Lecture 7: Supervised Learning Pt. 1

Linear Classifiers and Cross Validation INFO 1998: Introduction to Machine Learning



Announcements

- Assignment 6: Due tonight at 11:59pm
- Assignment 7: Due next Monday at 11:59pm
- **Mid-Semester Check-Ins:** Due today! Please get them done.
- **Next Lecture:** More Supervised Learning! (Decision Trees & Logistic Regression)



Agenda

- 1. Linear Classifiers
 - Linear Perceptron
 - Support Vector Machines (SVMs)
- 2. Cross Validation (K-Fold)



Linear Classifiers



Linear Classifiers

A linear classifier is a hyper plane that is used to classify our data points

A hyperplane is our **decision boundary** and our goal is to find the best hyper plane for our data.





Source

Linearly Separable

In this example, we cannot partition our dataset into yellow and purple with a linear decision boundary. This means that our data is <u>not</u> **linearly separable.**

Outliers are frequently the reason a data set is not linearly separable.





Perceptron Learning Algorithm

Goal: find a normal vector w that perfectly classifies all the points in our data set Algorithm:

Initialize classifier as some random hyperplane While there exists a misclassified point x: Adjust classifier slightly so that it classifies x correctly (or, is a little closer to classifying x correctly) End While

"Use your mistakes as your stepping stones"



History of the Linear Perceptron

Frank Rosenblatt was first to implement perceptron! \rightarrow Cornell professor and alum PHD '56 🐻

Gave him the title of 'Father of Deep Learning'

Deep Learning \rightarrow Neural Networks a.k.a. Multilayer Perceptrons





Limitations of Perceptron

The training algorithm will never terminate if your training dataset is not linearly separable 😔

> Is a great model to understand the intuition behind the training of a linear classifier: iteratively improve classifier by using misclassified points



Perceptron

Algorithm

```
Initialize \vec{w} = \vec{0}

while TRUE do

m = 0

for (x_i, y_i) \in D do

if y_i(\vec{w}^T \cdot \vec{x_i}) \leq 0 then

\vec{w} \leftarrow \vec{w} + y\vec{x}

m \leftarrow m + 1

end if

end for

if m = 0 then

break

end if

end while
```

// Initialize \vec{w} . $\vec{w} = \vec{0}$ misclassifies everything.

// Keep looping

- // Count the number of misclassifications, m
- // Loop over each (data, label) pair in the dataset,
- // If the pair $(\vec{x_i}, y_i)$ is misclassified
- // Update the weight vector \vec{w}

// Counter the number of misclassification

// If the most recent \vec{w} gave 0 misclassifications // Break out of the while-loop

// Otherwise, keep looping!



Lecture 7: Supervised Learning Pt. 1

Linear Classifiers and Cross Validation INFO 1998: Introduction to Machine Learning





Attendance!

Support Vector Machines



Classify (+) and (-)





Which Hyperplane?





Optimal Hyperplane





Support Vector Machine

Memory efficient Effective in a higher dimensions Slow calculation time



Maximal Margin Classifier

- We want to find a **separating** hyperplane
- Once we find candidates for the hyperplane, we try to maximize the margin, the normal distance from borderline points
 - Only Support Vectors matter





What if...





Which Decision Boundary is better?



Margins

Use cost function to penalize misclassified points

Choice of cost function makes margin "hard" vs. "soft"



Penalty of error: distance to hyperplane multiplied by error cost C.



Hard Margins

- High penalty value
- The hyperplane can be dictated by a single outlier





Soft Margins

- Used in non-linearly separable datasets
- Allow for misclassification
- Can account for "dirty" boundaries





Misclassification Penalty C





Kernels

- You cannot linearly divide the 2 classes on the *xy* plane at right
- Introduce new feature, z = x² + y²
 (radial kernel)
- Map 2 dimensional data onto 3 dimensional data. Now a hyperplane is easy to find.





Kernels





SVM has MANY Hyperparameters





Finding the Best Hyper Parameters



Curse of Dimensionality

Our search space for the optimal hyper-parameters increases **exponentially** as the number of hyper parameters we are considering increases



Overview

Perceptron	SVM
 A very simple model Will perform poorly if data is not linearly separable 	 More complex model because we have to choose the "penalty cost" associated with misclassifications Can transform feature space by choosing a Kernel







Cross Validation





Often used in practice with *k*=5 or *k*=10.

Create equally sized *k* partitions, or **folds**, of training data

For each fold:

- Treat the *k-1* other folds as training data.
- Test on the chosen fold.

The average of these errors is the validation error



Dataset

Suppose K = 5, 5-Fold CV









Calculate MSE = mse1





Calculate MSE = mse2





Calculate MSE = mse3



And so on





MSE = Avg(mse1...5)



Matters less how we divide up

Selection bias not present



Leave-1-Out Cross Validation

For each sample:

- Treat all other data as training data.
- Test on that one sample

The average of these errors is the validation error

Pro: Better on small datasets

Pro: More realistic (trained on most of the data)

Con: Takes longer to run



Coming Up

- Assignment 6: Due tonight at 11:59pm
- Assignment 7: Due next Monday at 11:59pm
- Mid-Semester Check-Ins: Due today! Please get them done ASAP
- **Next Lecture:** More Supervised Learning! (Decision Trees & Logistic Regression)

