

CSCI 3022-002 Intro to Data Science

Multiple Linear Regression

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Workflow for SLR so far:

1. Plot the data as a scatter plot
 - 1.1 Does **linearity** seem appropriate?
 - 1.2 Calculate $\hat{\beta}_0, \hat{\beta}_1$ and overlay the best-fit line $y = \beta_0 + \beta_1 X$.
2. Consider assumptions of SLR:
 - 2.1 Plot a histogram of the residuals: are they **normal**?
 - 2.2 Plot the residuals against x : are they changing?
3. Perform inference (on β s or on values of $Y|X$)

Announcements and Reminders

- ▶ Short HW for next week.
- ▶ NB day Friday.

Where we at?

Last time we talked about multiple linear regression. It's like simple linear regression, except now we can attempt to predict Y with a variety of things, including both different *features/predictors* X as well as transformations and augmentations of the original variables, like using x^2 .

Process: try to fix "problems" with assumptions. This means:

1. Plot the linear model
2. See if some predictors are redundant
3. Plot residuals of linear model, check for **normality, independence, structure**.
4. Hit model with a math-shaped stick to fix these problems.

Covariance

When two random variables X and Y are not independent, it is frequently of interest to assess how strongly they are related to one another.

Definition: *Covariance:*

The covariance between two rv's X and Y is defined as:

$$E[\underbrace{(X - \mu_X)}_{\text{X versus its mean}} \underbrace{(Y - \mu_Y)}_{\text{Y versus its mean}}]$$

If both variables tend to deviate in the same direction (both go above their means or below their means at the same time), then the covariance will be positive.

If the opposite is true, the covariance will be negative.

If X and Y are not strongly related, the covariance will be near 0.

Definition: *Correlation*

The *correlation* coefficient of X and Y , denoted by $\text{Corr}[X, Y]$ or just ρ , is the *unitless* measure of covariance defined by:

(-1) to (1) measure for X-Y relationship?

Multiple Linear Regression

Definition: Multiple Linear Regression

The multiple regression model is one where we allow each data point to have multiple characteristics (features/predictors) that we use to predict y . So for each data point we have p different X 's to predict y :

$$Y_i = \beta_0 + \beta_1 X_{1,i} + \beta_2 X_{2,i} + \cdots + \beta_p X_{p,i}$$

In matrix form:

$$\underline{Y} := [Y_1, Y_2, \dots, Y_n]^T$$

$$\underline{\beta} := [\beta_0, \beta_1, \dots, \beta_p]^T$$

$$X := \begin{bmatrix} 1 & X_{1,1} & \dots & X_{1,p} \\ 1 & X_{2,1} & \dots & X_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n,1} & \dots & X_{n,p} \end{bmatrix}$$

X data frame: design matrix

row: observation

columns: features

1st column: all 1's for β_0 intercept

So our model is $\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$. \underline{X} is called the **design** matrix. 4 (familiar) assumptions:

Multiple Linear Regression Estimators

Our estimators for the regression coefficients β rely on these assumptions, and look similar to those in SLR.

Multiple Linear Regression Estimators

Our estimators for the regression coefficients β rely on these assumptions, and look similar to those in SLR.

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix} = (X^T X)^{-1} X^T \underline{Y}$$

- intercept
partial regression coefficients
 \approx slopes

Multiple Linear Regression Estimators

$$X: \begin{pmatrix} 1 & 2 & 20 \\ 1 & 4 & 18 \end{pmatrix}$$

column · column

$$X^T X: \begin{pmatrix} 2 & 6 & 38 \\ 6 & 20 & 118 \\ 38 & 118 & 400 \end{pmatrix}$$

Our estimators for the regression coefficients β rely on these assumptions, and look similar to those in SLR.

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{pmatrix} = (X^T X)^{-1} X^T \underline{Y}$$

$$X^T X: \begin{bmatrix} 2 & 6 \\ 6 & 20 \end{bmatrix}$$

1: 2+1: 4
2: 2+4: 2

Each entry is
~ Cov(column_i, X column_j)

The $(X^T X)^{-1}$ bit corresponds to the $1/\sum (X_i - \bar{X})^2$ part from before, where the $X^T \underline{Y}$ part corresponds roughly to a covariance between X and Y .

Collinearity

The $(X^T X)^{-1}$ term in our regression coefficient errors is very similar to the "spread of X " term in the SLR coefficients. This time, however, it's a little nastier: it's the spread of X across *all* p dimensions of the predictors. Example:

Suppose we have roughly linear data, and we decide to fit the data with the model

$$y = \beta_0 + \beta_1 x + \beta_2 x^{1.000001} + \varepsilon$$

pred 1

pred 2

data
design

$$\begin{pmatrix} 1 & 2 & 2.00001 \\ 1 & 3 & 3.00001 \\ 1 & 5 & 5.00001 \\ \vdots & \vdots & \vdots \end{pmatrix}$$

Collinearity

The predictors $x_1 = x$ and $x_2 = x^{1.000001}$ are very highly related. In particular, if the *true* value of the regression line is $y = x$, the following functions are visually nearly identical:

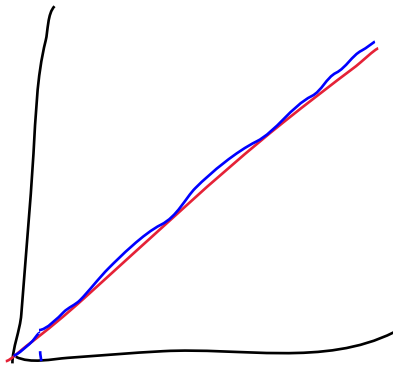
Collinearity

$$y = \beta_0 + \beta_1 x + \beta_2 x^{1.000001}$$

Handwritten annotations: $\beta_1 \rightarrow 0$ and $\beta_2 \rightarrow 1$ in blue ink. β_1 and β_2 are circled in red, with red arrows pointing to them from below.

The predictors $x_1 = x$ and $x_2 = x^{1.000001}$ are very highly related. In particular, if the *true* value of the regression line is $y = x$, the following functions are visually nearly identical:

1. $y = x$
2. $y = x^{1.000001}$



Collinearity

The predictors $x_1 = x$ and $x_2 = x^{1.000001}$ are very highly related. In particular, if the *true* value of the regression line is $y = x$, the following functions are visually nearly identical:

1. $y = x$
2. $y = x^{1.000001}$
3. $y = .5x + .5x^{1.000001}$
4. $y = 2x^{1.000001} - x$

Collinearity

The predictors $x_1 = x$ and $x_2 = x^{1.000001}$ are very highly related. In particular, if the *true* value of the regression line is $y = x$, the following functions are visually nearly identical:

1. $y = x$ 1,0
2. $y = x^{1.000001}$
3. $y = .5x + .5x^{1.000001}$
4. $y = 2x^{1.000001} - x$
5. $y = \underbrace{10^6 x - 999999 x^{1.000001}}$
6. $y = ax + bx^{1.000001}; \quad \forall a + b = 1.$

Collinearity

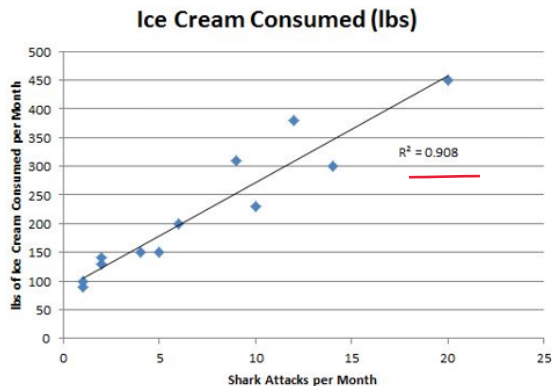
slope/coefficient for X relates
to the slope for $X^{(1.000000)}$

This is scary! The distribution of $\underline{\beta}$ has its own covariance, because the best choices for β_1 and β_2 may depend on each other. In the prior example, they would have a negative correlation of $-1!$.

In general, the interactions between coefficients is a function of the *linear independence* of the columns of the X matrix. In other words, we get a lot of negative effects if one predictor is describing one of the same things that we already have!

Correlations:

A SLR analysis of shark attacks vs ice cream sales at a Southern California beach indicates that there is a strong relationship between the two.



MLR:

→ season
of people at beach

Suppose we included both **temperature** and shark attacks as features in our model of ice cream sales. What would happen? Which one should we probably exclude, and why?

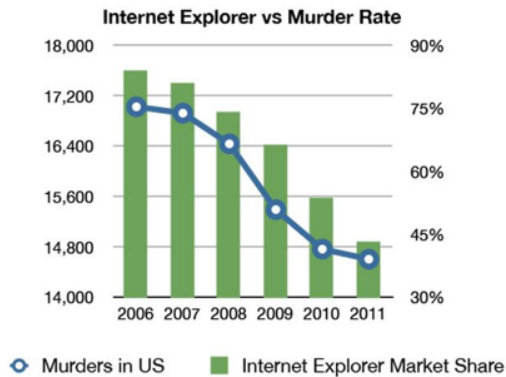
$X^T X$ is large: the columns:

shark attacks

temperature

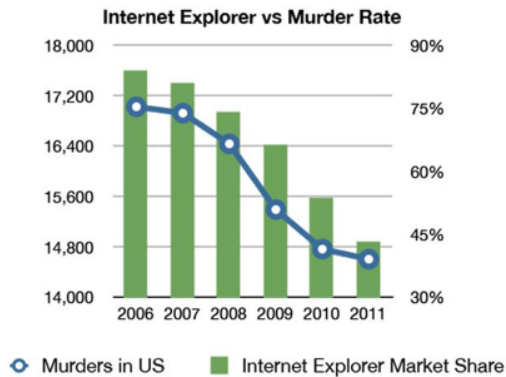
are related, so our estimates for both become confounded.

Correlations:



Why is this correlated?

Correlations:



Why is this correlated?



MLR Variable Selection

"(on the model)": Set slopes β
get CI & Hypotheses on β
check assumptions

The selection process of variables for a multiple linear regression is complicated! Typically, we begin by *running the model with all parameters included*. Then - **one at a time** - we discard the least useful columns of the X matrix. The following are *all* possible criteria for excluding a column.

1. Check your X -values for *redundant*, non-independent information. Discard an offending column, then repeat. The measure for this is a variance inflation factor. (high VIF = bad).
2. The estimator $\hat{\beta}_i$ for predictor column i is *not significantly different from zero* per a t -test. This means that in the presence of the other features, it's not really helping us predict y !
 H_0 : predictor i is useful : $\beta_i = 0$ check t_{stat} for β_i .
3. The model with predictor i has lower *adjusted R^2* than the model without it.
4. The variance captured by the model with predictor i has (~~lower adjusted R^2~~) than the model without it.
~~SSE~~
 test on Variance: 1) not a mean \Rightarrow no CLT
 2) Bootstrap
 3) new test!
 significantly less SSE

MLR Variable Selection

Typically we always include the first criteria, and then choose one of the other measures for improvement.

1. Discard x with high variance inflation factors.
2. Individual coefficient t -test significance.
3. Lower *adjusted* R^2 without a predictor.
4. The *variance captured* by the model with predictor i has lower *adjusted* R^2 than the model without it.

Sometimes you may even want to include non-significant predictors. Why? This may help the model be more *predictive*, even if we can't statistically point to which factors matter the most. A smaller model does a better job at **suggesting** causal relationships, but we never actually get true causality, so sometimes a more accurate prediction is all we want!

Special Cases: Variance

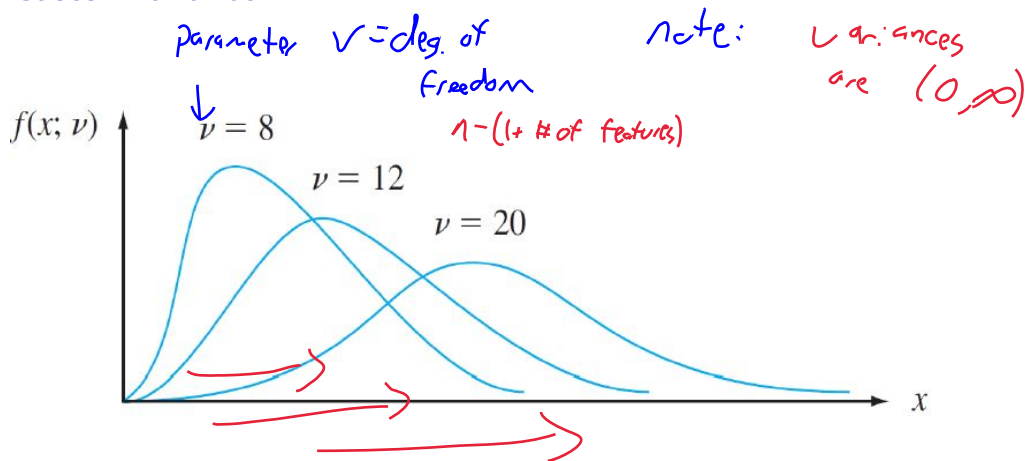
Definition: Chi-Squared: the distribution for variances of normal distributions

Let ν be a positive integer. The random variable X has a chi-squared distribution with parameter ν if the pdf of X is:

$$f(x) = \begin{cases} \frac{1}{2^{\nu/2}\Gamma(\nu/2)} x^{(\nu/2)-1} e^{-x/2} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

The parameter ν is called the number of degrees of freedom (df) of X . The symbol χ^2 is often used in place of “chi-squared.”

Special Cases: Variance



Special Cases: Variance

Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$. Then the random variable:

has a chi-squared (____) probability distribution with $n - 1$ df.

(In this class, we don't consider the case where the data is not normally distributed.)

Special Cases: Variance

Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$. Then the random variable:

$$\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} = \frac{(n-1)s^2}{\sigma^2}$$

normalize s^2
versus true σ^2

has a chi-squared (χ^2) probability distribution with $n - 1$ df.

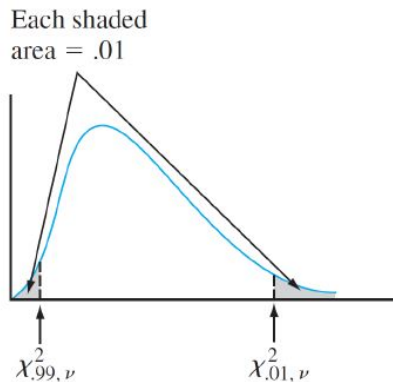
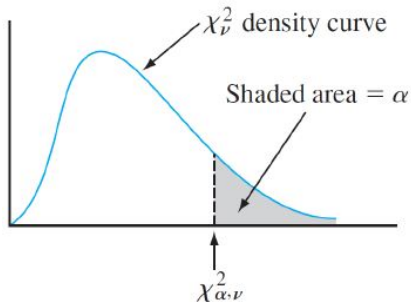
(In this class, we don't consider the case where the data is not normally distributed.)

stats. chi-squared: if variance = 1,
if variance = 14,

give me 95% CI on variance

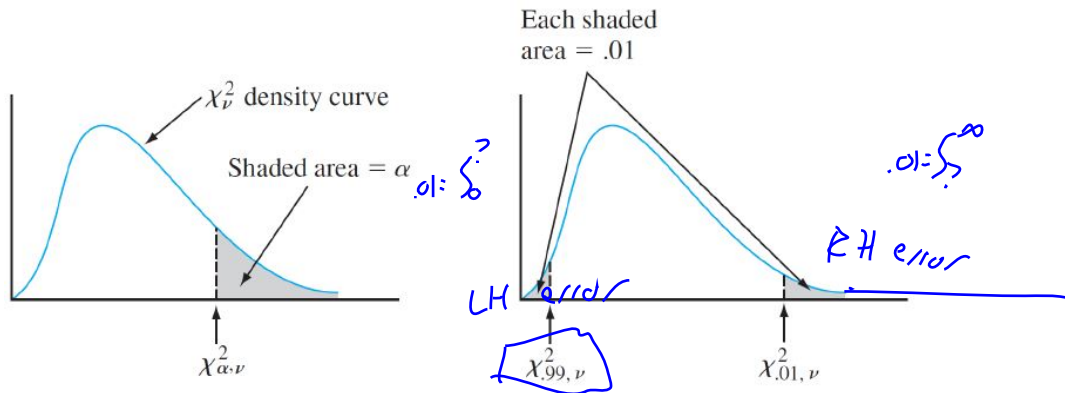
Special Cases: Variance

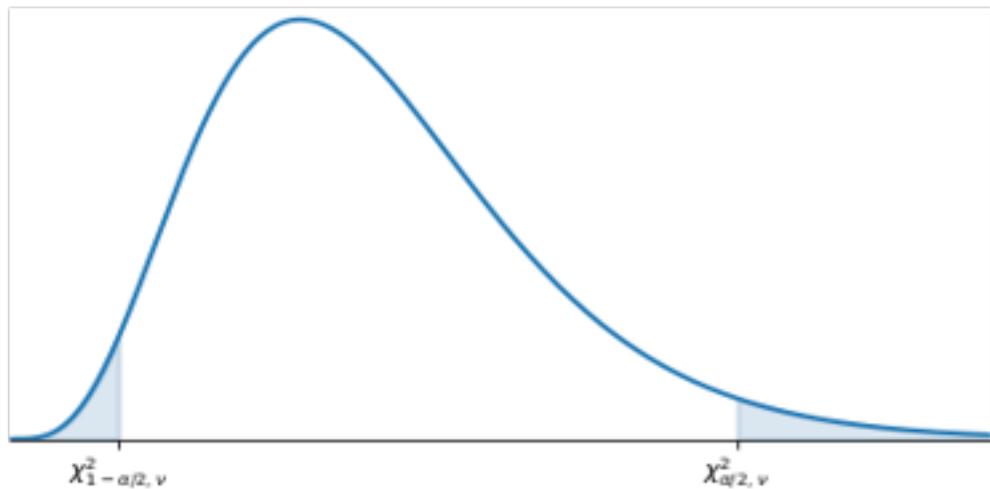
The chi-squared distribution is not symmetric, so these tables and functions contain values of _____ both for near 0 and 1.



Special Cases: Variance

The chi-squared distribution is not symmetric, so these tables and functions contain values of χ^2_{α} both for near 0 and 1.



Two tailed χ^2 

Special Cases: Variance

As a consequence:

$$1 - \alpha = P \left(\chi_{1-\alpha/2, n-1}^2 \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi_{\alpha/2, n-1}^2 \right)$$

Or, equivalently:

Thus we have a confidence interval for the variance. Taking square roots gives a CI for the standard deviation.

$$\text{Recall: } \left(\sqrt{\frac{1}{n}} \leq u \leq \sqrt{\frac{1}{n}} \right)$$

Special Cases: Variance

As a consequence:

$$\begin{aligned} 1 - \alpha &= P \left(\chi_{1-\alpha/2, n-1}^2 \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi_{\alpha/2, n-1}^2 \right) \\ &= P(1/\chi_{1-\alpha/2, n-1}^2 \geq \frac{\sigma^2}{(n-1)s^2} \geq 1/\chi_{\alpha/2, n-1}^2) \end{aligned}$$

Or, equivalently:

$$\left(\frac{(n-1)S^2}{\chi_{1-\alpha/2, n-1}^2}, \frac{(n-1)S^2}{\chi_{\alpha/2, n-1}^2} \right)$$

is a 100%(1 - α) CI for σ^2 .

Thus we have a confidence interval for the variance. Taking square roots gives a CI for the standard deviation.

A CI on Variance

Example: A large candy manufacturer produces packages of candy targeted to weigh 52g. The weight of the packages of candy is known to be normally distributed, but a QC engineer is concerned that the variation is too large. She selected 10 bags at random and weights them, for a sample variance of $4.2g^2$. Find a 95% CI for the variance and a 95% CI for the SD.

$$\alpha = .05, \quad \alpha/2 = .025 \quad n = 10 \quad s^2 = 4.2$$

A CI on Variance

Example: A large candy manufacturer produces packages of candy targeted to weigh 52g. The weight of the packages of candy is known to be normally distributed, but a QC engineer is concerned that the variation is too large. She selected 10 bags at random and weights them, for a sample variance of $4.2g^2$. Find a 95% CI for the variance and a 95% CI for the SD.

$$\alpha = .05, \quad \alpha/2 = .025 \quad n = 10 \quad s^2 = 4.2$$

$$\chi^2_{1-\alpha/2, n-1} = \chi^2_{.975, 9} = \text{stats.chi2.ppf}(0.025, 9) = 2.70$$

$$\chi^2_{\alpha/2, n-1} = \chi^2_{.025, 9} = \text{stats.chi2.ppf}(0.975, 9) = 19.02$$

A CI on Variance

Example: A large candy manufacturer produces packages of candy targeted to weigh 52g. The weight of the packages of candy is known to be normally distributed, but a QC engineer is concerned that the variation is too large. She selected 10 bags at random and weights them, for a sample variance of $4.2g^2$. Find a 95% CI for the variance and a 95% CI for the SD.

$$\alpha = .05, \quad \alpha/2 = .025 \quad n = 10 \quad s^2 = 4.2$$

$$\chi^2_{1-\alpha/2, n-1} = \chi^2_{.975, 9} = \text{stats.chi2.ppf}(0.025, 9) = 2.70$$

$$\chi^2_{\alpha/2, n-1} = \chi^2_{.025, 9} = \text{stats.chi2.ppf}(0.975, 9) = 19.02$$

$$\frac{(10-1)4.2}{19.02} < \sigma^2 \frac{(10-1)4.2}{2.70} \implies 1.99 < \sigma^2 < 14.0$$

$$\implies \sqrt{1.99} < \sigma < \sqrt{14.0}$$

Test for Equivalence of Variance

The F probability distribution has two parameters, denoted by ν_1 and ν_2 . The parameter ν_1 is called the numerator degrees of freedom, and ν_2 is the denominator degrees of freedom.

A random variable that has an F distribution cannot assume a negative value. The density function is complicated and will not be used explicitly, so it's not shown.

There is an important connection between an F variable and chisquared variables.

Test for Equivalence of Variance

If X_1 and X_2 are independent chi-squared rv's with ν_1 and ν_2 df, respectively, then the rv

can be shown to have an F distribution.

Recall that a chi-squared distribution was obtained by summing squared standard Normal variables (such as squared deviations for example). So a scaled ratio of two variances is a ratio of two scaled chi-squared variables.

Test for Equivalence of Variance

If X_1 and X_2 are independent chi-squared rv's with ν_1 and ν_2 df, respectively, then the rv

$$F = \frac{X_1/\nu_1}{X_2/\nu_2}$$

can be shown to have an F distribution.

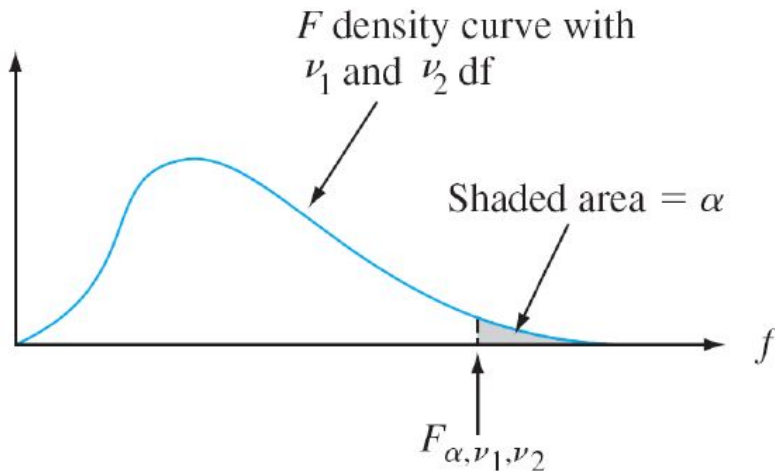
"F is the ratio (scaled) between two variances"

Recall that a chi-squared distribution was obtained by summing squared standard Normal variables (such as squared deviations for example). So a scaled ratio of two variances is a ratio of two scaled chi-squared variables.

$$\frac{\text{Var 1}}{\text{Var 2}} \approx \begin{cases} \text{large, if num is bigger} \\ 1, \text{ if } = \\ \text{small, if den is bigger} \end{cases}$$

Test for Equivalence of Variance

Figure below illustrates a typical F density function.:



Test for Equivalence of Variance

We use F_{α, ν_1, ν_2} for the value on the horizontal axis that captures of the area under the F density curve with ν_1 and ν_2 df in the upper tail.

The density curve is not symmetric, so it would seem that both upper- and lower-tail critical values must be tabulated. This is not necessary, though, because of the fact that

$$F_{1-\alpha, \nu_1, \nu_2} =$$

Test for Equivalence of Variance

We use F_{α, ν_1, ν_2} for the value on the horizontal axis that captures of the area under the F density curve with ν_1 and ν_2 df in the upper tail.

The density curve is not symmetric, so it would seem that both upper- and lower-tail critical values must be tabulated. This is not necessary, though, because of the fact that

$$F_{1-\alpha, \nu_1, \nu_2} = \frac{1}{F_{\alpha, \nu_1, \nu_2}}$$

For example, $F_{\underline{0.05}, 6, 10} = 3.22$ and $F_{\underline{0.95}, 10, 6} = 0.31 = 1/3.22$.

RH

LH

Test for Equivalence of Variance

A test procedure for hypotheses concerning the ratio σ_1^2/σ_2^2 is based on the following result.

Theorem:

Let X_1, X_2, \dots, X_m be a random sample from a normal distribution with variance σ_1^2 let Y_1, Y_2, \dots, Y_n be another random sample (independent of the X_i 's) from a normal distribution with variance σ_2^2 and let s_1^2 and s_2^2 denote the two sample variances. Then the rv

has an F distribution with $\nu_1 = m - 1$ and $\nu_2 = n - 1$.

Test for Equivalence of Variance

A test procedure for hypotheses concerning the ratio σ_1^2/σ_2^2 is based on the following result.

Theorem:

Let X_1, X_2, \dots, X_m be a random sample from a normal distribution with variance σ_1^2 let Y_1, Y_2, \dots, Y_n be another random sample (independent of the X_i 's) from a normal distribution with variance σ_2^2 and let s_1^2 and s_2^2 denote the two sample variances. Then the rv

$$F = \frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2}$$

has an F distribution with $\nu_1 = m - 1$ and $\nu_2 = n - 1$.

$$\frac{s_1^2}{s_2^2} \text{ predicts } \frac{\sigma_1^2}{\sigma_2^2}$$

Test for Equivalence of Variance

This theorem results from combining the fact that the variables $\frac{(n-1)s_2^2}{\sigma_2^2}$ and $\frac{(m-1)s_1^2}{\sigma_1^2}$ each have a chi-squared distribution with $n - 1$ and $m - 1$ df, respectively.

Because F involves a ratio rather than a difference, the test statistic is the ratio of sample variances.

The claim that $\sigma_1^2 = \sigma_2^2$ is then rejected if the ratio s_1^2/s_2^2 differs by too much from 1.

Test for Equivalence of Variance

Null hypothesis: H_0 :

Test statistic value:

<u>Alt Hypothesis</u>	<u>Rejection Region</u>
-----------------------	-------------------------

<u>p-value:</u>

Test for Equivalence of Variance

Null hypothesis: $H_0 : \sigma_1^2 = \sigma_2^2$

Test statistic value:

$$F = s_1^2 / s_2^2$$

Alt Hypothesis Rejection Region

p-value:

Test for Equivalence of Variance

Null hypothesis: $H_0 : \sigma_1^2 = \sigma_2^2$

Test statistic value:

$$F = s_1^2 / s_2^2$$

Alt Hypothesis	Rejection Region
$H_a : \sigma_1^2 > \sigma_2^2$	
$H_a : \sigma_1^2 < \sigma_2^2$	
$H_a : \sigma_1^2 \neq \sigma_2^2$	

p-value:

Test for Equivalence of Variance

Null hypothesis: $H_0 : \sigma_1^2 = \sigma_2^2$

Test statistic value:

$$F = s_1^2 / s_2^2$$

<u>Alt Hypothesis</u>	<u>Rejection Region</u>	<u>p-value:</u>
$H_a : \sigma_1^2 > \sigma_2^2$	$F_{stat} > F_{\alpha, m-1, n-1}$	
$H_a : \sigma_1^2 < \sigma_2^2$	$F_{stat} < F_{1-\alpha, m-1, n-1}$	
$H_a : \sigma_1^2 \neq \sigma_2^2$	$F_{stat} < F_{1-\alpha/2, m-1, n-1}$ OR $F_{stat} > F_{\alpha/2, m-1, n-1}$	

Test for Equivalence of Variance

Null hypothesis: $H_0 : \sigma_1^2 = \sigma_2^2$

Test statistic value:

$$F = s_1^2 / s_2^2$$

Alt Hypothesis	Rejection Region	p-value:
$H_a : \sigma_1^2 > \sigma_2^2$	$F_{stat} > F_{\alpha, m-1, n-1}$	$P(F_{m-1, n-1} > F_{stat})$
$H_a : \sigma_1^2 < \sigma_2^2$	$F_{stat} < F_{1-\alpha, m-1, n-1}$	$P(F_{m-1, n-1} < F_{stat})$
$H_a : \sigma_1^2 \neq \sigma_2^2$	$F_{stat} < F_{1-\alpha/2, m-1, n-1}$ OR $F_{stat} > F_{\alpha/2, m-1, n-1}$	(OR)

Test for Equivalence of Variance

Example: On the basis of data reported in the article “Serum Ferritin in an Elderly Population” (J. of Gerontology, 1979: 521–524), the authors concluded that the ferritin distribution in the elderly had a smaller variance than in the younger adults. (Serum ferritin is used in diagnosing iron deficiency.)

For a sample of 28 elderly men, the sample standard deviation of serum ferritin (mg/L) was $s_1 = 52.6$; for 26 young men, the sample standard deviation was $s_2 = 84.2$.

Does this data support the conclusion as applied to men? Use $\alpha = .01$.

Test for Equivalence of Variance

Example: On the basis of data reported in the article “Serum Ferritin in an Elderly Population” (J. of Gerontology, 1979: 521–524), the authors concluded that the ferritin distribution in the elderly had a smaller variance than in the younger adults. (Serum ferritin is used in diagnosing iron deficiency.)

For a sample of 28 elderly men, the sample standard deviation of serum ferritin (mg/L) was $s_1 = 52.6$; for 26 young men, the sample standard deviation was $s_2 = 84.2$.

Does this data support the conclusion as applied to men? Use $\alpha = .01$.

$$F_{27,25} = \frac{52.6}{84.2} = F_{stat}$$

Test for Equivalence of Variance

Example: On the basis of data reported in the article “Serum Ferritin in an Elderly Population” (J. of Gerontology, 1979: 521–524), the authors concluded that the ferritin distribution in the elderly had a smaller variance than in the younger adults. (Serum ferritin is used in diagnosing iron deficiency.)

For a sample of 28 elderly men, the sample standard deviation of serum ferritin (mg/L) was $s_1 = 52.6$; for 26 young men, the sample standard deviation was $s_2 = 84.2$.

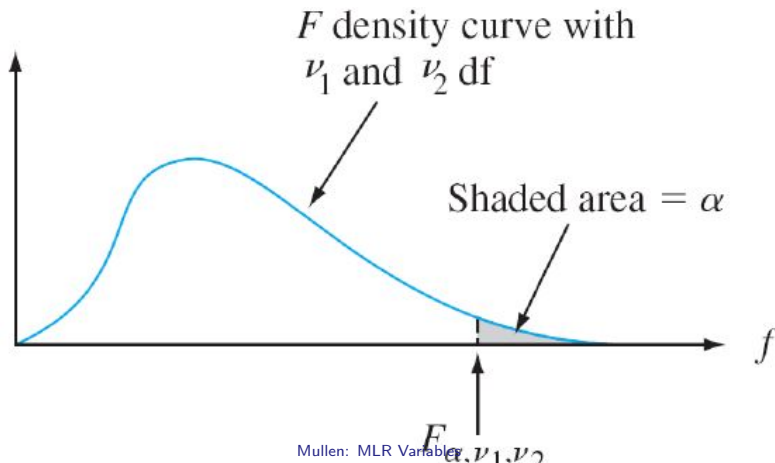
Does this data support the conclusion as applied to men? Use $\alpha = .01$.

$$F_{27,25} = \frac{52.6}{84.2} = F_{stat}$$

$$P(F_{27,25} \leq \frac{52.6}{84.2}) = \text{stats.f.cdf}(\frac{52.6}{84.2}, 27, 25) = .117 = p$$

Test for Equivalence of Variance

Recall: a typical F density function. When this thing took a value far from 1, we could conclude that the *ratio* being calculated had significantly different numerator from denominator. This is how we compared two variances.



The F-test for MLR

We use F statistics to compare variances. One way to compare linear models is to compare the variance in Y to the variance of your model: if your model is capturing a lot of the variance in Y , it's doing well!

In MLR we test the hypothesis:

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$$

Null Hyp: all features
are useless

which says that there is no useful linear relationship between y and any of the p predictors. We test against:

$$H_a : \text{any of the } B'_j\text{'s are nonzero.}$$

We could test each separately, but we would be committing the multiple comparisons fallacy. A better test is a joint test, and is based on a statistic that has an F distribution when H_0 is true.

The Full F

Null hypothesis:

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$$

Alternative hypothesis:

$$H_a : \text{at least one } \beta_j \neq 0.$$

Test statistic value:

$$F = \frac{SSR/(p+1)}{SST/(n-p+1)}$$

Rejection region for a level test: $f \geq F_{\alpha, p+1, n-(p+1)}$

The Partial F

Comparing variances also gives us another way - besides just adjusted R^2 - to compare between models.

Idea: compare the amount of variance captured by the larger model to the smaller model. If they're significantly different, we know the larger model is "adding" lots of information!

more features *less features*

As a hypothesis, this means testing that the parameters that are different between models are zero.

The Partial F

Suppose we wanted to know whether *some* subset of model parameters were zero:

Null:

Alternative:

The Partial F

Suppose we wanted to know whether *some* subset of model parameters were zero:

Full model: $\underline{Y} = \beta_0 + \beta_1 \underline{X_1} + \dots \beta_4 \underline{X_4}$

Reduced Model: $\underline{Y} = \beta_0 + \beta_2 \underline{X_2} + \beta_4 \underline{X_4}$

Null:

Alternative:

The Partial F

Suppose we wanted to know whether *some* subset of model parameters were zero:

Full model: $\underline{Y} = \beta_0 + \beta_1 \underline{X_1} + \dots \beta_4 \underline{X_4}$

Reduced Model: $\underline{Y} = \beta_0 + \beta_2 \underline{X_2} + \beta_4 \underline{X_4}$

Null: $H_0 : \beta_1 = \beta_3 = 0$

Alternative: Either/both of $\beta_1 \neq 0$ or $\beta_3 \neq 0$. Alternatively: the overall model captures significantly more variability in Y by including both β_1 and β_3 .

The Partial F

Suppose we wanted to know whether *some* subset of model parameters were zero:

Full model: $\underline{Y} = \beta_0 + \beta_1 \underline{X_1} + \dots \beta_4 \underline{X_4}$

Reduced Model: $\underline{Y} = \beta_0 + \beta_2 \underline{X_2} + \beta_4 \underline{X_4}$

Null: $H_0 : \beta_1 = \beta_3 = 0$

Alternative: Either/both of $\beta_1 \neq 0$ or $\beta_3 \neq 0$. Alternatively: the overall model captures significantly more variability in Y by including both β_1 and β_3 .

In many situations, one first builds a model containing p predictors and then wishes to know whether any of the predictors in a particular subset provide useful information about Y .

The Partial F

The test is carried out by fitting both the full and reduced models.

Because the full model contains not only the predictors of the reduced model but also some extra predictors, it should fit the data at least as well as the reduced model.

That is, if we let _____ be the sum of squared residuals for the full model and _____ be the corresponding sum for the reduced model, then _____

The Partial F

The test is carried out by fitting both the full and reduced models.

Because the full model contains not only the predictors of the reduced model but also some extra predictors, it should fit the data at least as well as the reduced model.

That is, if we let $\underline{SSE_{full}}$ be the sum of squared residuals for the full model and $\underline{SSE_{red}}$ be the corresponding sum for the reduced model, then $\underline{SSE_{full} < SSE_{red}}$

The Partial F

Intuitively, if _____ is a great deal smaller than _____, the full model provides a much better fit than the reduced model; the appropriate test statistic should then depend on the reduction _____ in unexplained variation.

Test statistic value:

Rejection region:

The Partial F

Intuitively, if $\underline{SSE_{full}}$ is a great deal smaller than $\underline{SSE_{red}}$, the full model provides a much better fit than the reduced model; the appropriate test statistic should then depend on the reduction $\underline{SSE_{red} - SSE_{full}}$ in unexplained variation.

Test statistic value:

$$F = \frac{(SSE_{red} - SSE_{full}) / (p - k)}{SSE_{full} / (n - (p + 1))}$$

Rejection region:

The Partial F

Intuitively, if \underline{SSE}_{full} is a great deal smaller than \underline{SSE}_{red} , the full model provides a much better fit than the reduced model; the appropriate test statistic should then depend on the reduction $\underline{SSE}_{red} - \underline{SSE}_{full}$ in unexplained variation.

Test statistic value:

$$F = \frac{(\underline{SSE}_{red} - \underline{SSE}_{full}) / (p - k)}{\underline{SSE}_{full} / (n - (p + 1))}$$

Rejection region: $f \geq F_{\alpha, p-k, n-(p+1)}$

Model Selection

So far, we have discussed a few of methods for finding the “best” model:

1. Comparison of adjusted R^2 .
2. F-test for model utility and F-test for determining significance of a subset of predictors.

There are other model selection techniques too:

3. Individual parameter t-tests.
4. Reduction of collinearity.
5. 'Best' transformations.
6. Forward/backward selection.

We will elaborate more on these and do some examples in the next day(s) of lecture.