

Akademia Ekonomiczno-Humanistyczna w Warszawie

SPRAWOZDANIE

INTELIGENTNA ANALIZA DANYCH

LAB3
ANALIZA TEKSTÓW,
PRZETRENOWANIE I OPTIMALIZACJA MODELU

28.11.2021

JOANNA PRAJZENDANC
36358

MIŁOSZ SAKOWSKI
36381

Spis treści

1. Cel i przebieg ćwiczenia.....	3
2. Definicje i założenia.....	3
2.1. Wyjaśnienie pojęć.....	3
3. Grid search.....	3
3.1. Zadanie #1.....	3
i. Treść polecenia.....	3
ii. Rozwiązanie.....	3
4. Przetrenowanie.....	11
4.1. Zadanie #2.....	11
i. Treść polecenia.....	11
ii. Rozwiązanie.....	11
iii. Wnioski.....	13
5. Analiza tekstu.....	14
5.1. Zadanie #3.....	14
i. Treść polecenia.....	14
ii. Rozwiązanie.....	14
5.2. Zadanie #4.....	18
i. Treść polecenia.....	18
ii. Rozwiązanie.....	18
5.3. Zadanie #5.....	19
i. Treść polecenia.....	19
ii. Rozwiązanie.....	19
5.4. Zadanie #6.....	21
i. Treść polecenia.....	21
ii. Rozwiązanie.....	21

1. Cel i przebieg ćwiczenia

Celem ćwiczenia było utrwalenie wiedzy w zakresie optymalizacji modelu (grid search), zapoznanie się z pojęciem przetrenowania modelu oraz sposobów analizy tekstów na podstawie gotowych wiadomości e-mail.

2. Definicje i założenia

2.1. Wyjaśnienie pojęć

W sprawozdaniu pojawiają się następujące pojęcia:

- ✧ grid search - sposób optymalizacji modelu polegający na zestawieniu ze sobą wszystkich możliwych konfiguracji zadanych parametrów
- ✧ przetrenowanie - stan modelu, który mimo idealnego dopasowania dla danych ze zbioru trenującego dopasowanie do danych ze zbioru testującego jest dużo niższe i niewystarczające

3. Grid search

3.1. Zadanie #1

i. Treść polecenia

Proszę pobrać dowolny zbiór danych ze strony <https://archive.ics.uci.edu/ml/index.php>

Następnie proszę zoptymalizować las metodą Grid Search (trzy parametry modelu mają zostać zoptymalizowane).

Jakie wartości parametrów są najlepsze i dlaczego?

ii. Rozwiązanie

A. Zbiór danych

```
import pandas as pd

fires_origin = pd.read_csv(
    "forestfires.csv", header=0, index_col=False)
d = {
    'X': fires_origin['X'].values,
    'Y': fires_origin['Y'].values,
    'temp': fires_origin['temp'].values,
    'RH': fires_origin['RH'].values,
    'wind': fires_origin['wind'].values,
    'rain': fires_origin['rain'].values,
    'area': fires_origin['area'].values
}
fires = pd.DataFrame(d)
print(fires.head())
```

	X	Y	temp	RH	wind	rain	area
0	7	5	8.2	51	6.7	0.0	0.0
1	7	4	18.0	33	0.9	0.0	0.0
2	7	4	14.6	33	1.3	0.0	0.0
3	8	6	8.3	97	4.0	0.2	0.0
4	8	6	11.4	99	1.8	0.0	0.0

Obraz 1: Ładowanie bazy danych z pliku

```

import statistics as stat
maximum = max(fires['temp'])
minimum = min(fires['temp'])
median = stat.median(fires['temp'])
diff = round((maximum - median)/4)
max_RH = max(fires['RH'])
min_RH = min(fires['RH'])
median_RH = stat.median(fires['RH'])
max_wind= max(fires['wind'])
min_wind= min(fires['wind'])
median_wind= stat.median(fires['wind'])

```

```

tp = fires['temp']
rh = fires['RH']
w = fires['wind']
serious = fires['temp'].copy()
for i, val in enumerate(fires['temp']):
    #print(val, rh[i], serious[i])
    if (val >= (median + diff) and rh[i] <= 30 and w[i] >= 3.0):
        serious.loc[i] = 4
    else:
        if (val >= median + diff):
            serious.loc[i] = 3
        else:
            if (val >= (median - diff) and rh[i] <= 50 and w[i] >= 4.0):
                serious.loc[i] = 3
            else:
                if (val >= median - 2*diff and rh[i] <= 50 and w[i] >= 5.0):
                    serious.loc[i] = 2
                else:
                    if (val >= median - 2*diff):
                        serious.loc[i] = 1
                    else:
                        serious.loc[i] = 0

```

Obraz 2: Algorytm klasyfikujący

```

x = pd.DataFrame(
    fires.values[:,0:11],
    columns =['X','Y','temp', 'RH', 'wind', 'rain', 'area']
)
y = pd.DataFrame(serious.values, columns =['SERIOUS'])

```

```

print(x.shape)
print(y.shape)

```

```

(517, 7)
(517, 1)

```

```

print(x[1:5])
print(y[1:5])

```

	X	Y	temp	RH	wind	rain	area
1	7.0	4.0	18.0	33.0	0.9	0.0	0.0
2	7.0	4.0	14.6	33.0	1.3	0.0	0.0
3	8.0	6.0	8.3	97.0	4.0	0.2	0.0
4	8.0	6.0	11.4	99.0	1.8	0.0	0.0

	SERIOUS
1	1.0
2	1.0
3	0.0
4	0.0

Obraz 3: Docelowy zbiór danych

B. Optymalizacja 1 parametru

```
# Grid Search - jedna zmienna
import numpy as np
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier
max_depth = np.array(
    [1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25])
model = RandomForestClassifier(n_estimators=2)
grid = GridSearchCV(estimator=model, param_grid=dict(max_depth=max_depth))
grid.fit(x, y.values.ravel())

GridSearchCV(estimator=RandomForestClassifier(n_estimators=2),
              param_grid={'max_depth': array([ 1,  2,  3,  4,  5,  6,  7,
        8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
        18, 19, 20, 21, 22, 23, 24, 25])})

print(grid)

GridSearchCV(estimator=RandomForestClassifier(n_estimators=2),
              param_grid={'max_depth': array([ 1,  2,  3,  4,  5,  6,  7,
        8,  9, 10, 11, 12, 13, 14, 15, 16, 17,
        18, 19, 20, 21, 22, 23, 24, 25])})

print(grid.cv_results_['params'])

[{'max_depth': 1}, {'max_depth': 2}, {'max_depth': 3}, {'max_depth': 4},
 {'max_depth': 5}, {'max_depth': 6}, {'max_depth': 7}, {'max_depth': 8},
 {'max_depth': 9}, {'max_depth': 10}, {'max_depth': 11}, {'max_depth': 12},
 {'max_depth': 13}, {'max_depth': 14}, {'max_depth': 15}, {'max_depth': 16},
 {'max_depth': 17}, {'max_depth': 18}, {'max_depth': 19}, {'max_depth': 20},
 {'max_depth': 21}, {'max_depth': 22}, {'max_depth': 23}, {'max_depth': 24},
 {'max_depth': 25}]

print(grid.cv_results_['mean_test_score'])

[0.57847274 0.65408887 0.71544063 0.79878641 0.85287528 0.81596341
 0.80477969 0.74665795 0.81643017 0.86273338 0.8626587  0.80070948
 0.83956311 0.85489171 0.80481703 0.84348394 0.83584765 0.79137416
 0.81803585 0.84150485 0.80830844 0.81435773 0.83943241 0.84156087
 0.8664115 ]
```

Obraz 4: Optymalizacja modelu - maksymalna głębokość drzew decyzyjnych

```
print(grid.best_params_)

{'max_depth': 7}
```

Obraz 5: Wynik optymalizacji modelu dla maksymalnej głębokości drzew decyzyjnych - najlepsze rozwiązanie

C. Optymalizacja 2 parametrów

```
# Grid Search - dwie zmienne
max_depth = np.array(
    [1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25])
n_estimators = np.array([2,3,4,5,6,7,8,9,10])
model = RandomForestClassifier()
grid = GridSearchCV(estimator=model,
                    param_grid=dict(
                        n_estimators=n_estimators,
                        max_depth=max_depth))
grid.fit(x, y.values.ravel())
print(grid.cv_results_['params'])
```

```
[{'max_depth': 1, 'n_estimators': 2}, {'max_depth': 1, 'n_estimators': 3}, {'max_depth': 1, 'n_estimators': 4}, {'max_depth': 1, 'n_estimators': 5}, {'max_depth': 1, 'n_estimators': 6}, {'max_depth': 1, 'n_estimators': 7}, {'max_depth': 1, 'n_estimators': 8}, {'max_depth': 1, 'n_estimators': 9}, {'max_depth': 1, 'n_estimators': 10}, {'max_depth': 2, 'n_estimators': 2}, {'max_depth': 2, 'n_estimators': 3}, {'max_depth': 2, 'n_estimators': 4}, {'max_depth': 2, 'n_estimators': 5}, {'max_depth': 2, 'n_estimators': 6}, {'max_depth': 2, 'n_estimators': 7}, {'max_depth': 2, 'n_estimators': 8}, {'max_depth': 2, 'n_estimators': 9}, {'max_depth': 2, 'n_estimators': 10}, {'max_depth': 3, 'n_estimators': 2}, {'max_depth': 3, 'n_estimators': 3}, {'max_depth': 3, 'n_estimators': 4}, {'max_depth': 3, 'n_estimators': 5}, {'max_depth': 3, 'n_estimators': 6}, {'max_depth': 3, 'n_estimators': 7}, {'max_depth': 3, 'n_estimators': 8}, {'max_depth': 3, 'n_estimators': 9}, {'max_depth': 3, 'n_estimators': 10}, {'max_depth': 4, 'n_estimators': 2}, {'max_depth': 4, 'n_estimators': 3}, {'max_depth': 4, 'n_estimators': 4}, {'max_depth': 4, 'n_estimators': 5}, {'max_depth': 4, 'n_estimators': 6}, {'max_depth': 4, 'n_estimators': 7}, {'max_depth': 4, 'n_estimators': 8}, {'max_depth': 4, 'n_estimators': 9}, {'max_depth': 4, 'n_estimators': 10}, {'max_depth': 5, 'n_estimators': 2}, {'max_depth': 5, 'n_estimators': 3}, {'max_depth': 5, 'n_estimators': 4}, {'max_depth': 5, 'n_estimators': 5}, {'max_depth': 5, 'n_estimators': 6}, {'max_depth': 5, 'n_estimators': 7}, {'max_depth': 5, 'n_estimators': 8}, {'max_depth': 5, 'n_estimators': 9}, {'max_depth': 5, 'n_estimators': 10}]]
```

Obraz 6: Optymalizacja modelu - maksymalna głębokość drzew decyzyjnych oraz ilość drzew decyzyjnych

```
print(grid.cv_results_['mean_test_score'])
```

```
[0.51460045 0.57645631 0.60541449 0.61712099 0.60944735 0.56092233
0.59002987 0.61314414 0.61486184 0.64826363 0.74902913 0.71581404
0.70617999 0.71603809 0.71796117 0.73144137 0.75614264 0.75828977
0.79107543 0.74673264 0.77976102 0.77968633 0.7448469 0.82593353
0.78541822 0.82203137 0.85309933 0.76773712 0.83370052 0.85870052
0.8317401 0.87050037 0.85888723 0.87038835 0.86663555 0.87236744
0.81239731 0.8801531 0.85517177 0.85306199 0.87442121 0.90147498
0.87821135 0.89759149 0.89361464 0.80660941 0.8742345 0.845295
0.89929052 0.90920463 0.89951456 0.89171023 0.91880134 0.92651232
0.80644137 0.88024645 0.9360717 0.88769604 0.90918596 0.94199029
0.94191561 0.93037715 0.93226288 0.79899178 0.87238611 0.89368932
0.89578043 0.92649365 0.89367065 0.9400112 0.91693428 0.94008588
0.81794249 0.89557506 0.88986184 0.9150112 0.89951456 0.94385736
0.94970127 0.94777819 0.93610904 0.81446975 0.88795743 0.89557506
0.91680358 0.9360717 0.92653099 0.93610904 0.93416729 0.94195295
0.78920836 0.88190814 0.90517177 0.90713219 0.91497386 0.94002987
0.92845407 0.93620239 0.95356609 0.84144884 0.85875653 0.90522778
0.90892457 0.91284541 0.93616505 0.92078043 0.95367812 0.93612771
0.8084578 0.85677745 0.84133682 0.93806945 0.93424197 0.93616505
0.9457991 0.94383869 0.95356609 0.83179612 0.88579164 0.89540702
0.90918596 0.91299477 0.95946602 0.95351008 0.94398805 0.95746826
0.79292382 0.8800224 0.89156087 0.93030246 0.92447722 0.94012323
0.92266617 0.94585512 0.9438947 0.8084578 0.88205751 0.8760829
0.94006721 0.89544436 0.93231889 0.94789022 0.92464526 0.91680358
0.81430172 0.85694548 0.86077296 0.93237491 0.89927184 0.93612771
0.92638163 0.94578043 0.92656833 0.81053025 0.89966393 0.8877894
0.92259149 0.89363331 0.91870799 0.92837939 0.9284354 0.92845407
0.8220127 0.85886856 0.89540702 0.9245519 0.89749813 0.91695295
0.93241225 0.94973861 0.93427931 0.81988424 0.88579164 0.90332338
0.92462659 0.9129761 0.92261016 0.93424197 0.92470127 0.94583645
0.78711725 0.90328603 0.87050037 0.90909261 0.91682226 0.92649365
0.94395071 0.93039582 0.93995519 0.82397311 0.88581031 0.88398058
0.88013443 0.93032114 0.95942868 0.94975728 0.9381068 0.94970127
0.85498506 0.88980583 0.87815534 0.92064974 0.92063107 0.94581777
0.91491785 0.93805078 0.94973861 0.81390963 0.89725541 0.87823002
0.91874533 0.91874533 0.93233757 0.93814414 0.95933532 0.94970127
0.87247946 0.88588499 0.87408514 0.92261016 0.93041449 0.9284354
0.91880134 0.93816281 0.94193428]
```

Obraz 7: Wyświetlenie średnich dopasowań modelu dla wszystkich kombinacji

```
print(grid.best_params_)
```

```
{'max_depth': 25, 'n_estimators': 9}
```

Obraz 8: Wynik optymalizacji modelu dla maksymalnej głębokości drzew decyzyjnych oraz ilości drzew - najlepsze rozwiązanie

D. Optymalizacja 3 parametrów

```
criterion = np.array(['gini', 'entropy'])
max_depth = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21])
n_estimators = np.array([2, 3, 4, 5, 6, 7, 8, 9, 10])
model = RandomForestClassifier()
grid = GridSearchCV(estimator=model, param_grid=dict(n_estimators=n_estimators,
max_depth=max_depth, criterion=criterion))
grid.fit(x, y.values.ravel())
print(grid.cv_results_['params'])
```

```
[{'criterion': 'gini', 'max_depth': 1, 'n_estimators': 2}, {'criterion': 'gini', 'max_depth': 1, 'n_estimators': 3}, {'criterion': 'gini', 'max_depth': 1, 'n_estimators': 4}, {'criterion': 'gini', 'max_depth': 1, 'n_estimators': 5}, {'criterion': 'gini', 'max_depth': 1, 'n_estimators': 6}, {'criterion': 'gini', 'max_depth': 1, 'n_estimators': 7}, {'criterion': 'gini', 'max_depth': 1, 'n_estimators': 8}, {'criterion': 'gini', 'max_depth': 1, 'n_estimators': 9}, {'criterion': 'gini', 'max_depth': 1, 'n_estimators': 10}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 2}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 3}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 4}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 5}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 6}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 7}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 8}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 9}, {'criterion': 'gini', 'max_depth': 2, 'n_estimators': 10}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 2}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 3}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 4}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 5}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 6}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 7}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 8}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 9}, {'criterion': 'gini', 'max_depth': 3, 'n_estimators': 10}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 2}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 3}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 4}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 5}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 6}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 7}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 8}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 9}, {'criterion': 'entropy', 'max_depth': 1, 'n_estimators': 10}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 2}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 3}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 4}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 5}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 6}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 7}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 8}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 9}, {'criterion': 'entropy', 'max_depth': 2, 'n_estimators': 10}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 2}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 3}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 4}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 5}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 6}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 7}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 8}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 9}, {'criterion': 'entropy', 'max_depth': 3, 'n_estimators': 10}]]
```

Obraz 9: Optymalizacja modelu - maksymalna głębokość drzew decyzyjnych, ilość drzew oraz sposób sprawdzania dopasowania podczas podziału w drzewach decyzyjnych


```
print(grid.cv_results_['mean_test_score'])
```

0.51051158	0.57615758	0.56282674	0.54570575	0.5627894	0.63030246
0.63840553	0.51039955	0.66896938	0.68246826	0.6496826	0.71751307
0.70982076	0.75423824	0.77957431	0.75218447	0.75054145	0.75263256
0.72522405	0.77561613	0.76023152	0.80261389	0.77959298	0.81250934
0.77397311	0.8471994	0.84729276	0.70026139	0.70789768	0.84914115
0.78732263	0.83762136	0.85494772	0.83967513	0.86848394	0.83384989
0.82025765	0.84144884	0.8568708	0.86269604	0.88203883	0.83956311
0.89930919	0.84919716	0.88394324	0.84917849	0.86842793	0.91103435
0.87410381	0.93420463	0.87826736	0.89749813	0.91499253	0.9071882
0.79301718	0.82796863	0.88000373	0.90321135	0.91872666	0.91290142
0.92651232	0.93614638	0.94971994	0.8434466	0.8993652	0.89352128
0.92262883	0.90905527	0.93037715	0.92645631	0.93622106	0.9554705
0.85089619	0.91114638	0.91110904	0.91678491	0.91295743	0.93220687
0.94004854	0.94002987	0.9575056	0.85868185	0.85110157	0.87414115
0.91687827	0.93808813	0.94583645	0.94002987	0.94393204	0.94779686
0.82410381	0.86454444	0.90134429	0.91680358	0.93614638	0.93045183
0.93228155	0.9690814	0.94193428	0.82191934	0.88033981	0.86842793
0.90143764	0.94006721	0.94578043	0.93618372	0.959354	0.94970127
0.82380508	0.88984317	0.90537715	0.89759149	0.92649365	0.92852875
0.93805078	0.9420463	0.94576176	0.81803585	0.88202016	0.88401792
0.93416729	0.918764	0.92839806	0.92458925	0.94970127	0.95746826
0.80683346	0.92063107	0.89361464	0.89559373	0.92266617	0.93806945
0.9342233	0.93032114	0.94772218	0.81209858	0.88771471	0.91103435
0.9092233	0.91480583	0.93039582	0.94970127	0.96325616	0.94960792
0.81992158	0.84746079	0.90911128	0.9110717	0.93224421	0.94008588
0.9342233	0.94195295	0.92457058	0.79292382	0.86476848	0.90319268
0.9264003	0.90511576	0.92656833	0.95160568	0.9304705	0.94383869
0.7928118	0.84113144	0.89932786	0.89759149	0.93230022	0.94010456
0.92063107	0.92852875	0.93424197	0.81629948	0.88790142	0.87440254
0.93228155	0.93612771	0.93032114	0.95356609	0.9438947	0.94583645
0.82007095	0.89372666	0.84551904	0.92462659	0.90916729	0.93228155
0.93808813	0.93801344	0.95554518	0.8279873	0.86648618	0.89940254
0.93808813	0.9148245	0.94391337	0.91112771	0.918764	0.94589246
0.79699403	0.89176624	0.86648618	0.93429798	0.88784541	0.94387603
0.94975728	0.93426064	0.94010456	0.84133682	0.8510829	0.91491785
0.92068708	0.91863331	0.93418596	0.95162435	0.95162435	0.95168036
0.75817774	0.88030246	0.90716953	0.93041449	0.88793876	0.94975728
0.93237491	0.94002987	0.94583645	0.61699029	0.52793129	0.61292009
0.51258402	0.49128081	0.62677371	0.60365945	0.62229276	0.59182226
0.60522778	0.74477222	0.69258775	0.67681105	0.74662061	0.79693801
0.7543876	0.74116878	0.72931292	0.69212099	0.78160941	0.73528753
0.81060493	0.78943241	0.78933906	0.79335325	0.8067401	0.81801718
0.81629948	0.8026699	0.82214339	0.84927184	0.80072816	0.83771471
0.87628827	0.87625093	0.8568708	0.76585138	0.85113891	0.87040702
0.87242345	0.88592233	0.90334205	0.89171023	0.88586632	0.90726288
0.87038835	0.87050037	0.8935773	0.89749813	0.91103435	0.92645631
0.92460792	0.93806945	0.92651232	0.77175131	0.80263256	0.91084765
0.91499253	0.90916729	0.9400112	0.90713219	0.93424197	0.92449589
0.86256535	0.85694548	0.91689694	0.89962659	0.91874533	0.94975728
0.92847274	0.91303211	0.95168036	0.82010829	0.87427184	0.89350261
0.92453323	0.95354742	0.93033981	0.93241225	0.93230022	0.9438947
0.85113891	0.85866318	0.89938387	0.89382001	0.91301344	0.91687827
0.93235624	0.93614638	0.93427931	0.86667289	0.88786408	0.90705751
0.92053771	0.93032114	0.92068708	0.91105302	0.96321882	0.93239358
0.79895444	0.8760829	0.88995519	0.92249813	0.91502987	0.93814414
0.9304705	0.95556385	0.95349141	0.80274459	0.86644884	0.91105302
0.91676624	0.91689694	0.93808813	0.92457058	0.94591113	0.94774085
0.77180732	0.86458178	0.89945855	0.91678491	0.9321882	0.9246826
0.93424197	0.92836072	0.93618372	0.75659074	0.81213592	0.89930919
0.90522778	0.8993652	0.93411128	0.92849141	0.94975728	0.94779686
0.82408514	0.88991785	0.89353996	0.92654966	0.93622106	0.94774085
0.96516057	0.95554518	0.9420463	0.80825243	0.86450709	0.92451456
0.90323002	0.90722554	0.93030246	0.93033981	0.94383869	0.93048917
0.83948842	0.85877521	0.89957058	0.90532114	0.91308813	0.93995519
0.92274085	0.92656833	0.9400112	0.85317401	0.87055639	0.87432786
0.90726288	0.93235624	0.91674757	0.93801344	0.92845407	0.97100448
0.81615011	0.89171023	0.89563107	0.93414862	0.93623973	0.90528379
0.91891337	0.93022778	0.93420463	0.81224795	0.87630695	0.90715086
0.91870799	0.94585512	0.93812547	0.95941001	0.96129574	0.92849141
0.80646004	0.88775205	0.90330471	0.93220687	0.91880134	0.93039582
0.9245519	0.93426064	0.94973861	0.76773712	0.86256535	0.84133682
0.90530246	0.91684093	0.95941001	0.91878267	0.9420463	0.9438947
0.79889843	0.87229276	0.9071882	0.92654966	0.92262883	0.95164302
0.91489918	0.93226288	0.92649365	0.76971621	0.86652353	0.85679612
0.92285288	0.92262883	0.9092233	0.93043316	0.93424197	0.92841673]

Obraz 10: Średnie wyniki dla wszystkich kombinacji podczas optymalizacji modelu

```
print(grid.best_params_)
```

```
{'criterion': 'entropy', 'max_depth': 20, 'n_estimators': 9}
```

Obraz 11: Wynik optymalizacji modelu dla maksymalnej głębokość drzew decyzyjnych, ilości drzew oraz sposobu sprawdzania dopasowania podczas podziału w drzewach decyzyjnych - najlepsze rozwiązanie

E. Najlepsze parametry dla wybranego modelu - zestawienie wyników optymalizacji modelu

Tabela 1: Zestawienie najlepszych parametrów w zależności od rodzaju optymalizacji modelu

Parametry optymalizacji	Najlepszy wynik: max_depth	Najlepszy wynik: n_estimators	Najlepszy wynik: criterion
<i>max_depth</i>	7	<i>default (2)</i>	<i>default („gini”)</i>
<i>max_depth, n_estimators</i>	25	9	<i>default („gini”)</i>
<i>max_depth, n_estimators, criterion</i>	20	9	„entropy”

Dla najprostszego modelu wystarczy las drzew decyzyjnych z dwoma drzewami o maksymalnej głębokości równej 7. Dla najbardziej dokładnego modelu potrzeba las drzew decyzyjnych z 9 drzewami o głębokości maksymalnej równej 20 oraz z kryterium dopasowywania ustawionym na „entropy”.

4. Przetrenowanie

4.1. Zadanie #2

i. Treść polecenia

Proszę zaprezentować zjawisko overfitingu na klasyfikatorze random forest.

ii. Rozwiązanie

A. Ładowanie bazy danych i bibliotek

```
import pandas as pd
from sklearn.ensemble import RandomForestClassifier

data = pd.read_csv("glass.data", header=None, index_col=0)
print(data.head())
x = pd.DataFrame(data.values[:,0:9], columns = [
    'RI', 'Na', 'Mg', 'Al', 'Si', 'K', 'Ca', 'Ba', 'Fe'])
y = pd.DataFrame(data.values[:,9], columns = ['Type'])

    1      2      3      4      5      6      7      8      9    10
0
1  1.52101  13.64  4.49  1.10  71.78  0.06  8.75  0.0  0.0  1
2  1.51761  13.89  3.60  1.36  72.73  0.48  7.83  0.0  0.0  1
3  1.51618  13.53  3.55  1.54  72.99  0.39  7.78  0.0  0.0  1
4  1.51766  13.21  3.69  1.29  72.61  0.57  8.22  0.0  0.0  1
5  1.51742  13.27  3.62  1.24  73.08  0.55  8.07  0.0  0.0  1

print(x[1:5])
print(y[1:5])

    RI      Na      Mg      Al      Si      K      Ca      Ba      Fe
1  1.51761  13.89  3.60  1.36  72.73  0.48  7.83  0.0  0.0
2  1.51618  13.53  3.55  1.54  72.99  0.39  7.78  0.0  0.0
3  1.51766  13.21  3.69  1.29  72.61  0.57  8.22  0.0  0.0
4  1.51742  13.27  3.62  1.24  73.08  0.55  8.07  0.0  0.0
Type
1    1.0
2    1.0
3    1.0
4    1.0

print(x.shape)
print(y.shape)

(214, 9)
(214, 1)

print(x[1:5])
print(y[1:5])

    RI      Na      Mg      Al      Si      K      Ca      Ba      Fe
1  1.51761  13.89  3.60  1.36  72.73  0.48  7.83  0.0  0.0
2  1.51618  13.53  3.55  1.54  72.99  0.39  7.78  0.0  0.0
3  1.51766  13.21  3.69  1.29  72.61  0.57  8.22  0.0  0.0
4  1.51742  13.27  3.62  1.24  73.08  0.55  8.07  0.0  0.0
Type
1    1.0
2    1.0
3    1.0
4    1.0
```

Obraz 12: Ładowanie bazy danych z pliku

```
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
clf = RandomForestClassifier()
```

Obraz 13: Dodanie bibliotek do trenowania modelu i przypisanie klasyfikatora drzewa losowego do zmiennej

B. Duże przetrenowanie

```
x_train, x_test, y_train, y_test = train_test_split(
    x, y, test_size=0.9, random_state=42)
```

```
print(x_train.shape)
print(y_train.shape)
print(x_test.shape)
print(y_test.shape)
```

```
(21, 9)
(21, 1)
(193, 9)
(193, 1)
```

```
clf = clf.fit(x_train, y_train.values.ravel())
y_train_pred_from_clf = clf.predict(x_train)
y_test_pred_from_clf = clf.predict(x_test)
print(accuracy_score(y_train, y_train_pred_from_clf))
print(accuracy_score(y_test, y_test_pred_from_clf))
```

```
1.0
0.5751295336787565
```

```
clf = clf.fit(x_train, y_train.values.ravel())
y_train_pred_from_clf = clf.predict(x_train)
y_test_pred_from_clf = clf.predict(x_test)
print(accuracy_score(y_train, y_train_pred_from_clf))
print(accuracy_score(y_test, y_test_pred_from_clf))
```

```
1.0
0.5544041450777202
```

```
clf = clf.fit(x_train, y_train.values.ravel())
y_train_pred_from_clf = clf.predict(x_train)
y_test_pred_from_clf = clf.predict(x_test)
print(accuracy_score(y_train, y_train_pred_from_clf))
print(accuracy_score(y_test, y_test_pred_from_clf))
```

```
1.0
0.6010362694300518
```

Obraz 14: Trenowanie modelu na bardzo małym zbiorze trenującym (10%) oraz wyniki dokładności modelu

- występuje duże przetrenowanie modelu

C. Małe przetrenowanie

```
x_train, x_test, y_train, y_test = train_test_split(
    x, y, test_size=0.1, random_state=42)
print(x_train.shape)
print(y_train.shape)
print(x_test.shape)
print(y_test.shape)
```

```
(192, 9)
(192, 1)
(22, 9)
(22, 1)
```

```
clf = clf.fit(x_train, y_train.values.ravel())
y_train_pred_from_clf = clf.predict(x_train)
y_test_pred_from_clf = clf.predict(x_test)
print(accuracy_score(y_train, y_train_pred_from_clf))
print(accuracy_score(y_test, y_test_pred_from_clf))
```

```
1.0
0.9545454545454545
```

```
clf = clf.fit(x_train, y_train.values.ravel())
y_train_pred_from_clf = clf.predict(x_train)
y_test_pred_from_clf = clf.predict(x_test)
print(accuracy_score(y_train, y_train_pred_from_clf))
print(accuracy_score(y_test, y_test_pred_from_clf))
```

```
1.0
0.9545454545454545
```

```
clf = clf.fit(x_train, y_train.values.ravel())
y_train_pred_from_clf = clf.predict(x_train)
y_test_pred_from_clf = clf.predict(x_test)
print(accuracy_score(y_train, y_train_pred_from_clf))
print(accuracy_score(y_test, y_test_pred_from_clf))
```

```
1.0
0.8636363636363636
```

Obraz 15: Trenowanie modelu na dużym zbiorze trenującym (90%) i wyniki dokładności

- występuje małe przetrenowanie modelu

iii. Wnioski

Bardzo dokładne trenowanie modelu nie jest wskazane, ponieważ może prowadzić do przetrenowania - model będzie się idealnie sprawdzał na danych ze zbioru trenującego, jednak będzie niewystarczająco skuteczny na danych ze zbioru testującego. Lepszym rozwiązaniem jest mniejsza dokładność na modelu trenującym, aby dokładność na danych testujących była zbliżona. Przetrenowaniu można zapobiec zwiększając zbiór trenujący oraz zmniejszając rygorystyczność warunków podczas trenowania modelu.

5. Analiza tekstu

5.1. Zadanie #3

i. Treść polecenia

Proszę powtórzyć ćwiczenia dla innego podzbioru kategorii, np. rec.autos, rec.motorcycles, rec.sport.baseball, rec.sport.hockey.

ii. Rozwiązanie

```
# http://qwone.com/~jason/20Newsgroups/

#1. Proszę powtorzyc cwiczenia dla innego podzbioru kategorii

#wczytanie zbioru trenujacego danych tekstowych
from sklearn.datasets import fetch_20newsgroups
categories = [
    'sci.crypt',
    'rec.sport.baseball',
    'rec.motorcycles',
    'sci.space']
train_dataset = fetch_20newsgroups(
    subset='train',
    categories=categories,
    shuffle=True,
    random_state=42
)
test_dataset = fetch_20newsgroups(
    subset='test',
    categories=categories,
    shuffle=True,
    random_state=42
)

#klasy
print(train_dataset.target_names)
print(test_dataset.target_names)

['rec.motorcycles', 'rec.sport.baseball', 'sci.crypt', 'sci.space']
['rec.motorcycles', 'rec.sport.baseball', 'sci.crypt', 'sci.space']

#liczba wiadomosci
print(len(train_dataset.data))
print(len(test_dataset.data))

2383
1585
```

Obraz 16: Załadowanie tekstów wiadomości e-mail

```
#wyswietlenie wiadomosci ze zbioru trenujacego o indeksie 180
print(train_dataset.data[180])
```

Obraz 17: Wyświetlenie wiadomości o indeksie 180

From: ajjb@adam4.bnsc.rl.ac.uk (Andrew Broderick)

Subject: Re: Solar Sail Data

Keywords: Solar Sail

Organization: Rutherford Appleton Lab, UK

Lines: 79

In article <1993Apr15.051746.29848@news.duc.auburn.edu> snydefj@eng.auburn.edu writes:

>

>I am looking for any information concerning projects involving Solar

> Sails

I was at an interesting seminar at work (UK's R.A.L. Space Science Dept.) on this subject, specifically on a small-scale Solar Sail proposed as a student space project. The guy giving the talk was keen to generate interest in the project. I'll type in the handout he gave out at the meeting. Here goes :

The Microlight Solar Sail

1. Introduction

The solar sail is a well-established concept. Harnessing the pressure of sunlight, a spacecraft would have unlimited range. In principle, such a vehicle could explore the whole Solar System with zero fuel consumption.

However it is more difficult to design a practical solar sail than most people realize. The pressure of sunlight is only about one kilogram per square kilometer. Deploying and controlling the large area of aluminized fabric which would be necessary to transport a 'conventional' type spacecraft is a daunting task. This is why, despite the potential of the idea, no such craft has actually been launched to date.

2. Design

Recent advances in microelectronics make possible a different concept: a tiny sail just a few metres in diameter which could be controlled purely by electronics, with no mechanical parts. Several attitude control methods are feasible: for example the pressure sunlight exerts on a panel of solar cells varies according to whether power is being drawn.

The key components of the craft will be a minute CCD camera developed at Edinburgh University which can act as both attitude sensor and data gathering device; solar cells providing ~1 watt power for control and communication; and a directional radio antenna etched onto the surface of the sail itself. Launched as a piggyback payload, the total cost of the mission can be limited to a few tens of thousands of dollars.

3. Missions

The craft would be capable of some ambitious missions. For example:
a) It could rendezvous with a nearby asteroid from the Apollo or Amor groups. Closeup pictures could be transmitted back to Earth at a low bit

rate.

b) It could be steered into a lunar polar orbit. Previously unobserved areas around the lunar poles could be viewed. By angling the sail to reflect sunlight downwards, polar craters whose bases never receive sunlight could be imaged. Bright reflections would confirm that volatiles such as water ice have become trapped in these locations.[Immensely valuable information for setting up a manned lunar base, BTW]

c) It could be sent to rendezvous with a small asteroid or comet nucleus. Impacting at low speed, a thin wire probe attached to the craft causes it to rebound while capturing a tiny sample in a sharp-edged tube, like performing a biopsy. Returning to Earth, the sail acts as an ideal re-entry parachute: load per unit area 20 gm/m² ensures that heat is reradiated so effectively that the sail temperature cannot exceed ~300 deg C. The material sample is recovered, enclosed in a small insulating container.

Contact: Colin Jack Tel. 0865-200447
Oxford Mathematical Designs, 131 High Street, Oxford OX1 4DH, England

This guy would love to hear from anyone interested in this project or seeking details or anything, and would be most happy to send you more information.

Andy

--

Andy Jonathan J. Broderick, | "I have come that they might have |
Rutherford Lab., UK | life, and have it to the full" |
Mail : ajjb@adam2.bnsc.rl.ac.uk | - Jesus Christ |

```
#wyswietlenie klasy wiadomosci o indeksie 180
print(train_dataset.target_names[train_dataset.target[180]])

sci.space
```

Obraz 18: Dopasowana kategoria wiadomości o indeksie 180

Wybrana wiadomość została przydzielona do kategorii „sci.space”, ponieważ zawiera słowa kluczowe typu „science”, „space” czy „Earth”.

```
from sklearn.feature_extraction.text import CountVectorizer
count_vect = CountVectorizer(analyzer='word', ngram_range=(1, 1))
X_train_counts = count_vect.fit_transform(train_dataset.data)
```

```
word_list = count_vect.get_feature_names_out()
count_list = X_train_counts.toarray().sum(axis=0)
word_count = dict(zip(word_list, count_list))
sorted(word_count.items(), key=lambda x: x[1], reverse=True)
```

```
[('the', 30530),
 ('to', 15205),
 ('of', 13501),
 ('and', 11909),
 ('in', 9861),
 ('is', 8380),
 ('that', 6952),
 ('it', 6349),
 ('for', 5829),
 ('on', 4651),
 ('from', 4554),
 ('you', 4510),
 ('edu', 4471),
 ('be', 4399),
 ('this', 4193),
 ('are', 3351),
 ('with', 3291),
 ('have', 3193),
 ('as', 3122),
 ('at', 2990),
 ('a', 2988),
 ('on', 2987),
 ('the', 2986),
 ('in', 2985),
 ('of', 2984),
 ('and', 2983),
 ('is', 2982),
 ('that', 2981),
 ('it', 2980),
 ('for', 2979),
 ('on', 2978),
 ('from', 2977),
 ('you', 2976),
 ('edu', 2975),
 ('be', 2974),
 ('this', 2973),
 ('are', 2972),
 ('with', 2971),
 ('have', 2970),
 ('as', 2969),
 ('at', 2968),
 ('a', 2967),
 ('on', 2966),
 ('the', 2965),
 ('in', 2964),
 ('of', 2963),
 ('and', 2962),
 ('is', 2961),
 ('that', 2960),
 ('it', 2959),
 ('for', 2958),
 ('on', 2957),
 ('from', 2956),
 ('you', 2955),
 ('edu', 2954),
 ('be', 2953),
 ('this', 2952),
 ('are', 2951),
 ('with', 2950),
 ('have', 2949),
 ('as', 2948),
 ('at', 2947),
 ('a', 2946),
 ('on', 2945),
 ('the', 2944),
 ('in', 2943),
 ('of', 2942),
 ('and', 2941),
 ('is', 2940),
 ('that', 2939),
 ('it', 2938),
 ('for', 2937),
 ('on', 2936),
 ('from', 2935),
 ('you', 2934),
 ('edu', 2933),
 ('be', 2932),
 ('this', 2931),
 ('are', 2930),
 ('with', 2929),
 ('have', 2928),
 ('as', 2927),
 ('at', 2926),
 ('a', 2925),
 ('on', 2924),
 ('the', 2923),
 ('in', 2922),
 ('of', 2921),
 ('and', 2920),
 ('is', 2919),
 ('that', 2918),
 ('it', 2917),
 ('for', 2916),
 ('on', 2915),
 ('from', 2914),
 ('you', 2913),
 ('edu', 2912),
 ('be', 2911),
 ('this', 2910),
 ('are', 2909),
 ('with', 2908),
 ('have', 2907),
 ('as', 2906),
 ('at', 2905),
 ('a', 2904),
 ('on', 2903),
 ('the', 2902),
 ('in', 2901),
 ('of', 2900),
 ('and', 2899),
 ('is', 2898),
 ('that', 2897),
 ('it', 2896),
 ('for', 2895),
 ('on', 2894),
 ('from', 2893),
 ('you', 2892),
 ('edu', 2891),
 ('be', 2890),
 ('this', 2889),
 ('are', 2888),
 ('with', 2887),
 ('have', 2886),
 ('as', 2885),
 ('at', 2884),
 ('a', 2883),
 ('on', 2882),
 ('the', 2881),
 ('in', 2880),
 ('of', 2879),
 ('and', 2878),
 ('is', 2877),
 ('that', 2876),
 ('it', 2875),
 ('for', 2874),
 ('on', 2873),
 ('from', 2872),
 ('you', 2871),
 ('edu', 2870),
 ('be', 2869),
 ('this', 2868),
 ('are', 2867),
 ('with', 2866),
 ('have', 2865),
 ('as', 2864),
 ('at', 2863),
 ('a', 2862),
 ('on', 2861),
 ('the', 2860),
 ('in', 2859),
 ('of', 2858),
 ('and', 2857),
 ('is', 2856),
 ('that', 2855),
 ('it', 2854),
 ('for', 2853),
 ('on', 2852),
 ('from', 2851),
 ('you', 2850),
 ('edu', 2849),
 ('be', 2848),
 ('this', 2847),
 ('are', 2846),
 ('with', 2845),
 ('have', 2844),
 ('as', 2843),
 ('at', 2842),
 ('a', 2841),
 ('on', 2840),
 ('the', 2839),
 ('in', 2838),
 ('of', 2837),
 ('and', 2836),
 ('is', 2835),
 ('that', 2834),
 ('it', 2833),
 ('for', 2832),
 ('on', 2831),
 ('from', 2830),
 ('you', 2829),
 ('edu', 2828),
 ('be', 2827),
 ('this', 2826),
 ('are', 2825),
 ('with', 2824),
 ('have', 2823),
 ('as', 2822),
 ('at', 2821),
 ('a', 2820),
 ('on', 2819),
 ('the', 2818),
 ('in', 2817),
 ('of', 2816),
 ('and', 2815),
 ('is', 2814),
 ('that', 2813),
 ('it', 2812),
 ('for', 2811),
 ('on', 2810),
 ('from', 2809),
 ('you', 2808),
 ('edu', 2807),
 ('be', 2806),
 ('this', 2805),
 ('are', 2804),
 ('with', 2803),
 ('have', 2802),
 ('as', 2801),
 ('at', 2800),
 ('a', 2799),
 ('on', 2798),
 ('the', 2797),
 ('in', 2796),
 ('of', 2795),
 ('and', 2794),
 ('is', 2793),
 ('that', 2792),
 ('it', 2791),
 ('for', 2790),
 ('on', 2789),
 ('from', 2788),
 ('you', 2787),
 ('edu', 2786),
 ('be', 2785),
 ('this', 2784),
 ('are', 2783),
 ('with', 2782),
 ('have', 2781),
 ('as', 2780),
 ('at', 2779),
 ('a', 2778),
 ('on', 2777),
 ('the', 2776),
 ('in', 2775),
 ('of', 2774),
 ('and', 2773),
 ('is', 2772),
 ('that', 2771),
 ('it', 2770),
 ('for', 2769),
 ('on', 2768),
 ('from', 2767),
 ('you', 2766),
 ('edu', 2765),
 ('be', 2764),
 ('this', 2763),
 ('are', 2762),
 ('with', 2761),
 ('have', 2760),
 ('as', 2759),
 ('at', 2758),
 ('a', 2757),
 ('on', 2756),
 ('the', 2755),
 ('in', 2754),
 ('of', 2753),
 ('and', 2752),
 ('is', 2751),
 ('that', 2750),
 ('it', 2749),
 ('for', 2748),
 ('on', 2747),
 ('from', 2746),
 ('you', 2745),
 ('edu', 2744),
 ('be', 2743),
 ('this', 2742),
 ('are', 2741),
 ('with', 2740),
 ('have', 2739),
 ('as', 2738),
 ('at', 2737),
 ('a', 2736),
 ('on', 2735),
 ('the', 2734),
 ('in', 2733),
 ('of', 2732),
 ('and', 2731),
 ('is', 2730),
 ('that', 2729),
 ('it', 2728),
 ('for', 2727),
 ('on', 2726),
 ('from', 2725),
 ('you', 2724),
 ('edu', 2723),
 ('be', 2722),
 ('this', 2721),
 ('are', 2720),
 ('with', 2719),
 ('have', 2718),
 ('as', 2717),
 ('at', 2716),
 ('a', 2715),
 ('on', 2714),
 ('the', 2713),
 ('in', 2712),
 ('of', 2711),
 ('and', 2710),
 ('is', 2709),
 ('that', 2708),
 ('it', 2707),
 ('for', 2706),
 ('on', 2705),
 ('from', 2704),
 ('you', 2703),
 ('edu', 2702),
 ('be', 2701),
 ('this', 2700),
 ('are', 2699),
 ('with', 2698),
 ('have', 2697),
 ('as', 2696),
 ('at', 2695),
 ('a', 2694),
 ('on', 2693),
 ('the', 2692),
 ('in', 2691),
 ('of', 2690),
 ('and', 2689),
 ('is', 2688),
 ('that', 2687),
 ('it', 2686),
 ('for', 2685),
 ('on', 2684),
 ('from', 2683),
 ('you', 2682),
 ('edu', 2681),
 ('be', 2680),
 ('this', 2679),
 ('are', 2678),
 ('with', 2677),
 ('have', 2676),
 ('as', 2675),
 ('at', 2674),
 ('a', 2673),
 ('on', 2672),
 ('the', 2671),
 ('in', 2670),
 ('of', 2669),
 ('and', 2668),
 ('is', 2667),
 ('that', 2666),
 ('it', 2665),
 ('for', 2664),
 ('on', 2663),
 ('from', 2662),
 ('you', 2661),
 ('edu', 2660),
 ('be', 2659),
 ('this', 2658),
 ('are', 2657),
 ('with', 2656),
 ('have', 2655),
 ('as', 2654),
 ('at', 2653),
 ('a', 2652),
 ('on', 2651),
 ('the', 2650),
 ('in', 2649),
 ('of', 2648),
 ('and', 2647),
 ('is', 2646),
 ('that', 2645),
 ('it', 2644),
 ('for', 2643),
 ('on', 2642),
 ('from', 2641),
 ('you', 2640),
 ('edu', 2639),
 ('be', 2638),
 ('this', 2637),
 ('are', 2636),
 ('with', 2635),
 ('have', 2634),
 ('as', 2633),
 ('at', 2632),
 ('a', 2631),
 ('on', 2630),
 ('the', 2629),
 ('in', 2628),
 ('of', 2627),
 ('and', 2626),
 ('is', 2625),
 ('that', 2624),
 ('it', 2623),
 ('for', 2622),
 ('on', 2621),
 ('from', 2620),
 ('you', 2619),
 ('edu', 2618),
 ('be', 2617),
 ('this', 2616),
 ('are', 2615),
 ('with', 2614),
 ('have', 2613),
 ('as', 2612),
 ('at', 2611),
 ('a', 2610),
 ('on', 2609),
 ('the', 2608),
 ('in', 2607),
 ('of', 2606),
 ('and', 2605),
 ('is', 2604),
 ('that', 2603),
 ('it', 2602),
 ('for', 2601),
 ('on', 2600),
 ('from', 2599),
 ('you', 2598),
 ('edu', 2597),
 ('be', 2596),
 ('this', 2595),
 ('are', 2594),
 ('with', 2593),
 ('have', 2592),
 ('as', 2591),
 ('at', 2590),
 ('a', 2589),
 ('on', 2588),
 ('the', 2587),
 ('in', 2586),
 ('of', 2585),
 ('and', 2584),
 ('is', 2583),
 ('that', 2582),
 ('it', 2581),
 ('for', 2580),
 ('on', 2579),
 ('from', 2578),
 ('you', 2577),
 ('edu', 2576),
 ('be', 2575),
 ('this', 2574),
 ('are', 2573),
 ('with', 2572),
 ('have', 2571),
 ('as', 2570),
 ('at', 2569),
 ('a', 2568),
 ('on', 2567),
 ('the', 2566),
 ('in', 2565),
 ('of', 2564),
 ('and', 2563),
 ('is', 2562),
 ('that', 2561),
 ('it', 2560),
 ('for', 2559),
 ('on', 2558),
 ('from', 2557),
 ('you', 2556),
 ('edu', 2555),
 ('be', 2554),
 ('this', 2553),
 ('are', 2552),
 ('with', 2551),
 ('have', 2550),
 ('as', 2549),
 ('at', 2548),
 ('a', 2547),
 ('on', 2546),
 ('the', 2545),
 ('in', 2544),
 ('of', 2543),
 ('and', 2542),
 ('is', 2541),
 ('that', 2540),
 ('it', 2539),
 ('for', 2538),
 ('on', 2537),
 ('from', 2536),
 ('you', 2535),
 ('edu', 2534),
 ('be', 2533),
 ('this', 2532),
 ('are', 2531),
 ('with', 2530),
 ('have', 2529),
 ('as', 2528),
 ('at', 2527),
 ('a', 2526),
 ('on', 2525),
 ('the', 2524),
 ('in', 2523),
 ('of', 2522),
 ('and', 2521),
 ('is', 2520),
 ('that', 2519),
 ('it', 2518),
 ('for', 2517),
 ('on', 2516),
 ('from', 2515),
 ('you', 2514),
 ('edu', 2513),
 ('be', 2512),
 ('this', 2511),
 ('are', 2510),
 ('with', 2509),
 ('have', 2508),
 ('as', 2507),
 ('at', 2506),
 ('a', 2505),
 ('on', 2504),
 ('the', 2503),
 ('in', 2502),
 ('of', 2501),
 ('and', 2500),
 ('is', 2499),
 ('that', 2498),
 ('it', 2497),
 ('for', 2496),
 ('on', 2495),
 ('from', 2494),
 ('you', 2493),
 ('edu', 2492),
 ('be', 2491),
 ('this', 2490),
 ('are', 2489),
 ('with', 2488),
 ('have', 2487),
 ('as', 2486),
 ('at', 2485),
 ('a', 2484),
 ('on', 2483),
 ('the', 2482),
 ('in', 2481),
 ('of', 2480),
 ('and', 2479),
 ('is', 2478),
 ('that', 2477),
 ('it', 2476),
 ('for', 2475),
 ('on', 2474),
 ('from', 2473),
 ('you', 2472),
 ('edu', 2471),
 ('be', 2470),
 ('this', 2469),
 ('are', 2468),
 ('with', 2467),
 ('have', 2466),
 ('as', 2465),
 ('at', 2464),
 ('a', 2463),
 ('on', 2462),
 ('the', 2461),
 ('in', 2460),
 ('of', 2459),
 ('and', 2458),
 ('is', 2457),
 ('that', 2456),
 ('it', 2455),
 ('for', 2454),
 ('on', 2453),
 ('from', 2452),
 ('you', 2451),
 ('edu', 2450),
 ('be', 2449),
 ('this', 2448),
 ('are', 2447),
 ('with', 2446),
 ('have', 2445),
 ('as', 2444),
 ('at', 2443),
 ('a', 2442),
 ('on', 2441),
 ('the', 2440),
 ('in', 2439),
 ('of', 2438),
 ('and', 2437),
 ('is', 2436),
 ('that', 2435),
 ('it', 2434),
 ('for', 2433),
 ('on', 2432),
 ('from', 2431),
 ('you', 2430),
 ('edu', 2429),
 ('be', 2428),
 ('this', 2427),
 ('are', 2426),
 ('with', 2425),
 ('have', 2424),
 ('as', 2423),
 ('at', 2422),
 ('a', 2421),
 ('on', 2420),
 ('the', 2419),
 ('in', 2418),
 ('of', 2417),
 ('and', 2416),
 ('is', 2415),
 ('that', 2414),
 ('it', 2413),
 ('for', 2412),
 ('on', 2411),
 ('from', 2410),
 ('you', 2409),
 ('edu', 2408),
 ('be', 2407),
 ('this', 2406),
 ('are', 2405),
 ('with', 2404),
 ('have', 2403),
 ('as', 2402),
 ('at', 2401),
 ('a', 2400),
 ('on', 2399),
 ('the', 2398),
 ('in', 2397),
 ('of', 2396),
 ('and', 2395),
 ('is', 2394),
 ('that', 2393),
 ('it', 2392),
 ('for', 2391),
 ('on', 2390),
 ('from', 2389),
 ('you', 2388),
 ('edu', 2387),
 ('be', 2386),
 ('this', 2385),
 ('are', 2384),
 ('with', 2383),
 ('have', 2382),
 ('as', 2381),
 ('at', 2380),
 ('a', 2379),
 ('on', 2378),
 ('the', 2377),
 ('in', 2376),
 ('of', 2375),
 ('and', 2374),
 ('is', 2373),
 ('that', 2372),
 ('it', 2371),
 ('for', 2370),
 ('on', 2369),
 ('from', 2368),
 ('you', 2367),
 ('edu', 2366),
 ('be', 2365),
 ('this', 2364),
 ('are', 2363),
 ('with', 2362),
 ('have', 2361),
 ('as', 2360),
 ('at', 2359),
 ('a', 2358),
 ('on', 2357),
 ('the', 2356),
 ('in', 2355),
 ('of', 2354),
 ('and', 2353),
 ('is', 2352),
 ('that', 2351),
 ('it', 2350),
 ('for', 2349),
 ('on', 2348),
 ('from', 2347),
 ('you', 2346),
 ('edu', 2345),
 ('be', 2344),
 ('this', 2343),
 ('are', 2342),
 ('with', 2341),
 ('have', 2340),
 ('as', 2339),
 ('at', 2338),
 ('a', 2337),
 ('on', 2336),
 ('the', 2335),
 ('in', 2334),
 ('of', 2333),
 ('and', 2332),
 ('is', 2331),
 ('that', 2330),
 ('it', 2329),
 ('for', 2328),
 ('on', 2327),
 ('from', 2326),
 ('you', 2325),
 ('edu', 2324),
 ('be', 2323),
 ('this', 2322),
 ('are', 2321),
 ('with', 2320),
 ('have', 2319),
 ('as', 2318),
 ('at', 2317),
 ('a', 2316),
 ('on', 2315),
 ('the', 2314),
 ('in', 2313),
 ('of', 2312),
 ('and', 2311),
 ('is', 2310),
 ('that', 2309),
 ('it', 2308),
 ('for', 2307),
 ('on', 2306),
 ('from', 2305),
 ('you', 2304),
 ('edu', 2303),
 ('be', 2302),
 ('this', 2301),
 ('are', 2300),
 ('with', 2299),
 ('have', 2298),
 ('as', 2297),
 ('at', 2296),
 ('a', 2295),
 ('on', 2294),
 ('the', 2293),
 ('in', 2292),
 ('of', 2291),
 ('and', 2290),
 ('is', 2289),
 ('that', 2288),
 ('it', 2287),
 ('for', 2286),
 ('on', 2285),
 ('from', 2284),
 ('you', 2283),
 ('edu', 2282),
 ('be', 2281),
 ('this', 2280),
 ('are', 2279),
 ('with', 2278),
 ('have', 2277),
 ('as', 2276),
 ('at', 2275),
 ('a', 2274),
 ('on', 2273),
 ('the', 2272),
 ('in', 2271),
 ('of', 2270),
 ('and', 2269),
 ('is', 2268),
 ('that', 2267),
 ('it', 2266),
 ('for', 2265),
 ('on', 2264),
 ('from', 2263),
 ('you', 2262),
 ('edu', 2261),
 ('be', 2260),
 ('this', 2259),
 ('are', 2258),
 ('with', 2257),
 ('have', 2256),
 ('as', 2255),
 ('at', 2254),
 ('a', 2253),
 ('on', 2252),
 ('the', 2251),
 ('in', 2250),
 ('of', 2249),
 ('and', 2248),
 ('is', 2247),
 ('that', 2246),
 ('it', 2245),
 ('for', 2244),
 ('on', 2243),
 ('from', 2242),
 ('you', 2241),
 ('edu', 2240),
 ('be', 2239),
 ('this', 2238),
 ('are', 2237),
 ('with', 2236),
 ('have', 2235),
 ('as', 2234),
 ('at', 2233),
 ('a', 2232),
 ('on', 2231),
 ('the', 2230),
 ('in', 2229),
 ('of', 2228),
 ('and', 2227),
 ('is', 2226),
 ('that', 2225),
 ('it', 2224),
 ('for', 2223),
 ('on', 2222),
 ('from', 2221),
 ('you', 2220),
 ('edu', 2219),
 ('be', 2218),
 ('this', 2217),
 ('are', 2216),
 ('with', 2215),
 ('have', 2214),
 ('as', 2213),
 ('at', 2212),
 ('a', 2211),
 ('on', 2210),
 ('the', 2209),
 ('in', 2208),
 ('of', 2207),
 ('and', 2206),
 ('is', 2205),
 ('that', 2204),
 ('it', 2203),
 ('for', 2202),
 ('on', 2201),
 ('from', 2200),
 ('you', 2199),
 ('edu', 2198),
 ('be', 2197),
 ('this', 2196),
 ('are', 2195),
 ('with', 2194),
 ('have', 2193),
 ('as', 2192),
 ('at', 2191),
 ('a', 2190),
 ('on', 2189),
 ('the', 2188),
 ('in', 2187),
 ('of', 2186),
 ('and', 2185),
 ('is', 2184),
 ('that', 2183),
 ('it', 2182),
 ('for', 2181),
 ('on', 2180),
 ('from', 2179),
 ('you', 2178),
 ('edu', 2177),
 ('be', 2176),
 ('this', 2175),
 ('are', 2174),
 ('with', 2173),
 ('have', 2172),
 ('as', 2171),
 ('at', 2170),
 ('a', 2169),
 ('on', 2168),
 ('the', 2167),
 ('in', 2166),
 ('of', 2165),
 ('and', 2164),
 ('is', 2163),
 ('that', 2162),
 ('it', 2161),
 ('for', 2160),
 ('on', 2159),
 ('from', 2158),
 ('you', 2157),
 ('edu', 2156),
 ('be', 2155),
 ('this', 2154),
 ('are', 2153),
 ('with', 2152),
 ('have', 2151),
 ('as', 2150),
 ('at', 2149),
 ('a', 2148),
 ('on', 2147),
 ('the', 2146),
 ('in', 2145),
 ('of', 2144),
 ('and', 2143),
 ('is', 2142),
 ('that', 2141),
 ('it', 2140),
 ('for', 2139),
 ('on', 2138),
 ('from', 2137),
 ('you', 2136),
 ('edu', 2135),
 ('be', 2134),
 ('this', 2133),
 ('are', 2132),
 ('with', 2131),
 ('have', 2130),
 ('as', 2129),
 ('at', 2128),
 ('a', 2127),
 ('on', 2126),
 ('the', 2125),
 ('in', 2124),
 ('of', 2123),
 ('and', 2122),
 ('is', 2121),
 ('that', 2120),
 ('it', 2119),
 ('for', 2118),
 ('on', 2117),
 ('from', 2116),
 ('you', 2115),
 ('edu', 2114),
 ('be', 2113),
 ('this', 2112),
 ('are', 2111),
 ('with', 2110),
 ('have', 2109),
 ('as', 2108),
 ('at', 2107),
 ('a', 2106),
 ('on', 2105),
 ('the', 2104),
 ('in', 2103),
 ('of', 2102),
 ('and', 2101),
 ('is', 2100),
 ('that', 2099),
 ('it', 2098),
 ('for', 2097),
 ('on', 2096),
 ('from', 2095),
 ('you', 2094),
 ('edu', 2093),
 ('be', 2092),
 ('this', 2091),
 ('are', 2090),
 ('with', 2089),
 ('have', 2088),
 ('as', 2087),
 ('at', 2086),
 ('a', 2085),
 ('on', 2084),
 ('the', 2083),
 ('in', 2082),
 ('of', 2081),
 ('and', 2080),
 ('is', 2079),
 ('that', 2078),
 ('it', 2077),
 ('for', 2076),
 ('on', 2075),
 ('from', 2074),
 ('you', 2073),
 ('edu', 2072),
 ('be', 2071),
 ('this', 2070),
 ('are', 2069),
 ('with', 2068),
 ('have', 2067),
 ('as', 2066),
 ('at', 2065),
 ('a', 2064),
 ('on', 2063),
 ('the', 2062),
 ('in', 2061),
 ('of', 2060),
 ('and', 2059),
 ('is', 2058),
 ('that', 2057),
 ('it', 2056),
 ('for', 2055),
 ('on', 2054),
 ('from', 2053),
 ('you', 2052),
 ('edu', 2051),
 ('be', 2050),
 ('this', 2049),
 ('are', 2048),
 ('with', 2047),
 ('have', 2046),
 ('as', 2045),
 ('at', 2044),
 ('a', 2043),
 ('on', 2042),
 ('the', 2041),
 ('in', 2040),
 ('of', 2039),
 ('and', 2038),
 ('is', 2037),
 ('that', 2036),
 ('it', 2035),
 ('for', 2034),
 ('on', 2033),
 ('from', 2032),
 ('you', 2031),
 ('edu', 2030),
 ('be', 2029),
 ('this', 2028),
 ('are', 2027),
 ('with', 2026),
 ('have', 2025),
 ('as', 2024),
 ('at', 2023),
 ('a', 2022),
 ('on', 2021),
 ('the', 2020),
 ('in', 2019),
 ('of', 2018),
 ('and', 2017),
 ('is', 2016),
 ('that', 2015),
 ('it', 2014),
 ('for', 2013),
 ('on', 2012),
 ('from', 2011),
 ('you', 2010),
 ('edu', 2009),
 ('be', 2008),
 ('this', 2007),
 ('are', 2006),
 ('with', 2005),
 ('have', 2004),
 ('as', 2003),
 ('at', 2002),
 ('a', 2001),
 ('on', 2000),
 ('the', 1999),
 ('in', 1998),
 ('of', 1997),
 ('and', 1996),
 ('is', 1995),
 ('that', 1994),
 ('it', 1993),
 ('for', 1992),
 ('on', 1991),
 ('from', 1990),
 ('you', 1989),
 ('edu', 1988),
 ('be', 1987),
 ('this', 1986),
 ('are', 1985),
 ('with', 1984),
 ('have', 1983),
 ('as', 1982),
 ('at', 1981),
 ('a', 1980),
 ('on', 1979),
 ('the', 1978),
 ('in', 1977),
 ('of', 1976),
 ('and', 1975),
 ('is', 1974),
 ('that', 1973),
 ('it', 1972),
 ('for', 1971),
 ('on', 1970),
 ('from', 1969),
 ('you', 1968),
 ('edu', 1967),
 ('be', 1966),
 ('this', 1965),
 ('are', 1964),
 ('with', 1963),
 ('have', 1962),
 ('as', 1961),
 ('at', 1960),
 ('a', 1959),
 ('on', 1958),
 ('the', 1957),
 ('in', 1956),
 ('of', 1955),
 ('and', 1954),
 ('is', 1953),
 ('that', 1952),
 ('it', 1951),
 ('for', 1950),
 ('on', 1949),
 ('from', 1948),
 ('you', 1947),
 ('edu', 1946),
 ('be', 1945),
 ('this', 1944),
 ('are', 1943),
 ('with', 1942),
 ('have', 1941),
 ('as', 1940),
 ('at', 1939),
 ('a', 1938),
 ('on', 1937),
 ('the', 1936),
 ('in', 1935),
 ('of', 1934),
 ('and', 1933),
 ('is', 1932),
 ('that', 1931),
 ('it', 1930),
 ('for', 1929),
 ('on', 1928),
 ('from', 1927),
 ('you', 1926),
 ('edu', 1925),
 ('be', 1924),
 ('this', 1923),
 ('are', 1922),
 ('with', 1921),
 ('have', 1920),
 ('as', 1919),
 ('at', 1918),
 ('a', 1917),
 ('on', 1916),
 ('the', 1915),
 ('in', 1914),
 ('of', 1913),
 ('and', 1912),
 ('is', 1911),
 ('that', 1910),
 ('it', 1909),
 ('for', 1908),
 ('on', 1907),
 ('from', 1906),
 ('you', 1905),
 ('edu', 1904),
 ('be', 1903),
 ('this', 1902),
 ('are', 1901),
 ('with', 1900),
 ('have', 1899),
 ('as', 1898),
 ('at', 1897),
 ('a', 1896),
 ('on', 1895),
 ('the', 1894),
 ('in', 1893),
 ('of', 1892),
 ('and', 1891),
 ('is', 1890),
 ('that', 1889),
 ('it', 1888),
 ('for', 1887),
 ('on', 1886),
 ('from', 1885),
 ('you', 1884),
 ('edu', 1883),
 ('be', 1882),
 ('this', 1881),
 ('are', 1880),
 ('with', 1879),
 ('have', 1878),
 ('as', 1877),
 ('at', 1876),
 ('a', 1875),
 ('on', 1874),
 ('the', 1873),
 ('in', 1872),
 ('of', 1871),
 ('and', 1870),
 ('is', 1869),
 ('that', 1868),
 ('it', 1867),
 ('for', 1866),
 ('on', 1865),
 ('from', 1864),
 ('you', 1863),
 ('edu', 1862),
 ('be', 1861),
 ('this', 1860),
 ('are', 1859),
 ('with', 1858),
 ('have', 1857),
 ('as', 1856),
 ('at', 1855),
 ('a', 1854),
 ('on', 1853),
 ('the', 1852),
 ('in', 1851),
 ('of', 1850),
 ('and', 1849),
 ('is', 1848),
 ('that', 1847),
 ('it', 1846),
 ('for', 1845),
 ('on', 1844),
 ('from', 1843),
 ('you', 1842),
 ('edu', 1841),
 ('be', 1840),
 ('this', 1839),
 ('are', 1838),
 ('with', 1837),
 ('have', 1836),
 ('as', 1835),
 ('at', 1834),
 ('a', 1833),
 ('on', 1832),
 ('the', 1831),
 ('in', 1830),
 ('of', 1829),
 ('and', 1828),
 ('is', 1827),
 ('that', 1826),
 ('it', 1825),
 ('for', 1824),
 ('on', 1823),
 ('from', 1822),
 ('you', 1821),
 ('edu', 1820),
 ('be', 1819),
 ('this', 1818),
 ('are', 1817),
 ('with', 1816),
 ('have', 1815),
 ('as', 1814),
 ('at', 1813),
 ('a', 1812),
 ('on', 1811),
 ('the', 1810),
 ('in', 1809),
 ('of', 1808),
 ('and', 1807),
 ('is', 1806),
 ('that', 1805),
 ('it', 1804),
 ('for', 1803),
 ('on', 1802),
 ('from', 1801),
 ('you', 1800),
 ('edu', 1799),
 ('be', 1798),
 ('this', 1797),
 ('are', 1796),
 ('with', 1795),
 ('have', 1794),
 ('as', 1793),
 ('at', 1792),
 ('a', 1791),
 ('on', 1
```

```
#macierz pomyłek
# 'sci.crypt', 'rec.sport.baseball', 'rec.motorcycles', 'sci.space'
from sklearn.metrics import confusion_matrix
confusion_matrix(test_dataset.target, predicted, labels=[0,1,2,3])

array([[389,  0,  7,  2],
       [ 1, 390,  3,  3],
       [ 1,  4, 391,  0],
       [ 1,  0, 20, 373]], dtype=int64)
```

Obraz 23: Macierz pomyłek dla wytrenowanego modelu

Na podstawie macierzy pomyłek widzimy, że model bardzo dobrze radzi sobie z dopasowywaniem kategorii do wiadomości; najbardziej myli „sci.crypt” z „rec.motorcycles” oraz „sci.space” z „rec.motorcycles”. Kategoria „rec.motorcycles” jest kategorią, którą najczęściej model błędnie wybiera.

5.2. Zadanie #4

i. Treść polecenia

Proszę zmienić użyty klasyfikator na drzewo decyzyjne i zbadać jakość wyników klasyfikacji.

ii. Rozwiązanie

```
from sklearn import tree
clf = Pipeline([
    ('vect', CountVectorizer(analyzer='word', ngram_range=(1, 1))),
    ('tfidf', TfidfTransformer()),
    ('clf', tree.DecisionTreeClassifier()),
])

#wytrenowanie pipeliny
clf.fit(train_dataset.data, train_dataset.target)

Pipeline(steps=[('vect', CountVectorizer()), ('tfidf', TfidfTransformer()),
                ('clf', DecisionTreeClassifier())])
```

Obraz 24: Trenowanie pipeliny z klasyfikatorem drzewa decyzyjnego - porównywanie na raz jednego słowa z jednym słowem

```
#proste testy wytrenowanego modelu
sentences = ['Asteroids are dangerous for Earth.',
             'I am dreaming about new motorcycle.',
             'I hate baseball games.',
             'Did you see my C implementations?']

predicted_categories = clf.predict(sentences)

for sentence, predicted_category in zip(sentences, predicted_categories):
    print('%r => %s' % (sentence, train_dataset.target_names[predicted_category]))

'Asteroids are dangerous for Earth.' => rec.sport.baseball
'I am dreaming about new motorcycle.' => rec.motorcycles
'I hate baseball games.' => rec.sport.baseball
'Did you see my C implementations?' => rec.sport.baseball
```

Obraz 25: Proste testy wytrenowanego modelu

Warto zauważyć, że przy zmianie klasyfikatora na drzewo decyzyjne model inaczej przypisuje kategorie już dla prostych testów wytrenowanego modelu - dwie z czterech wiadomości błędnie przypisuje do „rec.sport.baseball”.

```
#testy wytrenowanego modelu
import numpy as np
predicted = clf.predict(test_dataset.data)
np.mean(predicted == test_dataset.target)

0.8145110410094637
```

Obraz 26: Dokładność wytrenowanego modelu po zmianie klasyfikatora

Wyraźnie spada też dokładność wytrenowanego modelu.

```
#macierz pomyłek
# 'sci.crypt', 'rec.sport.baseball', 'rec.motorcycles', 'sci.space'
confusion_matrix(test_dataset.target, predicted, labels=[0,1,2,3])

array([[340, 28, 18, 12],
       [ 21, 329, 17, 30],
       [ 21, 38, 323, 14],
       [ 24, 45, 26, 299]], dtype=int64)
```

Obraz 27: Macierz pomyłek dla wytrenowanego modelu po zmianie klasyfikatora

Macierz pomyłek również wyraźnie się różni - model dużo częściej błędnie przypisuje kategorii - najczęściej myli się w wiadomościach z kategorii „sci.space”, a najczęściej błędnie przypisywaną kategorią jest „rec.sport.baseball”.

5.3. Zadanie #5

i. Treść polecenia

Proszę wytrenować klasyfikator na zbiorze testującym a przetestować na zbiorze trenującym.

ii. Rozwiązanie

```
test_dataset = fetch_20newsgroups(
    subset='train',
    categories=categories,
    shuffle=True,
    random_state=42)
train_dataset = fetch_20newsgroups(
    subset='test',
    categories=categories,
    shuffle=True,
    random_state=42)
#klasy
print(train_dataset.target_names)
print(test_dataset.target_names)
#liczba wiadomości
print(len(train_dataset.data))
print(len(test_dataset.data))

['rec.motorcycles', 'rec.sport.baseball', 'sci.crypt', 'sci.space']
['rec.motorcycles', 'rec.sport.baseball', 'sci.crypt', 'sci.space']
1585
2383
```

Obraz 28: Trenowanie modelu na zbiorze testującym a testowanie na zbiorze trenującym

```
[('the', 30530),
 ('to', 15205),
 ('of', 13501),
 ('and', 11909),
 ('in', 9861),
 ('is', 8380),
 ('that', 6952),
 ('it', 6349),
 ('for', 5829),
 ('on', 4651),
 ('from', 4554),
 ('you', 4510),
 ('edu', 4471),
 ('be', 4399),
 ('this', 4193),
 ('are', 3351),
 ('with', 3291),
 ('have', 3193),
 ('as', 3122),
```

Obraz 30: Trenowanie pipeline'u - porównywanie na raz jednego słowa z jednym słowem

```
#proste testy wytrenowanego modelu
sentences = ['Asteroids are dangerous for Earth.',
             'I am dreaming about new motorcycle.',
             'I hate baseball games.',
             'Did you see my C implementations?']

predicted_categories = clf.predict(sentences)

for sentence, predicted_category in zip(sentences, predicted_categories):
    print('%r => %s' % (sentence, train_dataset.target_names[predicted_category]))
```

```
'Asteroids are dangerous for Earth.' => sci.space
'I am dreaming about new motorcycle.' => rec.motorcycles
'I hate baseball games.' => rec.sport.baseball
'Did you see my C implementations?' => rec.motorcycles
```

Model przypisał błędnie tylko 1 wiadomość w prostych testach - zamiast „sci.crypt” oznaczył kategorię „rec.motorcycles”.

Obraz 32: Dokładność wytrenowanego modelu

Obraz 32: Dokładność wytrenowanego modelu

Model mimo wszystko posiada wysoką dokładność predykcji.

```
# 'sci.crypt', 'rec.sport.baseball', 'rec.motorcycles', 'sci.space'
confusion_matrix(test_dataset.target, predicted, labels=[0,1,2,3])

array([[576, 10, 7, 5],
       [ 5, 583, 8, 1],
       [ 1, 1, 592, 1],
       [ 2, 10, 21, 560]], dtype=int64)
```

Obraz 33: Macierz pomyłek wytrenowanego modelu

Model najczęściej myli „sci.space” z „rec.motorcycles” oraz „sci.crypt” z „rec.sport.baseball”. Kategoria „rec.motorcycles” jest najczęściej błędnie przypisywana.

5.4. Zadanie #6

i. Treść polecenia

Proszę ustawić w odpowiednim miejscu parametr `analyzer='char'` i sprawdzić jak zmieniają się wyniki.

ii. Rozwiązanie

```
count_vect = CountVectorizer(analyzer='char', ngram_range=(1, 1))
X_train_counts = count_vect.fit_transform(train_dataset.data)
word_list = count_vect.get_feature_names_out()
count_list = X_train_counts.toarray().sum(axis=0)
word_count = dict(zip(word_list, count_list))
sorted(word_count.items(), key=lambda x: x[1], reverse=True)
```

Obraz 34: Najczęściej pojawiające się znaki we wszystkich wiadomościach

```
[(' ', 611986),
 ('e', 356576),
 ('t', 274214),
 ('a', 244374),
 ('o', 235500),
 ('i', 226347),
 ('n', 216152),
 ('s', 208562),
 ('r', 192674),
 ('h', 130791),
 ('l', 128842),
 ('c', 112438),
 ('d', 107127),
 ('u', 93488),
 ('m', 85388),
 ('p', 73989),
 ('.', 70497),
 ('g', 65579),
 ('y', 62716),
 ('f', 60465),
 ('-', 58052),
 ('b', 55879),
 ('w', 52034),
 ('\n', 48768),
 ('v', 33854),
 (',', 32937),
 ('k', 29792),
 (':', 23327),
 ('1', 21961),
 ('>', 21654),
 ('0', 20665),
 ('"', 12544),
 (')', 12342),
 ('=', 12049),
 ('2', 11888),
 ('(', 11469),
 ('j', 11011),
 ('3', 10130),
 ('9', 9656),
 ('@', 9244),
 ('_', 9067),
 ('*', 8902),
 ('x', 8593),
 ('5', 8571),
 ('4', 7558),
 ('"', 7244),
 ('z', 6870),
 ('6', 6142),
 ('8', 6060),
 ('7', 5630),
 ('?', 5533),
 ('|', 5082),
 ('/', 4533),
 ('q', 3984),
 ('!', 3021),
 ('<', 2891),
 (';', 1978),
 ('~', 1898),
 ('\t', 1635),
 ('#', 1520),
 ('^', 1369),
 (']', 1287),
 ('[', 1273),
 ('+', 1206),
 ('$ ', 1172),
 ('`', 847),
 ('&', 820),
 ('\\', 720),
 ('%', 495),
 ('}', 257),
 ('{', 133),
 ('\xa0', 12),
 ('e', 9),
 ('e', 5),
 ('n', 5),
 ('c', 3),
 ('e', 3),
 ('\x1a', 1),
 ('e', 1)]
```



```

count_vect = CountVectorizer(analyzer='char', ngram_range=(5, 5))
X_train_counts = count_vect.fit_transform(train_dataset.data)
word_list = count_vect.get_feature_names()
count_list = X_train_counts.toarray().sum(axis=0)
word_count = dict(zip(word_list, count_list))
sorted(word_count.items(), key=lambda x: x[1], reverse=True)

```

```

[('----', 31948),
 (' the ', 27356),
 (' and ', 10438),
 ('====', 9864),
 ('ation', 7643),
 (' that', 6619),
 ('that ', 6043),
 ('____', 5571),
 ('*****', 5337),
 (' for ', 5231),
 ('n the', 4746),
 ('tion ', 4451),
 (' of t', 4195),
 ('f the', 4010),
 ('t the', 3967),
 (' this', 3902),
 ('of th', 3848),
 ('this ', 3710),
 (' with', 3511),
 ('----', 3402),
 ('----', 3392),
 ('----', 3382),
 ('----', 3372),
 ('----', 3362),
 ('----', 3352),
 ('----', 3342),
 ('----', 3332),
 ('----', 3322),
 ('----', 3312),
 ('----', 3302),
 ('----', 3292),
 ('----', 3282),
 ('----', 3272),
 ('----', 3262),
 ('----', 3252),
 ('----', 3242),
 ('----', 3232),
 ('----', 3222),
 ('----', 3212),
 ('----', 3202),
 ('----', 3192),
 ('----', 3182),
 ('----', 3172),
 ('----', 3162),
 ('----', 3152),
 ('----', 3142),
 ('----', 3132),
 ('----', 3122),
 ('----', 3112),
 ('----', 3102),
 ('----', 3092),
 ('----', 3082),
 ('----', 3072),
 ('----', 3062),
 ('----', 3052),
 ('----', 3042),
 ('----', 3032),
 ('----', 3022),
 ('----', 3012),
 ('----', 3002),
 ('----', 2992),
 ('----', 2982),
 ('----', 2972),
 ('----', 2962),
 ('----', 2952),
 ('----', 2942),
 ('----', 2932),
 ('----', 2922),
 ('----', 2912),
 ('----', 2902),
 ('----', 2892),
 ('----', 2882),
 ('----', 2872),
 ('----', 2862),
 ('----', 2852),
 ('----', 2842),
 ('----', 2832),
 ('----', 2822),
 ('----', 2812),
 ('----', 2802),
 ('----', 2792),
 ('----', 2782),
 ('----', 2772),
 ('----', 2762),
 ('----', 2752),
 ('----', 2742),
 ('----', 2732),
 ('----', 2722),
 ('----', 2712),
 ('----', 2702),
 ('----', 2692),
 ('----', 2682),
 ('----', 2672),
 ('----', 2662),
 ('----', 2652),
 ('----', 2642),
 ('----', 2632),
 ('----', 2622),
 ('----', 2612),
 ('----', 2602),
 ('----', 2592),
 ('----', 2582),
 ('----', 2572),
 ('----', 2562),
 ('----', 2552),
 ('----', 2542),
 ('----', 2532),
 ('----', 2522),
 ('----', 2512),
 ('----', 2502),
 ('----', 2492),
 ('----', 2482),
 ('----', 2472),
 ('----', 2462),
 ('----', 2452),
 ('----', 2442),
 ('----', 2432),
 ('----', 2422),
 ('----', 2412),
 ('----', 2402),
 ('----', 2392),
 ('----', 2382),
 ('----', 2372),
 ('----', 2362),
 ('----', 2352),
 ('----', 2342),
 ('----', 2332),
 ('----', 2322),
 ('----', 2312),
 ('----', 2302),
 ('----', 2292),
 ('----', 2282),
 ('----', 2272),
 ('----', 2262),
 ('----', 2252),
 ('----', 2242),
 ('----', 2232),
 ('----', 2222),
 ('----', 2212),
 ('----', 2202),
 ('----', 2192),
 ('----', 2182),
 ('----', 2172),
 ('----', 2162),
 ('----', 2152),
 ('----', 2142),
 ('----', 2132),
 ('----', 2122),
 ('----', 2112),
 ('----', 2102),
 ('----', 2092),
 ('----', 2082),
 ('----', 2072),
 ('----', 2062),
 ('----', 2052),
 ('----', 2042),
 ('----', 2032),
 ('----', 2022),
 ('----', 2012),
 ('----', 2002),
 ('----', 1992),
 ('----', 1982),
 ('----', 1972),
 ('----', 1962),
 ('----', 1952),
 ('----', 1942),
 ('----', 1932),
 ('----', 1922),
 ('----', 1912),
 ('----', 1902),
 ('----', 1892),
 ('----', 1882),
 ('----', 1872),
 ('----', 1862),
 ('----', 1852),
 ('----', 1842),
 ('----', 1832),
 ('----', 1822),
 ('----', 1812),
 ('----', 1802),
 ('----', 1792),
 ('----', 1782),
 ('----', 1772),
 ('----', 1762),
 ('----', 1752),
 ('----', 1742),
 ('----', 1732),
 ('----', 1722),
 ('----', 1712),
 ('----', 1702),
 ('----', 1692),
 ('----', 1682),
 ('----', 1672),
 ('----', 1662),
 ('----', 1652),
 ('----', 1642),
 ('----', 1632),
 ('----', 1622),
 ('----', 1612),
 ('----', 1602),
 ('----', 1592),
 ('----', 1582),
 ('----', 1572),
 ('----', 1562),
 ('----', 1552),
 ('----', 1542),
 ('----', 1532),
 ('----', 1522),
 ('----', 1512),
 ('----', 1502),
 ('----', 1492),
 ('----', 1482),
 ('----', 1472),
 ('----', 1462),
 ('----', 1452),
 ('----', 1442),
 ('----', 1432),
 ('----', 1422),
 ('----', 1412),
 ('----', 1402),
 ('----', 1392),
 ('----', 1382),
 ('----', 1372),
 ('----', 1362),
 ('----', 1352),
 ('----', 1342),
 ('----', 1332),
 ('----', 1322),
 ('----', 1312),
 ('----', 1302),
 ('----', 1292),
 ('----', 1282),
 ('----', 1272),
 ('----', 1262),
 ('----', 1252),
 ('----', 1242),
 ('----', 1232),
 ('----', 1222),
 ('----', 1212),
 ('----', 1202),
 ('----', 1192),
 ('----', 1182),
 ('----', 1172),
 ('----', 1162),
 ('----', 1152),
 ('----', 1142),
 ('----', 1132),
 ('----', 1122),
 ('----', 1112),
 ('----', 1102),
 ('----', 1092),
 ('----', 1082),
 ('----', 1072),
 ('----', 1062),
 ('----', 1052),
 ('----', 1042),
 ('----', 1032),
 ('----', 1022),
 ('----', 1012),
 ('----', 1002),
 ('----', 992),
 ('----', 982),
 ('----', 972),
 ('----', 962),
 ('----', 952),
 ('----', 942),
 ('----', 932),
 ('----', 922),
 ('----', 912),
 ('----', 902),
 ('----', 892),
 ('----', 882),
 ('----', 872),
 ('----', 862),
 ('----', 852),
 ('----', 842),
 ('----', 832),
 ('----', 822),
 ('----', 812),
 ('----', 802),
 ('----', 792),
 ('----', 782),
 ('----', 772),
 ('----', 762),
 ('----', 752),
 ('----', 742),
 ('----', 732),
 ('----', 722),
 ('----', 712),
 ('----', 702),
 ('----', 692),
 ('----', 682),
 ('----', 672),
 ('----', 662),
 ('----', 652),
 ('----', 642),
 ('----', 632),
 ('----', 622),
 ('----', 612),
 ('----', 602),
 ('----', 592),
 ('----', 582),
 ('----', 572),
 ('----', 562),
 ('----', 552),
 ('----', 542),
 ('----', 532),
 ('----', 522),
 ('----', 512),
 ('----', 502),
 ('----', 492),
 ('----', 482),
 ('----', 472),
 ('----', 462),
 ('----', 452),
 ('----', 442),
 ('----', 432),
 ('----', 422),
 ('----', 412),
 ('----', 402),
 ('----', 392),
 ('----', 382),
 ('----', 372),
 ('----', 362),
 ('----', 352),
 ('----', 342),
 ('----', 332),
 ('----', 322),
 ('----', 312),
 ('----', 302),
 ('----', 292),
 ('----', 282),
 ('----', 272),
 ('----', 262),
 ('----', 252),
 ('----', 242),
 ('----', 232),
 ('----', 222),
 ('----', 212),
 ('----', 202),
 ('----', 192),
 ('----', 182),
 ('----', 172),
 ('----', 162),
 ('----', 152),
 ('----', 142),
 ('----', 132),
 ('----', 122),
 ('----', 112),
 ('----', 102),
 ('----', 92),
 ('----', 82),
 ('----', 72),
 ('----', 62),
 ('----', 52),
 ('----', 42),
 ('----', 32),
 ('----', 22),
 ('----', 12),
 ('----', 2),
 ('----', 0)]

```

Obraz 35: Najczęściej pojawiające się 5 znakowe ciągi we wszystkich wiadomościach

```

clf = Pipeline([
    ('vect', CountVectorizer(analyzer='char', ngram_range=(1, 1))),
    ('tfidf', TfidfTransformer()),
    ('clf', MultinomialNB()),
])
#wytrenowanie pipeline
clf.fit(train_dataset.data, train_dataset.target)

Pipeline(steps=[('vect', CountVectorizer(analyzer='char')),
                 ('tfidf', TfidfTransformer()), ('clf', MultinomialNB())])

#proste testy wytrenowanego modelu
sentences = ['Asteroids are dangerous for Earth.',
             'I am dreaming about new motorcycle.',
             'I hate baseball games.',
             'Did you see my C implementations?']

predicted_categories = clf.predict(sentences)

for sentence, predicted_category in zip(sentences, predicted_categories):
    print('%r => %s' % (sentence, train_dataset.target_names[predicted_category]))

'Asteroids are dangerous for Earth.' => sci.crypt
'I am dreaming about new motorcycle.' => sci.crypt
'I hate baseball games.' => rec.sport.baseball
'Did you see my C implementations?' => sci.crypt

```

Obraz 36: Wytrenowanie pipeline dla porównania pojedynczych znaków

Przy porównaniu pojedynczych znaków model w prostych testach przypisał prawie wszystkie wiadomości do kategorii „sci.crypt”.

```

predicted = clf.predict(test_dataset.data)
np.mean(predicted == test_dataset.target)

0.3747634069400631

# 'sci.crypt', 'rec.sport.baseball', 'rec.motorcycles', 'sci.space'
confusion_matrix(test_dataset.target, predicted, labels=[0,1,2,3])

array([[116, 90, 161, 31],
       [ 31, 119, 231, 16],
       [ 29, 32, 296, 39],
       [ 27, 56, 248, 63]], dtype=int64)

```

Obraz 37: Dokładność i macierz pomyłek dla modelu porównującego pojedyncze znaki

Dokładność modelu jest bardzo niska. Model najczęściej myli „sci.crypt” z „rec.motorcycles” a kategorię „sci.space” prawie zawsze oznacza jako „rec.motorcycles”.

```
clf = Pipeline([
    ('vect', CountVectorizer(analyzer='char', ngram_range=(5, 5))),
    ('tfidf', TfidfTransformer()),
    ('clf', MultinomialNB()),
])
#wytrenowanie pipeline'u
clf.fit(train_dataset.data, train_dataset.target)

Pipeline(steps=[('vect', CountVectorizer(analyzer='char', ngram_range=(5, 5))),
                 ('tfidf', TfidfTransformer()), ('clf', MultinomialNB())])

#proste testy wytrenowanego modelu
sentences = ['Asteroids are dangerous for Earth.',
             'I am dreaming about new motorcycle.',
             'I hate baseball games.',
             'Did you see my C implementations?']

predicted_categories = clf.predict(sentences)

for sentence, predicted_category in zip(sentences, predicted_categories):
    print('%r => %s' % (sentence, train_dataset.target_names[predicted_category]))

'Asteroids are dangerous for Earth.' => sci.space
'I am dreaming about new motorcycle.' => rec.motorcycles
'I hate baseball games.' => rec.sport.baseball
'Did you see my C implementations?' => sci.crypt
```

Obraz 38: Trenowanie pipeline'u w modelu porównującym 5 znakowe ciągi

Model przypisał poprawnie wszystkie wiadomości w prostych testach.

```
predicted = clf.predict(test_dataset.data)
np.mean(predicted == test_dataset.target)

0.8845425867507887

# 'sci.crypt', 'rec.sport.baseball', 'rec.motorcycles', 'sci.space'
confusion_matrix(test_dataset.target, predicted, labels=[0,1,2,3])

array([[327,  4,  64,  3],
       [ 16, 344,  36,  1],
       [  6,  2, 385,  3],
       [  2,  1,  45, 346]], dtype=int64)
```

Obraz 39: Dokładność i macierz pomyłek dla modelu porównującego 5 znakowe ciągi

Model ma dość wysoką dokładność, jednak nie tak wysoką jak model porównujący ze sobą pojedyncze znaki. Najczęściej błędnie przypisuje „sci.crypt” jako „rec.motorcycles” oraz „sci.space” jako „rec.motorcycles”.