# **AIBridge**

Lecture 5

# Let's talk about the last lab!

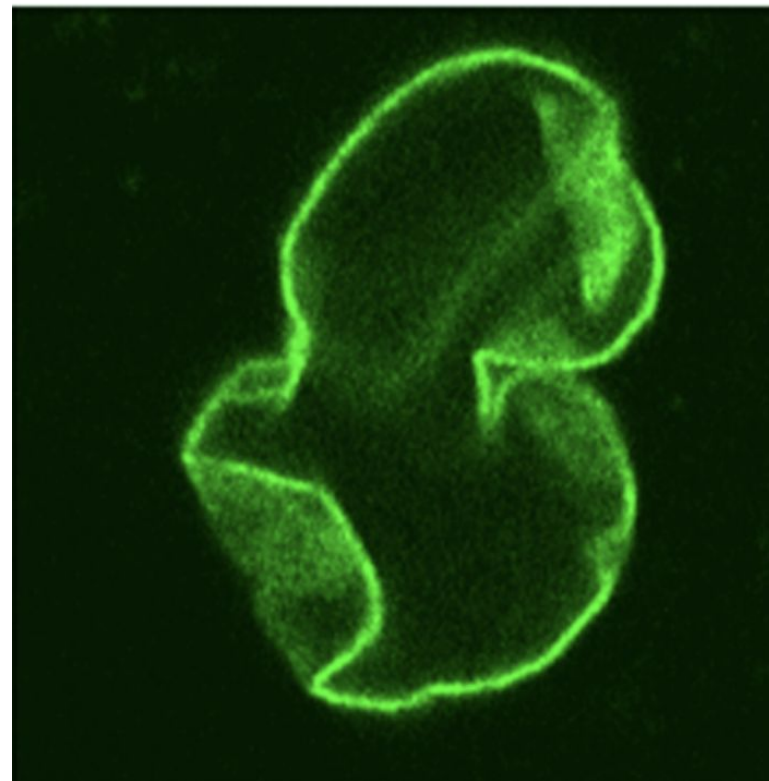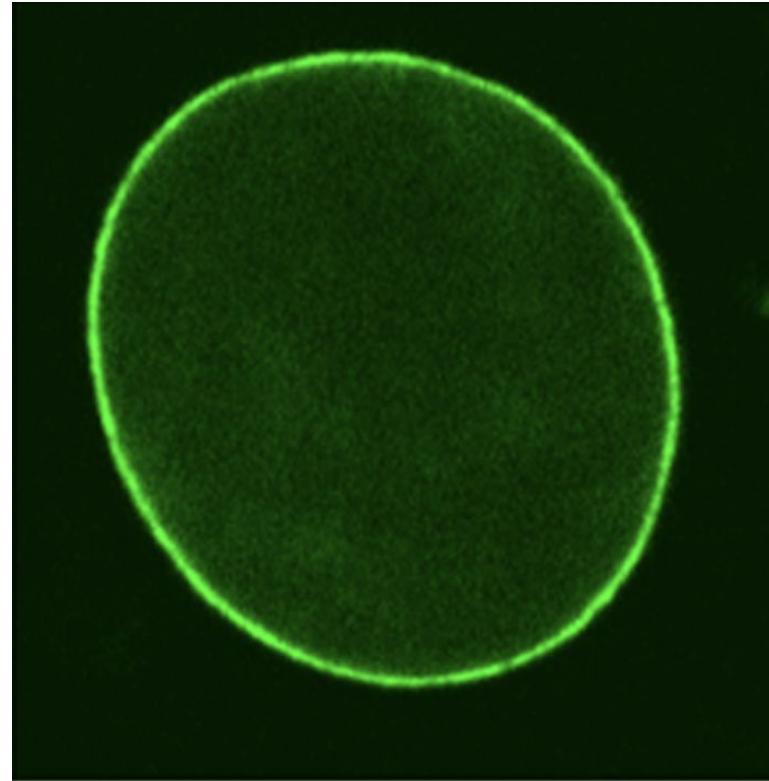# Let's talk about the last lab!

what does this even mean?

# What circumstances made the model fit better?
## worse?

# Accuracy

## "Why is it not enough?"

**Progeria affects ~159 patients in the US**

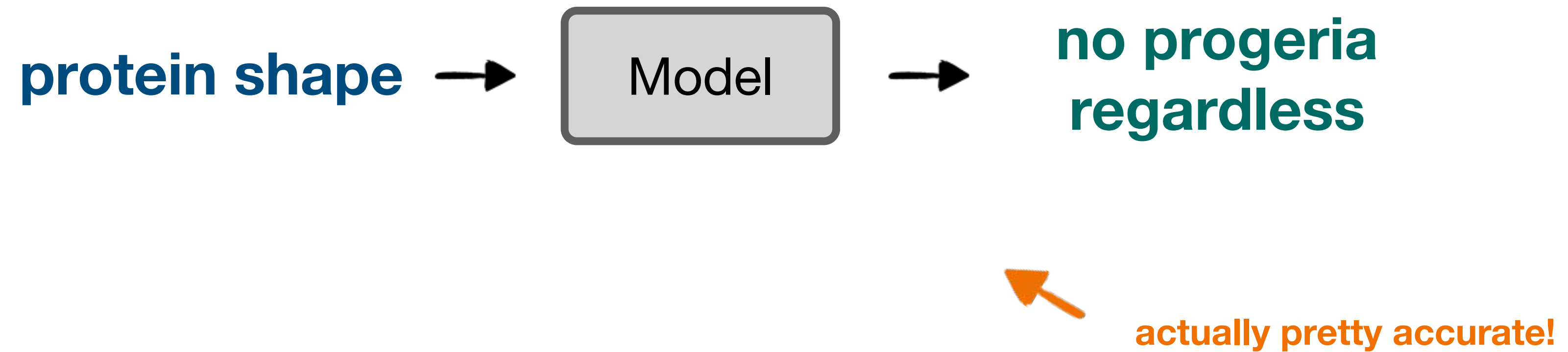we have a dataset of all American pediatric patients

Q: If my model predicts with 99.99% accuracy, is it good enough?

**protein shape** → Model → **progeria**
**(yes or no)**

**Progeria affects ~159 patients in the US**

we have a dataset of all American pediatric patients

**a proposed model:**

**protein shape** → Model → **no progeria regardless**
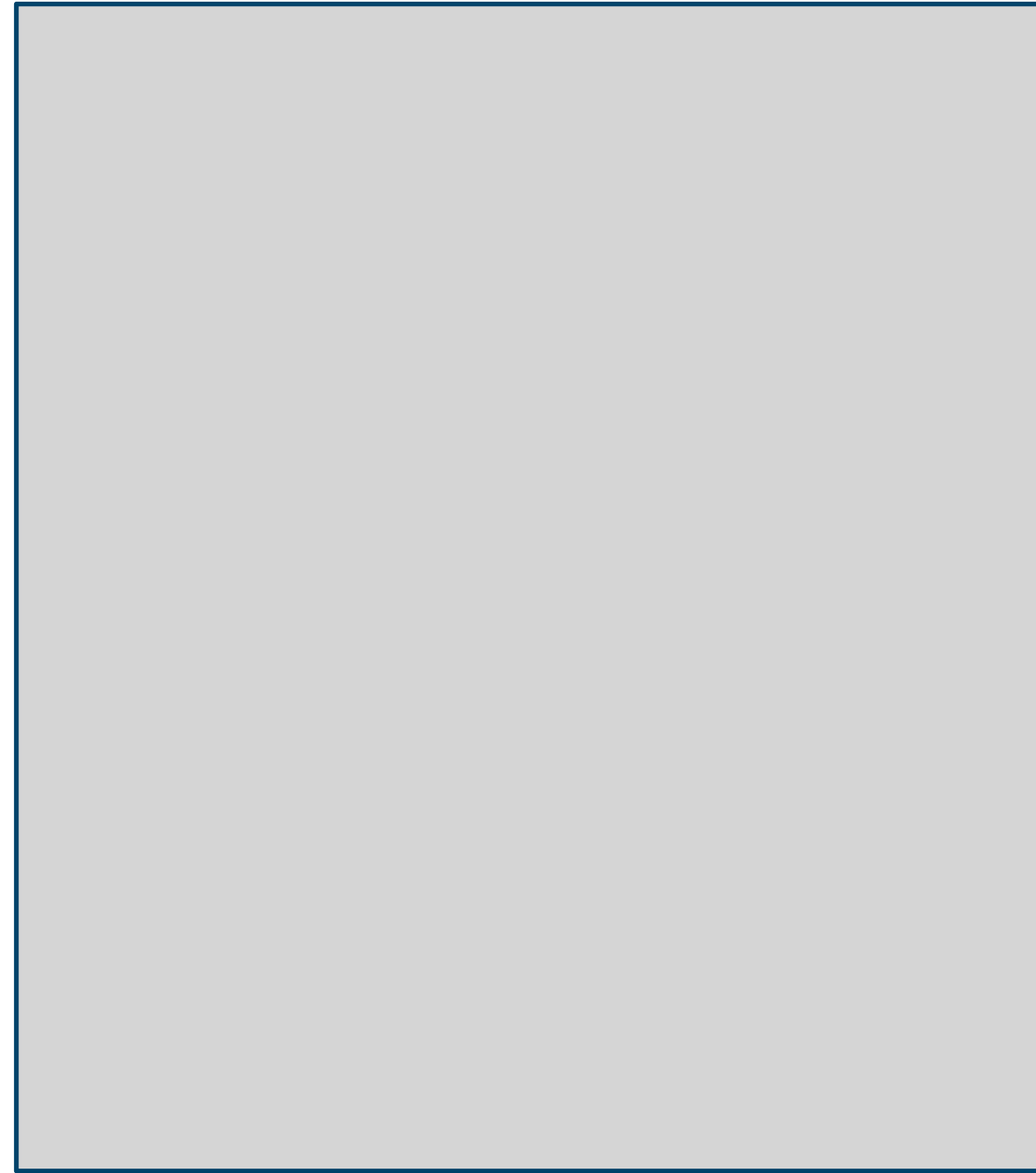
**actually pretty accurate!**

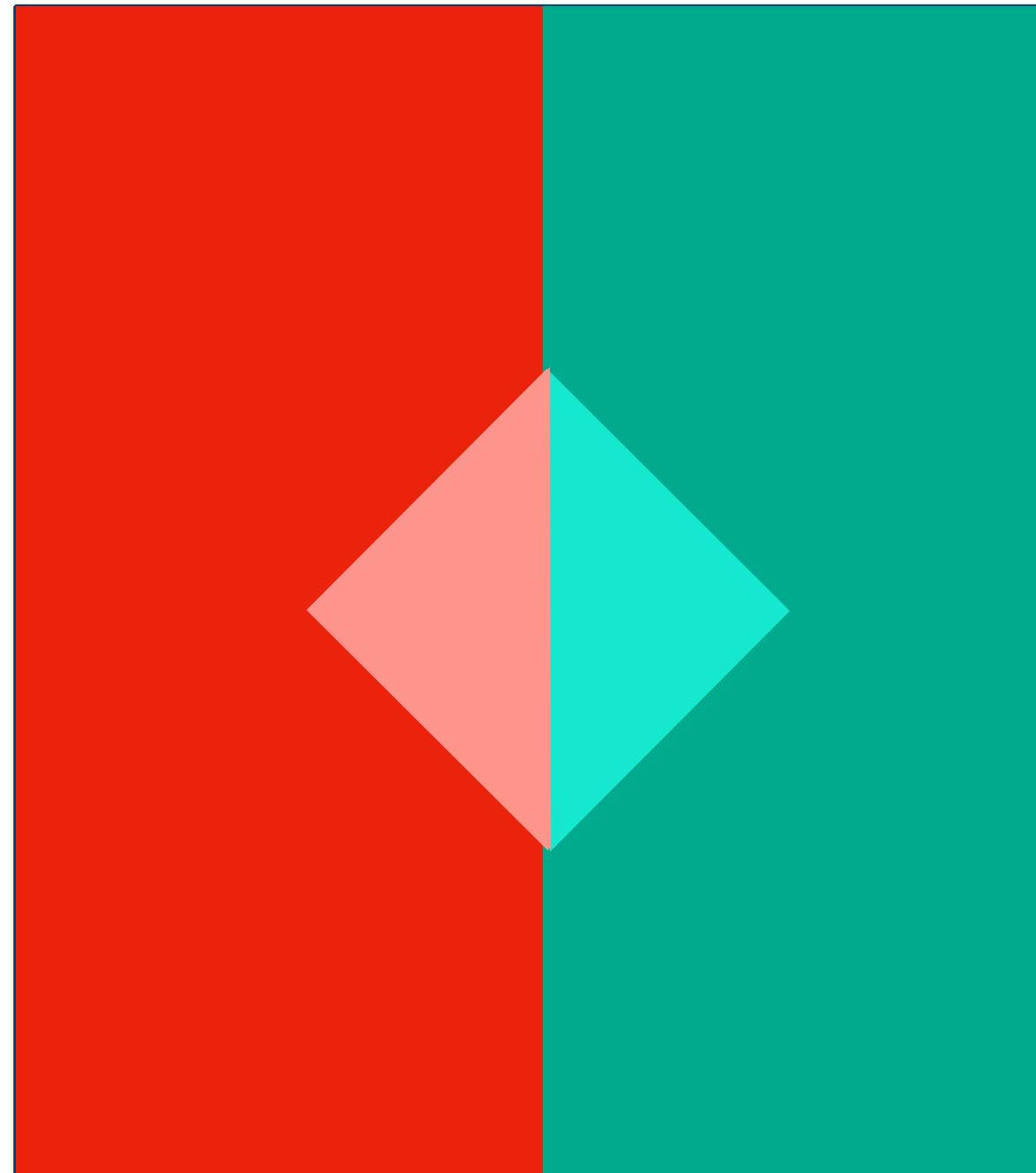**Progeria affects ~159 patients in the US**

we have a dataset of all American pediatric patients

# Accuracy , Precision, and Recall

"Selection space"
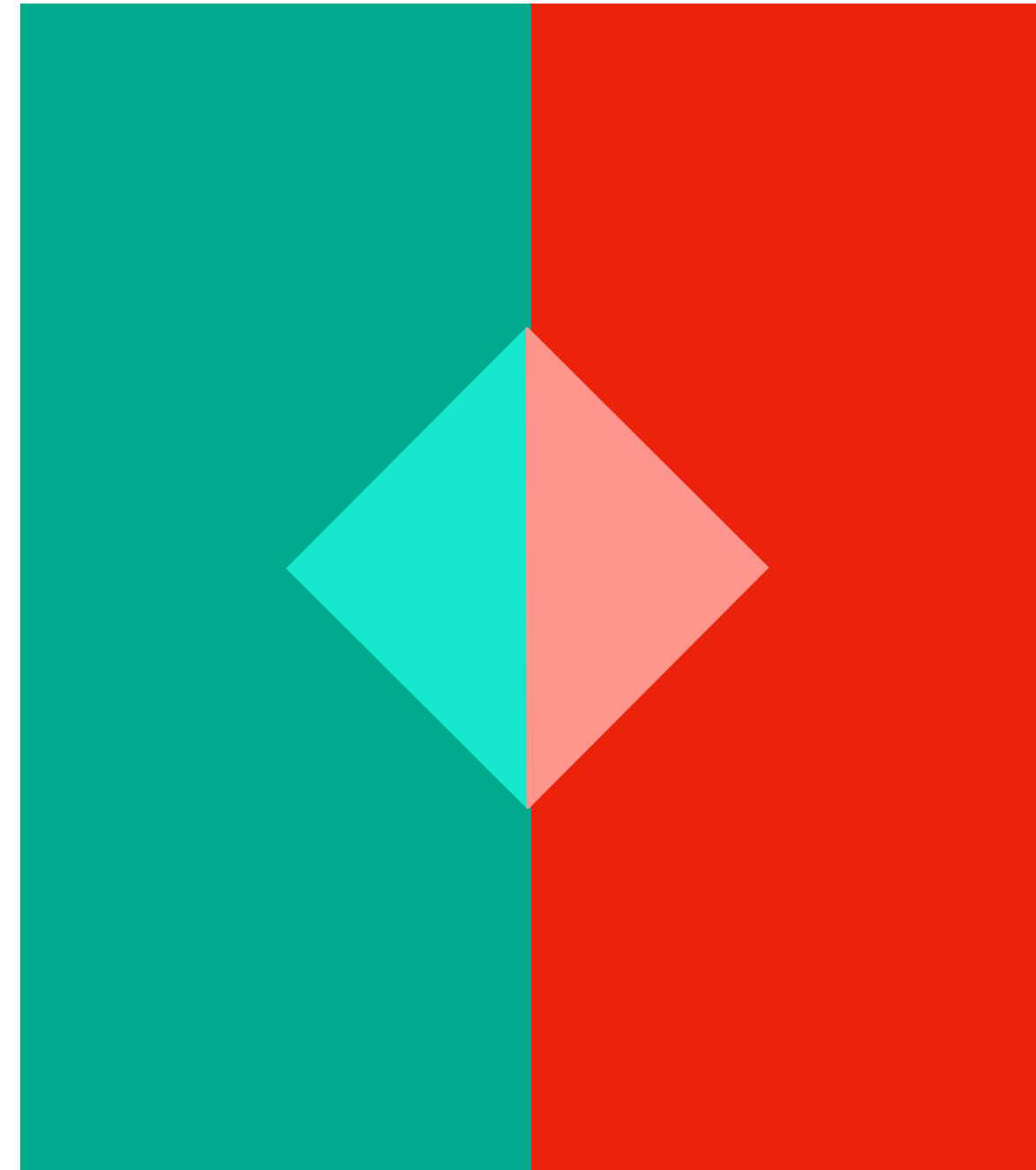
"Selection space"

Model selects **positive** and patient is **positive**

"Selection space"

1
1

Model selects **positive** and patient is **positive**

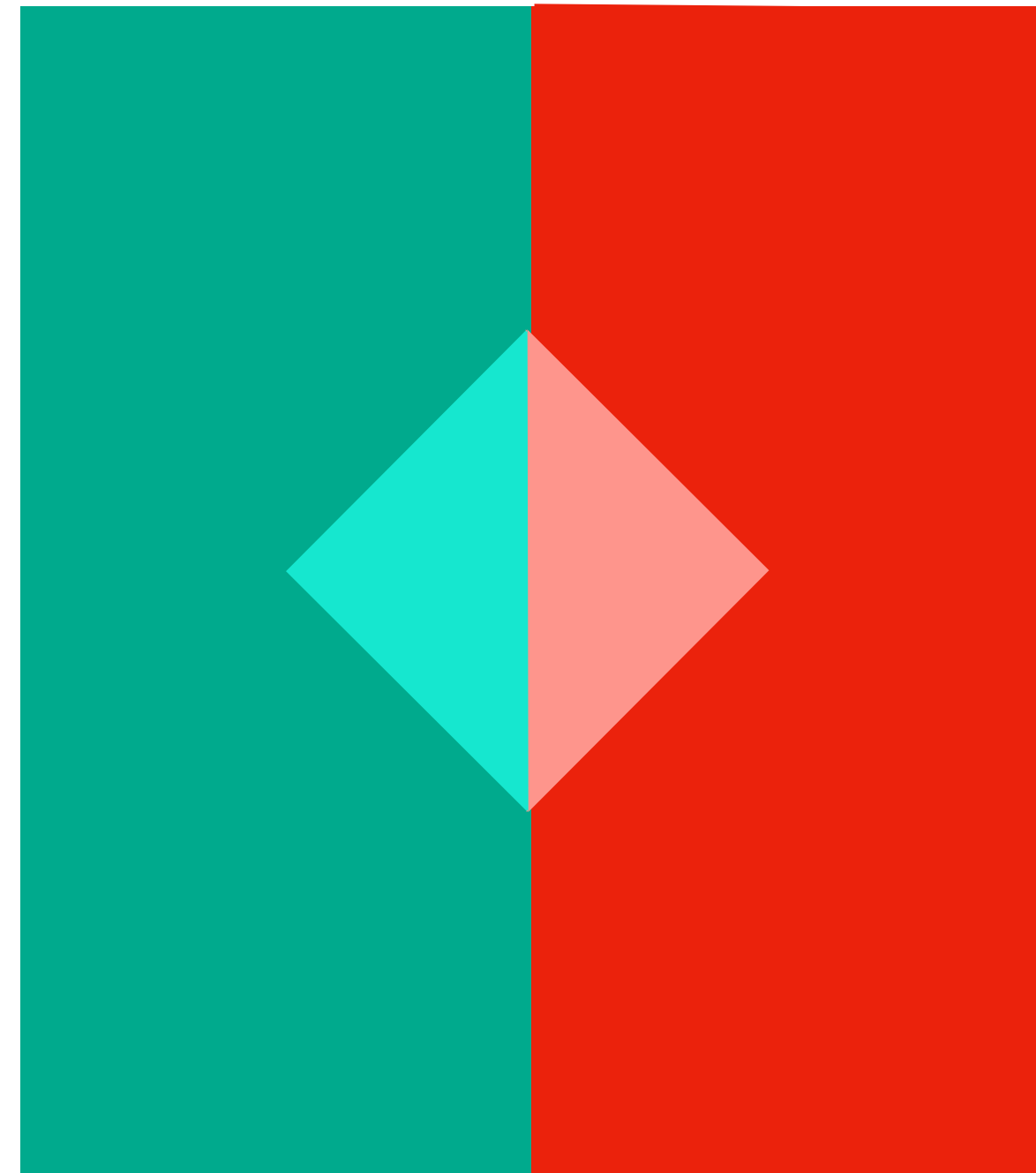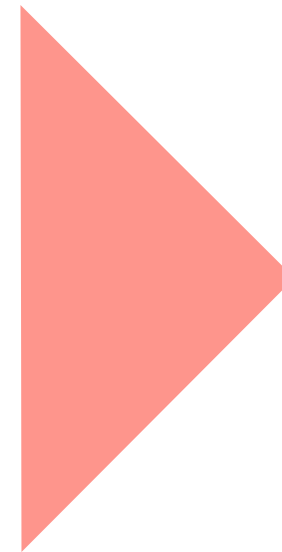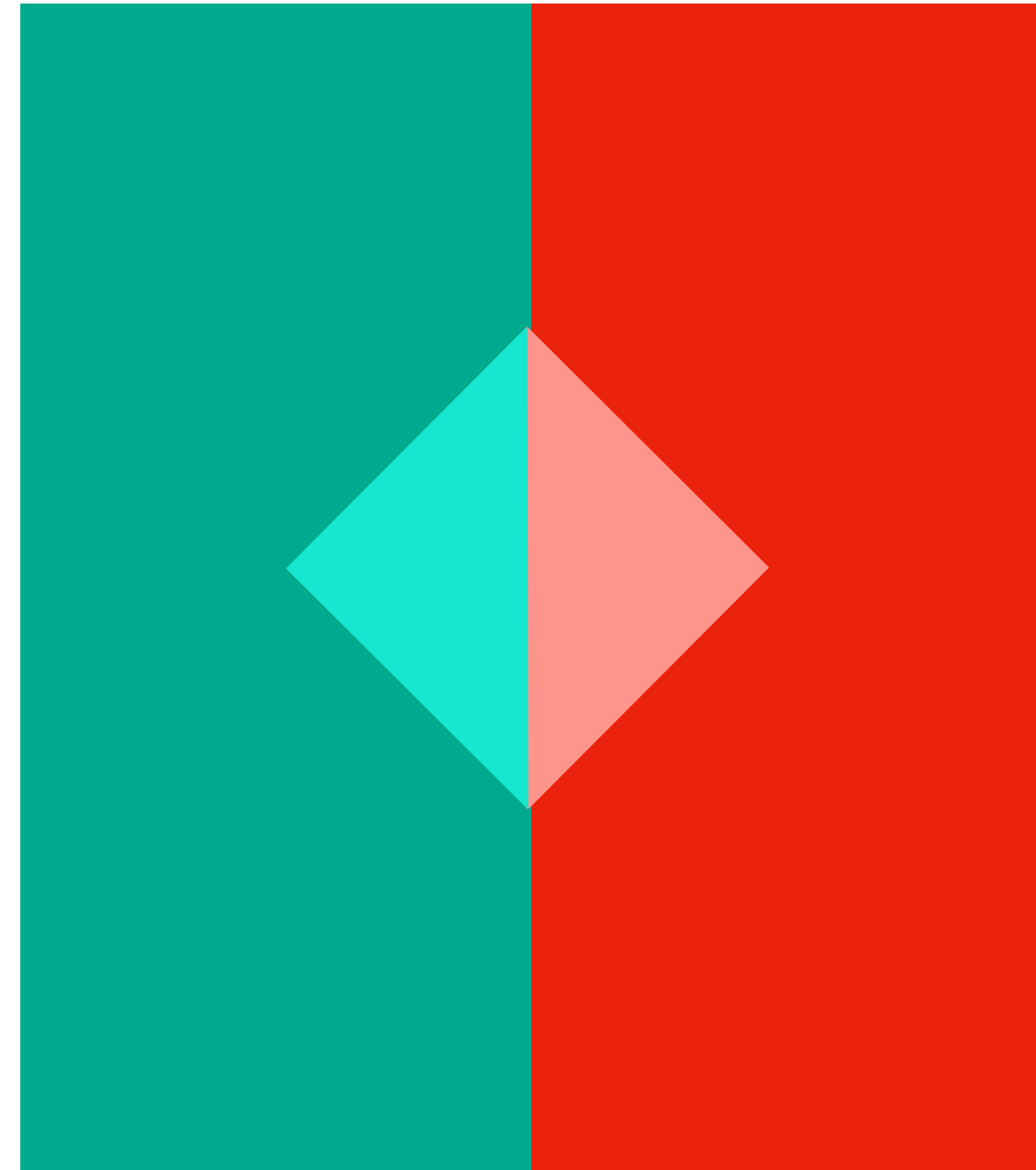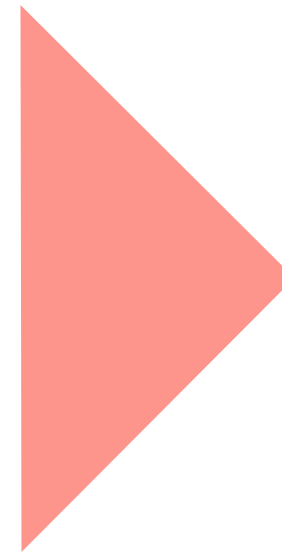Model selects **positive** and patient is **negative**



"Selection space"

Model selects **positive** and patient is **positive**

Model selects **positive** and patient is **negative**

Model selects **negative** and patient is **negative**
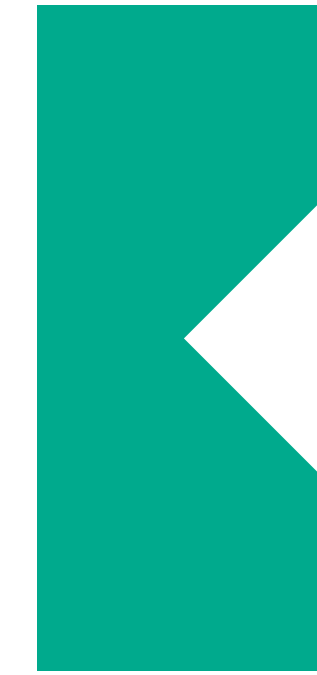
"Selection space"

13

TP: Model selects **positive** and patient is **positive**

FP: Model selects **positive** and patient is **negative**

"Selection space"
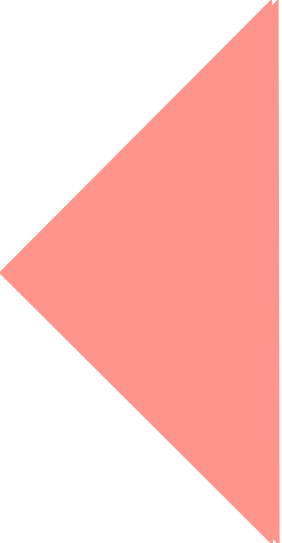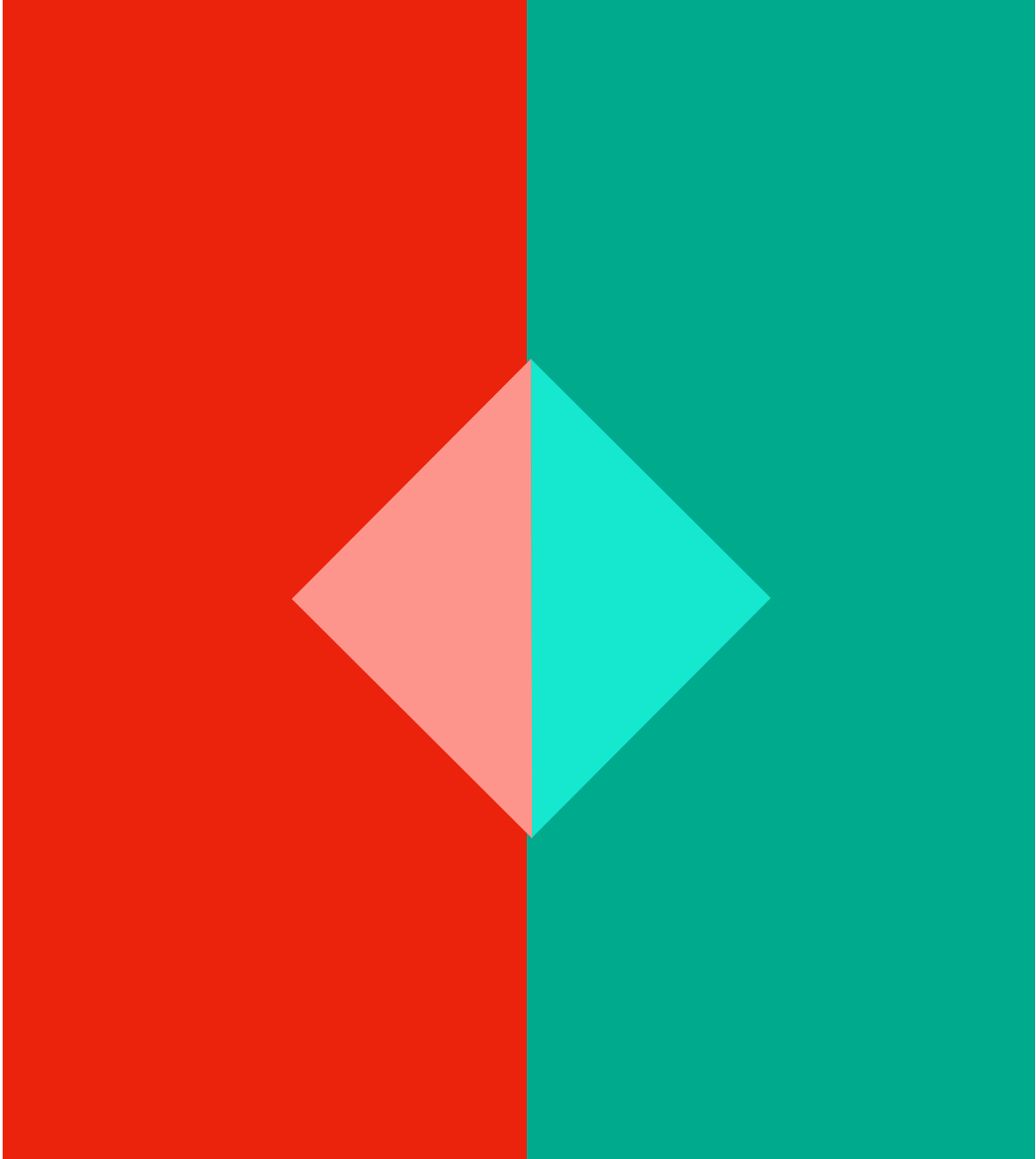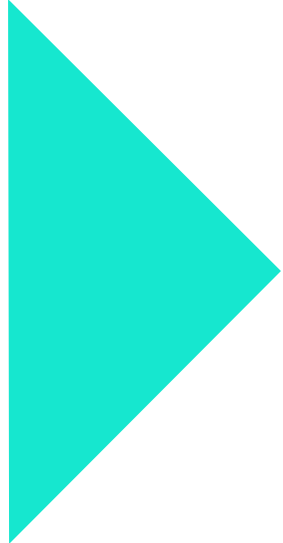
FN: Model selects **negative** and patient is **positive**

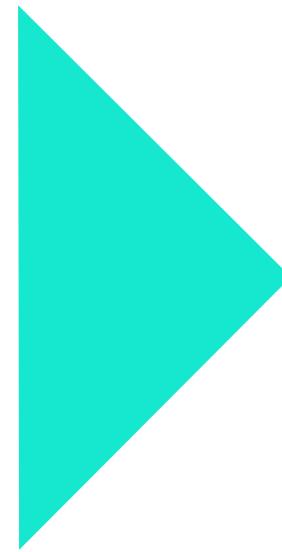TN: Model selects **negative** and patient is **negative**

# **T**RUE **P**OSITIVE

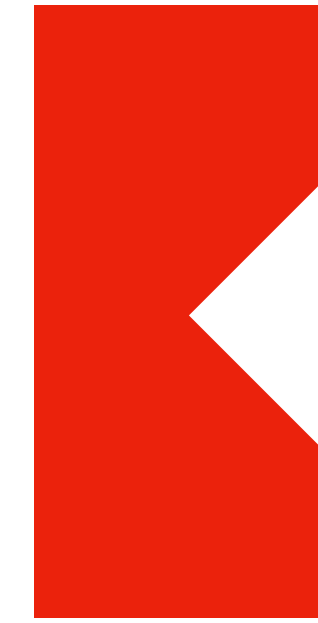TP: Model selects **positive** and patient is **positive**

# **F**ALSE **P**OSITIVE

FP: Model selects **positive** and patient is **negative**

# **F**ALSE **N**EGATIVE

FN: Model selects **negative** and patient is **positive**

# **T**RUE **N**EGATIVE

TN: Model selects **negative** and patient is **negative**
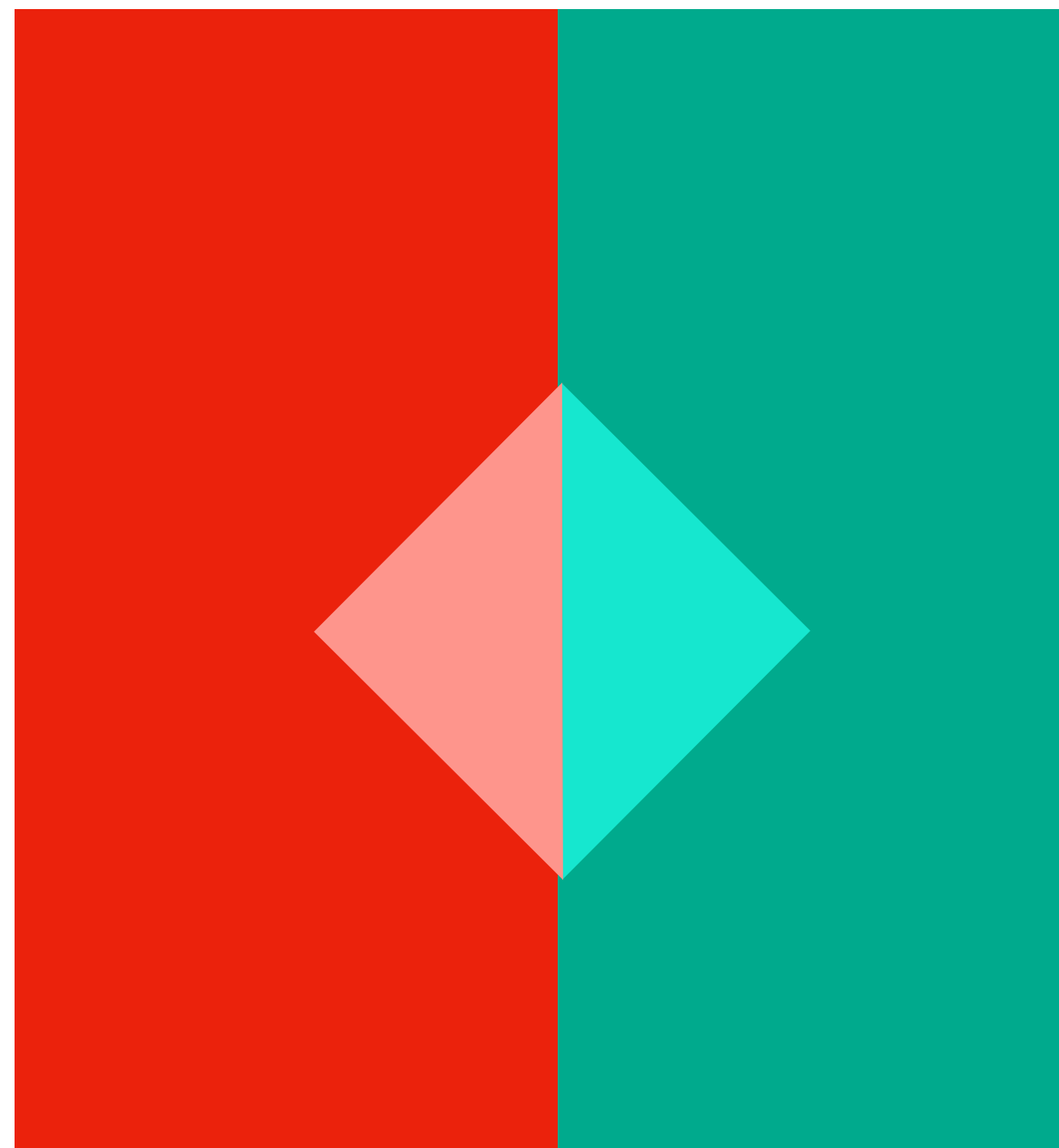
"Selection space"

**Accuracy**

Overall ability of model

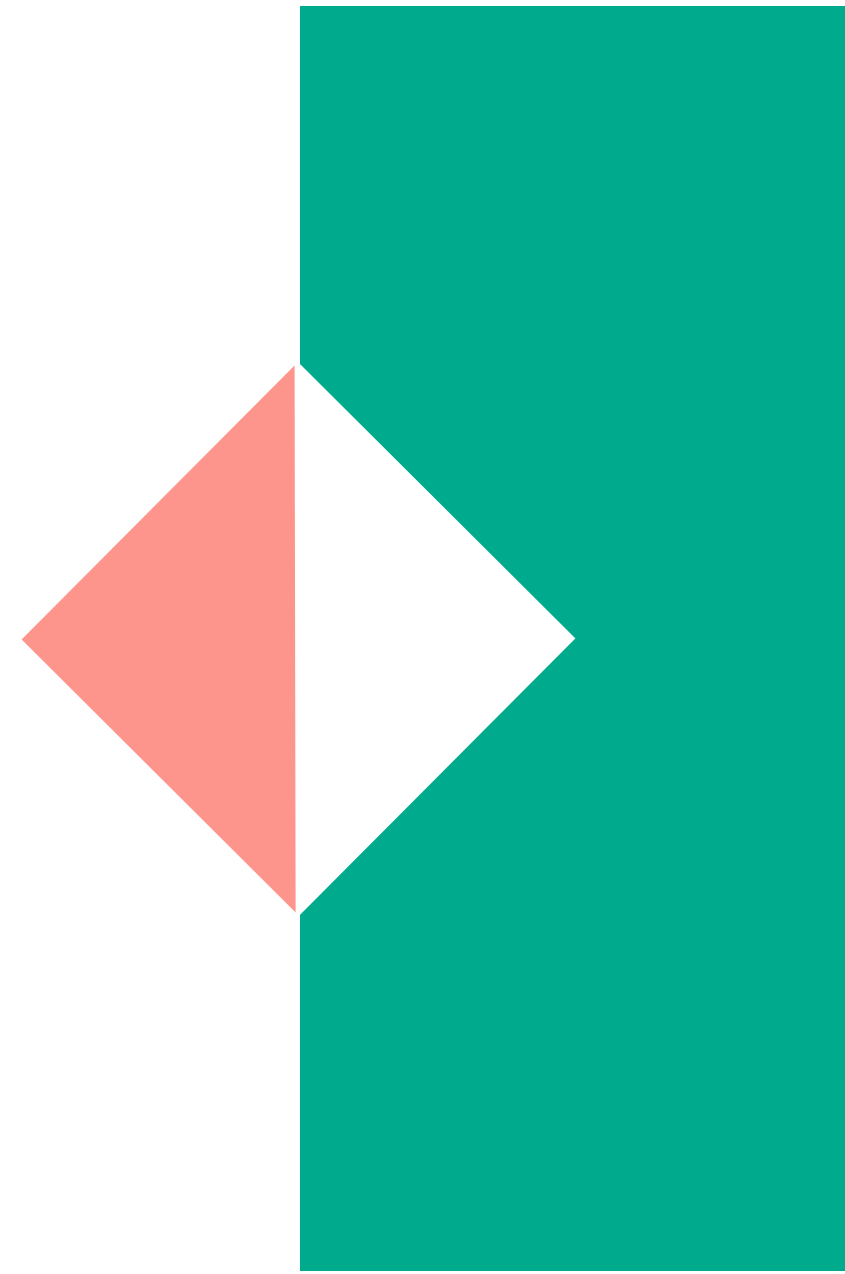"Number of cases where **we chose positive when patient is positive**

*and*

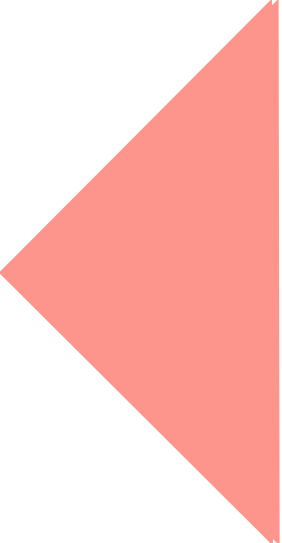Number of cases where **we chose negative when patient is negative**"

"Everything"

TP: Model selects **positive** and patient is **positive**

FP: Model selects **positive** and patient is **negative**

FN: Model selects **negative** and patient is **positive**

TN: Model selects **negative** and patient is **negative**

"Selection space"

**Accuracy**

Overall ability of model

"Number of cases where **we chose positive when patient is positive**"

**Precision**

Amount of selection that's actually correct.

"All selected **positive** by the model"
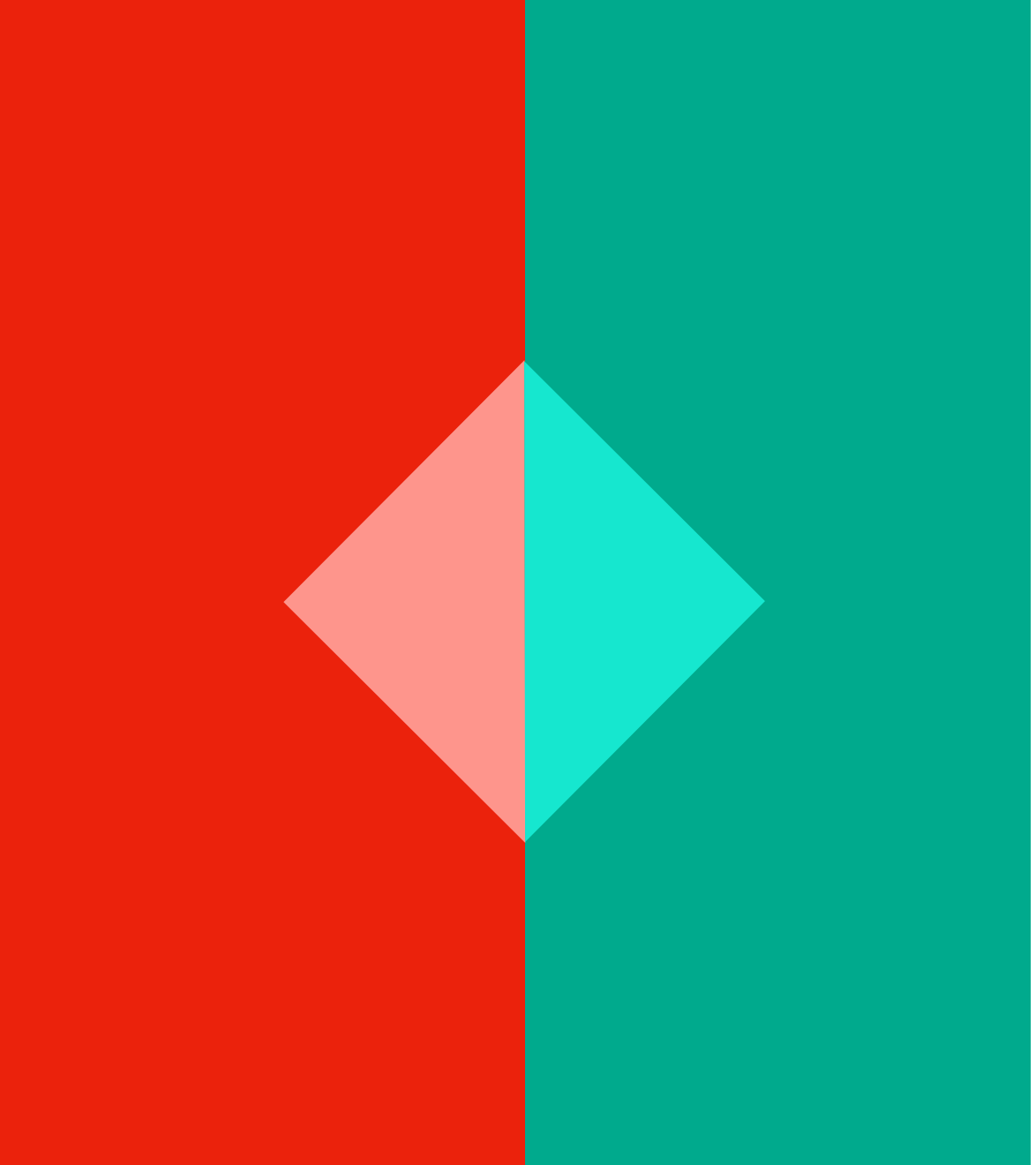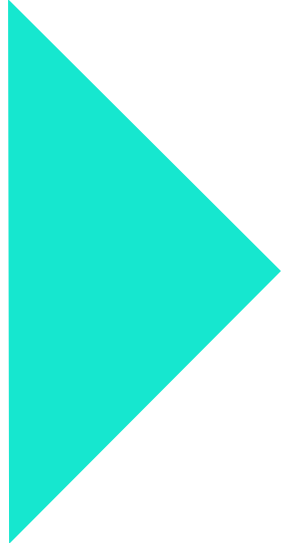
TP: Model selects **positive** and patient is **positive**

FP: Model selects **positive** and patient is **negative**

FN: Model selects **negative** and patient is **positive**

TN: Model selects **negative** and patient is **negative**

"Selection space"

**Accuracy**

Overall ability of model

**Precision**

Amount of selection that's actually correct.

19

"Number of cases where **we chose positive when patient is positive**"

**Recall**

Amount of what needs to be selected that is selected

"All cases that the patients are **positive**"

TP: Model selects **positive** and patient is **positive**

FP: Model selects **positive** and patient is **negative**

"Selection space"

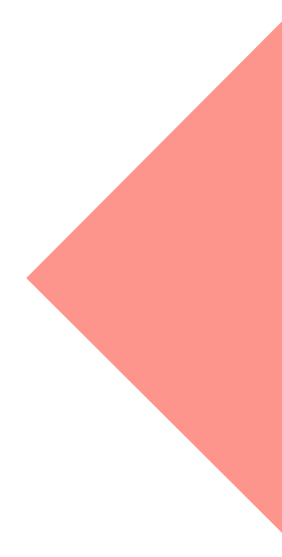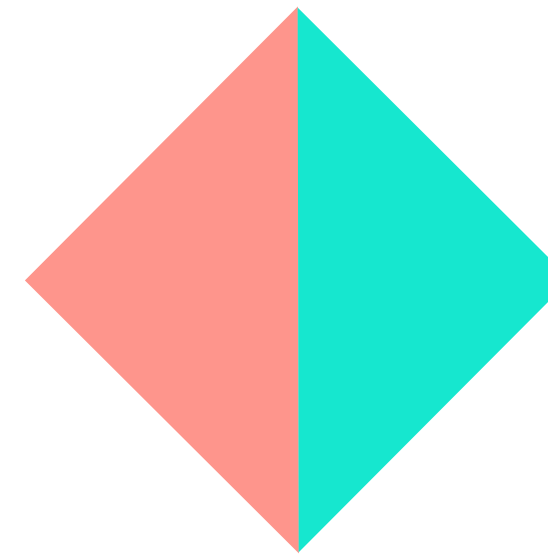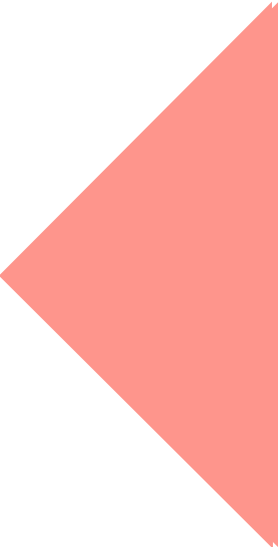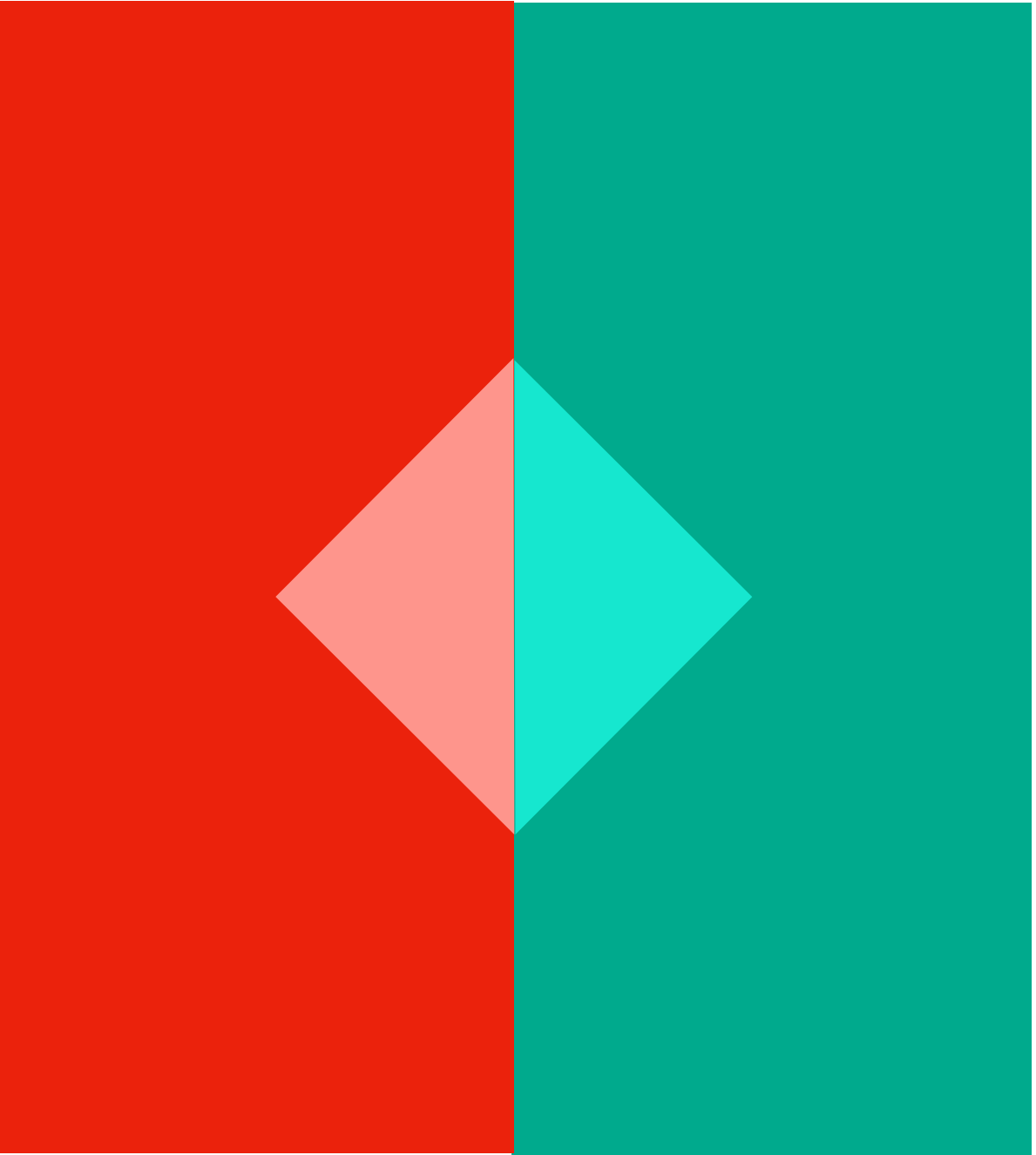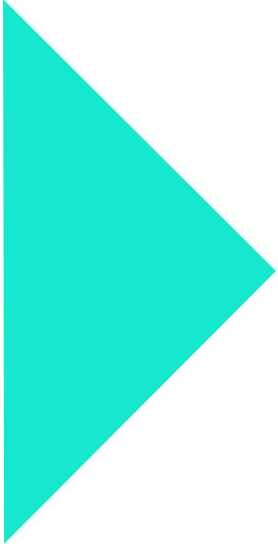FN: Model selects **negative** and patient is **positive**
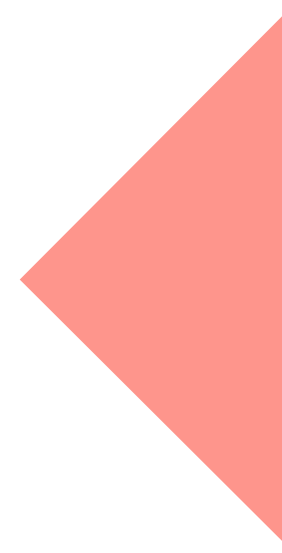
TN: Model selects **negative** and patient is **negative**

**Accuracy**

Overall ability of model

**Precision**

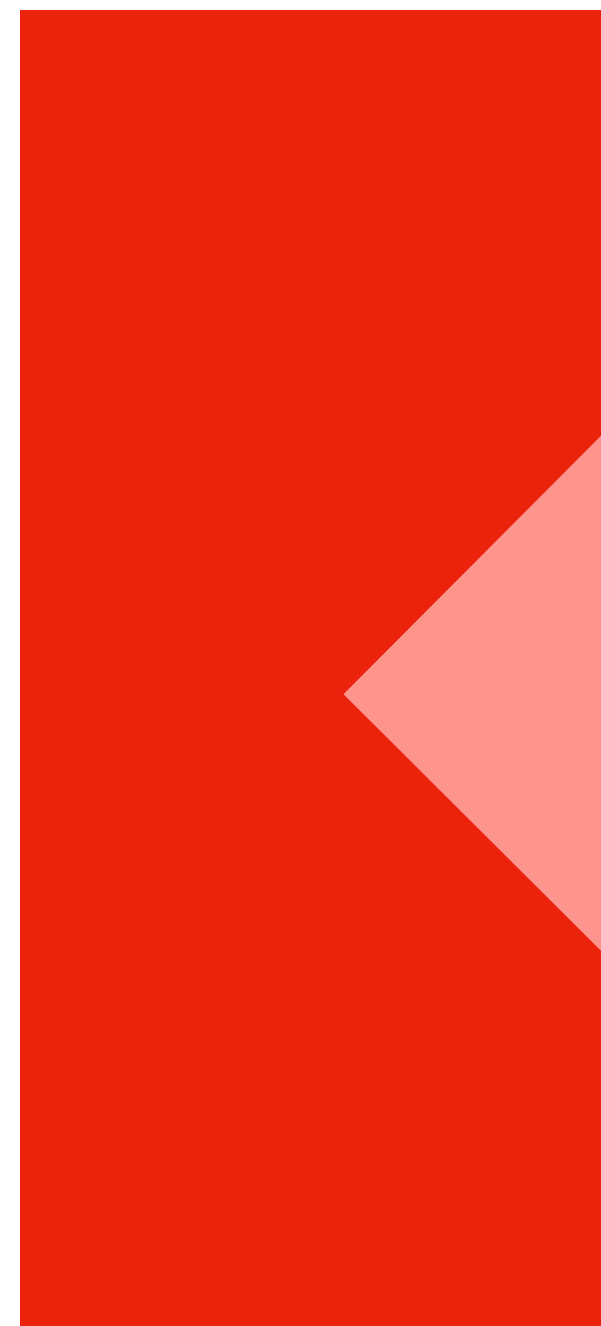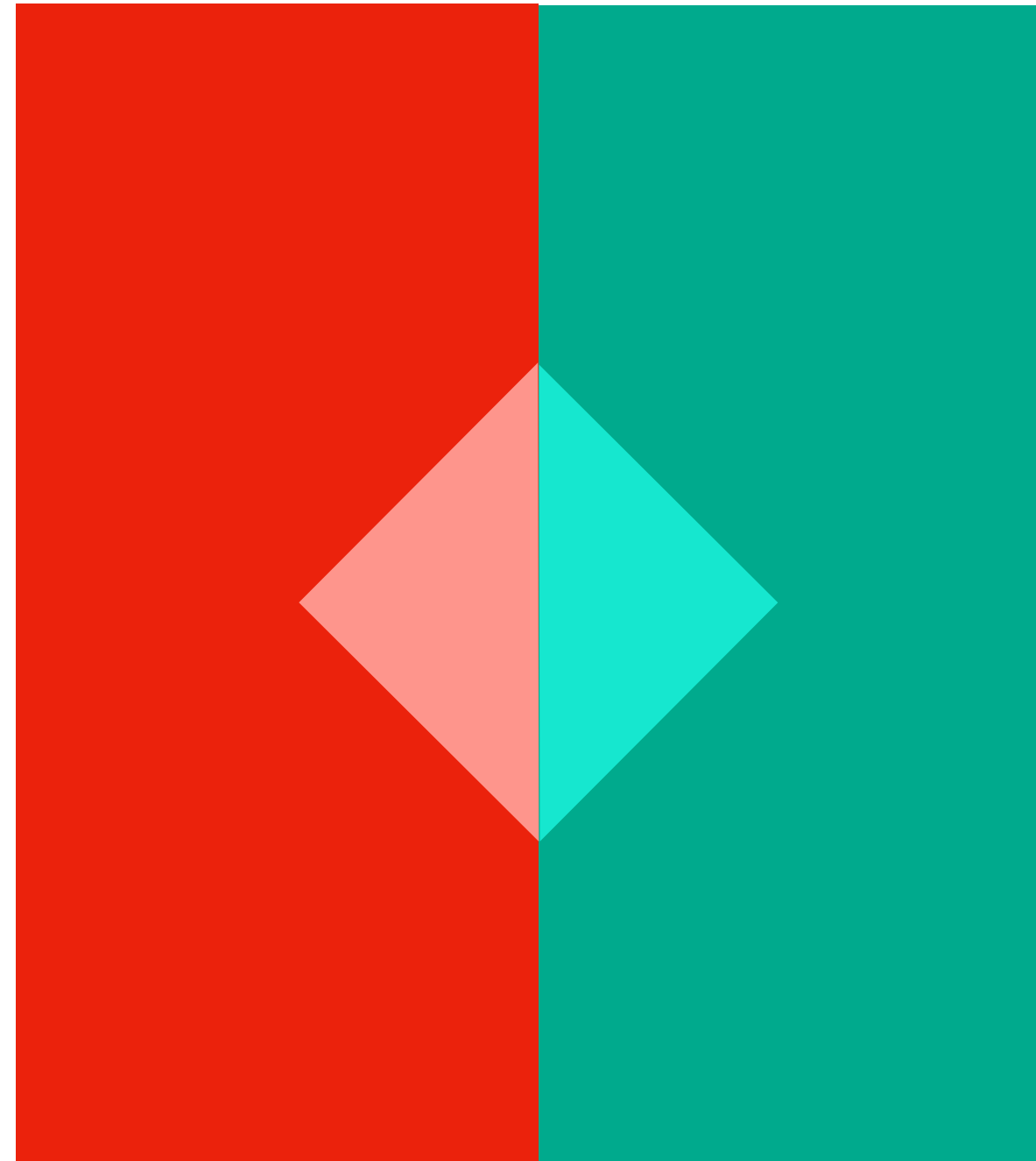Amount of selection that's actually correct.

**Recall**

Amount of what needs to be selected that is selected

"Selection space"

All positive samples

"Selection space"

All positive samples

All negative samples

"Selection space"

"Selection space"

All positive samples

All negative samples

"Selection space"

All positive samples

All negative samples

What our model selected as positive



"Selection space"

All positive samples

All negative samples

What our model selected as positive

What our model selected as negative

"Selection space"

All positive samples

All negative samples

"Number of things that we should select that we did select and the number of things that we shouldn't select that we didn't select."

**Accuracy**

Overall ability of model

"Out of everything"

"Selection space"

**Accuracy**

Overall ability of model

**Precision**

Amount of selection that's actually correct.

"Number of things that we **should** select that we *did* select"

"Number of things that we **should** select that we *did* select and the number of things that we **shouldn't** select that we *did* select."

"Selection space"

**Accuracy**

Overall ability of model

**Precision**

Amount of selection
that's actually correct.

**Recall**

Amount of what needs to
be selected that is selected

"Number of things that we
**should** select that we *did*
select"

"Number of things that we
**should** select in total"

"Selection space"

**Accuracy**

Overall ability of model

**Precision**

Amount of selection that's actually correct.

**Recall**

Amount of what needs to be selected that is selected

**T**RUE **P**OSITIVE

**F**ALSE **N**EGATIVE

**F**ALSE **P**OSITIVE

**T**RUE **N**EGATIVE

**Accuracy**

Overall ability of model

**Precision**

Amount of selection that's actually correct.

**Recall**

Amount of what needs to be selected that is selected

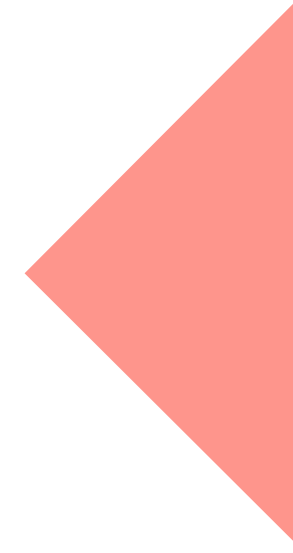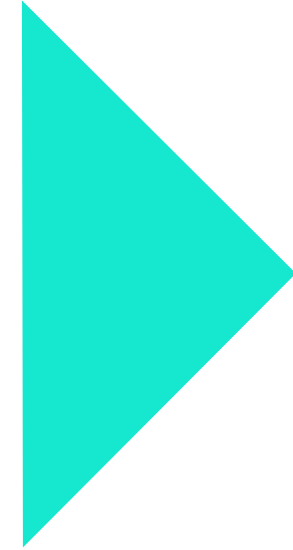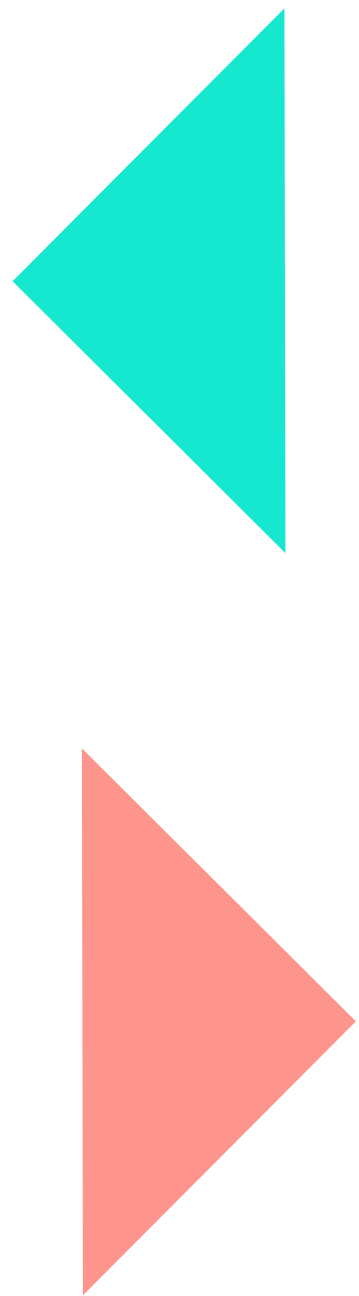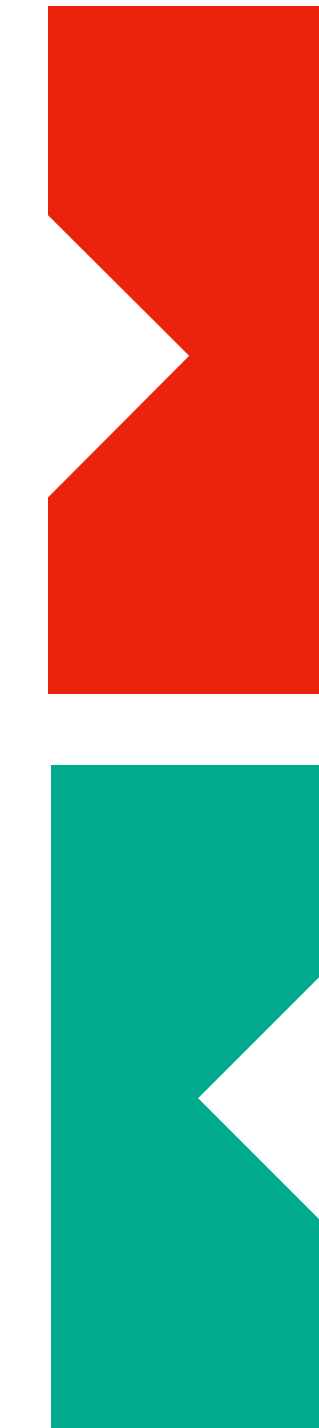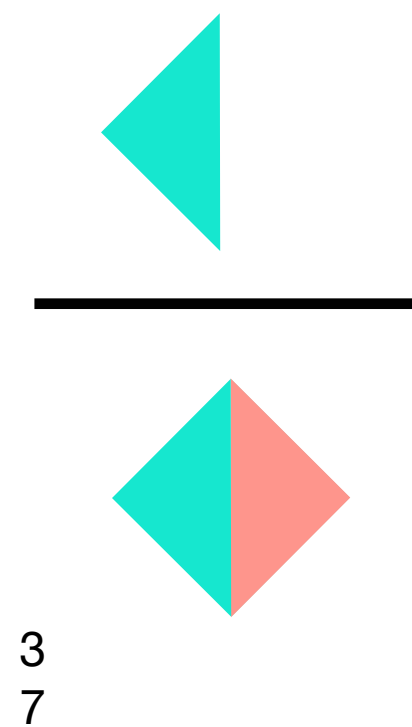| | Predicted condition | | | |
|---|---|---|---|---|
| **Total population** $= P + N$ | **Positive (PP)** | **Negative (PN)** | **Informedness**, bookmaker informedness (BM) $= TPR + TNR - 1$ | **Prevalence threshold** (PT) $= \dfrac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$ |
| **Actual condition** **Positive (P)** | **True positive (TP),** hit | **False negative (FN),** type II error, **miss, underestimation** | **True positive rate** (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \dfrac{TP}{P} = 1 - FNR$ | **False negative rate** (FNR), miss rate $= \dfrac{FN}{P} = 1 - TPR$ |
| **Negative (N)** | **False positive (FP),** type I error, **false alarm, overestimation** | **True negative (TN),** **correct rejection** | **False positive rate** (FPR), probability of false alarm, fall-out $= \dfrac{FP}{N} = 1 - TNR$ | **True negative rate** (TNR), specificity (SPC), selectivity $= \dfrac{TN}{N} = 1 - FPR$ |
| **Prevalence** $= \dfrac{P}{P + N}$ | **Positive predictive value** (PPV), precision $= \dfrac{TP}{PP} = 1 - FDR$ | **False omission rate** (FOR) $= \dfrac{FN}{PN} = 1 - NPV$ | **Positive likelihood ratio** (LR+) $= \dfrac{TPR}{FPR}$ | **Negative likelihood ratio** (LR−) $= \dfrac{FNR}{TNR}$ |
| **Accuracy** (ACC) $= \dfrac{TP + TN}{P + N}$ | **False discovery rate** (FDR) $= \dfrac{FP}{PP} = 1 - PPV$ | **Negative predictive value** (NPV) $= \dfrac{TN}{PN}$ $= 1 - FOR$ | **Markedness** (MK), deltaP (Δp) $= PPV + NPV - 1$ | **Diagnostic odds ratio** (DOR) $= \dfrac{LR+}{LR-}$ |
| **Balanced accuracy** (BA) $= \dfrac{TPR + TNR}{2}$ | **F$_1$ score** $= \dfrac{2\,PPV \times TPR}{PPV + TPR} = \dfrac{2TP}{2TP + FP + FN}$ | **Fowlkes–Mallows index** (FM) $= \sqrt{PPV \times TPR}$ | **Matthews correlation coefficient** (MCC) $= \sqrt{TPR \times TNR \times PPV \times NPV}$ $- \sqrt{FNR \times FPR \times FOR \times FDR}$ | **Threat score** (TS), critical success index (CSI), Jaccard index $= \dfrac{TP}{TP + FN + FP}$ |

https://en.wikipedia.org/wiki/Precision_and_recall

|  | **Predicted Positive** | **Predicted Negative** |
|---|---|---|
| **Actually Positive** | TP | FN |
| **Actually Negative** | FP | TN |

**Accuracy**

**Precision**

**Recall**

**Accuracy**

Overall ability of model

$$\frac{TP + TN}{Total}$$ **exactly zero**

**Precision**

Amount of selection that's actually correct.

$$\frac{TP}{TP + FP}$$

**Recall**

Amount of what needs to be selected that is selected

$$\frac{TP}{TP + FN}$$ **scaled properly!**
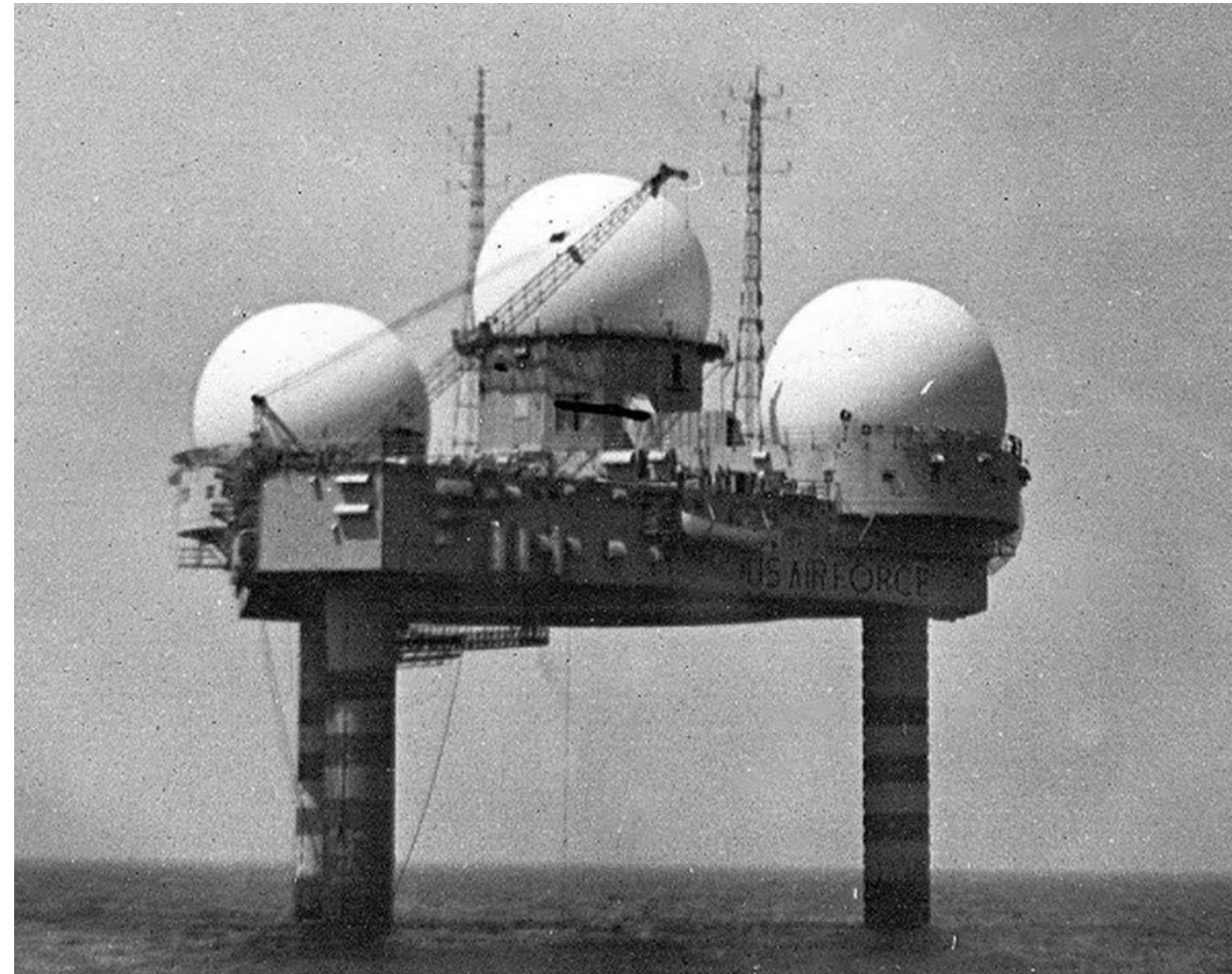
Model → **no progeria regardless**

**Progeria affects ~159 patients in the US**
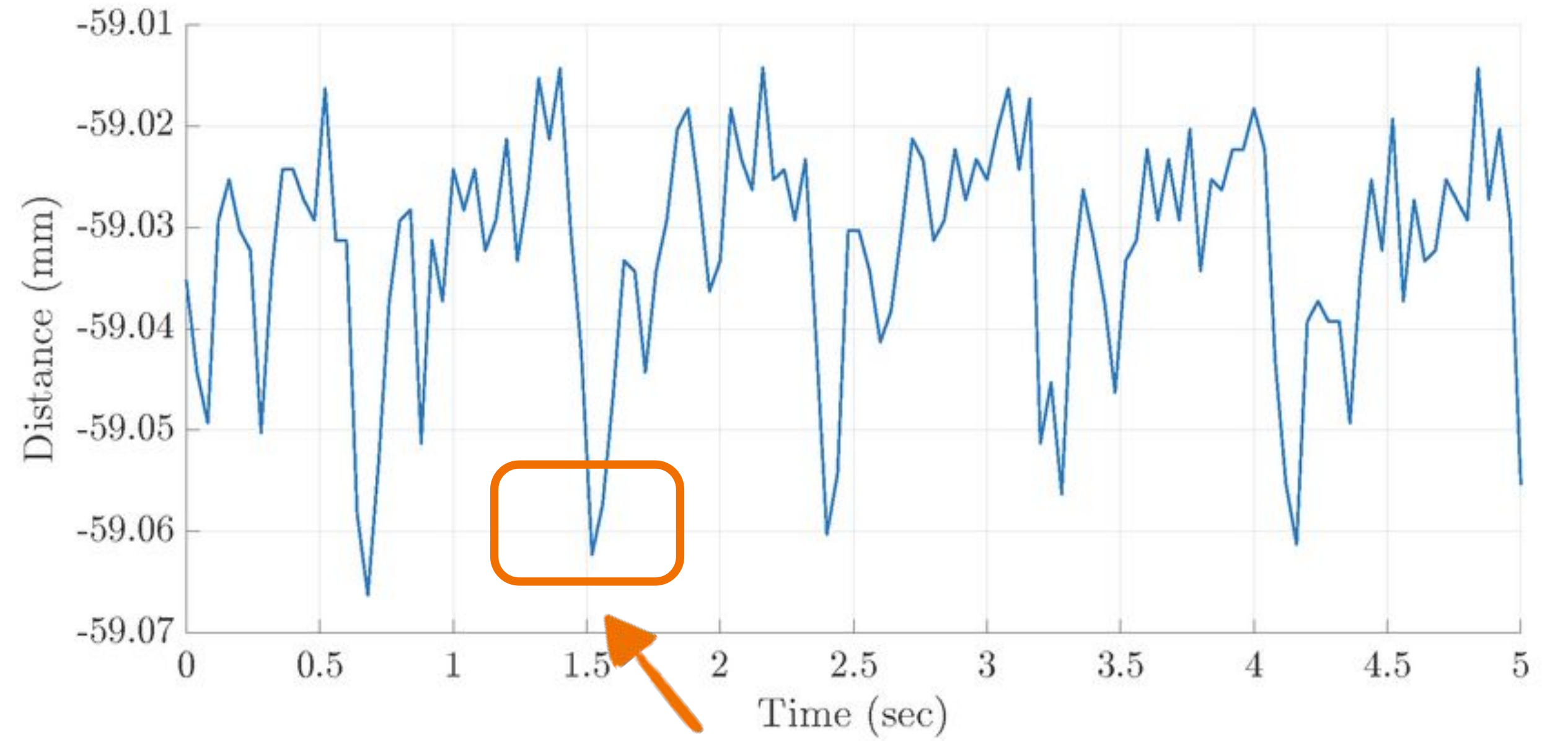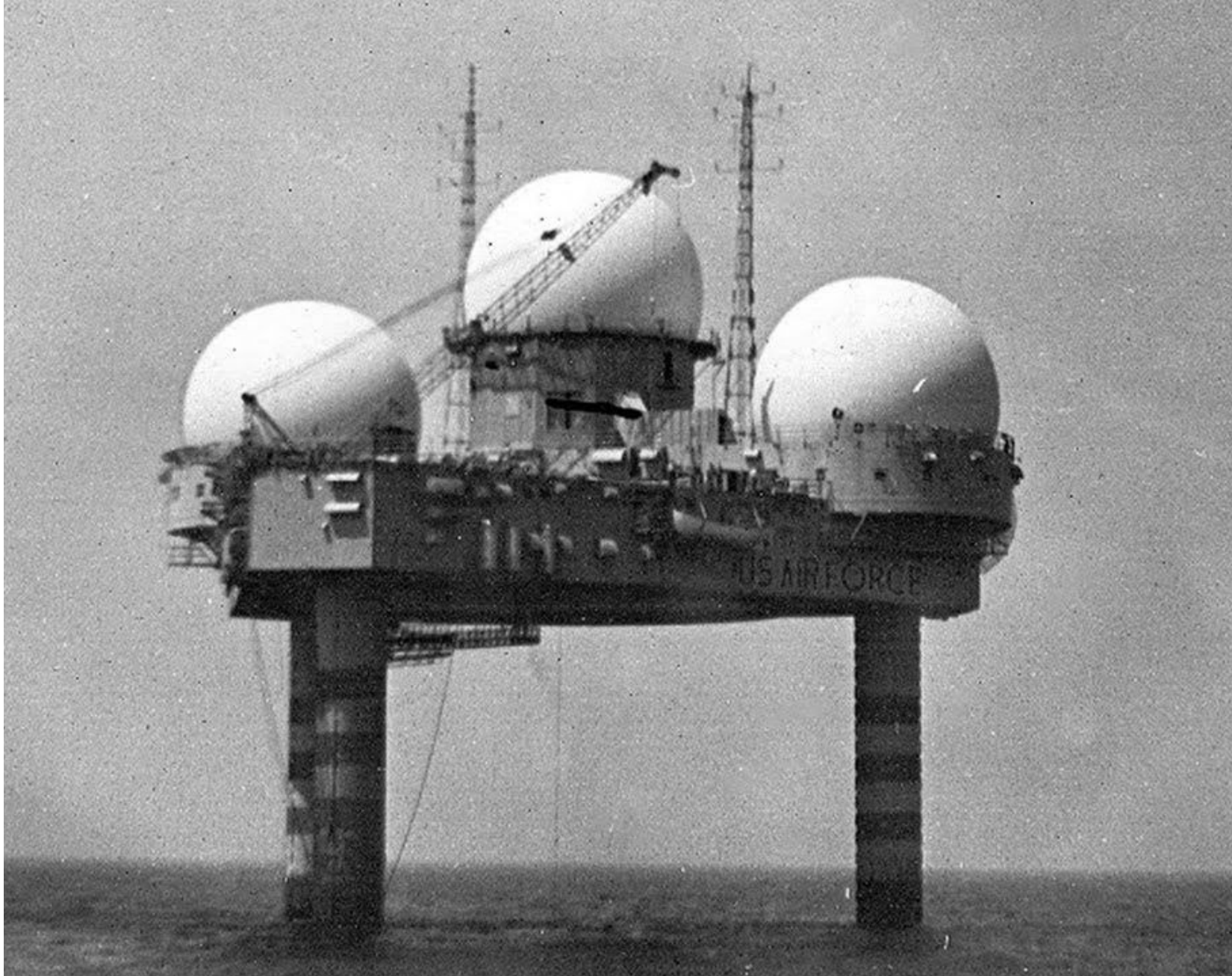
we have a dataset of all American pediatric patients
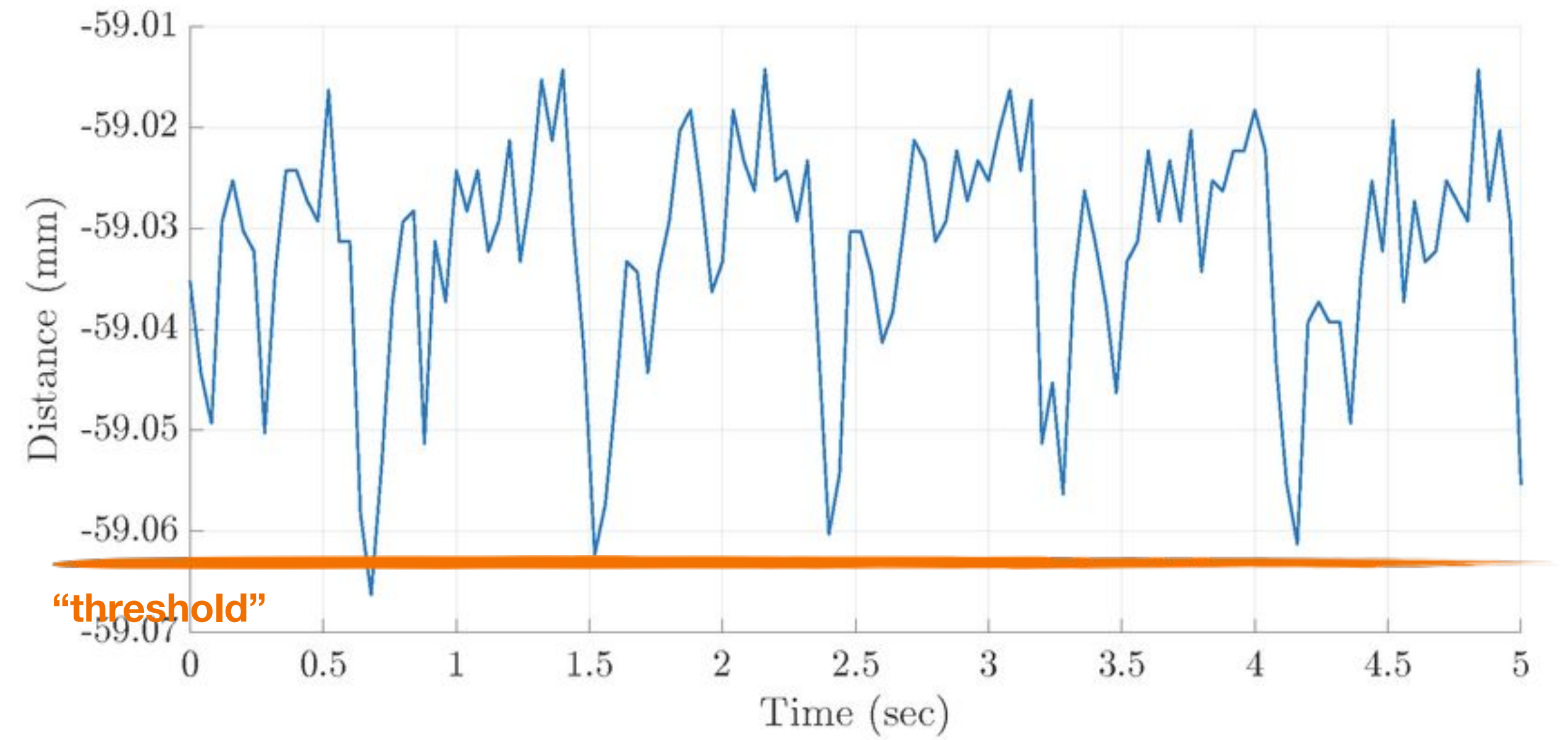
38

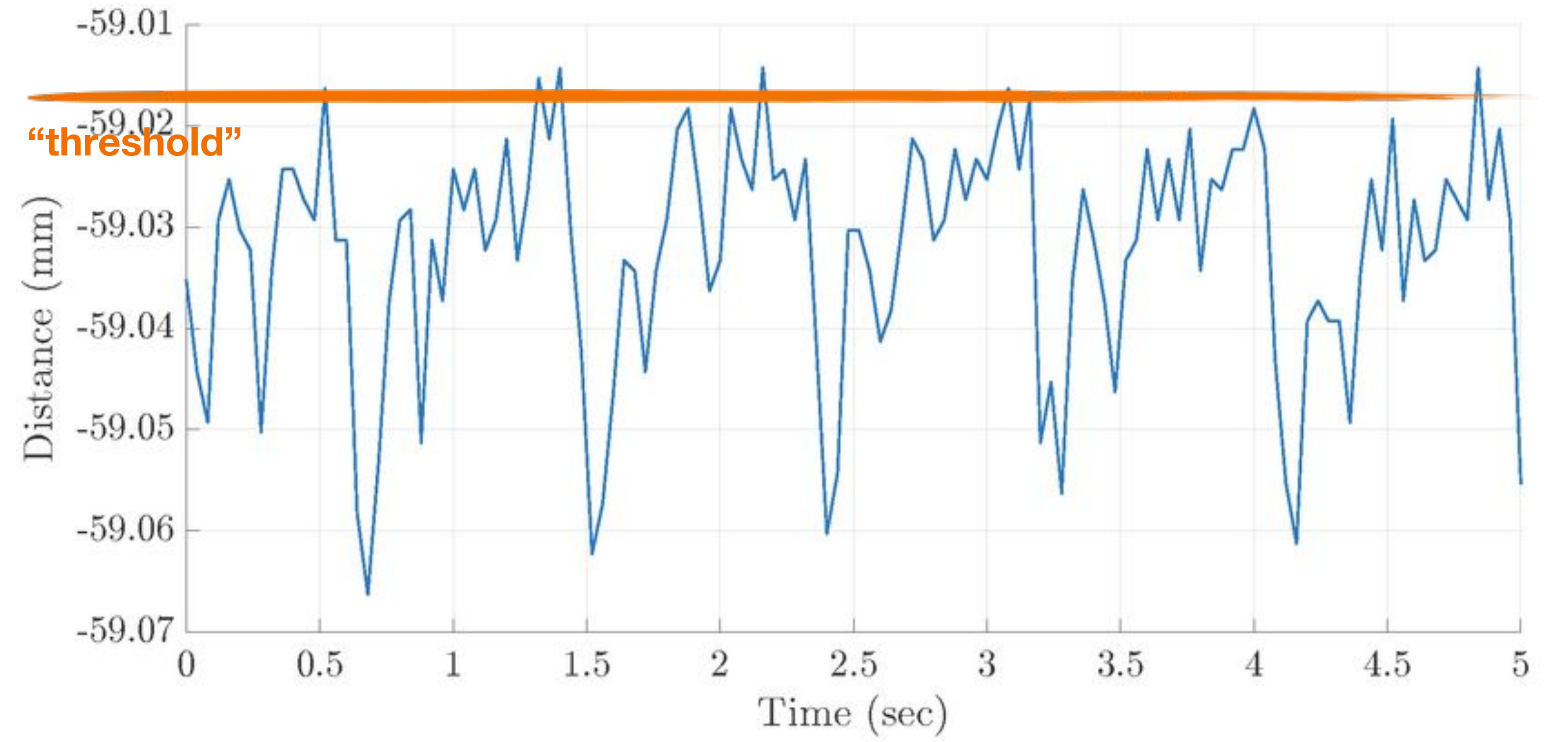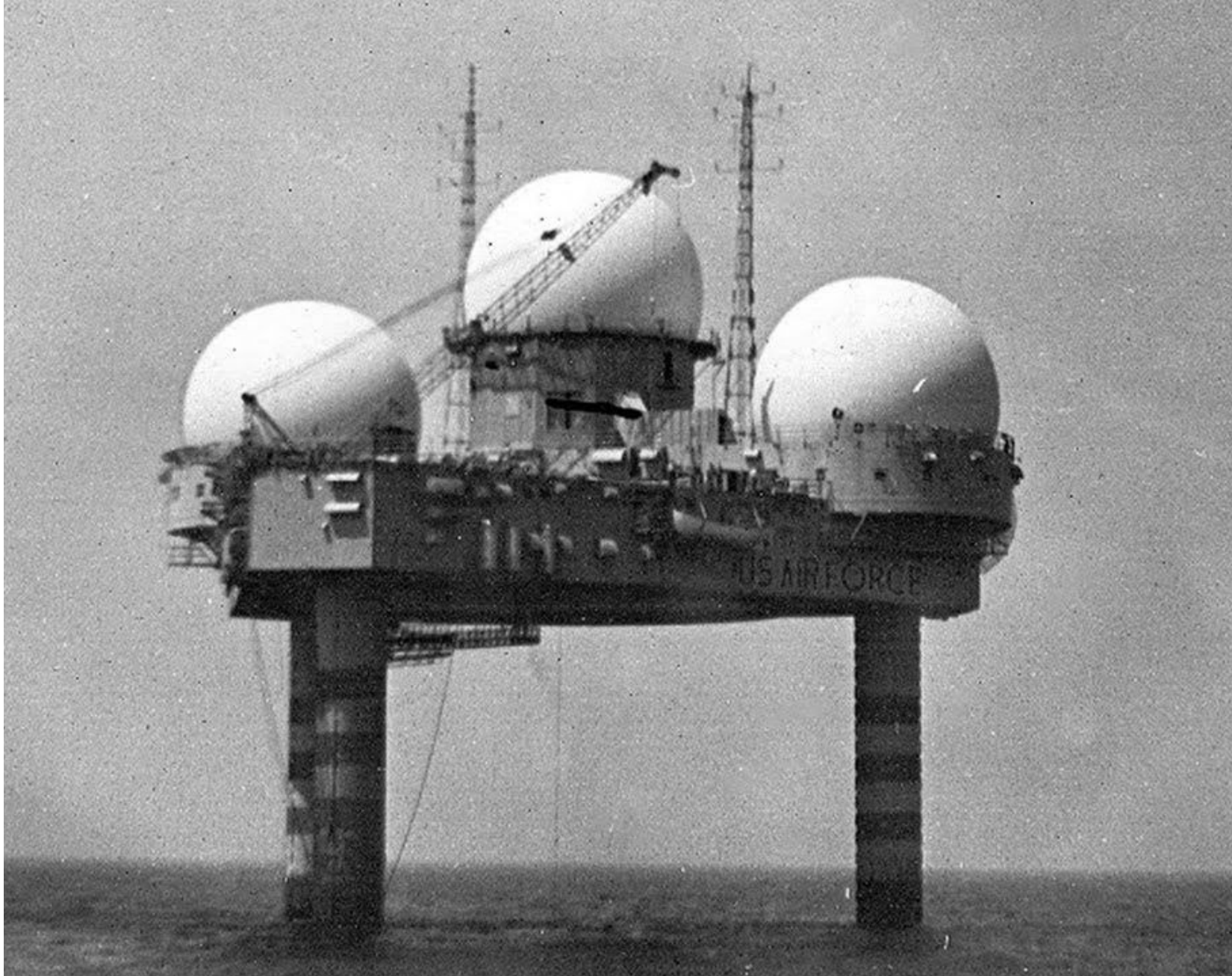# storytime!

# storytime!

# storytime!



does this count as detection?

# storytime!



"threshold"

high recall, low precision

42

# storytime!



"threshold"

high precision, low recall

43

# quantifying "threshold"

**quantifying "threshold"**

# ROC Curve!

**Great!**

**"Random"**

**True Positive**

# ROC Curve!

Receiver Operation Curve

**need lots of false positives before detecting a true positive**

**Awful.**

**False Positive**

■ **ROC Curve** quantify the amount of "error"/noise that is necessary for a classifier to make a good prediction

**Great!**

True Positive

False Positive

# AUC
area under [the ROC] curve

Q: how do you compare these points

■ **AUC** and also Precision-Recall Area Under Curve (PR AUC).

# what makes models fit better

**more** data

**balanced** data

**normalized** data

**quality** data

**more** data      **balanced** data      **normalized** data      **quality** data
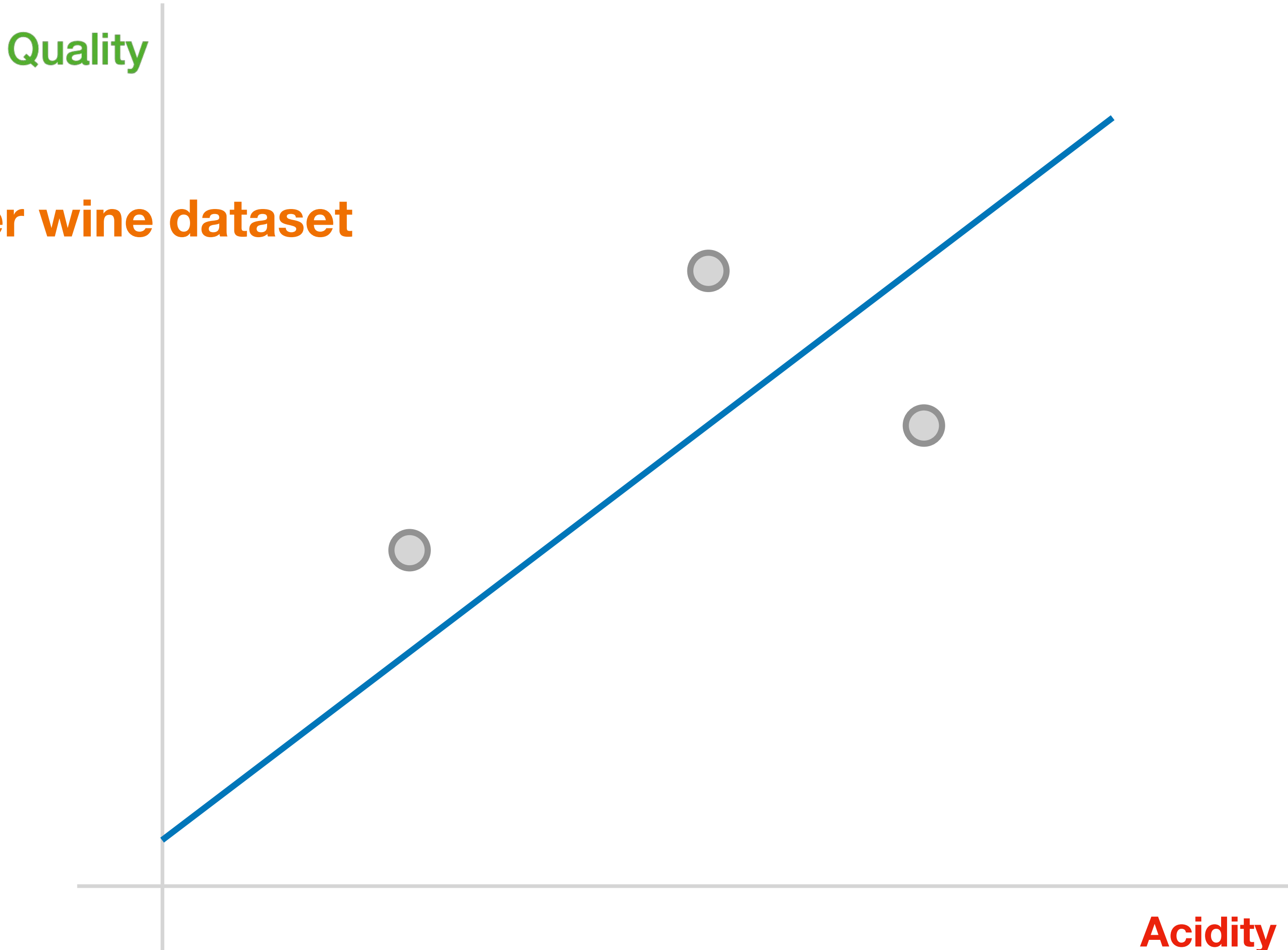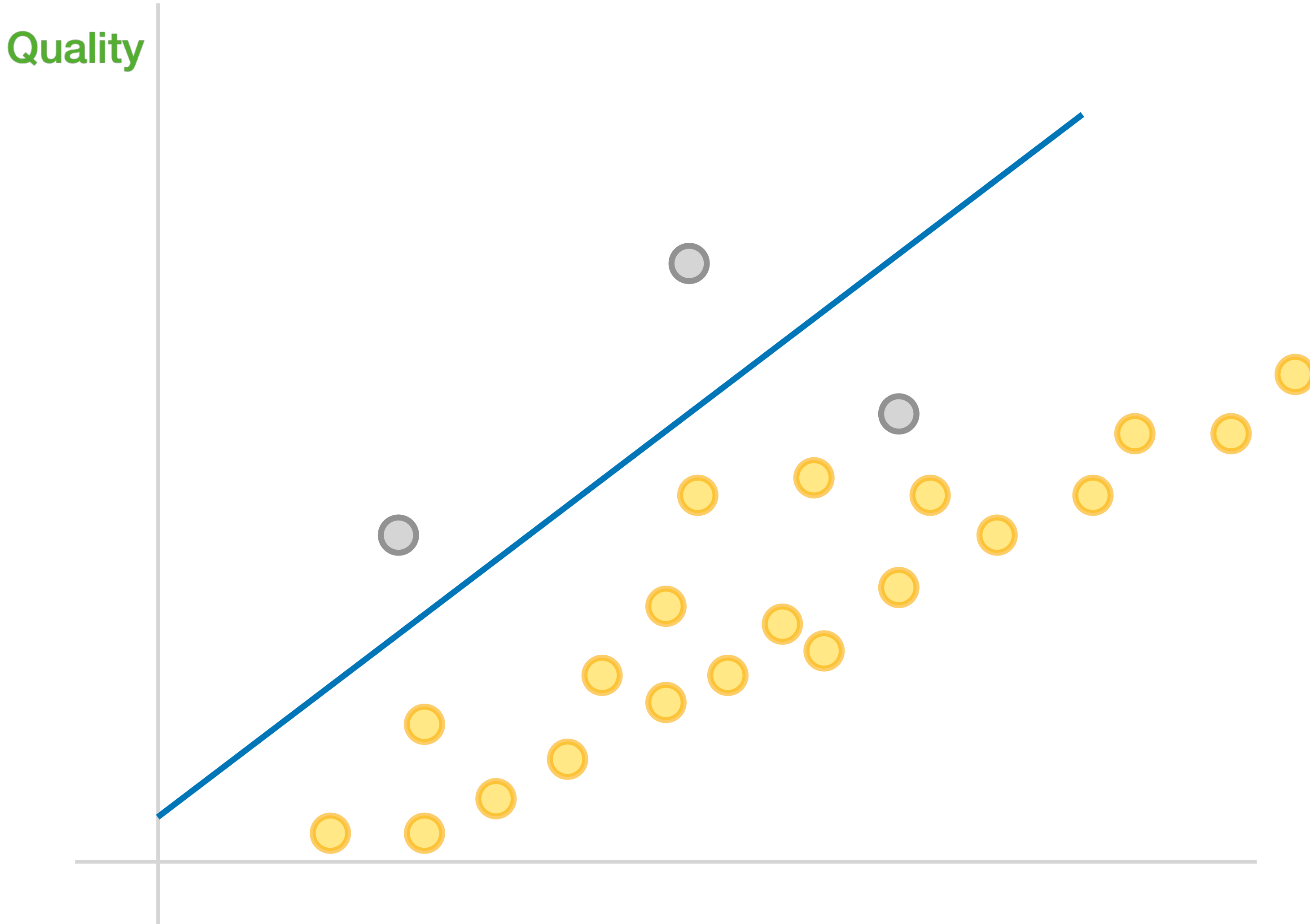
**more** data          **balanced** data          **normalized** data          **quality** data

# **more** data

**let's say we have a simpler wine dataset**
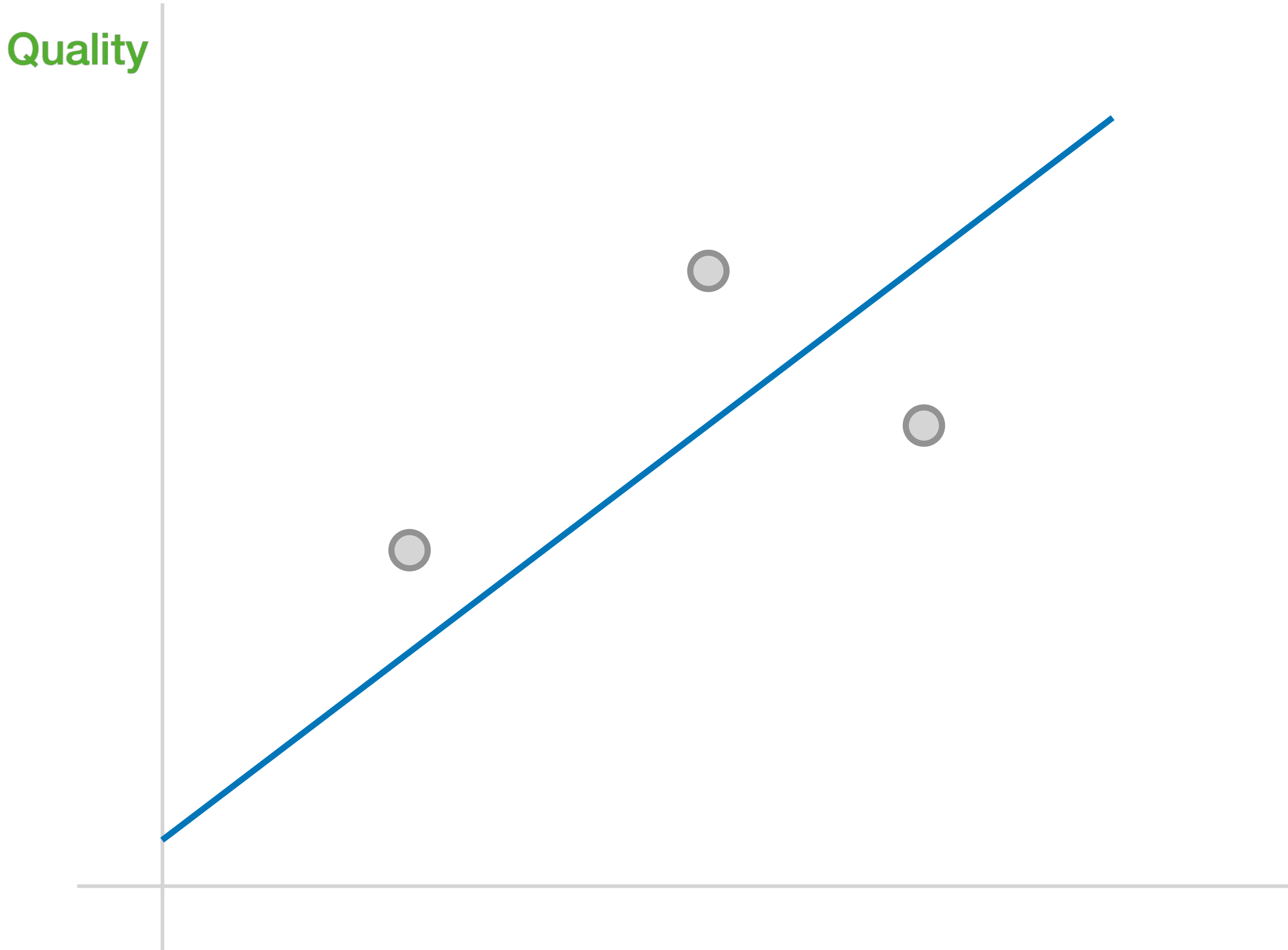
Quality  **on the y axis**
Acidity  on the x axis

Quality

Acidity

# **more** data

**Quality** on the y axis

Acidity on the x axis

Quality

# **more** data

**Quality** on the y axis

Acidity on the x axis

Quality

# **more** data

**Quality** on the y axis

Acidity on the x axis
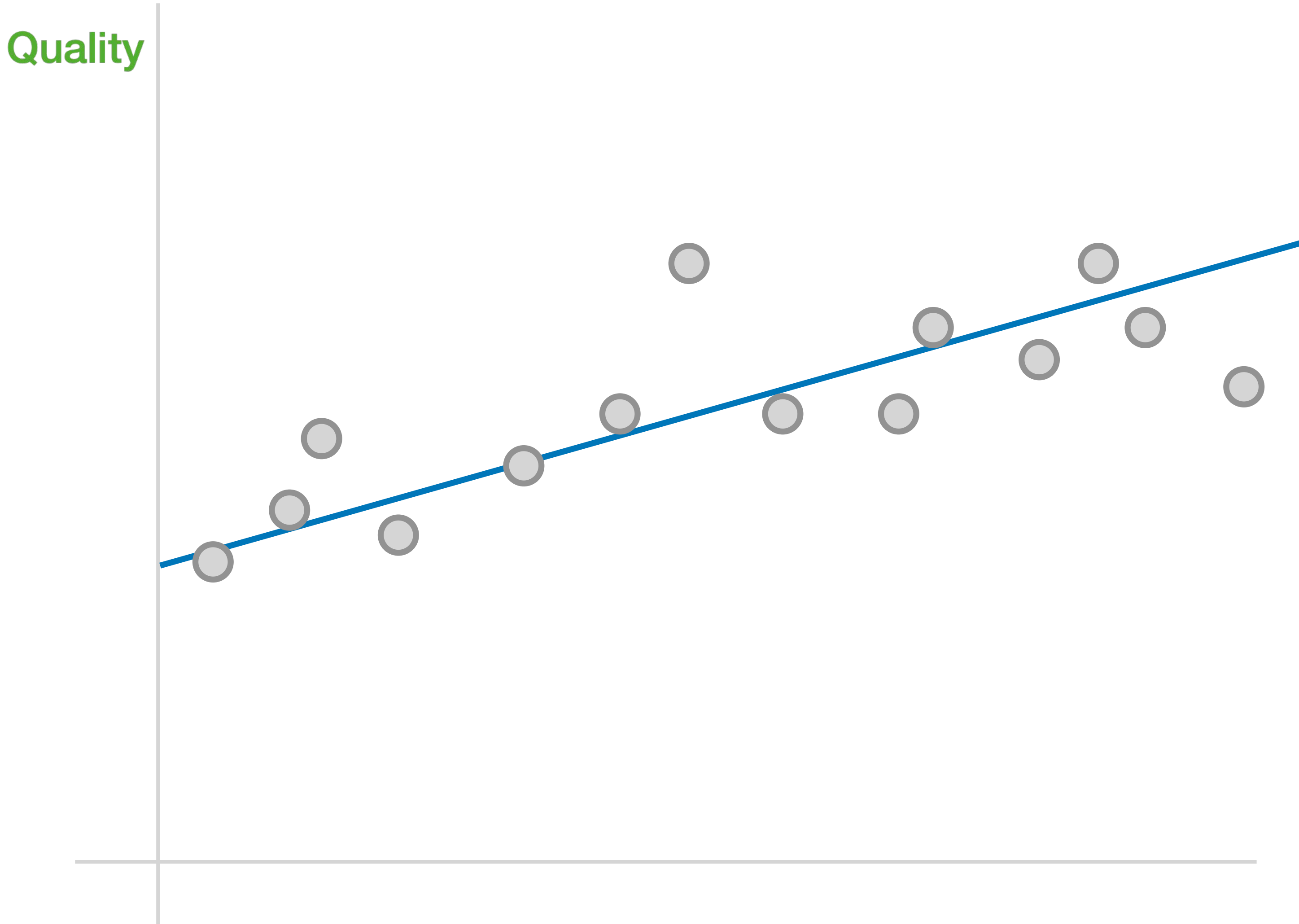
Quality

# **more** data

**Quality** on the y axis

Acidity on the x axis

Quality

# **more** data

Quality

Quality on the y axis
Acidity on the x axis

■ **use more data, get more accurate results**

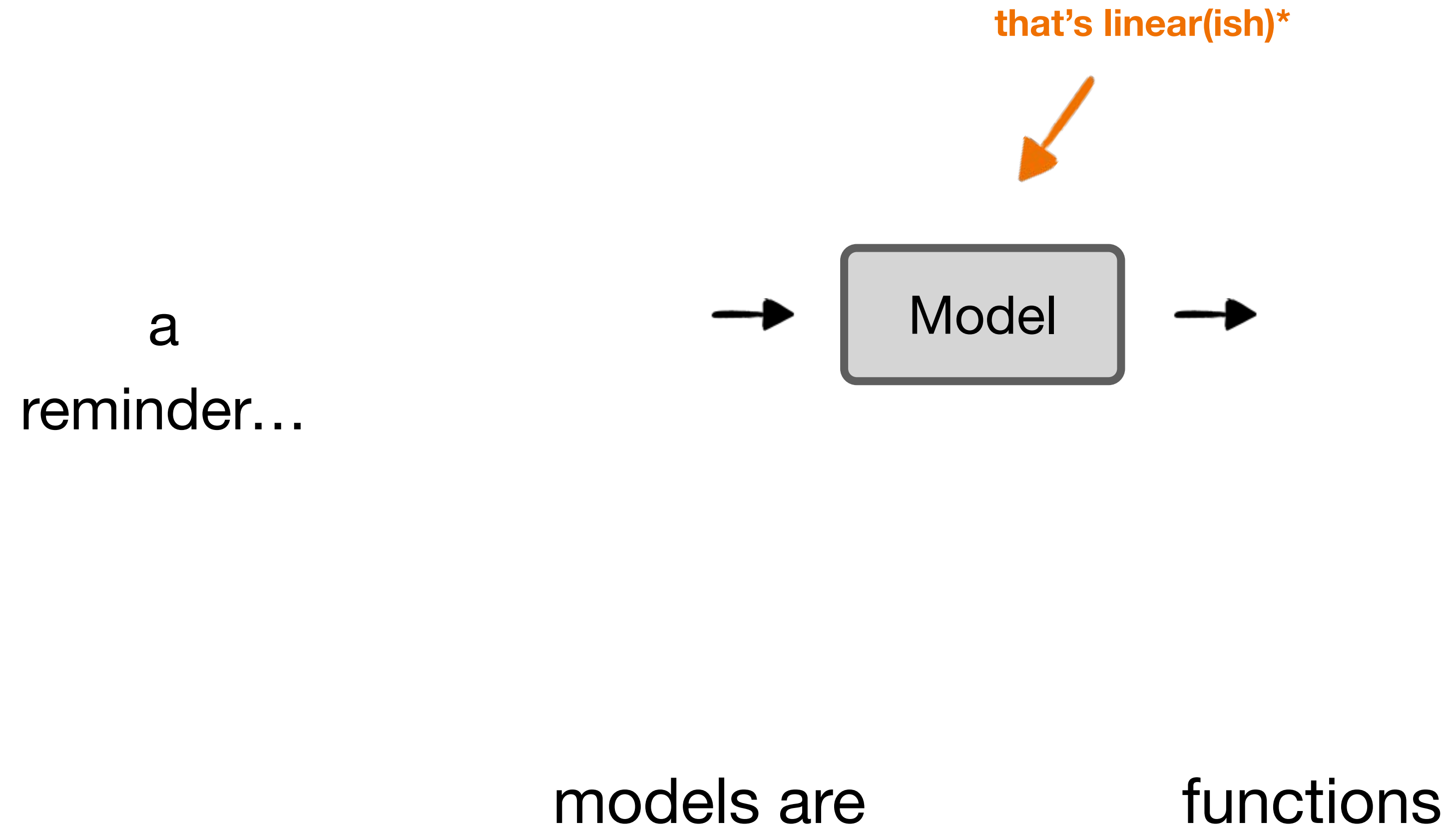**more** data          **balanced** data          **normalized** data          **quality** data

| Amber Colored | Angora Super Sweet | Black Ethiopian | Burbank Slicing |
| Dona | Sophie's Choice | White Bush | Ace 55 |

**more** data          **balanced** data          **normalized** data          **quality** data

**more**

data

that's linear(ish)*

→ Model →

a
reminder…

models are            functions

# more
## data

**that's linear(ish)\***

Quality
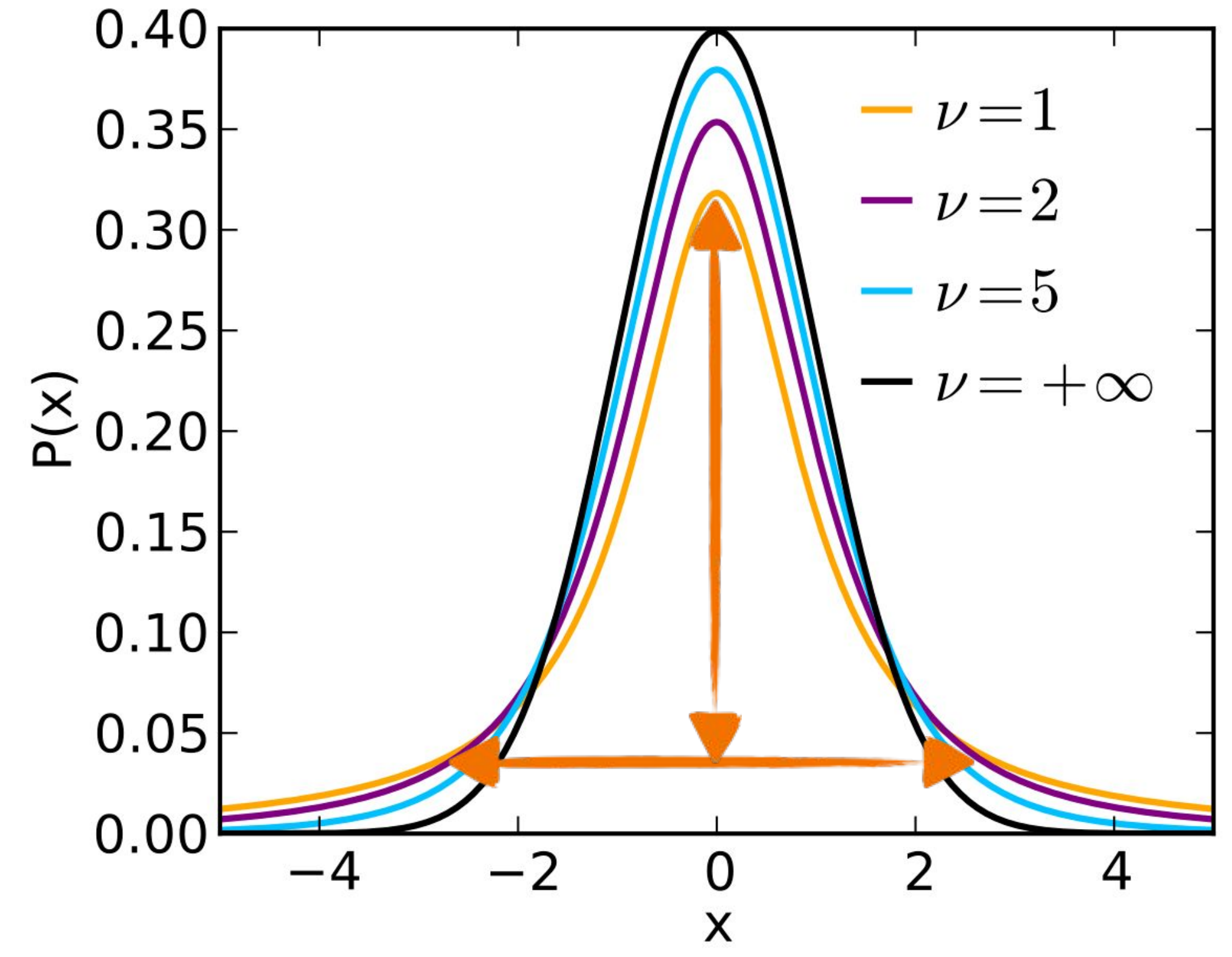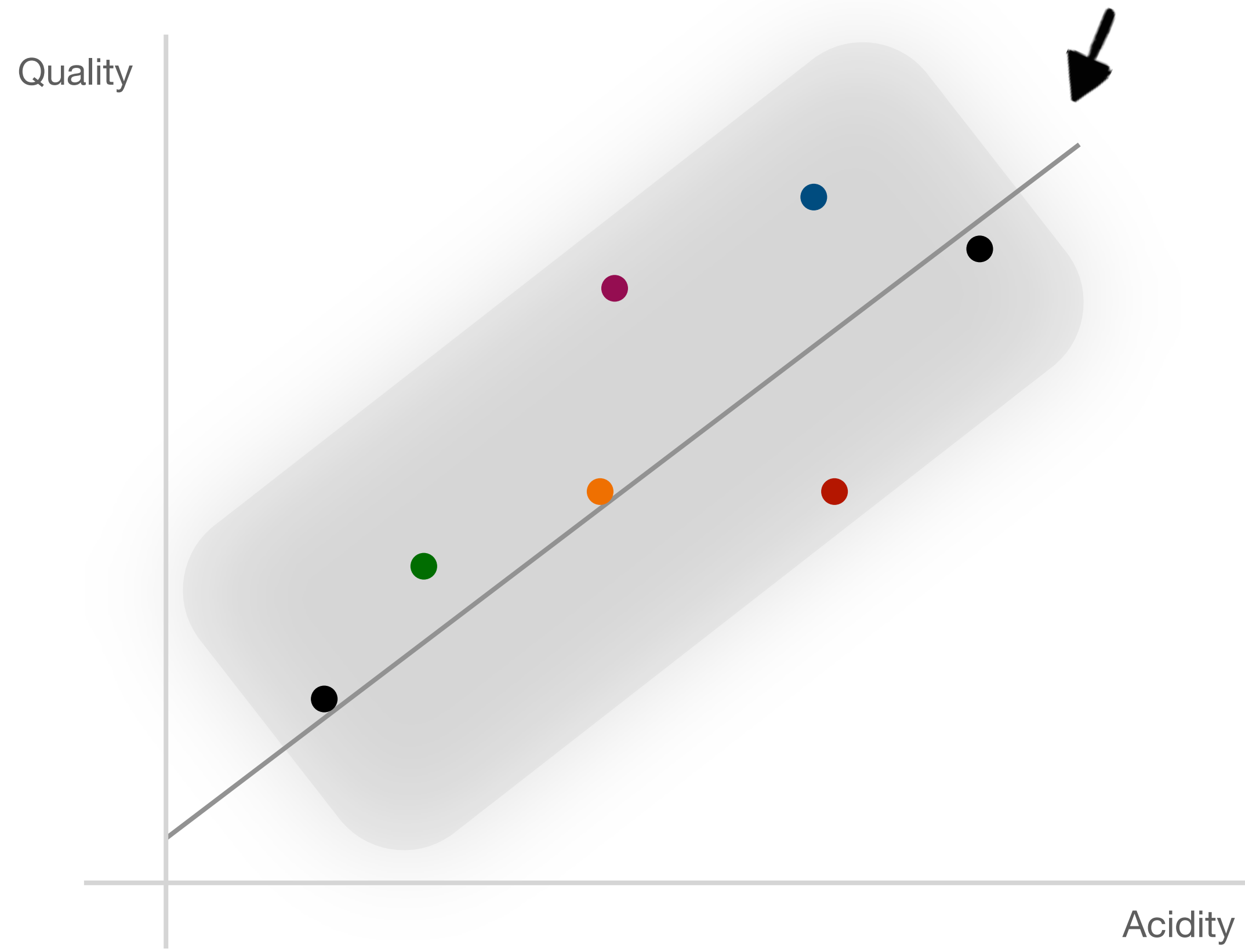
Acidity

# more
data

Quality

??? **still not super well described**

Acidity

# more
data

**more**
 data



wait… this is a t-test!

# more

data

# more
## data

**more**
data

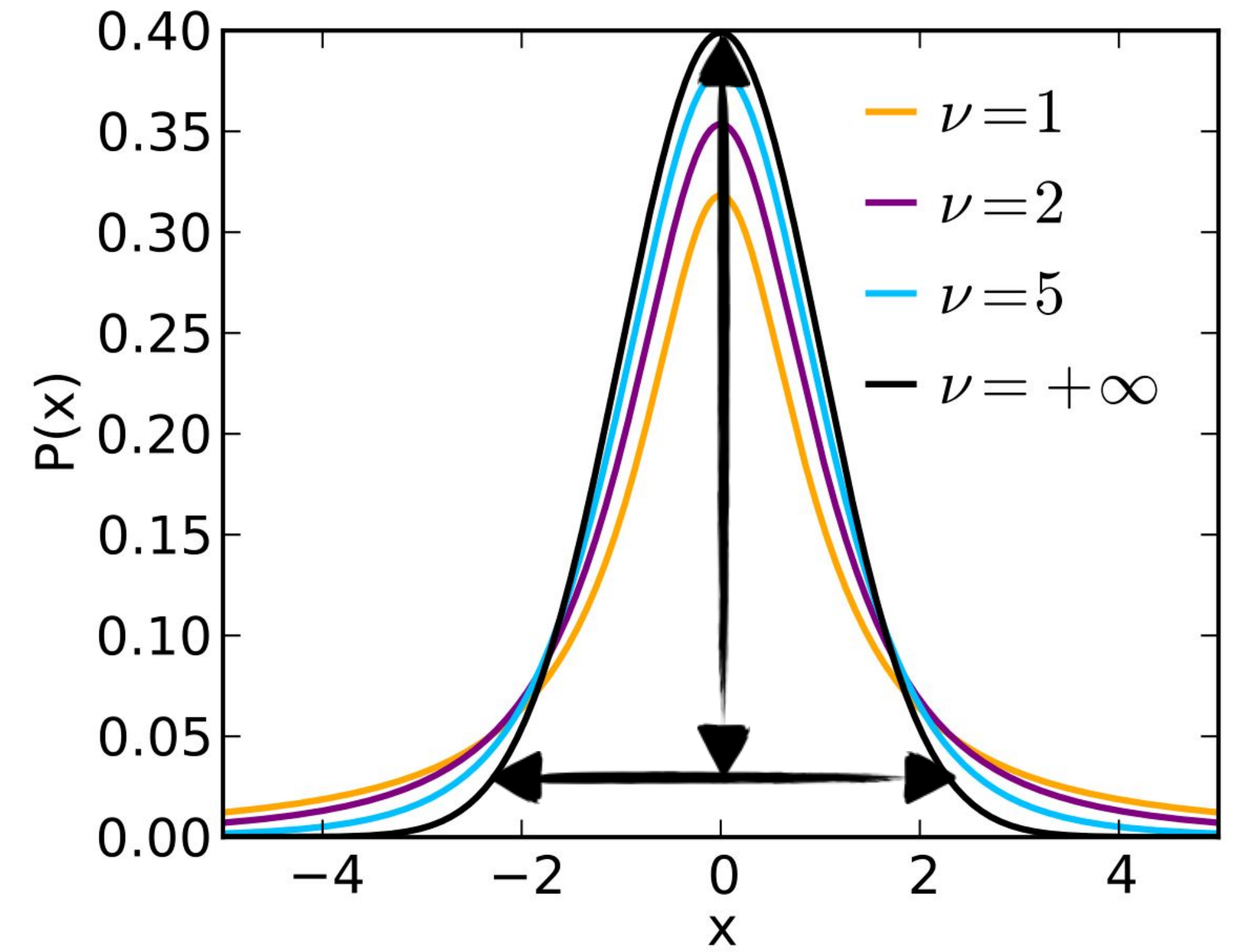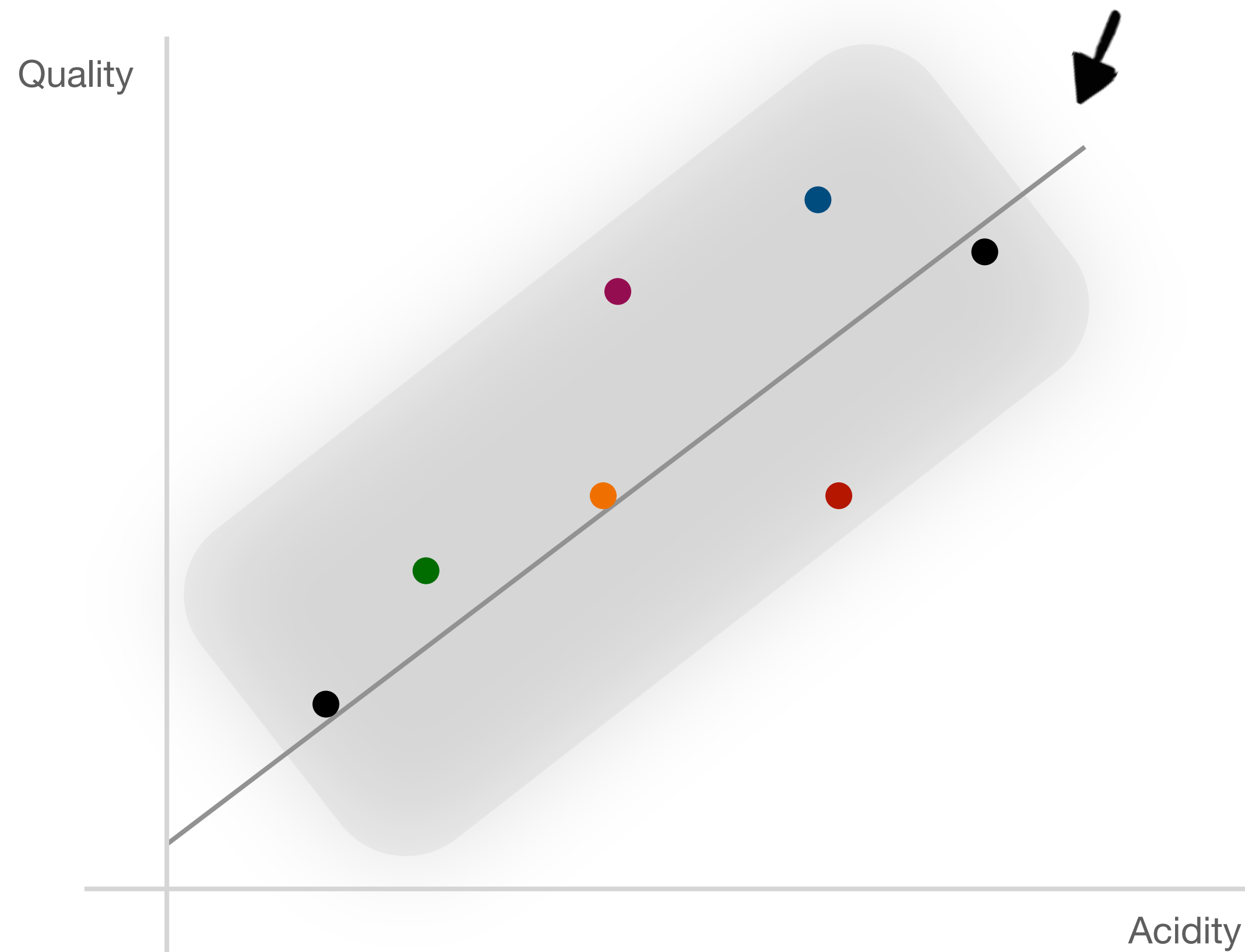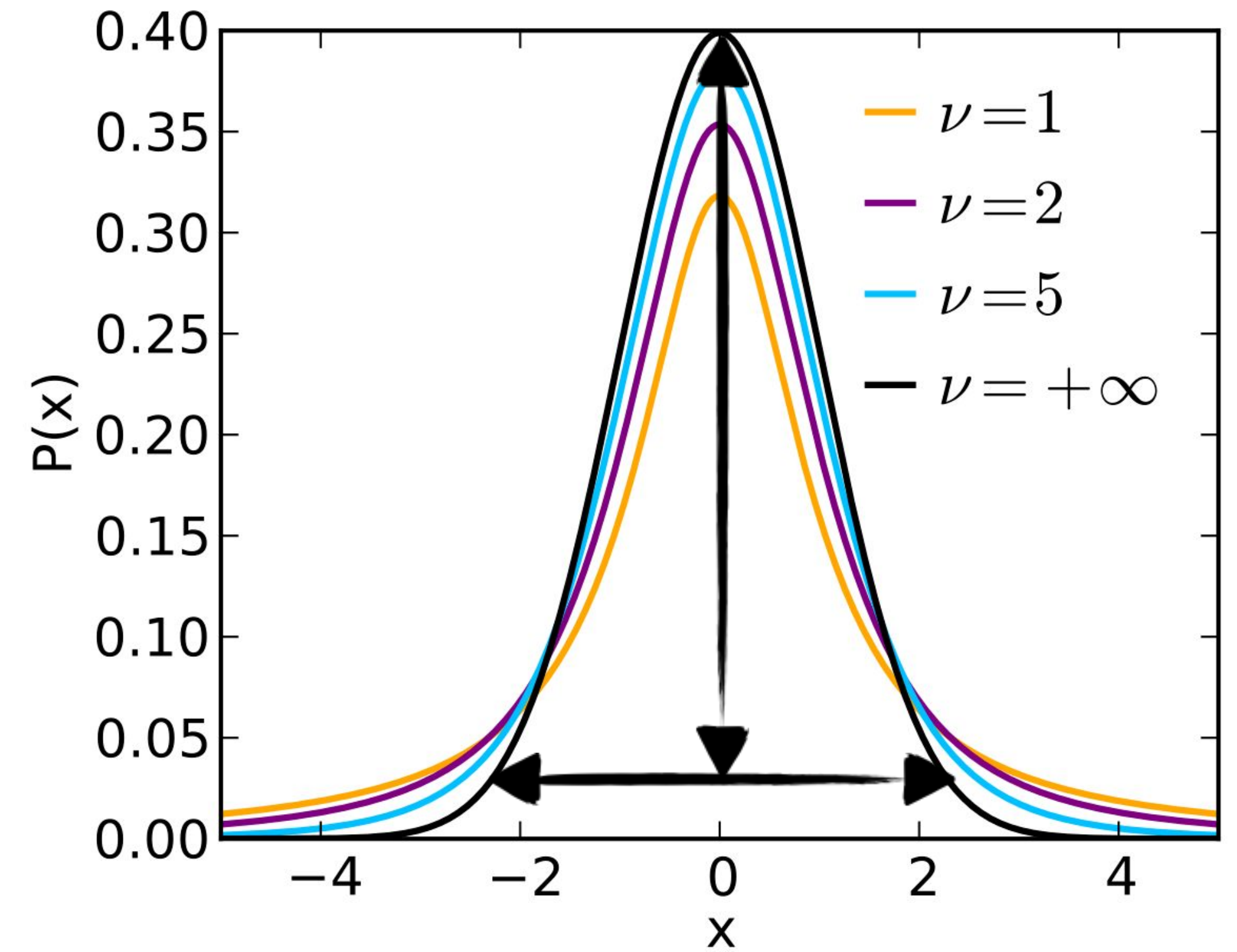■ **increased degrees-of-freedom** increases the probability of the population equaling sample

# more

data



■ ~~increased degrees-of-freedom increases the probability of the population equaling sample~~

■ **more data, better line**

more data **balanced** data normalized data quality data

**balanced** data

# Let's think about logistic functions!

# **balanced** data

| White | Red |
|-------|-----|

**in an ideal world**

**...but no**

# **balanced** data

**4898**          **1599**

White          Red

**What happens when we fit this dataset entirely?**

# **balanced** data

pretty well defined!

Red

White

Acidity

# **balanced** data

much more divergent

pretty well defined!

Red

White

Acidity

■ **balanced data, more accurate results**

more data **balanced** data normalized data **quality** data

more data          balanced data          **normalized** data          quality data

# normalized data

Acidity

Sulfur Dioxide

# **normalized** data

Acidity

**x - 3y > 50**

**True == Red**
**False == White**

**???**

Sulfur Dioxide
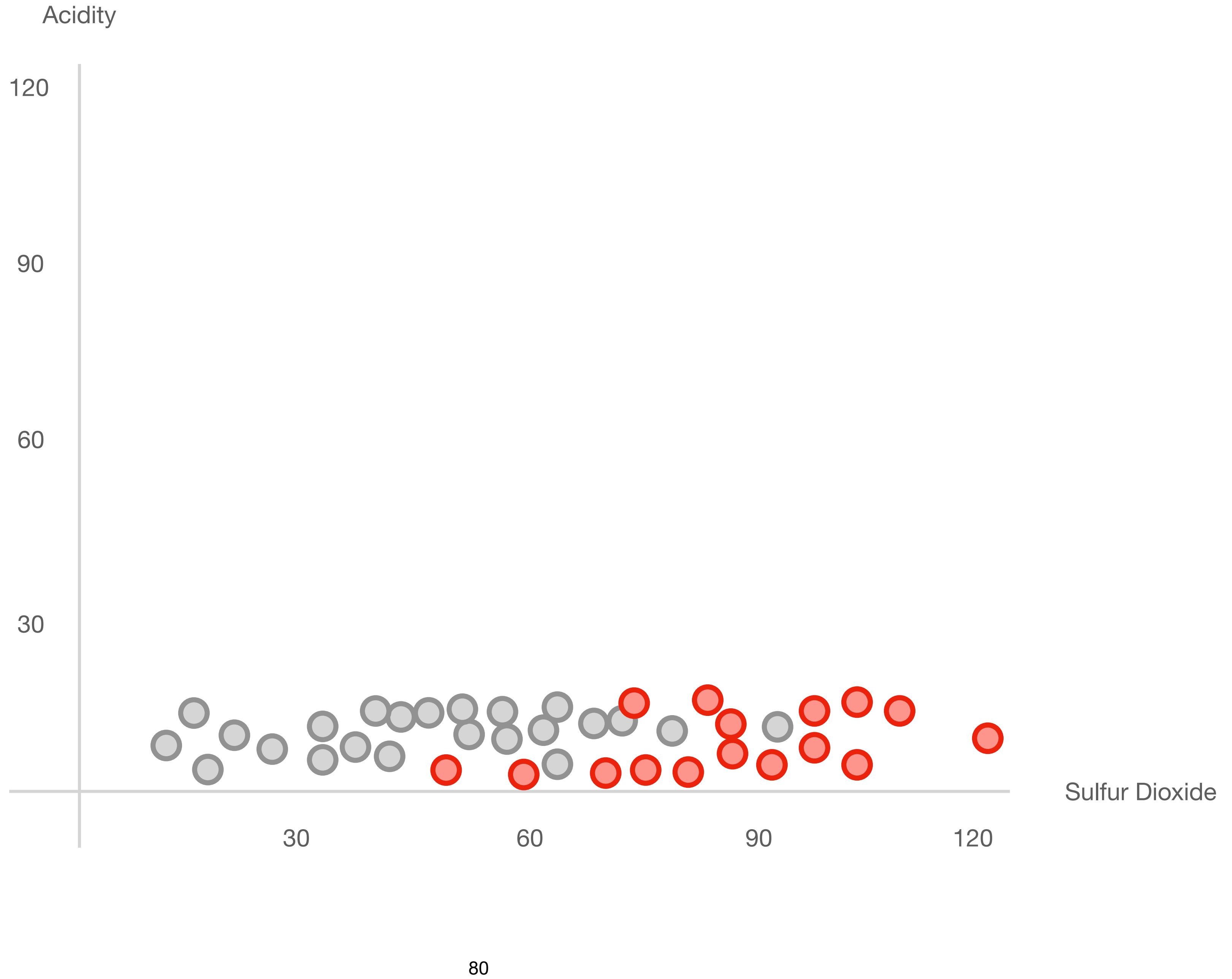
# **normalized** data

# **normalized** data

Acidity

Sulfur Dioxide

1.2x - y > 0.3

much better

■ **normalized data, better generalization, faster convergence**

# **normalized** data



Champagne

Red

White

Acidity

**??? how to fit a line**

84

# **normalized** data



Champagne

Red

White

Acidity

**??? how to fit a line**

# **normalized** data



Champagne

Red

White

Acidity

**??? how to fit a line**

# **normalized** data



Champagne

Red

White

Acidity

**??? how to fit a line**

# **normalized** data



White  Red  Champagne  Acidity

■ non-normalized data is hard to fit

# **normalized** data



■ ensure all features are internally normalized (same order of mag.)

more data          balanced data          **normalized** data          quality data

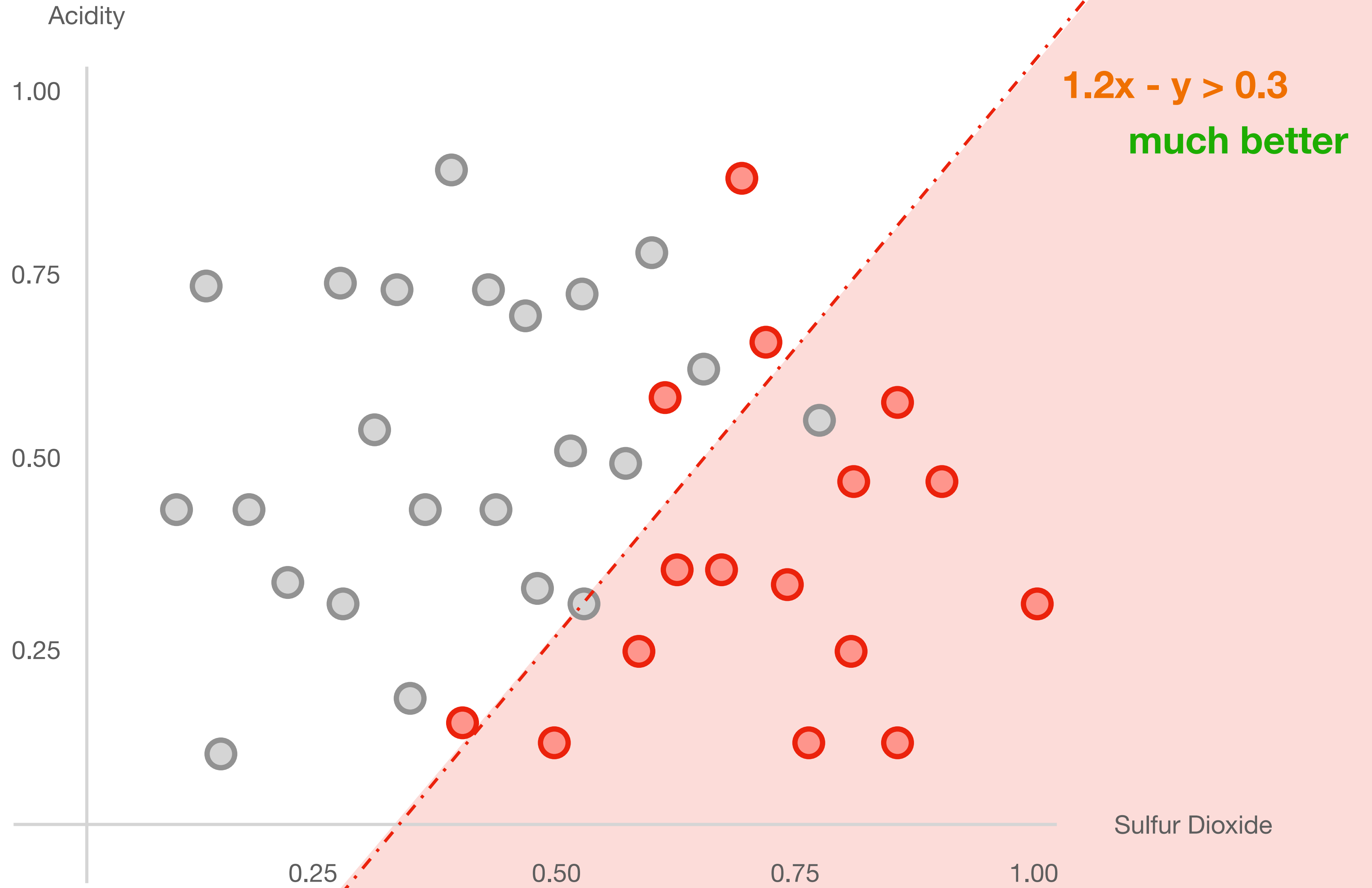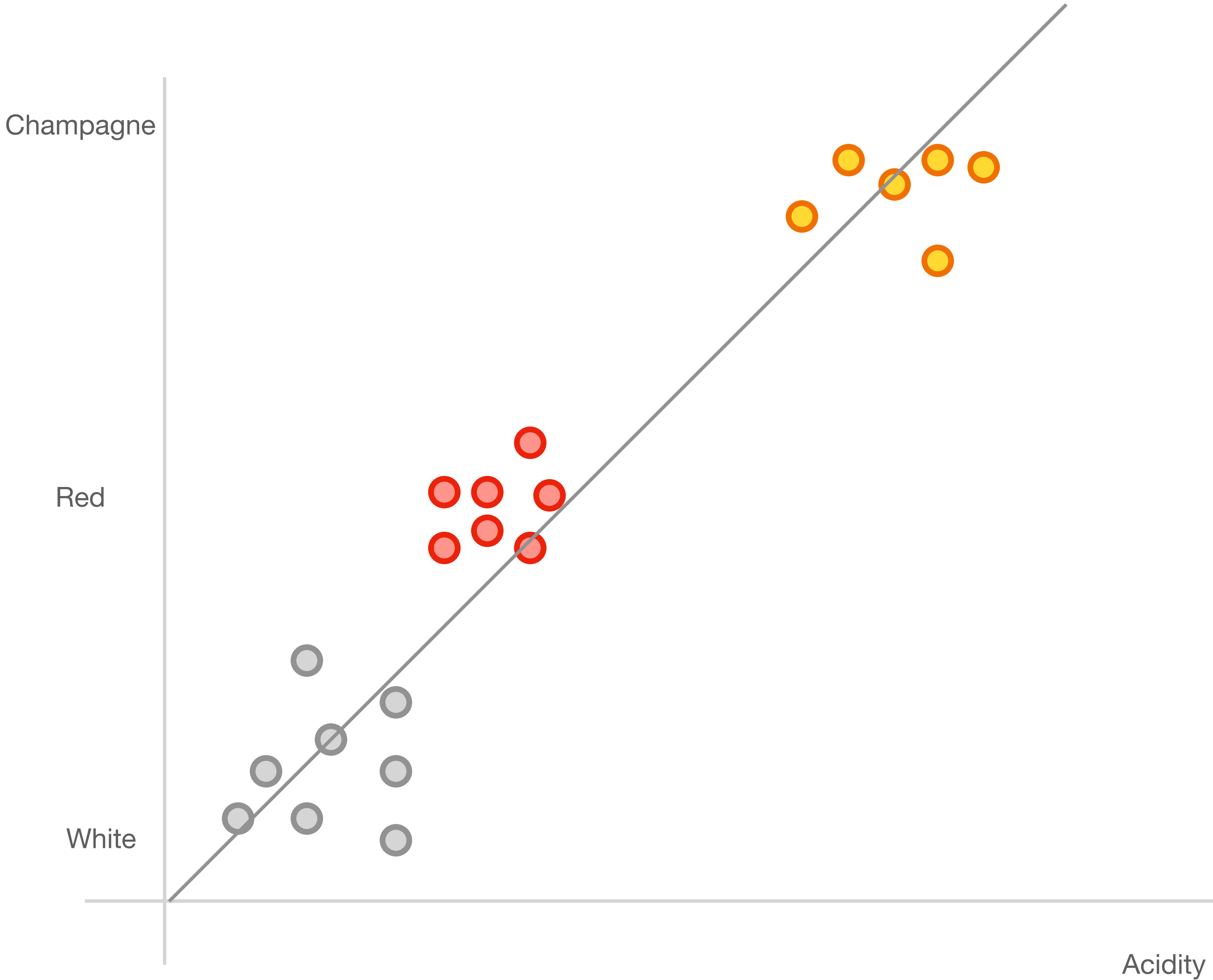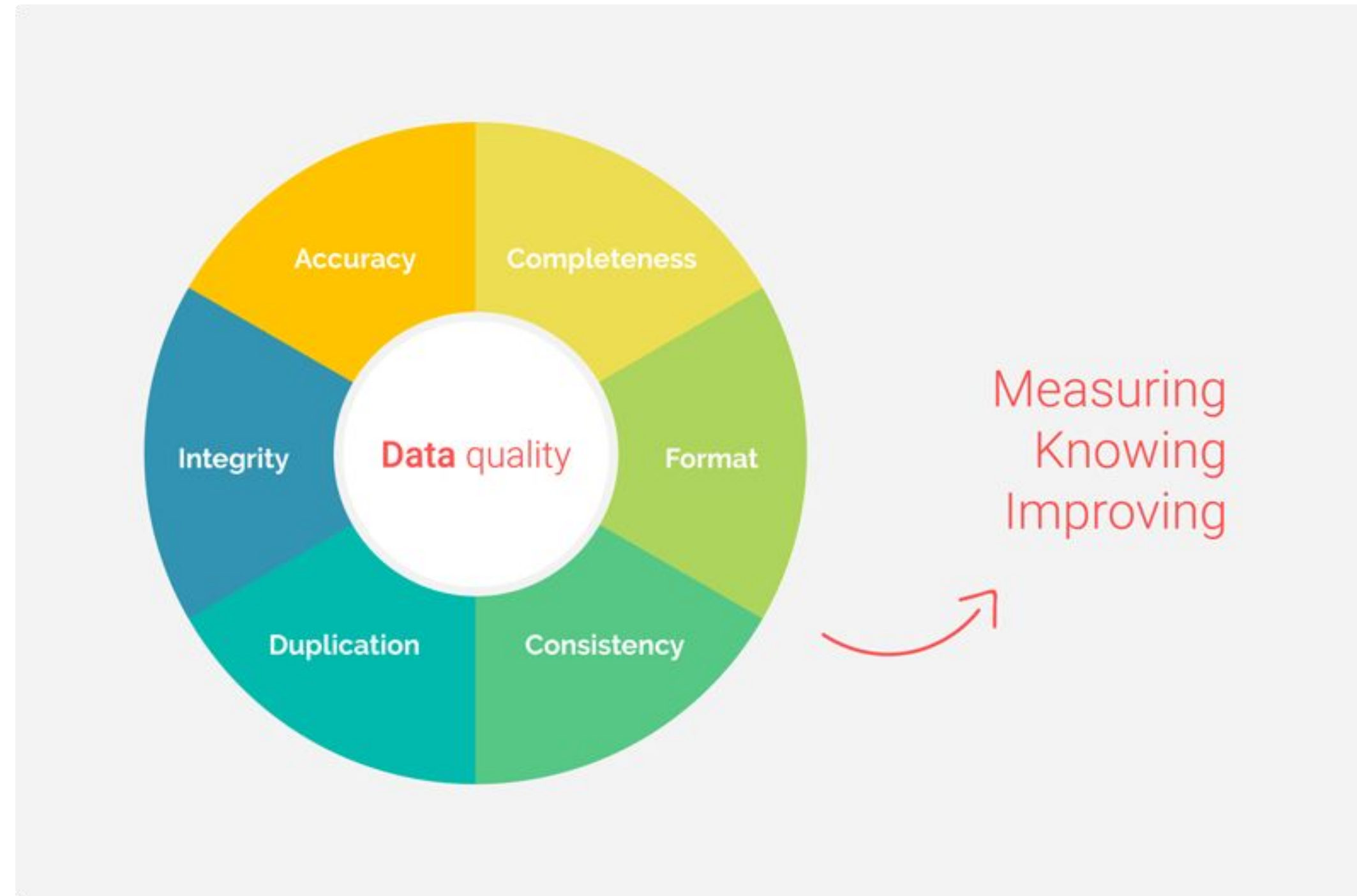more data          balanced data          normalized data          **quality** data

Image credit: Passionned Group

**more** data          **balanced** data          **normalized** data          **quality** data

# Missing Data

# Missing Data

Missing **completely at random**

Missing **at random**

Missing **not at random**

# Missing Data

remove

Use mean/most often

regression

**more** data          **balanced** data          **normalized** data

**let's clean some data!**