

Lecture 15

Multinomial logistic regression

$$\tilde{x} \in \mathbb{R}^D, \quad y \in \{1, \dots, C\}$$

If we define $(\tilde{x}, 1)$ for each datapoint,
 " \tilde{x}'

we can create a single weight matrix:

$$\underbrace{\tilde{a}}_{\text{dim} = C} = \underbrace{w \tilde{x} + b}_{\text{dim}} = \underbrace{\left(w \begin{array}{c|c} & b \\ \hline & \end{array} \right) \tilde{x}'}_{\text{" } w' \text{ "}} \begin{matrix} \text{dim} \\ \text{C} \times (D+1) \end{matrix}$$

$$\text{Recall that } p(y=c | \tilde{x}, \tilde{\theta}) = \frac{e^{a_c}}{\sum_{c'=1}^C e^{a_{c'}}}. \quad (*)$$

$$\text{Because } \sum_{c=1}^C p(y_n=c | \tilde{x}_n, \tilde{\theta}) = 1, \forall n$$

we do not need the last row in w' =>

=> we have $(C-1) \times (D+1)$ fitting prms.

We can use $p(y=c | \tilde{x}, \tilde{\theta}) = \frac{e^{a_c - a_\alpha}}{1 + \sum_{c'=1}^{C-1} e^{a_{c'} - a_\alpha}}$ $\stackrel{c=1, \dots, C-1}{\text{DBS}}$ => can use

Eg. (*) results in linear DBS => can use features $\tilde{g}(\tilde{x})$ to make non-linear DBS.

Objective function

$$-\mathcal{J}(\tilde{\theta}) = -\frac{1}{N} \log \prod_{n=1}^N \prod_{c=1}^C \mu_{n,c}^{y_{n,c}} \Leftrightarrow$$

$$\begin{cases} y_{n,c} = \mathbb{I}(y_n=c) \quad \text{one-hot encoding} \end{cases}$$

$$\begin{cases} \mu_{n,c} = \text{softmax}(\tilde{a}_n)_c \\ \tilde{a}_n = w \tilde{x}_n + b \end{cases}$$

$$\Leftrightarrow -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_{n,c} \log \mu_{n,c} = \frac{1}{N} \sum_{n=1}^N H(\vec{y}_n, \vec{\mu}_n),$$

where $H(\vec{p}, \vec{q}) = -\sum_{c=1}^C p_c \log q_c$.
cross-entropy

Optimization

First, consider $\frac{\partial \mu_c}{\partial a_j} = \frac{\partial}{\partial a_j} \left(\frac{e^{a_c}}{\sum_{c'=1}^C e^{a_{c'}}} \right) =$
 $= \frac{\delta_{cj} e^{a_c}}{\sum_{c'} e^{a_{c'}}} - \frac{e^{a_c}}{\left(\sum_{c'} e^{a_{c'}}\right)^2} e^{a_j} = \mu_c (\delta_{cj} - \mu_j)$.

Now, for a single datapoint \vec{x}' we have:

$$\frac{\partial \mathcal{L}}{\partial \vec{w}_j^T} (-\mathcal{L}(\vec{\theta})) = -\sum_{c=1}^C \frac{y_{c,j}}{\mu_c} \underbrace{\frac{\partial \mu_c}{\partial \vec{w}_j^T}}_{\text{class label}} \quad \textcircled{=} \\ \vec{w}_j^T = \text{vector of } \dim D+1 \quad a_j = \vec{w}_j^T \cdot \vec{x}' \quad \frac{\partial \mu_c}{\partial a_j} \frac{\partial a_j}{\partial \vec{w}_j^T} = \\ = \mu_c (\delta_{cj} - \mu_j) \vec{x}'$$

$$\textcircled{=} \sum_{c=1}^C y_c (\mu_j - \delta_{cj}) \vec{x}' = (\mu_j - y_j) \vec{x}'.$$

$\sum_c y_c = 1$

So, for all datapoints and all j (i.e., all classes), we obtain:

$$\underbrace{\vec{g}(\vec{\theta})}_{(D+1) \times (C-1) \text{ matrix}} = \frac{1}{N} \sum_{n=1}^N \vec{x}_n^T (\vec{\mu}_n - \vec{y}_n)^T$$

Likewise, we can compute the Hessian:
(for a single datapoint \tilde{x}')

$$\frac{\partial^2}{\partial \tilde{w}_k \partial \tilde{w}_j} (-J(\tilde{\theta})) = \frac{\partial}{\partial \tilde{w}_k} [(\mu_j - y_j) \tilde{x}'] = \\ = \mu_j (\delta_{jk} - \mu_k) \tilde{x}' \tilde{x}'^T.$$

Then for N datapoints we obtain:

$$H(\tilde{\theta}) = \underbrace{\frac{1}{N} \sum_{n=1}^N}_{c, c'} \mu_{n,c} (\delta_{cc'} - \mu_{n,c'}) \tilde{x}_n \tilde{x}_n^T$$

$(D+1) \times (D+1)$ matrix

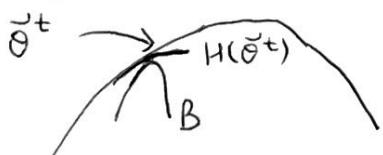
The entire $H(\tilde{\theta})$ is $(D+1)(C-1) \times (D+1)(C-1)$
and consists of $H_{c,c'}(\tilde{\theta})$ submatrices.

One can use grad or Newton's optimization.

One can also use bound optimization:

$$\text{consider } J(\tilde{\theta}) \approx J(\tilde{\theta}^t) + (\tilde{\theta} - \tilde{\theta}^t) \cdot \bar{g}(\tilde{\theta}^t) + \\ + \frac{1}{2} (\tilde{\theta} - \tilde{\theta}^t) \cdot B \cdot (\tilde{\theta} - \tilde{\theta}^t) \geq \\ \geq J(\tilde{\theta}^t) + (\tilde{\theta} - \tilde{\theta}^t) \cdot \bar{g}(\tilde{\theta}^t) + \frac{1}{2} (\tilde{\theta} - \tilde{\theta}^t) \cdot B \cdot (\tilde{\theta} - \tilde{\theta}^t),$$

where $H(\tilde{\theta}^t) \geq B$.



B is neg. def.

In our case (consider $G=2$ for simplicity), $H(\tilde{\theta}^t) = -\frac{1}{N} X^T S_t X$, where

$$S_t' = \text{diag} [\mu_{1,t}(1-\mu_{1,t}), \dots, \mu_{N,t}(1-\mu_{N,t})].$$

$$\text{Since } \mu_{i,t}(1-\mu_{i,t}) \leq \frac{1}{4},$$

$B = -\frac{1}{4N} X^T X$ will be more neg.-def. (i.e., it will curve down faster) than the 'real' curvature

Then we can use

$$\tilde{\theta}^{t+1} = \tilde{\theta}^t - 4N(X^T X)^{-1} \tilde{g}(\tilde{\theta}^t). \quad (**)$$

Since $(X^T X)^{-1}$ can be computed once, Eq. (**) should be more efficient than computing $N(X^T S_t X)^{-1}$ at each step.

MAP estimation

Recall that $\sum_{c=1}^G p(y_n=c | \tilde{x}_n, w') = 1$ leads to $(G-1)(D+1)$ indep. weights.

Now, let's use the original formula:

$$p(y=c | \tilde{x}, w') = \frac{e^{w'_c}}{\sum_{c'=1}^G e^{w'_c}}, \text{ which uses } G(D+1) \text{ weights.}$$

Introduce ℓ_2 regularization:

$$-\tilde{J}(\mathbf{w}') = -\sum_{n=1}^N \log p(y_n | \bar{x}_n, \mathbf{w}') + \frac{\lambda}{2} \sum_{c=1}^C \underbrace{\mathbf{w}_c^2}_{D+1 - \text{dim vector}}$$

$\lambda > 0$

Recall that

$$\frac{\partial}{\partial w_{cj}} (-\tilde{J}(\mathbf{w}')) = \sum_{n=1}^N (\mu_{nc} - y_{nc}) x'_{nj}.$$

$\uparrow j \uparrow$
 class label feature label
 $c=1, \dots, C$ $j=1, \dots, D+1$

Now, we need

$$\frac{\partial}{\partial w_{cj}} (-\tilde{J}(\mathbf{w}')) + \lambda w_{cj} = 0 \quad , \text{ or}$$

$$\lambda w_{cj} = \sum_{n=1}^N (y_{nc} - \mu_{nc}) x'_{nj}.$$

But then

$$\begin{aligned} \lambda \sum_{c=1}^C w_{cj} &= \sum_{c,n} (y_{nc} - \mu_{nc}) x'_{nj} = \\ &= \sum_n \left[\underbrace{\sum_c y_{nc}}_{1, \forall n} - \underbrace{\sum_c \mu_{nc}}_{1, \forall n} \right] x'_{nj} = 0. \end{aligned}$$

Thus, $\sum_{c=1}^C w_{cj} = 0$, $\forall j$ at the minimum \Rightarrow each ~~vector~~ column in the \mathbf{w}' matrix is constrained, we can use $C(D+1)$ weights.