

Question 1c) Use `[]` to select **Name** and **Year** in that order from the `baby_names` table.

Then repeat the same selection using the `.loc` notation instead.

```
In [13]: name_and_year1 = baby_names[['Name', 'Year']]
         name_and_year1.head()
```

```
Out[13]:
```

	Name	Year
0	Mary	1910
1	Annie	1910
2	Anna	1910
3	Margaret	1910
4	Helen	1910

```
In [14]: name_and_year2 = baby_names.loc[:, ['Name', 'Year']]
         name_and_year2.head()
```

```
Out[14]:
```

	Name	Year
0	Mary	1910
1	Annie	1910
2	Anna	1910
3	Margaret	1910
4	Helen	1910

Question 2a) A coin is flipped 10 times. How many possible outcomes have exactly 2 heads? Use LaTeX (not code) in the cell below to show all of your steps and fully justify your answer.

Note: In this class, you must always put your answer in the cell that immediately follows the question. DO NOT create any cells between this one and the one that says *Write your answer here, replacing this text.*

This is an example of combinatorics. The formula for combinatorics is

$$\binom{n}{k} = \frac{n!}{k! \cdot (n-k)!}.$$

In our case, $k = 2$ and $n = 10$. Using these values we then have

$$\binom{10}{2} = \frac{10!}{2! \cdot (10-2)!} = \frac{10!}{2 \cdot 8!} = \frac{10 \cdot 9}{2} = \frac{90}{2} = 45.$$

Question 2b) What is the probability that if I roll two 6-sided dice they add up to **at most** 9? Use LaTeX (not code) in the cell directly below to show all of your steps and fully justify your answer.

For this problem, there are at most 36 possibilities. We are seeking to find the combinations of die that do not exceed 9. That means, we are interested in the sums from 2 to 9. The following are the possible ways the sums can be computed:

- 2: (1,1)
- 3: (1,2), (2,1)
- 4: (1,3), (2,2), (3,1)
- 5: (1,4), (2,3), (3,2), (4,1)
- 6: (1,5), (2,4), (3,3), (4,2), (5,1)
- 7: (1,6), (2,5), (3,4), (4,3), (5,2), (6,1)
- 8: (2,6), (3,5), (4,4), (5,3), (6,2)
- 9: (3,6), (4,5), (5,4), (6,3)

Counting these up, there are 30 possibilities. This means the probability is then $30 / 36$. Formally,

$$\frac{30}{36} \approx 83.3\%.$$

Question 2c) Suppose you show up to a quiz completely unprepared. The quiz has 10 questions, each with 5 multiple choice options. You decide to guess each answer in a completely random way. What is the probability that you get exactly 3 questions correct? Use LaTeX (not code) in the cell directly below to show all of your steps and fully justify your answer.

To solve this problem, we need to use the binomial formula. Namely,

$$P(x = k) = \binom{n}{k} \cdot p^k \cdot (1 - p)^{n-k}.$$

In this formula:

- n = The number of trials (in our case, the number of questions).
- k = The number of successes (in our case, the number of answers that are correct).
- p = The probability of success for a trial (in our case, the probability of correctly selecting the right answer).

In this specific problem we can then see that our probability is

$$P(x = 3) = \binom{10}{3} \cdot \left(\frac{1}{5}\right)^3 \cdot \left(1 - \frac{1}{5}\right)^{10-3}.$$

Furthermore, the probability is then

$$P(x = 3) = \frac{10!}{3!(10-3)!} \cdot \left(\frac{1}{5}\right)^3 \cdot \left(\frac{4}{5}\right)^7 = 120 \cdot \frac{1}{125} \cdot \frac{16384}{78125} = \frac{393216}{1953125} \approx 0.201.$$

Question 3a) We commonly use sigma notation to compactly write the definition of the arithmetic mean (commonly known as the average):

$$\bar{x} = \frac{1}{n} (x_1 + x_2 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i$$

The *i*th *deviation from average* is the difference $x_i - \bar{x}$. Prove that the sum of all these deviations is 0 that is, prove that $\sum_{i=1}^n (x_i - \bar{x}) = 0$ (write your full solution in the box directly below showing all steps and using LaTeX).

Substituting in the arithmetic mean into our expression we have

$$\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n \left(x_i - \frac{1}{n} \sum_{j=1}^n x_j \right).$$

We then distribute the outer sum in our previous expression to get

$$\sum_{i=1}^n \left(x_i - \frac{1}{n} \sum_{j=1}^n x_j \right) = \sum_{i=1}^n x_i - \sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n x_j.$$

Looking at the term

$$\sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n x_j$$

we are effectively summing the arithmetic mean (average) n times. So the previous expression simplifies to

$$\sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n x_j = n \cdot \bar{x}.$$

Which, when we use our definition of the arithmetic mean, we get

$$n \cdot \frac{1}{n} \sum_{i=1}^n x_i = \sum_{i=1}^n x_i.$$

So, when we substitute this result into our previous expressions we have

$$\sum_{i=1}^n x_i - \sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n x_j = \sum_{i=1}^n x_i - \sum_{i=1}^n x_i = 0.$$

We can see from the above that this expression is indeed equal to 0.

Question 3b) Let x_1, x_2, \dots, x_n be a list of numbers. You can think of each index i as the label of a household, and the entry x_i as the annual income of Household i .

Consider the function

$$f(c) = \frac{1}{n} \sum_{i=1}^n (x_i - c)^2$$

In this scenario, suppose that our data points x_1, x_2, \dots, x_n are fixed and that c is the only variable.

Using calculus, determine the value of c that minimizes $f(c)$. You must use calculus to justify that this is indeed a minimum, and not a maximum.

Taking the derivative of $f(c)$ with respect to c we have

$$\frac{df}{dc} = \frac{1}{n} \sum_{i=1}^n \frac{d}{dc} (x_i - c)^2 = \frac{1}{n} \sum_{i=1}^n 2 \cdot (x_i - c) \cdot (-1) = \frac{1}{n} \sum_{i=1}^n -2 \cdot (x_i - c).$$

Since we are asked to show that the extrema that will arrive from this is a minima, we need to calculate the second derivative of $f(c)$ with respect to c . The second derivative is then

$$\frac{d^2f}{dc^2} = \frac{d}{dc} \left(\frac{1}{n} \sum_{i=1}^n -2 \cdot (x_i - c) \right) = \frac{-2}{n} \sum_{i=1}^n \frac{d}{dc} (x_i - c) = \frac{-2}{n} \sum_{i=1}^n (-1) = \frac{2}{n} \sum_{i=1}^n (1) = \frac{2}{n} \cdot (n) = 2.$$

Now, the extrema that can be found from the first derivative of $f(c)$ is then

$$0 = \frac{1}{n} \sum_{i=1}^n -2 \cdot (x_i - c) = -2 \cdot \sum_{i=1}^n (x_i - c) = \sum_{i=1}^n (x_i - c) \rightarrow \sum_{i=1}^n x_i - \sum_{i=1}^n c = 0 \rightarrow \sum_{i=1}^n x_i - n(c) = 0.$$

Rearranging the above and solving for c we then find

$$c = \frac{1}{n} \sum_{i=1}^n x_i$$

which is, the arithmetic mean (average). Since our second derivative is positive, this implies that our extrema is indeed a minimum.

Question 4b) I have a coin that lands heads with an unknown probability p .

Suppose I toss it 10 times and get the sequence TTTHTHHTTH.

If you toss this coin 10 times, the chance that you get the sequence above is a function of p . That function is called the *likelihood* of the sequence TTTHTHHTTH, so we will call it $L(p)$.

What is $L(p)$ for the sequence TTTHTHHTTH?

Write your answer using LaTeX below (i.e. your answer should be of the form: $L(p)$ =some function of p)

The probability of obtaining tails in this scenario is going to be $(1 - p)$. Using the logic of part 4a, $L(p)$ is then

$$L(p) = (1 - p)(1 - p)(1 - p)(p)(1 - p)(p)(p)(1 - p)(1 - p)(p).$$

Question 4c) Below is a section of code that will help you plot the function $L(p)$ that you defined above. Replace the ellipses with your function of p

```
In [27]: p = np.linspace(0, 1, 100)
         #This creates an array of 100 values equally spaced between 0 and 1

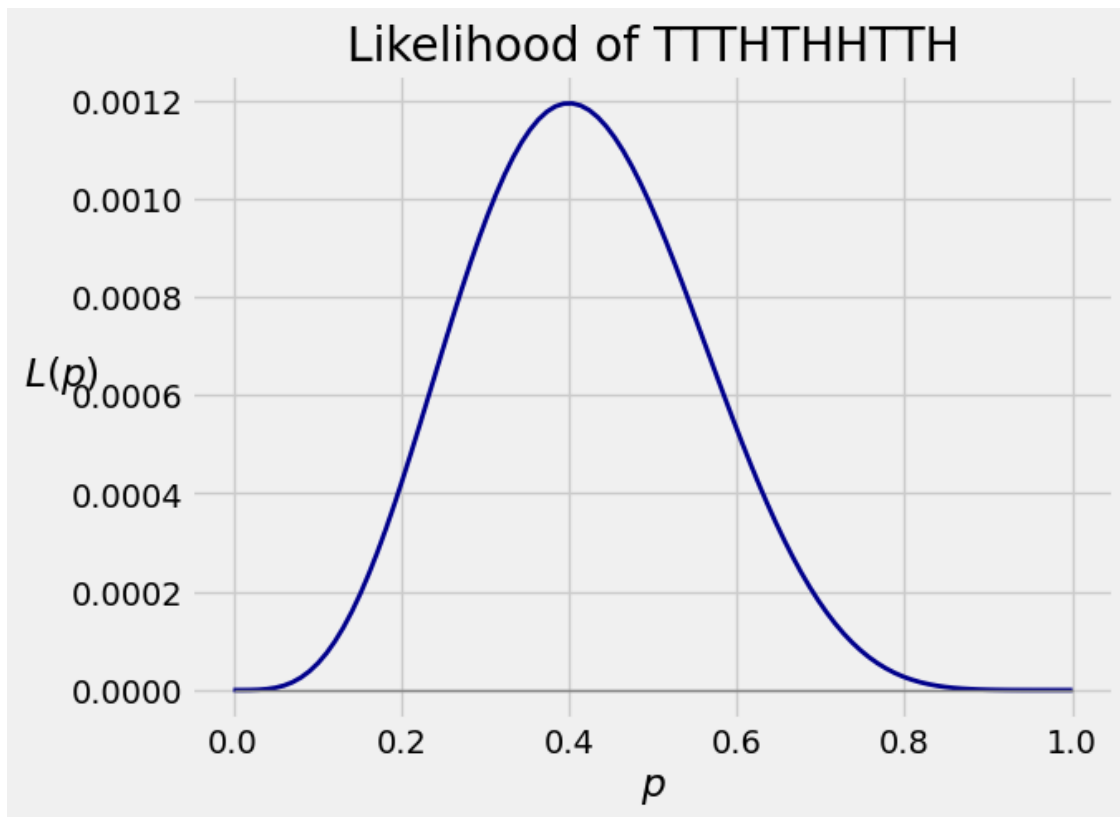
         likelihood = (1 - p)*(1 - p)*(1 - p)*(p)*(1 - p)*(p)*(p)*(1 - p)*(1 - p)*(p)

         plt.plot(p, likelihood, lw=2, color='darkblue')
         #This plots the likelihood function

         plt.plot([0, 1], [0, 0], lw=1, color='grey')
         #This plots a horizontal axis

         plt.xlabel('$p$')
         #This labels the x axis
         plt.ylabel('$L(p)$', rotation=0)
         #This labels the y-axis

         plt.title('Likelihood of TTTHTHHTTH');
         #This titles the plot
```



Question 4d) The value \hat{p} at which the likelihood function attains its maximum is called the *maximum likelihood estimate* (MLE) of p . Among all values of p , it is the one that makes the observed data most likely.

Using your plot above, what is the value of \hat{p} ?

Provide a simple interpretation of that value in terms of the data TTTHTHHTTH.

In our case, $\hat{p} = 0.4$. With $\hat{p} = 0.4$, this means that for the above sequence of tails and heads to most likely appear is when we have $p = 0.4$. Or in english, when we have a 60% chance of obtaining tails and a 40% chance of obtaining heads, this gives us the best chance of the aforementioned pattern happening.

Question 4e) Let's prove what you observed graphically above. That is, let's use calculus to find \hat{p} .

But wait before you start trying to find the value p where $L'(p) = 0$ (trust us, the algebra is not pretty...)

TIP:

The value \hat{p} at which the function L attains its maximum is the same as the value at which the function $\log(L)$ attains its maximum. To clarify, $\log(L)$ is the composition of \log and L : $\log(L)$ at p is $\log(L(p))$. Even though it doesn't make a difference for this problem, \log is now and forevermore the \log to the base e , not to the base 10.

This tip is hugely important in data science because many probabilities are products and the \log function turns products into sums. It's much simpler to work with a sum than with a product.

Armed with that tip use calculus to find \hat{p} . You don't have to check that the value you've found produces a max and not a min – we'll spare you that step.

Although you have provided me with a helpful hint, I am going to go the “not pretty” algebra route. We can rewrite the aforementioned $L(p)$ equation in a simpler form. $L(p)$ is then

$$L(p) = (1 - p)^6(p)^4.$$

Taking the derivative of this, by using the chain and product rules, we then find $L'(p)$ to be

$$L'(p) = 4(1 - p)^6(p)^3 - 6(1 - p)^5(p)^4 = (1 - p)^5(p)^3(4(1 - p) - 6p).$$

Since we are seeking to find the extrema of $L'(p)$, we need to find the value of p for when $L'(p) = 0$. Namely

$$\begin{aligned} 0 &= (1 - p)^5(p)^3(4(1 - p) - 6p) \\ &= (4(1 - p) - 6p) \\ &= 4 - 4p - 6p \\ &= 4 - 10p \\ 4 &= 10p \\ 0.4 &= p. \end{aligned}$$

The above implicates that $\hat{p} = 0.4$, which coincides with the aforementioned graph.

