

Machine Learning

Todd Davies

December 13, 2014

Introduction

Machine learning is concerned with creating mathematical "data structures" that allow a computer to exhibit behaviour that would normally require a human. Typical applications might be spam filtering, speech recognition, medical diagnosis, or weather prediction. The data structures we use (known as "models") come in various forms, e.g. trees, graphs, algebraic equations, probability distributions. The emphasis is on constructing these models automatically from data—for example making a weather predictor from a datafile of historical weather patterns. This course will introduce you to the concepts behind various Machine Learning techniques, including how they work, and use existing software packages to illustrate how they are used on data. The course has a fairly mathematical content although it is intended to be self-contained.

Aims

To introduce methods for extracting rules or learning from data, and provide the necessary mathematical background to enable students to understand how the methods work and how to get the best performance from them. This course covers basics of both supervised and unsupervised learning paradigms and is pitched towards any student with a mathematical or scientific background who is interested in adaptive techniques for learning from data as well as data analysis and modelling.

Additional reading

Introduction to machine learning (2nd edition)	Alpaydin, Ethem	2010
--	-----------------	------

Contents

- 1 Machine Learning
- 2 Nearest Neighbour Classifier
 - 2.1 Finding the distance between two n -dimensional points
 - 2.2 Computing the nearest neighbour
 - 2.3 Multiple nearest neighbours (K-NN)
 - 2.3.1 Overfitting

1 Machine Learning

Machine Learning is the creation of self-configuring data structures that allow a computer to do things that would be classed as ‘intelligent’ if a human did it.

Machine learning has been around in it’s infancy since the 40’s, where reasoning and logic were first studied by Claude Shannon and Kurt Godel. Steady progress was made with lots of funding through until the 70’s, where people then realised Artificial Intelligence was very hard, causing funding to dry up, and the term ‘AI Winter’ to be coined.

In the 80’s people started to look at biologically inspired algorithms such as neural networks and genetic algorithms, and there is more investment. This leads to the field of AI diverging into many other fields such as Computer Vision, NLP, Machine Learning etc. In the 00’s, ML begins to overlap with statistics, and the first *useful* applications emerge.

2 Nearest Neighbour Classifier

The Nearest Neighbour (NN) classifier is a simple implementation of a machine learning algorithm. We can give it a set of training data with each point labelled by a class, and then give it another datapoint, and that datapoint will be classified according to the data.

The premise is that you plot all the points of the training data on a graph, and when you want to classify another

datapoint, you simply plot it on the existing graph, and classify it by the classes of it's nearest neighbour(s).

One nice aspect of the NN classifier is that you can use it with as many features as you like with little extra effort. A standard implementation of the NN classifier might plot graphs in two dimensions, and thus will only be able to classify data with two features. However, if you want to include more (or less) features, you need to adjust the number of dimensions on your graph to match the number of features you're including. Then when you find the nearest neighbours on the graph, you find them in n -dimensional space rather than 2-dimensional space, where n is the number of dimensions on the graph.

2.1 Finding the distance between two n -dimensional points

In order to calculate which points in the training dataset are the nearest neighbours to the new data point, we can calculate the *Euclidean Distance* between the two points. In n dimensions, this how:

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \cdots + (p_{n-1} - q_{n-1})^2 + (p_n - q_n)^2}$$

This is actually a very simple computation. If you're given two arrays (containing the coordinates of your points) of equal lengths, then you can do something similar to Listing 1.

```
1 def euclideanDistance(c1: List[Double], c2: List[Double]) :  
    Double = {  
2   val delta: Set[Double] = for {(a,b) <- c1 zip c2} yield  
    math.pow(a-b, 2)  
3   Double = math.sqrt(delta.sum)  
4 }
```

Listing 1: Scala Euclidean Distance

The reason why using KNN on n-dimensions is that you can classify different types of data. A 2d KNN classifier could classify on two variables, such as weight and height, while a 256d classifier could classify a 16x16 monochrome image where each pixel is represented by a dimension.

2.2 Computing the nearest neighbour

Computing the nearest neighbour to a point is simply as easy as finding the point in the training data that has the smallest euclidean distance to it, as shown in Listing 3.

```
1 def nearestNeighbour(input: Point, existing: Set[Point]) :  
    Point = {  
2   existing.min((p: Point) => p.euclideanDistance(input))  
3 }
```

Listing 2: Scala Nearest Neighbour

Note:

Assume that there is a custom implementation of **Point** that exposes the method to find the euclidean distance between it and another point, similar to that in Listing 1

2.3 Multiple nearest neighbours (K-NN)

It's possible that more accurate classifications of points could be obtained by taking into account multiple close neighbours rather than one nearest one. In order to do this, you need to find the number of occurrences of a specific class in the top K nearest neighbours:

```
1 def kNearestNeighbour(k: Int, input: Point, existing: Set[
    Point]): Class = {
2   val topK: List[Class] = existing.toList.sorted(
3     (p: Point) => p.euclideanDistance(input)
4   ).take(k).map(_.class)
5   topK.groupBy(identity).maxBy(_._2.size)._1
6 }
```

Listing 3: Scala Nearest Neighbour

One must be careful when choosing a value of K for the classifier. As K increases proportional to the size of the dataset, then the number of incorrect classifications will increase.

Figure 1 shows two identical datasets that a KNN classification is being applied to. When $K = 3$, the correct classification of a square is made, and when $K = 5$, the wrong classification is made, simply because there are more circles than squares in the dataset.

With a nearest neighbour classifier, there is always one or more *boundaries* that defines which class a new item will be placed in, as shown in figure 2. Be aware that the **Decision Boundary** isn't always contiguous or a straight line.

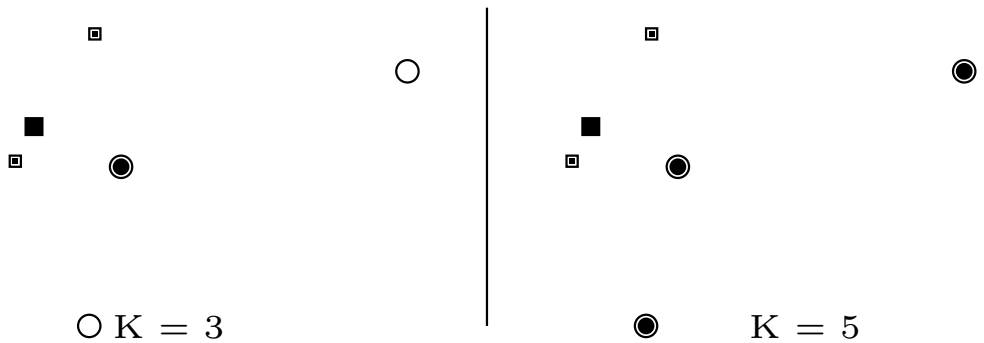


Figure 1: How K affects classification accuracy. The item being classified is the filled square and nearest neighbours are filled, and other elements are left unfilled. When $K=3$, it would be classified as square, when $K=5$, it'd be classified as a circle.

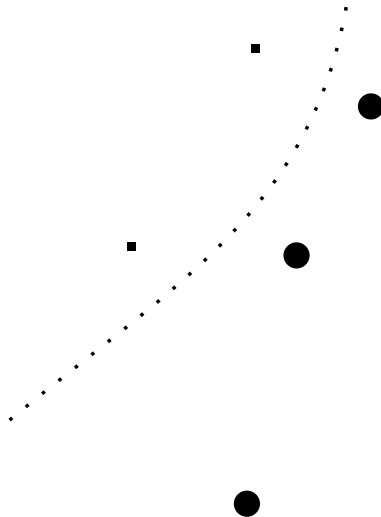


Figure 2: The dotted line represents the decision boundary for a KNN classifier

2.3.1 Overfitting

Overfitting is a problem with all ML algorithms, and occurs when an algorithm is trained on a small proportion of the dataset, that isn't representative of the whole data. It can also be when the algorithm is trained in such a way that it describes the noise or random error in the training set.

Avoiding overfitting is hard, but generally it involves choosing the right algorithm, the right parameters and the right training data.