

DS5220 - Homework 02.

KAVANA VENKATESH

1. Linear Regression.

$$\text{Given } \hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N (\theta^T \phi(x_i) - y_i)^2 + \lambda \|\theta - a\|_2^2$$

Using the definition of L_2 norm, we can write the above equation as,

$$\begin{aligned} \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N (\theta^T \phi(x_i) - y_i)^2 + \lambda \|\theta - a\|_2^2 \\ = \|\phi \hat{\theta} - y\|_2^2 + \lambda \|\hat{\theta} - a\|_2^2 \end{aligned}$$

$$\text{where } y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad \text{and} \quad \phi = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}$$

In order to find the closed form solution of the above equation, take derivative wrt θ and equate it to 0.

$$2 \phi^T (\phi \hat{\theta} - y) + 2 \lambda (\hat{\theta} - a) = 0.$$

$$\phi^T (\phi \hat{\theta} - y) = -\lambda (\hat{\theta} - a)$$

$$\phi^T \phi \hat{\theta} - \phi^T y = -\lambda \hat{\theta} + \lambda a$$

$$(\phi^T \phi + \lambda I) \hat{\theta} = (\phi^T y + \lambda a)$$

$$\Rightarrow \boxed{\hat{\theta} = (\phi^T \phi + \lambda I)^{-1} (\phi^T y + \lambda a)}$$

2. Given the Robust regression model

$$\min_{\theta} \sum_{i=1}^N \rho(y_i - \theta^T \phi(x_i))$$

2.a) Steps of Batch gradient descent in order to obtain the solution for θ .

$$\text{cost function} = \min_{\theta} \sum_{i=1}^N \begin{cases} \frac{1}{2} (y_i - \theta^T \phi(x_i))^2 \\ \delta |y_i - \theta^T \phi(x_i)| - \frac{1}{2} \delta^2 \end{cases}$$

$$= \min_{\theta} \sum_{i=1}^N \begin{cases} \frac{1}{2} (y_i - \phi^T(x_i) \theta)^2 & \text{if } |y_i - \theta^T \phi(x_i)| \leq \delta \\ \delta |y_i - \phi^T(x_i) \theta| - \frac{\delta^2}{2}, & |y_i - \theta^T \phi(x_i)| > \delta \end{cases}$$

Gradient of Huber loss will be,

$$\frac{\partial \text{cost}}{\partial \theta} = \frac{1}{N} \sum_{i=1}^N \begin{cases} -\phi(x_i) (y_i - \phi^T(x_i) \theta) & \text{if } |y_i - \theta^T \phi(x_i)| \leq \delta \\ -\delta \text{sign}(y_i - \phi^T(x_i) \theta) & \text{if } |y_i - \theta^T \phi(x_i)| > \delta \end{cases}$$

Sign will be +ve where $e \geq 0$.
-ve where $e < 0$.

→ Batch gradient descent

Step 1:- choose a start value θ .

Step 2:- we pick a value for learning rate ρ .

Step 3:- Let's iterate for $t=1, 2, \dots$ and update θ until convergence point.
$$\theta^{t+1} = \theta^t - \rho \nabla_{\theta} J(\theta).$$

By substituting this in the previous eq,

$$\theta^{t+1} = \theta^t - \rho \nabla_{\theta} J(\theta) = \theta^t - \rho \left(-\frac{1}{N} \sum_{i=1}^N \phi(x_i) (y_i - \phi^T(x_i) \theta) \right);$$

$$|y_i - \phi^T(x_i) \theta| \leq \delta.$$

$$\theta^{t+1} = \theta^t - \rho \left(-\frac{1}{N} \sum_{i=1}^N (+v_i) (\phi(x_i)) \right); y_i - \phi^T(x_i) \theta > \delta.$$

$$\theta^{t+1} = \theta^t - \rho \left(-\frac{1}{N} \sum_{i=1}^N (-v_i) (\phi(x_i)) \right); y_i - \phi^T(x_i) \theta < \delta.$$

Step 4:- Repeat Step 3 until the stopping criteria is met, which is change in norm of the \leq threshold value.
gradient

a.3) Stochastic gradient descent

Step 1:- choose the value of θ such that $\theta \in \mathbb{R}^{d+1}$.

Step 2:- pick a value for learning rate ρ .

Step 3: update θ by iterating for $i=1, 2$ until convergence point.

$$\theta^{t+1} = \theta^t - \eta \nabla_{\theta} J_g(x_i, \theta)$$

For random minibatch of size 1, the gradient descent changes as follows.

$$\nabla_{\theta} J_g(c) = \frac{d}{d\theta} J_g(c) = \begin{cases} -\frac{1}{2} \cdot 2 \phi(x_i) (y_i - \phi^T(x_i) \theta) ; & |y_i - \theta^T \phi(x_i)| \leq \delta \\ -\delta \cdot \text{sign}(y_i - \phi^T(x_i) \theta) \phi(x_i) ; & |y_i - \theta^T \phi(x_i)| > \delta \end{cases}$$

Step 4: we repeat Step 3 until convergence is met.

Change in norm of the gradient \leq threshold Value

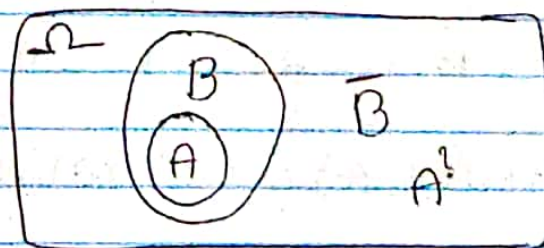
3. Given $P(A|B) + P(A|B^c) = 1 \quad \forall A, B \subseteq \Omega$
and $0 < P(A) < 1$

3.1) FALSE.

$P(A)$ can be written as $(A|\Omega) = P_\Omega(A)$.
This implies that by writing $P(A|B)$,
we are looking for probability of
an event A out of all the outcomes
of B .

$$\Rightarrow P(A|B) = P_B(A).$$

Once we look at P_B , we can not
move from P_B to $P_{\bar{B}}$. i.e., probabilities
in one sample space P_B can
not tell us anything about the
probabilities in another sample space
 $P_{\bar{B}}$.



\Rightarrow To disprove the above statement (given),
consider the popular example of
tossing a fair 6-sided die.

$$P(\text{Even}) = \frac{1}{2}, \quad P(\text{odd}) = \frac{1}{2}.$$

$$P(X=2 | \text{Even}) = \frac{1}{3} \quad \text{and}$$

$$P(X=2 | \overline{\text{Even}}) = P(X=2 | \text{odd})$$

Since it's a 2 event experiment and there are only 2 possibilities \Rightarrow
complement(even) = odd.

$$\text{Now, } P(X=2 | \overline{\text{Even}}) = P(X=2 | \text{odd}) = 0$$

$$P(X=2 | \text{Even}) + P(X=2 | \overline{\text{Even}}) = \frac{1}{3} + 0 \\ = \frac{1}{3} \neq 1$$

$$\text{Hence, } \boxed{P(A|B) + P(A|\overline{B}) \neq 1}$$

$$3.2) \quad P(B^c \cap (A \cup B)) + P(A^c \cup B) = 1$$

TRUE

Proof: $P(B^c \cap (A \cup B)) + P(A^c \cup B)$

$$\Rightarrow P((B^c \cap A) \cup (B^c \cap B)) + P(A^c \cup B) \quad \text{Distributive Law}$$

$$\Rightarrow P(B^c \cap A) + P(B^c \cap B) + P(A^c \cup B)$$

$$\text{But } P(B^c \cap B) = 0$$

$$\therefore P(B^c \cap A) + 0 + P(B \cup A^c)$$

$$\Rightarrow P(B^c \cap A) + P(B \cap A)^c$$

$$= \underline{\underline{1}} \quad (\because P(A) + P(A^c) = 1)$$

3.3) Given $P(A_1, A_2, \dots, A_n) = P(A_1)P(A_2|A_1) \dots$
 $\times P(A_n|A_1, \dots, A_{n-1})$.

By definition of the conditional probability,

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Consider the RHS of the equation,

$$P(A_1)P(A_2|A_1) \cdot P(A_3|A_1, A_2) \dots P(A_n|A_1, A_2, \dots, A_{n-1})$$

Now apply the definition of the conditional probability to the above eq.

$$P(A_1) \cdot \frac{P(A_2, A_1)}{P(A_1)} \cdot \frac{P(A_3, A_2, A_1)}{P(A_1, A_2)} \dots \frac{P(A_1, A_2, \dots, A_n)}{P(A_1, A_2, \dots, A_{n-1})}$$

We can notice that every term in the numerator and denominator gets cancelled except the last term.

$P(A_1, A_2, \dots, A_n)$, which is nothing but the RHS.

3-4) Given X and Y are discrete independent random variables, then
 $E[XY] = E[X]E[Y]$

For discrete random variables X and Y ,

$$E(XY) = \sum_i \sum_j x_i y_j p_{xy}(x_i, y_j)$$

$$= \sum_i \sum_j x_i y_j p_x(x_i) p_y(y_j)$$

Since for independent events,

$$p_{xy}(x_i, y_j) = p_x(x_i) p_y(y_j)$$

$$= \left(\sum_i x_i p_x(x_i) \right) \left(\sum_j y_j p_y(y_j) \right)$$

$$\boxed{E(XY) = E(X)E(Y)}$$

4. $p_\delta(X_i = x_i) = e^{-(\delta^2 + \delta x_i)}$

As the variables x_1, x_2, \dots, x_N are i.i.d,

a) $\mathcal{L}(\delta) = \prod_{i=1}^N p_\delta(X_i = x_i) = \prod_{i=1}^N e^{-(\delta^2 + \delta x_i)}$

is the likelihood function.

$$L(\delta) = e^{-(\delta^2 + \delta x_1)} \cdot e^{-(\delta^2 + \delta x_2)} \dots e^{-(\delta^2 + \delta x_n)}$$

$$= e^{-((\delta^2 + \delta^2 + \dots + \delta^2) + \delta(x_1 + x_2 + \dots + x_n))}$$

$$L(\delta) = e^{-(N\delta^2 + \delta \sum_{i=1}^N x_i)}$$

b). To find the log likelihood function, let's take \ln .

$$\ln(L(\delta)) = \ln(e^{-(N\delta^2 + \delta \sum_{i=1}^N x_i)})$$

$$= -N\delta^2 - \delta \sum_{i=1}^N x_i$$

Taking derivatives wrt δ ,

$$\frac{\partial}{\partial \delta} \ln(L(\delta)) = \frac{\partial}{\partial \delta} (-N\delta^2 - \delta \sum_{i=1}^N x_i)$$

$$= -2N\delta - \sum_{i=1}^N x_i$$

To find the maximum likelihood, let's equate the above eq. to 0.

$$\frac{\partial}{\partial \delta} \ln(L(\delta)) = -2N\delta - \sum_{i=1}^N x_i = 0$$

$$\Rightarrow -2N\delta = + \sum_{i=1}^N x_i$$

$$\delta^* = \frac{- \sum_{i=1}^N x_i}{2N}$$

5) Given $p(y=1|x) = \sigma(\omega^T x)$

5.1) $\omega^T x_n < 0.3$

In order to get the class of the function, let us use the Sigmoid function.

$$\sigma(\omega^T x_n) = \frac{1}{1 + e^{-\omega^T x_n}}$$

Let us now substitute the boundary of $\omega^T x$ i.e., 0.3

$$\sigma(0.3) = \frac{1}{1 + e^{-0.3}} = \frac{1}{1 + 0.74} \approx 0.57$$

WKT, for a Sigmoid function,

$$y = \begin{cases} 1 & ; \quad p(y=1|x) > 0.5 \\ 0 & ; \quad \text{Otherwise} \end{cases}$$

For $\sigma(z)$ to be equal to exactly $\frac{1}{2}$, $\omega^T x = 0$.

$$\Rightarrow \begin{cases} x_n \in \text{class 1} & \text{if } 0 \leq \omega^T x < 0.3 \\ x_n \in \text{class 0} & \text{if } \omega^T x < 0 \end{cases}$$

5.2) Given $\frac{1}{1+e^{\omega^T x}} = 0.7$

Since $\omega^T x$ is +ve, the function is written for class 0.

$$\therefore \frac{1}{1+e^{\omega^T x}} = 0.7$$

$$1+e^{\omega^T x} = \frac{1}{0.7} = \frac{1}{0.7}$$

$$e^{\omega^T x} = \frac{1}{0.7} - 1$$

$$\Rightarrow e^{\omega^T x} = \frac{0.3}{0.7}$$

$$\Rightarrow e^{-\omega^T x} = \frac{7}{3}$$

$$\therefore \frac{1}{1+e^{-\omega^T x}} = \frac{1}{1+\frac{7}{3}} = \frac{3}{10} = 0.3$$

Since $\sigma(z) = \frac{1}{1+e^{-\omega^T x}} = 0.3 < 0.5$,

$$\Rightarrow x_n \in \underline{\text{class 0}}$$

So, the model belongs to class 0 with 70% probability.

6) MLE estimates θ_j^y $j=0,1$ as well

6.1) as $\theta_{\bar{x}_j}^{x_l} | y_j$, $j=0,1$ and $\bar{x}_j = 0,1$ and
for $l=1,2$.

a) θ_j^y for $j=0,1$

$$P(y=0) = 3/7, \quad P(y=1) = 4/7$$

$$b) \quad P(x_1=0 | y=0) = 2/3 \quad P(x_1=0 | y=1) = 2/4$$

$$P(x_2=0 | y=0) = 2/3 \quad P(x_2=0 | y=1) = 2/4$$

$$P(x_1=1 | y=0) = 1/3 \quad P(x_1=1 | y=1) = 2/4$$

$$P(x_2=1 | y=0) = 1/3 \quad P(x_2=1 | y=1) = 2/4$$

$$6.2) \quad P(y=0 | x_1=0, x_2=1)$$

$$= P(y=0) \cdot P(x_1=0 | y=0) \cdot P(x_2=1 | y=0)$$

$$= 2/3 \cdot 1/3 \cdot 3/7 = \underline{\underline{2/21}}$$

6.3) Solution of $P(y=0 | x_1=0, x_2=1)$ without Naive Bayes is,

$$P(x_1=0, x_2=1 | y=0) \cdot P(y=0)$$

$$= \frac{1}{3} \cdot \frac{3}{7}$$

$$= \underline{\underline{\frac{1}{7}}}$$

Also, without the Naive Bayes assumption, we can directly look at the given table where in $x_1=0$ and $x_2=1$ co-exist and write the $P(y=0 | x_1=0, x_2=1)$ as,

$$P(y=0 | x_1=0, x_2=1) = \underline{\underline{\frac{1}{7}}}$$