

# 446 Section $\phi(7)$

TA: Yufei Zhang

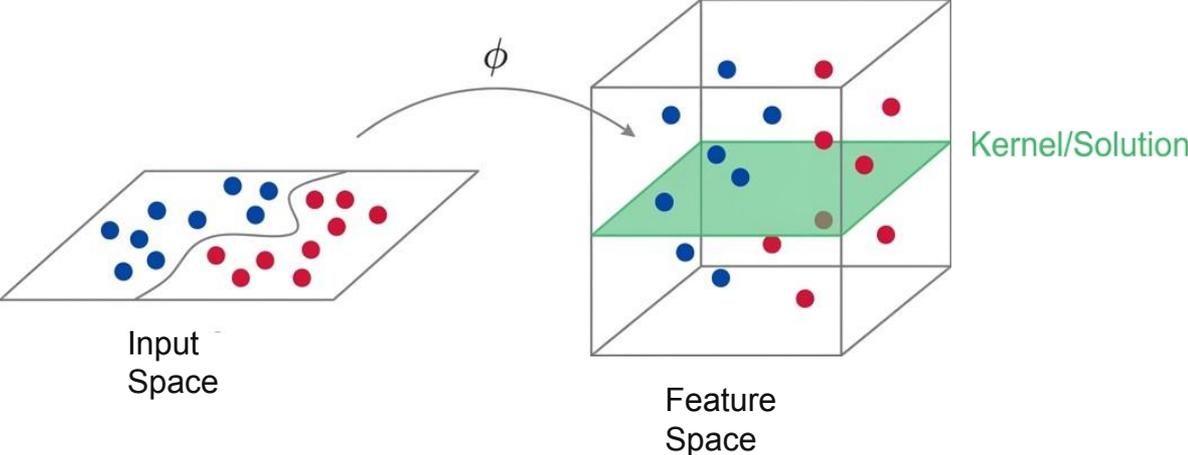
# Plans for today!

1. This
2. Reminders
3. Kernels
4. Problem 2a, 2b
5. PyTorch Colab Notebook (No Demo)

# Reminders

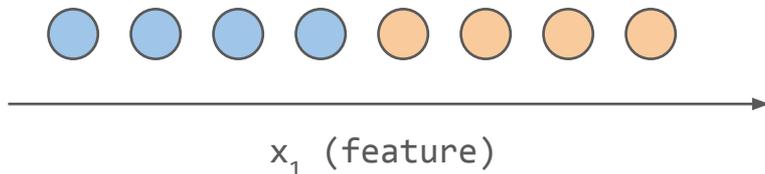
- HW3 due **Wed, Feb 25 @ 11:59pm**
- Midterm grades out!
  - Please only submit a regrade request if you have a *strong* case for extra points

# Kernels



# Motivation

Consider a 1D dataset with just feature  $x_1$ . The blue and orange dots denote two different classes.



Can you draw a horizontal line to separate these classes?

No

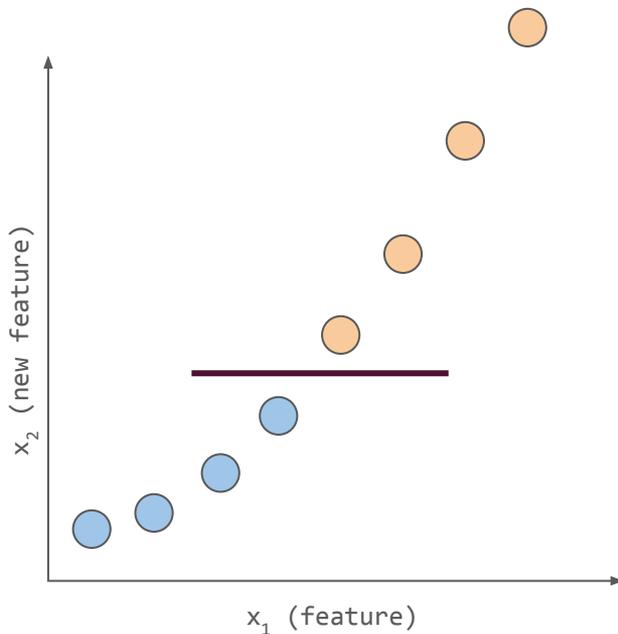
# Kernels

We need to create a new feature  $x_2$  from what we have (which is  $x_1$ )

$$\phi(x) = [\phi_1(x), \phi_2(x)]$$

$$\phi_1(x) = x$$

$$\phi_2(x) = x^2$$



Can you draw a horizontal line to separate these classes?

Yes!

# Problems with efficiency ← THIS IS WHY!

If we did things step by step...

1. Take our  $n$  data points, each one a  $d$  dimensional vector
2. To EACH data point, apply our kernel map ( $\phi(x)$ ), further blowing up the space/time complexity
3. Perform linear regression on the exploded set of features
4. Most likely fail...

Kernel trick!

# Ridge Regression Reminders:

Variables:

$$X \in \mathbb{R}^{n \times d}; y \in \mathbb{R}^n; w \in \mathbb{R}^d$$

Optimization objective:

$$\hat{w} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

Closed-form solution:

$$\hat{w} = (X^\top X + \lambda I)^{-1} X^\top y$$

Our goal is to get from **ridge regression** to **kernelized linear regression**

# Step 1: The Representer Theorem

The optimal weight vector  $\hat{w}$  lies in the span of the data points.

$$\hat{w} = X^T \alpha = \sum_{i=1}^n x_i \alpha_i$$

Where  $\alpha \in \mathbb{R}^n$  is a vector of coefficients.

This fact implies that **we can write  $w$  as a linear combination** of the data points ( $x_i$ ) and some scale factors ( $\alpha_i$ )  $\rightarrow$  useful for our proof later!

Let's do feature expansion!

From now on we are working with the *expanded* feature set:

$$X \rightarrow \phi(X)$$

NEW Optimization objective:

$$\hat{w}_\phi = \operatorname{argmin}_w \|y - \phi(X)w\|_2^2 + \lambda \|w\|_2^2$$

What if our feature set blows up so that  $d \gg n$ ? Or even blows up to infinite dimensions?

## Step 2: Substitution

$$X \rightarrow \phi(X)$$

$$\hat{w}_\phi = \operatorname{argmin}_w \|y - \phi(X)[w]\|_2^2 + \lambda \| [w] \|_2^2$$

$$\text{If } w = \phi(X)^T a$$

$$\hat{\alpha} = \operatorname{argmin}_\alpha \|y - \phi(X)[\phi(X)^T \alpha]\|_2^2 + \lambda \|[\phi(X)^T \alpha]\|_2^2$$

This is why we needed step 1. It makes this substitution possible **no matter the dimensionality** of the “blown up”  $X$ !

## Step 2: Define the Kernel Matrix

The **Kernel Function**:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

(Measures similarity  
between two points)

The **Kernel Matrix**:

$$\mathbf{K} \in \mathbb{R}^{n \times n}$$

where  $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$

How can we construct  
the kernel **matrix**  
using the kernel  
*function*?

## Step 2: Define the Kernel Matrix

The **Kernel Function**:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

(Measures similarity  
between two points)

The **Kernel Matrix**:

$$\mathbf{K} \in \mathbb{R}^{n \times n}$$

where  $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$

$$\mathbf{K} \begin{bmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \cdots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_1) & \cdots & \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_n) \end{bmatrix}$$

The Trick: We can compute  $\mathbf{K}$   
without ever calculating  $\phi$ .

$$K : \begin{matrix} & & & & n \\ & & & & \uparrow \\ n & \left[ \begin{array}{cccc} \phi(x_1)^T \phi(x_1) & \phi(x_2)^T \phi(x_1) & \dots & \phi(x_n)^T \phi(x_1) \\ \phi(x_1)^T \phi(x_2) & & & \\ \vdots & & & \\ \phi(x_1)^T \phi(x_n) & & & \phi(x_n)^T \phi(x_n) \end{array} \right] \end{matrix}$$

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

$$K_{ij} \in \mathbb{R}^{n \times n}, \quad K_{ij} = K(x_i, x_j)$$

The main takeaway is the formulation of the kernel matrix below.

$$\begin{matrix} & m \\ n & \left[ \begin{array}{c} \phi(x_1) \\ \vdots \\ \phi(x_n) \end{array} \right] \end{matrix} \quad \Downarrow \quad \begin{matrix} & n \\ m & \left[ \begin{array}{ccc} | & & | \\ \phi(x_1)^T & \dots & \phi(x_n)^T \\ | & & | \end{array} \right] \end{matrix}$$

$$K = \phi(x) \phi(x)^T : x \in \mathbb{R}^{n \times m}$$

## Substitute into Linear Regression

$$\text{If } \mathbf{K} = \phi(X)\phi(X)^\top$$

$$\begin{aligned}\hat{a} &= \underset{a}{\operatorname{argmin}} \|y - \phi(X)\phi(X)^\top a\|_2^2 + \lambda \|\phi(X)^\top a\|_2^2 \\ &= \underset{a}{\operatorname{argmin}} \|y - \phi(X)\phi(X)^\top a\|_2^2 + \lambda a^\top \phi(X)\phi(X)^\top a \\ &= \underset{a}{\operatorname{argmin}} \|y - \mathbf{K}a\|_2^2 + \lambda a^\top \mathbf{K}a\end{aligned}$$

**The dimension  $d$  has completely vanished.  
The problem now depends only on  $n$ .**

Now you are gonna take this all the way to the end!

# Step 3: Derive

## 2a. Kernelized Linear Regression

Recall that the definition of a kernel is the following:

**Definition 1.** A function  $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a *kernel* for a map  $\phi$  if  $K(x, x') = \phi(x) \cdot \phi(x') = \langle \phi(x), \phi(x') \rangle$  for all  $x, x'$ .

Consider regularized linear regression (without a bias, for simplicity). Our objective to find the optimal parameters  $\hat{w} = \arg \min_w L(w)$  for a dataset  $(x_i, y_i)_{i=1}^n$  that minimize the following loss function:

$$L(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$$

Note that from class, we know there is an optimal  $\hat{w}$  that lies in the span of the datapoints. Concretely, there exist  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$  such that  $\hat{w} = \sum_i \alpha_i x_i$ . Also recall from lecture that the expression of our loss function  $L(w)$  in terms of the kernel is:

$$L(w) = \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^T \mathbf{K}\alpha$$

Solve for the optimal  $\alpha$

## 2a. Kernelized Linear Regression

$$L(w) = \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda\alpha^T \mathbf{K}\alpha$$

1. **Gradient:** Take derivative of  $L(\alpha)$  w.r.t  $\alpha$ :

$$\nabla_{\alpha} L = -2\mathbf{K}(\mathbf{y} - \mathbf{K}\alpha) + 2\lambda\mathbf{K}\alpha$$

2. **First-Order Optimality:** Set to 0 and factor out  $K$ :

$$\mathbf{K}(\mathbf{K}\alpha - \mathbf{y} + \lambda\alpha) = 0$$

$$\mathbf{K}((\mathbf{K} + \lambda\mathbf{I})\alpha - \mathbf{y}) = 0$$

3. **Invert:** Assuming  $K$  is positive semi-definite (PSD):

$$(\mathbf{K} + \lambda\mathbf{I})\alpha = \mathbf{y}$$

4. **Solution:**

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

$$\nabla_{\alpha} L(w) = 0$$

$$-2\mathbf{K}(\mathbf{y} - \mathbf{K}\alpha) + 2\lambda\mathbf{K}\alpha = 0$$

$$-\mathbf{K}(\mathbf{y} - \mathbf{K}\alpha) + \lambda\mathbf{K}\alpha = 0$$

$$\mathbf{K}(\mathbf{K}\alpha - \mathbf{y} + \lambda\alpha) = 0$$

$$\mathbf{K}((\mathbf{K} + \lambda\mathbf{I})\alpha - \mathbf{y}) = 0$$

$$\mathbf{K}(\mathbf{K} + \lambda\mathbf{I})\alpha = \mathbf{K}\mathbf{y}$$

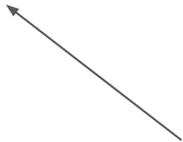
$$\hat{\alpha} = (\mathbf{K} + \lambda\mathbf{I})^{-1} \mathbf{y}$$

## Quick high-level recap

A kernel  $K(a, b)$  takes in  $d$  dimensional vectors  $a, b$  and gives us a scalar.

When you construct a kernel such that  $K(a, b) = \phi(a)^\top \phi(b)$

We can ultimately use  $K$  to efficiently apply  $\phi$  to our features.



Can have infinite  
dimensions

## STEP 4: PREDICT

Normal case:  $\hat{Y} = \hat{W} \cdot Z \quad : \quad Z \in \mathbb{R}^d$

Kernel case:  $\hat{Y} = \hat{W}_\phi \cdot \phi(z) \quad \hat{W}_\phi = \phi(x)^T a$

$$= \phi(x)^T a \phi(z)$$

$$= \sum_{i=1}^n a_i \phi(x_i)^T \phi(z)$$

$$= \sum_{i=1}^n a_i K(x_i, z)$$

Let's vectorize this!  $\longrightarrow$

You need some extra steps to predict a new datapoint, **but the predictive power compared to computational cost is well worth it!**

## 2b. Kernelized Linear Regression

$$L(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$$

$$L(w) = \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^T \mathbf{K}\alpha$$

Let us assume that we were using a linear kernel where  $\mathbf{K}_{ij} = x_i^T x_j$ . Suppose we have  $\mathbf{X}_{\text{test}}$  that we want to make prediction for after training on  $\mathbf{X}_{\text{train}}$ . Express the estimate  $\hat{\mathbf{Y}}$  in terms of  $\mathbf{K}_{\text{train}} = \mathbf{X}_{\text{train}} \mathbf{X}_{\text{train}}^T$ ,  $\mathbf{y}_{\text{train}}$ ,  $\mathbf{X}_{\text{train}}$  and  $\mathbf{X}_{\text{test}}$ . What would the general prediction formula look like if we are not using a linear kernel? Express the solution in terms of  $\mathbf{K}_{\text{train, test}}$  **Solution:**

## 2b. Kernelized Linear Regression

$$L(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$$

$$L(w) = \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^T \mathbf{K}\alpha$$

Let us assume that we were using a linear kernel where  $\mathbf{K}_{ij} = x_i^T x_j$ . Suppose we have  $\mathbf{X}_{\text{test}}$  that we want to make prediction for after training on  $\mathbf{X}_{\text{train}}$ . Express the estimate  $\hat{\mathbf{Y}}$  in terms of  $\mathbf{K}_{\text{train}} = \mathbf{X}_{\text{train}} \mathbf{X}_{\text{train}}^T$ ,  $\mathbf{y}_{\text{train}}$ ,  $\mathbf{X}_{\text{train}}$  and  $\mathbf{X}_{\text{test}}$ . What would the general prediction formula look like if we are not using a linear kernel? Express the solution in terms of  $\mathbf{K}_{\text{train, test}}$  **Solution:**

Express the estimate  $\hat{\mathbf{Y}}$ :

$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{X}_{\text{test}} \hat{\mathbf{w}} \\ &= \mathbf{X}_{\text{test}} \mathbf{X}_{\text{train}}^T \hat{\boldsymbol{\alpha}} \\ &= \mathbf{X}_{\text{test}} \mathbf{X}_{\text{train}}^T (\mathbf{K}_{\text{train}} + \lambda \mathbf{I})^{-1} \mathbf{y}_{\text{train}}\end{aligned}$$

General formula:

$$\begin{aligned}&= \mathbf{X}_{\text{test}} \mathbf{X}_{\text{train}}^T \hat{\boldsymbol{\alpha}} \\ &= \mathbf{K}_{\text{test, train}} \hat{\boldsymbol{\alpha}} \\ &\quad (m \times n)\end{aligned}$$

# Takeaways

1. **Efficiency:** Kernels allow us to operate in infinite dimensions without the computational cost ( $n$  vs  $d$ )
2. **Substitution:** We replace explicit features  $\phi(x)$  with pairwise dot products  $K$
3. **Result:** A regression model solved using only the data count  $n$

# Mercer's Conditions

Not just any function can be a Kernel. To ensure the optimization problem is convex and solvable, the Kernel Matrix  $\mathbf{K}$  must satisfy:

✓ **Symmetric:**

$$K_{ij} = K_{ji}$$

(The similarity between A and B is the same as B and A).

✓ **Positive Semi-Definite (PSD):**

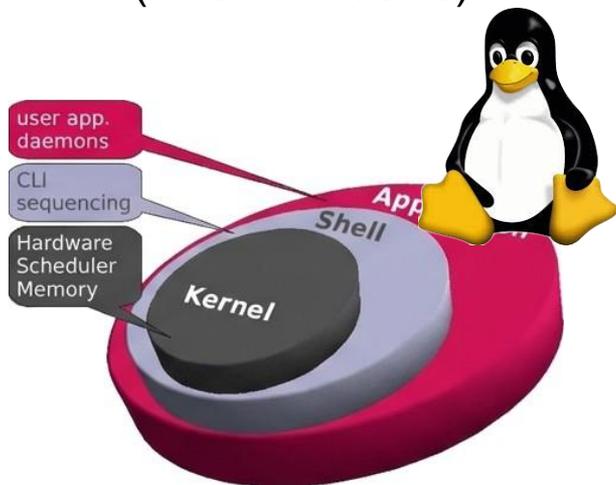
$$\mathbf{v}^T \mathbf{K} \mathbf{v} \geq 0 \quad \forall \mathbf{v}$$

(Ensures the loss surface is a bowl, not a saddle).

# A Tale of Three Kernels

Context matters.

OS Kernel  
(Linux/Windows)



Corn Kernel  
(*Zea mays*)



Image Kernel  
(Convolutions/CNNs)

131	162	232	84	91	207
104	-1	109	+1	237	109
243	-2	202	+2	135	26
185	-1	200	+1	61	225
157	124	25	14	102	108
5	155	16	218	232	249

**ML Kernel Method:** A similarity function  $K(x, z) = \langle \phi(x), \phi(z) \rangle$ .

Do this in your own time if you want a tutorial on **how to make a neural net** in PyTorch!

# PyTorch Colab