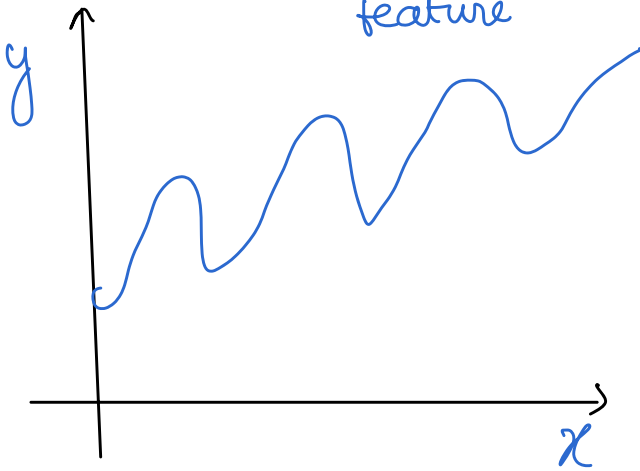


Kernel Methods

$$\underset{\text{attributes}}{x} \rightarrow \phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^4 \end{bmatrix}$$

feature

1d x mapped to higher dimensional space.
It can 4, 8 or ∞ .



Similar
for classification

Normally:

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \sum_{i=1}^n \underbrace{(y^{(i)} - \theta^T x^{(i)})}_{\text{Scalar}} \cdot x^{(i)} \quad \theta \in \mathbb{R}^d$$

with feature map

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \sum_{i=1}^n \underbrace{(y^{(i)} - \theta^T \phi(x^{(i)}))}_{\text{Scalar}} \phi(x^{(i)})$$

$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^p$

$\theta \in \mathbb{R}^p$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1^2 \\ x_1 x_2 \\ x_1 x_3 \\ \vdots \\ x_1^3 \\ x_1^2 x_2 \\ \vdots \end{bmatrix} \quad \begin{array}{l} \text{monomial terms of order } \leq 3 \\ p \approx O(d^3) \\ d^3 \rightarrow \text{features} \end{array}$$

Normal	$d = 1000$	$O(d)$
Feature Map	$(1000)^3$	$O(d^3)$

Claim:

$$\theta^{(t)} = \sum_{i=1}^n \beta_i^{(t)} \phi(x^{(i)}) \quad \rightarrow \text{every } \theta^{(t)} \text{ is linear comb}^n \text{ of } \phi(x)$$

evident from G.D. update rule

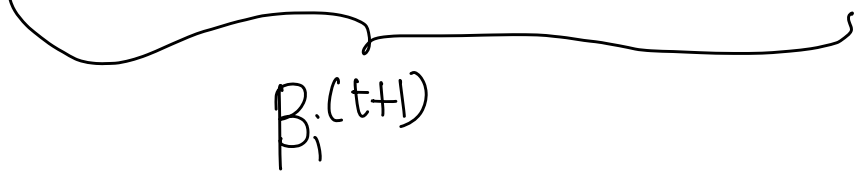
$$\theta^{(0)} = 0$$

$$\theta^{(1)} = \sum_{i=1}^n \underbrace{\alpha y^{(i)}}_{\beta_i^{(1)}} \phi(x^{(i)})$$

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \sum_{i=1}^n (y^{(i)} - \theta^{(t)T} \phi(x^{(i)})) \cdot \phi(x^{(i)})$$

$$= \sum_{i=1}^n \beta_i^{(t)} \phi(x^{(i)}) + \alpha \sum_{i=1}^n (y^{(i)} - (\sum_{j=1}^n \beta_j^{(t)} \phi(x^{(j)}))^T \phi(x^{(i)})) \cdot \phi(x^{(i)})$$

$$= \sum_{i=1}^n \left[\beta_i^{(t)} + \alpha (y^{(i)} - (\sum_{j=1}^n \beta_j^{(t)} \phi(x^{(j)}))^T \phi(x^{(i)})) \right] \phi(x^{(i)})$$


 $\beta_i^{(t+1)}$

$i \rightarrow i^{th}$ example

$\phi(x^{(i)})^T \phi(x^{(i)}) \rightarrow$ This can be precomputed as examples remain same.

$$\beta_i^{(t+1)} = \beta_i^{(t)} + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} \phi(x^{(j)})^T \phi(x^{(i)}))$$

Two elements of
the space

$$\text{Kernel} \triangleq k: \overset{!}{X} \cdot \overset{!}{X} = \mathbb{R}$$

$$\begin{aligned} K(x, z) &= \langle \phi(x), \phi(z) \rangle \\ &= \phi(x)^T \phi(z) \end{aligned}$$

$$x \in X, x \in \mathbb{R}^d \quad X = \mathbb{R}^d$$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1^3 \\ \vdots \end{bmatrix}$$

$$\begin{aligned} K(x, z) &= 1 + \langle x, z \rangle + \langle x, z \rangle^2 + \langle x, z \rangle^3 \\ &= \phi(x)^T \phi(z) \end{aligned} \quad \begin{array}{l} O(d^3) \text{ reduced} \\ \text{to } O(d) \end{array}$$

$$\phi: X \rightarrow \mathbb{R}^p$$

LINEAR REGRESSION (Kernalized)

<1> Precompute:

$$K_{ij} = K(x^{(i)}, x^{(j)}) = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$

($n^2 \rightarrow$ inner products)

$K \rightarrow$ stand for both kernel function & kernel matrix where kernel matrix is square symmetric matrix where dot product betⁿ all exs. pre evaluated.

<2> Loop:

$$\forall i \in \{1 \dots n\}$$

$$\beta_i^{(t+1)} = \beta_i^{(t)} + \alpha (y^{(i)} - \sum_{j=1}^n \beta_j^{(t)} K_{ij})$$

$$\beta^{(t+1)} = \beta^{(t)} + \alpha (\vec{y} - K \beta^{(t)}) \rightarrow \text{vectorized form}$$

Prediction

$$\beta^{(t)} \in \mathbb{R}^n$$

$$h_{\theta}(x) = \theta^T \phi(x)$$

$$= \sum_{i=1}^n \beta_i \phi(x^{(i)})^T \phi(x)$$

$$= \sum_{i=1}^n \beta_i K(x^{(i)}, x)$$

$x \rightarrow$ test example
 $x^{(i)} \rightarrow$ all training examples

Observations

$$\begin{aligned} \text{① Train: } \beta &:= \beta + \alpha(\vec{y} - K\beta) \\ \text{Test: } \hat{y} &= \sum_{i=1}^n K(x^{(i)}, x) \cdot \beta_i \end{aligned} \quad \left. \vphantom{\sum_{i=1}^n} \right\} \phi(x) \text{ does not appear}$$

② For Prediction
Need training example to be stored in memory
we give up $\phi(x)$

$$\begin{bmatrix} \beta^{(1)} \\ \vdots \\ \beta^{(n)} \end{bmatrix} \leftarrow \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \begin{matrix} x^{(i)} \in \mathbb{R}^d \\ -x^{(1)T} \\ \vdots \\ -x^{(n)T} \end{matrix} \quad \begin{bmatrix} \vec{y} \end{bmatrix}$$

$$\begin{array}{c} \boxed{\theta \in \mathbb{R}^d} \quad \theta^{(0)} \\ \downarrow \\ \boxed{ \text{ one comp per feature}} \quad \theta^{(1)} \\ \downarrow \\ \boxed{} \quad \theta^{(2)} \end{array}$$

$$\theta^{(t)} = \sum_{i=1}^n \beta_i^{(t)} x^{(i)}$$

$\beta \rightarrow 1$ per example (This is what allows us to scale to ∞ dimensions)

For logistic,

$$\beta^{(t+1)} = \beta^{(t)} + \alpha (\vec{y} - g(k\beta^{(t)}))$$

Kernel examples

$$x \in \mathbb{R}^d$$

$$1) K(x, z) = \langle x, z \rangle^2$$

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_d x_d \end{bmatrix}$$

First we come up
with $K(x, z)$ then
 $\phi(x)$

$$2) K(x, z) = (x^T z + c)^2$$

$$\phi(x) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ \vdots \\ \sqrt{2c} x_1 \\ \sqrt{2c} x_2 \end{bmatrix}$$

$$K(x, z) = \phi(x)^T \phi(z)$$

$K(x, z) \uparrow$ for similar (x, z)

\downarrow for dissimilar (x, z)

$$K(x, z) = \phi(x)^T \phi(z)$$

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right) \rightarrow \infty \text{ dimensional kernel vector}$$

↓
Popular kernel!!
Gaussian kernel

Necessary condⁿ for k to be a kernel

(1) k should be symmetric

$$K(x, z) = K(z, x)$$

(2) k is P.S.D.

$$\begin{aligned} Z^T K Z &= \sum_i \sum_j z_i k_{ij} z_j \\ &= \sum_k \left(\sum_i z_i \phi_k(x^{(i)}) \right)^2 \geq 0 \end{aligned}$$

Mercer Theorem

Let $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be a given function

For K to be kernel, it is necessary & sufficient for any $\{x^{(1)} \dots x^{(m)}\}$ the corresponding matrix $K_{ij} = K(x^{(i)}, x^{(j)})$ is P.S.D.

vector \longleftrightarrow functions
 matrices \longleftrightarrow operators

$$\begin{array}{c}
 \xrightarrow{x+h} \\
 K(x, z)
 \end{array}
 \begin{array}{c}
 \downarrow z+h \\
 \begin{bmatrix}
 \cdot & \cdot & \cdot & \cdot \\
 \cdot & \cdot & \cdot & \cdot \\
 \cdot & \cdot & \cdot & \cdot \\
 \cdot & \cdot & \cdot & \cdot \\
 \cdot & \cdot & \cdot & \cdot
 \end{bmatrix}
 \end{array}
 \begin{array}{c}
 x^{(1)} \\
 \vdots \\
 x^{(m)}
 \end{array}
 \begin{bmatrix}
 x^{(1)} & \dots & x^{(m)} \\
 \vdots & & \vdots
 \end{bmatrix}$$

① Construct ϕ

$$K(\cdot) = \phi^T \phi$$

O-R-

② Mercer's Theorem

$$\{x^{(1)} \dots x^{(m)}\}$$

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

is P.S.D.

$$\textcircled{3} \int \int f(x) K(x, x') f(x') dx dx' \geq 0$$

SVM

$$y^{(i)} \in \{+1, -\}$$

Parameters $\rightarrow w, b$

$$w \in \mathbb{R}^d \quad b \in \mathbb{R}$$

$$x \in \mathbb{R}^d$$

Functional
Margin

$$\text{Margin: } y^{(i)} (w^T x^{(i)} + b) > 0 \rightarrow \text{always}$$

$$w^T x + b > 0 \text{ for } y = +1$$

$$< 0 \text{ for } y = -1$$

Desire: Margin large

Calculate smallest margin & choose hyperplane which maximizes the smallest margin.

$$w \rightarrow 2w$$

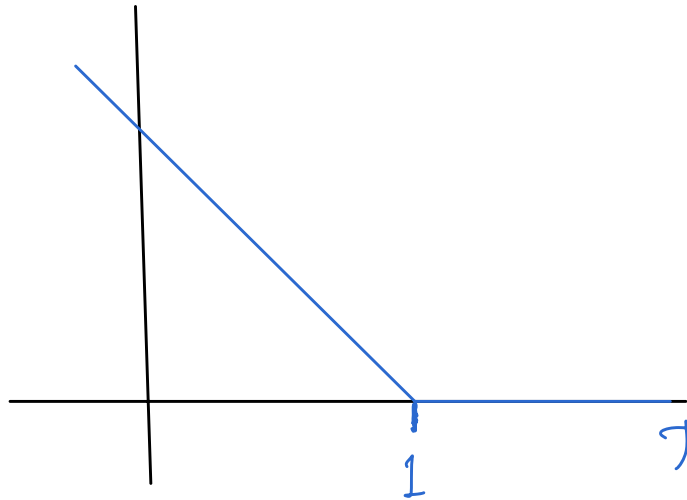
$$b \rightarrow 2b$$

Then margin doubles

So we can game the system,
so use Geometric margin

$$\text{Margin} = \gamma^{(i)}$$

$$\min_{w, b} \left(\sum_{i=1}^n \underbrace{\max[0, 1 - y^{(i)}(\omega^T x^{(i)} + b)]}_{\text{Hinge loss / SVM loss}} \right) + \underbrace{\frac{1}{C} \|\omega\|^2}_{\text{Penalizing}}$$



$$\min_{\xi, w, b} \quad \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{St. } y^{(i)} (\omega^T x^{(i)} + b) \geq 1 - \xi_i \quad \forall i \in \{1 \dots n\}$$

$$\xi_i \geq 0, \quad i = 1 \dots n \quad \text{Primal Convex Problem}$$

Dual Convex Problem

$$\max_x = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$0 \leq \alpha_i \leq C$$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0$$

We don't penalize b because then we don't give our algorithm to be close to the origin

Most α_i 's will be 0, very few non-zero
those set of examples \rightarrow non zero \rightarrow support vector, closest to margin

BAYESIAN METHODS