

Outline:

Kernel Methods

SVM

Feature map

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

$$x \rightarrow \phi(x)$$

fitting $h_\theta(x) = \theta^\top \phi(x)$ using gradient descent
with $\phi(x^{(i)})$ as features

$$\theta = 0$$

Loop

$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^\top \phi(x^{(i)})) \phi(x^{(i)})$$

Issue: $\theta, \phi(x) \in \mathbb{R}^P$ can be very high dimensional
runtime per iteration: $O(nP)$

Goal: improve to $O(n^2)$ per iteration

key observation:

θ can be represented as

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

\uparrow
scalar (n variables instead of P)

Proof by induction

At iteration 0

$$\theta = 0 = \sum_{i=1}^n 0 \cdot \phi(x^{(i)})$$

Suppose at iteration t, $\theta = \sum_{i=1}^n \beta_i \phi(x^{(i)})$

Next iteration:

$$\begin{aligned} \theta &:= \theta + \alpha \sum (y^{(i)} - \theta^\top \phi(x^{(i)})) \phi(x^{(i)}) \\ &\downarrow \\ &\sum \beta_i \phi(x^{(i)}) \end{aligned}$$

$$\hat{\theta} := \sum_{i=1}^n (\beta_i + \alpha(y^{(i)} - \theta^\top \phi(x^{(i)})) \phi(x^{(i)})$$

↑
new β_i

represent $\theta \in \mathbb{R}^P$ implicitly by $\beta \in \mathbb{R}^n$
works better $P \gg n$

Update rule for β

$$\begin{aligned}\beta_i &= \beta_i + \alpha(y^{(i)} - \theta^\top \phi(x^{(i)})) \\ &= \beta_i + \alpha(y^{(i)} - (\sum_{j=1}^n \beta_j \phi(x^{(j)})^\top) \phi(x^{(i)})) \\ &= \beta_i + \alpha(y^{(i)} - \sum_{j=1}^n \beta_j \langle \phi(x^{(j)}), \phi(x^{(i)}) \rangle)\end{aligned}$$

$\theta^\top = \left(\sum_{j=1}^n \beta_j \phi(x^{(j)}) \right)^\top$
 $= \sum_{j=1}^n \beta_j \phi(x^{(j)})^\top$

$$\alpha^\top b = \langle a, b \rangle$$

Observations

- ① $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ if i, j can be precomputed
- ② Often, $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ can be computed faster than $\mathcal{O}(P)$

$$\phi(u) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \\ x_1^2 \\ x_1 x_2 \\ \vdots \\ x_d^2 \\ x_1^3 \\ x_1^2 x_2 \\ \vdots \\ x_d^3 \end{bmatrix} \quad p = 1 + d + d^2 + d^3 = \mathcal{O}(d^3)$$

$$\langle x, z \rangle = \left(\sum_{i=1}^d x_i z_i \right)^2 = \sum_{i=1}^d \sum_{j=1}^d x_i x_j z_i z_j$$

$$\begin{aligned}\langle \phi(x), \phi(z) \rangle &= 1 + \sum_{i=1}^d x_i z_i + \sum_{i=1}^d \sum_{j=1}^d x_i x_j z_i z_j + \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d x_i x_j x_k z_i z_j z_k \\ &= 1 + \langle x, z \rangle + \langle x, z \rangle^2 + \langle x, z \rangle^3\end{aligned}$$

Can compute $\langle \phi(x), \phi(z) \rangle$ in $O(d)$ time!

$$K(x, z) = \langle \phi(x), \phi(z) \rangle \quad \text{Kernel}$$

$$K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

Algo:

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

for all $i, j \in \{1, \dots, n\}$

Set $\beta = 0$

Loop $\beta_L := \beta_L + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j K(x^{(i)}, x^{(j)}) \right)$

preprocessing: $O(n^2 d)$

Each iteration: $O(n^2)$ per iteration

Prediction: Given x , compute $\Theta^T \phi(x)$

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

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

$O(nd)$ time

$x^{(i)}$: i^{th} training example

x : new example

Deeper Observation

The algo only depends on $K(\cdot, \cdot)$

design of features \rightarrow design of Kernel functions

What kernels are valid ??

$$\exists \phi \text{ st. } K(x, z) = \langle \phi(x), \phi(z) \rangle$$

Necessary condition:
n data points $x^{(1)}, \dots, x^{(n)}$

$$\text{Kernel matrix } K \in \mathbb{R}^{n \times n} \quad K_{ij} = K(x^{(i)}, x^{(j)})$$

matrix K is positive semidefinite
 $K \succeq 0$

$$\boxed{k_{ij}}$$

$$K \succeq 0 \Leftrightarrow z^T K z \geq 0$$

$$z^T K z = \sum_i \sum_j z_i K_{ij} z_j$$

$$= \sum_i \sum_j z_i \phi(x^{(i)})^T \phi(x^{(j)}) z_j$$

$$= \sum_l \sum_j z_i \sum_l \phi_l(x^{(i)}) \phi_l(x^{(j)}) z_j$$

$$= \sum_l \left(\sum_i z_i \phi_l(x^{(i)}) \right)^2 \geq 0$$

$$\Rightarrow K \succeq 0$$

Thm (Mercer) K is a valid kernel fn (i.e. $K(z, z) = \phi(z)^T \phi(z)$)
iff for any $n < \infty$ and any $x^{(1)} \dots x^{(n)}$
the corresponding Kernel matrix K st. $K_{ij} = K(x^{(i)}, x^{(j)})$
is positive semidefinite

Other kernels:

$$K(x, z) = (x^T z + c)^2 = \langle \phi(x), \phi(z) \rangle$$

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

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

Gaussian Kernel: $K(x, z) = \exp\left(-\frac{\|x - z\|_2^2}{2\sigma^2}\right)$

$$K(x, z) = \langle \phi(x), \phi(z) \rangle$$

ϕ : infinite dimensional!

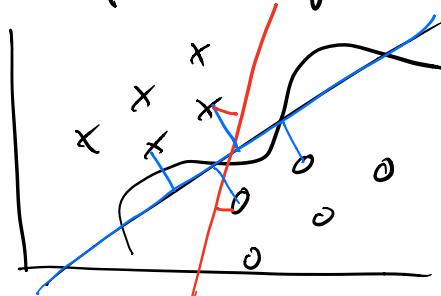
$$x_1, x_1 x_2, \dots, x_1^{16} x_2^{27}, \dots$$

Protein sequence classification
sequences of amino acids (d-T)

$$\begin{bmatrix} \text{AAAA} \\ \text{AAAB} \\ \vdots \\ \text{TTTT} \end{bmatrix} \rightarrow 20^4 \text{ dim vector}$$

Histograms: take min and sum up

SVMs for classification



linear:
 $\{x : w^T x + b = 0\}$

non-linear:
 $\{x : w^T \phi(x) + b = 0\}$

SVM: $\phi(x) = x$ for now

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

Warmup:

find w, b s.t.

$$y^{(i)} = 1 \quad w^T x^{(i)} + b > 0 \quad \textcircled{1}$$

$$y^{(i)} = -1 \quad w^T x^{(i)} + b < 0 \quad \textcircled{2}$$

Many such w, b

New goal: Among all w, b satisfying ① & ②

$$\text{Find } w, b \text{ st. } \max_{w, b} \min_{i \in \{1, \dots, n\}} \text{dist}(x^{(i)}, \text{boundary})$$

$$= \begin{cases} \frac{\mathbf{w}^T \mathbf{x}^{(i)} + b}{\|\mathbf{w}\|_2} & \text{if } \mathbf{w}^T \mathbf{x}^{(i)} + b \geq 0 \\ -\frac{(\mathbf{w}^T \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|_2} & \text{on positive side} \end{cases}$$

$$\max_{w, b} \quad \min_{i \in \{1, \dots, n\}} = \frac{y^{(i)}(w^T x^{(i)} + b)}{\|w\|_2}$$

Scale invariant $(w, b) \rightarrow (100w, 100b)$

wlog. we want to find w, b

$$\text{St. } \min_{i \in \{1, \dots, n\}} y^{(i)} (\mathbf{w}^\top \mathbf{x}^{(i)} + b) = 1 \quad (3)$$

$$\max \frac{1}{\|w\|_2} \equiv \min \|w\|_2$$

men || w||₂

$$\text{s.t. } \forall i \quad y^{(i)}(\omega^T x^{(i)} + b) \geq 1$$

$$\min \|w\|_2 \Rightarrow \min \frac{1}{2} \|w\|_2^2$$

$$\text{s.t. } \# \{ y^{(i)} (\omega^T x^{(i)} + b) \geq 1$$

Facts: ① optimal soln $w^* = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$ for some $\alpha_i \geq 0$

4

② The α in ④ is the optimizer of

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

$$\text{s.t. } \alpha_i \geq 0 \quad \sum \alpha_i y^{(i)} = 0$$

Kernelize: replace $\langle x^{(i)}, x^{(j)} \rangle$ with $K(x^{(i)}, x^{(j)})$
n variables

Fixing linearly separable assumption:

$$y^{(i)} (\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1 - \xi_i$$

$$\min \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i$$

