

5. Machine Learning

5.1 Machine learning

- ① 质量取决于 input data 的质量 ② model choice itself
- ② 进行步骤：① 选择第一个 training data set ② 选择一个 algorithm 来处理 training data set ③ with algorithm 来建立模型
- ④ 使用该模型
- ⑤ Supervised learning
 - ① 使用 labeled data ② training data: (\vec{x}_i, y_i) ③ 常用的：a new \vec{x} , predict the output y ④ Not all x of samples impact y of new data; All of samples impact y of new data.
 - ④ Unsupervised learning
 - ① 使用 data that does not have structure of objective answer
 - ② 常用目的：examine the information and identifying structure within it.

5.2 KNN

- 1. 步骤 ① 找 K 个最近 \vec{x} 的 training points $x_1 \dots x_K$
- ② 若 $x_i \in x$ 中有够多的属于 classifier C , 将 label $x \in C$
- 注：K 小：noise 对结果的影响较大 (overfitting)
K 大：破坏 KNN 的基本原理 (underfitting)
- 若 $K >$ 样本数，任一新 input 都有相同的 label。

注：KNN 的决策边界为非线性的

5.3 Logistic regression model

- 1. Logistic regression (二分类: $y \in \{0, 1\}$)

$$P(y=1|\vec{x}, \theta, b) = \frac{1}{1 + \exp(-(\theta^T \vec{x} + b))}$$

$$P(y=0|\vec{x}, \theta, b) = \frac{\exp(-(\theta^T \vec{x} + b))}{1 + \exp(-(\theta^T \vec{x} + b))} = 1 - P(y=1|\vec{x}, \theta, b)$$
- 注：对于多类问题： $y \in \{0, 1, \dots, k\}$

$P(y=i) = e^{\theta_i^T \vec{x} + b_i}$ $P(y=i|\vec{x}, \theta, b) = \frac{e^{\theta_i^T \vec{x} + b_i}}{1 + \sum_{j \neq i} e^{\theta_j^T \vec{x} + b_j}}$

注：选择回归得出的决策边界为线性的

2. MLE for logistic regression

根据 θ , 得得 likelihood of the labels 最大

$\max_{\theta} l(\theta, b) := \log \prod_i P(y_i|\vec{x}_i, \theta, b) = \sum_i \log P(y_i|\vec{x}_i, \theta, b)$

通常我们只考虑 averaged likelihood 最大化

$\max_{\theta, b} l(\theta, b) = \frac{1}{n} \sum_i \log P(y_i|\vec{x}_i, \theta, b)$

$l^0(\theta, b)$ is concave in (θ, b)

$\log P(y|\vec{x}, \theta, b) = (y-1)(\theta^T \vec{x} + b) - \log(1 + \exp(-\theta^T \vec{x} - b))$

由此要证明 $l(\theta, b) = \frac{1}{n} \sum_i \log P(y_i|\vec{x}_i, \theta, b)$ concave,

仅需证明 $\log(1 + \exp(-\theta^T \vec{x} - b))$ convex.

引理：若 $f(x)$ 是凸的，则 $f(g(x)) + (1-f(g(x)))f(h(x)) \geq f(g(x) + (1-f(g(x)))h(x))$

证明：令 $g(x) = f(x)$, 取 x , 依题意知 $f(x)$ 凸

若 $f(x)$ 是凸的，则取 x_1, x_2 , 依题意知 $f(x)$ 凸

即 $f(x_1) + (1-f(x_1))f(x_2) \geq f(x_1 + (1-f(x_1))x_2)$

即 $f(x_1) + (1-f(x_1))f(x_2) \geq f(x_1 + (1-f(x_1))x_2)$

$= f(x_1) + (1-f(x_1))x_2$

令 $(\theta, b) = (\theta_1, b_1 + t\vec{x}, b_2)$ 为任意单位向量

则 $l(\theta, b) = \log(1 + \exp(-\theta^T \vec{x} - b))$

$= \log(1 + \exp[-(\theta_1 - t)x_1 + b_1 + t\vec{x}^T \vec{x} + b_2])$

$= \log(1 + \exp[-(\theta_1 - t)x_1 + b_1 + t\vec{x}^T \vec{x} + b_2]) \geq 0$

\therefore No closed form solution maximizing $l(\theta, b)$

3. Gradient descent

$x^{(t+1)} = x^{(t)} + \alpha^{(t)} \frac{1}{m} \sum_i (y^{(i)} - 1)x^{(t)} + \frac{\exp(-\theta^T x^{(t)} - b^{(t)})}{1 + \exp(-\theta^T x^{(t)} - b^{(t)})} x^{(t)}$

注：几步 $\alpha^{(t)}$ 可以大一些，靠近极小值时小一些

若 $\alpha^{(t)}$ 为常数，可能会找不到极小值点

* At each iteration, we randomly choose a small batch of data points, and update using the stochastic gradient

$$\begin{aligned} \theta^{(t+1)} &= \theta^{(t)} + \alpha^{(t)} \frac{1}{m} \sum_i (y^{(i)} - 1)x^{(t)} + \frac{\exp(-\theta^T x^{(t)} - b^{(t)})}{1 + \exp(-\theta^T x^{(t)} - b^{(t)})} x^{(t)} \\ \theta^{(t+1)} &= \theta^{(t)} + \alpha^{(t)} \frac{1}{m} \sum_i (y^{(i)} - 1) + \frac{\exp(-\theta^T x^{(t)} - b^{(t)})}{1 + \exp(-\theta^T x^{(t)} - b^{(t)})} \end{aligned}$$

* B: the batch we use in each iteration

5.6 K-means clustering

1. Steps

1° 随机选取 k 个 cluster centers c^0, c^1, \dots, c^k

2° 分类并调整 cluster centers

① cluster assignment

将每个数据点 x^i 归入最近的 cluster centers

$c^j = \arg \min_{c^j} \text{dist}(x^i, c^j)$ (通常为 $\|x^i - c^j\|^2$)

② center adjustment

将 cluster centers 移至该类数据点的平均位置

$c^j = \frac{1}{|S_j|} \sum_{x^i \in S_j} x^i$ (在 S_j 中的点的总数)

3° 重复上一步直到 cluster centers 收敛或无改善

注：1° 不同的初始化可能得到不同的结果

2° algorithm 常会在 some iteration 之后停下来 (the breaking rule)

2. Similarity / Dissimilarity function

① Symmetry (对称性): $d(x, y) = d(y, x)$

② Positive separability: $d(x, y) \geq 0$ 当且仅当 $x = y$

③ Triangular inequality: $d(x, y) \leq d(x, z) + d(z, y)$

3. Distance functions

Minkowski distance: $d(x, y) = (\sum_{i=1}^n |x_i - y_i|^p)^{1/p}$

• Euclidean distance: $p=2$

• Manhattan distance: $p=1$, $d(x, y) = \sum_{i=1}^n |x_i - y_i|$

• inf-distance: $p=\infty$, $d(x, y) = \max_{i=1}^n |x_i - y_i|$

5.7 K-fold validation

1. Validation set approach

对于每一个可选的 model m : ① 用 training data 来 train model

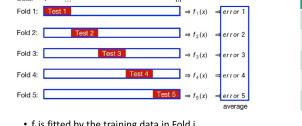
② 用 trained model 来预测 test data, 得到 prediction error (Err_m)

③ 通过 test error 算出最好的 model (min Err_m)

注：此方法通常的模型类型取决于 your separation of the data.
解决方法：重叠多次，每次使用不同的 training and test data.

2. K-fold validation

5-fold cross-validation (blank: training; red: test)



* i is fitted by the training data in Fold i .

* For each fold, use test data to test the prediction error.

* Use the average error as the model's prediction error.

注：K 的选取：通常挑选 K 和 10 , 以 $(\text{每一部分的占比}) = \frac{1}{K}$ 为宜
K 太大会 time-consuming, 若不考虑时间成本, K 选越大越好.

1. MLE: $X_1, \dots, X_n \sim \text{Ber}(\rho) \quad \max \log p$

$\log p(x, p) = p^x(1-p)^{1-x}$

$L(p) = \prod_{i=1}^n p^{x_i}(1-p)^{1-x_i} = p^{\sum x_i}(1-p)^{n-\sum x_i}$

$\ell(p) = \sum_{i=1}^n x_i \cdot \log p + (n-x_i) \cdot \log(1-p) \quad \ell'(p) = \sum_{i=1}^n x_i \cdot \frac{1}{p} + (n-x_i) \cdot \frac{-1}{1-p} = 0$

$\Rightarrow \sum_{i=1}^n x_i \cdot (1-p) + \left(\frac{n}{p}\right) x_i - n = 0 \Rightarrow \hat{p} = \frac{n}{\sum_{i=1}^n x_i}$

2. MLE

A geologist studied the composition of rocks in the shoreline area of Lake Michigan. He randomly selected 100 samples from the area, each containing 10 stones, and recorded the number of limestone stones in each sample. The geologist's data are as follows:

number of limestone in a sample	0	1	2	3	4	5	6	7	8	9	10
number of samples	0	1	6	7	23	26	21	12	3	1	0

number of limestone in a sample	0	1	2	3	4	5	6	7	8	9	10
number of samples	0	1	6	7	23	26	21	12	3	1	0

$\vec{x}_i = \# \text{limestone in sample } i$

Assume that these 100 observations are independent of each other. Find the maximum likelihood estimate (MLE) of the proportion p of limestone in the stones in this area.

$$\begin{aligned} L(p) &= \prod_{i=1}^{100} \left[p^{\vec{x}_i} (1-p)^{1-\vec{x}_i} \right] \\ &= \left[\prod_{i=1}^{100} \frac{1}{\vec{x}_i!} \right] \cdot p^{\sum \vec{x}_i} \cdot (1-p)^{100 - \sum \vec{x}_i} \\ &= \left[\prod_{i=1}^{100} \frac{1}{\vec{x}_i!} \right] \cdot p^{10 \times 100} \cdot (1-p)^{100 - 10 \times 100} \end{aligned}$$

11. Convex set

Some sets of probability distributions. Let x be a real-valued random variable with $\text{prob}(x = a_i) = p_i$, $i = 1, \dots, n$, where $a_1 < a_2 < \dots < a_n$. Of course $p \in \mathbb{R}^n$ lies in the standard probability simplex $P = \{p \mid 1^T p = 1, p \geq 0\}$. Which of the following conditions are convex in p ? (That is, for which of the following conditions is the set of $p \in P$ that satisfy the condition convex?)

- $\alpha \leq \mathbf{E} f(x) \leq \beta$, where $\mathbf{E} f(x)$ is the expected value of $f(x)$, i.e., $\mathbf{E} f(x) = \sum_{i=1}^n p_i f(a_i)$. (The function $f: \mathbb{R} \rightarrow \mathbb{R}$ is given.)
- $\text{prob}(x > \alpha) \leq \beta$.
- $\mathbf{E} |x|^3 \leq \alpha$.
- $\mathbf{E} x^2 \leq \alpha$.
- $\text{var}(x) \leq \alpha$, where $\text{var}(x) = \mathbf{E}(x - \mathbf{E} x)^2$ is the variance of x .
- $\text{var}(x) \geq \alpha$.
- $\text{quartile}(x) \geq \alpha$, where $\text{quartile}(x) = \inf\{\beta \mid \text{prob}(x \leq \beta) \geq 0.25\}$.
- $\text{quartile}(x) \leq \alpha$.

(a) 化为 $\vec{x} \leq \sum_{i=1}^n p_i a_i$ 为 two linear inequalities, 为凸

(b) 化为 $\sum_{i=1}^n p_i \leq \beta$ a linear inequality, 为凸

(c) 化为 $\sum_{i=1}^n p_i (a_i^3 - \alpha |a_i|) \leq 0$ a linear inequality, 为凸

(d) 化为 $\sum_{i=1}^n p_i a_i^2 \leq \alpha$ a linear inequality, 为凸

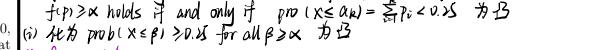
(e) 化为 $\sum_{i=1}^n p_i a_i^2 - (\sum_{i=1}^n p_i a_i) \leq \alpha$ 非凸, 反例:

取 $a_1 = 0, a_2 = 1, a_3 = \frac{1}{3}$, $p_1 = (1, 0, 0)$ 与 $p_2 = (0, 1, 1)$ 满足, 但 $p = (\frac{1}{2}, \frac{1}{2})$ 不满足

(f) $\text{prob}(x \leq \beta) \geq 0.25$ for all $\beta \in \mathbb{R}$ 为凸

(g) $\text{prob}(x \leq \beta) = b^T p + \beta^T A p \leq \alpha$ 为凸

(h) Let us denote $\text{quartile}(x) = f(p)$ to emphasize it is a function of p . The figure illustrates the definition. It shows the cumulative distribution for a distribution p with $f(p) = a_2$.



If $x \leq a_1$ always true. Otherwise, define $k = \max\{i \mid a_i < x\}$.

$f(p) \geq x$ holds if and only if $\text{prob}(x \leq a_k) = \frac{k}{n} p < 0.25$ 为凸

(i) If $\text{prob}(x \leq \beta) \geq 0.25$ for all $\beta \in \mathbb{R}$ 为凸

12. Convex set

Show that S_1 and S_2 are convex sets in $\mathbb{R}^{m \times n}$, then so is their partial sum

$$S = \{(x, y_1 + y_2) \mid x \in \mathbb{R}^m, y_1, y_2 \in \mathbb{R}^n, (x, y_1) \in S_1, (x, y_2) \in S_2\}.$$

Solution. We consider two points $(\bar{x}, \bar{y}_1 + \bar{y}_2)$, $(\tilde{x}, \tilde{y}_1 + \tilde{y}_2) \in S$, i.e., with

$$(\bar{x}, \bar{y}_1) \in S_1, \quad (\bar{x}, \bar{y}_2) \in S_2, \quad (\tilde{x}, \tilde{y}_1) \in S_1, \quad (\tilde{x}, \tilde{y}_2) \in S_2.$$

For $0 \leq \theta \leq 1$,

$$(\theta \bar{x} + (1-\theta)\tilde{x}, \theta \bar{y}_1 + (1-\theta)\tilde{y}_1) \in S_1, \quad (\theta \bar{x} + (1-\theta)\tilde{x}, \theta \bar{y}_2 + (1-\theta)\tilde{y}_2) \in S_2.$$

is in S because, by convexity of S_1 and S_2 ,

$$\begin{aligned} f(\theta \bar{x} + (1-\theta)\tilde{x}) &\leq \frac{n}{\sum_{i=1}^n x_i} \log \bar{x} \geq \frac{1}{\sum_{i=1}^n x_i} \log (\theta \bar{x} + (1-\theta)\tilde{x}) \\ \text{Similarly } f(\theta \bar{y}_1) &\leq \frac{n}{\sum_{i=1}^n y_i} \log \bar{y}_1 \leq \frac{1}{\sum_{i=1}^n y_i} \log (\theta \bar{y}_1 + (1-\theta)\tilde{y}_1); \end{aligned}$$

$$\text{Now } \lambda \frac{f(\theta \bar{x} + (1-\theta)\tilde{x})}{f(\theta \bar{x} + (1-\theta)\tilde{x}) + (1-\lambda)} \frac{f(\theta \bar{y}_1)}{f(\theta \bar{y}_1 + (1-\theta)\tilde{y}_1)} \leq \frac{1}{\sum_{i=1}^n x_i} \frac{\lambda \bar{x}_i + (1-\lambda)\tilde{x}_i}{\sum_{i=1}^n x_i + (1-\lambda)\tilde{x}_i}; = 1$$

As $f(\theta \bar{x} + (1-\theta)\tilde{x}) > 0$,

$$\lambda f(\bar{x}) + (1-\lambda)f(\tilde{x}) \leq f(\theta \bar{x} + (1-\theta)\tilde{x})$$

证 $f(x) = (x - \theta e_i)^2 + (y - \theta e_j)^2$ convex

$$\text{Let } g(\theta) = (x - \theta e_i - x_0)^2 + (y - \theta e_j - y_0)^2, \text{ Then } g'(0) = 2(x - \theta e_i - x_0) + 2(y - \theta e_j - y_0) e_i$$

$$g''(0) = 2 + 2\theta^2 e_i^2 \geq 2 > 0$$

where e_i denotes the i th element of the unit vector e .

13. Convexity

Show that $f(x_1, x_2) = (x_1^p + x_2^p)^{1/p}$ for $x_1, x_2 > 0, p < 1, p \neq 0$ is concave.

$$g(\theta) = f(x_1 + \theta e_1, x_2 + \theta e_2) = [(x_1 + \theta e_1)^p + (x_2 + \theta e_2)^p]^{1/p}. \text{ Then we have}$$

$$g'(\theta) = [(x_1 + \theta e_1)^{p-1} e_1 + (x_2 + \theta e_2)^{p-1} e_2]$$

$$= g(\theta)^{1-p} [(x_1 + \theta e_1)^{p-1} e_1 + (x_2 + \theta e_2)^{p-1} e_2]$$

$$g''(\theta) = (1-p)g(\theta)^{-p} g'(\theta) [(x_1 + \theta e_1)^{p-2} e_1 + (x_2 + \theta e_2)^{p-2} e_2]$$

$$= (1-p)g(\theta)^{1-2p} [(x_1 + \theta e_1)^{p-2} e_1^2 + (x_2 + \theta e_2)^{p-2} e_2^2]$$

$$= (1-p)g(\theta)^{1-2p} [(x_1 + \theta e_1)^{p-2} e_1^2 + (x_2 + \theta e_2)^{p-2} e_2^2]$$

$$= \frac{(1-p)}{g(\theta)} \left[\left(\sum_{i=1}^2 g(\theta)^{1-p} (x_i + \theta e_i)^{p-2} e_i^2 \right)^2 - \sum_{i=1}^2 g(\theta)^{2-p} (x_i + \theta e_i)^{p-2} e_i^2 \right]$$

$$\leq 0,$$

where we apply the Cauchy-Schwarz inequality ($\sum_{i=1}^n a_i b_i)^2 \leq (\sum_{i=1}^n a_i^2)(\sum_{i=1}^n b_i^2)$ with

$$a_i = g(\theta)^{-p/2} (x_i + \theta e_i)^{p/2}, b_i = g(\theta)^{1-p/2} (x_i + \theta e_i)^{p/2-1} e_i.$$

16. Convexity

For $f(x) = ax^3 + bx^2 + cx + d$ with $x \in \mathbb{R}$, find ranges of a, b, c, d to make $f(x)$ quasi-concave.

The first-order derivative of $f(x)$ with respect to x is $f'(x) = 3ax^2 + 2bx + c$.

- Case 1: $f'(x) \geq 0$ or $f'(x) \leq 0$ for all $x \in \mathbb{R}$. In this case, $f(x)$ is always monotone and thus quasi-concave.

- If $a = 0$, then $b = 0$.

- If $a \neq 0$, then $\Delta = 4b^2 - 12ac = 4(b^2 - 3ac) \leq 0$, i.e., $b^2 - 3ac \leq 0$.

- Case 2: $f'(x)$ will change its sign for $x \in \mathbb{R}$.

- If $a = 0$, then only when $b < 0$, $f(x)$ first increases up to x^* and then decreases, and thus $f(x)$ is quasi-concave.

- If $a \neq 0$ and $\Delta > 0$, i.e., $b^2 - 3ac > 0$, it is impossible that $f(x)$ is quasi-concave.

In summary, $f(x)$ is quasi-concave if and only if $a \neq 0, b^2 - 3ac \leq 0$ or $a = 0, b \leq 0$.

17. Convexity

Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex with $\text{dom } f = \mathbb{R}^n$, and bounded above on \mathbb{R}^n . Show that f is constant.

Assume f is not constant. $\exists x, y \in \mathbb{R}^n, f(x) < f(y)$

Define $g(t) = f(x+t(y-x))$

is also convex. $g(0) = f(x) < g(1) = f(y)$

If we assume $t \geq 1$, we will not use the convexity of $g(t)$ to get this's inequality:

$g(1) \leq \frac{t-1}{t}g(0) + \frac{1}{t}g(t)$

$\Rightarrow g(1) \geq t(g(1) - (t-1)g(0)) = g(0) + t(g(1) - g(0))$

As $t \rightarrow \infty$, $g(1) \rightarrow \infty$.

$f(x+t(y-x)) = g(t) \rightarrow \infty$

$f \rightarrow \infty$ as $t \rightarrow \infty$ contradiction of our condition.

We can say: f is constant. \square

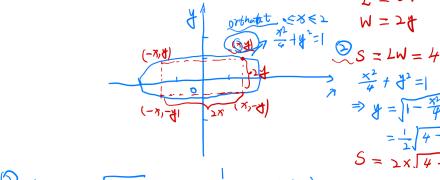
18. Convexity

A rectangle is to be inscribed in the ellipse $\frac{x^2}{4} + y^2 = 1$

$$\frac{x^2}{4} + y^2 = 1$$

What should the dimensions of the rectangle be to maximize its area?

What is the maximum area?



$$\textcircled{1} \quad L = 2x$$

$$\textcircled{2} \quad W = 2y$$

$$\textcircled{3} \quad S = LW = 4xy$$

$$\frac{x^2}{4} + y^2 = 1$$

$$y = \sqrt{1 - \frac{x^2}{4}} = \frac{\sqrt{4 - x^2}}{2}$$

$$S = 2x \cdot \frac{\sqrt{4 - x^2}}{2}$$

$$S(x) = 2\sqrt{4 - x^2} + 2x \cdot \frac{1}{2\sqrt{4 - x^2}} \cdot (-2x)$$

$$= \frac{8 - 4x^2}{\sqrt{4 - x^2}} = 0$$

$$\Rightarrow x = \sqrt{2}$$

Claim: For a continuous $S(x) = 2x\sqrt{4 - x^2}$ with the closed interval $[0, 2]$, then $S(x)$ can only attain its max or min at $x = 0$, $x = 2$, or $x = x^* (S(x^*) = 0)$

$$S(0) = 0, \quad S(2) = 0, \quad S(x^*) = S(\sqrt{2}) = 2\sqrt{2}\sqrt{2} = 4.$$

In conclusion,

$$L = 2x = 2\sqrt{2}, \quad W = 2y = \sqrt{2}$$

$$\text{Max}(S) = 4.$$

19. ML

We apply the gradient descent method to minimize the function $f(x) = (x-1)^2$.

- If the update is stuck in a loop, that is, $f(x^{(t+1)}) = f(x^{(t)})$ yet $|x^{(t+1)} - x^{(t)}| > \epsilon$, what can we say about the learning rate/step size $\alpha^{(t)}$?

- If the function $f(x)$ keeps descending as the update goes on, what can we say about the learning rate/step size $\alpha^{(t)}$?

- The updating scheme is

$$x^{(t+1)} = x^{(t)} - 2\alpha^{(t)}(x^{(t)} - 1).$$

It follows that

$$f(x^{(t+1)}) = f(x^{(t)}),$$

that is

$$(x^{(t)} - 2\alpha^{(t)}(x^{(t)} - 1) - 1)^2 = (x^{(t)} - 1)^2$$

which yields

$$\alpha^{(t)}(\alpha^{(t)} - 1)(x^{(t)} - 1)^2 = 0.$$

Since the update is stuck and does not converge, $x^{(t)} \neq 1$, we have $\alpha^{(t)} = 1$.

- If follows that

$$f(x^{(t+1)}) < f(x^{(t)}),$$

that is

$$(x^{(t)} - 2\alpha^{(t)}(x^{(t)} - 1) - 1)^2 < (x^{(t)} - 1)^2$$

which yields $\alpha^{(t)} < 1$.

