# Lecture 5: Hypothesis testing with the classical linear model

# Assumption MLR6: Normality

$$u \sim N(0, \sigma^2)$$
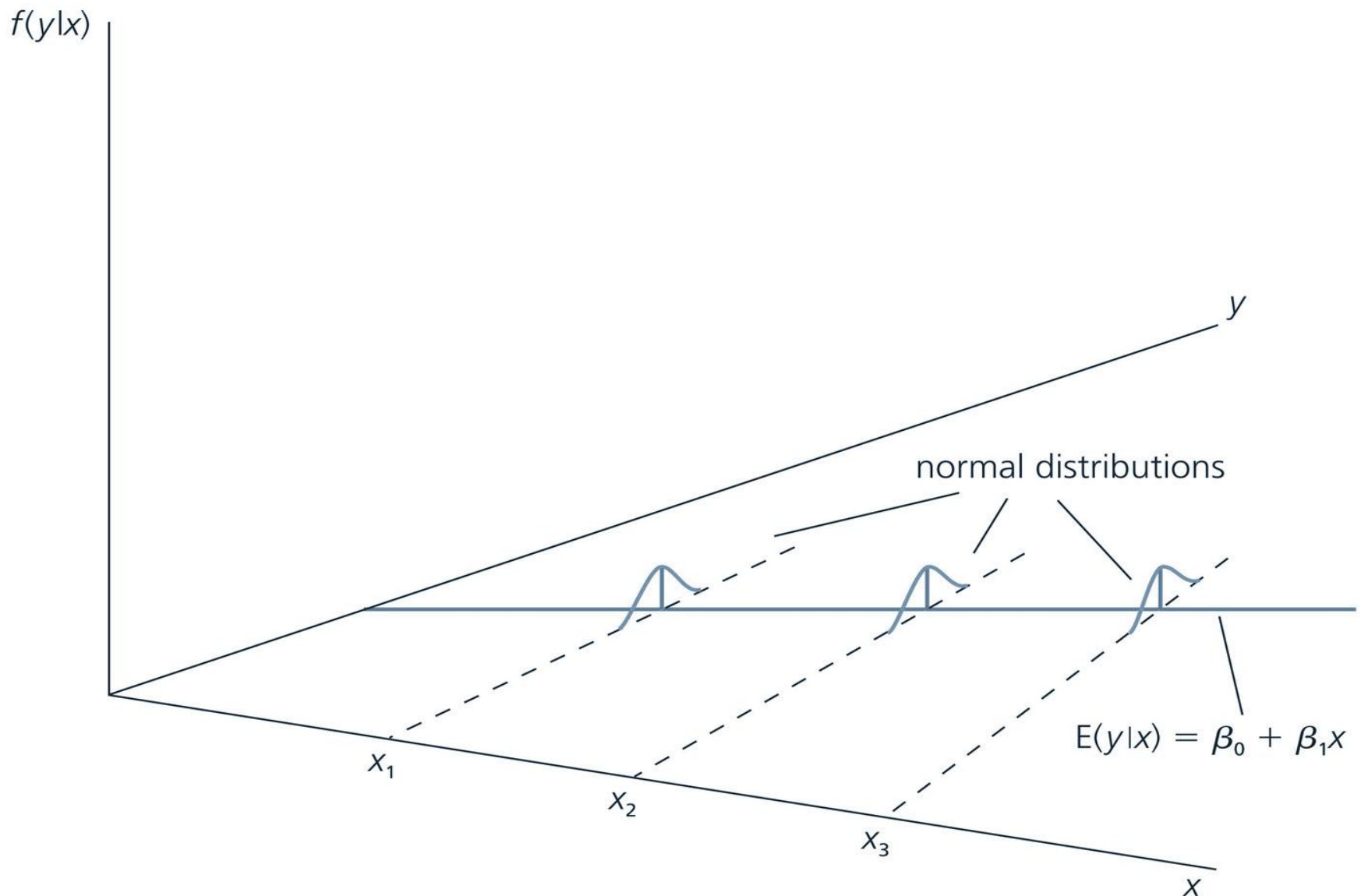
$$E(u \mid x_1, x_2, \ldots, x_k) = E(u) = 0$$

$$Var(u \mid x_1, x_2, \ldots, x_k) = Var(u) = \sigma^2$$

- MLR6 is not one of the Gauss-Markov assumptions. It's not necessary to assume the error is normally distributed in order to obtain the best linear unbiased estimator from OLS.

- MLR6 makes OLS the best unbiased estimator (linear or not), and allows us to conduct hypothesis tests.

**FIGURE 4.1**

**The homoskedastic normal distribution with a single explanatory variable.**



$f(y|x)$

$y$

normal distributions

$E(y|x) = \beta_0 + \beta_1 x$

$x_1$

$x_2$

$x_3$

$x$

# Assumption MLR6: Normality

- But if Y takes on only *n* distinct values, for any set of values of X, the residual can take on only *n* distinct values.
  - Non-normal errors: we can no longer trust hypothesis tests
  - Heteroscedasticity
- With a dichotomous Y, run a logit or probit model
- With an ordered categorial Y, ordered probit/logit
- With an unordered categorial Y, multinomial probit/logit
- With non-negative integer counts, poisson or negative binomial models
- But in chapter 5 we'll see large sample size can overcome this problem.

# Assumption MLR6: Normality

- If we assume that the error is normally distributed conditional on the x's, it follows:

$$\hat{\beta}_j \sim N(\beta_j, Var(\hat{\beta}_j)$$

$$(\hat{\beta}_j - \beta_j) / sd(\hat{\beta}_j) \sim N(0,1)$$

- In addition, any linear combination of beta estimates is normally distributed, and multiple estimates are jointly normally distributed.

# Assumption MLR6: Normality

- In practice, beta estimates follow the t-distribution:

$$(\hat{\beta}_j - \beta_j) / se(\hat{\beta}_j) \sim t_{n-k-1}$$

- Where *k* is the number of slope parameters, *k+1* is the number of unknown parameters (including intercept), *n* is the sample size, and *n-k-1* is the total degrees of freedom.

# Hypothesis testing:

1) State null and research hypotheses

2) Select significance level

3) Determine critical value for test statistic (decision rule for rejecting null hypothesis)

4) Calculate test statistic

5) Either reject or fail to reject (not "accept" null hypothesis)

# Hypothesis testing:

# Hypothesis Testing, example

- Go back to state poverty and homicide rates.

- Two-tailed vs. one-tailed test, which should we do? How do we read the output differently for the two tests?

$$H_0 : \beta = 0$$

$$H_1 : \beta \neq 0$$

two-tailed

$$H_0 : \beta \leq 0$$

$$H_1 : \beta > 0$$

one-tailed

# Hypothesis Testing, cont.

- We typically use two-tailed tests when we have no *a priori* expectation about the direction of a specific relationship. Otherwise, we use a one-tailed test.

- With poverty and homicide, a one-tailed test is justifiable.

- Hypothesis tests and confidence intervals in Stata regression output report t-test statistics, but it is important to understand what they mean because we don't always want to use exactly what Stata reports.

# Hypothesis Testing, cont.

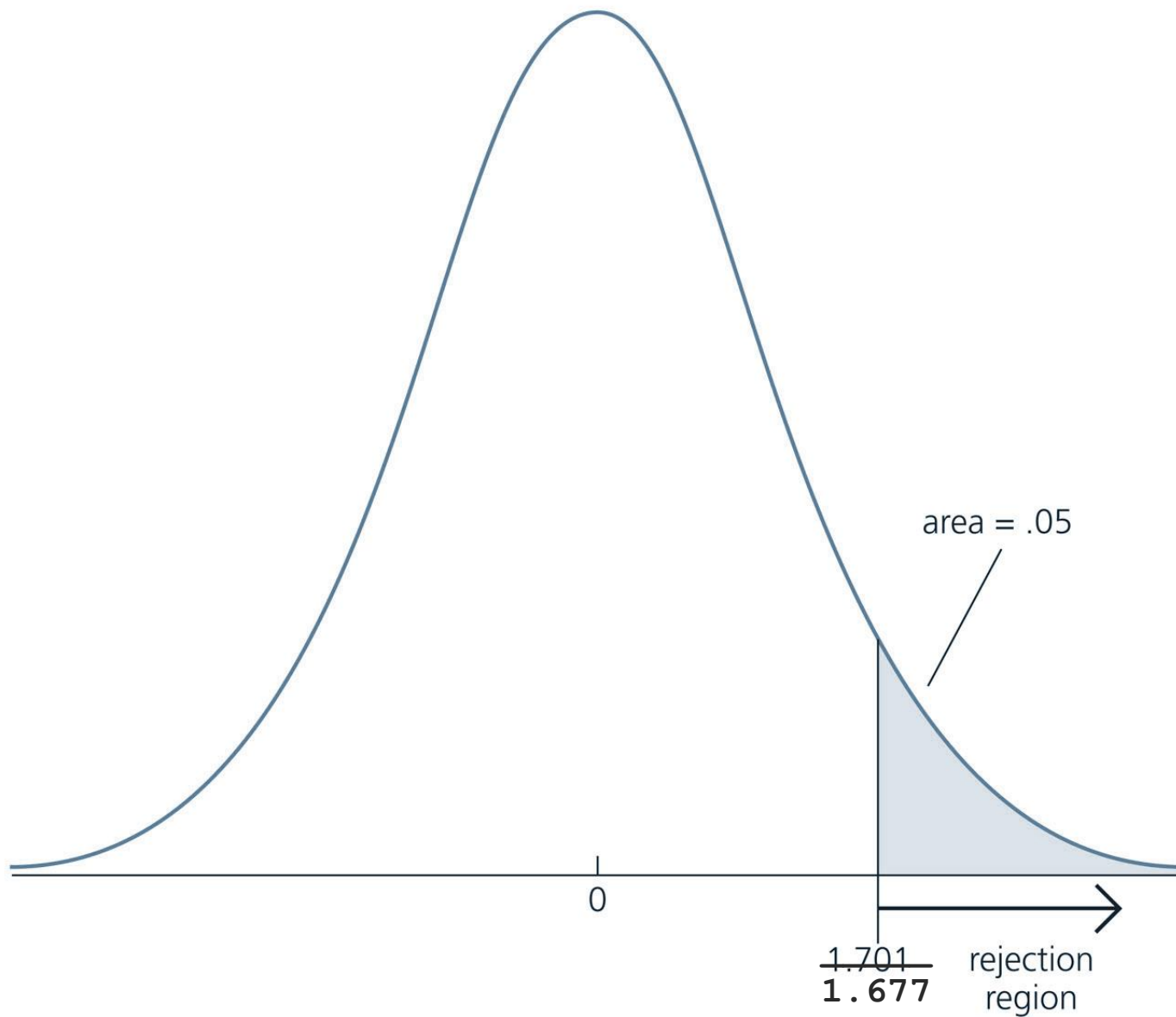- The alpha level is the chance that you will falsely reject the null hypothesis, Type 1 error.

- Step 2: select alpha level
  - Don't always use .05 alpha level.
  - Consider smaller alphas for very large samples, or when it's particularly important that you don't falsely reject the null hypothesis.
  - Use larger alphas for very small samples, or if it's not a big deal to falsely reject the null.

# Hypothesis Testing, cont.

- Step 3: determine critical value. This depends on the test statistic distribution, the alpha level and whether it's a one or two-tailed test.

- In a one-tailed t-test, with an alpha of .05, and a large sample size (>120), the critical value would be 1.64.

- But with 48 degrees of freedom (N-k-1), the critical value is ~1.68 (see Table 6.2, page 825).

FIGURE 4.2

**5% rejection rule for the alternative H$_1$: $\beta_j > 0$ with 28 df.** - - - **48** *df*



area = .05

0

1.701
**1.677**

rejection
region

# Hypothesis Testing, cont.

- To find critical t-statistics in Stata:

  - `. di invttail(48,.05)`

- You should look up these commands (ttail, invtail) and make sure you understand what they are doing. These are part of a larger class of density functions.

# Hypothesis Testing, cont.

- The test statistic is calculated as follows:

$$t_{\hat{\beta}} = \frac{\hat{\beta} - \beta_{H_0}}{\hat{\sigma}_{\hat{\beta}}}$$

- The null hypothesis value of beta is subtracted from our estimate and divided by its estimated standard error.

- This is compared to our pre-determined test statistic. If it's larger than 1.677, we reject the null.

# Hypothesis Testing, cont.

- Returning to the poverty and homicide rate example, we have an estimated beta of .475 and a standard error of .103.   If our null hypothesis is:

$$H_0 : \beta \leq 0$$

Then the test statistic is: $t.s. = \dfrac{.475 - 0}{.103} = 4.62$

- 4.62>1.677, so we reject the null hypothesis.

# Hypothesis Testing, cont.

- Stata also reports p-values in regression output. We would reject the null for any two-sided test where the alpha level is larger than the p-value. It's the area under the curve in a two-sided test.

- To find exact one-sided p values for the t-distribution in Stata:
  - `. di ttail(48,4.62)`
  - `.00001451`, so we would reject the null with any conventional alpha level

# Hypothesis Testing, warning

- Regression output always reports two-sided tests. You have to divide stata's outputted *p* values by 2 in order to get one-tailed tests, but make sure the coefficient is in the right direction!

- On the other hand, `ttail`, and `invttail` always report one-tailed values. Adjust accordingly.

- `Ttail` & `invttail` worksheet

...SALMON JELLY BEANS AND ACNE ($p > 0.05$).

...RED JELLY BEANS AND ACNE ($p > 0.05$).

...TURQUOISE JELLY BEANS AND ACNE ($p > 0.05$).

...MAGENTA JELLY BEANS AND ACNE ($p > 0.05$).

...YELLOW JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND NO LINK BETWEEN GREY JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND NO LINK BETWEEN TAN JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND NO LINK BETWEEN CYAN JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND A LINK BETWEEN GREEN JELLY BEANS AND ACNE ($p < 0.05$). WHOA!

WE FOUND NO LINK BETWEEN MAUVE JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND NO LINK BETWEEN BEIGE JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND NO LINK BETWEEN LILAC JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND NO LINK BETWEEN BLACK JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND NO LINK BETWEEN PEACH JELLY BEANS AND ACNE ($p > 0.05$).

WE FOUND NO LINK BETWEEN ORANGE JELLY BEANS AND ACNE ($p > 0.05$).

- What went wrong?
- What are the chances of finding at least one statistically significant variable at $p<.05$ when you are testing 20 variables?

`. di binomialtail(20,1,.05)` = .64

- There's a 64% chance of having at least one "statistically significant" result.
- This is the problem of multiple comparisons. How can you correct for this?

- The most common and simplest is the Bonferroni correction where you replace your original alpha level with alpha/k where k is the number of comparisons you make.

```
. di binomialtail(20,1,.05/20)
.04883012
```

# Confidence intervals

- Confidence intervals are related to hypothesis tests, but are interpreted much differently.

$$CI : \hat{\beta}_j \pm c \cdot se(\hat{\beta}_j)$$

- c is the t-value needed to obtain the correct % confidence interval. The 97.5% one-sided t-value is needed for a 95% confidence interval.

- Confidence intervals are always two-sided.

- Given the sample data, the confidence interval tells us, with X% confidence, that the *true* parameter falls within a certain range.

# Confidence intervals

- Going back to the homicide rate and poverty example, the estimated parameter for poverty was .475 with a standard error of .103.

- The 95% confidence interval, reported by Stata is .475+/-.103*2.01 = [.268,.682]

- So, with 95% confidence, the population value for the effect of poverty on homicide is between those two numbers.

- The 99% confidence interval will be wider in order to have greater confidence that the true value falls within that range:

  - .475+/-.103*2.68=[.199,.751]

# Example 4.2 (p. 126-8): student performance and school size

- Hypotheses: $H_0 : \beta_{enroll} \geq 0$

$$H_1 : \beta_{enroll} < 0$$

- Alpha: .05, one tailed, tcrit=-1.65
- Reject null hypothesis if t.s.<-1.65
- The estimated coefficient on enrollment, controlling for teacher compensation and staff:student ratio is -.00020 with a .00022 standard error.
- So the test statistic equals -.00020/.00022=-.91, fail to reject null.
- Functional form can change our conclusions! When school size is logged, we do reject the null.

# Other hypotheses about *β*

- We may want to test the hypothesis that β equals 1, or some other number besides zero. In this case, we proceed exactly as before, but t-statistic won't match the regression output. We subtract the hypothesized parameter size (now non-zero) from the parameter estimate before dividing by the standard error.

$$t_{\hat{\beta}} = \frac{\hat{\beta} - \beta_{H_0}}{\hat{\sigma}_{\hat{\beta}}}$$

- Stata, helpfully, will do this for us. After any regression type: "test varname=X", inserting the appropriate variable name and null parameter value.
- Example 4.4, p. 130-131

# Linear combinations of $\beta s$

- In section 4.4, Wooldredge goes through a detailed explanation of how to transform the estimated regression model in order to obtain $se(\beta_1 + \beta_2)$, which is necessary in order to directly test the hypothesis that $\beta_1 + \beta_2 = 0$.

- This method is correct, and it is useful to follow, I just prefer a different method after a regression model:

- "test x1+x2=0", replacing x1 and x2 with your variable names.

# Testing multiple linear restrictions

- Restricted model: Multiple restrictions are imposed on the data. (e.g. linearity, additivity, $X_j=0$, $X_j=X_k$, $X_j=3$, etc.)

- Unrestricted model: At least one of the above assumptions is relaxed, often by adding an additional predictor to the model.

- To test the null hypothesis, we conduct an F-test:

# F-test for restricted/unrestricted models

$$F\left(k_{UR} - k_R, n - k_{UR}\right) = \frac{\left(SSR_R - SSR_{UR}\right)/\left(k_{UR} - k_R\right)}{SSR_{UR}/\left(n - k_{UR}\right)}$$

- Where SSR refers to the residual sum of squares, and k refers to the number of regressors (including the intercept).

# F-test for restricted/unrestricted models, example

**Restricted model:**

```
. reg homrate poverty

      Source |       SS       df       MS              Number of obs =      50
-------------+------------------------------           F(  1,    48) =   21.36
       Model |  100.175656        1  100.175656        Prob > F      =  0.0000
    Residual |  225.109343       48  4.68977798        R-squared     =  0.3080
-------------+------------------------------           Adj R-squared =  0.2935
       Total |  325.284999       49  6.63846936        Root MSE      =  2.1656


------------------------------------------------------------------------------
     homrate |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
     poverty |    .475025   .1027807     4.62   0.000     .2683706    .6816795
       _cons |  -.9730529   1.279803    -0.76   0.451     -3.54627    1.600164
```

- **Why is this "restricted"? What restrictions are we imposing, and how might we test these?**

# F-test for restricted/unrestricted models, example

**Unrestricted model:**

```
. reg homrate poverty gradrate het

      Source |       SS       df       MS              Number of obs =      50
-------------+------------------------------           F(  3,    46) =   19.72
       Model |  183.012608     3  61.0042025           Prob > F      =  0.0000
    Residual |  142.272391    46  3.09287807           R-squared     =  0.5626
-------------+------------------------------           Adj R-squared =  0.5341
       Total |  325.284999    49  6.63846936           Root MSE      =  1.7587


------------------------------------------------------------------------------
     homrate |      Coef.   Std. Err.       t    P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
     poverty |   .3134171   .0922506      3.40   0.001     .1277263    .4991079
    gradrate |  -.0518914   .0404837     -1.28   0.206    -.1333809    .0295981
         het |   7.098508   2.174708      3.26   0.002     2.721047    11.47597
       _cons |   2.357754   3.913269      0.60   0.550    -5.519249    10.23476
------------------------------------------------------------------------------
```

- **Here, we have lifted the restriction that $\beta_{gradrate}=\beta_{het}=0$. The F-test of this restriction is calculated as follows:**

# F-test for restricted/unrestricted models, example

$$F\left(k_{UR} - k_R, n - k_{UR}\right) = \frac{\left(SSR_R - SSR_{UR}\right)/\left(k_{UR} - k_R\right)}{SSR_{UR}/\left(n - k_{UR}\right)}$$

$$F\left(4 - 2, 50 - 4\right) = \frac{\left(225.1 - 142.3\right)/\left(4 - 2\right)}{142.3/\left(50 - 4\right)}$$

$$F(2, 46) = \frac{41.4}{3.1} = 13.4, \left(p < .000\right)$$

- We can find the p-value in Stata using "`di Ftail(2,46,13.4)`" – or we can let Stata do all the calculations with "`test gradrate het`" after the unrestricted model. This is testing two restrictions jointly: `gradrate`=0 & `het`=0.

# F-test for restricted/unrestricted models, example

- This kind of test is appropriate when the difference between two models can be expressed as a set of restrictions or assumptions.

- It may take a little bit of imagination to recognize the "restrictions" in your restricted model.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_j x_j + u$$

- Examples:
  - $\beta_{j+1}=0$
  - the coefficient on the product of $x_1$ and $x_2$ is zero.

# F-test for restricted/unrestricted models, example

- One special type of restriction which is sometimes of interest in criminology, is that our models are the same across different groups. This follows the same logic, but is called a Chow test.

# F-test for restricted/unrestricted models, Chow test example

- The Chow test can be used in a couple different situations.
  - Completely different sets of data with the same variables
  - Sub-populations within one datset.
- Either way, we compare the SSR from a regression model with the two sets of data (or groups) pooled (restricted model), to the summed SSR from two separate regression models (unrestricted).
- What is the restriction in the restricted model?

# F-test for restricted/unrestricted models, Chow test example

Unrestricted model (two groups):

$$Y = \alpha_1 + \beta_1 X_1 + \delta_1 X_2 + \ldots + \gamma_1 X_k + \varepsilon$$

$$Y = \alpha_2 + \beta_2 X_1 + \delta_2 X_2 + \ldots + \gamma_2 X_k + \varepsilon$$

Restricted model (pooled):

$$Y = \alpha + \beta X_1 + \delta X_2 + \ldots + \gamma X_k + \varepsilon$$

$$\alpha_1 = \alpha_2, \beta_1 = \beta_2, \ldots \qquad \leftarrow \text{restrictions}$$

# F-test for restricted/unrestricted models, Chow test example

- Suppose we have a model for teen delinquency, but we think it differs for males and females. **Restricted model (note: we don't control for gender here):**

```
. reg dfreq1 age1 hisp black other msgrd sus1 r_wk biop1 smoke1

      Source |       SS       df       MS              Number of obs =    8669
-------------+------------------------------           F(  9,  8659) =   56.05
       Model |  86926.9336        9  9658.54818         Prob > F      =  0.0000
    Residual |  1492077.55     8659  172.315227         R-squared     =  0.0551
-------------+------------------------------           Adj R-squared =  0.0541
       Total |  1579004.48     8668    182.1648         Root MSE      =  13.127


------------------------------------------------------------------------------
      dfreq1 |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
        age1 |   .0392229   .1006811     0.39   0.697    -.1581361    .2365818
        hisp |   .3504168   .4325393     0.81   0.418     -.497463    1.198297
       black |  -.8410179   .3773185    -2.23   0.026    -1.580652   -.1013839
       other |  -.3523527   .4795827    -0.73   0.463    -1.292449    .5877435
       msgrd |   -.347948   .0939261    -3.70   0.000    -.5320656   -.1638304
        sus1 |   4.032407   .3474379    11.61   0.000     3.351346    4.713468
        r_wk |   .4253605   .1738242     2.45   0.014     .0846237    .7660973
       biop1 |  -.7830546    .300954    -2.60   0.009    -1.372996   -.1931132
      smoke1 |   3.773678   .3128116    12.06   0.000     3.160493    4.386864
       _cons |   2.279977   1.563863     1.46   0.145     -.785567    5.345521
------------------------------------------------------------------------------
```

# F-test for restricted/unrestricted models, Chow test example

- **Unrestricted model, part 1:**

```
. reg dfreq1 age1 hisp black other msgrd sus1 r_wk biop1 smoke1 if male==1

      Source |       SS           df       MS              Number of obs =     4436
-------------+----------------------------------           F(  9,   4426) =    30.13
       Model |  76722.9184         9   8524.76872           Prob > F        =   0.0000
    Residual |   1252212.9      4426    282.92203           R-squared       =   0.0577
-------------+----------------------------------           Adj R-squared   =   0.0558
       Total |  1328935.82      4435   299.647311           Root MSE        =    16.82

------------------------------------------------------------------------------
      dfreq1 |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
        age1 |   .1722392   .1807942     0.95   0.341    -.1822079    .5266863
        hisp |   .9256955    .778416     1.19   0.234    -.6003892     2.45178
       black |  -.9784988   .6809885    -1.44   0.151    -2.313577    .3565793
       other |  -.9274644     .86288    -1.07   0.283    -2.619141    .7642119
       msgrd |   -.322382    .165856    -1.94   0.052    -.6475428    .0027788
        sus1 |    4.10308   .5883441     6.97   0.000     2.949631    5.256529
        r_wk |   .4572663   .3020877     1.51   0.130    -.1349767    1.049509
       biop1 |  -1.566485   .5393016    -2.90   0.004    -2.623785   -.5091838
      smoke1 |   5.485458   .5611064     9.78   0.000     4.385409    6.585507
       _cons |   .7287458   2.803031     0.26   0.795    -4.766596    6.224088
------------------------------------------------------------------------------
```

# F-test for restricted/unrestricted models, Chow test example

- **Unrestricted model, part 2:**

```
. reg dfreq1 age1 hisp black other msgrd sus1 r_wk biop1 smoke1 if male==0

      Source |       SS          df       MS              Number of obs =     4233
-------------+------------------------------              F(  9,  4223) =   27.02
       Model | 12738.1839        9  1415.35377            Prob > F      =  0.0000
    Residual | 221189.499     4223   52.377338            R-squared     =  0.0545
-------------+------------------------------              Adj R-squared =  0.0524
       Total | 233927.682     4232  55.2759174            Root MSE      = 7.2372


------------------------------------------------------------------------------
      dfreq1 |      Coef.   Std. Err.       t    P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
        age1 |  -.0747234   .0793737    -0.94   0.347    -.2303376    .0808908
        hisp |  -.2843726   .3401383    -0.84   0.403    -.9512225    .3824774
       black |  -.5786616   .2963413    -1.95   0.051    -1.159646    .0023231
       other |   .1702421   .3771467     0.45   0.652    -.5691639     .909648
       msgrd |  -.1599254   .0776844    -2.06   0.040    -.3122276   -.0076232
        sus1 |   2.685044   .3034377     8.85   0.000     2.090147    3.279942
        r_wk |   .1609845   .1428534     1.13   0.260    -.1190833    .4410524
       biop1 |  -.2722414   .2388061    -1.14   0.254    -.7404268    .1959441
      smoke1 |   2.155191   .2475283     8.71   0.000     1.669905    2.640476
       _cons |   2.556738   1.240103     2.06   0.039     .1254827    4.987992
------------------------------------------------------------------------------
```

# F-test for restricted/unrestricted models, Chow test example

- **Chow test proceeds as follows:**

$$F\left(k_{UR} - k_R, n - k_{UR}\right) = \frac{\left(SSR_R - SSR_{UR}\right)/\left(k_{UR} - k_R\right)}{SSR_{UR}/\left(n - k_{UR}\right)}$$

$$F\left(20 - 10, 8669 - 20\right) = \frac{\left(1492078 - \left(1252212 + 221189\right)\right)/\left(20 - 10\right)}{\left(1252212 + 221189\right)/\left(8669 - 20\right)}$$

$$F(10, 8649) = \frac{18675/10}{1473402/8649} = 10.96, (p < .001)$$

- **Alternately, we could interact male with all other variables, run a fully interactive model, and . . .**

# F-test for restricted/unrestricted models, Chow test example

```
. reg dfreq1 age1 hisp black other msgrd sus1 r_wk biop1 smoke1 male mage1 mhisp mblack mother
    mmsgrd msus1 mr_wk mbiop1 msmoke1

      Source |       SS       df       MS              Number of obs =    8669
-------------+------------------------------           F( 19,  8649) =   32.63
       Model | 105602.083      19  5558.00435          Prob > F      =  0.0000
    Residual |  1473402.4    8649  170.355232          R-squared     =  0.0669
-------------+------------------------------           Adj R-squared =  0.0648
       Total | 1579004.48    8668    182.1648          Root MSE      =  13.052


------------------------------------------------------------------------------
      dfreq1 |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
        age1 |  -.0747234   .1431472    -0.52   0.602    -.355326    .2058792
        hisp |  -.2843726   .6134251    -0.46   0.643   -1.486832    .9180869
       black |  -.5786616   .5344391    -1.08   0.279    -1.62629    .4689663
       other |   .1702421   .6801683     0.25   0.802    -1.16305    1.503534
. . .
. . .
      mbiop1 |  -1.294243   .6005073    -2.16   0.031   -2.471381   -.1171058
     msmoke1 |   3.330267    .623581     5.34   0.000      2.1079    4.552634
       _cons |   2.556738   2.236474     1.14   0.253   -1.827285     6.94076
------------------------------------------------------------------------------

. test male mage1 mhisp mblack mother mmsgrd msus1 mr_wk mbiop1 msmoke1


      F( 10,  8649) =    10.96   <----- get the same answer!
           Prob > F =    0.0000
```

# F-test for restricted/unrestricted models, Chow test example

- We know that there are significant differences in average levels of delinquency between males and females.

- Part of the Chow test is that there is no difference in average levels between the two groups (same intercept).

- How would we test a modified Chow test where we allow males and females to have different levels of delinquency and just test if the effects of the covariates differ between the genders?

# F-test for restricted/unrestricted models, other uses

- This general test is used to calculate the overall F-statistic for every regression model. The restricted model is intercept only where all parameters are assumed to be zero.

- Interaction terms

# Stata's saved regression results

- After any regression:
    - "ereturn list" returns a list of all stored results
    - e(N): number of observations
    - e(mss): model sum of squares
    - e(df_m): model degrees of freedom
    - e(rss): residual sum of squares
    - e(df_r): residual degrees of freedom
    - e(F): F statistic
    - e(r2): r-squared
    - e(r2_a): adjusted r-squared
    - e(rmse): root mean squared error

# Next time:

Homework 6 Problems 4.2, 4.4, C4.6, C4.8

Answers posted – do not turn in.

Midterm: available today after class, due by 4:40pm 10/4, open to any non-interactive resource (books/notes/lectures/internet pages), but not other people.

Read: Wooldridge Chapter 5 (skim), Chapter 6