Problem 1: Hashing Distribution

To understand the distribution of collisions, let's start by supposing that there are $N$ spaces $m$ elements can be hashed into. Let the conditional expectation $E_m$ represent the expected value of $m$ elements being hashed into unique locations. Let 1 or 0 represent if a location is occupied or not. Then,

$$E_0 = 0$$
$$E_1 = 1 \qquad\qquad (E_m = \sum x_i P(x_i))$$
$$E_2 = E_1 + 1 * {}^{N-E_1}/N$$
$$E_3 = E_2 + 1 * {}^{N-E_2}/N$$
$$\vdots$$
$$E_m = E_{m-1} + 1 * {}^{N-E_{m-1}}/N$$

We can simplify the above by letting $c = {}^{N-1}/N$.

Then,
$$E_0 = 0$$
$$E_1 = 1$$
$$E_2 = cE_1 + 1 \qquad\qquad E_m = \sum_{i=0}^{m-1} c^i.$$
$$E_3 = cE_2 + 1$$
$$\vdots$$
$$E_m = cE_{m-1} + 1$$

$$* \quad \sum_{i=0}^{N} c^i = \frac{1 - c^{N+1}}{1 - c}.$$

Therefore, $E_m = \dfrac{1 - c^m}{1 - c} = N(1 - c^m)$.

We can make another approximation by using $e^{-1/N} \approx 1 - 1/N$ for $c$. Then, we get
$$E_m \approx N(1 - e^{-m/N}).$$

If we have $N = 512$ and we hash $m = 512$ elements, we should see:
$$N(1 - c^m) = 512\left(1 - ({}^{511}/512)^{512}\right) \approx 324.$$

This means that we had no collisions for 324 lists or spaces. Consequently, we should have around $512 - 324 = 188$ collisions. This matches with our simulation results.