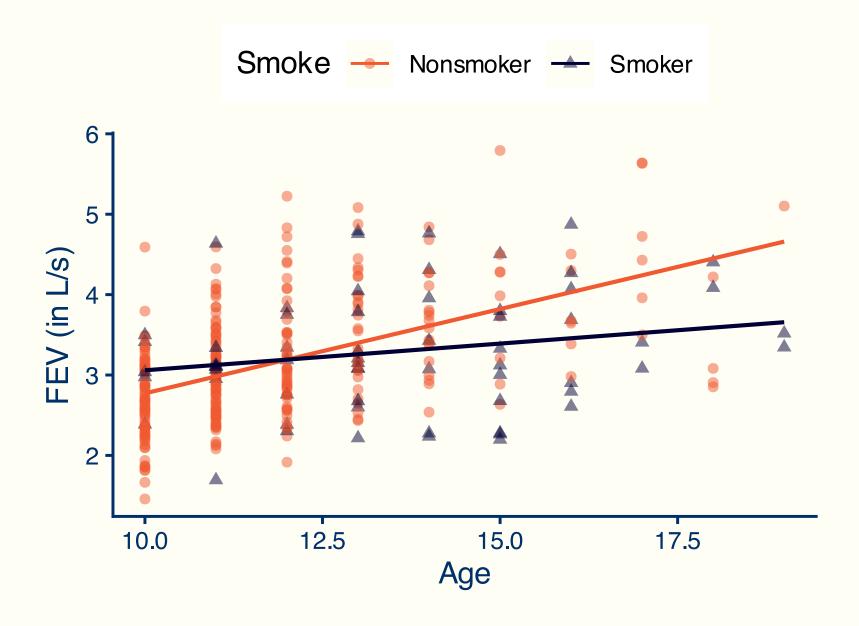
Additional Inferential Tools for MLR

Stat 230: Applied Regression Analysis

Linear Combinations

FEV example



FEV example

$$E(Y|X) = \beta_0 + \beta_1$$
age + β_2 smoke + β_3 age × smoke

Mean functions for each group

Smokers:
$$E(Y|X) = (\beta_0 + \beta_2) + (\beta_1 + \beta_3)$$
age

Non-smokers:
$$E(Y|X) = \beta_0 + \beta_1$$
 age

SEs for linear combinations

Estimate:

$$\widehat{\gamma} = c_1 \widehat{\beta}_i + c_2 \widehat{\beta}_j$$

Standard error:

$$SE(\widehat{\gamma}) = \sqrt{c_1^2 \{SE(\widehat{\beta}_i)\}^2 + c_2^2 \{SE(\widehat{\beta}_j)\}^2 + 2c_1c_2Cov(\widehat{\beta}_i, \widehat{\beta}_j)}$$

FEV example



Our standard table of coefficients doesn't contain the SEs we need to construct CIs or perform tests for these linear combinations.

term	estimate	std.error	statistic p.value
(Intercept)	0.675	0.251	2.684 0.008
Age	0.210	0.021	9.972 < 0.001
SmokeSmoker	1.720	0.563	3.053 0.002
Age:SmokeSmoker	-0.143	0.042	-3.397 <0.001

Covariance matrix

- The diagonal (top left to lower right) displays variances, $\{SE(\widehat{\beta}_j)\}^2$
- The off-diagonals display covariances, $Cov(\widehat{\beta}_i, \widehat{\beta}_j)$
- Covariance is like correlation, but not scaled to be between -1 and 1

```
(Intercept)
                                Age SmokeSmoker Age:SmokeSmoker
(Intercept)
                 0.063199 - 0.0052210
                                      -0.063199
                                                     0.0052210
                -0.005221 0.0004423
                                       0.005221
                                                    -0.0004423
Age
SmokeSmoker
               -0.063199 0.0052210
                                       0.317248
                                                    -0.0234032
Age: SmokeSmoker
              0.005221 -0.0004423
                                    -0.023403
                                                     0.0017798
```

SE for a linear combination

Smokers: $E(Y|X) = (\beta_0 + \beta_2) + (\beta_1 + \beta_3)$ age

Target: $\hat{\gamma} = \hat{\beta}_1 + \hat{\beta}_3$

```
(Intercept) Age SmokeSmoker Age:SmokeSmoker (Intercept) 0.063199 -0.0052210 -0.063199 0.0052210 -0.005221 0.0004423 0.005221 -0.0004423 SmokeSmoker -0.063199 0.0052210 0.317248 -0.0234032 Age:SmokeSmoker 0.005221 -0.0004423 -0.023403 0.0017798
```

$$SE(\widehat{\gamma}) = \sqrt{c_1^2 \{SE(\widehat{\beta}_1)\}^2 + c_2^2 \{SE(\widehat{\beta}_2)\}^2 + 2c_1c_2Cov(\widehat{\beta}_1, \widehat{\beta}_2)}$$

$$c_1 = 1, c_2 = 1$$

$$SE(\widehat{\beta_1})^2 = 0.0004423$$

$$SE(\widehat{\beta_2})^2 = 0.0017798$$

$$Cov(\widehat{\beta_1}, \widehat{\beta_2}) = -0.0004423$$

CI for a linear combination

Once we have $SE(\hat{\gamma})$, we can construct a CI as usual:

$$\widehat{\gamma} \pm t_{n-(p+1)}^* \cdot SE(\widehat{\gamma})$$

$$\hat{\gamma} = 0.0664$$
 $SE(\hat{\gamma}) = 0.0366$
 $df = n - (p + 1) = 345 - 4 = 341$
 $t_{341}^* \approx 1.649 \text{ for a } 90\% \text{ CI}$

Test for a linear combination

We can also test hypotheses about linear combinations of coefficients using a t-test:

$$H_0: \gamma = \gamma_0 \text{ vs. } H_a: \gamma \neq \gamma_0$$

$$t = \frac{\widehat{\gamma} - \gamma_0}{SE(\widehat{\gamma})}$$

Reference distribution: $t_{n-(p+1)}$

Example

- Work through the example on the handout with your neighbors
- Be prepared to share your strategy with the class

Extra Sums of Squares F-tests

Example: Used car prices

We have data on used car prices, including:

- Price: Price (in dollars)
- Mileage: Mileage (in thousands of miles)
- Make: Manufacturer (Buick, Cadillac, Chevrolet, Pontiac, SAAB, Saturn)

After some initial EDA, researchers want to fit an MLR model to predict log(price) using mileage and make.

Used car prices

Is Make a useful predictor of price, after accounting for Mileage?

term	estimate	std.error	statistic p.value
(Intercept)	10.090	0.034	293.975 < 0.001
Mileage	0.000	0.000	-7.359 <0.001
MakeCadillac	0.649	0.038	17.016 < 0.001
MakeChevrolet	-0.296	0.030	-9.809 < 0.001
MakePontiac	-0.144	0.033	-4.299 < 0.001
MakeSAAB	0.355	0.035	10.108 < 0.001
MakeSaturn	-0.397	0.041	-9.650 < 0.001

Review: Nested models

Model 1:

$$E(\log(Y)|X) = \beta_0 + \beta_1 \text{Mileage}$$

Model 2:

$$E(\log(Y)|X) = \beta_0 + \beta_1 \text{Mileage} + \beta_2 \text{Cadillac} + \beta_3 \text{Chevy} + \beta_4 \text{Pontiac} + \beta_5 \text{SAAB} + \beta_6 \text{Saturn}$$

Extra SS F-test:

Null hypothesis tells us how we can simplify Model 2 (full) to get Model 1 (reduced):

$$H_0: \beta_2 = \beta_3 = \beta_4 = \beta_5 = \beta_6 = 0$$

Extra Sums of Squares

Definition: the marginal reduction in the error sum of squares when one (or more) explanatory variable(s) is added to the model

F-test for comparing nested models:

$$F = \frac{(SSE_{full} - SSE_{reduced})/(df_{full} - df_{reduced})}{MSE_{full}}$$

ANOVA table from R

anova(car_lm) produces an ANOVA table for the MLR model, it gives **sequential sums of squares**, not the extra sums of squares we need for our F-test.

term	df	sumsq	meansq	statistic p.value
Mileage	1	2.959	2.959	50.941 < 0.001
Make	5	85.802	17.160	295.414 < 0.001
Residuals	797	46.297	0.058	NA NA

- Mileage row: extra SS for Mileage (compared to intercept-only model)
- Make row: extra SS for Make after accounting for Mileage
- Residuals row: SSE for the full model