

ECE 449. Assignment 1 . Conceptual Question

Q,

a). supervised learning.

use the book this customer bought before as training data,
with features like book type, language, number of words, price and so on.

b). reinforcement learning.

use the decision made every game as training data, like a sequence of
location where to move, also some feature to judge the decision efficiency
like whether win or lose, number of step used to win and so on.

c) both supervised or unsupervised.

it is a clustering problem essentially, with unsupervised learning. we can
effectively get some different group of movies with training data as some
basic features of movies like length, nation, language, director, title
and so on, and actually we don't know whether any one movie is what
type. as we have no label.

With supervised learning, we need training data above as well as all the type
labels each movie belong to.

d) unsupervised learning.

use basic features like style, rhythm, language and so on as training
data

e) both supervised learning

use data of previous customers and their
maximum debt.

Q₂

the key difference is : linear regression is a regression problem whose output is continuous while logistic regression, although named in regression, is actually a classification problem, which has discrete output.

linear regression example : Stock trend forecast

logistic regression example : spam mail detection.

Q₃:

$$(1,2), (2,1), (3,2) \quad \hat{y} = w_0 + w_1 x = H\hat{w}$$

with $\hat{w} = [w_0 \ w_1]^T$, $y = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$, $H = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 1 \end{bmatrix}$

so $H^T y = H^T H w$, $H^T H$ must be invertible.

$$\Rightarrow (H^T H)^{-1} H^T y = (H^T H)^{-1} (H^T H) w = w$$

$$w = (\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 1 \end{bmatrix}^T \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 1 \end{bmatrix})^{-1} \cdot \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 1 \end{bmatrix}^T \cdot \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 5 \\ 2 \\ 0 \end{bmatrix}$$

$$\text{so } w_0 = \frac{5}{3}, w_1 = 0.$$

Q₄:

$$(-2,1), (-1,0), (1,0), (2,2) \quad y = w_0 + w_1 x_1 + w_2 x_2^2$$

$$y = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 2 \end{bmatrix} \quad w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \quad H = \begin{bmatrix} 1 & -2 & 4 \\ 1 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix}$$

Q₅:

For Gradient descent, assume we fit our model with $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \dots$, $\theta = [\theta_0, \theta_1, \dots]^T$, then we analyze the loss by $J(\theta) = \min_{\theta} \theta^T \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$

start by some random guess of θ , we update θ_j ($j=0,1,2,\dots$) with

$$\theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$
, with α some learning rate defined until θ_j

convergence or through some iteration, one important thing is all θ_j should update simultaneously.

The different is Gradient Ascent algorithm are used to get max value, with the same idea of use gradient, but descent algorithm update θ by subtract but ascent use addition, which cause their different use.

Q6:

large learning rate:

Advantage: we may quickly converge that can shorter operation time

Disadvantage: as it jump large step each time, we may fail to converge or even diverge

Small learning rate:

Advantage: we can always find converge point

Disadvantage: we may get into some local min that are not good enough and take long operation time

Q7:

$$H(X) = - \sum_{i=1}^4 P_i \log_2 P_i$$

$$= -4 \times \frac{1}{4} \log_2 \frac{1}{4}$$

$$= 2.$$

Q8:

define X as different records of scores.

Y as different postgraduate admissions

$$X \text{ margin } (\frac{7}{16}, \frac{1}{4}, \frac{5}{32}, \frac{5}{32})$$

$$Y \text{ margin } (\frac{1}{4}, \frac{1}{4}, \frac{5}{8}, \frac{1}{8}).$$

$$\begin{aligned} H(Y|X) &= \sum_{x \in X, y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \\ &= \frac{3x_1}{8} \log \frac{\frac{7}{16}}{\frac{7}{8}} + \frac{1}{16} \log \frac{\frac{7}{16}}{\frac{1}{4}} + \frac{2x_1}{16} \log \frac{\frac{1}{4}}{\frac{5}{8}} + \frac{1}{8} \log \frac{\frac{1}{4}}{\frac{1}{8}} + \frac{2x_1}{32} \log \frac{\frac{5}{8}}{\frac{5}{16}} + \frac{4x_1}{16} \log \frac{\frac{5}{8}}{\frac{5}{16}} \\ &= \frac{3}{8} \log \frac{7}{2} + \frac{1}{16} \log 7 + \frac{1}{4} + \frac{1}{16} \log 5 + \frac{1}{8} \log \frac{5}{2} \\ &= 1.704 \end{aligned}$$

$$\begin{aligned}
H(x|Y) &= \sum_{x \in X, y \in Y} P(x,y) \log \frac{P(x)}{P(x,y)} = \frac{1}{8} \log \frac{\frac{1}{4}}{\frac{1}{8}} + \frac{1}{16} \log \frac{\frac{1}{4}}{\frac{1}{16}} + \frac{5^2}{32} \cdot \log \frac{\frac{1}{4}}{\frac{1}{5^2}} + \frac{1^2 \cdot 4}{16} \log \frac{\frac{1}{4}}{\frac{1}{16}} + \frac{1^2}{8} \log \frac{\frac{1}{4}}{\frac{1}{8}} \\
&\quad + \frac{1^2}{16} \log \frac{\frac{1}{3}}{\frac{1}{16}} + \frac{1}{8} \log \frac{\frac{1}{3}}{\frac{1}{8}} \\
&\approx \frac{1}{8} + \frac{1}{8} + \frac{3}{16} + \frac{1}{2} + \frac{1}{4} \log 3 + \frac{1}{8} \log 6 + \frac{1}{8} \\
&= 1.657
\end{aligned}$$

$$\begin{aligned}
I(X,Y) &= H(X) - H(X|Y) \\
&= -\frac{7}{16} \log \frac{7}{16} - \frac{1}{4} \log \frac{1}{4} - \frac{5}{32} \log \frac{5}{32} - \frac{5}{32} \log \frac{5}{32} - H(X|Y) \\
&= 0.202
\end{aligned}$$

Q9.

Follow this hint. first prove $h_i = x_i (x^T x)^{-1} x_i^T$

$$\begin{aligned}
x_i (x^T x)^{-1} x_i^T &= (1, x_i) \left[\begin{pmatrix} 1 & \cdots & 1 \\ x_0 & \cdots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_0 \\ \vdots & \vdots \\ 1 & x_0 \end{pmatrix} \right]^{-1} \begin{pmatrix} 1 \\ x_i \end{pmatrix} \\
&= (1, x_i) \left[\begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \right]^{-1} \begin{pmatrix} 1 \\ x_i \end{pmatrix} \\
&= (1, x_i) \left(\begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & N \end{pmatrix} \begin{pmatrix} 1 \\ x_i \end{pmatrix} \right) / (N \cdot \sum x_i - (\sum x_i)^2) \\
&= \frac{\sum x_i^2 - 2x_i \sum x_i + Nx_i^2}{N \sum x_i^2 - (\sum x_i)^2} \\
&= \frac{1}{N} + \frac{N \bar{x}^2 - 2N \bar{x} x_i + Nx_i^2}{N \sum x_i^2 - N \bar{x}^2} \\
&= \frac{1}{N} + \frac{(\bar{x} - x_i)^2}{\sum x_i^2 - N \bar{x}^2} = \frac{1}{N} + \frac{(\bar{x} - x_i)^2}{\sum (x_i^2 - \bar{x}^2)} \\
&= \frac{1}{N} + \frac{(\bar{x} - x_i)^2}{\sum (x_i^2 + \bar{x}^2) - 2x_i \bar{x}} = \frac{1}{N} + \frac{(\bar{x} - x_i)^2}{\sum (x_i - \bar{x})^2}
\end{aligned}$$

Then for $C_V = \frac{1}{N} \sum_{i=1}^N \left(\frac{y_i - \hat{y}_i}{1-h_i} \right)^2$, say x_i, y_i be X, Y without i^{th} row

$x_i^T x_i = (x^T x - x_i x_i^T)$, and with

$$(x_i^T x_i)^{-1} = (x^T x)^{-1} + \frac{(x^T x)^{-1} x_i x_i^T (x^T x)^{-1}}{1-h_i}$$

$$\begin{aligned}
\hat{\beta}_i &= (x_i^T x_i)^{-1} x_i^T y_i \\
&= \left((x^T x)^{-1} + \underbrace{(x^T x)^{-1} x_i x_i^T (x_i^T x_i)^{-1}}_{1-h_i} \right)^{-1} (x^T y - x_i^T y_i) \\
&= \hat{\beta} - \left(\underbrace{(x^T x)^{-1} x_i}_{1-h_i} \right)^{-1} (x_i^T (1-h_i) - x_i^T \hat{\beta} + h_i y_i) \\
&= \hat{\beta} - \underbrace{(x^T x)^{-1} x_i}_{1-h_i} (y_i - \hat{y}_i)
\end{aligned}$$

$$\begin{aligned}
\text{given } (y_i - \hat{y}_i)^2 &= \underbrace{(y_i - x_i^T \hat{\beta})^2}_{1-h_i} \\
&= (y_i - x_i^T (\hat{\beta} - \underbrace{(x^T x)^{-1} x_i (y_i - \hat{y}_i)}_{1-h_i}))^2 \\
&= (y_i - \hat{y}_i + \underbrace{h_i (y_i - \hat{y}_i)}_{1-h_i})^2 \\
&= \left(\frac{y_i - \hat{y}_i}{1-h_i} \right)^2 = \text{MSE}
\end{aligned}$$