

Problem 1 Votes

In this problem, we are a polling company trying to predict the vote in the 2016 Republican primary election in California. (In a party's primary election, people vote to determine the party's candidate for the actual election, which happens later.) Millions of people will vote in the election. We select a simple random sample of 200 people likely to vote in the Republican primary, and we ask who they will vote for.

We're interested in predicting the proportion of people who will vote for Donald Trump. Let's say that that proportion is *actually* .2. Our polling company doesn't know that, though.

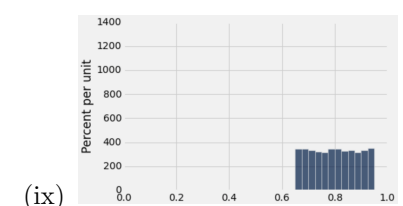
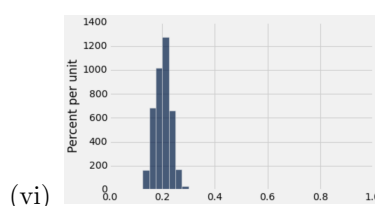
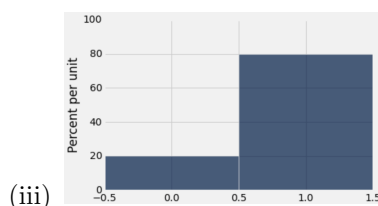
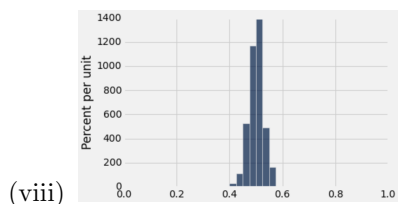
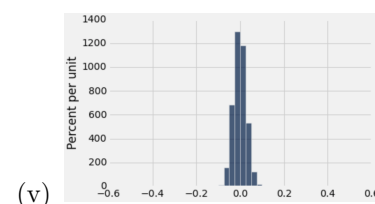
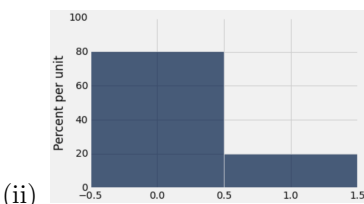
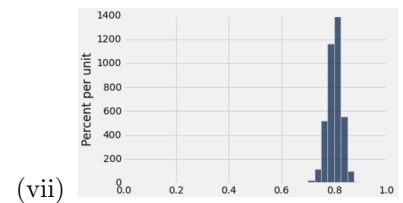
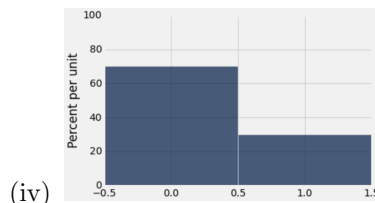
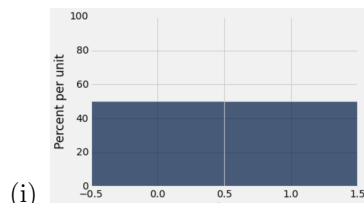
Suppose there is a table of all the electoral votes called `all_votes` and a table of our sample of 200 called `sample_votes`:

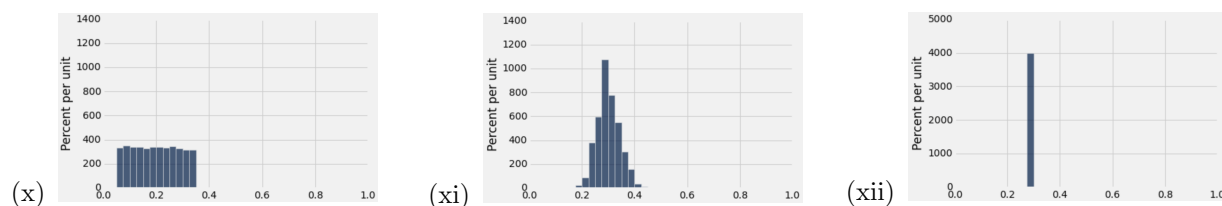
```
N = 200
all_votes = Table.read_table("all_votes.csv")
sample_votes = all_votes.sample(N)
sample_votes.show(10)
```

Name
Ted Cruz
Ted Cruz
Donald Trump
John Kasich
Marco Rubio
Ted Cruz
Marco Rubio
John Kasich
Ben Carson
Donald Trump

... (190 rows omitted)

For each of the scenarios below, we've described a histogram made from this data, or written some code that manipulates the data and makes a histogram. For each one, find in the list below *all* the histograms that would not be surprising results for the computations described.





Hint: Pay attention to the *mean* and *shape* of each histogram. The *spread* of each histogram might be hard to compute, but you can compare the spreads of similarly-shaped histograms to each other.

- (a) Suppose we mark each vote as a 1 if it's a vote for Trump and a 0 otherwise, and then make a histogram of those numbers from `all_votes`, using this code:

```
def one_if_trump(candidate):
    if candidate == "Donald Trump":
        return 1
    else:
        return 0

votes_with_trump_indicator = all_votes.with_column(
    "Trump",
    all_votes.apply(one_if_trump, "Name"))
votes_with_trump_indicator.hist("Trump", bins=[-0.5, 0.5, 1.5])
```

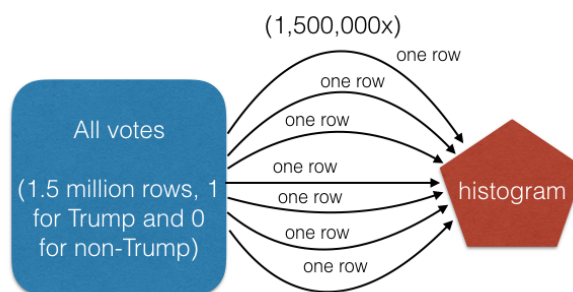


Figure 1: A visual depiction of (a).

- (b) We do the same thing, but instead of `all_votes`, we use the smaller sample our polling company actually sees.

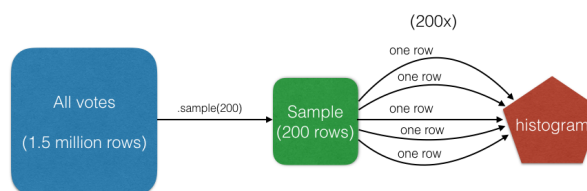


Figure 2: A visual depiction of (b).

- (c) Now we'd like to understand how our sample-based estimate of the proportion of Trump votes would "typically" (that is, across many re-runs of the sampling process) behave. We produce an empirical histogram of the proportions of Trump votes in 1000 different random samples of size 200. (If you're not sure how to write the code to do that, it would be a very useful exercise to try.)

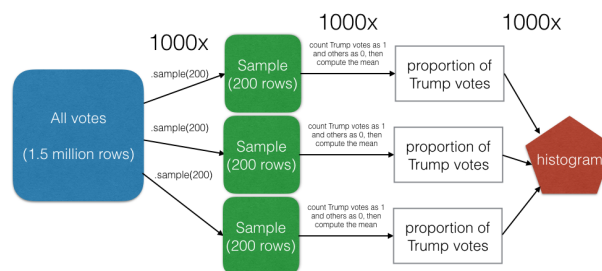


Figure 3: A visual depiction of (c).

- (d) Now we'd like to understand how the sample-based estimate of the proportion of *non-Trump* votes would typically behave. We produce an empirical histogram of the proportions of *non-Trump* votes in 1000 different random samples of size 200.
- (e) A political scientist on our staff believes that polls underestimate Trump's support because some Trump voters are reluctant to tell pollsters they're voting for Trump. (This is a real idea.) To correct for this, we estimate the *true proportion of Trump voters in California* by multiplying the proportion of Trump voters in our poll by 1.5. We produce an empirical histogram of these new estimated proportions of Trump votes (that is, the sample proportion of Trump votes times 1.5) from 1000 different random samples of size 200.

Answer:

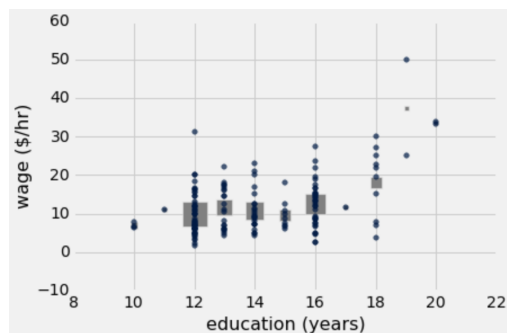
- (a) The only reasonable histogram is (ii), since 20% of the population are Trump voters.
- (b) The only reasonable histograms are (ii) and (iv). (i) and (iii) are possible but unlikely. That's not obvious unless you know how "typically variable" a sample of 200 is. The next question answers that question, and the histogram shows that getting half or 80% Trump voters is very unlikely, but getting 20% or 30% Trump voters is very possible.
- (c) The only reasonable histogram is (vi). That histogram is centered around 0.2 (since a proportion is just a mean, and the mean of a random sample is *on average* the mean of the population) and bell-shaped (which the Central Limit Theorem says it should be). So, for example, it can't be (v) or (viii) because it should be centered around 0.2, and it can't be (x) because it should be bell-shaped.
- (d) This is the same as the previous question, but centered around 0.8, since 0.8 of the population is non-Trump voters. So the only reasonable histogram is (vii), for the same reasons as in the previous question.
- (e) The only reasonable histogram is (xi). Take the histogram from (c) as a baseline and consider how the answer here should differ. Multiplying all the values in a distribution by a number doesn't change the shape of the distribution, so it should still be bell-shaped. The average of a bunch of numbers multiplied by 1.5 is 1.5 times their average, so the bell should be centered at $1.5 \times 0.2 = 0.3$. That leaves (xi) and (xii).
(xi) is shaped like (vi) but more spread out, and (xii) is like (vi) but much less spread out. So the remaining question is whether multiplying each value by 1.5 increases or decreases the distribution's spread. The answer is that it makes it more spread out. For example, the spread of the numbers $[-2, 0, 4]$ is twice the spread of the numbers $[-1, 0, 2]$. So the answer must be (xi).

Problem 2 Education

In a 1994 study, economists Ashenfelter and Krueger analyzed the *financial returns to education* – that is, how much getting more school increases income. They were interested in comparing twins, so they collected

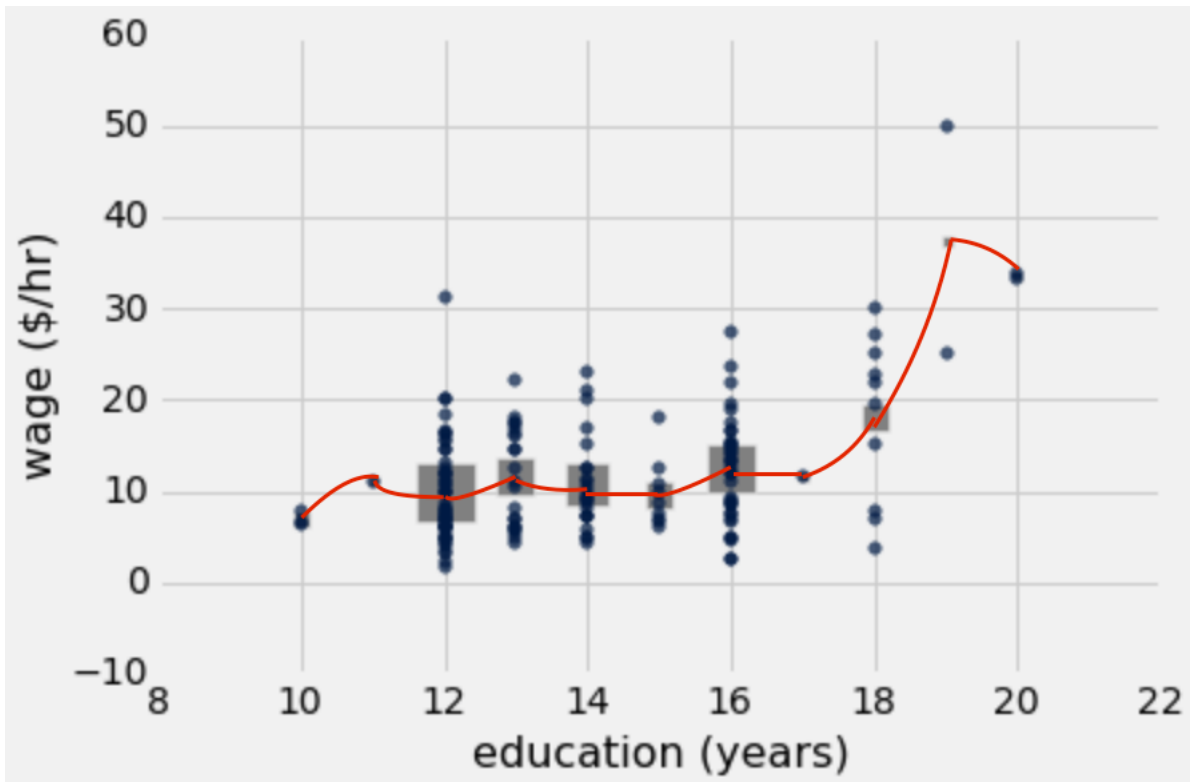
their data at the 16th Annual Twins Days Festival in Twinsburg, Ohio. We'll do a simplified version of their analysis.

Here's a graph of education (measured in years spent in school) and hourly wages for all the people in the study. The light squares denote the *mean* of each vertical slice. (The area of each square is proportional to the number of people in that slice.)



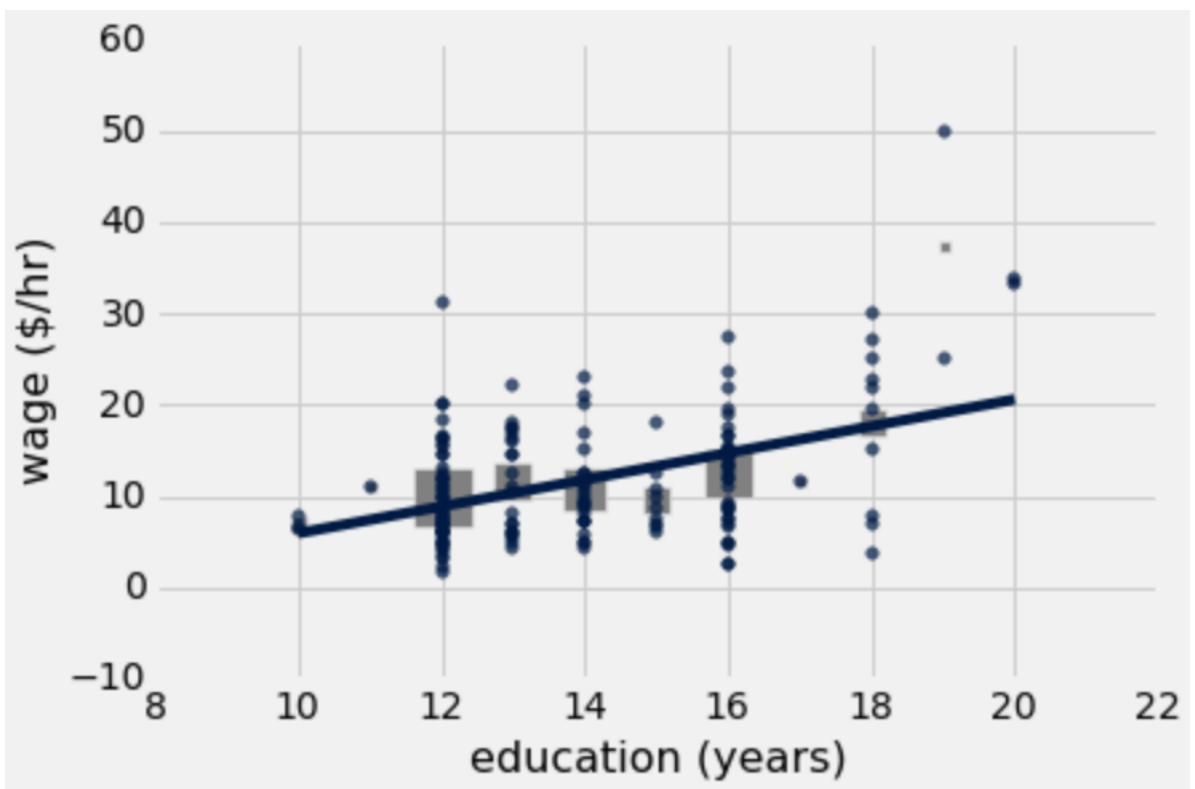
- (a) We'd like to model the relationship between education and hourly wages with a smooth curve. Suppose we don't restrict ourselves to fitting lines, and we allow any curve. Draw a curve that fits the data as well as possible. (If you'd like the goal to be stated more precisely: Draw a curve so that the squared vertical distance between the curve and each point is minimized.)
- (b) True or false: Your curve passes through the average of each vertical strip.
- (c) Now draw the *line* that fits the data as well as possible.
- (d) True or false: Your line passes through the average of each vertical strip.
- (e) Is the correlation between education and wages positive, negative, or zero? What's your best guess at the correlation?
- (f) Would you say there's a roughly linear relationship between education and wages in these data? If not, how would you characterize the relationship?
- (g) *[Optional]* It's important to note that these data only show an *association* between education and wages. Showing causation is much harder. For example, maybe we see this relationship because people with the perseverance to stay in school apply the same perseverance to get higher-paying jobs. That's an example of a *confounding factor*. Describe how knowing about the *association by itself* could be useful.
- (h) *[Optional]* Describe how knowing whether more education *causes* higher wages could be useful.
- (i) *[Optional]* Can you guess how Ashenfelter and Krueger used the fact that their dataset actually contained many pairs of *genetically-identical twins* to eliminate confounding factors like that?

Answer:



(a)

- (b) True. Since a curve can have only one value at each vertical strip, it can't fit all the data in a strip precisely. Instead, the best it can do is pass through the *mean* of each strip.



(c)

Notice that the line is pulled up a little toward the far-right values, but not all the way. That's because

there aren't very many people in the dataset with high education, so fitting those points is less important.

- (d) False. The means don't fall exactly on a line, so it's impossible for the line to go through each one. That means the line doesn't fit the data as well as the free curve. The simplicity of the line (versus our arbitrary-looking curve) is attractive, though.
- (e) Definitely positive. The line isn't a great fit, though, so our guess shouldn't be close to 1. The staff guess was .5. The actual correlation turned out to be 0.4383. (You don't need to be able to eyeball these things, but if you stare at enough ovals you might get a knack for it.)
- (f) Not really. It looks like wages are more-or-less flat except for two bumps: college graduates (typically 16 years of education) have higher wages than people with some or no college, and people with a bit of post-college education have higher wages than the college-educated.
- (g) Perhaps the government wants to design a social program to help people newly entering the workforce who will be earning low wages 10 years from now. It doesn't know how much people will earn in 10 years, but it does know how much education they have. Then it could target the program at people with low education as a *proxy* for low future wages. For this application we only need to *predict a value* rather than *predict the outcome of an intervention*, so it doesn't matter whether education causes low wages or is just associated with them.
- (h) Say we are deciding whether to implement a program to incentivize people to stay in school longer. If more schooling actually causes higher wages, then our program is more likely to actually help people.
- (i) It's plausible that identical twins are similar on some confounding factors like perseverance. So if we compare two twins with different education, the differences in their wages might actually be *caused* by the difference in education. *Note:* we haven't learned a technique to do that comparison yet.