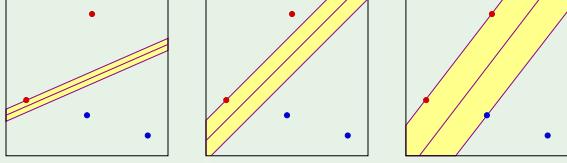


SVM

(Description)

Better linear separation

Linearly separable data



Different separating lines

Which is best?

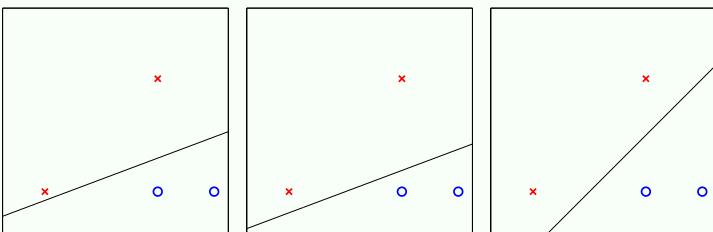
Two questions:

1. Why is bigger margin better?
2. Which \mathbf{w} maximizes the margin?

©  Creator: Yaser Abu-Mostafa • LFD Lecture 14

3/20

Which Separator Do You Pick?



©  Creator: Malik Magdon-Ismail

Maximizing the Margin: 3 /19

Robustness to noise →

Robustness to Noisy Data

Being robust to noise (measurement error) is good (remember regularization).

© Creator: Malik Magdon-Ismail

Maximizing the Margin: 4 /19

Thicker cushion means more robust →

Thicker Cushion Means More Robustness

We call such hyperplanes **fat**

© Creator: Malik Magdon-Ismail

Maximizing the Margin: 5 /19

Two crucial questions →

Two Crucial Questions

1. Can we efficiently find the fattest separating hyperplane?
2. Is a fatter hyperplane better than a thin one?

©  Creator: Malik Magdon-Ismail

Maximizing the Margin: 6 / 19

Pulling out the bias →

Separating The Data

Hyperplane $h = (b, \mathbf{w})$

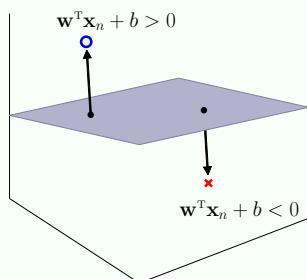
h separates the data means:

$$y_n(\mathbf{w}^T \mathbf{x}_n + b) > 0$$

By rescaling the weights and bias,

$$\min_{n=1,\dots,N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$$

(renormalize the weights so that the signal $\mathbf{w}^T \mathbf{x} + b$ is meaningful)

©  Creator: Malik Magdon-Ismail

Maximizing the Margin: 8 / 19

Distance to the hyperplane →

Distance to the Hyperplane

\mathbf{w} is normal to the hyperplane:

$$\mathbf{w}^T(\mathbf{x}_2 - \mathbf{x}_1) = \mathbf{w}^T\mathbf{x}_2 - \mathbf{w}^T\mathbf{x}_1 = -b + b = 0.$$

(because $\mathbf{w}^T\mathbf{x} = -b$ on the hyperplane)

Unit normal $\mathbf{u} = \mathbf{w}/\|\mathbf{w}\|$.

$$\begin{aligned} \text{dist}(\mathbf{x}, h) &= |\mathbf{u}^T(\mathbf{x} - \mathbf{x}_1)| \\ &= \frac{1}{\|\mathbf{w}\|} \cdot |\mathbf{w}^T\mathbf{x} - \mathbf{w}^T\mathbf{x}_1| \\ &= \frac{1}{\|\mathbf{w}\|} \cdot |\mathbf{w}^T\mathbf{x} + b| \end{aligned}$$

\mathbf{x}_1 is on the plane so $\mathbf{w}^T\mathbf{x}_1+b=0$

© Creator: Malik Magdon-Ismail
 Maximizing the Margin: 9 / 19
 Fatness of a separating hyperplane →

Fatness of a Separating Hyperplane

$$\text{dist}(\mathbf{x}, h) = \frac{1}{\|\mathbf{w}\|} \cdot |\mathbf{w}^T\mathbf{x} + b|$$

Fatness = Distance to the closest point

$\text{Since } |\mathbf{w}^T\mathbf{x}_n + b| = |y_n(\mathbf{w}^T\mathbf{x}_n + b)| = y_n(\mathbf{w}^T\mathbf{x}_n + b),$
 $\text{dist}(\mathbf{x}_n, h) = \frac{1}{\|\mathbf{w}\|} \cdot y_n(\mathbf{w}^T\mathbf{x}_n + b).$

$$\begin{aligned} \text{Fatness} &= \min_n \text{dist}(\mathbf{x}_n, h) \\ &= \frac{1}{\|\mathbf{w}\|} \cdot \min_n y_n(\mathbf{w}^T\mathbf{x}_n + b) && \leftarrow \text{separation condition} \\ &= \frac{1}{\|\mathbf{w}\|} && \leftarrow \text{the margin } \gamma(h) \end{aligned}$$

© Creator: Malik Magdon-Ismail
 Maximizing the Margin: 10 / 19
 Maximizing the margin →

The optimization problem

$$\text{Maximize} \frac{1}{\|\mathbf{w}\|}$$

$$\text{subject to } \min_{n=1,2,\dots,N} |\mathbf{w}^\top \mathbf{x}_n + b| = 1$$

$$\text{Notice: } |\mathbf{w}^\top \mathbf{x}_n + b| = y_n (\mathbf{w}^\top \mathbf{x}_n + b)$$

$$\text{Minimize } \frac{1}{2} \mathbf{w}^\top \mathbf{w}$$

$$\text{subject to } y_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad \text{for } n = 1, 2, \dots, N$$

SVM
(Solution - QP)

Example – Our Toy Data Set

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

$$y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$$

$$\begin{aligned} -b \geq 1 & \quad (i) \\ -(2w_1 + 2w_2 + b) \geq 1 & \quad (ii) \\ 2w_1 + b \geq 1 & \quad (iii) \\ 3w_1 + b \geq 1 & \quad (iv) \end{aligned}$$

(i) and (iii) gives $w_1 \geq 1$
(ii) and (iii) gives $w_2 \leq -1$
So, $\frac{1}{2}(w_1^2 + w_2^2) \geq 1$ ($b = -1, w_1 = 1, w_2 = -1$)

Optimal Hyperplane

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

margin: $\frac{1}{\|\mathbf{w}^*\|} = \frac{1}{\sqrt{2}} \approx 0.707$

For data points (i), (ii) and (iii) $y_n(\mathbf{w}^{*T} \mathbf{x}_n + b^*) = 1$

↑
Support Vectors

© Creator: Malik Magdon-Ismail Maximizing the Margin: 13 / 19 Quadratic programming →

Maximum Margin Hyperplane is QP

$$\underset{\mathbf{b}, \mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

subject to: $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \text{ for } n = 1, \dots, N.$

$$\underset{\mathbf{u} \in \mathbb{R}^q}{\text{minimize}} \quad \frac{1}{2} \mathbf{u}^T \mathbf{Q} \mathbf{u} + \mathbf{c}^T \mathbf{u}$$

subject to: $\mathbf{A} \mathbf{u} \geq \mathbf{a}$

Quadratic Programming

$$\underset{\mathbf{u} \in \mathbb{R}^q}{\text{minimize}} \quad \frac{1}{2} \mathbf{u}^T \mathbf{Q} \mathbf{u} + \mathbf{p}^T \mathbf{u}$$

subject to: $\mathbf{A} \mathbf{u} \geq \mathbf{c}$

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

$$y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$$

$$\begin{aligned} -b \geq 1 & \quad (i) \\ -(2w_1 + 2w_2 + b) \geq 1 & \quad (ii) \\ 2w_1 + b \geq 1 & \quad (iii) \\ 3w_1 + b \geq 1 & \quad (iv) \end{aligned}$$

$$\mathbf{Q} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \mathbf{p} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{A} = \begin{bmatrix} -1 & 0 & 0 \\ -1 & -2 & -2 \\ 1 & 2 & 0 \\ 1 & 3 & 0 \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Use your QP-solver to give
 $(b^*, w_1^*, w_2^*) = (-1, 1, -1)$

Constrained optimization

$$\text{Minimize} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{subject to} \quad y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \text{for } n = 1, 2, \dots, N$$

$$\mathbf{w} \in \mathbb{R}^d, \quad b \in \mathbb{R}$$

Lagrange? inequality constraints \implies KKT

©  Creator: Yaer Abu-Mostafa - LFD Lecture 14

11/20

What is New? Not quite the traditional optimization problem with equality constraints. The constraints here are inequality. KKT to the rescue.

Lagrange relaxation and KKT conditions

We consider the optimization problem

$$\begin{aligned} & \text{minimize} \quad f(\mathbf{x}) \\ & \text{s.t.} \quad g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \end{aligned}$$

where $f(\mathbf{x})$ and $g_i(\mathbf{x})$ are real valued functions.

If $\mathbf{g}(\mathbf{x}) = [g_1(\mathbf{x}) \quad \dots \quad g_m(\mathbf{x})]^T$ then (1) can be written

$$\begin{aligned} & \text{minimize} \quad f(\mathbf{x}) \\ & \text{s.t.} \quad \mathbf{g}(\mathbf{x}) \leq \mathbf{0}. \end{aligned}$$

The idea behind Lagrange relaxation is to put non-negative prices $y_i \geq 0$, on the constraints and then add these to the objective function.

This gives the (unconstrained) optimization problem:

$$\text{minimize} \quad f(\mathbf{x}) + \sum_{i=1}^m y_i g_i(\mathbf{x}) \quad (3)$$

which using $\mathbf{y} = [y_1 \quad \dots \quad y_m]^T$ can be written

$$\text{minimize} \quad f(\mathbf{x}) + \mathbf{y}^T \mathbf{g}(\mathbf{x}) \quad \boxed{L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^T \mathbf{g}(\mathbf{x})}$$

Theorem 1 (Weak duality). For an arbitrary $y \geq 0$ it holds that

$$\min_x L(x, y) \leq f(\hat{x}), \text{ where } \hat{x} \text{ is an optimal solution to (2).}$$

Proof: Since \hat{x} is a feasible solution to (2) it holds that $g(\hat{x}) \leq 0$. We get

$$\min_x L(x, y) \leq L(\hat{x}, y) = f(\hat{x}) + \underbrace{y^T g(\hat{x})}_{\geq 0} \leq f(\hat{x}). \quad (4)$$

- minimizing the Lagrange function provides lower bounds to the optimization problem (2).
- By an appropriate choice of y a good approximation of the optimal solution to (2) is searched for. In practical algorithms one tries to solve $\max_{y \geq 0} \min_x L(x, y)$. The next theorem gives conditions for the Lagrange multiplicator providing equality in (4).

Global optimality conditions

Theorem 2. If $(\hat{x}, \hat{y}) \in \mathbb{R}^n \times \mathbb{R}^m$ satisfies the conditions

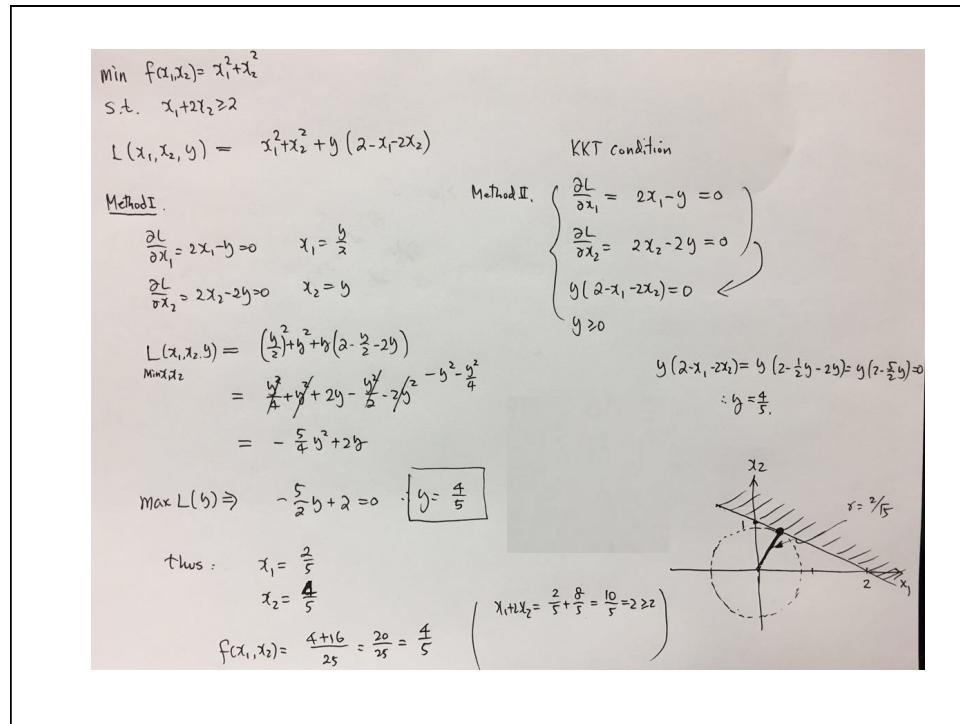
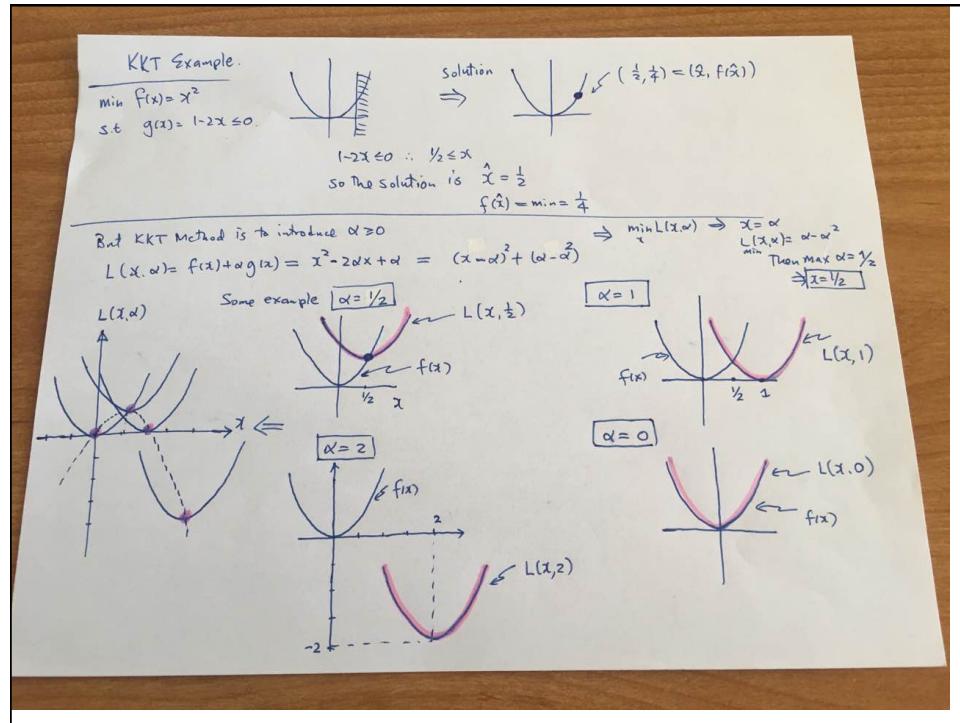
- (1) $L(\hat{x}, \hat{y}) = \min_x L(x, \hat{y})$,
- (2) $g(\hat{x}) \leq 0$,
- (3) $\hat{y} \geq 0$,
- (4) $\hat{y}^T g(\hat{x}) = 0$.

then \hat{x} is an optimal solution to (2).

Proof: If x is an arbitrary feasible solution to (2) it holds that $g(x) \leq 0$, which shows that

$$f(x) \geq f(x) + \hat{y}^T g(x) = L(x, \hat{y}) \geq L(\hat{x}, \hat{y}) = f(\hat{x})$$

where the first inequality follows from (3) and $g(x) \leq 0$, the second inequality follows from (1), and the last one from (4).



Lagrange formulation

$$\text{Minimize } \mathcal{L}(\mathbf{w}, \mathbf{b}, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1)$$

w.r.t. \mathbf{w} and b and maximize w.r.t. each $\alpha_n \geq 0$

$$\begin{aligned}\nabla_{\mathbf{w}} \mathcal{L} &= \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = \mathbf{0} \\ \frac{\partial \mathcal{L}}{\partial b} &= - \sum_{n=1}^N \alpha_n y_n = 0\end{aligned}$$

©  Creator: Yaser Abu-Mostafa - LFD Lecture 14

13/20

Substituting ...

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \quad \text{and} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

in the Lagrangian $\mathcal{L}(\mathbf{w}, \mathbf{b}, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1)$

we get $\mathcal{L}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{x}_n^T \mathbf{x}_m$

Maximize w.r.t. to $\boldsymbol{\alpha}$ subject to $\alpha_n \geq 0$ for $n = 1, \dots, N$ and $\sum_{n=1}^N \alpha_n y_n = 0$

©  Creator: Yaser Abu-Mostafa - LFD Lecture 14

14/20

The solution - quadratic programming

$$\min_{\alpha} \quad \frac{1}{2} \alpha^T \underbrace{\begin{bmatrix} y_1 y_1 \mathbf{x}_1^T \mathbf{x}_1 & y_1 y_2 \mathbf{x}_1^T \mathbf{x}_2 & \dots & y_1 y_N \mathbf{x}_1^T \mathbf{x}_N \\ y_2 y_1 \mathbf{x}_2^T \mathbf{x}_1 & y_2 y_2 \mathbf{x}_2^T \mathbf{x}_2 & \dots & y_2 y_N \mathbf{x}_2^T \mathbf{x}_N \\ \dots & \dots & \dots & \dots \\ y_N y_1 \mathbf{x}_N^T \mathbf{x}_1 & y_N y_2 \mathbf{x}_N^T \mathbf{x}_2 & \dots & y_N y_N \mathbf{x}_N^T \mathbf{x}_N \end{bmatrix}}_{\text{quadratic coefficients}} \alpha + \underbrace{(-\mathbf{1}^T) \alpha}_{\text{linear}}$$

subject to $\underbrace{\mathbf{y}^T \alpha = 0}_{\text{linear constraint}}$

$$\underbrace{\mathbf{0}}_{\text{lower bounds}} \leq \alpha \leq \underbrace{\infty}_{\text{upper bounds}}$$

© Creator: Yaser Abu-Mostafa - LFD Lecture 14

15/20

Still has to call a QP solver, but the problem is must simplified with a simple constraint.

QP hands us α

Solution: $\alpha = \alpha_1, \dots, \alpha_N$

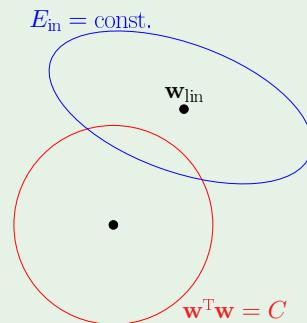
$$\Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

KKT condition: For $n = 1, \dots, N$

$$\alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1) = 0$$

We saw this before!

$\alpha_n > 0 \Rightarrow \mathbf{x}_n$ is a **support vector**



© Creator: Yaser Abu-Mostafa - LFD Lecture 14

16/20

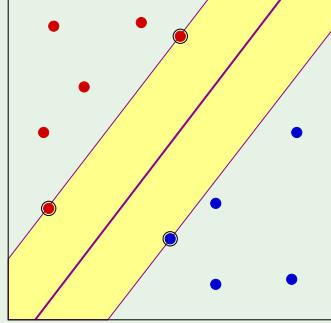
Support vectors

Closest \mathbf{x}_n 's to the plane: achieve the margin

$$\Rightarrow y_n (\mathbf{w}^\top \mathbf{x}_n + b) = 1$$

$$\mathbf{w} = \sum_{\mathbf{x}_n \text{ is SV}} \alpha_n y_n \mathbf{x}_n$$

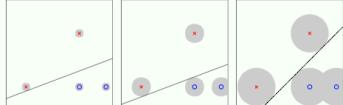
Solve for b using any SV:

$$y_n (\mathbf{w}^\top \mathbf{x}_n + b) = 1$$


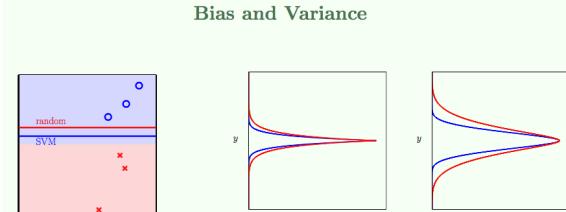
© Creator: Yaser Abu-Mostafa - LFD Lecture 14

17/20

Evidence that Larger Margin is Better



Bias and Variance



| | Random | SVM |
|-----------|--------|-------|
| bias | 0.02 | 0.015 |
| var | 0.059 | 0.038 |
| E_{out} | 0.079 | 0.053 |

-0.005

-0.021

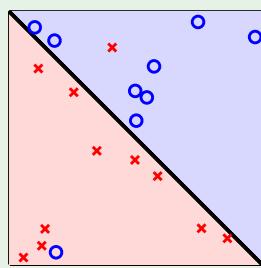
-0.026

13

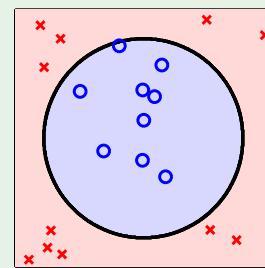
SVM (Extensions)

Two types of non-separable

slightly:



seriously:



SVM

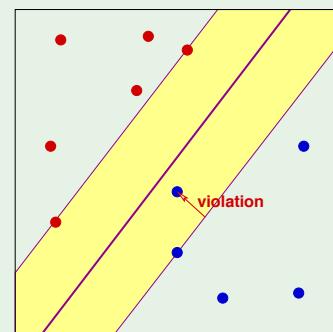
(Soft)

Error measure

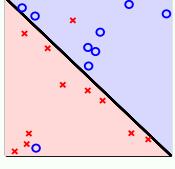
Margin violation: $y_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1$ fails

Quantify: $y_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n \quad \xi_n \geq 0$

$$\text{Total violation} = \sum_{n=1}^N \xi_n$$



Soft Margin SVM



tolerate error

$$\begin{aligned} & \underset{b, \mathbf{w}, \boldsymbol{\xi}}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ & \text{subject to: } y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \\ & \quad \xi_n \geq 0 \quad \text{for } n = 1, \dots, N \end{aligned}$$

Trades off 'soft in-sample error' $\sum_{n=1}^N \xi_n$ with weight norm $\frac{1}{2} \mathbf{w}^T \mathbf{w}$

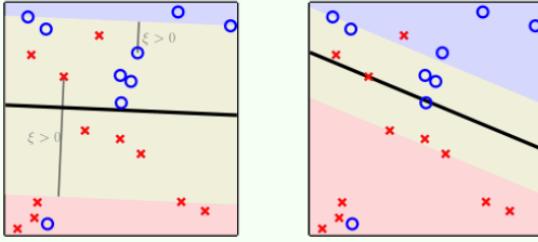
C plays the role of a regularization parameter ($\lambda \sim \frac{1}{C}$)

Choice of C is important - similar to choice of λ in regularization

← regularization

©  Creator: Malik Magdon-Ismail Overfitting and the Optimal Hyperplane: 13 / 17 Non-Separable Data →

Non-Separable Data



$C = 1$ $C = 500$

$$\begin{aligned} & \underset{b, \mathbf{w}, \boldsymbol{\xi}}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ & \text{subject to: } y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \\ & \quad \xi_n \geq 0 \quad \text{for } n = 1, \dots, N \end{aligned}$$

The new optimization

$$\text{Minimize} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n$$

$$\text{subject to} \quad y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \quad \text{for } n = 1, \dots, N$$

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

$$\mathbf{w} \in \mathbb{R}^d, \quad b \in \mathbb{R}, \quad \xi \in \mathbb{R}^N$$

©  Creator: Yaser Abu-Mostafa - LFD Lecture 15

16/20

Lagrange formulation

$$\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1 + \xi_n) - \sum_{n=1}^N \beta_n \xi_n$$

Minimize w.r.t. \mathbf{w} , b , and ξ and maximize w.r.t. each $\alpha_n \geq 0$ and $\beta_n \geq 0$

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = \mathbf{0}$$

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

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = C - \alpha_n - \beta_n = 0$$

©  Creator: Yaser Abu-Mostafa - LFD Lecture 15

17/20

and the solution is ...

$$\begin{aligned} \text{Maximize } \quad \mathcal{L}(\boldsymbol{\alpha}) &= \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{x}_n^\top \mathbf{x}_m \quad \text{w.r.t. to } \boldsymbol{\alpha} \\ \text{subject to } \quad 0 \leq \alpha_n \leq C \quad \text{for } n = 1, \dots, N \quad \text{and} \quad \sum_{n=1}^N \alpha_n y_n &= 0 \\ \implies \mathbf{w} &= \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \\ \text{minimizes } \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{n=1}^N \xi_n & \end{aligned}$$

© Creator: Yaser Abu-Mostafa - LFD Lecture 15

18/20

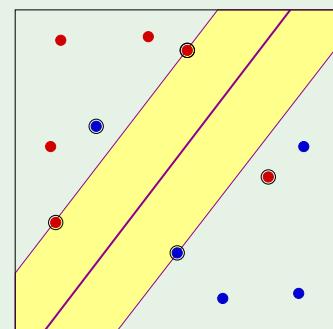
Types of support vectors

margin support vectors $(0 < \alpha_n < C)$

$$y_n (\mathbf{w}^\top \mathbf{x}_n + b) = 1 \quad (\xi_n = 0)$$

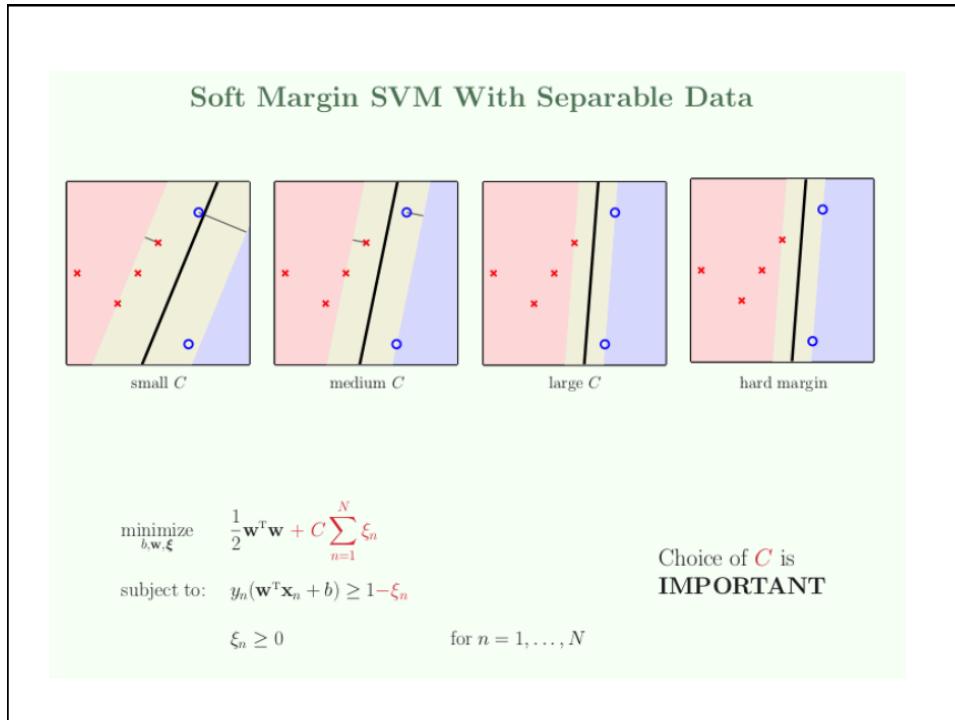
non-margin support vectors $(\alpha_n = C)$

$$y_n (\mathbf{w}^\top \mathbf{x}_n + b) < 1 \quad (\xi_n > 0)$$



© Creator: Yaser Abu-Mostafa - LFD Lecture 15

19/20



Two technical observations

1. **Hard margin:** What if data is not linearly separable?
“primal \longrightarrow dual” breaks down
2. **\mathcal{Z} :** What if there is w_0 ?
All goes to b and $w_0 \rightarrow 0$

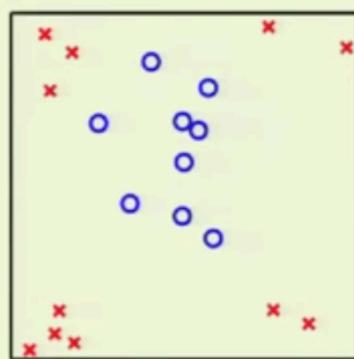
© Creator: Yaser Abu-Mostafa - LFD Lecture 15 20/20

SVM

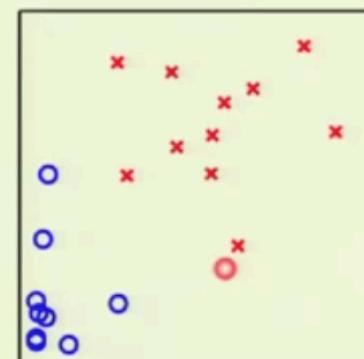
(Extension: Kernel)

\mathbf{z} instead of \mathbf{x}

$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{z}_n^\top \mathbf{z}_m$$



$\mathcal{X} \rightarrow \mathcal{Z}$



What do we need from the \mathcal{Z} space?

$$\mathcal{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{z}_n^\top \mathbf{z}_m$$

Constraints: $\alpha_n \geq 0$ for $n = 1, \dots, N$ and $\sum_{n=1}^N \alpha_n y_n = 0$

$$g(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{z} + b) \quad \text{need } \mathbf{z}_n^\top \mathbf{z}$$

$$\text{where } \mathbf{w} = \sum_{\mathbf{z}_n \text{ is SV}} \alpha_n y_n \mathbf{z}_n$$

$$\text{and } b: y_m (\mathbf{w}^\top \mathbf{z}_m + b) = 1 \quad \text{need } \mathbf{z}_n^\top \mathbf{z}_m$$

Generalized inner product

Given two points \mathbf{x} and $\mathbf{x}' \in \mathcal{X}$, we need $\mathbf{z}^\top \mathbf{z}'$

Let $\mathbf{z}^\top \mathbf{z}' = K(\mathbf{x}, \mathbf{x}')$ ([the kernel](#)) "inner product" of \mathbf{x} and \mathbf{x}'

Example: $\mathbf{x} = (x_1, x_2) \longrightarrow$ 2nd-order Φ

$$\mathbf{z} = \Phi(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$$

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= \mathbf{z}^\top \mathbf{z}' = 1 + x_1 x'_1 + x_2 x'_2 + \\ &\quad x_1^2 x'^2_1 + x_2^2 x'^2_2 + x_1 x'_1 x_2 x'_2 \end{aligned}$$

The trick

Can we compute $K(\mathbf{x}, \mathbf{x}')$ **without** transforming \mathbf{x} and \mathbf{x}' ?

Example: Consider $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^2 = (1 + x_1 x'_1 + x_2 x'_2)^2$

$$= 1 + x_1^2 x_1'^2 + x_2^2 x_2'^2 + \cancel{2x_1 x_1'} + \cancel{2x_2 x_2'} + \cancel{2x_1 x_1' x_2 x_2'}$$

This is an inner product!

$$(1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2)$$

$$(1, x_1'^2, x_2'^2, \sqrt{2}x_1', \sqrt{2}x_2', \sqrt{2}x_1'x_2')$$

© Creator: Yaser Abu-Mostafa - LFD Lecture 15

5/20

The Kernel $K(\cdot, \cdot)$ for a Transform $\Phi(\cdot)$

The Kernel tells you how to compute the inner product in \mathcal{Z} -space

$$K(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}') = \mathbf{z}^T \mathbf{z}'$$

Example: 2nd-order polynomial transform

$$\Phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}$$

$$K(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}') = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} x'_1 \\ x'_2 \\ x_1'^2 \\ \sqrt{2}x'_1x'_2 \\ x_2'^2 \end{bmatrix} \leftarrow O(d^2)$$

$$= x_1x'_1 + x_2x'_2 + x_1^2x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2x_2'^2$$

$$= \left(\frac{1}{2} + \mathbf{x}^T \mathbf{x}' \right)^2 - \frac{1}{4}$$

\uparrow
 computed quickly
 in d -space, in $O(d)$

© Creator: Malik Magdon-Ismail

Kernel Trick: 15 /18

Gaussian kernel →

The polynomial kernel

$\mathcal{X} = \mathbb{R}^d$ and $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$ is polynomial of order Q

The "equivalent" kernel $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^Q$

$$= (1 + x_1 x'_1 + x_2 x'_2 + \dots + x_d x'_d)^Q$$

Compare for $d = 10$ and $Q = 100$

Can adjust scale: $K(\mathbf{x}, \mathbf{x}') = (a \mathbf{x}^\top \mathbf{x}' + b)^Q$

We only need \mathcal{Z} to exist!

If $K(\mathbf{x}, \mathbf{x}')$ is an inner product in some space \mathcal{Z} , we are good.

Example: $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$

Infinite-dimensional \mathcal{Z} : take simple case

$$K(x, x') = \exp(-(x - x')^2)$$

$$= \exp(-x^2) \exp(-x'^2) \underbrace{\sum_{k=0}^{\infty} \frac{2^k (x)^k (x')^k}{k!}}_{\exp(2xx')}$$

The Gaussian Kernel is Infinite-Dimensional

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

Example: Gaussian Kernel in 1-dimension

$$\Phi(\mathbf{x}) = \begin{bmatrix} e^{-x^2} \sqrt{\frac{2^0}{0!}} \\ e^{-x^2} \sqrt{\frac{2^1}{1!}} x \\ e^{-x^2} \sqrt{\frac{2^2}{2!}} x^2 \\ e^{-x^2} \sqrt{\frac{2^3}{3!}} x^3 \\ e^{-x^2} \sqrt{\frac{2^4}{4!}} x^4 \\ \vdots \end{bmatrix} \quad K(x, x') = \Phi(x)^T \Phi(x') = \begin{bmatrix} e^{-x'^2} \sqrt{\frac{2^0}{0!}} \\ e^{-x'^2} \sqrt{\frac{2^1}{1!}} x' \\ e^{-x'^2} \sqrt{\frac{2^2}{2!}} x'^2 \\ e^{-x'^2} \sqrt{\frac{2^3}{3!}} x'^3 \\ e^{-x'^2} \sqrt{\frac{2^4}{4!}} x'^4 \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} e^{-x^2} \sqrt{\frac{2^0}{0!}} \\ e^{-x^2} \sqrt{\frac{2^1}{1!}} x \\ e^{-x^2} \sqrt{\frac{2^2}{2!}} x^2 \\ e^{-x^2} \sqrt{\frac{2^3}{3!}} x^3 \\ e^{-x^2} \sqrt{\frac{2^4}{4!}} x^4 \\ \vdots \end{bmatrix}$$

$$= e^{-x^2} e^{-x'^2} \sum_{i=0}^{\infty} \frac{(2xx')^i}{i!}$$

(infinite dimensional Φ)

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

© Creator: Malik Magdon-Ismail

Kernel Trick: 16 /18

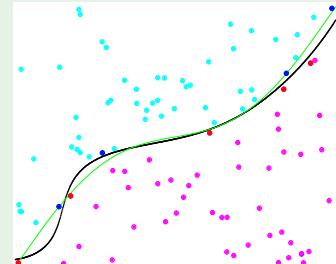
Bypass Z-space →

This kernel in action

Slightly non-separable case:

Transforming \mathcal{X} into ∞ -dimensional \mathcal{Z}

Overkill? Count the support vectors



© Creator: Yaser Abu-Mostafa - LFD Lecture 15

8/20

Kernel formulation of SVM

Remember quadratic programming? The only difference now is:

$$\begin{bmatrix} y_1y_1 K(\mathbf{x}_1, \mathbf{x}_1) & y_1y_2 K(\mathbf{x}_1, \mathbf{x}_2) & \dots & y_1y_N K(\mathbf{x}_1, \mathbf{x}_N) \\ y_2y_1 K(\mathbf{x}_2, \mathbf{x}_1) & y_2y_2 K(\mathbf{x}_2, \mathbf{x}_2) & \dots & y_2y_N K(\mathbf{x}_2, \mathbf{x}_N) \\ \dots & \dots & \dots & \dots \\ y_Ny_1 K(\mathbf{x}_N, \mathbf{x}_1) & y_Ny_2 K(\mathbf{x}_N, \mathbf{x}_2) & \dots & y_Ny_N K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

quadratic coefficients

Everything else is the same.

The final hypothesis

Express $g(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{z} + b)$ in terms of $K(-, -)$

$$\mathbf{w} = \sum_{\mathbf{z}_n \text{ is SV}} \alpha_n y_n \mathbf{z}_n \implies g(\mathbf{x}) = \text{sign} \left(\sum_{\alpha_n > 0} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b \right)$$

$$\text{where } b = y_m - \sum_{\alpha_n > 0} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}_m)$$

for any support vector ($\alpha_m > 0$)

How do we know that \mathcal{Z} exists ...

... for a given $K(\mathbf{x}, \mathbf{x}')$? valid kernel

Three approaches:

1. By construction
2. Math properties (*Mercer's condition*)
3. Who cares? ☺

Design your own kernel

$K(\mathbf{x}, \mathbf{x}')$ is a valid kernel iff

1. It is symmetric and 2. The matrix:
$$\begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \dots & K(\mathbf{x}_1, \mathbf{x}_N) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \dots & K(\mathbf{x}_2, \mathbf{x}_N) \\ \dots & \dots & \dots & \dots \\ K(\mathbf{x}_N, \mathbf{x}_1) & K(\mathbf{x}_N, \mathbf{x}_2) & \dots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

is **positive semi-definite**

for any $\mathbf{x}_1, \dots, \mathbf{x}_N$ (*Mercer's condition*)

The Kernel Allows Us to Bypass \mathcal{Z} -space

$x_n \in \mathcal{X}$

$\downarrow K(\cdot, \cdot)$

$g(\mathbf{x}) = \text{sign} \left(\sum_{\alpha_n^* > 0} \alpha_n^* y_n K(\mathbf{x}_n, \mathbf{x}) + b^* \right)$

$b^* = y_s - \sum_{\alpha_n^* > 0} \alpha_n^* y_n K(\mathbf{x}_n, \mathbf{x}_s)$

1: Input: \mathbf{X}, \mathbf{y} , regularization parameter C

2: Compute G: $G_{nm} = y_n y_m K(\mathbf{x}_n, \mathbf{x}_m)$.

3: Solve (QP):

$$\begin{array}{l} \text{minimize}_{\boldsymbol{\alpha}}: \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\alpha} - \mathbf{1}^T \boldsymbol{\alpha} \\ \text{subject to: } \mathbf{y}^T \boldsymbol{\alpha} = 0 \\ \mathbf{C} \geq \boldsymbol{\alpha} \geq \mathbf{0} \end{array} \rightarrow \boldsymbol{\alpha}^*$$

4: The final hypothesis is

$$g(\mathbf{x}) = \text{sign} \left(\sum_{\alpha_n^* > 0} \alpha_n^* y_n K(\mathbf{x}_n, \mathbf{x}) + b^* \right)$$

The Kernel-Support Vector Machine

| | |
|--|--|
| <u>Overfitting</u> <p>SVM Regression</p> <p>high $\tilde{d} \rightarrow$ complicated separator</p> <p>small # support vectors \rightarrow low effective complexity</p> <p>Can go to high (infinite) \tilde{d}</p> | <u>Computation</u> <p>Inner products with Kernel $K(\cdot, \cdot)$</p> <p>high $\tilde{d} \rightarrow$ expensive or infeasible computation</p> <p>kernel \rightarrow computationally feasible to go to high \tilde{d}</p> <p>Can go to high (infinite) \tilde{d}</p> |
|--|--|

© Creator: Malik Magdon-Ismail

Kernel Trick: 18 / 18

SVM

(Radial Basis Functions)

Learning From Data

Yaser S. Abu-Mostafa
California Institute of Technology

Lecture 16: Radial Basis Functions



Sponsored by Caltech's Provost Office, E&AS Division, and IST • Thursday, May 24, 2012



Outline

- RBF and nearest neighbors
- RBF and neural networks
- RBF and kernel methods
- RBF and regularization

 Creator: Yaser Abu-Mostafa - LFD Lecture 16

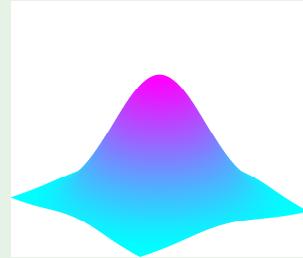
2/20

Basic RBF model

Each $(\mathbf{x}_n, y_n) \in \mathcal{D}$ influences $h(\mathbf{x})$ based on $\underbrace{\|\mathbf{x} - \mathbf{x}_n\|}_{\text{radial}}$

Standard form:

$$h(\mathbf{x}) = \sum_{n=1}^N w_n \underbrace{\exp(-\gamma \|\mathbf{x} - \mathbf{x}_n\|^2)}_{\text{basis function}}$$



 Creator: Yaser Abu-Mostafa - LFD Lecture 16

3/20

The learning algorithm

$$\text{Finding } w_1, \dots, w_N: \quad h(\mathbf{x}) = \sum_{n=1}^N w_n \exp\left(-\gamma \|\mathbf{x} - \mathbf{x}_n\|^2\right)$$

based on $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$

$E_{\text{in}} = 0: \quad h(\mathbf{x}_n) = y_n \text{ for } n = 1, \dots, N:$

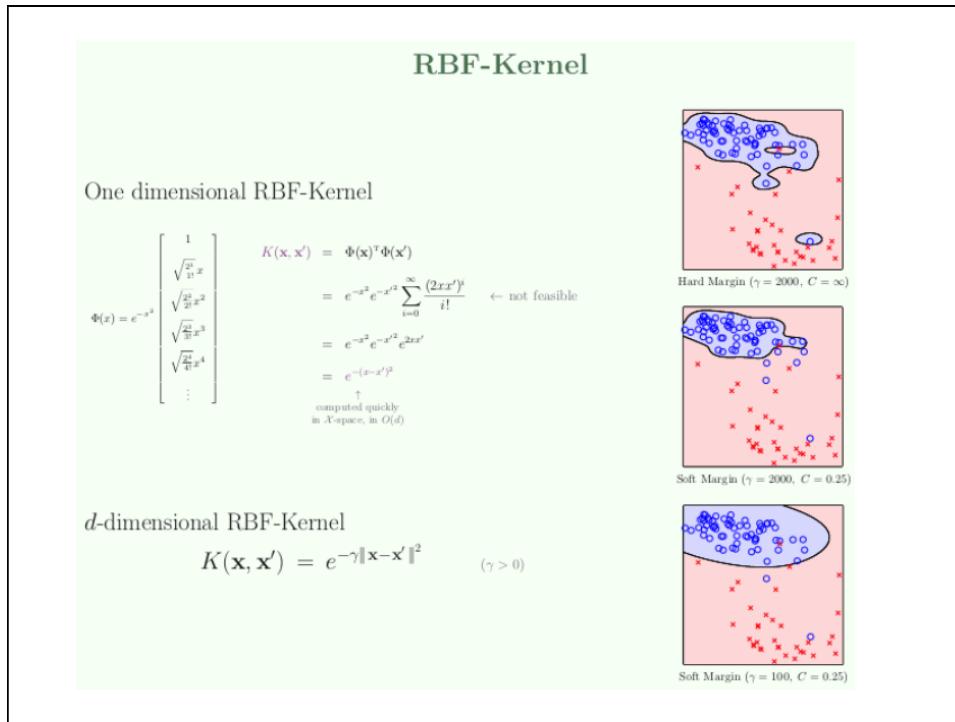
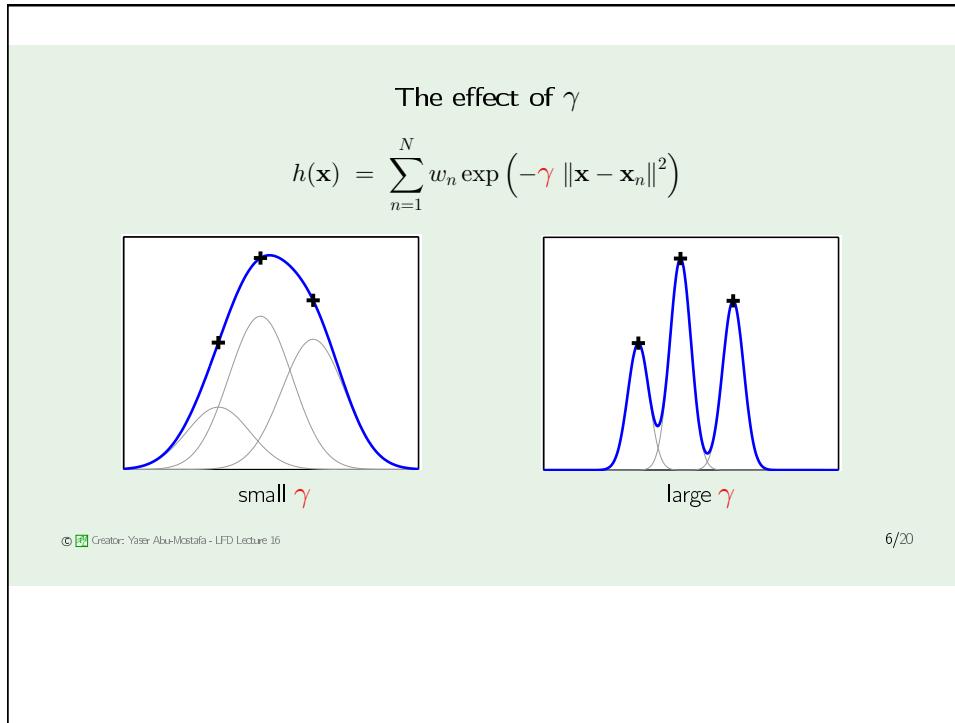
$$\sum_{m=1}^N w_m \exp\left(-\gamma \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) = y_n$$

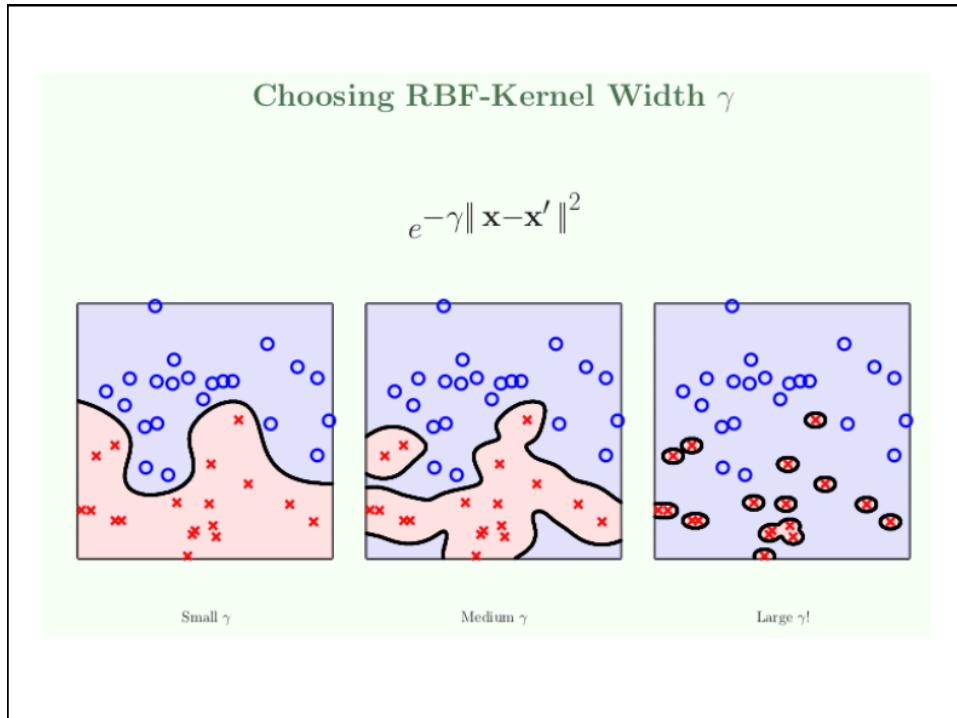
The solution

$$\sum_{m=1}^N w_m \exp\left(-\gamma \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) = y_n \quad N \text{ equations in } N \text{ unknowns}$$

$$\underbrace{\begin{bmatrix} \exp(-\gamma \|\mathbf{x}_1 - \mathbf{x}_1\|^2) & \dots & \exp(-\gamma \|\mathbf{x}_1 - \mathbf{x}_N\|^2) \\ \exp(-\gamma \|\mathbf{x}_2 - \mathbf{x}_1\|^2) & \dots & \exp(-\gamma \|\mathbf{x}_2 - \mathbf{x}_N\|^2) \\ \vdots & \vdots & \vdots \\ \exp(-\gamma \|\mathbf{x}_N - \mathbf{x}_1\|^2) & \dots & \exp(-\gamma \|\mathbf{x}_N - \mathbf{x}_N\|^2) \end{bmatrix}}_{\Phi} \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix}}_{\mathbf{w}} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\mathbf{y}}$$

If Φ is invertible, $\boxed{\mathbf{w} = \Phi^{-1}\mathbf{y}}$ “exact interpolation”





RBF for classification

$$h(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N w_n \exp \left(-\gamma \|\mathbf{x} - \mathbf{x}_n\|^2 \right) \right)$$

Learning: \sim linear regression for classification

$$s = \sum_{n=1}^N w_n \exp \left(-\gamma \|\mathbf{x} - \mathbf{x}_n\|^2 \right)$$

Minimize $(s - y)^2$ on \mathcal{D} $y = \pm 1$

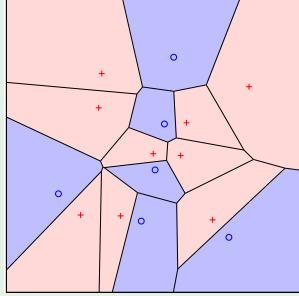
$$h(\mathbf{x}) = \text{sign}(s)$$

© Creator: Yaser Abu-Mostafa - LFD Lecture 16

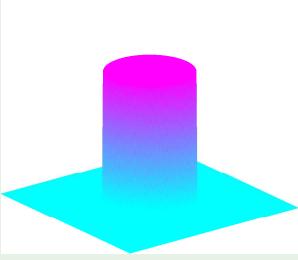
7/20

Relationship to nearest-neighbor method

Adopt the y value of a nearby point:



similar effect by a basis function:



© Creator: Yaes Abu-Mostafa - LFD Lecture 16

8/20

RBF with K centers

N parameters w_1, \dots, w_N based on N data points

Use $K \ll N$ centers: μ_1, \dots, μ_K instead of x_1, \dots, x_N

$$h(\mathbf{x}) = \sum_{k=1}^K w_k \exp\left(-\gamma \|\mathbf{x} - \mu_k\|^2\right)$$

1. How to choose the centers μ_k
2. How to choose the weights w_k

© Creator: Yaes Abu-Mostafa - LFD Lecture 16

9/20

Choosing the centers

Minimize the distance between \mathbf{x}_n and the **closest** center $\boldsymbol{\mu}_k$: *K*-means clustering

Split $\mathbf{x}_1, \dots, \mathbf{x}_N$ into clusters S_1, \dots, S_K

$$\text{Minimize} \sum_{k=1}^K \sum_{\mathbf{x}_n \in S_k} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

Unsupervised learning ⊕

NP-hard ⊗

©  Creator: Yaeser Abu-Mostafa - LFD Lecture 16

10/20

An **iterative** algorithm

Lloyd's algorithm: Iteratively minimize $\sum_{k=1}^K \sum_{\mathbf{x}_n \in S_k} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$ w.r.t. $\boldsymbol{\mu}_k, S_k$

$$\boldsymbol{\mu}_k \leftarrow \frac{1}{|S_k|} \sum_{\mathbf{x}_n \in S_k} \mathbf{x}_n$$

$$S_k \leftarrow \{\mathbf{x}_n : \|\mathbf{x}_n - \boldsymbol{\mu}_k\| \leq \text{all } \|\mathbf{x}_n - \boldsymbol{\mu}_\ell\|\}$$

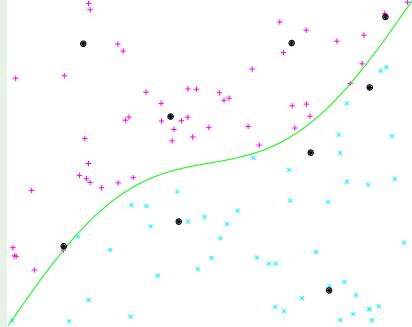
Convergence \longrightarrow **local minimum**

©  Creator: Yaeser Abu-Mostafa - LFD Lecture 16

11/20

Lloyd's algorithm in action

1. Get the data points
2. Only the inputs!
3. Initialize the centers
4. Iterate
5. These are your μ_k 's

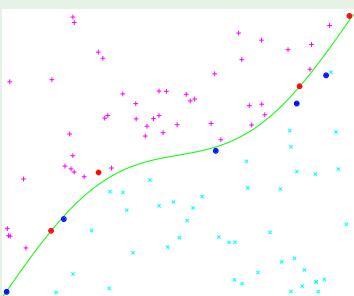


©  Creator: Yaser Abu-Mostafa • LFD Lecture 16

12/20

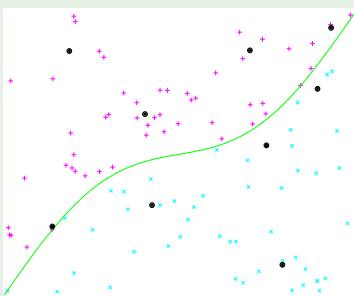
Centers versus support vectors

support vectors



©  Creator: Yaser Abu-Mostafa • LFD Lecture 16

RBF centers



13/20

Choosing the weights

$$\sum_{k=1}^K w_k \exp(-\gamma \|x_n - \mu_k\|^2) \approx y_n \quad N \text{ equations in } K < N \text{ unknowns}$$

$$\underbrace{\begin{bmatrix} \exp(-\gamma \|x_1 - \mu_1\|^2) & \dots & \exp(-\gamma \|x_1 - \mu_K\|^2) \\ \exp(-\gamma \|x_2 - \mu_1\|^2) & \dots & \exp(-\gamma \|x_2 - \mu_K\|^2) \\ \vdots & \vdots & \vdots \\ \exp(-\gamma \|x_N - \mu_1\|^2) & \dots & \exp(-\gamma \|x_N - \mu_K\|^2) \end{bmatrix}}_{\Phi} \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix}}_{w} \approx \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{y}$$

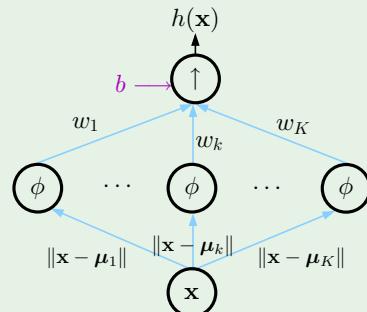
If $\Phi^T \Phi$ is invertible, $w = (\Phi^T \Phi)^{-1} \Phi^T y$ pseudo-inverse

RBF network

The “features” are $\exp(-\gamma \|x - \mu_k\|^2)$

Nonlinear transform depends on \mathcal{D}

\Rightarrow No longer a linear model



A bias term (b or w_0) is often added

