

Linear regression:  $y = \theta^T x_{new}$

Cost function:

$$J(\theta) = \sum_{i=1}^N (y^{(i)} - h_{\theta}(x^{(i)}))^2$$

$$= \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$$

also:  $J(\theta) = \left\| \begin{matrix} X \\ (N \times d+1) \end{matrix} \theta - \begin{matrix} Y \\ (d+1, 1) \end{matrix} \right\|_2^2$

$\hat{\theta} = \arg \min_{\theta} J(\theta) \quad \Downarrow \frac{\partial}{\partial \theta}$

$\Rightarrow \nabla J(\theta) = \frac{\partial J(\theta)}{\partial \theta} = 2X^T(X\theta - Y)$

$\Rightarrow X^T(X\theta - Y) = 0 \Rightarrow \hat{\theta} = (X^T X)^{-1} X^T Y$

method 2: (Gradient descent)

$\theta^{(t)} = \theta^{(t-1)} - \rho \left. \frac{\partial J}{\partial \theta} \right|_{\theta^{(t-1)}}$  (learning rate)

$\Rightarrow \theta^{(t)} = \theta^{(t-1)} - \rho X^T(X\theta^{(t-1)} - Y)$

method 3: stochastic Gradient descent:  
 $X$  is a sample of original  $[X]$

$X \Rightarrow \begin{bmatrix} \dots & 1 \\ x & 1 \\ \dots & 1 \end{bmatrix}$   
 $N(t+1)$

$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_d \\ b \end{bmatrix}_{(d+1) \times 1}$

Nonlinear features:  
 $x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \Rightarrow \phi(x) = \begin{pmatrix} x_1 \\ x_1^2 \\ x_2 \\ x_2^2 \\ \vdots \end{pmatrix}$

use  $\phi(x)$  instead of  $x$

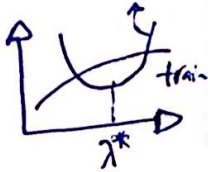
Ridge Regression:

$\hat{\theta}_{\lambda} = (X^T X + \lambda I)^{-1} X^T Y$   
 $\frac{\partial J(\theta)}{\partial \theta} = 0$

Set  $\lambda$ : k-fold cross validation



every time one  $D^{(i)}$  is hold-out set rest is train set. ...  
 find  $\lambda$  with min error:



Maximum Likelihood of observation  
 Training

1) Coin problem:

$x \in \{0, 1\}$   $P(x=1) = \theta$   
 $P(x=0) = 1 - \theta$   
 $P(x) = \theta^x (1 - \theta)^{1-x}$

$\max_{\theta} \log P(D|\theta) = \hat{\theta}_{ML} = \arg \max_{\theta} P(D|\theta)$

$\mathcal{L} = \log P(D|\theta) = \left( \sum_{i=1}^N x^{(i)} \right) \log \theta + \left( N - \sum_{i=1}^N x^{(i)} \right) \log (1 - \theta)$

$\Rightarrow \frac{\partial \mathcal{L}}{\partial \theta} = 0 \Rightarrow \hat{\theta}_{ML} = \frac{\sum_{i=1}^N x_i}{N}$

2) linear regression:

$P(y^{(i)} | x^{(i)}, \theta) = \mathcal{N}(y^{(i)}, \mu = x^{(i)T} \theta, \sigma^2)$   
 $= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)} - x^{(i)T} \theta)^2}{2\sigma^2}}$

$\Rightarrow \max_{\theta} \log P(D|\theta) = \hat{\theta}_{ML}$

$\mathcal{L} = \log P(D|\theta) \Rightarrow \hat{\theta} = \min_{\theta} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$

Generally: Cost function:  
 regularization:  $J(\theta) = \sum_{i=1}^N \ell(y^{(i)} - \theta^T x^{(i)}) + \lambda f(\theta)$

where  $\begin{cases} \ell(e) = e^2, f(\theta) = \|\theta\|_2^2 \rightarrow \text{Ridge} \\ \ell(e) = e^2, f(\theta) = \|\theta\|_1 \rightarrow \text{Lasso} \\ \ell(e) = |e|, \rightarrow \text{Robust} \end{cases}$



## MAP Estimation (Bayesian Est.)

$$\underbrace{P(\theta|D)}_{\text{Posterior}} \propto \underbrace{P(D|\theta)}_{\text{likelihood}} \underbrace{P(\theta)}_{\text{Prior}}$$

$$P(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)} \Rightarrow \begin{cases} \text{mean} = \frac{\alpha}{\alpha+\beta} \\ \text{mode} = \frac{\alpha-1}{\alpha+\beta-2} \end{cases}$$

Beta  $\alpha, \beta$

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} P(\theta|D) \rightarrow \text{small } N$$

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} P(D|\theta) \rightarrow \text{good for large } N$$

$$\begin{array}{c} \hat{\theta}_{\text{ML}} \quad \hat{\theta}_{\text{MAP}} \quad \hat{\theta}_{\text{Prior}} \\ \left\{ \begin{array}{l} N \rightarrow 0 \Rightarrow \hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{Prior}} \\ N \rightarrow \infty \Rightarrow \hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{ML}} \end{array} \right. \end{array}$$

for coin problem:  $P(x) = \theta^x (1-\theta)^{1-x}$

$$\hat{\theta}_{\text{MAP}} = \frac{N}{N+\alpha+\beta-2} \left( \frac{\sum x_i}{N} \right) + \frac{\alpha+\beta-2}{N+\alpha+\beta-2} \left( \frac{\alpha-1}{\alpha+\beta-2} \right)$$

$\hat{\theta}_{\text{ML}} \quad \hat{\theta}_{\text{Prior}}$

classification  $\begin{cases} \text{Generative Modeling: } P(x|y), P(y) \\ \text{Discriminative Modeling: } P(y|x) \end{cases}$

1) Generative:

$$\ell(\theta) = \sum_{i=1}^N \log P(x^{(i)}, y^{(i)} | \theta) =$$

$$= \sum_{i=1}^N \log P(y^{(i)} | \theta) + \sum_{i=1}^N \log P(x^{(i)} | y^{(i)}, \theta)$$

$$\hat{\theta}_j = \frac{\sum_{i=1}^N 1(y=j)}{\sum_{j=1}^K \sum_{i=1}^N 1(y=j)}$$

Naive Bayes with Laplace

$$\theta_{x_k|y} = \frac{\sum_{i=1}^N 1(y=j \wedge x_k = \bar{x}_k) + t}{\sum_{i=1}^N 1(y=j) + mt}$$

Smoothing

Naive Bayes:

$$P(x^{(i)} | y=j) = \prod_{k=1}^K P(x_k^{(i)} | y=j)$$

2) Discriminative Modeling: logistic regression

$$P(y=1|x) = \frac{1}{1 + e^{-w^T \phi(x)}}$$

$$P(y=0|x) = \frac{1}{1 + e^{w^T \phi(x)}}$$

$$\Rightarrow w = \arg \max_w \sum_{i=1}^N \log P(y^{(i)} | x^{(i)}, w) \rightarrow \ell(w)$$

$$\Rightarrow J(w) = -\ell(w)$$

$$\Rightarrow \frac{\partial J}{\partial w} = \sum_{i=1}^N -\phi(x^{(i)}) \left[ y^{(i)} - \frac{1}{1 + e^{-\phi(x^{(i)})^T w}} \right]$$

$$\frac{\partial J}{\partial w} = \sum_{i=1}^N -\phi(x^{(i)}) \left[ y^{(i)} - P(y=1|x^{(i)}, w) \right]$$

$$\Rightarrow \text{Gradient descent: } w^t = w^{t-1} - \rho \frac{\partial J}{\partial w}$$

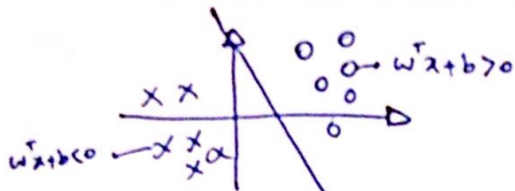
Neural Network update:

$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \frac{\partial J}{\partial z_i^{(l+1)}} \times \frac{\partial z_i^{(l+1)}}{\partial w_{ij}^{(l)}}$$

$$\frac{\partial J}{\partial b_i^{(l)}} = \frac{\partial J}{\partial z_i^{(l+1)}} \times \frac{\partial z_i^{(l+1)}}{\partial b_i^{(l)}}$$



# SVM



$$g(w^T x + b) = \begin{cases} +1 & w^T x + b > 0 \\ -1 & w^T x + b < 0 \end{cases}$$

$$\hat{\gamma}^{(i)} = y^{(i)} (w^T x^{(i)} + b)$$

functional margin

$$\gamma^{(i)} = \left( \frac{w^T x^{(i)} + b}{\|w\|_2} \right) y^{(i)}$$

geometric margin

$$\hat{\gamma} = \min_{i=1, \dots, N} \hat{\gamma}^{(i)}$$

overall functional margin

$$\gamma = \min_{i=1, \dots, N} \gamma^{(i)}$$

overall geometric margin

fix  $\hat{\gamma} = 1$

Primal SVM

$$\begin{cases} \min \|w\|_2^2 \\ \text{s.t. } y^{(i)} (w^T x^{(i)} + b) > 1 \quad \forall i=1, 2, \dots, N \end{cases}$$

generally:

$$\begin{cases} \min f(w) \\ \text{s.t. } g_i(w) \leq 0 \quad i=1, 2, \dots, k \\ h_i(w) = 0, \quad i=1, 2, \dots, L \end{cases}$$

$$\Rightarrow L(w, \alpha, \beta) \triangleq f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^L \beta_i h_i(w)$$

$$\Rightarrow \max_{\beta, \alpha \geq 0} L(w, \alpha, \beta) = f(w) \quad (\text{feasible})$$

if  $g_i \leq 0, h_i = 0$

$$\Rightarrow \min_w f(w) = \min_w \max_{\alpha \geq 0, \beta} L(w, \alpha, \beta)$$

$f_i(w), g_i(w)$  Convex,  $h_i(w)$  affine

$$\Rightarrow \max_{\alpha \geq 0, \beta} \min_w L = \min_w \max_{\alpha \geq 0, \beta} L(w, \alpha, \beta)$$

$$\Rightarrow \boxed{\min_w f(w) = \max_{\alpha \geq 0, \beta} \min_w L(w, \alpha, \beta)}$$

and

$$\frac{\partial L}{\partial w} = 0 \quad \frac{\partial L}{\partial \alpha} = 0$$

$$\frac{\partial L}{\partial \beta} = 0 \quad \alpha_i g_i(w) = 0 \quad g_i(w) \leq 0 \quad h_i(w) = 0$$

back to SVM primal:

$$\Rightarrow \frac{\partial L}{\partial w} = 0 \Rightarrow w^* = \sum_{i=1}^N \alpha_i^* y^{(i)} x^{(i)}$$

Dual SVM

$$\begin{aligned} \max_{\alpha_i, \alpha_i \geq 0} \quad & -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle + \sum_{i=1}^N \alpha_i \\ \text{s.t.} \quad & \alpha_i \geq 0 \quad \forall i=1, \dots, N \\ & \sum_{i=1}^N \alpha_i y^{(i)} = 0 \end{aligned}$$

$$b^* = - \frac{\max_{y^{(i)} = -1} w^{*T} x^{(i)} + \min_{y^{(i)} = 1} w^{*T} x^{(i)}}{2}$$

So for Test:

$$g(x) = \text{sgn}(w^T x + b^*) \quad \text{or} \quad k(x, x)$$

$$\text{or } g(x) = \text{sgn}\left(\sum_{i=1}^N \alpha_i^* y^{(i)} \langle x^{(i)}, x \rangle + b^*\right)$$

Kernels:  $K(\text{Point 1}, \text{Point 2}) = \text{number}$

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

\* if  $\phi(x) = \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix} \rightarrow K(x, z) = (x^T z)^2$

\* if  $\phi(x) = \begin{pmatrix} x_1^n \\ x_1^{n-1} x_2 \\ \vdots \\ x_2^n \end{pmatrix} \rightarrow \text{all monomials of degree } n$

$$\Rightarrow K(x, z) = (x^T z)^n$$

\* if  $\phi(x) = \begin{pmatrix} x_1^n \\ x_1^{n-1} x_2 \\ \vdots \\ x_2^n \end{pmatrix} \rightarrow \text{all monomials of deg. } n \text{ or less}$

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

\* if  $\phi(x) = \infty \text{ dim} \Rightarrow K(x, z) = \frac{-\|x - z\|_2^2}{2\sigma^2} e^{\frac{-\|x - z\|_2^2}{2\sigma^2}}$

So (No need) to calculate  $\phi(x), \phi(z)$

Kernel Matrix: for  $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$

$$\begin{bmatrix} K_{11} & \dots & K_{1N} \\ \vdots & & \vdots \\ K_{N1} & \dots & K_{NN} \end{bmatrix} \rightarrow K_{ij} = K(x^{(i)}, x^{(j)})$$

$K$  is kernel Matrix iff only if

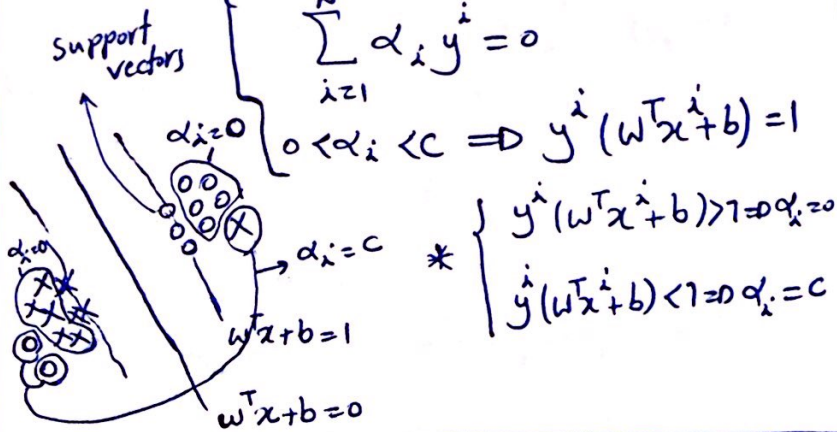
- 1)  $K = K^T$
- 2)  $K$  is PSD  $z^T K z \geq 0 \quad \forall z$



## SVM - Soft Margin:

Primal:  $\begin{cases} \min_{(w,b)} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N f_i \\ \text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1 - f_i, \forall i=1, \dots, N \\ f_i \geq 0, \forall i=1, \dots, N \end{cases}$

Dual:  $\begin{cases} \max_{\alpha_1, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y^{(i)} y^{(j)} k(x^{(i)}, x^{(j)}) \alpha_i \alpha_j \\ \text{s.t. } 0 \leq \alpha_i \leq C, \forall i=1, \dots, N \\ \sum_{i=1}^N \alpha_i y^{(i)} = 0 \end{cases}$



Ensemble Method:  $x \rightarrow \mathcal{P}(x)$  causes overfitting  
 $\hookrightarrow$  ensemble Methods are greater!

Boosting:

(1) For  $t=1, \dots, T$ : \* use weak classifier ( $h$ )

(2) if  $t=1$  then:

(3) Initialize weights  $D_{t,i} = \frac{1}{N}$

(4) else

(5) set weights to  $D_{t,i} \propto D_{t-1,i} e^{-\alpha_{t-1} y^{(i)} h_{t-1}(x^{(i)})}$

(6) End if

(7) Train a hypothesis  $h_t$  using  $D_t$

(8) Evaluate Training Error  $\epsilon_t = \sum_{i=1}^N D_{t,i} 1(h_t(x^{(i)}) \neq y^{(i)})$

(9) set  $\alpha_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$

(10) End for

(11)  $H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right) \Rightarrow$  Final classifier

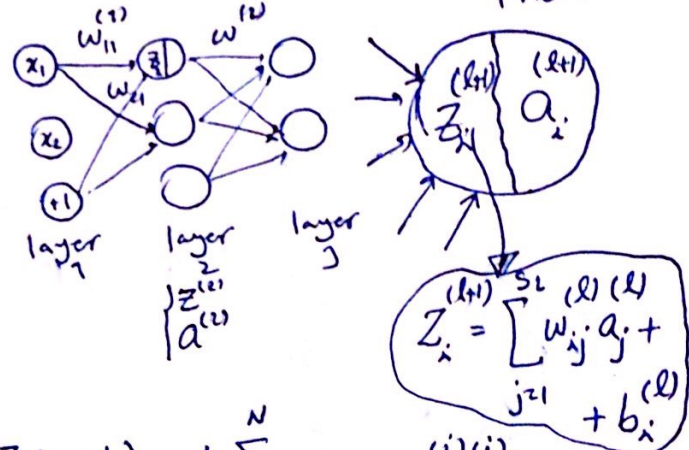
Bagging: ① No weight, random sampling with replacement

② uniform weight for hypothesis:  $\alpha_t = \frac{1}{T}$

in Boosting: (1)  $\epsilon_t \rightarrow 0 \Rightarrow \alpha_t \rightarrow \infty$ ,  $\epsilon_t = 0.5 \Rightarrow \alpha_t \rightarrow 0$

(2) if  $\begin{cases} h_{t-1}(x^{(i)}) y^{(i)} > 0 \xrightarrow{\text{correct}} D_{t-1,i} \times \text{small weight} \\ h_{t-1}(x^{(i)}) y^{(i)} < 0 \xrightarrow{\text{incorrect}} D_{t-1,i} \times \text{large weight} \end{cases}$

## Neural Network:



$$J_\lambda(w, b) = \frac{1}{N} \sum_{i=1}^N J(w, b, x^{(i)}, y^{(i)}) + \frac{\lambda}{2} \sum_{l=1}^{n-1} \|w^{(l)}\|_F^2$$

where  $J(w, b, x^{(i)}, y^{(i)}) = \left\| \frac{h(x^{(i)}; w, b)}{w, b} - y^{(i)} \right\|_2^2$

$\Rightarrow$  Gradient Descent:

$$w^{(l)} \leftarrow w^{(l)} - \alpha \frac{\partial J_\lambda}{\partial w} \bigg|_{w^{(l)}}$$

$$b^{(l)} \leftarrow b^{(l)} - \alpha \frac{\partial J_\lambda}{\partial b} \bigg|_{b^{(l)}}$$

Forward Propagation:  $\{w^{(l)}, b^{(l)}\}_{l=1}^{n-1}$  fixed

$$z^{(l+1)} = w^{(l+1)} a^{(l)} + b^{(l+1)}$$

$$a^{(l+1)} = \sigma(z^{(l+1)}); \forall l=1, 2, \dots, n-1$$

Back Propagation:  $\{z^{(l)}, a^{(l)}\}_{l=1}^n$  Computed

$$\delta^{(n)} = (a^{(n)} - y) \odot \sigma'(z^{(n)})$$

For  $l=n-1, n-2, \dots, 2, 1$

$$\textcircled{1} \text{ Set: } \delta^{(l)} = (w^{(l+1)} \delta^{(l+1)}) \odot \sigma'(z^{(l)})$$

$$\textcircled{2} \frac{\partial J}{\partial w} \bigg|_{w^{(l)}} = \delta^{(l+1)} a^{(l)T}$$

$$\textcircled{3} \frac{\partial J}{\partial b} \bigg|_{b^{(l)}} = \delta^{(l+1)}$$

$$\text{Update } w^{(l)} \leftarrow w^{(l)} - \alpha \left[ \frac{1}{N} \frac{\partial J}{\partial w} \bigg|_{w^{(l)}} + \lambda w^{(l)} \right]$$

$$b^{(l)} \leftarrow b^{(l)} - \alpha \left[ \frac{1}{N} \frac{\partial J}{\partial b} \bigg|_{b^{(l)}} \right]$$



## Principal Component Analysis (PCA)

(1) input :  $D = \{x^{(i)}\}_{i=1}^N = [X]$

(2) Output :  $[U]_{n \times d}$ ,  $[\mu]_{n \times 1}$ ,  $[y]_{d \times N}$

(3) Compute  $[\bar{x}] = \frac{1}{N} \sum_{i=1}^N x^{(i)}$

(4) Form :  $X = [x^{(1)} - \bar{x}, x^{(2)} - \bar{x}, \dots, x^{(N)} - \bar{x}]$

(5) Compute SVD of  $X$

$$X = U_X \Sigma_X V_X^T \quad U_X \in \mathbb{R}^{n \times n}$$

(6) Optimal  $U = U_X(:, 1:d)$

⑦ Principal Components :  $[y] = U^T (x - \bar{x})$

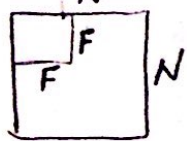
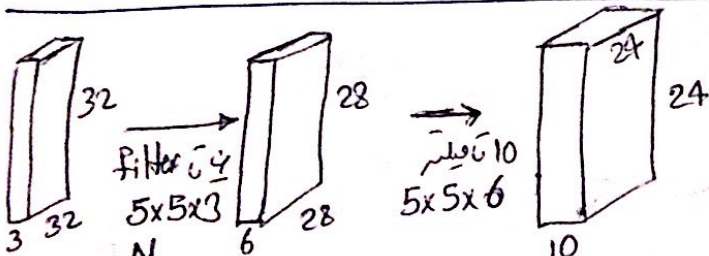
$$\text{PCA } \begin{cases} \min \frac{1}{2} \sum_{i=1}^N \|x^{(i)} - U y - \mu\|_2^2 \\ \text{s.t. } \sum_{i=1}^N y^{(i)} = 0, U U^T = I_d \end{cases}$$

⇒ Then we write  $L(\dots)$

$$\text{Then } \frac{\partial L}{\partial \mu} = 0 \Rightarrow \mu = \frac{1}{N} \sum_{i=1}^N x^{(i)}$$

$$\frac{\partial L}{\partial y_i} = 0 \Rightarrow y^{(i)} = U^T (x^{(i)} - \bar{x})$$

Plug in  $y^{(i)}$  &  $\mu$ , ... optz ✓✓



$$\Delta \text{ output size} = (N - F) / \text{stride} + 1$$

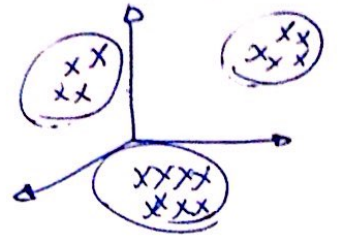
\* For every 0 pooling add 2 to N

→ input 32x32x3

$$10 \text{ filter } 5 \times 5 \times 3 \Rightarrow \# \text{Parameters} = (5 \times 5 \times 3 + 1) \times 10 = 760$$

bias

## Clustering



Cost Function :

$$J = \sum_{i=1}^N \min_{j=1,2,\dots,k} \|x^{(i)} - c_j\|_2^2$$

$$\Rightarrow \text{opt}_Z: \begin{cases} \min \sum_{i=1}^N \sum_{j=1}^k z_{ij} \|x^{(i)} - c_j\|_2^2 \\ \text{s.t. } z_{ij} \in \{0,1\}, \forall i,j; \sum_{j=1}^k z_{ij} = 1 \quad \forall i \end{cases}$$

where:  $z_{ij} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is assigned to } c_j \\ 0 & \text{else} \end{cases}$

Solve optz. with Coordinate Descent:

KMEANS

(1) Initialize  $\{c_j^{(0)}\}_{j=1}^k$  and  $\{z_{ij}^{(0)}\}_{i=1, j=1 \dots k}^N$

(2) For  $t=1,2,\dots$  do :

► For every  $i$

$$z_{ij}^{(t)} = \begin{cases} 1 & \text{if } j = \arg \min_{l=1,\dots,k} \|x^{(i)} - c_l^{(t-1)}\|_2 \\ 0 & \text{else} \end{cases}$$

► For every  $j$

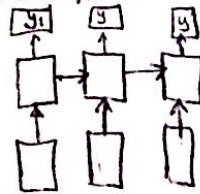
$$c_j^{(t)} = \frac{\sum_{i=1}^N 1(z_{ij}^{(t)} = 1) x^{(i)}}{\sum_{i=1}^N 1(z_{ij}^{(t)} = 1)}$$

(3) Output:  $\{c_j^{(t)}\}, \{z_{ij}^{(t)}\} \quad \forall i,j$

K-means Converges (Cost function always decreasing.)

$$J(\{z_{ij}^{(t)}\}, \{c_j^{(t)}\}) \geq J(\{z_{ij}^{(t+1)}\}, \{c_j^{(t+1)}\})$$

recurrent N.N



$$h_t = f_W(h_{t-1}, x_t)$$

ex:

$$h_t = \tanh(W_{hh} h_{t-1} + W_{hx} x_t)$$

$$y_t = W_{hy} h_t$$



$$\frac{P(A \cap B)}{P(B)} = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) P(B) = P(B|A) P(A) = P(A \cap B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$A, B \text{ indepe} \Rightarrow P(A \cap B) = P(A)P(B)$$

### Distributions

$$\textcircled{1} \text{ Gaussian: } f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2}$$

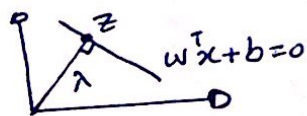
$$\textcircled{2} \text{ Bernoulli: } f(x; p) = p^x (1-p)^{1-x} \quad 0 < p < 1, x \in \{0, 1\}$$

$$\textcircled{3} \text{ Uniform: } f(x; a, b) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

$$\textcircled{4} \text{ Exponential Dist: } f(x, \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0, \lambda > 0 \\ 0 & x < 0 \end{cases}$$

$$\textcircled{5} \text{ Poisson Dist: } f(x, \lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \quad \lambda > 0$$

Distance from origin to hyperplane:



$$\lambda = \frac{|b|}{\|w\|_2}$$

$$\vec{z} = \lambda \frac{w}{\|w\|_2}$$

$$k(x, z) = c_1 k_1(n, z) + c_2 k_2(n, z) \quad \leftarrow \text{عج}$$

$$k_1(n, z) = \phi_1^T(x) \phi_1(z)$$

$$k_2(n, z) = \phi_2^T(x) \phi_2(z)$$

$$\Rightarrow k(x, z) = c_1 \phi_1^T(x) \phi_1(z) + c_2 \phi_2^T(x) \phi_2(z)$$

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

$$\text{where } \Phi = \begin{bmatrix} \sqrt{c_1} \phi_1(x) \\ \sqrt{c_2} \phi_2(x) \end{bmatrix} \quad \checkmark \checkmark$$

$$K(x, z) = k_1(n, z) k_2(n, z) \quad \leftarrow \text{عج}$$

$$k_1(n, z) \rightarrow \phi_1(x) = [f_1(n), f_2(n), \dots]$$

$$k_2(n, z) \rightarrow \phi_2(x) = [g_1(n), g_2(n), \dots]$$

$$K(n, z) = k_1(n, z) k_2(n, z) =$$

$$= \phi_1^T(x) \phi_1(z) \phi_2^T(x) \phi_2(z)$$

$$= \sum_{i=1}^{\infty} f_i(x) f_i(z) \sum_{j=1}^{\infty} g_j(x) g_j(z) =$$

$$= \sum_{i,j} (f_i(x) g_j(x)) (f_i(z) g_j(z))$$

$$= \sum_{i,j} c_{i,j}(x) c_{i,j}(z)$$

✓✓

Positive semi-definite ( $A_{nn}$ )

if  $\boxed{z^T A z \geq 0}$  for  $\forall z$

$$\Rightarrow \boxed{A = B^T B}$$

$$P(x_1, \dots, x_N) = \sum_{c_j} \prod_i P(x_i | c_j) P(c_j)$$

Naïve Bayes