

Problem 1

(1) set a random variable $X \in [-2, 2]$, $t=1$

$$\Pr[X \geq t] = \frac{1}{4} \geq 0 = \frac{E[X]}{t}.$$

so Markov's inequality doesn't hold for negative variable.

(2) For Chebyshev's Inequality: $\Pr[|X - E(X)| \geq t] \leq \frac{\text{Var}(X)}{t^2}$, $\forall t > 0$

$$\because E(\bar{X}) = E(X), \text{Var}(\bar{X}) = E\left[\left(\frac{1}{n} \sum_{i=1}^n x_i\right)^2\right] = \frac{\text{Var}(X)}{n} \text{ for } X \text{ i.i.d and } E(X)=0$$

$\therefore \Pr[|\bar{X} - E(X)| \geq t] \leq \frac{\text{Var}(X)}{nt^2}$, so large n and small variance imply better concentration.

(3) Hoeffding's inequality has tighter bound because it's decay exponentially $O(e^{-t^2})$

4). Only when all samples are i.i.d, we have $E_{\text{out}}(f) = E_{\text{sup}}[E_{\text{in}}(f)]$, which means we can use $E_{\text{in}}(f)$ as unbiased estimator of $E_{\text{out}}(f)$.

5). From $E_{\text{out}}(f) \leq E_{\text{in}}(f) + R(H) + \sqrt{\frac{\log \frac{1}{\delta}}{2n}}$ we know, more samples (bigger n), tighter bound, smaller error.

$$6). \text{ We know } \hat{R}_s(H) := E_{\epsilon} \left[\sup_{f \in H} \frac{1}{n} \sum_{i=1}^n \epsilon_i f_{\theta}(x_i) \right]$$

$$\text{when } |H|=1, \hat{R}_s(H) = E_{\epsilon} \left[\frac{1}{n} \sum_{i=1}^n \epsilon_i f_{\theta}(x_i) \right] = \sum_{i=1}^n f_{\theta}(x_i) E(\epsilon_i) = 0$$

$$\text{when } |H|=2^n, \hat{R}_s(H) = E_{\epsilon} \left[\sup_{f \in H} \frac{1}{n} \sum_{i=1}^n \epsilon_i f_{\theta}(x_i) \right] = E_{\epsilon} \left[\frac{1}{n} \sum_{i=1}^n 1 \right] = E[1] = 1$$

7) Firstly decide a ^{not so complex} hypothesis space H , then choose one optimization algorithm to get $f_{\hat{\theta}} \in H$, which minimize our designed loss function. So we can have small generalization error.

Problem 2. Define $h(s) := \sup_{f \in H} [E_{\text{out}}(f) - E_{\text{in}}(f; s)]$

Let $S = \{z_1, z_2, \dots, z_n\}$, $S' = \{z_1, \dots, z'_n, \dots, z_n\}$

$$h(s') - h(s) = \sup_{f \in H} [E_{\text{out}}(f) - E_{\text{in}}(f; s')] - \sup_{f \in H} [E_{\text{out}}(f) - E_{\text{in}}(f; s)]$$

$$\leq \sup_{f \in H} [E_{\text{in}}(f; s) - E_{\text{in}}(f; s')]$$

$$= \sup_{f \in H} \left[\frac{e(f(x_1), y_1) - e(f(x'_1), y'_1)}{n} \right]$$

$$\leq \frac{1}{n} \text{ since } e(f(x), y) \rightarrow [0, 1]$$

Apply McDiarmid's inequality to $h(S)$, we have

$$\Pr[h(S) - \mathbb{E}[h(S)] \geq t] \leq e^{-2nt^2}, \forall t > 0$$

Let $\delta = e^{-2nt^2}$, we then have with probability at least $1 - \delta$

$$h(S) \leq \mathbb{E}[h(S)] + \sqrt{\frac{\log \frac{1}{\delta}}{2n}}$$

since $h(S) = \sup_{f \in H} [E_{\text{out}}(f) - E_{\text{in}}(f)] \geq E_{\text{out}}(f) - E_{\text{in}}(f)$ for any $f \in H$

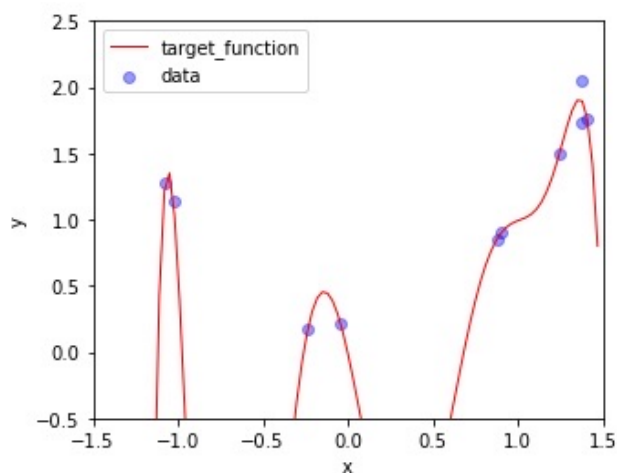
$$\text{we get } \forall f \in H, E_{\text{out}}(f) \leq E_{\text{in}}(f) + \mathbb{E}[h(S)] + \sqrt{\frac{\log \frac{1}{\delta}}{2n}}$$

$$\Rightarrow \mathbb{E}[h(S)] \leq R(H) \quad \text{so } E_{\text{out}}(f) \leq E_{\text{in}}(f) + R(H) + \sqrt{\frac{\log \frac{1}{\delta}}{2n}}$$

Problem 3

(a1): y is a $[10, 1]$ matrix, while X is a $[10, 9]$ vandermonde matrix.

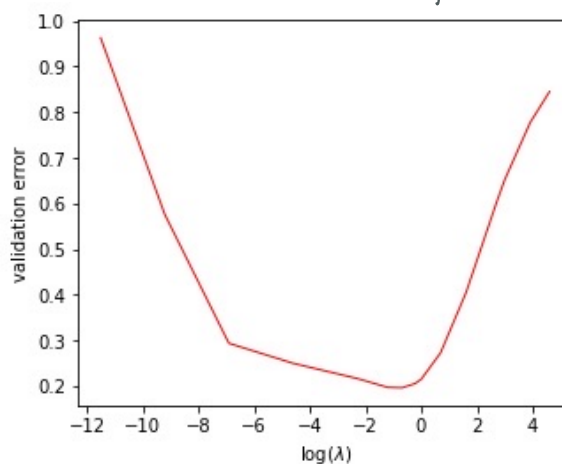
(a2):



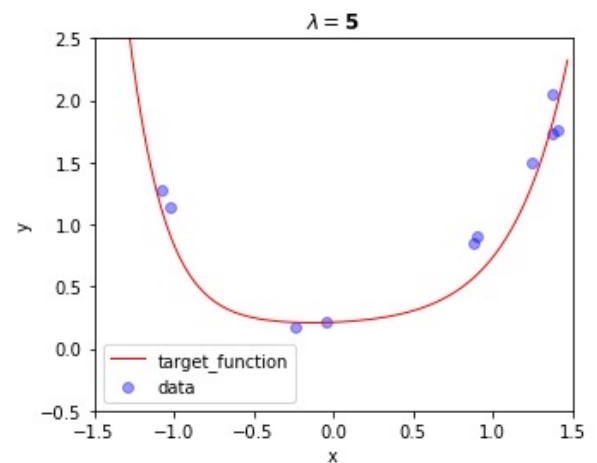
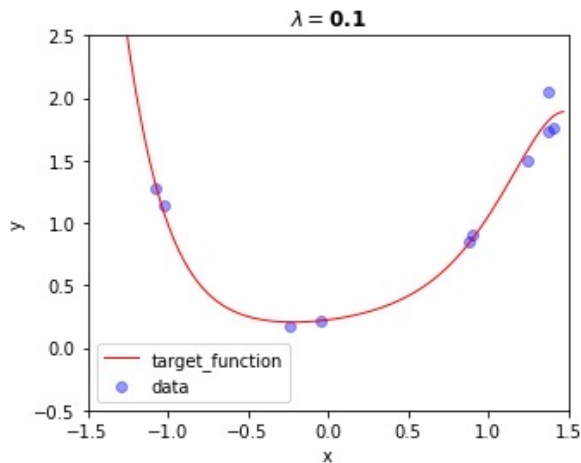
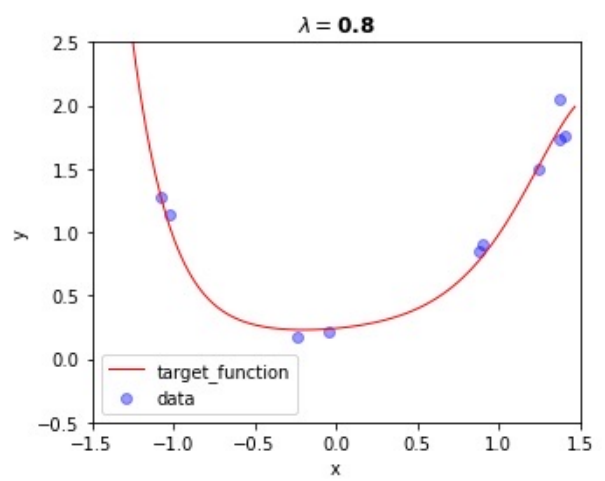
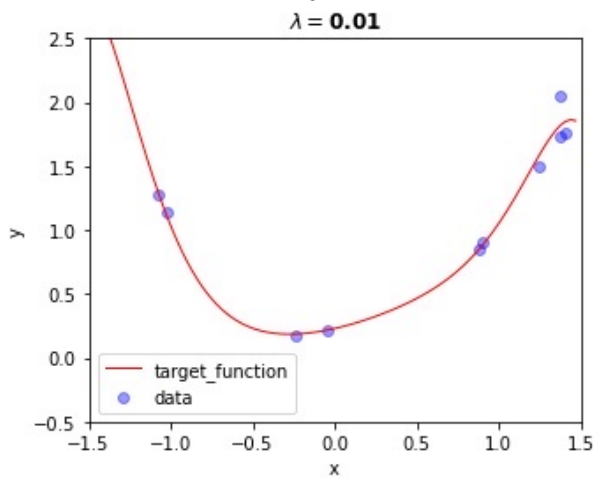
the overfitting is very obvious.

(a3) the test error is = 177.231208924576

(b1) the plot of validation error versus the value of λ is as follow.



(b2) the fitted curve using the four choice of λ is:



(b3) the test error for $\lambda = 0.01, 0.1, 0.8$, and 5 is $(1.21, 2.69, 4.11, 3.49)$

Problem 4

(a) the loss function can be define as:

$$U(\theta) := \sum_{n=1}^N [\theta_{y_n}^T X_n - \log \sum_{c=1}^{\text{Class}} \exp(\theta_c^T X_n)]$$

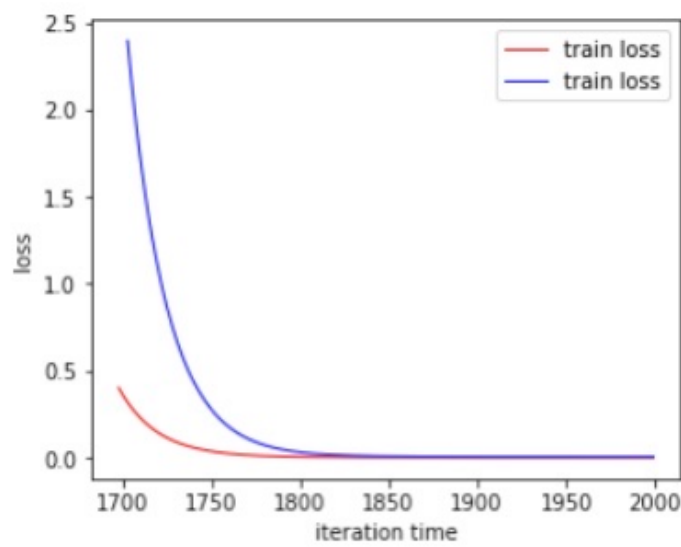
for the first term:

$$\frac{\partial \theta_{y_n}^T X_n}{\partial \theta_{y_n^{(j)}}} \theta_{y_n}^T X_n := [y_n = y'] X_n^{(j)}$$

for the second term:

$$\begin{aligned} \frac{\partial}{\partial \theta_{y_n^{(j)}}} \log \sum_{c=1}^{\text{Class}} \exp(\theta_c^T X_n) &:= \frac{\sum_{c=1}^{\text{Class}} \exp(\theta_c^T X_n) \times [y = y'] X_n^{(j)}}{\sum_{c=1}^{\text{Class}} \exp(\theta_c^T X_n)} \\ &= \sum_{c=1}^{\text{Class}} \left[\frac{\exp(\theta_c^T X_n)}{\sum_{c=1}^{\text{Class}} \exp(\theta_c^T X_n)} \times [y_n = y'] X_n^{(j)} \right] \\ &= \sum_{c=1}^{\text{Class}} [p(y=c | X_n) \times [y_n = y'] X_n^{(j)}] \\ &= p(y=y' | X_n) X_n^{(j)} \end{aligned}$$

Applying AGD we can get test loss and training loss decreasing as follow:



Then, we can further calculate the training error is 56.31%, and test accuracy is 43.58%. Compared with random probability, which is $\frac{1}{4} = 25\%$, multinomial logistic regression obviously has higher odds ratio.

Finally, comparing the corresponding parameters' norm, we find the most important feature is "RAM".

P.S: for more details, please see "p4.ipynb" in my folder.

1b) The learning problem can be formulated as:

$$\min_{\theta, Z} \ell(\theta, Z) := -\frac{1}{n} \cdot \mathbb{1}_{1 \times n} \cdot \{P \odot \log(h(X \cdot \theta \cdot \mathbb{1}_{1 \times K} - \mathbb{1}_{n \times 1} \cdot Z^T) \cdot M_{KK})\} \cdot \mathbb{1}_{K \times 1}$$

$$\theta = \mathbb{R}^{d \times 1}, \quad Z = \mathbb{R}^{K \times 1}$$

$$P = \mathbb{R}^{n \times K} \quad P_{ij} := \begin{cases} 1, & \text{if } f_i = j \\ 0, & \text{others} \end{cases}$$

$$M = \mathbb{R}^{K \times K} := \begin{pmatrix} 1 & -1 & 0 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

then, we get the gradient as:

$$\frac{\partial \ell(\theta, Z)}{\partial \theta} := -\frac{1}{n} \cdot (\mathbb{1}_{1 \times K} \otimes \mathbb{1}_{1 \times n}) \cdot \text{diag}(\text{vec}(P)) \cdot \partial_a \cdot \partial_{\log} \cdot (M^T \otimes \mathbb{I}_{n \times n}) \cdot \partial_h \cdot (\mathbb{1}_{K \times 1} \otimes X)$$

$$\frac{\partial \ell(\theta, Z)}{\partial Z} := \frac{1}{n} \cdot (\mathbb{1}_{1 \times K} \otimes \mathbb{1}_{1 \times n}) \cdot \text{diag}(\text{vec}(P)) \cdot \partial_a \cdot \partial_{\log} \cdot (M^T \otimes \mathbb{I}_{n \times n}) \cdot \partial_h \cdot (\mathbb{1}_{K \times 1} \otimes \mathbb{I}_{K \times K})$$

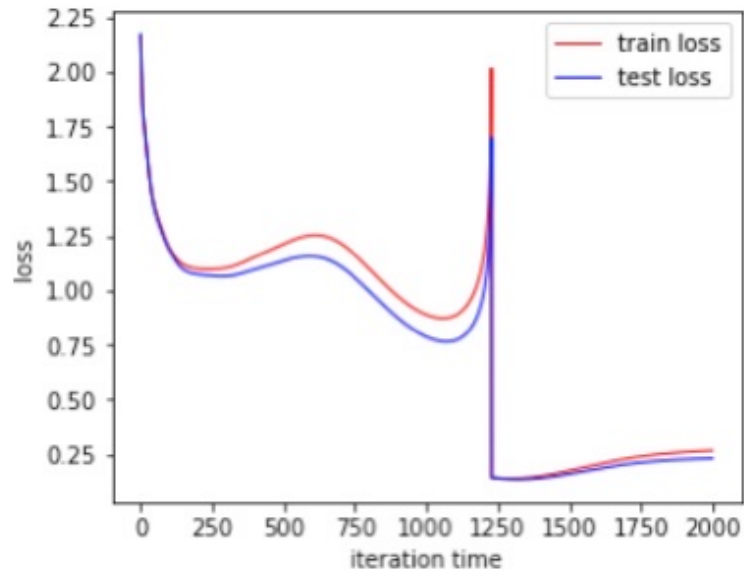
$$\text{where } H = h(X \cdot \theta \cdot \mathbb{1}_{1 \times K} - \mathbb{1}_{n \times 1} \cdot Z^T) \in \mathbb{R}^{n \times d}$$

$$Q_a := \text{diag}(\text{vec}(P)) \quad R^{nK \times nK}$$

$$Q_{\log} := \text{diag}(\text{vec}(H \otimes M)) \quad R^{nK \times nK}$$

$$Q_h := \text{diag}(\text{vec}(H')) \quad R^{nK \times nK}$$

By applying AGD, we can get loss decreasing as follow:



And the accuracy for train data is 76.66%, while the accuracy for train data is 73%. Both of them get improved greatly by ordinal logistic regression. So we can say, for ordinal label, it's better to use ordinal logistic regression than multinomial logistic regression.

As before, comparing the corresponding parameters norm, we find the most important feature is "RAM".