

# Semi-Supervised Learning

- Consider the problem of Prepositional Phrase Attachment.
  - Buy car with money ; buy car with wheel
- There are several ways to generate features. Given the limited representation, we can assume that all possible conjunctions of the 4 attributes are used. (15 feature in each example).
- Assume we will use **naïve Bayes** for learning to decide between  $[n,v]$
- Examples are:  $(x_1, x_2, \dots, x_n, [n,v])$

# Using naïve Bayes

- To use naïve Bayes, we need to use the data to estimate:

$P(n)$	$P(v)$
$P(x_1   n)$	$P(x_1   v)$
$P(x_2   n)$	$P(x_2   v)$
.....	
$P(x_n   n)$	$P(x_n   v)$

- Then, given an example  $(x_1, x_2, \dots, x_n, ?)$ , compare:

$$P(n | x) \sim P(n) P(x_1 | n) P(x_2 | n) \dots P(x_n | n)$$

and

$$P(v | x) \sim P(v) P(x_1 | v) P(x_2 | v) \dots P(x_n | v)$$

# Using naïve Bayes

- After seeing 10 examples, we have:

- $P(n) = 0.5; P(v) = 0.5$

$$P(x_1 | n) = 0.75; P(x_2 | n) = 0.5; P(x_3 | n) = 0.5; P(x_4 | n) = 0.5$$

$$P(x_1 | v) = 0.25; P(x_2 | v) = 0.25; P(x_3 | v) = 0.75; P(x_4 | v) = 0.5$$

- Then, given an example  $x = (1000)$ , we have:

$$P_n(x) \sim 0.5 \cdot 0.75 \cdot 0.5 \cdot 0.5 \cdot 0.5 = 3/64$$

$$P_v(x) \sim 0.5 \cdot 0.25 \cdot 0.75 \cdot 0.25 \cdot 0.5 = 3/256$$

- Now, assume that in addition to the 10 labeled examples, we also have 100 unlabeled examples.
- Will that help?

# Using naïve Bayes

- For example, what can be done with the example (1000) ?
  - We have an estimate for its label...
  - But, can we use it to improve the classifier (that is, the estimation of the probabilities that we will use in the future)?
- Option 1: We can make predictions, and believe them
  - Or some of them (based on what?)
- Option 2: We can assume the example  $x=(1000)$  is a
  - An **n-labeled** example with probability  $P_n(x)/(P_n(x) + P_v(x))$
  - A **v-labeled** example with probability  $P_v(x)/(P_n(x) + P_v(x))$
- Estimation of probabilities does not require working with integers!

# Using Unlabeled Data

The discussion suggests several algorithms:

1. Use a threshold. Chose examples labeled with high confidence. Label them  $[n, v]$ . Retrain.
2. Use fractional examples. Label the examples with fractional labels  $[p \text{ of } n, (1-p) \text{ of } v]$ . Retrain.

# Comments on Unlabeled Data

- Both algorithms suggested can be used iteratively.
- Both algorithms can be used with other classifiers, not only naïve Bayes. The only requirement – a robust confidence measure in the classification.
- There are other approaches to Semi-Supervised learning: See included papers (co-training; Yarowsky's Decision List/Bootstrapping algorithm; “graph-based” algorithms that assume “similar” examples have “similar labels”, etc.)
- What happens if instead of 10 labeled examples we start with 0 labeled examples?
- Make a Guess; continue as above; a version of EM

# EM

- EM is a **class of algorithms** that is used to estimate a probability distribution in the presence of missing attributes.
- Using it, requires an assumption on the underlying probability distribution.
- The algorithm can be very sensitive to this assumption and to the starting point (that is, the initial guess of parameters).
- In general, known to converge to a local maximum of the maximum likelihood function.

# Three Coin Example

- We observe a series of coin tosses generated in the following way:
- A person has three coins.
  - Coin 0: probability of Head is  $\alpha$
  - Coin 1: probability of Head  $p$
  - Coin 2: probability of Head  $q$
- Consider the following coin-tossing scenarios:



# Estimation Problems

- Scenario I: Toss one of the coins four times.

Observing HHTH

Question: Which coin is more likely to produce this sequence ?

- Scenario II: Toss coin 0. If Head – toss coin 1; o/w – toss coin 2

Observing the sequence HHHHT, THTHT, HHHHT, HHTTH

produced by Coin 0 , Coin1 and Coin2

Question: Estimate most likely values for  $p$ ,  $q$  (the probability of H in each coin) and the probability to use each of the coins ( $\alpha$ )

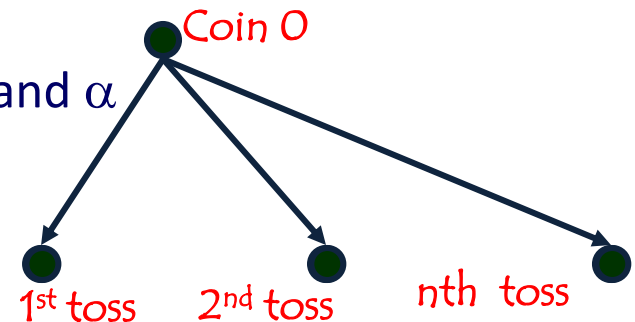
- Scenario III: Toss coin 0. If Head – toss coin 1; o/w – toss coin 2

Observing the sequence HHHT, HTHT, HHHT, HHTH

produced by Coin 1 and/or Coin 2

Question: Estimate most likely values for  $p$ ,  $q$  and  $\alpha$

There is no known analytical solution to this problem (general setting). That is, it is not known how to compute the values of the parameters so as to maximize the likelihood of the data.



# Key Intuition (1)

- If we knew which of the data points (HHHT), (HTHT), (HTTH) came from Coin1 and which from Coin2, there was no problem.
- Recall that the “simple” estimation is the **ML estimation**:
- Assume that you toss a  $(p, 1-p)$  coin  $m$  times and get  $k$  Heads  $m-k$  Tails.

$$\log[P(D|p)] = \log [ p^k (1-p)^{m-k} ] = k \log p + (m-k) \log (1-p)$$

- To maximize, set the derivative w.r.t.  $p$  equal to 0:

$$d \log P(D|p)/dp = k/p - (m-k)/(1-p) = 0$$

- Solving this for  $p$ , gives:  $p=k/m$

# Key Intuition (2)

- If we knew which of the data points (HHHT), (HTHT), (HTTH) came from Coin1 and which from Coin2, there was no problem.
- Instead, use an iterative approach for estimating the parameters:
  - Guess the probability that a given data point came from Coin 1 or 2; Generate fictional labels, weighted according to this probability.
  - Now, compute the most likely value of the parameters. [recall NB example]
  - Compute the likelihood of the data given this model.
  - Re-estimate the initial parameter setting: set them to maximize the likelihood of the data.  
(Labels  $\leftrightarrow$  Model Parameters)  $\leftrightarrow$  Likelihood of the data
- This process can be iterated and can be shown to converge to a local maximum of the likelihood function

# EM Algorithm (Coins) - I

- We will assume (for a minute) that we know the parameters  $\tilde{p}, \tilde{q}, \tilde{\alpha}$  and use it to estimate which Coin it is (Problem 1)
- Then, we will use this “label” estimation of the observed tosses, to estimate the **most likely** parameters
  - and so on...
- Notation:  $n$  data points; in each one:  $m$  tosses,  $h_i$  heads.
- What is the probability that the  $i$ th data point came from Coin1 ?
- **STEP 1 (Expectation Step):** (Here  $h=h_i$ )

$$\begin{aligned} P_1^i = P(\text{Coin1} | D^i) &= \frac{P(D^i | \text{Coin1}) P(\text{Coin1})}{P(D^i)} = \\ &= \frac{\tilde{\alpha} \tilde{p}^h (1 - \tilde{p})^{m-h}}{\tilde{\alpha} \tilde{p}^h (1 - \tilde{p})^{m-h} + (1 - \tilde{\alpha}) \tilde{q}^h (1 - \tilde{q})^{m-h}} \end{aligned}$$

# EM Algorithm

$$\begin{aligned}
 LL &= \sum_{i=1,n} \log \sum_{y=0,1} P(D^i, y \mid p, q, \alpha) = \\
 &= \sum_{i=1,n} \log \sum_{y=0,1} P(D^i \mid p, q, \alpha) \underbrace{P(y \mid D^i, p, q, \alpha)} = \\
 &= \sum_{i=1,n} \log E_{\mathbf{y}} P(D^i \mid p, q, \alpha) \geq \\
 &\geq \sum_{i=1,n} E_{\mathbf{y}} \log P(D^i \mid p, q, \alpha)
 \end{aligned}$$

Where the inequality is due to Jensen's Inequality.  
We maximize a lower bound on the Likelihood.

- Now, we would like to compute parameters that maximize it.
- We will maximize the log likelihood of (over  $n$  data points)
  - $LL = \sum_{i=1,n} \log P(D^i \mid p, q, \alpha)$
- But, one of the variables – the coin's name – is hidden. We can marginalize:
  - $LL = \sum_{i=1,n} \log \sum_{y=0,1} P(D^i, y \mid p, q, \alpha)$
- However, the sum is inside the log, making the solution difficult.
- Since the latent variable  $y$  is not observed, we cannot use the complete-data log likelihood. **Instead**, we use the expectation of the complete-data log likelihood under the posterior distribution of the latent variable to approximate  $\log p(D^i \mid p', q', \alpha')$
- We think of the likelihood  $\log P(D^i \mid p', q', \alpha')$  as a random variable that depends on the value  $y$  of the coin in the  $i^{\text{th}}$  toss. Therefore, instead of maximizing the LL we will maximize the expectation of this random variable (over the coin's name). [Justified using Jensen's Inequality; later & above]

# EM Algorithm

- We maximize the expectation (the coin name).

$$E[LL] = E\left[\sum_{i=1}^n \log P(D^i | p, q, \alpha)\right] \quad (1)$$

$$= \sum_{i=1}^n E[\log P(D^i | p, q, \alpha)] \quad (2)$$

$$= \sum_{i=1}^n p_1^i \log P(D^i | p, q, \alpha) + (1 - p_1^i) \log P(D^i | p, q, \alpha) \quad (3)$$

$$= \sum_{i=1}^n p_1^i \log \frac{P(D^i, 1 | p, q, \alpha)}{P(1 | D^i, p, q, \alpha)} + (1 - p_1^i) \log \frac{P(D^i, 0 | p, q, \alpha)}{P(0 | D^i, p, q, \alpha)} \quad (4)$$

$$= \sum_{i=1}^n p_1^i \log \frac{P(D^i, 1 | p, q, \alpha)}{p_1^i} + (1 - p_1^i) \log \frac{P(D^i, 0 | p, q, \alpha)}{1 - p_1^i} \quad (5)$$

$$= \sum_{i=1}^n p_1^i \log P(D^i, 1 | p, q, \alpha) - p_1^i \log p_1^i + (1 - p_1^i) \log P(D^i, 0 | p, q, \alpha) - (1 - p_1^i) \log (1 - p_1^i)$$

$$\begin{aligned} E[LL] &= E\left[\sum_{i=1, n} \log P(D^i | p, q, \alpha)\right] = \sum_{i=1, n} E[\log P(D^i | p, q, \alpha)] = \\ &= \sum_{i=1, n} P_1^i \log P(D^i, 1 | p, q, \alpha) + (1 - P_1^i) \log P(D^i, 0 | p, q, \alpha) \\ &\quad - \cancel{P_1^i \log P_1^i} - \cancel{(1 - P_1^i) \log (1 - P_1^i)} \end{aligned}$$

(Does not matter when we maximize)

- This is due to the linearity of the expectation and the random variable definition:

$$\begin{aligned} \log P(D^i, y | p, q, \alpha) &= \log P(D^i, 1 | p, q, \alpha) \quad \text{with Probability } P_1^i \\ &\quad \log P(D^i, 0 | p, q, \alpha) \quad \text{with Probability } (1 - P_1^i) \end{aligned}$$

# EM Algorithm (Coins) - IV

- Explicitly, we get:

$$\begin{aligned} E\left(\sum_i \log P(D^i | \tilde{p}, \tilde{q}, \tilde{\alpha})\right) &= \\ &= \sum_i P_1^i \log P(1, D^i | \tilde{p}, \tilde{q}, \tilde{\alpha}) + (1 - P_1^i) \log P(0, D^i | \tilde{p}, \tilde{q}, \tilde{\alpha}) = \\ &= \sum_i P_1^i \log(\tilde{\alpha} \tilde{p}^{h_i} (1 - \tilde{p})^{m - h_i}) + (1 - P_1^i) \log((1 - \tilde{\alpha}) \tilde{q}^{h_i} (1 - \tilde{q})^{m - h_i}) = \\ &= \sum_i P_1^i (\log \tilde{\alpha} + h_i \log \tilde{p} + (m - h_i) \log(1 - \tilde{p})) + \\ &\quad (1 - P_1^i) (\log(1 - \tilde{\alpha}) + h_i \log \tilde{q} + (m - h_i) \log(1 - \tilde{q})) \end{aligned}$$

# EM Algorithm (Cont.)

When computing the derivatives, notice  $P_1^i$  here is a constant; it was computed using the current parameters in the E step

- Finally, to find the most likely parameters, we maximize the derivatives with respect to  $\tilde{\mu}, \tilde{\alpha}$
- STEP 2: Maximization Step**
- (Sanity check: Think of the weighted fictional points)

$$\frac{dE}{d\tilde{\alpha}} = \sum_{i=1}^n \frac{P_1^i}{\tilde{\alpha}} - \frac{1-P_1^i}{1-\tilde{\alpha}} = 0 \quad \Rightarrow \quad \tilde{\alpha} = \frac{\sum P_1^i}{n}$$

$$\frac{dE}{d\tilde{\mu}} = \sum_{i=1}^n P_1^i \left( \frac{h_i}{\tilde{\mu}} - \frac{m-h_i}{1-\tilde{\mu}} \right) = 0 \quad \Rightarrow \quad \tilde{\mu} = \frac{\sum P_1^i \frac{h_i}{m}}{\sum P_1^i}$$

$$\frac{dE}{d\tilde{q}} = \sum_{i=1}^n (1-P_1^i) \left( \frac{h_i}{\tilde{q}} - \frac{m-h_i}{1-\tilde{q}} \right) = 0 \quad \Rightarrow \quad \tilde{q} = \frac{\sum (1-P_1^i) \frac{h_i}{m}}{\sum (1-P_1^i)}$$



# Models with Hidden Variables

- Now say we have two sets  $\mathcal{X}$  and  $\mathcal{Y}$ , and a joint distribution  $P(X, Y | \Theta)$

- If we had **fully observed data**,  $(X_i, Y_i)$  pairs, then

$$L(\Theta) = \sum_i \log P(X_i, Y_i | \Theta)$$

- If we have **partially observed data**,  $X_i$  examples, then

$$\begin{aligned} L(\Theta) &= \sum_i \log P(X_i | \Theta) \\ &= \sum_i \log \sum_{Y \in \mathcal{Y}} P(X_i, Y | \Theta) \end{aligned}$$

- The **EM (Expectation Maximization) algorithm** is a method for finding

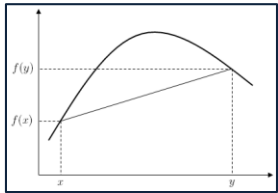
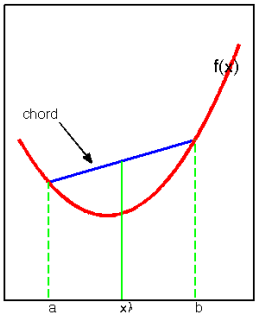
$$\Theta_{ML} = \operatorname{argmax}_{\Theta} \sum_i \log \sum_{Y \in \mathcal{Y}} P(X_i, Y | \Theta)$$

# EM: General Setting

- The EM algorithm is a general purpose algorithm for finding the maximum likelihood estimate in latent variable models.
- In the E-Step, we “fill in” the latent variables using the posterior, and in the M-Step, we maximize the expected complete log likelihood with respect to the complete posterior distribution.
  - Let  $D = (x_1, \dots, x_N)$  be the observed data, and
  - Let  $Z$  denote hidden random variables.
    - (We are not committing to any particular model.)
  - Let  $\theta$  be the model parameters. Then
- $\theta^* = \operatorname{argmax}_{\theta} p(x | \theta) = \operatorname{argmax}_{\theta} \sum_z p(x, z | \theta) =$
- $= \operatorname{argmax}_{\theta} \sum_z [p(z | \theta)p(x | z, \theta)]$
- This expression is called the complete log likelihood.

# EM: General Setting

Jensen's Inequality for convex functions:  
 $E(f(x)) \geq f(E(x))$   
But log is concave, so  
 $E(\log(x)) \leq \log(E(x))$



- To derive the **EM objective function**, we re-write the complete log likelihood function by multiplying it by  $q(z)/q(z)$ , where  $q(z)$  is an arbitrary distribution for the random variable  $z$ .
- $$\begin{aligned} \log p(x|\theta) &= \log \sum_z p(x, z | \theta) = \log \sum_z p(z|\theta) p(x|z, \theta) q(z)/q(z) = \\ &= \log E_q [p(z|\theta) p(x|z, \theta) / q(z)] \geq \\ &\geq E_q \log [p(z|\theta) p(x|z, \theta) / q(z)], \end{aligned}$$
- Where the inequality is due to Jensen's inequality applied to the *concave* function, log.
- We get the objective:  
$$L(\theta, q) = E_q [\log p(z|\theta)] + E_q [\log p(x|z, \theta)] - E_q [\log q(z)]$$
- The last component is an Entropy component; it is also possible to write the objective so that it includes a KL divergence (a distance function between distributions) of  $q(z)$  and  $p(z|x, \theta)$ .

Other  $q$ 's can be chosen [Samdani & Roth2012] to give other EM algorithms. Specifically, you can choose a  $q$  that chooses the most likely  $z$  in the E-step, and then continues to estimate the parameters (called Truncated EM, or Hard EM).

(Think back to the semi-supervised case)

- EM now continues iteratively, as a gradient ascent algorithm, where we choose  $q = p(z|x, \theta)$ .
- At the  $t$ -th step, we have  $q^{(t)}$  and  $\theta^{(t)}$ .
- E-Step:** update the posterior  $q$ , while holding  $\theta^{(t)}$  fixed:

$$q^{(t+1)} = \operatorname{argmax}_q L(q, \theta^{(t)}) = p(z|x, \theta^{(t)}).$$

- M-Step:** update the model parameters to maximize the expected complete log-likelihood function:

$$\theta^{(t+1)} = \operatorname{argmax}_\theta L(q^{(t+1)}, \theta)$$

To wrap it up, with the right  $q$ :

- $$L(\theta, q) = E_q \log [p(z|\theta) p(x|z, \theta) / q(z)]$$

$$= \sum_z p(z|x, \theta) \log [p(x, z|\theta) / p(z|x, \theta)] =$$

$$= \sum_z p(z|x, \theta) \log [p(x, z|\theta) p(x|\theta) / p(z, x|\theta)] =$$

$$= \sum_z p(z|x, \theta) \log [p(x|\theta)] = \log [p(x|\theta)] \sum_z p(z|x, \theta) = \log [p(x|\theta)]$$

- So, by maximizing the objective function, we are also maximizing the log likelihood function.

# The General EM Procedure

- Initially, the parameter  $\theta$  is set as  $\theta_0$
- In E step
  - We use the current parameter values  $\theta^{\text{old}}$  to find the posterior distribution of the latent variables given by  $p(Z|X, \theta^{\text{old}})$
  - Use  $p(Z|X, \theta^{\text{old}})$  to compute the expectation of the complete-data log likelihood  $\ln p(X, Z|\theta)$  under  $p(Z|X, \theta^{\text{old}})$

$$Q(\theta, \theta^{\text{old}}) = \sum_Z p(Z|X, \theta^{\text{old}}) \ln p(X, Z|\theta) \quad \text{E}$$

- In M step, we need to compute  $\theta^{\text{new}}$  which maximizes  $Q(\theta, \theta^{\text{old}})$

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}) \quad \text{M}$$

# EM Summary (so far)

- EM is a general procedure for learning in the presence of unobserved variables.
- We have shown how to use it in order to estimate the most likely density function for a mixture of (Bernoulli) distributions.
- EM is an iterative algorithm that can be shown to converge to a local maximum of the likelihood function.
- It depends on assuming a family of probability distributions.
- In this sense, it is a family of algorithms. The update rules you will derive depend on the model assumed.
- It has been shown to be quite useful in practice, when the assumptions made on the probability distribution are correct, but can fail otherwise.

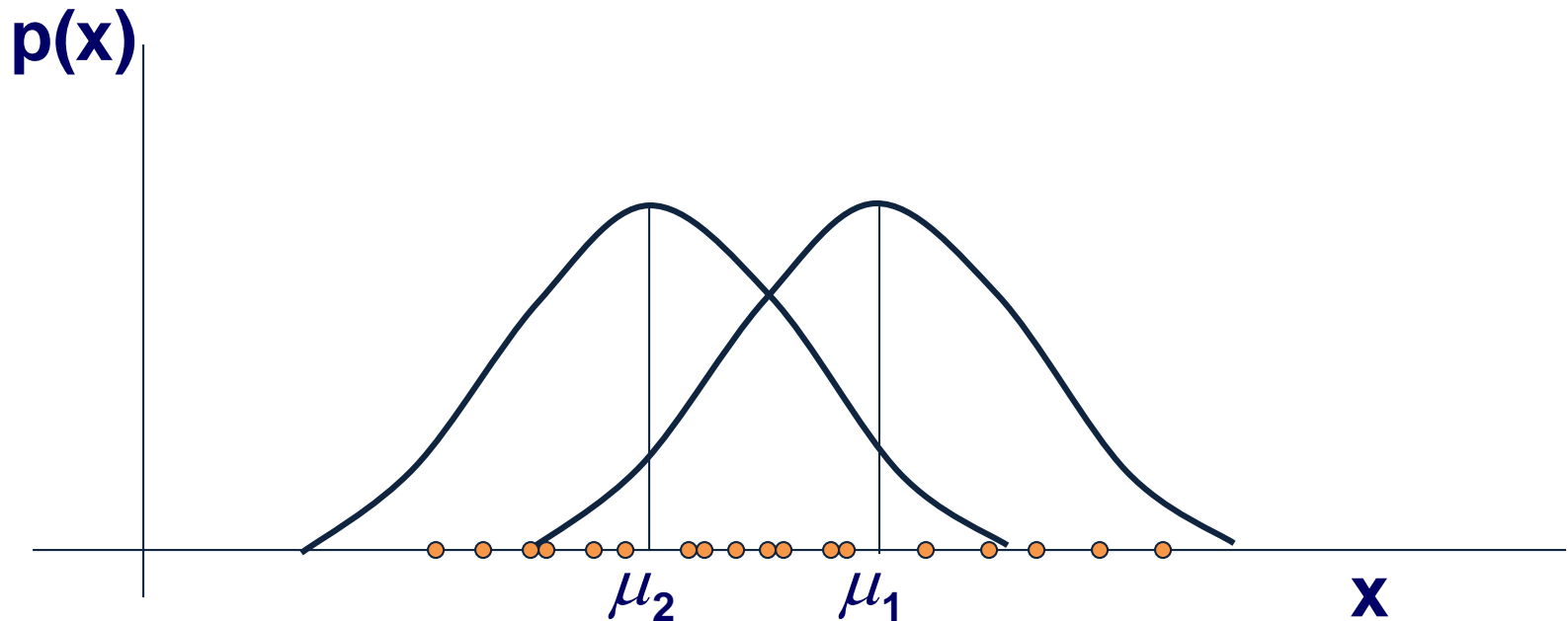
# EM Summary (so far)

- EM is a general procedure for learning in the presence of unobserved variables.
- The (family of ) probability distribution is known; the problem is to estimate its parameters
- In the presence of hidden variables, we can often think about it as a problem of a mixture of distributions – the participating distributions are known, we need to estimate:
  - Parameters of the distributions
  - The mixture policy
- Our previous example: Mixture of Bernoulli distributions

# Example: K-Means Algorithm

K-means is a **clustering** algorithm.

We are given data points, known to be sampled independently from a mixture of  $k$  Normal distributions, with means  $\mu_i, i=1,\dots,k$  and the same standard variation  $\sigma$





# Example: K-Means Algorithm

First, notice that if we knew that all the data points are taken from a normal distribution with mean  $\mu$ , finding its most likely value is easy.

$$p(\mathbf{x} \mid \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2} (\mathbf{x} - \mu)^2\right]$$

We get many data points,  $D = \{x_1, \dots, x_m\}$

$$\ln(L(D \mid \mu)) = \ln(P(D \mid \mu)) = \sum_i -\frac{1}{2\sigma^2} (\mathbf{x}_i - \mu)^2$$

Maximizing the log-likelihood is equivalent to minimizing:

$$\mu_{\text{ML}} = \operatorname{argmin}_{\mu} \sum_i (\mathbf{x}_i - \mu)^2$$

Calculate the derivative with respect to  $\mu$ , we get that the minimal point, that is, the most likely mean is  $\mu = \frac{1}{m} \sum_i \mathbf{x}_i$

# A mixture of Distributions

As in the coin example, the problem is that data is sampled from a mixture of  $k$  different normal distributions, and we do not know, for a given data point  $x_i$ , where is it sampled from.

Assume that we observe data point  $x_i$ ; what is the probability that it was sampled from the distribution  $\mu_j$  ?

$$\begin{aligned} P_{ij} = P(\mu_j | x_i) &= \frac{P(x_i | \mu_j)P(\mu_j)}{P(x_i)} = \frac{\frac{1}{k} P(x = x_i | \mu = \mu_j)}{\sum_{n=1}^k \frac{1}{k} P(x = x_i | \mu = \mu_n)} = \\ &= \frac{\exp\left[-\frac{1}{2\sigma^2} (x_i - \mu_j)^2\right]}{\sum_{n=1}^k \exp\left[-\frac{1}{2\sigma^2} (x_i - \mu_n)^2\right]} \end{aligned}$$

# A Mixture of Distributions

As in the coin example, the problem is that data is sampled from a mixture of  $k$  different normal distributions, and we do not know, for a given each data point  $x_i$ , where is it sampled from.

For a data point  $x_i$ , define  $k$  binary hidden variables,  $z_{i1}, z_{i2}, \dots, z_{ik}$ , s.t.  $z_{ij} = 1$  iff  $x_i$  is sampled from the  $j$ -th distribution.

$$\mathbf{E}[z_{ij}] = 1 \bullet \mathbf{P}(x_i \text{ was sampled from } \mu_j) +$$

$$0 \bullet \mathbf{P}(x_i \text{ was not sampled from } \mu_j) = \mathbf{P}_{ij}$$

$$\mathbf{E}[Y] = \sum_{y_i} y_i \mathbf{P}(Y = y_i)$$

$$\mathbf{E}[X + Y] = \mathbf{E}[X] + \mathbf{E}[Y]$$

# Example: K-Means Algorithms

Expectation: (here:  $h = \sigma, \mu_1, \mu_2, \dots, \mu_k$ )

$$p(\mathbf{y}_i | \mathbf{h}) = p(\mathbf{x}_i, \mathbf{z}_{i1}, \dots, \mathbf{z}_{ik} | \mathbf{h}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2} \sum_j \mathbf{z}_{ij} (\mathbf{x}_i - \mu_j)^2\right]$$

Computing the likelihood given the observed data  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  and the hypothesis  $h$  (w/o the constant coefficient)

$$\ln(P(Y | h)) = \sum_{i=1}^m -\frac{1}{2\sigma^2} \sum_j \mathbf{z}_{ij} (\mathbf{x}_i - \mu_j)^2$$

$$\begin{aligned} \mathbf{E}[\ln(P(Y | h))] &= \mathbf{E}\left[\sum_{i=1}^m -\frac{1}{2\sigma^2} \sum_j \mathbf{z}_{ij} (\mathbf{x}_i - \mu_j)^2\right] = \\ &= \sum_{i=1}^m -\frac{1}{2\sigma^2} \sum_j \mathbf{E}[\mathbf{z}_{ij}] (\mathbf{x}_i - \mu_j)^2 \end{aligned}$$

# Example: K-Means Algorithms

Maximization: Maximizing

$$Q(\mathbf{h} | \mathbf{h}') = \sum_{i=1}^m -\frac{1}{2\sigma^2} \sum_j \mathbf{E}[\mathbf{z}_{ij}] (\mathbf{x}_i - \mu_j)^2$$

with respect to  $\mu_j$  we get that:

$$\frac{dQ}{d\mu_j} = \mathbf{C} \sum_{i=1}^m \mathbf{E}[\mathbf{z}_{ij}] (\mathbf{x}_i - \mu_j) = \mathbf{0}$$

Which yields:

$$\mu_j = \frac{\sum_{i=1}^m \mathbf{E}[\mathbf{z}_{ij}] \mathbf{x}_i}{\sum_{i=1}^m \mathbf{E}[\mathbf{z}_{ij}]}$$

# Summary: K-Means Algorithms

Given a set  $D = \{x_1, \dots, x_m\}$  of data points,

guess initial parameters  $\sigma, \mu_1, \mu_2, \dots, \mu_k$

Compute (for all  $i, j$ )

$$p_{ij} = \mathbf{E}[z_{ij}] = \frac{\exp\left[-\frac{1}{2\sigma^2}(\mathbf{x}_i - \mu_j)^2\right]}{\sum_{n=1}^k \exp\left[-\frac{1}{2\sigma^2}(\mathbf{x}_i - \mu_n)^2\right]}$$

and a new set of means:

$$\mu_j = \frac{\sum_{i=1}^m \mathbf{E}[z_{ij}] \mathbf{x}_i}{\sum_{i=1}^m \mathbf{E}[z_{ij}]}$$

repeat to convergence

**Notice that this algorithm will find the best  $k$  means in the sense of minimizing the sum of square distance.**

# Summary: EM

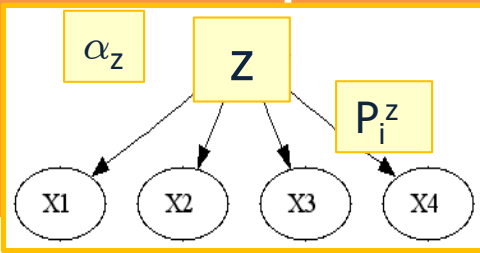
- EM is a general procedure for learning in the presence of unobserved variables.
- We have shown how to use it in order to estimate the most likely density function for a mixture of probability distributions.
- EM is an iterative algorithm that can be shown to converge to a local maximum of the likelihood function. Thus, might requires many restarts.
- It depends on assuming a family of probability distributions.
- It has been shown to be quite useful in practice, when the assumptions made on the probability distribution are correct, but can fail otherwise.
- As examples, we have derived an important clustering algorithm, the k-means algorithm and have shown how to use it in order to estimate the most likely density function for a mixture of probability distributions.

# More Thoughts about EM

- **Training:** a sample of data points,  $(x_0, x_1, \dots, x_n) \in \{0,1\}^{n+1}$
- **Task:** predict the value of  $x_0$ , given assignments to all  $n$  variables.

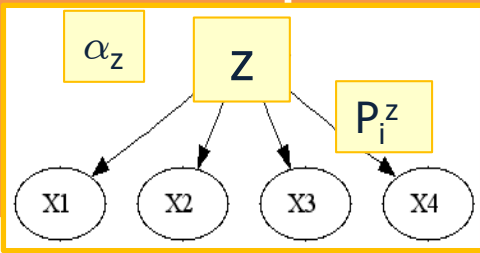


# More Thoughts about EM



- Assume that a set  $x^i \in \{0,1\}^{n+1}$  of data points is generated as follows:
- Postulate a hidden variable  $Z$ , with  $k$  values,  $1 \leq z \leq k$  with probability  $\alpha_z$ ,  $\sum_{1,k} \alpha_z = 1$
- Having randomly chosen a value  $z$  for the hidden variable, we choose the value  $x_i$  for each observable  $X_i$  to be 1 with probability  $p_i^z$  and 0 otherwise,  $[i = 0, 1, 2, \dots, n]$
- Training:** a sample of data points,  $(x_0, x_1, \dots, x_n) \in \{0,1\}^{n+1}$
- Task:** predict the value of  $x_0$ , given assignments to all  $n$  variables.

# More Thoughts about EM



- Two options:
- Parametric:** estimate the model using EM. Once a model is known, use it to make predictions.
  - Problem: Cannot use EM directly without an additional assumption on the way data is generated.
- Non-Parametric:** Learn  $X_0$  directly as a function of the other variables.
  - Problem: which function to try and learn?
- $X_0$  turns out to be a linear function of the other variables, when  $k=2$  (what does this mean?)

- Another important distinction to attend to is the fact that, once you estimated all the parameters with EM, you can answer many prediction problems e.g.,  $p(X_0, X_7, \dots, X_8 | X_1, X_2, \dots, X_n)$  while with Perceptron (say) you need to learn separate models for each prediction problem.

# EM

- Iterative procedure is defined as  $\Theta^t = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{t-1})$ , where

$$Q(\Theta, \Theta^{t-1}) = \sum_i \sum_{Y \in \mathcal{Y}} P(Y | X_i, \Theta^{t-1}) \log P(X_i, Y | \Theta)$$

- Key points:
  - Intuition: fill in hidden variables  $Y$  according to  $P(Y | X_i, \Theta)$
  - EM is guaranteed to converge to a local maximum, or saddle-point, of the likelihood function
  - In general, if

$$\operatorname{argmax}_{\Theta} \sum_i \log P(X_i, Y_i | \Theta)$$

has a simple (analytic) solution, then

$$\operatorname{argmax}_{\Theta} \sum_i \sum_Y P(Y | X_i, \Theta) \log P(X_i, Y | \Theta)$$

also has a simple (analytic) solution.

# The EM Algorithm

## Algorithm:

- Guess initial values for the hypothesis  $h = \sigma, \mu_1, \mu_2, \dots, \mu_k$
- Expectation: Calculate  $Q(h', h) = E(\text{Log } P(Y | h') | h, X)$   
using the current hypothesis  $h$  and the observed data  $X$ .
- Maximization: Replace the current hypothesis  $h$  by  $h'$ , that maximizes the Q function (the likelihood function)  
set  $h = h'$ , such that  $Q(h', h)$  is maximal
- Repeat: Estimate the Expectation again.