

NEWTON-RAPHSON ITERATIVE OPTIMIZATION

$$W^{(new)} = W^{(old)} - H^{-1} \nabla E(W^{(old)})$$

H is the Hessian matrix whose elements comprise the second derivatives of $E(W)$ w.r.t. W

$$H_{ij} = \frac{\partial^2 E}{\partial w_i \partial w_j}$$

1D Newton's method

$$x^{(new)} = x^{(old)} - g(x^{(old)}) / g'(x^{(old)}) \quad \text{finding } g=0$$

$$\hookrightarrow x^{(new)} = x^{(old)} - f'(x^{(old)}) / f''(x^{(old)}) \quad \text{finding } f'=0$$

Try to apply the Newton-Raphson method to linear regression as a practice

$$\nabla E(W) = \sum_{n=1}^N (W^T \phi_n - t_n) \phi_n = \bar{\Phi}^T \bar{\Phi} W - \bar{\Phi}^T t$$

$$H = \nabla \nabla E(W) = \sum_{n=1}^N \phi_n \phi_n^T = \bar{\Phi}^T \bar{\Phi}$$

Newton's update

$$W^{(new)} = W^{(old)} - (\bar{\Phi}^T \bar{\Phi})^{-1} \{ \bar{\Phi}^T \bar{\Phi} W^{(old)} - \bar{\Phi}^T t \}$$

$$= (\bar{\Phi}^T \bar{\Phi})^{-1} \bar{\Phi}^T t$$

$$\bar{\Phi} = \begin{bmatrix} \vdots \\ \phi_n^T \\ \vdots \end{bmatrix}$$

$N \times M$

we get the standard least-squares solution

Now try to apply the Newton-Raphson method to

$$E(W) = -\ln p(t|W) = -\sum_{n=1}^N \{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \}$$

$$\nabla E(W) = \sum_{n=1}^N (y_n - t_n) \phi_n = \bar{\Phi}^T (y - t)$$

$$H = \nabla \nabla E(W) = \sum_{n=1}^N y_n (1-y_n) \phi_n \phi_n^T = \bar{\Phi}^T R \bar{\Phi}$$

The Hessian is no longer constant ; it depends on W through the weighting matrix R .

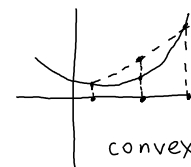
$$\left| \begin{array}{l} R \text{ } N \times N \text{ diagonal} \\ R_{nn} = y_n (1-y_n) \end{array} \right.$$

$$0 < y_n < 1$$

$$\Rightarrow v^T H v > 0 \text{ for an arbitrary } v$$

$$\Rightarrow H \text{ is positive definite}$$

The error function is a convex function of W and hence has a unique minimum.



$$W^{(new)} = W^{(old)} - (\bar{\Phi}^T R \bar{\Phi})^{-1} \bar{\Phi}^T (y - t)$$

$$= (\bar{\Phi}^T R \bar{\Phi})^{-1} \{ \bar{\Phi}^T R \bar{\Phi} W^{(old)} - \bar{\Phi}^T (y - t) \}$$

$$= (\bar{\Phi}^T R \bar{\Phi})^{-1} \bar{\Phi}^T R z$$

$$z = \bar{\Phi} W^{(old)} - R^{-1} (y - t) \quad \text{meaning?}$$

R depends on W . we need to update W iteratively

(Iterative Reweighted Least Squares, IRLS)

The weighting matrix R can be interpreted as variances

$$E[t] = \sigma(x) = y \quad \left| \begin{array}{l} \text{Bourneelli P.685} \\ p(t) = y^t (1-y)^{1-t} \end{array} \right.$$

$$\text{Var}[t] = E[t^2] - E[t]^2 = E[t] - E[t]^2 = y - y^2 = y(1-y)$$

$$\left\{ \begin{array}{l} t \in \{0, 1\} \\ t^2 = t \end{array} \right.$$

linearized problem in the space of $a = W^T \phi$

local linear approximation to logistic function

$$\left\{ \begin{array}{l} a_n(w) \approx a_n(w^{old}) + \frac{d a_n}{d y_n} \Big|_{w^{old}} (t_n - y_n) \\ = \phi_n^T w^{old} - \frac{(y_n - t_n)}{y_n(1-y_n)} = z_n \end{array} \right. \quad \boxed{\text{n-th element of } z}$$

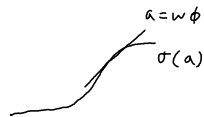
Meaning ?

z_n : as an effective target value in the space obtained by making a local linear approximation to the logistic sigmoid function around the current operating point w^{old}

$$\left\{ \begin{array}{l} \sigma = \frac{1}{1 + e^{-a}} \\ a = \ln \left(\frac{\sigma}{1-\sigma} \right) \\ \frac{da}{d\sigma} = \frac{1}{\sigma(1-\sigma)} \\ \sigma = y \quad \text{see (4.61) (4.88)} \end{array} \right.$$

(compared with least-squares solutions)

approximate y_n by a_n



MULTICLASS LOGISTIC REGRESSION

K classes

$$p(C_k | \phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k = w_k^T \phi \quad \left(\begin{array}{l} \text{softmax} \\ \text{posterior} \end{array} \right.)$$

discriminative approach

recall the 1-of-K coding scheme

the target value t_n for feature vector ϕ_n is

in the form $\begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$ only the k th component equals 1 $\phi_n \in C_k$

likelihood: $p(T | w_1, \dots, w_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k | \phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$

↑
N x K matrix

$$E(w_1, \dots, w_K) = -\ln p(T | w_1, \dots, w_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \quad (\text{cross entropy})$$

We want to find w_1, \dots, w_K that maximize E

To apply Newton-Raphson Method (IRLS), we need to compute ∇E and the Hessian $\nabla^2 E$

To begin with

$$\frac{\partial y_k}{\partial a_j} = y_k (I_{kj} - y_j)$$

$$I_{kj} = \begin{cases} 1, & j=k, \\ 0, & \text{otherwise.} \end{cases}$$

$$\frac{\partial y_k}{\partial a_k} = \frac{e^{a_k}}{\sum_i e^{a_i}} - \left(\frac{e^{a_k}}{\sum_i e^{a_i}} \right)^2$$

$$= y_k (1 - y_k)$$

$$\frac{\partial y_k}{\partial a_j} = -\frac{e^{a_k} e^{a_j}}{(\sum_i e^{a_i})^2}$$

$$= -y_k y_j \quad j \neq k$$

$$\frac{\partial E}{\partial y_{nk}} = -\frac{t_{nk}}{y_{nk}}$$

$$\frac{\partial E}{\partial a_{nj}} = \sum_{k=1}^K \frac{\partial E}{\partial y_{nk}} \frac{\partial y_{nk}}{\partial a_{nj}} = -\sum_{k=1}^K \frac{t_{nk}}{y_{nk}} y_{nk} (I_{kj} - y_{nj})$$

$$= -\sum_{k=1}^K t_{nk} (I_{kj} - y_{nj})$$

$$= -t_{nj} + \sum_{k=1}^K t_{nk} y_{nj}$$

$$= y_{nj} - t_{nj}$$

1-of-k coding
 $\forall n, \sum_k t_{nk} = 1$

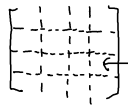
$$a_{nj} = w_j^T \phi_n$$

$$\nabla_{w_j} a_{nj} = \phi_n$$

compute the gradient w.r.t w_j

$$\nabla_{w_j} E(w_1, \dots, w_K) = \sum_{n=1}^N \frac{\partial E}{\partial a_{nj}} \nabla_{w_j} a_{nj} = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$

similarly the Hessian can be obtained by



$$\nabla_{w_k} \nabla_{w_j} E(w_1, \dots, w_K) = \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_n \phi_n^T$$

PROBIT REGRESSION

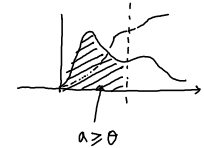
$$p(t=1|a) = f(a) \quad a = w^T \phi$$

for input ϕ_n , $a_n = w^T \phi_n$

Set the target by $\begin{cases} t_n=1, & \text{if } a_n \geq \theta, \\ t_n=0, & \text{otherwise.} \end{cases}$

If the value of θ is drawn from a density $p(\theta)$

$$\text{the } f(a) = \int_{-\infty}^a p(\theta) d\theta$$



suppose $p(\theta) \sim \mathcal{N}(\theta|0, 1)$

inverse probit function

$$\Psi(a) = \int_{-\infty}^a \mathcal{N}(\theta|0, 1) d\theta \quad \left(\begin{array}{l} \text{c.d.f. of} \\ \text{a Gaussian} \end{array} \right)$$

a related function:

$$\text{erf}(a) = \frac{2}{\sqrt{\pi}} \int_0^a \exp(-\theta^2) d\theta$$

$$\Psi(a) = \frac{1}{2} \left\{ 1 + \text{erf} \left(\frac{a}{\sqrt{2}} \right) \right\}$$

in MATLAB
 erf
 erfInv

The inverse probit function has a similar shape as the logistic sigmoid function.

We will use the inverse probit function in Bayesian logistic regression later.

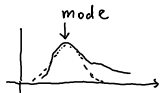
LAPLACE APPROXIMATION

2009年10月27日
下午 11:05

For prediction, we may integrate over the parameter W to get the predictive distribution. (Marginalization)

But the posterior distribution is no longer Gaussian in logistic regression.

How to resolve? Find a Gaussian centered at the mode of the posterior distribution as an approximation



Example:

$$p(z) = \frac{1}{Z} f(z) \quad Z = \int f(z) dz$$

z is a single continuous variable

① mode of $p(z)$: find z_0 st. $p'(z_0) = 0$

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

② Taylor expansion:

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2} a (z-z_0)^2 \quad \left. \begin{array}{l} \text{first derivative} \\ = 0 \end{array} \right\}$$

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

$$f(z) \simeq f(z_0) \exp \left\{ -\frac{a}{2} (z-z_0)^2 \right\}$$

normalized $\left\{ \begin{array}{l} z_0 \text{ must be} \\ \text{a local} \\ \text{maximum} \\ \Rightarrow a > 0 \end{array} \right.$

$$q(z) = \left(\frac{a}{2\pi} \right)^{\frac{1}{2}} \exp \left\{ -\frac{a}{2} (z-z_0)^2 \right\}$$

LAPLACE APPROXIMATION

2009年10月29日
上午 10:17

M -dimensional

$$p(\mathbf{Z}) = \frac{1}{Z} f(\mathbf{Z})$$

$$\ln f(\mathbf{Z}) \simeq \ln f(\mathbf{Z}_0) - \frac{1}{2} (\mathbf{Z}-\mathbf{Z}_0)^T \mathbf{A} (\mathbf{Z}-\mathbf{Z}_0)$$

$$\mathbf{A} = -\nabla \nabla \ln f(\mathbf{Z}) \Big|_{\mathbf{Z}=\mathbf{Z}_0} \quad \text{Hessian}$$

$$q(\mathbf{Z}) = \frac{|\mathbf{A}|^{\frac{1}{2}}}{(2\pi)^{\frac{M}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{Z}-\mathbf{Z}_0)^T \mathbf{A} (\mathbf{Z}-\mathbf{Z}_0) \right\}$$

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

mode $\mathbf{Z}_0 \Rightarrow$ local maximum $\Rightarrow \mathbf{A}$ positive definite

So, the Laplace approximation takes two steps:

- ① find the mode \mathbf{Z}_0 ;
- ② evaluate the Hessian at the mode.

more suitable for large data sets

Weakness: relies on a specific value of the variable, might fail to capture global properties.

BAYESIAN LOGISTIC REGRESSION

2009年10月29日

上午 10:40

Exact Bayesian inference for logistic regression is intractable.
 (The likelihood function comprises a product of logistic sigmoid function.)

Use Laplace approximation
 fitting a Gaussian centered at the mode of the posterior distribution

Begin with a Gaussian prior $p(w) = \mathcal{N}(w | m_0, S_0)$
 $\uparrow \nearrow$
 fixed hyperparameter

posterior $p(w | \mathbf{t}) \propto p(w) p(\mathbf{t} | w)$
 $\mathbf{t} = (t_1, \dots, t_N)^T$

$$\ln p(w | \mathbf{t}) = -\frac{1}{2} (w - m_0)^T S_0^{-1} (w - m_0) + \sum_{n=1}^N \{ t_n \ln y_n + (1 - t_n) \ln (1 - y_n) \} + \text{const}$$

(PRML 489)

$$y_n = \sigma(w^T \phi_n)$$

approximated by a Gaussian

↓
 what are the mean and the covariance?

2009年10月29日

上午 10:53

the mean : w_{MAP} Maximum a posteriori solution

We may use Newton-Raphson method to compute w_{MAP} as in maximum likelihood (regularized logistic regress vs. logistic regression)

the covariance :

$$S_N^{-1} = -\nabla \nabla \ln p(w | \mathbf{t}) = S_0^{-1} + \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T$$

so $q(w) = \mathcal{N}(w | w_{\text{MAP}}, S_N)$

is our approximation to the posterior distribution

Predictive Distribution (for class C_1)

Given a new feature vector $\phi(x)$
 marginalizing w.r.t. the posterior distribution $p(w | \mathbf{t})$

$$p(C_1 | \phi, \mathbf{t}) = \int p(C_1 | \phi, w) p(w | \mathbf{t}) dw \approx \int \sigma(w^T \phi) q(w) dw$$

Still intractable since $\sigma(w^T \phi)$ is the logistic sigmoid

2009年10月29日
上午 11:15

An interesting technique

$$a = w^T \phi$$

we write

$$\sigma(w^T \phi) = \int \delta(a - w^T \phi) \sigma(a) da$$

δ is the Dirac delta function

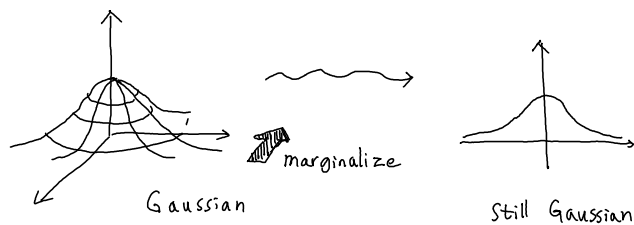
$\sigma(w^T \phi)$ is expressed in the form of integration

$$\begin{aligned} \text{Therefore, } \int \sigma(w^T \phi) q(w) dw & \\ &= \iint \delta(a - w^T \phi) \sigma(a) q(w) da dw \\ &= \int \sigma(a) \int \delta(a - w^T \phi) q(w) dw da \\ &= \int \sigma(a) p(a) da \end{aligned}$$

$$\text{where } p(a) = \int \delta(a - w^T \phi) q(w) dw$$

observe that $\delta(a - w^T \phi)$ imposes a linear constraint on w
 $p(a)$ is actually a marginal distribution from the joint
distribution $q(w)$ by integrating out all directions
orthogonal to ϕ .

$q(w)$ is a Gaussian $\xrightarrow{\text{marginalization}}$ still a Gaussian



What is the mean and variance of $p(a)$

2009年10月29日
上午 11:37

$$\begin{aligned} \mu_a &= E[a] = \int p(a) a da \\ &= \iint \delta(a - w^T \phi) q(w) a dw da \\ &= \int q(w) \int \delta(a - w^T \phi) a da dw \\ &= \int q(w) w^T \phi dw \\ &= w_{MAP}^T \phi = m_N^T \phi \\ s_a^2 &= \text{var}[a] = \int p(a) \{a^2 - E[a]^2\} da \\ &= \iint \delta(a - w^T \phi) q(w) \{a^2 - E[a]^2\} dw da \\ &= \int q(w) \int \delta(a - w^T \phi) \{a^2 - E[a]^2\} da dw \\ &= \int q(w) \{ (w^T \phi)^2 - E[w^T \phi]^2 \} dw \\ &= \phi^T S_N \phi \end{aligned}$$

the predictive distribution becomes

$$\begin{aligned} p(C_i | \#) &\approx \int \sigma^T(w^T \phi) q(w) dw \\ &= \int \sigma(a) p(a) da \\ &= \int \sigma(a) \mathcal{N}(a | \mu_a, s_a^2) da \end{aligned}$$

meaning?

2009年10月29日

下午 12:02

$$p(C_i | \mathbf{t}) \simeq \int \sigma(a) \mathcal{N}(a | \mu_a, s_a^2) da$$

convolution of a Gaussian with logistic sigmoid
cannot be evaluated analytically

By approximation again

What kind of function might look like a logistic sigmoid?
recall the inverse probit function (c.d.f of Gaussian)

$$\int \Psi(\lambda a) \mathcal{N}(a | \mu, \sigma^2) da = \Psi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{\frac{1}{2}}}\right)$$

requires a lot of work but is doable $\left| \lambda^2 = \frac{\pi}{8} \right.$

So if we approximate $\sigma(a)$ by the inverse probit function $\Psi(\lambda a)$, we can write down the integration analytically.

$$p(C_i | \phi, \mathbf{t}) \simeq \int \sigma(a) \mathcal{N}(a | \mu_a, s_a^2) da \simeq \sigma(\kappa(s_a^2) \mu_a)$$
$$\kappa(s_a^2) = (1 + \pi s_a^2 / 8)^{-\frac{1}{2}}$$

observe that if we have $p(C_i | \phi, \mathbf{t}) = 0.5$ as the decision boundary due to $\mu_a = 0$, which means $\mathbf{w}_{\text{MAP}}^T \phi = 0$,
So the boundary will be the same as the one obtained by the MAP solution. \Rightarrow equal prior

if $\mu_a \neq 0$, $p(C_i | \phi, \mathbf{t}) \neq 0.5$ more general