# Lecture September 24

- Short note about Project scaling:

  - OLS, with intercept or not, no change in MSE

    Do it yourself.

    $X \Rightarrow X - \text{mean}(X)$

    $y \Rightarrow y - \text{mean}(y)$

    $\beta_0 = \text{mean}(y)$

  - Ridge: regularization term

    $$\lambda \|\beta\|_2^2$$

    $$= \lambda \sum_{j=0}^{p-1} \beta_j^2$$

    $\beta_0$ not included.

- same with Lasso.

$$X = \begin{bmatrix} 1 & x_0^0 & x_0^0 & - - & x_0 \\ 1 & x_1^1 & x_1^2 & - \cdots & \\ 1 & 1 & & & \\ 1 & 1 & & & \\ 1 & 1 & & & \end{bmatrix}$$

if you keep the intercept column, when comparing own code with sklearn fit_intercept = False.

— Coding

Python

— numpy

— Pandas

— CV $\begin{cases} k = 5 \quad \boxed{\text{Train} \mid - \mid - \mid \cdots \mid \text{Test}} \\ k \text{ Fold in sklearn.} \end{cases}$

———————— * ————————

Classification &
Logistic regression

— Linear Regression

$$y = f(x) + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2)$$

$$\tilde{y} = X\beta$$

$$E[y] = X\beta$$

$$\text{var}[y] = \sigma^2$$

$$y \sim N(X\beta, \sigma^2)$$

— Binary classification

$$y_i = 1 \implies p(y_i = 1 | x_i \beta)$$

$$y_i = 0 = 1 - p(y_i = 1 | x_i \beta)$$

$$y = p(x) + \varepsilon$$

$\varepsilon \sim$ Binomial distribution

$y, \varepsilon$ are iid

$$D = \{ (x_0, y_0), (x_1 y_1) \ldots (x_{n-1} y_{n-1}) \}$$

$$P(D|\beta)$$

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\arg\max} \; P(D|\beta)$$

$$C(\beta) = - \log P(D|\beta)$$

$$P(D|\beta) = \prod_{i=0}^{n-1} P(y_i = 1|\beta)^{y_i} \underbrace{\left(1 - P(y_i = 1|\beta)\right)}_{P(y_i = 0|\beta)}^{1-y_i}$$
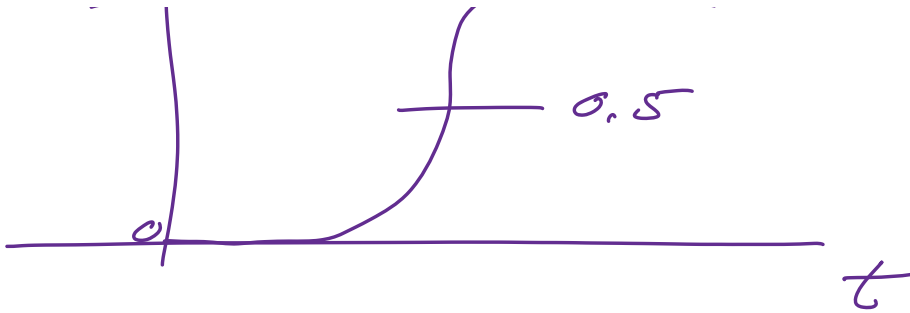
$$P(y_i = 1|x_i\,\beta) = \frac{e^{t(x_i\,\beta)}}{1 + e^{t(x_i\,\beta)}}$$

$$= \frac{e^t}{1 + e^t}$$

$$1 - P(y_i = 1|x_i\,\beta) = \frac{1}{1 + e^t}$$

$$0 \le P(y_i|x_i\,\beta) \le 1$$

if

$$t = t(x, \beta) = \beta_0 + \beta_1 x$$

$$\left( t = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots \beta_{p-1} x_{ip} \right)$$

$$C(\beta) = - \sum_{i=0}^{n-1} \left\{ y_i \log P_i' + (1 - y_i) \log(1 - P_i') \right\}$$

$$\left( P_i = p(y_i = 1 | x_i \beta) \right)$$

$$P_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$

$$C(\beta) = - \sum_{i=0}^{n-1} \left\{ y_i (\beta_0 + \beta_1 x_i - \log(1 + e^{\beta_0 + \beta_1 x_i})) - (1 - y_i) \log(1 + e^{\beta_0 + \beta_1 x_i}) \right\}$$

$$\frac{\partial C(\beta)}{\partial \beta_0} = 0 = - \sum (y_i - P_i)$$

$$\frac{\partial C(\beta)}{\partial \beta_1} = 0 = - \sum x_i (y_i - P_i)$$

$$\frac{\partial C}{\partial \beta} = 0 = \boxed{-X^T(y - P)}$$

$$P(x_i, y_i | \beta)$$

$$\frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$

non-linear
depedence on $\beta$

$$y_i = [1, 1] \quad \rightsquigarrow \quad \tilde{y}_i = [-1, 1]$$

$$\bigcirc^2 \qquad\qquad T$$

$$\frac{\partial L}{\partial \beta \partial \beta^T} = X W X$$

$$0 \le p_i \le 1$$

$$\left( \begin{array}{l} W_{ii} = \overset{g_i=1}{\overbrace{\frac{p_i}{}(1-p_i)}} \\ w_{ij} = 0 \quad i \ne j \quad i \ne j \end{array} \right)$$

$$W \ge 0$$

In Linear Regression

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = X^T X = H$$

$$(\text{Hessian})$$

$$\boxed{g = -X^T(y-P) = 0}$$

$$f(s) = 0$$

Newton-Raphson

Taylor expand

$$f(s) = f(x) + (s-x)f'(x) + \frac{(s-x)^2}{} f''$$

$$\approx f(x) + (s-x) f'(x) = 0$$

$$S = x - f(x)/f'(x)$$

suggests an iterative procedure :

$$S \Rightarrow x_{m+1} = x_m - f(x_m)/f'(x_m)$$

stop when $|x_{m+1} - x_m| \leq \delta$

$$\delta \sim 10^{-10}$$

$$\underline{g(\beta)} = -x^T(y - p(\beta)) = 0$$

Generalization of
Newton-Raphson to
more than one variable
applied to $g(p) = 0$

$$\beta_{m+1} = \beta_m - g(\beta_m)/H(\beta)$$

$$\frac{\partial g}{\partial p^T} = \frac{\partial g^T}{\partial \beta} = H$$

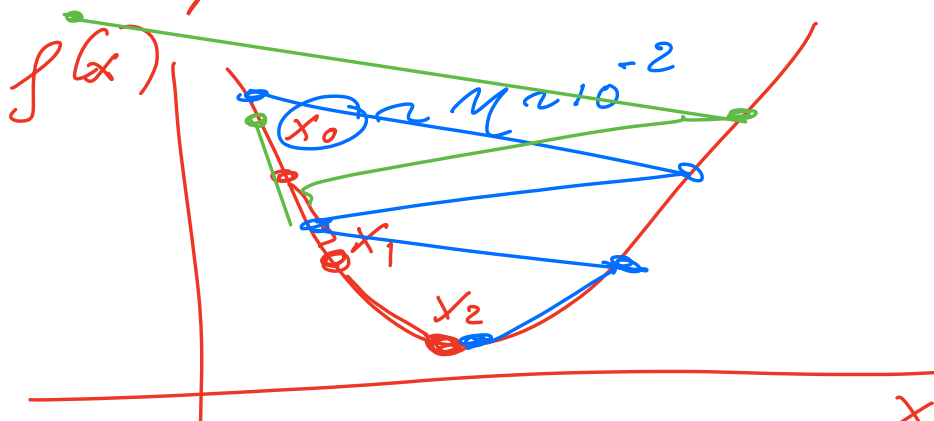$$\beta_{n+1} = \beta_n + \frac{1}{H} X^T (y - p(\beta_n))$$

$$\frac{1}{H} \longrightarrow \quad \eta = \text{learning rate}$$

$$\boxed{\beta_{n+1} = \beta_n - \eta\, g(\beta_n)}$$

$\eta$ is a parameter

$$\eta = [10^{-5}, 10^{-9}, \dots]$$

gradient descent