

1. What are the pros and cons between the regression hypotheses based on the objective

$$\arg_{\theta} \min \sum_{t=1}^N \left[r^{(t)} - h(\mathbf{x}^{(t)}; \theta) \right]^2$$

and $\arg_{\theta} \min \sum_{t=1}^N |r^{(t)} - h(\mathbf{x}^{(t)}; \theta)|$ respectively? (Hint: consider the training complexity and prediction accuracy)

A1: 令

$$\arg_{\theta} \min \sum_{t=1}^N \left[r^{(t)} - h(\mathbf{x}^{(t)}; \theta) \right]^2 \quad \text{為 1 式 Standard deviation}^*N$$

令

$$\arg_{\theta} \min \sum_{t=1}^N |r^{(t)} - h(\mathbf{x}^{(t)}; \theta)| : \quad \text{為 2 式 Mean absolute deviation}^*N$$

使用 Mean absolute deviation 的好處是：

單一計算較容易，不需要進行平方的運算

且較直覺，簡單，易於了解。

使用 Mean absolute deviation 的壞處是：

無法微分，因此，在求此線性方程組近似解時，會較困難。

使用 Standard deviation 的好處是：

可微分，因此容易找到最佳近似解 Eg: $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r}$

如果有解，則為最佳解。

使用 Standard deviation 的壞處是：

衡量誤差是根據其平方值，可能有較大誤差。

2. In logistic regression, show that $l(\beta) = \sum_{t=1}^N \left\{ y^{(t)} \beta^\top \tilde{x}^{(t)} - \log \left(1 + e^{\beta^\top \tilde{x}^{(t)}} \right) \right\}$.

A2:因為

$$\begin{aligned} \underline{l(\beta)} &= \log \prod_{t=1}^N p(\mathbf{x}^{(t)}, r^{(t)} | \beta) \\ &= \log \prod_{t=1}^N P(r^{(t)} | \mathbf{x}^{(t)}, \beta) p(\mathbf{x}^{(t)} | \beta) \\ &\propto \log \prod_{t=1}^N \pi(\mathbf{x}^{(t)}; \beta)^{q^{(t)}} (1 - \pi(\mathbf{x}^{(t)}; \beta))^{(1-q^{(t)})} \end{aligned}$$

$$\begin{aligned} l(\beta) &= \sum_{t=1}^N q^t * \log \pi(x^t; \beta) + (1 - q^t) * \log(1 - \pi(x^t; \beta)) \\ &= \sum_{t=1}^N \log(1 - \pi(x^t; \beta)) + q^t * [\log \pi(x^t; \beta) - \log(1 - \pi(x^t; \beta))] \end{aligned}$$

$$\begin{aligned} \because \log(1 - \pi(x^t; \beta)) &= \log\left(1 - \frac{1}{1 + e^{-B^T x}}\right) = \log\left(\frac{e^{-B^T x}}{1 + e^{-B^T x}}\right) = \log\left(\frac{1}{\frac{1}{e^{-B^T x}} + 1}\right) \\ &= -\log(1 + e^{B^T x}) \end{aligned}$$

$$\begin{aligned} \because q^t * [\log \pi(x^t; \beta) - \log(1 - \pi(x^t; \beta))] &= q^t \left[\log \frac{\pi(x^t; \beta)}{1 - \pi(x^t; \beta)} \right] = q^t \left[\log \frac{\frac{1}{1 + e^{-B^T x}}}{1 - \frac{1}{1 + e^{-B^T x}}} \right] \\ &= q^t \left[\log\left(\frac{1}{1 + e^{-B^T x} - 1}\right) \right] = q^t \left[\log\left(\frac{1}{e^{-B^T x}}\right) \right] = -q^t [\log e^{-B^T x}] = q^t B^T x \end{aligned}$$

$$\therefore l(\beta) = \sum_{t=1}^N \{-\log(1 + e^{B^T x}) + q^t B^T x\}$$

得證

3. Read Appendix C on the definitions of convex set and functions.

- (a) Show that the intersection of convex sets, $\bigcap_{i \in \mathbb{N}} C_i$ where $C_i \subseteq \mathbb{R}^n$, is convex.
(b) Show that the log-likelihood function for logistic regression, $l(\beta)$, is concave.

A3:

(a) 依照定理

Given a convex set $C_1, C_2 \subseteq \mathbb{R}^n$,

- Scaling: $\beta C = \{\beta \mathbf{x} : \mathbf{x} \in C\}$ is convex for any $\beta \in \mathbb{R}$
- Sum: $C_1 + C_2 = \{\mathbf{x}_1 + \mathbf{x}_2 : \mathbf{x}_1 \in C_1, \mathbf{x}_2 \in C_2\}$ is convex

設 $\mathbf{x}, \mathbf{y} \in C_1 \cap C_2$ 且 $0 \leq \alpha \leq 1$

$\therefore \mathbf{x}, \mathbf{y}$ 屬於 C_1 與 C_2

因(依題目設定, C_1, C_2 皆為 convex set)

$\therefore \alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in C_1 \cap C_2$ <-using 定理 Scaling and Sum

證明 $C_1 \cap C_2$ 為凸集合

重複上述動作 證明 $C_1 \cap C_2 \cap C_3$ 為凸集合

令此交集為 C_{123}

重複上述動作直到 證明 $C_1 \cap C_2 \cap \dots \cap C_n$ 為凸集合

得證.

(b)(符號使用與 a 小題不同) 依照 concave 定義

A real-valued function f on an interval (or, more generally, a convex set in vector space) is said to be concave if, for any \mathbf{x} and \mathbf{y} in the interval and for any t in $[0, 1]$,

$$f((1-t)\mathbf{x} + (t)\mathbf{y}) \geq (1-t)f(\mathbf{x}) + (t)f(\mathbf{y}).$$

令 $\gamma \alpha + (1 - \gamma) \beta \in C$ (convex set)

對於任意 $\alpha, \beta \in C$ 且 $0 \leq \gamma \leq 1$

$$l(\beta) = \sum_{t=1}^N \left\{ y^{(t)} \beta^T \tilde{\mathbf{x}}^{(t)} - \log(1 + e^{\beta^T \tilde{\mathbf{x}}^{(t)}}) \right\}$$

$$l(\gamma \alpha + (1 - \gamma) \beta) = \sum_{t=1}^N \left\{ y^{(t)} (\gamma \alpha + (1 - \gamma) \beta)^T \tilde{\mathbf{x}}^{(t)} - \log(1 + e^{(\gamma \alpha + (1 - \gamma) \beta)^T \tilde{\mathbf{x}}^{(t)}}) \right\}$$

$$\gamma l(\alpha) + (1 - \gamma) l(\beta) =$$

$$\begin{aligned} & \sum_{t=1}^N \left\{ \gamma y^{(t)} \alpha^T \tilde{\mathbf{x}}^{(t)} - \log(1 + e^{(\gamma \alpha)^T \tilde{\mathbf{x}}^{(t)}}) \right\} + \sum_{t=1}^N \left\{ (1 - \gamma) y^{(t)} \beta^T \tilde{\mathbf{x}}^{(t)} - \log(1 + e^{(\gamma \beta)^T \tilde{\mathbf{x}}^{(t)}}) \right\} \\ &= \sum_{t=1}^N \left\{ y^{(t)} \{ \gamma \alpha^T + (1 - \gamma) \beta^T \} \tilde{\mathbf{x}}^{(t)} - \gamma \{ \log(1 + e^{(\gamma \alpha)^T \tilde{\mathbf{x}}^{(t)}}) \} - (1 - \gamma) \{ \log(1 + e^{(\gamma \beta)^T \tilde{\mathbf{x}}^{(t)}}) \} \right\} \\ &= \sum_{t=1}^N \left\{ y^{(t)} \{ \gamma \alpha^T + (1 - \gamma) \beta^T \} \tilde{\mathbf{x}}^{(t)} - \{ \log(1 + e^{(\gamma \alpha)^T \tilde{\mathbf{x}}^{(t)}})^{\gamma} + \log(1 + e^{(\gamma \beta)^T \tilde{\mathbf{x}}^{(t)}})^{(1-\gamma)} \} \right\} \end{aligned}$$

顯然

$$l(\gamma \alpha + (1 - \gamma) \beta) \geq \gamma l(\alpha) + (1 - \gamma) l(\beta)$$

所以 $l(\beta)$ is concave

4. Consider the locally weighted linear regression problem with the following objective:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \frac{1}{2} \sum_{i=1}^N l^{(i)}(\mathbf{w}^\top \begin{bmatrix} 1 \\ \mathbf{x}^{(i)} \end{bmatrix} - r^{(i)})^2$$

local to a given instance \mathbf{x}' whose label will be predicted, where $l^{(i)} = \exp(-\frac{(\mathbf{x}' - \mathbf{x}^{(i)})^2}{2\tau^2})$ for some constant τ .

(a) Show that the above objective can be written as the form

$$(\mathbf{X}\mathbf{w} - \mathbf{r})^\top \mathbf{L}(\mathbf{X}\mathbf{w} - \mathbf{r}).$$

Specify clearly what \mathbf{X} , \mathbf{r} , and \mathbf{L} are.

(b) Give a close form solution to \mathbf{w} . (Hint: recall that we have $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{r}$ in linear regression when $l^{(i)} = 1$ for all i)

(c) Suppose that the training examples $(\mathbf{x}^{(i)}, r^{(i)})$ are i.i.d. samples drawn from some joint distribution with the marginal:

$$p(r^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^{(i)}}} \exp\left(-\frac{(r^{(i)} - \mathbf{w}^\top \begin{bmatrix} 1 \\ \mathbf{x}^{(i)} \end{bmatrix})^2}{2\sigma^{(i)2}}\right)$$

where $\sigma^{(i)}$'s are constants. Show that finding the maximum likelihood of \mathbf{w} reduces to solving the locally weighted linear regression problem above. Specify clearly what the $l^{(i)}$ is in terms of the $\sigma^{(i)}$'s.

A4: Local Weight Linear Regression
$$\sum_i l(\mathbf{x}^{(i)}; \mathbf{x}') (r^{(i)} - \mathbf{w}^\top \begin{bmatrix} 1 \\ \mathbf{x}^{(i)} \end{bmatrix})^2$$

(a) 自行定義符號 \mathbf{l} 為 $n \times 1$ 矩陣 \mathbf{l}_i 為 $l(\mathbf{x}^{(i)}; \mathbf{x}')$, $\mathbf{L} = (\sqrt{\mathbf{l}})(\sqrt{\mathbf{l}})^\top$ 為 $n \times n$ 矩陣

$$\begin{aligned} \text{原式可寫成} &= \left((\sqrt{\mathbf{l}})^\top (\mathbf{X}\mathbf{w} - \mathbf{r}) \right)^\top \left((\sqrt{\mathbf{l}})^\top (\mathbf{X}\mathbf{w} - \mathbf{r}) \right) = (\mathbf{X}\mathbf{w} - \mathbf{r})^\top (\sqrt{\mathbf{l}})(\sqrt{\mathbf{l}})^\top (\mathbf{X}\mathbf{w} - \mathbf{r}) \\ &= (\mathbf{X}\mathbf{w} - \mathbf{r})^\top \mathbf{L}(\mathbf{X}\mathbf{w} - \mathbf{r}) \quad \text{得證} \end{aligned}$$

(b) $(\mathbf{X}\mathbf{w} - \mathbf{r})^\top \mathbf{L}(\mathbf{X}\mathbf{w} - \mathbf{r})$, 對 \mathbf{W} 作微分, 並令其一階微分為零

$$\because \nabla((\mathbf{X}\mathbf{w} - \mathbf{r})^\top \mathbf{L}(\mathbf{X}\mathbf{w} - \mathbf{r})) = (\mathbf{X}^\top)(\mathbf{L} + \mathbf{L}^\top)(\mathbf{X}\mathbf{w} - \mathbf{r}) = 2(\mathbf{X}^\top)(\mathbf{L})(\mathbf{X}\mathbf{w} - \mathbf{r}) = 0$$

$$\therefore \mathbf{w} = (\mathbf{X}^\top \mathbf{L} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{L} \mathbf{r}$$

(c) $\because P(\mathbf{W}|\mathbf{X}) = \frac{P(\mathbf{X}|\mathbf{W}) * P(\mathbf{W})}{P(\mathbf{X})} \propto P(\mathbf{X}|\mathbf{W})$, and we assume $P(\mathbf{W})$ all same, ' (\mathbf{X}, \mathbf{r}) are i.i.d.'

\because by maximizing likelihood \mathbf{W}

$$\therefore \log P(\mathbf{X}|\mathbf{W}) = \log\left(\prod_{t=1}^N P(\mathbf{x}^t, r^t | \mathbf{W})\right) = \log\left(\prod_{t=1}^N P(r^t | \mathbf{x}^t, \mathbf{W}) * P(\mathbf{x}^t | \mathbf{W})\right)$$

$$\propto \log\left(\prod_{t=1}^N \frac{1}{\sqrt{2\pi\sigma^t}} e^{-\frac{(r^t - \mathbf{w}^t \begin{bmatrix} 1 \\ \mathbf{x}^t \end{bmatrix})^2}{2(\sigma^t)^2}}\right) = \sum_{t=1}^N -\log(\sqrt{2\pi\sigma^t}) - \frac{(r^t - \mathbf{w}^t \begin{bmatrix} 1 \\ \mathbf{x}^t \end{bmatrix})^2}{2(\sigma^t)^2}$$

$$\propto \sum_{t=1}^N -\frac{(r^t - \mathbf{w}^t \begin{bmatrix} 1 \\ \mathbf{x}^t \end{bmatrix})^2}{2(\sigma^t)^2} = -\frac{1}{2} \sum_{t=1}^N \frac{1}{(\sigma^t)^2} (r^t - \mathbf{w}^t \begin{bmatrix} 1 \\ \mathbf{x}^t \end{bmatrix})^2, \quad l^t = \frac{1}{(\sigma^t)^2} \quad \text{得證.}$$

Coding problem

- (d) Implement a linear regressor (see the spec for more details) on the provided 1D dataset. Plot the data and your fitted line. (Hint: don't forget the intercept term)
- (e) Implement 4 locally weighted linear regressors (see the spec for more details) on the same dataset with $\tau = 0.1, 1, 10$, and 100 respectively. Plot the data and your 4 fitted curves (for different x 's within the dataset range).
- (f) Discuss what happens when τ is too small or large.