

ML CS229.

PS 1.

1. Linear Classifiers (Logistic Regression & GDA)

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(g(\theta^T x^{(i)})) + (1-y^{(i)}) \log(1-g(\theta^T x^{(i)}))$$

1: $\nabla_{\theta} J(\theta)$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \frac{1}{g(\theta^T x^{(i)})} - (1-y^{(i)}) \frac{1}{1-g(\theta^T x^{(i)})} \right].$$

$$g'(\theta^T x^{(i)}) \cdot (\theta^T x^{(i)})' = \underline{g(\theta^T x^{(i)}) (1-g(\theta^T x^{(i)}))} \cdot x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \cdot (1-g(\theta^T x^{(i)})) \cdot (1-y^{(i)}) \cdot g(\theta^T x^{(i)}) \right] x_j^{(i)}$$

$$\Rightarrow \nabla_{\theta_j} J(\theta) = +\frac{1}{m} \sum_{i=1}^m [g(\theta^T x^{(i)}) - y^{(i)}] x_j^{(i)}$$

$$\Rightarrow \nabla_{\theta} J(\theta) = \frac{1}{m} (g(\theta^T x) - y)^T X$$

$$H_{jk} = \frac{\partial^2 J(\theta)}{\partial \theta_j \partial \theta_k} = \frac{1}{m} \sum_{i=1}^m g(\theta^T x^{(i)}) [1-g(\theta^T x^{(i)})] x_j^{(i)} x_k^{(i)} = \frac{1}{\exp(-(\theta^T x + \theta_0)) + 1}$$

$$\rightarrow H = \frac{1}{m} \cdot X^T g(X\theta) (1-g(X\theta)) X$$

$$Z^T X^T \cdot \underbrace{g(X\theta) \cdot (1-g(X\theta))}_{\text{ER} \sim [0,1]} \cdot Z = k(Z^T X)^T \cdot (X \cdot Z) = k(Z^T X) \cdot X^T \cdot Z = k(Z^T X)^2 \geq 0$$

$$\rightarrow Z^T Z = \frac{1}{m} \sum_{i=1}^m \underbrace{g(\theta^T x^{(i)}) [1-g(\theta^T x^{(i)})]}_{\text{row } \alpha} x_j^{(i)} x_k^{(i)} z_j \cdot z_k$$

$$= \frac{1}{m} \sum_i \sum_j \sum_k z_j \cdot x_k^{(i)} \cdot \cancel{x_j^{(i)}} \cdot z_k = \frac{1}{m} \sum_i (Z^T X^{(i)})^2 \geq 0$$

$\therefore H \geq 0$

(c) GDA, Gaussian Discriminant Analysis Note)

Bayes' Theorem $P(y=1|x) = \frac{P(x|y=1; \mu_1, \Sigma)}{P(x)}$

$\therefore P(y=1|x) = \frac{P(x|y=1; \mu_1, \Sigma) P(y=1; \phi)}{P(x|y=0; \mu_0, \Sigma) P(y=0; \phi) + P(x|y=1; \mu_1, \Sigma) P(y=1; \phi)}$

$= \frac{1}{\frac{P(x|y=0; \mu_0, \Sigma) P(y=0; \phi)}{P(x|y=1; \mu_1, \Sigma) P(y=1; \phi)} + 1}$

$= \frac{1}{\frac{\exp(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0)) (1-\phi)}{\exp(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)) \phi} + 1}$

$= \frac{1}{\frac{\frac{1}{\phi} \exp(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0)) + \frac{1}{1-\phi} (x-\mu_1)^T \Sigma^{-1} (x-\mu_1)}{\exp(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0)) + \exp(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1))} + 1}$

$= \frac{1}{\exp\left(\frac{\log \frac{\phi}{1-\phi}}{2} + \frac{1}{2} (x-\mu_0)^T \Sigma^{-1} (x-\mu_0) + (x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right) + 1}$

$= \exp\left(\frac{\mu_0 - \mu_1}{2} \cdot \frac{x + \log \frac{\phi}{1-\phi}}{\sqrt{\phi(1-\phi)}} + \frac{\mu_0^2 - \mu_1^2}{2 \Sigma}\right) + 1$

(d) maximize $L(\phi, \mu_0, \mu_1, \Sigma) =$

$$\text{① } P(x^{(i)}|y^{(i)}; \phi) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}})\right\}$$

where $\mu_{y^{(i)}} = \begin{cases} \mu_0, & y^{(i)}=0 \\ \mu_1, & y^{(i)}=1 \end{cases} \Rightarrow \mu_{y^{(i)}} = \{y^{(i)}=0\} \mu_0 + \{y^{(i)}=1\} \mu_1$

$$\text{② } P(y^{(i)}; \phi) = \begin{cases} \phi & y=1 \Rightarrow P(y^{(i)}; \phi) \\ 1-\phi & y=0. \end{cases} \quad \{y^{(i)}=1\} \cup \{y^{(i)}=0\}$$

$$\therefore L = \sum_{i=0}^m \log P(x^{(i)}|y^{(i)}; \phi) + \sum_{i=0}^m \log P(y^{(i)}; \phi)$$

plug ① ② into L:

$$= -m \left(\log(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}} \right) - \frac{1}{2} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) +$$

$$\sum_{i=1}^m \{y^{(i)}=1\} \cdot \log \phi + \left(\sum_{i=1}^m \{y^{(i)}=0\} \right) \log(1-\phi)$$

$$= -m \left[\frac{d}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma| \right] + \dots$$

$$\Delta \frac{\partial L}{\partial \phi} = \frac{1}{\phi} \cdot \sum_{i=1}^m \{y^{(i)}=1\} + \frac{1}{1-\phi} \left[m - \sum_{i=1}^m \{y^{(i)}=1\} \right]$$

$$\frac{\partial L}{\partial \phi} = 0 \Rightarrow \phi = \frac{1}{m} \Delta = \frac{1}{m} \sum_{i=1}^m \{y^{(i)}=1\}$$

$$\Delta \frac{\partial L}{\partial \Sigma} = -\frac{1}{2} \cdot \Sigma^{-1} + \frac{1}{2} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})^T (x^{(i)} - \mu_{y^{(i)}}) \cdot \Sigma^{-1}$$

$$\frac{\partial L}{\partial \Sigma} = 0 \Rightarrow \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})^T (x^{(i)} - \mu_{y^{(i)}})$$

$$\frac{\partial l}{\partial \mu_0} = \frac{\partial L}{\partial \mu_{y^{(i)}}} \cdot \frac{\partial \mu_{y^{(i)}}}{\partial \mu_0} = \sum_{i=1}^m \frac{\partial}{\partial \mu_0} \left\{ -\frac{1}{2} \sum_{i=1}^m [x^{(i)T} - 2\mu_{y^{(i)}} x^{(i)} + \mu_{y^{(i)}}^2] \right\}$$

$$= \sum_i^m \left(-\frac{1}{2} \sum_{i=1}^m [-2x^{(i)} + 2\mu_{y^{(i)}}] \right) \quad \{ y^{(i)} = 0 \}$$

$$= \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}}) \cdot \{ y^{(i)} = 0 \}$$

maximize, hence

$$\frac{\partial l}{\partial \mu_0} = 0 \Rightarrow \mu_0 = \frac{\sum_{i=1}^m \{ y^{(i)} = 0 \} \cdot x^{(i)}}{\sum_{i=1}^m \{ y^{(i)} = 0 \}}$$

$$\text{By symmetry: } \mu_0 = \frac{\sum_{i=1}^m \{ y^{(i)} = 1 \} \cdot x^{(i)}}{\sum_{i=1}^m \{ y^{(i)} = 1 \}}$$

3. Incomplete, Positive - Only Labels $x \rightarrow t \quad \{ y \}$

(a). conditionally independent. By def.

$$P(y^{(i)}=1 | t^{(i)}=1, x^{(i)}) = P(t^{(i)}=1 | x^{(i)}) \cdot p(x^{(i)})$$

$$= p(t=1 | y=1, x).$$

$$= p(t=1 | y=1, x) \cdot p(y=1 | x) \cdot p(x)$$

$$\Delta p(t=1 | x) = p(y=1 | x) \cdot p(t=1 | y=1, x) \cdot p(x)$$

Given that $p(t^{(i)}=1 | y^{(i)}=1) = 1$

t & x conditionally independent

$$= \frac{p(y=1 | x)}{2}$$

$$\text{where } \alpha = p(y^{(i)}=1 | t^{(i)}=1, x^{(i)}) = p(y^{(i)}=1 | t^{(i)}=1)$$

(b) \sqrt{t} positive (labeled) examples in V

$$\text{i.e. } V_t = \{ x^{(i)} \in V \mid y^{(i)}=1 \}$$

$$\therefore p(t^{(i)}=1 | x^{(i)}) \approx 1 \text{ when } x \in V_t$$

Given that $\alpha = p(y^{(i)}=1 | t^{(i)}=1)$

$$h(x^{(i)}) = p(y^{(i)}) = 1 | x^{(i)} = \alpha \cdot p(t=1 | x) \approx \alpha$$

Take $\bar{\alpha}$ or weighted average.

$$\frac{p(y|x)}{p(x|y)}$$

3. Poisson Regression

(a). Poisson Regression.

$$p(y; \lambda) = \frac{\lambda^y e^{-\lambda}}{y!} = \frac{\exp(-\lambda + y \ln \lambda)}{y!}$$

Exponential Family: P_{22}

$$p(y; \eta) = b(y) \cdot \frac{\exp(\eta^T T(y) - \psi(\eta))}{1}$$

$b(y) = \frac{1}{y!}$
sufficient statistic

$T(y) = y$
natural parameter / canonical

$$\times \begin{cases} \eta = \ln \lambda \\ \psi(\eta) = \lambda = e^\eta \end{cases} \quad \text{partition function} \quad = e^\eta \quad \text{sigmoid function } P_{22}$$

cb) GLM. Generalized Linear Model

$$\begin{aligned} \text{p26} \quad T(y) &= y & h(x) &= E(y|x) \\ \text{canonical response fun} \quad h(x) &= E(y|x) = \lambda = e^{\theta^T x} & \theta^T x &= \frac{T(y; \eta)}{\eta} \end{aligned}$$

η : natural parameter

$$\text{(c). } \log p(y^{(i)} | x^{(i)}; \theta) = \log \frac{1}{y^{(i)}!} \exp[\theta^T x^{(i)} y^{(i)} - e^{\theta^T x^{(i)}}]$$

$$= -\log y^{(i)}! + \theta^T x^{(i)} y^{(i)} - e^{\theta^T x^{(i)}}$$

$$\frac{\partial \log p(\cdots)}{\partial \theta_j} = y^{(i)} \cdot x_j^{(i)} + e^{\theta^T x^{(i)}} \cdot x_j^{(i)}$$

$$= [y^{(i)} - e^{\theta^T x^{(i)}}] \cdot x_j^{(i)}$$

derive. \Downarrow back

p19. Stochastic gradient update ascent rule:

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

$$h_\theta(x^{(i)}) = \lambda = e^{\theta^T x^{(i)}}$$

LMS

4. Convexity of GLM

(a) $E[Y|X; \theta] = \int y p(y; \eta) dy \quad \text{①}$

$$\frac{\partial}{\partial \eta} \int p(y; \eta) dy = 0.$$

Hence,

$$\begin{aligned} \frac{\partial}{\partial \eta} \int p(y; \eta) dy &= \int \frac{\partial}{\partial \eta} p(y; \eta) dy \\ &= \int b(y) \cdot \exp(\eta y - \alpha(\eta)) \left[y - \frac{\partial \alpha(\eta)}{\partial y} \right] dy \\ &= \int p(y; \eta) \cdot \left[y - \frac{\partial \alpha(\eta)}{\partial y} \right] dy \\ &= E[Y|X; \eta] - \frac{\int p(y; \eta) dy \frac{\partial \alpha(\eta)}{\partial y}}{\int p(y; \eta) dy} \\ \Rightarrow E[Y; \eta] &= \frac{\partial \alpha(\eta)}{\partial y}. \end{aligned}$$

(b). $\frac{\partial}{\partial \eta} \int y p(y; \eta) dy = \frac{\partial^2 \alpha(\eta)}{\partial \eta^2}$

$$\frac{\partial}{\partial \eta} \int y p(y; \eta) dy = \int y \frac{\partial}{\partial \eta} p(y; \eta) dy \quad \text{the same way}$$

$$\begin{aligned} &= \int y^2 p(y; \eta) dy - \frac{\partial \alpha(\eta)}{\partial \eta} \cdot \frac{\int y \cdot p(y; \eta) dy}{m} \\ &= E[Y^2; \eta] - E[Y; \eta]^2 \\ &= \text{Var}[Y; \eta] = \text{Var}[Y|X; \eta]. \end{aligned}$$

(c) p13 negative log likelihood. P13

$$\begin{aligned} \text{negative log likelihood } l(\theta) &= - \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= - \sum_{i=1}^m \left[\log b(y^{(i)}) + \theta^T x^{(i)} y^{(i)} - \alpha(\theta^T x^{(i)}) \right]. \end{aligned}$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = \sum_{i=1}^m [\alpha'(\theta^T x^{(i)}) - y^{(i)}] x_j^{(i)}$$

$$H_{jk} = \frac{\partial^2 l(\theta)}{\partial \theta_j \partial \theta_k} = \sum_{i=1}^m \alpha''(\theta^T x^{(i)}) \cdot x_j^{(i)} \cdot x_k^{(i)}$$

$$\text{PSD} \Leftrightarrow Z^T H Z = \sum_i \sum_j \sum_k \alpha''(\theta^T x^{(i)}) \frac{x_j^{(i)} x_k^{(i)} z_j z_k}{(x^{(i)} \cdot Z)^2} > 0$$

5. Locally weighted Linear Regression.

case i. $w \in \mathbb{R}^m$

$$(X\theta - y)^T W (X\theta - y)$$

$$\begin{array}{c} i=j \\ \vdots \\ i \\ \vdots \\ m \end{array} \xrightarrow{\text{normal equation}} \begin{array}{c} X^T \\ \vdots \\ X \\ \vdots \\ X^T \end{array} \quad \begin{array}{c} w \\ \vdots \\ w \\ \vdots \\ w \end{array} \quad \begin{array}{c} Y \\ \vdots \\ Y \\ \vdots \\ Y \end{array}$$

$$\therefore w_{ij} = \begin{cases} \frac{1}{m} w^{(i)} & i=j \\ 0 & i \neq j. \end{cases}$$

$$\text{ii. } \nabla_{\theta} J(\theta) = \nabla_{\theta} (X\theta - y)^T W (X\theta - y)$$

$$= \nabla_{\theta} (\theta^T X^T - y^T) W (X\theta - y) \quad \text{tr } A = \text{tr } A^T.$$

$$= \nabla_{\theta} (\theta^T X^T W X \theta - y^T W X \theta - \underbrace{\theta^T X^T W y}_{\text{from eqn 1}} - y^T W y) \quad \cancel{\theta^T W y}$$

$$= \nabla_{\theta} (\theta^T X^T W X \theta - 2y^T W X \theta) \quad \sim R$$

$$= \nabla_{\theta} \text{tr} (\theta^T X^T W X \theta - 2y^T W X \theta) \quad \cancel{\theta^T W X} / 2 \cancel{\theta^T W X}$$

$$= X^T W X \theta + (X^T W X)^T \theta - 2(y^T W X)^T \quad \cancel{X^T \cdot W^T X}$$

$$(A^T W^T = W) \text{ by symmetry}$$

$$\nabla_{\theta} \text{tr} AB^T C = B^T A^T C + B A^T C \quad \cancel{B^T A^T C}$$

$$= 2 X^T W X \theta - 2 X^T W y$$

$$\nabla_{\theta} J(\theta) = 0 \Rightarrow \theta = (X^T W X)^{-1} X^T W y$$

$$\text{iii. } l(\theta) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta)$$

$$= \sum_{i=1}^m -\log(\sqrt{2\pi} \sigma^{(i)}) - \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = \sum_{i=1}^m \frac{y^{(i)} - \theta^T x^{(i)}}{(\sigma^{(i)})^2} \cancel{x_j^{(i)}}$$

(b) underfitting

(c) bandwidth parameter

$$T = 0.05$$

Python. src.

1.(b).

$x.$ shape

$m, n = x.$ shape. $(800, 3)$

$\begin{cases} m = x.$ shape [0] \\ $n = x.$ shape [1] \end{cases}

| | | | |
|--|-----|---|---|
| | 0 | 1 | 2 |
| | 1 | | |
| | | 1 | |
| | | | 1 |
| | 800 | | |

✓ Logistic regression (Newton's method).

$$\theta = \theta - H^{-1} \nabla_{\theta} L(\theta)$$

where: np.linalg.inv()

$$g(x) = \frac{1}{1+e^{-\theta^T x}}$$

$$H = g(X) \cdot [1 - g(X)] \cdot X^T \cdot X / m$$

$$|\theta_k' - \theta_k| < \epsilon$$

Prediction

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$$

1. (e) Logistic vs GDA

DB1: Logistic, GDA (Not Gaussian Distribution)

(h). Box-Cox transformation.

2. (e). decision boundary: $P(t^{(i)} = 1 | x^{(i)}) = \frac{1}{2}$.

$$\Rightarrow \frac{1}{2} P(y^{(i)} = 1 | x^{(i)}) = \frac{1}{2}$$

$$\frac{1}{2} \frac{1}{1+e^{-\theta^T x}} = \frac{1}{2}$$

$$\theta^T x + \ln\left(\frac{\theta^2}{2} - 1\right) = 0.$$

add to θ_0 .

3. (d). $\hat{x} = \frac{2 \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}}{m}$ LMS

~~step-size~~ $x = x_1 \dots x_4$ $x^T = x_1 \dots x_4$ $y, y^{(i)} (2500, 1)$
 $x_1 (1, 2500, v)$ $x_2 (2, 2500, v)$ $x_3 (3, 2500, v)$ $x_4 (4, 2500, v)$

if np.linalg.norm(theta - theta, ord=1) < self.eps:
break

$$w^{(i)} = \exp\left(-\frac{\|x^{(i)} - \bar{x}\|^2}{2\sigma^2}\right)$$

$$\theta = (x^T W x)^{-1} x^T W y$$

$$y' = \theta^T x.$$