

# the rise of the machines

A primer to machine learning and predictive analytics using Azure ML



Beat Schwegler

head in the cloud feet on the ground

Twitter: @cloudbeatsch

Blog: <http://cloudbeatsch.com>

ALL MEN

SHOULD  
BE



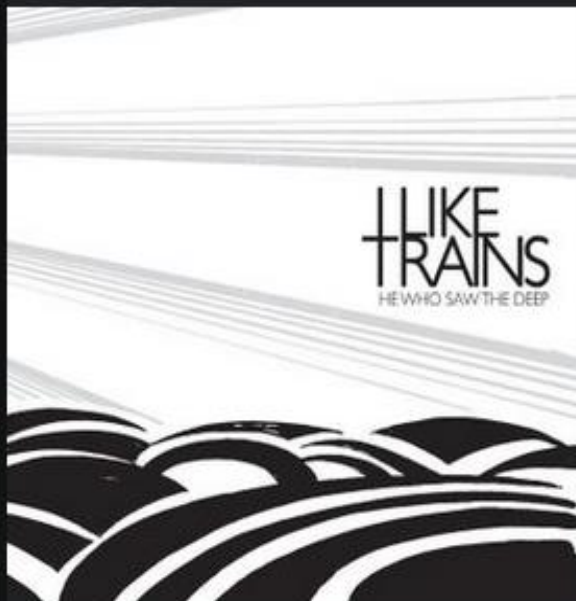
If you like **Francis International Airport**, we recommend **Naked Lunch**.



**Naked Lunch**

409 FOLLOWERS

You listened to **Engineers**. Check out **I Like Trains**.



**We Saw The Deep**

I Like Trains

You listened to **Francis International Airport**. Here's an album you might like.



**Everyone Is Having Fun**

Jack Beauregard

People who listen to **William Fitzsimmons** are also listening to **Joseph Arthur**.



**Joseph Arthur**

13,486 FOLLOWERS

Like **Headphone** and **Balthazar**? Check out **dEUS**.



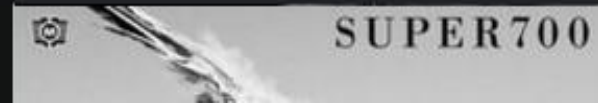
If you like **Most Unpleasant Men**, we recommend **The Secret Love Parade**.



If you like **Collapse Under The Empire**, we recommend **ef**.



You listened to **Francis International Airport**. Check out **Super700**.



| predictive analytics  
| predictions based on models

bing

New York, NY (NYC) to Los Angeles, CA (LAX)

Mon, 12/03 - Fri, 12/21 · 1 adult · Economy · [Change search](#)



Tip: **Buy** · Fares rising \$50+ · 80%+ Confidence

| predictive maintenance  
| fix it before it breaks

| predictive assistant  
| make predictions personal

fundamentals of machine learning

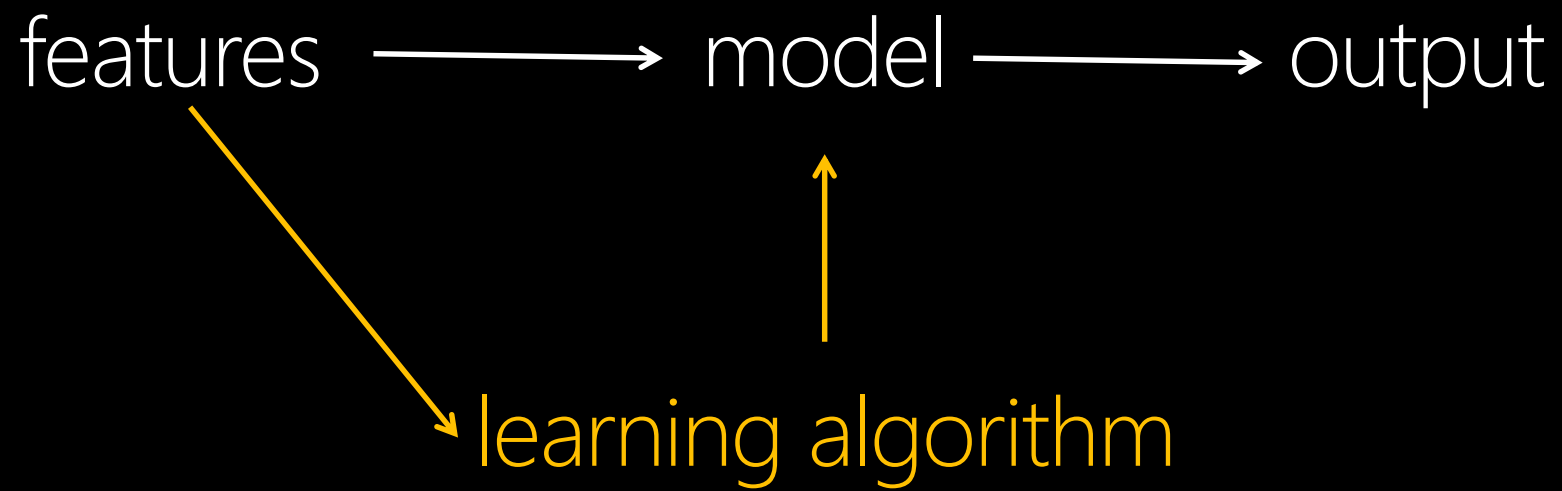
```
graph TD; A[fundamentals of machine learning] --> B[data science workflow]; B --> C[maml experiments];
```

data science workflow

maml experiments



machine learning  
algorithms and systems that improve  
their performance with experience



| features describe the domain  
| e.g. income, age, education, ...

| labels augment the learning data  
| e.g. this photo contains a man

unsupervised learning

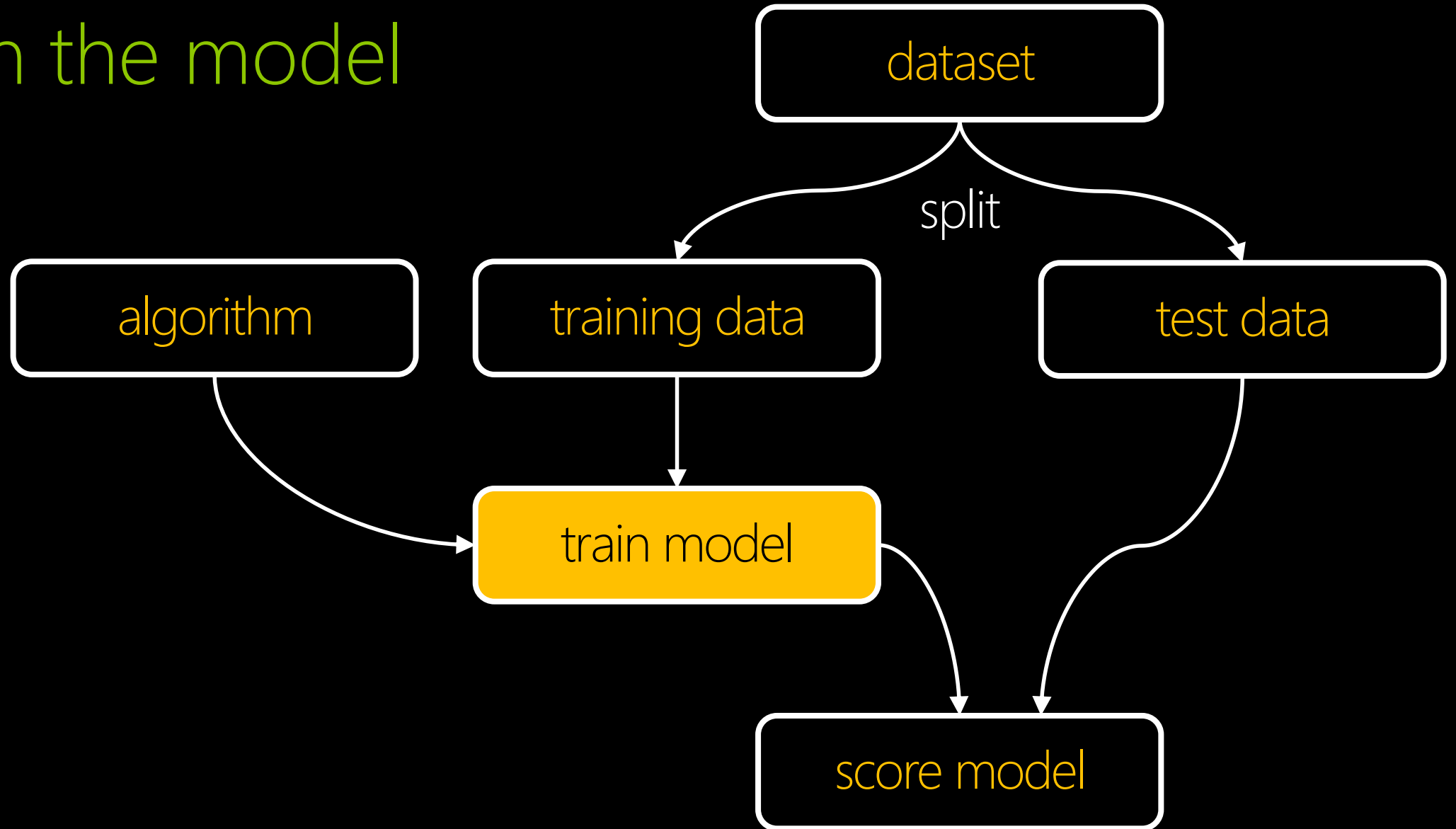
discovering clusters and associations using

unlabeled data

| supervised learning  
| using labeled data to train the model



# train the model



different tasks require different algorithms  
finding similar companies k-means clustering  
hourly bike rental prediction regression tree  
credit risk prediction decision tree

machine learning tools

open source: R, mahout, python, weka ...

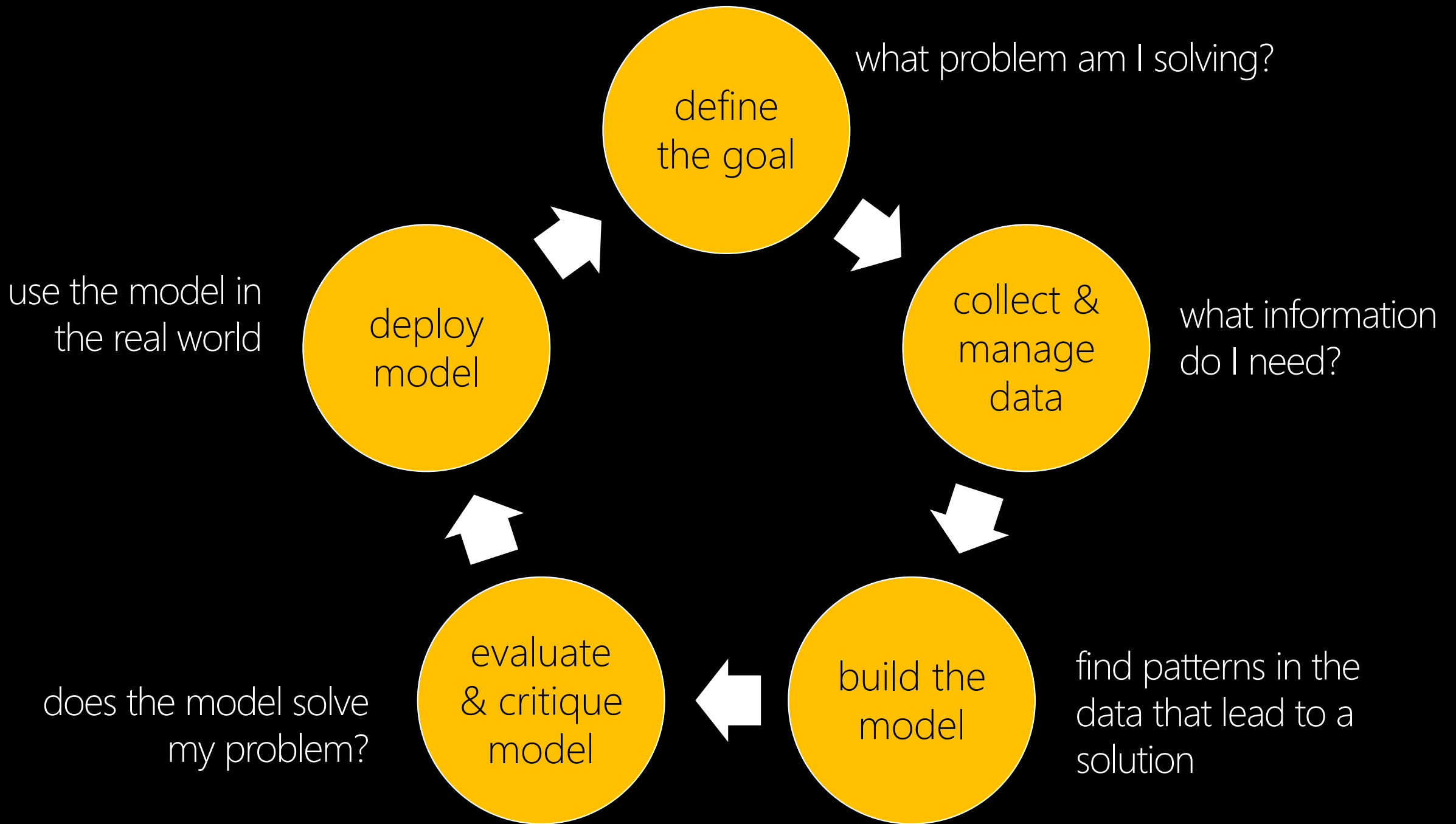
commercial: Azure ML, SAS, SPSS, MATLAB ...

fundamentals of machine learning

```
graph TD; A[fundamentals of machine learning] --> B[data science workflow]; B --> C[maml experiments]; style B fill:#ffff00,stroke:#fff,stroke-width:2px;
```

data science workflow

maml experiments



ask the right question!  
when will the clutch fail?

versus

what is the probability that the clutch fails  
within the next 3 months?



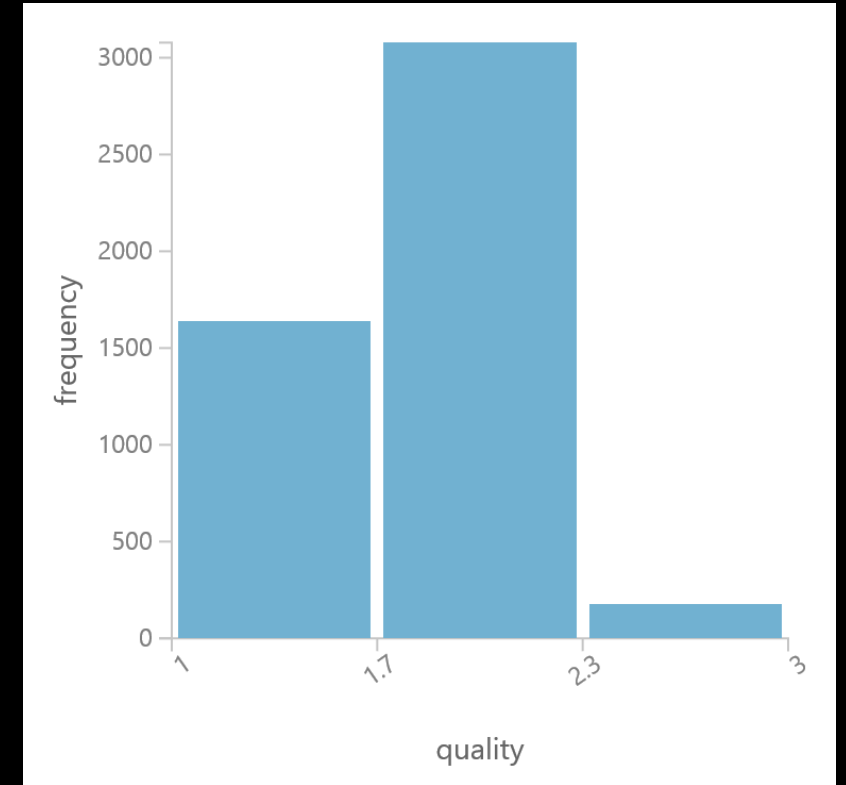
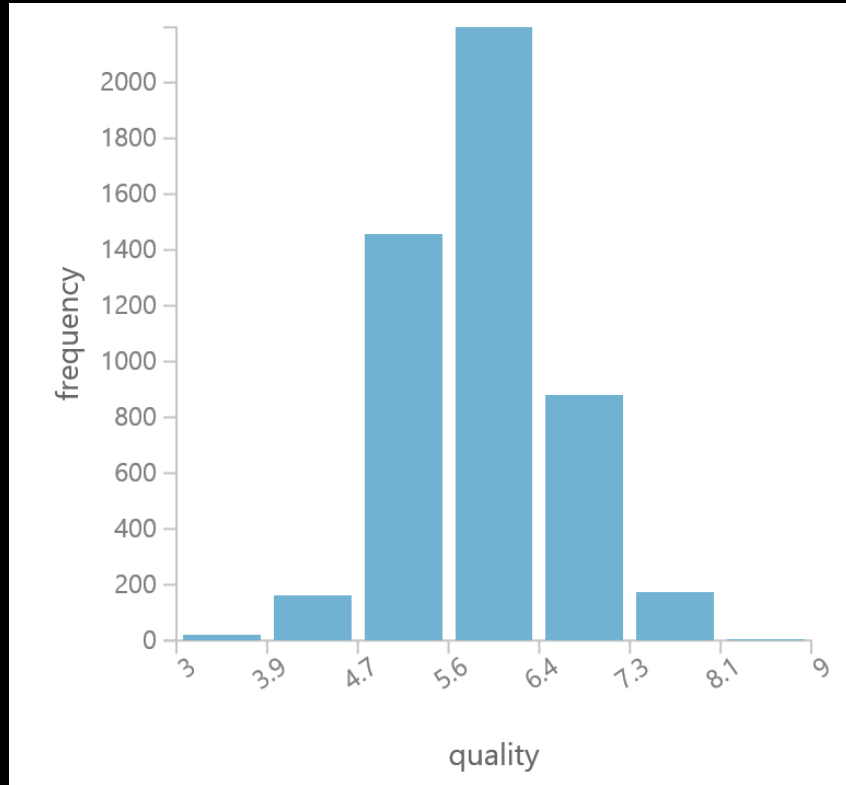
feature construction

transform e.g. cumulative time dependent

clipping outliers e.g. faulty measurements

normalize e.g. avoid feature dominance

quantize e.g. create categories (great/bad)



e.g. quantize wine quality ratings into  
bad/good/great buckets

train & test the model

better data often beats better algorithms

more training data  $\neq$  better model

does the algorithm correlate features?

does it support online learning?

deciding on algorithms

linear regression to predict

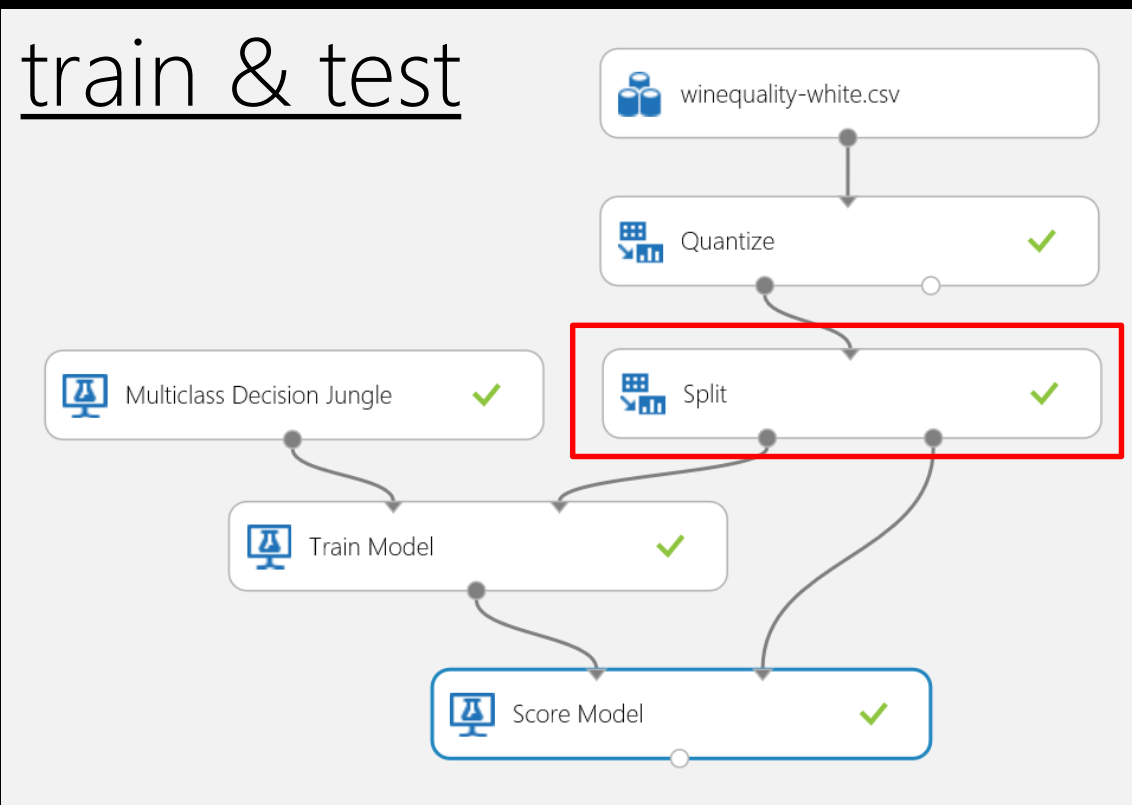
decision trees

to classify classification tree (2 & multi)

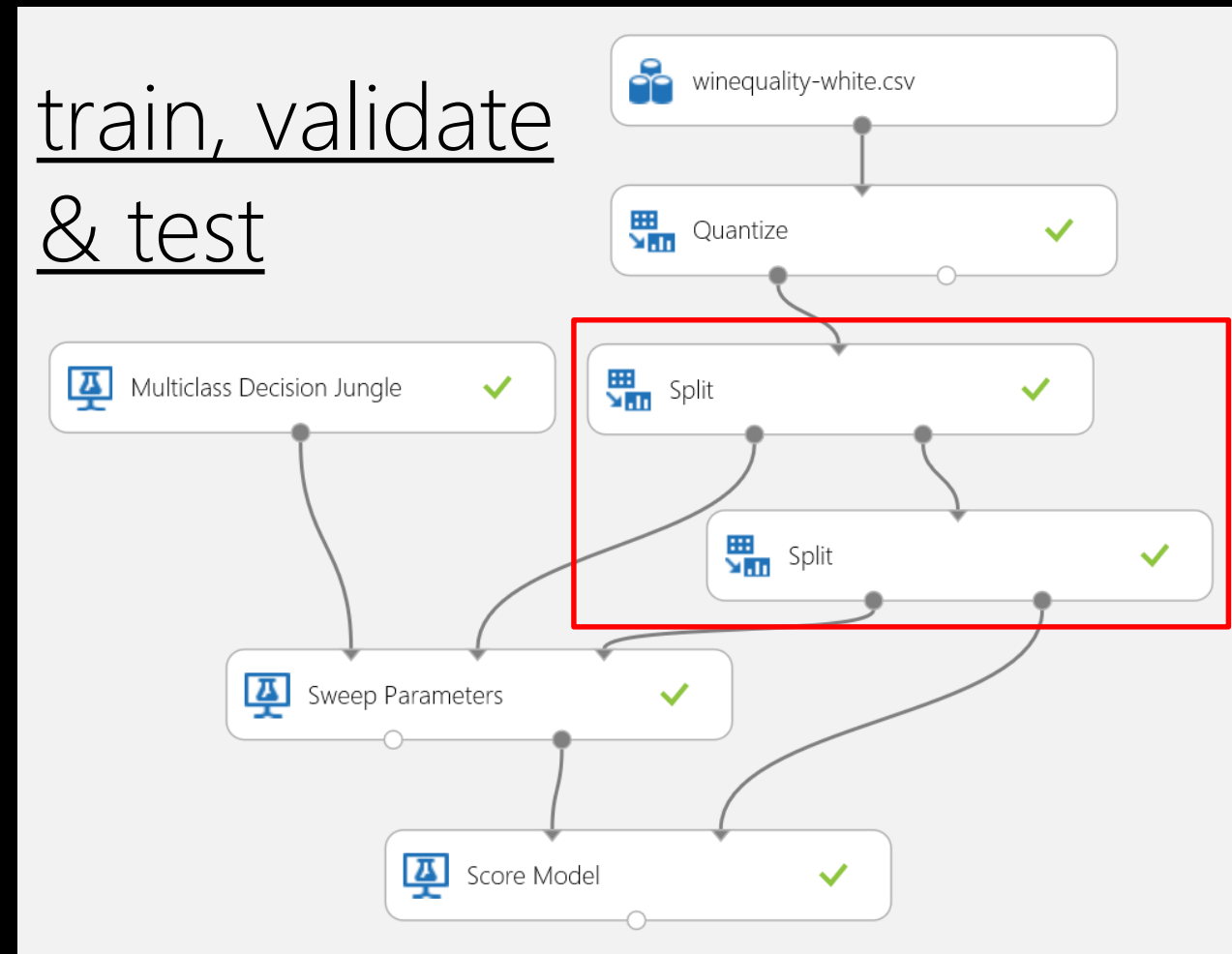
to predict regression tree

# 2-way partitioning versus 3-way partitioning

## train & test



## train, validate & test

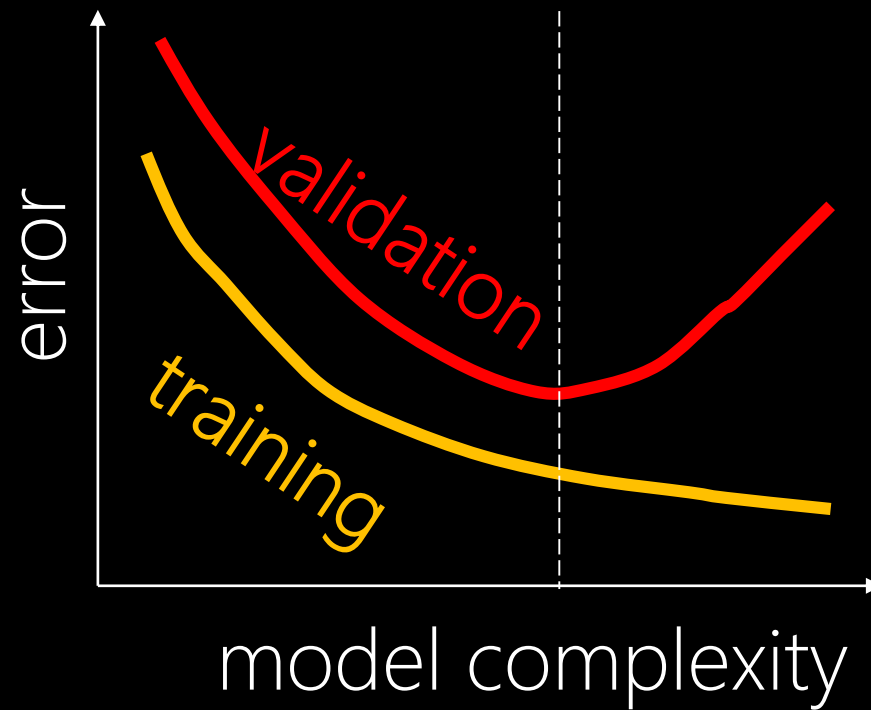


# cross validation - divide dataset into n folds

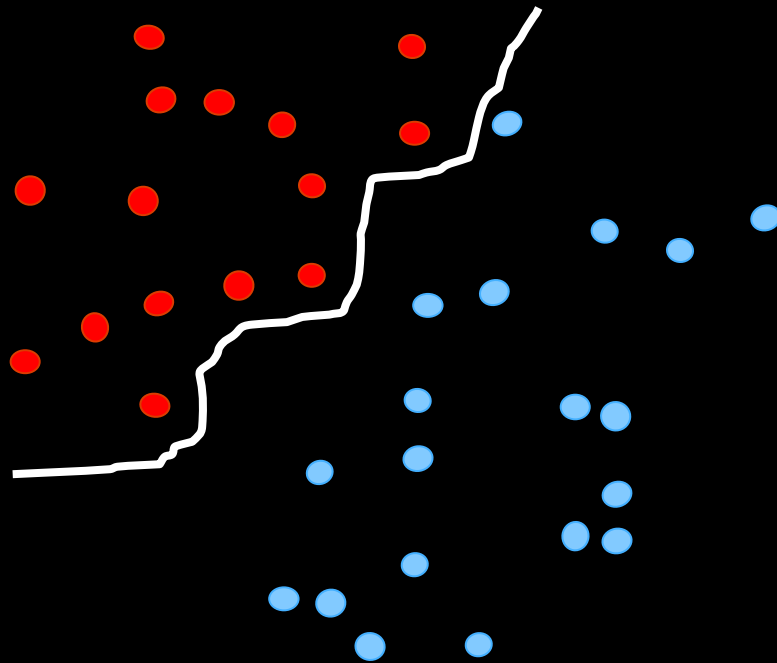
	A	B	C	D	E	F
Iteration 1	Test	Train	Train	Train	Train	Train
Iteration 2	Train	Test	Train	Train	Train	Train
Iteration 3	Train	Train	Test	Train	Train	Train
Iteration 4	Train	Train	Train	Test	Train	Train
Iteration 5	Train	Train	Train	Train	Test	Train
Iteration 6	Train	Train	Train	Train	Train	Test



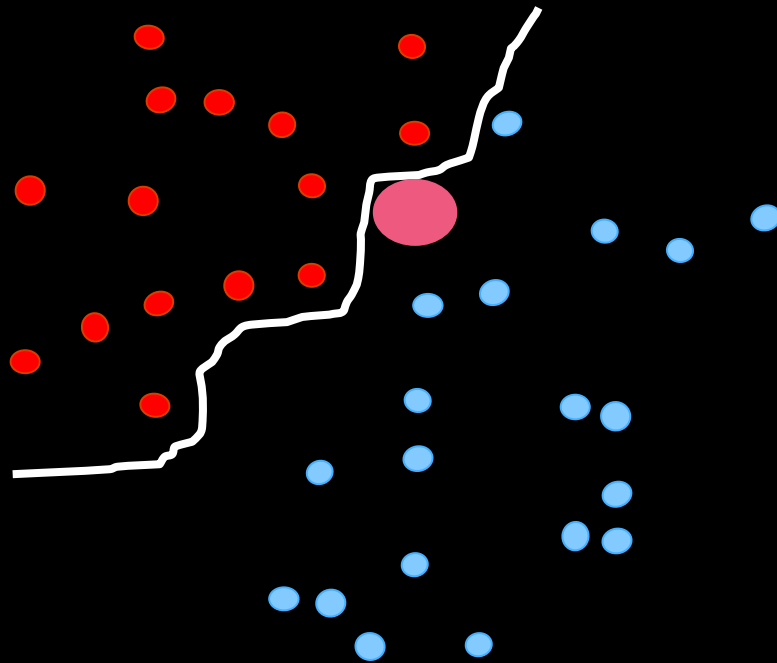
trade off between training error & validation performance



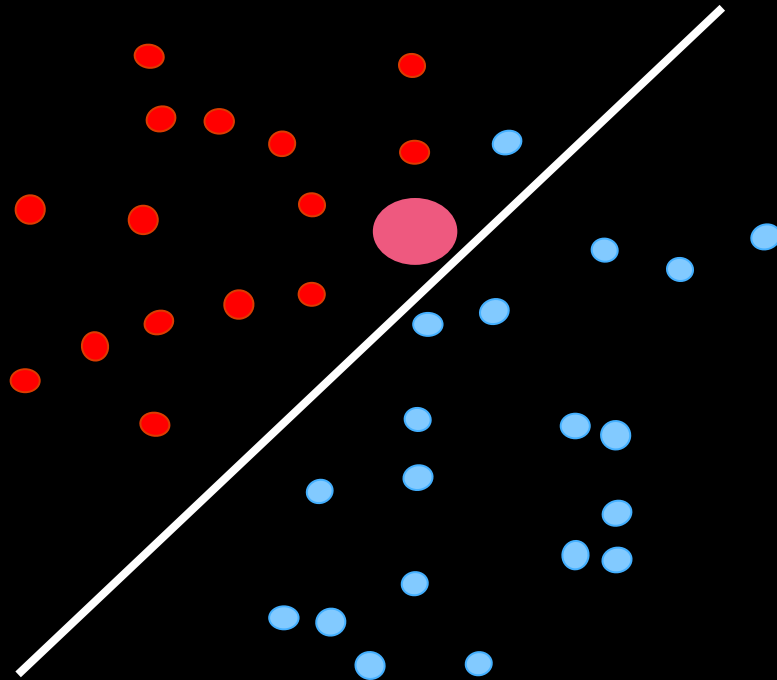
# over fitting and generalization



# over fitting and generalization



# over fitting and generalization



True Positive	False Negative	Accuracy	Precision	Threshold		AUC
<b>621</b>	<b>129</b>	<b>0.733</b>	<b>0.790</b>	<b>0.5</b>		<b>0.758</b>
False Positive	True Negative	Recall	F1 Score			
<b>165</b>	<b>185</b>	<b>0.828</b>	<b>0.809</b>			

accuracy:  $\frac{TP+TN}{P+N}$

→ 1 means all predictions are correct (over fitting)

precision:  $\frac{TP}{TP+FP}$

→ 1 means all P are P (but increased # of FN)

recall:  $\frac{TP}{TP+FN}$

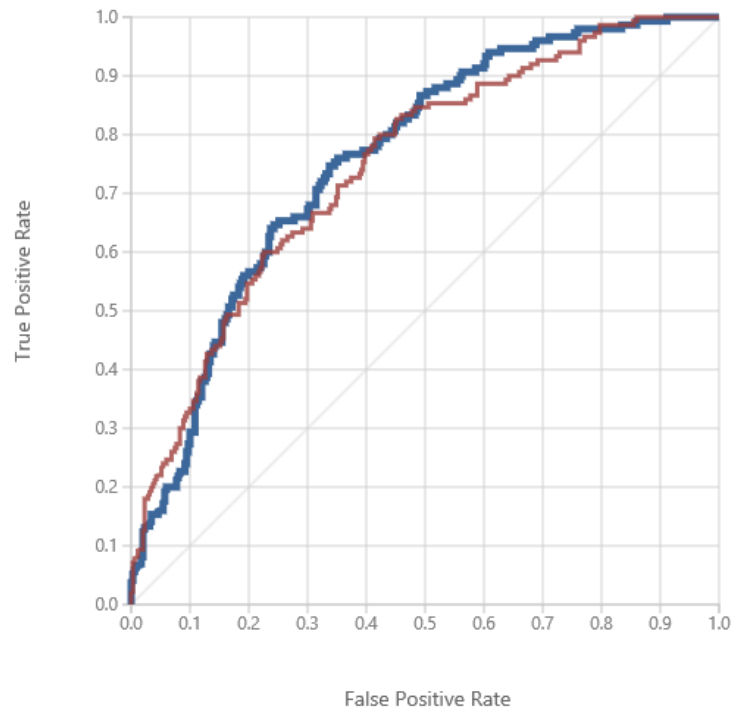
→ 1 means no FN (just another name for TP rate)

F1 score:  $\frac{precision+recall}{2}$

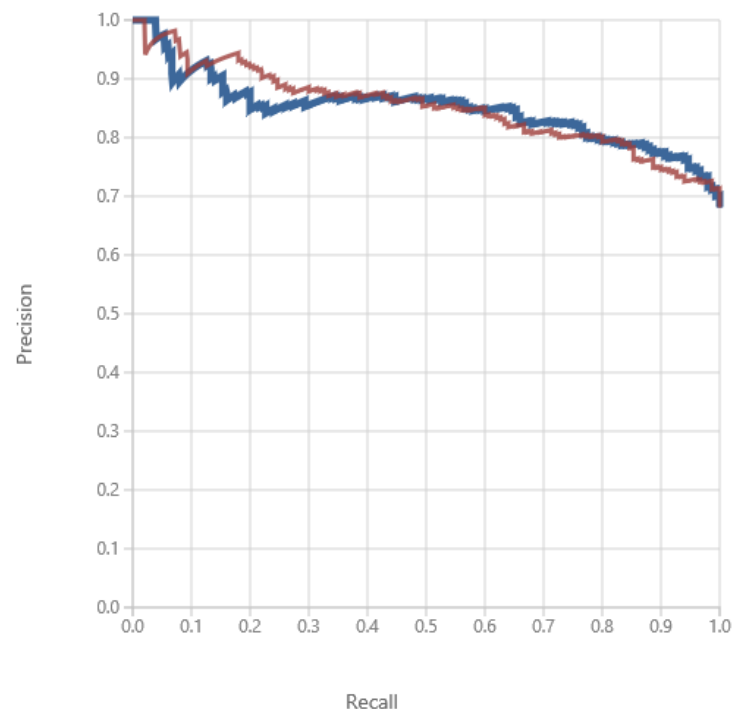
→ the closer to 1, the better

# evaluate the model

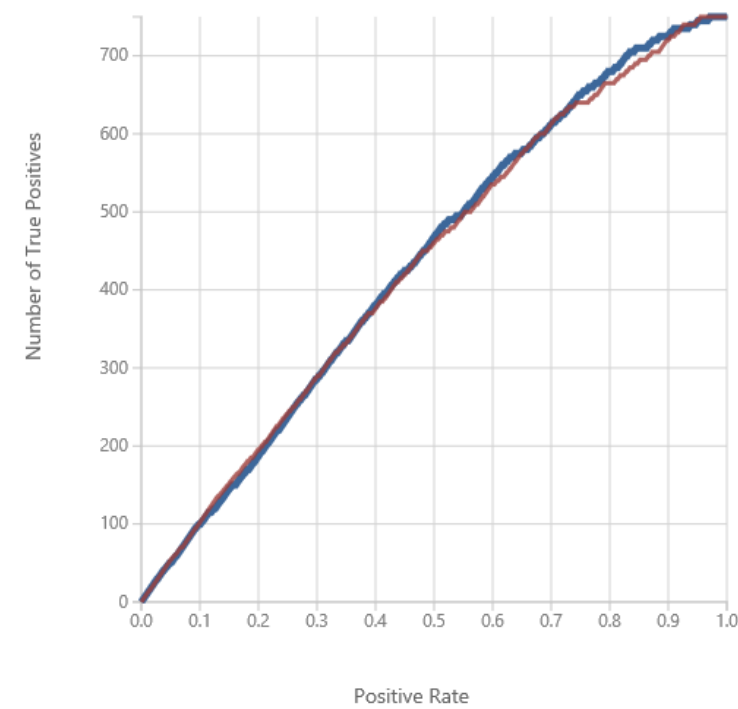
## ROC (Receiver Operating Characteristic)



## Precision/Recall



## Lift



True Positive	False Negative	Accuracy
<b>621</b>	<b>129</b>	<b>0.733</b>
False Positive	True Negative	Recall
<b>165</b>	<b>185</b>	<b>0.828</b>

Precision
<b>0.790</b>
F1 Score
<b>0.809</b>



AUC  
**0.758**



# put the model in production as a Web Service

```
{
  "Id": "wineTypeScore",
  "Instance":
  {
    "FeatureVector":
    {
      "pH": "3.5",
      "alcohol": "10.7"
    },
    "GlobalParameters": {}
  }
}
```

predict  
→

Scored Labels	Scored Probabilities
white	0.999882
red	0.242678
white	0.533091
white	0.971664
white	0.971664
white	0.533091
white	0.6598
white	0.999882
red	0.242678
white	0.888038

fundamentals of machine learning

```
graph TD; A[fundamentals of machine learning] --> B[data science workflow]; B --> C[maml experiments];
```

data science workflow

maml experiments

References to the wine dataset:

P. Cortez, A. Cerdeira, F. Almeida, T. Matos and J. Reis.  
Modeling wine preferences by data mining from physicochemical properties.  
In Decision Support Systems, Elsevier, 47(4):547-553. ISSN: 0167-9236.

Available at:

[@Elsevier] <http://dx.doi.org/10.1016/j.dss.2009.05.016>

[Pre-press (pdf)] <http://www3.dsi.uminho.pt/pcortez/winequality09.pdf>

[bib] <http://www3.dsi.uminho.pt/pcortez/dss09.bib>

in closing...