

Titanik Kazasının Veri Madenciliği Yöntemleri İle İncelenmesi

Investigation of Titanic Accident by Data Mining Methods

Fatih Afşin
Teknoloji Fakültesi Yazılım Mühendisliği / Fırat Üniversitesi
Elazığ,Türkiye

ÖZET

RMS(Royal Mail Ship) Titanik'in batması , tarihin en azılı gemi kazalarından biridir . 15 Nisan 1912'de ilk seferinde Titanik bir buz dağına çarparak battı ve 2224 yolcusundan 1514 yolcu ve mürettebat yaşamını yitirdi. Batması imkansız denilen bu gemide bu denli bir can kaybının olmasının sebeplerinden birisi ise yolcu ve mürettebat için yeterli sayıda cankurtaran olmamasıydı . Efsaneleşmiş bu trajik kaza uluslararası toplumu şoke etti ve gemiler için daha iyi güvenlik düzenlemelerine yol açtı. Hayatta kalmak için bazı şans unsurları bulunsa da,bazı insanların(kadın,çocuk ve üst sınıf gibi) diğer insanlara nazaran hayatta kalma olasılığı daha yüksekti. Günümüzde veri madenciliği yöntemlerinin gelişmesi ile özellikle hangi yolcuların enkazdan kurtulduğunu tahmin etmek için makine öğrenmesi araçlarından 5 tanesini (Desicion Tree , SVM , Random Forest , KNN,Logistic Regression) algoritmalarını uygulayarak elde ettiğimiz başarımlar incelenmiştir.

Investigation of Titanic Accident by Data Mining Method

ABSTRACT

The RMS(Royal Mail Ship) Titanic sank is one of the shipwecks at least in history . On the first time on April 15,1912, the Titanic slammed into an iceberg and sank and 2224 crew died.One of the reasons why this ship, which is called impossible to sink,caused such a loss of life was that there were not enough lifeguards fort he passengers and crew. This legendary tragic accident shocked the international community and led to better safety arrangements for ships. While there were some chansenes for survival,some people(such as women,children,and upper classes) were more likely to survive than others. Today,with the development of data minning methods,our performance has been examined by applying 5 of the machine learning tools(Desicion Tree,SVM, Random Forest,KNN,Logistic Regression) to predict which passengers escaped the wreckage.

1.GİRİŞ

Titanik , White Star Line şirketine ait Oympic sınıfı bir transatlantik yolcu gemisiydi. Harland and Wolff tersanelerinde üretilmiş yapımı tamamlandığında dünyanın en büyük buharlı yolcu gemisiydi . 15 Nisan 1912 gecesi daha ilk seferinde bir uz dağına çarpmış ve yaklaşık iki saat kırk dakika içinde Kuzey Atlantik'in buzlu sularına gömülmüştür.Batışı 1514 kişinin ölüyle sonuçlanmış ve en büyük deniz felaketlerinden biri olarak tarihe geçmiştir.[1] Bu çalışmamızın veri setini bizlere Kaggle sunuyor . Kaggle veri bilimcilerin ,veri setlerinde oluşan problemleri çözmeye çalıştıkları global bir yarışma platformudur. Kaggle'ın bizlere sunmuş olduğu titanik veri setini (VS) tanımakla başlayalım Eğitim setini (train.csv),makine öğrenim modelimizi oluşturmak için kullanacağız . Modelimiz , yolcuların cinsiyet ve sınıfı gibi "özelliklere" dayanacaktır. Survived : hayatta kalma (0=Hayır,1=Evet) , Pclass : bilet sınıfı (1=1., 2=2., 3=3.sınıf) , Sex : Cinsiyet , Sibsp : titanik'teki

kardeş/eşsayısı , Parch : Titanik'teki aynı aileye ait ebeveynlerin/çocukların sayısı , Ticket : bilet numarası , Fare : ödenen ücret , Cabin = kabin numarası , Embarked ise yolcuların bindikleri limanlardır . Bu çalışmada kullandığımız VS 'den test ve eğitim setleri oluşturup algoritmalarımız üzerindeki başarımlarını ölçüp algoritmalarımızı karşılaştıracaktır. Çalışmamız da sırasıyla VS üzerinde boş ,eksik (nul ,NaN) değerleri veri ön işleme adımlarından sınıflarına göre medyan değerlerini veya ortalama değerlerini uygulayarak veri setimizde kayıp verilerimizden kurtulup , verinin görselleştirilmesi , Kullanacağımız makine öğrenme algoritmalarının performanslarını arttırmak için Feature enginering (özellik mühendisliği) aşmasını yapıp modellerimizi eğitimini yapıp başarımlarını karşılaştırdık .

2.VERİ MADENCİLİĞİ (DATA MINING)

Veri madenciliği , büyük miktardaki verinin içinden geleceği tahmin edilmesinde yardımcı olacak

anlamli ve yararlı baęlantı ve kuralların veya tahminlerin bilgisayar programlarının aracılığıyla aranması ve analizidir . Ayrıca veri madencilięi , ok byk miktardaki verilerin iindeki iliřkileri inceleyerek aralarındaki baęlantıyı bulmaya yardımcı olan ve veri tabanı sistemleri ierisinde gizli kalmıř bilgilerin ekilmesini saęlayan veri analiz teknięidir. Bu iřlemlerin uygulama alanları olduka geniřtir . Bu alanlar ierisinde veri grsellięi , yapay sinir aęları , istatistik , yapay ğrenme vb. gibi disiplinler bulunmaktadır. Makine ğrenmesinde sayısal verilerin tahmin edilmesi iin Prediction Algorithms kullanılır . Sayısal olmayan yani kategorik verilerin tahmini iin ise Classification(Sınıflandırma) kullanılır.

2.1 Bazı Makine ğrenmesi Algoritmaları

Bu alıřmamızda kullandığımız 5 makine ğrenmesi algoritmalarını tanıyalım ;

Logistic Regression : Lojik regresyon sınıflandırma iřlemi yapmaya yarayan bir regresyon yntemidir .Kategorik veya sayısal verilerin sınıflandırılmasında kullanılır .Baęımlı deęiřkenin yani sonucun sadece 2 farklı deęer alabilmesi durumunda alıřır (Evet/Hayır , Erkek/Kadın, řiřman/Zayıf vs.) . [3]

SVM: (Support Vector Machine) Lojik regresyon ile benzer bir sınıflandırma algoritmasıdır . Her ikisinde iki sınıfı ayıran en iyi izgiyi bulmaya alıřırlar. Algoritma izlecek doęrunun iki sınıfında elemanlarına en uzak yerden geecek řekilde ayarlanmasını saęlar . Hibir parametre almayan bir sınıflayıcıdır. SVM aynı zamanda doęrusal ve doęrusal olmayan verileri de sınıflandırabilir ancak genellikle verileri doęrusal olarak sınıflandırmaya alıřır. [2]

Decision Tree: Karar aęacı algoritması , veri madencilięi sınıflandırma algoritmalarından biridir . nceden tanımlanmıř bir hedef deęiřkene sahiptir. Yapıları itibariyle en tepeden ařaęı inen bir strateji sunmaktadır. Yani basit karar verme adımları uygulanarak byk miktardaki kayıtları , ok kk kayıt gruplarına blerek kullanılan bir yapıdır .[4]

Random Forest : Rastgele ormanlar veya rastgele karar ormanları , sınıflandırma , regresyon ve dięer grevler iin , eęitim ařamasında ok sayıda karar aęacı oluřturarak problemin tipine gre sınıf veya sayı tahmini yapan bir toplu ğrenme yntemidir. [7]

KNN : K-En yakın komřu algoritması Thomas Cover tarafından sınıflandırma ve regresyon iin nerilen parametrik olmayan bir yntemdir. Her iki durumda da , girdi zellik alanındaki en yakın

rneklerinden oluřur. ıktı K-NN sınıflandırma veya regresyon iin kullanılıp kullanılmadıęına baęlıdır. [6]

2.2 Veri Madencilięinde Karřılařılan Problemler

Veri madencilięi uygulamalarında karřılařılabilecek problemler řunlardır.

Artık Veri: Artık veri problemde istenilen sonucu elde etmek iin kullanılan rneklem kmesindeki gereksiz niteliklerdir. Bu durum pek ok iřlem arasında karřımıza ıkabilir .

Boř Veri: Bir veri setinde boř deęer , zellikler arasında herhangi bir nitelięin deęeri olabilir . Boř deęer ,tanımı gereęi kendisi de dahil olmak zere hibir deęere eřit olmayan veridir.

Dinamik veri: Kurumsal evrim ii veri tabanları dinamiktir ve ierięi srekli olarak deęiřir. Bu durum, bilgi keřfi metotları iin nemli sakıncalar doęurmaktadır.

Eksik veri: Veri kmesinin byklęnden ya da doęasından kaynaklanmaktadır.

Eksik veriler olduęunda yapılması gerekenler řunlardır:

- Eksik veri ieren kayıt veya kayıtlar ıkarılabilir.
- Deęiřkenin ortalaması eksik verilerin yerine kullanılabilir.
- Var olan verilere dayalı olarak en uygun deęer kullanılabilir.

Grltl ve Kayıp Deęerler: Veri giriři veya veri toplanması esnasında oluřan sistem dıřı hatalara grlt denir. Byk veri tabanlarında pek ok nitelięin deęeri yanlış olabilir. Veri toplanması esnasında oluřan hatalara lmden kaynaklanan hatalar da dāhil olmaktadır. Bu hataların sonucu olarak birok nitelięin deęeri yanlış olabilir ve bu yanlışlardan dolayı veri madencilięi amacına tam olarak ulařmayabilir.

Sınırlı Bilgi: Veri tabanları genel olarak basit ğrenme iřlerini saęlayan zellik veya nitelikleri sunmak gibi veri madencilięi dıřındaki amalar iin hazırlanmıřlardır. Bu yzden, ğrenme grevini kolaylařtıracak bazı zellikler bulunmayabilir.[8]

3.Verit Madencilięi Ařamaları

Veri Madencilięi Yntemi ile Titanik Kazasının Arařtırılması veri madencilięi uygulaması ařaęıda ařamalar halinde verilmiřtir.

3.1 Verilerin Temizlenmesi

VS verilerinde bir sapma , anormal bir değer olup olmadığının tespiti için veri kalitesi incelenmiştir.

3.1.1 Artık Verilerin Temizlenmesi

Bu çalışmamızda modelimizi eğitmeden önce Passenger Id ve Cabin adlı feature'leri (Nitelik,özellik) veri setimizde modelimizin eğitimi sırasında hata ve sorun ile karşılaşmamak için silinmiştir

3.1.2 Boş Değerlerin Doldurulması

VS' de boş değer olarak Embarked yani yolcuların hangi limandan bindiğini gösteren nitelik içerisinde saptanmış olup bu durumdan kurtulmak için embarked değeri ile bilete ödenen fiyat arasında bir ilişki bulup boş değer saptanan kayıtlarda ödenen ücrete göre yolcunun embarked değeri girilmiştir.

3.1.3 Eksik Verilerin Tamamlanması

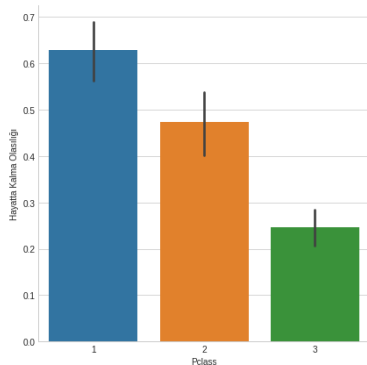
Bu çalışmamızda VS üzerinde bir'den fazla nitelik'te eksik veri saptanmış olup bu verileri VS değerlerini bozmamak için bazı niteliklerde ortalamasını bazı niteliklerde ise medyan ortanca değeri yazılmıştır.

3.2 Verilerin Görselleştirilmesi

VS de bulunan verileri grafik halinde sunulmuştur.

3.2.1 Sınıflara Göre Hayatta Kalma

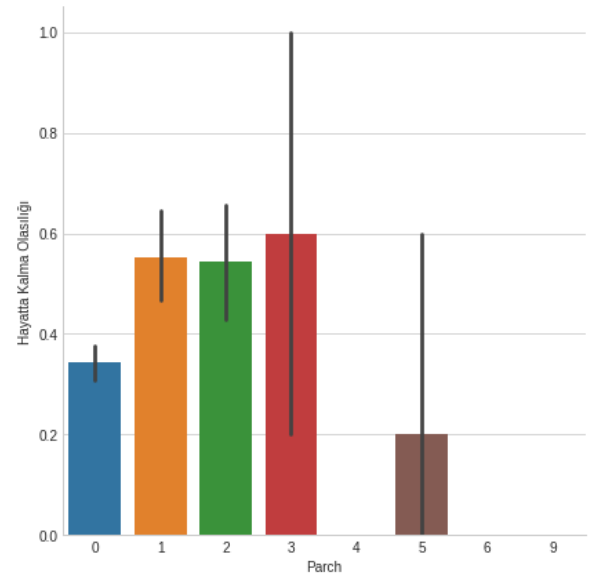
Pclass yolcularımızın seyahat ettiği sınıflardır Şekil-1 'de gördüğümüzde 1.sınıfta hayatta kalma oranı yüksek olduğudur. Burada değerlerimizin ve özelliklerimizin çok olmayışından ötürü modelimiz için yeni bir özellik oluşturulmasına gerek duyulmamıştır .



Şekil-1

3.2.2 Aile Sayısının Hayatta Kalma Oranları

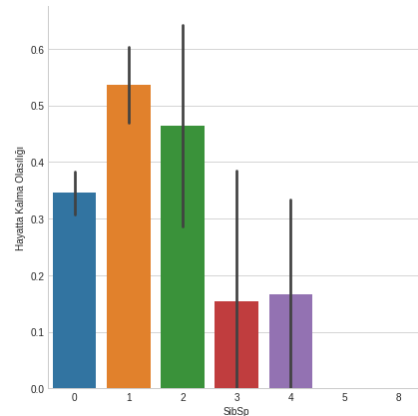
Şekil 2'deki grafikte bar plotların üstünde olan çizgiler şunu gösterir ; örnek olarak 3'ü ele alırsak ortalama olarak 3 bireye sahip ailelerin hayatta kalma olasılığı %60 dır ama ortadaki çizgi bize 3 kişiye sahip ailenin hayatta kalma oranını 0.2-1.0 aralığında olabileceğini de gösteriyor . Burada hayatta kalma olasılığı çok değişken ve büyük bir aralığa sahip .Model geliştirirken bunu göz önüne alıp veri setimizi Parch ile SibSp yi 3 ten aşağısı ve yukarısı olarak birleştirip standart sapmamızın daha tutarlı olmasını sağladık .



Şekil-2

3.2.3 Ebeveyn Çocuk Sayısına Göre Hayatta Kalma

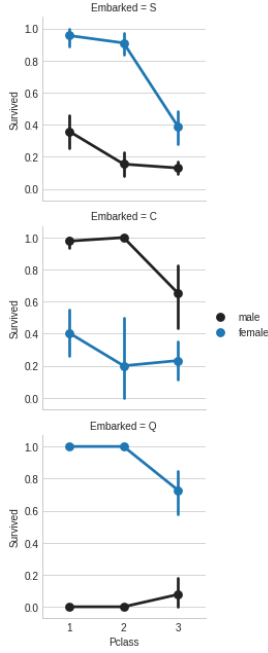
Şekil-3 te gördüğümüz eğer 2 den fazla SibSp değerine yani ebeyn çocuk sayısına sahip bir ailenin hayatta kalma olasılığı sert bir şekilde azalıyor. SibSp==0 veya 1 veya 2 olan grubun hayatta kalma oranları daha fazla. Bu kategoriye göre yeni bir feature(özellik) oluşturduk.



Şekil-3

3.2.4 Yolcuların Bindikleri Liman-Cinsiyet için Hayatta Kalma Oranları

Şekil-4 te görüldüğü gibi kadın yolcular erkek yolculara göre daha fazla hayatta kalmışlardır. Erkeklerin ise bindikleri limana göre C limanında binenlerin hayatta kalma olasılığı daha yüksek olduğu görülmüştür. Bu durum sınıflandırma yaparken doğrudan modelimizde kullanılmıştır.



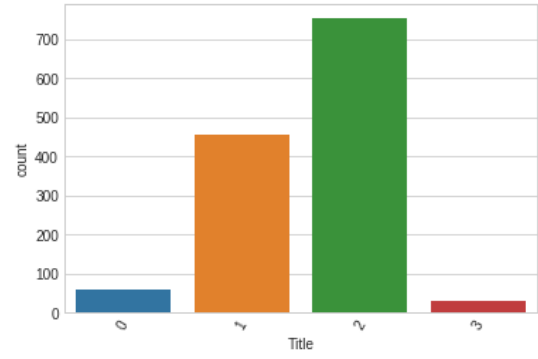
Şekil-4

3.3 Öz nitelik Seçimi

Makine öğreniminde öz nitelik çıkarımı büyük bir veri kümesini açıklamak için gereken kaynak miktarını azaltmayı içerir [9]. Veri setimizde 3 farklı öz nitelik değeri oluşturulmuştur.

3.3.1 Yolcu isimlerinin Başlıklara Bölünmesi

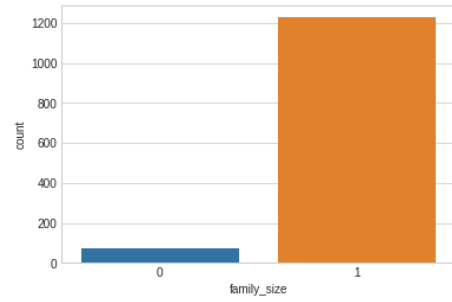
Şekil-5 teki gibi Name niteliğini “Master = 0”, “Miss, Ms, Mile, MRS = 1”, “Mr = 2”, “Mrs = 3” olarak 4 başlığa indirgeyip title adlı yeni bir öz nitelik oluşturulmuştur.



Şekil-5

3.3.2 Aile Niteliğinin İndirgenmesi

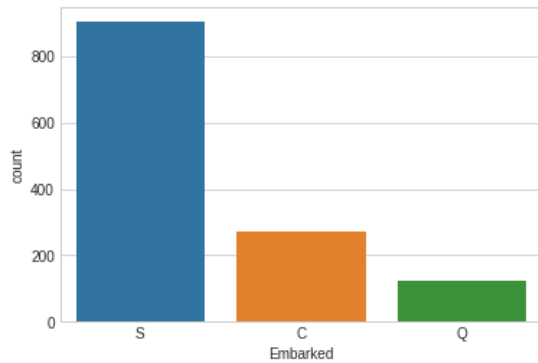
Aile öz niteliği incelendiğinde aile sayısının 5 ten büyük olduğu durumlarda hayatta kalma oranlarının sert bir şekilde düşmesinden ötürü bu niteliği aile sayısı 5 ten küçük ve büyük olmak üzere Şekil-6 'da olduğu gibi family_size adında yeni bir öz niteliğe dönüştürülmüştür.



Şekil-6

3.3.3 Liman Niteliğinin İndirgenmesi

Yolcuların bindikleri limanlara göre Şekil-7 'de olduğu gibi Embarked adı altında “S”, “C”, “Q” formatına dönüştürülmüştür.



Şekil -7

4. Model'in Oluşturulması

Öncelikle modelimizden iki farklı test veri setini oluşturup bu veri setlerimizi makine öğrenmesi algoritmalarımızda teker teker başarımlarını test edilmiştir.

4.1 Test Veri Setinin Oluşturulması

Modelimizi oluşturmak için hayatta kalma niteliğini baz alarak rastgele bir şekilde veri setimizi %33 oranında test veri setlerine böldük .

```
X_train 590
X_test 291
y_train 590
y_test 291
test 418
```

Şekil-8

4.2 Makine Öğrenmesi Algoritmalarının Uygulanması

Sırası ile Desicion Tree , SVC, Random Forest, Logistic Regression,K-NN algoritmaları uygulandı . Bu modellerin içerisinde bulunan parametrelerin en iyi parametresini (Hyper Paramater) Grid Search yöntemine göre arayıp , bulduğumuz parametrelerin en iyi değerlerini karşılaştırırken Cross Validation yöntemini uyguladık.

```
Fitting 10 folds for each of 250 candidates, totalling 2500 fits

[Parallel(n_jobs=1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=1)]: Done 88 tasks | elapsed: 1.6s
[Parallel(n_jobs=1)]: Done 2500 out of 2500 | elapsed: 6.8s finished
[Parallel(n_jobs=1)]: Using backend LokyBackend with 4 concurrent workers.

0.8355932283389831
Fitting 10 folds for each of 28 candidates, totalling 280 fits

[Parallel(n_jobs=1)]: Done 120 tasks | elapsed: 1.7s
[Parallel(n_jobs=1)]: Done 280 out of 280 | elapsed: 4.3s finished
[Parallel(n_jobs=1)]: Using backend LokyBackend with 4 concurrent workers.

0.7983850847457627
Fitting 10 folds for each of 54 candidates, totalling 540 fits

[Parallel(n_jobs=1)]: Done 42 tasks | elapsed: 5.4s
[Parallel(n_jobs=1)]: Done 192 tasks | elapsed: 20.6s
[Parallel(n_jobs=1)]: Done 442 tasks | elapsed: 50.5s
[Parallel(n_jobs=1)]: Done 540 out of 540 | elapsed: 1.0min finished

0.840677966101695
Fitting 10 folds for each of 14 candidates, totalling 140 fits

[Parallel(n_jobs=1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=1)]: Done 123 tasks | elapsed: 0.8s
[Parallel(n_jobs=1)]: Done 140 out of 140 | elapsed: 1.0s finished
/opt/conda/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:940: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)
[Parallel(n_jobs=1)]: Using backend LokyBackend with 4 concurrent workers.

0.8283389838508474
Fitting 10 folds for each of 40 candidates, totalling 400 fits

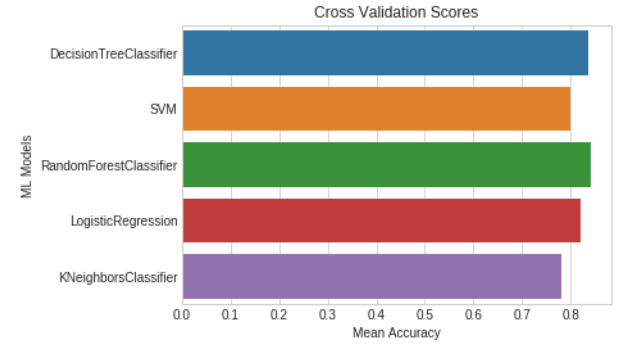
[Parallel(n_jobs=1)]: Done 280 tasks | elapsed: 1.2s

0.7796610169491525

[Parallel(n_jobs=1)]: Done 400 out of 400 | elapsed: 1.6s finished
```

Şekil-9

Kullandığımız modellerdeki başarımların oranlarının Cross Validation skorlarını tablo üzerinde gösterimi Şekil-10 verilmiştir.



Şekil-10

4.Sonuç

Titanik veri setimiz üzerinde görselleştirme ,veri ön işleme ve yeni öznitelikleri oluşturulmasından sonra 5 farklı makine öğrenimi modeli eğitmeye başladık . Bu modellerin içerisinde bulunan parametrelerin en iyi parametresini (Hyper Paramater) Grid Search yöntemine göre arayıp , bulduğumuz parametrelerin en iyi değerlerini karşılaştırırken Cross Validation yöntemini uyguladık ve sonuç olarak 0.840677966101695 yani %84 başarımlarında Random Forest algoritması ile başarımlar oranı en yüksek algoritmamızın Random Forest olduğunu gördük.

Veri görselleştirme kısmında özellikle cinsiyeti kadın olan ve üst sınıf insanların hayatta kalma oranlarının yüksek olduğunu ve üst sınıfta olan insanların aynı limandan bindiklerini görmüş olduk. Model başarımlarını yükseltebilmek için daha kapsamlı öznitelikler çıkartılıp , gürültülü özellikleri tanımlayıp bu verileri kaldırarak başarımlar oranı iyileştirilebilir.

KAYNAKLAR (REFERENCES)

- [1]Wiki,RMS Titanic,Wikipedia
- [2]Kadir Ulgen ,(2016)Makine Öğrenmesinde Karar Ağaçları, Medium
- [3]Ekrem Hatipoğlu , (2018)logisctic regression ,part-8,Medium
- [4] Ekrem Hatipoğlu, (2018)machine-learning, classification ,supportvector ,Medium
- [5]Wiki , Rastgele Orman, Wkipedia.Org
- [6]Wiki , K-nearest_neighbors_algorithm ,Wikipedia.Org
- [7]Kadir Alan ,(2020) veri madenciliği yöntemleri , Medium
- [8]Levent Sabah,Hüseyin Bayraktar ,(2020) Veri madenciliği birliktelik kuralları yöntemi kullanılarak binaların risk durumlarının belirlenmesi,Dergi Park
- [9]Wiki,öznitelik çıkarımı,Wikipedia