

M.L.  $\Rightarrow$  Error Minimization (or) Likelihood Maximization

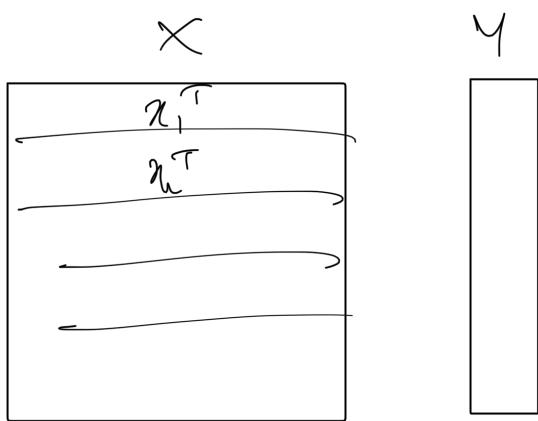
$$\text{Likelihood function} \Rightarrow P(x|\mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

For this to be a probability function,

$$\int_{-\infty}^{\infty} P(x|\mu, \sigma^2) dx = 1$$

$$\begin{aligned} \text{MLE} &\Rightarrow P(D|\theta) \\ \text{MAP} &\Rightarrow P(\theta|D) \\ &\Leftarrow P(\theta|\theta) \propto P(\theta) \end{aligned}$$

If prior is also included, it is MAP.



Maximum Likelihood Estimation

$\Rightarrow$  learn  $\theta$  such that

$P(Y|X, \theta)$  is maximized

Likelihood Function

Coin Toss

Bernoulli Distribution  $\Rightarrow n=1$

Likelihood is like a score which tells how likely a dataset is generated by a model.

$$m^n (1-m)^{(1-n)} \Rightarrow n=0 \text{ or } 1$$

i.i.d  $\Rightarrow$  independent identically distributed

$\bigcup n$  coin tosses

$$P(X|\mu) = \prod_{i=1}^n P(x_i|\mu) \Rightarrow \text{Joint Likelihood}$$

When many P's are multiplied  $\Rightarrow$  answer is small.  $\Rightarrow$  Underflow error

So use log likelihood.

$$\text{Joint log likelihood} = \log(P(X(\mu))) = \log\left(\prod_{i=1}^n P(x_i|\mu)\right)$$

$$= \sum_{i=1}^N \log P(n_i|\mu)$$

$$\sum_{i=1}^N \log(P(n_i|\mu)) = \sum_{i=1}^N \log(\mu^{x_i} (1-\mu)^{1-x_i})$$

Find  $\mu$  such that this is maximized

$$= \sum_{i=1}^N [x_i \log \mu + (1-x_i) \log(1-\mu)] \Rightarrow JLL$$

$$\frac{\partial(JLL)}{\partial \mu} = \frac{\partial}{\partial \mu} \left[ \left( \sum_{i=1}^N x_i \right) \log \mu + \left( \sum_{i=1}^N (1-x_i) \right) \log(1-\mu) \right]$$

$$\sum_{i=1}^N x_i = N_H$$

$$= \frac{\partial}{\partial \mu} \left[ N_H \log \mu + (N - N_H) \log(1-\mu) \right]$$

$$= \frac{N_H}{\mu} - \frac{N - N_H}{1-\mu} = \frac{N_H - N_H \mu - N \mu + N_H \mu}{\mu(1-\mu)}$$

$$\frac{\partial(JLY)}{\partial \mu} = \frac{N_H - N \mu}{\mu(1-\mu)}$$

Let  $\frac{\partial \text{JLL}}{\partial \mu} = 0$

$$\frac{N_H - N \times \mu}{\mu(1-\mu)} = 0 \Rightarrow N_H - N \times \mu = 0$$

$$\mu = \frac{N_H}{N}$$

We know this formula. But it was formally derived like this

## Dice

"K" faces

$x$  = k-dimensional vector where 1st outcome is 1 & others are zeros. ( $i$ -indexed)

(Or)

$n$  = a no. b/w 1...k.

$P(x|\mu) \Rightarrow$  Here  $\mu$  is a k-dimensional vector where  $\mu[i]$  represents probability of dice showing  $i$ .

$$P(x|\mu) = \mu_1^{I(n=1)} \mu_2^{I(n=2)} \mu_3^{I(n=3)} \dots \mu_k^{I(n=k)}$$

$$P(n|\mu) = \prod_{i=1}^k \mu_i^{I(n=i)} \xrightarrow{\text{or}} \underbrace{\mu_1 \mu_2 \dots \mu_k}_{\text{kth element of vector "n".}}$$

$$\text{If } n = [0 \ 1 \ 0 \ 0] \Rightarrow I(n=1) = I(n=3) = I(n=4) = 0 \\ \therefore I(n=2) = 1$$

Let there be  $N$  samples like these,

$$x_1, x_2, x_3, \dots, x_N$$

$x_n = n^{\text{th}}$  sample.



$$J.C. = P(X|\mu) = \prod_{n=1}^N P(x_n|\mu)$$

$$\begin{aligned}
 J.L.L &= \log(P(X|\mu)) = \sum_{n=1}^N \log(P(x_n|\mu)) \\
 &= \sum_{n=1}^N \log \left[ \frac{\prod_{i=1}^k \mu_i^{x_{ni}}}{\prod_{i=1}^k \mu_i^{N_i}} \right] = \sum_{n=1}^N \sum_{i=1}^k x_{ni} \times \log \mu_i \\
 &\quad - \lambda \left[ \sum_{i=1}^k \mu_i - 1 \right]
 \end{aligned}$$

(Lagrangian  $\hookrightarrow$  constraint  
 As, sum of probs. = 1.

$$\frac{\partial (J.L)}{\partial \mu_i} = \sum_{n=1}^N \sum_{i=1}^k \frac{x_{ni}}{\mu_i} - \lambda = 0 \Rightarrow \frac{N_i}{\mu_i} = \lambda$$

$$\lambda = \frac{N_i}{\mu_i} \quad (\text{or}) \quad \mu_i = \frac{N_i}{\lambda}$$

$$\sum_{n=1}^N \sum_{i=1}^k x_{ni} = ?$$

$$\sum_{i=1}^k \mu_i = 1 \Rightarrow \frac{\sum_{i=1}^k N_i}{\lambda} = 1 \Rightarrow \lambda = \sum_{i=1}^k N_i \Rightarrow \boxed{\lambda = N}$$

Posterior  $\propto$  likelihood \* prior

$$P(\theta|X, \alpha) \propto P(X|\theta) P(\theta|\alpha)$$

$\hookrightarrow$  hyperparameters for  $\theta$

For coin toss,

$$P(X|\theta) = P(X|\mu) = \prod_{i=1}^N \mu^{n_i} (1-\mu)^{l-n_i}$$

$$\beta(\theta|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \times \theta^{a-1} \times [1-\theta]^{b-1} \Rightarrow \text{Beta function}$$

For bernoulli,

$$\beta(\mu|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \times \mu^{N_1+a-1} \times (1-\mu)^{N-N_1+b-1}$$

Conjugate prior  $\Rightarrow$  When prior doesn't change its functional form when converted to posterior

For dice game,

Dirichlet distribution  $\Rightarrow$  Distribution of distribution.  
- Distribution of prob. vectors over discrete space.

$$P(\theta|x,\alpha) \propto P(x|\theta) \times P(\theta|\alpha)$$

$$\propto \prod_{i=1}^N \Gamma(n_i + \alpha_i)$$

$$\propto \prod_{i=1}^N \prod_{j=1}^k \alpha_j^{n_{ij}}$$

$$k \text{ faces} \Rightarrow \sum_{i=1}^k \alpha_i = 1$$

$$\times \alpha_1, \alpha_2, \alpha_3, \dots, \alpha_k \geq 0$$

$$[\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_k] \approx \sum_{i=1}^k \alpha_i$$

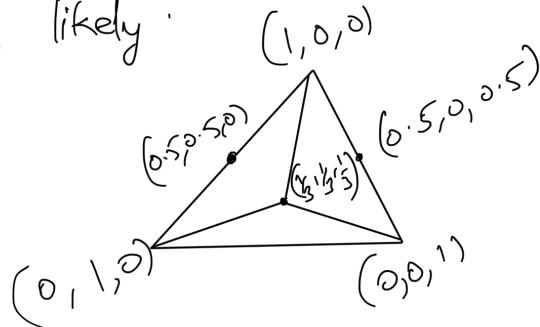
$\Gamma \Rightarrow$  gamma function

(extension of factorial function for real nos.)

$$\text{Dir}(\alpha) = \frac{\prod_{i=1}^k \alpha_i}{\prod_{i=1}^k \Gamma(\alpha_i)} \times \prod_{i=1}^k \alpha_i^{\alpha_i - 1}$$

$$\text{Posterior} \Rightarrow \frac{\prod_{i=1}^k (\alpha_i + N_i)}{\prod_{i=1}^k (\alpha_i + N_i)} \times \prod_{i=1}^k \alpha_i^{N_i + \alpha_i - 1}$$

$\alpha = 1 \Rightarrow$  Every point in the simplex of  $n - D$  data are equally likely.



$\alpha \uparrow \Rightarrow$  More uniform distribution

$\alpha \downarrow \Rightarrow$  Some events less likely, while, some are more likely.

↳ prior. Decides which kind of samples generated.

Naive Bayes  $\Rightarrow$  Probabilistic classifier

SVM  $\Rightarrow$  Non-probabilistic classifier

$\downarrow$   
feature conditional probabilities are independent of each other.

$$\hat{P}(y_i | x, \theta) \Rightarrow \hat{y} = \max_{i=1 \dots L} (\hat{P}(y_i | x, \theta)) \Rightarrow \text{Probabilistic classifier}$$

Generative classifier (approach)  
Discriminative approach

}

$\Rightarrow$  lead

$\hookrightarrow$  Directly modeling conditional probabilities

> Works on joint distribution  $\Rightarrow$  by product is classifier

$$P(y, X) = P(y) \times P(X|y)$$

↓  
prior      ↓  
likelihood

$$P(X|y) = \prod_{i=1}^K P(x_i|y)$$

Generative models are easy to implement. Work well on smaller data sets.

$$J.L. = \prod_{n=1}^D P(y_n, x_n) = \prod_{n=1}^D P(y_n) \times P(x_n|y_n)$$

If  $W = [0 \ 1 \ 0 \ 0 \ | \ 0 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ \dots]$  → No. of words in Corpus

$\downarrow$  word-0 absent     $\downarrow$  word-1 present

$\downarrow$  word-2 absent     $\downarrow$  no. of times this occurred

&  $x = [0 \ 1 \ 0 \ 0 \ 3 \ 0 \ 2 \ 0 \ 1 \ 5 \ 9 \ 0 \ \dots]$

$$P(x|y) = \prod_{i=1}^K P(w_i|y)^{x_{ni}}$$

So, if  $|V|=30,000$ , then you need 30,000 parameters to model it. May overfit.

$$J.L. = \prod_{n=1}^D P(y_n) \times \prod_{i=1}^K P(w_i|y_n)^{x_{ni}}$$

Assume words are conditionally independent given  $y$ .

$$P(y=l) = \frac{1}{m} \sum_{i=1}^m I(y_{il}=1)$$

$L \times K \Rightarrow$  model size

$L^{th}$  row &  $i^{th}$  column  $\Rightarrow P(w_i|y=l)$

$$P(w_i|y=l) = \frac{N_{li}}{N_l} \rightarrow \text{No. of times } i^{th} \text{ feature occurs in class } l$$

$\Rightarrow$  Total no. of occurrences of all features in class  $l$ .

Like try to build a model of what malignant tumors look like & another model of what benign tumors look like.

Generative

- Naive Bayes
- Builds separate models to characterize each class.

- Learns  $P(x|y)$ .  
Also  $P(y)$ .

Then Compute  $P(y=k|n)$

$$P(y=k|n) = P(n|y=k) P(y=k) / P(n)$$

$P(n) = \sum_{k=1}^K P(n|y=k) P(y=k)$

What are the features like when  $y=0$ .

Discriminative

- logistic regression
- Builds a model to separate the classes.
- Learns  $P(y|x)$  directly.

# Logistic Regression Binary classification

$$P(y|x, \theta)$$

$$y = \theta^T x$$

$$\hookrightarrow (-\infty, \infty)$$

↓ map to

$[0, 1] \Rightarrow$  (use sigmoid function)

$$\sigma(n) = \frac{1}{1+e^{-n}}$$

$$y = \sigma(\theta^T x) \Rightarrow P(y|x, \theta) = \frac{1}{1+e^{-\theta^T x}}$$

- Doesn't assume conditional independence of conditional prob. of features.

$\hat{y}_i \Rightarrow$  predicted.

$$E_p = (y_i - \hat{y}_i)^2$$

$$E = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$\min_{\theta} E \rightarrow$  to fit the model

$$\min_{\theta} E = \min_{\theta} \frac{1}{N} \sum_{n=1}^N [y_n - \sigma(\theta^T x_n)]^2$$

$$\frac{\partial E}{\partial \theta} = \frac{1}{N} \sum_{n=1}^N 2 [y_n - \sigma(\theta^T x_n)] \times \sigma(\theta^T x_n) \times (1 - \sigma(\theta^T x_n))$$

$$\frac{\partial}{\partial \theta} (\theta^T b) = \begin{bmatrix} \frac{\partial (\theta^T b)}{\partial \theta_1} \\ \frac{\partial (\theta^T b)}{\partial \theta_2} \\ \vdots \end{bmatrix}$$

$x \quad x_n$   
→ vector

$k=1$

i.e. estimates the probability distribution of the dataset.  
→ calculates the probability by counting.

$$a^T b = a_1 b_1 + a_2 b_2 + a_3 b_3 + \dots + a_k b_k$$

$$\frac{\partial(a^T b)}{\partial a_1} = b_1, \quad \frac{\partial(a^T b)}{\partial a_2} = b_2, \quad \dots$$

$$\therefore \begin{bmatrix} \frac{\partial(a^T b)}{\partial a_1} \\ \frac{\partial(a^T b)}{\partial a_2} \\ \vdots \\ \frac{\partial(a^T b)}{\partial a_k} \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{bmatrix} = b$$

Derive  $\Rightarrow \frac{\partial}{\partial x} (x^T A x)$

possible only if  $A$  is a matrix; if  $x$  is a vector.

		$y$	$x_1$	$x_2$
1			1.5	2.8
0			5	6
1			3.4	6
0			7.6	8

$$\omega^T$$

-1	1.5	0.5
----	-----	-----

$\downarrow$        $\downarrow$        $\downarrow$

$\theta_0$      $\theta_1$      $\theta_2$

$$\omega = \omega - \alpha \frac{\partial E}{\partial \omega}$$

Demonstrate gradient descent for one step. Take  $\alpha = 1$ .

$$\chi_0 = 1$$

$$\text{loss function} = L = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

$$w^{t+1} = w^t - \alpha \times \frac{\partial L}{\partial w}$$

$$L = \frac{1}{4} \left[ (1 - \sigma(2.65))^2 + (0 - \sigma(9.5))^2 + (1 - \sigma(7.1))^2 + (0 - \sigma(14.4))^2 \right]$$

$$L = \frac{1}{4} \left[ (1 - 0.934)^2 + (0 - 0.999)^2 + (1 - 0.999)^2 + (0 - 0.9999)^2 \right] = 0.5$$

$$\frac{\partial L}{\partial w} = \frac{1}{N} \sum_{n=1}^N 2 [y_n - \sigma(w^T x_n)] \times \sigma(w^T x_n) \times (1 - \sigma(w^T x_n)) \times x_n$$

$$\frac{\partial L}{\partial w} = \frac{1}{4} \times 2 \times \left[ ((1 - \sigma(-1 + 1.5 \times 1.5 + 0.5 \times 2.8)) \times \chi_1 - (-1 + 1.5 \times 1.5 + 0.5 \times 2.8)) \times \chi_2 \times ((1 - \sigma(-1 + 1.5 \times 1.5 + 0.5 \times 2.8))) \times \chi_3 + (0 - \dots - ) \times \chi_4 \right]$$

$$\sigma(w^T x_1) = \sigma(2.65) = 0.934$$

$$\sigma(w^T x_2) = \sigma(9.5) = 0.999$$

$$\sigma(w^T x_3) = \sigma(7.1) = 0.999$$

$$\sigma(w^T x_4) = \sigma(14.4) = 0.999$$

$$\frac{\partial L}{\partial w} = \frac{1}{2} \times \left[ ((1 - 0.934) \times 0.934 \times (1 - 0.934) \times \begin{bmatrix} 1.5 \\ 2.8 \end{bmatrix} \right. \\ + ((0 - 0.999) \times 0.999 \times (1 - 0.999) \times \begin{bmatrix} 5 \\ 6 \end{bmatrix} \left. \right] \\ + ((1 - 0.999) \times 0.999 \times (1 - 0.999) \times \begin{bmatrix} 3.4 \\ 6 \end{bmatrix} \left. \right] \\ + ((0 - 0.999) \times 0.999 \times (1 - 0.999) \times \begin{bmatrix} 7.6 \\ 8 \end{bmatrix} \left. \right]$$

$$= \frac{1}{2} \times \left\{ 0.004068 \times \begin{bmatrix} 1 \\ 1.5 \\ 2.8 \end{bmatrix} - 0.000998 \begin{bmatrix} 1 \\ 5 \\ 6 \end{bmatrix} \right. \\ \left. + 0.999 \times 10^{-6} \begin{bmatrix} 1 \\ 3.4 \\ 6 \end{bmatrix} \right. \\ \left. - 0.000998 \begin{bmatrix} 1 \\ 7.6 \\ 8 \end{bmatrix} \right\} \\ = \frac{1}{2} \times \begin{bmatrix} 0.006 \\ 0.0114 \end{bmatrix} - \begin{bmatrix} 0.005 \\ 0.006 \end{bmatrix} + \begin{bmatrix} 3.4 \times 10^{-6} \\ 6 \times 10^{-6} \end{bmatrix} \\ - \begin{bmatrix} 0.0025848 \\ 0.008 \end{bmatrix} \\ = \begin{bmatrix} 0.00207196 \\ -0.0065 \\ -0.0026 \end{bmatrix} = \begin{bmatrix} 0.00103649 \\ -0.00325 \\ -0.0013 \end{bmatrix}$$

$$\omega = \omega - \alpha \times \frac{\partial L}{\partial w} \Rightarrow \omega = \begin{bmatrix} -1 \\ 1.5 \\ 0.5 \end{bmatrix} - \begin{bmatrix} 0.00103649 \\ -0.00325 \\ -0.0013 \end{bmatrix}$$

$$\omega = \begin{bmatrix} -1.0036249 \\ 1.50325 \\ 0.5013 \end{bmatrix}$$

MSE is not used as optimization objective of logit regression. Instead we compute KL divergence b/w  $y$  &  $\hat{y}$ .

$$KL(\vec{P} \parallel \vec{Q}) = \sum_{i=1}^L p_i \log \frac{p_i}{q_i}$$

$$\begin{aligned} L &= \frac{1}{N} \sum_{n=1}^N KL(y_n \parallel \hat{y}_n) = \frac{1}{N} \sum_{n=1}^N \sum_{l=1}^L y_{nl} \log \frac{y_{nl}}{\hat{y}_{nl}} \\ &= \frac{1}{N} \sum_{n=1}^N -\log \left[ p(y_n=l|x_n, w) \right] \end{aligned}$$

$$p(y_n=l|x_n, w) = \sigma(w^\top x_n)$$

Softmax ( $l, n, w$ ) =  $\frac{e^{w_l^\top x}}{\sum_{i=1}^L e^{w_i^\top x}}$

$\downarrow$

$\begin{bmatrix} w_1^\top \\ w_2^\top \\ \vdots \\ w_L^\top \end{bmatrix}$

$\downarrow$

$p(y=l|n)$

$\Downarrow P_l$

$\exp(\text{Score 1}) \Rightarrow w_1^\top x$

$\exp(\text{Score 2}) \Rightarrow w_2^\top x$

$\exp(\text{Score 3}) \Rightarrow w_3^\top x$

$\vdots$

$\exp(\text{Score } l) \Rightarrow w_l^\top x$

Some may be negative.  
So apply  $\exp$

$\hookrightarrow$  This gives  $[P_1 \ P_2 \ \dots \ P_L]$

$$KL((y_1 \ y_2 \ \dots \ y_L) \parallel (P_1 \ P_2 \ \dots \ P_L))$$

$$= \sum_{l=1}^L y_l \log \frac{y_l}{P_l} = -\log p(y|x, w)$$

$$= \sum_{l=1}^L y_l \log y_l - \sum_{l=1}^L y_l \log P_l$$

$\hookrightarrow$  Minimize this  $\Rightarrow \log p(y|x, w)$   
increases

$$d_{f_m}(u_v) = \frac{vdu - udv}{\sqrt{v^2}}$$

$$\frac{\partial P_l}{\partial w_{CK}} = \frac{\partial [\text{softmax}(l, n, w)]}{\partial w_{CK}} = \frac{\partial}{\partial w_{CK}} \cdot \frac{e^{w_l^T n}}{\sum_{i=1}^L e^{w_i^T n}} = \frac{(\Sigma) \times \frac{\partial}{\partial w_{CK}} (e^{w_l^T n})}{(\Sigma)^2}$$

$\frac{\partial}{\partial w_{CK}} (e^{w_l^T n})$   
 $= -x e^{w_l^T n} \times e^{w_C^T n} \times n_K$   
 $\quad \left( \sum_{i=1}^L e^{w_i^T n} \right)^2$   
 $+ \frac{I(l=c) \times e^{w_l^T n} \times n_C}{\sum_{i=1}^L e^{w_i^T n}}$   
 $= -P_l \times P_c \times n_K + I(l=c) \times P_l \times n_C$   
 $= P_l (-P_c + I(l=c)) n_K$

$$\boxed{\frac{\partial P_l}{\partial w_C} = P_l (I(l=c) - P_c) n}$$

$$E = KL((y_1, y_2, \dots, y_n) || (p_1, p_2, \dots, p_n))$$

$$= \sum_{l=1}^L y_l \log \frac{y_l}{p_l} = -\log(p(y/n, w))$$

$$E = \frac{1}{N} \sum_{n=1}^N KL(y_n || p_n)$$

*One hot vector*

$y_{nl} \log y_{nl} - y_{nl} \log p_{nl}$   
 Here  $w_{CK}$ ,  
 so not considered

$$w_{CK}^{t+1} = w_{CK}^t - \alpha \frac{\partial E}{\partial w_{CK}} \Rightarrow \text{Gradient Descent.}$$

$$\frac{\partial E}{\partial w_{CK}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial w_{CK}} \left[ \sum_{l=1}^L y_{nl} \times \log \frac{y_{nl}}{p_{nl}} \right] = -\frac{1}{N} \sum_{n=1}^N \sum_{l=1}^L I(y_{nl}=1) \frac{\partial}{\partial w_{CK}} (\log(p_{nl}))$$

$$= -\frac{1}{N} \sum_{n=1}^N \sum_{l=1}^L I(y_{nl}=1) \times \frac{1}{p_{nl}} \times \frac{\partial}{\partial w_{CK}} (p_{nl})$$

$$= -\frac{1}{N} \sum_{n=1}^N \sum_{l=1}^L I(y_{nl}=1) \times \frac{1}{p_{nl}} \times p_{nl} \times (I(c=1) - P_c) n_K$$

$$\frac{\partial E}{\partial w_{ck}} = -\frac{1}{N} \sum_{n=1}^N \sum_{l=1}^L I(y_{nl}=1) (I(c=l) - p_{nc}) \times x_{nk}$$

$$I(c=l) = 1, \text{ if } c = l$$

No. of parameters in logistic =  $L + L \times K$   
 ↳ bias

No. of parameters in naive bayes =  $L + L \times K$

changing the gradient  
 as  $w_{clc}^{+f'} = w_{clc}^+ - \alpha \left( \frac{\partial E}{\partial w_{ck}} H^{-1} \right)$   
 ↳ prior prob.  $H = \frac{\partial^2 E}{\partial w_{ck}^2} \Rightarrow$  Hermitian

→ Naive Bayes assumes conditional probs. to be independent.  
 So logistic fares better than naive bayes.

Indirect way of solving the classification problem  
 by naive bayes also makes it less suitable to classification  
 when compared to logistic regression.

Generative algorithms make more assumptions; whereas discriminative learning algs make lesser assumptions. If assumptions are wrong, then D.A.'s fair better.

G.A.  $\Rightarrow$  computationally efficient than D.A.'s  
 But more data present  $\Rightarrow$  D.A.'s preferred

↓  
 Helps in  
 faster  
 convergence  
 ↓  
 Called  
 E-BFGS

## E-M Algorithm

Suppose there are documents from three classes Cricket, Hockey & Football. But they aren't labelled.

$$\begin{bmatrix} & \\ & \\ & \vdots \\ & \end{bmatrix} \begin{bmatrix} T_1 \\ T_2 \\ \vdots \\ T_n \end{bmatrix}$$

E: Desired posterior probs. for each document

$$\begin{pmatrix} P_{11} & P_{12} & \dots & \dots & P_{1K} \\ P_{21} & P_{22} & \dots & \dots & P_{2K} \\ P_{31} & P_{32} & \dots & \dots & P_{3K} \end{pmatrix}$$

M: Use posterior prob. to derive updated parameters.

## Solve AI's EM example

Eg: Let there be two dice with three faces (A, B, C)

$$\text{Dice 1} \Rightarrow \begin{matrix} A & B & C \\ 0.3 & 0.4 & 0.3 \end{matrix} \quad \text{Update } P(A/\text{Dice 1}), \dots$$

Dice 2  $\Rightarrow 0.1, 0.4, 0.5$

C C A B B C  $\Rightarrow T_1$   
 A B A A A B  $\Rightarrow T_2$   
 A C C B C A  $\Rightarrow T_3$   
 A A B B C B  $\Rightarrow T_4$   
 A C B C B C  $\Rightarrow T_5$

If sequence of throws is given, then probability is the multiplication of prior.

If sequence is not given, then apply binomial distribution, multinomial distribution, etc.

	P(D1)	P(D2)	No. of A's attr. to D1	No. of B's attr. to D1	No. of C's attr. to D1	No. of A's attr. to D2	B's to D2	C's to D2
T1	0.3932	0.607	0.3932	0.7864	1.1796	0.607	1.24	1.82
T2	0.988	0.0122	3.952	1.976	0	0.0488	0.0245	0
T3	0.66	0.34	1.32	0.66	1.98	0.68	0.34	1.02
T4	0.84	0.16	1.68	2.52	0.84	0.32	0.48	0.16
T5	0.39	0.61	0.39	0.78	1.17	0.61	1.22	1.83

$$P(A/D1) = \frac{\text{No. of A's attr. to D1}}{\text{Total no. of rolls with D1}} = \frac{7.7352}{19.6272} = 0.394$$

$$P(A/D2) = \frac{\text{No. of A's attr. to D2}}{\text{Total no. of rolls with D2}} = \frac{2.2658}{10.3752} = 0.2183$$

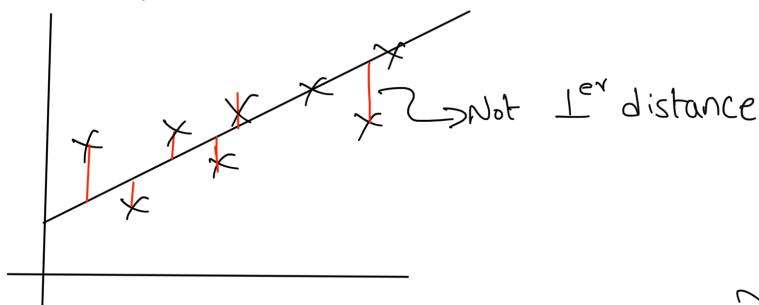
# Linear Regression

$y$  = Real-valued variable

$X \Rightarrow N \times K$  data matrix

$$\hat{y} = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_K x_K = \alpha^T x$$

$$E = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2$$



2 dimensions

$$E = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2 = \frac{1}{N} \sum_{n=1}^N (y_n - m x_n - c)^2$$

$$= d_1 m^2 + d_2 m c + d_3 c^2 + d_4 m + d_5 c + d_6$$

$\underbrace{\sum_{n=1}^N x_n}_{\text{m}}$        $\underbrace{\sum_{n=1}^N y_n^2}_{\text{c}}$

$$\frac{\partial E}{\partial m} = 2d_1 m + d_2 c + d_4 = 0$$

$$\frac{\partial E}{\partial c} = d_2 m + 2d_3 c + d_5 = 0$$

$$d_2 \times (2d_1 m + d_2 c + d_4) = 0$$

$$2d_1 \times (d_2 m + 2d_3 c + d_5) = 0$$

$$\Rightarrow \frac{2d_1 d_2 m + d_2^2 c + d_2 d_4}{2d_1 d_2 m + 4d_1 d_3 c + 2d_1 d_5} = 0$$

$$\Rightarrow c = \frac{2d_5 d_1 - d_4 d_2}{d_2^2 - 4d_1 d_3}$$

$$m = \frac{d_2 d_5 - 2 d_1 d_4}{4 d_1 d_3 - d_2^2}$$

$$m = \frac{\left( \sum x_n \right) \times \left( \sum y_n \right) - N \left[ - \sum x_n y_n \right]}{N \left( \sum x_n^2 \right) - \left( \sum x_n \right)^2}$$

$$m = \frac{N \sum x_n y_n - \left( \sum x_n \right) \left( \sum y_n \right)}{N \left( \sum x_n^2 \right) - \left( \sum x_n \right)^2} \quad \Rightarrow \text{2nd order}$$

K-d  $\Rightarrow E = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2 = \frac{1}{N} (y - x\alpha)^T (y - x\alpha)$

$$E = \frac{1}{N} (y^T y - (x\alpha)^T y - y^T (x\alpha) + (x\alpha)^T (x\alpha))$$

$$E = \frac{1}{N} (y^T y - 2 y^T (x\alpha) + \alpha^T (x^T x) \alpha)$$

$$\frac{\partial E}{\partial \alpha} = \frac{1}{N} \left[ 2(x^T x \alpha) - 2(y^T x)^T \right] = 0$$

$\frac{\partial}{\partial \alpha} (x^T x \alpha) = 2 x^T x$   
 $\frac{\partial}{\partial \alpha} (y^T x)^T = 2 y^T x$   
 $\frac{\partial}{\partial \alpha} (\alpha^T \alpha) = 2 \alpha$

$$\alpha = (x^T x)^{-1} x^T y$$

$$\text{Regularizer} \Rightarrow \beta \alpha^T \alpha$$

$$\begin{aligned}
 \frac{\partial (\beta \alpha^T \alpha)}{\partial \alpha} &= \frac{\partial (\beta \alpha^T I \alpha)}{\partial \alpha} \\
 &= \beta I 2 I \alpha \\
 &= 2 \beta \alpha
 \end{aligned}$$

$$\frac{\partial E \text{ with regularizer}}{\partial \alpha} \Rightarrow 2 x^T x \alpha - 2 (y^T x)^T + 2 \beta \alpha = 0$$

$$(x^T x + \beta I) \alpha = x^T y$$

$$\alpha = (x^T x + \beta I)^{-1} x^T y$$

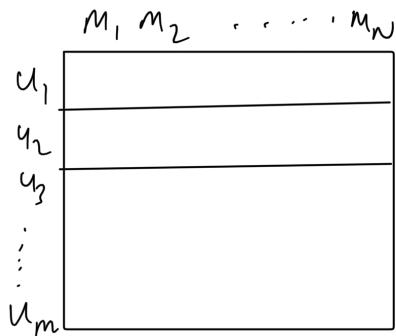
Entropy

# Recommender Systems

large dataset, but sparse.

① Latent factor model

↳ Each user vector to a lower dimension



→ Recommender Systems are based on one of two strategies ⇒ content filtering & collaborative filtering.

Read Netflix paper. Read eigen vectors, SVD etc.

Read Recommender System with SVD

## Problem

	M <sub>1</sub>	M <sub>2</sub>	M <sub>3</sub>	M <sub>4</sub>	M <sub>5</sub>	⇒ indicates no rating
U <sub>1</sub>	1	0	4	1	0	
U <sub>2</sub>	0	0	4	4	2	
U <sub>3</sub>	1	2	0	0	1	
U <sub>4</sub>	0	0	5	2	5	
U <sub>5</sub>	0	4	0	2	5	

P → weights for user  
Q → Movie features

$$Q = \begin{bmatrix} -0.5 & -0.91 \\ 0.04 & -0.9 \\ 0.23 & -0.5 \\ 0.89 & -0.53 \\ -0.92 & -0.91 \end{bmatrix} \quad P = \begin{bmatrix} 0.57 & -0.48 \\ -0.35 & -0.27 \\ 0.17 & -0.84 \\ 0.49 & -0.37 \\ 0.59 & -0.67 \end{bmatrix}$$

User1 & User3 have rated Movie1.

First fix  $P$  & find  $\theta$ .  
Then fix  $\theta$  & find  $P$ .

$$P = \begin{bmatrix} 0.57 & -0.48 \\ 0.17 & -0.84 \end{bmatrix}, R = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

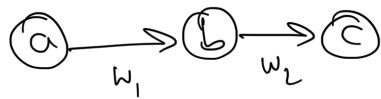
$$Q = (P^T P)^{-1} P^T R$$

$$\theta = \begin{bmatrix} \frac{5553}{10000} & \frac{5001}{10000} \\ \frac{5001}{10000} & \frac{1469}{10000} \end{bmatrix}^{-1} P^T R$$

$$Q = \begin{bmatrix} \frac{15045775}{1971098} \\ -\frac{4596125}{657366} \end{bmatrix} = \begin{bmatrix} 7.63 \\ -6.99 \end{bmatrix} \hookrightarrow q_1$$

After updating  $Q$ , use the updated(new)  $Q$  to solve for  $P$ .

# Neural Networks



$t = \text{target}$

$$E = (c - t)^2$$

$$b = f(w_1 a)$$

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial c} \times \frac{\partial c}{\partial g} \times \frac{\partial g}{\partial w_2}$$

$$c = g(w_2 b)$$

$$= 2(c - t) \times g(w_2 b) (1 - g(w_2 b)) \times b$$

$$f = \text{Tanh}$$

$$g = \text{sigmoid}$$

$$\begin{aligned} \frac{\partial E}{\partial w_1} &= \frac{\partial E}{\partial c} \times \frac{\partial c}{\partial g} \times \frac{\partial g}{\partial b} \times \frac{\partial b}{\partial f} \\ &\quad \times \frac{\partial (w_1 a)}{\partial w_1} \end{aligned}$$

- Not convex optimization
- So at best we get local optima only.

$y = f(a), y \approx \hat{y}$ ,  
To approximate "f",  
that function "f",  
we use a complex  
neural network, that  
adjusts the weights  
to make  $y \approx \hat{y}$ .

$$\begin{aligned} &= 2(c - t) \times g(w_2 b) (1 - g(w_2 b)) \\ &\quad \times w_2 \times (1 - \tanh^2(w_1 a)) \\ &\quad \times a \end{aligned}$$

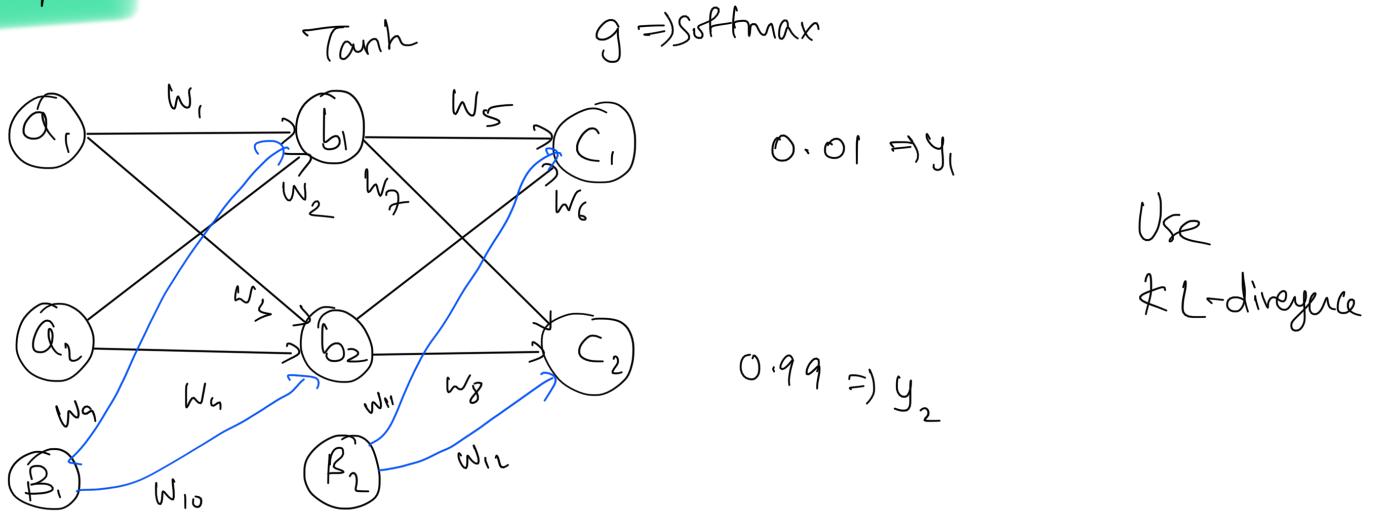
- Binary classification  $\Rightarrow$  Sigmoid Activation  $\Rightarrow$  Near to o/p layer

Tanh  $\Rightarrow$  Middle layers.

( $\hookrightarrow$  -1 to 1  
symmetric about 0.)

$$w_2^{t+1} = w_2^t - \alpha \left( \frac{\partial E}{\partial w_2} \right)$$

## Example



$$E = \sum_{i=1}^2 y_i \log \frac{y_i}{c_i}$$

$$b_1 = \tanh((w_1 a_1 + w_2 a_2) + w_9 \beta_1)$$

$$b_2 = \tanh((w_3 a_1 + w_4 a_2) + w_{10} \beta_1)$$

$$[c_1 \ c_2] = \text{softmax} ((w_5 b_1 + w_6 b_2) + w_{11} \beta_2, (w_7 b_1 + w_8 b_2) + w_{12} \beta_2)$$

$$E = y_1 \log \frac{y_1}{c_1} + y_2 \log \frac{y_2}{c_2}$$

$$E = y_1 \log y_1 - y_1 \log c_1 + y_2 \log y_2 - y_2 \log c_2$$

$$c_1 = \frac{e^{w_5 b_1 + w_6 b_2 + w_{11} \beta_2}}{e^{w_5 b_1 + w_6 b_2 + w_{11} \beta_2} + e^{w_7 b_1 + w_8 b_2 + w_{12} \beta_2}}$$

$$c_2 = \frac{e^{w_7 b_1 + w_8 b_2 + w_{12} \beta_2}}{e^{w_5 b_1 + w_6 b_2 + w_{11} \beta_2} + e^{w_7 b_1 + w_8 b_2 + w_{12} \beta_2}}$$

$$\begin{aligned} \frac{\partial E}{\partial w_5} &= \frac{\partial E}{c_1} \times \frac{\partial c_1}{\partial w_5} + \frac{\partial E}{c_2} \times \frac{\partial c_2}{\partial w_5} \\ &= \frac{-y_1}{c_1} \times \frac{(e^{w_5 \dots} + e^{w_7 \dots}) \times (e^{w_5 \dots} \times b_1) - (e^{w_5 \dots})(e^{w_7 \dots} \times b_1)}{(e^{w_5 \dots} + e^{w_7 \dots})^2} \\ &\quad - \frac{y_2}{c_2} \times \frac{(e^{w_7 \dots} + e^{w_8 \dots}) \times 0 - (e^{w_7 \dots})(e^{w_8 \dots} \times b_2)}{(e^{w_7 \dots} + e^{w_8 \dots})^2} \end{aligned}$$

$$\frac{\partial E}{\partial w_5} = -\frac{y_1}{c_1} \times (c_1 - c_1^2 b_1) + \frac{y_2}{c_2} \times c_1 c_2 b_1$$

$$\frac{\partial E}{\partial w_5} = -y_1 \times b_1 \times (1 - c_1) + y_2 c_1 b_1 \quad \text{---}$$

$$\frac{\partial E}{\partial w_6} = \frac{\partial E}{\partial c_1} \times \frac{\partial c_1}{\partial w_6} + \frac{\partial E}{\partial c_2} \times \frac{\partial c_2}{\partial w_6}$$

$$= -\frac{y_1}{c_1} \times \frac{(e^{w_5} + e^{\tilde{w}_5})(e^{w_5} \times b_2) - (e^{\tilde{w}_5})(e^{w_5} \times b_2)}{(e^{\tilde{w}_5} + e^{\tilde{w}_5})^2}$$

$$-\frac{y_2}{c_2} \times \frac{(e^{w_7} + e^{\tilde{w}_7})(e^{w_7} \times b_1) - (e^{\tilde{w}_7})(e^{w_7} \times b_1)}{(e^{\tilde{w}_7} + e^{\tilde{w}_7})^2}$$

$$\frac{\partial E}{\partial w_6} = -\frac{y_1}{c_1} \times (c_1 b_2 - c_1^2 b_2) + \frac{y_2}{c_2} \times c_1 c_2 b_2$$

$$\frac{\partial E}{\partial w_6} = -y_1 b_2 (1 - c_1) + y_2 c_1 b_2 \quad \text{---}$$

$$\frac{\partial E}{\partial w_7} = \frac{\partial E}{\partial c_1} \times \frac{\partial c_1}{\partial w_7} + \frac{\partial E}{\partial c_2} \times \frac{\partial c_2}{\partial w_7}$$

$$= -\frac{y_1}{c_1} \times \frac{(e^{w_7} + e^{\tilde{w}_7})(e^{w_7} \times b_1) - (e^{\tilde{w}_7})(e^{w_7} \times b_1)}{(e^{w_7} + e^{\tilde{w}_7})^2}$$

$$-\frac{y_2}{c_2} \times \frac{(e^{w_7} + e^{\tilde{w}_7})(e^{w_7} \times b_1) - (e^{\tilde{w}_7})(e^{w_7} \times b_1)}{(e^{w_7} + e^{\tilde{w}_7})^2}$$

$$\frac{\partial E}{\partial w_7} = + \frac{y_1}{c_1} \times c_1 c_2 b_1 - \frac{y_2}{c_2} \times (c_2 b_1 - c_1^2 b_1)$$

$$\frac{\partial E}{\partial w_7} = y_1 c_2 b_1 - y_2 b_1 (1 - c_2) \quad \text{---}$$

$$\begin{aligned}\frac{\partial E}{\partial w_8} &= \frac{\partial E}{\partial c_1} \times \frac{\partial c_1}{\partial w_8} + \frac{\partial E}{\partial c_2} \times \frac{\partial c_2}{\partial w_8} \\ &= -\frac{y_1}{c_1} \times \frac{(e^n + e^r) \times 0 - (e^{ws..}) (e^{w_7 ..} \times b_2)}{(e^n + e^r)^2} \\ &\quad - \frac{y_2}{c_2} \times \frac{(e^n + e^r) (e^{w_7 ..} \times b_2) - (e^{ws..}) (e^{w_7 ..} \times b_2)}{(e^n + e^r)^2} \\ &= + \frac{y_1}{c_1} c_1 c_2 b_2 - \frac{y_2}{c_2} (c_2 b_2 - c_1^2 b_2)\end{aligned}$$

$$\frac{\partial E}{\partial w_8} = y_1 c_2 b_2 - y_2 b_2 (1 - c_2) \quad \text{---}$$

$$\begin{aligned}\frac{\partial E}{\partial w_{11}} &= \frac{\partial E}{\partial c_1} \times \frac{\partial c_1}{\partial w_{11}} + \frac{\partial E}{\partial c_2} \times \frac{\partial c_2}{\partial w_{11}} \\ &= -\frac{y_1}{c_1} \times \frac{(e^n + e^r) (e^{ws..} \times \beta_2) - (e^{ws..}) (e^{ws..} \times \beta_2)}{(e^n + e^r)^2} \\ &\quad - \frac{y_2}{c_2} \times \frac{(e^n + e^r) \times 0 - (e^{ws..}) (e^{ws..} \times \beta_2)}{(e^n + e^r)^2} \\ &= -\frac{y_1}{c_1} (c_1 \beta_2 - c_1^2 \beta_2) + \frac{y_2}{c_2} c_1 c_2 \beta_2\end{aligned}$$

$$\frac{\partial E}{\partial w_{11}} = -y_1 \beta_2 (1 - c_1) + y_2 c_1 \beta_2 \quad \text{---}$$

$$\begin{aligned}\frac{\partial E}{\partial w_{12}} &= \frac{\partial E}{\partial c_1} \times \frac{\partial c_1}{\partial w_{12}} + \frac{\partial E}{\partial c_2} \times \frac{\partial c_2}{\partial w_{12}} \\ &= -\frac{y_1}{c_1} \times \frac{(e^N + e^N)(0) - (e^{w_5}) (e^{w_7} \times \beta_2)}{(e^N + e^N)^2} \\ &\quad - \frac{y_2}{c_2} \times \frac{(e^N + e^N) (e^{w_7} \times \beta_2) - (e^{w_7}) (e^{w_5} \times \beta_2)}{(e^N + e^N)^2} \\ &= +\frac{y_1}{c_1} \times c_1 c_2 \beta_2 - \frac{y_2}{c_2} \times (c_2 \beta_2 - c_1^2 \beta_2)\end{aligned}$$

$$\frac{\partial E}{\partial w_{12}} = y_1 c_2 \beta_2 - y_2 \beta_2 (1 - c_2)$$

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial c_1} \times \frac{\partial c_1}{\partial b_1} \times \frac{\partial b_1}{\partial w_1} + \frac{\partial E}{\partial c_2} \times \frac{\partial c_2}{\partial b_1} \times \frac{\partial b_1}{\partial w_1}$$

$$\frac{\partial c_1}{\partial b_1} = \frac{(e^N + e^N) (e^{w_5} \times w_5) - (e^{w_5}) (e^{w_5} \times w_5 + e^{w_7} \times w_7)}{(e^N + e^N)^2}$$

$$\frac{\partial c_1}{\partial b_1} = c_1 w_5 - c_1 (c_1 w_5 + c_2 w_7)$$

$$\frac{\partial c_2}{\partial b_1} = \frac{(e^N + e^N) (e^{w_7} \times w_7) - (e^{w_7}) (e^{w_5} \times w_5 + e^{w_7} \times w_7)}{(e^N + e^N)^2}$$

$$\frac{\partial c_2}{\partial b_1} = c_2 w_7 - c_2 (c_1 w_5 + c_2 w_7)$$

$$\frac{\partial b_1}{\partial w_1} = (1 - \tanh^2(w_1 a_1 + w_2 a_2 + w_9 \beta_1)) \times a_1$$

$$\frac{\partial R}{\partial w_1} = -\frac{y_1}{c_1} \times (c_1 w_5 - c_1 (c_1 w_5 + c_2 w_7)) \times \frac{\partial b_1}{\partial w_1} - \frac{y_2}{c_2} \times (c_2 w_7 - c_2 (c_1 w_5 + c_2 w_7)) \times \frac{\partial b_1}{\partial w_1}$$

$$\frac{\partial E}{\partial w_1} = \left( -y_1 (w_5 - c_1 w_5 - c_2 w_7) - y_2 (w_7 - c_1 w_5 - c_2 w_7) \right) \times (1 - \tanh^2(w_1 a_1 + w_2 a_2 + w_9 \beta_1)) \times a_1$$

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial c_1} \times \frac{\partial c_1}{\partial b_1} \times \frac{\partial b_1}{\partial w_2} + \frac{\partial E}{\partial c_2} \times \frac{\partial c_2}{\partial b_1} \times \frac{\partial b_1}{\partial w_2}$$

$$\frac{\partial E}{\partial w_2} = \left( -y_1 (w_5 - c_1 w_5 - c_2 w_7) - y_2 (w_7 - c_1 w_5 - c_2 w_7) \right) \times (1 - \tanh^2(w_1 a_1 + w_2 a_2 + w_9 \beta_1)) \times a_2$$

$$\frac{\partial E}{\partial w_9} = \left( -y_1 (w_5 - c_1 w_5 - c_2 w_7) - y_2 (w_7 - c_1 w_5 - c_2 w_7) \right) \times (1 - \tanh^2(w_1 a_1 + w_2 a_2 + w_9 \beta_1)) \times \beta_1$$

$$\frac{\partial E}{\partial w_3} = \frac{\partial E}{\partial c_1} \times \frac{\partial c_1}{\partial b_2} \times \frac{\partial b_2}{\partial w_3} + \frac{\partial E}{\partial c_2} \times \frac{\partial c_2}{\partial b_2} \times \frac{\partial b_2}{\partial w_3}$$

$$\frac{\partial c_1}{\partial b_2} = \frac{(e^{w_5} + e^{w_7})(e^{w_5} \times w_6) - (e^{w_7}) (e^{w_5} \times w_6 + e^{w_7} \times w_8)}{(e^{w_5} + e^{w_7})^2}$$

$$\frac{\partial c_1}{\partial b_2} = c_1 w_6 - c_1 (c_1 w_6 + c_2 w_8)$$

$$\frac{\partial c_2}{\partial b_2} = \frac{(e^{w_5} + e^{w_7})(e^{w_7} \times w_8) - (e^{w_7}) (e^{w_5} \times w_6 + e^{w_7} \times w_8)}{(e^{w_5} + e^{w_7})^2}$$

$$\frac{\partial c_2}{\partial b_2} = c_2 w_8 - c_2 (c_1 w_6 + c_2 w_8)$$

$$\frac{\partial b_2}{\partial w_3} = (1 - \tanh^2(w_1 a_1 + w_2 a_2 + w_{10} \beta_1)) \times a_1$$

$$\frac{\partial E}{\partial w_3} = -\frac{y_1}{c_1} \times (c_1 w_6 - c_1 (c_1 w_6 + c_2 w_8)) \times \frac{\partial b_2}{\partial w_3} - \frac{y_2}{c_2} \times (c_2 w_8 - c_2 (c_1 w_6 + c_2 w_8)) \times \frac{\partial b_2}{\partial w_3}$$

$$\frac{\partial E}{\partial w_3} = \left( -y_1(w_6 - c_1 w_6 - c_2 w_8) - y_2(w_8 - c_1 w_6 - c_2 w_8) \right) \times (1 - \tanh^2(w_3 a_1 + w_4 a_2 + w_{10} \beta_1)) \times a_1$$

$$\frac{\partial E}{\partial w_4} = \left( -y_1(w_6 - c_1 w_6 - c_2 w_8) - y_2(w_8 - c_1 w_6 - c_2 w_8) \right) \times (1 - \tanh^2(w_3 a_1 + w_4 a_2 + w_{10} \beta_1)) \times a_2$$

$$\frac{\partial E}{\partial w_{10}} = \left( -y_1(w_6 - c_1 w_6 - c_2 w_8) - y_2(w_8 - c_1 w_6 - c_2 w_8) \right) \times (1 - \tanh^2(w_3 a_1 + w_4 a_2 + w_{10} \beta_1)) \times \beta_1$$

## Gaussian Distribution

$$x_1, x_2, \dots, x_N, x \sim P(x|\mu, \sigma^2)$$

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i, \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

$$P(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi} \sigma} \times e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

random variable  
location of the peak  
width of the curve  
how far  $x$  is from  $\mu$  [how likely  $x$  is]

$x$ 's closer to  $\mu$  have higher probability.

## Multivariate

$$\underbrace{[x_1, x_2, \dots, x_D]}^T \rightarrow x \text{ in } D\text{-dimensional space}$$

↳ If this is a mobile rating for each feature.

Eg:  $[3, 4, 2, 1]$   
 ↳ price  
 ↳ camera  
 ↳ battery  
 ↳ display

Model these vectors as multivariate gaussian dist. as features may be correlated. → Price depends on display, camera, etc.

- Multivariate gaussian dist. is useful when modelling numerical data with multiple dimensions & all dimensions are continuous valued.

$$P(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}-\boldsymbol{\mu})}$$

trace of a matrix =  $\sum_{i=1}^D \lambda_i$

multiplication of eigen values of  $\Sigma$ .

also a  
D-dimensional vector

covariance matrix of size  $D \times D$ .

Mahalanobi's distance

is a measure of the distance between a point and a distribution. It is a generalization of the Euclidean distance, taking into account the correlations b/w the variables. It is calculated as the square root of the difference b/w the vector representation of an observation & the mean of the distribution, scaled by the covariance matrix of the distribution.

Commonly used in ML to detect outliers, classification & quantifying the similarity b/w two points. The distance can be interpreted as the no. of standard deviations separating the point from the mean of the distribution, taking into account the correlation b/w the variables.

When  $\Sigma^{-1} = I \Rightarrow$  then  $\Delta^2 = (\mathbf{x}-\boldsymbol{\mu})^T (\mathbf{x}-\boldsymbol{\mu})$

$\hookrightarrow$  Euclidean Distance

→ Usually in gaussian dist.  $\Sigma^{-1}$  is symmetric (As  $\Sigma$  is symmetric) & real-valued.

→ If  $A$  is a real-valued matrix,

$$A = \frac{1}{2} (A + A^T) + \frac{1}{2} (A - A^T)$$

$\downarrow$  Symmetric (B)  $\Rightarrow M^T = M$

$\downarrow$  Anti-symmetric (C)  $\Rightarrow M^T = -M$

$$B_{ij} = A_{ij} + A_{ji}$$

$\downarrow$  ' + ', so  $B_{ji}$  is also same

$$C_{ij} = A_{ij} - A_{ji}$$

$\downarrow$   $C_{ji} = -C_{ij}$

$\downarrow$   $A_{ij} = -A_{ji}$  & all diagonal elements are 0

→ If  $A$  is anti-symmetric,  $\mathbf{x}^T A \mathbf{x} = 0$ , for any  $\mathbf{x}$

$\mathbf{x}^T A \mathbf{x} = \sum_{i,j \in D} x_i \cdot x_j \cdot A_{ij} \rightarrow$

$A_{ji} = -A_{ij} \rightarrow$

$x_i \cdot x_j \cdot A_{ij} + x_j \cdot x_i \cdot A_{ji} = 0$

So, we assume  $\Sigma^{-1}$  is symmetric in  $(x-\mu)^T \Sigma^{-1} (x-\mu)$ .

For this not to be 0.

$$\sum_{D \times D} x_i u_i = \lambda_i u_i \quad \begin{array}{l} \text{eigen vector} \\ \text{eigen value} \end{array} \quad [\Sigma = \text{full-rank matrix}]$$

eigen vector of a matrix is a vector which when multiplied by that matrix, will not change its direction but only gets multiplied by a scalar.

$$U = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_D \end{bmatrix} \Rightarrow \text{eigen vectors}$$

For any  $u_i \neq u_j$ ,  $u_i^T u_j = 0$

$U U^T = U^T U = I$

$U^T U = I \quad \text{if } u_i \text{ & } u_j \text{ are orthogonal}$

$U^T U = I \quad \text{if } u_i \text{ & } u_j \text{ are from same eigen values}$

because diagonal elements represent  $u_i^T u_i = 1$  & off-diagonal represent  $u_i^T u_j = 0$

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ 0 & 0 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & \cdots & \lambda_D \end{bmatrix}$$

$$\Sigma = \sum_{i=1}^D \lambda_i u_i u_i^T \quad \begin{array}{l} \text{row vector} \\ \text{column vector} \end{array}$$

$$\Sigma^{-1} = U \Lambda^{-1} U^T$$

$$= \sum_{i=1}^D \frac{1}{\lambda_i} u_i u_i^T$$

$$u_i = [a_1 \ a_2 \ \dots \ a_n]$$

$$u_i^T = [a_1 \ a_2]$$

$$u_i u_i^T = \begin{bmatrix} a_1^2 & a_1 a_2 \\ a_2 a_1 & a_2^2 \end{bmatrix}$$

does it depend on the rank of a matrix?

Yes, the number of eigenvectors that a matrix has can depend on the dimension of the matrix, the algebraic multiplicities of the eigenvalues, and whether or not the eigenvalues are repeated.

If a square matrix  $A$  has distinct eigenvalues, then its rank is equal to the number of non-zero eigenvalues. In this case,  $A$  has linearly independent eigenvectors, where  $n$  is the dimension of  $A$ .

However, if  $A$  has repeated eigenvalues, the eigenvectors may not be linearly independent, in which case the rank of  $A$  can be less than  $n$ . In this case,  $A$  has fewer than  $n$  independent eigenvectors.

For example, a rank-deficient matrix, i.e., a matrix with rank less than its dimension, can have fewer than linearly independent eigenvectors. A zero matrix is an extreme example of a rank-deficient matrix that has infinitely many eigenvectors corresponding to the eigenvalue zero.

This is because eigenvectors of a symmetric matrix corresponding to different eigenvalues are orthogonal.

To see why this is the case, let  $u$  and  $v$  be eigenvectors of a symmetric matrix  $A$  corresponding to distinct eigenvalues  $\lambda$  and  $\mu$ , respectively. Then, we have:

$$Au = \lambda u$$

$$Av = \mu v$$

Taking the dot product of both sides of the above equations, we get:

$$u^T Av = \lambda u^T v$$

Multiplying both sides by  $v^T$ , we get:

$$v^T Av = \lambda v^T u$$

Since  $A$  is symmetric, we have:

$$v^T A u = u^T A v$$

Substituting this into the previous equation, we get:

$$\lambda v^T u = \mu v^T u$$

Since  $\lambda$  and  $\mu$  are distinct, we have  $\lambda - \mu \neq 0$ . Therefore, we can divide both sides of the above equation by  $\lambda - \mu$  to get:

$$v^T u = 0$$

No, it's general, as it's not defined for the zero eigenvalues.

However, there are some special cases where the eigenvalue of zero is the same as the eigenvalue of  $\lambda$ . One such case is when a symmetric matrix with positive definite entries has a zero eigenvalue. In this case, the zero eigenvalue is the sum of all the other eigenvalues.

Specifically, if  $\lambda$  is an eigenvalue of a real symmetric matrix  $A$ , then  $\lambda$  is also an eigenvalue of  $A^T$ . If  $\lambda$  is an eigenvalue of  $A$  with multiplicity  $k$ , then  $\lambda$  is also an eigenvalue of  $A^T$  with multiplicity  $k$ . Specifically, if  $\lambda$  is an eigenvalue of  $A$  with multiplicity  $k$ , then  $\lambda$  is also an eigenvalue of  $A^T$  with multiplicity  $k$ .

$$D^2 = (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \rightarrow \text{Mahalanobis Distance}$$

$$= (\mathbf{x} - \boldsymbol{\mu})^\top \left[ \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^\top \right] (\mathbf{x} - \boldsymbol{\mu}) = \mathbf{z}^\top \left[ \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^\top \right] \mathbf{z}$$

$$= \sum_{i=1}^D \frac{1}{\lambda_i} \underbrace{\mathbf{z}^\top \mathbf{u}_i \mathbf{u}_i^\top \mathbf{z}}_{\substack{1 \times D \\ D \times 1 \\ 1 \times D \\ D \times 1}} \hookrightarrow \text{Scalar} = y$$

$$= \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \hookrightarrow \text{Indexed by eigen vector } \mathbf{i}$$

$$P(\mathbf{n}/\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2} \sum_{i=1}^D \frac{y_i^2}{\lambda_i}}$$

$y = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu})$  acting as a transformation matrix  
 $\hookrightarrow D \times 1 \text{ vector}$

$$= \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \times e^{-\frac{1}{2} \sum_{i=1}^D \frac{y_i^2}{\lambda_i}}, \quad y_i \sim N(y_0, \lambda_i)$$

$\hookrightarrow$  also a Gaussian

$$\text{Ex: } \mathbf{x} = [3 \ 4 \ 2 \ 1]$$

$$\text{Let correspond } \mathbf{y} = [0.2 \ 0.8 \ -0.2 \ 0.6]$$

$$\text{Let } \boldsymbol{\mu} = [2 \ 3 \ 0.5 \ 0.5]$$

$$(\mathbf{x} - \boldsymbol{\mu}) = [1 \ 1 \ 1.5 \ 0.5]$$

$$\therefore \mathbf{U}(\mathbf{x} - \boldsymbol{\mu}) = \mathbf{y} = \underline{\hspace{10em}}$$

→ there is some correlation b/w features of  $\mathbf{x}$ ,  
 $0.2$  doesn't correspond to only one dimension, but encapsulates correlation  
of all features. Each  $y_i$  is independent.

$$P(y_0/\lambda) = \frac{1}{\sqrt{2\pi} \sqrt{\lambda}} e^{-\frac{1}{2} \frac{y_0^2}{\lambda}}$$

$\downarrow \sigma^2$

$$\mu = \int_{-\infty}^{\infty} x P(x) dx \rightarrow \text{for univariate}$$

$$\sigma^2 = E[(x-\mu)^2] = E(x^2 + \mu^2 - 2\mu x) = E[x^2] + \mu^2 - 2\mu E[x]$$

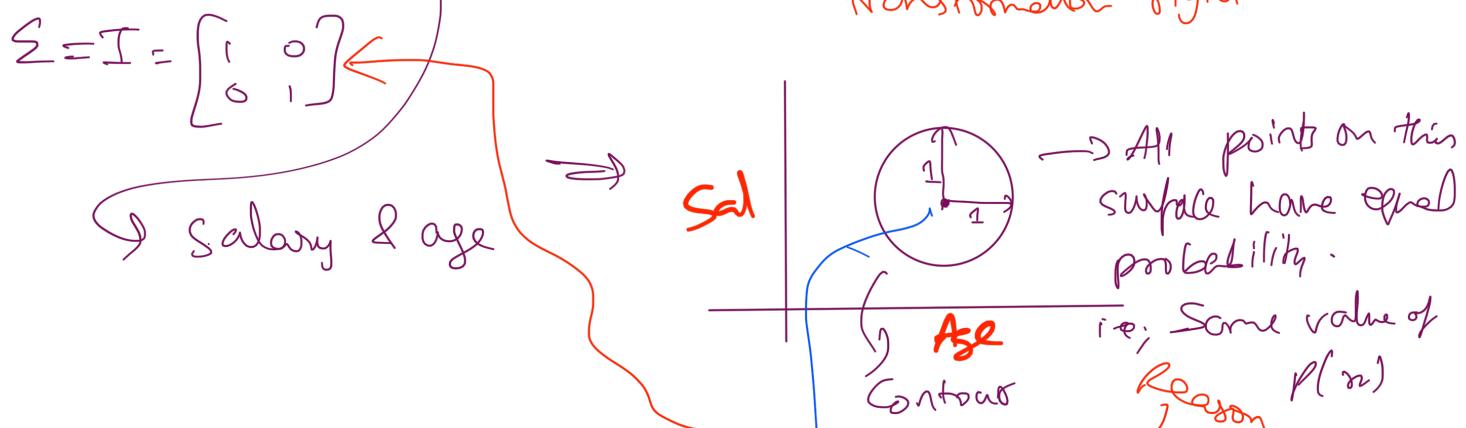
$$= E[x^2] + \mu^2 - 2\mu^2$$

$$\sigma^2 = E[x^2] - \mu^2 = E[x^2] - (E[x])^2$$

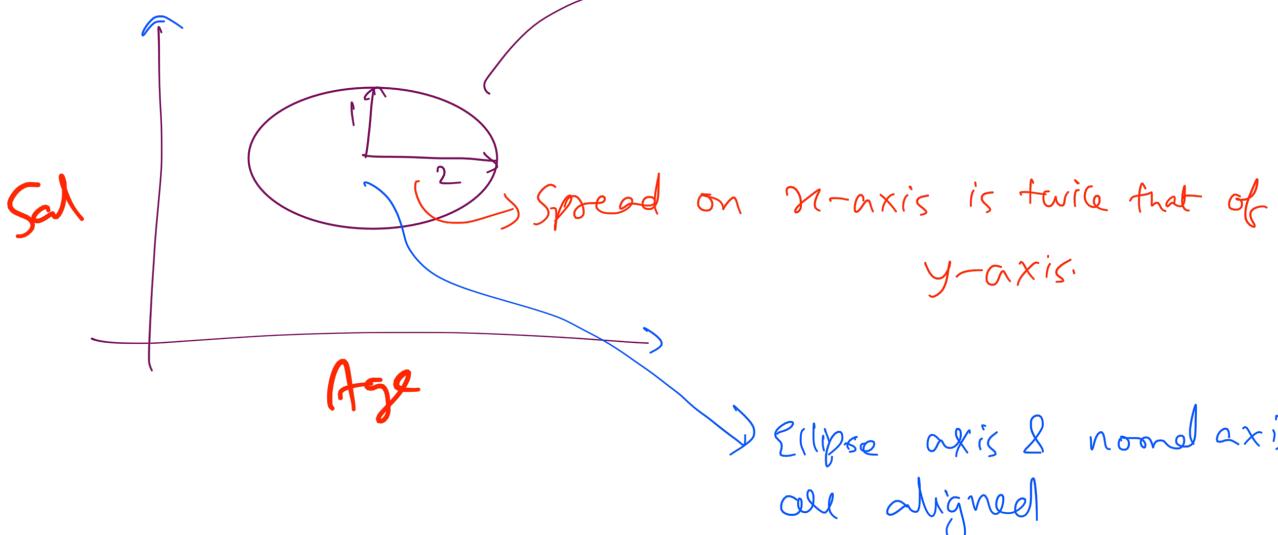
$$U = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_p^T \end{bmatrix}, \quad y = U \begin{bmatrix} x - \mu \\ D \times 1 \\ D \times D \\ D \times 1 \end{bmatrix}$$

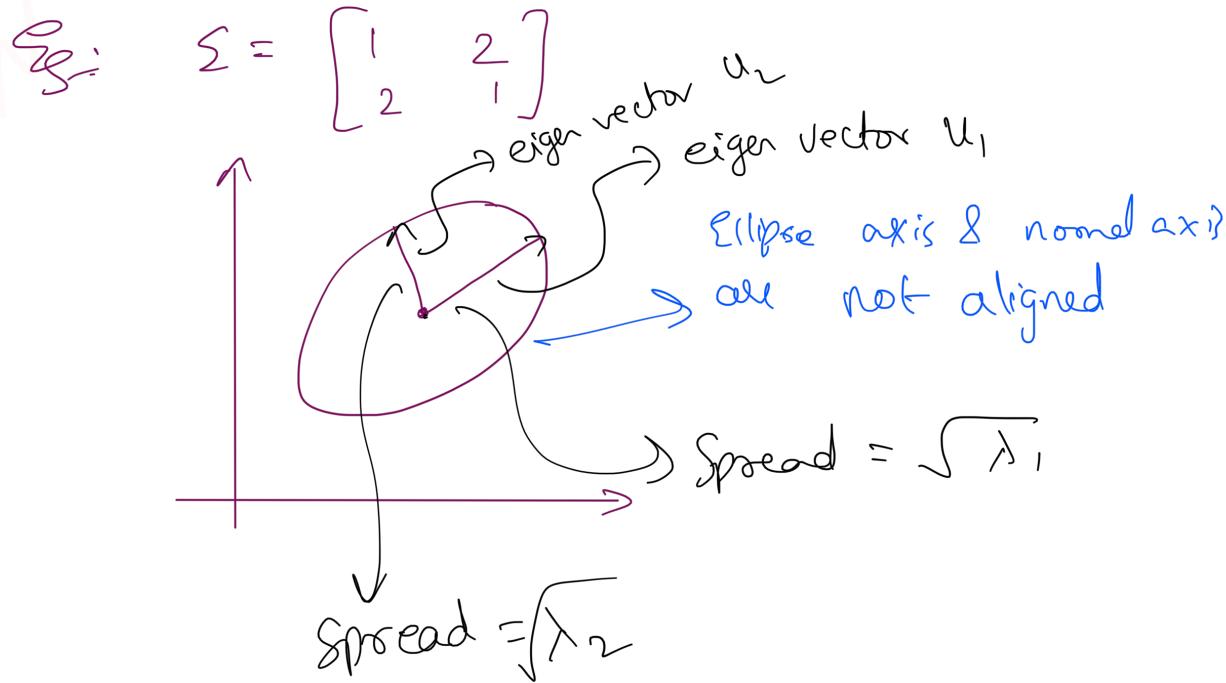
- Transformation preserves length  
 of any vector  
 - And also angle b/w any  
 two vectors.

Eg: Two-variate space.



$$\text{Now, } \Sigma = I = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$





$$\text{Multivariate} \Rightarrow P(\mathbf{n}) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (\mathbf{n}-\mu)^T \Sigma^{-1} (\mathbf{n}-\mu)}$$

$$\begin{aligned} E(\mathbf{n}) &= \int_{-\infty}^{\infty} \mathbf{n} P(\mathbf{n}) d\mathbf{n} = \int_{-\infty}^{\infty} (\mathbf{d} + \mu) \times \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} \mathbf{d}^T \Sigma^{-1} \mathbf{d}} d\mathbf{d} \\ &= \underbrace{\int_{-\infty}^{\infty} \mu P(\mathbf{n}) d\mathbf{n}}_{= \mu \times 1 \Rightarrow \text{as}} + \underbrace{\int_{-\infty}^{\infty} \mathbf{d} P(\mathbf{d}) d\mathbf{d}}_{=0} \quad \text{as } P(\mathbf{d}) \text{ is gaussian} \text{ so, } -\int_{-\infty}^{\infty} \mathbf{d} P(\mathbf{d}) = 0 \\ &= \mu \times 1 \Rightarrow \int_{-\infty}^{\infty} P(\mathbf{n}) d\mathbf{n} = 1 \quad \text{This is an odd function, so odd func} = 0 \end{aligned}$$

$$E(\mathbf{n}) = \mu$$

$$E(\mathbf{n}^2) = E(\mathbf{n}\mathbf{n}^T)$$

if  $\mathbf{n} = \begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ n_D \end{bmatrix}$

$\mathbf{n}\mathbf{n}^T = \begin{bmatrix} n_1^2 & n_1n_2 & \dots & n_1n_D \\ n_2n_1 & n_2^2 & \dots & n_2n_D \\ \vdots & \vdots & \ddots & \vdots \\ n_Dn_1 & n_Dn_2 & \dots & n_D^2 \end{bmatrix}$

All possible multiplications

$$\begin{aligned}
 E(nn^T) &= \int_{-\infty}^{\infty} nn^T p(n) dn = \int_{-\infty}^{\infty} [d + \mu] [d + \mu]^T p(d) dd \\
 &= \int_{-\infty}^{\infty} dd^T p(d) dd + 2 \int_{-\infty}^{\infty} d n^T p(d) dd + \mu \mu^T \int_{-\infty}^{\infty} p(d) dd
 \end{aligned}$$

↓  
 0 ⇒ as odd function  
 ↓ 1  
 as  $\int_{-\infty}^{\infty} p(n) dn = 1$

$$E(nn^T) = \int_{-\infty}^{\infty} dd^T p(d) dd + \mu \mu^T$$

$$\begin{aligned}
 \Rightarrow & U^T y = U^T U d \\
 & U^T y = I d = d
 \end{aligned}$$

W.C.T.  $y = U[n - \mu] = Ud \Rightarrow U^T y = d$

$$U^T = \begin{bmatrix} u_1 & u_2 & \dots & u_D \\ u_{11} & u_{12} & \dots & u_{1D} \\ u_{21} & u_{22} & \dots & u_{2D} \\ \vdots & \vdots & & \vdots \\ u_{D1} & u_{D2} & \dots & u_{DD} \end{bmatrix}, y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_D \end{bmatrix}$$

$$d = U^T y = \sum_{i=1}^D u_i y_i \xrightarrow{\substack{\text{scalar} \\ \text{vector}}}$$

$$dd^T = \left[ \sum_{i=1}^D u_i y_i \right] \left[ \sum_{j=1}^D u_j y_j \right]^T = \sum_{i=1}^D u_i y_i \sum_{j=1}^D u_j^T y_j$$

$$dd^T = \sum_{i=1}^D \sum_{j=1}^D y_i y_j u_i u_j^T$$

$$E(nn^T) = \int_{-\infty}^{\infty} dd^T p(d) dd + \mu \mu^T$$

$$\begin{aligned}
 &= \int_{-\infty}^{\infty} \sum_{i=1}^D \sum_{j=1}^D y_i y_j u_i u_j^T p(d) dd + \mu \mu^T \\
 &= \sum_{i=1}^D \sum_{j=1}^D \mu_i \mu_j^T \int_{-\infty}^{\infty} y_i y_j p(y) dy
 \end{aligned}$$

$$u_i^\top u_j = 0, \text{ but } u_i^\top u_j \neq 0, \text{ as } u_i = u_j = D \neq 1 \text{ & } u_i^\top = u_j^\top = 1 \times D$$

$$\frac{D}{\prod_{i=1}^D} e^{-\frac{1}{2} \sum_{i=1}^D \frac{y_i^2}{\lambda_i}}$$

$$E(nn^\top) = \sum_{i=1}^D \sum_{j=1}^D u_i^\top u_j \int_{-\infty}^{\infty} y_i y_j \left[ \prod_{k=1}^D P(y_k) dy_k \right]$$

when  $i \neq j$ ,

$$\int_{-\infty}^{\infty} y_i P(y_i) dy_i \times \int_{-\infty}^{\infty} y_j P(y_j) dy_j \times \prod_{\substack{k=1 \\ k \neq i, j}}^D P(y_k) dy_k$$

↳ 0      ↳ 0      ↳ 1

$\hookrightarrow$  Bring out  $P(y_i) dy_i$  &  $P(y_j) dy_j$  as it is multiplication

when  $i = j$ ,

$$E(nn^\top) = \sum_{i=1}^D u_i^\top u_i \int_{-\infty}^{\infty} y_i^2 P(y) dy$$

$$= \sum_{i=1}^D u_i^\top u_i \int_{-\infty}^{\infty} y_i^2 P(y_i) dy_i \times \prod_{k=1}^{D-1} \int_{-\infty}^{\infty} P(y_k) dy_k$$

$$= \underbrace{\sum_{i=1}^D u_i^\top u_i}_{\text{In Univariate Case,}} \lambda_i + \mu \mu^\top$$

$$\sigma^2 = E[y_i^2] - (E[y_i])^2$$

$$E(nn^\top) = \Sigma + \mu \mu^\top$$

$$\Sigma = E(nn^\top) - \mu \mu^\top$$

$$\lambda_i = \int_{-\infty}^{\infty} y_i^2 P(y_i) dy_i - 0$$

As,  $P(y_i)$  has  $\sqrt{\lambda_i}$  standard deviation and 0 as mean

Deriving conditional & marginal distribution of a gaussian distributed normal variable

Eg- random variable (gaussian dist.)

= [Battery, Camera, Display, pricing, hardware]

(Assume quality of all these features is available in numerical form)

Example  $\Rightarrow n = [5, 4, 3, 3, 3] \Rightarrow$  multivariate gaussian distributed variable

$$n \sim P(m(\mu, \Sigma)) = N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Let  $P = \Sigma^{-1}$   
Precision Matrix

If these features are joint distributed according to gaussian distribution, what is the form of conditional probability of some features, given other features?

$\rightarrow P(x_r/x_c)$ , where  $x_c$  = constant/evidence features

Conditional distribution given joint distribution is gaussian  
 $x_r$  = random/obj features

Here,  $n = [5, 4, \underbrace{3, 3, 3}_{\downarrow n_r \quad \downarrow n_c}]$ ,  $P = \begin{bmatrix} 2 & 3 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} P_{rr} & P_{rc} \\ P_{cr} & P_{cc} \end{bmatrix}, P_{rc} = P_{cr}^T$

$$\begin{bmatrix} P_{rr} & P_{rc} \\ P_{cr} & P_{cc} \end{bmatrix}$$

Precision Matrix

Now,  $-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)$

$$= -\frac{1}{2} \left[ \underbrace{(x_r - \mu_r)^T}_{1 \times 2}, \underbrace{(x_c - \mu_c)^T}_{1 \times 3} \right] \begin{bmatrix} P_{rr} & P_{rc} \\ P_{cr} & P_{cc} \end{bmatrix} \begin{bmatrix} (x_r - \mu_r) \\ (x_c - \mu_c) \end{bmatrix} \rightarrow 2 \times 1 \\ \rightarrow 3 \times 1$$

$$= -\frac{1}{2} \left[ (\boldsymbol{x}_r - \boldsymbol{\mu}_r)^T P_{rr} + (\boldsymbol{x}_c - \boldsymbol{\mu}_c)^T P_{cc} , \underbrace{(\boldsymbol{x}_r - \boldsymbol{\mu}_r)^T P_{rc} + (\boldsymbol{x}_c - \boldsymbol{\mu}_c)^T P_{cr}}_{(1 \times 2)} \right] \downarrow$$

$1 \times 2$

$1 \times 1$

$$\times \begin{bmatrix} (\boldsymbol{x}_r - \boldsymbol{\mu}_r) \\ (\boldsymbol{x}_c - \boldsymbol{\mu}_c) \end{bmatrix} \rightarrow 2 \times 1 \\ \rightarrow 3 \times 1$$

① involves quadratic term of  $\boldsymbol{x}_r$

② involves linear term of  $\boldsymbol{x}_r$

$$= -\frac{1}{2} \left[ (\boldsymbol{x}_r - \boldsymbol{\mu}_r)^T P_{rr} (\boldsymbol{x}_r - \boldsymbol{\mu}_r) + (\boldsymbol{x}_c - \boldsymbol{\mu}_c)^T P_{cr} (\boldsymbol{x}_r - \boldsymbol{\mu}_r) \right. \\ \left. + (\boldsymbol{x}_r - \boldsymbol{\mu}_r)^T P_{rc} \times (\boldsymbol{x}_c - \boldsymbol{\mu}_c) + (\boldsymbol{x}_c - \boldsymbol{\mu}_c)^T P_{cc} \times (\boldsymbol{x}_c - \boldsymbol{\mu}_c) \right]$$

$$-\frac{1}{2} (\boldsymbol{x} - \boldsymbol{\mu})^T P (\boldsymbol{x} - \boldsymbol{\mu}) = -\frac{1}{2} \boldsymbol{x}^T P \boldsymbol{x} + \cancel{\frac{1}{2} \boldsymbol{x}^T P \boldsymbol{\mu}} + \boldsymbol{\mu}^T P \boldsymbol{x}$$

Gaussian Dist. form. So we can

say that  $P(\boldsymbol{x}_r | \boldsymbol{x}_c)$  is of gaussian dist. form only.

$$= -\frac{1}{2} \left[ (\boldsymbol{x}_r^T P_{rr} \boldsymbol{x}_r) - 2 \cancel{\boldsymbol{x}_r^T P_{rr} \boldsymbol{\mu}_r} + \boldsymbol{\mu}_r^T P_{rr} \boldsymbol{\mu}_r \right. \\ \left. + (\boldsymbol{x}_c - \boldsymbol{\mu}_c)^T P_{cr} \boldsymbol{x}_r + \boldsymbol{x}_r^T P_{rc} (\boldsymbol{x}_c - \boldsymbol{\mu}_c) \right. \\ \left. + (\boldsymbol{x}_c - \boldsymbol{\mu}_c)^T P_{cr} \cancel{- (-\boldsymbol{\mu}_r)} \right. \\ \left. + (-\boldsymbol{\mu}_r)^T P_{rc} (\boldsymbol{x}_c - \boldsymbol{\mu}_c) \right. \\ \left. + (\boldsymbol{x}_c - \boldsymbol{\mu}_c)^T P_{cc} (\boldsymbol{x}_c - \boldsymbol{\mu}_c) \right]$$

Aim: To find  $\boldsymbol{\mu}$  &  $\Sigma$  for  $P(\boldsymbol{x}_r | \boldsymbol{x}_c)$

As when we talk about a Gaussian dist., we have to know about its  $\boldsymbol{\mu}$  &  $\Sigma$ .

$$= -\frac{1}{2} \left[ (\boldsymbol{x}_r^T P_{rr} \boldsymbol{x}_r) - 2 \cancel{\boldsymbol{x}_r^T P_{rr} \boldsymbol{\mu}_r} + 2 \cancel{\boldsymbol{x}_r^T P_{rc} (\boldsymbol{x}_c - \boldsymbol{\mu}_c)} \right. \\ \left. + \text{const} + \text{const} + \text{const} \right]$$

As here no  $P_{rr}$ , so  $P_{rr}^{-1} P_{rr} \times P_{rr}^{-1} = I$

$$= -\frac{1}{2} \left[ \boldsymbol{x}_r^T P_{rr} \boldsymbol{x}_r \right] + \boldsymbol{x}_r^T P_{rr} \left[ \boldsymbol{\mu}_r - P_{rr}^{-1} P_{rc} (\boldsymbol{x}_c - \boldsymbol{\mu}_c) \right] \\ + \text{const.}$$

$$\therefore P(\boldsymbol{x}_r | \boldsymbol{x}_c) = N \left( \boldsymbol{\mu}_{r/c}, \Sigma_{r/c} \right) \rightarrow 2 \times 2 \text{ matrix}$$

Here,  $2 \times 1 \rightarrow 2 \times 1$

$$\Sigma_{r/c} = P_{rr}^{-1}$$

$$M_{r/c} = M_r - P_{rr}^{-1} P_{rc} (n_c - \mu_c)$$

Cause, it is of the form  $(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$n_r^T$$

$$P_{rr}$$

$$n_r$$

→ Cause,  $\mu$  is whatever is filled by a linear term in "n\_r".

$$\Sigma = \begin{bmatrix} \Sigma_{rr} & \Sigma_{rc} \\ \Sigma_{cr} & \Sigma_{cc} \end{bmatrix}$$

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}$$

$$A^{-1} = A - BD^{-1}C$$

$$\text{So, } P_{rr}^{-1} = P_{rr} - P_{rc} P_{cc}^{-1} P_{cr}$$

$$P = \Sigma^{-1}$$

$$P = \begin{bmatrix} \Sigma_{rr} & \Sigma_{rc} \\ \Sigma_{cr} & \Sigma_{cc} \end{bmatrix}^{-1} = \begin{bmatrix} P_{rr} & P_{rc} \\ P_{cr} & P_{cc} \end{bmatrix}$$

Generally,

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CMBD^{-1} \end{bmatrix}$$

$$M = (A - BD^{-1}C)^{-1}$$

$$\therefore P_{rr} = (P_{rr} - P_{rc} P_{cc}^{-1} P_{cr})^{-1}$$

$$\therefore P_{rr}^{-1} = P_{rr} - P_{rc} P_{cc}^{-1} P_{cr}$$

## Marginal Gaussian Distribution

If joint distribution,  $P(x_r, x_c)$  is gaussian, then the conditional distribution,  $P(x_r | x_c)$  is also gaussian.

Now, find,

$$P(x_r) = \int P(x_r, x_c) dx_c$$

↳ also Gaussian

As we have to differentiate w.r.t.  $x_c$ , we need to consider the terms containing  $x_c$  from the same eqn. used in the conditional distribution

$$-\frac{1}{2} \left[ \underbrace{(x_r - \mu_r)^T P_{rr} (x_r - \mu_r)}_{\text{const}} + (x_c - \mu_c)^T P_{cc} (x_c - \mu_c) + (x_r - \mu_r)^T P_{rc} \times (x_c - \mu_c) + (x_c - \mu_c)^T \times P_{rc} \times (x_r - \mu_r) \right]$$

$$= -\frac{1}{2} \left[ \underbrace{x_c^T P_{cr} x_r}_{\text{cancel}} - \underbrace{x_c^T P_{cr} \mu_r}_{\text{cancel}} + \underbrace{x_r^T P_{rc} x_c}_{\text{cancel}} - \underbrace{\mu_r^T P_{rc} x_c}_{\text{cancel}} + x_c^T P_{cc} x_c - \mu_c^T P_{cc} x_c + \text{const} \right]$$

$$= -\frac{1}{2} \left[ \underbrace{x_c^T P_{cr} x_r}_{\text{cancel}} - 2 \underbrace{x_c^T P_{cr} \mu_r}_{\text{cancel}} + \underbrace{x_r^T P_{rc} x_c}_{\text{cancel}} + \underbrace{\mu_r^T P_{rc} x_c}_{\text{cancel}} - \underbrace{x_c^T P_{cc} \mu_c}_{\text{cancel}} - \underbrace{\mu_c^T P_{cc} x_c}_{\text{cancel}} \right]$$

$$= -\frac{1}{2} \left[ 2 x_c^T P_{cr} x_r + x_c^T P_{cc} x_c - 2 x_c^T P_{cr} \mu_r - 2 x_r^T P_{rc} \mu_c \right]$$

$$= -\frac{1}{2} [x_c^T P_{cr} x_r - x_c^T P_{cr} \mu_r + x_c^T P_{cr} \mu_r + x_c^T P_{cc} \mu_c]$$

$$= -\frac{1}{2} \boldsymbol{x}_c^T P_{cc} \boldsymbol{x}_c + \boldsymbol{x}_c^T (P_{cc} \boldsymbol{m}_c - P_{cr} (\boldsymbol{n}_r - \boldsymbol{m}_r))$$

$$\text{let } m = [P_{cc} \boldsymbol{m}_c - P_{cr} (\boldsymbol{n}_r - \boldsymbol{m}_r)]$$

$$\rightarrow = -\frac{1}{2} \boldsymbol{x}_c^T P_{cc} \boldsymbol{x}_c + \boldsymbol{x}_c^T m$$

$$= -\frac{1}{2} (\boldsymbol{x}_c - P_{cc}^{-1} m)^T P_{cc} (\boldsymbol{x}_c - P_{cc}^{-1} m)$$

$$+ \frac{1}{2} m^T P_{cc}^{-1} m$$

↳ standard quadratic form of a Gaussian dist.

$\therefore$  When integrating over  $P(c)$ , that integral over  $\boldsymbol{x}_c$  will take the form,

$$\int e^{-\frac{1}{2} (\boldsymbol{x}_c - P_{cc}^{-1} m)^T P_{cc} (\boldsymbol{x}_c - P_{cc}^{-1} m)} d\boldsymbol{x}_c$$

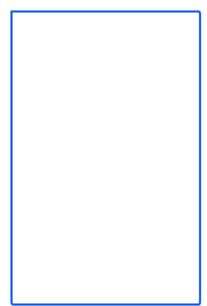
↳ integral over unnormalize gaussian.

So, result = reciprocal of the normalization coeff.

## Semi-Supervised Learning



$D_L \rightarrow$  limited.



This is also a version of  
EM-algorithm

- Some labelled data, remaining unlabeled data. Build
- Build a model using the labelled data.  $\Rightarrow$  sub-optimal ( $\theta^0$ )
- Then build an optimal classifier using the entire data  
 $\hookrightarrow \theta^1$
- Then re-label  $D_U$  using  $\theta^1$ ,  $\Rightarrow$  you get  $\theta^2$
- Do this till changes are not significant.

why does this work?

Find out.

↳ as some features that are not in  $D_L$  are in  $D_U$ . How do these features get weights?  $\rightarrow$  Brief: There unknown features co-occur with some features that are already in  $D_L$ .

- This will work well only if Naive Bayes classifier is used. Not so good when logistic regression, SVM, Gaus neural nets are used as they are discriminative.

$I(f_i \in n_i)$   $\Rightarrow$  sets all  $x_i$  such that  $f_i$  is in  $n_i$

$$KL(P(y|n_i) || P_0(y|n_i))$$

model output

ground truth  $\rightarrow [0.95, 0.05]$

for all documents containing feature  $f_i$ .

$$F \sum_{i=1}^{ID_L} \sum_{n=1} P_i(y|l_{mn}) \parallel P_0(y|m_n)$$

**labelled feature learning is weak-supervised**

**labelled features**  $\Rightarrow$  Use if you think the presence of some features indicates a particular class.

→ let there be a feature vector matrix  $X$  containing 2000 documents of which not all are labelled.

- let the word excellent be present in 100 documents & we assume excellent to be a positive label indicator  $\Rightarrow [0.95 \ 0.05] \Rightarrow$  labelled feature  
 $y = +ve \quad y = -ve \rightarrow x$  contains excellent as a word  
 $P(y|x \in "excellent") = [0.95, \ 0.05]$

∴ the loss function for this specific feature will contain only these 100 documents.

Assume there are 10 positive labelled features, 10 negative labelled features, & 10 neutral labelled features.

If no. of words in vocab is 1 lakh,

then model size  $\Rightarrow 3 \times 1$  lakh  $\Rightarrow L \times k \rightarrow$  vocab. size  
 $\underbrace{d}_{\text{pos, -ve, neutral}}$   $\downarrow$   $\text{no. of labels}$

Shouldn't be  $3 \times 30$  as to include only the labelled features, as, other features also contribute to the label given to the document.

$$P(y/x \in \text{"excellent"}) = \begin{bmatrix} y=2 & y=1 & y=0 \\ 0.90 & 0.05 & 0.05 \end{bmatrix} \Rightarrow \hat{P}_k$$

$\hookrightarrow$  probability vector

$$KL(y \parallel \hat{P}_o(y/x)) \Rightarrow \text{Loss function for supervised classification}$$

$\hookrightarrow$  one-hot vector

$$\hat{P}_o(y/I(f_k \in x)) = \frac{1}{\sum_{n=1}^N I(f_k \in x_n)} \sum_{n=1}^N I(f_k \in x_n) P_o(y/x_n)$$

$\hookrightarrow$  count of docs. containing  $f_k$

parameter matrix  
of size  $L \times K$

$$\begin{array}{c} \text{Ex: } \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.9 & 0.05 & 0.05 \\ 0.7 & 0.2 & 0.1 \end{bmatrix} \\ \hline [0.8 & 0.11 & 0.09] \end{array} \Rightarrow \begin{array}{l} \text{let these be the} \\ \text{labels given by our} \\ \text{model to docs containing } f_k \end{array}$$

$\hookrightarrow$  average the 3 vectors.  
 We aren't working on per document classification here, as the ground truth for labelled features. Say that all docs containing it must have prob. dist. similar to it. So we work on avg. of all docs containing that feature.

Now, if ground truth is  $[0.90 \ 0.05 \ 0.05]$

Then for feature  $f_k \Rightarrow KL([0.90 \ 0.05 \ 0.05], [0.8 \ 0.11 \ 0.09])$

Using labelled features, we reduce the time taken to get gold labels

$$(LSS) \Rightarrow \sum_{k=1}^K KL\left(\tilde{P}_k(y|I(f_k \in n)) || \hat{P}_{k\theta}(y|I(f_k \in n))\right)$$

Gradient update equations

$$L = \sum_{k=1}^K \sum_y \tilde{P}_k(y|I(f_k \in n)) \log \frac{\tilde{P}_k(y|I(f_k \in n))}{\hat{P}_{k\theta}(y|I(f_k \in n))}$$

$$L = \sum_{k=1}^K \sum_y \left[ \tilde{P}_k(y|I(f_k \in n)) \log \tilde{P}_k(y|I(f_k \in n)) - \tilde{P}_k(y|I(f_k \in n)) \log \hat{P}_{k\theta}(y|I(f_k \in n)) \right]$$

$$L = \sum_{k=1}^K \sum_{n=1}^N \sum_{l=1}^M \left[ \tilde{P}_k(y=l|I(f_k \in n)) \log \tilde{P}_k(y=l|I(f_k \in n)) - \tilde{P}_k(y=l|I(f_k \in n)) \log \hat{P}_{k\theta}(y=l|I(f_k \in n)) \right]$$

$$\frac{\partial P_\theta^{(y=l/n)}}{\partial \theta_{ij}} = p_l (-p_i + I(l=i)) n_j \quad \Rightarrow p_l = \frac{e^{w_l^T n}}{\sum_y e^{w_y^T n}}$$

$$\begin{aligned} \frac{\partial L}{\partial \theta_{ij}} &= - \sum_{k=1}^K \sum_{n=1}^N \sum_{l=1}^M \underbrace{\tilde{P}_k(y=l|I(f_k \in n))}_{\hat{P}_{k\theta}(y=l|I(f_k \in n))} \times \hat{P}_{k\theta}(y=l|I(f_k \in n)) \\ &\quad \times \left( -\hat{P}_{k\theta}(y=i|I(f_k \in n)) + I(l=i) \right) \\ &\quad \times n_{nj} \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \theta_{ij}} (L) &= - \sum_{k=1}^K \sum_{n=1}^N \sum_{l=1}^M I(f_k \in n) \times \hat{P}_k(y=l|n) \times \\ &\quad \left( -\hat{P}_{k\theta}(y=i|n) + I(l=i) \right) \times n_{nj} \end{aligned}$$

# K-means Clustering

Optimization Objective  $\Rightarrow$

$$\sum_{c=1}^C N_c \sum_{i=1}^n \|x_i - \mu_c\|_2^2$$

no. of datapoints in cluster "c"      Within cluster sum of square  
  ↳ L<sub>2</sub>-norm

## Algorithm

① Initialize K random points in a D-dimensional space randomly.  $(\mu_1, \mu_2, \dots, \mu_K) \Rightarrow$  Cluster centroids.

- ② For each data point  $x_1, \dots, x_N$ , assign it to nearest cluster center.
- ③ Calculate new cluster centers.

$$\mu_{1\dots K}^{t+1} = \text{mean of data points in respective clusters.}$$

## Probabilistic K-means

While assigning cluster to a data point, assign probabilities of this data point belonging to each cluster.

Model seedvalued  $\Rightarrow$  Gaussian dist.

Model text  $\Rightarrow$  Multinomial dist.

Model arrival times  $\Rightarrow$  Poisson

What kind of distribution to use for a particular problem?

→ Assume mixture of that particular distribution in real world.

## Gaussian Mixture Models

## Expectation-Maximization

↳ used when something is unobserved.

Assume "c" mixtures. Mixture model has  $\mu_1, \Sigma_1, \mu_2, \Sigma_2, \dots, \mu_c, \Sigma_c$

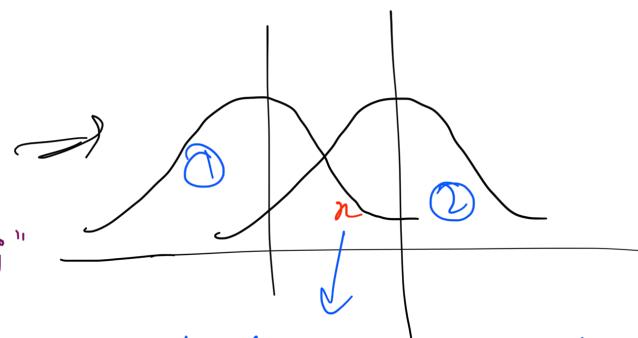
$$x \sim P(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2} \Sigma^{-1}(x-\mu)^T (x-\mu)}$$

↳ for single gaussian

$x \sim$  Mixture of gaussians

$$= \sum_{i=1}^C \pi_i \times N(x|\mu_i, \Sigma_i)$$

↳ weight of gaussian "i"



prob of  $n$  is weighted

sum of likelihoods that  $x$  come from ① & ②.

→ No. of parameters if  $C=2$  &  
dimensions of data points is 5 ( $D=5$ ),

$$\mu_1 \Rightarrow 5 \times 1$$

$$\mu_2 \Rightarrow 5 \times 1$$

$$\Sigma_1 \Rightarrow 5 \times 5 \Rightarrow \text{Max} \Rightarrow 62 \text{ parameters}$$

$$\Sigma_2 \Rightarrow 5 \times 5$$

$$\pi_1 \Rightarrow 1$$

$$\pi_2 \Rightarrow 1 \text{ or } (1 - \pi_1)$$

## E-M steps

① Initialization  $\Rightarrow \pi = [\pi_1, \pi_2, \pi_3, \dots, \pi_k] \rightarrow k \times 1$

$$\mu = [\mu_1, \mu_2, \dots, \mu_k] \rightarrow K \times D$$

Each mean is a  $D$ -dimensional vector

$$\Sigma_1, \Sigma_2, \dots, \Sigma_k \Rightarrow K \times D \times D$$

② E-step  $\Rightarrow$  Calculate the probability of each point  $x_i$  to be from each of the mixture components. O/p is a  $k$ -dimensional probability vector, for each point.

③ M-step  $\Rightarrow$  Recalculate  $\mu_{1..k}, \Sigma_{1..k} \Delta \pi$ .

$$\mu_k = \frac{1}{N_k} \sum_{i \in X} p_{ik} x_i$$

$$N_k = \sum_{i=1}^N p_{ik}$$

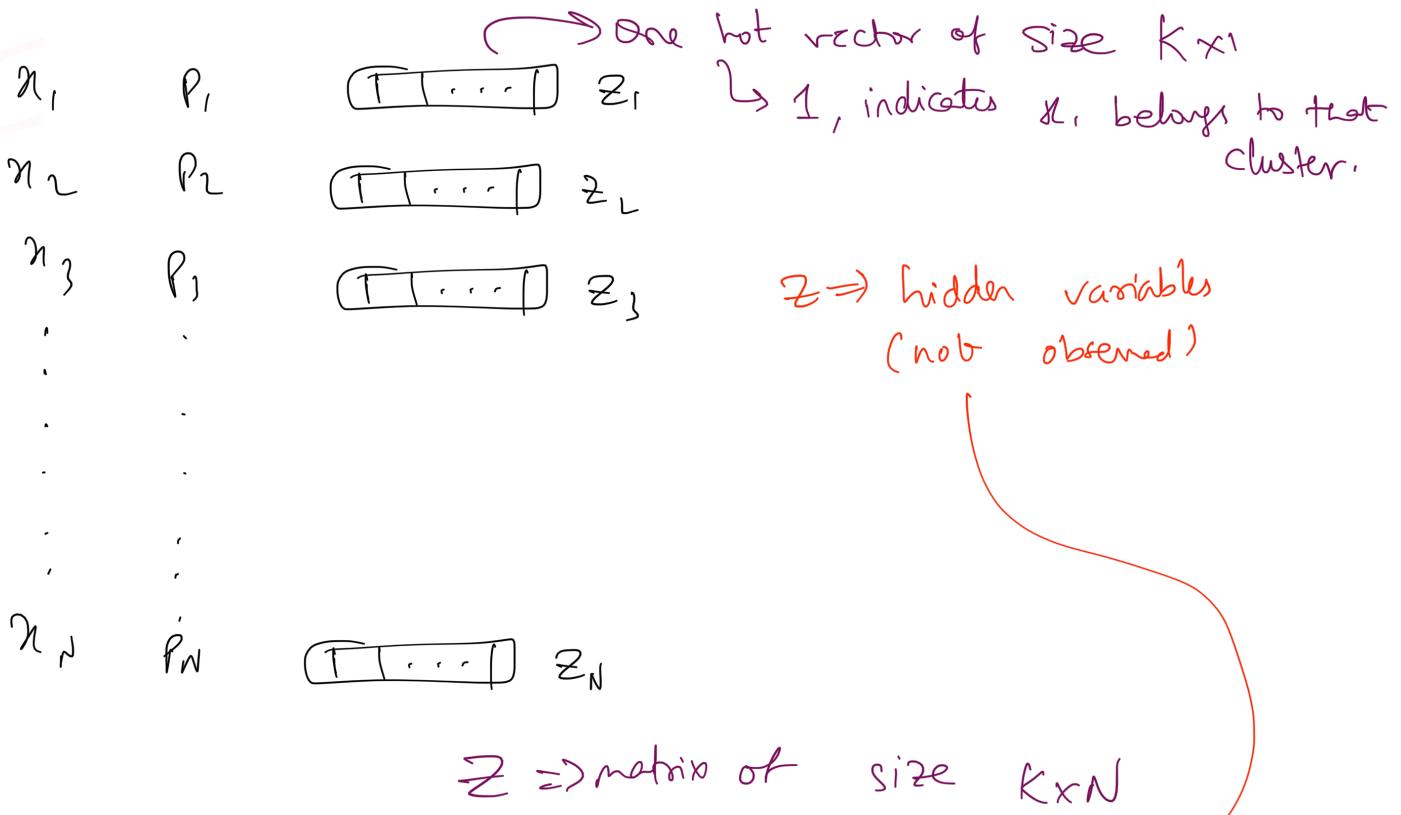
$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^N p_{ik} \times [x_i - \mu_k]^2$$

$$P_{ik} = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{k'=1}^K \pi_{k'} N(x_i | \mu_{k'}, \Sigma_{k'})}$$

$$P(x|\theta) = \sum_{k=1}^K \pi_{ik} N(x | \mu_k, \Sigma_k)$$

prob of  $x_i$  coming from  $k^{th}$  gaussian dist.

↳ likelihood function for  $x$  coming from this entire mixture of gaussians



$$P(z_{ic}=1/\theta, x) = \text{Posterior}$$

$$\ln P(x/\theta) = \ln \sum_z P(x, z/\theta)$$

$$P(x/\theta) = \sum_z P(x, z/\theta)$$

$\rightarrow$  Can apply gradient descent but using EM is advised as  
 these type of equations (log summation over sub-components)  
 do not have closed-form solutions

$\ln P(x, z/\theta) \Rightarrow$  Complete log likelihood  
 $\downarrow$   
 $\begin{array}{l} \text{observed variable} \\ \text{unobserved variable} \end{array}$

take expectation of this under posterior prob. distribution.



$$E_p \ln P(x, z/\theta)$$

$\downarrow$   
 because we do not know the identity of the cluster of each  $x_i$ .  
 Same like in coin toss where we don't know the identity of coin A / coin B in E-M coin toss problem.

log summation solution is difficult. So, do it in e-m steps & fill the unobserved variables, to form complete log likelihood function which is easier to solve (using posterior probs.)

$$\frac{\partial}{\partial M_i} \left( \ln P(x|\theta) \right) = \frac{\partial}{\partial M_i} \left[ \ln \sum_{k=1}^K \pi_k N(n|m_k, \varepsilon_k) \right]$$

Gradient Descent

$$= \frac{1}{\sum_{k=1}^K \pi_k N(n|m_k, \varepsilon_k)} \times \pi_i \times N(n|m_i, \varepsilon_i) \times -\frac{2}{2} \times (x - M_i)^T \varepsilon_i^{-1} (-)$$

$$\hookrightarrow p_i$$

$$= p_i \times (x - M_i)^T \varepsilon_i^{-1}$$

$$\ln P(x|\theta) = \ln \sum_z p(x, z|\theta)$$

$$L(q, \theta) + k L(q || p(z|x, \theta))$$

$\hookrightarrow$  Lower Bound

Tries to match  $q$  with the posterior prob.

$\rightarrow$  Then fixes  $q$  & tries to optimize  $\theta$ .

$$x_1, x_2, \dots, x_N \Rightarrow x \sim N(\mu, \Sigma)$$

$\downarrow$   
 $\partial x^1$

$\downarrow$   
 $\partial x^1$

$\downarrow$   
 $\partial x^D$

Lamda

Given  $N$  Data points, likelihood of dataset,  $X$ ,

$$X = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}_N$$

$$P(X|\theta) = \prod_{n=1}^N P(x_n|\theta)$$

$$\theta = \mu, \Sigma$$

$$= \prod_{n=1}^N N(x_n|\mu, \Sigma)$$

$$\log \text{likelihood} = \log P(X|\theta) = \sum_{n=1}^N \log \frac{1}{(2\pi)^D} + \log \frac{1}{\det(\Sigma)} + \log \det^{-1}(\Sigma)$$

Derivative  $\xrightarrow{\text{then}}$  we get,

$$\mu = \frac{1}{N} \sum_{n=1}^N x_n, \quad \Sigma = \frac{1}{N} \sum_{n=1}^N [x_n - \mu] [x_n - \mu]^T$$

$\underbrace{\quad}_{D \times 1}$      $\underbrace{\quad}_{1 \times D}$  (outer product)     $\underbrace{\quad}_{D \times D}$

Multinomial

$$x_1, x_2, \dots, x_N, \quad x \sim \text{Multinomial}(n|\theta), \quad n \text{ is } K\text{-dimensional}$$

$$\sum_{k=1}^K \theta_k = 1 \quad \& \quad \theta_k \geq 0$$

Each  $x$  is a one-hot vector

$$P(X|\theta) = \prod_{n=1}^N P(x_n|\theta)$$

$$= \prod_{n=1}^N \prod_{k=1}^K \theta_k^{x_{nk}}$$

$$\log P(X|\theta) = \sum_{n=1}^N \sum_{k=1}^K x_{nk} \log \theta_k \rightarrow \left[ \sum_{k=1}^K \theta_k - 1 \right]$$

$$\theta_k = \frac{n_k}{N}$$

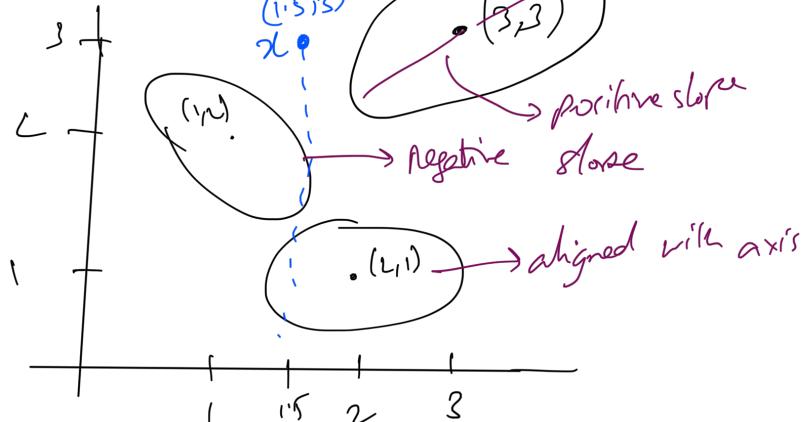
## EM Mixture model

$$x_1, x_2, \dots, x_N, x \sim p(x|\theta)$$

$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 1 & -0.8 \\ 0.8 & 1 \end{bmatrix}$$



Assume data is from these three gaussians.

$$x_1, x_2, \dots, x_N \rightarrow x \sim p(x|\theta)$$

Three scores from each gaussian for each point  $x$ .

$$\rightarrow S_1, S_2, S_3$$

$$\text{Probs} \rightarrow \underbrace{\frac{S_1}{S_1+S_2+S_3}}_{P_1}, \underbrace{\frac{S_2}{S_1+S_2+S_3}}_{P_2}, \underbrace{\frac{S_3}{S_1+S_2+S_3}}_{P_3}$$

↳ Prob. of  $x$  comes from GMML.

$$\sum_{i=1}^3 P_i = 1$$

→ We are able to calculate these probabilities as we know the parameters.

$$x_1, x_2, \dots, x_N \rightarrow x \sim P(x|\theta)$$

$z_1=2$     $z_2=1$     $\dots$     $z_N=1$

→ We don't have these tags.

→ If we atleast had these tags we would have calculated the parameters of each individual GMM.

So we use EM algo to solve this GMM problem.

$$p(x|\theta) = \sum_{k=1}^K \pi_k \times N(x|\mu_k, \Sigma_k) \rightarrow \text{Mixture model likelihood}$$

↓ weight of component k

↓ likelihood of a single data point.

$\sum_{i=1}^K \pi_i = 1$ , & each  $\pi_i > 0$

$$P(X|\theta) = \prod_{n=1}^N p(x_n|\theta)$$

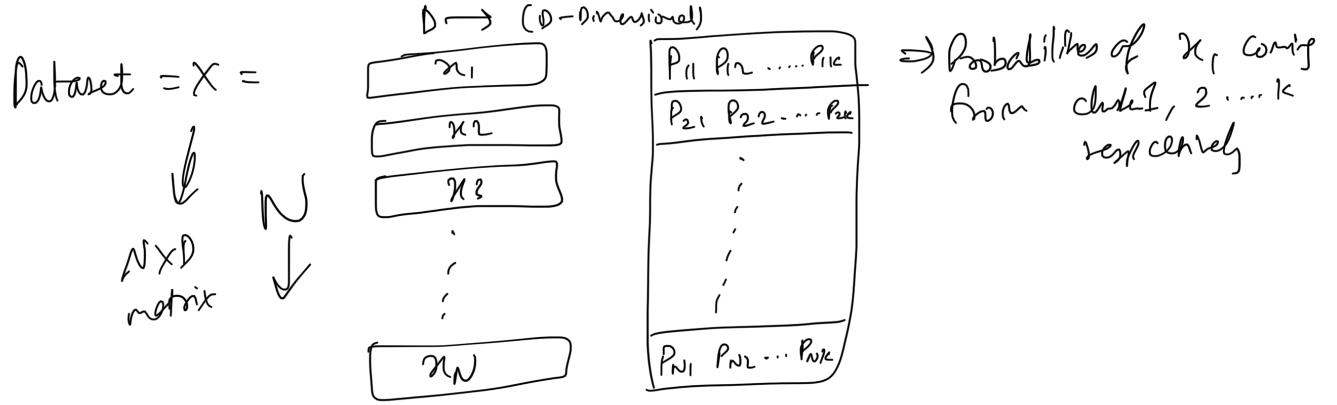
$$\text{No. of parameters} = K \times [D + D \times D] + K$$

$$P(X|\theta) = \prod_{n=1}^N \sum_{k=1}^K \pi_k N(x_n|\mu_k, \Sigma_k)$$

$$\log P(X|\theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k N(x_n|\mu_k, \Sigma_k)$$

$$\frac{\partial}{\partial \mu_k} (\ln - \mu_k)$$

$$\begin{aligned} \frac{\partial \log P(X|\theta)}{\partial \mu_k} &= \sum_{n=1}^N \frac{1}{\sum_{k'=1}^K \pi_{k'} N(x_n|\mu_{k'}, \Sigma_{k'})} \times \pi_k N(x_n|\mu_k, \Sigma_k) \\ &= \sum_{n=1}^N \frac{\pi_k N(x_n|\mu_k, \Sigma_k)}{\sum_{k'=1}^K \pi_{k'} N(x_n|\mu_{k'}, \Sigma_{k'})} \times (x_n - \mu_k)^T \Sigma_k^{-1} (-1) \\ &\quad \downarrow \frac{\partial n^T A_n}{\partial n} = 2A_n (\sigma^2) \\ &\quad \uparrow \frac{\partial \Sigma_k^{-1}}{\partial n} \end{aligned}$$



$$\frac{\partial \log P(x|\theta)}{\partial \mu_k} = \sum_{n=1}^N \frac{\pi_k N(x_n/\mu_k, \Sigma_k)}{\sum_{k'=1}^K \pi_{k'} N(x_n/\mu_{k'}, \Sigma_{k'})} \times (x_n - \mu_k)^T \Sigma_k^{-1}$$

$\hookrightarrow p_{nk}$

$$O = \sum_{n=1}^N p_{nk} [x_n - \mu_k]^T \underbrace{\Sigma_k^{-1}}_{\text{neglected as it is common for every data point}}$$

$$O = \sum_{n=1}^N p_{nk} x_n^T - \sum_{n=1}^N p_{nk} \mu_k^T$$

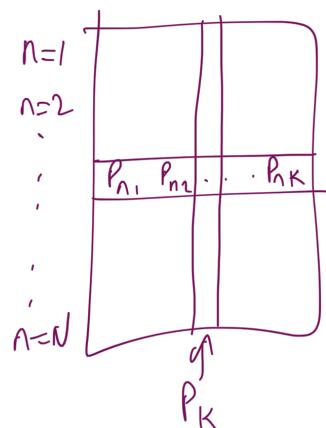
$$\Rightarrow \sum_{n=1}^N p_{nk} x_n^T = \sum_{n=1}^N p_{nk} \mu_k^T$$

$$\Rightarrow \sum_{n=1}^N p_{nk} x_n = \mu_k \sum_{n=1}^N p_{nk} \quad (\text{Transpose on both sides})$$

$$\Rightarrow \mu_k = \frac{\sum_{n=1}^N p_{nk} x_n}{\sum_{n=1}^N p_{nk}}$$

let  $N_k = \sum_{n=1}^N p_{nk}$

$$\mu_k = \frac{\sum_{n=1}^N p_{nk} x_n}{N_k}$$



$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N p_{nk} [x_n - \mu_k] [x_n - \mu_k]^T$$

$$\frac{\partial \log P(x|\theta)}{\partial \pi_K} = \sum_{n=1}^N \log \frac{N(n_n | \mu_K, \Sigma_K)}{\sum_{k'=1}^K \pi_{k'} N(n_n | \mu_{k'}, \Sigma_{k'})} \times \frac{\pi_K}{\sum_{k'=1}^K \pi_{k'}} \quad \text{Multiply w.m.d. by } \pi_K$$

$$0 = \frac{1}{\pi_K} \sum_{n=1}^N p_{nk} - \lambda \quad \begin{aligned} & - \lambda \left( \sum_{i=1}^K \pi_i - 1 \right) \\ & \frac{\partial}{\partial \pi_K} (\pi_1 + \pi_2 + \dots + \pi_K) = 1 \end{aligned}$$

$$\pi_K = \frac{\sum_{n=1}^N p_{nk}}{\lambda} \quad \& \quad \lambda = N$$

$$\therefore \pi_K = \frac{N_k}{N}$$

Mixture proportion

When we are estimating  $\mu_K$ , we assume  $p_{nk}$  as it is not a function involving  $\mu_K$ ; which isn't true.

So, we can use these formulas only if we fix  $\mu_K$  &  $\Sigma_K$ . (like in the EM algo), else there is no closed form solution.

$z \Rightarrow$  latent variable  $\Rightarrow$  indicate the component associated to  $x$   
 ↳ Not observed.

Dataset  $\Rightarrow x_1, x_2, \dots, x_N \Rightarrow$  observed data

$z_1, z_2, \dots, z_N \Rightarrow$  latent variable

MLE  $\Rightarrow \theta^* = \underset{\theta}{\operatorname{argmax}} \ln P(x|\theta)$

$z \Rightarrow K$ -dimension vector  $\Rightarrow$  if  $z_{ik}=1 \Rightarrow x$  is drawn from  $k^{th}$  gaussian component.

$P(z_k=1|x, \theta) \Rightarrow$  Posterior probability of a data points comp = 1

$$\hookrightarrow = \frac{\pi_K N(x | \mu_K, \Sigma_K)}{\sum_{k'=1}^K \pi_{k'} N(x | \mu_{k'}, \Sigma_{k'})}$$

$P(z|x, \theta) \Rightarrow$  Posterior over latent variables.

$$\begin{aligned} P(z|x, \theta) &= \prod_{n=1}^N P(z_n|x, \theta) \\ &= \prod_{n=1}^N \prod_{k=1}^K P(z_{nk}=1|x, \theta) \end{aligned}$$

$$\ln P(x|\theta) = \ln \sum_z P(x, z|\theta) \Rightarrow \text{Maximize this objective w.r.t } \theta.$$

But there is a  $\sum_z$  inside log which does not lead to a closed form solution.

So we maximize  $E[\ln P(x, z|\theta)]$ , w.r.t.  $P(z|x, \theta)$

$$Q(\theta, \theta^{old}) \rightarrow E[\ln P(x, z|\theta)] \text{ w.r.t. } P(z|x, \theta^{old})$$

① E-step consists of calculating  $P(z|x, \theta^{old})$  for dataset  $X$ .

② M-step  $\Rightarrow Q(\theta, \theta^{old}) = E_{P(z|x, \theta^{old})} [\ln P(x, z|\theta)]$

Find  $\theta$  maximizing  $Q$  function

Instead of maximizing  $\ln \sum_z P(x, z|\theta)$ , i.e.  $\ln P(x|\theta)$  we maximize  $Q(\theta, \theta^{old})$

-  $Q(\theta, \theta^{old})$  will keep the summation of all "z", but outside the logarithm.

- Easy to find if the components are known, i.e. summation is outside the log.

For GMM, M-step

$$\mu_k^{new} = \frac{1}{N_k} \sum_{n=1}^N P(z_{nk}=1 | x_n, \theta^{old}) \times x_n$$

$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N P(z_{nk}=1 | x_n, \theta^{old}) \times [x_n - \mu_k][x_n - \mu_k]^T$$

$$\pi_{ik}^{new} = \frac{N_k}{N}, \quad N_k = \sum_{n=1}^N P(z_{nk}=1 | x_n | \theta^{old})$$

$$P(x, z | \theta) = P(z | x, \theta) \times P(x | \theta) \Rightarrow \text{Bayes rule}$$

$$\ln P(x, z | \theta) = \ln P(z | x, \theta) + \ln P(x | \theta)$$

$$\ln P(x | \theta) = \ln P(x, z | \theta) - \ln P(z | x, \theta)$$

$q(z) \Rightarrow$  arbitrary distribution over latent variables

$$z_n \rightarrow z_{n1}, z_{n2}, \dots, z_{nk}$$

0.4      0.2      0.4

$$\mathbb{E}_{q(z)} [\ln P(x | \theta)] = \mathbb{E}_{q(z)} [\ln [P(x, z | \theta)]] - \mathbb{E}_{q(z)} [\ln P(z | x, \theta)]$$

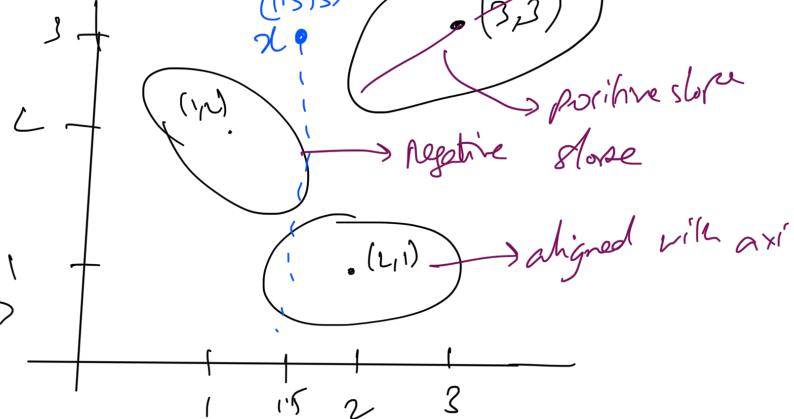
$$\ln P(x | \theta) = \mathbb{E}_{q(z)} \left[ \ln \frac{P(x, z | \theta)}{q(z)} \right] + \mathbb{E}_{q(z)} \left[ \ln \frac{q(z)}{P(z | x, \theta)} \right]$$

$$\ln (P(x | \theta)) \leq \underbrace{\mathcal{L}(z, \theta)}_{\text{lower bound}} + \overbrace{\mathbb{E}_{q(z)} [\ln \frac{q(z)}{P(z | x, \theta)}]}^{\geq 0 \text{ (always)}}$$

$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 1 & -0.8 \\ 0.8 & 1 \end{bmatrix}$$



Taking a point  $x = [1.5, 3]$ , we get the probability density function as

$[0.011, 0.015, 0.019]$	& after normalization,	
$\downarrow$	$\downarrow$	$\downarrow$
Caussian 1	Caussian 2	Caussian 3
$[0.25, 0.33, 0.41]$		

$$\ln \sum_{k=1}^K \pi_k \times P(x | \mu_k, \Sigma_k)$$

Let  $\pi = [\pi_1, \pi_2, \pi_3] = \text{initial}$

$$\ln \frac{1}{3} (0.011 + 0.015 + 0.019) \Rightarrow \text{Original objective}$$

↳ Not maximizable due to the  $\sum$  inside log

$$P(x, z_1 = 1 | \theta) \quad P(x, z_2 = 1 | \theta) \quad P(x, z_3 = 1 | \theta)$$

$$E[\ln 0.011] \times 0.25 + \ln[0.015] \times 0.33 + \ln[0.019] \times 0.41$$

$P(z_m | \theta)$    
 Represents joint dist. of observed & latent variable  
 weighted by corresponding posterior probability

$$P(x, z | \theta) = P(z | x, \theta) \times P(x | \theta)$$

$$\ln P(x, z | \theta) = \ln P(z | x, \theta) + \ln P(x | \theta)$$

$$\ln P(x | \theta) = \ln P(x, z | \theta) - \ln P(z | x, \theta)$$

↳ posterior

$q(z) \Rightarrow$  Arbitrary probability dist. over latent variable

$$z_n \rightarrow z_{n1}, z_{n2}, \dots, z_{nk}$$

$0.4 \quad 0.2 \quad 0.4$

$$E_{q(z)} \left[ \ln p(x|\theta) \right] = E_{q(z)} \left[ \ln p(x, z|\theta) \right] - E_{q(z)} \left[ \ln p(z|x, \theta) \right]$$

↙ No,  $z$  in this term. So remains same.

$$\ln p(x|\theta) = E \left[ \ln \left( \frac{p(x, z|\theta)}{q(z)} \right) \right] - E \left\{ \ln \frac{q(z)}{p(z|x, \theta)} \right\}$$

The effect of this division, is that we have added entropy of  $q(z)$  ( $+H(q)$ ). So subtracting entropy from second term will keep it same.

$$\ln p(x|\theta) = L(q, \theta) + KL(q||p)$$

↙ Lower bound ↘ always  $\geq 0$

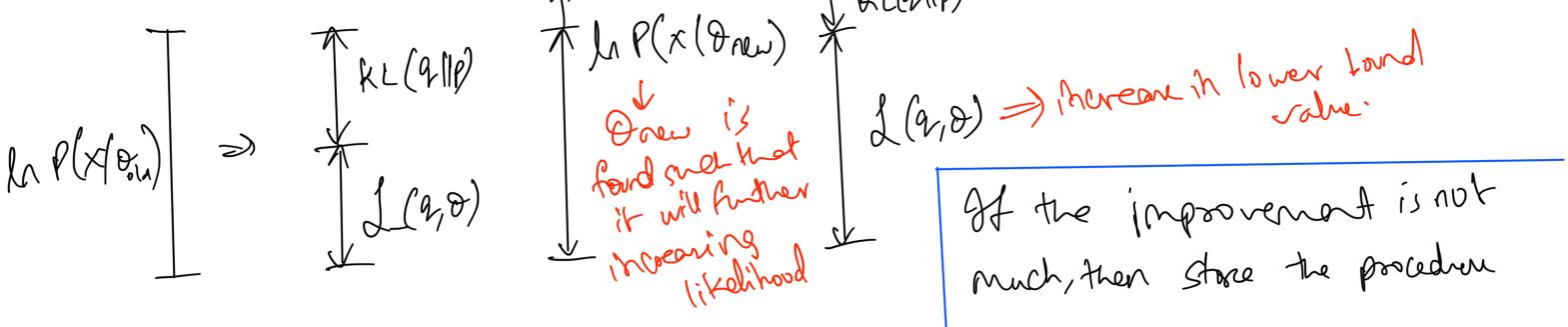
optimized w.r.t  $q$  in E-step & w.r.t.  $\theta$  in M-step

→ When  $KL(q||p)$  is more, then  $\ln p(x|\theta)$  &  $L(q, \theta)$  are more distant with respect to each other.  
↓ vice-versa.

So, when  $KL(q||p)$  is closer to zero,  $\ln p(x|\theta)$  &  $L(q, \theta)$  will have similar like values; ie. a tight lower bound. i.e. if it is exactly equal to LNS.

$$\theta' = \underset{\theta}{\operatorname{argmax}} L(q, \theta)$$

$\rightarrow p(z|x, \theta)$



## Mixture of Bernoulli Distribution

CMM  $\Rightarrow$  mixtures of distributions over continuous variables

$\Rightarrow$  mixtures of distributions over discrete binary variables.

↳ latent class analysis.

Consider a set of  $D$  binary variables,  $x_i \Rightarrow i=1 \text{ to } D$ , with corresponding  $\mu_i$ .

$$x = (x_1, \dots, x_D)^T \text{ & } \mu = (\mu_1, \dots, \mu_D)^T$$

$$P(x|\mu) = \prod_{i=1}^D \mu_i^{x_i} (1-\mu_i)^{1-x_i}$$

$$E[x] = \mu$$

$$\text{Cov}[x] = \text{diag}\{\mu_i(1-\mu_i)\}$$

Now, consider a finite mixture of these distributions given by,

$$P(x|\mu, \pi) = \sum_{k=1}^K \pi_k P(x|\mu_k), \text{ & } P(x|\mu_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1-\mu_{ki})^{1-x_i}$$

$$E[x] = \sum_{k=1}^K \pi_k \mu_k$$

$$\text{Cov}[x] = \sum_{k=1}^K \pi_k \left\{ \sum_{i=1}^D \mu_{ki} \mu_{ki}^T \right\} - E[x] E[x]^T$$

$\hookrightarrow = \text{diag}\{\pi_k \mu_k (1-\mu_k)\}$

$$\log \text{likelihood} = \ln P(x|\mu, \pi) = \sum_{n=1}^N \ln \sum_{k=1}^K \pi_k P(x_n|\mu_k)$$

$\downarrow$

$\{\text{inside log. Use EM.}$

latent variable  $z \Rightarrow$  one hot  $K-d$  vector.

$$P(x|z, \mu) = \prod_{k=1}^K P(x|z_k, \mu_k)^{z_{ik}}$$

$$P(z|\pi) = \prod_{k=1}^K \pi_k^{z_{ik}}$$

$\hookrightarrow$  prior probs-

$$\ln P(x, z|\mu, \pi) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ \ln \pi_k + \sum_{i=1}^D x_{ni} \ln \mu_{ki} + (1-x_{ni}) \ln (1-\mu_{ki}) \right\}$$

$$E_z \left[ \ln p(x, z | \mu, \pi) \right] = \sum_{n=1}^N \sum_{k=1}^K r(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^D x_{ni} \ln \mu_{ki} + (1-x_{ni}) \ln (1-\mu_{ki}) \right\}$$

$$\begin{aligned} \gamma_{nk} &= E[z_{nk}] = \frac{\sum_{z_{nk}} z_{nk} [\pi_k p(x_n | \mu_k)]}{\sum_{z_{nj}} [\pi_j p(x_n | \mu_j)]^{z_{nj}}} \\ &= \frac{\pi_k p(x_n | \mu_k)}{\sum_{j=1}^K \pi_j p(x_n | \mu_j)} \end{aligned}$$

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$

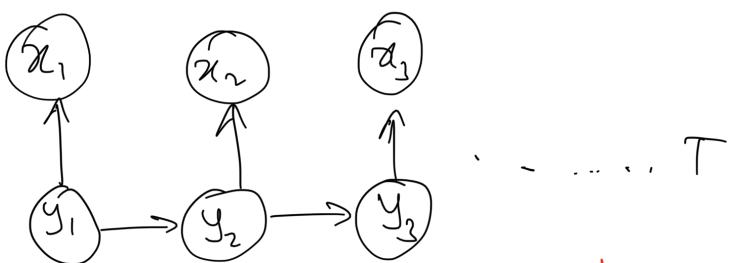
$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n$$

$$\pi_k = \frac{N_k}{N}$$

Hidden Markov Model

→ used for seq2seq classification

like NER tagging, for tagging

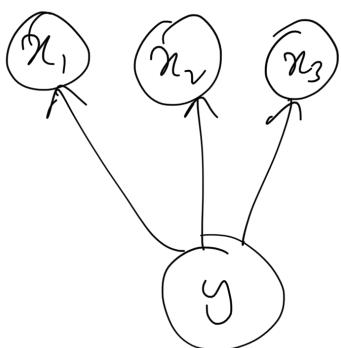


$$\Rightarrow P(x, y) = P(y_1) \times P(x_1 | y_1) \times \prod_{i=2}^T P(x_i | y_i) \times P(y_i | y_{i-1})$$

$$= \prod_{i=1}^T P(x_i | y_i) \times P(y_i | y_{i-1})$$

↳ Size of model =  $L \times k + L \times L$

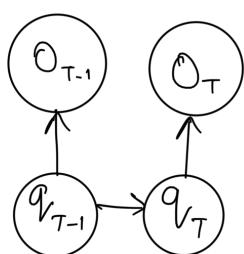
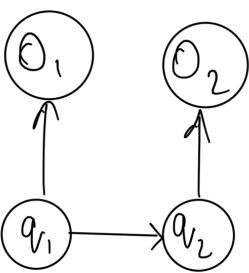
for naive Bayes  $\Rightarrow$



$$\Rightarrow P(x, y) = \frac{P(y) \times P(x | y)}{\prod_{i=1}^k P(x_i | y)}$$

Model Size  $\Rightarrow L \times k + L$

A tutorial of HMM by Robinson



$P(o_1 = 1 | q_0) = \pi_1$ ,  
 $P(o_1 = 2 | q_0) = \pi_2$

$\pi = [\pi_1, \pi_2, \dots, \pi_N]$

initial probabilities

$q_t = 1, \dots, N$   
 ↳ represents state at time "t"

$P(o, q | \lambda) \rightarrow$  To find  
 ↳  $\lambda = \{A, B, \pi\}$

$N \Rightarrow$  no. of distinct states

$$A = N \begin{bmatrix} & \\ & \\ & \end{bmatrix}, \quad B = N \begin{bmatrix} & \\ & \\ & \end{bmatrix} \rightarrow p(i|j)$$

$M =$  no. of distinct observation symbols

$$P(O|\lambda) = \sum_{q \in Q} P(O, q|\lambda)$$

↳ No. of possible state sequences =  $N^T$

$\checkmark$   $T$  observations, with each having  $N$  possible states.

→ so computationally expensive

Brute force  $\Rightarrow$  Compute prob. of all state sequences

$$\text{Joint probability} \Rightarrow P(O, q|\lambda) = \pi[q_0] \times \prod_{t=1}^T P(q_t | q_{t-1}) \times P(O_t | q_t)$$

Forward probability

$$\alpha_t(i) = P(O_1, O_2, \dots, O_t, q_t = i | \lambda)$$

$$\text{Initialisation} \Rightarrow \alpha_1(i) = P(O_1, q_1 = i | \lambda)$$

$$\alpha_1(i) = \pi_i \times B_i(O_1) \quad \text{→ } i^{\text{th}} \text{ row & } O_1^{\text{th}} \text{ column of } B.$$

prob. of being in state  $i$  initially & observing  $O_1$

$$\alpha_t(i) \circ$$

$$\alpha_{t+1}(i) = \left[ \sum_{j=1}^N \alpha_t(j) \times A_{ji} \right] \times B_i(O_{t+1})$$

$$\alpha_t(2) \circ$$

$$\overset{i}{\alpha_{t+1}(i)}$$

Computational effort =  $T N^2$

$$\alpha_t(N) \circ$$

$$P(O/\lambda) = \sum_{i=1}^N \alpha_T(i) \quad , \quad \alpha_T(i) = P(O_1, O_2, \dots, O_T, q_T=i/\lambda)$$

## Backward Probability

$$\beta_t(i) = P(O_{t+1}, \dots, O_T / q_t=i, \lambda)$$

$$\beta_T(i) = 1$$

$$\beta_t(i) = \sum_{j=1}^N A_{ij} \times \beta_{t+1}(j) B_j(O_{t+1})$$

$$\alpha_t(i) = P(O_1, O_2, \dots, O_t, q_t=i/\lambda)$$

$$\beta_t(i) = P(O_{t+1}, \dots, O_T / q_t=i, \lambda)$$

$$P(O_1, O_2, \dots, O_T, q_t=i/\lambda) = \alpha_t(i) \times \beta_t(i)$$

$$P(q_t=i/O, \lambda) = \frac{\alpha_t(i) \times \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \times \beta_t(j)} = r_t(i)$$

The most probable state sequence given the observed sequences.

Viterbi Algorithm  $\rightarrow$  For most probable state sequence  
 solves inference problem.

EM Algorithm  $\rightarrow$  To train, to find parameters.

$$\mathcal{S}_t(i) = \pi_i \times \beta_i(\theta_t)$$

$$\mathcal{S}_{t+1}(j) = \max_{i=1\dots N} \mathcal{S}_t(i) \times A_{ij} \times \beta_j(\theta_{t+1})$$

$$\psi_{t+1}(j) = \operatorname{argmax} \mathcal{S}_t(i) A_{ij} \times \beta_j(\theta_{t+1})$$

# HMM not class

Notations  $\Rightarrow$  time instants associated with state changes,  
 $t=1, 2, 3, \dots$

State at  $t=t \Rightarrow q_t$

$a_{ij} = P(q_t = S_j / q_{t-1} = S_i) \Rightarrow$  transition probability from state  $S_i$  at  $q_{t-1}$  to  $S_j$  at  $q_t$ .

$N = \text{no. of states}$

$$a_{ij} \geq 0 \quad \& \quad \sum_{j=1}^N a_{ij} = 1$$

$\downarrow$

$A = \text{transition prob matrix, size } N \times N$

$\pi \Rightarrow \text{initial state probability vector, size } N \times 1$

$$\pi_i = P(q_1 = S_i) \quad \xrightarrow{\text{say } S_i}$$

$\Rightarrow$  Given that the model is in a known state, what is the probability that it stays in that same state for exactly "d" days?

$$\Omega = \{S_1, S_2, S_3, \dots, S_d, S_{d+1}, \dots, S_N\}$$

$$P(\Omega/\text{Model}) = \pi_i \times (a_{ii})^{d-1} \times (1 - a_{ii}) = p_i(d)$$

$\hookrightarrow 1$  (as it is given that at  $t=1 \Rightarrow q_t = S_i$ )

$\Rightarrow$  discrete probability density function of duration d in state i.

Based on  $p_i(d)$ , we can calculate the expected no. of observations (duration) in a state, conditioned on starting in that state w,

$$\bar{d}_i = \sum_{d=1}^{\infty} d P_i(d)$$

$$= \sum_{d=1}^{\infty} d (a_{ii})^{d-1} (1-a_{ii}) = \frac{1}{1-a_{ii}}$$

$\therefore$  Expected no. of consecutive days in  $S_i$  given

$$\text{start is } S_i = \frac{1}{1-a_{ii}}$$

→ Elements of a HMM

①  $N \Rightarrow$  no. of states in the model.

$S = \{S_1, S_2, \dots, S_N\}$  & state at time  $t \Rightarrow q_t$   
 $(N \text{ can be the vocabulary size})$

②  $M \Rightarrow$  no. of distinct observation symbols per state.

Congress to the physical obj of the system being modeled.

$V = \{V_1, V_2, \dots, V_M\}$  ( $M$  can be the no. of NER tags)

③ State transition prob. dist.  $A = \{a_{ij}\}$  of size  $N \times N$ .

$$a_{ij} = P[q_t = S_j | q_{t-1} = S_i]$$

④ Observation symbol prob. dist. in state  $j$ ,  $B = \{b_j(k)\}$  of size  $N \times M$ . (also called emission prob from hidden obj)

$$b_j(k) = P[V_k \text{ at } t | q_t = S_j]$$

⑤ Initial state distribution,  $\pi = \{\pi_i\}$  of size  $N \times 1$

$$\pi_i = P(q_1 = S_i)$$

→ Given appropriate values of  $N, M, A, B, \pi$ , HMM can be used as a generator to give an observation sequence

$$O = O_1, O_2, \dots, O_T,$$

where each  $O_t$  is a symbol from  $V$ ,  $T$  is the number of observations in the sequence.

Steps for generating

- ① Choose initial state  $q_1 = s_1$  according to  $\pi$ .
- ② Set  $t = 1$
- ③ Choose  $O_t = v_k$  acc. to symbol prob. dist state  $q_i \Rightarrow b_i(v_k)$
- ④ Transit to new state  $q_{t+1} = s_j$  according to state prob. dist for state  $s_i \Rightarrow a_{ij}$
- ⑤ Set  $t = t+1$  & return to step 3 if  $t \leq T$ .  
else terminate.

$$\lambda = [A, B, \pi]$$

Parameter set of the model

C. The Three Basic Problems for HMMs<sup>5</sup>

Given the form of HMM of the previous section, there are three basic problems of interest that must be solved for the model to be useful in real-world applications. These problems are the following:

Problem 1: Given the observation sequence  $O = O_1, O_2, \dots, O_T$ , and a model  $\lambda = (A, B, \pi)$ , how do we efficiently compute  $P(O|\lambda)$ , the probability of the observation sequence, given the model?

Problem 2: Given the observation sequence  $O = O_1, O_2, \dots, O_T$ , and the model  $\lambda$ , how do we choose a corresponding state sequence  $Q = q_1, q_2, \dots, q_T$  which is optimal in some meaningful sense (i.e., best "explains" the observations)?

Problem 3: How do we adjust the model parameters  $\lambda = (A, B, \pi)$  to maximize  $P(O|\lambda)$ ?

→ Evaluation problem, i.e. given a model & a sequence, how likely is the sequence generated by the model.  
(OR)

Scoring how well a given model suits a given obs. seqn.

→ Uncovering the hidden part of the model i.e. to find correct state sequence.

→ Optimize the model parameters to best describe how a given obs. seq. comes about.

To fix ideas, consider the following simple isolated word speech recognizer. For each word of a  $W$ -word vocabulary, we want to design a separate  $N$ -state HMM. We represent the speech signal of a given word as a time sequence of coded spectral vectors. We assume that the coding is done using a spectral codebook with  $M$  unique spectral vectors; hence each observation is the index of the spectral vector closest (in some spectral sense) to the original speech signal. Thus, for each vocabulary word, we have a training sequence consisting of a number of repetitions of sequences of codebook indices of the word (by one or more talkers). The first task is to build individual word models. This task is done by using the solution to Problem 3 to optimally estimate model parameters for each word model. To develop an understanding of the physical meaning of the model states, we use the solution to Problem 2 to segment each of the word training sequences into states, and then study the properties of the spectral vectors that lead to the observations occurring in each state. The goal here would be to make refinements on the model (e.g., more states, different codebook size, etc.) so as to improve its capability of modeling the spoken word sequences. Finally, once the set of  $W$  HMMs has been designed and optimized and thoroughly studied, recognition of an unknown word is performed using the solution to Problem 1 to score each word model based upon the given test observation sequence, and select the word whose model score is highest (i.e., the highest likelihood).

high probability of belonging to the Gaussian component centered at its location, and the responsibility of that component for that data point will be close to 1. On the other hand, the responsibilities of the other components will be close to 0 because their probability density functions are becoming very narrow and are centered at other locations. As a result, the responsibilities of other data points will approach zero except for the data point that is closest to the Gaussian component that is collapsing onto it.

In a hmm ner tagger, what are the input, hidden and output states

In a Hidden Markov Model (HMM) NER Tagger, the input state is the sequence of words in a sentence. The output state is the sequence of tags assigned to each word in the sentence, which correspond to named entities such as person names, locations, and organizations.

The hidden state is the sequence of unknown, underlying states that generate the output sequence. These hidden states are often referred to as the "hidden Markov chain" or "latent variable". In the context of NER, the hidden states represent the true labels for each word, which are not directly observable and must be inferred from the observed data.

## Solution to Problem 1 $P(O|\lambda)$

Brute force  $\Rightarrow$  Enumerating every possible state sequence of length  $T$  (no. of observations)

$$Q = q_1 q_2 \dots q_T \Rightarrow \text{If there are } "N" \text{ states, no. of such sequences} = N^T$$

$$\Rightarrow P(Q | \lambda) = \prod_{t=1}^T P(q_t | q_{t-1}, \lambda)$$

$$O = o_1 o_2 \dots o_T$$

We assumed statistical independence of observations,

$$P(Q | \lambda) = b_{q_1}(o_1) \times b_{q_2}(o_2) \times \dots \times b_{q_T}(o_T)$$

$\downarrow$   
emission probability of token represented by  $o_t$  at time  $t=1$  from hidden state  $q_t$ .

$$P(Q | \lambda) = \pi_{q_1} \times a_{q_1 q_2} \times a_{q_2 q_3} \times \dots \times a_{q_{T-1} q_T}$$

Joint probability of  $O$  &  $Q$ , i.e. prob. that  $O$  &  $Q$  occur simultaneously  $\Rightarrow P(O, Q|\lambda) = P(O|Q, \lambda) \times P(Q|\lambda)$

The probability of  $O$  given  $\lambda$ , is obtained by summing the joint prob. over all possible state sequences  $q$ , giving;

$$P(O|\lambda) = \sum_{\text{all } Q} P(O|Q, \lambda) \times P(Q)$$

$$= \sum_{q_1 q_2 \dots q_T} \pi_{q_1} \cdot b_{q_1}(o_1) \times a_{q_1 q_2} b_{q_2}(o_2) \times \dots \times a_{q_{T-1} q_T} b_{q_T}(o_T)$$

Interpretation of above computation & Computational Complexity

The interpretation of the computation in the above equation is the following. Initially (at time  $t = 1$ ) we are in state  $q_1$  with probability  $\pi_{q_1}$ , and generate the symbol  $O_1$  (in this state) with probability  $b_{q_1}(O_1)$ . The clock changes from time  $t$  to  $t + 1$  ( $t = 2$ ) and we make a transition to state  $q_2$  from state  $q_1$  with probability  $a_{q_1 q_2}$  and generate symbol  $O_2$  with probability  $b_{q_2}(O_2)$ . This process continues in this manner until we make the last transition (at time  $T$ ) from state  $q_{T-1}$  to state  $q_T$  with probability  $a_{q_{T-1} q_T}$  and generate symbol  $O_T$  with probability  $b_{q_T}(O_T)$ .

A little thought should convince the reader that the calculation of  $P(O|\lambda)$ , according to its direct definition (17) involves on the order of  $2T \cdot N^T$  calculations, since at every  $t = 1, 2, \dots, T$ , there are  $N$  possible states which can be reached (i.e., there are  $N^T$  possible state sequences), and for each such state sequence about  $2T$  calculations are required for each term in the sum of (17). (To be precise, we need  $(2T - 1)N^T$  multiplications, and  $N^T - 1$  additions.) This calculation is computationally unfeasible, even for small values of  $N$  and  $T$ ; e.g., for  $N = 5$  (states),  $T = 100$  (observations), there are on the order of  $2 \cdot 100 \cdot 5^{100} \approx 10^{72}$  computations

and observations are the input words as they are observed

Picture a HMM like neural network, where hidden states are analogous to neurons. Each input goes to a specific hidden state from which a specific output is generated. hidden states are usually the outputs, like NER & POS tags. While neurons are just hidden computational units.

Efficient Solution  $\Rightarrow$  Forward - Backward Procedure

to model  $P(O|\lambda)$

$$\alpha_t(i) = \underbrace{P(O_1, O_2, \dots, O_t, q_t = s_i | \lambda)}_{\text{Probability of the partial observation sequence } O_1, O_2, \dots, O_t \text{ (until time } t\text{)}}$$

& state  $s_i$  at time time  $t$ .

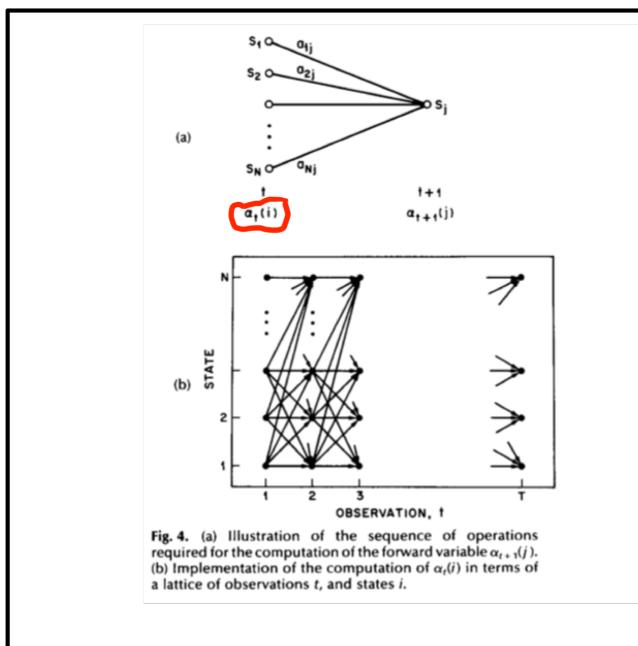
### ① Initialization

$$\alpha_1(i) = \underbrace{\prod_{j=1}^N b_j(O_1)}_{P(O_1)} \quad \text{for all } i, \text{ such that } 1 \leq i \leq N$$

This step initializes the forward probabilities as the joint probabilities of state  $s_i$  & initial observation  $O_1$ .

### ② Induction

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) \times a_{ij} \right] \times b_j(O_{t+1}), \quad 1 \leq t \leq T-1, \quad 1 \leq j \leq N$$



\* We only need  $\alpha_T(i)$  for all possible  $i$  to solve problem 1, but  $\beta_T(i)$  is introduced next, just because it is used in solving problem 2.

- Figure shows how state  $S_j$  can be reached at time  $t+1$  from the  $N$  possible states  $S_i, 1 \leq i \leq N$  at time  $t$ .
- As  $\alpha_t(i)$  is the prob. of joint event that  $O_1, O_2, \dots, O_b$  are observed and the state at time  $t$  is  $S_i$ , the product  $\alpha_t(i) \times a_{ij}$  is then the probability of the joint event that  $O_1, O_2, \dots, O_b$  are observed & state  $S_j$  is reached at time  $t+1$  via state  $S_i$  at time  $t$ .
- Summing this product over all " $N$ " possible states, results in the prob. of  $S_j$  at time  $t+1$  with all accompanying previous partial observations.
- Once this is done,  $\alpha_{t+1}(j)$  is obtained by accounting for  $O_{t+1}$  in state  $j$ , i.e.  $b_j(O_{t+1})$ .
- The computation of  $\alpha_{t+1}(j)$  is done for all states  $j$ ,  $1 \leq j \leq N$ ; for a given time  $t$ .

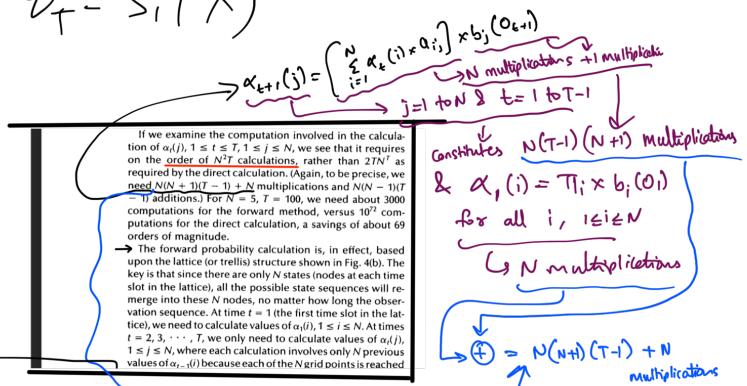
### ③ Termination

$$P(O/x) = \sum_{i=1}^N \alpha_T(i)$$

- Step 3 gives  $P(O/x)$ , because by definition,

$$\alpha_T(i) = P(O, O_2, \dots, O_T, S_T = S_i/x)$$

→ Computational Complexity  $\Rightarrow N^2 T$



→ In the same manner, consider a backward variable  $\beta_t(i)$  defined as

$$\beta_t(i) = P(O_{t+1} O_{t+2} \dots O_T / q_t = s_i / \lambda)$$

↓  
Probability of the partial observation sequence from  $t+1$  to the end, given state  $s_i$  at time  $t$  & model  $\lambda$ .

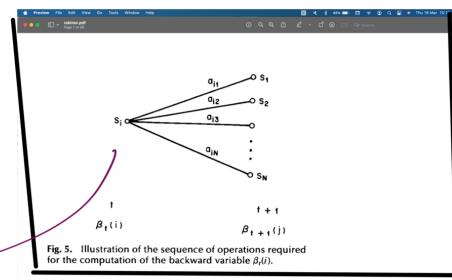
## ① Initialization

$$\beta_T(i) = 1, \quad 1 \leq i \leq N$$

## ② Induction

$$\beta_t(i) = \sum_{j=1}^N \alpha_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$

$$, t = T-1, T-2, \dots, 1 \quad \& \quad 1 \leq i \leq N$$



In order to have been at state  $s_i$  at time  $t$ , and to account for the observation sequence from time  $t+1$  on, you have to consider all possible states  $s_j$  at time  $t+1$ , accounting for the transitions from  $s_i$  to  $s_j$  as well as the observation  $O_{t+1}$  in state  $j$  & then account for the remaining partial observation from state  $j$ .

Computational Complexity  $\Rightarrow \underline{N^2 T}$

## Solution to Problem 2

$$\rightarrow P(O|\theta, \lambda)$$

- Problem 1 has an exact solution which can be given.  
 But there are several possible ways to find the "optimal" state sequence. The difficulty lies with the definition of the optimal state sequence.

One such optimality criterion is to choose the states  $q_t$  which are individually most likely.  $\rightarrow$  a principle or standard by which something is judged

↳ This optimality criterion maximizes the expected no. of correct individual states.

To implement this solution, we define the variable,

$$\gamma_t(i) = P(q_t = s_i | O, \lambda)$$

→ probability of being in state  $s_i$  at time  $t$ , given the observation sequence  $O$  & model  $\lambda$ .

$$\gamma_t(i) = \frac{\alpha_t(i) \times \beta_t(i)}{\sum_{i=1}^N \alpha_t(i) \times \beta_t(i)} = \frac{P(O_1, O_2, \dots, O_T, q_t = s_i | \lambda)}{P(O | \lambda)}$$

accounts for partial observation sequence  $O_1, O_2, \dots, O_T$  & state  $s_i$  at  $t$

→ accounts for partial obs. seq.  $O_{t+1}, \dots, O_T$  given state  $s_i$  at  $t$ .

→ Normalization factor makes  $\gamma_t(i)$  a probability measure such that  $\sum_{i=1}^N \gamma_t(i) = 1$

- Using  $\gamma_t(i)$ , we solve for the individually most likely state  $q_t$  at time  $t$ , as,

$$q_t = \operatorname{argmax}_{1 \leq i \leq N} [\gamma_t(i)], \text{ for all } t, 1 \leq t \leq T$$

- Although the above optimality criterion maximizes the expected no. of correct states (by choosing the most likely state for each  $t$ ), there could be some problems with the resulting state sequence.
- This is when some state transitions may have zero probability, i.e.  $a_{ij}=0$  for some  $i \neq j$ .
- So the "optimal" state sequence may not even be a valid state sequence, as the optimality criterion above just choose the most likely state at every instant, without regard to the probability of occurrence of sequences of states.
- So, we need to modify the criterion.  
 For example, we could solve for the state sequence that maximizes the expected number of correct pairs of states  $(q_t, q_{t+1})$  or triples of states  $(q_t, q_{t+1}, q_{t+2})$ , etc.
- These criteria might be reasonable for some applications, but the most widely used criterion is to find the single best state sequence (path), i.e. to maximize  $P(Q|O, \lambda)$  which is equivalent to maximizing  $P(Q, O|\lambda)$ .  
 ↳ Formal technique for this is the Viterbi Algorithm based on Dynamic Programming.

**Viterbi Algorithm** → To find the single best sequence,  $Q = \{q_1, q_2, \dots, q_T\}$  for the given observation sequence  $O = \{o_1, o_2, \dots, o_T\}$ , we define,

$$S_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1, q_2, \dots, q_t = i, o_1, o_2, \dots, o_t / \lambda)$$

the best score (highest prob.) along a single path, at time  $t$ , which accounts for first " $t$ " observation & ends at state  $q_i$ .

By induction we have,  $\delta_{t+1}(j) = \left[ \max_i \delta_t(i) a_{ij} \right] \times b_j(o_{t+1})$

To actually retrieve the state sequence, we need to keep track of the argument which maximized for each  $t$  &  $j$ . Store these in the array  $\psi_t(j)$

### ① Initialization

$$\delta_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

$$\psi_1(i) = 0$$

### ② Recursion

$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(o_t), \quad 2 \leq t \leq T \quad \& \quad 1 \leq j \leq N$$

→ At each instant of time, for each hidden state, calculate prob till then. i.e. you are keeping track of all hidden states possible from where  $s_j$  is reached at time "t-1" at time "t".

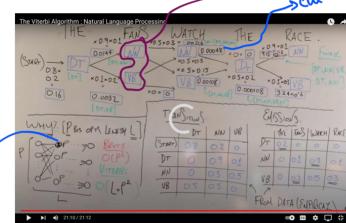
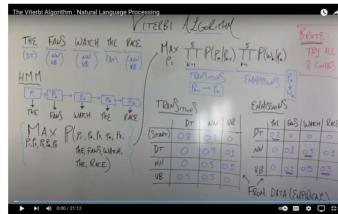
$$\psi_t(j) = \underset{1 \leq i \leq N}{\text{argmax}} [\delta_{t-1}(i) a_{ij}], \quad 2 \leq t \leq T \quad \& \quad 1 \leq j \leq N$$

↓ index which gives max value

### ③ Termination

$$P^* = \max_{1 \leq i \leq N} [\delta_T(i)]$$

$$q_T^* = \underset{1 \leq i \leq N}{\text{argmax}} [\delta_T(i)]$$



### ④ Path (state seq.) backtracking

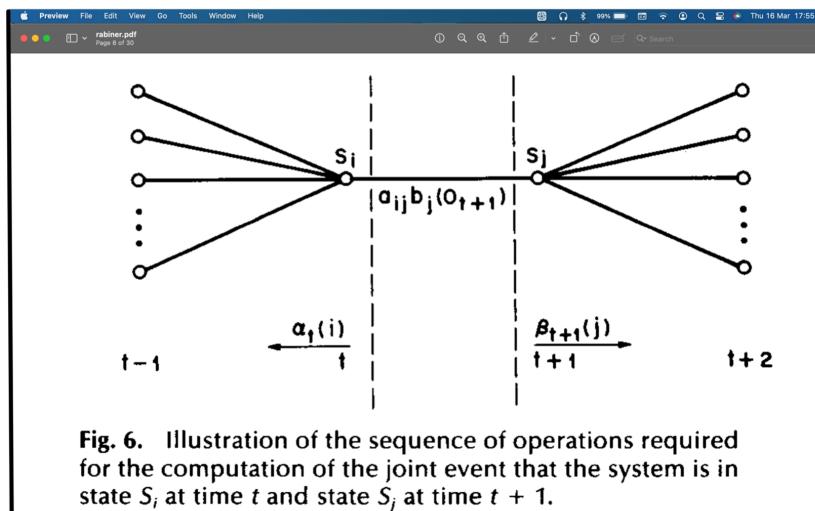
$$q_t^* = \psi_{t+1}[q_{t+1}^*], \quad t = T-1, T-2, \dots, 1$$

→ Similar to forward prob. Calculation. Major difference is max over summation in the forward prob. calculation.

Computational complexity of Viterbi  $\Rightarrow O(N^2 T)$

## Solution to Problem 3

- Most difficult problem of HMM's is to determine a method to adjust the model parameters  $(A, B, \pi)$  to maximize the prob. of observation sequence given to the model.  
 Global optimum  $\rightarrow$  ?  $\Rightarrow$  No optimal way  
 Local optimum  $\rightarrow$  Yes  $\Rightarrow$  Baum-Welch method based on EM.
- In order to describe the procedure for re-estimation i.e. iterative update & improvement of HMM parameters, we first define  $\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$   $\curvearrowright$   
 The seq. of events leading to the condition here are illustrated below.



- From the definition of forward & backward probs, i.e.  $\alpha_t(i)$  &  $\beta_t(i)$ , we can express  $\xi_t(i, j)$  as,

$$\begin{aligned}
 \xi_t(i, j) &= \frac{\alpha_t(i) \times a_{ij} \times b_j(O_{t+1}) \times \beta_{t+1}(j)}{P(O|\lambda)} \\
 &= \frac{\alpha'_t(i) \times a_{ij} \times b_j(O_{t+1}) \times \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) \times a_{ij} \times b_j(O_{t+1}) \times \beta_{t+1}(j)} \Rightarrow P(q_t=S_i, q_{t+1}=S_j, O|\lambda) \\
 &\hookrightarrow P(O|\lambda) \Rightarrow \text{can also be accounted by } \sum_{i=1}^N \alpha_t(i)
 \end{aligned}$$

$\gamma_t(i)$  = Prob. of being in State  $S_i$  at time  $t$ ; given  $O$  &  $X$ .

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

- If we sum over  $\gamma_t(i)$ , we get a quantity which can be interpreted as the expected no. of times that state  $S_i$  is visited, (or) equivalently, the expected no. of transitions made from state  $S_i$ , if we exclude  $t=T$  from the summation.

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected no. of transition from } S_i$$

- Similarly summation over  $\xi_t(i, j)$  over  $t$ , gives expected no. of transitions from state  $S_i$  to state  $S_j$ .

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{expected no. of transitions from } S_i \text{ to } S_j$$

### Reestimation formulas

$$\bar{\gamma}_i = \text{expected frequency in state } S_i \text{ at time } t=1 = \gamma_1(i)$$

$$\bar{\alpha}_{ij} = \frac{\text{expected no. of transitions from } S_i \text{ to } S_j}{\text{Total expected no. of transitions from } S_i} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$$\bar{b}_j(k) = \frac{\text{expected no. of times in state } S_j \text{ & observing } V_k}{\text{Total expected no. of times in state } S_j}$$

$$= \sum_{t=1}^T \gamma_t(j)$$

Such that

$$O_t = V_k$$

$$\sum_{t=1}^T \gamma_t(j)$$

- It has been proven that,  $\bar{\lambda} = \{\bar{A}, \bar{B}, \bar{\pi}\}$  is a model such that,
 $P(O|\bar{\lambda}) \geq P(O|\lambda)$ 

Equal only if  $\bar{\lambda} = \lambda$ , ie  $\lambda$  is a critical point of the likelihood function.

→ Obs. seq. is more likely to have been produced from model  $\bar{\lambda}$ ; when compared to  $\lambda$ .
- Repeatedly doing this will eventually end in a maximum likelihood estimate of the HMM. (local maxima only)

$$Q(\lambda, \bar{\lambda}) = \sum_Q P(Q|O, \lambda) \log [P(O, Q|\bar{\lambda})]$$

$$\max_{\bar{\lambda}} [Q(\lambda, \bar{\lambda})] \Rightarrow P(O|\bar{\lambda}) \geq P(O|\lambda)$$

→ Important point to note is that, at every update

$$\sum_{i=1}^N \bar{\pi}_i = 1 \quad \& \quad \sum_{i=1}^N \bar{a}_{ij} = 1 \quad \& \quad \sum_{k=1}^M \bar{b}_i(k) = 1$$

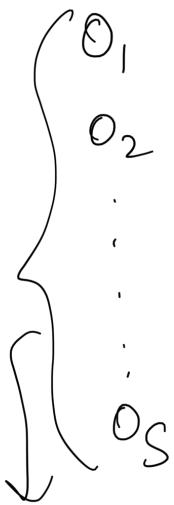
This can be achieved using lagrange multipliers.

By using lagrange multipliers, it is shown that  $P$  is maximized when,

$$\bar{\pi}_i = \pi_i \frac{\partial P}{\partial \bar{\pi}_i} \quad , \quad \bar{a}_{ij} = \frac{a_{ij} \frac{\partial P}{\partial a_{ij}}}{\sum_{k=1}^N \bar{a}_{ik} \frac{\partial P}{\partial \bar{a}_{ik}}} \quad \& \quad \bar{b}_i(k) = b_i(k) \frac{\partial P}{\partial b_i(k)}$$

$$\bar{b}_i(k) = b_i(k) \frac{\partial P}{\partial b_i(k)} \overrightarrow{\frac{\sum_{l=1}^M b_l(l) \frac{\partial P}{\partial b_l(l)}}{}}$$

## Dataset



"S" no. of sequences, each of length "T"

$O_s = \text{Sentence / Sequence } i^*$

$$P(O, q | \lambda) = \prod_{t=1}^T P(O_t | q_t) \times P(q_t | q_{t-1})$$

↗ parameters come from matrix B  
 ↗ parameters come from matrix A & vector  $\pi$

$$= \pi_{q_1} \times B_{q_1}(O_1) \times \prod_{t=2}^T P(O_t | q_t) \times P(q_t | q_{t-1})$$

$$P(O | \lambda) = \sum_{q \in Q} P(O, q | \lambda) = \sum_{q \in Q} P(O | q, \lambda) \times P(q | \lambda)$$

$\hookrightarrow N^T$  possibilities

$$\sum_{s=1}^S \log P(O_s | \lambda) = \log \sum_{q \in Q} P(O_s, q | \lambda)$$

E-step  $\Rightarrow P(q | O_s, \lambda^t)$

M-step  $\Rightarrow E_{P(q | O_s, \lambda)} [\log P(O_s, q | \lambda)]$

Expectation  
 ↓  
 using posterior dist. to maximize  $\log P(O_s, q | \lambda)$

$$\log P(O_{1:T} | \lambda) = \log \pi_{q_1} + \left[ \prod_{t=1}^T P(O_t | q_t) \right] \left[ \prod_{t=2}^T P(q_t | q_{t-1}) \right]$$

$\hookrightarrow_B$        $\hookrightarrow_A$

$$E_{P(\hat{q}|O, \lambda^{old})} [\log \pi_i] + \sum_{t=1}^T E_{P(O_t|O, \lambda^{old})} [\log P(O_t | q_t)] + \sum_{t=2}^T E_{P(q_t|O, \lambda^{old})} [\log P(q_t | q_{t-1})]$$

$$\downarrow \sum_{i=1}^N [\log \pi_i] \ell(q_1=i|O, \lambda^{old}) - g \left( \sum_{i=1}^N \pi_i - 1 \right)$$

$\hookrightarrow$  lagrangian multiplier

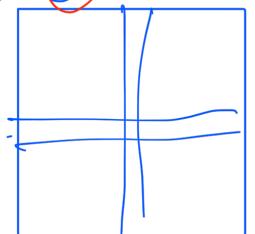
$\frac{\partial}{\partial \pi_i}$

$$\frac{\gamma_i(i)}{\pi_i} - g = 0 \Rightarrow \pi_i = \frac{\gamma_i(i)}{g} \quad \& \quad g = \sum_{i=1}^N \gamma_i(i)$$

$$B_{ij} = P(O=j | q=i)$$

M

$$\boxed{\pi_i = \frac{\gamma_i(i)}{\sum_{i=1}^N \gamma_i(i)}}$$



$$E_{P(O|O, \lambda^{old})} \left[ \sum_{t=1}^T \log P(O_t | q_t) \right] - \sum_{i=1}^N h_i \left( \sum_{j=1}^M B_{ij} - 1 \right)$$

Each row is a probability vector

$$= \sum_{t=1}^T \sum_{i=1}^N \left[ (\log P(O_t | q_t = i)) \right] P(q_t = i | O, \lambda^{old})$$

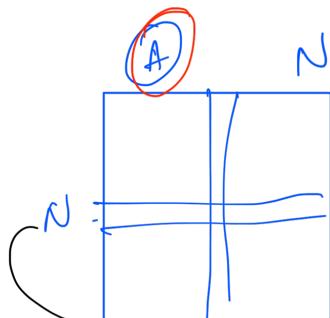
$$E_{P(q|O, \lambda^{old})} \left[ \sum_{m=1}^M \sum_{t=1}^T \sum_{i=1}^N \log P(O_t | q_t = i) \times I(O_t = m) \right]$$

$$- \sum_{i=1}^N h_i \left( \sum_{j=1}^M B_{ij} - 1 \right)$$

$$\frac{\partial}{\partial B_{ij}} = \sum_{t=1}^T \frac{I(O_t = m)}{B_{im}} \gamma_t(i) - h_i = 0$$

$$B_{im} = \frac{\sum_{t=1}^T I(\theta_t = m) \times \gamma_t(i)}{h_i}$$

$$h_i = \sum_{m=1}^M \sum_{t=1}^T I(\theta_t = m) \times \gamma_t(i)$$



$$\frac{\partial}{\partial A_{ij}}$$

$$\sum_{t=2}^T E_{P(q_t | O, \lambda^{old})} [\log P(q_t / q_{t-1})] - \sum_{i=1}^N k_i \left[ \sum_{j=1}^N A_{ij} - 1 \right]$$

$$\sum_{t=1}^{T-1} \sum_{i=1}^N \sum_{j=1}^N \left[ \log A_{q_t=i, q_{t+1}=j} \right] \times P(q_t=i, q_{t+1}=j | O, \lambda)$$

$\xi_t(i, j)$

This is because summation is from  $t=1$  here

$$\sum_{t=1}^{T-1} \frac{I(q_t=i) \ I(q_{t+1}=j) \ \xi_t(i, j)}{A_{ij}} - k_i = 0$$

$$A_{ij} = \frac{\sum_{t=1}^{T-1} I(q_t=i) \times I(q_{t+1}=j) \ \xi_t(i, j)}{k_i}$$

$$\begin{aligned} k_i &= \sum_{j=1}^N \sum_{t=1}^{T-1} I(q_t=i) \ I(q_{t+1}=j) \ \times \xi_t(i, j) \\ &= \sum_{t=1}^{T-1} \gamma_t(i) \end{aligned}$$

$$Y_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

Weak-supervised learning  $\Rightarrow$  learning from ~~labelled~~ features

K-means, CMM, general EM

HMM  $\Rightarrow$  parameter estimation

End  
Ki  
Sadko

Transfer learning using spectral learning

SVM, PCA, SVD  $\rightarrow$  EM for PLSA

Mixture models, Latent Dirichlet Allocation,  
Gibbs Sampling, Variational Bayes

## Graphical Models

$V \Rightarrow$  Vocab

Docs.  
 $\rightarrow w_{11} \quad w_{1n_1} \quad w_{12} \quad \dots \quad \dots \quad w_{1N_1}$   
 $\rightarrow w_{21} \quad w_{2n_2} \quad w_{22} \quad \dots \quad \dots \quad \dots \quad w_{2N_2}$   
 $\vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots$   
 $\rightarrow w_{D1} \quad w_{D2} \quad w_{D3} \quad \dots \quad \dots \quad w_{DN_D}$

index of word from Vocab  $V$

$w_{ij} = j^{\text{th}}$  word in  $i^{\text{th}}$  document - (index of this word)

Each doc. of diff. length

$D = \text{no. of docs.}$

$N_D = \text{no. of words in doc. } D$

Multinomial  $\Rightarrow$  Multi( $d|\theta$ )  
 $\hookrightarrow [\theta_1, \theta_2, \dots, \theta_{N_D}]$   
 $\hookrightarrow \text{prob. of all words}$

$$\text{Multi}(d/e) = \prod_{n=1}^{N_d} \theta_{w_n}$$

$$\text{Joint likelihood} = \prod_{d=1}^D \prod_{n=1}^{N_d} \theta_{w_{dn}}$$

$$\text{Joint log likelihood} = \sum_{d=1}^D \sum_{n=1}^{N_d} \log \theta_{w_{dn}} - \lambda \left[ \sum_{w=1}^{|V|} \theta_w - 1 \right]$$

(Lagrange multiplier)

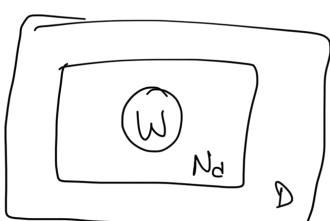
$$\prod_{n=1}^{N_d} \theta_{w_n} = \prod_{w=1}^{|V|} \theta_w^{N_{dw}} \quad [N_{dw} = \text{no. of times "w" occurs in doc. } d]$$

$$\therefore JLL = \sum_{d=1}^D \sum_{w=1}^{|V|} \log \theta_w^{N_{dw}} - \lambda \left[ \sum_{w=1}^{|V|} \theta_w - 1 \right]$$

$$\frac{\partial JLL}{\partial \theta_w} = \sum_{d=1}^D \frac{N_{dw}}{\theta_w} - \lambda = 0$$

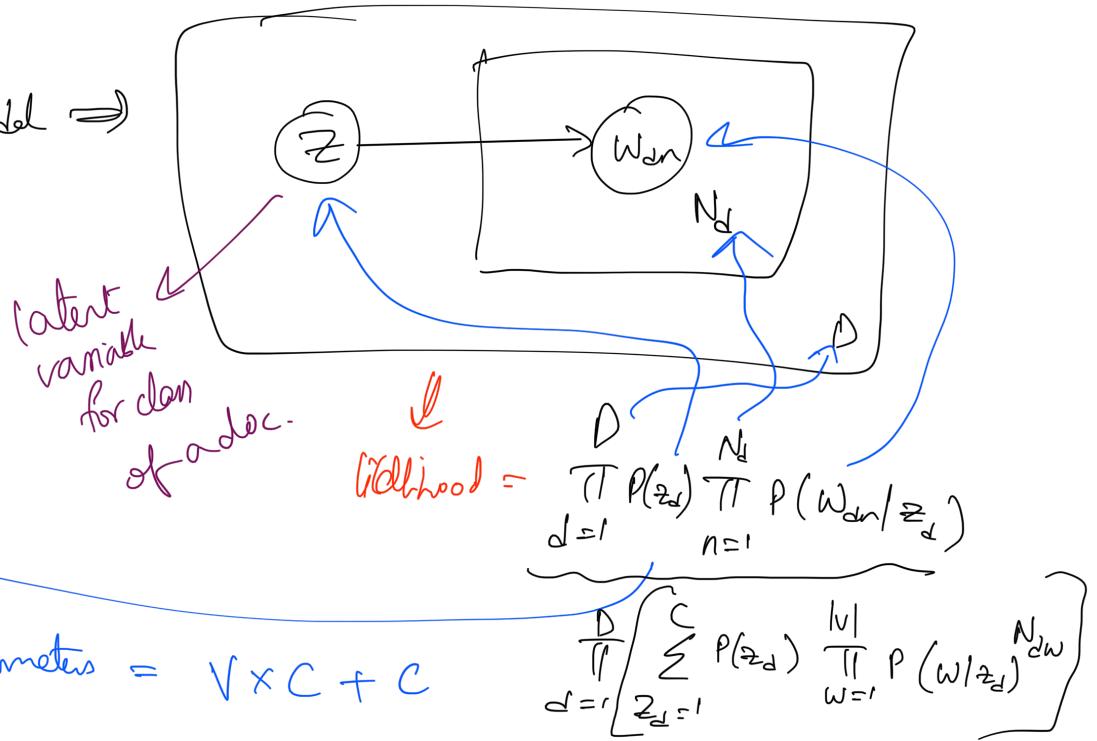
$$\therefore \theta_w = \frac{\sum_{d=1}^D N_{dw}}{\lambda} \Rightarrow \lambda = N \rightarrow \text{Total no. of tokens}$$

$$\therefore \theta_w = \frac{\sum_{d=1}^D N_{dw}}{N}$$



Read plate Notation  $\Rightarrow$

for mixture model  $\Rightarrow$



$$\text{No. of Parameters} = V \times C + C$$

$$JLL = \sum_{d=1}^D \left[ \log \sum_{z_d=1}^C P(z_d) \times \frac{1}{W} \prod_{w=1}^W P(w | z_d)^{N_{dw}} \right]$$

$$E\text{-step: } P(z_d | w_d, \theta^{old})$$

$$M\text{-step: } E_{P(z_d | w_d, \theta^{old})} \left[ \sum_{d=1}^D \log P(z_d, w_d | \theta) \right]$$

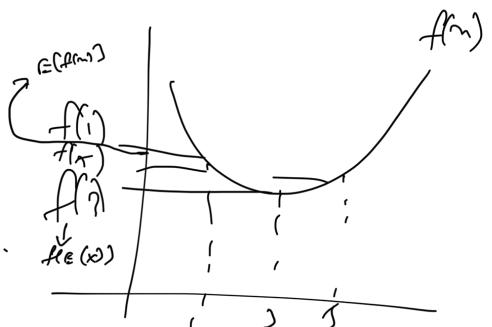




## Jensen's Inequality

Let  $f$  be a convex function

$\hookrightarrow$  i.e.  $f''(n) \geq 0$ , for all  $n$ .



Let  $X$  be a random variable,

$$\text{Then } f(E[X]) \leq E[f(X)]$$

If  $f(n)$  is strictly convex,  $\rightarrow f''(n) > 0$

$$\text{then } E[f(X)] = f(E[X])$$

$$\begin{aligned} p(n=1) &= \frac{1}{L} \\ p(n=5) &= \frac{1}{L} \\ \text{mean} &= f(E[n]) \\ &= 3 \end{aligned}$$

$$E[n] = \frac{1+5}{2} = 3$$

$$\begin{aligned} E[f(n)] &= \frac{1}{2} \times f(1) + \frac{1}{2} \times f(5) \\ &= 3 \end{aligned}$$

$\Rightarrow$  If  $f$  is a concave function

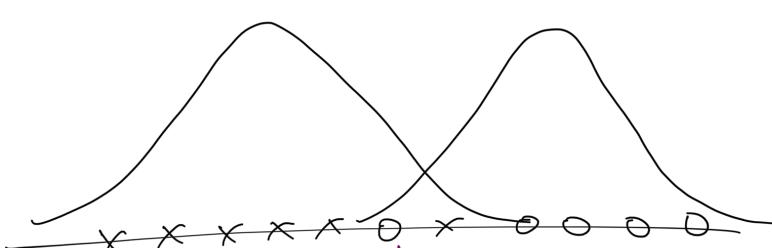
$\hookrightarrow f''(n) \leq 0$ , for all  $n$

$$\& f(E[X]) \geq E[f(X)]$$

If  $f$  is convex  $\Rightarrow -f$  is concave

EM

1-D Example



If we know the class then we can calculate the  $M_k$  &  $S_k$  for each cluster. Else if, we know the  $M_k$  &  $S_k$ , we can predict the classes

for each  $n$ .

The problem here is that we don't know both.  
Classic chicken & eggs problem.

Mixture of Gaussian Model.

Suppose that a latent (hidden/unobserved) random variable  $Z$  &  $x^{(i)}, z^{(i)}$  are distributed,

$$p(x^{(i)}, z^{(i)}) = p(x^{(i)}|z^{(i)}) \times p(z^{(i)})$$

where  $z^{(i)} \sim \text{Multinomial}(\phi)$   
 $\phi$  parameters

$$x^{(i)}|z^{(i)}=j \sim N(\mu_j, \Sigma_j) \quad z \in \{1, 2, \dots, k\}$$

If we know the  $z^{(i)}$ 's we use MLE:

$$\ell(\phi, \mu, \Sigma) = \sum_{i=1}^m \log p(x^{(i)}, z^{(i)} | \phi, \mu, \Sigma)$$

$$\phi_j = \frac{1}{m} \sum_{i=1}^m 1 \times I(z^{(i)}=j)$$

$$\mu_j = \frac{\sum_{i=1}^m 1 \times I(z^{(i)}=j) x^{(i)}}{\sum_{i=1}^m 1 \times I(z^{(i)}=j)}$$

E-step (true value of  $z^{(i)}$ 's)

$w_j^{(i)}$  is how much  
 $z^{(i)}$  is assigned to  
the Gaussian "j"

Set  $w_j^{(i)} = P(z^{(i)}=j | x^{(i)} | \phi, \mu, \Sigma)$

$$= P(x^{(i)} | z^{(i)}=j) \times P(z^{(i)}=j) \rightarrow z \sim \text{Multinomial}(\phi)$$

$$\sum_{l=1}^k P(x^{(i)} | z^{(i)}=l) P(z^{(i)}=l)$$

$\sqrt{(\mu_j, \Sigma_j)}$

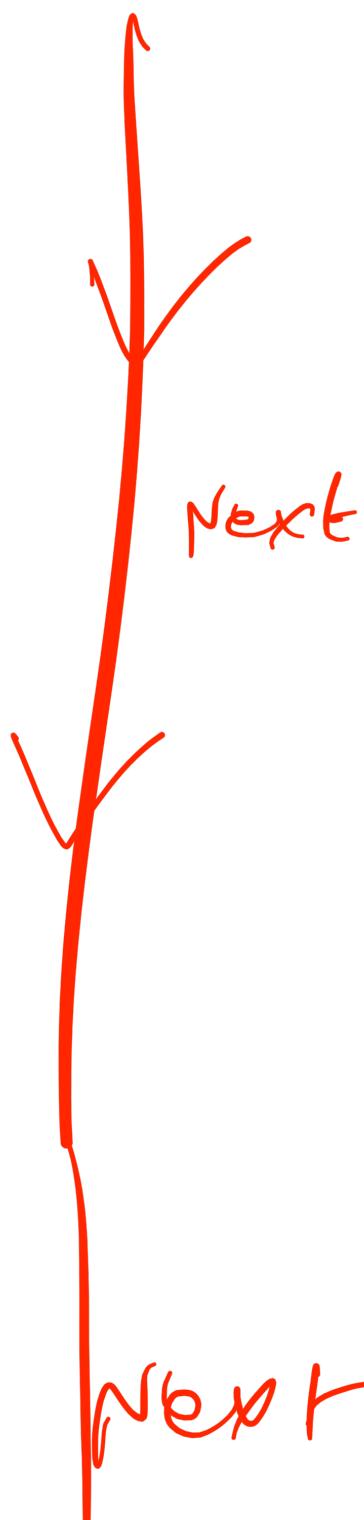
M-step

$$\phi_j = \frac{1}{m} \sum_{i=1}^m w_j^{(i)}$$

$$\mu_j = \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}}$$

$$w_j^{(i)} = \mathbb{E} [1 \times \mathbb{I}(z^{(i)} = j)]$$

E<sup>M</sup>  $\Rightarrow$  soft clustering assignment.



enr

Used in anomaly detection.

## Density estimation

Model  $P(n) \Rightarrow$  Estimate the density from which the dataset points came from.  
 heat  
 vibration  
 Anomaly  
 if  $P(n) < \epsilon \Rightarrow$  flag anomaly.

Mixture of gaussians for each feature.

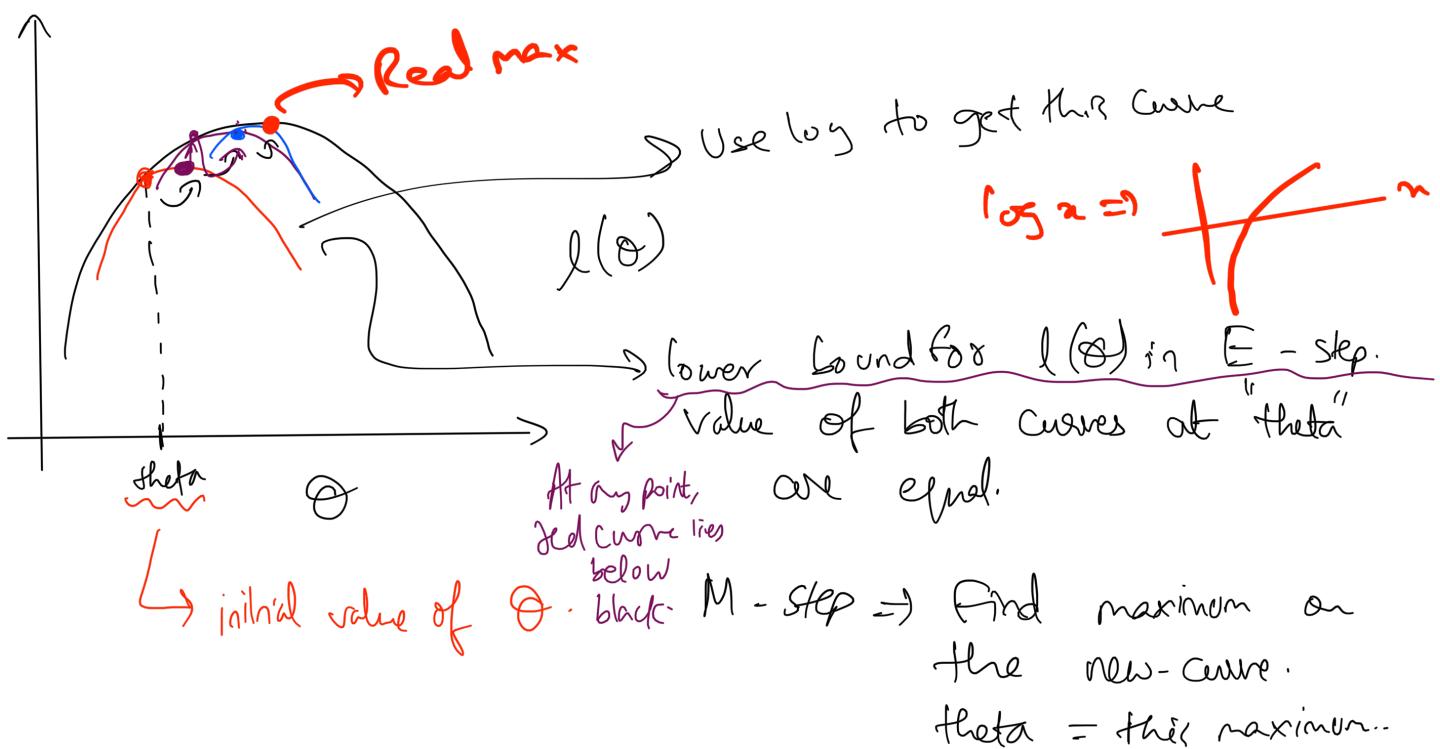
Read ML notes  $\Rightarrow$  for Anomaly detection

Have model for  $P(n, z | \theta)$

Only observe  $X$   $\{x^{(1)}, \dots, x^{(m)}\}$

$$l(\theta) = \sum_{i=1}^m \log P(x^{(i)}, z^{(i)} | \theta)$$

$$= \sum_{i=1}^m \log \sum_{z^{(i)}} P(n^{(i)}, z^{(i)} | \theta)$$



$$\max_{\theta} \sum_i \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)} | \theta) \rightarrow l(\theta)$$

$$= \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})}$$

where  $Q_i(z^{(i)})$  is a probability distribution

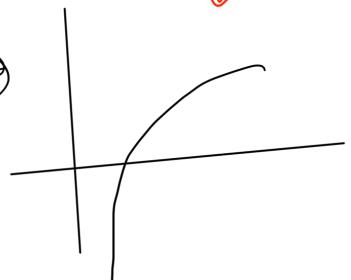
$$\text{i.e. } \sum_{z^{(i)}} Q_i(z^{(i)}) = 1$$

$$\geq \sum_i \log \mathbb{E}_{z^{(i)} \sim Q_i} \left[ \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})} \right]$$

$$= \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})}$$

i.e.,  $f(E[x]) \geq E[f(x)]$

( $\log x \geq$ )



If  $z = \{1, \dots, 10\}$   
Roll a 10-sided dice

then

$$\mathbb{E}[g(z)] = \sum_z P(z) g(z)$$

$$\mathbb{E}[z] = \sum_z P(z) \cdot z$$

→ On a given iteration of EM (with params  $\theta$ )

we want,  
 $\log \mathbb{E}_{z^{(i)} \sim Q_i} \left[ \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})} \right] = \mathbb{E}_{z^{(i)} \sim Q_i} \left[ \log \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})} \right]$

For this to hold, we need

$$\frac{P(x^{(i)}, z^{(i)})}{Q_i(z^{(i)})} = \text{Constant}$$

Set  $\sum_{z^{(i)}} Q_i(z^{(i)}) = 1$

$$Q_i(z^{(i)}) = \frac{P(x^{(i)}, z^{(i)} | \theta)}{\sum_{z^{(i)}} P(x^{(i)}, z^{(i)} | \theta)}$$

After derivation,

$$Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}, \theta)$$

E-step  $\Rightarrow$  Set  $Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}, \theta)$

M-step  $\Rightarrow \theta = \underset{\theta}{\operatorname{argmax}} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)} | \theta)}{Q_i(z^{(i)})}$

## General

⇒ The gamma function, denoted by  $\Gamma_n$  is a generalization of the factorial func. for non-integer values of  $n$ .

$$\Gamma_n = \int_0^\infty t^{(n-1)} e^{-t} dt$$

⇒ A simplex is geometric object defined by a set of  $n+1$  vertices in  $n$ -dimensional space.

Eg = triangle in 2-d space, defined by 3 vertices.  
tetrahedron in 3-d space, defined by 4 vertices.

⇒ Dirichlet distribution is a probability dist. on the simplex. It is generalization of  $\beta$ -dist. in  $n$ -dimensions.

$$f(m|\alpha) = \frac{\Gamma(\alpha_0)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K m_i^{\alpha_i - 1}$$

$$\alpha_0 = \sum_{i=1}^K \alpha_i$$

Dirichlet is used as a prior for categorical data, where the simplex represents the set of possible categorical outcomes.

⇒ Multinomial Distribution

$$P(X=n) = \frac{n!}{x_1! x_2! x_3! \dots x_K!} \times P_1^{x_1} P_2^{x_2} \dots P_K^{x_K}$$

