

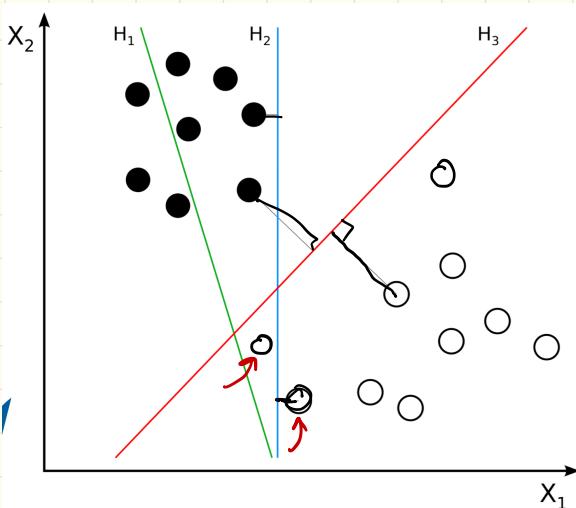
Topic 8

Support vector machines (SVM)

Plan:

1. Motivation for SVM
2. How to solve SVM optimisation problem?
3. Duality (a proof of the method of solving SVM)

Section I : Motivation and derivation.



Binary classific.

Previously
linear regression →

→ logistic func
→ classification

H₁ - wrong

H₂ - OK

H₃ - OK

Idea: maximise the total/average distance between points and a decision boundary.

1.2 Simple linear algebra/geometry

Assumptions: binary classification

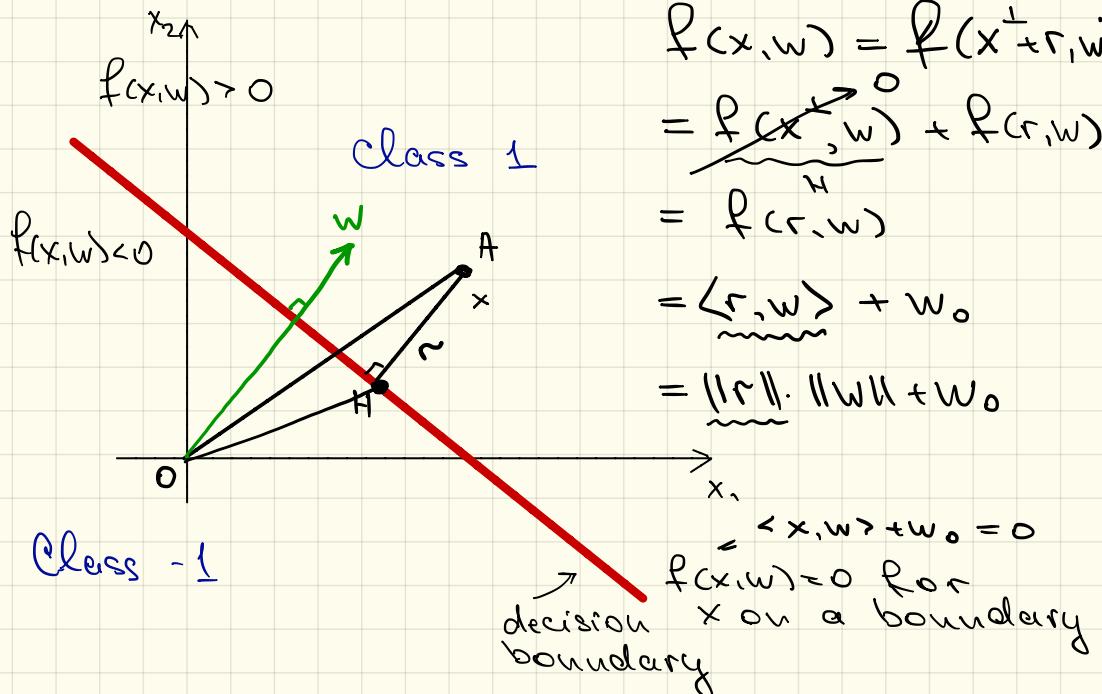
(based on linear regression), class labels $\{-1, 1\}$.

Start with $\{x^{(i)}, y_i\}_{i=1}^s$.

$$f(x, w) = \langle \tilde{x}, \tilde{w} \rangle \stackrel{\in \mathbb{R}^{d+1}}{\substack{\text{extended}}} = w_0 + \langle x, w \rangle$$

if $f(x, w) \geq 0 \Rightarrow$ label x with 1

if $f(x, w) < 0 \Rightarrow$ label x with -1



Goal express r in terms of $f(x, w)$.

$$AH \parallel \bar{w}$$

H is a projection of A on a line
 $\Rightarrow \overrightarrow{OH} = x^\perp$

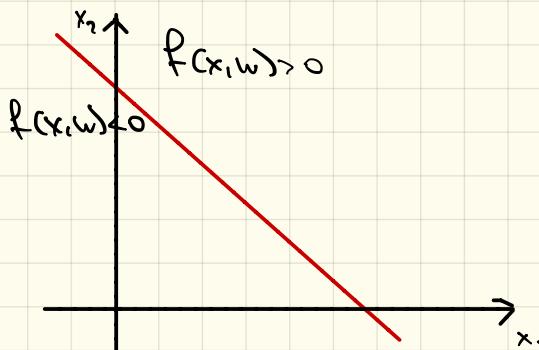
Any x : $x = x^\perp + r$
 $\hookrightarrow \in \text{line}$ $\rightarrow \perp \text{ to the line.}$

$$\|r\| = \frac{|f(x, w) - w_0|}{\|w\|}$$

max

$$\Rightarrow \boxed{\max \frac{f(x, w)}{\|w\|}}$$

What about classification?



$$\begin{cases} x^{(i)}, y_i \end{cases}_{i=1}^s$$

$$\{ -1, 1 \}$$

We want to have $f(x^{(i)}, w) > 0$ if

$y_i = 1$ and $f(x^{(i)}, w) < 0$ if

$y_i = -1$. $\Leftrightarrow y_i \cdot f(x^{(i)}, w) \geq 0$

" \Leftrightarrow " $y_i \cdot f(x^{(i)}, w) \geq 1$.

Conclusion: $\max \sum_{i=1}^s \frac{|f(x^{(i)}, w)|}{\|w\|} = |y_i|$

with $y_i \cdot f(x^{(i)}, w) \geq 1$ for $\forall i$.

Reformulate: $\min \|w\|$ with
the same constraints

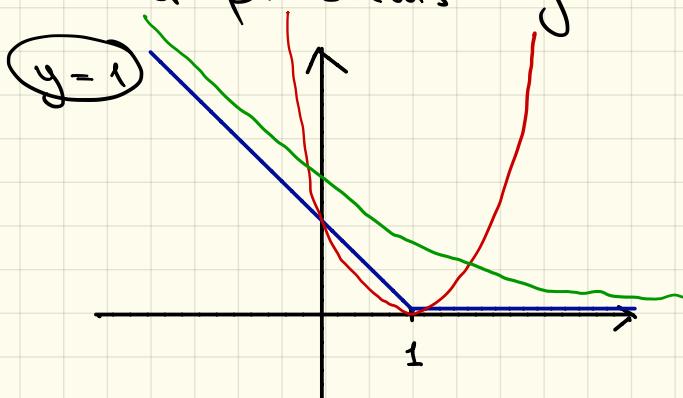
$$\min_w \left\{ \sum_{i=1}^s \max\{0, 1 - y_i f(x^{(i)}, w)\} + \frac{\alpha}{2} \|w\|^2 \right\}$$

Support vector machine optimisation problem.

$$f(x^{(i)}, w) = \langle \tilde{x}^{(i)}, \tilde{w} \rangle = w_0 + \langle x^{(i)}, w \rangle$$

$$H(z) = \max(0, 1 - yz) \text{ - Hinge-loss}$$

In our problems $y = \pm 1$



- MSE(z) = $(1-yz)^2$
- Logistic(z)
 $= \log(1+e^{-yz})$

2. Solving SVM optimisation

$$\min_w \left\{ \sum_{i=1}^s \max(0, 1 - y_i f(x^{(i)}, w)) + \frac{\alpha}{2} \|w\|^2 \right\}$$

$= L(w)$

Problems: ① $\max(0, 1 - yz)$ is non-differentiable.

② rewrite everything in terms of X, y, w

$$L(w) = \sum_{i=1}^s \max(0, 1_i - (Xw)_i) + \frac{\alpha}{2} \|w\|^2$$

$$1 = \underbrace{(1, 1, \dots, 1)}_s^\top$$

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ 0 \\ \vdots \\ y_s \end{pmatrix}$$

$$\max(0, z) = \max_{\lambda \in [0, 1]} \lambda z$$

$$L(w) = \sum_{i=1}^s \max_{\lambda_i \in [0, 1]} \left\{ \lambda_i \cdot (1_i - (Xw)_i) + \frac{\alpha}{2} \|w\|^2 \right\}$$

$$L(w) = \sum_{i=1}^s \max_{\lambda_i \in [0, 1]} \left\{ \lambda_i \cdot \underbrace{\left(\mathbb{1}_i - (\mathbf{Y} \mathbf{X}_w)_i \right)}_{\text{margin}} \right\} + \frac{\alpha}{2} \|w\|^2$$

$$= \boxed{\max_{\Lambda} \sum_{i=1}^s \left(\Lambda (\mathbb{1} - \mathbf{Y} \mathbf{X}_w)_i \right)_i + \frac{\alpha}{2} \|w\|^2}$$

$$\Lambda = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & 0 \\ 0 & \ddots & \lambda_s \end{pmatrix}$$

SVM problem:

$$\min_w \max_{\Lambda} \sum_{i=1}^s \left(\Lambda (\mathbb{1} - \mathbf{Y} \mathbf{X}_w)_i \right)_i + \frac{\alpha}{2} \|w\|^2$$

$$\hat{w} = \arg \min_w \max_{\Lambda} \sum_{i=1}^s \left(\Lambda (\mathbb{1} - \mathbf{Y} \mathbf{X}_w)_i \right)_i + \frac{\alpha}{2} \|w\|^2$$

Crucial idea: change $\min \leftrightarrow \max$
order !

Example: $f(x,y) = \sin(cx+y)$

$$\min_x \max_y f(x,y) = \min_x 1 = 1$$

$$\max_y \min_x f(x,y) = \max_y -1 = -1$$

we can't change min-max order here.

Assume we can change min-max in SVM optimisation problem.

$$\max_\lambda \min_w \underbrace{\sum_{i=1}^n \lambda (1 - y_i X_w)_+ + \frac{\alpha}{2} \|w\|^2}_{L(w, \lambda)}$$

Observation:

- $L(w, \lambda)$ is differentiable in w and λ .
- $L(w, \lambda)$ is convex in w , is concave in λ .

$L(w, \Lambda)$ is minimised in w 's when $\nabla_w L(w, \Lambda) = 0$.

$$L(w, \Lambda) = \sum_{i=1}^s \lambda_i (1 - y_i \langle \mathbf{x}^{(i)}, w \rangle) + \frac{\alpha}{2} \|w\|^2$$

$$\nabla_w L(w, \Lambda) = \sum_{i=1}^s -\lambda_i y_i \mathbf{x}^{(i)} + \alpha w = 0$$

$$w = \frac{1}{\alpha} \sum_{i=1}^s \lambda_i y_i \mathbf{x}^{(i)} = \frac{1}{\alpha} X^T Y \underbrace{\lambda}_{\text{---}}$$

$$\max_{\Lambda \in \mathbb{R}^{s \times s}} L(\underbrace{\frac{1}{\alpha} X^T Y \lambda}_{\text{---}}, \Lambda)$$

$$= \max_{\Lambda} \left\{ \sum_{i=1}^s \left(\Lambda \left(1 - \frac{1}{\alpha} \mathbf{Y} X \mathbf{X}^T Y \lambda \right) \right) + \frac{\alpha}{2} \frac{1}{\alpha^2} \|X^T Y \lambda\|^2 \right\}$$

$$= \max_{\Lambda} \left\{ \left\langle \lambda, \mathbf{1} \right\rangle - \frac{1}{2} \sum_{i=1}^s \left\langle \Lambda \mathbf{Y} \mathbf{X} \mathbf{X}^T Y \lambda, \mathbf{1} \right\rangle + \frac{1}{2\alpha} \|X^T Y \lambda\|^2 \right\}$$

$$= \max_{\lambda \in [0, \mathbb{R}^s]} \left\{ \langle \lambda, \mathbf{1} \rangle - \frac{1}{2} \sum_{i=1}^s (\underbrace{\lambda Y X^T Y \lambda}_{(Y X^T Y \lambda)_i})_i + \frac{1}{2\alpha} \|X^T Y \lambda\|^2 \right\}$$

$$= \max_{\lambda \in [0, \mathbb{R}^s]} \left\{ \langle \lambda, \mathbf{1} \rangle - \frac{1}{2} \|X^T Y \lambda\|^2 + \frac{1}{2\alpha} \|X^T Y \lambda\|^2 \right\}$$

$$= \max_{\lambda \in [0, \mathbb{R}^s]} \left\{ \langle \lambda, \mathbf{1} \rangle - \frac{1}{2\alpha} \|X^T Y \lambda\|^2 \right\}$$

L(λ)

$L(\lambda)$ is differentiable, concave

$$\nabla_{\lambda} L(\lambda) = \mathbf{1} - \frac{1}{2\alpha} \nabla_{\lambda} \langle Y X X^T Y \lambda, \lambda \rangle$$

$$= \mathbf{1} - \frac{1}{2} Y X X^T Y \lambda$$

Bad idea: $\lambda = \alpha (Y X X^T Y)^{-1} \mathbf{1}$.
 Restrictions !

Initial $\max_{x \in \mathbb{R}^n}$ problem can be written as

$$\max_{\lambda \in \mathbb{R}^s} \{L(\lambda) - x_{[0,1]^s}(\lambda)\}$$

Proximal gradient ascent
(max)

We can find λ 's by iterative procedure:

$$\lambda^{(k+1)} = \text{proj}_{[0,1]^s} \left(\lambda^{(k)} + \tau \left(\mathbf{I} - \frac{1}{2} \mathbf{X} \mathbf{X}^T \mathbf{Y} \right) \lambda^{(k)} \right)$$

Algorithm:

- Take $\lambda = (0, \dots, 0)$ and repeat the above steps until "convergence".
- Calculate $\hat{w} = \frac{1}{2} \mathbf{X}^T \mathbf{Y} \lambda$.

Recap: SVM problem

$$\hat{w} = \arg \min_w \left\{ \sum_{i=1}^S \max(0, 1 - y_i \langle x^{(i)}, w \rangle) + \frac{\lambda}{2} \|w\|^2 \right\}$$

$$\max(0, z) = \max_{\lambda \in [0, 1]} \lambda z$$

$$L(w, \lambda) = \sum_{i=1}^S \lambda_i (1 - y_i \langle x^{(i)}, w \rangle) + \frac{\lambda}{2} \|w\|^2$$

$$\hat{w} = \arg \min_w \max_{\lambda} L(w, \lambda)$$

$$\Leftrightarrow \max_{\lambda} \underbrace{\min_w L(w, \lambda)}_{\cdot}$$

$$\hat{w} = \frac{1}{\lambda} X^T Y$$

$\hat{\lambda} \longrightarrow \hat{w}$

Duality

Question: Let $f(x, y)$ be a real valued function. What could be proven about $\min_x \max_y f(x, y)$ and $\max_y \min_x f(x, y)$?

Statement 1: (Max-Min inequality)

$$\max_y \min_x f(x, y) \leq \min_x \max_y f(x, y)$$

for every y
we look for
 $x(y)$ such that

$$f(x(y), y) \leq f(x, y) \text{ for any } x.$$

↓ max over y

* we find
 y^* such that

$$f(x(y^*), y^*) \geq f(x(y), y)$$

$$\begin{aligned} & \overbrace{\quad\quad\quad}^{\substack{f(x, y(x)) \\ \geq f(x, y) \\ \forall y.}} \end{aligned}$$

Our goal to have

$$\max_y \min_x f(x,y) = \min_x \max_y f(x,y)$$

Theorem (von Neumann , 1928)

Let $X \subset \mathbb{R}^m$, $Y \subset \mathbb{R}^n$ be compact, convex sets. If $f: X \times Y \rightarrow \mathbb{R}$ is continuous and convex-concave, that means

$f(\cdot, y): X \rightarrow \mathbb{R}$ - convex

$f(x, \cdot): Y \rightarrow \mathbb{R}$ - concave

then

$$\min_{x \in X} \max_{y \in Y} f(x,y) = \max_{y \in Y} \min_{x \in X} f(x,y)$$

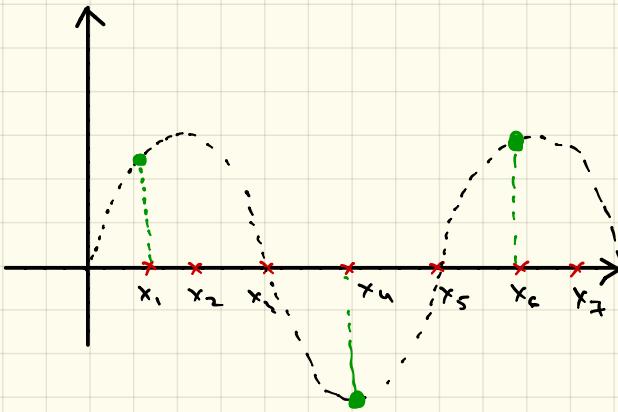
Topic 9: Interpolation with graphs.

Another method of solving classification problems.

Goal: Having a list of data inputs with partial list of outputs classify the rest of inputs based on their similarities

Problem: $\{x_i\}_{i \in I_1} \subset X$
 \uparrow not even numbers / vectors
for $i \in I_2 \subset I_1$ we are given class labels y_i . To find y_i for $i \in I_1$

Example:



$$I_1 = \{1, 2, \dots, 7\}$$

$$I_2 = \{1, 4, 6\}$$

+ additional information
about x_i 's.



$$f(x) \text{ for } x \in I_1$$

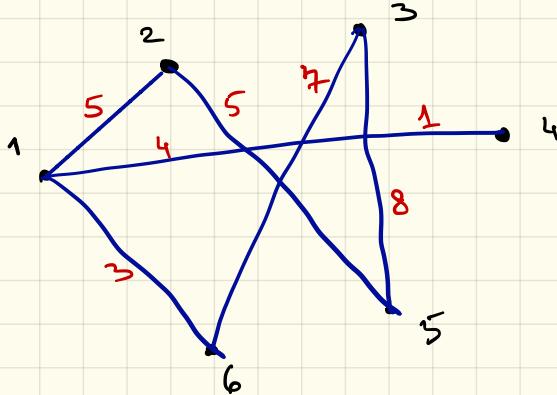
What is this additional information?

In this case this will be a similarity between points.

9.1 Graphs as a tool for similarity representation.

Graphs = nodes and edges.

We will consider undirected weighted graphs.



- edges don't have directions
- every edge (u,v) is enriched with the weight w_{uv} .

Def The incidence matrix of a graph $G = (V, E)$

is a matrix of a size $|E| \times |V|$ of the form

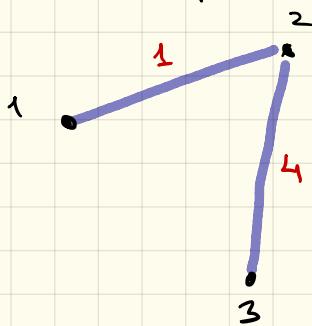
$$M_{\text{ev}} = \begin{cases} -\sqrt{w_{ij}}, & e = (i, j), v = i \\ \sqrt{w_{ij}}, & e = (i, j), v = j \\ 0, & e = (i, j), v \neq i, j \end{cases}$$

↑ edge ↓ node

The Laplacian matrix of G is given by

$$L = M^T M$$

Example:



$$V = \begin{pmatrix} v_1 & v_2 & v_3 \\ 1, & 2, & 3 \end{pmatrix}$$

$$E = \left(\begin{matrix} (1, 2) & ; & (2, 3) \\ e_1 & & e_2 \end{matrix} \right)$$

$$M = \begin{matrix} e_1 & \begin{pmatrix} v_1 & v_2 & v_3 \\ -1 & 1 & 0 \\ 0 & -2 & 2 \end{pmatrix} \\ e_2 & \end{matrix}$$

$$L = M^T M = \begin{pmatrix} -1 & 0 \\ 1 & -2 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} -1 & 1 & 0 \\ 0 & -2 & 2 \end{pmatrix}$$

$$= \begin{pmatrix} 1 & -1 & 0 \\ -1 & 5 & -4 \\ 0 & -4 & 4 \end{pmatrix} \quad \text{Wuv}$$

$$L \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} \approx \begin{pmatrix} \tilde{f}'(x_1) \\ \vdots \\ \tilde{f}'(x_n) \end{pmatrix}$$

↑ total weight

↑ weighted derivative.