



linear predictor $w^T x$
 $w^T \phi(x)$
 \uparrow feature expansion.

Feature Expansion.

① $x \in \mathbb{R}^d$ "Quadratic Expansion"

$$\phi(x) = (1, \sqrt{2} x_1, \sqrt{2} x_2, \dots, \sqrt{2} x_d, x_1^2, \dots, x_d^2, \sqrt{2} x_1 x_2, \sqrt{2} x_1 x_3, \dots, \sqrt{2} x_{d-1} x_{d-2})$$

$$\phi(x)^T \phi(x') = \underline{(1 + x^T x')^2}$$

better than $O(d^2)$? $O(d)$ linear time.

② "Products of all subsets"

$$\phi(x) = \left(\prod_{i \in S} x_i \right)_{S \subseteq [d]}$$

$$(x_1, x_2, \dots, x_d, x_1 x_2, \dots, x_{d-1} x_d, \dots, x_1 \dots x_d) \in \mathbb{R}^{2^d}$$

$$\phi(x)^T \phi(x) = \prod_{i=1}^d (1 + \tilde{\kappa}_i \tilde{\kappa}_i')$$

How hard?
 $O(2^d)$

\downarrow
 $O(d)$ time

③ $\phi(x)^T \phi(x) = \exp\left(-\frac{\|x-x'\|_2^2}{2\sigma^2}\right)$

"Gaussian Kernel"

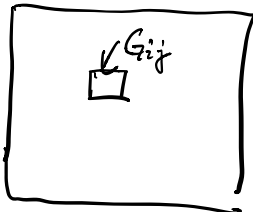
$O(d)$ linear

$$\phi(x) \in \mathbb{R}^{\infty}$$

$$\text{In } x \in \mathbb{R}, \phi(x) = \exp\left(\frac{-x^2}{2\sigma^2}\right) \left(1, \frac{x}{\sigma}, \frac{1}{2!} \left(\frac{x}{\sigma}\right)^2, \dots\right)$$

Kernel

A kernel function $K: \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ is a symmetric function such that for any $x_1, \dots, x_n \in \mathcal{X}$,

the Gram matrix G  $G_{ij} = K(x_i, x_j)$

is positive semidefinite.

$$\rightarrow \underline{\phi^T G \phi \geq 0, \quad \forall \phi \in \mathbb{R}^n.}$$

How to show $K(\cdot, \cdot)$ is a kernel?

Easy way: Find ϕ mapping such that

$$\underline{\phi(x)^T \phi(x')} = K(x, x').$$

new features.

"similarity measure"

↳ implies the definition

$$G_{ij} = \phi(x_i)^T \phi(x_j)$$

$$G = \underbrace{\Phi^T \Phi}_{\text{p.s.d.}}, \quad \Phi = [\phi(x_1), \dots, \phi(x_n)]$$

$$\forall f, \quad f^T \Phi^T \Phi f = (\Phi f)^2 \geq 0.$$

Examples:

① $K(x, x') = \underline{x^T x'}$
when is this useful?
 $n \ll d$

② Polynomial: $K(x, x') = (1 + x^T x')^k$

③ Gaussian Kernel/
Radial Basis Function (RBF)

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

.....
other examples in the note.

Dual SVM

$$\begin{aligned} \max_{\alpha, \lambda} \quad & \sum_i \lambda_i - \frac{1}{2} \sum_{i, j \in [n]} \lambda_i \lambda_j y_i y_j \underbrace{\phi^T(x_i) \phi(x_j)}_{k(x_i, x_j)} \\ \text{s.t.} \quad & 0 \leq \lambda_i \leq C. \end{aligned}$$

For RBF, w lives in \mathbb{R}^∞ ?

$$\underline{w^*} = \sum_{i=1}^n y_i \lambda_i^* \underline{\phi(x_i)}$$

For a new test point x .

$$w^{*T} \phi(x) = \sum_{i=1}^n y_i \lambda_i^* \underbrace{\phi(x_i)^T \phi(x)}_{k(x_i, x)}$$

kernel \leftrightarrow "similarity"

Build new kernels from old kernels.

$$k(x, y) = c k_1(x, y), \text{ for } c > 0.$$

kernel. kernel

HW₂ {

$$k(x, y) = \underbrace{k_1(x, y)}_{\text{kernel.}} + \underbrace{k_2(x, y)}_{\text{kernels}}$$

$$k(x, y) = k_1(x, y) k_2(x, y)$$

other examples in notes.

Kernel Ridge Regression

$$A = \begin{bmatrix} \leftarrow x_1 \rightarrow \\ \vdots \\ \leftarrow x_n \rightarrow \end{bmatrix} \quad b = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$\text{Ridge Reg} : \hat{w} = (A^T A + \lambda I)^{-1} A^T b$$

$$\begin{aligned} \text{Alternative form} : & \quad (A^T A + \lambda I)^{-1} A^T \\ & = A^T (G + \lambda I_n)^{-1} \quad \left(\begin{array}{l} \text{See note} \\ \text{for a proof} \end{array} \right) \\ & \quad \downarrow \\ & \quad G \in \mathbb{R}^{n \times n}, \quad G_{ij} = x_i^T x_j \end{aligned}$$

|

$$\hat{w} = A^T \underbrace{(G + \lambda I_n)^{-1}}_v b$$

\Downarrow
 $G_{ij} = \phi(x_i)^T \phi(x_j)$

$$\hat{w} = A^T v = \sum_{i=1}^n v_i x_i$$

New test point x

prediction $x^T \hat{w} = \sum_{i=1}^n v_i \underbrace{x^T x_i}_{k(x, x_i)}$

\downarrow
 $\phi(x)^T$

$$\phi(x)^T \hat{w} = \sum_{i=1}^n v_i \underbrace{\phi(x)^T \phi(x_i)}_{k(x, x_i)}$$