

# MTH 221

Fundamentals of Machine Learning

Batuhan Bardak

**Lecture 7:** Support Vector Machines

**Date:** 14.11.2023

# Announcements

- Midterm date
  - December 1, 2023 from 18:40 to 20:40
  - Informatics Institute, II-02 and II-03
  - You are responsible for the topics we have covered and discussed in class until the exam (including the content shared by the guests in the Special Guest lecture and this paper as well <https://homes.cs.washington.edu/~pedrod/papers/cacm12.pdf>)

# Announcements

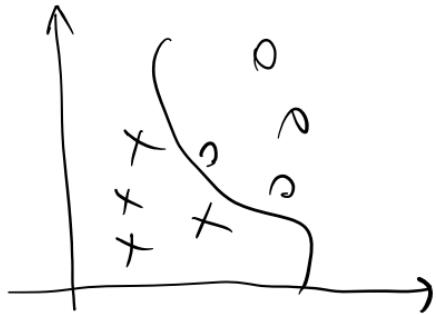
- Special Guest Lecture, On November 21, 2023
  - Barkın Saritaş, Senior Sales Strategy, Planning & Operations Manager @ Github
    - Transforming Data into Decisions: The Anatomy of a Corporate Analytics Project
  - Yunuscan Koçak, Machine Learning Scientist @ Booking
    - From Jupyter to Production: A Pragmatic Guide to Deploying Machine Learning Models
  - Salim Tütüncü, Senior Solutions Architect @ AWS
    - Generative AI: How to Leverage AWS ML Services to Drive Value

# Announcements

- Project Progress Report Deadline
  - December 3, 2023 by 23:59
  - The outline should be
    - EDA
    - Experimental Reesults
    - Future Work
    - At least 3 pages long in PDF Format
    - Send to TA while keeping me in the CC.

# Recap Logistic Regression

\* classification.

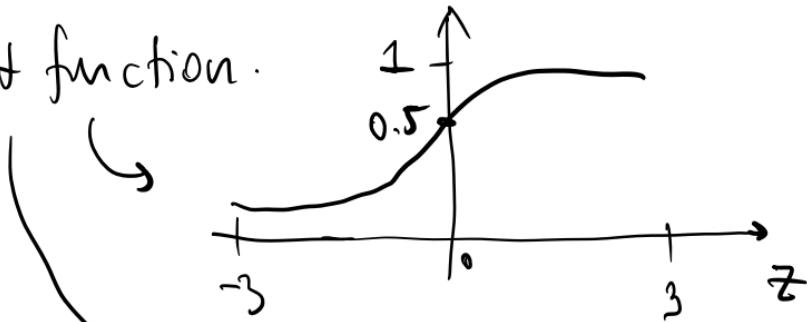


$$f_{\vec{w}, b}(\vec{x})$$

$$\vec{z} = \vec{w}\vec{x} + b$$

$$g(z) = \frac{1}{1+e^{-z}}$$

sigmoid function.



$$g(z) = \frac{1}{1+e^{-z}}$$

# Recap Gradient Descent

↗ linear.  
 ↗ new.  
 ↗ logistic.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \left( f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)} \right)^2 \rightarrow \text{linear reg.}$$

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) \right]$$



$\min_w J(w, b)$

Repeat until converge

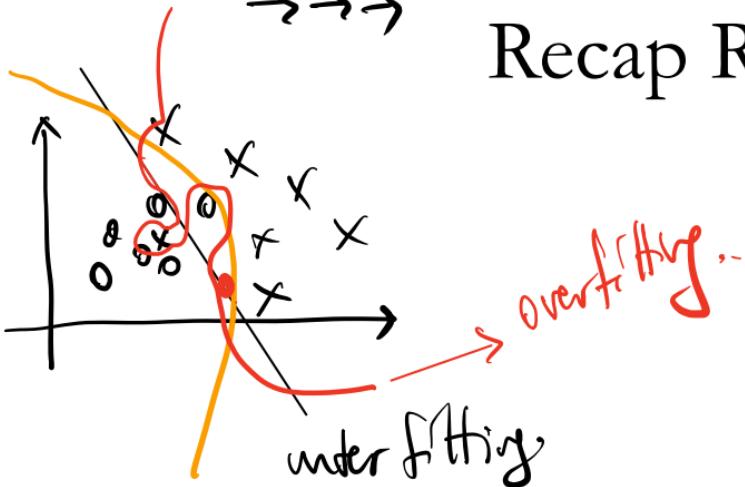
$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$w_J \rightarrow \rightarrow \rightarrow$$

}

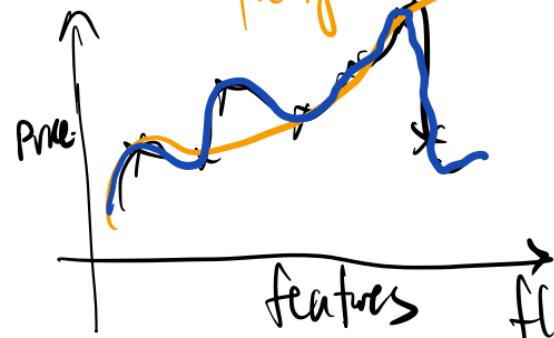
↓ learning rate

## Recap Regularization



- Collect more training data.
- Select features ~ 1 - 100  
feature selection algorithms.

- Regularization.  
→ Reduce the size of parameters  $w_j$ .



$$f(x) = 28x - 385x^2 + 174x^4 + 100$$

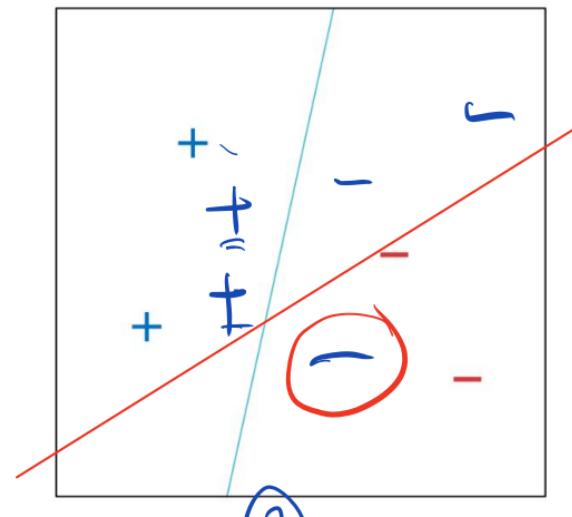
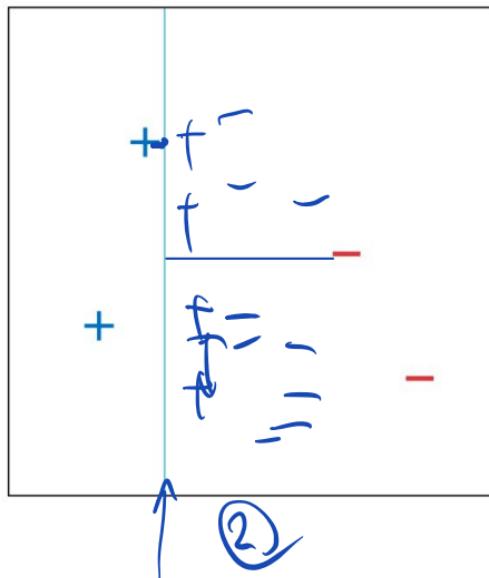
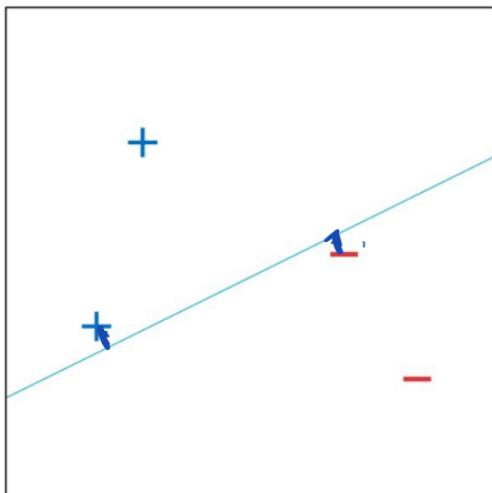
0

# Plan for today

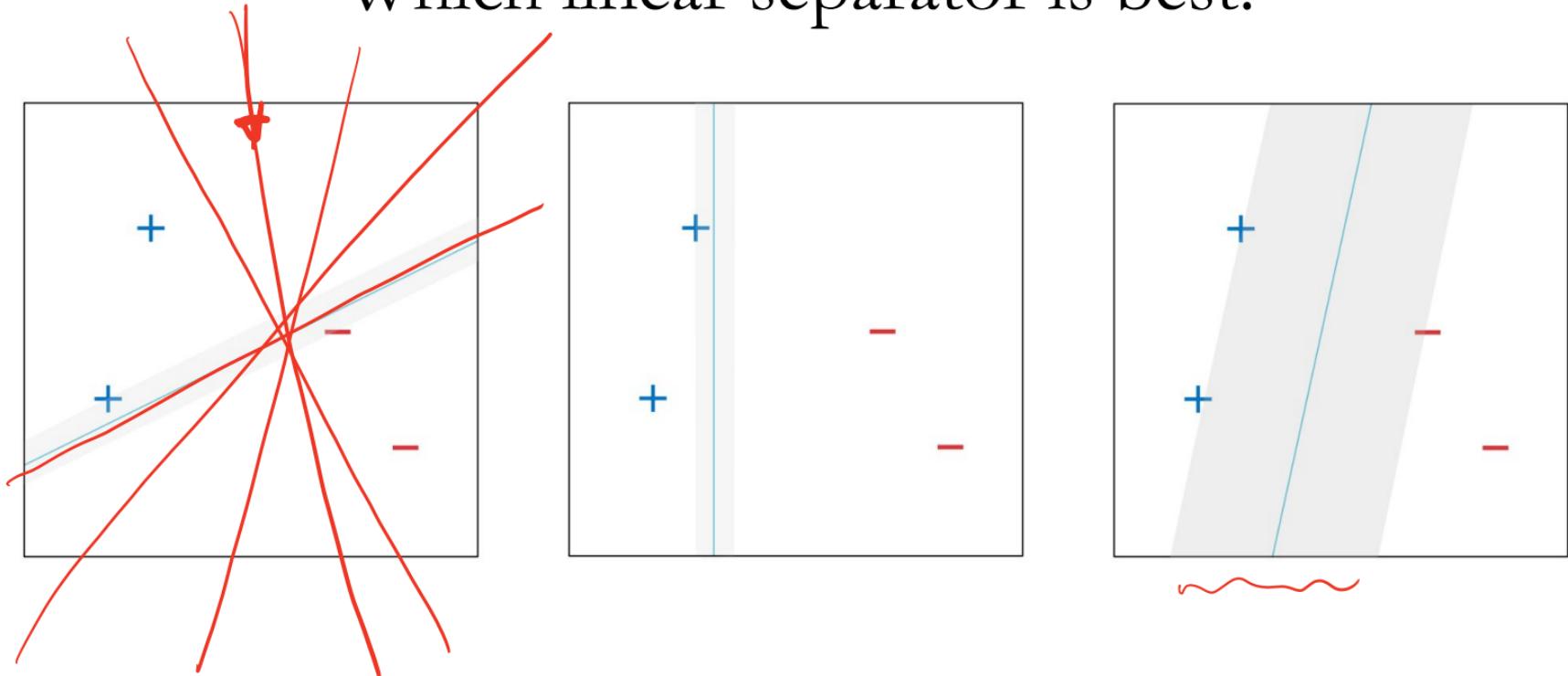
- SVM for linearly separable classes
- SVM for linearly inseparable classes
- SVM for nonlinear decision boundaries
  - Kernel functions

A hand-drawn diagram illustrating the SVM optimization problem. It shows a blue line representing the decision boundary, with a vector  $w$  originating from it. A point  $x$  is shown above the line, with its distance to the line labeled  $\gamma(w, b)$ . The formula for  $\gamma(w, b)$  is given as  $\gamma(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x_i) - y_i)^2$ . The diagram is divided into two regions by the line: Region I (labeled "I fit data") contains points  $x_1, x_2, \dots, x_m$ , while Region II (labeled "II regularizer pen") contains points  $x_{m+1}, x_{m+2}, \dots, x_n$ . Below the line, there are two arrows: one pointing up labeled  $\lambda \uparrow \rightarrow \text{underfit.} > 0$  and one pointing down labeled  $\lambda \downarrow \rightarrow \text{overfit}$ .

# Which linear separator is best?



# Which linear separator is best?



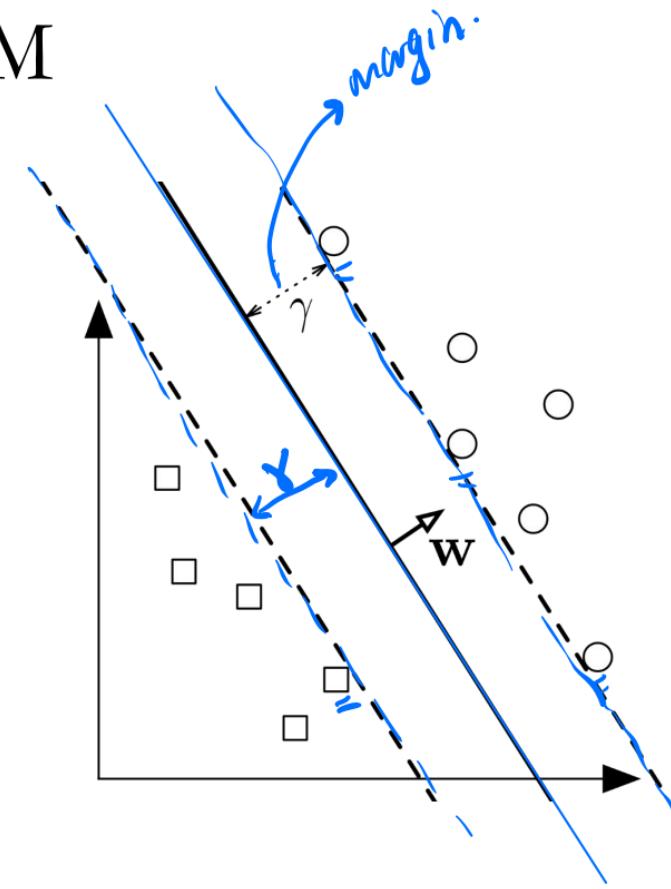
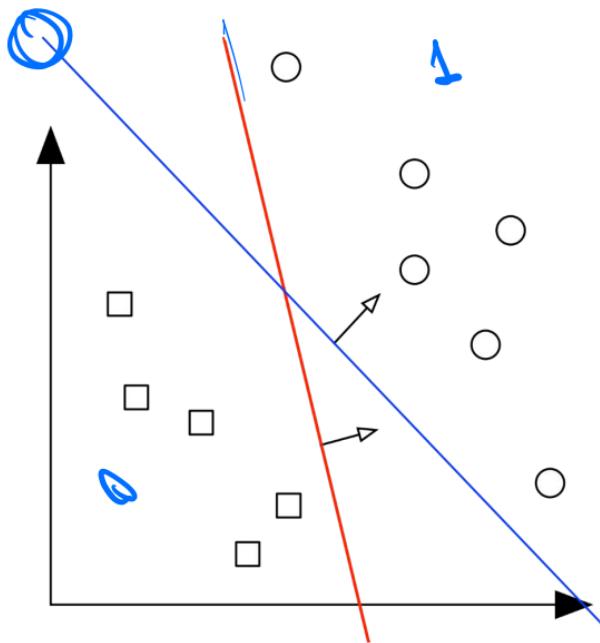
# SVM

- The Support Vector Machine (SVM) is a linear classifier that can be viewed as an extension of the Perceptron.
- The Perceptron guaranteed that you find a hyperplane if it exists. The SVM finds the **maximum margin** separating hyperplane.

# Maximal Margin Linear Separators

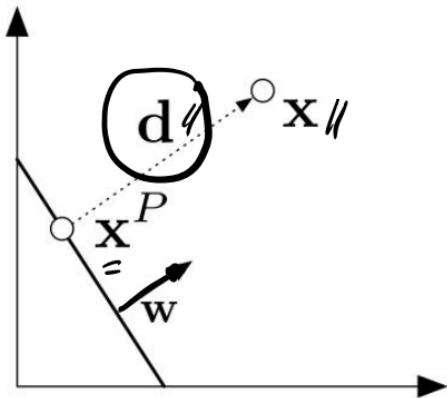
- The margin of a linear separator is the distance between it and the nearest training data point.
- Questions:
  - How can we efficiently find a maximal-margin linear separator?
  - Why are linear separators with larger margins better?
  - What can we do if the data is not linearly separable?

# SVM



# Geometric Margin

Hyperplane defined by  $(w, b)$ , i.e.,  
 $\{x : w^T x + b = 0\}$



Fact 1.  $x - x^P$  is parallel to  $w$ :

$$x - x^P = \alpha w$$

Fact 2.  $x^P$  is on the hyperplane:

$$w^T x^P + b = 0$$

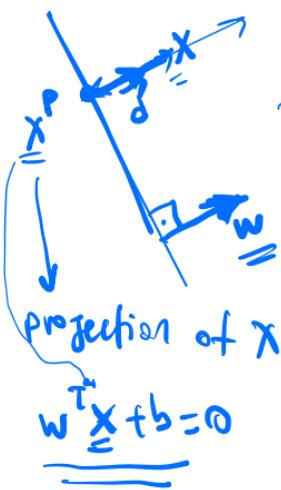
Fact 1 + fact 2 implies:

$$w^T(x - \alpha w) + b = 0 \rightarrow \alpha = (w^T x + b) / \|w\|_2^2$$

Final step:

$$d = \|x - x^P\|_2 = \|\alpha w\|_2 = \frac{|w^T x + b|}{\|w\|_2}$$

$H_{w,b} = \{x : w^T x + b = 0\}$  find the distance b/w point "x" and hyperplane.



## Geometric Margin

\* "f" must be the rescaled version of

$$\begin{aligned} w^T x_p + b &= 0 \\ &= w^T (\vec{x} - \vec{d}) + b \\ &= w^T (\vec{x} - \alpha \vec{w}) + b = 0 \end{aligned}$$

$$\begin{aligned} \vec{d} &= \alpha \vec{w} \\ \alpha &= \frac{w^T x + b}{w^T w} \end{aligned}$$

- length  $\rightarrow$  norm of  $d$

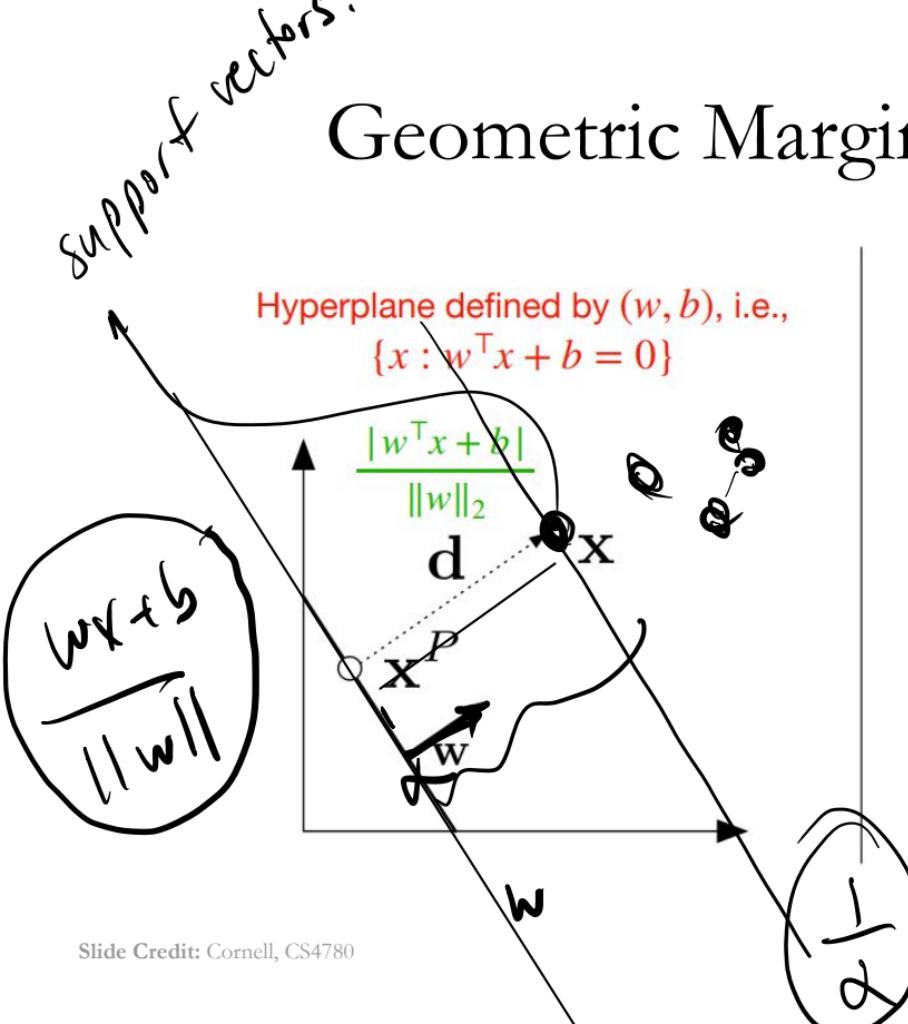
distance of a point  $x$  to the hyperplane

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2}$$

$$\begin{aligned} \|d\|_2 &= \sqrt{d^T d} \\ &= \sqrt{\alpha^2 w^T w} \\ &= \alpha \sqrt{w^T w} \end{aligned}$$

$$\begin{aligned} d &= \frac{w^T x + b}{\sqrt{w^T w}} \\ &= \frac{w^T x + b}{\|w\|_2} \end{aligned}$$

# Geometric Margin is Scale Invariant



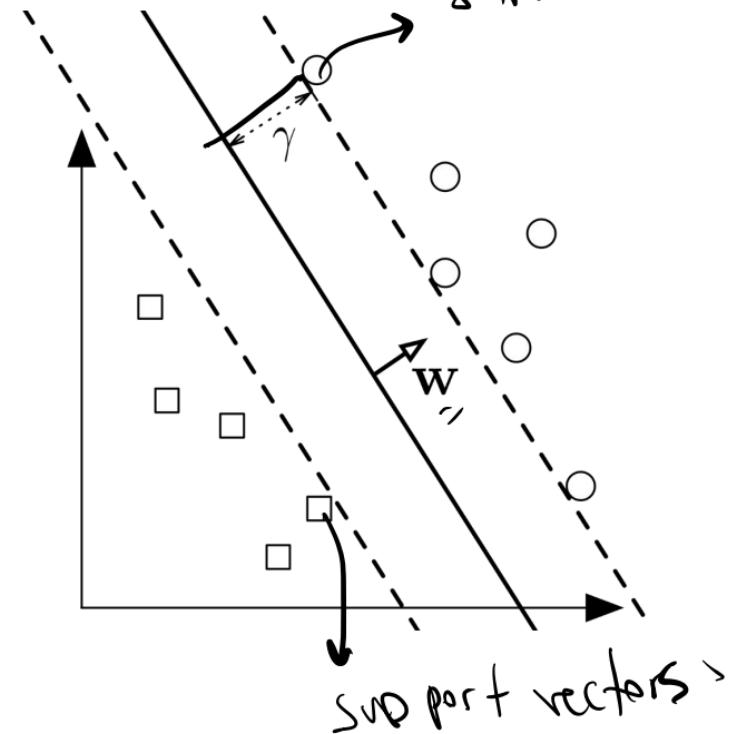
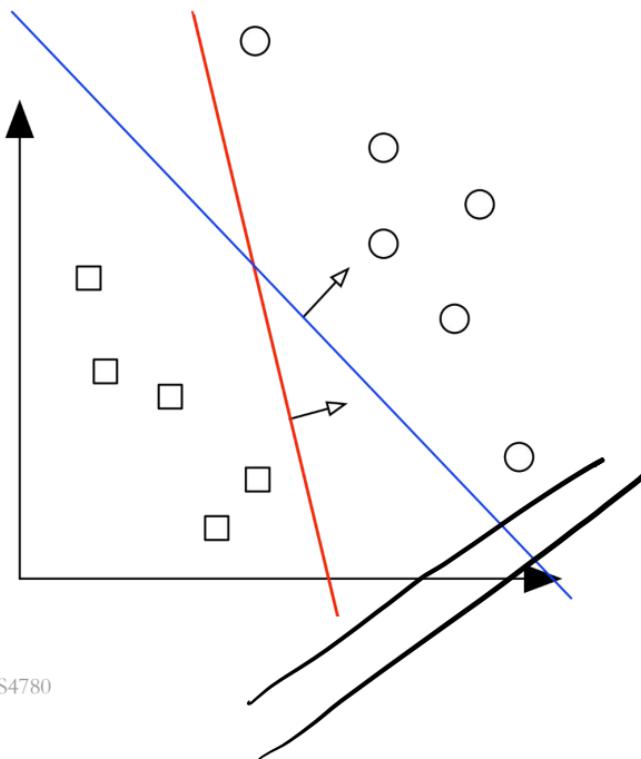
We scale  $(w, b)$  by a constant  $\gamma \in \mathbb{R}^+$

Q: is the hyperplane defined by  
 $(\gamma w, \gamma b)$  different?

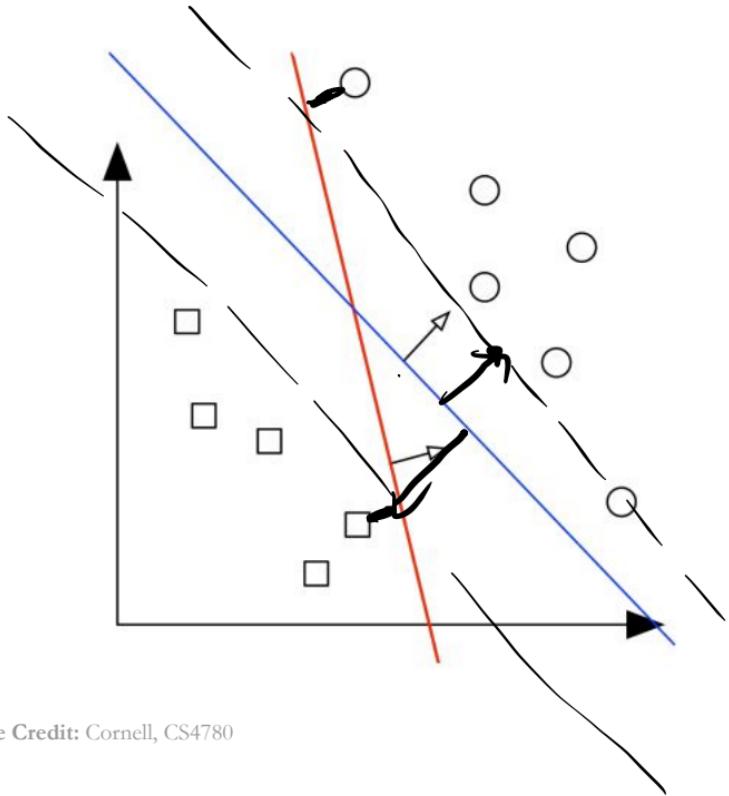
Q: does the margin change?

Hyperplane & Geometric margin are  
scale invariant!

# Which linear classifier is better?



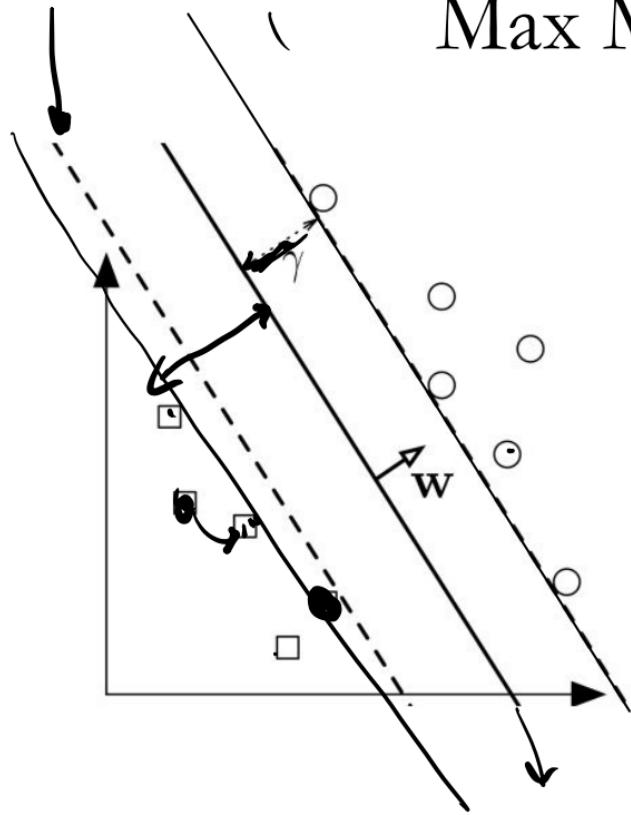
# Max Margin Classifier



The Goal of SVM:

Find a hyperplane that has the largest  
Geometric margin

# Max Margin Classifier



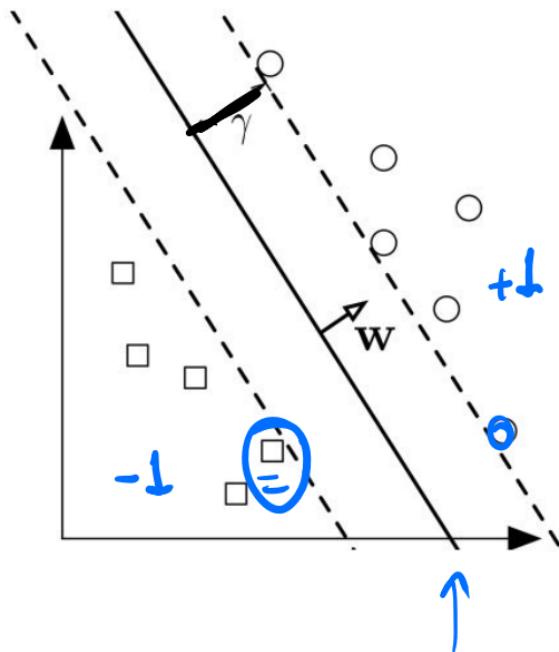
Given a linearly separable dataset  $\{x_i, y_i\}_{i=1}^n$ , the minimum geometric margin is defined as

$$\gamma(w, b) := \min_{x_i \in \mathcal{D}} \frac{|x_i^T w + b|}{\|w\|_2}$$

Goal: we want to find  $(w, b)$  s.t. it separates the data, and maximizes  $\gamma(w, b)$

$$\max_{w, b} \gamma(w, b)$$

# Max Margin Classifier



We want to find  $(w, b)$  s.t. it separates the data, and maximizes  $\gamma(w, b)$

Constraint

$$\max_{w,b} \underline{\gamma(w, b)}$$

s.t.  $\forall i, y_i(w^\top x_i + b) \geq 0$

Plug in the def of  $\gamma(w, b)$ :

$$\max_{w,b} \frac{1}{\|w\|_2} \min_{x_i} |w^\top x_i + b|$$

$$\text{s.t. } \forall i, y_i(w^\top x_i + b) \geq 0$$

# SVM for separable data: Max Margin Classifier

We want to find  $(w, b)$  s.t. it separates the data, and maximize  $\gamma(w, b)$

$$\min_{w, b} w^T w$$

$$\text{s.t. } \forall i, y_i(w^T x_i + b) \geq 0$$

$$\min_i |w^T x_i + b| = 1$$

Recall that margin & hyperplane is scale invariant

For any  $(w, b)$ , we can always scale it by some constant to have

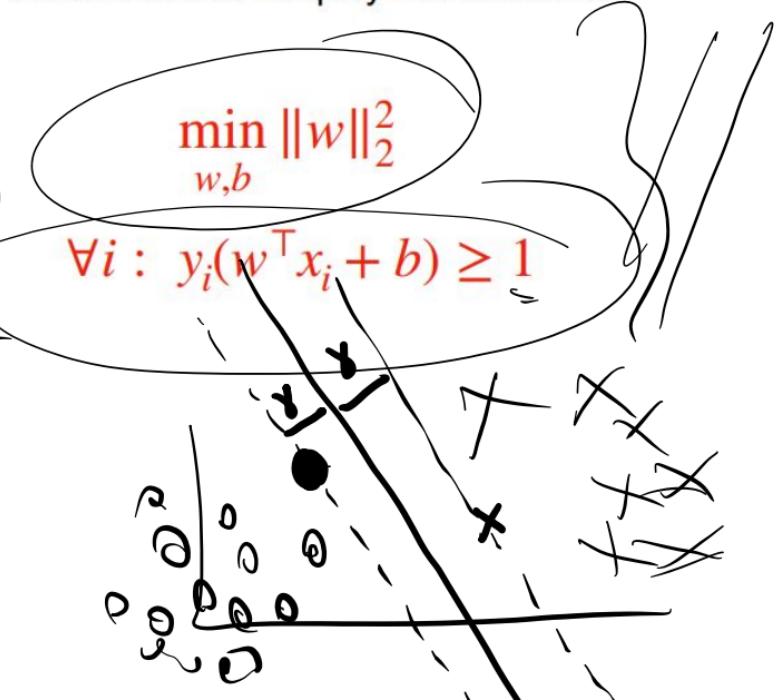
$$\min_{x_i} |w^T x_i + b| = 1$$

Without loss of generality, let's just focus on such  $(w, b)$  pairs with  $\min_{x_i} |w^T x_i + b| = 1$

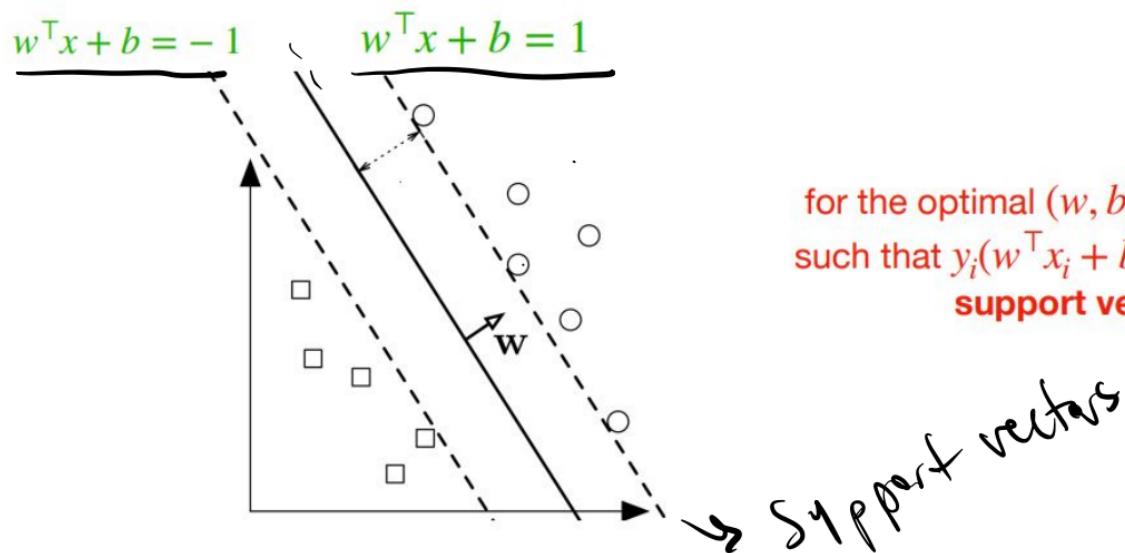
# SVM for separable data: Max Margin Classifier

$$\begin{aligned} & \min_{w,b} \|w\|_2^2 \\ \text{s.t. } & \forall i, y_i(w^\top x_i + b) \geq 0 \\ & \underbrace{\min_i |w^\top x_i + b| = 1} \end{aligned}$$

We can further simplify the constraint

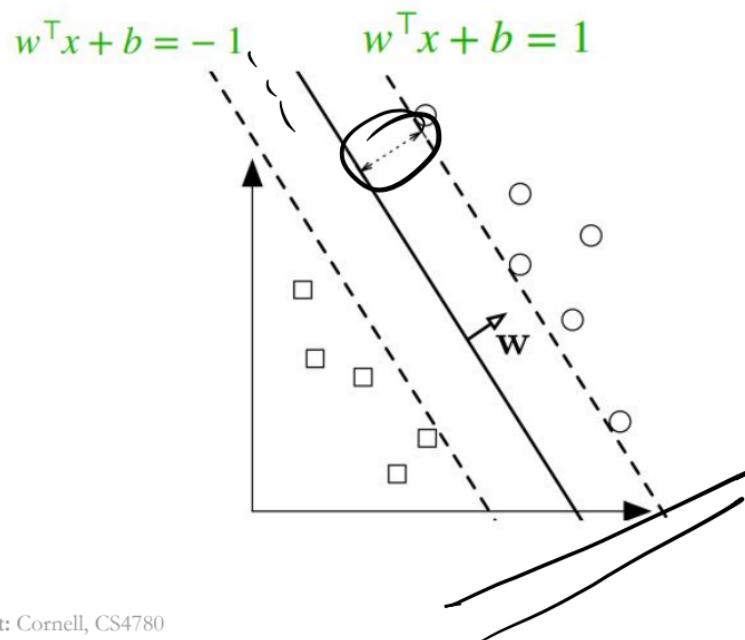


# Max Margin Classifier



for the optimal  $(w, b)$  pair, points  $x_i$  such that  $y_i(w^T x_i + b) = 1$  are called support vectors

# Support Vectors



for the optimal  $(w, b)$  pair, points  $x_i$  such that  $y_i(w^\top x_i + b) = 1$  are called **support vectors**

# SVM for non-separable data

If data is not linearly separable, then **there is no**  $(w, b)$   
can satisfy  $\forall i : y_i(w^\top x_i + b) \geq 1$



= **infeasible.**

Idea: introducing slack variables to relax the constraint, i.e., find  $(w, b, \xi_i)$ , s.t,

$$\begin{aligned} \forall i : y_i(w^\top x_i + b) &\geq 1 - \xi_i, & \text{soft constraints} \\ \xi_i &\geq 0, \forall i \end{aligned}$$

Q: does this always has a feasible solution?

# SVM for non-separable data

Idea: introducing slack variables to relax the constraint, i.e., find  $(w, b, \xi_i)$ , st,

$$\forall i : y_i(w^\top x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad \xi_i(x_i)$$

We still want our margin to be somewhat large, i.e., we want slack variables to be as small as possible

$$\min_{w,b,\xi} \|w\|_2^2 + c \sum_{i=1}^n \xi_i \quad \text{Penalizing large slacks}$$

$$\forall i : y_i(w^\top x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

$$w^\top x_i \geq 1$$

# SVM for non-separable data

$$\min_{w,b,\xi_i} \|w\|_2^2 + c \sum_{i=1}^n \xi_i$$

Penalizing large slacks

---

$$\forall i : y_i(w^\top x_i + b) \geq 1 - \xi_i, \xi_i \geq 0$$

We can turn this constrained opt to a unconstraint opt w/ a single objective.

Q: For any fixed  $(w, b)$  pair, how to set  $\xi_i$ , such that the obj is minimized?

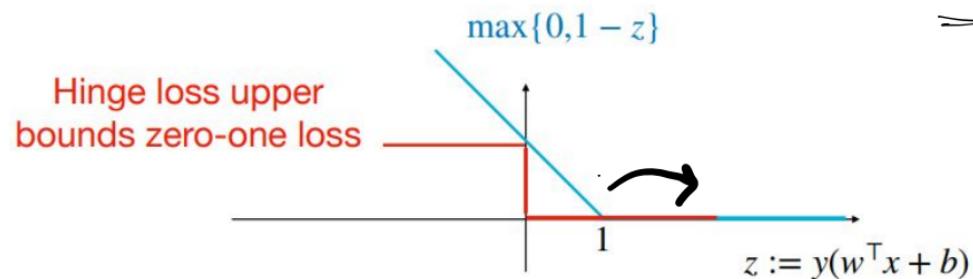
---

A: set  $\xi_i = \max\{0, 1 - y_i(w^\top x_i + b)\}$

# SVM for non-separable data

$$\min_{w,b} \|w\|_2^2 + c \sum_{i=1}^n \max \{0, 1 - y_i(w^\top x_i + b)\}$$

Hinge loss



Hinge loss starts penalizing when functional margin falls below 1

# SVM for non-separable data

$$\min_{w,b} \|w\|_2^2 + c \sum_{i=1}^n \max \{0, 1 - y_i(w^\top x_i + b)\}$$

Trades off  $\|w\|_2^2$  and functional margins over data

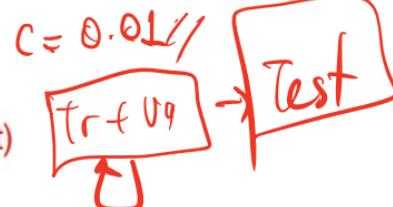
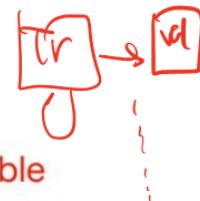
When  $c \rightarrow +\infty$ :  
 forcing  $y_i(w^\top x_i + b) \geq 1$  for as many data points as possible

When  $c \rightarrow 0^+$ :

The solution  $w \rightarrow \mathbf{0}$  (i.e., we do not care about hinge loss part)



$c = 0.05$



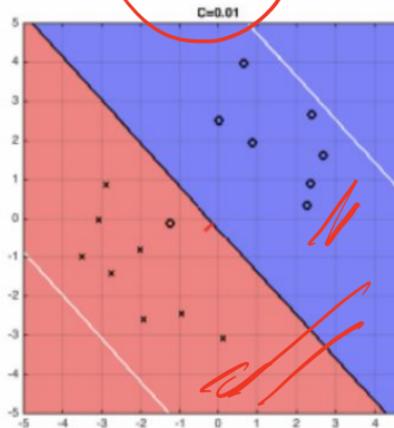
what if  $\gamma$

# SVM for non-separable data

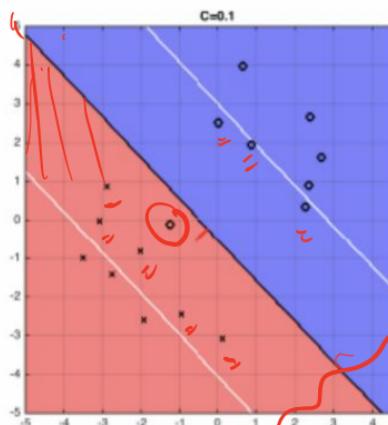
$$\min_{w,b} \|w\|_2^2 + c \sum_{i=1}^n \max \{0, 1 - y_i(w^\top x_i + b)\}$$

Trades off  $\|w\|_2^2$  and functional margins over data

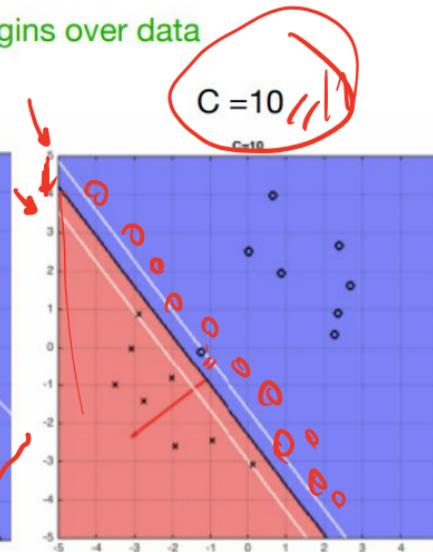
$C = 0.01$



$C = 1$



$C=10$



# Summary

## 1. SVM for linearly separable data

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

$$\forall i : y_i(w^\top x_i + b) \geq 1$$

## 2. SVM for non-separable data

$$\min_{w,b} \|w\|_2^2 + c \sum_{i=1}^n \max \{0, 1 - y_i(w^\top x_i + b)\}$$

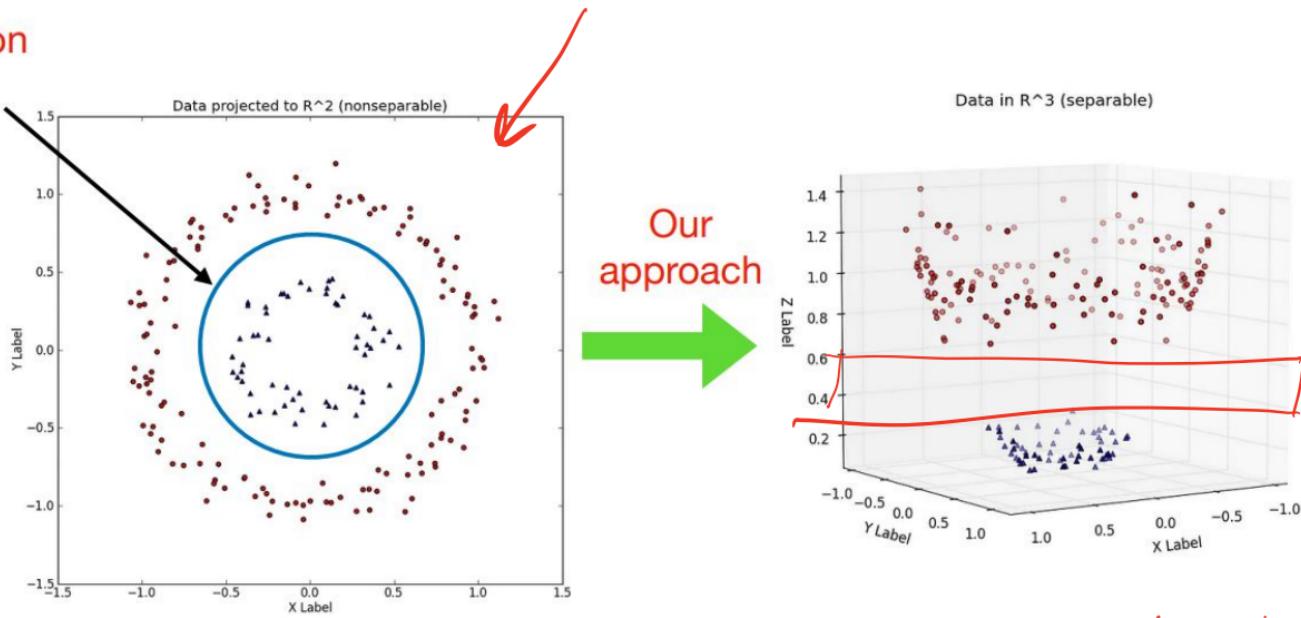
Hinge loss

# Kernel SVMs

- Go from linear models to more powerful nonlinear ones.
- Keep convexity (ease of optimization).
- Generalize the concept of feature engineering.

# Kernels

Goal: Non-linear decision boundary



Slide Credit: Cornell, CS4780

$\phi$

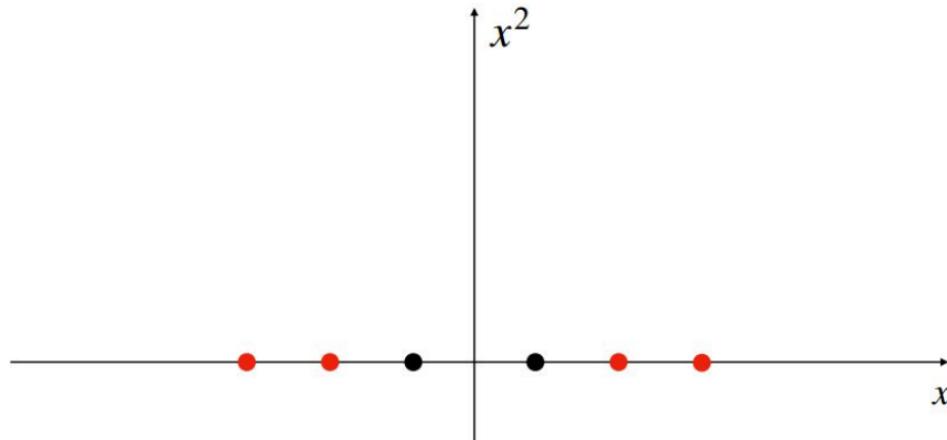
$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

Kernel trick

# Feature mapping

Define  $\phi(\mathbf{x}) \in \mathbb{R}^m$  as a feature mapping (often  $m > d$ )

Ex 1:  $\mathbf{x} \in \mathbb{R}$ ,  $\phi(\mathbf{x}) = [x, x^2]^\top \in \mathbb{R}^2$



# Feature mapping

Define  $\phi(\mathbf{x}) \in \mathbb{R}^m$  as a feature mapping (often  $m > d$ )

Ex 2: quadratic  
feature mapping  $\phi$

$$\boxed{\mathbf{x} = [x_1, x_2]^\top,}$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]^\top$$

= \_\_\_\_\_

# Feature mapping

Define  $\phi(\mathbf{x}) \in \mathbb{R}^m$  as a feature mapping (often  $m > d$ )

Ex 2: cubic feature  
mapping  $\phi$

$$\mathbf{x} = [x_1, x_2]^\top,$$

---

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_1^3, x_2^3, x_1 x_2^2, x_1^2 x_2]^\top$$

Q: in general, for  $\mathbf{x} \in \mathbb{R}^d$ , and a p-th order polynomial feature  $\phi$ , what's the dimension of  $\phi(\mathbf{x})$ ?

$$\text{at least } \binom{d}{p}$$

**Dim of  $\phi(\mathbf{x})$  can be very large!**

# Feature mapping

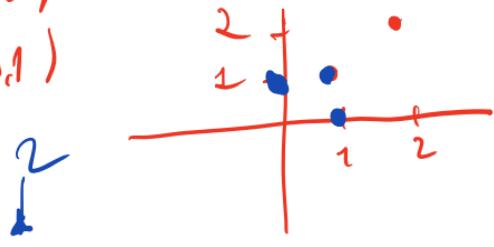
Define  $\phi(\mathbf{x}) \in \mathbb{R}^m$  as a feature mapping (often  $m > d$ )

Ex 2: cubic feature  
mapping  $\phi$

$$\mathbf{x} = [x_1, x_2]^\top,$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_1^3, x_2^3, x_1 x_2^2, x_1^2 x_2]^\top$$

class = 1 :  $(1, 1)$   $(2, 2)$   
class = 2 :  $(1, 0)$   $(0, 1)$

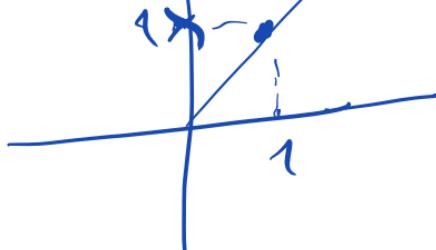


Q: in general, for  $\mathbf{x} \in \mathbb{R}^d$ , and a p-th order polynomial feature  $\phi$ , what's the dimension of  $\phi(\mathbf{x})$ ?

at least  $\binom{d}{p}$

Dim of  $\phi(\mathbf{x})$  can be very large!  
 $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x} \cdot \mathbf{x}')^2$

$$\mathbf{z} = (z_1, z_2, z_3) = (x_1, \sqrt{2}x_1, x_2)$$



## Example of Kernels

Point (1,1) becomes:

$$(1^2, \sqrt{1+1}, 1, 1^2) = (1, \sqrt{2}, 1)$$

$$k_{\text{linear}}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

$$k_{\text{poly}}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d$$

$$k_{\text{rbf}}(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

(1,0)

$$(1^2, \sqrt{2}, 1, 0, 0^2) = (1, 0, 0)$$

$$k_{\text{sigmoid}}(\mathbf{x}, \mathbf{x}') = \tanh(\gamma \mathbf{x}^T \mathbf{x}' + r)$$

$$k_{\cap}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^p \min(x_i, x'_i)$$

- If  $k$  and  $k'$  are kernels, so are  $k + k'$ ,  $kk'$ ,  $ck'$ , ...

# Polynomial Kernel vs Features

$$k_{\text{poly}}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d$$

Primal vs Dual Optimization

Explicit polynomials → compute on n\_samples \* n\_features \*\* d

Kernel trick → compute on kernel matrix of shape n\_samples \* n\_samples

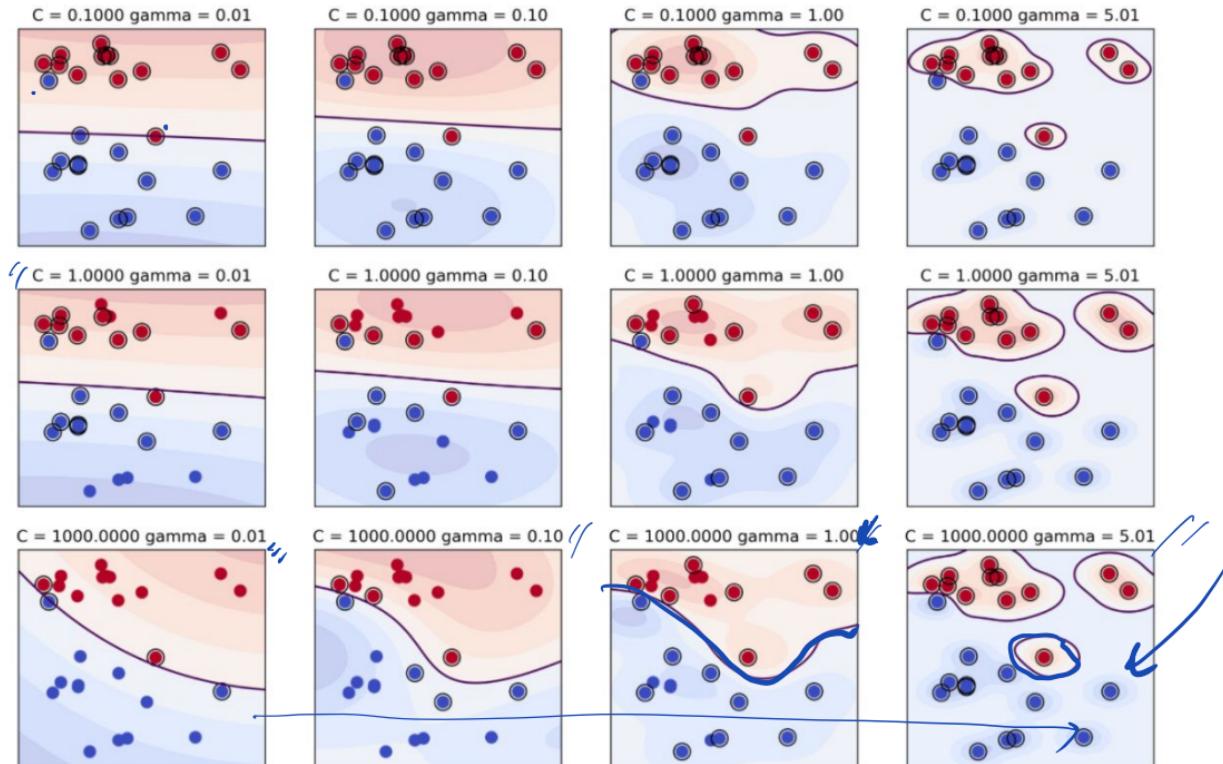
For a single feature:

$$(x^2, \sqrt{2}x, 1)^T (x'^2, \sqrt{2}x', 1) = x^2 x'^2 + 2x x' + 1 = (x x' + 1)^2$$

# Kernels in Practice

- Dual coefficients less interpretable
- Long runtime for “large” datasets (100k samples)
- Real power in infinite-dimensional spaces: rbf!
- Rbf is “universal kernel” - can learn (aka overfit) anything.

# Parameters for RBF Kernels



## Next Class:

Decision Trees & Ensemble Learning

Novel

Not -

- <sup>last</sup> lecture
- Review / hands on example
- mistakes
-