Plan:

| | |
|---|---|
| 11/13 | important talk |
| 11/18 | go over for test. |
| 11/20 | Test 2 |
| 11/25 | Go over test 2 |
| 12/2 | } Presentations |
| 12/4 | } & projects |

## Privacy issues

Privacy for the data that is used

__Problem.__ We have data that we actually use, but we need to perserve privacy.

__We have.__
- statistical Data base : You are interested in statistics and no in individual records.
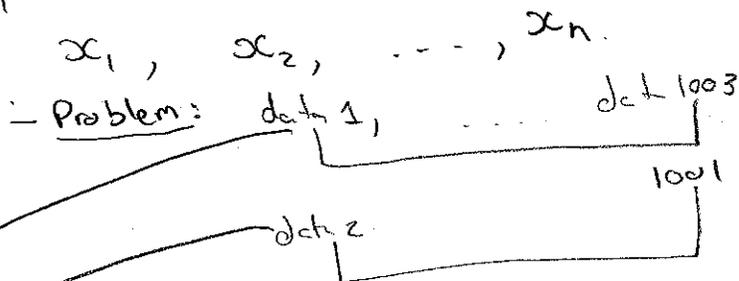- usual databases. objective get individual records

e.g.,
+ Census database :— good source for correlations
+ Medical database :—
(Need to get data without disclosing sensitive info.)

&ast; Many traditional ways to protect privacy.
+ limit size of the sample below.

$$x_1, \quad x_2, \quad \ldots, \quad x_n.$$

— __Problem:__ data 1, $\ldots$ data 1003

1001

$$\frac{S_1 + \ldots + S_K}{K} = \bar{S}_1$$

$$\frac{S_1 + \ldots + S_{K-1}}{K-1} = \bar{S}_2$$

data 2

$$\bar{S}_1 \cdot K = S_1 + \ldots + S_{K-1} + S_K$$

$$\bar{S}_2(K-1) = S_1 + \ldots + S_{K-1}$$

$$\bar{S}_1 \cdot K - \bar{S}_2(K-1) = \ldots S_K$$

|  |  |  |  |  |
|---|---|---|---|---|
| 10 | 100 | 80 | 50 | $\longrightarrow 60 = \bar{S_1}$ |
| $S_1$ | $S_2$ | $S_3$ | $S_4$ | $\longleftarrow \dfrac{S_1 + S_2 + S_3}{3}$ |

$$\frac{10 + 100 + 80 + 50}{4} = 60$$

$$\frac{10 + 100 + 80}{3} = \frac{190}{3} = 63.335$$

$$S_1 + S_2 + S_3 + S_4 = \bar{S_1} \cdot 4 = 240$$
$$S_1 + S_2 + S_3 = \bar{S_2} \cdot 3 = 190$$

Another trick.

you
always
get the

Same
random #.

$$\begin{cases} S_. + \Delta S_1 = S_1 \\ S_. + \Delta S_2 = S_2 \\ \quad \cdots \\ S_. + \Delta S_n = S_n \end{cases}$$

Radically different approach.

- Don't keep exact data
- Keep only bounds.  [20, 30], [30, 40], ...,
    Salary : [10, 30], [20, 100]    [100, 300]

+ Privacy is preserved.
− How to compute statistical characteristics.

| If we know $S_i$ | We only know |
|---|---|

$$S = \frac{S_1 + \ldots + S_n}{n}$$

$$S_i \in [\underline{S_i}, \overline{S_i}]$$

$$V = \frac{1}{n} \sum_{i}^{n} (S_i - S_{av})^2$$

$$\underline{S_{av}} = \frac{\underline{S_1} + \ldots + \underline{S_n}}{n}$$

$$= \frac{1}{n} \sum_{i=1}^{n} S_i^2 - \left( \frac{1}{n} \sum_{i=1}^{n} S_i \right)^2$$

$$\overline{S_{av}} = \frac{\overline{S_1} + \ldots + \overline{S_n}}{n}$$

In this case the fcn is monotonic.

$$y = f(x_1, \ldots, x_n)$$

Previously: $[x_i - \Delta_i, x_i + \Delta_i]$

We know: $x_i \in [\underline{x_i}, \overline{x_i}]$

Want range $\{ y = f(x_1, \ldots, x_n) : x_i \in [\underline{x_i}, \overline{x_i}] \}$



$\underline{x_i}$   $\overline{x_i}$

$\widetilde{x_i} - \Delta_i$   $\widetilde{x_i} + \Delta$

$$\Delta_i = \frac{\overline{x_i} - \underline{x_i}}{2}$$ half width.

Use Monte Carlo when $f$ is easy.

# Monte Carlo Method

☑ fewer calls to $f$

— more computations $\sum_{i=1}^{n} |f(\tilde{x}_1, \ldots, \tilde{x}_{i-1}, \tilde{x}_i + \Delta_i, \tilde{x}_{i+1}, \ldots, \tilde{x}_n) - \tilde{y}|$

⚠ Question

    — Here is the characteristic
    — Here is the Input
    — Find the values.

| interval | mid point | half width |
|---|---|---|
| [10, 20] | 15 | 5 |
| [20, 30] | 25 | 5 |
| [60, 70] | 65 | 5 |

$S_1$      $S_2$