

# CSCI 567: Machine Learning

Vatsal Sharan  
Spring 2026

Lecture 6, February 20

# Administrivia

- HW2 out, due in less than 1 week.
- Exam 1 in 2 weeks, more details on Ed

Recap

# Regularized least squares

We looked at regularized least squares with non-linear basis:

$$\begin{aligned}\mathbf{w}^* &= \underset{\mathbf{w}}{\operatorname{argmin}} F(\mathbf{w}) \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} (\|\Phi\mathbf{w} - \mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_2^2) \\ &= (\Phi^T\Phi + \lambda\mathbf{I})^{-1} \Phi^T\mathbf{y}\end{aligned}$$

$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_n)^T \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

This solution operates in the space  $\mathbb{R}^M$  and  $M$  could be huge (and even infinite).

# Regularized least squares solution: Another look

We realized that we can write,

$$\mathbf{w}^* = \Phi^T \boldsymbol{\alpha} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)$$

Thus the least square solution is **a linear combination of features of the datapoints!**

We calculated what  $\boldsymbol{\alpha}$  should be,

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

where  $\mathbf{K} = \Phi \Phi^T \in \mathbb{R}^{n \times n}$  is the **kernel matrix**.

# Kernel trick

The prediction of  $w^*$  on a new example  $x$  is

$$w^{*\top} \phi(x) = \sum_{i=1}^n \alpha_i \phi(x_i)^\top \phi(x)$$

Therefore, *only inner products in the new feature space matter!*

Kernel methods are exactly about computing inner products *without explicitly computing  $\phi$* . The exact form of  $\phi$  is inessential; *all we need to do is know the inner products  $\phi(x)^\top \phi(x')$* .

# The kernel trick: Example 1

Consider the following polynomial basis  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ :

$$\phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

What is the inner product between  $\phi(\mathbf{x})$  and  $\phi(\mathbf{x}')$ ?

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 = (\mathbf{x}^T \mathbf{x}')^2 \end{aligned}$$

Therefore, *the inner product in the new space is simply a function of the inner product in the original space.*

# Kernel functions

**Definition:** a function  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is called a *kernel function* if there exists a function  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^M$  so that for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ ,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

## Popular kernels:

1. Polynomial kernel

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^M$$

for  $c \geq 0$  and  $M$  is a positive integer.

2. Gaussian kernel or Radial basis function (RBF) kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}\right) \quad \text{for some } \sigma > 0.$$

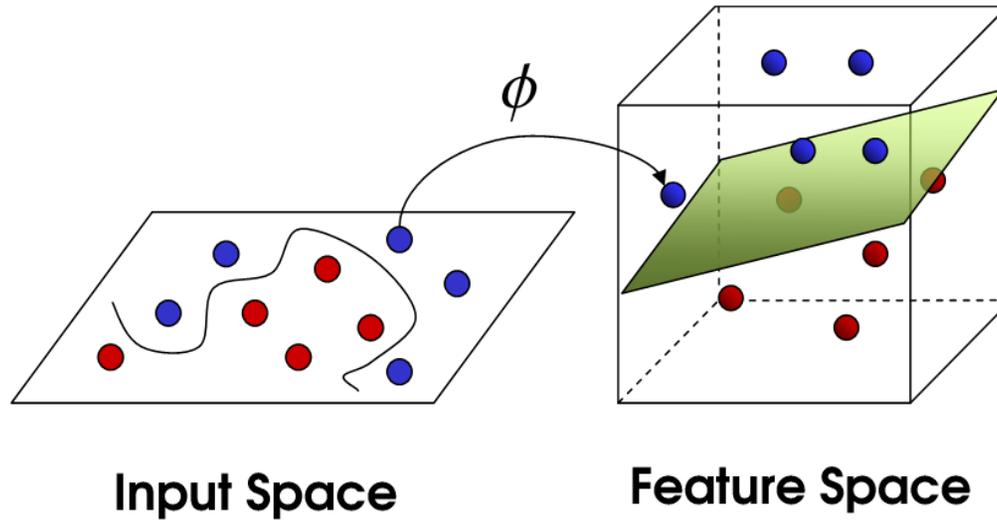
# Prediction with kernels

As long as  $\mathbf{w}^* = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)$ , prediction on a new example  $\mathbf{x}$  becomes

$$\mathbf{w}^{*\top} \phi(\mathbf{x}) = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}).$$

This is known as a **non-parametric method**. Informally speaking, this means that there is no fixed set of parameters that the model is trying to learn (remember  $\mathbf{w}^*$  could be infinite). Nearest-neighbors is another non-parametric method we have seen.

# Classification with kernels



Similar ideas extend to the classification case, and we can predict using  $\text{sign}(\mathbf{w}^T \phi(\mathbf{x}))$ .  
Data may become linearly separable in the feature space!

We'll see this today.

A diagram illustrating a Support Vector Machine (SVM) for a two-class problem. It shows a set of data points in a 2D space, divided into two classes: blue circles and red circles. A solid gray line represents the decision boundary, which is the hyperplane that separates the two classes. Two dashed gray lines represent the margins, which are parallel to the decision boundary and define the regions of maximum separation from the support vectors. The support vectors are the data points that lie on the margins. The text "Support vector machines (SVMs)" is overlaid on the diagram in a large, bold, red font.

**Support vector  
machines (SVMs)**

# Why study SVM?

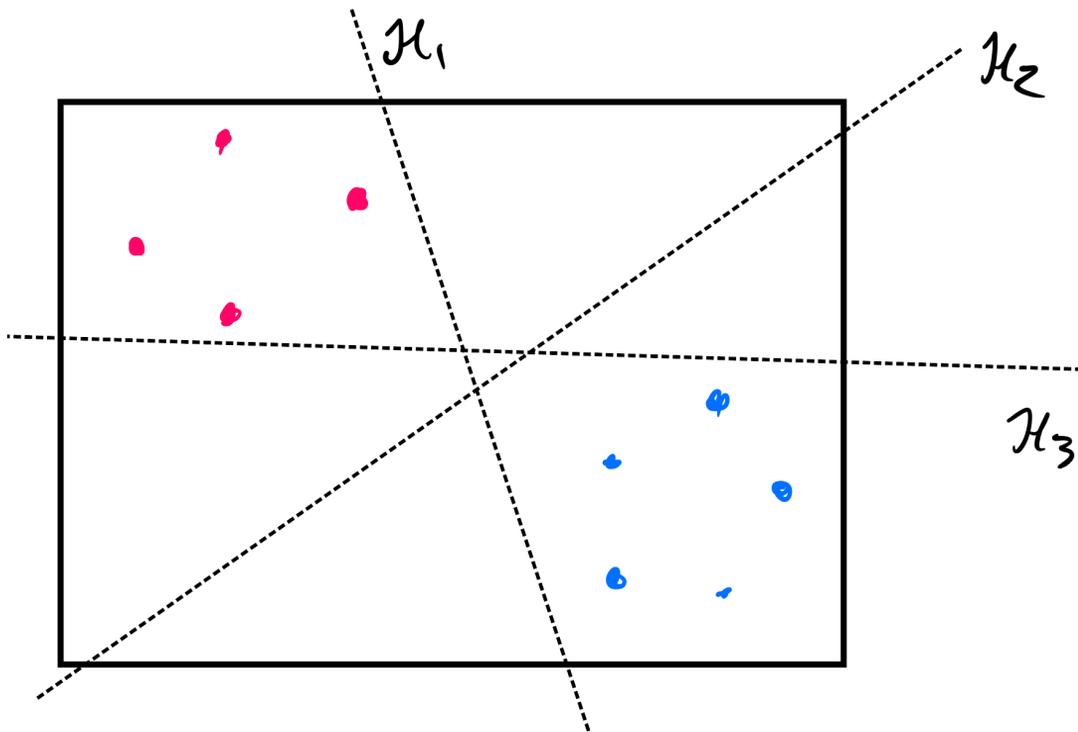
- One of the most commonly used classification algorithms
- Allows us to explore the concept of *margins* in classification
- Works well with the kernel trick
- Strong theoretical guarantees

We focus on **binary classification** here.

The *function class for SVMs is a linear function on a feature map  $\phi$  applied to the datapoints*:  $\text{sign}(\mathbf{w}^T \phi(\mathbf{x}) + b)$ . Note, the bias term  $b$  is taken separately for SVMs, you'll see why.

# Margins: separable case, geometric intuition

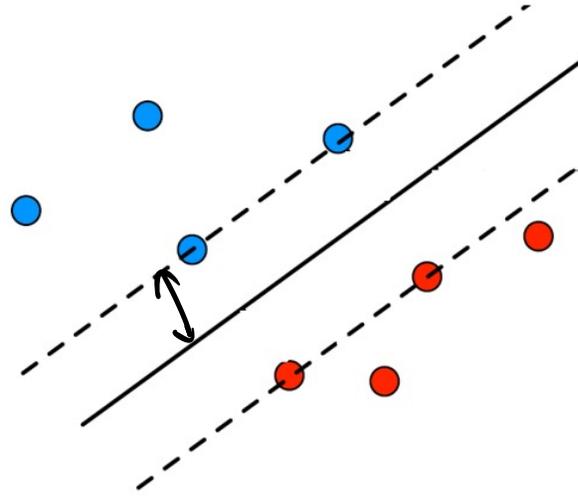
When data is **linearly separable**, there are infinitely many hyperplanes with zero training error:



Which one should we choose?

# Margins: separable case, geometric intuition

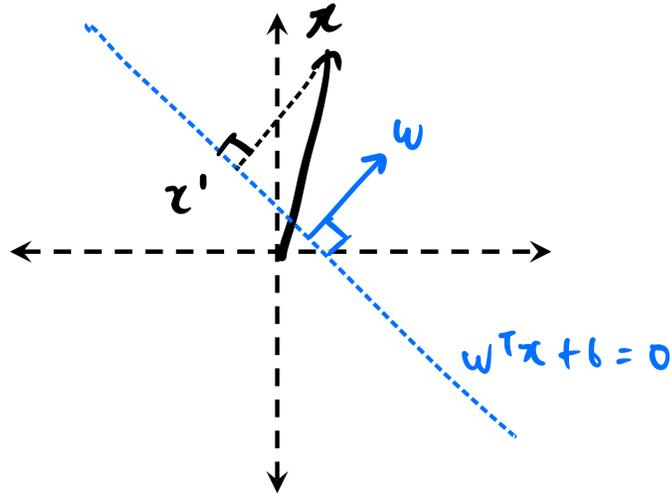
The further away the separating hyperplane is from the datapoints, the better.



**Margin for linearly separable data:** Distance from the hyperplane to the point closest to the hyperplane.

# Formalizing geometric intuition: Distance to hyperplane

What is the **distance** from a point  $x$  to a hyperplane  $\{x : w^T x + b = 0\}$ ?



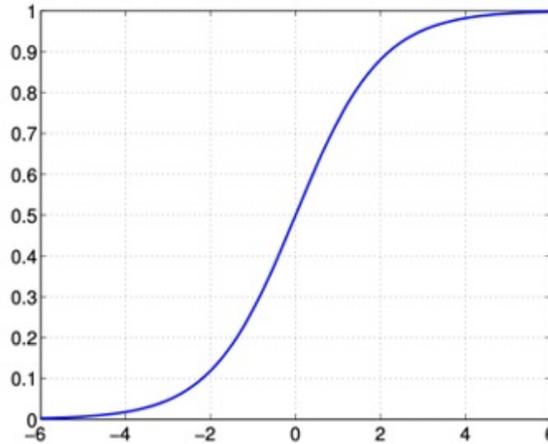
Assume the **projection** is  $x' = x - \beta \frac{w}{\|w\|_2}$ , then

$$0 = w^T \left( x - \beta \frac{w}{\|w\|_2} \right) + b = w^T x - \beta \|w\| + b \implies \beta = \frac{w^T x + b}{\|w\|_2}.$$

Therefore the distance is  $\|x - x'\|_2 = |\beta| = \frac{|w^T x + b|}{\|w\|_2}$ .

For a hyperplane that correctly classifies  $(x, y)$ , the distance becomes  $\frac{y(w^T x + b)}{\|w\|_2}$ .

# Margins: functional motivation



$$\Pr[y = 1 \mid \mathbf{x}; \mathbf{w}] = \sigma(y(\mathbf{w}^T \mathbf{x} + b)) = \frac{1}{1 + \exp(-y(\mathbf{w}^T \mathbf{x} + b))}$$

If  $y = 1$ , want  $\mathbf{w}^T \mathbf{x} + b \gg 0$

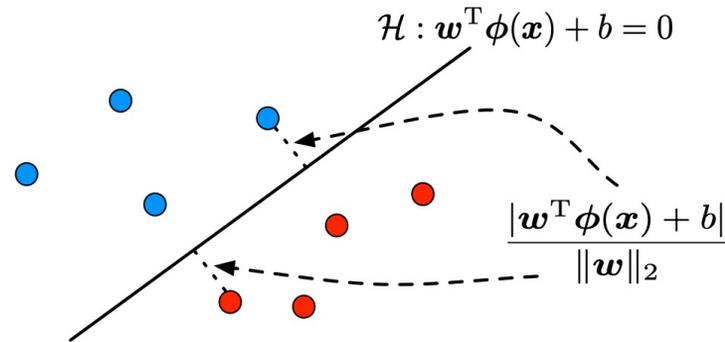
If  $y = -1$ , want  $\mathbf{w}^T \mathbf{x} + b \ll 0$

$\therefore$  want  $y(\mathbf{w}^T \mathbf{x} + b) \gg 0$

# Maximizing margin

**Margin:** the *smallest* distance from all training points to the hyperplane

$$\text{MARGIN OF } (\mathbf{w}, b) = \min_i \frac{y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)}{\|\mathbf{w}\|_2}$$



The intuition “**the further away the better**” translates to solving

$$\max_{\mathbf{w}, b} \min_i \frac{y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)}{\|\mathbf{w}\|_2} = \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \min_i y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)$$

# Maximizing margin, rescaling

**Note:** rescaling  $(w, b)$  by multiplying both by some scalar does not change the hyperplane.

Decision boundary:  $w^T \phi(x) + b = 0 \iff (10^6 w)^T \phi(x) + 10^6 b = 0$

multiply original  $(w, b)$  by  $\frac{1}{\min_i (y_i (w^T \phi(x_i) + b))}$

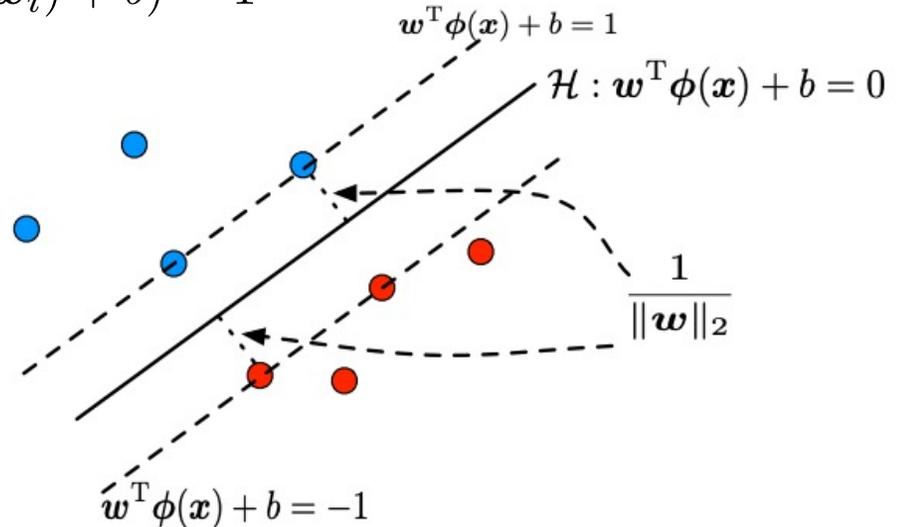
We can thus always scale  $(w, b)$  s.t.  $\min_i y_i (w^T \phi(x_i) + b) = 1$

The margin then becomes

MARGIN OF  $(w, b)$

$$= \frac{1}{\|w\|_2} \min_i y_i (w^T \phi(x_i) + b)$$

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



# SVM for separable data: “Primal” formulation

For a separable training set, we aim to solve

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \quad \text{s.t.} \quad \min_i y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) = 1$$

(this is non-convex)

This is equivalent to

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$
$$\text{s.t.} \quad y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \quad \forall i \in [n]$$

Minimizing a convex  
function with convex  
constraints is convex

SVM is thus also called *max-margin* classifier. The constraints above are called *hard-margin* constraints.

# General non-separable case

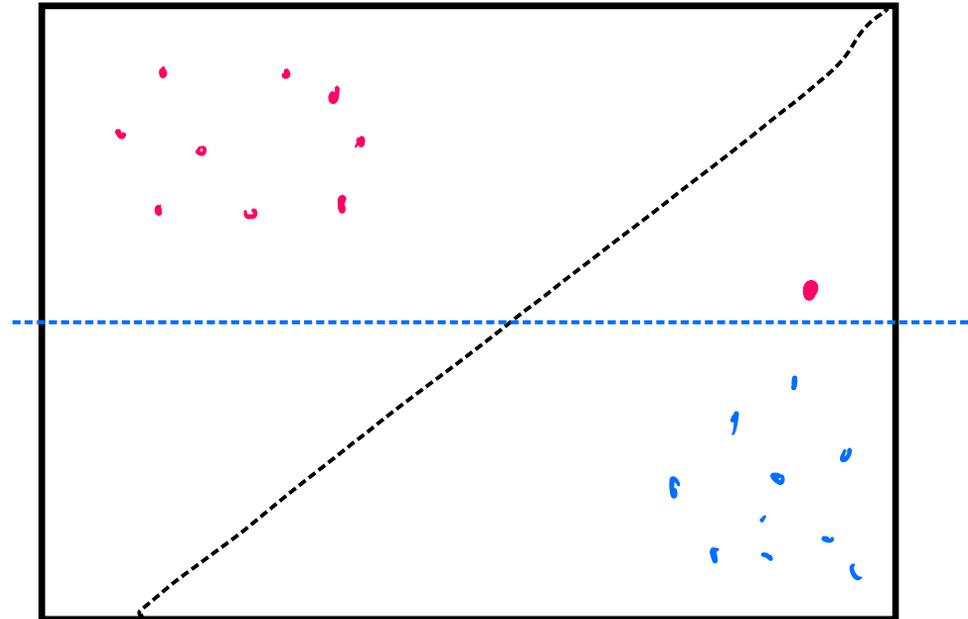
If data is not linearly separable, the previous constraint

$$y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \quad \forall i \in [n]$$

is obviously *not feasible*. What is the right thing to do?

cannot even  
match  
 $\text{sign}(\mathbf{w}^T \phi(\mathbf{x}_i) + b) = y_i$   
 $\forall i \in [n]$

Forcing classifier  
to classify all  
datapoints correctly  
may be bad



## General non-separable case

If data is not linearly separable, the previous constraint  $y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \forall i \in [n]$  is not feasible. And more generally, forcing classifier to always classify all datapoints correctly may not be the best idea.

To deal with this issue, we relax the constraints to  $\ell_1$  **norm soft-margin** constraints:

$$\begin{aligned} y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) &\geq 1 - \xi_i, \quad \forall i \in [n] \\ \iff 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) &\leq \xi_i, \quad \forall i \in [n] \end{aligned}$$

where we introduce **slack variables**  $\xi_i \geq 0$ .

Recall the hinge loss:  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$ . In our case,  $z = y(\mathbf{w}^T \phi(\mathbf{x}) + b)$ .

## Aside: Why $\ell_1$ penalization?



Don't  
just minimize  
the  
objective function

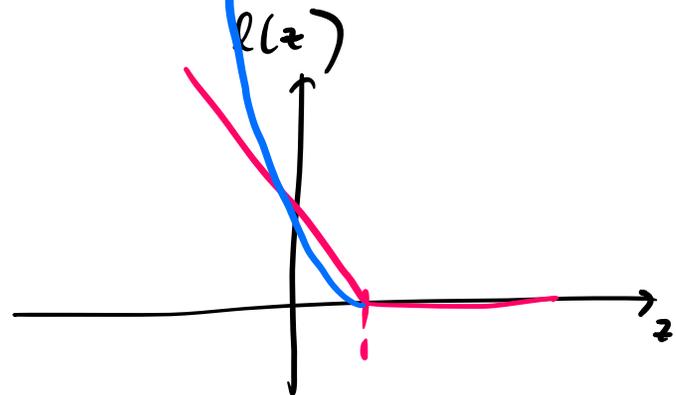


Feel the  
objective  
function!

## Aside: Why $\ell_1$ penalization?

hinge loss  $l(z) = \max(0, 1 - z)$

squared hinge loss  $l(z) = (\max(0, 1 - z))^2$



What would be different?

$z^2$  grows much faster than  $z$

squared hinge loss would really penalize getting some predictions very wrong

## Aside: Why $\ell_1$ penalization?

Because of this, absolute value loss can be more robust to outliers in data compared to squared loss.

A 1-D regression example: mean vs. median

If I have  $x_1, x_2, \dots, x_n$ :

What is  $w_{\ell_2}^* = \arg \min_w \sum_i (x_i - w)^2$ ?  $w_{\ell_2}^* = \frac{\sum_i x_i}{n}$

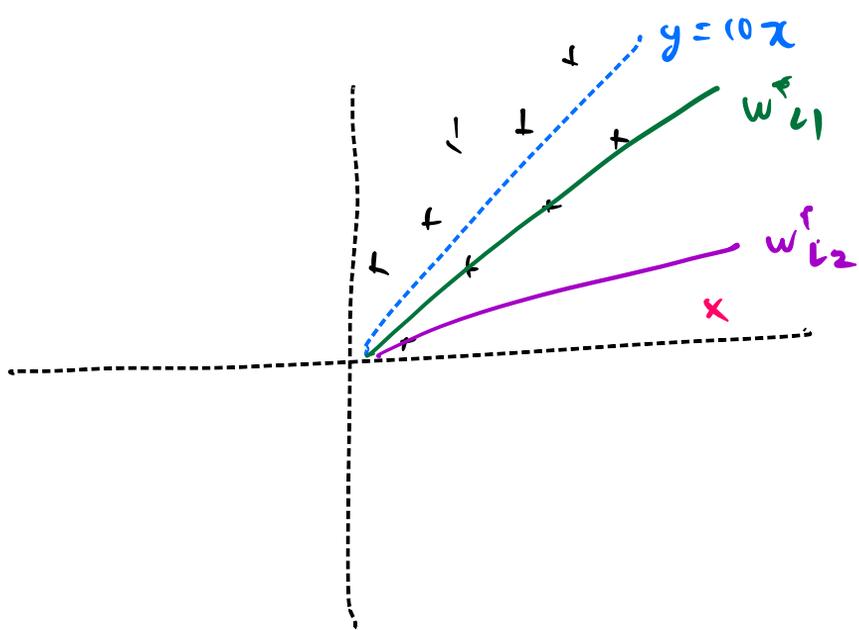
What is  $w_{\ell_1}^* = \arg \min_w \sum_i |x_i - w|$ ?  $w_{\ell_1}^* = \text{median}(x_1, \dots, x_n)$

*Median is more robust to outliers than mean.*

## Aside: Why $\ell_1$ penalization?

For 1-D regression

$$y = 10x + \text{noise}$$



consider

$$w_{L2}^* = \arg \min_w \sum_i (y_i - wx_i)^2$$

$$w_{L1}^* = \arg \min_w \sum_i |y_i - wx_i|$$

## Back to SVM: General non-separable case

If data is not linearly separable, the constraint  $y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \forall i \in [n]$  is not feasible.

To deal with this issue, we relax the constraints to  $\ell_1$  **norm soft-margin** constraints:

$$y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \forall i \in [n]$$

where we introduce **slack variables**  $\xi_i \geq 0$ .

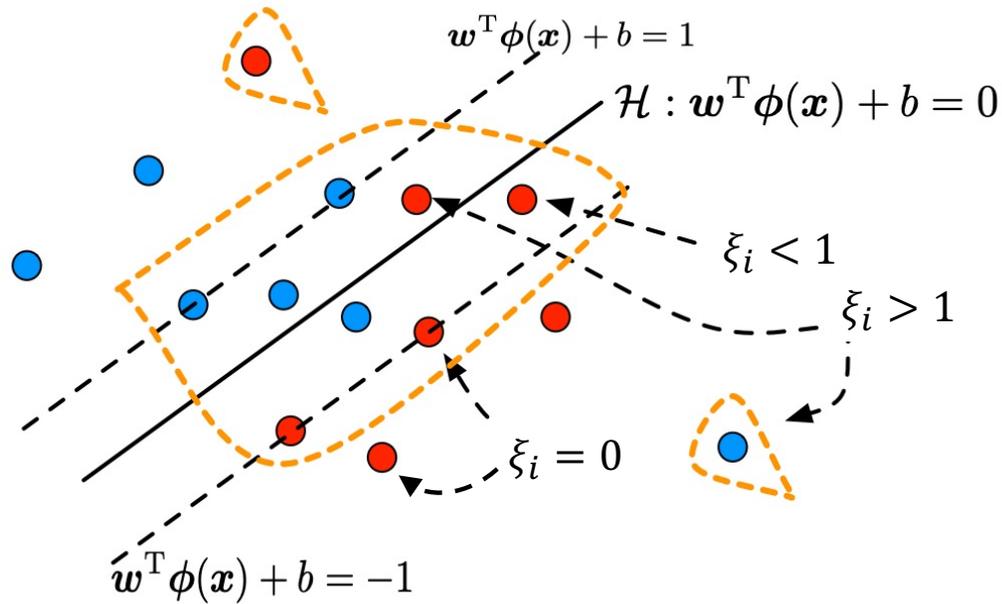
# SVM General Primal Formulation

We want  $\xi_i$  to be as small as possible. The objective becomes

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \forall i \in [n] \\ & \xi_i \geq 0, \quad \forall i \in [n] \end{aligned}$$

where  $C$  is a hyperparameter to balance the two goals.

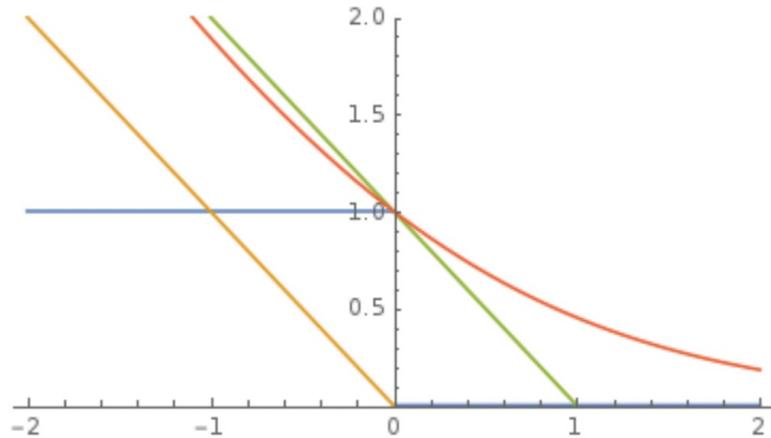
# Understanding the slack conditions



- when  $\xi_i^* = 0$ , point is classified correctly and satisfies large margin constraint.
- when  $\xi_i^* < 1$ , point is classified correctly but does not satisfy large margin constraint.
- when  $\xi_i^* > 1$ , point is misclassified.

# Primal formulation: Another view

In one sentence: **linear model with  $\ell_2$  regularized hinge loss**. Recall:



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\} \rightarrow$  Perceptron
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z)) \rightarrow$  logistic regression
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\} \rightarrow$  **SVM**

# Primal formulation: Another view

For a linear model  $(\mathbf{w}, b)$ , this means

$$\min_{\mathbf{w}, b} \sum_i \max \{0, 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- recall  $y_i \in \{-1, +1\}$
- a nonlinear mapping  $\phi$  is applied
- the bias/intercept term  $b$  is used explicitly (why is this done?)

*What is the relation between this formulation and the one which we just saw before?*

# Equivalent forms

## The formulation

$$\min_{\mathbf{w}, b, \{\xi_i\}} C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \leq \xi_i, \quad \forall i \in [n]$$

$$\xi_i \geq 0, \quad \forall i \in [n]$$

In order to minimize  $\sum_i \xi_i$

we should set  $\xi_i$  to be as small as possible:

is equivalent to

$$\min_{\mathbf{w}, b, \{\xi_i\}} C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } \max \{0, 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)\} = \xi_i, \quad \forall i \in [n]$$

## Equivalent forms

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \max \{0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)\} = \xi_i, \quad \forall i \in [n] \end{aligned}$$

is equivalent to

$$\min_{\mathbf{w}, b} C \sum_i \max \{0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)\} + \frac{1}{2} \|\mathbf{w}\|_2^2$$

and

$$\min_{\mathbf{w}, b} \sum_i \max \{0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

with  $\lambda = 1/C$ . *This is exactly minimizing  $\ell_2$  regularized hinge loss!*

# Optimization

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \forall i \in [n] \\ & \xi_i \geq 0, \quad \forall i \in [n]. \end{aligned}$$

- it is a convex (in fact, a **quadratic**) problem
- thus can apply any convex optimization algorithms, e.g. SGD
- there are **more specialized and efficient** algorithms
- but usually we apply kernel trick, which requires solving the *dual problem*

A diagram illustrating a Support Vector Machine (SVM) classification. It shows two classes of data points: blue circles and red circles. A solid gray line represents the decision boundary, which is the line that best separates the two classes. Two dashed gray lines represent the margins, which are the boundaries of the regions that are furthest from the decision boundary while still containing all data points of their respective classes. The text "SVMs: Dual formulation & Kernel trick" is overlaid on the diagram in a large, bold, red font.

**SVMs:  
Dual formulation  
& Kernel trick**

Recall SVM formulation for separable case,

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } y_i (w^T \phi(x_i) + b) \geq 1 \quad \forall i \in [n]$$

Can we use the kernel trick??

Can we show that  $w^*$  is a linear combination of feature vectors  $\phi(x_i)$ ??

# How did we show this for regularized least squares?

By setting the gradient of  $F(\mathbf{w}) = \|\Phi\mathbf{w} - \mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_2^2$  to be  $\mathbf{0}$ :

$$\Phi^T(\Phi\mathbf{w}^* - \mathbf{y}) + \lambda\mathbf{w}^* = \mathbf{0}$$

we know

$$\mathbf{w}^* = \frac{1}{\lambda}\Phi^T(\mathbf{y} - \Phi\mathbf{w}^*) = \Phi^T\boldsymbol{\alpha} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)$$

Thus the least square solution is **a linear combination of features of the datapoints!**

# Is optimal predictor linear combination of feature vectors?

$$F(\mathbf{w}) = \underbrace{\sum_{i=1}^n \max \{0, 1 - y_i (\mathbf{w}^\top \phi(\mathbf{x}_i) + b)\}}_{\text{hinge loss}} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

This is a convex problem. Therefore, gradient descent (GD) will find a minimizer with any initialization (for some suitable step size).

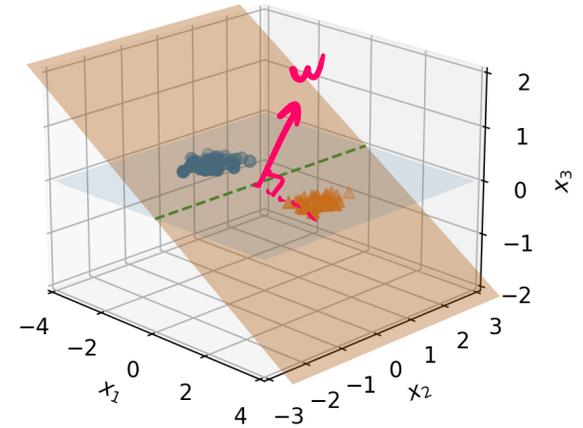
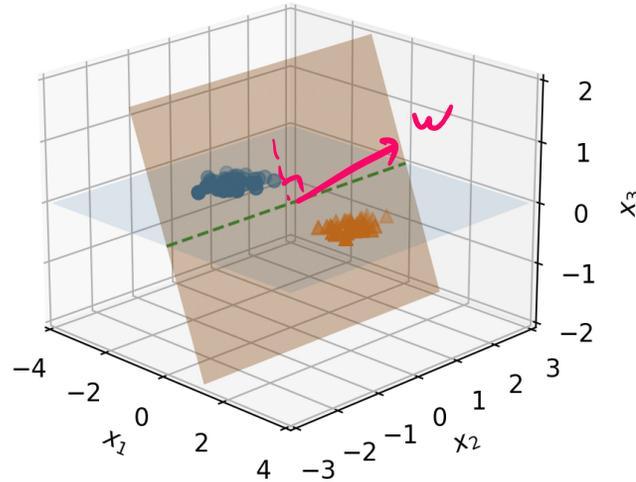
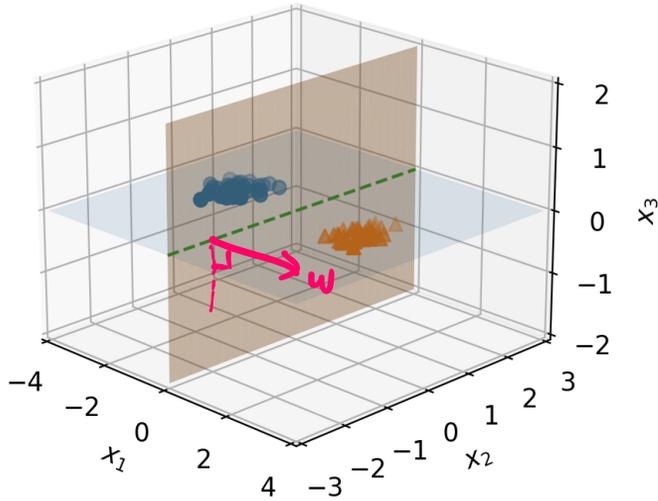
*Exercise:* Suppose we do GD with initialization  $\mathbf{w}^{(0)} = \mathbf{0}$ . Show that gradient descent iterates on  $F(\mathbf{w})$  at any time step  $t$  satisfy  $\mathbf{w}^{(t)} = \sum_{i=1}^n \alpha_i^{(t)} y_i \phi(\mathbf{x}_i)$  for some  $\alpha_i^{(t)}$ .

For the SVM problem,  $\mathbf{w}^* = \sum_{i=1}^n d_i^* y_i \phi(\mathbf{x}_i)$

We can also geometrically understand why  $w^*$  should lie in the span of the data:

-ve

+ve



$$w_3 =$$

$$w_3 > 0$$

$$w_3 \gg 0$$

all datapoints have  $\tau_3 = 0$

# Kernelizing SVM

If  $w = \sum_{i=1}^n d_i y_i \phi(x_i)$ , how can we use this?

$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } y_j (w^T \phi(x_j) + b) \geq 1, \quad \forall j \in [n].$$

→

$$\min_{d, b} \frac{1}{2} \left\| \sum_i d_i y_i \phi(x_i) \right\|_2^2$$

$$\text{s.t. } y_j \left( \left( \sum_i d_i y_i \phi(x_i) \right)^T \phi(x_j) + b \right) \geq 1$$

This is equivalent to

$$\min_{d, b} \frac{1}{2} \sum_i \sum_j d_i d_j y_i y_j \phi(x_i)^T \phi(x_j)$$

$$\text{s.t. } y_j \left( \sum_i d_i y_i \phi(x_i)^T \phi(x_j) + b \right) \geq 1 \quad \forall i \in [n]$$

# SVM: Dual form for separable case

With some optimization theory (Lagrange duality, not covered in this class), we can show this is equivalent to,

$$\begin{aligned} \max_{\{\alpha_i\}} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \geq 0, \quad \forall i \in [n] \end{aligned}$$

# SVM: Dual form for separable case

Using the kernel function  $k$  for the mapping  $\phi$ , we can kernelize this!

$$\begin{aligned} \max_{\{\alpha_i\}} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \geq 0, \quad \forall i \in [n] \end{aligned}$$

No need to compute  $\phi(x)$ . This is also a **quadratic program** and many efficient optimization algorithms exist.

# SVM: Dual form for the general case

For the primal for the general (non-separable) case:

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \forall i \in [n] \\ & \xi_i \geq 0, \quad \forall i \in [n]. \end{aligned}$$

The dual is very similar,

$$\begin{aligned} \max_{\{\alpha_i\}} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i \leq C, \quad \forall i \in [n]. \end{aligned}$$

# Prediction using SVM

How do we predict given the solution  $\{\alpha_i^*\}$  to the dual optimization problem?

Remember that,

$$\mathbf{w}^* = \sum_i \alpha_i^* y_i \phi(\mathbf{x}_i) = \sum_{i:\alpha_i^* > 0} \alpha_i^* y_i \phi(\mathbf{x}_i)$$

A point with  $\alpha_i^* > 0$  is called a “**support vector**”. Hence the name SVM.

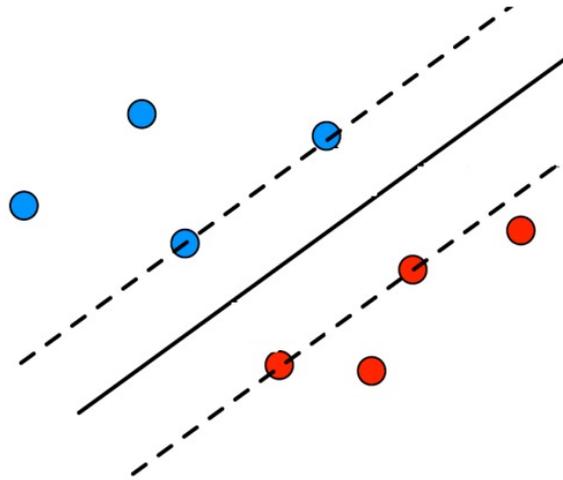
To make a prediction on any datapoint  $\mathbf{x}$ ,

$$\begin{aligned} \text{sign} \left( \mathbf{w}^{*\top} \phi(\mathbf{x}) + b^* \right) &= \text{sign} \left( \sum_{i:\alpha_i^* > 0} \alpha_i^* y_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}) + b^* \right) \\ &= \text{sign} \left( \sum_{i:\alpha_i^* > 0} \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}) + b^* \right). \end{aligned}$$

All we need now is to identify  $b^*$ .

# How to compute bias term $b^*$ ?

We will only consider the separable case (general case is not too different):



It can be shown (we will not cover in class), that in the separable case the support vectors lie on the margin.

$$y_i (\mathbf{w}^* \cdot \phi(x_i) + b^*) = 1 \quad \Rightarrow \quad y_i^2 (\mathbf{w}^* \cdot \phi(x_i) + b^*) = y_i$$

$$\Rightarrow (\mathbf{w}^*) \cdot \phi(x_i) + b^* = y_i$$

$$\Rightarrow b^* = y_i - \mathbf{w}^* \cdot \phi(x_i) = y_i - \sum_{j: \alpha_j^* > 0} \alpha_j^* y_j k(x_j, x_i), \quad \text{for any } i \text{ s.t. } \alpha_i^* > 0$$



**SVMs:  
Understanding  
them further**

The diagram illustrates a linear Support Vector Machine (SVM) in a 2D space. It shows two classes of data points: blue circles and red circles. A solid gray line represents the decision boundary, which is the hyperplane that separates the two classes. Two dashed gray lines represent the margins, which are parallel to the decision boundary and define the maximum distance from the decision boundary to the nearest data points of each class. The blue points are located above the decision boundary, and the red points are located below it. The margins are wider for the blue class than for the red class.

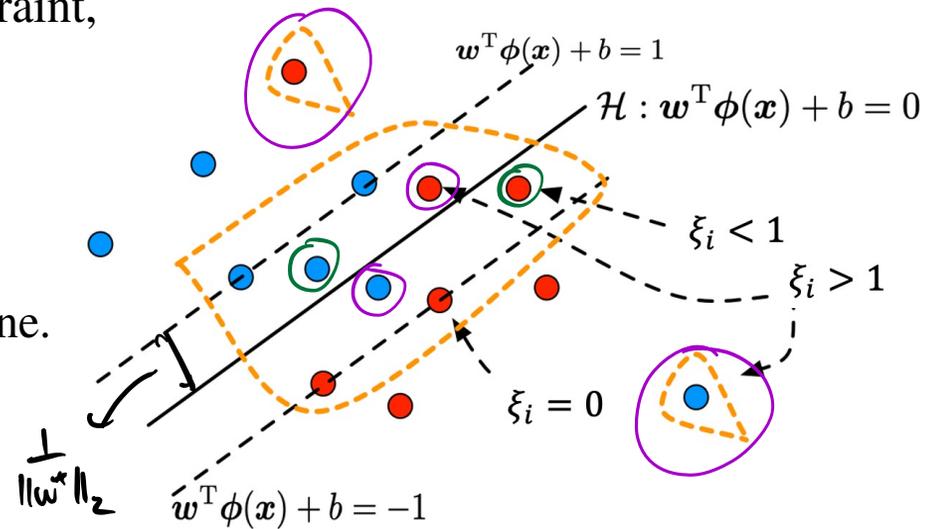
# Understanding support vectors

Support vectors are  $\phi(\mathbf{x}_i)$  such that  $\alpha_i^* > 0$ .

They are the set of points which satisfy one of the following:

- (1) they are tight with respect to the large margin constraint,
- (2) they do not satisfy the large margin constraint,
- (3) they are misclassified.

- when  $\xi_i^* = 0$ ,  $y_i(\mathbf{w}^{*\top} \phi(\mathbf{x}_i) + b^*) = 1$ , and thus the point is  $1/\|\mathbf{w}^*\|_2$  away from the hyperplane.
- when  $\xi_i^* < 1$ , the point is classified correctly but does not satisfy the large margin constraint.
- when  $\xi_i^* > 1$ , the point is misclassified.



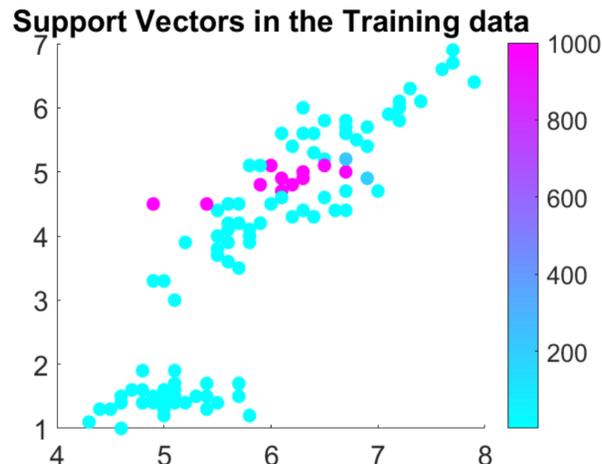
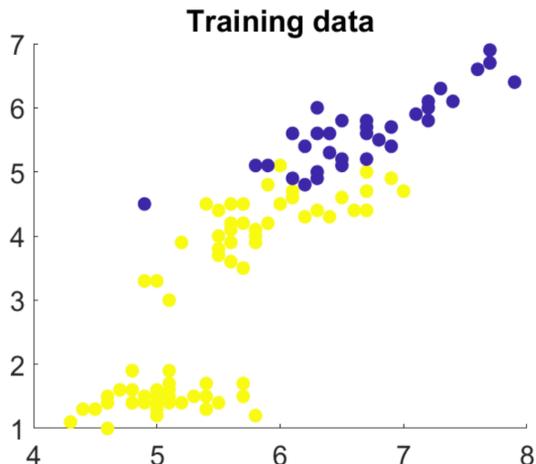
Support vectors (circled with the orange line) are the only points that matter!

# Understanding support vectors

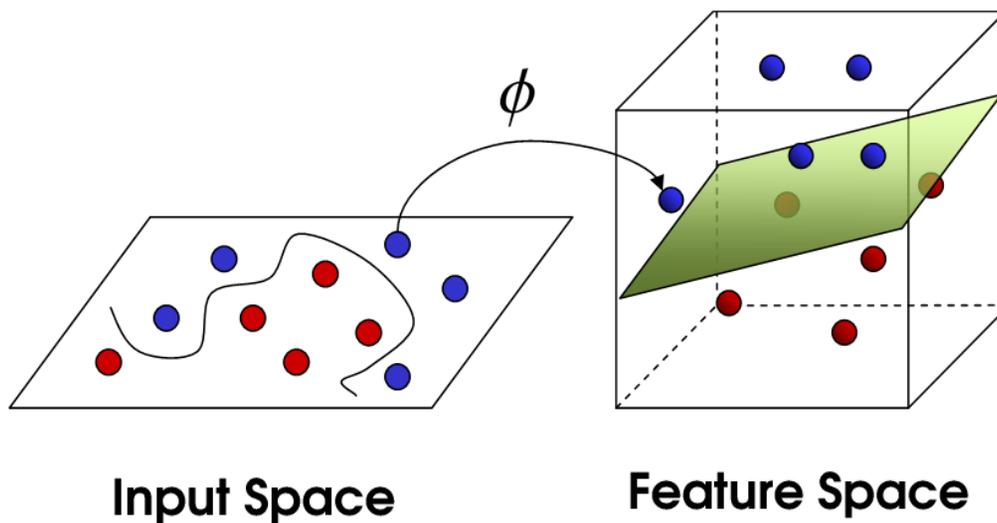
One potential drawback of kernel methods: **non-parametric**, need to potentially keep all the training points.

$$\text{sign} \left( \mathbf{w}^{*\text{T}} \phi(\mathbf{x}) + b^* \right) = \text{sign} \left( \sum_{i=1}^n \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}) + b^* \right).$$

For SVM though, very often  $\#\text{support vectors} = |\{i : \alpha_i^* > 0\}| \ll n$ .

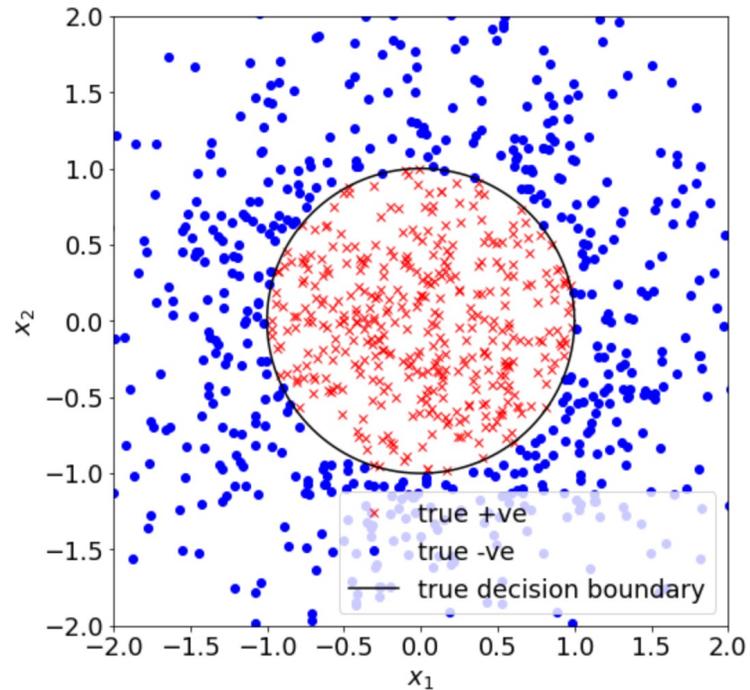


# Examining the effect of kernels



Data may become linearly separable when lifted to the high-dimensional feature space!

# Polynomial kernel: example



Switch to Colab

# Gaussian kernel: example

## Gaussian kernel or Radial basis function (RBF) kernel

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

for some  $\sigma > 0$ . This is also parameterized as,

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma\|\mathbf{x} - \mathbf{x}'\|_2^2)$$

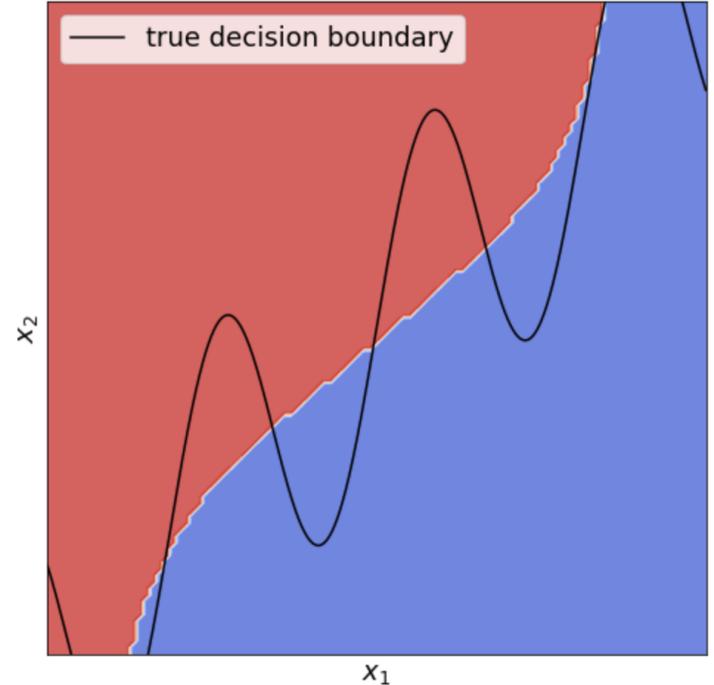
for some  $\gamma > 0$ .

What does the decision boundary look like?

What is the effect of  $\gamma$ ?

Note that the prediction is of the form

$$\text{sign}\left(\mathbf{w}^{*\top}\phi(\mathbf{x}) + b^*\right) = \text{sign}\left(\sum_{i:\alpha_i^* > 0} \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}) + b^*\right).$$



Switch to Colab

# SVM: Summary of mathematical forms

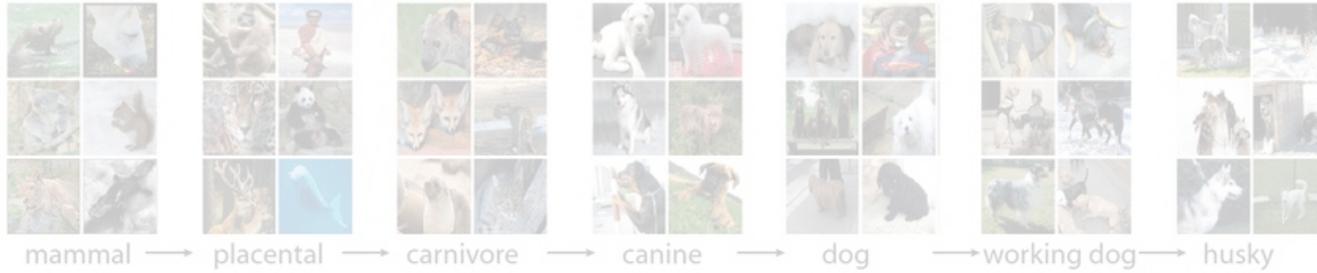
SVM: **max-margin linear classifier**

**Primal** (equivalent to minimizing  $\ell_2$  regularized hinge loss):

$$\begin{aligned} \min_{\mathbf{w}, b, \{\xi_i\}} \quad & C \sum_i \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \forall i \in [n] \\ & \xi_i \geq 0, \quad \forall i \in [n]. \end{aligned}$$

**Dual** (kernelizable, reveals what training points are support vectors):

$$\begin{aligned} \max_{\{\alpha_i\}} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j \boldsymbol{\phi}(\mathbf{x}_i)^\top \boldsymbol{\phi}(\mathbf{x}_j) \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i \leq C, \quad \forall i \in [n]. \end{aligned}$$



# Multiclass classification

# Setup

Recall the setup:

- input (feature vector):  $\mathbf{x} \in \mathbb{R}^d$
- output (label):  $y \in [C] = \{1, 2, \dots, C\}$
- goal: learn a mapping  $f : \mathbb{R}^d \rightarrow [C]$

**Examples:**

- recognizing digits ( $C = 10$ ) or letters ( $C = 26$  or  $52$ )
- predicting weather: sunny, cloudy, rainy, etc
- predicting image category: ImageNet dataset ( $C \approx 20K$ )

# Linear models: Binary to multiclass

Step 1: *What should a linear model look like for multiclass tasks?*

Note: a linear model for binary tasks (switching from  $\{-1, +1\}$  to  $\{1, 2\}$ )

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ 2 & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

can be written as

$$\begin{aligned} f(\mathbf{x}) &= \begin{cases} 1 & \text{if } \mathbf{w}_1^T \mathbf{x} \geq \mathbf{w}_2^T \mathbf{x} \\ 2 & \text{if } \mathbf{w}_2^T \mathbf{x} > \mathbf{w}_1^T \mathbf{x} \end{cases} \\ &= \operatorname{argmax}_{k \in \{1, 2\}} \mathbf{w}_k^T \mathbf{x} \end{aligned}$$

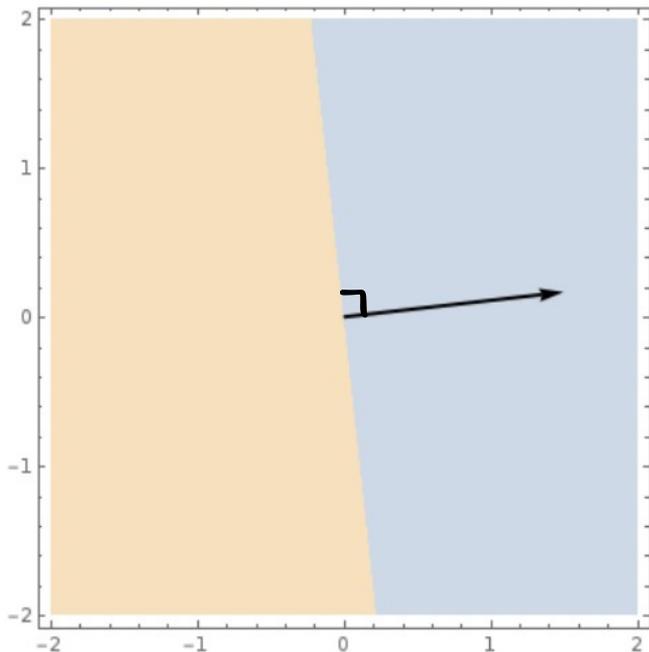
$$\mathbf{w}_1 \neq \mathbf{w}_2 \quad \text{s.t.}$$

$$\mathbf{w}_1 - \mathbf{w}_2 = \mathbf{w}$$

for any  $\mathbf{w}_1, \mathbf{w}_2$  s.t.  $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$

Think of  $\mathbf{w}_k^T \mathbf{x}$  as **a score for class  $k$** .

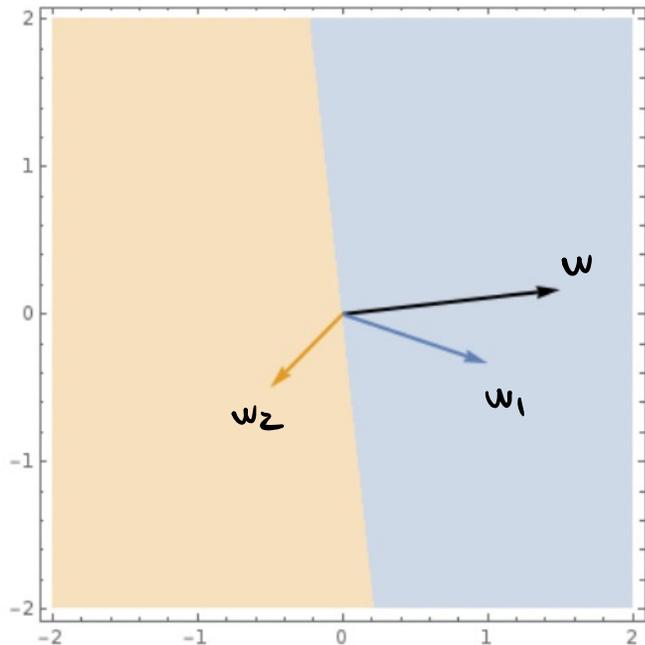
## Linear models: Binary to multiclass



$$w = \left(\frac{3}{2}, \frac{1}{6}\right)$$

- **Blue class:**  
 $\{x : w^T x \geq 0\}$
- **Orange class:**  
 $\{x : w^T x < 0\}$

# Linear models: Binary to multiclass



$$w = \left(\frac{3}{2}, \frac{1}{6}\right) = w_1 - w_2$$

$$w_1 = \left(1, -\frac{1}{3}\right)$$

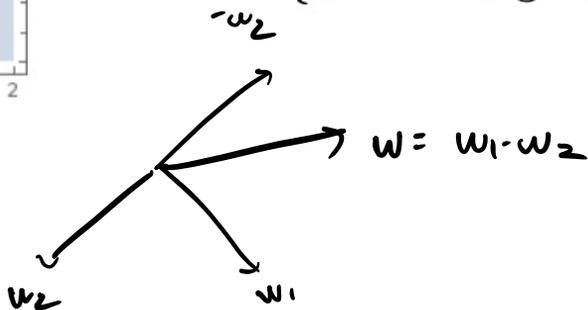
$$w_2 = \left(-\frac{1}{2}, -\frac{1}{2}\right)$$

- Blue class:

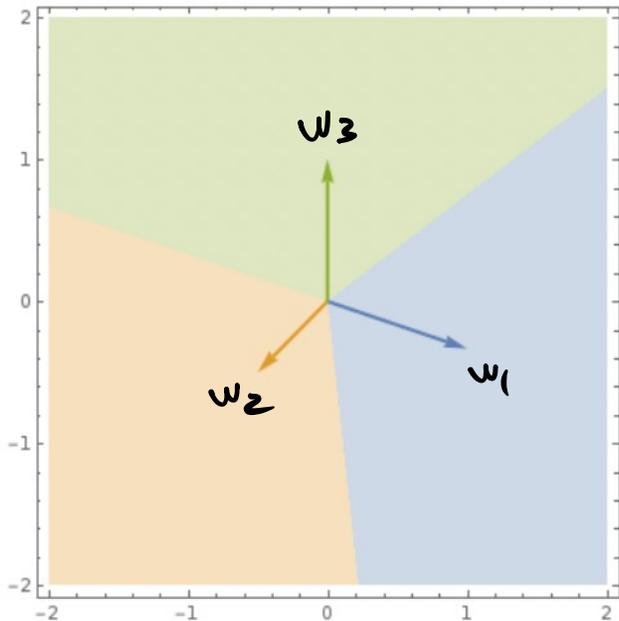
$$\{x : 1 = \operatorname{argmax}_k w_k^T x\}$$

- Orange class:

$$\{x : 2 = \operatorname{argmax}_k w_k^T x\}$$



# Linear models: Binary to multiclass



$$w_1 = (1, -\frac{1}{3})$$
$$w_2 = (-\frac{1}{2}, -\frac{1}{2})$$
$$w_3 = (0, 1)$$

- Blue class:  
 $\{x : 1 = \operatorname{argmax}_k w_k^T x\}$
- Orange class:  
 $\{x : 2 = \operatorname{argmax}_k w_k^T x\}$
- Green class:  
 $\{x : 3 = \operatorname{argmax}_k w_k^T x\}$

## Function class: Linear models for multiclass classification

$$\mathcal{F} = \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x} \mid \mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^d \right\}$$
$$= \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} (\mathbf{W} \mathbf{x})_k \mid \mathbf{W} \in \mathbb{R}^{C \times d} \right\}$$

Next, lets try to generalize the loss functions. Focus on the logistic loss today.

