

TP 2

YUNLONG JIAO, 23/01/2019

Ex 1. Derive analytical solutions to the following optimization problems. Then implement them in R/Python.

$$D = \left\{ (x_i, y_i) \right\}_{i=1}^n$$

$x_i \in \mathbb{R}^p, y_i \in \mathbb{R}$

① OLS:  $\min_{\beta \in \mathbb{R}^p} \frac{1}{n} \|y - X\beta\|^2 := J(\beta)$

$y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$   
 $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_m^T \end{pmatrix} \in \mathbb{R}^{m \times p}$

$$\frac{\partial J}{\partial \beta} = -\frac{2}{n} X^T (y - X\beta) = 0 \quad \text{Model: } y_i = \beta^T x_i + \varepsilon_i$$

$$(X^T X) \beta = X^T y$$

If  $(X^T X)$  invertible.

$$\hat{\beta}^{OLS} = (X^T X)^{-1} X^T y$$

⇒ Prediction:

$$x^* \in \mathbb{R}^p$$

$$\hat{y}^* = (x^*)^T \hat{\beta}^{OLS} = (x^*)^T (X^T X)^{-1} X^T y$$

② Ridge:  $\min_{\beta} \frac{1}{n} \|y - X\beta\|^2 + \lambda \|\beta\|^2 \quad (\lambda > 0)$

$$\frac{\partial J}{\partial \beta} = -\frac{2}{n} X^T (y - X\beta) + 2\lambda \beta = 0$$

$$\hat{\beta}^{Ridge} = (X^T X + \lambda n I)^{-1} X^T y \quad \text{(Ex)}$$

$$= X^T (X X^T + \lambda n I)^{-1} y$$

⇒ Prediction:

$$x^* \in \mathbb{R}^p$$

$$\hat{y}^* = (x^*)^T \hat{\beta}^{Ridge}$$

$$= (x^*)^T (X^T X + \lambda n I)^{-1} X^T y$$

$$\stackrel{\text{(Ex)}}{=} \underbrace{(x^*)^T X^T}_{K_{x^*, X}} \underbrace{(X X^T + \lambda n I)^{-1} y}_{K_{X X^T + \lambda n I}}$$

$$K_{x^*, X} \quad K_{X X^T + \lambda n I}$$

GP

### ③ Weighted Ridge Regression (WRR)

$$\min \frac{1}{n} \sum_i w_i (y_i - \beta^T x_i)^2 + \lambda \|\beta\|^2 \quad (w_i > 0)$$

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

$$W = \text{diag}(w_1, \dots, w_n)$$

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

$$= \frac{1}{n} \|\tilde{y} - \tilde{X}\beta\|^2 + \lambda \|\beta\|^2 \quad \tilde{y} = W^{\frac{1}{2}}y, \quad \tilde{X} = W^{\frac{1}{2}}X$$

Due to Ridge ②.

$$\hat{\beta}^{W\text{-Ridge}} = (X^T W X + \lambda n I)^{-1} X^T W y$$

Ex 2: (IRLS for <sup>Ridge</sup> Logistic Regression.)

(90min)

$$\mathcal{D} = \{(\vec{x}_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^p, y_i \in \{-1, +1\}$$

$$X \in \mathbb{R}^{n \times p}, y \in \{-1, +1\}^n$$

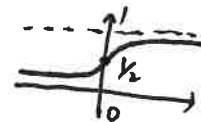
Model:  $\mathbb{P}(y_i | \vec{x}_i, \beta) = \begin{cases} \sigma(\beta^T x_i) & \text{if } y_i = +1 \\ 1 - \sigma(\beta^T x_i) & \text{if } y_i = -1 \end{cases} = \sigma(y_i \beta^T x_i)$

where: sigmoid

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

$$\sigma(-a) = 1 - \sigma(a)$$

$$\sigma'(a) = \sigma(a) \sigma(-a)$$



LR,  $\min \frac{1}{n} \sum_i [-\log \mathbb{P}(y_i | x_i, \beta)] + \lambda \|\beta\|^2 \stackrel{\text{Gaussian Prior}}{:=} J(\beta)$

$$= -\frac{1}{n} \sum_i \log \sigma(y_i \beta^T x_i) + \lambda \|\beta\|^2$$

① Derive:  $\nabla J(\beta), \nabla^2 J(\beta)$

$$\nabla_{\beta} J(\beta) = -\frac{1}{n} \sum \frac{\sigma(y_i \beta^T x_i) \sigma(-y_i \beta^T x_i)}{\sigma(y_i \beta^T x_i)} y_i x_i + 2\lambda \beta$$

$$= -\frac{1}{n} \sum y_i \sigma(-y_i \beta^T x_i) x_i + 2\lambda \beta$$

$$:= -\frac{1}{n} X^T P(\beta) y + 2\lambda \beta$$

$$P(\beta) = \text{diag}[\sigma(-y_i \beta^T x_i)]_n$$

$$(x_1, \dots, x_n) \setminus \left( \begin{matrix} y \\ \vdots \\ y \end{matrix} \right)_n \quad \sqrt{2}$$

$$\nabla^2 J(\beta) = \frac{1}{n} \sum \underbrace{\sigma(y_i \beta^T x_i) \sigma(-y_i \beta^T x_i)}_{\text{diag}[\sigma(\beta^T x_i) \sigma(-\beta^T x_i)]} x_i x_i^T + 2\lambda I$$

$$=: \frac{1}{n} X^T W(\beta) X + 2\lambda I$$

(x<sub>1</sub>, ..., x<sub>n</sub>) \setminus (||)

$$W(\beta) = \text{diag}[\sigma(\beta^T x_i) \sigma(-\beta^T x_i)]_n$$

each step of.

② Show that the <sup>Raphson</sup> Newton solver of LRR is equivalent to WRR.

Gradient Descent:

$$\min_x f(x) \quad f \in C^1$$

- Initialize  $x := x^{(0)}$
  - For  $k=0, \dots, L$ , or until convergence (\*)
- $$x^{(k+1)} = x^{(k)} - t \nabla f(x^{(k)})$$
- step size

at  $k$ -th step:

$$f(x) \approx f(x^{(k)}) + (x - x^{(k)})^T \nabla f(x^{(k)}) + \frac{1}{2t} \|x - x^{(k)}\|^2 := f_{g_1}(x)$$

$$\min_x f_{g_1}(x) \quad (\text{around } x^{(k)})$$

$$\Rightarrow x^{(k+1)} = x^{(k)} - t \nabla f(x^{(k)})$$

Newton Raphson solver:

$$\min f(x) \quad f \in C^2 \quad (Ex 1.3)$$

- Initialize  $x := x^{(0)}$
- For  $k=0, \dots, L$ , or until convergence (\*)

$$x^{(k+1)} = x^{(k)} - (\nabla^2 f(x^{(k)}))^{-1} \nabla f(x^{(k)})$$

at  $k$ -th step:

$$f(x) \approx f(x^{(k)}) + (x - x^{(k)})^T \nabla f(x^{(k)}) + \frac{1}{2} (x - x^{(k)})^T \nabla^2 f(x^{(k)}) (x - x^{(k)}) := f_{g_2}(x)$$

$$\min_x f_{g_2}(x) \quad (\text{around } x^{(k)})$$

$$\Rightarrow x^{(k+1)} = x^{(k)} - (\nabla^2 f(x^{(k)}))^{-1} \nabla f(x^{(k)})$$

(\*) Stopping criteria way too simplified here...!

(+) Newton is rarely used, consider Conjugate Gradient (~~Truncated~~ Newton), BFGS, (Quasi-Newton)

$$\beta^{\text{new}} \leftarrow \beta^{\text{old}} - [\nabla^2 J(\beta^{\text{old}})]^{-1} \nabla J(\beta^{\text{old}})$$

LRR:  $\min J(\beta) := -\frac{1}{n} \sum_i \log \sigma(y_i \beta^T x_i)$

at  $k$ -th step:

$$\underbrace{(J(\beta) \approx)}_{f_{g_1}} J(\beta) \approx J(\beta^{(k)}) + (\beta - \beta^{(k)})^T \nabla J(\beta^{(k)}) + \frac{1}{2} (\beta - \beta^{(k)})^T \nabla^2 J(\beta^{(k)}) (\beta - \beta^{(k)})$$

where  $\beta^T \nabla J(\beta^{(k)}) = -\frac{1}{n} \beta^T X^T P y + 2\lambda \beta^T \beta^{(k)}$

$$-\beta^T \nabla^2 J(\beta) \beta^{(k)} = -\frac{1}{n} \beta^T X^T W X \beta^{(k)} - 2\lambda \beta^T \beta^{(k)}$$

$$\frac{1}{2} \beta^T \nabla^2 J(\beta^{(k)}) \beta = \frac{1}{2n} \beta^T X^T W X \beta + \lambda \beta^T \beta$$

$$P := P(\beta^{(k)})$$

$$W := W(\beta^{(k)})$$

$$\Rightarrow 2 J_g(\beta) = \frac{1}{n} \beta^T X^T W X \beta - \frac{2}{n} \beta^T X^T W \underbrace{(X \beta^{(k)} + W^{-1} P y)}_{=: z} + 2\lambda \|\beta\|^2 + C$$

$$= \frac{1}{n} (z - X\beta)^T W (z - X\beta) + 2\lambda \|\beta\|^2 + C$$

③ Implement

solve LRR (...) in R/Python.

Set:

$$[f_i] = [(\beta^{(k)})^T x_i] \quad i=1, \dots, n.$$

$$[W_i] = \text{diag} [\sigma(f_i) \sigma(-f_i)]$$

$$[P_i] = \text{diag} [\sigma(-y_i f_i)]$$

$$[z_i] = [f_i + y_i P_i W_i^{-1}] = [f_i + y_i / \sigma(y_i f_i)]$$

Call: solve WRR (z, X, W, 2λ) at each update.