

Foundations of Machine Learning

CentraleSupélec — Fall 2017

10. Support Vector Machines

Chloé-Agathe Azencott

Centre for Computational Biology, Mines ParisTech
chloe-agathe.azencott@mines-paristech.fr

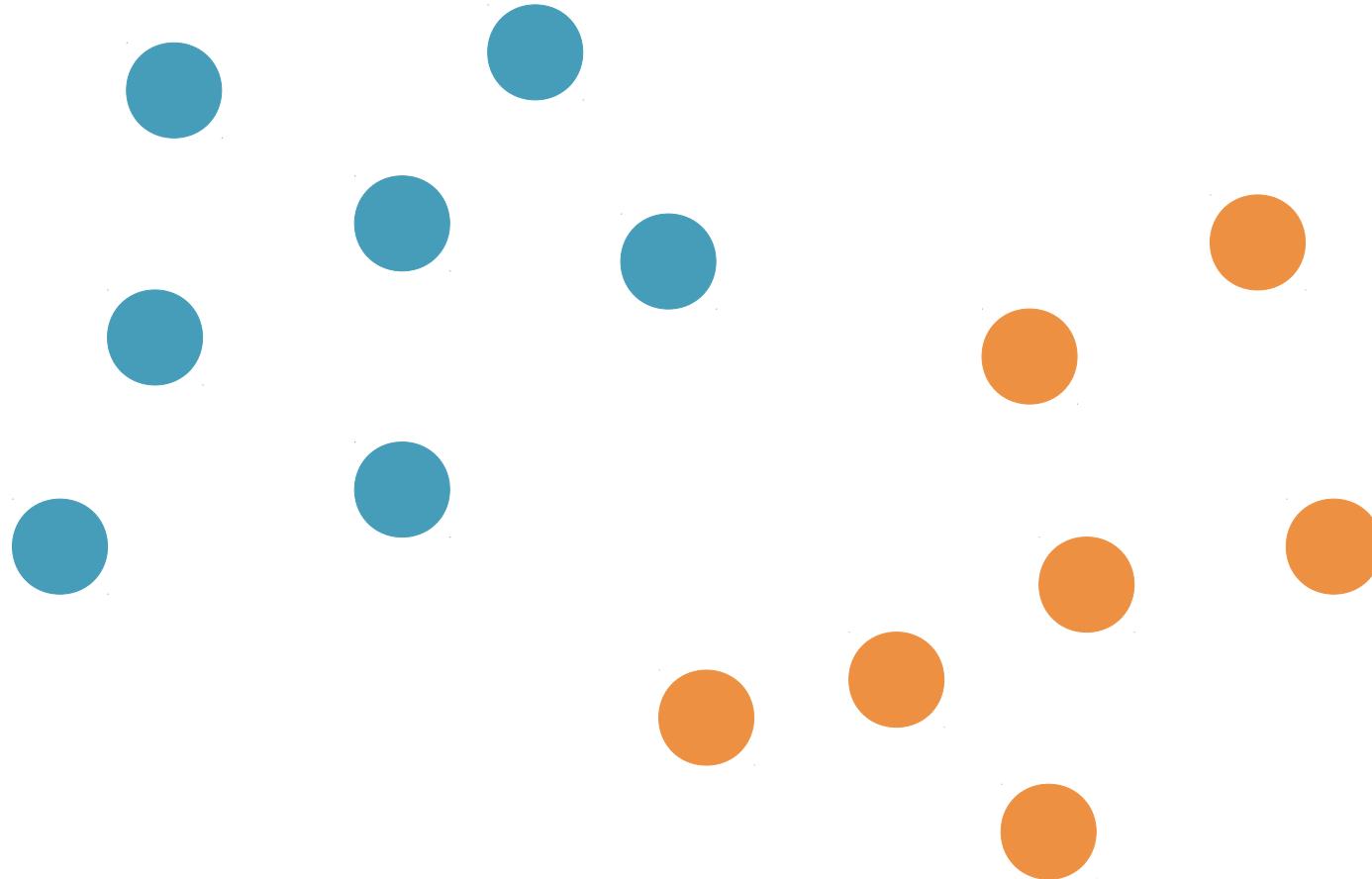


Learning objectives

- Define a **large-margin classifier** in the separable case.
- Write the corresponding **primal** and **dual** optimization problems.
- Re-write the optimization problem in the case of **non-separable data**.
- Use the **kernel trick** to apply soft-margin SVMs to **non-linear** cases.
- Define kernels for **real-valued data, strings, and graphs**.

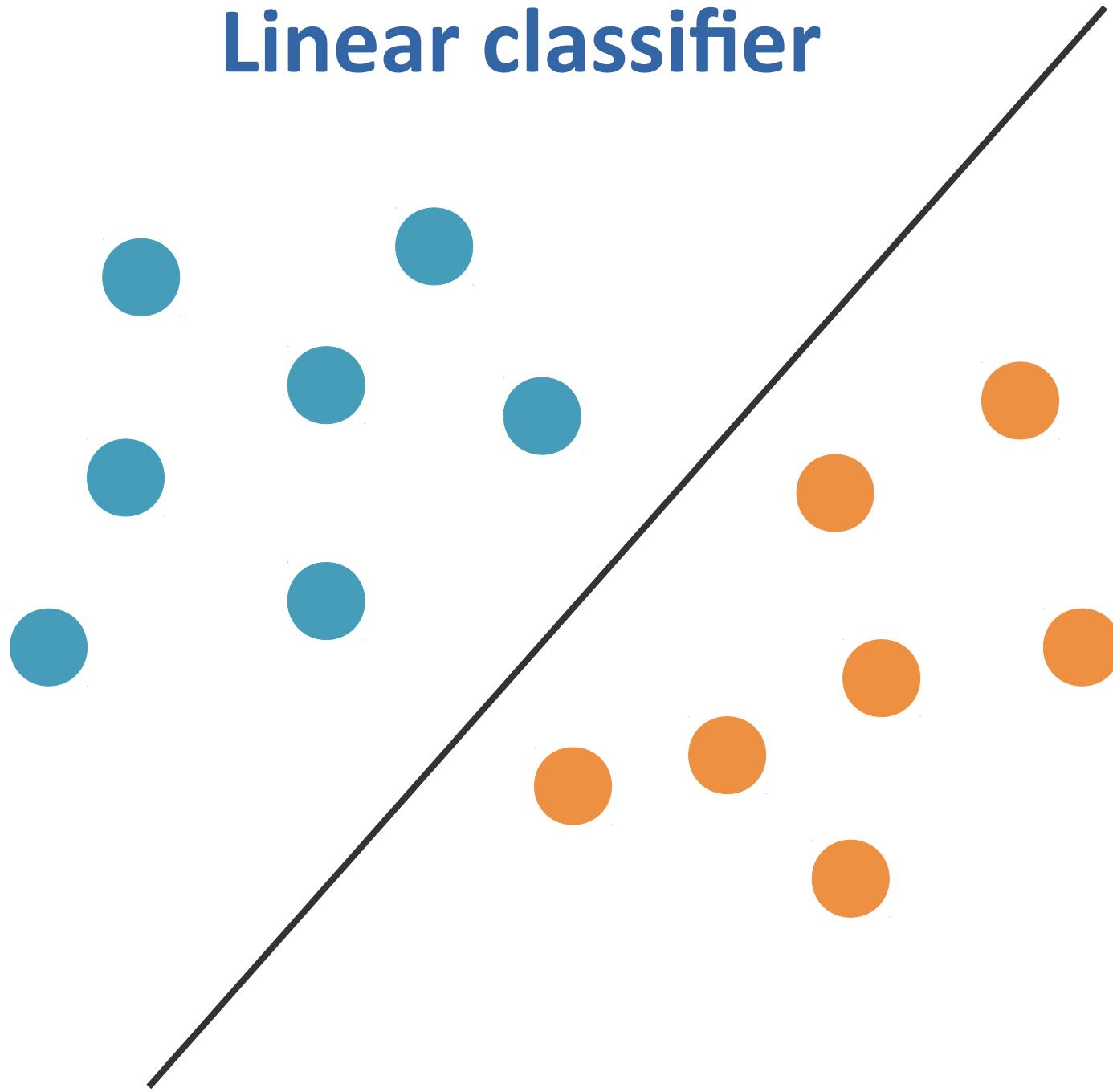
The linearly separable case: hard-margin SVMs

Linear classifier

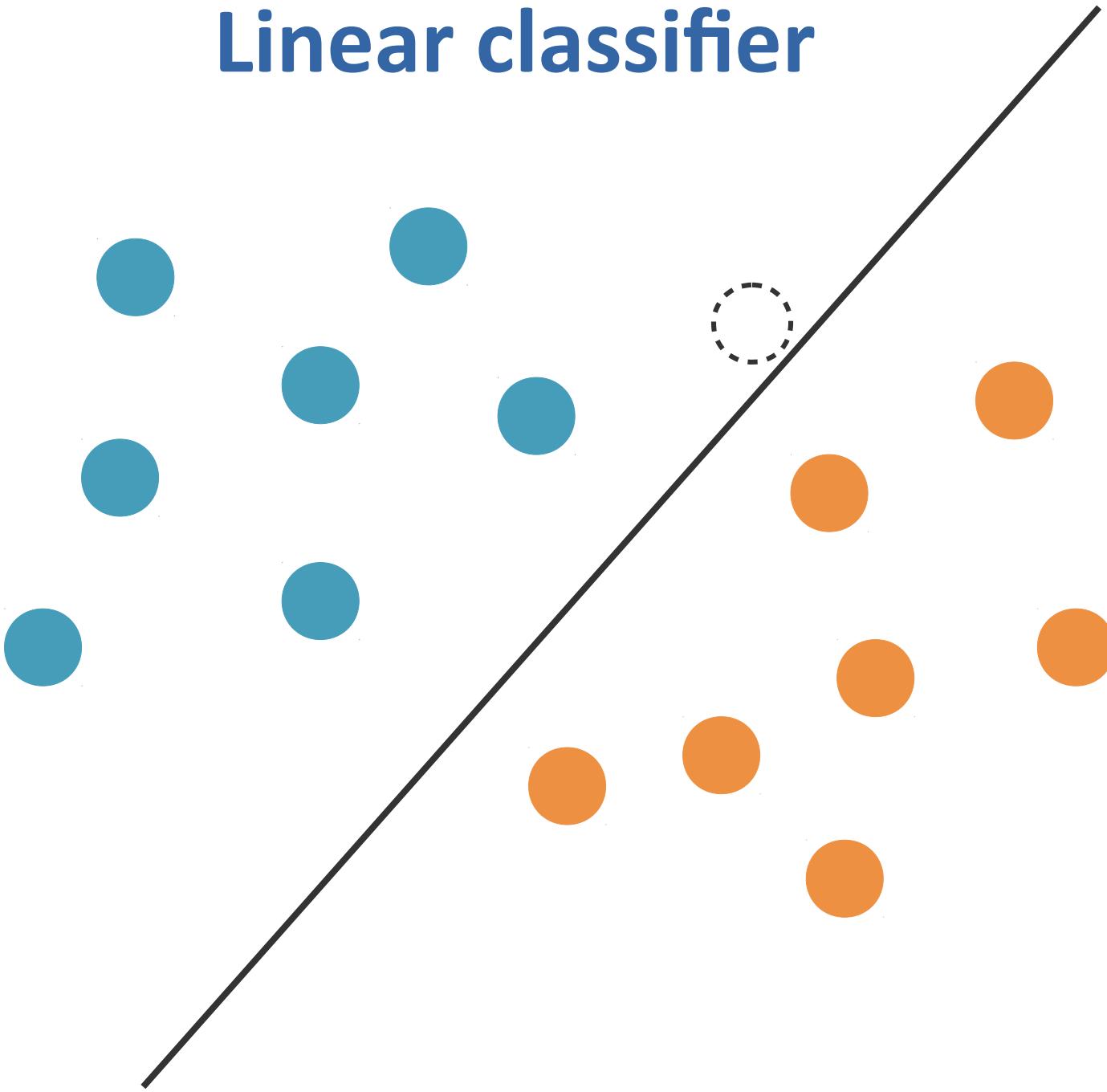


Assume data is **linearly separable**:
there exists a line that separates + from -

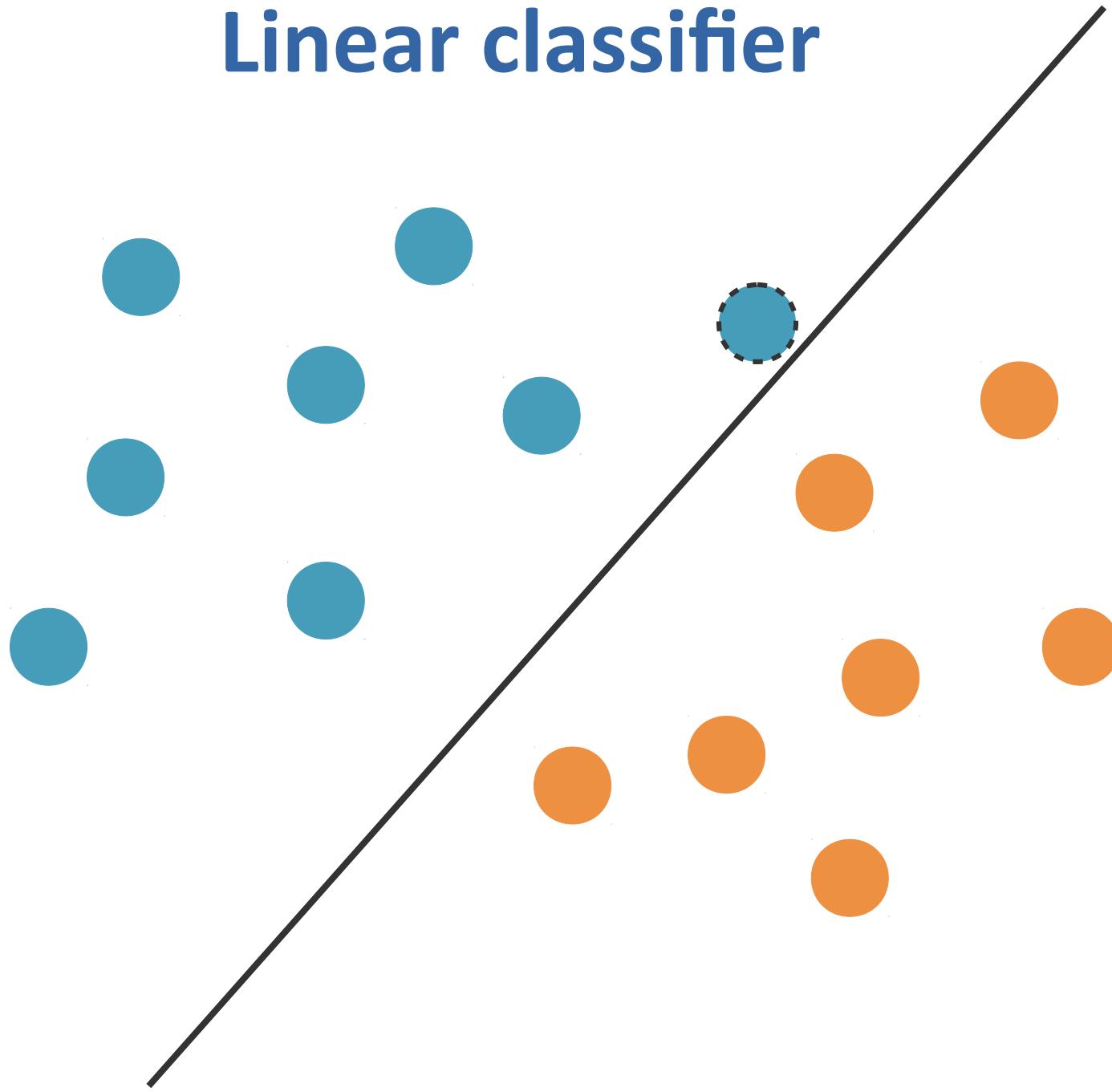
Linear classifier



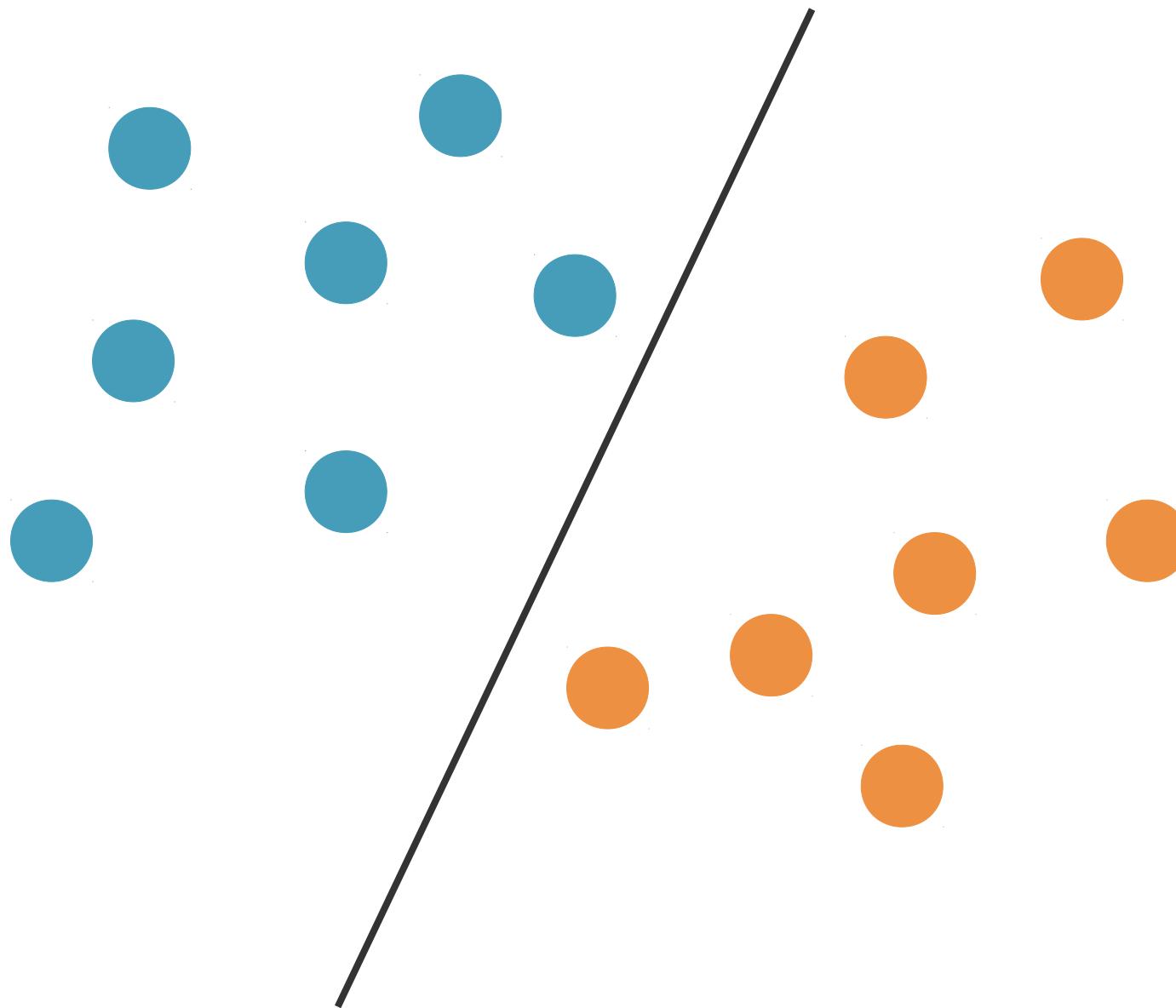
Linear classifier



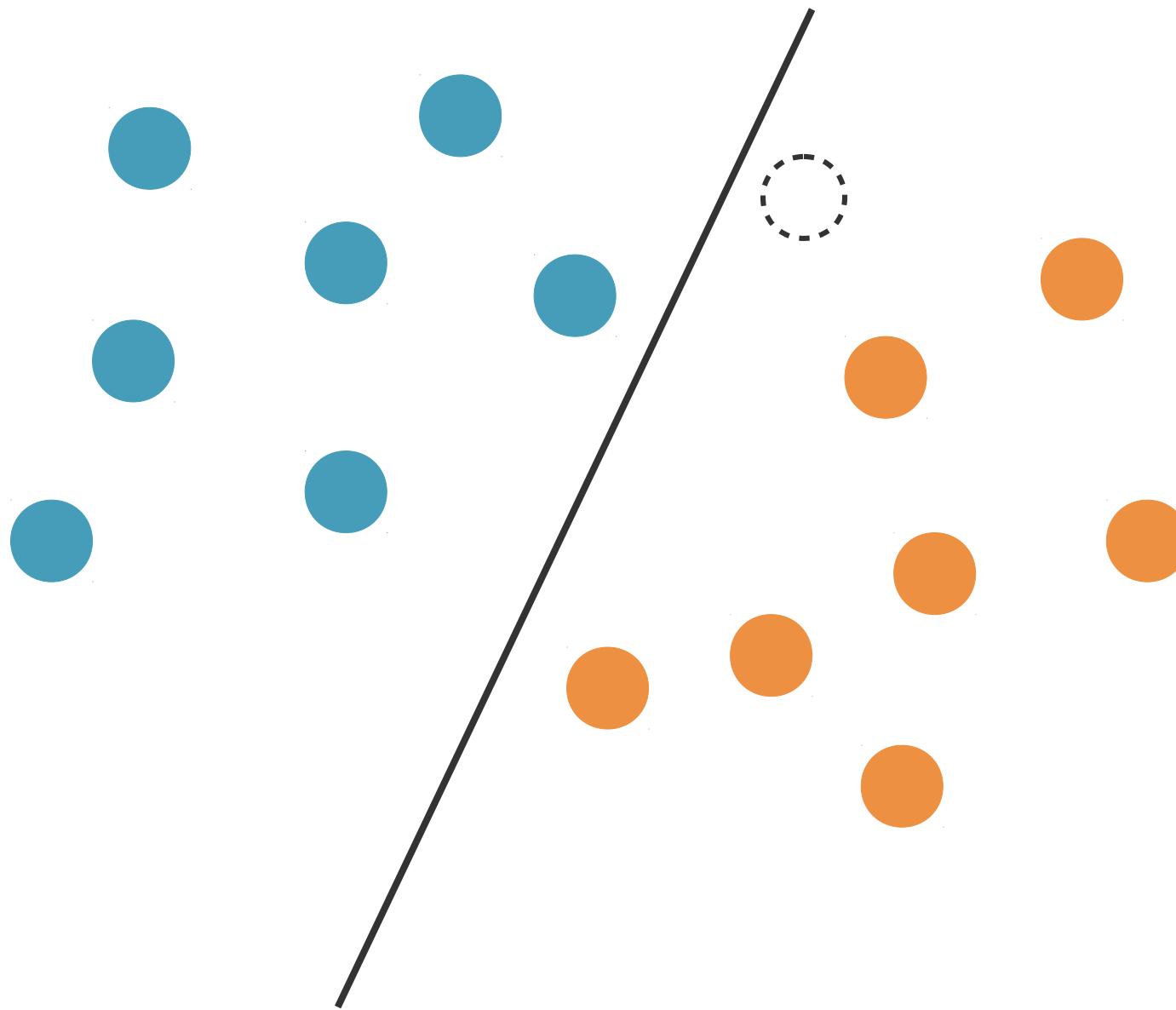
Linear classifier



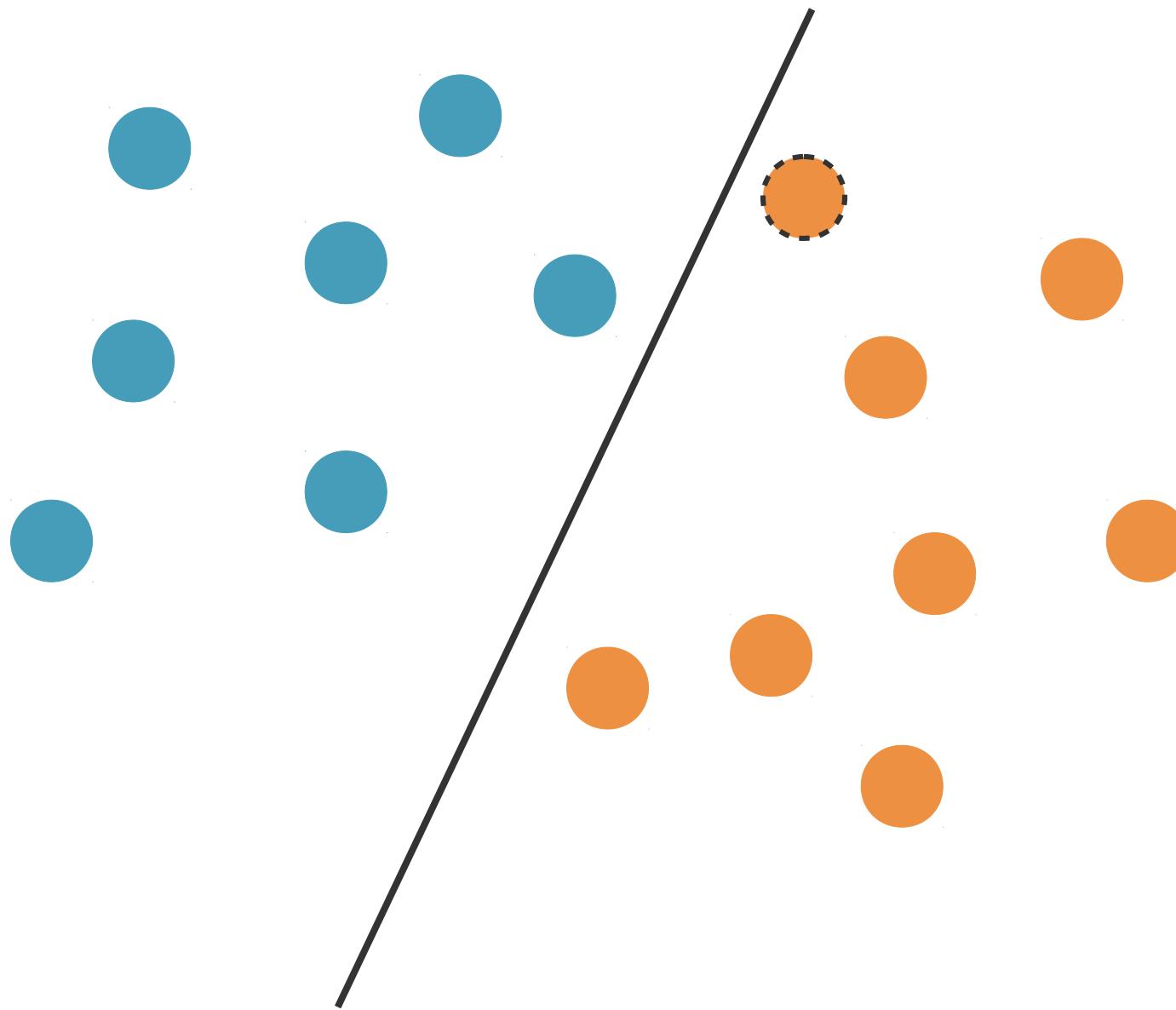
Linear classifier



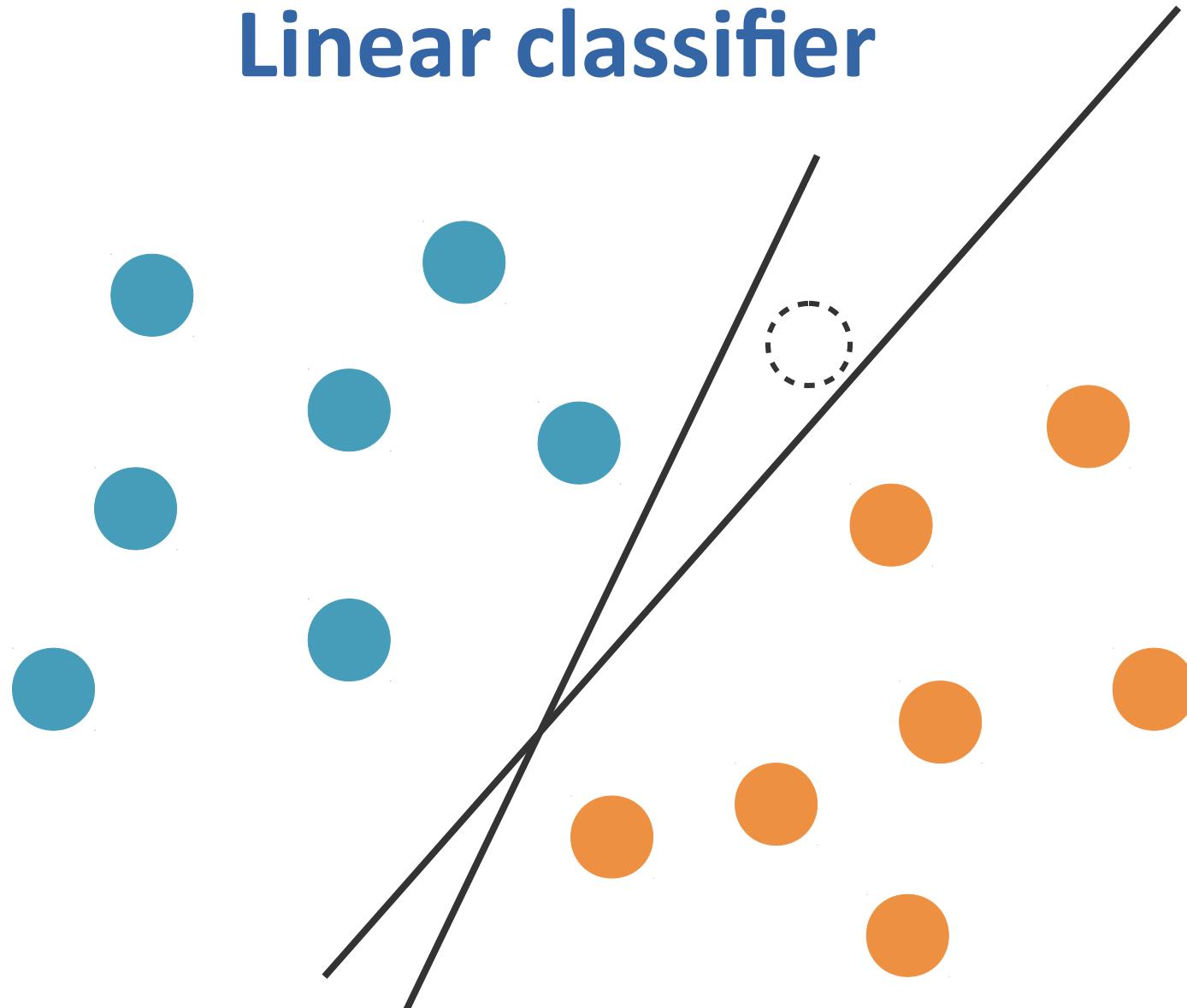
Linear classifier



Linear classifier

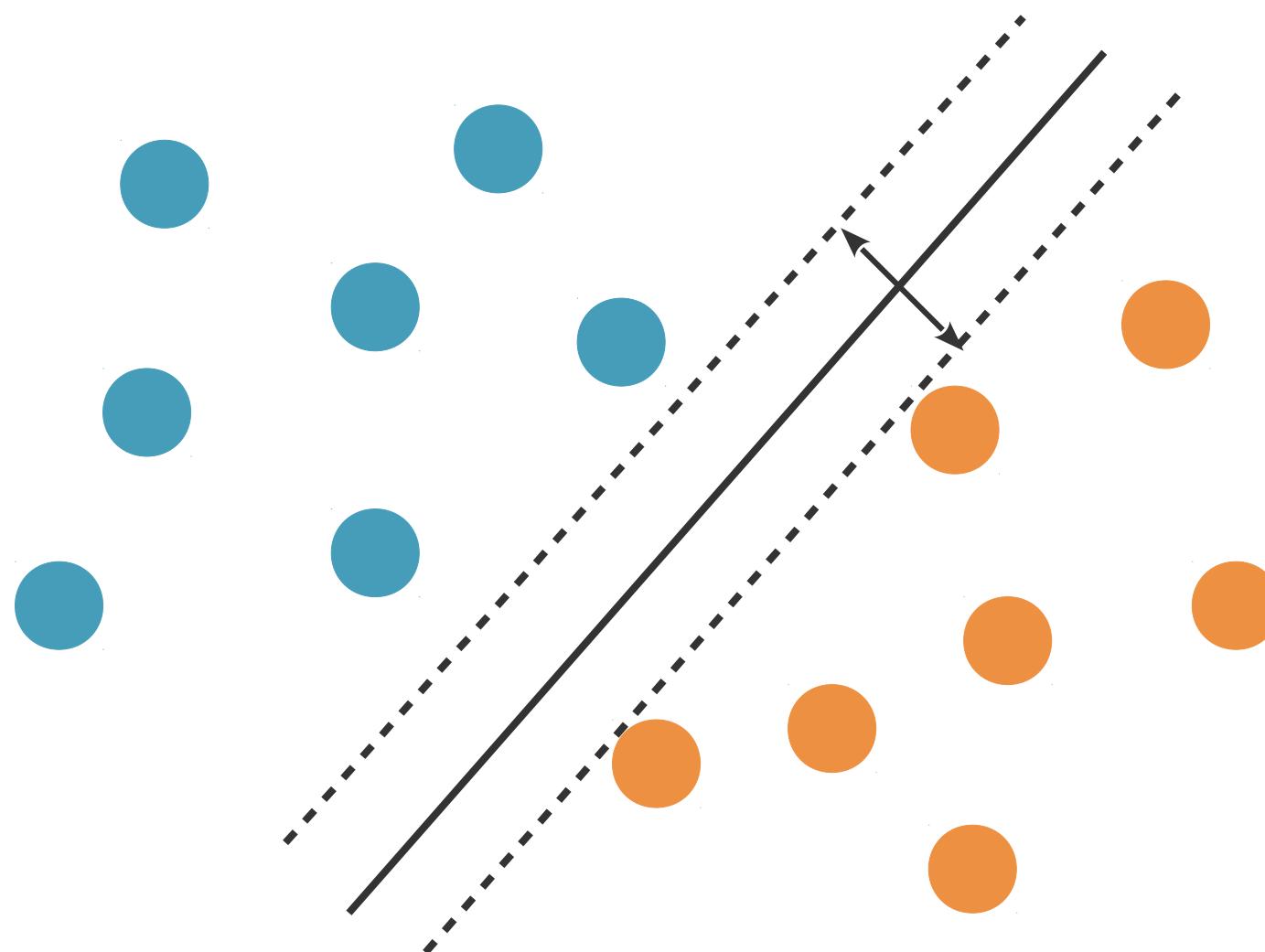


Linear classifier



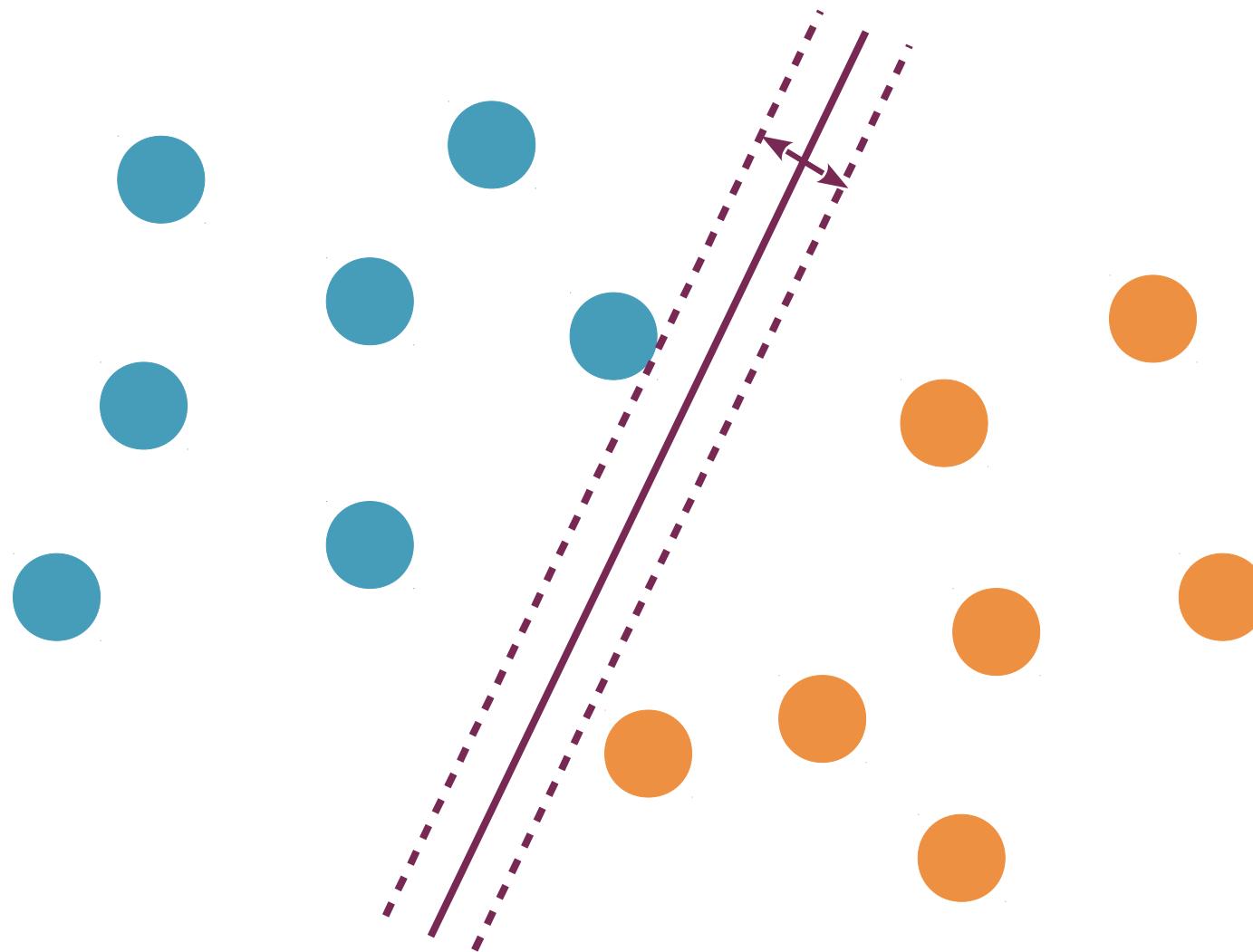
Which one is better?

Margin of a linear classifier

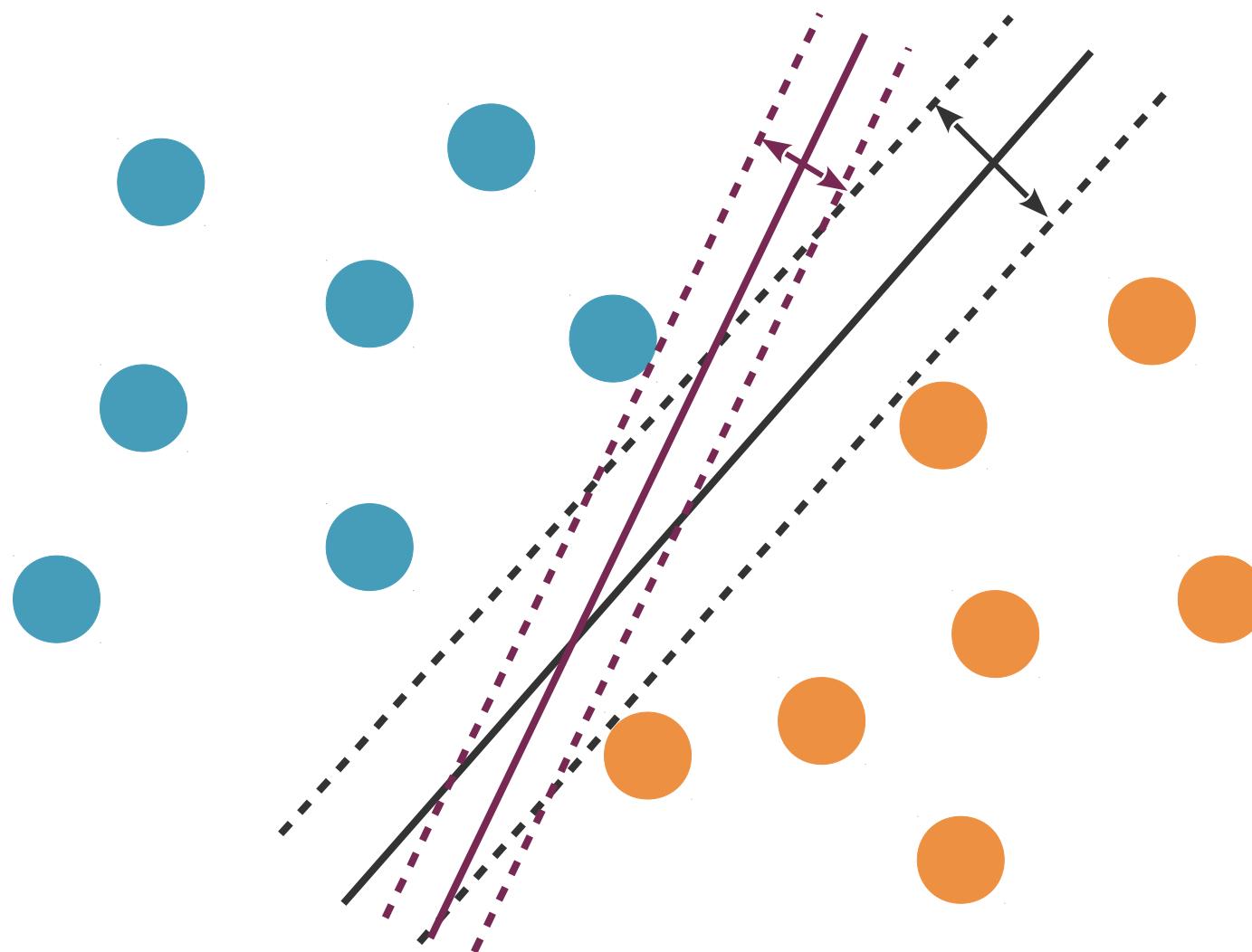


Margin: Twice the distance from the separating hyperplane to the closest training point.

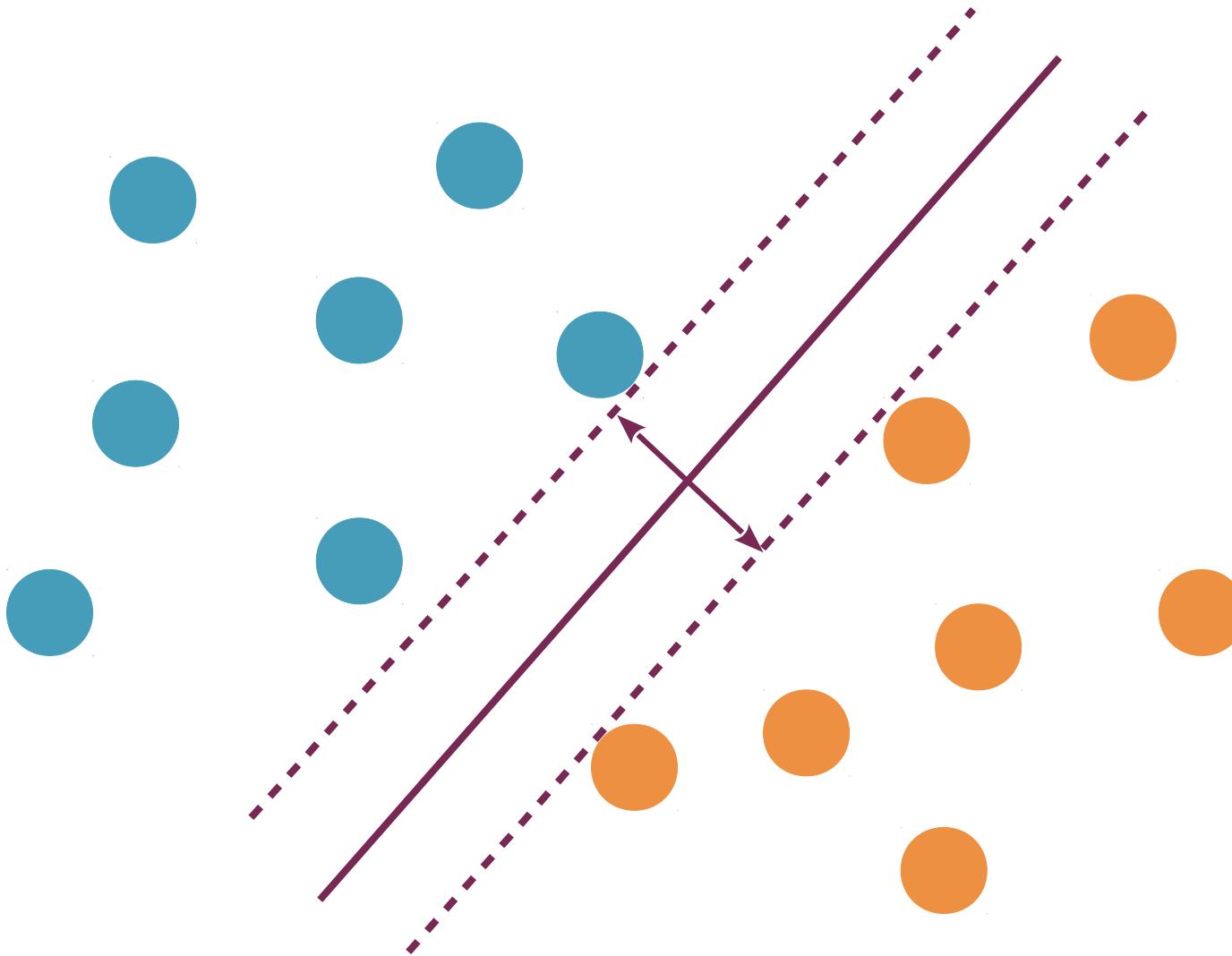
Margin of a linear classifier



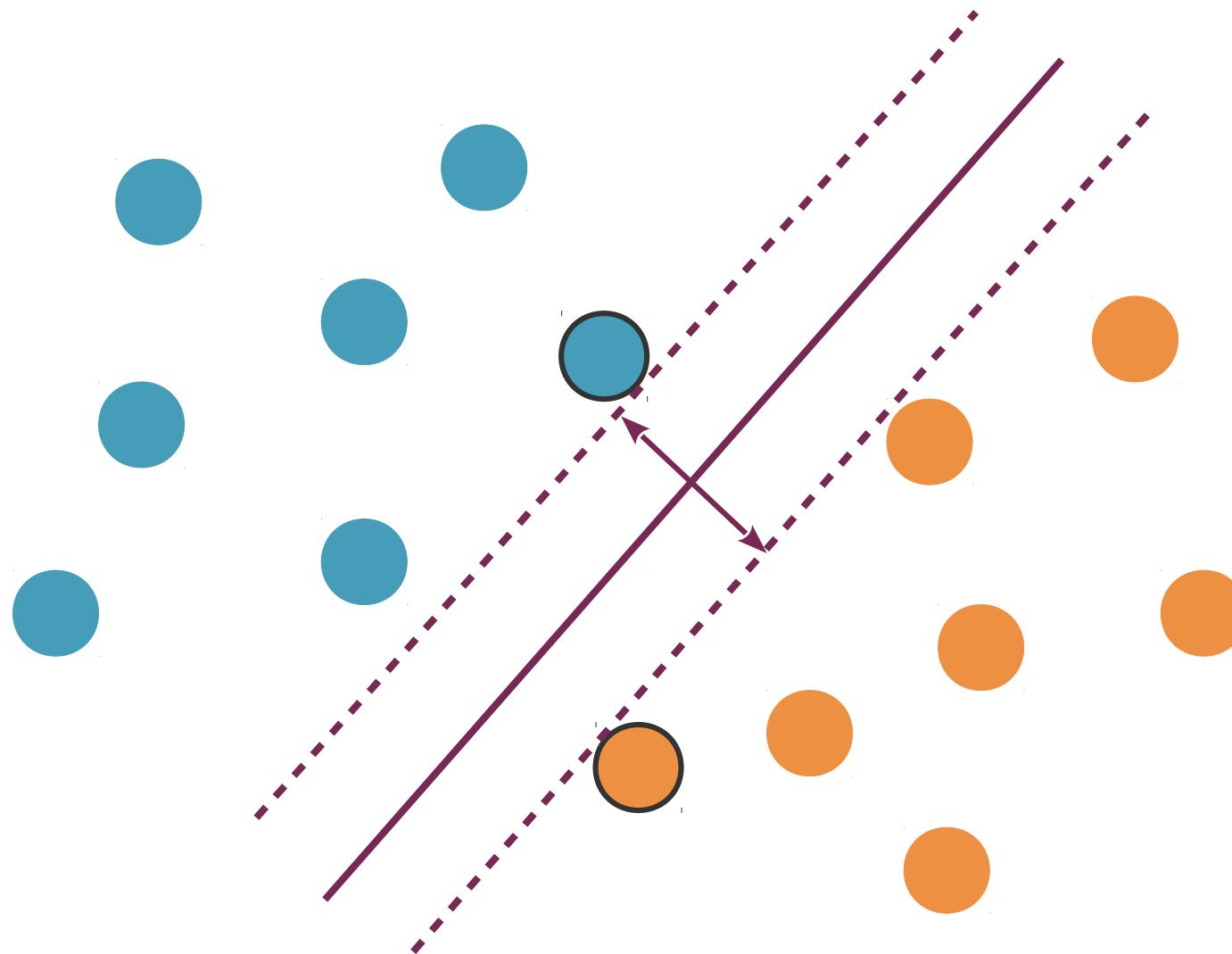
Margin of a linear classifier



Largest margin classifier: Support vector machines



Support vectors

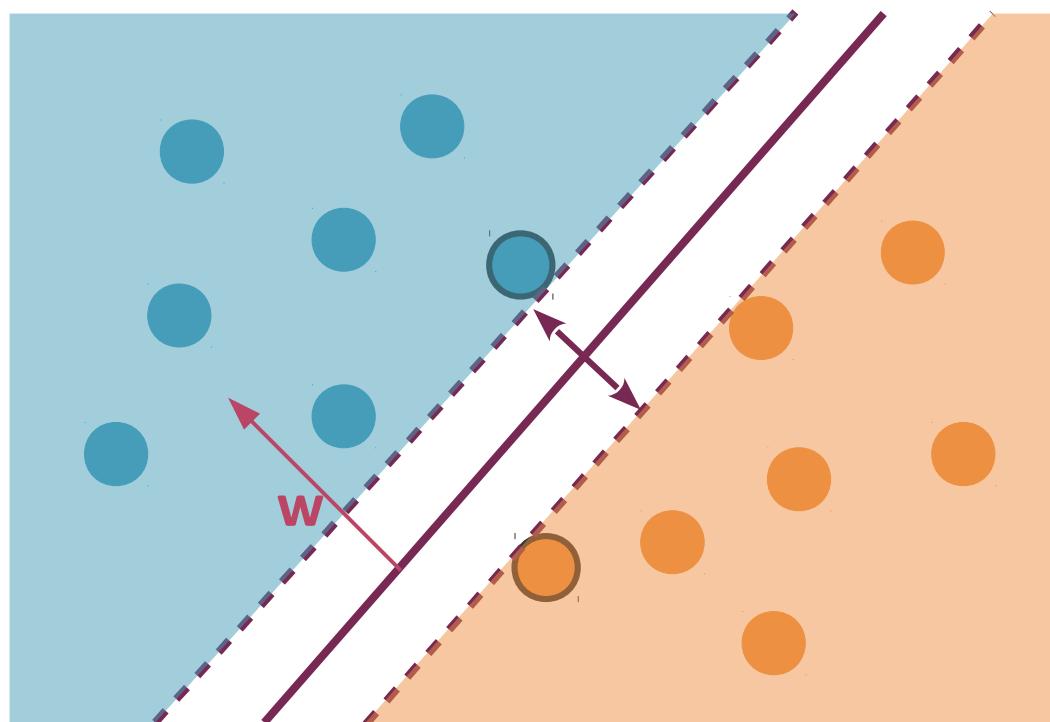


Formalization

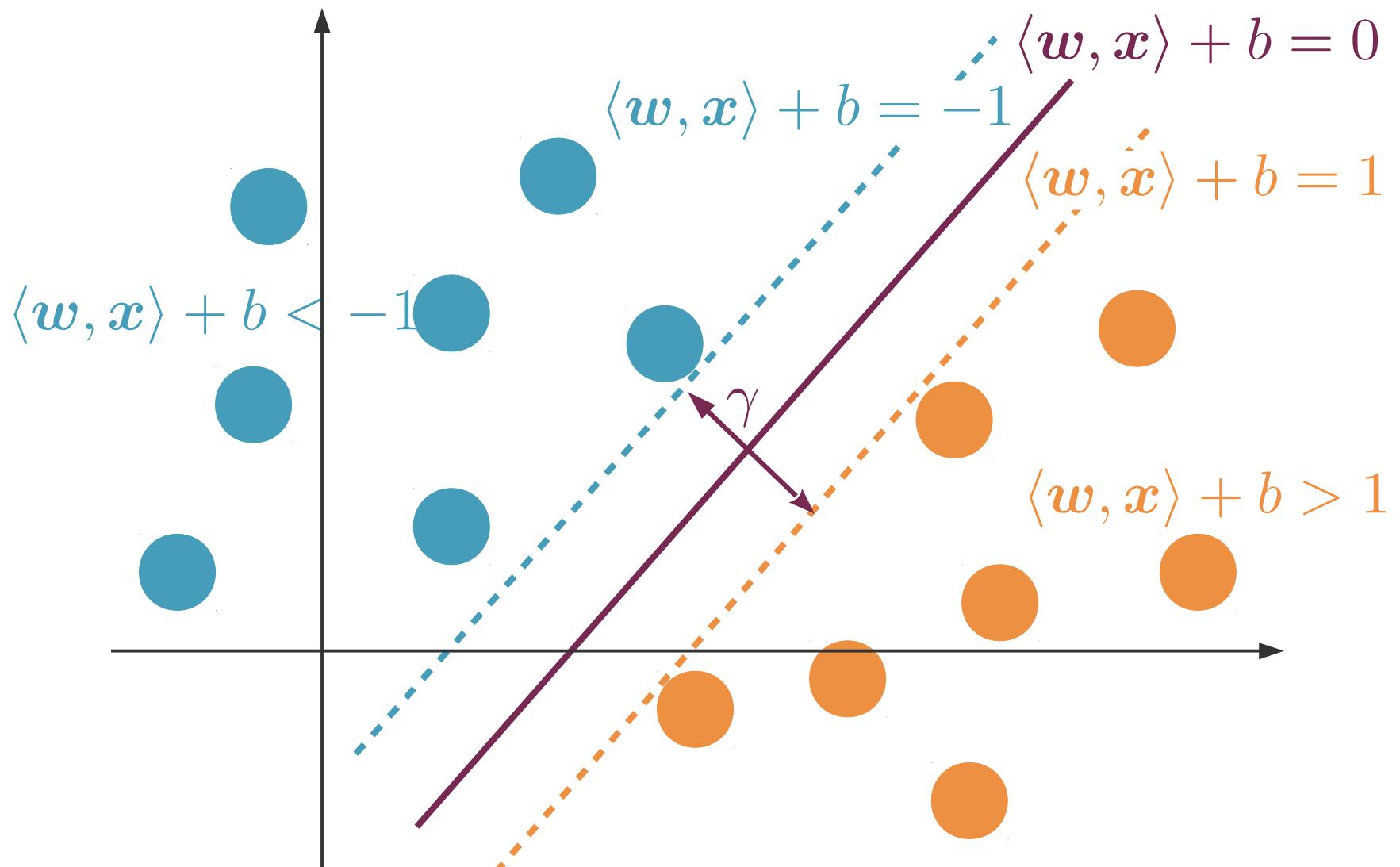
- **Training set**

$$\mathcal{D} = \{(x^1, y^1), \dots, (x^n, y^n)\} \quad x^i \in \mathbb{R}^p \quad y^i \in \{-1, +1\}$$

- What are the equations of the 3 parallel hyperplanes?
- How is the “blue” region defined? The “orange” one?

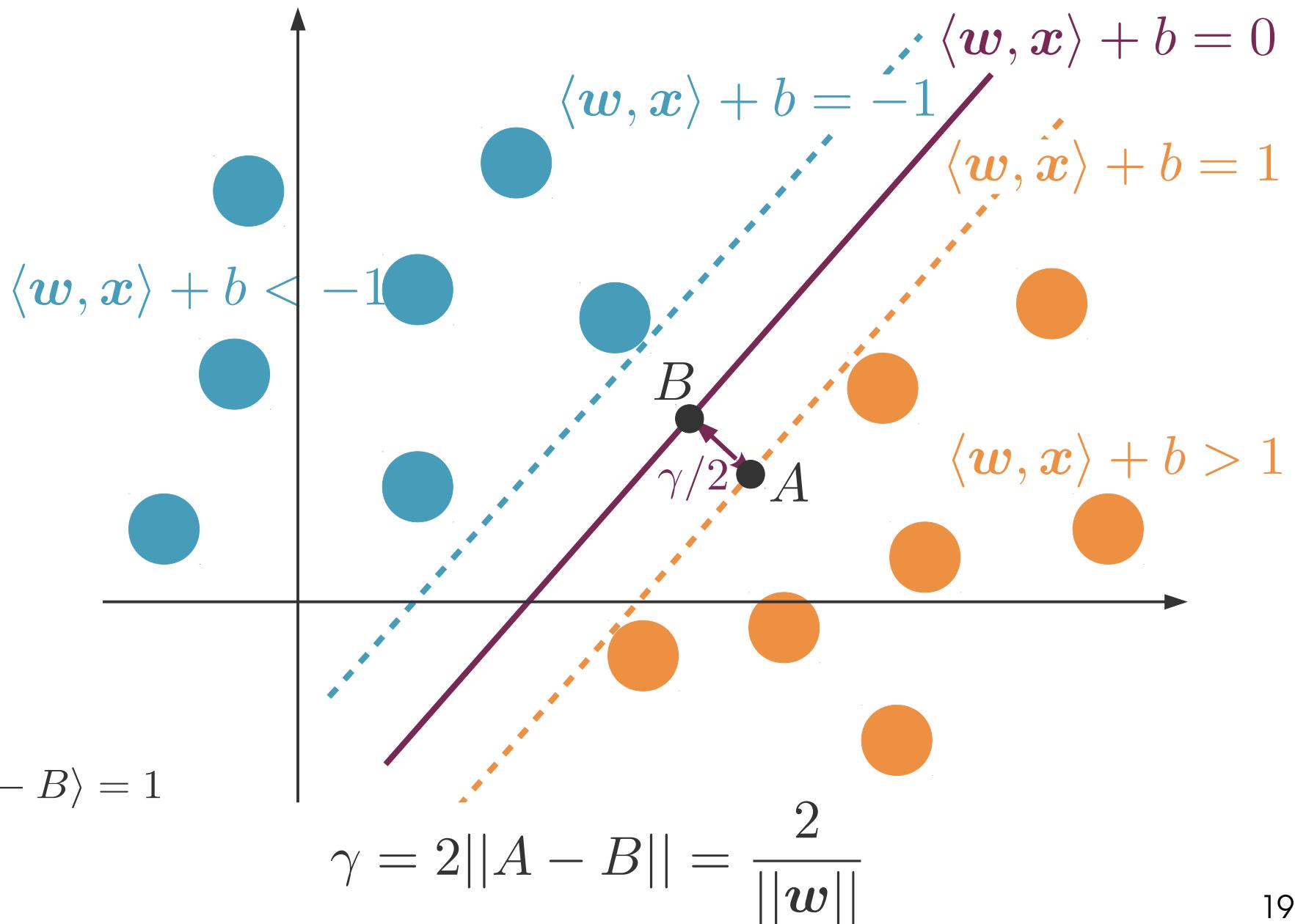


Largest margin hyperplane



What is the size of the margin γ ?

Largest margin hyperplane



Optimization problem

- Training set

$$\mathcal{D} = \{(x^1, y^1), \dots, (x^n, y^n)\} \quad x^i \in \mathbb{R}^p \quad y^i \in \{-1, +1\}$$

- Assume the data to be linearly separable

$$\exists(\mathbf{w}, b) \in \mathbb{R}^p \times \mathbb{R} \text{ s.t. } \begin{cases} \langle \mathbf{w}, x^i \rangle + b > 0 & \text{if } y^i = +1 \\ \langle \mathbf{w}, x^i \rangle + b < 0 & \text{if } y^i = -1 \end{cases}$$

- Goal: Find (\mathbf{w}^*, b^*) that define the hyperplane with largest margin.

Optimization problem

- Margin maximization:

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2$$

- Correct classification of the training points:

- For positive examples:

$$y^i = 1 \text{ and } \langle \mathbf{w}, \mathbf{x}^i \rangle + b \geq 1$$

- For negative examples:

$$y^i = -1 \text{ and } \langle \mathbf{w}, \mathbf{x}^i \rangle + b \leq -1$$

- Summarized as



Optimization problem

- Margin maximization:

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2$$

- Correct classification of the training points:

- For positive examples:

$$y^i = 1 \text{ and } \langle \mathbf{w}, \mathbf{x}^i \rangle + b \geq 1$$

- For negative examples:

$$y^i = -1 \text{ and } \langle \mathbf{w}, \mathbf{x}^i \rangle + b \leq -1$$

- Summarized as: $y^i \cdot (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) \geq 1$

- Optimization problem:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^p, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ s.t. } y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1 \geq 0 \quad \forall i \in \{1, \dots, n\}$$

Optimization problem

- Find (w^*, b^*) that minimize $\frac{1}{2}||w||^2$
under the n constraints $y^i(\langle w, x^i \rangle + b) - 1 \geq 0$



Optimization problem

- Find (\mathbf{w}^*, b^*) that minimize $\frac{1}{2} \|\mathbf{w}\|^2$ under the n constraints $y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1 \geq 0$
- We introduce one **dual variable** α_i for each constraint (i.e. each training point)
- **Lagrangian:**

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1).$$

$$\mathbf{w} \in \mathbb{R}^p \quad \boldsymbol{\alpha} \in \mathbb{R}_+^n \quad b \in \mathbb{R}$$

Lagrange dual of the SVM

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1)$$

- **Lagrange dual function:** $q : \mathbb{R}^n \rightarrow \mathbb{R}$

$$q(\alpha) = \inf_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha)$$

- **Lagrange dual problem:**

$$\max_{\alpha} \inf_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1)$$

- **Strong duality:** Under **Slater's conditions**, the optimum of the primal is the optimum of the dual.

The function to optimize is convex and the equality constraints are affine.

Minimizing the Lagrangian of the SVM

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1)$$

- $L(\mathbf{w}, b, \alpha)$ is **convex quadratic** in \mathbf{w} and minimized for



Minimizing the Lagrangian of the SVM

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1)$$

- $L(\mathbf{w}, b, \alpha)$ is **convex quadratic in \mathbf{w}** and minimized for:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, b^*, \alpha) = 0 \Leftrightarrow \mathbf{w}^* - \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i = 0 \Leftrightarrow \mathbf{w}^* = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i.$$

Minimizing the Lagrangian of the SVM

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1)$$

- $L(\mathbf{w}, b, \alpha)$ is **convex quadratic in \mathbf{w}** and minimized for:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, b^*, \alpha) = 0 \Leftrightarrow \mathbf{w}^* - \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i = 0 \Leftrightarrow \mathbf{w}^* = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i.$$

- $L(\mathbf{w}, b, \alpha)$ is **affine in b** . Its minimum is $-\infty$ except if. 

Minimizing the Lagrangian of the SVM

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1)$$

- $L(\mathbf{w}, b, \alpha)$ is **convex quadratic in \mathbf{w}** and minimized for:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, b^*, \alpha) = 0 \Leftrightarrow \mathbf{w}^* - \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i = 0 \Leftrightarrow \mathbf{w}^* = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i.$$

- $L(\mathbf{w}, b, \alpha)$ is **affine in b** . Its minimum is $-\infty$ except if:

$$\nabla_b L(\mathbf{w}^*, b^*, \alpha) = 0 \Leftrightarrow \sum_{i=1}^n \alpha_i y^i = 0.$$

SVM dual problem

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1)$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i$$

$$L(\mathbf{w}^*, b, \boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - b \sum_{i=1}^n \alpha_i y^i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$\sum_{i=1}^n \alpha_i y^i = 0$$

$$L(\mathbf{w}^*, b, \boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

- **Lagrange dual function:**

$$\begin{aligned} q(\boldsymbol{\alpha}) &= \inf_{\mathbf{w} \in \mathbb{R}^p, b \in \mathbb{R}} L(\mathbf{w}, b, \boldsymbol{\alpha}) \\ &= \begin{cases} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle & \text{if } \sum_{i=1}^n \alpha_i y^i = 0 \\ -\infty & \text{otherwise.} \end{cases} \end{aligned}$$

- **Dual problem:** maximize $q(\boldsymbol{\alpha})$ subject to $\boldsymbol{\alpha} \geq 0$.

Maximizing a quadratic function under box constraints can be solved efficiently using dedicated software.

Optimal hyperplane

- Once the optimal α^* is found, we recover (\mathbf{w}^*, b^*)

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y^i \mathbf{x}^i$$

$$b^* = \frac{\max_{i:y^i=-1} \langle \mathbf{w}^*, \mathbf{x}^i \rangle + \min_{i:y^i=+1} \langle \mathbf{w}^*, \mathbf{x}^i \rangle}{2}$$

- Determining b^* :

- Closest positive point to the separating hyperplane: $\arg \min_{i:y^i=+1} \langle \mathbf{w}^*, \mathbf{x}^i \rangle$ verifies $\langle \mathbf{w}^*, \mathbf{x} \rangle + b^* = 1$

- Closest negative point to the separating hyperplane: $\arg \max_{i:y^i=-1} \langle \mathbf{w}^*, \mathbf{x}^i \rangle$ verifies $\langle \mathbf{w}^*, \mathbf{x} \rangle + b^* = -1$

- The **decision function** is hence:

$$\begin{aligned} f^*(\mathbf{x}) &= \langle \mathbf{w}^*, \mathbf{x} \rangle + b^* \\ &= \sum_{i=1}^n \alpha_i^* y^i \langle \mathbf{x}^i, \mathbf{x} \rangle + b^* \end{aligned}$$

Lagrangian

- minimize $f(w)$ under the constraint $g(w) \geq 0$

$$f(w) = \frac{1}{2} \|w\|^2$$

$$g(w) = y(\langle w, x \rangle + b) - 1$$

abusive notation: $g(w, b)$

How do we write this in terms of the gradients of f and g ?

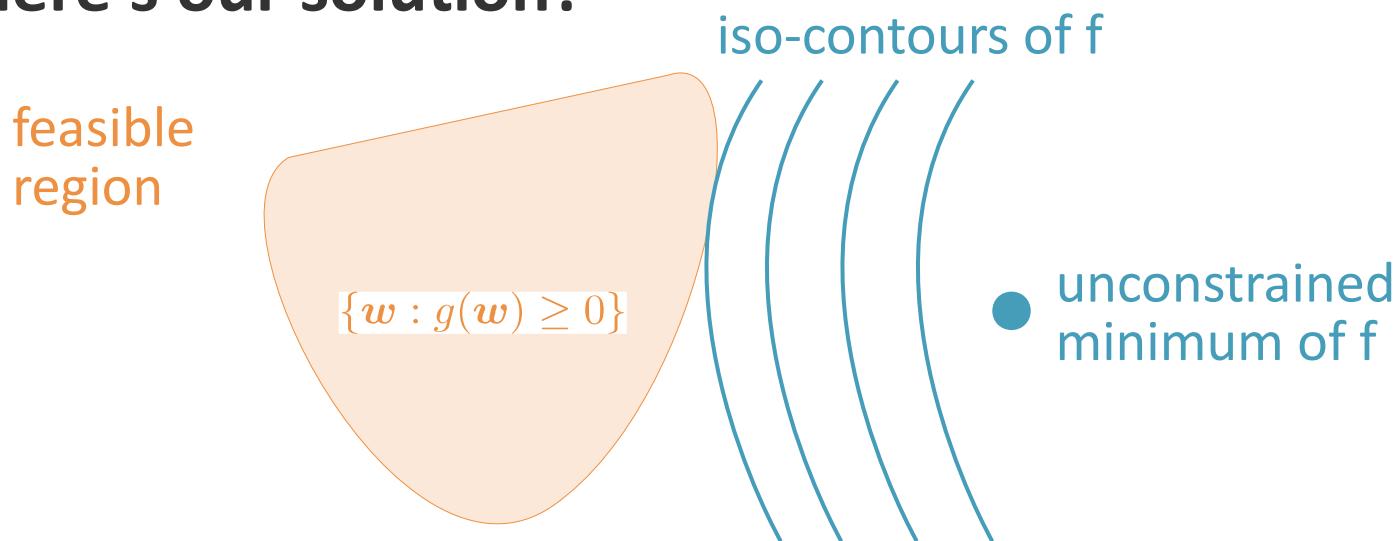
Lagrangian

- minimize $f(\mathbf{w})$ under the constraint $\mathbf{g}(\mathbf{w}) \geq 0$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$$

$$g(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b) - 1$$

If the minimum of $f(\mathbf{w})$ doesn't lie in the feasible region, where's our solution?



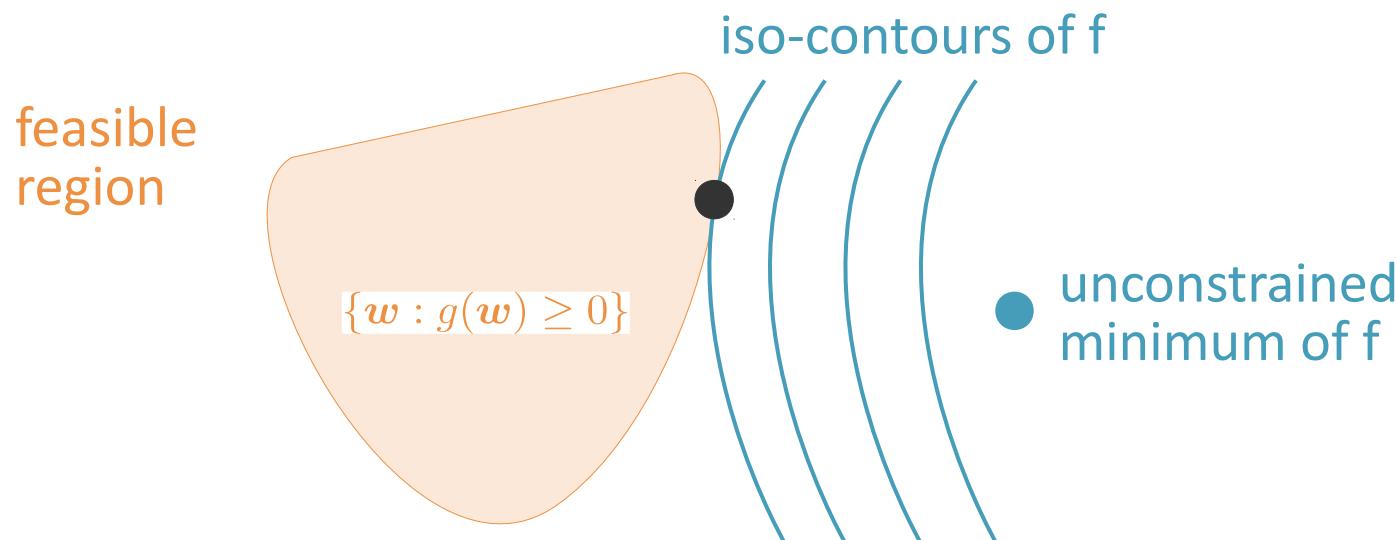
Lagrangian

- minimize $f(\mathbf{w})$ under the constraint $g(\mathbf{w}) \geq 0$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$$

$$g(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b) - 1$$

How do we write this in terms of the gradients of f and g ?



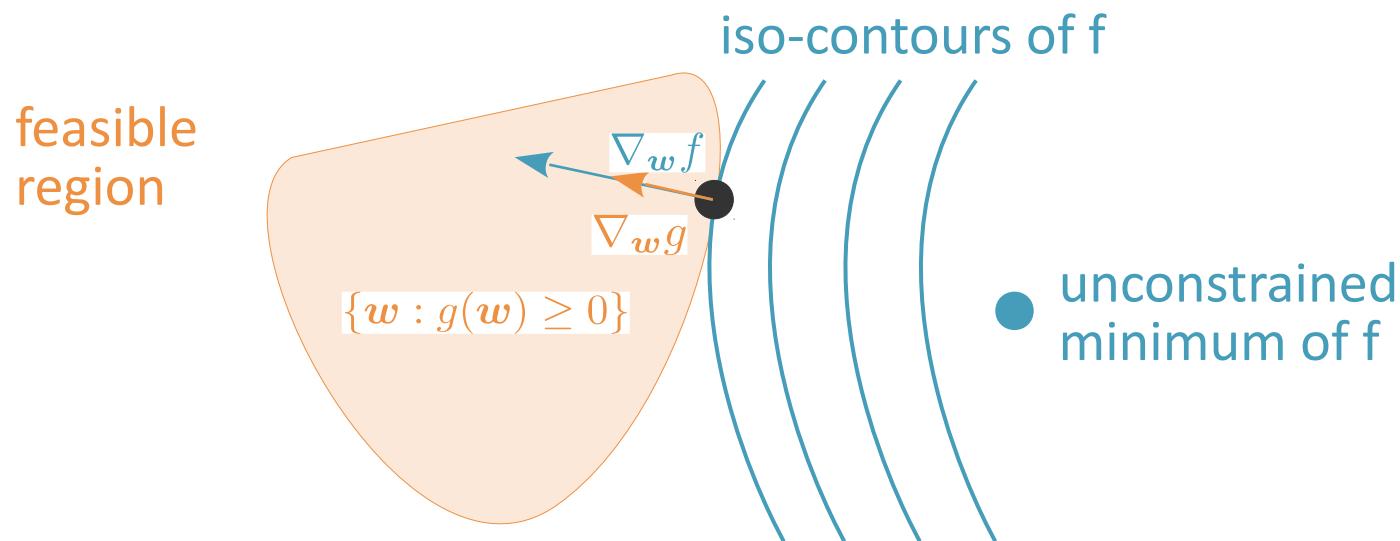
Lagrangian

- minimize $f(\mathbf{w})$ under the constraint $g(\mathbf{w}) \geq 0$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$$

$$g(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b) - 1$$

How do we write this in terms of the gradients of f and g ?



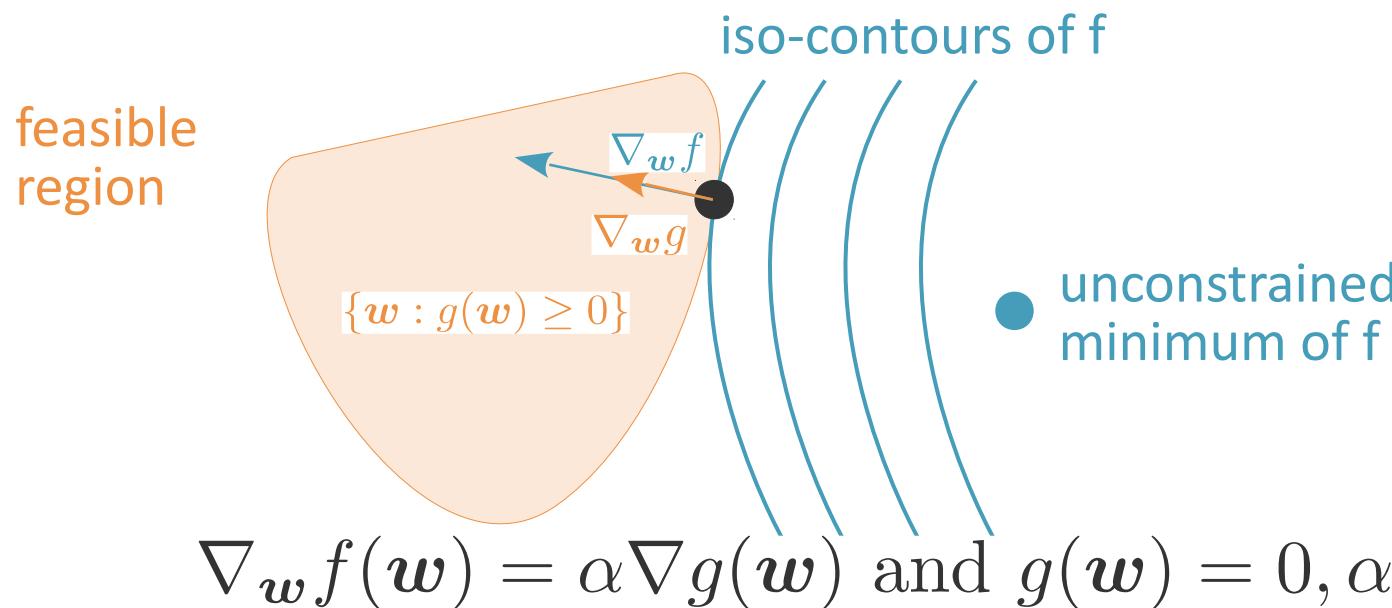
Lagrangian

- minimize $f(\mathbf{w})$ under the constraint $g(\mathbf{w}) \geq 0$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$$

$$g(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b) - 1$$

How do we write this in terms of the gradients of f and g ?



Lagrangian

- minimize $f(\mathbf{w})$ under the constraint $g(\mathbf{w}) \geq 0$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \quad g(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b) - 1$$

Case 1: the unconstrained minimum lies in the feasible region.

Case 2: it does not.

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = \alpha \nabla g(\mathbf{w}) \text{ and } g(\mathbf{w}) = 0, \alpha > 0$$

How do we summarize both cases?

Lagrangian

- minimize $f(\mathbf{w})$ under the constraint $\mathbf{g}(\mathbf{w}) \geq \mathbf{0}$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \quad g(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b) - 1$$

Case 1: the unconstrained minimum lies in the feasible region.

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = 0 \text{ and } g(\mathbf{w}) \geq 0$$

Case 2: it does not.

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = \alpha \nabla g(\mathbf{w}) \text{ and } g(\mathbf{w}) = 0, \alpha > 0$$

- Summarized as:

$$\begin{cases} \nabla_{\mathbf{w}}(f(\mathbf{w}) - \alpha g(\mathbf{w})) = 0 \\ \alpha g(\mathbf{w}) = 0 \end{cases} \quad \text{and } \alpha \geq 0.$$

Lagrangian

- minimize $f(\mathbf{w})$ under the constraint $g(\mathbf{w}) \geq 0$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \quad g(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b) - 1$$

$$\begin{cases} \nabla_{\mathbf{w}}(f(\mathbf{w}) - \alpha g(\mathbf{w})) = 0 \\ \alpha g(\mathbf{w}) = 0 \end{cases} \quad \text{and } \alpha \geq 0.$$

Lagrangian: $L(\mathbf{w}, \alpha) = f(\mathbf{w}) - \alpha g(\mathbf{w})$

α is called the **Lagrange multiplier**.

Lagrangian

- minimize $f(\mathbf{w})$ under the constraints $\mathbf{g}_i(\mathbf{w}) \geq 0$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \quad g_i(\mathbf{w}) = y_i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1$$

$$\begin{cases} \nabla_{\mathbf{w}}(f(\mathbf{w}) - \alpha_i g_i(\mathbf{w})) = 0 \\ \alpha_i g_i(\mathbf{w}) = 0 \end{cases} \quad \text{and } \alpha_i \geq 0.$$

How do we deal with n constraints?

Lagrangian

- minimize $f(\mathbf{w})$ under the constraints $\mathbf{g}_i(\mathbf{w}) \geq 0$

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \quad g_i(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1$$

$$\begin{cases} \nabla_{\mathbf{w}}(f(\mathbf{w}) - \alpha_i g_i(\mathbf{w})) = 0 \\ \alpha_i g_i(\mathbf{w}) = 0 \end{cases} \quad \text{and } \alpha_i \geq 0.$$

Use n Lagrange multipliers

- Lagrangian:

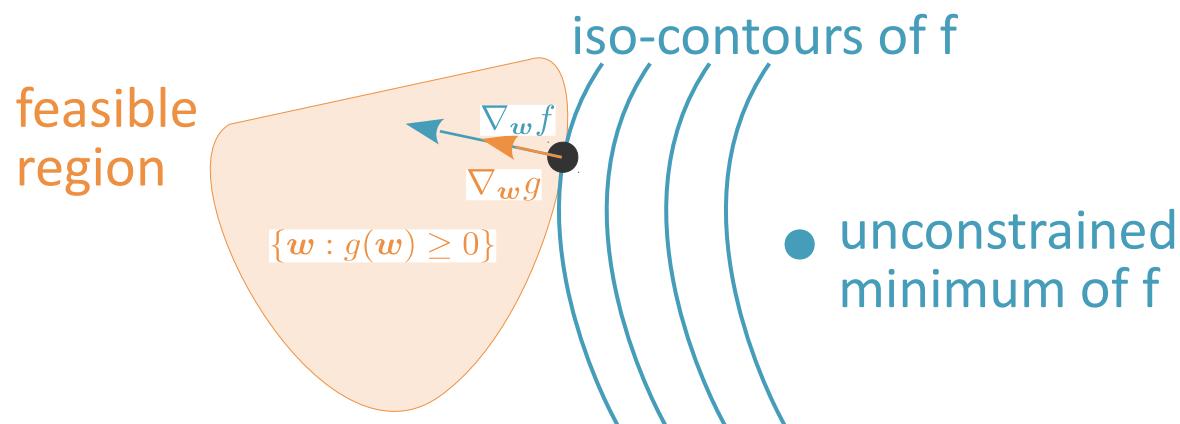
$$L(\mathbf{w}, \boldsymbol{\alpha}) = f(\mathbf{w}) - \sum_{i=1}^n \alpha_i g_i(\mathbf{w})$$

Support vectors

- **Karun-Kush-Tucker conditions:**

Either $\alpha_i = 0$ (case 1) or $g_i=0$ (case 2)

$$g_i(\mathbf{w}) = y(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1 \quad \begin{cases} \nabla_{\mathbf{w}}(f(\mathbf{w}) - \alpha_i g_i(\mathbf{w})) = 0 \\ \alpha_i g_i(\mathbf{w}) = 0 \end{cases} \quad \text{and } \alpha_i \geq 0.$$



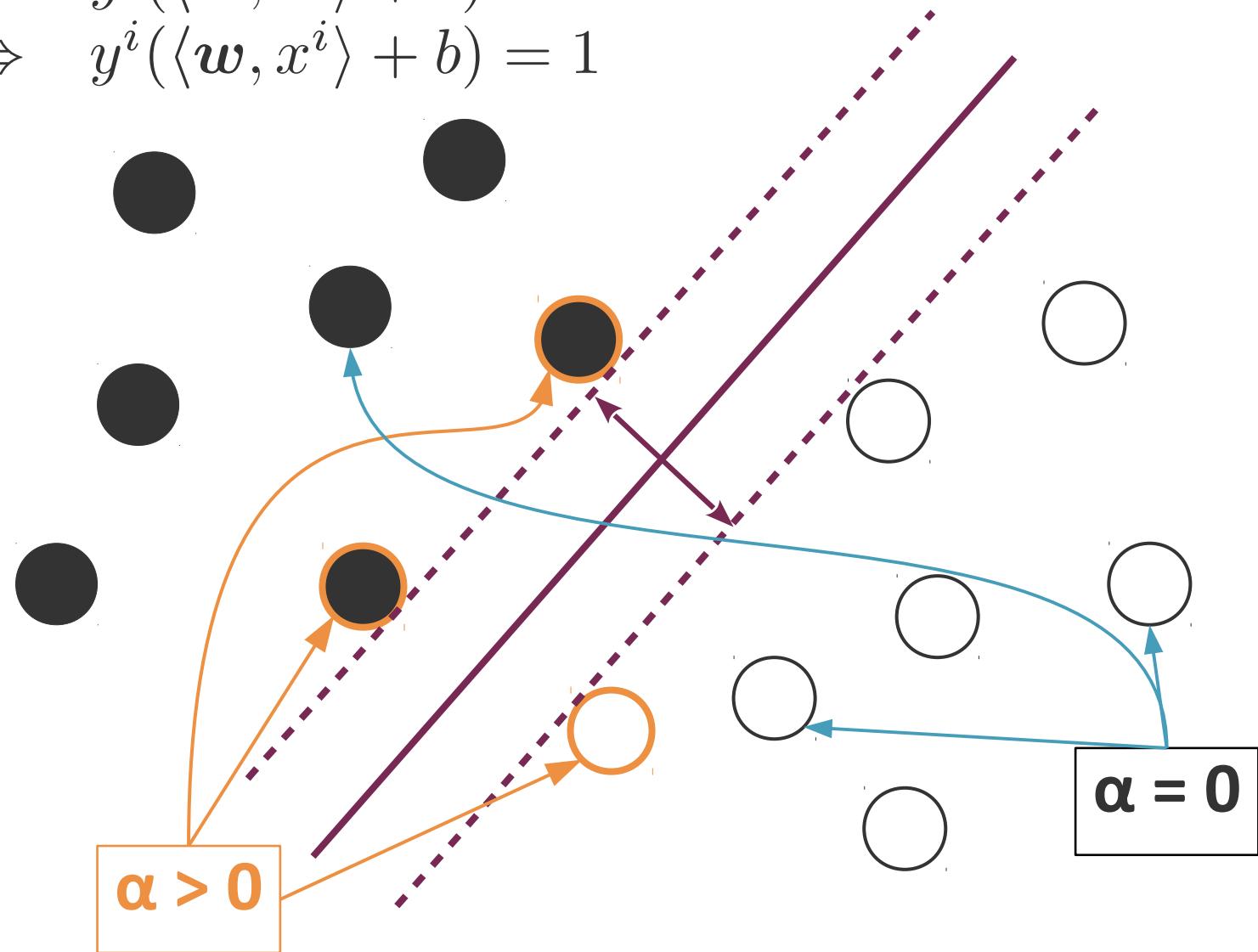
$$\textbf{Case 1: } \alpha_i = 0 \Rightarrow y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) > 1$$

$$\textbf{Case 2: } \alpha_i > 0 \Rightarrow y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) = 1$$

Support vectors

$$\alpha_i = 0 \Rightarrow y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) > 1$$

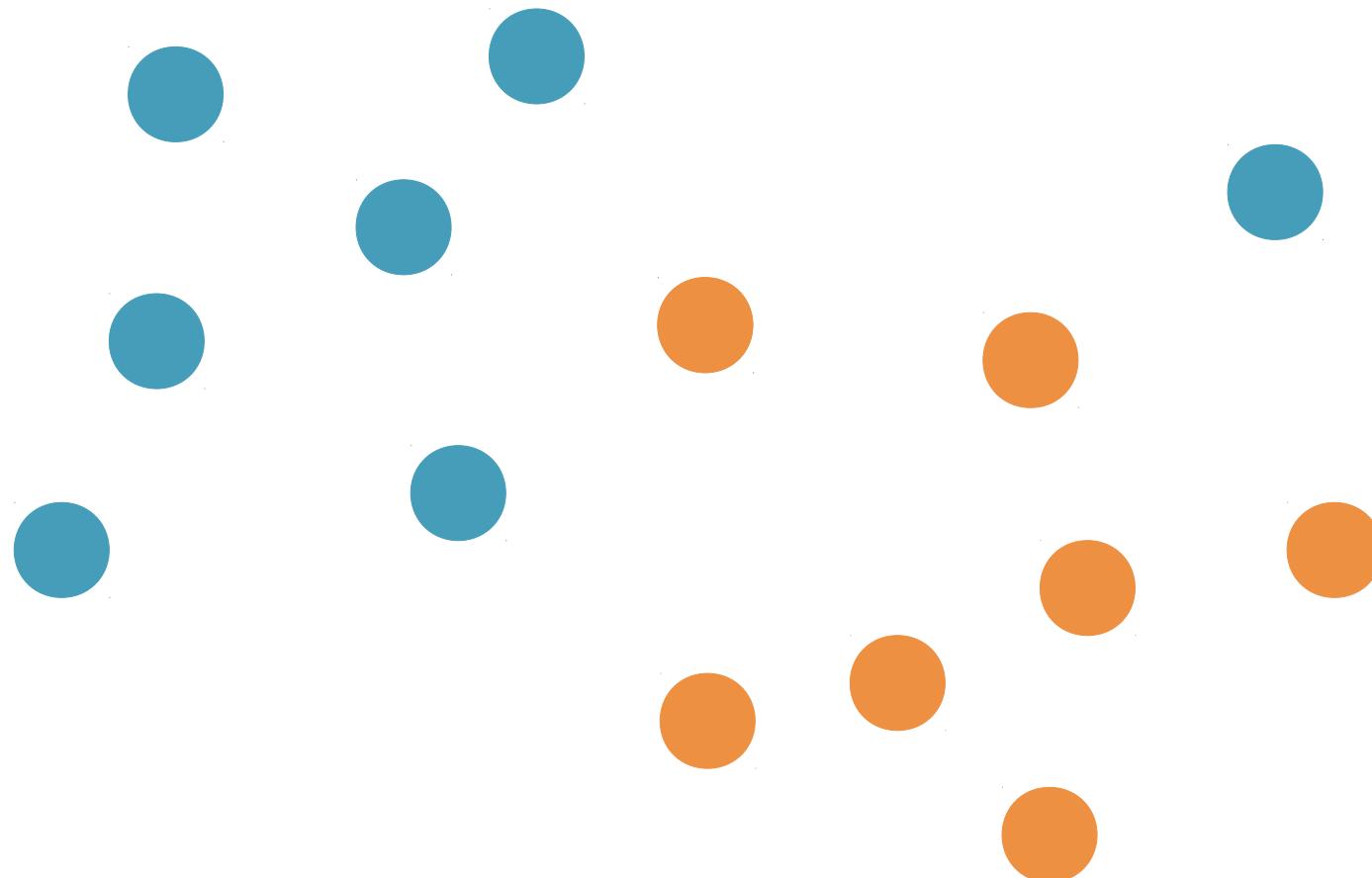
$$\alpha_i > 0 \Rightarrow y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) = 1$$



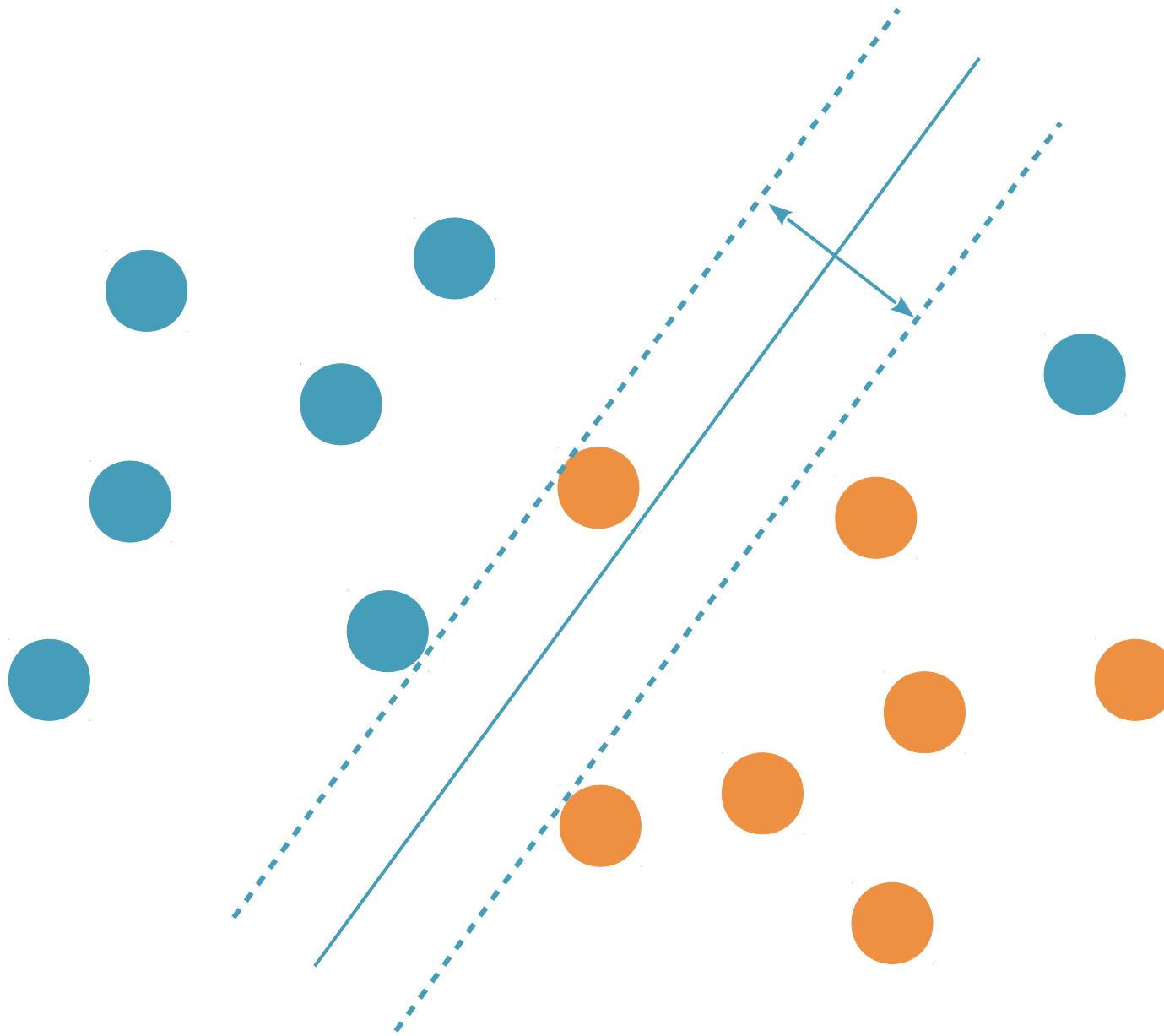
The non-linearly separable case: soft-margin SVMs.

Soft-margin SVMs

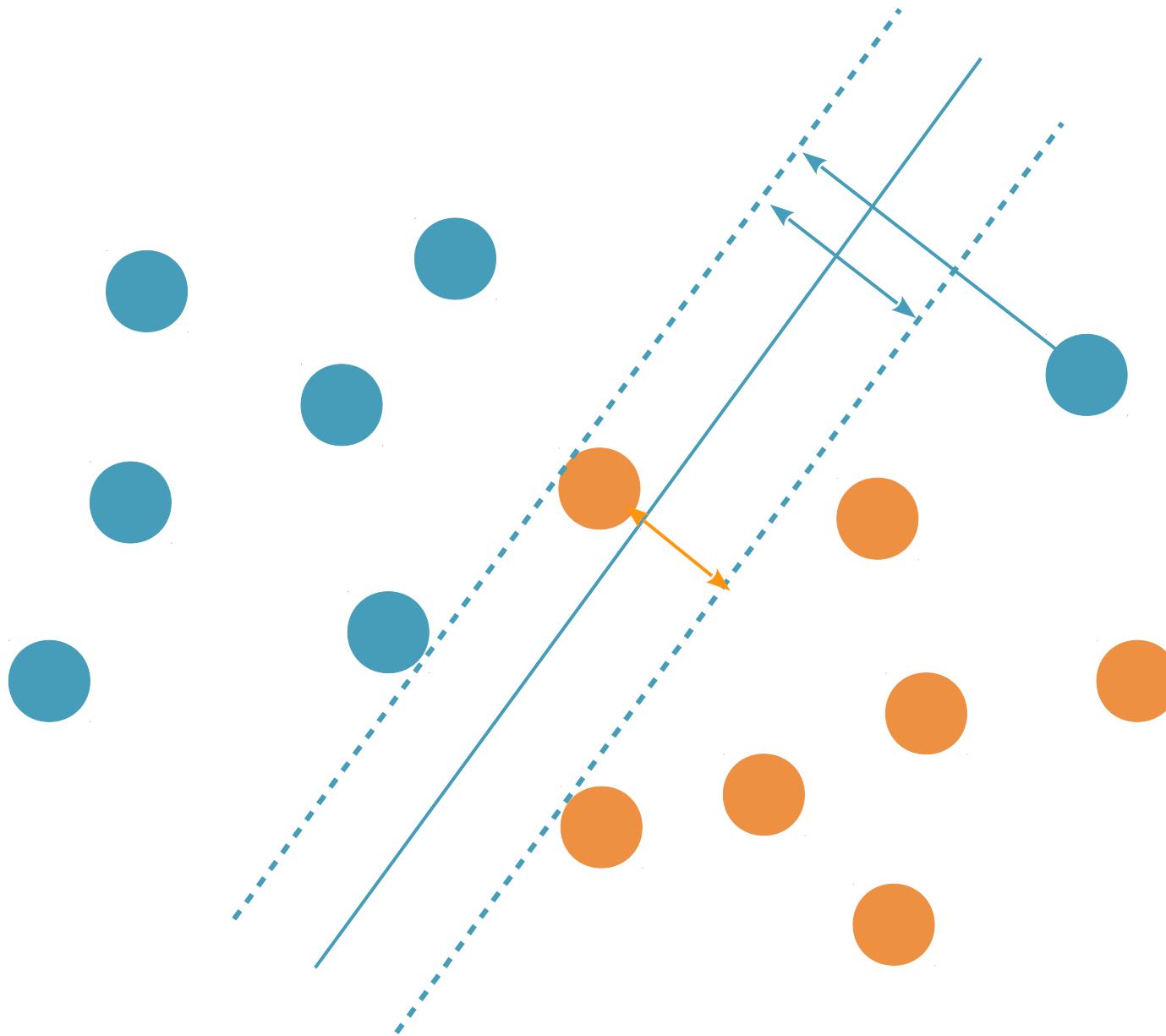
What if the data are not linearly separable?



Soft-margin SVMs



Soft-margin SVMs



Soft-margin SVMs

- Find a trade-off between **large margin** and **few errors**.

$$\min_f \left(\frac{1}{\text{margin}(f)} + C \times \text{error}(f) \right)$$

What does this remind you of?

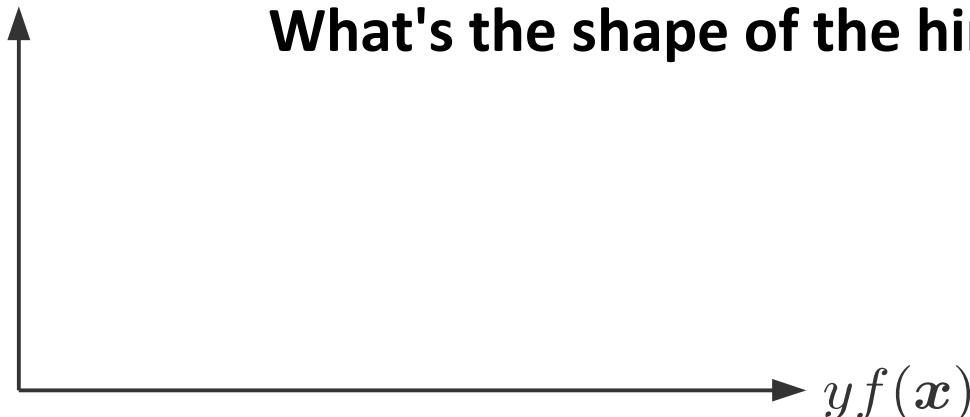
SVM error: hinge loss

- We want for all i: $y^i f(\mathbf{x}^i) \geq 1$

- **Hinge loss** function:

$$\begin{aligned} l_{\text{hinge}}(u, y) &= \max(1 - yu, 0) \\ &= \begin{cases} 0 & \text{if } yu \geq 1 \\ 1 - yu & \text{otherwise} \end{cases} \end{aligned}$$

$$l_{\text{hinge}}(y, f(\mathbf{x}))$$



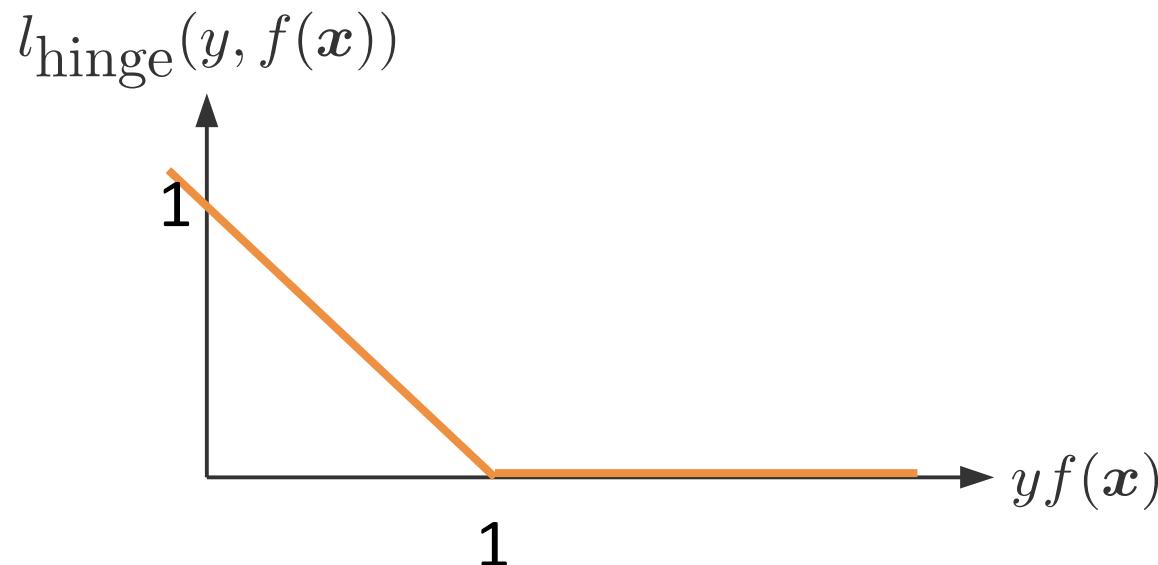
What's the shape of the hinge loss?

SVM error: hinge loss

- We want for all i: $y^i f(\mathbf{x}^i) \geq 1$

- **Hinge loss** function:

$$\begin{aligned} l_{\text{hinge}}(u, y) &= \max(1 - yu, 0) \\ &= \begin{cases} 0 & \text{if } yu \geq 1 \\ 1 - yu & \text{otherwise} \end{cases} \end{aligned}$$



Soft-margin SVMs

- Find a trade-off between **large margin** and **few errors**.

$$\min_f \left(\frac{1}{\text{margin}(f)} + C \times \text{error}(f) \right)$$

- **Error:**

$$\begin{cases} 0 & \text{if } y(\langle \mathbf{w}, \mathbf{x} \rangle + b) \geq 1 \\ 1 - y(\langle \mathbf{w}, \mathbf{x} \rangle + b) & \text{otherwise.} \end{cases}$$

- The **soft-margin SVM** solves:

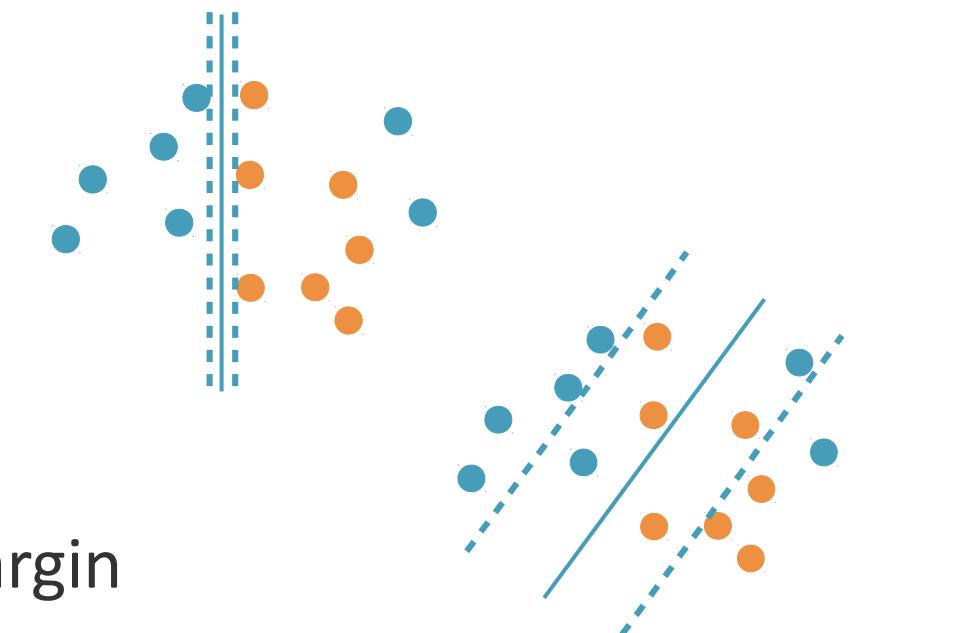
$$\arg \min_{\mathbf{w}, b} \left(\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b)) \right)$$

The C parameter

$$\min_f \left(\frac{1}{\text{margin}(f)} + C \times \text{error}(f) \right)$$

- **Large C**

makes few errors



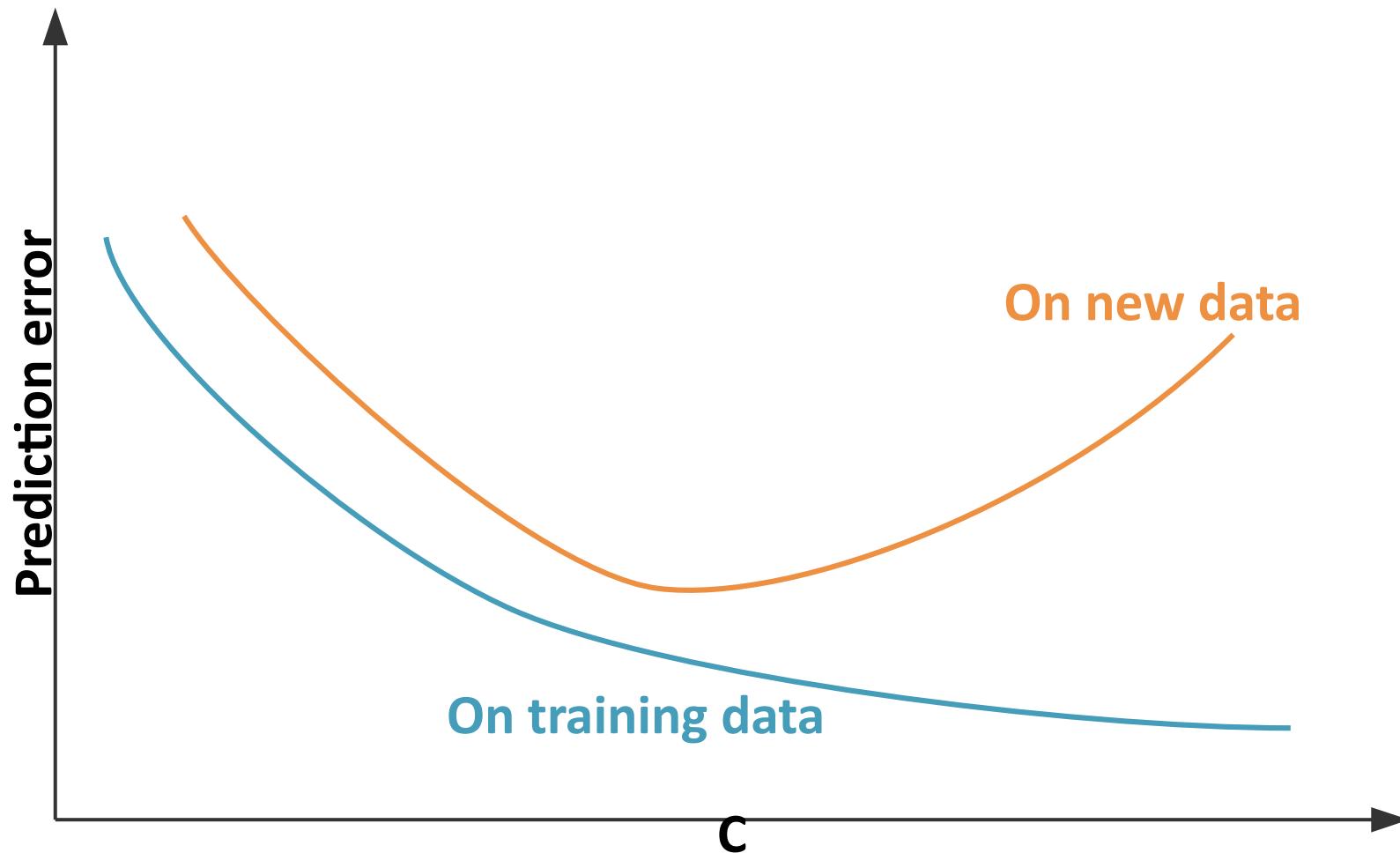
- **Small C**

ensures a large margin

- **Intermediate C**

finds a tradeoff

It is important to control C



Slack variables

$$\arg \min_{\mathbf{w}, b} \left(\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b)) \right)$$

is equivalent to:

$$\begin{array}{lll} \arg \min & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s. t. } y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) & \geq & 1 - \xi_i \\ \xi_i & \geq & 0 \quad \forall i \end{array}$$

slack variable:
distance btw y.f(x) and 1

Lagrangian of the soft-margin SVM

- **Primal**

$$\begin{array}{lll} \arg \min & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s. t. } & y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \quad \forall i \end{array}$$

- **Lagrangian**

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \mathbf{r}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1 + \xi_i) - \sum_i r_i \xi_i$$

- Min the Lagrangian (partial derivatives in \mathbf{w} , b , $\boldsymbol{\xi}$)

$$\begin{aligned} \mathbf{w} &= \sum_i \alpha_i y^i \mathbf{x}^i \\ \sum_i \alpha_i y^i &= 0 \\ \alpha_i &= C - r_i \end{aligned}$$

- **KKT conditions**

$$\begin{aligned} \alpha_i (y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) - 1 + \xi_i) &= 0 \\ r_i \xi_i &= 0 \end{aligned}$$

Dual formulation of the soft-margin SVM

- **Dual:** Maximize

$$q(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

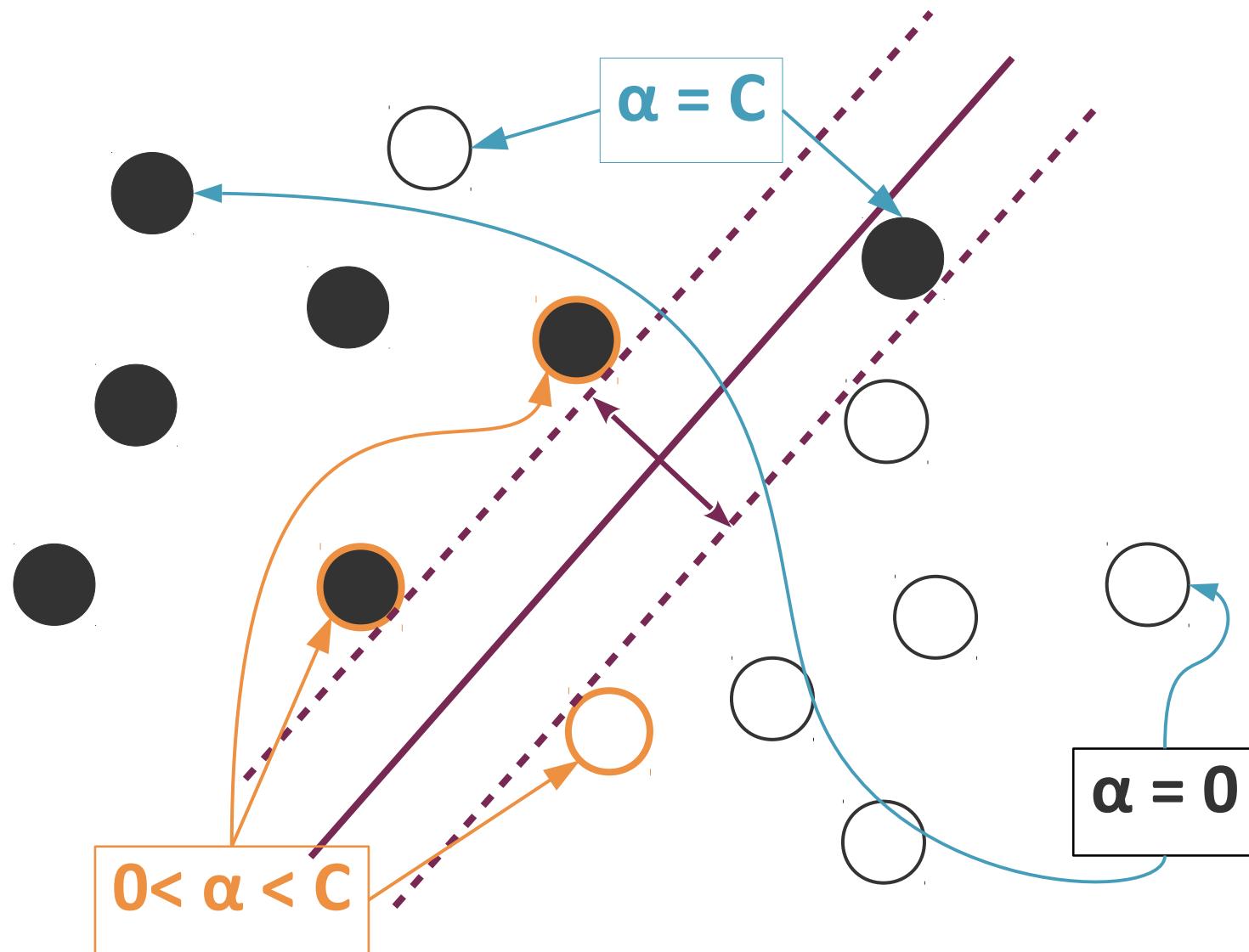
- under the constraints

$$\begin{cases} 0 \leq \alpha_i \leq C & \text{for } i = 1, \dots, n \\ \sum_{i=1}^n \alpha_i y^i = 0 \end{cases}$$

- **KKT conditions:**

$$\begin{array}{lll} \alpha_i = 0 & \Rightarrow & y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) > 1 & \text{"easy"} \\ \alpha_i = C & \Rightarrow & y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) < 1 & \text{"hard"} \\ 0 < \alpha_i < C & \Rightarrow & y^i (\langle \mathbf{w}, \mathbf{x}^i \rangle + b) = 1 & \text{"somewhat hard"} \end{array}$$

Support vectors of the soft-margin SVM



Primal vs. dual

- What is the dimension of the primal problem?

$$\arg \min_{\mathbf{w}, b} \left(\sum_{i=1}^n l_{\text{hinge}}(\langle \mathbf{w}, \mathbf{x}^i \rangle + b, y^i) + \lambda \|\mathbf{w}\|^2 \right)$$

- What is the dimension of the dual problem?

$$\arg \max_{\boldsymbol{\alpha}} L(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$0 \leq \alpha_i \leq C \text{ and } \sum_{i=1}^n \alpha_i y^i = 0$$

Primal vs. dual

- Primal: (\mathbf{w}, b) has **dimension $(p+1)$** .

$$\arg \min_{\mathbf{w}, b} \left(\sum_{i=1}^n l_{\text{hinge}}(\langle \mathbf{w}, \mathbf{x}^i \rangle + b, y^i) + \lambda \|\mathbf{w}\|^2 \right)$$

Favored if the data is **low-dimensional**.

- Dual: α has **dimension n** .

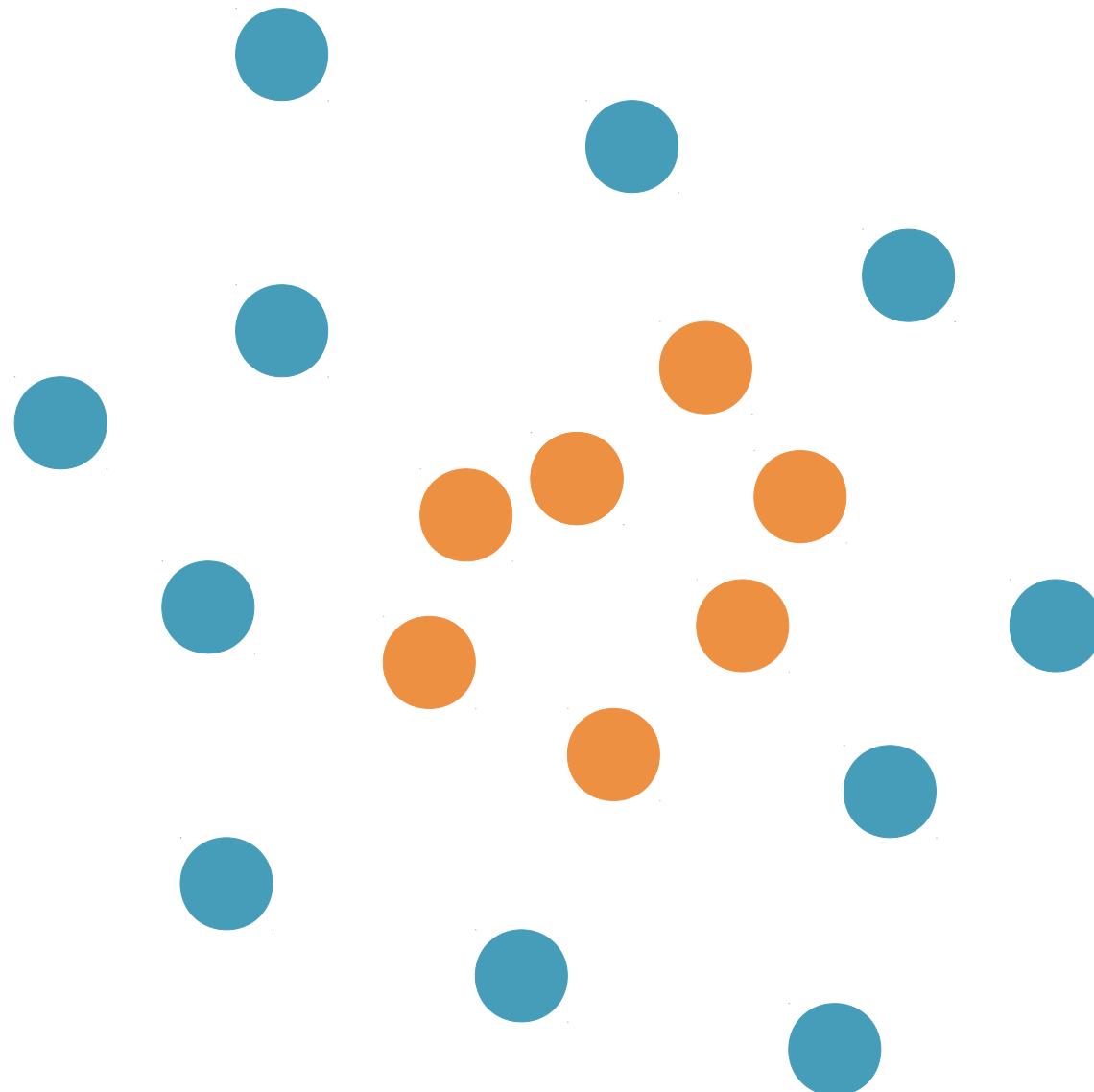
$$\arg \max_{\alpha} L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

$$0 \leq \alpha_i \leq C \text{ and } \sum_{i=1}^n \alpha_i y^i = 0$$

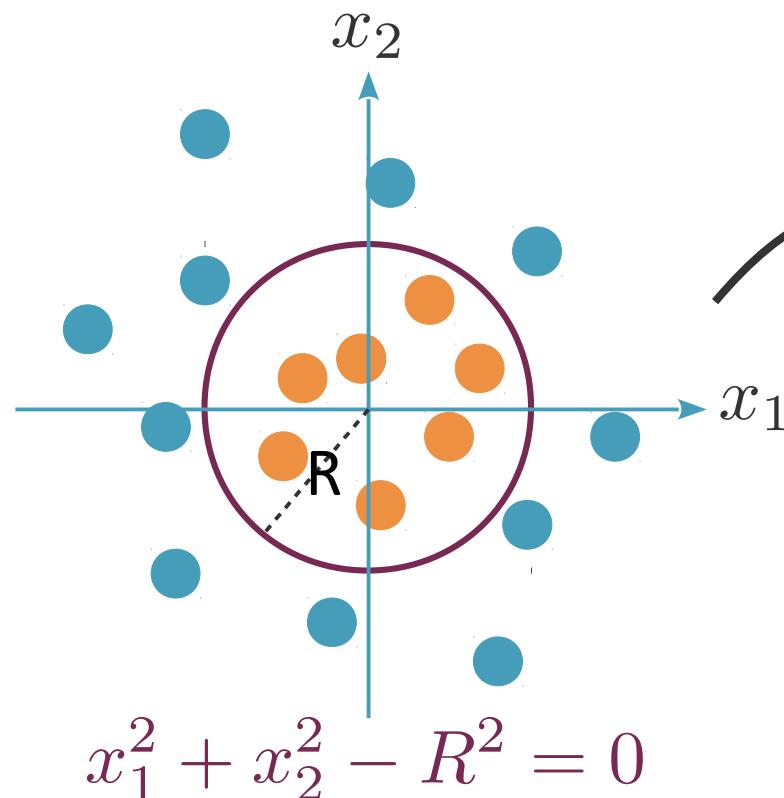
Favored if there is **little data** available.

The non-linear case: kernel SVMs.

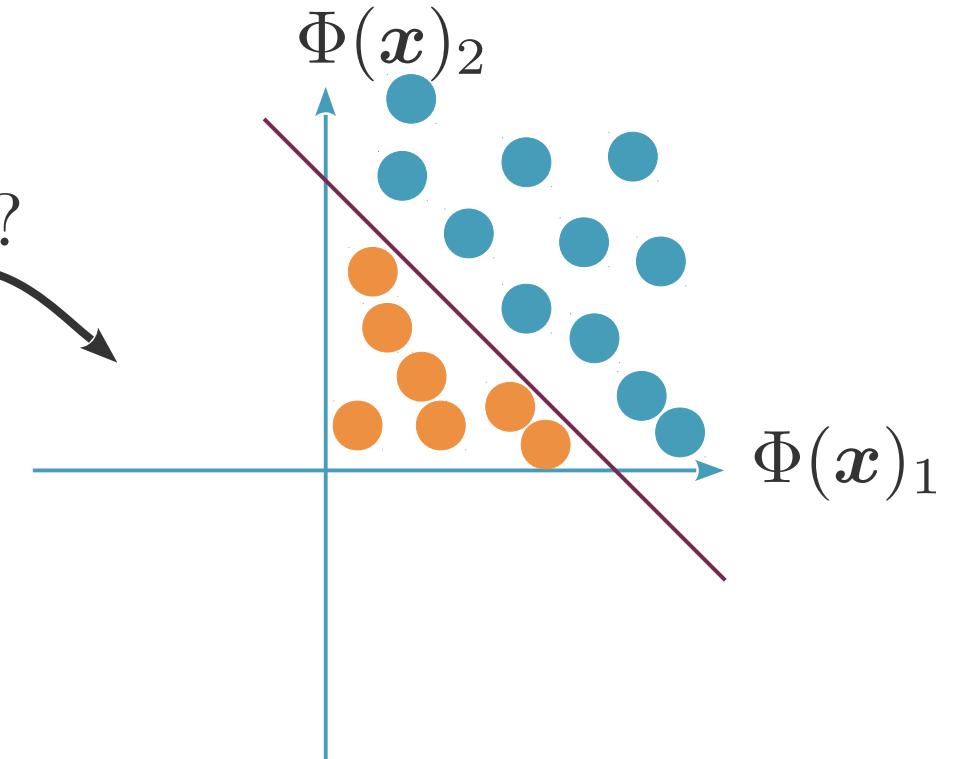
Non-linear SVMs



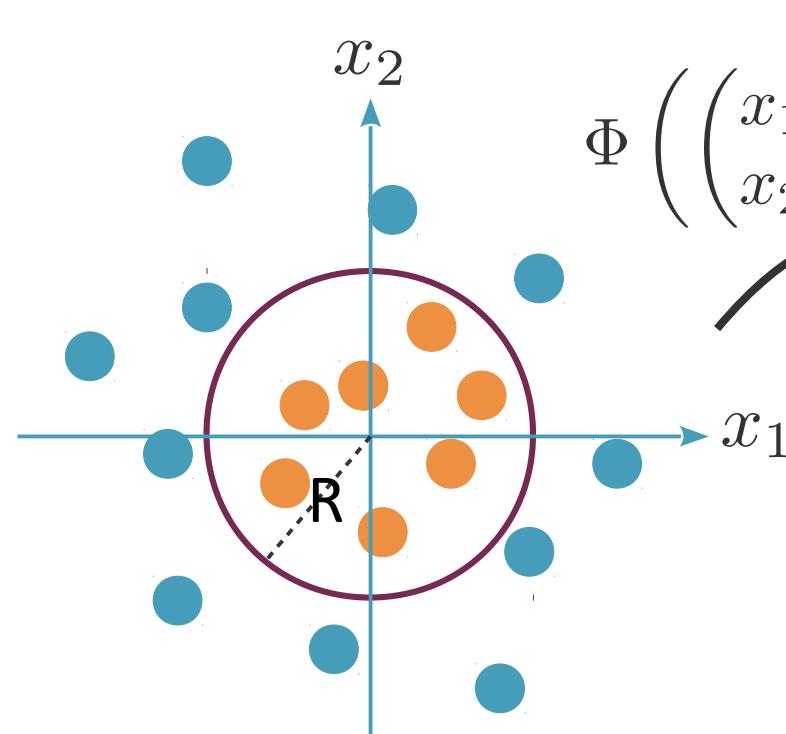
Non-linear mapping to a feature space



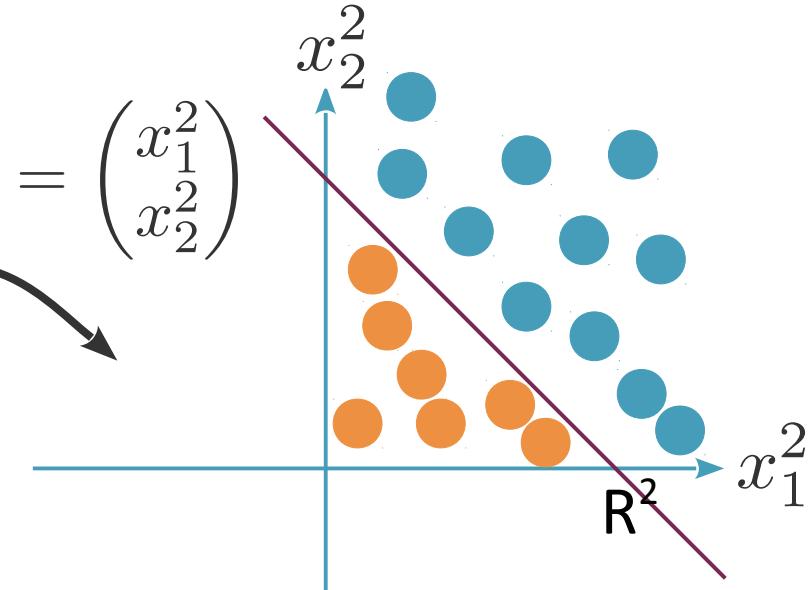
$\Phi?$



Non-linear mapping to a feature space



$$\Phi \left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right) = \begin{pmatrix} x_1^2 \\ x_2^2 \end{pmatrix}$$



SVM in the feature space

- **Train:**

$$\arg \max_{\alpha} L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \Phi(\mathbf{x}^i), \Phi(\mathbf{x}^j) \rangle_{\mathcal{H}}$$

under the constraints

$$0 \leq \alpha_i \leq C \text{ and } \sum_{i=1}^n \alpha_i y^i = 0$$

- **Predict** with the decision function

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y^i \langle \Phi(\mathbf{x}^i), \Phi(\mathbf{x}) \rangle_{\mathcal{H}} + b^*$$

Kernels

For a given mapping

$$\Phi : \mathcal{X} \mapsto \mathcal{H}$$

from the space of objects \mathcal{X} to some Hilbert space \mathcal{H} , the **kernel** between two objects x and x' is the inner product of their images in the feature spaces.

$$\forall x, x' \in \mathcal{X}, K(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$$

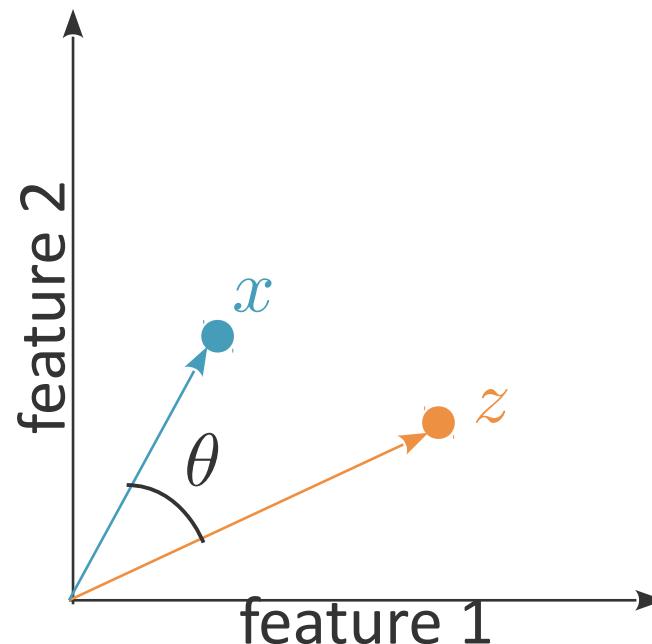
$$K(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$$

- E.g. $K(x, x') = \left\langle \begin{pmatrix} x_1^2 \\ x_2^2 \end{pmatrix}, \begin{pmatrix} {x'_1}^2 \\ {x'_2}^2 \end{pmatrix} \right\rangle = x_1^2 {x'_1}^2 + x_2^2 {x'_2}^2$
- **Kernels allow us to formalize the notion of similarity.**

Dot product and similarity

- Normalized dot product = cosine

$$\rho(\mathbf{x}, \mathbf{z}) = \frac{\sum_{j=1}^p x_j z_j}{\sqrt{\sum_{j=1}^p x_j^2} \sqrt{\sum_{j=1}^p z_j^2}} = \frac{\langle \mathbf{x}, \mathbf{z} \rangle}{\|\mathbf{x}\| \cdot \|\mathbf{z}\|} = \cos \theta$$



Kernel trick

- Many linear algorithms (in particular, linear SVMs) can be performed in the feature space H **without explicitly computing the images $\phi(x)$** , but instead by computing kernels $K(x, x')$
- It is sometimes easy to compute kernels which correspond to large-dimensional feature spaces: **$K(x, x')$ is often much simpler to compute than $\phi(x)$.**

SVM in the feature space

- **Train:**

$$\arg \max_{\alpha} L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \Phi(\mathbf{x}^i), \Phi(\mathbf{x}^j) \rangle_{\mathcal{H}}$$

under the constraints

$$0 \leq \alpha_i \leq C \text{ and } \sum_{i=1}^n \alpha_i y^i = 0$$

- **Predict** with the decision function

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y^i \langle \Phi(\mathbf{x}^i), \Phi(\mathbf{x}) \rangle_{\mathcal{H}} + b^*$$

SVM with a kernel

- **Train:**

$$\arg \max_{\alpha} L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j K(\mathbf{x}^i, \mathbf{x}^j)$$

under the constraints

$$0 \leq \alpha_i \leq C \text{ and } \sum_{i=1}^n \alpha_i y^i = 0$$

- **Predict** with the decision function

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y^i K(\mathbf{x}^i, \mathbf{x}) + b^*$$

Which functions are kernels?

- A function $K(\mathbf{x}, \mathbf{x}')$ defined on a set X is a **kernel** iff it exists a Hilbert space H and a mapping $\phi: X \rightarrow H$ such that, for any \mathbf{x}, \mathbf{x}' in X :

$$K(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_H$$

- A function $K(\mathbf{x}, \mathbf{x}')$ defined on a set X is **positive definite** iff it is **symmetric** and satisfies:

$$\forall N \in \mathbb{N}, \forall (\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N) \in \mathcal{X}^N \text{ and } (a_1, a_2, \dots, a_N) \in \mathbb{R}^N$$

$$\sum_{i=1}^N \sum_{j=1}^N a_i a_j K(\mathbf{x}^i, \mathbf{x}^j) \geq 0.$$

- Theorem [Aronszajn, 1950]: **K is a kernel iff it is positive definite.**

Positive definite matrices

- Have a **unique Cholesky decomposition** $K = LL^\top$
L: lower triangular, with positive elements on the diagonal
- **Sesquilinear form is an inner product**

$$x, x' \mapsto x^\top K x'$$

- **conjugate symmetry** $x^\top K x' = (x'^\top K x)^*$
- **linearity in the first argument**

$$a(x^\top K x') = (ax)^\top K x'$$

$$(x + z)^\top K x' = x^\top K x' + z^\top K x'$$

- **positive definiteness**

$$x^\top K x' \geq 0$$

$$x^\top K x = 0 \Rightarrow x = 0$$

Polynomial kernels

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2 \quad \Phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix} \in \mathbb{R}^3$$

Compute $K(\mathbf{x}, \mathbf{x}')$



Polynomial kernels

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2 \quad \Phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix} \in \mathbb{R}^3$$

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= x_1^2 {x'_1}^2 + 2x_1x_2x'_1x'_2 + x_2^2 {x'_2}^2 \\ &= \langle \mathbf{x}, \mathbf{x}' \rangle^2 \end{aligned}$$

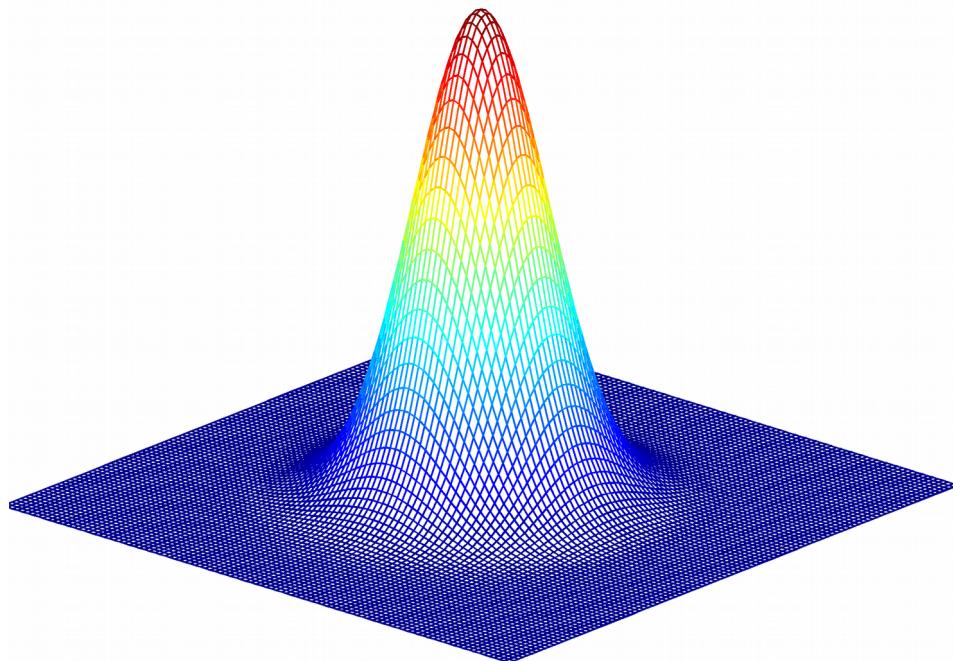
More generally, for $\mathcal{X} = \mathbb{R}^p$

$$K(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle + 1)^d$$

is an inner product in a feature space of all monomials of degree up to d .

Gaussian kernel

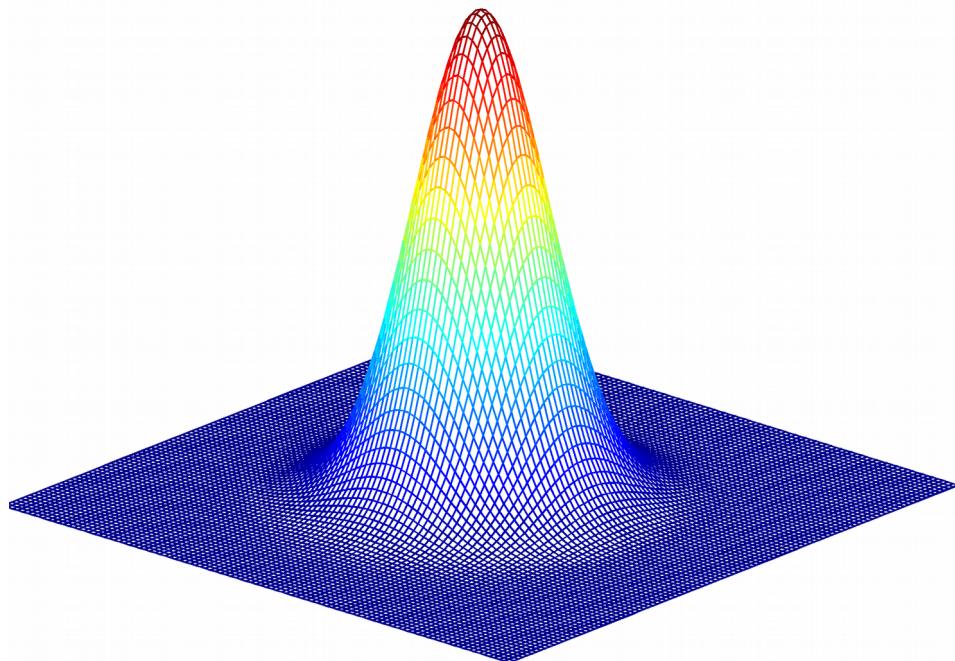
$$K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$



**What is the dimension of
the feature space?**

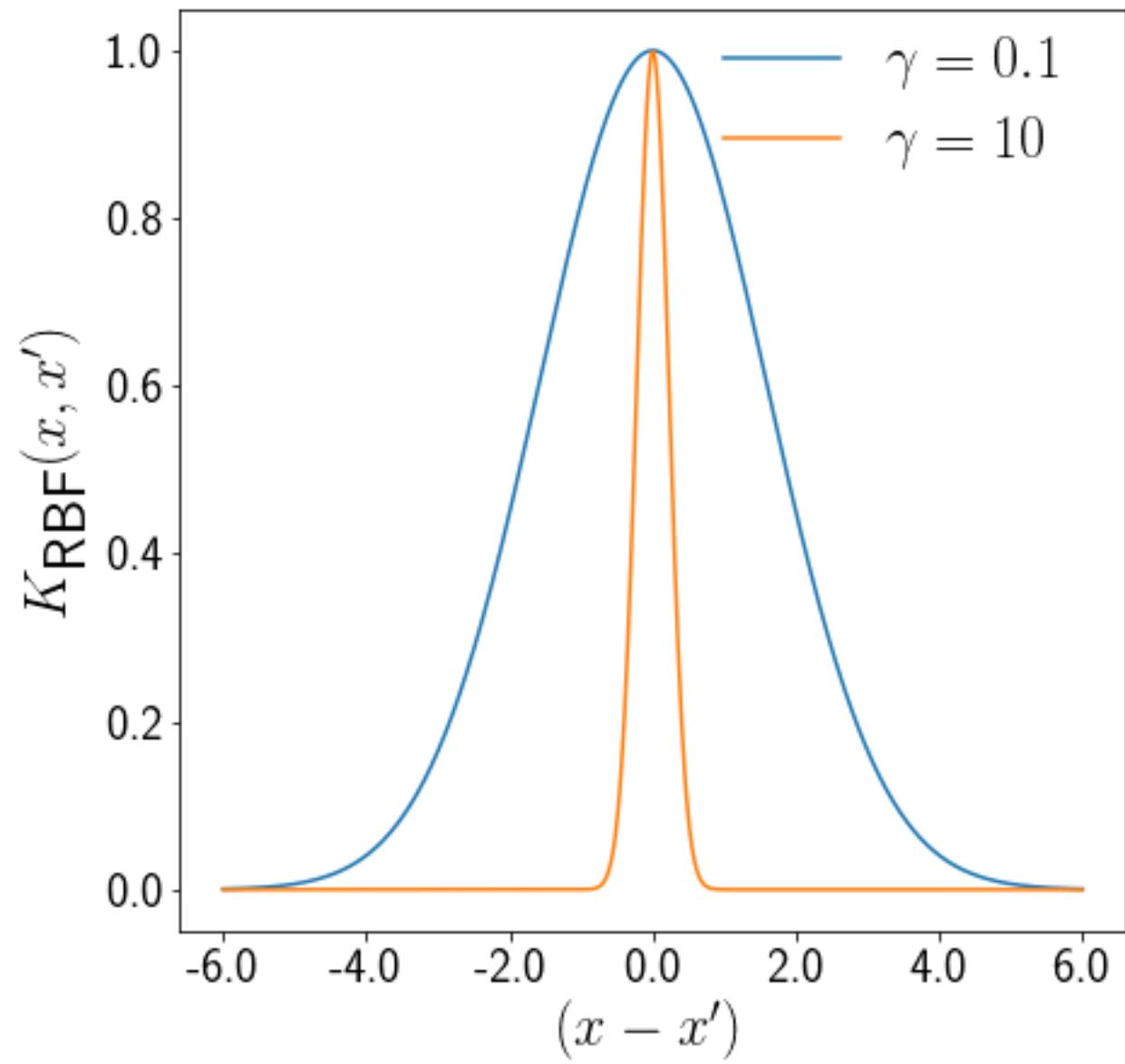
Gaussian kernel

$$K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

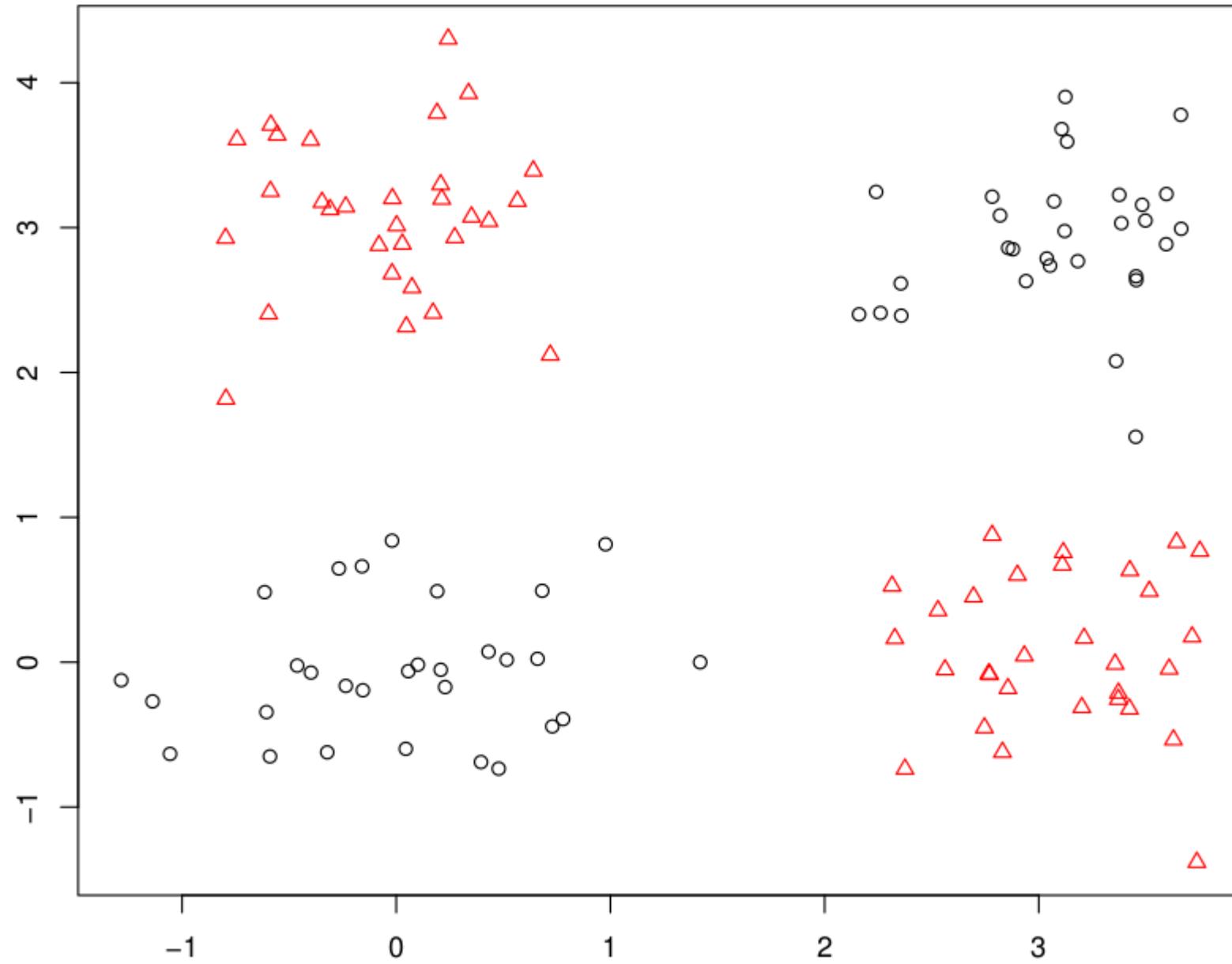


The feature space has
infinite dimension.

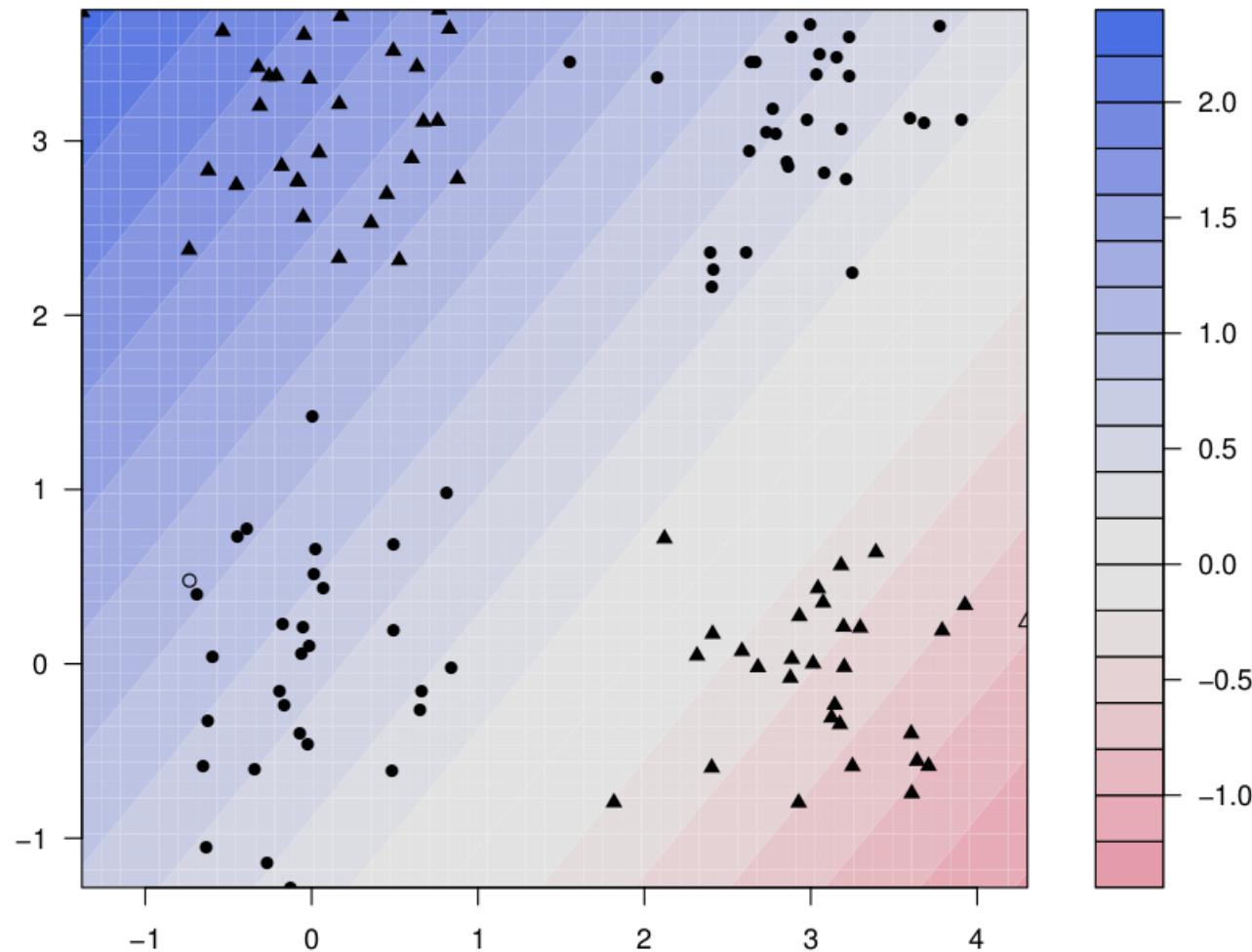
$$\begin{aligned} K(x, x') &= \exp\left(\frac{-1}{2\sigma^2}\|x\|^2\right) \exp\left(\frac{1}{\sigma^2}\langle x, x' \rangle\right) \exp\left(\frac{-1}{2\sigma^2}\|x'\|^2\right) \\ &= f(x) \sum_{r=0}^{+\infty} \frac{\langle x, x' \rangle^r}{\sigma^{2r} r!} f(x') \end{aligned}$$



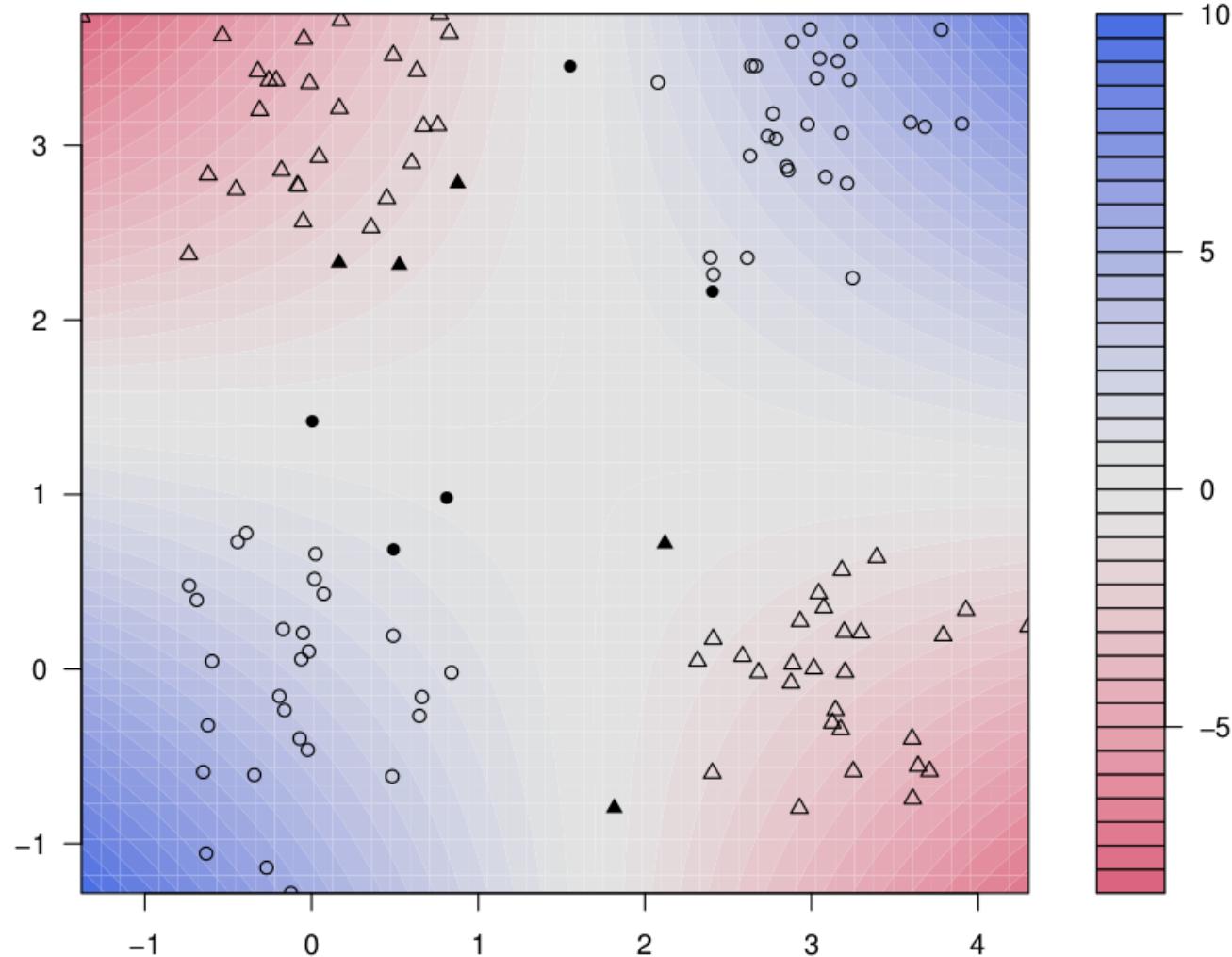
Toy example



Toy example: linear SVM



Toy example: polynomial SVM (d=2)



Kernels for strings

Protein sequence classification

Goal: predict which proteins are secreted or not, based on their sequence.

- Secreted proteins:

MASKATLLLAFATLLFATCIARHQQRQQQQNQCQLQNIEA...

MARSSLFTFLCLAVFINGCLSQIEQQSPWEFQGSEVW...

MALHTVLIMLSLLPMLEAQNPEHANITIGEPITNETLGWL...

...

- Non-secreted proteins:

MAPPSSVFAEVVPQAQPVLVFKLIADFREDPDPRKVNLGVG...

MAHTLGLTQPNSTEPHKISFTAKEIDVIEWKGDIILVVG...

MSISESYAKEIKTAFRQFTDFPIEGEQFEDFLPIIGNP ..

...

Substring-based representations

- Represent strings based on the presence/absence of substrings of fixed length.

$$\Phi(x) = \{\Phi_u(x)\}_{u \in \mathcal{A}^k}$$

The diagram illustrates the components of a substring-based representation. It shows the formula $\Phi(x) = \{\Phi_u(x)\}_{u \in \mathcal{A}^k}$. Two arrows point to specific parts of the formula: one from a question mark to the term $\Phi_u(x)$, and another from the text "Strings of length k" to the set \mathcal{A}^k .

Substring-based representations

- Represent strings based on the presence/absence of substrings of fixed length.

$$\Phi(x) = \{\Phi_u(x)\}_{u \in \mathcal{A}^k}$$

- Number of occurrences of u in x : **spectrum kernel** [Leslie et al., 2002].

Substring-based representations

- Represent strings based on the presence/absence of substrings of fixed length.

$$\Phi(x) = \{\Phi_u(x)\}_{u \in \mathcal{A}^k}$$

- Number of occurrences of u in x: **spectrum kernel** [Leslie et al., 2002].
- Number of occurrences of u in x, up to m mismatches: **mismatch kernel** [Leslie et al., 2004].

Substring-based representations

- Represent strings based on the presence/absence of substrings of fixed length.

$$\Phi(x) = \{\Phi_u(x)\}_{u \in \mathcal{A}^k}$$

- Number of occurrences of u in x: **spectrum kernel** [Leslie et al., 2002].
- Number of occurrences of u in x, up to m mismatches: **mismatch kernel** [Leslie et al., 2004].
- Number of occurrences of u in x, allowing gaps, with a weight decaying exponentially with the number of gaps: **substring kernel** [Lohdi et al., 2002].

Spectrum kernel

$$K(\mathbf{x}, \mathbf{x}') = \sum_{u \in \mathcal{A}^k} \Phi_u(\mathbf{x}) \Phi_u(\mathbf{x}')$$

- **Implementation:**

- Formally, a sum over $|\mathcal{A}^k|$ terms
- **How many non-zero terms in $\Phi(\mathbf{x})$?**



Spectrum kernel

$$K(\mathbf{x}, \mathbf{x}') = \sum_{u \in \mathcal{A}^k} \Phi_u(\mathbf{x}) \Phi_u(\mathbf{x}')$$

- **Implementation:**

- Formally, a sum over $|\mathcal{A}^k|$ terms
- At most $|\mathbf{x}| - k + 1$ non-zero terms in $\Phi(\mathbf{x})$
- Hence: Computation in $O(|\mathbf{x}| + |\mathbf{x}'|)$

- **Prediction** for a new sequence \mathbf{x} :

$$f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b$$



Write $f(\mathbf{x})$ as a function of only $|\mathbf{x}| - k + 1$ weights.

Spectrum kernel

$$K(\mathbf{x}, \mathbf{x}') = \sum_{u \in \mathcal{A}^k} \Phi_u(\mathbf{x}) \Phi_u(\mathbf{x}')$$

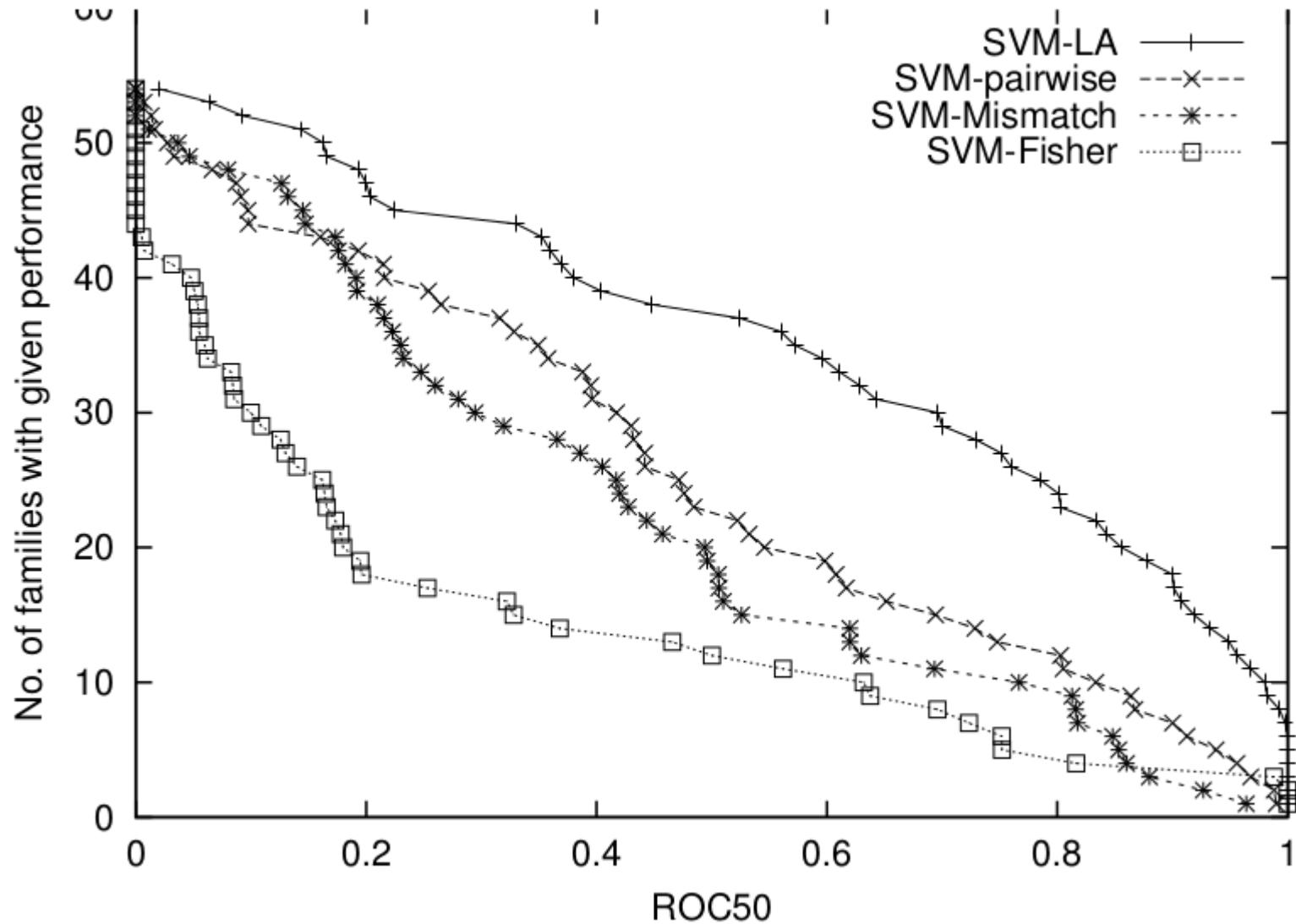
- **Implementation:**

- Formally, a sum over $|\mathcal{A}^k|$ terms
- At most $|\mathbf{x}| - k + 1$ non-zero terms in $\Phi(\mathbf{x})$
- Hence: Computation in $O(|\mathbf{x}| + |\mathbf{x}'|)$

- **Fast prediction** for a new sequence \mathbf{x} :

$$\begin{aligned} f(\mathbf{x}) &= \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b \\ &= \sum_{u \in \mathcal{A}^k} w_u \Phi_u(\mathbf{x}) + b \\ &= \sum_{j=1}^{|\mathbf{x}|-k+1} w_{x_j x_{j+1} \dots x_{j+k-1}} + b \end{aligned}$$

The choice of kernel matters

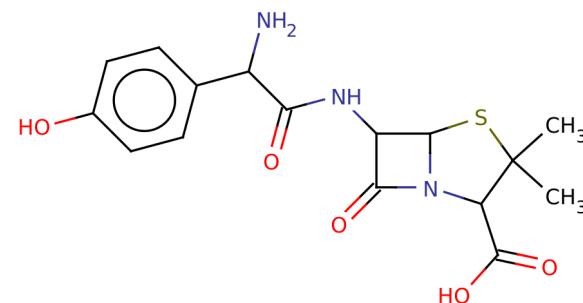


Performance of several kernels on the SCOP superfamily recognition kernel [Saigo et al., 2004]

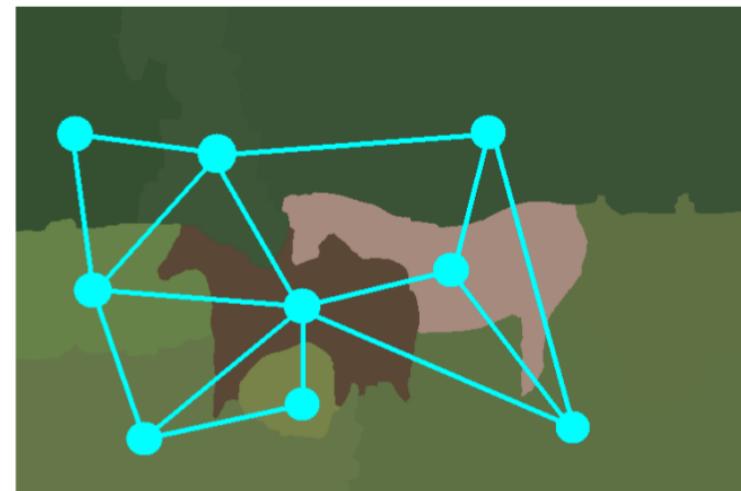
Kernels for graphs

Graph data

- Molecules

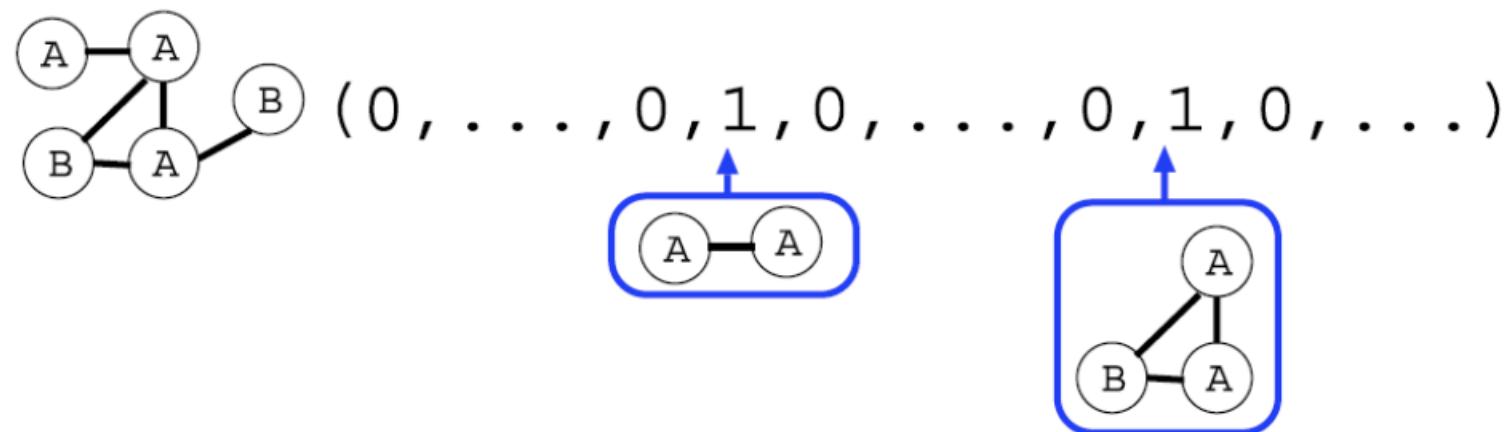


- Images



[Harchaoui & Bach, 2007]

Subgraph-based representations



Tanimoto & MinMax

- The Tanimoto and MinMax similarities are kernels

$$s(\mathbf{x}^1, \mathbf{x}^2) = \frac{\sum_{j=1}^p (\mathbf{x}_j^1 \text{ AND } \mathbf{x}_j^2)}{\sum_{j=1}^p (\mathbf{x}_j^1 \text{ OR } \mathbf{x}_j^2)}$$

$$s(\mathbf{x}^1, \mathbf{x}^2) = \frac{\sum_{j=1}^p \min(\mathbf{x}_j^1, \mathbf{x}_j^2)}{\sum_{j=1}^p \max(\mathbf{x}_j^1, \mathbf{x}_j^2)}$$

Which subgraphs to use?

- **Indexing by all subgraphs...**
 - Computing all subgraph occurrences is NP-hard.
 - Actually, finding whether a given subgraph occurs in a graph is NP-hard in general.

A Quasipolynomial Time Algorithm for Graph Isomorphism: The Details

Posted on November 12, 2015 by j2kun

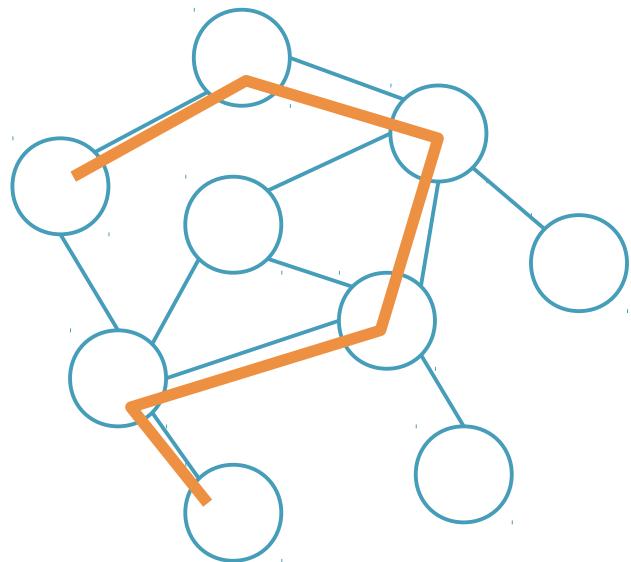
Update 2015-11-16: Laci has [posted the talk on his website](#). It's an hour and a half long, and I encourage you to watch it if you have the time 😊

<http://jeremykun.com/2015/11/12/a-quasipolynomial-time-algorithm-for-graph-isomorphism-the-details/>

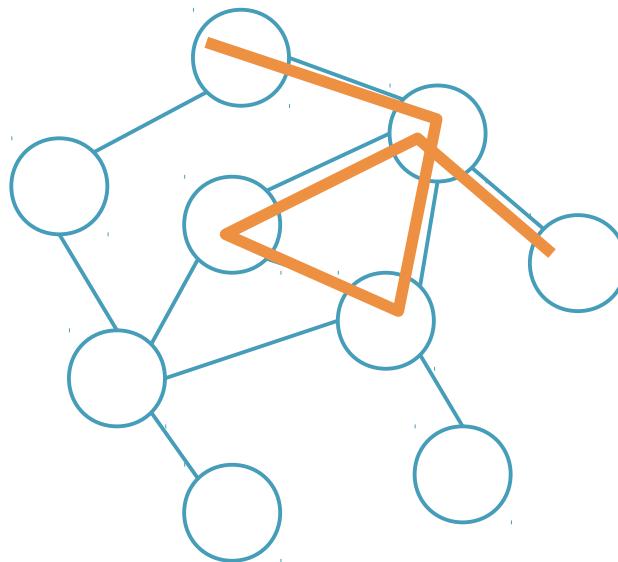
Which subgraphs to use?

- **Specific subgraphs** that lead to computationally efficient indexing:
 - Subgraphs selected based on **domain knowledge**
E.g. chemical fingerprints
 - All **frequent subgraphs** [Helma et al., 2004]
 - All **paths** up to length k [Nicholls 2005]
 - All **walks** up to length k [Mahé et al., 2005]
 - All **trees** up to depth k [Rogers, 2004]
 - All **shortest paths** [Borgwardt & Kriegel, 2005]
 - All **subgraphs up to k vertices (graphlets)** [Shervashidze et al., 2009]

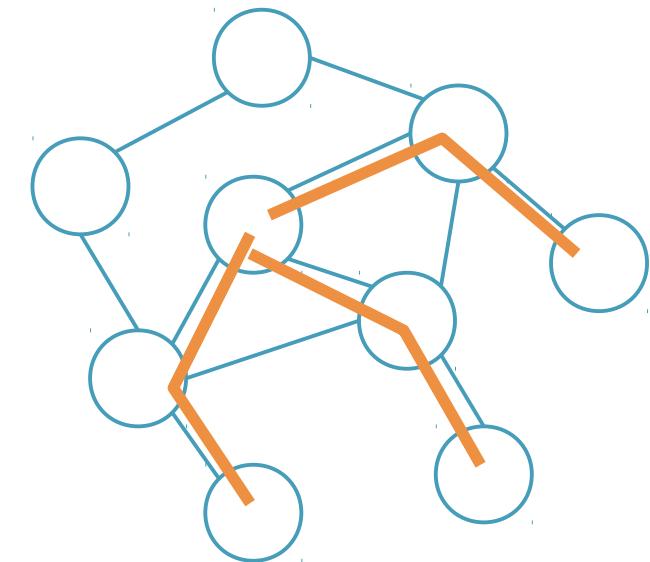
Which subgraphs to use?



Path of length 5

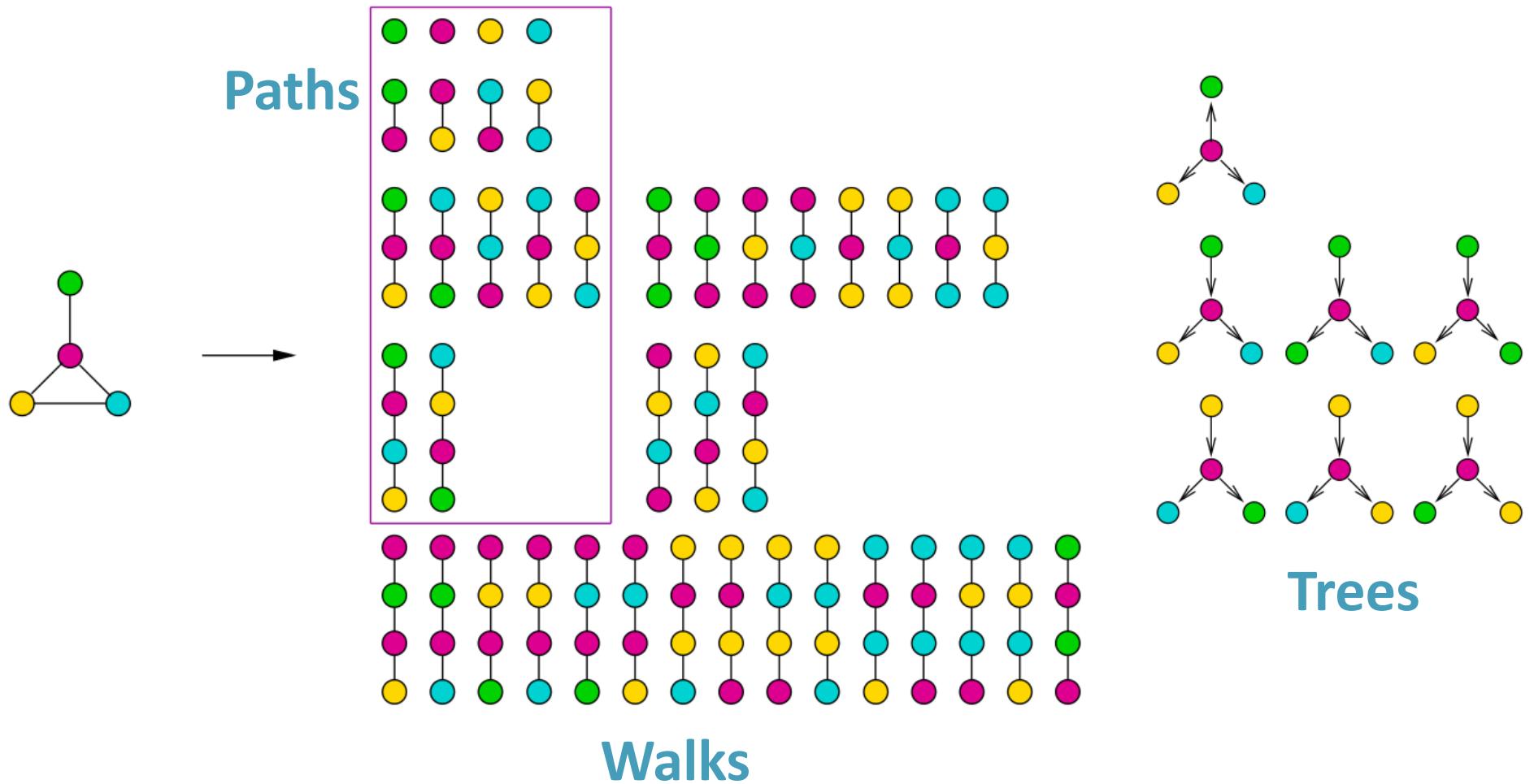


Walk of length 5



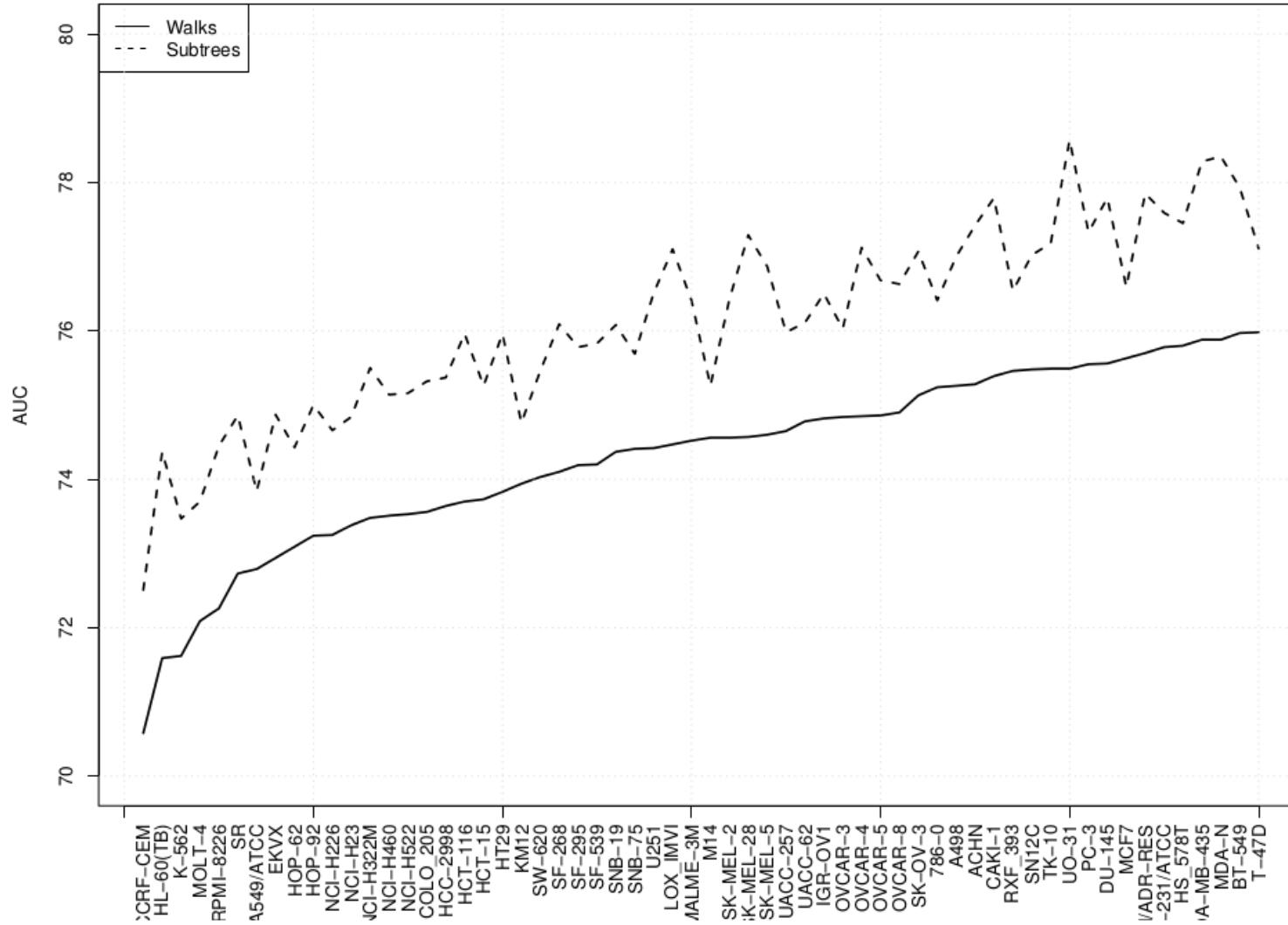
Tree of depth 2

Which subgraphs to use?



[Harchaoui & Bach, 2007]

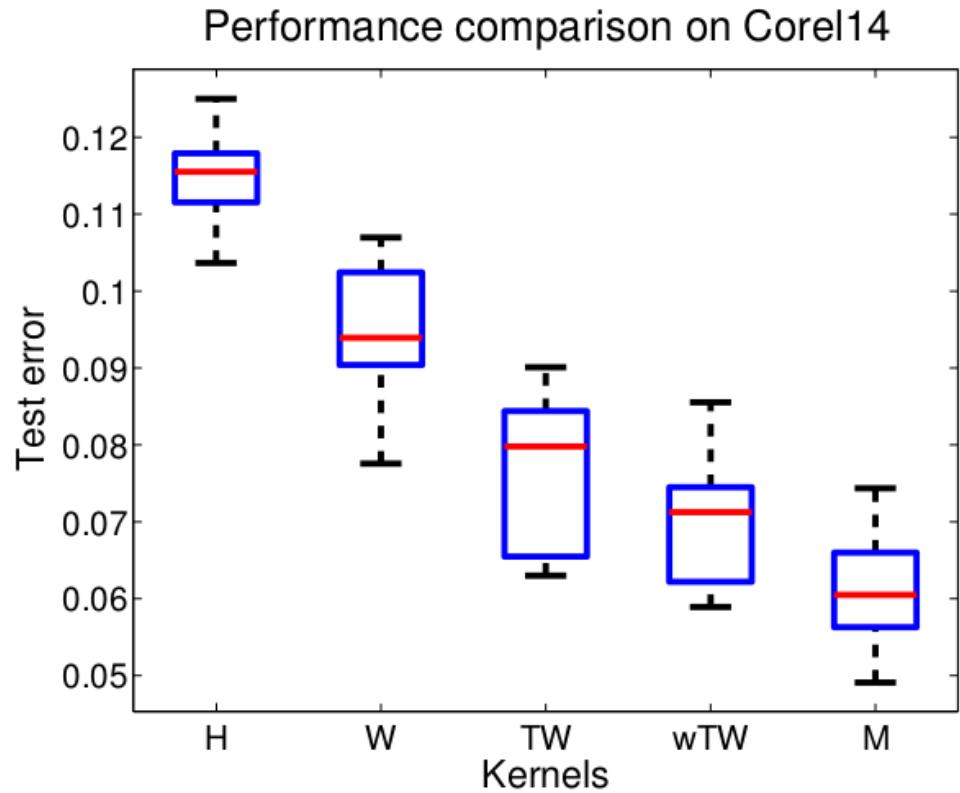
The choice of kernel matters



Predicting inhibitors for 60 cancer cell lines [Mahé & Vert, 2009]

The choice of kernel matters

- COREL14: 1400 natural images, 14 classes
- **Kernels:** histogram (H), walk kernel (W), subtree kernel (TW), weighted subtree kernel (wTW), combination (M).



Summary

- Linearly separable case: **hard-margin SVM**
- Non-separable, but still linear: **soft-margin SVM**
- Non-linear: **kernel SVM**
- Kernels for
 - real-valued data
 - strings
 - graphs.

- *A Course in Machine Learning.*
http://ciml.info/dl/v0_99/ciml-v0_99-all.pdf
 - **Soft-margin SVM** : Chap 7.7
 - **Kernel SVM**: Chap 11.1 – 11.6
- *The Elements of Statistical Learning.*
<http://web.stanford.edu/~hastie/ElemStatLearn/>
 - **Separating hyperplane**: Chap 4.5.2
 - **Soft-margin SVM**: Chap 12.1 – 12.2
 - **Kernel SVM**: Chap 12.3
 - String kernels: Chap 18.5.1
- *Learning with Kernels*
<http://agbs.kyb.tuebingen.mpg.de/lwk/>
 - **Soft-margin SVM**: Chap 1.4
 - **Kernel SVM**: Chap 1.5
 - **SVR**: Chap 1.6
 - **Kernels**: Chap 2.1
- *Convex Optimization*
<https://web.stanford.edu/~boyd/cvxbook/>
 - **SVM optimization** : Chap 8.6.1

Practical matters

- **Preparing for the exam**
 - Previous exams with solutions on the course website
- **Next week: special session! 2 x 1.5 hrs**
 - Introduction to **artificial neural networks**
 - Introduction to **deep learning** and **Tensorflow** (J. Boyd)
Jupyter notebook will be available for download
 - **Deep learning for bioimaging** (P. Naylor)

Lab

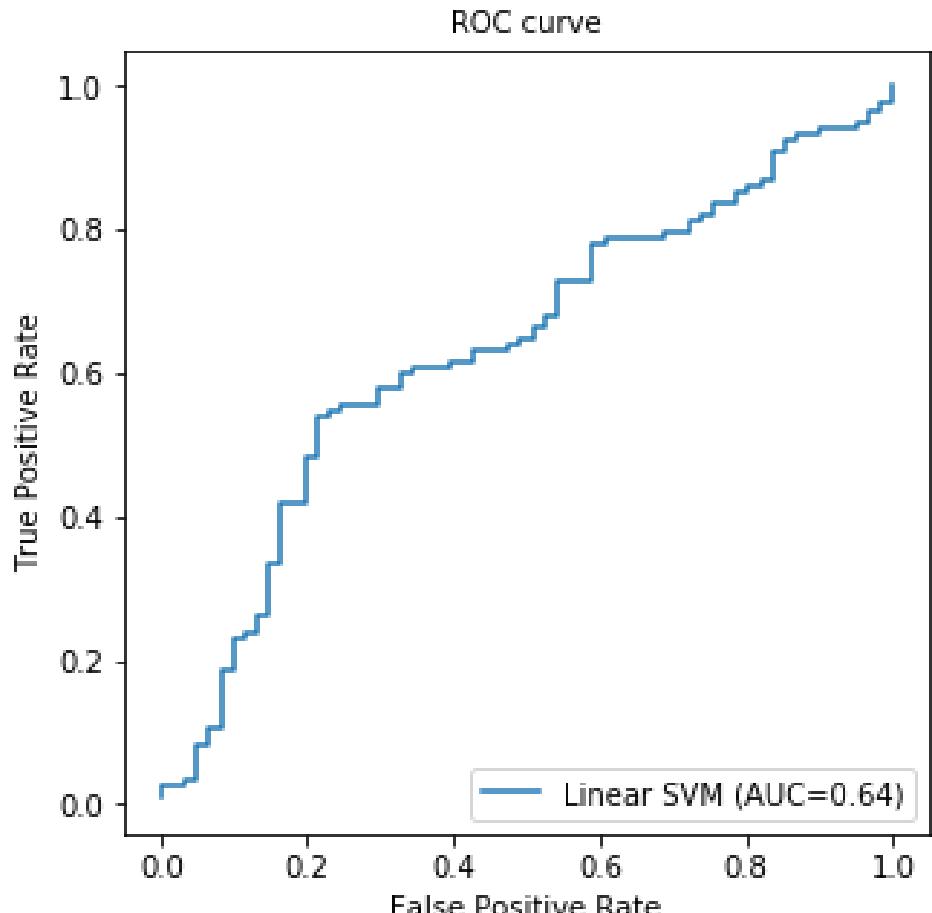
- Redefining cross_validate

```
def cross_validate(design_matrix, labels, classifier, cv_folds):  
    ...  
  
    pred = np.zeros(labels.shape) # Hold all predictions, in correct order.  
    for tr, te in cv_folds:  
        # Restrict data to train/test folds  
        Xtr = design_matrix[tr, :]  
        ytr = labels[tr]  
        Xte = design_matrix[te, :]  
  
        # Fit classifier  
        classifier.fit(Xtr, ytr)  
  
        # Compute decision function on test data  
        # TODO  
        yte = clf.decision_function(Xte)  
  
        # Update pred  
        pred[te] = yte  
  
    return pred
```

Linear SVM

The data is not easily separated by a hyperplane

Support vectors are either correctly classified points that support the margin or errors.



```
: # TODO
print('No. support vectors class 0: %d' % clf.n_support_[0])
print('No. support vectors class 1: %d' % clf.n_support_[1])
print('No. training samples: %d' % X.shape[0])
```

No. support vectors class 0: 43
No. support vectors class 1: 62
No. training samples: 183

Many support vectors suggest the data is not easy to separate and there are many errors.

Linear kernel matrix

```
kmatrix = X.dot(X.T) # TODO

# heatmap + color map
fig, ax = plt.subplots(figsize=(5, 5))
plot = ax.imshow(kmatrix, cmap=plt.cm.PuRd)

# set axes boundaries
ax.set_xlim([0, X.shape[0]]) ; ax.set_ylim([0, X.shape[0]])

# flip the y-axis
ax.invert_yaxis() ; ax.xaxis.tick_top()

# plot colorbar to the right
plt.colorbar(plot, pad=0.1, fraction=0.04)
```

No visible pattern.

Dark lines correspond to vectors with highest magnitude.



Linear kernel matrix (after feature scaling)

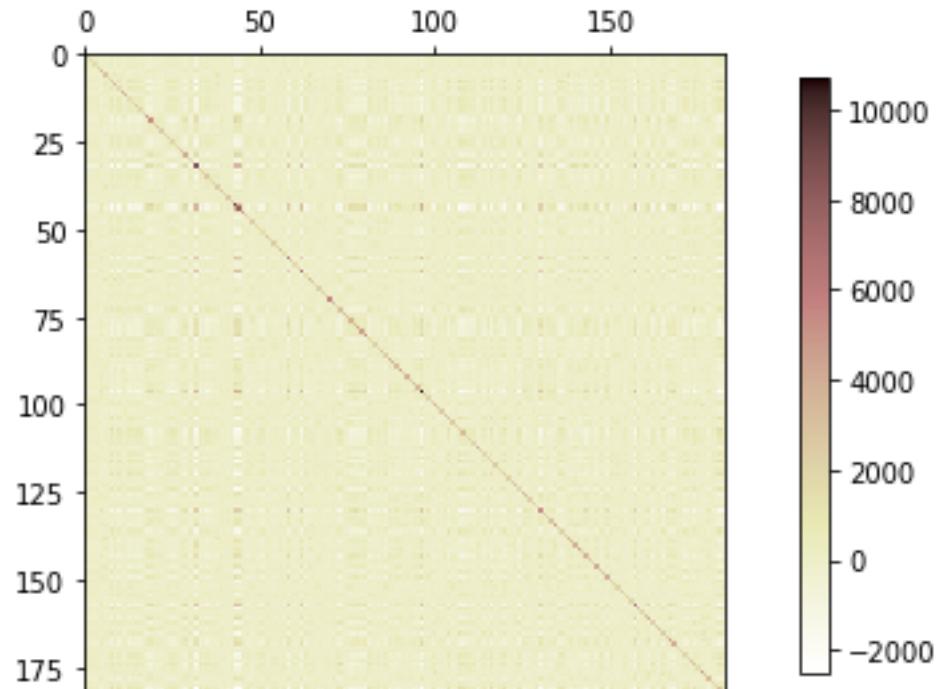
```
# TODO
X_scaled = preprocessing.scale(X)
kmatrix = X_scaled.dot(X_scaled.T)

# heatmap + color map
fig, ax = plt.subplots(figsize=(5, 5))
plot = ax.imshow(kmatrix, cmap=plt.cm.pink_r)
```

The kernel values are on a smaller scale than previously.

The diagonal emerges (the most similar sample to an observation is itself).

Many small values.



```

def cross_validate_with_scaling(design_matrix, labels, classifier, cv_folds):

    pred = np.zeros(labels.shape)
    for tr, te in cv_folds:
        # Restrict data to train/test folds
        Xtr = design_matrix[tr, :]
        ytr = labels[tr]
        Xte = design_matrix[te, :]

        # Create scaler object
        # TODO
        sc = preprocessing.StandardScaler()

        # Fit the scaler and transform training data
        # TODO
        Xtr = sc.fit_transform(Xtr)

        # Transform test data
        # TODO
        Xte = sc.transform(Xte)

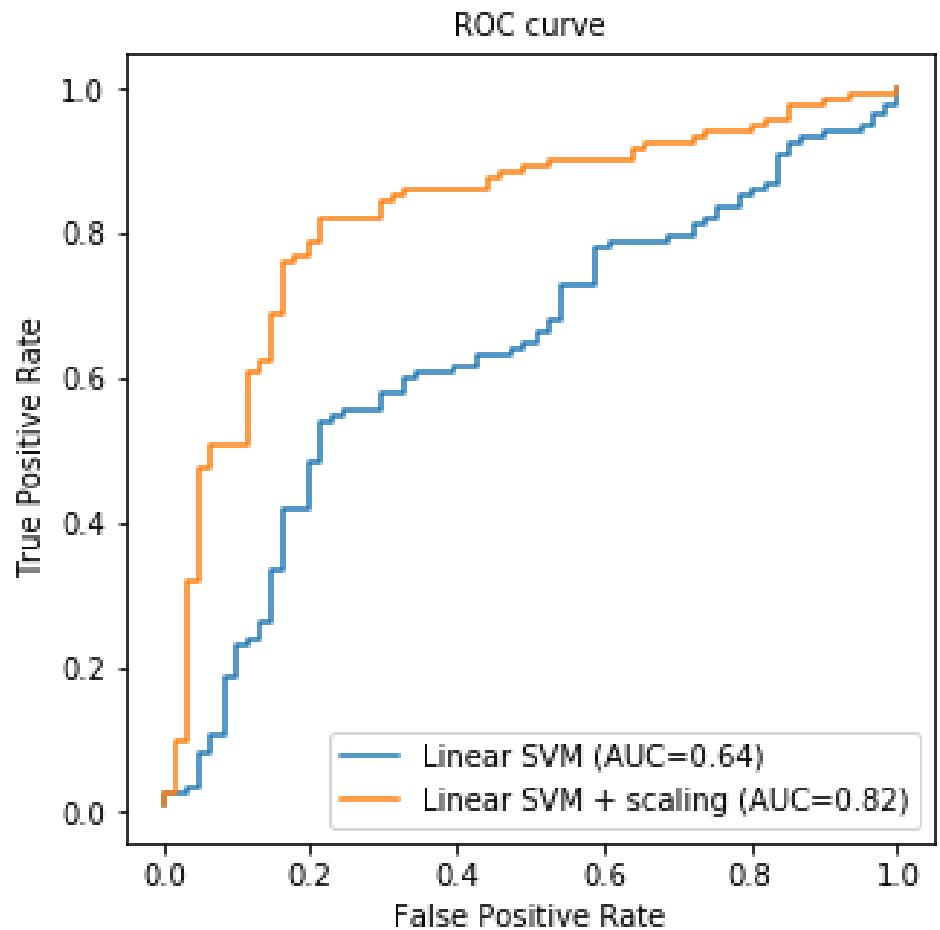
        # Fit classifier
        classifier.fit(Xtr, ytr)

        # Compute decision function on test data
        # TODO
        yte = classifier.decision_function(Xte)

        # Update pred
        # TODO
        pred[te] = yte

    return pred

```



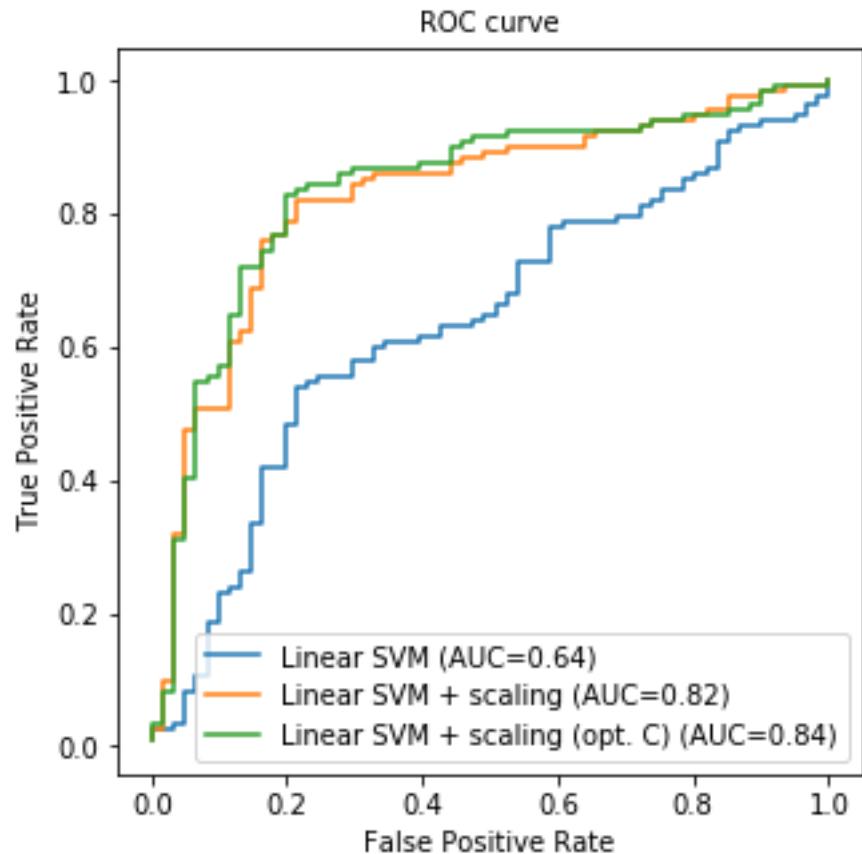
Linear SVM with optimal C

```
from sklearn.model_selection import GridSearchCV
parameters_dict = {'C': np.logspace(-4, 2, 7)}

clf_C = GridSearchCV(svm.SVC(kernel='linear'),
                     parameters_dict, scoring='roc_auc')

y_pred_linear_scaled = cross_validate_with_scaling(X, y, clf_C, folds)
```

An SVM classifier with optimized C



On each pair (tr, te):

- scaling factors are computing on Xtr
- Xtr, Xte are scaled accordingly
- for each value of C:
 - an SVM is cross-validated on Xtr_scaled
 - the best of these SVM is trained on the full Xtr_scaled and applied to Xte_scaled (this produces one prediction per data point of X)

Polynomial kernel SVM

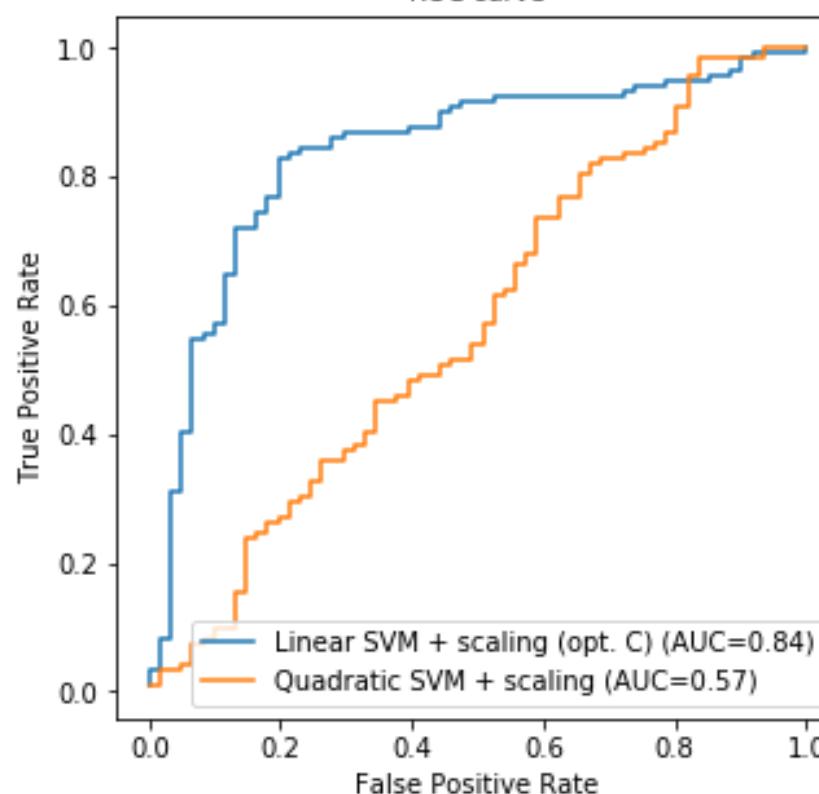
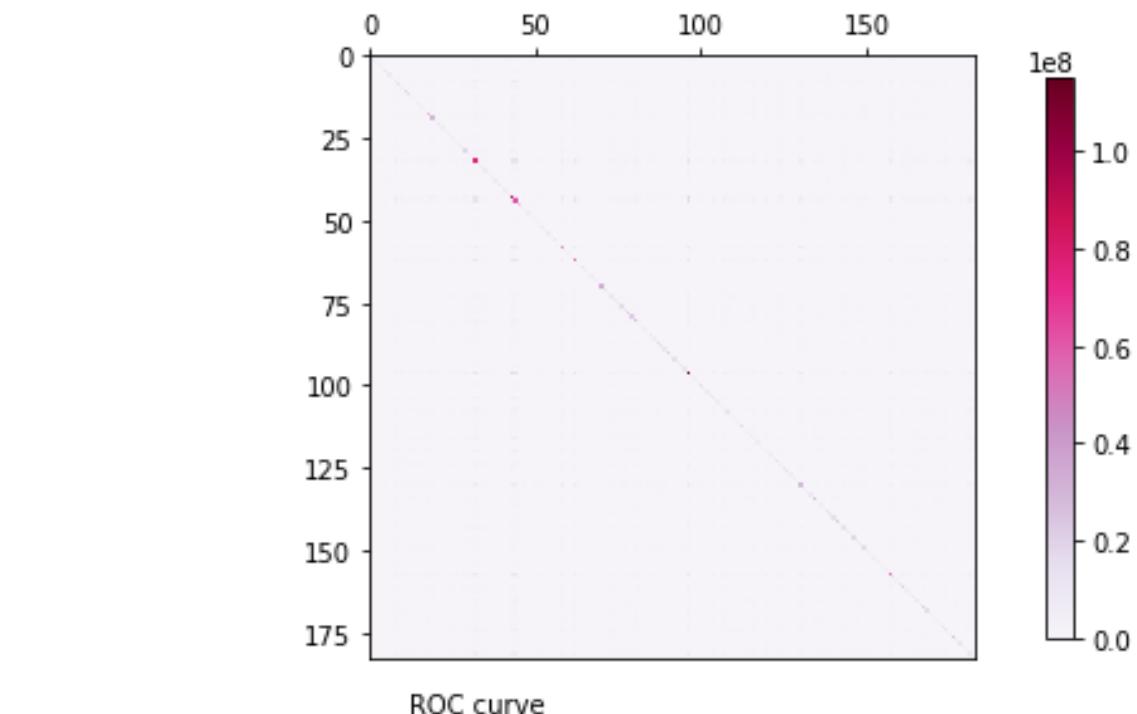
Polynomial kernel with
 $r=0$
 $d=2$

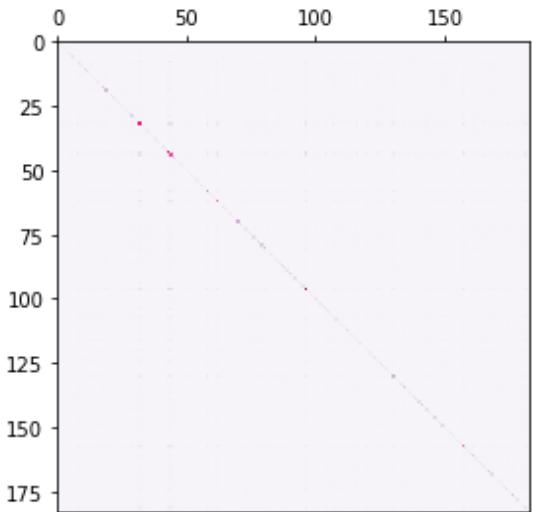
Computed on X_{scaled}

The matrix is really close to identity, nothing can be learned.

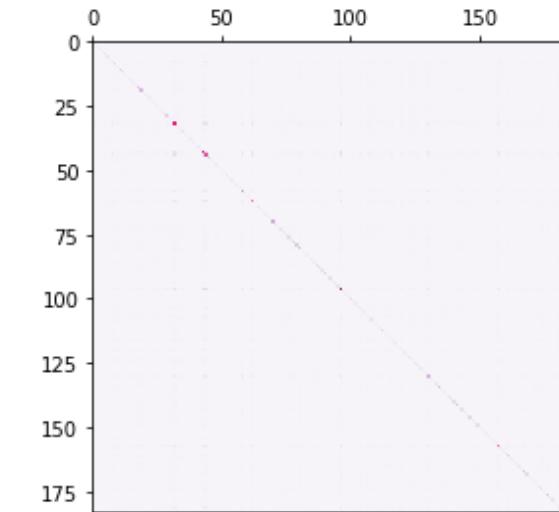
This gets worse if you increase d .

Changing r can give us a more reasonable matrix.

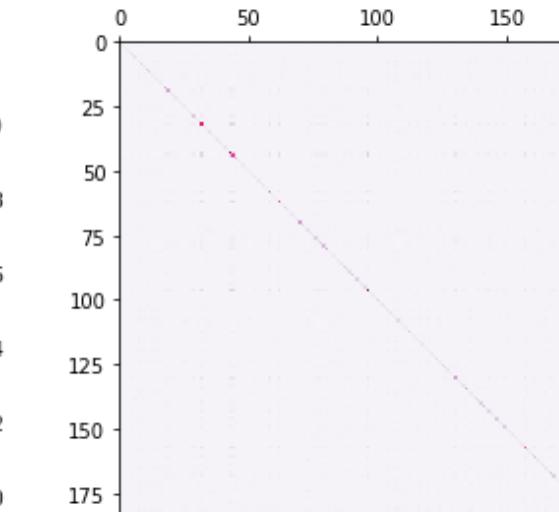




$r=10$

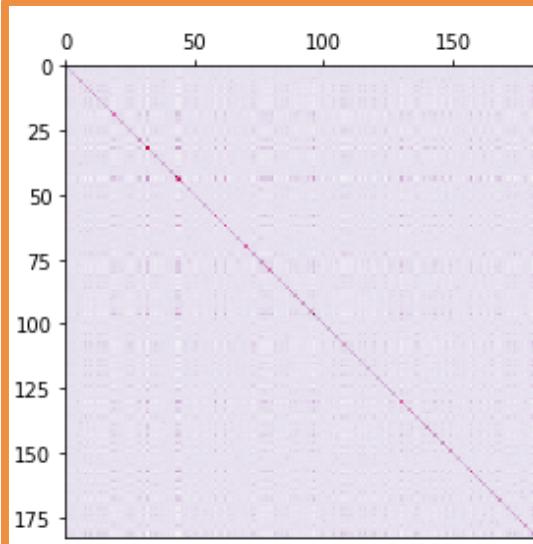


$r=100$

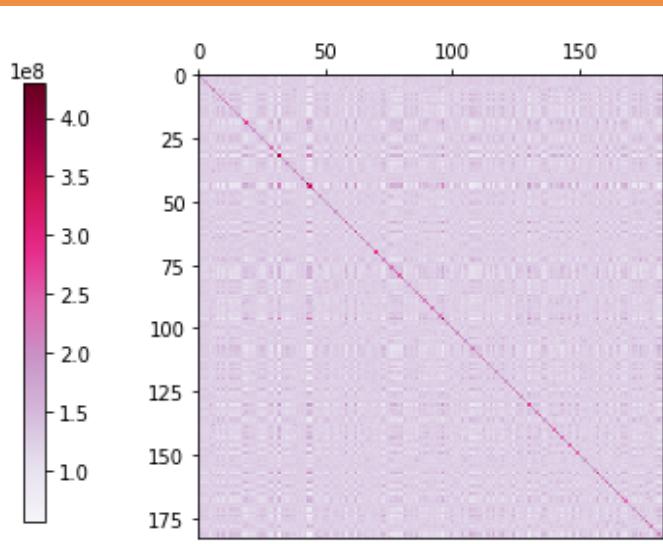


$r=1000$

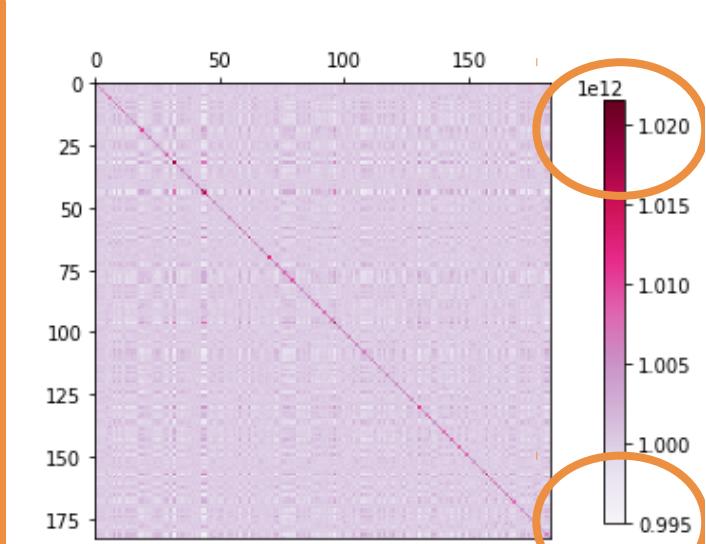
The kernel matrix is almost the identity matrix



$r=10000$

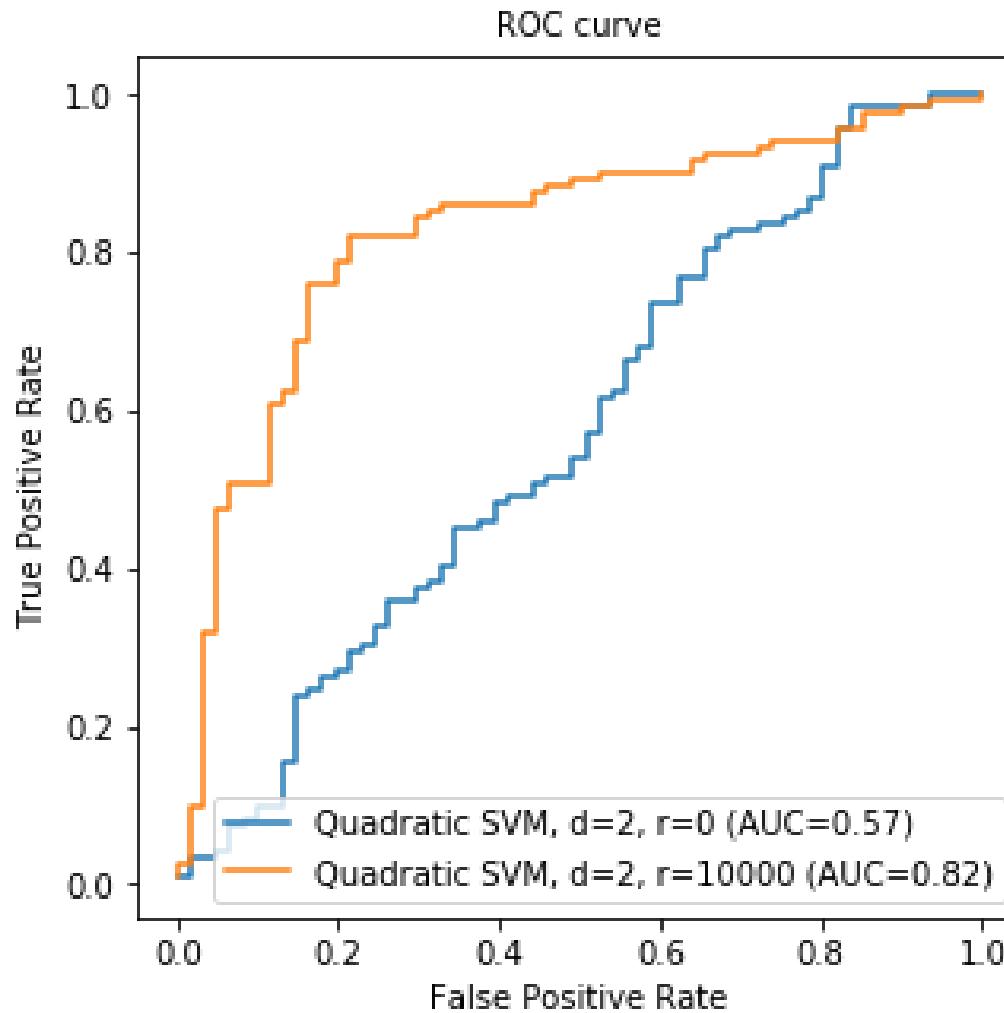


$r=100000$



$r=1000000$
Almost all 1s

Reasonable range of values for r



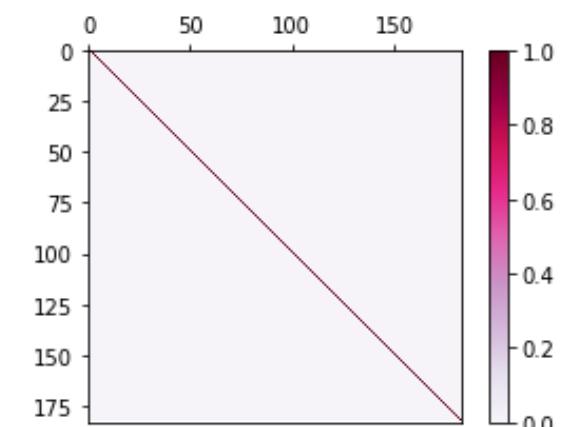
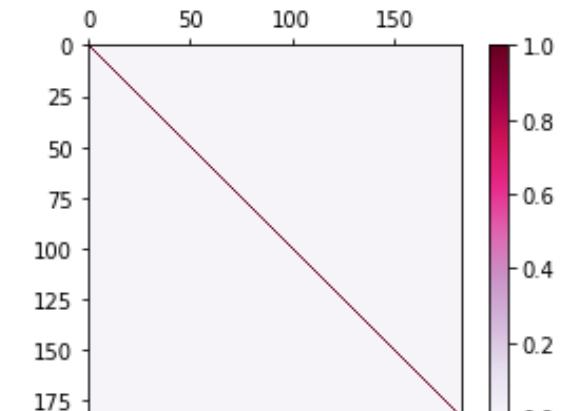
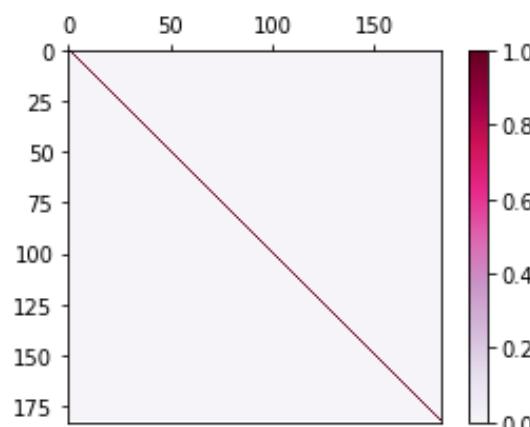
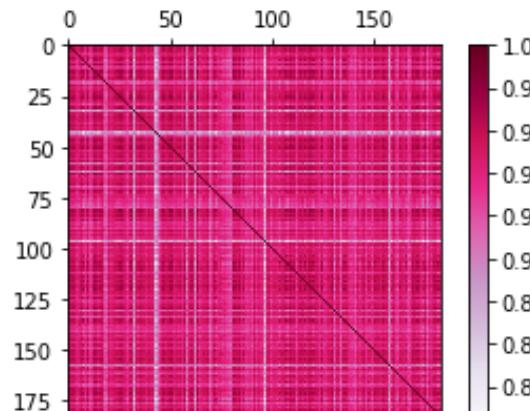
- For a fair comparison with the linear kernel, cross-validate C and r.
- For r, use a logspace between 10000 and 100000 based on your observation of the kernel matrix.

Gaussian kernel SVM

- What values of gamma should we use? Start by spreading out values.

```
for i, gamma in enumerate([1e-5, 1e-2, 1e2, 1e5]):  
    ax = fig.add_subplot(2, 2, i + 1)  
    pairwise_sq_dists = rbf_kernel(X_scaled, gamma=gamma)
```

- When $\gamma > 1e-2$, the kernel matrix is close to the identity.
- When $\gamma = 1e-5$, the kernel matrix is getting close to a matrix of all 1s.
- If we choose γ much smaller, the kernel matrix is going to be so close to a matrix of all 1s the SVM won't learn well.

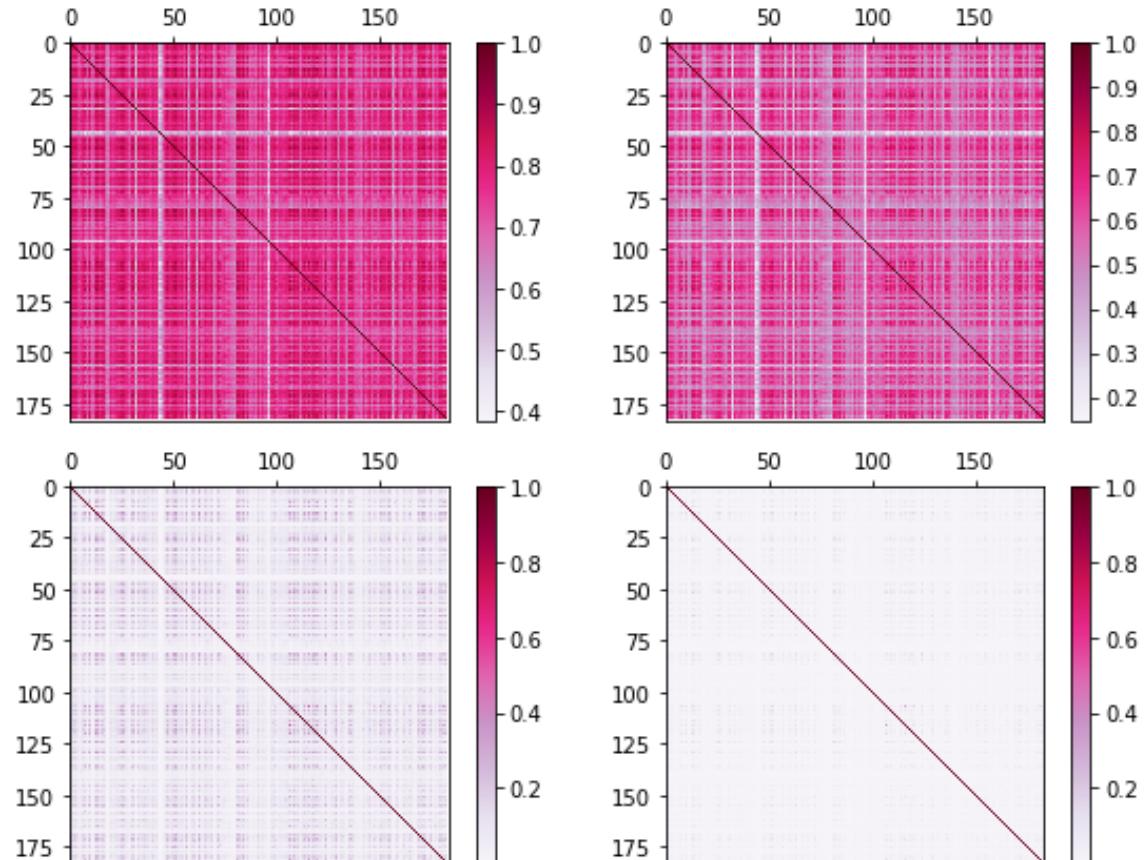


Gaussian kernel SVM

- What values of gamma should we use? Start by spreading out values.

```
for i, gamma in enumerate([5e-5, 1e-4, 5e-4, 1e-3]):  
    ax = fig.add_subplot(2, 2, i + 1)  
    pairwise_sq_dists = rbf_kernel(X_scaled, gamma=gamma)
```

- The kernel matrix is more reasonable when gamma is between 5e-5 and 5e-4.



Gaussian kernel SVM

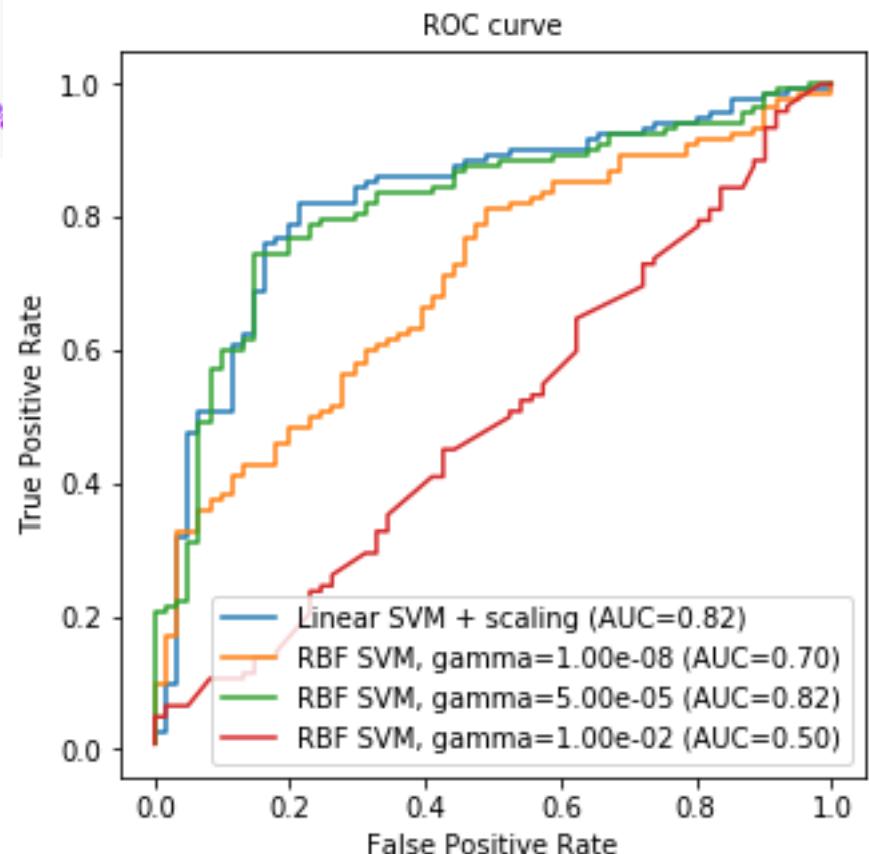
```
fig, ax = plt.subplots(figsize=(5, 5))
ax.plot(fpr2, tpr2,
         label='Linear SVM + scaling (AUC=%2f)' % auc2)

for g in [1e-8, 5e-5, 1e-2]:
    clf = svm.SVC(kernel='rbf', gamma=g)
    ypred_rbf = cross_validate_with_scaling(X, y, clf, folds)

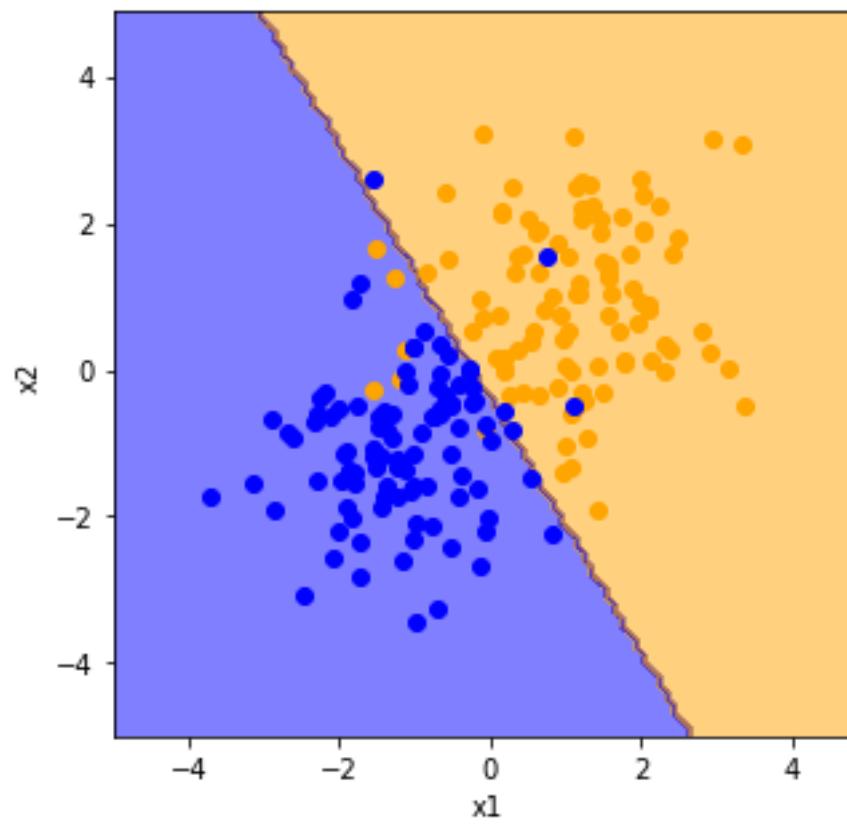
    fpr5, tpr5, thresholds5 = metrics.roc_curve(y, ypred_rbf)
    auc5 = metrics.auc(fpr5, tpr5)

    plt.plot(fpr5, tpr5,
              label='RBF SVM, gamma=%2e (AUC=%2f)' % (g, auc5))
```

- The best performance we obtain is indeed for a gamma of 5e-5.
- To fairly compare to the linear SVM, one should cross-validate C.



Linear SVM decision boundary



Quadratic SVM decision boundary

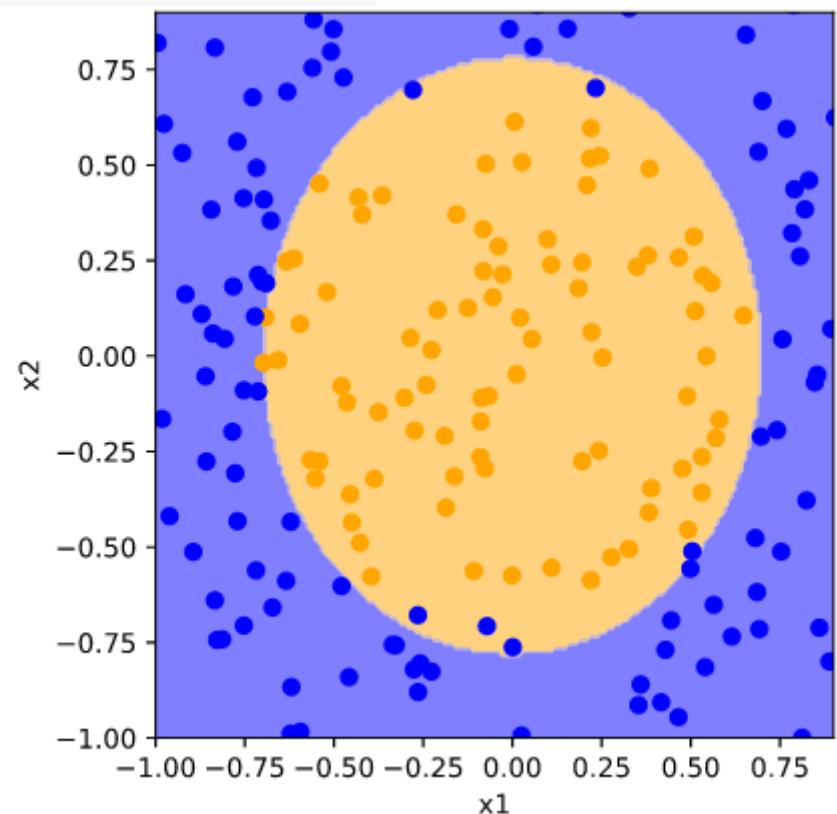
```
# create mesh coordinates
lim = 1
xx, yy = np.meshgrid(np.arange(-lim, +lim, 0.01), np.arange(-lim, +lim, 0.01))
zz_poly = np.zeros(xx.shape)

# create mesh of predictions
for i in range(xx.shape[0]):
    for j in range(xx.shape[1]):
        zz_poly[i, j] = clf_poly.predict(np.array([xx[i, j], yy[i, j]]))

# visualise mesh as contour plot
fig = plt.figure(figsize=(10, 5))

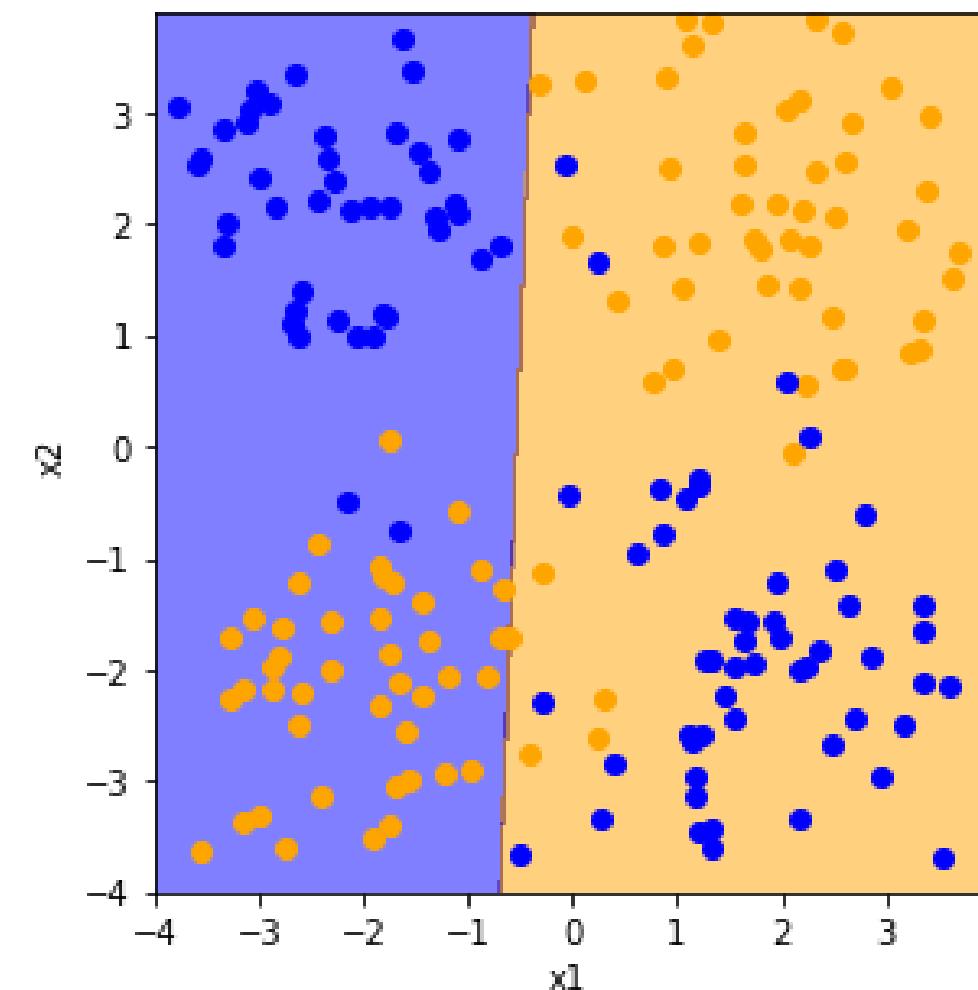
ax = fig.add_subplot(1, 2, 2)
ax.set_xlabel('x1') ; ax.set_ylabel('x2')
ax.set_xlim([-lim, +lim-0.1]) ; ax.set_ylim([-lim, +lim-0.1])
ax.contourf(xx, yy, zz_poly, alpha=0.5, colors=('blue', 'orange'))

# plot training data
ax.scatter(X[class_pos, 0], X[class_pos, 1], color='orange')
ax.scatter(X[class_neg, 0], X[class_neg, 1], color='blue')
```



Separating XOR

linear kernel, accuracy=0.48



RBF kernel, accuracy=0.98

