

Zaawansowane metody uczenia maszynowego

Wykład 1

Mateusz Serocki

Politechnika Gdańska

October 18, 2023

Agenda

Kilka słów o mnie

- Aktualnie starszy inżynier uczenia maszynowego w NIKE
- Main skillset: Python/AWS
- 5.5 roku doświadczenia zawodowego
- Main topics: Regresja, Klasyfikacja, Prognoza, MLOps
- LinkedIn: <https://www.linkedin.com/in/mateuszserocki/>
- Email: Mateusz.Serockiii@gmail.com

Redukcja wymiarowości

Definicja

Redukcja wymiarowości to transformacja danych z przestrzeni wielowymiarowej do przestrzeni niskowymiarowej, tak aby reprezentacja niskowymiarowa zachowała pewne znaczące właściwości oryginalnych danych, idealnie zbliżone do ich wewnętrznego wymiaru.

Interpretacja słowna

Innymi słowami staramy się zmniejszyć wymiarowość naszego zbioru danych przy równoczesnym zachowaniu maksymalnej ilości informacji jaka z tego zbioru pochodzi.

Przykład redukcji wymiarów

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	0.0
1	4.9	3.0	1.4	0.2	0.0
2	4.7	3.2	1.3	0.2	0.0
3	4.6	3.1	1.5	0.2	0.0
4	5.0	3.6	1.4	0.2	0.0
...
145	6.7	3.0	5.2	2.3	2.0
146	6.3	2.5	5.0	1.9	2.0
147	6.5	3.0	5.2	2.0	2.0
148	6.2	3.4	5.4	2.3	2.0
149	5.9	3.0	5.1	1.8	2.0

Figure: Full dataset

	feature1	feature2	target
0	5.80	5.20	0.0
1	5.60	5.00	0.0
2	5.35	4.80	0.0
3	5.35	4.70	0.0
4	5.70	5.10	0.0
...
145	9.30	7.85	2.0
146	8.80	7.25	2.0
147	9.10	7.50	2.0
148	8.90	7.35	2.0
149	8.45	6.80	2.0

150 rows × 3 columns

Figure: Reduced dataset

Wady redukcji wymiarów

Redukcja wymiarów ma też negatywne skutki takie jak:

- Zmniejszenie ilości informacji
- Ryzyko usunięcia potencjalnie informatywnej zmiennej
- Całkowita lub częściowa utrata interpretowalności zmiennych
- Kolejny krok potrzebny do przygotowania danych

Zalety redukcji wymiarów

Pozytywne aspekty redukcji wymiarów to:

- Redukcja szumu pochodzącego ze zmiennych
- Brak zmiennych które nie mają uzasadnionego wpływu na wynik
- Mniejsza moc obliczeniowa potrzebna do utworzenia modelu
- Możemy zrezygnować z pobierania niektórych zmiennych jeżeli na tym etapie sadzimy że są nieistotne

Wariancja

Teoria

W teorii prawdopodobieństwa i statystyce wariancja jest kwadratem odchylenia od średniej zmiennej losowej. Wariancje często definiuje się także jako kwadrat odchylenia standardowego.

Zastosowanie

W praktyce wariancja odzwierciedla zmienność danej cechy, niska wariancja może nieść za sobą niską informatywność. Warto rozważyć usunięcie takiej cechy jeżeli nasz zbiór jest bardzo duży. Koniecznie musimy sprawdzić wpływ tego usunięcia na wynik modelu.

Korelacja

Typy korelacji

- Pearson - standardowa korelacja liniowa
- Spearman - korelacja rankingowa
- Kendal - korelacja rankingowa dla grup z powtarzającymi się wartościami
- more...

Przykład liczenia korelacji Pearsona

[Link](#)

Teoria dot. korelacji spearmana

[Link](#)

Teoria dot. korelacji Kendalla

[Link](#)

Biblioteka python do analizy zbioru danych

Pandas profiler

Biblioteka służy do szybkiej analizy danych, pomaga nam w zauważeniu potencjalnych problemów wewnątrz naszych danych, oraz może służyć do wstępnej analizy naszego zbioru

```
import pandas as pd
data = pd.read_csv('path_to_your_data')
profile = data.profile_report(
    title = "Pandas Profiling Report",
    correlations = {
        "pearson": {"calculate": True},
        "spearman": {"calculate": True},
        "kendall": {"calculate": True}
    },
)
profile.to_file("pandas_profiled.html")
```

Braki danych

Sposoby na radzenie sobie z brakiem danych

- Usuniecie
- Inputacja
- Podstawienie

Przykład 1

Przykładem algorytmu który radzi sobie z brakami danych jest XGBoost który wykorzystuje średnia wartość danej zmiennej zamiast wartości NULL (brak)

Przykład 2

Przykładem algorytmu który nie radzi sobie z brakami danych jest regresja liniowa (i pochodne) które w przypadku napotkania wartości null zwróca błąd processowania

Przykład analizy braku danych

```
import sklearn.datasets
import pandas as pd
import numpy as np

iris = sklearn.datasets.load_iris(as_frame=True)
data = pd.DataFrame(data= np.c_[iris['data'],
                                iris['target']],
                    columns= iris['feature_names'] + ['target'])
data.isnull().sum()
```

Usuwanie zmiennych

Code example

Przykład usunięcia kolumny za pomocą python - pandas

```
import pandas as pd
```

```
data = pd.read_csv('path_to_your_data')
```

```
new_data = data.drop('column_name', axis=1)
```

```
# or
```

```
data.drop('column_name', axis=1, inplace=True)
```

Standaryzacja zmiennych

Zapamiętaj

Każdy algorytm który działa na podstawie liczenia odległości między różnymi zmiennymi wymaga od nas wykonania standaryzacji

Typy standaryzacji

- MinMaxScaler
- StandardScaler

Przykład

Na wykładzie.

Pytanie z *

Czy w przypadku gdy wszystkie nasze cechy posiadają rozkład binarny, standaryzacja jest wymagana?

Principal component analysis

Analiza głównych składowych (PCA) to popularna technika analizy dużych zbiorów danych zawierających dużą liczbę wymiarów/cech na obserwacje, zwiększająca interpretowalność danych przy jednoczesnym zachowaniu maksymalnej ilości informacji i umożliwiającą wizualizację danych wielowymiarowych

Przykład redukcji wymiarów

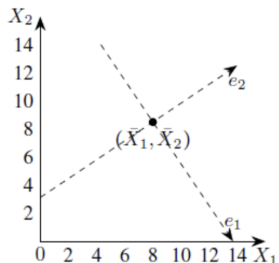


Figure: Rzut wektorów własnych

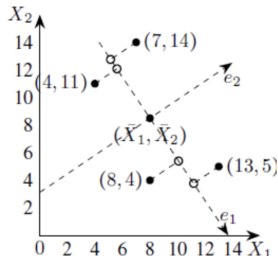


Figure: Mapowanie na podstawie wektorów własnych

RandomForest

Opis algorytmu

Losowe lasy lub losowe lasy decyzyjne to metoda uczenia się zespołowego do klasyfikacji, regresji i innych zadań, która polega na konstruowaniu wielu drzew decyzyjnych w czasie szkolenia.

Zastosowanie

- Regresja
- Klasyfikacja

Dokładne działanie algorytmu

Na wykładzie.

RandomForest

Code example

```
import sklearn.datasets
import sklearn.ensemble as ens
import pandas as pd

iris = sklearn.datasets.load_iris(as_frame=True)
X_clf = pd.DataFrame(data=iris.data, columns=iris.feature_name)
y_clf = iris.target

model_clf = ens.RandomForestClassifier()
model_clf.fit(X=X_clf, y=y_clf)
model_clf.predict(X_clf)

boston = sklearn.datasets.load_boston()
X_reg = pd.DataFrame(data=boston.data, columns=boston.feature_name)
y_reg = boston.target

model_reg = ens.RandomForestRegressor()
model_reg.fit(X=X_reg, y=y_reg)
model_reg.predict(X_reg)
```

Least Absolute Shrinkage and Selection Operator

Opis algorytmu

LASSO jest metoda analizy regresji, która przeprowadza zarówno selekcję zmiennych oraz ich regularyzację w celu zwiększenia dokładności przewidywania i interpretowalności modelu statystycznego, który generuje.

Historia

Została wprowadzona przez Roberta Tibshiraniego w 1996r. Warto w tym miejscu zaznaczyć, że w naszej pracy korzystamy ze zmodyfikowanej metody LASSO. Pierwotnie autor użył metody najmniejszych kwadratów do estymacji parametrów $\hat{\beta}$, natomiast my wykorzystujemy logarytm największej wiarygodności. Metoda ta została zmodyfikowana w 2006 r.

LASSO

Niech \mathbf{X} oznacza macierz predyktorów, a \mathbf{y} wektor odpowiedzi. Dla danego parametru t współczynniki $\hat{\beta}^{\text{LASSO}} = (\hat{\beta}_1, \dots, \hat{\beta}_k)$ są rozwiązaniem poniższej optymalizacji z ograniczeniem:

$$\min\{-L(\mu; \mathbf{y})\} \quad \text{dla} \quad \sum_{j=1}^k |\beta_j| \leq t \quad (1)$$

gdzie:

– L jest funkcja największej wiarygodności.

Wzór możemy zapisać w równoważnej formie Lagrange'a:

$$\min\left\{-L(\mu; \mathbf{y}) + \lambda \cdot \sum_{j=1}^k \|\beta_j\|\right\} \quad (2)$$