Notes de Cours de Machine
Learning     M1 MINT

---

Slide 2 : régression logistique et softmax.

Logistique   $\mathbb{P}(Y=1/X=x) = \sigma(\langle x,w\rangle + b)$

Softmax   $\mathbb{P}(Y=k/X=x) = \dfrac{\exp(\langle x,w_k\rangle + b_k)}{\sum_{k=1}^{K} \exp(\langle x,w_k\rangle + b_k)}$

$\sum_{k=1}^{K} \mathbb{P}(Y=k/X=x) = 1$

Cas où $K=2$ classification binaire.

$\mathbb{P}(Y=1/X=x) = \dfrac{\exp(\langle x,w_1\rangle + b_1)}{\exp(\langle x,w_{-1}\rangle + b_{-1}) + \exp(\langle x,w_1\rangle + b_1)}$

$= \dfrac{\exp(\langle x,w_1\rangle + b_1 - \langle x,w_{-1}\rangle - b_{-1})}{1 + \exp(\langle x,w_1\rangle + b_1 - \langle x,w_{-1}\rangle - b_1)}$

$= \dfrac{\exp(\langle x,w_1 - w_{-1}\rangle + b_1 - b_1)}{1 + \exp(\langle x,w_1 - w_{-1}\rangle + b_1 - b_{-1})}$

Si on pose $w = w_1 - w_{-1}$ et $b = b_1 - b_{-1}$

$= \dfrac{\exp(\langle x,w\rangle + b)}{1 + \exp(\langle x,w\rangle + b)} = \sigma(\langle x,w\rangle + b).$

Conclusion : la softmax généralise la sigmoïde.

Slide 22

$$\frac{\mathbb{P}(Y=1 \mid X=x)}{\mathbb{P}(Y=-1 \mid X=x)} = \text{odds}.$$

$$\mathbb{P}(Y=1) = \frac{1}{4}$$
$$\mathbb{P}(Y=-1) = 3/4$$
$$\text{odds} = \frac{1/4}{3/4} = \frac{1}{3}.$$

$$\mathbb{P}(Y=1) = 1 \quad \mathbb{P}(Y=-1) = 0$$
$$\text{odds} = +\infty$$

$$\mathbb{P}(Y=1) = \frac{1}{2}$$
$$\mathbb{P}(Y=-1) = \frac{1}{2}$$

$$\mathbb{P}(Y=1) = 0 \quad \mathbb{P}(Y=-1) = 1 \qquad \text{odds} = 1$$
$$\text{odds} = 0.$$

Dans le cas de la régression logistique

$$\mathbb{P}(Y=1 \mid X=x) = \sigma(\langle x, w \rangle + b). \quad \sigma : x \to \frac{e^x}{1+e^x}$$
$$= \frac{1}{1+e^{-x}}$$

$$\text{Odds} = \frac{\sigma(\langle x, w \rangle + b)}{1 - \sigma(\langle x+w \rangle + b)}$$

$$= \frac{e^{\langle x, w \rangle + b} / 1+e^{\langle x, w \rangle + b}}{1 - \frac{e^{\langle x, w \rangle + b}}{1 + e^{\langle x, w \rangle + b}}}$$

$$= \frac{\dfrac{e^{\langle x, w \rangle + b}}{1+e^{\langle x, w \rangle + b}}}{\dfrac{1+e^{\langle x, w \rangle + b} - e^{\langle x, w \rangle + b}}{1+e^{\langle x, w \rangle + b}}}$$

$$= e^{\langle x, w \rangle + b}$$

$i_1$ et $i_2$ dont les features sont égales
sauf la $j$   $x_{i_1}^j = x_{i_2}^j + 1$
$$x_{i_1}^k = x_{i_2}^k \qquad \forall k \in [\![1, d]\!] \setminus \{j\}$$

$$\frac{\text{odds}(i_1)}{\text{odds}(i_2)} = \frac{\exp(\langle x_{i_1}, w\rangle + b)}{\exp(\langle x_{i_2}, w\rangle + b)}$$

$$= \exp\left(\langle x_{i_1}, w\rangle + b - \langle x_{i_2}, w\rangle - b\right)$$

$$= \exp\left(\langle x_{i_1} - x_{i_2}, w\rangle\right) = \exp\left(\sum_{k=1}^{d} (x_{i_1}^k - x_{i_2}^k) w_k\right)$$

$$= \exp\left((x_{i_1}^{j} - x_{i_2}^{j}) w_j\right) = \exp(w_j).$$

on choisit $\hat{Y}_+ = 1$

si $\quad \mathbb{P}(Y=1 / X=x) \geqslant \mathbb{P}(Y=-1/X=x)$

$\Longleftrightarrow \quad \dfrac{e^{\langle x, w\rangle + b}}{1 + e^{\langle x, w\rangle + b}} \geqslant \dfrac{1}{1 + e^{\langle x, w\rangle + b}}$

$\Longleftrightarrow \quad e^{\langle x, w\rangle + b} \geqslant 1$

$\Longleftrightarrow \quad \langle x, w\rangle + b \geqslant 0$

$\Longleftrightarrow \quad \langle x, w\rangle \geqslant -b \qquad \longrightarrow$ règle de classification linéaire

l'espace des features est partagé par un hyperplan d'équation $\langle x, w\rangle + b = 0$.

On veut calculer $-\frac{1}{n} \log \mathcal{L}$ ← vraisemblance.

Vraisemblance $= \mathcal{L}(y_1, \ldots, y_n, x_1, \ldots, x_n ; w, b) = \mathcal{L}$

$$= \prod_{i=1}^{n} \mathbb{P}(Y = y_i / X = x_i)$$

$$\mathbb{P}(Y = 1 / X = x_i) = \sigma(\langle x_i, w\rangle + b) = \frac{e^{\langle x_i, w\rangle + b}}{1 + e^{\langle x_i, w\rangle + b}}$$

$$\mathbb{P}(Y = -1 / X = x_i) = 1 - \sigma(\langle x_i, w\rangle + b)$$

$$= \frac{1}{1 + e^{\langle x_i, w\rangle + b}}$$

Probit
$\phi(\langle x_i, w\rangle + b)$

$$= \prod_{i=1}^{n} \left( \frac{e^{\langle x_i, w\rangle + b}}{1 + e^{\langle x_i, w\rangle + b}} \right)^{\mathbb{1}(y_i = 1)} \left( \frac{1}{1 + e^{\langle x_i, w\rangle + b}} \right)^{\mathbb{1}(y_i = -1)}$$

$$= \prod_{i=1}^{n} \left( \frac{1}{e^{-(\langle x_i, w\rangle + b)} + 1} \right)^{\mathbb{1}(y_i = 1)} \left( \frac{1}{1 + e^{\langle x_i, w\rangle + b}} \right)^{\mathbb{1}(y_i = -1)}$$

$$= \prod_{i=1}^{n} \frac{1}{1 + e^{-y_i(\langle x_i, w\rangle + b)}}$$

$$\log \mathcal{L} = \sum_{i=1}^{n} \log \frac{1}{1 + e^{-y_i(\langle x_i, w\rangle + b)}}$$

$$= - \sum_{i=1}^{n} \log \left( 1 + e^{-y_i \langle x_i, w\rangle + b} \right)$$

Empirical loss $= -\frac{1}{n} \log \mathcal{L} = \frac{1}{n} \sum_{i=1}^{n} \log \left( 1 + e^{-y_i \langle x_i, w\rangle + b} \right)$

Définir les "bons" $\hat{\omega}$ et $\hat{b}$ au max. de vraisemblance comme

$$(\hat{\omega}, \hat{b}) \in \underset{\omega \in \mathbb{R}^d, b \in \mathbb{R}}{\text{argmin}} \frac{1}{n} \sum_{i=1}^{n} \log\left(1 + e^{-y_i(\langle x_i, \omega \rangle + b)}\right)$$

problème d'optimisation $\overset{\text{strict.}}{\vee}$ convexe et différentiable.

$$\nabla f(\omega) = \frac{1}{n} \sum_{i=1}^{n} \ell'(y_i, \langle x_i, \omega \rangle) x_i$$

$$\nabla^2 f(\omega) = \frac{1}{n} \sum_{i=1}^{n} \ell''(y_i, \langle x_i, \omega \rangle) x_i x_i^\top$$

$f$ est convexe $\Longleftrightarrow$ $\ell$ l'est pour tout $y$. $\quad y' \mapsto \ell(y, y')$

$L$ - régulière

$$\|\nabla f(\omega) - \nabla f(\omega')\|_2 \leq L \|\omega - \omega'\|_2.$$

si $f$ est 2 fois diff.

$L$ - régulière $\Longleftrightarrow$ $\sigma_{max}(\nabla^2 f(\omega)) \leq L$.

logistique $\ell(y_i, \langle x_i, \omega \rangle) = \log\left(1 + e^{-y_i \langle x_i, \omega \rangle}\right)$

$$\ell'(y_i, \langle x_i, \omega \rangle) = \frac{\overset{\textcircled{-}}{} y_i x_i e^{\langle x_i, \omega \rangle}}{1 + e^{-y_i \langle x_i, \omega \rangle}}$$

$$= y_i \left( \sigma(y_i \langle x_i, \omega \rangle) - 1 \right) x_i$$

$$\ell''\left(y_i, \langle x_i, \omega \rangle\right) = y_i \left(\sigma'\left(y_i \langle x_i, \omega \rangle\right)\right) \quad x_i$$

$$\sigma(x) = \frac{e^x}{1+e^x}$$

$$\sigma'(x) = \frac{e^x(1+e^x) - e^x e^x}{(1+e^x)^2}$$

$$= \frac{e^x + e^{2x} - e^{+2x}}{(1+e^x)^2} = \frac{e^x}{1+e^x} \cdot \frac{1}{1+e^x}$$

$$= \sigma(x)\left(1 - \sigma(x)\right)$$

$$= \underset{\underset{1}{\shortparallel}}{y_i^2} \; \sigma\left(y_i \langle x_i, \omega \rangle\right)\left(1 - \sigma\left(y_i \langle x_i, \omega \rangle\right)\right) x_i x_i^T$$

$$\nabla^2 f(\omega) = \frac{1}{n} \sum_{i=1}^{n} \ell''\left(y_i, \langle x_i, \omega \rangle\right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \underbrace{\sigma\left(y_i \langle x_i, \omega \rangle\right)\left(1 - \sigma\left(y_i \langle x_i, \omega \rangle\right)\right)}_{\sigma\left(y_i \langle x_i, \omega \rangle\right) = \mathbb{P}\left(Y = y_i / X = x_i\right)} x_i x_i^T$$

$x \longmapsto \sigma(x)\left(1 - \sigma(x)\right)$ son maximum vaut $\frac{1}{4}$

$$\lambda_{max}\left(\nabla^2 f(\omega)\right) \le \frac{1}{4n} \underbrace{\lambda_{max}\left(\sum_{i=1}^{n} x_i x_i^T\right)}_{\phantom{x}}$$

$$\llcorner \text{logistic.}$$

Si $f$ est $L$-régulière.
$$f(w) \leq f(w') + \langle \nabla f(w'), w - w' \rangle + \frac{L}{2} \|w - w'\|_2^2.$$

A l'itération:
$$f(w) \leq f(w^k) + \langle \nabla f(w^k), w - w^k \rangle + \frac{L}{2} \|w - w^k\|_2^2.$$

descente de gradient $w^{k+1} = \text{argmin} \;(1)$

On veut minimiser $f(w) + g(w)$

$$f(w) + g(w) \leq f(w^k) + \langle \nabla f(w^k), w - w^k \rangle + \frac{L}{2} \|w - w^k\|_2^2$$
$$+ g(w)$$

$$w^{k+1} = \underset{w}{\text{argmin}} \; \cancel{f(w^k)} + \langle \nabla f(w^k), w - w^k \rangle$$
$$+ \frac{L}{2} \underbrace{\|w - w^k\|_2^2} + g(w)$$

$$\frac{L}{2} \| w - (w^k - \frac{1}{L} \nabla f(w^k)) \|_2^2$$
$$= \frac{L}{2} \|w - w^k\|_2^2 \; \cancel{+} \; \cancel{2} \frac{\cancel{L}}{\cancel{2}} \langle w - w^k, \frac{1}{L} \nabla f(w^k) \rangle$$
$$+ \frac{L}{2} \frac{1}{L^2} \| \nabla f(w)^k \|_2^2$$

$$w^{k+1} = \underset{w}{\text{argmin}} \; \frac{L}{2} \| w - (w^k - \frac{1}{L} \nabla f(w^k)) \|_2^2 + g(w).$$
$$= \underset{w}{\text{argmin}} \; \frac{1}{2} \| w - (w^k - \frac{1}{L} \nabla f(w^k)) \|_2^2 + \frac{1}{L} g(w)$$

Si $g = 0$
on retombe sur la
descente de gradient

$g : \mathbb{R}^d \to \mathbb{R}$ convexe (pas forcément diff.)
on définit son opérateur proximal.

$$\text{prox}_g(w) = \underset{w' \in \mathbb{R}^d}{\arg\min} \left\{ \frac{1}{2} \|w - w'\|_2^2 + g(w') \right\}$$

: prox du LASSO.

$z \in \mathbb{R}, \quad z' \in \mathbb{R}$

$$\underset{z'}{\arg\min} \quad \frac{1}{2}(z' - z)^2 + d|z'|$$

sur $\mathbb{R}_+$   $z' - z + d$.   le minimum est atteint pour $z' = z - d$.

sur $\mathbb{R}_-$   $z' - z - d$.

le minimum est atteint $z' = z + d$.

ce minimum est $\leq 0$. seulement si $z \leq -d$.

ce minimum est $\geq 0$ seulement si $z \geq d$.

Sur $[-d, d]$ le minimum est atteint en $0$.

$$z^* = \underset{}{\arg\min} \frac{1}{2}(z' - z)^2 + d|z'| = \begin{cases} z - d & \text{si } z \geq d \\ z + d & \text{si } z \leq -d \\ 0 & \text{si } z \in [-d, d] \end{cases}$$

$$\boxed{z^*(z) = \text{sign}(z)\,(|z| - d)_+}$$

si $z \geq d$.   $z^*(z) = 1 \cdot (z - d)_+ = z - d$

si $z \leq -d$   $z^*(z) = (-1)(-z - d)_+ = z + d$.

si $z \in [-d, d]$   $z^*(z) = \binom{1}{-1}\underset{\underset{\leq 0}{\leq d}}{(|z| - d)_+} = 0$.

A une itération $k$ $w^k, b^k$.

$$f(w,b) \leq f(w^k, b^k) + \left\langle \nabla_{w,b} f(w^k, b^k), \begin{pmatrix} w - w^k \\ b - b^k \end{pmatrix} \right\rangle$$
$$+ g(w)$$
$$+ \frac{L}{2} \left\| \begin{pmatrix} w - w^k \\ b - b^k \end{pmatrix} \right\|_2^2 + g(w) \quad \textcircled{$\bigstar$}$$

$$\nabla_{w,b} f(w^k, b^k) = \begin{pmatrix} \nabla_w f(w^k, b^k) \\ \nabla_b f(w^k, b^k) \end{pmatrix} \in \mathbb{R}^{d+1}$$

qu'on minimise
pour obtenir
$w^{k+1}, b^{k+1}$.

$w \in \mathbb{R}^d, b \in \mathbb{R}$

$$\begin{pmatrix} w \\ b \end{pmatrix} = \begin{pmatrix} w_1 \\ \vdots \\ w_d \\ b \end{pmatrix} \in \mathbb{R}^{d+1} \qquad \left\| \begin{pmatrix} w \\ b \end{pmatrix} \right\|_2^2 = \sum_{j=1}^{d} w_j^2 + b^2.$$
$$= \|w\|_2^2 + b^2.$$

$\textcircled{$\bigstar$} = \left\langle \nabla_w f(w^k, b^k), w - w^k \right\rangle + \nabla_b f(w^k, b^k)(b - b^k)$

$$+ \frac{L}{2} \|w - w^k\|_2^2 + \frac{L}{2}(b - b^k)^2 + g(w).$$

$$w^{k+1} = \arg\min_{w} \left\langle \nabla_w f(w^k, b^k), w - w^k \right\rangle + \frac{L}{2} \|w - w^k\|_2^2$$
$$+ g(w)$$

$$b^k = \arg\min_{b} \nabla_b f(w^k, b^k)(b - b^k) + \frac{L}{2}(b - b^k)^2$$

Conclusion: l'update $w^{k+1}$ est donné
par la descente de
gradient proximale

l'update $b^{k+1}$ est donné
par la descente de gradient.