IDS 702: MODULE 1.2

INTRODUCTION TO MULTIPLE LINEAR REGRESSION

DR. OLANREWAJU MICHAEL AKANDE



MULTIPLE LINEAR REGRESSION

Multiple linear regression (MLR) assumes the following distribution for a response variable y_i given p potential covariates/predictors/features x_i = (x_{i1}, x_{i2}, ..., x_{ip}).

 $y_i=eta_0+eta_1x_{i1}+eta_2x_{i2}+\ldots+eta_px_{ip}+\epsilon_i; \ \ \epsilon_i\stackrel{iid}{\sim}\mathcal{N}(0,\sigma^2), \ \ i=1,\ldots,n.$

• We can also write the model as:

$$egin{aligned} y_i \stackrel{iid}{\sim} \mathcal{N}(eta_0 + eta_1 x_{i1} + eta_2 x_{i2} + \ldots + eta_p x_{ip}, \sigma^2). \ p(y_i | oldsymbol{x}_i) &= \mathcal{N}(eta_0 + eta_1 x_{i1} + eta_2 x_{i2} + \ldots + eta_p x_{ip}, \sigma^2). \end{aligned}$$

- MLR assumes that the conditional average or expected value of a response variable is a linear function of potential predictors.
- Note that the linearity is in terms of the "unknown" parameters (intercept and slopes).
- Just like in SLR, MLR also assumes values of the response variable follow a normal curve within any combination of predictors.



MLR

- Just as we had under SLR, here each β_j represents the true "unknown" value of the parameter, while $\hat{\beta}_j$ represents the estimate of β_j .
- Similarly, y_i represents the true value of the response variable, while \hat{y}_i represents the predicted value. That is,

$${\hat y}_i={\hateta}_0+{\hateta}_1x_{i1}+{\hateta}x_{i2}+\ldots+{\hateta}x_{ip}.$$

- Also, the residuals e_i are our estimates of the true "unobserved" errors ϵ_i . Thus,

$$e_i=y_i-\left|{\hateta}_0+{\hateta}_1x_{i1}+{\hateta}x_{i2}+\ldots+{\hateta}x_{ip}
ight|=y_i-{\hat y}_i.$$

- Since the e_i's estimate the e_i's, we expect them to also be independent, centered at zero, and have constant variance.
- We will get into this more under model assessment.

MLR: ESTIMATION

 Estimated coefficients are found by taking partial derivatives of the sum of squares of the errors

$$\sum_{i=1}^n \left(y_i - \left[eta_0 + eta_1 x_{i1} + eta_2 x_{i2} + \ldots + eta_p x_{ip}
ight]
ight)^2,$$

with respect to each parameter, that is, $\beta_0, \beta_1, \ldots, \beta_p$.

- This is the ordinary least squares (OLS) method.
- Resulting formulas are a bit messy to write down in this form.
- However, there is a very nice matrix algebra representation as we will see soon.



MLR: ESTIMATION

- An alternative derivation uses maximum likelihood estimation (MLE).
- First, not that if each Y_i , with $i=1,\ldots,n$, follows the normal distribution $Y_i\sim \mathcal{N}(\mu,\sigma^2)$, then the likelihood is

$$egin{aligned} L(\mu,\sigma^2|y_1,\ldots,y_n) &= \prod_{i=1}^n \left(2\pi\sigma^2
ight)^{-rac{1}{2}} e^{-rac{1}{2\sigma^2}(y_i-\mu)^2} \ &= \left(2\pi\sigma^2
ight)^{-rac{n}{2}} e^{-rac{1}{2\sigma^2}\sum\limits_{i=1}^n (y_i-\mu)^2}. \end{aligned}$$

• So that for MLR, the likelihood is

$$L(eta_0,eta_1,\dots,eta_p,\sigma^2|y_1,\dots,y_n) = ig(2\pi\sigma^2ig)^{-rac{n}{2}} e^{-rac{1}{2\sigma^2}\sum\limits_{i=1}^n (y_i - [eta_0+eta_1x_{i1}+\dots+eta_px_{ip}])^2}.$$

- To get the MLEs, take the log of the likelihood, differentiate with respect to each parameter in $(\beta_0, \beta_1, \ldots, \beta_p, \sigma^2)$, and set to zero.
- Again, resulting formulas for $(\beta_0, \beta_1, \dots, \beta_p)$ are a bit messy to write down in this form.



MLR: ESTIMATION

• The MLE for σ^2 (work it out to convince yourself) is

$$\hat{\sigma}_{ ext{MLE}}^2 = rac{1}{n}\sum_{i=1}^n \left(y_i - \left[\hat{eta}_0 + \hat{eta}_1 x_{i1} + \ldots + \hat{eta}_p x_{ip}
ight]
ight)^2
onumber \ = rac{1}{n}\sum_{i=1}^n \left(y_i - \hat{y}_i
ight)^2 = rac{1}{n}\sum_{i=1}^n e_i^2.$$

- However, the MLE is biased. That is, $\mathbb{E}[\hat{\sigma}_{\mathrm{MLE}}^2] \neq \sigma^2$.
- Therefore, we often used the following "unbiased" estimator for σ^2 .

$$\hat{\sigma}^2 = s_e^2 = rac{1}{n-(p+1)}\sum_{i=1}^n \left(y_i - \hat{y}_i
ight)^2 = rac{1}{n-(p+1)}\sum_{i=1}^n e_i^2.$$

• Most software packages will estimate s_e^2 automatically.



MLR: MATRIX REPRESENTATION

Let

$$oldsymbol{y} = egin{bmatrix} y_1 \ y_2 \ dots \ y_n \end{bmatrix} oldsymbol{X} = egin{bmatrix} 1 & x_{11} & x_{12} & \ldots & x_{1p} \ 1 & x_{21} & x_{22} & \ldots & x_{2p} \ dots & dots$$

• Then, we can write the MLR model as

$$oldsymbol{y} = oldsymbol{X}oldsymbol{eta} + oldsymbol{\epsilon}; \;\;oldsymbol{\epsilon} \sim \mathcal{N}(0,\sigma^2oldsymbol{I}).$$

- The OLS and MLE estimates of all (p+1) coefficients (intercept plus p slopes) is then given by

$$\hat{oldsymbol{eta}} = ig(oldsymbol{X}^Toldsymbol{X}ig)^{-1}oldsymbol{X}^Toldsymbol{y}.$$

Ideally, n should be bigger than p. Why?

There are many ways around the p>n problem. If there is time, we may look at some options.

IDS 702

MLR: MATRIX REPRESENTATION

The predictions can then be written as

$$\hat{oldsymbol{y}} = oldsymbol{X} \hat{oldsymbol{eta}} = oldsymbol{X} \left[oldsymbol{\left(X^T oldsymbol{X}ig)^{-1} oldsymbol{X}^T oldsymbol{y}}
ight] = \left[oldsymbol{X} oldsymbol{\left(X^T oldsymbol{X}ig)^{-1} oldsymbol{X}^T
ight] oldsymbol{y}.$$

• The residuals can be written as

$$oldsymbol{e} = oldsymbol{y} - \hat{oldsymbol{y}} = oldsymbol{y} - \left[oldsymbol{X}ig(oldsymbol{X}^Toldsymbol{X}ig)^{-1}oldsymbol{X}^Tig]oldsymbol{y} = \left[oldsymbol{1}_n - oldsymbol{X}ig(oldsymbol{X}^Toldsymbol{X}ig)^{-1}oldsymbol{X}^Tig]oldsymbol{y}$$

where $\mathbf{1}_n$ is a matrix of ones

• The n imes n matrix

$$oldsymbol{H} = oldsymbol{X}ig(oldsymbol{X}^Toldsymbol{X}ig)^{-1}oldsymbol{X}^T$$

is often called the projection matrix or the hat matrix.

• We will see some important features of the elements of *H* soon.



MLR: MATRIX REPRESENTATION

In matrix form,

$$s_e^2 = \sum_{i=1}^n rac{\left(y_i - \hat{y}_i
ight)^2}{n - (p+1)} = rac{(m{y} - m{X}\hat{m{eta}})^T(m{y} - m{X}\hat{m{eta}})}{n - (p+1)} = rac{m{e}^Tm{e}}{n - (p+1)}.$$

- The variance of the OLS estimates of all (p+1) coefficients (intercept plus p slopes) is

$$\mathbb{V}\left[\hat{oldsymbol{eta}}
ight] = \sigma^2ig(oldsymbol{X}^Toldsymbol{X}ig)^{-1}$$

 Notice that this is a covariance matrix; the square root of the diagonal elements give us the standard errors for each β_j, which we can use for hypothesis testing and interval estimation.

What are the off-diagonal elements?

- When estimating $\mathbb{V}[\hat{m{eta}}]$, plug in s_e^2 as an estimate of σ^2 .
- Now that we have a basic introduction, we are ready see how to fit MLR models.

WHAT'S NEXT?

Move on to the readings for the next module!

