

Problem (1)

Amin Asadi

- a) In the first one the train-test partition is 50% train & 50% test. In the second one the train portion is 95% and the test portion is 5%. So because in the second one the train set is larger than first one, the second model is more generalized and its parameters estimation have less variance. However since the test set is smaller in the second one, the performance results variance is larger in the second one.

The reason for the 90% accuracy being more than the 80% accuracy might be the more generalized model in the second one, but the accuracy of the second model may be very different on another test set (because as we noted above, it has larger variance) and hence the 90% accuracy is not reliable. Also increasing the train set size from 200 to 380 has not effected to much in the performance in the second model. It could have been better to increase train set size for example by 100 and instead have a bit larger train set in the second model.

- b) Statistical Inference = Probability⁻¹

+ In probability we are concerned about the properties of data from a specified probability distribution

For example if we know that Random Variables X_1, \dots, X_5 have a standard Gaussian distribution and they are iid, what is the probability of X being more than 0.5 ($P(X > 0.5)$)

+ In Statistical Inference we are concerned about finding the properties of the distribution of a specified set of data.

For example suppose that we know $X_1 = 1$ and $X_2 = 2, \dots$ and $X_5 = 3$. We are concerned about finding parameter (mean) θ of the distribution $N(\theta, 1)$ such the $X_1, \dots, X_5 \stackrel{iid}{\sim} N(\theta, 1)$

We use Inference mainly in Hypothesis testing, Estimation and Confidence Intervals.

Problem 2:

- a) Suppose that we have a function F which is dependent on one or more variables. The goal of gradient descent is to find the local minimum of F . This algorithm is based on the fact that, at any point x that F is differentiable on a neighbourhood of x , the fastest decrease in F , occurs when we go from x proportional to negative of gradient of F at point x . We can find the local minimum of F by iteratively taking such steps.

$$\text{Here } F = J(\theta) = J(w; b) = \frac{1}{2} \sum_{i=1}^q (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\left[\frac{\partial J}{\partial b} \right] = \frac{1}{2} \sum_{i=1}^q 2(h_{\theta}(x^{(i)}) - y^{(i)}) \left(\frac{\partial (w^T x^{(i)} + b)}{\partial b} \right) (1 - \tanh^2(w^T x^{(i)} + b))$$

$$= \sum_{i=1}^q (h_{\theta}(x^{(i)}) - y^{(i)}) (1 - \tanh^2(w^T x^{(i)} + b))$$

$$\left[\frac{\partial J}{\partial w_j} \right] = \frac{1}{2} \sum_{i=1}^q 2(h_{\theta}(x^{(i)}) - y^{(i)}) \left(\frac{\partial h_{\theta}(x^{(i)})}{\partial w_j} \right) = \sum_{i=1}^q (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} (1 - \tanh^2(w^T x^{(i)} + b))$$

$$w_j = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad x_j^{(i)}$$

- b) If learning Rate $= \alpha$:

$$\text{update rules: } b := b - \alpha \frac{\partial J}{\partial b}, \quad w_j := w_j - \alpha \frac{\partial J}{\partial w_j}$$

If the learning rate is small then we will take small steps towards the minimum but we are sure that we keep getting closer and closer to that point (It will converge to local minimum). On the other hand, if the learning rate is too high, we will take bigger steps and thus we will get to the minimum in less steps but there is high chance that we overshoot the local minimum and repeatedly fall into other side of the point and hence do not converge. So the learning rate should be tuned as a hyper parameter.

further explanation of the derivative in part a:

$$\begin{aligned}
 \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{1}{2} \sum_{i=1}^q 2(h_{\theta}(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_j} (h_{\theta}(x^{(i)}) - y^{(i)}) \\
 &= \frac{1}{2} \sum_{i=1}^q 2(h_{\theta}(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_j} (\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_m x_m^{(i)} + b - y^{(i)}) \\
 &= \sum_{i=1}^q (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} (1 - \tanh^2(\omega^T x + b)) \times \frac{d}{dx}(\tanh(x)) \Big|_{\omega^T x + b}
 \end{aligned}$$

Problem 3:

We prove that there is a unique polynomial of degree k or less that passes through these $k+1$ points:

let $p(x) = a_k x^k + a_{k-1} x^{k-1} + \dots + a_1 x + a_0$ be a polynomial
($a_i \geq 0$)

If we substitute the $k+1$ points in $p(x)$ we will have a system of $k+1$ linear equations:

$k+1$ points: $(x_0, y_0), \dots, (x_k, y_k)$

$$a_k x_0^k + a_{k-1} x_0^{k-1} + \dots + a_0 = y_0$$

\vdots

$$a_k x_k^k + a_{k-1} x_k^{k-1} + \dots + a_0 = y_k$$

$$\underbrace{\begin{bmatrix} x_0^k & x_0^{k-1} & \dots & x_0 & 1 \\ x_1^k & x_1^{k-1} & \dots & x_1 & 1 \\ \vdots & \vdots & & \vdots & \vdots \\ x_k^k & x_k^{k-1} & & x_k & 1 \end{bmatrix}}_X \underbrace{\begin{bmatrix} a_k \\ a_{k-1} \\ \vdots \\ a_0 \end{bmatrix}}_A = \underbrace{\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_k \end{bmatrix}}_Y$$

X is a Vandermonde Matrix with $\det(X) = \prod_{1 \leq i < j \leq k+1} (x_j - x_i)$

because $k+1$ points have distinct $x \rightarrow \det(X) \neq 0$

$\rightarrow A = X^{-1}Y \rightarrow$ So $P(x)$ is unique.

Obviously if p had a degree more than k , the number of variables would be more than number of equations and hence ~~it~~^{the system} would have more than one solutions \rightarrow more than one polynomials

Problem 4:

suppose

The X_i are fixed

$$\beta_1 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x}) Y_i - \bar{Y} \sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$Y = b_0 + b_1 x$
 $\text{Var}(Y_i) = \sigma^2$

$$= \frac{\sum (x_i - \bar{x}) Y_i - \bar{Y} \sum x_i + \bar{Y} \cdot n \cdot \frac{\sum x_i}{n}}{\sum (x_i - \bar{x})^2} = \frac{\sum \overset{K_i}{(x_i - \bar{x})} Y_i}{\sum (x_i - \bar{x})^2} = \sum K_i Y_i \quad (1)$$

$$\sum_i K_i = \sum_i \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2} = \frac{1}{\sum (x_i - \bar{x})^2} \sum_i (x_i - \bar{x}) = 0 \quad (2)$$

$$\sum_i K_i^2 = \sum_i \left(\frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2} \right)^2 = \frac{\sum (x_i - \bar{x})^2}{(\sum (x_i - \bar{x})^2)^2} = \frac{1}{\sum (x_i - \bar{x})^2} \quad (3)$$

$$\sum_i K_i x_i = \frac{\sum_i x_i^2 - x_i \bar{x}}{\sum (x_i - \bar{x})^2} = \frac{\sum_i x_i^2 - 2x_i \bar{x} + \bar{x}^2 + x_i \bar{x} - \bar{x}^2}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\sum (x_i - \bar{x})^2 + \bar{x} \sum_i (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = 1 \quad (4)$$

$$E(\beta_1) = E(\sum K_i Y_i) = \sum K_i E(Y_i) = \sum K_i (b_0 + \beta_1 x_i)$$

$$\stackrel{(2)(4)}{=} \beta_0 \sum_i K_i + \beta_1 \sum_i K_i x_i = b_1 \quad (5)$$

$$\text{Var}(\beta_1) = \text{Var}(\sum K_i Y_i) = \sum K_i^2 \text{Var}(Y_i) + \sum_{i \neq j} K_i K_j \text{Cov}(Y_i, Y_j)$$

$$\stackrel{(3)}{=} \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

$$\text{Cov}(\beta_0, \beta_1) = E[(\beta_0 - E(\beta_0))(\beta_1 - E(\beta_1))] \quad (6) \text{ By Definition}$$

$$E(\beta_0) = \bar{Y} - E(\beta_1) \bar{x} \stackrel{(5)}{=} \bar{Y} - b_1 \bar{x} \quad (7)$$

$$(7) \rightarrow \beta_0 - E(\beta_0), \beta_0 - (\bar{Y} - \beta_1 \bar{x}) = \bar{Y} - \beta_1 \bar{x} - \bar{Y} + \beta_1 \bar{x} = -\bar{x}(\beta_1 - b_1) \quad (8)$$

$$(7), (8) \Rightarrow \text{Cov}(\beta_0, \beta_1) = E[(-\bar{x}(\beta_1 - b_1))(\beta_1 - b_1)]$$

$$= E(-\bar{x}(\beta_1 - b_1)^2) = -\bar{x} E(\beta_1 - b_1)^2$$

$$= -\bar{x} \text{Var}(\beta_1) = -\bar{x} \frac{\sigma^2}{\sum_i (x_i - \bar{x})^2}$$

$$\beta_0, \beta_1 \rightarrow \text{independent} \rightarrow \text{Cov}(\beta_0, \beta_1) = 0 \rightarrow -\bar{x} \frac{\sigma^2}{\sum_i (x_i - \bar{x})^2} = 0$$

$$\rightarrow \begin{cases} \bar{x} = 0 \rightarrow \text{میانگین داده ها صفر باشد} \\ \sigma^2 = 0 \rightarrow \text{واریانس خوجی ها صفر باشد} \\ \text{خوجی ها برابر باشند} \end{cases}$$

Problem 5:

$$a) \sigma^2 = \text{Var}(Y) = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n-1}$$

درجه آزادی را $n-1$ فرض کرده ایم

$$\bar{Y} = \frac{\sum_{i=1}^8 Y_i}{n} = \frac{40+41+43+42+44+42+43+42}{8} = 42.125$$

$$\rightarrow \sigma^2 = \frac{\sum_{i=1}^8 (Y_i - 42.125)^2}{7} \approx 1.55$$

$$\bar{x} = \frac{0.5+1+\dots+4}{8} = 1.5$$

$$\beta_1 = \frac{\sum_{i=1}^8 (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^8 (x_i - \bar{x})^2} = \frac{\sum_{i=1}^8 (Y_i - 42.125)(x_i - 1.5)}{\sum_{i=1}^8 (x_i - 1.5)^2} \approx 0.54$$

$$\beta_0 = \bar{Y} - \beta_1 \bar{x} = 42.125 - 0.54 \times 1.5 \approx 40.89$$

Based on (1) in problem 4 solution: $\begin{cases} \beta_1 = \sum K_i Y_i \\ K_i = \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2} \end{cases}$

By variance sum law: $\text{Var}(\beta_1) = \sum_{i=1}^n K_i^2 \text{Var}(Y_i) + \sum_i \sum_{j \neq i} K_i K_j \text{Cov}(Y_i, Y_j)$

$$\rightarrow \text{Var}(\beta_1) = \sum_{i=1}^n K_i^2 \text{Var}(Y_i) = \sigma^2 \sum_{i=1}^n K_i^2 = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

independency \leftarrow

$$= \frac{1.55}{\sum_{i=1}^8 (x_i - 1.5)^2} = \frac{1.55}{10.5} \approx 0.147$$

$$\beta_0 = \bar{Y} - \beta_1 \bar{x} \rightarrow \text{Var}(\beta_0) = \text{Var}(\bar{Y}) + \text{Var}(\beta_1 \bar{x}) - 2\text{Cov}(\bar{Y}, \bar{x} \beta_1)$$

$$= \frac{\sigma^2}{n} + \bar{x}^2 \text{Var}(\beta_1) - 2\bar{x} \text{Cov}(\bar{Y}, \beta_1)$$

$$\text{Cov}(\bar{Y}, \beta_1) = E[(\bar{Y} - E(\bar{Y}))(\beta_1 - E(\beta_1))] = 0$$

$$\rightarrow \text{Var}(\beta_0) = \frac{\sigma^2}{n} + \bar{x}^2 \text{Var}(\beta_1) = \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{\sum (x_i - \bar{x})^2} = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right)$$

$$\rightarrow \text{Var}(\beta_0) = \frac{1.55}{8} + (1.5)^2 (0.147) = 0.52$$

$$b) \text{Cov}(\beta_0, \beta_1) = \frac{-\bar{x} \sigma^2}{\sum (x_i - \bar{x})^2} = \frac{-1.5 \times 1.55}{10.5} \approx -0.22$$

$$\text{Cor}(\beta_0, \beta_1) = \frac{\text{Cov}(\beta_0, \beta_1)}{\sqrt{\text{Var}(\beta_0) \text{Var}(\beta_1)}} = \frac{-0.22}{\sqrt{0.52 \times 0.147}} = -2.87$$

Problem 8:

Suppose for example that we have built a classifier which predicts if someone has cancer or not. We have tested the classifier with 1000 samples and the confusion matrix is:

Actual \ Predict	Positive	Negative
Positive	1	1
Negative	0	998

Accuracy: This metric measures that how much the model was good in predicting the correct class. $\text{Accuracy} = \frac{\# \text{ correct predictions}}{\# \text{ total predictions}}$

$$\text{Here accuracy} = \frac{1+998}{1+1+0+998} = \frac{999}{1000} = 99.9\% \text{ which is very high}$$

Precision: This metric measures that what fraction of the positive predictions are actually positive.

$$\text{precision} = \frac{\# \text{ true positive}}{\# \text{ total positive predictions}} = \frac{1}{0+1} = 100\%$$

Recall: This metric measures that what fraction of actual positives were predicted correctly, i.e. positive.

$$\text{recall} = \frac{\# \text{ true positive}}{\# \text{ total actual positives}} = \frac{1}{1+1} = 50\%$$

In this example although accuracy is very high (99.9%) but the important thing is that samples which actually have cancer (actual positives) must not be predicted negative, otherwise their lives are put in danger. If an actual negative is predicted positive is bad but it does not have as much cost as low recall. So in our example high recall is the most important.

When we want to design a spam email classifier, precision is the most important metric because we do not want an important email to be predicted spam. But if a spam email is predicted non-spam it does not have a high cost. In other examples where the general correctness of classification is important and false predictions are not of a too much cost, the accuracy is the most important metric.