

[illegible]

# 01. Support Vector Machines

## Winter 2020/2021

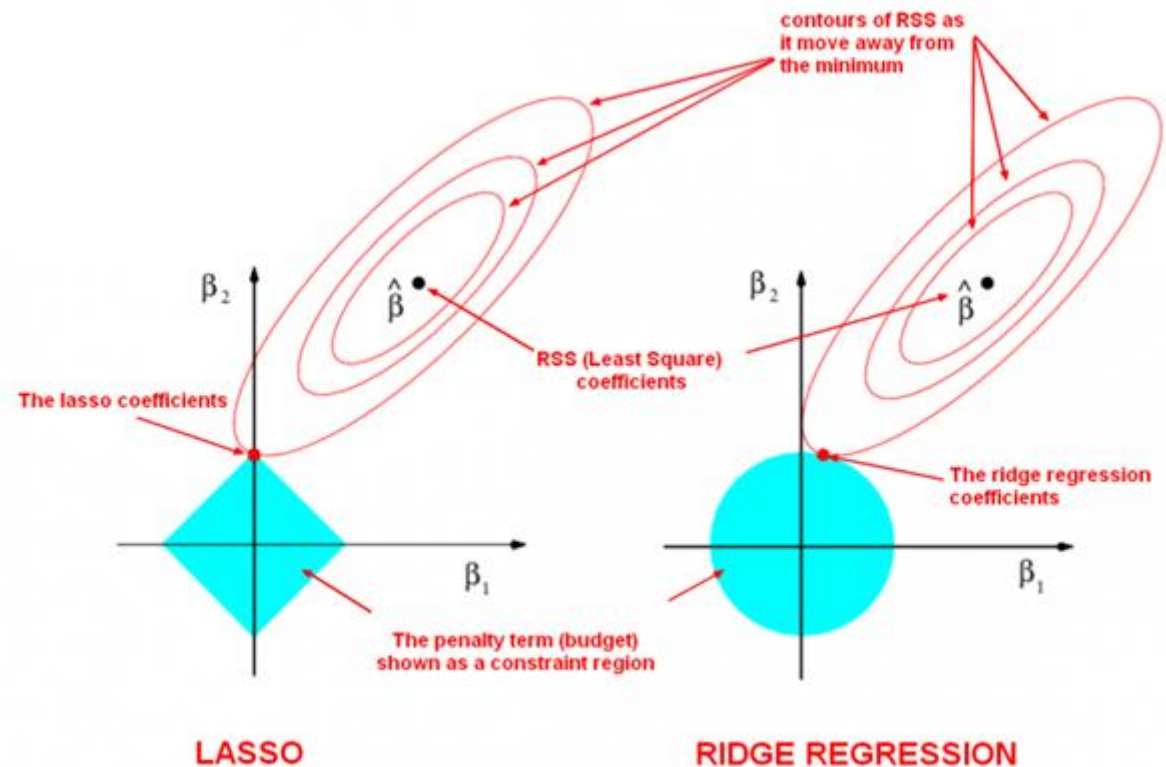


## Recap: regularized regression

—In LASSO regression we use

$$E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} |\mathbf{w}|$$

—LASSO regression can be used for variable selection



Source: <https://www.quora.com/How-would-you-describe-the-difference-between-linear-regression-lasso-regression-and-ridge-regression>



# Agenda

- Introduction
- Learning problem & linear classification
- Linear models: regression & logistic regression
- Non-linear transformation, overfitting & regularization
- **Support Vector Machines and kernel learning**
- Neural Networks: shallow [and deep]
- Theoretical foundation of supervised learning
- Unsupervised learning

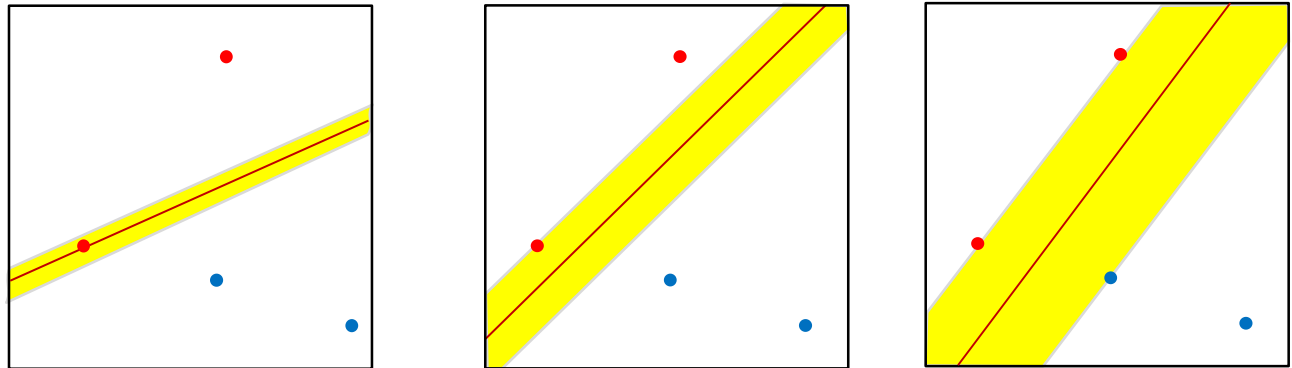


# Better linear separation

Linearly separable data

Different separating lines

What means best?



Two questions:

- Why is bigger margin better?
- Which **w** maximizes the margin?



## Finding $\mathbf{w}$ with large margin

Let  $\mathbf{x}_n$  be the nearest data point to the plane  $\mathbf{w}^\top \mathbf{x} = 0$ . How far is it?

2 preliminary technicalities:

– Normalize  $\mathbf{w}$ :

$$|\mathbf{w}^\top \mathbf{x}_n| = 1$$

– Pull out  $w_0$ :

$$\mathbf{w} = (w_1, \dots, w_d) \text{ as } b$$

– The plane is now defined by  $\mathbf{w}^\top \mathbf{x} + b = 0$   
(no  $x_0$ )



## Computing the distance

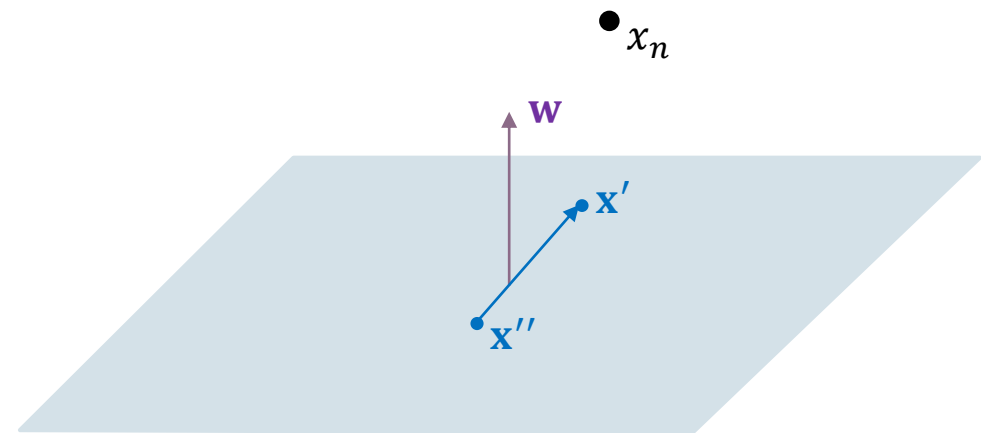
The distance between  $x_n$  and the plane  $\mathbf{w}^\top \mathbf{x} + b = 0$  where  $|\mathbf{w}^\top x_n + b| = 1$

The vector  $\mathbf{w}$  is  $\perp$  to the plane in the  $\mathcal{X}$  space:

Take  $\mathbf{x}'$  and  $\mathbf{x}''$  on the plane

$$\mathbf{w}^\top \mathbf{x}' + b = 0 \quad \text{and} \quad \mathbf{w}^\top \mathbf{x}'' + b = 0$$

$$\Rightarrow \mathbf{w}^\top (\mathbf{x}' - \mathbf{x}'') = 0$$





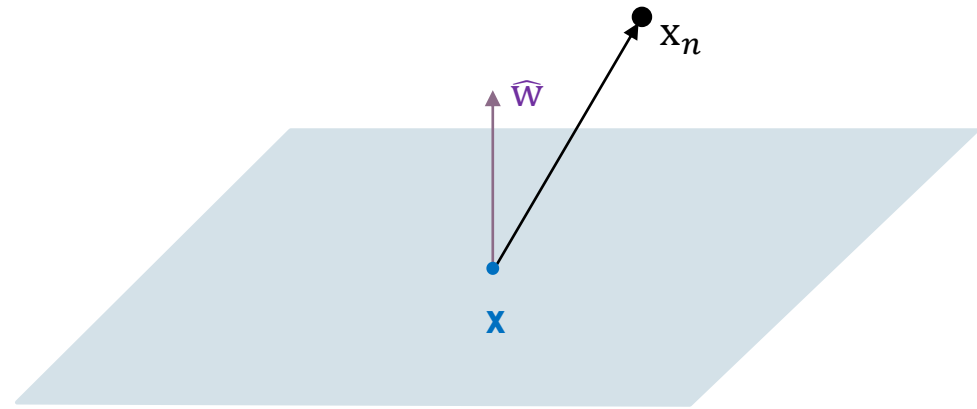
and the distance is ...

Distance between  $x_n$  and the plane:

Take any point  $x$  on the plane

Projection  $x_n - x$  on  $\mathbf{w}$

$$\hat{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|} \Rightarrow \text{distance} = |\hat{\mathbf{w}}^\top (x_n - x)|$$



$$\text{distance} = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^\top x_n - \mathbf{w}^\top x| = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^\top x_n + b - \mathbf{w}^\top x - b| = \frac{1}{\|\mathbf{w}\|}$$



# The optimization problem

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

$$\text{subject to } \min_n |\mathbf{w}^\top \mathbf{x}_n + b| = 1$$

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

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

$$\text{subject to } y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1$$

for  $n = 1, 2, \dots, N$

How can we solve this optimization problem?





## Lagrange formulation (using Karush-Kuhn-Tucker)

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

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

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

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



## Substituting ...

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

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

we get

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

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



## The solution – quadratic programming

$$\min_{\alpha} \frac{1}{2} \underbrace{\alpha^T \begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 & \cdots & y_1 y_N x_1^T x_N \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 & \cdots & y_2 y_N x_2^T x_N \\ \vdots & \vdots & \ddots & \vdots \\ y_N y_1 x_N^T x_1 & y_N y_2 x_N^T x_2 & \cdots & y_N y_N x_N^T x_N \end{bmatrix} \alpha}_{\text{quadratic coefficients}} + \underbrace{(-1^T) \alpha}_{\text{linear}}$$

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

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



## Quadratic Programming finds the $\alpha$ 's

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

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

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

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

That leads to the conclusion

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



# Support vectors

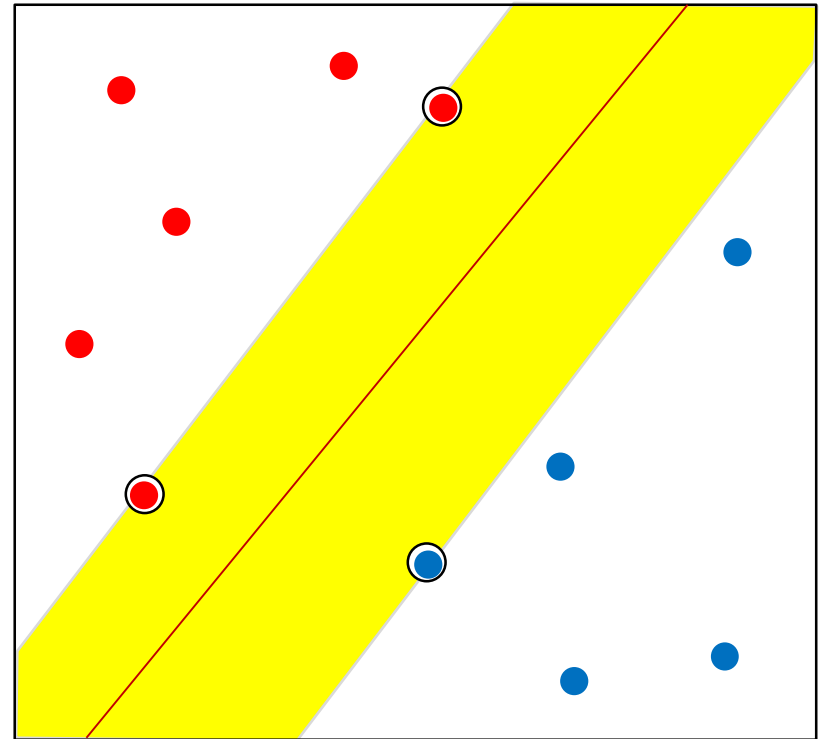
Closest  $\mathbf{x}_n$ 's to the plane: achieve the margin

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

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

Solve for  $b$  using any SV:

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





## An interesting insight

— If we can express  $\mathbf{w}$  in terms of a linear combinations of  $\mathbf{x}_n$ :

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

— ... then we can express the decision function as

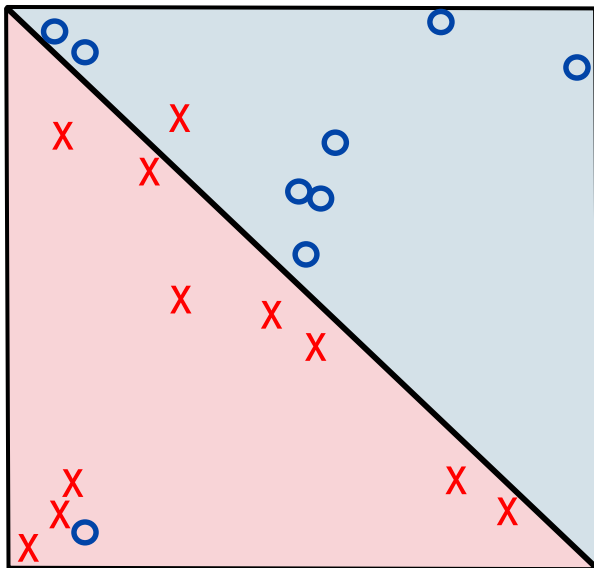
$$h(\mathbf{x}) = \text{sign} \left( \sum_{\mathbf{x}_n \text{ is SV}} \alpha_n y_n \mathbf{x}_n^T \mathbf{x} + b \right)$$

... so it also depends on the linear combination

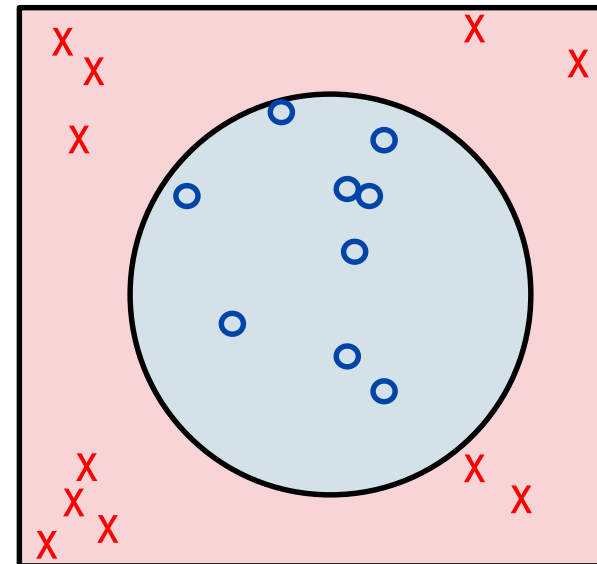


# What, if data is not linearly seperable?

slightly:



seriously:



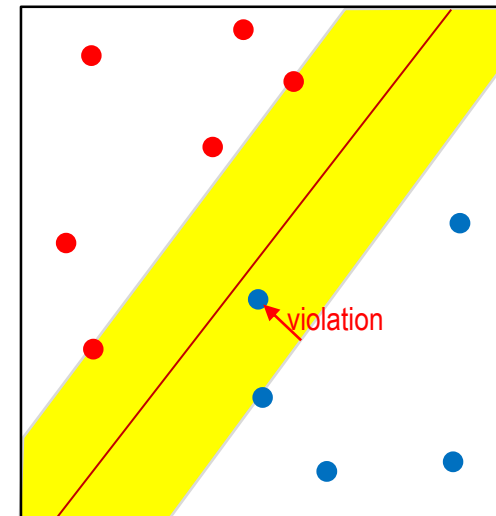


# Accepting margin violations

Margin violation:  $y_n(w^\top x_n + b) \geq 1$  fails

Require:  $y_n(w^\top x_n + b) \geq 1 - \xi_n$  with  $\xi_n \geq 0$

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







## The new optimization

Minimize  $\frac{1}{2} \mathbf{w}^\top \mathbf{w} + c \sum_{n=1}^N \xi_n$

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

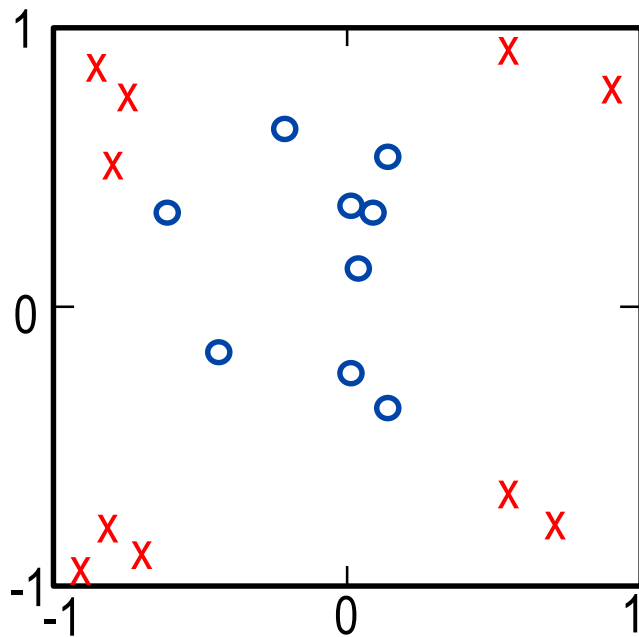
and  $\xi_n \geq 0$  for  $n = 1, \dots, N$

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

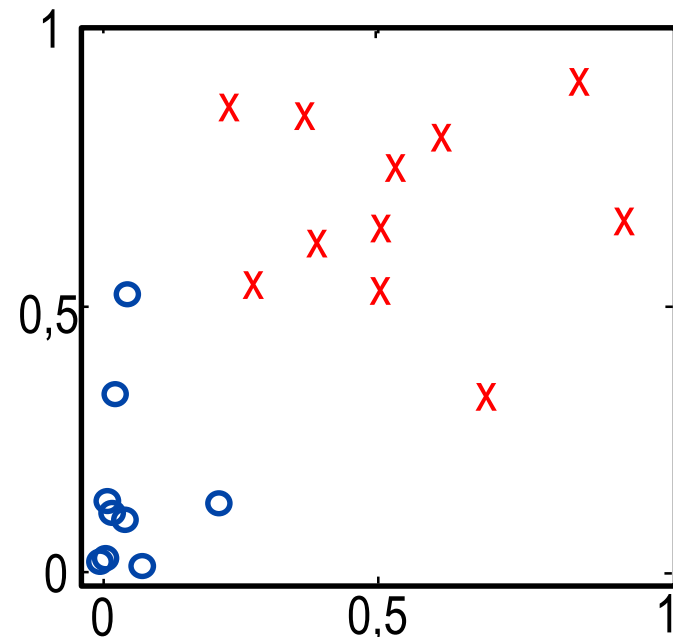


# Using non-linear transformations

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



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





## 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  $\sum_{n=1}^N \alpha_n y_n = 0$

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

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

and  $b : y_m (\mathbf{w}^\top \mathbf{z}_m + b) = 1$  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:  $\mathbf{x} = (x_1, x_2) \rightarrow 2^{\text{nd}}\text{-order } \Phi$

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

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{z}^\top \mathbf{z}' = 1 + x_1 x_1' + x_2 x_2' + x_1^2 x_1'^2 + x_2^2 x_2'^2 + x_1 x_1' x_2 x_2'$$



## The final hypothesis

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

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

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

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



# Design your own kernel

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

1. It is symmetric and 2. The matrix

$$\begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \cdots & K(\mathbf{x}_1, \mathbf{x}_N) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \cdots & K(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & K(\mathbf{x}_N, \mathbf{x}_2) & \cdots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

is positive semi-definite

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



# Backup



# Lagrange formulation

$$\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) =$$

$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + \mathcal{C} \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n (\mathbf{y}_n (\mathbf{w}^\top \mathbf{x}_n + b) - 1 + \xi_n) - \sum_{n=1}^N \beta_n \xi_n$$

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

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

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

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





and the solution is ...

Maximize  $\mathcal{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{x}_n^\top \mathbf{x}_m$  w.r.t. to  $\alpha$

subject to  $0 \leq \alpha_n \leq \mathcal{C}$  for  $n = 1, \dots, N$  and  $\sum_{n=1}^N \alpha_n y_n = 0$

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

$$\text{minimizes } \frac{1}{2} \mathbf{w}^\top \mathbf{w} + \mathcal{C} \sum_{n=1}^N \xi_n$$



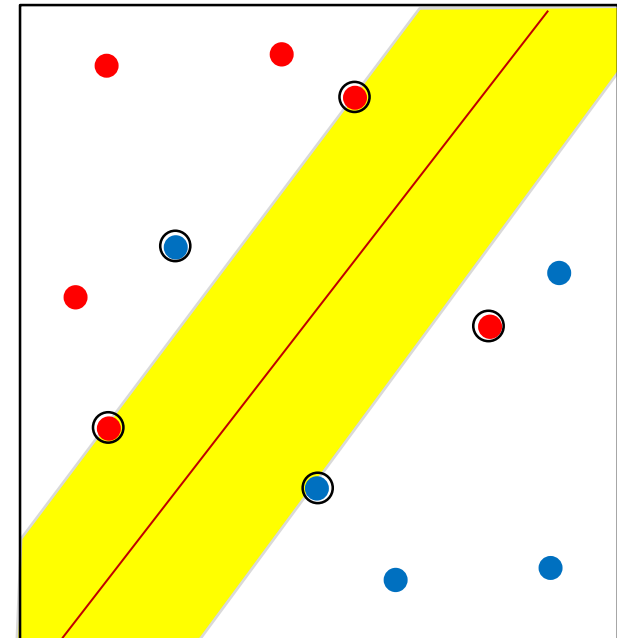
# Types of support vectors

**margin** support vectors  $(0 < \alpha_n < \mathcal{C})$

$$y_n(w^\top x_n + b) = 1 \quad (\xi_n = 0)$$

**non-margin** support vectors  $(\alpha_n = \mathcal{C})$

$$y_n(w^\top x_n + b) < 1 \quad (\xi_n > 0)$$





## The Kernel trick

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

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

$$= 1 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x'_1 x_2 x'_2$$

This is an inner product!

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

$$(1, x'^2_1, x'^2_2, \sqrt{2}x'_1, \sqrt{2}x'_2, \sqrt{2}x'_1x'_2)$$



## 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

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

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



## Kernel formulation of SVM

Remember quadratic programming? The only difference is now:

$$\underbrace{\begin{bmatrix} y_1 y_1 K(x_1, x_1) & y_1 y_2 K(x_1, x_2) & \cdots & y_1 y_N K(x_1, x_N) \\ y_2 y_1 K(x_2, x_1) & y_2 y_2 K(x_2, x_2) & \cdots & y_2 y_N K(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ y_N y_1 K(x_N, x_1) & y_N y_2 K(x_N, x_2) & \cdots & y_N y_N K(x_N, x_N) \end{bmatrix}}_{\text{quadratic coefficients}}$$

Everything else is the same.