

Yapay Zeka - odev 2

Parsa Kazerooni - 19011915 @ YTU - Prof.Dr. Mehmet Fatih Amasyalı

Bu proje, YTU Yapay Zeka dersi için hazırlanmıştır. Web Scraping ile Hollanda'nın en büyük supermarket'i olan [Albert Heijn](#) üzerinden yemek ürünleri ve beslenme detayları veri olarak çekilmiştir. Sonra supervised learning algoritmaları ile bu veriler üzerinden ne kadar sağlıklı bir yemek olduğunu sınıflandırıcı modeller tasarlanmıştır. Sağlık ölçümü, Fransa'dan kaynaklanan [Nutri-Score](#) yöntemidir. Sonra bu modellerin performansları karşılaştırılmıştır.

Kurulum

requirements.txt dosyası içerisindeki kütüphanelerin kurulumu için:

```
pip install -r requirements.txt
```

1. Scraping

[scraping](#) klasörü içerisindeki [populate_urls.py](#) dosyası, farklı ürün kategorilerinden ürünlerin url'lerini çekmektedir. Bu url'ler daha sonra [populate_products.py](#) dosyasında ürünlerin detaylarını çekmek için kullanılmaktadır. Bu dosyaların çekmiş olduğu veriler [data](#) klasöründe bulunmaktadır.

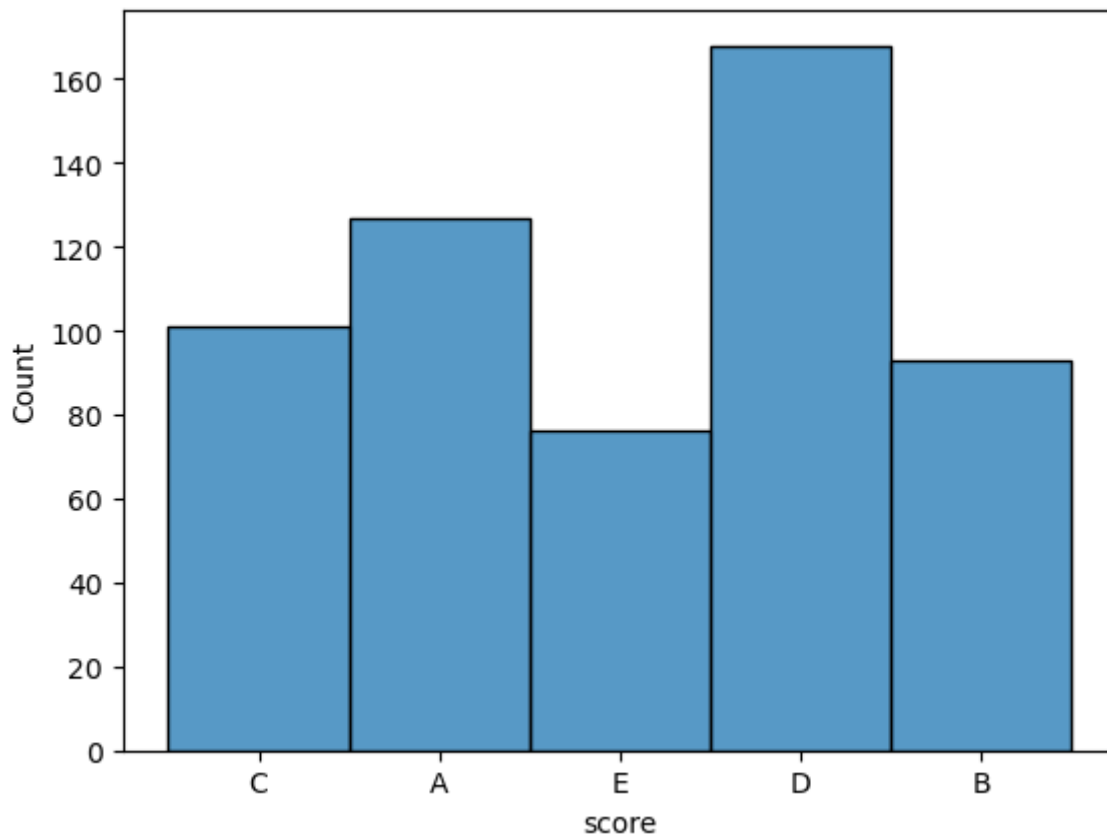
```
python scraping/populate_urls.py  
python scraping/populate_products.py
```

2. Data

Automation işlemi, DDoS saldırısı gibi algılanmaması için veriler daha küçük parçalar ve aralıklar halinde çekilmiştir. Bu yüzden verilerin birleştirilmesi gerekmektedir. [data](#) klasöründe bulunan verileri birleştirmek için [merge_data.py](#) dosyası kullanılmaktadır. Bu dosya, verileri birleştirip [data/ah.csv](#) dosyasına kaydetmektedir.

```
python data/merge_csvs.py
```

Yaklaşık 600 farklı ürünün, 20 farklı beslenme özelliği bulunmaktadır. Bu özelliklerin bir kısmı çok az sayıda bulunmakta olduğu için, sonradan analize dahil edilmemiştir. Bu yüzden, 20 özellikten 10 tanesi seçilmiştir.

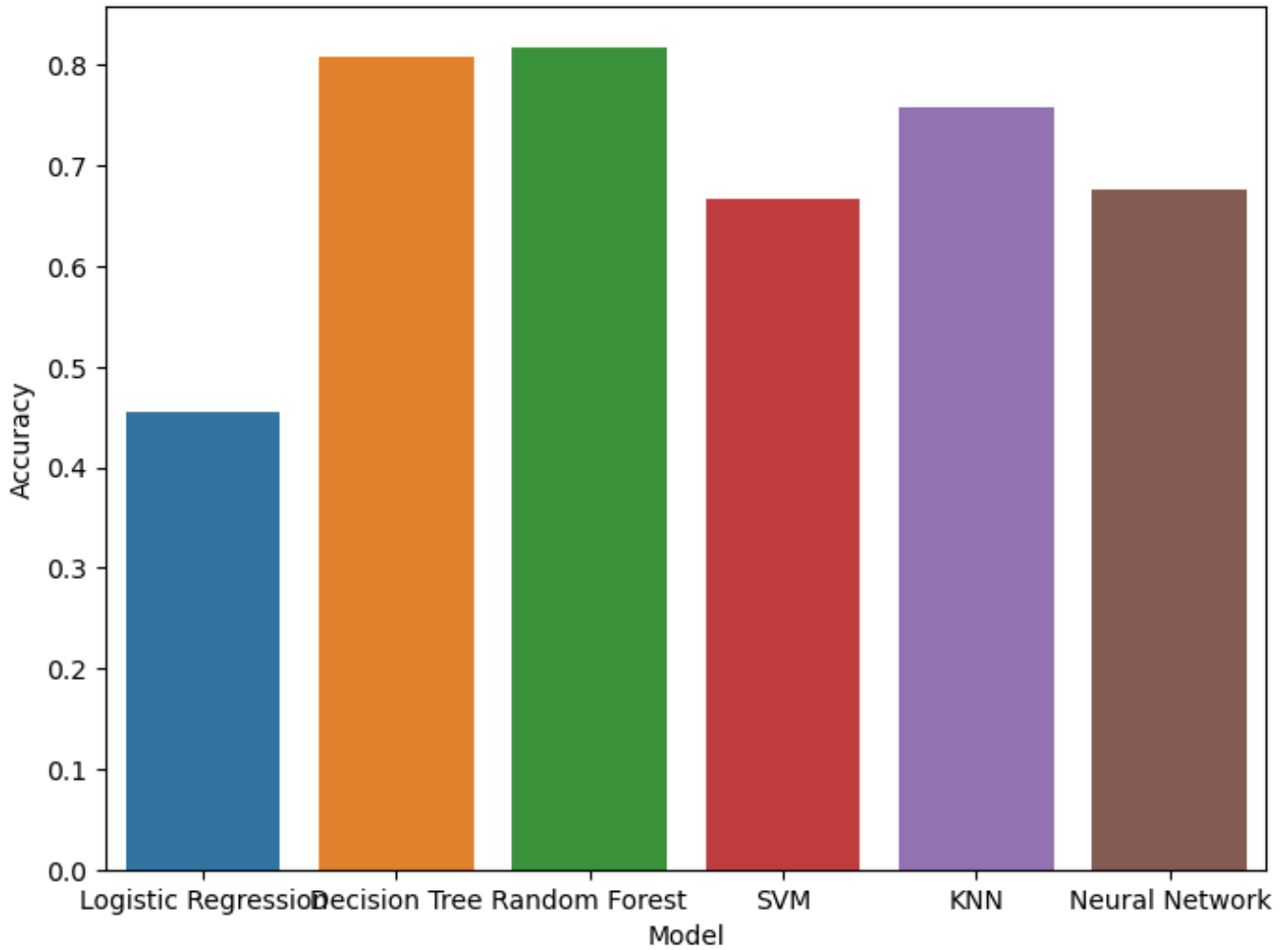


3. Notebook

[notebook.ipynb](#) dosyasi, verilerin islenmesi ve algoritmaların eğitilmesi ve karşılaştırılmasını içerir.

4. Sonuclar

6 tane farklı sınıflandırıcı model, Logistic Regression, Decision Tree, Random Forest, KNN, SVM ve Neural Network, kullanılmıştır. Bu modellerin performansları karşılaştırılmıştır. Sonuçlar aşağıdaki gibidir:



Daha iyi performans saglamak icin, Normalization ve feature selection uygulanmistir.

10-fold cross validation ile modellerin performanslari karsilastirilmistir. Her modelin score gecmisleri kaydedilmistir. ve ona gore t-testi yapilmistir. Sonuclar asagidaki gibidir:

```
Logistic Regression vs Decision Tree: t=-10.98, p=0.00
Logistic Regression vs Random Forest: t=-11.64, p=0.00
Logistic Regression vs SVM: t=-6.20, p=0.00
Logistic Regression vs KNN: t=-8.82, p=0.00
Logistic Regression vs Neural Network: t=-6.35, p=0.00
Decision Tree vs Random Forest: t=-0.29, p=0.77
Decision Tree vs SVM: t=4.24, p=0.00
Decision Tree vs KNN: t=1.53, p=0.14
Decision Tree vs Neural Network: t=3.90, p=0.00
Random Forest vs SVM: t=4.65, p=0.00
Random Forest vs KNN: t=1.85, p=0.08
Random Forest vs Neural Network: t=4.28, p=0.00
SVM vs KNN: t=-2.55, p=0.02
SVM vs Neural Network: t=-0.25, p=0.80
KNN vs Neural Network: t=2.26, p=0.04
```

References

1. <https://www.ah.nl/>
2. <https://en.wikipedia.org/wiki/Nutri-Score>
3. <https://machinelearningmastery.com/k-fold-cross-validation/>