

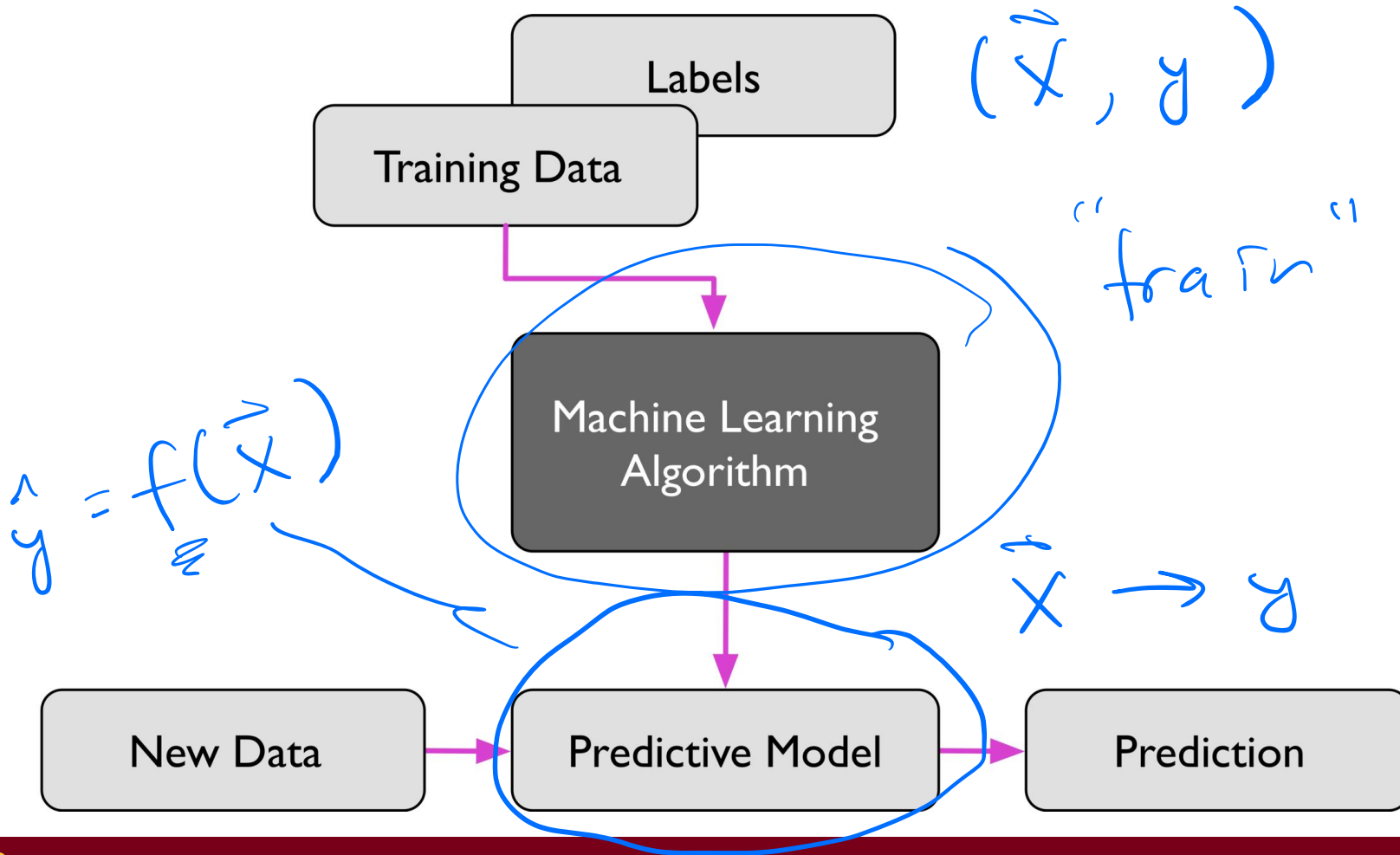
Three core types of ML problems

* Supervised \rightarrow data is labeled (\hat{y}) *
 $\hookrightarrow \hat{y} = f(\vec{x})$

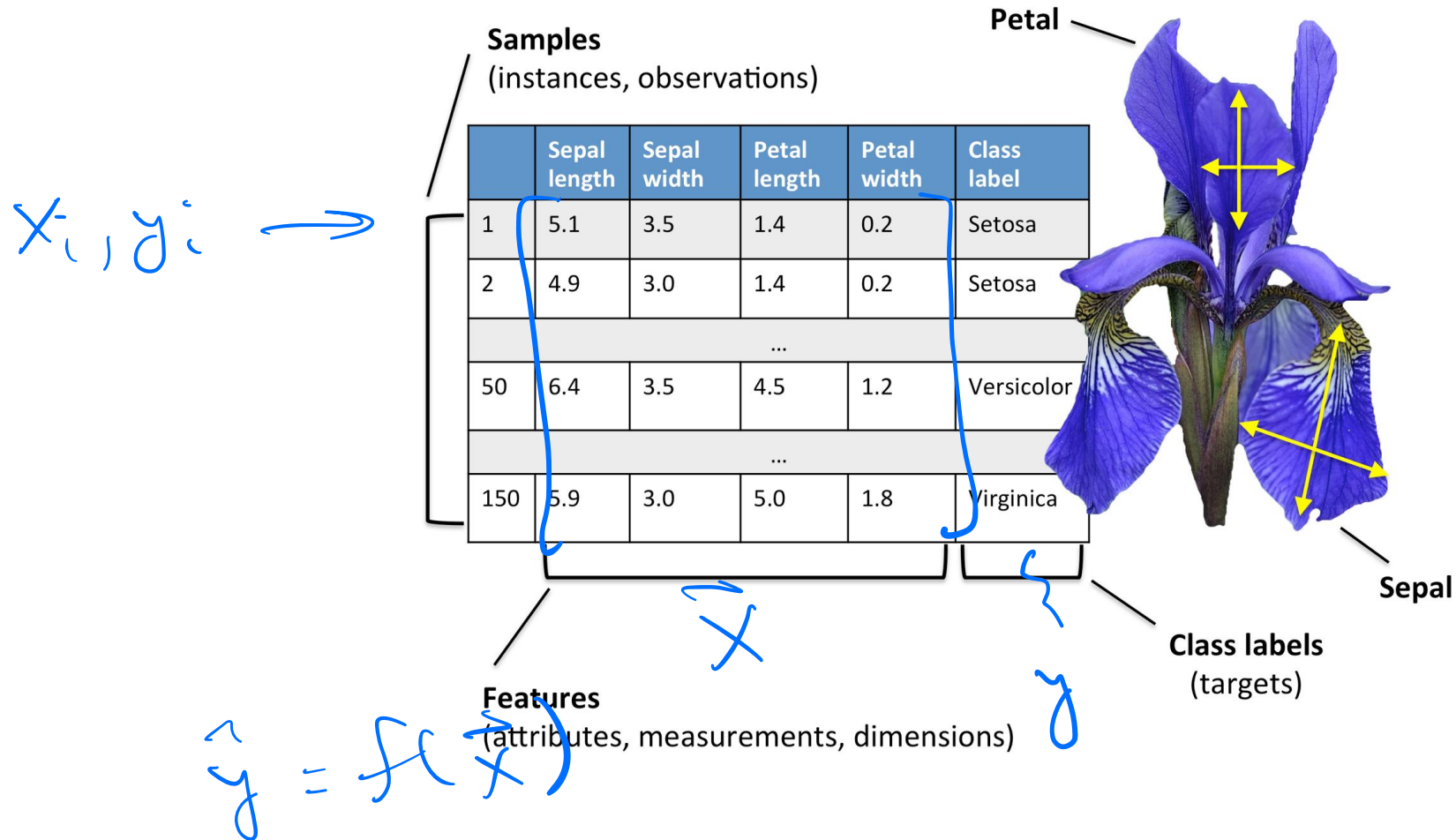
Unsupervised \rightarrow data is unlabeled
 \hookrightarrow find structure in \vec{x}

Reinforcement \rightarrow learn "actions" max
reward.

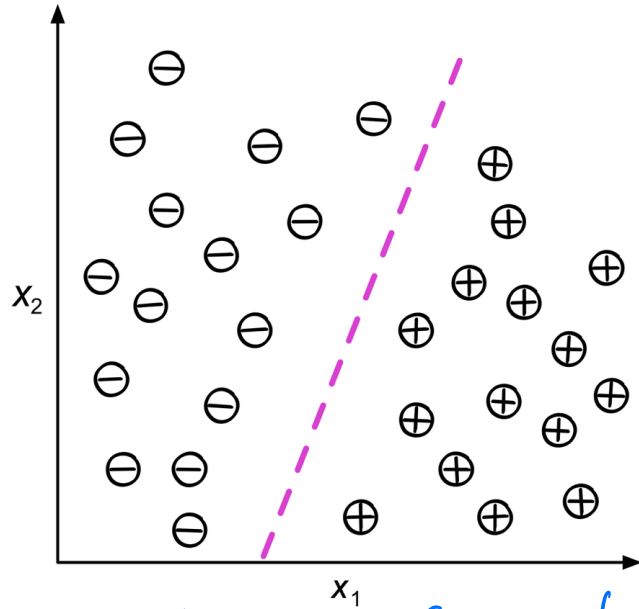
Generic supervised learning problem



Generic supervised learning problem

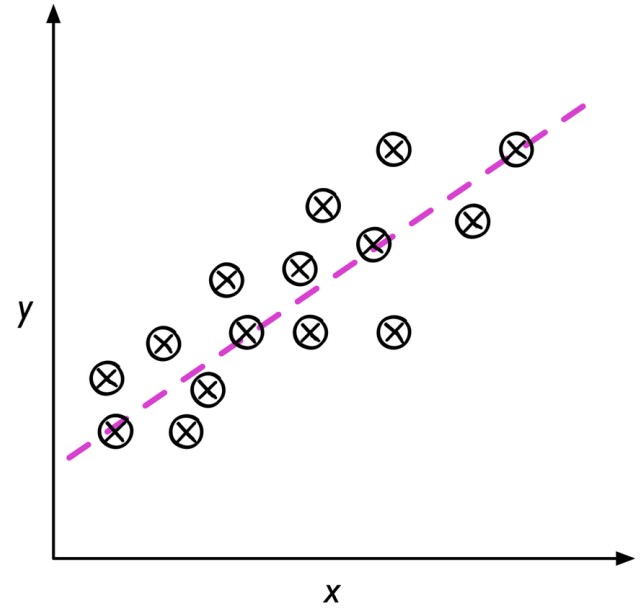


Classification vs regression



Classification
(discrete) \leftarrow

y



\rightarrow Regression
(continuous)

Why ML models fail

$\hat{y} = f(\vec{x}) \hat{=} y$, Failure = \hat{y} is not close enough y

Causes

- ↳ \vec{x} is not expressive enough (bad features)
- ↳ y is not long enough (not enough data)
- ↳ f can't map $\vec{x} \rightarrow y$ (wrong model)
- ↳ \vec{x} and/or y are biased, noisy, incomplete (bad data)

How much data is enough?

#data needed \propto

- 1) how weak \vec{X} is
- 2) how close we need \hat{y} and y to be
- 3) how broad (\vec{X}, y) are

very rough guide

↳ well-structured data + good features

↳ $\sim 10^2 +$

↳ weak structure,

↳ $\sim 10^4 +$

basic features, complex problem



What about features?

What are attributes of **good**/**bad** features?

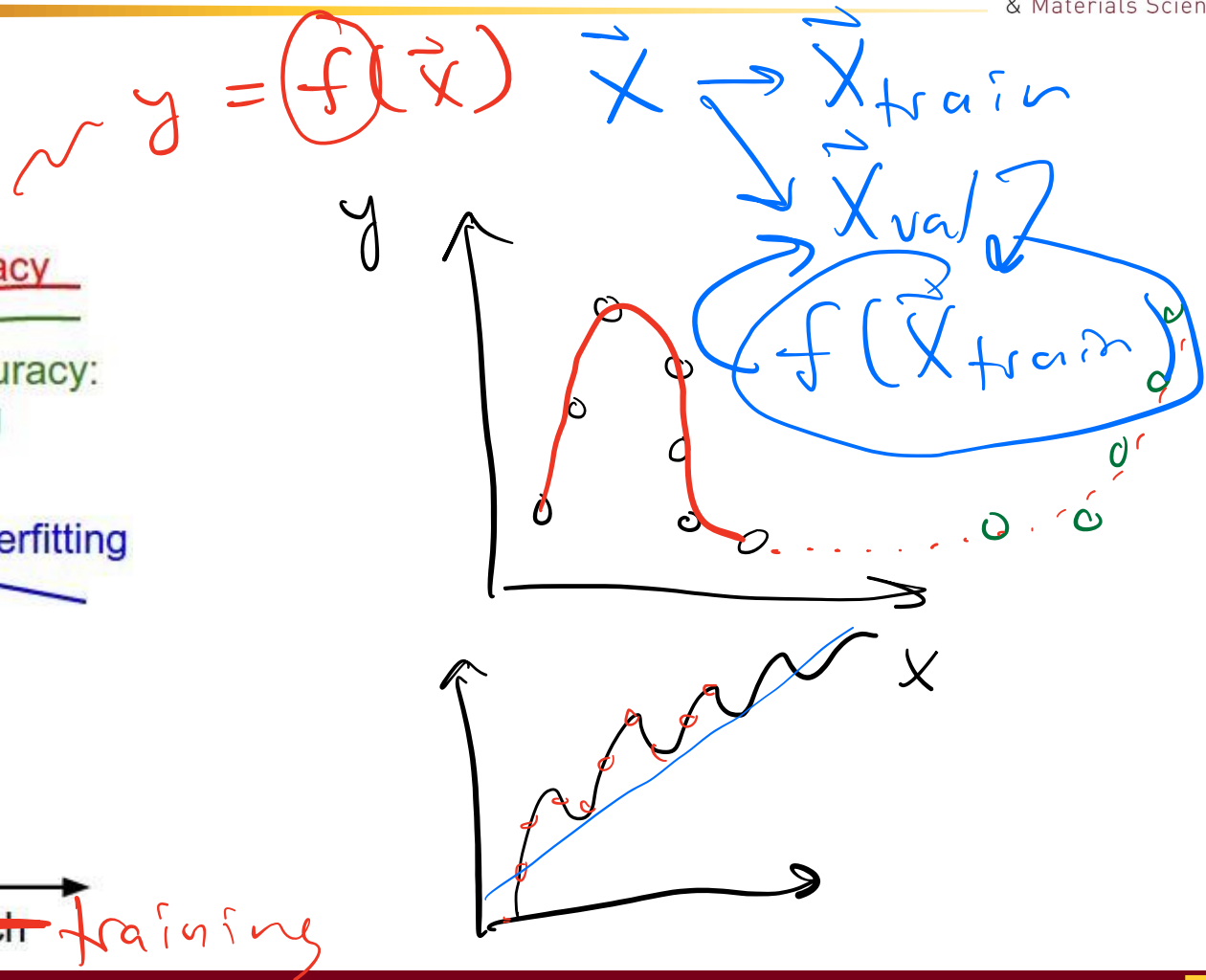
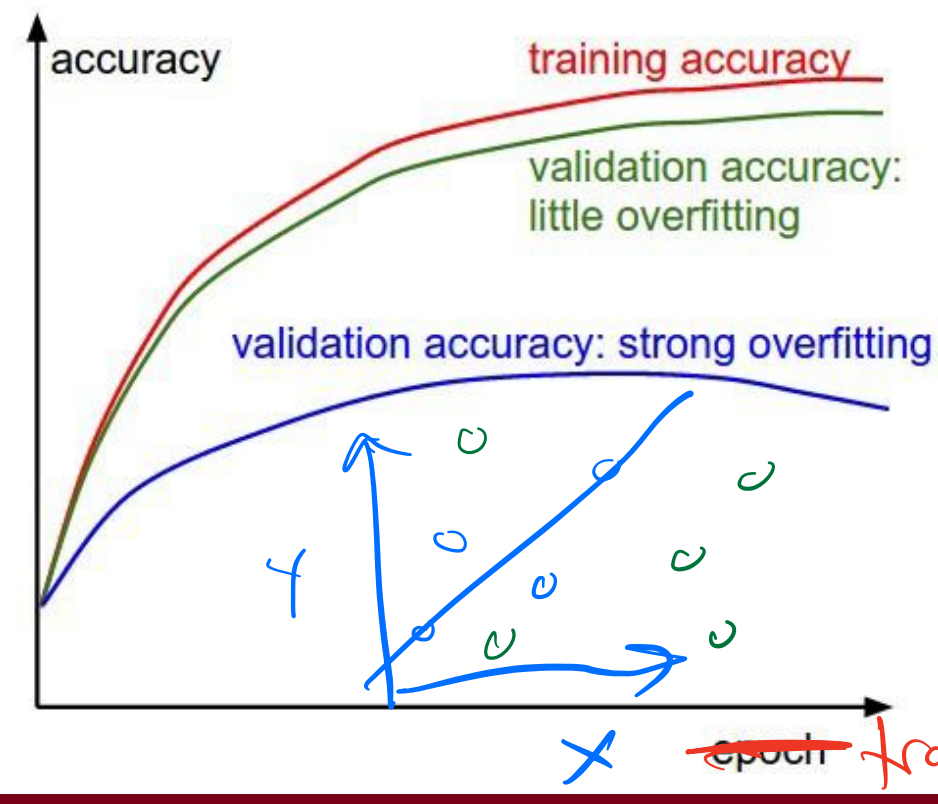
Good

- Strong covariance w/ target.
- Cheap to calculate
- Complementary.
- Interpretable
- Well defined.

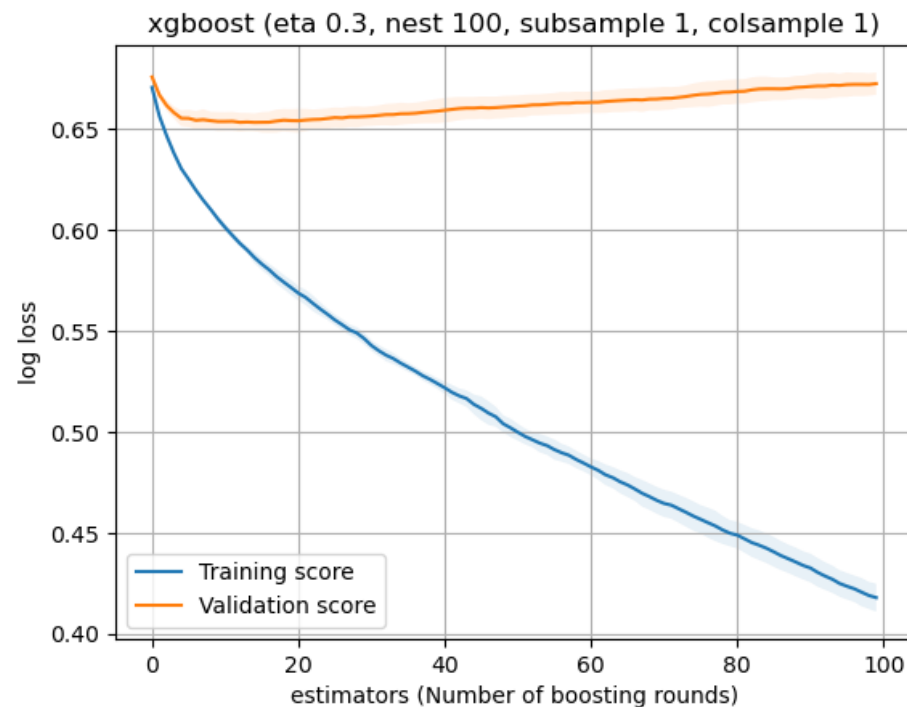
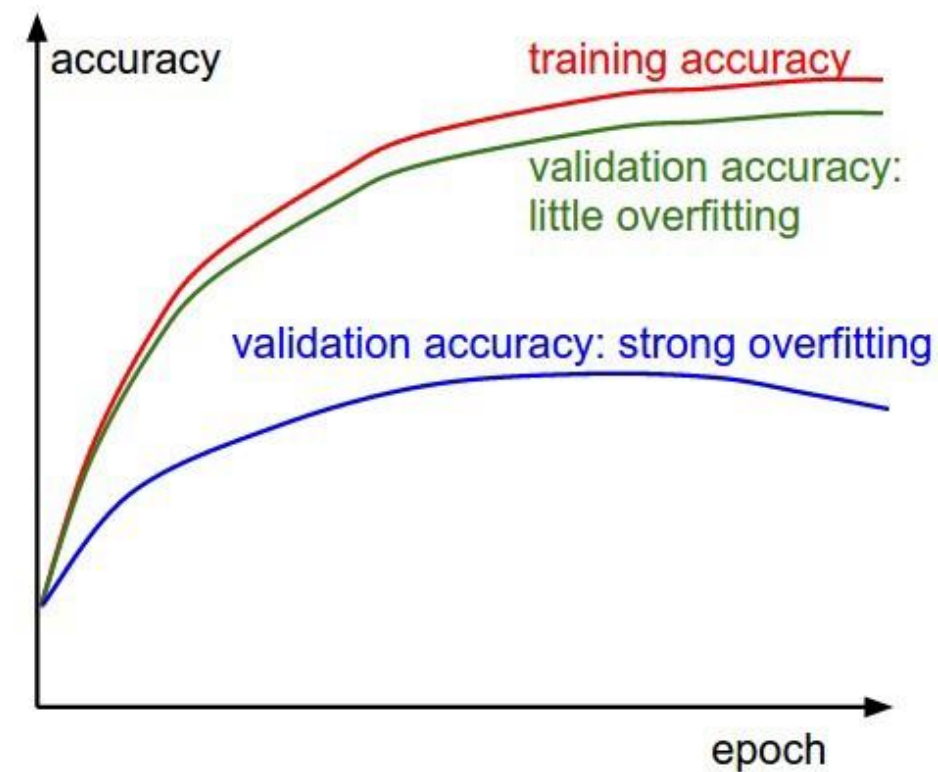
Bad

- (- Irrelevant)
- Expensive
- Correlated
- Sparse
- Opaque

Validation is essential!

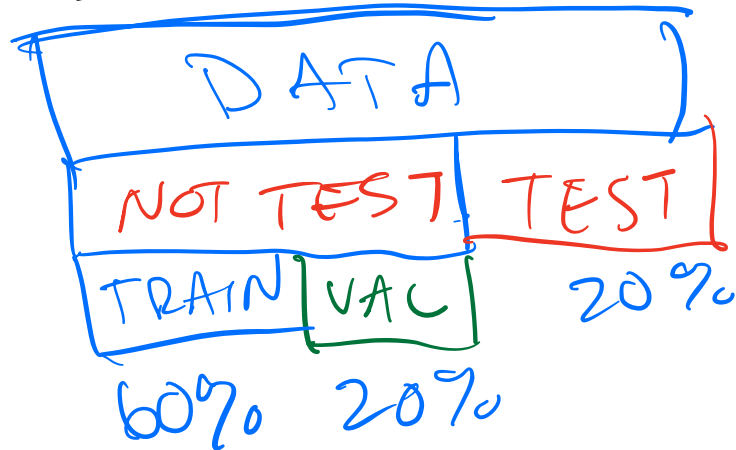


Validation is essential!



How to approach an ML problem

- 1) Research question
- 2) Hypothesis
- 3) Data collection
- 4) Analysis



- 4a) Preprocess data
- 4b) Develop training & validation protocol.
- 4c) Train models
- 4d) Evaluate
- 4e) Test

Domain knowledge (your background) is key!

What is the right property to predict?

What is the right way to assess performance?

Which data points are compatible?

Which features might be informative?

...