

# Sınıflandırma için Makine Öğrenimi

Hüseyin Cem ARAS – İbrahim Berkay ÖZ

Danışman: Doktor Öğretim Üyesi Özgül VUPA ÇİLENGİROĞLU  
Çalışma Ortağı: SANEM Plastik Tasarım Merkezi Müdürü Alim Fatih KILINÇ



## AMACIMIZ

İzmir ili Torbalı ilçesinde plastik masa örtüsü üretimi gerçekleştiren SANEM PLASTİK adlı fabrikadan alınan 2018-2019 yılları arasındaki veriler ile makine öğrenimi algoritmalarını kullanarak model oluşturulup, modele göre fabrikanın verimliliğini arttırmaktır.

## OEE

Overall Equipment Effectiveness” kısaltmasıdır. Türkçe çevirisi Toplam Ekipman Etkinliği’dir. OEE bütün ekipmanların ne ölçüde kullanıldığına işaret eden bir TPM(total productive maintenance) hesabıdır. Arızalar, ekipman ayarları, duruşlar, çalışma hızındaki azalmalar, iskartalar ve yeniden işlem gibi kayıplar üzerine düşer.

OEE [%] = Kullanılabilirlik Oranı x Performans Oranı x Kalite Oranı

Kullanılabilirlik Oranı [%]: Ekipmana ait sebeplerden (arıza, ayarlamadan kaynaklanan duruş süresi vs.) kaynaklanan kullanılabilirlik miktarını gösterir.

Performans Oranı [%]: Çalışma hızlarında tasarımla belirlenmiş hızlara göre düşüşleri ve birkaç saniyelik duruşları hesap eder.

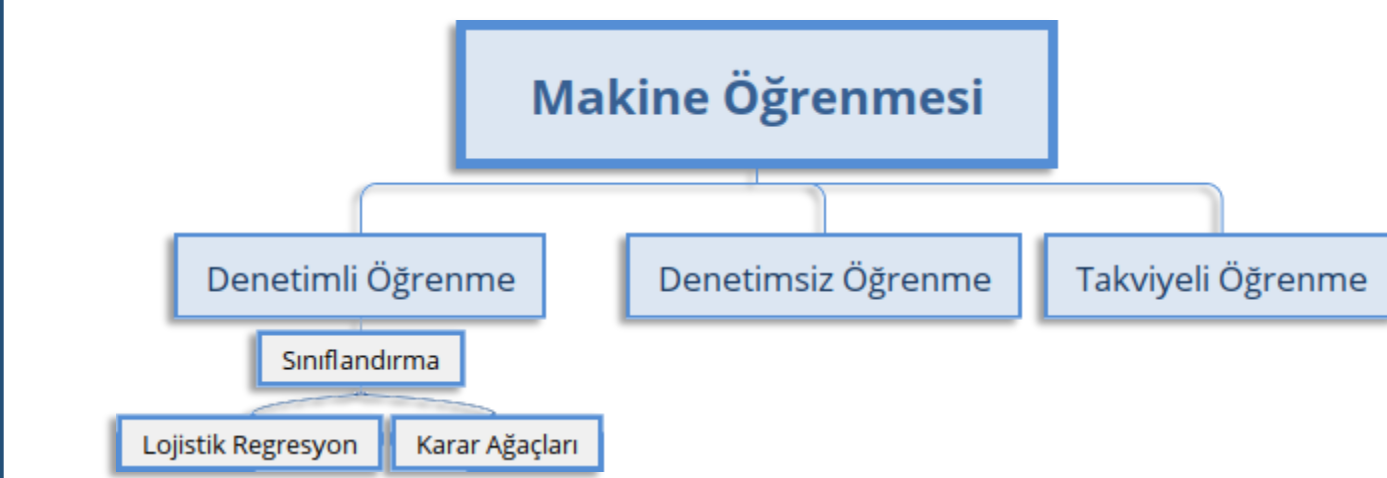
Kalite Oranı [%]: Toplam işlenen parçaların iskarta ve yeniden işlem kayıplarının yüzdesidir.

## VERİ SETİ

İzmir ili Torbalı ilçesinde plastik masa örtüsü üretimi gerçekleştiren SANEM PLASTİK adlı fabrikadan alınan 2018-2019 yılları arasındaki veriler (n=198, train verisi=163) ile makine öğrenimi algoritmalarını kullanarak model oluşturulup, modele göre fabrikanın verimliliğini arttırmaktır.

OEE, üretim, saat, kayıp metre, değişkenleri sürekli, hafta ve gün değişkenleri kategorik değişkenlerdir. Model için sürekli değişkenlerimiz, medyana göre bölünüp O&1 olarak kategorik değişken haline getirilmiştir. Modellerimizde OEE değişkeni bağımlı değişken iken diğer değişkenler bağımsız olarak ele alınmıştır. Modellerin amacı bağımsız değişkenlerin, OEE puanına nasıl etkileyeceğini araştırmaktır.

## KULLANILAN YÖNTEMLER



### Lojistik Regresyon:

Normallik varsayımının bozulması nedeniyle doğrusal regresyon analizine alternatif olmaktadır. Doğrusal regresyon analizinde bağımlı değişkenin değeri, lojistik regresyon analizinde ise bağımlı değişkenin alabileceği değerlerden birinin gerçekleşme olasılığı tahmin edilmektedir. Temel olarak lojistik regresyonda bağımsız değişkenler ile iki ya da çok sınıflı kategorik bağımlı değişken arasındaki ilişkinin tanımlanması için matematiksel modelleme yapmak amaçlanmaktadır.

$$P_i = \frac{\exp(\sum_{k=0}^p \beta_k x_{ik})}{1 + \exp(\sum_{k=0}^p \beta_k x_{ik})}$$

### CART Algoritması:

CART algoritması, ağaç yapısına dayalı olarak sınıflandırma ve regresyon modellerinin türetilmesi için yaygın olarak kullanılan bir istatistiksel prosedürdür. CART ağaç modeli, tek değişkenli ikili kararların bir hiyerarşisini içerir.

CART ağaçları, kesin bir heterojenlik (impurity) ölçüsüne bağlı olarak düğümlere ayrılmış iki değerli (binary) ağaçlardır ve bu nedenle de sonuçta homojen dallar oluşmaktadır. Ağacın hedefi benzer veya aynı çıktı değerlerine sahip olma eğiliminde olan alt gruplar yaratmaktır. CART modelleri için bölünmelerin bulunmasında kullanılan dört farklı heterojenlik ölçüsü mevcuttur. Kategorik hedef değişkenler için Gini. Twoing veya (sıralayıcı hedef değişkenleri için) sıralı Twoing sürekli hedef değişkenler için ise en küçük kareli sapma kullanılabilir.

$$g(t) = 1 - \sum_j p^2(j/t)$$

## KAYNAKLAR

- [https://lean.org.tr/oe-nedir/, Can YÜKSELEN]
- Atalay M., Çelik E., 2017, Artificial Intelligence And Machine Learning Applications in Big Data Analysis
- Özkaya A., 2012, Makine Öğrenmesi ile Ürün Sınıflandırma İncelemesi
- KALAYCI E., 2018, Comparison of machine learning techniques for classification of phishing web sites
- Abbas G., Shirali M., Moloud V., Amal S., 2018, Predicting the outcome of occupational accidents by CART and CHAID methods at a steel factory in Iran
- Yan-yan S., Ying L., 2015, Decision tree methods: applications for classification and prediction
- Antipov E., Pokryshevskaya E., 2009, Applying CHAID for logistic regression diagnostics and classification accuracy improvement

## GITHUB

Projemize ait dosyalara Github üzerinden erişebilirsiniz:  
<https://github.com/hcemaras/Siniflandirma-icin-Makine-Ogrenmesi>

## LOJİSTİK REGRESYON

### Full model

```

> full.model = glm(OEE ~. - tecrube ,family=binomial, data=train)
> summary(full.model)

Call:
glm(formula = OEE ~. - tecrube, family = binomial, data = train)

Deviance Residuals:
    Min       1Q   Median       3Q      Max 
-2.7868  -0.4047   0.2040   0.6068   1.4250 

Coefficients:
(Intercept)  -0.14138    0.54440   -0.260    0.79509 
uretim1       3.90538    0.80857    4.830 1.37e-06 *** 
saat1        -2.21958    0.86006   -2.581  0.00986 ** 
kayip_metre1 -1.90456    1.00137   -1.902  0.05718 . 
gun2         0.09838    0.73493    0.134  0.89351 
gun3        -0.20531    0.75528   -0.272  0.78575 
gun4        -0.03542    0.73268   -0.048  0.96144 
gun5         1.07417    0.82050    1.309  0.19048 
gun6         1.56849    0.87277    1.797  0.07231 . 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 214.49  on 162  degrees of freedom
Residual deviance: 120.85  on 154  degrees of freedom
AIC: 138.85

Number of Fisher Scoring iterations: 6
  
```

Veri setimizdeki tüm veriler ile kurulan lojistik regresyon modelidir. görüldüğü üzere üretim, saat anlamlı çıkmıştır. Kayıp metre değişkeni 0.05 sınırında olduğundan ‘.’ ile işaretlenmiştir.

Yapılan Korelasyon analizine göre kayıp metre ve saat değişkeni ile ilişkili bulunmuştur (0.48). Bu yüzden kayıp metrenin modele alınmasına karar verilmiştir.

En iyi modeli oluşturmak için stepwise yöntemi kullanılmıştır.

### En iyi model

```

> full.model = glm(OEE ~. - tecrube -kayip_metre ,family=binomial, data=train)
> best.model <- stepAIC(full.model, direction = "both",trace = FALSE)
> summary(best.model)

Call:
glm(formula = OEE ~ uretim + saat, family = binomial, data = train)

Deviance Residuals:
    Min       1Q   Median       3Q      Max 
-2.5993  -0.2757   0.2635   0.6622   1.2583 

Coefficients:
(Intercept)  0.2805    0.3149    0.891    0.373 
uretim1      3.0630    0.6392    4.792 1.65e-06 *** 
saat1       -3.5317    0.6397   -5.521 3.38e-08 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 214.49  on 162  degrees of freedom
Residual deviance: 131.67  on 160  degrees of freedom
AIC: 137.67
  
```

Bu durumda en iyi modelde üretim, saat değişkenleri %95 anlamlılık düzeyinde istatistiksel olarak anlamlı çıkmıştır. Daha sonra OR, CI ve Wald testine bakılmıştır.

Wald test: Chi-squared test:  
X2 = 31.2, df = 2, P(> X2) = 1.6e-07

Kategorik değişkenlerle kurulan lojistik regresyon en iyi modelindeki Wald değerlerine bakıldığında, tüm p değerleri 0.05’den küçük olduğu için katsayılar %95 güven seviyesinde istatistiksel olarak anlamlı bulunmuştur.

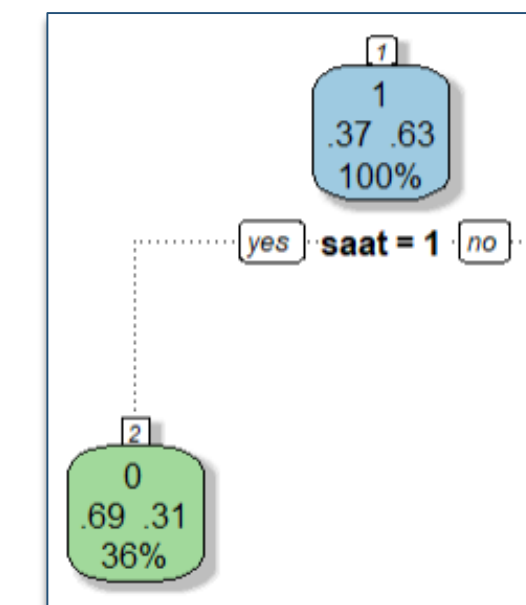
Ayrıca OR ve güven aralıklarına bakıldığında, OR’a ait güven aralıklarının 1’i içermediği görülmektedir. Bu yüzden OR’ler anlamlıdır ve yorumlanabilir.

Üretimi yüksek olanların, üretimi düşük olanlara göre Toplam Ekipman Etkinliğinin (OEE) yüzdesinin yüksek olma şansı yaklaşık olarak 21 katıdır. Saat değişkeninde ise saati az olanların, yüksek olanlara göre OEE yüzdesinin yüksek olma şansı ise 1/0.29 oranından yaklaşık olarak 3 katı olduğu elde edilmektedir.

## CART ALGORİTMASI

### En iyi model

CART ağacımız Gini katsayısına göre hesaplanıp, ilk olarak saat kategorik değişkenine ayrılmıştır. Bu durumda OEE’nin yorumlanmasında en önemli değişken olarak saat değişkeni olduğu söylenebilir.



Saati ortalamadan az olanlar ise, sağ tarafa ayrılmıştır. 163 birimin %64 oranla 104 birim bu düğüme düşmüştür. 104 birimin %18’i yani 19 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 85 tanesi ise OEE’si 1 olan sınıfta yer almıştır.

## YENİ MODEL

Fabrika verileri ile kurulan lojistik regresyon modelinde saat ve üretim değişkenleri modelde anlamlı çıkmıştır. CART modelinde ise en önemli değişkenin saat olduğu tespit edilmiştir. Daha sonra CART modelinde ağaç derinliği kontrol edildiğinde 2. düğümün üretim değişkenine ayrıldığı gözlemlenmiştir. Fakat üretim değişkeni tek başına yeterli ayrıştırmayı yapamadığından ağaçta kırılmıştır.

Ancak modelde değişken sayısını artırıp modelin nasıl değiştiğini görmek için verilerimize simülasyon çalışması ile makineyi çalıştıran kişinin tecrübesi (10 yıldan az&10 yıl ve daha fazla olarak belirlenmiş kategorik değişken) eklenmiştir. Bu yeni değişkene göre modeller kurulmuştur.

## LOJİSTİK REGRESYON

### Tecrübe, Üretim, Saat (TÜS)

Burada en uygun modelde tecrübe, üretim, saat değişkenleri %95 anlamlılık düzeyinde istatistiksel olarak anlamlı çıkmıştır.

Daha sonra OR, CI ve Wald testine bakılmıştır.

```

> model2 = glm(OEE ~uretim+ saat + tecrube ,family=binomial, data=train)
> summary(model2)

Call:
glm(formula = OEE ~ uretim + saat + tecrube, family = binomial, data = train)

Deviance Residuals:
    Min       1Q   Median       3Q      Max 
-2.1667  -0.3807   0.1177   0.5443   1.9852 

Coefficients:
(Intercept)  -1.2441    0.4895   -2.541    0.011 * 
uretim1       3.4907    0.7087    4.926 8.41e-07 *** 
saat1        -4.0670    0.7303   -5.569 2.56e-08 *** 
tecrube1      2.7232    0.5543    4.913 8.98e-07 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 214.49  on 162  degrees of freedom
Residual deviance: 100.22  on 159  degrees of freedom
AIC: 108.22
  
```

Wald test: Chi-squared test:  
X2 = 24.3, df = 2, P(> X2) = 5.4e-06

Kategorik değişkenlerle kurulan en uygun lojistik regresyon modelindeki Wald değerlerine bakıldığında, tüm p değerleri 0.05’den küçük olduğu için katsayılar %95 güven seviyesinde istatistiksel olarak anlamlı bulunmuştur.

Ayrıca OR ve güven aralıklarına bakıldığında, OR’a ait güven aralıklarının 1’i içermediği görülmektedir. Bu yüzden OR’ler anlamlıdır ve yorumlanabilir.

(Intercept) uretim1 saat1 tecrube1  
0.28821447 32.81030772 0.01712815 15.22833475

10 Yıllık Tecrübesi olanların, olmayanlara göre OEE yüzdesinin yüksek olma şansı yaklaşık olarak 15 kattır.

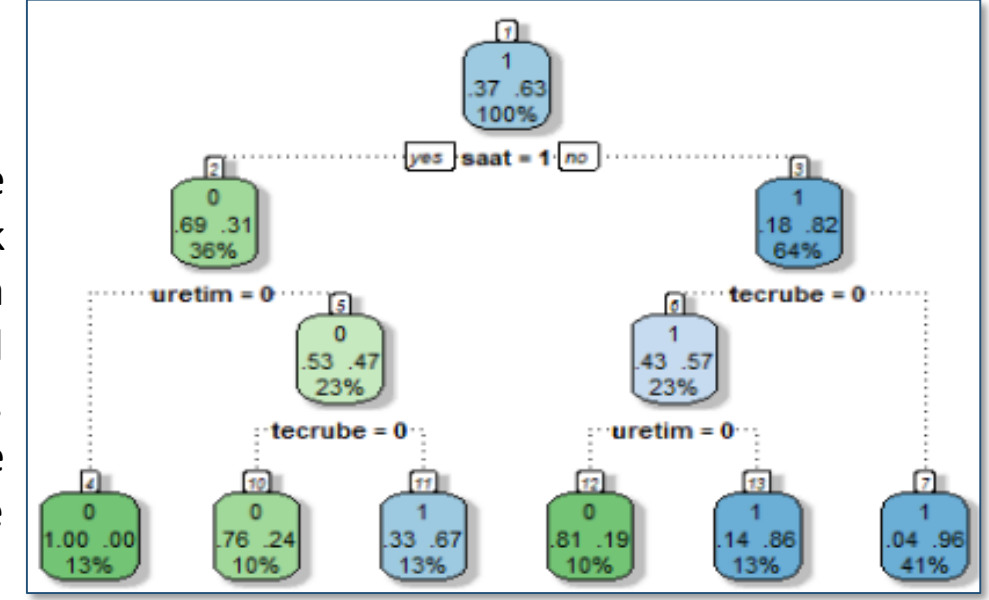
Üretimi yüksek olanların, üretimi düşük olanlara göre OEE yüzdesinin yüksek olma şansı yaklaşık olarak 33 kattır.

Saati az olanların saati yüksek olanlara göre OEE yüzdesinin yüksek olma şansı ise 1/0.17 yaklaşık olarak 6 kattır.

## CART ALGORİTMASI

### Tecrübe, Üretim, Saat (TÜS)

CART ağacımız Gini katsayısına göre hesaplanıp, ilk olarak saat kategorik değişkeni ile ayrılmıştır. Saati ortalama veya ortalamadan yüksek olanlar, sol tarafa ayrılıp, 2. düğüme yer almıştır. 163 birimlik “train” verisi üzerinde hesaplandığı için bu 2. düğüme %36 oranla 59 birim düşmüştür.



2. düğüm daha sonra üretim değişkenine göre ikiye ayrılmıştır. Üretimi ortalama ya da ortalamadan yüksek olanlar sağ tarafa ayrılıp 5. düğümü oluşturmuşlardır. 5. düğüme toplam verinin %23’üne denk gelen 38 birim düşmüştür. Bu 38 birimin %53’ü yani 20 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 18 tanesi ise OEE’si 1 olan sınıfta yer almıştır. 2. düğüme üretimi ortalamadan az olanlar sol tarafa ayrılıp 4. düğümü oluşturmuşlardır. 4. düğüme toplam verinin %13’üne denk gelen 21 birim düşmüştür. Bu 21 birimin %100’ü yani tamamı OEE’si 0 olan sınıfta yer almıştır. OEE’si 1 olan sınıfta yer alan birim yoktur. 5. düğüm daha sonra tecrübe değişkenine göre ikiye ayrılmıştır. On yıllık tecrübesi olanlar sol tarafa ayrılıp 10. düğümü, olmayanlar sağ tarafa ayrılıp 11. düğümü oluşturmuşlardır. 10. düğüme toplam verinin %10’una denk gelen 16 birim düşmüştür. Bu 16 birimin %76’sı yani 12 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 4 tanesi ise OEE’si 1 olan sınıfta yer almıştır. 11. düğüme ise toplam verinin %13’üne denk gelen 21 birim düşmüştür. Bu 21 birimin %33’ü yani 7 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 14 tanesi ise OEE’si 1 olan sınıfta yer almıştır. Saati ortalamadan az olanlar ise, sağ tarafa ayrılmıştır. 163 birimin %64 oranla 104 birim bu düğüme yani 3. düğüme düşmüştür. 104 birimin %18’i yani 19 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 85 tanesi ise OEE’si 1 olan sınıfta yer almıştır. 3. düğüm daha sonra tecrübe değişkenine göre ayrılmıştır. On yıllık tecrübesi olmayanlar sol tarafa ayrılıp 6. düğümü oluşturmuştur. 6. düğüme toplam verinin %23’üne denk gelen 38 birim düşmüştür. Bu 38 birimin %43’ü yani 16 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 22 tanesi ise OEE’si 1 olan sınıfta yer almıştır. 3. düğümden on yıllık tecrübesi olanlar sağ tarafa ayrılıp 7. düğümü oluşturmuşlardır. 7. düğüme toplam verinin %41’ine denk gelen 67 birim düşmüştür. Bu 67 birimin %4’ü yani 3 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 64 tanesi ise OEE’si 1 olan sınıfta yer almıştır. 6. düğüm daha sonra üretim değişkenine göre ikiye ayrılmıştır. Üretimi az olanlar sol tarafa ayrılıp 12. düğümü, çok olanlar sağ tarafa ayrılıp 13. düğümü oluşturmuşlardır. 12. düğüme toplam verinin %10’una denk gelen 16 birim düşmüştür. Bu 16 birimin %81’i yani 13 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 3 tanesi ise OEE’si 1 olan sınıfta yer almıştır. 13. düğüme ise toplam verinin %13’üne denk gelen 21 birim düşmüştür. Bu 21 birimin %14’ü yani 3 tanesi OEE’si 0 olan sınıfta yer almıştır. Kalan 18 tanesi ise OEE’si 1 olan sınıfta yer almıştır.

## SONUÇ

İzmir ili Torbalı ilçesinde plastik masa örtüsü üretimi gerçekleştiren SANEM Plastik adlı fabrikadan 2018-2019 yılları arasında alınan 198 veri ile makine öğrenimi algoritmalarını kullanarak saat ve üretim ve değişkenlerinin OEE puanı üzerindeki etkileri, lojistik regresyon ve CART algoritmaları üzerinden modeller kurulmuştur. Kurulan lojistik regresyon modelinde saat ve üretim değişkenleri modelde anlamlı çıkmıştır. CART modelinde ise en önemli değişkenin saat olduğu tespit edilmiştir.

Modellerimizin iki büyük kısıtlaması vardır. Biri değişken sayısının diğeri ise örneklem sayısının azlığıdır. Bu amaçla modellemeye değişken sayısını artırıp modelin nasıl değiştiğini görmek için verilerimize simülasyon çalışması ile makineyi çalıştıran kişinin tecrübesi (10 yıldan az & 10 yıl ve daha fazla olarak belirlenmiş kategorik değişken) eklenmiştir. Son eklenen değişken ile tecrübe, üretim, ve saat değişkenleri lojistik regresyon modelimizde anlamlı, CART modelinde ise yine en önemli değişkenin saat olduğu, daha sonra tecrübe en son olarak da üretim olduğu bulunmuştur.

Bu çalışma ile üretim yapan fabrikaların verisinde makine öğrenimi kapsamında yer alan sınıflandırma modellerinden lojistik regresyon ve CART modellerinin uygulanabilirliği tespit edilmiştir. Ayrıca bu tip verilerde değişik sınıflandırma algoritmaları da kullanılabilir ve algoritmaların performansları da karşılaştırılabilir.

Fabrikadaki uzmanlarla yapılan ortak çalışmalarla bu tespitin farklı veri türlerinde ve farklı modellerle yapılması planlanmalıdır. Böylece sanayi-üniversite işbirliği gerçekleşmesi sağlanabilir.