

1. Given that $R_{tr}(\beta) = \frac{1}{N} \sum_{i=1}^N (y_i - \beta^T x_i)^2$
 $R_{te}(\beta) = \frac{1}{M} \sum_{i=1}^M (\hat{y}_i - \beta^T \hat{x}_i)^2$

A regression model with p parameters

$$y = XB + \epsilon, \text{ where } \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

solving $y = XB + \epsilon$.

I get the least square estimate of parameter β

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

$$\text{we know that } E(\hat{\beta}) = \beta$$

Let $(x_1, y_1) \dots (x_N, y_N)$ be the training data.

$(\hat{x}_1, \hat{y}_1) \dots (\hat{x}_M, \hat{y}_M)$ be the test data

$$E(y) = E(X\hat{\beta}) + E(\epsilon)$$

$$E(y) = X\beta$$

So expectation of y under both training data and test data give the same result.

$$\text{when } \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$y_i \sim N(x_i \beta, \sigma^2)$$

$$\hat{y}_i \sim N(\hat{x}_i \beta, \sigma^2)$$

We know that $s_i^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \hat{\beta}^T x_i)^2$ is unbiased

s_i^2 is estimation of σ^2 .

$$E(s_i^2) = E\left(\frac{1}{N-1} \sum_{i=1}^N (y_i - \hat{\beta}^T x_i)^2\right)$$

we already know that

$$\begin{aligned} E(R_{tr}(\hat{\beta})) &= E\left(\frac{1}{N} \sum_{i=1}^N (y_i - \hat{\beta}^T x_i)^2\right) \\ &= \frac{N-1}{N} E\left(\frac{1}{N-1} \sum_{i=1}^N (y_i - \hat{\beta}^T x_i)^2\right) \\ &= \frac{N-1}{N} E(s_i^2) \\ &= \frac{N-1}{N} \sigma^2 \end{aligned}$$

Similarly, let $s_i^2 = \frac{1}{M-1} \sum_{i=1}^M (\hat{y}_i - \hat{\beta}^T \hat{x}_i)^2$

which is an unbiased estimation of σ^2

$$E(s_i^2) = E\left(\frac{1}{M-1} \sum_{i=1}^M (\hat{y}_i - \hat{\beta}^T \hat{x}_i)^2\right)$$

$$\begin{aligned} E(R_{te}(\hat{\beta})) &= E\left(\frac{1}{M} \sum_{i=1}^M (\hat{y}_i - \hat{\beta}^T \hat{x}_i)^2\right) \\ &= \frac{M-1}{M} E\left(\frac{1}{M-1} \sum_{i=1}^M (\hat{y}_i - \hat{\beta}^T \hat{x}_i)^2\right) \\ &= \frac{M-1}{M} \sigma^2 \end{aligned}$$

$$\text{Since } E(R_{tr}(\hat{\beta})) = \frac{N-1}{N} \sigma^2$$

$$\text{and } E(R_{te}(\hat{\beta})) = \frac{M-1}{M} \sigma^2$$

and we have $N \geq M \geq p$

$$\text{we have } \frac{N-1}{N} \leq \frac{M-1}{M}$$

$$\text{thus } E(R_{tr}(\hat{\beta})) \leq E(R_{te}(\hat{\beta}))$$

2. Set $h(t) = f(x + t(y-x))$ for $t \in (0,1)$

$$h'(t) = \nabla f(x + t(y-x))^T (y-x)$$

Use equation 1.1

$$f(x) - h'(0) = f(x) - \nabla f(x)^T (y-x) = \nabla f(x)^T (y-x)$$

$$\begin{aligned}
 &= (\nabla f(x) - \nabla f(x+t(y-x)))^T \cdot (y-x) \\
 &\geq mt \|x-y\|_2^2 \\
 h'(t) &\geq h'(0) + mt \|x-y\|_2^2
 \end{aligned}$$

use FTC, I have

$$h(1) - h(0) = \int_0^1 h'(t) dt$$

$$f(y) - f(x) \geq \int_0^1 h'(0) + mt \|x-y\|_2^2 dt$$

$$\leq h'(0) + \frac{m}{2} \|x-y\|_2^2$$

$$\geq \nabla f(x)^T (y-x) + \frac{m}{2} \|x-y\|_2^2$$

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|x-y\|_2^2$$

Thus I proved that (1) is equivalent to

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|y-x\|_2^2$$

(b) assume $g(x) = f(x) - \frac{m}{2} \|x\|_2^2$, where $\nabla g(x) = \nabla f(x) - mx$.

Therefore I get

$$\begin{aligned}
 (\nabla g(x) - \nabla g(y))^T (x-y) &= (\nabla f(x) - \nabla f(y) - \frac{m}{2}(x-y))^T (x-y) \\
 &= (\nabla f(x) - \nabla f(y))^T (x-y) - (m(x-y))^T (x-y) \\
 &\geq 0
 \end{aligned}$$

Thus $g(x)$ is convex in domain of f .

Therefore $\nabla^2 g(x) \geq 0$ and I get $\nabla^2 f(x) - mI \geq 0$
 which means that I proved that (1) is equivalent
 to $\nabla^2 f(x) \geq mI$

3(a) $P_r(q_2 = \text{Happy}) = 0.8$

(b) $P_r(O_1 = \text{frown}) = 0.8 \times 0.1 + 0.2 \times 0.5$
 $= 0.08 + 0.1 = 0.18$

(c) $P_r(q_2 = \text{Happy} | O_1 = \text{frown}) = \frac{0.8 \times 0.1}{0.18} = \frac{4}{9}$

(d) $P_r(O_{100} = \text{yell}) = P_r(O_{100} = \text{yell} \cap q_{100} = \text{happy})$
 $+ P_r(O_{100} = \text{yell} \cap q_{100} = \text{angry})$
 $= 0.2$

(e) $q_1 = \text{happy}$

$P(f_2 | h_1) = 0.8 \times 0.1 = 0.08$
 $P(f_2 | a_1) = 0.2 \times 0.5 = 0.1$ so $q_2 = \text{angry}$

$P(f_3 | h_2) = 0.2 \times 0.1 = 0.02$ so $q_3 = \text{angry}$
 $P(f_3 | a_2) = 0.8 \times 0.5 = 0.4$

the most likely sequence is happy, angry, angry.

4.(a) The optimization problem now becomes $R(\theta) + \lambda J(\theta)$

which equals $\sum_{i=1}^n \sum_{k=1}^K (y_{ik} - f_k(x_i))^2 + \lambda (\sum_{k,m} B_{km}^2 + \sum_{m,i} \alpha_{mi}^2)$

taking the gradient of $R(\theta) + \lambda J(\theta)$

I get $\frac{\partial R + \partial \lambda J(\theta)}{\partial B_{km}} = \frac{\partial R}{\partial B_{km}} + \frac{\partial \lambda J}{\partial B_{km}}$
 $= 2k_i z_{mi} + 2\lambda \sum_{km} B_{km}$

$$\text{and } \frac{\partial R + \partial \lambda J(\theta)}{\partial a_{m1}} = \sum_i m_i x_{i1} + 2\lambda \sum_{m1} a_{m1}$$

According to the problem, a gradient update at the

$(r+1)$ st iteration has the form

$$B_{km}^{(r+1)} = B_{km}^{(r)} - \frac{\partial R}{\partial B_{km}^{(r)}} \quad a_{m1}^{(r+1)} = a_{m1}^{(r)} - \frac{\partial R}{\partial a_{m1}}$$

taking the new gradient update, $\frac{\partial R + \partial \lambda J(\theta)}{\partial B_{km}}$ and $\frac{\partial R + \partial \lambda J(\theta)}{\partial a_{m1}}$ into the gradient update.

$$\text{I get } B_{km}^{(r+1)} = B_{km}^{(r)} - \frac{\partial R}{\partial B_{km}^{(r)}} - \frac{\partial \lambda J(\theta)}{\partial B_{km}^{(r)}}$$

$$= B_{km}^{(r)} + 2(y_{ik} - f_k(x_i)) g'_k(B_k^T z_i) z_{mi} - 2\lambda \sum_{km} B_{km}^{(r)}$$

$$a_{m1}^{(r+1)} = a_{m1}^{(r)} + \sum_{k=1}^K 2(y_{ik} - f_k(x_i)) g'_k(B_k^T z_i) B_{km}^{(r)} (a_{m1}^{(r)} x_{i1}) - 2\lambda \sum_{m1} a_{m1}^{(r)}$$

Above are the gradient update for this regularized problem

3) Stochastic gradient descent replaces the actual gradient, which is calculated by the entire dataset, with an estimate. The estimate is calculated from a randomly subset of the data. This could largely reduce the computational burden, when n is large, achieving faster iterations in trade for a lower convergence rate