

## SL: PUTTING IT ALL TOGETHER

### Test information

- The second test is next week Monday (12/1)
- In-class, closed book, closed notes
- Similar to previous test
- Not comprehensive – starts at HMM
- List of topics posted on Piazza

## Test information

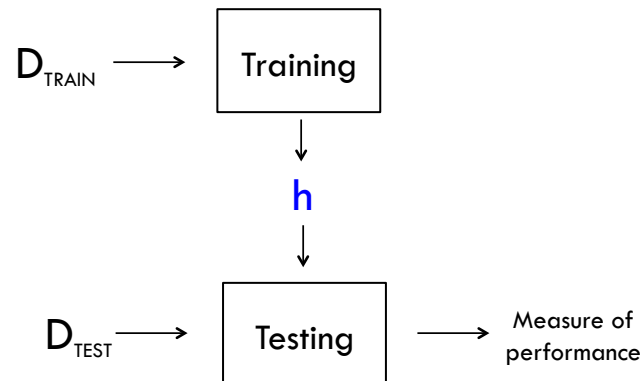
- Covered
  - ▣ HMM, filtering, smoothing, particle filtering
  - ▣ Supervised learning, naïve Bayes
  - ▣ Decision trees, neural networks, support vector machines
  - ▣ Clustering
- Not Covered
  - ▣ Prediction, Most likely explanation, Viterbi Algorithm
  - ▣ Won't ask you to derive Backpropagation/SVMs
  - ▣ No calculator needed

## Today

- Reading
    - ▣ AIMA 18.4 (Cross-validation)
    - ▣ AIMA 18.10-18.11 (Ensembles)
  - Goals
    - ▣ Step 1: Formulating the problem
    - ▣ Step 2: Exploring the data
    - ▣ Step 3: Feature Selection
    - ▣ Step 4: Training
    - ▣ Step 5: Testing
- } The first 4 steps are not necessarily done in a strict linear progression

## Overview

$$D = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, N\} \quad \text{where} \quad f(\mathbf{x}_i) = y_i$$



## Step 1: Formulate the problem

- What quantity are you predicting?
  - ▣ real-valued, categorical, structure?
  - ▣ Changing over time?
  - ▣ Classification
    - Binary classification? Multi-class classification?
    - Singly-labeled? Multi-labeled?
    - For multi-labeled classification tasks, how correlated are the labels?
- What data do you have?
  - ▣ Where to get labeled data? (Amazon mechanical turk)
  - ▣ How much labeled data?
  - ▣ What is the quality of the labeled data?
  - ▣ Are the labels learnable given the data?
  - ▣ Is the distribution of labels in the data skewed/imbalanced?

## Guiding Principles

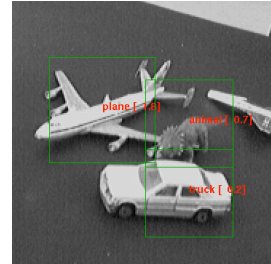
- Unsupervised learning as a surrogate for supervised learning...is a headache. Just get more data
- Reproducibility
- Think of how you would justify each decision you made
- Start simple and iterate

## Reducing multi-class to binary task

One-vs-All	<table border="1"> <tbody> <tr><td>x1</td><td>c1</td></tr> <tr><td>x2</td><td>c3</td></tr> <tr><td>x3</td><td>c1</td></tr> <tr><td>x4</td><td>c2</td></tr> </tbody> </table> <p style="text-align: center; color: blue;">original training data</p>	x1	c1	x2	c3	x3	c1	x4	c2	<table border="1"> <tbody> <tr><td>x1</td><td>1</td></tr> <tr><td>x2</td><td>-1</td></tr> <tr><td>x3</td><td>1</td></tr> <tr><td>x4</td><td>-1</td></tr> </tbody> </table> <p style="text-align: center; color: blue;">c1 vs. all</p>	x1	1	x2	-1	x3	1	x4	-1	<table border="1"> <tbody> <tr><td>x1</td><td>-1</td></tr> <tr><td>x2</td><td>-1</td></tr> <tr><td>x3</td><td>-1</td></tr> <tr><td>x4</td><td>1</td></tr> </tbody> </table> <p style="text-align: center; color: blue;">c2 vs. all</p>	x1	-1	x2	-1	x3	-1	x4	1	<table border="1"> <tbody> <tr><td>x1</td><td>-1</td></tr> <tr><td>x2</td><td>1</td></tr> <tr><td>x3</td><td>-1</td></tr> <tr><td>x4</td><td>-1</td></tr> </tbody> </table> <p style="text-align: center; color: blue;">c3 vs. all</p>	x1	-1	x2	1	x3	-1	x4	-1
x1	c1																																			
x2	c3																																			
x3	c1																																			
x4	c2																																			
x1	1																																			
x2	-1																																			
x3	1																																			
x4	-1																																			
x1	-1																																			
x2	-1																																			
x3	-1																																			
x4	1																																			
x1	-1																																			
x2	1																																			
x3	-1																																			
x4	-1																																			
One-vs-One	<table border="1"> <tbody> <tr><td>x1</td><td>c1</td></tr> <tr><td>x2</td><td>c3</td></tr> <tr><td>x3</td><td>c1</td></tr> <tr><td>x4</td><td>c2</td></tr> </tbody> </table> <p style="text-align: center; color: blue;">original training data</p>	x1	c1	x2	c3	x3	c1	x4	c2	<table border="1"> <tbody> <tr><td>x1</td><td>1</td></tr> <tr><td>x3</td><td>1</td></tr> <tr><td>x4</td><td>-1</td></tr> </tbody> </table> <p style="text-align: center; color: blue;">c1 vs. c2</p>	x1	1	x3	1	x4	-1	<table border="1"> <tbody> <tr><td>x1</td><td>1</td></tr> <tr><td>x2</td><td>-1</td></tr> <tr><td>x3</td><td>1</td></tr> </tbody> </table> <p style="text-align: center; color: blue;">c1 vs. c3</p>	x1	1	x2	-1	x3	1	<table border="1"> <tbody> <tr><td>x2</td><td>-1</td></tr> <tr><td>x4</td><td>1</td></tr> </tbody> </table> <p style="text-align: center; color: blue;">c2 vs. c3</p>	x2	-1	x4	1								
x1	c1																																			
x2	c3																																			
x3	c1																																			
x4	c2																																			
x1	1																																			
x3	1																																			
x4	-1																																			
x1	1																																			
x2	-1																																			
x3	1																																			
x2	-1																																			
x4	1																																			

## Multi-label Classification

- Each example can be labeled with multiple labels
  - ▣ Don't confuse this with multi-class classification!
  - ▣ Common for document classification or object recognition
- One-vs-all
- One classifier for every possible combination of labels
  - ▣ Combinatorial explosion
  - ▣ Limited training data



## Step 2: Exploratory Data Analysis

- Look at the data. It's surprising how often we forget to actually do this!
- **Exploratory Data Analysis** (EDA) is a statistical mindset
  - ▣ Box plots, histograms, scatter plots, mean, mode, deviations
  - ▣ Can guide the modeling process by
    - give you insight into the data
    - help (in)validate your assumptions
    - detect outliers

## Step 3: Feature Selection

- What features should I use?
  - ▣ Dimensionality reduction if exist time/space constraints
  - ▣ Reduce noise in the data (irrelevant or redundant features)
- Dimensionality reduction
  - ▣ Principal component analysis (PCA)
  - ▣ Singular value decomposition (SVD)
  - ▣ Canonical correlation analysis (CCA)
- Regularization
  - ▣ Use every feature but penalize classifiers that are overly complex

$$\text{Error}(w) = \sum_{i=1}^N (y_i - h_w(x_i)) + \lambda \|w\|^2$$

encourages sparse weight vectors

## Other tricks

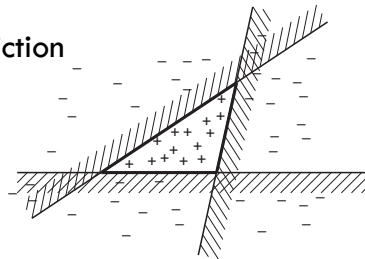
- Scale input features
- Transform features
  - ▣ e.g., take log
- Higher-order features
  - ▣ e.g., product of features
- Again, EDA can help guide this process

## Step 4: Training

- Pick your classifier
  - ▣ Decision tree, perceptron, neural network, SVM, linear regression, logistic regression, random forests, ensembles, Gaussian process regression, hidden Markov models, conditional random field, Bayesian networks,...
  - ▣ Bagging or Boosting
- Your choice is informed by all of the previous steps
- Often there are parameters that must be tuned...

## Ensembles of Classifiers

- An ensemble of classifiers – A group of classifiers whose predictions are combined to produce one final prediction
- Benefits
  - ▣ Harder to make a wrong prediction
  - ▣ More expressive hypothesis



## Boosting

- Learn a series of **weak classifiers**
- Weight each weak classifier to create a final strong classifier
- Often the weight for each classifier is proportional to its accuracy
- **AdaBoost** (Freund and Schapire 1995)

## Bagging

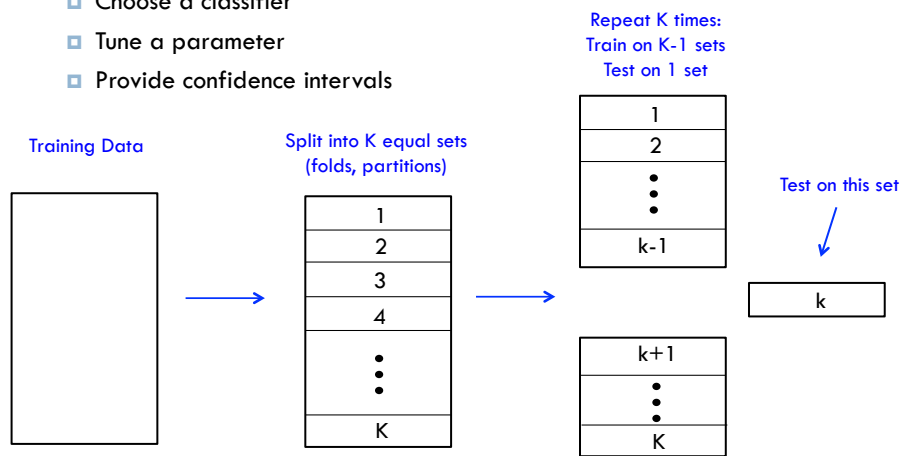
- Short for “Bootstrap aggregating”
- Given training set  $D$ 
  - Generate  $M$  new training sets  $D_i$  where  $|D_i| < |D|$  by sampling from  $D$  with replacement
  - This is a statistical technique known as **bootstrapping**
  - Train a classifier on each of the  $M$  new training sets
  - Combine output of  $M$  classifiers using averaging or voting
- Random Forests (Breimen, 2001)
  - Bagged decision trees



## Cross Validation

### □ K-fold cross validation

- Choose a classifier
- Tune a parameter
- Provide confidence intervals



## Step 5: Testing

- We have a final hypothesis
- We now use our hypothesis to predict on new (unseen) examples from the test set.
  - There's no going back and tweaking the classifier based on its test set performance!
- Where do these new unseen examples come from?
  - External source
  - Set aside from training data

## Binary Classification: Measures of Performance

- Let  $D_{\text{TEST}} = \{(x_i, y_i) \mid i=1 \dots N\}$  be our test set and  $\{h_i\}$  be the set of predicted values
- The contingency table is given by:

	$y = 1$	$y = 0$
$h = 1$	TP	FP
$h = 0$	FN	TN

- TP is the number of *true positives*
- FP is the number of *false positives*
- FN is the number of *false negatives*
- TN is the number of *true negatives*

## Binary Classification: Measures of Performance

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F_1\text{-score} = 2 \cdot \frac{\text{Prec} \cdot \text{Recall}}{\text{Prec} + \text{Recall}}$$

	$y = 1$	$y = 0$
$h = 1$	TP	FP
$h = 0$	FN	TN

Contingency Table

## Binary Classification: Measures of Performance

$$\text{Accuracy} = \frac{7 + 8}{7 + 8 + 2 + 3} = \frac{15}{20} = .75$$

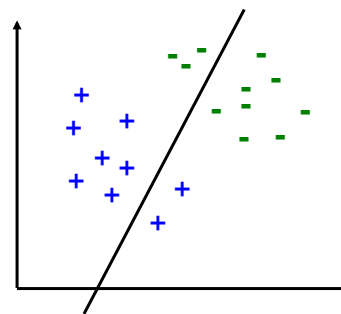
$$\text{Precision} = \frac{7}{7 + 3} = .70$$

$$\text{Recall} = \frac{7}{7 + 2} = .78$$

$$F_1\text{-score} = 2 \left( \frac{.70 \cdot .78}{.70 + .78} \right) = 2 \left( \frac{.546}{1.48} \right) = .74$$

	y = 1	y = 0
h = 1	7	3
h = 0	2	8

Contingency Table



## Multi-class Classification: Measures of performance

- Evaluate each label separately using a “one-vs-all” approach

- ▣ Macro-averaging

- Compute the measure (precision, recall,  $F_1$ ) for each class
    - Average across all C classes
    - Gives equal weight to all classes

- ▣ Micro-averaging

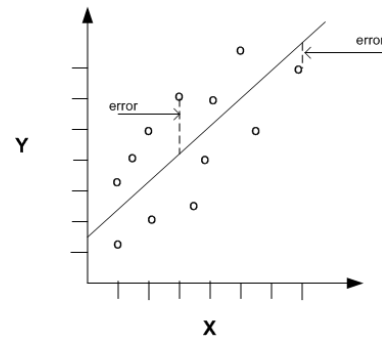
- Pool the TP, FP, FN, TN for all C classes
    - Compute the measure (precision, recall, F1)
    - Weighted towards performance of most likely class

	$y_c = 1$	$y_c = 0$
$h_c = 1$	$TP_c$	$FP_c$
$h_c = 0$	$FN_c$	$TN_c$

Contingency Table

## Regression: Measures of performance

- Mean-squared error
- Root mean-squared error
- Mean absolute error
- Mean absolute percentage
- ...



## Summary

- Overview
  - Step 1: Formulate the problem
  - Step 2: Explore the data
  - Step 3: Feature Selection
  - Step 4: Training
  - Step 5: Testing