

CDA Homework 3

Author: Joseph D Ciaravino

1.1 - see attachment

Conceptual Questions

1.1

PDF:

$$f(x; k, \theta) = \frac{1}{\Gamma(k)\theta^k} x^{k-1} e^{-\frac{x}{\theta}}, \quad x \geq 0$$

Log Likelihood = $L(k, \theta) = \prod_{i=1}^N f(x_i; k, \theta)$

$$\ell(k, \theta) = (k-1) \sum_{i=1}^N \ln(x_i) - \sum_{i=1}^N \frac{x_i}{\theta} - Nk \ln(\theta) - N \ln(\Gamma(k))$$

$$\frac{\partial \ell}{\partial \theta} = 0 \Rightarrow \boxed{\frac{1}{KN} \sum_{i=1}^N x_i}$$

1.2 Please compare the pros and cons of KDE over histogram, and give atleast one advantage and disadvantage to each.

KDE

Pro - flexible (various smoothing kernel functions) - kernel bandwidth can change smoothness as well.

Cons - Poor in terms of memory requirement. All data will need to be in memory ($m \times n$), expensive computationally.

Histogram

Pros - Can easily expand for more data (placing in bins is computation). Parameters increase with 1/bin size which is an advantage over KDE in that it can be more memory efficient in certain circumstances, depending on bin size.

Cons - output depends on bins size. Not efficient for high-dimensional data

1.3

Shown below are how Bayes rule in part(A) is used to form the basis for constructing the E-step in EM for GMM. The components in Section (A) are linked to components in fully formed calculation in Section (B)

(Conceptual Questions)

1.3

A) $P(z|x) = \frac{P(x|z)P(z)}{P(x)}$ = $\frac{P(x,z)}{\sum_{z'} P(x,z')}$

\downarrow likelihood

\downarrow prior

\downarrow posterior

Posterior in e-step

$P(z|K) = \pi'_k$

B) $\pi'_k = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{k=1..K} \pi_k N(x_i | \mu_k, \Sigma_k)}$

\uparrow prior

\uparrow likelihood

Normalization constant

$\left\{ \sum_{z'} \pi'_{z'} N(x | \mu_{z'}, \Sigma_{z'}) \right\}$

Normalization constant
circle over all values

2.A

\Leftarrow

2B.

from question 1.3

$$E\text{-step} \quad t_k^i = \frac{\pi_k N(x^i | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k N(x^i | \mu_k, \Sigma_k)}$$

Closed form:

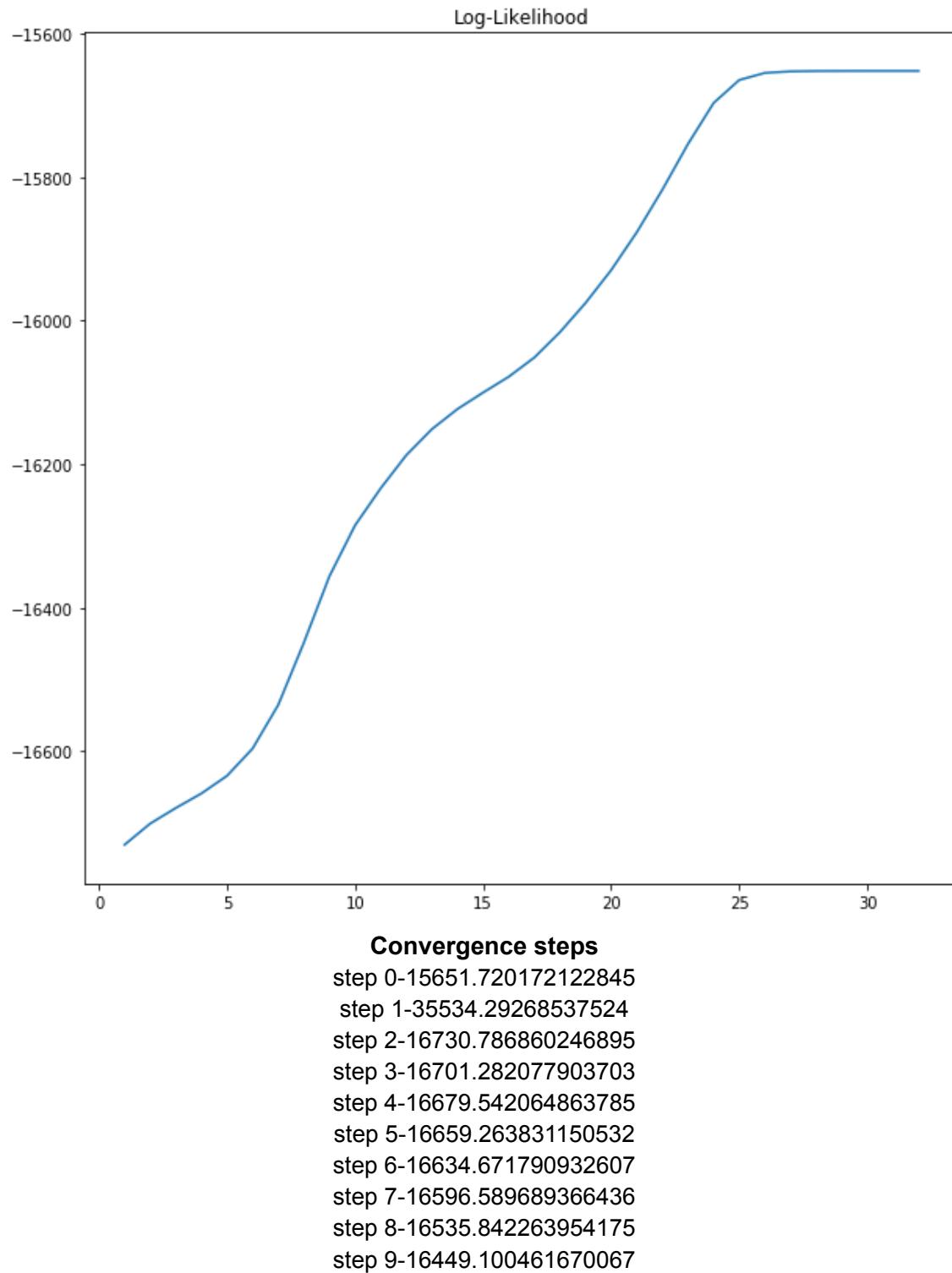
$$t_k^i = \frac{\pi_k |\Sigma_k|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)\right)}{\sum_{k=1}^K \left[\pi_k |\Sigma_k|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)\right) \right]}$$

M-step source: Sides EM Algorithm

$$\pi_k = \frac{\sum_{i=1}^n t_k^i}{n} \quad \mu_k = \frac{\sum_{i=1}^n t_k^i x^i}{\sum_{i=1}^n t_k^i} \quad \Sigma_k = \frac{\sum_{i=1}^n t_k^i (x^i - \mu_k)(x^i - \mu_k)^T}{\sum_{i=1}^n t_k^i}$$

2.B

As shown below the Log-likelihood function converges after a reasonable number of steps and time (in this instance 32) decreasing the log-likelihood until the function is maximized.

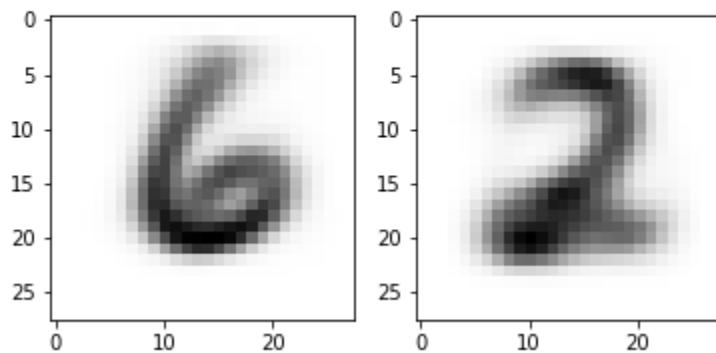


```
step 10-16356.310514457975
step 11-16285.398108684658
step 12-16233.76097593973
step 13-16187.542795014568
step 14-16151.274063243598
step 15-16123.161782855062
step 16-16100.215349168291
step 17-16078.156998488743
step 18-16051.362895303719
step 19-16015.97885868341
step 20-15975.251266938812
step 21-15929.44656565035
step 22-15876.666302886388
step 23-15817.224481862124
step 24-15753.475177275322
step 25-15696.398958787231
step 26-15664.364618231984
step 27-15654.515252597856
step 28-15652.379628648454
step 29-15651.89065830424
step 30-15651.765713163051
step 31-15651.732126330218
step 32-15651.722902182253
training covered
```

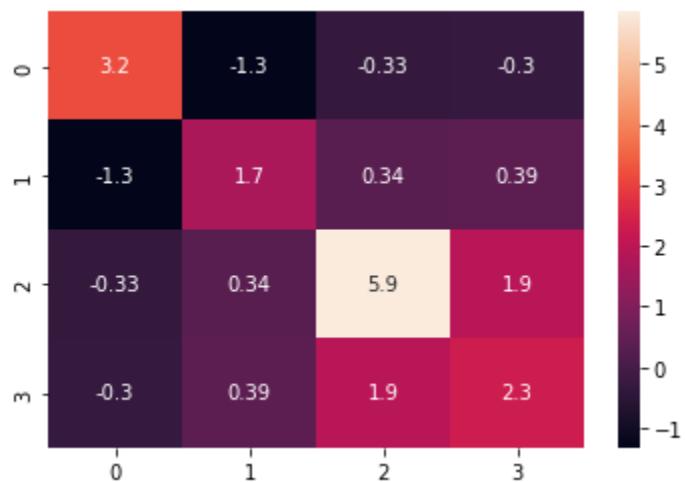
2C.

Mean Images

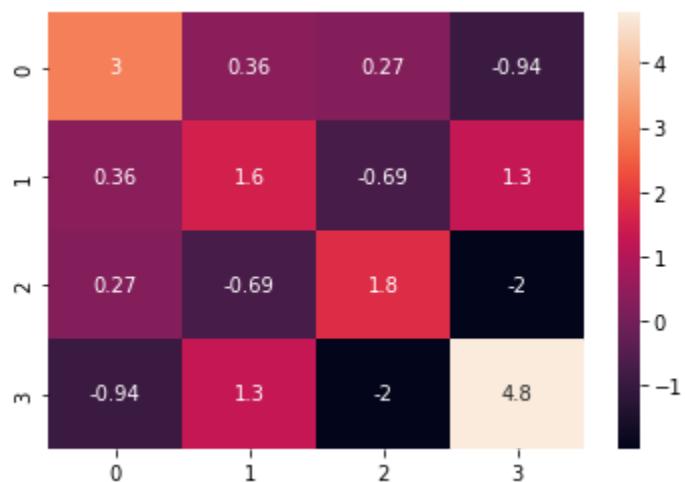
Mean Images 6 and 2



First Covariance Matrix



Second Covariance Matrix



Weights

Component 1: 0.5134041139402375
Component 2: 0.48659588605976245

2.D

Confusion Matrix - EM GMM

	precision	recall	f1-score	support
0	0.94	0.99	0.96	975
1	0.99	0.93	0.96	1015
accuracy			0.96	1990
macro avg	0.96	0.96	0.96	1990
weighted avg	0.96	0.96	0.96	1990

\

Confusion Matrix - Kmeans Algorithm

	precision	recall	f1-score	support
0	0.95	0.93	0.94	1044
1	0.93	0.94	0.93	946
accuracy			0.94	1990
macro avg	0.94	0.94	0.94	1990
weighted avg	0.94	0.94	0.94	1990

The precision for 2 labels is slightly higher for k-means than the GMM, however GMM outperforms for the 6 digit by a greater margin. Overall the GMM computationally doesn't take much longer than kmeans and seems to have greater accuracy.

Additionally, GMM is a “soft” classifier as opposed to kmeans hard classification using Euclidean distance, so by simply inferring labels using tau, we are missing some extra information on bayesian probability.