# Exercise 1
## Solutions

David Phelz

9/12/2016

David Puelz

# Linear Regression

- Consider $y = X\beta + \varepsilon$

$$y \in \mathbb{R}^N \quad, \quad X \in \mathbb{R}^{N \times p}, \quad \beta \in \mathbb{R}^p$$

- Principal of weighted least squares:

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^{N} \frac{w_i}{2} (y_i - x_i^T \beta)^2$$

A) (i) $\sum_{i=1}^{N} \frac{w_i}{2} (y_i - x_i^T \beta)^2$

$$= \frac{1}{2} \sum_{i=1}^{N} (y_i - x_i^T \beta) w_i (y_i - x_i^T \beta)$$

$$= \frac{1}{2} \sum_{i=1}^{N} (y_i w_i y_i - 2 y_i w_i x_i^T \beta + x_i^T \beta w_i x_i^T \beta)$$

$$= \frac{1}{2} (y^T W y - 2 y^T W X \beta + (X\beta)^T W X \beta)$$

$$= \frac{1}{2} (y - X\beta)^T W (y - X\beta)$$

(ii) Need to find minimum of

$$f(\beta) = y^T W y - 2\beta^T X^T W y + \beta^T X^T W X \beta$$

$$\Rightarrow \frac{\partial f(\beta)}{\partial \beta} = -2 X^T W y + 2 X^T W X \beta \quad \text{(matrix derivative)}$$

FOC is: $\dfrac{\partial f(\beta)}{\partial \beta} = 0 = -2 X^T W y + 2 X^T W X \hat{\beta}$

$\Rightarrow \hat{\beta}$ must satisfy:

$$0 = -2 X^T W y + 2 X^T W X \hat{\beta}$$

$$X^T W y = (X^T W X) \hat{\beta} \quad \blacksquare$$

(iii) Let's look at the QR factorization.

Theorem: Given problem

$$\min_{\beta} \| y - X\beta \|_2 \quad \text{, its solution set is}$$

⎫
⎬ one approach
⎭

equivalent to $R\beta = Q^T y$ where $Q$ is
the solution set of

orthonormal, $R$ is upper triangular, and both come from the QR factorization of $X$.

A second approach would be to solve the set of normal equations: (FOC from least-squares)

$$(X^T W X)\beta = X^T W y \qquad \text{where the}$$

unknown is $y$.

Write: $C = X^T W X$

$$d = X^T W y$$

$$L L^T = C \qquad \leftarrow \text{Cholesky factorization}$$

$$\text{(L upper triangular)}$$

Algorithm: Solve $L z = d$, forward substitution
for $z$ $\qquad\qquad\qquad$ since $L$ is triangular.

Solve $\qquad\qquad\qquad\qquad$ backward substitution
for $\beta$ $\quad L^T \beta = z$ ,

---

Class notes:

3 standard factorizations

1) Cholesky (LU) $\quad$⎫$\longrightarrow$ fast, unstable (turfethan)
$\qquad\qquad\qquad\qquad\qquad$ book.
2) QR $\qquad\qquad\qquad$ ⎬$\longrightarrow$ fast, stable.
3) SVD $\qquad\qquad\qquad$⎭$\longrightarrow$ matrix that is close to
$\qquad\qquad\qquad\qquad\qquad\qquad$ rank deficient.

$\downarrow$

Focus on QR:

$$W^{1/2}X = QR \qquad (\text{where} \quad W^{1/2}W^{1/2} = W)$$

$N \times N$    $N \times P$     $N \times P$

square
$P \times P$
orthonormal
columns

"upper
$\triangle$"

reduced QR
factorization.

looking back at normal equations:

$$\Rightarrow X^T W Y = X^T W X \beta$$

$$\Rightarrow X^T W^{1/2} W^{1/2} Y = X^T W^{1/2} W^{1/2} X \beta$$

$$\Rightarrow (QR)^T W^{1/2} Y = (QR)^T QR \beta$$

$$\Leftrightarrow \cancel{R^T} Q^T W^{1/2} Y = \cancel{R^T} Q^T QR \beta$$

$$Q^T W^{1/2} Y = R \beta \qquad \longleftarrow \quad \text{now we have an upper triangular system!!!}$$

# Generalized Linear Models,

Let $y_i \sim Bin(m_i, w_i)$

with $w_i(\beta) = \dfrac{1}{1 + \exp\{-x_i^T \beta\}}$.

A) likelihood is:

$$\ell(\beta) = -\log\left\{ \prod_{i=1}^{N} p(y_i | \beta) \right\}$$

where $p(y_i | \beta) = \binom{m_i}{y_i} w_i^{y_i} (1-w_i)^{m_i-y_i}$

$$\Rightarrow \ell(\beta) = -\log\left\{ \prod_{i=1}^{N} \binom{m_i}{y_i} w_i^{y_i} (1-w_i)^{m_i-y_i} \right\}$$

$$\propto -\sum_{i=1}^{N} \left[ y_i \log w_i + (m_i - y_i) \log(1-w_i) \right]$$

let's define

$$\ell_i(\beta) = y_i \log w_i + (m_i - y_i) \log(1-w_i)$$

so that

$$\ell(\beta) = -\sum_{i=1}^{N} \ell_i(\beta)$$

→ for a general function, $f(x)$

$$\frac{d}{dx} \log f(x) = \frac{f'(x)}{f(x)}.$$

→ Similarly,

$$\nabla_\beta \log w_i(\beta) = \frac{\nabla_\beta w_i(\beta)}{w_i(\beta)}.$$

Solving for gradient of $w_i(\beta)$:

$$\nabla_\beta w_i(\beta) = \frac{-\exp\{\bar{x}_i^T \beta\}}{(1+\exp\{-x_i^T \beta\})^2} x_i = w_i(1-w_i) x_i$$

$$\implies \frac{\nabla_\beta w_i(\beta)}{w_i(\beta)} = \frac{-\exp\{-x_i^T \beta\}}{1+\exp\{-x_i^T \beta\}} x_i$$

$$= -(1-w_i) x_i$$

Now, looking back at the likelihood, we have

$$\nabla l_i(\beta) = \nabla(y_i \log w_i) + (m_i - y_i)\log 1 - w_i)$$

$$= -y_i(1-w_i)x_i + (m_i - y_i)w_i x_i$$

$$= -y_i x_i + y_i w_i x_i + m_i w_i x_i - y_i w_i x_i$$

$$= -(y_i - m_i w_i)x_i \qquad \}$$

where $\quad \nabla_\beta \log\left(1 - w_i(\beta)\right) = \quad w_i x_i \quad$ by similar argument to above.

Therefore,

$$\nabla \ell(\beta) = -\sum_{i=1}^{N} \nabla \ell_i(\beta)$$

$$= -\sum_{i=1}^{N} \left(y_i - m_i w_i\right) x_i$$

$$= -X^T S \qquad X \text{ matrix}, \qquad S \text{ is vector with } N \text{ components \& } i^{th} \text{ component } \{y_i - m_i w_i\}$$

c) let $\beta_0 \in \mathbb{R}^P$, we call that we write

$\ell(\beta)$ as

$$\ell(\beta) = -\sum_{i=1}^{N} \ell_i(\beta)$$

with $\ell_i(\beta) = y_i \log w_i + (m_i - y_i) \log(1 - w_i)$

Taylor's theorem states (approximating $\ell(\beta)$ near $\beta_0$)

$$\ell(\beta) \approx \ell(\beta_0) + \underbrace{\nabla \ell(\beta_0)^T (\beta - \beta_0)}_{\text{1st order term}} + \underbrace{\frac{1}{2}(\beta - \beta_0)^T \nabla^2 \ell(\beta_0)(\beta - \beta_0)}_{\text{2nd order term}}$$

let's solve for the Hessian, $\nabla^2 \ell(\beta_0)$.

Recall:

$$\nabla \ell(\beta) = -X^T S \quad , \quad S \text{ vector with}$$
$$\text{i}\underline{\text{th}} \text{ component}$$
$$y_i - m_i w_i$$
$$= -X^T(y - Mw) \quad , \quad M \text{ is matrix}$$
$$\text{(diagonal) with}$$
$$= -X^T y + X^T Mw \quad \quad M_{ii} = m_i.$$

Now, we aim to calculate

$$\nabla^2 \ell(\beta) = \nabla(-X^T y + X^T Mw)$$

$$= \nabla(X^T Mw)$$

$$= X^T M \nabla w$$

where the gradient of a vector field is a second order tensor (ie: matrix).

$$\nabla w = \frac{\partial w_i}{\partial \beta_j} e_i \otimes e_j \qquad \text{where } \{e_i\}_{i=1}^N \text{ is standard basis for } \mathbb{R}^N.$$

$$= \begin{bmatrix} \frac{\partial w_1}{\partial \beta_1} & \cdots & \frac{\partial w_1}{\partial \beta_N} \\ \vdots & & \vdots \\ \frac{\partial w_N}{\partial \beta_1} & & \frac{\partial w_N}{\partial \beta_N} \end{bmatrix}$$

$$= \begin{bmatrix} \nabla_\beta w_1(\beta) & \cdots & \nabla_\beta w_N(\beta) \end{bmatrix}^T$$

We've shown!

$$\nabla_\beta w_i(\beta) = w_i(1-w_i) x_i$$

$$\Longrightarrow$$

$$\nabla w = \begin{bmatrix} w_1(1-w_1)x_1 & \cdots & w_N(1-w_N)x_N \end{bmatrix}^T$$

$$= \tilde{W} X \qquad , \qquad W \text{ diagonal matrix with } W_{ii} = w_i(1-w_i)$$

Therefore

So simple.

$$\boxed{\nabla^2 \ell(\beta) = \nabla(X^T M w) = X^T M \nabla w = X^T M \tilde{W} X}$$

So, by Taylor's theorem

$$\ell(\beta) \approx$$

$$\ell(\beta_0) - (X^T S)^T (\beta - \beta_0) + \frac{1}{2}(\beta - \beta_0)^T X^T A X (\beta - \beta_0)$$

$$\begin{cases} \text{with} : A = M\tilde{W} \\ \tilde{W}_{ii} = w_i(1 - w_i) \\ M_{ii} = m_i \end{cases}$$

$$= $$

$$\ell(\beta_0) - S^T(X\beta - X\beta_0) + \frac{1}{2}(X\beta - X\beta_0)^T A (X\beta - X\beta_0)$$

$$= $$

$$\ell(\beta_0) - S^T(X\beta - \tilde{z}) + \frac{1}{2}(X\beta - \tilde{z})^T A (X\beta - \tilde{z})$$

$$\begin{cases} \tilde{z} = X\beta_0. \end{cases}$$

completing
the
square!
w.r.t to
$\underline{X\beta - \tilde{z}}$

$$\propto \frac{1}{2}\left(X\beta - \tilde{z} - A^{-1}S\right)^T A \left(X\beta - \tilde{z} - A^{-1}S\right) + c$$

Flipping
around
signs

$$= \frac{1}{2}\left(\{\tilde{z} + A^{-1}S\} - X\beta\right)^T A \left(\{\tilde{z} + A^{-1}S\} - X\beta\right) + c$$

$$= \frac{1}{2}(z - X\beta)^T A (z - X\beta) + c$$

$$z = X\beta_0 + (M\tilde{W})^{-1}S \quad , \quad A = M\tilde{W}$$