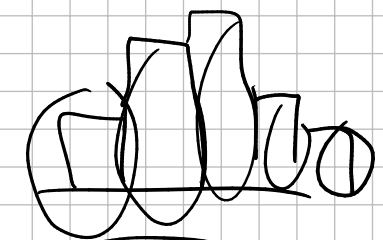


histogram, box plot



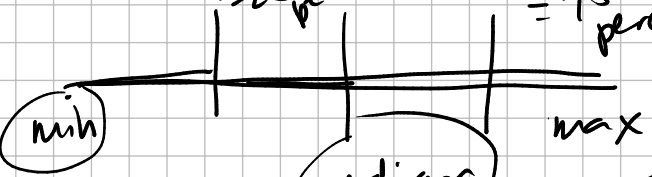
data plenty of observations

numerical summaries

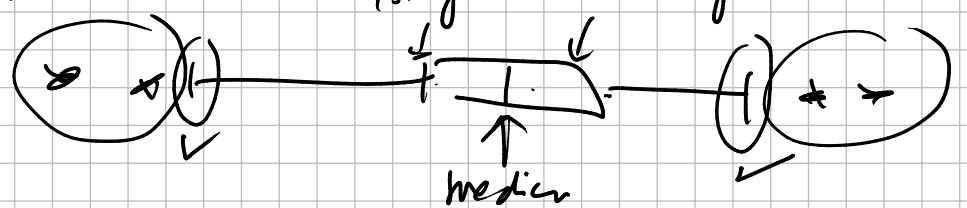
reduce dimensions

1st quartile = 25th percentile
3rd quartile = 75th percentile

lower limit
upper limit
count



1, 2, 3, 4, 5
1, 2, 3, 4



1, 2, 3, 4, 5
1000

2.5

continuous variables

histograms / box plot / Some of numerical summaries

might not be useful for categorical variables } discrete variables
dummy variables

mean (average)

$$\frac{1}{n} \sum_{t=1}^n X_t$$

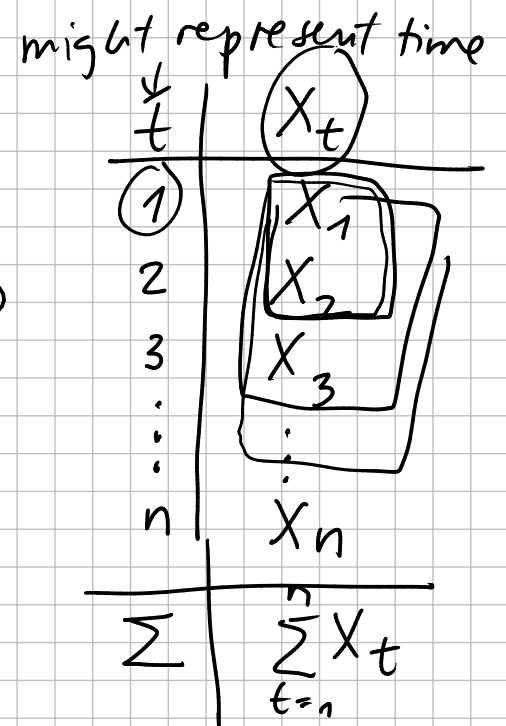
X_t variable indexed by t
subscript

PERNUM
PERWT

n ← ending point
 \sum = summation
 $t=1$ → starting point

ending point

$$\sum_{t=1}^{(n)} X_t = X_{(t=1)} + X_{(t=2)} + X_{(t=3)} + \dots + X_{(t=n)}$$



$$\text{mean} = \bar{X} = \frac{1}{n} \sum_{t=1}^{(n)} X_t = \text{simple average}$$

$$\sum_{t=3}^{(n-2)} X_t = X_3 + X_4 + \dots + X_{n-2}$$

moving average

trimmed mean

Survey:

$$\sum_{t=3}^6 (-1)^t X_t = (-1)^3 X_3 + (-1)^4 X_4 + (-1)^5 X_5 + (-1)^6 X_6$$

$$= (-1)^3 \cdot 3 + (-1)^4 \cdot 4 + (-1)^5 \cdot 5 + (-1)^6 \cdot 6 = 2$$

$$X_t = t$$

$$\sum_{t=3}^6 X_t = X_3 + X_4 + X_5 + X_6 = 2 + 2 + 2 + 2 = 8$$

Suppose $X_t = c$, where c is some constant, for all t

$$\sum_{t=3}^6 X_t = 4 \cdot c \quad \text{if } X_t = c \text{ for all } c$$

$$\begin{aligned} \sum_{t=3}^6 X_t &= \sum_{t=3}^6 [2] = 2 \sum_{t=3}^6 (1) \\ &= 2(1+1+1+1) \\ &= 8 \end{aligned}$$

$$\sum_{t=3}^6 X_t = X_3 + X_4 + X_5 + X_6 = 2y_3 + 2y_4 + 2y_5 + 2y_6$$

Suppose $X_t = 2y_t$

$$\begin{aligned} &= 2(y_3 + y_4 + y_5 + y_6) \\ &= 2 \sum_{t=3}^6 y_t \end{aligned}$$

In other words,

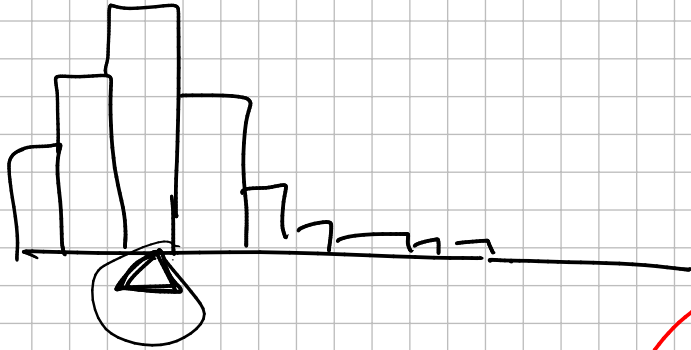
$$\sum_{t=3}^6 2y_t = 2 \sum_{t=3}^6 y_t$$

$$\sum_{t=3}^6 (X_t)^2 = X_3^2 + X_4^2 + X_5^2 + X_6^2$$

NOTE: $\left(\sum_{t=3}^6 X_t \right)^2 = \left[X_3 + X_4 + X_5 + X_6 \right]^2$

generally,
not
equal!

Where in the histogram can you find the mean?

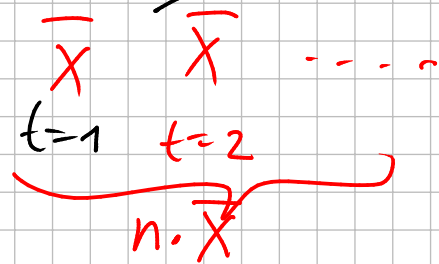


$$\bar{X} = \frac{1}{n} \sum_{t=1}^n X_t$$

$$0 = \left(\frac{1}{n} \sum_{t=1}^n X_t \right) - \bar{X}$$

any $X_t!$ $\Rightarrow \left(\frac{1}{n} \sum_{t=1}^n X_t \right) - \left(\frac{1}{n} \sum_{t=1}^n \bar{X} \right)$

$$0 = \frac{1}{n} \sum_{t=1}^n X_t - \frac{1}{n} \sum_{t=1}^n \bar{X} = \frac{1}{n} \sum_{t=1}^n (X_t - \bar{X})$$



average of $(X_t - \bar{X})$'s = 0
deviation from the mean

how every observing X differs from \bar{X}

negative deviations

$$X_t < \bar{X}$$

$$X_t > \bar{X}$$

positive deviations

sd() std dev of n variable numerical

X Name	X Company	X Industry	...

sd(execComp \$ TotalCompMil) -
sd(execComp \$ SalaryThou) -
apply

std dev = $\sqrt{\frac{1}{n} \sum_{t=1}^n (X_t - \bar{X})^2}$ or sometimes $\sqrt{\frac{1}{n-1} \sum_{t=1}^n (X_t - \bar{X})^2}$

$(X_t - \bar{X})^2$ $(X_t - \bar{X})^2 (0.01)^2$

It does not make sense to propose a measure of spread

based on $\frac{1}{n} \sum_{t=1}^n (X_t - \bar{X})$!
always 0.

alternative measure of spread

absolute deviation

$$\frac{1}{n} \sum |X_t - \bar{X}|$$

$|1-3| = 3$
 $|3| = 3$

$$\text{variance} = \frac{1}{n} \sum (X_t - \bar{X})^2$$

changes the measurement of (X_t)
(units)

X_t measured in dollars

$(X_t - \bar{X})^2$ measured in dollars²

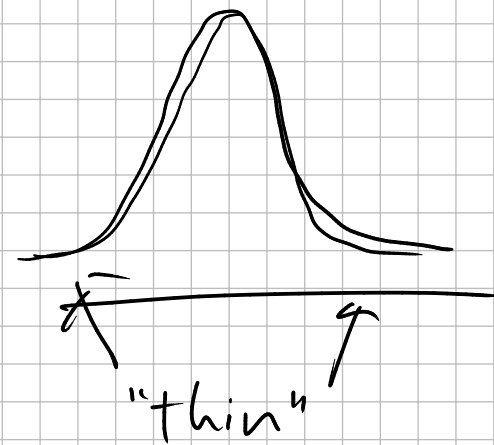
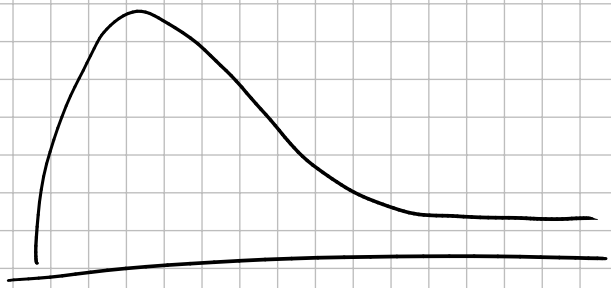
mean (exec Comp & Total Comp Mi) → see what happens!

461300	461.3
- 696798	- 696.7
-----	-----
255498 dollars	255.4 thousand dollars

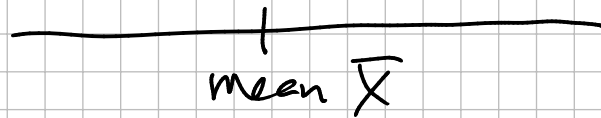
$$\frac{X_t - \bar{X}}{sd} = \text{number} \quad \text{no unit/unitless}$$

$$\pm 3$$

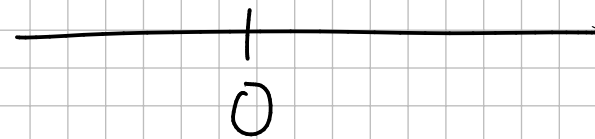
$$X_t - \bar{X} = \text{number} * sd \Rightarrow X_t = \bar{X} + \text{number} * sd$$



$$\frac{X_t - \bar{X}}{sd}$$



original scale



standardization

z-score

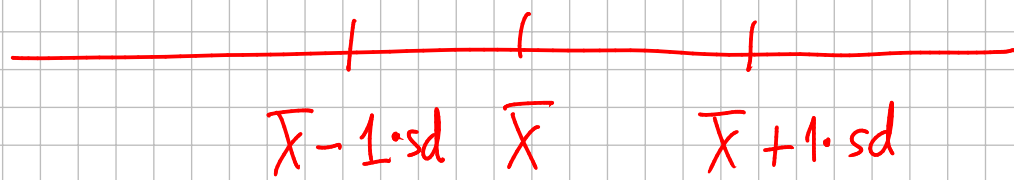
mean of $\left(\frac{X_t - \bar{X}}{sd} \right)$

mean $\left(\text{abs}(z_{Sal}) > 1 \right)$

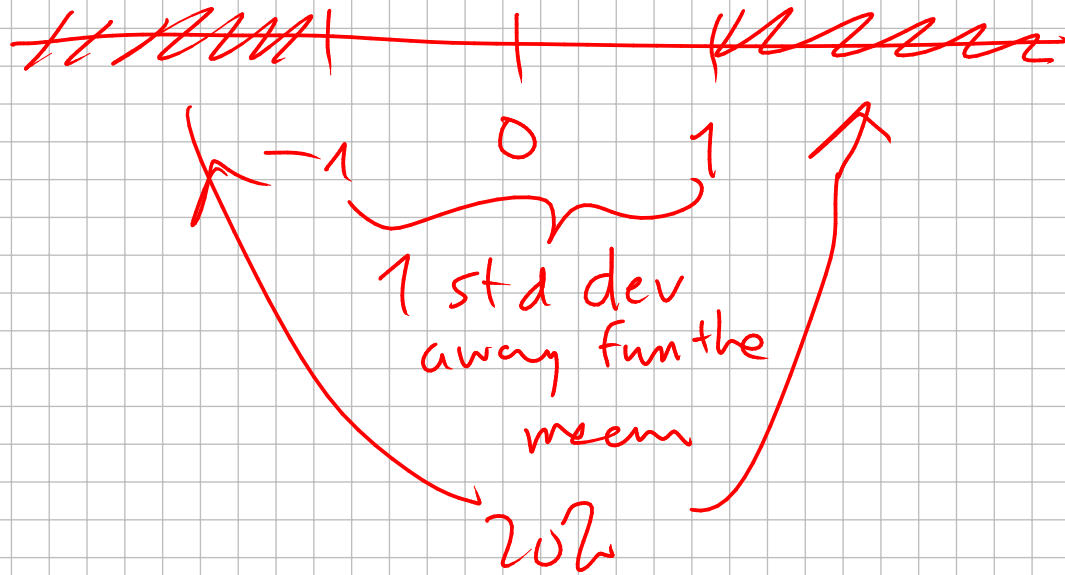
ISO 1 z-scores

? proportion

$\text{abs}(z_{Sal}) > 1$	z-scores	0.5	0.6	-1.7	1.8	...
	1.2					
	TRUE	FALSE	FALSE	TRUE	TRUE	



original scale
dollars



standardized scale
no units

$6.29e+02$ in R

$6.29 * 10^2$ scientific notation
 $= 629$

$$X_t \rightarrow X_t - \bar{X}$$

translation

$$X_t \rightarrow \frac{X_t - \bar{X}}{sd}$$

combination of translation & scaling

$$F = \left[\frac{9}{5} \right] C + 32$$

$X_t \rightarrow \log X_t$ natural logarithmic transformation

weird for economics algebra

$$\log X_t = \log_e X_t$$
$$\log X_t = \log_{10} X_t$$

strange, you note!

$$X_t \rightarrow X_t^2$$

$$X_t \rightarrow \frac{(X_t - \bar{X})^2}{(X_t - \bar{X})(X_t - \bar{X})}$$

nonlinear transformations

$$\log_b X = ?$$

↳ exponent

$$\log_b X = y \Leftrightarrow$$

$$b^y = X$$

$$\log_b x_2 + \log_b x_1 = \log_b (x_1 \cdot x_2)$$

$$\log_b X_2 - \log_b X_1 = \log_b \left(\frac{X_2}{X_1} \right)$$

$$\log_b X^a = a \log_b X$$

$$\log_{\underbrace{100}_{10^2}} 50 = y$$

$$\& \log_{10} 5 = 0.7$$

$$\Rightarrow 10^{0.7} = 5$$

$$100^y = 50 \Rightarrow (10^2)^y = 50 \Rightarrow 10^{2y} = 50$$

$$\Rightarrow \log_{10} 10^{2y} = \log_{10} 50$$

$$2y \log_{10} 10 = \log_{10} 50$$

$$2y = \log_{10} 50$$

$$y = \frac{\log_{10} 50}{2}$$

$$= \frac{1.7}{2} = 0.85$$

$\log 0$ does not exist
in \mathbb{R} - INF

$$10^? = 10$$

$$\log_{10} 50 = \log_{10} (5 \cdot 10)$$

$$= \underbrace{\log_{10} 5}_{0.7} + \underbrace{\log_{10} 10}_1$$

NOTE

$$\log_b (X_1 + X_2) = ? \neq \log_b X_1 + \log_b X_2$$

original scale
data

transformations



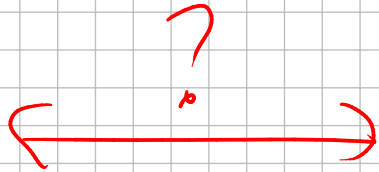
reduce ↓

visualization

numerical summaries

data

↓ reduce



$$\textcircled{10 \wedge 6} = 10^6$$

1e06 * 10⁶

(scientific notation)

a^x

a^x