

## Mašinsko učenje – Domaći 4

### Random forest

S obzirom da *random forest* ima tendenciju da *overfitt*-uje podatke, potrebno je izvršiti regularizaciju hiperparametara modela. Najbolji način da se ovo uradi je pomoću *grid search*-a.

Pomoću *random grid search*-a mogu se zadati opsezi vrednosti hiperparametara od kojih će on na slučajan način birati kombinacije i računati K-fold kros validaciju za svaku kombinaciju. Na kraju će biti sačuvan najbolji model i njegovi parametri. Velika prednost *random grid search*-a je to što se na brz način mogu pronaći solidni parametri jer ne razmatra sve moguće kombinacije već samo one koje je slučajno odabrao.

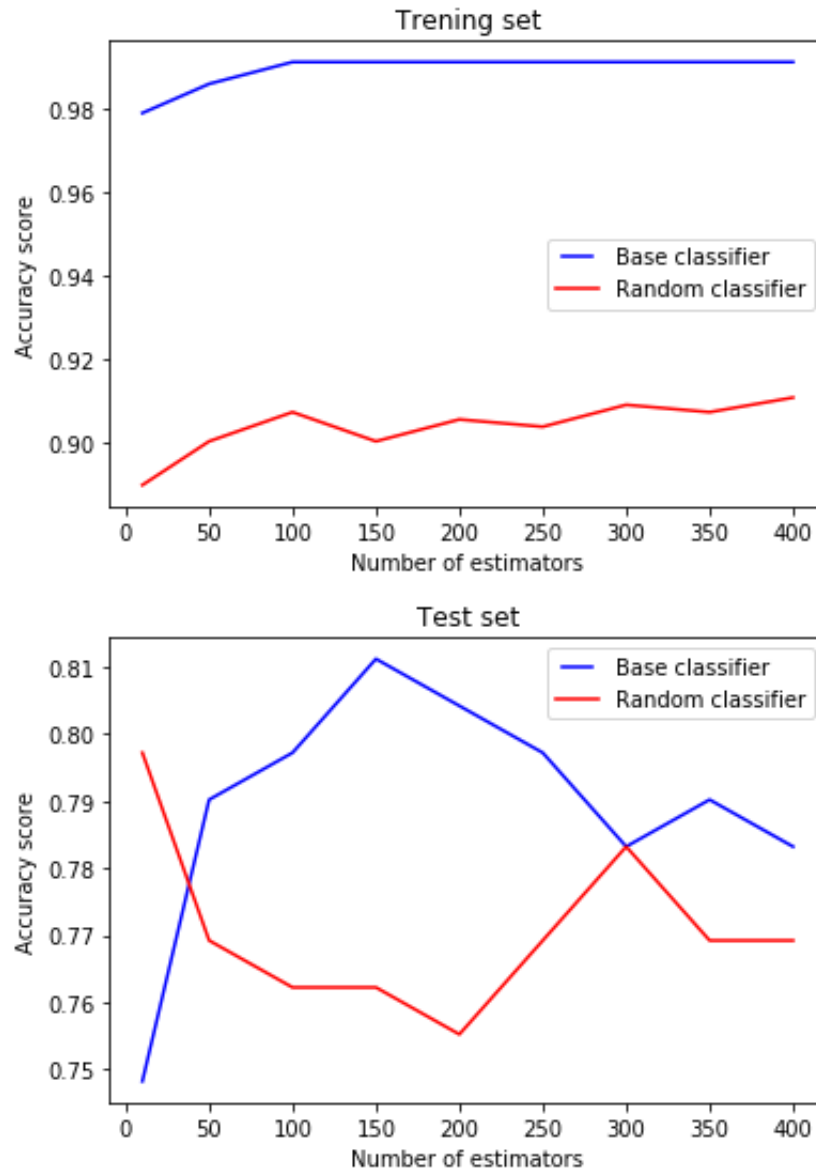
Nakon primene *random grid search*-a može se primeniti *grid search* sa vrednostima parametara blizu onih koje je pronašao *random grid search*. Na taj način se ne ispituju samo slučajne kombinacije hiperparametara već sve moguće kombinacije, i moguće je pronaći najbolji mogući model.

Parametri koji su isprobani za *random grid search*:

```
{'bootstrap': [True, False],  
 'max_depth': [10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110, None],  
 'max_features': [2, 4, 6, 'auto'],  
 'min_samples_leaf': [1, 2, 5],  
 'min_samples_split': [2, 5, 10],  
 'n_estimators': [10, 50, 100, 150, 200, 250, 300, 350, 400]}
```

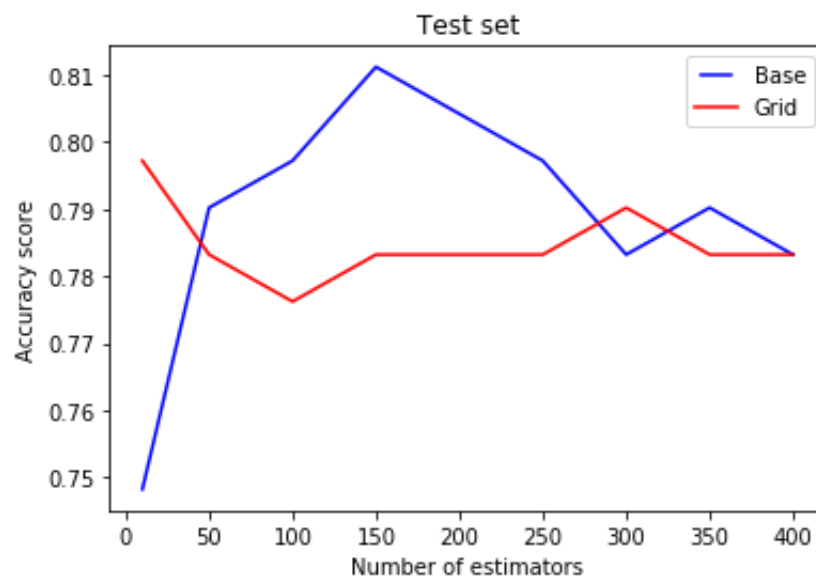
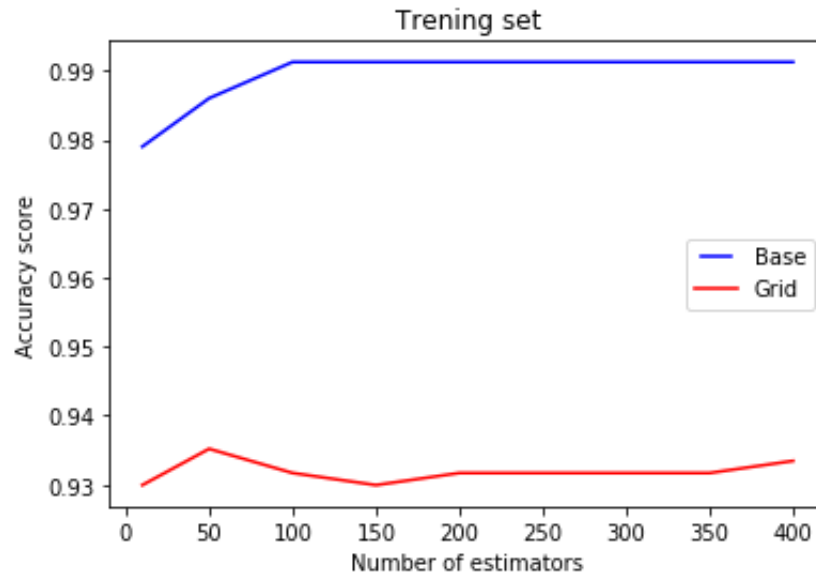
Svaki put kada se pozove ova metoda mogu se dobiti različiti parametri modela, s obzirom da se uvek na slučajan način biraju kombinacije hiperparametara.

Na sledećoj slici mogu se videti tačnosti za *base* klasifikator bez *tunovanih* hiperparametara i za klasifikator pronađen pomoću *random grid search*-a. Iako deluje da *base* klasifikator ima bolju tačnost, ako se poredi tačnost na trening setu, može se uočiti da *base* klasifikator *overfitt*-uje podatke.



Nakon analize performansi modela pronadjenog *random grid search*-om, nije lose pokusati i dalje tunovanje hiperparametara pomocu standardnog *grid search*-a.

Vrednosti parametara razmatrane u *grid search*-u su u okolini vrednosti parametara najboljeg moguceg modela koji je pronadjen *random grid search*-om.

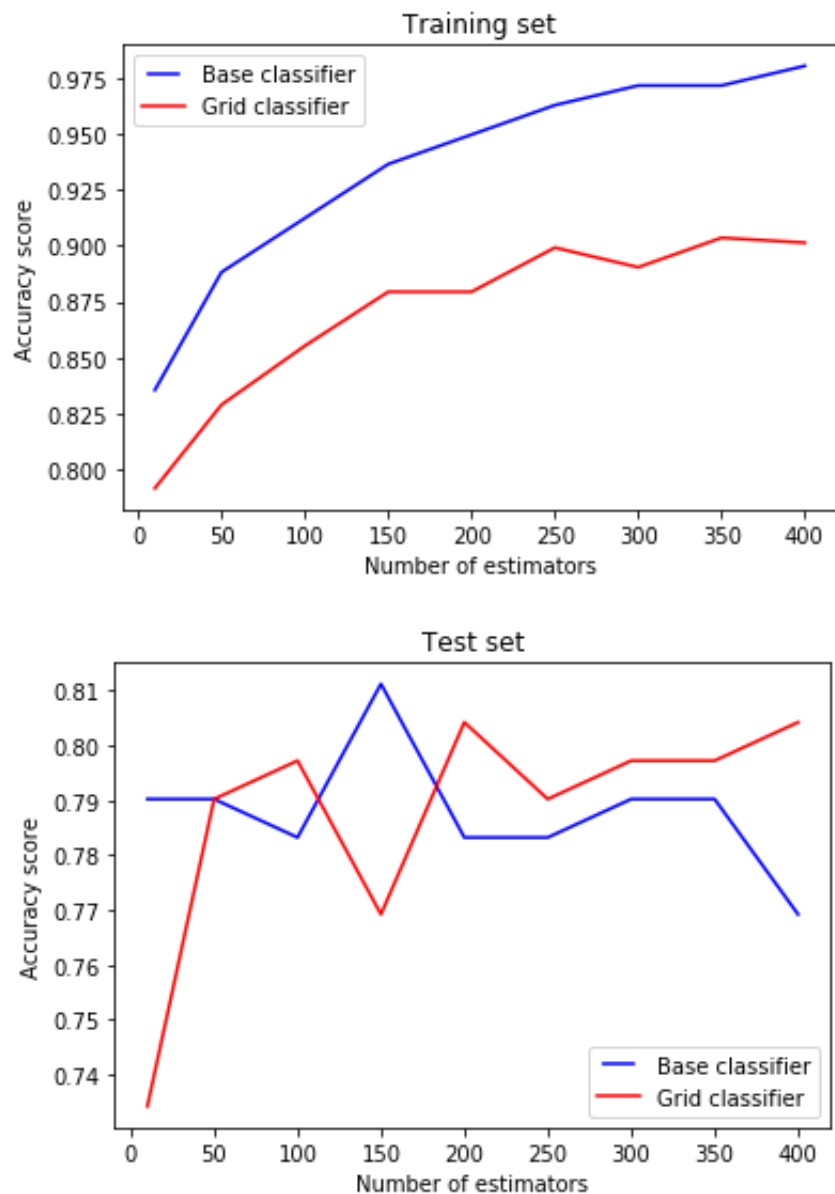


## Gradient Boosting Classifier

Kod ovog klasifikatora bitan parametar koji se pojavljuje, koga nema *Random Forest* klasifikator, je *learning rate*. Parametri koji su razmatrani u *grid search*-u su:

```
{'learning_rate': [0.01, 0.1, 0.5],
 'max_depth': [1, 2, 5, 10],
 'max_features': [2, 4, 6],
 'min_samples_leaf': [1, 2, 3],
 'min_samples_split': [2, 4, 6],
 'n_estimators': [5, 10, 50, 100, 200],
 'subsample': [0.6, 0.8, 1]}
```

Tacnost na treningu i testu za *base* klasifikatora bez zadavanja hiperparametara i najboljeg klasifikatora pronadjenog *grid search*-om prikazana je na sledecim slikama:



Na prvoj slici moze se videti da s povecavanjem broja estimatora raste i tacnost klasifikacije na trening setu. S obzirom da tacnost klasifikacije na test setu ostaje ista, moze se zakljuciti da dolazi do *overfitting*-

a. Pojava *overfitting*-a je znatno manja kod klasifikatora pronađenog pomoću *grid search*-a u odnosu na *base* klasifikator.