



NYU

Center for
Data Science

Week 08.1: Search and MinHash

DS-GA 1004: Big Data

Finding items in a large collection

- **Search** and **recommendation** rely on similarity calculation
- User provides a “**query**”, e.g.:
 - Search string
 - Example document
 - Latent representation
- System returns a list of matching “**documents**” from the database

Examples

- Text search: search string \Leftrightarrow the web (documents)
- Recommender systems: user representation \Leftrightarrow item representation
- Reverse image search: photo \Leftrightarrow library
- Copyright detection: uploaded video \Leftrightarrow all of youtube

Basic approach

- Given a query q
- For each document d in collection
 - Compute **similarity**(q, d)
- Order collection by decreasing **similarity**
- Return top k documents

How can we do
this efficiently?

Problems of scale

- Size of collection (N)
- Dimension of representation
- More generally: complexity of computing similarity
- *Can we do better than brute force search?*

Approximate search

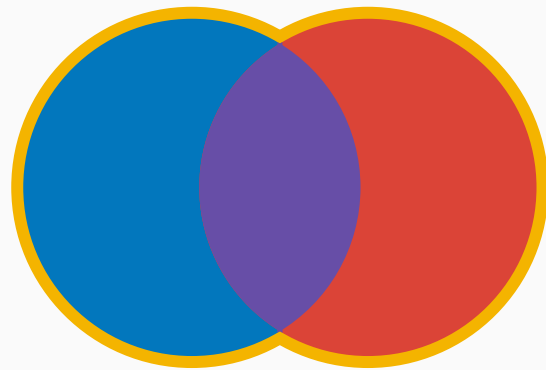
- Use a **fast method** to identify $n \ll N$ **candidate nearest neighbors**
- Use **true similarity** on **candidate set** to discard any false positives
- **Fast method** usually requires some kind of data structure
 - Search time should be strictly sub-linear: $n \sim o(N)$
 - e.g. $\log(N)$ or \sqrt{N}

MinHash

[Broder, 1997]

Jaccard similarity for sets

- Items are represented as **sets**, e.g.:
 - $A = \{\text{words contained in document } A\}$
 - $A = \{\text{users who interacted with item } A\}$
- **Jaccard similarity** is the ratio
$$J(A, B) = |A \cap B| / |A \cup B|$$
- $D(A, B) = 1 - J(A, B)$ is the **Jaccard distance metric**



MinHash

- Fix a random ordering π of the elements
- Imagine a table of set memberships
- For each set \mathbf{S} , its hash is:

$$h(\mathbf{S} \mid \pi) = \min \{k \mid \pi(k) \in \mathbf{S}\}$$

- Index of the first (permuted) item belonging to set \mathbf{S}

$\pi(k)$		Doc 1	Doc 2	...
1	Item 3	1	0	
2	Item 75	0	0	
3	Item 21	0	1	
4	Item 1	0	0	
5	Item 2004	0	1	

$$h(\text{Doc 1} \mid \pi) = 1$$

$$h(\text{Doc 2} \mid \pi) = 3$$

Permutation indexing

A = {"T.rex", "Stegosaurus", "PDP-11"}

B = {"Apples", "Bananas", "Pine cones"}

C = {"T.rex", "Bananas", "Penguins"}

D = {"Apples", "Turtles", "Pine cones"}

$\pi(\mathbf{k})$		A	B	C	D
1	PDP-11	1	0	0	0
2	Penguins	0	0	1	0
3	Pine cones	0	1	0	1
4	Turtles	0	0	0	1
5	Apples	0	1	0	1
6	T.rex	1	0	1	0
7	Bananas	0	1	1	0
8	Stegosaurus	1	0	0	0

Permutation indexing

$A = \{\text{"T.rex", "Stegosaurus", "PDP-11"}\} \rightarrow h(A|\pi) = 1$

$B = \{\text{"Apples", "Bananas", "Pine cones"}\} \rightarrow h(B|\pi) = 3$

$C = \{\text{"T.rex", "Bananas", "Penguins"}\} \rightarrow h(C|\pi) = 2$

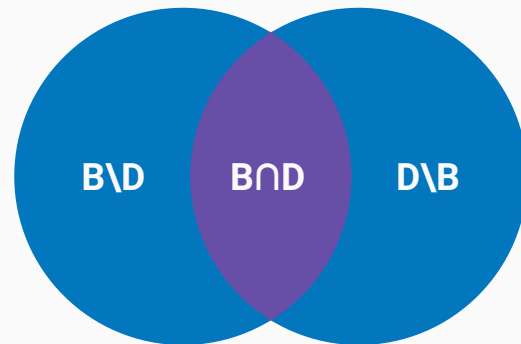
$D = \{\text{"Apples", "Turtles", "Pine cones"}\} \rightarrow h(D|\pi) = 3$

Hash collision is more likely when sets overlap.
Let's analyze this!

$\pi(k)$		A	B	C	D
1	PDP-11	1	0	0	0
2	Penguins	0	0	1	0
3	Pine cones	0	1	0	1
4	Turtles	0	0	0	1
5	Apples	0	1	0	1
6	T.rex	1	0	1	0
7	Bananas	0	1	1	0
8	Stegosaurus	1	0	0	0

$$P[h(S_1) = h(S_2)] = \text{Jaccard}(S_1, S_2)$$

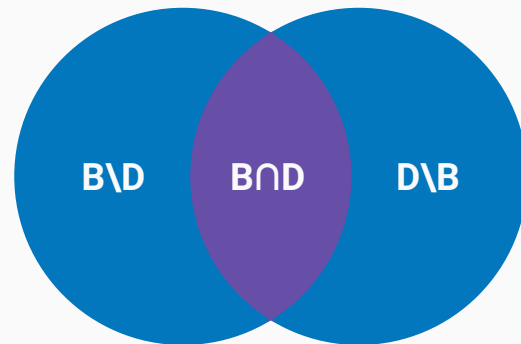
- For two sets S_1 and S_2 , there are three types of rows:
 - Type 1:** $\pi(k) \in S_1 \cap S_2$
 - Type 2:** $\pi(k) \in S_1 \Delta S_2$
 - Type 3: $\pi(k) \notin S_1 \cup S_2$
- Collision** \Leftrightarrow **type 1 row** occurs before all **type 2 rows**



$\pi(k)$		A	B	C	D
1	PDP-11	1	0	0	0
2	Penguins	0	0	1	0
3	Pine cones	0	1	0	1
4	Turtles	0	0	0	1
5	Apples	0	1	0	1
6	T.rex	1	0	1	0
7	Bananas	0	1	1	0
8	Stegosaurus	1	0	0	0

$$P[h(S_1) = h(S_2)] = \text{Jaccard}(S_1, S_2)$$

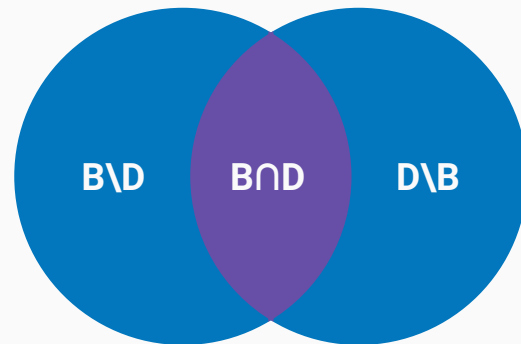
- For two sets S_1 and S_2 , there are three types of rows:
 - Type 1:** $\pi(k) \in S_1 \cap S_2$
 - Type 2:** $\pi(k) \in S_1 \Delta S_2$
 - Type 3: $\pi(k) \notin S_1 \cup S_2$
- Collision** \Leftrightarrow **type 1 row** occurs before all **type 2 rows**
- $P[\text{Collision}] = (\# \text{ Type 1}) / (\# \text{ Type 1} + \# \text{ Type 2})$



	$\pi(k)$	A	B	C	D
1	PDP-11	1	0	0	0
2	Penguins	0	0	1	0
3	Pine cones	0	1	0	1
4	Turtles	0	0	0	1
5	Apples	0	1	0	1
6	T.rex	1	0	1	0
7	Bananas	0	1	1	0
8	Stegosaurus	1	0	0	0

$$P[h(S_1) = h(S_2)] = \text{Jaccard}(S_1, S_2)$$

- For two sets S_1 and S_2 , there are three types of rows:
 - Type 1:** $\pi(k) \in S_1 \cap S_2$
 - Type 2:** $\pi(k) \in S_1 \Delta S_2$
 - Type 3: $\pi(k) \notin S_1 \cup S_2$
- Collision** \Leftrightarrow **type 1 row** occurs before all **type 2 rows**
- $P[\text{Collision}] = (\# \text{ Type 1}) / (\# \text{ Type 1} + \# \text{ Type 2})$
 $= |S_1 \cap S_2| / (|S_1 \cap S_2| + |S_1 \Delta S_2|)$



	$\pi(k)$	A	B	C	D
1	PDP-11	1	0	0	0
2	Penguins	0	0	1	0
3	Pine cones	0	1	0	1
4	Turtles	0	0	0	1
5	Apples	0	1	0	1
6	T.rex	1	0	1	0
7	Bananas	0	1	1	0
8	Stegosaurus	1	0	0	0

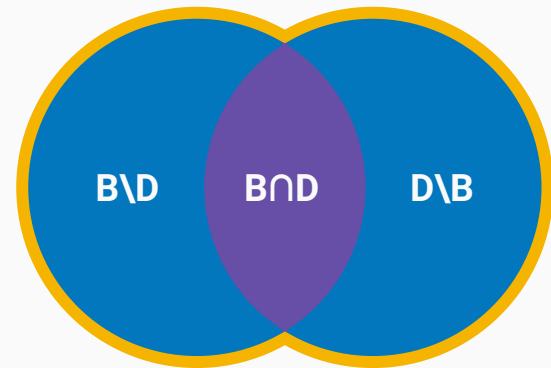
$$P[h(S_1) = h(S_2)] = \text{Jaccard}(S_1, S_2)$$

- For two sets S_1 and S_2 , there are three types of rows:

- **Type 1:** $\pi(k) \in S_1 \cap S_2$
- **Type 2:** $\pi(k) \in S_1 \Delta S_2$
- Type 3: $\pi(k) \notin S_1 \cup S_2$

- Collision** \Leftrightarrow **type 1 row** occurs before all **type 2 rows**

- $P[\text{Collision}] = (\# \text{ Type 1}) / (\# \text{ Type 1} + \# \text{ Type 2})$
 $= |S_1 \cap S_2| / (|S_1 \cap S_2| + |S_1 \Delta S_2|)$
 $= |S_1 \cap S_2| / (|S_1 \cup S_2|)$
 $= J(S_1, S_2) \blacksquare$



	$\pi(k)$	A	B	C	D
1	PDP-11	1	0	0	0
2	Penguins	0	0	1	0
3	Pine cones	0	1	0	1
4	Turtles	0	0	0	1
5	Apples	0	1	0	1
6	T.rex	1	0	1	0
7	Bananas	0	1	1	0
8	Stegosaurus	1	0	0	0

Monte carlo approximation

- **B** and **D** had **P[collision]** = $\frac{1}{2}$ for **single permutation** π
- But we want the **probability of collision** over **all choices** of π
- **Idea:**
 - Generate m **random permutations** $\pi_1, \pi_2, \dots, \pi_m$
 - **Count hash collisions** between A and B over all π_i 's
 - $J(\mathbf{A}, \mathbf{B}) \approx \# \text{ collisions} / m$

$\pi(\mathbf{k})$		A	B	C	D
1	PDP-11	1	0	0	0
2	Penguins	0	0	1	0
3	Pine cones	0	1	0	1
4	Turtles	0	0	0	1
5	Apples	0	1	0	1
6	T.rex	1	0	1	0
7	Bananas	0	1	1	0
8	Stegosaurus	1	0	0	0



$h(\mathbf{S} \pi_i)$	A	B	C	D
π_1	0	1	2	1
π_2	2	1	2	5
π_3	0	0	3	6
π_4	9	5	1	1
...				

MinHash signatures

Searching with MinHash

- User provides query q
- **Candidates** $\leftarrow \{ \}$
- For each π_i :
 - Compute $h(q \mid \pi_i)$
 - Candidates \leftarrow Candidates + **colliding documents** $S : h(q \mid \pi_i) = h(S \mid \pi_i)$
- Return **Candidates** ordered by
 - # Collisions (estimated similarity score), or
 - Or full Jaccard similarity (more accurate, but slower)

$h(S \mid \pi_i)$	A	B	C	D
π_1	0	1	2	1
π_2	2	1	2	5
π_3	0	0	3	6
π_4	9	5	1	1
...				

Note that we don't to compare to the full collection!

Only those that collide.

Extension to bags / Ruzicka approximation

[Chen, Philbin, Zisserman 2008]

- Idea: reduce bags to sets by uniquely identifying each repetition

{dog}	→ {dog ₁ }
{dog, dog}	→ {dog ₁ , dog ₂ }
{dog, dog, dog}	→ {dog ₁ , dog ₂ , dog ₃ }

- Jaccard on expanded sets = **Ruzicka similarity** on original bags

$$R(\mathbf{A}, \mathbf{B}) = \frac{\sum_i \min(\mathbf{A}[i], \mathbf{B}[i])}{\sum_j \max(\mathbf{A}[j], \mathbf{B}[j])}$$

Improving over word counts

- Word n-grams:
 - “T.rex stomps loudly” → {“T.rex”, “stomps”, “loudly”, “T.rex stomps”, “stomps loudly”}
- Character shingles:
 - “T.rex stomps loudly” → {“T.rex st”, “.rex sto”, “rex stom”, “ex stomp”, ...}

Efficient approximation

- Permuting the entire universe is **expensive**
 - ... and it would not support growing item sets!
- Instead, replace *permutations* π_i with *hashes* H_i
 - A permutation is a **perfect hash**: distinct elements cannot collide
 - Approximate this by an **imperfect hash**: distinct elements *may collide*
- As long as H_i are unlikely to collide, this can still work

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	∞	∞	∞	∞
H_2	∞	∞	∞	∞

Signature table is initialized to infinity for each entry

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0 ∞	∞	∞	∞
H_2	0 ∞	∞	∞	∞

A, H_1 : $0 < \infty \rightarrow \text{update}$
A, H_2 : $0 < \infty \rightarrow \text{update}$

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0	∞	1 $\leftrightarrow \infty$	∞
H_2	0	∞	2 $\leftrightarrow \infty$	∞

C, H_1 : $1 < \infty \rightarrow \text{update}$
C, H_2 : $2 < \infty \rightarrow \text{update}$

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0	2 ∞	1	2 ∞
H_2	0	4 ∞	2	4 ∞

B, H_1 : $2 < \infty \rightarrow$ update
B, H_2 : $4 < \infty \rightarrow$ update
D, H_1 : $2 < \infty \rightarrow$ update
D, H_2 : $4 < \infty \rightarrow$ update

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0	2	1	2
H_2	0	4	2	0 4

D, H_1 : $3 > 2 \rightarrow$ no update
D, H_2 : $0 < 4 \rightarrow$ update

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0	0 2	1	0 2
H_2	0	1 3	2	0

B, H_1 : $0 < 2 \rightarrow$ update
B, H_2 : $1 < 3 \rightarrow$ update
D, H_1 : $0 < 2 \rightarrow$ update
D, H_2 : $1 > 0 \rightarrow$ no update

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0	0	1	0
H_2	0	1	2	0

A, H_1 : $1 > 0 \rightarrow$ no update
A, H_2 : $3 > 0 \rightarrow$ no update
C, H_1 : $1 = 1 \rightarrow$ no update
C, H_2 : $3 > 2 \rightarrow$ no update

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0	0	1	0
H_2	0	1	2	0

B, H_1 : $2 > 0 \rightarrow$ no update
B, H_2 : $5 > 1 \rightarrow$ no update
C, H_1 : $2 > 1 \rightarrow$ no update
C, H_2 : $5 > 2 \rightarrow$ no update

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0	0	1	0
H_2	0	1	2	0

A, H_1 : $3 > 0 \rightarrow$ no update

A, H_2 : $1 > 0 \rightarrow$ no update

MinHash example

x	$H_1(x)$	$H_2(x)$	A	B	C	D
PDP-11	0	0	1			
Penguins	1	2			1	
Pine cones	2	4		1		1
Turtles	3	0				1
Apples	0	1		1		1
T.rex	1	3	1		1	
Bananas	2	5		1	1	
Stegosaurus	3	1	1			

	A	B	C	D
H_1	0	0	1	0
H_2	0	1	2	0

Collisions:

- $H_1: A \equiv B \equiv D \neq C$
- $H_2: A \equiv D \neq B \neq C$

Failure modes of MinHash

- Permutation MinHash:
 - **Collisions** are **more likely** when a **small set of items** are shared across **many documents**
 - ⇒ “Stop-words” can be deadly! *“The”, “and”, “or”, etc...*
- **Hashing approximation** only makes things **worse**
 - Two distinct items can now hash to the same value
 - Collision probability only increases with the hashing approximation
 - **High collision likelihood** ⇒ **large candidate sets** ⇒ **slow retrieval**

How can we reduce the size of candidate sets, but maintain high recall?

... come back for part 2!