

CSCI 567: Machine Learning

Vatsal Sharan
Fall 2022

Lecture 5, Sep 22

Administrivia

- HW2 due in about a week.
- Quiz 1 in 2 weeks.

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

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

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

Kernel methods are exactly about computing inner products *without explicitly computing ϕ* . The exact form of ϕ is inessential; *all we need to do is know the inner products $\phi(x)^T \phi(x')$* .

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 $\phi(\mathbf{x}')$?

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

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 $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

Popular kernels:

1. Polynomial kernel

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M$$

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

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

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

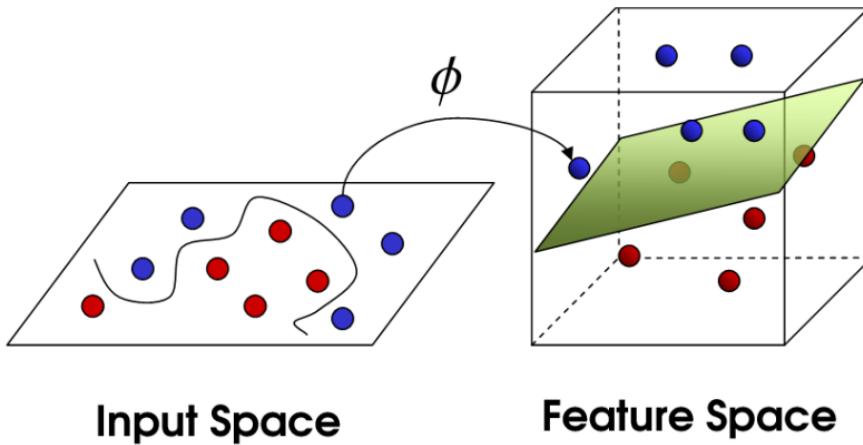
Prediction with kernels

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.

Support vector machines (SVMs)



1.1 Why study SVM?

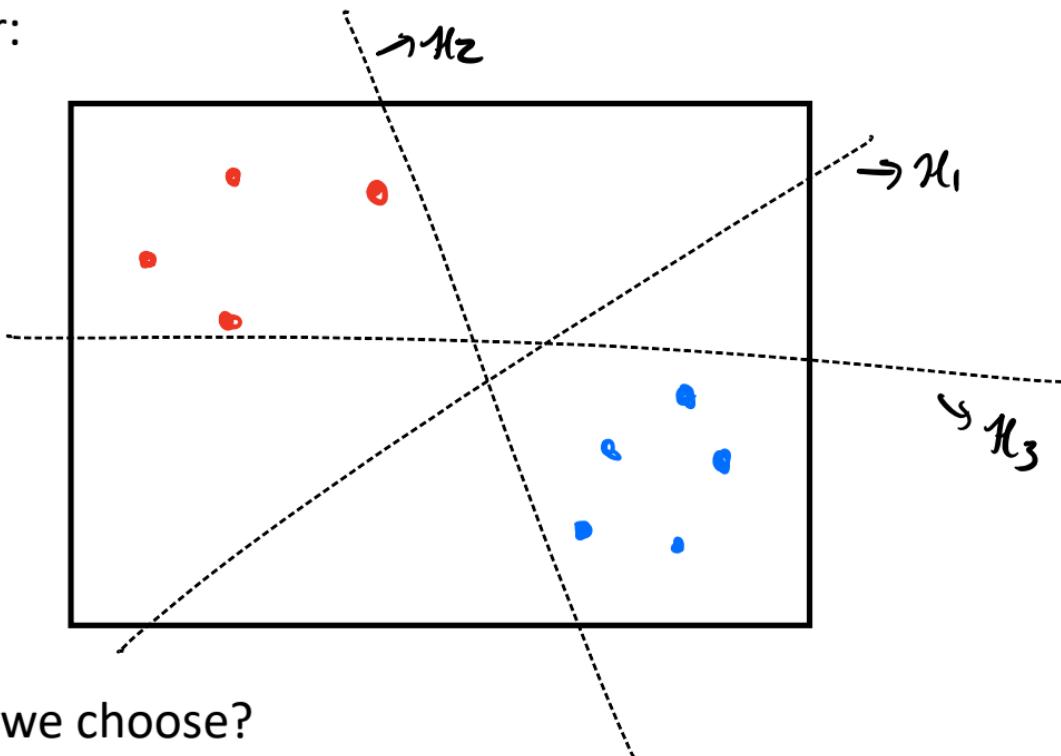
- One of the most commonly used classification algorithms
- Allows us to explore the concept of *margins* in classification
- Works well with the kernel trick
- Strong theoretical guarantees

We focus on **binary classification** here.

The *function class for SVMs is a linear function on a feature map ϕ 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.

1.2 Margins: separable case, geometric intuition

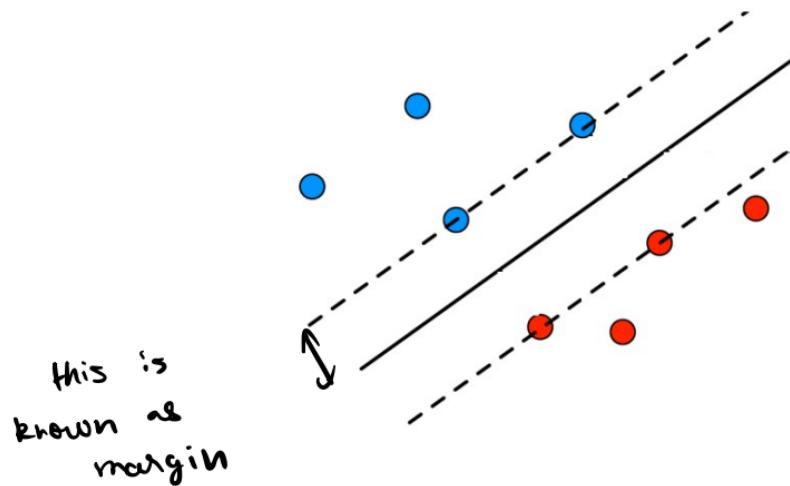
When data is **linearly separable**, there are infinitely many hyperplanes with zero training error:



Which one should we choose?

1.2 Margins: separable case, geometric intuition

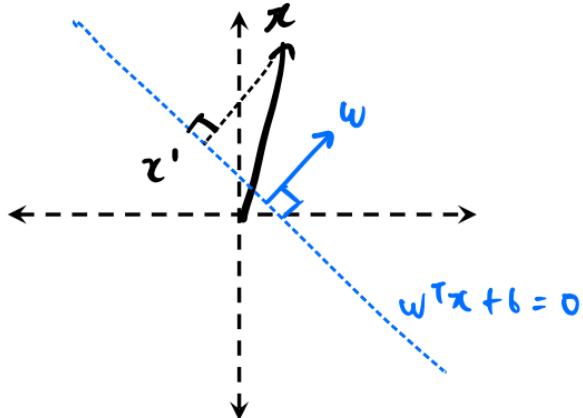
The further away the separating hyperplane is from the datapoints, the better.



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

1.2 Formalizing geometric intuition: Distance to hyperplane

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



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

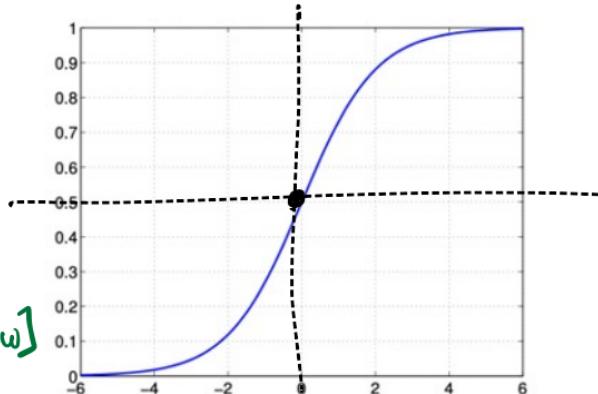
$$0 = \mathbf{w}^T \left(x - \beta \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \right) + b = \mathbf{w}^T x - \beta \|\mathbf{w}\| + b \implies \beta = \frac{\mathbf{w}^T x + b}{\|\mathbf{w}\|_2}.$$

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

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

$$\text{sign}(\mathbf{w}^T x + b) = y$$

1.2 Margins: functional motivation



This should be $\Pr[y = 1 | x; w]$

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

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

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

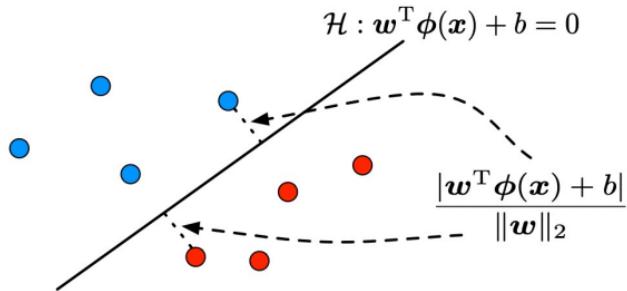
∴ want $y(w^T x + b) \gg 0$

1.3 Maximizing margin

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

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

data
 $\left\{ (\mathbf{x}_i, y_i) \right\}, i \in \{1, \dots, n\}$



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

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

1.3 Maximizing margin, rescaling

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

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

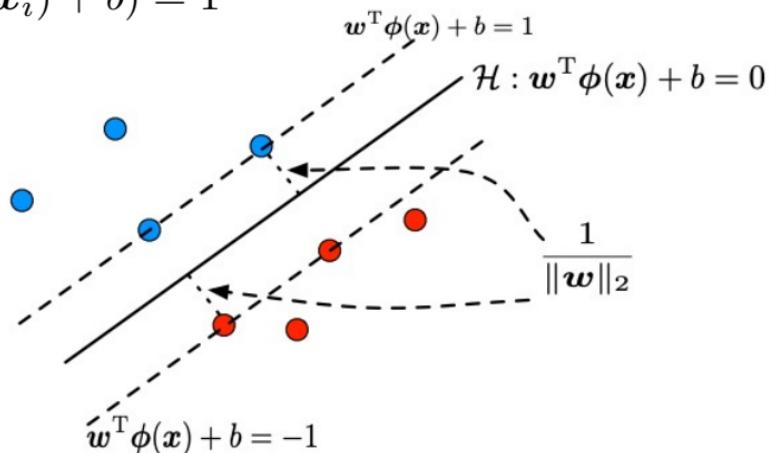
We can thus always scale (w, b) s.t. $\min_i y_i(w^T \phi(x_i) + b) = 1$

The margin then becomes

MARGIN OF (w, b)

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

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



1.4 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

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

This is convex! Minimizing a
convex function
with convex constraint
is convex.

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

1.5 General non-separable case

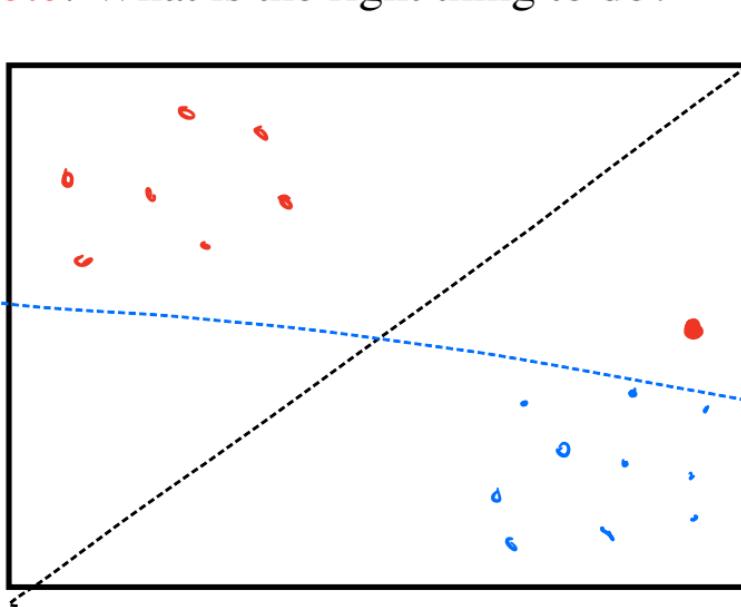
If data is not linearly separable, the previous constraint

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

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

Can't even match
 $\text{sign}(\mathbf{w}^T \phi(\mathbf{x}) + b)$
& $y \neq \pm 1 \in [n]$,
if not linearly
separable.

Even if data is
linearly separable,
should we always
separate it?



forcing classifier to
classifying all datapoints
correctly might not
be good.

1.5 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.

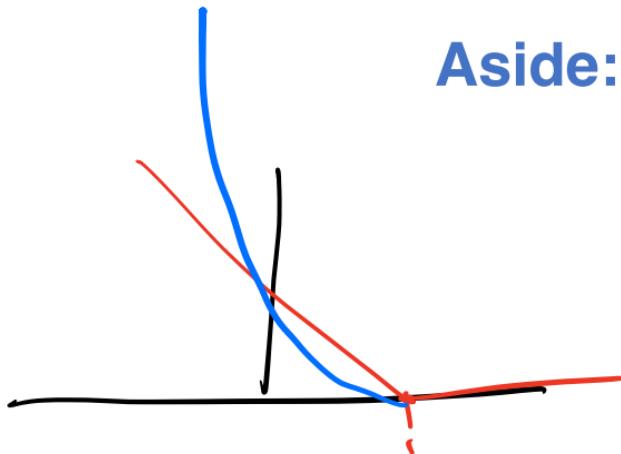
To deal with this issue, we relax the constraints to ℓ_1 norm soft-margin constraints:

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

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 ℓ_1 penalization?



$$\text{hinge loss } \ell(z) = \max(0, 1-z)$$

$$\text{squared hinge loss } \ell(z) = \max(0, 1-z)^2$$

what would be different?

χ^2 grows much faster than χ .

Squared hinge loss would really penalize getting some predictions wrong.

Aside: Why ℓ_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 $w_{L2}^* = \underset{w}{\operatorname{argmin}} \sum_i (x_i - w)^2$? $w_{L2}^* = \frac{\sum x_i}{n}$

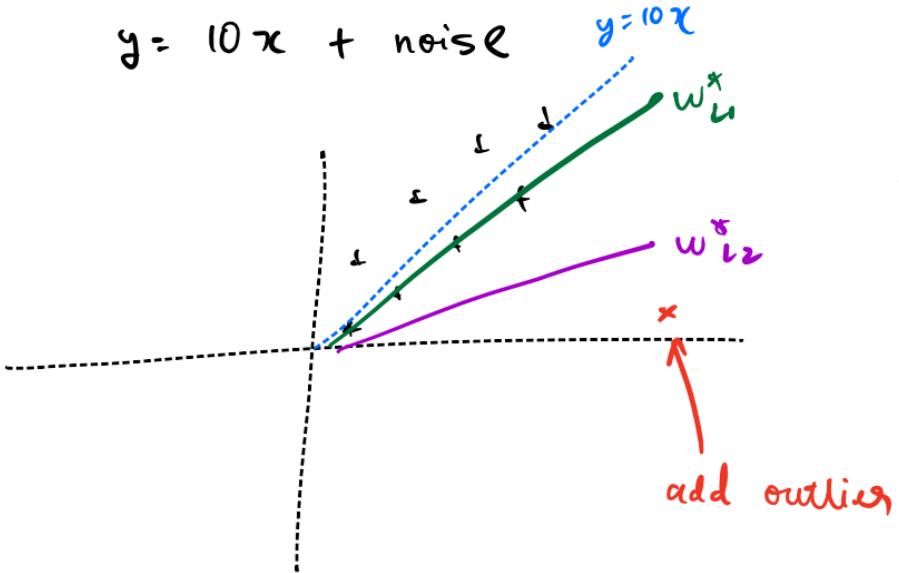
median is more
robust to outliers
than mean

What is $w_{L1}^* = \underset{w}{\operatorname{argmin}} \sum_i |x_i - w|$? $w_{L1}^* = \operatorname{median}(x_1, \dots, x_n)$

Aside: Why ℓ_1 penalization?

for 1-D regression.

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



consider

$$w_{L2}^* = \underset{i}{\operatorname{argmin}} \sum (y_i - w x_i)^2$$

$$w_{L1}^* = \underset{i}{\operatorname{argmin}} \sum |y_i - w x_i|$$

1.5 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 ℓ_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 $\xi_i \geq 0$.

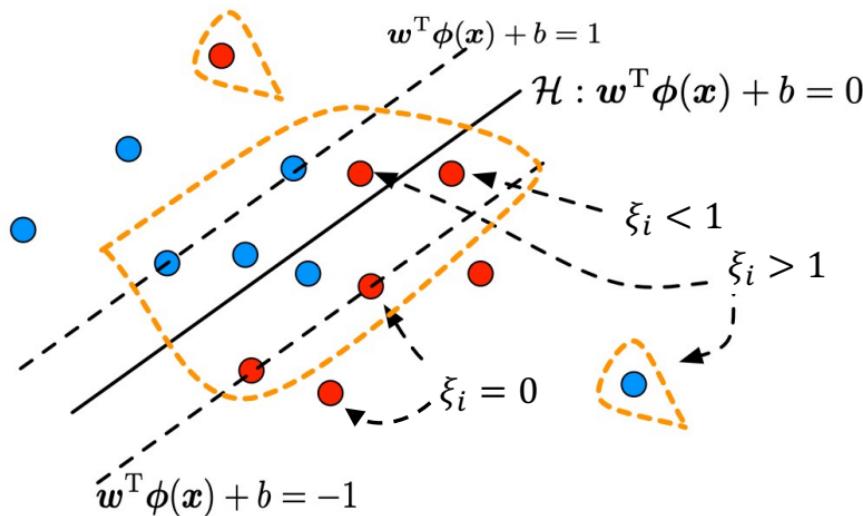
1.5 SVM General Primal Formulation

We want ξ_i to be as small as possible. The objective becomes

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

where $\textcolor{red}{C}$ is a hyperparameter to balance the two goals.

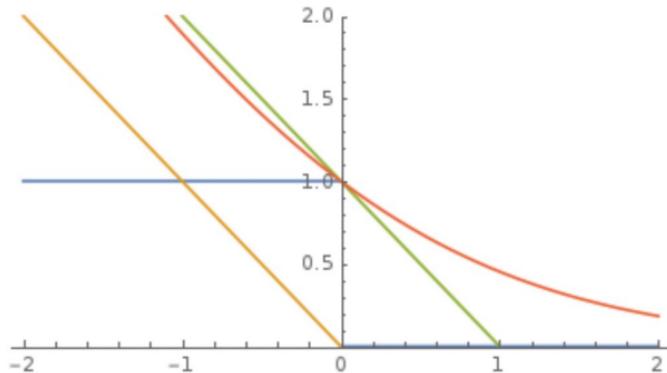
1.6 Understanding the slack conditions



- 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.

1.7 Primal formulation: Another view

In one sentence: **linear model with ℓ_2 regularized hinge loss**. Recall:



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

1.7 Primal formulation: Another view

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

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

- recall $y_i \in \{-1, +1\}$
- a nonlinear mapping ϕ 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?

1.7 Equivalent forms

The formulation

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

In order to minimize $\sum_i \xi_i$,

we should set ξ_i to be as small as possible, which is:

is equivalent to

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

1.7 Equivalent forms

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

is equivalent to

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

and

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

with $\lambda = 1/C$. *This is exactly minimizing ℓ_2 regularized hinge loss!*

1.8 Optimization

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

- 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*



SVMs:
Dual formulation
& Kernel trick

A scatter plot illustrating a linear SVM model. The plot features several blue circles (labeled 'A') and red circles (labeled 'B') scattered across a white background. A solid gray line, representing the decision boundary, runs diagonally from the bottom-left towards the top-right. Two dashed gray lines, representing the margin, extend from the decision boundary. The distance between the closest blue and red points and the decision boundary is labeled 'Margin'. The text 'SVMs:' is positioned above the plot, followed by 'Dual formulation' and '& Kernel trick' on separate lines.

Recall SVM formulation for separable case :

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

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

(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!**

2.1 Kernelizing SVM

(Claim: For the SVM problem, $w^* = \sum_i \alpha_i y_i \phi(x_i)$)

Informal Proof:

formulation as a linear model with ℓ_2 regularized hinge loss.

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

brace under max term labeled "hinge loss"

This is a convex problem \therefore GD will find a minimizer with any initialization (for some appropriate step size).

Recall $\ell_{\text{hinge}}(z) = \max(0, 1 - z)$

$$\frac{\partial F(\omega)}{\partial \omega} = \sum_{i=1}^n \left(\frac{\partial \ell_{\text{hinge}}(z)}{\partial z} \Bigg| \begin{array}{l} (-y_i \phi(x_i)) \\ z = y_i(\omega^\top \phi(x_i) + b) \end{array} \right) + \lambda \omega^{(t)}$$

$$\omega^{(0)} \leftarrow 0$$

$$\omega^{(t+1)} \leftarrow \omega^{(t)} - \eta \sum_{i=1}^n \left(\frac{\partial \ell_{\text{hinge}}(z)}{\partial z} \Bigg| \begin{array}{l} (-y_i \phi(x_i)) \\ z = y_i(\omega^\top \phi(x_i) + b) \end{array} \right) + \lambda \omega^{(t)}$$

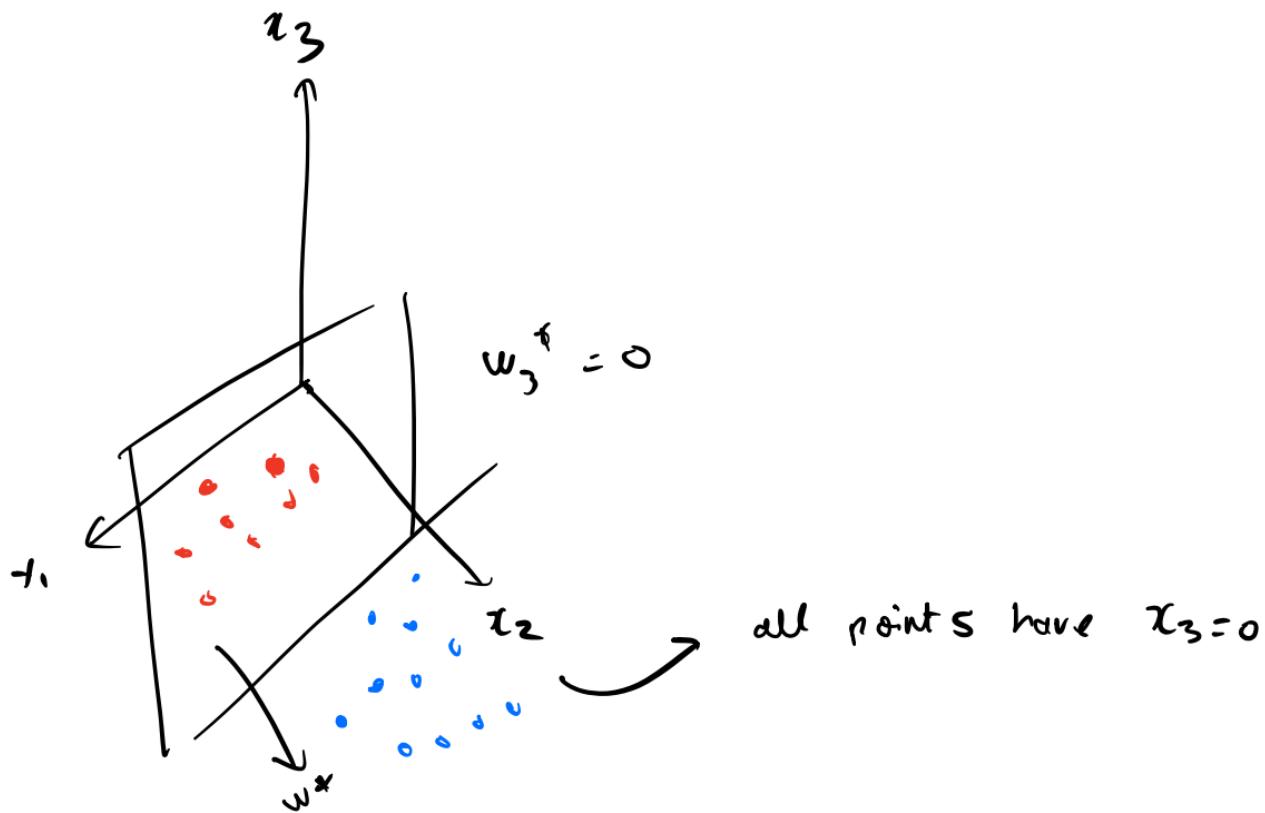
$\therefore w^{(t)}$ always lie in span of $\phi(x_i)$

$$w^{(t)} = \sum d_i^{(t)} y_i \phi(x_i) + t, \text{ for some } d_i^{(t)}$$

$$\therefore w^* = \sum d_i^* y_i \phi(x_i) \text{ for some } d_i^*$$

■

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



If $w = \sum_i d_i y_i \phi(x_i)$, how can we use this?

$$\begin{array}{ll} \min_{w,b} & \frac{1}{2} \|w\|_2^2 \\ \text{s.t. } & y_j(w^T \phi(x_j) + b) \geq 1, \quad \forall j \in [n]. \end{array} \rightarrow \begin{array}{ll} \min_{d,b} & \frac{1}{2} \left\| \sum_i d_i y_i \phi(x_i) \right\|_2^2 \\ \text{s.t. } & y_j \left(\sum_i d_i y_i \phi(x_i)^T \phi(x_j) + b \right) \geq 1 \end{array}$$

This is equivalent to,

$$\begin{array}{ll} \min_{d,b} & \frac{1}{2} \sum_i \sum_j d_i d_j y_i y_j \phi(x_i)^T \phi(x_j) \\ \text{s.t. } & y_j \left(\sum_i d_i y_i \phi(x_i)^T \phi(x_j) + b \right) \geq 1 \quad \forall i \in [n]. \end{array}$$

2.2 SVM: Dual form for separable case

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

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

2.2 SVM: Dual form for separable case

Using the kernel function k for the mapping ϕ , 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 $\phi(\mathbf{x})$. This is also a **quadratic program** and many efficient optimization algorithms exist.

2.3 SVM: Dual form for the general case

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

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

The dual is very similar,

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

2.4 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.

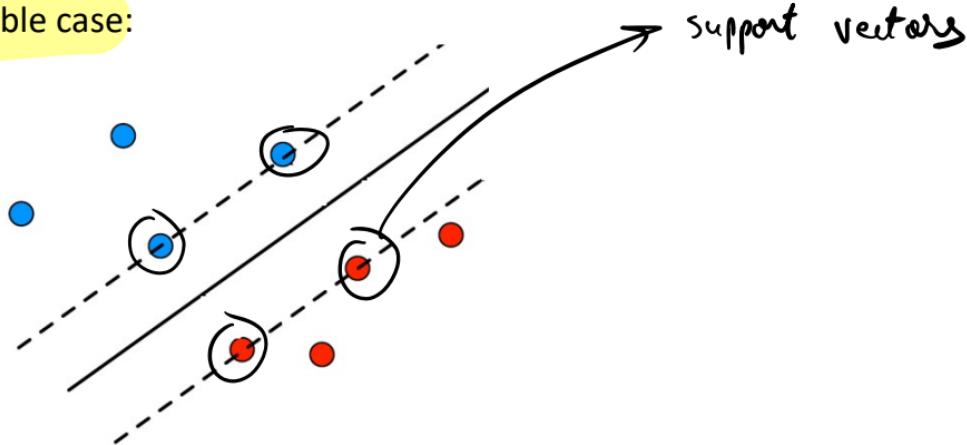
To make a prediction on any datapoint \mathbf{x} ,

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

All we need now is to identify b^* .

2.5 Bias term b^*

First, let's consider the separable case:



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

$$\begin{aligned} y_i (\hat{w}^\top \phi(x_i) + b^*) &= 1 \Rightarrow y_i^2 (\hat{w}^{*\top} \phi(x_i) + b^*) = y_i \\ \Rightarrow \hat{w}^{*\top} \phi(x_i) + b^* &= y_i \\ \Rightarrow b^* &= y_i - \hat{w}^{*\top} \phi(x_i) \quad \text{for any } i \text{ s.t. } d_i > 0. \end{aligned}$$

2.5 Bias term b^*

General (non-separable case):

For any support vector $\phi(\mathbf{x}_i)$ with $0 < \alpha_i^* < C$, it can be shown that $1 = y_i(\mathbf{w}^{*\top}\phi(\mathbf{x}_i) + b^*)$ (i.e. that support vector lies on the margin). Therefore, as before,

$$b^* = y_i - \mathbf{w}^{*\top}\phi(\mathbf{x}_i) = y_i - \sum_{j=1}^n \alpha_j^* y_j k(\mathbf{x}_j, \mathbf{x}_i).$$

In practice, often *average* over all i with $0 < \alpha_i^* < C$ to stabilize computation.

With α^* and b^* in hand, we can make a prediction on any datapoint \mathbf{x} ,

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



SVMs: Understanding them further

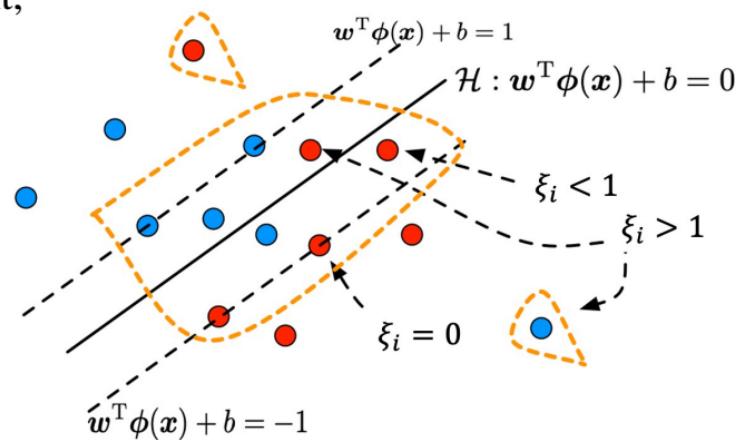
3.1 Understanding support vectors

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 $\xi_i^* = 0$, $y_i(\mathbf{w}^{*T}\phi(\mathbf{x}_i) + b^*) = 1$,
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.



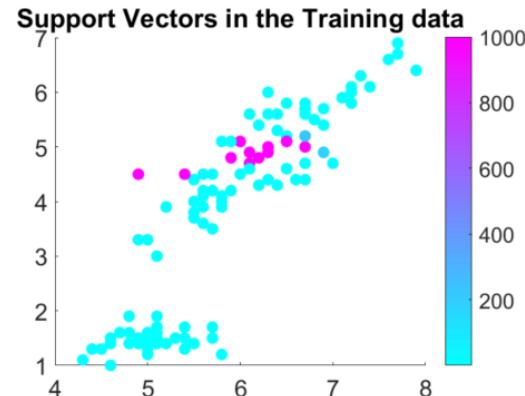
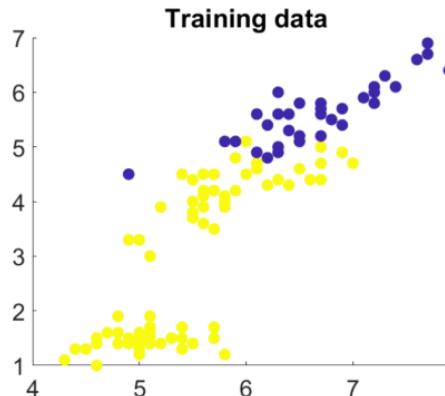
Support vectors (circled with the orange line) are the only points that matter!

3.1 Understanding support vectors

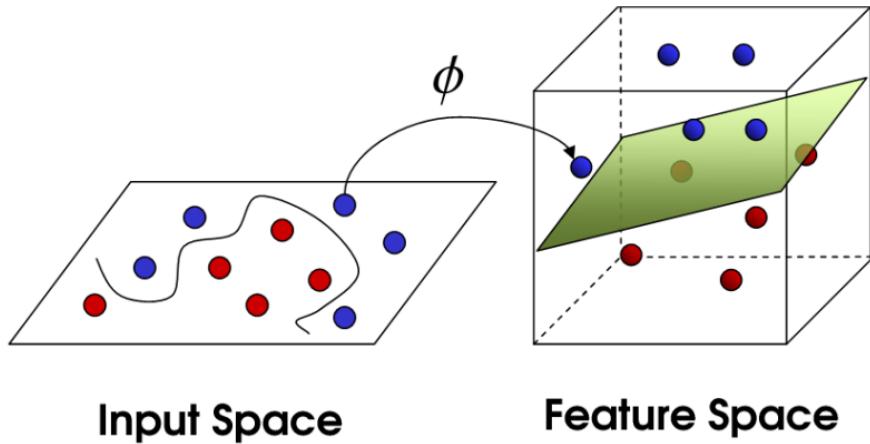
One potential drawback of kernel methods: **non-parametric**, need to potentially keep all the training points.

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

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

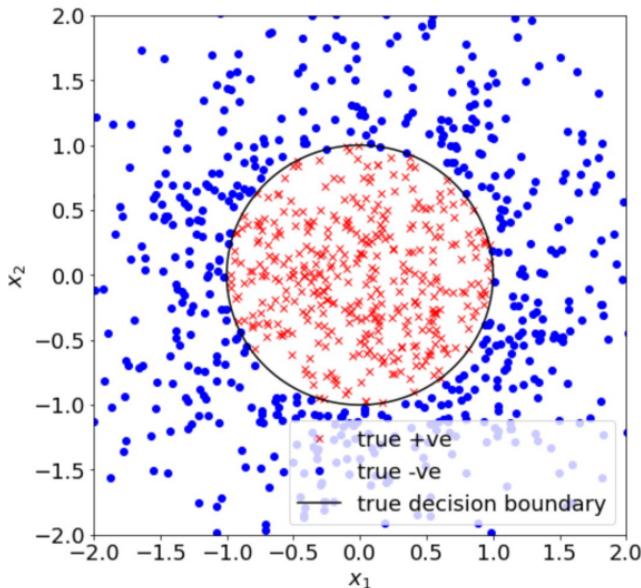


3.2 Examining the effect of kernels



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,

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

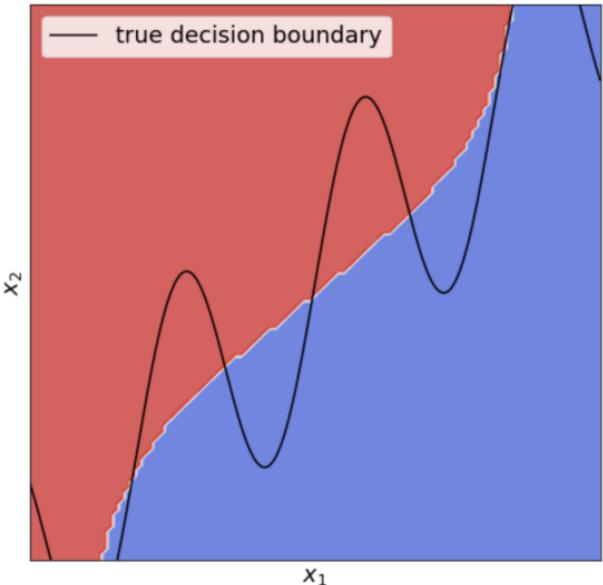
for some $\gamma > 0$.

What does the decision boundary look like?

What is the effect of γ ?

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 ℓ_2 regularized hinge loss):

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

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

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