# Assignment 4

Foundations of Machine Learning

IIT-Hyderabad

Sept-Dec 2021

Submitted by:

Ankita Jain

BM21MTECH14001

## 1. Non-Uniform Weights in Linear Regression

**1. Given:**

A dataset : $(x_n, t_n) ; n = 0, 1, \ldots N$

Non negative weighing factor $g_n > 0$ for each data point

Resultant error function:

$$E_D(w) = \frac{1}{2} \sum_{n=1}^{N} g_n \left(t_n - w^T \phi(x_n)\right)^2 \quad -(1)$$

$\phi(\cdot) \rightarrow$ any representation of data

**(a)** To find $w^*$ that minimizes $E_D(w)$

**Soln** Let's first evaluate the derivative of the above error function $-(1)$ w.r.t. $w$

$$\frac{\partial E_D(w)}{\partial w} = \frac{1}{2} \sum_{n=1}^{N} g_n \left(t_n - w^T \phi(x_n)\right) \cdot 2 \cdot \left(-\bar{\phi}(x_n)\right)$$

$\boxed{\phi w^T \phi(x) = \phi^T(x) w}$ hence

Equating the derivative to 0,

$$\sum_{n=1}^{N} g_n \left(t_n - w^T \phi(x_n)\right) \phi^T(x_n) = 0$$

Substituting

$$\sum_{n=1}^{N} g_n t_n \phi^T(x_n) - \sum_{n=1}^{N} g_n w^T \phi(x_n) \phi^T(x_n) = 0$$

$\Phi = \begin{bmatrix} \phi_{11} & \phi_{12} \cdots \\ \phi_{21} \\ \vdots \end{bmatrix}$

$t = [t_1 \cdots t_n]$

Substituting $\sqrt{g_n}\, \phi(x_n) = \phi'(x_n)$ (feature mapping) $\sqrt{g_n}\, t_n = t'_n$

$$\sum_{n=1}^{N} t'_n \phi'(x_n) - w^T \sum_{n=1}^{N} \phi'(x_n) \phi'(x_n)^T = 0$$

(Putting the above in vectorized form)

$$\phi'(x_n) t - w^T \phi'(x_n)^T \phi'(x_n) = 0$$

$$\therefore w = (\phi^T \phi)^{-1} \phi^T t$$

(b) To find two alternatives of the above weighted sum of squares error function in terms of

(i) Data dependent noise variance

Data dependent noise variance is indicated by $\beta$ in the following equation.

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^{N} \{t_n - w_{ML}^T \phi(x_n)\}^2 \quad \text{—②}$$

also known as maximizing the log likelihood function w.r.t. noise parameter $\beta$

Now, if we substitute $r_{ML} = \ell n \beta^{-1}$; our eqn ② becomes:

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_n \ell n \{t_n - w\phi(x_n)\}$$

Hence we, see this is one of the interpretation in terms of data independent noise variance

(ii) ... Data Points

(b) Replicated Data Points

the variable $g_n$ which has been termed as a non negative weight for an individual data point $(x_n, t_n)$

can also be the number of times a single data point has been replicated

i.e the effective number of times the data point has been repeated it sam can be treated the same as a weight attached to a specific data point

# 2. Bayes Optimal Classifier:

2. Given 5 hypothesis that could guide a robot to move either Forward (F)

Left (L)

or Right (R)

| $P(h_i|D)$ | $P(F|h_i)$ | $P(L|h_i)$ | $P(R|h_i)$ |
|---|---|---|---|
| 0.4 | 1 | 0 | 0 |
| 0.2 | 0 | 1 | 0 |
| 0.1 | 0 | 0 | 1 |
| 0.1 | 0 | 1 | 0 |
| 0.2 | 0 | 1 | 0 |

## MAP Estimate:

MAP hypothesis is defined as follows:

From the table above

$$h_{MAP} \equiv \underset{h \in H}{argmax}\, P(h|D)$$

where
$D \to data$
$h \to hypothesis$

$$= \underset{h \in H}{argmax}\, (0.4, 0.2, 0.1, 0.1, 0.2)$$

$= h \to h_1 \to$ Robot should move Forward

$$= \underset{h \in H}{argmax}\, \frac{P(D|h)P(h)}{P(D)}$$

$$= \underset{h \in H}{argmax}\, P(D|h)P(h)$$

# Bayes Optimal Estimate

Bayes Optimal Classification is defined as:

$$h_{BO} \equiv \underset{h \in H}{\text{argmax}} \sum_{h \in H} P(V_j | h_i) P(h_i | D)$$

where $V \rightarrow$ set of possible classification of the examples

From the table provided provided:

$$\sum_{h \in H} P(F | h_i) P(h_i | D) = \overset{h_1}{(1 * 0.4)} = 0.4$$

$$\sum_{h \in H} P(L | h_i) P(h_i | D) = \overset{h_2}{\underset{+}{1 * 0.2}} = 0.5$$

$$\overset{h_4}{(1 * 0.1)}$$
$$+$$
$$\overset{h_5}{(1 * 0.2)}$$

$$\sum_{h \in H} P(R | h_i) P(h_i | D) = \overset{h_3}{(1 * 0.1)} = 0.1$$

$$\therefore h_{BO} = \text{argmax} (0.4, 0.5, 0.1)$$

Bayes Optimal recommends the robot to move Left

As we see both recommendations are not same

No other classification method using the same hypothesis space & prior knowledge can outperform Bayes Optimal Classifier on the average

# 3. VC-Dimension:

**3.** Given: 1D data $x$ where Hypothesis

$H$ is parametrized by $\{p, q\}$ &

$$h_{p,q} = \begin{cases} 1 & ; \text{if } p < x \leq q \\ 0 & ; \text{otherwise} \end{cases}$$

To find VC dimension of $H$

**Soln** Let's assume VC = 2
then for any 2 values, $H$ must
be able to shatter this subset of $x$

$x = \{x_1, x_2\}$ where $p < x_1, x_2 \leq q$

possible dichotomies

|  | | Label | | |
|-----|---|---|---|---|
| $x_1$ | 0 | 0 | 1 | 1 |
| $x_2$ | 0 | 1 | 0 | 1 |

$\therefore$ In this all $x_1, x_2 < q$ &

$x_1, x_2 > p$

the $H$ can shatter one subset
of $x$ Hence VC dimension is
atleast 2

Now, let's assume VC=3

then H must be able to shatter any one
subset of x $(x_1, x_2, x_3)$ toprove the
above.

Labels

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $x_1$ | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| $x_2$ | 0 | | 0 | 1 | 1 | 0 | 0 | 1 | 1 |
| $x_3$ | 0 | | 1 | 0 | 1 | 0 | 1 | 0 | 1 |

∴ the condition for hypothesis
to be true, $p \leq x < q$

i.e if we assume, $x_1 < x_2 < x_3$
           any
then $\{x_1, x_3\}$ will always
contain $x_2$, hence not all
the dichotomies can be realized
as marked above

Hence the VC dimension for the
given hypothesis is 2

# 4. Regularizer:

**4. Given:** $D$ dimensional data $\{x_1, \ldots, x_D\}$

Linear model : $y(x, w) = w_0 + \sum_{k=1}^{D} w_k x_k$

$N$ such data samples

MSE : $E(w) = \frac{1}{2} \sum_{i=1}^{N} \left( y(x_i, w) - t_i \right)^2$

Suppose Gaussian Noise $N(0, \sigma^2)$ is added to each input variable $x_k$

**To find :** Relation b/w minimizing MSE averaged over noisy data & minimizing standard MSE averaged over noise free data with a $l_2$ regularization term, with bias $w_0$ omitted

**Solution** Let the model with noisy inputs be defined as:

$$\hat{y}'(x, w) = w_0 + \sum_{k=1}^{D} w_k (x_{nk} + \varepsilon_{nk})$$

where $\varepsilon_{nk} \to$ noise added to input except bias
$$\to N(0, \sigma^2)$$

Reducing the above eqn

$$\hat{y}(x, w) = w_0 + \sum_{u=1}^{D} w_k x_{nu} + \sum_{u=1}^{D} \varepsilon_{nu}$$

$$\hat{y}(x, w) = y(x, w) + \sum_{u=1}^{D} \varepsilon_{nu}$$

Substituting this in MSE for noisy input:

$$\hat{E} = \frac{1}{2} \sum_{n=1}^{N} \left( \hat{y}_n - t_n \right)^2$$

$$= \frac{1}{2} \sum_{n=1}^{N} \left\{ (\hat{y}_n)^2 - 2 y_n t_n + t_n^2 \right\}$$

$$\hat{E} = \frac{1}{2} \sum_{n=1}^{N} \left\{ \left( y_n + \sum_{k=1}^{D} w_k e_{nk} \right)^2 - 2 \left( y_n + \sum_{k=1}^{D} w_k e_{nk} \right) t_n + t_n^2 \right\}$$

$$\hat{E} = \frac{1}{2} \sum_{n=1}^{N} \left\{ y_n^2 + 2 y_n \sum_{k=1}^{D} w_k e_{nk} + \left( \sum_{k=1}^{D} e_k w_k e_{nk} \right)^2 - 2 y_n t_n - 2 t_n \sum_{k=1}^{D} w_k e_{nk} + t_n^2 \right\}$$

Finding Averaged MSE

$$E[\hat{E}] = E\left\{ \frac{1}{2} \sum_{n=1}^{N} \left( y_n^2 - 2 y_n t_n + t_n^2 \right) \right\} \rightarrow 0$$

$$+ E\left\{ \sum_{n=1}^{N} \left( 2 y_n \sum_{k=1}^{D} w_k e_k \right) \right\} \quad \therefore e_k = 0 \text{ mean}$$

$$- E\left\{ \frac{1}{2} \sum_{k=1}^{N} 2 t_n \sum_{k=1}^{D} w_k e_{nk} \right\} \rightarrow 0 \quad \therefore e_k = 0 \text{ mean}$$

$$+ E\left\{ \frac{1}{2} \sum_{k=1}^{D} \left( \sum_{k=1}^{D} w_k e_k \right)^2 \right\}$$

$$= E\{E\} + \frac{1}{2} \sum_{n=1}^{N} E\left\{ \left( \sum_{k=1}^{D} w_k e_{wk} \right)^2 \right\}$$

$$= E\{E\} + \frac{1}{2} \sum_{k=1}^{N} \sum_{k=1}^{D} w_k^2 \underbrace{E(e_{wk}^2)}_{\sigma^2}$$

$$= E\{E\} + \frac{1}{2} \sum_{k=1}^{D} \sigma^2 \sum_{k=1}^{D} w_k^2$$

$\therefore$ the rel b/w noisy input MSE & std MSE is:

$$\boxed{E[\hat{E}] = \frac{1}{2} \sum_{n=1}^{N} \left( y(x_n, w) - t_n \right)^2 + \frac{\sigma^2}{2} \sum_{k=1}^{D} w_k^2}$$

where $\hat{E}$ is <u>MSE due to noisy input</u>

1st term is due to std MSE
2nd term Regularizes Input

# 5. Logistic Regression:

```
Cost after 0 epoch is:  0.6793333930895056
Solution to 5b(i)
The updated value of w,b at the end of the epoch is:
 [1.45104757 0.49269336] -1.0073819919017042


 Solution to 5b(ii)
The logistic model P(y=1|x1,x2) is:
 0.577540312892111
Corresponding cross entropy function =
 0.6770734435545663
```
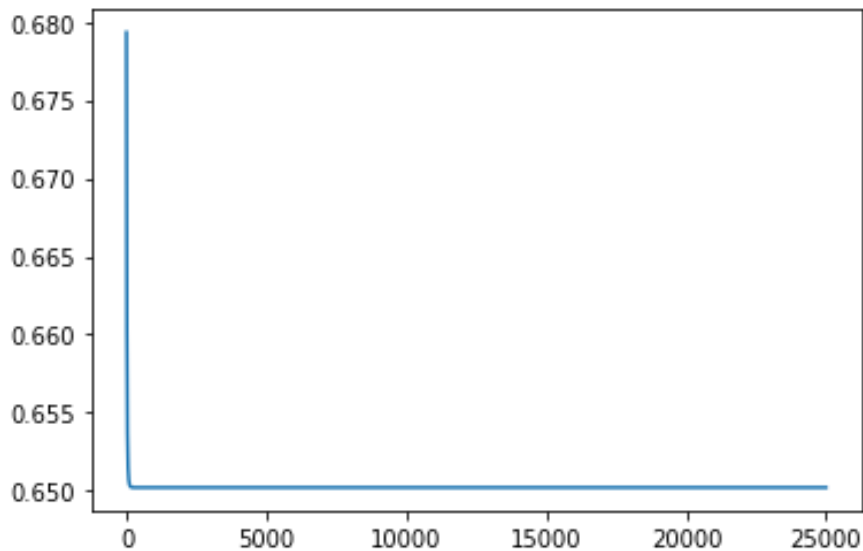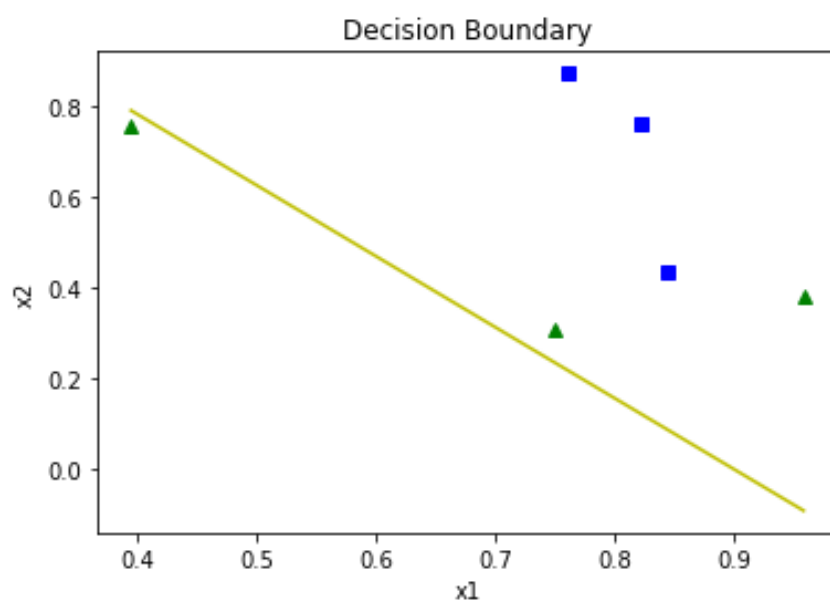


```
The solution to 5b(iii)

Accuracy:  66.66666666666667
Precision:  60.0
Recall:  100.0
```

# 6. Kaggle - Taxi Fare Prediction:

The following two models performed best.

- o An ensemble of kNN and lightGBM
- o LightGBM model

| Name | Submitted | Wait time | Execution time | Score |
|------|-----------|-----------|----------------|-------|
| sub1.csv | just now | 1 seconds | 0 seconds | 2.97710 |

Complete

Jump to your position on the leaderboard ▾

Your most recent submission

| Name | Submitted | Wait time | Execution time | Score |
|------|-----------|-----------|----------------|-------|
| sub2.csv | just now | 1 seconds | 0 seconds | 2.98912 |

Complete

Jump to your position on the leaderboard ▾

A brief on the models selected and why they performed better.

1. LightGBM

   Faster training speed and higher efficiency: Light GBM use histogram based algorithm i.e binning information  Lower memory usage: Replaces continuous values to discrete bins which result in lower memory.

2. Ensemble

   There are two main reasons to use an ensemble over a single model, and they are related; they are: Performance: An ensemble can make better predictions and achieve better performance than any single contributing model. Robustness: An ensemble reduces the spread or dispersion of the predictions and model performance