

Mathematics of Big Data, I

Lecture 7: Unsupervised Learning, K-mean Clustering, Gaussian Mixture, Jensen's inequality, and EM

Weiqing Gu

Professor of Mathematics

Director of the Mathematics Clinic

Harvey Mudd College

Summer 2017

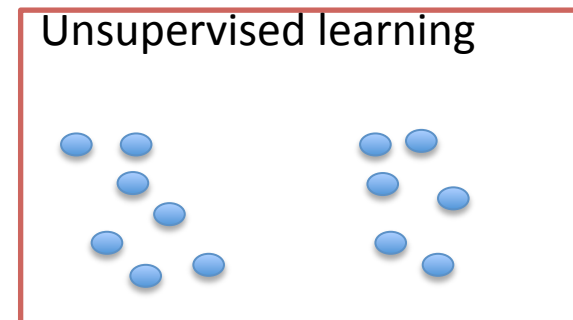
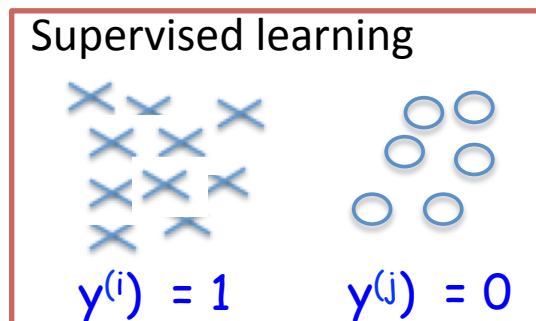
Today's topics

- The k-means clustering algorithm
- Mixtures of Gaussians
- Jensen's inequality
- The EM (Expectation-Maximization) Algorithm

What is a clustering problem?

A clustering problem is an unsupervised learning problem

- Given a training set $\{x^{(1)}, \dots, x^{(m)}\}$, here each $x^{(i)}$ is in \mathbf{R}^n .
- Goal: want to group the data into a few cohesive “clusters.”
- Note: the difference between unsupervised learning and supervised learning is that no labels $y^{(i)}$ are given.



In general, if only $\{x^{(1)}, \dots, x^{(m)}\}$ given for a problem, but no labels $y^{(i)}$ are given, then the problem is an unsupervised learning problem!

The k-means clustering algorithm

1. Initialize cluster centroids $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$ randomly.

2. Repeat until convergence: {

For every i , set

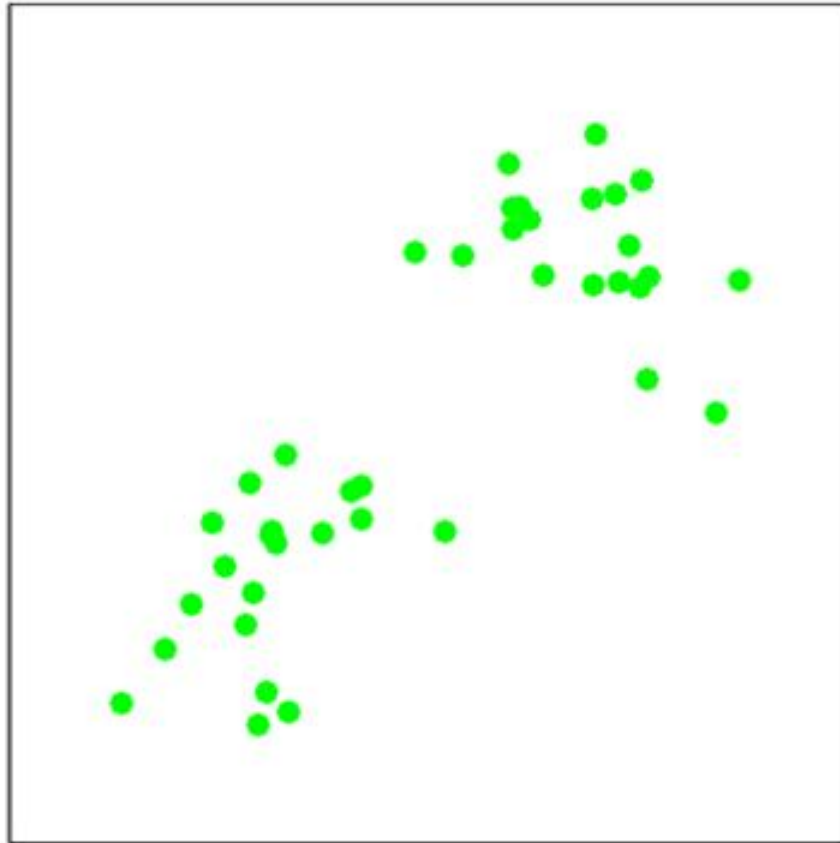
$$c^{(i)} := \arg \min_j ||x^{(i)} - \mu_j||^2.$$

For each j , set

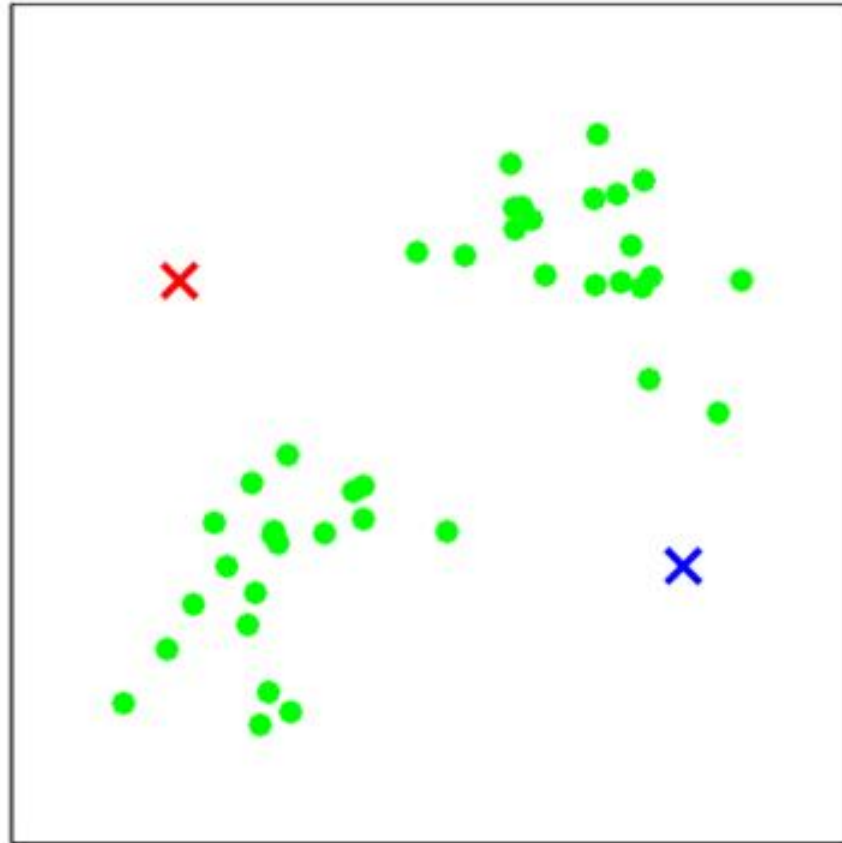
$$\mu_j := \frac{\sum_{i=1}^m 1\{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m 1\{c^{(i)} = j\}}.$$

}

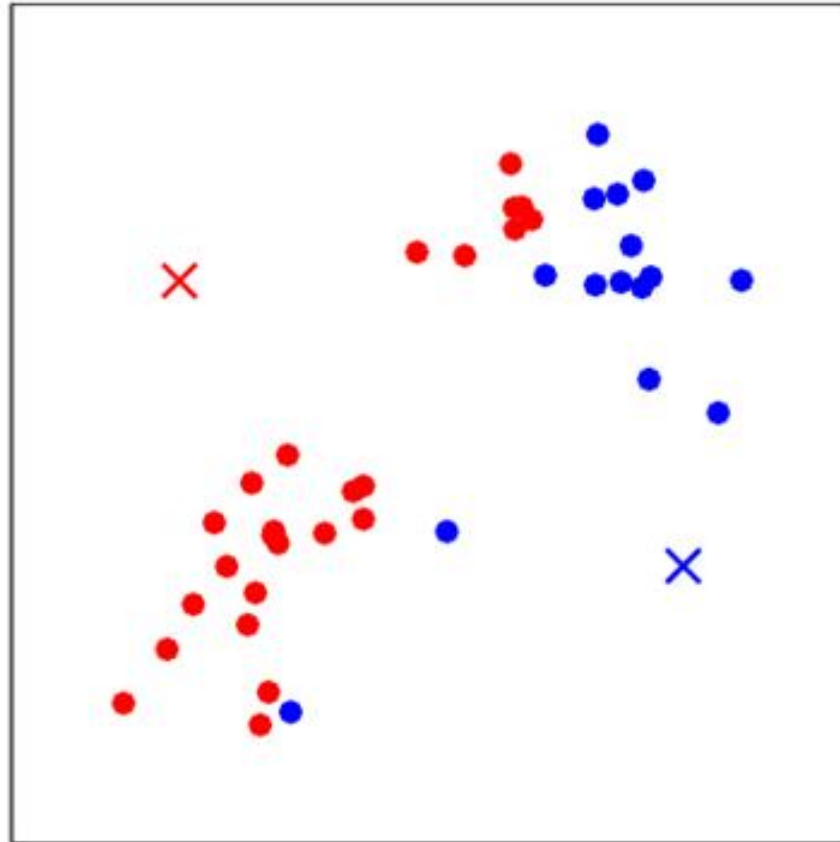
Example of K-mean clustering



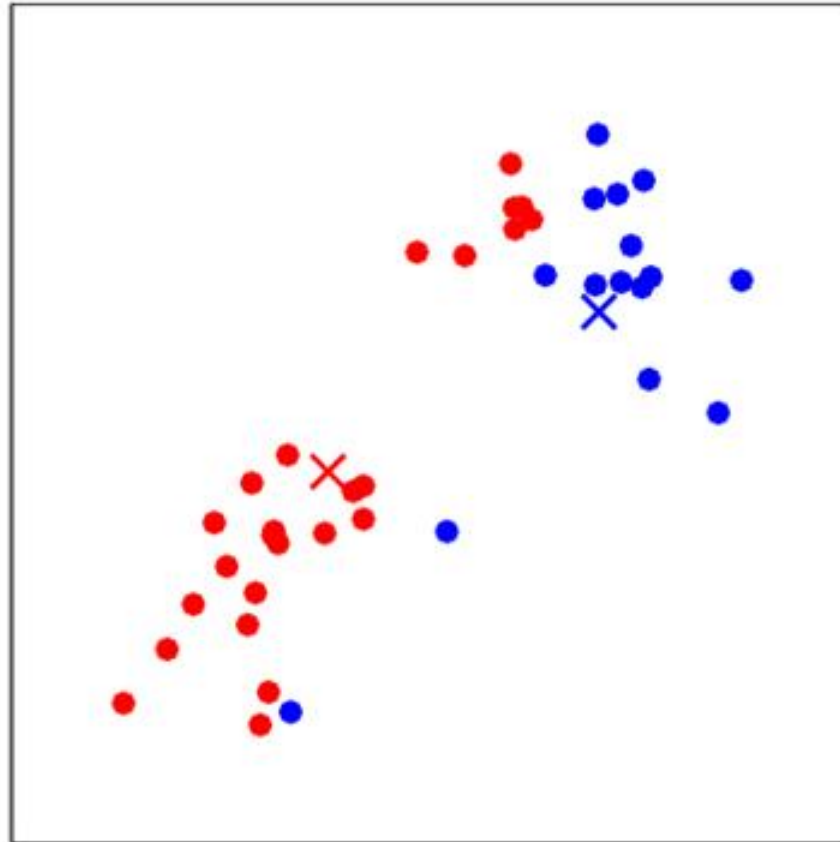
Example of K-mean clustering



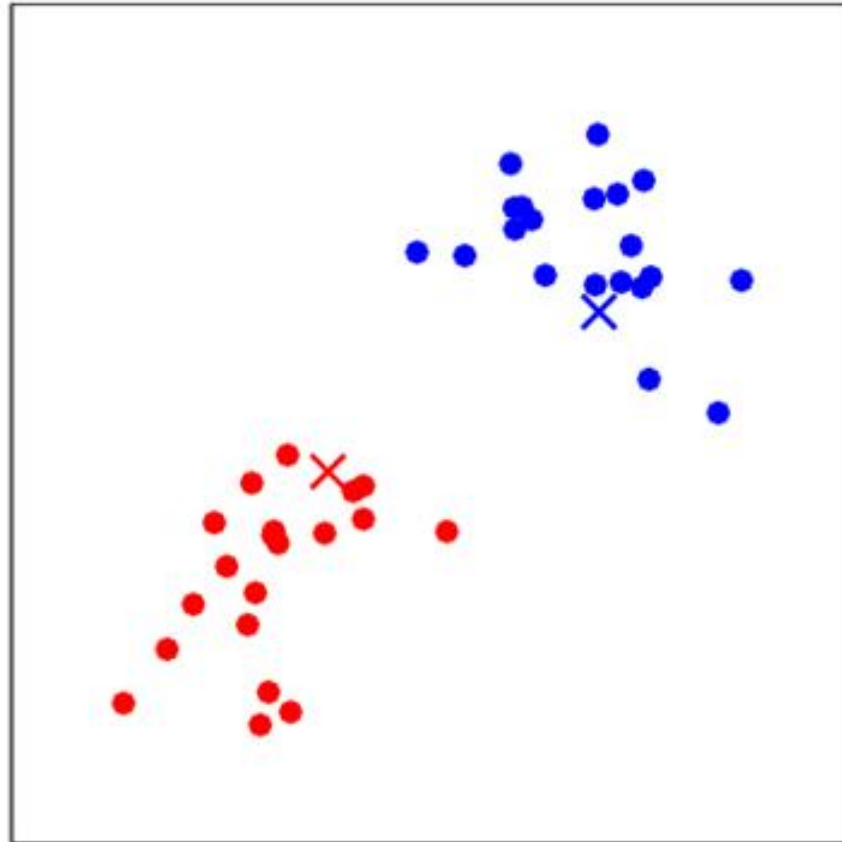
Example of K-mean clustering



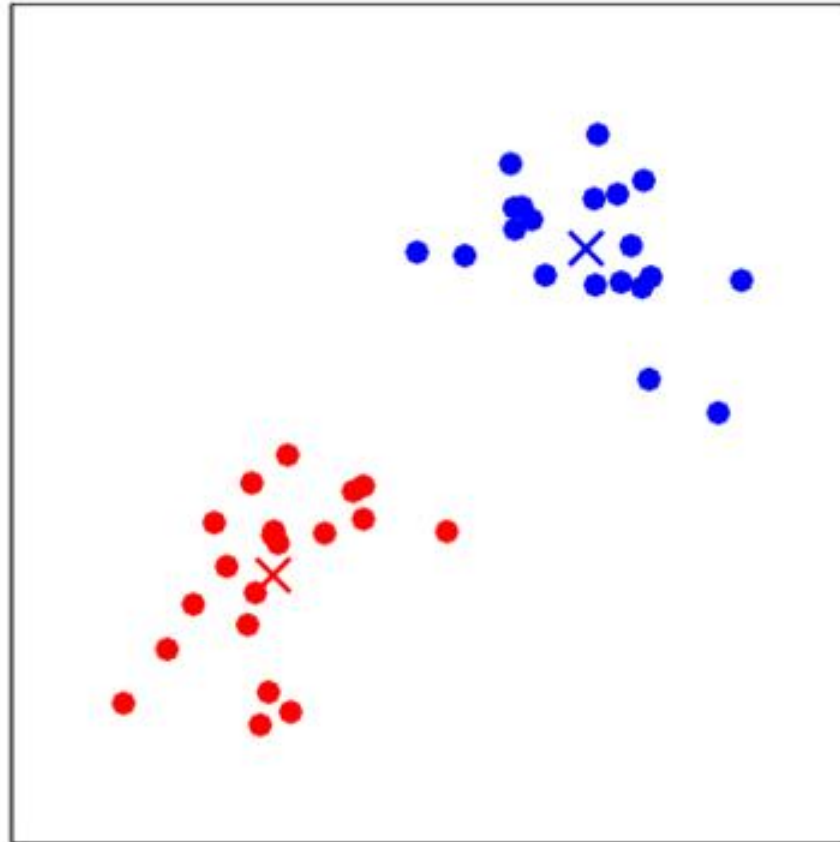
Example of K-mean clustering



Example of K-mean clustering



Example of K-mean clustering



Example of K-mean clustering

- In the Figure above for K-means algorithm: Training examples are shown as dots, and cluster centroids are shown as crosses.
- (a) Original dataset.
- (b) Random initial cluster centroids (in this instance, not chosen to be equal to two training examples).
- (c-f) Illustration of running two iterations of k-means.
- In each iteration, we assign each training example to the closest cluster centroid (shown by “painting” the training examples the same color as the cluster centroid to which is assigned); then we move each cluster centroid to the mean of the points assigned to it. (Best viewed in color.)
- Images courtesy Michael Jordan.

Q: Is the k-means algorithm guaranteed to converge?

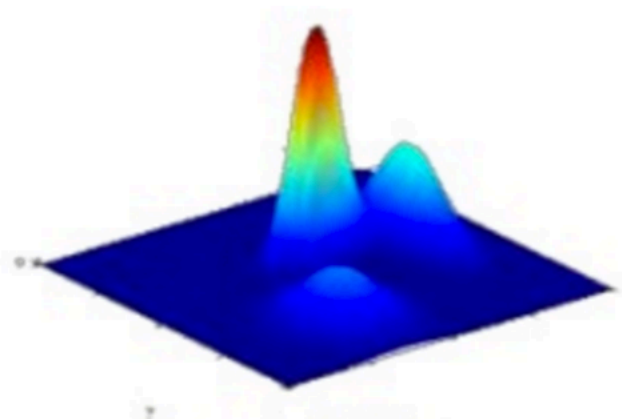
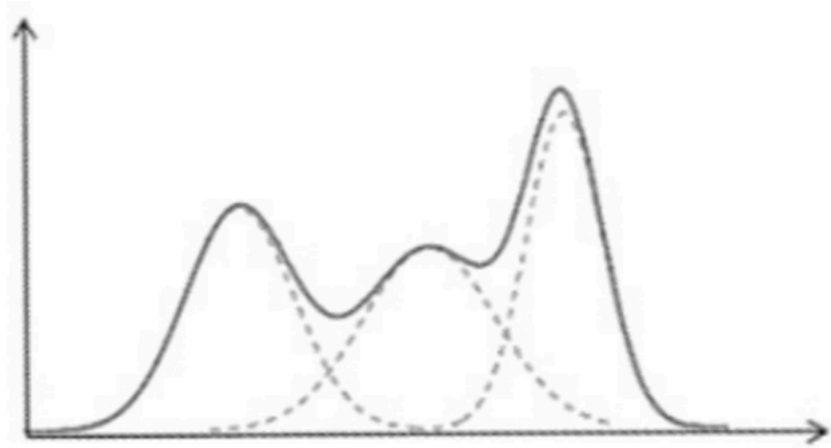
- Yes it is, but it might converge to a local optimization point instead of a global one in the following sense.
- Define the distortion function to be:
$$J(c, \mu) = \sum_{i=1}^m ||x^{(i)} - \mu_{c(i)}||^2$$
- Here J measures the sum of squared distances between each training example $x^{(i)}$ and the cluster centroid $\mu_{c(i)}$ to which it has been assigned.
- It can be shown that k-means is exactly coordinate descent on J .
- This means that the inner-loop of k-means repeatedly minimizes J with respect to c while holding μ fixed, and then minimizes J with respect to μ while holding c fixed.
- Thus, J must monotonically decrease, and the value of J must converge. (Usually, this implies that c and μ will converge too.)
- In theory, it is possible for k-means to oscillate between a few different clusterings—i.e., a few different values for c and/or μ —that have exactly the same value of J , but this almost never happens in practice.)

Note: The distortion function J is non-convex, so no global minimum is guaranteed.

- That is to say the coordinate descent on J is not guaranteed to converge to the global minimum: k-means can be susceptible to local optima.
- But very often k-means will work fine and come up with very good clusterings despite this.
- Try heuristic method if you are worried about getting stuck in bad local minima:
- Run k-means many times (using different random initial values for the cluster centroids μ_j).
- Then, out of all the different clusterings found, select the one that gives the lowest distortion $J(c, \mu)$.

Change gears:

Multivariate Gaussian Mixture Model



$$p(\boldsymbol{\theta}) = \sum_{i=1}^K \phi_i \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

*Mimic linear combination of vectors,
here is a convex combination.*

$$\phi_j \geq 0, \sum_{j=1}^K \phi_j = 1$$

where the i^{th} vector component is characterized by normal distributions weights ϕ_i , means $\boldsymbol{\mu}_i$ and covariance matrices $\boldsymbol{\Sigma}_i$.

Gaussian Mixture Models

Like K-Means, GMM clusters have centers.

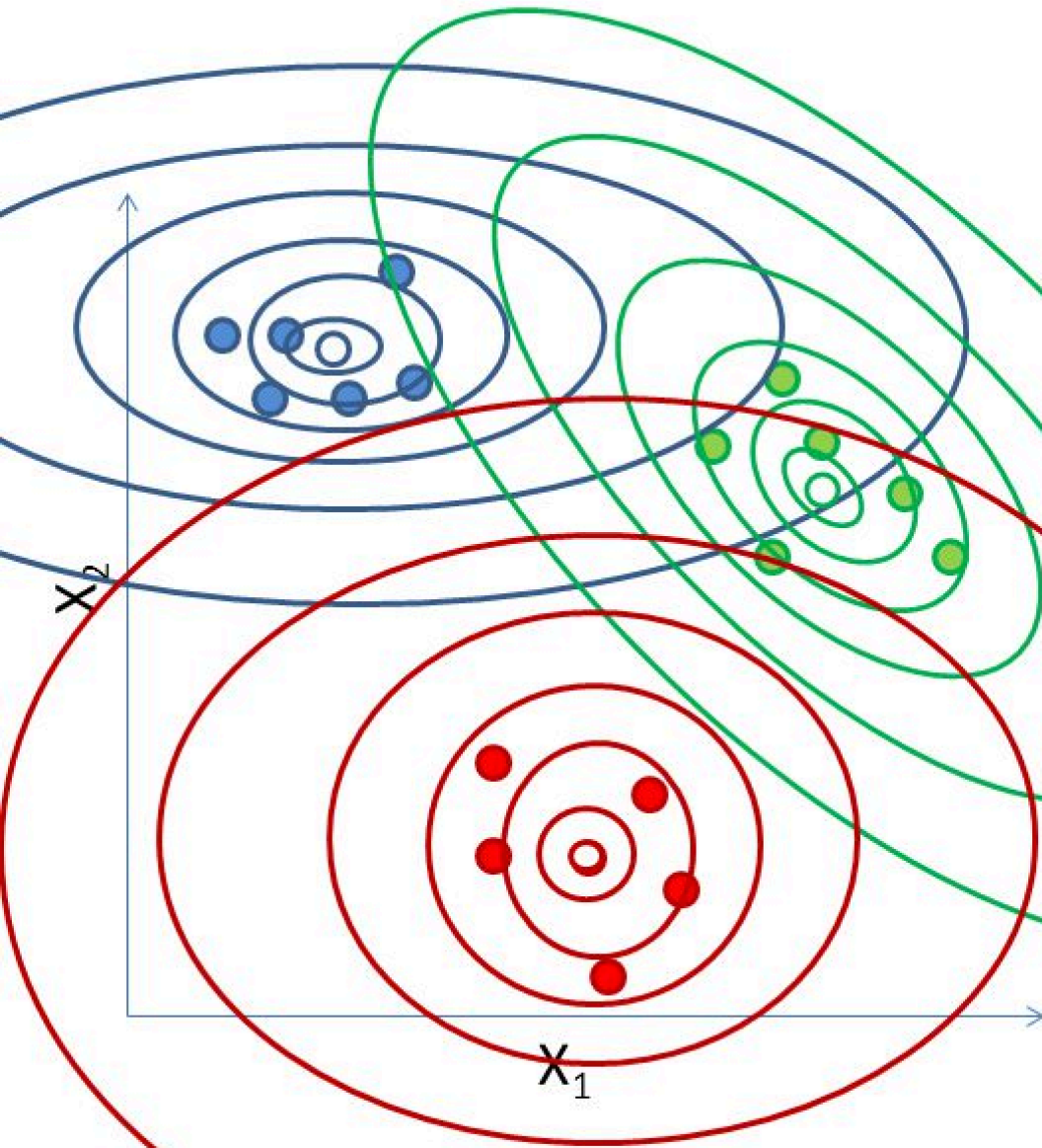
In addition, they have probability distributions that indicate the probability that a point belongs to the cluster.

These ellipses show “level sets”: lines with equal probability of belonging to the cluster.

Notice that green points still have SOME probability of belonging to the blue cluster, but it's much lower than the blue points.

This is a more complex model than K-Means: distance from the center can matter more in one direction than another.

A soft clustering methods. Sign more responsibility to one of the Gaussians.



How Gaussian mixture model and Expectation-Maximization (EM) related?

Key: the posterior distribution $p(\boldsymbol{\theta}|\mathbf{x})$ is *also* a Gaussian mixture mode !

$$p(\boldsymbol{\theta}|\mathbf{x}) = \sum_{i=1}^K \tilde{\phi}_i \mathcal{N}(\tilde{\boldsymbol{\mu}}_i, \tilde{\boldsymbol{\Sigma}}_i)$$

with new parameters $\tilde{\phi}_i$, $\tilde{\boldsymbol{\mu}}_i$ and $\tilde{\boldsymbol{\Sigma}}_i$ that are updated using the EM algorithm.

What is an EM algorithm?

EM = **Expectation-Maximization**

The EM (Expectation-Maximization) Algorithm

- Expectation of what?
- Maximization of what?
- *Work out details with students on the board.*

EM (Expectation-Maximization) Algorithm

- Given a training set $\{x^{(1)}, \dots, x^{(m)}\}$
- Note: since we are in the unsupervised learning setting, so these points do not come with any labels.
- Goal: Model the data by specifying a joint distribution $p(x^{(i)}, z^{(i)}) = p(x^{(i)} | z^{(i)})p(z^{(i)})$.

Here, $z^{(i)} \sim \text{Multinomial}(\phi)$ $\phi_j \geq 0, \sum_{j=1}^k \phi_j = 1$
and the parameter ϕ_j gives $p(z^{(i)} = j)$

and $x^{(i)} | z^{(i)} = j \sim \mathcal{N}(\mu_j, \Sigma_j)$

k denote the number of values that the $z^{(i)}$'s can take on.

Work out details with students on the
board

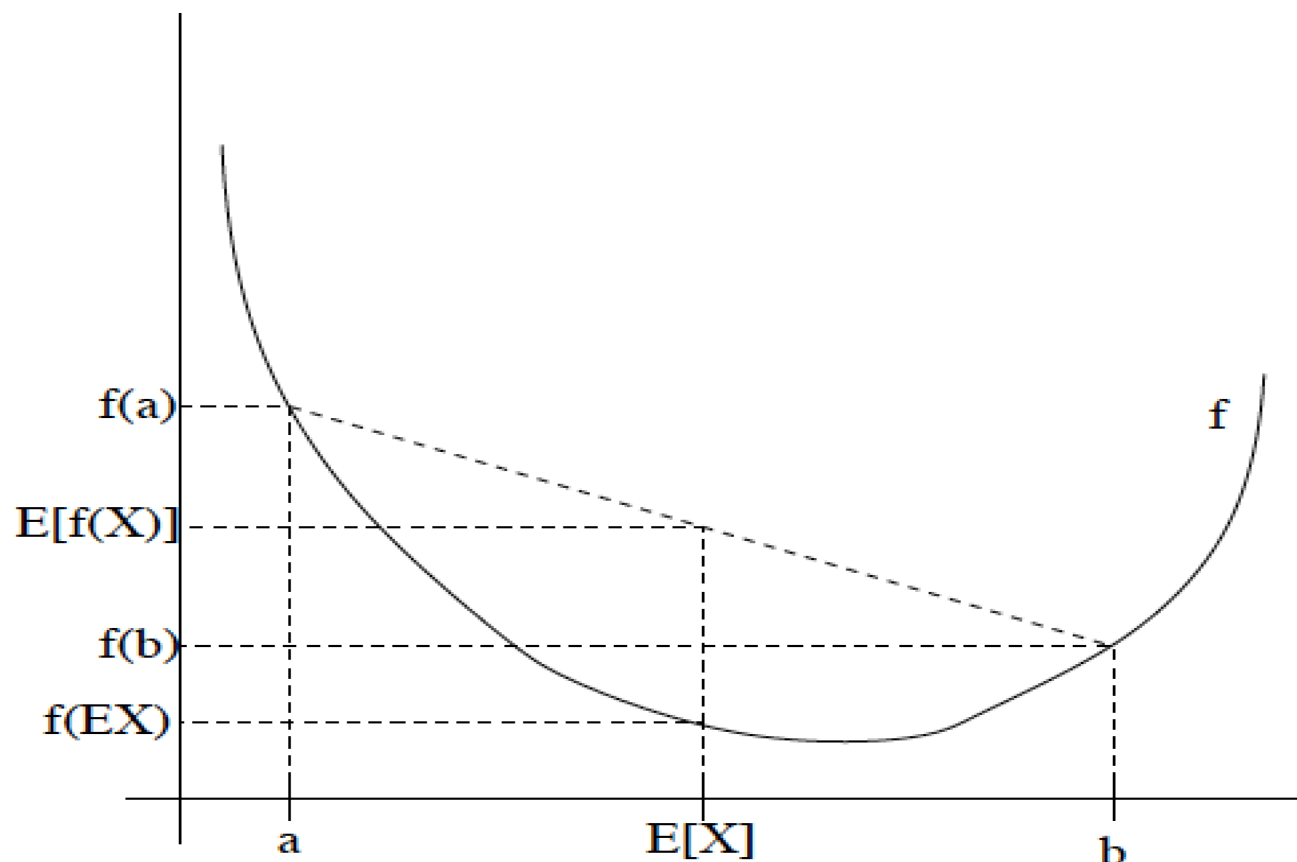
- The parameters of our model are thus ϕ , μ and Σ . To estimate them, write:

Jensen's inequality

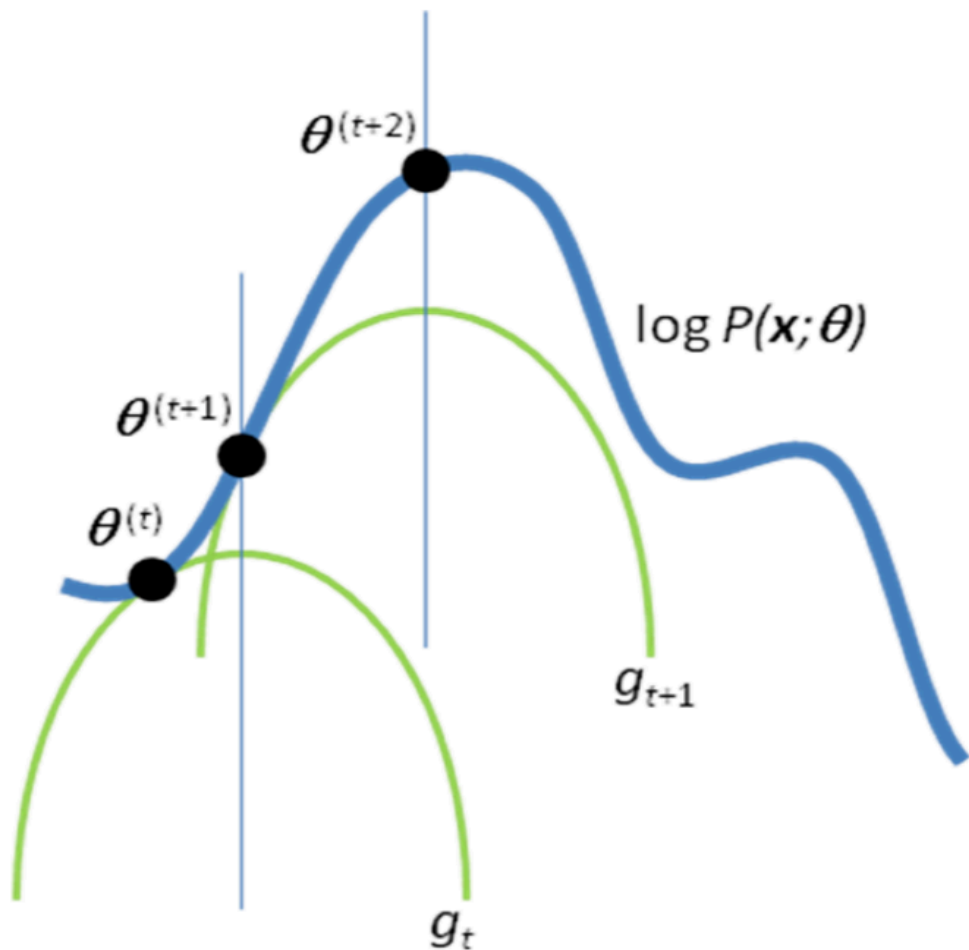
Theorem. Let f be a convex function, and let X be a random variable. Then:

$$E[f(X)] \geq f(EX).$$

Moreover, if f is strictly convex, then $E[f(X)] = f(EX)$ holds true if and only if $X = E[X]$ with probability 1 (i.e., if X is a constant).

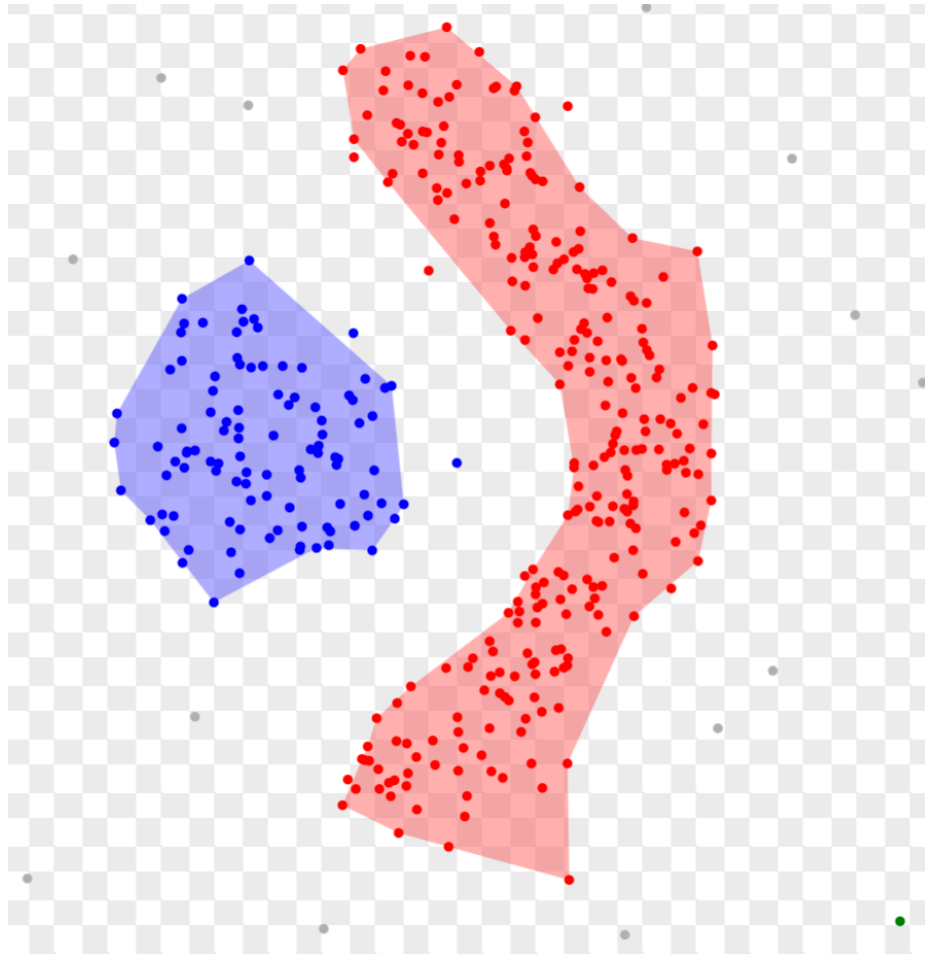


Geometry of EM algorithm

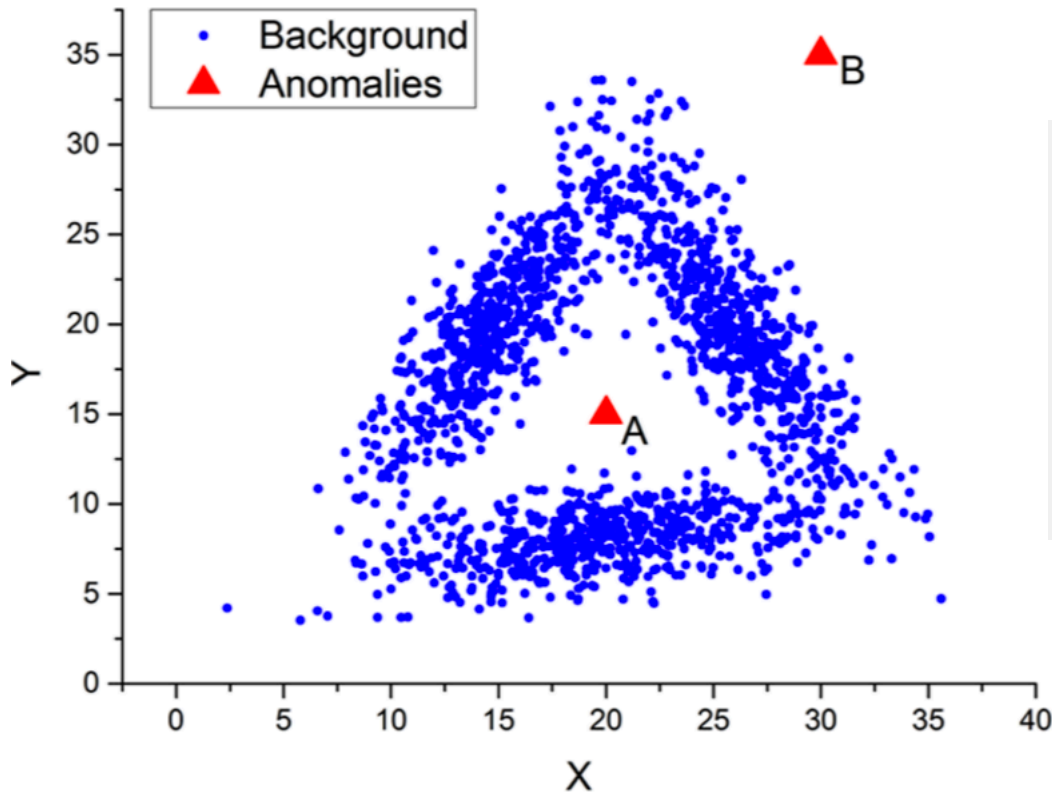


Supplementary Figure 1 Convergence of the EM algorithm. Starting from initial parameters $\theta^{(t)}$, the E-step of the EM algorithm constructs a function g_t that lower-bounds the objective function $\log P(x; \theta)$. In the M-step, $\theta^{(t+1)}$ is computed as the maximum of g_t . In the next E-step, a new lower-bound g_{t+1} is constructed; maximization of g_{t+1} in the next M-step gives $\theta^{(t+2)}$, etc.

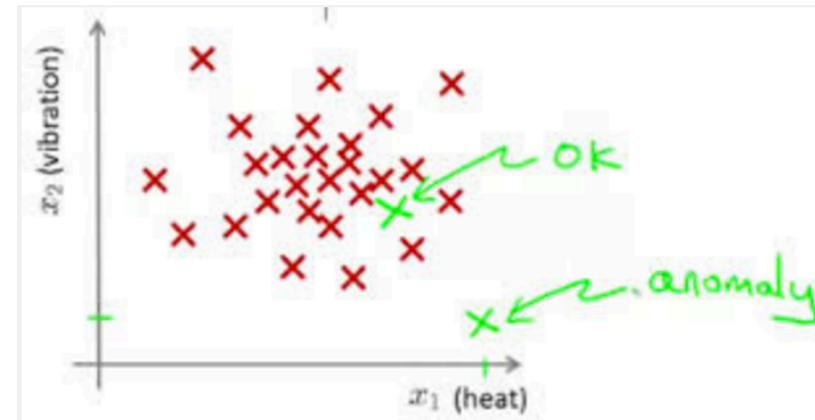
Density estimation using EM Algorithm



Anomaly Detection using Density Estimation

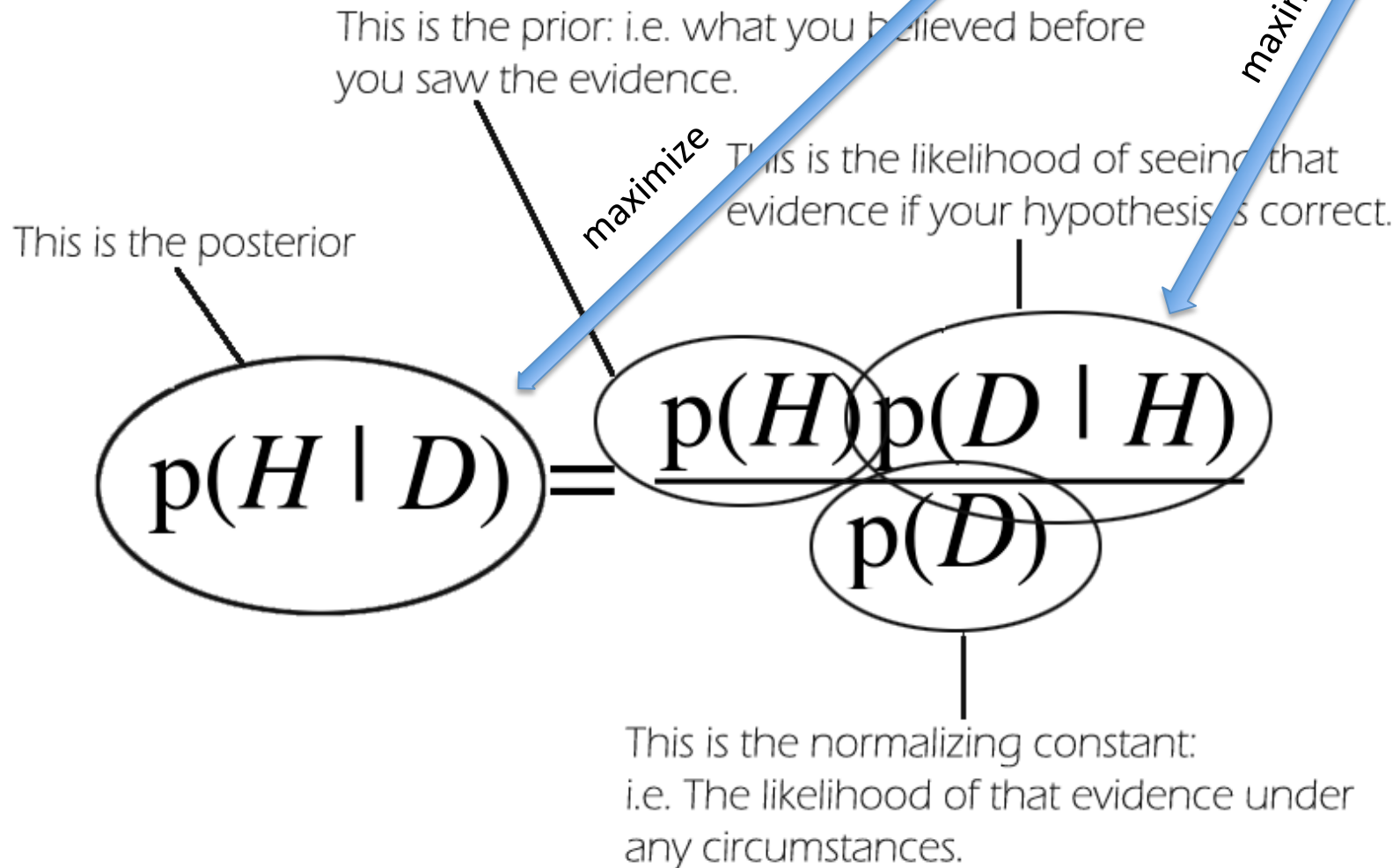


(a)



Recall: Maximum Likelihood Estimation (MLE) and Maximum a Posterior (MAP)

In both cases, the data D is given.



Recall: MLE =Maximum Likelihood Estimation

- In statistics, maximum likelihood estimation (MLE) is a method of estimating the parameters of a statistical model given observations, by finding the parameter values that maximize the likelihood of making the observations given the parameters.

Recall: MLE = Maximum Likelihood Estimate

Assume that we want to estimate an unobserved population parameter θ on the basis of observations x . Let f be the **sampling distribution** of x , so that $f(x \mid \theta)$ is the probability of x when the underlying population parameter is θ . Then the function:

$$\theta \mapsto f(x \mid \theta)$$

is known as the **likelihood function** and the estimate:

$$\hat{\theta}_{\text{ML}}(x) = \arg \max_{\theta} f(x \mid \theta)$$

is the maximum likelihood estimate of θ .

Recall: MAP

- Maximum a posteriori (MAP) estimation is a model of posterior distribution.
- The MAP can be used to obtain a point estimate of an unobserved quantity on the basis of empirical data.

Now assume that a prior distribution g over θ exists. This allows us to treat θ as a **random variable** as in **Bayesian statistics**. We can calculate the **posterior distribution** of θ using **Bayes' theorem**:

$$\theta \mapsto f(\theta \mid x) = \frac{f(x \mid \theta) g(\theta)}{\int_{\vartheta \in \Theta} f(x \mid \vartheta) g(\vartheta) d\vartheta}$$

where g is density function of θ , Θ is the domain of g .

The method of maximum a posteriori estimation then estimates θ as the **mode** of the posterior distribution of this random variable:

$$\hat{\theta}_{\text{MAP}}(x) = \arg \max_{\theta} f(\theta \mid x) = \arg \max_{\theta} \frac{f(x \mid \theta) g(\theta)}{\int_{\vartheta} f(x \mid \vartheta) g(\vartheta) d\vartheta} = \arg \max_{\theta} f(x \mid \theta) g(\theta).$$

The denominator of the posterior distribution (so-called **marginal likelihood**) does not depend on θ