

homework-01

September 4, 2023

1 Homework 1

1.1 References

- Lectures 1 through 4 (inclusive).

1.2 Instructions

- Type your name and email in the “Student details” section below.
- Develop the code and generate the figures you need to solve the problems using this notebook.
- For the answers that require a mathematical proof or derivation you should type them using latex. If you have never written latex before and you find it exceedingly difficult, we will likely accept handwritten solutions.
- The total homework points are 100. Please note that the problems are not weighed equally.

```
[2]: import matplotlib.pyplot as plt
%matplotlib inline
import matplotlib_inline
matplotlib_inline.backend_inline.set_matplotlib_formats('svg')
import seaborn as sns
sns.set_context("paper")
sns.set_style("ticks")

import numpy as np
import scipy
import scipy.stats as st
import urllib.request
import os

def download(
    url : str,
    local_filename : str = None
):
    """Download a file from a url.

    Arguments
    url          -- The url we want to download.
    local_filename -- The filename to write on. If not
                     specified
```

```

"""
if local_filename is None:
    local_filename = os.path.basename(url)
urlretrieve(url, local_filename)

```

1.3 Student details

- **First Name:** Dominic
- **Last Name:** Ferro
- **Email:** dferro@purdue.edu
- **Used generative AI to complete this assignment (Yes/No):** No
- **Which generative AI tool did you use (if applicable)?:**

1.4 Problem 1

Disclaimer: This example is a modified version of the one found in a 2013 lecture on Bayesian Scientific Computing taught by Prof. Nicholas Zabaraz. I am not sure where the original problem is coming from.

We are tasked with assessing the usefulness of a tuberculosis test. The prior information I is:

The percentage of the population infected by tuberculosis is 0.4%. We have run several experiments and determined that: + If a tested patient has the disease, then 80% of the time the test comes out positive. + If a tested patient does not have the disease, then 90% of the time the test comes out negative.

To facilitate your analysis, consider the following logical sentences concerning a patient:

A: The patient is tested and the test is positive.

B: The patient has tuberculosis.

A. Find the probability that the patient has tuberculosis (before looking at the result of the test), i.e., $p(B|I)$. This is known as the base rate or the prior probability. **Answer:**

$p(B|I) = 0.4\%$. Before looking at the test, we only know that this patient is a member of the population.

B. Find the probability that the test is positive given that the patient has tuberculosis, i.e., $p(A|B, I)$. **Answer:**

$p(A|BI) = 80\%$. This comes from the first condition of the givens.

C. Find the probability that the test is positive given that the patient does not have tuberculosis, i.e., $p(A|\neg B, I)$. **Answer:**

$p(A|\neg B, I) = 1 - p(\neg A|\neg B, I) = 1 - 0.9 = 0.1$. There is a 10% chance that the test is positive given that the patient does not have tuberculosis.

D. Find the probability that a patient that tested positive has tuberculosis, i.e., $p(B|A, I)$. **Answer:**

$$p(B|A, I) = \frac{p(B, A|I)}{p(A|I)} = \frac{p(A|B, I)p(B|I)}{p(A|I)}$$

We need $p(A|I)$, the probability of any patient receiving a positive test, which is the probability of a patient having tuberculosis and testing positive plus the probability of a patient not having tuberculosis and testing positive:

$$p(A|I) = p(A, B|I) + p(A, \neg B|I) = p(A|B, I)p(B|I) + p(A|\neg B, I)p(\neg B|I)$$

$$p(A|I) = 0.8 * 0.004 + 0.1 * 0.996 = 0.1028$$

Thus,

$$p(B|A, I) = \frac{p(A|B, I)p(B|I)}{p(A|I)} = \frac{0.8*0.004}{0.1028} = 0.03113$$

There is a 3.113% chance that a patient that tested positive has tuberculosis.

E. Find the probability that a patient that tested negative has tuberculosis, i.e., $p(B|\neg A, I)$. Does the test change our prior state of knowledge about the patient? Is the test useful? **Answer:**

$$p(B|\neg A, I) = \frac{p(\neg A|B, I)p(B|I)}{p(\neg A|I)} = \frac{0.2*0.004}{1-0.1028} = 0.00089. \text{ There is a 0.089\% chance that a patient with a negative test has tuberculosis.}$$

The test does change our prior state of knowledge about the patient. If there is a positive test, the probability of the patient having tuberculosis increases from 0.4% to 3.1%. If there is a negative test, the probability decreases to 0.089%.

This test is not very useful as there is initially a low probability of the patient having tuberculosis. A positive test does not indicate a high likelihood of the patient actually having tuberculosis. While we gain some knowledge from the test, it is hardly enough to be considered useful.

F. What would a good test look like? Find values for

$$p(A|B, I) = p(\text{test is positive}|\text{has tuberculosis}, I),$$

and

$$p(A|\neg B, I) = p(\text{test is positive}|\text{does not have tuberculosis}, I),$$

so that

$$p(B|A, I) = p(\text{has tuberculosis}|\text{test is positive}, I) = 0.99.$$

There are more than one solutions. How would you pick a good one? Thinking in this way can help you set goals if you work in R&D. If you have time, try to figure out whether or not there exists such an accurate test for tuberculosis **Answer:**

$$\text{Using the expression: } p(B|A, I) = \frac{p(A|B, I)p(B|I)}{p(A|B, I)p(B|I) + p(A|\neg B, I)p(\neg B|I)}$$

We can find that $p(A|B, I) = 0.999$ and $p(A|\neg B, I) = 0.00004$ gives $p(B|A, I) = 0.99013$

It is most important to minimize the number of false positives. The factor $p(A|\neg B, I)p(\neg B|I)$ is the term that impacts the final probability the most. A quick google search says that the false positive rate in tuberculosis tests due to cross contamination is 2% roughly three orders of magnitude larger than the idealized false positive rate.

1.5 Problem 2 - Practice with discrete random variables

Consider the Categorical random variable:

$$X \sim \text{Categorical}(0.3, 0.1, 0.2, 0.4),$$

taking values in $\{0, 1, 2, 3\}$. Find the following (you may use `scipy.stats.rv_discrete` or do it by hand):

A. The expectation $\mathbb{E}[X]$.

Answer:

```
[3]: # Categorical Probabilites
ps = [0.3, 0.1, 0.2, 0.4]
# Corresponding values
xs = [0, 1, 2, 3]
# Create the categorical rv
X = st.rv_discrete(name="P2_Categorical", values=(xs,ps))

# Expectation
print(f"E[X] = {X.expect():.2f}")
```

$\mathbb{E}[X] = 1.70$

B. The variance $\mathbb{V}[X]$.

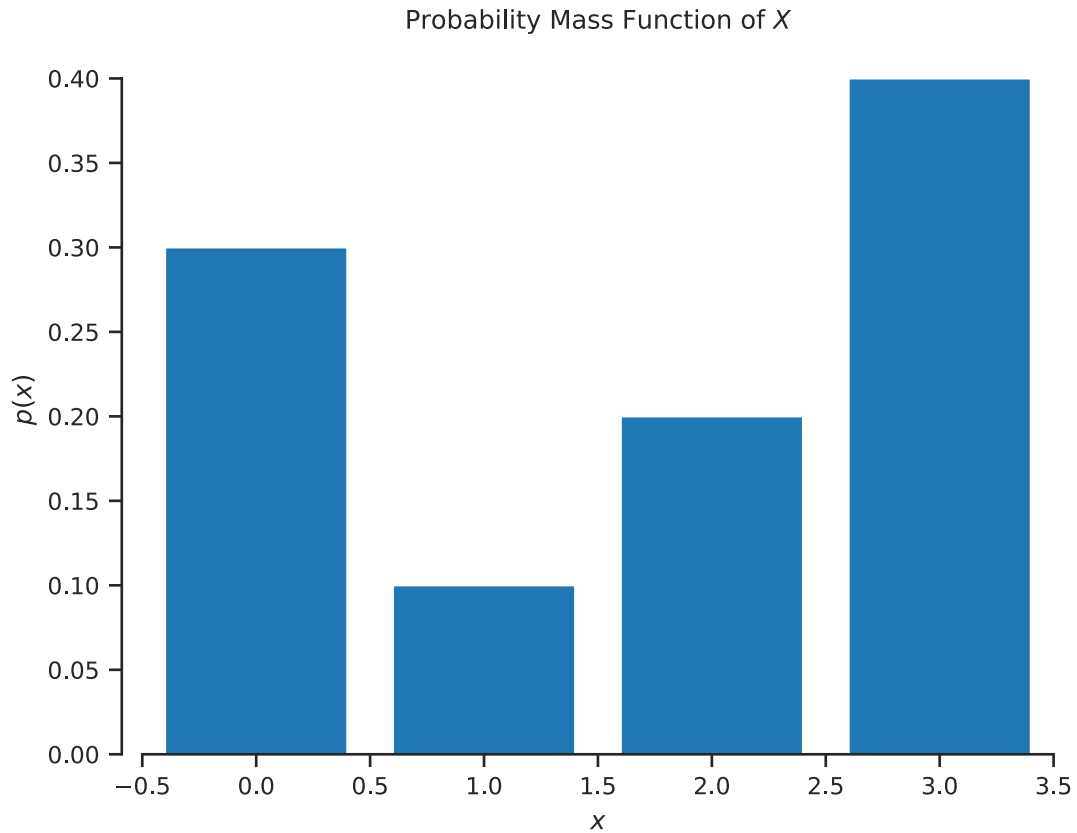
Answer:

```
[4]: # Variance
print(f"V[X] = {X.var():.2f}")
```

$\mathbb{V}[X] = 1.61$

C. Plot the probability mass function of X .

```
[5]: fig, ax = plt.subplots()
ax.bar(xs, X.pmf(xs))
ax.set_xlabel("$x$")
ax.set_ylabel("$p(x)$")
ax.set_title("Probability Mass Function of $X$")
sns.despine(trim=True)
```



D. Find the probability that X is in $\{0, 2\}$.

Answer:

```
[6]: px_02 = 0
      for ii in range(0,3):
          px_02 += X.pmf(ii)

      print(f"p(x in {0,2}) = {px_02:.2f}")
```

$p(x \text{ in } (0, 2)) = 0.60$

E. Find $\mathbb{E}[4X + 3]$.

Answer:

$$\mathbb{E}[4X + 3] = 4\mathbb{E}[X] + 3$$

```
[7]: print(f"E[4X + 3] = {4*X.expect() + 3:.2f}")
```

$$\mathbb{E}[4X + 3] = 9.80$$

F. Find $\mathbb{V}[4X + 3]$.

Answer:

$$\mathbb{V}[4X + 3] = \mathbb{V}[4X] + \mathbb{V}[3] = 16\mathbb{V}[X]$$

```
[8]: print(f"V[4X+3] = {16*X.var():.2f}")
```

$$\mathbb{V}[4X+3] = 25.76$$

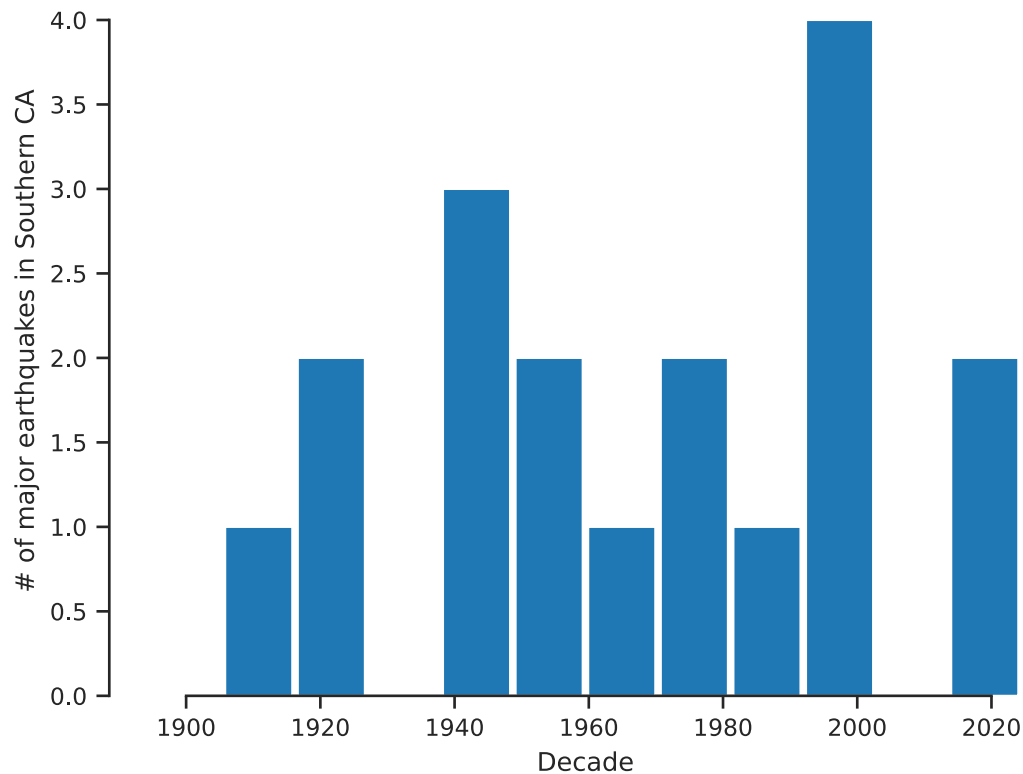
1.6 Problem 3 - Predicting the probability of major earthquakes in Southern California

The [San Andreas fault](#) extends through California forming the boundary between the Pacific and the North American tectonic plates. It has caused some of the major earthquakes on Earth. We are going to focus on Southern California and we would like to assess the probability of a major earthquake, defined as an earthquake of magnitude 6.5 or greater, during the next ten years.

A. The first thing we are going to do is go over a [database of past earthquakes](#) that have occurred in Southern California and collect the relevant data. We are going to start at 1900 because data before that time may be unreliable. Go over each decade and count the occurrence of a major earthquake (i.e., count the number of orange and red colors in each decade). We have done this for you.

```
[24]: eq_data = np.array([
    0, # 1900-1909
    1, # 1910-1919
    2, # 1920-1929
    0, # 1930-1939
    3, # 1940-1949
    2, # 1950-1959
    1, # 1960-1969
    2, # 1970-1979
    1, # 1980-1989
    4, # 1990-1999
    0, # 2000-2009
    2 # 2010-2019
])
fig, ax = plt.subplots(dpi=150)
ax.bar(np.linspace(1900, 2019, eq_data.shape[0]), eq_data, width=10)
ax.set_xlabel('Decade')
ax.set_ylabel('# of major earthquakes in Southern CA')
#plt.legend(loc="best", frameon=False)
sns.despine(trim=True);
```

WARNING:matplotlib.legend:No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



B. The [Poisson distribution](#) is a discrete distribution with values $\{0, 1, 2, \dots\}$ which is commonly used to model the number of events occurring in a certain time period. It is the right choice when these events are happening independently and the probability of any event happening over a small period of time is constant. Let's use the Poisson to model the number of earthquakes X occurring in a decade. We write:

$$X \sim \text{Poisson}(r),$$

where r is the *rate parameter* of Poisson. The rate is the number of events per time period. Here, r is the number of earthquakes per decade. Using the data above, we can set the rate as the empirical average of the observed number of earthquakes per decade:

```
[10]: r = np.mean(eq_data)
      print('r = {0:1.2f} major earthquakes per decade'.format(r))
```

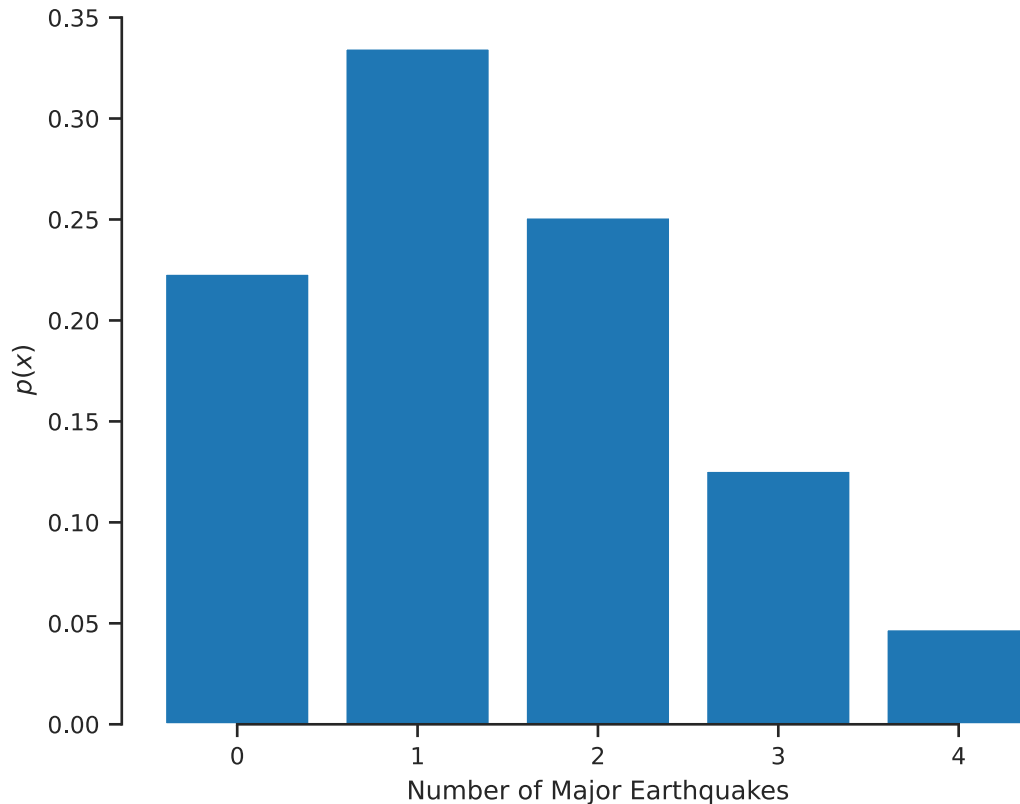
```
r = 1.50 major earthquakes per decade
```

Strictly speaking, **this is not how you should be calibrating models!!!** We will learn about the **right** way (which uses Bayes' rule) in the subsequent lectures. But it will do for now as the answer you would get using the **right** way is, for this problem, almost the same. Let's define a Poisson distribution using `scipy.stats.poisson` (see documentation [here](#)):

```
[11]: X = st.poisson(r)
```

A. Plot the probability mass function of X .

```
[31]: fig, ax = plt.subplots(dpi=300)
      ax.bar(eq_data, X.pmf(eq_data))
      ax.set_xlabel("Number of Major Earthquakes")
      ax.set_ylabel("$p(x)$")
      sns.despine(trim=True)
```



B. What is the probability that at least one major earthquake will occur during the next decade?

Answer:

```
[65]: print(f"p(x>0) = {1-X.pmf(0):0.3f}")
```

$p(x>0) = 0.777$

C. What is the probability that at least one major earthquake will occur during the next two decades? Hint: Consider two independent and identical copies of X , say X_1 and X_2 . And consider their sum $Y = X_1 + X_2$. Read [this](#) about the sum of two independent Poisson distributions.

Answer:

```
[66]: Y = st.poisson(r+r)
      print(f"p(x>0) = {1-Y.pmf(0):0.3f}")
```


$p(x>0) = 0.950$

D. What is the probability that at least one major earthquake will occur during the next five decades? **Answer:**

```
[67]: Y_D = st.poisson(5*r)
      print(f"p(x>0) = {1-Y_D.pmf(0):0.3f}")
```

$p(x>0) = 0.999$

1.7 Problem 4 - Failure of a mechanical component

Assume that you designing a gear for a mechanical system. Under normal operating conditions the gear is expected to fail at a random time. Let T be a random variable capturing the time the gear fails. What should the probability density of T look like?

Here are some hypothetical data to work with. Suppose that we took ten gears and we worked them until failure. The failure times (say in years) are as follows:

```
[16]: time_to_fail_data = np.array(
      [
          10.5,
          7.5,
          8.1,
          8.4,
          11.2,
          9.3,
          8.9,
          12.4
      ]
    )
```

Why does each gear fail at different times? There are several sources of uncertainty. The most important are:

- Manufacturing imperfections.
- Different loading conditions.

If this was a controlled fatigue experiment, then we could eliminate the second source of uncertainty by using exactly the same loading conditions.

Now, we are going to fit a probability density function to these data. Which one should we use? Well, new gears do not fail easily. So, the probability density function of T should be close to zero for small T . As time goes by, the probability density should increase because various things start happening to the material, e.g., crack formation, fatigue, etc. Finally, the probability density must again start going to zero as time further increases because nothing lasts forever... A probability distribution that is commonly used to model this situation is the [Weibull](#). We are going to fit some fail time data to a Weibull distribution and then you will have to answer a few questions about failing times.

The Weibull has parameters and we are going to fit them to the available data. The method we are going to use is called the *maximum likelihood method*. We haven't really talked about this, and it

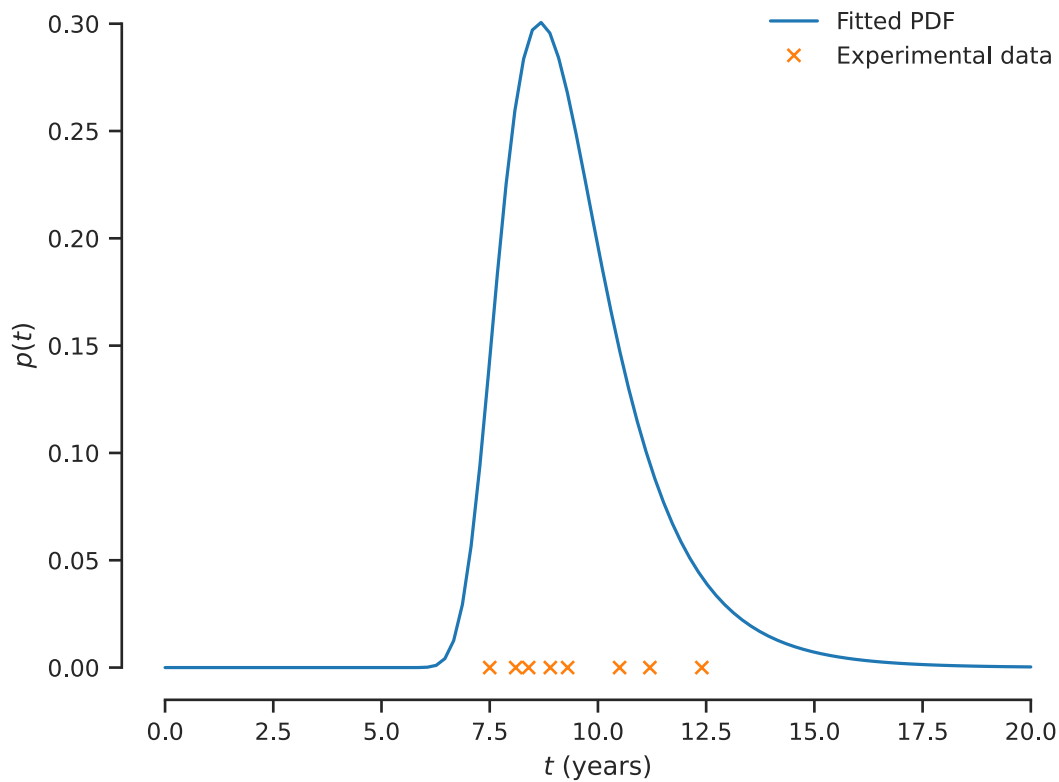
is not important to know what it is to do this homework problem. We will learn about maximum likelihood in later lectures. Here is how we fit the parameters using `scipy.stats`:

```
[17]: fitted_params = st.exponweib.fit(time_to_fail_data, loc=0)
      T = st.exponweib(*fitted_params)
      print(f"Fitted parameters: {fitted_params}")
```

Fitted parameters: (448.066965711728, 0.7099665338918923, 3.4218808260575804, 0.41627831297126994)

Let's plot the fitted Weibul PDF and the data we used:

```
[18]: fig, ax = plt.subplots()
      ts = np.linspace(0.0, 20.0, 100)
      ax.plot(
          ts,
          T.pdf(ts),
          label="Fitted PDF"
      )
      ax.plot(
          time_to_fail_data,
          np.zeros_like(time_to_fail_data),
          "x",
          label="Experimental data"
      )
      ax.set_xlabel(r"$t$ (years)")
      ax.set_ylabel(r"$p(t)$")
      plt.legend(loc="best", frameon=False)
      sns.despine(trim=True);
```



Now you have to answer a series of questions about the random variable T that we just fitted.

A. Find the mean fail time and its variance. Hint: Do not integrate anything by hand. Just use the functionality of `scipy.stats`.

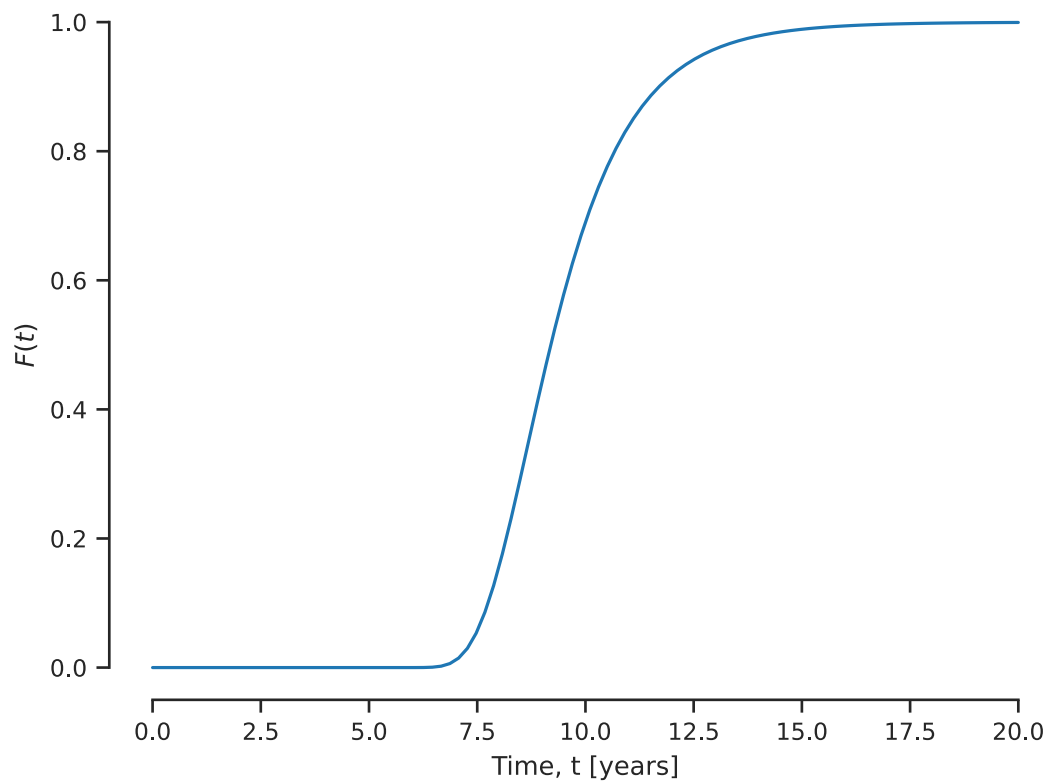
```
[38]: t_mean = T.expect()
      t_var = T.var()
      print(f"E[T] = {t_mean:.2f}")
      print(f"V[T] = {t_var:.2f}")
```

E[T] = 9.53

V[T] = 2.88

B. Plot the cumulative distribution function $F(t) = P(T \leq t)$ of T .

```
[40]: fig, ax = plt.subplots(dpi=300)
      ax.plot(ts, T.cdf(ts))
      ax.set_xlabel('Time, t [years]')
      ax.set_ylabel('$F(t)$')
      sns.despine(trim=True)
```

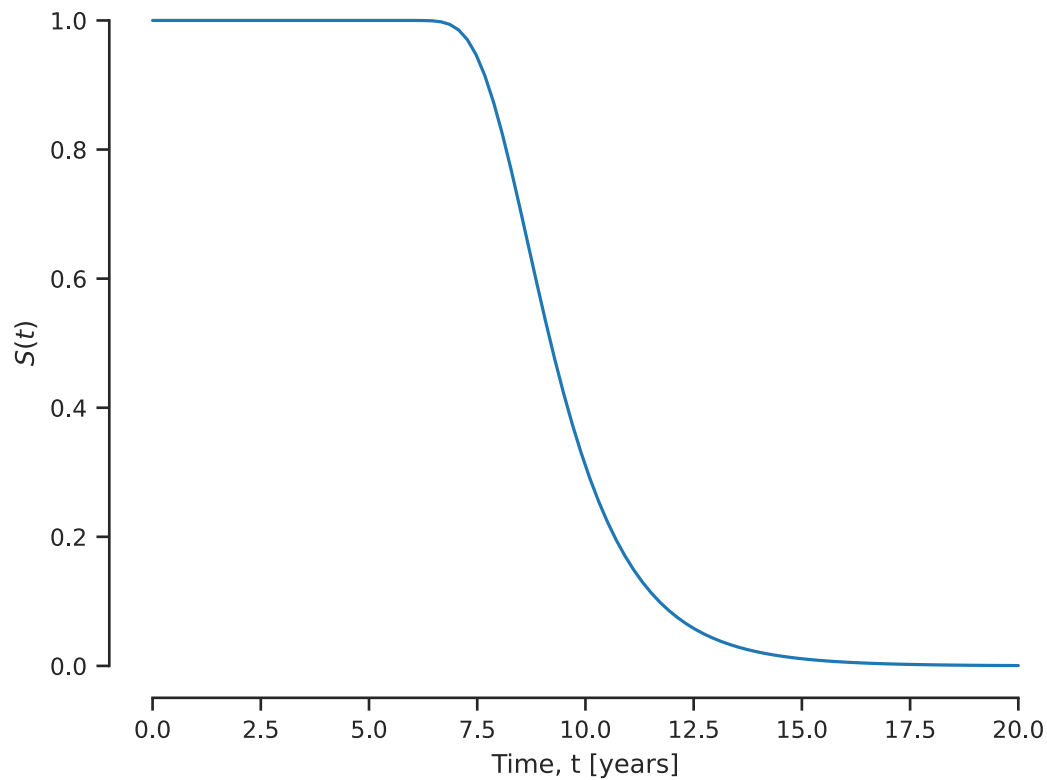


C. Plot the probability that gear survives for more than t as a function of t . That is, plot the function:

$$S(t) = p(T > t).$$

Hint: First connect $S(t)$ to the cumulative distribution function $F(t)$ of T .

```
[42]: S = 1 - T.cdf(ts)
fig, ax = plt.subplots(dpi=300)
ax.plot(ts, S)
ax.set_xlabel('Time, t [years]')
ax.set_ylabel('$S(t)$')
sns.despine(trim=True)
```



D. Find the probability that the gear lasts anywhere between 8 and 10 years.

```
[51]: print(f"p(8<=T<=10) = {(T.cdf(10) - T.cdf(8)):.3f}")
```

$p(8 \leq T \leq 10) = 0.534$

E. Find the time t^* such that the probability that the gear fails before t^* is 0.01.

```
[63]: nn = int(1e5)
tspan = np.linspace(0, 20, nn)
for t in tspan:
    if T.cdf(t) >= 0.01:
        tstar = t
        break
print(f"At time t* = {tstar:.3f}, the probability that the gear fails is_
↳ p(T<t*) = {T.cdf(tstar):.3f}")
```

At time $t^* = 6.975$, the probability that the gear fails is $p(T < t^*) = 0.010$