

PSTAT 126

Regression Analysis

Laura Baracaldo

Lecture 6

Multiple Linear Regression

Multiple Linear Regression Models (MLR)

Consider the linear regression model with p predictors:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$$

For each $i = 1, \dots, n$ we assume $\epsilon_i \stackrel{i.i.d}{\sim} N(0, \sigma^2)$. Thus $y_i \stackrel{ind}{\sim} N(\mu_i, \sigma^2)$, with $\mu_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$.

Multiple Linear Regression Models (MLR)

Matrix Representation: By stacking up all the observations, we can obtain the matrix representation of the MLR model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

This implies $\mathbf{y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$.

Estimation

LS and ML Estimation of β : Provided $\mathbf{X}^T \mathbf{X}$ is non-singular:

$$\hat{\beta}_{LSE} = \hat{\beta}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

So that $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{H}\mathbf{y}$, with $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ the *Projection Matrix* or *Hat Matrix*, which corresponds to the orthogonal projection onto the column space of \mathbf{X} .

Estimation of σ^2 :

$$\hat{\sigma}^2 = \frac{\hat{\epsilon}^T \hat{\epsilon}}{n - p^*} = \frac{SSR}{n - p^*}$$

Algebraic Properties of the Projection Matrix H

- H is symmetric: $H^T = H$.
- H is idempotent: $H^2 = H$.
- If X is a $n \times p^*$ matrix with $\text{rank}(X) = p^*$, then $\text{rank}(H) = p^*$.
- X is invariant under H : $HX = X$
- The eigenvalues of H consist of p^* ones, and $n - p^*$ zeros.

Algebraic Properties of the Projection Matrix H

Note that the residuals of the MLR model can be written as:

$$\hat{\epsilon} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{H}\mathbf{y} = (\mathbf{I} - \mathbf{H})\mathbf{y} = \mathbf{M}\mathbf{y}$$

- $\mathbf{M} = (\mathbf{I} - \mathbf{H})$ is symmetric: $\mathbf{M}^T = \mathbf{M}$.
- \mathbf{M} is idempotent: $\mathbf{M}^2 = \mathbf{M}$
- $\text{rank}(\mathbf{M}) = n - p^*$.
- The eigenvalues of \mathbf{M} consist of $n - p^*$ ones, and p^* zeros.
- $\hat{\epsilon}$ and \mathbf{X} are orthogonal: $\hat{\epsilon}^T \mathbf{X} = \mathbf{0}$. Under normality this is equivalent to statistical independence.
- \mathbf{M} and \mathbf{H} are orthogonal: $\mathbf{MH} = \mathbf{0}$. This implies $\hat{\epsilon}$ and $\hat{\beta}$ are independent (Under normality).

Inference on β and σ^2

- Distribution of $\hat{\beta}$:

$$\begin{aligned}E(\hat{\beta}) &= E \left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \right] \\&= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E(\mathbf{y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta\end{aligned}$$

$$\begin{aligned}V(\hat{\beta}) &= V \left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \right] \\&= \left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right] V(\mathbf{y}) \left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right]^T = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2\end{aligned}$$

$$Var(\hat{\beta}_j) = \sigma^2 \left[\mathbf{X}^T \mathbf{X} \right]_{jj}^{-1} \quad Cov(\hat{\beta}_i, \hat{\beta}_j) = \sigma^2 \left[\mathbf{X}^T \mathbf{X} \right]_{ij}^{-1}$$

$$\Rightarrow \hat{\beta} \sim N_{p^*}(\beta, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2).$$

Gauss-Markov Theorem: $\hat{\beta}$ is the best linear unbiased estimate (BLUE) of β .

Inference on β and σ^2

- **Distribution of $\hat{\sigma}^2$.** Note that:

$$\hat{\epsilon} = M\mathbf{y} = M(\mathbf{X}\beta + \epsilon) = M\epsilon$$

Lemma

If A is a symmetric and idempotent $n \times n$ real matrix and $Z \sim N(0, I_n)$ is a random vector of n independent standard normal variables, then $Z^T A Z$ has chi-squared(r) distribution, r being the trace of A .

From the lemma we can prove that $\frac{\hat{\epsilon}^T \hat{\epsilon}}{\sigma^2} \sim \chi^2_{(n-p^*)}$ since:

$$\frac{\hat{\epsilon}^T \hat{\epsilon}}{\sigma^2} = \frac{\epsilon^T M \epsilon}{\sigma^2}$$

Species Example

```
Beta.hat<- solve(crossprod(X))%*(t(X)%*y);t(Beta.hat)
```

```
##      Intercept      Area Elevation      Scruz      Adjacent
## [1,]  7.075377 -0.02397793 0.3195734 -0.2393552 -0.07484842
```

```
sigma.hat <- sqrt(sum(fit1$residuals^2)/(fit1$df.residual)); sigma.hat
```

```
## [1] 59.74333
```

```
XtX.inverse <- solve(crossprod(X)); XtX.inverse
```

```
##      Intercept      Area      Elevation      Scruz      Adjacent
## Intercept  9.849524e-02  3.879513e-05 -1.589754e-04 -4.059337e-04  2.421334e-05
## Area      3.879513e-05  1.296222e-07 -2.439943e-07  1.163061e-07  4.003452e-08
## Elevation -1.589754e-04 -2.439943e-07  7.323822e-07 -1.457015e-07 -1.471043e-07
## Scruz     -4.059337e-04  1.163061e-07 -1.457015e-07  7.594187e-06 -1.368476e-08
## Adjacent  2.421334e-05  4.003452e-08 -1.471043e-07 -1.368476e-08  7.747375e-08
```

```
Beta.hat.SE <- sigma.hat*sqrt(diag(XtX.inverse)); Beta.hat.SE
```

```
##      Intercept      Area      Elevation      Scruz      Adjacent
## [1,] 18.74981844  0.02150944  0.05112795  0.16463800  0.01662902
```

Species Example

```
##
## Call:
## lm(formula = Species ~ Area + Elevation + +Scruz + Adjacent,
##     data = gala)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -111.637  -34.930   -7.864   33.432  182.524
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.07538    18.74982   0.377 0.709093
## Area        -0.02398     0.02151  -1.115 0.275554
## Elevation     0.31957     0.05113   6.250 1.54e-06 ***
## Scruz        -0.23936     0.16464  -1.454 0.158434
## Adjacent     -0.07485     0.01663  -4.501 0.000136 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 59.74 on 25 degrees of freedom
## Multiple R-squared:  0.7658, Adjusted R-squared:  0.7284
## F-statistic: 20.44 on 4 and 25 DF,  p-value: 1.39e-07
```

Model Performance

- ① Goodness of fit.
- ② Estimators Quality.
- ③ Usefulness of predictors.
- ④ Prediction accuracy.

Goodness of fit

- Residual Standard Error (RSE) $\hat{\sigma}$

$$\hat{\sigma} = \sqrt{\frac{\hat{\epsilon}^T \hat{\epsilon}}{n - p^*}}$$

The smaller $\hat{\sigma}$ the better. But, how small?

Determination coefficient R^2

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{SSR}{SST} = 1 - \frac{\mathbf{y}^T \mathbf{M} \mathbf{y}}{\mathbf{y}^T \mathbf{M}_1 \mathbf{y}}$$

Where $\mathbf{M}_1 = (\mathbf{I} - \mathbf{1}(\mathbf{1}^T \mathbf{1})^{-1} \mathbf{1}^T)$

It can be proved that:

$$R^2 = \text{cor}(\mathbf{y}, \hat{\mathbf{y}})$$

Hypothesis tests for a set of predictors

We want to assess the usefulness of the predictors in the prediction of the response. For instance, let's suppose we want to test whether there is at least one predictor useful in the prediction of the response:

$$H_0 : \beta_1 = \dots = \beta_p = 0 \quad H_1 : \beta_j \neq 0 \text{ for at least one } j \in 1, \dots, p.$$

This can be tested by comparing the full model: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ to the null model: $\mathbf{y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$.

General Hypothesis Test

Consider two models M_1 and M_2 , such that $c(M_1) \subset c(M_2)$. Where $c(M)$ denotes the column space spanned by the predictors matrix of model M . This means that predictors included in model M_1 is a subset of predictors included in M_2 . How to decide which model is better?

- If there is not much difference in how well both models fit the data, we go with the *smaller* model M_1 .
- If the *larger* model M_2 fits better because of the additional variables, we go with model M_2 .

General Hypothesis Test

We have learned that we can evaluate the performance of a model by calculating the SSR . Small values of SSR will imply a better fit of the model. For $c(M_1) \subset c(M_2)$ it is always true that:

$$SSR_{M_2} \leq SSR_{M_1}$$

Let's consider the difference $SSR_{M_1} - SSR_{M_2}$. If this difference is small, there is not substantial difference between the fit of the larger model and the smaller model. Thus, we go with M_1 . If this difference is significant, the fit of the smaller model is substantially worse than the fit of the large model. Therefore, we go with model M_2 .

F-Test

We want to test the hypothesis:

$$H_0 : SSR_{M_1} = SSR_{M_2} \quad H_1 : SSR_{M_1} > SSR_{M_2}$$

We consider the F-Statistic:

$$\begin{aligned} F &= \frac{(SSR_{M_1} - SSR_{M_2}) / (df_{M_1} - df_{M_2})}{SSR_{M_2} / df_{M_2}} \\ &= \frac{\mathbf{y}^T (\mathbf{H}_2 - \mathbf{H}_1) \mathbf{y} / (df_{M_1} - df_{M_2})}{\mathbf{y}^T (\mathbf{I} - \mathbf{H}_2) \mathbf{y} / df_{M_2}} \sim F_{(df_{M_1} - df_{M_2}, df_{M_2})} \end{aligned}$$

With \mathbf{H}_1 and \mathbf{H}_2 the projection matrices for models M_1 and M_2 respectively. $F \sim F_{(df_{M_1} - df_{M_2}, df_{M_2})}$, since $\frac{\mathbf{y}^T (\mathbf{H}_2 - \mathbf{H}_1) \mathbf{y}}{\sigma^2} \sim \chi^2_{df_{M_1} - df_{M_2}}$ and $\frac{\mathbf{y}^T (\mathbf{I} - \mathbf{H}_2) \mathbf{y}}{\sigma^2} \sim \chi^2_{df_{M_2}}$

Case 1: Global F-Test

The global F-test:

$$H_0 : \beta_1 = \dots = \beta_p = 0 \quad H_1 : \beta_j \neq 0 \text{ for at least one } j \in 1, \dots, p$$

Which is equivalent to:

$H_0 : M_0$ and M_F fit the data similarly well , $H_1 : M_F$ fits the data better

Where M_0 is the null model (with no predictors): $\mathbf{y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$, and M_F is the full model (with p predictors): $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$.

Case 1: Global F-Test

Total Sum of Squares:

$$SSR_{M_0} = \mathbf{y}^T(\mathbf{I} - \mathbf{H}_0)\mathbf{y} = \sum_{i=1}^n (y_i - \bar{y})^2, \text{ with } df = n - 1.$$

Residual Sum of Squares:

$$SSR_{M_F} = \mathbf{y}^T(\mathbf{I} - \mathbf{H}_F)\mathbf{y} = \sum_{i=1}^n (y_i - \hat{y}_i)^2, \text{ with } df = n - p^*.$$

Regression Sum of Squares:

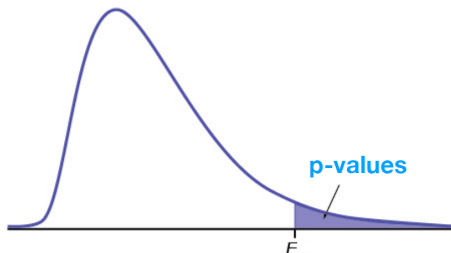
$$\Rightarrow SSR_{M_0} - SSR_{M_F} = \mathbf{y}^T(\mathbf{H}_F - \mathbf{H}_0)\mathbf{y} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2, \text{ with } df = p.$$

Recall $\mathbf{H}_0 = \mathbf{X}_0(\mathbf{X}_0^T \mathbf{X}_0)^{-1} \mathbf{X}_0^T = \mathbf{1}(\mathbf{1}^T \mathbf{1})^{-1} \mathbf{1}^T$ and
 $\mathbf{H}_F = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

Case 1: Global F-Test

For a significance level α , we reject H_0 if $F > F_{(1-\alpha; df_{M_1} - df_{M_2}, df_{M_2})}$ or equivalently, if $p\text{-value} < \alpha$.

Recall: p -value: Probability of obtaining tests results as least as extreme as the observed results.



Case 2: F-test for a pair of predictors

$$H_0 : \beta_l = \beta_k = 0 \quad H_1 : \beta_l \neq 0 \vee \beta_k \neq 0$$

Which is equivalent to:

$H_0 : M_1$ and M_F fit the data similarly well , $H_1 : M_1$ fits the data better

Where M_1 is the model that does not include predictors X_l nor X_k and M_F is the full model (with p predictors).

Case 2: F-test for a pair of predictors

$$SSR_{M_1} = \mathbf{y}^T(\mathbf{I} - \mathbf{H}_1)\mathbf{y}, \text{ with } df = n - (p^* - 2) = n - p + 1.$$

$$SSR_{M_F} = \mathbf{y}^T(\mathbf{I} - \mathbf{H}_F)\mathbf{y}, \text{ with } df = n - p^* = n - p - 1.$$

$$\Rightarrow SSR_{M_1} - SSR_{M_F} = \mathbf{y}^T(\mathbf{H}_F - \mathbf{H}_1)\mathbf{y}, \text{ with } df = 2.$$

Recall $H_1 = \mathbf{X}_1(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T$, with \mathbf{X}_1 the matrix without predictors X_l and X_k and $H_F = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.

Species Example: $\hat{\sigma}$ and R^2

```
fit1<- lm( Species ~ Area+Elevation+ + Scrutz+ Adjacent, data=gala)
RSE<- sqrt(sum(fit1$residuals^2)/fit1$df.residual);RSE ## With Formula
```

```
## [1] 59.74333
```

```
sigma(fit1) ##With lm function
```

```
## [1] 59.74333
```

```
y.hat<- fit1$fitted.values
R2<- cor(y.hat,y)^2;R2 ## With Formula
```

```
## [1] 0.7658462
```

```
summary(fit1)$r.squared ##With lm function
```

```
## [1] 0.7658462
```

Species Example: Global F-Test

$H_0 : \beta_{Area} = \beta_{Elevation} = \beta_{Scruz} = \beta_{Adjacent} = 0$ $H_1 : \beta_j \neq 0$ for at least one $j \in \{Area, Elevation, Scrutz, Adjacent\}$

```
fullmodel<- lm( Species ~ Area+Elevation+ Scrutz+ Adjacent, data=gala)
nullmodel <- lm(Species~1, data=gala)
anova1<-anova(nullmodel, fullmodel);anova1
```

```
## Analysis of Variance Table
##
## Model 1: Species ~ 1
## Model 2: Species ~ Area + Elevation + Scrutz + Adjacent
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      29 381081
## 2      25  89232  4    291850 20.442 1.39e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
pval<- 1-pf(anova1$F[2],4,25);pval
```

```
## [1] 1.389845e-07
```

Species Example: F-test for a pair of predictors

$$H_0 : \beta_{Area} = \beta_{Scruz} = 0 \quad H_1 : \beta_{Area} \neq 0 \vee \beta_{Scruz} \neq 0$$

```
fullmodel<- lm( Species ~ Area+Elevation+ Scruz+ Adjacent, data=gala)
Model1 <- lm(Species~Elevation+ Adjacent, data=gala)
anova2<-anova(Model1, fullmodel);anova2
```

```
## Analysis of Variance Table
```

```
##
```

```
## Model 1: Species ~ Elevation + Adjacent
```

```
## Model 2: Species ~ Area + Elevation + Scruz + Adjacent
```

```
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)
```

```
## 1      27 100003
```

```
## 2      25  89232  2      10771 1.5089 0.2406
```

```
pval<- 1-pf(anova2$F[2],2,25);pval
```

```
## [1] 0.2406149
```