

April 28, 2020 kernel methods

$$\|x_i - x_j\|_2^2 = \cancel{\|x_i\|_2^2} - \cancel{\|x_j\|_2^2} + 2x_i^T x_j$$

$x_i^T x_j \approx$ how similar x_i, x_j are

node2vec
t-SNE
PCA

$$k(x_i, x_j) = x_i^T x_j \quad (\text{similarity})$$

$$k(x_i, x_j) = -\|x_i - x_j\|_2^2$$

Example: squared exponential / Gaussian

$$k(x_i, x_j) = \sigma^2 \exp\left(-\frac{1}{2} \|x_i - x_j\|_2^2 / \ell^2\right)$$

Example: polynomial (degree 2)

$$k(x_i, x_j) = \underline{\underline{(c + x_i^T x_j)^2}}$$

$$(c + \gamma^T z)^2 = \left(c + \sum_{k=1}^d \gamma_k z_k \right)^2$$

$$= c^2 + 2c \sum \gamma_k z_k + \left(\sum \gamma_k z_k \right)^2$$

$$= c \cdot c + \sum (\sqrt{2c} \gamma_k) (\sqrt{2c} z_k) + \sum_{k=1}^d \gamma_k^2 z_k^2 + \sum_{i < j} (\sqrt{2} \gamma_i \gamma_j) (\sqrt{2} z_i z_j)$$

$$\phi(x) = (c, \sqrt{2c} x_1, \dots, \sqrt{2c} x_d, x_1^2, \dots, x_d^2, \sqrt{2} x_1 x_2, \dots, \sqrt{2} x_{d-1} x_d)$$

$$\Rightarrow = \phi(y)^T \phi(z) \quad \phi: x \in \mathbb{R}^d \rightarrow \phi(x) \in \mathbb{R}^{1+d+d+\binom{d}{2}}$$

ϕ is called a feature map

What about Gaussian?

$$k(x_i, x_j) = \exp\left(-\frac{1}{2} \|x_i - x_j\|_2^2\right)$$

$$= 1 - \frac{r^2}{2} + \frac{r^4}{8} - \frac{r^6}{48} + \dots$$

$$r = \|x_i - x_j\|_2$$

$$\exp\left(-\frac{1}{2} r^2\right)$$

Key idea: $k(x_i, x_j)$ often easy to compute
 $k(x_i, x_j) = x_i^T x_j$ $k(x_i, x_j) = \exp(-\|x_i - x_j\|_2^2)$

Key idea: if our method only relies on $x_i^T x_j$ as measure of similarity, then we can think of using $k(x_i, x_j)$
"kernel trick"

Example: kernel PCA

$$n \times \begin{matrix} d \\ \boxed{X} \end{matrix} = n \times \begin{matrix} d \\ \boxed{U} \end{matrix} \begin{matrix} d \\ \boxed{\Sigma} \end{matrix} \begin{matrix} d \\ \boxed{V^T} \end{matrix}$$

k

columns of X (features) have zero mean

PCs: V_k

Projection: $X V_k = U_k \Sigma_k$

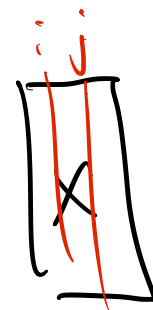
$$d \quad X^T X = V \Sigma U^T U \Sigma V^T = V \Sigma^2 V^T$$

PCA: first k eigenvectors of $X^T X$

sample covariance



$$(X^T X)_{ij} = \underline{x_i^T x_j}$$



$$X^T X \Rightarrow K \quad K_{ij} = k(x_i, x_j)$$

$$k(x, y) = k(y, x)$$

$$K \succeq 0$$

drop-in replacement?

Symmetric

\bar{V}_k = first k eigenvectors of \bar{K}

$$\bar{K} = (I - \frac{1}{n} \mathbf{1}\mathbf{1}^T) K (I - \frac{1}{n} \mathbf{1}\mathbf{1}^T)$$

zero near

Function approx.

$$\underline{x_1, \dots, x_n} \in \mathbb{R}^d \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$$

demographic data in county \xRightarrow{f} election (% Biden / Trump)

$$f \approx \hat{f} \quad \hat{f}(x) = x^T c^*$$

$$c^* = \underset{c}{\operatorname{arg\,min}} \sum \|Xc - f_x\|_2^2$$

$\begin{matrix} x_1^T \\ \vdots \\ x_n^T \end{matrix}$

$\begin{matrix} f(x_1) \\ \vdots \\ f(x_n) \end{matrix}$

Basis functions ϕ_1, \dots, ϕ_d $\phi_j(x) = \underline{x_j}$

$$\begin{aligned} \hat{f} &= \sum c_j \phi_j & \hat{f}(x) &= \sum c_j \phi_j(x) = \sum c_j x_j \\ & & &= c^T x = x^T c \end{aligned}$$

More complex space

$$\psi_j: \mathbb{R}^d \rightarrow \mathbb{R} \quad j = 1, \dots, m$$

$$\hat{f} = \sum_{j=1}^m c_j \psi_j \quad \hat{f}(x) = \sum_{j=1}^m c_j \psi_j(x)$$

$$\min_c \|\Psi c - f_x\|_2^2$$

$$\begin{matrix} \psi_1(x_1) & \dots & \psi_m(x_1) \\ \vdots & & \vdots \\ \psi_1(x_n) & \dots & \psi_m(x_n) \end{matrix}$$

Key idea: ψ_j can be nonlinear (feature maps)
still get OLS

$$\bullet \quad n \geq m \implies c^* = (\Psi^T \Psi)^{-1} \Psi^T f_x$$

$$\star \quad n \leq m \implies \text{HW1} \quad \Psi c^* = f_x$$

$$\min_c \|c\|_2^2$$

$$\text{s.t. } \Psi c = f_X$$

$$\Rightarrow c^* = \Psi^T (\Psi \Psi^T)^{-1} f_X$$

$$\hat{f} = \sum c_j^* \psi_j \quad \hat{f}(z) = \sum c_j^* \psi_j(z) = \Psi(z)^T c^*$$

$$= \Psi(z)^T \Psi^T (\Psi \Psi^T)^{-1} f_X$$

$$\begin{matrix} \text{m} \\ \left[\psi_1(z) \dots \psi_m(z) \right] \\ \text{m} \end{matrix} \Psi^T$$

$$\begin{matrix} \downarrow \\ K_{XX} \end{matrix} \quad (K_{XX})_{ij} = \Psi(x_i)^T \Psi(x_j) = k(x_i, x_j)$$

$$\begin{aligned} \triangleq K_{ZX} \quad (K_{ZX})_j &= \sum_l \psi_l(z) \psi_l(x_j) \\ &= \Psi(z)^T \Psi(x_j) \\ &= k(z, x_j) \end{aligned}$$

$$\hat{f}(z) = \underline{K_{ZX}} \underline{K_{XX}^{-1}} \underline{f_X} \quad K_{XX} \rightarrow 0$$

Vector space H with orthonormal basis ψ_1, \dots, ψ_m

Approx. f with $\hat{f} \in H$ $\langle \psi_i, \psi_j \rangle_H = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$

$$\hat{f} = \sum c_j \psi_j$$

$$\begin{aligned} \|\hat{f}\|_H^2 &= \left\langle \sum_j c_j \psi_j, \sum_i c_i \psi_i \right\rangle_H \\ &= \sum_{i,j} c_i c_j \langle \psi_i, \psi_j \rangle_H \\ &= \sum_i c_i^2 \langle \psi_i, \psi_i \rangle_H = \sum_i c_i^2 = \|c\|_2^2 \end{aligned}$$

$$\min_{\hat{f}} \|\hat{f}\|_H^2 \quad \text{s.t.} \quad \hat{f}_x = f_x$$

example:  $\langle \psi_i, \psi_j \rangle_H = \int_{S^1} \psi_i(x) \psi_j(x) dx$

Point approx $k_y \in H$ $k_y(x) = \sum_{j=1}^{\infty} \psi_j(y) \psi_j(x)$
 $k_y = \sum \psi_j(y) \psi_j$

$g \in H$ $g = \sum a_i \psi_i$ $g(y) = \sum a_i \psi_i(y)$

$\langle g, k_y \rangle_H = \langle \sum a_i \psi_i, \sum \psi_j(y) \psi_j \rangle_H$

$= \sum_{j=1}^{\infty} a_i \psi_j(y) \langle \psi_i, \psi_j \rangle_H$

$= \sum_{i=1}^{\infty} a_i \psi_i(y) \langle \psi_i, \psi_i \rangle_H$

$= g(y)$ "reproducing"

need to be more careful in infinite dims

H is a Reproducing Kernel Hilbert Space (RKHS)
 if $g(y) = \langle g, k_y \rangle_H$ for any $g \in H$

feature maps orthonormal basis ψ_1, \dots, ψ_m for H

\Rightarrow RKHS

Need: inner product $\langle \psi(x_i), \psi(x_j) \rangle_{\ell_2}$

Other direction: RKHS \Rightarrow kernel function

kernel \Rightarrow RKHS