

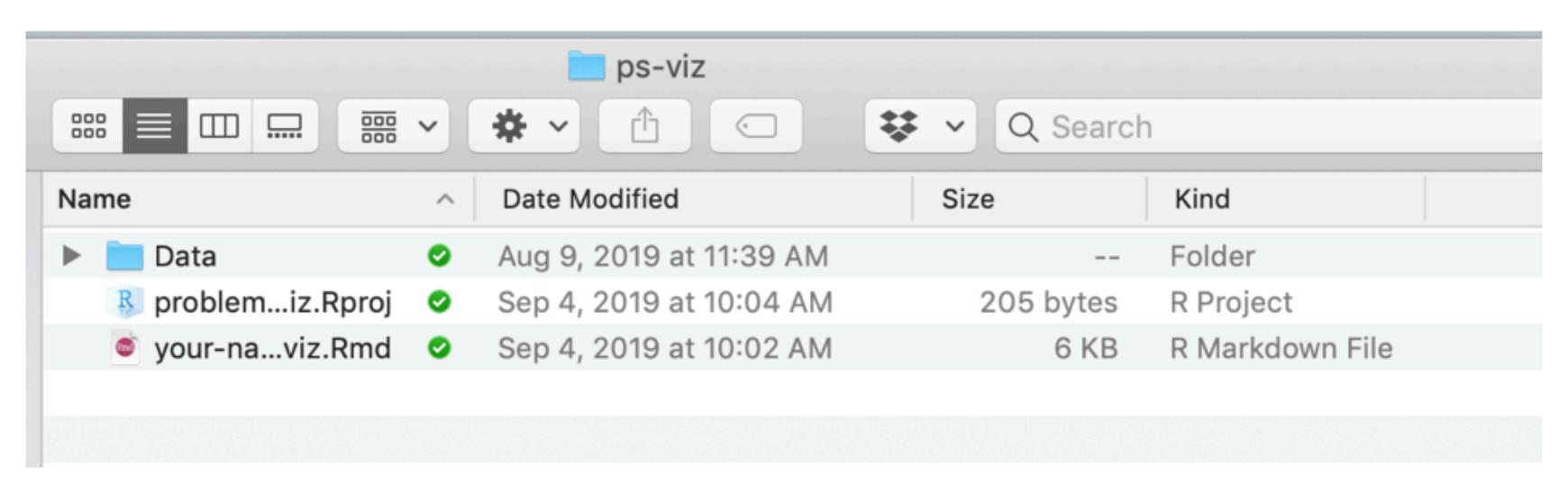
TODAY'S AGENDA

- (1) Regression Recap
- (2) Regression with categorical variables
- "Explaining" outcomes

Loading data

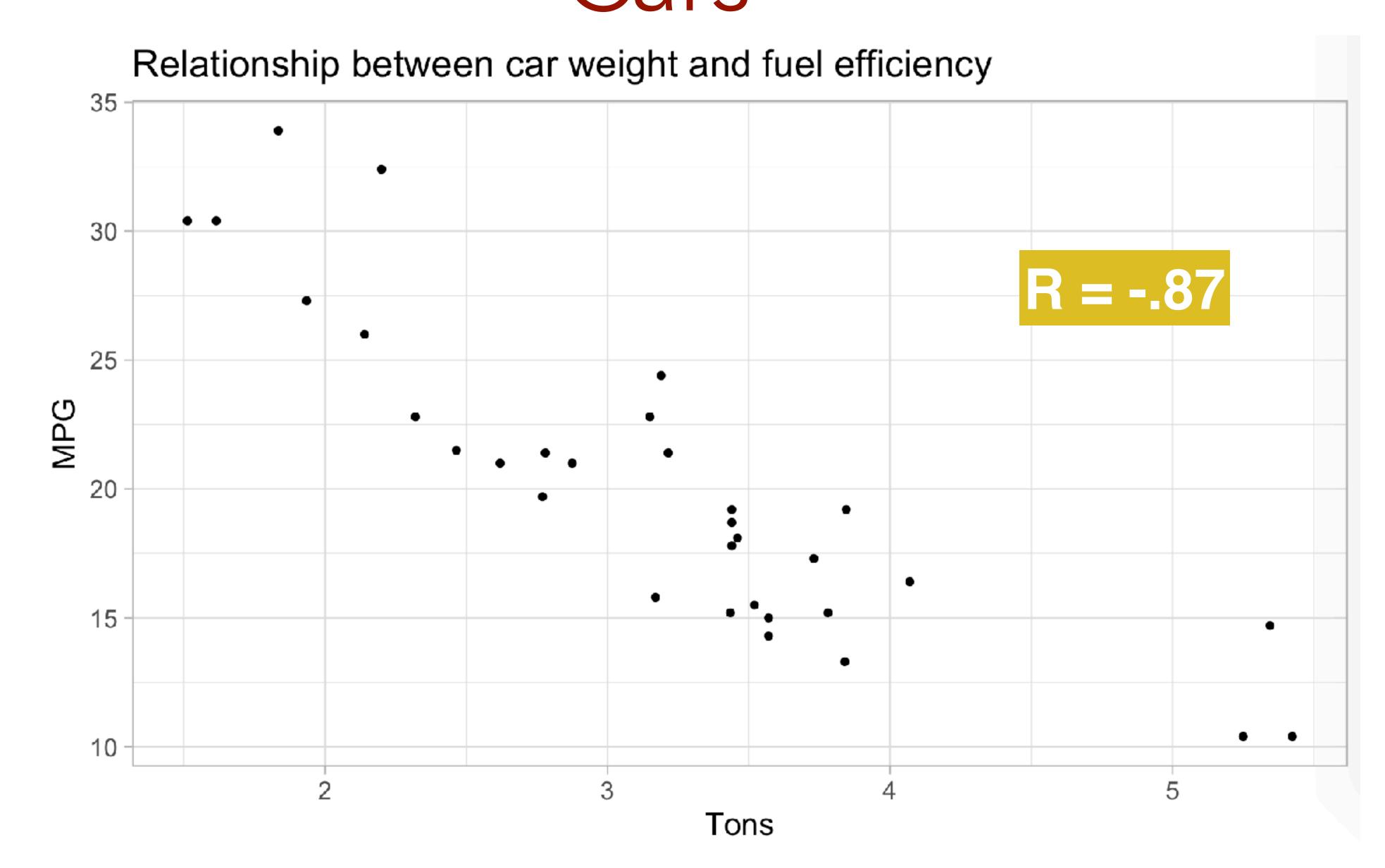
Need to be in RProject

Need file path for data and right extension



```
# read data
movies = read_csv("Data/movies.csv")
```

Cars



Linear models

$$\hat{y} = mx + b$$
 $\hat{y} = \beta_0 + \beta_1 x_1 + \epsilon$

Y Outcome variable

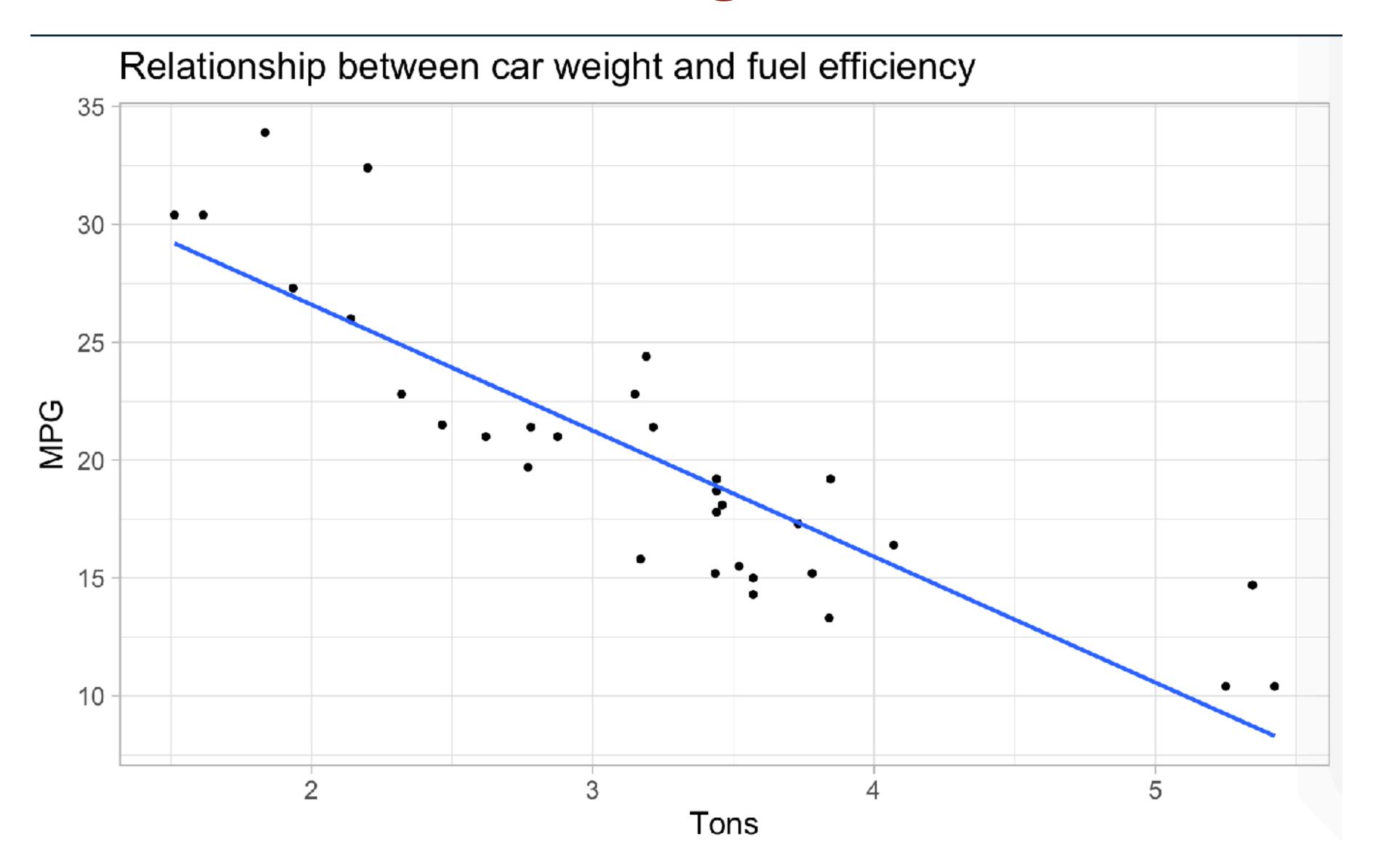
 \hat{x} x_1 Explanatory variable

 $\hat{\beta}_0$ Y intercept

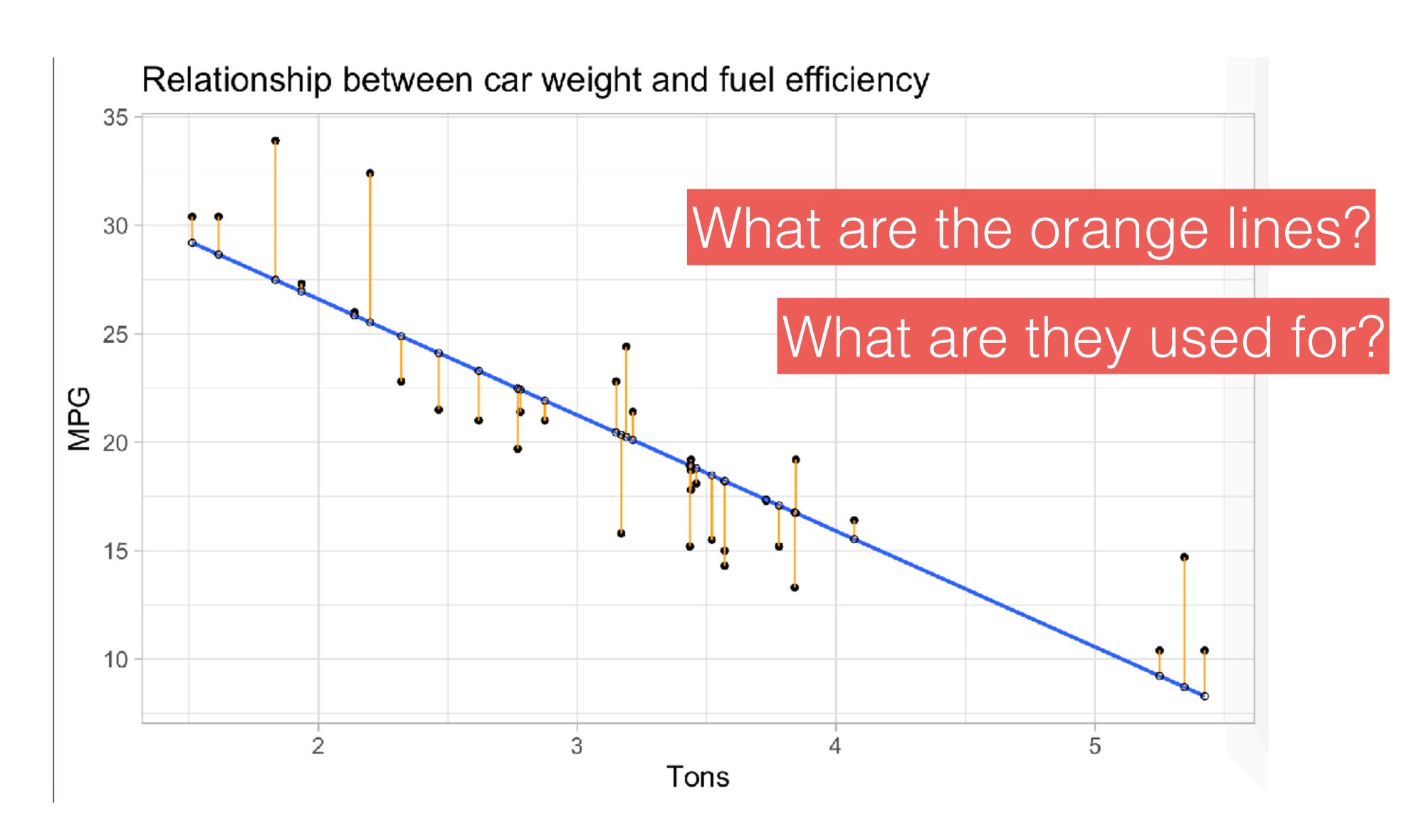
 $\hat{\beta}_1$ Slope

 $\hat{\epsilon}$ Error (residual)

Drawing lines



Drawing lines



Residuals

$$m\hat{p}g = \beta_0 + \beta_1 weight + \epsilon$$

mpg =
$$37.29 - 5.34(wt) + \epsilon$$

What is mpg_hat?

What are residual units in?

```
model_cars = lm(mpg ~ wt, data = mtcars)

## using moderndive
fitted = get_regression_points(model_cars)
```

```
wt mpg_hat residual
     ΙD
          mpg
  <int> <dbl> <dbl>
                      <dbl>
                              <dbl>
                      23.3
                             -2.28
         21
               2.62
         21
               2.88
                      21.9
                             -0.92
                             -2.09
         22.8
              2.32
                      24.9
                             1.30
         21.4
               3.22
                      20.1
                             -0.2
         18.7
               3.44
                       18.9
                             -0.693
         18.1
               3.46
                       18.8
                             -3.90
      7 14.3 3.57
                       18.2
                              4.16
      8 24.4 3.19
                       20.2
      9 22.8 3.15
                       20.4
                              2.35
     10 19.2 3.44
                       18.9
                              0.3
10
```

Regression with categories

Regression with categorical variables

Numeric variables

(continuous)

Numbers

Categorical variables

(Factors)

Words

You tell me

Income

Weight

Gender GDP per capita 18-25 26-34

Strongly Agree

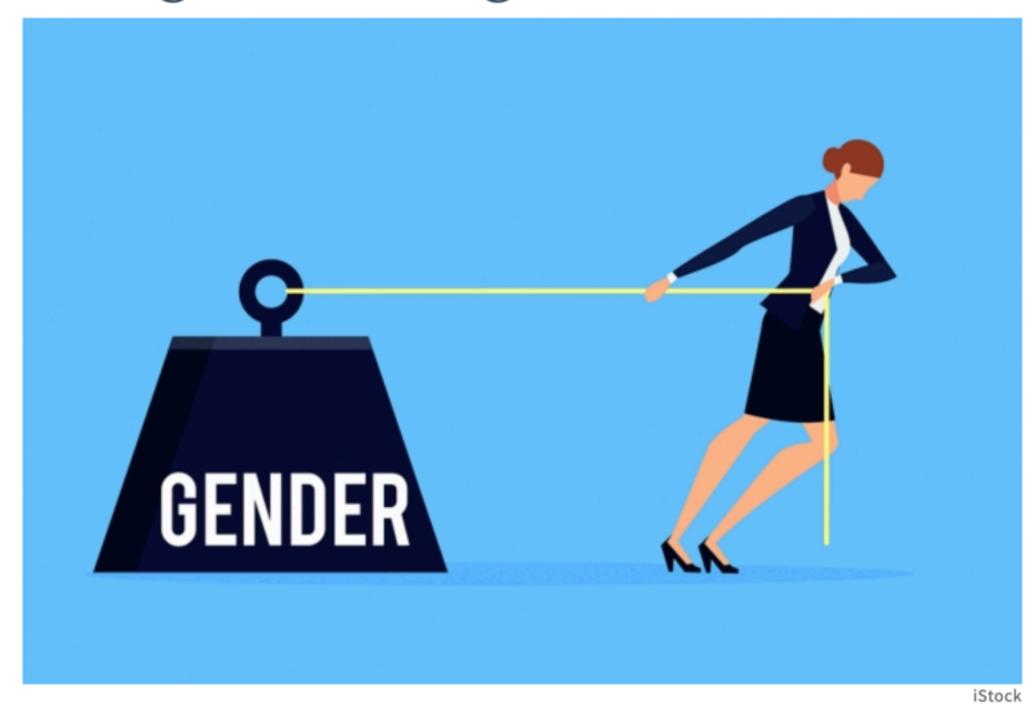
Strongly disagree

Bechdel test

Life Expectancy

Bias in student evaluations

Why We Must Stop Relying on Student Ratings of Teaching



News & Views Careers Events Reports & Data

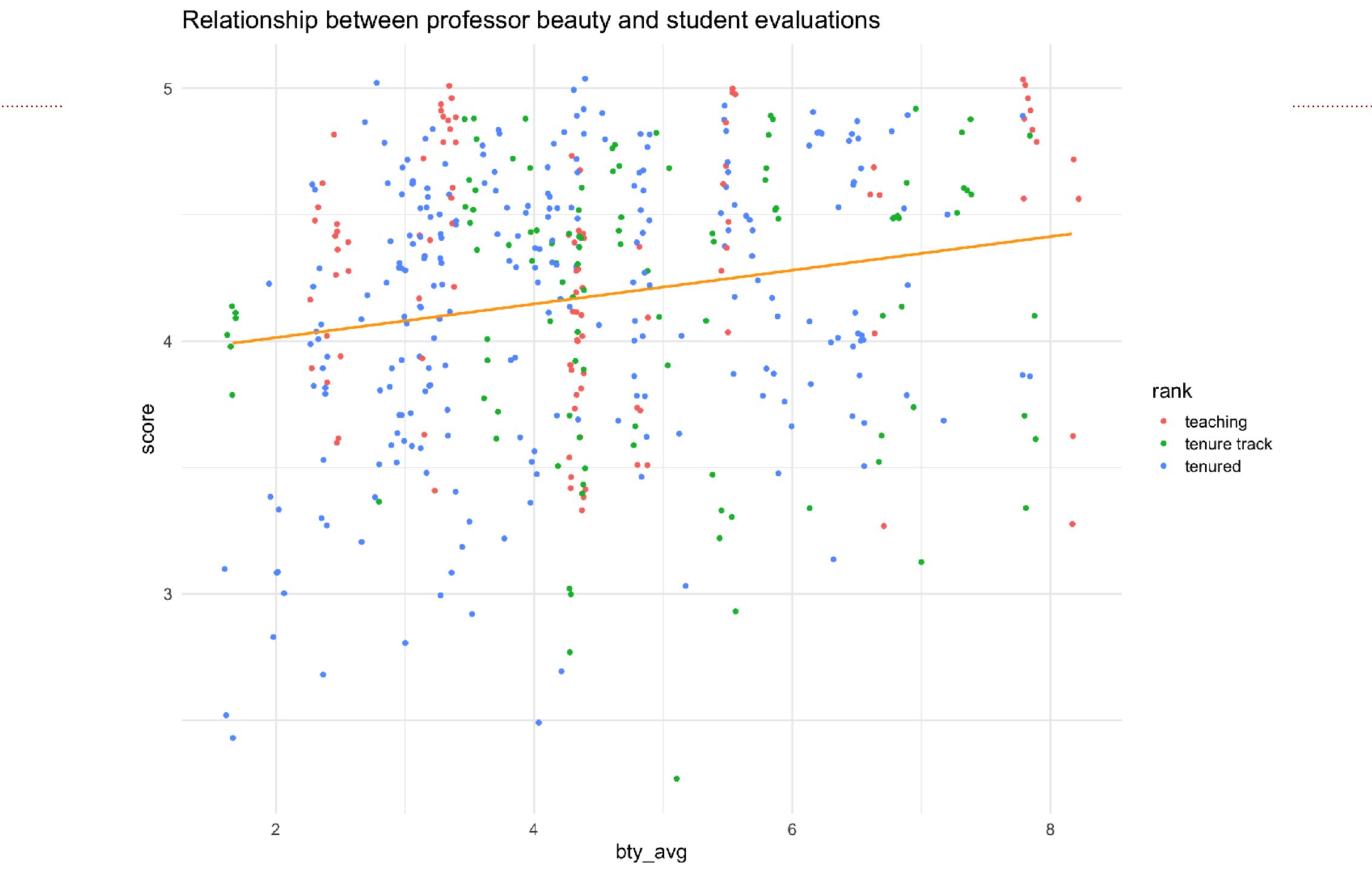
Trending: New Admissions Era? Winning Tenure, in Court



#News #1

#Teaching And Learning

Teaching Evals: Bias and Tenure



Rank differences

```
## average by tenure
evals %>%
group_by(rank) %>%
summarise(avg = mean(score))
```

```
rank avg
<fct> <fct> <dbl>
1 teaching 4.28
2 tenure track 4.15
3 tenured 4.14
```

```
model_rank = lm(score ~ rank, data = evals)

# get model output
```

get_regression_table(model_rank)

```
estimate std_error statistic p_value lower_ci upper_ci
 term
                  <dbl>
                                  <dbl> <dbl>
                          <dbl>
                                                 <dbl>
 <chr>
                                                        <dbl>
            4.28
1 intercept
                          0.054 79.9
                                         0 4.18
                                                        4.39
2 ranktenure track -0.13
                          0.075 \quad -1.73 \quad 0.084
                                                -0.277
                                                        0.017
3 ranktenured
             -0.145
                                         0.023
                          0.064 -2.28
                                                -0.27
                                                        -0.02
```

score = $\alpha + \beta_1$ (rank_{tenure} track) + β_2 (rank_{tenured}) + ϵ

score = 4.28 - 0.13(rank_{tenure track}) - 0.15(rank_{tenured}) + ϵ

Score for teaching faculty

score = 4.28 - 0.13(rank_{tenure track}) - 0.15(rank_{tenured}) + ϵ

score =
$$4.28 - 0.13 \times 0 - 0.15 \times 0 + \epsilon$$

score = 4.28

Score for tenured faculty

score = 4.28 - 0.13(rank_{tenure track}) - 0.15(rank_{tenured}) + ϵ

score =
$$4.28 - 0.13 \times 0 - 0.15 \times 1 + \epsilon$$

score
$$= 4.28 - 0.15$$

4.13

Regression coefficients

Averages

```
term estimate std_
<chr> <chr> intercept 4.28
ranktenure track -0.13
ranktenured -0.145
```

```
rank avg
<fct> <fct> <dbl>
1 teaching 4.28
2 tenure track 4.15
3 tenured 4.14
```

Why do you think teaching = intercept?

Factors

Baseline category/intercept = lowest level in a factor variable

```
[457] tenure track tenure track
[463] tenure track
Levels: teaching tenure track tenured
```

When categorical variable *isn't* a factor, R picks one, typically first in alphabetical order

But you should set the factor a value that makes sense!

Interpretation

On average, y is \(\beta \) units larger (or smaller) in xn, compared to x0

On average, student evaluations are .13 points lower for tenure-track faculty than for teaching faculty

Compared to teaching faculty, tenured faculty get, on average, . 15 fewer points on teaching evals

A useful analogy: sliders and switches



score =
$$3.88 + 0.07$$
(bty_avg) + ϵ

This is a **SLOPE**



score =
$$4.28 - 0.13$$
(rank_{tenure track}) - 0.15 (rank_{tenured}) + ϵ

This is an OFFSET

Your turn

Regress score against a different categorical variable in *evals*

Interpret results for us

Most negative residual?

Most positive residual?

Interpreting residuals

TABLE 5.9: Regression points (First 10 out of 142 countries)

What does Afghanistan having a residual of -27 tell us?

country	lifeExp	continent	lifeExp_hat	residual
Afghanistan	43.8	Asia	70.7	-26.900
Albania	76.4	Europe	77.6	-1.226
Algeria	72.3	Africa	54.8	17.495
Angola	42.7	Africa	54.8	-12.075
Argentina	75.3	Americas	73.6	1.712
Australia	81.2	Oceania	80.7	0.515
Austria	79.8	Europe	77.6	2.180
Bahrain	75.6	Asia	70.7	4.907
Bangladesh	64.1	Asia	70.7	-6.666
Belgium	79.4	Europe	77.6	1.792