Non-Linearizing Linear Methods¹¹²

- A simple non-linear extension of linear methods while keeping advantages of linear methods:
 - Map the original data to a feature space by a non-linear transformation
 - Run linear algorithm in the feature space



PCA in Feature Space

Non-linear transformation: $\phi: \boldsymbol{x} \to \boldsymbol{f}$

Feature vectors: $\{\boldsymbol{f}_i \mid \boldsymbol{f}_i = \phi(\boldsymbol{x}_i)\}_{i=1}^n$

Centered feature vectors:

$$\boldsymbol{f}_i \longleftarrow \boldsymbol{f}_i - \frac{1}{n} \sum_{j=1}^n \boldsymbol{f}_j$$

Eigenproblem: $C_f \psi = \lambda \psi$ $C_f = \frac{1}{n} \sum_{i=1}^n f_i f_i^\top$ PCA Projection: $g_i = B_{PCA} f_i$

 \mathcal{M}

$$\boldsymbol{B}_{PCA} = (\boldsymbol{\psi}_1 | \boldsymbol{\psi}_2 | \cdots | \boldsymbol{\psi}_m)^\top$$

Example

 $\blacksquare d = 2$





Example (cont.)

116

Run PCA in feature space.



Example (cont.)

Pull the results back to input space.



Non-linear PCA describes the original data much better than linear PCA.

PCA in High-Dimensional ¹¹⁹ Feature Space

- If $\dim(\mathcal{H})$ is high, description ability of non-linear PCA would be better.
- However, when dim(H) is very large, PCA in feature space is computationally demanding.
- We need a trick to reduce computational burden.

Dual Eigenproblem

$$C_f \psi = \lambda \psi$$

Since
$$C_f = \frac{1}{n} \sum_{i=1}^n f_i f_i^\top$$
,
 $\exists \alpha, \ \psi = \sum_{j=1}^n \alpha_j f_j$

Eigenproblem is

$$\frac{1}{n} \sum_{i,j=1}^{n} \alpha_j \langle \boldsymbol{f}_i, \boldsymbol{f}_j \rangle \boldsymbol{f}_i = \lambda \sum_{i=1}^{n} \alpha_i \boldsymbol{f}_i$$

Dual Eigenproblem (cont.)

121

$$\frac{1}{n}\sum_{i,j=1}^{n}\alpha_{j}\langle \boldsymbol{f}_{i},\boldsymbol{f}_{j}\rangle\boldsymbol{f}_{i}=\lambda\sum_{i=1}^{n}\alpha_{i}\boldsymbol{f}_{i}$$

In matrix form, $FK\alpha = \lambda' F\alpha$ $F \equiv (f_1 | f_2 | \cdots | f_n)$ $K_{i,j} \equiv \langle f_i, f_j \rangle \quad \lambda' \equiv n\lambda$

Solution of the above eigenproblem is given by $K\alpha = \lambda' \alpha$

 $(\lambda'_1 \ge \lambda'_2 \ge \cdots \ge \lambda'_n) \qquad (\langle \boldsymbol{\alpha}^{(i)}, \boldsymbol{\alpha}^{(j)} \rangle = \delta_{i,j})$

Kernel Trick $K\alpha = \lambda' \alpha$

- Given K, solving dual is faster than primal for $\dim(\mathcal{H}) > n$.
- However, calculating K is still painful.
- "Kernel trick": For some transformation $\phi(\boldsymbol{x})$,

$$\exists K(\boldsymbol{x}, \boldsymbol{x}') \ s.t. \ K(\boldsymbol{x}_i, \boldsymbol{x}_j) = \langle \boldsymbol{f}_i, \boldsymbol{f}_j \rangle$$

$$\boldsymbol{f}_i = \phi(\boldsymbol{x}_i)$$

122

 $\mathbf{K}(\mathbf{x}, \mathbf{x}')$:Kernel function

Kernels

$$K(\boldsymbol{x}_i, \boldsymbol{x}_j) = \langle \boldsymbol{f}_i, \boldsymbol{f}_j \rangle = \langle \phi(\boldsymbol{x}_i), \phi(\boldsymbol{x}_j) \rangle$$

- Rather than directly defining $\phi(x)$, we implicitly specify $\phi(x)$ by K(x, x').
- Kernel matrix: $K_{i,j} \equiv K(x_i, x_j)$
- Implicit mapping $\phi(x)$ exists if K(x, x') is s.t.
 - $oldsymbol{K}$ is symmetric: $oldsymbol{K}^ op=oldsymbol{K}$
- *K* is positive semi-definite: ∀*y*, ⟨*Ky*, *y*⟩ ≥ 0
 Such K(x, x') is often called the Mercer kernel or reproducing kernel.

Examples of Kernels

Polynomial kernel:

$$K(\boldsymbol{x}, \boldsymbol{x}') = \langle \boldsymbol{x}, \boldsymbol{x}' \rangle^c \qquad c \in \mathbb{N}$$

• $d = 2 \quad c = 2$
 $\langle \boldsymbol{x}, \boldsymbol{x}' \rangle^2 = (ss' + tt')^2 \qquad \boldsymbol{x} = (s, t)^\top$
 $= sss's' + 2ss'tt' + ttt't'$
 $\mathbf{f} = \phi(\boldsymbol{x}) = (s^2, \sqrt{2st}, t^2)^\top$
 $\dim(\mathcal{H}) = 3$

• In general, $\dim(\mathcal{H}) = \begin{pmatrix} c+d-1\\ c \end{pmatrix}$

Examples of Kernels (cont.) ¹²⁵

Gaussian kernel:

$$K(x, x') = \exp(-||x - x'||^2/c)$$
 $c > 0$

 $\dim(\mathcal{H}) = \infty$

Playing with Kernel Trick $\langle \boldsymbol{f}_i, \boldsymbol{f}_j \rangle = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$

Exercise: Prove the following • $\|\boldsymbol{f}_i - \boldsymbol{f}_j\|^2 = K(\boldsymbol{x}_i, \boldsymbol{x}_i) - 2K(\boldsymbol{x}_i, \boldsymbol{x}_j) + K(\boldsymbol{x}_j, \boldsymbol{x}_j)$

•
$$\cos \theta = \frac{K(\boldsymbol{x}_i, \boldsymbol{x}_j)}{\sqrt{K(\boldsymbol{x}_i, \boldsymbol{x}_i)K(\boldsymbol{x}_j, \boldsymbol{x}_j)}}$$

For Gaussian kernel

•
$$\|\boldsymbol{f}_i - \boldsymbol{f}_j\|^2 = 2 - 2K(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

•
$$\cos \theta = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$$



Centering in Feature Space ¹²⁷

For implicit feature map, how do we center $\{f_i\}_{i=1}^n$?

$$oldsymbol{f}_i \longleftarrow oldsymbol{f}_i - rac{1}{n}\sum_{j=1}^n oldsymbol{f}_j$$

We only need centered kernel matrix.Centered kernel matrix is given by

$$K \leftarrow -K - NK - KN + NKN$$

 $N_{i,j} = 1/n$

Homework

128

Prove that the centered kernel matrix is given by

 $K \leftarrow K - NK - KN + NKN$ $N_{i,j} = 1/n$