

In the absence of any prior knowledge about the shape of  $p(x)$  we can model it as a Gaussian. ✓

Many times we can identify that the given observation  $x$  has been generated through a structured process.

That is, the distribution over the features  $p(x)$  has an underlying structure.

If our model for  $p(x)$  can capture this structure, then the model will be closely following (approximating) the underlying data generating distribution and will likely perform well.

An underlying structure to the data generating process would imply the presence of additional hidden (not observed) properties of data that describe (define) that structure.

These hidden attributes (properties) are the latent attributes.

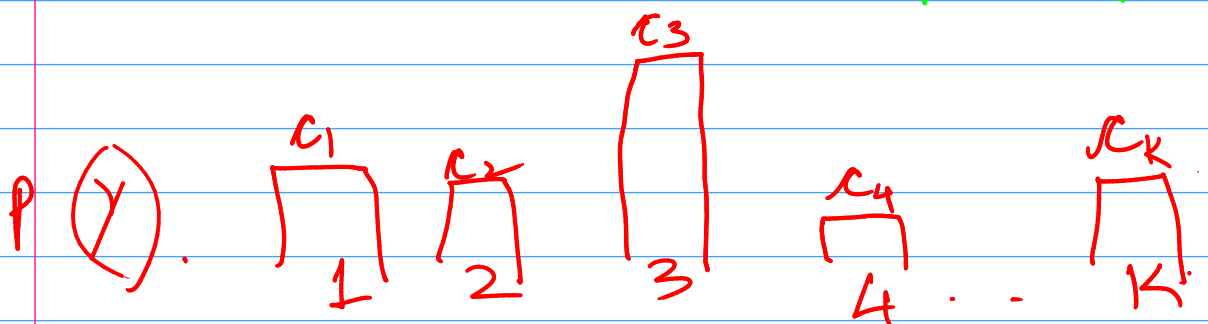
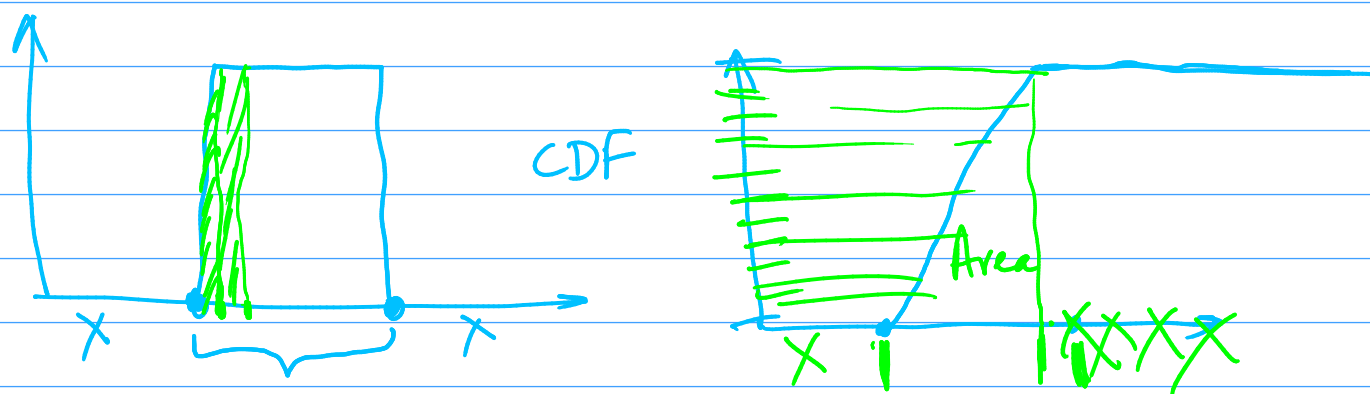
We don't know the values assigned to the latent attributes. These are unknown parameters of our model.

We treat them as uncertain (random) variables which follow an unknown distribution.



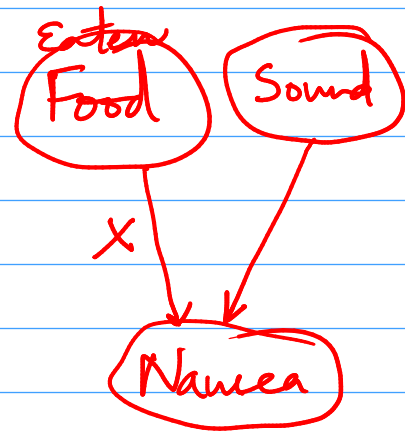
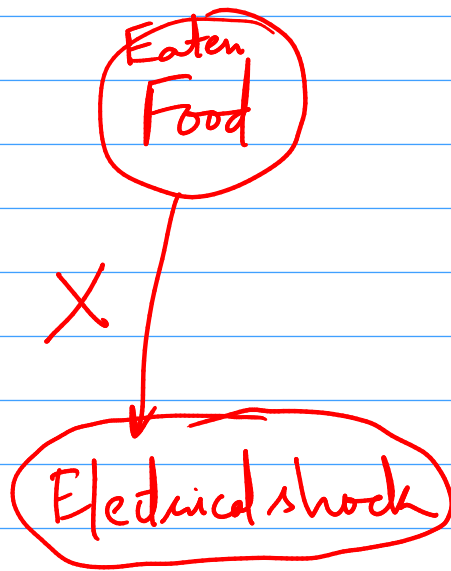
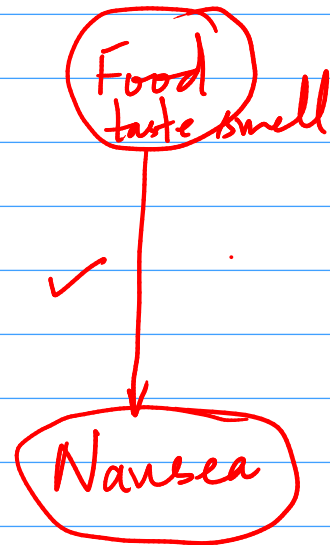
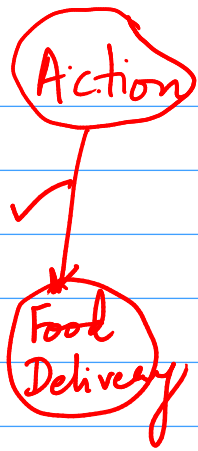
Given a  $n \sim U$  uniform distribution

Cumulative distribution



$$\sum_y c_y = 1.$$

Absurd  
Useless  
Hypothesis



Inductive bias.

Inductive learning: <sup>learning from</sup> Examples  
(training learning)



A simple form of structure is categorized examples <sup>(observation)</sup>  
 That is, the features that we observe belong to some category  $Y = y_i$ .  
 (Not a target variable) <sup>Latent variable</sup>  $Y$  <sup>Category</sup>

Such a structure is denoted graphically as hidden  $Y$   $\downarrow$   $X$  <sup>feature observed</sup>  
 $Y$  is a random variable taking discrete values.  
 $Y \in [1, K] \quad 1, 2, 3, \dots, K.$  <sup>labels</sup>

$X$  is a random variable denoting the observation vector  $\underline{x}_i$   
 $\underline{x} \in \mathbb{R}^d$

To generate an observation, we first sample from a categorical variable  $Y = y_i$  and then given the category  $y_i$  we sample an observation  $\underline{x}_i$  from the conditional distribution  $P(X = \underline{x}_i | Y = y_i)$  ✓

Random variable  $Y$  has a prior distribution associated with it  $P(Y = y_i)$  ✓. Since  $Y$  is discrete valued, we can have a categorical distribution to describe  $P(Y = y)$ .

Using such a structured model, the probability of observing example  $\underline{x}_i$  can be written as  $P(X = \underline{x}_i, Y = y_i)$

unsupervised  $P(X = \underline{x}_i) = \sum_{y_i=1}^K P(Y = y_i) P(X = \underline{x}_i | Y = y_i)$   
 i.e. summing out the joint distribution over the random variable  $Y$ .

$$P(a) = \sum_b P(a, b)$$

$$= \sum_{y_i=1}^K \underbrace{C_{y_i}}_{\text{Gaussian}} \frac{1}{\sqrt{(2\pi)^d |\Sigma_{y_i}|}}$$

$$\exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu}_{y_i})^T \Sigma_{y_i}^{-1} (\underline{x} - \underline{\mu}_{y_i})\right)$$

Multi variate Gaussian  $d$ -dim space.

Here  $P\{X = \underline{x}_i | Y = y_i\}$  has been modelled as a Gaussian distribution for every value of  $Y = y_i$ .

Since  $Y$  can take  $k$  values, there will be  $k$  such Gaussian distributions, with parameters  $(\underline{\mu}_{y_i}, \underline{\Sigma}_{y_i})$  ✓

While referring to the parameters of these Gaussians, we drop the example specific subscript  $i$ , and simply write

$(\underline{\mu}_y, \underline{\Sigma}_y)$   
 $d \times 1$  vector  
 $d \times d$  matrix

Though we have introduced a structure for modelling the generative distribution  $P(X)$ , is there any advantage in assuming this structure? incorporating prior knowledge

And, how do we solve for the parameters  $\Rightarrow$  better learning.

$\bigcirc Y$   $P(Y)$  prior  
 $\bigcirc X$   $P(X|Y)$  likelihood of  $x$ ..  
 Where  $\underline{\pi}_y$  are the parameters of the categorical distribution  $P\{Y=y\}$  and  $\underline{\mu}_y, \underline{\Sigma}_y$  are the parameters of the conditional distribution  $P\{X=\underline{x} | Y=y\}$  modeled as a multivariate

Gaussian.

We formulate our solution  $\underline{\theta}$  to be the one that maximizes the likelihood of observing the examples in the training set  $S = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_m\}$  MLE

The likelihood of the training set  $S$

$$\underline{P}_{\underline{\theta}}[S] = \prod_{i=1}^m \underline{P}_{\underline{\theta}}[X = \underline{x}_i]$$

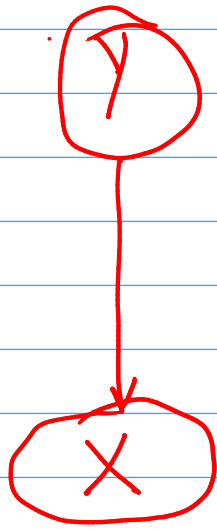
$\underline{x}_i$

independence assumption



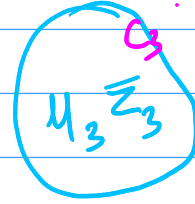
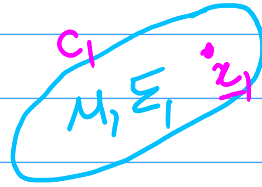
# Gaussian Mixture Model (GMM)

Assume  
value  $K=5$



$P(X)$   
is maximized.

$$\sum_y c_y = 1$$



Hyper parameter  
 $K=5$

$K=3$

$K=8$

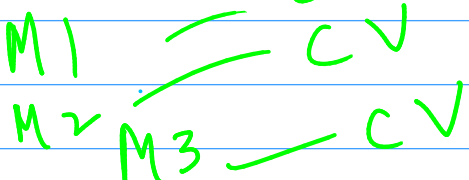


$P(X)$

Try different values  
values of  $K$ .  
check the fit quality

Penalty: Complexity of the model.

Assessment  
trade off score.  $(\text{fit quality} + \lambda \text{ complexity})$   
(error)  
(penalty).



(training (source data)).

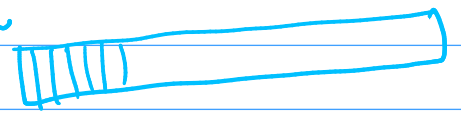
K-fold

cross validation

Leave one out cross validation

$$S = \{x_1, \dots, x_m\}$$

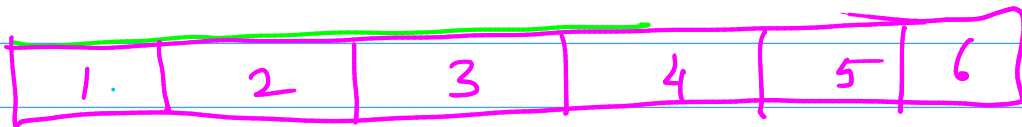
$K=m$



K subsets

$K=6$

Model



Model fixed

Keep aside one subset

Error

Session 1

Train(2,3,4,5,6)

Test(1)  $e_1$

2.

Train(1,3,4,5,6)

Test(2)  $e_2$

3

Train(1,2,4,5,6)

Test(3)  $e_3$

$\vdots$

Test

$K$

Performance:

Avg( $e_i$ ).

(Generalization performance)

$S: \left( \begin{array}{cc} \text{70\% Train} & \text{30\% Test} \\ S_1 & S_2 \end{array} \right)$  one-time partition.

Error( $S_2$ )

For tuning hyper parameters

$K=3$

KFCV — Test performance

$K=4$

KFCV — "

$K=5$

]

To prevent numeric underflow, we take the log likelihood.

$$L(\underline{\theta}) = \log \left( \prod_{i=1}^m P_{\underline{\theta}}[X = \underline{x}_i] \right)$$

likelihood function.

$$= \sum_{i=1}^m \log P_{\underline{\theta}}[X = \underline{x}_i]$$

Joint distribution.  
(given by the structural model).

$$\arg \max_{\underline{\theta}} L(\underline{\theta}) = \sum_{i=1}^m \log \left( \sum_{y_i=1}^k P_{\underline{\theta}} \left[ \underbrace{X = \underline{x}_i}_{\text{obs}} , \underbrace{Y = y_i}_{\text{hidden}} \right] \right)$$

i.e. we write  $P[X = \underline{x}_i]$  as a marginalization of the latent variable  $Y$  over the joint distribution of observed and latent variables.   
 (summing out)

The optimal parameters for the model are the ones that maximize the likelihood function.

$\theta$ : hyperparameter

$$\underline{\theta}^* = \arg \max_{\underline{\theta}} L(\underline{\theta}) = \arg \max_{\underline{\theta}} \sum_{i=1}^m \log \left( \sum_{y_i=1}^k P_{\underline{\theta}}[X = \underline{x}_i | Y = y_i] \right)$$

exp  
Gaussian

$\underline{\theta}$ : hyperparameter

This maximization is hard because the log acts on the summation.

Therefore, we maximize  $L(\underline{\theta})$  by formulating a surrogate function  $G(\underline{Q}, \underline{\theta})$  which takes in another parameter.   
 a matrix  $\underline{Q}$  of size  $m \times k$

Every row of the matrix  $\underline{Q}$  is a categorical distribution over  $K$  discrete values.  
 $\sum \text{row entries} = 1$

$$\underline{Q} = \begin{matrix} & \begin{matrix} 1 & 2 & \dots & k \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ m \end{matrix} & \begin{bmatrix} & & & \\ & & & \\ & & & \\ \textcolor{red}{i} & & & \\ & & & \\ & & & \\ m & p_1^m & p_2^m & p_3^m & p_k^m \end{bmatrix} \end{matrix}$$

$m \times k$   
posterior  
 $p(Y = y_i | X = \underline{x}_i)$



We are uncertain about the hidden variables.  
 We cannot observe them.  
 We capture this uncertainty by defining a posterior distribution over  $Y = y_i$  given every training example  $X = x_i$ .



$$P[Y = y_i] = \text{prior}$$

$$P[Y = y_i | X = x_i] : \text{posterior}$$



$$P[X = x_i | Y = y_i] : \text{likelihood}$$

✓ The matrix  $Q$  captures the posterior distribution

$$G(Q, \theta)$$

$$Q_{i,y} = P[Y = y_i | X = x_i]$$

responsibility  
 that the  $y_i^{\text{th}}$  Gaussian  
 takes for the  
 feature  $x_i$

The surrogate function  $G(Q, \theta)$  is formulated as

$$G(Q, \theta) = F(Q, \theta) - \sum_{i=1}^m \sum_{y=1}^k Q_{i,y} \log Q_{i,y}$$

$$F(Q, \theta) = \sum_{i=1}^m \sum_{y=1}^k Q_{i,y} \log \left( P_{\theta} [X = x_i, Y = y_i] \right)$$

Posterior distr  
 over  $Y = y_i$

likelihood of  
 complete data:  
 observed + hidden  
 random variables

✓  $\therefore$  the function  $F(Q, \theta)$  gives

the expectation of the log likelihood of the complete data with respect to the posterior distribution over the latent (hidden) variable  $Y$ .