Kötü amaçlı URL'leri Algılayan Sistem Tasarımı

Hazırlayanlar

17110131006-Mehmet OKYAY

18110131505-Abdourazak ALİ EGUEH

Projenin Amacı

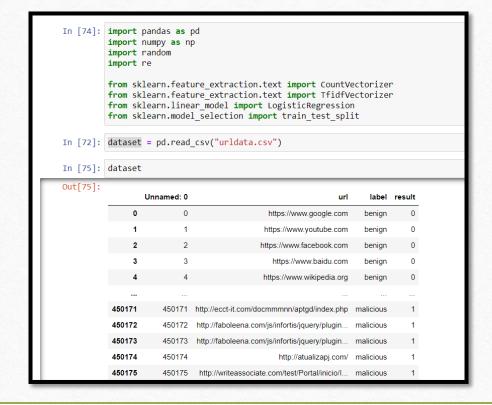
• Projemiz bilgisayar ve ağ güvenliğinde herhangi bir internet sitesinin güvenilir olup olmadığını tahmin eden bir makine öğrenmesi uygulamasıdır. Bu uygulama ile veri setinde bulunan güvenilir veya güvenilmeyen siteler makine öğrenmesi yolu ile eğitilip yeni listede olmayan yeni bir domain adresinin güvenilip güvenilemeyeceği konusunda fikir sahibi olmamızı sağlamaktadır.

Projenin Hazırlanma Aşamaları

- Veri Toplama
- Veri önişleme
- Verinin eğitilmesi
- Verinin test Edilmesi
- Eğitimin Değerlendirilmesi ve Geliştirilmesi

Veriyi Toplama

 Projemizde kullanılan veri seti Kaggle üzerinden indirilip kullanılmıştır.
 Veri setimiz Unnamed, url, label, result olmak üzere 4 niteliğe sahip ve 450176 adet kayıttan oluşmaktadır.



Veri Önişleme

- Veri setindeki label alanındaki değerleri result değerlerine göre yeniden etiketlendi.
- Veri setinde null değere sahip bir kayıtın olup olmadığına bakılarak eğitim modelinde hata oranının artmasına neden olabilecek değerler tespit edildi.

In [66]: dataset.replace({'label':{'malicious':'kötü_amaçli', 'benign':'iyi_amaçli'}},inplace=True)

```
Unnamed: 0 url label result

0 False False False False
1 False False False False
2 False False False False
3 False False False False
4 False False False False
4 False False False False
5 False False False False
6 False False False False
7 False False False False
8 False False False False
9 False False False False
1 False False False False False
1 False False False False False
1 False False False False False False
1 False False False False False False
1 False False False False False False False False
2 False ```

## Veri Önişleme

- Web sitelerinin domain adresleri ile güvenilir veya güvenilmez site ilişkisini kurabilmek için url adreslerindeki geçen kelimeleri gruplarını ayırmak gerekir. Yani url de geçen kelime ile label arasında bir bağlantı kurulacaktır. Bunun için Pythonda bir fonksiyon tanımladık.
- Ayrıca bu bağlantıyı kurmak için re(regular expression) kütüphanesi tanımlandı.

```
def makeTokens(f):
 tkns_BySlash = str(f.encode('utf-8')).split('/')—*# slah belirteçları yap.
 total_Tokens = []
 for i in tkns_BySlash:
 tokens = str(i).split('-')-*# - ile bölündüğünde / yapma işlemi
 tkns_ByDot = []
 for j in range(0,len(tokens)):
 temp_Tokens = str(tokens[j]).split('.')*# noktayla böldükten sonra belitrteçler yapar.
 tkns_ByDot = tkns_ByDot + temp_Tokens
 total_Tokens = total_Tokens + tokens + tkns_ByDot
 total_Tokens = list(set(total_Tokens))-*#gereksiz belirteçler kaldırılıyor.
 if 'com' in total_Tokens:
 total_Tokens.remove('com')-*#.com uzantılar birçok kez kaldırıldıği için direk kaldırıldı.
 return total_Tokens
```

## Eğitim Modeli Oluşturma

• Bağımlı(y) ve bağımsız değişkenler(X) veri setinden ilgili sütunları seçilerek alındı ve daha sonra metini ayıklama fonksiyonu TF x IDF Skorlama Modeline ile kelimelerin sıklık ve başkınlık oranları ile bağımlı değişken arasındaki ilişki kurmak amacıyla modelin eğitimi lojistik regresyon algoritması ile gerçekleştirildi.

```
y = dataset["label"]

url_list = dataset["url"]

vectorizer = TfidfVectorizer(tokenizer=makeTokens)

X = vectorizer.fit_transform(url_list)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

logit = LogisticRegression()

logit.fit(X_train, y_train)

LogisticRegression()
```

#### Projenin Test Edilmesi

• Veri setinin eğitilmesinden sonra tahmin işlemleri için random olarak belirlenen web sitelerinin test edilmesi aşamasıdır.

```
X predict = ["https://bm.ksu.edu.tr/", "https://www.kervaniho.meb.k12.tr/",
"https://www.w3schools.com/python/numpy/default.asp","https://www.udemy.com/"]
X predict
['https://bm.ksu.edu.tr/',
 'https://www.kervaniho.meb.k12.tr/',
 'https://www.w3schools.com/python/numpy/default.asp',
 'https://www.udemy.com/']
X predict = vectorizer.transform(X predict)
New predict = logit.predict(X predict)
print(New predict)
['iyi_amaçli' 'iyi_amaçli' 'iyi_amaçli' 'iyi_amaçli']
X predict1 = ["www.buyfakebillsonlinee.blogspot.com",
"www.unitedairlineslogistics.com",
"www.stonehousedelivery.com",
"www.silkroadmeds-onlinepharmacy.com"]
X predict1 = vectorizer.transform(X predict1)
New predict1 = logit.predict(X predict1)
print(New predict1)
['kötü amaçli' 'kötü amaçli' 'kötü amaçli']
```

## Projenin Değerlendirilmesi

• Yaptığımız projenin %99'un üzerinde bir başarı ile web sitelerinin domain adreslerine bakarak bu sitenin güvenli veya güvenli olmayan şeklinde tahminde bulunması bilgisayar ve ağ güvenliğinde bilgi güvenliği açısından yorumlandığında gizlilik ilkesine uygudur

```
print("Başarı: ",logit.score(X_test, y_test))
Başarı: 0.9948020791683326
```