

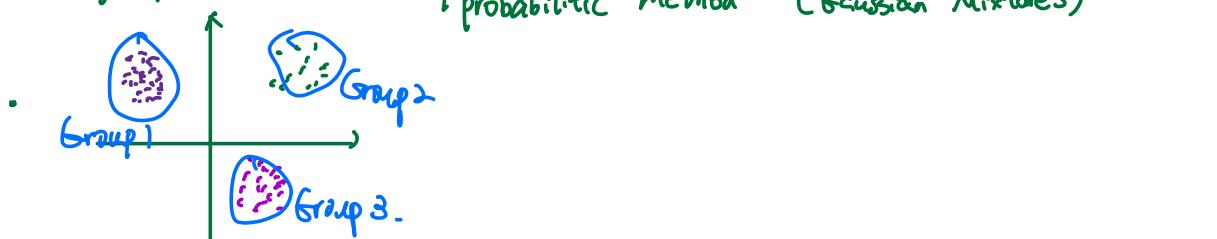
Chapter 4 Clustering, Mixture model and Expectation Maximization

Lecture 1

1. What is clustering?
2. K-means problem and cost function
3. K-means algorithm
4. Guarantees, convergence and local optima
5. How to pick K?

• What is clustering?

- Unsupervised. Often used in exploratory data analysis.
- Many approaches exists.
 - { deterministic method (K-means)
 - probabilistic method (Gaussian Mixtures)



• K-means Problem statement + Cost function

Given:
 $X_{1:N}$: a dataset $\begin{bmatrix} X_1^T \\ \vdots \\ X_N^T \end{bmatrix}_{NxD}$ $X = []_{Dx1}$
 K : designed number of clusters, must be ≥ 1

Goal:

Find locations of K cluster centers $M_{1:K} \rightarrow []_{Dx1}$ $M = \begin{bmatrix} M_1^T \\ \vdots \\ M_K^T \end{bmatrix}_{KxD}$
Such that if we compute the distance from each data point x_n to its nearest cluster center, and sum up the total distance - this total distance is minimized.

Cost function:

$$\begin{aligned} J'(X_{1:N}, M_{1:K}) &= \sum_{n=1}^N \min_{k \in \{1, 2, \dots, K\}} \text{dist}(X_n, M_k) \\ &= \sum_{n=1}^N \min_{k \in \{1, 2, \dots, K\}} \sum_{d=1}^D (X_{nd} - M_{kd})^2 \end{aligned}$$

$$= \sum_{n=1}^N \min_{k \in \{1, 2, \dots, K\}} (x_n - \mu_k)^T (x_n - \mu_k)$$

- Equivalent Formulation with Assignment Variables

Let's indicate which cluster (indexed by $k \in \{1, 2, \dots, k\}$) is closest to data point n with new variable

r_n is a one-hot indicator vector of size k $r = []_{n \times k}$

Example: If $k=4$, the possible r_n values are:

x_n 's closest cluster

4
3
2
1

$$r_n = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

The cost function can be rewritten as:

$$\begin{aligned} J(x_{1:N}, \mu_{1:k}, r_{1:N}) &= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \text{dist}(x_n, \mu_k) \\ &= \sum_{n=1}^N \sum_{k=1}^K r_{nk} (x_n - \mu_k)^T (x_n - \mu_k) \end{aligned}$$

Previous: $= \sum_{n=1}^N \min_{k \in \{1, 2, \dots, K\}} (x_n - \mu_k)^T (x_n - \mu_k)$

等价是相同的，先显式地写出来！

K-Means Coordinate Descent Algorithm

- Goal: Find cluster locations $\mu_{1:k}$ and assignments $r_{1:N}$ that minimize the cost
- Opt problem: $\mu_{1:k}^*, r_{1:N}^* = \operatorname{arg\min} J(x_{1:N}, \mu_{1:k}, r_{1:N})$
- How to solve?
 - Gradient descent? No GD needs all continuous Variable
 - Coordinate descent? Yes!

Iteratively update some vars while holding other fixed:

Each iteration has 2 steps:

Step 1. Fix $\mu = \mu^{t-1}$. Update r

$$r^t = \underset{r}{\operatorname{arg\,min}} J(X_{1:N}, \mu^{t-1}, r) \rightarrow \text{Assignment step}$$

Step 2. Fix $r = r^t$. Update μ

$$\mu^t = \underset{\mu}{\operatorname{arg\,min}} J(X_{1:N}, \mu, r^t) \rightarrow \text{Cluster location step}$$

Assignment Update Step

- Aim: $\underset{r_{1:N}}{\operatorname{arg\,min}} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \operatorname{dist}(x_n, \mu_k)$

- Notice: No interaction between r_n terms. Each just in one term of sum.

Can solve for each r_n value separately:

$$= \min_{r_1} f(x_1, r_1) + \min_{r_2} f(x_2, r_2) + \dots + \min_{r_N} f(x_N, r_N)$$

$$\Rightarrow r_n^* = \underset{r_n \in \text{onehot}[K]}{\operatorname{arg\,min}} \sum_k r_{nk} \operatorname{dist}(x_n, \mu_k)$$

$$= \text{One hot indicator of } k^*, k^* = \underset{k \in \{1, 2, \dots, K\}}{\operatorname{arg\,min}} \operatorname{dist}(x_n, \mu_k)$$

Cluster Location Update Step:

Aim: $\sum_{n=1}^N \sum_{k=1}^K r_{nk} (x_n - \mu_k)^T (x_n - \mu_k)^T$

$$\Rightarrow \sum_{k=1}^K \sum_{n=1}^N r_{nk} (x_n - \mu_k)^T (x_n - \mu_k)^T$$

For each k , we want to solve:

$$\underset{\mu \in \mathbb{R}^D}{\operatorname{arg\,min}} \sum_n r_{nk} (x_n - \mu_k)^T (x_n - \mu_k)$$

Take the second derivative:

$$\nabla_{\mu_k} \ell(\mu_k) = 0$$

set μ_k to mean location of

$$\mu_k^* = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}} \quad \checkmark \text{ all data assigned to cluster } k$$

K-Means Algorithm

Input: $X_{1:N}$: dataset

K : num clusters

K : num clusters

$\mu_{1:k}^0$: initial locations
can be randomly drawn
can be picked manually
just need to be distinct

can be randomly drawn
can be picked manually
just need to be distinct

Procedure:

for iteration t in $1, 2, \dots$ until converged

step 1
update
 $r_{1:N}$

for n in 1,2,...N:

Update $r_n^t \leftarrow \text{to_onehot}(k^*)$, $k^* = \underset{k \in 1-K}{\operatorname{argmin}} \text{dist}(x_n, \mu_k)$

step 2
update μ_{IK}

for k in 1, 2, ..., K:

$$\text{Update } \mu_k^t \leftarrow \frac{\sum r_{nk}^t x_n}{\sum r_{nk}^t}$$

Return: $\mu_{1:k}^*$ optimal cluster locations
 $r_{1:N}^*$ optimal assignment indicators

Guarantees + Assessing Convergence

key idea of coordinate descent: Each step is guaranteed to improve the cost function
(or at convergence, will stay the same)

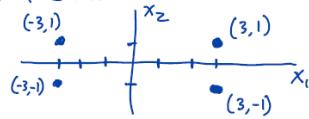
$$\dots \geq J(x_{1:N}, \mu^T, r^T)$$

after iter T
step 2

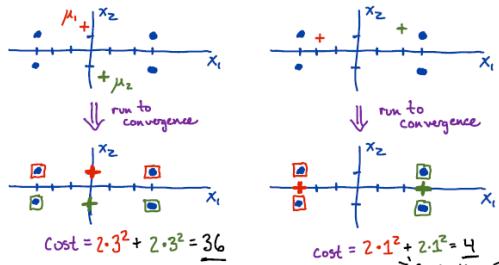
Issue : local optima (try kmeans++)

Issue : Local Optima

Consider $N=4$ dataset



Two possible initializations



For a given K , if you run K-means twice, might get arbitrarily worse solution from one vs. other.

To be robust, try many initializations and pick best solution

How to select number of cluster k ?

How to select number of clusters K ?

Not easy to do with training set alone.

Can always improve cost function by adding more clusters up to limit $K=N$, where each data point assigned to own cluster. Will have cost = 0!

Two options

(1) Pick K that minimizes holdout set cost

$$K^* = \underset{K \in \{2, 3, \dots\}}{\operatorname{argmin}} J(x_{1:N}^{\text{val}}, \mu_{1:K})$$

Might favor very large K

$$\underset{\mu_{1:K}}{\operatorname{argmin}} J(x_{1:N}^{\text{tr}}, \mu_{1:K})$$

(2) Add a "complexity penalty"

$$K^* = \underset{K \in \{2, 3, \dots\}}{\operatorname{argmin}} J(x_{1:N}, \mu_{1:K}) + \lambda \text{penalty}(K)$$

now need to pick $\lambda > 0$ and penalty function

Lecture 2. Gaussian Mixture Models

(1) Why mixture models

(2) Gaussian mixture model two views: without / with hidden assignments

(3) Estimating posterior over assignment indicators for one example

(4) Estimating parameters of GMM given dataset (logSumExp trick)

(5) Problems with ML estimation for GMMs

Why mixture models?

Why Mixture Models?

Common data analysis problem:

We observe many examples of a diverse population

Examples: many patients in hospital
have vital signs taken

many DNA sequences obtained
from cells found near a tumor

We have good domain knowledge that examples come from some
"types" or "groups" or "clusters"

Examples: patients: disease A vs. disease B vs ... disease Z
cells: healthy vs. cancer vs. cancer w/ extra mutation

But, we just observe raw data, not cluster indicators.

Goal: Represent many clusters to get flexible model
that fits data well.

Mixture model idea

Observed Data

cluster 1

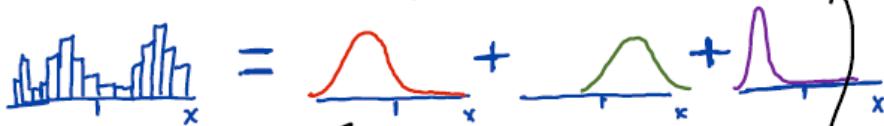
45 %

cluster 2

40 %

cluster 3

15 %

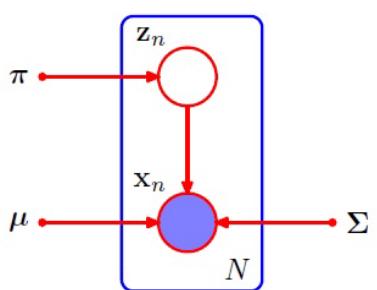


Each "cluster" has two parts:

cluster-specific
PDF over data

appearance probability
in overall population

Gaussian mixture models (GMM) Graphical model



Parameters:

$$z = [z_1, \dots, z_k]$$

$$\pi = [\pi_1, \pi_2, \dots, \pi_k] \quad \pi_k \in \Delta^k$$

$$\mu = [\mu_1, \mu_2, \dots, \mu_k] \quad \mu_k \in \mathbb{R}^D$$

$$\Sigma = [\Sigma_1, \dots, \Sigma_k], \quad \Sigma_k \in \mathbb{R}^{D \times D}$$

Aim: Define a distribution over real valued vectors via a mixture of k separate Gaussian distributions.

$$z = [z_1, z_2, \dots, z_k] \rightarrow \text{one-hot vector size } k$$

$$z_k = \begin{cases} 1 & \text{if we assign cluster } k \text{ to generate our data } x \\ 0 & \text{otherwise} \end{cases}$$

$z \sim \text{Cat}(\pi_{1:k})$: assign to cluster k with prob π_k

$$p(z) = \prod_{k=1}^K \pi_k^{z_k} \quad \pi_k \in \Delta^k$$

The GMM model can be defined as:

$$\begin{aligned} p(x, z) &= p(z) \cdot p(x|z) \\ &= \text{Cat PMF}(z|\pi_{1:k}) \cdot \prod_{k=1}^K \text{NMPDF}(x|\mu_k, \Sigma_k)^{z_k} \\ &= \prod_{k=1}^K \pi_k^{z_k} \cdot N(x|\mu_k, \Sigma_k)^{z_k} \end{aligned}$$

Computing posterior over assignments

Aim: Given a GMM with known parameters π, μ, Σ and a single data example $x \in \mathbb{R}^D$, which of the k clusters generated this x ?

Can write this as a posterior:

$$p(z_k=1 | x) = \frac{p(x|z_k=1) \cdot p(z_k=1)}{p(x)} = \frac{p(x|z_k=1) \cdot p(z_k=1)}{\sum_z p(z) p(x|z)}$$

$$= \frac{N(x|\mu_k, \Sigma_k) \cdot \pi_k}{\sum_{l=1}^k \pi_l N(x|\mu_l, \Sigma_l)} \quad \text{GMM PDF}$$

Thus, can write posterior as a categorical distribution:

$$q_{LB}(x) = \text{Cat PMF } [r_1, r_2, \dots, r_k]$$

$$\text{where } r_k = \frac{1}{\text{GMM PDF}(x)} \pi_k \cdot N(x|\mu_k, \Sigma_k)$$

ML Estimation of Parameters

Given: N data examples $X_{1:N}$ s.t. $X_n \in \mathbb{R}^D$

k : number of assumed cluster

Goal: Estimate value of $6mn$ parameters:

$$k \pi$$

$$k \times D \mu_{ik}$$

$$k \times D \times D \Sigma_{ik}$$

Method: Maximize likelihood of N observed examples $X_{1:N}$

★ LogsumExp trick to avoid over/underflow

$$\pi^*, \mu^*, \Sigma^* = \arg \max_{\pi, \mu, \Sigma} \frac{\sum_{i=1}^N \log \text{GMM PDF}(X_i | \pi, \mu, \Sigma)}{N}$$

$$= \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k N(X_i | \mu_k, \Sigma_k) \right)$$

$$= \text{logsumExp} \left([\log \pi_1, \dots, \log \pi_K] + \right. \\ \left. [\log N(X_i | \mu_1, \Sigma_1) + \dots + \log N(X_i | \mu_K, \Sigma_K)] \right)$$

More:

`scipy.special.logsumexp`

$$\log \sum_k e^{a_k}$$

$$m = \max_k a_k$$

$$\log \left[e^m \left(\sum_k e^{a_k - m} \right) \right]$$

$$m + \log \sum_k e^{a_k - m} \leq 0$$

Example:

logsum

$$\log \left(\underbrace{e^{1000}}_{\approx 1000 + \epsilon} + \underbrace{e^0}_{\text{np.exp}(0)} + \underbrace{e^0}_{\text{np.exp}(0)} \right)$$

$$\approx 1000 + \epsilon$$

$$\text{np.log} \left(\underbrace{\text{np.exp}(1000)}_{+\infty} + \underbrace{\text{np.exp}(0)}_{\text{np.exp}(0)} + \underbrace{\text{np.exp}(0)}_{\text{np.exp}(0)} \right)$$

$$\log \left(e^{1000} \left(e^0 + e^{-1000} + e^{-1000} \right) \right)$$

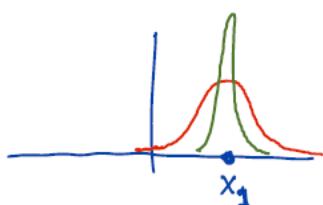
$$\log(e^{1000}) + \log \left(\underbrace{e^0 + e^{-1000} + e^{-1000}}_{\approx 1} \right)$$

Problem with ML Estimation.

Problems with ML Estimation

As we've seen before, ML estimation can have problems with poor generalization performance especially when data set size is small.

Consider $K=1$ GMM (basically, a single Gaussian) for $N=1$ example with $D=1$ dimension. $\Sigma_k = \sigma_k^2$ scalar variance



ML goal:

$$\max_{\mu_1 \in \mathbb{R}, \sigma_1^2 > 0} \log \text{NormPDF}(x_1 | \mu_1, \sigma_1^2)$$

Two solutions shown

$$\mu_1 = x_1, \sigma_1^2 = 1$$

$$\mu_1 = x_1, \sigma_1^2 = \frac{1}{100}$$

Problem: We can make this point's PDF value $\rightarrow +\infty$ as $\sigma_1^2 \rightarrow 0^+$

Consider more general GMMs:

Problems arise when $N > 1$ and $K \geq 2$,

when one cluster with $\sigma_k^2 \rightarrow 0$ lands on top of one point, and other clusters explain other examples



Problem summary: ML can yield $-\infty$ PDF

- Want all clusters to have non-zero variance to have hope of generalizing well.
- Want finite PDF values even in training set for sensible model selection.

Ways to Avoid this Problem

(1) Constrain variances away from zero

$$\sigma^2 > \varepsilon \text{ not } \sigma^2 > 0$$

for some $\varepsilon >> 0$, say 0.01

(2) Penalized ML

$$\max \log p(x_{1:N}) - \lambda \text{penalty}(\sigma)$$

Add penalty term to objective
with weight $\lambda > 0$

Penalty term gives high "cost" to
problematic σ values

(3) Do MAP instead of ML

Pick a prior distribution $p(\sigma)$
that favors variances far from zero

$$\max \log p(x_{1:N}) + \log p(\sigma)$$

Figure 9.7 Illustration of how singularities in the likelihood function arise with mixtures of Gaussians. This should be compared with the case of a single Gaussian shown in Figure 1.14 for which no singularities arise.

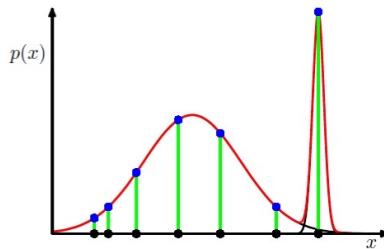
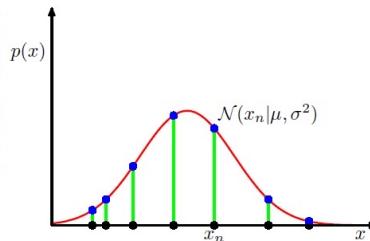


Figure 1.14 Illustration of the likelihood function for a Gaussian distribution, shown by the red curve. Here the black points denote a data set of values $\{x_n\}$, and the likelihood function given by (1.53) corresponds to the product of the blue values. Maximizing the likelihood involves adjusting the mean and variance of the Gaussian so as to maximize this product.



Lecture 3. EM algorithm and gradient Descent

1. Penalized ML optimization problem for GMMs
2. Gradient Descent for GMMs
3. Derivation of Coordinate Descent for GMMs
4. The EM Coordinate Descent Algorithm for GMMs

Penalized ML optimization problem for GMMs

- Parameters:

$\pi \in \mathbb{R}^k$: k -length vector with non-negative entries that sum to one.

$\mu \in \mathbb{R}^D$: μ_k is D -length vector

$\Sigma \in \mathbb{R}^{D \times D}$: Σ_k is $D \times D$ symmetric, PD matrix

- Object:

$$\begin{aligned}\pi^*, \mu^*, \Sigma^* &= \underset{\pi, \mu, \Sigma}{\operatorname{argmax}} \sum_{i=1}^N \log \text{GMM PDF}(x_i | \pi, \mu, \Sigma) \\ &= \underset{\pi, \mu, \Sigma}{\operatorname{argmin}} - \sum_{i=1}^N \log \text{GMM PDF}(x_i | \pi, \mu, \Sigma)\end{aligned}$$

- Object with penalty:

Goal: Do ML, but avoid pathology where one cluster has variance shrink to 0 to get "infinite" likelihood

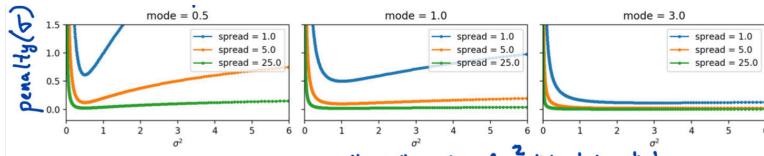
Instead of "full" $D \times D$ covariance, let's assume $\Sigma_k = \begin{bmatrix} \sigma_{k1} & & \\ & \sigma_{k2}^2 & \\ & & \ddots & \\ & & & \sigma_{kD}^2 \end{bmatrix}$

Then we can define penalty over each dimension's

Std. deviation scalar:

$$\text{penalty}(\sigma) = \frac{1}{m^2 s} \log \sigma + \frac{1}{2} \frac{1}{m \cdot s} \frac{1}{\sigma^2}$$

- Visualization and interpretation



$M > 0$ is the "mode", value of σ^2 that is least penalized
 $S > 0$ is the "spread", high values \rightarrow larger region around mode favored
 Note: σ^2 near zero is always highly penalized

Can view this penalty as closely related to a Gamma distribution prior on the variance. Connects penalized ML to MAP estimation

We emphasize this is one possible penalty.
 The diagonal covariance assumption is for simplicity.

- Object after penalty:

$$\pi^*, \mu^*, \Sigma^* = \arg \min_{\pi, \mu, \Sigma} - \sum_{i=1}^n \log \text{GMMPOF}(x_i | \pi, \mu, \Sigma) + \lambda \sum_{k=1}^K \sum_{d=1}^D \text{Penalty}(\theta_k)$$

- use gradient descent:

Sketch: Input: initial values π^0, μ^0, Σ^0

Procedure: for iter t in 1, 2, ... until converged:

$$\begin{aligned}\pi^t &\leftarrow \pi^{t-1} - \varepsilon \nabla_{\pi} \ell(\pi^{t-1}, \mu^{t-1}, \Sigma^{t-1}) \\ \mu^t &\leftarrow \mu^{t-1} - \varepsilon \nabla_{\mu} \ell(\pi^{t-1}, \mu^{t-1}, \Sigma^{t-1}) \\ \Sigma^t &\leftarrow \Sigma^{t-1} - \varepsilon \nabla_{\Sigma} \ell(\pi^{t-1}, \mu^{t-1}, \Sigma^{t-1})\end{aligned}$$

△ Problem:

Both π and Σ are constrained

π must sum to one and be non-negative

Σ must be symmetric and PD

The GD update may produce values that violate constraints

$$\pi^t \leftarrow \pi^{t-1} - \varepsilon \cdot g^{t-1} \text{ will not keep } \pi^t \in \Delta^K$$

- Solution: Reparameterize! Find unconstrained parameterization

that maps to our constraint version

For π :

Let $\pi = [\pi_1, \pi_2, \dots, \pi_k]$ be a "probability vector". That means π defines a discrete PMF over k choices. Mathematically, π must have non-negative entries that sum to one.

Define transform function $\pi(\cdot) : \mathbb{R}^k \rightarrow \Delta^k$ (softmax)

$$\pi([s_1, \dots, s_k]) = \left[\frac{e^{s_1}}{Z}, \frac{e^{s_2}}{Z}, \dots, \frac{e^{s_k}}{Z} \right]$$

$$\text{where } Z = \sum_{i=1}^k e^{s_i}$$

For δ

Let δ_k be a vector that must be positive

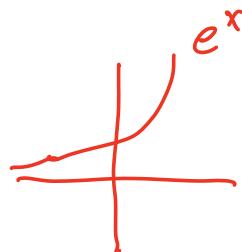
$$\delta_k = [\delta_{k1}, \dots, \delta_{kD}]$$

Option 1:

Define $t_k \in \mathbb{R}^D$, and use transformation $\delta(\cdot)$

$$\delta(t_k) \triangleq [e^{t_{k1}}, e^{t_{k2}}, \dots, e^{t_{kD}}]$$

by definition, for any input t , $\delta(t)$ should have positive entries.



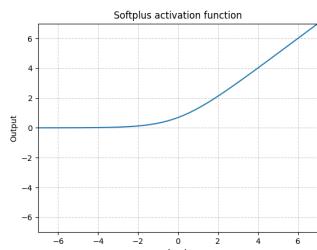
Option 2:

Softplus function

$$\text{softplus}(x) = \log(1 + e^x)$$

$$\delta(t_k) \triangleq [\text{softplus}(t_{k1}), \text{softplus}(t_{k2}), \dots, \text{softplus}(t_{kD})]$$

Option 2 might be better in GD, because less sensitive to small changes in input t .



• Summary for GD for GMMs

Gradient Descent for GMMs.

Parameters : $s: s_{1:k}, s_k \in \mathbb{R}$ $\mu: \mu_{1:k}, \mu_k \in \mathbb{R}^D$ $t: t_{1:k}, t_k \in \mathbb{R}^D$ all unconstrained real values

Goal: $s^*, \mu^*, t^* = \underset{\substack{s \in \mathbb{R} \\ \mu \in \mathbb{R}^D \\ t \in \mathbb{R}^D}}{\operatorname{argmin}} \alpha^{\text{PML}}(s, \mu, t)$
 then set $\pi^* \leftarrow \pi(s^*), \sigma^* \leftarrow \sigma(t^*)$. π^*, μ^*, σ^* Return

$$\text{Loss: } \alpha^{\text{PML}}(s, \mu, t) = \alpha^{\text{PML}}(\pi(s), \mu, \sigma(t))$$

$$\begin{aligned} \text{Grad: } \nabla_s \alpha^{\text{PML}} &= \nabla_s \pi(s) \cdot \nabla_\pi \alpha^{\text{PML}}(\pi(s), \dots) \\ \nabla_t \alpha^{\text{PML}} &= \nabla_t \sigma(t) \cdot \nabla_\sigma \alpha^{\text{PML}}(\dots, \sigma(t)) \end{aligned} \quad \text{chain rule!}$$

GD is now possible!

Can do gradients by hand or via automatic differentiation

Keep in Mind:

- still need to select step size $\epsilon > 0$ (as in any GD)
- loss has many local optima. May need many runs of GD to find a decent solution

Coordinate descent method iteration + derivation

- Objective:

$$\pi^*, \mu^*, \Sigma^* = \arg \min - \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \text{MVNPDF}(x_n | \mu_k, \Sigma_k)$$

- Idea:

Analytically consider an optimal set of parameters π^*, μ^*, Σ^* .

To be optimal, they must satisfy these three equations.

$$\begin{aligned} \nabla_{\pi_k} L^{ML}(\pi^*, \mu^*, \Sigma^*) &= \vec{0}_K \\ \nabla_{\mu_k} L^{ML}(\pi^*, \mu^*, \Sigma^*) &= \vec{0}_D \\ \nabla_{\Sigma_k} L^{ML}(\pi^*, \mu^*, \Sigma^*) &= \vec{0}_{D \times D} \end{aligned}$$

expand out

$$\begin{aligned} \nabla_{\mu_k} L^{ML} &= \sum_{n=1}^N \nabla_{\mu_k} [\log \text{GMM PDF}(x_n | \pi, \mu, \Sigma)] \\ &= \sum_{n=1}^N \frac{1}{\text{GMM PDF}(x_n)} \cdot \nabla_{\mu_k} \text{GMM PDF}(x_n | \pi, \mu, \Sigma) \\ &= \sum_{n=1}^N \frac{1}{\text{GMM PDF}(x_n)} \cdot \nabla_{\mu_k} \left[\sum_{l=1}^K \pi_l C(\Sigma_k) \cdot e^{\frac{1}{2} (x_n - \mu_k)^T \Sigma^{-1} (x_n - \mu_k)} \right] \\ &= \sum_{n=1}^N \frac{1}{\text{GMM PDF}(x_n)} \underbrace{\pi_k C(\Sigma_k) \cdot \nabla_{\mu_k} e^{\frac{1}{2} (x_n - \mu_k)^T \Sigma^{-1} (x_n - \mu_k)}}_{P(Z_{nk}=1 | x_n)} \\ &= \sum_{n=1}^N \frac{1}{\text{GMM PDF}(x_n)} \underbrace{P(Z_{nk}=1 | x_n)}_{\pi_k C(\Sigma_k) \cdot e^{\frac{1}{2} (x_n - \mu_k)^T \Sigma^{-1} (x_n - \mu_k)}} \cdot \nabla_{\mu_k} \left[\sum_{l=1}^K \pi_l C(\Sigma_k) \cdot e^{\frac{1}{2} (x_n - \mu_k)^T \Sigma^{-1} (x_n - \mu_k)} \right] \end{aligned}$$

Define $\gamma_{nk} = \frac{P(Z_{nk}=1 | x_n)}{P(x_n)} = P(Z_{nk}=1 | x_n) \rightarrow \text{posterior}$

$$= \sum_{n=1}^N \gamma_{nk} - \Sigma_k^{-1} (x_n - \mu_k)$$

Insight: Can simplify use soft "responsibilities" $\gamma_{nk} \in \Delta^k$
 Fix these, and become easy to solve

$$\sum_{n=1}^N \gamma_{nk} \sum_k (\mathbf{x}_n - \mathbf{m}_k) = 0$$

$$\sum_k \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \mathbf{m}_k) = 0$$

$$\boxed{\mathbf{m}_k^* = \frac{\sum_n \gamma_{nk} \mathbf{x}_n}{\sum_n \gamma_{nk}}}$$

Interpretation: weighted empirical mean
 of data assigned to cluster k

Similarly:

$$\pi_k^* = \frac{1}{N} \sum_{n=1}^N \gamma_{nk} \quad \text{(PRML 9.22)} \quad \begin{matrix} \text{fraction of all data points} \\ \text{assigned to cluster k} \end{matrix}$$

$$\Sigma_k^* = \frac{1}{\sum_n \gamma_{nk}} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \mathbf{m}_k^*) (\mathbf{x}_n - \mathbf{m}_k^*)^T \quad \begin{matrix} \text{(9.19)} \\ \text{weighted empirical covariance of} \\ \text{data assigned to cluster k} \end{matrix}$$

EM algorithm for GMM

- Algorithm for estimating the parameters of a GMM using (penalized) maximum likelihood

Coordinate Descent algorithm for GMMs

Similar to K Means

Minimizes the negative ML loss function \mathcal{L}^{ML}

Known as Expectation-Maximization or "E-M"
for reasons discussed in next class

Input: $X_{1:N}$ dataset
 $\pi^0, \mu_{1:k}^0, \Sigma_{1:k}^0$ initial parameters

Procedure: for iter $t \in 1, 2, \dots$ until converged:

E step
update responsibilities

for example $n \in 1, 2, \dots N$:
for cluster $k \in 1, 2, \dots K$:
 $\tilde{\gamma}_{nk} \leftarrow \pi_k^{t+1} \text{MVNormPDF}(x_n | \mu_k^{t+1}, \Sigma_k^{t+1})$
 $\gamma_n \leftarrow \left[\frac{\tilde{\gamma}_{n1}}{\sum_k \tilde{\gamma}_{nk}}, \frac{\tilde{\gamma}_{n2}}{\sum_k \tilde{\gamma}_{nk}}, \dots, \frac{\tilde{\gamma}_{nK}}{\sum_k \tilde{\gamma}_{nk}} \right]$

M step
update parameters

for cluster $k \in 1, 2, \dots K$:
 $\pi_k^t \leftarrow \frac{1}{N} \sum_n \tilde{\gamma}_{nk}$
 $\mu_k^t \leftarrow \frac{1}{\sum_n \tilde{\gamma}_{nk}} \sum_n \tilde{\gamma}_{nk} x_n$
 $\Sigma_k^t \leftarrow \frac{1}{\sum_n \tilde{\gamma}_{nk}} \sum_n \tilde{\gamma}_{nk} (x_n - \mu_k^t)^T (x_n - \mu_k^t)$

could do with diagonal parameters
 $\Sigma_k = [\sigma_1^2 \dots \sigma_K^2]$
can adjust to penalty term

Will converge to a fixed point where π^*, μ^*, Σ^* is local optima of \mathcal{L}^{ML}

» Lecture 4 A new view of EM
 (as a principle, general purpose optimization algorithm)

Topics: Recap: EM for GMMs
GMM as latent variable model

Idea 1: Complete likelihood easier than incomplete

Idea 2: Expectation of complete likelihood also easy

Idea 3: Can formulate objective function that:

- Principled: \geq a lower bound of incomplete likelihood

- Tractable: Uses expectation of complete likelihood

Idea 4: EM is coordinate ascent optimization applied to this objective

• Recap EM for GMMs

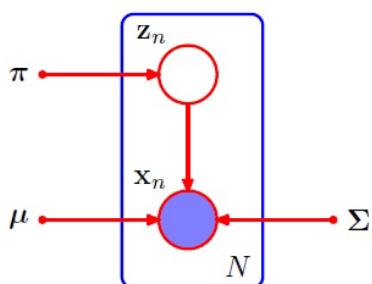
Big idea:

Coordinate ascent algorithm: each substep updates subset of variables (coordinates)

Open question:

- What principles let us use EM?
- What objectives are we optimizing?

• Step back: Latent Variable Models



z_n is an one-hot indicates, cluster assigned to example n $z_n \in \text{onehot}(k)$

z_n is latent or "hidden", we cannot observe it directly.

$$p(z_n) = \text{Cat}(\pi_1, \pi_2, \dots, \pi_k) = \prod_{k=1}^K \pi_k^{z_{nk}}$$

$$p(x_n | z_n) = \prod_{k=1}^K \text{NormPDF}(x_{nk} | \mu_k, \sigma_k^2)^{z_{nk}}$$

(z_n is a one-hot, so only one term in each product will be used)

$$\text{Joint: } P(X_n, Z_n) = P(Z_n) \cdot P(X_n | Z_n)$$

$$\text{Posterior: } P(Z_n | X_n) = \frac{P(X_n, Z_n)}{P(X_n)} = \frac{P(X_n, Z_n)}{\sum_{Z_n} P(X_n, Z_n)}$$

Assume $Z_K = 1$

$$f_{nk} = P(\underline{Z_{nk}=1} | X_n) = \frac{\pi_k \text{Norm}(X_n | \mu_k, \sigma_k^2)}{\sum_{\ell=1}^K \pi_\ell \text{Norm}(X_n | \mu_\ell, \sigma_\ell^2)}$$

$$P(Z_n | X_n) = \text{Cat}(f_{n1}, f_{n2}, \dots, f_{nK})$$

$$\begin{aligned} f_{nk} &\geq 0 \\ \sum_{k=1}^K f_{nk} &= 1 \end{aligned}$$

Idea 1: Complete likelihood is easier to optimize

Incomplete likelihood: (marginal)

$$P(X_n) = \sum_k \pi_k N(X_n | \mu_k, \sigma_k^2)$$

$$\text{ML: } \pi^*, \mu^*, \sigma^* = \underset{\pi, \mu, \sigma}{\operatorname{argmax}} \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k N(X_n | \mu_k, \sigma_k^2) \right)$$

Take derivative of μ_k :

$$\log(f(\mu_1) + f(\mu_2) + \dots + f(\mu_K))$$

Very complicated! $\mu_1, \mu_2, \dots, \mu_K$ is interactive

Complete likelihood (Joint)

$$\max_{\pi, \mu, \sigma^2} \sum_n \log P(X_n, Z_n)$$

$$\Rightarrow \sum_n \log \sum_{k=1}^K \pi_k^{Z_{nk}} \cdot N(X_n | \mu_k, \sigma_k^2)^{Z_{nk}}$$

$$\Rightarrow \sum_n \sum_k Z_{nk} \log \pi_k + Z_{nk} \log N(X_n | \mu_k, \sigma_k^2)$$

Not interactive. easier to optimize.

Idea 2. Can optimize expectation of complete likelihood

Suppose we have a distribution over z that we thought was accurate:

$$q(z_n) = \text{Cat}(r_{n1}, r_{n2}, \dots, r_{nK})$$

$$\mathbb{E}_q[z_{nk}] = r_{nk}$$

Suppose we wanted to compute

$$\mathbb{E}_{q(z)} [\log p(x, z)]$$

$$p(z_n=1) = r_{nk}$$

$$p(z_n=0) = 1 - r_{nk}$$

$$\text{mean: } (x_{nk} + 0 \times (1 - r_{nk})) = r_{nk}$$

$$= \mathbb{E}_{q(z)} \left[\sum_n \sum_k z_{nk} \log \pi_k + z_{nk} \log N(x_n | \mu_k, \sigma_k^2) \right]$$

$$= \sum_n \sum_k \underbrace{\mathbb{E}_q[z_{nk}]}_{r_{nk}} \log \pi_k + \mathbb{E}_q[z_{nk}] \log N(x_n | \mu_k, \sigma_k^2)$$

Punchline: Expectations of complete log likelihood are easy to evaluate.
easy to optimize for π, μ, σ^2

Idea 3. We can develop an objective

(a) A lower bound on $\log p(x)$, so can be interpreted in principled way as optimizing likelihood

(b) uses expectations of complete likelihood, so it is tractable

$$\log p(x_n) = \log p(x_n) + \sum_z q(z_n)$$

Notation

$$\sum_z q(z) = 1$$

$$= \sum_z \log p(x_n) \cdot q(z_n)$$

$$E(x) = \sum_i x_i p_i$$

$$= \mathbb{E}_{q(z_n)} \log p(x_n)$$

$$\text{Bayes: } p(z_n|x_n) = \frac{p(x_n|z_n)}{p(x_n)}$$

$$= \mathbb{E}_{q(z_n)} \log p\left(\frac{p(x_n, z_n)}{p(z_n|x_n)} \cdot \frac{q(z_n)}{q(z_n)}\right)$$

$$= \mathbb{E}_{q(z_n)} \log p\left(\frac{p(x_n, z_n)}{q(z_n)} \cdot \frac{q(z_n)}{p(z_n|x_n)}\right)$$

$$= \mathbb{E}_{q(z_n)} \left(\log p\left(\frac{p(x_n, z_n)}{q(z_n)}\right) + \log \left(\frac{q(z_n)}{p(z_n|x_n)}\right) \right)$$

$$\begin{aligned}
 &= E_{q(z_n)} \log p\left(\frac{p(x_n, z_n)}{q(z_n)}\right) + E_{q(z_n)} \left(\log \frac{q(z_n)}{p(z_n | x_n)} \right) \\
 &= E_{q(z_n)} \log p\left(\frac{p(x_n, z_n)}{q(z_n)}\right) + \text{entropy of } q(z_n) \\
 \log p(x_n) &= \underbrace{E_{q(z_n)} \log p(x_n, z_n)}_{\substack{\text{expected complete} \\ \text{likelihood}}} - \underbrace{E_{q(z_n)} \log q(z_n)}_{\substack{\text{entropy} \\ \text{of } q(z_n)}} \quad \left| \begin{array}{l} \text{KL divergence:} \\ \text{KL}(P||Q) = \sum p(x) \log \frac{p(x)}{q(x)} \\ \text{KL} \geq 0 \\ \text{entropy:} \\ H(X) = -\sum_i p(x_i) \log p(x_i) \end{array} \right.
 \end{aligned}$$

We have a principled lower bound of the incomplete likelihood!

$$\log p(x_n) \geq \boxed{E_{q(z_n)} \log \frac{p(x_n, z_n)}{q(z_n)}} \quad \text{Evidence Lower Bound (ELBO)}$$

$$L(x_n, r_n, \pi, \mu, \sigma) = E_{q(z_n)} \log p(x_n, z_n) - E_{q(z_n)} \log q(z_n)$$

We can easily calculate L given

| | |
|--|--|
| $\begin{cases} x & \text{data} \\ r_n & \text{prob vector} \\ \pi & \text{weights} \\ \mu & \text{means} \\ \sigma & \text{variances} \end{cases}$ | $\left. \begin{array}{l} \text{parameters} \\ \text{fmm} \end{array} \right\}$ |
|--|--|

Goal:

$$\Rightarrow \sum_k (r_{nk} \log \pi_k + r_{nk} \log \text{Norm PDF}(x_n | \mu_k, \sigma_k^2) - r_{nk} \log r_{nk})$$

Idea 4: EM is coordinate ascent optimization applied to the objective.

New View of EM:

Optimizing an objective \mathcal{L}
that is lower bound of
incomplete likelihood

$$\log p(x|\pi, \mu, \Sigma) \geq \sum_n \mathcal{L}(x_n, r_n, \pi, \mu, \Sigma)$$

E step: Visit every example n ,
find $q(z_n|r_n)$ that
maximizes \mathcal{L}
given current params
 π, μ, Σ

$$r_n = \underset{r_n \in \Delta^K}{\operatorname{argmax}} \mathcal{L}(x_n, r_n, \pi, \mu, \Sigma)$$

Visit each K :

M step: Find point estimate of
each param π, μ, Σ
that maximizes whole dataset objective

$$\pi, \mu, \Sigma = \underset{\pi, \mu, \Sigma}{\operatorname{argmax}} \sum_n \mathcal{L}(x_n, r_n, \pi, \mu, \Sigma)$$

EM (coordinate descent) guarantee improvement!

Question: How good is L as a bound?

Earlier, we saw

$$\log p(x_n) = \mathcal{L}(x_n, r_n, \pi, \mu, \Sigma) + \text{KL}(q(z) \parallel p(z|x))$$

KL always ≥ 0

Recall KL equals 0 ONLY when $q(z) = p(z|x)$

Earlier, we said

$$p(z_n|x_n) = \text{Cat}(y_{n1}, y_{n2}, \dots, y_{nK})$$

where $y_{nk} = \frac{\pi_k \text{Norm}(x_n | \mu_k, \Sigma_k)}{\sum_l \pi_l \text{Norm}(x_n | \mu_l, \Sigma_l)}$

So if we set $r_n = y_n$, then $\text{KL term} = 0.0$ and bound is tight

Turns out, this is optimal E step update,

$$\log p(x_n) = \mathcal{L}(x_n, y_n, \pi, \mu, \Sigma)$$

with equality

EM for Gaussian Mixtures

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters (comprising the means and covariances of the components and the mixing coefficients).

1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \quad (9.23)$$

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (9.24)$$

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T \quad (9.25)$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad (9.26)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}). \quad (9.27)$$

4. Evaluate the log likelihood

$$\ln p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad (9.28)$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

Page 15

The General EM Algorithm

Given a joint distribution $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ over observed variables \mathbf{X} and latent variables \mathbf{Z} , governed by parameters $\boldsymbol{\theta}$, the goal is to maximize the likelihood function $p(\mathbf{X}|\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$.

1. Choose an initial setting for the parameters $\boldsymbol{\theta}^{\text{old}}$.

2. **E step** Evaluate $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$.

3. **M step** Evaluate $\boldsymbol{\theta}^{\text{new}}$ given by

$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) \quad (9.32)$$

where

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}). \quad (9.33)$$

4. Check for convergence of either the log likelihood or the parameter values.
If the convergence criterion is not satisfied, then let

$$\boldsymbol{\theta}^{\text{old}} \leftarrow \boldsymbol{\theta}^{\text{new}} \quad (9.34)$$

and return to step 2.