

Question 1— coursera

Question 2

半監督學習 (semi-supervised learning) 在隱私保護上的應用。例如銀行業者希望利用機器學習的方法評估客戶的信用進行信貸風險管理。這一情境下如果利用所有資料進行建模就有可能涉及到客戶的個人隱私問題。半監督學習只需要利用一部分已標記的數據和大量未標記的數據進行建模，這樣我們可以避免使用一些客戶敏感資訊 (例如身份訊息、居住地、手機號碼等)，更關注某一類客戶在行為上有什麼共性。

3. $\because f$ 在 $D = \{(x_n, y_n)\}_{n=1}^N$ 上 $f(x_n) = y_n$, where $y_n \in Y = \{-1, +1\}$

f 在 test input $\{x_{N+1}, x_{N+2}, x_{N+3} \dots x_{N+L}\}$ 上 $f(x_{N+L}) = -1$ or $+1$

\therefore 对于 test set, f 产生 $\{y_{N+1}, y_{N+2}, y_{N+3} \dots y_{N+L}\}$ 有 2^L 种可能的情况

$$E_{OTS}(g, f) = \frac{1}{L} \sum_{l=1}^L [g(x_{N+l}) \neq f(x_{N+l})] \quad \text{where } A(D) = g$$

$$\begin{cases} g(x_{N+l}) \neq f(x_{N+l}), & OTS_l = 1 \\ g(x_{N+l}) = f(x_{N+l}), & OTS_l = 0 \end{cases} \quad \text{OTS: off-Training-Set error.}$$

\therefore 存在 $f(x_{N+l})$ 使得 $\sum_{l=1}^L OTS_l = \sum_{l=1}^L [g(x_{N+l}) \neq f(x_{N+l})]$ 的值为 $[0, L]$ 中的任意一个整数

\therefore 当 $\sum_{l=1}^L OTS_l = \sum_{l=1}^L [g(x_{N+l}) \neq f(x_{N+l})] = n$ 时, where $n \in [0, L]$

满足上述情况的 f 有 C_n^L 个.

$\therefore f$ are equally likely in probability.

$$\therefore P\left(\sum_{l=1}^L OTS_l = n\right) = \frac{C_n^L}{2^L}$$

$$\therefore E_f \{E_{OTS}(A(D), f)\} = \sum E_{OTS}(A(D), f) \cdot P$$

二项式定理

$$(x+y)^n = \sum_{k=0}^n C_k^n x^{n-k} y^k = \sum_{k=0}^n C_k^n x^k y^{n-k} = \frac{1}{L} \cdot 0 \cdot \frac{C_0^L}{2^L} + \frac{1}{L} \cdot 1 \cdot \frac{C_1^L}{2^L} + \dots + \frac{1}{L} \cdot L \cdot \frac{C_L^L}{2^L}$$

$$(1+x)^n = \sum_{k=0}^n C_k^n x^k = C_0^n + x C_1^n + x^2 C_2^n + \dots + x^n C_n^n = \frac{1}{L} \cdot \frac{1}{2^L} \cdot (0 \cdot C_0^L + 1 \cdot C_1^L + \dots + L \cdot C_L^L)$$

$$\frac{d}{dx} (1+x)^n = n(1+x)^{n-1} = \frac{1}{L} \cdot \frac{1}{2^L} \cdot (L \cdot 2^{L-1}) \quad \downarrow \text{见左侧.}$$

$$= C_1^n + 2x C_2^n + \dots + n x^{n-1} C_n^n$$

$$\text{当 } x=1 \text{ 时 } n \cdot 2^{n-1} = 1 \cdot C_1^n + 2 \cdot C_2^n + \dots + n C_n^n$$

$$= \frac{1}{2}$$

\therefore 对于任意 A , 上式与 A 无关 $E_f \{E_{OTS}(A(D), f)\} = \frac{1}{2} \quad \#$

R08944052 斯晓宇

4. 设 $P(A)$: 选中 A 类骰子的机率

$P(D)$: 选中 D 类骰子的机率

$$\therefore \text{green } 1 \in A \text{ or } D \quad P(A) \perp P(B) \perp P(C) \perp P(D)$$

$$\therefore P(\text{green } 1) = P(A \cup D) = P(A) + P(D) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

$$\therefore P(\text{pick five green } 1\text{'s}) = \left(\frac{1}{2}\right)^5 = \frac{1}{32}$$

5.

	green	
	Even	Odds
A		1, 3, 5
B	2, 4, 6	
C	4, 6	5
D	2	1, 3

✓ green 1 \in A or D

green 2 \in B or D

✓ green 3 \in A or D

• green 4 \in B or C

green 5 \in A or C

• green 6 \in B or C

抽 5 个骰子,

某一种数字全为绿色的机率。

选中某一数字全为绿色会有 4 种可能的情况

A or D — 1, 3

B or D — 2

B or C — 4, 6

A or C — 5

$$P(\text{"some number" that is purely green}) = [P(A \cup D)]^5 + [P(B \cup D)]^5 + [P(B \cup C)]^5 +$$

$$[P(A \cup C)]^5 - [P(A)]^5 - [P(B)]^5 - [P(C)]^5 - [P(D)]^5$$

$[P(A \cup D)]^5, [P(A \cup C)]^5$ 都包含 5 次全为 A 类的情况, 故 $- [P(A)]^5$

依此类推需要减去 5 次全为某个类别的机率

$$\begin{aligned} \therefore P(\text{"some number" that is purely green}) &= \left(\frac{1}{2}\right)^5 + \left(\frac{1}{2}\right)^5 + \left(\frac{1}{2}\right)^5 + \left(\frac{1}{2}\right)^5 - \left(\frac{1}{4}\right)^5 - \left(\frac{1}{4}\right)^5 - \left(\frac{1}{4}\right)^5 - \left(\frac{1}{4}\right)^5 \\ &= \frac{4}{2^5} - \frac{4}{4^5} = \frac{4}{2^5} - \frac{4}{2^{10}} = \frac{4(2^5 - 1)}{2^{10}} = \frac{31}{256} \end{aligned}$$

Findings 在下一页!

5. Findings:

可以把 problem 4-5 抽象为一个 Feasibility of Learning 的问题.

Bag \rightarrow 巨大的 data set

A, B, C, D \rightarrow 四种 data 类型

dice $\begin{cases} \text{绿色的数字} & h(x) \neq f(x) \\ \text{橙色的数字} & h(x) = f(x) \end{cases}$

1, 2, 3, 4, 5, 6 \rightarrow hypothesis $h_1, h_2, h_3, h_4, h_5, h_6$

pick 5 dice from the bag \rightarrow sample $N=5$.

依题意对于任意一个 h_n $E_{out}(h_n) = \frac{1}{2}$

当 $N=5$ 且抽到的是 5 个 green dice, 此时 $E_{in} = 1$, $|E_{out} - E_{in}| = 0.5$

所以该情况可视为 Bad D

Question 4 得到: $P_D[\text{Bad D for } h_n] = \frac{1}{32} = \frac{8}{256}$

Question 5 得到: $P_D[\text{Bad D}] = \frac{31}{256} \leq 4 \cdot P_D[\text{Bad D for } h_n] = \frac{32}{256}$

根据 Hoeffding's Inequality:

$$\begin{aligned} \text{Bound of Bad Data } P_D[\text{Bad D}] &= P_D[\text{Bad D for } h_1 \text{ or Bad D for } h_2 \dots \text{ or Bad D for } h_6] \\ &\leq P_D[\text{Bad D for } h_1] + P_D[\text{Bad D for } h_2] + \dots + P_D[\text{Bad D for } h_6] \\ &= 6 \cdot P_D[\text{Bad D for } h_n] \end{aligned}$$

所以理论上 $P_D[\text{Bad D}] \leq 6 \cdot P_D[\text{Bad D for } h_n]$

但在 Question 5 中 $P_D[\text{Bad D}] \leq 4 \cdot P_D[\text{Bad D for } h_n]$

我认为这是因为 Question 5 中只有 4 种情况

所以使得有效计算 Bound of Bad Data 的

hypothesis 变为只有 4 种。

$\begin{cases} A \text{ or } D \rightarrow h_1, h_3 \\ B \text{ or } D \rightarrow h_2 \\ B \text{ or } C \rightarrow h_4, h_6 \\ A \text{ or } C \rightarrow h_5 \end{cases}$

说明在 hypothesis set 中有一些 h 是相似的,

找到实际有效的 hypothesis 可以更精确估计 Bound.

Question6

The average number of updates is 39.22113676731794

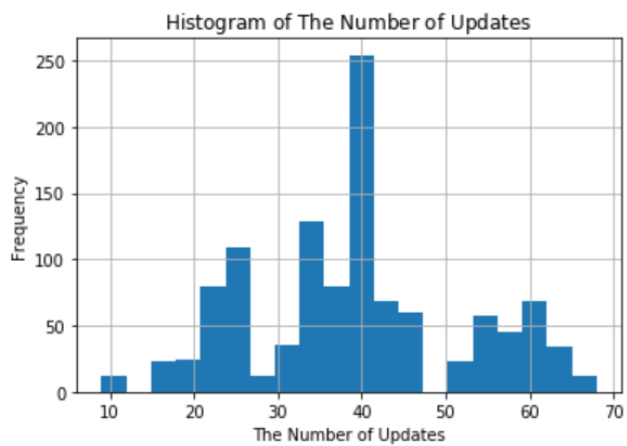
jupyter hw1_6 最后检查: 5 分钟前 (自动保存)

File Edit View Insert Cell Kernel Widgets Help Snippets

代码

```
#Plot histogram to show number of updates
plt.figure()
plt.hist(save_update, bins=20)
plt.ylabel('Frequency')
plt.xlabel('The Number of Updates')
plt.title(r'$\mathrm{Histogram\ of\ The\ Number\ of\ Updates}$')
plt.grid(True)
plt.show()
```

The average number of updates is 39.22113676731794



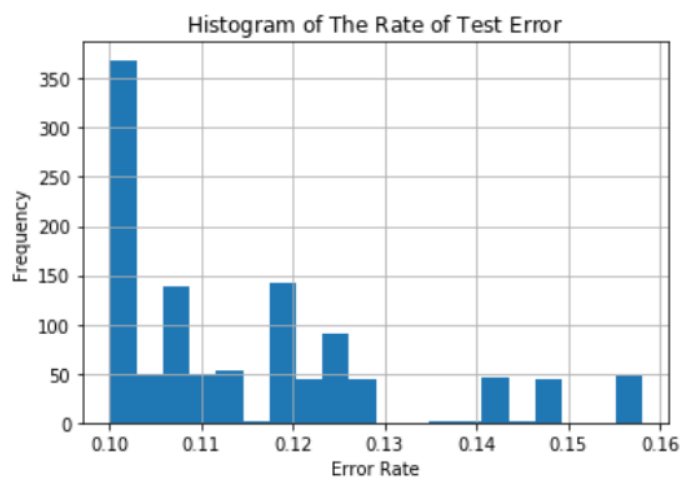
Question7

The average error rate(pocket) on the test set is 0.1149

```
In [49]: avg_error = mean(eval)
print("The average error rate on the test set is ", avg_error)

#Plot histogram to show number of test error
plt.figure()
plt.hist(eval, bins=20)
plt.ylabel('Frequency')
plt.xlabel('Error Rate')
plt.title(r'\mathrm{Histogram\ of\ The\ Rate\ of\ Test\ Error}')
plt.grid(True)
plt.show()
```

The average error rate on the test set is 0.11490230905861457

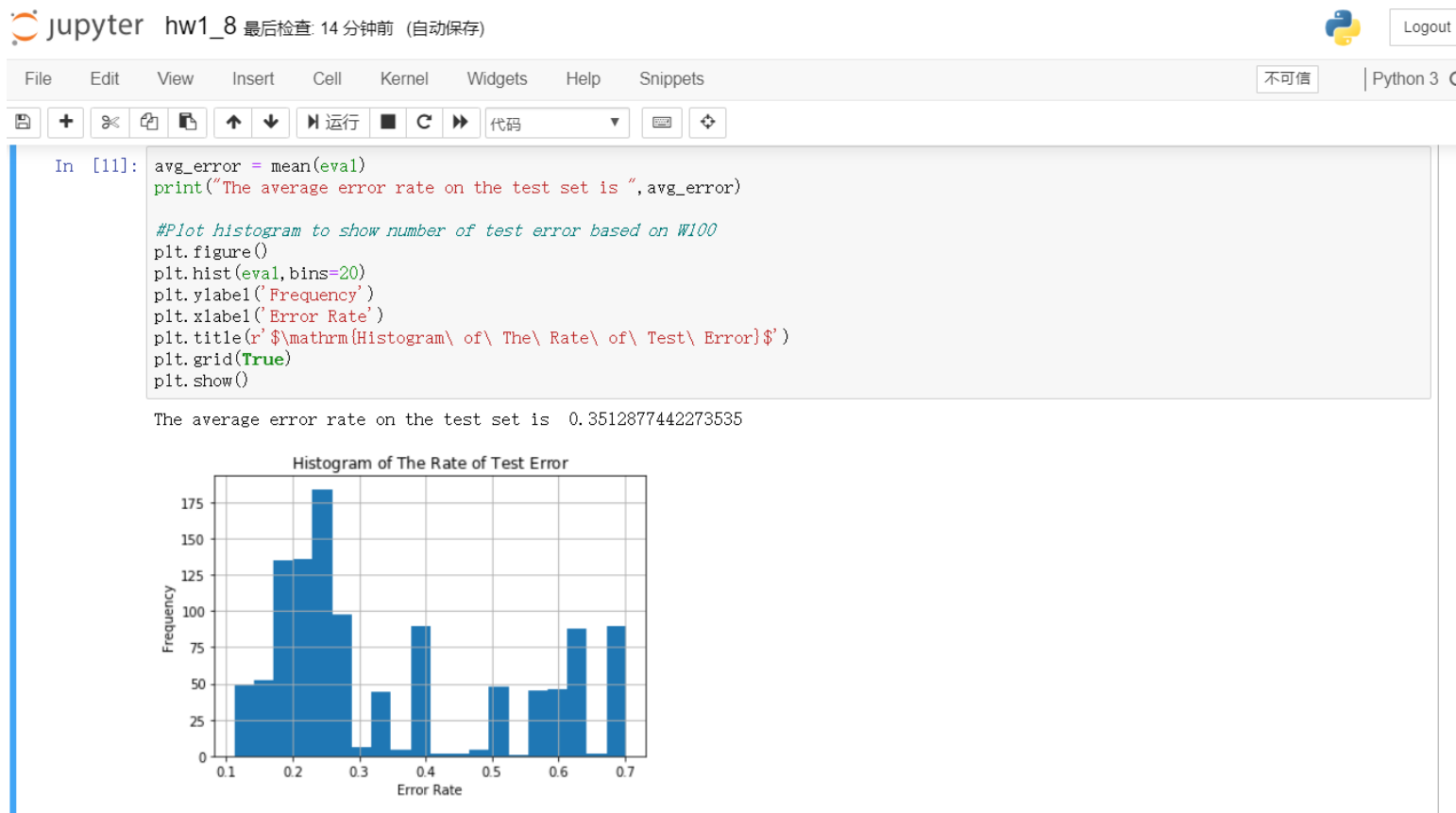


Question8

The average error rate(w_{100}) on the test set is 0.35128

Findings:

實驗 Pocket PLA 演算法，第 7 題 return $w_{\text{hat}}(\text{pocket vector})$ ，第 8 題 return w_{100} 。 $w_{\text{hat}}(\text{pocket vector})$ 保存每次更新使得錯分率更小的權重，相較於直接使用更新到 100 次的權重 w_{100} 。在有限的 iteration 和 update 次數中，使用 $w_{\text{hat}}(\text{pocket})$ 的平均錯誤率=0.115 顯著小於 w_{100} 的平均錯誤率=0.351 可以驗證使用 pocket 演算法能得到最佳的分類線，雖然仍有分錯的點，但已經是最少的了。



實驗Pocket PLA演算法，第7題return $w_{\text{hat}}(\text{pocket vector})$ ，第8題return w_{100} 。 $w_{\text{hat}}(\text{pocket vector})$ 保存每次更新使得錯分率更小的權重，相較於直接使用更新到100次的權重 w_{100} 。在有限的iteration和update次數中，使用 $w_{\text{hat}}(\text{pocket})$ 的平均錯誤率=0.115 顯著小於 w_{100} 的平均錯誤率=0.351 可以驗證使用pocket演算法能得到最佳的分類線，雖然仍有分錯的點，但已經是最少的了。

9. It doesn't work

R08944052 斯晓宇

Refer to page 16 in Lecture 2

$$R^2 = \max_n \|X_n\|^2 \quad \rho = \min_n y_n \frac{w_f^T}{\|w_f\|} X_n$$

The number of mistake corrections $T \leq \frac{R^2}{\rho^2}$

如果 X_n 线性缩小 10 倍, X_n 的模长 $\|X_n\|$ 也会缩小 10 倍.

$$R' = \max_n \left\| \frac{X_n}{10} \right\| \quad \rho' = \min_n y_n \frac{w_f^T}{\|w_f\|} \frac{X_n}{10}$$

$$\frac{R'^2}{\rho'^2} = \frac{(R/10)^2}{(\rho/10)^2} = \frac{R^2}{\rho^2} \geq T$$

\therefore 线性缩放 X_n 不会对收敛速度有所影响.