

Machine Learning in Python

Rohith Mohan

GradQuant

Spring 2018



Chet Haase

@chethaase

Follow

A Machine Learning algorithm walks into a bar.

The bartender asks, "What'll you have?"

The algorithm says, "What's everyone else having?"

8:24 AM - 1 Nov 2017

Interviewer: What's your biggest strength?

Me: I'm an expert in machine learning.

Interviewer: What's $9 + 10$?

Me: Its 3.

Interviewer: Not even close. It's 19.

Me: It's 16.

Interviewer: Wrong. Its still 19.

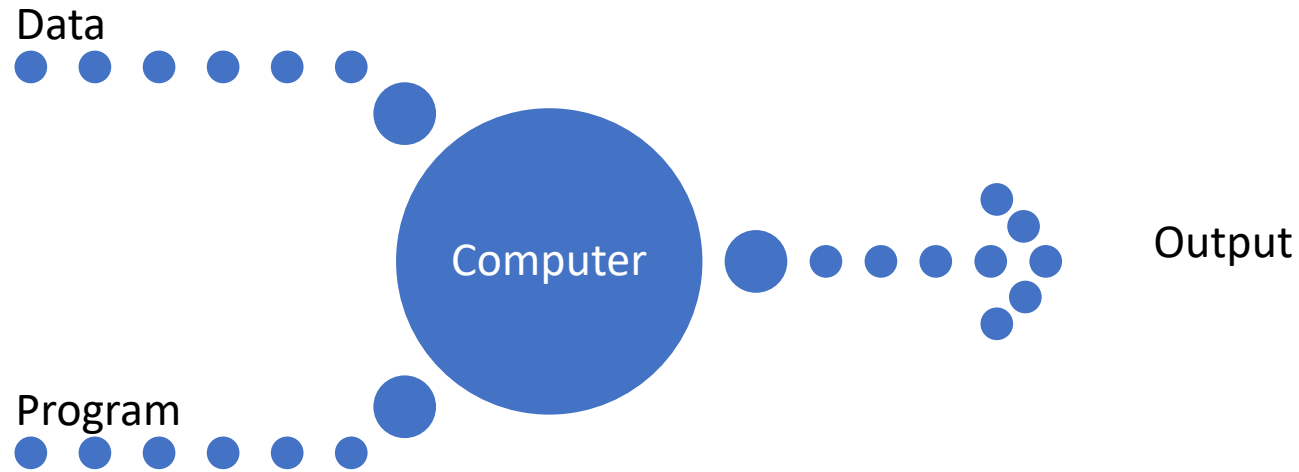
Me: It's 18.

Interviewer: No, it's 19.

Me: it's 19.

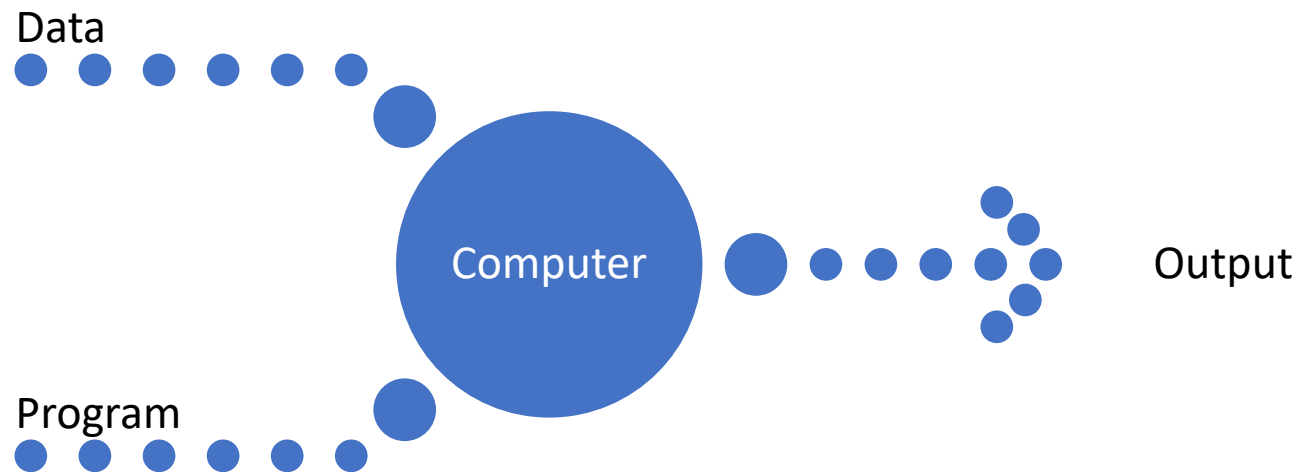
Interviewer: You're hired

Traditional Programming



- Getting computers to program themselves
- Coding is the bottleneck, let data dictate programming

Machine Learning



Formal Definitions

- Arthur Samuel (1959)
 - “Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.”
 - Created a program for computer to play itself in checkers (10000s games) and learn at IBM
- Tom Mitchell (1998)
 - “Well-posed Learning Problem: A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E .”

Machine Learning

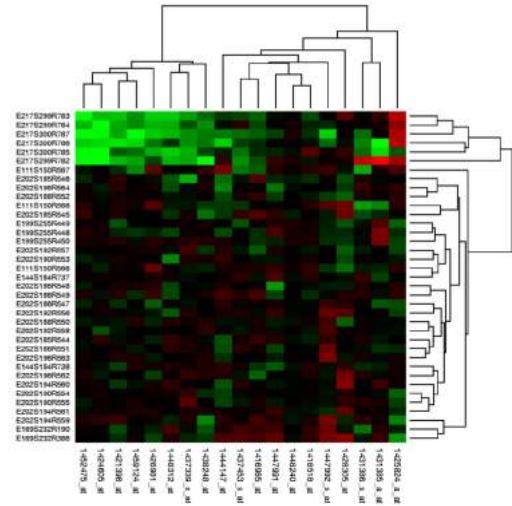
- Developed out of initial work in Artificial Intelligence (AI)
- Increased availability of large datasets and advances in computing architecture boosted usage in recent times

Overview [\[edit \]](#)

Decade ↕	Summary ↕
<1950s	Statistical methods are discovered and refined.
1950s	Pioneering machine learning research is conducted using simple algorithms.
1960s	Bayesian methods are introduced for probabilistic inference in machine learning. ^[1]
1970s	'AI Winter' caused by pessimism about machine learning effectiveness.
1980s	Rediscovery of backpropagation causes a resurgence in machine learning research.
1990s	Work on machine learning shifts from a knowledge-driven approach to a data-driven approach. Scientists begin creating programs for computers to analyze large amounts of data and draw conclusions – or “learn” – from the results. ^[2] Support vector machines (SVMs) and recurrent neural networks (RNNs) become popular.
2000s	Kernel methods grow in popularity, ^[3] and competitive machine learning becomes more widespread. ^[4]
2010s	Deep learning becomes feasible, which leads to machine learning becoming integral to many widely used software services and applications.

Usage

Natural Language Processing
+ Computer Vision



Mining and clustering
gene expression data to
identify individuals

These recommendations are based on items you own and more.

view: All | New Releases | Gaming_Seq

- The Singularity Is Near: When Humans Transcend Biology**
by Ray Kurzweil (Sep 26, 2006)
Average Customer Review: (112)
In Stock
List Price: \$16.00
Price: \$12.24
37 used & new from \$9.97
- Fantastic Voyage: Live Long Enough to Live Forever**
by Ray Kurzweil (Sep 27, 2005)
Average Customer Review: (63)
Available from [these sellers](#).
17 used & new from \$11.36
- Army of Darkness**
DVD ~ Ian Abercrombie (Aug 19, 1998)
Average Customer Review: (529)
In Stock
List Price: \$12.98
Price: \$8.99
82 used & new from \$2.99

Recommendation algorithms



Reproducing human
behavior (True AI)

<https://www.irishnews.com/magazine/science/2018/01/01/news/12-of-the-biggest-scientific-breakthroughs-of-2017-that-might-just-change-the-world-1222695/>

<http://www.idownloadblog.com/2016/05/12/google-translate-offline-mode/>

<https://www.flickr.com/photos/theadamclarke/2589233355> <https://de.wikipedia.org/wiki/Genexpressionsanalyse>

Common steps in ML workflow

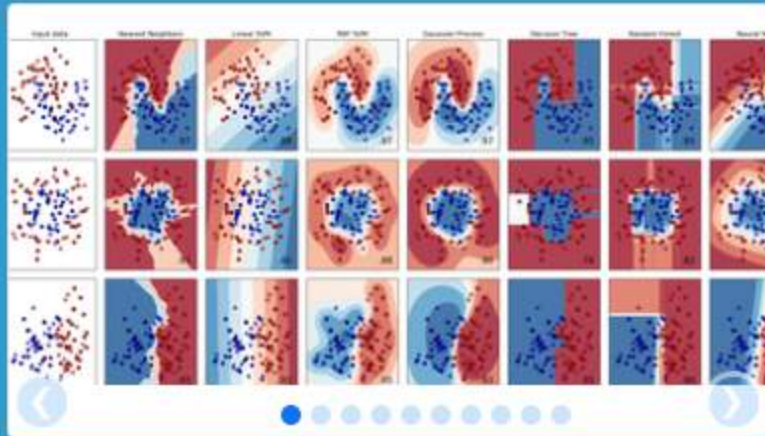
- Collect data (various sources, UCI data repository, news orgs, Kaggle)
- Prepare data (exploratory analysis, feature selection, regularization)
- Selecting and training model (train and test datasets, what model?)
- Evaluating model (accuracy, precision, ROC curves, F1 score)
- Optimizing performance (change model, # of features, scaling)

scikit-learn



Home Installation Documentation ▾ Examples

Google Custom Search



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Preprocessing

- Clean data and deal with missing values, etc.
- Feature scaling - rescaling features to be more sensible
- Standardization - getting various features into similar range (e.g. -1 to 1)
 - Square footage of a house (100s of ft) vs # of rooms (1-5)

```
>>> from sklearn import preprocessing
>>> import numpy as np
>>> X_train = np.array([[ 1., -1.,  2.],
...                    [ 2.,  0.,  0.],
...                    [ 0.,  1., -1.]])
>>> X_scaled = preprocessing.scale(X_train)

>>> X_scaled
array([[ 0. ...., -1.22....,  1.33....],
       [ 1.22....,  0. ...., -0.26....],
       [-1.22....,  1.22...., -1.06....]])
```

```
>>> X_scaled.mean(axis=0)
array([ 0.,  0.,  0.])

>>> X_scaled.std(axis=0)
array([ 1.,  1.,  1.] )
```

Preprocessing

- Clean data and deal with missing values, etc.
- Feature scaling - rescaling features to be more sensible
- Standardization - getting various features into similar range (e.g. -1 to 1)
 - Square footage of a house (100s of ft) vs # of rooms (1-5)
- Normalization – scaling to some standard (e.g. subtract mean & divide by SD)
- Many others (regularization, imputation, generating polynomial features, etc.)

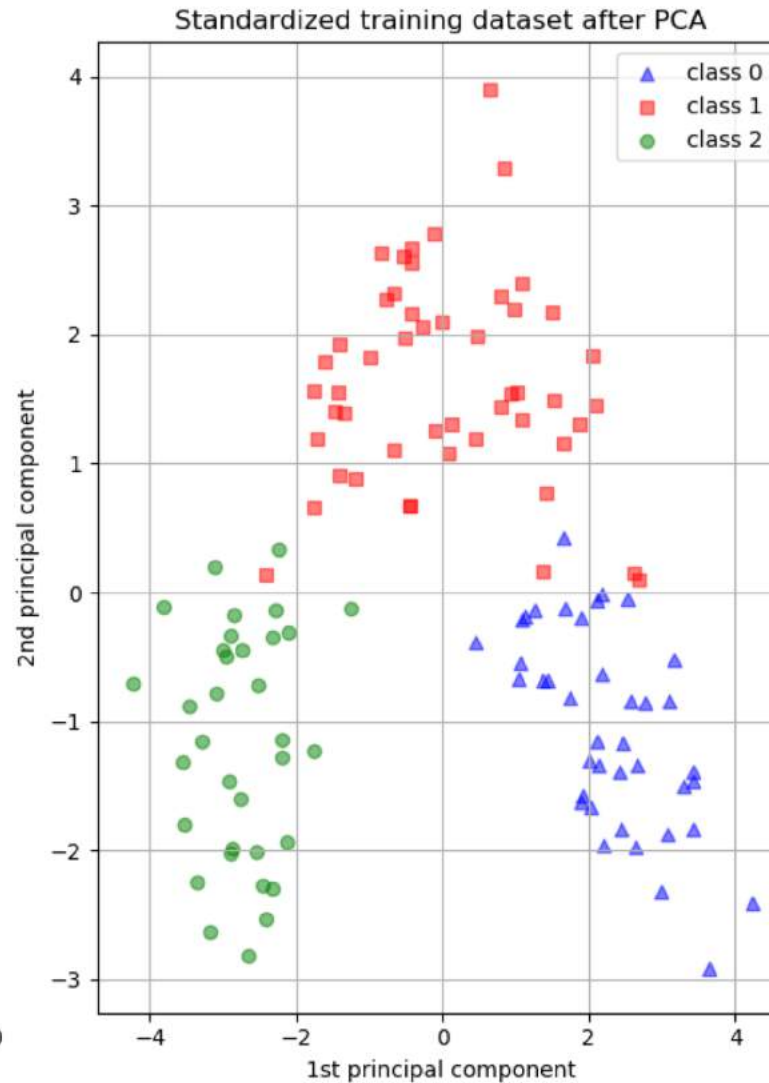
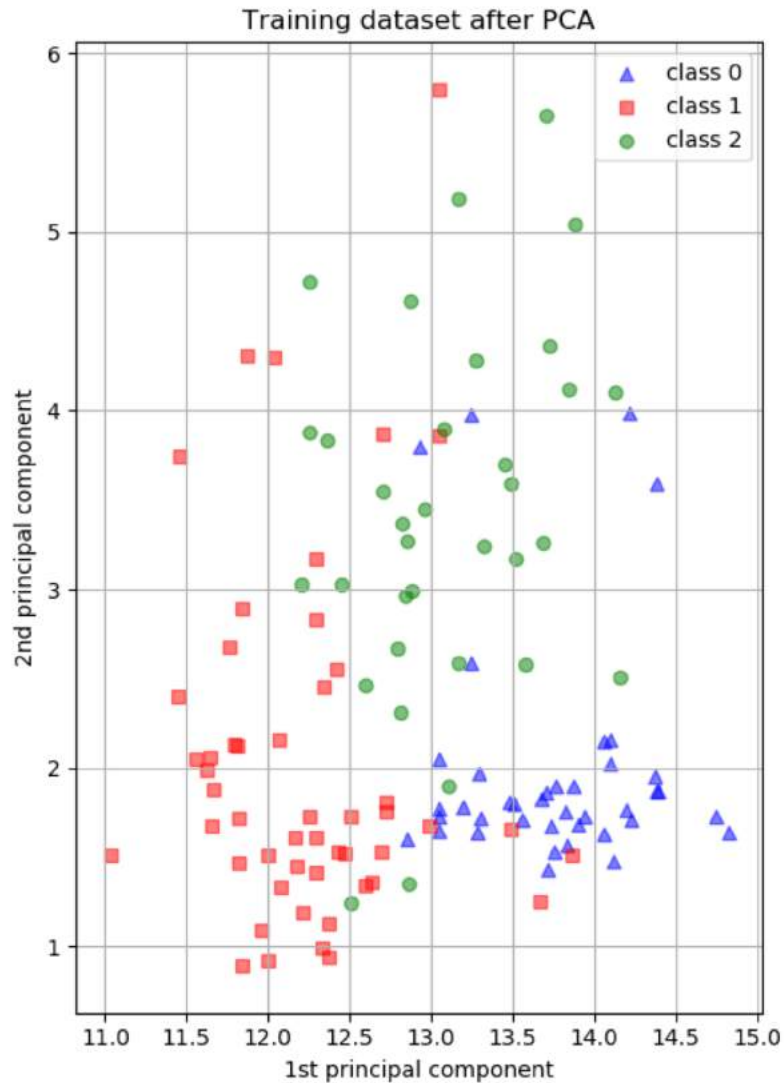
Preprocessing

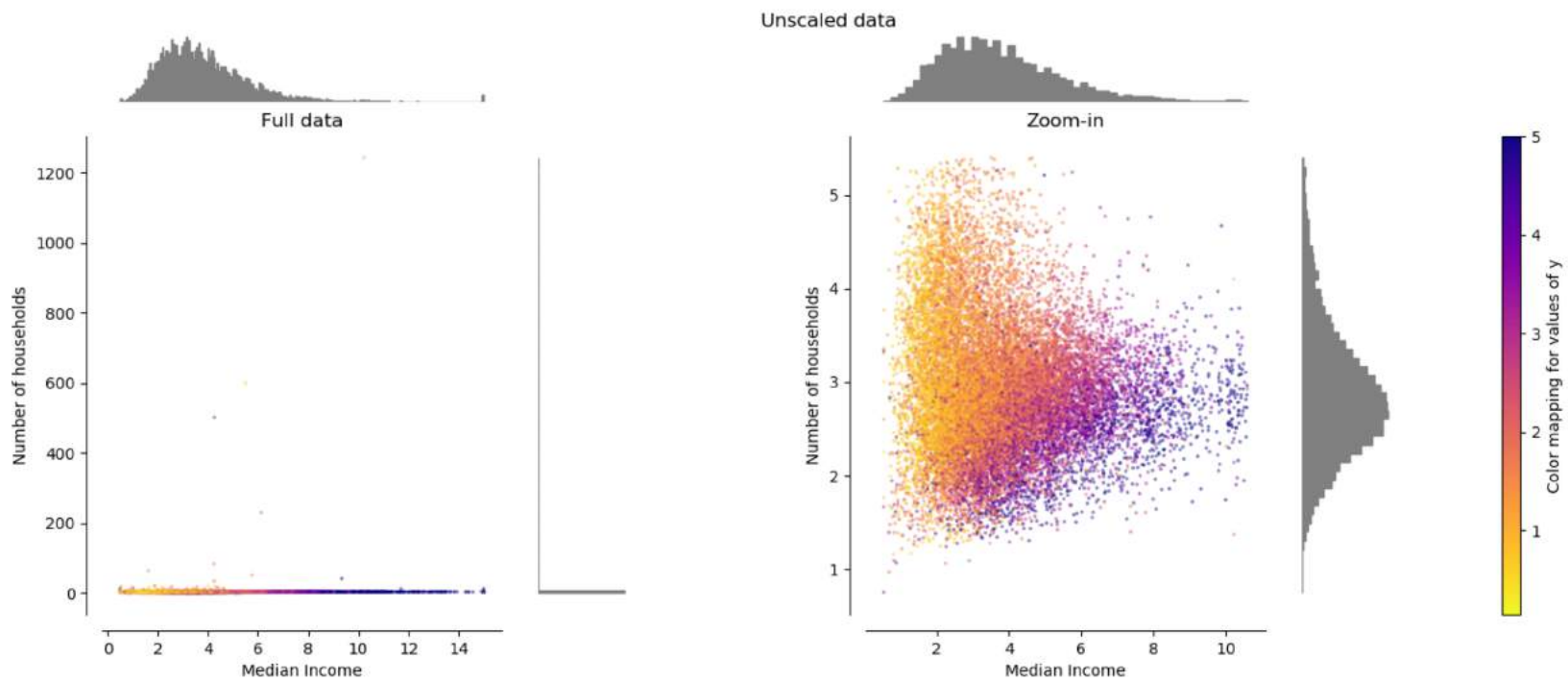
- Clean data and deal with missing values, etc.
- Feature scaling - rescaling features to be more sensible
- Standardization - getting various features into similar range (e.g. -1 to 1)
 - Square footage of a house (100s of ft) vs # of rooms (1-5)
- Normalization – scaling to some standard (e.g. subtract mean & divide by SD)

```
>>> X = [[ 1., -1.,  2.],
...      [ 2.,  0.,  0.],
...      [ 0.,  1., -1.]]
>>> X_normalized = preprocessing.normalize(X, norm='l2')
```

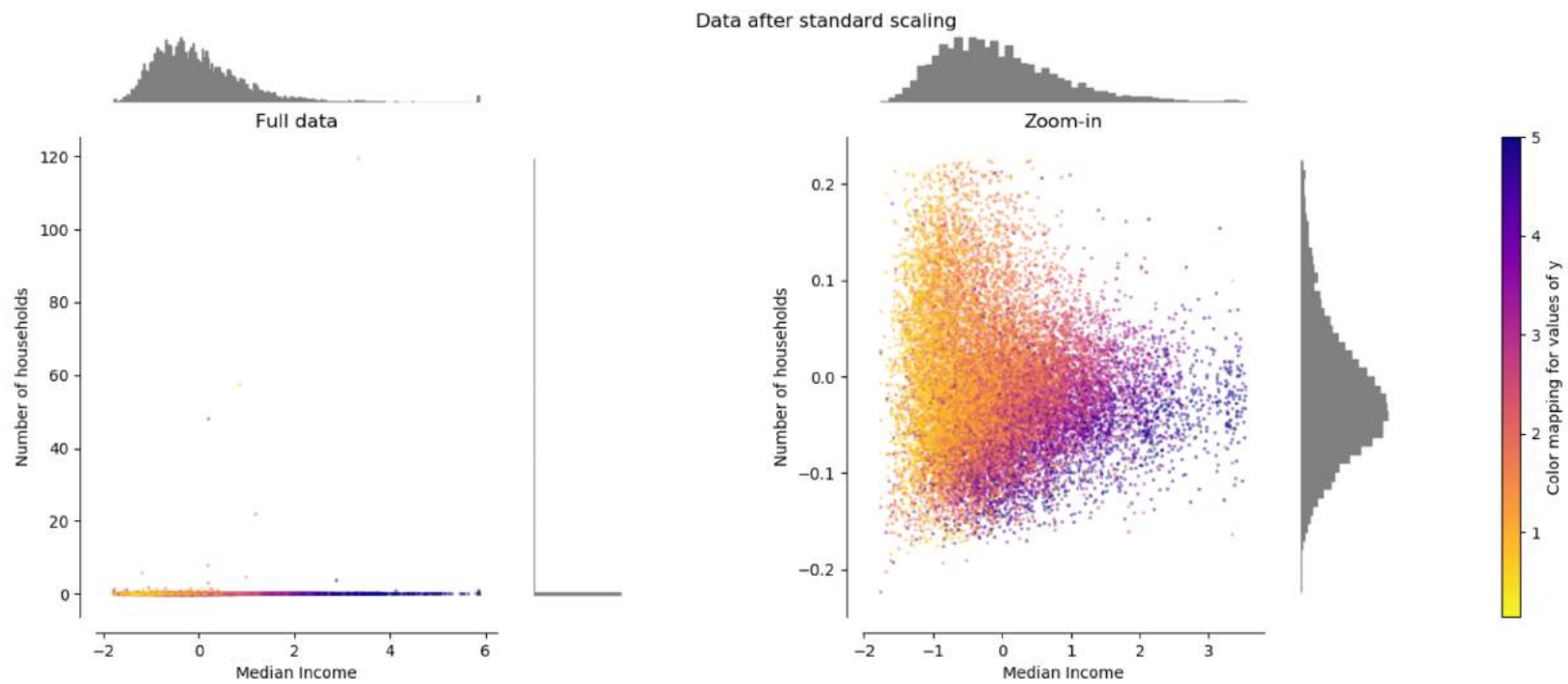
```
>>> X_normalized
array([[ 0.40..., -0.40...,  0.81...],
       [ 1. ....,  0. ....,  0. ....],
       [ 0. ....,  0.70..., -0.70...]])
```

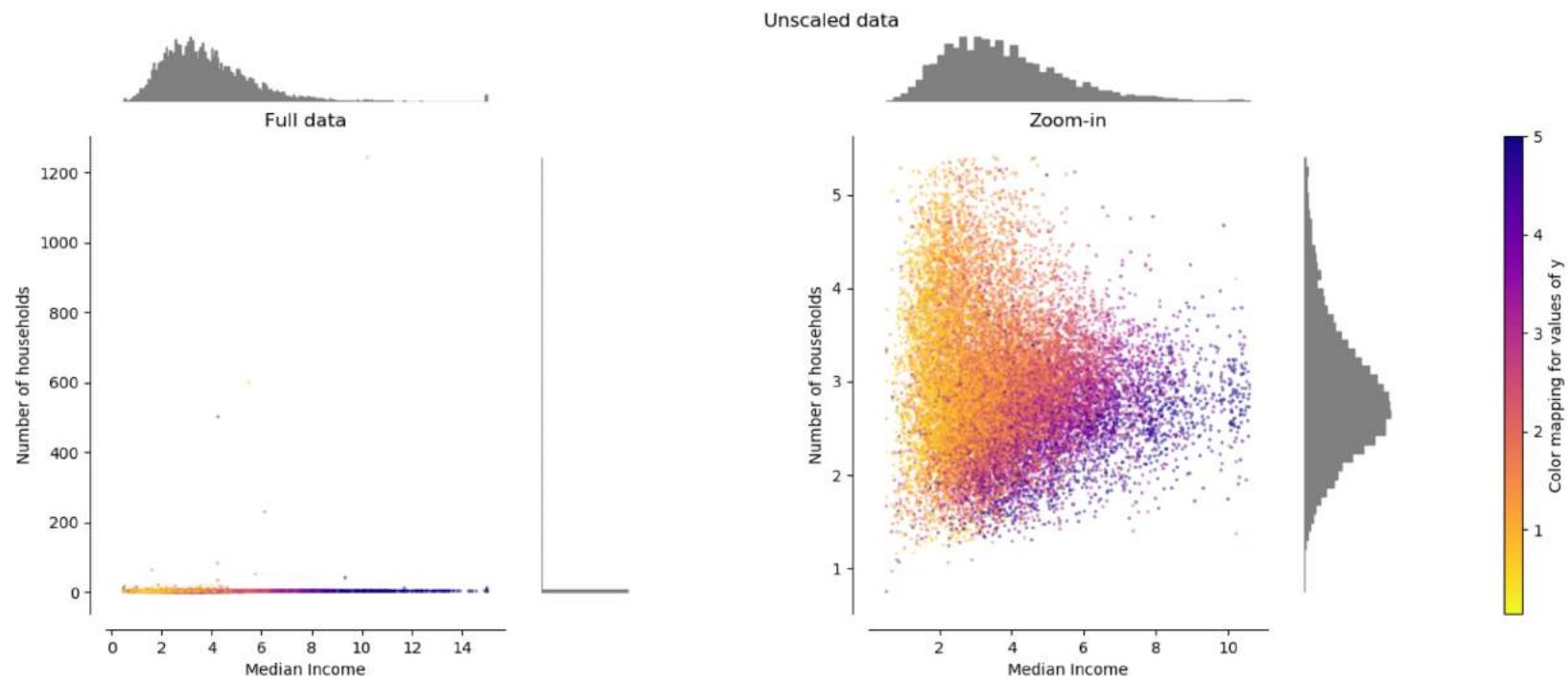
Importance of feature scaling



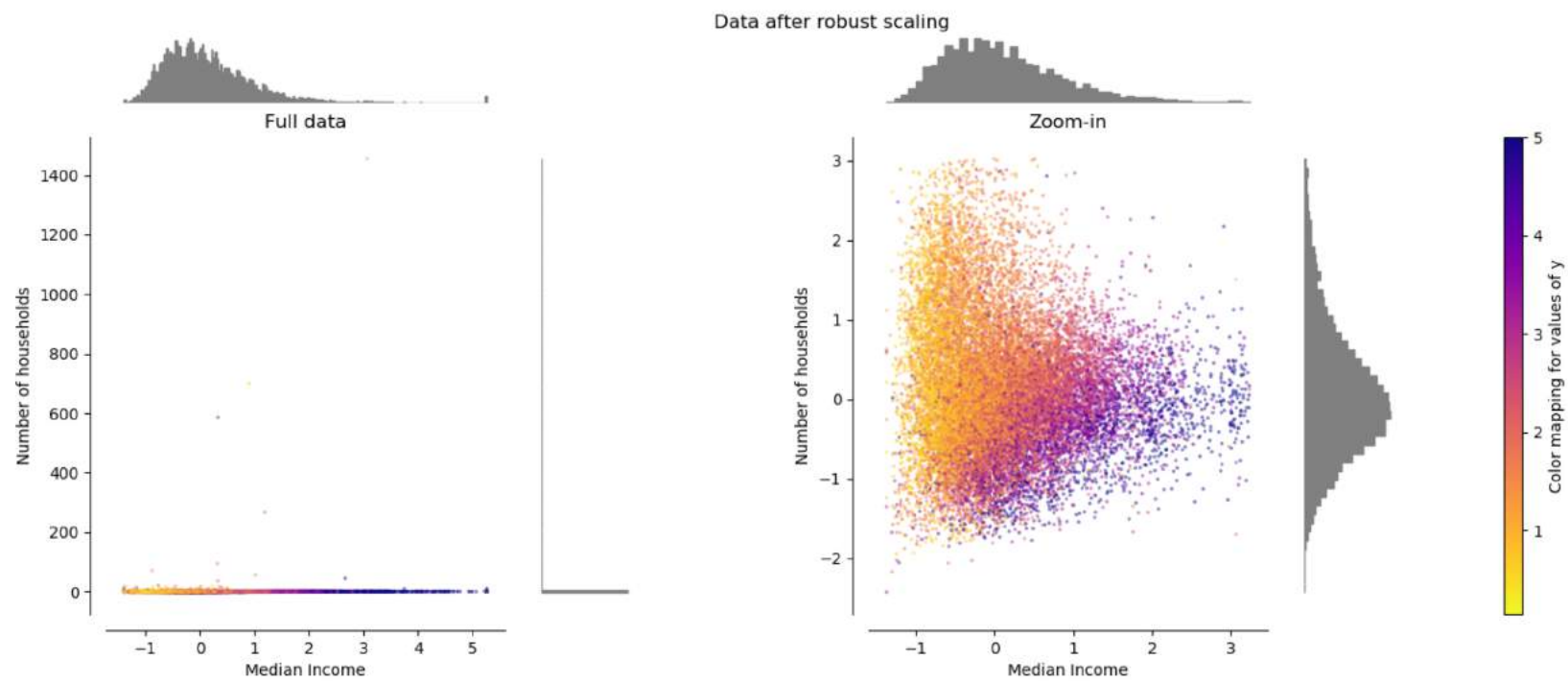


StandardScaler





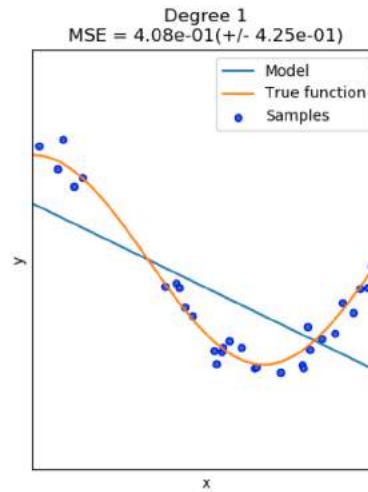
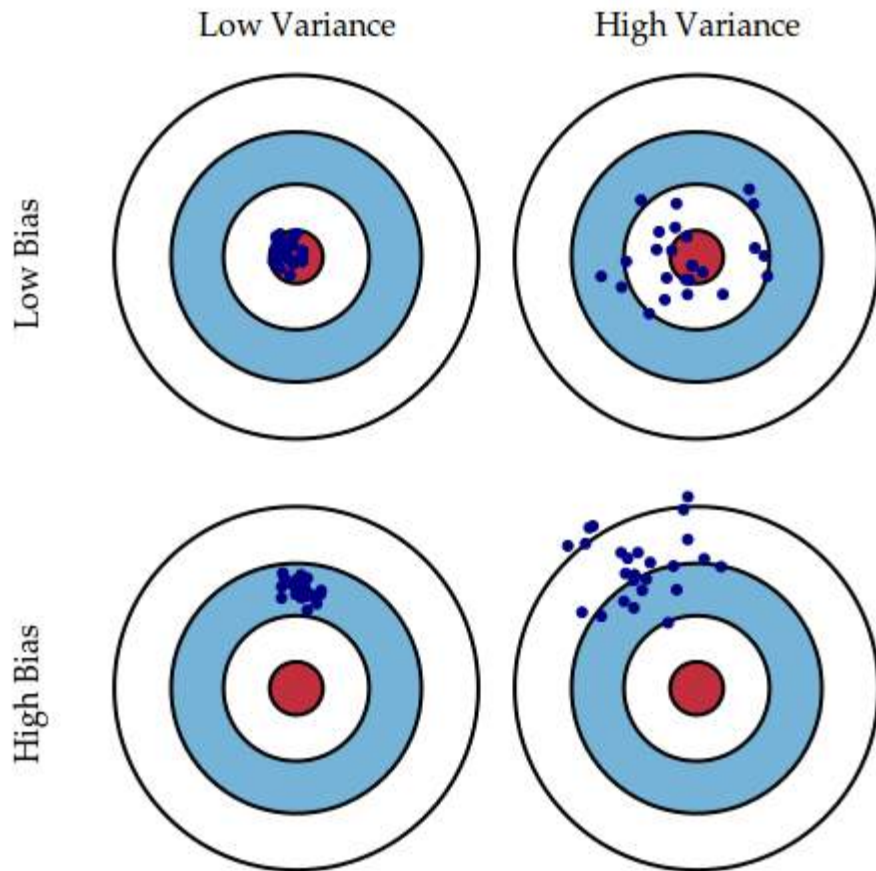
RobustScaler



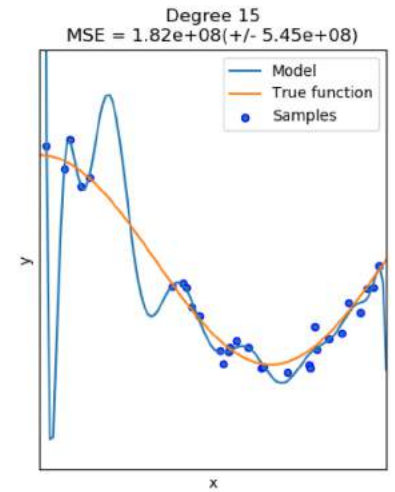
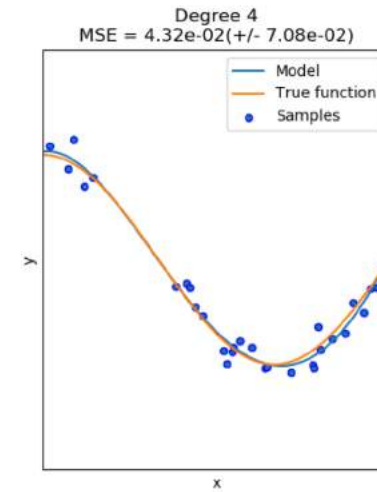
Train Test (Cross Validate?)

- Why do we need to split up our datasets?
 - Overfitting
- Split dataset
 - Train – for training your model on
 - Test – evaluate performance of model
 - Usually 40% for testing is enough
- Validation set?
- Cross-validation
 - Split up training set into subsets and evaluate performance (can be more computationally expensive but conserves data)
- Hyper-parameter tuning

Bias-variance tradeoff

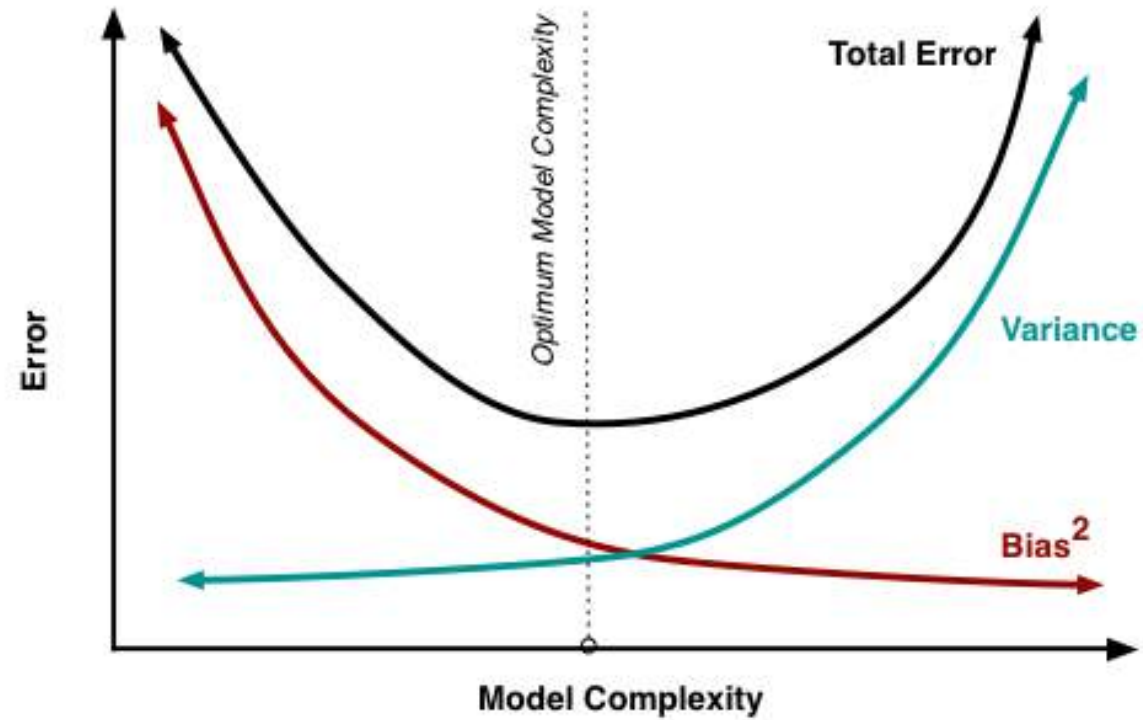


Underfitting
High Bias



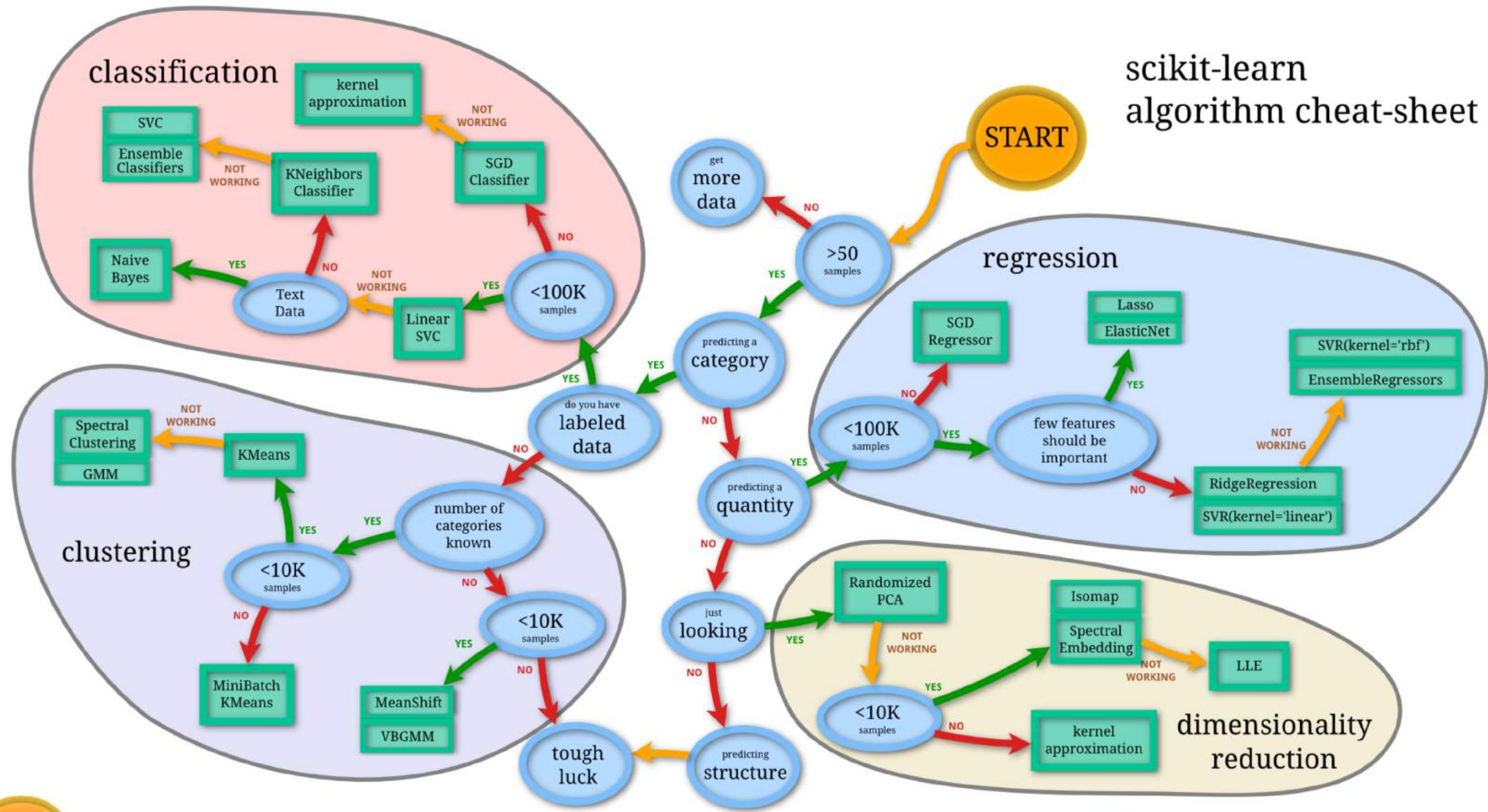
Overfitting
High Variance

Bias-variance tradeoff



How to select a model?

scikit-learn algorithm cheat-sheet



Supervised vs Unsupervised Learning

- Supervised
 - Regression, classification
 - Input variables, output variable, learn mapping of input to output
- Unsupervised
 - Clustering, association, etc.
 - No correct answers and no teacher
- Semi-supervised
 - Partially labeled dataset of images
 - Mixing both techniques is what occurs in real-world

Regression

- Linear regression (OLS)

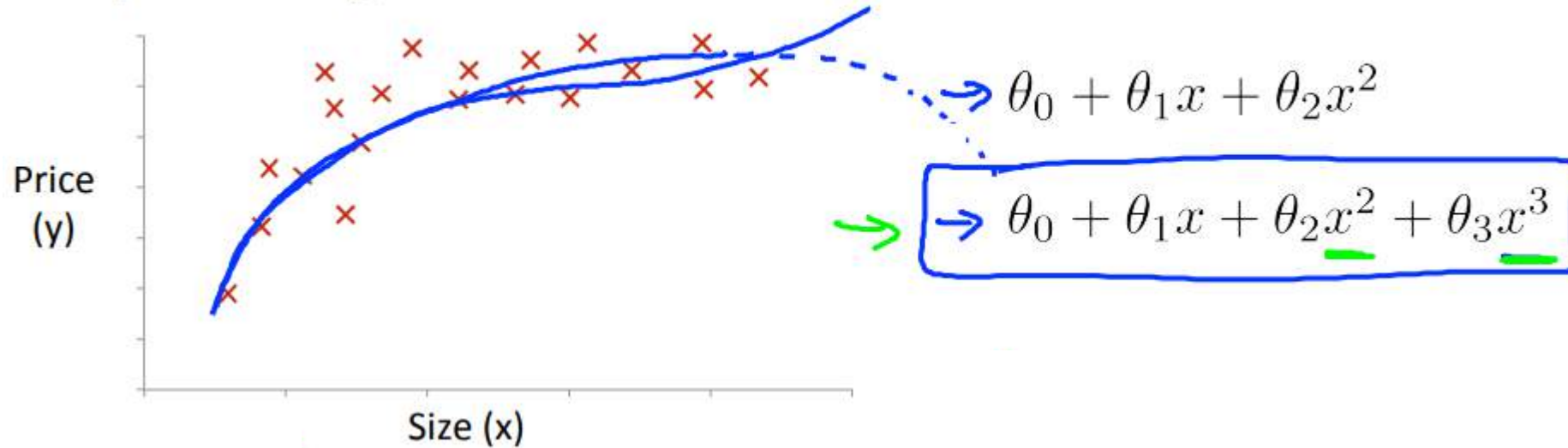
$$Y = \beta_0 + \sum_{j=1..p} \beta_j X_j + \varepsilon$$

- Prediction
- Multiple variables/features
 - Feature selection

```
>>> from sklearn.feature_selection import VarianceThreshold
>>> X = [[0, 0, 1], [0, 1, 0], [1, 0, 0], [0, 1, 1], [0, 1, 0], [0, 1, 1]]
>>> sel = VarianceThreshold(threshold=(.8 * (1 - .8)))
>>> sel.fit_transform(X)
array([[0, 1],
       [1, 0],
       [0, 0],
       [1, 1],
       [1, 0],
       [1, 1]])
```

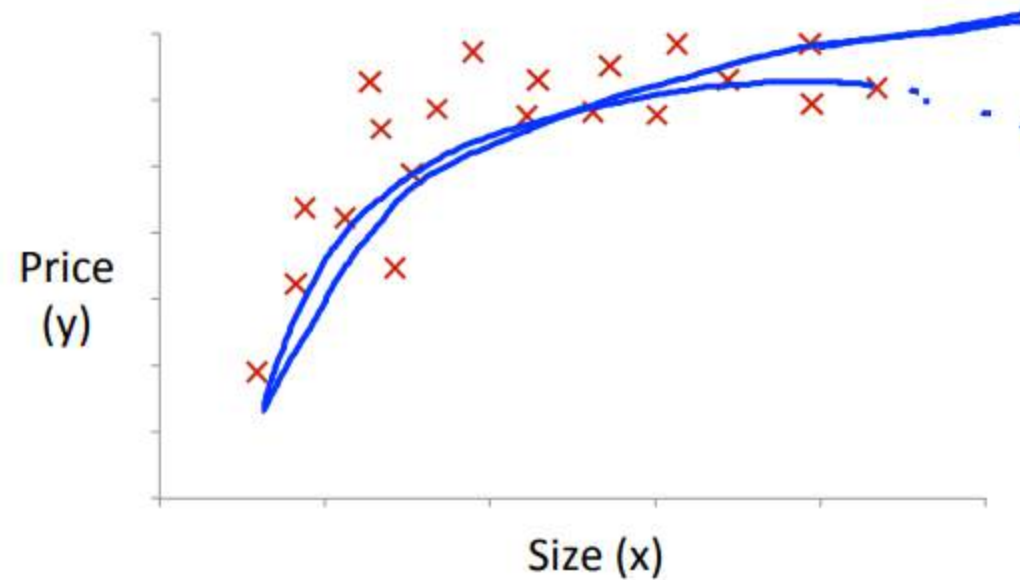
Feature Selection

Polynomial regression



Feature Selection

Choice of features



$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

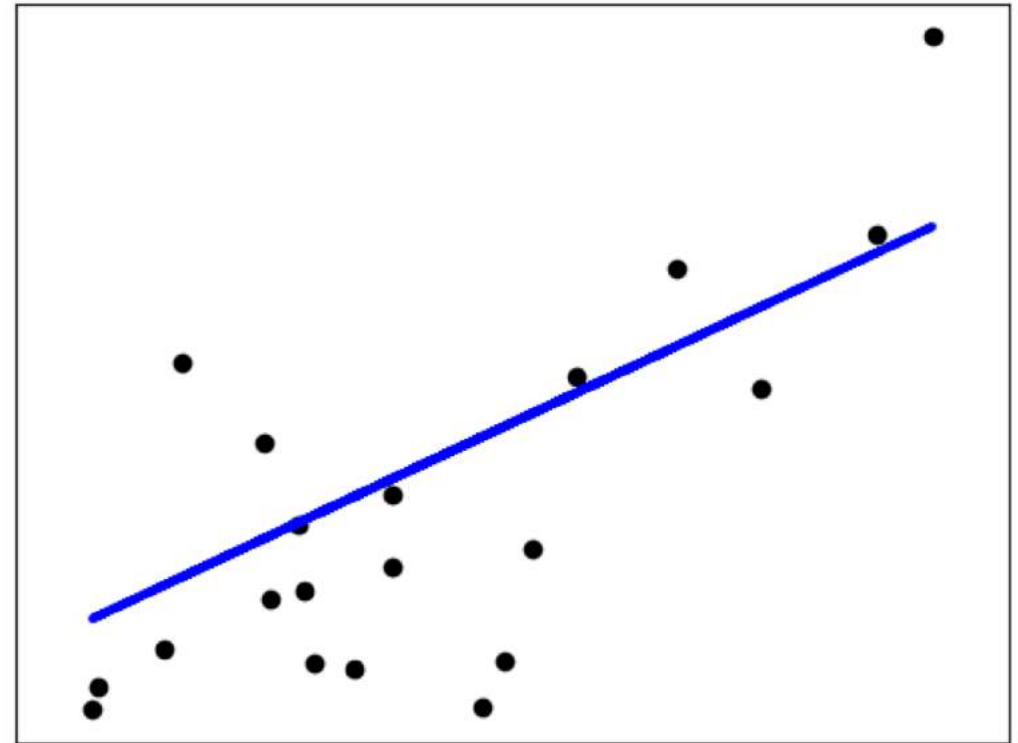
$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2\sqrt{(\text{size})}$$

Regression

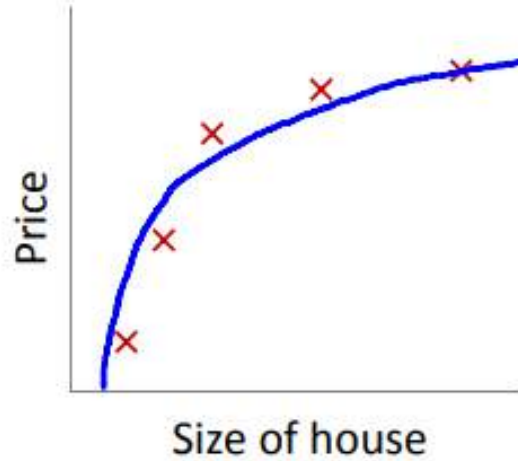
- Linear regression (OLS)

$$Y = \beta_0 + \sum_{j=1..p} \beta_j X_j + \varepsilon$$

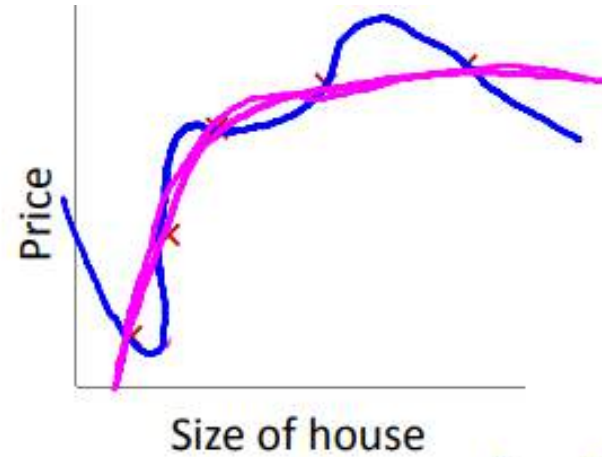
- Prediction
- Multiple variables/features?
 - Feature selection
 - Length, width of a house (area?)
 - Regularization



Regularization

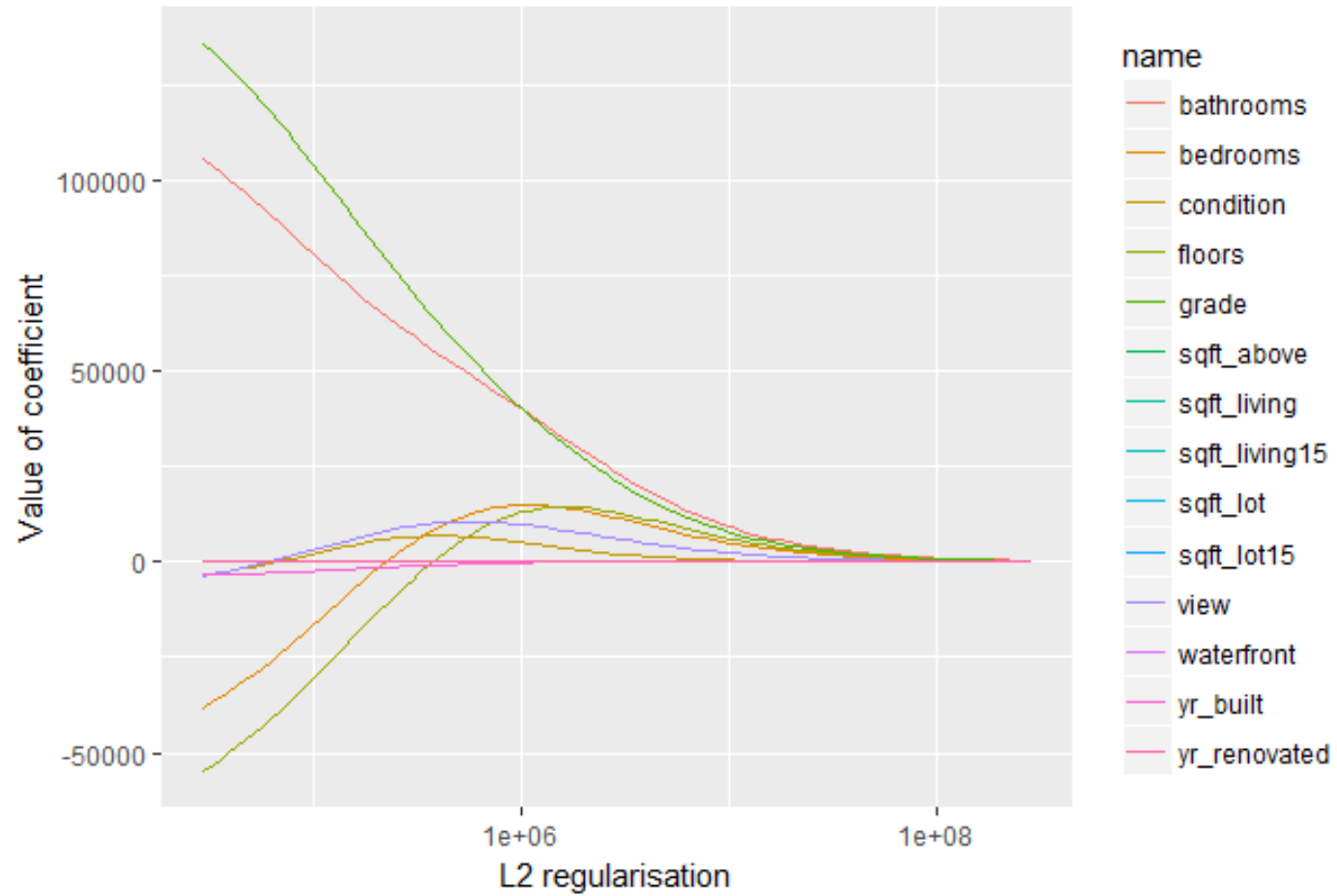


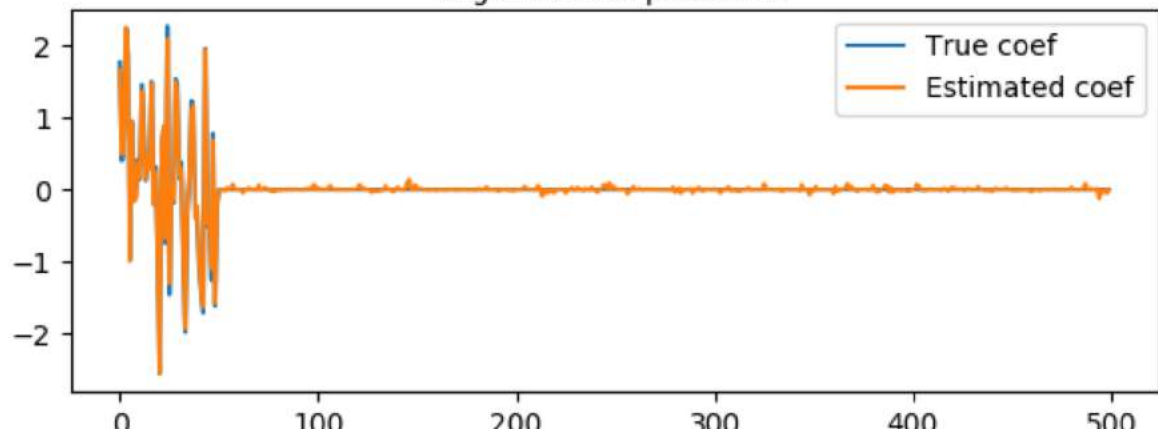
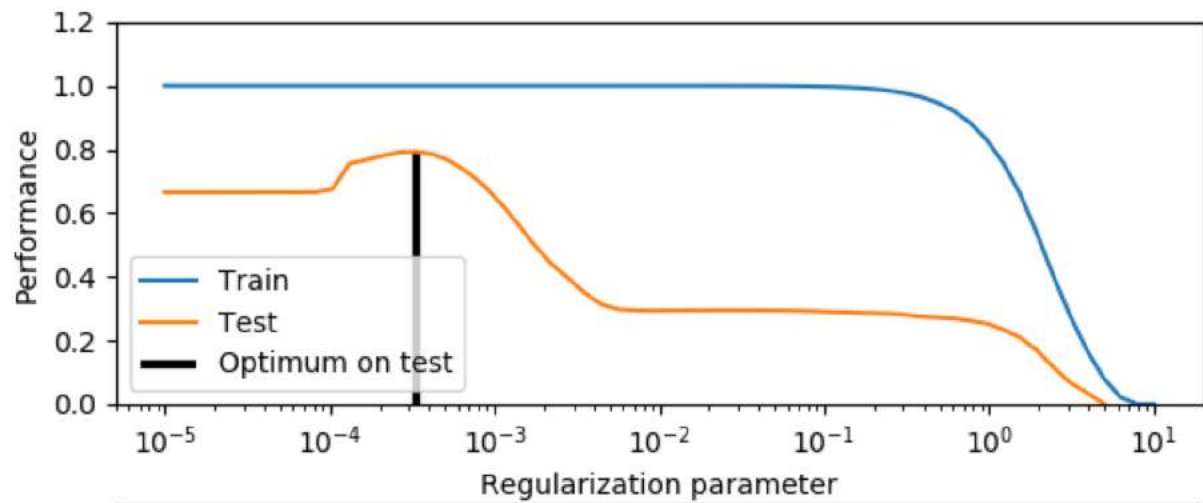
$$\theta_0 + \theta_1 x + \theta_2 x^2$$



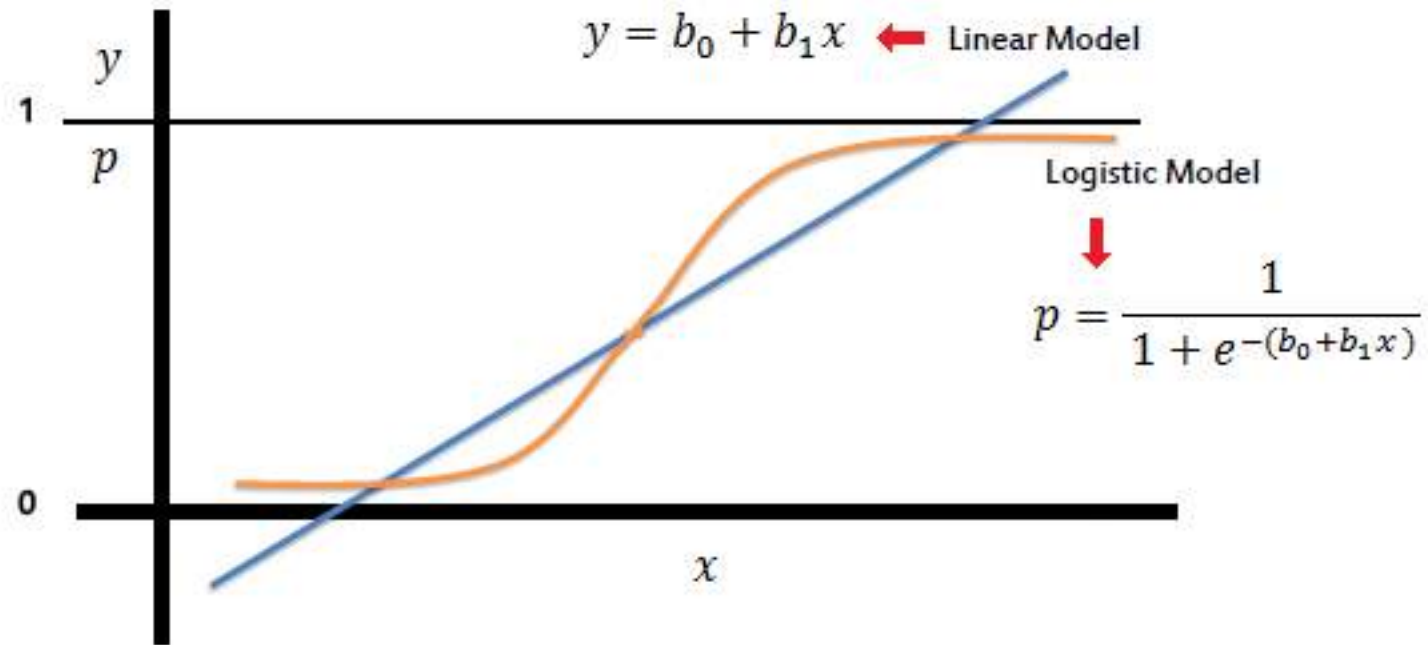
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \cancel{\theta_3 x^3} + \cancel{\theta_4 x^4}$$

Regularization



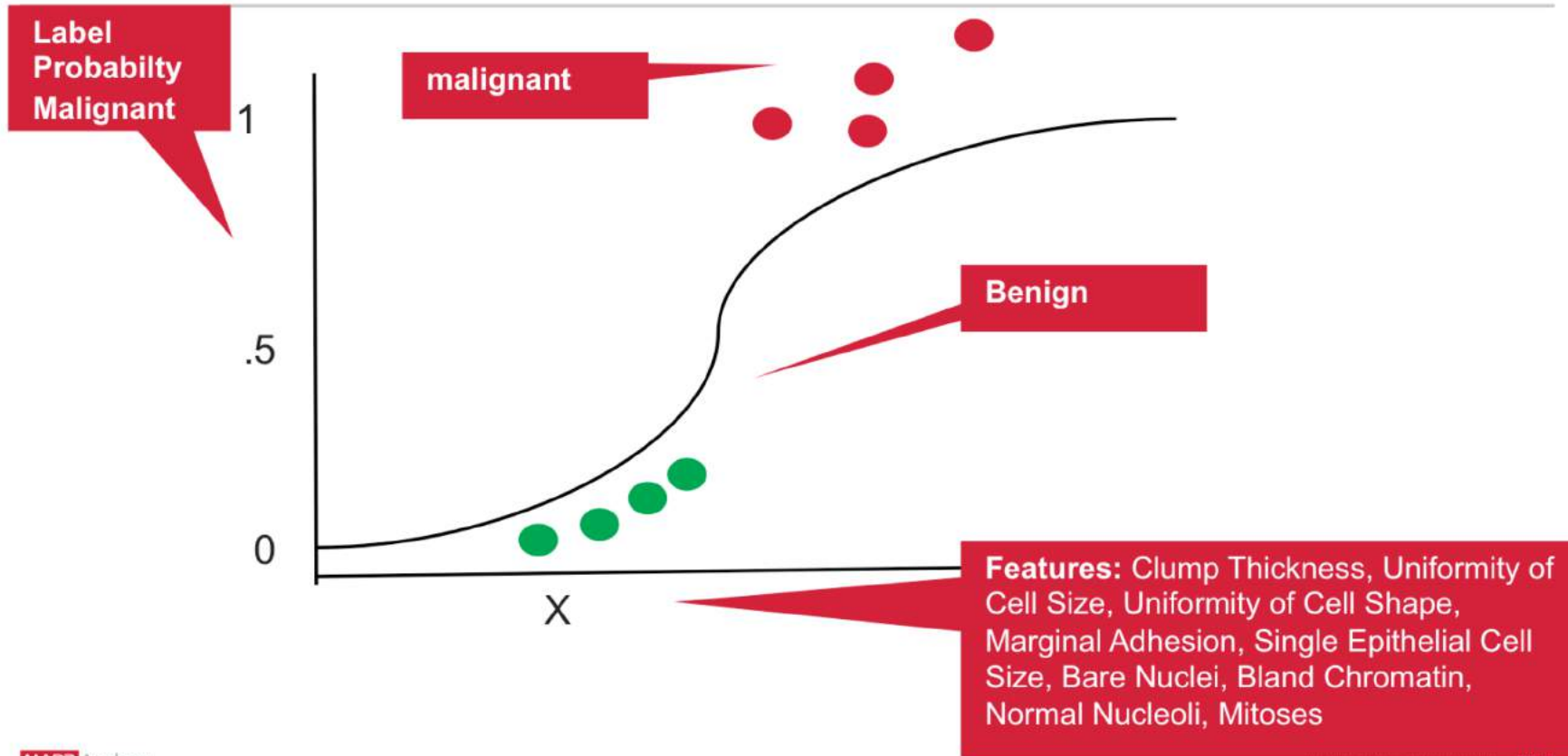


Classification – Logistic Regression

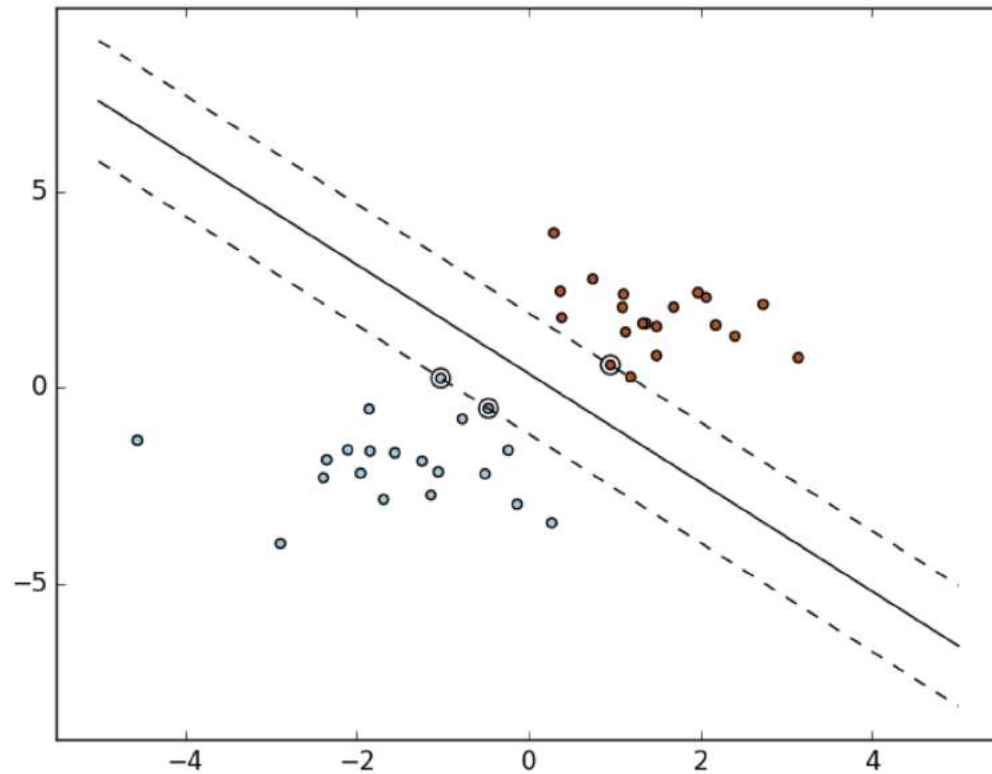


Classification – Logistic Regression

Breast Cancer Logistic Regression Example

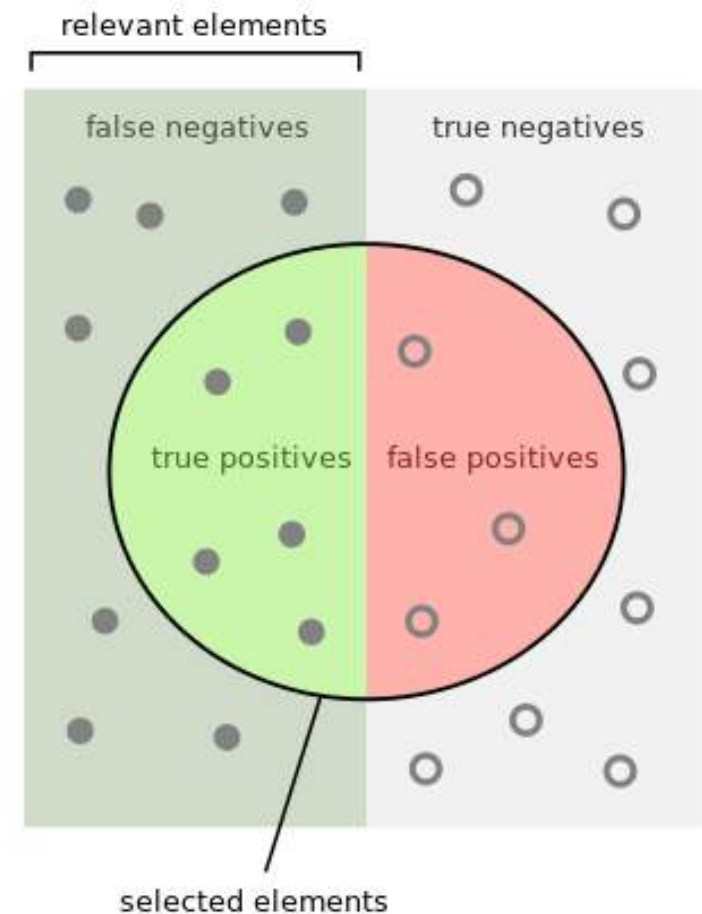


Classification – SVM



Evaluating Performance

- Accuracy – how many predictions are correct on a dataset?
 - Can be a flawed metric
- Precision and Recall



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

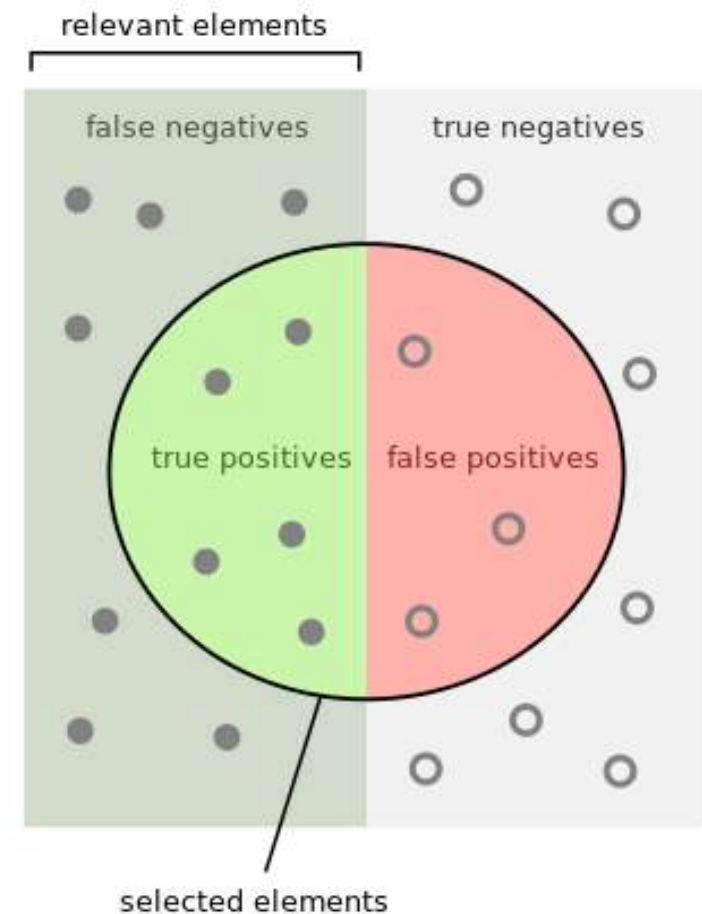
How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Evaluating Performance

- Accuracy – how many predictions are correct on a dataset?
 - Can be a flawed metric
- Precision and Recall
- ROC curves
- F1 score

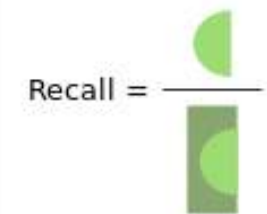
$$F_1 = \frac{2}{\frac{1}{\text{recall}} + \frac{1}{\text{precision}}} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$



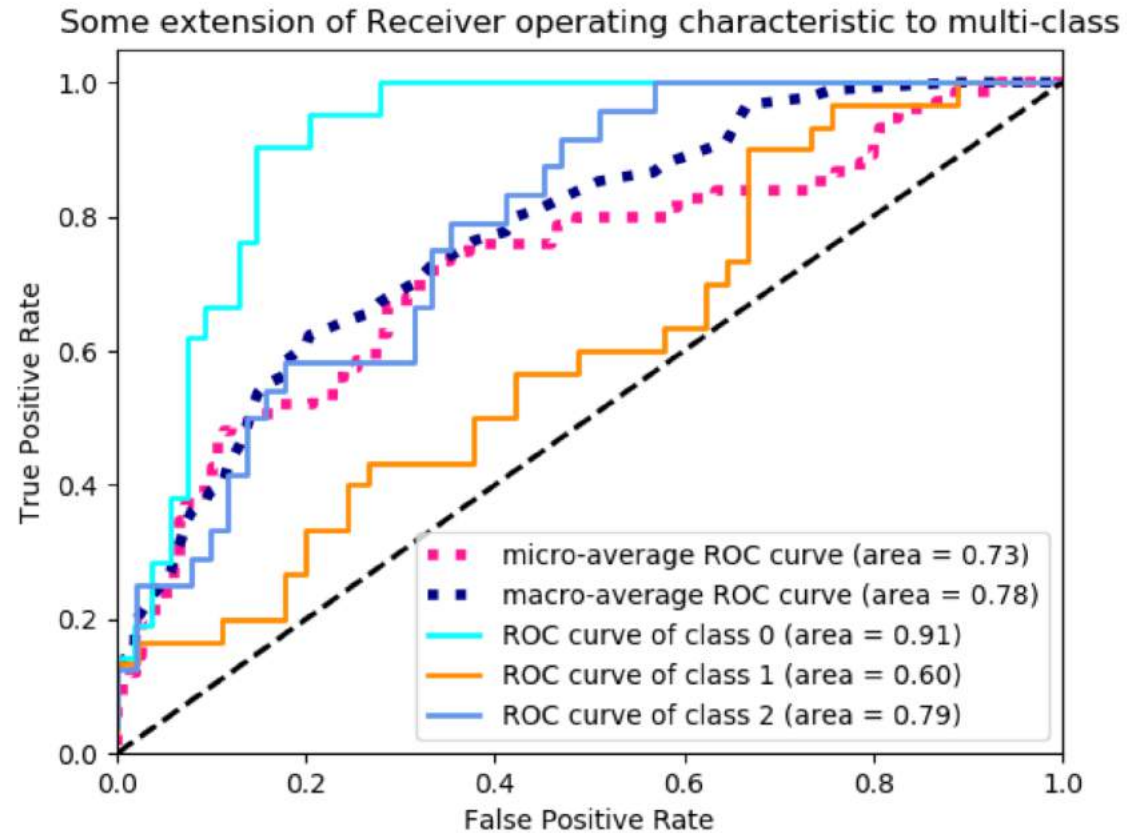
How many selected items are relevant?



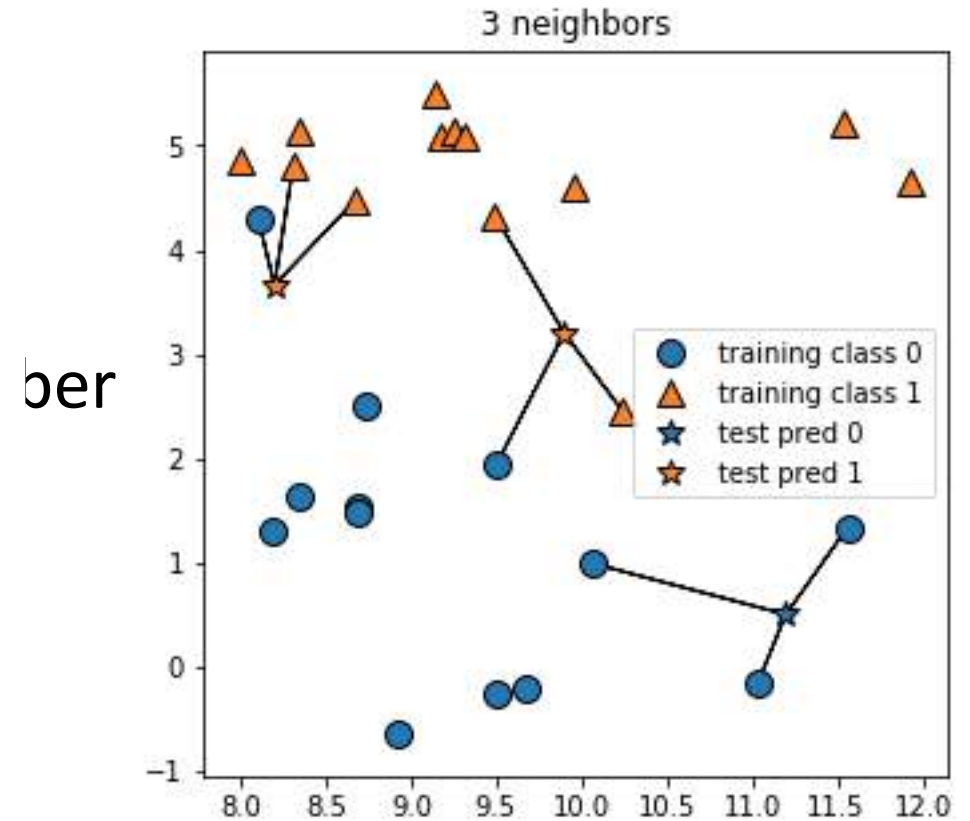
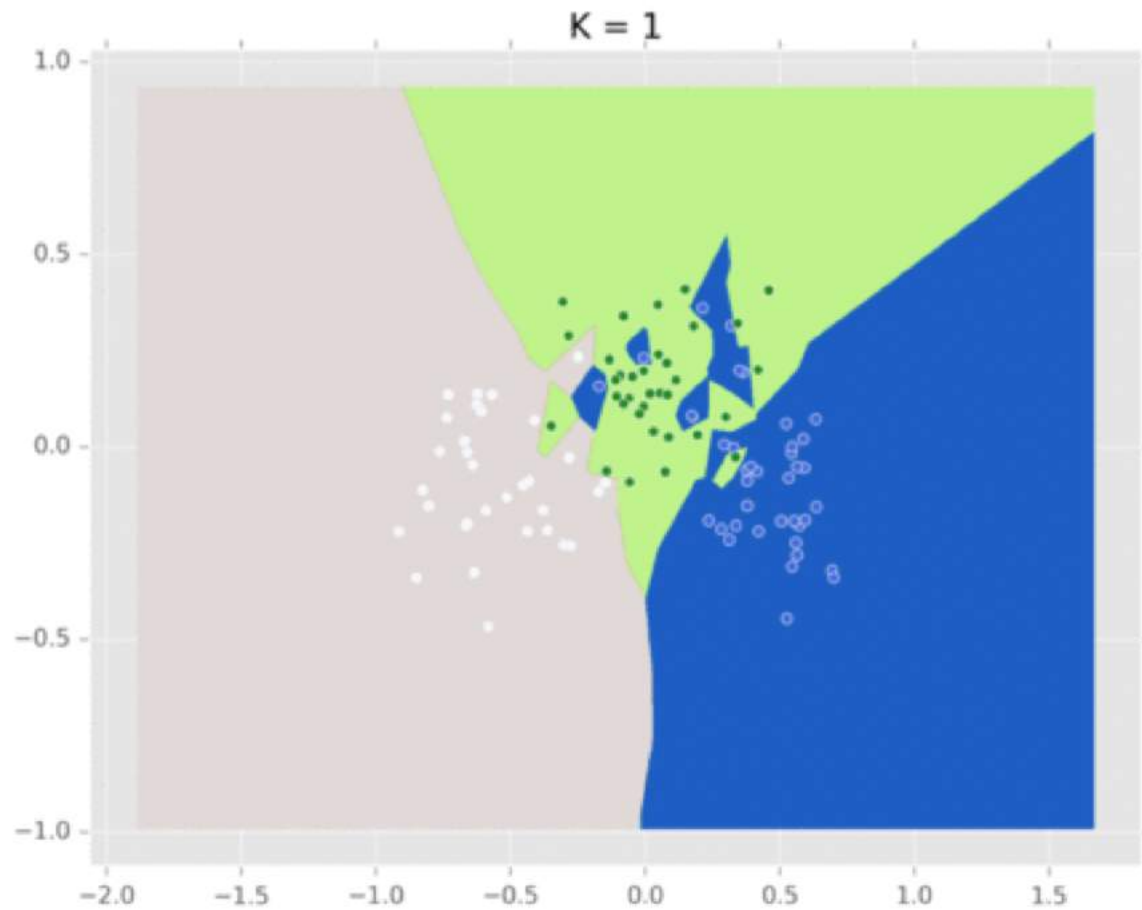
How many relevant items are selected?



Evaluating Performance



Classification - K-Nearest Neighbors



Clustering

- Unsupervised learning
- Can help you understand structure of your data
- Various types of clustering: K-means, Hierarchical, Ward

K-means

- Randomly choose k centroids
- Form clusters around it
- Take mean of cluster to identify new centroid
- Repeat until convergence

