

Kernel Ridge Regression & The Kernel trick

Least Squares:

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

(Ridge Reg)

closed form sol:

$$w^* = (X^T X + \lambda I_d)^{-1} X^T y$$

$X \in \mathbb{R}^{N \times d}$

Proof: $P \in \mathbb{R}^{N \times N}$ and $Q \in \mathbb{R}^{N \times N}$

$$P(QP + I) = PQP + P = (PQ + I)P$$

Assume that $(QP + I)$ and $(PQ + I)$ are invertible:

$$(QP + I)^{-1}P = P(PQ + I)^{-1}$$

then $P = X^T$ and $Q = \frac{I}{\lambda}$

$$\Rightarrow (\frac{1}{\lambda} X^T X + I)^{-1} X^T = X^T (\frac{1}{\lambda} X X^T + I)^{-1}$$

$$\Rightarrow (X^T X + \lambda I)^{-1} X^T = X^T (X X^T + \lambda I)^{-1}$$

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix}_{N \times D}$$

feature
cols

$$X^T = \begin{pmatrix} x_1 & x_2 & \dots & x_N \end{pmatrix}_{D \times N}$$

feature
rows

Alternative:

We will show:

$$w^* = X^T \underbrace{(X X^T + \lambda I_N)^{-1}}_{{D \times N} \quad {N \times N}} y \quad (\text{Ridge Regression})$$

Why the alternative form:

1) Complexity:

- $(X^T X + \lambda I)^{-1} X^T y = O(d^3 + Nd^2)$
- $X^T (X X^T + \lambda I_N)^{-1} y = O(N^3 + dN^2)$

2) Structure:

$$w^* = X^T \alpha^* \quad \text{where } \alpha^* = (X X^T + \lambda I_N)^{-1} y \quad (N \times 1)$$

$$w^* = \sum \alpha_i x_i$$

w^* is a comb. of the feature vectors
 $w^* \in$ the space spanned by the feature vectors
(column space of X^T)

Representation theorem:

For any loss function \mathcal{L} , if $w^* = \min_w \sum_{i=1}^N \mathcal{L}(x_i^T w, y_i) + \frac{\lambda}{2} \|w\|_2^2$
then there exists $\alpha^* \in \mathbb{R}^N$ s.t. $w^* = X^T \alpha^*$

$$w^* = \arg \min_w \frac{1}{2} \|y - Xw\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

$$\alpha^* = \arg \min \frac{1}{2} \alpha^T (X X^T + \lambda I) \alpha - \alpha^T y$$

$$\Rightarrow w^* = X^T \alpha^*$$

why is this of interest?

• Complexity (N vs D): if $N \gg D \rightarrow$ Normal sol:

• if $N \ll D \rightarrow$ Alt. sol:

• $K = X X^T$: Kernel Matrix

→ Kernel trick

Kernel function:

$$K = \begin{pmatrix} x_1^T x_1 & x_1^T x_2 & \dots \\ x_2^T x_1 & \dots & \dots \end{pmatrix}_{N \times N} = (x_i^T x_j)_{i,j}$$

Feature map: $\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$
with $d \ll D$; $x_i \sim (\Phi(x_i))_{j=1}^D$

with $d \ll D$; $x_i \sim (\Phi(x_i))_{j=1}^D$

$$K = \begin{pmatrix} \phi(x_1)^T \phi(x_1) & \dots & \phi(x_1)^T \phi(x_n) \\ \vdots & \ddots & \vdots \\ \phi(x_n)^T \phi(x_1) & \dots & \phi(x_n)^T \phi(x_n) \end{pmatrix} \in \mathbb{R}^{N \times N}$$

Kernelfkt II:

$$K(x, x') = \phi^T(x) \phi(x')$$

(bupper) $\xrightarrow{\quad} O(d)$

Kerneltrick:

$$K = (K(x_i, x_j))_{ij}$$

Ex:

- 1) $K(x, x') = x^T x' = \phi^T(x) \phi(x') \Rightarrow \phi(x) = x$
- 2) $x, x' \in \mathbb{R}, K(x, x') = (x \cdot x')^2 \Rightarrow \phi(x) = x^2, \phi^T(x) \phi(x') = x^2 \cdot x'^2 = (x \cdot x')^2$
- 3) $x, x' \in \mathbb{R}^3, K(x, x') = (x_1 x'_1 + x_2 x'_2 + x_3 x'_3)^2 = \phi^T(x) \phi(x)$
 $x = (x_1, x_2, x_3), x' = (x'_1, x'_2, x'_3) \Rightarrow K(x, x') = (x_1 x'_1)^2 + (x_2 x'_2)^2 + (x_3 x'_3)^2 + 2(x_1 x'_1 x'_2 + x_1 x'_1 x'_3 + x_2 x'_2 x'_3)$
 $\phi(x) = (x_1^2 \ x_2^2 \ x_3^2 \ \ln x_1 x_2 \ \ln x_1 x_3 \ \ln x_2 x_3)^T \in \mathbb{R}^6$
 $K(x, x') = \phi^T(x) \phi(x)$ (Polynomial Kernel)

4) Radial basis fkt:

$$x, x' \in \mathbb{R}^d, K(x, x') = e^{-(x-x')^T(x-x')} \quad O(d)$$

2. f. assume $x, x' \in \mathbb{R}$

$$\begin{aligned} K(x, x') &= e^{-(x-x')^2} \\ &= e^{-x^2 - x'^2 + 2x \cdot x'} \\ &= e^{-x^2} e^{-x'^2} e^{2x \cdot x'} \\ &= e^{-x^2} e^{-x'^2} \sum_{k=0}^{\infty} \frac{x^k x'^k}{k!} \quad \text{Tayl. exp.} \end{aligned}$$

$$\phi(x) = e^{-x^2} \left(\dots, \frac{\sqrt{2}}{\sqrt{k!}} x^k, \dots \right)$$

$$\phi^T(x) \phi(x) = K(x, x') \quad (\text{inf. dim.})$$

Building new kernels:
from old ones

• K_1, ϕ_1 and K_2, ϕ_2 :

$$k(x, x') = \alpha K_1(x, x') + \beta K_2(x, x') \quad \alpha, \beta \geq 0$$

$$= \alpha \phi_1^T(x) \phi_1(x') + \beta \phi_2^T(x) \phi_2(x')$$

$$= (\sqrt{\alpha} \phi_1(x), \sqrt{\beta} \phi_2(x)) \begin{pmatrix} \sqrt{\alpha} \phi_1(x') \\ \sqrt{\beta} \phi_2(x') \end{pmatrix}$$

$$= \phi^T(x) \phi(x')$$

$$\phi(x) = \begin{pmatrix} \sqrt{\alpha} \phi_1(x) \\ \sqrt{\beta} \phi_2(x) \end{pmatrix} \in \mathbb{R}^{d_1+d_2}$$

$$\cdot k(x, x') = k_1(x, x') \cdot k_2(x, x')$$

$$\phi \in \mathbb{R}^{d_1 \times d_2}, \quad K(x, x') = \phi_1^T(x) \phi_1(x') \phi_2^T(x) \phi_2(x') \\ = \phi^T(x) \phi(x')$$

$$\phi(x) = \begin{pmatrix} (\phi_1(x))_1 & \phi_2(x) \\ \vdots & \\ (\phi_1(x))_{d_1} & \phi_2(x) \end{pmatrix} \in \mathbb{R}^{d_1 \times d_2}$$

$$\phi^T(x) \phi(x') = \sum_{i=1}^{d_1} (\phi_1(x))_i \phi_2(x')^T (\phi_1(x'))_i \phi_2(x')$$

$$= \sum_{i=1}^{d_1} (\phi_1(x))_i (\phi(x'))_i \phi_2^T(x) \phi_2(x')$$

$$= \phi_1^T(x) \phi_1(x') \phi_2^T(x) \phi_2(x') = K(x, x')$$

• There exist a lot of other rules.

Mercer's condition:

Given k (kernel), how do we know it $\exists \alpha \phi$ s.t.:

$$K(x, x') = \phi^T(x) \phi(x')$$

Mercer's condition:

$$\cdot k(x, x') = k(x', x) \quad \forall x, x'$$

$$\cdot \forall N \geq 0, \#(x_i)_{i=1}^N \quad K = (k(x_i, x_j))_{i,j=1}^N \quad K \text{ must be positive semi-definite (PSD)}$$

<p><u>Predicting using Kernel:</u></p> $\phi^T(x) = \begin{pmatrix} \phi(x_1) & \dots & \phi(x_n) \end{pmatrix}_{D \times N}$	<p>$w^* \rightsquigarrow \phi^T(x) w^*$. D can be very large</p> <p>How to predict only using the Kernel:</p> $\phi^T(x) w^* = \sum_N \phi^T(x) \phi(x)^T \alpha^*$ $\phi^T(x) w^* = \sum_{i=1}^N k(x, x_i) \alpha_i \rightarrow much\ smaller$ <p><u>Rq:</u> $y = \phi^T(x) w^*$ linear pred in the concatenated space $= f_{w^*}(x)$ non-linear in the original space. (\checkmark)</p>
---	---