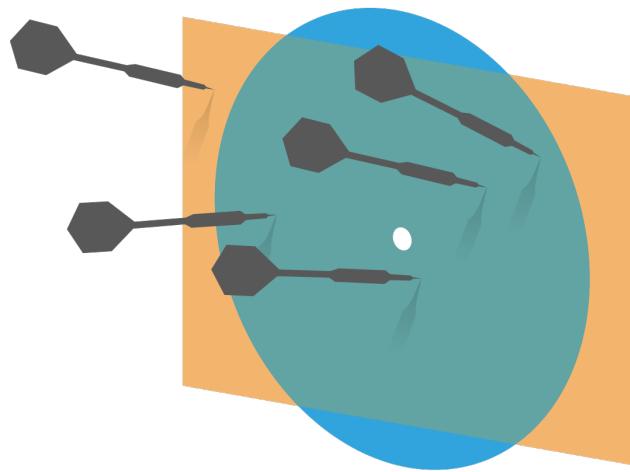


Probability and Statistics for Computer Science



“Correlation is not Causation”
but Correlation is so beautiful!

Credit: wikipedia

Last time

Mean

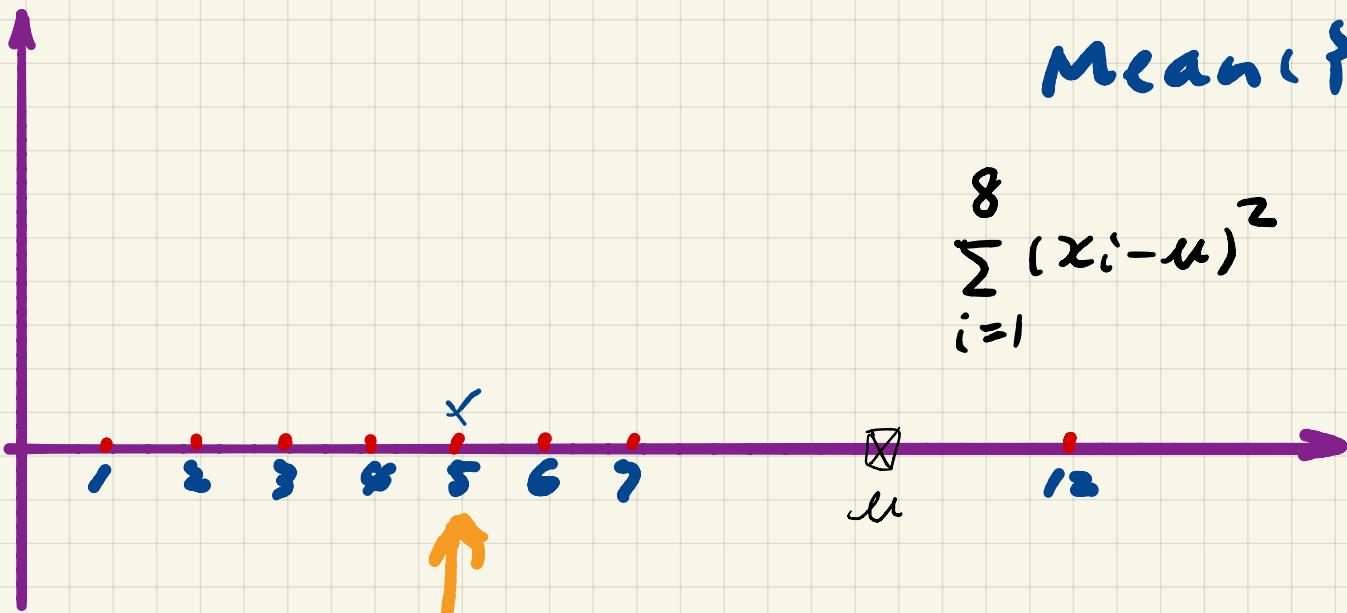
Standard deviation

Variance

Standardizing data

$$\{x_i\} \quad i \in [1, 8]$$

$$\{x_i\} = 1, 2, 3, 4, 5, 6, 7, 12$$

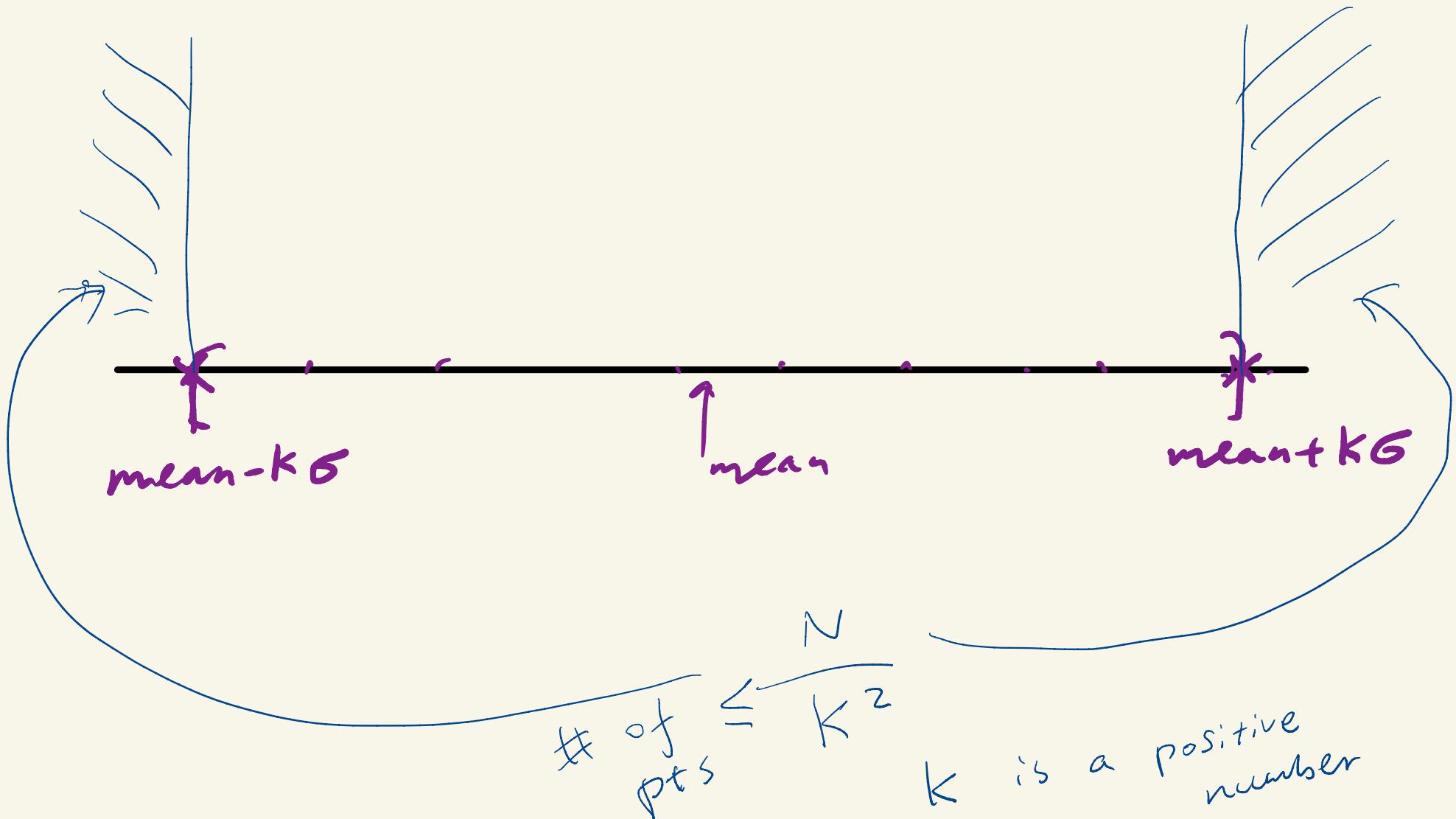


$$\text{Mean}(\{x_i\}) = 5$$

$$\sum_{i=1}^8 (x_i - \mu)^2$$

$$\underset{\mu}{\operatorname{argmin}} \left(\sum_i (x_i - \mu)^2 \right)$$

$$\{x_i\} \quad i \in [1, N]$$



Q: Estimate the range of data in standard coordinates

✳ The interval [-5, 5] covers x% data,
choose the closest estimate of x.

A. 80%

B. 99

C. 96

D. 96%

$$\hat{x}_i = \frac{x_i - \text{mean}(\{x_i\})}{\text{std}(\{x_i\})}$$

Objectives

- ✳ Median, Percentile, Mode, IQR,
- ✳ Scatter plots for relationships
- ✳ Correlation Coefficient
- ✳ Other Visualization for *relationships*
Heatmap, 3D bar, Time series plots,

Median

- ✳ We first sort the data set $\{x_i\}$
- ✳ Then *if* the number of items N is **odd**
median = middle item's value
if the number of items N is **even**
median = mean of middle 2 items' values

Properties of Median

- Scaling data scales the median

$$\text{median}(\{k \cdot x_i\}) = k \cdot \text{median}(\{x_i\})$$

- Translating data translates the median

$$\text{median}(\{x_i + c\}) = \text{median}(\{x_i\}) + c$$

Percentile

- ✿ k^{th} percentile is the value relative to which $k\%$ of the data items have smaller or equal numbers
- ✿ Median is roughly the 50^{th} percentile

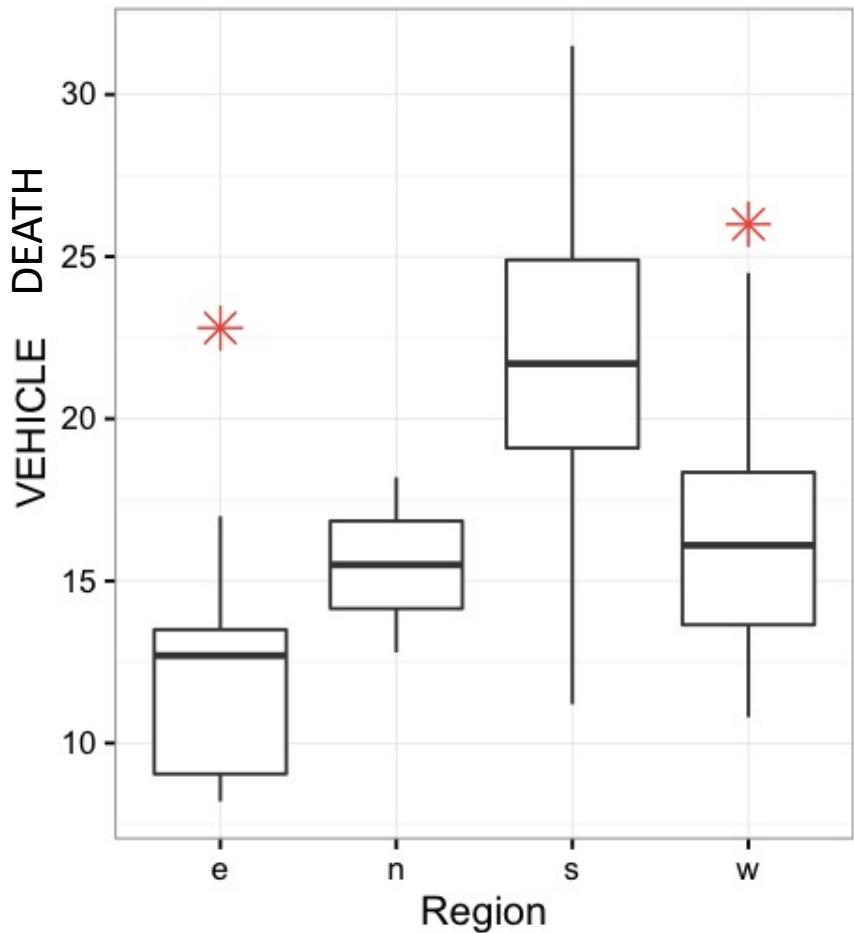
Interquartile range

- ✿ $iqr = (75\text{th percentile}) - (25\text{th percentile})$
- ✿ Scaling data scales the interquartile range
$$iqr(\{k \cdot x_i\}) = |k| \cdot iqr(\{x_i\})$$
- ✿ Translating data does **NOT** change the interquartile range
$$iqr(\{x_i + c\}) = iqr(\{x_i\})$$

Box plots

- ✿ Boxplots
- ✿ Simpler than histogram
- ✿ Good for outliers
- ✿ Easier to use for comparison

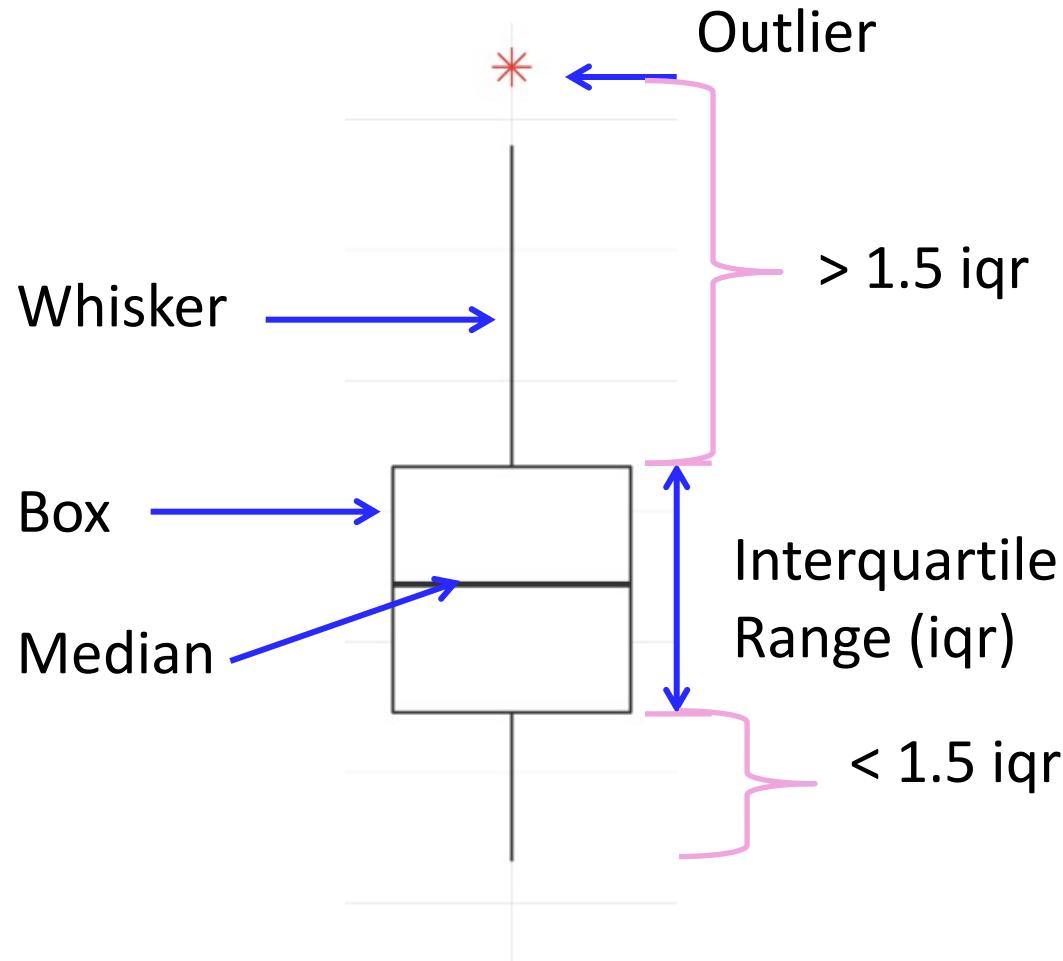
Vehicle death by region



Data from
<https://www2.stetson.edu/~jrasp/data.htm>

Boxplots details, outliers

How to
define
outliers?
(the default)



Sensitivity of summary statistics to outliers

- mean and standard deviation are very sensitive to outliers
- median and interquartile range are not sensitive to outliers

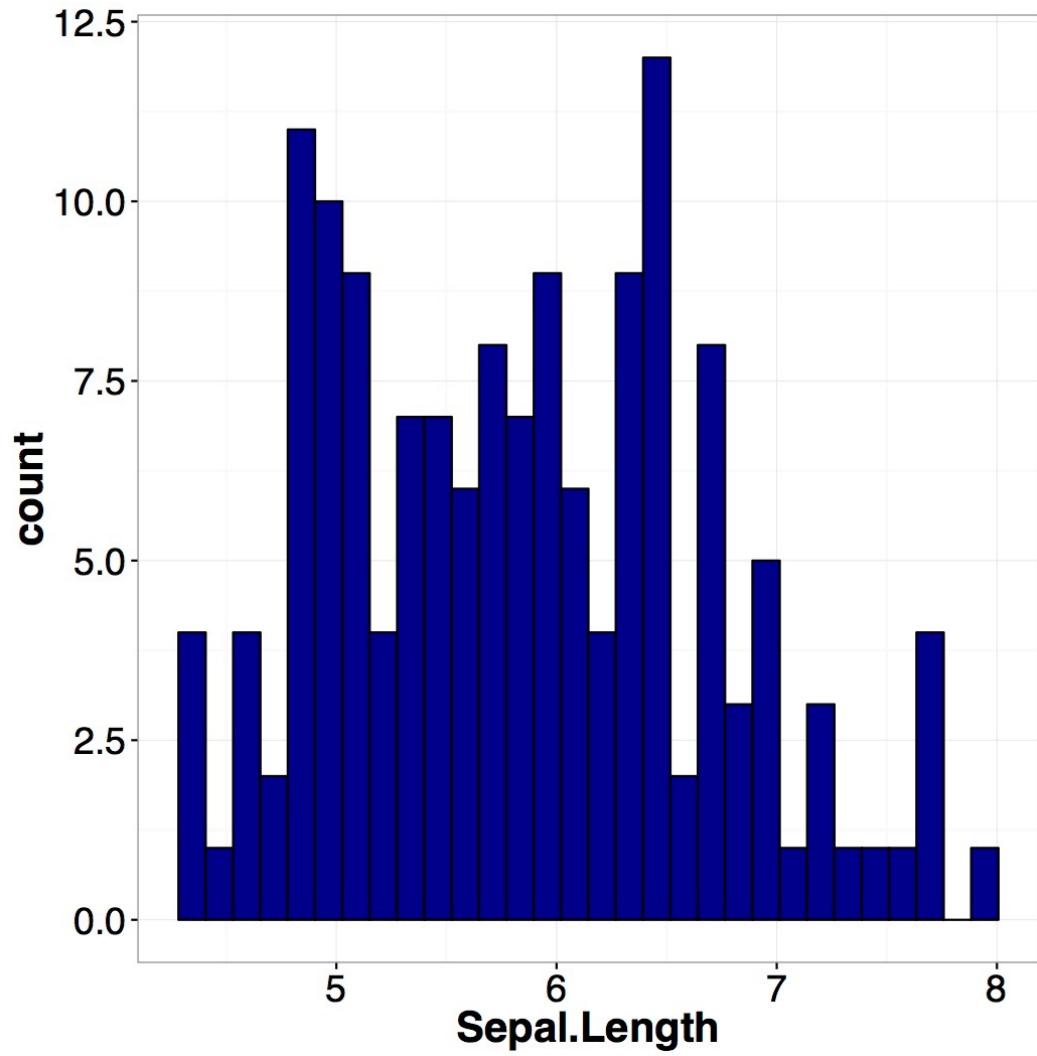
Modes

- ✳ Modes are peaks in a histogram
- ✳ If there are more than 1 mode, we should be curious as to why

Multiple modes

- ✿ We have seen
the “iris” data
which looks to
have several
peaks

Data: “iris” in R

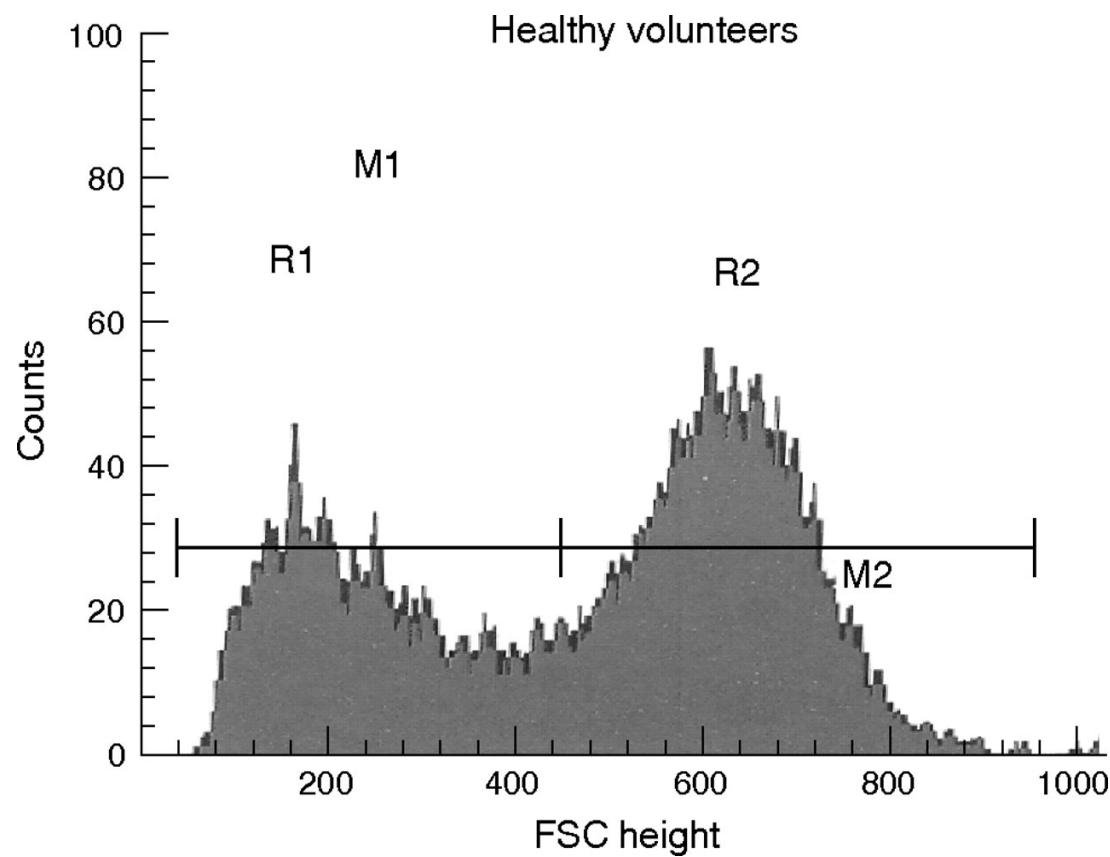


Example Bi-modes distribution

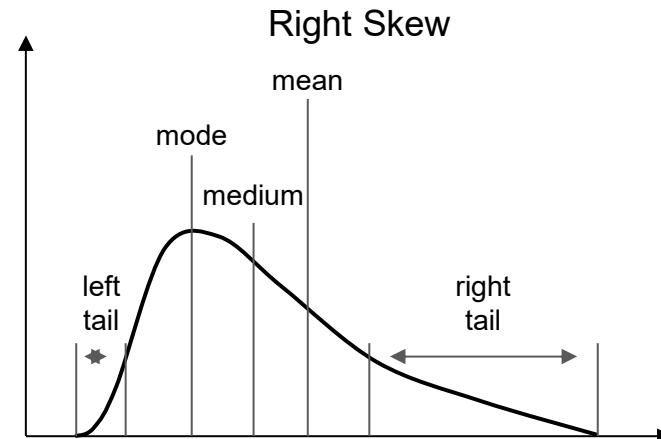
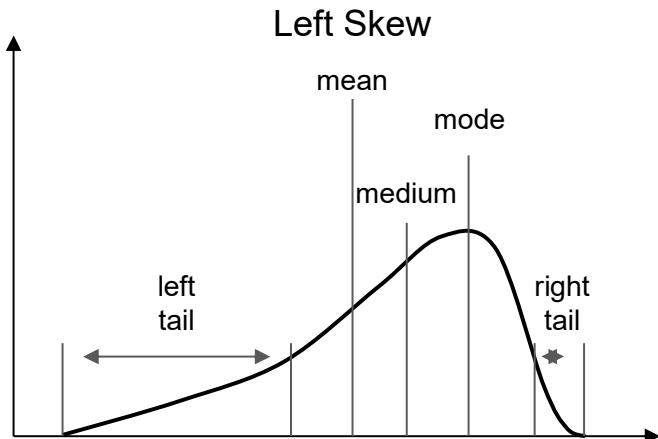
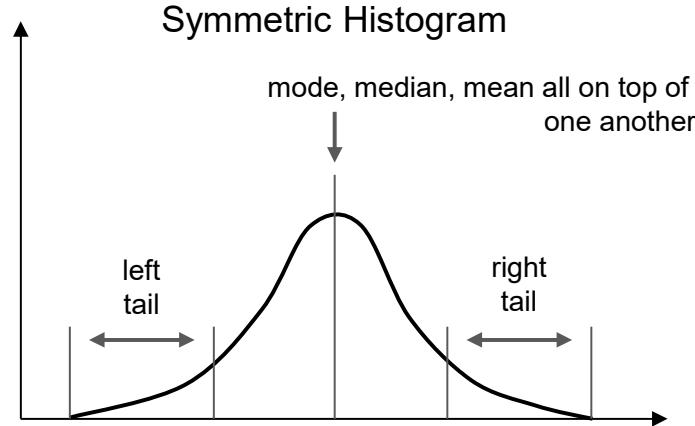
- Modes may indicate multiple populations

Data: Erythrocyte cells in healthy humans

Piagnerelli, JCP 2007



Tails and Skews



Relationship between data features

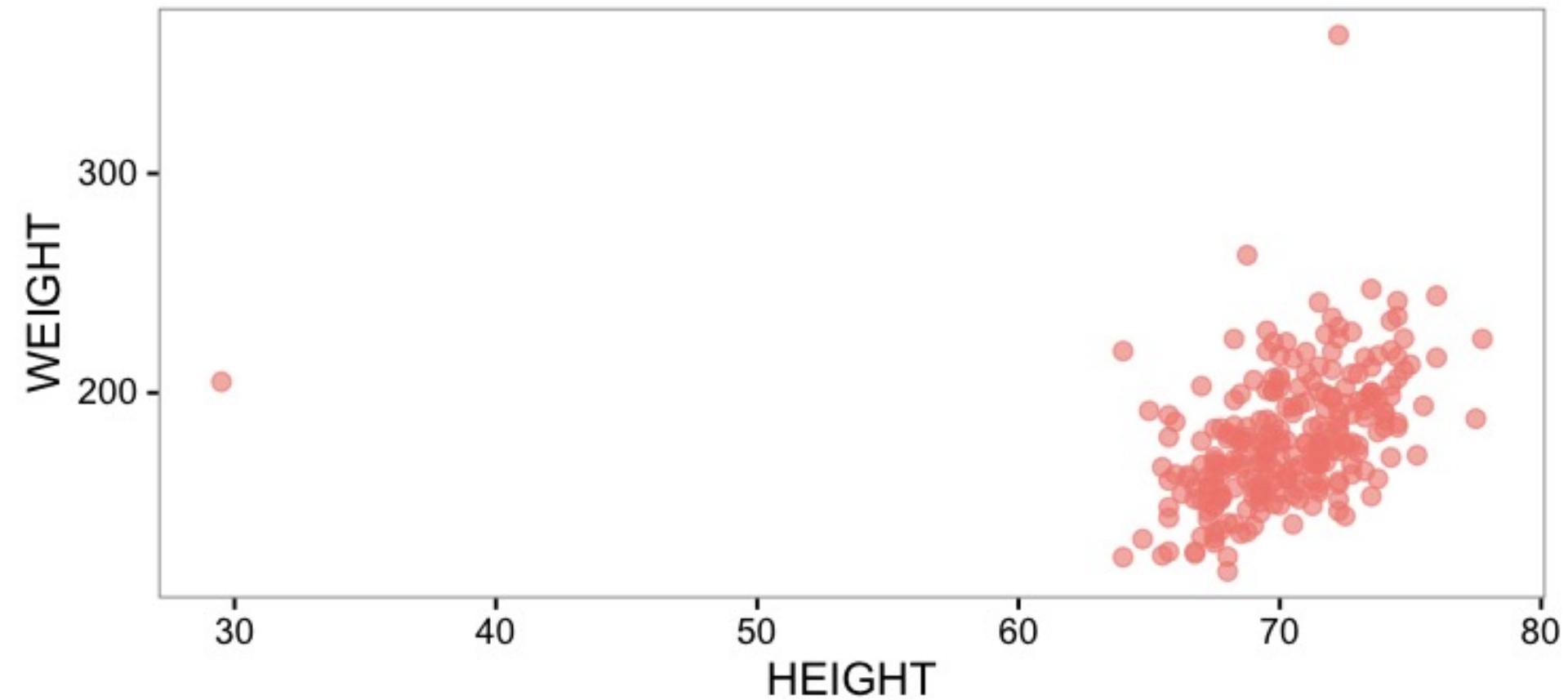
- Example: Does the weight of people relate to their height?

IDNO	BODYFAT	DENSITY	AGE	WEIGHT	HEIGHT
1	12.6	1.0708	23	154.25	67.75
2	6.9	1.0853	22	173.25	72.25
3	24.6	1.0414	22	154.00	66.25
4	10.9	1.0751	26	184.75	72.25
5	27.8	1.0340	24	184.25	71.25
6	20.6	1.0502	24	210.25	74.75
7	19.0	1.0549	26	181.00	69.75
8	12.8	1.0704	25	176.00	72.50
9	5.1	1.0900	25	191.00	74.00
10	12.0	1.0722	23	198.25	73.50

- x : HEIGHT, y: WEIGHT

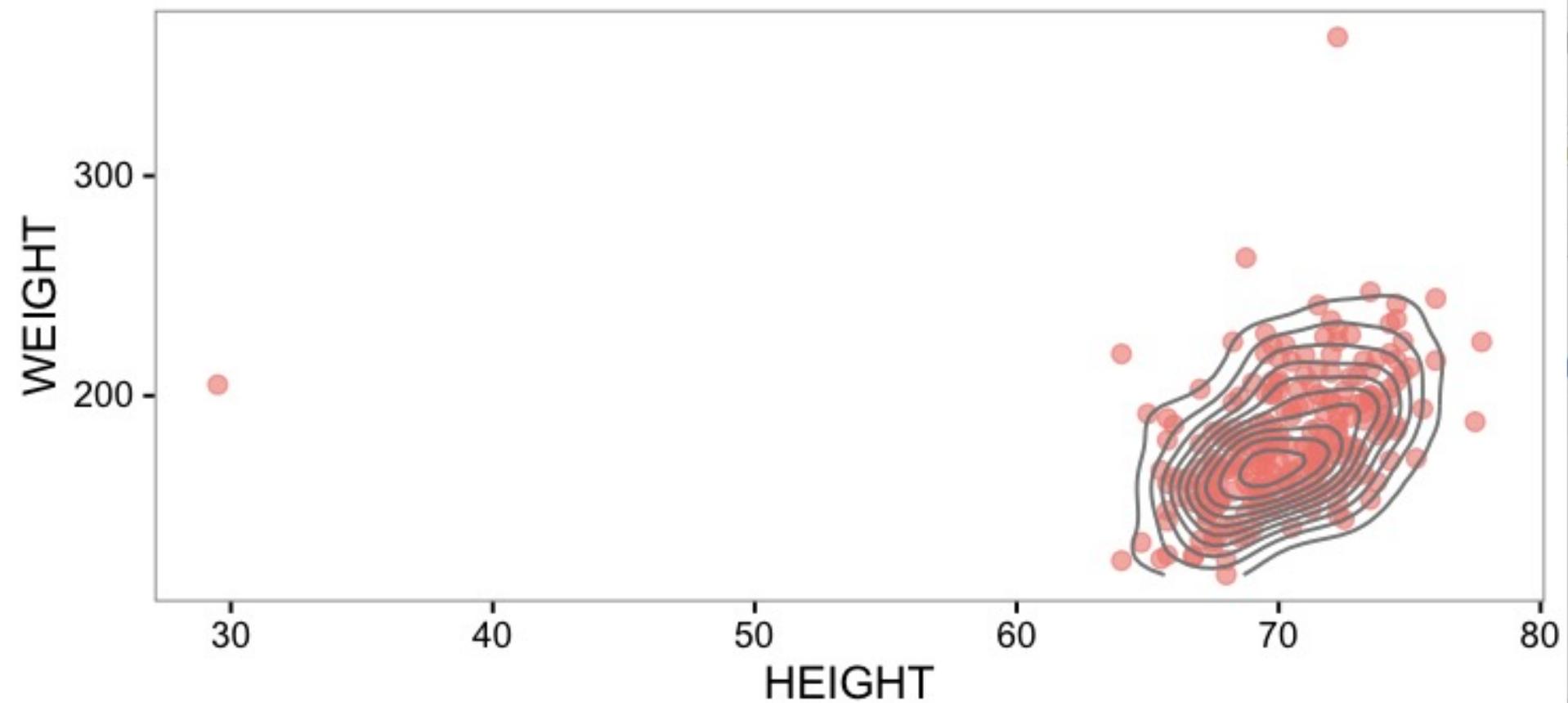
Scatter plot

Body Fat data set



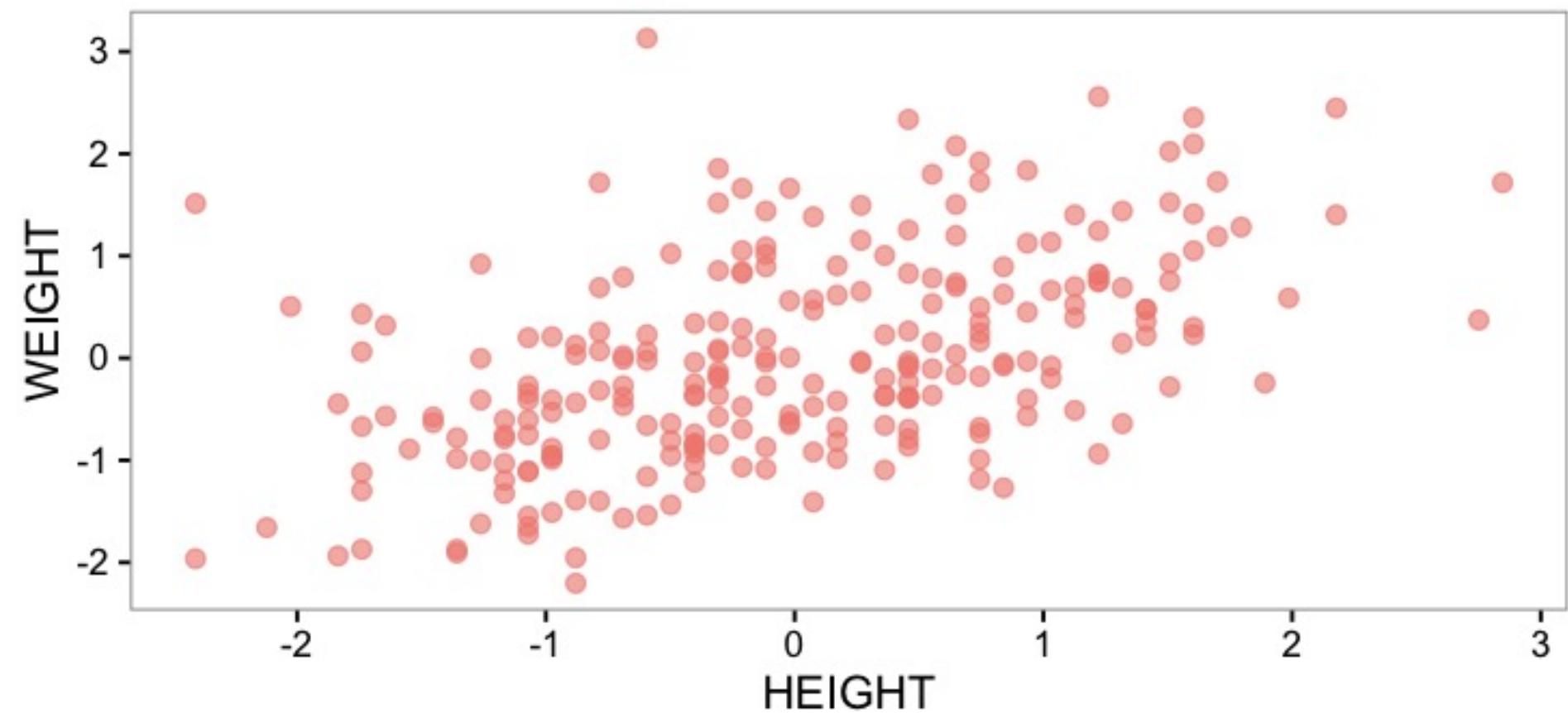
Scatter plot

✳️ Scatter plot with density



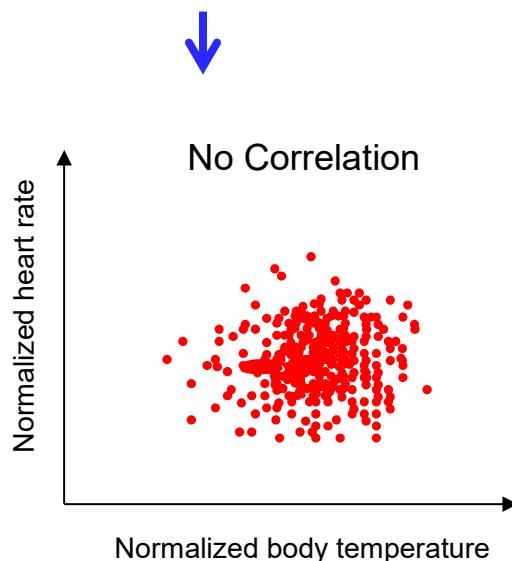
Scatter plot

Removed of outliers & standardized

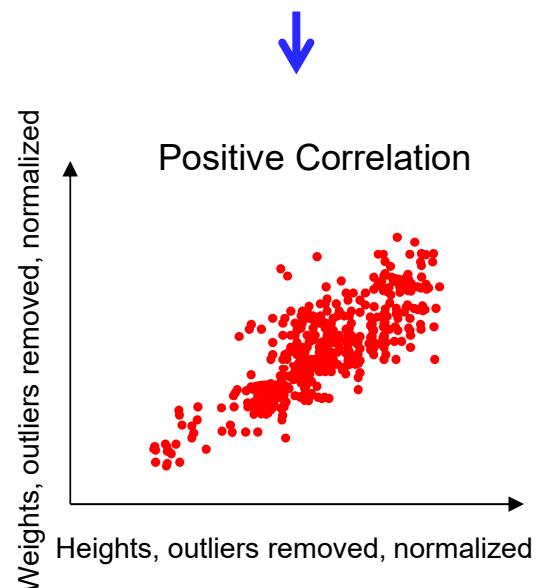


Correlation seen from scatter plots

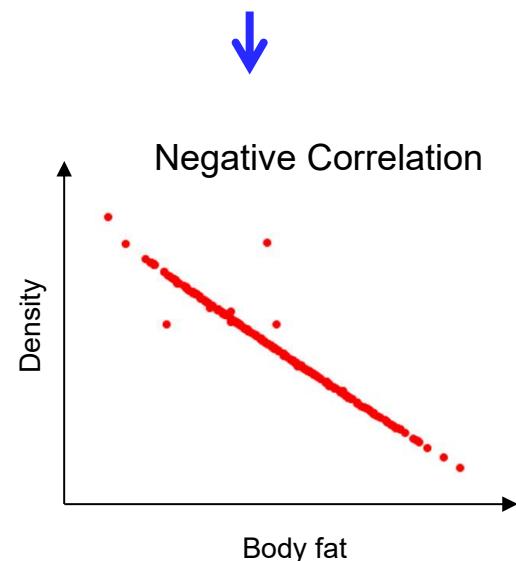
Zero
Correlation



Positive
correlation



Negative
correlation



What kind of Correlation?

- ✳ Line of code in a database and number of bugs
- ✳ Frequency of hand washing and number of germs on your hands
- ✳ GPA and hours spent playing video games
- ✳ earnings and happiness

Correlation Coefficient

Given a data set $\{(x_i, y_i)\}$ consisting of items $(x_1, y_1) \dots (x_N, y_N)$,

Standardize the coordinates of each feature:

$$\hat{x}_i = \frac{x_i - \text{mean}(\{x_i\})}{\text{std}(\{x_i\})} \quad \hat{y}_i = \frac{y_i - \text{mean}(\{y_i\})}{\text{std}(\{y_i\})}$$

Define the correlation coefficient as:

$$\text{corr}(\{(x_i, y_i)\}) = \frac{1}{N} \sum_{i=1}^N \hat{x}_i \hat{y}_i$$

Correlation Coefficient

$$\hat{x}_i = \frac{x_i - \text{mean}(\{x_i\})}{\text{std}(\{x_i\})} \quad \hat{y}_i = \frac{y_i - \text{mean}(\{y_i\})}{\text{std}(\{y_i\})}$$

$$\text{corr}(\{(x_i, y_i)\}) = \frac{1}{N} \sum_{i=1}^N \hat{x}_i \hat{y}_i$$

$$= \text{mean}(\{\hat{x}_i \hat{y}_i\})$$

Q: Correlation Coefficient

- ✳️ Which of the following describe(s) correlation coefficient correctly?
 - A. It's unitless
 - B. It's defined in standard coordinates
 - C. Both A & B

$$\text{corr}(\{(x_i, y_i)\}) = \frac{1}{N} \sum_{i=1}^N \hat{x}_i \hat{y}_i$$

A visualization of correlation coefficient

<https://rpsychologist.com/d3/correlation/>

In a data set $\{(x_i, y_i)\}$ consisting of items
 $(x_1, y_1) \dots (x_N, y_N)$,

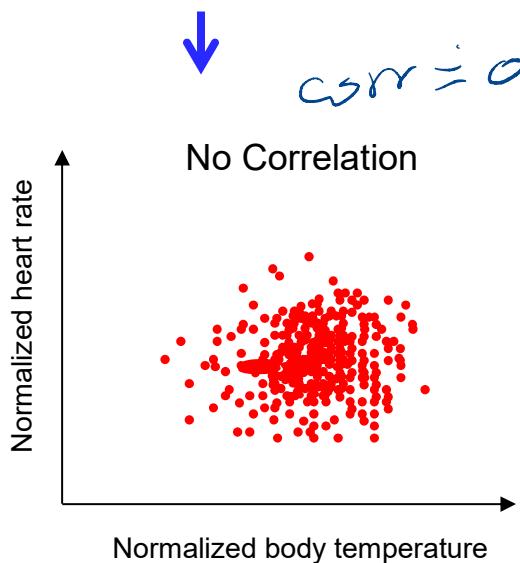
$corr(\{(x_i, y_i)\}) > 0$ shows positive correlation

$corr(\{(x_i, y_i)\}) < 0$ shows negative correlation

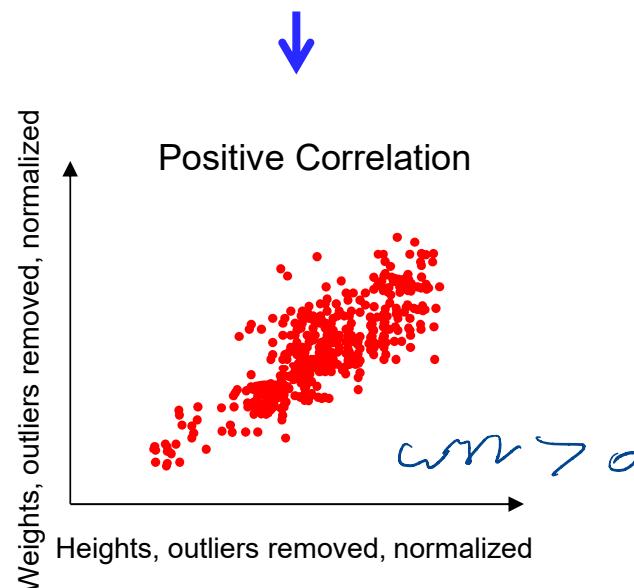
$corr(\{(x_i, y_i)\}) = 0$ shows no correlation

Correlation seen from scatter plots

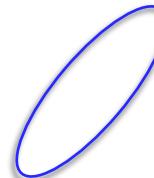
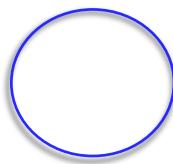
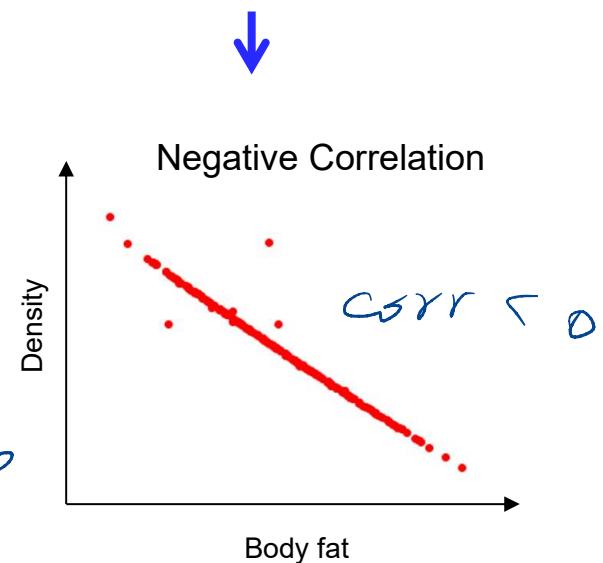
Zero
Correlation



Positive
correlation



Negative
correlation



The Properties of Correlation Coefficient

- ✿ The correlation coefficient is symmetric

$$\text{corr}(\{(x_i, y_i)\}) = \text{corr}(\{(y_i, x_i)\})$$

- ✿ Translating the data does **NOT** change the correlation coefficient

The Properties of Correlation Coefficient

- Scaling the data may change the sign of the correlation coefficient

$$\text{corr}(\{(a x_i + b, c y_i + d)\})$$

$$= \text{sign}(a c) \text{corr}(\{(x_i, y_i)\})$$

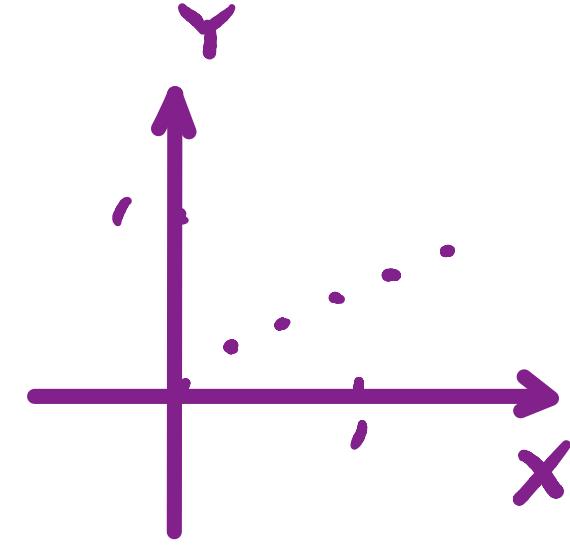
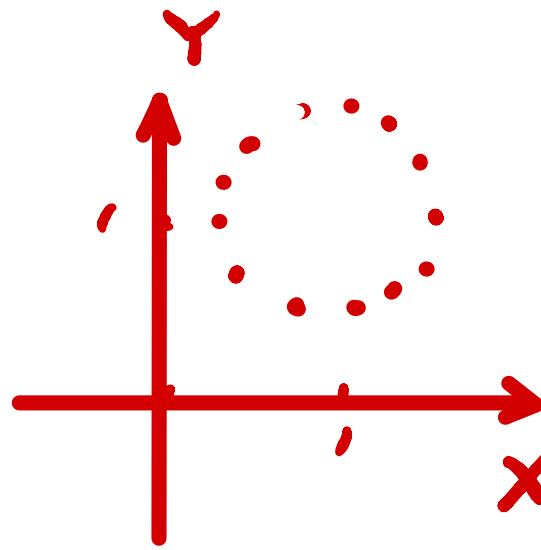
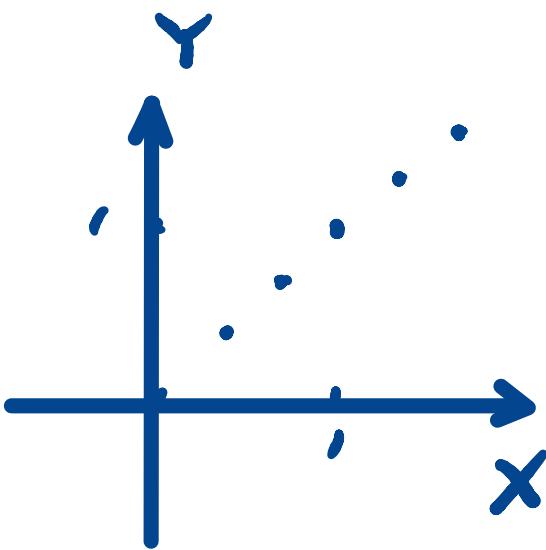
The Properties of Correlation Coefficient

- ★ The correlation coefficient is bounded within [-1, 1]

$\text{corr}(\{(x_i, y_i)\}) = 1$ if and only if $\hat{x}_i = \hat{y}_i$

$\text{corr}(\{(x_i, y_i)\}) = -1$ if and only if $\hat{x}_i = -\hat{y}_i$

Q. Which of the following has correlation coefficient equal to 1?



- A. Left and right
- B. Left
- C. Middle

Review and
finish it at
home.

$$y = ax ; \quad a > 0$$

$$\hat{y} = \frac{ax - \text{mean}\{y\}}{\text{std}\{y\}}$$

$$= \frac{ax - a\text{mean}\{x\}}{a\text{std}\{x\}}$$

$$= \frac{x}{1}$$

$$\text{corr} = \frac{\sum \hat{x}_i \hat{y}_i}{N}$$

Concept of Correlation Coefficient's bound

- The correlation coefficient can be written as

$$\text{corr}(\{(x_i, y_i)\}) = \frac{1}{N} \sum_{i=1}^N \hat{x}_i \hat{y}_i$$

$$\text{corr}(\{(x_i, y_i)\}) = \sum_{i=1}^N \frac{\hat{x}_i}{\sqrt{N}} \frac{\hat{y}_i}{\sqrt{N}}$$

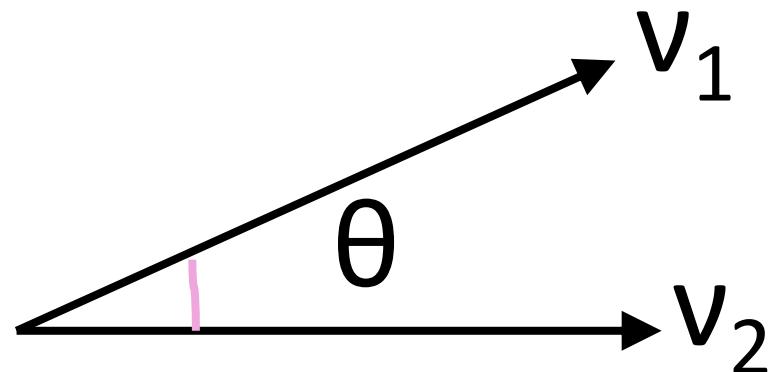
- It's the inner product of two vectors

$$\left\langle \frac{\hat{x}_1}{\sqrt{N}}, \dots, \frac{\hat{x}_N}{\sqrt{N}} \right\rangle \text{ and } \left\langle \frac{\hat{y}_1}{\sqrt{N}}, \dots, \frac{\hat{y}_N}{\sqrt{N}} \right\rangle$$

Inner product

- Inner product's geometric meaning:

$$|\nu_1| |\nu_2| \cos(\theta)$$



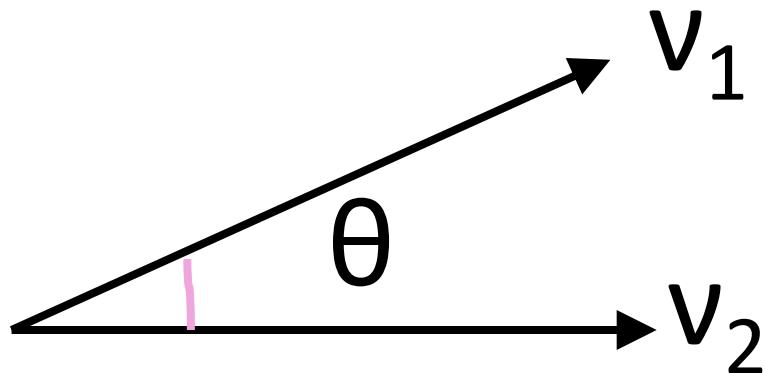
- Lengths of both vectors

$$\nu_1 = \left\langle \frac{\widehat{x_1}}{\sqrt{N}}, \dots, \frac{\widehat{x_N}}{\sqrt{N}} \right\rangle \quad \nu_2 = \left\langle \frac{\widehat{y_1}}{\sqrt{N}}, \dots, \frac{\widehat{y_N}}{\sqrt{N}} \right\rangle$$

are 1

Bound of correlation coefficient

$$|corr(\{(x_i, y_i)\})| = |\cos(\theta)| \leq 1$$



$$v_1 = \left\langle \frac{\widehat{x_1}}{\sqrt{N}}, \dots, \frac{\widehat{x_N}}{\sqrt{N}} \right\rangle \quad v_2 = \left\langle \frac{\widehat{y_1}}{\sqrt{N}}, \dots, \frac{\widehat{y_N}}{\sqrt{N}} \right\rangle$$

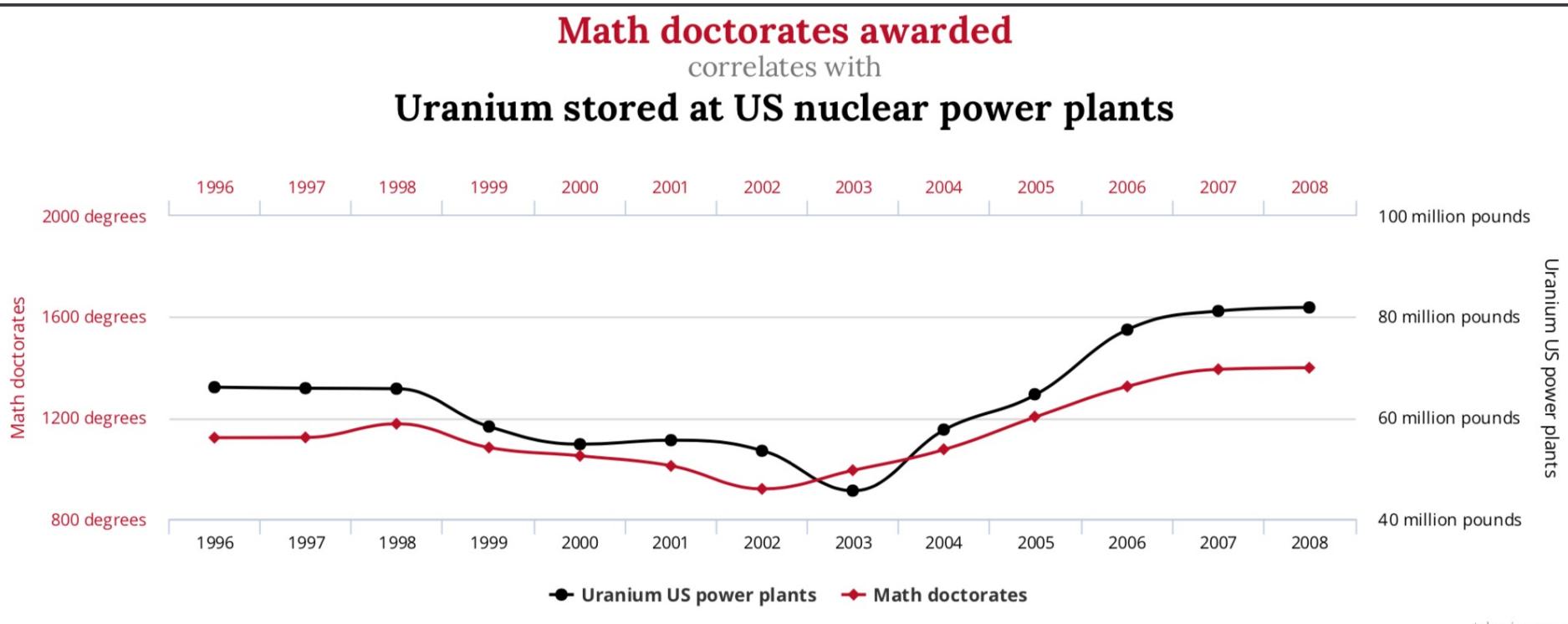
The Properties of Correlation Coefficient

- ★ Symmetric
- ★ Translating invariant
- ★ Scaling only may change sign
- ★ bounded within $[-1, 1]$

Using correlation to predict



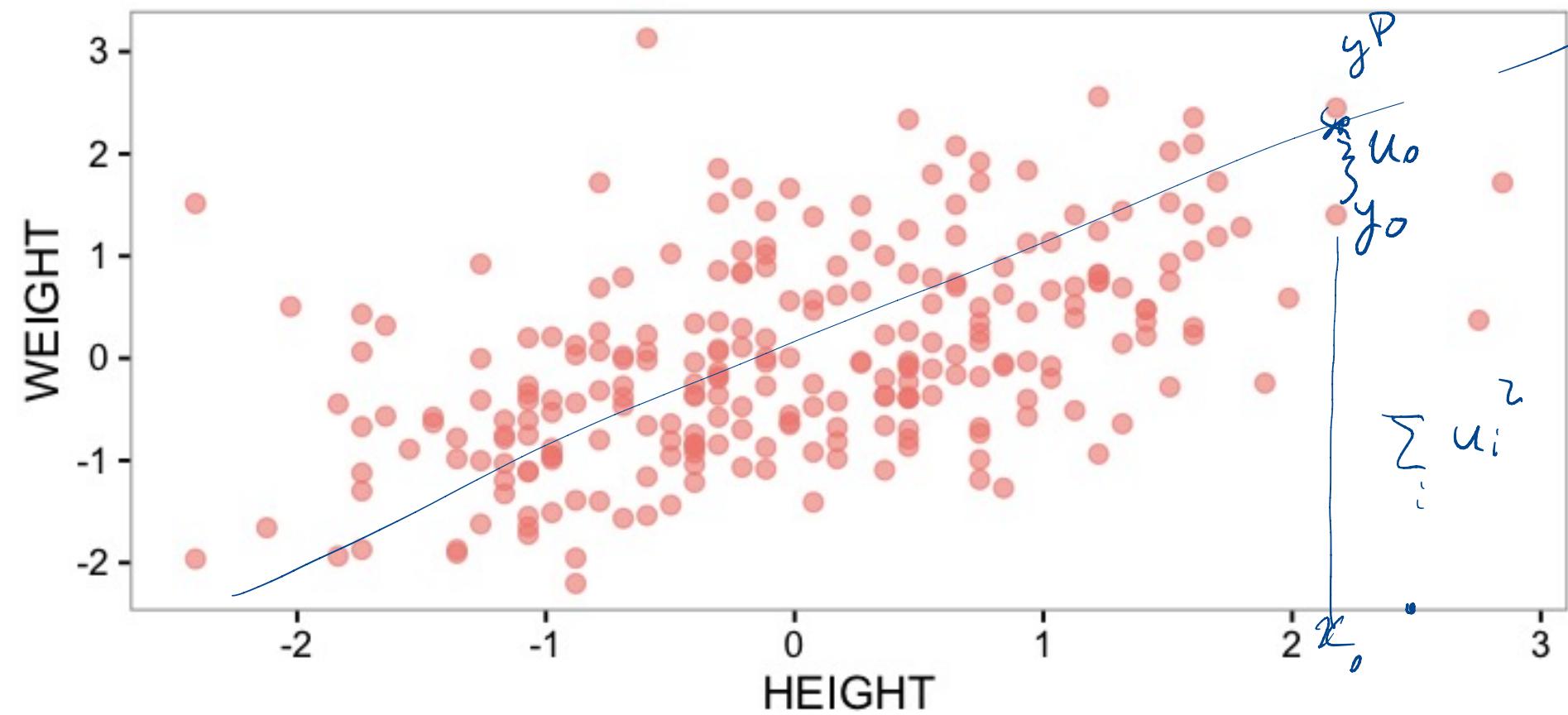
Caution! Correlation is NOT Causation



Credit: Tyler Vigen

How do we go about the prediction?

Removed of outliers & standardized



Using correlation to predict

- Given a correlated data set $\{(x_i, y_i)\}$
we can predict a value y_0^p that goes
with x_0 a value
- In standard coordinates $\{(\hat{x}_i, \hat{y}_i)\}$
we can predict a value \hat{y}_0^p that goes
with \hat{x}_0 a value

Q:

✳️ Which coordinates will you use for the predictor using correlation?

- A. Standard coordinates
- B. Original coordinates
- C. Either

Linear predictor and its error

- ✳ We will assume that our predictor is linear

$$\hat{y}^p = a \hat{x} + b$$

- ✳ We denote the prediction at each \hat{x}_i in the data set as \hat{y}_i^p

$$\hat{y}_i^p = a \hat{x}_i + b$$

- ✳ The error in the prediction is denoted u_i

$$u_i = \hat{y}_i - \hat{y}_i^p = \hat{y}_i - a \hat{x}_i - b$$

Require the mean of error to be zero

We would try to make the mean of error equal to zero so that it is also centered around 0 as the standardized data:

Require the variance of error is minimal

Require the mean of error to be zero

We would try to make the mean of error equal to zero so that it is also centered around 0 as the standardized data: $\text{mean}(\{u_i\}) = 0 \quad \checkmark$

$$\begin{aligned}\text{mean}(\{u_i\}) &= \text{mean}(\{\hat{y} - \hat{y}^P\}) \\ &= \text{mean}(\{\hat{y} - a\hat{x} - b\}) \\ \text{mean}(\{Kx+b\}) &= \text{mean}(\{\hat{y}\}) - \text{mean}(\{\hat{x}\}) \\ &= K \text{mean}(\{x\}) + b\end{aligned}$$

$$= k \text{mean}(\{x\}) + b$$

$$-b = 0$$

$$b = 0$$

Require the variance of error is minimal

minimize $\text{var}(\{u_i\})$

$$\begin{aligned}\text{var}(\{u_i\}) &= \text{mean}(\{u_i - \text{mean}(\{u_i\})\}^2) \\ &= \text{mean}(\{u_i^2\}) \\ &= \text{mean}(\{(\hat{y} - \hat{y}^P)^2\}) \quad u_i = \hat{y} - \hat{y}^P \\ &= \text{mean}(\{(\hat{y} - a\hat{x})^2\}) \quad = \hat{y} - a\hat{x} \\ &= \text{mean}(\{\hat{y}^2 - 2a\hat{x}\hat{y} + a^2\hat{x}^2\}) \\ &= \cancel{\text{mean}(\{\hat{y}^2\})} - 2a\text{mean}(\{\hat{x}\hat{y}\}) + a^2\text{mean}(\{\hat{x}^2\})\end{aligned}$$

Require the variance of error is minimal

minimize $\text{var}(\{u_i\})$

$$\text{var}(\{u_i\}) = \text{mean}(\{u_i - \text{mean}(\{u_i\})\}^2)$$
$$= \text{mean}(\{u_i^2\})$$

$$= \text{mean}(\{(\hat{y} - \hat{y}^P)^2\}) \quad u_i = \hat{y} - \hat{y}^P$$
$$= \text{mean}(\{(\hat{y} - a\hat{x})^2\}) \quad = \hat{y} - a\hat{x}$$
$$\therefore b = 0$$

$$= \text{mean}(\{\hat{y}^2 - 2a\hat{x}\hat{y} + a^2\hat{x}^2\})$$

$$= \text{mean}(\{\hat{y}^2\}) - 2a \text{mean}(\{\hat{x}\hat{y}\}) + a^2 \text{mean}(\{\hat{x}^2\})$$

$$\text{mean}(\{\hat{y}^2\})$$

$$= \text{mean}(\{(\hat{y} - 0)^2\})$$

$$= \text{mean}(\{(\hat{y} - \text{mean}(\hat{y}))^2\})$$

$$= \text{var}(\{\hat{y}\}) = 1$$

\downarrow
 $\text{corr} \rightarrow r$

Require the variance of error is minimal

$$\begin{aligned} \text{var}(\{u\}) &= \text{mean}(\{\hat{y}^2\}) - 2a \text{mean}(\{\hat{x}\hat{y}\}) \\ &\quad + a^2 \text{mean}(\{\hat{x}^2\}) \\ &= 1 - 2a \text{mean}(\{\hat{x}\hat{y}\}) + a^2 \\ &= 1 - 2a \text{corr}(\{x, y\}) + a^2 \quad \text{var}(\{u\}) \\ r &= \text{corr}(\{x, y\}) \quad \text{argmin}(\downarrow) \\ &= 1 - 2ar + a^2 \quad a = ? \end{aligned}$$
$$\frac{d \text{var}(\{u\})}{da} = -2r + 2a = 0$$

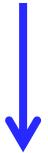
Require the variance of error is minimal

$$\hat{y}^P = \alpha \hat{x} + b$$
$$= r \hat{x}$$

$$\begin{aligned}a &= r \\b &= 0\end{aligned}$$

Here is the linear predictor!

$$\hat{y}^p = r \hat{x}$$



Correlation coefficient

Prediction Formula

✳ In standard coordinates

$$\widehat{y}_0^p = r \widehat{x}_0 \text{ where } r = \text{corr}(\{(x_i, y_i)\})$$

✳ In original coordinates

$$\frac{y_0^p - \text{mean}(\{y_i\})}{\text{std}(\{y_i\})} = r \frac{x_0 - \text{mean}(\{x_i\})}{\text{std}(\{x_i\})}$$

Root-mean-square (RMS) prediction error



Given $\text{var}(\{u_i\}) = 1 - 2ar + a^2$
& $a = r$

$$\text{var}(\{u_i\}) = 1 - r^2$$

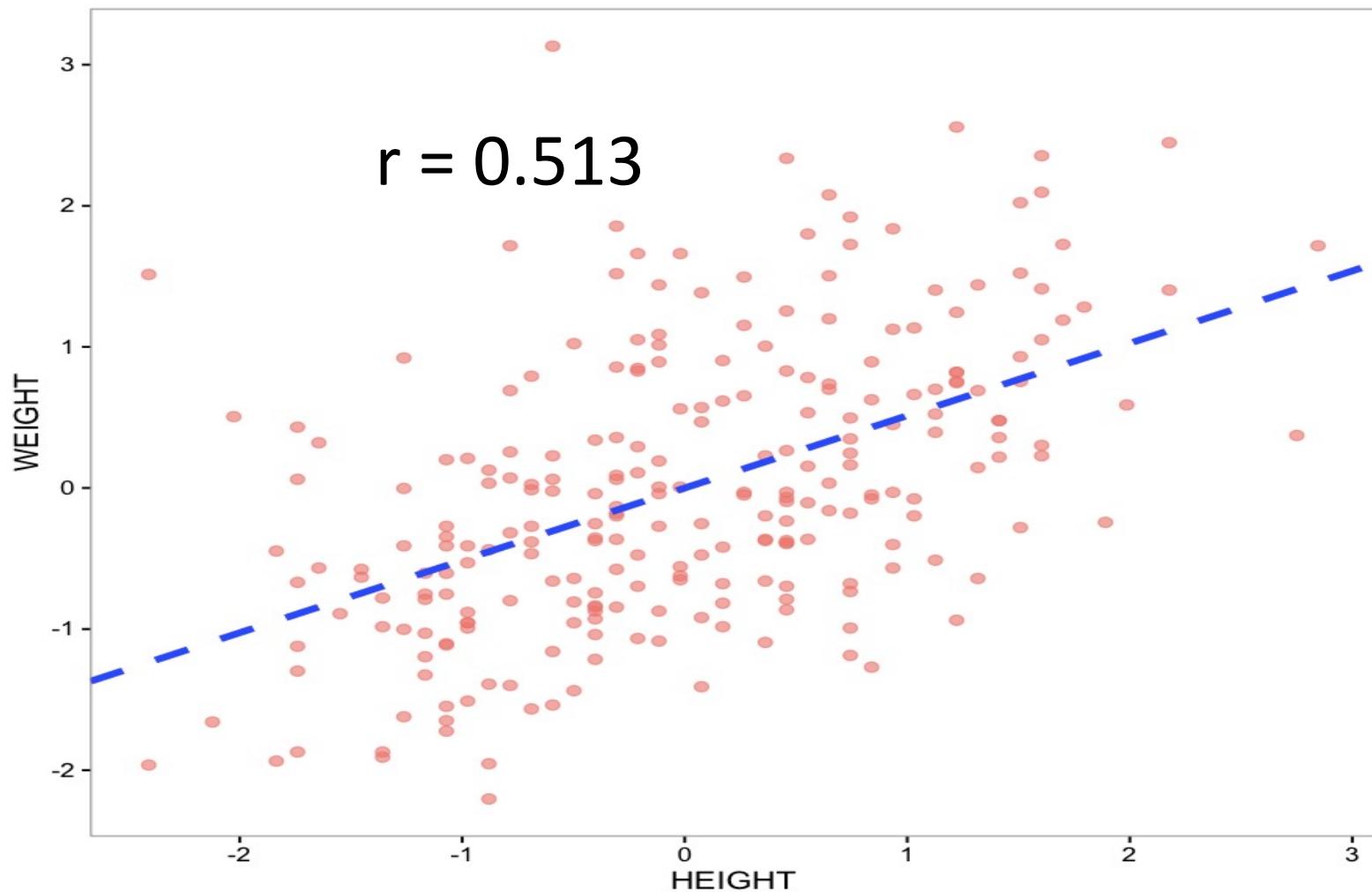


$$\begin{aligned}\text{RMS error} &= \sqrt{\text{mean}(\{u_i^2\})} \\ &= \sqrt{\text{var}(\{u_i\})} \\ &= \sqrt{1 - r^2}\end{aligned}$$

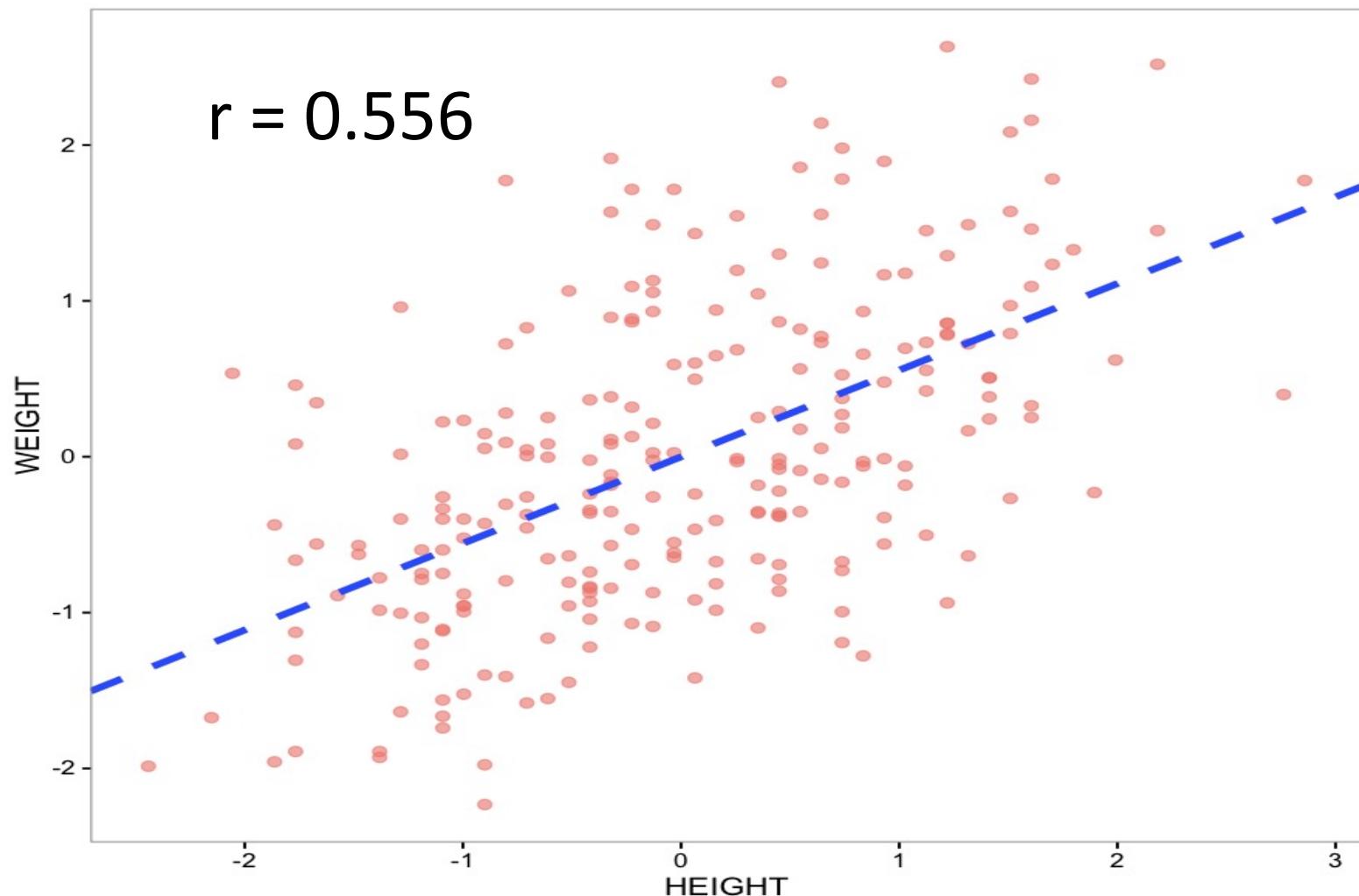
See the error through simulation

<https://rpsychologist.com/d3/correlation/>

Example: Body Fat data

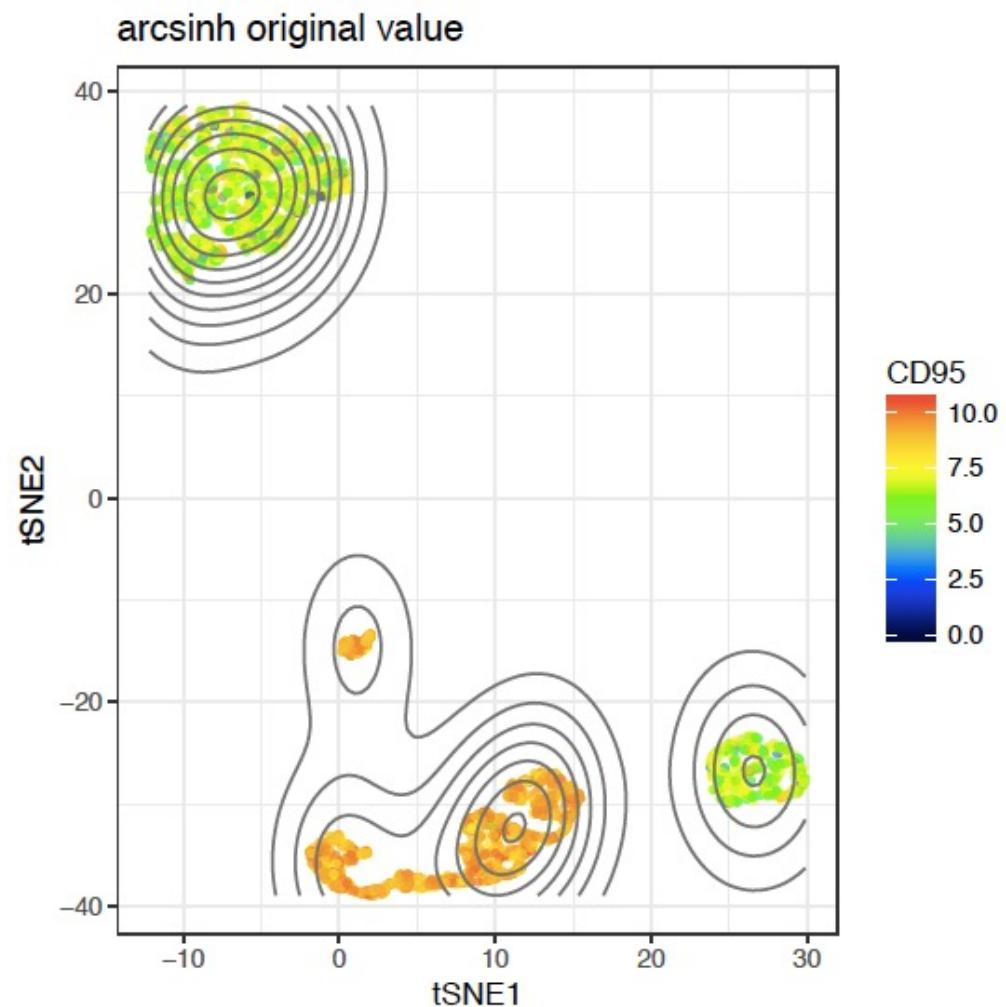


Example: remove 2 more outliers

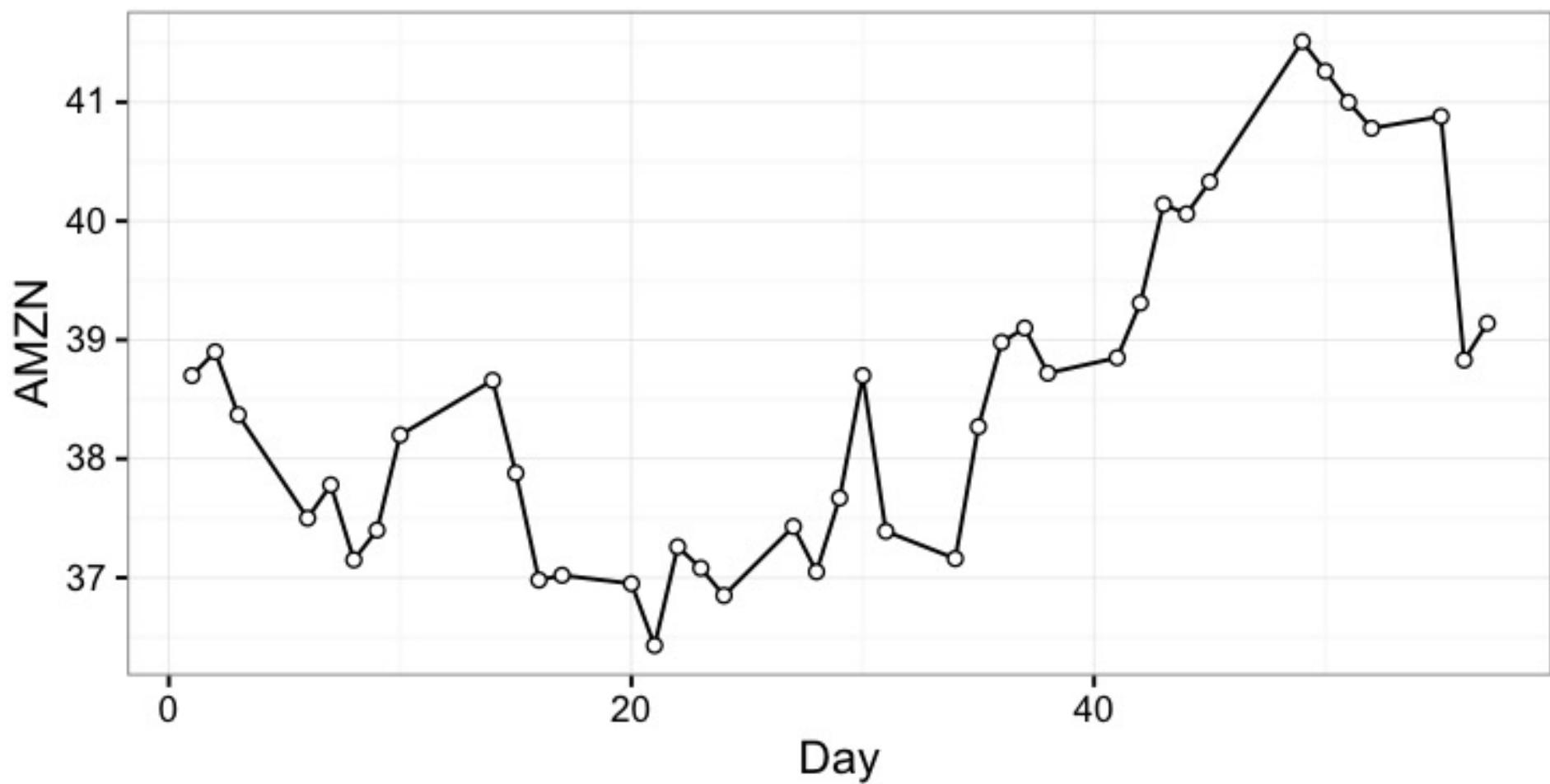


Scatter plot

✳️ Coupled with heatmap to show a 3rd feature



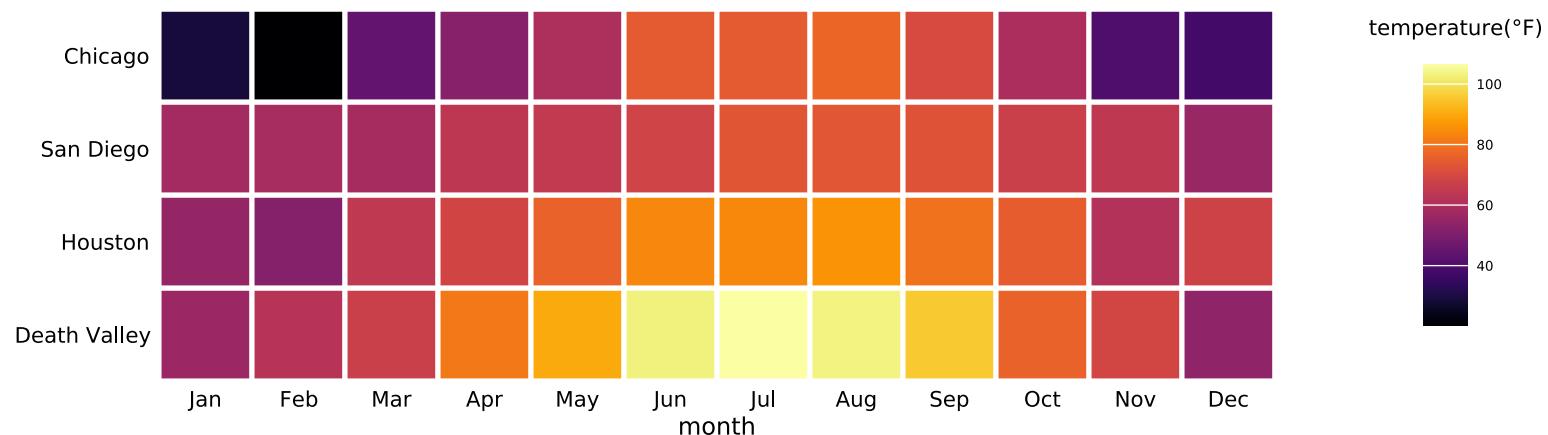
Time Series Plot: Stock of Amazon



Heatmap

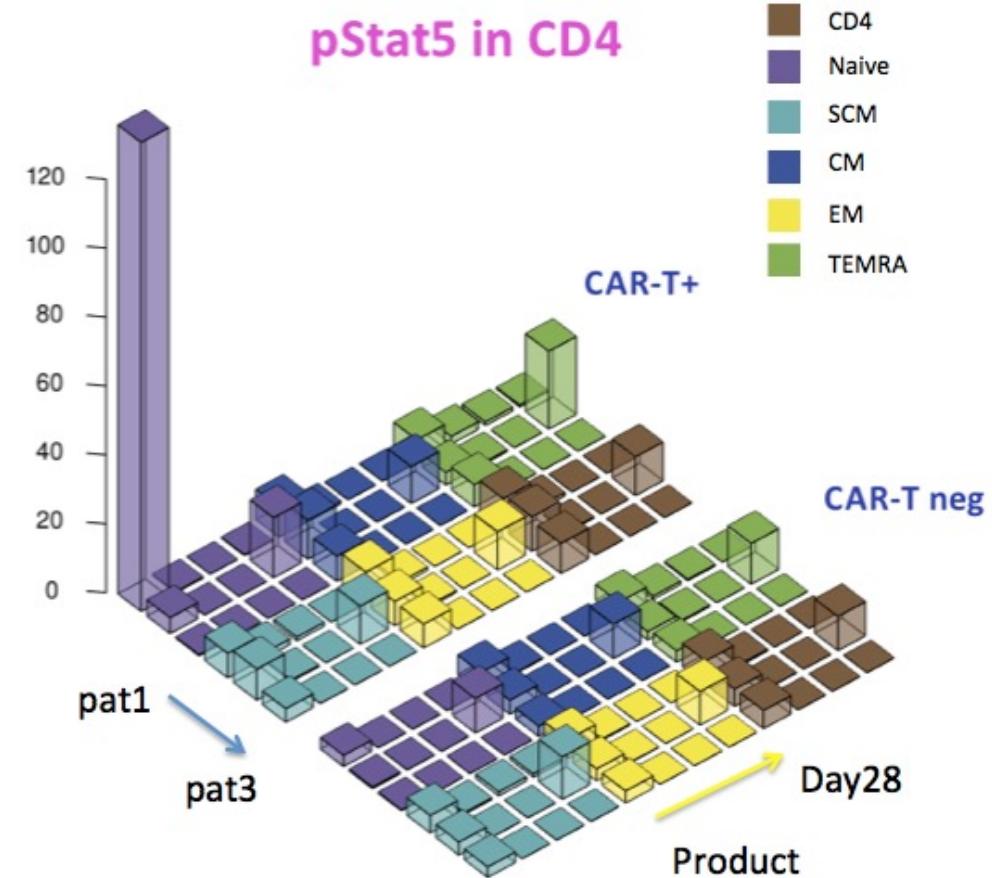
✿ Display matrix of data via gradient of color(s)

Summarization of 4 locations' annual mean temperature by month



3D bar chart

✳️ Transparent
3D bar chart
is good for
small # of
samples
across
categories



Assignments

- ✿ Quiz1 open at 4:30pm today on PL
- ✿ Finish reading Chapter 2 of the textbook
- ✿ Work on the Week 2 module on Canvas
- ✿ Next time: Probability a first look

Additional References

- ★ Charles M. Grinstead and J. Laurie Snell
"Introduction to Probability"
- ★ Morris H. Degroot and Mark J. Schervish
"Probability and Statistics"

See you next time

*See
You!*

