Naive Bayes: simple boundary, apply Bayes' theorem with the "naive" assumption of independence between every pair of features

1. efficient

http://scikit-learn.org/stable/modules/naive\_bayes.html

sensitivity: Pr (positive| have it)

**specitivity**: Pr (negative| do not have it)

<u>Support Vector Machine (SVM)</u>: maximize the margin (distance to the nearest point) -> robustness

**Kernel trick:**  $(x, y) \rightarrow (x_1, \ldots, x_m)$  not separable  $\rightarrow$  separable

SVM to prevent over-fitting:

- 1. Kernel
- 2. <u>C</u>: control trade-off between smoothing decision boundary and classifying training points correctly
- 3. gamma

Linear kernel only gives a linear decision boundary.

Entropy: how a decision tree decides where to split the data

- 1. measure of impurity in a bunch of examples
- 2. find the threshold where the split of data is as pure as possible

**Information Gain:** entropy (parent) - weighted average of entropy (children)

- 1. Decision tree maximizes information gain
- 2. at the start of the data, entropy (parent) is 1

# **Algorithms:**

- 1. K-nearest neighbor
- 2. random forest (ensemble method) -> fit classifier to sub-tree and average result

3. boosted decision tree (adaboost) (ensemble method) -> fit a base classifier and then adjust to fit the mis-classified data

## **More data > Fine-Tuned Algorithm!!**

### **Types of data:**

- 1. numerical
- 2. categorical number (limited number of discrete values)
- 3. time-series (temporal, data, time)
- 4. text data

### **Outlier Rejection:**

- 1. train all the data
- 2. remove the data with highest residual error -> 10%
- 3. train the data gain
- 4. maybe repeat 1 to 3

Feature Scaling:  $x' = \frac{x - x_{min}}{x_{max} - x_{min}} \in [0, 1]$ , prone to outlier

Bag of Words: take the words and count the frequency Stop words (low-information words occur frequently)

## **Feature Selection**

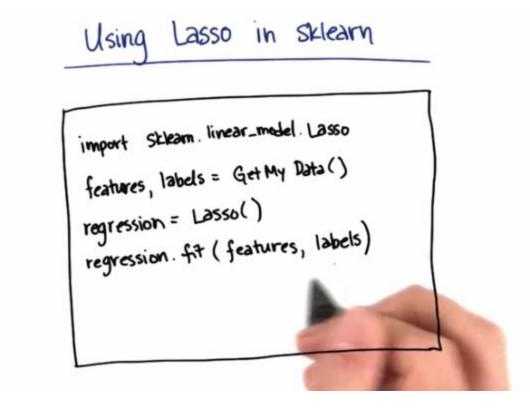
- 1. select best features
- 2. add new features

<u>Univariate Feature Selection:</u> treats each feature independently and asks how much power it gives you in classifying or regressing.

sklearn -> feature selection

from sklearn.feature selection import \*\*\*

Lasso Regression: min SSE + lambda \* |\beta| (|\beta| is # of features)



<u>Difference between Feature Selection and PCA</u>: PCA makes composite features, while feature selection eliminates some features.

**variance**: the spread of a data distribution

<u>Principal Component</u>: the direction with the largest variance

```
def doPCA():
    from sklearn.decomposition import PCA
    pca = PCA(n_components=2)
    pca.fit(data)
    return pca

pca = doPCA()
print pca.explained_variance_ratio_
first_pc = pca.components_[0]
second_pc = pca.components_[1]

transformed_data = pca.transform(data)
for ii, jj in zip(transformed_data, data):
    plt.scatter( first_pc[0]*ii[0], first_pc[1]*ii[0], co r="r")
    plt.scatter( second_pc[0]*ii[1], second_pc[1]*ii[1], co or="c")
plt.xlabel("bonus")
plt.xlabel("bonus")
plt.ylabel("long-term incentive")
plt.show()
```

## K-fold cross validation:

- 1. run K times
- 2. choose one as test set, the others as train sets

Accuracy is not good for skewed class (only very few data points are of interest)

#### Confusion Matrix:

(2 by 2 matrix, actual positive/negative v.s. predict positive/negative) -> <u>for</u> skewed class

- 1. **false alarm**: true negative, predict positive (n by n matrix, where n is # of features)
  - 1. **recall**: if true, the probability of identifying it to be true
  - 2. **precision**: if predicted to be true, the probability of it to be true

### positive for actual, true or false for prediction

PREDICTED Ariel Sharon 1 1 [ 13 4 1] Colin Powell 0 55 0] 1 8 2] Donald Rumsfeld 25 0 123 1] George W Bush 1 7 Gerhard Schroeder 14 4] 2 1 10 **Hugo Chavez** 0] Tony Blair 26]] Precision = true positives

true positives + false pos.

true positives + false negatives

true positives + false negatives Positives X TRUE POSITIVES 0 FALSE NECTIVES X FALSE