

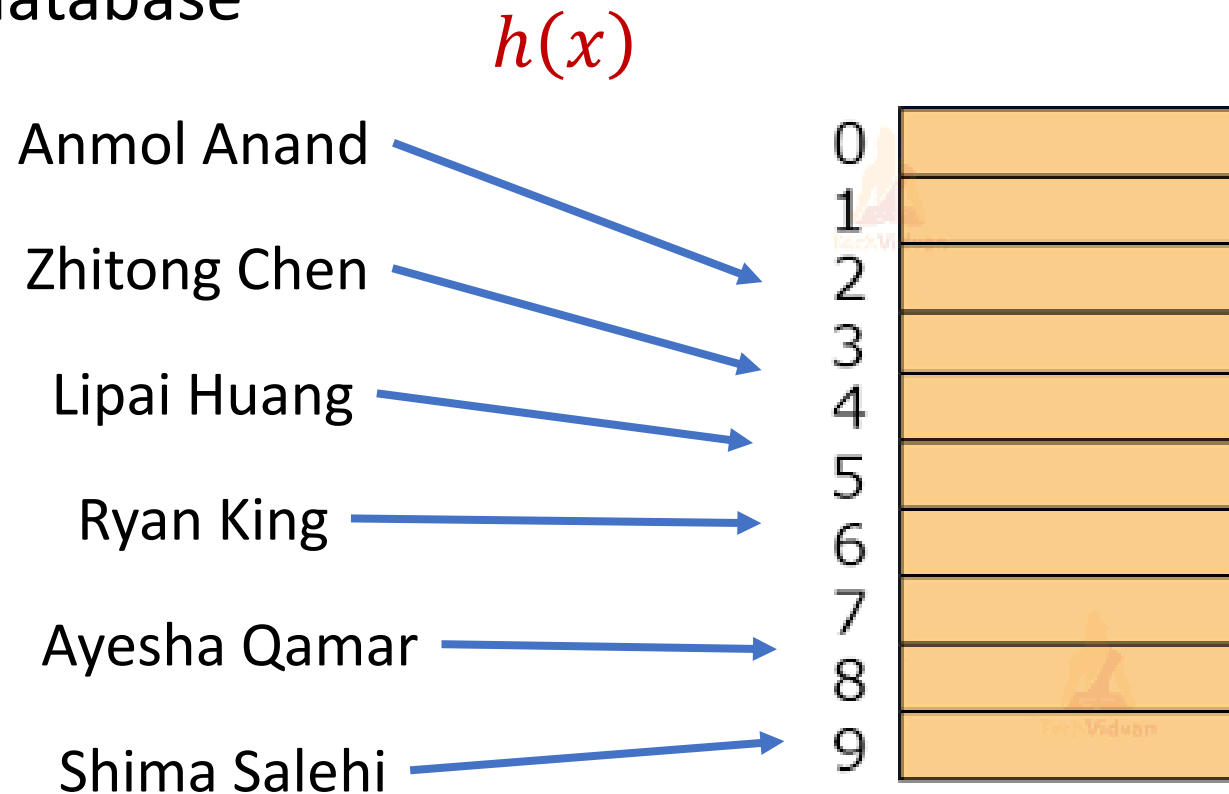
# CSCS 689: Special Topics in Modern Algorithms for Data Science

## Lecture 3

Samson Zhou

# Last Time: Hashing

- Hashing is a method to quickly map items from a universe to a location in a database



# Last Time: Birthday Paradox

- Suppose we have a fair  $n$ -sided die. “On average”, how many times should we roll the die before we see a repeated outcome among the rolls? Example: 1, 5, 2, 4, 5
- $\Theta(1)$
- $\Theta(\log n)$
- $\Theta(\sqrt{n})$
- $\Theta(n)$

# Last Time: Birthday Paradox

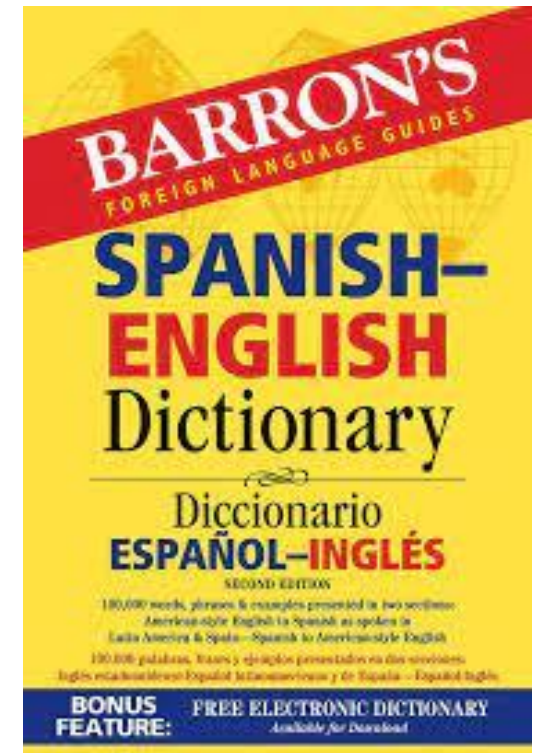
- Suppose we have a fair  $n$ -sided die. “On average”, how many times should we roll the die before we see a repeated outcome among the rolls? Example: 1, 5, 2, 4, 5
- $\Theta(1)$
- $\Theta(\log n)$
- $\Theta(\sqrt{n})$
- $\Theta(n)$

# Future

- **Next Monday:** Sign up for LaTeX scribe note slots
- **Today:** Meet your classmates (1)
- **Next Monday:** Meet your classmates (2), receive and consider list of potential projects/groups
- **Next Wednesday:** Discuss potential project groups
- **Next Friday:** Email me the members/group name

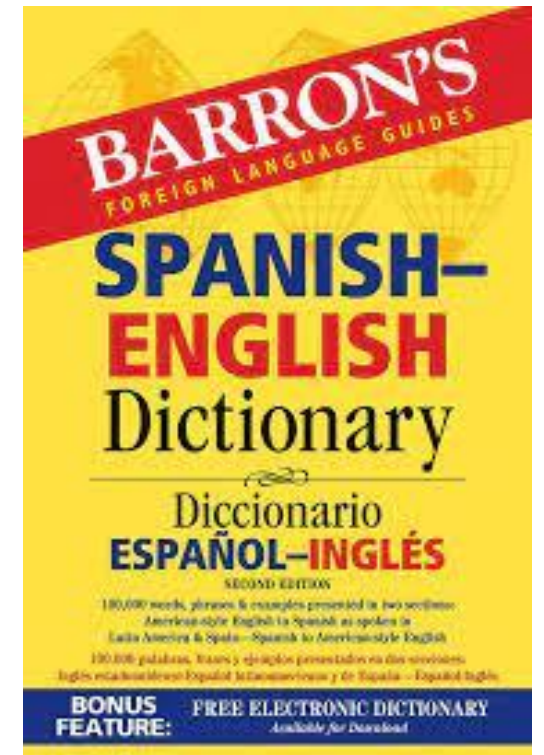
# Case Study

- We are trying to learn a new language on an app, which claims to have a database of *1 million words*
- Each time we ask the app, it gives us a random word in the database
- We want to verify the claim



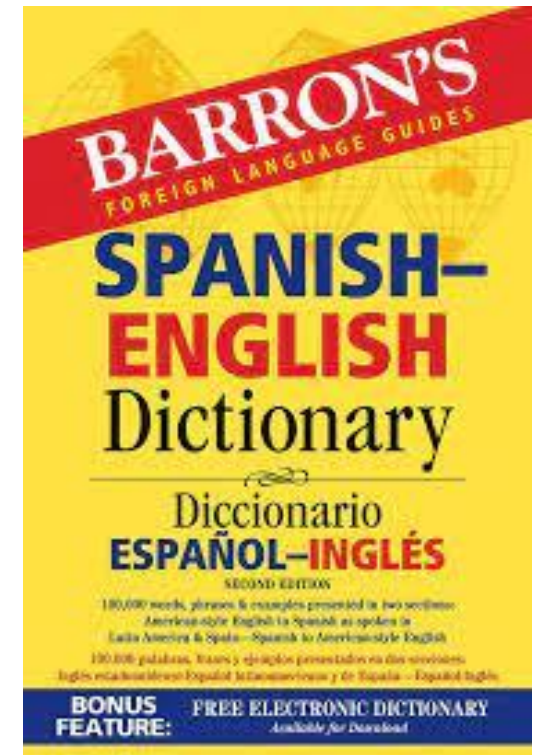
# Case Study

- We could use the app until we see 1 million unique words, but that would take at least *1 million checks*
- Instead, we use the app for *1000 times* and count the number of pairwise duplicates
- If there are many duplicates, the database is probably not very large



# Case Study

- We use the app for  $k$  times and count the number of pairwise duplicates
- If we see the same word on the 3-rd time, the 100-th time, and the 205-th time, there are 3 pairwise duplicates: (3, 100), (3, 205), (100, 205)





# Expected Value

- The expected value of a random variable  $X$  over  $\Omega$  is:

$$E[X] = \sum_{x \in \Omega} \Pr[X = x] \cdot x$$

- The “average value of the random variable”
- Linearity of expectation:  $E[X + Y] = E[X] + E[Y]$

# Linearity of Expectation

- Linearity of expectation:  $E[X + Y] = E[X] + E[Y]$

$$E[X + Y] = \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot (x + y)$$

# Linearity of Expectation

- Linearity of expectation:  $E[X + Y] = E[X] + E[Y]$

$$\begin{aligned} E[X + Y] &= \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot (x + y) \\ &= \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot x + \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot y \end{aligned}$$

# Linearity of Expectation

- Linearity of expectation:  $E[X + Y] = E[X] + E[Y]$

$$\begin{aligned} E[X + Y] &= \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot (x + y) \\ &= \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot x + \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot y \\ &= \sum_{x \in \Omega_X} x \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] + \sum_{y \in \Omega_Y} y \sum_{x \in \Omega_X} \Pr[X = x, Y = y] \end{aligned}$$

# Linearity of Expectation

- Linearity of expectation:  $E[X + Y] = E[X] + E[Y]$

$$\begin{aligned} E[X + Y] &= \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot (x + y) \\ &= \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot x + \sum_{x \in \Omega_X} \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] \cdot y \\ &= \sum_{x \in \Omega_X} x \sum_{y \in \Omega_Y} \Pr[X = x, Y = y] + \sum_{y \in \Omega_Y} y \sum_{x \in \Omega_X} \Pr[X = x, Y = y] \\ &= \sum_{x \in \Omega_X} x \cdot \Pr[X = x] + \sum_{y \in \Omega_Y} y \cdot \Pr[Y = y] = E[X] + E[Y] \end{aligned}$$

# Expected Value

- Suppose we roll a 6-sided die
- Let  $X$  be the outcome of the roll
- What is  $E[X]$ ?

# Birthday Paradox

- Suppose we have a fair  $n$ -sided die that we roll  $k = 1, 2, 3, 4, \dots$  times. What is the probability we see a repeated outcome among the rolls?

$$\left(1 - \frac{0}{n}\right) \left(1 - \frac{1}{n}\right) \dots \left(1 - \frac{k-1}{n}\right)$$

# Birthday Paradox, Revisited

- Suppose we have a fair  $n$ -sided die that we roll  $k = 1, 2, 3, 4, \dots$  times. What is the expected number of pairwise collisions among the rolls?
- Let  $X_i$  be the number of pairwise collisions on the  $i$ -th roll
- We have  $E[X_i] = \frac{i-1}{n}$



# Birthday Paradox, Revisited

- Let  $X$  be the number of pairwise collisions after  $k$  rolls
- What is  $E[X]$ ?

# Birthday Paradox, Revisited

- Let  $X$  be the number of pairwise collisions after  $k$  rolls

$$\begin{aligned} E[X] &= E[X_1 + \cdots + X_k] \\ &= E[X_1] + \cdots + E[X_k] \\ &= \frac{0}{n} + \cdots + \frac{k-1}{n} \\ &= \frac{k(k-1)}{n} \end{aligned}$$

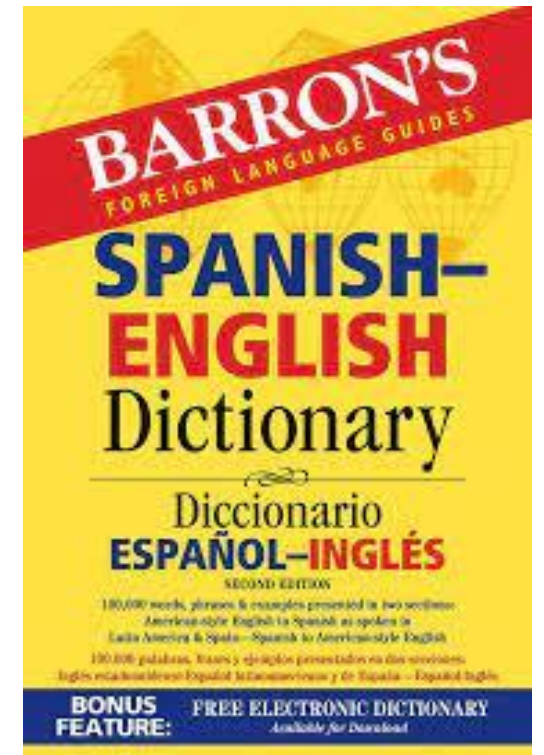
# Birthday Paradox, Revisited

- $E[X] = \frac{k(k-1)}{n}$
- $\frac{(k-1)^2}{n} \leq E[X] \leq \frac{k^2}{n}$
- $k = \sqrt{n} + 1$  implies  $E[X] \geq 1$
- $k = \frac{\sqrt{n}}{2}$  implies  $E[X] \leq \frac{1}{4}$

# Case Study

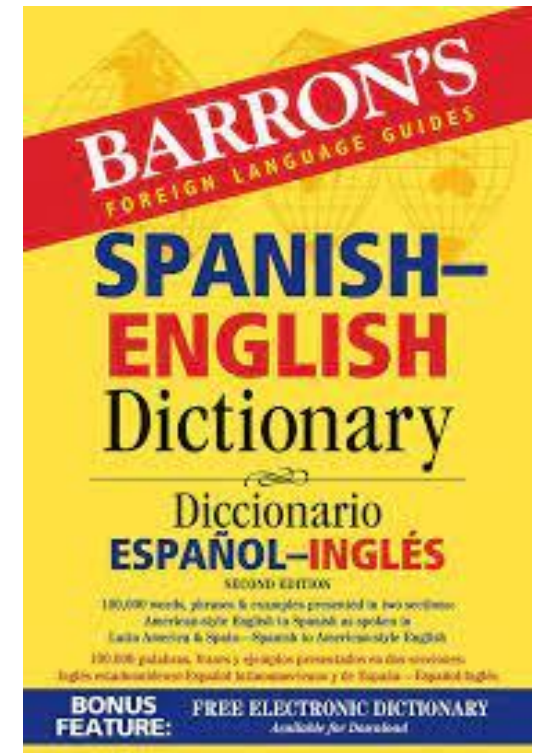
- We use the app for  $k = 1000$  times and count the number of pairwise duplicates

- If the database contains *1 million words*, the expected number of pairwise duplicates is  $E[X] = \frac{k(k-1)}{n} < 0.5$



# Case Study

- If the database contains *1 million words*, the expected number of pairwise duplicates is  $E[X] = \frac{k(k-1)}{n} < 0.5$
- ...We see **20** duplicates
- We think the claim is incorrect, but how can we be sure?



# Concentration Inequalities

- Concentration inequalities bound the probability that a random variable is “far away” from its expectation
- Often used in understanding the performance of statistical tests, the behavior of data sampled from various distributions, and for our purposes, the guarantees of randomized algorithms

# Markov's Inequality

- Let  $X \geq 0$  be a non-negative random variable. Then for any  $t > 0$ :

$$\Pr[X \geq t \cdot E[X]] \leq \frac{1}{t}$$

# Proof of Markov's Inequality

- Let  $X \geq 0$  be a non-negative random variable. Then for any  $t > 0$ :

$$\begin{aligned} E[X] &= \sum_{x \in \Omega} \Pr[X = x] \cdot x \\ &= \sum_{x \geq t \cdot E[X]} \Pr[X = x] \cdot x + \sum_{x < t \cdot E[X]} \Pr[X = x] \cdot x \\ &\geq \sum_{x \geq t \cdot E[X]} \Pr[X = x] \cdot x \\ &\geq t \cdot E[X] \sum_{x \geq t \cdot E[X]} \Pr[X = x] \\ &= t \cdot E[X] \cdot \Pr[X \geq t \cdot E[X]] \end{aligned}$$



# Birthday Paradox

- Suppose we have a fair  $n$ -sided die that we roll  $k = 1, 2, 3, 4, \dots$  times. What is the probability we see a repeated outcome among the rolls?

$$\left(1 - \frac{0}{n}\right) \left(1 - \frac{1}{n}\right) \dots \left(1 - \frac{k-1}{n}\right)$$

# Birthday Paradox, Revisited

- Suppose we have a fair  $n$ -sided die that we roll  $k = 1, 2, 3, 4, \dots$  times. What is the expected number of pairwise collisions among the rolls?
- Let  $X_i$  be the number of pairwise collisions on the  $i$ -th roll
- We have  $E[X_i] = \frac{i-1}{n}$

# Birthday Paradox, Revisited

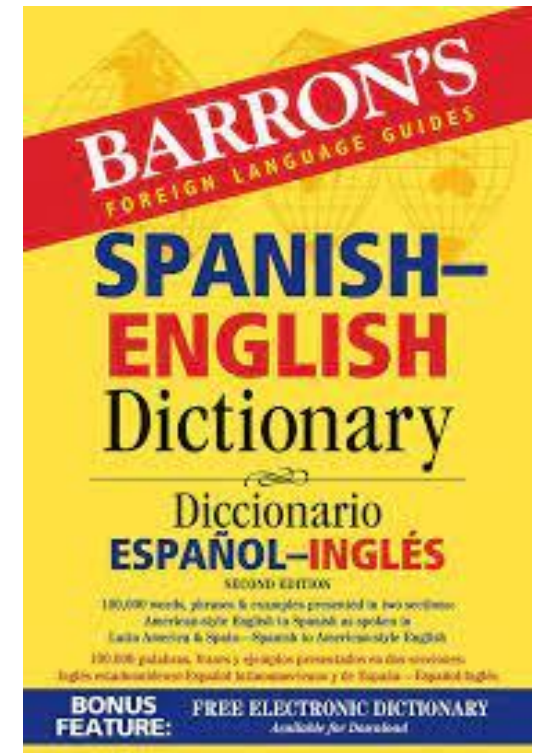
- $E[X] = \frac{k(k-1)}{n}$
- $\frac{(k-1)^2}{n} \leq E[X] \leq \frac{k^2}{n}$
- $k = \sqrt{n} + 1$  implies  $E[X] \geq 1$
- $k = \frac{\sqrt{n}}{2}$  implies  $E[X] \leq \frac{1}{4}$

# Birthday Paradox, Revisited

- $E[X] = \frac{k(k-1)}{n}$
- $\frac{(k-1)^2}{n} \leq E[X] \leq \frac{k^2}{n}$
- $k = \sqrt{n} + 1$  implies  $E[X] \geq 1$
- $k = \frac{\sqrt{n}}{2}$  implies  $E[X] \leq \frac{1}{4}$ ,  $\Pr[X \geq 1] \leq \frac{1}{4}$ ,

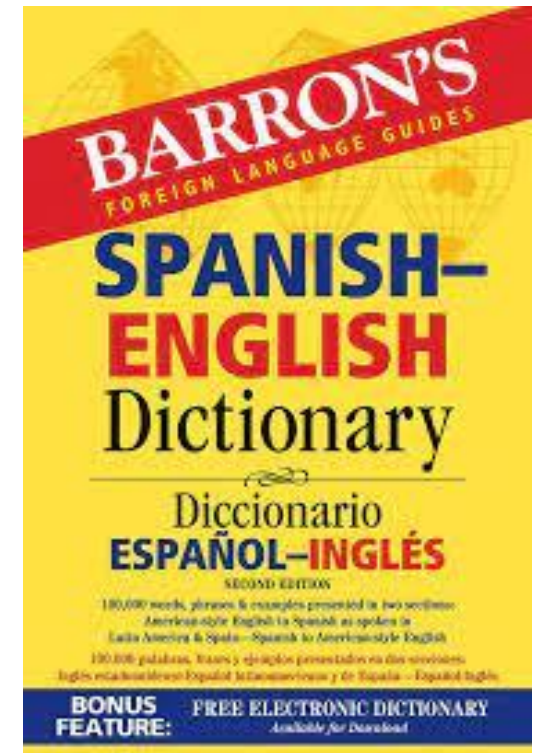
# Case Study

- If the database contains *1 million words*, the expected number of pairwise duplicates is  $E[X] = \frac{k(k-1)}{n} < 0.5$
- ...We see **20** duplicates
- We think the claim is incorrect, but how can we be sure?



# Case Study

- If the database contains *1 million words*, the expected number of pairwise duplicates is  $E[X] = \frac{k(k-1)}{n} < 0.5$
- ...We see **20** duplicates
- $\Pr[X \geq 20] \leq \frac{1}{40}$



# CSCS 689: Special Topics in Modern Algorithms for Data Science

## Lecture 4

Samson Zhou

## Trivia Question #3 (Max Load)

- Suppose we have a fair  $n$ -sided die that we roll  $n$  times. “On average”, what is the largest number of times any outcome is rolled? Example: 1, 5, 2, 4, 1, 3, 1 for  $n = 7$
- $\Theta(1)$
- $\tilde{\Theta}(\log n)$
- $\tilde{\Theta}(\sqrt{n})$
- $\tilde{\Theta}(n)$



## Trivia Question #4 (Coupon Collector)

- Suppose we have a fair  $n$ -sided die. “On average”, how many times should we roll the die before we all possible outcomes among the rolls? Example: 1, 5, 2, 4, 1, 3, 1, 6 for  $n = 6$
- $\Theta(n)$
- $\Theta(n \log n)$
- $\Theta(n\sqrt{n})$
- $\Theta(n^2)$

# Moments

- For  $p > 0$ , the  $p$ -th moment of a random variable  $X$  over  $\Omega$  is:

$$E[X^p] = \sum_{x \in \Omega} \Pr[X = x] \cdot x^p$$

# Variance

- The variance of a random variable  $X$  over  $\Omega$  is:

$$\text{Var}[X] = E[X^2] - (E[X])^2$$

- Linearity of variance for *independent* random variables:  $\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$
- “How far numbers are from the average”

# Variance

- Suppose  $X$  takes the value  $1$  with probability  $\frac{1}{2}$  and takes the value  $-1$  with probability  $\frac{1}{2}$
- What is  $E[X]$ ?
- What is  $\text{Var}[X]$ ?

# Variance

- Suppose  $Y$  takes the value  $100$  with probability  $\frac{1}{2}$  and takes the value  $-100$  with probability  $\frac{1}{2}$
- What is  $E[Y]$ ?
- What is  $\text{Var}[Y]$ ?

# Chebyshev's Inequality

- Let  $X$  be a random variable with expected value  $\mu := E[X]$  and variance  $\sigma^2 := \text{Var}[X]$

$$\Pr[|X - \mu| \geq k\sigma] \leq \frac{1}{k^2}$$

- “What is the probability a random variable is far away from its average?”