

Unit 1: Linear models: Estimation

Eugene Katsevich

September 7, 2021

1 Types of predictors; interpreting linear model coefficients (Agresti 1.2)

The types of predictors x_j (e.g. binary or continuous) has less of an effect on the regression than the type of response, but it is still important to pay attention to the former.

Intercepts. It is common to include an *intercept* in a linear regression model, a predictor x_0 such that $x_{i0} = 1$ for all i . When an intercept is present, we index it as the 0th predictor. The simplest kind of linear model is the *intercept-only model* or the *one-sample model*:

$$y = \beta_0 + \epsilon. \quad (1)$$

The parameter β_0 is the mean of the response.

Binary predictors. In addition to an intercept, suppose we have a binary predictor $x_1 \in \{0, 1\}$ (e.g. $x_1 = 1$ for patients who took blood pressure medication and $x_1 = 0$ for those who didn't). This leads to the following linear model:

$$y = \beta_0 + \beta_1 x_1 + \epsilon. \quad (2)$$

Here, β_0 is the mean response (say blood pressure) for observations with $x_1 = 0$ and $\beta_0 + \beta_1$ is the mean response for observations with $x_1 = 1$. Therefore, the parameter β_1 is the difference in mean response between observations with $x_1 = 1$ and $x_1 = 0$. This parameter is sometimes called the *effect* or *effect size* of x_1 , though a causal relationship might or might not be present. The model (2) is sometimes called the *two-sample model*, because the response data can be split into two “samples”: those corresponding to $x_1 = 0$ and those corresponding to $x_1 = 1$.

Categorical predictors. A binary predictor is a special case of a categorical predictor: A predictor taking two or more discrete values. Suppose we have a predictor $w \in \{w_0, w_1, \dots, w_{C-1}\}$, where $C \geq 2$ is the number of categories and w_0, \dots, w_{C-1} are the *levels* of w . E.g. suppose $\{w_0, \dots, w_{C-1}\}$ is the collection of U.S. states, so that $C = 50$. If we want to regress a response on the categorical predictor w , we cannot simply set $x_1 = w$ in the context of the linear regression (2). Indeed, w does not necessarily take numerical values. Instead, we need to add a predictor x_j for each of the levels of w . In particular, define $x_j \equiv \mathbb{1}(w = w_j)$ for $j = 1, \dots, C - 1$ and consider the regression

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_{C-1} x_{C-1} + \epsilon. \quad (3)$$

Here, category 0 is the *base category*, and β_0 represents the mean response in the base category. The coefficient β_j represents the difference in mean response between the j th category and the base category.

Quantitative predictors. A quantitative predictor is one that can take on any real value. For example, suppose that $x_1 \in \mathbb{R}$, and consider the linear model

$$y = \beta_0 + \beta_1 x_1 + \epsilon. \quad (4)$$

Now, the interpretation of β_1 is that an increase in x_1 by 1 is associated with an increase in y by β_1 . We must be careful to avoid saying “an increase in x_1 by 1 *causes* y to increase by β_1 ” unless we make additional causal assumptions. Note that the units of x_1 matter. If x_1 is the height of a person, then the value and the interpretation of β_1 changes depending on whether that height is measured in feet or in meters.

Ordinal predictors. There is an awkward category of predictor in between categorical and continuous called *ordinal*. An ordinal predictor is one that takes a discrete number of values, but these values have an intrinsic ordering, e.g. $x_1 \in \{\text{small}, \text{medium}, \text{large}\}$. It can be treated as categorical at the cost of losing the ordering information, or as continuous if one is willing to assign quantitative values to each category.

Multiple predictors. A linear regression need not contain just one predictor (aside from an intercept). For example, let’s say x_1 and x_2 are two predictors. Then, a linear model with both predictors is

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon. \quad (5)$$

When there are multiple predictors, the interpretation of coefficients must be revised somewhat. For example, β_1 in the above regression is the effect of an increase in x_1 by 1 *while holding x_2 constant* or *while adjusting for x_2* or *while controlling for x_2* . If y is blood pressure, x_1 is a binary predictor indicating blood pressure medication taken and x_2 is sex, then β_1 is the effect of the medication on blood pressure while controlling for sex. In general, the coefficient of a predictor depends on what other predictors are in the model. As an extreme case, suppose the medication has no actual effect, but that men generally have higher blood pressure and higher rates of taking the medication. Then, the coefficient β_1 in the single regression model (2) would be nonzero but the coefficient in the multiple regression model (5) would be equal to zero. In this case, sex acts as a *confounder*.

Interactions. Note that the multiple regression model (5) has the built-in assumption that the effect of x_1 on y is the same for any fixed value of x_2 (and vice versa). In some cases, the effect of one variable on the response may depend on the value of another variable. In this case, it’s appropriate to add another predictor called an *interaction*. Suppose x_1 is quantitative (e.g. years of job experience) and x_2 is binary (e.g. sex, with $x_2 = 1$ meaning male). Then, we can define a third predictor x_3 as the product of the first two, i.e. $x_3 = x_1 x_2$. This gives the regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon. \quad (6)$$

Now, the effect of adding another year of job experience is β_1 for females and $\beta_1 + \beta_3$ for males. The coefficient β_3 is the difference in the effect of job experience between males and females.

2 Model matrices, model vectors spaces, and identifiability (Agresti 1.3-1.4)

The matrix \mathbf{X} is called the *model matrix* or the *design matrix*. Concatenating the linear model equations across observations give us an equivalent formulation:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}; \quad \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}, \text{ Var}[\boldsymbol{\epsilon}] = \sigma^2 \mathbf{I}_n$$

or

$$\mathbb{E}[\mathbf{y}] = \mathbf{X}\boldsymbol{\beta} = \boldsymbol{\eta}.$$

As $\boldsymbol{\beta}$ varies in \mathbb{R}^p , the set of possible vectors $\boldsymbol{\eta} \in \mathbb{R}^n$ is defined

$$C(\mathbf{X}) \equiv \{\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^p\}.$$

$C(\mathbf{X})$, called the *model vector space*, is a subspace of \mathbb{R}^n : $C(\mathbf{X}) \subseteq \mathbb{R}^n$. Since

$$\mathbf{X}\boldsymbol{\beta} = \beta_1 \mathbf{x}_{*1} + \cdots + \beta_p \mathbf{x}_{*p},$$

the model vector space is the column space of the matrix \mathbf{X} .

The *dimension* of $C(\mathbf{X})$ is the rank of \mathbf{X} , i.e. the number of linearly independent columns of \mathbf{X} . If $\text{rank}(\mathbf{X}) < p$, this means that there are two different vectors $\boldsymbol{\beta}$ and $\boldsymbol{\beta}'$ such that $\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta}'$. Therefore, we have two values of the parameter vector that give the same model for \mathbf{y} . This makes $\boldsymbol{\beta}$ *not identifiable*, and makes it impossible to reliably determine $\boldsymbol{\beta}$ based on the data. For this reason, we will generally assume that $\boldsymbol{\beta}$ is *identifiable*, i.e. $\mathbf{X}\boldsymbol{\beta} \neq \mathbf{X}\boldsymbol{\beta}'$ if $\boldsymbol{\beta} \neq \boldsymbol{\beta}'$. This is equivalent to the assumption that $\text{rank}(\mathbf{X}) = p$. Note that this cannot hold when $p > n$, so for the majority of the course we will assume that $p \leq n$. In this case, $\text{rank}(\mathbf{X}) = p$ if and only if \mathbf{X} has *full-rank*.

As an example when $p \leq n$ but when $\boldsymbol{\beta}$ is still not identifiable, consider the case of a categorical predictor. Suppose the categories of w were $\{w_1, \dots, w_{C-1}\}$, i.e. the baseline category w_0 did not exist. In this case, the model (3) would not be identifiable because $x_0 = 1 = x_1 + \cdots + x_{C-1}$ and thus $x_{*0} = 1 = x_{*1} + \cdots + x_{*,C-1}$. Indeed, this means that one of the predictors can be expressed as a linear combination of the others, so \mathbf{X} cannot have full rank. A simpler way of phrasing the problem is that we are describing $C - 1$ intrinsic parameters (the means in each of the $C - 1$ groups) with C model parameters. There must therefore be some redundancy. For this reason, if we include an intercept term in the model then we must designate one of our categories as the baseline and exclude its indicator from the model.

3 Least squares estimation (Agresti 2.1.1, 2.7.1)

Now, suppose that we are given a dataset (\mathbf{X}, \mathbf{y}) . How do we go about estimating $\boldsymbol{\beta}$ based on this data? The canonical approach is the *method of least squares*:

$$\hat{\boldsymbol{\beta}} \equiv \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2. \quad (7)$$

The quantity

$$\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 = \|\mathbf{y} - \hat{\boldsymbol{\mu}}\|^2 = \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 \quad (8)$$

is called the *residual sum of squares* (*RSS*), and it measures the lack of fit of the linear regression model. We therefore want to choose $\hat{\boldsymbol{\beta}}$ to minimize this lack of fit. Note that if $\boldsymbol{\epsilon}$ is assumed to be $N(0, \sigma^2 \mathbf{I}_n)$, then the least squares solution would also be the maximum likelihood solution. Indeed, for $y_i \sim N(\mu_i, \sigma^2)$, the log-likelihood is

$$\log \left[\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y_i - \mu_i)^2}{2\sigma^2} \right) \right] = \text{constant} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_i)^2.$$

Letting $L(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$, we can do some calculus to derive that

$$\frac{\partial}{\partial \boldsymbol{\beta}} L(\boldsymbol{\beta}) = -\mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \quad (9)$$

Setting this vector of partial derivatives equal to zero, we arrive at the *normal equations*:

$$-\mathbf{X}^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0 \iff \mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{y}. \quad (10)$$

If \mathbf{X} is full rank, the matrix $\mathbf{X}^T\mathbf{X}$ is invertible and we can therefore conclude that

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}. \quad (11)$$

Now that we have derived the least squares estimator, we can compute its bias and variance. To obtain the bias, we first calculate that

$$\mathbb{E}[\hat{\boldsymbol{\beta}}] = \mathbb{E}[(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}] = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbb{E}[\mathbf{y}] = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta}.$$

Therefore, the least squares estimator is unbiased. To obtain the variance, we compute

$$\begin{aligned} \text{Var}[\hat{\boldsymbol{\beta}}] &= \text{Var}[(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}] \\ &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\text{Var}[\mathbf{y}]\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} \\ &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T(\sigma^2\mathbf{I}_n)\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} \\ &= \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}. \end{aligned}$$

According to the Gauss-Markov theorem, this covariance matrix computed above is the smallest (in the sense of positive semidefinite matrices) among all linear unbiased estimates of $\boldsymbol{\beta}$.

4 Least squares estimation R demo (Agresti 2.6)

The R demo will be based on the `ScotsRaces` data from the textbook. Data description (quoted from the textbook):

“Each year the Scottish Hill Runners Association publishes a list of hill races in Scotland for the year. The table below shows data on the record time for some of the races (in minutes). Explanatory variables listed are the distance of the race (in miles) and the cumulative climb (in thousands of feet).”

```
library(tidyverse)
```

```
# read the data into R
scots_races = read_tsv("../data/ScotsRaces.dat", col_types = "cddd")
scots_races

## # A tibble: 35 x 4
##   race                distance climb  time
##   <chr>              <dbl> <dbl> <dbl>
## 1 GreenmantleNewYearDash    2.5  0.65  16.1
## 2 Carnethy5HillRace         6    2.5   48.4
## 3 CraigDunainHillRace       6    0.9   33.6
## 4 BenRhaHillRace           7.5  0.8   45.6
## 5 BenLomondHillRace         8    3.07  62.3
## 6 GoatfellHillRace         8    2.87  73.2
```

```
## 7 BensofJuraFellRace      16    7.5  205.  
## 8 CairnpappleHillRace     6     0.8  36.4  
## 9 ScoltyHillRace          5     0.8  29.8  
## 10 TraprainLawRace        6     0.65 39.8  
## # ... with 25 more rows
```

```
# Exploration
```

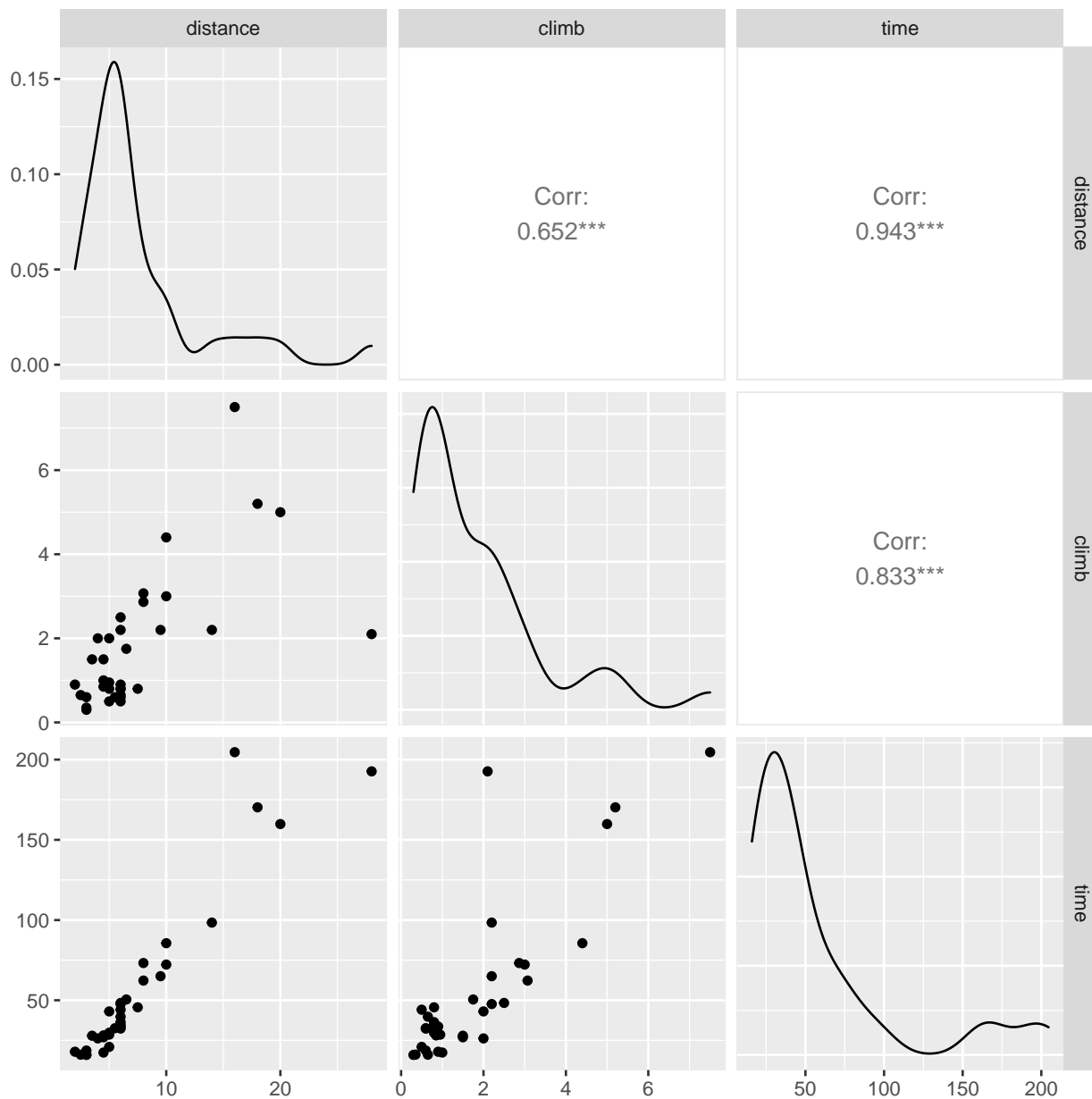
```
# pairs plot
```

```
GGally::ggpairs(scots_races %>% select(-race))
```

```
## Registered S3 method overwritten by 'GGally':
```

```
##   method from
```

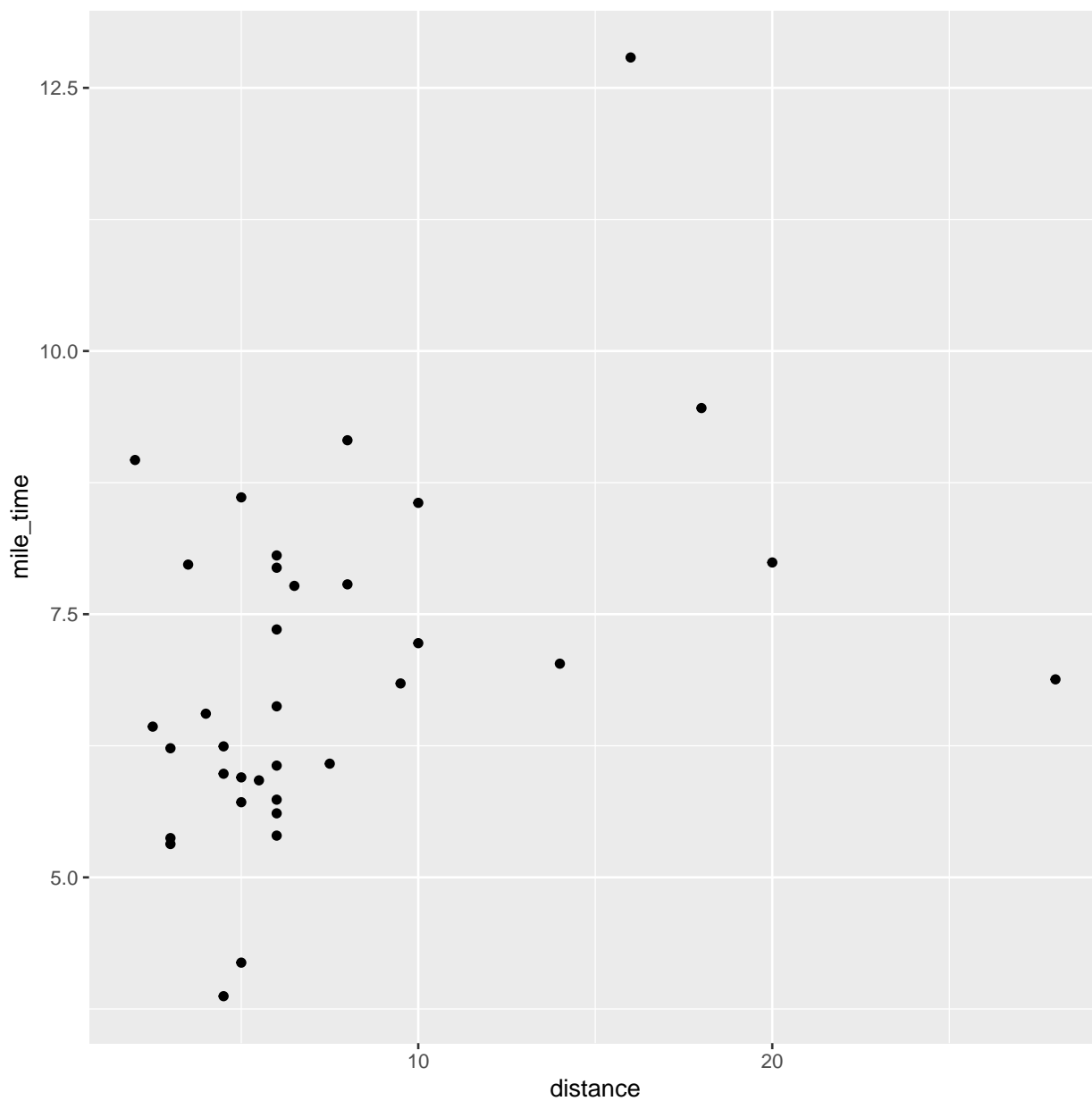
```
##   +.gg   ggplot2
```



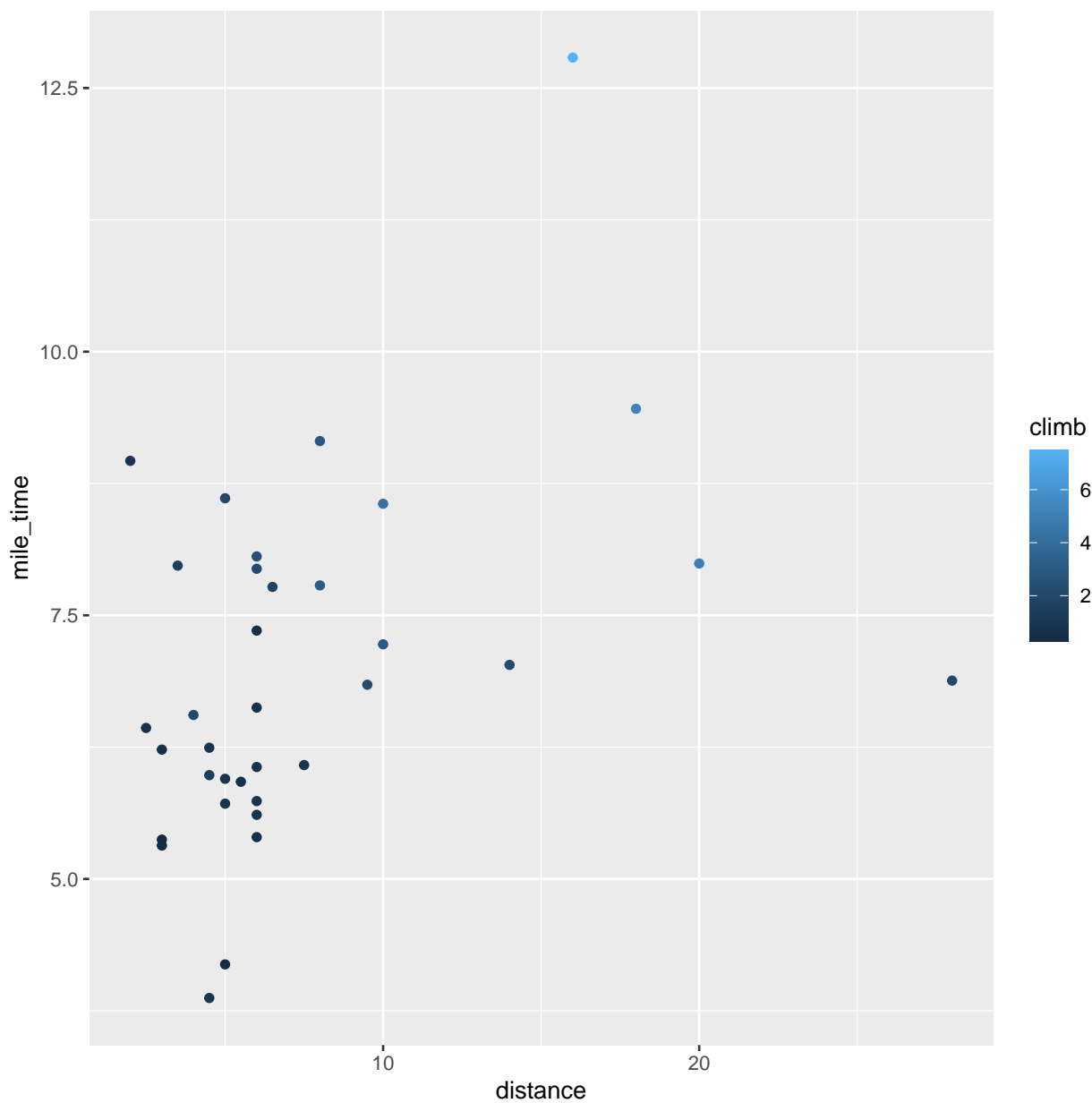
Q: What are the typical ranges of the variables?
 # Q: What are the relationships among the variables?

mile time versus distance

```
scots_races %>%
  mutate(mile_time = time/distance) %>%
  ggplot(aes(x = distance, y = mile_time)) +
  geom_point()
```



```
scots_races %>%  
  mutate(mile_time = time/distance) %>%  
  ggplot(aes(x = distance, y = mile_time, label = race,  
             color = climb)) +  
  geom_point()
```



```
# Q: How does mile time vary with distance?
# Q: What races deviate from this trend?
# Q: How does climb play into it?
```

```
# Linear model
```

```
# Q: What is the effect of an extra mile of distance on time?
lm_fit = lm(time ~ distance + climb, data = scots_races)
summary(lm_fit)
```

```
##
## Call:
```



```
## lm(formula = time ~ distance + climb, data = scots_races)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.654  -4.842   1.110   4.667  27.762
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -13.1086     2.5608  -5.119 1.41e-05 ***
## distance       6.3510     0.3578  17.751 < 2e-16 ***
## climb        11.7801     1.2206   9.651 5.37e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.734 on 32 degrees of freedom
## Multiple R-squared:  0.9717, Adjusted R-squared:  0.97
## F-statistic: 549.9 on 2 and 32 DF,  p-value: < 2.2e-16
```

```
# Linear model with interaction
```

```
# Q: What is the effect of an extra mile of distance on time
# for a run with low climb?
```

```
# Q: What is the effect of an extra mile of distance on time
# for a run with high climb?
```

```
lm_fit_int = lm(time ~ distance*climb, data = scots_races)
summary(lm_fit_int)
```

```
##
## Call:
## lm(formula = time ~ distance * climb, data = scots_races)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -23.197  -2.797   0.628   2.243  18.963
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -0.7672     3.9058  -0.196  0.84556
## distance       4.9623     0.4742  10.464 1.07e-11 ***
## climb         3.7133     2.3647   1.570  0.12650
## distance:climb  0.6598     0.1743   3.786  0.00066 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.338 on 31 degrees of freedom
```

```
## Multiple R-squared:  0.9807, Adjusted R-squared:  0.9788
## F-statistic: 524.1 on 3 and 31 DF,  p-value: < 2.2e-16
```

5 Linear regression as orthogonal projection (Agresti 2.2, 2.3, 2.4.2, 2.4.3, 2.4.4)

Let's think about the mapping $\mathbf{y} \mapsto \hat{\boldsymbol{\mu}} = \mathbf{X}\hat{\boldsymbol{\beta}} \in C(\mathbf{X})$. We claim that this mapping is an *orthogonal projection*. Geometrically it makes sense, since we define $\hat{\boldsymbol{\beta}}$ so that $\hat{\boldsymbol{\mu}} \in C(\mathbf{X})$ is as close to \mathbf{y} as possible. The shortest path between a point and a plane is the perpendicular. One way of seeing this is to show that $\mathbf{v}^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0$ for each $\mathbf{v} \in C(\mathbf{X})$. Since the columns $\{\mathbf{x}_{*1}, \dots, \mathbf{x}_{*p}\}$ of \mathbf{X} form a basis for $C(\mathbf{X})$, it suffices to show that $\mathbf{x}_{*j}^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0$ for each $j = 1, \dots, p$. This is a consequence of the normal equations $\mathbf{X}^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0$ derived in (10).

To derive the projection matrix corresponding to this orthogonal projection, we write

$$\hat{\boldsymbol{\mu}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{H}\mathbf{y}, \quad (12)$$

where

$$\mathbf{H} \equiv \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \quad (13)$$

is called the *hat matrix*. This is the orthogonal projection matrix onto $C(\mathbf{X})$. Recall that a matrix \mathbf{P} is an orthogonal projection onto a subspace \mathbf{W} if for all $\mathbf{v} \in \mathbf{W}$ we have $\mathbf{P}\mathbf{v} = \mathbf{v}$ and for all $\mathbf{v} \in \mathbf{W}^\perp$ we have $\mathbf{P}\mathbf{v} = 0$. We can check for example the first of these conditions by noting that if $\mathbf{v} \in C(\mathbf{X})$, then $\mathbf{v} = \mathbf{X}\boldsymbol{\beta}$ for some $\boldsymbol{\beta} \in \mathbb{R}^p$. Therefore, we have

$$\mathbf{H}\mathbf{v} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta} = \mathbf{v}.$$

A simple example of \mathbf{H} can be obtained by considering the intercept-only regression.

One consequence of this observation is that the fitted values $\hat{\boldsymbol{\mu}}$ depend on \mathbf{X} only through $C(\mathbf{X})$. As we will see in Homework 1, there are many different model matrices \mathbf{X} leading to the same model space. Essentially, this reflects the fact that there are many different bases for the same vector space. Consider for example changing the units on the columns of \mathbf{X} . It can be verified that not just the fitted values $\hat{\boldsymbol{\mu}}$ but also the predictions on a new set of features remain invariant to reparametrization (this follows from parts (a) and (b) of Homework 1 Problem 1). Therefore, while reparametrization can have a huge impact on the fitted coefficients, it has no impact on the predictions of linear regression.

The orthogonality property of least squares, together with the Pythagorean theorem, leads to the following fundamental relationship. Let's say that $S \subset \{1, \dots, p\}$ is a subset of the predictors. First regress \mathbf{y} on \mathbf{X} to get $\hat{\boldsymbol{\beta}}$ as usual. Then, we consider the *partial model matrix* \mathbf{X}_{*S} obtained by selecting only the columns in S . Regression \mathbf{y} on \mathbf{X}_{*S} results in $\hat{\boldsymbol{\beta}}_S$ (note: $\hat{\boldsymbol{\beta}}_S$ is not necessarily obtained from $\hat{\boldsymbol{\beta}}$ by extracting the coefficients corresponding to S). Now, consider the three points $\mathbf{y}, \mathbf{X}\hat{\boldsymbol{\beta}}, \mathbf{X}_{*S}\hat{\boldsymbol{\beta}}_S \in \mathbb{R}^n$. Since $\mathbf{X}\hat{\boldsymbol{\beta}}$ and $\mathbf{X}_{*S}\hat{\boldsymbol{\beta}}_S$ are both in $C(\mathbf{X})$, it follows by the orthogonal projection property that $\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is orthogonal to $\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}_{*S}\hat{\boldsymbol{\beta}}_S$. In other words, these three points form a right triangle. By the Pythagorean theorem, we conclude that

$$\|\mathbf{y} - \mathbf{X}_{*S}\hat{\boldsymbol{\beta}}_S\|^2 = \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 + \|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}_{*S}\hat{\boldsymbol{\beta}}_S\|^2. \quad (14)$$

We will rely on this fundamental relationship throughout this course.

For now, we can extract a few consequences of the relationship (14). As a starting point, consider the case when $S = \{0\}$, i.e. the partial model is the intercept-only model. In this case, $\mathbf{X}_{*S} = \mathbf{1}_n$ and $\hat{\beta}_S = \bar{y}$. Therefore, equation (14) implies that

$$\|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2 = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 + \|\mathbf{X}\hat{\beta} - \bar{y}\mathbf{1}_n\|^2. \quad (15)$$

Equivalently, we can rewrite this equation as follows:

$$\text{TSS} \equiv \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 \equiv \text{SSR} + \text{SSE}. \quad (16)$$

- R^2 , and what happens as you add more predictors
- Special case: One-way ANOVA

6 Correlation, multiple correlation, and R^2 (Agresti 2.1.3, 2.4.6)

TBD.

7 Collinearity, adjustment, and partial correlation (Agresti 2.2.4, 2.5.6, 2.5.7)

TBD.

8 R demo

TBD.