

L8

- ★ Midterm
- ★ Prep: SVMs & kernels
- ★ Multiplexer perception
- ★ Toward neural network

Midterm: 3/31 Next Tuesday

- ★ close book
 - ★ 90 minutes
 - ★ 12-hour window
 - ★ self-enforced (honor code)
 - ★ lecture 1-8
 - ★ conceptual, algorithm & analysis, code
 - ★ ~6 questions
 - ★ 1-page cheat sheet
 - ★ posted as an assignment, you can scan and upload solution.
-

Support vector machines

- Hinge loss + regularizer
- optimize via dual
- gives a solution that is robust

kernel method

- Nonlinear separable datasets
- kernel trick

$$X = (x^{(1)}, x^{(2)}, \dots, x^{(d)})$$

$$\phi \rightarrow (x^{(1)}, x^{(2)}, \dots, x^{(d)}, \\ x^{(1)}x^{(2)}, \dots, x^{(1)}x^{(d)}, \\ x^{(1)}x^{(2)}x^{(3)}, \dots \\ \exp(x^{(1)}) \dots)$$

$$x \rightarrow \boxed{\phi} \rightarrow \phi(x)$$

"kernel feature mapping"

Challenge: $\phi(x)$ can be very high dimensional

idea: Observe that in most linear models, we access the data via inner products

$\langle \phi(x), w \rangle \rightarrow$ really what we need

$\langle \phi(x_i), \phi(x_j) \rangle \rightarrow$ Dot products in kernel space

$k(x_i, x_j)$ "kernel inner product"

Examples: ① $k(x_i, x_j) = \langle x_i, x_j \rangle$ linear kernel

② $k(x_i, x_j) = (\langle x_i, x_j \rangle)^2$ or $(1 + \langle x_i, x_j \rangle)^2$
quadratic kernel

③ $k(x_i, x_j) = (1 + \langle x_i, x_j \rangle)^p$ polynomial kernel?

④ $k(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|_2^2) = \exp(-\frac{\|x_i - x_j\|_2^2}{6\gamma})$
Gaussian kernel, or Radial basis function / **RBF**

Rules of kernel design

① Efficient computation $k(x_i, x_j)$

should be easy computable

② Symmetry $k(x_i, x_j) = k(x_j, x_i)$

③ $k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$ for some ϕ for all x_i, x_j
 $= (c + \langle \phi(x_i), \phi(x_j) \rangle)^r$

Guaranteed to exist via Mercer's Theorem

[statement: Given any n data points x_1, \dots, x_n

Form the symmetric $n \times n$ matrix M

$$M_{ij} = k(x_i, x_j)$$

Then the eigenvalues of M have to be non-negative

i.e. M has to be positive semi-definite]

经验法则

Thumb rules for constructing kernels

* $k(x, z) = x^T z$ is a kernel

* $k(x, z) = \text{poly}(k(x, z))$ is a kernel

* $k(x, z) = a k_1(x, z) + b k_2(x, z)$ is a kernel

* $k(x, z) = \exp(k(x, z))$ is a kernel

$k(x, z) = f(x) k(x, z) f(z)$ is a kernel

Exercise : show using above rules that

$k(x, z) = \exp(-\|x - z\|_2^2)$ is a valid kernel

Instantiate in a concrete example \rightarrow perceptron

kernel Perceptron

Input : $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

Output : $\text{Sign}(\langle w, \phi(x_i) \rangle) = y_i$

Algorithm :

0. initialize $w_0 \leftarrow 0$

1. Repeat :

For each $(x_i, y_i) \in S$

if $\text{Sign}(\langle w_t, \phi(x_i) \rangle) \neq y_i$

$w_{t+1} \leftarrow w_t + y_i \phi(x_i)$

if no change in w_t after sweeping through data set . exit

until $t \leq \text{max epochs}$

$\langle w_t, \phi(x_i) \rangle \rightarrow ?$ kernel inner product

trick : observe : $w_0 = 0$

$w_1 = y_i \phi(x_i)$ for some i

$w_2 = y_1 \phi(x_1) + y_1' \phi(x_1')$

In general, $W \in \sum_{i \in S} \alpha_i y_i \phi(x_i)$

α_i encodes # times its data point has been visited.

\therefore just a list of α_i 's $\in \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix}$

During test $\text{sign}(W, \phi(x))$

$$= \text{sign} \left(\left\langle \sum_{i \in S} \alpha_i y_i \phi(x_i), \phi(x) \right\rangle \right)$$

$$= \text{sign} \left(\sum_{i \in S} \alpha_i y_i \left\langle \phi(x_i), \phi(x) \right\rangle \right)$$

$$= \text{sign} \left(\sum_{i \in S} \alpha_i y_i K(\phi(x_i), \phi(x)) \right)$$

_____ model

Other algorithms [linear regression, SVMs, ridge regression, ecc.] can all be kernelized.

Summary:

- Establish that final model is a linear combination of data points

$$W = \sum_{i \in S} \alpha_i y_i x_i$$

- Ensure that intermediate updates can be implemented via inner products.
- Replace all occurrence of $\langle x_i, x_j \rangle$ to $K(x_i, x_j)$

kernel nearest neighbours ?

Train: $(x_1, y_1) \dots (x_n, y_n)$

Test: x

$$f(x) = \sum_{i=1}^n y_i \delta(x, x_i)$$

$\delta \rightarrow$ delta function $= \begin{cases} 1 & \text{if } x_i \text{ is nearest of } x \\ 0 & \text{otherwise} \end{cases}$

kernel perceptron/SVM.

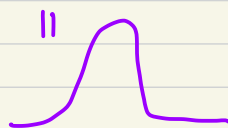
$$f(x) = \text{sgn} \left(\sum_{i=1}^n \alpha_i y_i k(x, x_i) \right)$$

结果更取决于
nearest label.

$x = x_i$ 时
 $k(x, x_i) = 1$

几乎是 KNN 的

If we choose $k(x, x_i) = \exp(-\gamma \|x - x_i\|^2)$



As limit $\gamma \rightarrow \infty$,

$$k(x, x_i) \rightarrow \delta(x, x_i)$$

所以有这种
关系

linear
models

kernel

nearest neighbour

Neural Network

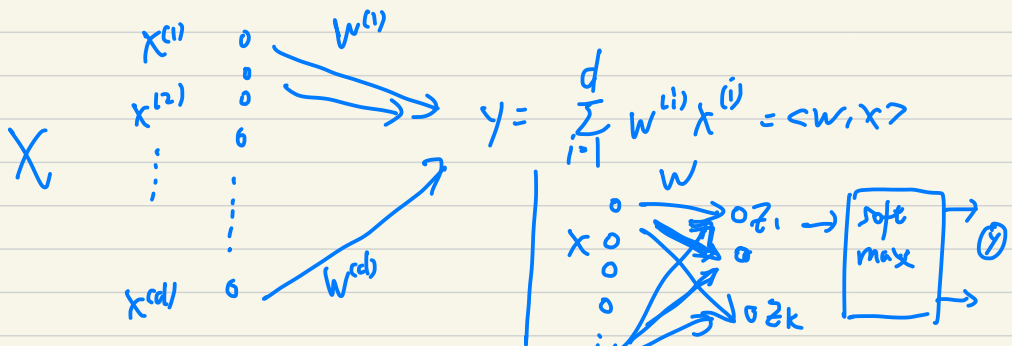
Step back

eg linear regression

- choose a linear model: $\langle w, x \rangle$
- compare with provided label via loss function $\ell(\cdot, \cdot)$ eg.

$$L(y, \langle w, x \rangle) = \frac{1}{2} \|y - w^T x\|_2^2$$

- If this is non-zero, update via grad descent



eg. logistic regression (k-class)

- linear model $\langle w_1, x \rangle$
 $\langle w_2, x \rangle$
 \vdots
 $\langle w_k, x \rangle$

- Softmax: $\hat{y} = \text{softmax}(z)$

$$\hat{y}_i = \frac{\exp(\langle w_i, z \rangle)}{\sum_{i=1}^k \exp(\langle w_k, z \rangle)}$$

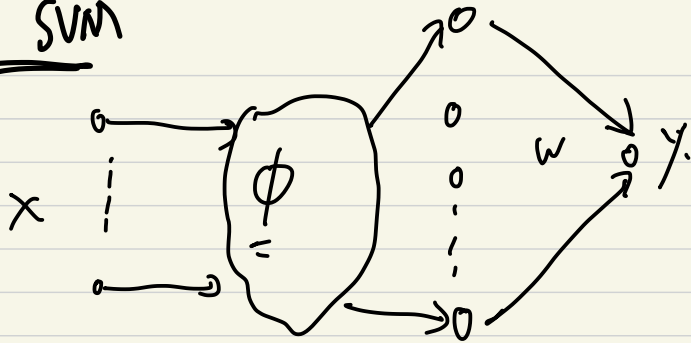
- cross entropy: $L(y, \hat{y}_i)$

$$= - \sum_{i=1}^k y_i \log(\hat{y}_i)$$

$$y = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

- if loss is non-zero, update w . $\langle \phi(x), w \rangle$

kernel SVM



kernel methods: hand chosen

Neural methods : learned feature maps.