**Alin Tomescu**, http://people.csail.mit.edu/~alinush
6.867 Machine learning | Prof. Tommi Jaakkola | Week 9, Thursday, October 31st, 2013| Lecture 16
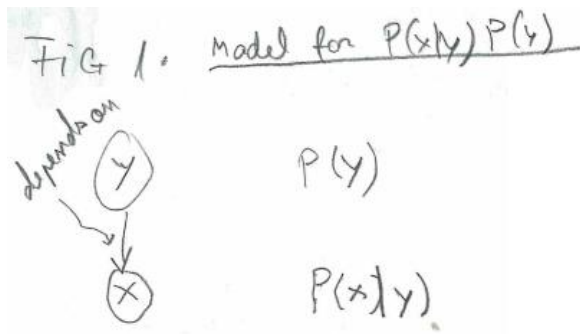
# Lecture 16

**Project grading criteria:** How you describe what you achieved, show what you understand (practical / theoretical aspect of the material). It does not matter if you succeeded in your classification task

**2 topics that will be on final:** SVM with kernels, EM algorithm.

$$P(x, y; \theta) = \sum_{y=1}^{k} P(x \mid y; \theta) \, P(y; \theta)$$
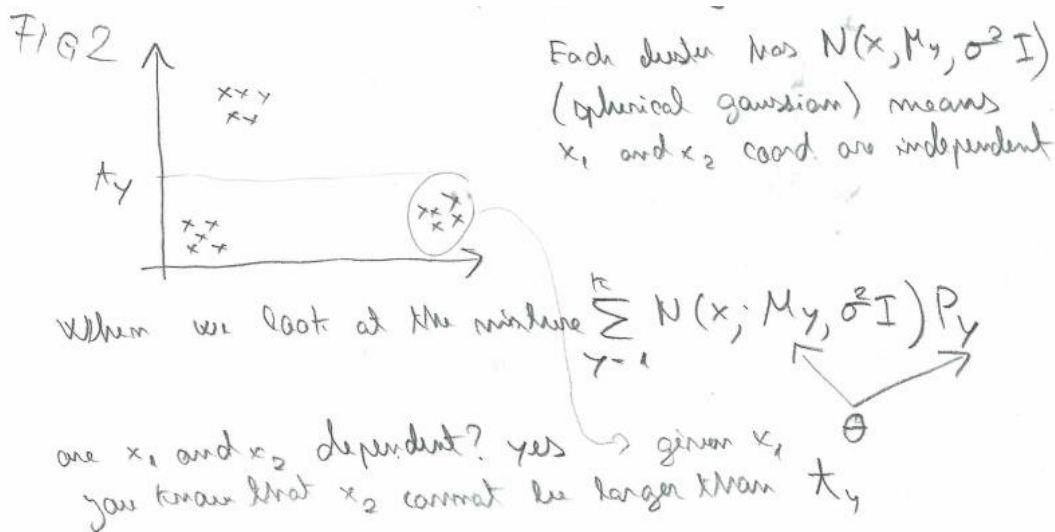
A simple graphical model for this will look like:



This is a "compact description for generating the values" and this is how you sample/generate values:

- $y$ does not depend on anything, so just sample
- $x$ depends on $y$, so you can only sample it after you sampled an $y$

## The EM algorithm

**First step:** Initialize. We will see this is quite important, since it defines a starting point from which we start refining.

$$\theta^{[0]}$$

**E-step:** Performs the *data completion.* In each step, we map each data point $x^{(i)}$ to $k$ different data points $(x^{(i)}, y)$ but not with equal weighting (see previous $q$ function).

$$q^{[m]}(y \mid i) = P(y \mid x^{(i)}; \theta^{[m]}) = \frac{P(x^{(i)} \mid y; \theta^{[m]}) P_y^{[m]}}{P(x^{(i)}; \theta^{[m]})} = \frac{P(x^{(i)} \mid y; \theta^{[m]}) P_y^{[m]}}{\sum_{y'=1}^{k} P(x^{(i)} \mid y'; \theta^{[m]}) P_{y'}^{[m]}}$$

We do this for all $y = 1, \dots, k$ and for the entire data set $i = 1, \dots, n$.

**Note:** Now we kind of have $nk$ points after the $x^{(i)} \to (x^{(i)}, y)$ mapping.

**M-step:** Simply takes that data set and maximizes the log-likelihood of it.

$$\theta^{[m+1]} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^{n} \sum_{y=1}^{k} q^{[m]}(y \mid i) \log \left[ N\left( x^{(i)}; \mu_y^{[m+1]}, \sigma^2 I \right) P_y^{[m+1]} \right], \theta = \left[ \mu_y^{[m+1]}, P_y^{[m+1]} \right]$$
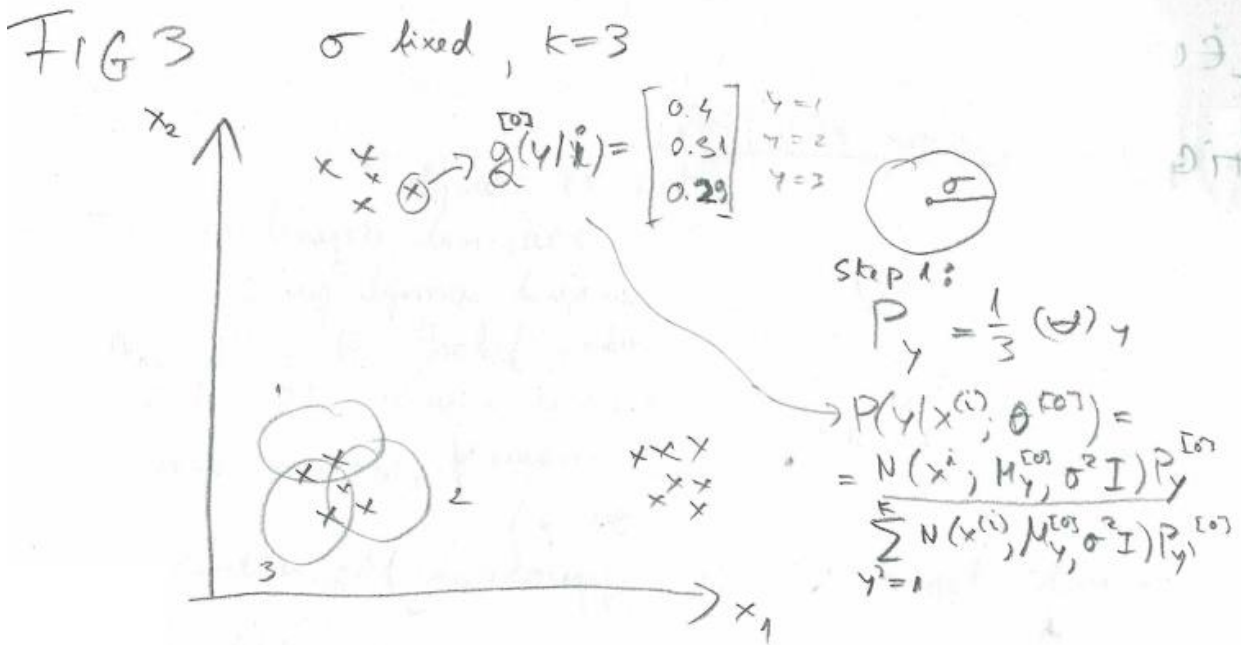
If we compute the argmax, we obtain:

$$\hat{p}_y^{[m+1]} = \frac{\sum_{i=1}^{n} q^{[m]}(y \mid i)}{n}, i = 1 \dots k$$

$$\hat{\mu}_y^{[m+1]} = \frac{1}{\sum_{i=1}^{n} q^{[m]}(y \mid i)} \sum_{i=1}^{n} q^{[m]}(y \mid i) \cdot x^{(i)}, i = 1 \dots k$$

**Note:** $\sum_{y=1}^{k} q(y \mid i) = 1$

## Analyzing the EM algorithm

Let's see how starting from a poor model might lead to improvement and how important the initialization step is.

**Initialization step:** Pick probabilities for each $y$, pick mean for the Gaussians for each $y$.
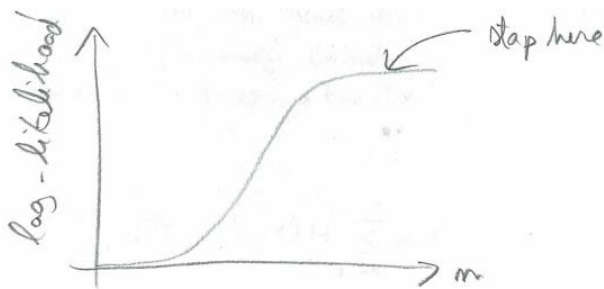
$$P_y = \frac{1}{3}, \forall y$$

$$\mu_y = \cdots$$

To pick the means we could just pick random points for each cluster.

**E-step**: compute $q^{[m]}(y \mid i) = P(y \mid x^{(i)}; \theta^{[m]})$

$$q^{[0]}(y \mid i) = P(y \mid x^{(i)}; \theta^{[0]})$$

**M-step**: In the above example, the Gaussian $y = 1$ only looks at $q(1, i)$ and goes towards the points that have higher $q(1, i)$ values. The other Gaussians proceed similarly.

**Question:** When do you stop? **Answer:** It generally takes an infinite amount of steps to converge to the maximum likelihood. The likelihood grows slowly initially and then it starts increasing as the Gaussians "see" the points and then it starts growing slowly again.
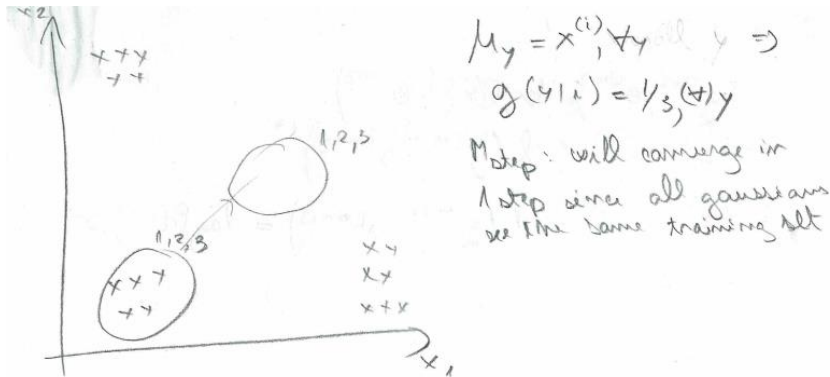
**Note:** You will always converge in the EM algorithm, but not necessarily to the best solution. Only guaranteed to converge to a locally optimal solution.

**Question:** What will happen with different initialization, why is it important?
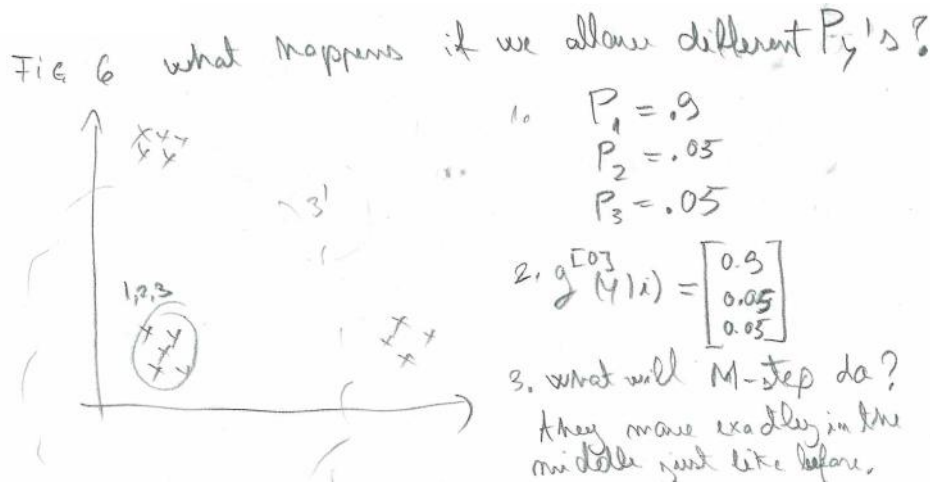
## Case 1: $\mu_y = x^{(i)}, \forall y$ and $P_y = p, \forall y$ and fixed $\sigma$

In this case the EM algorithm will converge in one step to a non-optimal solution.



## Case 2: $\mu_y = x^{(i)}, \forall y$ and $P_{y^{(i)}} \neq P_{y^{(j)}}$ and fixed $\sigma$

In this case, the same solution is computed.



## Case 3: $\mu_y = x^{(i)}, \forall y$ and $P_{y^{(i)}} \neq P_{y^{(j)}}$ and different $\sigma$

**Note:** These are not the only ways to initialize badly. Plenty of other bad ways…

If I allow different $\sigma$'s, the Gaussian with the higher $\sigma$ will "see" points further away "better" and will move towards these points more than the other Gaussians

# EM algorithm as a maximization problem

We will show that the E-step and M-step both maximize for the log-likelihood.

What is this lower bound that I am constantly maximizing and how it relates to the log-likelihood?

If we can maximize $P(x^{(i)}; \theta)$, we can also maximize over the sum:

$$P(x^{(i)}; \theta) = \sum_{y=1}^{k} P(x^{(i)}, y; \theta)$$

$$P(x^{(i)}, y; \theta) = P(x^{(i)} \mid y; \theta)P(y; \theta) = P(y \mid x^{(i)}; \theta)P(x^{(i)}; \theta)$$

$$l(q, \theta) = \sum_{y=1}^{k} q(y \mid i) \log P(x^{(i)}, y; \theta) + \sum_{y} q(y \mid i) \log \frac{1}{q(y \mid i)}$$

We added this (Shanon) entropy term (measure of uncertainty) $\sum_y q(y \mid i) \log \frac{1}{q(y \mid i)}$ as a penalty.

**Claim 1:** For any parameters, $l(q, \theta) \leq \log P(x^{(i)}; \theta)$

**Claim 2:** If I start from parameters $\theta^{[m]}$, $l(q^{[m]}, \theta^{[m]}) = \log P(x^{(i)}; \theta^{[m]})$, when $q^{[m]}(y \mid i) = P(y \mid x^{(i)}; \theta^{[m]})$

$$l(q^{[m]}, \theta^{[m]}) = \sum_{y=1}^{k} q^{[m]}(y \mid i) \left( \frac{\log P(x^{(i)}, y; \theta^{[m]})}{q^{[m]}(y \mid i)} \right) = \sum_{y=1}^{k} q^{[m]}(y \mid i) \left( \frac{\log P(x^{(i)}; \theta^{[m]})P(y \mid x^{(i)}; \theta^{[m]})}{q^{[m]}(y \mid i)} \right)$$

$$= \sum_{y=1}^{k} q^{[m]}(y \mid i) \left( \frac{\log P(x^{(i)}; \theta^{[m]})q^{[m]}(y \mid i)}{q^{[m]}(y \mid i)} \right) = \sum_{y=1}^{k} q^{[m]}(y \mid i) \log P(x^{(i)}; \theta^{[m]}) = \log P(x^{(i)}; \theta^{[m]})$$

FIG 7:

$$\ell\left(g^{[m]}, \theta^{[m]}\right) = \log P(x^{(i)}; \theta^{[m]})$$

M-step $\longrightarrow \leq \ell\left(g^{[m]}, \theta^{[m+1]}\right)$

E-step $\longrightarrow \leq \ell\left(g^{[m+1]}, \theta^{[m+1]}\right) = \log P(x^{(i)}; \theta^{[m+1]})$

$$l(q^{[m]}, \theta^{[m]}) = \log P(x^{(i)}; \theta^{[m]}) \leq l(q^{[m]}, \theta^{[m+1]}) \leq l(q^{[m+1]}, \theta^{[m+1]}) = \log P(x^{(i)}; \theta^{[m+1]})$$