

Intro to Topological Data Analysis

Statistics and Dimension Reduction

Chunyin Siu, Feb 14, 2025

Intro to Topological Data Analysis

Statistics and Dimension Reduction

Chunyin Siu, Feb 14, 2025

Intro to **Topological Data Analysis**

Statistics and Dimension Reduction

Chunyin Siu, Feb 14, 2025

What is Data Analysis

What is Data Analysis

1000110110001011000110111001011011100111001
00110001001110110001110100111011100010110101
00110111000101100011011100101101110011100100
11000100111011000111010011101110001011010110
00110111000101100011011100101101110011100100
11000100111011000111010011101110001011010110
00011011100010110001101110010110111001110010
01100010011101100011101001110111000101101011
10001101110001011000110111001011011100111001
00110001001110110001110100111011100010110101
11011100010110001101110010110111001110010001
00010011101100011101001110111000101101011010
00110111000101100011011100101101110011100100
11000100111011000111010011101110001011010110
00110111000101100011011100101101110011100100

Get Data

Sérgio Valle Duarte <https://commons.wikimedia.org/wiki/File:BinaryData.jpg>

What is Data Analysis

```
10001101100010110001101100101101100111001  
00110001001110110001110100111011100010110101  
00110110001011000110110010110110011100100  
11000100111011000111010011101110001011010110  
0011011100010110001101110010110110011100100  
11000100111011000111010011101110001011010110  
00011011100010110001101110010110110011100100  
01100010011101100011101001110111000101101011  
10001101110001011000110111001011011100111001  
0011000100111011000111010011101110001011010110  
11011100010110001101110010110111001110010001  
00010011101100011101001110001011010110100  
0011011100010110001101110010110110011100100  
11000100111011000111010011101110001011010110  
0011011100010110001101110010110110011100100
```

```
15 const LOCALE = globalThis.navigator.language  
16  
17 const div = document.body.appendChild(document.createElement('div'))  
18 const list = div.appendChild(document.createElement('ol'))  
19  
20 const dayNames = new Map()  
21  
22 for (let i = 0; i < 7; ++i) {  
23   const d = Temporal.PlainDate.From({  
24     year: Temporal.Now.plainDateISO().year,  
25     month: 1,  
26     day: i + 1,  
27   })  
28   dayNames.set(d.dayOfWeek, d.toLocaleString(LOCALE, { weekday: 'long' }))  
29 }  
30  
31 for (const num of [...dayNames.keys()].sort((a, b) => a - b)) {  
32   list.appendChild(Object.assign(  
33     document.createElement('li'),  
34     { textContent: dayNames.get(num) },  
35   ))  
36 }  
37  
38 }  
39
```

Get Data

Sérgio Valle Duarte <https://commons.wikimedia.org/wiki/File:BinaryData.jpg>

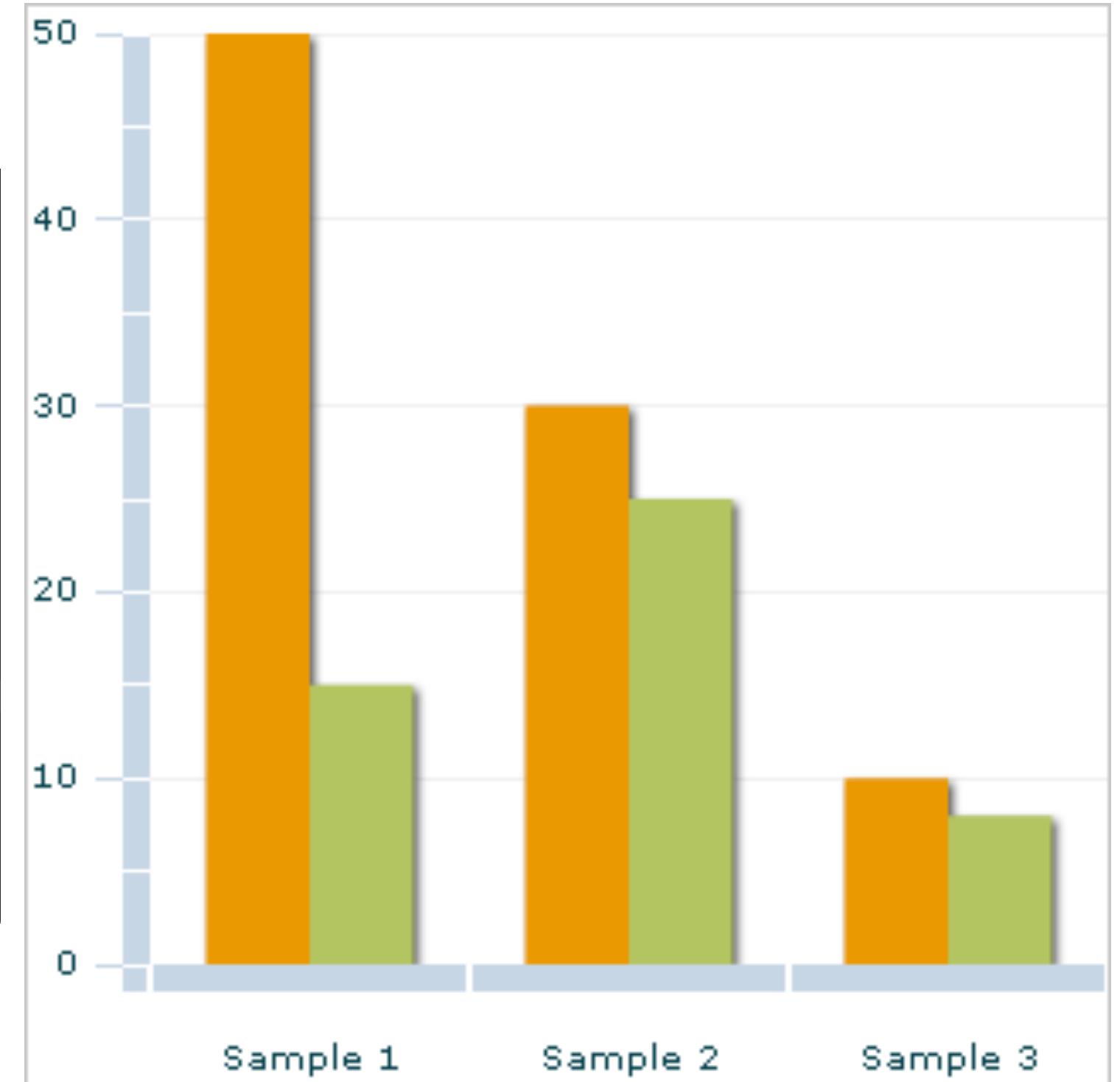
Analyze Data

Lionel Rowe https://commons.wikimedia.org/wiki/File:JavaScript_code.png

What is Data Analysis

```
10001101100010110001101100101101100111001  
0011000100111011000111010011101100010110101  
00110110001011000110110010110110011100100  
11000100111011000111010011101110001011010110  
0011011100010110001101110010110110011100100  
11000100111011000111010011101110001011010110  
000110111000101100011011100101101011001110010  
01100010011101100011101001110111000101101011  
10001101110001011000110111001011011100111001  
0011000100111011000111010011101110001011010110  
11011100010110001101110010110111001110010001  
0001001110110001110100111000101101011010  
00110111000101100011011100101101110011100100  
11000100111011000111010011101110001011010110  
00110111000101100011011100101101110011100100
```

```
15 const LOCALE = globalThis.navigator.language  
16  
17 const div = document.body.appendChild(document.createElement('div'))  
18 const list = div.appendChild(document.createElement('ol'))  
19  
20 const dayNames = new Map()  
21  
22 for (let i = 0; i < 7; ++i) {  
23   const d = Temporal.PlainDate.From({  
24     year: Temporal.Now.plainDateISO().year,  
25     month: 1,  
26     day: i + 1,  
27   })  
28   dayNames.set(d.dayOfWeek, d.toLocaleString(LOCALE, { weekday: 'long' }))  
29 }  
30  
31 for (const num of [...dayNames.keys()].sort((a, b) => a - b)) {  
32   list.appendChild(Object.assign(  
33     document.createElement('li'),  
34     { textContent: dayNames.get(num) },  
35   ))  
36 }  
37  
38 }  
39
```



Get Data

Sérgio Valle Duarte <https://commons.wikimedia.org/wiki/File:BinaryData.jpg>

Analyze Data

Lionel Rowe https://commons.wikimedia.org/wiki/File:JavaScript_code.png

Get Nice Plots

https://commons.wikimedia.org/wiki/File:Adobe_Flex_ColumnChart.png

What is Data Analysis



Get Data

<https://banknxt.com/2021/12/08/tarot-cards-meanings/>

What is Data Analysis



Get Data

<https://banknxt.com/2021/12/08/tarot-cards-meanings/>



Analyze Data

Lionel Rowe https://commons.wikimedia.org/wiki/File:JavaScript_code.png

What is Data Analysis



Get Data

<https://banknxt.com/2021/12/08/tarot-cards-meanings/>



Analyze Data

Lionel Rowe https://commons.wikimedia.org/wiki/File:JavaScript_code.png



Get Nice Plots

<https://slashandsroll.com/when-to-use-the-celestial-cross-spread/>

**Explaining data variability
while acknowledging the effect of random
fluctuations**

Explaining data variability
while acknowledging the effect of **random**
fluctuations

Explaining data variability
while acknowledging the effect of **random**
fluctuations

BELIEF

Four Levels of Statistical Belief

Four Levels of Statistical Belief

Descriptive

mean, variance, etc

Four Levels of Statistical Belief

Frequentist

parameter estimation, p-value

Descriptive

mean, variance, etc

Four Levels of Statistical Belief



Four Levels of Statistical Belief



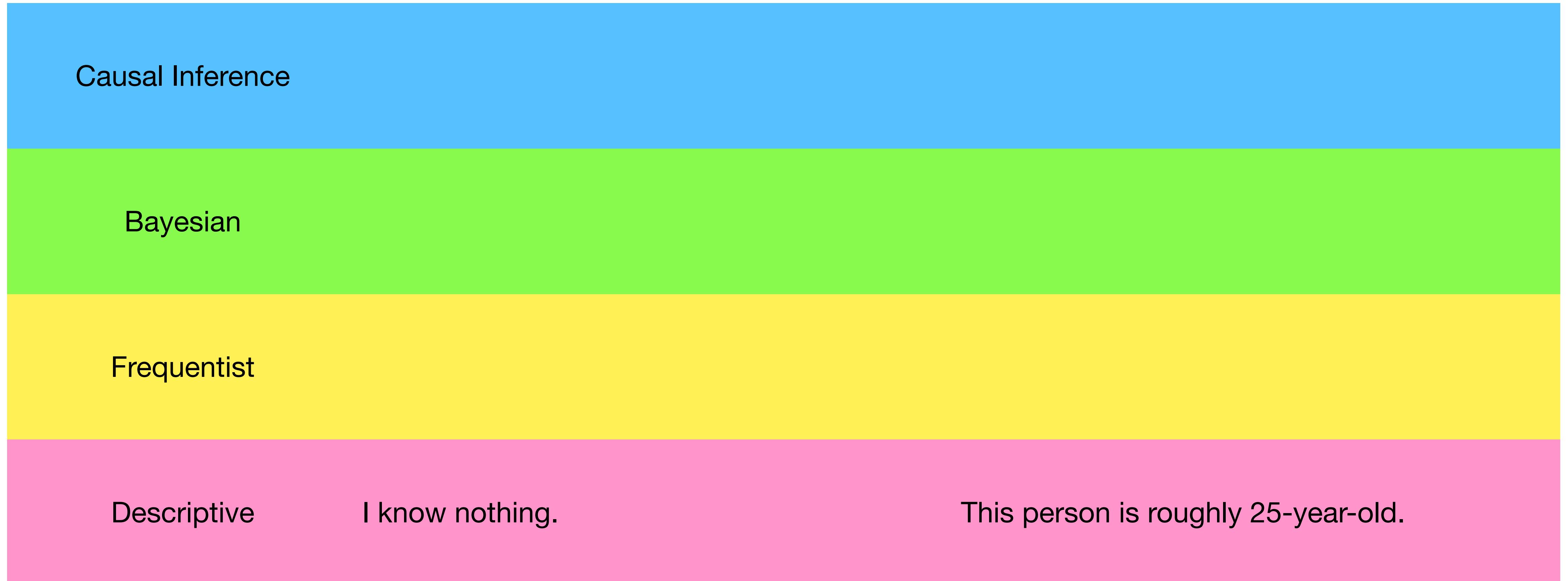
Four Levels of Statistical Belief



Rick? Or Morty?



Rick? Or Morty?



Rick? Or Morty?

Causal Inference

Bayesian

Frequentist

I know all the possible worlds.

This person is either Rick, or Morty.

Descriptive

I know nothing.

This person is roughly 25-year-old.

Rick? Or Morty?

Causal Inference

Bayesian

I know all possible worlds,
and I know the probability of each world.

This person was found in Mortytown.

Frequentist

I know all the possible worlds.

This person is either Rick, or Morty.

Descriptive

I know nothing.

This person is roughly 25-year-old.

Rick? Or Morty?

Causal Inference

I know all parallel universes,
even though only one can materialize.

Will Jessica make Morty happy?

Bayesian

I know all possible worlds,
and I know the probability of each world.

This person was found in Mortytown.

Frequentist

I know all the possible worlds.

This person is either Rick, or Morty.

Descriptive

I know nothing.

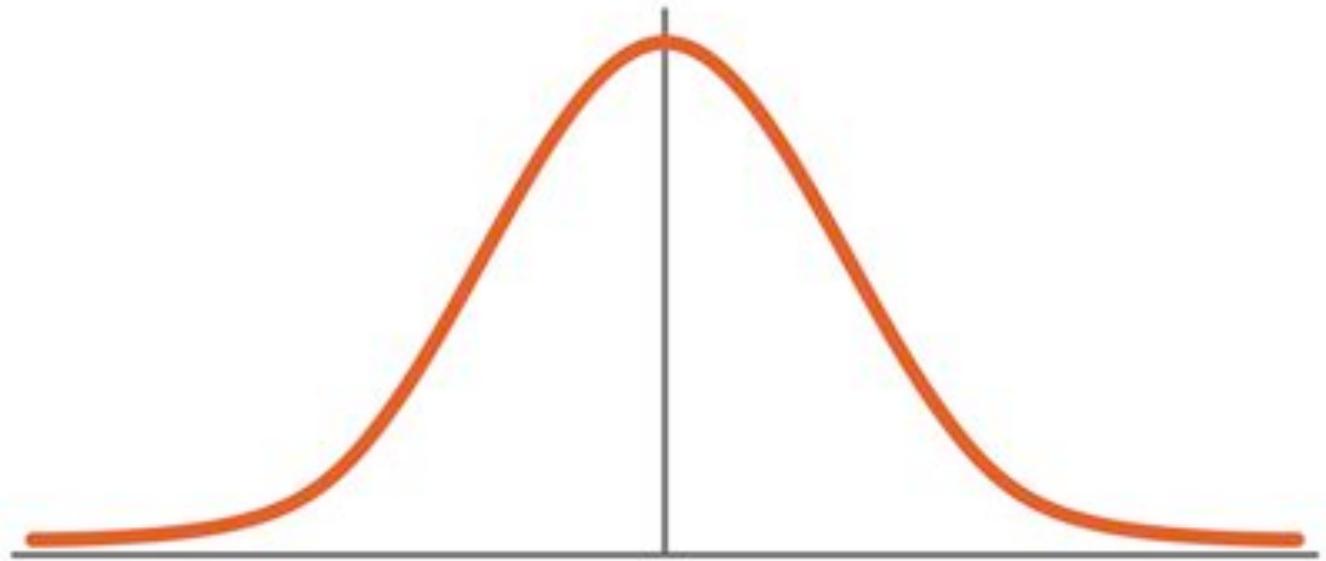
This person is roughly 25-year-old.

In the Citadel

- mean age = 42
- sd = 28

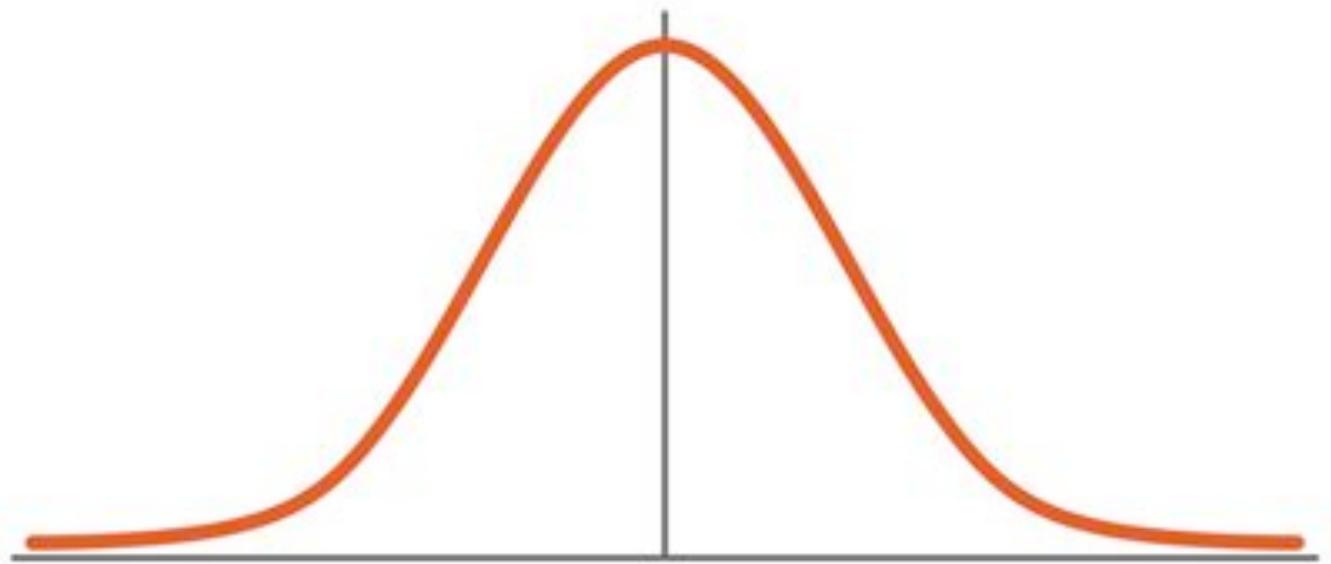
In the Citadel

- mean age = 42
- sd = 28



In the Citadel

- mean age = 42
- sd = 28



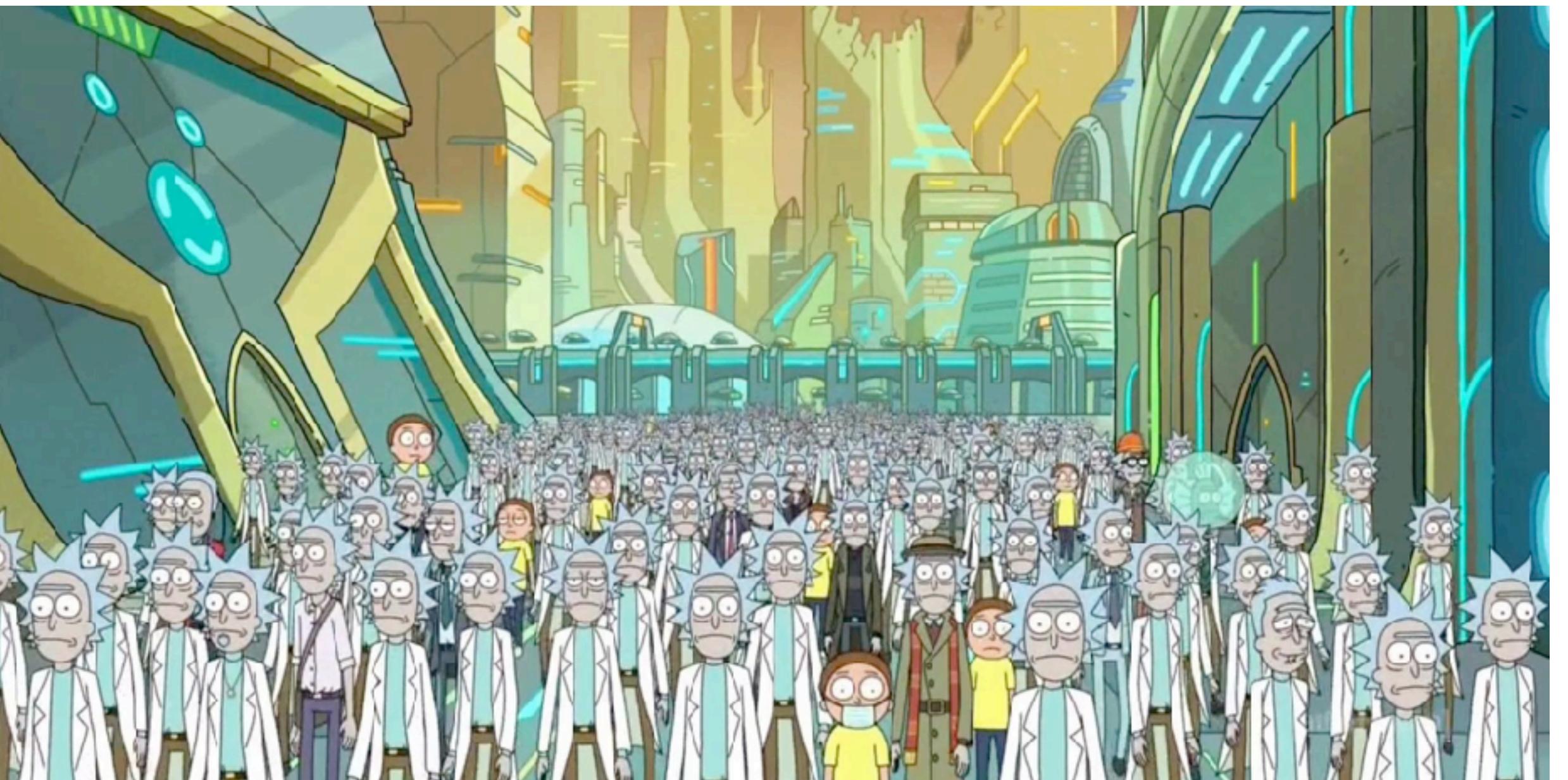
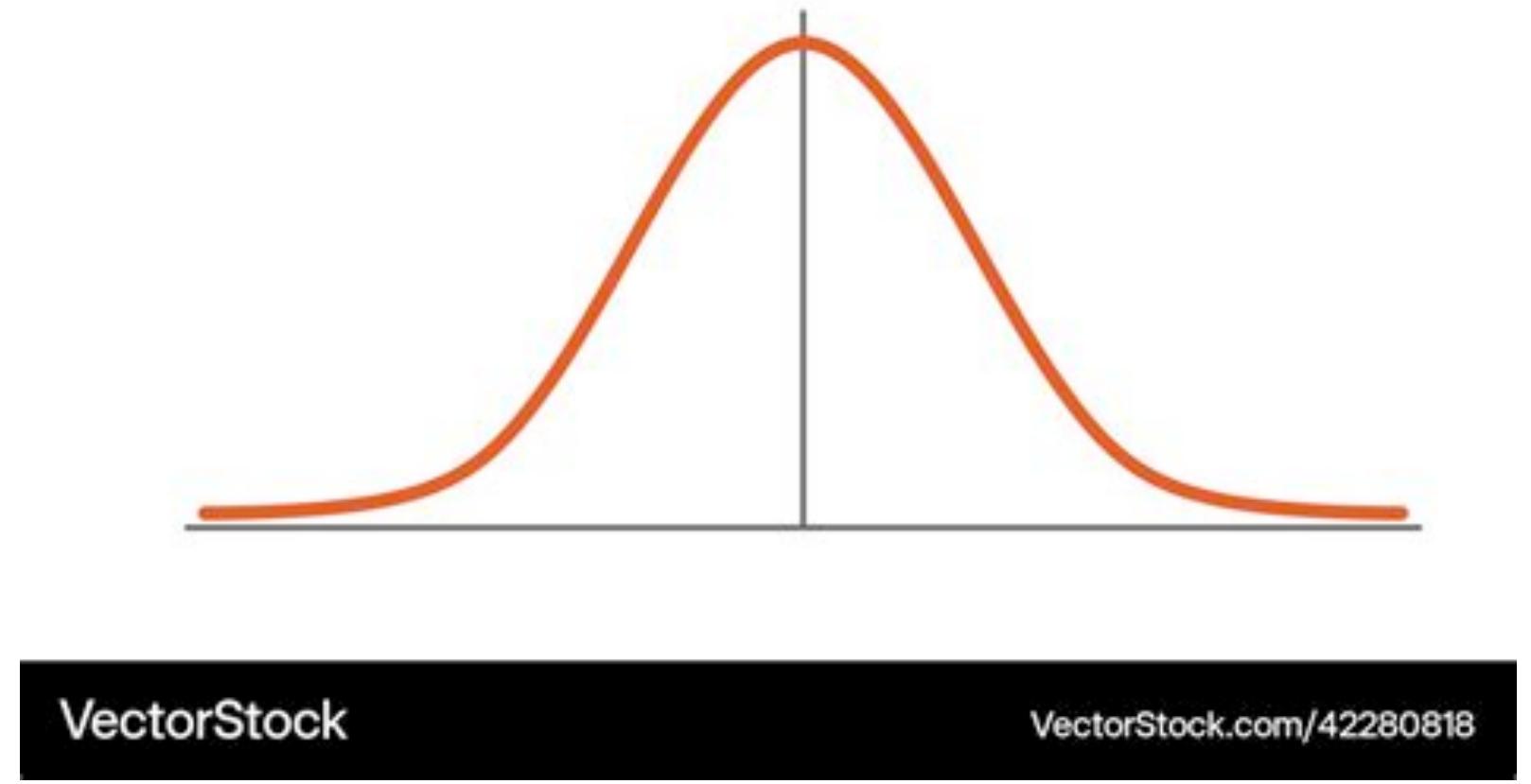
VectorStock

VectorStock.com/42280818

- age of 38% of the population in $[28, 56]$

In the Citadel

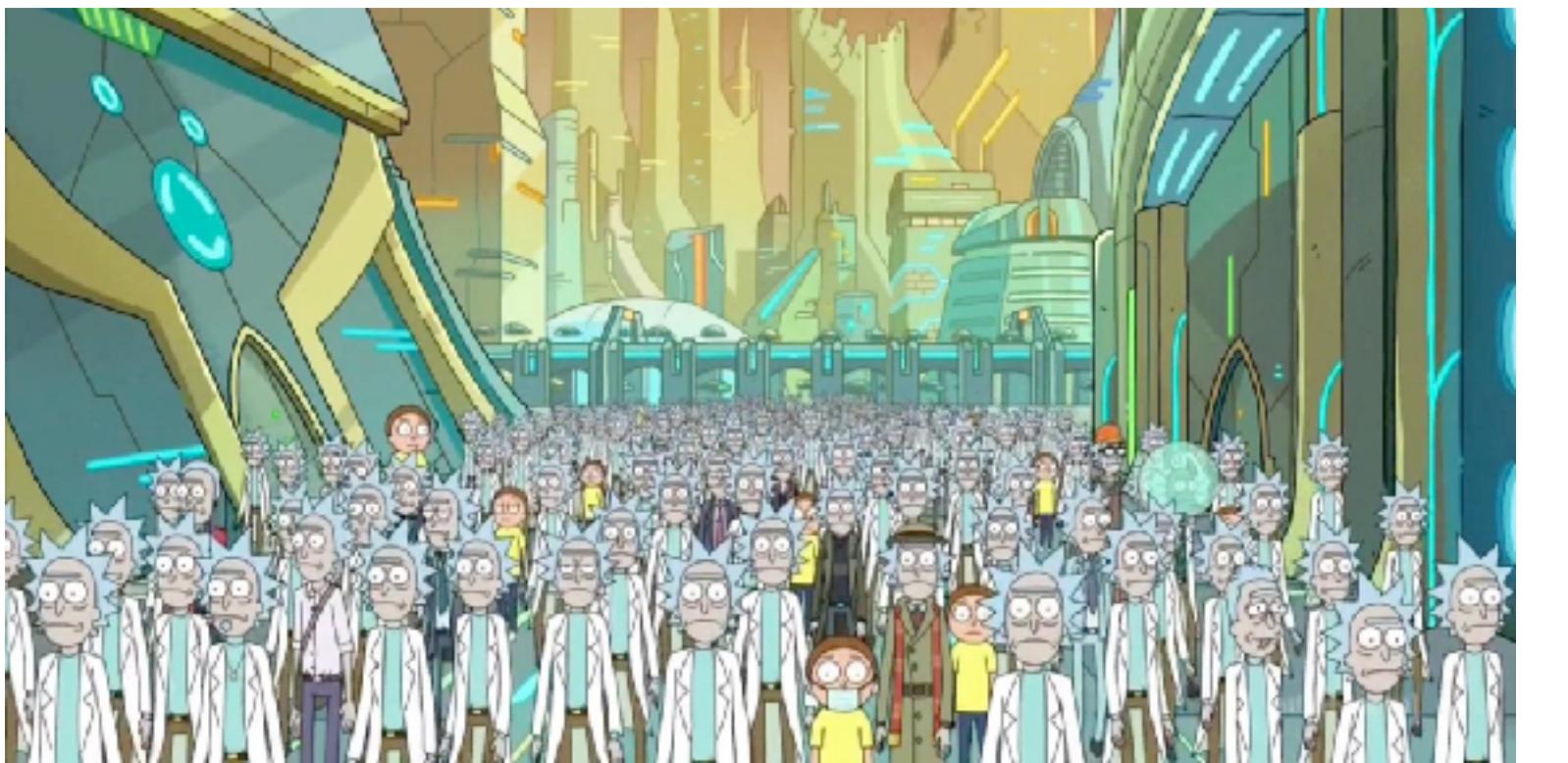
- mean age = 42
- sd = 28



- age of 38% of the population in $[28, 56]$

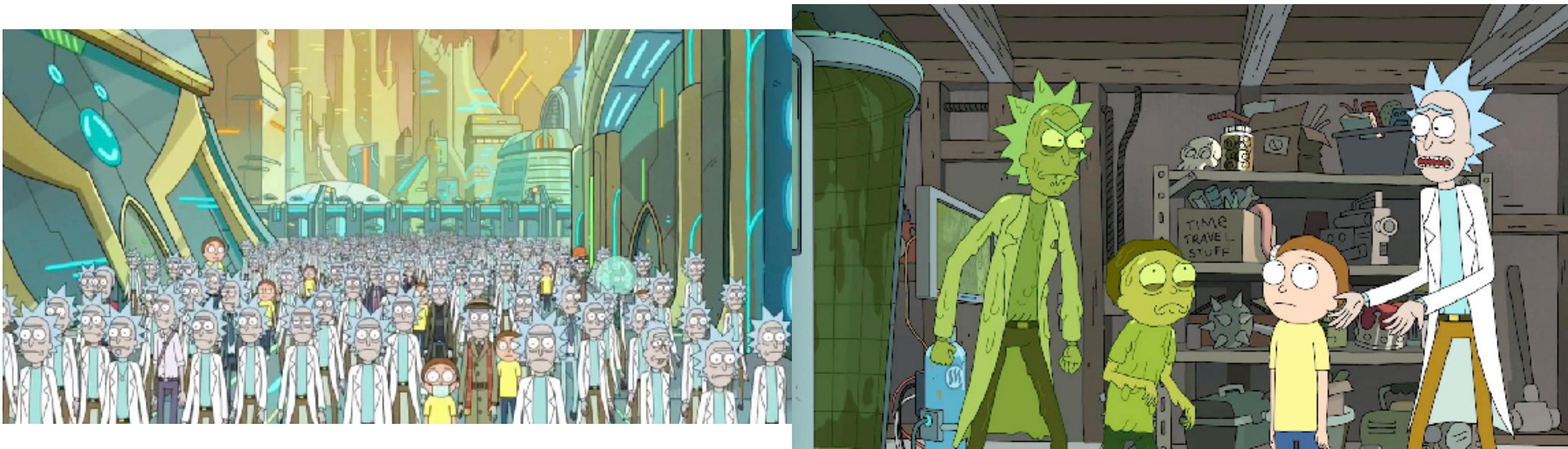
Only Descriptive Statistics Works Here

The cult of mean



mean \neq mode

The cult of mean



mean \neq mode

$(\text{positive} + \text{negative})/2 = 0$

The cult of mean



mean \neq mode

$(\text{positive} + \text{negative})/2 = 0$

$(\text{signal} + \text{noise})/2 = \text{weak signal}$

The cult of $p < 0.05$

The cult of $p < 0.05$

p-value = probability of positive result if it is actually wrong

The cult of $p < 0.05$

p-value = probability of positive result if it is actually wrong

Assumptions

The cult of $p < 0.05$

p-value = probability of positive result if it is actually wrong

Assumptions

A non-Rick is not necessarily a Morty

The cult of $p < 0.05$

p-value = probability of positive result if it is actually wrong

Assumptions	Number of Tests
A non-Rick is not necessarily a Morty	

The cult of $p < 0.05$

p-value = probability of positive result if it is actually wrong

Assumptions	Number of Tests
A non-Rick is not necessarily a Morty	For 100 wrong claims, you get 5 positive results.

The cult of $p < 0.05$

p-value = probability of positive result if it is actually wrong

Assumptions	Number of Tests
A non-Rick is not necessarily a Morty	<p>For 100 wrong claims, you get 5 positive results.</p> <p>Multiple hypothesis testing / FDR control comes at a price.</p>

The cult of $p < 0.05$

p-value = probability of positive result if it is actually wrong

Assumptions	Number of Tests	Effect Size
A non-Rick is not necessarily a Morty	<p>For 100 wrong claims, you get 5 positive results.</p> <p>Multiple hypothesis testing / FDR control comes at a price.</p>	

The cult of $p < 0.05$

p-value = probability of positive result if it is actually wrong

Assumptions	Number of Tests	Effect Size
A non-Rick is not necessarily a Morty	For 100 wrong claims, you get 5 positive results. Multiple hypothesis testing / FDR control comes at a price.	0.0000000001 > 0, so?

Explaining data variability
while acknowledging the effect of **random**
fluctuation

BELIEF

ATOM

[Wasserstein, Schirm, Lazar, 2019]

ATOM

[Wasserstein, Schirm, Lazar, 2019]

- Be **Thoughtful** about your belief.

ATOM

[Wasserstein, Schirm, Lazar, 2019]

- Be **Thoughtful** about your **belief**.
- Be **Open** about your **belief**.

ATOM

[Wasserstein, Schirm, Lazar, 2019]

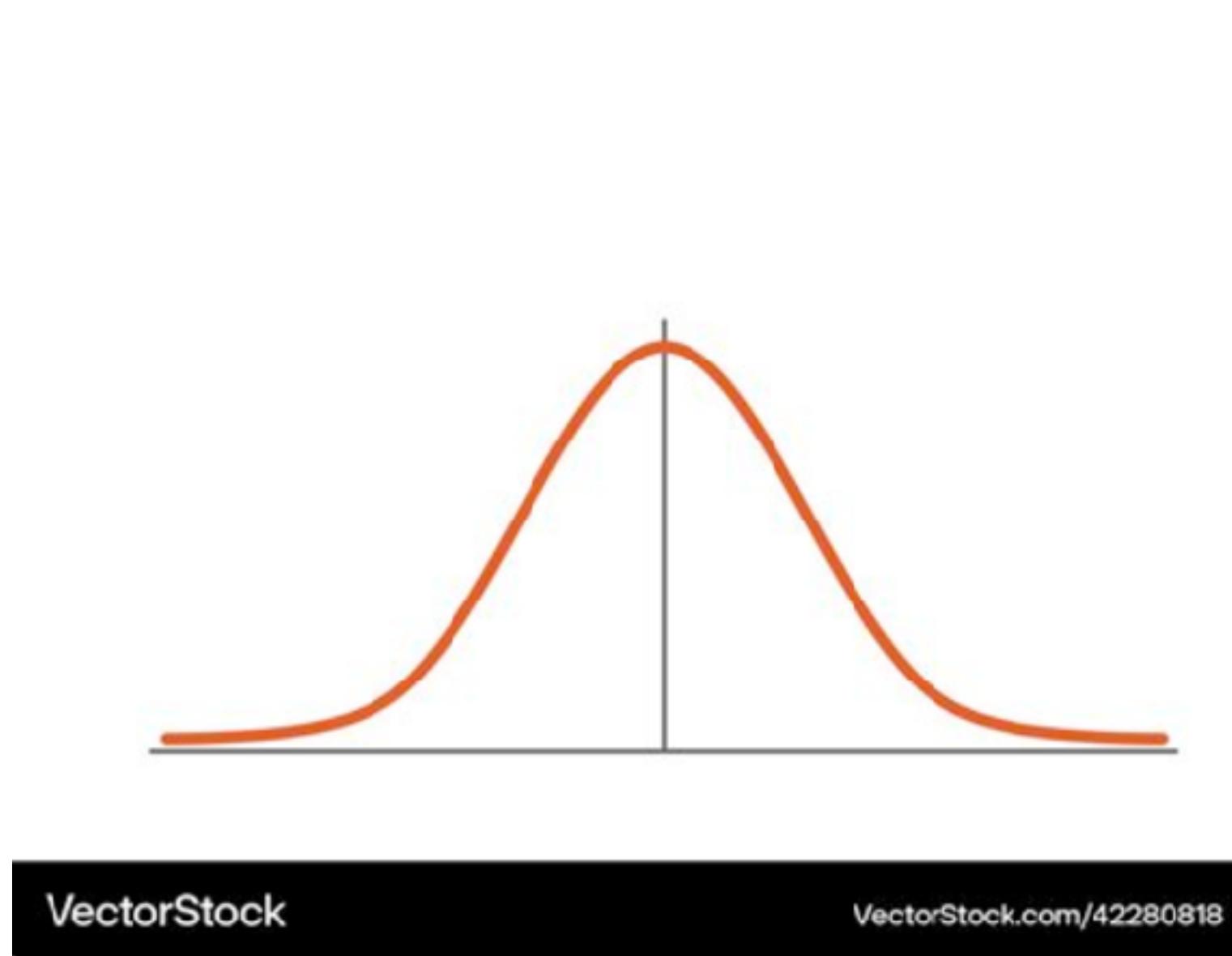
- Be **Thoughtful** about your belief.
- Be **Open** about your belief.
- Be **Modest** about your results.

ATOM

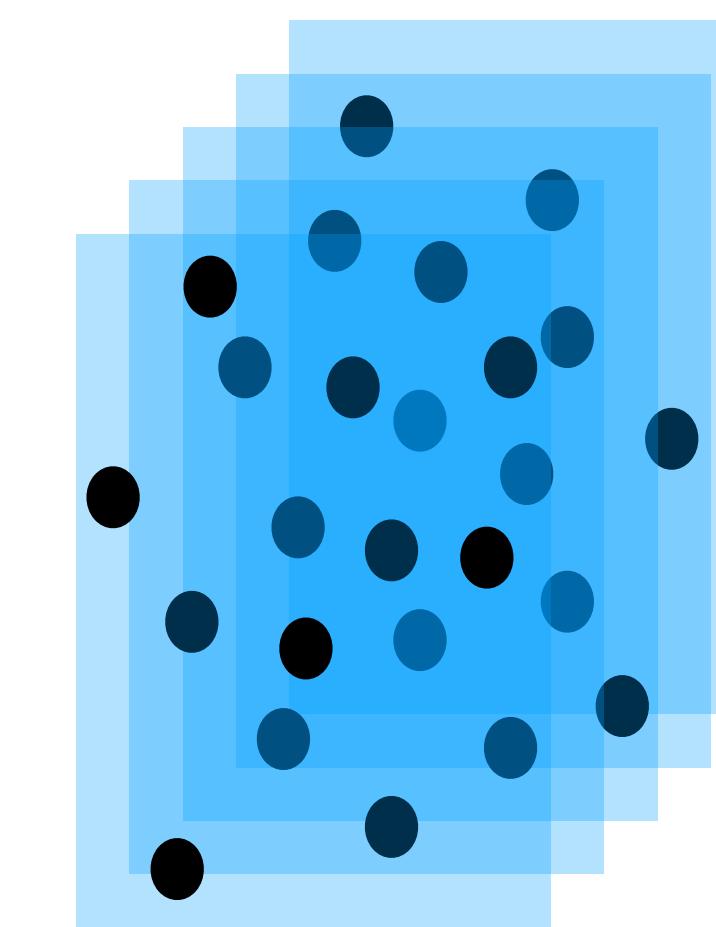
[Wasserstein, Schirm, Lazar, 2019]

- **Accept** randomness.
- Be **Thoughtful** about your belief.
- Be **Open** about your belief.
- Be **Modest** about your results.

How to Gauge Randomness

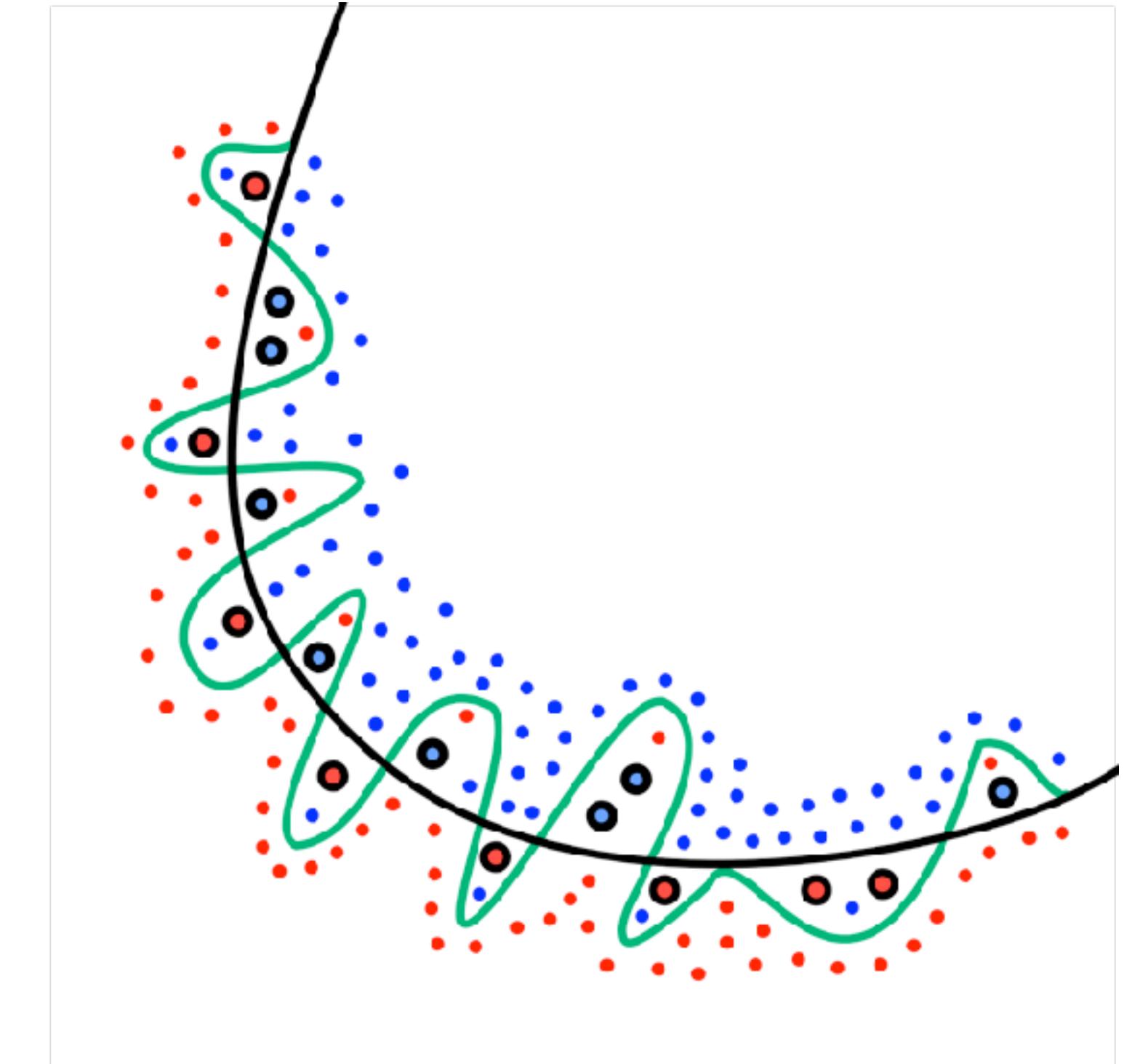


Normal Approximation



Bootstrap

Lionel Rowe https://commons.wikimedia.org/wiki/File:JavaScript_code.png

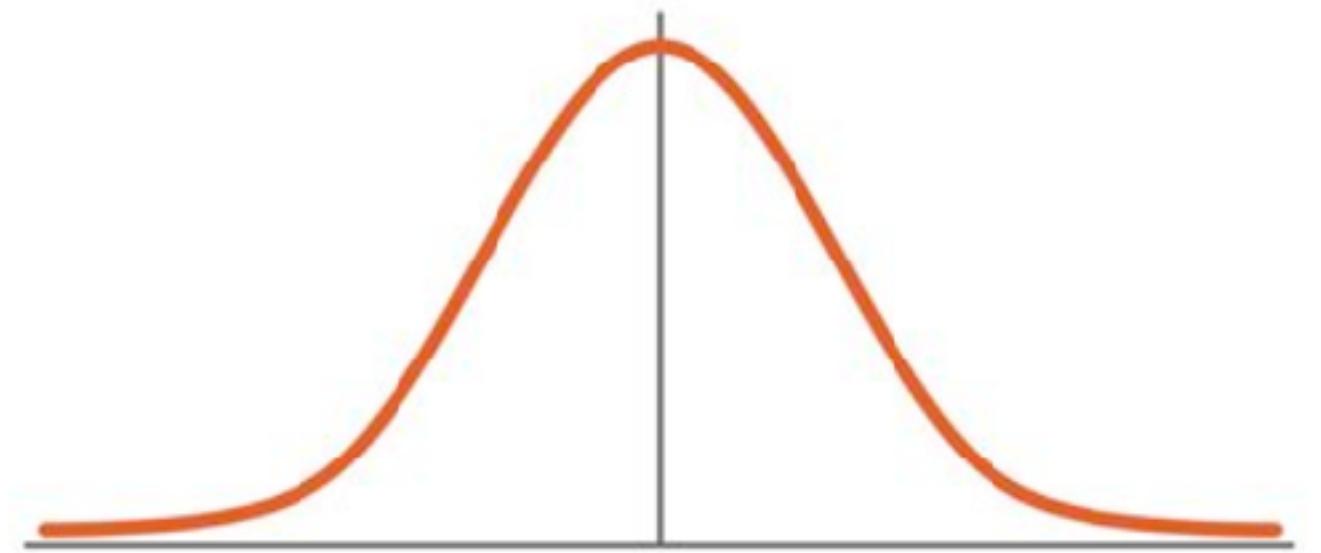


Cross Validation

Chabacano <https://commons.wikimedia.org/wiki/File:Overfitting.svg>

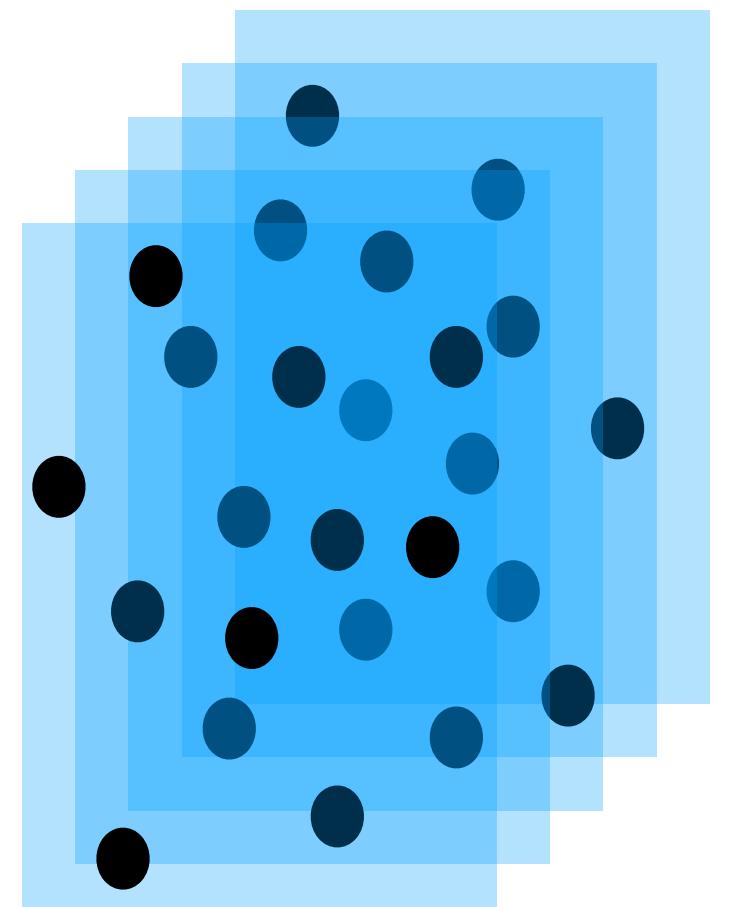
Normal Approximation

- when you compute mean and standard deviation and nothing else



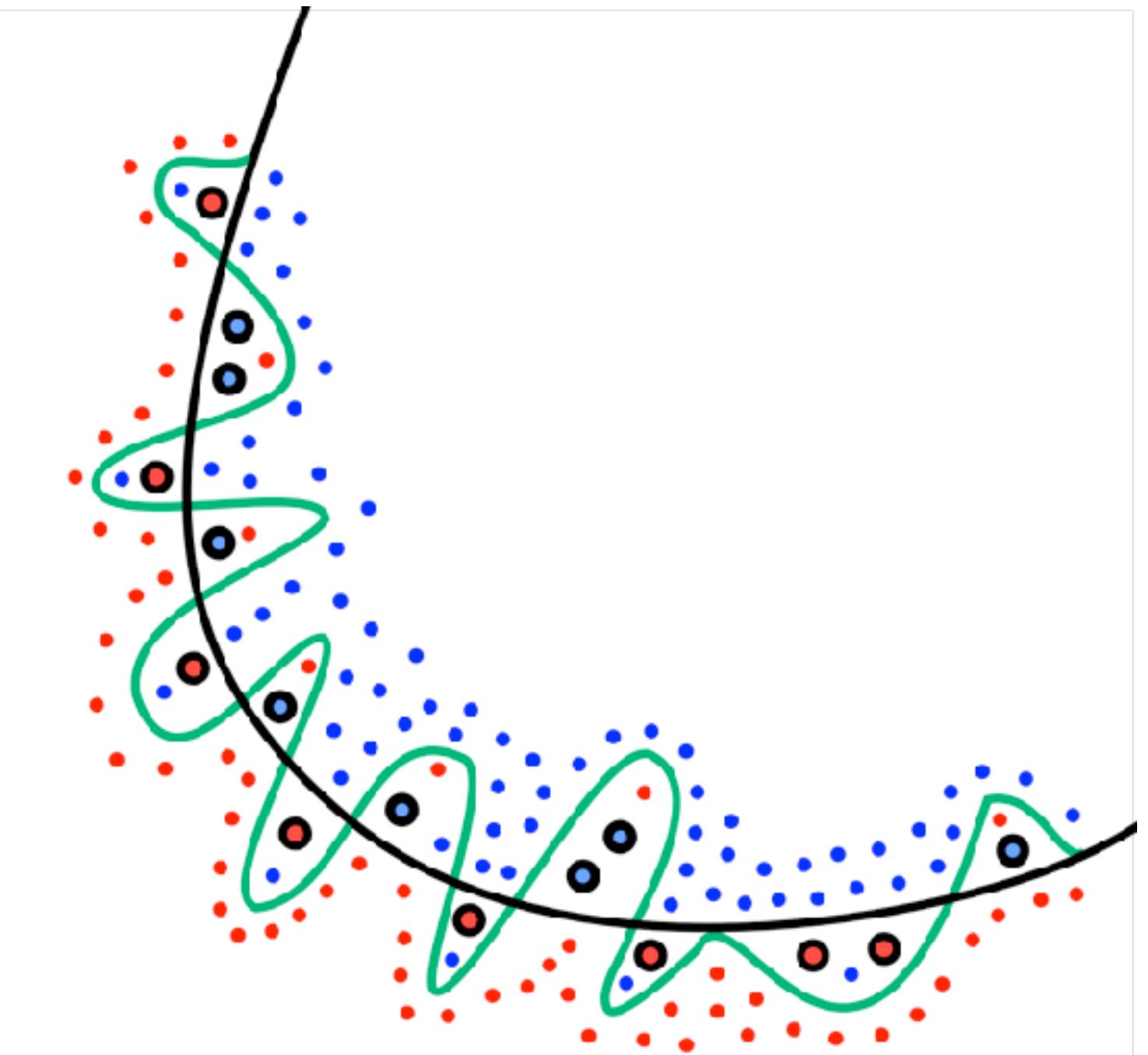
Bootstrap / Subsampling

- Bootstrap
 - Generate many **same-sized** samples by **resampling** from the data set
- Subsampling
 - Generate many small **subsamples**
 - Look at the resultant distribution of the sample means



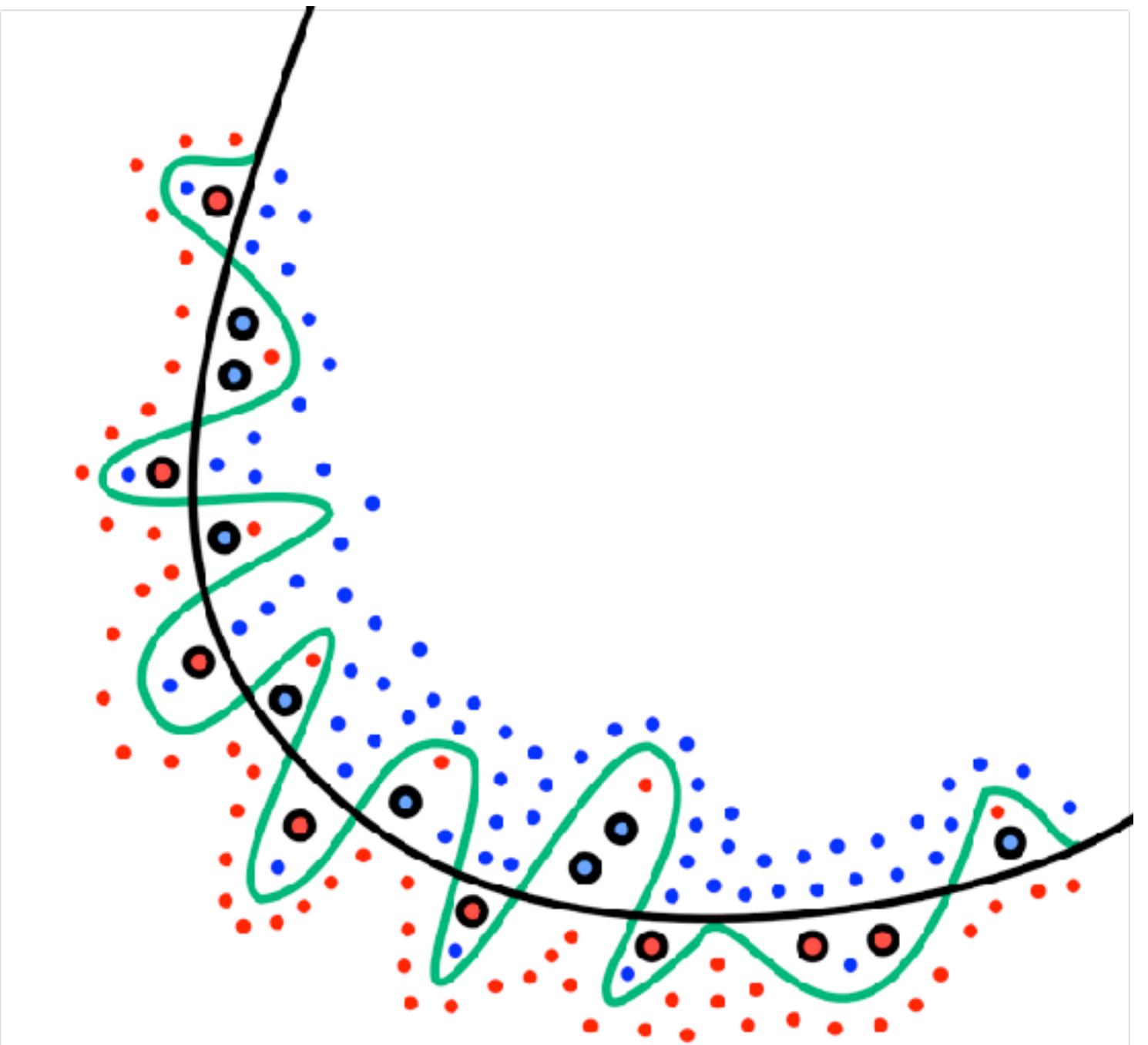
Cross-Validation

- Maybe this is neither Rick nor Morty

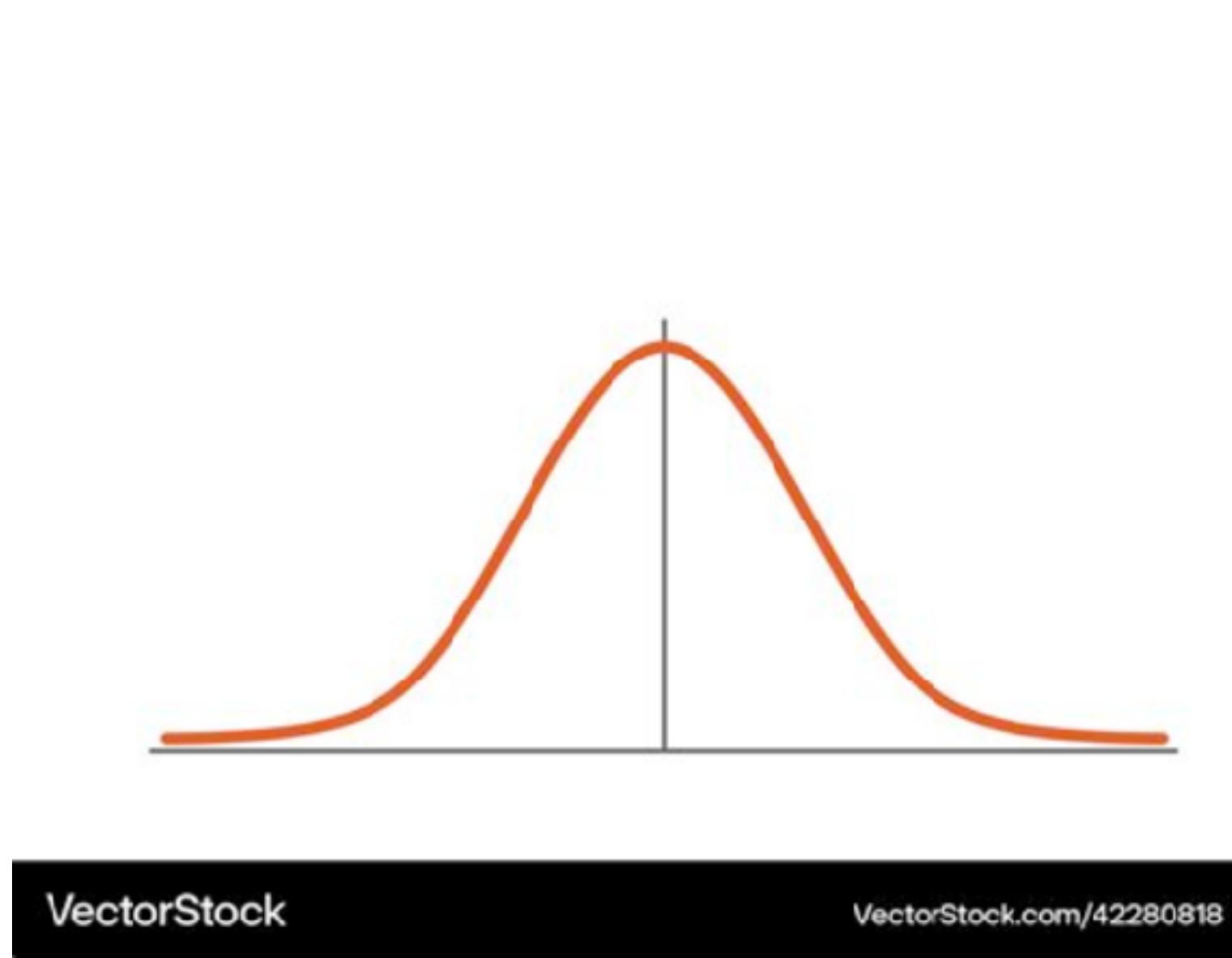


Cross-Validation

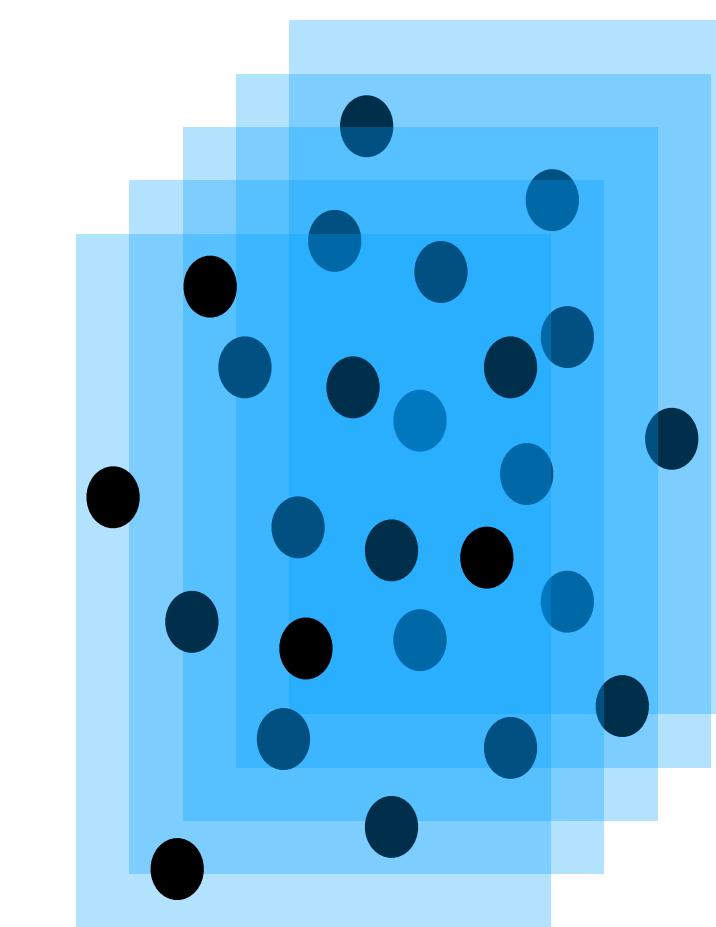
- Maybe this is neither Rick nor Morty
- But all models are wrong...



How to Gauge Randomness

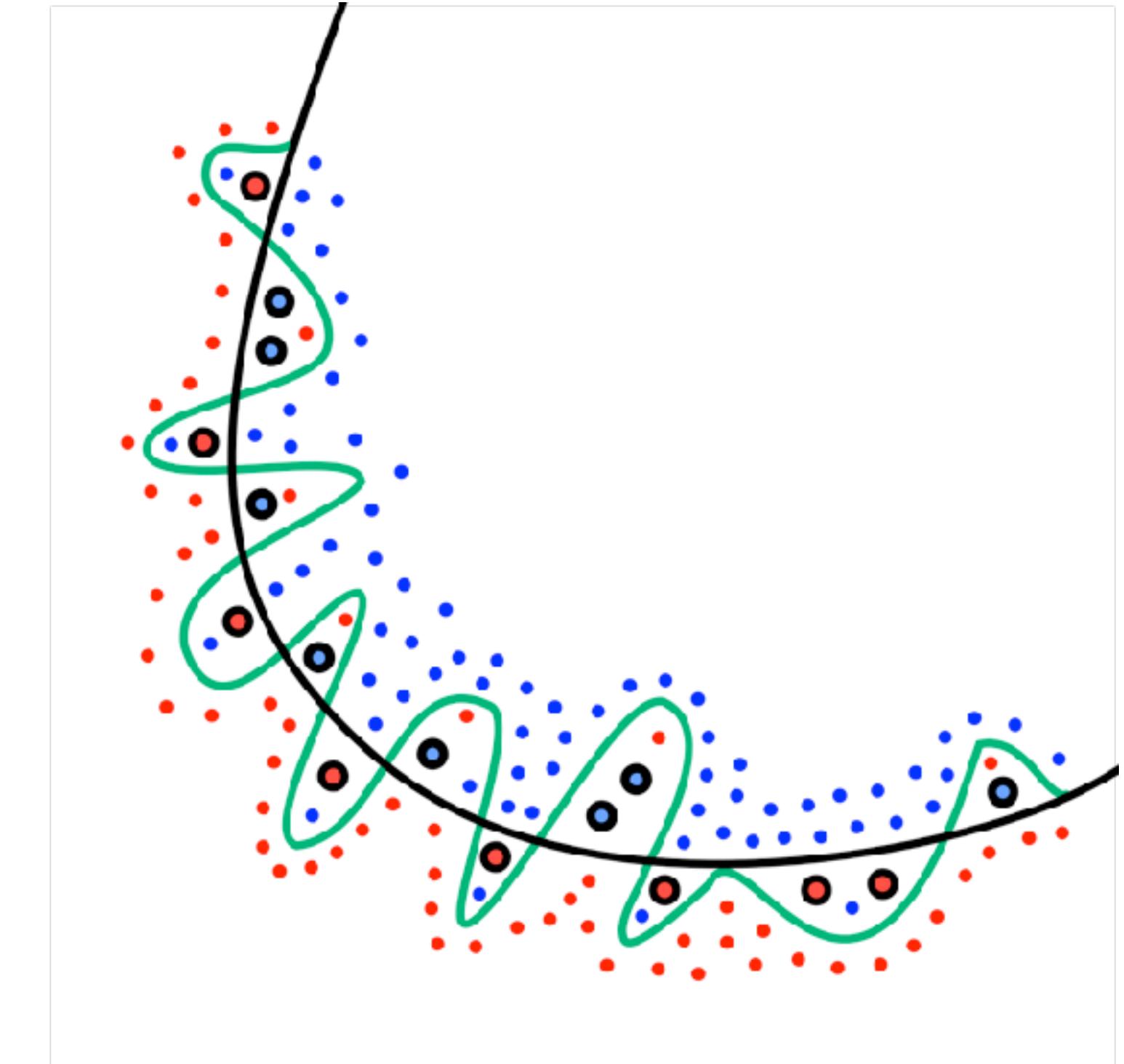


Normal Approximation



Bootstrap

Lionel Rowe https://commons.wikimedia.org/wiki/File:JavaScript_code.png



Cross Validation

Chabacano <https://commons.wikimedia.org/wiki/File:Overfitting.svg>

Dimension Reduction

Dimension Reduction

Dimension Reduction



Walt Disney Pictures / Barry Wetcher

Enchanted, 2007

Dimension Reduction

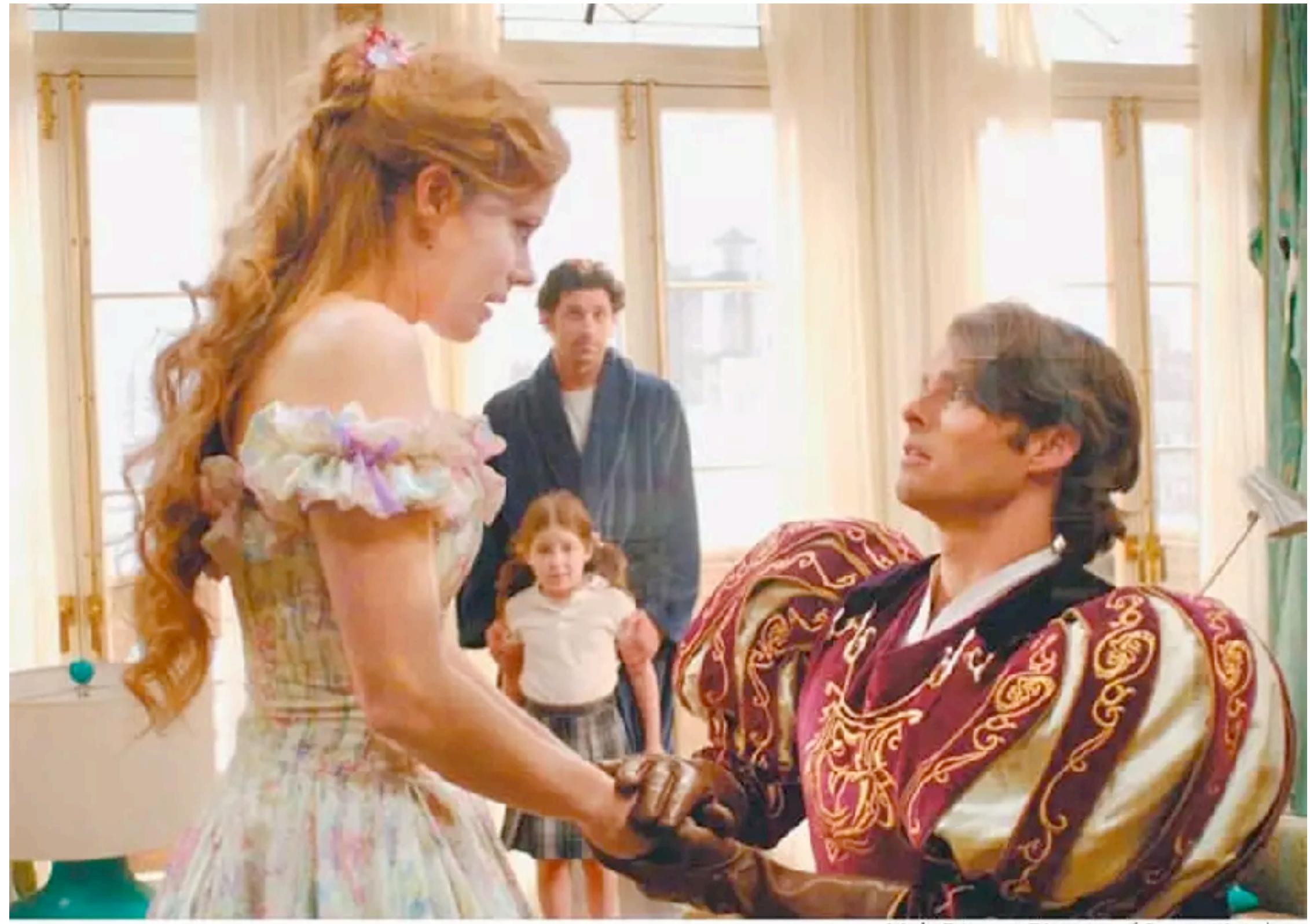


Walt Disney Pictures / Barry Wetcher



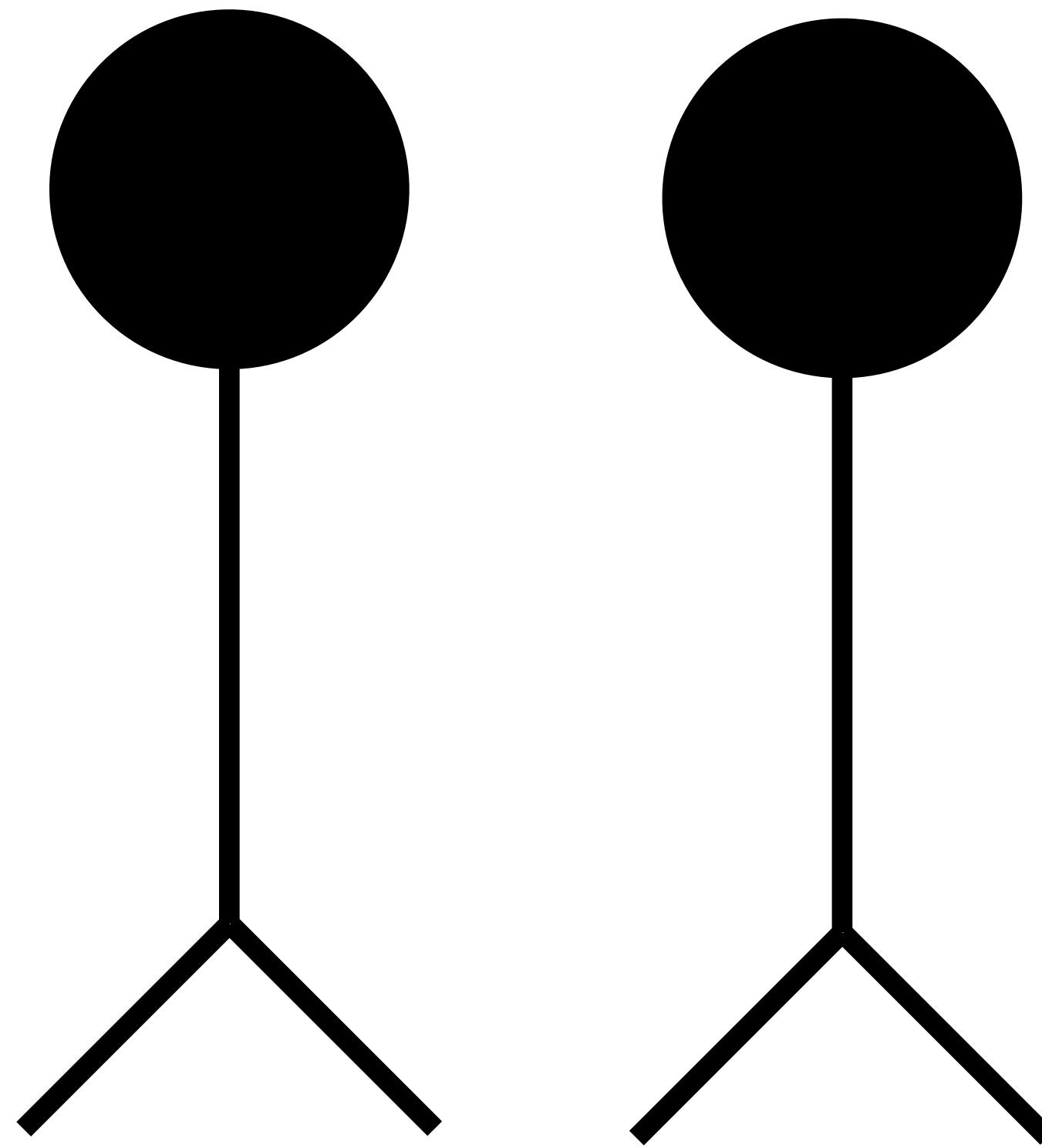
Enchanted, 2007

Input, Output



Walt Disney Pictures / Barry Wetcher

Enchanted, 2007



Different Methods

Principal Component
Analysis

Multi-Dimensional Scaling

UMAP

Mapper

Different Methods

Principal Component
Analysis

Preserve linear information

Multi-Dimensional Scaling

UMAP

Mapper

Different Methods

Principal Component
Analysis

Preserve linear information

Multi-Dimensional Scaling

Preserve distances as far as possible

UMAP

Mapper

Different Methods

Principal Component
Analysis

Preserve linear information

Multi-Dimensional Scaling

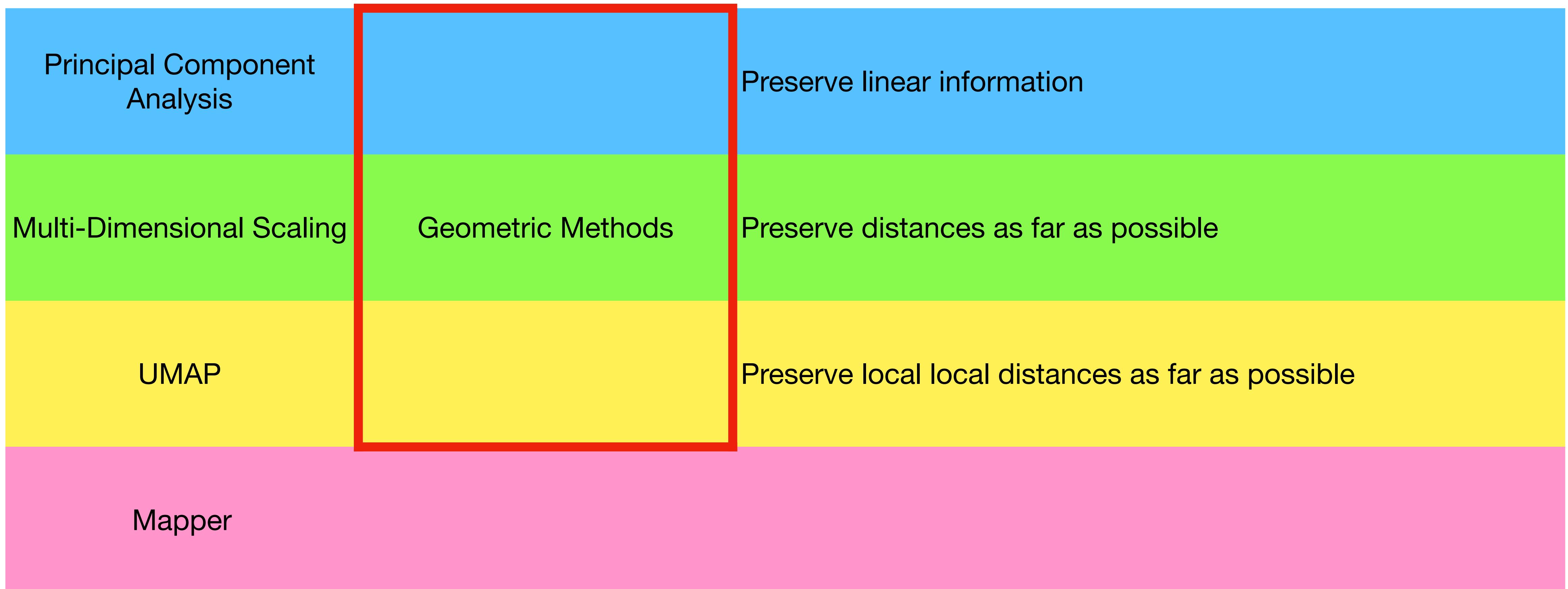
Preserve distances as far as possible

UMAP

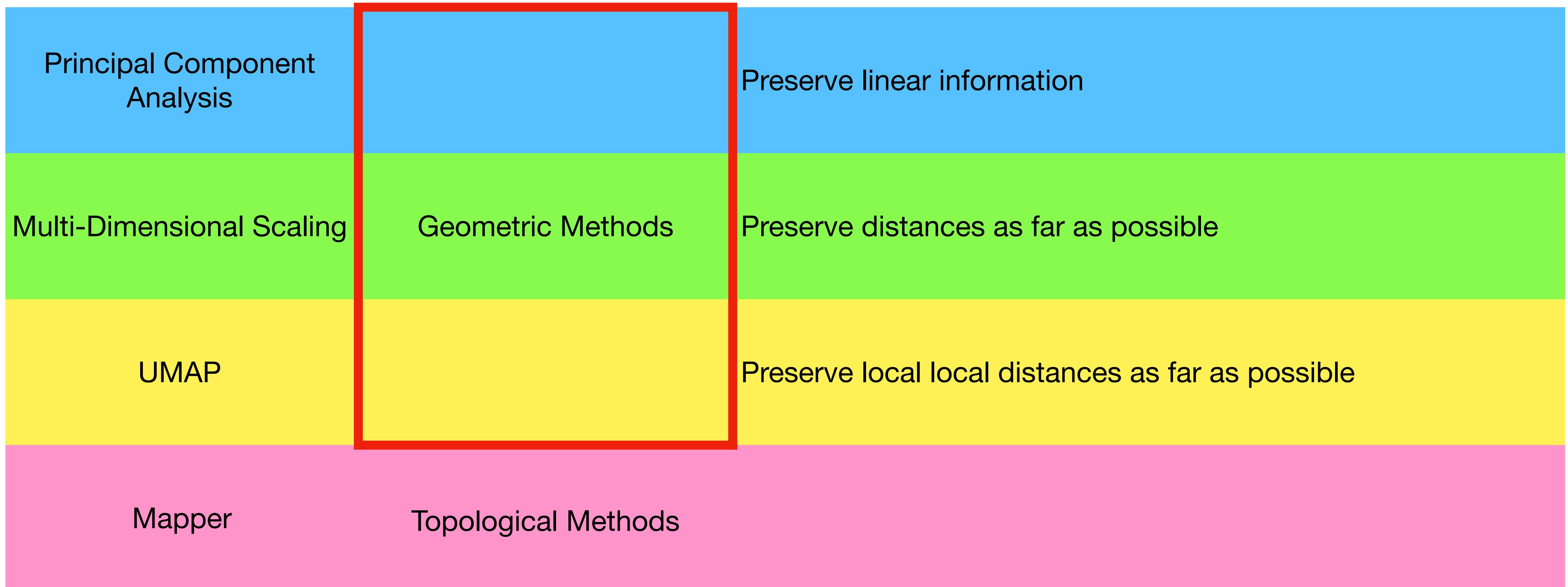
Preserve local local distances as far as possible

Mapper

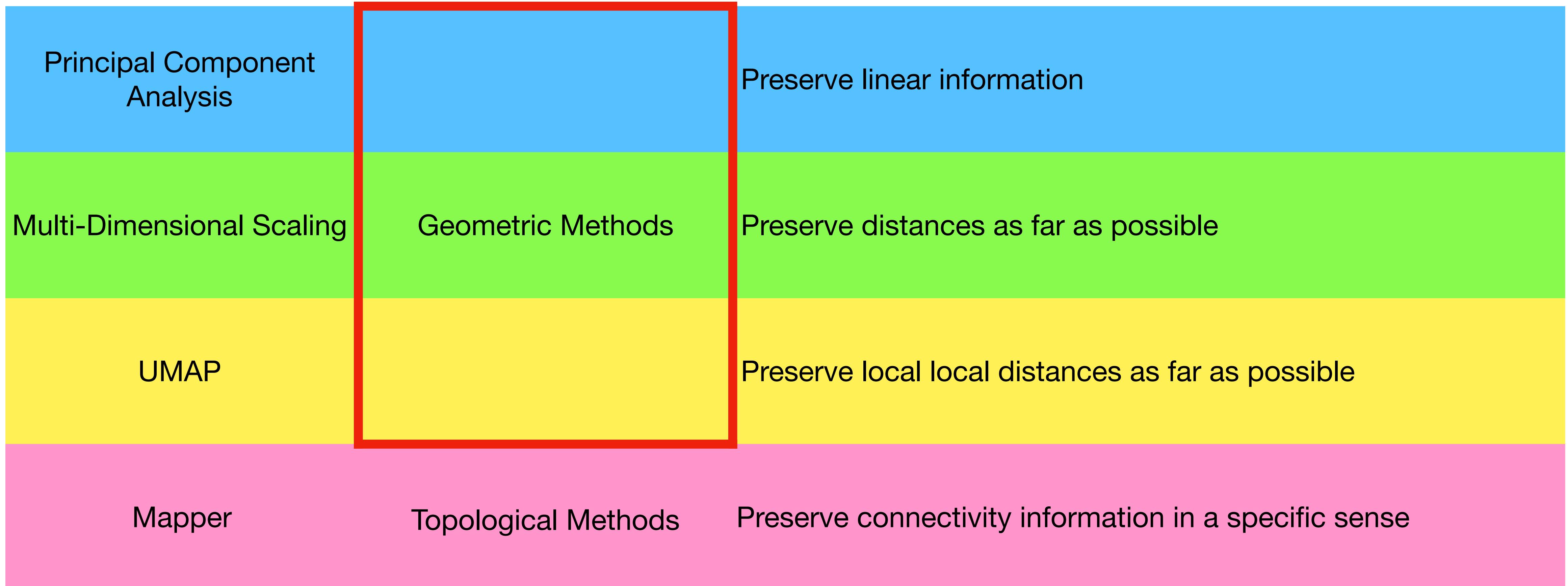
Different Methods



Different Methods



Different Methods



As good as possible...

As good as possible...

- Nearby points in high-dim remains nearby in low-dim.
- Reverse?

As good as possible...

- Nearby points in high-dim remains nearby in low-dim.
- Reverse?



Cinderella III: A Twist in Time (2017)

As good as possible...

- Nearby points in high-dim remains nearby in low-dim.
- Reverse needs cheating



Alice in the Wonderland (1951)

Mapper

- Nearby points in high-dim remains nearby in low-dim.
- Reverse guaranteed by using graph topology to transcend Euclidean geometry



Alice in the Wonderland (1951)

Mapper

- Nearby points in high-dim remains nearby in low-dim.
- Reverse guaranteed by using graph topology to transcend Euclidean geometry
- Price: loss of metric control



Alice in the Wonderland (1951)

Mapper

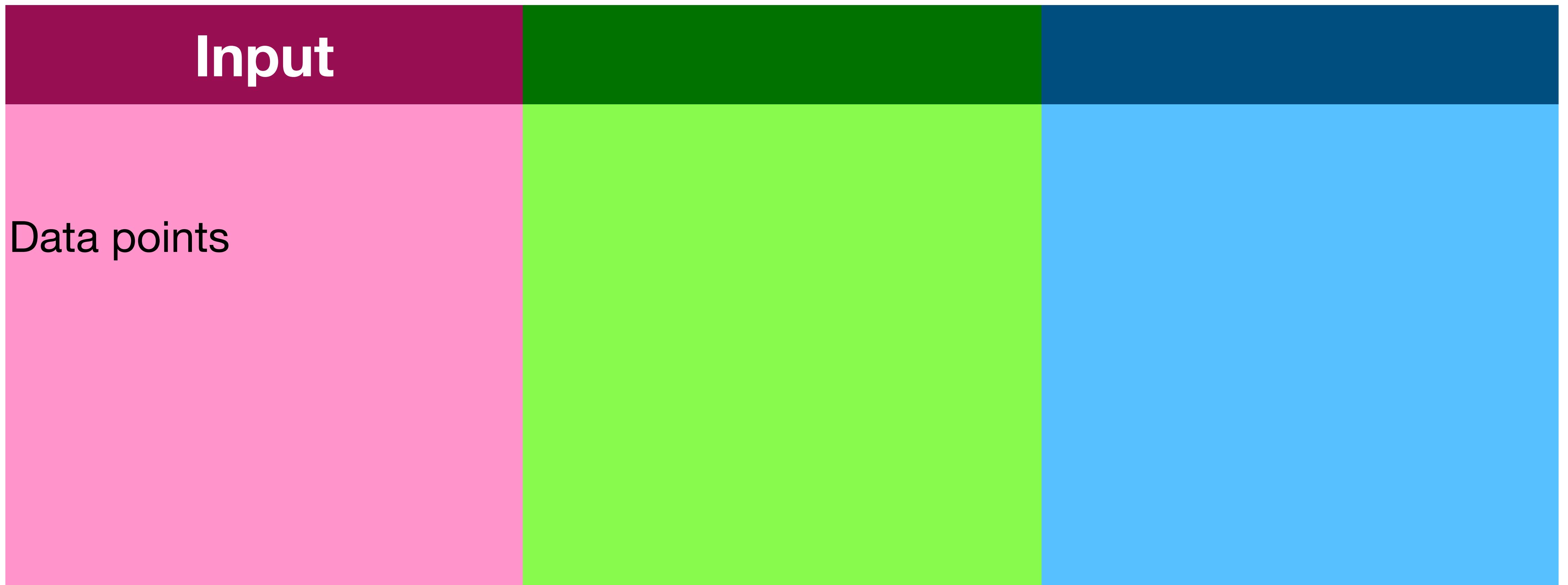
- Nearby points in high-dim remains nearby in low-dim.
- Reverse guaranteed by using graph topology to transcend Euclidean geometry
- Price: loss of metric control
- Get to keep: connectivity



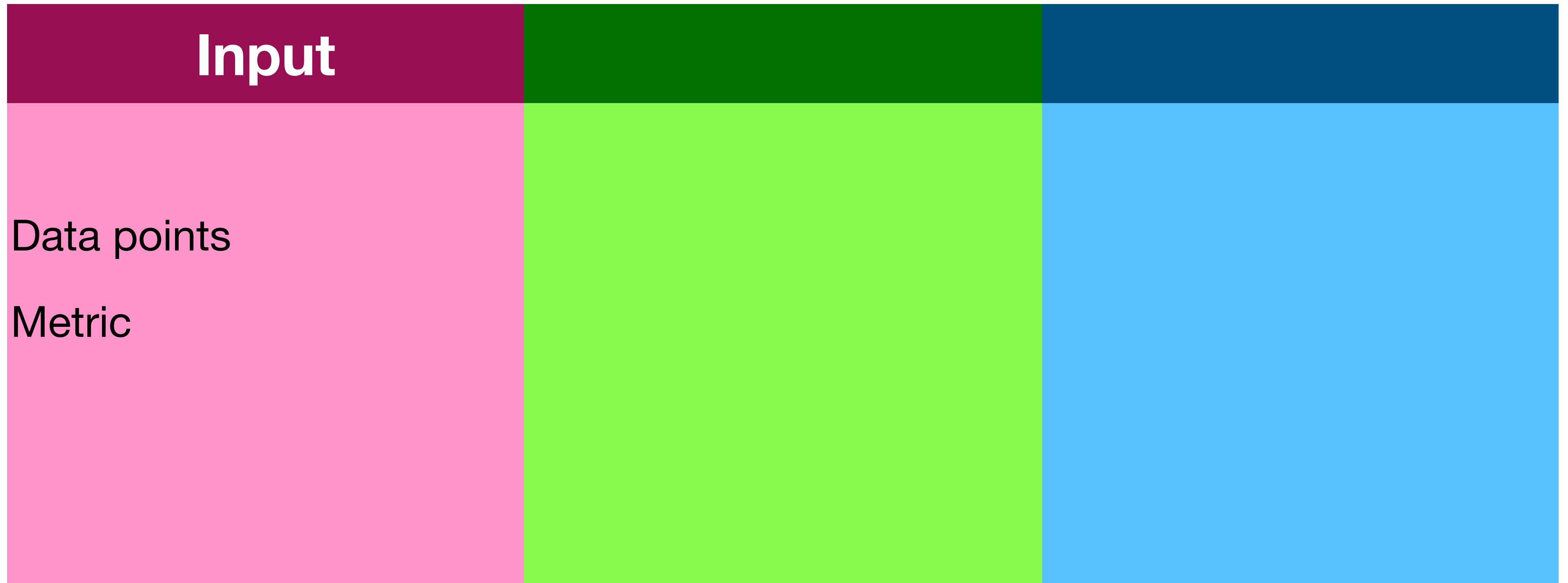
Alice in the Wonderland (1951)

What is Mapper?

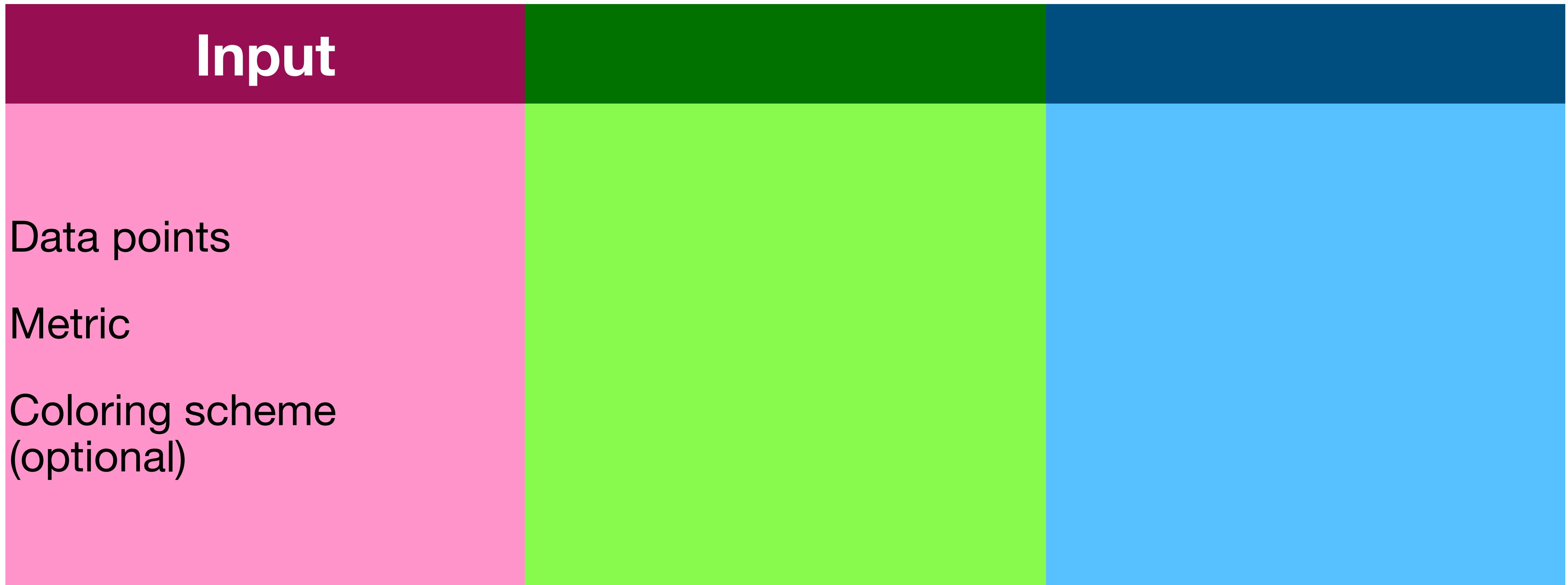
What is Mapper?



What is Mapper?



What is Mapper?



What is Mapper?

Input		Output
<p>Data points</p> <p>Metric</p> <p>Coloring scheme (optional)</p>		<p>A (colored) graph, where</p>

What is Mapper?

Input		Output
Data points Metric Coloring scheme (optional)		A (colored) graph, where nodes are clusters of data points

What is Mapper?

Input		Output
Data points Metric Coloring scheme (optional)		A (colored) graph, where nodes are clusters of data points edges are nonempty intersection of clusters

What is Mapper?

Input	Parameters	Output
Data points	Discretization parameters	A (colored) graph, where nodes are clusters of data points
Metric		
Coloring scheme (optional)		edges are nonempty intersection of clusters

What is Mapper?

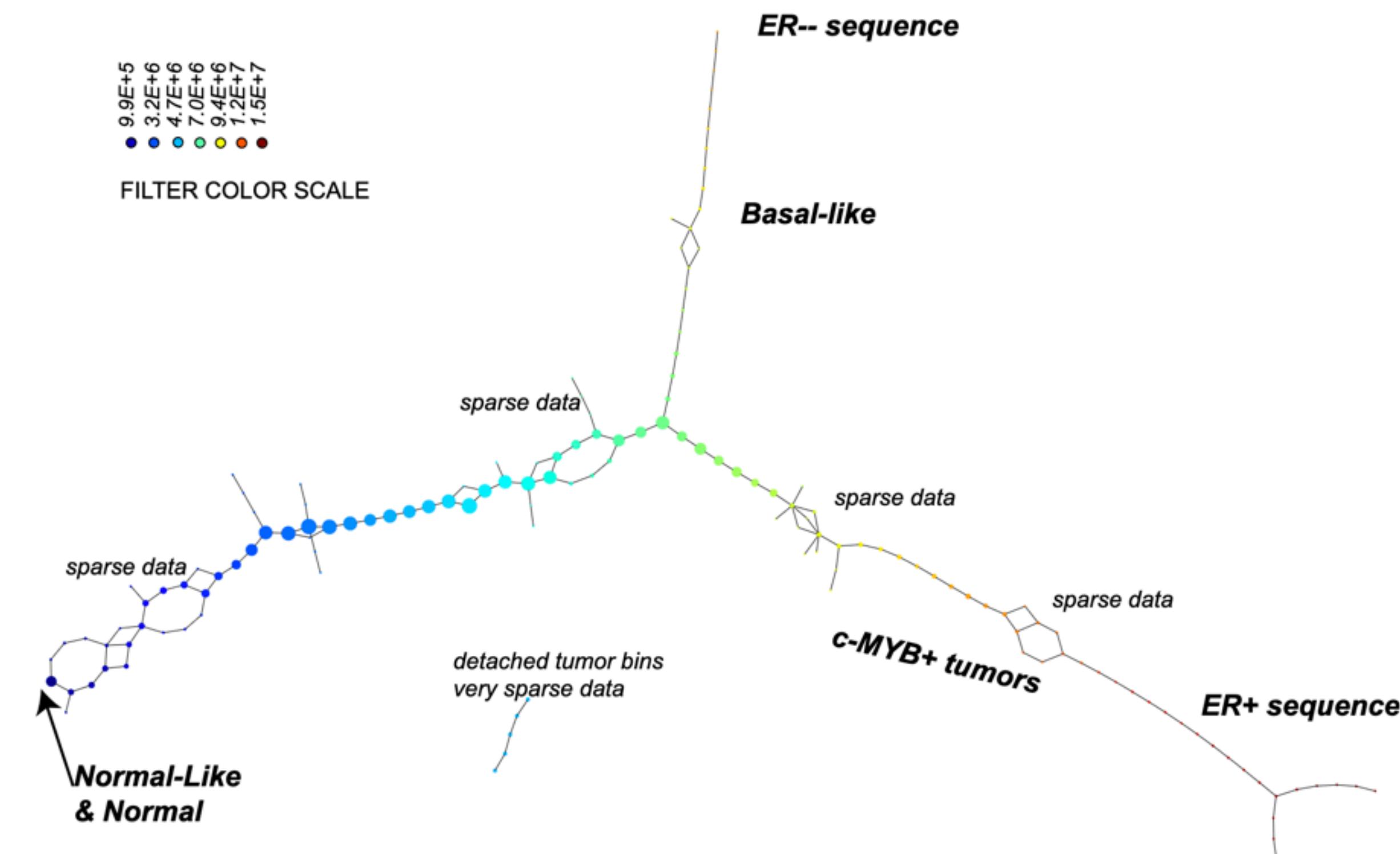
Input	Parameters	Output
Data points	Discretization parameters	A (colored) graph, where nodes are clusters of data points
Metric	Lens (if lens-based)	edges are nonempty intersection of clusters
Coloring scheme (optional)		

What is Mapper?

Input	Parameters	Output
Data points	Discretization parameters	A (colored) graph, where nodes are clusters of data points
Metric	Lens (if lens-based)	
Coloring scheme (optional)	Subroutine parameters	edges are nonempty intersection of clusters

Nicolau et al, 2011

- data: Transcriptional microarray data
- sample: 295 tumors
- features: 262 genes
- lens: Normal component to the linear subspace of healthy tissues



What is Mapper Good For?

proximity guarantee, and...?

What is Mapper Good For?

proximity guarantee, and...?

- Distinguish points that are indistinguishable through the lens

What is Mapper Good For?

proximity guarantee, and...?

- Distinguish points that are indistinguishable through the lens
- Good for sub-typing

What is Mapper Good For?

proximity guarantee, and...?

- Distinguish points that are indistinguishable through the lens
- Good for sub-typing
- Finding out relationship between them

Limitations of Mapper

Limitations of Mapper

- cannot retain metric information

Limitations of Mapper

- cannot retain metric information
- choices of lens

Limitations of Mapper

- cannot retain metric information
- choices of lens
- parameter tuning

Limitations of Mapper

- cannot retain metric information
- choices of lens
- parameter tuning
- difficulty of clustering

Intro to **Topological Data Analysis**

Statistics and Dimension Reduction

Chunyin Siu, Feb 14, 2025