# CIS 5200: MACHINE LEARNING KERNELS



Content here draws from material by Jake/Shivani (UPenn), Christopher De Sa (Cornell)





# Surbhi Goel

#### Spring 2023

### OUTLINE - TODAY

Recap of SVMs
Function Maps
Kernel Functions
Kernelization
Demo!

### SVM - SOFT-MARGIN

#### **Primal:**

min w,b

such that

 $\frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^{m} \xi_i$ i=1 $\xi_i \geq 0, \forall i \in [m]$ 

#### **Dual:**



such that

max

α



# SOFT-SVM - LOSS MINIMIZATION VIEW



Is equivalent to the following loss minimization problem for  $C = \frac{1}{2\lambda m}$ :

$$\min_{w,b} \frac{1}{m} \sum_{i=1}^{m} \max(0, 1 - y_i(w^{\mathsf{T}}x_i + b))$$

 $\ell_2$ -regularized hinge loss minimization

$$\|_{2}^{2} + C \sum_{i=1}^{m} \xi_{i}$$

$$x_{i} + b) \geq 1 - \xi_{i}, \forall i \in [m]$$

$$y_{i} \in [m]$$

 $+ \lambda \|w\|^2$ 

n



4

### NON-SEPARABLE



#### What can we do if data is like this?

# FEATURE MAP - MAPTO HIGHER DIMENSIONS

Map data into to a higher dimensional space using feature map  $\phi$  $x \mapsto \phi(x)$ 

What features should we use?



## FEATURE MAP - MAPTO HIGHER DIMENSIONS

Consider the following feature map:

 $\phi(x) = \begin{vmatrix} x_2 \\ x_1^2 \end{vmatrix}$  $x_2^2$ 

Is the data linearly separable in this feature space?

- Let  $w = [0,0,1,1]^T$  and  $b = r^2$ , then we have
  - $w^{\mathsf{T}}\phi(x) + b = x_1^2 + x_2^2 r^2$

7



### FEATURE MAP - LINEAR TO NON-LINEAR



**Predictor function:** Linear functions  $w^{\top}\phi(x) + b$ 



**Training data:**  $\{(x_1, y_1), \dots, (x_m, y_m)\} \mapsto \{(\phi(x_1), y_1), \dots, (\phi(x_m), y_m)\}$ 

### FEATURE MAP - CHALLENGE

Consider the following feature map:

 $\phi(x) = \begin{array}{c} \vdots \\ x_d \\ x_1^2 \\ x_1 x \end{array}$ 

1 ' x<sub>1</sub>

 $x_1 x_2$ 

 $\mathcal{X}_{d}$ 



# What is the dimension of this map? $D = (d+1)^2$

What if we take all monomials to degree r?

$$D = (d+1)^r$$

This is huge!



### RECALL - SOFT-SVM

max α

such that



We only need to compute inner products  $\phi(x_i)^{\top}\phi(x_i)$ 

#### With feature map:





### KERNELTRICK - EXAMPLE

 $\phi(x) = \begin{vmatrix} 1 \\ x_1 \\ \vdots \\ x_d \\ x_1^2 \\ x_1^2 \\ x_1 x_2 \end{vmatrix}$  $x_1 x_2$  $\frac{1}{x_d^2}$ 



#### Let's compute inner product:

### $\phi(x)^{\top}\phi(x') = 1 + x^{\top}x' + (x^{\top}x')^2$

What is the computational cost of this?

# KERNEL FUNCTIONS

A kernel is a function  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  that satisfies: for some feature map  $\phi$  that maps  ${\mathcal X}$  to some inner product space  ${\mathcal V}$ . For data  $x_1, \ldots, x_m$  the kernel matrix  $K \in \mathbb{R}^{m \times m}$ 

A kernel k is valid if for any  $x_1, \ldots, x_m$ : K is symmetric and positive semi-definite

# $k(x, x') = \langle \phi(x), \phi(x') \rangle$

# $K_{ii} = k(x_i, x_j)$

 $K = K^{\top} \qquad \qquad \text{For any } x, x^{\top} K x \ge 0$ All eigenvalues are non-negative



#### KERNELS - EXAMPLES

**\*** Linear:

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

# **\* Polynomial**: for degree *r* $k(x, x') = (1 + x^{T}x')^{r}$

\* Gaussian/Random Basis Function (RBF): for some parameter  $\sigma > 0$  $k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$ 

### LETS KERNELIZE!

- training points,  $w = \sum \alpha_i x_i$ i=1
- Replace  $x_i^T x_i \rightarrow k(x_i, x_i)$  everywhere



#### Show that the solution to your problem lies in the span of the

# Rewrite the algorithm and the predictor so that all training or test points are only accessed in inner-products $(x_i^{\mathsf{T}}x_i)$ with other points

#### EXAMPLE - SOFT-SVM



### EXAMPLE - KERNEL SVM





$$gn\left(\sum_{i=1}^{m} \alpha_i y_i k(x_i, x) + b\right)$$

### EXAMPLE - PERCEPTRON

#### Algorithm 1: Perceptron

Initialize  $w_1 = 0 \in \mathbb{R}^d$ for t = 1, 2, ... do if  $\exists i \in [m] \ s.t. \ y_i \neq \text{sign} (w_t^\top x \text{else output } w_t$ end

Update is always in the feature space

Can we write the algorithm in terms of  $\alpha$ ?

#### if $\exists i \in [m] \ s.t. \ y_i \neq \text{sign}\left(w_t^\top x_i\right)$ then update $w_{t+1} = w_t + y_i x_i$

e, so 
$$w_* = \sum_{i=1}^m \alpha_i x_i$$
 for some  $\alpha \in \mathbb{R}^m$ 

### EXAMPLE - KERNEL PERCEPTRON

Algorithm 2: Perceptron - Dual Initialize  $\alpha_1 = 0 \in \mathbb{R}^d$ for t = 1, 2, ... do if  $\exists i \in \{1, \ldots, m\}$  s.t.  $y_i \neq \operatorname{sign}\left(\sum_{j=1}^m \sum_{j=1}^m \sum_{$ else output  $\alpha_t$  $\mathbf{end}$ 

Now we can kernelize this since it only depends on inner products!

$$=_1 \alpha_{tj} x_j^{\top} x_i$$
 then update  $\alpha_{(t+1)i} = \alpha_{ti} + y_i$ 



### EXAMPLE - RIDGE REGRESSION





$$+ \lambda \|w\|_{2}^{2} = \frac{1}{m} \|Y - Xw\|^{2} + \lambda \|w\|_{2}^{2}$$

Can w be expressed as a linear combination of the input datapoints? Proof by contradiction!

We have  $w = \sum \alpha_i x_i = X^T \alpha$  for some  $\alpha$ 

i=1

#### EXAMPLE - RIDGE REGRESSION



# Each element of $XX^{\dagger}$ is an inner product $x_i^{\dagger}x_j$ for some $i, j \in [m]$

Prediction is  $w^{T}x$ 

$$x = \sum_{i=1}^{m} \alpha_i x_i^{\mathsf{T}} x = \alpha^{\mathsf{T}} X x$$

### EXAMPLE - KERNEL RIDGE REGRESSION



Prediction is  $w^{\top}\phi(x) =$ 

 $k_x = [k(x,$ 



#### Here $K_{ij} = k(x_i, x_j)$ is the kernel/gram matrix

$$= \sum_{i=1}^{m} \alpha_i k(x_i, x) = \alpha^{\mathsf{T}} k_x \text{ where}$$
$$= x_1) \dots k(x, x_m) ]^{\mathsf{T}}$$



DEMO

By folks at Cornell CS

### POWER OF KERNELS

- training points,  $w = \sum \alpha_i x_i$ i=1
- Replace  $x_i^T x_i \rightarrow k(x_i, x_i)$  everywhere for a valid kernel k



#### Show that the solution to your problem lies in the span of the

There is a general theorem called the Representer Theorem which tells us when this is true

Rewrite the algorithm and the predictor so that all training or test points are only accessed in inner-products  $(x_i^{\mathsf{T}}x_i)$  with other points

Super Powerful!

### CHALLENGE

# \* How do we choose a good feature map $\phi$ ? \* Feature map is the same for all inputs!

#### Can learn the feature map itself $\rightarrow$ deep learning!