

TP4 - solution

1. Le gradient de perte par rapport à  $a_i$  est donné en dérivant la formule suivante par rapport à  $W$  :

$$E_D(W) = -\sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_{\vec{w}_i}(\vec{x})$$

Sachant que  $y_{\vec{w}_i}(\vec{x}_n) = \frac{e^{a_i}}{\sum_c e^{a_c}}$  et  $a_i = \vec{w}_i^T \phi$

Nous pouvons trouver  $\frac{\partial E(w)}{\partial w}$  par la règle de dérivée en chaîne. Ainsi nous pouvons décortiquer le problème tel qu'illustré ici :

$$\frac{\partial E(w)}{\partial w} = \frac{\partial E(w)}{\partial y_{\vec{w}_i}(\vec{x}_n)} \frac{\partial y_{\vec{w}_i}(\vec{x}_n)}{\partial w} = \frac{\partial E(w)}{\partial y_{\vec{w}_i}(\vec{x}_n)} \frac{\partial y_{\vec{w}_i}(\vec{x}_n)}{\partial a_i} \frac{\partial a_i}{\partial w}$$

Et si nous simplifions, nous pouvons effectuer le travail suivant :

$$\frac{\partial E(w)}{\partial A} \cdot \frac{\partial A}{\partial B} \cdot \frac{\partial B}{\partial w}$$

Où :

$$A = y_{\vec{w}_i}(\vec{x})$$

$$\text{Et } B = a_i$$

Ainsi

$$1) \frac{\partial E(w)}{\partial A}$$

$$\frac{\partial E(w)}{\partial A} = \frac{\partial -\sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_{\vec{w}_i}(\vec{x}_n)}{\partial y_{\vec{w}_i}(\vec{x}_n)}$$

<=>

$$\frac{\partial -\sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln A}{\partial A}$$

par définition

<=>

$$-\sum_{n=1}^N \sum_{k=1}^K \frac{t_{kn}}{A}$$

par dérivée

$$2) \frac{\partial A}{\partial B}$$

$$\frac{\partial A}{\partial B} = \frac{\partial \frac{e^{a_i}}{\sum_c e^{a_c}}}{\partial a_i}$$

<=>

$$\frac{\partial \frac{e^B}{\sum_c e^B}}{\partial B}$$

par définition

<=>

$$\begin{aligned}
& \frac{e^B \sum_c e^B - e^B e^B}{(\sum_c e^B)^2} && \text{selon } \frac{\partial g(x)}{\partial x} = \frac{g'(x)f(x) - g(x)f'(x)}{f^2(x)} \\
\langle \Rightarrow \rangle & \frac{e^B (\sum_c e^B - e^B)}{\sum_c e^B} && \text{par associativité} \\
\langle \Rightarrow \rangle & \frac{e^B}{\sum_c e^B} \left( \frac{\sum_c e^B - e^B}{\sum_c e^B} \right) && \text{arithmétique} \\
\langle \Rightarrow \rangle & \frac{e^B}{\sum_c e^B} \left( \frac{\sum_c e^B}{\sum_c e^B} - \frac{e^B}{\sum_c e^B} \right) && \text{arithmétique} \\
\langle \Rightarrow \rangle & \frac{e^B}{\sum_c e^B} \left( I_n - \frac{e^B}{\sum_c e^B} \right) && \text{arithmétique} \\
\langle \Rightarrow \rangle & A(I_n - A) && \text{par définition} \\
3) \quad \frac{\partial B}{\partial w} & && \\
\langle \Rightarrow \rangle & \frac{\partial B}{\partial w} = \frac{\partial \bar{w}_i^T \vec{x}}{\partial w} && \\
& \vec{x} && \text{par dérivation}
\end{aligned}$$

Il ne reste qu'à remettre les solutions trouvés dans la définition de dérivée en chaine plus haut :

$$\begin{aligned}
& \frac{\partial E(w)}{\partial A} \cdot \frac{\partial A}{\partial B} \cdot \frac{\partial B}{\partial w} \\
\langle \Rightarrow \rangle & - \sum_{n=1}^N \sum_{k=1}^K \frac{t_{kn}}{A} \cdot A(I_n - A) \cdot \vec{x}_n && \text{par définition} \\
\langle \Rightarrow \rangle & - \sum_{n=1}^N \sum_{k=1}^K t_{kn} (I_n - A) \cdot \vec{x}_n && \text{arithmétique} \\
\langle \Rightarrow \rangle & \sum_{n=1}^N \sum_{k=1}^K t_{kn} (A - I_n) \cdot \vec{x}_n && \text{arithmétique} \\
\langle \Rightarrow \rangle & \sum_{n=1}^N \sum_{k=1}^K t_{kn} (y_{\bar{w}_i}(\vec{x}_n) - I_n) \cdot \vec{x}_n && \text{par définition} \\
\langle \Rightarrow \rangle & \sum_{n=1}^N \sum_{k=1}^K (t_{kn} y_{\bar{w}_i}(\vec{x}_n) - t_{kn}) \cdot \vec{x}_n && \text{par distributivité} \\
\langle \Rightarrow \rangle & \sum_{n=1}^N (y_{\bar{w}_i}(\vec{x}_n) - t_{ni}) \cdot \vec{x}_n && \text{par définition d'un « one-hot »} \\
& \text{CQFD}
\end{aligned}$$

2.

- a. Une distribution de vraisemblance est la distribution de probabilité conditionnelle d'observer un paramètre dans une classe donnée, ex :  $p(x|t)$  soit la probabilité d'observer  $x$  étant donné  $t$ . Pour la calculer, nous pourrions supposer que chaque classe de véhicule suit une distribution gaussienne et ainsi calculer une gaussienne par classe. Ainsi, dans chaque classe  $t$ , nous aurons une formule gaussienne qui nous permet d'en connaître la répartition  $x$ .
  - b. La distribution à priori est la distribution de probabilité d'observer une catégorie divisée par le nbr de véhicule total, ex :  $p(t)$ . Pour chaque catégorie, nous pourrions calculer le nombre de véhicule divisé par le nombre total de véhicule.
  - c. Oui, puisque les voitures sport sont de même catégorie et ont plus de chance de consommer tous beaucoup d'essence et peu de chance qu'une de ces voitures consomme très peu d'essence.
3. Si nous partons des principes de forces qui s'appliquent sur un corps en mouvement, nous avons les paramètres suivants : la masse d'un objet, la gravité qui agit sur celui-ci, la pression du sol sur ce corps, la pente du sol sur cette masse ainsi que la friction entre la masse et le sol. Ainsi, l'équilibre complexe ainsi formé peut être résumé dans les fonctions suivantes :

$$a = \frac{(v_t - v_{t-\delta})}{\delta}$$
$$v = \frac{(d_t - d_{t-\delta})}{\delta}$$
$$\delta = \frac{d_t - d_{t-\delta}}{v_t - v_{t-\delta}}$$

où

et où

a : accélération  
v : vitesse  
d : position  
 $\delta$  : variation de pente de vitesse

Le formule de descente de gradient de type momentum se définissent comme suit :

$$w_{t+1} = w_t - \eta \nabla E_{\vec{x}_n}(w_t)$$
$$v_{t+1} = \rho v_t - \nabla E_{\vec{x}_n}(w_t)$$
$$w_{t+1} = w_t - \eta v_{t+1}$$

Où

w : un certain poids  
v : vitesse  
 $\eta$  : constante de contrôle de vitesse initiale  
 $\rho$  : constante de contrôle de variation de vitesse  
t : temps donné  
 $\nabla E_{\vec{x}_n}$  : gradient

Pour mettre en valeur les liens entre la descente de gradient de type momentum et la position, vitesse et accélération, il suffit de prendre les équations de vitesse et d'accélération, d'isoler  $v_t$  et  $d_t$  et d'étudier les similarités. Nous obtiendrons les équations suivantes :

$$\begin{aligned}d_t &= d_{t-\delta} + v\delta \\v_t &= v_{t-\delta} + a\delta\end{aligned}$$

Par rapport aux formules du momentum suivantes :

$$\begin{aligned}w_{t+1} &= w_t - \eta v_{t+1} \\v_{t+1} &= \rho v_t - \nabla E_{\vec{x}_n}(w_t)\end{aligned}$$

- Les  $w_t$  tout comme  $d_t$  peuvent d'abord être interprétés comme la hauteur de la masse dans une pente en un temps  $t$ .
- $\eta$  peut être vue comme étant la variation de l'itération à laquelle on regarde la vitesse.
- $\rho$  peut être vu comme la friction qui ralentit la masse
- $\nabla E_{\vec{x}_n}(w_t)$  et l'accélération ( $a$ ) peuvent être vu comme la force qui fait varier la position de la masse dans la pente (la gravité).
- La différences des signes + et - s'explique par la direction d'application des variations de vitesses et de gradient sur la masse. Dans le gradient, + signifie que le poid remonte alors que dans la vitesse, un  $\delta +$  signifie que la masse accélère vers le bas et prend plus de distance.

C'était les similitudes entre la descente de gradient de type momentum et la position, vitesse et accélération.