## CSCI 567: Machine Learning

Vatsal Sharan Spring 2024

Lecture 5, February 9



#### Administrivia

- HW2 out, due in less than 2 weeks.
- Exam 1 in 3 weeks, more details next week

# Recap

#### Regularized least squares

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

$$egin{aligned} oldsymbol{w}^* &= rgmin_{oldsymbol{w}} F(oldsymbol{w}) \ &= rgmin_{oldsymbol{w}} \left( \| oldsymbol{\Phi} oldsymbol{w} - oldsymbol{y} \|_2^2 + \lambda \| oldsymbol{w} \|_2^2 
ight) \ &= \left( oldsymbol{\Phi}^T oldsymbol{\Phi} + \lambda oldsymbol{I} 
ight)^{-1} oldsymbol{\Phi}^T oldsymbol{y} \end{array} egin{align*} oldsymbol{\Phi} = \left( egin{align*} oldsymbol{\phi}(oldsymbol{x}_1)^T \\ \vdots \\ oldsymbol{\phi}(oldsymbol{x}_2)^T \\ \vdots \\ oldsymbol{\phi}(oldsymbol{x}_n)^T \end{array} 
ight), \quad oldsymbol{y} = \left( egin{align*} y_1 \\ y_2 \\ \vdots \\ y_n \end{array} 
ight) \ &= \left( oldsymbol{\Phi}^T oldsymbol{\Phi} + \lambda oldsymbol{I} \right)^{-1} oldsymbol{\Phi}^T oldsymbol{y} \end{aligned}$$

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,

$$oldsymbol{w}^* = oldsymbol{\Phi}^{\mathrm{T}} oldsymbol{lpha} = \sum_{i=1}^n lpha_i oldsymbol{\phi}(oldsymbol{x}_i)$$

Thus the least square solution is a linear combination of features of the datapoints! We calculated what  $\alpha$  should be,

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

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

#### **Kernel trick**

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

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

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  $\phi(\mathbf{x})^T \phi(\mathbf{x}')$ .

#### The kernel trick: Example 1

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

$$m{\phi}(m{x}) = \left(egin{array}{c} x_1^2 \ \sqrt{2}x_1x_2 \ x_2^2 \end{array}
ight)$$

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

$$\phi(\mathbf{x})^{\mathsf{T}}\phi(\mathbf{x}') = x_1^2 x_1'^2 + 2x_1 x_2 x_1' x_2' + x_2^2 x_2'^2$$
$$= (x_1 x_1' + x_2 x_2')^2 = (\mathbf{x}^{\mathsf{T}} \mathbf{x}')^2$$

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

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

#### **Popular kernels:**

1. Polynomial kernel

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

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

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

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

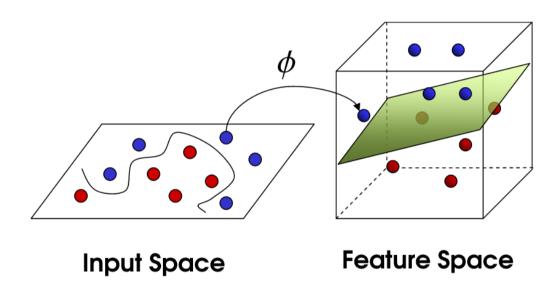
#### **Prediction with kernels**

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

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

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  $w^*$  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  $sign(\boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}))$ . Data may become linearly separable in the feature space!

We'll see this today.

# Support vector machines (SVMs)

#### 1.1 Why study SVM?

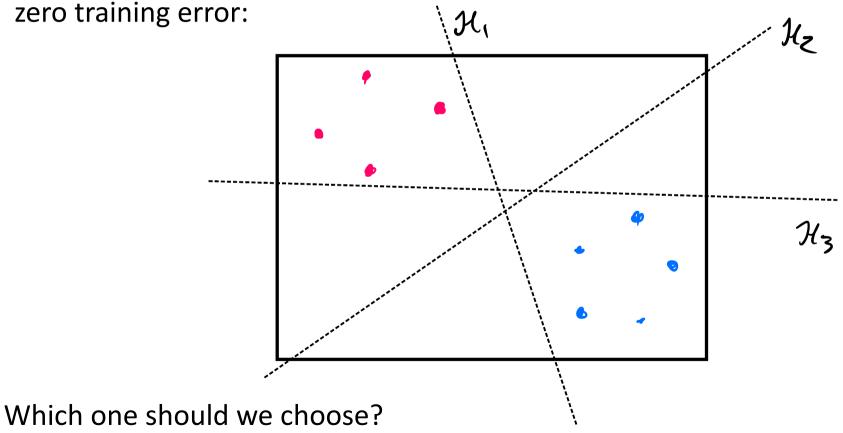
- One of the most commonly used classification algorithms
- Allows us to explore the concept of *margins* in classification
- Works well with the kernel trick
- Strong theoretical guarantees

We focus on **binary classification** here.

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

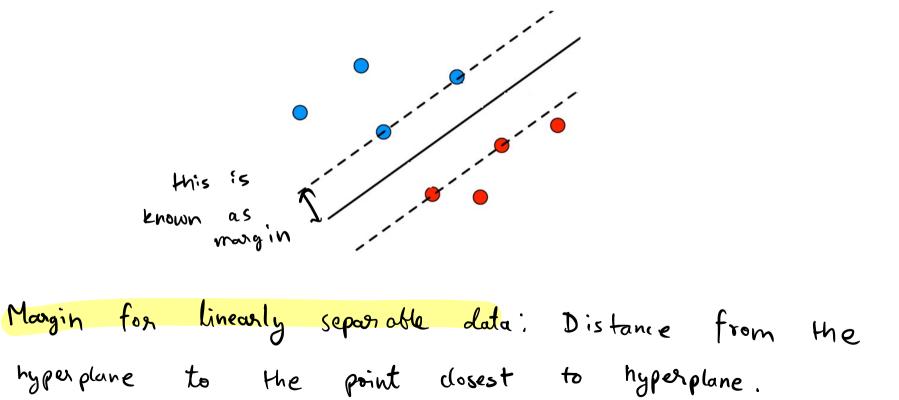
#### 1.2 Margins: separable case, geometric intuition

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



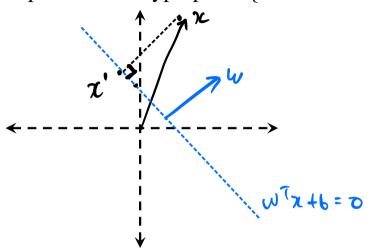
#### 1.2 Margins: separable case, geometric intuition

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



#### 1.2 Formalizing geometric intuition: Distance to hyperplane

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



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

$$0 = \boldsymbol{w}^{\mathrm{T}} \left( \boldsymbol{x} - \beta \frac{\boldsymbol{w}}{\|\boldsymbol{w}\|_{2}} \right) + b = \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x} - \beta \|\boldsymbol{w}\| + b \implies \beta = \frac{\boldsymbol{w}^{\mathrm{T}} \boldsymbol{x} + b}{\|\boldsymbol{w}\|_{2}}.$$

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

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

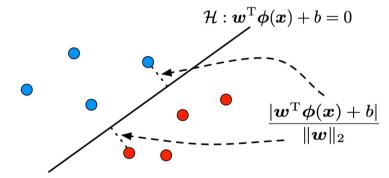
#### 1.2 Margins: functional motivation

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

#### 1.3 Maximizing margin

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

MARGIN OF 
$$(\boldsymbol{w},b) = \min_i \frac{y_i(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b)}{\|\boldsymbol{w}\|_2}$$
 dot  $\{\boldsymbol{\psi}_i, \boldsymbol{\psi}_i\}$ 



The intuition "the further away the better" translates to solving

$$\max_{\boldsymbol{w},b} \min_{i} \frac{y_i(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b)}{\|\boldsymbol{w}\|_2} = \max_{\boldsymbol{w},b} \frac{1}{\|\boldsymbol{w}\|_2} \min_{i} y_i(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b)$$

#### 1.3 Maximizing margin, rescaling

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

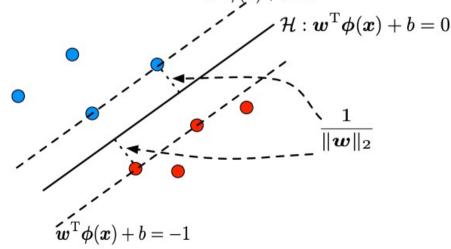
Decision boundary: 
$$w^T \phi(x) + b = 0 \iff (w^T \phi(x) + b) = 0$$

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

The margin then becomes

MARGIN OF 
$$(\boldsymbol{w}, b)$$

$$= \frac{1}{\|\boldsymbol{w}\|_2} \min_{i} y_i(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b)$$
$$= \frac{1}{\|\boldsymbol{w}\|_2}$$



#### 1.4 SVM for separable data: "Primal" formulation

For a separable training set, we aim to solve

$$\max_{m{w},b} rac{1}{\|m{w}\|_2}$$
 s.t.  $\min_i y_i(m{w}^{\mathrm{T}}m{\phi}(m{x}_i) + b) = 1$ 

This is equivalent to

$$\min_{m{w},b} \ rac{1}{2} \|m{w}\|_2^2$$
 This is convariant. Set.  $y_i(m{w}^{ ext{T}}m{\phi}(m{x}_i)+b) \geq 1, \ orall \ i \in [n]$ 

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

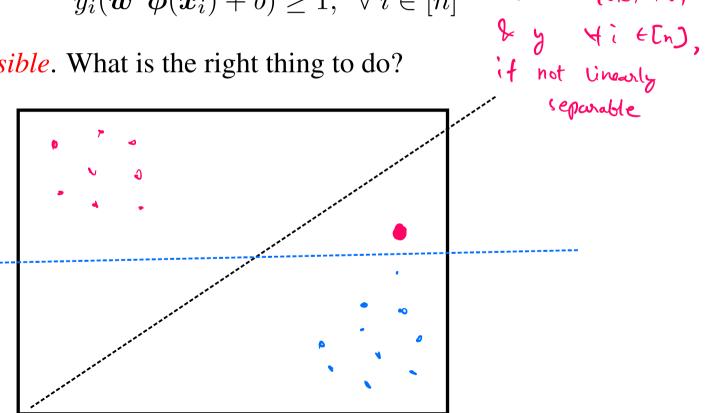
#### 1.5 General non-separable case

If data is not linearly separable, the previous constraint

arable, the previous constraint 
$$y_i(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_i) + b) \geq 1, \ \ \forall \ i \in [n]$$

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

even if data is separable, separable, separable it?



#### 1.5 General non-separable case

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

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

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

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

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



Aside: Why  $\ell_1$  penalization?

hinge loss 
$$l(z) = ma_{+}(0, 1-z)$$

Squared hinge loss  $l(z) = (ma_{+}(0, 1-z))^{2}$ 

The grows much faster than  $\pi$ !

Squared hinge loss would pendize where predictions much more.

#### 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

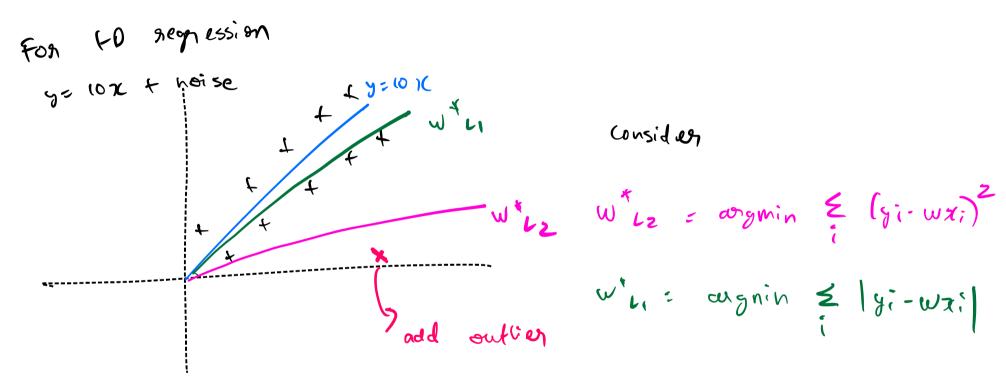
If I have  $x_1, x_2, \ldots, x_n$ :

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

Median is more robust to outliers than mean.

### Aside: Why $\ell_1$ penalization?



#### 1.5 Back to SVM: General non-separable case

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

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

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

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

#### 1.5 SVM General Primal Formulation

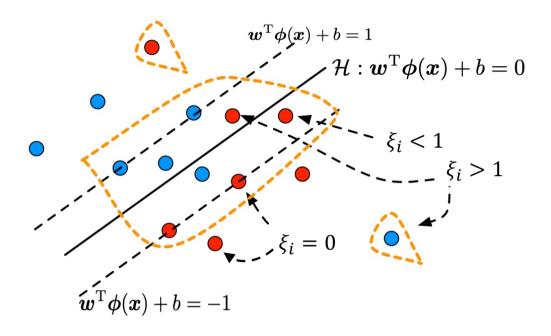
We want  $\xi_i$  to be as small as possible. The objective becomes

$$\min_{\boldsymbol{w},b,\{\boldsymbol{\xi_i}\}} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2 + C \sum_{i} \boldsymbol{\xi_i}$$
s.t.  $y_i(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b) \ge 1 - \boldsymbol{\xi_i}, \ \forall i \in [n]$ 

$$\boldsymbol{\xi_i} \ge 0, \ \forall i \in [n]$$

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

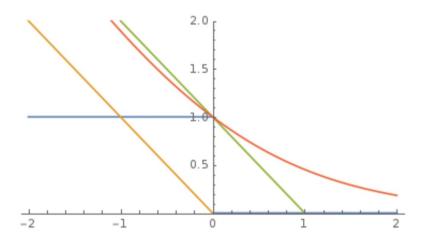
#### 1.6 Understanding the slack conditions



- 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  $\ell_2$  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  $\ell_{\text{hinge}}(z) = \max\{0, 1-z\} \rightarrow \mathbf{SVM}$

#### 1.7 Primal formulation: Another view

For a linear model (w, b), this means

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

- 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?

#### 1.7 Equivalent forms

#### The formulation

$$\min_{\boldsymbol{w},b,\{\xi_i\}} C \sum_i \xi_i + \frac{1}{2} \|\boldsymbol{w}\|_2^2$$
 s.t.  $1 - y_i(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b) \leq \xi_i, \quad \forall \, i \in [n]$  
$$\xi_i \geq 0, \quad \forall \, i \in [n]$$
 In order to minimize  $\boldsymbol{\xi}_i$  as found set  $\boldsymbol{\xi}_i$ , to be as small as possible is equivalent to

$$\min_{\boldsymbol{w},b,\{\xi_i\}} C \sum_i \xi_i + \frac{1}{2} \|\boldsymbol{w}\|_2^2$$
s.t. 
$$\max \left\{ 0, 1 - y_i(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b) \right\} = \xi_i, \ \forall i \in [n]$$

#### 1.7 Equivalent forms

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

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

#### is equivalent to

$$\min_{\boldsymbol{w},b} C \sum_{i} \max \left\{ 0, 1 - y_i(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b) \right\} + \frac{1}{2} \|\boldsymbol{w}\|_2^2$$

and

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

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

#### 1.8 Optimization

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

- 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

Recall the SVM for separable case: min I lolle w, b s.t. 4: (w) \$(xi) +1) 71 #i & Cn) (an we use kernel trick? (an we show that we is a linear combination of \$ (xi)?

#### How did we show this for regularized least squares?

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

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

we know

$$oldsymbol{w}^* = rac{1}{\lambda} oldsymbol{\Phi}^{ ext{T}}(oldsymbol{y} - oldsymbol{\Phi} oldsymbol{w}^*) = oldsymbol{\Phi}^{ ext{T}} oldsymbol{lpha} = \sum_{i=1}^n lpha_i oldsymbol{\phi}(oldsymbol{x}_i)$$

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

#### 2.1 Kernelizing SVM

$$\frac{\partial F(\omega)}{\partial \omega} = \frac{2}{i\epsilon!} \left( \frac{\partial \ln \log(2)}{\partial z} \Big|_{z=y_i(\omega^T \phi(x_i) + \delta)} \frac{\partial z}{\partial \omega} + \lambda \omega \right)$$

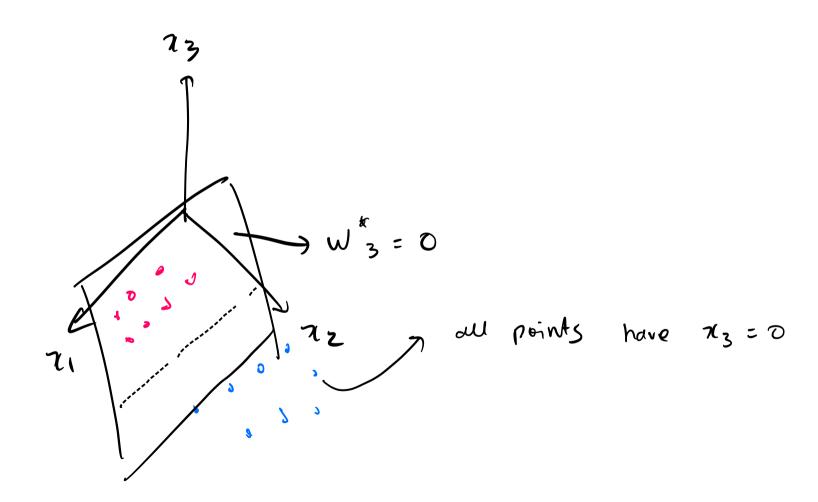
$$\frac{\partial F(\omega)}{\partial \omega} = \frac{2}{i\epsilon!} \left( -1(2\epsilon) \Big|_{z=y_i(\omega^T \phi(x_i) + \delta)} (y_i \phi(x_i) + \delta) \right) + \lambda \omega$$

$$w^{(\delta)} = 0$$

$$w^{(t+1)} = w^{(t)} \eta \left( \sum_{i=1}^{m} -2! \left( y_i w^{T} \phi(x_i) + i \leq l \right) |y_i \phi(x_i) \rangle + \lambda w^{(t)} \right)$$

-. W always lies in span of  $\phi(x_i)$ .  $w^{(t)} = \sum_{i} u^{(t)} y_i \phi(x_i) + t$ , for some  $d_i^{(t)}$   $w^{(t)} = \sum_{i} d_i^{(t)} y_i \phi(x_i) + t$ , for some  $d_i^{(t)}$ 

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



If w= Zdi yi 
$$\phi(\chi_i)$$
,

$$\min_{\boldsymbol{w},b} \frac{1}{2} \|\boldsymbol{w}\|_{2}^{2}$$
s.t.  $y_{j}(\boldsymbol{w}^{T}\phi(\boldsymbol{x}_{j}) + b) \geq 1, \ \forall j \in [n].$ 

$$S.t. \ y_{j}(\boldsymbol{w}^{T}\phi(\boldsymbol{x}_{j}) + b) \geq 1, \ \forall j \in [n].$$

$$S.t. \ y_{j}(\boldsymbol{z}^{T}\phi(\boldsymbol{x}_{j}) + b) \geq 1, \ \forall j \in [n].$$

This is equivalent to

rain

L Z Z di di yi yi yi 
$$\phi(x_i)^{T}\phi(x_i)$$

S.t. yi (\geq aiyi \alpha(x\_i)^{T}\alpha(x\_i) \di) >1. \Hi \in Inj

### 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,

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

### 2.2 SVM: Dual form for separable case

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

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

No need to compute  $\phi(x)$ . 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:

$$\min_{\boldsymbol{w},b,\{\xi_i\}} C \sum_{i} \xi_i + \frac{1}{2} \|\boldsymbol{w}\|_2^2$$
s.t.  $y_i(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}_i) + b) \ge 1 - \xi_i, \quad \forall i \in [n]$ 

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

The dual is very similar,

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

### 2.4 Prediction using SVM

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

Remember that,

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

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

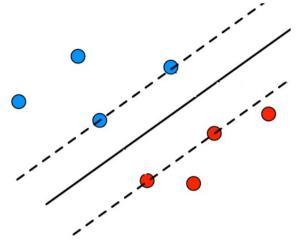
To make a prediction on any datapoint x,

$$\operatorname{sign}\left(\boldsymbol{w}^{*\mathsf{T}}\phi(\boldsymbol{x}) + b^{*}\right) = \operatorname{sign}\left(\sum_{i:\alpha_{i}^{*}>0} \alpha_{i}^{*}y_{i}\phi(\boldsymbol{x}_{i})^{\mathsf{T}}\phi(\boldsymbol{x}) + b^{*}\right)$$
$$= \operatorname{sign}\left(\sum_{i:\alpha_{i}^{*}>0} \alpha_{i}^{*}y_{i}k(\boldsymbol{x}_{i},\boldsymbol{x}) + b^{*}\right).$$

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.

$$y_{i} \left( w^{*} \right) \left( \psi(x_{i}) + \psi(x_{i}) \right) = 1$$
 =>  $y_{i}^{2} \left( w^{*} \right) \left( \psi(x_{i}) + \psi(x_{i}) \right) = y_{i}^{2}$   
=>  $w^{*} \left( \psi(x_{i}) \right) + \psi(x_{i}) +$ 

### 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  $y_i(\mathbf{w}^{*T}\phi(\mathbf{x}_i) + b^*) = 1$  (i.e. that support vector lies on the margin). Therefore, as before,

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

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

With  $\alpha^*$  and  $b^*$  in hand, we can make a prediction on any datapoint x,

$$\operatorname{sign}\left(\boldsymbol{w}^{*T}\phi(\boldsymbol{x}) + b^{*}\right) = \operatorname{sign}\left(\sum_{i:\alpha_{i}^{*}>0} \alpha_{i}^{*}y_{i}k(\boldsymbol{x}_{i},\boldsymbol{x}) + b^{*}\right).$$

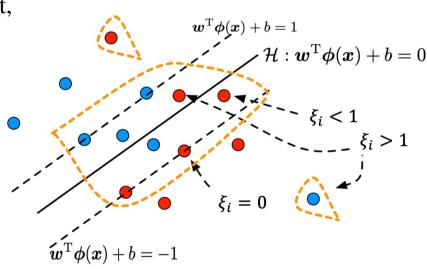
# • SVMs: Understanding them further

# 3.1 Understanding support vectors

Support vectors are  $\phi(x_i)$  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 contraint,
- (2) they do not satisfy the large margin contraint,
- (3) they are misclassified.
- when  $\xi_i^* = 0$ ,  $y_i(\boldsymbol{w}^{*T}\boldsymbol{\phi}(\boldsymbol{x}_i) + b^*) = 1$ , and thus the point is  $1/\|\boldsymbol{w}^*\|_2$  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.



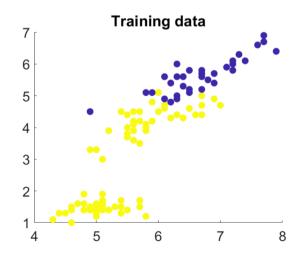
Support vectors (circled with the orange line) are the only points that matter!

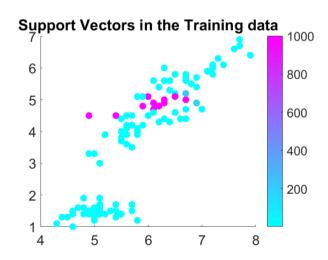
### 3.1 Understanding support vectors

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

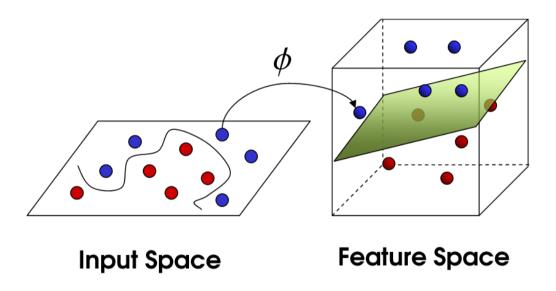
$$\operatorname{sign}\left(\boldsymbol{w}^{*T}\phi(\boldsymbol{x}) + b^{*}\right) = \operatorname{sign}\left(\sum_{i=1}^{n} \alpha_{i}^{*}y_{i}k(\boldsymbol{x}_{i}, \boldsymbol{x}) + b^{*}\right).$$

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



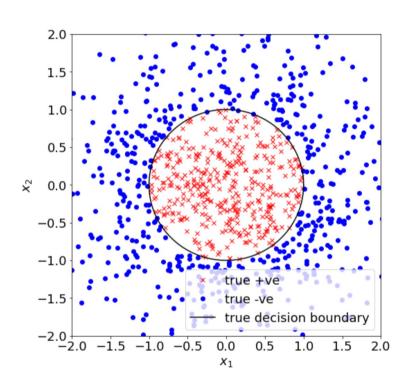


### 3.2 Examining the effect of kernels



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

# Polynomial kernel: example



# Gaussian kernel: example

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

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

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

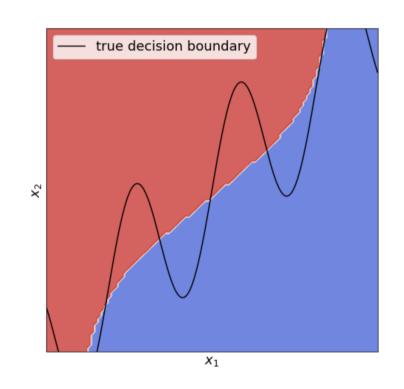
$$k(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\gamma \|\boldsymbol{x} - \boldsymbol{x}'\|_{2}^{2}\right)$$

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

$$\operatorname{sign}\left(\boldsymbol{w}^{*\mathsf{T}}\phi(\boldsymbol{x})+b^{*}\right)=\operatorname{sign}\left(\sum_{i:\alpha^{*}>0}\alpha_{i}^{*}y_{i}k(\boldsymbol{x}_{i},\boldsymbol{x})+b^{*}\right).$$



Switch to Colab

### **SVM: Summary of mathematical forms**

**SVM:** max-margin linear classifier

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

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

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

$$\max_{\{\alpha_i\}} \quad \sum_{i} \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j \phi(\boldsymbol{x}_i)^{\mathsf{T}} \phi(\boldsymbol{x}_j)$$
s.t. 
$$\sum_{i} \alpha_i y_i = 0 \quad \text{and} \quad 0 \le \alpha_i \le C, \quad \forall i \in [n].$$