

ColumbiaX: Machine Learning

Lecture 16

Prof. John Paisley

Department of Electrical Engineering
& Data Science Institute

Columbia University

↗ focus on this lecture

SOFT CLUSTERING VS HARD CLUSTERING MODELS

↳ which we saw with k-means 2 lectures ago.

HARD CLUSTERING MODELS

Review: K-means clustering algorithm

Given: Data x_1, \dots, x_n , where $x \in \mathbb{R}^d$

Goal: Minimize $\mathcal{L} = \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}\{c_i = k\} \|x_i - \mu_k\|^2$. * over c_i, μ_k

- Iterate until values no longer changing

1. Update c : For each i , set $c_i = \arg \min_k \|x_i - \mu_k\|^2$

2. Update μ : For each k , set $\mu_k = (\sum_i x_i \mathbb{1}\{c_i = k\}) / (\sum_i \mathbb{1}\{c_i = k\})$

~~Then when we allocate datapoints to nearest cluster, we update cluster centroids by simply averaging.~~
K-means is an example of a *hard clustering* algorithm because it assigns each observation to only one cluster. (after it converges and outputs the final clustering.) And we don't have any uncertainty. **

In other words, $c_i = k$ for some $k \in \{1, \dots, K\}$. There is no accounting for the “boundary cases” by hedging on the corresponding c_i .

* objective function:

an objective that looks like this where we sum over each data point. The squared distance to the cluster centroid that it was assigned to as encoded by the parameters C_i .

So X_i is assigned to one of capital K clusters and C_i indexes what cluster that is.

So we sum up the total number of square distances.

And that's our objective that we're trying to minimize

** What else we want to make?

However you could imagine that on the boundary of this clusters for example a data point that is could go either way.

We might want want encode our uncertainty about the clustering of that data point with some sort of a weight.

Whereas the points that are right in the middle of the cluster, we would want to be more confident about

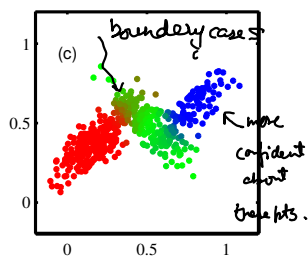
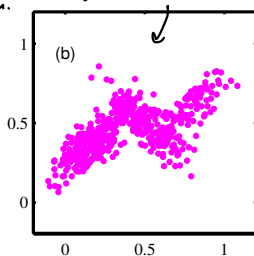
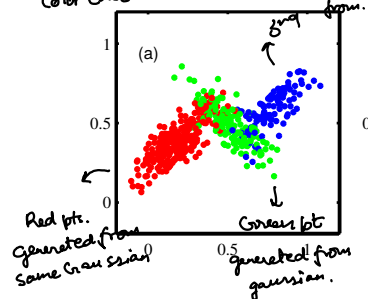
This leads to concept of a soft clustering algorithm.

SOFT CLUSTERING MODELS

A soft clustering algorithm breaks the data across clusters intelligently.

color encodes which Gaussian it came from.

what we see



(left) True cluster assignments of data from three Gaussians.

(middle) The data as we see it.

(right) A soft-clustering of the data accounting for borderline cases.

WEIGHTED K-MEANS (SOFT CLUSTERING EXAMPLE)

(considering this in the context of k-means leads to weighted k-means. Modification here)

Weighted K-means clustering algorithm (means from hard to soft clustering algorithm)

Given: Data x_1, \dots, x_n , where $x \in \mathbb{R}^d$ * ② $-\xi_i H(\phi_i)$

Goal: Minimize $\mathcal{L} = \sum_{i=1}^n \sum_{k=1}^K \phi_i(k) \frac{\|x_i - \mu_k\|^2}{\beta}$ But here divide by β . over ϕ_i and μ_k

Conditions: $\phi_i(k) > 0$ and $\sum_{k=1}^K \phi_i(k) = 1$. Set parameter $\beta > 0$. * *

► Iterate the following , $H(\phi_i) = \text{entropy}$

1. Update ϕ : For each i , update the word allocation weights

$$\phi_i(k) = \frac{\exp\{-\frac{1}{\beta}\|x_i - \mu_k\|^2\}}{\sum_j \exp\{-\frac{1}{\beta}\|x_i - \mu_j\|^2\}}, \text{ for } k = 1, \dots, K$$

if x_i is very far away from μ_k , this value will be small.

2. Update μ : For each k , update μ_k with the weighted average

$$\mu_k = \frac{\sum_i x_i \phi_i(k)}{\sum_i \phi_i(k)}$$

each data point weight by the probs. of coming from that cluster.

② Instead of multiplying it by an indicator that says x_i was assigned to cluster k , we multiply it by a weight $\phi_i(k)$.

$\phi_i(k)$ → data point we are looking at. → cluster we are considering.

These ϕ s now can be viewed as probability weights. So each value of $\phi_i(k) \geq 0$.

And sum of them is equal to 1 (over k).

→ And so in this case, we're breaking each data point up into different clusters according to how probable we think that cluster is in generating the data point we're looking at.

** For a particular iteration, we look through each of the data points and instead of assigning them to each cluster we assign responsibilities across the clusters according to weight ϕ .

And so in this case, we're breaking each data point up into different clusters according to how probable we think that cluster is in generating the data point we're looking at.

③ So this is how we take proximity to the cluster centroids into account. To break the responsibility for that given data point across the clusters according to how close that data point is to each cluster.

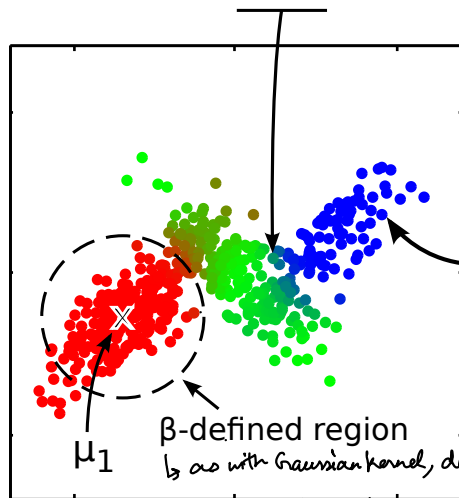
Difference:

So you can imagine that if these ϕ 's are pure indicators meaning the vector of all zeroes except for a one in the dimension indicating what cluster x_i came from.

That would correspond to a hard clustering algorithm and we would simply be summing up all the points in the k th cluster and dividing by the total number of points. Now we are softening that by taking fractions of points and splitting them across clusters

SOFT CLUSTERING WITH WEIGHTED K-MEANS

$\phi_i = 0.75$ on green cluster & 0.25 blue cluster



The border points are split across clusters to account for the fact that we're less certain which cluster they come from.

What does it mean to be close and what does it mean to be far away? A parameter we have to set.

Difference k-means and mixture models:

k-means;

With soft k-means we took each data point and we, instead of assigning it to one and only one cluster, we broke that data point across clusters using a vector of weights.

And we loosely use the term probability vector, even though there were no probabilistic assumptions.

MIXTURE MODELS

Mixture Models;

So now we're gonna look at a modeling framework where we do make these sorts of probabilistic assumptions, by defining probability Distributions.

PROBABILISTIC SOFT CLUSTERING MODELS

Probabilistic vs non-probabilistic soft clustering

The weight vector ϕ_i is *like* a probability of x_i being assigned to each cluster.

A **mixture model** is a probabilistic model where ϕ_i actually *is* a probability distribution according to the model.

Mixture models work by defining:

(we explicitly define a probability distribution on it.)

- ▶ A prior distribution on the cluster assignment indicator c_i
- ▶ A likelihood distribution on observation x_i given the assignment c_i

So, c_i is going to come back into our model, and we're going to define a probability distribution.

So, we're now using probabilistic definition to put these constraints on our parameters.

Intuitively we can connect a mixture model to the Bayes classifier:

- ▶ Class prior \rightarrow cluster prior. This time, we *don't* know the "label" **
- ▶ Class-conditional likelihood \rightarrow cluster-conditional likelihood ***

* It's the prior on a particular data point coming from a particular cluster, except this time now, we don't know the label for each observation.

**

Previously the cluster, each cluster was like a class, and we knew the label for each observation, and therefore we knew which cluster it came from. Now we don't know that information

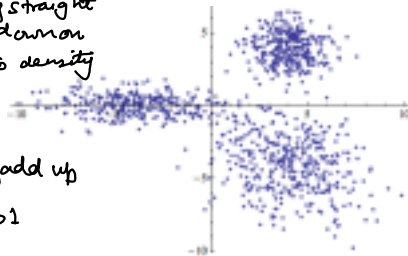
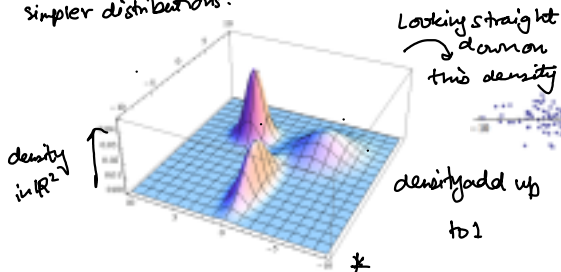
*** Conditional distribution similarity:

And similarly with the bayes classifier we had the class conditional distributions, so given what class it comes from we have a specific class distribution on the data. Now we have cluster conditional distributions, so condition on a data point coming from a particular cluster we have. A specific distribution on that data point for that, coming from that cluster.

MIXTURE MODELS

weighted combination of simpler distributions.

- generative model, we're going to define a probability distribution on the data and try to learn parameters of that distribution.



(a) A probability distribution on \mathbb{R}^2 .

(b) Data sampled from this distribution.

And now we want to propose a Gaussian mixture model. And then learn the resulting means, covariances and mixing weights on these 3 different Gaussians.

Before introducing math, some key features of a mixture model are:

1. It is a generative model (defines a probability distribution on the data)
2. It is a weighted combination of simpler distributions.
 - ▶ Each simple distribution is in the same distribution family (i.e., a Gaussian).
 - ▶ The “weighting” is defined by a discrete probability distribution.

* These three distributions are actually 3 different Gaussian distributions with different means and covariances. We weight each of these and put them together according to the probability of particular Gaussian to get our final density.

^{to define 2}
MIXTURE MODELS [↗] define a generative process. (same with all bayesian models. If we want to define what's it doing we have to know how data is generated from that model.)

Generating data from a mixture model

Data: x_1, \dots, x_n , where each $x_i \in \mathcal{X}$ (can be complicated, but think $\mathcal{X} = \mathbb{R}^d$)

Model parameters: A K -dim distribution π and parameters $\theta_1, \dots, \theta_K$.
^{*}
^{on each of those clusters.}

Generative process: For observation number $i = 1, \dots, n$, ^{1st cluster parameter}

1. Generate cluster assignment: $c_i \stackrel{iid}{\sim} \text{Discrete}(\pi)$ ^{**} $\Rightarrow \text{Prob}(c_i = k | \pi) = \pi_k$. ^{*}

2. Generate observation: $x_i \sim p(x | \theta_{c_i})$. ⁰
^{↳ class of probability distributions.}

Some observations about this procedure:

- ▶ First, each x_i is randomly assigned to a cluster using distribution π .
- ▶ c_i indexes the cluster assignment for x_i
 - ▶ This picks out the index of the parameter θ used to generate x_i .
 - ▶ If two x 's share a parameter, they are clustered together. (generated from the same probability distribution.)

- * if we assume there are k different clusters, we have k dimensional probability distribution on each of those clusters, and also the parameters for the clusters.
- ** Every single data point that we are going to assume is generated independently with same discrete π distribution.
- *** Meaning that the probability that we assign any given data point to the k^{th} cluster given π is equal to π_k .
- o Then, when we assign it to the cluster, we pick the parameter to generate the data point according to the indexing of the cluster. And so we use the same exact probability distribution.

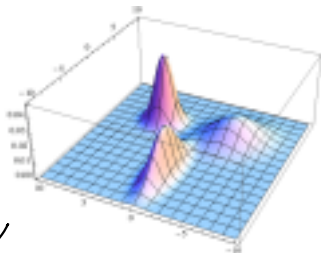
Distribution of p :

So we're going to assume that P is the same for all clusters.

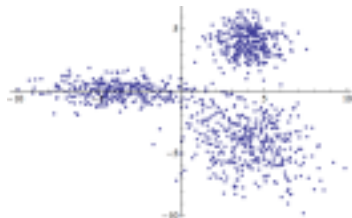
So it's, for example, a multivariate Gaussian or something more complicated than that.

The only thing that's differing across clusters is what parameter we use for that cluster. ($\theta - \mu, \Sigma$ for each)

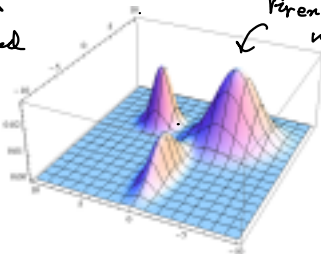
MIXTURE MODELS (Gaussian setting to see the impact of weights only.)



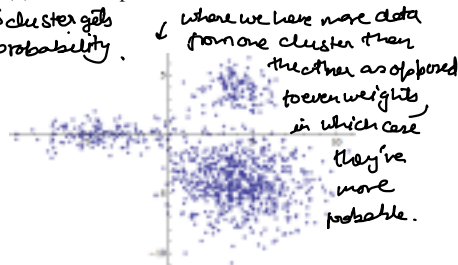
(a) Uniform mixing weights



(b) Data sampled from this distribution.



(c) Uneven mixing weights



(d) Data sampled from this distribution.

same Gaussian but we have changed weights

Here, this cluster gets more probability.

where we have more data from one cluster than the other as opposed to even weights, in which case they're more probable.

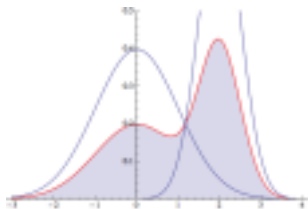
Gaussians define the regions & the weights define how dense the data is in each region.

GAUSSIAN MIXTURE MODELS

ILLUSTRATION

Gaussian mixture models are mixture models where $p(x|\theta)$ is Gaussian.

Mixture of two Gaussians



The red line is the density function.

$$\pi = [0.5, 0.5]$$

$$(\mu_1, \sigma_1^2) = (0, 1)$$

$$(\mu_2, \sigma_2^2) = (2, 0.5)$$

Influence of mixing weights



The red line is the density function.

$$\pi = [0.8, 0.2]$$

$$(\mu_1, \sigma_1^2) = (0, 1)$$

$$(\mu_2, \sigma_2^2) = (2, 0.5)$$

GAUSSIAN MIXTURE MODELS (GMM)

The model

Parameters: Let π be a K -dimensional probability distribution and (μ_k, Σ_k) be the mean and covariance of the k th Gaussian in \mathbb{R}^d .

Generate data: For the i th observation, We've defined the hypothesis of how the data got to us by defining the generative model. Now,

1. Assign the i th observation to a cluster, $c_i \sim \text{Discrete}(\pi)$
2. Generate the value of the observation, $x_i \sim N(\mu_{c_i}, \Sigma_{c_i})$ that defines a joint likelihood on the data

Definitions: $\mu = \{\mu_1, \dots, \mu_K\}$ and $\Sigma = \{\Sigma_1, \dots, \Sigma_K\}$. given some parameters.

Goal: We want to learn π , μ and Σ .

Once we have the data, our goal is go back and infer what are the parameters that could have generated that data. (Inference problem)

GAUSSIAN MIXTURE MODELS (GMM)

Maximum likelihood

Objective: Maximize the likelihood over model parameters π , μ and Σ by treating the c_i as auxiliary data using the EM algorithm.

↳ additional hidden data that we introduce. So we want to integrate out c_i .

$$p(x_1, \dots, x_n | \pi, \mu, \Sigma) = \prod_{i=1}^n p(x_i | \pi, \mu, \Sigma) = \prod_{i=1}^n \sum_{k=1}^K p(x_i, c_i = k | \pi, \mu, \Sigma)$$

i.e. assumption simplification (under the first product)
no c_i here (we integrated it out) (under the second product)

The summation over values of each c_i “integrates out” this variable.

why? c_i is Θ_L (last lecture)

We can't simply take derivatives with respect to π , μ_k and Σ_k and set to zero to maximize this because there's no closed form solution. (can't solve it exactly).

We could use gradient methods, but EM is cleaner.

Goal to maximize the joint likelihood of data given only π , μ & Σ , where we have summed the parameters c_i out of the model,

(because we can get closed form updates.)

* So CI , the cluster assignment now is going to be the additional hidden data that we introduced, and then do. So we want to integrate out CI . Remember, from the Algorithm that we discussed last time. What that means is that we wanted to maximize the likelihood of all of the data, x_1 through x_n , given π , μ , and σ .

* *

But now we integrate or in this case, because CI can only take a discrete and finite number of values. We sum over all those possible values to integrate them out.

- Joint likelihood of the i^{th} observation and cluster assignment which cluster the i^{th} observation came from.

EM ALGORITHM

Q: Why not instead just include each c_i and maximize $\prod_{i=1}^n p(x_i, c_i | \pi, \mu, \Sigma)$ since (we can show) this is easy to do using coordinate ascent?

*

A: We would end up with a hard-clustering model where $c_i \in \{1, \dots, K\}$.
Our goal here is to have soft clustering, which the sum effectively does.

EM and the GMM

We will not derive everything from scratch. However, we can treat c_1, \dots, c_n as the auxiliary data that we integrate out. *(latent data, missing data in our model, that we want to integrate over.)*

Therefore, we use EM to

$$\text{maximize} \quad \sum_{i=1}^n \ln p(x_i | \pi, \mu, \Sigma) \quad \text{by using} \quad \sum_{i=1}^n \ln p(x_i, c_i | \pi, \mu, \Sigma)$$

marginal likelihood where c_i is gone *log of this augmented / extended joint likelihood.*

Let's look at the outlines of how to derive this.

*

Because when we do that, when we write out the augmented joint likelihood, where we include C_i as a parameter we want to do a point estimate over, we take the log of this, and we take derivatives and solve. What we would find is that we can optimize that very easily, we can get closed form updates for every parameter in the model including C_i . So that is something that can be done

But,

what we would find by doing that is we would end up with a hard clustering model. Meaning that, when we wanted to update a particular cluster assignment, C_i , the maximizing that C_i would lead to assigning the data point to the most probable cluster. So the maximum likelihood solution, or the map solution, actually, for C_i , is to assign the i th data point to the most probable cluster. And so we would have a hard clustering algorithm.

Goal;

Our goal here is to integrate out the effect of the cluster assignment. Which we're going to see when we do that in the Framework. It's going to lead to a soft clustering algorithm that we anticipate should be more accurate and should be better than the hard clustering algorithm because we've integrated out one of the unknowns in our model. So we have fewer potential local optimal solutions that we could converge to.

THE EM ALGORITHM AND THE GMM (How the equation is modified for this model structure?)

From the last lecture, the generic EM objective is

We introduced the additional variable to our model Θ_2 and integrate it out.

And we do a 2-step procedure:

maximize the log of probability of the x given parameter Θ_1

$$\ln p(x|\theta_1) = \int q(\theta_2) \ln \frac{p(x, \theta_2|\theta_1)}{q(\theta_2)} d\theta_2 + \int q(\theta_2) \ln \frac{q(\theta_2)}{p(\theta_2|x, \theta_1)} d\theta_2$$

$K L(\theta_2 | p)$

The EM objective for the Gaussian mixture model is

Take this equality, and we just write it using the Gaussian mixture model parameters:

Conditional posterior distribution.

think as this

$$\sum_{i=1}^n \ln p(x_i|\pi, \mu, \Sigma) = \sum_{i=1}^n \sum_{k=1}^K q(c_i = k) \ln \frac{p(x_i, c_i = k|\pi, \mu, \Sigma)}{q(c_i = k)}$$

In next slide, we are going to ignore it not involve the parameters we're interested in.

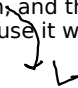
Θ_1 corresponds to π
 Θ_2 corresponds to c_i

$$\sum_{i=1}^n \sum_{k=1}^K q(c_i = k) \ln \frac{q(c_i = k)}{p(c_i|x_i, \pi, \mu, \Sigma)}$$

it's k diverges

Because c_i is discrete, the integral becomes a sum.

*
remember the procedures were to set Q equal to P .
Given that updated Q distribution,
calculate this expectation, and then maximize this expectation over θ_1 .
And θ_2 is gone because it was integrated out.



* *

So, the log of the likelihood of all of the data, given π , μ and σ
because the data are assumed to be IID, that turns into the log of the product,
which turns into the sum of the log of each individual likelihood.
So we have a sum over all of the data of the log of the probability of each point
given our Gaussian mixture model parameters.

○ And because c_i is a discrete variable the integral is going to turn in to a sum.

So when you have continuous
variables that you integrate when you have discrete ones you sum.

so this integral, overall, the values that θ_2 can take turns into a sum
overall the values that C_i can take.

EM SETUP (ONE ITERATION)

(Set KL Divergence = 0)

you do this previously k

First: Set $q(c_i = k) \leftarrow p(c_i = k | x_i, \pi, \mu, \Sigma)$ using Bayes rule:

$$p(c_i = k | x_i, \pi, \mu, \Sigma) \propto p(c_i = k | \pi) p(x_i | c_i = k, \mu, \Sigma)$$

likelihood of the i^{th} observation

conditional posterior assigning it to k^{th} cluster

given that it was assigned to cluster.

We can solve the posterior of c_i given π, μ and Σ :

$$q(c_i = k) = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_i | \mu_j, \Sigma_j)} \Rightarrow \phi_i(k)$$

IInd part after we update each of these key distributions to

↑ serving same purpose as weighted k-means algo. Now we have a firm probabilistic interpretation of what we're doing.

E-step: Take the expectation using the updated q 's

Sum over each observation,

$$Q = \sum_{i=1}^n \sum_{k=1}^K \phi_i(k) \ln p(x_i, c_i = k | \pi, \mu_k, \Sigma_k) + \text{constant w.r.t. } \pi, \mu, \Sigma$$

(a)

M-step: Maximize Q with respect to π and each μ_k, Σ_k .

(a)

weight we assign to k^{th} cluster.

We can expand this, write out the Gaussian form of it. And target this $f^{(n)}$ that we want to maximise.

what does this mean?

* The prior probability of assigning a data point to cluster k is just equal to π_k .

** The likelihood of i^{th} observation coming from cluster k is equal to Gaussian
→ where we evaluate the likelihood of x_i given the μ & Σ of k^{th} cluster.

→ Normalization

And then to turn that in to a probability distribution we simply normalize the right-hand side.

So this is a proportionality, to turn this into an equality we normalize this thing by dividing over the sum of all possible values that z_i can take, meaning we sum this numerator over all values k .

So we take this divided by the sum over all possible values for the cluster assignment.

* * * Where we evaluate the log of the likelihood of the data coming from each cluster.

So, that's this term the log and the likelihood of the data point coming from k^{th} cluster plus the log of the prior probability of it coming from the k^{th} cluster.

So that's equivalent to the log of the joint likelihood of the observation coming from the k^{th} cluster.

And it coming from, any point coming from the k^{th} cluster.

?

← joint?

M-STEP CLOSE UP

* the log of likelihood of data given π, μ and Σ where we integrate out c .

Aside: How has EM made this easier?

Original objective function:

$$\mathcal{L} = \sum_{i=1}^n \ln \sum_{k=1}^K p(x_i, c_i = k | \pi, \mu_k, \Sigma_k) = \sum_{i=1}^n \ln \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k).$$

Annotations:
 - An arrow points from the handwritten note "the log of likelihood of data given π, μ and Σ where we integrate out c ." to the term $p(x_i, c_i = k | \pi, \mu_k, \Sigma_k)$.
 - An arrow points from the handwritten note "how?" to the term $\sum_{k=1}^K$.
 - An arrow points from the handwritten note "integrate out c " to the term c_i .
 - An arrow points from the handwritten note "the log of likelihood of data given π, μ and Σ where we integrate out c ." to the term \ln .

The log-sum form makes optimizing π , and each μ_k and Σ_k difficult.

Using EM here, we have the M-Step:

Before we had a logsum, we turned it into a sum of log.

$$Q = \sum_{i=1}^n \sum_{k=1}^K \phi_i(k) \underbrace{\{\ln \pi_k + \ln N(x_i | \mu_k, \Sigma_k)\}}_{\ln p(x_i, c_i = k | \pi, \mu_k, \Sigma_k)} + \text{constant w.r.t. } \pi, \mu, \Sigma$$

The sum-log form is easier to optimize. We can take derivatives and solve.

* And this log sum,
this sum appearing inside of the log is what makes it difficult.
If you took the derivative of this with respect to μ or σ or π you would see that this sum inside is
what's the thing that's making everything complicated for us.

↓ why

* If we now take the derivative of this and try to maximize it with respect to π , μ and σ , we'll find that it's very easy to do that.
And so that's what we do.
We actually take the derivative of this thing with respect to μ , the means, the covariances, and also π , conditioned with the constraint that π is a probability distribution.
And we maximize this function with respect to those three unknown parameters.
And we'll find that we get a closed form solution.

EM FOR THE GMM

Algorithm: Maximum likelihood EM for the GMM

Given: x_1, \dots, x_n where $x \in \mathbb{R}^d$

Goal: Maximize $\mathcal{L} = \sum_{i=1}^n \ln p(x_i | \pi, \mu, \Sigma)$.
Log of marginal likelihood of each data point over π, μ, Σ where there is no cluster assignment at all to be found because it is integrated out.

► Iterate until incremental improvement to \mathcal{L} is “small”

1. **E-step:** For $i = 1, \dots, n$, set *

(Prob. of i^{th} data pt. coming from k^{th} cluster.)

$$\phi_i(k) = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_i | \mu_j, \Sigma_j)}, \quad \text{for } k = 1, \dots, K$$

how did you get these expressions?

2. **M-step:** For $k = 1, \dots, K$, define $n_k = \sum_{i=1}^n \phi_i(k)$ and update the values

Then we update 3 parameters of the model:

(prior prob. of the k^{th} cluster.)

$$\pi_k = \frac{n_k}{n}, \quad \mu_k = \frac{1}{n_k} \sum_{i=1}^n \phi_i(k) x_i, \quad \Sigma_k = \frac{1}{n_k} \sum_{i=1}^n \phi_i(k) (x_i - \mu_k)(x_i - \mu_k)^T$$

soft-covariance matrix

Comment: The updated value for μ_k is used when updating Σ_k .

empirical distribution over the clusters acc. to our assignments (our soft assignments)

1. We perform soft-clustering of the data by taking each data point and calculate conditional posterior probability of that data point coming from each of the k different clusters.

And we can do that this way **, where calculate the prior probability of the k^{th} cluster times the likelihood of the observation given it comes from

k^{th} cluster. And then divide by the sum of that over each cluster,

o sum of the expected no. of points that we're going to see coming from the k^{th} cluster acc to the current iteration

oo

For the mean, we then take each data point.

We multiply it by the probability of that data point coming from the K^{th} cluster.

We sum it up and

divide by the expected total number points coming from that cluster.

Mean calculation same (difference from k means).

So notice that this actually is identical to weighted k means.

This is exactly the same as weighted k means.

The only difference is the way that we calculated this ϕ .

So this ϕ , the way we calculated ϕ with the Gaussian mixture model, used this equation.

If you look at the way we calculated ϕ for weighted k means, it used a slightly different equation.

[what is difference in ϕ]

But then once we have the soft clustering of the i^{th} observation of each observation, we can calculate the mean exactly in the same way.

Covariance matrix

But then in addition to that with the Gaussian mixture model, we get a covariance matrix for the K th cluster, which is also like the weighted empirical covariance of the data coming from the K th cluster.

So from μ_k here, this is the mean of the K th cluster.

We use the updated μ here.

We calculate this outer product.

We multiply it by the probability that the i th point came from the K th cluster to begin with.

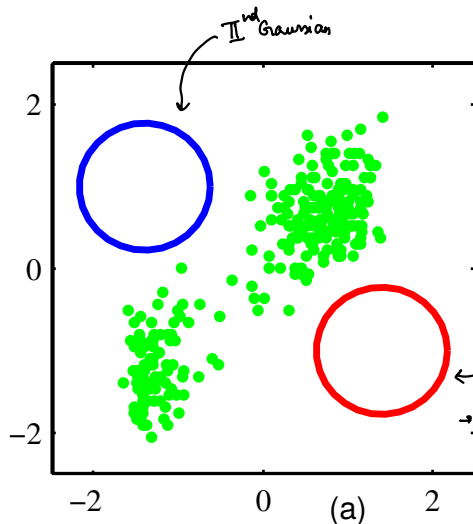
Sum that up over each data point.

And then divide by the expected total number of points.

So this is a soft covariance matrix.

GAUSSIAN MIXTURE MODEL: EXAMPLE RUN

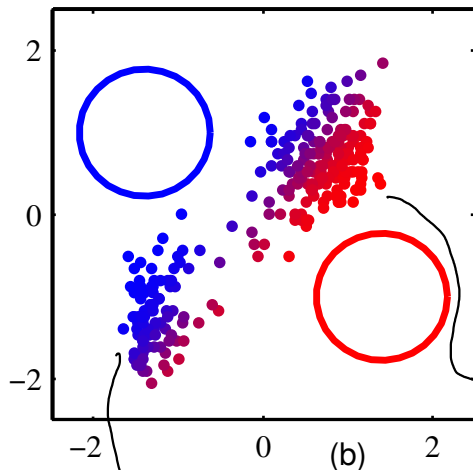
In this case, we define 2 Gaussians.



A random initialization

→ Circle indicates spherical covariance for each of these Gaussians.
Also, they are equally probable, the wt. is not being shown on the prob. of each of these clusters.

GAUSSIAN MIXTURE MODEL: EXAMPLE RUN



Iteration 1 (E-step)

Assign data to clusters

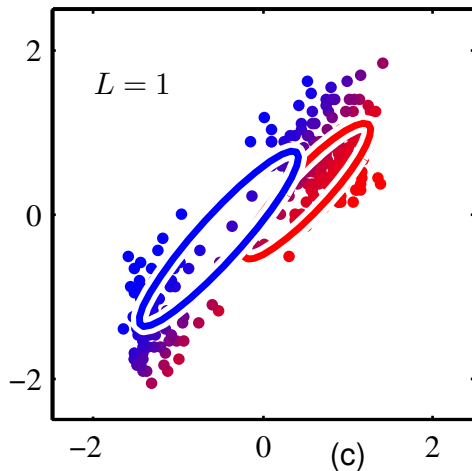
We take data point and calculate the posterior prob. of that point coming from each of these clusters.

color code based on the weight of prob. so this point is entirely red to indicate its almost entirely coming from this cluster.

in the middle, its 50/50 on either of these 2 clusters.

Diff. from k-means. Everything on 1 side, entirely red, and everything on other side entirely blue.

GAUSSIAN MIXTURE MODEL: EXAMPLE RUN

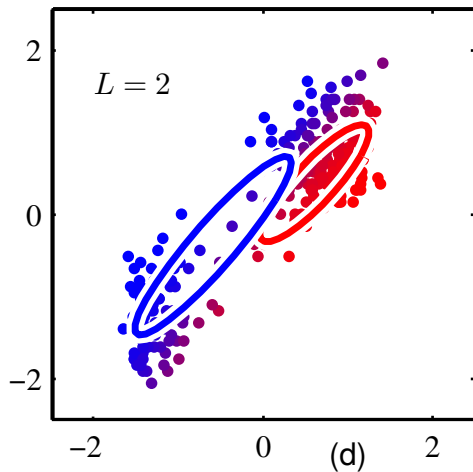


Iteration 1 (M-step)

Update the Gaussians

We calculate the empirical weighted mean and weighted covariance.

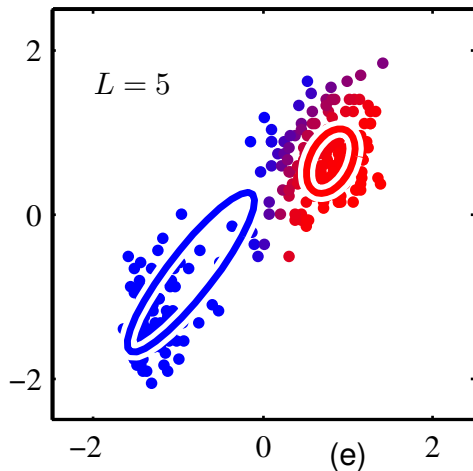
GAUSSIAN MIXTURE MODEL: EXAMPLE RUN



Iteration 2

Assign data to clusters
and update the Gaussians

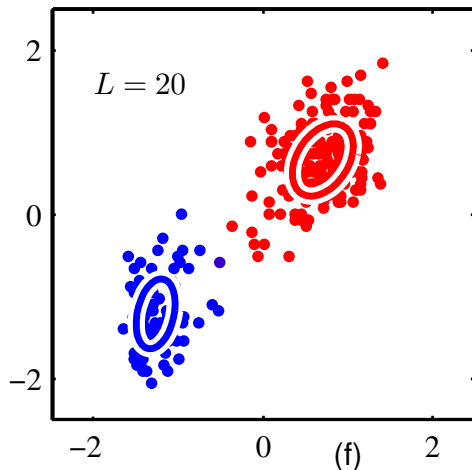
GAUSSIAN MIXTURE MODEL: EXAMPLE RUN



Iteration 5 (skipping ahead)

Assign data to clusters
and update the Gaussians

GAUSSIAN MIXTURE MODEL: EXAMPLE RUN



Iteration 20 (convergence)

Assign data to clusters
and update the Gaussians

[Turns out for this particular dataset,
we're very confident in our clustering.]

And we also get a probability of each of
the Gaussians. So, the Gaussian has a
certain prob. based on the empirical
probability of data coming from that cluster.

GMM AND THE BAYES CLASSIFIER

Intuitively they feel similar, but there are significant differences.

with k clusters

where we had a class-specific multi-variate Gaussian density

The GMM feels a lot like a K -class Bayes classifier, where the label of x_i is

$$\text{label}(x_i) = \arg \max_k \pi_k N(x_i | \mu_k, \Sigma_k).$$

class prior probability of data point coming from class k times the likelihood of that particular data point.

► π_k = class prior, and $N(\mu_k, \Sigma_k)$ = class-conditional density function.

Also, ► We learned π , μ and Σ using maximum likelihood here too. (And had similar looking updates for those parameters, except we didn't have an algorithm in that case. We just simply gave an equation for MLE for π, μ, Σ .)

For the Bayes classifier, we could find π , μ and Σ with a single equation because the class label was *known*. Compare with the GMM update: (where we have an algorithm.)

$$\pi_k = \frac{n_k}{n}, \quad \mu_k = \frac{1}{n_k} \sum_{i=1}^n \phi_i(k) x_i \quad \Sigma_k = \frac{1}{n_k} \sum_{i=1}^n \phi_i(k) (x_i - \mu_k)(x_i - \mu_k)^T$$

They're almost identical. But since $\phi_i(k)$ is changing we have to update these values. With the Bayes classifier, " ϕ_i " encodes the label, so it was known.

* But what's changed:

The key thing that's changed is that with Bayes classifier we already had the cluster assignments. So the clusters were like the classes and we had already assigned each of the data points to their specific classes. Because we had labels for them in our dataset. So we didn't need to ask which cluster/class any data point came from. We simply assigned it to its specific class.


GMM:

Now, with clustering, with this GMM which has a very similar form. We suddenly don't know for any data point which cluster it came from.

What that tells you:

→ Conditioned on the cluster assignment we can do MLE exactly. But when we don't have the cluster assignment, we have to iterate between updating the cluster assignment and then updating the model parameters.

We have to run an iterative algorithm because we can't optimize both simultaneously.

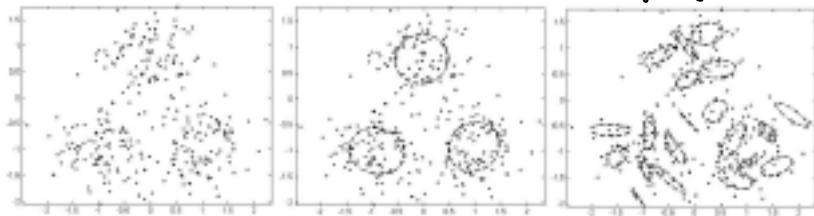
* * $\phi_i(x)$ is like encoding a probability of what label the i^{th} data point has, if you were to view a cluster like a label. For that data point we're inferring both the labels and class-specific densities simultaneously. 

Although, it's not typical to think of this sort of unsupervised mixture modelling as a classification problem. We don't even usually discuss the clusters as being classes but intuitively the math is very similar.

CHOOSING THE NUMBER OF CLUSTERS

no. of clusters you give it is the no. of clusters it learns (similarity with kmeans)

overfitting with GMM.



Maximum likelihood for the Gaussian mixture model can overfit the data. It will learn as many Gaussians as it's given.

There are a set of techniques for this based on the Dirichlet distribution.

A Dirichlet prior is used on π which encourages many Gaussians to disappear (i.e., not have any data assigned to them).



EM FOR A GENERIC MIXTURE MODEL

Generalized to mixture models using different distributions besides multi-variate Gaussians.

Algorithm: Maximum likelihood EM for mixture models

Given: Data x_1, \dots, x_n where $x \in \mathcal{X}$ ^{some parameter set} ^{i.i.d. distribution}

Goal: Maximize $\mathcal{L} = \sum_{i=1}^n \ln p(x_i | \pi, \theta)$, where $p(x | \theta_k)$ is problem-specific. ^{cluster specific likelihood. *} ^{mixing weights}

- Iterate until incremental improvement to \mathcal{L} is “small”

1. **E-step:** For $i = 1, \dots, n$, set ^{prior probability of a point coming from k^{th} cluster} $\phi_i(k) = \frac{\pi_k p(x_i | \theta_k)}{\sum_j \pi_j p(x_i | \theta_j)}$, for $k = 1, \dots, K$
^{conditional posterior probability of i^{th} observation coming from the k^{th} cluster.} ^{likelihood of i^{th} data point coming from the k^{th} cluster **}

^{similar from before:} $\left[\begin{array}{l} \text{sum the probabilities of all of the data coming from that} \\ \text{cluster. Get the expected no. of points we're gonna see from} \\ \text{that cluster.} \end{array} \right]$

2. **M-step:** For $k = 1, \dots, K$, define $n_k = \sum_{i=1}^n \phi_i(k)$ and set ^{prior prob. of any pt coming from k^{th} cluster.} $\pi_k = \frac{n_k}{n}$, ^{parameter to the distribution for k^{th} cluster.} $\theta_k = \arg \max_{\theta} \sum_{i=1}^n \phi_i(k) \ln p(x_i | \theta)$ ^{likelihood of i^{th} pt given θ . ** *}

(same exact thing)

Comment: Similar to generalization of the Bayes classifier for any $p(x | \theta_k)$.

* of x given the k^{th} cluster is going to be generic. With the Gaussian mixture model, this was a multivariate Gaussian. And μ_k was k^{th} mean and Σ_k covariance for that Gaussian.

Now this is any prob. distribution and Θ_k is the set of parameters for that distribution.

** This is now different from before. Before this was multivariate Gaussian, now it's whatever distribution we've defined on x ; according to the problem we're looking at.

*** So for example, before this was a Gaussian, where we wanted to maximize this thing over the mean and the covariance.

And we took the derivative and we found that we could do that in closed form.

Now, depending on what this distribution is, we might need to run an algorithm to update to maximize this, or maybe we can take the derivatives and set to zero and solve. But we'll have a different equation given the model that we're considering.

Given the densities that the likelihoods that we're using for the model we're considering.