

CS M146: Introduction to Machine Learning

Kernels

Aditya Grover



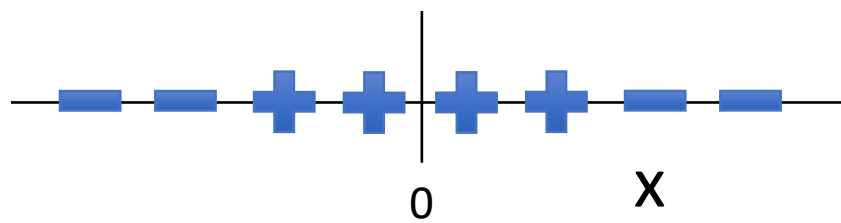
<https://aditya-grover.github.io/>



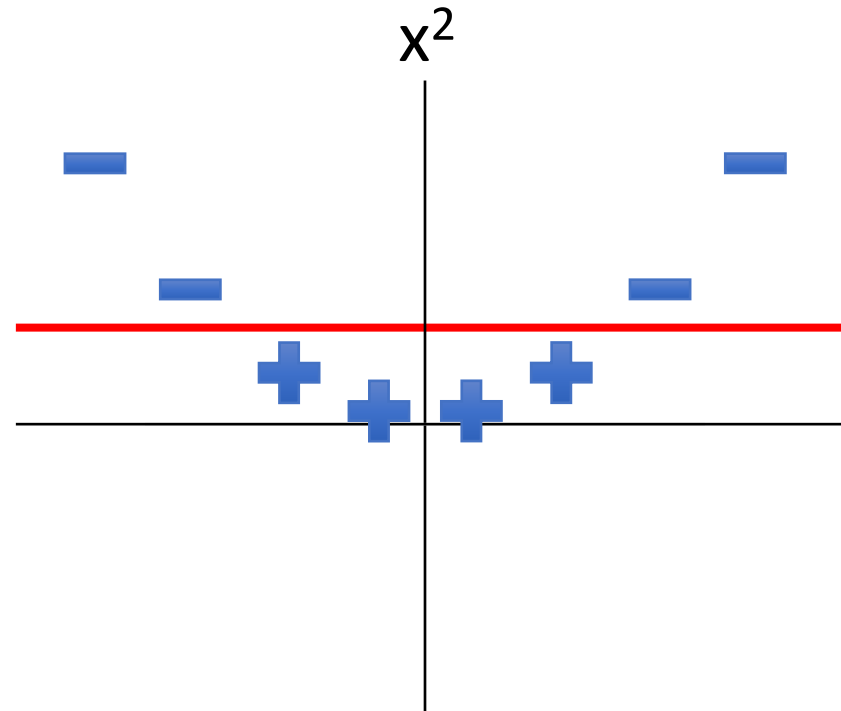
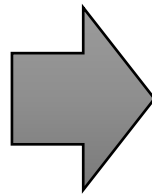
@adityagrover_

Revisiting Basis Functions

Can we separate the + and - examples using a linear decision boundary?

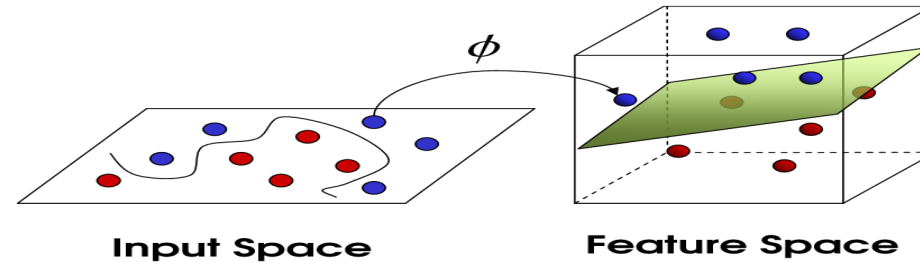


No



Yes

Mapping into a New Feature Space



- **Feature map/basis function:** Any function $\phi: \mathcal{X} \rightarrow \hat{\mathcal{X}}$ that maps input attributes to a feature space
- For example, with $\mathbf{x} \in \mathbb{R}^2$, we can define:
$$\phi(\mathbf{x}) = \phi([x_1, x_2]) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$
 - Here, $\mathcal{X} \equiv \mathbb{R}^2$ and $\hat{\mathcal{X}} \equiv \mathbb{R}^6$
 - For simplicity of notation, assume there is no bias
- **Key advantage:** Can still use a hypothesis class linear in θ
$$h_{\theta}(\mathbf{x}) = \theta^T \phi(\mathbf{x})$$

Learning With Feature Maps

- Gradient descent update rule

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}) - y^{(i)}) \phi(\mathbf{x}^{(i)})$$

- Each update requires computing $\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)})$
 - What if $\phi(\mathbf{x})$ is very high dimensional?
- E.g., consider all polynomial features with degree ≤ 3
$$\phi(\mathbf{x}) = \phi([x_1, x_2]) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_1^3, x_2^3, x_1^2 x_2, x_1 x_2^2]$$
 - **PRO:** Very expressive feature space, $\hat{\mathcal{X}} \equiv \mathbb{R}^{10}$
 - **CON 1:** Prone to overfitting (need to use regularization!)
 - **CON 2:** High compute/memory requirements for gradient descent

Revisiting Regularized Linear Models

- Hypothesis for linear models

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x})$$

- ℓ_2 regularized loss for linear models:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n L(\boldsymbol{\theta}^T \mathbf{x}^{(i)}, y^{(i)}) + \lambda \|\boldsymbol{\theta}_{1:d}\|_2^2$$

where

- **Linear Regression:** $g(z) = z$ and $L(\boldsymbol{\theta}^T \mathbf{x}^{(i)}, y^{(i)}) = \frac{1}{2} (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2$
- **Perceptron:** $g(z) = \text{sign}(z)$ and $L(\boldsymbol{\theta}^T \mathbf{x}^{(i)}, y^{(i)}) = \max(0, -y^{(i)} \boldsymbol{\theta}^T \mathbf{x}^{(i)})$
- **Logistic Regression:** $g(z) = \text{sigmoid}(z)$ and $L(\boldsymbol{\theta}^T \mathbf{x}^{(i)}, y^{(i)}) = -y^{(i)} \log(\text{sigmoid}(\boldsymbol{\theta}^T \mathbf{x}^{(i)})) - (1 - y^{(i)}) \log(1 - \text{sigmoid}(\boldsymbol{\theta}^T \mathbf{x}^{(i)}))$

Linear Models With Feature Maps

- Hypothesis for linear models **with feature maps**:

$$h_{\theta}(\mathbf{x}) = g(\boldsymbol{\theta}^T \phi(\mathbf{x}))$$

- ℓ_2 regularized loss for linear models **with feature maps**:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n L(\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}), y^{(i)}) + \lambda \|\boldsymbol{\theta}_{1:d}\|_2^2$$

Can train and evaluate linear models with expressive feature maps efficiently using **Representer Theorem**

Representer Theorem

- **Representer Theorem** (special case): For any $\lambda > 0$, there exists some real-valued coefficients $\beta_i \in \mathbb{R}$ such that the minimizer of the regularized empirical risk $J(\boldsymbol{\theta})$ can be expressed as:

$$\hat{\boldsymbol{\theta}} = \sum_{i=1}^n \beta_i \phi(\mathbf{x}^{(i)}) \quad (1)$$

- **Proof:** Since $\hat{\boldsymbol{\theta}}$ is a stationary point,

$$\nabla_{\boldsymbol{\theta}} J(\hat{\boldsymbol{\theta}}) = 0 \quad (2)$$

$$\text{Since } \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n L'(\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}), y^{(i)}) \phi(\mathbf{x}^{(i)}) + 2\lambda \boldsymbol{\theta} \quad (3)$$

$$\text{Combining (2) and (3), we get } \hat{\boldsymbol{\theta}} = -\frac{1}{2\lambda n} \sum_{i=1}^n L'(\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}), y^{(i)}) \phi(\mathbf{x}^{(i)}) \quad (4)$$

Letting $\beta_i = -\frac{1}{2\lambda n} L'(\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}), y^{(i)})$ in (4) gives (1) and finishes the proof

Reparametrized Risk

- Key immediate benefit of Representer Theorem is that it allows us to **reparameterize** the (regularized) empirical
- Empirical risk for linear models with ℓ_2 regularization:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n L(\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}), y^{(i)}) + \lambda \|\boldsymbol{\theta}\|_2^2 \quad (1)$$

- By Representer Theorem, we can plug $\hat{\boldsymbol{\theta}} = \sum_{i=1}^n \beta_i \phi(\mathbf{x}^{(i)})$ in (1) to obtain an equivalent reparameterized objective \tilde{J} w.r.t. $\boldsymbol{\beta}$

$$\begin{aligned} \tilde{J}(\boldsymbol{\beta}) &= \frac{1}{n} \sum_{i=1}^n L\left(\left(\sum_{j=1}^n \beta_j \phi(\mathbf{x}^{(j)})\right)^T \phi(\mathbf{x}^{(i)}), y^{(i)}\right) + \lambda \left\| \sum_{i=1}^n \beta_i \phi(\mathbf{x}^{(i)}) \right\|_2^2 \\ &= \frac{1}{n} \sum_{i=1}^n L\left(\sum_{j=1}^n \beta_j \phi(\mathbf{x}^{(j)})^T \phi(\mathbf{x}^{(i)}), y^{(i)}\right) + \lambda \sum_{i=1}^n \sum_{j=1}^n \beta_i \beta_j \phi(\mathbf{x}^{(j)})^T \phi(\mathbf{x}^{(i)}) \end{aligned}$$

Reparametrized Risk

- $\tilde{J}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n L(\sum_{j=1}^n \beta_j \phi(\mathbf{x}^{(j)})^T \phi(\mathbf{x}^{(i)}), y^{(i)}) + \lambda \sum_{i=1}^n \sum_{j=1}^n \beta_i \beta_j \phi(\mathbf{x}^{(j)})^T \phi(\mathbf{x}^{(i)})$

- Define a **kernel function** $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \phi(\mathbf{x}^{(j)})^T \phi(\mathbf{x}^{(i)})$

$$\tilde{J}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n L(\sum_{j=1}^n \beta_j k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}), y^{(i)}) + \lambda \sum_{i=1}^n \sum_{j=1}^n \beta_i \beta_j k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

- **Notes:**

- $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is fixed during training and can be precomputed for all i, j
- During training, we only need to optimize over $\boldsymbol{\beta}$. In contrast, earlier we needed to compute $\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)})$ at every iteration during training
- Is that it? No, more savings in line via the **kernel trick**!

Kernel Trick

- Consider the **kernel function** $k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^T \phi(\mathbf{v})$
 - Which is easier to compute: $\phi(\mathbf{u})$ or $k(\mathbf{u}, \mathbf{v})$?
 - Naively, for $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k$, both require $O(k)$ time.
- Kernel trick:** Sometimes computing $k(\mathbf{u}, \mathbf{v})$ is much more efficient
- E.g., $\phi(x) = \phi([x_1, x_2]) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_1^3, x_2^3, x_1^2 x_2, x_1 x_2^2]$
 $k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^T \phi(\mathbf{v})$

$$\begin{aligned} &= 1 + \sum_{a=1}^2 u_a v_a + \sum_{a,b \in \{1,2\}} u_a u_b v_a v_b + \sum_{a,b,c \in \{1,2\}} u_a u_b u_c v_a v_b v_c \\ &= 1 + \sum_{a=1}^2 u_a v_a + (\sum_{a=1}^2 u_a v_a)^2 + (\sum_{a=1}^2 u_a v_a)^3 \\ &= 1 + \mathbf{u}^T \mathbf{v} + (\mathbf{u}^T \mathbf{v})^2 + (\mathbf{u}^T \mathbf{v})^3 \end{aligned}$$

- Computing $\mathbf{u}^T \mathbf{v}$ (and consequently $k(\mathbf{u}, \mathbf{v})$) only requires $O(d)$ time.

Kernel Trick

- Consider the **kernel function** $k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^T \phi(\mathbf{v})$
 - Which is easier to compute: $\phi(\mathbf{u})$ or $k(\mathbf{u}, \mathbf{v})$?
 - Naively, for $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k$, both require $O(k)$ time.
- **Kernel trick:** Sometimes computing $k(\mathbf{u}, \mathbf{v})$ is much more efficient
- E.g., $\phi(x) = \phi([x_1, x_2]) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_1^3, x_2^3, x_1^2 x_2, x_1 x_2^2]$
 $k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^T \phi(\mathbf{v})$

$$\begin{aligned} &= 1 + \sum_{a=1}^2 u_a v_a + \sum_{a,b \in \{1,2\}} u_a u_b v_a v_b + \sum_{a,b,c \in \{1,2\}} u_a u_b u_c v_a v_b v_c \\ &= 1 + \sum_{a=1}^2 u_a v_a + (\sum_{a=1}^2 u_a v_a)^2 + (\sum_{a=1}^2 u_a v_a)^3 \\ &= 1 + \mathbf{u}^T \mathbf{v} + (\mathbf{u}^T \mathbf{v})^2 + (\mathbf{u}^T \mathbf{v})^3 \end{aligned}$$

- Computing $\mathbf{u}^T \mathbf{v}$ (and consequently $k(\mathbf{u}, \mathbf{v})$) only requires $O(d)$ time.

Kernelized Linear Models

- Kernel Trick: "Only compute what you need"
- We can apply the kernel trick to train and evaluate linear models with basis functions/feature maps → **Kernelized Linear Models**
- **Definition:** Given a kernel function k and a set of points $S = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \dots, \mathbf{x}^{(m)}\}$, we can define a **kernel matrix** $K \in \mathbb{R}^{m \times m}$ as:

$$K_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

Kernelized Linear Models

- **Step 1 (Pre-processing):** Pre-compute Kernel Matrix $K \in \mathbb{R}^{n \times n}$ for n training points

- **Step 2 (Kernelized Training):** Update β by gradient descent over reparameterized risk

$$\tilde{J}(\beta) = \frac{1}{n} \sum_{i=1}^n L\left(\sum_{j=1}^n \beta_j K_{ij}, y^{(i)}\right) + \lambda \sum_{i=1}^n \sum_{j=1}^n \beta_i \beta_j K_{ij}$$

- **Step 3 (Kernelized Evaluation):** Make predictions on a test point x by computing its **weighted** similarity with all training points :

$$h_{\theta}(x) = g(\theta^T \phi(x)) = g\left(\sum_{i=1}^n \beta_i k(x^{(i)}, x)\right)$$

Example: Kernalized Ridge Regression

- Standard objective

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}))^2 + \lambda \|\boldsymbol{\theta}\|_2^2$$

- Gradient Update

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \left(\frac{1}{n} \sum_{i=1}^n (\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}) - y^{(i)}) \phi(\mathbf{x}^{(i)}) + \lambda \boldsymbol{\theta} \right)$$

- Kernalized objective

$$\tilde{J}(\boldsymbol{\beta}) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \sum_{j=1}^n \beta_j K_{ij})^2 + \lambda \boldsymbol{\beta}^T \mathbf{K} \boldsymbol{\beta}$$

- Kernalized Gradient Update

$$\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} - \alpha \left(\frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^n \beta_j K_{ij} - y^{(i)} \right) \mathbf{K}^{(i)} + \lambda \sum_{i=1}^n \beta_i \mathbf{K}^{(i)} \right)$$

where $\mathbf{K}^{(i)} = \mathbf{K}_{i,:}^T$ (transpose of i -th row of \mathbf{K})

Valid Kernels

- Given a feature map ϕ , we can define a kernel function $k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^T \phi(\mathbf{v})$
- Intuitively, a kernel function defines a notion of similarity between datapoints
 - E.g., $K_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is the similarity between $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$
- There are many possible notions of similarity
 - Euclidean distance, cosine similarity
 - Make your own similarity kernel! E.g., $\log \|\mathbf{x}^{(i)}\|_2 - \log \|\mathbf{x}^{(j)}\|_2$
- Given an arbitrary kernel function $k(\mathbf{u}, \mathbf{v})$, does it imply a feature map ϕ ?

Valid Kernels

- **Definition:** A kernel is **valid** if it can be induced via a feature map ϕ
- How to determine if k is a valid kernel?
- **Method 1** (First Principles): To prove k is valid, show that there exists a function $\phi: \mathcal{X} \rightarrow \hat{\mathcal{X}}$
$$k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^T \phi(\mathbf{v}) \text{ for all } \mathbf{u}, \mathbf{v} \in \mathcal{X}$$
 - Non-trivial to apply in practice!
- **Method 2:** Use **Mercer's Theorem**
- **Method 3:** Use kernel composition rules

Mercer's Theorem

- **Mercer's Theorem:** A kernel function k is valid **if and only if** over for *any* set of $m > 0$ points, the corresponding kernel matrix K is **symmetric** and **positive semi-definite (PSD)**
- *Note:* The m points are not necessarily the training set
- **Symmetric:** $K_{ij} = K_{ji}$ for all i, j
 - Equivalently, $K = K^T$
- **Positive semi-definite:** For all $\mathbf{z} \in \mathbb{R}^m$, $\mathbf{z}^T K \mathbf{z} \geq 0$
 - Equivalently, all eigenvalues of K are real and positive

Example

- Prove $k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^2$ is a valid kernel using Mercer's theorem

- **Symmetry:**

- $k(\mathbf{v}, \mathbf{u}) = (\mathbf{v}^T \mathbf{u})^2 = (\mathbf{u}^T \mathbf{v})^2 = k(\mathbf{u}, \mathbf{v})$

- $K_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = K_{ji}$

- **PSD:** For any $\mathbf{z} \in \mathbb{R}^m$, we have:

$$\begin{aligned} \mathbf{z}^T \mathbf{K} \mathbf{z} &= \sum_{i=1}^m \sum_{j=1}^n z_i z_j K_{ij} \\ &= \sum_{i=1}^m \sum_{j=1}^n z_i z_j (\mathbf{x}^{(i)T} \mathbf{x}^{(j)})^2 \\ &= \sum_{i=1}^m \sum_{j=1}^n [(\sqrt{z_i} \mathbf{x}^{(i)})^T (\sqrt{z_j} \mathbf{x}^{(j)})]^2 \\ &\geq 0 \end{aligned}$$

Kernel Composition Rules

- Valid kernels can be **composed** with each other using many common algebraic operations to give valid kernels
- **Composition Rules:** If $k_1(\mathbf{u}, \mathbf{v})$ and $k_2(\mathbf{u}, \mathbf{v})$ are valid kernels, then the following are also valid kernels:
 - $k_1(\mathbf{u}, \mathbf{v}) + k_2(\mathbf{u}, \mathbf{v})$ is a valid kernel
 - $k_1(\mathbf{u}, \mathbf{v})k_2(\mathbf{u}, \mathbf{v})$ is a valid kernel
 - $ck_1(\mathbf{u}, \mathbf{v})$ for any positive real $c > 0$
 - $g(\mathbf{u})k_1(\mathbf{u}, \mathbf{v})g(\mathbf{v})$ for any real-valued function g
 - $\exp k_1(\mathbf{u}, \mathbf{v})$

Summary

Kernels

Motivation: Featurize inputs to arbitrary high dimensions

Representer Theorem and Kernel Trick

Ensures we can compute loss and gradients efficiently using kernels

Identifying Valid Kernels

First Principles, Mercer's Theorem, Kernel Compositions