We are uncertain about the hidden variables. $Y$

We cannot observe them.

We capture this uncertainty by defining a posterior distribution over $Y = y_i$ given every training example $X = x_i$.

$P[Y=y]$

$P[Y = y_i] = $ prior ✓

$y_i$ is the sampled value

$P[Y = y_i \mid X = x_i]$ : posterior

$P[Y = y_i]$

$X = x_i$

$P[X = x_i \mid Y = y_i]$ : likelihood ✓

Generating an observation by sampling the likelihood function

The matrix $Q$ captures the posterior distribution

$$G(Q, \theta)$$

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

responsibility that the $y^{th}$ Gaussian takes for the feature $x_i$

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

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

for every row $i$ : Entropy of the cat distr

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

obs    latent

Maxim

Expectation . . . .

Posterior distr over $Y = y_i$ for $X = x_i$

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 $Y$.

## Given Array (Data)

$$a_1 \; a_2 \; a_3 \cdots a_d$$

## Distribution

$$p_1 \; p_2 \cdots p_d$$
$$\sum_i p_i = 1$$

## Expectation of the array w.r.t. the distribution

$$\sum_{i=1}^{d} p_i a_i$$

Mean of a data array $=$ Expect Uniform distribution w.r.t Uniform distr
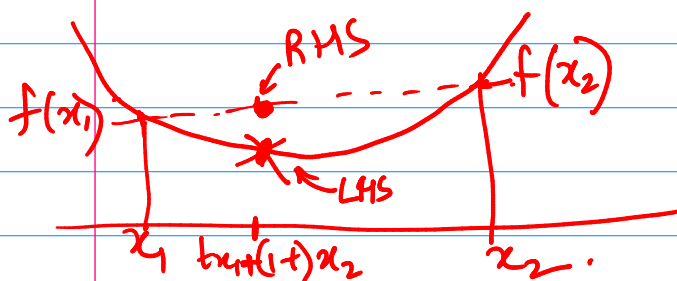
$$\sum_i \frac{1}{d} a_i = \frac{1}{d} \sum a_i$$

## Complete data log likelihood

$$x_i \downarrow \begin{array}{c} i \end{array} \left[ \begin{array}{cccc} \overset{y=1}{O} & \overset{y=2}{O} & \overset{y=3}{O} & O \\ & & & \\ \log\left( P_\theta\left[ X = x_i, Y = y_i \right] \right) & & & \end{array} \right] \quad m \times d$$

## Posterior over the latent variable $y$

$$x_i \downarrow \left[ \begin{array}{c} P[Y = y_i \mid X = x_i] \\ Q_{iy} \\ Q_i \end{array} \right] \quad m \times d$$

$$Q = \quad x_i \downarrow \begin{array}{c} x_1 \\ x_2 \\ x. \end{array} \left[ \begin{array}{ccccc} \overset{1}{\frac{1}{p_1}} & \overset{2}{\frac{1}{p_2}} & \overset{3}{\frac{1}{p_3}} & \overset{4}{=} & = \\ = & = & = & = & = \end{array} \right. \left. \begin{array}{c} \Sigma = 1. \\ \Sigma = 1. \end{array} \right.$$

## Entropy

$$\sum_i p_i \log \frac{1}{p_i}$$

## Jensen's Inequality

$$\underbrace{f\left( t\,x_1 + (1-t)\,x_2 \right)}_{LHS} \leq \underbrace{t\,f(x_1) + (1-t)\,f(x_2)}_{RHS}$$



If $f$ is a convex function

likelihood of $\underline{x}_i$.  $Y = y_i$  $\underline{x}_i \sim P(X|Y=y_i)$. $(\underline{x}_i \ y_i)$

$\underline{x}_i$ Posterior $P(Y=y|X=\underline{x}_i)$

The $G(Q, \underline{\theta})$ can be simplified as:

$$F(Q, \theta)$$

$$G(Q, \underline{\theta}) = \sum_{i=1}^{m} \sum_{y_i=1}^{k} Q_{i,y_i} \log P_{\underline{\theta}}[X=\underline{x}_i, Y=y_i] - $$

$$Q = \begin{bmatrix} \overset{y}{\overrightarrow{\circ \circ \circ \circ \circ \circ}} \\ i \\ \downarrow \end{bmatrix} Q_{i,y}$$

$Y=y_i \ (Y=y)$

$$\sum_{i=1}^{m} \sum_{y_i=1}^{k} Q_{i,y_i} \log Q_{i,y_i}$$

entropy posterior over latent variable

$\underline{x}_i$

$y_i \equiv y$

$P[y|X=\underline{x}_i]$

since it is implied that the random variable $Y$ takes the value $y_i$ in the context of example $x_i$ we write $y$ instead of $y_i$
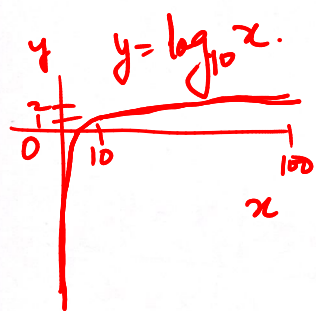
$$\log\frac{A}{B} = \log A - \log B$$

$$G(Q, \underline{\theta}) = \sum_{i=1}^{m} \sum_{y=1}^{k} Q_{i,y} \log \frac{P_{\underline{\theta}}[X=\underline{x}_i, Y=y]}{Q_{i,y}}$$

$Q_{i,y}$ : responsibility taken by the $y^{th}$ category for generating $\underline{x}_i$

✓ The upper bound of $G(Q, \underline{\theta})$ is $L(\underline{\theta})$.

$$G(Q, \theta) = \sum_{i=1}^{m} \sum_{y=1}^{k} Q_{i,y} \log \frac{P_{\underline{\theta}}[X=\underline{x}_i, Y=y_i]}{Q_{i,y}}$$

$$y = \log_{10} x.$$
$$0 \quad 10 \quad 100 \quad x$$

$$\leq \sum_{i=1}^{m} \log\left( \sum_{y=1}^{k} Q_{i,y} \frac{P_{\underline{\theta}}[X=\underline{x}_i, Y=y_i]}{Q_{i,y}} \right)$$

$f$

Jensen's inequality

$$= \sum_{i=1}^{m} \log\left( \sum_{y=1}^{k} P_{\underline{\theta}}[X=\underline{x}_i, Y=y_i] \right)$$

Summing out the latent vars

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

observed data log likelihood.

$$= L(\underline{\theta})$$

$$\therefore \quad G(Q, \underline{\theta}) \leq L(\underline{\theta}).$$

The maximization procedure involves alternating between

two steps

**E-step**

I $\quad Q^{(t+1)} = \underset{Q}{\arg\max} \; G(Q, \underline{\theta}^{(t)})$ $\quad$ Compute $Q^{(t+1)}$ with elements

closed form $\quad Q_{i,y}^{(t+1)} = P_{\underline{\theta}^{(t)}}[Y = y_i \mid X = x_i]$

**M-step**

II $\quad \underline{\theta}^{(t+1)} = \underset{\underline{\theta}}{\arg\max} \; G(Q^{(t+1)}, \underline{\theta})$ $\quad \underline{\theta}^{(t+1)} = \underset{\underline{\theta}}{\arg\max} \; F(Q^{t+1}, \underline{\theta})$

**EM algorithm**

$Y \longrightarrow X$

We know the property of $G$ function that $G(Q, \underline{\theta}) \leq L(\underline{\theta})$

that means $\quad G(Q, \underline{\theta}^{(t)}) \leq L(\underline{\theta}^{(t)})$

Now, if we substitute a $Q$ matrix with elements

$$Q_{i,y} = P_{\underline{\theta}}[Y = y_i \mid X = x_i] \quad \text{in } G(Q, \underline{\theta})$$

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

$$y_i \equiv y$$

$$= \sum_{i=1}^{m} \sum_{y_i=1}^{k} Q_{i,y} \log \frac{P_{\underline{\theta}}\{X = x_i, Y = y_i\}}{Q_{i,y}}$$

$$= \sum_{i=1}^{m} \sum_{y_i=1}^{k} P_{\underline{\theta}}[Y = y_i \mid X = x_i] \log \frac{P_{\underline{\theta}}[X = x_i, Y = y_i]}{P_{\underline{\theta}}[Y = y_i \mid X = x_i]}$$

$$= \sum_{i=1}^{m} \sum_{y_i=1}^{k} P_{\underline{\theta}}[Y = y_i \mid X = x_i] \log P_{\underline{\theta}}[X = x_i]$$

$$= \sum_{i=1}^{m} \log P_{\underline{\theta}}[X = x_i] \sum_{y_i=1}^{k} P_{\theta}[Y = y_i \mid X = x_i]$$

$$= \sum_{i=1}^{m} \log P_{\theta}[X = x_i] = L(\underline{\theta}).$$