

Hashing

- SHA-1 hashing (recommended)

$$\text{hash}(\underbrace{\text{concat}(\text{salt}, x)}_{\text{string}}) \rightarrow [m]$$

- Multiplicative Hashing

$$h_a(x) = \lfloor m \cdot \text{fac}(x \cdot a) \rfloor$$

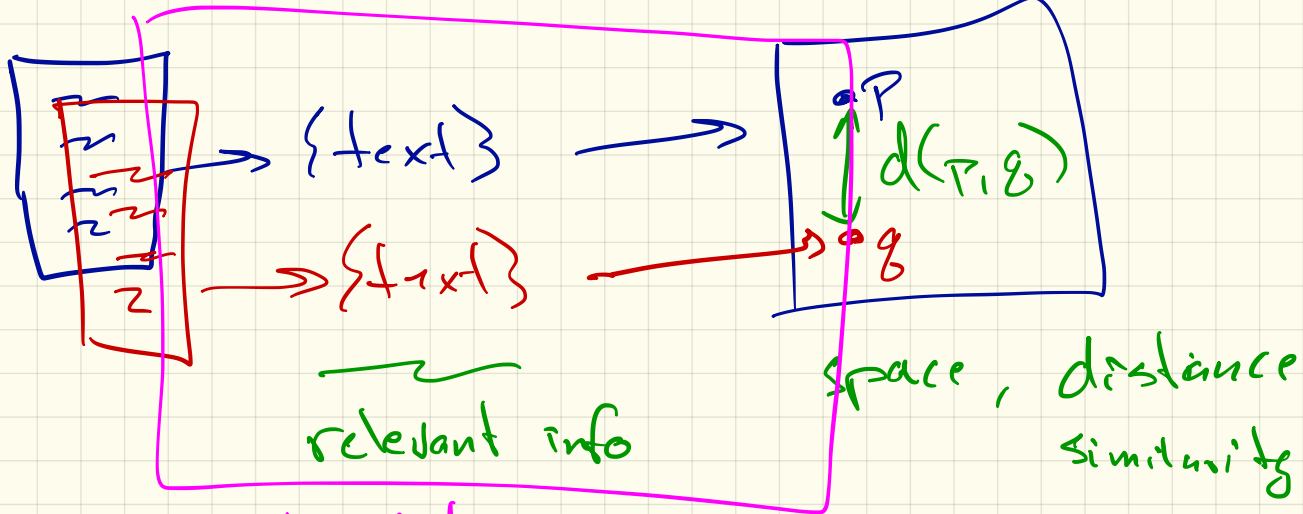
↑
deterministic

- ~~Modular $h(x) = x \pmod{m}$~~ don't

Distances between document (text)

Jaccard Distance, k -Grams

- Given 2 homeworks: did the plagiarizer close to copies?
- Given keyword in Google, which webpages are similar? Are 2 pages duplicates.
- emails? is it spam?



Modeling

Today

P, q are sets

Jaccard Distance

Alto-Vista

Common
 $P, q \in \mathbb{R}^d$

$P = (p_1, p_2, \dots, p_n)$

Sets

$$\{a, b, c\} = \{a, a, b, c\}$$

$$= \{c, b, a\}$$

$$\neq \{a, b\}$$

$$\neq \{a, b, c, d\}$$

Distance

$d(A, B)$

small \leftrightarrow A, B
close

large \leftrightarrow far

$0 \iff A = B$

$[0, \infty)$

$d(A, B) = 1 - s(A, B)$ (Jaccard)

$d(A, B) = \sqrt{s(A, A) + s(B, B) - 2s(A, B)}$

Similarity

$s(A, B)$

small \leftrightarrow A, B
far

large \leftrightarrow close

$1 \iff A = B$

$[0, 1]$

Jaccard Distance / Similarity

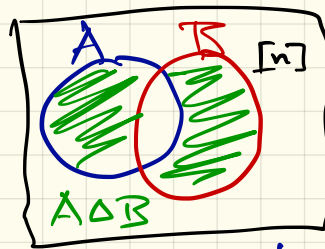
$$A = \{0, 1, 2, 5, 6\} \cup \{3\}$$

$$B = \{0, 2, 3, 5, 7, 9\}$$

$$JS(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|\{0, 2, 5\}|}{|\{0, 1, 2, 3, 5, 6, 7, 9\}|}$$

$$JD(A, B) = 1 - JS(A, B) = \frac{3}{8} = 0.375$$

Generalized Set Distances



- Hamming Sim

$$\text{Ham}(A, B) = \frac{|A \cap B| + |A \cup B|}{|A \cap B| + |A \cup B| + |A \Delta B|} = 1 - \frac{|A \Delta B|}{|[n]|}$$

- Andberg Sim

$$\text{Andb}(A, B) = \frac{|A \cap B|}{|A \cup B| + |A \Delta B|}$$

$$\text{SS}(A, B)$$

$$= S_{1,0,0,1}(A, B)$$

- Dice $(A, B) = \frac{2|A \cap B|}{|A| + |B|}$

$$S_{x,y,z,z'}(A, B) = \frac{x|A \cap B| + y|A \cup B| + z|A \Delta B|}{x|A \cap B| + y|A \cup B| + z'|A \Delta B|}$$

L3: Jaccard Similarity and k -Grams

Jeff M. Phillips

January 17, 2018

k-Grams with Words

k-gram

2-word gram

I am Sam.

Sam I am.

I do not like green eggs and ham.

I do not like them, Sam I am.

{ [I am], [am Sam], [Sam Sam],

k -Grams with Words

I am Sam.

Sam I am.

I do not like green eggs and ham.

I do not like them, Sam I am.

Words $k = 1$:

{[I], [am], [Sam], [do], [not], [like], [green],
[eggs], [and], [ham], [them]}

k -Grams with Words

I am Sam.

Sam I am.

I do not like green eggs and ham.

I do not like them, Sam I am.

Words $k = 1$:

{[I], [am], [Sam], [do], [not], [like], [green],
[eggs], [and], [ham], [them]}

Words $k = 2$:

{[I am], [am Sam], [Sam Sam], [Sam I], [am I], [I
do], [do not], [not like], [like green], [green
eggs], [eggs and], [and ham], [ham I], [like them],
[them Sam]}

k-Grams with Characters

$k = 3$

I am Sam.

Sam I am.

Characters $k = 3$:

{[iam], [ams], [msa], [sam], [ami], [mia]}

Modeling

- word vs. char
- value of k
- punctuation
- white space
- Capitalization

k -Grams with Characters

I am Sam.

Sam I am.

Characters $k = 3$:

{[iam], [ams], [msa], [sam], [ami], [mia]}

Characters $k = 4$:

{[iams], [amsa], [msam], [sams], [sami], [amia],
[miam]}

k -Grams and Jaccard

D_1 : I am Sam.

D_2 : Sam I am.

D_3 : I do not like green eggs and ham.

D_4 : I do not like them, Sam I am.

Words $k = 2$:

{[I am], [am Sam], [Sam Sam], [Sam I], [am I], [I do], [do not], [not like], [like green], [green eggs], [eggs and], [and ham], [like them], [them Sam]}

k-Grams and Jaccard

D_1 : [I am], [am Sam]

D_2 : [Sam I], [I am]

D_3 : [I do], [do not], [not like], [like green]
[green eggs], [eggs and], [and ham]

D_4 : [I do], [do not], [not like], [like them], [them Sam]
[Sam I], [I am]

k-Grams and Jaccard

D_1 : [I am], [am Sam]

D_2 : [Sam I], [I am]

D_3 : [I do], [do not], [not like], [like green]
[green eggs], [eggs and], [and ham]

D_4 : [I do], [do not], [not like], [like them], [them Sam]
[Sam I], [I am]

Jaccard Similarity: $JS(A, B) = \frac{|A \cap B|}{|A \cup B|}$

k-Grams and Jaccard

D_1 : [I am], [am Sam]

D_2 : [Sam I], [I am]

D_3 : [I do], [do not], [not like], [like green]
[green eggs], [eggs and], [and ham]

D_4 : [I do], [do not], [not like], [like them], [them Sam]
[Sam I], [I am]

Jaccard Similarity: $JS(A, B) = \frac{|A \cap B|}{|A \cup B|}$

$$JS(D_1, D_2) = 1/3 \approx 0.333$$

k-Grams and Jaccard

D_1 : [I am], [am Sam]

D_2 : [Sam I], [I am]

D_3 : [I do], [do not], [not like], [like green]
[green eggs], [eggs and], [and ham]

D_4 : [I do], [do not], [not like], [like them], [them Sam]
[Sam I], [I am]

Jaccard Similarity: $JS(A, B) = \frac{|A \cap B|}{|A \cup B|}$

$$JS(D_1, D_2) = \frac{1}{3} \approx 0.333$$

$$JS(D_1, D_3) = 0 = 0.0$$

k-Grams and Jaccard

D_1 : [I am], [am Sam]

D_2 : [Sam I], [I am]

D_3 : [I do], [do not], [not like], [like green]
[green eggs], [eggs and], [and ham]

D_4 : [I do], [do not], [not like], [like them], [them Sam]
[Sam I], [I am]

Jaccard Similarity: $JS(A, B) = \frac{|A \cap B|}{|A \cup B|}$

$$JS(D_1, D_2) = \frac{1}{3} \approx 0.333$$

$$JS(D_1, D_3) = \frac{0}{8} = 0.0$$

$$JS(D_1, D_4) = \frac{1}{8} = 0.125$$

k-Grams and Jaccard

D_1 : [I am], [am Sam]

D_2 : [Sam I], [I am]

D_3 : [I do], [do not], [not like], [like green]
[green eggs], [eggs and], [and ham]

D_4 : [I do], [do not], [not like], [like them], [them Sam]
[Sam I], [I am]

Jaccard Similarity: $JS(A, B) = \frac{|A \cap B|}{|A \cup B|}$

$$JS(D_1, D_2) = 1/3 \approx 0.333$$

$$JS(D_1, D_3) = 0 = 0.0$$

$$JS(D_1, D_4) = 1/8 = 0.125$$

$$JS(D_2, D_3) = 0 = 0.0$$

$$JS(D_2, D_4) = 2/7 \approx 0.286$$

$$JS(D_3, D_4) = 3/11 \approx 0.273$$

Bag-of-Words Model

[Sam | am] ← document D_1

$$D_1 \rightarrow P \in \mathbb{R}^m$$

$$P = (P_1, P_2, P_3, \dots, P_m)$$

Annotations: $P_1 = 2$ (labeled "2" = Sam), $P_3 = 1$ (labeled "3" = T), $P_5 = 1$ (labeled "5" = am)

$P_i = \#$ occurrences of word " i "

Annotations: "Sam" points to the first row, "2" points to the second element of the first row

$$P = \begin{pmatrix} 0 & 2 & 1 & 1 & 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 2 & 1 & 0 & 1 & 0 & 0 & \dots & 0 \end{pmatrix}$$

Continuous Bag of Words

I am Sam Sam I am I do not like green eggs and ham I
do not like them Sam I am