



Feature maps and Kernels

CE-477: Machine Learning - CS-828: Theory of Machine Learning
Sharif University of Technology
Fall 2024

Fatemeh Seyyedsalehi

Feature maps

- ▶ What if the output variable y can be more accurately represented as a non-linear function of x ?
 - ▶ We need a more expressive family of models than linear models
 - ▶ A stronger hypothesis space
- ▶ Example:
 - ▶ We start by considering setting cubic functions:
$$y = w_0 + w_1x + w_2x^2 + w_3x^3$$
 - ▶ We can view the cubic function as a linear function over the a different set of feature variables

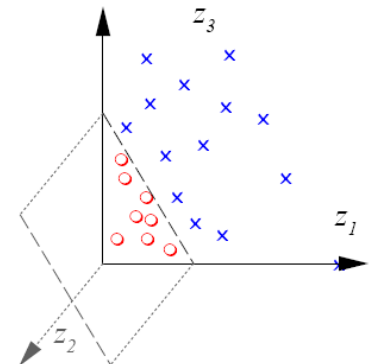
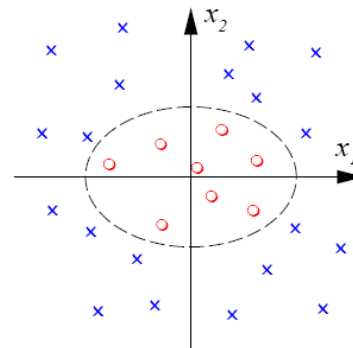
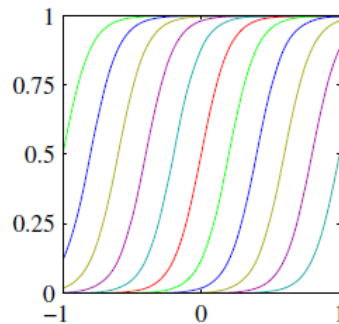
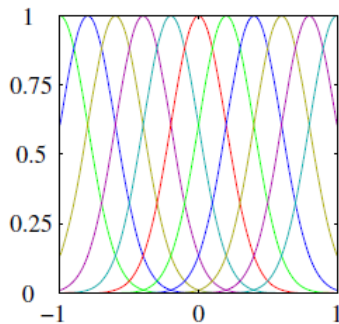
Feature maps

- Concretely, let the following function

$$\phi: \mathbb{R} \rightarrow \mathbb{R}^4 \quad \phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \in \mathbb{R}^4$$

- Thus, a cubic function of the variable x can be viewed as a linear function over the variables (x) .

$$y = w_0 + w_1x + w_2x^2 + w_3x^3 = \mathbf{w}^T \phi(x)$$



Revisiting SSE update rule

- ▶ The gradient descent algorithm for fitting the following model

$$\mathbf{w}^T \phi(x)$$

- ▶ The update rule

$$\mathbf{w} = \mathbf{w} + \eta \sum_{i=1}^N (y^i - \mathbf{w}^T x^i) x^i$$

- ▶ Consider a feature map as: $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$

Revisiting SSE update rule

- ▶ New update rule:

$$\mathbf{w} = \mathbf{w} + \eta \sum_{i=1}^N (y^i - \mathbf{w}^T \phi(x^i)) \phi(x^i)$$

- ▶ The corresponding stochastic gradient descent update rule

$$\mathbf{w} = \mathbf{w} + \eta (y^i - \mathbf{w}^T \phi(x^i)) \phi(x^i)$$

- ▶ Gradient update above becomes computationally expensive when the features map is high-dimensional.

Revisiting SSE update rule

- ▶ $d = 1000 \rightarrow$ each update needs computing and storing a 1000000000 dimensional vector

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \in \mathbb{R}^4 \longrightarrow \phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \\ x_1^2 \\ x_1 x_2 \\ \dots \\ x_1^3 \\ x^2 x_2 \\ \dots \end{bmatrix}$$

- ▶ Solution: kernel method

Revisiting SSE update rule

- ▶ The main observation is that at any time, \mathbf{w} can be represented as a linear combination of the following vectors

$$\phi(x^1), \dots, \phi(x^n)$$

- ▶ Showing this inductively:

- ▶ Initialization:

$$\mathbf{w} = 0 = \sum_{i=1}^n 0 \cdot \phi(x^i)$$

- ▶ Assume at some point we have:

$$\mathbf{w} = \sum_{i=1}^n \beta_i \cdot \phi(x^i)$$

- ▶ Then: ...

Revisiting SSE update rule

► Then:

$$\begin{aligned}\mathbf{w} &= \mathbf{w} + \eta \sum_{i=1}^N (y^i - \mathbf{w}^T \phi(x^i)) \phi(x^i) \\ &= \sum_{i=1}^n \beta_i \cdot \phi(x^i) + \eta \sum_{i=1}^N (y^i - \mathbf{w}^T \phi(x^i)) \phi(x^i) \\ &= \sum_{i=1}^N (\underbrace{\beta_i + \eta(y^i - \mathbf{w}^T \phi(x^i))}_{\text{new } \beta_i}) \phi(x^i)\end{aligned}$$

Revisiting SSE update rule

- ▶ Replacing \mathbf{w} as:

$$\mathbf{w} = \sum_{i=1}^n \beta_i \cdot \phi(x^i)$$

- ▶ The update rule is:

$$\forall i \in \{1, \dots, n\} \quad \beta_i = \beta_i + \eta \left(y^i - \left(\sum_{j=1}^n \beta_j \cdot \phi(x^j) \right)^T \phi(x^i) \right)$$

- ▶ The inner product of two feature map:

$$\phi(x^j)^T \phi(x^i): \langle \phi(x^j), \phi(x^i) \rangle$$

- ▶ Two benefits of this new update rule:
 - ▶ We can pre-compute these pairs
 - ▶ They does not necessarily require computing ϕ s explicitly.

Revisiting SSE update rule

- ▶ Prediction time:

$$\mathbf{w}^T \phi(x) = \sum_{i=1}^n \beta_i \cdot \phi(x^i)^T \phi(x) = \sum_{i=1}^n \beta_i K(x^i, x)$$

- ▶ All we need to know about the feature map is encapsulated in the corresponding kernel function.
- ▶ The kernel function corresponding to a feature map

$$\begin{aligned} \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R} \\ K(x, z) &= \langle \phi(x), \phi(z) \rangle \end{aligned}$$

Mercer Theorem

- ▶ Kernel trick → Extension of many well-known algorithms to kernel-based ones
 - ▶ By substituting the dot product with the kernel function
 - ▶ $k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}')$
 - ▶ $k(\mathbf{x}, \mathbf{x}')$ shows the dot product of \mathbf{x} and \mathbf{x}' in the transformed space.
 - ▶ kernel functions K can indeed be efficiently computed, with a cost proportional to d (the dimensionality of the input) instead of p .
- ▶ Idea: when the input vectors appears only in the form of dot products, we can use kernel trick
 - ▶ Solving the problem without explicitly mapping the data
 - ▶ Explicit mapping is expensive if $\boldsymbol{\phi}(\mathbf{x})$ is very high dimensional

Mercer Theorem

- ▶ Construct kernel functions directly
 - ▶ Ensure that it is a valid kernel
 - ▶ Corresponds to an inner product in some feature space.
- ▶ Kernels as similarity metrics:
 - ▶ We can think of $k(\mathbf{x}, \mathbf{z})$ as some measurement of how similar are $\phi(\mathbf{x})$ and $\phi(\mathbf{z})$, or of how similar are \mathbf{x} and \mathbf{z} .
- ▶ We need a way to test whether a kernel is valid without having to construct $\phi(\mathbf{x})$

Mercer Theorem

- ▶ Let $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be given. Then for K to be a valid kernel, it is **necessary and sufficient** that for any $\{x^1, \dots, x^N\}$ the corresponding kernel matrix is **symmetric positive semi-definite**.
- ▶ The Gram matrix $\mathbf{K}_{N \times N}: K_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$
 - ▶ For a set of training points $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix}$$

Kernel examples

- ▶ Polynomial kernel

$$\begin{aligned} K(x, z) &= (x^T z + c)^2 \\ &= \sum_{i,j=1}^d (x_i x_j)(z_i z_j) + \sum_{i=1}^d (\sqrt{2c} x_i)(\sqrt{2c} z_i) + c^2. \end{aligned}$$

- ▶ Gaussian kernels

- ▶ Corresponds to an infinite dimensional feature map

$$K(x, z) = \exp \left(-\frac{\|x - z\|^2}{2\sigma^2} \right).$$

Kernel examples

- ▶ Kernels also can be defined on general types of data
 - ▶ Kernel functions do not need to be defined over vectors
 - ▶ Just we need a symmetric positive definite matrix
- ▶ Thus, many algorithms can work with general (non-vectorial) data
 - ▶ Kernels exist to embed strings, trees, graphs, ...
 - ▶ E.g. for sets:
$$k(A, B) = 2^{|A \cap B|}$$
 - ▶ E.g. for graphs: random walk kernel, counts the number of paths in random walks of two graphs

References

- ▶ [1]: Andrew Ng, Machine learning, Stanford (main_notes.pdf)