

NON-PARAMETRIC CLASSIFIERS

K-NEAREST NEIGHBORS

- > For each test point, find the K nearest samples in the training data: a distance metric such as the 2-norm $\|x_{\text{test}} - x_i\|_2^2$, $i=1, \dots, N$ may be used.
- > Classify the point according to the majority of their class labels.

Classification Rule (Naive)

① $\vec{d} = \|x - x_i\|$ $i=1, \dots, N$

② Assume the target function (that assigns labels) is

$$g: \mathbb{R}^n \rightarrow \{1, \dots, C\}$$

Find position of smallest elements of \vec{d} : $\{d_j\}, j \in T$, where

$T \subseteq \{1, \dots, N\}$ is the set of the positions of the smallest K entries of \vec{d} .

③ Define set $S \triangleq \{(x_{T(1)}, g(x_{T(1)})), \dots, (x_{T(K)}, g(x_{T(K)}))\}$ where $T(i)$ is the

i^{th} element of set T . The class of x_{test} is

$$g(x_{\text{test}}) = \text{mode}\{g(x_{T(1)}), \dots, g(x_{T(K)})\}$$

Rule of thumb: $k < \sqrt{N}$, N being the no. of training samples

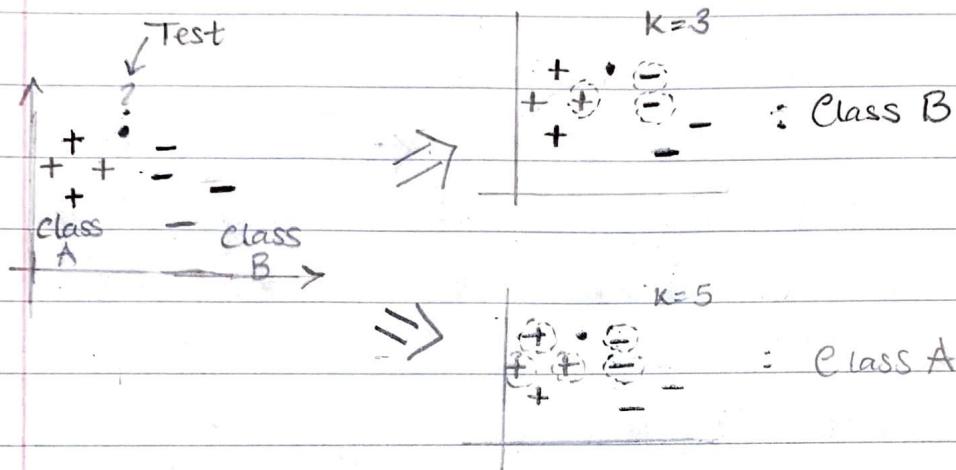


Figure above illustrates how the selection of k may affect classifier

- ① In general, the KNN scheme is intuitive and simple to implement. But it may become very expensive for large datasets. Notwithstanding it is easy to parallelize.

① PARZEN WINDOW (PROBABILITY DENSITY ESTIMATION)

In this approach, we split our domain of interest (n -dimensional space) into hypercubes (bins or intervals). For each bin, we 'count' or measure how many samples (per total samples counted) fall into that bin; and divide by the volume of the hypercube.

This approach helps build a probability density estimate over the domain:

$$P_{\text{est}}(x) = \frac{K_N/N}{\Delta x} \quad K_N: \text{number of training samples in}$$

the interval Δx

N : total number of samples counted

K_N/N : frequency

Δx : volume, size of interval

In particular, we can write this probability estimate using a windowing function

$$P_N(x) = \frac{1}{N} \sum_{i=1}^N \frac{1}{V_N} \varphi\left(\frac{x-x_i}{h_N}\right)$$

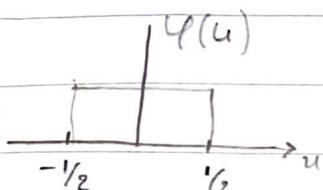
Window function:

measures if or how much
 x_i is in the bin with
width h_N centered at x

h_N : width of bin

$V_N = h_N^d$: volume of bin

$$\varphi(u) = \begin{cases} 1 & |u_i| < \frac{1}{2} \quad i=1, \dots, n \\ 0 & \text{otherwise} \end{cases}$$



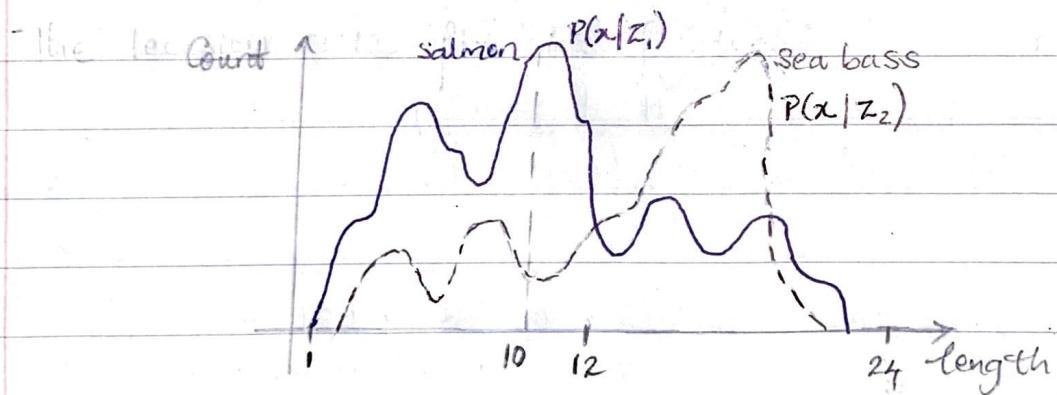
Note that we may also use a normal distribution window

$$\varphi(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2}$$

and define $h_N = h_1/\sqrt{N}$ where h_1 becomes the parameter at our disposal. This definition is to ensure that

$$\lim_{n \rightarrow \infty} V_n = 0$$

the estimate converges to the true probability density function



Decision Rule: $\operatorname{argmax}_{i \in \{1, \dots, c\}} p(x|w_i)$

PARAMETRIC CLASSIFIERS

We would like to have parametric methods for classification that summarize the data with a fixed number of parameters that does not scale with the number of training data available. We focus on what we call discriminant functions.

BINARY CLASSIFICATION

Given training data $\{x_i, z_i\}$ for $i=1, \dots, N$, $x_i \in \mathbb{R}^n$, $z_i \in \{-1, 1\}$ we want a discriminant function

$$g(x_i) \begin{cases} > 0 & \text{if } z_i = 1 \\ < 0 & \text{if } z_i = -1 \end{cases}$$

i.e. $z_i g(x_i) > 0$ implies correct classification

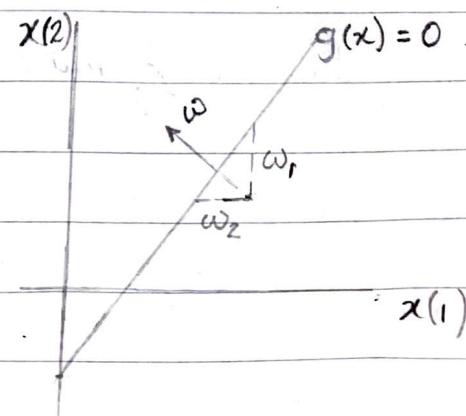
LINEAR CLASSIFIER

If we assume that the classes are linearly separable then we can write

$$g(\hat{x}) = \hat{\omega}^\top \hat{x} + \hat{w}_0$$

In 2D, the contour $g(x) = 0$ is a line

$$\text{since } x(1) = -\frac{\omega_1}{\omega_2}, x(2) = -\frac{b}{\omega_2}$$

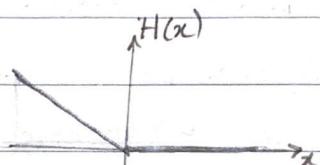


In n -dimensional space, it is a hyperplane.

We can write $\omega = \begin{bmatrix} \hat{\omega} \\ \omega_0 \end{bmatrix}$, $x = \begin{bmatrix} \hat{x} \\ 1 \end{bmatrix}$, so that $g(x) = \omega^T x$

Recall that we want $z_i g(x_i) > 0$, so we are only interested in cases where

$$\underbrace{\max(0, -z_i g(x_i))}_{\text{Perception loss function: generates misclassification in one term}} > 0$$



for when $z_i g(x_i) > 0$: $\max(0, -z_i g(x_i)) = 0$

only when $z_i g(x_i) < 0$: $\max(0, -z_i g(x_i)) > 0$
ie misclassification

PERCEPTRON ALGORITHM

We want to minimize misclassifications. This is equivalent to minimizing

$$J(\omega) = \sum_{i=1}^N \max(0, -z_i \omega^T x_i)$$

which we can re-write as

$$J(\omega) = - \sum_{i \in M} z_i \omega^T x_i$$

where $M \subset \{1, \dots, N\}$ is the set of misclassified samples

① To find ω that minimizes the cost function J , i.e
$$\arg \min_{\omega} J(\omega)$$

we use gradient descent.

Algorithm

for $K = 1, 2, \dots$

 update $M_K = \{i \mid z_i \omega^T x_i < 0\}$ for $\omega = \omega^{(K)}$

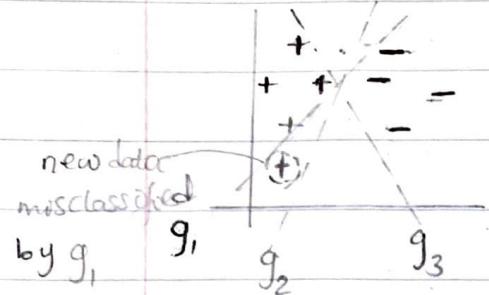
Gradient Descent

$$\omega^{(K+1)} \leftarrow \omega^{(K)} + \alpha_K \sum_{i \in M_K} z_i x_i$$

End for

This algorithm typically converges for linearly separable data.
Note that, in fact, J is convex.

SUPPORT VECTOR MACHINE



But what is the best hyperplane $g(x, \omega)$ that separates the data: if the data is linearly separable, there will likely be several (that are optimal based on the perceptron algorithm), but it may quickly mis-classify on unseen (test) data.

So we want an hyperplane (characterized with the orthogonal vector ω) that will maintain a margin between the closest element of a class and the dividing hyperplane.

In fact, $g(x, \omega)$ should be such that it is equidistant from the closest vector (in terms of projection onto ω) in class

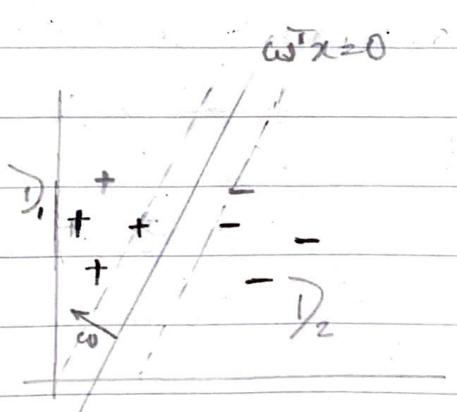
$$\mathcal{D}_1 \triangleq \{i \mid z_i = 1\} \text{ and the closest vector in class } \mathcal{D}_1 \triangleq \{i \mid z_i = 1\}$$

So we can think of 3 parallel hyperplanes

$$\hat{\omega}^T x + \omega_0 = d \text{ Supporting } \mathcal{D}_1$$

$$\hat{\omega}^T x + \omega_0 = 0 \text{ dividing hyperplane}$$

$$\hat{\omega}^T x + \omega_0 = -d \text{ Supporting } \mathcal{D}_2$$



We can normalize by dividing by d to get

$$\hat{\omega}^T x = 1 \text{ (for } \mathcal{D}_1\text{)} \quad \hat{\omega}^T x = -1 \text{ (for } \mathcal{D}_2\text{)}$$

$$\text{where } \omega = \frac{1}{d} \begin{bmatrix} \tilde{\omega} \\ \omega_0 \end{bmatrix}, \quad x = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

So,

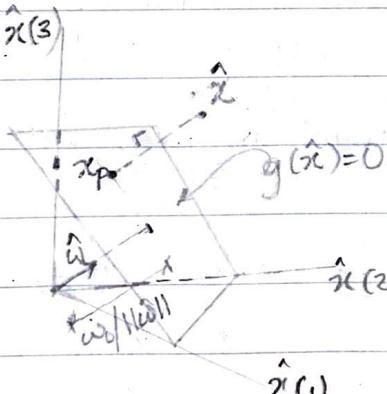
$$\tilde{\omega}^T x_i \geq 1 \quad \text{if } z_i = 1$$

$$\tilde{\omega}^T x_i \leq 1 \quad \text{if } z_i = -1$$

which can be written as

$$z_i \tilde{\omega}^T x_i \geq 1 \quad \forall i$$

Now for any $\hat{x} \in \mathbb{R}^n$, we can write



$$\hat{x} = x_p + r \frac{\hat{\omega}}{\|\hat{\omega}\|}$$

$$\text{since } g(x_p) = 0$$

$$g(\hat{x}) = \hat{\omega}^T \hat{x} + \hat{\omega}_0$$

$$= \omega^T x_p + r \|\omega\| + \omega_0$$

$$= r \|\hat{\omega}\|$$

$$r = \frac{g(\hat{x})}{\|\hat{\omega}\|}$$

$$\text{The distance of the hyperplane from } \hat{x}=0 = \frac{\omega_0}{\|\hat{\omega}\|}$$

It follows that the distance from any vector x_i on the plane

$$w^T x_i = 1 \text{ from any vector } x_i \text{ on the plane } w^T x_i = 1 \text{ vs } \frac{2}{\|w\|}$$

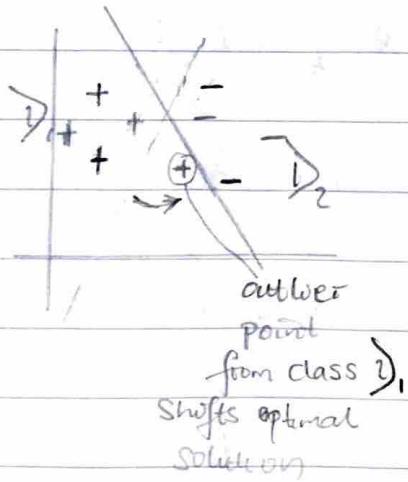
We want to maximize this distance while ensuring that the classification criterion ($x_i w^T x_i \geq 1 + \xi_i$) holds:

$$\min \|w\|^2$$

$$\text{s.t. } x_i w^T x_i \geq 1 + \xi_i \quad \forall i \in \{1, N\}$$

This is a quadratic problem with linear constraints, hence there is a unique minimum.

But the constraint ($x_i w^T x_i \geq 1$) may limit us from getting 'good' margins



To make the constraint softer, we introduce slack variables $\xi_i \geq 0$

if $0 < \xi_i \leq 1$ the point is between margin and the correct side of hyperplane.

This is called margin violation

if $\xi_i > 1$, point is misclassified.

The optimization problem now becomes

Equivalent expression

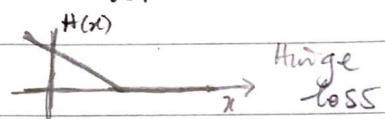
$$\min_{\omega, \xi} \|\omega\|_2^2 + C \sum_{i=1}^N \xi_i$$

$$z_i \omega^\top x_i \geq 1 - \xi_i \quad \forall i$$

$$\xi_i = \max(0, 1 - z_i \omega^\top x_i)$$

Hinge loss function

$$\min_{\omega} \|\omega\|_2^2 + C \sum_{i=1}^N \max(0, 1 - z_i \omega^\top x_i)$$



large C makes constraint hard to ignore

Small C allows constraint to be easily ignored \rightarrow large margin

In SVM's, it is common to first transform or map x to a higher dimension

$$y_i = \varphi(x_i)$$

generally nonlinear mapping

in expectation that in the higher dimensional space, a separating hyperplane will exist

MULTICLASS

1. In multiclass cases,

$$g_i(x) = \omega_i^T x$$

We assign x to z_i if $g_i(x) > g_j(x) \forall j \neq i$

We can write the multi-class SVM problem as (without slack variables)

$$\max_{\omega_1, \dots, \omega_c} \sum_k \|\hat{\omega}_k\|_2^2 + C \sum_{i=1}^N \xi_i$$

$$\text{st. } \hat{\omega}_{z_i}^T x_i - \hat{\omega}_k^T x_i \geq 1 - \xi_i$$

$$\forall k \in \{1, 2, \dots, c\}, k \neq z_i$$

$$\forall i$$

