

## PGM HOMEWORK2 REPORT

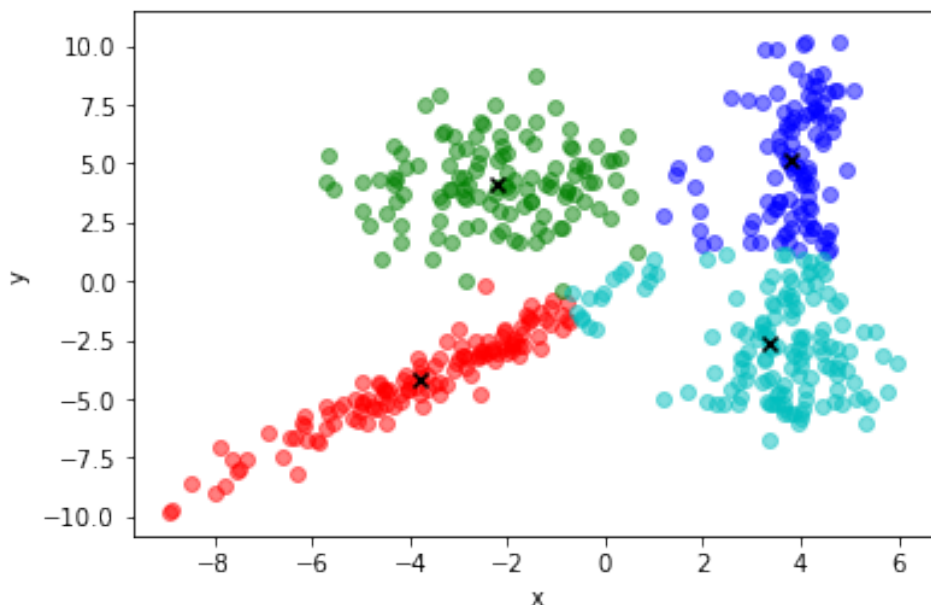
Please, see the scans in the end of the report for theoretical questions.

### 4 Implementation - Gaussian mixtures :

a) For the implementation of the first question:  $K\_means(data,k)$  which is initialized randomly by  $k$  different points from data and returns a list containing in the first element the clusters found and in the second their centers. The function  $distortion(kmeans)$  takes the  $k\_means$  results as arguments and return the distortion of the classes found.

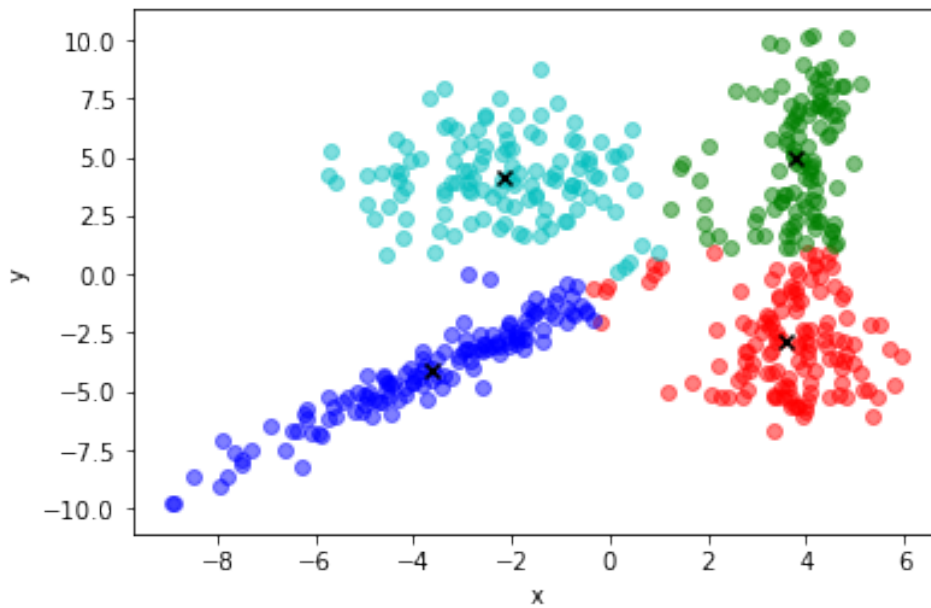
To plot the clusters and theirs centers, I implemented the function  $plt\_kmeans(cluster,center,color)$  which plot the cluster with color color and its center in black.

Figure 1: K-means algorithm Result for a random initialization  
Distortion = 1108.46



Cluster centers are:  $\begin{bmatrix} -3.78 & -4.22 \\ -2.24 & 4.16 \\ 3.8 & 5.1 \\ 3.36 & -2.66 \end{bmatrix}$

Figure 2: K-means algorithm Result for a second random initialization  
Distortion = 1102.55



Cluster centers are:  $\begin{bmatrix} 3.6 & -2.89 \\ -2.16 & 4.11 \\ -3.64 & -4.05 \\ 3.79 & 5. \end{bmatrix}$

Table 1: Comparison between different results of k-means algorithms on train data

Clusters' centers	Distortion	Clusters' centers	Distortion
$\begin{bmatrix} 3.6 & -2.89 \\ -2.16 & 4.11 \\ -3.64 & -4.05 \\ 3.79 & 5. \end{bmatrix}$	1102.54	$\begin{bmatrix} -2.14 & 3.97 \\ -3.72 & -4.18 \\ 3.79 & 5. \\ 3.57 & -2.88 \end{bmatrix}$	1103.92
$\begin{bmatrix} -3.78 & -4.22 \\ -2.24 & 4.16 \\ 3.8 & 5.1 \\ 3.36 & -2.66 \end{bmatrix}$	1108.46	$\begin{bmatrix} -2.24 & 4.13 \\ -3.82 & -4.27 \\ 3.8 & 5.1 \\ 3.34 & -2.64 \end{bmatrix}$	1109.42
$\begin{bmatrix} 3.8 & 5.1 \\ -3.66 & -4.07 \\ 3.48 & -2.7 \\ -2.24 & 4.24 \end{bmatrix}$	1105.84	$\begin{bmatrix} -3.78 & -4.22 \\ 3.36 & -2.71 \\ -2.24 & 4.16 \\ 3.8 & 5.03 \end{bmatrix}$	1107.88

## **Rhoulam Adil**

From the different results, we could conclude that k-means algorithm get trapped in local minimas and we have to run it several time with random initialization and preserve the best clustering which has the minimal distortion. However, we can notice that for the results in the table 1 the distortion value is in the range [1102,1110] which gives an error of  $8/1106 = 0.7\%$  which is not a large precision error. We can notice also that the centers are approximately the same (differ by the first value after comma).

**b)** For the theoretical question, please see the scans in the end of the end of the report. For the implementation of the EM algorithm for Gaussian Mixture when the covariance matrices are isotropic (proportional to the identity matrix), I used the functions: *initializationStep(data,k)*: Which returns initialized means of gaussians and weights with K\_means algorithm and coefficients of covariances with a strict positive random value, *stepE(data,means,cov,weights,k)* which returns the responsibilities (posterior probabilities) and the current log-likelihood, *stepM(data,resp,k)* which returns the updated means, variances and weights given the responsabilities. Then implemented the EM algorithm with function *algorithmEM(data,k)* which returns the means, covariances, weights and responsibilities when the log-likelihood converges.

The graphical Representation of the training data, the centers, as well as 90% percentage of covariance ellipses learnt by EM algorithm is shown in Figure 4.

For the representation of the latent variables for all data points with the parameters learned by EM, I implemented the function *clusters(data,k)* which assigns each point to the cluster with the highest value of responsibility. The Figure 5 shows the clusters with different colors.

Figure4: Training data plotted with centers and 90 percent of the isotropic covariance matrices ellipse learnt by EM

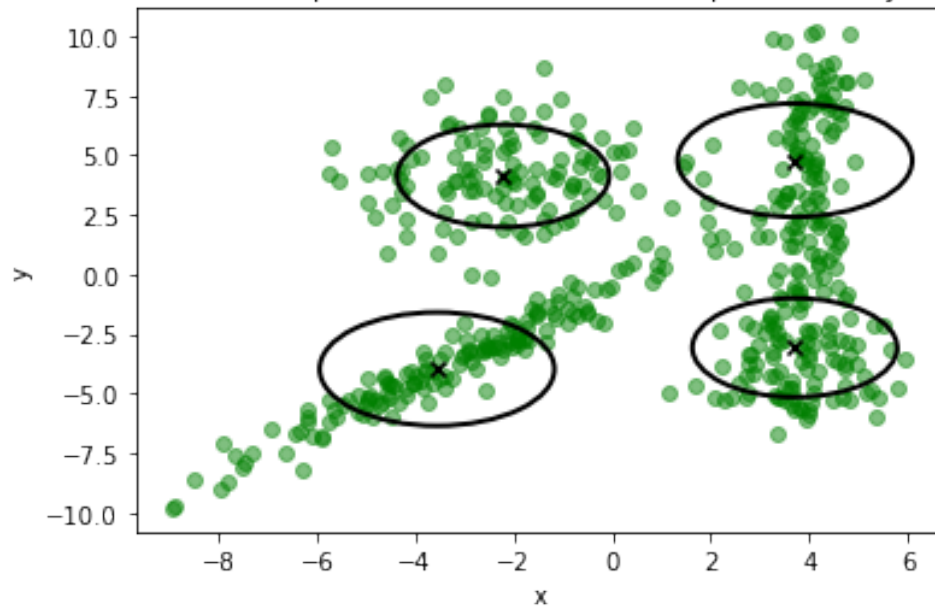
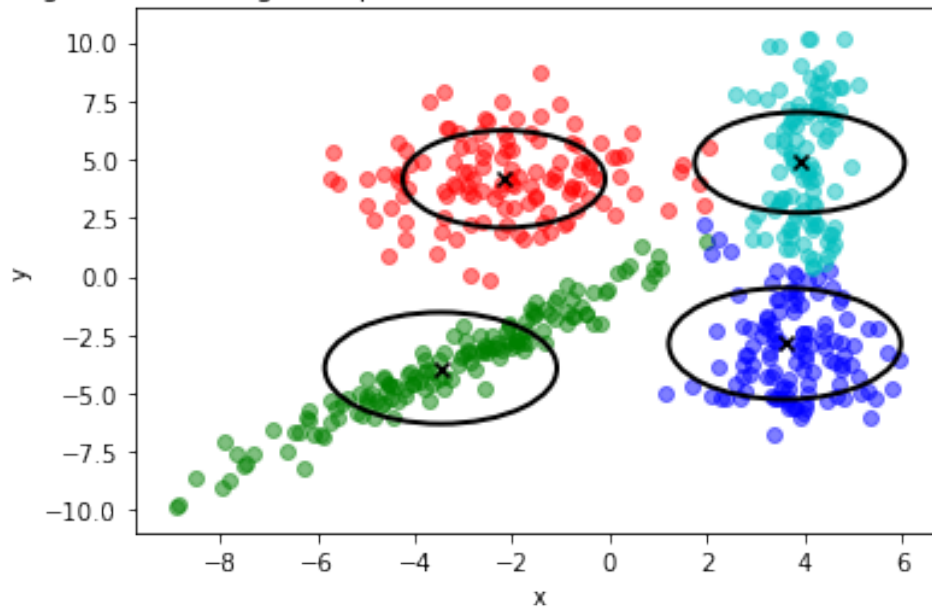


Figure 5: Training data plotted with different colors for latent variables

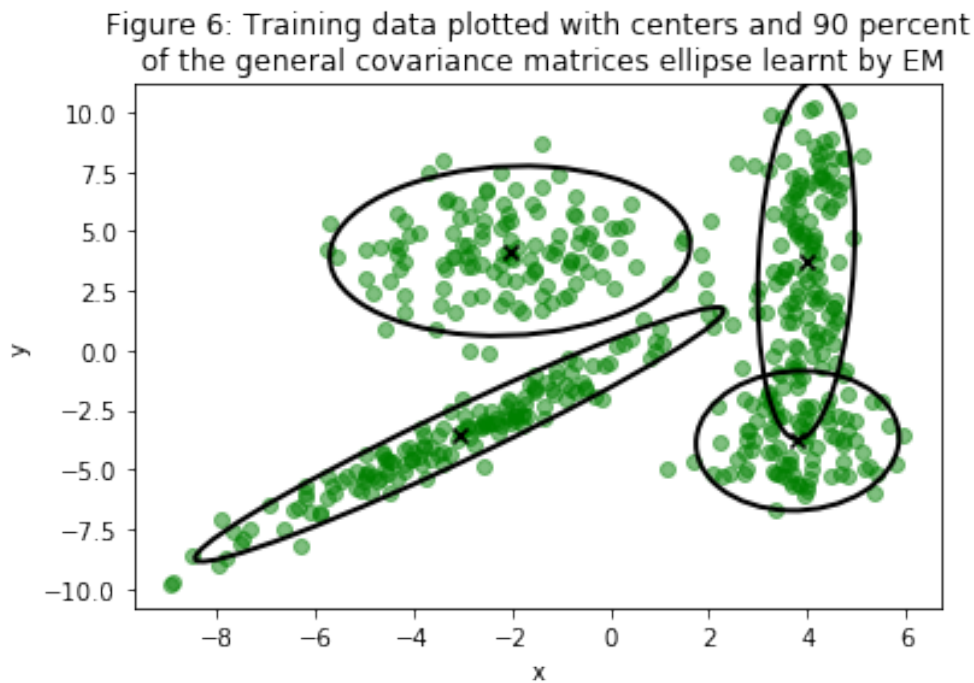


**c).** For the implementation of the EM algorithm for Gaussian Mixture when the covariance matrices are in general form, I used the functions:  
`intializationStepGeneralCase(data,k)` Which returns initialized means of gaussians, weights and covariances with K\_means algorithm. `stepE(data,means,cov,weights,k)`

**Rhoulam Adil**

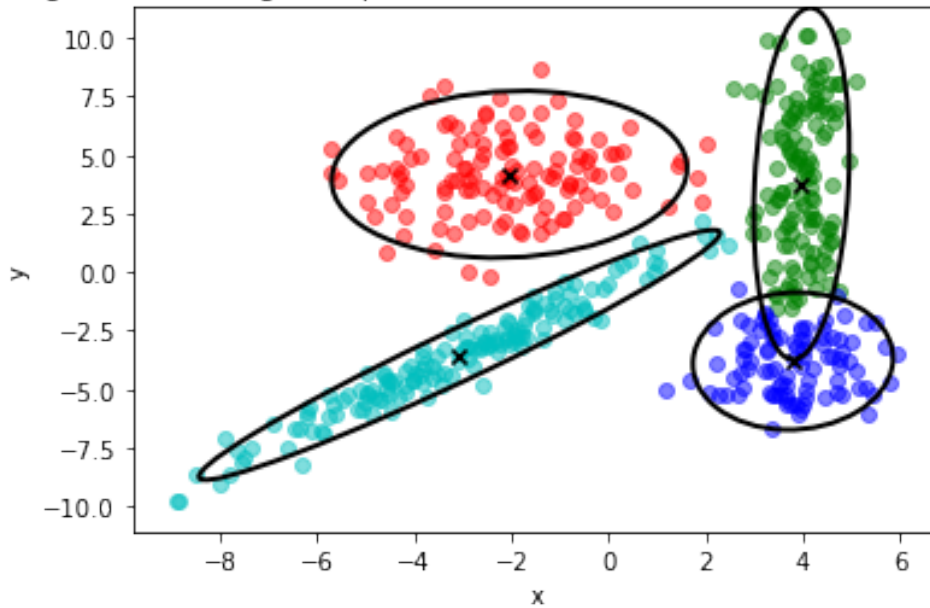
the same one as for the isotropic case and `stepMGeneralCase(data,resp,k)` which returns the updated means, covariance matrices and weights given the responsibilities. Then implemented the EM algorithm with function `algorithmEMgeneralCase(data,k)` which returns the means, covariances, weights and responsibilities when the log-likelihood converges.

The graphical Representation of the training data, the centers, as well as 90% percentage of covariance ellipses learnt by EM algorithm is shown in Figure 6.



For the representation of the latent variables for all data points with the parameters learned by EM, I implemented the function `clustersEMgeneralCase(data,k)` which assigns each point to the cluster with the highest value of responsibility. The Figure 7 shows the clusters with different colors.

Figure 7: Training data plotted with different colors for latent variables



d) The comparison of maximum log-likelihood between different models and data (train and test) is show in the following table:

Table 2: Comparison of Maximum Log-likelihood of Train data and Test data

	Maximum Log-Likelihood on Train data	Maximum Log-Likelihood on Test data
<b>Isotropic covariance matrices</b>	-4469.33	-4341.51
<b>General covariance matrices</b>	-3358.18	-3404.92

We can notice that the maximum log-likelihood when using general covariance matrices is larger than the one using isotropic covariance matrices, but if we compare the elapsed times while convergence of the two models ( 34.39 seconds on test-data with isotropic model and 48.83 seconds on test-data with general model), we find that the second one converges rapidly than the first, this was expected because the isotropic covariances has only 1 parameter to learn while in the general case we have  $d(d+1)/2$  with  $d$ : the dimension of data. Moreover, by analyzing the figures 4, 5, 6 and 7, we can say that the model of isotropic covariances isn't suitable for our data.



## HOME WORK 2 - PGM

### 1) Conditional independence and factorizations:

1) if  $x \perp\!\!\!\perp y | z$  then  $p(x, y | z) = p(x | z) p(y | z)$

$$p(x | y, z) = \frac{p(x, y | z)}{p(y | z)} = \frac{p(x, y | z) p(z)}{p(y | z) p(z)} = \frac{p(x | z) p(y | z) p(z)}{p(y | z) p(z)} = p(x | z)$$

then  $p(x | y, z) = p(x | z)$

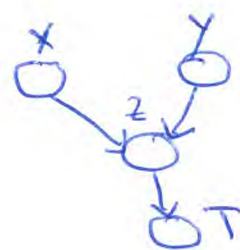
if  $p(x | y, z) = p(x | z)$ :

$$p(x, y | z) = \frac{p(x | y, z) p(y | z)}{p(z)} = \frac{p(x | z) p(y | z)}{p(z)} = \frac{p(x | z) p(y | z) p(z)}{p(z)} = p(x | z) p(y | z)$$

then  $p(x, y | z) = p(x | z) p(y | z)$

2) - let  $p \in \mathcal{L}(G)$

then  $p = p(x) p(y) p(z | x, y) p(T | z)$



-  $x \perp\!\!\!\perp y | T$ ?

using d-separation:

- if  $c$  is the chain from  $x, y$   
we have  $z \in c$  and  $(x, z, y)$  is a v-structure  
and  $z \notin S$  ( $S = \{T\}$  set that is observed)  
but  $T$  is a descendant of  $z$  and it's observed  
then  $x \not\perp\!\!\!\perp y | T$

3)- a) let  $X \perp Y | Z = 0$ ;  $X \perp Y | Z = 1$  and  $Z \in \{0, 1\}$

such that  $P(Z=1) = q \in [0, 1]$

let  $j \in \{0, 1\}$ :

$$P(X, Y, Z=j) = P(Z=j) P(X, Y | Z=j) \stackrel{①}{=} P(Z=j) P(X | Z=j) P(Y | Z=j) \\ = P(Z=j) \times \frac{P(X) P(Z=j | X)}{P(Z=j)} \times \frac{P(Y) P(Z=j | Y)}{P(Z=j)} \quad ③$$

From another side:

$$P(X, Y, Z=j) = P(X, Y) P(Z=j | X, Y) \stackrel{②}{=} P(X) P(Y) P(Z=j | X, Y) \quad ④$$

From ③ and ④; we deduce:

$$P(Z=j) P(Z=j | X, Y) = P(Z=j | X) P(Z=j | Y)$$

if we denote  $x = P(Z=1 | X)$ ,  $y = P(Z=1 | Y)$

then:

$$\begin{cases} q P(Z=1 | X, Y) = xy \\ (1-q) P(Z=0 | X, Y) = (1-x)(1-y) \end{cases}$$

$$(1-q) q P(Z=1 | X, Y) = (1-q) q (1 - P(Z=0 | X, Y)) \\ = (1-q) q - q (1-q) P(Z=0 | X, Y) \\ = (1-q) q - q (1-x)(1-y)$$

$$\text{then } (1-q) xy = (1-q) q - q (1-x)(1-y)$$

$$\Leftrightarrow xy - qxy = q - q^2 - q(1-x-y+xy)$$

$$\Leftrightarrow xy - qxy - q + q^2 + q - qx - qy + qxy = 0$$

$$\Leftrightarrow xy - qx - qy + q^2 = 0$$

$$\Leftrightarrow x(y-q) = q(y-q)$$

$$\Leftrightarrow (x-q)(y-q) = 0 \Leftrightarrow x-q=0 \text{ or } y-q=0$$

$$\Leftrightarrow \text{or } P(Z=1 | X) = P(Z=1)$$

$$\Leftrightarrow \text{or } \begin{cases} P(Z=1 | Y) = P(Z=1) \\ P(Z=0 | X) = P(Z=0) \text{ and } P(Z=0 | Y) = P(Z=0) \end{cases}$$

$$\Leftrightarrow \text{or } \begin{cases} X \perp Z \\ Y \perp Z \end{cases}$$



b) Let  $(X, Y, Z)$  random variables in a finite space we can enumerate  $Z(\Omega)$  such that:

$$Z(\Omega) = \{z_1, \dots, z_n\}$$

Let  $A \subset Z(\Omega)$   $A \cup A^c = Z(\Omega)$  and  $A \cap A^c = \emptyset$

$(A, A^c)$  partition of  $Z(\Omega)$ :  $q_A = P(Z \in A)$  and

$$1 - q_A = q_{A^c} = P(Z \in A^c)$$

$$\begin{cases} Z' = 0 \Leftrightarrow Z \in A \\ Z' = 1 \Leftrightarrow Z \in A^c \end{cases}$$

$Z'$  is a binary variable, we will have the same result as the question before. Let reproduce the result.

$$P(X, Y, Z \in A) = P(X)P(Y)P(Z \in A | X, Y)$$

$$\text{and } P(X, Y, Z \in A) = P(Z \in A) \frac{P(Z \in A | X)P(X)}{P(Z \in A)} \frac{P(Z \in A | Y)P(Y)}{P(Z \in A)}$$

$$\text{then } P(Z \in A | X, Y) P(Z \in A) = \underbrace{P(Z \in A | X)}_{x_A} \underbrace{P(Z \in A | Y)}_{y_A}$$

$$q_A P(Z \in A | X, Y) = x_A y_A$$

$$(1 - q_A) P(Z \in A^c | X, Y) = (1 - x_A)(1 - y_A)$$

We do calculation, we find:

$$(x_A - q_A)(y_A - q_A) = 0 \quad \forall A \in Z(\Omega)$$

$$\text{then } \begin{cases} x_A = q_A \\ \text{or} \\ y_A = q_A \end{cases} \quad \forall A \in Z(\Omega)$$

$$\Leftrightarrow \begin{cases} P(Z \in A | X) = P(Z \in A) \\ \text{or} \\ P(Z \in A | Y) = P(Z \in A) \end{cases} \quad \forall A \in Z(\Omega)$$

$$\Leftrightarrow \begin{cases} P(Z = z_i | X) = P(Z = z_i) \\ \text{or} \\ P(Z = z_i | Y) = P(Z = z_i) \end{cases}$$

$$A = \{z_i\} \quad \forall i \in \{1, \dots, n\}$$



Let  $\mathcal{J} = \{i \in \{1, \dots, n\} \mid P(z=z_i | X) = P(z=z_i) \}$

Then:  $\forall i \notin \mathcal{J}: P(z=z_i | Y) = P(z=z_i)$   
 we should show that  $\mathcal{J} = \emptyset$  or  $\mathcal{J} = \{1, \dots, n\}$

$$\sum_{i \in \mathcal{J}} P(z=z_i) + \sum_{i \notin \mathcal{J}} P(z=z_i) = \sum_{i \in \mathcal{J}} P(z=z_i | X) + \sum_{i \notin \mathcal{J}} P(z=z_i | Y)$$

$$1 = \sum_{i \in \mathcal{J}} P(z=z_i | X) + \sum_{i \notin \mathcal{J}} P(z=z_i | Y)$$

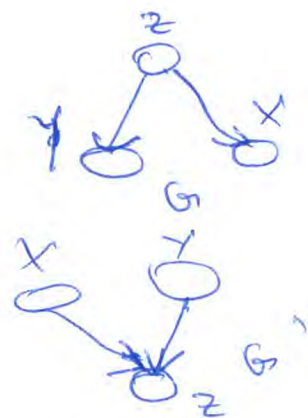
$$\Rightarrow \sum_{i \in \mathcal{J}} (P(z=z_i | X) - P(z=z_i | Y)) = 0$$

$$\Rightarrow \sum_{i \in \mathcal{J}} \frac{P(z=z_i | X)}{P(X)} - \frac{P(z=z_i | Y)}{P(Y)} = 0$$

$$\Rightarrow \sum_{i \in \mathcal{J}} P(Y) P(z=z_i | X) - P(z=z_i | Y) P(X) = 0$$

$$\sum_{i \in \mathcal{J}} P(z=z_i | X) = \sum_{i \in \mathcal{J}} P(z=z_i | Y)$$

I didn't have time to think more about it.  
 but maybe by using  $|A|=2, |A|=3 \dots$  we  
 could find the independence, because it seems  
 logical that  $X \perp\!\!\!\perp Z$  or  $Y \perp\!\!\!\perp Z$ . If we suppose  
 that  $X \not\perp\!\!\!\perp Z$  and  $Y \not\perp\!\!\!\perp Z$  then  
 the graphs  $G$  and  $G'$  are equivalent  
 since  $L(G) = L(G')$  but they don't  
 have the same ~~to~~ independence  
 since a result shows that two graphs  
 are Markov equivalent iff they have  
 the same links (edges without orientation)  
 and the same set of v structures



The demo is given by Pent and Verma but  
I didn't find it on the web.



## 2 - Distributions factorizing in a graph:

1)  $G = (V, E)$  a DAG ;  $i \rightarrow j \in E$  is a covered edge (i.e.  $\pi_j = \pi_i \cup \text{hit}$ ).

$G' = (V, E')$  with  $E' = E \setminus \{i \rightarrow j\} \cup \{j \rightarrow i\}$

let show  $L(G) = L(G')$ :

first  $\pi_j = \pi_i \cup \text{hit} \Leftrightarrow$

$E' = E \setminus \{i \rightarrow j\} \cup \{j \rightarrow i\} \Rightarrow$

the first graph is an  $G$   
and the second in  $G'$ .

to show that  $L(G) = L(G')$

it's sufficient to show that  
since the edge that relates  $i$  and  $j$  is the  
only one that changed orientation

let  $p \in L(G)$ .  $p = p(\pi_i) p(x_j | \pi_i, x_i) p(x_i | \pi_i)$

$$p(x_j | \pi_j) = p(x_j | x_i, \pi_i) = \frac{p(x_j, x_i | \pi_i)}{p(x_i | \pi_i)}$$

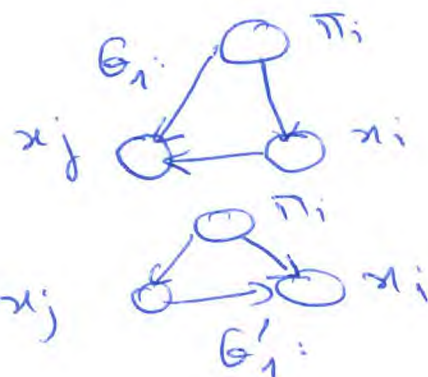
$$= \frac{p(x_i | x_j, \pi_i) p(x_j | \pi_i)}{p(x_i | \pi_i)}$$

$$\text{then } p(x_j | x_i, \pi_i) p(x_i | \pi_i) = p(x_i | x_j, \pi_i) p(x_j | \pi_i)$$

and:  $p = p(\pi_i) p(x_i | x_j, \pi_i) p(x_j | \pi_i) \in L(G')$

By equivalence  $L(G) = L(G')$

then  $L(G) = L(G')$ .

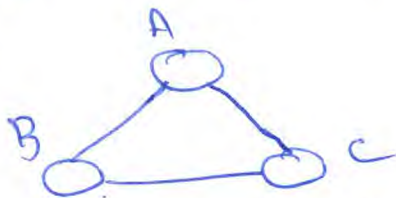




2) Let  $G$  be a directed tree and  $G'$  it's corresponding undirected tree

First let show that the maximal cliques in  $G'$  are composed of only two nodes if  $|V| \geq 2$

Let suppose that there exists a clique of 3 nodes  $A, B, C$  in  $G'$  then it takes the form



then since  $G'$  is the

undirected tree corresponding to  $G$ , there exists an orientation of this clique that is in  $G$ .

We can see that all the  $3^2 = 9$  possible orientations gives a V-structure; but the directed tree  $G$  doesn't contain V-structures then it's an absurd. then:  $|\text{cliques}_{\max}(G')| \leq 2$

- if  $G' = (\text{root}, \emptyset)$  then  $\text{cliques}(G') = 1$  and  $h(G') = h(G)$  obviously.

- if  $|V| \geq 2$  then there exists a node that have a parent then  $|\text{cliques}_{\max}(G')| \geq 2$ .

finally  $|\text{cliques}_{\max}(G')| = 2$

without loss of generality let  $(x_1 \rightarrow x_n)$

be the nodes of  $G$  in a chronological order:

$\rho(G) = \prod_{i=1}^n \rho(|\pi_i|/\pi_i)$  and  $|\pi_i| = 1$  (cardinal of the sets of parents of  $x_i$  is equal to one by the definition of a tree).

- for  $G'$  the only cliques that exists are  $(\pi_i, x_i) \quad \forall i \in [1, n]$



$$\text{so } p \in L(G') \Leftrightarrow p = \frac{1}{Z} \prod_{c \in C} \psi_c(x_c) \\ = \frac{1}{Z} \prod_{i=1}^n \psi_{C_i}(x_{C_i})$$

with:  $C_i = \{x_1, \dots, x_n\}$

by identifying  $\psi_{C_i}(x_{C_i}) = p(x_i | \pi_i)$

$$Z = \sum_{x_1, \dots, x_n} \prod_{i=1}^n \psi_{C_i}(x_{C_i}) \\ = \sum_{x_1, \dots, x_n} \prod_{i=1}^n p(x_i | \pi_i) \\ = \prod_{i=1}^n \sum_{x_i} p(x_i | \pi_i) = 1$$

we find that:  $p = \prod_{i=1}^n p(x_i | \pi_i) \in L(G)$

then  $L(G) = L(G')$  by equivalence.

### 3-Entropy and Mutual Information:

1)  $X$  a discrete RV with  $|X| = k$

the entropy:  $H(X) = - \sum_{x \in X} P_X(x) \log P_X(x)$

with  $P_X(x) = P(X=x)$

$$D(p||q) = \begin{cases} +\infty & \text{if } \exists x \in X, \text{ st } q(x)=0 \text{ and } p(x) \neq 0 \\ \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} & \text{otherwise,} \end{cases}$$

a) let consider  $g: [0, +\infty[ \rightarrow \mathbb{R}$   
 $x \mapsto -x \log(x)$  [ $0 \log 0 = 0$ ]

$$g'(x) = -\log(x) + 1 \Rightarrow g''(x) = -\frac{1}{x} < 0 \quad \forall x \in ]0, +\infty[$$

then  $g$  is concave.

$$H(X) = - \sum_{x \in X} P_X(x) \log P_X(x) = + \sum_{x \in X} g(P_X(x))$$

$$\geq \sum_{x \in X} P_X(x) g(1) = 0 \quad \text{because } g(1) = 0$$

Jensen inequality

$$H(X) = 0 \Rightarrow - \sum_{x \in X} P_X(x) \log P_X(x) = 0$$

$$P_X(x) \log P_X(x) \geq 0 \quad \forall x \in X \quad \text{because } P_X(x) \in [0, 1]$$

$$\text{then } P_X(x) \log P_X(x) = 0 \quad \forall x \in X$$

$$\Rightarrow P_X(x) = 0 \text{ or } P_X(x) = 1 \quad \forall x \in X$$

$$\text{and because } \sum_{x \in X} P_X(x) = 1$$

$$\text{then } \exists! x_0 \in X \text{ st } P_X(x_0) = 1 \text{ and } \forall x \neq x_0, P_X(x) = 0$$

$$\Leftarrow \text{reciprocally: if } \exists! x_0 \in X \text{ st } P_X(x_0) = 1 \text{ and } P_X(x) = 0 \quad \forall x \neq x_0$$

$$\text{then } H(X) = - \sum_{x \in X} P_X(x) \log P_X(x) = 0$$



finally, the equality holds if and only if  $X$  is a constant with probability 1.

b) let  $P_X \sim X$  and  $q \sim U([1, k])$

$$\begin{aligned} D(P_X \| q) &= \sum_{x \in X} P_X(x) \log \left( \frac{P_X(x)}{1/k} \right) \\ &= -H(X) + \sum_{x \in X} P_X(x) \log(k) \\ &= -H(X) + \log(k) \end{aligned}$$

$$\boxed{D(P_X \| q) = -H(X) + \log(k)}$$

c) Using the fact that  $D(P_X \| q) \geq 0$ , we deduce that  $H(X) \leq \log(k)$

2)  $(X_1, X_2)$  two RVs defined over the finite set  $X_1 \times X_2$ .  $P_{12}$ ,  $P_1$  and  $P_2$  denote respectively the joint distribution of  $(X_1, X_2)$ , the marginals of  $X_1, X_2$ .

$$I(X_1, X_2) = \sum_{(x_1, x_2) \in X_1 \times X_2} P_{12}(x_1, x_2) \log \left( \frac{P_{12}(x_1, x_2)}{P_1(x_1) P_2(x_2)} \right)$$

$$a) \quad I(X_1, X_2) = - \sum_{(x_1, x_2) \in X_1 \times X_2} P_{12}(x_1, x_2) \log \frac{P_1(x_1) P_2(x_2)}{P_{12}(x_1, x_2)}$$

$$= + \sum_{(x_1, x_2) \in X_1 \times X_2} P_{12}(x_1, x_2) f \left( \frac{P_1(x_1) P_2(x_2)}{P_{12}(x_1, x_2)} \right)$$

with  $f: \mathbb{R}_+^{++} \rightarrow \mathbb{R}$  ;  $f'(x) = \frac{1}{x^2} > 0$   
 $x \mapsto -\log(x)$

then  $f$  is convex and using Jensen Inequality

$$I(X_1, X_2) \geq f \left( \sum_{(x_1, x_2) \in X_1 \times X_2} P_{12}(x_1, x_2) \frac{P_1(x_1) P_2(x_2)}{P_{12}(x_1, x_2)} \right)$$

$$= f \left( \sum_{x_1 \in X_1} P_1(x_1) \cdot \sum_{x_2 \in X_2} P_2(x_2) \right) = f(1 \cdot 1) = 0$$

then  $I(x_1, x_2) \geq 0$

b)

$$I(x_1, x_2) = \sum_{(x_1, x_2) \in X_1 \times X_2} p_{12}(x_1, x_2) \log(p_{12}(x_1, x_2)) - \sum_{(x_1, x_2) \in X_1 \times X_2} p_{12}(x_1, x_2) \log(p_1(x_1))$$

$$\log(p_1(x_1)) = \sum_{(x_1, x_2) \in X_1 \times X_2} p_{12}(x_1, x_2) \log(p_1(x_1))$$

$$\begin{aligned} \sum_{(x_1, x_2) \in X_1 \times X_2} p_{12}(x_1, x_2) \log(p_1(x_1)) &= \sum_{x_1 \in X_1} \left( \sum_{x_2 \in X_2} p_{12}(x_1, x_2) \right) \log(p_1(x_1)) \\ &= \sum_{x_1 \in X_1} p_1(x_1) \log(p_1(x_1)) = -H(x_1) \end{aligned}$$

the same for the last term

then  $I(x_1, x_2) = -H(x_1, x_2) + H(x_1) + H(x_2)$

c) we have  $I(x_1, x_2) \geq 0$  then

$$H(x_1, x_2) \leq H(x_1) + H(x_2)$$

$H(x_1, x_2)$  is maximal when  $I(x_1, x_2) = 0$  which means

$$H(x_1, x_2) = H(x_1) + H(x_2)$$

$$I(x_1, x_2) = 0 \Leftrightarrow \sum_{(x_1, x_2) \in X_1 \times X_2} p_{12}(x_1, x_2) \log \left( \frac{p_1(x_1) p_2(x_2)}{p_{12}(x_1, x_2)} \right) = 0$$

$$\sum_{(x_1, x_2) \in X_1 \times X_2} \log \left( \frac{p_1(x_1) p_2(x_2)}{p_{12}(x_1, x_2)} \right) \quad \left( \text{Equality holds in the Jensen inequality} \right)$$

we have:

$$I(x_1, x_2) = D(p_{12} \| p_1 p_2)$$

with  $D$  is the Kullback-Leibler over  $(x_1, x_2)$

$$D(p_{12} \| p_1 p_2) = 0 \Leftrightarrow p_{12} = p_1 p_2$$

then the joint distribution of maximal entropy is given by  $p_{12} = p_1 p_2$  which means  $x_1$  and  $x_2$  are independent



#### 4) Gaussian Mixtures:

b) Let  $D = (x_1, \dots, x_N)$  be  $N$  iid observations of  $x$  that have the distribution  $p(x|z_k) = N(x|\mu_k, \Sigma_k)$  with  $p(z) = \prod_{k=1}^K \pi_k z_k$  ( $\sum_{k=1}^K \pi_k = 1$  and  $\pi_k \geq 0 \forall k \in \{1, \dots, K\}$ )

$$p(x) = \sum_{k=1}^K p(x|z_k=1) p(z_k) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k)$$

$$\log p(D|\pi, \Sigma, \mu) = \log \prod_{i=1}^N p(x_i|\pi, \Sigma, \mu)$$

$$\left[ \log p(D|\pi, \Sigma, \mu) = \sum_{i=1}^N \log \left( \sum_{k=1}^K \pi_k N(x_i|\mu_k, \Sigma_k) \right) \right]$$

If we denote  $\gamma(z_k) = p(z_k=1|x_i)$  then by Bayes formula we deduce that:

$$\gamma(z_k) = \frac{\pi_k N(x_i|\mu_k, \Sigma_k)}{\sum_{p=1}^K \pi_p N(x_i|\mu_p, \Sigma_p)}$$

and considering that the covariance matrices are diagonal:  $\Sigma_k = \sigma_k^2 \mathbf{I}$ ;  $\sigma_k^2 > 0 \forall k \in \{1, \dots, K\}$

calculating the M-step:

- Maximum likelihood:

•) For  $\mu_k$ :  $\sum_{i=1}^N \frac{\pi_k N(x_i|\mu_k, \Sigma_k)}{\sum_{p=1}^K \pi_p N(x_i|\mu_p, \Sigma_p)} \sigma_k^2 (x_i - \mu_k) = 0$

$$\frac{\partial \log p(D)}{\partial \mu_k} = 0 \Rightarrow \sum_{i=1}^N \frac{\pi_k N(x_i|\mu_k, \Sigma_k)}{\sum_{p=1}^K \pi_p N(x_i|\mu_p, \Sigma_p)} \sigma_k^2 (x_i - \mu_k) = 0$$

$$\Rightarrow \sum_{i=1}^N \gamma(z_k) \sigma_k^2 (x_i - \mu_k) = 0$$

$$\Rightarrow \hat{\mu}_k = \frac{\sum_{i=1}^N \gamma(z_k) x_i}{\sum_{i=1}^N \gamma(z_k)}$$

•) For  $\sigma_k^2 = v_k$ :

$$\frac{\partial \log p(D)}{\partial v_k} = \sum_{i=1}^N \frac{\pi_k}{\sum_{p=1}^K \pi_p N(x_i|\hat{\mu}_p, v_p)}$$

$$\frac{\partial N(x_i|\hat{\mu}_p, v_p)}{\partial v_k} = 0$$



$$N(x_i | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{D/2} \det(\Sigma_k I)^{1/2}} \exp \left\{ -\frac{1}{2\Sigma_k} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k) \right\}$$

then:

$$\begin{aligned} \frac{\partial N(x_i)}{\partial \Sigma_k} &= -\frac{D}{2} \Sigma_k^{-\frac{D}{2}-1} \frac{1}{(2\pi)^{D/2}} \exp \left\{ -\frac{1}{2\Sigma_k} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k) \right\} \\ &\quad + \frac{1}{2\Sigma_k^2} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k) N(x_i | \mu_k, \Sigma_k) \\ &= -\frac{D}{2} \Sigma_k^{-1} N(x_i | \mu_k, \Sigma_k) + \frac{1}{2\Sigma_k^2} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k) N(x_i | \mu_k, \Sigma_k) = 0 \end{aligned}$$

$$\frac{\partial \log p(D)}{\partial \Sigma_k} = 0 \Rightarrow \sum_{i=1}^N \sigma(z_{ik}) \left( -\frac{D}{2\Sigma_k} + \frac{1}{2\Sigma_k^2} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k) \right) = 0$$

$$\Rightarrow \Sigma_k = \frac{\sum_{i=1}^N \sigma(z_{ik}) (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k)}{D \sum_{i=1}^N \sigma(z_{ik})}$$

• For  $\pi_k$ : we should take into account the condition that  $\sum_{k=1}^K \pi_k = 1$

We have to maximize:

$$\log p(D | \mu, \pi, \Sigma) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

taking the gradient of this quantity by  $\pi_k$  we get:

$$\sum_{i=1}^N \frac{N(x_i | \mu_k, \Sigma_k)}{\sum_{p=1}^K \pi_p N(x_i | \mu_p, \Sigma_p)} + \lambda = 0$$

then multiplying by  $\pi_k$  and taking the sum over  $k$  we get:

$$N + \lambda = 0 \Rightarrow \lambda = -N$$

then

$$\sum_{i=1}^N \frac{1}{\pi_k} \sigma(z_{ik}) = N$$

$$\Rightarrow \pi_k = \frac{\sum_{i=1}^N \sigma(z_{ik})}{N}$$

$\pi_k, \mu_k$  have the same form as in the case of general covariance matrix.