

Unsupervised methods: no "training data", i.e. labelled examples from particular classes

- ex:
- k-means (entire well-separated, compact clusters)
 - DBSCAN: look for connected high density regions
 - Spectral clustering: approximate the best graph cut

Supervised learning: Some (perhaps a lot) of labelled instances are available.

- Problems:
- regress a function $f(x)$ given samples $\{(x_i, f(x_i))\}_{i=1}^n$.
 - build a classifier $C(x) \in \{1, \dots, K\}$ given labeled examples $\{(x_i, C(x_i))\}$,

• Of course, classification is a special case of regression.

• Many ways to do this; classical topics in statistics.

• We will focus on classification, though similar ideas will apply to regression.

• We will focus on three methods:

- 1) Nearest neighbor classification: simple, inelegant, suffers obvious limitations
- 2) Support vector machines: more complicated, elegantly handle nonlinear data through the "kernel trick"
- 3) Neural networks: extremely complicated, poorly understood theoretically, amazing results in some settings

Remark: Battle between (kernel) SVM and neural networks is ongoing...


Neural networks have the decisive upper hand at the moment.

Let $\{(x_i, y_i)\}_n$ be labelled data:
 - $x_i \in \mathbb{R}^D$
 - $y_i \in \{1, \dots, K\}$

The nearest neighbor classifier sets for a new point $x \in \mathbb{R}^D$

The label $C(x) = \text{mode} \{x_i^{NN}\}_{i=1}^{K_{NN}}$, where
 1) K_{NN} is some integer ≥ 1
 2) $\{x_i^{NN}\}_{i=1}^{K_{NN}}$ are the K_{NN} nearest neighbors of x , i.e. $\{x_i^{NN}\}_{i=1}^{K_{NN}}$ are the K_{NN} minimum of $\|x - x_i\|_2$.

This is simple in an important way: nothing to do before assigning a label - "lazy" classifier

ex:  If $K_{NN} = 4 \Rightarrow C(x) = 2$
 $K_{NN} = 8 \Rightarrow C(x) = 1$

Choosing K_{NN} is important, but it's the only parameter.

It can, moreover, be learned through a general process in supervised learning called cross validation.

Indeed, one can take the training data $\{x_i\}_{i=1}^n$ and split it into a training and validation set. (3)

For k -NN classification, the whole training set can be used for cross-validation. Choose k that maximizes accuracy on the training set, say by testing x_i on $\{x_j\}_{j \neq i}$ for all $i=1, \dots, n$ (leave one out cross-validation).

Remark: While simple (and often effective), k -NN classification suffers from at least two serious problems.

Problem 1: Computation. Suppose we fix $x \in \mathbb{R}^D$. Who among the n training points are the k_{nn} nearest neighbors? Well, for each x_i , need to compute the L^2 distance in \mathbb{R}^D . Need to do this k_{nn} times $\rightarrow O(n \cdot D \cdot k_{nn})$.

Indexing structures (kd-trees, case trees) allow us to speed up this.

If you want $k_{nn} \rightarrow \sqrt{n}$, it may be faster to just sort the pairwise distances $O(D \log(n))$.

Problem 2: Curse of Dimensionality. KNN-classification (4)

tries to use pairwise distances to make decisions. So, there should be some information in the pairwise distances for this to work.

• But, in high dimensions, pairwise distances are uninformative!

• Intuitively, for the unit ball (volume = 1) or $[0,1]^D$, most of the volume concentrates near the boundary.

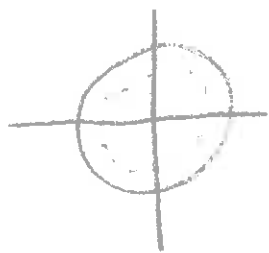
• Exercise: Let $B_r^D = \{x \in \mathbb{R}^D \mid \|x\|_2 \leq r\}$ be the radius r sphere in \mathbb{R}^D . Then the surface area of B_r^D is

$$SA(B_r^D) = \frac{D \pi^{D/2}}{\Gamma(D/2 + 1)} r^{D-1}, \quad \Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$$

• So, to go from surface area to volume, integrate over r :

$$\begin{aligned} \text{Vol}(B_r^D) &= \int_0^r \frac{D \pi^{D/2}}{\Gamma(D/2 + 1)} \rho^{D-1} d\rho \\ &= \frac{\pi^{D/2}}{\Gamma(D/2 + 1)} \int_0^r D \rho^{D-1} d\rho \end{aligned}$$

So, how much volume is between the sphere of radius 1 (unit sphere) and the sphere of radius $(1-\epsilon)$? in \mathbb{R}^2 : (5)



$$\text{vol}(B_1^2) - \text{vol}(B_{1-\epsilon}^2)$$

$$= \pi \cdot 1^2 - \pi (1-\epsilon)^2$$

$$= \pi [1 - (1-\epsilon)^2]$$

(P) vol = area
SA = perimeter
in \mathbb{R}^2

As a proportion of the volume of B_1^2 :

$$\frac{\text{vol}(B_1^2) - \text{vol}(B_{1-\epsilon}^2)}{\text{vol}(B_1^2)}$$

$$= 1 - (1-\epsilon)^2$$

$$= 2\epsilon - \epsilon^2$$

How about in \mathbb{R}^D ? $\text{vol}(B_1^D) - \text{vol}(B_{1-\epsilon}^D)$

$$= \frac{\pi^D}{\Gamma(\frac{D}{2}+1)} \int_{1-\epsilon}^1 \rho^D d\rho$$

$$= \frac{\pi^D}{\Gamma(\frac{D}{2}+1)} [1 - (1-\epsilon)^D]$$

(6)

As a proportion of total volume:

$$\frac{\text{vol}(B_{\epsilon}^D) - \text{vol}(B_{1-\epsilon}^D)}{\text{vol}(B_{\epsilon}^D)}$$
$$= 1 - (1-\epsilon)^D$$

• So, if $D=100$, $\epsilon = \frac{1}{2}$, then $1 - \left(\frac{1}{2}\right)^{100} \approx 1$

• This concentration of measure of phenomenon makes KNN good when D is large, without a lot of samples.