

Slovenská technická univerzita v Bratislave
Fakulta elektrotechniky a informatiky

Neurónové siete

Zadanie č.1

Obsah

| | | |
|---|--|----|
| 1 | Načítanie a spracovanie dát | 3 |
| 2 | Zhodnotenie náročnosti problému | 6 |
| 3 | Trénovanie, vyhodnotenie siete a nájdenie dobrých parametrov | 6 |
| 4 | Grid search | 12 |
| 5 | Tabulka experimentov | 13 |

Zoznam obrázkov

| | | |
|----|--|----|
| 1 | Stlpec pH pred škálovaním. | 4 |
| 2 | Stlpec pH po škálovaní. | 4 |
| 3 | Zobrazenie pitnosti vody v našom datasete. | 5 |
| 4 | DummyClassifier strategy. | 5 |
| 5 | Klasifikácia pomocou log. regresie. | 6 |
| 6 | MLPClassifier klasifikačný report. | 7 |
| 7 | Konfuzna matica pre tréningový set pri solveri adam. | 8 |
| 8 | Konfuzna matica pre tréningový set pri solveri SGD. | 9 |
| 9 | Konfuzna matica pre validačný set pri solveri adam. | 10 |
| 10 | Konfuzna matica pre validačný set pri solveri SGD. | 11 |
| 11 | Tabulka experimentov. | 13 |
| 12 | Konfúzna matica pre náš najlepší výsledok. | 13 |
| 13 | Priebeh trénovania pre náš najlepší výsledok. | 14 |

1 Načítanie a spracovanie dát

Prvým krokom je načítanie údajov pomocou knižnice **pandas** do nášho DataFrame-u. Škálovanie dát sme uskutočnili pomocou knižnice **MinMaxScaler**.

Priemerná hodnota pred normalizáciou - water train.

| | | |
|----|-----------------|--------------|
| 1 | ph | 7.041691 |
| 2 | Hardness | 196.310660 |
| 3 | Solids | 21756.102873 |
| 4 | Chloramines | 7.163238 |
| 5 | Sulfate | 335.505887 |
| 6 | Conductivity | 424.118360 |
| 7 | Organic_carbon | 14.260028 |
| 8 | Trihalomethanes | 66.115472 |
| 9 | Turbidity | 4.027384 |
| 10 | Potability | 0.390244 |

Št. odchýlka pred normalizáciou - water train.

| | | |
|----|-----------------|-------------|
| 1 | ph | 1.595465 |
| 2 | Hardness | 34.804883 |
| 3 | Solids | 8796.578151 |
| 4 | Chloramines | 1.621486 |
| 5 | Sulfate | 41.991894 |
| 6 | Conductivity | 80.942866 |
| 7 | Organic_carbon | 3.257755 |
| 8 | Trihalomethanes | 17.125366 |
| 9 | Turbidity | 0.794728 |
| 10 | Potability | 0.488177 |

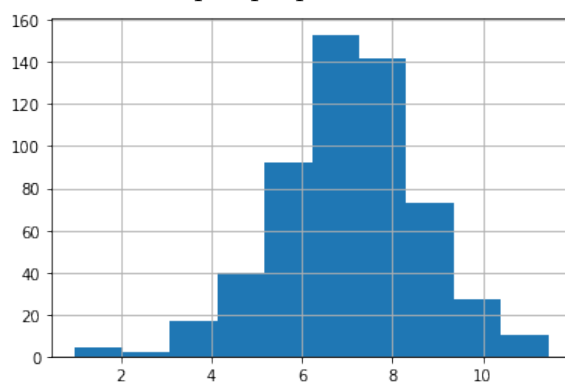
Priemerná hodnota PO normalizácií - water train.

| | | |
|----|-----------------|----------|
| 1 | ph | 0.578573 |
| 2 | Hardness | 0.474704 |
| 3 | Solids | 0.351937 |
| 4 | Chloramines | 0.441357 |
| 5 | Sulfate | 0.514041 |
| 6 | Conductivity | 0.435472 |
| 7 | Organic_carbon | 0.467340 |
| 8 | Trihalomethanes | 0.548045 |
| 9 | Turbidity | 0.468048 |
| 10 | Potability | 0.390244 |

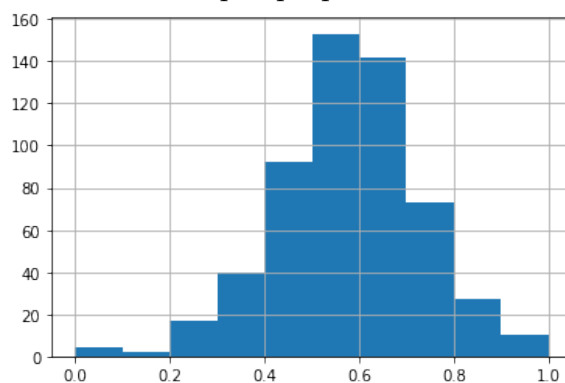
Št. odchýlka PO normalizácií - water train.

| | | |
|----|-----------------|----------|
| 1 | ph | 0.152533 |
| 2 | Hardness | 0.144171 |
| 3 | Solids | 0.144428 |
| 4 | Chloramines | 0.157565 |
| 5 | Sulfate | 0.145519 |
| 6 | Conductivity | 0.200047 |
| 7 | Organic_carbon | 0.165337 |
| 8 | Trihalomethanes | 0.143558 |
| 9 | Turbidity | 0.155906 |
| 10 | Potability | 0.488177 |

Obr. 1: Stlpec ph pred škalovaním.

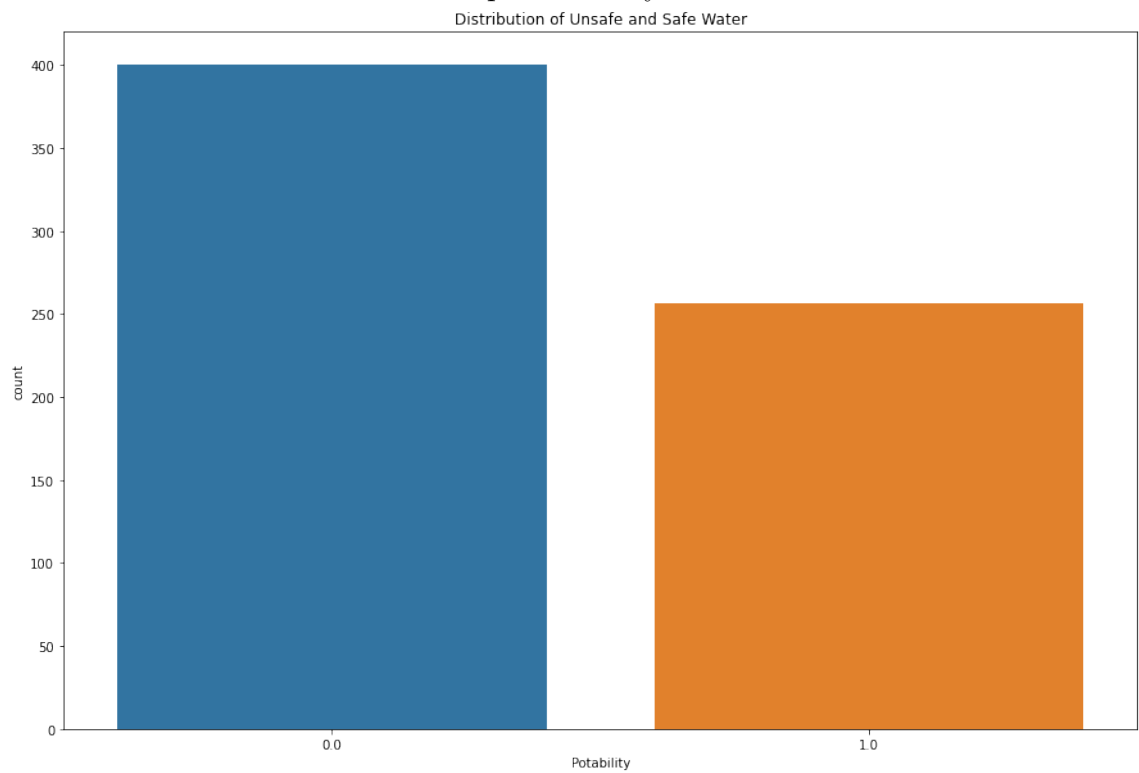


Obr. 2: Stlpec ph po škalovaní.



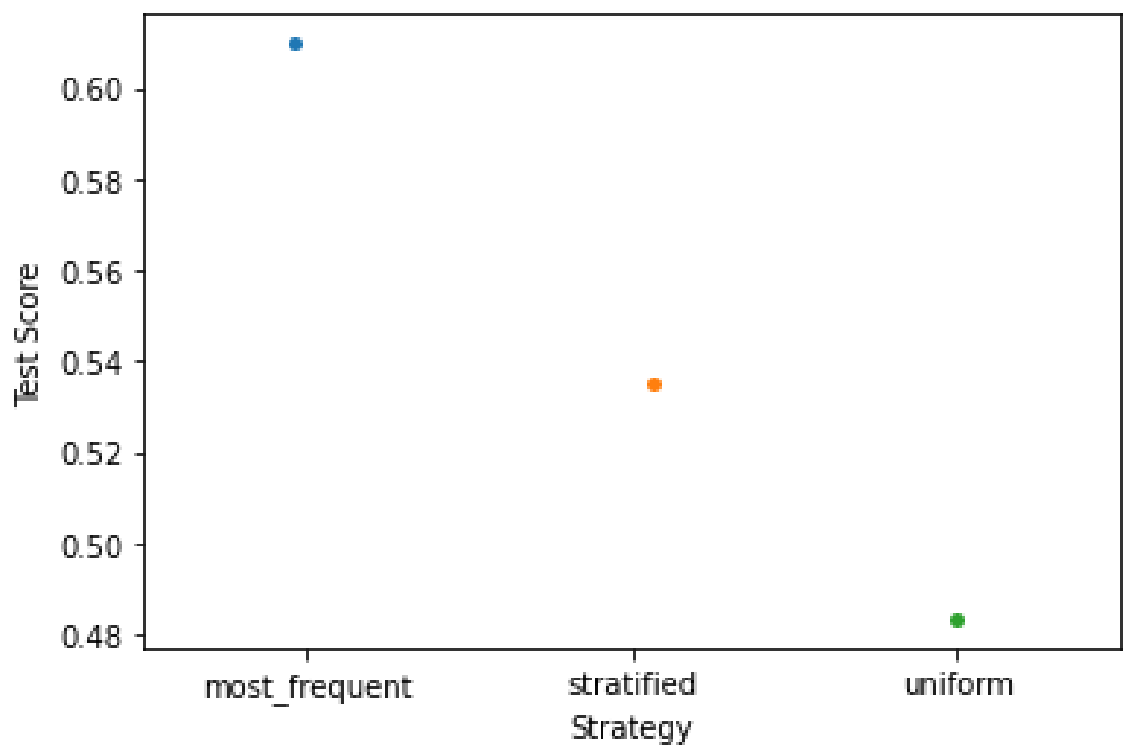
Výsledky máme rovnake avšak, práve po normalizácií (od 0 po 1) dát nám umožňuje modelu konvergovať k lepším hmotnostiam a následne vedie k presnejšiemu modelu.

Obr. 3: Zobrazenie pitnosť vody v našom datasete.



Podľa tohto grafu vieme určiť, že náš dataset obsahuje väčšinu nepitnej vody (označené číslom 0).

Obr. 4: DummyClassifier strategy.

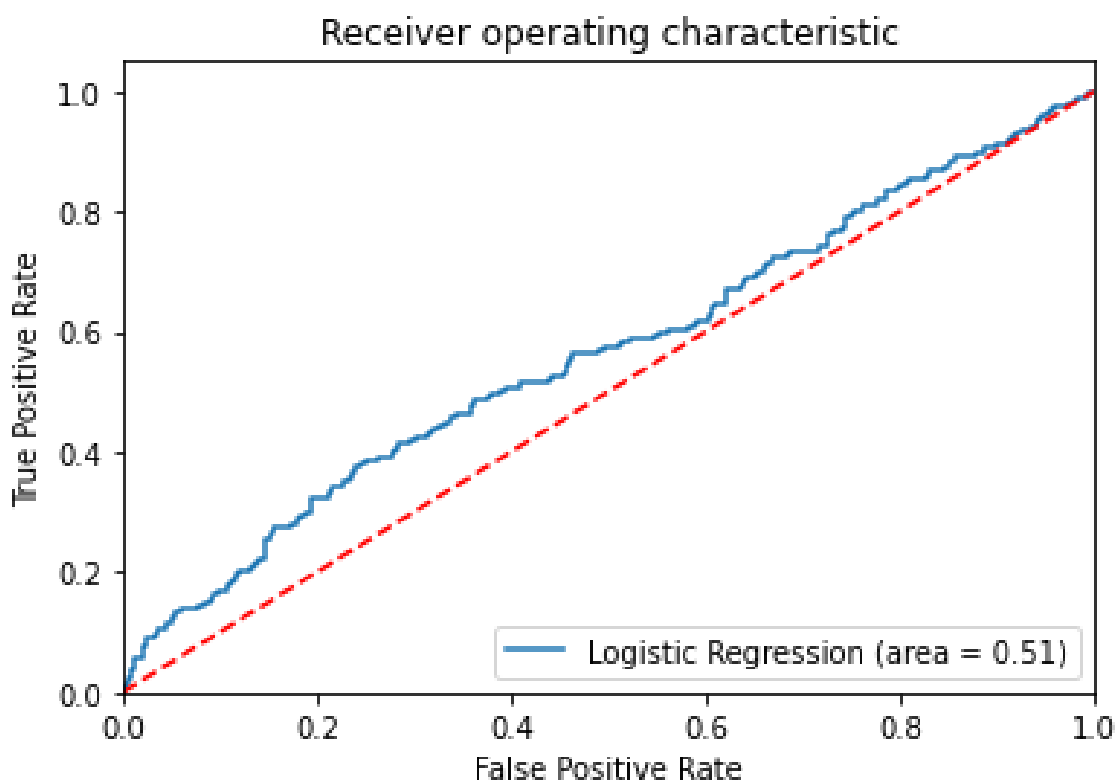


2 Zhodnotenie náročnosti problému

Pri stratégií najfrekvencovanejši klasifikátor vždy predpovedá najčastejšie označenie triedy v tréningových údajoch to nám vyšlo 60.1 percent. Pri stratégií stratifikovaný nám vyšla presnosť približne 53.5 percent, ktorý generuje náhodné predpovede rešpektovaním distribúcie tried tréningovej množiny. Pri stratégií uniform nám generuje predpovede rovnomerne náhodne kde nám vyšlo 48.3 percent.

Klasifikácia pomocou logistickej regresie na testovacej množine nám presnosť vyšla 62 percent.

Obr. 5: Klasifikácia pomocou log. regiesie.



Bodkovaná čiara predstavuje ROC krivku čisto náhodného klasifikátora; dobrý klasifikátor zostáva čo najďalej od tejto čiary (smerom k ľavému hornému rohu).

3 Trénovanie, vyhodnotenie siete a nájdenie dobrých parametrov

```
1 x_train, x_val, y_train, y_val = train_test_split(
```



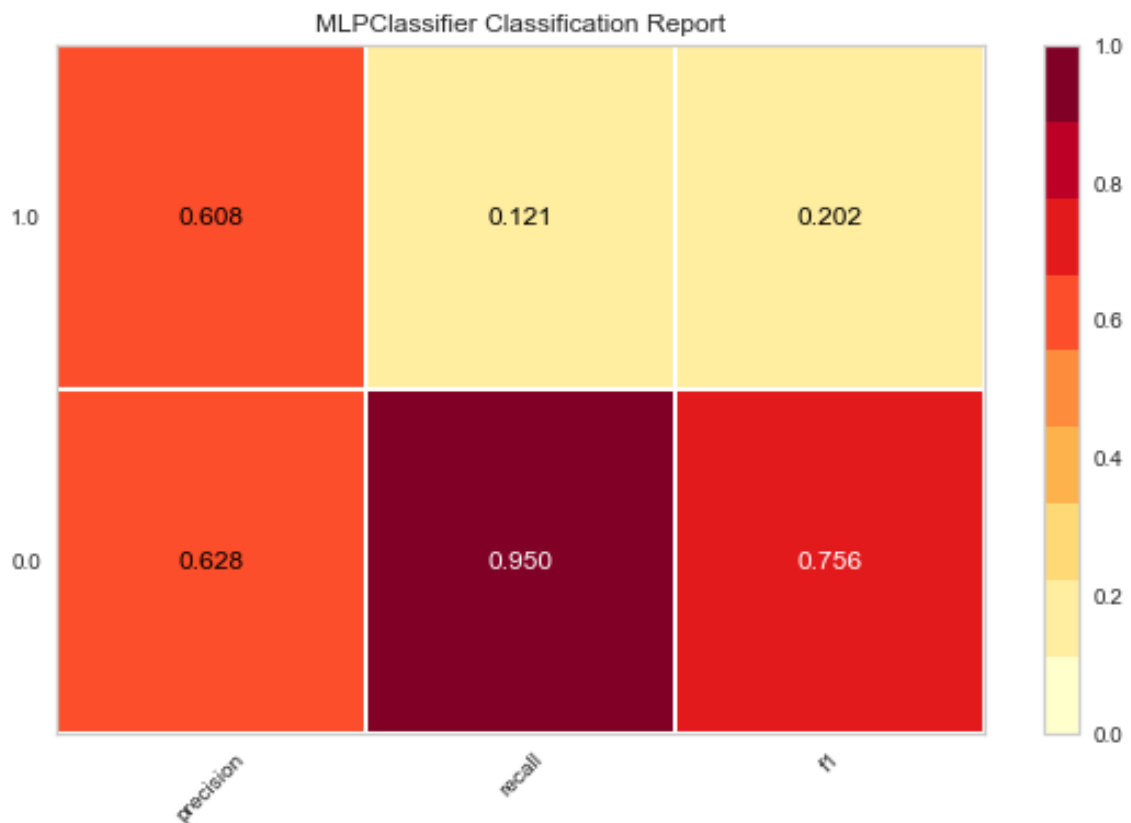
```

2     x_train, y_train, test_size=0.2, random_state=1)
3
4 mlpc = MLPClassifier(activation='tanh', solver='adam', alpha=0.0001,
    learning_rate_init= 0.01, max_iter=50, hidden_layer_sizes=(12, 6)
    , random_state=0)
5 mlpc = mlpc.fit(x_val,y_val)
6 result = mlpc.predict(x_test)

```

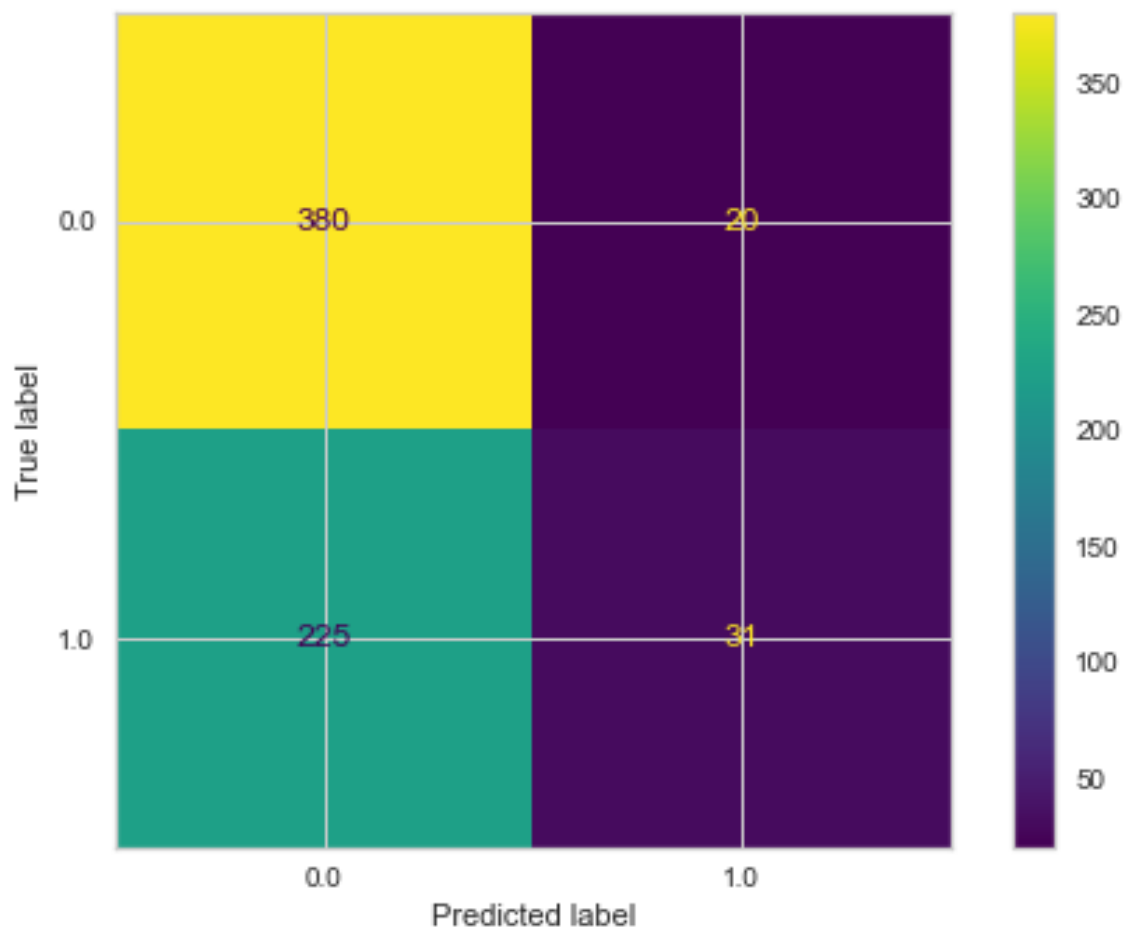
Na tréovanie sme použili triedu MLPClassifier s aktivačnou funkciou - the hyperbolic tan function, solver je použitý defaultne adam. Skryté vrstvy s počtom 12 a 6. Počet iterácií sme nastavili na 50. Validačné dáta sme rozdelili z tréovacej množniny v pomere 80:20.

Obr. 6: MLPClassifier klasifikačný report.

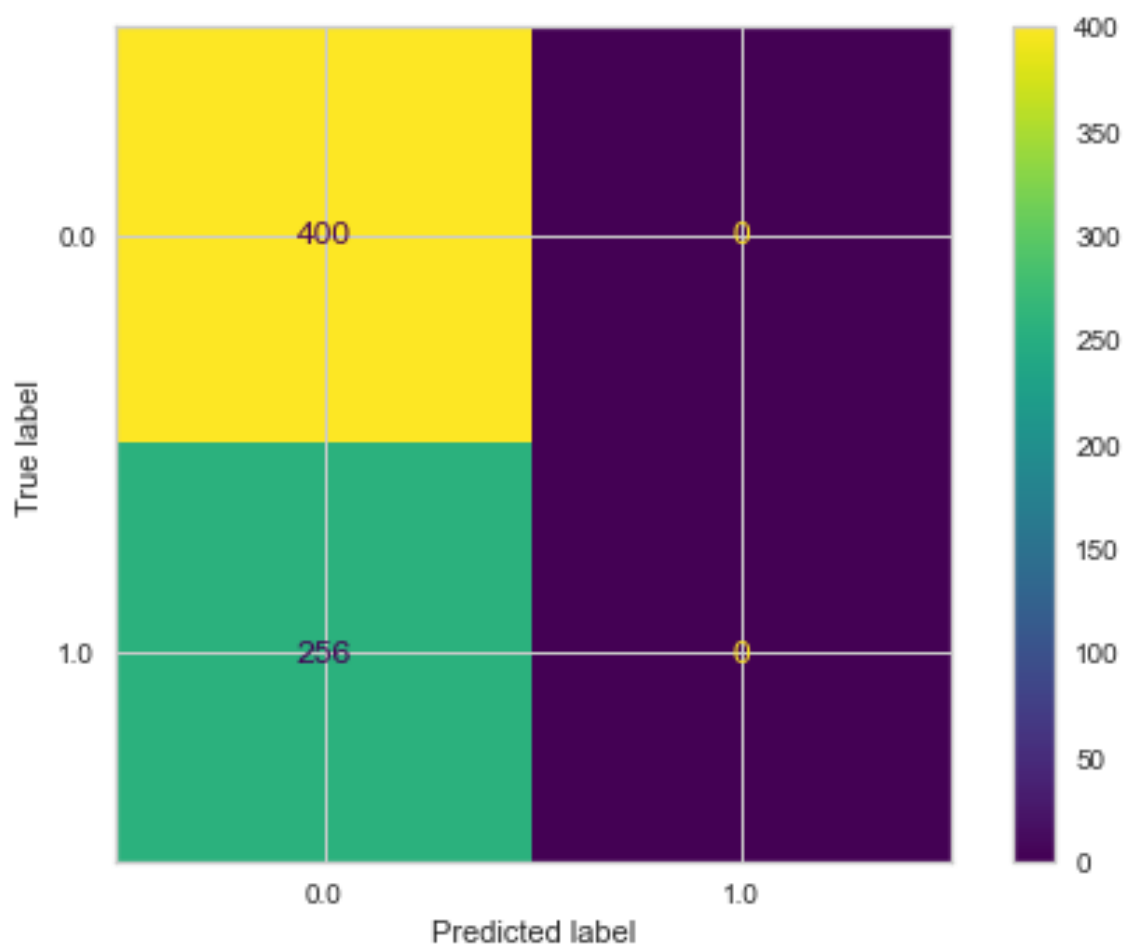


Presnosť pri tejto klasifikácii nám vyšla 62 percent.

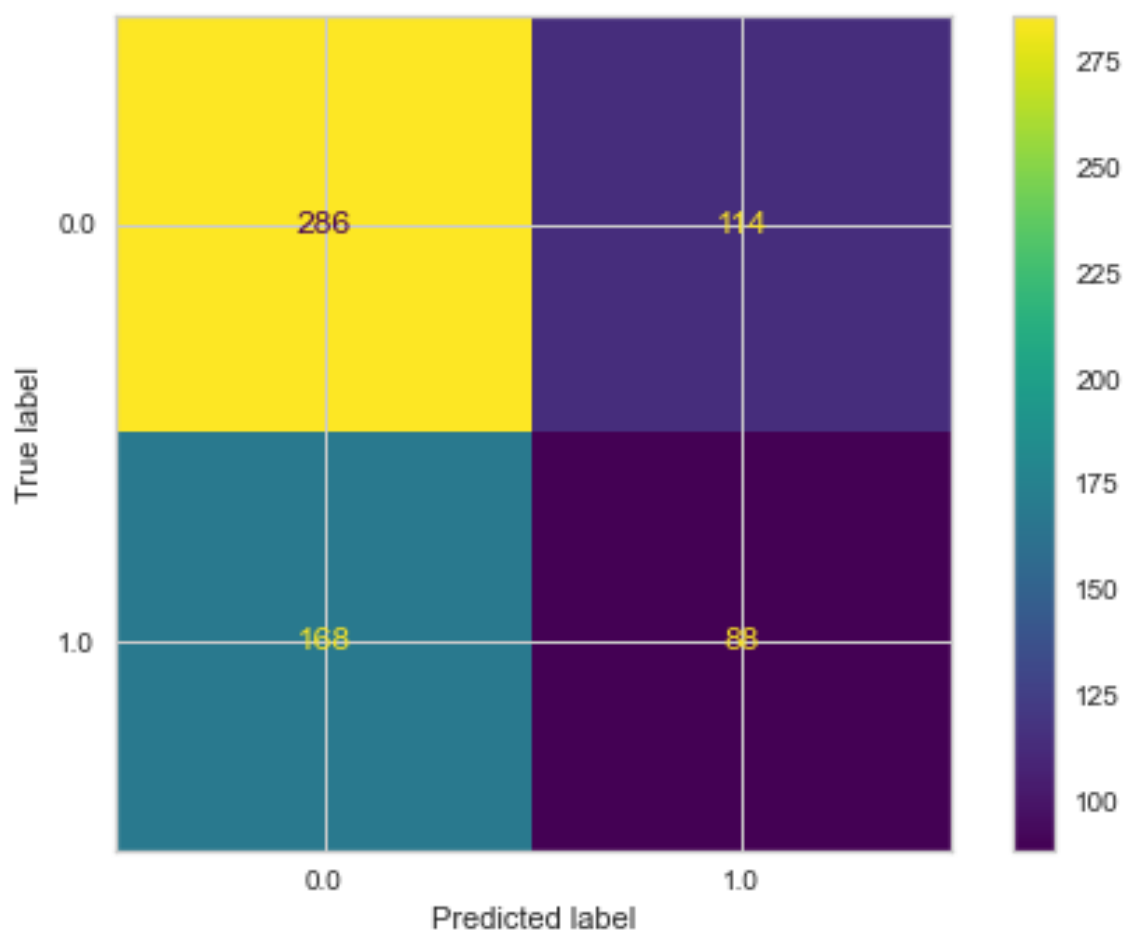
Obr. 7: Konfuzna matica pre treningovy set pri solveri adam.



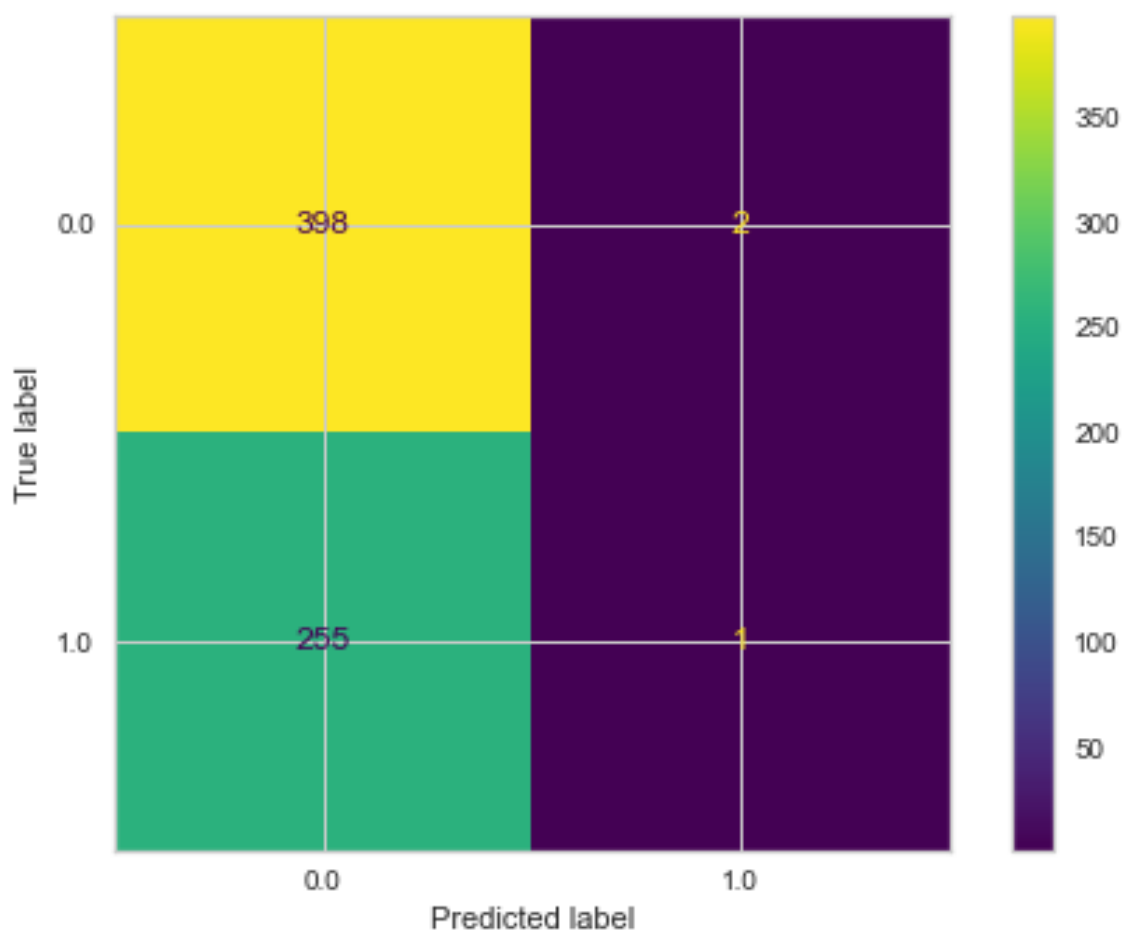
Obr. 8: Konfuzna matica pre treningovy set pri solveri SGD.



Obr. 9: Konfuzna matica pre validačný set pri solveri adam.



Obr. 10: Konfuzna matica pre validačný set pri solveri SGD.



Použili sme solver SGD. Max iter. 50 a skrytými vrstvami 18, 12 a 8. Pri validačnej množine nám vyšla presnosť 54 percent a pri tréningovej množine nám vyšla presnosť siete 62.59 percent.

Zdá sa, že lepšie výsledky sme dosiahli pre tréningovú množinu aj pomocou solvera SGD. Môžeme to vidieť na obr. číslo 10 pri konfúznej matici a pri celkovej úspešnosti 62.6 percent.

```

1 a = cross_val_score(mlpc, x_train, y_train, cv=5,
2                       scoring='accuracy').mean()*100
3 print("%f " % a)

```

Skúšali sme meniť rôzne nastavenia ako sú solvery, počet skrytých vrstiev, počet iterácií či rýchlosť učenia. Lepšie výsledky sme získali keď sme nastavili solver SGD. Pri nastavení rýchlosti učenia, aktivačnej funkcie sme získali najvyššiu presnosť a to 71.6 percent.

4 Grid search

Teraz sa pozrime, ako použiť GridSearchCV na zlepšenie presnosti nášho modelu.

```
1 param_grid = {'C': [0.1, 1, 10, 100],
2               'gamma': [1, 0.1, 0.01, 0.001, 0.0001],
3               'gamma': ['scale', 'auto'],
4               'kernel': ['linear']}
5 grid = GridSearchCV(SVC(), param_grid, refit = True, verbose = 3,
6                     n_jobs=-1)
7 # fitting the model for grid search
8 grid.fit(x_train, y_train)
9 # print best parameter after tuning
10 print(grid.best_params_)
11 grid_predictions = grid.predict(x_test)
12
13 # print classification report
14 print(classification_report(y_test, grid_predictions))
```

```
1 {'C': 0.1, 'gamma': 'scale', 'kernel': 'linear'}
2           precision    recall  f1-score   support
3
4      0.0         0.61      1.00      0.76       400
5      1.0         0.00      0.00      0.00       256
6
7      accuracy              0.61       656
8      macro avg           0.30      0.50      0.38       656
9      weighted avg           0.37      0.61      0.46       656
```

Môže sa zdať, že 'C': 100, 'gamma': 'scale', 'kernel': 'linear' sú najlepšie hodnoty pre hyperparametre pre model SVM. No nemusí to tak byť. Ale pre akýkoľvek iný súbor údajov môže mať model SVM rôzne optimálne hodnoty pre hyperparametre, ktoré môžu zlepšiť jeho výkon.

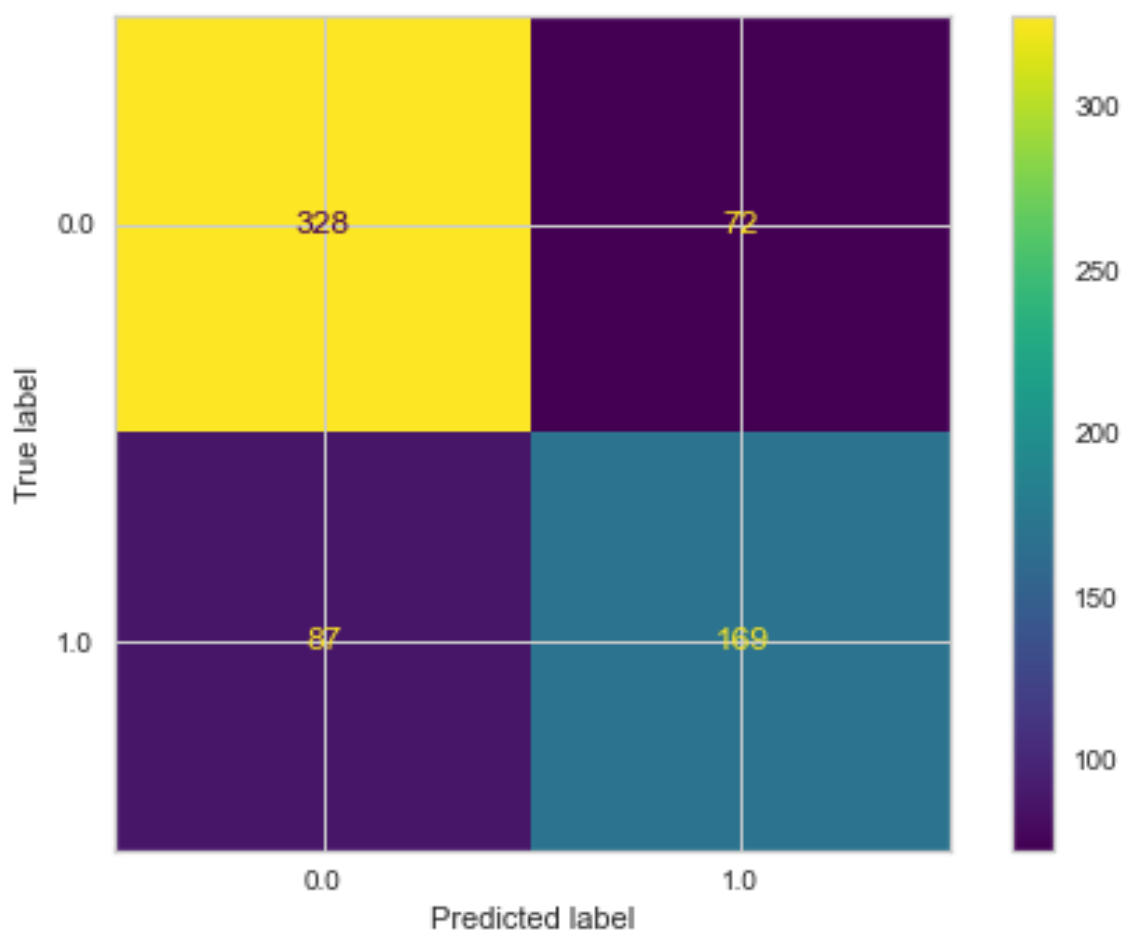
5 Tabulka experimentov

Obr. 11: Tabulka experimentov.

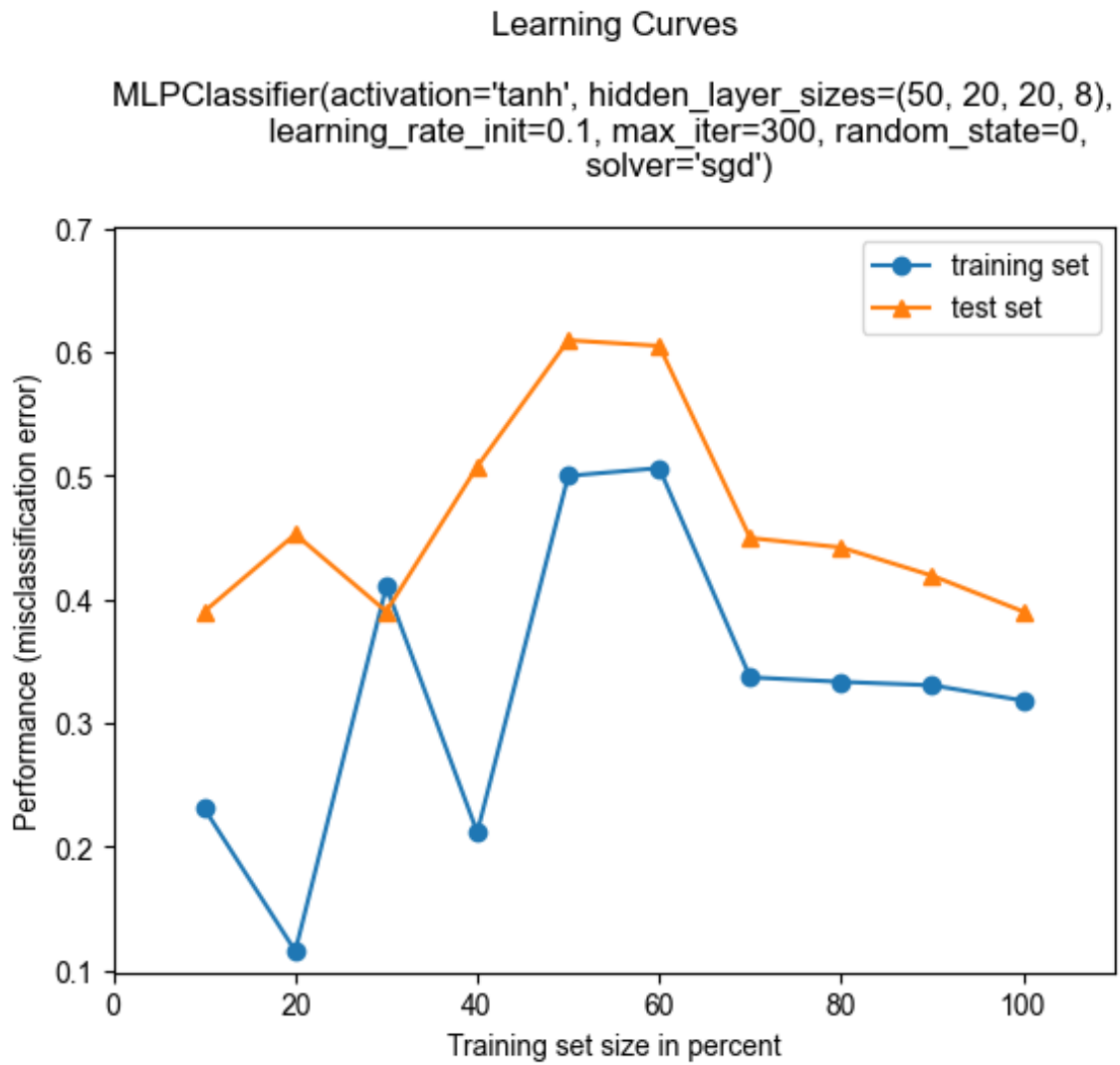
| set | ativation | solver | learning rate | max iter | hidden layers | accuracy % |
|-------|-----------|--------|---------------|----------|---------------|------------|
| train | tanh | sgd | 0.01 | 50 | 18,12,8 | 62.5 |
| train | relu | sgd | 0.1 | 100 | 50,30,8 | 66 |
| train | relu | sgd | 0.1 | 150 | 50,30,8 | 71.6 |
| train | relu | sgd | 0.1 | 200 | 50,30,8 | 73 |
| train | relu | sgd | 0.1 | 300 | 50,30,8 | 72.8 |
| train | relu | sgd | 0.1 | 300 | 50,20,20, 8 | 69.8 |
| train | tanh | sgd | 0.1 | 300 | 50,20,20, 8 | 75.7 |
| train | tanh | adam | 0.1 | 300 | 50,20,20, 8 | 60 |

Skúšali sme rôzne nastavenia parametrov a najlepší výsledok trénovania pomocou MLPClassifier sme získali 75.7 percent.

Obr. 12: Konfúzna matica pre náš najlepší výsledok.



Obr. 13: Priebeh tréovania pre náš najlepší výsledok.



Literatúra

https://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPClassifier.html