# lecture 7 : Operations on data sketches

$M_1 = (|A, m_1)$  $\qquad$ $M_2 = (|B, m_2)$

$$\begin{vmatrix} (1,\lambda_1) \\ (2,\lambda_2) \\ (3,\lambda_3) \\ (1,\lambda_1) \\ \vdots \end{vmatrix} \qquad\qquad \begin{vmatrix} (3,\lambda_3) \\ (4,\lambda_4) \\ (5,\lambda_5) \\ \vdots \end{vmatrix}$$

|A

$\begin{matrix} (1,\lambda_1) \\ (2,\lambda_2) \end{matrix}$ ... $\quad (3,\lambda_3) \quad (1,\lambda_4)$

$(5,\lambda_5)$  |B

$\downarrow$ $\qquad\qquad$ $\downarrow$

sketch A $\qquad$ sketch B

$|A|_w = \lambda_1 + \lambda_2 + \lambda_3$

---

$A = (A_1, A_2, \ldots, A_m)$ ,

$B = (B_1, B_2, \ldots, B_m)$ ,

$A_n = \min\{S_1^{(h)}, S_2^{(h)}, S_3^{(h)}\} = Exp(\lambda_1 + \lambda_2 + \lambda_3)$

$B_k = \min\{S_3^{(h)}, S_4^{(h)}, S_5^{(h)}\} = Exp(\lambda_3 + \lambda_4 + \lambda_5)$

$S_i^{(h)} \sim Exp(\lambda_i)$

$S_i^{(h)} = \dfrac{\ln(h(i \| k))}{-\lambda_i}$

# SUM

- $A, B \Rightarrow |A \cup B|_w = ?$ , $|X|_w = \sum\limits_{(i, \lambda_i) \in X} \lambda_i$

- Let $S_A = \{ S_i^{(k)} : (i, \lambda_i) \in A \}$,

  $S_B = \{ S_i^{(k)} : (i, \lambda_i) \in B \}$ and value of $k$ is unimportant it's just the number of experiment.

- Note that $\min\{ S_A \cup S_B \} = \min\{ \min\{S_A\}, \min\{S_B\} \}$

  $\underbrace{\min\{S_1, S_2, S_3\}}\in \underbrace{\min\{S_3, S_4, S_5\}}$

  thus sketch $C = (\min\{A_1, B_1\}, \dots \min\{A_m, B_m\})$ is exactly the same sketch as the sketch we would get by observing elements of $A \cup B$.

  generally $\min\{ S_1^{(1)}, S_2^{(1)}, S_3^{(1)}, S_4^{(1)}, S_5^{(1)} \}$

  $\Downarrow$

  $\min\{ \min\{S_1^{(1)}, S_2^{(1)}, S_3^{(1)}\}, \min\{S_3^{(1)}, S_4^{(1)}, S_5^{(1)}\}\}$

- Since sketch $C$ represents $A \cup B$, we can use estimator from the previous lecture:

$$\hat{\lambda} := \frac{m-1}{\sum\limits_{k=1}^{m} c_k} = \frac{m-1}{\sum\limits_{k=1}^{m} \min\{A_k, B_k\}}$$

- $E[\hat{\lambda}] = \lambda \rightarrow E\left[ \frac{m-1}{\sum \min\{A_k, B_k\}} \right] = |A \cup B|_w$

$$\rightarrow SE[\dots] = \frac{1}{\sqrt{m-2}}$$

This can be generized to any number of sets:

$$\min\{ S_A \cup S_B \cup S_C \dots \} = \min\{ \min\{S_A\}, \min\{S_B\}, \min\{S_C\} \dots \}$$

$\downarrow$

$$E\left[\frac{m-1}{\sum \min\{A_H, B_H, C_H \dots\}}\right] = |A \cup B \cup C \cup \dots|_w$$

---

## INTERSECTION

- Jaccard similarity $\quad J(A,B) = \dfrac{|A \cap B|}{|A \cup B|}$

- Weighted Jaccard similarity $\quad J_w(A,B) = \dfrac{|A \cap B|_w}{|A \cup B|_w}$

- Note that $|A \cap B|_w = |A \cup B|_w \cdot J_w(A,B)$

## Lemma 6

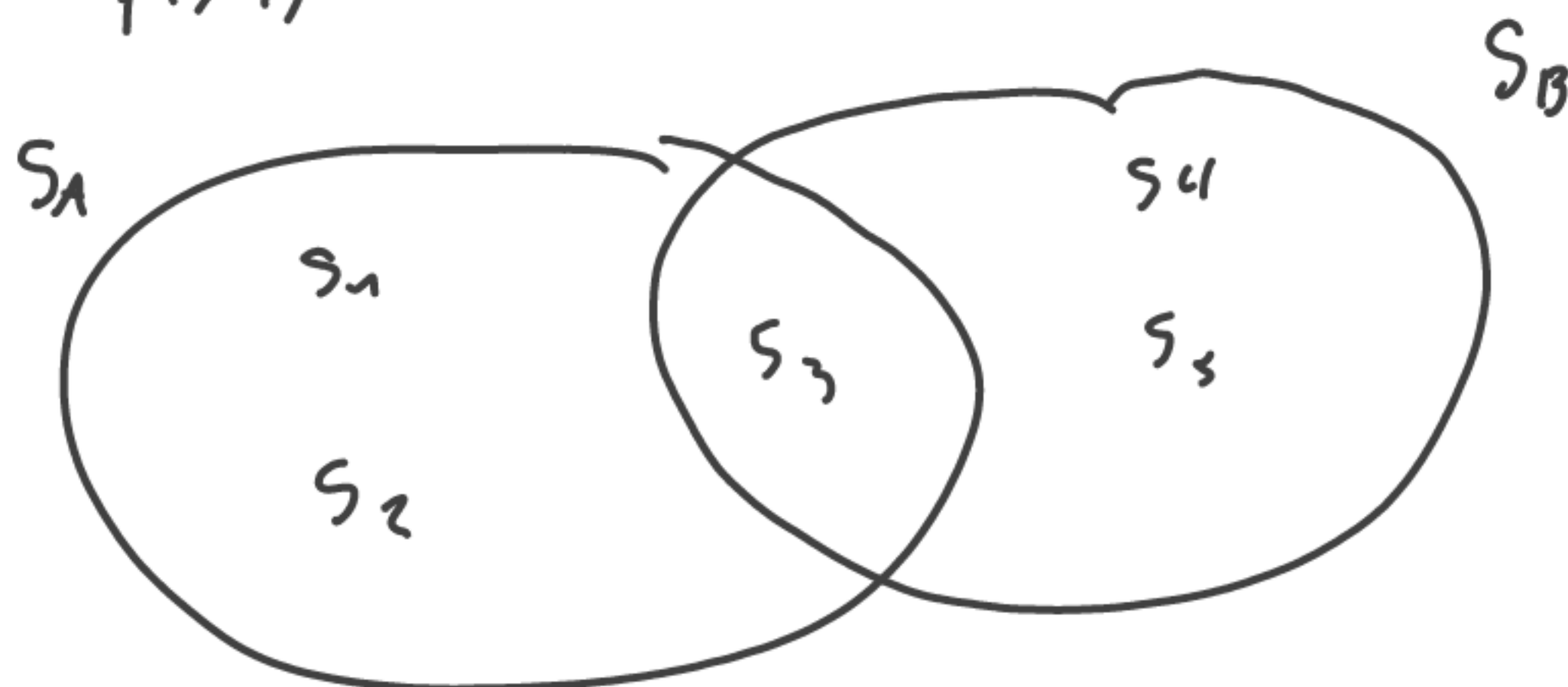$A \rightsquigarrow A = (A_1, A_2, \dots A_m), \quad B \rightsquigarrow B = (B_1, B_2, \dots B_m)$

$$\left(\forall_k \in \{1,2,3,\dots m\}\right)\left(P_r[A_k = B_k] = J_w(A,B)\right)$$

the probability that $A$ and $B$ are equal on some position is like computing $J_w$.

proof:

Fix $k$ and let $S_A = \{ S_i^{(k)} : (i, \lambda_i) \in A \}$, $S_B = \{ S_i^{(k)} : (i, \lambda_i) \in B \}$

$S_i^{(k)} \sim Exp(\lambda_i)$



$A_k = \min\{S_A\}$    $B_k = \min\{S_B\}$

to have $A_k = B_k$ the minimum must be in the <u>intersection</u>

$$\Pr[A_{1k} = B_{1k}] = \Pr\left[\min\{S_A \cap S_B\} < \min\{(S_A \cup S_B) \setminus (S_A \cap S_B)\}\right] \underset{=}{\overset{\textcircled{8}}{}}$$

- $\min\{S_A \cap S_B\} \sim \text{Exp}(|A \cap B|_w)$    $\overset{\text{"}\lambda_3}{}$

- $\min\{(S_A \cup S_B) \setminus (S_A \cap S_B)\} \sim \text{Exp}(|(A \cup B) \setminus (A \cap B)|_w)$

also we know that

$$\mathcal{X} \sim \text{Exp}(x) \quad , \quad Y \sim \text{Exp}(y) \quad \Rightarrow \quad \Pr[\mathcal{X} < Y] = \frac{x}{x+y} \quad (\text{ex.}) \qquad \textcircled{**}$$

$$\underset{=}{} \frac{|A \cap B|_w}{|A \cap B|_w + |(A \cup B) \setminus (A \cap B)|_w} = \frac{|A \cap B|_w}{|A \cup B|_w} = J_w(A, B)$$

---

## Conclusion

- To estimate $J_w(A, B)$ it is enough to estimate $\Pr[A_k = B_k]$

- Let $\bar{J}_w(A, B) := \dfrac{\sum\limits_{k=1}^{m} \mathbb{1}_{A_n = B_k}}{m}$ , $\mathbb{1}_{A_n = B_k} = \begin{cases} 1 & \text{if } A_k = B_k \\ 0 & \text{if } A_k \neq B_k \end{cases}$

- $E[\bar{J}_w(A, B)] = \dfrac{\sum\limits_{k=1}^{m} E[\mathbb{1}_{A_k = B_n}]}{m} = \dfrac{\sum\limits_{k=1}^{m} \Pr[A_n = B_k]}{m} = \dfrac{m \cdot \Pr(A_k = B_k)}{m}$

$$= J_w(A, B)$$

- $\bar{I}(A, B) := \dfrac{m - 1}{\sum\limits_{k=1}^{m} \min\{A_k, B_k\}} \cdot \dfrac{\sum\limits_{k=1}^{m} \mathbb{1}_{A_k = B_n}}{m}$

<span style="color:magenta">↖  ↑<br>are independent</span><br>thus we get: $E[\bar{I}(A, B)] = |A \cap B|_w$

generally:

$$\cdot \ \overline{I}(A|B,C,\dots) := \underbrace{\dfrac{m-1}{\sum\limits_{K=1}^{m} \min\{A_n,B_n,C_n,\dots\}}}_{|A \cap B \cap C \cap \dots|_w} \cdot \underbrace{\dfrac{\sum\limits_{K=1}^{m} \mathbb{1}_{A_n = B_n = C_n = \dots}}{m}}_{|A \cup B \cup C \dots|_w} =$$



$$\boxed{\begin{array}{c} J_w(A,B,C,\dots) \\ \shortparallel \\[4pt] \dfrac{|A \cap B \cap C \cap \dots|_w}{|A \cup B \cup C \cup \dots|_w} \end{array}}$$

$$\cdot \ \mathbb{E}\big[\overline{I}(A,B,C\dots)\big] = |A \cap B \cap C \dots|_w \ \leftarrow \text{unbiased}$$

---

# COMPLEMENT

$$\cdot \ \widetilde{D}(A,B) := \underbrace{\dfrac{m-1}{\sum\limits_{K=1}^{m} \min\{A_n,B_n\}}}_{|A \cup B|_w} \cdot \underbrace{\dfrac{\sum\limits_{K=1}^{m} \mathbb{1}_{A_n < B_n}}{m}}_{\dfrac{|A \setminus B|_w}{|A \cup B|_w}}$$

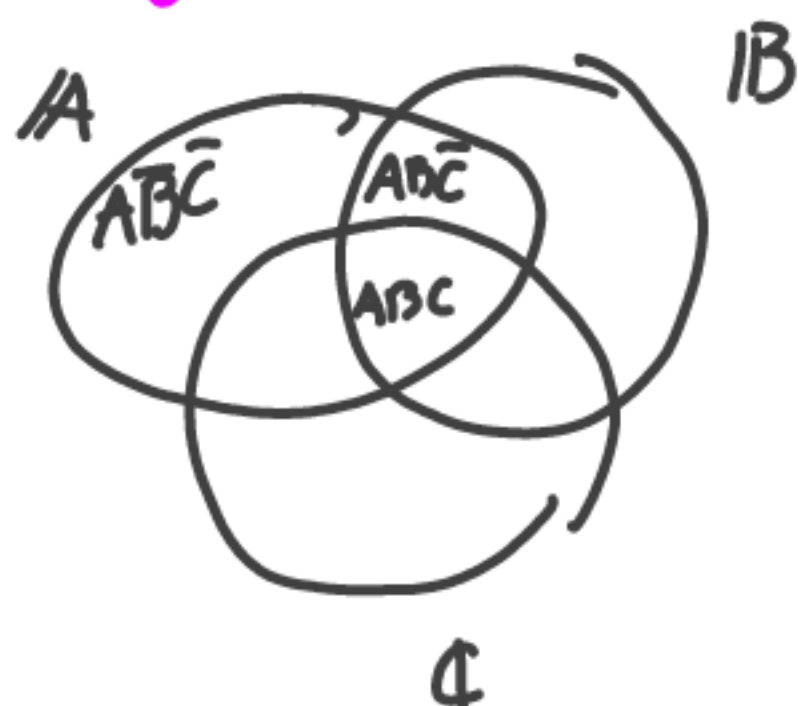$$\cdot \ \mathbb{E}\big[\widetilde{D}(A,B)\big] = |A \setminus B|_w$$

# Any sequence of operations

$A, B, C, \ldots$ can be simulated on sketches $A, B, C, \ldots$.

① Find disjunctive normal form (DNF),

e.g. $(A \setminus B) \cup (A \wedge B \wedge C) = A \cap \bar{B} \cap \bar{C} \cup A \cap B \cap \bar{C} \cup A \cap B \cap C$

② Estimate each conjuction seperetely and sum up all estimates.



③ Based on the results we have we can estimate each conjuction

e.g. $|A B C|_w \rightarrow$ we use $\bar{I}(A, B, C)$

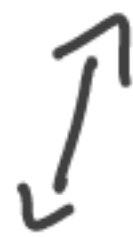$|A \bar{B} \bar{C}|_w = |A \setminus (B \cup C)|$

we need to count positions in the sketch such that

$$m' = |\{k : A_k < \min\{B_k, C_k\}\}|$$

$\bullet\ E\left[\dfrac{m-1}{\sum \min\{A_k, B_k, C_k\}} \cdot \dfrac{m'}{m}\right] = |A \bar{B} \bar{C}|_w$

$$|AB\bar{C}\,D\bar{E}\bar{F}| = |ABD\bar{C}\bar{E}\bar{F}| = |ABD \setminus (C \cup E \cup F)|$$

$$m' = |\{k : A_k = B_k = D_k < \min(C_k, E_k, F_k)\}|$$

in general:

$$E\left[\frac{m-1}{\sum\limits_{k=1}^{m} \min(A_k, B_k, C_k, D_k, E_k, F_k)} \cdot \frac{m'}{m}\right] = |AB\bar{C}D\bar{E}\bar{F}|_w$$