# Lab 7: Regression

Welcome to Lab 7!

Today we will get some hands-on practice with linear regression.

```
In [1]:   # Run this cell, but please don't change it.

          # These lines import the Numpy and Datascience modules.
          import numpy as np
          from datascience import *

          # These lines do some fancy plotting magic.
          import matplotlib
          %matplotlib inline
          import matplotlib.pyplot as plt
          plt.style.use('fivethirtyeight')
          import warnings
          warnings.simplefilter('ignore', FutureWarning)
```

# 1. How Faithful is Old Faithful?

(Note: clever title comes from here, a statistical discussion of Old Faithful.)

Old Faithful is a geyser in Yellowstone National Park in the central United States. It's famous for erupting on a fairly regular schedule. You can see a video below.

```
In [23]:  # For the curious: this is how to display a YouTube video in a
          # Jupyter notebook.  The argument to YouTubeVideo is the part
          # of the URL (called a "query parameter") that identifies the
          # video.  For example, the full URL for this video is:
          #   https://www.youtube.com/watch?v=wE8NDuzt8eg
          from IPython.display import YouTubeVideo
          YouTubeVideo("wE8NDuzt8eg")
```

Out[23]:



Old Faithful Geyser eruption Yellowsto...

Some of Old Faithful's eruptions last longer than others. When it has a long eruption, there's generally a longer wait until the next eruption.

If you visit Yellowstone, you might want to predict when the next eruption will happen, so you can see the rest of the park and come to see the geyser when it happens. Today, we will use a dataset on eruption durations and waiting times to see if we can make such predictions accurately with linear regression.

The dataset has one row for each observed eruption. It includes the following columns:

- **duration**: Eruption duration, in minutes
- **wait**: Time between this eruption and the next, also in minutes

Run the next cell to load the dataset.

```
In [2]:  faithful = Table.read_table("~/DS_113_S23/Labs/Lab_7/faithful.csv")
         faithful
```

Out[2]:

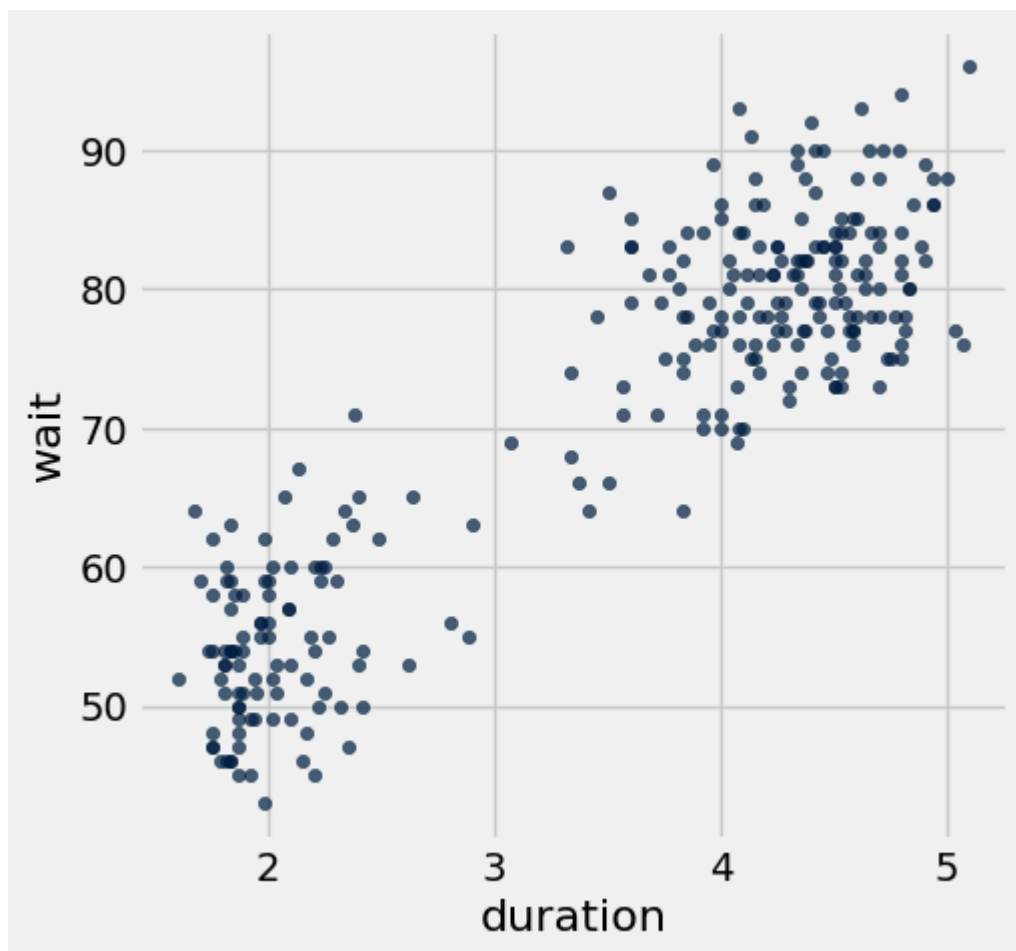| duration | wait |
| --- | --- |
| 3.6 | 79 |
| 1.8 | 54 |
| 3.333 | 74 |
| 2.283 | 62 |
| 4.533 | 85 |
| 2.883 | 55 |
| 4.7 | 88 |
| 3.6 | 85 |
| 1.95 | 51 |
| 4.35 | 85 |

... (262 rows omitted)

We would like to use linear regression to make predictions, but that won't work well if the data aren't roughly linearly related. To check that, we should look at the data.

## Question 1

Make a scatter plot of the data. It's conventional to put the column we will try to predict on the vertical axis and the other column on the horizontal axis.

```
In [7]:  faithful.scatter(0,1)
```

## Question 2

Look at the scatter plot. Are eruption duration and waiting time roughly linearly related? Is the relationship positive, as we claimed earlier?

The relationship is positive. We can see that for a longer eruption duration, most of the data show that the waiting time is also longer.

We're going to continue with the provisional assumption that they are linearly related, so it's reasonable to use linear regression to analyze this data.

We'd next like to plot the data in standard units. Recall that, if `nums` is an array of numbers, then

```
(nums - np.mean(nums)) / np.std(nums)
```

...is an array of those numbers in standard units.

## Question 3

Compute the mean and standard deviation of the eruption durations and waiting times. **Then** create a table called `faithful_standard` containing the eruption durations and waiting times in standard units. (The columns should be named `"duration (standard units)"` and `"wait (standard units)"`.

In [10]:
```python
duration_mean = np.mean(faithful["duration"])
duration_std = np.std(faithful["duration"])
wait_mean = np.mean(faithful["wait"])
wait_std = np.std(faithful["wait"])

duration_standard = (faithful["duration"] - duration_mean) / duration_std
wait_standard = (faithful["wait"] - wait_mean) / wait_std

faithful_standard = Table().with_columns(
    "duration (standard units)", duration_standard,
    "wait (standard units)", wait_standard)
faithful_standard
```

Out[10]:

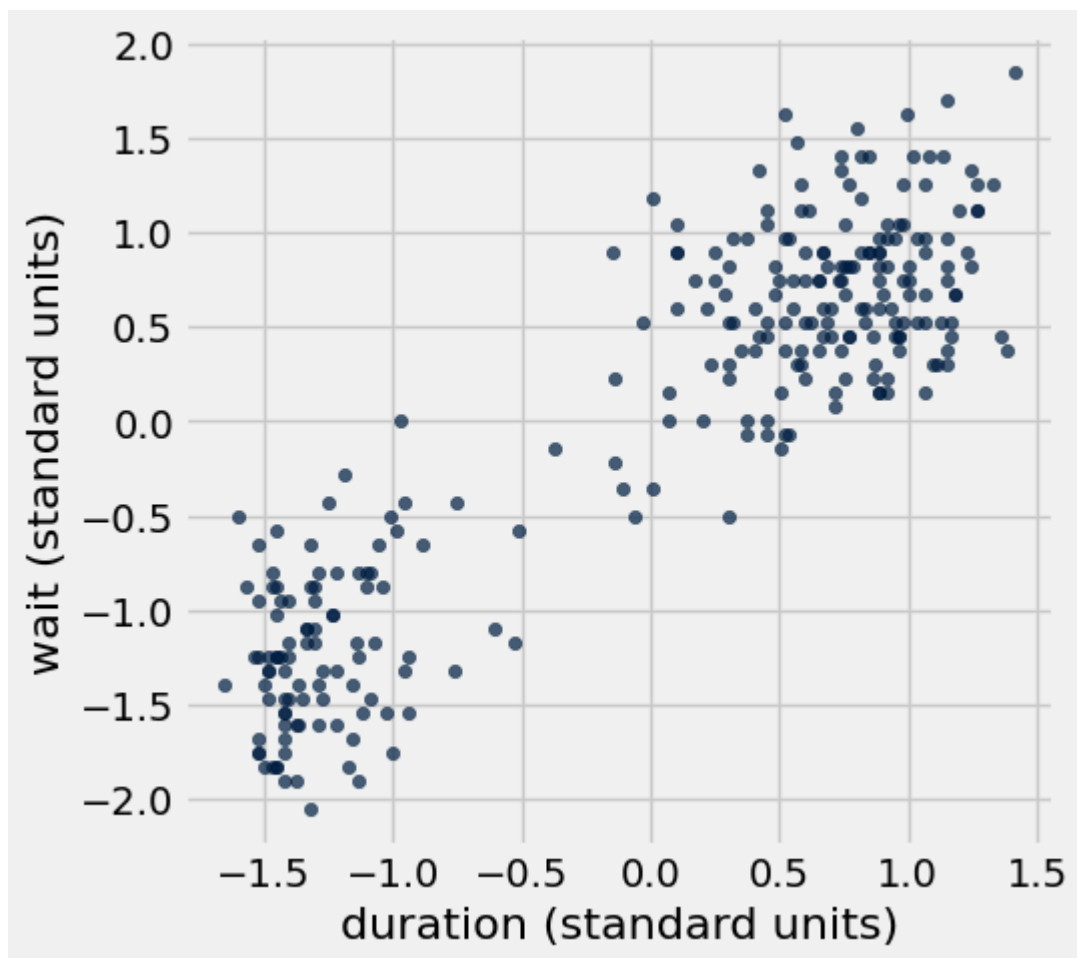| duration (standard units) | wait (standard units) |
| --- | --- |
| 0.0984989 | 0.597123 |
| -1.48146 | -1.24518 |
| -0.135861 | 0.228663 |
| -1.0575 | -0.655644 |
| 0.917443 | 1.03928 |
| -0.530851 | -1.17149 |
| 1.06403 | 1.26035 |
| 0.0984989 | 1.03928 |
| -1.3498 | -1.46626 |
| 0.756814 | 1.03928 |

... (262 rows omitted)

## Question 4

Plot the data again, but this time in standard units.

In [11]:
```python
faithful_standard.scatter(0,1)
```

You'll notice that this plot looks exactly the same as the last one! The data really are different, but the axes are scaled differently. (The method `scatter` scales the axes so the data fill up the available space.) So it's important to read the ticks on the axes.

## Question 5

Among the following numbers, which would you guess is closest to the correlation between eruption duration and waiting time in this dataset? Why?

- -1
- 0
- 1

1 ; We can see the positive relationship in the scatter plot and they are gathered around each other.

## Question 6

Compute the correlation `r`.

```
In [12]:  r = (faithful_standard.column(0) * faithful_standard.column(1)).mean()
          r
```

```
Out[12]:  0.90081116832181318
```

# 2. The regression line

Recall that the correlation is the slope of the regression line when the data are put in standard units.
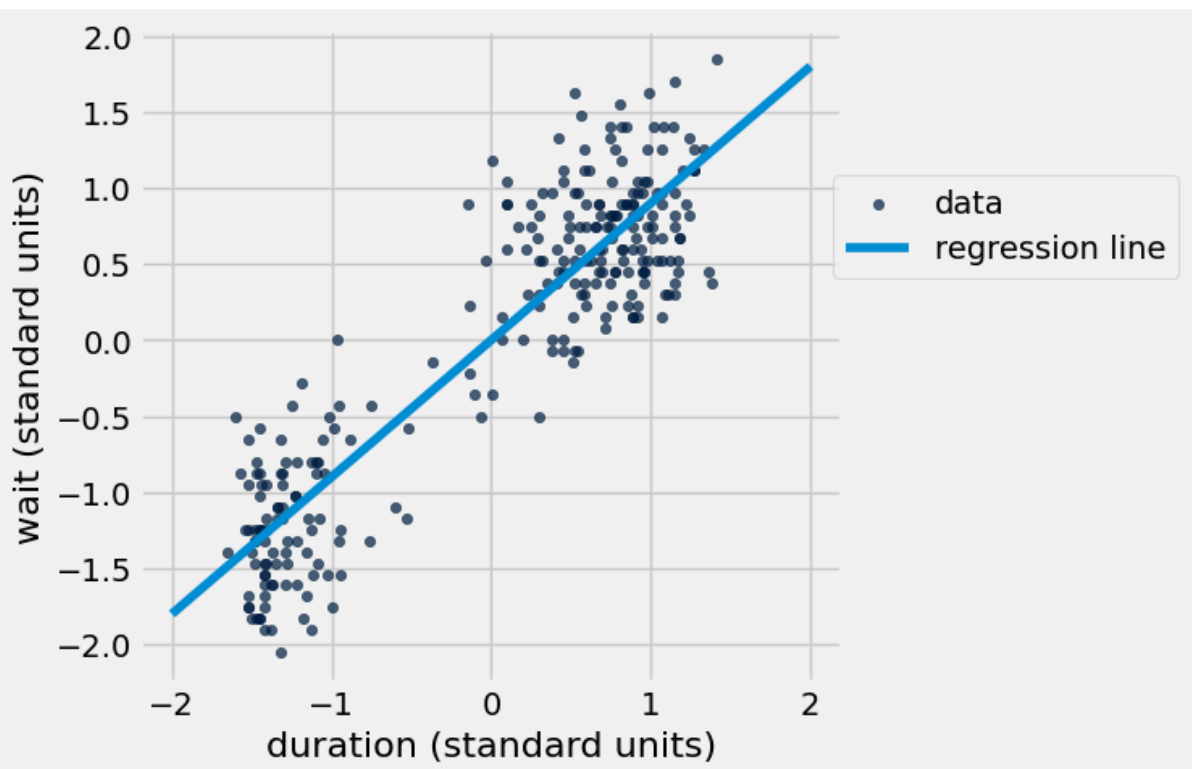
The next cell plots the regression line in standard units:

$$\text{waiting time (standard units)} = r \times \text{eruption duration (standard units)}.$$

Then, it plots the original data again, for comparison.

```
In [13]: def plot_data_and_line(dataset, x, y, point_0, point_1):
             """Makes a scatter plot of the dataset, along with a line passing throug
             dataset.scatter(x, y, label="data")
             plt.plot(make_array(point_0.item(0), point_1.item(0)), make_array(point_
             plt.legend(bbox_to_anchor=(1.5,.8))

         plot_data_and_line(faithful_standard, "duration (standard units)", "wait (st
```



How would you take a point in standard units and convert it back to original units? We'd have to "stretch" its horizontal position by `duration_std` and its vertical position by `wait_std`.

**That means the same thing would happen to the slope of the line.**

## Question 1

What is the slope of the regression line in original units?

```
In [14]: slope = r * faithful["wait"].std()/faithful["duration"].std()
         slope
```

```
Out[14]: 10.729641395133529
```

We know that the regression line passes through the point `(duration_mean, wait_mean)`. You might recall from high-school algebra that the equation for the line is

therefore:

$$\text{waiting time} - \texttt{wait\_mean} = \texttt{slope} \times (\text{eruption duration} - \texttt{duration\_mean})$$

After rearranging that equation slightly, the intercept turns out to be:

```
In [15]:  intercept = slope*(-duration_mean) + wait_mean
          intercept
```

```
Out[15]:  33.474397022753344
```

# 3. Investigating the regression line

The slope and intercept tell you exactly what the regression line looks like. To predict the waiting time for an eruption, multiply the eruption's duration by `slope` and then add `intercept`.

## Question 1

Compute the predicted waiting time for an eruption that lasts 2 minutes, and for an eruption that lasts 5 minutes.

```
In [16]:  two_minute_predicted_waiting_time = 2*slope + intercept
          five_minute_predicted_waiting_time = 5*slope + intercept

          # Here is a helper function to print out your predictions
          # (you don't need to modify it):
          def print_prediction(duration, predicted_waiting_time):
              print("After an eruption lasting", duration,
                    "minutes, we predict you'll wait", predicted_waiting_time,
                    "minutes until the next eruption.")

          print_prediction(2, two_minute_predicted_waiting_time)
          print_prediction(5, five_minute_predicted_waiting_time)
```
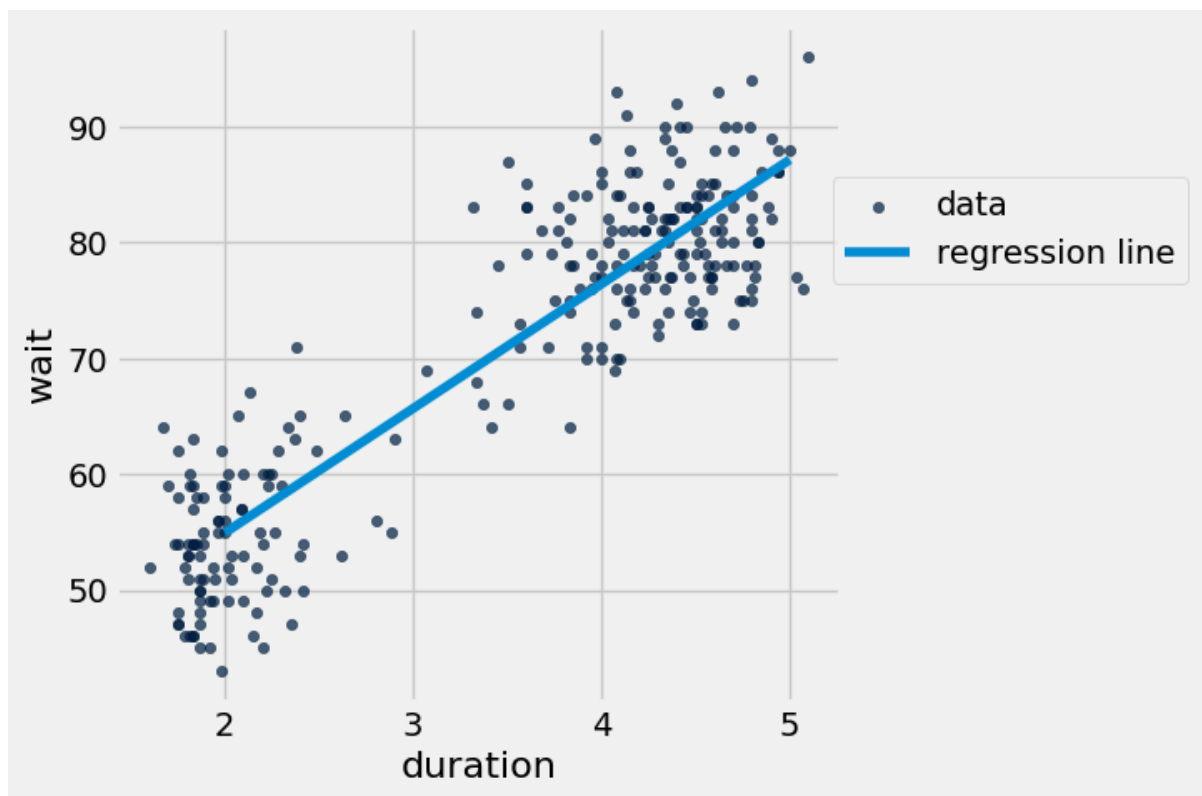
```
After an eruption lasting 2 minutes, we predict you'll wait 54.933679813 mi
nutes until the next eruption.
After an eruption lasting 5 minutes, we predict you'll wait 87.1226039984 m
inutes until the next eruption.
```

The next cell plots the line that goes between those two points, which is (a segment of) the regression line.

```
In [17]:  plot_data_and_line(faithful, "duration", "wait", make_array(2, two_minute_pr
```

## Question 2

Make predictions for the waiting time after each eruption in the `faithful` table. (Of course, we know exactly what the waiting times were! We are doing this so we can see how accurate our predictions are.) Put these numbers into a column in a new table called `faithful_predictions`. Its first row should look like this:

|duration|wait|predicted wait| |-|-|-| |3.6|79|72.1011|

*Hint:* Your answer can be just one line. There is no need for a `for` loop; use array arithmetic instead.

```
In [18]: faithful_predictions = faithful.with_column("predicted wait", faithful["dura
         faithful_predictions
```

Out[18]:

| duration | wait | predicted wait |
|---|---|---|
| 3.6 | 79 | 72.1011 |
| 1.8 | 54 | 52.7878 |
| 3.333 | 74 | 69.2363 |
| 2.283 | 62 | 57.9702 |
| 4.533 | 85 | 82.1119 |
| 2.883 | 55 | 64.408 |
| 4.7 | 88 | 83.9037 |
| 3.6 | 85 | 72.1011 |
| 1.95 | 51 | 54.3972 |
| 4.35 | 85 | 80.1483 |

... (262 rows omitted)

## Question 3

How close were we? Compute the *residual* for each eruption in the dataset. The residual is the difference (not the absolute difference) between the actual waiting time and the predicted waiting time. Add the residuals to `faithful_predictions` as a new column called `"residual"`, naming the resulting table `faithful_residuals`.

*Hint:* Again, your code will be much simpler if you don't use a `for` loop.

```
In [19]: faithful_residuals = faithful_predictions.with_column("residual",
                                                 faithful_predictions["v
         faithful_residuals
```

Out[19]:

| duration | wait | predicted wait | residual |
|---|---|---|---|
| 3.6 | 79 | 72.1011 | 6.89889 |
| 1.8 | 54 | 52.7878 | 1.21225 |
| 3.333 | 74 | 69.2363 | 4.76371 |
| 2.283 | 62 | 57.9702 | 4.02983 |
| 4.533 | 85 | 82.1119 | 2.88814 |
| 2.883 | 55 | 64.408 | -9.40795 |
| 4.7 | 88 | 83.9037 | 4.09629 |
| 3.6 | 85 | 72.1011 | 12.8989 |
| 1.95 | 51 | 54.3972 | -3.3972 |
| 4.35 | 85 | 80.1483 | 4.85166 |

... (262 rows omitted)

Here is a plot of the residuals you computed. Each point corresponds to one eruption. It shows how much our prediction over- or under-estimated the waiting time.

```
In [20]: faithful_residuals.scatter("duration", "residual", color="r")
```
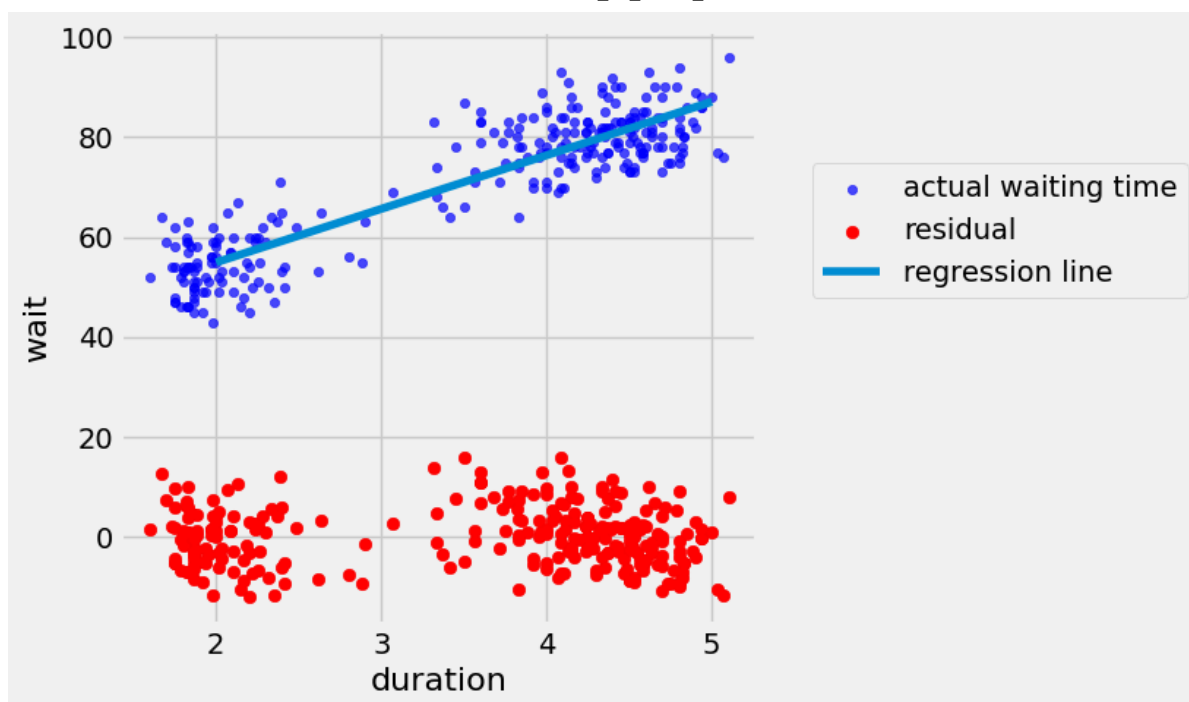
There isn't really a pattern in the residuals, which confirms that it was reasonable to try linear regression. It's true that there are two separate clouds; the eruption durations seemed to fall into two distinct clusters. But that's just a pattern in the eruption durations, not a pattern in the relationship between eruption durations and waiting times.

# 4. How accurate are different predictions?

Earlier, you should have found that the correlation is fairly close to 1, so the line fits fairly well on the training data. That means the residuals are overall small (close to 0) in comparison to the waiting times.

We can see that visually by plotting the waiting times and residuals together:

```
In [21]:  faithful_residuals.scatter("duration", "wait", label="actual waiting time",
          plt.scatter(faithful_residuals.column("duration"), faithful_residuals.column
          plt.plot(make_array(2, 5), make_array(two_minute_predicted_waiting_time, fiv
          plt.legend(bbox_to_anchor=(1.7,.8));
```

However, unless you have a strong reason to believe that the linear regression model is true, you should be wary of applying your prediction model to data that are very different from the training data.

## Question 1

In `faithful`, no eruption lasted exactly 0, 2.5, or 60 minutes. Using this line, what is the predicted waiting time for an eruption that lasts 0 minutes? 2.5 minutes? An hour?

In [22]:
```python
zero_minute_predicted_waiting_time = 0*slope + intercept
two_point_five_minute_predicted_waiting_time = 2.5*slope + intercept
hour_predicted_waiting_time = 60*slope + intercept

print_prediction(0, zero_minute_predicted_waiting_time)
print_prediction(2.5, two_point_five_minute_predicted_waiting_time)
print_prediction(60, hour_predicted_waiting_time)
```

```
After an eruption lasting 0 minutes, we predict you'll wait 33.4743970228 m
inutes until the next eruption.
After an eruption lasting 2.5 minutes, we predict you'll wait 60.2985005106
minutes until the next eruption.
After an eruption lasting 60 minutes, we predict you'll wait 677.252880731
minutes until the next eruption.
```

**Question 2.** Do you believe any of these values are reliable predictions? If you don't believe some of them, say why.

Prediction for 2.5 minutes and 60 minutes duration seems to reliable, although we are not sure of the 60 minutes duration predition. However, for 0 minute duration, we cannot say the the predicted waiting time is about 33.47 minutes since the eruptions did not happened and we are still waiting from the previous eruption until now.