

4.

The wage equation is given as

1.

$$W_i = \beta_0 + \beta_1 X_i + u_i$$

where  $W_i$  is the wage of worker  $i$

$X_i$  is the labour market experience of worker  $i$

A)

We suspect that the intercept is different for men and women. To test this suspicion the model is changed as

$$W_i = \beta_0 + \beta_1 X_i + \beta_2 M + u_i$$

where  $M$  is a dummy variable that can take values 1 and 0. It takes value 1 for man and 0 for woman.

Now,

$$W_i = \beta_0 + \beta_1 X_i + \beta_2 + u_i, \text{ for man } M=1$$

and

$$= \beta_0 + \beta_1 X_i + u_i, \text{ for woman } M=0$$

If the value of  $\beta_2$  is significant then our suspicion is correct and intercept is different for men and women. If  $\beta_2$  is zero then our suspicion is incorrect.

B) We suspect that the slope is different for men and women. To test this suspicion the model is changed as

$$W_i = \beta_0 + (\beta_1 + \beta_2 M) X_i + \mu_i$$

where  $M$  is a dummy variable that can take values 1 and 0. It takes value 1 for man and 0 for woman.

Now,

$$W_i = \beta_0 + (\beta_1 + \beta_2) X_i + \mu_i, \text{ for man } M=1$$

and

$$= \beta_0 + \beta_1 X_i + \mu_i, \text{ for woman } M=0$$

If the value of  $\beta_2$  is significant then our suspicion is correct and the slope is different for men and women. If  $\beta_2$  is zero then our suspicion is incorrect.



DELTA Pg No. 1  
c) We suspect that the relationship between wage and experience resembles an upward slope. To test this suspicion the model is changed as

$$W_i = \beta_0 + \beta_1 X_i^2 + \mu_i$$

We square the  $X_i$  term of the model to make it always positive.

Now if the value of  $W_i$  obtained is positive then it is an upward slope else it is downward.

Since the terms  $\beta_0$  and  $\mu_i$  can also be negative we can square them too. The updated equation becomes

$$W_i = \beta_0^2 + \beta_1 X_i^2 + \mu_i^2$$

Since all the other terms are positive the value obtained on RHS will help us determine if the slope is upward (~~not~~ positive)



2. L2 regularization or Ridge regression promotes smaller coefficients. It is used to prevent the model from overfitting. When we use a model with higher complexity the coefficients of the model become very large and the model starts to overfit and doesn't perform too well on the testing set. To prevent this we use ridge regression. in which a penalty term is added to the cost function

$$J(\theta) = \frac{1}{M} \sum_{i=1}^M ( \theta^T \phi(x_i) - y_i )^2$$

$$+ \lambda \sum_{j=1}^p \theta_j^2$$

penalty term.

Adding this term, reduces the value of coefficients and prevents overfitting

3. Suppose the parameters are given by  $D = \{(x_i, y_i) : i \in [1, N]\}$

It is a linear model  $y = \theta^T x + c$  and  $\theta$  are the weights

Using Bayes Theorem we can write

$$\frac{P(\theta)}{P(D)} = \frac{P(D|\theta) P(\theta)}{P(D)}$$

Taking log on both sides.

$$\log P\left(\frac{\theta}{D}\right) = \log\left(\frac{D}{\theta}\right) + \log P(\theta) - \log P(D)$$

Using maximum a-posterior Inference the solution is given by

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \log\left(\frac{\theta}{D}\right)$$

$$= \underset{\theta}{\operatorname{argmax}} \left( \log\left(\frac{D}{\theta}\right) + \log P(\theta) - \log P(D) \right)$$

Using MAP we can ignore  $-\log P(D)$

$$= \underset{\theta}{\operatorname{argmax}} \left( \log\left(\frac{D}{\theta}\right) + \log P(\theta) \right) \quad \text{--- (1)}$$



log likelihood is given by

$$\log P(\frac{D}{\theta}) = \log P(Y|X, \theta)$$

$$= \log \prod_{n=1}^N P(y_n | x_n, \theta)$$

$$= \sum_{n=1}^N \log \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(y_n - \theta^T x_n)^2}{2\sigma^2}}$$

$$= \sum_{n=1}^N \left( \log \frac{1}{\sqrt{2\pi}\sigma^2} - \frac{(y_n - \theta^T x_n)^2}{2\sigma^2} \right)$$

$$= \sum_{n=1}^N \left( \frac{-1}{2} \log(2\pi\sigma^2) - \frac{(y_n - \theta^T x_n)^2}{2\sigma^2} \right)$$

②

Prior is given by

$$P(\theta) = \frac{1}{(2\pi)^{D/2}} e^{-\frac{\lambda}{2} \theta^T \theta}$$

$\lambda$  is a positive scalar

$$\log P(\theta) = \log \left( \frac{1}{(2\pi)^{D/2}} e^{-\frac{\lambda}{2} \theta^T \theta} \right)$$

$$= -\frac{D}{2} \log(2\pi) - \frac{\lambda}{2} \theta^T \theta \quad \text{--- (3)}$$

Substituting values of ② and ③ in ① and ignoring the constant values that do not depend on  $\theta$  we get

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} \left( \frac{-1}{2} \theta^T \theta - \sum_{n=1}^N \frac{(y_n - \theta^T x_n)^2}{2\sigma} \right)$$

Changing the signs to positive and  $\operatorname{argmax}$  to  $\operatorname{argmin}$  (~~as~~ ~~in~~)

$$= \underset{\theta}{\operatorname{argmin}} \left( \frac{1}{2} \theta^T \theta + \sum_{n=1}^N \frac{(y_n - \theta^T x_n)^2}{2\sigma} \right)$$

This is the equation for L2 regularization.