Now that we have our pmf, we want to find the cdf. This is done by doing the following

$$F_Y(y) = P(Y \le y) = \sum_{i=1}^{y} (1-p)^{i-1} p = p \left[ \frac{1 - (1-p)^y}{1 - (1-p)} \right] = 1 - (1-p)^y \qquad (10.6)$$

by the geometric series.

To prove this is a cdf we show that as $y$ goes to $-\infty$ we get 0 (which is true because $F_Y(y)$ is defined to be 0 for $y < 0$, as it goes to $\infty$ we get 1 (which is true because $1 - p$ is a number less than 1 and raised to infinity goes to 0), it is monotonically increasing, and lastly we show it is right continuous by $F_Y(y + \epsilon) = F_Y(y)$ as $\epsilon$ goes to 0.

1.2 Say we have a standard normal distribution that we take draws from $X_1, X_2, ...$ until we draw a value that is greater than some value $p$, in which case we stop sampling. What is the expected value of the draws?

The first step is to write down what we want. The way to think of this is that we want to know the expected value of $X$ (where $X$ is a normal random variable), but with the condition that $X < p$. This can be written as $E[X|X < p]$. Using the standard normal distribution we have:

$$
\begin{aligned}
E[X|X < p] &= \int_{-\infty}^{\infty} x f_X(x|X < p) dx \qquad \text{(by definition of expectation)} \\
&= \int_{-\infty}^{\infty} x \frac{f_X(x, X < p)}{P(X < p)} dx \qquad \text{(by definition of conditional probability)} \\
&= \frac{1}{P(X < p)} \int_{-\infty}^{p} x \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} dx \qquad \text{(limit integral bounds to match conditional prob.)} \\
&= \frac{1}{P(X < p)\sqrt{2\pi}} \int_{-\infty}^{p} x e^{\frac{-x^2}{2}} dx
\end{aligned}
$$

The quantity $P(X < p)$ is given by the cdf of the standard normal distribution. The integral, using u-substitution, is found to be:

$$
\begin{aligned}
\int_{-\infty}^{p} x e^{\frac{-x^2}{2}} dx &= \int_{-\infty}^{\frac{-p^2}{2}} -e^u du \\
&= -e^{\frac{-p^2}{2}} + e^{-\infty} \\
&= -e^{\frac{-p^2}{2}}
\end{aligned}
\qquad (10.7)
$$

where $u = \frac{-x^2}{2}$.

All together we have $\frac{-e^{\frac{-p^2}{2}}}{P(X<p)\sqrt{2\pi}}$. If $p = 1$ for example then we have

$$\frac{-e^{\frac{-1}{2}}}{P(X < 1)\sqrt{2\pi}} \approx \frac{-e^{\frac{-1}{2}}}{0.84 * \sqrt{2\pi}} \approx -0.288 \qquad (10.8)$$

One question I had is if there is a difference between taking the expectation of a sequence from the distribution and stopping once we have a value that is greater than $p$ vs. taking the

expectation of a set sample with variables that are greater than $p$ thrown out.

From a simulation perspective there is no difference if we think of the set sample as a bunch of sequences tied together, broken up by the values greater than $p$. In this light, the sequences are essentially defined by the realizations that are greater than $p$. So even though thinking about this process from a sequence point of view, I feel it is safe to conclude that we are just taking the normal random sample process approach, looking for the expectation of $X$ conditioned on the event that $X > p$.

The code for this simulation is given below:

```python
from scipy.stats import norm
import numpy as np

mean_rvs = []
N = 10000
p = 1
for i in range(N):
    rvs = []
    less = True
    while less:
        rv = norm.rvs(size=1) # sample one normal, standardized random variable
        if rv < p: # if draw is less than p, add to sequence
            rvs.append(rv[0])
        else: # if draw is greater than p, stop sequence
            less=False
    if len(rvs)!=0: # take mean of sequence, as long as there is as least one random variable in sequence
        mean_rvs.append(np.mean(rvs))
```

where the mean of "mean_rvs" is around -0.288.

1.3 Find the probability that two points on a unit line have a distance less than 0.5.

$$
\begin{aligned}
P(|Y - X| < 0.5) &= P(-0.5 < Y - X < 0.5) \\
&= P(Y - X < 0.5) - P(Y - X < -0.5) \\
&= P(Y < X + 0.5) - P(Y < X - 0.5)
\end{aligned}
\tag{10.9}
$$

We can think of $Y, X$ having a joint uniform distribution on the unit square. This gives us every possible combination of these two variables. If we think of it in this context and relate the two variables together then we can deal with the probability terms above.

The first part of the last expression above is going to be the area (since this is the uniform distribution) below the line $Y = X + 0.5$ on the unit square. This ends up being $\frac{7}{8}$. The second part will be the area below the line $Y < X - 0.5$ which ends up being $\frac{1}{8}$. Therefore we have a difference of $\frac{6}{8}$ or 0.75.