

Multiclass logistic regression Lecture 9

Now consider $K > 2$:

$$p(C_k | \vec{x}) \equiv y_k(\vec{x}) = \frac{e^{a_k}}{\sum_{j=1}^K e^{a_j}}$$

$$a_k = \vec{w}_k^T \cdot \vec{x}$$

The goal is to determine $\{\vec{w}_k\}$ directly.

Use $\vec{t}_n = \{ \dots t_{n,1} \dots t_{n,k} \dots t_{n,K} \}$
 "1 if $\vec{t}_n \in C_k$, all other entries = 0

$$T = \left(\begin{matrix} & \underbrace{\hspace{10em}}_K \\ & \vdots \\ & \vdots \end{matrix} \right) \left. \vphantom{\begin{matrix} & \underbrace{\hspace{10em}}_K \\ & \vdots \\ & \vdots \end{matrix}} \right\} N \quad T_{nk} = t_{nk}$$

Then

$$\mathcal{L} = p(T | \vec{w}_1, \dots, \vec{w}_K) = \prod_{n=1}^N \prod_{k=1}^K \underbrace{p(C_k | \vec{y}_n)}_{y_{nk}}^{t_{nk}} \quad \text{①}$$

$$\text{②} \quad \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

$$E(\vec{w}_1, \dots, \vec{w}_K) = -\log \mathcal{L} = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_{nk}$$

Now, consider

$$\underbrace{\frac{\partial E}{\partial \vec{w}_j}}_{\text{vector derivative}} = - \sum_n \sum_k t_{nk} \frac{1}{y_{nk}} \frac{\partial y_{nk}}{\partial \vec{w}_j} \quad \text{③}$$

$$\frac{\partial y_{nk}}{\partial a_j} \frac{\partial a_j}{\partial \vec{w}_j}$$

\vec{y}_n

note that $\frac{\partial a_{j'}}{\partial \vec{w}_j} = 0, j' \neq j$

$$\frac{\partial y_k}{\partial a_j} = y_k \delta_{kj} - \frac{e^{a_k}}{(\sum_{j'} e^{a_{j'}})^2} e^{a_j} =$$

$$= y_k \delta_{kj} - y_k y_j.$$

$$\square - \sum_{n,k} \frac{t_{nk}}{y_{nk}} [y_{nk} \delta_{kj} - y_{nk} y_{nj}] \vec{y}_n =$$

$$= \sum_{n,k} [t_{nk} y_{nj} - t_{nk} \delta_{kj}] \vec{y}_n = \sum_n [y_{nj} - \underline{t_{nj}}] \vec{y}_n.$$

$$\uparrow$$

$$\sum_k t_{nk} = 1, \forall n.$$

MK x MK matrix

Moreover,

$$H = \frac{\partial^2 E}{\partial \vec{w}_k \partial \vec{w}_j} = \sum_n \left[\frac{\partial y_{nj}}{\partial \vec{w}_k} \right] \vec{y}_n =$$

$$= \sum_n [y_{nj} \delta_{kj} - y_{nk} y_{nj}] \underbrace{\vec{y}_n \vec{y}_n^T}_{M \times M \text{ matrix}}$$

Can show that $\vec{u}^T H \vec{u} > 0 \Rightarrow$ unique
minimum
positive definite

Thus can do NR algorithm again.

Probit regression

If desired, $\sigma(a)$ can be replaced by

$$\Phi(a) = \int_{-\infty}^a d\theta \mathcal{N}(\theta|0,1) \leftarrow \text{cumulative gaussian, or probit function}$$

Indeed, consider $\sigma(a) = \frac{1}{1+e^{-a}}$ vs. $\Phi(\lambda a)$.

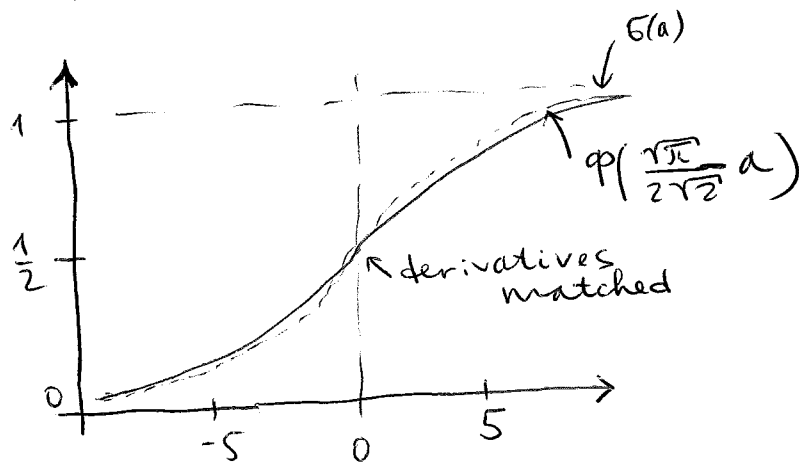
$$\left. \frac{d\sigma(a)}{da} \right|_{a=0} = \sigma(0)(1-\sigma(0)) = \frac{1}{4}$$

$$\text{But } \left. \frac{d\Phi(\lambda a)}{da} \right|_{a=0} = \left. \frac{d\Phi(\lambda a)}{d(\lambda a)} \lambda \right|_{a=0} = \frac{\lambda}{\sqrt{2\pi}} e^{-(\lambda a)^2} \Big|_{a=0} = \frac{\lambda}{\sqrt{2\pi}}$$

Match derivatives at $a=0$:

$$\frac{\lambda}{\sqrt{2\pi}} = \frac{1}{4} \Rightarrow \lambda^2 = \frac{\pi}{8}$$

So, $\sigma(a) \approx \Phi\left(\frac{\sqrt{\pi}}{2\sqrt{2}} a\right)$:



Can repeat the ML analysis for the logistic regression with probit regression \Rightarrow similar results in practice.

Laplace approximation

(saddle-point approx'n)

Consider

$$p(z) = \frac{f(z)}{Z}$$

$$\int dz p(z) = 1 \Rightarrow Z = \int dz f(z).$$

$$f(z) \approx f(z_0) e^{-\frac{A}{2}(z-z_0)^2}$$

$$\Rightarrow \left. \frac{d \log f(z)}{dz} \right|_{z_0} = 0 \text{ as well}$$

$$\uparrow A = - \left. \frac{d^2 \log f(z)}{dz^2} \right|_{z_0}$$

$$\left. \frac{df(z)}{dz} \right|_{z_0} = 0$$

@ max

$$\log f(z) \approx \log f(z_0) - \frac{1}{2} A (z-z_0)^2$$

Then $p(z) \Rightarrow q(z) = \sqrt{\frac{A}{2\pi}} e^{-\frac{A}{2}(z-z_0)^2} = \mathcal{N}(z|z_0, A^{-1})$

 Likewise, $f(\vec{z}) \approx f(\vec{z}_0) e^{-\frac{1}{2}(\vec{z}-\vec{z}_0)^T A (\vec{z}-\vec{z}_0)}$

stationary point

$$p(\vec{z}) \Rightarrow q(\vec{z}) = \frac{|A|^{1/2}}{(2\pi)^{D/2}} e^{-\frac{1}{2}(\vec{z}-\vec{z}_0)^T A (\vec{z}-\vec{z}_0)} =$$

$$= \mathcal{N}(\vec{z}|\vec{z}_0, A^{-1})$$

$$A = - \left. \nabla_{\vec{z}} \nabla_{\vec{z}} \log f(\vec{z}) \right|_{\vec{z}_0}$$

$D \times D$ matrix

Note that

$$Z = \int d\vec{z} f(\vec{z}) \approx f(\vec{z}_0) \frac{(2\pi)^{D/2}}{|A|^{1/2}}$$

$$\approx f(\vec{z}_0) \int \frac{d\vec{z}}{\sqrt{|A|}} e^{-\frac{1}{2}(\vec{z}-\vec{z}_0)^T A (\vec{z}-\vec{z}_0)} \text{ saddle-point approx'n}$$

Model comparison

Consider a set of models $\{M_i\}$ with
prms $\{\tilde{\theta}_i\}$. Define likelihood
 $p(\mathcal{D} | \tilde{\theta}_i, M_i)$, then

$$\underbrace{p(\mathcal{D} | M_i)}_{\text{model evidence}} = \int d\tilde{\theta}_i p(\mathcal{D} | \tilde{\theta}_i, M_i) p(\tilde{\theta}_i | M_i)$$

Under saddle-point approximation,

$$f(\tilde{x}) \Rightarrow p(\mathcal{D} | \tilde{\theta}_i, M_i) p(\tilde{\theta}_i | M_i) :$$

$$p(\mathcal{D} | M_i) \approx p(\mathcal{D} | \tilde{\theta}_{i, \text{MAP}}, M_i) p(\tilde{\theta}_{i, \text{MAP}} | M_i) \times$$

$$\times \frac{(2\pi)^{M/2}}{|A_i|^{1/2}}$$

where M is the # model prms and A_i is the Hessian of $\log p(\mathcal{D} | \tilde{\theta}_i, M_i) p(\tilde{\theta}_i | M_i)$ at $\tilde{\theta}_{i, \text{MAP}}$.
in model M_i [strictly speaking, $M \rightarrow M_i$]

$$A_i = - \vec{\nabla}_{\tilde{\theta}} \vec{\nabla}_{\tilde{\theta}} \left[\log(p(\mathcal{D} | \tilde{\theta}_i, M_i) p(\tilde{\theta}_i | M_i)) \right] \Big|_{\tilde{\theta}_{i, \text{MAP}}} =$$

$$Z p(\tilde{\theta}_i | \mathcal{D}, M_i), \text{ s.t. } \log(\dots) = \log p(\tilde{\theta}_i | \mathcal{D}, M_i) + \log Z$$

$$\log(\dots) = \log p(\tilde{\theta}_i | \mathcal{D}, M_i) + \log Z$$

$$= - \vec{\nabla}_{\tilde{\theta}} \vec{\nabla}_{\tilde{\theta}} \left[\log p(\tilde{\theta}_i | \mathcal{D}, M_i) \right] \Big|_{\tilde{\theta}_{i, \text{MAP}}} =$$

posterior

Finally,

$$\left[\log p(\mathcal{D} | M_i) \approx \log p(\mathcal{D} | \tilde{\theta}_{i, \text{MAP}}, M_i) + \log p(\tilde{\theta}_{i, \text{MAP}} | M_i) + \frac{M}{2} \log(2\pi) - \frac{1}{2} \log |A_i| \right]$$

"penalty" for model complexity

Moreover, assume that priors are given by
 $p(\vec{\theta}_i | M_i) = \mathcal{N}(\vec{\theta}_i | \vec{0}, \alpha^{-1} \mathbb{I}) \stackrel{\text{as before}}{=} \frac{\alpha^{M/2}}{(2\pi)^{M/2}} e^{-\frac{1}{2} \vec{\theta}_i^T (\alpha \mathbb{I}) \vec{\theta}_i}$

Then
$$A_i = \underbrace{-\vec{\nabla} \vec{\nabla} [\log p(\mathcal{D} | \vec{\theta}_i, M_i)]}_{\equiv H} \Big|_{\vec{\theta}_i, \text{MAP}}$$

$$\underbrace{-\vec{\nabla} \vec{\nabla} [\log p(\vec{\theta}_i | M_i)]}_{\alpha \mathbb{I}} \Big|_{\vec{\theta}_i, \text{MAP}} = H + \alpha \mathbb{I} \approx H$$

if α small
or we have lots of data
(N is large)

We have:
$$-\log p(\mathcal{D} | \vec{\theta}_i, M_i) \xrightarrow{\text{note that}} \sum_{n=1}^N \dots \text{ and } \uparrow \text{ as } N \uparrow$$

$$\log p(\mathcal{D} | M_i) \approx \log p(\mathcal{D} | \vec{\theta}_i, \text{MAP}, M_i) - \frac{M}{2} \log(2\pi) + \frac{M}{2} \log \alpha + \frac{M}{2} \log(2\pi)$$

$$- \frac{\alpha}{2} \vec{\theta}_i, \text{MAP}^T \vec{\theta}_i, \text{MAP} - \frac{1}{2} \log |H| \log \text{in the large-} N \text{ limit}$$

correction
same argument as above (α small and/or N large)

Finally,

$$H = \sum_{n=1}^N H_n = N \underbrace{\langle H \rangle}_{\frac{1}{N} \sum_n H_n} \Rightarrow \log |H| = \log |N \langle H \rangle| =$$

$$= \log [N^M \langle H \rangle] =$$

\uparrow this assumes that the rank of H is M

$$= M \log N + \log \langle H \rangle$$

We obtain:

$$\left[\log p(\mathcal{D} | M_i) \approx \log p(\mathcal{D} | \vec{\theta}_i, \text{MAP}, M_i) - \frac{M}{2} \log N \right]$$

Bayesian information criterion (BIC)

complexity penalty