

RÉSUMÉ THÉORIQUE

1 Préalables

- Dériver les équations normales de la régression linéaire au moyen de l'algèbre linéaire

On cherche \vec{x} qui solutionne $\mathbf{A}\vec{x} = \vec{b}$. Comme l'équation n'admet pas de solution pour \vec{x} , on cherche plutôt $\vec{\hat{x}}$ qui solutionne $\mathbf{A}\vec{\hat{x}} = \vec{p}$, avec l'objectif de minimiser l'erreur $\vec{e} = \vec{p} - \vec{b}$. Puisque $\vec{p} \in C(\mathbf{A})$ et que \vec{e} est minimal quand \vec{e} est perpendiculaire à $C(\mathbf{A})$, on conclut, par le théorème fondamental de l'algèbre linéaire, que $\vec{e} \in N(\mathbf{A}^T)$:

$$\begin{aligned}\mathbf{A}^T \vec{e} &= \vec{0} \\ \mathbf{A}^T (\mathbf{A}\vec{\hat{x}} - \vec{b}) &= \vec{0} \\ \vec{\hat{x}} &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \vec{b}\end{aligned}$$

- Démontrer que $\mathbf{A}^T \mathbf{A}$ est inversible ssi les colonnes de \mathbf{A} sont indépendantes

Si $\mathbf{A}\vec{x} = \vec{0}$, alors $\mathbf{A}^T \mathbf{A}\vec{x} = \vec{0}$; si $\mathbf{A}^T \mathbf{A}\vec{x} = \vec{0}$ alors $\vec{x}^T \mathbf{A}^T \mathbf{A}\vec{x} = 0$, $\|\mathbf{A}\vec{x}\| = 0$, $\vec{x} = \vec{0}$: $\mathbf{A}^T \mathbf{A}$ et \mathbf{A} ont le même espace nul. Si les colonnes de \mathbf{A} sont indépendantes, $N(\mathbf{A})$ ne contient que le vecteur nul, $N(\mathbf{A}^T \mathbf{A})$ ne contient que le vecteur nul, les colonnes de $\mathbf{A}^T \mathbf{A}$ sont indépendantes et donc $\mathbf{A}^T \mathbf{A}$ est inversible.

- Écrire les trois premiers termes de l'expansion en série de Taylor de $f(\vec{x})$ où $\vec{x} \in \mathbb{R}^n$

$$f(\vec{x}) = f(\vec{x}_0) + (\vec{x} - \vec{x}_0)^T \left. \nabla f(\vec{x}) \right|_{\vec{x}_0} + \frac{1}{2} (\vec{x} - \vec{x}_0)^T \left. H_f(\vec{x}) \right|_{\vec{x}_0} (\vec{x} - \vec{x}_0) + \dots$$

$\left. \nabla f(\vec{x}) \right|_{\vec{x}_0}$: gradient de $f(\vec{x})$ évalué à \vec{x}_0

$\left. H_f(\vec{x}) \right|_{\vec{x}_0}$: matrice hessienne de $f(\vec{x})$ évaluée à \vec{x}_0

- Exprimer la *distance de Mahalanobis* de $\vec{\mu}$ à \vec{x} en termes des valeurs propres de Σ

On veut montrer : $\Delta^2 = (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})$

Soit \vec{u}_i tel que $\Sigma \vec{u}_i = \lambda_i \vec{u}_i$

$$\Sigma = \sum_{i=1}^D \lambda_i \vec{u}_i \vec{u}_i^T$$

$$\Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \vec{u}_i \vec{u}_i^T$$

Notons $y_i = \vec{u}_i^T (\vec{x} - \vec{\mu})$

$$= (\vec{x} - \vec{\mu})^T \vec{u}_i$$

$$\Delta^2 = (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})$$

$$= (\vec{x} - \vec{\mu})^T \left(\sum_{i=1}^D \frac{1}{\lambda_i} \vec{u}_i \vec{u}_i^T \right) (\vec{x} - \vec{\mu})$$

$$= \sum_{i=1}^D \frac{1}{\lambda_i} (\vec{x} - \vec{\mu})^T \vec{u}_i \vec{u}_i^T (\vec{x} - \vec{\mu})$$

$$= \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$$

2 Théorie de l'information

- Démontrer l'inégalité de Gibbs au moyen de l'inégalité de Jensen

Inégalité de Gibbs : $D_{KL}(P||Q) \geq 0$

$$\begin{aligned} D_{KL}(P||Q) &= \sum_{x \in \mathcal{A}_X} P(x) \log \frac{P(x)}{Q(x)} \\ &= \sum_{x \in \mathcal{A}_X} P(x) \left(-\log \frac{Q(x)}{P(x)} \right) \\ &= E \left[-\log \frac{Q(X)}{P(X)} \right] \end{aligned}$$

$-\log$ est une fonction convexe

par l'inégalité de Jensen, $E[f(X)] \geq f(E[X])$:

$$\begin{aligned} D_{KL}(P||Q) &\geq -\log \left(E \left[\frac{Q(X)}{P(X)} \right] \right) \\ &= -\log \left(\sum_{x \in \mathcal{A}_X} P(x) \frac{Q(x)}{P(x)} \right) \\ &= -\log \left(\sum_{x \in \mathcal{A}_X} Q(x) \right) = -\log(1) \end{aligned}$$

$$D_{KL}(P||Q) \geq 0$$

3 Régression linéaire

- Pour le modèle de la régression linéaire, énoncer les équations pour : la distribution de la variable cible ; les distributions *a priori* et *a posteriori* du vecteur des pondérations ; la distribution prédictive de la variable cible dans le cas où $\mathbf{S}_0 = \alpha \mathbf{I}$ et $\vec{m}_0 = 0$

$$\begin{aligned}
 p(t | \vec{x}, \vec{w}) &= \mathcal{N}(t | y(\vec{x}, \vec{w}), \beta^{-1}) \\
 p(\vec{w}) &= \mathcal{N}(\vec{w} | \vec{m}_0, \mathbf{S}_0) \\
 p(\vec{w} | \vec{t}, \mathbf{X}) &= \mathcal{N}(\vec{w} | \vec{m}_N, \mathbf{S}_N) \\
 \text{où } \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta \Phi^T \Phi \\
 \vec{m}_N &= \mathbf{S}_N (\mathbf{S}_0^{-1} \vec{m}_0 + \beta \Phi^T \vec{t}) \\
 p(t | \vec{x}, \vec{t}, \mathbf{X}) &= \mathcal{N}(t | y(\vec{x}, \vec{m}_N), \sigma_N^2(\vec{x})) \\
 \text{où } \sigma_N^2(\vec{x}) &= \frac{1}{\beta} + \vec{\phi}(\vec{x})^T \mathbf{S}_N \vec{\phi}(\vec{x})
 \end{aligned}$$

- Établir la correspondance entre le paramètre de régularisation dans la fonction d'erreur des moindres carrés et le paramètre de précision de la distribution a priori du vecteur des pondérations

$$\begin{aligned}
 p(\vec{w} | \vec{t}) &\propto \prod_{n=1}^N p(t_n | \vec{x}_n, \vec{w}) p(\vec{w}) \\
 &= \prod_{n=1}^N \mathcal{N}(t_n | y_n, \beta^{-1}) \mathcal{N}(\vec{w} | 0, \alpha^{-1}) \text{ où } y_n = y(\vec{x}_n, \vec{w}) \\
 &= k \prod_{n=1}^N \exp\left(-\frac{\beta}{2}(t_n - y_n)^2\right) \exp\left(-\frac{\alpha}{2}\vec{w}^T \vec{w}\right) \\
 E(\vec{w}) &= -\ln p(\vec{w} | \vec{t}) = \frac{\beta}{2} \sum_{n=1}^N (t_n - y_n)^2 + \frac{\alpha}{2} \vec{w}^T \vec{w} + K \\
 E(\vec{w}) \text{ est minimum quand } &\frac{\beta}{2} \sum_{n=1}^N (t_n - y_n)^2 + \frac{\alpha}{2} \vec{w}^T \vec{w} \text{ est minimum} \\
 \text{quand } &\frac{1}{2} \sum_{n=1}^N (t_n - y_n)^2 + \frac{1}{2} \frac{\alpha}{\beta} \vec{w}^T \vec{w} \text{ est minimum}
 \end{aligned}$$

Le paramètre de régularisation λ est donc équivalent à $\frac{\alpha}{\beta}$

- Au moyen d'une notation sans omission, formuler la triple intégrale de la distribution *prédictive* obtenue dans le cadre d'un traitement bayésien complet (portant sur \vec{w} , α , β) du modèle de régression linéaire

$$p(t | \vec{x}, \vec{t}, \mathbf{X}) = \iiint p(t | \vec{x}, \vec{w}, \beta) p(\vec{w} | \vec{t}, \mathbf{X}, \alpha) p(\alpha, \beta | \vec{t}, \mathbf{X}) d\vec{w} d\alpha d\beta$$

4 Réseaux neuronaux

- Énoncer la relation qui lie une sortie y_k au vecteur d'entrée \vec{x} (à D composantes) dans un réseau neuronal à une couche cachée composée de M unités

$$y_k = f(a_k) = f\left(\sum_{j=0}^M w_{kj}^{(2)} \sigma\left(\sum_{i=0}^D w_{ji}^{(1)} x_i\right)\right)$$

- Dans un réseau neuronal, pour les modèle de régression, de classification binaire et de classification multinomiale, énoncer : le modèle pour la variable cible ; la formulation de y_k d'une unité de sortie en fonction de son activation a_k ; la fonction d'erreur

$$p(t|\vec{x}, \vec{w}) = \mathcal{N}(t|y(\vec{x}, \vec{w}), \beta^{-1})$$

$$y(\vec{x}, \vec{w}) = f(a) = a$$

$$E(\vec{w}) = \frac{\beta}{2} \sum_{n=1}^N \{y(\vec{x}_n, \vec{w}) - t_n\}^2$$

$$p(t|\vec{x}, \vec{w}) = y(\vec{x}, \vec{w})^t \{1 - y(\vec{x}, \vec{w})\}^{1-t}$$

$$y(\vec{x}, \vec{w}) = f(a) = \sigma(a) = \frac{1}{1 + \exp(-a)}$$

$$E(\vec{w}) = - \sum_{n=1}^N t_n \ln \{y(\vec{x}_n, \vec{w})\} + (1 - t_n) \ln \{1 - y(\vec{x}_n, \vec{w})\}$$

$$p(\vec{t}|\vec{x}, \vec{w}) = \prod_{k=1}^K \{y_k(\vec{x}, \vec{w})\}^{t_k}$$

$$y_k(\vec{x}, \vec{w}) = f(a_k) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$E(\vec{w}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln \{y_k(\vec{x}_n, \vec{w})\}$$

où t_{nk} est la k^e composante de la cible \vec{t} de la n^e observation

- Pour un réseau neuronal à K sorties, avec une couche cachée composée de M unités, dériver les équations de $\nabla_{\vec{w}} E_n$ calculé par rétropropagation

$$\begin{aligned}\frac{\partial E_n}{\partial w_{ji}^{(l)}} &= \frac{\partial E_n}{\partial a_j^{(l)}} \frac{\partial a_j^{(l)}}{\partial w_{ji}^{(l)}} = \frac{\partial E_n}{\partial a_j^{(l)}} z_i^{(l-1)} \text{ où } z_i^{(1)} = h(a_i^{(1)}), z_i^{(0)} = x_i \\ &= \delta_j^{(l)} z_i^{(l-1)} \text{ où } \delta_j^{(l)} \equiv \frac{\partial E_n}{\partial a_j^{(l)}}\end{aligned}$$

$\delta_k^{(2)}$ se calcule directement à partir de la fonction d'erreur.

Si la fonction f dans $y_k = f(a_k)$ est l'inverse de la fonction de liaison canonique de la distribution de $t|\vec{x}$:

$$\delta_k^{(2)} = y_k - t_k$$

$$\begin{aligned}\text{Quelle que soit la forme de la fonction } f : \delta_j^{(1)} &= \frac{\partial E_n}{\partial a_j^{(1)}} = \sum_{k=1}^K \frac{\partial E_n}{\partial a_k^{(2)}} \frac{\partial a_k^{(2)}}{\partial a_j^{(1)}} \\ &= \sum_{k=1}^K \delta_k^{(2)} \frac{\partial a_k^{(2)}}{\partial a_j^{(1)}}\end{aligned}$$

$$\text{De la somme } a_k^{(2)} = w_{k0}^{(2)} h(a_0^{(1)}) + w_{k1}^{(2)} h(a_1^{(1)}) + w_{k2}^{(2)} h(a_2^{(1)}) + \dots$$

$$\text{un seul terme lie } a_k^{(2)} \text{ et } a_j^{(1)} : w_{kj}^{(2)} h(a_j^{(1)})$$

$$\frac{\partial a_k^{(2)}}{\partial a_j^{(1)}} = w_{kj}^{(2)} \frac{dh(a)}{da} \Big|_{(a_j^{(1)})}$$

$$\delta_j^{(1)} = h'(a_j^{(1)}) \sum_{k=1}^K w_{kj}^{(2)} \delta_k^{(2)}$$

5 Méthodes noyau

- Résumer la *représentation noyau* du modèle de régression linéaire des moindres carrés régularisés (Ridge) en indiquant : la définition reliant le nouveau vecteur des paramètres \vec{a} et l'ancien vecteur \vec{w} ; la fonction noyau ; la relation entre \vec{w} et \vec{a} qui produit une valeur minimum pour la fonction d'erreur ; la solution de \vec{a} qui minimise la fonction d'erreur ; la solution de régression $y(\vec{x})$ en termes de la fonction noyau et du vecteur des paramètres \vec{a}

$$\Phi \vec{w} = \vec{t} - \lambda \vec{a}$$

$$k(\vec{x}_n, \vec{x}_m) = \vec{\phi}(\vec{x}_n)^T \vec{\phi}(\vec{x}_m)$$

$$\vec{w} = \Phi^T \vec{a}$$

$$\vec{a} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \vec{t}$$

$$\text{où } \mathbf{K} = \Phi \Phi^T$$

$$= K_{nm} = k(\vec{x}_n, \vec{x}_m)$$

$$y(\vec{x}) = \vec{a}^T \vec{k}(\vec{x})$$

$$\text{où } \vec{k}(\vec{x}) = \Phi \vec{\phi}(\vec{x})$$

$$k_n(\vec{x}) = k(\vec{x}_n, \vec{x})$$

- Donner une formulation complète du modèle de Nadaraya-Watson

Au moyen d'une densité de probabilité modélisée par un estimateur de Parzen :

$$p(\vec{x}, t) = \frac{1}{N} \sum_{n=1}^N f(\vec{x} - \vec{x}_n, t - t_n)$$

on déduit la fonction de régression $y(\vec{x}) = \mathbb{E}[t | \vec{x}]$

$$y(\vec{x}) = \sum_{n=1}^N k(\vec{x}, \vec{x}_n) t_n$$

$$\text{où } k(\vec{x}, \vec{x}_n) = \frac{g(\vec{x} - \vec{x}_n)}{\sum_{m=1}^N g(\vec{x} - \vec{x}_m)}$$

$$g(\vec{x}) \equiv \int_{-\infty}^{\infty} f(\vec{x}, t) dt$$

- Résumer le modèle de la régression linéaire par *processus gaussien* en indiquant : les trois hypothèses du modèle ; les trois équations intermédiaires qui relèvent de l'approche par processus gaussien ; la prédictive :

$$\text{modèle : } y(\vec{x}) = \vec{w}^T \vec{\phi}(\vec{x})$$

$$p(\vec{w}) = \mathcal{N}(\vec{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

$$p(t | y) = \mathcal{N}(t | y, \beta^{-1})$$

$$\text{approche processus gaussien : } p(\vec{y}) = \mathcal{N}(\vec{y} | \mathbf{0}, \mathbf{K})$$

$$\text{où } \mathbf{K} = \frac{1}{\alpha} \mathbf{\Phi} \mathbf{\Phi}^T$$

$$K_{nm} = \frac{1}{\alpha} \vec{\phi}(\vec{x}_n)^T \vec{\phi}(\vec{x}_m) \\ = k(\vec{x}_n, \vec{x}_m)$$

$$p(\vec{t} | \vec{y}) = \mathcal{N}(\vec{t} | \vec{y}, \beta^{-1} \mathbf{I})$$

$$p(\vec{t}) = \mathcal{N}(\vec{t} | \mathbf{0}, \mathbf{C})$$

$$\text{où } \mathbf{C} = \mathbf{K} + \beta^{-1} \mathbf{I}$$

$$C_{nm} = k(\vec{x}_n, \vec{x}_m) + \beta^{-1} \delta_{nm}$$

$$\text{prédictive : } p(t | \vec{x}, \vec{t}, \mathbf{X}) = \mathcal{N}(t | m(\vec{x}), \sigma^2(\vec{x}))$$

$$\text{où } m(\vec{x}) = \vec{k}^T \mathbf{C}^{-1} \vec{t}$$

$$\sigma^2(\vec{x}) = c - \vec{k}^T \mathbf{C}^{-1} \vec{k}$$

$$k_n = k(\vec{x}_n, \vec{x})$$

$$c = k(\vec{x}, \vec{x}) + \beta^{-1}$$