





```
[29]: from time import time
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import SGDClassifier
from sklearn.tree import DecisionTreeClassifier

def train_pred_eval(X_train=X_train, y_train=y_train, X_test=X_test, y_test=y_test, random_state=0):

    # Instantiate
    svm_model = SVC(random_state=random_state)
    sgd_model = SGDClassifier(random_state=random_state)
    knn_model = KNeighborsClassifier()
    decision_model = DecisionTreeClassifier(random_state=random_state)

    learners = [svm_model, sgd_model, knn_model, decision_model]
    cross_score = []
    fbeta = {}
    train_time = {}

    # Train
    for learner in learners:

        # Learn
        start = time()
        learner.fit(X_train, y_train)
        end = time()

        # Predict
        learner_name = learner.__class__.__name__
        y_pred = learner.predict(X_test)

        # Scores
        train_time[learner_name] = end-start
        fbeta[learner_name] = fbeta_score(y_test, y_pred, beta=0.5)
        cross_score[learner_name] = np.mean(cross_val_score(learner, X_test, y_test))

    print(f"Train time : {train_time}")
    print(f"Cross Val : {cross_score}")
    print(f"F0.5 Score : {fbeta}")

train_pred_eval()
```

train time : {'SVC': 106.65814685821533, 'SGDClassifier': 0.45618605613708496, 'KNeighborsClassifier': 0.013034820556464625, 'DecisionTreeClassifier': 0.8054039478302002}  
Cross Val : {'SVC': 0.838553178551686, 'SGDClassifier': 0.8436705362078497, 'KNeighborsClassifier': 0.8110558319513543, 'DecisionTreeClassifier': 0.80829187396395158}  
F0.5 Score : {'SVC': 0.70325993939153713, 'SGDClassifier': 0.711759504862953, 'KNeighborsClassifier': 0.6581325301204819, 'DecisionTreeClassifier': 0.6539823008849557}

Looking at the results above, Show that SVC perform well on testing data, but it take a lot of time to train, considering both metrics and time I suggest choosing the SGDClassifier model. It runs faster, predicts well enough and more importantly it is simple

## Feature Selection

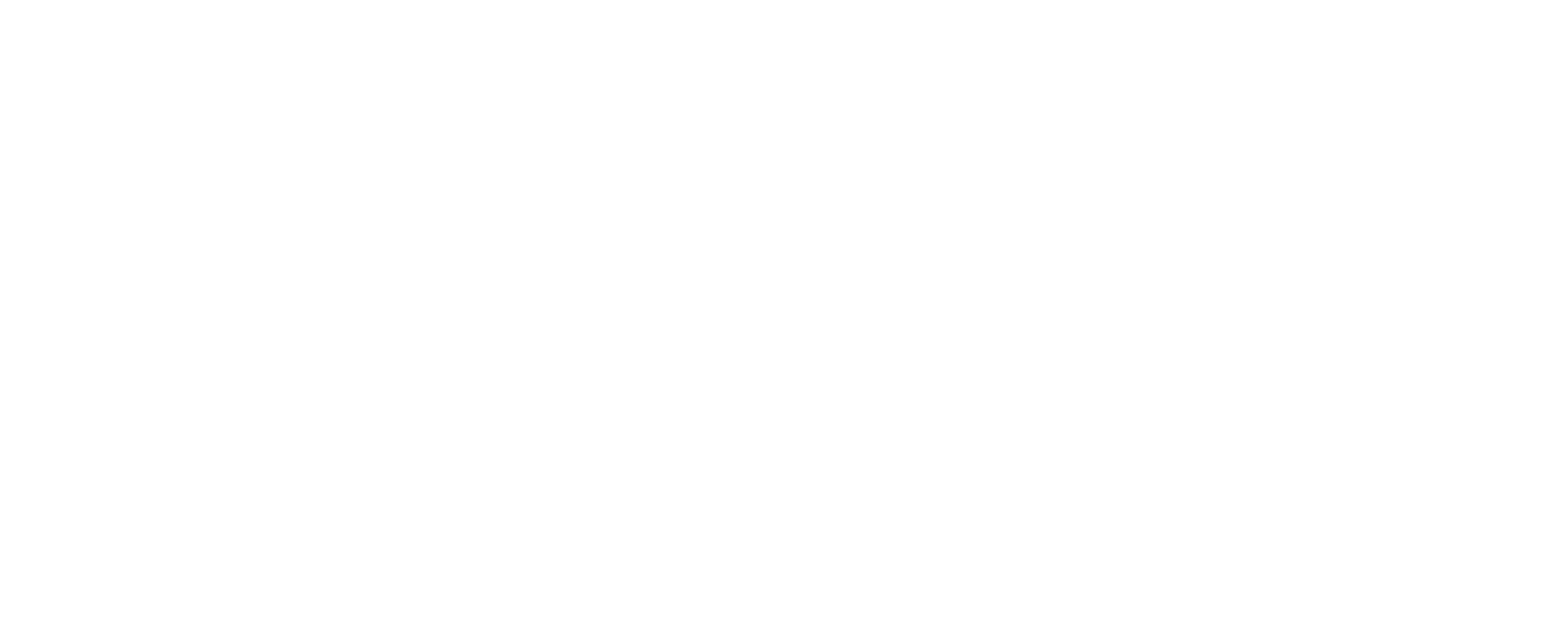
```
In [30]: # Using Random Forest
from sklearn.ensemble import RandomForestClassifier
# train the supervised model on the training set using .fit(X_train, y_train)
model = RandomForestClassifier()
model.fit(X_train, y_train)

# Extract the feature importances using .feature_importances_
importances = model.feature_importances_
feature_importances = pd.DataFrame(model.feature_importances_, index=model.feature_names_in_, columns=['Importance'])

top_ten_features = feature_importances[0:10]
top_ten_features
```

Out[30]:

	importance
age	0.240368
hours-per-week	0.115853
capital-gain	0.107676
marital-status, Married-civ-spouse	0.091896
education-num	0.064642
capital-loss	0.038426
marital-status, Never-married	0.036294
occupation, Exec-managerial	0.020426
relationship, Not-in-family	0.018981
sex, Male	0.018023



## In Progress

```
In [31]:
```