

Theoretical Basics of Big Data Analytics and Real Time Computation Algorithms *Spring 2016*

- What is **Big Data**?
 - Why is it **new and important**?
 - What **tools** do we need to deal with it?
- Big collections of data:
 - Retail store inventory and transactions
 - Cell company records
 - Weather records

Databases: store and manage data

- **Digitized data** - routinely collected

- Shopping transactions
- Search queries
- Internet traffic
- Readings from sensors

Often tossed away or stored and never used

- **Recently: something special**

- New kinds of information
- Was not anticipated
- Emerges when BIG

- **Valuable **hidden** information**

- not visible
- has to be extracted,
- processed

- **Big Data** is about:

- New information
- Specific tools

Such definition is

- is too vague
- relies on the notion of *information*

Information?

“Logic and Information” *Keith Devlin*

- What is *Iron*?
 - Go to Iron age
 - Ask Ironsmith
 - Examples: raw, processing, things
- Result: unsatisfactory
 - No frame of reference
 - Need to know: molecular structure...
- Now we are in **Information Age**
 - Know it exists
 - Have examples
 - Definitions in special cases
- Entering **Big Data (sub) Age**

Examples of Big Data

“BIG DATA: A Revolution That Will Transform How We Live, Work, and Think”

*Viktor Mayer-Schönberger
and Kenneth Cukier*

- **Target - Detecting:** a woman is pregnant
 - Two dozen products used as proxies
 - Estimate pregnancy stages, due date
 - Send relevant coupons
- **Correlation-based techniques**
 - Predict mechanical failures
 - Things break down *gradually*
 - Sensors + correlation analysis:
 - * Whirling motor
 - * Excessive heat
 - * ...

• UPS

- Replacing parts: 2-3years
 - * Inefficient
 - Predictive analysis
 - * Monitoring individual parts
 - Predictions made automatically
 - Based on:
 - * Great number of cases
 - * Correlation analysis
 - * No complex models
- Modern cars
- Lots of sensors
 - * Temperature, Vibrations, Voltages...
 - Use *complex* models of prediction
 - Information tossed away - no learning
 - Imagine: transmitted, collected, and analyzed...

• H1N1 Virus 2009

- Only hope: to slow its spread
- Need to know: where it is
- US Center for Diseases Control (CDC):
 - * Doctors: to inform of new flu cases
 - * Week or two out of date
 - * \Rightarrow Delays blinded health agencies

Few weeks before H1B1:

- * Google: paper in “Nature”
- * Predict spread if winter flu by looking what people were *searching*
- * 3 Bil. queries a day

- **Google “learning” technique**
 - CDC data for 2003-2008
 - Correlations:
 - * Search queries (50M most common)
 - * Flu spread
 - Result:
 - * Combination of 45 search terms
 - * + Math model
 - * = Strong correlation with official figures
 - So, in 2009 Google - more timely indicator
- No need in:
- * mouth swabs
 - * contacting doctors...
- Instead: huge amount of data
- * Too Big
 - * Too Noisy

- **Unexpected data in existing collections:**
 - Too Big
 - Too Noisy
- **Arranging new studies**

Aspirin and orange juice vs Cancer

 - **Standard way:**
 - * Specific tests
 - * Time
 - * Low confidence (small amount of data)
 - **Big Data way:**
 - * Digitized med. records
 - * Shopping transactions
 - * Search queries
 - * ... (lots of other data)

Other Big Data challenges

- **LHC Large Hadron Collider**
 - **150 Mil.** sensors
 - **40 Mil.** times per second
 - Only **0.001%** saved
 - **25 PB** in 2012 (1PB=1000TB)
 - If all recorded:
 - * **150 Mil. PB/year**
 - * = **200** x all other sources in the world
- **Modern Aircraft**
 - **100,000** sensors
 - Only **3 GB** in an hour flight

Seems not Big, but

- Monitoring in **real time**
- **Combinations** of readings
- In **dynamics**
- Need to make very fast **predictions**

⇒ Big Data challenge

- **Digitized Media Streaming**

- Large volumes
- But: Nothing is hidden

⇒ Not considered as Big Data

- **Information in Big Data**

- Hidden
- Requires special tools

Analogy:

- Rare mineral
- Nuclear fusion energy

Big Data Manipulations: Basic Steps

- **Extract** pieces of information (probably from distributed sources)
- **Unify** - transform to “canonical” form
 - Compact
 - Easy to handle
 - Contains sufficient information
- **Combine** pieces
- **Update** when new info arrives
- **Utilize** - Decision making

Simplest Example

x - object of interest (unknown value)

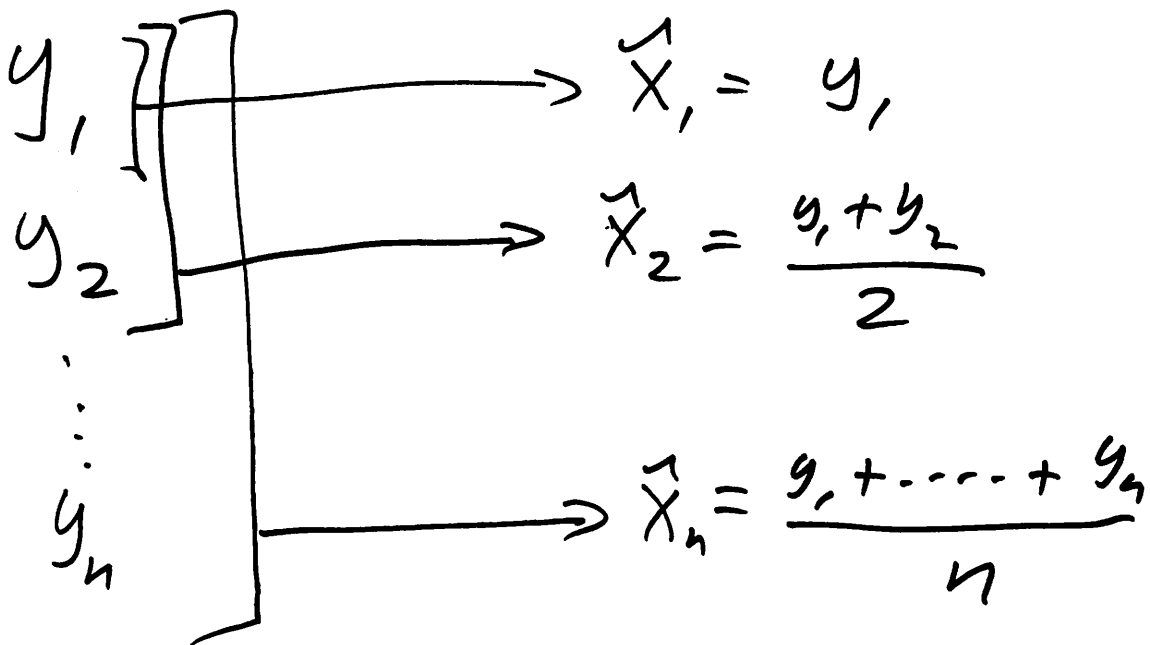
Observations:

$$y_i = x + \varepsilon_i, \quad i = 1, \dots, n$$

ε_i - i.i.d. random values, $E\varepsilon_i = 0$.

A good estimate of x :

$$\hat{x} = \frac{1}{n} \sum_{i=1}^n y_i$$



Updating \hat{x}

n : have \hat{x}_n , get y_{n+1}

$$\hat{x}_{n+1} = \frac{1}{n+1} \sum_{i=1}^{n+1} y_i = \frac{1}{n+1} \left(\underbrace{\sum_{i=1}^n y_i}_{= n \hat{x}_n} + y_{n+1} \right)$$

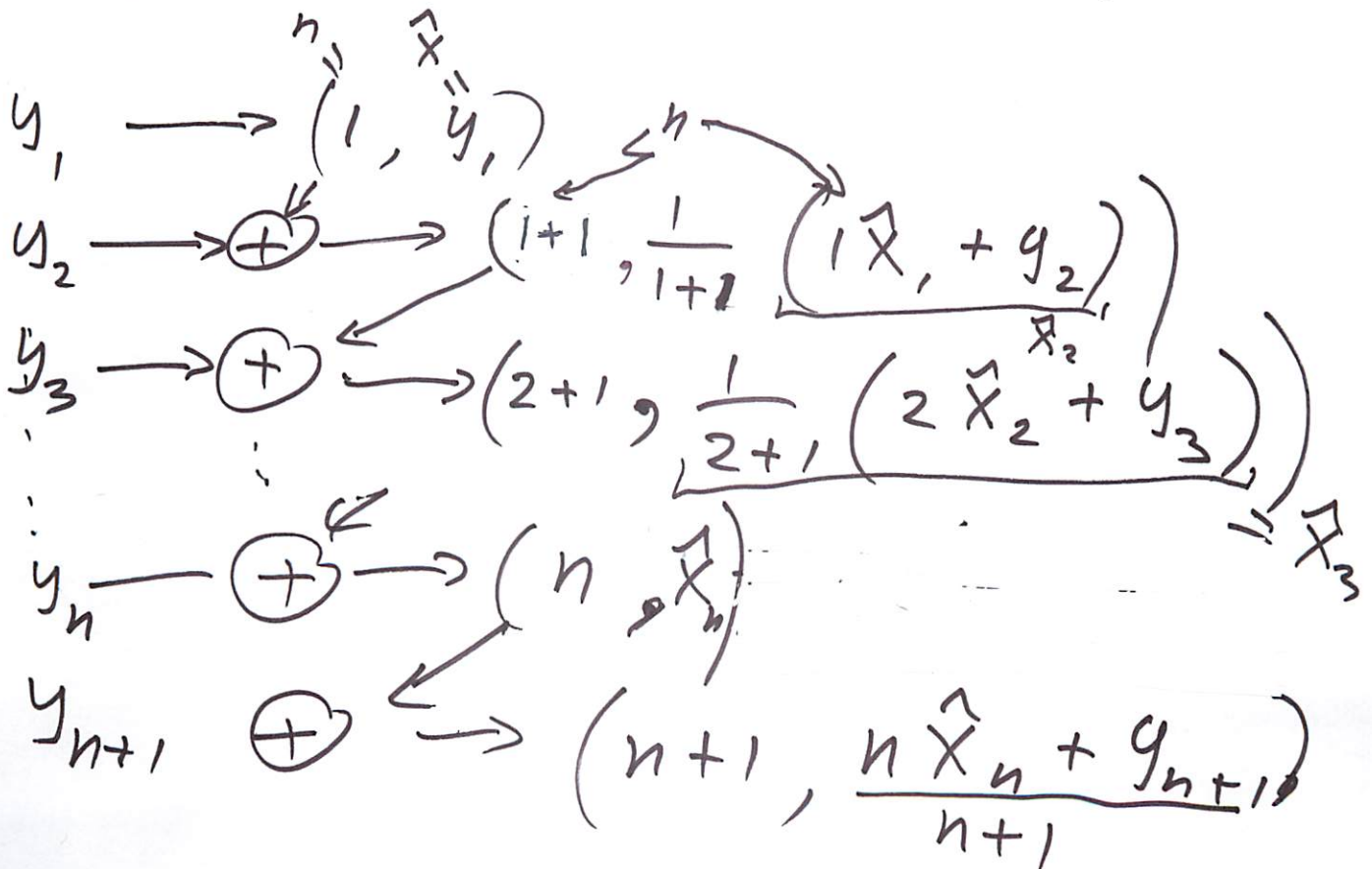
$$= \frac{1}{n+1} (n \hat{x}_n + y_{n+1})$$

or

$$\hat{x}_{n+1} = \hat{x}_n + \frac{1}{n+1} (y_{n+1} - \hat{x}_n)$$

In addition to (\hat{x}_n) need to keep n .

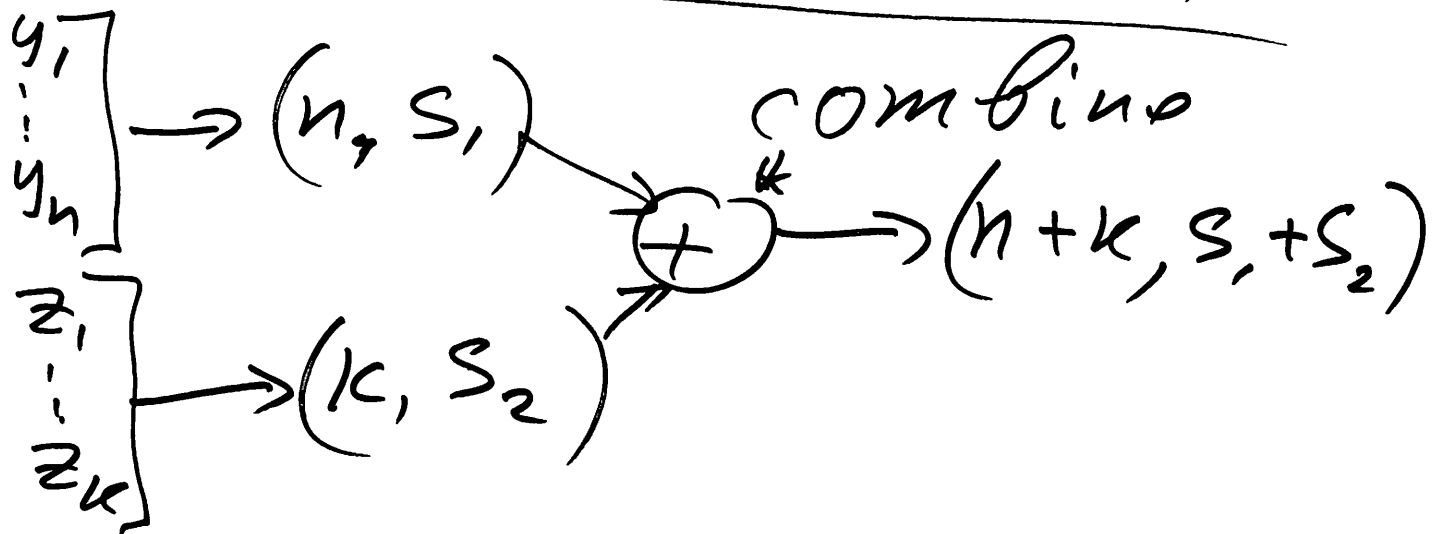
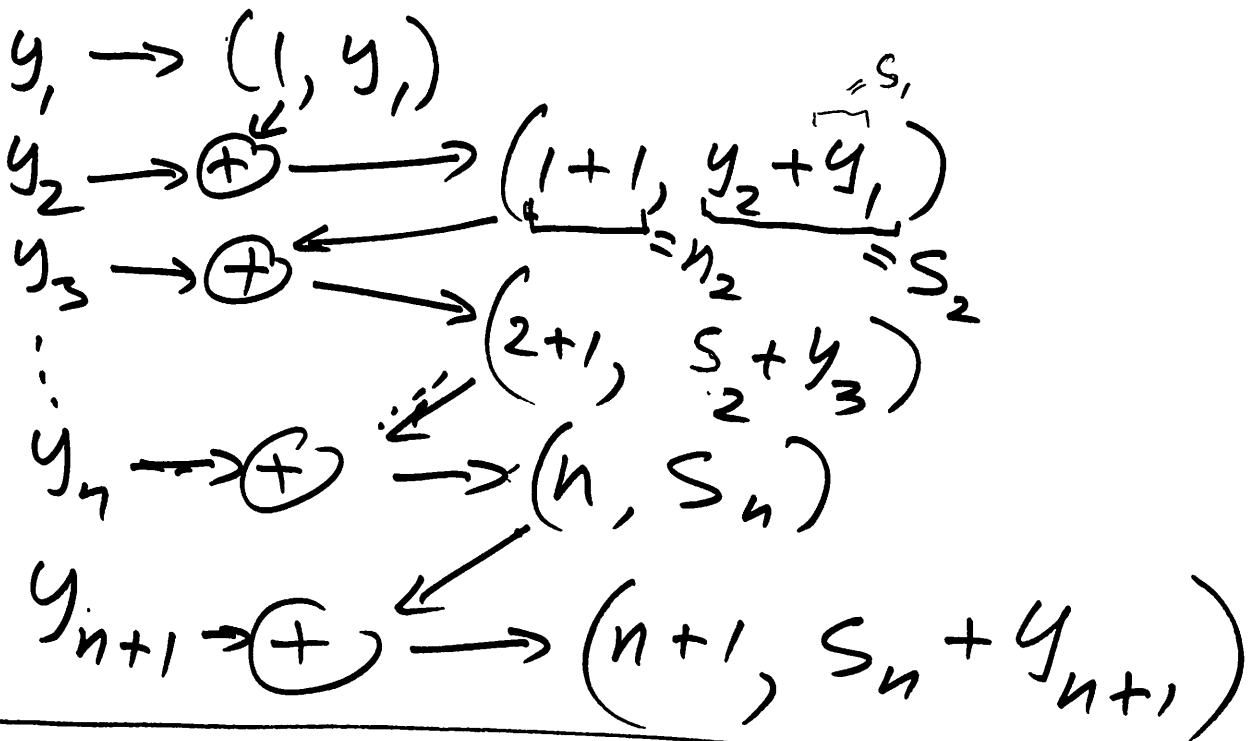
"Explicit" form of information: (n, \hat{x}_n) .



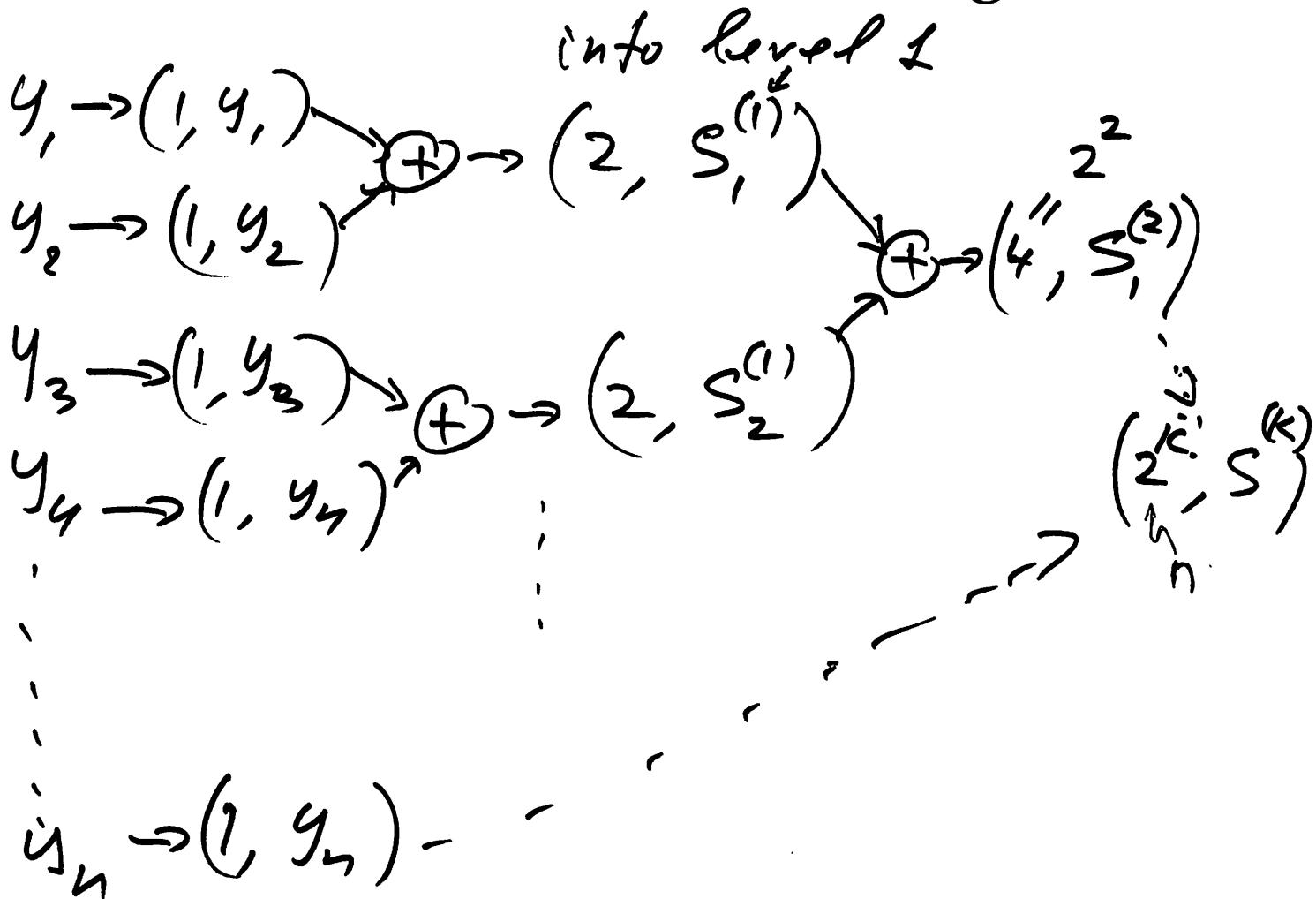
Updating "Canonical" Information (n, S)

n - number of readings, $S = \sum_{i=1}^n y_i$

$$(n, S) \Rightarrow \hat{X} = S/n$$



Concurrent Combining



$$n = 2^k \quad k = \log_2 n$$

if $n = 1000$

$k = 10$

$n = 1M$

$k = 20$

$n = 1B$

$k = 30$

Precision of \hat{x}

$$y_i = x + \varepsilon_i$$

$$\varepsilon_i - \text{i.i.d.}, \quad E\varepsilon_i = 0$$

$$\text{Var}(\varepsilon_i) = E(\varepsilon_i - E\varepsilon_i)^2 = E\varepsilon_i^2 = \underline{\underline{\sigma^2}}$$

$$\hat{x} = \frac{1}{n} \sum_{i=1}^n y_i - \text{unbiased est. of } x, \text{ i.e. } E\hat{x} = x$$

$$E\hat{x} = E\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n} E \sum_{i=1}^n (x + \varepsilon_i)$$

$$= \frac{1}{n} \sum_{i=1}^n x + \frac{1}{n} E \sum_{i=1}^n \varepsilon_i = x$$

$$\underbrace{\hspace{10em}}_{\substack{= \sum_{i=1}^n E\varepsilon_i \\ \downarrow \\ = 0}}$$

$$\text{Var}(\hat{x}) = E(\hat{x} - x)^2 = E\left(\frac{1}{n} \sum_{i=1}^n (x + \varepsilon_i) - x\right)^2$$

$$= E\frac{1}{n^2} \left(\sum_{i=1}^n (\cancel{x} + \varepsilon_i - \cancel{x})\right)^2 = \frac{1}{n^2} \sum_{i,j=1}^n E\varepsilon_i \varepsilon_j$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \underline{\underline{\frac{\sigma^2}{n}}}$$

$$\left(\sum \varepsilon_i\right)^2 = \sum_i \varepsilon_i \cdot \sum_j \varepsilon_j$$

ind $\Rightarrow E\varepsilon_i \varepsilon_j = \underbrace{E\varepsilon_i}_{=0} \cdot \underbrace{E\varepsilon_j}_{=0} = 0$ if $i \neq j$

$$\text{Var}(\hat{x}) = \frac{\sigma^2}{n} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

If we collect canonical info (n, S) :

$$(n, S) \Rightarrow \hat{x} = \frac{S}{n}, \quad \text{Var}(\hat{x}) = \frac{\sigma^2}{n}.$$

(n, S) is sufficient to obtain \hat{x} and its variance, but only when σ^2 is known.

*Suppose σ^2 is **not** known*

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{x})^2$$

- Unbiased estimate of σ^2 .

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{x})^2 &= \overbrace{\sum_{i=1}^n y_i^2}^T - 2 \overbrace{\sum_{i=1}^n y_i \cdot \hat{x}}^{=S} + n\hat{x}^2 && \hat{x} = \frac{S}{n} \\ &= T - 2S \frac{S}{n} + n \left(\frac{S}{n} \right)^2 = T - \frac{S^2}{n} \end{aligned}$$

$$T = \sum_{i=1}^n y_i^2$$

New canonical information (n, S, T) :

$$n = \sum_{i=1}^n y_i^0, \quad S = \sum_{i=1}^n y_i^1, \quad T = \sum_{i=1}^n y_i^2$$

$$\widehat{\sigma^2} = \frac{1}{n-1} \left(T - \frac{S^2}{n} \right)$$

Estimate of the variance of \widehat{x} :

$$V = \widehat{\text{Var}(\widehat{x})} = \frac{\widehat{\sigma^2}}{n} = \frac{1}{n(n-1)} \left(T - \frac{S^2}{n} \right)$$

$$(n, S, T) \Rightarrow \widehat{x} = \frac{S}{n}, \quad V = \frac{1}{n(n-1)} \left(T - \frac{S^2}{n} \right)$$

Updating can. info:

$$(n, S, T) \xrightarrow{y} \oplus \longrightarrow (n+1, S+y, T+y^2)$$

Combining can. info:

$$\begin{array}{l} (n_1, S_1, T_1) \\ (n_2, S_2, T_2) \end{array} \xrightarrow{\oplus} \longrightarrow (n_1+n_2, S_1+S_2, T_1+T_2)$$

Info in **explicit form** $(n, \underline{\hat{x}}, \underline{V})$:
 Have (n, \hat{x}_n, V_n) , receive y_{n+1}

$$\hat{x}_{n+1} = \hat{x}_n + \frac{y_{n+1} - \hat{x}_n}{n+1}$$

$$\hat{\sigma}_{n+1}^2 = \frac{n-1}{n} \hat{\sigma}_n^2 + \frac{(y_{n+1} - \hat{x}_n)^2}{n+1}$$

(n) = V_n

$$V_{n+1} = \frac{\hat{\sigma}_{n+1}^2}{n+1} = \frac{n-1}{n+1} V_n + \left(\frac{y_{n+1} - \hat{x}_n}{n+1} \right)^2$$

Updating explicit info:

$$(n, \hat{x}, V) \xrightarrow{y} \left(n+1, \hat{x}_n + \frac{y_{n+1} - \hat{x}_n}{n+1}, V \right)$$

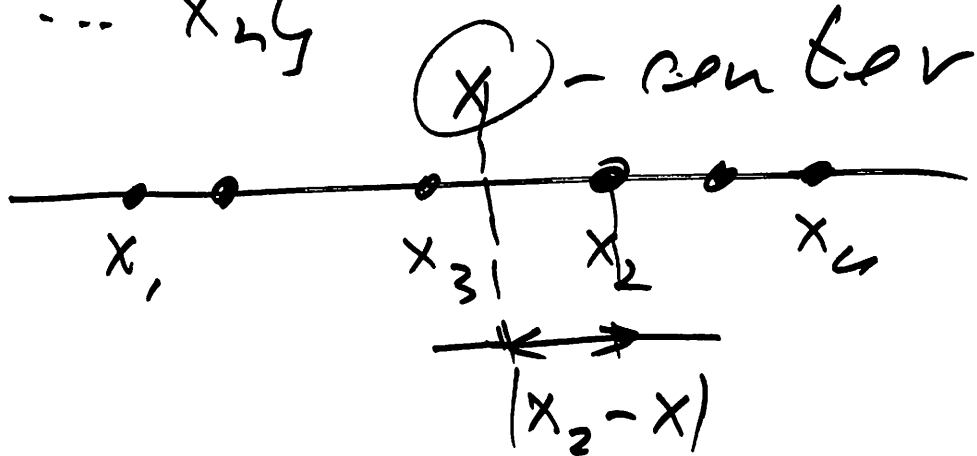
Explicit form for single observation:

$$y \rightarrow (1, y, n?)$$

$$V = \frac{1}{n(n-1)} \sum_{i=1}^n (y_i - \bar{x})^2$$

\Rightarrow Information in explicit form *may not exist*

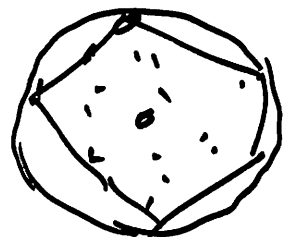
Center of
a set of points
 $\{x_1, \dots, x_n\}$



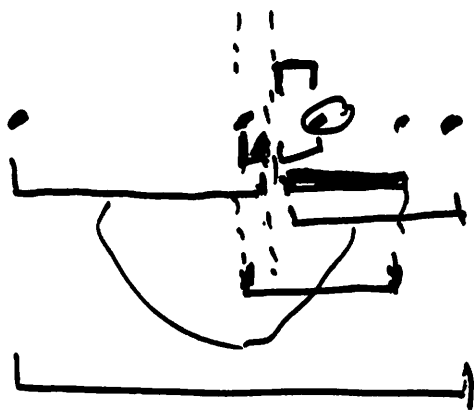
$$a) \sum_{i=1}^n |x - x_i| \sim \min_x$$

$$b) \sum_{i=1}^n (x - x_i)^2 \sim \min_x$$

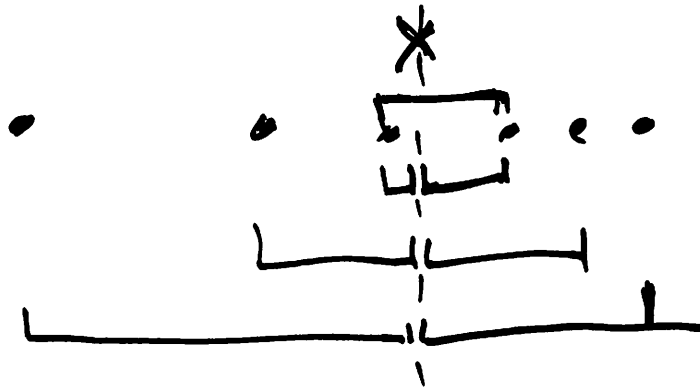
$$c) \max_{i=1, \dots, n} |x - x_i| \sim \min_x$$



a.)

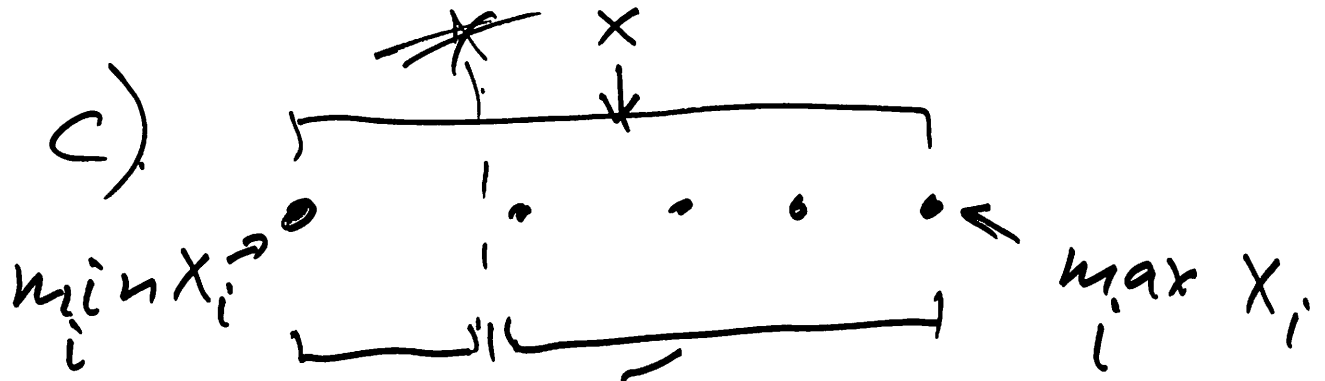


median
(odd)



even
not unique

$$\begin{aligned}
 \text{b)} \quad \sum_{i=1}^n (x - x_i)^2 &= \sum_{i=1}^n \left(\underbrace{(x - \bar{x}) + (\bar{x} - x_i)} \right)^2 = \\
 &= \underbrace{n(x - \bar{x})^2} + 2 \sum_{i=1}^n \underbrace{(x - \bar{x})(\bar{x} - x_i)}_{=0} \\
 &\quad + \underbrace{\sum_{i=1}^n (\bar{x} - x_i)^2}_{= \text{const}} \\
 &= \underbrace{n(x - \bar{x})^2}_{\text{min if } x = \bar{x}} + \sum_{i=1}^n (\bar{x} - x_i)^2 \\
 \Rightarrow \text{center is } \underline{\underline{\text{mean}}}
 \end{aligned}$$



$$\min_i \max |x - x_i|$$

x

$$x = \frac{1}{2} (\min_i x_i + \max_i x_i)$$

call it "middle"