

# Lecture 15:

## Coordinate Descent (continued)

---

- How to solve non-smooth optimization like Lasso?

$$\hat{w}_{\text{Lasso}} = \arg \min_{w \in \mathbb{R}^d} \underbrace{\|y - Xw\|_2^2 + \lambda \|w\|_1}_{f(w)}$$

W

# Coordinate descent for Lasso

- let us apply coordinate descent on Lasso, which minimizes
$$\text{minimize}_w \mathcal{L}(w) + \lambda \|w\|_1 = \|\mathbf{X}w - \mathbf{y}\|_2^2 + \lambda \|w\|_1$$
- the goal is to derive an **analytical rule** for updating  $w_j^{(t)}$ 's
- let us first write the update rule explicitly for  $w_1^{(t)}$ 
  - first step is to write the loss in terms of  $w_1$

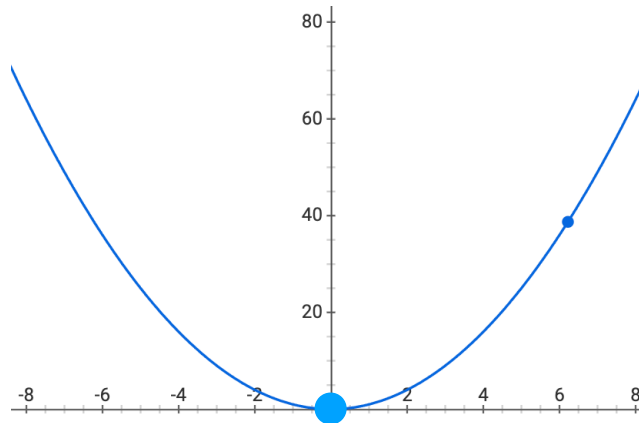
$$\left\| \mathbf{X}[:,1]w_1 - (\mathbf{y} - \mathbf{X}[:,2:d]w_{2:d}) \right\|_2^2 + \lambda \left( |w_1| + \underbrace{\|w_{2:d}\|_1}_{\text{constant}} \right)$$

- hence, the coordinate descent update boils down to

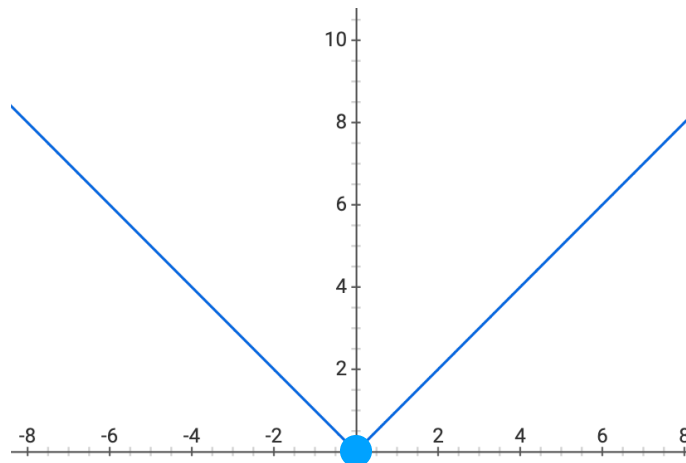
$$w_1^{(t)} \leftarrow \arg \min_{w_1} \underbrace{\left\| \mathbf{X}[:,1]w_1 - (\mathbf{y} - \mathbf{X}[:,2:d]w_{2:d}^{(t-1)}) \right\|_2^2 + \lambda |w_1|}_{f(w_1)}$$

# How do we find the minima?

- for **convex differentiable** functions, the minimum is achieved at points where gradient is zero



- for **convex non-differentiable** functions, the minimum is achieved at points where sub-gradient includes zero



## Finding the minima for $(aw_1 - b)^2 + \lambda |w_1|$

- the minimizer  $w_1^{(t)}$  is when zero is included in the sub-gradient

$$\partial f(w_1) = \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

## Finding the minima for $(aw_1 - b)^2 + \lambda |w_1|$

- the minimizer  $w_1^{(t)}$  is when zero is included in the sub-gradient

$$\partial f(w_1) = \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

## Finding the minima for $(aw_1 - b)^2 + \lambda |w_1|$

- the minimizer  $w_1^{(t)}$  is when zero is included in the sub-gradient

$$\partial f(w_1) = \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

## Finding the minima for $(aw_1 - b)^2 + \lambda |w_1|$

- considering all three cases, we get the following update rule by setting the sub-gradient to zero

$$w_1^{(t)} \leftarrow \begin{cases} \frac{b}{a} - \frac{\lambda}{2a^2} & \text{for } 2ab > \lambda \\ 0 & \text{for } -\lambda \leq 2ab \leq \lambda \\ \frac{b}{a} + \frac{\lambda}{2a^2} & \text{for } \lambda < -2ab \end{cases} \iff \frac{-\lambda}{2a^2} \leq \frac{b}{a} \leq \frac{\lambda}{2a^2}$$

# How do we find the minimizer?

- the minimizer  $w_1^{(t)}$  is when zero is included in the sub-gradient

$$\partial f(w_1) = \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

- case 1:

- $2a(aw_1 - b) + \lambda = 0$  for some  $w_1 > 0$

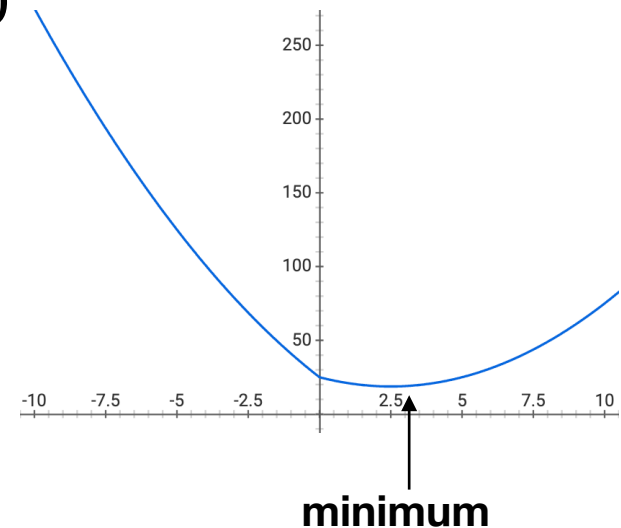
- this happens when

$$w_1 = \frac{-\lambda + 2ab}{2a^2} > 0$$

- hence,

$$w_1^{(t)} \leftarrow \frac{b}{a} - \frac{\lambda}{2a^2},$$

if  $\lambda < 2ab$





- case 2:

- $2a(aw_1 - b) - \lambda = 0$  for some  $w_1 < 0$

- this happens when

$$w_1 = \frac{\lambda + 2ab}{2a^2} < 0$$

- hence,

$$w_1^{(t)} \leftarrow \frac{b}{a} + \frac{\lambda}{2a^2},$$

if  $\lambda < -2ab$

- case 3:

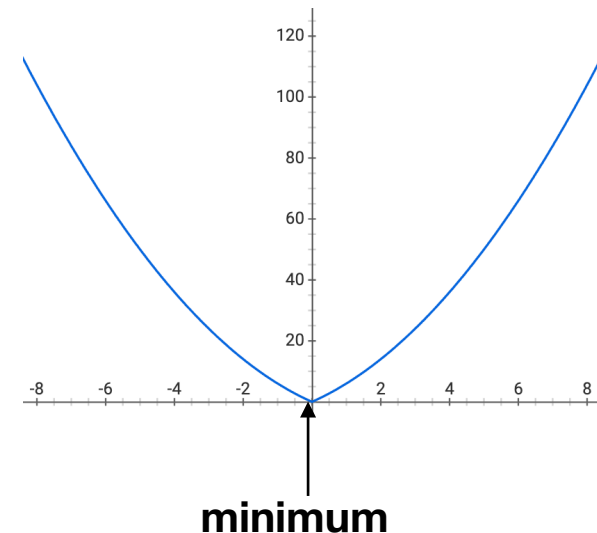
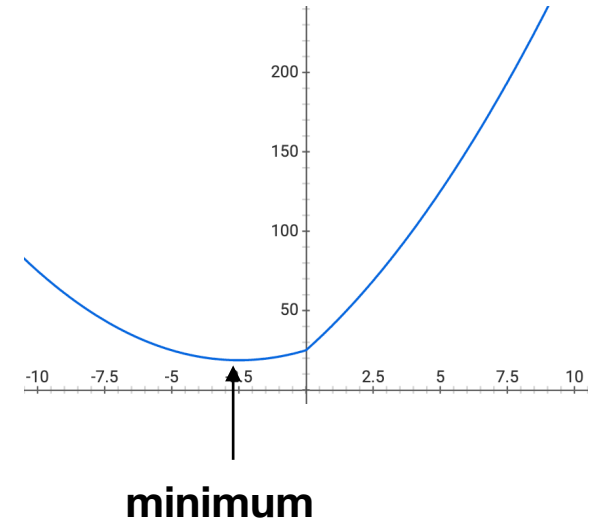
- $0 \in [-2ab - \lambda, -2ab + \lambda]$

- and  $w_1 = 0$

- hence,

$$w_1^{(t)} \leftarrow 0,$$

if  $-\lambda \leq 2ab \leq \lambda$

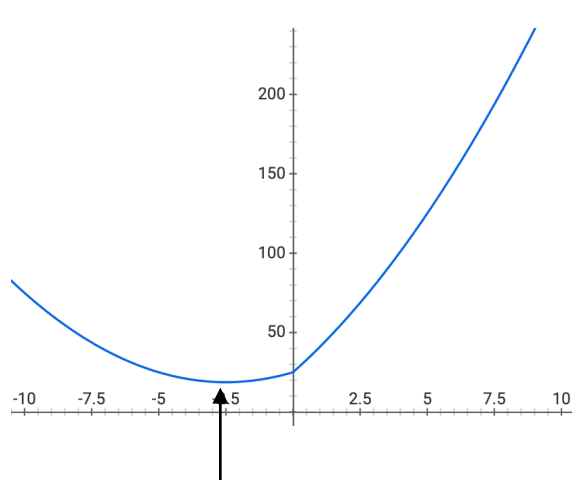


# Coordinate descent on Lasso

- considering all three cases, we get the following update rule by setting the sub-gradient to zero

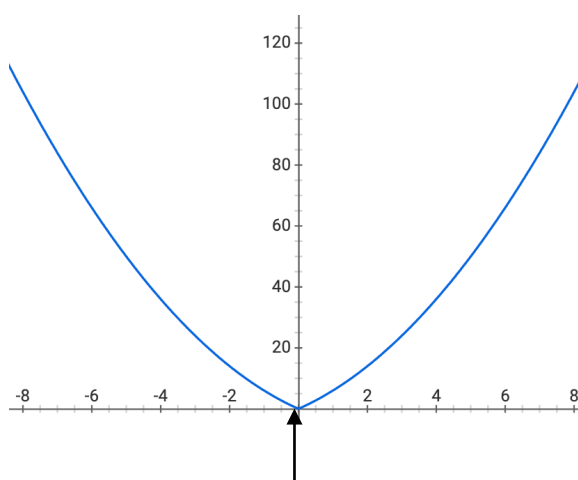
$$w_1^{(t)} \leftarrow \begin{cases} \frac{b}{a} - \frac{\lambda}{2a^2} & \text{for } 2ab > \lambda \\ 0 & \text{for } -\lambda \leq 2ab \leq \lambda \\ \frac{b}{a} + \frac{\lambda}{2a^2} & \text{for } \lambda < -2ab \end{cases}$$

• where  $a = \sqrt{\mathbf{X}[:, 1]^T \mathbf{X}[:, 1]}$ , and  $b = \frac{\mathbf{X}[:, 1]^T (\mathbf{y} - \mathbf{X}[:, 2:d] w_{-1})}{\sqrt{\mathbf{X}[:, 1]^T \mathbf{X}[:, 1]}}$

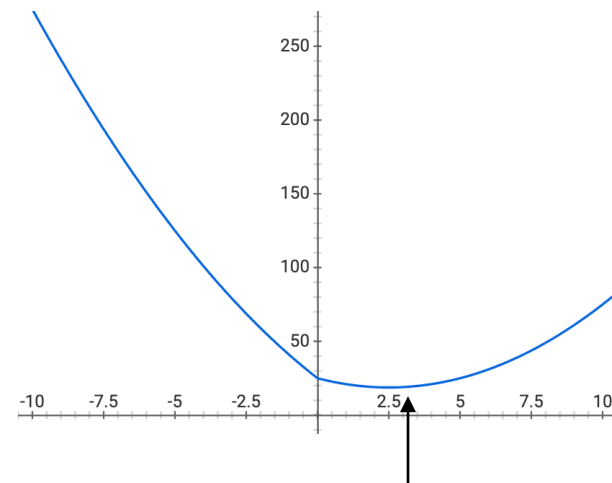


10

minimum



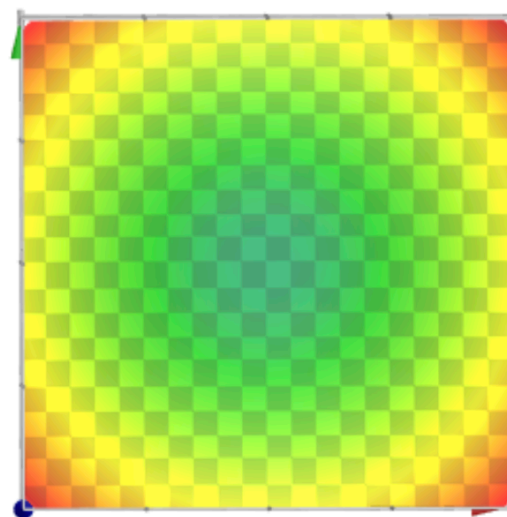
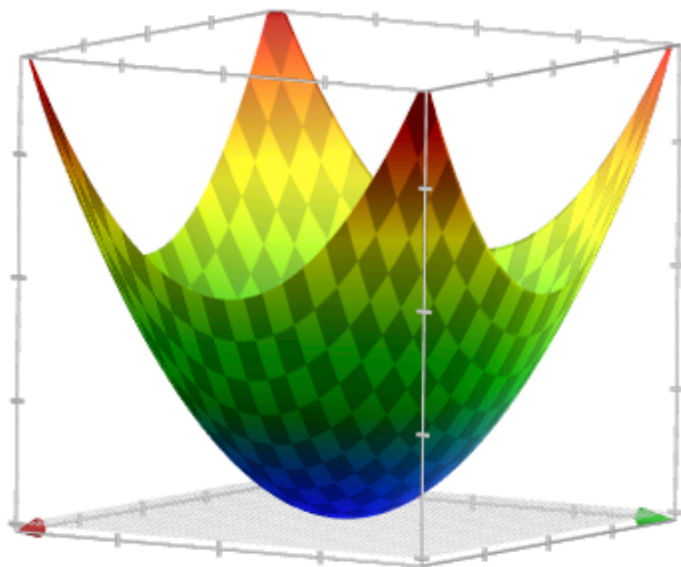
minimum



minimum

# When does coordinate descent work?

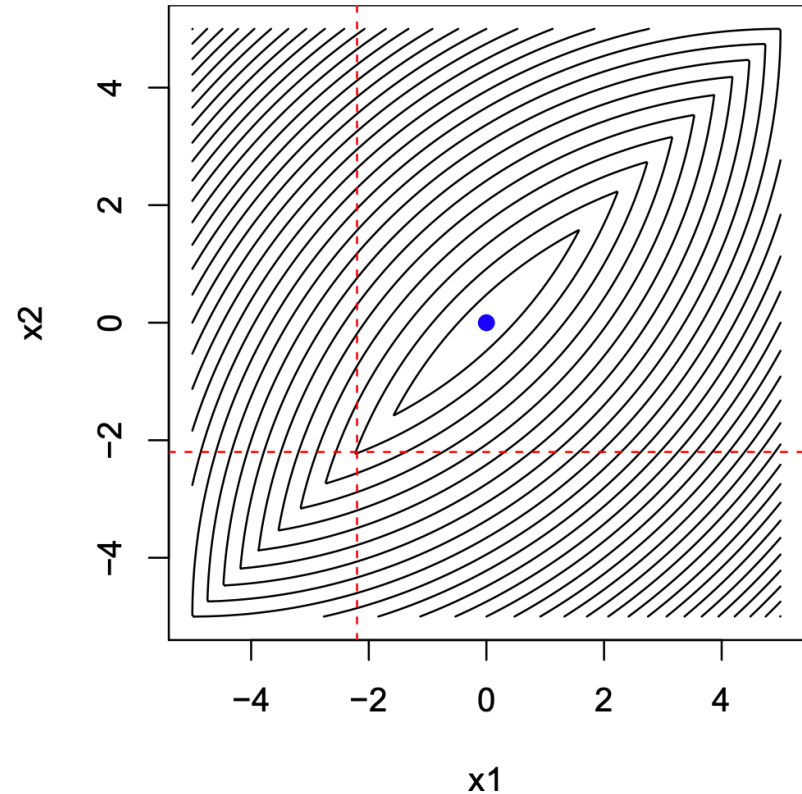
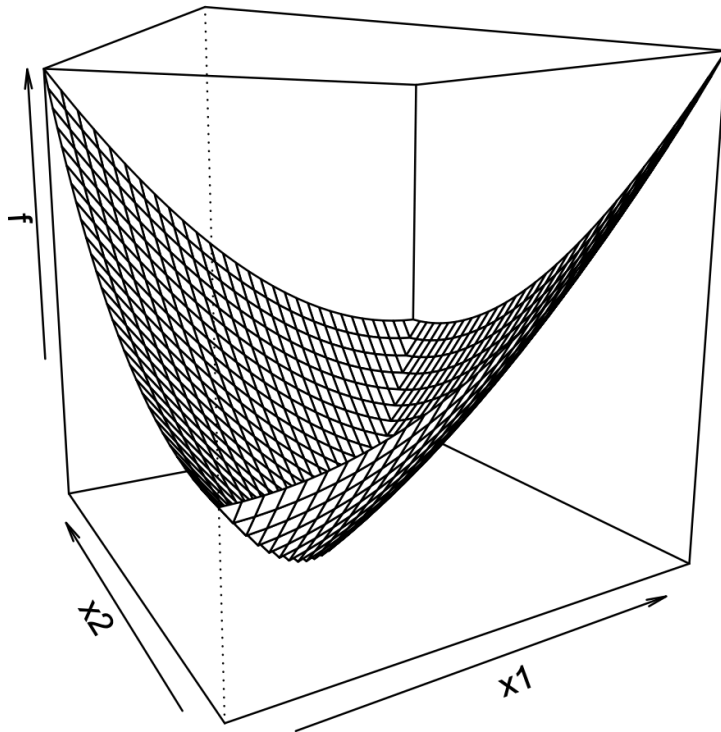
- Consider minimizing a **differentiable convex** function  $f(x)$ , then coordinate descent converges to the global minima



- when coordinate descent has stopped, that means  $\frac{\partial f(x)}{\partial x_j} = 0$  for all  $j \in \{1, \dots, d\}$
- this implies that the gradient  $\nabla_x f(x) = 0$ , which happens only at minimum

# When does coordinate descent work?

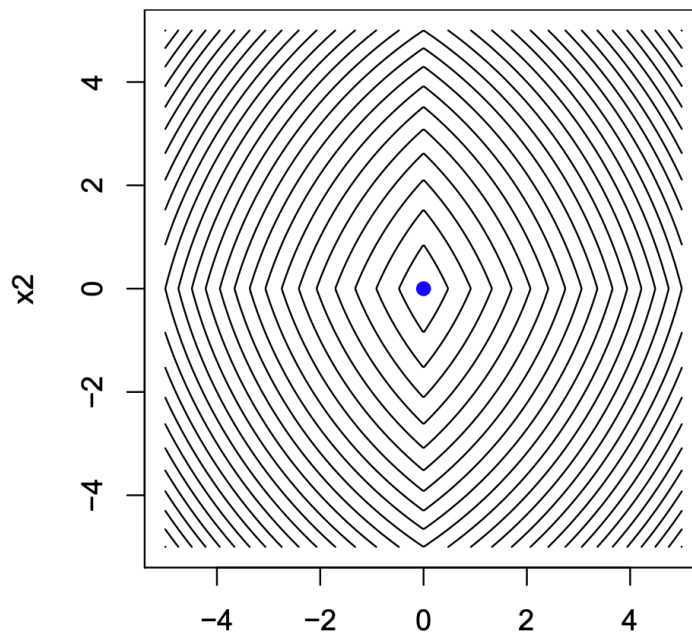
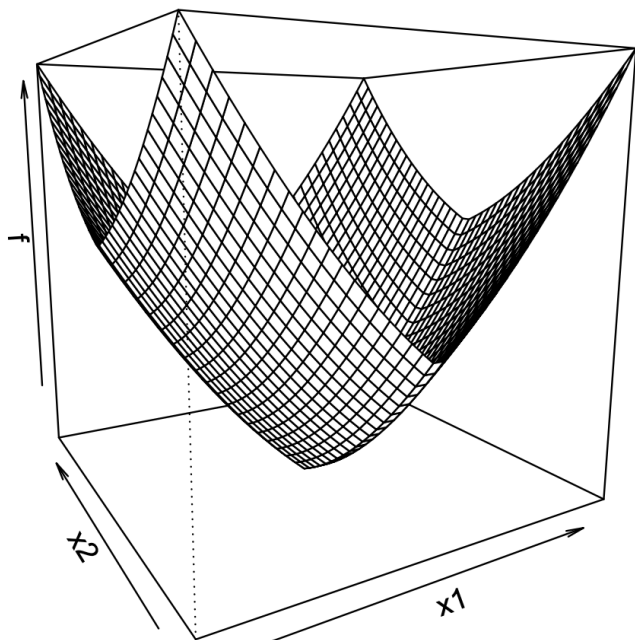
- Consider minimizing a **non-differentiable convex** function  $f(x)$ , then coordinate descent can get stuck



$$f(x_1, x_2) = (3x_1 + 4x_2 + 1)^2 + \lambda |x_1 - x_2|$$

# When does coordinate descent work?

- then how can coordinate descent find optimal solution for Lasso?
- consider minimizing a **non-differentiable convex** function but has a structure of  $f(x) = g(x) + \sum_{j=1}^d h_j(x_j)$ , with differentiable convex function  $g(x)$  and coordinate-wise non-differentiable convex functions  $h_j(x_j)$ 's, then coordinate descent converges to the global minima



$$f(x_1, x_2) = (3x_1 + 4x_2 + 1)^2 + \lambda|x_1| + \lambda|x_2|$$

# Questions?

---

# Lecture 16:

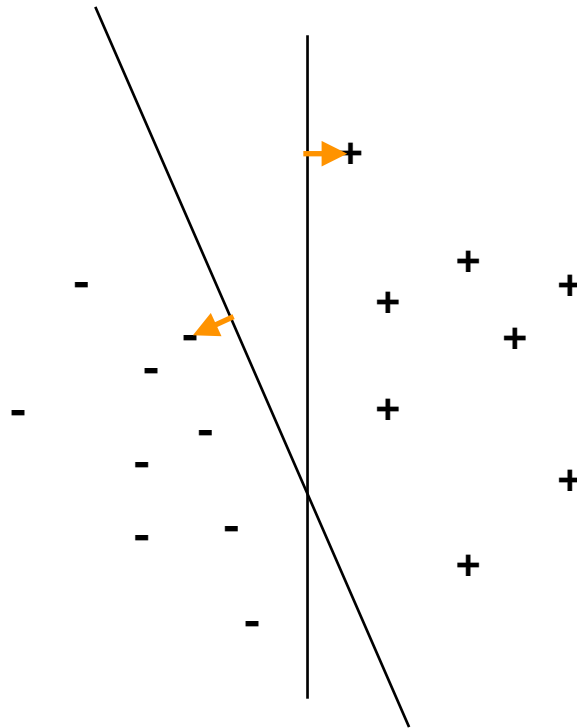
# Support Vector Machines

---



# How do we choose the best linear classifier?

- informally, **margin** of a set of examples to a decision boundary is the distance to the closest point to the decision boundary
- for linearly separable datasets, **maximum margin** classifier is a natural choice
- large margin implies that the decision boundary can change without losing accuracy, so the learned model is more robust against new data points





# Geometric margin

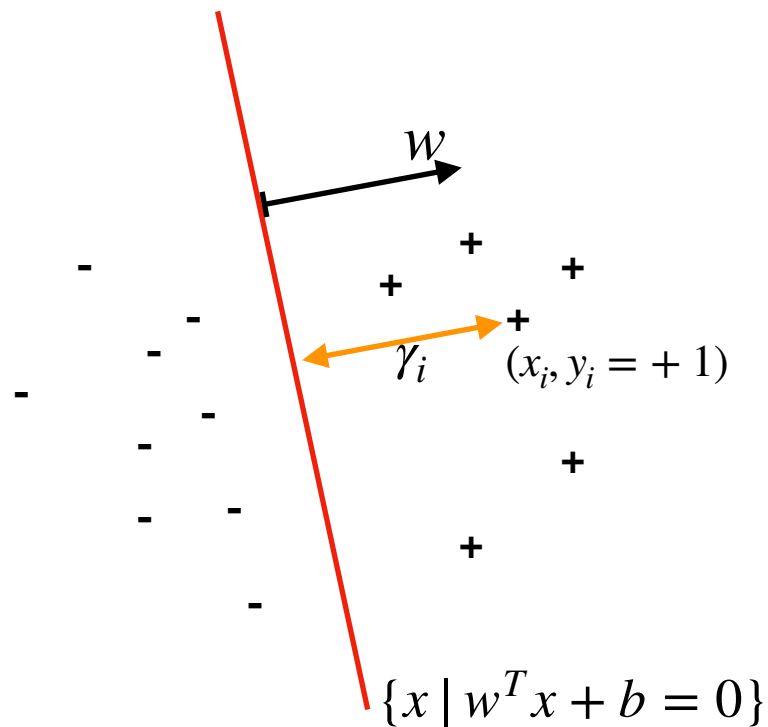
- given a set of training examples  $\{(x_i, y_i)\}_{i=1}^n$
- and a linear classifier  $(w, b) \in \mathbb{R}^d \times \mathbb{R}$
- such that the decision boundary is a separating hyperplane  $\{x \mid \underbrace{b + w_1x[1] + w_2x[2] + \dots + w_dx[d]}_{w^T x + b} = 0\},$

which is the set of points that are orthogonal to  $w$  with a shift of  $b$

- we define **functional margin** of  $(b, w)$  with respect to a training example  $(x_i, y_i)$  as the distance from the point  $(x_i, y_i)$  to the decision boundary, which is

$$\gamma_i = y_i \frac{(w^T x_i + b)}{\|w\|_2}$$

(The proof is on the next slide)



# Geometric margin

- the distance  $\gamma_i$  from a hyperplane  $\{x \mid w^T x + b = 0\}$  to a point  $x_i$  can be computed geometrically as follows
- We know that if you move from  $x_i$  in the negative direction of  $w$  by length  $\gamma_i$ , you arrive at the line, which can be written as

$$\left( x_i - \frac{w}{\|w\|_2} \gamma_i \right) \text{ is in } \{x \mid w^T x + b = 0\}$$

- so we can plug the point in the formula:

$$w^T \left( x_i - \frac{w}{\|w\|_2} \gamma_i \right) + b = 0$$

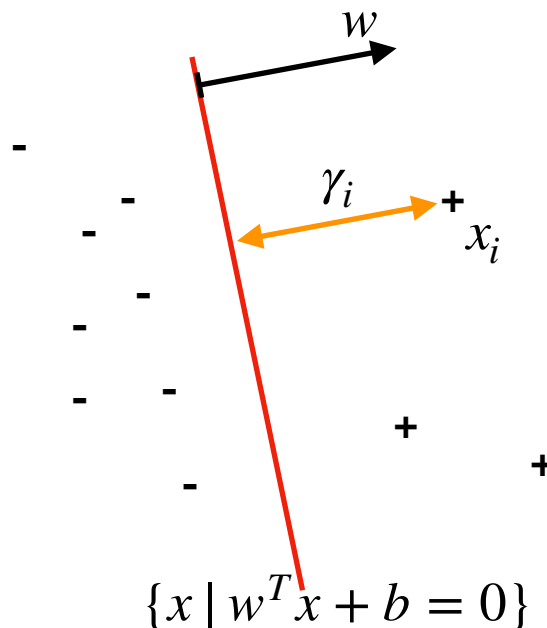
which is

$$w^T x_i - \frac{\|w\|_2^2}{\|w\|_2} \gamma_i + b = 0$$

and hence

$$\gamma_i = \frac{w^T x_i + b}{\|w\|_2},$$

and we multiply it by  $y_i$  so that for negative samples we use the opposite direction of  $-w$  instead of  $w$

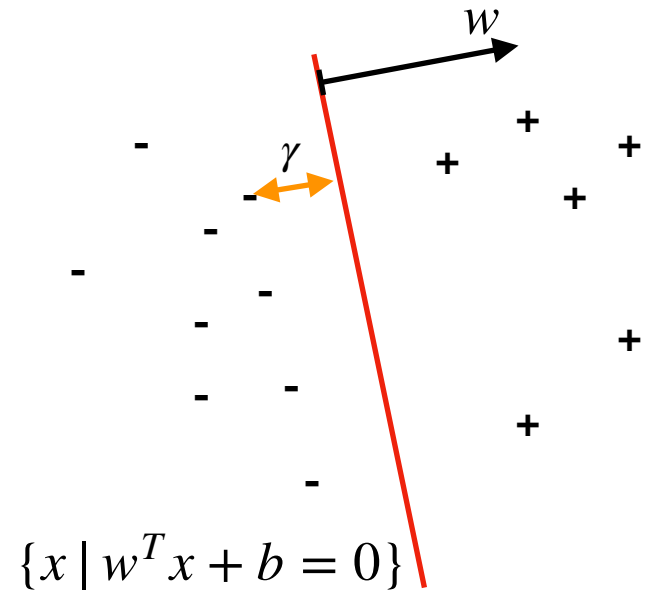


# Geometric margin

- the **margin** with respect to a set is defined as

$$\gamma = \min_{i=1}^n \gamma_i$$

- among all linear classifiers, we would like to find one that has the **maximum margin**



# Maximum margin classifier

- we propose the following optimization problem:

$$\text{maximize}_{w \in \mathbb{R}^d, b \in \mathbb{R}, \gamma \in \mathbb{R}} \quad \gamma$$

(maximize the margin)

$$\text{subject to} \quad \frac{y_i(w^T x_i + b)}{\|w\|_2} \geq \gamma \quad \text{for all } i \in \{1, \dots, n\}$$

(s.t.  $\gamma$  is a lower bound on the margin)

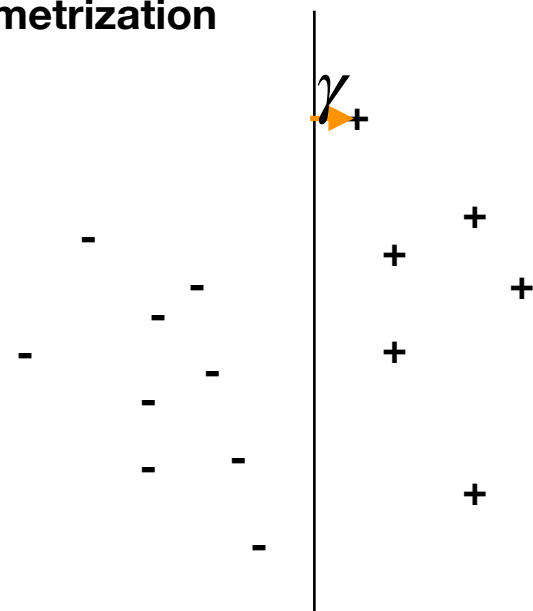
- if we fix  $(w, b)$ , the optimal solution of the optimization is the margin
- together with  $(w, b)$ , this finds the classifier with the maximum margin
- note that this problem is **scale invariant** in  $(w, b)$ , i.e. changing a  $(w, b)$  to  $(2w, 2b)$  does not change either the feasibility or the objective value, hence the following reparametrization is valid
- the above optimization looks difficult, so we transform it using **reparametrization**

$$\text{maximize}_{w \in \mathbb{R}^d, b \in \mathbb{R}, \gamma \in \mathbb{R}} \quad \gamma$$

$$\text{subject to} \quad \frac{y_i(w^T x_i + b)}{\|w\|_2} \geq \gamma \quad \text{for all } i \in \{1, \dots, n\}$$

$$\|w\|_2 = \frac{1}{\gamma}$$

- Because of scale invariance, the optimal solution does not change, as the solutions to the original problem did not depend on  $\|w\|_2$ , and only depends on the direction of  $w$



- maximize $_{w \in \mathbb{R}^d, b \in \mathbb{R}, \gamma \in \mathbb{R}} \quad \gamma$

$$\text{subject to } \frac{y_i(w^T x_i + b)}{\|w\|_2} \geq \gamma \text{ for all } i \in \{1, \dots, n\}$$

$$\|w\|_2 = \frac{1}{\gamma}$$

- the above optimization still looks difficult, but can be transformed into

$$\text{maximize}_{w \in \mathbb{R}^d, b \in \mathbb{R}} \quad \frac{1}{\|w\|_2} \quad \text{(maximize the margin)}$$

$$\text{subject to } \frac{y_i(w^T x_i + b)}{\|w\|_2} \geq \frac{1}{\|w\|_2} \text{ for all } i \in \{1, \dots, n\} \text{ (now } \frac{1}{\|w\|_2} \text{ plays the role of a lower bound on the margin)}$$

which simplifies to

$$\text{minimize}_{w \in \mathbb{R}^d, b \in \mathbb{R}} \quad \|w\|_2^2$$

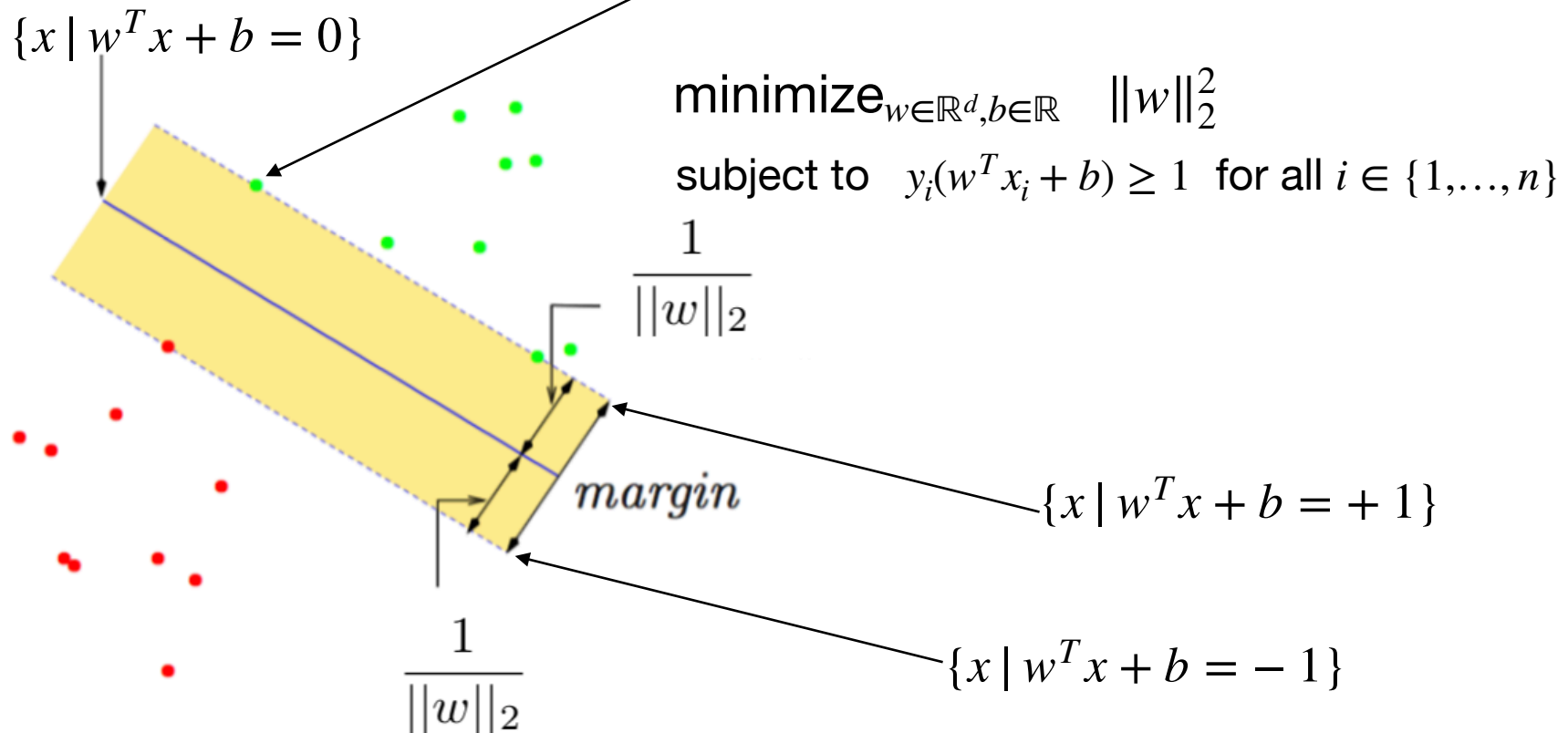
$$\text{subject to } y_i(w^T x_i + b) \geq 1 \text{ for all } i \in \{1, \dots, n\}$$

- this is a **quadratic program with linear constraints**, which can be easily solved

- once the optimal solution is found, the margin of that classifier  $(w, b)$  is  $\frac{1}{\|w\|_2}$

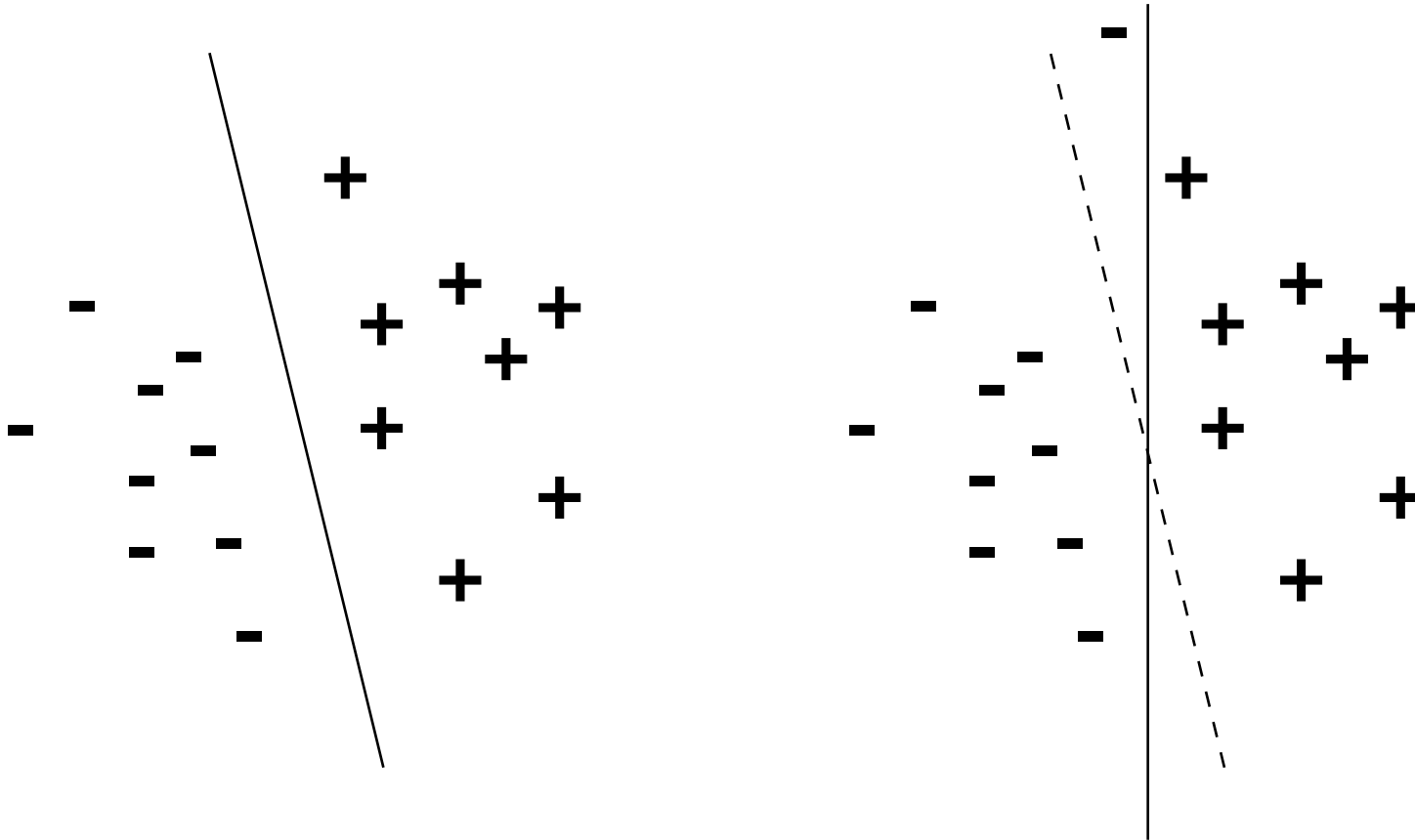
# What if the data is not separable?

- we cheated a little in the sense that the reparametrization of  $\|w\|_2 = \frac{1}{\gamma}$  is possible only if the margins are positive, i.e. the data is linearly separable with a positive margin
- otherwise, there is no feasible solution
- the examples at the margin are called **support vectors**



# Two issues

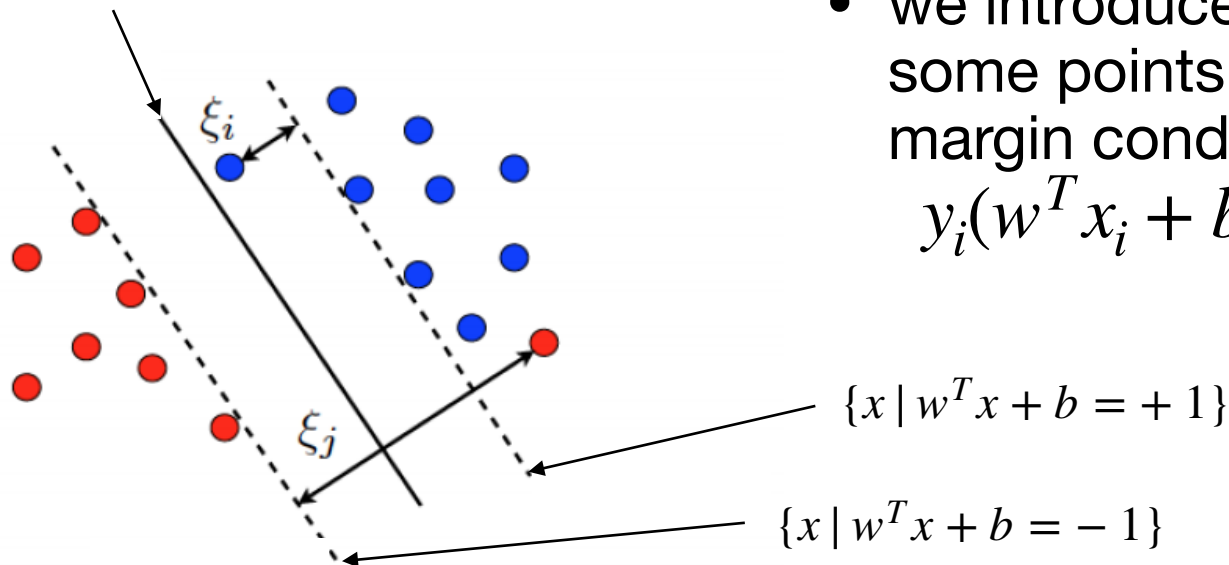
- max-margin formulation we proposed is sensitive to outliers



- it does not generalize to non-separable datasets

# What if the data is not separable?

$$\{x \mid w^T x + b = 0\}$$



- we introduce slack so that some points can violate the margin condition

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$

- this gives a new optimization problem with some positive constant  $c \in \mathbb{R}$

$$\text{minimize}_{w \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad \|w\|_2^2 + c \sum_{i=1}^n \xi_i$$

$$\text{subject to} \quad y_i(w^T x_i + b) \geq 1 - \xi_i \quad \text{for all } i \in \{1, \dots, n\}$$

$$\xi_i \geq 0 \quad \text{for all } i \in \{1, \dots, n\}$$

the (re-scaled) margin (for each sample) is allowed to be less than one, but you pay  $c\xi_i$  in the cost, and  $c$  balances the two goals: maximizing the margin for most examples vs. having small number of violations



# Support Vector Machine

- for the optimization problem

$$\text{minimize}_{w \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad \|w\|_2^2 + c \sum_{i=1}^n \xi_i$$

$$\text{subject to} \quad y_i(w^T x_i + b) \geq 1 - \xi_i \quad \text{for all } i \in \{1, \dots, n\}$$

$$\xi_i \geq 0 \quad \text{for all } i \in \{1, \dots, n\}$$

notice that at optimal solution,  $\xi_i$ 's satisfy

- $\xi_i = 0$  if margin is big enough  $y_i(w^T x_i + b) \geq 1$ , or
- $\xi_i = 1 - y_i(w^T x_i + b)$ , if the example is within the margin  $y_i(w^T x_i + b) < 1$

- so one can write

- $\xi_i = \max\{0, 1 - y_i(w^T x_i + b)\}$ , which gives

$$\text{minimize}_{w \in \mathbb{R}^d, b \in \mathbb{R}} \quad \frac{1}{c} \|w\|_2^2 + \sum_{i=1}^n \max\{0, 1 - y_i(w^T x_i + b)\}$$

# Sub-gradient descent for SVM

- SVM is the solution of

$$\text{minimize}_{w \in \mathbb{R}^d, b \in \mathbb{R}} \quad \frac{1}{c} \|w\|_2^2 + \sum_{i=1}^n \max\{0, 1 - y_i(w^T x_i + b)\}$$

- as it is non-differentiable, we solve it using sub-gradient descent
- which is exactly the same as gradient descent, except when we are at a non-differentiable point, we take one of the sub-gradients instead of the gradient (recall sub-gradient is a set)
- this means that we can take (a generic form derived from previous page)

$$\partial_w \ell(w^T x_i + b, y_i) = \mathbf{I}\{y_i(w^T x_i + b) \leq 1\}(-y_i x_i)$$

and apply

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \left( \sum_{i=1}^n \mathbf{I}\{y_i((w^{(t)})^T x_i + b^{(t)}) \leq 1\}(-y_i x_i) + \frac{2}{c} w^{(t)} \right)$$

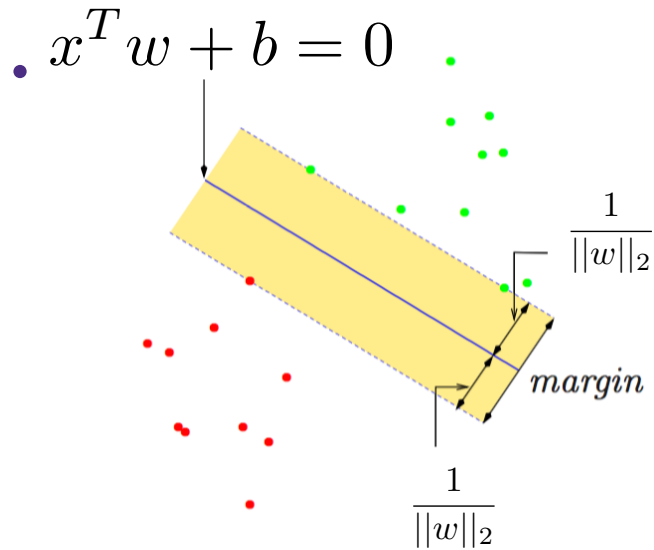
$$b^{(t+1)} \leftarrow b^{(t)} - \eta \sum_{i=1}^n \mathbf{I}\{y_i((w^{(t)})^T x_i + b^{(t)}) \leq 1\}(-y_i)$$

# Lecture 17: Kernels

---



# What if the data is not linearly separable?



some points don't satisfy margin constraint:

$$\min_{w,b} \|w\|_2^2$$

$$y_i(x_i^T w + b) \geq 1 \quad \forall i$$

Two options:

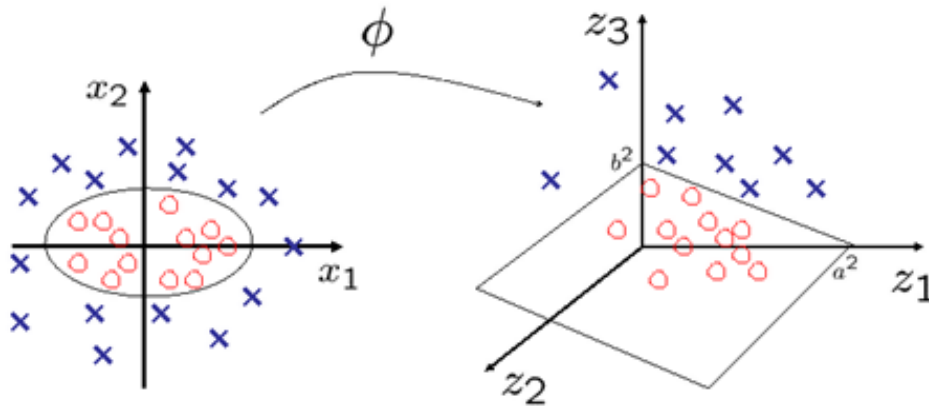
1. Introduce slack to this optimization problem (Support Vector Machine)
2. **Lift to higher dimensional space (Kernels)**

# What if the data is not linearly separable?

- Use features, for example,

$$x = (x_1, x_2) \in \mathbb{R}^2$$

$$z = (x_1, x_2, \sqrt{x_1^2 + x_2^2}) \in \mathbb{R}^3$$



- In high dimensional feature space, it is easier to linearly separate different classes

# Creating Features

- Feature mapping  $h : \mathbb{R}^d \rightarrow \mathbb{R}^p$  maps original data into a rich and high-dimensional feature space (usually  $d \ll p$ )

For example, in  $d=1$ , one can use

$$h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix} = \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^p \end{bmatrix}$$

For example, for  $d>1$ , one can generate vectors  $\{u_j\}_{j=1}^p \subset \mathbb{R}^d$

and define features:

$$h_j(x) = (u_j^T x)^2$$

$$h_j(x) = \frac{1}{1 + \exp(u_j^T x)}$$

$$h_j(x) = \cos(u_j^T x)$$

**Feature space can get really large really quickly!**

# Degree-d Polynomials

---

# How do we deal with high-dimensional lifts/data?

## A fundamental trick in ML: use kernels

A function  $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a *kernel* for a map  $\phi$  if  $K(x, x') = \phi(x) \cdot \phi(x')$  for all  $x, x'$ .

This notation is for dot product (which is the same as inner product)

- So, if we can represent our
  - training algorithms and
  - decision rules for prediction
- as functions of dot products of feature maps (i.e.  $\{\phi(x) \cdot \phi(x')\}$ )  
and if we can find a kernel for our feature map such that

$$K(x, x') = \phi(x) \cdot \phi(x')$$

then we can avoid explicitly computing (high-dimensional)  $\{\phi(x)\}$



# Linear Regression as Kernels

---

- Consider Ridge regression:  $\hat{w} = \arg \min_{w \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}w\|_2^2 + \lambda \|w\|_2^2$

## Dot-product of polynomials

$\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) =$  polynomials of degree exactly  $d$

$$d = 1 : \phi(u) = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \langle \phi(u), \phi(v) \rangle = u_1 v_1 + u_2 v_2$$

$$d = 2 : \phi(u) = \begin{bmatrix} u_1^2 \\ u_2^2 \\ u_1 u_2 \\ u_2 u_1 \end{bmatrix} \quad \langle \phi(u), \phi(v) \rangle = u_1^2 v_1^2 + u_2^2 v_2^2 + 2u_1 u_2 v_1 v_2$$

## Dot-product of polynomials

$\Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) =$  polynomials of degree exactly  $d$

$$d = 1 : \phi(u) = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \langle \phi(u), \phi(v) \rangle = u_1 v_1 + u_2 v_2$$

$$d = 2 : \phi(u) = \begin{bmatrix} u_1^2 \\ u_2^2 \\ u_1 u_2 \\ u_2 u_1 \end{bmatrix} \quad \langle \phi(u), \phi(v) \rangle = u_1^2 v_1^2 + u_2^2 v_2^2 + 2u_1 u_2 v_1 v_2$$

**Feature space can get really large really quickly!**

General  $d$  : Dimension of  $\phi(u)$  is roughly  $p^d$  if  $u \in \mathbb{R}^p$

Feature expansion can be written **implicitly**  $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^p$

# Examples of Kernels

- **Polynomials of degree exactly d**

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^p$$

- **Polynomials of degree up to d**

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^p$$

- **Gaussian (squared exponential) kernel**

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

- **Sigmoid**

$$K(u, v) = \tanh(\gamma \cdot u^T v + r)$$

# The Kernel Trick

---

**Pick a kernel  $K$**

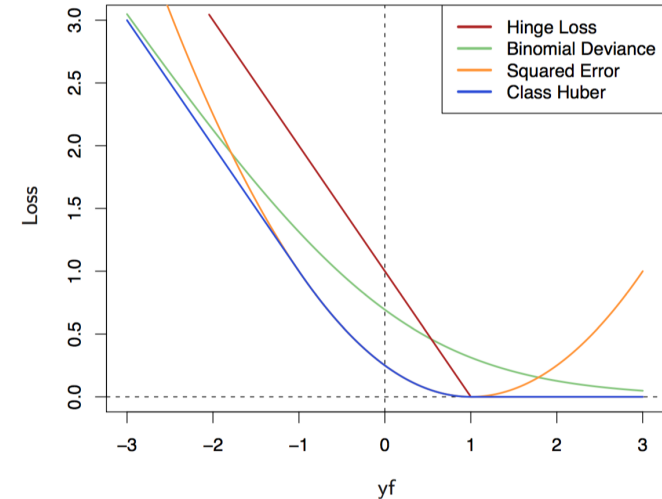
**For a linear predictor, show  $w = \sum_i \alpha_i x_i$**

**Change loss function/decision rule to only access data through dot products**

**Substitute  $K(x_i, x_j)$  for  $x_i^T x_j$**

# Loss Functions

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$



- Loss functions:

$$\sum_{i=1}^n \ell_i(w)$$

Squared error Loss:  $\ell_i(w) = (y_i - x_i^T w)^2$

Logistic Loss:  $\ell_i(w) = \log(1 + \exp(-y_i x_i^T w))$

0/1 loss:  $\ell_i(w) = \mathbb{I}[\text{sign}(y_i) \neq \text{sign}(x_i^T w)]$

Hinge Loss:  $\ell_i(w) = \max\{0, 1 - y_i x_i^T w\}$

# The Kernel Trick for regularized least squares

---

$$\hat{w} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda ||w||_w^2$$

There exists an  $\alpha \in \mathbb{R}^n$ :  $\hat{w} = \sum_{i=1}^n \alpha_i x_i$

# The Kernel Trick for regularized least squares

$$\hat{w} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \|w\|_w^2$$

There exists an  $\alpha \in \mathbb{R}^n$ :  $\hat{w} = \sum_{i=1}^n \alpha_i x_i$

$$\begin{aligned} \hat{\alpha} &= \arg \min_{\alpha} \sum_{i=1}^n (y_i - \sum_{j=1}^n \alpha_j \langle x_j, x_i \rangle)^2 + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \langle x_i, x_j \rangle \\ &= \arg \min_{\alpha} \sum_{i=1}^n (y_i - \sum_{j=1}^n \alpha_j K(x_i, x_j))^2 + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j) \\ &= \arg \min_{\alpha} \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^T \mathbf{K}\alpha \end{aligned}$$

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$



# Why regularization?

---

Typically,  $\mathbf{K} \succ 0$ . What if  $\lambda = 0$ ?

$$\hat{\alpha} = \arg \min_{\alpha} ||\mathbf{y} - \mathbf{K}\alpha||_2^2 + \lambda \alpha^T \mathbf{K} \alpha$$

# Why regularization?

Typically,  $\mathbf{K} \succ 0$ . What if  $\lambda = 0$ ?

$$\hat{\alpha} = \arg \min_{\alpha} ||\mathbf{y} - \mathbf{K}\alpha||_2^2 + \lambda \alpha^T \mathbf{K} \alpha$$

Unregularized kernel least squares can (over) fit **any data!**

$$\hat{\alpha} = \mathbf{K}^{-1} \mathbf{y}$$

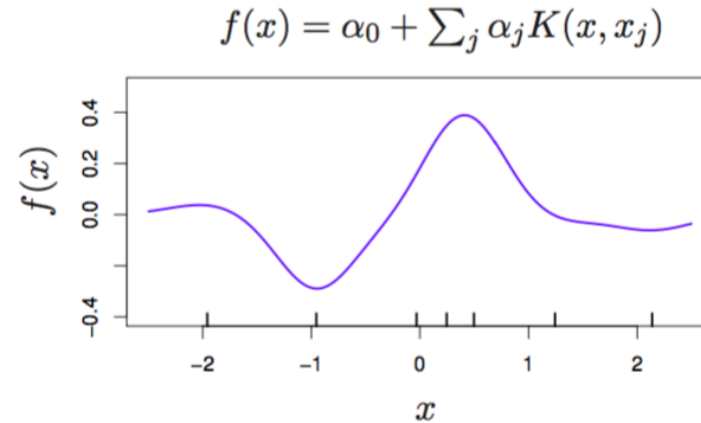
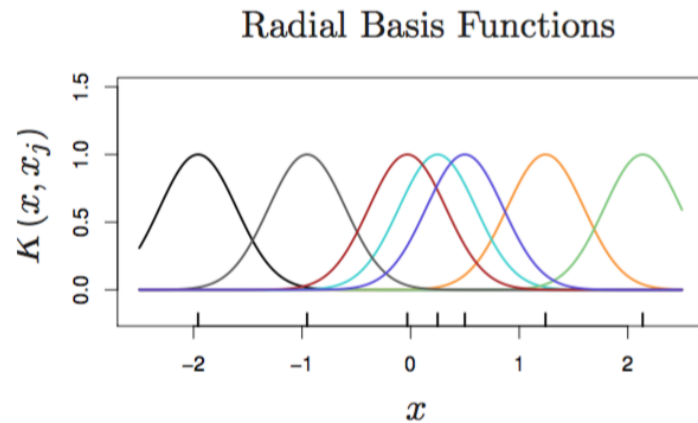
# The Kernel Trick for SVMs

---

# RBF Kernel

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

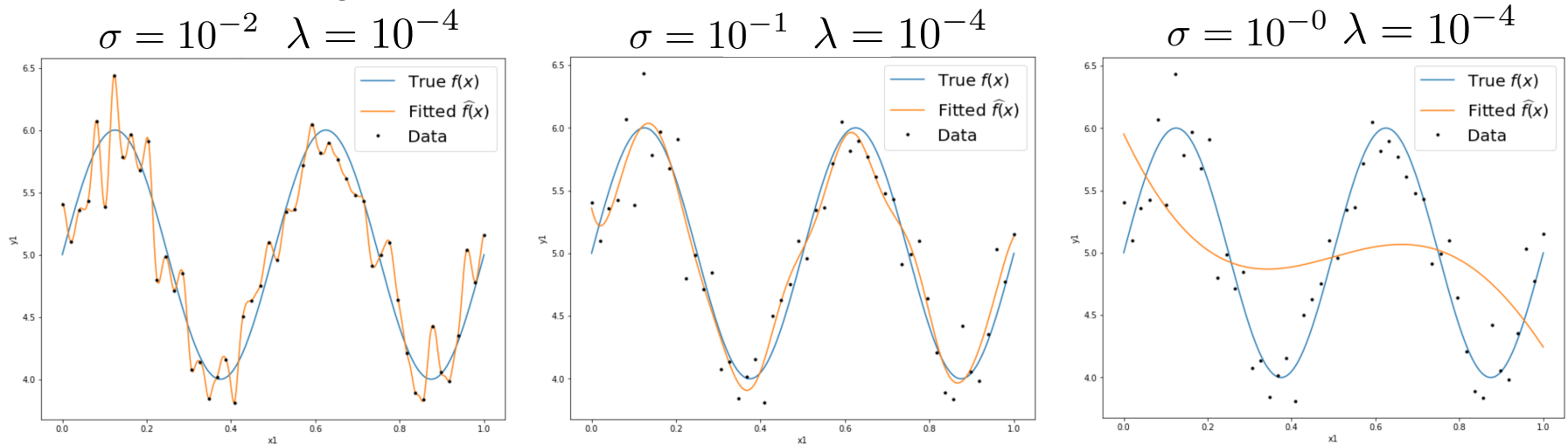
This is like weighting “bumps” on each point



# RBF Kernel

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

The bandwidth sigma has an enormous effect on fit:

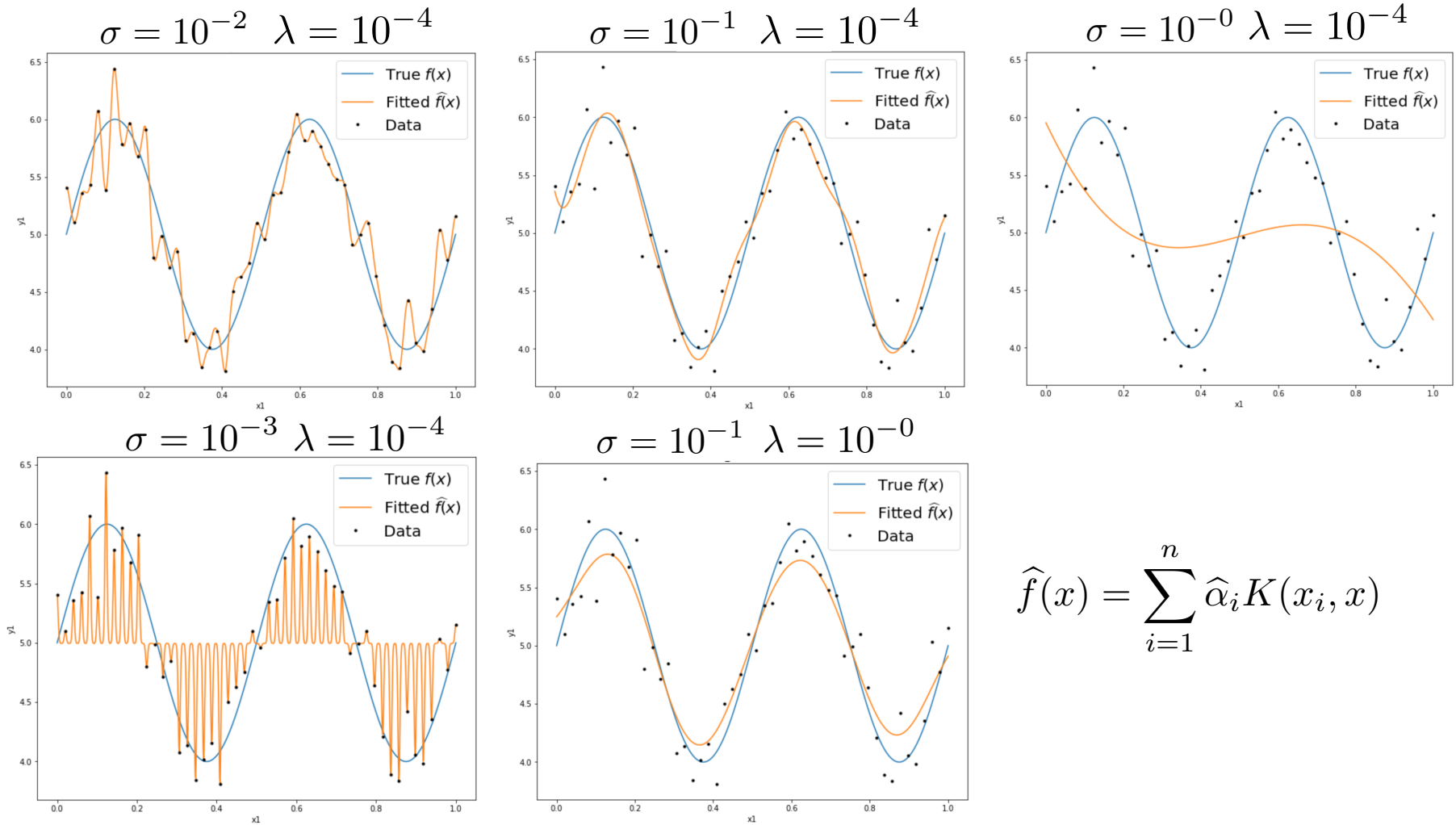


$$\hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i K(x_i, x)$$

# RBF Kernel

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

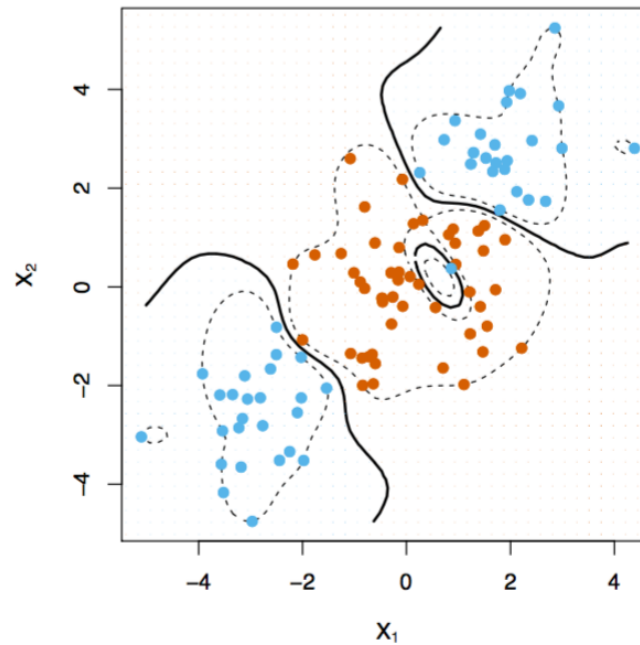
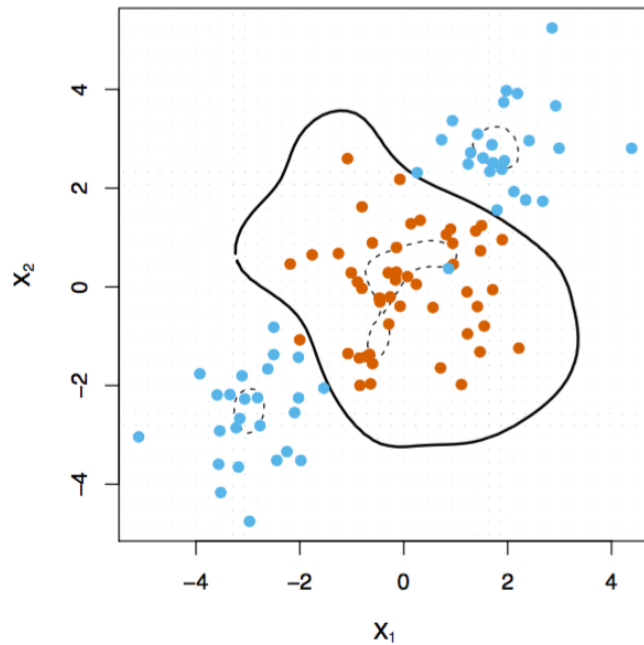
The bandwidth sigma has an enormous effect on fit:



$$\hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i K(x_i, x)$$

# RBF kernel and random features

$$\hat{w} = \sum_{i=1}^n \max\{0, 1 - y_i(b + x_i^T w)\} + \lambda \|w\|_2^2$$
$$\min_{\alpha, b} \sum_{i=1}^n \max\{0, 1 - y_i(b + \sum_{j=1}^n \alpha_j \langle x_i, x_j \rangle)\} + \lambda \sum_{i,j=1}^n \alpha_i \alpha_j \langle x_i, x_j \rangle$$



# RBF kernel and random features

---

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

If  $n$  is very large, allocating an  $n$ -by- $n$  matrix is tough.



# RBF kernel and random features

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

If  $n$  is very large, allocating an  $n$ -by- $n$  matrix is tough.

$$2 \cos(\alpha) \cos(\beta) = \cos(\alpha + \beta) + \cos(\alpha - \beta)$$

$$e^{jz} = \cos(z) + j \sin(z)$$

$$\phi(x) = \begin{bmatrix} \sqrt{2} \cos(w_1^T x + b_1) \\ \vdots \\ \sqrt{2} \cos(w_p^T x + b_p) \end{bmatrix}$$

$$w_k \sim \mathcal{N}(0, 2\gamma I)$$

$$b_k \sim \text{uniform}(0, \pi)$$

[Rahimi, Recht NIPS 2007]  
“NIPS Test of Time Award, 2018”

# String Kernels

Example from Efron and Hastie, 2016

Amino acid sequences of different lengths:

x1

IPTSALVKETLALLSTHRTLLIANETLRIPVPVHKNHQLCTEEIFQGIGTLESQTVQGGTV  
ERLFKNLSLIKKYIDGQKKKCGEERRRVNQFLDY**LQE**FLGVMNTEWI

x2

PHRRDLCSRSIWLARKIRSDLTALTESYVKHQGLWSELTEAER**LQEN**LQAYRTFHVLLA  
RLLEDQQVHFTPTGDFHQAIHTLLLQVAAFAYQIEELMILLEYKIPRNEADGMLFEKK  
LWGLKV**LQE**LSQWTVRSIHDLRFISSHQTGIP

All subsequences of length 3 (of possible 20 amino acids)  $20^3 = 8,000$

$$h_{\text{LQE}}^3(x_1) = 1 \text{ and } h_{\text{LQE}}^3(x_2) = 2.$$

# Fixed Feature V.S. Learned Feature

---