

Homework 1 PGM

Sharone Dayan, Michael Sutton

October 2017

1 Learning in Discrete Graphical Model

Considérons le modèle suivant :

Soit Z et X deux variables discrètes qui prennent respectivement M et K valeurs tel que :

$$\begin{aligned} p(Z = m) &= \pi_m & \forall m \in \llbracket 1, \dots, M \rrbracket \\ p(X = k | Z = m) &= \theta_{mk} & \forall k \in \llbracket 1, \dots, K \rrbracket, \forall m \in \llbracket 1, \dots, M \rrbracket \end{aligned}$$

Données : On a un échantillon $(z^{(1)}, x^{(1)}), \dots, (z^{(N)}, x^{(N)})$ d'observations i.i.d du couple (Z, X) .

Posons $Y = (Y_1, \dots, Y_M)$ avec $Y_m = \mathbb{1}_{Z=m}$. Posons $B = \begin{pmatrix} B_1 \\ \vdots \\ B_M \end{pmatrix}$ avec $(B_m)_k = \mathbb{1}_{X=k|Z=m}$, on a alors :

$$p(Y_m = m) = p(Z = m) = \pi_m \quad \forall m \in \llbracket 1, \dots, M \rrbracket$$

$$p(Y = y) = \prod_{m=1}^M \pi_m^{y_m}$$

L'évènement $\{Y = k\}$ correspond à $\{Y_k = 1 \text{ et } Y_l = 0 \ \forall l \neq k\}$, $Y \in [0, 1]^K$.

Et concernant B et Z on a les relations :

$$p(B_{mk} = 1) = p((B_m)_k = 1) = p(X = k | Z = m) = \theta_{mk} \quad \forall k \in \llbracket 1, \dots, K \rrbracket, \forall m \in \llbracket 1, \dots, M \rrbracket$$

$$p(B = b) = \prod_{m=1}^M \prod_{k=1}^K \theta_{mk}^{b_{mk}}$$

mettre un texte comme pour Y decrivant B

Ecrivons la fonction de vraisemblance :

$$\begin{aligned}
L((\pi, \theta)) &= p_{(\pi, \theta)}((Z^{(1)}, X^{(1)}) = (z^{(1)}, x^{(1)}), \dots, (Z^{(N)}, X^{(N)}) = (z^{(N)}, x^{(N)})) \\
&= \prod_{n=1}^N p_{(\pi, \theta)}((Z^{(n)}, X^{(n)}) = (z^{(n)}, x^{(n)})) \quad \text{par independance} \\
&= \prod_{n=1}^N p(X^{(n)} = x^{(n)} | Z^{(n)} = z^{(n)}) p(Z^{(n)} = z^{(n)}) \\
&= \prod_{n=1}^N p(B^{(n)} = b^{(n)}) p(Y^{(n)} = y^{(n)}) \\
&= \prod_{n=1}^N \prod_{m=1}^M \prod_{k=1}^K \theta_{mk}^{b_{mk}^{(n)}} \prod_{m=1'}^M \pi_{m'}^{y_{m'}^{(n)}} \\
&= \prod_{n=1}^N \prod_{m=1}^M \prod_{k=1}^K \theta_{mk}^{b_{mk}^{(n)}} \pi_m^{y_m^{(n)}}
\end{aligned}$$

En remarquant que les cas ou $\pi_i = 0$ (resp. $\theta_{mk} = 0$) entrainerai que $y_i = 0$ sur toute les obesrvation (resp. $B_{mk} = 0$ sur toute les observations), les valeurs correspondantes dans le produit valent 1 et "disparaissent". On en conclu donc que la log vraisemblance est bien defini et vaut :

$$\begin{aligned}
l((\pi, \theta)) &= \sum_{n=1}^N \sum_{k=1}^K \sum_{m=1}^M b_{mk}^{(n)} \log(\theta_{mk}) + y_m^{(n)} \log(\pi_m) \\
&= \sum_{m=1}^M \sum_{k=1}^K \sum_{n=1}^N b_{mk}^{(n)} \log(\theta_{mk}) + y_m^{(n)} \log(\pi_m)
\end{aligned}$$

En notant $n_m = \sum_{n=1}^N y_m^{(n)}$, qui correspond au nombre d'observations de Z prenant la valeur j , et en notant $n_{mk} = \sum_{n=1}^N b_{mk}^{(n)}$ qui correspond au nombre d'observations de (Z, X) prenant la valeur (m, k) on a :

$$l((\pi, \theta)) = \sum_{m=1}^M \sum_{k=1}^K n_{mk} \log(\theta_{mk}) + n_m \log(\pi_m)$$

L'objectif est donc de maximiser la fonction de log vraisemblance $l((\pi, \theta))$ sous les contrainte que $\sum_{m=1}^M \pi_m = 1$ et $\forall m \sum_{k=1}^K \theta_{mk} = 1$, ou en d'autres termes :

$$\min_{\substack{\pi_m \geq 0 \\ \theta_{mk} \geq 0}} f((\pi, \theta)) = -l((\pi, \theta)) = \sum_{m=1}^M \sum_{k=1}^K -n_{mk} \log(\theta_{mk}) - n_m \log(\pi_m) \quad \text{s.c.} \quad 1^T \pi = 1 \in \mathbb{R} \quad \text{and} \quad 1^T \theta = 1 \in \mathbb{R}^M$$

Le Lagrangien de ce problème donne :

$$\mathcal{L}(\pi, \theta, \lambda) = -l((\pi, \theta)) + \lambda_0 \left(\sum_{m=1}^M \pi_m - 1 \right) + \sum_{m=1}^M (\lambda_j \left(\sum_{k=1}^K \theta_{mk} - 1 \right))$$

Il est évident que $n_i \geq 0$ car $y_i \geq 0$ donc f est convexe comme somme de fonction convexe. De plus l'ensemble $\{\pi_m \geq 0, \theta_{mk} \geq 0, \forall m \in \llbracket 1, \dots, M \rrbracket \forall k \in \llbracket 1, \dots, K \rrbracket\}$ est convexe, il s'agit d'un problème d'optimisation

convexe. Les contraintes son linéaires, et **il existe $\pi_1, \pi_2, \dots, \pi_M$ tq $\sum_{i=1}^M \pi_i = 1$** , donc d'après la qualification de contraintes de Slater, le problème a la propriété de forte dualité. Ainsi :

$$\min_{\substack{\pi_m \geq 0 \\ \theta_{mk} \geq 0}} f((\pi, \theta)) = \max_{\lambda} \min_{\substack{\pi_m \geq 0 \\ \theta_{mk} \geq 0}} \mathcal{L}(\pi, \theta, \lambda)$$

Je suis pas sur que c'est phrase soit la bonne justification Comme $L(\pi, \lambda)$ est convexe par rapport à π , le minimum se trouve en annulant la dérivée de $L(\pi, \lambda)$ par rapport à π . Ainsi, on obtient :

$$\frac{\partial \mathcal{L}}{\partial \pi_i} = -\frac{K n_i}{\pi_i} + \lambda = 0 \quad \forall i \in \llbracket 1, \dots, M \rrbracket$$

Donc :

$$\pi_i \lambda = K n_i \quad \forall i \in \llbracket 1, \dots, M \rrbracket$$

En substituant cette égalité à la contrainte $\sum_{i=1}^M \pi_i = 1$, on obtient $\lambda = K \sum_{i=1}^M n_i$, d'où $\lambda = kN$ ($\neq 0$) avec N le nombre d'observations.

On obtient finalement :

$$\hat{\pi}_i = \frac{K n_i}{\lambda} = \frac{\sum_{j=1}^n y_i^{(j)}}{n} = \frac{n_i}{n} \quad \forall i \in \llbracket 1, \dots, M \rrbracket$$

Nous cherchons donc à trouver :

$$\hat{\pi}_{ML} = \operatorname{argmax}_{\pi_m} L(z_1, \dots, z_M; \pi_m)$$