

Banka müşterisi kayıp analizi

Abdullah Jamous

Siber Güvenlik ve Kriptografi Bölümü, Yıldız Teknik Üniversitesi, Türkiye

Özet İçinde bulunduğumuz çağda, veri analizi ve makine öğrenmesi tüm kuruluşlar için en önemli önceliktir. Nesnelerin İnterneti (IoT), büyük veri ve bulut bilgi işlem gibi İnternet teknolojilerinin hızla gelişmesiyle birlikte, bireyler, hükümet yetkilileri ve ordu, veri analizinde sorunlarla karşı karşıyadır. Üstel veri büyüme hızı göz önüne alındığında, araştırmacılar için zorlu bir görev, akıllı şehirler tasarlarken büyük miktarda veriyi etkin bir şekilde nasıl yönetecekleridir. Önerilen sonuçlar, %60 ile %80 arasında değişen sonuç doğruluk değerleri ile yüksek düzeyde bir belirsizlik elde etti.

Anahtar Kelimeler Makine Öğrenmesi, Naive Bayes, k-en yakın komşuluk, karar ağaçları.

I. Giriş

Kimin kalacağını ve gelecekte kimleri kaybedebileceğimizi tahmin edebilmek için müşteri kaybı analizi yapmak çok önemlidir. Makine öğrenimi kullanılarak yapılan veri analizi, geri çekilmeye maruz kalan ve katkıda bulunmaya çalışan müşterilerin örüntüsünün tahmin edilmesine katkıda bulunmuştur. Bu müşterilerle ilgilenmek ve karşılaştıkları sorunları bilmek ve Bu sorunlardan kurtulmak için bankalar, bu analizleri yaptıktan sonra, mümkün olan en fazla sayıda müşterinin kendilerinde kalmasını sağlamak için performans verimliliğini ve iş kalitesini artırabilir.

II. Veri Kümesi

Veri setini kaggle.com'dan aldık, veri setinde 10.000 satır ve 12 sütun var, sütunlar arasında sütunun kendisi ile ilgili verileri içerdiği için sildiğimiz bir sütun var, o da credit_id ve ayrıca iki sütunumuz var. kategorik veri içerenleri x algoritmasını kullanarak sayısal verilere çevireceğiz, ayrıca eksik veri var mı kontrol edeceğiz, eksik veri olan satırları sileceğiz, gerekirse veriyi indireceğiz

III. Sınıflandırma Modelleri

Bu deneyde kullanılan yöntemler kısaca bahsetmek istiyorum:

I. k-en yakın komşuluk

KNN algoritması sınıflandırılmak istenen bir veriyi daha önceki verilerle olan yakınlık ilişkisine göre sınıflandıran bir algoritmadır. Örneğin; $k = 3$ olarak alırsak, yeni gelen verinin eski verilere olan uzaklıkları ölçülür ve en yakın 3 tanesi belirlenir. K sayısı genelde 3 ile 7 arasında seçilir. K sayısı tek sayı olması gerekiyor. $K=5$ seçilirse $k=3$ socundan daha başarılı çıkıyor.

II. Naive Bayes

Naïve Bayes sınıflandırması olasılık ilkelerine göre tanımlanmış bir dizi hesaplama ile, sisteme sunulan verilerin sınıfını yani kategorisini tespit etmeyi amaçlar. Bir eleman için her durumun olasılığını hesaplar ve olasılık değeri en yüksek olana göre sınıflandırır. Az bir eğitim verisiyle çok başarılı işler çıkartabilir.

III.karar ağaçları (decision trees)

Amaç veri öğeleri basit kurallar kuralları bu kuralları öğrenerek bir değişkenin değerini tahmin eden modeli oluşturmaktır. Algoritma eksik değerleri desteklemez makineyi eğitmeden önce eksik değerleri (eksik değer) hesaplamamız gerekir. Rastgele Orman görünümü karar ağacı (Karar ağacı) gibi hem sınıflandırma (Sınıflandırma) hem de regresyon (Regresyon) problemlerinde kullanılabilir. Çalışma mantığı birden fazla karar ağacı oluşturur. Bir sonuç üreteceği zaman bu karar ağaçlarındaki ortalama değer alınır ve sonuç üretilir.

IV. Deneysel Analiz

Veri kümesinde yapılan ön işlem : Eğitimde kullanılmayan sütunları silmek ve 1000 satır veri seçmek

▼ Delete unused columns and take part of the data

```
[4] del df['customer_id']

remove_n = 9000
drop_indices = np.random.choice(df.index, remove_n, replace=False)
df = df.drop(drop_indices)
```

Kategorik veri içeren iki sütunumuz var, bu veriyi sayısal veriye dönüştürmek için one-hot-encoding algoritmasını kullandım.

```

on-hot-encoding algorithm used for category columns

[5]: exchanger_rates.head(10)
exchanger_rates.head(10)
exchanger_rates.tail(10)
exchanger_rates.info()
exchanger_rates.describe()

Out[5]:
country_code  country_name  country_currency  country_flag  gender  gender_age  age  balance  products_number  credit_card  action  source  estimated_salary  share
0  99  990  1.0  0.0  0.0  0.0  1.0  0.0  2  103594.00  1  1  1  78904.12  0
1  99  990  0.0  1.0  0.0  0.0  1.0  0.0  3  70349.40  1  2  0  17007.24  0
2  99  470  1.0  0.0  0.0  0.0  1.0  0.0  4  17180.00  0  0  0  3940.22  0
3  99  470  0.0  0.0  0.0  1.0  1.0  0.0  4  17180.00  0  0  0  3940.22  0
4  91  900  0.0  1.0  0.0  0.0  1.0  0.0  3  146604.07  2  0  0  84624.57  0
5  91  900  0.0  1.0  0.0  0.0  1.0  0.0  3  146604.07  2  0  0  84624.57  0

```

Eksik veri olup olmadığını da kontrol ettim.

▼ Check for missing data

```
[3] df.isnull().values.any()
```

False

Ondan sonra datayı parçalara ayırdım, 400 satır eğitmek için, 300 satır doğrulamak için, 300 satır test yapmak için

```

# Division of samples

X_train, y_train, X_valid, y_valid, X_test, y_test = train_valid_test_split(X, target = 'churn', train_size=0.4, valid_size=0.1, test_size=0.5)

print(X_train.shape), print(y_train.shape)
print(X_valid.shape), print(y_valid.shape)
print(X_test.shape), print(y_test.shape)

(4000, 11)
(4000,)
(1600, 11)
(1600,)
(4000, 11)
(4000,)
(None, None)

```

İkinci aşama doğrulama aşamasıdır:
Knn algoritması
doğrulama aşaması

[illegible]

Test aşaması

K-Nearest Neighbors (KNN) with test data

```
[12] y_pred_3 = knn3.predict(X_test)
     y_pred_5 = knn5.predict(X_test)
```

```
from sklearn.metrics import accuracy_score
print("Accuracy with k=3 :", accuracy_score(y_test, y_pred_3)*100)
print("Accuracy with k=5 :", accuracy_score(y_test, y_pred_5)*100)
```

Accuracy with k=3 : 67.33333333333333
Accuracy with k=5 : 70.0

Önemli gördüğüm şey k değeri ne kadar büyük olsa o kadar iyi başarı oranı alabiliriz.

Naive Bayes algoritması:

Doğrulama aşaması

```
[14]: from sklearn.model_selection import train_test_split
      from sklearn.naive_bayes import GaussianNB

      gnb = GaussianNB()
      y_pred = gnb.fit(X_train, y_train).predict(X_val)
      print('Number of mislabeled points out of a total %d points: %d' % (X_test.shape[0], (y_val != y_pred).sum()))
      print('Accuracy : ', (X_test.shape[0] - (y_val != y_pred).sum())/(X_test.shape[0])*100)

      Number of mislabeled points out of a total 300 points: 64
      Accuracy : 78.66666666666666
```

Test aşaması

[illegible]

Decision trees

Doğrulama aşaması

[illegible]

Test aşaması

```

gccnu9cl:  _e"333333333333333
    bu7ur("gccnu9cl:  _w6rl7c2"9ccnu9cl"zco6(λ"7e2r' λ"bu6q),799)
[13] λ"bu6q = c7t'bu6q7c7(λ"7e2r)
q6c7j0i0n 7r6e2 w777 7e2f d979

```

Gördüğüm kadarıyla öğretmiş olduğum algoritmalarından en iyi sonuç veren Naive Bayes algoritmasıdır.

IV. Öğrenme modelinizin başarısın

Başarı oranı 60% lardan 80%lere kadar.