

# SPR Day 19

Algorithms for estimating the parameters of a Gaussian mixture model using (penalized) maximum likelihood

Readings: Bishop PRML Sec. 9.2.1  
ML for GMMs  
Sec. 9.2.2  
EM for GMMs

Outline: (1) Maximum Likelihood (ML) and Penalized ML Optimization Problems

not in Bishop PRML → (2) Gradient descent methods  
intuition + derivation

(3) Coordinate descent methods  
intuition + derivation

(4) Expectation-Maximization algorithm

# ML Estimation: Problem Statement

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

$K$ : number of assumed clusters

Goal: Estimate values of all GMM parameters

total size  $K$   
 $K \times D$   
 $K \times D \times D$

$$\begin{aligned} \pi &= \pi_{1:K} : K\text{-length vector with non-negative entries that sum to one} \\ \mu &= \mu_{1:K} : \mu_k \text{ is } D\text{-length vector of reals} \\ \Sigma &= \Sigma_{1:K} : \Sigma_k \text{ is } D \times D \text{ sym, pos definite matrix} \end{aligned}$$

Method: Maximize likelihood of  $N$  observed examples  $X_{1:N}$   
under model where each  $x_n \stackrel{iid}{\sim} \text{GMM}(\pi, \mu, \Sigma)$

Objective:  $\pi^*, \mu^*, \Sigma^* = \underset{\substack{\pi \in \Delta^K \\ \mu: \mu_k \in \mathbb{R}^D \\ \Sigma: \Sigma_k \in \mathbb{R}^{D \times D} \\ \text{sym} \\ \text{pos def}}}{\text{argmax}} \sum_{n=1}^N \log \text{GMMPDF}(x_n | \pi, \mu, \Sigma)$

equivalently as a min problem:

$$\pi^*, \mu^*, \Sigma^* = \underset{\substack{\pi \in \Delta^K \\ \mu: \mu_k \in \mathbb{R}^D \\ \Sigma: \Sigma_k \in \mathbb{R}^{D \times D}}}{\text{argmin}} \underbrace{-1 \cdot \sum_{n=1}^N \log \text{GMMPDF}(x_n | \pi, \mu, \Sigma)}_{\mathcal{L}^{ML}(\pi, \mu, \Sigma)}$$

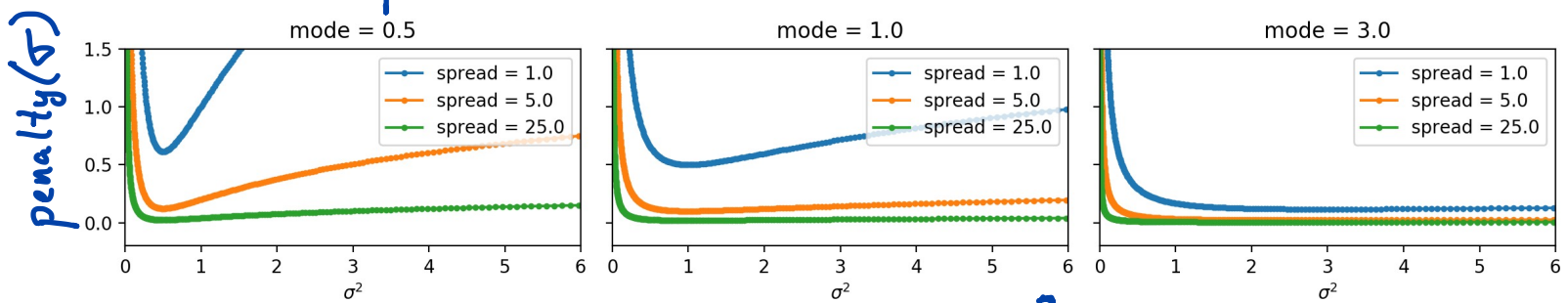
Loss function to minimize:  $\mathcal{L}^{ML}$

# Penalty on Variance Parameter

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  $\Sigma_k = \begin{bmatrix} \sigma_{k1}^2 & & & \\ & \sigma_{k2}^2 & & \\ & & \dots & \\ & & & \sigma_{kD}^2 \end{bmatrix}$   
 call this a "diagonal" covariance.

Then we can define penalty over each dimension's *std. deviation* scalar

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


$m > 0$  is the "mode", value of  $\sigma^2$  that is least penalized  
 $s > 0$  is the "spread", higher values  $\rightarrow$  larger region around mode is favored  
 Note:  $\sigma^2$  near zero is always highly penalized

Can view is 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.

# Penalized ML Estimation: Problem Statement

$$\hat{\pi}^*, \hat{\mu}^*, \hat{\Sigma}^* = \underset{\substack{\pi \in \Delta^K \\ \mu: \mu_k \in \mathbb{R}^D \\ \sigma: \sigma_k \in \mathbb{R}_+^D}}{\operatorname{argmin}} \underbrace{-\sum_{n=1}^N \log \text{GMM PDF}(x_n | \pi, \mu, \sigma)}_{\alpha^{\text{PML}}(\pi, \mu, \sigma)} + \lambda \sum_{k=1}^K \sum_{d=1}^D \text{penalty}(\sigma_{kd})$$

Loss function to minimize:  $\alpha^{\text{PML}}$

We'll use this in CP3!

# Gradient Descent as an algorithm to estimate GMM Parameters

Goal: Find values of  $\pi^*$ ,  $\mu^*$ ,  $\Sigma^*$  that minimize either:

- ML loss  $d^{ML}$
- penalized ML loss  $d^{PML}$

Idea: All parameters have continuous real domains, so let's use gradient descent

Sketch: Input: initial values  $\pi^0, \mu^0, \Sigma^0$

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

$$\pi^t \leftarrow \pi^{t-1} - \epsilon \nabla_{\pi} d(\pi^{t-1}, \mu^{t-1}, \Sigma^{t-1})$$
$$\mu^t \leftarrow \mu^{t-1} - \epsilon \nabla_{\mu} d(\pi^{t-1}, \mu^{t-1}, \Sigma^{t-1})$$
$$\Sigma^t \leftarrow \Sigma^{t-1} - \epsilon \nabla_{\Sigma} d(\pi^{t-1}, \mu^{t-1}, \Sigma^{t-1})$$

Problem: Both  $\pi$  and  $\Sigma$  are constrained

- $\pi$  must sum to one and be non-negative
- $\Sigma_k$  must be a symmetric, pos def matrix

The GD update may produce values that violate constraints

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

Solution: Reparameterize. Find unconstrained parameterization that maps to our constrained version.

# Unconstrained Parameterization for a Probability Vector

Let  $\pi = [\pi_1, \pi_2, \dots, \pi_K]$  be a "probability vector"

That means  $\pi$  defines a discrete PMF over  $K$  choices.

Mathematically,  $\pi$  must have non-negative entries + sum to one.

## Unconstrained Parameterization:

Let  $s \in \mathbb{R}^K$  be an unconstrained vector of reals

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

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

aka softmax

$$Z = \sum_{l=1}^K e^{s_l}$$

by definition, output vector  
is a probability vector

# An Unconstrained Parameterization for a Positive Vector

Let  $\sigma_k$  be a vector that must be positive.  
 $\sigma_k = [\sigma_{k1} \dots \sigma_{kD}]$

## OPTION 1

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

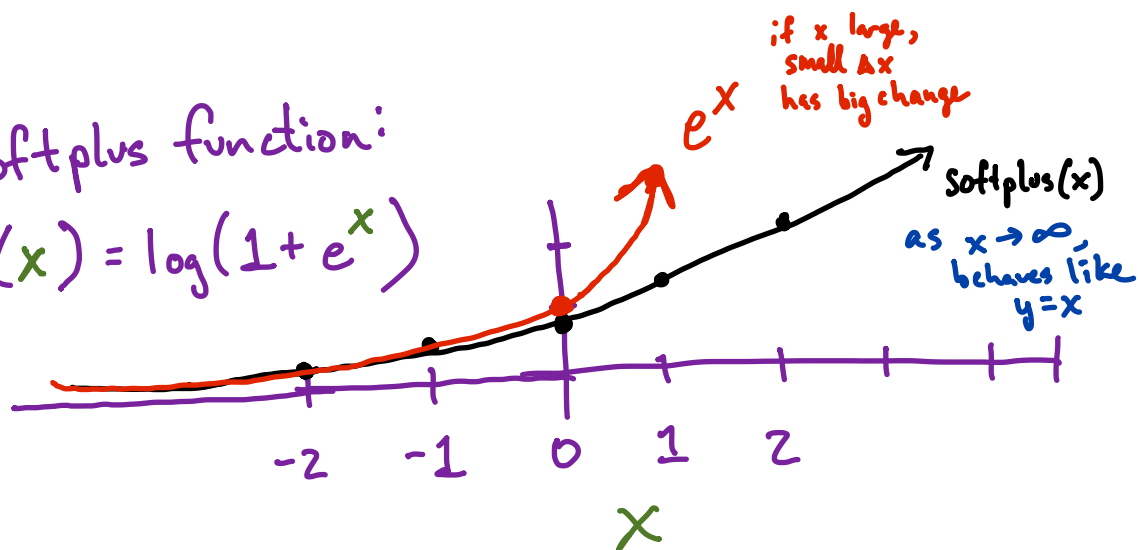
$$\sigma(t_k) \triangleq [e^{t_{k1}}, e^{t_{k2}}, \dots, e^{t_{kD}}]$$

by definition, for any input  $t$ ,  
 $\sigma(t)$  will have all positive entries

## OPTION 2

Define softplus function:

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



Give  $t_k \in \mathbb{R}^D$ , we define transformation  $\sigma(\cdot)$ :

$$\sigma(t_k) = [\text{softplus}(t_{k1}), \dots, \text{softplus}(t_{kD})]$$

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

# 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:  $s^*, \mu^*, t^* = \operatorname{argmin} \alpha^{\text{TPML}}(s, \mu, t)$   
then set  $\pi^* \leftarrow \pi(s^*), \sigma^* \leftarrow \sigma(t^*)$ . Return  $\pi^*, \mu^*, \sigma^*$

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

Grad:  $\nabla_s \alpha^{\text{RPML}} = \nabla_s \pi(s) \cdot \nabla_{\pi} \alpha^{\text{PML}}(\pi(s), \dots)$   
 $\nabla_t \alpha^{\text{RPML}} = \nabla_t \sigma(t) \cdot \nabla_{\sigma} \alpha^{\text{PML}}(\dots, \sigma(t))$   
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



# Towards a Coordinate Descent Algorithm to Estimate Parameters

Recall our ML objective (ignore penalty for simplicity here)

$$\pi^*, \mu^*, \Sigma^* = \underset{\substack{\pi \text{ valid} \\ \mu: \text{GMM} \\ \Sigma \text{ params}}}{\text{argmin}} \underbrace{-1 \cdot \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \text{MVNPDF}(x_n | \mu_k, \Sigma_k)}_{d^{ML}}$$

Idea: Analytically consider an optimal set of parameters  $\pi^*, \mu^*, \Sigma^*$ .  
To be optimal, they must satisfy these blue eqns.

$$\nabla_{\pi} d^{ML}(\pi^*, \mu^*, \Sigma^*) = \vec{0}_K \quad \leftarrow \text{all zeros shape } (K,)$$

$$\nabla_{\mu_k} d^{ML}(\pi^*, \mu^*, \Sigma^*) = \vec{0}_D$$

$$\nabla_{\Sigma_k} d^{ML}(\pi^*, \mu^*, \Sigma^*) = \vec{0}_{D,D}$$

Let's expand this out:

$$\nabla_{\mu_k} d^{ML} = \sum_{n=1}^N \nabla_{\mu_k} [\log \text{GMMPDF}(x_n | \pi, \mu, \Sigma)]$$

$$= \sum_{n=1}^N \frac{1}{\text{GMMPDF}(x_n)} \cdot \left[ \sum_{j=1}^K \pi_j \cdot \nabla_{\mu_k} N(x_n | \mu_j, \Sigma_j) \right]$$

$$= \sum_{n=1}^N \frac{1}{p(x_n)} \cdot \pi_k N(x_n | \mu_k, \Sigma_k) \cdot \nabla_{\mu_k} \left[ -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right]$$

$$\gamma_{nk} = \frac{p(z_{nk}=1 | x_n)}{p(x_n)} = \frac{p(x_n, z_{nk}=1)}{p(x_n)}$$

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

**Insight:**

Can simplify using soft "responsibilities"  $\gamma_n \in \Delta^K$  fix these, and becomes easy to solve.

Continuing to solve blue equations

$$\nabla_{\pi} d^{ML}(\pi^*, \mu^*, \Sigma^*) = \vec{0}_K$$

$$\nabla_{\mu_k} d^{ML}(\pi^*, \mu^*, \Sigma^*) = \vec{0}_D$$

$$\nabla_{\Sigma_k} d^{ML}(\pi^*, \mu^*, \Sigma^*) = \vec{0}_{D,D}$$

Can write all gradients  
in terms of

optimal parameters  $\pi^*, \mu^*, \Sigma^*$

and responsibilities

$\gamma_n \in \Delta^K$   
(treat as additional variable)

For example, solving for  $\mu$  gives

$$\sum_{n=1}^N \gamma_{nk} \Sigma_k^{-1} (x_n - \mu_k) = \vec{0}$$

$$\Sigma_k^{-1} \left( \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k) \right) = \vec{0}$$

$$\sum_{n=1}^N \gamma_{nk} (x_n - \mu_k) = \vec{0}$$

multiply both sides by  $\Sigma_k$

$$\mu_k^* = \frac{\sum_n \gamma_{nk} x_n}{\sum_n \gamma_{nk}}$$

(PRML Eq. 9.19)

Interpretation

weighted empirical mean of data assigned to cluster  $k$

Similar derivations yield

$$\pi_k^* = \frac{1}{N} \sum_{n=1}^N \gamma_{nk}$$

(PRML 9.22)

fraction of all data points assigned to cluster  $k$

$$\Sigma_k^* = \frac{1}{\left(\sum_n \gamma_{nk}\right)} \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k^*) (x_n - \mu_k^*)^T$$

(PRML 9.19)

weighted empirical covariance of data assigned to cluster  $k$

# Coordinate Descent algorithm for GMMs

Similar to K Means

Minimizes the negative ML loss function  $d^{ML}$

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  $k \in 1, 2, \dots, K$ :

$$\tilde{\gamma}_{nk} \leftarrow \pi_k^{t-1} \text{MVNormPDF}(x_n | \mu_k^{t-1}, \Sigma_k^{t-1})$$
$$\gamma_n \leftarrow \left[ \frac{\tilde{\gamma}_{n1}}{\sum_l \tilde{\gamma}_{nl}}, \frac{\tilde{\gamma}_{n2}}{\sum_l \tilde{\gamma}_{nl}}, \dots, \frac{\tilde{\gamma}_{nK}}{\sum_l \tilde{\gamma}_{nl}} \right]$$

M step  
update parameters

for cluster  $k \in 1, 2, \dots, K$ :

$$\pi_k^t \leftarrow \frac{1}{N} \sum_n \gamma_{nk}$$
$$\mu_k^t \leftarrow \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} x_n$$
$$\Sigma_k^t \leftarrow \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (x_n - \mu_k^t)(x_n - \mu_k^t)^T$$

could do with diagonal parametrization  
 $\Sigma_k^t = [\sigma_1^2, \dots, \sigma_D^2]$   
can adjust to penalty term

Will converge to a fixed point where  $\pi^*, \mu^*, \Sigma^*$  is local optima of  $d^{ML}$