

# Visualisation avec les cartes topologiques catégorielles

Mustapha LEBBBAH, Fouad BADRAN, Sylvie THIRIA  
CNAM, paris 6, Renault

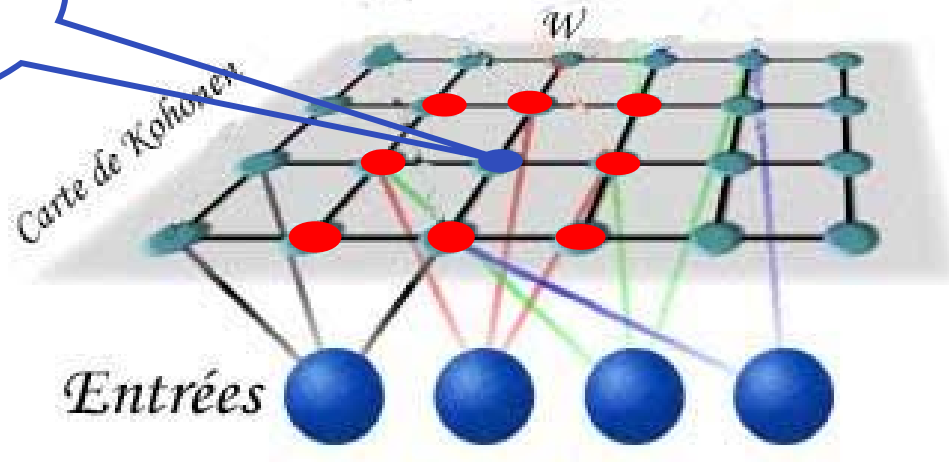
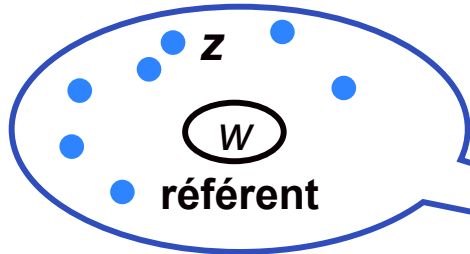
CLERMONT FERRAND 20/01/2004

# PLAN

- Cartes topologiques et Données qualitatives
- Carte topologique binaire: **BTM**
- Carte topologique probabiliste: **CTM**
- Validation du CTM
- Conclusion



# Cartes topologiques et données qualitatives

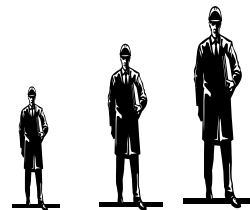


$$K(\delta(c_1, c_2))$$

Fonction noyau

•Données qualitatives

Taille:



Sexe:





# Variables qualitatives et codage

Taille: **P**etit, Moyen, **G**rand

Petit, **M**oyen, **G**rand

Petit, Moyen, **G**rand



Ordinale

1 0 0

1 1 0

1 1 1

Couleur : **r**ouge, vert, bleu

rouge, **v**ert, bleu

rouge, vert, **b**leu



Disjonctif

1 0 0

0 1 0

0 0 1



# Le modèle BTM

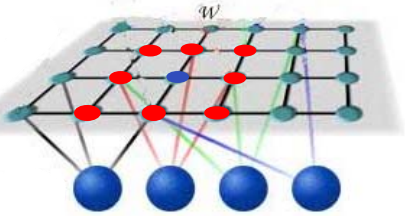
Goveart 88, Marchetti 89

Distance de Hamming

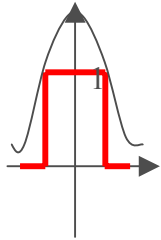


100111101010

Additif, Disjonctif



$H(.,.)$



**BTM**

1	1	1	1	1	0	0	1	1	1	0	0
1	1	0	1	1	1	1	1	0	0	0	0
1	1	1	1	0	0	0	1	1	1	1	1
1	1	1	1	1	1	0	1	1	0	0	0
1	1	1	1	1	1	0	1	1	1	1	0
1	0	0	1	1	1	0	1	1	1	0	0

$K_1$

$K_2$

$K_2$

$K_3$

$K_4$

...

Fonction de coût

$$E(W) = \sum_{z_i \in APP} \sum_{r \in C} K(\delta(c, r)) H(z_i, W_r)$$

1 1 1 1 0 1 0 1 0 1 0 0

Centre médian

**Minimisation utilisant les nuées dynamiques**



# Cartes Topologiques et Probabilité

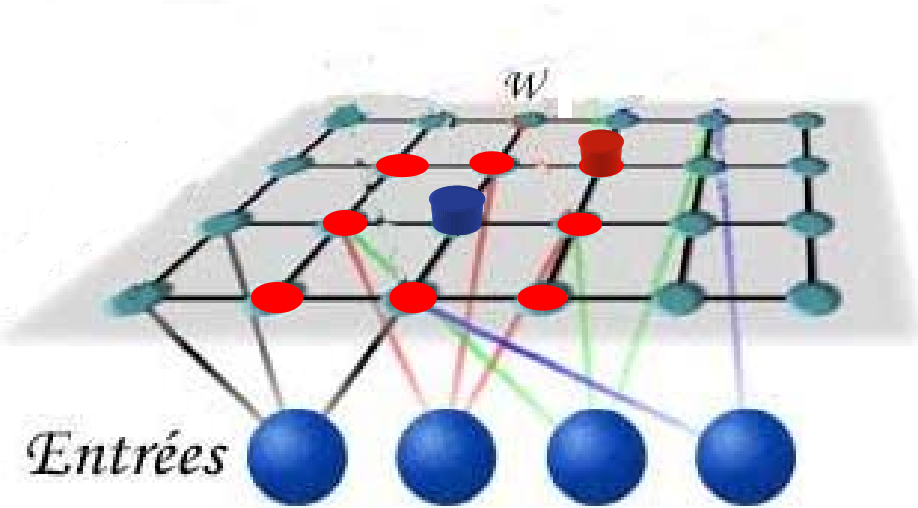
$$p(\mathbf{z}/\mathbf{c}_1)$$

$$p(\mathbf{c}_2)$$

$$p(\mathbf{c}_1/\mathbf{c}_2) = \frac{K(\delta(\mathbf{c}_1, \mathbf{c}_2))}{\sum_{\mathbf{c}}}$$

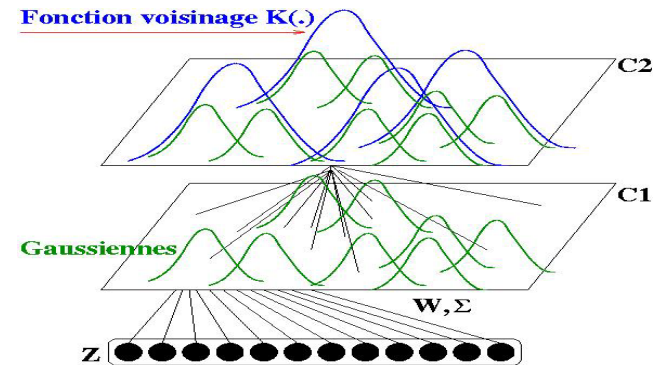
$$p(\mathbf{z}) = \sum_{\mathbf{c}_2} p(\mathbf{c}_2) p_{\mathbf{c}_2}(\mathbf{z})$$

$$p_{\mathbf{c}_2}(\mathbf{z}) = \sum_{\mathbf{c}_1} p(\mathbf{c}_1/\mathbf{c}_2) p(\mathbf{z}/\mathbf{c}_1)$$



$$p(\mathbf{z}/\mathbf{c}_1) = \text{Loi Gaussienne sphérique}$$

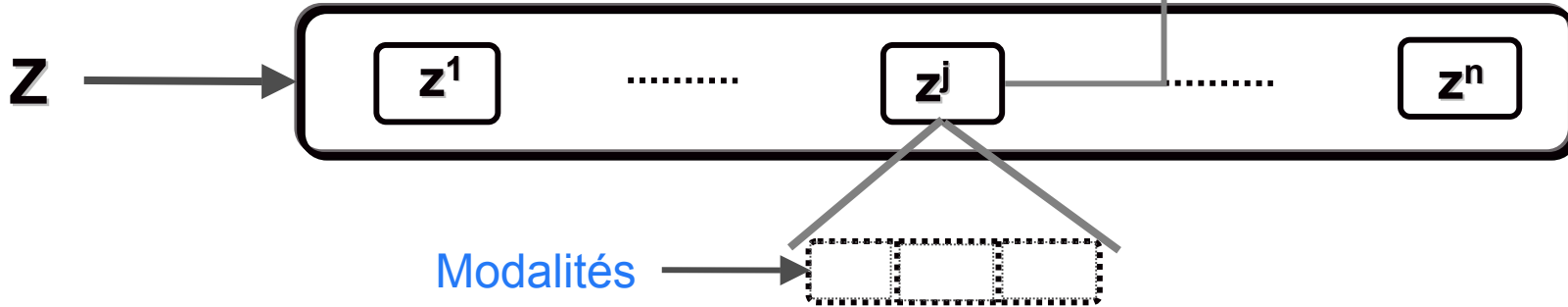
$$p(\mathbf{z}) = \text{Mélange de mélanges locales de lois Gaussiennes}$$





# Le Modèle CTM

Variable catégorielle



Variables indépendantes →  $p(\mathbf{z}/\mathbf{c}) = \prod_{j=1..n} p(z^j/\mathbf{c})$

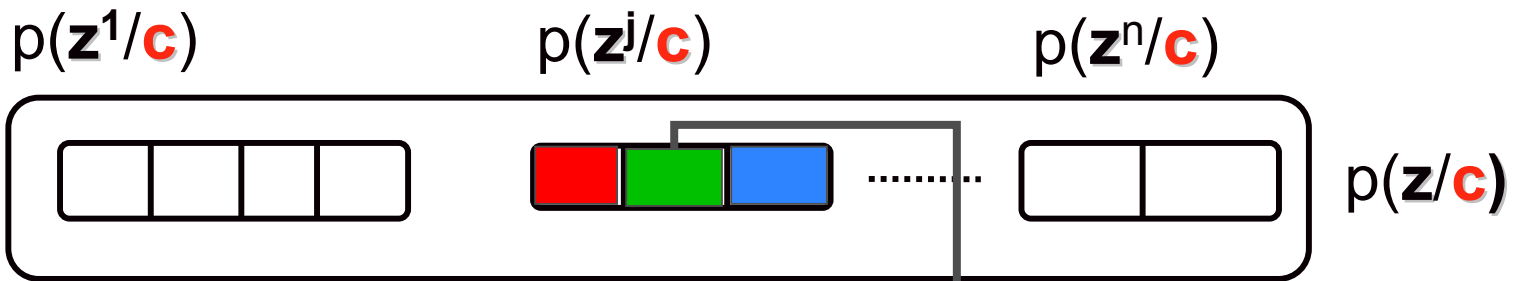


Table de probabilités

$$p(\mathbf{z}) = \sum_{c_2} p(c_2) \sum_{c_1} p(c_1/c_2) p(\mathbf{z}/c_1)$$

⚡ EM

CTM= Carte Topologique+ données catégorielles



# Formules obtenues:

$$p(\mathbf{c}_2) = \frac{\sum_z p(\mathbf{c}_2/z)}{\sum_{\mathbf{c}_2} \sum_z \sum_{\mathbf{c}_1} p(\mathbf{c}_1, \mathbf{c}_2/z)}$$

Probabilité a priori

$$p(z^{j,\text{mod}}/\mathbf{c}) = \frac{\sum_{z \in z_j, \text{mod}} p(\mathbf{c}_1/z)}{\sum_z p(\mathbf{c}_1/z)}$$

Une composante de la table de probabilités unidimensionnelle

$$p(\mathbf{c}_1/z) = \sum_{\mathbf{c}_2} p(\mathbf{c}_1, \mathbf{c}_2/z)$$

$$p(\mathbf{c}_2/z) = \sum_{\mathbf{c}_1} p(\mathbf{c}_1, \mathbf{c}_2/z)$$

$$p(\mathbf{c}_1, \mathbf{c}_2/z) = \frac{p(\mathbf{c}_2) \overset{\text{Fonction de voisinage}}{K(\delta(\mathbf{c}_1, \mathbf{c}_2))} p(z/\mathbf{c}_1)}{\sum_{\mathbf{c}_2} p(\mathbf{c}_2) p(z/\mathbf{c}_2)}$$

Fonction de voisinage



Prise en compte de l'ordre topologique





# Validation du CTM

## Base réelle belge

1106 assurés, 9 variables, Bon / Mauvais conducteur



- Utilité (Privé, Professionnelle)
- Sexe (Homme, Femme, Véhicule de Société )
- Langue (Français, Autre)
- Age (Vieux, Moyen, Jeune)
- Localisation (Capitale, Province)
- Bonus (1,2)
- Police (86, Autre)
- Puissance (Grande, Petite)
- Age Véhicule (Ancien, Nouveau)



$p(z/c)$

U	S	Lg	Ag	Lo	B	Po	Pu	AV
Pv	?	?	?	?	?	?	?	?
Pf	?	?	?	?	?	?	?	?
	?		?					

20 paramètres à estimer

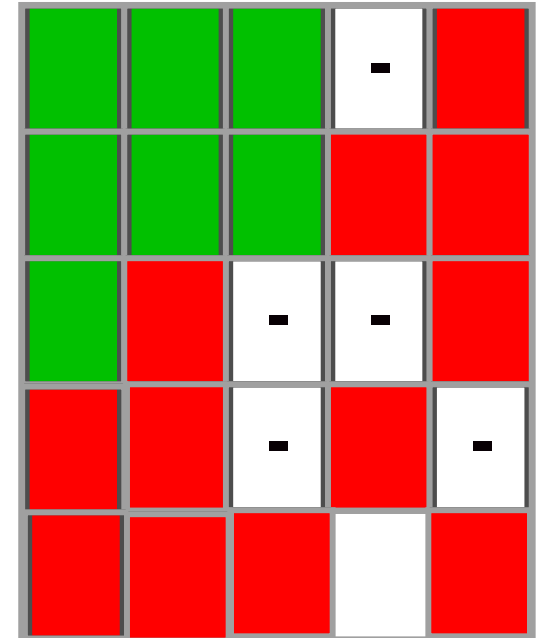


# Visualisation multi-dimensionnelle

H V - An	- J - An	H J - An	F J Pt Nou	- J Gr -
H - - -	- J - An	H J Gr An	F J Pt Nou	VS J Gr Nou
H V Gr -	H - Gr -	VS J Pt Nou	VS J Pt Nou	VS J Gr Nou
- - Gr -	H - Gr -	VS J Pt An	VS J Pt An	VS J Gr Nou
H M Gr -	- M - -	F M Pt -	F M Pt Nou	F M Pt Nou

Représentation des 4 variables

(Sexe, Age, Puissance, Age du Véhicule,)



Affectation avec la probabilités a posteriori  $p(c/z)$

Etiquetage

bon/mauvais



# Base de Sémométrie

## Base réelle

**1128 individus, 70 Mots (variables)**

7 notes sémiométriques (modalités)

- Univocité sémantique
- Stabilité sémantique
- Non-consensualité

ACHETER  
AMBITION  
ANGOISSE  
ARGENT  
ASTUCIEUX  
AUDACE  
BIJOU  
CADEAU

.....

Mot

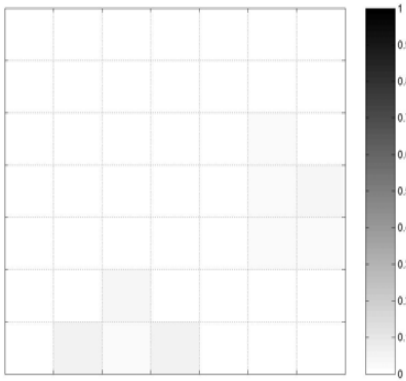


**7 paramètres à estimer**

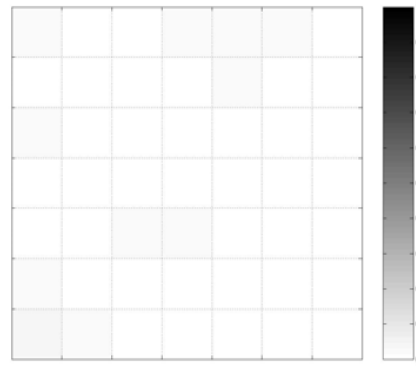


# Répartition de probabilité du mot : ACHETER

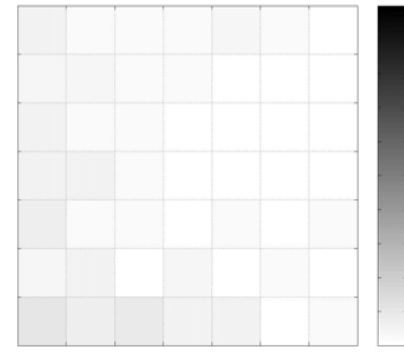
note:1..7



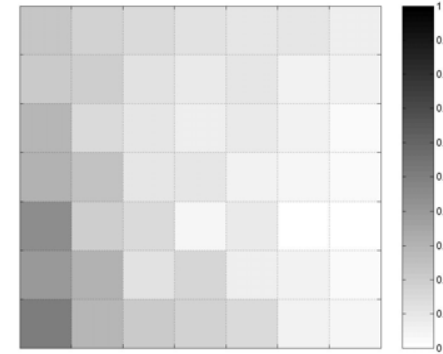
1



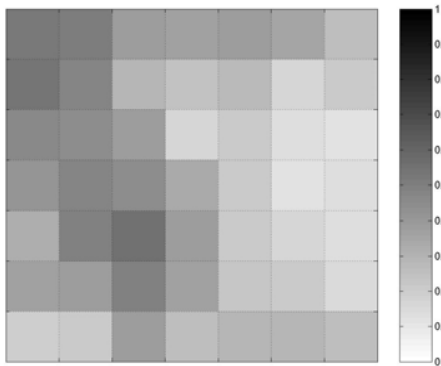
2



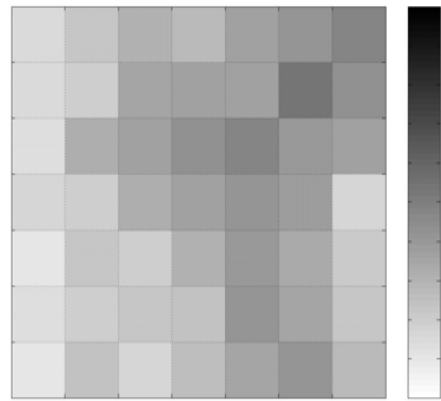
3



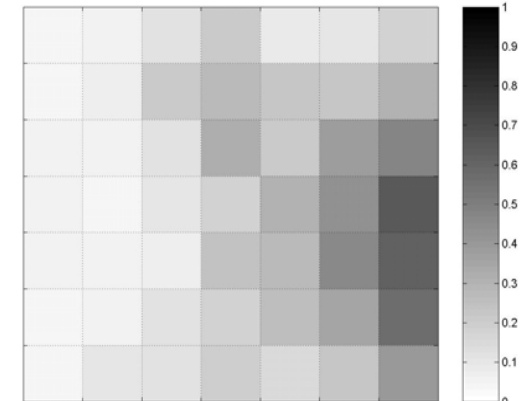
4



5



6

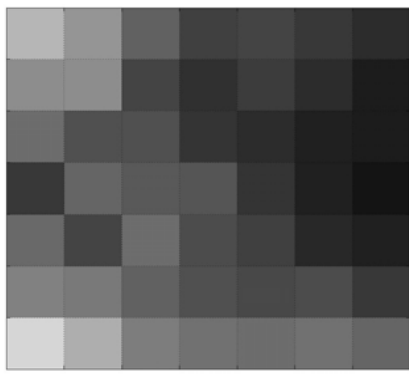


7

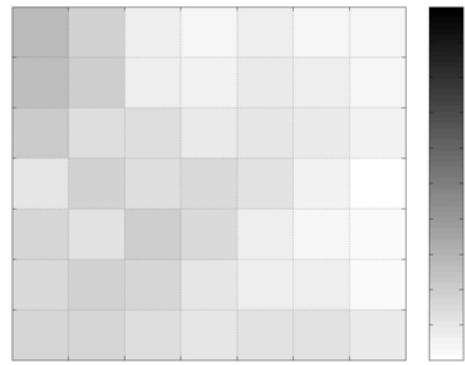


# Répartition de probabilité du mot : Mort

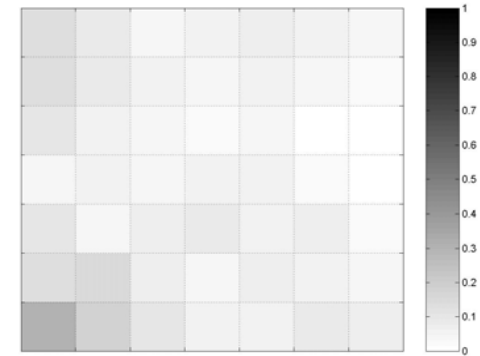
note: 1..7



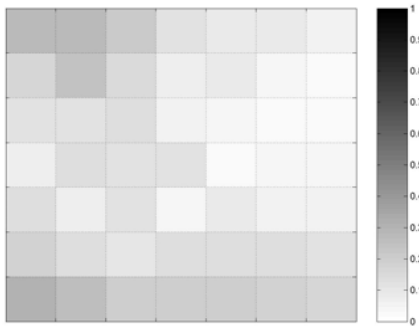
1



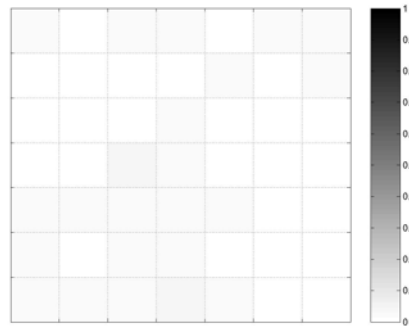
2



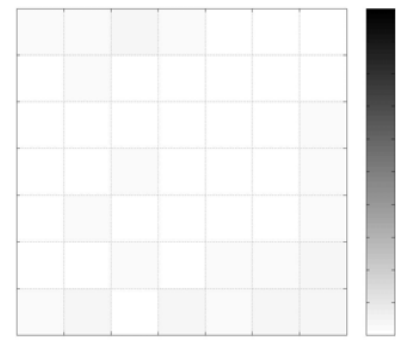
3



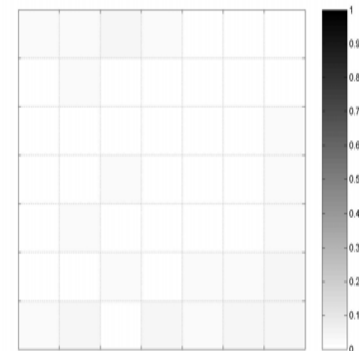
4



5



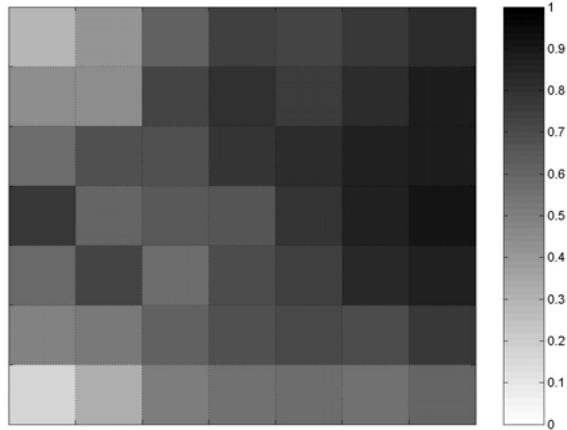
6



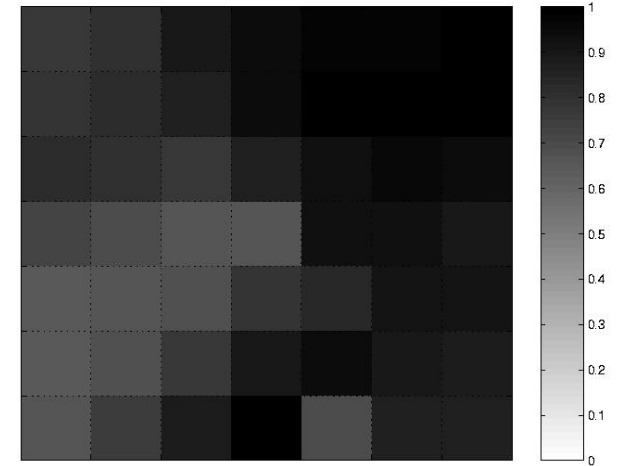
7



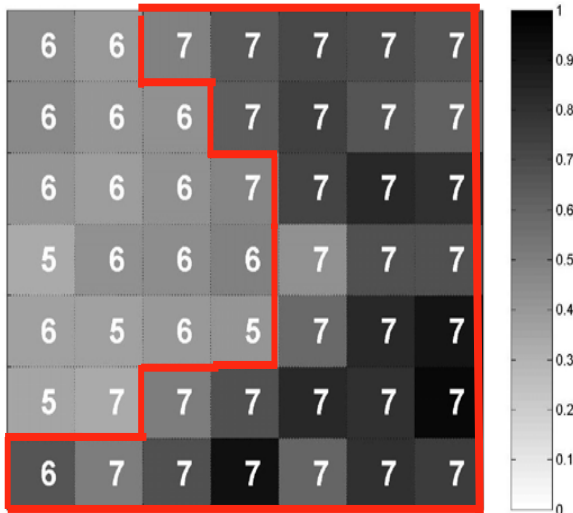
# Répartition de la probabilité Maximale



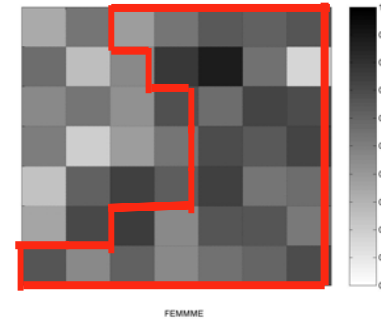
Mort (1)



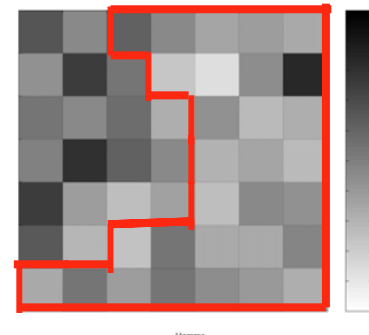
Guerre (1)



Fleur



♀  
Femme



♂  
Homme  
14



# Conclusion

## CTM

- Etendre le champs d'application des cartes topologiques
- Réaliser une partition de données en des sous-ensembles homogènes
- Outils de visualisation
- Utiliser comme classifieur



# Perspectives

- Utiliser une loi adaptée aux données qualitatives codées en additives
- Adapter la CAH à ce type de données
- Evolution des états sur la carte (HMM)

