause FARMSCO2_COMDY	128	0.14086	0.8887
FARMSUSGG10YR_INDEX does not Granger Cause FARMSEURUSD_CUR	128	0.58422	0.5702
FARMSEURUSD_CUR does not Granger Cause FARMSUSGG10YR_INDEX	128	0.58312	0.5825
FARMSEURUSD_CUR does not Granger Cause FARMSCO1_COMDY	128	0.11911	0.9368
FARMSEURUSD_CUR does not Granger Cause FARMSCO2_COMDY	128	1.5787	0.2122
FARMSEURUSD_CUR does not Granger Cause FARMSUSGG10YR_INDEX	128	4.13002	0.0198
FARMSUSGG10YR_INDEX does not Granger Cause FARMSEURUSD_CUR	128	0.15428	0.8757
FARMSUSGG10YR_INDEX does not Granger Cause FARMSCO1_COMDY	128	4.13002	0.0198
FARMSUSGG10YR_INDEX does not Granger Cause FARMSCO2_COMDY	128	2.58772	0.1078
FARMSUSGG10YR_INDEX does not Granger Cause FARMSEURUSD_CUR	128	0.57374	0.5825

JOHANSEN EŞBÜTÜNLÜME TESTİ

Trace Test	Trace	Statistic	Prob. >
Trace Test indicates 5 cointegrating eq(s) at the 0.05 level			
* denotes rejection of the hypothesis at the 0.05 level			
*MacKinnon-Haug-Michelis (1999) p-values			

Kısa dönemli ilişkileri test etmek için Granger-Nedensellik testi uygulanmıştır. Nedensellik analizinin amacı değişkenler arasındaki ilişkinin varlığını test etmek ve ilişkinin yönünü belirlemektir. Granger nedensellik testi sonucunda para arzı değişiminin S&P 500 endeksi değişimi üzerinde anlamlı bir nedensel etkisi bulunmuştur.

HATA DÜZELTME MODELİ

Variable	Coefficient	Std. Error	t-Statistic	Prob.
FARMSCO1_COMDY	0.00000			
FARMSCO2_COMDY	0.00000			
FARMSEURUSD_CUR	0.00000			
FARMSUSGG10YR_INDEX	0.00000			
C	0.00000			

Uzun dönemdeki katsayıların olduğu denkleme aşağıdaki gibidir.
SPX = 33.03 CO1 + 0.83 ECORUSN + 710.73 EURUSD + 966.22 USGG10YR +115.86

SONUÇ

Bu çalışmada, seçili makroekonomik değişkenlerden para arzı, petrol fiyatı, döviz kuru ve faiz oranının S&P 500 endeksi üzerindeki etkisi araştırılmıştır. Çalışmada yapılan birim kök testleri sonucunda verilerin ilk halinde durağan olmadığı fakat birinci dereceden farkları alındığında durağan olduğu sonucuna ulaşılmıştır. Sonrasında çoklu regresyon modeli kurulumu ve modelin anlamlılığını test edilmiştir ve model anlamlı çıkmıştır. S&P 500 endeksine etki eden bağımsız değişkenlerin petrol fiyatı, para arzı, döviz kuru ve faiz oranı olduğu sonucuna varılmıştır. Granger nedensellik testi sonucunda para arzı değişiminin S&P 500 endeksi değişimini üzerinde anlamlı bir nedensel etkisi bulunmuştur. Johansen eşbütünlüme testi sonucunda %5 seviyesinde bağımlı (S&P 500 endeksi) ve bağımsız değişkenler (petrol fiyatı, para arzı, döviz kuru ve faiz oranı) arası uzun dönemli ilişkinin var olduğu sonucuna ulaşılmıştır. Kısa dönem dalgalanmalarına kadar strede uzun dönem dengisini yakaladığına da hata düzeltme modeli ile bulunmuştur.

KAYNAKÇA

[1] Sirucek, M. (2012). Macroeconomic variables and stock market: A review
[2] Dolajic, B. (2022). Seçili Makroekonomik Göstergelerin Borsa İstanbul (Bist-100) Endeksine Etkisi
[3] Aygün, M. (2022). Makroekonomik Değişkenler ile Borsa İstanbul Arasındaki Nedensellik İlişkisi: Toda-Yamamoto Testi

İSTATİSTİK BÖLÜMÜ

Kredi Kartı Başvurusu Veri Setinde Keşifsel Veri Analizi ve Makine Öğrenmesi Uygulaması

Öğrenci: Beyzagül GÖKSU - 18023031

Danışman: Prof.Dr. Ali Hakan BÜYÜKLÜ

Finansal hizmetlerin giderek dijitalleştiği günümüzde, kredi kartları bireylerin günlük harcamalarını yönetmeleri ve finansal ihtiyaçları için kolay bir ödeme aracı sağlamaları açısından önemli bir rol oynamaktadır. Bankalar ve finansal kuruluşlar, kredi kartı başvurularını değerlendirirken risk analizi yaparak başvuru sahiplerine kredi kartı verme veya reddetme kararı vermektedirler. Bu süreçte, teknolojinin de gelişmesiyle geleneksel yöntemlerin yanı sıra veri analitiği ve makine öğrenmesi gibi teknikler doğru ve hızlı kararlar almak için kullanılmaya başlanmıştır.

Bu çalışma, keşifsel veri analizi ve makine öğrenmesi tekniklerini kullanarak kredi kartı başvurularını incelemeyi hedeflemektedir. Kredi kartı başvurusu veri seti üzerinde gerçekleştirilen bu analizler, başvuruların kabul edilme etkenlerinin anlaşılmasına ve onaylanma veya reddedilme olasılıklarının tahmin edilmesine yardımcı olacaktır. Bununla birlikte, bu çalışma, finansal kuruluşlara ve kredi kartı sağlayıcılara; başvuru sürecini optimize etme, risk yönetimini iyileştirme ve müşteri memnuniyetini artırma potansiyeli sunmaktadır.

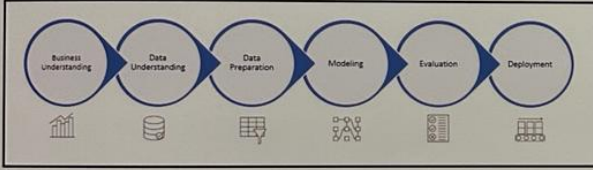
Kredi Kartı Nedir?

Kredi limiti, bir kredi kartı veya bir kredi hattı gibi kredi ürünleri için belirlenir. Kredi kartı limiti, kredi kartı sahibinin belirli bir dönemde kullanabileceği maksimum tutarı ifade eder. Kredi hattı limiti ise, müşterinin belirli bir süre boyunca kullanabileceği maksimum kredi miktarını ifade eder. Kredi limiti, müşterinin finansal durumuna ve risk profiline bağlı olarak belirlenir. Kredi limiti aynı zamanda müşterinin kredi notu, gelir düzeyi ve ödeme geçmiş gibi faktörlerle ilişkilidir. Kredi limiti, müşterinin kredi kullanma kapasitesini gösterir ve müşteri bu limit dahilinde kredi alabilir veya harcama yapabilir.

Bankalarda Risk Yönetimi

Saunders ve Cornett (2014), banka yönetiminin hissedarlar için getirileri artırmayı amaçladığını ve bunun riskin artmasıyla gerçekleştiğini belirtmektedir. Bankaların karşılaştığı birçok risk bulunmaktadır. Bunlar kredi riski, faiz oranı riski, piyasa riski, bilanço dışı risk, teknoloji ve operasyonel risk, döviz riski, ülke veya egemen riski, likidite riski, likidite riski ve iflas riski olarak sıralanabilir. Bu riskleri etkin bir şekilde yöneterek, bankalar daha iyi performans gösterebilirler.

Keşifçi Veri Analizi



Veri Keşfi: İlk adım veri setinin anlaşılması ve keşfedilmesidir. Bu adımda, verilerin kaynakları, toplama yöntemleri ve veri noktalarının neyi temsil ettiği gibi bilgilere dikkat edilir.

Veri Temizliği: Veri setindeki hataları, eksik değerleri ve anormal verileri tespit etmek ve bunları düzeltmek için veri temizliği yapılır. Bu adımda, verilerin tutarlılığı ve doğruluğu sağlanır. Eksik veriler, aykırı değerler veya çelişkili veriler gibi problemler belirlenir ve uygun bir şekilde ele alınır.

Veri Görselleştirme: Veri setinin keşfedilmesi ve analizi için görselleştirmeler kullanılır. Grafikler, histogramlar, dağılım eğrileri ve kutu grafikleri gibi görsel araçlar kullanılarak veri setindeki desenler ve ilişkiler daha iyi anlaşılır. Görselleştirmeler, veri setindeki eğilimleri ve varyasyonları göstererek analistlere daha fazla içgörü sağlar.

İlişkilerin ve Desenlerin Analizi: Bu adımda, veri setindeki değişkenler arasındaki ilişkileri ve desenleri keşfetmek için istatistiksel analiz yöntemleri kullanılır. Korelasyon analizi, regresyon analizi, kümeleme ve faktör analizi gibi teknikler kullanılarak veriler arasındaki bağlantılar ve desenler ortaya çıkarılır. Bu adım, veri setindeki önemli faktörleri belirlemek için kullanılır.

Makine Öğrenmesi Hakkında

Makine öğrenmesi, belirli bir problemi çözmek için birçok farklı algoritma ve teknik kullanır. Bu algoritmalar, veriyeye dayalı örüntüleri bulmak ve modele öğrenme yeteneği kazandırmak için kullanılır. Makine öğrenmesi modelleri, genellikle eğitim veri setleri üzerinde eğitilir, ardından bu öğrenme sürecinden elde edilen bilgilerle yeni verileri tahmin etmek veya sınıflandırmak için kullanılır. Bu çalışmada makine öğrenmesi modellerinden LightGBM ve CatBoosting öne çıkmıştır.

LGBM: Ağaç tabanlı bir modeldir ve özellikle sınıflandırma ve regresyon problemlerinde kullanılır. Gelişmiş bir özellik olan "leaf-wise" büyüme stratejisini kullanarak ağaçları hızlı bir şekilde oluşturur. Bu strateji, enformasyon kazancı maksimize edilerek ağaç yapısının derinleşmesini sağlar.

CatBoosting: Hızlı ve ölçeklenebilir bir algoritmadır ve büyük veri setlerinde ve yüksek boyutlu özelliklerde etkili bir şekilde çalışabilir. Overfitting dirençlidir ve doğru hiper parametre ayarlarıyla yüksek performans sağlayabilir.

Veri Hakkında

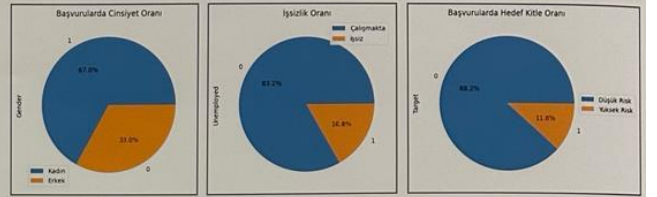
Araştırmada kullanılan veri seti Kaggle web sitesinden alınmıştır.

Veri setinin adı "Credit Card Approval Prediction"dir. İki ayrı veri setinden oluşmaktadır.

application_record.csv: 438.557 müşterinin kişisel verilerini içermektedir.

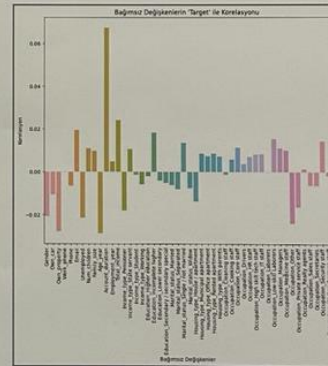
credit_record.csv: Belli müşterilerin 12 aylık kredi durumlarını içermektedir.

Her müşterinin kredi durumları bilinmediğinden çalışmada kullanılmak üzere kredi durumları bilinen müşteriler üzerinden veri setleri birleştirildi. Veri işlendikten sonra çalışma 36449 müşterinin kişisel verileri ve kredi skorları ile devam etti.



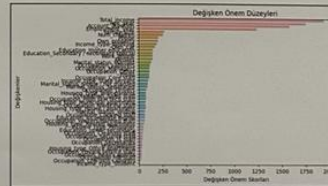
Çalışmada müşterilerin %67'si Kadın, %33'ü erkektir. Çoğunluk iş sahibi iken %16.8'i çalışmamakta dolayısıyla bu müşterilerin kredi başvurusunun reddedilmesi beklenmektedir. Kredi skoru olarak belirlenen Hedef değişkenine göre yüksek risk grubunda olan %11.8 müşteri reddedilmiştir. Yani işsiz olmasına rağmen kredi başvurusu kabul edilen %5 müşteri vardır.

Makine Öğrenmesi Uygulamaları

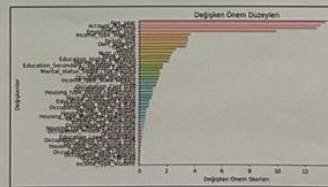


En iyi pozitif ilişkilerin sırasıyla hesap yaşı, yıllık gelir, e-posta ve yüksek öğrenim olduğu gözlemlenmektedir. En güçlü negatif ilişkiler ise yaş, gayrimenkul sahipliği, işsizlik ve cinsiyet olarak gözlemleniyor.

Sonuç olarak çalışmada kullanılan makine öğrenmesi modellerinin en ilişkili değişkenler olarak hesap yaşı, yıllık gelir, e-posta ve yüksek öğrenim değişkenlerini öngörmesi bekleniyor.



LGBM: Modelin doğruluk skoru 0.8853'tür. Yıllık gelir, yaş, çalışma süresi ve hesap süresi değişkenlerini modeli en çok etkileyen değişkenler olarak öngörmüştür.



CatBoosting: Modelin doğruluk skoru 0.8855'tür. LGBM modeli gibi yıllık gelir, yaş, çalışma süresi ve hesap süresi değişkenlerini modeli en çok etkileyen değişkenler olarak öngörmüştür.

KAYNAKÇA

- 'Credit Card Approval Prediction by Using Machine Learning Techniques' M. P. C. Pei University of Colombo School of Computing
- 'Predicting Credit Card Approvals using Machine Learning' by Lopa Nayak, Aman Sangal
- Credit Card Approval Prediction Data from Kaggle



YFU FEN EDEBİYAT FAKULTESİ

İSTATİSTİK BÖLÜMÜ

TÜRKİYE'DE DEPREM BİLGİSİ VE BİLİNCİ ÜZERİNE BİR ARAŞTIRMA

EMİRHAN BAYRAK 18023033

Danışman: Dr. Öğr. Üyesi Doğan YILDIZ

Deprem Dünya genelinde en tehlikeli ve en çok zarar veren doğal afetlerden biridir. Birçok ülke gibi Türkiye de deprem riski yüksek bir ülkedir ve sık sık depremlerle karşı karşıya kalmaktadır. Bu durum, Türkiye'de deprem bilincinin önemini artırmaktadır. Bu tezde, Türkiye'deki deprem bilincini ölçülmektedir. Bu amaçla, farklı yaş gruplarından ve farklı eğitim düzeylerinden katılımcılara toplamda 18 soruluk bir anket çalışması yapılmıştır. Ankette, deprem hakkındaki bilgi düzeyi, deprem hazırlığı yapma eğilimi, deprem anında nasıl davranılması gerektiği konusunda bilgi sahibi olma durumu gibi konulara yer verilmiştir. Ankete katılan 192 kişinin cevapları, SPSS programı kullanılarak analiz edilmiştir. Analizde, Cronbach testi ile güvenilirlik testi edildi, aynı zamanda çapraz tablolar kullanılmıştır ve bu çapraz tablolar Ki-Kare testi ile desteklenmiştir. Bu yöntemlerle sorular arasındaki ilişki ilişkiler arasında hipotezler test edilmiştir. Analiz sonuçları elde edilen SPSS çıktıları üzerinden yorumlanmıştır.

DEPREM

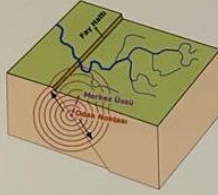
Depremler, dünya tarihinde sayısız felakete neden olmuştur. İnsanlar, depremlerin doğrudan ve dolaylı etkileri yoluyla binlerce yıl boyunca zarar görmüş ve yaralanmışlardır. Depremler, binaların ve diğer yapıların yıkılmasına, su kaynaklarının kirlenmesine, elektrik ve ulaşım ağlarının hasar görmesine ve hatta ekonomik çöküntülere neden olabilir.

DEPREM TÜRLERİ

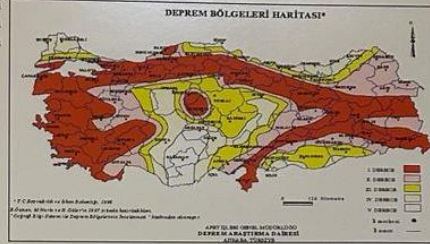
Günümüzde bilinen 7 farklı deprem türü vardır bunlar tektonik depremler, volkanik depremler, çökme depremler, yer kayması depremleri, heyelan depremleri, denizaltı depremleri, ve yapı depremleridir. Dünyadaki depremlerin yaklaşık %95'ini "Tektonik Depremler" oluşturur.

Tektonik Depremleri açıklayacak olursak;

Yer kabuğundaki levhaların sürtünmesi veya birbirinden uzaklaşması sonucu meydana gelen depremlerdir.



DEPREM VE TÜRKİYE



Türkiye, Kuzey Anadolu Fay Hattı, Doğu Anadolu Fay Hattı ve Batı Anadolu Fay Hattı gibi önemli fay hatları üzerinde yer alır. Bu fay hatları, sık sık depremlerin meydana geldiği bölgelerdir ve büyük depremlere neden olabilirler.

Ayrıca, Marmara Denizi bölgesi, Türkiye'nin en yüksek deprem riskine sahip bölgelerinden biridir. Bu bölgelerdeki depremler, önemli hasarlara ve kayıplara yol açabilir.

ANKET SONUÇLARI VE CRONBACH TESTİ SONUCU

DEMOGRAFİK SORULARIN SONUÇLARI				DEPREM İLE İLGİLİ SORULARIN SONUÇLARI			
Sorular	Şıklar	N	%	Sorular	Şıklar	N	%
Cinsiyetiniz nedir?	Erkek	103	46,4%	Depremlerin nasıl meydana geldiğini biliyor musunuz?	Evet, tam olarak biliyorum	11	5,7%
	Kadın	89	39,6%		Biraz biliyorum	73	37,9%
Yaş aralığınız nedir?	18-25	135	70,3%	Hüç bilmiyorum	2	2,6%	
	26-35	29	15,1%	Kendinizi ve ailenizi bir depreme hazırlamak için neler yapmaya çalıştığınızı biliyor musunuz?	Evet, her şeyi biliyorum	85	44,3%
	36-50	12	6,3%		Biraz şey biliyorum	99	51,4%
	50 ve üzeri	16	8,3%	Hiçbir şey biliyorum	8	4,2%	
Eğitim durumunuz nedir?	İlkokul	5	2,6%	Deprem anında yapılması gereken işlemleri biliyor musunuz?	Evet, tam olarak biliyorum	89	46,3%
	Ortaokul	11	5,7%	Biraz biliyorum	86	44,8%	
	Lise	30	15,6%	Hiçbir şey biliyorum	7	3,6%	
	Üniversite	141	73,4%	Deprem olabileceği bir bölgede yaşıyor musunuz?	Evet, depremler olabileceği bir bölgede yaşıyorum	139	72,4%
	Üniversite Lisansüstü	5	2,6%	Hayır, depremler olabileceği bir bölgede yaşamıyorum	5	2,6%	
Çalışma durumunuz nedir?	Kamu Sekreteri	11	5,7%	Bilinmeyen depreme dayanıklı olduğuna düşüncünüz var mıdır?	Evet	104	54,2%
	Özel Sekreter	62	32,3%		Hayır	88	45,8%
	Öğrenci	99	51,6%	Bilinmeyen depremleri için belediye ve deprem dayanıklılık testi yaptırdığınız mıdır?	Evet	52	27,1%
	Evlü Hanımı	7	3,6%		Hayır	140	72,9%
	Kendi işinin patronuyum	3	1,6%	Yaşadığınız bölge için belediye ve depreme dayanıklı olduğuna düşüncünüz var mıdır?	Evet	62	32,3%
	Çalışmıyorum	10	5,2%		Hayır	130	67,7%
Gelir durumunuz nedir?	0-4.500	97	50,5%	Deprem sırasında meydana gelen depremleri biliyor musunuz?	Evet	119	62,1%
	4.501-13.500	39	20,3%	Bilinmeyen depremleri biliyor musunuz?	Evet	75	39,1%
	13.501-20.000	41	21,4%	Deprem sırasında ve sonrasında yapılması gerekenleri biliyor musunuz?	Evet, sahibim	57	29,7%
	20.001-35.000	11	5,7%	Hayır, sahip değilim	135	70,3%	
	35.001 ve üzeri	4	2,1%	Deprem sırasında, binaya zarar vermişse, binanın yeniden yapıldığını biliyor musunuz?	Evet, bildim	131	68,2%
Yaşadığınız konutun tipi nedir?	Müstakil	17	8,9%	Deprem sırasında, binanın yeniden yapıldığını biliyor musunuz?	Evet, bildim	91	47,4%
	Apartman	133	69,3%	Hayır, bildim	131	68,2%	
	Site	33	17,2%	Hayır, bildim	131	68,2%	
	Yurt/Lojman	9	4,7%	Evet, tam olarak biliyorum	72	37,5%	
Hemenindeki kişi sayısı nedir?	1	9	4,7%	Biraz biliyorum	82	42,7%	
	2	33	17,2%	Hiçbir şey biliyorum	38	19,8%	
	3	44	22,9%				
	4	68	35,4%				
	5+	38	19,8%				

Reliability Statistics

Reliability Statistics	Value	df
Cronbach's Alpha	0,869	192
Cronbach's Alpha Based on Standardized Items	0,864	192
N of Items	18	

Bu anket için Cronbach alfa katsayısı hesaplandığında, 0,669 olarak bulunmuştur. Bu değer, anketin iç tutarlılık düzeyinin kabul edilebilir bir seviyede olduğunu göstermektedir. Literatürde, Cronbach alfa katsayısının 0,70 veya üzeri olması genellikle kabul edilebilir bir iç tutarlılık işareti olarak kabul edilirken, bazı durumlarda 0,60 veya daha düşük değerler de kabul edilebilir olabilmektedir. Dolayısıyla, 0,669 Cronbach alfa değeri, anketin güvenilirliğini dair olumlu bir kanıt sunmaktadır.

ÇAPRAZ TABLO ANALİZİ VE Kİ-KARE TESTİ

Çapraz tablolar, verilerdeki değişkenlerin dağılımını anlamak ve aralarındaki ilişkiyi kontrol etmek için Ki-Kare testiyi destekler. Bu doğrultuda anket çalışmasında demografik sorular ile deprem ile ilgili sorular arasında çapraz tablo Kare testleri uygulandı ve bu uygulamalara örnek vermem gerekirse; 4. soru "Deprem olabilecek bir bölgede yaşıyor musunuz?" ile 9. soru olan "Deprem sırasında, deprem sırasında ve sonrasında yapılması gerekenleri anlatan bir acil durum çantası kullanıyor musunuz?" arasındaki ilişkiyi kontrol etmemiz gerektiği ortaya çıktı.

H₀: Soru4 ile Soru9 arasında ilişki yoktur.

H₁: Soru4 ile Soru9 arasında ilişki vardır.

$\alpha = 0,05$

Karar kuralına göre bakıldığında p-değerimiz (0,05) güvenilirlik düzeyimizden küçük olduğundan dolayı H₀'ı reddetmeliyiz. 4. soru 9. soru arasında ilişki vardır. Deprem bölgesinde yaşayan insanların sayısı bir değişim acil durum çantasına sahip olan insanların sayı etkileyeceğini söyleyebiliriz.

SONUÇ VE ÖNERİ

Anket çalışması, Türkiye'deki deprem bilincinin değerlendirilmesi amacıyla yapılmıştır. Katılımcılar deprem hakkındaki bilgileri sahipleriyle paylaşarak deprem tedbir alma konusunda yeterli kalmamıştır. Cinsiyet, yaş ve eğitim düzeyi gibi demografik faktörler, deprem bilinci ve tedbir alma davranışları üzerinde etkilidir. Deprem bilincini artırmak için eğitim programları bilgilendirme kampanyaları önemlidir. Bu çalışma, deprem bilincini artırmak ve toplumun daha iyi hazırlanması sağlamak için eğitim programları ve bilgilendirme kampanyalarının gerekliliğini vurgulamaktadır.



SİSMOGRAF

Deprem ölçümleri, sismograflar adı verilen cihazlar kullanılarak yapılır. Sismograflar, deprem sırasında yer kabuğundaki hareketleri kaydeden hassas aletlerdir. Sismografların temel işlevi, deprem anındaki titreşimleri ölçmek ve kaydetmektir. Sismograflar, yerin altındaki kayaların hareketlerini ölçmek için kullanılan özel bir dizi sensöre sahiptir. Bu sensörler, yer kabuğundaki herhangi bir titreşimi ölçmek için tasarlanmıştır.

Sismografin çalışma prensibi oldukça basittir: sensörler, deprem sırasında yer kabuğundaki titreşimleri algılar ve bu titreşimleri diğt üzerine çizerek kaydeder. Deprem ölçümleri, sismografların kaydettiği verilerin analiziyle yapılır. Bu veriler, depremin büyüklüğü ve merkez üssünün konumu hakkında bilgi sağlar. Deprem ölçümleri, genellikle Richter ölçeği adı verilen ölçekte ifade edilir ve bu ölçeğe bir artış, depremin şiddetinin 10 kat arttığını gösterir.

DEPREM ÖNCESİNDE YAPILMASI GEREKENLER

Deprem öncesinde kendimizi güven altına tutmak için yapılabilecek ve alınabilecek birçok tedbir vardır. Bunlardan bazıları şu şekilde sıralayabiliriz;

- Deprem sigortası yaptırmak:** Ev, iş yeri ve diğer mülkler için deprem sigortası yaptırmak önemlidir.
- Acil durum çantası hazırlamak:** Deprem sırasında veya sonrasında kullanılmak üzere bir acil durum çantası hazırlamak faydalı olabilir. Bu çanta içinde gıda, su, ilaçlar, ilk yardım malzemeleri, çakmak, el feneri ve bataryalı gibi temel eşyalar bulunmalıdır.
- Yüksek yerlerde güvenliğini sağlamak:** Deprem sırasında eşyaların düşmesini önlemek için, dolap ve rafların sabitlenmesi gerekmektedir.
- Acil durum planı yapmak:** Ailenizle veya ev arkadaşlarınızla bir acil durum planı yapmak, herkesin deprem sırasında nerede planacağı ve nasıl iletişim kuracağı konusunda bilgi sahibi olmasını sağlayacaktır.
- Bina güvenliğini kontrol etmek:** Deprem öncesinde bina güvenliğini kontrol etmek önemlidir. Bina güçlendirme işlemleri yapılması gereken yapılar varsa, bu çalışmaların yapılması gerekmektedir.

ACİL DURUM ÇANTASININ HAZIRLANMASI

de bulundurulacak acil durum çantası, doğal afetler gibi beklenmedik olaylarda hızlı bir şekilde hareket etmenize yardımcı olur: kritik durumlarda hayatta kalmayı sağlayabilir. Acil durum çantasında bulunması gerekenleri şu şekilde sıralayabiliriz: su, gıda, c yarıdım kitleri, hijyen malzemeleri, giyim, ışılandırma ve iletişim malzemeleri. Bu malzemeler yaşanan bölgeye ve kişinin istra ihtiyaçlarına göre düzenlenebilir. Bu ürünlerden mümkün olduğunca 3 ile 5 gün arasında tüketilebilecek miktarda olmalıdır.

DEPREM ANINDA YAPILMASI GEREKENLER

Deprem anında bireyin depremi hissetmesi, hissettikten sonra buna tepki vermesi ve hissettikten sonra buna vereceği tepkiyle birlikte yapması gereken önemli davranışlar vardır. Bunları şöyle sıralayabiliriz; sakin olmak, hemen korunmak, kapı ve pencere anlarına yakın olmak, asansör kullanmamak ve acil çıkış kapılarına yakın durmak gibi birkaç örnek verilebilir.

DEPREM SONRASINDA YAPILMASI GEREKENLER

Deprem sonrasında yapılması gerekenler, insanların hayatta kalmalarına, yaralanmalarının önlenmesine ve hasarın azaltılmasına yardımcı olacak önlemleri içerir. Deprem sonrasında yapılması gerekenler şu şekildedir; kendinizi ve diğerlerinizi güvende tutun, ararlara yardım edin, elektrik, gaz ve su bağlantılarını kesin, haberleri takip edin, yapısal hasarın kontrolü, geçici barınak sağlayın ve acil durum çantası kullanın.

KAYNAKÇA

- 1] Freed, A. M. (2007). Afterslip (and only afterslip) following the 2004 Parkfield, California, earthquake. Geophysical Research Letters, 34(6).
- 1] Akkaş, H. H. (2023). Deprem Bilinci. Kahramanmaraş Merkezli Depremler Sonrası İçin, 23.



YTU FEN EDEBİYAT FAKÜLTESİ İSTATİSTİK BÖLÜMÜ

KURAL BAZLI VE DENETİMSİZ ÖĞRENME YÖNTEMLERİ İLE MÜŞTERİ SEGMENTASYONU

Mertcan DEMİREL 20023603

Danışman: Prof. Dr. Ersoy ÖZ

ÖZET

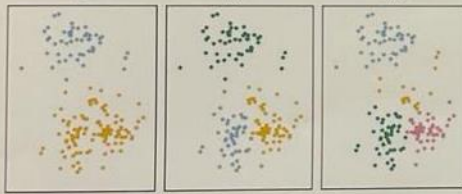
Bu tez çalışmasında kullanılan veri seti, İngiltere merkezli çevrimiçi işlem veri setidir. Veri seti, 01/12/2010 ile 09/12/2011 tarihleri arasındaki tüm işlemleri kapsamaktadır ve toptancı müşterileri içermektedir. Veri setinde 37 farklı ülke bulunmaktadır. Ancak, daha sağlıklı sonuçlar alabilmek için veri seti, "İngiltere", "Fransa", "Almanya", "İsviçre", "İspanya", "Portekiz", "Belçika" ülkeleri baz alınarak indirilmiştir. Bu veriler, müşteri segmentasyonu için kullanılmıştır ve k-ortalama ve hiyerarşik kümeleme yöntemleri Python programlama dili üzerinde geliştirilmiştir. Her bir ülke için ortalama 5 ya da 6 ideal kümeleme sayısı belirlenmiş ve ülkeler bu sayıya göre kümelendirilmiştir. Segmentasyon sonucunda, bu çalışma Streamlit kütüphanesi kullanılarak bir internet sitesi olarak yayımlanmıştır. Bu çalışma, müşterileri segmente ederek hangi kümeye yatırım yapılması gerektiğini belirlemeye ve sadık müşteri grubunu bulmaya yardımcı olabilecek değerli bir kaynak olabilir.

RFM ANALİZİ

RFM analizi, "Recency" (yenilik), "Frequency" (sıklık) ve "Monetary" (para) kelimelerinin baş harflerinden oluşan bir analiz yöntemidir. Bu yöntem, müşteri satın alma davranışlarını üç kategoriye ayırarak değerlendirir. Yenilik, müşterinin son satın alma işleminden itibaren geçen süreyi ifade eder ve yeni satın almalar daha değerli kabul edilir. Sıklık, müşterinin belirli bir zaman diliminde kaç kez satın alma yaptığını gösterirken, para ise müşterinin işletmeye yaptığı toplam harcamayı ifade eder. RFM analizi, müşteri davranışlarını anlamak için kullanılan ve işletmelere pazarlama stratejilerini optimize etmede yardımcı olan bir yöntemdir.

K-ORTALAMA KÜMELEME

K-ortalama kümeleme yöntemi, verileri iki aşamada kümelere ayırmak için çalışır. İlk olarak, veriler rastgele seçilen merkezi noktalara atanır ve her veri noktası en yakın merkezi noktaya atanır. Bu adım, verilere en yakın merkezi noktaya atanan kümelere oluşmasını sağlar. Daha sonra, her kümeye ait verilerin ortalaması alınarak yeni merkezi noktalar belirlenir. Bu adım, verilerin yeniden düzenlenmesiyle yeni bir kümeleme işlemi yapılır ve merkezi noktaların konumu ve kümelere içeriği değişerek daha doğru bir kümeleme sağlanır.



HİYERARŞİK KÜMELEME

Hiyerarşik kümeleme yöntemi, aglomeratif ve bölücü yaklaşımlar olmak üzere iki farklı şekilde uygulanabilir. Aglomeratif yaklaşımda, veriler başlangıçta ayrı kümelere yer alır ve bu kümeler birleştirilerek daha büyük kümeler oluşturulur. Bölücü yaklaşımda ise, tüm veriler bir kümede yer alır ve bu küme alt kümeler halinde bölünür. Hiyerarşik kümeleme yöntemi, verilerin gruplandırılmasında dendrogram adı verilen bir yapıyı kullanır. Bu yöntem, verilerin doğrusal olmayan yapısı nedeniyle k-ortalama yönteminden daha doğru sonuçlar üretebilir, ancak büyük veri kümeleri için yüksek hesaplama maliyeti gibi bazı dezavantajları da bulunmaktadır.

Dendrogramlar, hiyerarşik kümeleme yöntemlerinde kullanılan bir görselleştirme aracıdır ve veri noktalarının kümeleme sürecindeki birleşmelerini gösterir. Hiyerarşik kümeleme, benzerlik veya uzaklık ölçütlerine dayanarak veri noktalarını adım adım birleştirerek veya bölerek kümeler oluşturur. Dendrogramlar, bu birleşim ve bölünme adımlarını göstererek, küme sayısı ve yapıları hakkında bilgi sağlayarak kümeleme sonuçlarını analiz etmek ve uygun küme sayısını belirlemek için kullanışlı bir araç olarak kullanılır.

VERİ SETİ

Irvine, Kaliforniya Üniversitesi'ndeki Makine Öğrenimi ve Akıllı Sistemler Merkezi'nden elde edilen "Online Retail" adı bir veri seti kullanılmıştır. Veri seti, toplamda 541.908 veri içeren 8 değişken içerir. Bunlar; "InvoiceNo", "StockCode", "Description", "Quantity", "InvoiceDate", "UnitPrice", "CustomerID", "Country" adındaki değişkenlerdir.

KAYNAKÇA

- [1] Kotler, P., & Keller, K. L. (2016). Marketing Management (14th edition). Shanghai: Shanghai People's Publishing House.
- [2] Mitchell, T. M. (2007). Machine learning (Vol. 1). New York: McGraw-hill.

KÜMELEME ANALİZİ

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
from sklearn.preprocessing import StandardScaler

# Veri setini yükleyelim
df = pd.read_csv('data.csv')

# Veriyi temizleyelim
df.dropna(inplace=True)

# Veriyi ölçeklendirelim
scaler = StandardScaler()
df_scaled = scaler.fit_transform(df)

# K-ortalama kümeleme yapalım
kmeans = KMeans(n_clusters=5, random_state=0)
clusters = kmeans.fit_predict(df_scaled)

# Silüvyo skorunu hesaplayalım
silhouette_scores = silhouette_score(df_scaled, clusters)

# Silüvyo grafiğini çizelim
plt.figure(figsize=(10, 8))
sns.scatterplot(x=df_scaled[:, 0], y=df_scaled[:, 1], hue=clusters)
plt.title('K-ortalama Kümeleme Sonuçları')
plt.show()
```

Veri setindeki en son alışveriş tarihine 2 gün eklenerek analiz tarihi belirlendi ve "today_date" değişkenine saklandı. Ardından RFM metrikleri hesaplandı ve skorlar 1-5 aralığına indirilerek karşılaştırılabilir hale getirildi. RFM skorlarının kullanılacağı segmentasyon için hesaplamalar yapıldı. Örneğin, R=1, F=2, M=3 ise RFM skoru "123" olacaktır. Bu işlemlerin ardından, her iki yöntemi de kullanabilmek için veri seti alt kümelere ayrıldı ve bu alt kümeler [0,1] aralığına normalize edilerek diğer algoritmalara uygun hale getirildi.

K-ORTALAMA KÜMELEME

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
from sklearn.preprocessing import StandardScaler

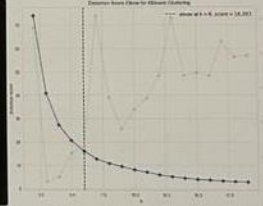
# Veri setini yükleyelim
df = pd.read_csv('data.csv')

# Veriyi ölçeklendirelim
scaler = StandardScaler()
df_scaled = scaler.fit_transform(df)

# K-ortalama kümeleme yapalım
kmeans = KMeans(n_clusters=5, random_state=0)
clusters = kmeans.fit_predict(df_scaled)

# Silüvyo skorunu hesaplayalım
silhouette_scores = silhouette_score(df_scaled, clusters)

# Silüvyo grafiğini çizelim
plt.figure(figsize=(10, 8))
sns.scatterplot(x=df_scaled[:, 0], y=df_scaled[:, 1], hue=clusters)
plt.title('K-ortalama Kümeleme Sonuçları')
plt.show()
```



Optimal küme sayısını belirlemek için k-ortalama yöntemi kullanılarak 2'den 20'ye kadar değişen küme sayılarına göre "elbow" yöntemi grafiği oluşturuldu. Elde edilen optimal sonuçlar temel alınarak bir model oluşturuldu ve bu modele dayanarak RFM skorları elde edildi. Sonuçları yorumladığımızda, "recency" değeri en düşük olan küme grubunun sadık müşterilere ait olduğunu söyleyebiliriz. Örneğin, 1. küme en düşük "recency" değerine sahip olup ortalama olarak "29" gün önce alışveriş yapmıştır. Ayrıca, bu küme toplamda "5" ürün satın almış ve bu ürünlerin toplam değeri "903" sterlin olarak kaydedilmiştir.

HİYERARŞİK KÜMELEME

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import HierarchicalClustering
from sklearn.metrics import silhouette_score
from sklearn.preprocessing import StandardScaler

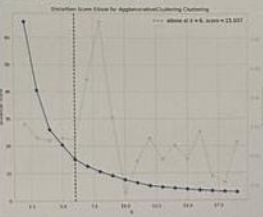
# Veri setini yükleyelim
df = pd.read_csv('data.csv')

# Veriyi ölçeklendirelim
scaler = StandardScaler()
df_scaled = scaler.fit_transform(df)

# Hiyerarşik kümeleme yapalım
hclust = HierarchicalClustering(n_clusters=5, linkage='ward', random_state=0)
clusters = hclust.fit_predict(df_scaled)

# Silüvyo skorunu hesaplayalım
silhouette_scores = silhouette_score(df_scaled, clusters)

# Silüvyo grafiğini çizelim
plt.figure(figsize=(10, 8))
sns.scatterplot(x=df_scaled[:, 0], y=df_scaled[:, 1], hue=clusters)
plt.title('Hiyerarşik Kümeleme Sonuçları')
plt.show()
```



Optimal kümeleme sayısı "elbow" yöntemiyle belirlenmiş ve bu sonuca göre bir model oluşturularak RFM skorları elde edilmiştir. Sonuçların yorumlanmasına göre, en düşük "recency" değerine sahip olan küme grubu sadık müşterilere aittir. Örneğin, 1. küme en düşük "recency" değerine sahip olup ortalama olarak "12" gün önce alışveriş yapmıştır. Ayrıca, bu küme 7 adet ürün satın almış ve bu ürünlerin toplam değeri 3624 sterlin olarak kaydedilmiştir.

SONUÇ VE ÖNERİ

Bu çalışma, RFM skorlarına dayalı olarak k-ortalama ve hiyerarşik kümeleme yöntemlerinin müşteri segmentasyonunda kullanılmasını araştırmaktadır. Analizler, RFM skorlarıyla birlikte bu iki yöntemin müşterilerin segmentlerine ayrılmasında etkili ve anlamlı sonuçlar sağladığını göstermektedir. K-ortalama algoritması benzer RFM skorlarına sahip müşterileri gruplandırmak için kullanılırken, hiyerarşik kümeleme yöntemi RFM skorlarına dayalı olarak müşterileri bir ağaç yapısıyla hiyerarşik olarak sınıflandırmak için kullanılmıştır. Sonuç olarak, RFM skorlarına dayalı olarak bu iki yöntemin birlikte kullanılması, müşteri segmentasyonunda önemli bir ilerleme kaydetmektedir ve şirketlere daha etkili pazarlama stratejileri oluşturma fırsatı sunmaktadır.



Özet

İnsan sayısı bilindiği üzere her geçen gün artmaktadır ve bu artışın oluşturduğu çeşitli maliyetler vardır. Bunlardan biri de trafiktir. Trafik akışının tıkanması veya tamamen durması, insanların kaybettiği vakit, artan karbon emisyonu ve kullanılan yakıt miktarında artış gibi birçok olumsuz etkisi vardır. Bu olumsuzlukları minimuma indirmekle uğraşan; karayolları genel müdürlükleri, belediyeler, ulaştırma mühendisleridir ve ihtiyaç uydukları en önemli verilerden birisi trafik akışı verileridir. Trafik akışı verileri birçok yöntem ile elde edilebilir ancak az maliyeti, kullanım esnekliği gibi faktörlerden dolayı kameralar en iyi çözümler arasındadır. Bu çalışmada kamera kayıtlarından trafik akışı verisi üretilmeye çalışılacaktır.

Trafik Akışı Verisi Nedir

Trafik akışı verisi trafikte bulunan araçların sayısının bazı kırımlara göre sayılması ve ayrıca araç hızlarını da içeren bir veridir. Bu kırımlar araçların tipleri; araba, kamyon, motosiklet vb. ve araçların yol üstündeki rotalarının bilgisidir. Aşağıda yapılan çalışmanın sonucundan bir örnektir. Çalışmada gidiş ve geliş olan bir yol üzerinden bu veri üretilmiştir.

Cls	Direction	Kmph
2	Left	140,10
2	Right	134,44
2	Left	113,54
7	Right	88,93
2	Right	134,76

- **Cls** : araçların tiplerini içermektedir 2 araba 7 kamyon olmak üzere.
- **Direction** : Araçların rota bilgilerini içerir. Kayıt alındığı konuma göre 'Left' gidiş 'Right' geliş tir.
- **Kmph**: Araçların hızlarını içerir.

Kullanılan Algoritma

Bu çalışmada yapay sinir ağı ile çalışan Yolov5_StrongSort_OsNet Algoritmasından faydalanılmıştır. Algoritma one-stage (tek adım) nesne tanıma algoritmasının olan Yolov5'in üzerine nesne takip sisteminin inşa edilmesi ile oluşturulmuştur. Tıpkı nesne takip de olduğu gibi nesne tanıma ve sınıflama kısımları daha önceden eğitilerek geliştirilmiş bir algoritmadır. Yolov5'e bakacak olursak CNN (Evrişimli Sinir Ağı) tabanlı, nesne tespiti yapan ve canlı görüntüleri işleyebilecek çalışma hızına sahip bir algoritmadır. Sonuç olarak eliminizdeki araç nesne tanıma, tahmin ve tespit konusunda başarılı olsa da trafik akışı verisi için çeşitli eklemeler yapılmalıdır.

Algoritmanın Özelleştirilmesi

1. Giriş ve Çıkış Tespiti

Öncelikle trafik görüntüsü kaydında araçların yol üzerinde giriş ve çıkış piksellerinin sayısal olarak tespit edilmesi gerekmektedir. Bu işlem görüntünün bir karesi üzerinden fotoğrafta gözüktüğü gibi gerçekleştirilmiştir.



Giriş ve çıkışlar harici algoritma çalıştırılarak canlı olarak video kaydında ekranda oynatılmaktadır. Bu pencereye yapılan hesaplamaların gözükmesi için görüntüye giren araçlar için sayıçlar eklenecektir. Bu sayıçların ve yazıların konularında hesaplanmıştır.

2. Sayım

Bu işlem için aşağıda gözüktüğü gibi if-else blokları oluşturulmuştur. Bu bloklar algoritmanın ürettiği nesne tespiti ve takibi sonucunda elde edilen aracın merkez piksel koordinatlarının daha önceden belirlenen koordinatlar içerisinde olması halinde çeşitli ham verilerin tutulmasını sağlayan kodları tetiklenmesini ve araç ile ilgili ham verilerin tutulmasını sağlar.

```
def count_obj(box_w, h, id, cls):
    global ccl, cti, ccr, ctr

    if id not in control:
        if int(box[1] + int(box[3]-box[1])/2) > (h-370):
            if int(box[0] + (box[2]-box[0])/2) > 50 and int(box[0] + (box[2]-box[0])/2) < 615:

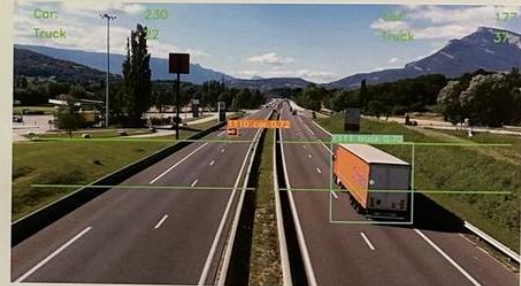
                control.append(id)
                dctst['TimeSt'].append(time.time())
                dctst['ID'].append(id)
                dctst['cls'].append(cls)
                dctst['Direction'].append('Left')

                if cls==2:
                    ccl+=1
                elif cls==7:
                    cti+=1
```

3. Verilerin İşlenmesi

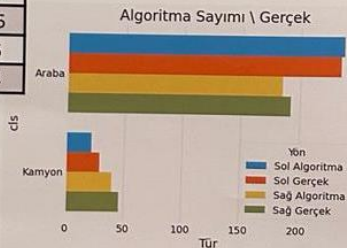
Tutulmuş ham verilerin içerisinde aracın rotası ve sınıfı istenilen halde gelmektedir ancak aracın hızı için ek hesaplamalara ihtiyaç vardır. Araçların belirlenen giriş ve çıkış alanlarına girdikten sonra bilgisayarın saat bilgisi çekilmekte ve aradaki farktan piksel saniye cinsinden hız hesaplanmaktadır. Bu hız verisi belirlenen bölgenin uzunluğu bilinmediğinden farazi olarak hesaplanmıştır.

Sonuçlar



Algoritma çalışırken oluşa çıktığı şekildedeki gibi gözükmektedir. Yolda bulunan araçlar sadece araba kamyonlardan oluşmaktadır. Kamyon olarak ayrılan araçlar ile tırlar 'truck' olarak tek bir seviyede değerlendirilmiştir. Ağıdaki görsellerde algoritma ile gözlem yolu sayım arasındaki fark gözükmektedir.

Araç Sınıfı	Yön	Algoritma Sayımı	Gerçek Sayım
Araba	Sağ	185	194
	Sol	238	235
Kamyon	Sağ	22	26
	Sol	40	44



KAYNAKÇA

- [1] Yolov5_StrongSort_OsNet https://github.com/hdnh2006/Yolov5_DeepSort_Pytorch
- [2] Yardım, M. S., & Akyıldız, G., (2005). Akıllı Ulaştırma Sistemleri ve Türkiye'deki Uygulamalar . 6. Ulaştırma Kongresi (pp.405-414). İstanbul, Turkey
- [3] <https://towardsdatascience.com/artificial-neural-networks-for-total-beginners-d8cd07abaae4>



ÖZET

Online alışveriş, internetin herhangi bir aracı kullanılarak gerçekleştirilen alışveriş şeklidir. E-ticaret ya da çevrimiçi alışveriş gibi isimlerle de anılan bu pazarlama şekli, Türkiye'de özellikle son 5 yılda popüler bir hale gelmiştir. Online alışveriş, reklamlıcılığın, pazarlamanın, alım-satım ve ödeme işlemlerinin tamamının yapıldığı platformları kapsadığından, klasik pazarlamadaki kavramlara ek olarak pek çok kavramı beraberinde getirmiştir. Bu çalışmada, online alışveriş ile ilgili sorular içeren 20 soruluk bir anket hazırlanmıştır. Ankete katılan 123 kişinin cevapları, SPSS programı kullanılarak analiz edilmiştir. Sorulara verilen cevaplara, çapraz tablolar yardımıyla bir ön analiz yapılmıştır. Sonrasında ki-kare testi kullanılarak seçilen sorular arasında bir ilişkinin olup olmadığı test edilmiştir. Analiz sonuçları elde edilen SPSS çıktıları üzerinden yorumlanmıştır.

ONLINE ALIŞVERİŞİN TANIMI ve ÖZELLİKLERİ

Online alışverişin basit bir tabirle "internet üzerinden satın alma işlemi" olarak açıklayabiliriz. Günümüzde internet kullanımının giderek artması ile birlikte, online alışveriş de hızla yaygınlaşmıştır. İnternetin sunduğu kolaylık, zaman ve maliyet tasarrufu, tüketicilerin internet üzerinden alışveriş yapmasını teşvik etmektedir.

Online alışverişin özelliklerini; 7/24 açık olması, ürünlerin detaylı bir şekilde sunulabilmesi, ürünlerin fiyat ve özelliklerinin karşılaştırılabilmesi, ödeme kolaylığı ve ürünlerin ev veya iş yerine teslim edilebilmesi olarak açıklayabiliriz. Online alışverişte ödeme yöntemleri de oldukça çeşitlidir. Kredi kartı, banka havalesi, kapıda ödeme ve benzeri birçok ödeme yöntemi kullanılabilir. Bunlara ek olarak, online alışveriş siteleri, binlerce ürün seçeneği sunar ve tüketiciler, mağaza gezerek bu kadar geniş bir seçeneği bulamazlar. Tüm bunları online alışverişin avantajları olarak da gösterebiliriz.



ONLINE ALIŞVERİŞİN GELİŞİMİ ve TARİHİ



Online alışverişin tarihçesi, 1980'lerin sonlarına kadar uzanmaktadır. Bu dönemde, online alışverişin temelinde telefon, televizyon veya posta siparişi ile çalışan şirketler bulunmaktaydı. İnternetin yaygınlaşmasıyla birlikte online alışveriş de önemli bir değişim geçirdi. İlk e-ticaret web siteleri, 1990'ların ortalarında açılmaya başladı.

1994 yılında, bir bilgisayar bilimcisi olan Phil Brandenberger, Pizza Hut'un web sitesi üzerinden pizza siparişi veren ilk müşteri olarak kaydedildi. Bu olay, online alışverişin doğuşu olarak kabul edilmektedir. Aynı yıl içinde, Jeff Bezos, Amazon.com'u kurmuştur. Amazon, başlangıçta sadece kitap satan bir web sitesi olarak faaliyet gösterse de, daha sonra elektronik, ev gereçleri, giyim gibi birçok farklı kategoride ürün satışına başlamıştır. 2000'li yılların başlarına kadar, online alışverişin gelişimi yavaş bir şekilde ilerlemiştir. 2000'lerin ortalarında, çevrimiçi alışveriş yapmak için gerekli olan altyapı ve teknolojik gelişmeler hızlanmaya başlamıştır. Mobil cihazların yaygınlaşması, e-ticaret web sitelerinin daha kullanıcı dostu ve erişilebilir hale gelmesi gibi faktörler, online alışverişin hızlı bir şekilde yayılmasına yol açmıştır. Günümüzde de, online alışverişin gelişimi her geçen gün hız kazanmaktadır. Özellikle, COVID-19 pandemisi nedeniyle, online alışverişe olan talep daha da artmıştır. Tüketiciler, evden çıkmadan ürünleri sipariş ederek ihtiyaçlarını karşılamaktadırlar. Ayrıca, birçok perakende şirketi, online alışverişin yaygınlaşmasıyla birlikte dijital pazarlama ve reklamlık stratejilerine de ağırlık vermeye başlamıştır.

SOSYAL MEDYANIN ONLINE ALIŞVERİŞE ETKİSİ

Sosyal medya, online alışveriş deneyimini daha kişisel hale getirmiştir. Birçok online alışveriş sitesi, müşterilerin alışveriş geçmişlerine ve tercihlerine dayalı öneriler sunar. Sosyal medya ise bu önerilerin daha da kişiselleştirilmesine olanak sağlar. Örneğin, bir tüketicinin sosyal medya hesabında beğendiği veya takip ettiği bir marka, kendi ürünlerini sosyal medya üzerinden önererek tüketicinin alışveriş kararını etkileyebilir.



Sosyal medya, tüketicilerin satın alma kararlarını etkilemek için markaların kullandığı bir pazarlama aracıdır. Birçok marka, sosyal medyada hedef kitlelerine özel reklamlar yayımlayarak, tüketicilerin ilgisini çekmeye ve satın alma kararlarını etkilemeye çalışır. Sosyal medya platformları ayrıca influencer marketing adı verilen bir pazarlama stratejisi için de önemli bir rol oynamaktadır. İşletmeler, sosyal medya fenomenlerini kullanarak ürünlerini tanıtabilir ve onların takipçileri aracılığıyla daha da geniş bir kitleye ulaşabilirler. Bu yöntem, tüketicilerin ürünlerle daha olumlu bir şekilde yaklaşmalarına yardımcı olabilmektedir.

ANKET SONUÇLARI

Soru	İstatistik	N	%
Cinsiyetiniz nedir?	Erkek	49	39,8
	Kadın	74	64,2
Yaş aralığınız nedir?	18-25	80	65,0
	26-35	12	9,8
	36-50	11	8,9
	50 ve üstü	20	16,3
	Hiçbiri	1	0,8
Eğitim durumunuz nedir?	Ortaokul	3	2,4
	Lise	18	14,6
	Üniversite	93	75,6
	Lisansüstü	8	6,5
Çalışma durumunuz nedir?	Kamu Sektörü	9	7,3
	Özel Sektör	26	21,1
	Öğrenci	67	54,5
	Ev Hanımı	8	6,5
	Kendi işinin patronuyum	1	0,8
	Çalışmıyorum	12	9,8

Soru	İstatistik	N	%
Günlük internet kullanım süreniz nedir?	1-3 saat	12	9,8
	3-5 saat	34	27,6
	5-10 saat	40	32,5
	10 saat ve üstü	9	7,3
Sosyal medya hesabınız var mı?	Evet	120	97,6
	Hayır	3	2,4
Sosyal medya hesabınız var ise nerelerde kullanıyorsunuz?	Akış Telefon	113	91,8
	İkişer	1	0,8
	Diğer	1	0,8
	Sosyal medya hesabınız yok	3	2,4
Apađdaki sosyal medya platformlarından en çok hangisini kullanıyorsunuz?	Facebook	7	5,7
	Twitter	19	15,4
	Instagram	79	64,2
	YouTube	13	10,6
	LinkedIn	0	0,0
	Diğer	2	1,6
	Sosyal medya hesabınız yok	3	2,4
İnternette bir ürün/hizmet satın alırken sosyal medyadan etkilenmişleriniz olduğunu düşünür müsünüz?	Evet	93	75,6
	Hayır	27	22,0
	Sosyal medya hesabınız yok	3	2,4
Sosyal medya platformlarından yapılan tanımlan kampanyaların, klasik medya araçlarından (TV ve gazete gibi) daha etkili buluyor musunuz?	Evet	101	82,1
	Hayır	19	15,4
	Sosyal medya hesabınız yok	3	2,4

Soru	İstatistik	N	%
İnternet üzerinden alışveriş yaptığınız sıklığına ne derseniz?	Evet	123	100,0
En son ne zaman internet üzerinden alışveriş yaptınız?	Son 1 hafta içerisinde	67	54,5
	Son 1 ay içerisinde	37	30,1
	Son 3 ay içerisinde	18	14,6
	Son 1 yıl içerisinde	5	4,1
	İnternet üzerinden alışveriş yapmıyorum	0	0,0
Ne sıklıkla internet üzerinden alışveriş yapıyorsunuz?	Hiçbiri	0	0,0
	Hiçbiri bir	13	10,6
	Haftada bir	14	11,4
	Ayda bir	47	38,2
	Ayda birden fazla	32	26,0
	Yıla bir	17	13,8
	İnternet üzerinden alışveriş yapmıyorum	0	0,0
İnternette gördüğünüz ürün reklamlarını tıklayarak satın alma gerçekleşti mi?	Evet	47	38,2
	Hayır	76	61,8
	İnternet üzerinden alışveriş yapmıyorum	0	0,0
İnternette gördüğünüz bir ürünün indirime girmiş olmasının, o ürünü almanıza etkili miydi?	Evet	108	87,8
	Hayır	13	10,6
	İnternet üzerinden alışveriş yapmıyorum	0	0,0
İnternet üzerinden alışveriş hangi sektörler için kullanıyorsunuz? (Birden fazla seçmek için lütfen her satıra bir kutucuğu işaretleyiniz)	Elektronik	70	56,9
	Giyim Sektörü	100	81,3
	Eğlence Sektörü	40	32,5
	Elektronik Sektör	61	49,6
	Kozmetik Sektörü	52	42,3
	Diğer	26	21,1
	İnternet üzerinden alışveriş yapmıyorum	0	0,0
İnternet üzerinden alışveriş yaparken dikkatli hangisine dikkat edersiniz?	Ürünlerin Fiyatı	70	56,9
	Ürünlerin Kalitesi	99	80,5
	Ticariletiler	44	35,8
	Teslimat Süresi	0	0,0
	İnternet üzerinden alışveriş yapmıyorum	0	0,0
Apađdaki alışveriş sitelerinden hangisini veya hangilerini tercih ediyorsunuz? (Birden fazla seçmek için lütfen her satıra bir kutucuğu işaretleyiniz)	Trendyol	106	86,2
	Hepsiburada	66	53,2
	Getir	13	10,6
	Amazon	54	43,9
	YeniDünya	53	43,1
	NI1	15	12,2
	Çoksepetim	27	22,0
	Diğer	27	22,0
İnternet üzerinden alışveriş yaptığınız bölgede	İnternet üzerinden alışveriş yapmıyorum	0	0,0
	Reklam Bankası Kurum	115	93,5
	Kapıda Ödeme	8	6,5
İnternet üzerinden alışveriş yaptığınız bölgede	İnternet üzerinden alışveriş yapmıyorum	0	0,0
	Sosyal Medya	30	24,4
	Mobil Uygulamalar	88	71,5
	Diğer	5	4,1

HİPOTEZLER ve Kİ-KARE TESTİ YORUMU

Chi-Square Tests	Value	df	Asymptotic Significance (2-sided)
Pearson Chi-Square	20,141*	16	.214
Likelihood Ratio	17,581	16	.349
Linear-by-Linear Association	1,099	1	.294
N of Valid Cases	123		

a. 19 cells (76.0%) have expected count less than 5. The minimum expected count is .11.

Chi-Square Tests	Value	df	Asymptotic Significance (2-sided)
Pearson Chi-Square	13,899*	6	.031
Likelihood Ratio	13,627	6	.034
Linear-by-Linear Association	4,266	1	.039
N of Valid Cases	123		

a. 7 cells (58.3%) have expected count less than 5. The minimum expected count is .45.

Ki-Kare İlişki Hipotezleri

H₀: Sorular arasında herhangi bir ilişki yoktur.

H₁: Sorular arasında bir ilişki vardır.

İLİŞKİ BULUNAN VE İLİŞKİ BULUNMAYAN SORULAR

Ki-kare analizine göre aralarında ilişki bulunan sorular şunlardır;	Ki-kare analizine göre aralarında ilişki bulunmayan sorular şunlardır;
<ul style="list-style-type: none"> Cinsiyet ile Sosyal Medyanın Alışverişe Etkisi Arasındaki İlişki Yaş ile Günlük İnternet Kullanım Sıklığı Arasındaki İlişki Cinsiyet ile Sosyal Medyadaki Reklamların Etkisi Arasındaki İlişki Cinsiyet ile İnternette Yapılan Son Alışverişin Yapılma Zamanı Arasındaki İlişki Yaş ile Ürün Satışı için Kullanılan Platformlar Arasındaki İlişki Sosyal Medya Kullanımı ile İnternette Yapılan Son Alışverişin Yapılma Zamanı Arasındaki İlişki Sosyal Medyanın Alışverişe Etkisi ile İndirimlerin Ürünü Satın Almadaki Etkisi Arasındaki İlişki 	<ul style="list-style-type: none"> Yaş ile Sosyal Medyanın Alışverişe Etkisi Arasındaki İlişki Çalışma Durumu ile Sosyal Medya Kullanımı Arasındaki İlişki Çalışma Durumu ile En Çok Kullanılan Sosyal Medya Platformu Arasındaki İlişki Cinsiyet ile Tercih Edilen Ödeme Yöntemleri Arasındaki İlişki Eğitim Durumu ile İnternet Üzerinden Yapılan Alışverişin Sıklığı Arasındaki İlişki Sosyal Medyanın Alışverişe Etkisi ile İnternet Sitelerindeki Reklamların Alışverişe Etkisi Arasındaki İlişki İnternette Yapılan Son Alışverişin Yapılma Zamanı ile İnternette Alışveriş Yaparken Dikkat Edilen Öncelikler Arasındaki İlişki

KAYNAKÇA

- [1] Jarvenpaa, S. L., & Todd, P. A. (1996). Consumer reactions to electronic shopping on the World Wide Web. International Journal of electronic commerce, 1(2), 59-88.
- [2] Park, C. H., & Kim, Y. G. (2003). Identifying key factors affecting consumer purchase behavior in an online shopping context. International journal of retail & distribution management.
- [3] TOPAL, I., & TEMİZKAN, V. (2016). Tüketicilerin mobil sosyal medya kullanımının marka



Dünyada birçok ülkenin kamu harcamalarında önemli bir yere sahip olan askeri harcamalar ulusal güvenliği sağlama açısından ülkelerin önemli harcama kalemleri arasında yer almaktadır. Bu doğrultuda yapılan bu çalışmanın amacı savunma harcamaları bakımından önemli büyüklüğe sahip olan ülkelerin gayri safi yurt içi hasılları ile askeri harcamalar değişkenlerinin hangi modellere uygun olduğunu ve değişkenler arasındaki ilişkiyi tespit etmektir. Çalışmada panel veri analizi yapılmasına imkan sağlayan Stata programı kullanılarak yapılan analizlerde öncelikle veri seti için uygun olan model tespit edilmeye çalışılmış ve Rassel Etkiler Modelinin uygun olduğu sonucuna ulaşılmıştır. Yatay kesit bağımlılığını araştırmak için Pesaran CD testine başvurulmuş ve test sonucunda yatay kesit bağımlılığının olduğu tespit edilmiş olup İkinci Nesil Panel Birim Kök Testi olan Pesaran CADF Testi uygulanmıştır. Sonuç olarak, serilerin durağan olduğu ve seçilen ülkelerde gayri safi yurt içi hasıla ile askeri harcamalar arasında pozitif bir nedensellik ilişkisinin olduğu tespit edilmiştir.

PANEL VERİ VE PANEL VERİ ANALİZİ

- Ekometrik analizlerde herhangi bir konuda hem zamana göre hem de birimlere göre analiz yapılması gerektiğinde, genellikle bu analizler zamana ve birime göre ayrı ayrı yapılmaktadır. Zamana göre yapılan analizler zaman serileri analizi olmakta, birimlere göre yapılan analizler ise yatay kesit analizi olmaktadır. Zaman serileri ve yatay kesit analizinin birleştirilmesini ve uygun modellerin test edilmesini sağlayan yöntemlere panel veri analizi denilmektedir.
- Panel veriler kullanılarak oluşturulan panel veri modelleri yardımıyla ekonomik ilişkilerin tahmin edilmesi yöntemine "Panel Veri Analizi" denilmektedir. Bu analizde genelde, yatay kesit biriminin (N) sayısının dönem sayısından fazla olduğu durumlarda çalışılmaktadır.

PANEL VERİNİN ÖZELLİKLERİ

- Herhangi bir yatay kesitte araştırma konusu olan birimlerin (firmalar, ülkeler vb.) davranışlarını etkileyen sayısız ölçülemez açıklayıcı değişken vardır. Bu değişkenlerin dışlanması sapmalı tahminlere neden olmaktadır. Benzer bir durum mikro birimlerin davranışlarını hep aynı yönde ancak her bir zaman döneminde farklı bir şekilde etkileyen zaman serisi değişkenlerinin dışlanması halinde de geçerlidir. Panel veri bu problem giderilmesine olanak tanınmaktadır.
- Panel veri bir dönemden diğerine meydana gelen değişim ile mikro birimler arasındaki farklılığı birleştirilerek suretiyle değişkenlik meydana getiren çoklu doğrusallığı azaltmaktadır.

PANEL VERİ ANALİZİNİN AVANTAJLARI

- Zaman serisi ve yatay kesit analizi ile kıyaslandığında panel veri analizi, araştırmacıya daha geniş bir veri seti ile çalışma imkânı sunar.
- Yeterli bir zaman uzunluğunda, değişim dinamiklerinin çalışmasında panel veri analizi yatay kesit ve zaman serisi analizlerine göre daha avantajlı bir yöntemdir.

PANEL VERİ ANALİZİNİN DEZAVANTAJLARI

- Her birim için zaman serisi boyutunun kısa olabilmesi.
- Yatay kesit ve zaman serisi gözlemleri arasında meydana gelen parametre farklılıklarının göz önüne alınmadığı durumlarda birtakım sapmaların ortaya çıkması ve bu durumun parametrelerin tutarsız ve anlamlı olmayan tahminlerine sebep olması.

PANEL VERİ ANALİZİ MODEL TÜRLERİ

- Sabit Etkiler Modeli:** Sabit etki modeli oluşturulan katsayılar, birimlere ve/veya zamana göre farklılık göstermektedirler. Bu modelin temelinde, birimler arası farklılıkların modelde yer alan sabit terimdeki farklılıklar aracılığıyla yok edilebileceğini düşüncesi yatmaktadır. Bundan dolayı her birimi temsil edecek farklı sabit terimler kullanılmaktadır. Bu modelde yatay kesit birimleri arasındaki farklılıklar sabit terimdeki farklılık açıkladığı ve kükü/gölgü değişken yardımıyla tahmin edildiği için model kükü değişkenli model olarak da bilinmektedir.
- Rassel Etkiler Modeli:** Sabit etki modelinde birimler arası farklılığın sabit olduğu ve bunun sabit terimdeki farklılıklarla giderileceği varsayılmaktadır. Ancak bazı çalışmalarda birimlerin tesadüfi olarak seçildiği durumlar olabilmektedir. Bu birimler arası farklılıklar da tesadüfi olduğu için rassel (tesadüfi) etkiler modeli sabit etki modeline alternatif olarak geliştirilmiş bir yöntemdir. Rassel etkiler modelinin avantajı, zaman değişkenlerinin modele ilave edilebilmesidir. Rassel etkiler modelinde birim ve zaman etkileri tesadüfi değişken olarak modelde hata teriminin bileşeni olarak analize dahil edilmektedir.

VERİ TİPLERİ

Stata'da veri tipi çok önemlidir. Çalıştığımız veri sürekli bir veri ise "float" olarak kaydedilmesi gerekir. Stata'da tam sayılar için sadece "byte", "int" ve "long" kullanılır. "float", yedi basamaklı sayı olarak depolanmaktadır. Sayıların büyüklüğü önemli değildir. 9 basamaklı sayılar veya daha azı için "long", 9 basamaktan daha fazlası ise "double" olarak depolanabilmektedir.

ANALİZLER

Oluşturduğumuz veri setinde başlangıçta toplam 1450 adet veri bulunmaktaydı. Zaman serisi bakımından 1997-2021 arası dönem ele alınmış olup, yatay kesit serisi bakımından Avrupa ve Orta Asya ülkelerindeki 49 ülke ele alınmıştır. Eksik veriler içerisindeki dolay veri setinden 2020-2021 yıllarını ve çeşitli ülkeleri çıkartarak veri setini dengeli (balanced) hale getirdiğimizde veri setinde toplam 966 adet veri bulunmaktaydı. Veri setindeki bağımlı değişken Askeri Harcamalar (Military Expenditures (mex)) (current USD), bağımsız değişken Gayri Safi Yurt İçi Hasıla (gross domestic product (gdp)) (current US\$) olmaktadır.

- Veri setini panel veri olarak tanımladıktan sonra sağ tarafta bulunan cıktıda yer alan strongly balanced ifadesi panel veri analizine başlayabilmemize olanak sağlar.

```

statst countryid time
-----
panel variable: countryid (strongly balanced)
time variable: time, 1997 to 2019
data in units
-----

```

- Sağ tarafta yer alan cıktıda veri setinin özeti hakkında olacak olursak veri setindeki değişkenleri, toplam gözlem sayısını, değişkenlerimizin ortalaması, standart sapması ve minimum-maximum değerlerini görebiliriz.

```

summarize
-----
Variable      Obs      Mean      Std. Dev.      Min.      Max.
-----
countryid     966      2008      0.000000      1997      2019
time          966      4.74111     7.484111     1.744000     9.744000
gdp           966      7.808000    1.484000    5.007700     9.804000
mex           966      22.8       12.1000     1.000000     41.000000
-----

```

- Sağ tarafta yer alan cıktıda ise veri setinde bulunan değişkenlerin detaylı özeti ve her bir değişken için toplam gözlem sayısını, yatay kesit birim sayısı ile zaman birimi sayısını görebiliriz.

```

estat sby num
-----
Variable      Mean      Std. Dev.      Min.      Max.      Observations
-----
gdp            7.808000    1.484000    5.007700     9.804000     966
mex            22.800000  12.100000    1.000000     41.000000    966
-----

```

KAYNAKÇA

- Dr. Öğr. Üyesi Elif Tuna Uygulamalı Panel Veri Analizi Ders notları
- Doc. Dr. Ferda Yelen Tatoglu Panel Veri Ekometrisi
- Hsiao, C., 2003, Analysis of Panel Data, Cambridge University Press, United Kingdom.
- Baltagi, B. H., 2005, Econometric Analysis of Panel Data, Third Edition, John Wiley & Sons Inc, England.

- Klasik Model ile Rassel Etkiler Modeli arasında tercih yapmak için LM testini kullanırız. Test hipotezleri; H0: Tesadüfi etkiler sıfırdır. H1: Tesadüfi etkiler sıfır değildir, şeklindedir. Yan tarafta bulunan cıktıda yer alan Prob>chibar2=0.0000 ifadesi neticesinde H0 hipotezini reddetmekle herhangi bir hata yapmayacağımızı söyleriz ve bu nedenle H0 hipotezi reddedilir. Dolayısıyla tesadüfi etkiler sıfır değildir, rassel etkiler modeli uygundur sonucuna ulaşılır.

```

estatlm
-----
Residual and Panel (strongly balanced) test for random effects
Test of H0: rho = 0 (no random effects)
Test of H1: rho > 0 (random effects)
-----
Number of obs = 966
Number of panels = 49
Number of time periods = 23
-----
F(1, 916) = 12.4614
Prob > F = 0.0000
-----

```

- Rassel Etkiler Modeli ile Sabit Etkili Model arasında seçim yapmak için Hausman Testi uygulayacağız. Test hipotezleri; H0: Tesadüfi etki modeli uygundur. H1: Tesadüfi etki modeli uygun değildir, şeklindedir. Yan tarafta bulunan cıktıda yer alan Prob>chibar2=0.3618 ifadesi neticesinde %95 güven düzeyinde H0 hipotezini reddetmekle hata yaparız ve H0 hipotezi reddedilmez. Dolayısıyla tesadüfi etkiler sıfır değildir, rassel etkiler modeli uygundur sonucuna ulaşılır.

```

hausman
-----
Test of H0: rho = 0 (no random effects)
Test of H1: rho > 0 (random effects)
-----
Number of obs = 966
Number of panels = 49
Number of time periods = 23
-----
F(1, 916) = 0.3618
Prob > F = 0.5482
-----

```

- Yapılan iki test neticesinde veri setimize uygun olan modelin Rassel Etkiler Modeli olduğu sonucuna ulaşılmıştır.

RASSAL ETKİLER MODELİ

Rassel Etkiler Modelini tahmin etmek için çeşitli kodlar çalıştırarak sağ tarafta bulunan cıktıyı elde ederiz. Cıktıdan da görüleceği üzere (corr(L, x)=0 ifadesi) Rassel Etkiler Modeli hata terimleri ile açıklayıcı değişkenler arasındaki korelasyonun 0 olduğunu varsayar. Modelin overall R-square değeri 0.82'dir. Cıktıda gdp değişkeninin katsayısının anlamlı olup sabit katsayısının anlamsız olduğu görülmektedir. rho değeri rassel etkiler modelinin varyansının %27.8'inin birimler arasındaki farklılıktan kaynaklandığını söylemektedir. Prob>chi2=0.0000 ifadesi modelin tüm katsayılarını anlamlı olduğunu göstermektedir.

```

xtreg, re
-----
Number of obs = 966
Number of panels = 49
Number of time periods = 23
-----
F(1, 916) = 12.4614
Prob > F = 0.0000
-----

```

Ş ZAMANLI KORELASYON

Ş zamanlı korelasyon, panel veri modelinde yatay kesit birimleri arasında bağımlılık olmadığını ifade eder. Eğer bu bağımlılık mevcut ise analizlerimiz hatalı sonuçlar verecektir. Breusch ve Pagan LM Testi ve Pesaran CD Testi ile test edilir.

PESARAN CD TESTİ

- Pesaran CD Testi kısa panel verilerde (N>T ise) kullanılır. Dolayısıyla bizim panelimiz için uygundur. Test hipotezleri; H0: Yatay kesit bağımlılığı yoktur. H1: Yatay kesit bağımlılığı vardır, şeklindedir. Cıktıda yer alan Pr = 0.0000 ifadesi, H0 hipotezini reddetmekle hata yapmayacağımızı anlamına gelmektedir. Dolayısıyla H0 hipotezi reddedilir, değişkenlerimiz arasında yatay kesit bağımlılığı vardır.

```

xtcd, pesaran abs
-----
Pesaran's test of cross sectional independence = 41.438, Pr = 0.0000
Average absolute value of the off-diagonal elements = 0.574
-----

```

- Pesaran CD Testi sonucunda modelimizde bulunan değişkenler arasında yatay kesit bağımlılığı olduğu sonucuna ulaşılmış olup yatay kesit bağımlılığını konu alan İkinci Nesil Panel Birim Kök Testi olan Pesaran CADF testine başvurulmuştur.

PESARAN CADF TESTİ

- Pesaran CADF Testi'ni gdp değişkenine uygularken test hipotezleri; H0: İngdp değişkeni durağan değildir. H1: İngdp değişkeni durağandır, şeklindedir. Yan tarafta bulunan cıktıda yer alan P-value=0.0000 ifadesi neticesinde H0 hipotezi reddedilir ve İngdp değişkeni durağandır sonucuna ulaşılır.

```

psarandf, ingdp
-----
Pesaran's CADF Test: Engle
Cross-sectional average of fixed period detrended and extreme t-ratios truncated
autoregressive residuals
-----
Number of obs = 966
Number of panels = 49
Number of time periods = 23
-----
F(1, 916) = 12.4614
Prob > F = 0.0000
-----

```

- Pesaran CADF Testi'ni mex değişkenine uygularken test hipotezleri; H0: İngdp değişkeni durağan değildir. H1: İngdp değişkeni durağandır, şeklindedir. Yan tarafta bulunan cıktıda yer alan P-value=0.0000 ifadesi neticesinde H0 hipotezi reddedilir ve İnmex değişkeni durağandır sonucuna ulaşılır.

```

psarandf, mex
-----
Pesaran's CADF Test: Engle
Cross-sectional average of fixed period detrended and extreme t-ratios truncated
autoregressive residuals
-----
Number of obs = 966
Number of panels = 49
Number of time periods = 23
-----
F(1, 916) = 12.4614
Prob > F = 0.0000
-----

```

SONUÇLAR

- Çalışma süresince yapılan analizler neticesinde panel veri setimize için en uygun modelin Rassel (Tesadüfi) Etkiler Modeli olduğu sonucuna ulaşılmıştır.
- Gayri safi yurt içi hasıla ile askeri harcamalar arasında pozitif bir nedensellik ilişkisi olduğu sonucuna ulaşılmıştır.
- Pesaran CD Testi'nin sonucuna göre yatay kesit bağımlılığı vardır.
- Değişkenlerimizde yatay kesit bağımlılığı olduğundan dolayı İkinci Nesil Panel Birim Kök Testleri uygulanmıştır. İkinci Nesil Panel Birim Kök Testlerini uygulayarak yapılan analizler sonucunda değişkenlerimizin durağan olduğu sonucuna ulaşılmıştır.



YTU FEN EDEBİYAT FAKULTESİ İSTATİSTİK BÖLÜMÜ

DERİN ÖĞRENME İLE FİLM YORUMLARINDAN DUYGU ANALİZİ

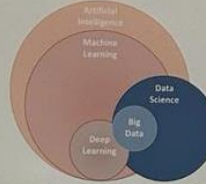
Alpay TUNCER 18023040

Danışman: Doç. Dr. Gülder KEMALBAY

Bilgisayar teknolojisindeki değişiklikler hayatın her alanında yenilikleri getirmiştir. Verilerin farklı tiplerinin derinlenip saklanıp işlenmesi ve sosyal medyanın kullanımı geleneksel medya araçlarına bir alternatif oluşturur. Metin madenciliği, yapılandırılmamış metinleri işlemek için kullanılır. Duygu analizi, metindeki duyguları belirlemeyi amaçlar. Filmler, herkesin ilgisini çekebilen ve çeşitli amaçlar için kullanılabilen eğlence araçlarıdır. Filmlerdeki sonuçlar, duygu analizinin film yorumlarını otomatik olarak analiz etmek, yorumlar ile verilen puanların arasındaki bağlantıyı anlamak ve genel izleyici duyarlılığını belirlemek için değerli bir araç olduğunu göstermektedir. Bu, sinema endüstrisinin film ticari başarısını önceden tahmin etmek ve izleyici eğilimlerini daha iyi anlamak için bu tür analizleri nasıl kullanabileceğine dair potansiyel yolları tartışır.

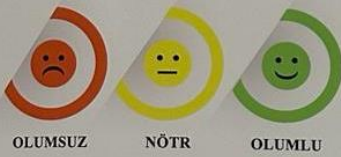
Metin madenciliği, doğal dil işleme (NLP) ve makine öğrenmesi uygulamaları için büyük veri setlerinden değerli bilgiler çıkarılması sürecidir. Bu, metinlerde belirli bilgi parçalarını, örüntüleri, duygusal durumları veya diğer önemli verileri tanımlamak ve analiz etmek için kullanılan bir tekniktir.

Doğal dil işleme (NLP), insanların doğal dillerini anlama, yorumlama, oluşturma ve manipüle etme yeteneğine sahip bilgisayar sistemlerinin geliştirilmesiyle ilgilene bilgisayar bilimi dalıdır. Bilgisayarların metin veya sesli dil verilerini anlamalarını, yanıtlarını, etkileşimlerini otomatikleştirme ve insan dilinin karmaşıklığını ve anlamlarını modelleme amacını taşır.



Makine öğrenmesi, bilgisayarların belirli bir görevi yerine getirirken performanslarını geliştirmek için veriden bağımsız olarak öğrenme yeteneklerine sahip olmalarını sağlayan bir yapay zeka dalıdır. Bu süreç, deneyim ve örüntülerin analizine dayalı algoritma geliştirme ile otomatik model oluşturma ve sürekli iyileştirme sağlar.

Derin öğrenme, genellikle çok katmanlı yapay sinir ağıları kullanarak karmaşık örüntülerin verilerden öğrenilmesiyle ilgilene bir makine öğrenmesi alt dalıdır. Bu, görüntü ve ses tanıma, metin anlama, çeviri, oyun oynama ve diğer karmaşık görevlerde üstün performans sağlama yeteneğine sahip modellerin oluşturmaktadır.



Duygu analizi, metinlerdeki duygusal tonları, tutumları ve yargıları belirlemek için kullanılan bir yapay zeka tekniğidir. Yani metinden insan duygusunu çıkarır. İngilizce'deki "Sentiment Analysis" teriminin Türkçe karşılığıdır. Genellikle doğal dil işleme, metin analizi, hesaplama dilbilim ve biyometrik ölçümler gibi teknikleri içerir.

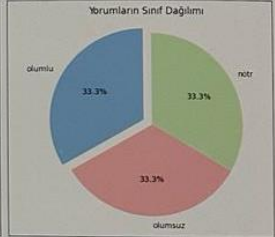
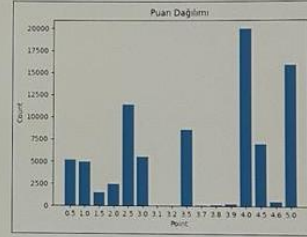
Duygu analizi, metinlerin içerisinde bulunan subjektif bilgileri ve bu bilgilerin duygusal yükünü sınıflandırmak ve anlamak için kullanılan bir yapay zeka tekniğidir. Bu teknik, doğal dil işleme (NLP), metin analizi ve hesaplama dilbilimini kullanarak metinlerde bulunan pozitif, negatif veya nötr duyguları belirler. Duygu analizi genellikle üç türde yapılır: polarite tabanlı (olumlu, olumsuz, nötr), duygu tabanlı (mutluluk, kızgınlık, üzüntü vb.) ve konu tabanlı (belli bir konu veya konular hakkında insanların duygularını analiz etme).

Duygu analizinin ham maddesi kullanılan dildir. Bu yüzden araştırmacının analizi yaptığı dile hakim olması analizinin başarısı için son derece önemlidir. Literatürde özellikle İngilizce çalışmaların fazlalığı göze çarpmaktadır, bundaki temel neden dilin yapısının basit olmasıdır. Türkçe gibi sondan eklemeli ve karmaşık yapıya sahip dillerde ise duygu analizi yapmak çok daha zordur. Türkçede bir kök; aldığı yapım ve çekim ekleri ile çok farklı anlamlara bürünebilmektedir.

Bir film, seri halinde gösterilen hareketli görüntülerden meydana gelen bir eserdir. Filmlerin oluşumu, gerçek insanların ve objelerin kamera ile çekimi veya animasyon, özel efektler gibi teknolojik yöntemler ile bu unsurların yaratılması yoluyla gerçekleşir.

Film eleştirisinin amacı filmleri çözümlenmek ve değerlendirmektir. Akademik film eleştirisi ve gazeteci film eleştirisi olmak üzere iki ana kategoriye ayrılır, genellikle bu iki türden birini izler. Fakat Bu çalışmada bireysel kullanıcıların IMDb internet sitesi üzerinden yaptığı tarafsız yorumlar kullanılmıştır.

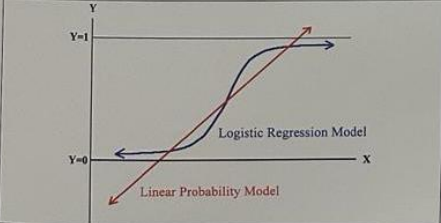
IMDb (Internet Movie Database), filmler, televizyon dizileri, TV programları, video oyunları ve internet tabanlı içerikler hakkında bilgi sunan online bir platformdur. Bu platform, oyuncular, yapımcılar, ekip, biyografiler, özetler, trivia, derecelendirmeler ve yorumlar gibi sinema ve televizyon yapımlarıyla ilgili bilgileri içerir. Uygulamada kullanılan veri seti, hazır verisetleri ve örnek çalışmaların olduğu Kaggle sitesinden elde edilmiştir. Elde edilen veri seti "Turkish Movie Sentiment Analysis Dataset" ismi ile Kaggle sitesinden kolayca bulunabilir. Veriler ile ilgili görselleştirmeler aşağıdaki gibidir.



Lojistik regresyon, bir bağımlı değişkenin iki kategorik değeri (örneğin, evet/hayır, 1/0) arasındaki ilişkiyi analiz etmek için kullanılan bir istatistiksel modeldir.

Lojistik regresyon, bağımlı değişkenin olasılık dağılımını tahmin etmek için bağımsız değişkenlerin lineer kombinasyonunu kullanır.

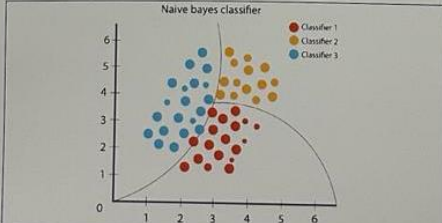
Sonuç olarak, lojistik regresyon, sınıflandırma problemlerinde kullanılan bir modeldir ve doğrusal olmayan ilişkileri yakalamak için lojistik fonksiyonu ile bağımlı değişkenin dönüşümünü kullanır.



Naive Bayes, makine öğrenmesinde sınıflandırma problemleri için kullanılan bir olasılık tabanlı bir algoritmadır.

Naive Bayes, Bayes teoremine dayanır ve sınıflandırma yaparken bağımsız değişkenler arasındaki ilişkisizlik varsayımını kullanır.

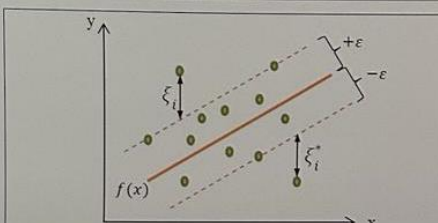
Bu algoritma, veri setindeki özelliklerin bir sınıfa ait olma olasılıklarını hesaplayarak yeni bir veri noktasını en olası sınıfa atar ve bu nedenle genellikle metin sınıflandırması gibi doğal dil işleme problemlerinde etkilidir.



Destek Vektör Makineleri (SVM), makine öğrenmesinde sınıflandırma ve regresyon problemleri için kullanılan bir algoritmadır.

SVM, veri noktalarını birer vektör olarak temsil eder ve bu vektörlerin sınıflarını birbirinden en iyi şekilde ayıran bir hiperdüzlemi bulmaya çalışır.

SVM, veri noktalarını sınıflandıran maksimum marjinal sınıflandırma prensibini kullanır ve aynı zamanda çekirdek fonksiyonları aracılığıyla doğrusal olmayan ilişkileri de yakalayabilir.

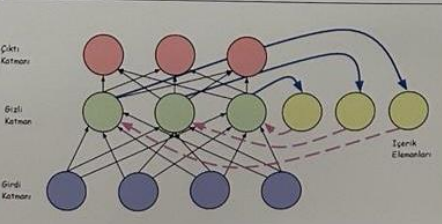


Yinelemeli Sinir Ağı (RNN), makine öğrenmesinde tekrar eden yapıların modellenmesi için kullanılan bir yapay sinir ağı türüdür.

RNN'ler, girdi verilerinin yani sıra önceki çıktılardan da geri besleme döngüsü aracılığıyla kullanarak, geçmiş bilgiyi hafızasında tutma yeteneğine sahiptir.

Bu özellik, zaman serileri, doğal dil işleme, konuşma tanıma gibi problemlerde zaman bağımlılıklarını ve uzun vadeli bağımlılıkları modellemek için etkili bir yöntem sunar.

RNN'lerin en yaygın varyasyonları LSTM (Uzun Kısa Vadeli Hafıza) ve GRU (Geriye Dönük Birim) olarak bilinir, bu varyasyonlar RNN'lerin uzun vadeli bağımlılıkları daha iyi yakalayabilmesini sağlamak için geliştirilmiştir.



Classification Report:				
	Precision	Recall	f1-Score	Support
Olumsuz	0.89	0.68	0.77	4602
Nötr	0.53	0.82	0.65	4253
Olumlu	0.71	0.75	0.72	4892
Accuracy			0.72	13747
Macro Avg	0.71	0.75	0.71	13747
Weighted Avg	0.78	0.72	0.73	13747

Bu çalışma, IMDb'deki film yorumlarından duygu analizi ile duygu tahminleme yapmayı amaçlamıştır. Uygulanan modeller arasında RNN modeli en yüksek (0.72) doğruluk oranını göstermiştir. Çalışmadaki ilk amacımız Türkçe dilini hazırlanmış Türkçe Lemmatizer ile çalışmak ve literatürdeki diğer çalışmalardan farklı olup olmadığını test etmektir, ayrıca bu veriye en uygun modeli tespit etmek, bu modelin performansını incelemek ve film yorumlarından duygu analizi yapabileceği kapasitesini ölçmektir. Bu şekil bir veri setine RNN modelinin uygun olabileceğini ve çalışmada kullanılan diğer modellerden daha başarılı olduğunu söyleyebiliriz. Olumlu ve olumsuz yorumların tahmini, nötr yorumlara göre daha yüksektir çünkü olumlu ve olumsuz yorumlarda belirli kelimeler öne çıkmaktadır. Veriler bireysel kullanıcılardan kontrolsüz alındığı ve verilen puanların yorumlarla doğru eşleşemeyeceğini göz önünde bulduğumuzda alınan doğruluk oranının yeterli olduğunu düşünebiliriz. Ancak, bu çalışmada kullanılan veri setinin, özellikle olumsuz yorumların daha az olduğu bir yapıya sahip olduğunu gözlemledik.

KAYNAKÇA

[1] Kızılkaya, Y. M. (2018). Duygu Analizi ve Sosyal Medya Alanında Uygulama (Doktora Tezi). Uludağ Üniversitesi. Sosyal Bilimler Enstitüsü, Bursa. <https://acikerisim.uludag.edu.tr/> (Temmuz, 2018).

[2] Elmas, Ş. Ş. (2019). Sosyal Medya Mesajlarının Veri Madenciliği Yöntemi ile Duygu Analizi (Sivas İli Örneği) (Yüksek Lisans Tezi). Sivas Cumhuriyet Üniversitesi. Sosyal Bilimler Enstitüsü, Sivas. <https://acikerisim.cumhuriyet.edu.tr/> (Ocak, 2019).

[3] <https://aws.amazon.com/tr/what-is/sentiment-analysis/>



Öznitelik Seçimi ve Boosting Algoritmalarını Kullanarak Makine Öğrenmesi Üzerinden Tahmin Modellemesi Geliştirme

Hazırlayan: 18023012 - Hürkan Üstündağ
Bitirme Çalışması Danışmanı: Doç.Dr.Öyküm Esra Yiğit
YILDIZ TEKNİK ÜNİVERSİTESİ İSTATİSTİK BÖLÜMÜ

Özet:

Günümüzde, veri biliminin ve yapay zeka teknolojilerinin gelişimi, özellikle makine öğrenmesi alanındaki yenilikler, daha önce çözülmemiş karmaşık sorunları çözmeye ve belirsizliği yönetme potansiyeline sahiptir. Boosting algoritmaları ve öznitelik seçimi teknikleri bu yenilikler arasındadır. Boosting algoritmaları sayesinde modelin performansı iyi derecelere yükseltilir ve öznitelik seçimi teknikleri sayesinde de model daha doğru ve etkili bir şekilde çalışabilir. Spaceship Titanic adlı veri seti bir yolcu uzay gemisidir ve amaç hangi yolcuların kaybolduğunu tahmin edebilmektir. Bu tahmin işlemi veri seti uygun bir şekilde temizlenerek bahsi geçen teknikler uygulanarak gerçekleştirilmiştir ve model performansının bu durumdan pozitif bir şekilde etkilendiği bu analiz sonucunda da yine görülmektedir.

Öznitelik Seçimi

Öznitelik seçimi, bir makine öğrenmesi modelinin eğitim sürecinde en önemli adımlardan biridir. Modelin veriler üzerinde nasıl eğileceğini ve hangi özelliklerin modelin tahminlerinde önemli bir rol oynayacağını belirlemek, genel model performansını ve öğrenme sürecinin etkinliğini önemli ölçüde etkiler. Modelin performansını birkaç farklı yolla artırabilir. İlk olarak, gereksiz veya ilgisiz özelliklerin kaldırılması, modelin eğitim süresini hızlandırır. Bu, özellikle büyük veri setlerinde ve yüksek boyutlu veri setlerinde önemlidir, çünkü bu tür veri setlerinde özniteliklerin sayısı genellikle verinin kendisinden daha fazla olabilir. İkincisi, öznitelik seçimi modelin anlaşılabilirliğini artırabilir. Az sayıda özellik, modelin çalışma mekanizmasını daha kolay anlamayı ve yorumlamayı sağlar. Son olarak, öznitelik seçimi, modelin genelleme yeteneğini artırabilir, yani modelin yeni, görülmemiş veriler üzerinde daha iyi performans göstermesini sağlar.

Öznitelik Seçimi Teknikleri

Filtre teknikleri

Bu yöntemler genellikle hızlı ve etkilidir çünkü makine öğrenmesi modelleri tarafından kullanılan özelliklerin bir ön değerlendirilmesini yaparlar. Bunlar, özelliklerin her birini bağımsız olarak değerlendiren istatistiksel tekniklerdir.

1) Bilgi Kazancı

Bir özelliğin hedef değişkenin belirsizliğini ne kadar azalttığını ölçer. Özellikle karar ağaçları gibi algoritmalarda kullanılır, çünkü ağaç dallarına noktaları belirlerken hangi özelliklerin en bilgilendirici olduğunu belirlemeye yardımcı olur.

2) Ki-Kare Testi

Kategorik özellikler için kullanılır ve hedef değişken ile olan bağımsızlığı test eder. Ki-Kare testi, bir özelliğin ve hedef değişkenin bağımsız olduğunu belirleme olasılığını ölçer, bu da özelliğin hedef sınıfı tahmin etmek için yararlı olabileceğini gösterir.

3) Fisher Skoru

Özelliklerin ayrılabilişliliğini ölçer. Fisher skoru yüksek olan özellikler hedef değişkeni daha iyi ayırabilir. Fisher Skoru, bir özelliğin hedef sınıflar arasında ne kadar iyi ayrım yaptığını ölçer, bu da özelliğin hedef sınıfı tahmin etmek için ne kadar yararlı olabileceğini gösterir.

4) Korelasyon Katsayısı

Bu, özelliklerin hedef değişken ile olan doğrusal korelasyonunu ölçer. Yüksek bir korelasyon katsayısı, bir özelliğin hedef değişken ile güçlü bir ilişkisi olduğunu ve dolayısıyla tahmin için yararlı olabileceğini gösterir.

5) Varyans Eşiği

Bu yöntem, düşük varyanslı özellikleri (yani verileri az değişen özellikleri) çıkarır. Çünkü düşük varyans genellikle az bilgi anlamına gelir. Bu ölçüm, bir özelliğin değerlerinin ortalama değerinden ne kadar farklı olduğunu ölçer.

Sarmal Metodlar

Bu metodlar, özellik seçimini bir arama problemi olarak ele alır. Özelliklerin farklı kombinasyonları belirlenir ve her biri için ayrı değerlendirilir. Kapak metodları genellikle zaman alıcıdır çünkü birçok farklı özellik kombinasyonunu değerlendirmek gerekmektedir. Ancak, bunlar genellikle daha doğru sonuçlar verir çünkü özelliklerin seçimi belirli bir modelin performansına dayalıdır.

1) İleri Doğru Seçim Tekniği

Bu yöntem, boş bir özellik seti ile başlar ve model performansını en çok artıran özellikleri tek tek ekler.

2) Geriye Doğru Seçim Tekniği

Bu yöntem, tüm özellikler ile başlar ve model performansını en az etkileyen özellikleri tek tek çıkarır.

3) Geri Arama Tekniği

Bu yöntem, modelin önemli özellikler olan özellikleri teker teker çıkarır ve en iyi altkümenin hangi özelliklerden oluştuğunu belirler.

Kaynakça

[1] (Liu, H., Motoda, H. (eds) Feature Selection for Knowledge Discovery and Data Mining, 1998)

[2] <https://www.kaggle.com/code/michalbrezk/xgboost-classifier-and-hyperparameter-tuning-85>

[3] <https://www.analyticsvidhya.com/blog/2020/02/4-boosting-algorithms-machine-learning/>

Gömülü Metodlar

Bu metodlar, modelin eğitimi sırasında özellik seçimini gerçekleştirir. Modelin karmaşıklığına bir ceza uygulayarak modelin aşırı uyumunu kontrol ederler.

1) L1 Düzenleme (Lasso)

Regresyon katsayılarını sınırlayan ve bu sınırlama sonucu bazı özelliklerin katsayılarını tamamen sıfıra çeken bir yöntemdir. Bu nedenle Lasso, modelin karmaşıklığını kontrol ederken aynı zamanda otomatik bir özellik seçimi gerçekleştirir. Bu yönüyle, Lasso, özellik seçimini modelin eğitimi sırasında gömülü bir şekilde gerçekleştiren bir özellik seçim tekniğidir.

2) L2 Düzenleme (Ridge)

Ridge regresyonu, regresyon katsayılarını sınırlar, ancak onları tamamen sıfıra çekmez. Bu, Ridge regresyonunun aşırı uyumu önlemesine yardımcı olur ve modelin genelleme kabiliyetini artırır. Ancak, Ridge, katsayıları tamamen sıfıra çekmediği için, tüm özellikleri modele dahil eder ve bu nedenle otomatik bir özellik seçimi sağlamaz. Ancak, özelliklerin katsayılarının büyüklüğüne bir sınırlama getirerek, önemli özelliklerin modelde daha belirgin hale gelmesine yardımcı olabilir.

3) Rassel Orman

Her özelliğin sınıflandırmaya performansına katkısını ölçer ve önemsiz özellikleri belirler

Boosting Algoritmaları

Boosting, makine öğrenmesinde bir topluluk öğrenme tekniği olup genellikle sınıflandırma ve regresyon problemlerini çözmeye amacıyla kullanılır. Boosting, birden çok zayıf öğrenici (genellikle karar ağaçları) kullanarak bir güçlü öğrenici oluşturmayı amaçlar. Her bir öğrenici, önceki öğrenicinin hatalarını düzeltmeye odaklanır. Boosting yöntemleri genellikle genel bir modelin performansını artırmak için kullanılır ve aşırı öğrenmeyi engeller yeteneği ile bilinir.

Boosting Algoritma Çeşitleri

AdaBoost (Adaptive Boosting)

AdaBoost, en popüler boosting algoritmalarından biridir. Bu yöntem, iterasyonları boyunca zayıf sınıflandırıcıları birleştirir ve her bir sınıflandırıcının ağırlığını hatalarına bağlı olarak ayarlar. AdaBoost, genellikle aşırı öğrenme eğilimini azaltır ve genel model performansını artırır.

Gradient Boosting

AdaBoost'un bir geliştiğidir ve her iterasyonda modelin hatalarını azaltmak için gradient descent algoritmasını kullanır. Genellikle ağaç tabanlı zayıf öğrenicilerle kullanılır ve güçlü bir tahmin modeli oluşturur.

XGBoost (Extreme Gradient Boosting)

XGBoost, Gradient Boosting'in hızlandırılmış ve daha etkili bir versiyonudur. Bu algoritma hem hız hem de performansı ile popüler olmuştur ve birçok Kaggle yarışmasında kazanan modellerde kullanılmıştır.

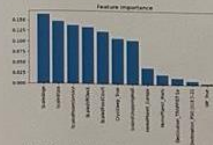
LightGBM

Microsoft tarafından geliştirilen LightGBM, Gradient Boosting'in bir başka uygulamasıdır. Daha hızlı eğitim süresi daha verimli hafıza kullanımı ile bilinir. Ayrıca, büyük veri setleriyle çalışma yeteneği ile de dikkat çeker.

CatBoost

Yandex tarafından geliştirilen CatBoost, özellikle kategorik verilerle çok iyi çalışabilen bir gradient boosting algoritmasıdır. Bu özellik, onu diğer boosting algoritmalarından ayırır.

Uygulama



Feature	Importance
Survived	0.12
Sex	0.11
Age	0.10
SibSp	0.09
Pclass	0.08
Embarked	0.07
Fare	0.06
Smoker	0.05
FamilySize	0.04
IsAlone	0.03
AgeBand	0.02
SexMale	0.01
EmbarkedC	0.01
EmbarkedQ	0.01
EmbarkedS	0.01
Age0-17	0.01
Age18-25	0.01
Age26-35	0.01
Age36-45	0.01
Age46-55	0.01
Age56-65	0.01
Age65+	0.01

Veri temizlendikten sonra öznitelik seçim tekniklerinden rassel orman tekniği uygulanmıştır ve bu sayede en önemli değişkenlerin hangileri olduğu ilk şekilde görülmektedir. Daha sonra bağımlı değişkeni en az etkileyen değişkenler çıkarıldı ve model tahmin edildi. 2. şekil bir lojistik regresyon modellemesidir ve gözüktüğü üzere 0.75'lik bir performansa sahiptir fakat 3. şekilde de görüldüğü üzere Xgboost uygulandığında modelde artış gerçekleşmiştir ve 0.78'lik bir skora ulaşılmıştır.



İSTATİSTİK BÖLÜMÜ

MÜŞTERİ YAŞAM BOYU DEĞERİ TAHMİNİNİN PAZARLAMAYA ETKİSİ

Aleyna ATAN 18023067

Danışman: Prof.Dr.Gülhayat GÖLBAŞI ŞİMŞEK

ÖZET

Bu tez çalışmasında müşteri yaşam boyu değerinin şirketlerin yapabilecekleri pazarlama stratejileri ile müşterileri portföylerinde uzun vadede nasıl tutabilecekleri üzerine bir inceleme yapılmıştır. Bu çalışma esnasında müşteri yaşam boyu değeri tanımı, hesaplama yöntemi ve bu değerın uygulama alanları hakkında bilgiler verilmiştir. Aynı zamanda bir müşterinin bir işletme için, müşteriyi elde tutmanın ekonomik açıdan ne kadar faydalı olabileceği hakkında bir hesaplama yapılmıştır. Bu hesaplama yapılırken müşteri o işletmeden ne kadar sıklıkla satın almaya başladığına, hangi periyotlarda satın alımı yapıldığına, kaç defa satın almaya başladığına ve bu satın almaların parasal değerine ait veriler RFM metoduna tabii tutulmuştur. Sonraki aşamada müşterinin parasal değeri, yani toplam fiyatı bulunmuştur. BG/NBD modelini kullanarak her müşteri için satın almış sayısını tahmini yapılmıştır. Bu tahminler de yapıldıktan sonra Gama-Gama modeli ile müşteri yaşam boyu değerini tahmin etmek için gereken parasal değerin tahmini yani her işlem için en olası değer tahmini yapılmıştır.

MÜŞTERİ YAŞAM BOYU DEĞERİ

Müşteri yaşam boyu değeri (Customer Lifetime Value, CLTV), bir müşterinin bir şirketle bağı süresince sağlanmış olacağı toplam maddi değeri ifade eder. Müşteri yaşam boyu değeri, bir şirketin müşterisi ile kurduğu ilişki süresince müşterinin yapacağı satın almalar, vermiş olduğu hizmetlerde elde edeceği gelirleri ve diğer sağladığı tüm değerlerin toplamı olacak şekilde hesaplanır.

MÜŞTERİ YAŞAM BOYU DEĞERİNİN HESAPLANMASI

Müşteri yaşam boyu değeri (CLTV), bir işletme müşterisinin sağlayacağı tahmini değeri hesaplamak için kullanılan bir değerdir. Müşteri yaşam boyu değeri, hesaplaması aşağıda verilen formül üzerinden yapılır:

Müşteri Yaşam Boyu Değeri (CLTV) = (Müşteri Değeri / Müşteri Kaybetme Oranı) * Kar Marjı

- Müşteri Değeri = Ortalama Sipariş Değeri * Satın Alma Sıklığı
- Ortalama Sipariş Değeri = Toplam Fiyat / Toplam İşlem
- Satın Alma Sıklığı = Toplam İşlem / Toplam Müşteri Sayısı
- Müşteri Kaybetme Oranı = 1 - Birden Fazla Alışveriş Yapan Müşteri Oranı
- Birden Fazla Alışveriş Yapan Müşteri Oranı = Birden Fazla Alışveriş Yapan Müşteri Sayısı / Toplam Müşteri Sayısı
- Kar Marjı = Toplam Fiyat * 0.10 (10 değişkenlik olabilir).

MÜŞTERİ YAŞAM BOYU DEĞERİ HESAPLAMA MODELLERİ

- Berger ve Nasr'ın müşteri yaşam boyu değeri hesaplama modeli
- Bruhn'un müşteri yaşam boyu değeri hesaplama modelleri
- Collings ve Baxter'in müşteri yaşam boyu değeri hesaplama modeli
- Gelbrich ve Wünschmann'ın müşteri yaşam boyu değeri hesaplama modeli
- RFM modeli

RFM analizi, pazarlamada özellikle direkt pazarlamada en çok bilinen ve uygulanan müşteri bölümlenme yöntemlerinden biridir. RFM; Recency, Frequency, Monetary kelimelerinin baş harflerinden oluşur. Recency, müşterinin son işleminin güncelliğine, frequency işlem sıklığına, monetary de müşterinin harcadığı toplam parayı ifade etmektedir. RFM analizi 1970'li yıllardan beridir kullanılan bir yöntemdir (Bery ve Linoff, 2004: 447-70). Mevcut müşterilerimiz arasında yeni bir teklife cevap vermeye hazır olanı hangisidir sorusuna cevap veren bir tekniktir. Bu teknik genellikle doğrudan pazarlamada kullanılır.

PAZARLAMA

4P	4C
Ürün (Product)	Tüketiciye Sağlanan Değer (Customer Value)
Fiyat (Price)	Tüketiciye Maliyeti (Cost to Customer)
Fiziksel Dağıtım (Place)	Kolaylık (Convenience)
Tutundurma (Promotion)	İletişim (Communication)

Pazarlama, bir şirketin ürünlerini veya hizmetlerini hedef kitleye tanıtmak, müşteri talebini oluşturmak ve satışları arttırmak için kullanılan stratejilerin ve faaliyetlerin bir bütünüdür. Pazarlama, müşteri ihtiyaçlarını anlamak, ürün veya hizmetin değerini iletmek, müşteri ilişkilerini yönetmek ve rekabet avantajı elde etmek gibi amaçları içerir.

PAZARLAMA VE MÜŞTERİ YAŞAM BOYU DEĞERİ (CLTV) İLİŞKİSİ

Pazarlama ve müşteri yaşam boyu değeri arasında yakın bir ilişki bulunmaktadır. Pazarlama, bir şirketin ürünlerini veya hizmetlerini hedef kitleye tanıtmak, marka bilinirliğini arttırmak ve müşteri talebini oluşturmak için kullanılan stratejilerin ve faaliyetlerin bir bütünüdür. Müşteri yaşam boyu değeri ise, bir müşterinin şirketle olan ilişkisi boyunca sağladığı toplam gelirin, maliyetin ve karlılığın ölçüsüdür. Pazarlama faaliyetleri, müşteri yaşam boyu değerini arttırmak için önemli bir rol oynar. İyi bir pazarlama stratejisi, müşteriye değer sunan ürünlerin veya hizmetlerin tanıtımını, marka bağlılığını artırıcı faaliyetleri ve müşteri deneyimini iyileştirme çabalarını içerir. Bu şekilde, müşterilerin şirketle olan ilişkileri güçlenir ve müşteri sadakati oluşturulur.

PAZARLAMADA TEMEL YAKLAŞIMLAR

Pazarlamadaki yaklaşımları temel olarak:

- İlişki pazarlaması yaklaşımı,
- Müşteri odaklı yaklaşım,
- Değer temelli yaklaşım,
- Pazar yönelimli yaklaşım diye maddelemek mümkündür.

UYGULAMA

Müşteri yaşam boyu değerini hesaplamak şirketler açısından oldukça önemli bir hesaplamadır. Müşterinin yaşam boyu değerini bulmak buna göre hareket etmek müşterileri seçmek ve müşteriye özel iletişim stratejilerine ilişkin karar vermek için iyi bir temel oluşturur. Müşteri yaşam boyu değeri, pazarlama yapacağımız analizlerden biri RFM analizidir. RFM analizi ve müşteri yaşam boyu değeri, pazarlama stratejileri ve müşteri ilişkileri yönetimi için önemli araçlardır. Bu analizler, şirketlerin müşterilerini daha iyi anlamalarına ve onlara daha kişiselleştirilmiş deneyimler sunmalarına yardımcı olur.

Bu uygulamada kullanılan veri seti, University of California Irvine bünyesindeki Machine Learning Repository'den alınmıştır. Veri kümesi, İngiltere merkezli bir e-ticaret şirketinin 01/12/2010 ve 09/12/2011 tarihleri arasında gerçekleşen, 541.910 adet online alışveriş işlemlerini içermektedir.

VERİ ANALİZİ

Veri setinde tüm düzeltmeler yapıldıktan sonrasında müşteri yaşam boyu değerini (CLTV) tahmin etmek için bazı tahminler yapılmıştır. Bu tahminler, geçmiş satın alma davranışlarını analiz etmeye yarar ve en iyi müşterileri bulmak için kullanılan bir pazarlama tekniğidir. Müşterilerin ne kadar sıklıkla alışveriş yaptıklarını, şu ana kadar harcadıkları toplam tutarı, işletmeden en son ne zaman alışveriş yaptıklarını gibi benzer bilgileri içerir. Bu değerler şunlardır; sıklık, yenilik, T, parasal değerdir.

Veri setinde tüm düzeltmeler ve ön analizler yapıldıktan sonra, BG (Beta, Geo) kullanılarak, BG/NBD modelini yeni veri çerçevesine uygulandı ve her bir müşteri için satın almış sayısını tahmini edildi. BG/NBD modelinin özeti çıkarıldı.

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	Total Price
0	59836	891234 WHITE HANGING HEART T-LIGHT HOLDER	6.0	2010-12-01 08:26:00	2.55	17920.0	United Kingdom	15.30
1	59836	71053 WHITE METAL LANTERN	6.0	2010-12-01 08:26:00	3.39	17920.0	United Kingdom	20.34
2	59836	84408 CREAM CUPID HEARTS COAT HANGER	6.0	2010-12-01 08:26:00	2.75	17920.0	United Kingdom	16.50
3	59836	84293 KNITTED UNION FLAG HOT WATER BOTTLE	6.0	2010-12-01 08:26:00	3.39	17920.0	United Kingdom	20.34
4	59836	84202 RED WOOLLY HOTTIE WHITE HEART	6.0	2010-12-01 08:26:00	3.39	17920.0	United Kingdom	20.34

Veri seti ön işleme sokulduktan sonrasında müşteri yaşam boyu değerini (CLTV) tahmin etmek için bazı tahminler yapılmıştır. Bu tahminler, geçmiş satın alma davranışlarını analiz etmeye yarar ve en iyi müşterileri bulmak için kullanılan bir pazarlama tekniğidir. Müşterilerin ne kadar sıklıkla alışveriş yaptıklarını, şu ana kadar harcadıkları toplam tutarı, işletmeden en son ne zaman alışveriş yaptıklarını gibi benzer bilgileri içerir. Bu değerler şunlardır; sıklık, yenilik, T, parasal değerdir.

Veri setinde tüm düzeltmeler ve ön analizler yapıldıktan sonra, BG (Beta, Geo) kullanılarak, BG/NBD modelini yeni veri çerçevesine uygulandı ve her bir müşteri için satın almış sayısını tahmini edildi. BG/NBD modelinin özeti çıkarıldı.

CustomerID	frequency	recency	T	monetary_value
12747.0	10.0	367.0	369.0	375.725000
12748.0	112.0	373.0	373.0	257.314911
12749.0	3.0	210.0	213.0	999.106667
12820.0	3.0	323.0	326.0	256.573333
12823.0	4.0	222.0	296.0	252.450000

Parasal değerlerin hesaplanması için Unit Price ve Quantity değişkenleri çarpılarak yeni bir özellik oluşturuldu. Bu yeni özellik ile her bir müşteri için toplam fiyat bulundu.

Müşteri yaşam boyu değerinin veri setindeki 6 ay için tahmin edilebilmesine ilişkin tüm modeller hazırlandı. BG/NBD ve Gama Gama modelleri kullanılarak Müşteri Yaşam Boyu Değeri tahmin edildi.

coef	se(coef)	lower 95% bound	upper 95% bound	
r	2.381731e+00	1.013873e-01	2.183012e+00	2.580451e+00
alpha	1.120749e+02	5.267568e+00	1.017504e+02	1.223993e+02
a	2.301100e-14	3.668831e-10	-7.190678e-10	7.191138e-10
b	2.398753e-05	3.819882e-01	-7.486728e-01	7.487208e-01

Müşteri yaşam boyu değeri tahmini sonrası, her müşterinin segmentine göre ürün sunmak, alt segmentte bulunan müşterilerin yaşam boyu değerlerini arttırmak amacıyla pazarlama planı oluşturmak ve müşteri kazanımı için maliyeti azaltmak amacı ile üst segmentlere odaklanmak amacı ile müşteriler yaşam boyu değerlerine göre segmentlere ve gruplara ayrıldı.

Segment	frequency	recency	T	monetary_value	expected_purc_6_months	6_Months_CLV
Hibernating	3.174713	221.337931	292.245977	147.152900	2.533847	384.667903
Need Attention	4.027650	240.140553	284.440092	265.378345	2.989161	745.573232
Loyal Customers	5.592166	242.944700	276.027650	365.362234	3.759721	1234.753807
Champions	11.108046	261.354023	280.475862	584.771128	5.994943	3095.935737

SONUÇ VE ÖNERİLER

Tahmin sonuçlarına bakıldığında, şirketlerin müşterilerinin yaşam boyu değerlerini tahmin etmelerinin gelecek dönemler için çok fazla efor sarf etmeksizin gelirlerini arttırmalarını sağlayacağını ortaya koymaktadır. Bu sonuçlardan yola çıkılarak, şirketler ilişkilerini devam ettirmesini bekledikleri müşterilerinin sadakatini sağlamak ve arttırmak için pek çok farklı yol izleyebilir. Müşterileri, yaşam boyu değerlerine göre gruplayabilir ve bu gruplar için özel kampanyalar uygulayabilir. Müşterilerin yaratacakları fayda göz önüne alınarak yapılacak kampanyalarla müşterilere özel faydalar sunulması sağlanırken bu sayede şirket gereksiz maliyetlerden kaçınılmış olacaktır. Aynı zamanda yaşam boyu müşteri değeri farklı analizlere de giridi olabilir. Her müşterinin firmaya ne kadar fayda sağlayacağını bilmek firmanın gelecek yatırım kararlarında pazarlama faaliyetlerine kadar pek çok avantaj sağlayacaktır.

KAYNAKÇA

- [1] YAPRAKLI, T. Ş., & Keser, E. (2008). MÜŞTERİ YAŞAM BOYU DEĞERİNİN ANALİZİ: BİR SAHA ARAŞTIRMASI. Atatürk Üniversitesi Sosyal Bilimler Enstitüsü Dergisi, 12(2), 483-503.
- Tran, K. G., Nguyen, V. H., & Ho, T. Customer segmentation analysis and customer lifetime value prediction using Pareto/NBD and Gamma-Gamma model.
- Yıldırım, Y. O. PAZARLAMADA DEĞİŞEN GÜÇLÜ DİJİTAL DÖNEME GEÇİŞ VE PAZARLAMA 4.0 YAKLAŞIMLARI.



ÖZET

Bu çalışmada, deprem sonrası travma düzeyi ve deprem risk algısının işte üretkenlik üzerindeki etkisini belirlemek için bir anket çalışması yapılmıştır. Ankette, katılımcıların deprem sonrası yaşadıkları travma düzeyi ve deprem risk algısı ölçülmüştür. Ayrıca, katılımcıların işte üretkenliklerini belirlemek için sorular sorulmuştur. Çalışma kapsamında, veri analiz yöntemleri kullanılmıştır. Veri analiz yöntemleri olarak kişisel değişkenleri incelemek için betimleyici istatistiksel yöntemler ve araştırma değişkenleri arasındaki etkiyi incelemek için çoklu regresyon analizi kullanılmıştır. Çalışma kapsamında, deprem sonrası travma düzeyinin, depresyon-anksiyete-stres seviyesinin, deprem risk algısının ve iş tatmininin; işte üretkenlik üzerinde anlamlı bir etkisi olduğu belirlenmiştir.

DEPREM SONRASI TRAVMA VE DEPREM RİSK ALGISI

Depremzedelerin yaşadıkları afet sonucunda bazı psikolojik sorunlarla karşı karşıya kaldıkları söylenebilir. Travmatik olaylar olarak adlandırılan bu sorunlar, depremzedelerin deprem sonrası uzama ile etkileşimini gerektirmektedir. Risk, olumsuz olaylar nedeniyle ortaya çıkabilecek sonuçların (kayıp veya hasar) olasılığıdır.

DASÖ-21

DASÖ, hasta tarafından doldurulabilen, anksiyete ve depresyonun temel belirtilerini içeren, yüksek psikometrik standartları karşılayan, anksiyete ve depresyonu birbirinden ayırt edebilen bir ölçek oluşturmak amacıyla geliştirilmiştir.

İŞTE ÜRETKENLİK VE İŞ TATMİNİ

İşte üretkenlik, hem işin sonuçları hem de sonuçlara ulaşma sürecindeki kritik davranışları içerir. İş tatmini, "özelliklerinin değerlendirilmesinden kaynaklanan bir işe olumlu bir duygu" olarak tanımlanmaktadır.

ANKET SONUÇLARI VE GÜVENİLİRLİK ANALİZİ

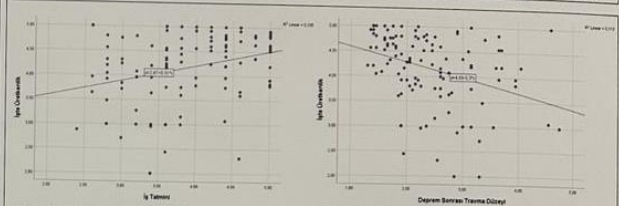
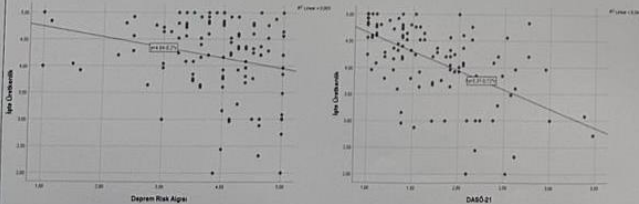
DEMOGRAFİK SORULARIN SONUÇLARI			DEĞİŞKENLERE YÖNELİK BETİMLEYİCİ İSTATİSTİKLER								
Soru	İçerik	Oran (%)	Değişken	Min.	Maks.	Ort.	SS				
Cinsiyetiniz nedir?	Erkek	60	652,2	Deprem Sonrası Travma Düzeyi	1	5	2,41	0,78			
	Kadın	55	647,8		Deprem Risk Algısı	1	5	3,91	0,89		
	Hesabı	0	0,0			DASÖ-21	1	5	1,68	0,57	
	Ortaokul	1	10,9				İşte Üretkenlik	1	5	4,26	0,7
	Lise	1	10,9					İş Tatmini	1	5	3,85
Üniversite	88	1018,5	Güvenilirlik Analizi Bulguları								
Üniversite	22	118,1		Değişkenler	Cronbach's Alpha				Madde Sayısı		
Çalışma durumunuz?	Özel Sektör	89		1077,4	Deprem Sonrası Travma Düzeyi	0,916			8		
	Emekli	2		14,7	Deprem Risk Algısı	0,95	20				
	Öğrenci	2		14,7	DASÖ-21	0,943	21				
	Serbest Meslek Çalışmıyorum	5	14,3	İşte Üretkenlik	0,96	25					
	Çalışmıyorum	1	10,9	İş Tatmini	0,776	5					
Çalışma yılınız?	0-5 yıl	74	1048,3	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	6-10 yıl	15	113	Çoklu Regresyon Analizi Sonuç Tablosu							
	11-15 yıl	14	112,2	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	16-21 yıl	4	105,5	Çoklu Regresyon Analizi Sonuç Tablosu							
	22+ yıl	8	107	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
İşyerindeki istihdam türünüz?	Yeni Zamanlar	95	1080,9	Çoklu Regresyon Analizi Sonuç Tablosu							
	Yeni Zamanlar	16	110,7	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	Serbest Çalışan	4	105,5	Çoklu Regresyon Analizi Sonuç Tablosu							
	ÖZG	54	1047	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	Edisyon/İşleme Ofisi	17	1123,5	Çoklu Regresyon Analizi Sonuç Tablosu							
İşyerindeki çalışma şekliniz?	Andevis Bölgesi	19	1126,5	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	Marmara Bölgesi	2	14,7	Çoklu Regresyon Analizi Sonuç Tablosu							
	Karadeniz Bölgesi	89	1077,4	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	İç Anadolu Bölgesi	0	0,0	Çoklu Regresyon Analizi Sonuç Tablosu							
	Doğu Anadolu Bölgesi	0	0,0	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
Yaşadığınız bölge?	Sarıyer/Beşiktaş Bölgesi	0	0,0	Çoklu Regresyon Analizi Sonuç Tablosu							
	0-5 yaşlık	17	1048,8	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	6-10 yaşlık	24	1020,9	Çoklu Regresyon Analizi Sonuç Tablosu							
	11-15 yaşlık	17	1048,8	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	16-21 yaşlık	15	113,7	Çoklu Regresyon Analizi Sonuç Tablosu							
Ortaokulunuz bina kaç yaşlık?	22+ yaşlık	41	1037,7	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							
	11-21 yaşlık	1	10,9	Çoklu Regresyon Analizi Sonuç Tablosu							
	Hiçbir şey yok	1	10,9	Normal dağılım grafikleri; tahminlere ait hatalar normal dağılmaktadır.							

ÇOKLU REGRESYON ANALİZİ

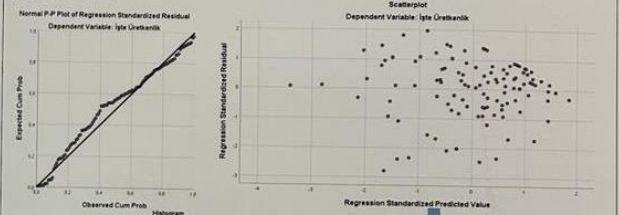
H0: Deprem Sonrası Travma Düzeyinin, Depresyon-Anksiyete-Stres Seviyesinin, Deprem Risk Algısının, İş Tatmininin; İşte Üretkenlik üzerinde etkisi yoktur.

- H0a: Deprem Sonrası Travma Düzeyi ve İşte Üretkenlik arasında ilişki yoktur.
- H0b: Depresyon-Anksiyete-Stres Seviyesi ve İşte Üretkenlik arasında ilişki yoktur.
- H0c: Deprem Risk Algısı ve İşte Üretkenlik arasında ilişki yoktur.
- H0d: İş Tatmini ve İşte Üretkenlik arasında ilişki yoktur.

Bağımsız değişkenler ile bağımlı değişken arasındaki doğrusallık kontrolü dağılım grafiği (scatter plot) ile sağlanmıştır.



Dağılım grafiklerinde bulunan düz çizgiler değişkenler arasında doğrusallık olduğunu göstermektedir.



Kalıntı grafiği; noktaların grafik üzerinde oldukça yayıldığı ve uç değer olmadığı görülmektedir. Dolayısıyla eş varyanslılık vardır.

Değişkenler	β	Standart Hata	Beta	t	P	VIF
Sabit*	4,744	,364	-	13,038	,000	-
DASÖ Seviyesi	-,729	,184	-,175	-6,235	,000	1,858
Travma Düzeyi	,159	,099	,638	1,640	,104	2,371
İş Tatmini	,205	,129	,036	2,954	,004	1,047
Deprem Risk Algısı	-,132	,120	-,142	-1,954	,053	1,496

*İşe Üretkenlik (Bağımlı Değişken)
*R = ,649 *R² = ,421 *Adjusted R² = ,399

Değişkenler	İşte Üretkenlik	Deprem Risk Algısı	DASÖ-21	İş Tatmini	Deprem Sonrası Travma Düzeyi
İşte Üretkenlik	1,000	-,263	-,593	,312	-,335
Deprem Risk Algısı	-,263	1,000	,328	,045	,569
DASÖ-21	-,593	,328	1,000	-,174	,662
İş Tatmini	,312	,045	-,174	1,000	-,044
Deprem Sonrası Travma Düzeyi	-,335	,569	,662	-,044	1,000

KORELASYON TABLOSU

SONUÇ

Bu çalışmada, deprem sonrası travma düzeyi ve deprem risk algısının işte üretkenlik üzerindeki etkisini belirlemek için bir anket çalışması yapılmıştır. H0 hipotezi reddedilmiştir, yani etki anlamlıdır. Bağımsız değişkenlerin bağımlı değişken üstündeki etkilerine bakıldığında, depresyon-stres-anksiyete seviyesi ve iş tatmini değişkenlerinin etkisi olduğu görülmüş fakat deprem sonrası travma düzeyinin işte üretkenlik üzerinde etkisi için yeterli kanıt bulunamamıştır. Deprem risk algısı değişkeni için ise marjinal düzeyde etkisi görülmüş olup bu değişkenin etkisini araştırmak için çalışma geliştirilmesi önerilmektedir.

KAYNAKÇA

- [1] Clark, L. A., & Watson, D. (1991). Tripartite model of anxiety and depression: Psychometric evidence and taxonomic implications. *J Abnorm Psychol*, 100(3), 316-336
- [2] Seçer, İ. (2015). Spss ve Lisrel ile Pratik Veri Analizi, 2. Baskı. ss: 28. Anı Yayıncılık, Ankara.
- [3] Endicott J, Nee J (1997) Endicott Work Productivity Scale (EWPS): A new measure to



İSTATİSTİK BÖLÜMÜ

ARIMA İle Türkiye Buğday Verimi Tahmini

Tuğba ÇERİK 17023062

Danışman: Doç. Dr. Gülder KEMALBAY

ÖZET

Tarım Mahsulleri Ofisi tarafından 2021 yılında yayımlanan 1938- 2021 yılları arasında kaydedilen Türkiye buğday verimi veri seti ile geleneksel doğrusal zaman serisi modellerinden biri olan ARIMA yöntemi kullanılacaktır. Araştırmanın amacına ulaşmak için buğday verim veri seti, ARIMA modeli ile elde edilen Otoregresif parametre AR mertebesi kadar gecikmeli bağımsız değişken kullanılacaktır. Uygunluk (fitness) fonksiyonu için MAE (mean absolute error; ortalama mutlak hata) kriteri seçilecektir. ARIMA tahmin modeli oluşturulacaktır.

BUĞDAY

İnsanların tüketim alışkanlıklarını belirlemeye yönelik yapılan araştırmalar sonucunda, pirinçten sonra en yoğun olarak tüketilen tahıl ürününün buğday olduğu belirlenmiştir. Bu çalışmada veri seti olarak, TMO (2021) tarafından yayımlanan 1938-2021 yılları arası buğday verisi kullanılmıştır.

Tablo 1. Türkiye'de Her On Yıllık Periyot için Buğday Üretim Miktarı

YILLAR	ÜRETİM (Ton)
1941	3.483.147
1951	5.600.000
1961	7.000.000
1971	13.500.000
1981	17.000.000
1991	20.400.000
2001	19.000.000
2011	21.800.000
2021	17.650.000

ARIMA

Box-Jenkins yöntemi olarak bilinen ARIMA, zaman serisi verilerini tahmin etmek için kullanılan popüler istatistiksel bir modeldir. Verideki geçmiş verileri kullanarak noktalar arasındaki oto korelasyonu açıklamayı hedefleyen, basit ve esnek yapıda bir yöntemdir. Modelin en önemli özelliği, minimum sayıda parametre ile kurulmasıdır.

- AR (oto regresyon) gözlem değerleri ve gecikmeli veriler arasındaki bağımlı ilişkiyi kullanır.
- I (entegre), gözlemlerin farkı alınarak zaman serisi durağan hale getirilir.
- MA (hareketli ortalama) bir zaman serisinin herhangi bir dönemdeki gözlem değeri ile hata terimlerinin lineer olarak ifade edildiği yöntemdir.

Matematiksel olarak ARIMA modeli Denklem (1) gibidir:

$$\Delta d Z_t = c + (\theta_1 \Delta d Z_{t-1} + \dots + \theta_p \Delta d Z_{t-p}) - (\beta_1 et-1 + \dots + \beta_q et-1) + et \quad (1)$$

DENEYSEL SONUÇLAR



Şekil 1. Yıllık buğday verimi zaman serisi grafiği

Ayırılabilir bazı desenler olduğu görülmektedir. Serinin, ortalama da durağan olmadığı yani ortalamasının zaman içinde sabit ilermediği ve zamana bağlı olarak farklılık gösterdiği görülmektedir.

Tablo 2. Verim serisinin birinci farkı için ADF test istatistiği

Model	Level								
	kesmesiz ve trendsiz			kesmeli ve trendsiz			kesmeli ve trendli		
Kritik	%1	%5	%10	%1	%5	%10	%1	%5	%10
Değerler	-2.6	-1.95	-1.61	-3.51	-2.89	-2.58	-4.04	-3.45	-3.15
DF Test İstatistiği	-3.7497			-7.9045			-7.9414		

DF test istatistikleri, %1, %5 ve %10 anlamlılık düzeyleri için kritik değerlerden küçüktür. DF Test İstatistiği, kritik değerlerden küçük ise H0 hipotezi reddedilir. Seri birim kök içermektedir. İlk farkı alınan seri durağan hale getirilmiştir. Bu nedenle, $\tau_1 = \Delta y_t$, serinin bir birim kök içerdiği H_0 hipotezi reddedilir, yani $\Delta Y_t = (1 - L)Y_t$ serisi birinci farkta durağandır. Bu nedenle, ARIMA modelimizin entegre kısmı (I) 1'ye eşit olacak, yani d=1 olacaktır.

Tablo 3. Aday AICc ARIMA(p,d,q) modelleri için değerler

Model	AICc	Model	AICc
ARIMA(2,1,2)	614.394	ARIMA(0,1,0)	628.5677
ARIMA(0,1,0)	629.9381	ARIMA(1,1,1)	610.9329
ARIMA(1,1,0)	617.4868	ARIMA(0,1,2)	610.9345

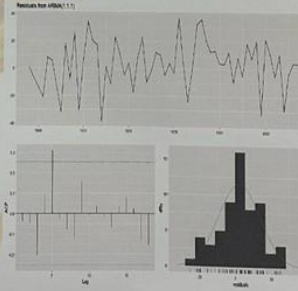
Tablo 4. ARIMA model sonuçları

	AR1	MA1
Coefficients	0.0388	-0.5689
Standard Error	0.0000	0.0000
p-value	0.0000	0.0000

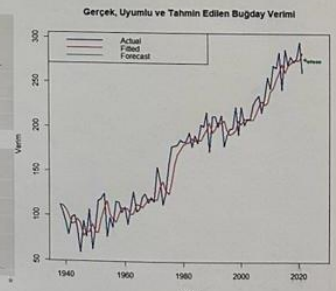
Verilen çıktıya dayanarak, ARIMA(1,1,1) modelini temsil eden denklem şu şekilde ifade edilebilir:

$$r_t = 0.0388r_{t-1} - 0.5689r_{t-1} + e_t$$

Bu denklem, bir zaman serisi verisindeki değerlerin önceki değerleri ve hata terimi arasındaki ilişkiyi açıklar. İfade edilen model, birinci dereceden otoregresif (AR) ve birinci dereceden hareketli ortalama (MA) bileşenlerini içermektedir. ARIMA(1,1,1) modelinin, veri serisinin durağan hale getirildikten sonra (1 kez fark alındıktan sonra) uygulanan bir regresyon modeli olduğunu ifade eder.



Şekil 2. ARIMA(1,1,1) modeli için artık analizi



Şekil 3. Actual, Fit ve tahmin grafiği

Şekil 12'de verilen artık analizi, ACF grafiğinde önemli pikler olmadığını ve artıkların yaklaşık olarak normal dağılımına göstermektedir. Seride otokorelasyon gözlenmemektedir. Bununla birlikte, artıkların varyansı zaman içinde sabit gibi görünmemektedir.

SONUÇ

Bu çalışmanın sonuçları, ARIMA(1,1,1) modelinin buğday verimi tahmini için kullanışlı bir araç olduğunu göstermektedir. Bu tür çalışmalar, on birinci kalkınma planı gibi, tarım sektöründe verimliliği artırmak amacıyla ve gelecekteki buğday üretimini daha iyi planlamak için önemli bir bilgi sağlayabilir.

KAYNAKÇA

- [1] KEMALBAY, G., & KORKMAZOĞLU, O. B. (2021). Sarıma-arch versus genetic programming in stock price prediction. Sigma Journal of Engineering and Natural Sciences, 39(2), 110-122.
- [2] T.C. Cumhurbaşkanlığı Strateji ve Bütçe Başkanlığı (2019). "On Birinci Kalkınma Planı (2019-2023)", <http://www.sbb.gov.tr/wpcontent/uploads/2019/07/On-Birinci-Kalkinma-Planı.pdf> [Erişim Tarihi: 27.06.2021].
- [3] Bakanlık, T. (2019). Tarımsal Ekonomi ve Politika Geliştirme Enstitüsü. Tarım Ürünleri Piyasası, Buğday, Ocak, 2021.



ÖZET

Tezimin özetine ünlü Amerikalı İstatistikçi W. Edwards Deming'in bir sözüyle başlamak istiyorum; "İstatistikler, gerçekliğin dilidir; anlamak için dikkatli bir şekilde dinlemeli ve doğru bir şekilde yorumlamalıyız. Bu tezde belirtilen konu üzerinde SPSS kullanılarak bazı istatistiksel testler yapılarak alakalılık derecesine bakılıp bu testlerin yorumları yapılmıştır. İlk önce spor ve sporun faydaları hakkında kısa bir bilgi verilip sonrasında yapılan testler sonucu bulgular ve sonuç ile öneriler verilmiştir.

ÜNİVERSİTE ÖĞRENCİLERİNDE SPOR YAPMANIN AKADEMİK BAŞARIYA

ETKİSİ

GİRİŞ

Spor, önceden belirlenmiş kurallara göre bireysel veya takım halinde yapılan, genellikle rekabete dayalı yarışma ve kişisel eğlence veya mükemmelliğe ulaşmak için yapılan fiziksel veya zihinsel aktivite. Sporları kabaca homo sapiens türünün medeniyete ulaşmadan önce doğayla veya diğer canlılarla yaptığı fiziksel mücadelelerin günümüzdeki medeni karşılığı olarak da tanımlayabiliriz.

Spor yapmanın bazı faydaları da aşağıdaki gibidir;

Kasları güçlendirir. Vücut dayanıklılığını artırır, esneklik verir. Kemikleri güçlendirir ve kemik yoğunluğunun artmasını sağlar. Vücut enerjisinin artmasını sağlar, zindelik verir. Kilonuzu dengede tutmanızı sağlar ve eklemere binen yükü azaltır. Depresyonu ve kaygıyı engeller. Uykuyu düzenler. Rahat ve derin bir uyku sağlar. Kalp hastalıklarını önler. Tansiyonu düzenler. Yüksek kolesterol riskini azaltır. Daha kaliteli bir yaşam sunar.

Kalp rahatsızlıkları risklerini azaltma, tansiyonu normal tutma, kandaki Triglycreid oranını sabit tutma, kolesterolü düzenlenme, kalp-damar dolaşımını geliştirme gibi etkileriyle sağlık sorunları yaratacak durumları engeller.

ÇALIŞMANIN AMACI

Bu araştırmanın hedefi, üniversite öğrencilerinde sporun akademik başarıya etkisini incelemek ve sporun akademik performans üzerindeki olası faydalarını belirlemektir. Araştırma, sporla ilgili faktörlerin öğrencilerin ders başarıları ve notları üzerindeki ilişkisini anlamak için yapılmaktadır. Bu bilgiler, sporun öğrencilerin akademik performansını desteklemede nasıl kullanılabilirliği konusunda stratejiler geliştirmek için kullanılabilir. Bu araştırma, üniversite öğrencilerinde sporun akademik başarı üzerindeki etkisini belirlemek için önemli bir ışık tutacaktır. Araştırma sonuçları, sporun akademik performansı nasıl etkileyebileceği, sporun öğrencilerin ders başarıları ve notları üzerindeki ilişkisi gibi konuları aydınlatacaktır. Ayrıca, sporun öğrencilerin konsantrasyon, bilişsel fonksiyonlar ve stres yönetimi gibi becerilerini nasıl geliştirebileceğini göstererek, üniversiteler ve eğitimciler için sporun akademik başarıyı desteklemek için kullanılabilir bir araç olduğunu vurgulayacaktır.

MATERYAL VE METODLAR

Verileri sorularını özel olarak hazırlanmış ve toplamda 150 kişiye ulaşan bir anket hazırlayarak toplanmıştır. Soruların ana konuları yapılan sporun çeşidi, akademik başarı ve demografik bilgidir. Elde edilen veriler SPSS(İstatistik paket programı) ile analiz edilmiştir.

BULGULAR

Katılımcıların spor yapma durumları ile genel akademik not ortalamaları arasındaki ilişkiye bakılmıştır. Bu ilişkiye göre GANO puanı 1-2 arasında olan 7 kişi (%50) spor yapıyorken 7 kişi (%50) spor yapmamaktadır, GANO puanı 2-3 arasında olan 58 kişi (%77.3) spor yapıyorken 17 kişi (%22.7) spor yapmamaktadır, GANO puanı 3-4 arasında olan 45 kişi (%73.8) spor yapıyorken 16 kişi (%26.2) spor yapmamaktadır.

Katılımcıların kendi okul başarı değerlendirmeleri ile spor yapmaya başladıktan sonraki akademik başarı durumları arasındaki ilişki analiz edilmiştir. Elde edilen sonuçlara göre; okul başarısını 'çok iyi' olarak değerlendiren öğrencilerin 41'i (%73.2) spor yapmaya başladıktan sonra akademik başarısının arttığını savunurken 8'i (%14.3) değişmediğini veya azaldığını düşünmektedir. Okul başarısını 'idare eder' olarak değerlendiren öğrencilerin 34'ü (%38.6) spor yapmaya başladıktan sonra akademik başarısının arttığını savunurken 27'si (%30.7) değişmediğini veya azaldığını düşünmektedir. Okul başarısını 'çok kötü' olarak değerlendiren öğrencilerden 0'ı (%0) başarısının arttığını savunurken 3'ü (%50) başarısının değişmediğini veya azaldığını düşünmektedir. Katılımcıların spor yapma durumları ile kendi okul başarı değerlendirmeleri arasındaki ilişki analiz edilmiştir. Bulgulara göre okul başarısını 'çok iyi' olarak değerlendiren öğrencilerin; 47'si (%83.9) spor yaparken 9'u (%16.1) spor yapmamaktadır. 'idare eder' olarak değerlendiren öğrencilerin 60'ı (%68.2) spor yaparken 28'i (%31.8) spor yapmamaktadır. 'Çok kötü' olarak değerlendiren öğrencilerin 3'ü (%50) spor yaparken 3'ü (%50) spor yapmamaktadır.

Katılımcıların üniversite sınavına hazırlanırken düzenli olarak spor yapma durumları ile şu anki okul başarı değerlendirmeleri arasındaki ilişkiye bakılmıştır. Bu ilişkiye göre okul başarısını; 'çok iyi' olarak değerlendiren öğrencilerin 32'i (%57.1) üniversite sınavına hazırlanırken düzenli olarak spor yapmışken 24'ü (%42.9) yapmamıştır; 'idare eder' olarak değerlendiren öğrencilerin 29'u (%33) düzenli olarak spor yapmışken 59'u (%67) yapmamıştır; 'çok kötü' olarak değerlendiren öğrencilerin 1'i (%16.7) düzenli olarak spor yapmışken 5'i (%83.3) spor yapmamıştır.

Katılımcıların okulda göstermiş olduğu performanstan memnun olma durumları ile spor yapma durumları arasındaki ilişki analiz edilmiştir. Elde edilen sonuçlara göre; okuldaki performansından memnun olan öğrencilerin 91'i (%81.3) spor yapıyorken 21'i (%18.8) spor yapmamaktadır. Okuldaki performansından memnun olmayan öğrencilerin 10'u (%50) spor yapıyorken 9'u (%50) spor yapmamaktadır.

Katılımcıların okulda göstermiş olduğu performanstan memnun olma durumları ile spor yapmaya başladıktan sonraki akademik durumları arasındaki ilişki incelenmiştir. Bu ilişkiye göre; okuldaki performansından memnun olan öğrencilerin 71'inin (%63.4) spor yapmaya başladıktan sonra akademik başarıları artmışken 23'ünün (%20.5) akademik başarıları değişmemiş veya azalmıştır. Okuldaki performansından memnun olmayan öğrencilerin 4'ünün (%10.5) spor yapmaya başladıktan sonra akademik başarıları artmışken 15'inin (%39.5) akademik başarıları değişmemiş veya azalmıştır.

Katılımcıların okulda göstermiş olduğu performanstan memnun olma durumları ile kaç yıldır spor yaptıkları arasındaki ilişki analiz edilmiştir. Elde edilen bulgulara göre; okuldaki performansından memnun olan öğrencilerin 12'si (%10.7) 0-1 yıldır, 17'si (%15.2) 1-2 yıldır, 23'ü (%20.5) 3-4 yıldır ve 42'si (%37.5) 4 yıldan uzun süredir spor yapmaktadır. Okuldaki performansından memnun olmayan öğrencilerin 1'i (%2.6) 0-1 yıldır, 3'ü (%7.9) 1-2 yıldır, 3'ü (%7.9) 3-4 yıldır ve 11'i (%28.9) 4 yıldan uzun süredir spor yapmaktadır.

Katılımcıların okulda göstermiş oldukları performanstan memnun olma durumları ile haftada kaç gün antrenman yaptıkları arasındaki ilişki incelenmiştir. Elde edilen bulgulara göre; haftada 1-3 gün spor yapan öğrencilerin 51'i (%83.6) okul performansından memnunken 10'u (%16.4) memnun değildir. Haftada 3-5 gün antrenman yapan öğrencilerin 18'i (%81.8) okul performansından memnunken 4'ü (%18.2) memnun değildir.

SONUÇ VE ÖNERİLER

Bu çalışmada sporun akademik başarı üzerindeki etkisi ortaya çıkarılmıştır. Buna göre GANO puanı yükseldikçe spor yapan öğrenci oranının da yükseldiği görülmüştür. Bir başka bulguya göre okul başarısını çok iyi olarak değerlendiren öğrencilerinin büyük çoğunluğu spor yapmaya başladıktan sonra akademik başarılarının arttığını düşünmektedir. Düzenli ve aktif olarak spor yapan öğrencilerin ise kendilerini daha fazla başarılı olarak gördüğü tespit edilmiştir. Üniversite sınavına hazırlanırken düzenli olarak spor yapan öğrenciler şu an kendilerini daha fazla başarılı olarak değerlendirmektedir. Aynı zamanda şu an spor yapan öğrenciler okuldaki başarılarından daha fazla memnun oldukları saptanmıştır. Üniversitede kendisini başarılı olarak değerlendiren öğrencilerin anlamlı olarak çoğunluğu aktif olarak spor yapmaktadır. Öğrencilerin kendilerini akademik bağlamda başarılı görmeleri ile kaç yıldır spor yaptıkları arasında da anlamlı ilişki bulunmuştur. Buna göre spor yapma süresi arttıkça kendilerini başarılı olarak değerlendiren öğrencilerin oranı da artmaktadır. Elde edilen bu sonuçlara göre, okulda memnun olan öğrencilerin spor yapma oranının daha yüksek olduğu ve spor yapmaya başladıktan sonra akademik başarılarının arttığını görülmüştür. Ayrıca, düzenli olarak spor yapan öğrencilerin genellikle okuldaki performanslarından daha memnun oldukları gözlemlenmiştir.

Sporun akademik başarıyı olumlu yönde etkilediği yapılan testlerin sonuçları doğrultusunda ortaya çıkmıştır. Buna göre; kurumlar ve eğitim sistemleri, sporun akademik başarıyı desteklemesine olanak sağlamak için uygun ortamlar ve olanaklar sunmalıdır. Spor tesislerine erişim sağlamak, spor kulüpleri ve takımları oluşturmak, spor etkinliklerini teşvik etmek gibi adımlar atılabilir. Ayrıca, spor yapmayı teşvik eden programlar ve politikalar geliştirilmeli ve öğrencilere spor yapma fırsatları sunulmalıdır. Böylece öğrencilerin spor yapmaları teşvik edilebilir.

KAYNAKÇA

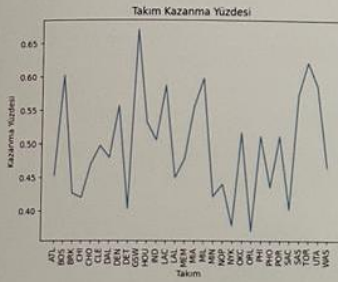
- <https://tr.wikipedia.org/wiki/Spor>
- <https://hthayat.haberturk.com/saglik/egzersiz/haber/1011758-spor-yapmanin->



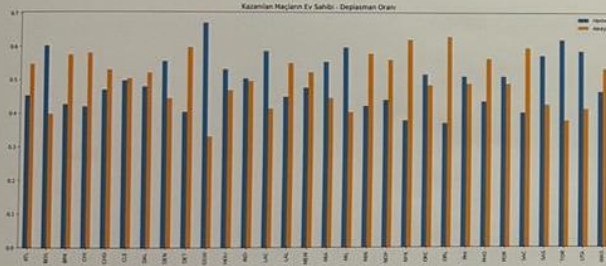
ÖZET

Bu çalışmada, NBA'nın 7 sezonuna altı maç istatistiklerini ve sonuçlarını içeren bir veri seti kullanılmıştır. Temel Bileşenler Regresyonu (PCR) ve Kısmi En Küçük Kareler Regresyonu (PLSR) yöntemleri kullanılarak bu veriler analiz edilmiş ve Destek Vektör Makinesi (SVM) algoritmasıyla sonuçlar tahmin edilmiştir. PCR ve PLSR yöntemleri kullanılarak boyut indirgeme işlemi gerçekleştirilmiş ve ardından SVM algoritmasıyla maç sonuçlarının tahmin edilmesi için bir model oluşturulmuştur. SVM, sınıflandırma ve regresyon problemlerinde etkili bir şekilde kullanılabilen bir makine öğrenimi algoritmasıdır. PCR, çok sayıda bağımsız değişkenin olduğu durumlarda boyut indirgeme sağlamak için kullanılan bir yöntemdir. PCR, veri setindeki varyansın büyük bir kısmını açıklayan bileşenleri belirler ve bunları regresyon analizinde kullanır. Bu şekilde, model karmaşıklığını azaltarak veri setinin özünü yakalamaya çalışır. PLSR hem bağımsız değişkenlerin hem de bağımlı değişkenin boyutunu azaltmayı hedefleyen bir yöntemdir. PLSR, bağımlı değişkenle en yüksek korelasyona sahip bileşenleri bulmaya odaklanır ve bu bileşenleri kullanarak regresyon analizi yapar. Bu sayede, önemli bilgileri korurken modelin karmaşıklığını azaltır.

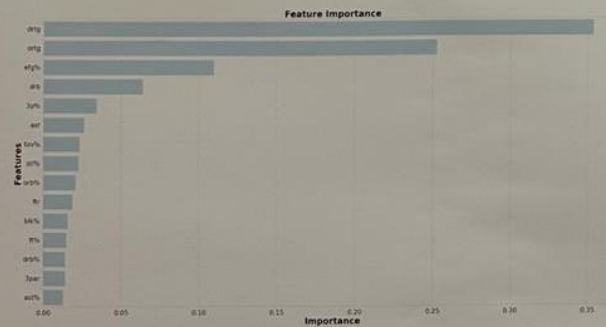
MAÇ SONUÇLARININ ANALİZİ



Yandaki grafikte tüm takımların 7 sezon boyunca oynadıkları maçların kazanma yüzdeleri gösterilmektedir. Tüm bu sezonlar boyunca en başarılı takım Golden State Warriors takımıdır. En başarısız takım ise Orlando Magic takımıdır.



Yukarıdaki grafikte ele alınan 7 sezonluk maçlarda takımların kazandıkları maçlar ev sahibi ve deplasman olacak şekilde oranlanmıştır. Bu grafik sonucunda ev sahibi olarak en çok kazanma oranına sahip olan takım "GSW" takımıdır. "ORL" takımı ise kazandığı maçları daha çok deplasmandayken kazanmıştır.



Random Forest algoritması ile tüm değişkenlerin değişken önemlilikleri hesaplanmıştır. Bir takımın maçı kazanmasına en çok etki eden 5 değişken ve oranları aşağıdaki gibidir:

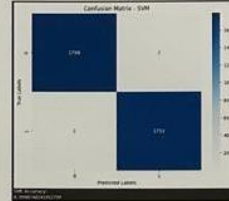
- ♦ drtg(defansif reyting): 0.353
- ♦ ortg(ofansif reyting): 0.252
- ♦ efg%(efektif atış yüzdesi): 0.109
- ♦ drb(defansif rebound sayısı): 0.064
- ♦ 3p%(3 sayı yüzdesi): 0.034

KAYNAKÇA

- [1] dataquestio. (2021). NBA Games Walkthrough: Get Data [Jupyter Notebook]. GitHub.
- [2] Weiner, J. (2019). Predicting the outcome of NBA games with machine learning. Towards Data Science
- [3] Wu, L., & Lee, Y. H. (2016). Modeling and forecasting the outcomes of NBA basketball games. Journal of Forecasting, 35(7), 646-661. doi: 10.1002/for.2411

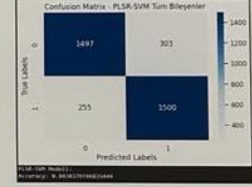
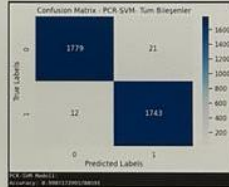
TAHMİN SONUÇLARI

TEMEL BİLEŞENLER İLE SVM TAHMİN MODELİ



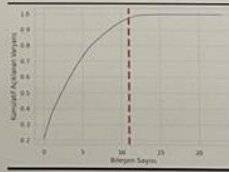
Support Vector Machine (SVM) modeli kullanılarak elde edilen sonuçlar, doğruluk oranı ile gösterilmiştir. Confusion matrix ve buna dayalı ısı haritası, modelin sınıflandırma performansını daha ayrıntılı olarak göstermektedir. Doğruluk oranı, y_test veri setindeki gerçek etiketlerle yapılan tahminlerin yüzdesini ifade etmektedir. Bu modelde veri setinde boyut indirgeme yapılmadan ham veri ile tahminleme yapılmıştır. Tahmin sonucunda 0.9988 accuracy değeri çıkmıştır

TÜM TEMEL VE KISMİ BİLEŞENLER KULLANILARAK SVM TAHMİN MODELİ

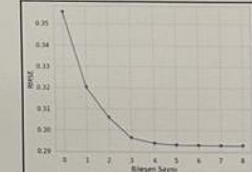


Yukarıdaki grafiklerde kurulan modellerin sonuçları confusion matrix olarak gösterilmektedir. Tüm bileşenler kullanıldığında PCR-SVM modelinde accuracy değeri 0.9907 çıkmıştır. PLSR-SVM modelinde ise accuracy değeri 0.8430 olarak elde edilmiştir.

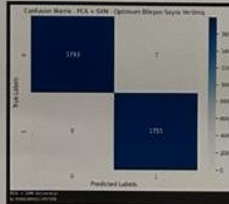
OPTIMUM BİLEŞEN SAYISI İLE SVM TAHMİN MODELİ



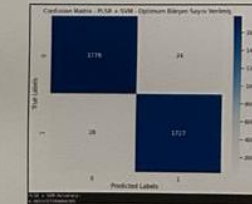
PCA ile veri modellemede kullanılmak üzere bileşen sayısı seçmek için grafik oluşturulmuştur. Sonuç olarak 11 bileşenin veri setinin 0.91 ile büyük bir kısmını açıkladığı bulunmuştur ve bileşen sayısı olarak "n_component" parametresinde 11 kullanılacaktır.



Tüm bileşenlerin RMSE değerleri hesaplanarak grafikte gösterilmiştir. Optimum bileşen sayısı seçilirken RMSE değeri en düşük olan bileşen seçilmektedir. Burada 9. Bileşen en düşük RMSE değerine sahip olduğu için seçilmiştir.



Optimum bileşen sayısı kullanıldığında PCR-SVM modelinde accuracy değeri 0.9980 çıkmıştır.



Optimum bileşen sayısı kullanıldığında PLSR-SVM modelinde accuracy değeri 0.9853 çıkmıştır.

SONUÇ

Bu çalışmada, boyut indirgeme yöntemleri SVM ile birleştirilerek model performansının değerlendirilmesi yapılmıştır. Optimum bileşen sayısının seçilmesi, daha güvenilir tahminler elde etmek amacıyla model performansını artırmada yardımcı olmuştur.