

Assignment 1

Please fill out the relevant cells below according to the instructions. When done, save the notebook and export it to PDF, upload both the `.ipynb` and the PDF file to Canvas.

Group Members

Group submission is highly encouraged. If you submit as part of group, list all group members here. Groups can comprise up to 4 students.

- Vanessa Roser
- Dedan Campbell

Problem 1: Central Limit Theorem (2pts)

Use `scipy.stats` to draw N samples from the uniform and the Cauchy distribution. Confirm whether the mean μ of these samples (which is itself a RV) has a distribution $p(\mu)$ that converges to a normal distribution when $N \rightarrow \infty$.

A simple way of testing for normality of the distribution of means is the **68-95-99.7 rule**, i.e. you expect that there are only about 5% of the means (of a draw of N samples) that deviate from $\text{mean}(\mu)$ by more than $2\sqrt{\text{var}(\mu)}$.

Visualization can be helpful but is itself not a sufficient confirmation of normality!

```
In [1]: from scipy import stats
import numpy as np
import matplotlib.pyplot as plt

#mean and variance of the uniform and cauchy distribution
#uniform
mu_u = 0.5
var_u = 1/12
#cauchy
mu_c = 0
var_c = 1

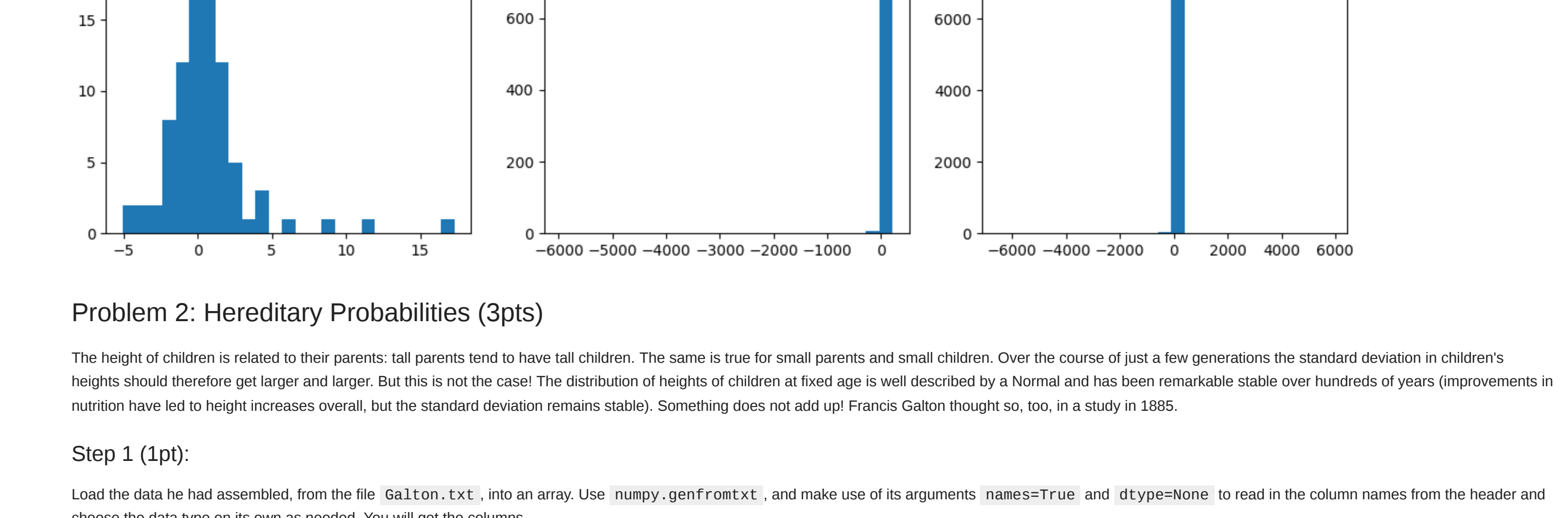
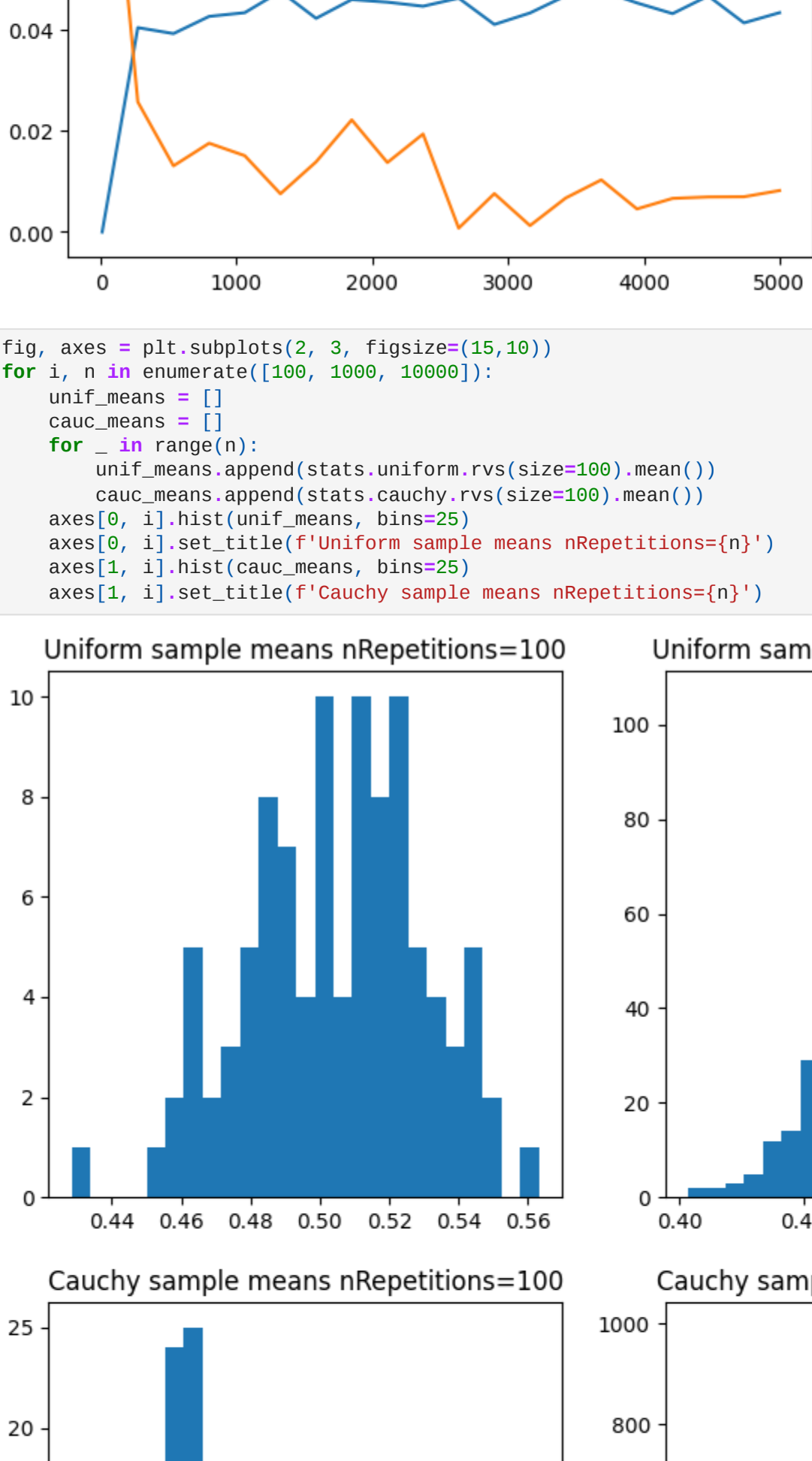
#draw N samples from the uniform and cauchy distribution
def convergence_test(n, size=100, dist='uniform'):
    """
    Sample from the specified distribution n times.
    """
    nus = np.zeros(n)
    for i in range(n):
        if dist=='uniform':
            samp = stats.uniform.rvs(size=size)
            mus[i] = samp.mean()
        elif dist=='cauchy':
            samp = stats.cauchy.rvs(size=size)
            mus[i] = samp.mean()
        else:
            raise ValueError('Invalid distribution')
    upper = np.mean(mus)+2*np.std(mus)
    lower = np.mean(mus)-2*np.std(mus)
    return (mus>upper) | (mus<lower), sum() / n

# Vary the number of samples
n_samples = np.linspace(10, 5000, 20)
ratio_u = np.zeros_like(n_samples)
ratio_c = np.zeros_like(n_samples)

for i, n in enumerate(n_samples):
    ratio_u[i] = convergence_test(n, size=100, dist='uniform')
    ratio_c[i] = convergence_test(n, size=50, dist='cauchy')

plt.plot(n_samples, ratio_u, label='uniform distribution convergence')
plt.plot(n_samples, ratio_c, label='Cauchy distribution convergence')
plt.axhline(0.95, color='r', ls='--', label='0.95')
plt.legend()
plt.title('Since the Cauchy distribution has neither a mean nor a variance, the central limit theorem does not apply.')
print('NOTE: Since the Cauchy distribution has neither a mean nor a variance, the central limit theorem does not apply.')

% samples >2SD from mean
```



Problem 2: Hereditary Probabilities (3pts)

The height of children is related to their parents: tall parents tend to have tall children. The same is true for small parents and small children. Over the course of just a few generations the standard deviation in children's heights should therefore get larger and larger. But this is not the case! The distribution of heights of children at fixed age is well described by a Normal and has been remarkable stable over hundreds of years (improvements in nutrition have led to height increases overall, but the standard deviation remains stable). Something does not add up! Francis Galton thought so, too, in a study in 1885.

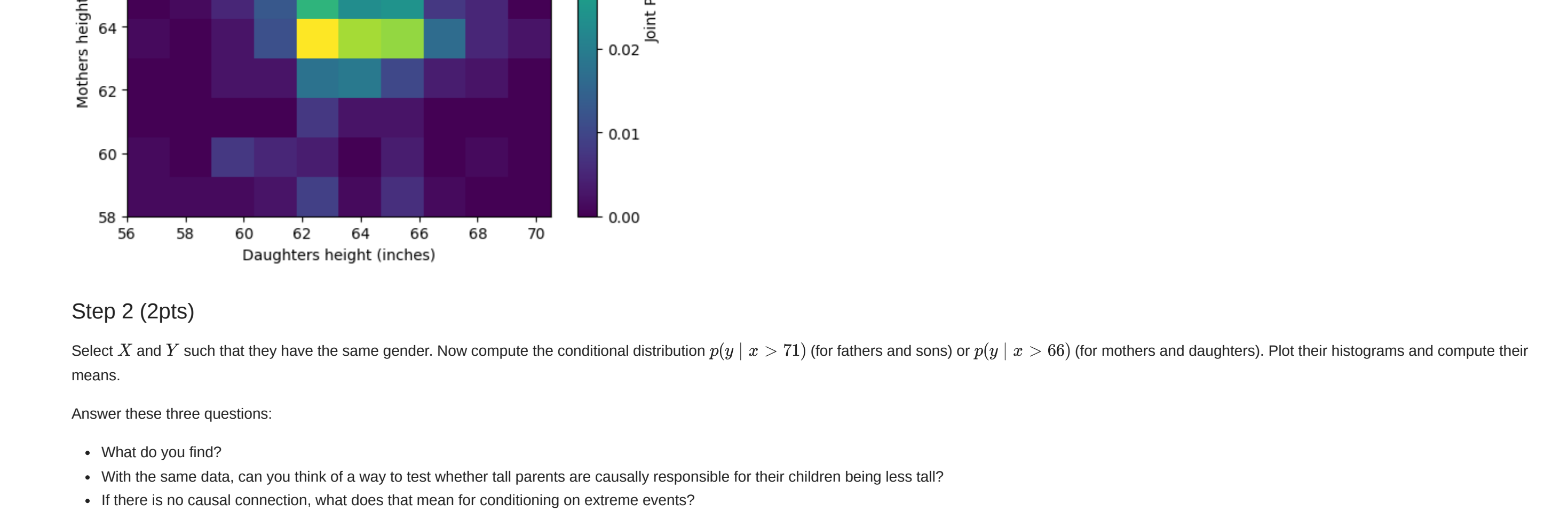
Step 1 (1pt)

Load the data he had assembled, from the file `Galton.txt`, into an array. Use `numpy.genfromtxt`, and make use of its arguments `names=True` and `dtype=None` to read in the column names from the header and choose the data type on its own as needed. You will get the columns

- `Family`: The family that the child belongs to, labeled from 1 to 204 and 136A
- `Father`: The father's height, in inches
- `Mother`: The mother's height, in inches
- `Gender`: The gender of the child, male (M) or female (F)
- `Height`: The height of the child, in inches (presumably fully grown)
- `Kids`: The number of kids in the family of the child

Make a visualization of the joint distribution of X , the parent's height (pick either father or mother), and Y , the children's height (pick either son or daughter).

Tip: The `matplotlib.pyplot.hist2d` is useful. Don't forget labels and units.



Step 2 (2pts)

Select X and Y such that they have the same gender. Now compute the conditional distribution $p(y | x > 71)$ (for fathers and sons) or $p(y | x > 66)$ (for mothers and daughters). Plot their histograms and compute their means.

Answer these three questions:

- What do you find?
- With the same data, can you think of a way to test whether tall parents are causally responsible for their children being less tall?
- If there is no causal connection, what does that mean for conditioning on extreme events?



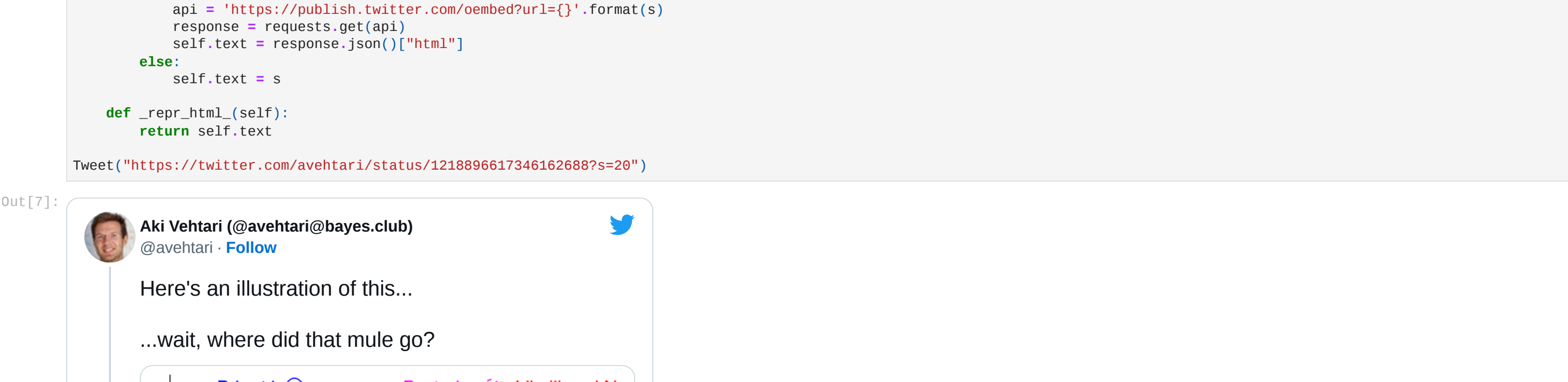
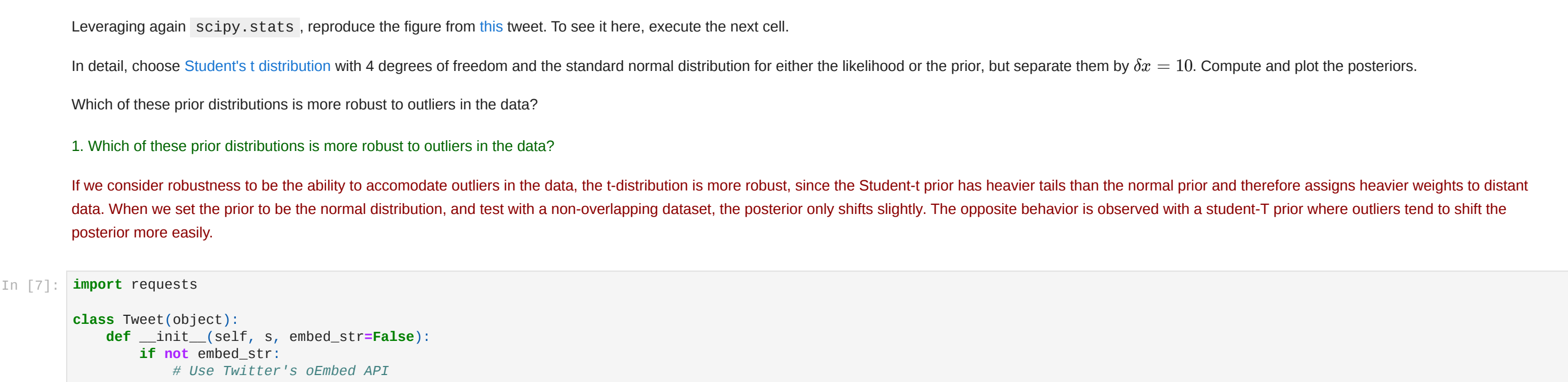
Problem 3: Likelihood vs Prior (1pt)

Leveraging again `scipy.stats`, reproduce the figure from this tweet. To see it here, execute the next cell.

In detail, choose **Student's t** distribution with 4 degrees of freedom and the standard normal distribution for either the likelihood or the prior, but separate them by $\Delta\mu = 10$. Compute and plot the posteriors.

Which of these prior distributions is more robust to outliers in the data?

If we consider robustness to be the ability to accommodate outliers in the data, the t-distribution is more robust, since the Student-t prior has heavier tails than the normal prior and therefore assigns heavier weights to distant data. When we set the prior to be the normal distribution, and test with a non-outlying dataset, the posterior only shifts slightly. The opposite behavior is observed with a student-t prior where outliers tend to shift the posterior more easily.



Problem 4: Hubble was no Bayesian (4pts)

...but you can be!

In 1929, Edwin Hubble published a seminal paper, in which he compared the radial velocity of astronomical objects (i.e. how fast these objects move towards or away from us) with their distance. The former can be done pretty precisely with spectroscopy, the latter is much more uncertain.

He saw that the velocity increases with distance and speculated that this could be the sign of a cosmological expansion. This led cosmologist to believe in the Big Bang theory.

Step 0:

Load the data from the file `hubble.txt` into an array with `numpy.genfromtxt`, and make again use of the arguments `names` and `dtype`. You should get 6 columns

- `GAT`, `NUMBER`: These two combined give you the name of the galaxy.
- `R`: distance in Mpc
- `V`: radial velocity in km/s
- `RA`, `DEC`: equatorial coordinates of the galaxy

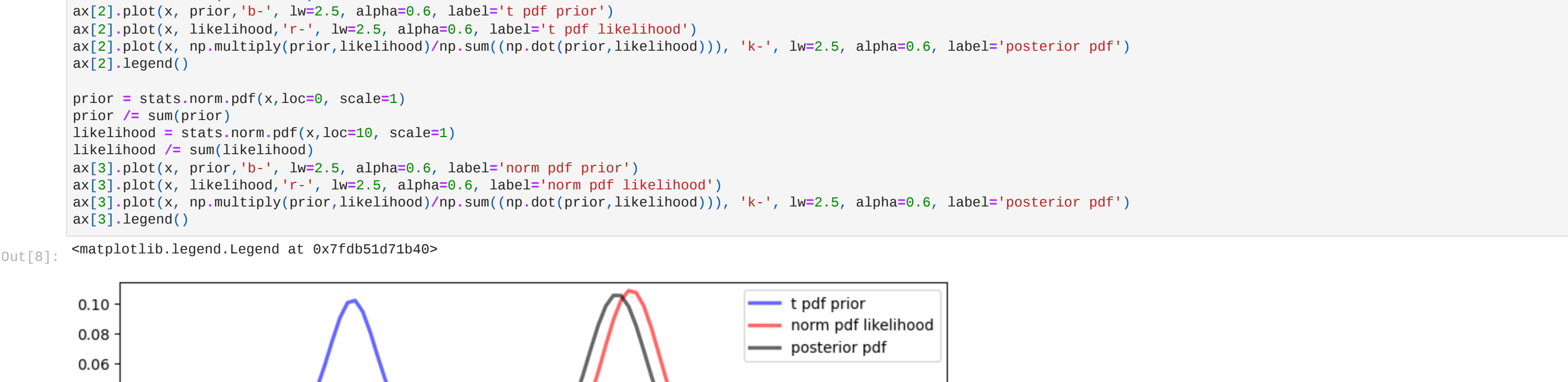
Make a scatter plot of R vs V (that means the independent variable is V). Don't forget labels and units...



Step 1 (1pt)

Use linear regression to determine the MLE of the slope b for the line $R = bV$. This is a linear model with **no intercept**. Print the MLE. Then, create a new version of the scatter plot by adding the MLE line.

Tip: You don't need measurement uncertainties (there aren't any in Hubble's data) to determine the MLE.

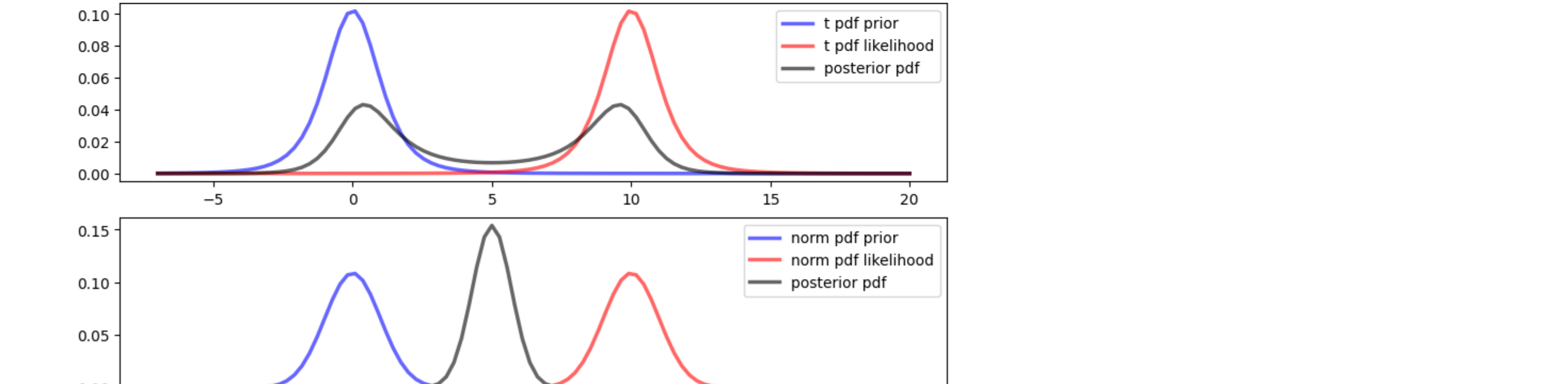


Step 2 (2pts)

The full Gaussian likelihood of the linear regression problem has a term for the intercept a , slope b , and uncertainty σ of R . We will assume that the uncertainties of all data points are identical. Adopt uninformative priors for all of the parameters $\theta = (a, b, \sigma)$.

Compute the log posterior on a reasonably fine grid of (a, b, σ) , picking suitable limits for every parameter. Then marginalize out σ and plot the log posterior for the remaining parameters (a, b) .

Tip: The function `scipy.special.logsumexp` is useful.



Step 3 (1pt)

Use the function `sample_2d` below to draw samples from the 2D array of the posterior of (a, b) . Create a final version of the scatter plot by adding the lines that correspond to these posterior draws.

Tip: When plotting, set the transparency `alpha` to values < 1 , so that multiple draws of the same parameter pair become visually more important.

