

Examen Apprentissage Statistique  
Master Informatique UPMC – spécialité DAC

31-01-2017

Notes de cours autorisées

## Skip Gram avec embeddings gaussiens

1. En ½ page, le principe du modèle skip gram pour les embeddings de mots.

Au lieu de représenter les mots par des vecteurs, on va les représenter par une distribution Gaussienne : à chaque mot on associe une moyenne et une matrice de variance-covariance. Le modèle est plus riche et permet de représenter des incertitudes liées à la représentativité des données. On note  $w$  un terme du dictionnaire, et  $c(w)$  un contexte de  $w$ . Par la suite, ce contexte sera réduit à 1 terme uniquement. Au mot  $w_i$  on associera la représentation  $z_i \sim \mathcal{N}(z; \mu_i, \Sigma_i)$ ,  $z_i \in \mathbb{R}^n$ . Le modèle skip-gram emploie le produit scalaire pour mesurer la similarité de deux représentations vectorielles. Le produit scalaire sera remplacé ici par une similarité entre deux distributions. On considère par la suite deux alternatives pour cela.

2. Similarité symétrique : noyau produit de probabilités.

Pour deux distributions  $f$  et  $g$ , le noyau, qui étend le produit scalaire vectoriel, est défini par  $K(f, g) = \int_{z \in \mathbb{R}^n} f(z)g(z)dz$ . On note pour simplifier  $K(\mathcal{N}(z; \mu_i, \Sigma_i), \mathcal{N}(z; \mu_j, \Sigma_j)) = K(\mathcal{N}_i, \mathcal{N}_j)$ . On emploiera comme similarité  $S(\mathcal{N}_i, \mathcal{N}_j) = K(\mathcal{N}_i, \mathcal{N}_j)$ . On a alors (résultat admis) :

$$\log K(\mathcal{N}_i, \mathcal{N}_j) = -\frac{1}{2} \log |\Sigma_i + \Sigma_j| - \frac{1}{2} (\mu_i - \mu_j)^T (\Sigma_i + \Sigma_j)^{-1} (\mu_i - \mu_j) - \frac{n}{2} \log 2\pi$$

Où  $|\Sigma|$  est le déterminant de  $\Sigma$ .

Calculer  $\frac{\partial \log K(\mathcal{N}_i, \mathcal{N}_j)}{\partial \mu_i}$  et  $\frac{\partial \log K(\mathcal{N}_i, \mathcal{N}_j)}{\partial \Sigma_i}$  (pour la première on peut utiliser la dérivée d'un scalaire par rapport à un vecteur, pour la seconde, on passera par la dérivée de la fonction scalaire par rapport aux composantes  $\Sigma_{kl}$  de  $\Sigma$ ).

3. Similarité asymétrique : Distance de Kulback-Leibler

La deuxième alternative est la distance de Kulback-Leibler entre deux distributions :  $D_{KL}(\mathcal{N}_j || \mathcal{N}_i) = \int_z (\log \frac{\mathcal{N}_j(z)}{\mathcal{N}_i(z)}) \mathcal{N}_j(z) dz$ , on emploiera comme similarité  $S(\mathcal{N}_i, \mathcal{N}_j) = -D_{KL}(\mathcal{N}_j || \mathcal{N}_i)$

3.1 Montrer :

$$\int_z (\log \mathcal{N}_i(z)) \mathcal{N}_j(z) dz = -\frac{1}{2} (n \log 2\pi + \log |\Sigma_i| + \text{Tr}(\Sigma_i^{-1} \Sigma_j) + (\mu_i - \mu_j)^T \Sigma_i^{-1} (\mu_i - \mu_j))$$

3.2 En déduire :

$$D_{KL}(\mathcal{N}_j, \mathcal{N}_i) = \frac{1}{2} (\text{Tr}(\Sigma_i^{-1} \Sigma_j) + (\mu_i - \mu_j)^T \Sigma_i^{-1} (\mu_i - \mu_j) - n - \log \frac{|\Sigma_j|}{|\Sigma_i|})$$

3.3 Calculer  $\frac{\partial D_{KL}(\mathcal{N}_j, \mathcal{N}_i)}{\partial \mu_i}$ ,  $\frac{\partial D_{KL}(\mathcal{N}_j, \mathcal{N}_i)}{\partial \Sigma_i}$ ,  $\frac{\partial D_{KL}(\mathcal{N}_j, \mathcal{N}_i)}{\partial \Sigma_j}$

4. Algorithme d'apprentissage. Etant donné un terme  $w$  et un terme de contexte  $c$ , on va optimiser par gradient stochastique une fonction de coût qui a la forme suivante :

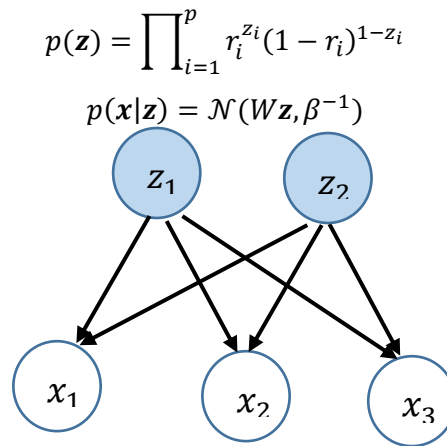
$L(w, c_p, c_n) = \max(0, 1 - S(w, c_p) + S(w, c_n))$  où  $c_p$  est un terme de contexte « positif », i.e. qui co-occure avec  $w$  dans le corpus et  $c_n$  est un terme de contexte négatif.  $S(\cdot)$  dénote l'une des deux similarités introduites ci-dessus et par abus de notation  $S(w, c)$  désigne la similarité des représentations correspondantes.

4.1 Expliquer cette fonction de coût.

4.2 Donner un algorithme de gradient stochastique pour chacune des deux similarités, qui optimise cette fonction de coût.

## Codage sparse binaire

On étudie un modèle génératif illustré sur la figure ci dessous. Une donnée  $x \in R^n$  est générée par le modèle suivant :



$z_i \in \{0,1\}$  suit une loi de Bernoulli de paramètre  $r_i$ ,  $\mathbf{z} \in \{0,1\}^p$ ,  $W \in R^{n \times p}$ ,  $\beta^{-1}$  est une matrice diagonale  $n \times n$  dite matrice de précision (c'est l'inverse de la matrice de covariance). On note que le modèle implique  $p(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^n p(x_i|\mathbf{z})$  (indépendance conditionnelle) et  $p(\mathbf{z}) = \prod_{i=1}^p p(z_i)$ . On veut estimer les paramètres du modèle par maximum de vraisemblance. Pour cela on a besoin de calculer les probabilités a posteriori  $p(\mathbf{z}|\mathbf{x})$  qui n'ont pas d'expression analytique simple. On va utiliser à la place une approximation variationnelle de champ moyen.

1. Optimiser la vraisemblance demande de connaître les  $p(\mathbf{z}|\mathbf{x})$ . Pour s'en convaincre on va regarder la dérivée de  $p(\mathbf{x})$  par rapport à un paramètre  $r_i$ . Montrer :

$$\frac{\partial}{\partial r_i} \log p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}) \frac{\frac{\partial}{\partial r_i} p(\mathbf{z})}{p(\mathbf{z})} = E_{\mathbf{z} \sim p(\mathbf{z}|\mathbf{x})} \left[ \frac{\partial}{\partial r_i} \log p(\mathbf{z}) \right]$$

2.  $p(\mathbf{z}|\mathbf{x})$  ne peut pas être calculé de façon analytique. L'approximation variationnelle consiste à l'approximer par une distribution simple  $q(\mathbf{z}|\mathbf{x})$  qui permettra d'aboutir à une expression analytique, et à optimiser une borne de la log vraisemblance. Montrer :

$$\log p_{\theta}(\mathbf{x}) = V_L(\theta, \phi; \mathbf{x}) + D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x}))$$

avec

$$V_L(\theta, \phi; \mathbf{x}) = E_{q(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]$$

Où  $\theta$  et  $\phi$  désignent respectivement les paramètres du modèle  $(\mathbf{r}, W, \beta^{-1})$  et de la distribution  $q$  (définis plus bas).

Comme  $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \geq 0$ ,  $V_L(\theta, \phi; \mathbf{x})$  constitue une borne inférieure de  $p_\theta(\mathbf{x})$  et c'est elle que l'on va optimiser.

Montrer que  $V_L(\theta, \phi; \mathbf{x})$  se réécrit sous la forme

$$V_L(\theta, \phi; \mathbf{x}) = -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) + E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$$

3. Pour obtenir une forme simple de  $q_\phi(\mathbf{z}|\mathbf{x})$ , on va faire l'hypothèse de champ moyen :

$$q_\phi(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^p q(z_i|\mathbf{x})$$

Puisque  $z_i \in \{0,1\}$ ,  $q(z_i|\mathbf{x})$  sera modélisé par une distribution de Bernoulli de paramètre  $q(z_i = 1|\mathbf{x}) = \hat{z}_i$

3.1 Montrer

$$V_L(\theta, \phi; \mathbf{x}) = E_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left[ \sum_{i=1}^p \log p(z_i) + \sum_{i=1}^n \log p(x_i|\mathbf{z}) - \sum_{i=1}^p \log q(z_i|\mathbf{x}) \right]$$

3.2 En introduisant les paramètres des distributions, montrer que cette expression se développe en :

$$V_L(\theta, \phi; \mathbf{x}) = \sum_{i=1}^p \hat{z}_i (\log r_i - \log \hat{z}_i) + (1 - \hat{z}_i) (\log(1 - r_i) - \log(1 - \hat{z}_i)) \\ + \frac{1}{2} \sum_{i=1}^n \left[ \log \frac{\beta_i}{2\pi} - \beta_i \left( x_i^2 - 2x_i W_{i,:} \hat{\mathbf{z}} + \sum_{j=1}^p \left( \sum_{k=1}^p W_{i,j} W_{i,k} \hat{z}_j \hat{z}_k \right) \right) \right]$$

Où  $W_{i,:}$  représente la ligne  $i$  de la matrice  $W$  et  $W_{i,j}$  l'élément  $i, j$

Cette expression montre que la borne a une expression que l'on peut optimiser.

3.3 proposer un algorithme pour apprendre les paramètres  $\theta$  et  $\phi$

3.4 A quoi sert ce modèle et comment l'utiliser ?

### Formules utiles

$A, B$  et  $X$  sont des matrices qui ont les bonnes dimensions pour se multiplier,  $x$  est un scalaire,  $Tr(A)$  est la trace de  $A$  et  $|A|$  son déterminant,  $\mathbf{x}$  est un vecteur,  $x$  est un scalaire.

### Trace de matrices

$$Tr(A + B) = Tr(A) + Tr(B), Tr(AB) = Tr(BA), x^T x = Tr(xx^T), x^T Ax = Tr(xAx^T) = Tr(Axx^T)$$

### Dérivation matricielle

La dérivée d'un scalaire par rapport à un vecteur est le vecteur des dérivées du scalaire par rapport à chacune des composantes (le gradient habituel)  $\left[\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_n}\right]$ . De même la dérivée d'un scalaire par rapport à une matrice est la matrice des dérivées de ce scalaire par rapport à chaque élément de la

matrice 
$$\begin{pmatrix} \frac{\partial y}{\partial x_{11}} & \dots & \frac{\partial y}{\partial x_{p1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{1q}} & \dots & \frac{\partial y}{\partial x_{pq}} \end{pmatrix}.$$

La dérivée d'un vecteur par rapport à un scalaire est le vecteur des dérivées des composantes par rapport à ce scalaire  $\left[\frac{\partial y_1}{\partial x}, \dots, \frac{\partial y_n}{\partial x}\right]^T$ . La dérivée d'une matrice par rapport à un scalaire est la matrice

des dérivées des composantes 
$$\begin{pmatrix} \frac{\partial y_{11}}{\partial x} & \dots & \frac{\partial y_{1n}}{\partial x} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_{m1}}{\partial x} & \dots & \frac{\partial y_{mn}}{\partial x} \end{pmatrix}.$$

### Dérivées de scalaires

On suppose que le scalaire dépend des coefficients de la matrice ou l'inverse suivant les cas.

$$\begin{aligned} \frac{\partial}{\partial A} Tr(AB) &= B^T, \frac{\partial}{\partial A} \ln|A| = (A^{-1})^T = (A^T)^{-1}, \frac{\partial \ln|A|}{\partial x} = Tr(A^{-1} \frac{\partial A}{\partial x}), \frac{\partial Tr(A)}{\partial x} = Tr(\frac{\partial A}{\partial x}), \frac{\partial \mathbf{x}^T A \mathbf{x}}{\partial x} = \\ (A + A^T)\mathbf{x}, \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial x} &= \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial x} = \mathbf{a}, \frac{\partial \mathbf{x}^T A^{-1} \mathbf{y}}{\partial A} = -(A^{-1})^T \mathbf{x} \mathbf{y}^T (A^{-1})^T, \frac{\partial}{\partial A} Tr(X^T A^{-1} Y) = -(A^{-1} Y X^T A^{-1})^T \end{aligned}$$

### Dérivées de matrices par rapport à un scalaire

On suppose que les coefficients de la matrice dépendent du scalaire.

$$\frac{\partial A^{-1}}{\partial x} = -A^{-1} \frac{\partial A}{\partial x} A^{-1}, \frac{\partial AB}{\partial x} = \frac{\partial A}{\partial x} B + A \frac{\partial B}{\partial x}$$

### Espérance et trace

$$E_x[A\mathbf{x}] = A E_x[\mathbf{x}], E_x[Tr(A\mathbf{x})] = Tr(A E_x(\mathbf{x}))$$

