# Math 423/533: The Main Theoretical Topics

- sample size $n$, data index $i$
- number of predictors, $p$ ($p = 2$ for simple linear regression)
- $y_i$: response for individual $i$
- $\mathbf{x}_i = (x_{i1}, \ldots, x_{ip})^\top$ – $(1 \times p)$ row vector
- $\mathbf{X}$ - $(n \times p)$ matrix containing all predictors for all individuals $i = 1, \ldots, n$.
- $\mathbf{y} = (y_1, \ldots, y_n)^\top$ – $(n \times 1)$ column vector
- $Y_i$ and $\mathbf{Y}$: random variables corresponding to responses

## Linear model assumptions

For $i = 1, \ldots, n$,

$$\mathbb{E}_{Y_i|\mathbf{x}_i}[Y_i|\mathbf{x}_i] = \mathbf{x}_i\beta = \sum_{j=1}^{p} \beta_j x_{ij}$$

and

$$\mathbb{V}\mathrm{ar}_{Y_i|\mathbf{x}_i}[Y_i|\mathbf{x}_i] = \sigma^2$$

where

$$\beta = (\beta_1, \ldots, \beta_p)^\top$$

is the $(p \times 1)$ vector of regression coefficients, and $\sigma^2 > 0$ is the error variance.

We assume also that $Y_1, \ldots, Y_n$ are independent given $\mathbf{x}_1, \ldots, \mathbf{x}_n$.

## Linear model assumptions

In vector form

$$\mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}|\mathbf{X}] = \mathbf{X}\beta \qquad (n \times 1)$$

and

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}|\mathbf{X}] = \sigma^2\mathbf{I}_n \qquad (n \times n).$$

This is equivalent to a model specification of

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

where

$$\mathbb{E}_{\epsilon|\mathbf{X}}[\epsilon|\mathbf{X}] = \mathbf{0}_n \qquad\qquad \mathbb{V}\mathrm{ar}_{\epsilon|\mathbf{X}}[\epsilon|\mathbf{X}] = \sigma^2\mathbf{I}_n.$$

We usually consider including the 'special' predictor

$$x_{i0} \equiv 1 \qquad i = 1, \ldots, n$$

and specify the model

$$\mathbb{E}_{Y_i|\mathbf{x}_i}[Y_i|\mathbf{x}_i] = \mathbf{x}_i\beta = \beta_0 + \sum_{j=1}^{p} \beta_j x_{ij} = \sum_{j=0}^{p} \beta_j x_{ij}$$

This model has $p + 1$ $\beta$ parameters.

We will let $p$ count the total number of predictors, including the intercept term.

We specify

$$\mathbb{E}_{Y_i|\mathbf{x}_i}[Y_i|\mathbf{x}_i] = \beta_0 + \beta_1 x_{i1} = [1 \; x_{i1}] \left[ \begin{array}{c} \beta_0 \\ \beta_1 \end{array} \right] = \mathbf{x}_i \beta$$

with $p = 2$ parameters in the regression model.

This model posits a straight line relationship between $x$ and $y$.

On the basis of data $(x_{i1}, y_i)$, $i = 1, \ldots, n$, we choose the line of best fit according to the least squares principle. We estimate parameters $\beta = (\beta_0, \beta_1)^\top$ by $\widehat{\beta}$ where

$$\widehat{\beta} = \arg\min_\beta \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1})^2 = \arg\min_\beta S(\beta)$$

where we may also write, in vector form,

$$S(\beta) = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta)$$

We achieve the minimization by calculus.

## The Normal Equations

We solve

$$\frac{\partial S(\beta)}{\partial \beta_0} = -2 \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_{i1}) = 0$$

$$\frac{\partial S(\beta)}{\partial \beta_1} = -2 \sum_{i=1}^{n} x_{i1}(y_i - \beta_0 - \beta_1 x_{i1}) = 0$$

These two equations can be written

$$n\beta_0 + \beta_1 \sum_{i=1}^{n} x_{i1} = \sum_{i=1}^{n} y_i$$

$$\beta_0 \sum_{i=1}^{n} x_{i1} + \beta_1 \sum_{i=1}^{n} x_{i1}^2 = \sum_{i=1}^{n} x_{i1} y_i$$

or, in matrix form

$$(\mathbf{X}^\top \mathbf{X})\beta = \mathbf{X}^\top \mathbf{y}$$

## The Normal Equations

These equations are the 'Normal Equations'. If the symmetric $p \times p = 2 \times 2$ matrix

$$\mathbf{X}^\top \mathbf{X}$$

is non-singular, then we may write the solution

$$\widehat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

which yields a $p \times 1 = 2 \times 1$ vector of least squares estimates.

Explicitly, we have

$$\left[ \begin{array}{c} \widehat{\beta_0} \\ \widehat{\beta_1} \end{array} \right] = \left[ \begin{array}{cc} n & \sum x_{i1} \\ \sum x_{i1} & \sum x_{i1}^2 \end{array} \right]^{-1} \left[ \begin{array}{c} \sum y_i \\ \sum x_{i1} y_i \end{array} \right]$$

# Estimates

$$
\begin{bmatrix} n & \sum_{i=1}^{n} x_{i1} \\ \sum_{i=1}^{n} x_{i1} & \sum_{i=1}^{n} x_{i1}^2 \end{bmatrix}^{-1} = \frac{1}{n \sum_{i=1}^{n} x_{i1}^2 - \{\sum_{i=1}^{n} x_{i1}\}^2} \begin{bmatrix} \sum_{i=1}^{n} x_{i1}^2 & -\sum_{i=1}^{n} x_{i1} \\ -\sum_{i=1}^{n} x_{i1} & n \end{bmatrix}
$$

We write

$$
\begin{aligned}
S_{xx} &= \sum_{i=1}^{n} x_{i1}^2 - \left\{\sum_{i=1}^{n} x_{i1}\right\}^2 = \sum_{i=1}^{n} (x_{i1} - \bar{x}_1)^2 \\
S_{xy} &= \sum_{i=1}^{n} x_{i1} y_i - \frac{1}{n} \left\{\sum_{i=1}^{n} x_{i1} \sum_{i=1}^{n} y_i\right\} = \sum_{i=1}^{n} y_i (x_{i1} - \bar{x}_1)
\end{aligned}
$$

Thus, after some algebra

$$\begin{aligned} \widehat{\beta}_0 &= \overline{y} - \widehat{\beta}_1 \overline{x}_1 \\ \widehat{\beta}_1 &= \frac{S_{xy}}{S_{xx}} \end{aligned}$$

Define for $i = 1, \ldots, n$,

$$e_i = y_i - (\widehat{\beta}_0 + \widehat{\beta}_1 x_{i1}) = y_i - \widehat{y}_i$$

- $e_i$ – $i$th residual
- $\widehat{y}_i$ – $i$th fitted value.

# Statistical properties of least squares estimators

It is evident from the formula

$$\widehat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{A}\mathbf{y}$$

say, where

$$\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

that the least squares estimates are merely linear combinations of the observed responses $\mathbf{y} = (y_1, \ldots, y_n)^\top$.

Specifically in the simple linear regression

$$\widehat{\beta}_0 = \sum_{i=1}^{n} \left( \frac{1}{n} - \bar{x}_1 c_i \right) y_i \qquad\qquad \widehat{\beta}_1 = \sum_{i=1}^{n} c_i y_i$$

where, for $i = 1, \ldots, n$,

$$c_i = \frac{x_{i1} - \bar{x}_1}{S_{xx}}.$$

In random variable form, we have the estimators

$$\widehat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \mathbf{A}\mathbf{Y}$$

and thus, under the model assumptions

$$\mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}|\mathbf{X}] = \mathbf{X}\beta$$

and

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}|\mathbf{X}] = \sigma^2 \mathbf{I}_n$$

we can study distributional properties of the estimators.

We have, using elementary properties of expectation and variance,

$$\mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\widehat{\beta}|\mathbf{X}] = \beta \qquad (p \times 1)$$

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\widehat{\beta}|\mathbf{X}] = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} \qquad (p \times p)$$

with $p = 2$. Explicitly

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\widehat{\beta}_0|\mathbf{X}] = \sigma^2 \frac{\sum x_{i1}^2}{n S_{xx}} = \sigma^2 \left( \frac{1}{n} + \frac{(\bar{x}_1)^2}{S_{xx}} \right)$$

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\widehat{\beta}_1|\mathbf{X}] = \frac{\sigma^2}{S_{xx}}$$

For the simple linear regression

1. $\sum\limits_{i=1}^{n} e_i = \mathbf{1}_n^\top \mathbf{e} = 0.$

2. $\sum\limits_{i=1}^{n} x_{i1} e_i = \underset{\sim}{x}_1^\top \mathbf{e} = 0$

3. $\sum\limits_{i=1}^{n} y_i = \sum\limits_{i=1}^{n} \widehat{y}_i.$

4. $\sum\limits_{i=1}^{n} \widehat{y}_i e_i = 0,$

that is, the observed residual vector $\mathbf{e}$ is orthogonal to the observed $n \times 1$ vectors

$$\underset{\sim}{\mathbf{x}}_1 = (x_{11}, \ldots, x_{n1})^\top \qquad \text{and} \qquad \widehat{\mathbf{y}} = (\widehat{y}_1, \ldots, \widehat{y}_n)^\top.$$

These results arise as $\widehat{\beta}$ solves

$$\mathbf{X}^{\top}(\mathbf{y} - \mathbf{X}\widehat{\beta}) = \mathbf{0}_p$$

and the results above follow immediately.

# Estimating $\sigma^2$

Let

$$
\begin{aligned}
\mathrm{SS_{Res}} &= \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n}(y_i - \widehat{y_i})^2 \\
&= \sum_{i=1}^{n} y_i^2 - n(\bar{y})^2 - \widehat{\beta}_1 S_{xy} \\
&= \mathrm{SS_T} - \widehat{\beta}_1 S_{xy}
\end{aligned}
$$

say, where

$$
\mathrm{SS_T} = \sum_{i=1}^{n} y_i^2 - n(\bar{y})^2 = \sum_{i=1}^{n}(y_i - \bar{y})^2
$$

We study the statistical properties of the random variable

$$
\sum_{i=1}^{n}(Y_i - \widehat{Y_i})^2 = \sum_{i=1}^{n}(Y_i - \mathbf{x}_i\widehat{\beta})^2 = (\mathbf{Y} - \mathbf{X}\widehat{\beta})^{\top}(\mathbf{Y} - \mathbf{X}\widehat{\beta})
$$

where

$$\widehat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

is the vector of least squares estimators.

But

$$\widehat{\mathbf{Y}} = \mathbf{X}\widehat{\beta} = \mathbf{X}(\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

say, where

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}$$

is the 'hat matrix'.

We can show that $\mathbf{H}$ is symmetric, and that

$$\mathbf{H}^{\top}\mathbf{H} = \mathbf{H}$$

so $\mathbf{H}$ is idempotent.

Now consider the simpler model where dependence on $x_{i1}$ is omitted, and we merely have an intercept term. Predictions in this model use the $(n \times 1)$ matrix

$$\mathbf{X} = (1, 1, \ldots, 1)^\top = \mathbf{1}_n$$

yielding the corresponding hat matrix

$$\mathbf{H}_1 = \mathbf{1}_n (\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$$

which is merely the $(n \times n)$ matrix with all elements equal to $1/n$.

We have that

$$\text{SS}_{\text{Res}} = (\mathbf{Y} - \mathbf{X}\widehat{\beta})^\top (\mathbf{Y} - \mathbf{X}\widehat{\beta}) = \mathbf{Y}^\top (\mathbf{I}_n - \mathbf{H}) \mathbf{Y}$$

where the $(n \times n)$ matrix $(\mathbf{I}_n - \mathbf{H})$ is symmetric and idempotent.

Now, we have the sum of squares decomposition

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \widehat{y}_i)^2 + \sum_{i=1}^n (\widehat{y}_i - \bar{y})^2$$

or

$$\text{SS}_{\text{T}} = \text{SS}_{\text{Res}} + \text{SS}_{\text{R}}$$

## Estimating $\sigma^2$ (cont.)

Similarly to the previous result we have

$$SS_T = \mathbf{Y}^\top (\mathbf{I}_n - \mathbf{H}_1)\mathbf{Y}$$

and

$$SS_R = \mathbf{Y}^\top (\mathbf{H} - \mathbf{H}_1)\mathbf{Y}$$

yielding the representation

$$\mathbf{Y}^\top (\mathbf{I}_n - \mathbf{H}_1)\mathbf{Y} = \mathbf{Y}^\top (\mathbf{I}_n - \mathbf{H})\mathbf{Y} + \mathbf{Y}^\top (\mathbf{H} - \mathbf{H}_1)\mathbf{Y}$$

where the $(n \times n)$ matrices $(\mathbf{I}_n - \mathbf{H}_1)$ and $(\mathbf{H} - \mathbf{H}_1)$ are also symmetric and idempotent.

## Estimating $\sigma^2$ (cont.)

Using the result for the expectation of a quadratic form that if $\mathbf{V}$ is a $k$-dimensional random vector with

$$\mathbb{E}[\mathbf{V}] = \mu \qquad \mathbb{V}\mathrm{ar}[\mathbf{V}] = \Sigma$$

then for $k \times k$ matrix $\mathbf{A}$, we have

$$\mathbb{E}[\mathbf{V}^\top \mathbf{A} \mathbf{V}] = \mathrm{trace}(\mathbf{A}\Sigma) + \mu^\top \mathbf{A} \mu$$

it follows that

$$\mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}^\top (\mathbf{I}_n - \mathbf{H})\mathbf{Y}] = (n - p)\sigma^2$$

Hence an unbiased estimator of $\sigma^2$ is

$$\widehat{\sigma}^2 = \frac{\mathrm{SS_{Res}}}{n - p} = \mathrm{MS_{Res}}$$

with $p = 2$.

Using similar methods

$$\mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\mathrm{SS}_{\mathrm{R}}|\mathbf{X}] = \mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}^{\top}(\mathbf{H} - \mathbf{H}_1)\mathbf{Y}|\mathbf{X}] = (p-1)\sigma^2 + \beta_1^2 S_{xx}$$

and

$$\mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\mathrm{SS}_{\mathrm{T}}|\mathbf{X}] = \mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}^{\top}(\mathbf{I}_n - \mathbf{H}_1)\mathbf{Y}|\mathbf{X}] = (n-1)\sigma^2 + \beta_1^2 S_{xx}$$

We have that

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\widehat{\beta}|\mathbf{X}] = \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}$$

which is estimated by

$$\widehat{\sigma}^2(\mathbf{X}^\top\mathbf{X})^{-1}.$$

The standard errors of the estimators are estimated by the square roots of the diagonal elements of this matrix; denote them by

$$\mathrm{e.s.e}(\widehat{\beta}_j) \quad j = 0, 1.$$

# Hypothesis Testing

We can formulate hypothesis tests for the parameters provided we make the normality assumption

$$\epsilon | \mathbf{X} \sim \text{Normal}(\mathbf{0}_n, \sigma^2 \mathbf{I}_n).$$

For $j = 0, 1$, to test

$$
\begin{aligned}
\text{H}_0 &: \quad \beta_j = 0 \\
\text{vs} \quad \text{H}_1 &: \quad \beta_j \neq 0
\end{aligned}
$$

we use the test statistic

$$t_j = \frac{\widehat{\beta}_j}{\text{e.s.e}(\widehat{\beta}_j)}.$$

If $H_0$ is true, we have by standard distributional results that corresponding statistic

$$T_j \sim \text{Student}(n - p)$$

with $p = 2$. We reject $H_0$ at significance level $\alpha$ if

$$|t_j| > t_{\alpha/2, n-p}$$

where $t_{\alpha, \nu}$ is the $1 - \alpha$ quantile of the Student-t distribution with $\nu$ degrees of freedom.

A $(1 - \alpha) \times 100\%$ confidence interval for $\beta_j$ is

$$\widehat{\beta_j} \pm t_{\alpha/2, n-p} \times \text{e.s.e}(\widehat{\beta_j}) \qquad j = 0, 1.$$

The $R^2$ statistic is defined by

$$R^2 = \frac{SS_R}{SS_T} = 1 - \frac{SS_{Res}}{SS_T}$$

and is a measure of the global adequacy of $x$ as a predictor of $y$.

The adjusted $R^2$ statistic is defined by

$$R^2_{Adj} = 1 - \frac{SS_{Res}/(n-p)}{SS_T/(n-1)}$$

and is a measure that acknowledges that $SS_{Res}$ decreases in expectation as $p$ increases.

Residual plots are used to assess 'local' model adequacy.

If the model assumptions are correct, then the residual plots

- $e_i$ vs $i$
- $e_i$ vs $x_{i1}$
- $e_i$ vs $\widehat{y}_i$

should be 'patternless' that is, should not exhibit systematic patterns in either mean-level or variability.

The residuals should form a horizontal 'band' around zero, with equal variability around zero everywhere.

Predictions from the model at value of $x$ are formed by using the estimated regression coefficients; at $x = x_{i1}$ observed in the sample, we have the prediction equal to the fitted value

$$\widehat{y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 x_{i1}.$$

In vector form, we have

$$\widehat{\mathbf{y}} = \mathbf{X}\widehat{\beta}.$$

At $x = x_1^{\text{new}}$, we have the prediction

$$\widehat{y}^{\text{new}} = \widehat{\beta}_0 + \widehat{\beta}_1 x_1^{\text{new}}.$$

In the random variable form we have predictions

$$\widehat{\mathbf{Y}} = \mathbf{X}\widehat{\beta} = \mathbf{H}\mathbf{Y}$$

so that

$$\mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\widehat{\mathbf{Y}}|\mathbf{X}] = \mathbf{X}\beta$$

and

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\widehat{\mathbf{Y}}|\mathbf{X}] = \sigma^2\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top = \sigma^2\mathbf{H}$$

Therefore a $(1-\alpha) \times 100\%$ **confidence interval** for the prediction at $x = x_{i1}$ is

$$\widehat{y}_i \pm t_{\alpha/2, n-p} \times \sqrt{\widehat{\sigma^2 h_{ii}}}$$

where $h_{ii}$ is the $(i,i)$th diagonal element of $\mathbf{H}$.

For a prediction at $x = x_1^{\text{new}}$, we have that

$$\mathbb{V}\text{ar}_{\mathbf{Y}|\mathbf{X}}[\widehat{Y}^{\text{new}}|\mathbf{X}] = \sigma^2 x^{\text{new}}(\mathbf{X}^\top\mathbf{X})^{-1}(x^{\text{new}})^\top = \sigma^2 h^{\text{new}}$$

and a $(1 - \alpha) \times 100\%$ **confidence interval** for the prediction at $x = x_1^{\text{new}}$ is

$$\widehat{y}^{\text{new}} \pm t_{\alpha/2, n-p} \times \sqrt{\widehat{\sigma}^2 h^{\text{new}}}$$

A **prediction interval** at $x = x_1^{\text{new}}$ incorporates the random variation that is present in the observations. Let

$$\widehat{Y}_{\text{O}}^{\text{new}} = \widehat{Y}^{\text{new}} + \epsilon^{\text{new}}$$

where $\epsilon^{\text{new}}$ is a zero mean, variance $\sigma^2$ random residual error, independent of all other random quantities. Then

$$\begin{aligned}
\mathbb{V}\text{ar}_{\mathbf{Y}|\mathbf{X}}[\widehat{Y}_{\text{O}}^{\text{new}}|\mathbf{X}] &= \mathbb{V}\text{ar}_{\mathbf{Y}|\mathbf{X}}[\widehat{Y}^{\text{new}}|\mathbf{X}] + \mathbb{V}\text{ar}_{\mathbf{Y}|\mathbf{X}}[\epsilon^{\text{new}}|\mathbf{X}] \\
&= \sigma^2 h^{\text{new}} + \sigma^2 \\
&= \sigma^2(1 + h^{\text{new}}).
\end{aligned}$$

Thus a $(1 - \alpha) \times 100\%$ **prediction interval** for the prediction at $x = x_1^{\text{new}}$ is

$$\widehat{y}^{\text{new}} \pm t_{\alpha/2, n-p} \times \sqrt{\widehat{\sigma}^2(1 + h^{\text{new}})}$$

The sums-of-squares decomposition

$$SS_T = SS_{Res} + SS_R$$

forms the basic component of the Analysis of Variance (ANOVA) as it describes how overall observed variability in response $y$ ($SS_T$) is decomposed into

- a component corresponding to the residual errors ($SS_{Res}$) and
- a component corresponding to the regression ($SS_R$).

Under the assumption of Normality of residual errors,

$$\epsilon | \mathbf{X} \sim \text{Normal}(\mathbf{0}_n, \sigma^2 \mathbf{I}_n),$$

and the hypothesis that $\beta_1 = \mathbf{0}$, we have the result that for the sums-of-squares random variables

$$
\begin{aligned}
\frac{\text{SS}_{\text{T}}}{\sigma^2} &= \frac{\mathbf{Y}^\top (\mathbf{I}_n - \mathbf{H}_1) \mathbf{Y}}{\sigma^2} &\sim\; \chi^2_{n-1} \\
\frac{\text{SS}_{\text{Res}}}{\sigma^2} &= \frac{\mathbf{Y}^\top (\mathbf{I}_n - \mathbf{H}) \mathbf{Y}}{\sigma^2} &\sim\; \chi^2_{n-p} \\
\frac{\text{SS}_{\text{R}}}{\sigma^2} &= \frac{\mathbf{Y}^\top (\mathbf{H} - \mathbf{H}_1) \mathbf{Y}}{\sigma^2} &\sim\; \chi^2_{p-1}
\end{aligned}
$$

with $\text{SS}_{\text{Res}}$ and $\text{SS}_{\text{R}}$ independent.

## The Analysis of Variance (cont.)

Consequently we can show that under the hypothesis, the random variable

$$F = \frac{SS_R/(p-1)}{SS_{Res}/(n-p)}$$

has a Fisher-F distribution with $p-1$ and $n-p$ degrees of freedom

$$F \sim \text{Fisher}(p-1, n-p).$$

We can construct a test of $H_0 : \beta_1 = 0$ based on this result: we reject $H_0$ at significance level $\alpha$ if

$$F > F_{\alpha, p-1, n-p}$$

where $F_{\alpha, \nu_1, \nu_2}$ is the $(1-\alpha)$ quantile of the Fisher-F distribution with $\nu_1$ and $\nu_2$ degrees of freedom

This test is equivalent to the test of $H_0 : \beta_1 = 0$ based on the $t$-statistic; we have that

$$t_1^2 = \left\{ \frac{\widehat{\beta}_1}{\text{e.s.e}(\widehat{\beta}_1)} \right\}^2 = F$$

and the two-tailed test based on $t_1$ is equivalent to the one-tailed test based on $F$.

## The Analysis of Variance (cont.)

The ANOVA table arranges the requires information in tabular form:

| Source | SS | df | MS | F |
|---|---|---|---|---|
| Regression | $SS_R$ | $p - 1$ | $SS_R/(p-1)$ | $F$ |
| Residual | $SS_{Res}$ | $n - p$ | $SS_{Res}/(n-p)$ | |
| Total | $SS_T$ | $n - 1$ | | |

# Maximum Likelihood Estimation

Under the assumption of Normality of residual errors,

$$\epsilon | \mathbf{X} \sim \text{Normal}(\mathbf{0}_n, \sigma^2 \mathbf{I}_n),$$

so that

$$\mathbf{Y} | \mathbf{X} \sim \text{Normal}(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n),$$

we may consider using a maximum likelihood (ML) procedure to estimate the model parameters $(\beta, \sigma^2)$.

The likelihood function for data $\mathbf{y}$ is

$$L(\beta, \sigma^2) = \left( \frac{1}{2\pi\sigma^2} \right)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \right\}$$

The log-likelihood is

$$
\begin{aligned}
\ell(\beta, \sigma^2) &= -\frac{n}{2}\log\sigma^2 - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^\top(\mathbf{y} - \mathbf{X}\beta) + \text{constant} \\
&= -\frac{n}{2}\log\sigma^2 - \frac{1}{2\sigma^2}S(\beta) + \text{constant}.
\end{aligned}
$$

We seek to maximize this function with respect to $\beta$ and $\sigma^2$.

It is evident that, in terms of $\beta$, the maximum value of $\ell(\beta, \sigma^2)$ is attained (for any $\sigma^2$) when $S(\beta)$ is minimized. Thus the maximum likelihood estimate of $\beta$ is identical to the least squares estimate.

To maximize over $\sigma^2$, note that

$$\frac{\partial \ell(\beta, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} S(\beta).$$

Equating to zero, and setting $\beta = \widehat{\beta}$, we have the solution

$$
\begin{aligned}
\widehat{\sigma}^2_{\text{ML}} &= \frac{1}{n} S(\widehat{\beta}) = \frac{1}{n}(\mathbf{y} - \mathbf{X}\widehat{\beta})^\top (\mathbf{y} - \mathbf{X}\widehat{\beta}) \\
&= \frac{1}{n} \sum_{i=1}^{n} (y_i - \widehat{y}_i)^2 \\
&= \frac{\text{SS}_{\text{Res}}}{n}.
\end{aligned}
$$

We may decide to treat the predictor $x$ as a random variable, and make a bivariate Normal distribution assumption for the data pairs $(x_i, y_i), i = 1, \ldots, n$.

We specify that

$$\left[ \begin{array}{c} X \\ Y \end{array} \right] \sim \text{Normal} \left( \left[ \begin{array}{c} \mu_X \\ \mu_Y \end{array} \right], \left[ \begin{array}{cc} \sigma_X^2 & \sigma_{XY} \\ \sigma_{XY} & \sigma_Y^2 \end{array} \right] \right)$$

or

$$\left[ \begin{array}{c} X \\ Y \end{array} \right] \sim \text{Normal} \left( \mu, \Sigma \right)$$

Writing

$$\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

the joint density is

$$f_{X,Y}(x,y) = \frac{1}{2\pi} \frac{1}{\sigma_X \sigma_Y \sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \right.$$

$$\left. \left[ \left( \frac{x - \mu_X}{\sigma_X} \right)^2 - \frac{2\rho(x - \mu_X)(y - \mu_Y)}{\sigma_X \sigma_Y} + \left( \frac{y - \mu_Y}{\sigma_Y} \right)^2 \right] \right\}$$

- $\sigma_{XY}$ is the covariance parameter
- $\rho$ is the correlation parameter.

We may factorize the joint density

$$f_{X,Y}(x,y) = f_{Y|X}(y|x)f_X(x)$$

where

$$X \sim \text{Normal}(\mu_X, \sigma_X^2)$$

and

$$Y|X = x \sim \text{Normal}\left(\mu_Y + \frac{\sigma_{XY}}{\sigma_X^2}(x - \mu_X), \sigma_Y^2 - \frac{\sigma_{XY}^2}{\sigma_X^2}\right)$$

or equivalently

$$Y|X = x \sim \text{Normal}\left(\mu_Y + \rho\frac{\sigma_Y}{\sigma_X}(x - \mu_X), \sigma_Y^2(1 - \rho^2)\right)$$

Equating the conditional expectation of $Y$ given $X = x$

$$\mathbb{E}_{Y|X}[Y|x] = \mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(x - \mu_X)$$

with a simple linear regression

$$\mathbb{E}_{Y|X}[Y|x] = \beta_0 + \beta_1 x$$

we identify

$$\begin{aligned}
\beta_0 &= \mu_Y - \mu_X \rho \frac{\sigma_Y}{\sigma_X} \\
\beta_1 &= \rho \frac{\sigma_Y}{\sigma_X}
\end{aligned}$$

The sample correlation is

$$r = \frac{\sum\limits_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2 \sum\limits_{i=1}^{n}(y_i - \bar{y})^2}} = \frac{S_{xy}}{\sqrt{S_{xx}SS_T}}$$

so that

$$\widehat{\beta}_1 = \left(\frac{SS_T}{S_{xx}}\right)^{1/2} r$$

and

$$r^2 = \frac{SS_R}{SS_T} = R^2.$$

# Multiple Linear Regression

The simple linear regression model extends to include multiple predictors

$$\mathbb{E}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}|\mathbf{X}] = \mathbf{X}\beta$$

in a straightforward way.

- Least squares estimation, inference, testing prediction etc. proceeds in precisely the same way.
- Model checking using residuals and $R^2$ also proceeds as for simple linear regression.
- $F$-testing procedures for assessing the utility of including all the predictors proceed as before, and there is a natural extension to the use of ANOVA tables in regression.

However, now we may also consider individual *t*-tests for the coefficients, from a single fit. Furthermore we may consider simultaneous tests for collections of parameters using the *F*-test.

- We split $\mathbf{X} = [\mathbf{X}^{(1)}\ \mathbf{X}^{(2)}]$ and $\beta = [\beta^{(1)}\ \beta^{(2)}]$ and consider the model
$$\mathbf{Y} = \mathbf{X}\beta + \epsilon = \mathbf{X}^{(1)}\beta^{(1)} + \mathbf{X}^{(2)}\beta^{(2)} + \epsilon.$$

- We compare a simpler model with a more complex model, where the simpler model is obtained by hypothesizing that some $\beta$s are zero, that is

$$\text{Full Model} \quad : \quad \mathbf{Y} = \mathbf{X}\beta + \epsilon$$

$$\text{Reduced Model} \quad : \quad \mathbf{Y} = \mathbf{X}^{(1)}\beta^{(1)} + \epsilon$$

that is, under the reduced model we assume $\beta^{(2)} = \mathbf{0}_r$.

In general the parameter estimates arising from the two models will be **different**; that is

- $\beta^{(1)}$ estimates from the Full model will in general not be equal to $\beta^{(1)}$ estimates from the Reduced model, as in the Full model

$$\widehat{\beta} = \left[ \begin{array}{c} \widehat{\beta}^{(1)} \\ \widehat{\beta}^{(2)} \end{array} \right] = \left[ \begin{array}{cc} \{\mathbf{X}^{(1)}\}^{\top}\mathbf{X}^{(1)} & \{\mathbf{X}^{(1)}\}^{\top}\mathbf{X}^{(2)} \\ \{\mathbf{X}^{(2)}\}^{\top}\mathbf{X}^{(1)} & \{\mathbf{X}^{(2)}\}^{\top}\mathbf{X}^{(2)} \end{array} \right]^{-1} \left[ \begin{array}{c} \{\mathbf{X}^{(1)}\}^{\top}\mathbf{y} \\ \{\mathbf{X}^{(2)}\}^{\top}\mathbf{y} \end{array} \right]$$

and in the reduced model

$$\widehat{\beta}^{(1)} = (\{\mathbf{X}^{(1)}\}^{\top}\mathbf{X}^{(1)})^{-1}\{\mathbf{X}^{(1)}\}^{\top}\mathbf{y}.$$

We modify the earlier sums of squares decomposition

$$\sum_{i=1}^{n}(y_i - \bar{y})^2 = \sum_{i=1}^{n}(y_i - \widehat{y}_i)^2 + \sum_{i=1}^{n}(\widehat{y}_i - \bar{y})^2$$

$$SS_T = SS_{Res} + SS_R$$

to

$$\sum_{i=1}^{n}y_i^2 = \sum_{i=1}^{n}(y_i - \widehat{y}_i)^2 + \sum_{i=1}^{n}\widehat{y}_i^2$$

$$\overline{SS}_T = SS_{Res} + \overline{SS}_R$$

Note that

$$\mathrm{SS_T} = \overline{\mathrm{SS}}_\mathrm{T} - n\{\bar{y}\}^2$$

and

$$\mathrm{SS_R} = \overline{\mathrm{SS}}_\mathrm{R} - n\{\bar{y}\}^2.$$

If $\overline{SS}_R(\beta)$ and $\overline{SS}_R(\beta^{(1)})$ denote the regression sums of squares from the Full and Reduced models respectively, the *extra sum of squares* due to $\beta^{(2)}$ in the presence of $\beta^{(1)}$ is

$$\overline{SS}_R(\beta^{(2)}|\beta^{(1)}) = \overline{SS}_R(\beta) - \overline{SS}_R(\beta^{(1)}).$$

This facilitates the *F*-test of the null hypothesis

$$H_0 \: : \: \beta^{(2)} = \mathbf{0}_r$$

that is, that the Reduced model is an adequate simplification of the Full model.

We perform a partial $F$-test using

$$F = \frac{(\overline{SS}_R(\beta^{(2)}|\beta^{(1)}))/r}{MS_{Res}(\beta)} = \frac{(SS_{Res}(\beta^{(1)}) - SS_{Res}(\beta))/r}{SS_{Res}(\beta)/(n-p)}$$

which is distributed as

$$\text{Fisher}(r, n-p)$$

if $H_0$ is true.

The test may be used sequentially: for example, if parameter $\beta = (\beta_0, \beta_1, \beta_2, \beta_3)^\top$, then we have

$$\overline{SS}_R(\beta_0, \beta_1, \beta_2, \beta_3) = \overline{SS}_R(\beta_0) + \overline{SS}_R(\beta_1|\beta_0) + \overline{SS}_R(\beta_2|\beta_0, \beta_1) + \overline{SS}_R(\beta_3|\beta_0, \beta_1, \beta_2)$$

and also

$$\overline{SS}_R(\beta_1, \beta_2, \beta_3|\beta_0) = \overline{SS}_R(\beta_1|\beta_0) + \overline{SS}_R(\beta_2|\beta_0, \beta_1) + \overline{SS}_R(\beta_3|\beta_0, \beta_1, \beta_2)$$

We also have

$$\overline{SS}_R(\beta_1, \beta_2, \beta_3|\beta_0) \equiv SS_R(\beta_1, \beta_2, \beta_3).$$

This latter decomposition allows us to test

- whether $X_1$ is worth including in the model, then
- whether $X_2$ is worth including, when $X_1$ is already included, then
- whether $X_3$ is worth including, when $X_1$ and $X_2$ are already included.

**Note:** If the $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ blocks are **orthogonal**

$$\{\mathbf{X}^{(1)}\}^{\top}\mathbf{X}^{(2)} = \mathbf{0}_{p-r,r}$$

then

- $\beta$ estimates from Full and Reduced model are equal;
- we have the identity

$$\overline{\mathrm{SS}}_{\mathrm{R}}(\beta^{(1)}|\beta^{(2)}) = \overline{\mathrm{SS}}_{\mathrm{R}}(\beta^{(1)}).$$

# The General Linear Hypothesis

The general linear hypothesis specifies a more general type of constraint on the model, namely

$$\mathrm{H_0} \ : \ \mathbf{A}\beta = \mathbf{0}$$

for some ($m \times p$) matrix $\mathbf{A}$. If the hypothesis places $r$ linearly independent constraints on $\beta$ to make the Reduced model, and if

$$\mathrm{SS_H} = \mathrm{SS_{Res}}(\mathrm{Reduced}) - \mathrm{SS_{Res}}(\mathrm{Full})$$

then the test statistic

$$F = \frac{\mathrm{SS_H}/r}{\mathrm{SS_{Res}}(\mathrm{Full})/(n - p)}$$

is distributed as Fisher($r, n - p$) if $\mathrm{H_0}$ is true.

When testing each of the $\beta$s in turn using *t*-tests, we should account for multiple testing by controlling the *familywise Type I error rate* using *multiple testing corrections*:

- Bonferroni correction.

Confidence ellipsoids allow for simultaneous confidence statements to be made about multiple $\beta$ parameters; the region

$$\mathbf{b} : \frac{(\widehat{\beta} - \mathbf{b})^\top (\mathbf{X}^\top \mathbf{X})(\widehat{\beta} - \mathbf{b})}{p \mathrm{MS}_{\mathrm{Res}}(\beta)} \leq F_{\alpha, p, n-p}$$

defines a region in $\mathbb{R}^p$ that exhibits $(1 - \alpha) \times 100\%$ confidence.

*Multicollinearity* corresponds to dependence/correlation amongst the predictors:

- can lead to inflation of the variance of estimators if present;
- can be measured by variance inflation factors;
- equivalent to $R^2$ measures constructed for the predictors;
- can be computed from the correlation matrix from the predictors.

# Special Types of Predictors

- polynomial terms;
- factor predictors: discrete predictors measured on a nominal (non-numerical) scale;
- interactions: an interaction is a modification of the effect of one predictor in the presence of another.
- higher-order interactions.

Important issues include

- how to represent factor predictors using indicator functions;
- how to count parameters;
- the interpretation of interactions;
- removing the intercept.

Our goal is to find the simplest possible model that adequately explains the observed response data.

- Forward selection;
- Backward elimination;
- Stepwise elimination;

# Model Selection Criteria

- $R^2$-based methods;
- largest $R^2_{\text{Adj}}$ equivalent to minimum $\text{MS}_{\text{Res}}$;
- Mallows's $C_p$;
- AIC;
- BIC;

- **Over-simplified model:** if the chosen model omits important predictors, then
    - the resulting estimates will be biased, have **lower** variance, and can have **higher** mean-squared error, depending on the magnitude of the effect of the omitted predictors;
    - the estimate of $\sigma^2$ will be too high;
    - predictions will have higher mean squared error.

  If the omitted predictors are orthogonal to the included ones, then the bias is removed.

- **Over-complex model:** if the chosen model includes unnecessary predictors, then
    - the resulting estimates will be unbiased, have **higher** variance, and **higher** mean-squared error;
    - the estimate of $\sigma^2$ will be unbiased;
    - predictions will have bias, higher variance and higher mean squared error.

- PRESS residuals/statistic;
- leave-one-out/deletion residuals;
- leverage;
- influence in estimation and prediction;
- outliers;
- influence diagnostics.

# Generalizing Least Squares

We may relax the variance assumption

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}|\mathbf{X}] = \sigma^2 \mathbf{I}_n$$

to

$$\mathbb{V}\mathrm{ar}_{\mathbf{Y}|\mathbf{X}}[\mathbf{Y}|\mathbf{X}] = \sigma^2 \mathbf{V}$$

where $\mathbf{V}$ is a square, symmetric, non-singular matrix. This allows for

- unequal variances in the conditional response distribution;
- correlation amongst the residual errors $\epsilon$.

This leads to the amended least squares objective function

$$(\mathbf{y} - \mathbf{X}\beta)^\top \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta)$$

This leads to the estimates

$$\widehat{\beta} = (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^\top \mathbf{X}^\top \mathbf{V}^{-1} \mathbf{y}$$

and the properties of the estimators, predictions etc. follow in a straightforward fashion.

A special case is *weighted least squares*, where

$$\mathbf{V} = \operatorname{diag}(1/w_1, 1/w_2, \ldots, 1/w_n)$$

which accommodates unequal variances.

We may use transformations of the response $y_i$ to a different form to make the linear regression model assumptions more appropriate:

- log;
- square root;
- Box-Cox transform.