

Data Graduate Program N°2 Cdiscount Academy

Lundi 20/09 – Mercredi 23/10



SOMMAIRE

- SVM



SVM

Doing really well with linear decision surfaces



SVM – STRENGTH



- Good generalization
 - in theory
 - in practice
- Works well with few training instances
- Find globally best model
- Efficient algorithms
- Amenable to the kernel trick



SVM - NOTATION



To better match notation used in SVMs
...and to make matrix formulas simpler

We will drop using superscripts for the i^{th} instance

i^{th} instance

$$\boldsymbol{x}^{(i)} \rightarrow \mathbf{x}_i$$

Bold denotes vector



i^{th} instance label

$$y^{(i)} \rightarrow y_i$$

Non-bold denotes scalar

j^{th} feature of i^{th} instance

$$x_j^{(i)} \rightarrow x_{ij}$$

SVM – LINEAR SEPARATORS

- Training instances

$$\mathbf{x} \in \mathbb{R}^{d+1}, x_0 = 1$$

$$y \in \{-1, 1\}$$

- Model parameters

$$\boldsymbol{\theta} \in \mathbb{R}^{d+1}$$

- Hyperplane

$$\boldsymbol{\theta}^\top \mathbf{x} = \langle \boldsymbol{\theta}, \mathbf{x} \rangle = 0$$

- Decision function

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^\top \mathbf{x}) = \text{sign}(\langle \boldsymbol{\theta}, \mathbf{x} \rangle)$$



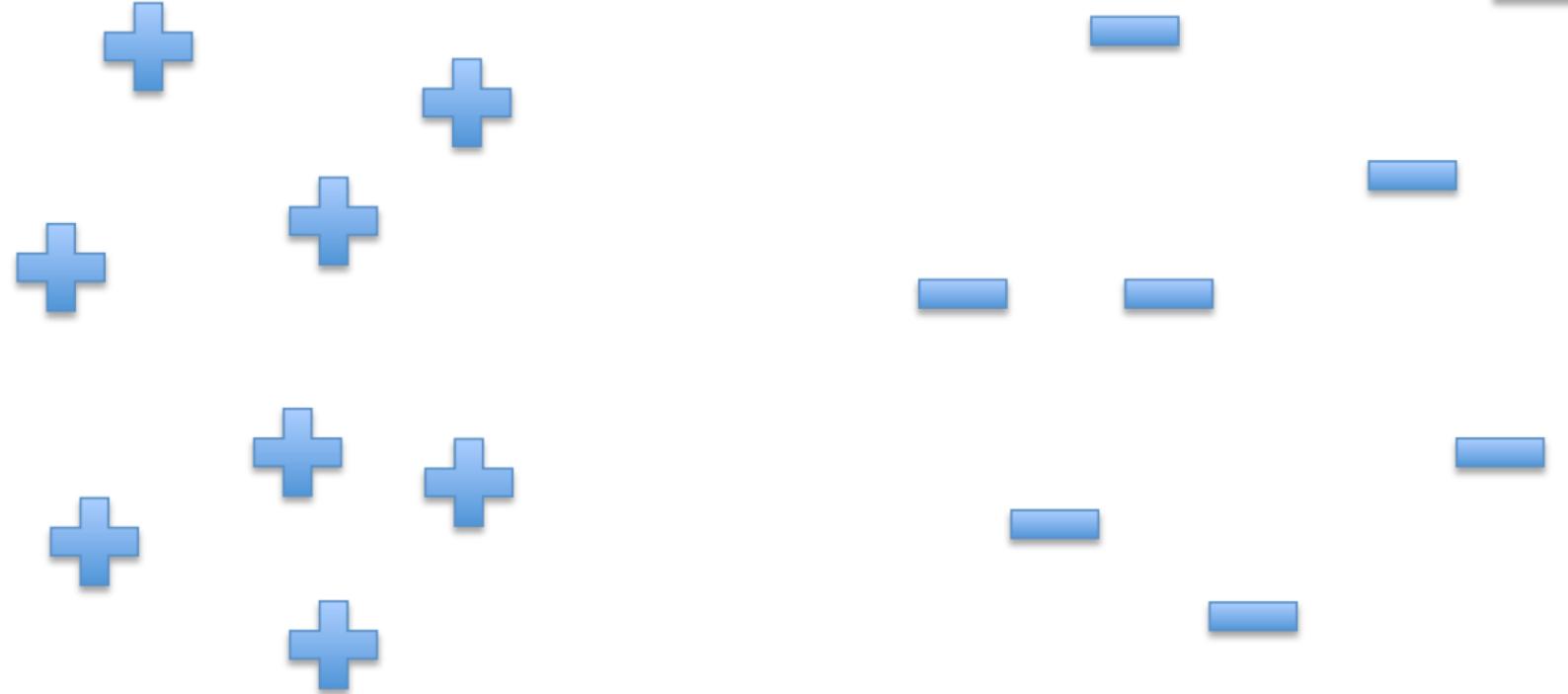
Recall:

Inner (dot) product:

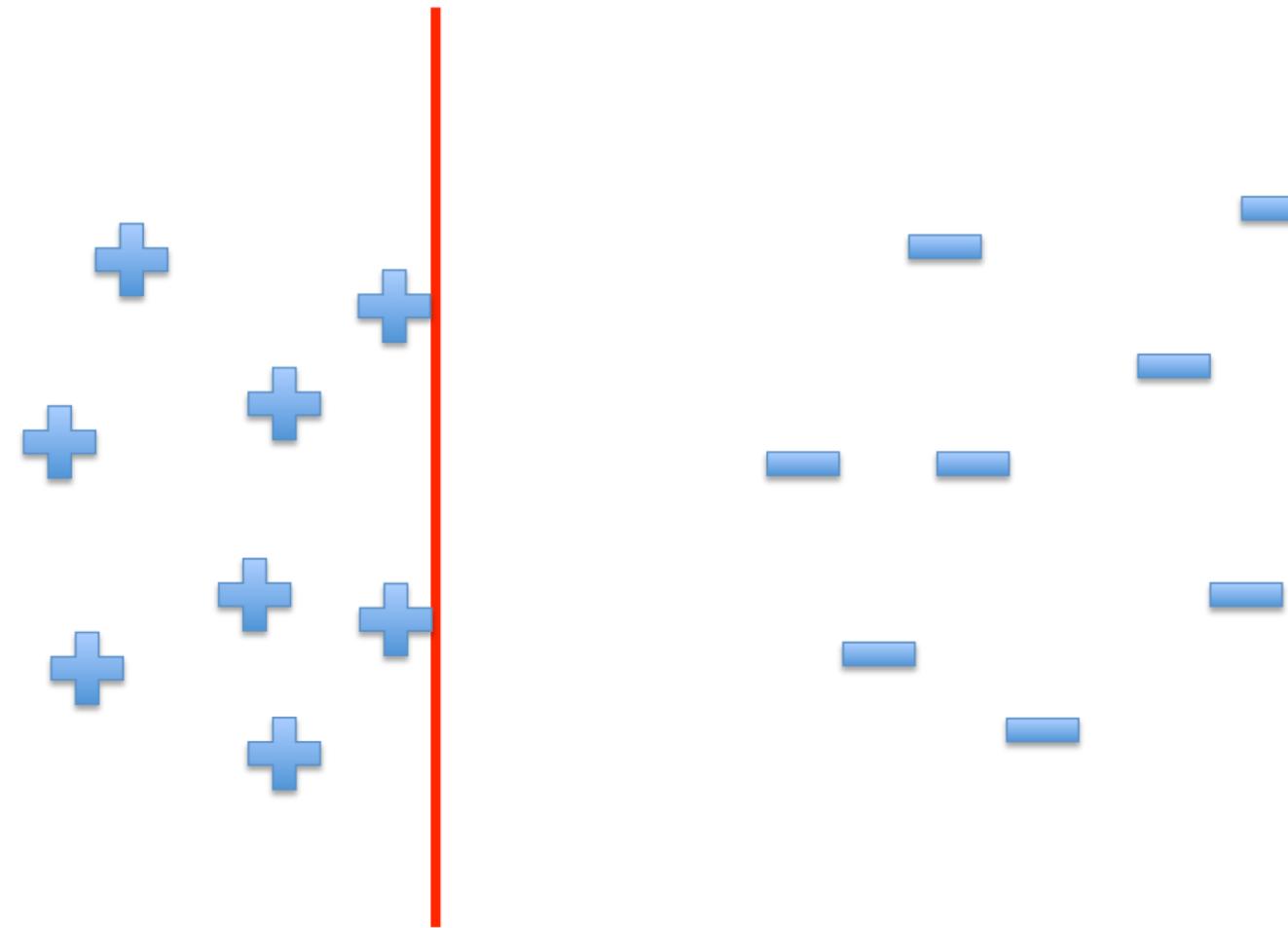
$$\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u} \cdot \mathbf{v} = \mathbf{u}^\top \mathbf{v}$$

$$= \sum_i u_i v_i$$

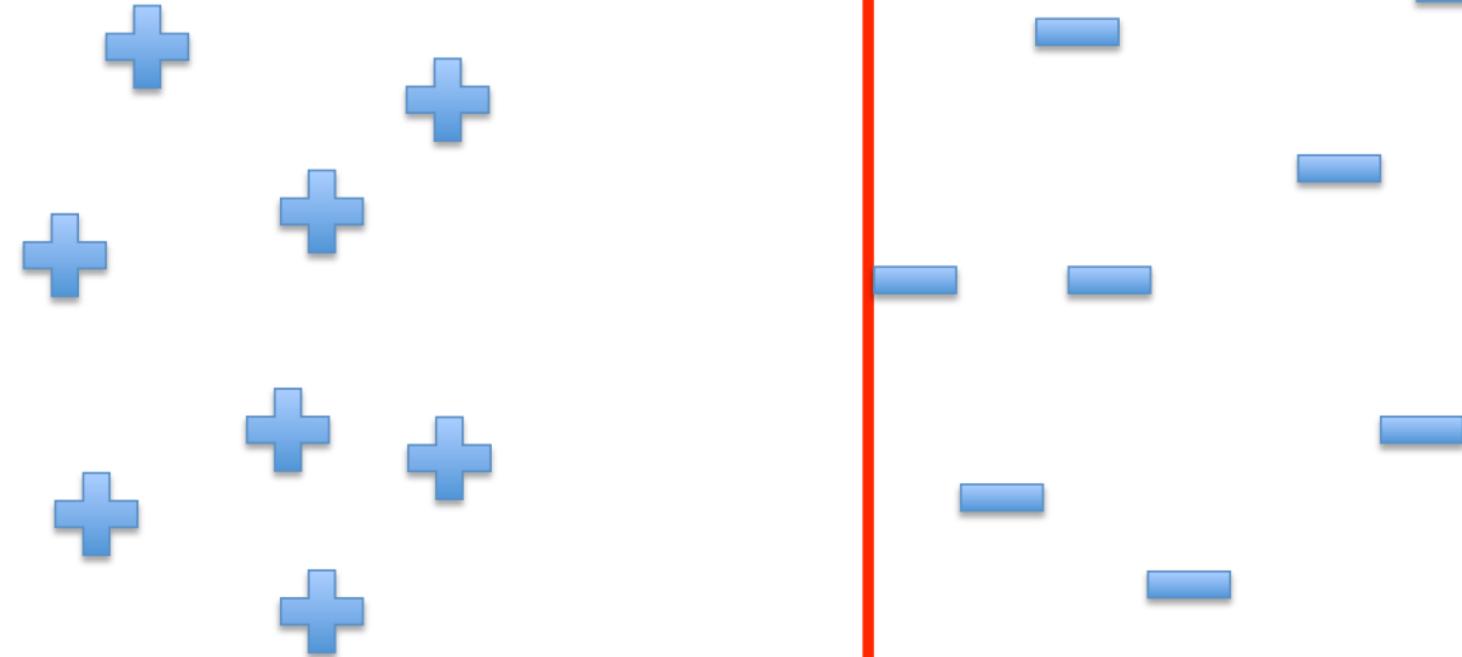
SVM – INTUITIONS



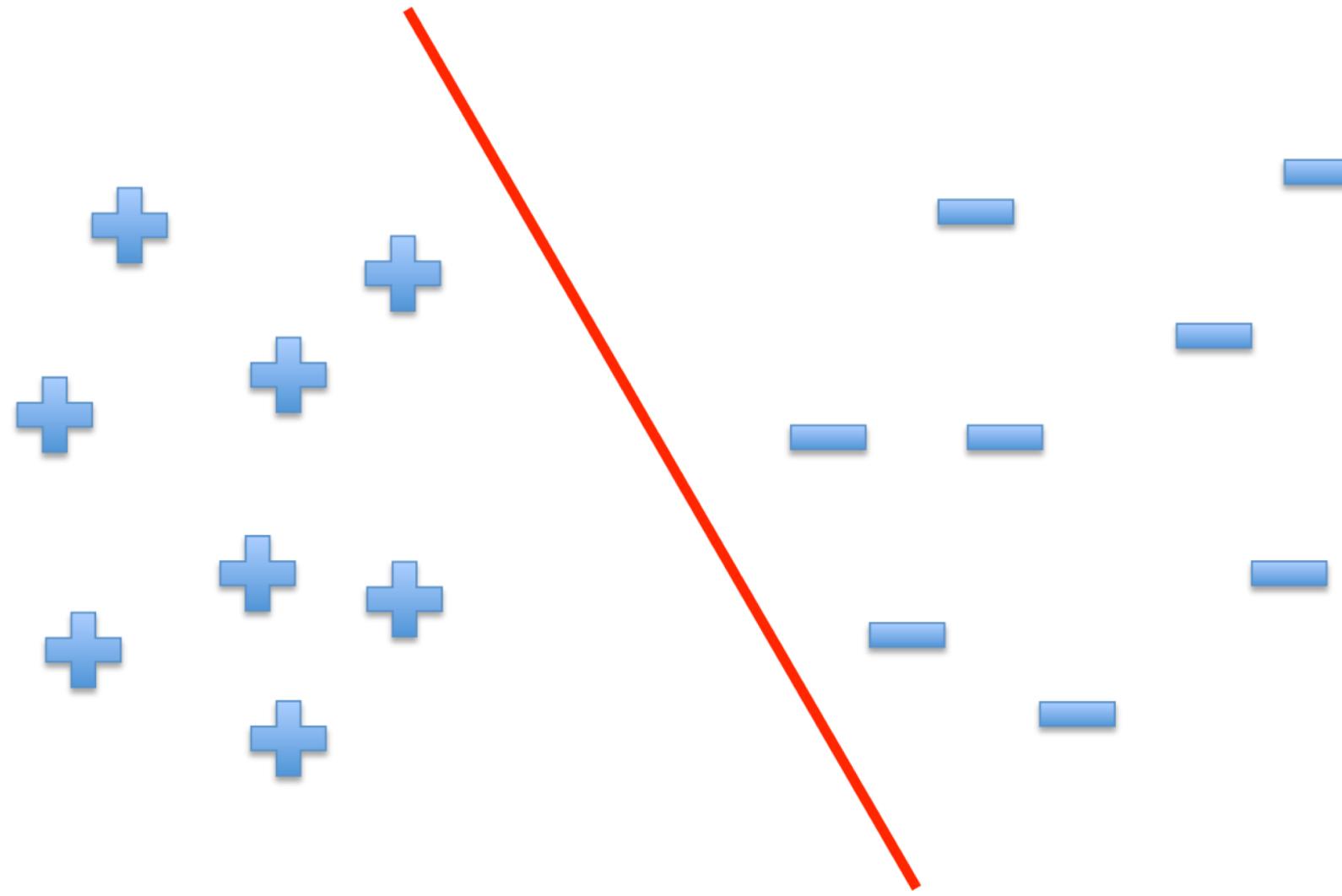
SVM – INTUITIONS



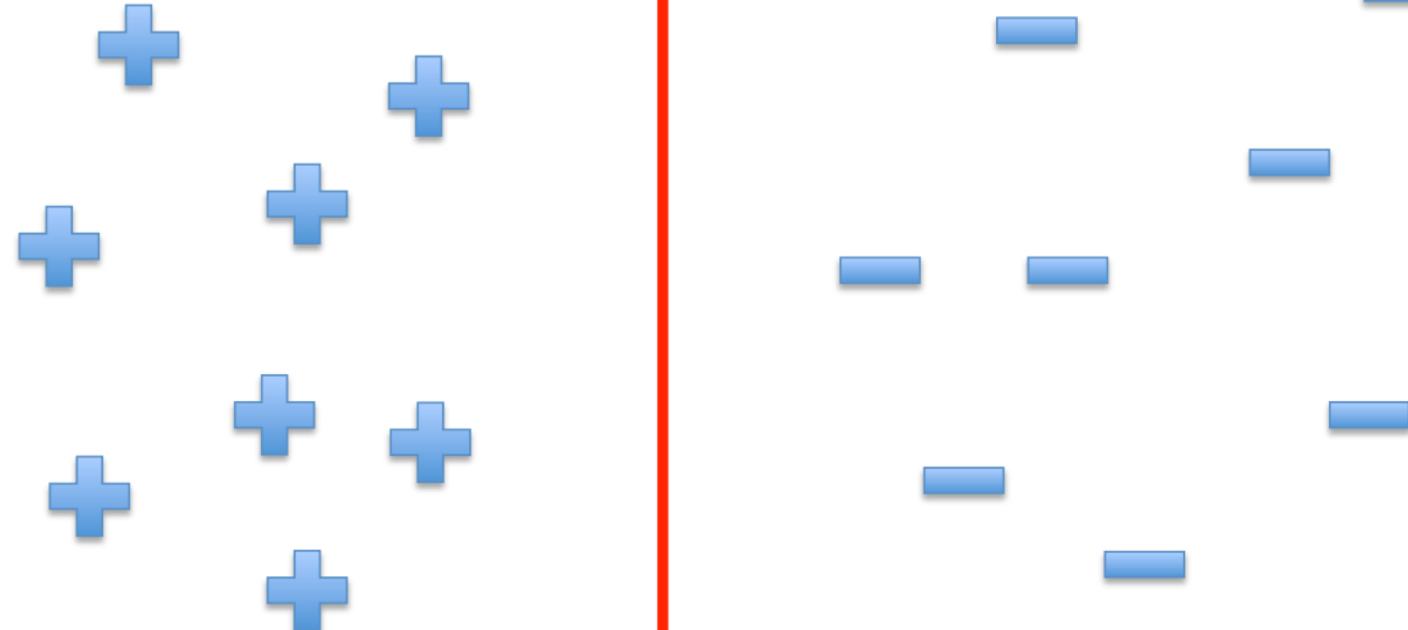
SVM – INTUITIONS



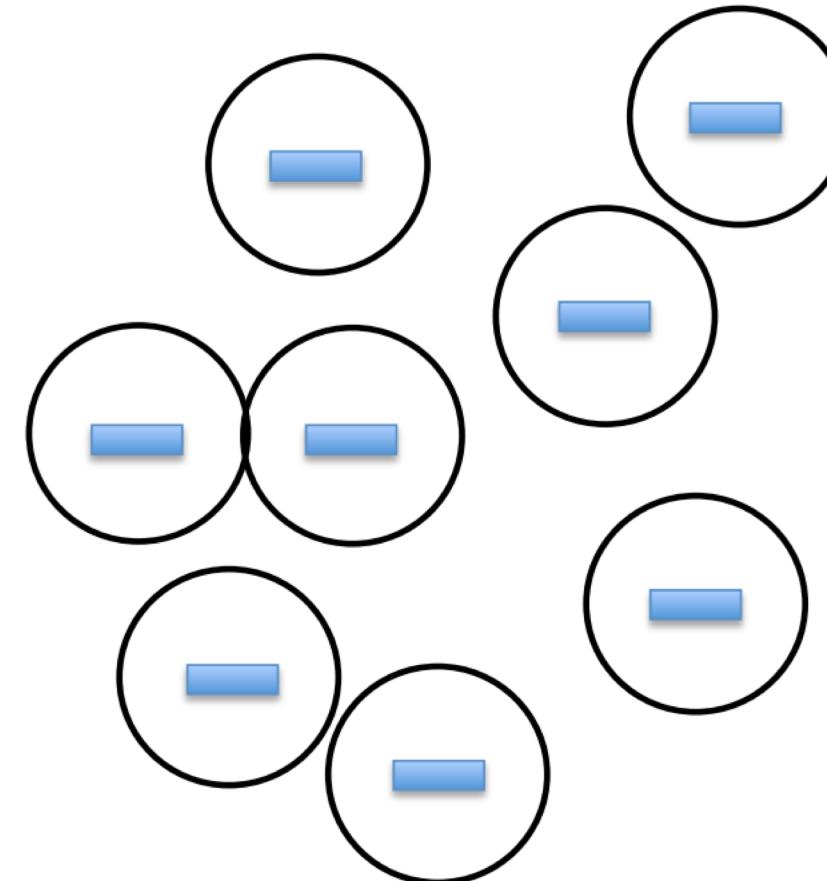
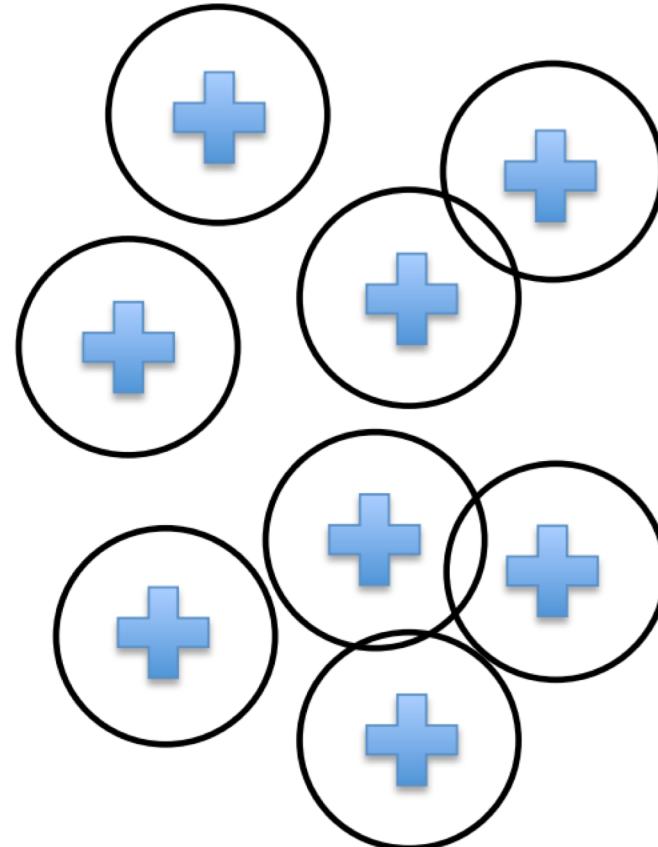
SVM – INTUITIONS



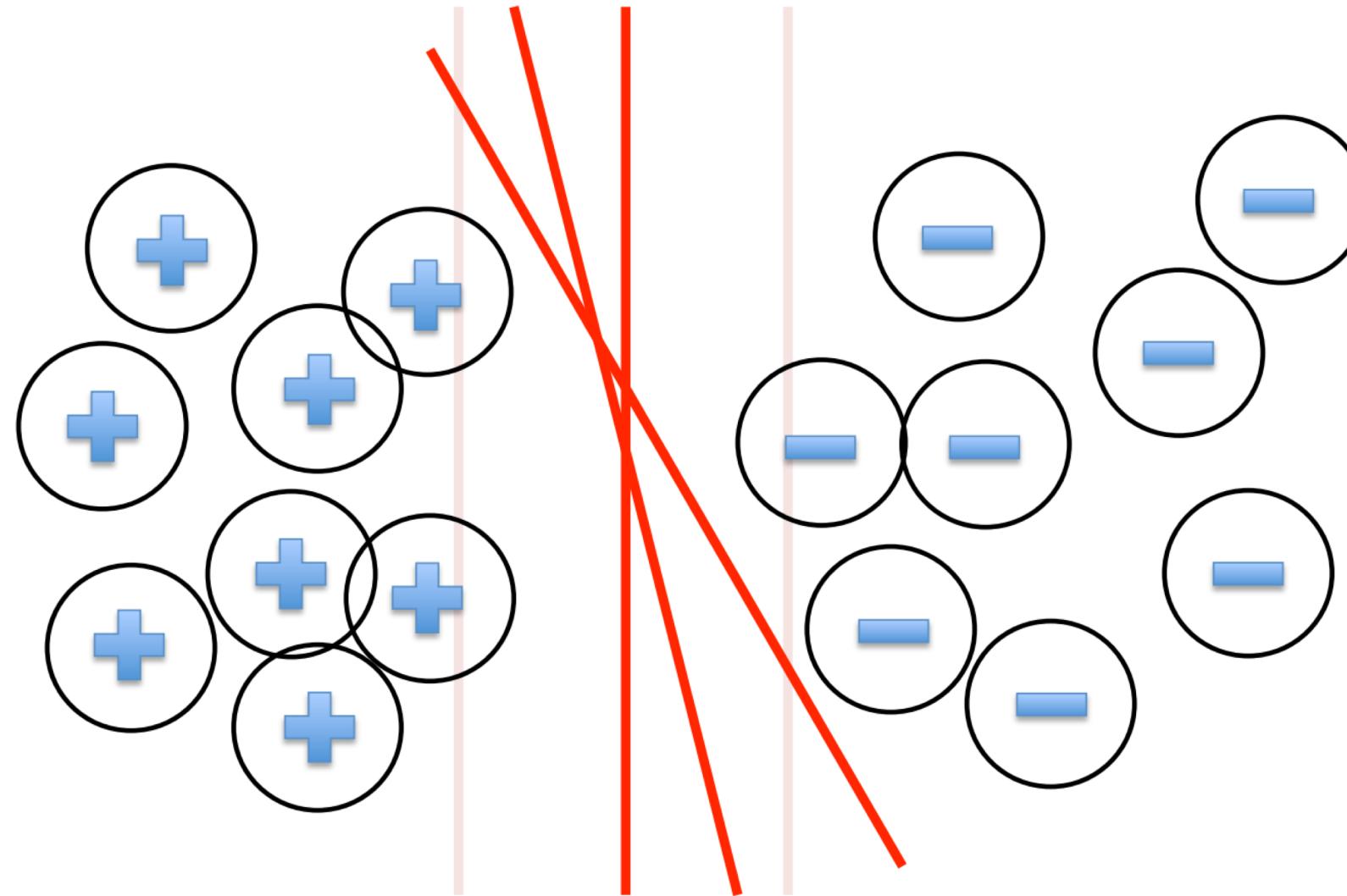
SVM – INTUITIONS : A « GOOD » SEPARATOR



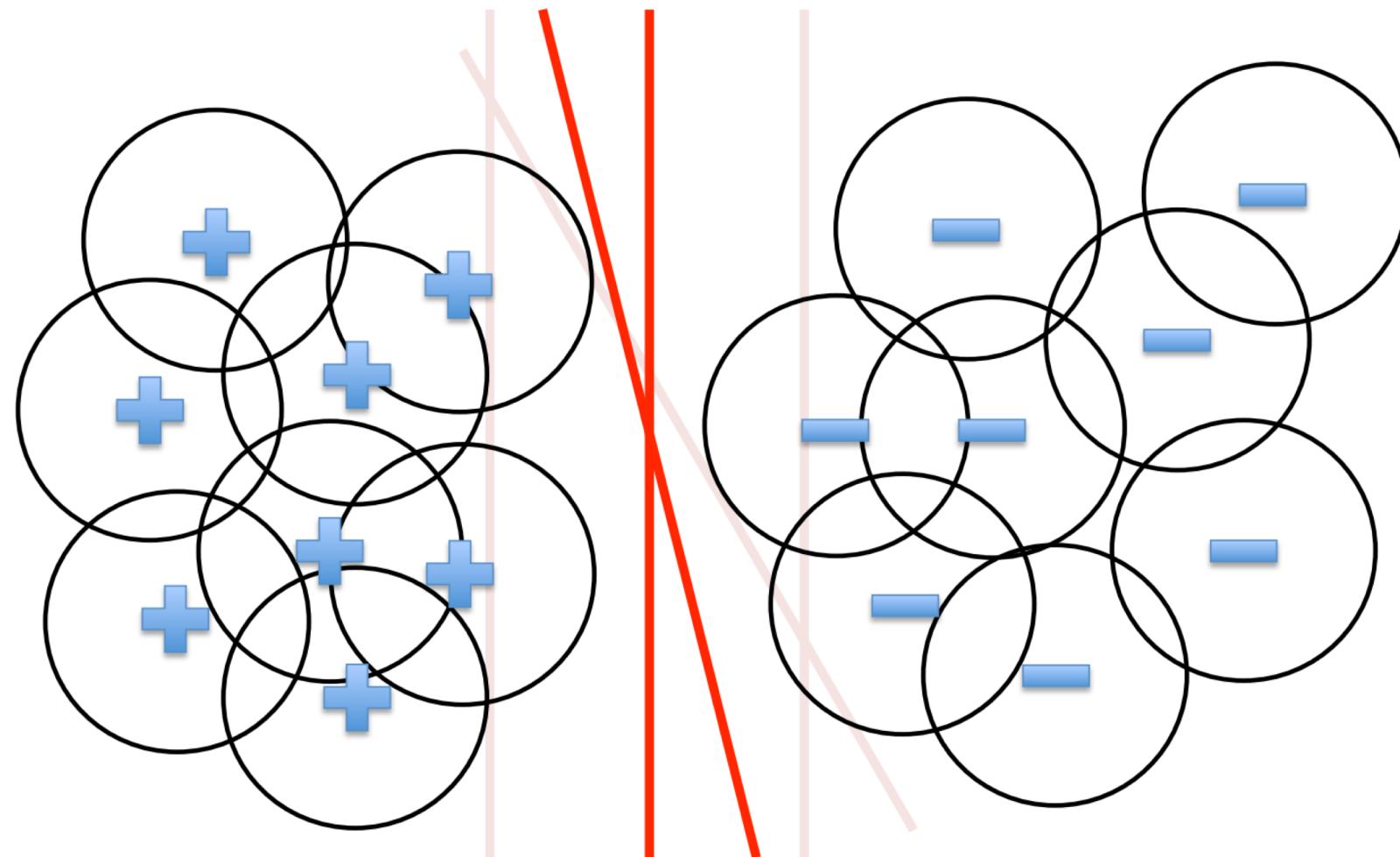
SVM – NOISE IN THE OBSERVATIONS



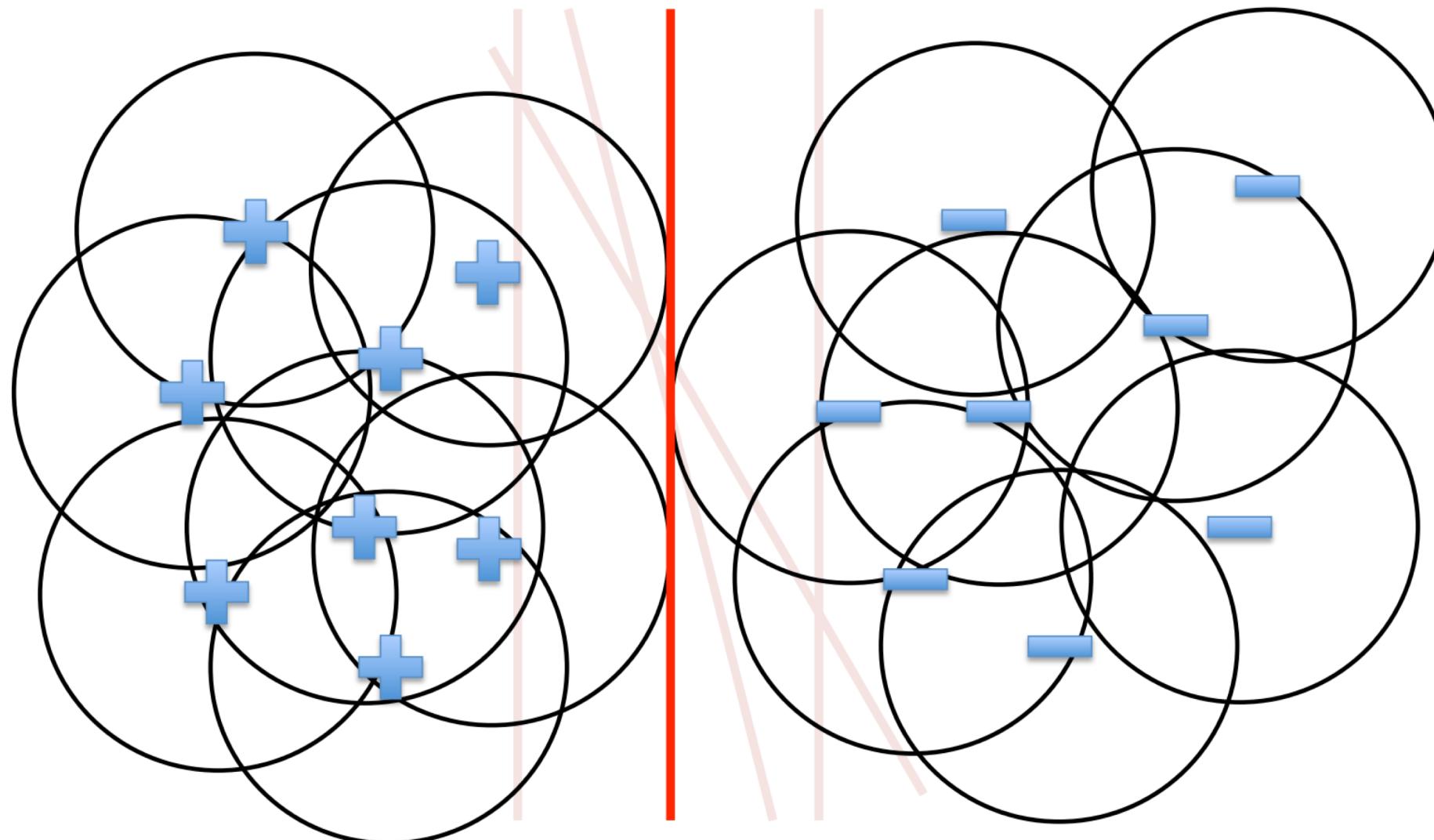
SVM – RULING OUT SOME SEPARATORS



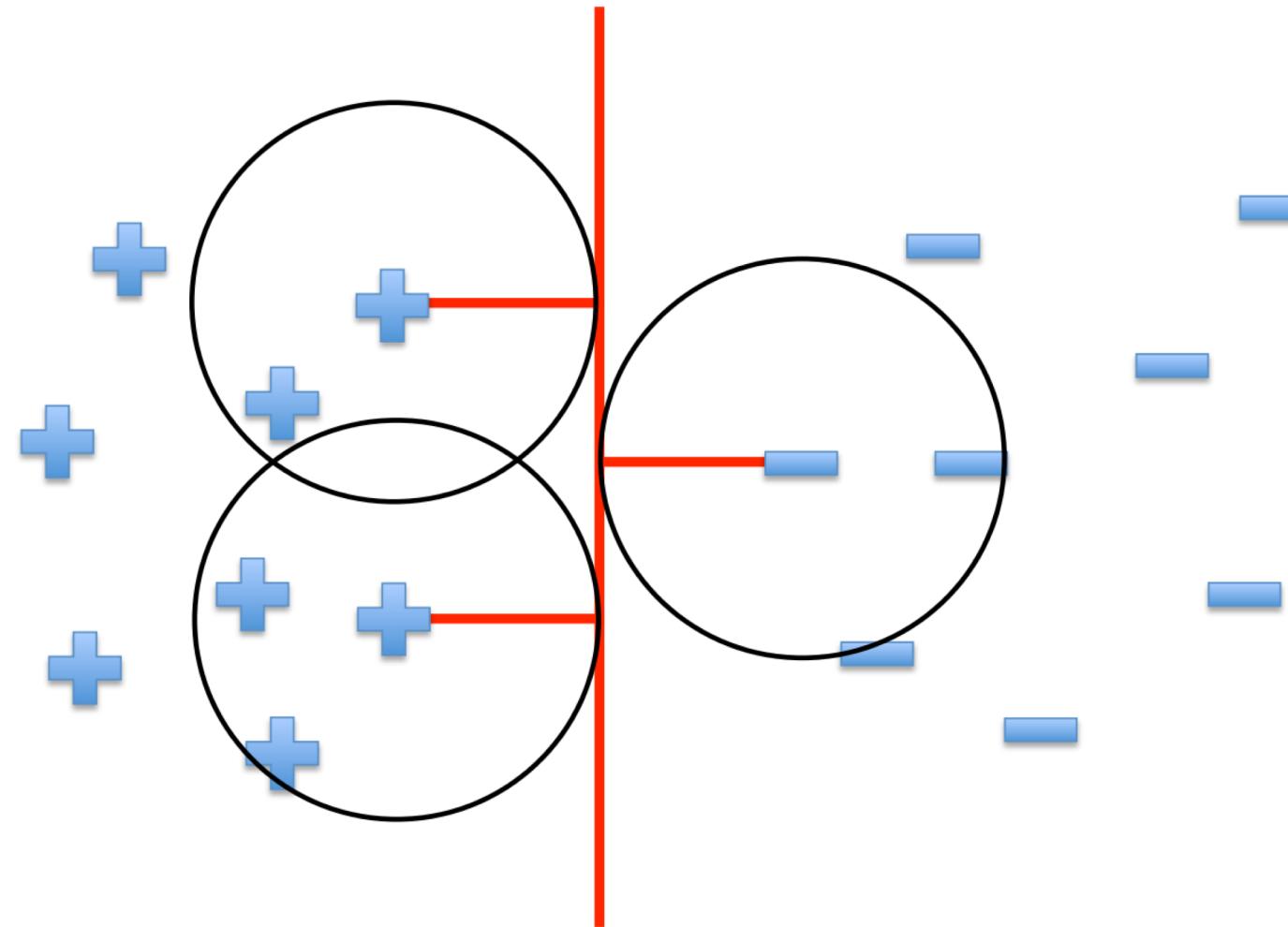
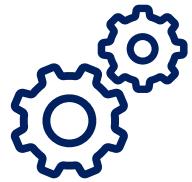
SVM – LOTS OF NOISE



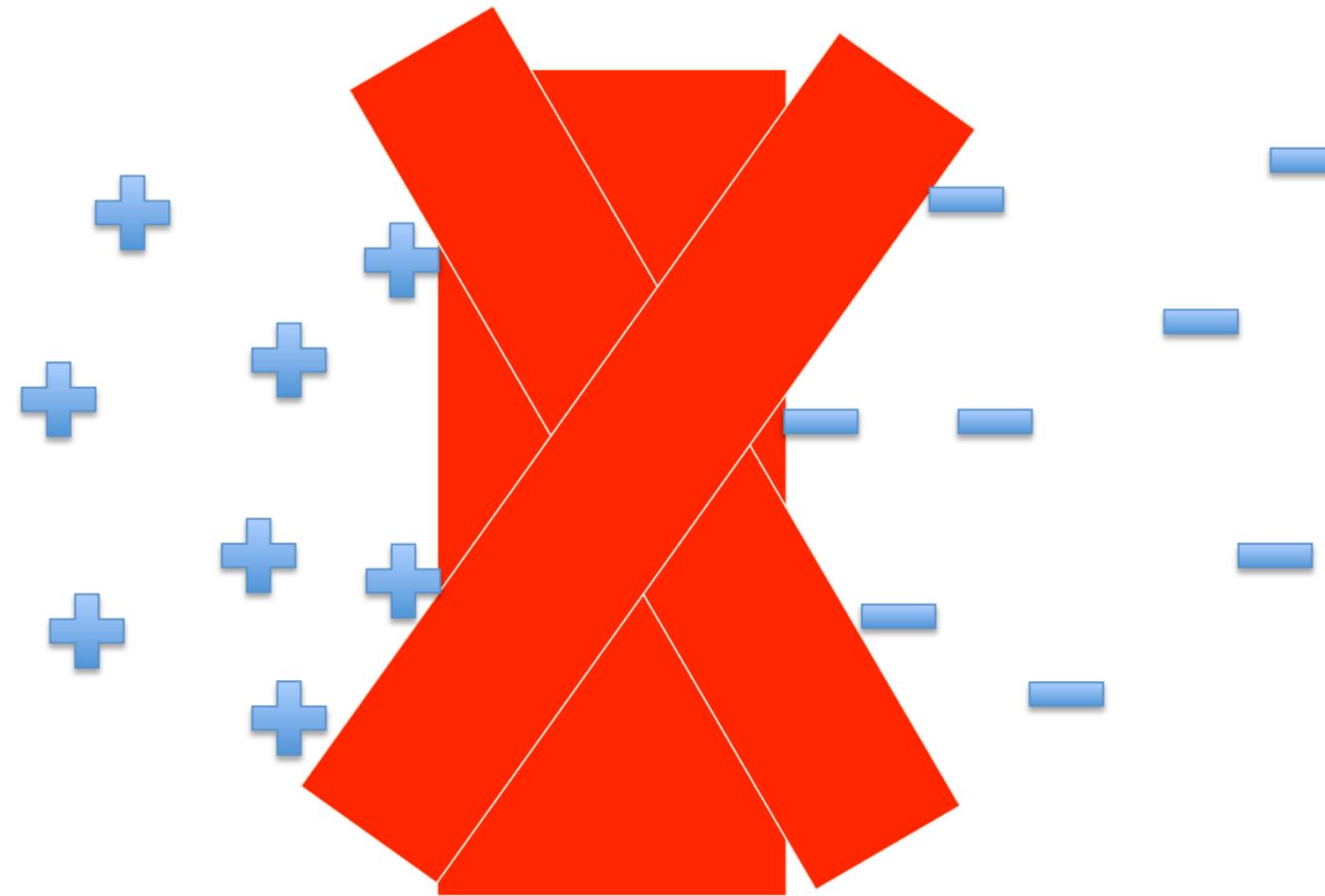
SVM – ONLY ONE SEPARATOR REMAINS



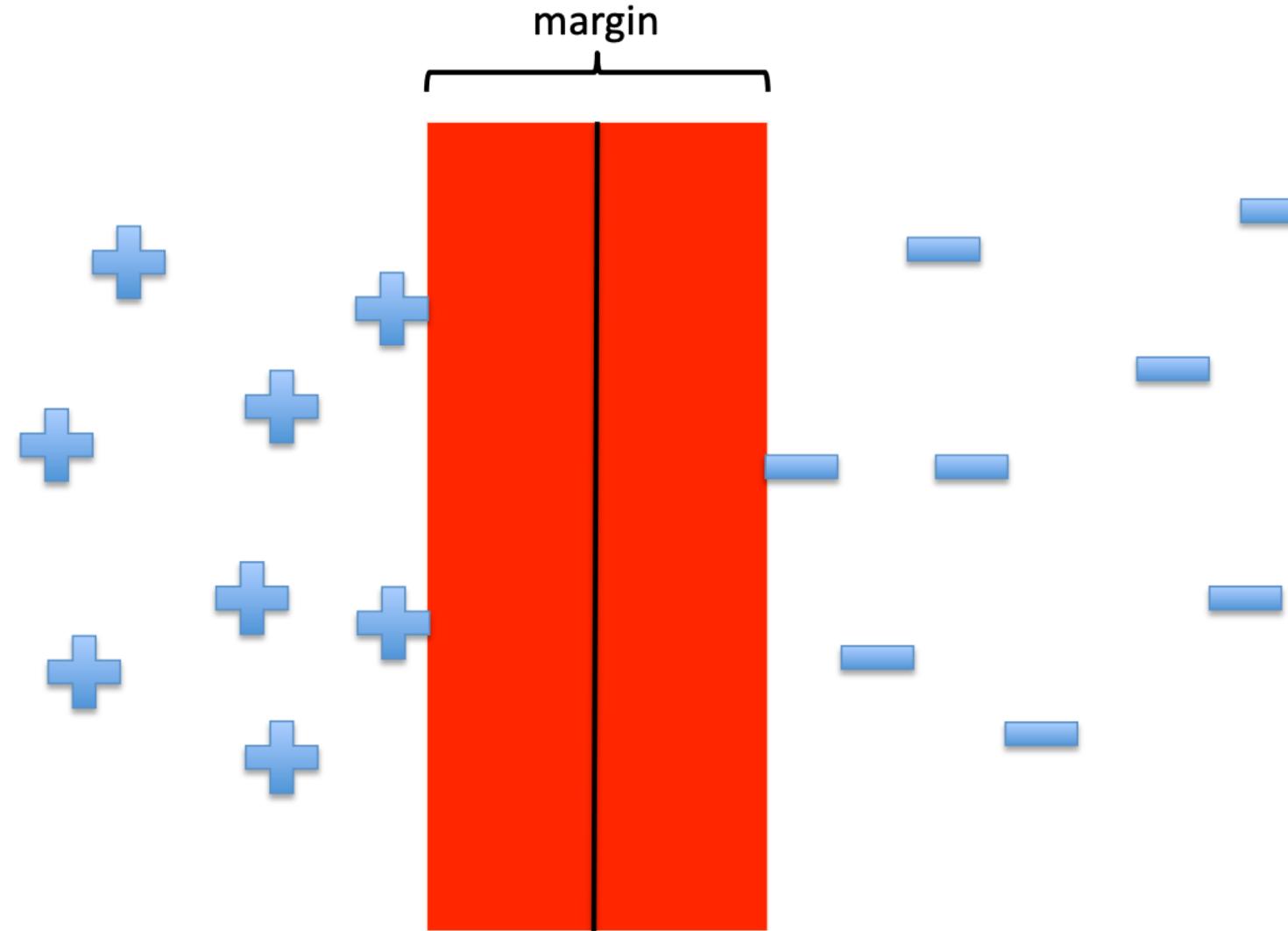
SVM – MAXIMIZING THE MARGIN



SVM – « FAT SEPARATORS »



SVM – « FAT SEPARATORS »



SVM – WHY MAXIMIZE MARGIN



Increasing margin reduces *capacity*

- i.e., fewer possible models



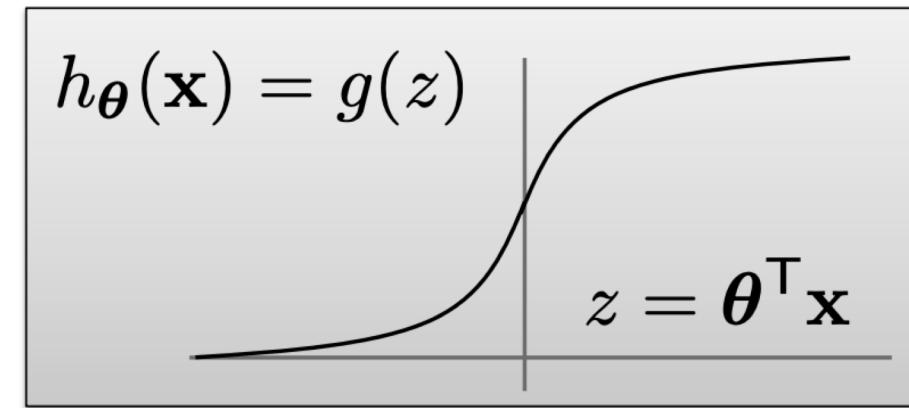
Lesson from Learning Theory:

- If the following holds:
 - H is sufficiently constrained in size
 - and/or the size of the training data set n is large, then low training error is likely to be evidence of low generalization error

ALTERNATIVE VIEW OF LOGISTIC REGRESSION



$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$



If $y = 1$, we want $h_{\theta}(\mathbf{x}) \approx 1$, $\theta^T \mathbf{x} \gg 0$

If $y = 0$, we want $h_{\theta}(\mathbf{x}) \approx 0$, $\theta^T \mathbf{x} \ll 0$



$$J(\theta) = - \sum_{i=1}^n [y_i \log h_{\theta}(\mathbf{x}_i) + (1 - y_i) \log (1 - h_{\theta}(\mathbf{x}_i))]$$

$\min_{\theta} J(\theta)$ $\underbrace{\text{cost}_1(\theta^T \mathbf{x}_i)}$ $\underbrace{\text{cost}_0(\theta^T \mathbf{x}_i)}$

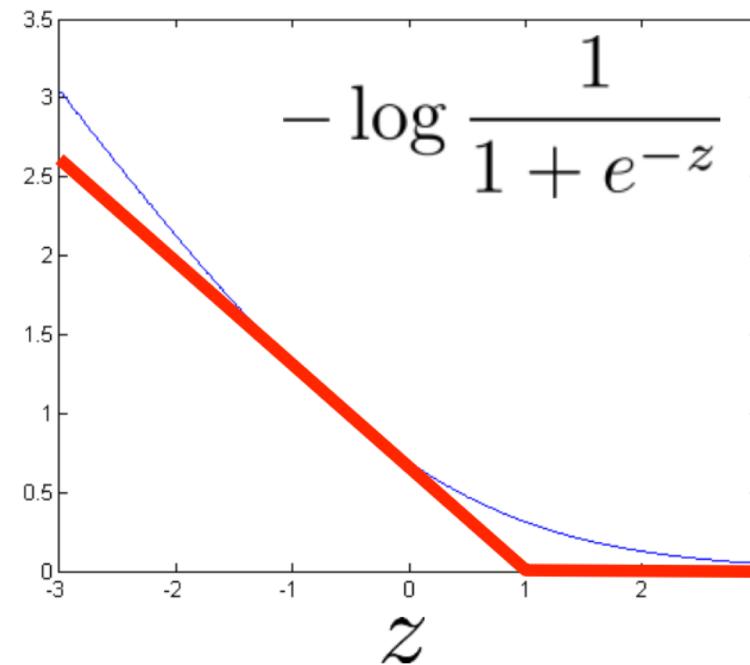
ALTERNATIVE VIEW OF LOGISTIC REGRESSION

Cost of example: $-y_i \log h_{\theta}(\mathbf{x}_i) - (1 - y_i) \log (1 - h_{\theta}(\mathbf{x}_i))$

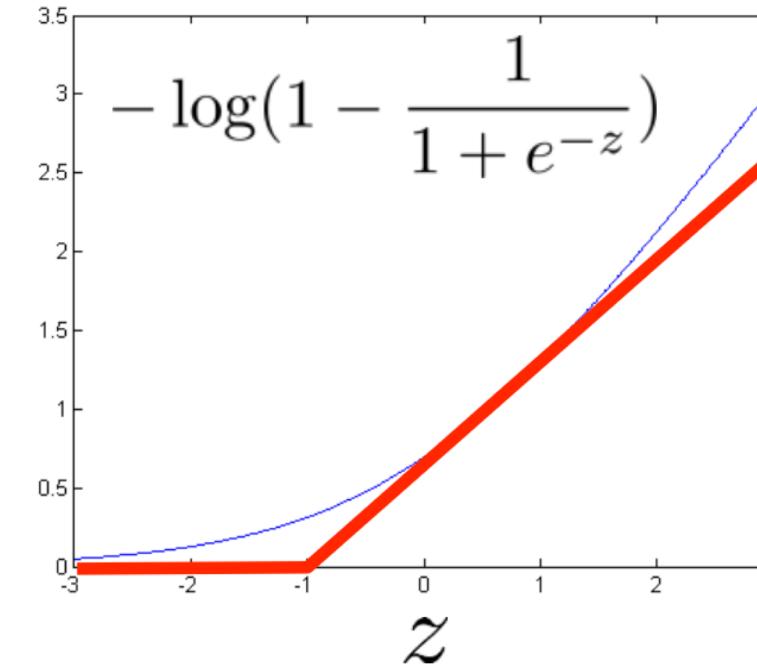


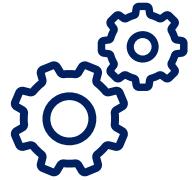
$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \quad z = \theta^T \mathbf{x}$$

If $y = 1$ (want $\theta^T \mathbf{x} \gg 0$):



If $y = 0$ (want $\theta^T \mathbf{x} \ll 0$):





Logistic Regression:

$$\min_{\theta} -\sum_{i=1}^n [y_i \log h_{\theta}(\mathbf{x}_i) + (1 - y_i) \log (1 - h_{\theta}(\mathbf{x}_i))] + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$



Support Vector Machines:

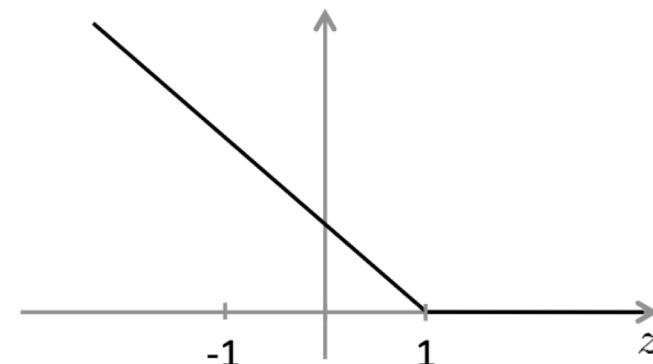
$$\min_{\theta} C \sum_{i=1}^n [y_i \text{cost}_1(\boldsymbol{\theta}^T \mathbf{x}_i) + (1 - y_i) \text{cost}_0(\boldsymbol{\theta}^T \mathbf{x}_i)] + \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

You can think of C as similar to $\frac{1}{\lambda}$

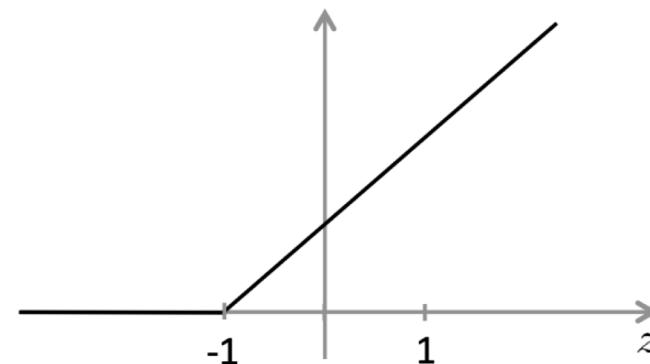


$$\min_{\theta} C \sum_{i=1}^n [y_i \text{cost}_1(\theta^\top \mathbf{x}_i) + (1 - y_i) \text{cost}_0(\theta^\top \mathbf{x}_i)] + \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

If $y = 1$ (want $\theta^\top \mathbf{x} \geq 1$):



If $y = 0$ (want $\theta^\top \mathbf{x} \leq -1$):



$$\ell_{\text{hinge}}(h(\mathbf{x})) = \max(0, 1 - y \cdot h(\mathbf{x}))$$



$$\min_{\theta} C \sum_{i=1}^n [y_i \text{cost}_1(\theta^\top \mathbf{x}_i) + (1 - y_i) \text{cost}_0(\theta^\top \mathbf{x}_i)] + \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

$y = 1 / 0$

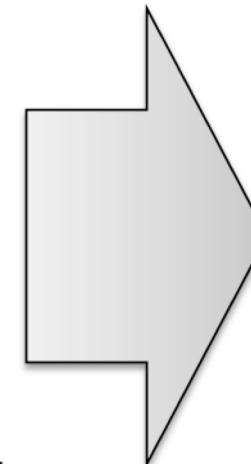
with $C = 1$

$y = +1 / -1$



$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

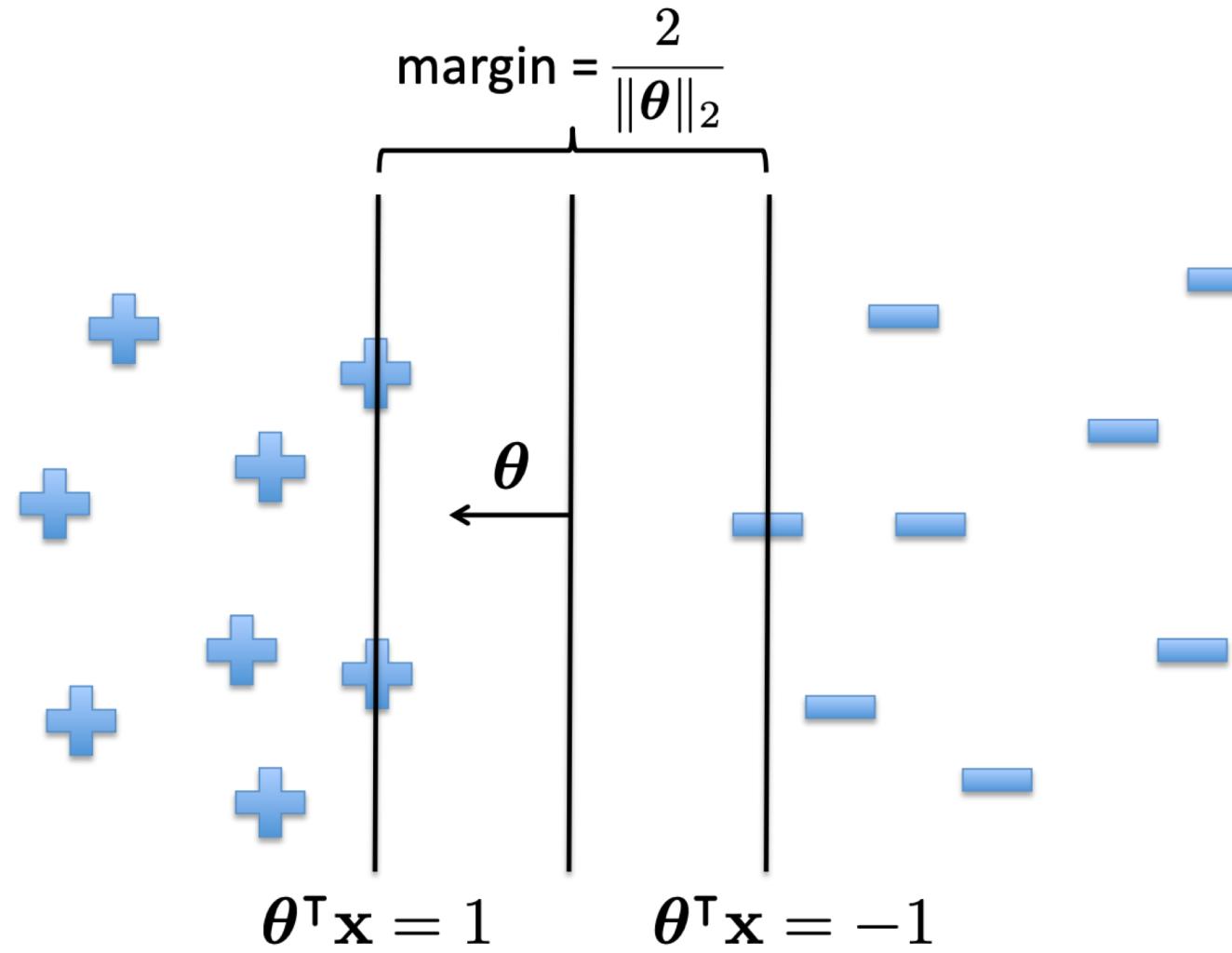
s.t. $\theta^\top \mathbf{x}_i \geq 1 \quad \text{if } y_i = 1$
 $\theta^\top \mathbf{x}_i \leq -1 \quad \text{if } y_i = -1$



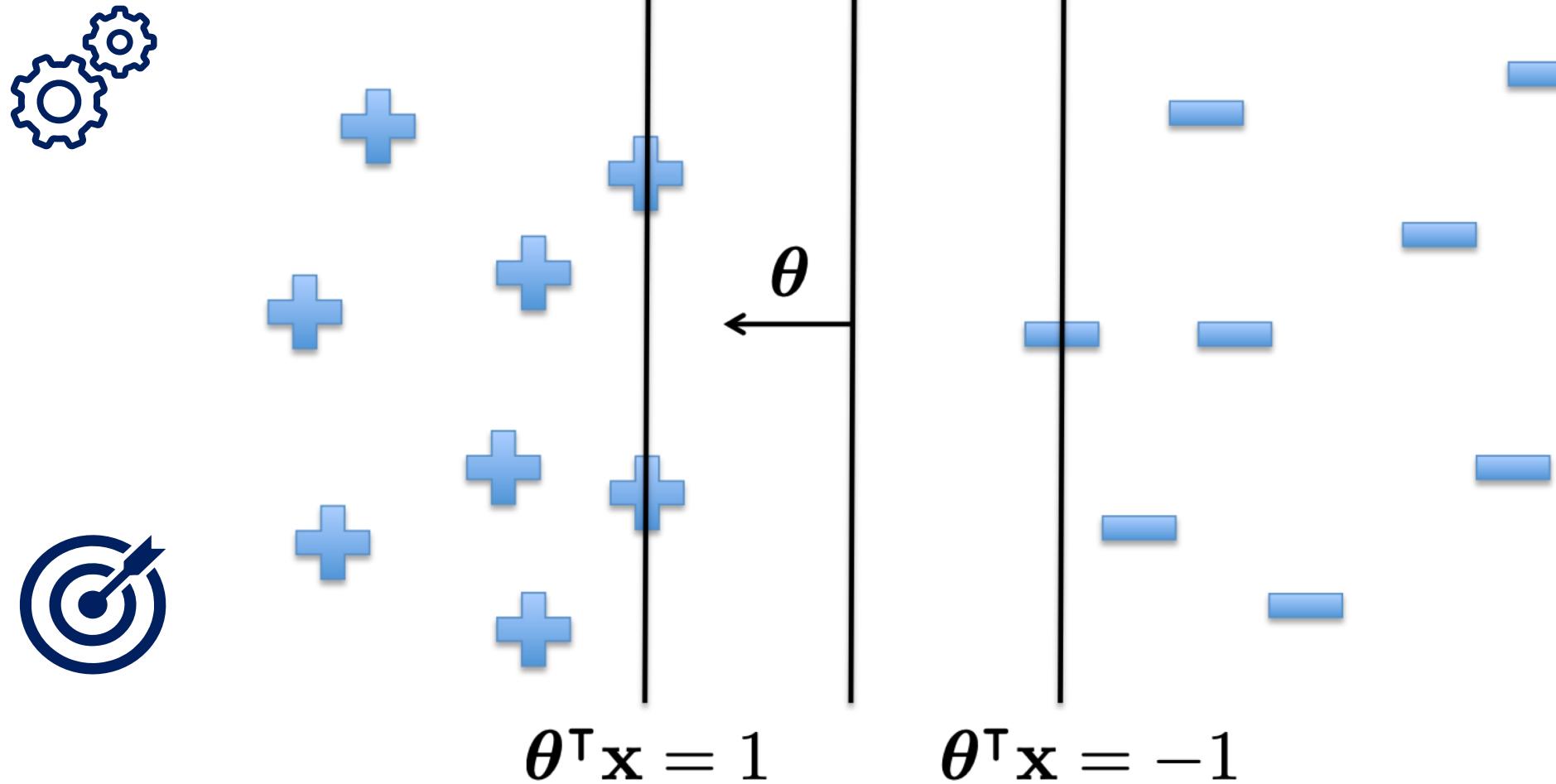
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

s.t. $y_i(\theta^\top \mathbf{x}_i) \geq 1$

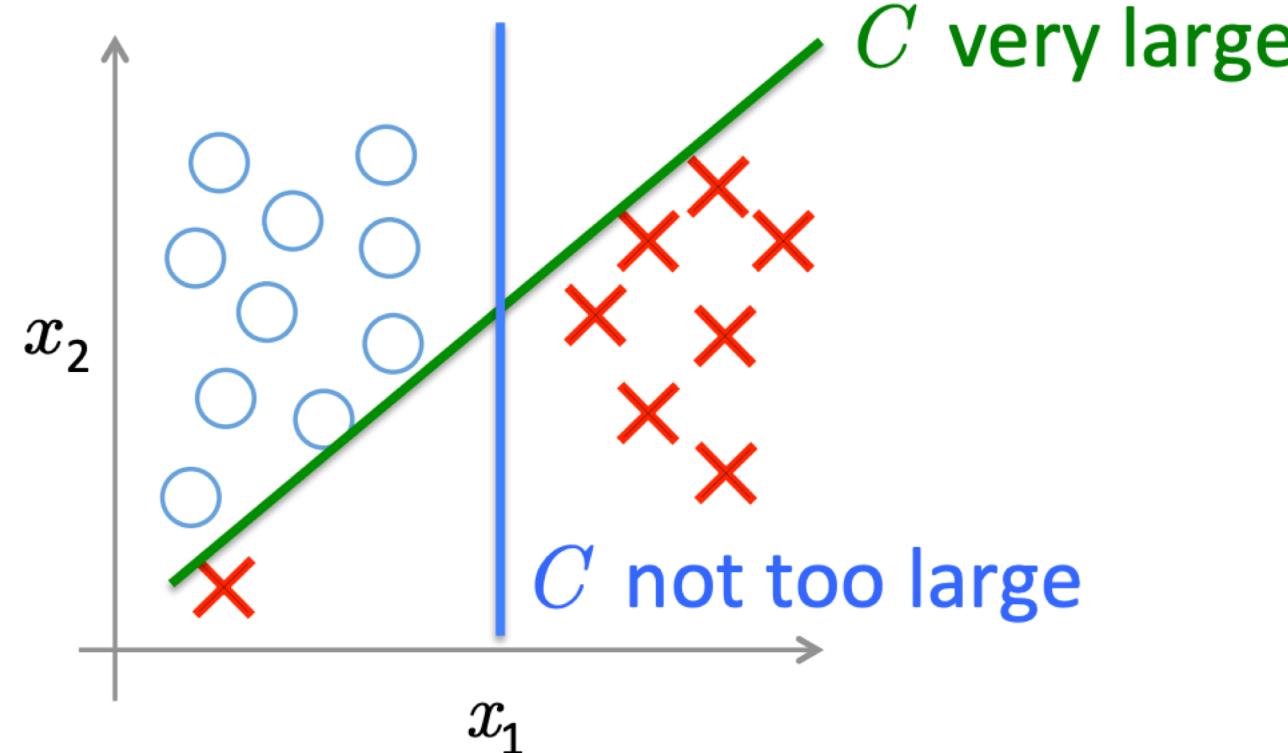
SVM : MAXIMUM MARGIN HYPERPLANE



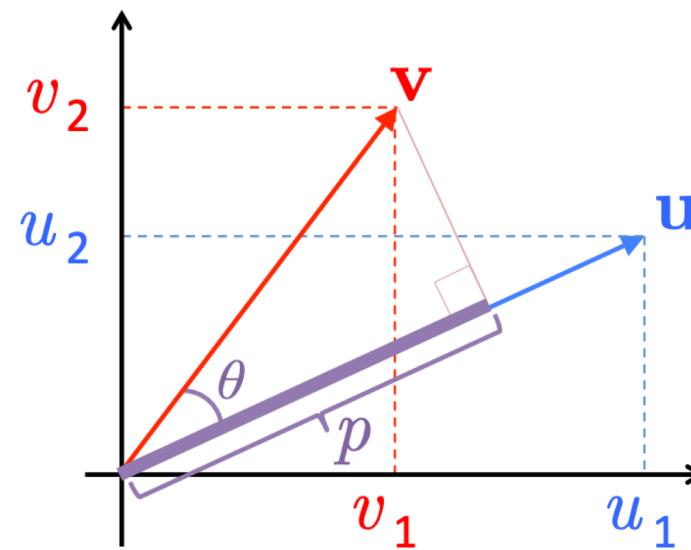
SVM : SUPPORT VECTORS



SVM : LARGE MARGIN CLASSIFIER IN PRESENCE OF OUTLIERS



SVM : VECTOR INNER PRODUCT



$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{aligned}\|\mathbf{u}\|_2 &= \text{length}(\mathbf{u}) \in \mathbb{R} \\ &= \sqrt{u_1^2 + u_2^2}\end{aligned}$$



$$\begin{aligned}\mathbf{u}^\top \mathbf{v} &= \mathbf{v}^\top \mathbf{u} \\ &= u_1 v_1 + u_2 v_2 \\ &= \|\mathbf{u}\|_2 \|\mathbf{v}\|_2 \cos \theta \\ &= p \|\mathbf{u}\|_2 \quad \text{where } p = \|\mathbf{v}\|_2 \cos \theta\end{aligned}$$

SVM : UNDERSTANDING THE HYPERPLANE

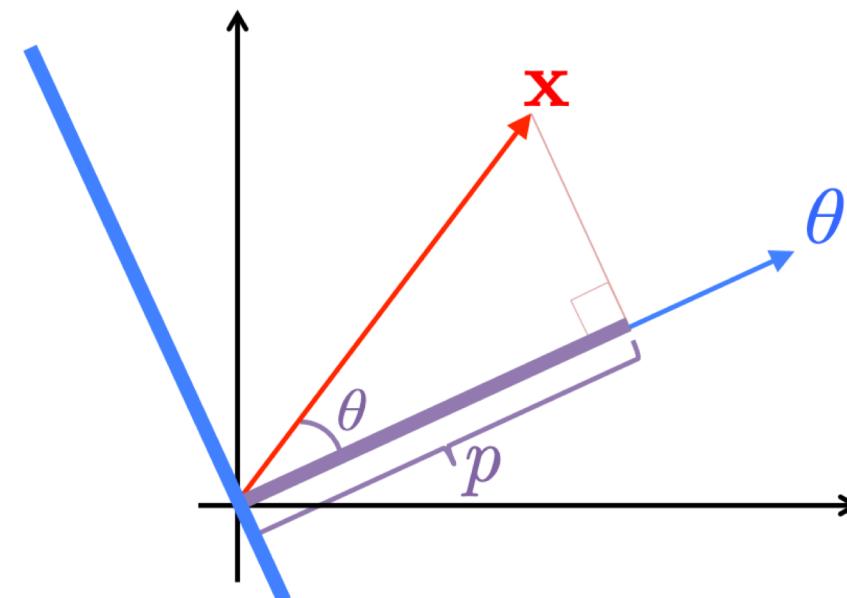


$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

s.t. $\theta^\top \mathbf{x}_i \geq 1 \quad \text{if } y_i = 1$

$\theta^\top \mathbf{x}_i \leq -1 \quad \text{if } y_i = -1$

Assume $\theta_0 = 0$ so that the hyperplane is centered at the origin, and that $d = 2$



$$\begin{aligned}\theta^\top \mathbf{x} &= \|\theta\|_2 \underbrace{\|\mathbf{x}\|_2 \cos \theta}_p \\ &= p\|\theta\|_2\end{aligned}$$

SVM : MAXIMIZING THE MARGIN



$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

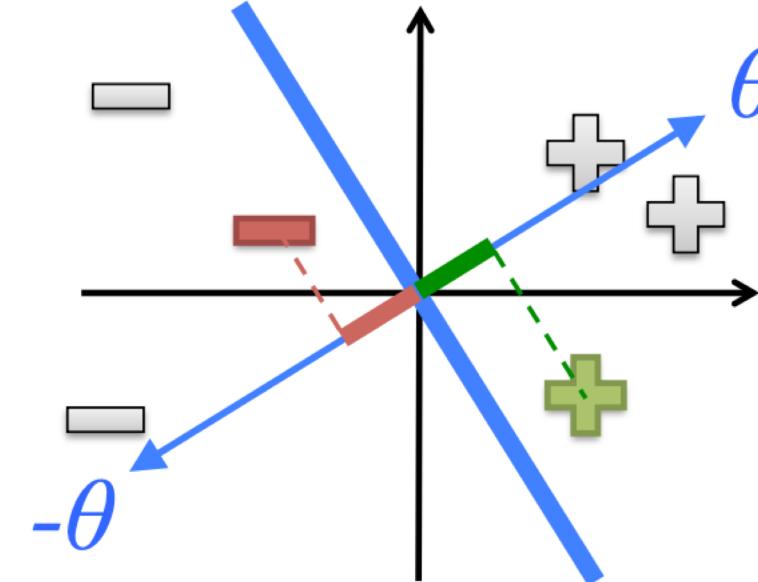
$$\begin{aligned} \text{s.t. } \theta^\top \mathbf{x}_i &\geq 1 & \text{if } y_i = 1 \\ \theta^\top \mathbf{x}_i &\leq -1 & \text{if } y_i = -1 \end{aligned}$$



Assume $\theta_0 = 0$ so that the hyperplane is centered at the origin, and that $d = 2$

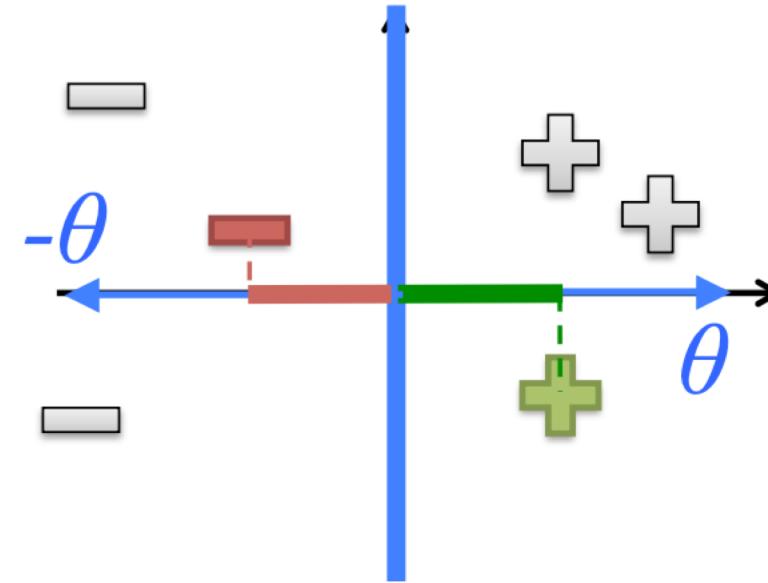
Let p_i be the projection of \mathbf{x}_i onto the vector θ

SVM : MAXIMIZING THE MARGIN



Since p is small, therefore $\|\theta\|_2$ must be large to have $p\|\theta\|_2 \geq 1$ (or ≤ -1)

SVM : MAXIMIZING THE MARGIN



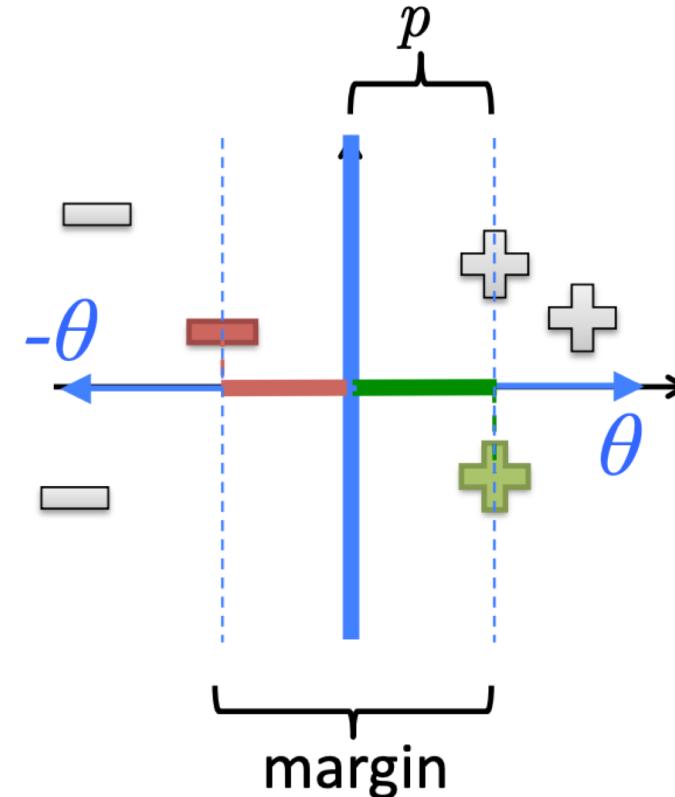
Since p is larger, $\|\theta\|_2$ can be smaller
in order to have $p\|\theta\|_2 \geq 1$ (or ≤ -1)



SVM : SIZE OF THE MARGIN

For the support vectors, we have $p\|\theta\|_2 = \pm 1$

- p is the length of the projection of the SVs onto θ



Therefore,

$$p = \frac{1}{\|\theta\|_2}$$

$$\text{margin} = 2p = \frac{2}{\|\theta\|_2}$$

SVM : THE DUAL PROBLEM

The primal SVM problem was given as


$$\begin{aligned} \min_{\theta} \quad & \frac{1}{2} \sum_{j=1}^d \theta_j^2 \\ \text{s.t.} \quad & y_i(\theta^\top \mathbf{x}_i) \geq 1 \quad \forall i \end{aligned}$$

Can solve it more efficiently by taking the Lagrangian dual



- Duality is a common idea in optimization
- It transforms a difficult optimization problem into a simpler one
- Key idea: introduce slack variables α_i for each constraint
 - α_i indicates how important a particular constraint is to the solution

SVM : THE DUAL PROBLEM

- The Lagrangian is given by



$$L(\theta, \alpha) = \frac{1}{2} \sum_{j=1}^d \theta_j^2 - \sum_{i=1}^n \alpha_i (y_i \theta^\top \mathbf{x} - 1)$$

$$\text{s.t. } \alpha_i \geq 0 \quad \forall i$$

- We must minimize over θ and maximize over α
- At optimal solution, partials w.r.t θ 's are 0



Solve by a bunch of algebra and calculus ...
and we obtain ...

SVM : THE DUAL REPRESENTATION



$$\begin{aligned} \text{Maximize } J(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{s.t. } \alpha_i &\geq 0 \quad \forall i \\ \sum_i \alpha_i y_i &= 0 \end{aligned}$$



The decision function is given by

$$h(\mathbf{x}) = \text{sign} \left(\sum_{i \in \mathcal{SV}} \alpha_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b \right)$$

$$\text{where } b = \frac{1}{|\mathcal{SV}|} \sum_{i \in \mathcal{SV}} \left(y_i - \sum_{j \in \mathcal{SV}} \alpha_j y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right)$$

SVM : UNDERSTANDING THE DUAL



$$\text{Maximize } J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{s.t. } \alpha_i \geq 0 \quad \forall i$$

$$\sum_i \alpha_i y_i = 0$$



Balances between the weight of constraints for different classes

Constraint weights (α_i 's) cannot be negative

SVM : UNDERSTANDING THE DUAL



$$\text{Maximize } J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$
$$\text{s.t. } \alpha_i \geq 0 \quad \forall i$$
$$\sum \alpha_i y_i = 0$$

Points with different labels increase the sum

Points with same label decrease the sum

Measures the similarity between points



Intuitively, we should be more careful around points near the margin

SVM : UNDERSTANDING THE DUAL



$$\begin{aligned} \text{Maximize } J(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{s.t. } \alpha_i &\geq 0 \quad \forall i \\ \sum_i \alpha_i y_i &= 0 \end{aligned}$$



In the solution, either:

- $\alpha_i > 0$ and the constraint is tight ($y_i(\theta^\top \mathbf{x}_i) = 1$)
 - point is a support vector
- $\alpha_i = 0$
 - point is not a support vector

SVM : EMPLOYING THE SOLUTION

- Given the optimal solution α^* , optimal weights are


$$\theta^* = \sum_{i \in SVs} \alpha_i^* y_i \mathbf{x}_i$$

- In this formulation, have *not* added $x_0 = 1$

- Therefore, we can solve one of the SV constraints

$$y_i (\theta^* \cdot \mathbf{x}_i + \theta_0) = 1$$



to obtain θ_0

- Or, more commonly, take the average solution over all support vectors

SVM : WHAT IF DATA ARE NOT LINEARLY SEPARABLE?



- Cannot find θ that satisfies $y_i(\theta^\top \mathbf{x}_i) \geq 1 \quad \forall i$
- Introduce slack variables ξ_i

$$y_i(\theta^\top \mathbf{x}_i) \geq 1 - \xi_i \quad \forall i$$



- New problem:

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2 + C \sum_i \xi_i$$

$$\text{s.t. } y_i(\theta^\top \mathbf{x}_i) \geq 1 - \xi_i \quad \forall i$$

SVM : STRENGTHS



- Good generalization in theory
- Good generalization in practice
- Work well with few training instances
- Find globally best model
- Efficient algorithms
- Amenable to the kernel trick ...



SVM : WHAT IF SURFACE IS NON-LINEAR?

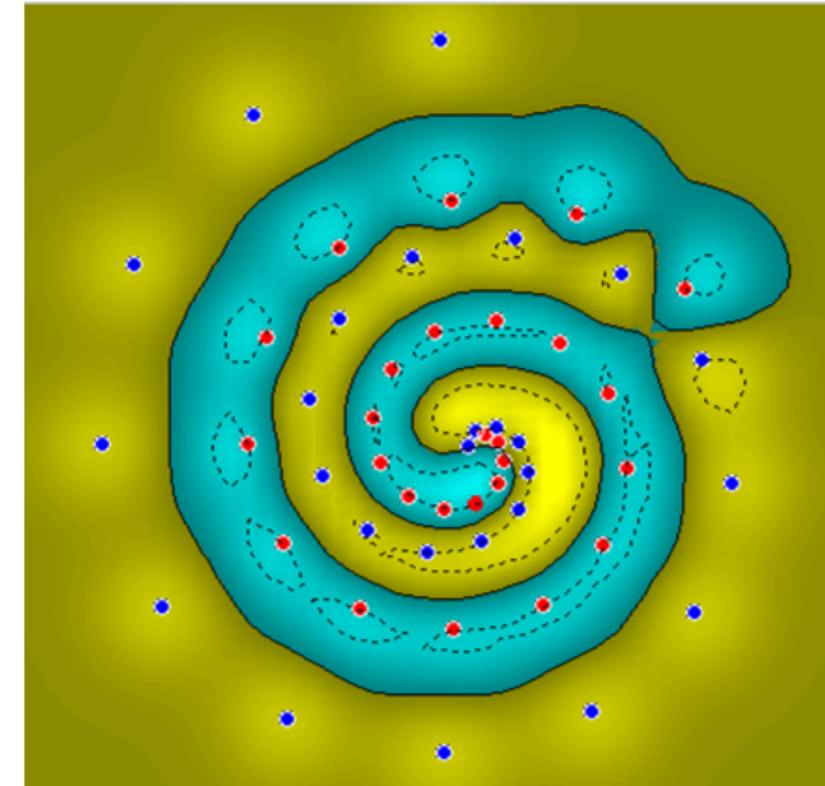
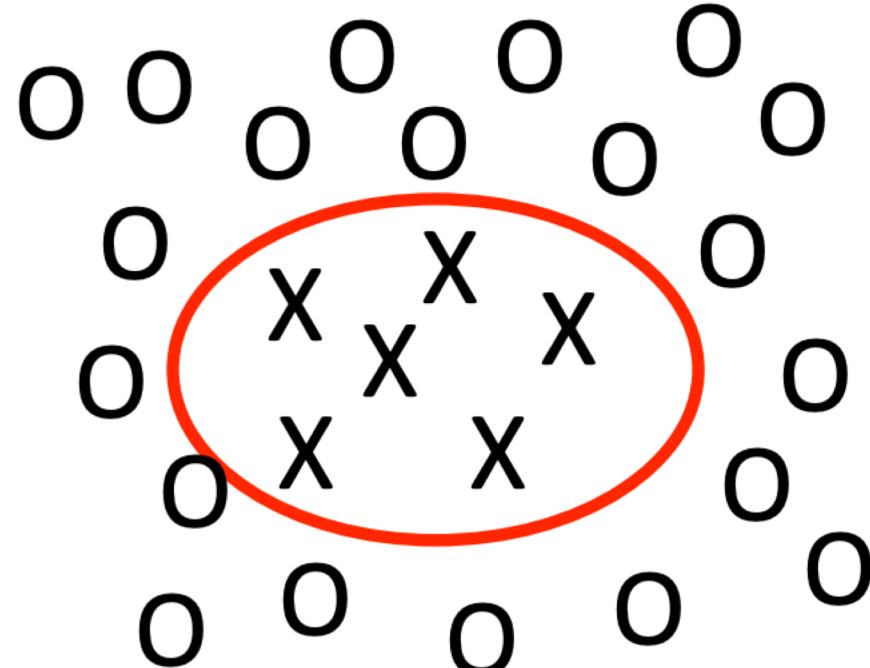


Image from <http://www.atrandomresearch.com/iclass/>

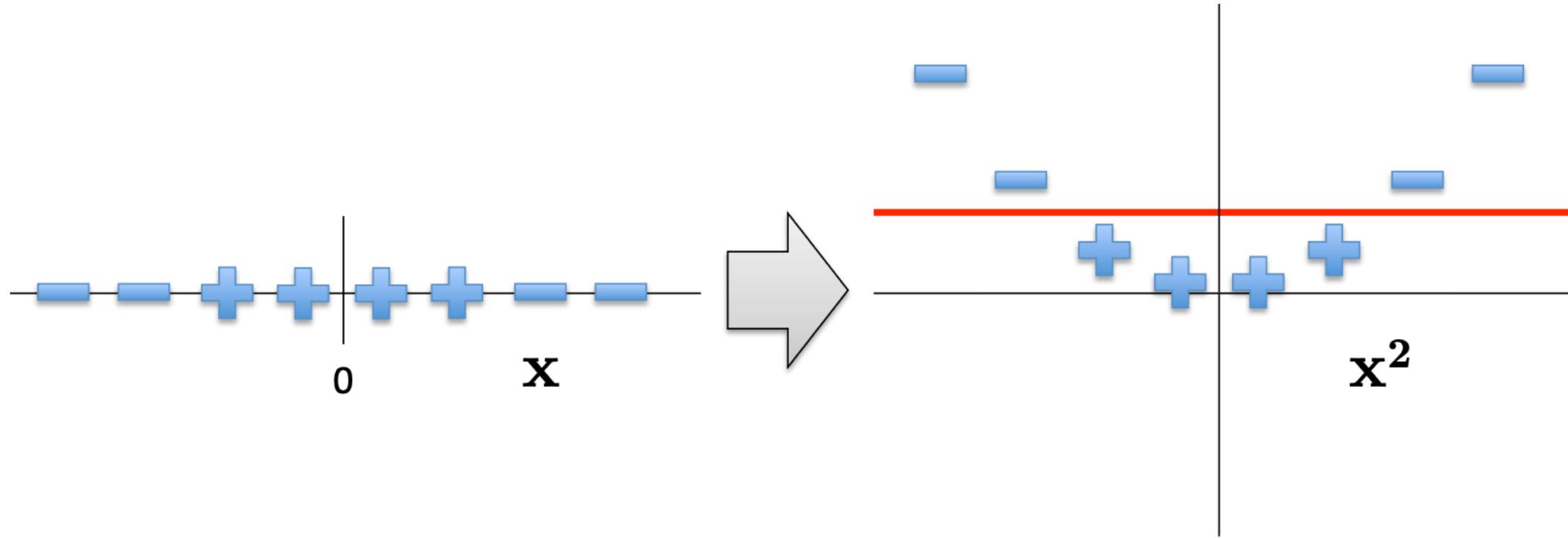


Kernerl Methods

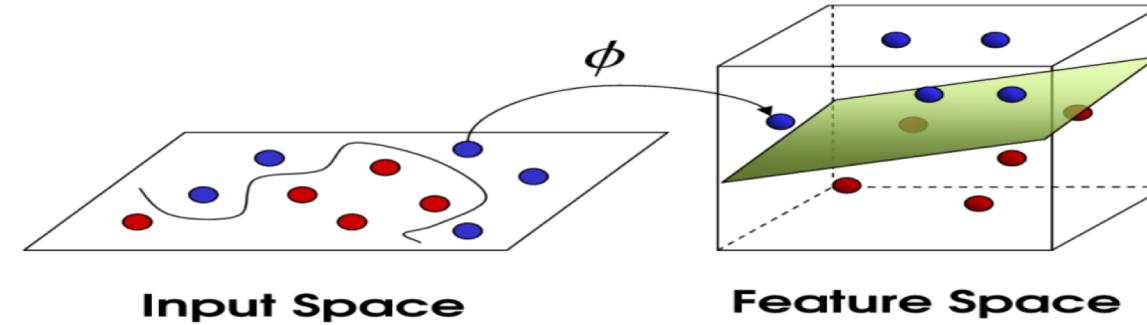
Making the Non-Linear Linear



KERNEL METHODS: WHEN LINEAR SEPARATORS FAIL



KERNEL METHODS: MAPPING INTO A NEW FEATURE SPACE



$$\Phi : \mathcal{X} \mapsto \hat{\mathcal{X}} = \Phi(\mathbf{x})$$

- For example, with $\mathbf{x}_i \in \mathbb{R}^2$
$$\Phi([x_{i1}, x_{i2}]) = [x_{i1}, x_{i2}, x_{i1}x_{i2}, x_{i1}^2, x_{i2}^2]$$
- Rather than run SVM on \mathbf{x}_i , run it on $\Phi(\mathbf{x}_i)$
 - Find non-linear separator in input space
- What if $\Phi(\mathbf{x}_i)$ is really big?
- Use kernels to compute it implicitly!



KERNEL METHODS



- Find kernel K such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$$

- Computing $K(\mathbf{x}_i, \mathbf{x}_j)$ should be efficient, much more so than computing $\Phi(\mathbf{x}_i)$ and $\Phi(\mathbf{x}_j)$



- Use $K(\mathbf{x}_i, \mathbf{x}_j)$ in SVM algorithm rather than $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- Remarkably, this is possible!

THE POLYNOMIAL KERNEL



Let $\mathbf{x}_i = [x_{i1}, x_{i2}]$ and $\mathbf{x}_j = [x_{j1}, x_{j2}]$

Consider the following function:

$$\begin{aligned} K(\mathbf{x}_i, \mathbf{x}_j) &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle^2 \\ &= (x_{i1}x_{j1} + x_{i2}x_{j2})^2 \\ &= (x_{i1}^2x_{j1}^2 + x_{i2}^2x_{j2}^2 + 2x_{i1}x_{i2}x_{j1}x_{j2}) \\ &= \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle \end{aligned}$$



where

$$\Phi(\mathbf{x}_i) = [x_{i1}^2, x_{i2}^2, \sqrt{2}x_{i1}x_{i2}]$$

$$\Phi(\mathbf{x}_j) = [x_{j1}^2, x_{j2}^2, \sqrt{2}x_{j1}x_{j2}]$$

THE POLYNOMIAL KERNEL



- Given by $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle^d$
 - $\Phi(\mathbf{x})$ contains all monomials of degree d



- Useful in visual pattern recognition
 - Example:
 - 16x16 pixel image
 - 10^{10} monomials of degree 5
 - Never explicitly compute $\Phi(\mathbf{x})$!
- Variation: $K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + 1)^d$
 - Adds all lower-order monomials (degrees 1,...,d)!

THE KERNEL TRICK



“Given an algorithm which is formulated in terms of a positive definite kernel K_1 , one can construct an alternative algorithm by replacing K_1 with another positive definite kernel K_2 ”

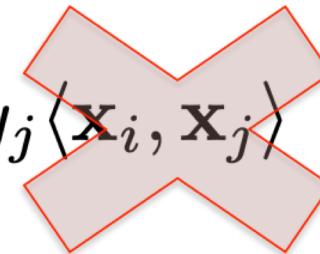


➤ SVMs can use the kernel trick

INCORPORATING KERNELS INTO SVM



$$J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$



$$J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$



$$\text{s.t. } \alpha_i \geq 0 \quad \forall i$$

$$\sum_i \alpha_i y_i = 0$$

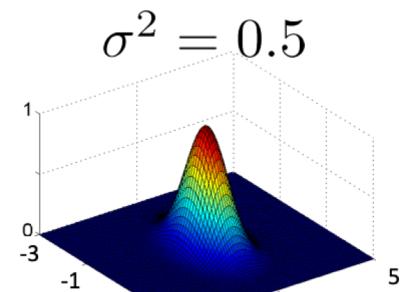
THE GAUSSIAN KERNEL

- Also called Radial Basis Function (RBF) kernel

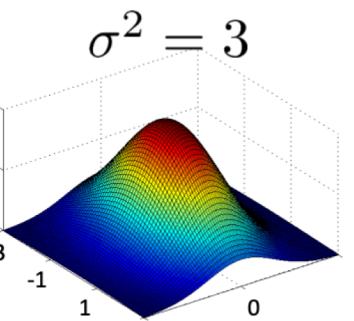
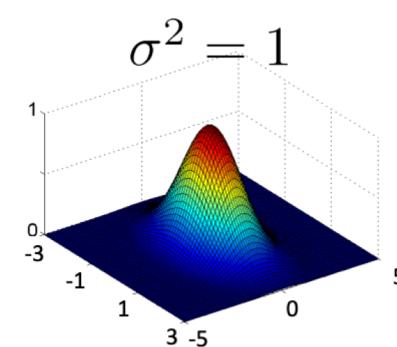


$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

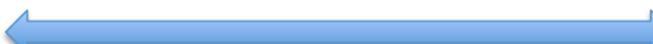
- Has value 1 when $\mathbf{x}_i = \mathbf{x}_j$
- Value falls off to 0 with increasing distance
- Note: Need to do feature scaling before using Gaussian Kernel



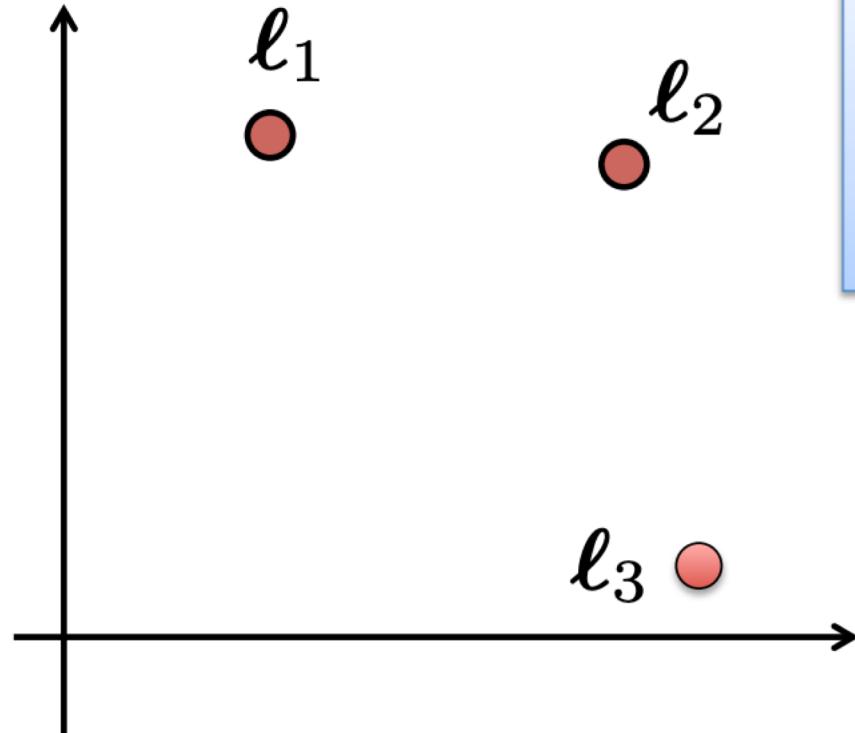
lower bias,
higher variance



higher bias,
lower variance



THE GAUSSIAN KERNEL EXAMPLE



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

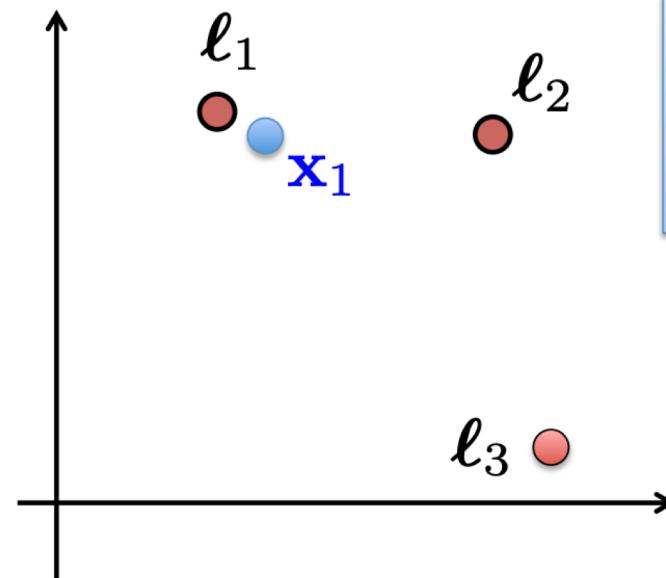
Imagine we've learned that:

$$\theta = [-0.5, 1, 1, 0]$$

Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \ell_1) + \theta_2 K(\mathbf{x}, \ell_2) + \theta_3 K(\mathbf{x}, \ell_3) \geq 0$



THE GAUSSIAN KERNEL EXAMPLE



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:

ℓ_3

$$\boldsymbol{\theta} = [-0.5, 1, 1, 0]$$

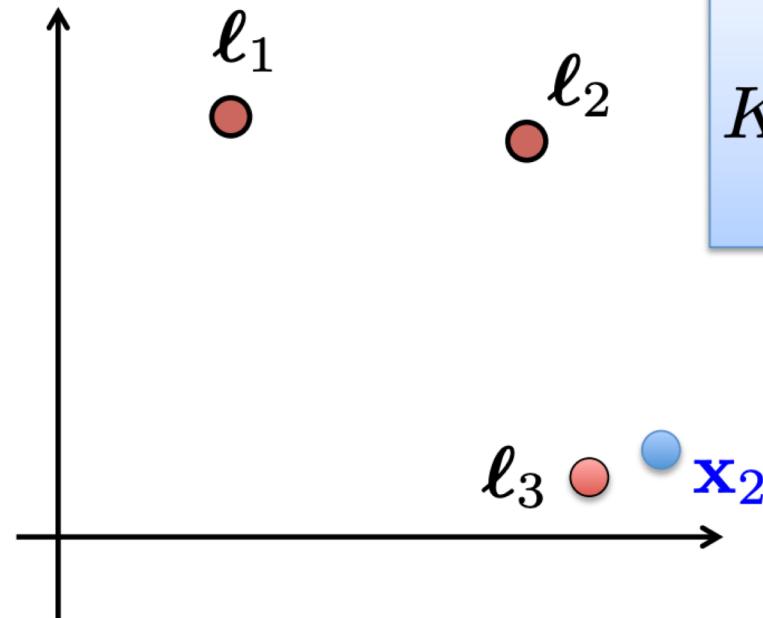
Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \ell_1) + \theta_2 K(\mathbf{x}, \ell_2) + \theta_3 K(\mathbf{x}, \ell_3) \geq 0$



- For \mathbf{x}_1 , we have $K(\mathbf{x}_1, \ell_1) \approx 1$, other similarities ≈ 0

$$\begin{aligned}\theta_0 + \theta_1(1) + \theta_2(0) + \theta_3(0) \\ = -0.5 + 1(1) + 1(0) + 0(0) \\ = 0.5 \geq 0, \text{ so predict +1}\end{aligned}$$

THE GAUSSIAN KERNEL EXAMPLE



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:

$$\ell_3 \bullet \mathbf{x}_2$$

$$\boldsymbol{\theta} = [-0.5, 1, 1, 0]$$

Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \ell_1) + \theta_2 K(\mathbf{x}, \ell_2) + \theta_3 K(\mathbf{x}, \ell_3) \geq 0$

- For \mathbf{x}_2 , we have $K(\mathbf{x}_2, \ell_3) \approx 1$, other similarities ≈ 0

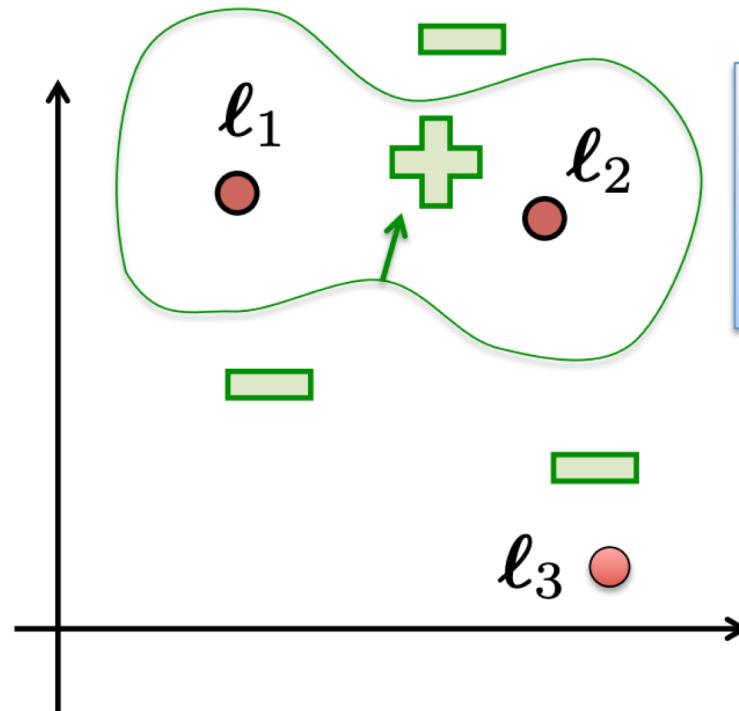
$$\theta_0 + \theta_1(0) + \theta_2(0) + \theta_3(1)$$

$$= -0.5 + 1(0) + 1(0) + 0(1)$$

$= -0.5 < 0$, so predict -1



THE GAUSSIAN KERNEL EXAMPLE



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:

$$\boldsymbol{\theta} = [-0.5, 1, 1, 0]$$



Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \ell_1) + \theta_2 K(\mathbf{x}, \ell_2) + \theta_3 K(\mathbf{x}, \ell_3) \geq 0$

Rough sketch of decision surface

OTHER KERNELS

- Sigmoid Kernel


$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\alpha \mathbf{x}_i^\top \mathbf{x}_j + c)$$

- Neural networks use sigmoid as activation function
- SVM with a sigmoid kernel is equivalent to 2-layer perceptron

- Cosine Similarity Kernel


$$K(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{x}_i^\top \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$$

- Popular choice for measuring similarity of text documents
- L₂ norm projects vectors onto the unit sphere; their dot product is the cosine of the angle between the vectors

OTHER KERNELS

- Chi-squared Kernel



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp \left(-\gamma \sum_k \frac{(x_{ik} - x_{jk})^2}{x_{ik} + x_{jk}} \right)$$

- Widely used in computer vision applications
- Chi-squared measures distance between probability distributions
- Data is assumed to be non-negative, often with L_1 norm of 1



- String kernels
- Tree kernels
- Graph kernels

BONUS : THE MATH BEHIND KERNELS



What does it *mean* to be a kernel?

- $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$ for some Φ

What does it *take* to be a kernel?

- The Gram matrix $G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
 - Symmetric matrix
 - Positive semi-definite matrix:
 $\mathbf{z}^T G \mathbf{z} \geq 0$ for every non-zero vector $\mathbf{z} \in \mathbb{R}^n$



Establishing “kernel-hood” from first principles is non-trivial

A FEW GOOD KERNELS...



- Linear Kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- Polynomial kernel $K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + c)^d$
– $c \geq 0$ trades off influence of lower order terms
- Gaussian kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$
- Sigmoid kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\alpha \mathbf{x}_i^\top \mathbf{x}_j + c)$



Many more...

- Cosine similarity kernel
- Chi-squared kernel
- String/tree/graph/wavelet/etc kernels

PRACTICAL ADVICE FOR APPLYING SVMS

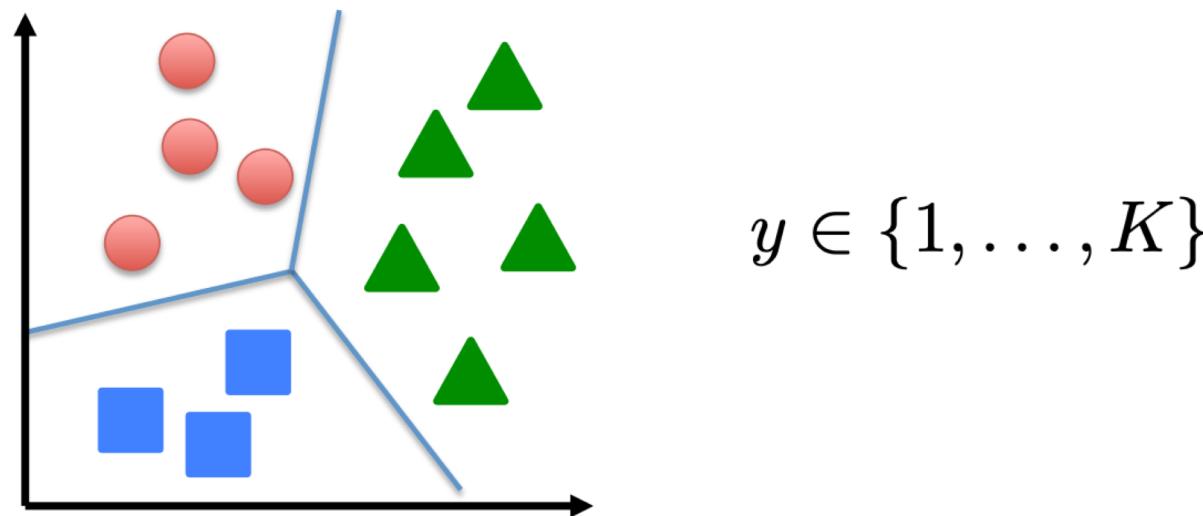
- Use SVM software package to solve for parameters
 - e.g., SVMlight, libsvm, cvx (fast!), etc.
- Need to specify:
 - Choice of parameter C
 - Choice of kernel function
 - Associated kernel parameters



$$\text{e.g., } K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + c)^d$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

MULTI-CLASS CLASSIFICATION WITH SVMS



- Many SVM packages already have multi-class classification built in
- Otherwise, use one-vs-rest
 - Train K SVMs, each picks out one class from rest, yielding $\theta^{(1)}, \dots, \theta^{(K)}$
 - Predict class i with largest $(\theta^{(i)})^\top \mathbf{x}$

SVMS VS LOGISTIC REGRESSION

n = # training examples d = # features



If d is large (relative to n) (e.g., $d > n$ with $d = 10,000$, $n = 10-1,000$)

- Use logistic regression or SVM with a linear kernel

If d is small (up to 1,000), n is intermediate (up to 10,000)

- Use SVM with Gaussian kernel

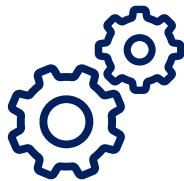
If d is small (up to 1,000), n is large (50,000+)

- Create/add more features, then use logistic regression or SVM without a kernel



Neural networks likely to work well for most of these settings, but may be slower to train

OTHER SVM VARIATIONS



- nu SVM
 - nu parameter controls:
 - Fraction of support vectors (lower bound) and misclassification rate (upper bound)
 - E.g., $\nu = 0.05$ guarantees that $\geq 5\%$ of training points are SVs and training error rate is $\leq 5\%$
 - Harder to optimize than C-SVM and not as scalable
 - SVMs for regression
 - One-class SVMs
 - SVMs for clustering
- ...



CONCLUSION



- SVMs find optimal linear separator
 - The kernel trick makes SVMs learn non-linear decision surfaces
-
- Strength of SVMs:
 - Good theoretical and empirical performance
 - Supports many types of kernels
-
- 
- Disadvantages of SVMs:
 - “Slow” to train/predict for huge data sets (but relatively fast!)
 - Need to choose the kernel (and tune its parameters)



www.keyrus.com

Shriman TIWARI

Tech Lead/Manager Data Science

Mobile: +33 (0)6 49 71 80 68
shriman.tiwari@keyrus.com

KEYRUS