CS 541 a3
Ghenxu Wang
10457625

# Gradient Calculation:

Suppose $x \in R^d$ and $y \in R$ are known.
Calculate the gradient of the following functions.

○ Sigmoid function: $f(w) = \dfrac{1}{1 + e^{-x \cdot w}}$

Since $\left(\dfrac{1}{u}\right)' = \dfrac{-u'}{u^2}$

Using Chain Rule:

$$\dfrac{\partial f(w)}{\partial w} = \dfrac{-e^{-xw} \cdot (-x)}{(1 + e^{-x \cdot w})^2} = \dfrac{x \cdot e^{-xw}}{(1 + e^{-xw})(1 + e^{-xw})}$$

$$= \dfrac{x}{(1 + e^{-xw})} \cdot \dfrac{1 + e^{-xw} - 1}{(1 + e^{-xw})} = x \cdot \dfrac{1}{1 + e^{-xw}} \cdot \left(1 - \dfrac{1}{1 + e^{-xw}}\right)$$

$$= x \cdot f(w) \cdot (1 - f(w))$$

○ Logistic loss $F(w) = \log(1 + e^{-y \cdot x \cdot w})$

Since $(\log u)' = \dfrac{1}{u}$

Using Chain Rule :

$$\frac{\partial}{\partial w} \log(1 + e^{-y \cdot x \cdot w})$$

$$= \frac{1}{1 + e^{-y \cdot x \cdot w}} \cdot e^{-y \cdot x \cdot w} \cdot (-y \cdot x)$$

$$= \frac{-y \cdot x}{e^{y \cdot x \cdot w} + 1}$$

# Linear Regression :

1. Since we have, for vector $z \in \mathbb{R}^d$, $z^T \cdot z = \sum_{i=1}^{d} z_i^2$

$\therefore \quad F(w) = \frac{1}{2} \| y - Xw \|_2^2$

$= \frac{1}{2} (y - X \cdot w)^T \cdot (y - Xw)$

$\nabla F(w) = \nabla \frac{1}{2} ( y^T \cdot y - y^T \cdot Xw - w^T \cdot x^T \cdot y + w^T x^T \cdot X \cdot w )$

$= \frac{1}{2} ( 0 - y^T x - x^T \cdot y + x^T X \cdot w + w^T \cdot x^T X )$

$= \frac{1}{2} ( -2 x^T y + 2 x^T X \cdot w )$

$= x^T X \cdot w - x^T \cdot y$

Let $\nabla F(w) = 0$, we have :

$$X^T X \cdot w = X^T y$$

When $n > d$, $(X^T X)^{-1}$ exists :

$$\Rightarrow \quad w = (X^T X)^{-1} \cdot X^T \cdot y$$

For Hessian matrix of $F(w)$:

$$H(F) = \frac{\partial^2 \bar{F}(w)}{\partial w \cdot \partial w^T} = \frac{\partial}{\partial w^T} \left( \frac{\partial F(w)}{\partial w} \right) = \frac{\partial}{\partial w^T} \cdot \nabla \bar{F}(w)$$

$$= \frac{\partial}{\partial w^T} \left( x^T x \cdot W - x^T y \right) = x^T x$$

$$= \begin{bmatrix} \sum\limits^{n} x_{i1}^2 & \sum\limits^{n} x_{i1} \cdot x_{i2} & \cdots & & \sum\limits^{n} x_{i1} \cdot x_{id} \\ \sum\limits^{n} x_{i2} \cdot x_{i1} & \sum\limits^{n} x_{i2}^2 & \cdots & & \sum\limits^{n} x_{i2} \cdot x_{id} \\ \vdots & \vdots & \ddots & & \vdots \\ \sum\limits^{n} x_{id} \cdot x_{i1} & \sum\limits^{n} x_{id} \cdot x_{i2} & \cdots & & \sum\limits^{n} x_{id}^2 \end{bmatrix}_{d \times d}$$

The matrix $x^T x \geq 0$, is positive semidefinite, which shows that $F(w)$ is a convex program.

2.

When we using the least squares formulation, it is equal to calculate the maximum likelihood estimation for the samples.

we can prove it:

The samples are $(x_i, y_i)$, the prediction is $\hat{y}_i|_w$, then we have $y = \hat{y} + \varepsilon$

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

$\therefore \quad y - \hat{y} \sim N(0, \sigma^2)$

$\Rightarrow \quad y \sim N(\hat{y}, \sigma^2)$

The likelihood function is:

$$L(w) = P(y|x; w) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - \hat{y})^2}{2\sigma}\right)$$

$$\Rightarrow \log L(w) = n \log \frac{1}{\sqrt{2\pi}} + \sum_{i=0}^{n} -\frac{(y_i - \hat{y}_i)^2}{2\sigma}$$

If we want the maximum of $\log L(w)$, $\sum_{i=0}^{m}(y - \hat{y})^2$ is minimal.

2. The least squares formulation will larger progress for getting to the optimum W. then $||y- \hat{y}||_2^{100}$

So the least squares will faster when we iterate $w^t = w^{t-1} - \eta \cdot \nabla (w^{t-1})$ to get $W_{opt}$. .

3.

When the rank of $X^T X$ : $Rank(X^T X) = d$, then

$H(\bar{F}) = X^T X \geqslant m \ (m > 0)$. $\bar{F}(w)$ is strongly-convex.

otherwise, is not strongly-convex.