

11. Mixture Models and EM

Mixture Models:

Recap K-means clustering

$\underline{x} \in \mathbb{R}^N$
data space

M clusters: q_1, q_2, \dots, q_M

M prototypes (centroids): $\underline{w}_q \in \mathbb{R}^N$

assignment variables for each point

$$w_q^{(\alpha)} := \begin{cases} 1 & \text{point assigned to cluster } q \\ 0 & \text{otherwise} \end{cases}$$

$$\sum_{q=1}^M w_q^{(\alpha)} = 1 \quad \text{a point has to be assigned to something.}$$

slide #4

measure "size" of a cluster post-hoc

↑
radius of the circles

cluster "size" does not tell us how "important" a cluster is.

high density vs. low density clusters

more ← important → less

Whatever is generating the data generates more points in cluster q_1 than q_2

We can measure this density post-hoc but what about an algorithm that does this as it's clustering?

↓
This algorithm should be able to capture densities

Remember density estimation?

Mixture Models

Combine (A) density estimation
with (B) clustering

Many ways of describing the same thing:

- cluster the data and estimate the density of each cluster.
- increase the resolution of your density estimation such that you end up with a mode around each cluster.
- assume M sources are generating points in the data find a measure for:

* 1. the probability of source $q \in \{1, \dots, M\}$ generating any point. Are you contributing to the data?

** 2. the probability of a point coming from source q_1 vs. q_2 vs. q_M who generated this point?

$$\underline{x} \in \mathbb{R}^N \sim P(\underline{x})$$

observations

joint unknown distribution of the data

Mixture models express $P(\underline{x})$ as a linear combination of components, ~~each~~

$$P(\underline{x}) = \sum_{q=1}^M \underbrace{P(\underline{x} | q)}_{**} \underbrace{P(q)}_*$$

slide #6 Example of Gaussian mixture models
choice of basis function

What happens now?

$P(\underline{x})$ is unknown, we want to estimate it.

~~Remember~~

Pick a model class: Gaussian Mixture model
parametric

Recap from density estimation:

slide #9

Objective: maximize likelihood $P(\underline{x}^{(a)}) \stackrel{!}{=} \max$

corresponds to minimizing neg. log likelihood

maximizing, minimizing
no real difference

depends on the cost funct.

any max can be turned into min.

easier to optimize
sums are easier to handle
than products.

$$ET = -\ln P =$$

$$= -\sum_x \ln \sum_q P(x^{(a)}|q)P(q) \stackrel{!}{=} \min$$

slide #10

We've actually seen a special case
of this before: soft-clustering

slides #15, #16, #17

Optimization

expressions for the different partial derivatives of all the parameters of the individual Gaussian components we are trying to find.

Lagrangian multiplier method required because of constraint $\sum_{q=1}^M P(q) = 1$

ensures we don't want useless bystanding components.

each component has to contribute something to the dataset. ~~Not necessarily to every point.~~

side effect

Not necessarily to every point.

degenerate solutions "collapse"

smaller M actually better. ← Actually we

slide #19

just want to avoid a component being only useful in explaining a very very small number of points.

The EM algorithm

Expectation Maximization

An iterative method that alternates between two steps for approximately maximizing the likelihood function

Recap MLE

Problem: estimate distribution parameters (e.g. μ, Σ of Gaussian mixture models)

Approach:

likelihood function $\hat{P}(\{\underline{x}^{(\alpha)}\} | \underline{w}) = \prod \hat{P}(\underline{x}^{(\alpha)} | \underline{w})$
↑ observations ↑ set of parameters \underline{w}
measures that the points come from this estimated distribution.

The true distribution $P(\underline{x})$ is unknown.

All we can do is find a $\hat{P}(\underline{x} | \underline{w})$ that produces non-zero probabilities for every observation $\underline{x}^{(\alpha)}$ and assigns more mass to more frequently occurring sample regions.

slide #9

Optimization

use the log to maximize sum instead of product.

$$\ln \hat{P}(\{\underline{x}^{(\alpha)}\} | \underline{w}) = \sum_{\alpha=1}^P \ln \hat{P}(\underline{x}^{(\alpha)} | \underline{w})$$

$\stackrel{!}{=} \max_{\underline{w}}$

OR minimize negative log likelihood

$$-\ln \hat{P}(\{\underline{x}^{(\alpha)}\} | \underline{w}) = -\sum_{\alpha=1}^P \ln \hat{P}(\underline{x}^{(\alpha)} | \underline{w})$$

A scenario to justify approximation of MLE via EM

MLE with latent variable models is too costly.

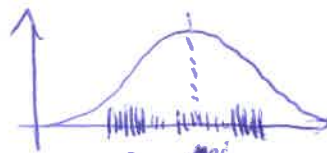
? → see example (e.g. density estimation for clustered data)

Assignment variables are unknown.

Example scenario for latent variable models:

What is the height distribution of students ~~in a school~~ at the university?

Approach 1: parametric fit ⇒ fit a Gaussian



latent variables

two groups

maybe gender,
a lot of basketball players

There are hidden causes in these observations.

We want to 1. increase the resolution of our density estimation

2. estimate density around each unknown group

mixture models



$$\hat{P}(\underline{x}^{(i)} | \underline{w})$$

see slide #25

we need clustering

slide #25

So what is the EM algorithm?

A solution to the chicken-or-egg problem

we need to know who is assigned
to which group to estimate the
density of that group

But we don't know the assignments.

~~To know the assignments we need a~~

But to get the assignment we need
the density of that group....

What do we do:

evenly
~~equally~~

[randomly assign samples to groups
estimate the density of each group

use the density to ~~as~~ "re-" assign
points to groups

use the current assignments to
update the densities

Sounds a lot like K-means.

K-mean is a special case of EM.

Now we need reassurance that alternating between the
steps improves both our assignments AND our density estimates.

slide 175

x

visible variables

m

latent (hidden) variables

w

parameters we want to estimate

slide 176

$$+ \ln P(\{ \underline{x}^{(\alpha)} \}, \{ \underline{m}^{(\alpha)} \} | \underline{w}) \stackrel{!}{=} \max$$

But w depends on m!
 and m is unknown →

How do ~~ALL~~ x, m and w interact with one another?

$$\begin{aligned} P(\underline{x} | \underline{w}) &= \prod_{\alpha} P(\underline{x}^{(\alpha)} | \underline{w}) = \\ &= \prod_{\alpha} \sum_{\underline{m}} P(\underline{x}^{(\alpha)}, \underline{m} | \underline{w}) \end{aligned}$$

marginalizing over m (summing out m) recovers the likelihood.

from the product rule: $P(a, b) = P(a|b)P(b)$

$$P(\underline{x}, \underline{m} | \underline{w}) = P(\underline{m} | \underline{x}, \underline{w}) \cdot P(\underline{x} | \underline{w})$$

$$P(\underline{x}, \underbrace{\underline{m}}_{\sim} | \underline{w}) = P(\underline{x} | \underline{m}, \underline{w}) \cdot P(\underline{m} | \underline{w})$$

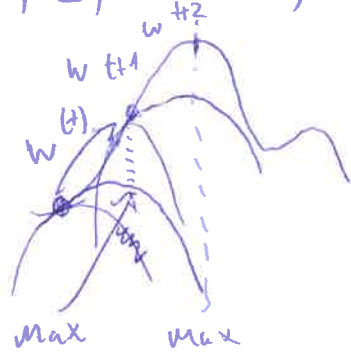
$$\begin{aligned} P(\underline{x} | \underline{w}) &= \sum_{\underline{m}} P(\underline{x}, \underline{m} | \underline{w}) \\ &= \sum_{\underline{m}} P(\underline{m} | \underline{w}) P(\underline{x} | \underline{m}, \underline{w}) \end{aligned}$$

posterior:

$$P(\underline{m} | \underline{x}, \underline{w}) = \frac{P(\underline{x}, \underline{m} | \underline{w})}{P(\underline{x} | \underline{w})}$$

posterior of hidden m given observed x

$$P(\underline{m} | \underline{x}, \underline{w}) P(\underline{x} | \underline{w}) = P(\underline{x} | \underline{m}, \underline{w}) P(\underline{m} | \underline{w})$$



show convergence plot

$$P(\underline{x} | \underline{w}) = \frac{P(\underline{x} | \underline{m}, \underline{w}) P(\underline{m} | \underline{w})}{P(\underline{m} | \underline{x}, \underline{w})}$$

Define a distribution $P(\underline{m})$ and insert it above:

$$P(\underline{x} | \underline{w}) = \frac{P(\underline{x} | \underline{m}, \underline{w}) P(\underline{m} | \underline{w})}{P(\underline{m})} \cdot \frac{P(\underline{m})}{P(\underline{m} | \underline{x}, \underline{w})}$$

take logarithm on both sides:

$$\ln P(\underline{x} | \underline{w}) = \ln \left(\frac{P(\underline{x} | \underline{m}, \underline{w}) P(\underline{m} | \underline{w})}{P(\underline{m})} \right) + \ln \left(\frac{P(\underline{m})}{P(\underline{m} | \underline{x}, \underline{w})} \right)$$

Compute the expectation w.r.t. $P(\underline{m})$:

$$\mathbb{E}[\ln P(\underline{x} | \underline{w})] = \int P(\underline{m}) \ln \left(\frac{P(\underline{x} | \underline{m}, \underline{w}) P(\underline{m} | \underline{w})}{P(\underline{m})} \right) d\underline{m} + \int P(\underline{m}) \ln \left(\frac{P(\underline{m})}{P(\underline{m} | \underline{x}, \underline{w})} \right) d\underline{m}$$

$\leftarrow D_{KL}(P(\underline{m}) || \frac{P(\underline{m})}{P(\underline{m} | \underline{x}, \underline{w})}) \geq 0$

Jensen's inequality

For convex functions $E[f(x)] \geq f(E[x])$

for concave $E[f(x)] \leq f(E[x])$

The logarithm function is concave

1st derivative of $\ln(a) = \frac{1}{a}$ strictly decreases when x increases

$$E[\ln P(\underline{x}|\underline{w})] = E[\mathcal{L}(P(\underline{m}), \underline{w})] + D_{KL} \left(\begin{matrix} P(\underline{m}) \\ P(\underline{m}|\underline{x}, \underline{w}) \end{matrix} \right)$$

$\ln P(\underline{x}|\underline{w})$ has a lower bound that is

~~$$\int P(\underline{m}) \ln P(\underline{m}) d\underline{m}$$~~

$$\mathcal{L}(P(\underline{m}), \underline{w}) = \int P(\underline{m}) \ln \left(\frac{P(\underline{x}|\underline{m}, \underline{w}) P(\underline{m}|\underline{w})}{P(\underline{m})} \right) d\underline{m}$$

E-step using $\underline{w}^{\text{old}}$:

Maximize $\mathcal{L}(P(\underline{m}), \underline{w}^{\text{old}})$ w.r.t $P(\underline{m})$
keeping $\underline{w}^{\text{old}}$ fixed

↓ solution:

$$\mathcal{L}(P(\underline{m}), \underline{w}^{\text{old}}) \stackrel{!}{=} \max$$

when $D_{KL} \rightarrow 0$

this happens when $P(\underline{m}) = P(\underline{m}|\underline{x}, \underline{w}^{\text{old}})$

M-step

keep $P(\underline{m})$ fixed

maximize $\mathcal{L}(P(\underline{m}), \underline{w})$ w.r.t $\underline{w} \rightarrow \underline{w}^{\text{new}}$

We either have a higher lower bound \mathcal{L} or the same because we already reached the maximum.

increasing $\mathcal{L} \rightarrow$ increase of likelihood yay!
but not by the same amount.

the difference is explained by the D_{KL} term. ≥ 0
increase in likelihood is actually larger.

Int ~~substitution~~

~~P(m)~~

if we are keeping $P(\underline{m})$ fixed between
E and M steps we are

actually measuring $P(\underline{m})$ as a function of $\underline{w}^{\text{old}}$.

Therefore:

$$P(\underline{m}) := P(\underline{m} | \underline{x}, \underline{w}^{\text{old}})$$

From this follows:

~~L(m)~~

$$\mathcal{L}(P(\underline{m} | \underline{x}, \underline{w})^{\text{old}}, \underline{w}) =$$

$$= \sum_{\underline{m}} P(\underline{m} | \underline{x}, \underline{w}^{\text{old}}) \ln P(\underline{x} | \underline{m}, \underline{w})$$

no old for this one.

$$- \sum_{\underline{m}} P(\underline{m} | \underline{x}, \underline{w}^{\text{old}}) \ln P(\underline{m} | \underline{x}, \underline{w}^{\text{old}})$$

entropy of $P(\underline{m} | \underline{x}, \underline{w}^{\text{old}})$

independent of $\underline{w} \rightarrow \text{constant}$.

$$) := \cancel{Q}(\underline{w}, \underline{w}^{\text{old}})$$

the ~~step~~ are