$\text{Def}^n$: Functional margin of a hyperplane $(w, b)$

w.r.t. $(x^{(i)}, y^{(i)})$ is:

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b)$$

If $y^{(i)} = 1$. want $w^T x^{(i)} + b \gg 0$
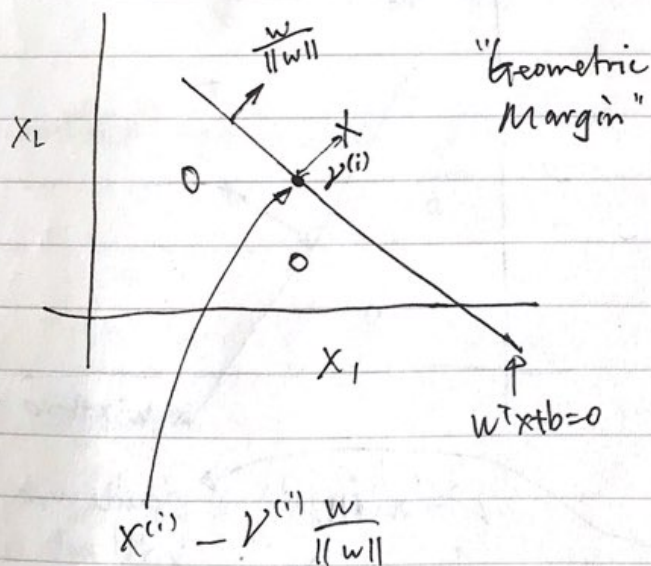
If $y^{(i)} = -1$. want $w^T x^{(i)} + b \ll 0$

If $y^{(i)}(w^T x^{(i)} + b) > 0$. then classified $(x^{(i)}, y^{(i)})$ correctly.

$$\hat{\gamma} = \min_i \hat{\gamma}^{(i)} \qquad (\text{worst case})$$

$w \to 2w$

$b \to 2b$.

then $\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b)$ doubles. so: $\|w\| \overset{\text{set}}{=} 1$.



"Geometric Margin"

$x^{(i)} - \gamma^{(i)} \dfrac{w}{\|w\|}$

$w^T \left( x^{(i)} - \gamma^{(i)} \dfrac{w}{\|w\|} \right) + b = 0$.

$w^T x^{(i)} + b = \gamma^{(i)} \dfrac{w^T w}{\|w\|} = \gamma^{(i)} \|w\|$

$$\gamma^{(i)} = \left( \dfrac{w}{\|w\|} \right)^T x^{(i)} + \dfrac{b}{\|w\|}$$

More generally, geometric margin

$$\gamma^{(i)} = y^{(i)} \left( \dfrac{w^T}{\|w\|} x + \dfrac{b}{\|w\|} \right).$$

If $\|w\| = 1$:

$$\hat{\gamma}^{(i)} = \gamma^{(i)}$$

$$\boxed{\gamma^{(i)} = \dfrac{\hat{\gamma}^{(i)}}{\|w\|}}$$

Geometric margin:

$$\gamma = \min_i \gamma^{(i)}$$

Maximum classifier:

$$\max_{\gamma, w, b} \gamma$$
$$s.t. \quad y^{(i)}(w^T x^{(i)} + b) \geq \gamma$$
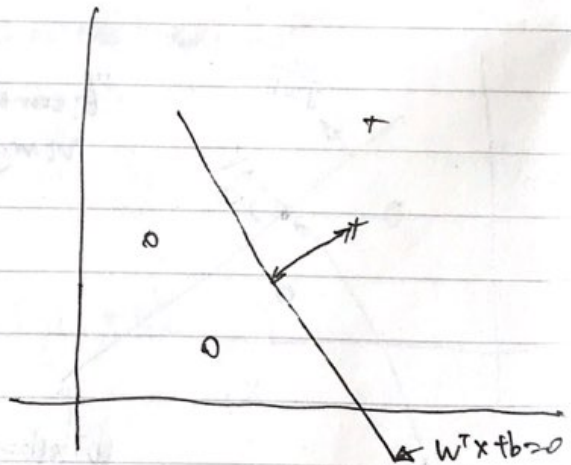$$\|w\| = 1 \qquad (\forall i)$$

$w \rightarrow 10w$

$b \rightarrow 10b$

does not change Geometric margin

---

[Lecture 7]

- Optimal Margin Classifier
- Primal/Dual optimization
  problem (KKT)
- SVM dual
- Kernels

$$h_{w,b}(x) = g(w^T x + b)$$
$$g(z) = \begin{cases} 1 & if \ z \geq 0 \\ 0 & otherwise \end{cases}$$
$$y \in \{-1, 1\}$$



$\leftarrow w^T x + b = 0$

Func. Margin: $\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b)$

Geo. Margin: $\gamma^{(i)} = y^{(i)}\left(\frac{w^T}{\|w\|} x^{(i)} + \frac{b}{\|w\|}\right)$

equivalent edition of
the optimal problem:

$\gamma = \min_i \gamma^{(i)}$ → double $w, b$

$\hat{\gamma} = \min_i \hat{\gamma}^{(i)}$ will not change the position of

$\|w\| = 1$

$|w_1| = 1$ the hyperplane

$w_1^2 + |w_1| = 17$

**#1:** $\max\limits_{\nu,w,b} \nu$

$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \nu \quad (i=1,...,m)$

$\|w\| = 1 \quad \leftarrow \text{Non-convex.}$

**#2:** $\max\limits_{\hat{\nu},w,b} \dfrac{\hat{\nu}}{\|w\|}$

$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq \hat{\nu}$

$\dfrac{\hat{\nu}}{\|w\|} = \nu$

$\boxed{\hat{\nu}=1} \leftarrow \text{impose this constraint}$

$\text{then } \min\limits_{i} y^{(i)}(w^T x^{(i)} + b) = 1.$

**#3:** $\min\limits_{w,b} \|w\|^2 \quad \leftarrow \text{max } \dfrac{1}{\|w\|}$

$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1$



**Convex optimization:**

$\min\limits_{w} f(w)$

$\text{s.t. } h_i(w) = 0, \quad i=1,...,l.$

**Lagrangian:**

$\mathcal{L}(w,\beta) = f(w) + \sum\limits_{i=1}^{l} \beta_i h_i(w)$

Lagrange Multipliers

$\dfrac{\partial \mathcal{L}}{\partial w} \overset{\text{set}}{=} 0 \qquad \dfrac{\partial \mathcal{L}}{\partial \beta} \overset{\text{set}}{=} 0.$

QP: quadratic programming software
gradient descent algo.

$h(w) = \begin{pmatrix} h_1(w) \\ h_2(w) \\ \vdots \\ h_l(w) \end{pmatrix} = \vec{0}$

for $w^*$ to be a solution, it is necessary that

$\exists \beta^*$ st. $\dfrac{\partial f(w^*, \beta^*)}{\partial w} = 0$. $\dfrac{\partial f(w^*, \beta^*)}{\partial \beta} = 0$.

## Primal Problem

$\min\ f(w)$

St. $g_i(w) \leq 0$, $i = 1 \dots k$ $\quad$ ("$\vec{g}(w) \leq \vec{0}$")

$\quad h_j(w) = 0$. $j = 1 \dots \ell$ $\quad$ ("$\vec{h}(w) = \vec{0}$")

Lagrangian:

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{j=1}^{\ell} \beta_j h_j(w)$$

Define: $\theta_p(w) = \max_{\substack{\alpha \geq 0 \\ \beta}} \mathcal{L}(w, \alpha, \beta)$

Consider:

$$p^* = \min_{w} \max_{\substack{\alpha \geq 0 \\ \beta}} \mathcal{L}(w, \alpha, \beta) = \min_{w} \theta_p(w)$$

$\theta_p(w)$ : - If $g_i(w) > 0$

$\qquad$ then $\theta_p(w) = \infty$

$\qquad$ - If $h_i(w) \neq 0$

$\qquad$ then $\theta_p(w) = \infty$

$\qquad$ - Otherwise, $\theta_p(w) = f(w)$

Then $\theta_p(w) = \begin{cases} f(w), & \text{if constraints satisfied } (g, h) \\ \infty, & \text{otherwise} \end{cases}$

So. $\min_{w} \theta_p(w) = $ original problem

## Dual Problem

$$\theta_D(\alpha, \beta) = \min_w \mathcal{L}(w, \alpha, \beta)$$

$$d^* = \max_{\substack{\alpha \geq 0 \\ \beta}} \min_w \mathcal{L}(w, \alpha, \beta) = \max_{\substack{\alpha \geq 0 \\ \beta}} \theta_D(\alpha, \beta)$$

$$d^* \leq p^* \qquad \boxed{\max \min (\dots) \leq \min \max (\dots)}$$

↑ general result for any function

$$\max_{y \in \{0,1\}} \left( \underbrace{\min_{x \in \{0,1\}} 1\{x = y\}}_{0} \right) \leq \min_{x \in \{0,1\}} \left( \underbrace{\max_{y \in \{0,1\}} 1\{x = y\}}_{1} \right)$$

Some conditions: $d^* = p^*$ $\boxed{\begin{array}{l}\text{for optimal margin classifier.} \\ \text{dual problem has better character.}\end{array}}$

---

let f be convex (Hessian $H \geq 0$)

suppose $h_i$ is affine ($h_i(w) = a_i^T w + b$)

and suppose $g_i$ is (strictly) feasible. ($\exists w$ st. $\forall_i \ g_i(w) < 0$)

Then $\exists w^*, \alpha^*, \beta^*$ st. $w^*$ solve primal

$\alpha^*, \beta^*$ solve dual, and $p^* = d^* = \mathcal{L}(w^*, \alpha^*, \beta^*)$

$\left( \begin{array}{l}
\dfrac{\partial}{\partial w} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0 \\[2mm]
\boxed{\alpha_i^* \ g_i(w^*) = 0} \quad \longleftarrow \text{KKT complementary condition} \\[2mm]
\dfrac{\partial}{\partial \beta} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0 \qquad \text{KKT: short for} \\[2mm]
g_i(w^*) \leq 0 \qquad\qquad\qquad \text{Karush-Kuhn-Tucker} \\[2mm]
\alpha_i^* \geq 0
\end{array} \right.$

KKT condition:

If $\alpha_i > 0 \implies g_i(w^*) = 0$ ($g_i(w)$ is an "active" constraint).

Lagrange multiplier: $\alpha_i, \beta_i \xrightarrow{\text{in SVM problem}} \alpha_i$

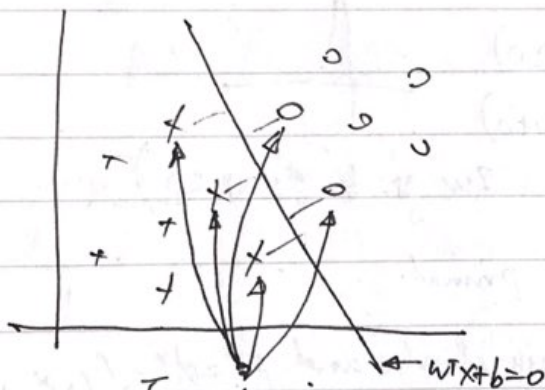Parameters : $w \xrightarrow{\text{in SVM problem}} w, b$ (notation change

$$\min \frac{1}{2}\|w\|^2$$

$$\text{S.t.} \quad \underbrace{y^{(i)}(w^T x^{(i)} + b) \geq 1}, \quad i = 1 \dots m.$$

$$g_i(w, b) = -y^{(i)}(w^T x^{(i)} + b) + 1 \leq 0$$

$$\alpha_i > 0 \implies g_i(w, b) = 0 \quad \text{(active constraint)}$$
$$\iff (x^{(i)}, y^{(i)}) \text{ has functional margin } 1.$$



Func. Margin = 1
(usually $\alpha \neq 0$)
"Support Vectors"

$w^T x + b = 0$

Support vectors : $\alpha_i > 0$

$\alpha_i = 0$ for most non-support vectors

$$\min \quad \mathcal{L}(w,b,\alpha) = \tfrac{1}{2}\|w\|^2 - \sum_{i=1}^{m}\alpha_i\left(y^{(i)}(w^T x^{(i)}+b)-1\right)$$

Dual problem:

$$\theta_D(\alpha) = \min_{w,b} \mathcal{L}(w,b,\alpha)$$

$$\nabla_w \mathcal{L}(w,b,\alpha) = w - \sum_{i=1}^{m}\alpha_i y^{(i)} x^{(i)} \overset{\text{set}}{=} 0 \quad \Rightarrow \quad \boxed{w = \sum_{i=1}^{m}\alpha_i y^{(i)} x^{(i)}}$$

$$\frac{\partial}{\partial b}\mathcal{L}(w,b,\alpha) = -\sum_{i=1}^{m}\alpha_i y^{(i)} \overset{\text{set}}{=} 0 \;.$$

$$\mathcal{L} = \tfrac{1}{2}\underbrace{w^T w}_{} - \sum_{i=1}^{m}\alpha_i\left(y^{(i)}(w^T x^{(i)}+b)-1\right)$$

$$\hookrightarrow \left(\sum_{i=1}^{m}\alpha_i y^{(i)} x^{(i)}\right)^T \left(\sum_{j=1}^{m}\alpha_j y^{(j)} x^{(j)}\right)$$

$$\mathcal{L} = \tfrac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)}\rangle$$

$$- \sum_{i=1}^{m}\sum_{j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)}\rangle + \sum_{i=1}^{m}\alpha_i$$

$$= \underbrace{\sum_{i=1}^{m}\alpha_i - \tfrac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j \boxed{\langle x^{(i)}, x^{(j)}\rangle}}_{W(\alpha)}$$

$$\max_{\alpha}\; \theta_D(\alpha)$$

- If $\sum_i y^{(i)}\alpha_i \neq 0$
$$\theta_D(\alpha) = -\infty$$

- If $\sum_i y^{(i)}\alpha_i = 0$
then $\theta_D(\alpha) = W(\alpha)$

Dual Problem:

$$\max \; W(\alpha)$$
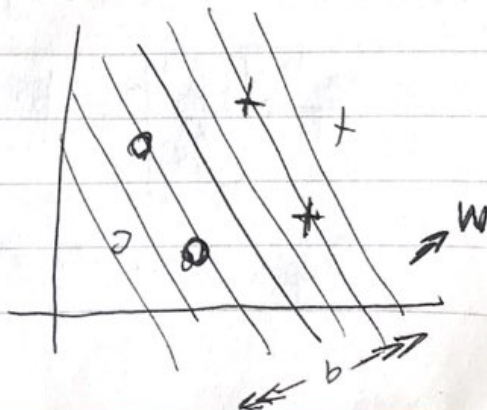$$\text{St.} \quad \alpha_i \geq 0$$
$$\sum_i y^{(i)}\alpha_i = 0$$

Solve for $\alpha$. (later)

$$w = \sum_{i=1}^{m}\alpha_i y^{(i)} x^{(i)}$$

$$b = \frac{\max\limits_{i:\, y^{(i)}=1} w^T x^{(i)} + \min\limits_{i:\, y^{(i)}=0} w^T x^{(i)}}{2}$$

$$w = \sum_i \alpha_i y^{(i)} x^{(i)}$$

$$h_{w,b}(x) = g(w^T x + b)$$
$$\uparrow \text{ threshold function}$$

$$w^T x + b = \sum_{i=1}^{m} \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b.$$

Kernels   $x^{(i)}$ —— very high dimentional $(x^{(i)} \in \mathbb{R}^{\infty})$
$\langle x^{(i)}, x^{(j)} \rangle$ efficiently computed

Lecture 8

SVM

- kernels
- soft margin
- SMO

$$\min \frac{1}{2} \|w\|^2$$
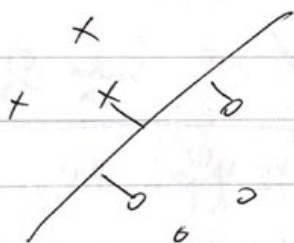$$\text{st. } y^{(i)}(w^T x^{(i)} + b) \geq 1$$

Dual Problem:

$$\max \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j \, y^{(i)} y^{(j)} \boxed{\langle x^{(i)}, x^{(j)} \rangle}$$

$$\text{st. } \alpha_i \geq 0$$
$$\sum_i y_i \alpha_i = 0$$
$$w = \sum_i \alpha_i y^{(i)} x^{(i)}$$



$$h_{w,b}(x) = g(w^T x + b)$$
$$= g\left( \sum_i \alpha_i y^{(i)} \boxed{\langle x^{(i)}, x \rangle} + b \right)$$

Have   $x \in \mathbb{R}$ living area

$$x \xrightarrow{\Phi} \begin{bmatrix} x \\ x^2 \\ x^3 \\ x^4 \end{bmatrix} = \phi(x)$$

Replace $\langle x^{(i)}, x^{(j)} \rangle$ with $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$

$\phi(x)$ – very high dim : cannot compute efficiently

$K(x^{(i)}, x^{(j)}) = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ : sometimes compute efficiently.

$\underbrace{\quad\quad\quad}$
$x, z \in \mathbb{R}^n$

$$K(x, z) = (x^T z)^2 = \left(\sum_{i=1}^{n} x_i z_i\right)\left(\sum_{j=1}^{n} x_j z_j\right) = \sum_{i=1}^{n} \sum_{j=1}^{n} (x_i x_j)(z_i z_j)$$

$$= (\phi(x))^T \phi(z)$$

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \\ \sqrt{2} x_1 \\ \sqrt{2} x_2 \\ \sqrt{2} x_3 \\ c \end{bmatrix}$$

Need $O(n^2)$ to compute $\phi(x)$.

But need $O(n)$ time to compute $K(x, z)$

$K(x, z) = (x^T z + c)^2$

Generally: $K(x, z) = (x^T z + c)^d$

$\rightarrow \binom{n+d}{d}$ features of all monomials up to degree $d$.

$x \longmapsto \phi(x) \qquad z \longmapsto \phi(z)$

$\langle \phi(x), \phi(z) \rangle$

$x, z$ similar : $\langle \phi(x), \phi(z) \rangle$ maybe large?
---- dissimilar : ⌢ ⌢ ⌢ ⌢ ... small.

How to choose kernel function?

33

$K(x, z)$ — large if $x, z$ similar

small if $x, z$ dissimilar

$$K(x, z) = e^{-\frac{\|x - z\|^2}{2\sigma^2}} \quad \rightarrow \text{invalid}?$$

$\exists \phi$ s.t. $K(x, z) = \langle \phi(x), \phi(z) \rangle$ ?

Suppose $K$ is a kernel. Let $\{x^{(1)}, \ldots, x^{(m)}\}$ be given

let $K \in \mathbb{R}^{m \times n}$

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

Then for any vector $z \in \mathbb{R}^m$. $\boxed{z^T K z}$

$$z^T K z = \sum_i \sum_j z_i K_{ij} z_j$$

$$= \sum_i \sum_j z_i \, \phi(x^{(i)})^T \phi(x^{(j)}) \, z_j$$

$$= \sum_i \sum_j z_i \sum_k \left(\phi(x^{(i)})\right)_k \left(\phi(x^{(j)})\right)_k z_j$$

$$= \sum_k \sum_i \sum_j z_i \left(\phi(x^{(i)})\right)_k \left(\phi(x^{(j)})\right)_k z_j$$

$$= \sum_k \left(\sum_i z_i \, \phi(x^{(i)})_k\right)^2 \geq 0$$

$K \geq 0$. (positive semi-definite).   And the inverse
                                        judgement is true

Theorem (Mercer): Let $K(x, z)$ be given, then $K$ is valid

(is (Mercer) Kernel ($\therefore$ e. $\exists \phi$ st. $K(x,z) = \phi(x)^T \phi(z)$))

iff. for all $\{x^{(1)} \ldots x^{(m)}\}$ ($m < \infty$) the kernel matrix

$K \in \mathbb{R}^{m \times m}$ is symmetric positive semidefinite.

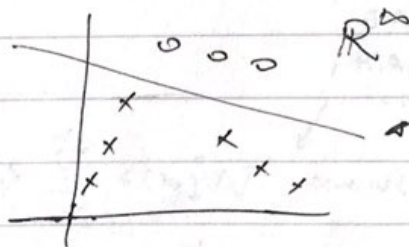eg: $K(x, x) = -1$   not valid. $\therefore$ $-1 \neq \phi(x)^T \phi(x)$

Choose $K(x, z) = e^{-\frac{\|x - z\|^2}{2\sigma^2}}$ (Gaussian Kernel)

or $(x^T z + \epsilon)^d$, or etc.

Replace $\langle x^{(i)}, x^{(j)} \rangle$ with $K(x^{(i)}, x^{(j)})$ ✓

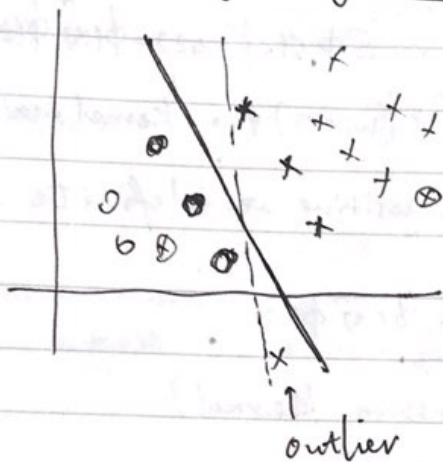$$x^{(i)} \longmapsto \phi(x^{(i)})$$ ✗

$x \in \mathbb{R}$



$\mathbb{R}^\infty$

→ have linear classifier in high-dim. space

$\langle x^{(i)}, x^{(j)} \rangle \to K(x^{(i)}, x^{(j)})$ for any algo. which can be written as inner-product form

$\langle x^{(i)}, x^{(j)} \rangle \to K(x^{(i)}, x^{(j)})$ for any algo. which can be written as inner-product form

35

# $L_1$ norm soft margin SVM



$$\min_{w, b, \xi_i} \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{m} \xi_i$$

$$\text{St. } y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0, \quad i = 1, \dots, m.$$

penalty

If $y^{(i)}(w^T x^{(i)} + b) > 0 \Rightarrow$ classified correctly.

$$\mathcal{L}(w, b, \xi, \alpha, r) = \frac{1}{2}\|w\|^2 + C \sum_i \xi_i$$
$$- \sum_{i=1}^{m} \alpha_i \left( y^{(i)}(w^T x^{(i)} + b) - 1 + \xi_i \right) - \sum_{i=1}^{m} r_i \xi_i$$
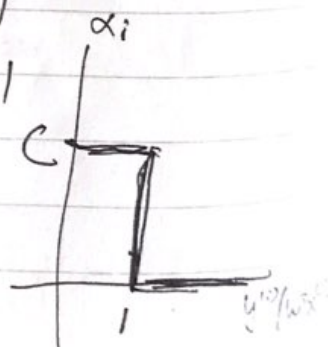
SOME
MATH

$$\max \ W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j \, y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle.$$

$$\text{St. } \sum_{i=1}^{m} y^{(i)} \alpha_i = 0$$

$$0 \leq \alpha_i \leq C. \quad i = 1, \dots, m$$

convergence criteria (derive from KKT condition):

$$\begin{cases} \alpha_i = 0 & \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1 \\ \alpha_i = C & \Rightarrow y^{(i)}(w^T x^{(i)} + b) \leq 1 \\ 0 < \alpha_i < C & \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1 \end{cases}$$

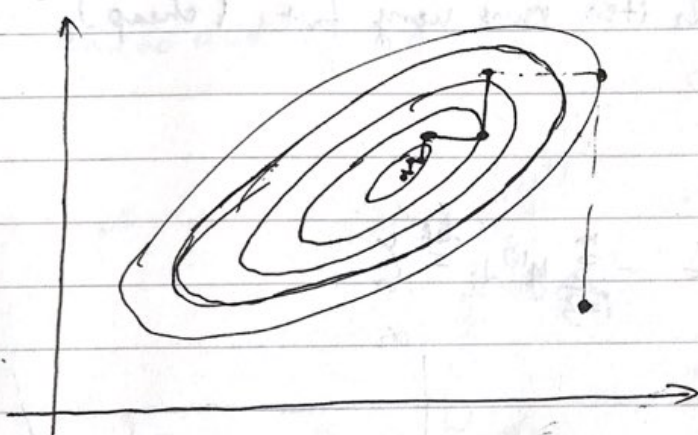Digression:

coordinate ascent (another opt. question)

maxs $W(\alpha_1, \ldots, \alpha_m)$ (no constraints on $\alpha_i$'s)

Repeat {
    For $i=1$ to $m$:

      Hold everything except $\alpha_i$ fixed.

}     $\alpha_i := \underset{\hat{\alpha}_i}{\arg\max} \, W(\alpha_1, \ldots \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1} \ldots \alpha_m)$

$W(\alpha_1, \ldots, \alpha_m)$

high dim: not fixed order $\alpha_1, \ldots \alpha_m$

    heuristic value function decide which is chosen to vary (change)

    * more steps to converge (compared to Newton's method)
    * inner loops executes very quick.

## SMO

Coordinate ascent cannot work directly on SVM dual opt. problem.

∴ constraints : $\sum_{i=1}^{m} y^{(i)} \alpha_i = 0$.

∴ Change 2 $\alpha_i$'s at a time.

SMO Algo. is due to

John Platt @ Microsoft

Outline:

Select $\alpha_i, \alpha_j$ (heuristic)

Hold all $\alpha_i$'s fixed except $\alpha_i, \alpha_j$

Optimize $W(\alpha)$ w.r.t. $\alpha_i \alpha_j$ s.t. constraints $\leftarrow$    (\*) key step

$$\left[ \begin{array}{l} \text{SMO do extremely efficiently:} \\ \qquad \text{Although many iterations (large number)} \\ \qquad \text{But each iter. runs very fast (cheap)} \end{array} \right.$$

Update $\alpha_1, \alpha_2$

Know $\sum_i y^{(i)} \alpha_i = 0$.

$$y^{(1)} \alpha_1 + y^{(2)} \alpha_2 = -\sum_{i=3}^{m} y^{(i)} \alpha_i \overset{def}{=} \zeta$$
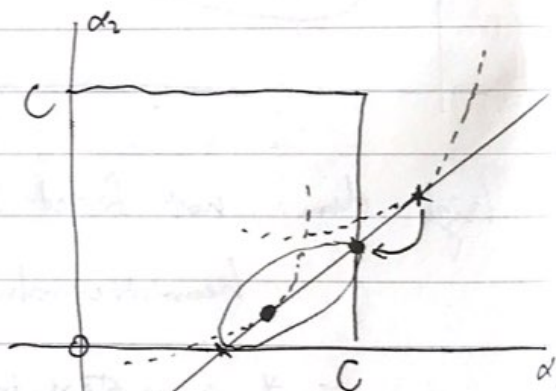
$$0 \le \alpha_i \le C.$$

$$W(\alpha_1, \alpha_2, \alpha_3 \dots \alpha_m)$$
$$= W\left(\frac{\zeta - y^{(2)} \alpha_2}{y^{(1)}}, \alpha_2, \alpha_3, \dots \alpha_m\right)$$
$$= a\alpha_2^2 + b\alpha_2 + c$$
(quadratic func.)

optimize quadratic func.

in inner loop operates very efficiently.

$$y^{(1)} \alpha_1 + y^{(2)} \alpha_2 = \zeta$$
$$\Rightarrow \alpha_1 = \frac{\zeta - y^{(2)} \alpha_2}{y^{(1)}}$$

Compute $b$ is not hard. Do it after class.



38

## Applications of SVM

① Handwriter's Digit Recognition

$$K(x,y) = (x^T y)^d \quad \text{or} \quad e^{-\frac{\|x-y\|^2}{2\sigma^2}}$$

SVM is comparable with best neural networks

$(0 \leq x \leq 0$

$x \in \mathbb{R}^{100}$

② Classify protein seq's

amino acid seq. $(A . \sim Z)$

BAJTSIAJBAJTAU

$\phi(x) = ?$ (hard problem...)

$$\phi(x) \in \mathbb{R}^{(20^4)}$$

$$= \mathbb{R}^{160000}$$

$$
\begin{array}{l|c}
AAAA & 0 \\
AAAB & 0 \\
AAAC & 0 \\
\vdots & \vdots \\
AAAZ & 0 \\
AABA & 0 \\
\vdots & \vdots \\
BAJT & 2 \quad \leftarrow \text{occur 2 times} \\
\vdots & \vdots \\
TSIA & 1 \\
\vdots & \vdots \\
ZZZZ & 0
\end{array}
\Bigg\} = \phi(x)
$$

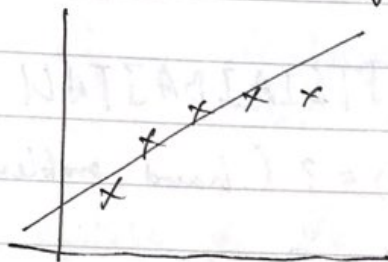use Dynamic Programming (DP):

compute $\phi(x)^T \phi(z)$.

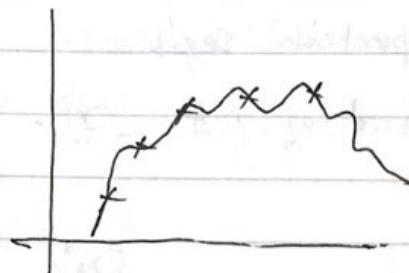Basic (Applicative) Part of this course is over.

NEXT:    Understanding these algorithms
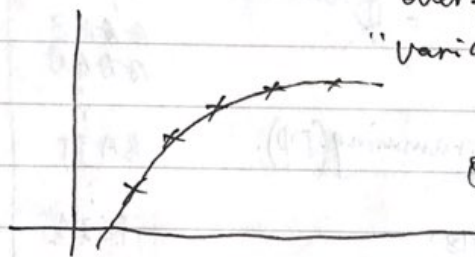
39

## Lecture 9 — learning theory

- Bias / Variance
- Empirical Risk Minimization
- Union Bound / Hoeffding Inequality
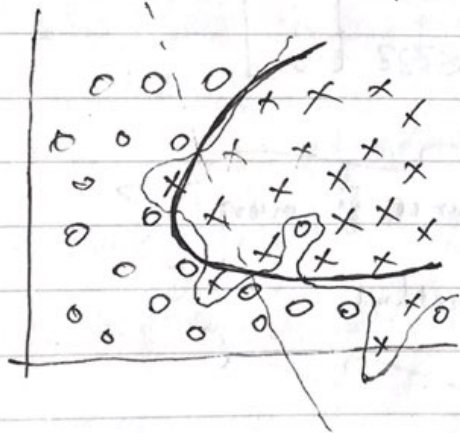- Uniform convergence.

$\theta_0 + \theta_1 x$
"underfit"
"bias"

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
"overfit"
"variance"

$\theta_0 + \theta_1 x + \theta_2 x^2$
"fit"
"balanced"

$h_\theta(x) = g(\theta_0 + \theta_1 x + \cdots + \theta_n x^4)$

Linear Classification:

$h_\theta(x) = g(\theta^T x)$

$g(z) = 1\{z \geq 0\}$      (note $y \in \{0,1\}$)

$S = \{(x^{(i)}, y^{(i)})\}_{i=1}^m$ (training set)     $(x^{(i)}, y^{(i)}) \sim \mathcal{D}$

Training error of $h_\theta$:

$$\hat{\varepsilon}(h_\theta) = \hat{\varepsilon}_S(h_\theta) = \frac{1}{m} \sum_{i=1}^m 1\{h_\theta(x^{(i)}) \neq y^{(i)}\}$$

ERM:    $\hat{\theta} = \arg\min_\theta \hat{\varepsilon}_S(h_\theta)$      $\leftarrow$ SVM. logistic regression.
                                       is approx. of ERM.

                                      (ERM is more general)

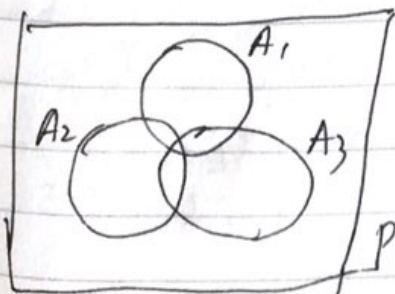Hypothesis class $\mathcal{H} = \{h_\theta : \theta \in \mathbb{R}^{n+1}\}$.

            $h_\theta : X \longmapsto \{0,1\}$.

     ERM:    $\hat{h} = \arg\min_{h \in \mathcal{H}} \hat{\varepsilon}_S(h)$

Generalization error:

$$\varepsilon(h) = P_{(x,y) \in \mathcal{D}}(h(x) \neq y).$$

---

Union Bound:    let $A_1, A_2 \ldots A_k$ be $k$ event.
                                      (not necessarily independent)

                   Then

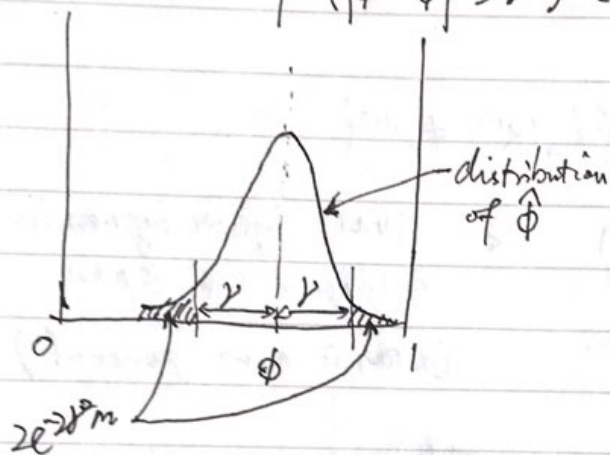$$P(A_1 \cup A_2 \cup \cdots \cup A_k) \leq P(A_1) + P(A_2) + \cdots + P(A_k)$$



     "or"

$$P(A_1 \cup A_2 \cup A_3) \leq P(A_1) + P(A_2) + P(A_3)$$

**Hoeffding Inequality :** Let $z_1, \ldots z_m$ be $\sim$ I.I.D

Bernoulli $(\phi)$ random variable $(P(z_i = 1) = \phi)$.

Let $\hat{\phi} = \frac{1}{m} \sum_{i=1}^{m} z_i$ and let any $\gamma > 0$ be fixed.

Then
$$P\left(|\hat{\phi} - \phi| > \gamma\right) \leq 2 e^{-2\gamma^2 m}$$



distribution of $\hat{\phi}$

$2e^{-2\gamma^2 m}$

$0$   $\phi$   $1$

$\gamma$   $\gamma$

**Central Limit Theorem :** works only m $\bar{\imath}$s large.

## The case of finite $\mathcal{H}$

$\mathcal{H} = \{h_1, h_2, \ldots h_k\}$   $k$ hypothesis

$\hat{h} = \underset{h \in \mathcal{H}}{\arg\min} \; \hat{\varepsilon}_s(h_i)$

Strategy:   (1) $\hat{\varepsilon} \approx \varepsilon$

(2) Show bound on $\varepsilon(\hat{h})$.

(1) **Fix** any $h_j \in \mathcal{H}$.

Define $Z_i = \mathbb{1}\{h_j(x^{(i)}) \neq y^{(i)}\} \in \{0,1\} \sim$ Bernoulli

$P(Z_i = 1) = \varepsilon(h_j)$. (All $Z_i \sim$ I.I.D)

$\hat{\varepsilon}(h_j) = \frac{1}{m}\sum_{i=1}^{m} Z_i = \frac{1}{m}\sum_{i=1}^{m} \mathbb{1}\{h_j(x^{(i)}) \neq y^{(i)}\}$

$\underbrace{\qquad}_{\text{mean } \varepsilon(h_j)}$

By Hoeffding Ineq:

$$P\left(\left|\varepsilon(h_j) - \hat{\varepsilon}(h_j)\right| \geq \nu\right) = 2 e^{-2\nu^2 m}.$$

$A_j = $ event that $\left|\varepsilon(h_j) - \hat{\varepsilon}(h_j)\right| > \nu$.

$P(A_j) \leq 2e^{-2\nu^2 m}.$

$P\left(\exists\, h_j \in \mathcal{H} \text{ st. } \left|\varepsilon(h_j) - \hat{\varepsilon}(h_j)\right| > \nu\right)$

$= P(A_1 \cup A_2 \cup \cdots \cup A_k) \leq \sum_{i=1}^{k} P(A_i) \leq \sum_{i=1}^{k} 2e^{-2\nu^2 m} = 2ke^{-2\nu^2 m}.$

( 1 − both sides ):

$P\left(\nexists\, h_j \in \mathcal{H} \text{ st. } \left|\varepsilon(h_j) - \hat{\varepsilon}(h_j)\right| > \nu\right)$

$= P\left(\forall\, h_j \in \mathcal{H} \text{ st. } \left|\varepsilon(h_j) - \hat{\varepsilon}(h_j)\right| \leq \nu\right) \geq 1 - 2ke^{-2\nu^2 m}$

So w.p. $1 - 2ke^{-2\nu^2 m}$, $\hat{\varepsilon}(h)$ will be within $\nu$ of $\varepsilon(h)$ for all $h \in \mathcal{H}$.

"with probability"

"Uniform Convergence"

Given $\nu$ and $\delta$, what is $m$?

$$\delta = 2Ke^{-2\nu^2 m}, \text{ solve for } m.$$

So long as $\quad m \geq \frac{1}{2\nu^2} \log \frac{2k}{\delta}$,

then w.p. $1-\delta$, we have $|\mathcal{E}(h) - \hat{\mathcal{E}}(h)| \leq \nu$ for all $h \in \mathcal{F}$

"Sample complexity" bound.

$\forall k, \log k \leq 30$
in practical

"Error bound"

Solve for $\nu$ for fixed $m, \delta$.

w.p. $1-\delta$, we have that $\forall h \in \mathcal{H}$.

$$|\hat{\mathcal{E}}(h) - \mathcal{E}(h)| \leq \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}} \qquad \text{"}\nu\text{"}$$

(2) Let assume $\forall h \in \mathcal{H}$. $|\mathcal{E}(h) - \hat{\mathcal{E}}(h)| \leq \nu$. ①

Can we prove something about $\mathcal{E}(\hat{h})$.

$$\hat{h} = \arg\min_{h \in \mathcal{H}} \hat{\mathcal{E}}(h) \qquad ②$$

$$h^* = \arg\min_{h \in \mathcal{H}} \mathcal{E}(h) \qquad ③$$

$$\mathcal{E}(\hat{h}) \leq \hat{\mathcal{E}}(\hat{h}) + \nu \qquad -- \text{ by } ①$$
$$\leq \hat{\mathcal{E}}(h^*) + \nu \qquad -- \text{ by } ②$$
$$\leq \mathcal{E}(h^*) + \nu + \nu$$
$$= \mathcal{E}(h^*) + 2\nu.$$

$\hat{\mathcal{E}}(h)$ — training error of $h$

$\mathcal{E}(h)$ — generalization error of $h$.

44

**Theorem.** Let $|\mathcal{H}| = k$ and let $m, \delta$ be fixed, then w.p. $1-\delta$

$$\varepsilon(\hat{h}) \leq \min_{h \in \mathcal{H}} \varepsilon(h) + 2\sqrt{\frac{1}{2m} \log \frac{2k}{\delta}} \qquad (\ast)$$

$\underbrace{\qquad}_{\text{"bias"}} \quad \underbrace{}_{\varepsilon(h^*)} \qquad \underbrace{}_{\text{"variance"}}$

Set $\gamma = \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}}$, we know ① holds

w.p. $1-\delta$. (i.e. $|\varepsilon(h) - \hat{\varepsilon}(h)| \leq \gamma \quad \forall h \in \mathcal{H}$)

which implies $(\ast)$

$\mathcal{H} \subseteq \mathcal{H}'$
$\uparrow$ linear func. $\quad \uparrow$ quadratic func.

"underfit" high bias

"overfit" high variance

error

generalization error

training error

model complexity (Degree of polynomial. size of $\mathcal{H}$ etc.)

($\tau$ in locally-weighted ~~linear~~ linear regression)

**Corollary.** Let $|\mathcal{H}| = k$, let any $\delta, \gamma$ be fixed.

Then form
$$\varepsilon(\hat{h}) \leq \min_{h \in \mathcal{H}} \varepsilon(h) + 2\gamma$$

w.p. $1-\delta$, it suffices that

$$m \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta} = O\left(\frac{1}{\gamma^2} \log \frac{k}{\delta}\right) \leftarrow \text{important when generalize to infinite hypothesis class } \mathcal{H}.$$

45

VC dimension

Model Selection

  - Cross validation

  - feature selection

Bayesian statistics & Regularization

Let $|\mathcal{H}| = k$, and let $\gamma, \delta$ be fixed.

Then for $\varepsilon(\hat{h}) \leq \min_{h \in \mathcal{H}} \varepsilon(h) + 2\gamma$ w.p. $1-\delta$,

it suffices that

$$m \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta} = O\left(\frac{1}{\gamma^2} \log \frac{k}{\delta}\right)$$

* generalize to infinite hypothesis class:

Say $\mathcal{H}$ is parameterized by $d$ real numbers (linear boundary

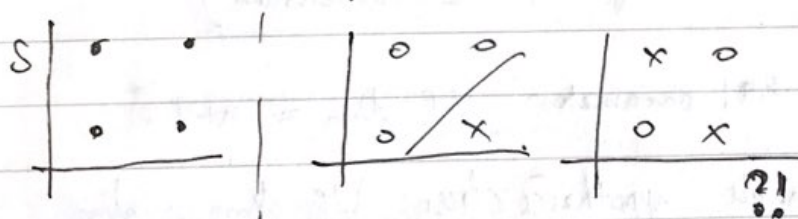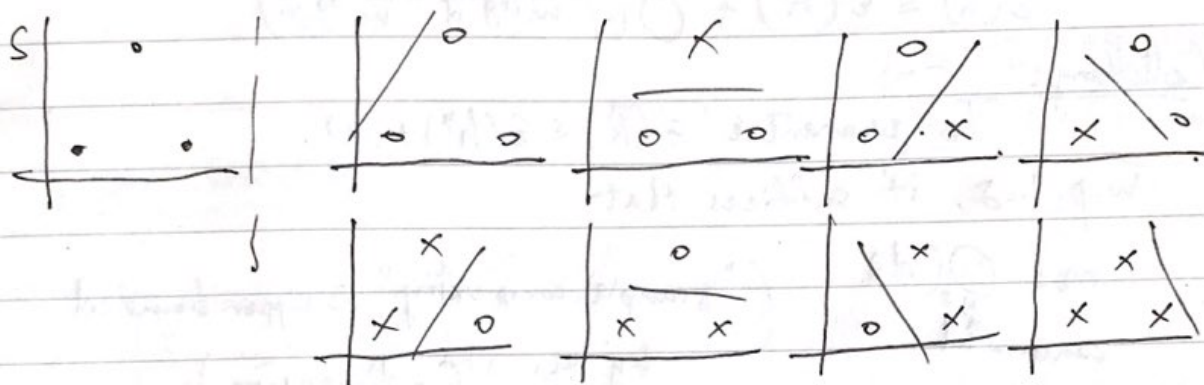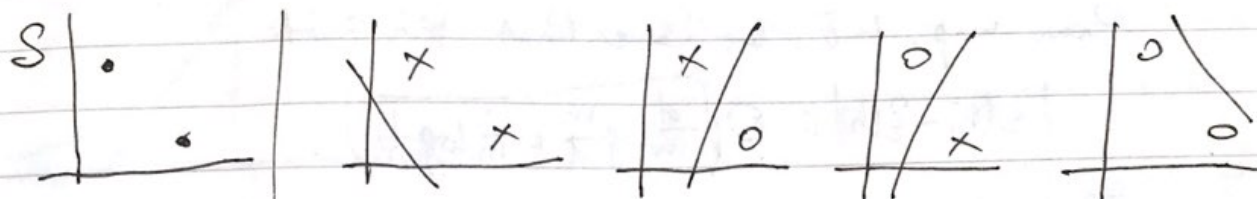$\rightarrow 64d$ bits in computer. digital

$k = |\mathcal{H}| = 2^{64d}$.

Sufficed that: $m \geq O\left(\frac{1}{\gamma^2} \log \frac{2^{64d}}{\delta}\right)$

$$= O\left(\frac{d}{\gamma^2} \log \frac{1}{\delta}\right).$$

          intuitively

Definition: Given a set $S = \{x^{(1)} \ldots x^{(d)}\}$. more forma[l]

we say $\mathcal{H}$ <u>shatters</u> $S$ if $\mathcal{H}$ can <u>realize any labeling</u>

on it.

$\mathcal{H} = \{$linear classifiers in 2D$\}$

$S$

$S$

$S$

?!

Definition: The Vapnik-Chervonenkis dimension of $\mathcal{H}$
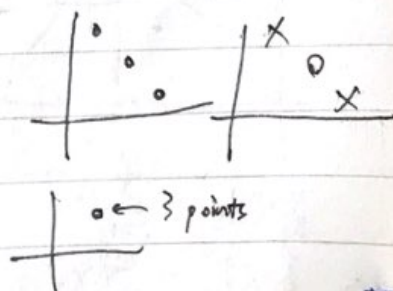(i.e. $VC(\mathcal{H})$) is the size of the largest set shattered by $\mathcal{H}$

E.g. $\mathcal{H} = \{$linear classifiers in 2D$\}$.

$VC(\mathcal{H}) = 3$.

* Only need exist one set of 3 points can be shattered.

More generally, in $n$ dimensions,

$VC(\{$linear classifiers in $n$ dim$\}) = n+1$.

$\leftarrow$ 3 points

47

**Theorem**. Let $\mathcal{H}$ be given and let $VC(\mathcal{H})=d$.

then w.p. $1-\delta$, we have that $\forall h \in \mathcal{H}$:

$$|\varepsilon(h) - \hat{\varepsilon}(h)| \leq O\left(\sqrt{\frac{d}{m}\log\frac{m}{d} + \frac{1}{m}\log\frac{1}{\delta}}\right).$$

Thus, w.p. $1-\delta$, are also have

$$\varepsilon(\hat{h}) \leq \varepsilon(h^*) + O\left(\sqrt{\frac{d}{m}\log\frac{m}{d} + \frac{1}{m}\log\frac{1}{\delta}}\right)$$

**Collollary:**

To gharantee $\varepsilon(\hat{h}) \leq \varepsilon(h^*) + 2\nu$,

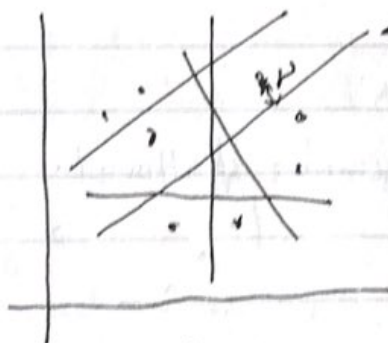w.p. $1-\delta$, it suffices that

$$m = \underset{\text{const}}{\underset{\uparrow \uparrow}{O}}_{\nu,\delta}(d)$$

("sample complexity" is upper bounded by the VC-dimension)

Eg. logistic, $n+1$ parameter $VC$ dim $= n+1$.

\* For most reasonable hypothesis classes. VC-dim usually <u>linear</u> in $\#$ parameters of the model/hypothesis.

\* Class of linear <u>seperators</u> with large margin actually has low Vc-dim.



If $\|x^{(i)}\|_2 \leq R$.

$$VC(\mathcal{H}) \leq \left\lceil\frac{R^2}{4\gamma^2}\right\rceil + 1.$$

SVM automatically find low VC-dim classifier.

$$\|x\|_2^2 = \sum_{i=1}^{n} x_i^2$$

$$\|x\|_2^2 = \sum_{i=1}^{\infty} x_i^2 \quad (\text{converge condition}).$$