COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Andrew McGregor

Lecture 7

- Estimate $\#$ distinct elements $d$ in stream $x_1, \ldots, x_n \in U$

- Estimate # distinct elements $d$ in stream $x_1, \ldots, x_n \in U$
- **Basic Algorithm:**
  - Let $h_1, h_2, \ldots, h_k : U \to [0,1]$ be random hash functions.
  - Compute $s = \frac{1}{k} \sum s_i$ where $s_i = \min_{j \in [n]} h_i(x_j)$
  - Return $\hat{d} = 1/s - 1$

- Estimate # distinct elements $d$ in stream $x_1, \ldots, x_n \in U$
- **Basic Algorithm:**
  - Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions.
  - Compute $s = \frac{1}{k} \sum s_i$ where $s_i = \min_{j \in [n]} h_i(x_j)$
  - Return $\hat{d} = 1/s - 1$
- **Analysis:** $\mathbb{E}[s] = \frac{1}{d+1}$ and $\mathrm{Var}[s] \le \frac{1}{k(d+1)^2}$. Setting $k = O(\epsilon^{-2})$ ensures $(1 - \epsilon)\mathbb{E}[s] \le s \le (1 + \epsilon)\mathbb{E}[s]$ with probability at least $3/4$ and

$$(1 - \epsilon)\mathbb{E}[s] \le s \le (1 + \epsilon)\mathbb{E}[s] \Rightarrow (1 - 4\epsilon)d \le \hat{d} \le (1 + 4\epsilon)d$$

- Estimate # distinct elements $d$ in stream $x_1, \ldots, x_n \in U$
- **Basic Algorithm:**
  - Let $h_1, h_2, \ldots, h_k : U \to [0, 1]$ be random hash functions.
  - Compute $s = \frac{1}{k} \sum s_i$ where $s_i = \min_{j \in [n]} h_i(x_j)$
  - Return $\hat{d} = 1/s - 1$
- **Analysis:** $\mathbb{E}[s] = \frac{1}{d+1}$ and $\text{Var}[s] \leq \frac{1}{k(d+1)^2}$. Setting $k = O(\epsilon^{-2})$ ensures $(1 - \epsilon)\mathbb{E}[s] \leq s \leq (1 + \epsilon)\mathbb{E}[s]$ with probability at least $3/4$ and

$$(1 - \epsilon)\mathbb{E}[s] \leq s \leq (1 + \epsilon)\mathbb{E}[s] \Rightarrow (1 - 4\epsilon)d \leq \hat{d} \leq (1 + 4\epsilon)d$$

- **Median Trick:** If an algorithm returns a sufficiently accurate numerical answer with probability at least $3/4$, run it $O(\log(1/\delta))$ times and take the median answer. This will have the required accuracy with probability at least $1 - \delta$.

Our algorithm uses continuous valued fully random hash functions.

Our algorithm uses continuous valued fully random hash functions. Can't be implemented...

Our algorithm uses continuous valued fully random hash functions. Can't be implemented...

- The idea of using the minimum hash value of $x_1, \ldots, x_n$ to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.

Our algorithm uses continuous valued fully random hash functions.
Can't be implemented...

- The idea of using the minimum hash value of $x_1, \ldots, x_n$ to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.
- Flajolet-Martin (LogLog) algorithm and HyperLogLog.

Our algorithm uses continuous valued fully random hash functions. Can't be implemented...

- The idea of using the minimum hash value of $x_1, \ldots, x_n$ to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.
- Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| | |
|---|---|
| $h(x_1)$ | 1010010 |
| $h(x_2)$ | 1001100 |
| $h(x_3)$ | 1001110 |
| $\vdots$ | |
| $h(x_n)$ | 1011000 |

Our algorithm uses continuous valued fully random hash functions. Can't be implemented...

- The idea of using the minimum hash value of $x_1, \ldots, x_n$ to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.
- Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| | |
|---|---|
| $h(x_1)$ | **1010010** |
| $h(x_2)$ | **1001100** |
| $h(x_3)$ | **1001110** |
| $\vdots$ | |
| $h(x_n)$ | **1011000** |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

Our algorithm uses continuous valued fully random hash functions. Can't be implemented…

- The idea of using the minimum hash value of $x_1, \ldots, x_n$ to estimate the number of distinct elements naturally extends to when the hash functions map to discrete values.
- Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| | |
|---|---|
| $h(x_1)$ | **1010010** |
| $h(x_2)$ | **1001100** |
| $h(x_3)$ | **1001110** |
| | . . . |
| $h(x_n)$ | **1011000** |

Estimate # distinct elements based on maximum number of trailing zeros **m**.
The more distinct hashes we see, the higher we expect this maximum to be.

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| $h(x_1)$ | **1010010** |
|----------|-------------|
| $h(x_2)$ | **1001100** |
| $h(x_3)$ | **1001110** |
| ⋮ | |
| $h(x_n)$ | **1011000** |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| | |
|---|---|
| $h(x_1)$ | **1010010** |
| $h(x_2)$ | **1001100** |
| $h(x_3)$ | **1001110** |
| | |
| | ⋮ |
| | |
| $h(x_n)$ | **1011000** |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

With $d$ distinct elements, roughly what do we expect **m** to be?

a) $O(1)$    b) $O(\log d)$    c) $O(\sqrt{d})$    d) $O(d)$

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| | |
|---|---|
| $h(x_1)$ | **101001**0 |
| $h(x_2)$ | **10011**00 |
| $h(x_3)$ | **1001110** |
| | |
| $\vdots$ | |
| | |
| $h(x_n)$ | **1011**000 |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

With $d$ distinct elements, roughly what do we expect **m** to be?

a) $O(1)$    b) $O(\log d)$    c) $O(\sqrt{d})$    d) $O(d)$

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| | |
|---|---|
| $h(x_1)$ | **101001**0 |
| $h(x_2)$ | **1001**100 |
| $h(x_3)$ | **1001110** |
| | |
| ⋮ | ⋮ |
| | |
| $h(x_n)$ | **1011**000 |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

With $d$ distinct elements, roughly what do we expect **m** to be?

$$\Pr(\mathbf{h}(x_i) \text{ has } x \text{ trailing zeros}) =$$

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| $h(x_1)$ | **1010010** |
|---|---|
| $h(x_2)$ | **1001100** |
| $h(x_3)$ | **1001110** |
| $\vdots$ | |
| $h(x_n)$ | **1011000** |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

With $d$ distinct elements, roughly what do we expect **m** to be?

$$\Pr(\mathbf{h}(x_i) \text{ has } x \text{ trailing zeros}) = \frac{1}{2^x}$$

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| | |
|---|---|
| $h(x_1)$ | **101001**0 |
| $h(x_2)$ | **1001**1**00** |
| $h(x_3)$ | **1001110** |
| | |
| | . |
| | . |
| | . |
| | |
| | |
| $h(x_n)$ | **1011**000 |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

With $d$ distinct elements, roughly what do we expect **m** to be?

$$\Pr(\mathbf{h}(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}}$$

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| $h(x_1)$ | **101001**<span style="color:blue">**0**</span> |
|---|---|
| $h(x_2)$ | **1001**<span style="color:blue">**100**</span> |
| $h(x_3)$ | **1001110** |
| | |
| | ⋮ |
| | |
| $h(x_n)$ | **1011**<span style="color:blue">**000**</span> |

Estimate # distinct elements based on maximum number of trailing zeros $\mathbf{m}$.

With $d$ distinct elements, roughly what do we expect $\mathbf{m}$ to be?

$$\Pr(\mathbf{h}(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| $\mathbf{h}(x_1)$ | **1010010** |
| $\mathbf{h}(x_2)$ | **1001100** |
| $\mathbf{h}(x_3)$ | **1001110** |
| | $\vdots$ |
| $\mathbf{h}(x_n)$ | **1011000** |

Estimate # distinct elements based on maximum number of trailing zeros $\mathbf{m}$.

With $d$ distinct elements, roughly what do we expect $\mathbf{m}$ to be?

$$\Pr(\mathbf{h}(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

So with $d$ distinct hashes, expect to see 1 with $\log d$ trailing zeros. Expect $\mathbf{m} \approx \log d$.

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| | |
|---|---|
| $h(x_1)$ | **1010010** |
| $h(x_2)$ | **1001100** |
| $h(x_3)$ | **1001110** |
| $\vdots$ | |
| $h(x_n)$ | **1011000** |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

With $d$ distinct elements, roughly what do we expect **m** to be?

$$\Pr(\mathbf{h}(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

So with $d$ distinct hashes, expect to see 1 with $\log d$ trailing zeros. Expect **m** $\approx \log d$. **m** takes $\log \log d$ bits to store.

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| $\mathbf{h}(x_1)$ | **1010010** |
|---|---|
| $\mathbf{h}(x_2)$ | **1001100** |
| $\mathbf{h}(x_3)$ | **1001110** |
| $\vdots$ | |
| $\mathbf{h}(x_n)$ | **1011000** |

Estimate # distinct elements based on maximum number of trailing zeros $\mathbf{m}$.

With $d$ distinct elements, roughly what do we expect $\mathbf{m}$ to be?

$$\Pr(\mathbf{h}(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

So with $d$ distinct hashes, expect to see 1 with $\log d$ trailing zeros. Expect $\mathbf{m} \approx \log d$. $\mathbf{m}$ takes $\log \log d$ bits to store.

**Total Space:** $O\left(\frac{\log \log d}{\epsilon^2} + \log d\right)$ for an $\epsilon$ approximate count.

Flajolet-Martin (LogLog) algorithm and HyperLogLog.

| $h(x_1)$ | **101001**<span style="color:blue">0</span> |
| $h(x_2)$ | **1001**<span style="color:blue">00</span> |
| $h(x_3)$ | **1001110** |
| $\vdots$ | |
| $h(x_n)$ | **1011**<span style="color:blue">000</span> |

Estimate # distinct elements based on maximum number of trailing zeros **m**.

With $d$ distinct elements, roughly what do we expect **m** to be?

$$\Pr(h(x_i) \text{ has } \log d \text{ trailing zeros}) = \frac{1}{2^{\log d}} = \frac{1}{d}.$$

So with $d$ distinct hashes, expect to see 1 with $\log d$ trailing zeros. Expect **m** $\approx \log d$. **m** takes $\log \log d$ bits to store.

**Total Space:** $O\left(\frac{\log \log d}{\epsilon^2} + \log d\right)$ for an $\epsilon$ approximate count.

**Note:** Careful averaging of estimates from multiple hash functions.

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\text{space used } = O\left(\frac{\log \log d}{\epsilon^2} + \log d\right)$$

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\text{space used } = O\left(\frac{\log\log d}{\epsilon^2} + \log d\right)$$

$$= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} + \lceil \log_2 d \rceil \text{ bits}[1]$$

---

[1]. 1.04 is the constant in the HyperLogLog analysis. Not important!

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\text{space used } = O\left(\frac{\log \log d}{\epsilon^2} + \log d\right)$$

$$= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} + \lceil \log_2 d \rceil \text{ bits}^{[1]}$$

$$= \frac{1.04 \cdot 5}{.02^2} + 30 = 13030 \text{ bits} \approx 1.6 \text{ } kB!$$

---

[1]. 1.04 is the constant in the HyperLogLog analysis. Not important!

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\text{space used } = O\left(\frac{\log \log d}{\epsilon^2} + \log d\right)$$

$$= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} + \lceil \log_2 d \rceil \text{ bits}^1$$

$$= \frac{1.04 \cdot 5}{.02^2} + 30 = 13030 \text{ bits} \approx 1.6 \text{ kB!}$$

**Mergeable Sketch:** Consider the case (essentially always in practice) that the items are processed on different machines.

---

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\text{space used} = O\left(\frac{\log \log d}{\epsilon^2} + \log d\right)$$

$$= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} + \lceil \log_2 d \rceil \text{ bits}[1]$$

$$= \frac{1.04 \cdot 5}{.02^2} + 30 = 13030 \text{ bits} \approx 1.6 \text{ kB}!$$

**Mergeable Sketch:** Consider the case (essentially always in practice) that the items are processed on different machines.

• Given data structures (sketches) $HLL(x_1, \ldots, x_n)$, $HLL(y_1, \ldots, y_n)$ it is easy to merge them to give $HLL(x_1, \ldots, x_n, y_1, \ldots, y_n)$.

---

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\text{space used} = O\left(\frac{\log\log d}{\epsilon^2} + \log d\right)$$

$$= \frac{1.04 \cdot \lceil \log_2 \log_2 d \rceil}{\epsilon^2} + \lceil \log_2 d \rceil \text{ bits}^1$$

$$= \frac{1.04 \cdot 5}{.02^2} + 30 = 13030 \text{ bits} \approx 1.6 \text{ kB}!$$

**Mergeable Sketch:** Consider the case (essentially always in practice) that the items are processed on different machines.

- Given data structures (sketches) $HLL(x_1, \ldots, x_n)$, $HLL(y_1, \ldots, y_n)$ it is easy to merge them to give $HLL(x_1, \ldots, x_n, y_1, \ldots, y_n)$. How?

---

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

Using HyperLogLog to count 1 billion distinct items with 2% accuracy:

$$\text{space used } = O\left(\frac{\log\log d}{\epsilon^2} + \log d\right)$$

$$= \frac{1.04 \cdot \lceil\log_2\log_2 d\rceil}{\epsilon^2} + \lceil\log_2 d\rceil \text{ bits}[1]$$

$$= \frac{1.04 \cdot 5}{.02^2} + 30 = 13030 \text{ bits} \approx 1.6 \text{ kB!}$$

**Mergeable Sketch:** Consider the case (essentially always in practice) that the items are processed on different machines.

- Given data structures (sketches) $HLL(x_1, \ldots, x_n)$, $HLL(y_1, \ldots, y_n)$ it is easy to merge them to give $HLL(x_1, \ldots, x_n, y_1, \ldots, y_n)$. How?

- Set the maximum # of trailing zeros to the maximum in the two sketches.

---

1. 1.04 is the constant in the HyperLogLog analysis. Not important!

Questions on distinct elements counting?

**Jaccard Index:** A similarity measure between two sets.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\# \text{ shared elements}}{\# \text{ total elements}}.$$



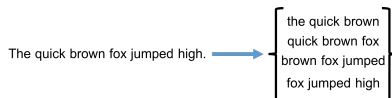Natural measure for similarity between bit strings – interpret an $n$ bit string as a set, containing the elements corresponding the positions of its ones. $J(x, y) = \frac{\# \text{ shared ones}}{\text{total ones}}$.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\#\ \text{shared elements}}{\#\ \text{total elements}}.$$

**Want Fast Implementations For:**

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\text{\# shared elements}}{\text{\# total elements}}.$$

**Want Fast Implementations For:**

- **Near Neighbor Search:** Have a database of $n$ sets/bit strings and given a set $A$, want to find if it has high Jaccard similarity to anything in the database. $\Omega(n)$ time with a linear scan.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\# \text{ shared elements}}{\# \text{ total elements}}.$$

**Want Fast Implementations For:**

- **Near Neighbor Search:** Have a database of $n$ sets/bit strings and given a set $A$, want to find if it has high Jaccard similarity to anything in the database. $\Omega(n)$ time with a linear scan.

- **All-pairs Similarity Search:** Have $n$ different sets/bit strings and want to find all pairs with high Jaccard similarity. $\Omega(n^2)$ time if we check all pairs explicitly.

Will speed up via randomized locality sensitive hashing.

**Document Similarity:**

- E.g., detecting plagiarism, copyright infringement, spam.

- Use Shingling + Jaccard similarity.

The quick brown fox jumped high. $\longrightarrow$
$$\begin{bmatrix} \text{the quick brown} \\ \text{quick brown fox} \\ \text{brown fox jumped} \\ \text{fox jumped high} \end{bmatrix}$$

**Document Similarity:**

- E.g., detecting plagiarism, copyright infringement, spam.
- Use Shingling + Jaccard similarity. ($n$-grams, $k$-mers)



The quick brown fox jumped high. →
the quick brown
quick brown fox
brown fox jumped
fox jumped high

ATAGCCGTAGT →
ATAG  CCGT
TAGC  CGTA
AGCC  GTAG
GCCG  TAGT

Online recommendation systems are often based on **collaborative filtering**. Simplest approach: find similar users and make recommendations based on those users.

Online recommendation systems are often based on **collaborative filtering**. Simplest approach: find similar users and make recommendations based on those users.



- Twitter: represent a user as the set of accounts they follow. Match similar users based on the Jaccard similarity of these sets. Recommend that you follow accounts followed by similar users. Netflix: look at sets of movies watched. Amazon: look at products purchased, etc.

**Entity Resolution Problem:** Want to combine records from multiple data sources that refer to the same entities.

## APPLICATION: ENTITY RESOLUTION

**Entity Resolution Problem:** Want to combine records from multiple data sources that refer to the same entities.

- E.g. data on individuals from voting registrations, property records, and social media accounts. Names and addresses may not exactly match, due to typos, nicknames, moves, etc.

- Still want to match records that all refer to the same person using all pairs similarity search.

## APPLICATION: ENTITY RESOLUTION

**Entity Resolution Problem:** Want to combine records from multiple data sources that refer to the same entities.

- E.g. data on individuals from voting registrations, property records, and social media accounts. Names and addresses may not exactly match, due to typos, nicknames, moves, etc.

- Still want to match records that all refer to the same person using all pairs similarity search.

See Section 3.8.2 of *Mining Massive Datasets* for a discussion of a real world example involving 1 million customers. Naively this would be $\binom{1000000}{2} \approx 500$ billion pairs of customers to check!

Many applications to spam/fraud detection. E.g.

Many applications to spam/fraud detection. E.g.

- **Fake Reviews**: Very common on websites like Amazon.
  Detection often looks for (near) duplicate reviews on similar
  products, which have been copied. 'Near duplicate' measured
  with shingles + Jaccard similarity.

Many applications to spam/fraud detection. E.g.

- **Fake Reviews**: Very common on websites like Amazon. Detection often looks for (near) duplicate reviews on similar products, which have been copied. 'Near duplicate' measured with shingles + Jaccard similarity.
- **Lateral phishing**: Phishing emails sent to addresses at a business coming from a legitimate email address at the same business that has been compromised.
  - One method of detection looks at the recipient list of an email and checks if it has small Jaccard similarity with any previous recipient lists. If not, the email is flagged as possible spam.

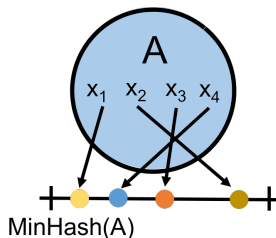**Goal:** Speed up Jaccard similarity search (near neighbor and all-pairs similarity search).

**Goal:** Speed up Jaccard similarity search (near neighbor and all-pairs similarity search).

**Strategy:** Use random hashing to map each set to a very compressed representation. Jaccard similarity can be estimated from these representations.

**Goal:** Speed up Jaccard similarity search (near neighbor and all-pairs similarity search).

**Strategy:** Use random hashing to map each set to a very compressed representation. Jaccard similarity can be estimated from these representations.

**MinHash(A):** [Andrei Broder, 1997 at Altavista]

- Let $\mathbf{h} : U \to [0, 1]$ be a random hash function

- $\mathbf{s} := 1$

- For $x_1, \ldots, x_{|A|} \in A$

  - $\mathbf{s} := \min(\mathbf{s}, \mathbf{h}(x_k))$

- Return $\mathbf{s}$

**Goal:** Speed up Jaccard similarity search (near neighbor and all-pairs similarity search).

**Strategy:** Use random hashing to map each set to a very compressed representation. Jaccard similarity can be estimated from these representations.

**MinHash(A):** [Andrei Broder, 1997 at Altavista]

- Let $\mathbf{h} : U \to [0, 1]$ be a random hash function

- $\mathbf{s} := 1$

- For $x_1, \ldots, x_{|A|} \in A$

  - $\mathbf{s} := \min(\mathbf{s}, \mathbf{h}(x_k))$

- Return $\mathbf{s}$



MinHash(A)

**Goal:** Speed up Jaccard similarity search (near neighbor and all-pairs similarity search).

**Strategy:** Use random hashing to map each set to a very compressed representation. Jaccard similarity can be estimated from these representations.

**MinHash(A):** [Andrei Broder, 1997 at Altavista]

- Let $\mathbf{h} : U \to [0, 1]$ be a random hash function

- $\mathbf{s} := 1$

- For $x_1, \ldots, x_{|A|} \in A$

  - $\mathbf{s} := \min(\mathbf{s}, \mathbf{h}(x_k))$

- Return $\mathbf{s}$



MinHash(A)

Identical to our distinct elements sketch!

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

- Since we are hashing into the continuous range $[0, 1]$, we will never have $\mathbf{h}(x) = \mathbf{h}(y)$ for $x \neq y$ (i.e., no spurious collisions)
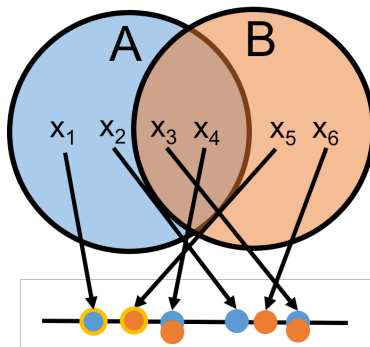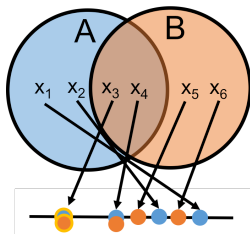
For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

- Since we are hashing into the continuous range $[0, 1]$, we will never have $\mathbf{h}(x) = \mathbf{h}(y)$ for $x \neq y$ (i.e., no spurious collisions)

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

- Since we are hashing into the continuous range $[0, 1]$, we will never have $\mathbf{h}(x) = \mathbf{h}(y)$ for $x \neq y$ (i.e., no spurious collisions)

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

- Since we are hashing into the continuous range $[0, 1]$, we will never have $\mathbf{h}(x) = \mathbf{h}(y)$ for $x \neq y$ (i.e., no spurious collisions)

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

- Since we are hashing into the continuous range $[0, 1]$, we will never have $\mathbf{h}(x) = \mathbf{h}(y)$ for $x \neq y$ (i.e., no spurious collisions)

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

- Since we are hashing into the continuous range $[0, 1]$, we will never have $\mathbf{h}(x) = \mathbf{h}(y)$ for $x \neq y$ (i.e., no spurious collisions)

- $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.

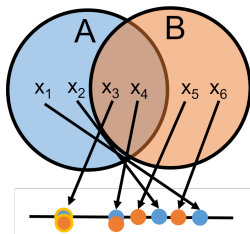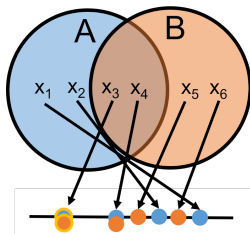For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



$\Pr(MinHash(A) = MinHash(B)) = ?$

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



$$\Pr(MinHash(A) = MinHash(B)) = \frac{|A \cap B|}{\text{total \# items hashed}}$$

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

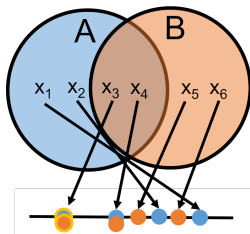**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



$$\Pr(MinHash(A) = MinHash(B)) = \frac{|A \cap B|}{\text{total \# items hashed}}$$
$$= \frac{|A \cap B|}{|A \cup B|}$$

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



$$\Pr(MinHash(A) = MinHash(B)) = \frac{|A \cap B|}{\text{total } \# \text{ items hashed}}$$
$$= \frac{|A \cap B|}{|A \cup B|} = J(A, B).$$