

Lecture 8

Sampling

$$X_y = \begin{cases} 1 & \text{if } y \in D \\ 0 & \text{oth} \end{cases}$$

$$E(X_y) = \frac{|D|}{|D|}$$

Sampling without replacement.

$$\int_{y,z} = \begin{cases} 1 & y=z \\ 0 & \text{oth} \end{cases}$$

Q: What does $y=z$ means?

1° $y=z$

numerically $y=z$

$y, z \in \mathbb{R}$

not necessarily.

2° $y=z$ if they refer to the same piece of D

How to assign identity to D

- pick any order
- index in order is id

val	1	1	2	3	4	5
id	1	2	3	4	5	6

$$E(X_y X_z)$$

$$\sum_{y \in D} \sum_{z \in D} \left(\int_{y,z} + \dots \right) > ?$$

From now on \mathcal{D} will have relational structure

$$\mathcal{D} = \{R_1, \dots, R_k\}$$

R_i is a relation with a tuple identity (primary key)

$$f(\mathcal{D}) = f(R_1, \dots, R_k)$$

A	B	ID
1	2	1
1	2	2
3	4	3

$$\int_{t, t'} = \begin{cases} 1 & \text{if } t \text{ is the same as } t' \\ 0 & \text{otherwise} \end{cases}$$

Previous example

$$\mathcal{D} = \{R\}$$

$$R(A)$$

$$f(\mathcal{D}) = \sum_{t \in R} t.A$$

$$f(\mathcal{D}) = \sum_{t \in \mathcal{D}} t$$

$$x = f(\mathcal{D}) = f(\mathcal{D}') = \sum_{t \in R'} t.A$$

$$x = \alpha \cdot \sum_{t \in R} t.A$$

$$x_t = \begin{cases} 1 & \text{if } t \in R' \\ 0 & \text{otherwise} \end{cases}$$

$$\mathcal{D}' = \{R'\} \quad R'(A)$$

R' is a sample of R sample without replacement

$$\alpha = \frac{|R|}{|R'|}$$

$$E(x) = \alpha \sum_{t \in R} E(x_t) \cdot t.A$$

$$E(x^2) = \alpha \sum_{t \in R} \sum_{t' \in R} E(x_t x_{t'}) \cdot t.A \cdot t'.A$$

Problem 2

$$D = \langle R(A) \rangle$$

$$\langle \cdot \rangle = \frac{|R|^{-1}}{|R|^{-1}}$$

$$f(D) = \sum_{t \in R} g(t, A)$$

where $g(x)$ is some fct.

$$\forall t \in R, X_t = \begin{cases} 1 & \text{if } t \in R' \\ 0 & \text{oth.} \end{cases}$$

$$E[X_t] = \frac{|R'|}{|R|} = \frac{1}{\alpha}$$

$$\begin{aligned} \text{if } t \in R \quad E[X_t X_{t'}] &= P(t \in R' \wedge t' \in R') \\ &= P(t \in R') \cdot P(t' \in R' | t \in R') \\ &= \frac{1}{\alpha} \begin{cases} \frac{1}{\alpha} & \text{if } t \neq t' \\ 1 & \text{oth.} \end{cases} \end{aligned}$$

$$\begin{aligned} E[X_t X_{t'}] &= \frac{1}{\alpha} \left(\sum_{t, t'} \dots + (1 - \sum_{t, t'} \dots) \frac{1}{\alpha} \right) \\ &= \frac{1}{\alpha} \left(\frac{1}{\alpha} + \sum_{t, t'} \dots \left(1 - \frac{1}{\alpha}\right) \right) \end{aligned}$$

X_t + prep. has nothing to do with $g(\cdot)$

$$D \quad X = \alpha \cdot \sum_{t \in R} g(t, A) = \alpha \sum_{t \in R} X_t g(t, A)$$

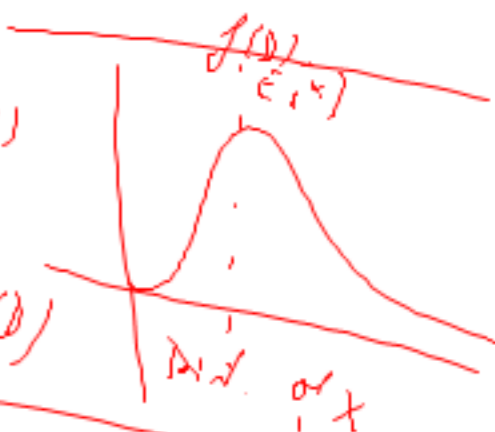
$$E[X] = \alpha \sum_{t \in R} E[X_t] \cdot g(t, A) = \sum_{t \in R} g(t, A) = f(D)$$

Example:

g encodes reliability indicator:

$$g(A) = \begin{cases} A & \text{if } A \geq 10 \\ 0 & \text{oth.} \end{cases}$$

SELECT SWs A FROM IP WHERE A >= 10



$$\begin{aligned} E[X^2] &= \alpha^2 \sum_{t \in R} \sum_{t' \in R} E[X_t X_{t'}] \cdot g(t, A) \cdot g(t', A) \\ &= \alpha^2 \frac{1}{\alpha} \left(\sum_{t \in R} g(t, A) \right)^2 - \sum_{t \in R} g(t, A)^2 \\ &= \alpha \frac{1}{\alpha} (1 - \frac{1}{\alpha}) \sum_{t \in R} g(t, A) \end{aligned}$$