

Systemy Sztucznej Inteligencji

dokumentacja projektu

Obliczanie prawdopodobieństwa wystąpienia udaru na podstawie danych pacjenta

Wiktor Machoń, Piotr Kołodziejski

29 maja 2022

Część I

Opis programu

Celem projektu jest zbudowanie klasyfikatora w oparciu o Naiwny Klasyfikator Bayesa. Klasyfikator ten ma za zadanie ocenić czy dany pacjent, ze swoimi parametrami jest w grupie zagrożonej dostaniem udaru.

Instrukcja obsługi

Program został napisany w języku programowania Python przy użyciu Jupyter notebook. Program (plik z rozszerzeniem .ipynb) należy uruchomić z poziomu aplikacji Jupyter Notebook. Plik .csv z danymi powinien znajdować się w tym samym folderze, co plik rozwiązania.

Dodatkowe informacje

Aby klasyfikator uznać za działający, jego trafność predykcji musi być na poziomie co najmniej 85 procent.

Część II

Opis działania

Naiwny klasyfikator Bayesa tworzymy w oparciu o Twierdzenie Bayesa, które opisujemy poniższym równaniem:

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)} \quad (1)$$

W naszym programie używamy zmiennych X (dla danych wejściowych) oraz y (dla danych wyjściowych), więc nasz wzór przyjmuje postać:

$$P(y|X) = \frac{P(X|y) * P(y)}{P(X)} \quad (2)$$

Przyjmujemy naiwne założenie, że zmienne są od siebie niezależne, co implikuje, że:

$$P(X|y) = P(x_1|y) * P(x_2|y) * \dots * P(x_n|y) \quad (3)$$

Prawdopodobieństwo obliczamy dla y , więc $P(X)$ jest stałą, co oznacza, że możemy usunąć ją z równania otrzymując przybliżenie:

$$P(X|y) = P(y) * \prod_{i=1}^n P(x_i|y) \quad (4)$$

Korzystając z tego wzoru musimy wyznaczyć maksymalną wartość y .

Algorytm

Implementacja Naiwnego Klasyfikatora Bayesa w oparciu o Twierdzenie Bayesa, opisane wyżej. Dane z pliku .csv wczytujemy korzystając w biblioteki pandas, tworząc tym samym dataframe. Zamieniamy wartości tekstowe na liczbowe oraz dzielimy zbiór na treningowy i testowy, tak że ten pierwszy to 70 procent zbioru. Przed podziałem policzone zostały średnie i odchylenia stadardowe dla każdej klasy.

Data: Dane wejściowe dataframe X

Result: Klasy kolejnych wierszy - lista Predictions

Wczytaj dataframe. $Predictions = [...]$;

foreach $i \in X.index$ **do**

 Zainicjuj listę $ClassLikelihood$;

 Przypisz obecny wiersz jako obiekt do zmiennej $instance = X.loc[i]$;

foreach $cls \in classes$ **do**

 Utwórz listę $FeatureLikelihoods$;

 Dodaj do listy $FeatureLikelihoods$ $\log(prior[cls])$;

foreach $col \in xTrain.columns$ **do**

$data = instance[col]$;

$mean = means[col].loc[cls]$;

$variance = var[col].loc[cls]$;

 Obliczyć gęstość prawdopodobieństwa

$Likelihood = Normal(data, mean, variance)$

if $Likelihood \neq 0$ **then**

$Likelihood = np.log(Likelihood)$;

else

$Likelihood = 1/len(xTrain)$;

end

end

$TotalLikelihood = sum(FeatureLikelihoods)$;

$ClassLikelihood.append(TotalLikelihood)$;

end

$MaxIndex = ClassLikelihood.index(max(ClassLikelihood))$;

$Prediction = classes[MaxIndex]$;

$Predictions.append(Prediction)$;

end

Algorithm 1: Algorytm drukowania informacji o liczbie parzystej/nieparzystej.

Algorytm odrzucania najmniej istotnych kolumn

Data: Dane wejściowe dataframe y , X

Result: Nazwy dwóch kolumn, których odrzucenie będzie skutkowało zwiększeniem dokładności

Utwórz pusty słownik $colsAndAcc$;

foreach $col1 \in X.columns$ **do**

$x1 = X.drop(col1, axis = 1)$;

foreach $col2 \in x1.columns$ **do**

$x2 = x1.drop(col2, axis = 1)$;

$prediction = Predict(x2)$;

$colsAndAcc[str(col1 + ' and ' + col2)] = round(Accuracy(y, prediction), 5)$;

end

end

Zwróć maksymalną wartość ze słownika $colsAndAcc$ oraz jej klucz

Algorithm 2: Algorytm odrzucania najmniej istotnych kolumn

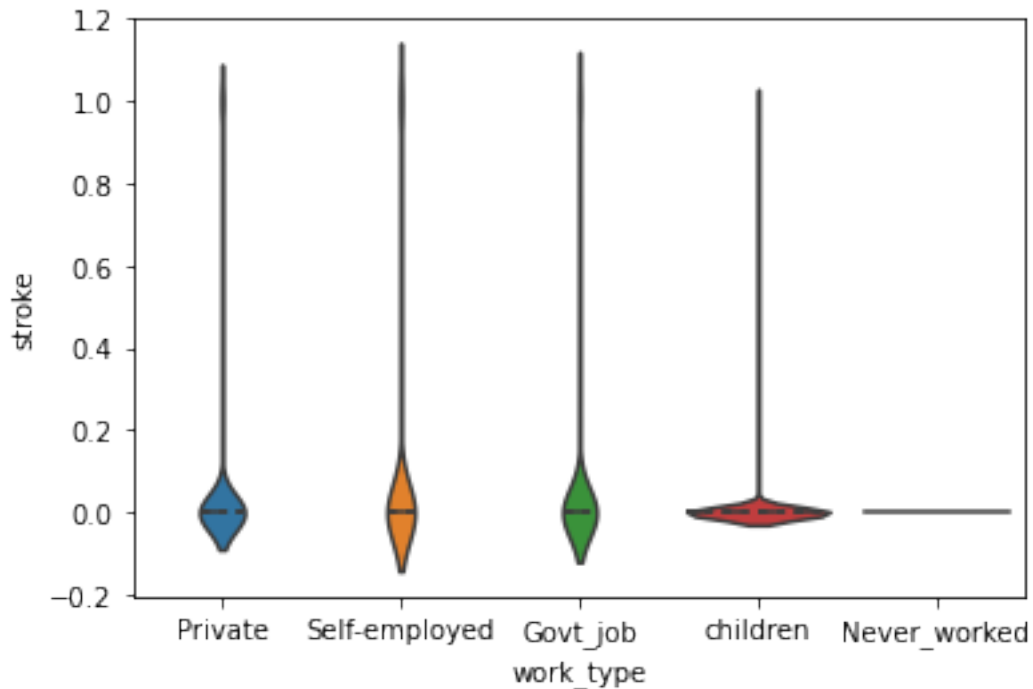
Implementacja

W pliku z danymi występowały dane tekstowe, więc żeby móc je wykorzystać zostały one przekształcone w dane numeryczne. Gdy dla danej klasy występowały tylko dwie wartości to zostały one przekształcone na wartości 0 lub 1. Gdy wartości było więcej to dodawano kolejne liczby naturalne. Jeżeli nie ma konkretnej informacji (jak np. *Other* w klasie *gender*) wartości takiej zostało przypisane -1.

```
1  def ToNumData(X):
2      for col in X.columns:
3          if col == "gender":
4              for row_nr in range(len(X[col])):
5                  if(X[col].iloc[row_nr] == "Male"):
6                      X[col].iloc[row_nr] = 0
7                  elif (X[col].iloc[row_nr] == "Female"):
8                      X[col].iloc[row_nr] = 1
9                  elif (X[col].iloc[row_nr] == "Other"):
10                     X[col].iloc[row_nr] = -1
11             elif col == "ever_married":
12                 for row_nr in range(len(X[col])):
13                     if(X[col].iloc[row_nr] == "Yes"):
14                         X[col].iloc[row_nr] = 1
15                     elif (X[col].iloc[row_nr] == "No"):
16                         X[col].iloc[row_nr] = 0
17             elif col == "Residence_type":
18                 for row_nr in range(len(X[col])):
19                     if(X[col].iloc[row_nr] == "Rural"):
20                         X[col].iloc[row_nr] = 1
21                     elif (X[col].iloc[row_nr] == "Urban"):
22                         X[col].iloc[row_nr] = 0
23             elif col == "smoking_status":
24                 for row_nr in range(len(X[col])):
25                     if(X[col].iloc[row_nr] == "formerly smoked"):
26                         X[col].iloc[row_nr] = 1
27                     elif (X[col].iloc[row_nr] == "never smoked"):
28                         X[col].iloc[row_nr] = 0
29                     elif (X[col].iloc[row_nr] == "smokes"):
30                         X[col].iloc[row_nr] = 2
31                     elif (X[col].iloc[row_nr] == "Unknown"):
32                         X[col].iloc[row_nr] = -1
33             elif col == "work_type":
34                 for row_nr in range(len(X[col])):
35                     if(X[col].iloc[row_nr] == "Never_worked"):
36                         X[col].iloc[row_nr] = 0
37                     elif (X[col].iloc[row_nr] == "Govt_job"):
38                         X[col].iloc[row_nr] = 1
39                     elif (X[col].iloc[row_nr] == "Private"):
40                         X[col].iloc[row_nr] = 2
41                     elif (X[col].iloc[row_nr] == "children"):
42                         X[col].iloc[row_nr] = -1
43                     elif (X[col].iloc[row_nr] == "Self-employed"):
44                         X[col].iloc[row_nr] = 3
45
46     return X
```

Testy

Sprawdzenie jak różne formy zatrudnienia mają się do możliwości wystąpienia udaru.



Rysunek 1: Forma zatrudnienia a możliwość wystąpienia udaru

Eksperymenty

W ramach sprawdzenia i głębszej analizy danych sprawdzono jakie kolumny mają wpływ na wynik klasyfikacji. Gdy kolumna ta nie miała wpływu została odrzucana. W ten sposób efektywność wzrosła z 86 do 92 procent.

```
1
2 def colSelection(y,X):
3     i = 1
4     cols_and_acc = dict()
5     for col1 in X.columns:
6         x1 = X.drop(col1, axis = 1)
7         for col2 in x1.columns:
8             x2 = x1.drop(col2, axis = 1)
9             prediction = Predict(x2)
10            cols_and_acc[str(col1 + " and " + col2)] = round(Accuracy(y,
11                                                                prediction), 5)
12            i+=1
13 print("Accuracy without columns "+max(cols_and_acc, key=cols_and_acc
14 .get)+" : "+str(cols_and_acc[max(cols_and_acc, key=cols_and_acc.
15 get)]))
```

```
colSelection(yVal,xVal)
```

Accuracy without columns hypertension and avg_glucose_level: 0.92532

Rysunek 2: Po odrzuceniu powyższych kolumn klasyfikator ma wyższą skuteczność

Pełen kod aplikacji

```
1 import numpy as np
2 import random as rd
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib as plt
6
7 strokes = pd.read_csv('healthcare-dataset-stroke-data.csv')
8 strokes.info()
9 strokes.describe()
10 strokes.head()
11
12 sns.violinplot(y='stroke', x='work_type', data=strokes, inner='quartile
    ')
13
14 sns.pairplot(strokes, hue='stroke', markers='+')
15
16 def ToNumData(X):
17     for col in X.columns:
18         if col == "gender":
19             for row_nr in range(len(X[col])):
20                 if(X[col].iloc[row_nr] == "Male"):
21                     X[col].iloc[row_nr] = 0
22                 elif (X[col].iloc[row_nr] == "Female"):
23                     X[col].iloc[row_nr] = 1
24                 elif (X[col].iloc[row_nr] == "Other"):
25                     X[col].iloc[row_nr] = -1
26
27         elif col == "ever_married":
28             for row_nr in range(len(X[col])):
29                 if(X[col].iloc[row_nr] == "Yes"):
30                     X[col].iloc[row_nr] = 1
31                 elif (X[col].iloc[row_nr] == "No"):
32                     X[col].iloc[row_nr] = 0
33
34         elif col == "Residence_type":
35             for row_nr in range(len(X[col])):
36                 if(X[col].iloc[row_nr] == "Rural"):
37                     X[col].iloc[row_nr] = 1
38                 elif (X[col].iloc[row_nr] == "Urban"):
39                     X[col].iloc[row_nr] = 0
40
41         elif col == "smoking_status":
```



```

42         for row_nr in range(len(X[col])):
43             if(X[col].iloc[row_nr] == "formerly smoked"):
44                 X[col].iloc[row_nr] = 1
45             elif (X[col].iloc[row_nr] == "never smoked"):
46                 X[col].iloc[row_nr] = 0
47             elif (X[col].iloc[row_nr] == "smokes"):
48                 X[col].iloc[row_nr] = 2
49             elif (X[col].iloc[row_nr] == "Unknown"):
50                 X[col].iloc[row_nr] = -1
51
52         elif col == "work_type":
53             for row_nr in range(len(X[col])):
54                 if(X[col].iloc[row_nr] == "Never_worked"):
55                     X[col].iloc[row_nr] = 0
56                 elif (X[col].iloc[row_nr] == "Govt_job"):
57                     X[col].iloc[row_nr] = 1
58                 elif (X[col].iloc[row_nr] == "Private"):
59                     X[col].iloc[row_nr] = 2
60                 elif (X[col].iloc[row_nr] == "children"):
61                     X[col].iloc[row_nr] = -1
62                 elif (X[col].iloc[row_nr] == "Self-employed"):
63                     X[col].iloc[row_nr] = 3
64
65     return X
66
67 #Dropping rows where bmi is NaN
68 strokes.dropna(subset= ["bmi"], inplace=True)
69
70 #Changing non-numeric data to numeric types
71 strokes = ToNumData(strokes)
72 strokes["gender"] = pd.to_numeric(strokes["gender"])
73 strokes["ever_married"] = pd.to_numeric(strokes["ever_married"])
74 strokes["Residence_type"] = pd.to_numeric(strokes["Residence_type"])
75 strokes["smoking_status"] = pd.to_numeric(strokes["smoking_status"])
76 strokes["work_type"] = pd.to_numeric(strokes["work_type"])
77
78 strokesTrain = strokes.sample(frac=0.7, random_state=1)
79 strokesVal = strokes.drop(strokesTrain.index)
80
81 yTrain = strokesTrain["stroke"]
82 xTrain = strokesTrain.drop("stroke", axis = 1).drop("id", axis = 1)
83
84 yVal = strokesVal["stroke"]
85 xVal = strokesVal.drop("stroke", axis = 1).drop("id", axis = 1)
86
87 means = strokesTrain.groupby(["stroke"]).mean()
88 var = strokesTrain.groupby(["stroke"]).var()
89 classes = np.unique(strokesTrain["stroke"].tolist())
90 prior = (strokesTrain.groupby(["stroke"]).count()/len(strokesTrain)).
91         iloc[:,0]
92
93 def Normal(n, mu, var):
94     sd = np.sqrt(var)
95     pdf = (np.e ** (-0.5 * ((n - mu)/sd) ** 2)) / (sd * np.sqrt(2 * np.
96         pi))

```

```

95     return pdf # pdf - probability density function
96
97 def Predict(X):
98     Predictions = []
99
100    for i in X.index: # Loop through each instances
101        ClassLikelihood = []
102        instance = X.loc[i]
103
104        for cls in classes: # Loop through each class
105
106            FeatureLikelihoods = []
107            FeatureLikelihoods.append(np.log(prior[cls])) # Append log
108                prior of class 'cls'
109
110            for col in X.columns: # Loop through each feature
111
112                data = instance[col]
113
114                mean = means[col].loc[cls] # Find the mean of column '
115                    col' that are in class 'cls'
116                variance = var[col].loc[cls] # Find the variance of
117                    column 'col' that are in class 'cls'
118
119                Likelihood = Normal(data, mean, variance)
120
121                if Likelihood != 0:
122                    Likelihood = np.log(Likelihood) # Find the log-
123                        likelihood evaluated at x
124                else:
125                    Likelihood = 1/len(X)
126
127                FeatureLikelihoods.append(Likelihood)
128
129            TotalLikelihood = sum(FeatureLikelihoods) # Calculate
130                posterior
131            ClassLikelihood.append(TotalLikelihood)
132
133            MaxIndex = ClassLikelihood.index(max(ClassLikelihood)) # Find
134                largest posterior position
135            Prediction = classes[MaxIndex]
136            Predictions.append(Prediction)
137
138    return Predictions
139
140 PredictTrain = Predict(xTrain)
141 PredictVal = Predict(xVal)
142
143 def Accuracy(y, prediction):
144
145     # Function to calculate accuracy
146     y = list(y)
147     prediction = list(prediction)
148     score = 0

```

```

144     for i, j in zip(y, prediction):
145         if i == j:
146             score += 1
147
148     return score / len(y)
149
150 round(Accuracy(yTrain, PredictTrain), 5)
151
152 round(Accuracy(yVal, PredictVal), 5)
153
154 def ColumnsSelection(y,X):
155     for col in X.columns:
156         x = X.drop(col, axis = 1)
157         prediction = Predict(x)
158         print(f"Accuracy without column {col}: {str(round(Accuracy(y,
159                                     prediction)*100, 5))}%")
160
161 ColumnsSelection(yTrain,xTrain)
162
163 ColumnsSelection(yVal,xVal)
164
165 def colSelection(y,X):
166     i = 1
167     cols_and_acc = dict()
168     for col1 in X.columns:
169         x1 = X.drop(col1, axis = 1)
170         for col2 in x1.columns:
171             x2 = x1.drop(col2, axis = 1)
172             prediction = Predict(x2)
173             cols_and_acc[str(col1 + " and " + col2)] = round(Accuracy(y,
174                                     prediction), 5)
175             i+=1
176
177     print(f"Accuracy without columns {max(cols_and_acc, key=cols_and_acc
178                                     .get)}: {str(cols_and_acc[max(cols_and_acc, key=cols_and_acc.get)
179                                     ])}")
180
181 colSelection(yVal,xVal)

```
