# Machine Learning from Data
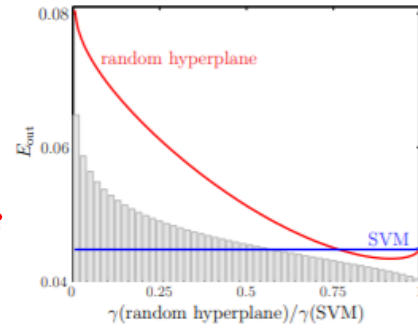
Lecture 25: Spring 2021

# Today's Lecture
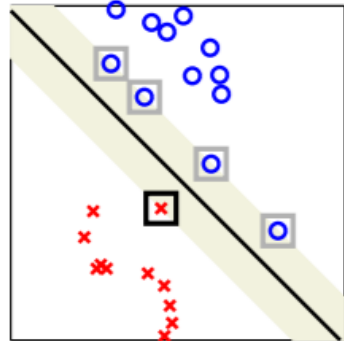
- The Kernel Trick ✓

# Large Margin is Better

## Controling Overfitting



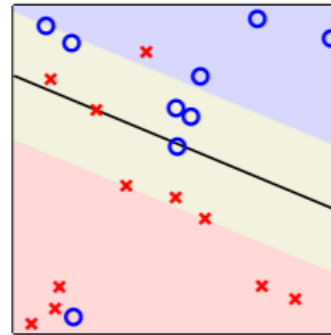$$E_{out}$$

random hyperplane

SVM

$$\gamma(\text{random hyperplane})/\gamma(\text{SVM})$$

**Theorem.** $d_{vc}(\gamma) \leq \left\lceil \dfrac{R^2}{\gamma^2} \right\rceil + 1$



$$E_{cv} \leq \frac{\#\ \text{support vectors}}{N}$$

## Non-Separable Data



$$\text{minimize}_{b,\mathbf{w},\xi} \quad \frac{1}{2}\mathbf{w}^{\mathsf{T}}\mathbf{w} + C\sum_{n=1}^{N}\xi_n$$

$$\text{subject to:} \quad y_n(\mathbf{w}^{\mathsf{T}}\mathbf{x}_n + b) \geq 1 - \xi_n$$

$$\xi_n \geq 0 \qquad\qquad \text{for } n = 1,\dots,N$$



$\Phi_2 + \text{SVM}$     $\Phi_3 + \text{SVM}$     $\Phi_3 + \text{pseudoinverse algorithm}$

Complex hypothesis that does not overfit because it is 'simple', controlled by only a few support vectors.

**1. Original data**
$$\mathbf{x}_n \in \mathcal{X}$$

**2. Transform the data**
$$\mathbf{z}_n = \Phi(\mathbf{x}_n) \in \mathcal{Z}$$

$\Phi$

**4. Classify in $\mathcal{X}$-space**
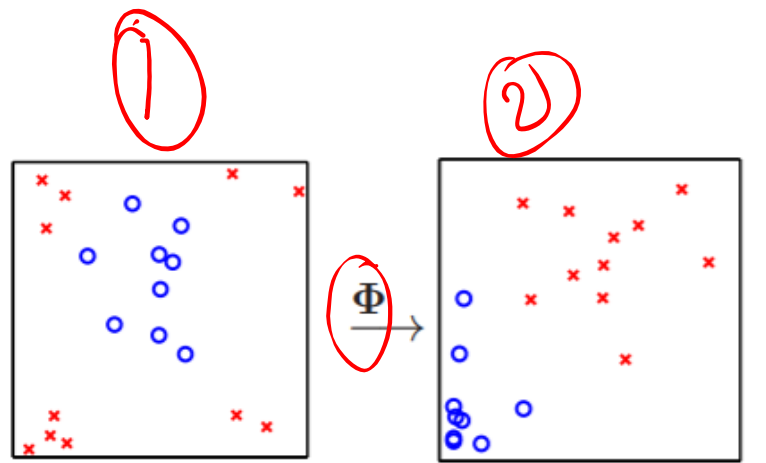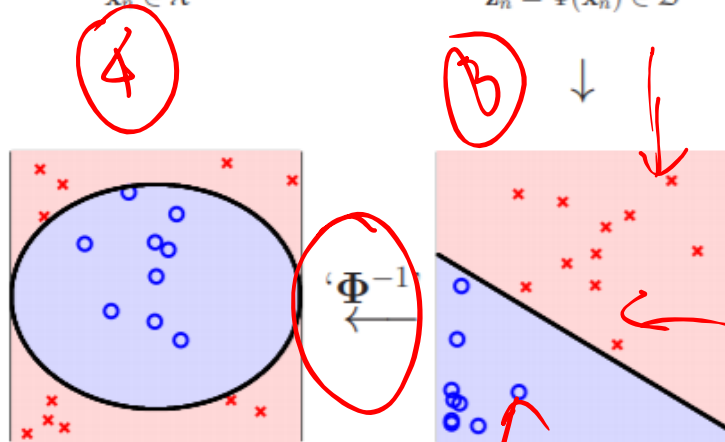$$g(\mathbf{x}) = \tilde{g}(\Phi(\mathbf{x})) = \text{sign}(\tilde{\mathbf{w}}^\mathsf{T}\Phi(\mathbf{x}))$$

**3. Separate data in $\mathcal{Z}$-space**
$$\tilde{g}(\mathbf{z}) = \text{sign}(\tilde{\mathbf{w}}^\mathsf{T}\mathbf{z})$$

$`\Phi^{-1}`$

$\mathcal{X}$-space is $\mathbb{R}^d$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N$

$y_1, y_2, \ldots, y_N$

no weights

$d_{vc} = d + 1$

$g(\mathbf{x}) = \text{sign}(\tilde{\mathbf{w}}^\mathsf{T}\Phi(\mathbf{x}))$

$\mathcal{Z}$-space is $\mathbb{R}^{\tilde{d}}$

$$\mathbf{z} = \Phi(\mathbf{x}) = \begin{bmatrix} 1 \\ \Phi_1(\mathbf{x}) \\ \vdots \\ \Phi_{\tilde{d}}(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ \vdots \\ z_{\tilde{d}} \end{bmatrix}$$

$\mathbf{z}_1, \mathbf{z}_2, \ldots, \mathbf{z}_N$

$y_1, y_2, \ldots, y_N$

$$\tilde{\mathbf{w}} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{\tilde{d}} \end{bmatrix}$$

$d_{vc} = d + 1$

Have to **transform** the data to the $\mathcal{Z}$-space.

# Today's Lecture

- How to use nonlinear transforms without physically transforming data to Z-space?

$$\begin{cases} \min_{\underline{\omega, b}} & \frac{1}{2} \omega^T \omega \qquad\qquad\qquad\qquad\qquad\qquad \underline{QP} \\ \text{s.t.} & y_n(\omega^T x_n + b) \geq 1 \qquad \forall n = 1, 2 \dots \underline{\underline{N}} \end{cases}$$

$\rightarrow d + 1$ optimization variables, $N$ constraints.

$\rightarrow$ Transformed Space $\rightarrow$ Large no. of features.

## $\underline{PRIMAL} \iff \underline{\underline{DUAL}}$

$\downarrow$

lagrange's multiplier. $\rightarrow \alpha$

$$\min_{\alpha} \left\{ \frac{1}{2} \sum_{n=1}^{M} \sum_{m=1}^{N} \alpha_n \alpha_m \, y_n y_m \left( \underbrace{x_n^T x_m}_{\uparrow} \right) - \sum_i \alpha_i \right\}$$

$$s.t. \quad \sum_n \alpha_n y_n = 0$$

$$\alpha_n \geq 0 \qquad \qquad \qquad QP$$

$$\text{Output} \longrightarrow \alpha_n (\underline{\alpha_n^*})$$

$$\omega^* = \sum_{n=1}^{N} \alpha_n^* y_n x_n \quad \Big\} \text{classification}$$

Pick some points

$$\alpha_s^* > 0 \{ Sv \}$$

$$b^* = y_s - (\omega^*)^T x_s$$

**Dual:** N optimization variables, $N+1$ constraints (simple)

Does not depend on dimensions.

$$d \longleftrightarrow N$$

## Lagrangian

$$\Rightarrow L(\omega, b, \alpha) = \frac{1}{2}\omega^T\omega + \sum_{n=1}^{N} \alpha_n(1 - y_n(\omega^T x_n + b))$$

Lagrangian $\alpha_n \geqslant 0$

$\underline{\min} L$ w.r.t. $\omega, b$ ✓

$\underline{\underline{\max}} L$ w.r.t. $\alpha's (\alpha_n \geqslant 0)$ ✓

or

"Proof": $1 - y_n(\omega^T x_n + b) \leq 0 - ①$

OR $1 - y_n(\omega^T x_n + b) > 0 - ②$ for each $n$

$$\alpha_n(1 - y_n(\omega^T x_n + b)) > \underline{\underline{0}}$$

$$\alpha_n \longrightarrow \text{large} \longrightarrow \alpha \longrightarrow L \to \infty \times$$

$$\underbrace{1 - y_n(w^T x_n + b)}_{} \leq 0 \quad \Longrightarrow \underbrace{y_n(w^T x_n + b) \geq 1}_{} \}$$

$$\left(1 - y_n(w^T x_n + b)\right) = 0 \lor \text{ or } \alpha_n = 0$$

$$\alpha_n \underline{\underline{\quad}}\left(1 - y_n(w^T x_n + b)\right) \cancel{\leq} 0 \} \qquad \alpha_n \geq 0$$

$$\searrow_{-ve}$$

## Summarize

$$\alpha_n \left(1 - y_n(w^T x_n + b)\right) = 0 \checkmark$$

$$y_n(w^T x_n + b) = 1 \quad \checkmark \longrightarrow x_n \text{ is } \underline{\underline{SV}}$$

$$\alpha_n = 0 \checkmark$$

$$L^* = \frac{1}{2} w^T w \longleftarrow$$

$$y_n(w^T x_n + b) \geqslant 1$$

What happens when we solve this problem?

$$\frac{\partial L}{\partial w} \Rightarrow w - \sum_{m=1}^{N} \alpha_n y_n x_n = 0 \Rightarrow \underline{\underline{w}} = \sum_{n=1}^{N} \alpha_n y_n x_n \Bigg\}$$

$$\frac{\partial L}{\partial b} \Rightarrow -\sum_{n=1}^{N} \alpha_n y_n = 0 \Rightarrow \sum_{n=1}^{N} \alpha_n y_n = 0$$

$$\frac{1}{2} w^T w$$

$$L(\alpha) = \to \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} \alpha_n \alpha_m y_n y_m \left( x_n^T x_m \right) + \sum \alpha_n$$

$$\to -1 \sum_{n=1}^{N} \alpha_n y_n x_n^T \sum_{m=1}^{N} \alpha_m y_m x_m$$

$$L(\alpha) = \left( -\frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} \alpha_n \alpha_m y_n y_m \left( x_n^T x_m \right) + \sum \alpha_n \right)$$

$$\max \quad L(\alpha)$$
$$\text{s.t} \quad \alpha \geq 0$$
$$\sum_{n=1}^{N} \alpha_n y_n = 0$$

$$\min \quad L(\alpha)$$
$$\text{s.t.} \quad \sum_{n=1}^{N} \alpha_n y_n = 0$$
$$\alpha \geq 0$$

$$\Uparrow$$

$$\underline{\text{DUAL}} \longrightarrow \text{Primal}$$

$$\underline{\text{End Proof}}$$

Solving $\longrightarrow$ $\alpha_n^*$

$$w_n^* = \sum_{n=1}^{N} \alpha_n^* y_n x_n$$

$$b = y_s - w^{*T} x_s \longleftarrow$$

$$\alpha_n \left( 1 - y_n \left( w^T x_n + b \right) \right) = 0$$
$$\text{If } \alpha_n \neq 0, \alpha_s > 0$$
$$y_s \left( w^T x_s + b \right) = 1$$

Summary

$$\min_{u} \quad u^{\top} Q u + p^{\top} u$$

$$\text{s.t.} \quad \underline{A u \geq C}$$

$$\left[ X_S X_S^{\top} \right]_{NM} = \begin{bmatrix} y_1 x_1^{\top} \\ y_2 x_2^{\top} \\ \vdots \\ y_N x_N^{\top} \end{bmatrix} \Bigg|_n \quad \Big[ y_1 x_1^{\top} \cdots y_N x_N^{\top} \Big] \Bigg\|_{\underset{A}{}}$$

$$= y_n x_n^{\top} \ y_m x_m \Big\} G_M$$



$$x = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \quad y = \begin{bmatrix} -1 \\ -1 \\ +1 \\ +1 \end{bmatrix}$$

$$x_S = \begin{bmatrix} 0 & 0 \\ -2 & -2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \quad \text{signed.}$$

$$G_M = X_S X_S^T \Rightarrow Q$$

$$P = -1^T$$

$$A \text{ (constraints)}$$

$$= \begin{bmatrix} y^T \\ -y^T \\ I \end{bmatrix}$$

$$\begin{cases} y^T \alpha = 0 \rightarrow y^T \alpha \geq 0 \\ \qquad\qquad -y^T \alpha \geq 0 \\ \alpha \geq 0 \rightarrow I\alpha \geq 0 \end{cases}$$

$$C = \begin{bmatrix} 0 \\ 0 \\ 0_N \end{bmatrix}$$

Got QPAC $\longrightarrow$ QP Solver $\longrightarrow u^*$

$$\omega^* \quad b^*$$

$$\omega^* = \sum_n \alpha_n^* y_n x_n \qquad \alpha^* = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ 1 \\ 0 \end{bmatrix}$$

$$= \frac{1}{2} \cdot -1 \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \frac{1}{2} -1 \begin{bmatrix} 2 \\ 2 \end{bmatrix} + 1.1 \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$b^* = -1 - 1 \begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} = -1 \Big\}$$

$$g(x) = \text{sign}(x_1 - x_2 - 1)$$

Primal Version $g(x) = \text{sign}(w^{*T}x + b^{*})$

Dual Version: $g(x) = \text{sign}\left(\sum\limits_{n=1}^{N} \alpha_n^{*} y_n \underline{(x_n^T x)} + b^{*}\right)$

where $b^{*} = y_s - \sum\limits_{n=1}^{N} \alpha_n^{*} y_n \left(\underline{x_n^T x_s}\right)$

dot product

$\alpha_n = 0$

$\rightarrow$ Most of $\alpha$'s are zero, $\alpha_s$ (need) $\longrightarrow$ huge saving in computation

$\rightarrow$ Don't need all the data.

$y, \alpha$

$\rightarrow [G]_{nm} \longrightarrow x_n$ dotted into $x_m \longrightarrow \alpha_n^{*} \, x_n x_m$

$$G_{nm} \rightarrow N \times N \text{ matrix of dot products} \} \alpha_n^*$$

$$y \rightarrow N \times 1$$

## INNER PRODUCT ALGORITHMS

1) $\underbrace{K(x, x')}_{\text{Kernel}} \longrightarrow x \cdot x'$

$$\min_{\alpha} \sum_{nm} \alpha_n \alpha_m y_n y_m K(x_n, x_m) - \sum_n \alpha_n$$

$$\text{s.t } \sum_n \alpha_n y_n = 0$$

$$\alpha_n \geq 0$$

$$\downarrow$$

$$\alpha_n^*$$

$$g(x) = \text{sign}\left(\sum_{n=1}^{\alpha} \alpha_r^* \, K(x_{n}, x) + b^*\right)$$

$$b^* = y_s - \sum_{n=1}^{N} \alpha_n y_n \, K(x_n, x_s)$$

$$\alpha_s > 0$$

→ Dot product regardless of the space.

$$x_1 \; x_2 \cdots x_N \longrightarrow z_1, z_2 \cdots z_N$$

Target $y_1 \; y_2 \cdots y_N \longrightarrow y_1 \; y_2 \cdots y_N$ ?

$$\underline{\underline{z_n^T z_m}}$$

Test point $\longrightarrow \underline{\underline{z_n^T z}}$ ?

$$\longrightarrow K(\underbrace{x_n \; x_m}_{}) \xrightarrow{\;\;\checkmark\;\;} z_n^T z_n \;\Big\}$$

$$K(x_n, x) \xrightarrow{\;\;\checkmark\;\;} z_n^T z$$

$$K(x_n \; x_S) \xrightarrow{\;\;\checkmark\;\;} z_n^T z_S$$

## Example 1

① 2nd order polynomial transform.

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \longrightarrow z = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ \sqrt{2}\, x_1 x_2 \\ x_2^2 \end{bmatrix}$$

$$z \cdot z' = \Phi(x)^T \cdot \phi(x')$$

$$z' = \begin{bmatrix} x_1' \\ x_2' \\ x_1'^2 \\ \sqrt{2}\, x_1' x_2' \\ x_2'^2 \end{bmatrix}$$

$$z \cdot z' = x_1 x_1' + x_2 x_2' + x_1^2 (x_1')^2 + x_2^2 (x_2')^2 + 2 x_1 x_2 x_1' x_2'$$

$$= \left( x_1 x_1' + x_2 x_2' + \frac{1}{2} \right)^2 - \frac{1}{4}$$

$$z \cdot z' = \left( x \cdot x' + \frac{1}{2} \right)^2 - \frac{1}{4} \Rightarrow K(x, x')$$

$\downarrow$ x.space

$$K(x, x') = (x \cdot x' + \tfrac{1}{\xi})^Q - C$$

This kernel $\Longrightarrow$ get dot product in Z-space without actually going thru.

1) Kernel $\longrightarrow$ dot products!

2) Huge computational saving.

Separable data

C (penetration)

# Infinite dimensional Space

$$x \longrightarrow e^{-x^2} \begin{pmatrix} \sqrt{\frac{2}{0!}} \, x^0 \\ \sqrt{\frac{2}{1!}} \, x^1 \\ \sqrt{\frac{2^3}{3!}} \, x^3 \\ \vdots \\ \sqrt{\frac{2^k}{k!}} \, x^k \\ \vdots \end{pmatrix} = Z$$

$$x' \longrightarrow e^{-x'^2} \begin{pmatrix} \sqrt{\frac{2}{0!}} \, x'^0 \\ \sqrt{\frac{2}{1!}} \, x'^1 \\ \vdots \\ \sqrt{\frac{2}{k!}} \, x'^k \end{pmatrix}$$

$$Z \cdot Z' = e^{-x^2} e^{-(x')^2} \left[ \frac{2^0}{0!}(x x')^0 + \frac{2^1}{1!}(x^2 x'^1)^1 + \cdots \right]$$

$$\underbrace{\qquad\qquad}_{e^{2xx'}}$$

$$= e^{-(x-x')^2}$$

$$K(x, x') = e^{-(x-x')^2} \implies e^{-r(\|x-x'\|)^2}$$

$\underbrace{\hphantom{K(x, x')}}$ gaussian kernel

Kernel $\longrightarrow$ important.

Any kernel ( <u>symmetric</u>, <u>+ve definite</u> )

<u>Infinite dimensions</u>.

1) Computationally feasible $\rightarrow$ Inf. dim
$K(x, x')$

non-falsifiable
$E_{in}$ VS $E_{out}$

2) Regularize (A lot!)

3) Maximally Regularize $\rightarrow$ Hyperplane.

4) Small no. of SVs.

# Thanks!