



Evaluation Metrics & Model Selection



Table of Content

What will We Learn Today?

1. Train test split
2. Evaluation Metrics
3. Hyperparameter Tuning





There are two questions:

1. how well is my model doing?, is this model good or not?
2. how do we improve it based on these metrics?

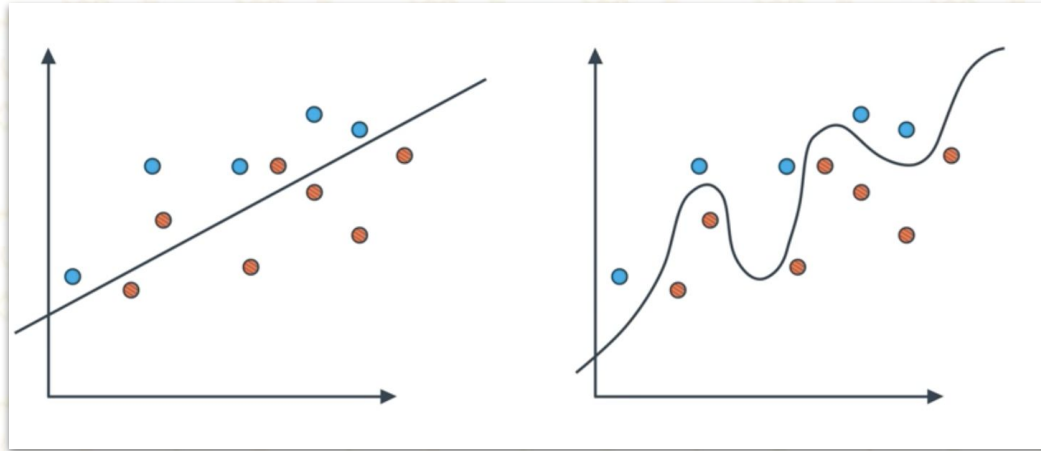
Train test split



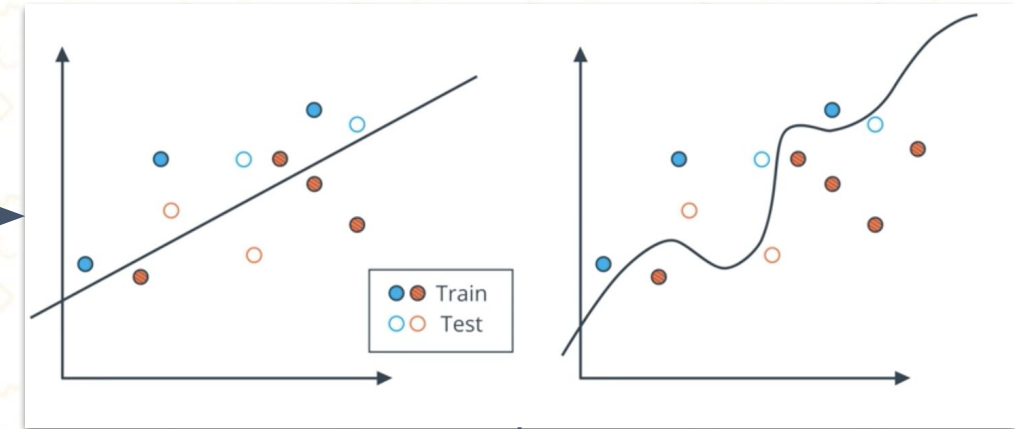


Train test split

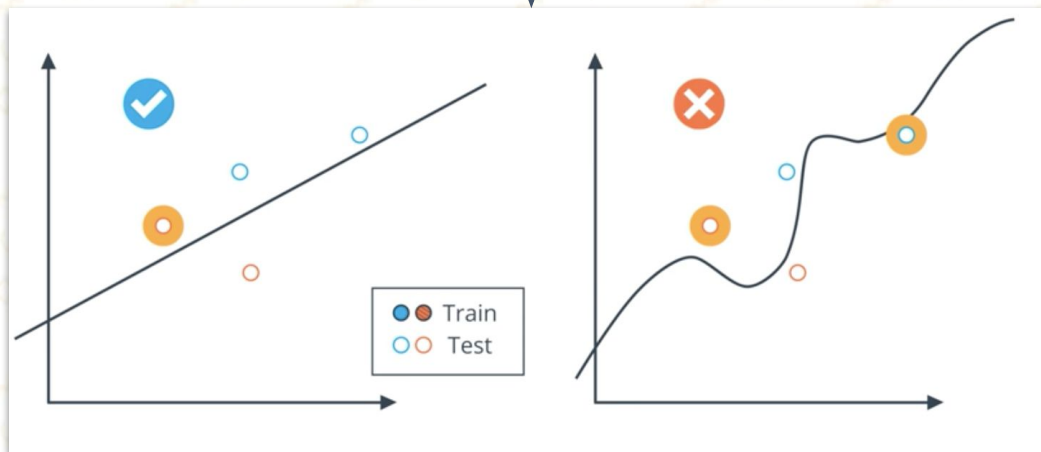
Which model is better?



Split them into training and testing



source: udacity machine learning nanodegree



```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test =  
train_test_split(X, y, test_size = 0.25)
```

Evaluation Metrics

How well is my model doing? How do we improve the model based on metrics?





Confusion Matrix

The confusion matrix is a simple table that stores these **four values**



PATIENTS

DIAGNOSIS

	Diagnosed Sick	Diagnosed Healthy
Sick	1000 True Positives	200 False Negatives
Healthy	800 False Positives	8000 True Negatives

10, 000 PATIENTS

source: udacity machine learning nanodegree

Two types of error:

1. **Type 1 error** (false positive), this is when we misdiagnose a healthy patient as sick
2. **Type 2 error** (false negative), this is when we misdiagnose a sick patient as healthy



Accuracy

Out of all target data, how many did we classify correctly?



PATIENTS

DIAGNOSIS

	Diagnosed Sick	Diagnosed Healthy
Sick	1000 True Positives	200 False Negatives
Healthy	800 False Positives	8000 True Negatives

$$Accuracy = \frac{(TruePositive + TrueNegative)}{(TP + FP + FN + TN)}$$

$$= \frac{1.000 + 8.000}{10.000}$$

$$= 90\%$$

10, 000 PATIENTS

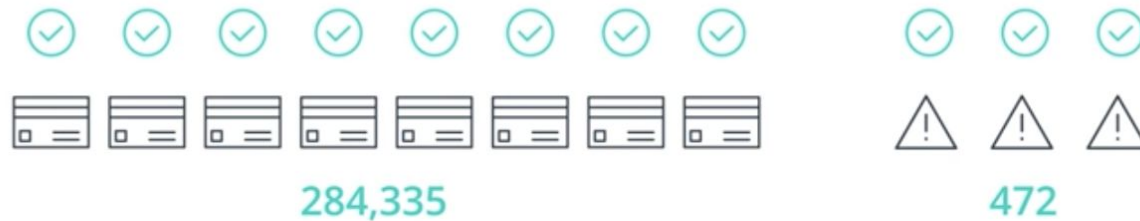
source: udacity machine learning nanodegree

```
from sklearn.metrics import accuracy_score  
accuracy_score(y_true, y_pred)
```




Accuracy may not be the best metric

◦ CREDIT CARD FRAUD



MODEL: ALL TRANSACTIONS ARE GOOD.

$$\text{ACCURACY} = \frac{284,335}{284,887} = 99.83\%$$

	Guessed Positive	Guessed Negative
Positive	True Positives	False Negatives
Negative	False Positives	True Negatives

source: udacity machine learning nanodegree

$$\text{Accuracy} = \frac{(\text{TruePositive} + \text{TrueNegative})}{(TP + FP + FN + TN)}$$

So, what's the next?



Precision, Recall & F1 Score

False positive: when the patient is healthy, and you diagnose them as sick

False negative: when the patient is sick, and you diagnose them as healthy

Which one you choose?







$$\text{Precision} = \frac{(\text{TruePositives})}{(TP+FP)}$$

Out of all the points predicted to be positive, how many of them were actually positive?

$$\text{Recall} = \frac{(\text{TruePositive})}{(TP+FN)}$$

Out of the points that are labeled positive, how many of them were correctly predicted as positive?

F1 score: the average between precision and recall

		Diagnosis	
		 DIAGNOSED SICK	 DIAGNOSED HEALTHY
Patients	 SICK	True Positive	 FALSE NEGATIVE
	 HEALTHY	 FALSE POSITIVE	True Negative

source: udacity machine learning nanodegree

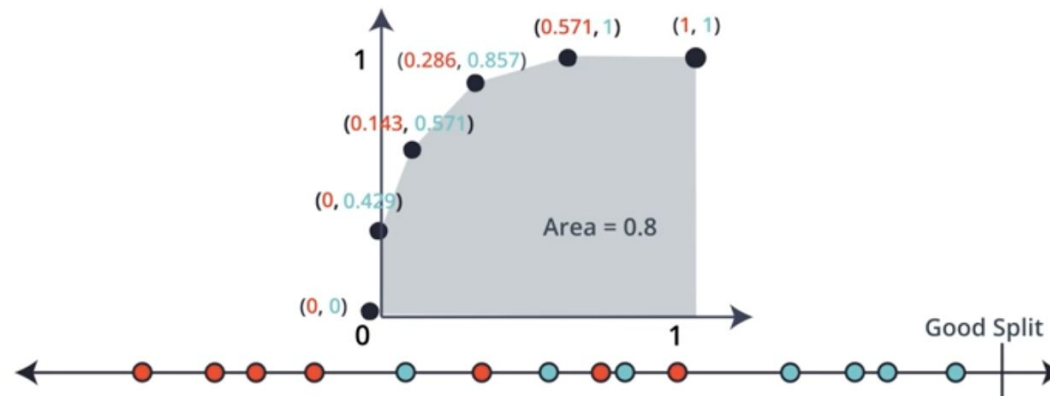





ROC Curve

Receiver Operator Characteristic curve

$$\text{True Positive Rate} = \frac{\text{TRUE POSITIVES}}{\text{ALL POSITIVES}} \quad \text{Sensitivity}$$

$$\text{False Positive Rate} = \frac{\text{FALSE POSITIVES}}{\text{ALL NEGATIVES}} \quad \text{Specificity}$$



Diagnosis		
	DIAGNOSED SICK	DIAGNOSED HEALTHY
Patients	 SICK True Positive	 FALSE NEGATIVE
	 HEALTHY FALSE POSITIVE	True Negative

source: udacity machine learning nanodegree



Regression Case

There are three metrics to see how well our model in the regression case:

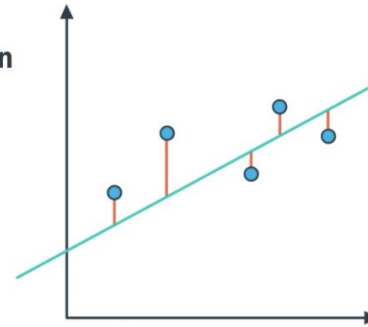
1. Mean absolute error, is just calculated with the mean absolute error function

```
from sklearn.metrics import mean_absolute_error
from sklearn.linear_model import LinearRegression

classifier = LinearRegression()
classifier.fit(X,y)

guesses = classifier.predict(X)

error = mean_absolute_error(y, guesses)
```



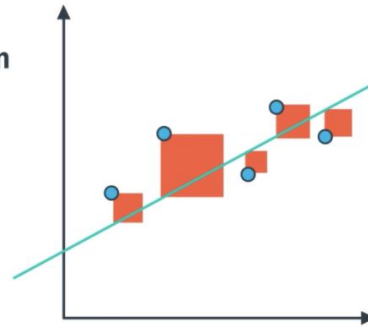
2. Mean squared error add the squares of the distance between the points and the line

```
from sklearn.metrics import mean_squared_error
from sklearn.linear_model import LinearRegression

classifier = LinearRegression()
classifier.fit(X,y)

guesses = classifier.predict(X)

error = mean_squared_error(y, guesses)
```

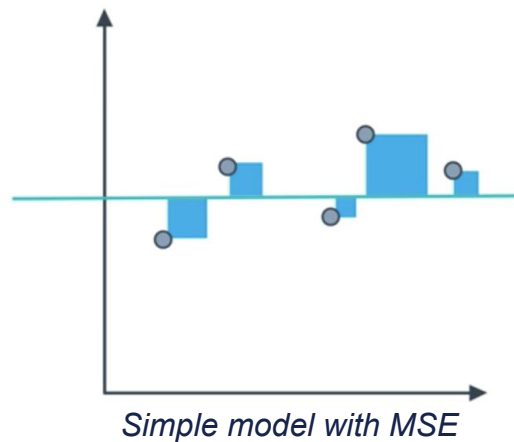




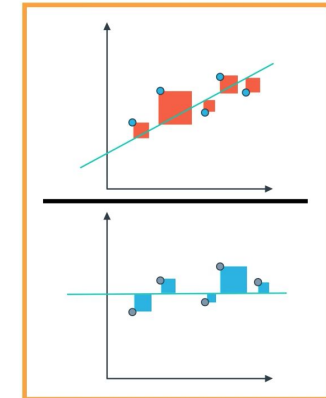
Regression Case

There are three metrics to see how well our model in the regression case:

3. R2 score is a metric based on comparing our model to the simplest possible model



$$\text{R}^2 = 1 -$$



There are two results:

1. **Bad Model**, the error should be similar, R2 score should be close to 0
2. **Good Model**, the mean squared error from our model should be smaller than the simple model, R2 score should be close to 1

```
from sklearn.metrics import r2_score  
r2_score(y_true, y_pred)
```

Hyperparameter Tuning



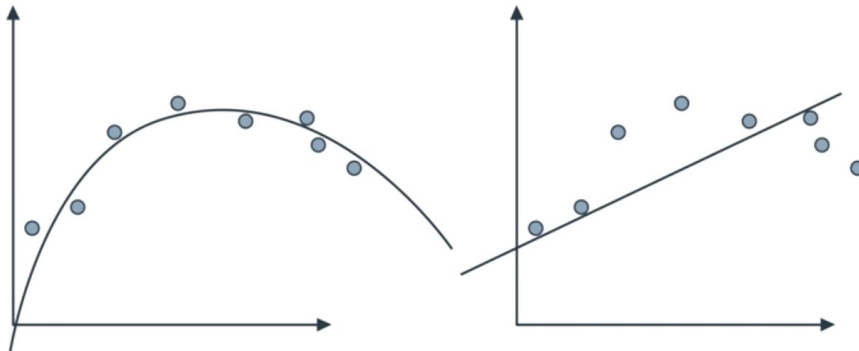


Types of Error

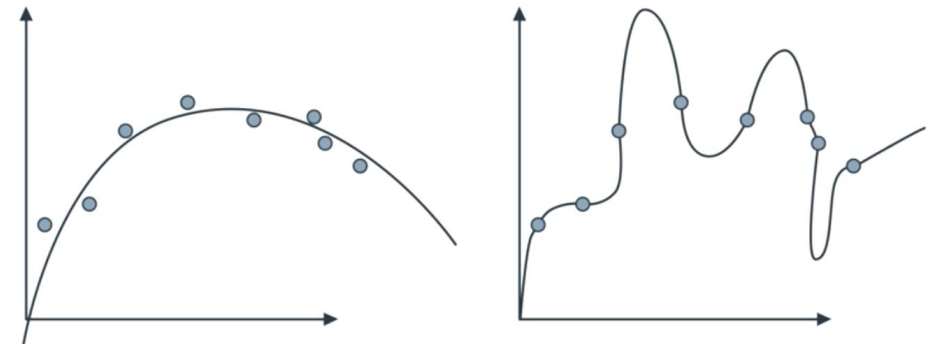
There are two types of error:

1. **Underfitting**, when we oversimplify the problem
2. **Overfitting**, when we overcomplicate the problem

○ UNDERFITTING
Error due to bias
This model will not do well in the training set



○ OVERFITTING
Error due to variance
This model does great in the training set

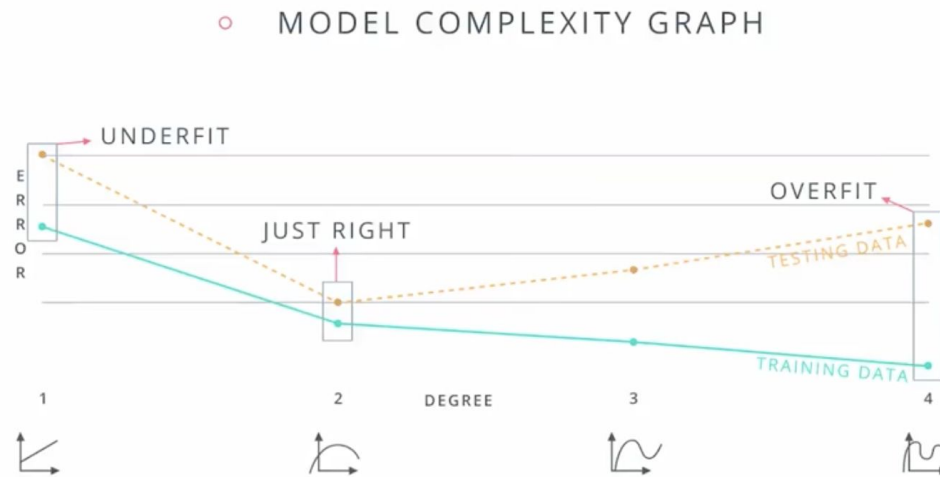


source: udacity machine learning nanodegree



Cross Validation

Instead of having training and testing set, we can use cross validation that will be used for **making decision**



source: udacity machine learning nanodegree



Never use your testing set for training



K-Fold - Cross Validation

Even if we have separated the data into three sets, it is better not to throw away the data that can be useful for training our algorithm.

with K-Fold Cross Validation

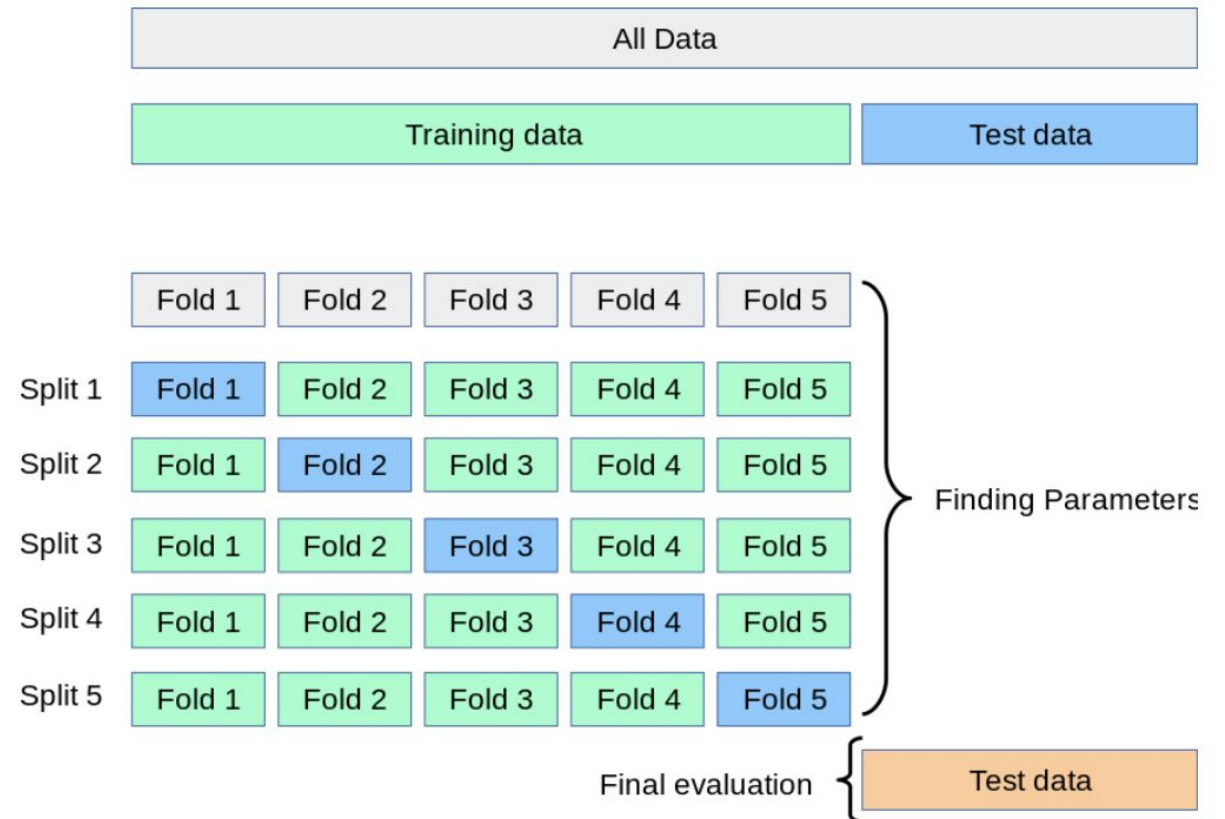
- Break our data into K buckets.
- Then, average the results to get a final model

```
from sklearn.model_selection import KFold
```

```
X = np.array([[1, 2], [3, 4], [1, 2], [3, 4]])  
y = np.array([1, 2, 3, 4])
```

```
kf = KFold(n_splits=10, shuffle = True)
```

source: scikit-learn.org





Grid Search CV

Finding the best parameter to use it as our model.

◦ GRID SEARCH CROSS VALIDATION

Kernel \ C	Linear	Polynomial
0.1	 F1 SCORE = 0.5	 F1 SCORE = 0.2
1	 F1 SCORE = 0.8	 F1 SCORE = 0.4
10	 F1 SCORE = 0.6	 F1 SCORE = 0.6

source: udacity machine learning nanodegree



```
from sklearn.model_selection import GridSearchCV  
from sklearn.metrics import make_scorer  
from sklearn.metrics import f1_score
```

```
scorer = make_scorer(f1_score)
```

```
parameters = {'max_depth':[2,4,6,8],  
              'min_samples_leaf':[2,4,6,8]}
```

```
grid_obj = GridSearchCV(clf, parameters,  
                        scoring=scorer)
```

```
grid_fit = grid_obj.fit(X, y)
```

**Thank
YOU**

