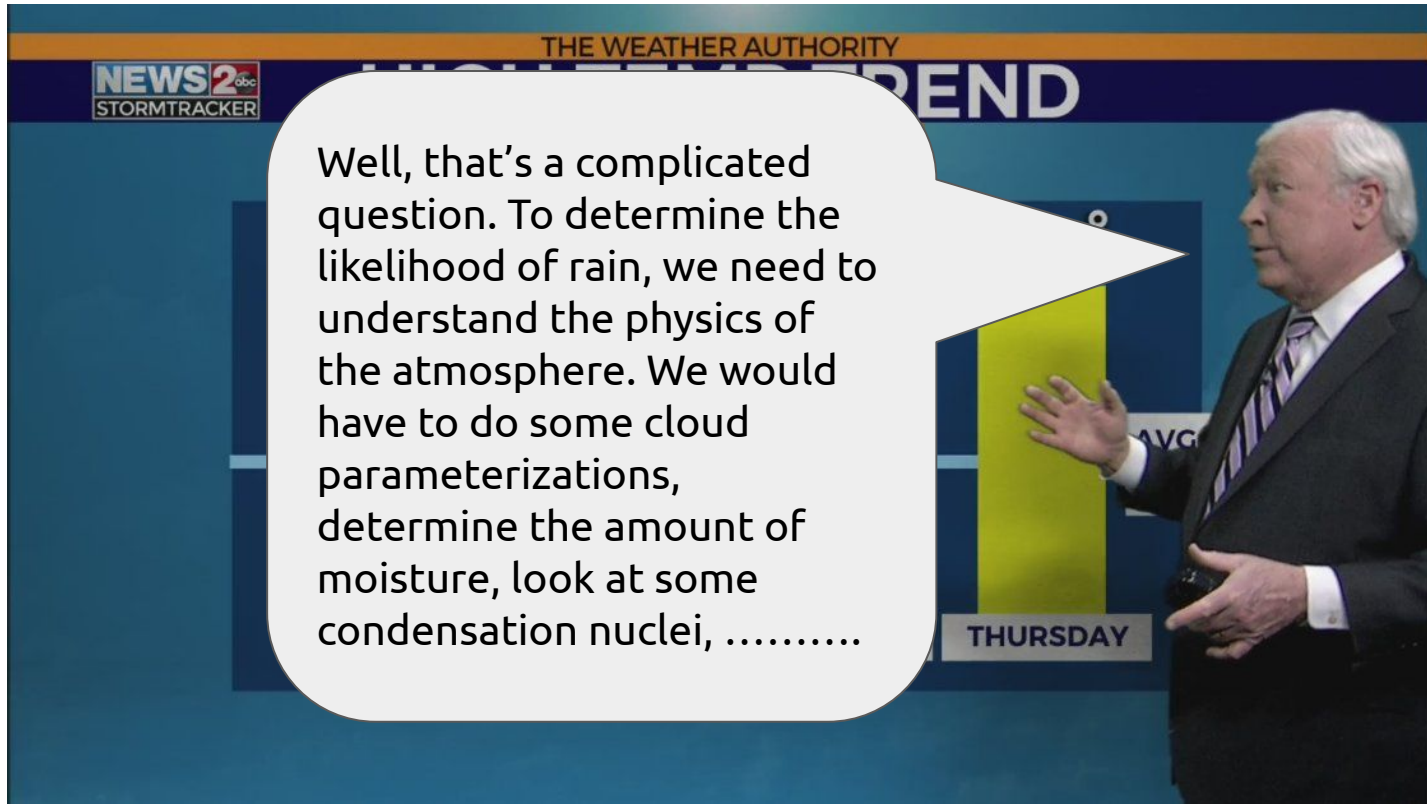# Introduction to Supervised Learning

Question - Is it going to rain today? Do I need my umbrella?

# Question - Is it going to rain today? Do I need my umbrella?

Question - Is it going to rain today? Do I need my umbrella?

You look outside and it looks like this:

Question - Is it going to rain today? Do I need my umbrella?

You look outside and it looks like this:

Question - Is it going to rain today? Do I need my umbrella?

You look outside and it looks like this:

Question - Is it going to rain today? Do I need my umbrella?
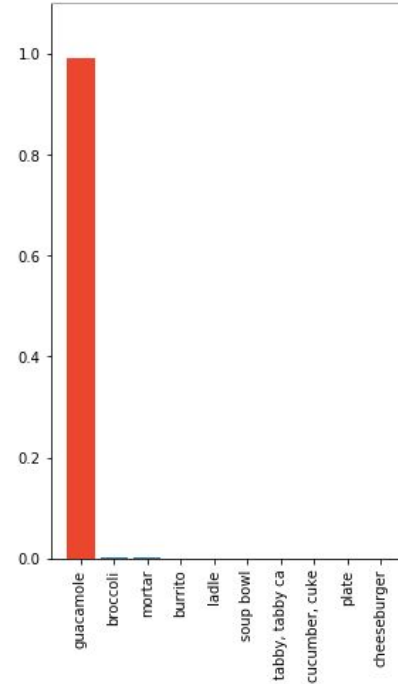
You look outside and it looks like this:



**Moral:** Given a good set of predictor variables and enough "experience" (training data), we can often make good predictions.

Before we get started, a word of caution: machine learning is more than just dumping a lot of data into a magic black box to make predictions.

Before we get started, a word of caution: machine learning is more than just dumping a lot of data into a magic black box to make predictions.

STAT+

# IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal documents show

By Casey Ross ✔ and Ike Swetlitz  July 25, 2018

Reprints

https://www.statnews.com/2018/07/25/ibm-watson-recommended-unsafe-incorrect-treatments/

Before we get started, a word of caution: machine learning is more than just dumping a lot of data into a magic black box to make predictions.
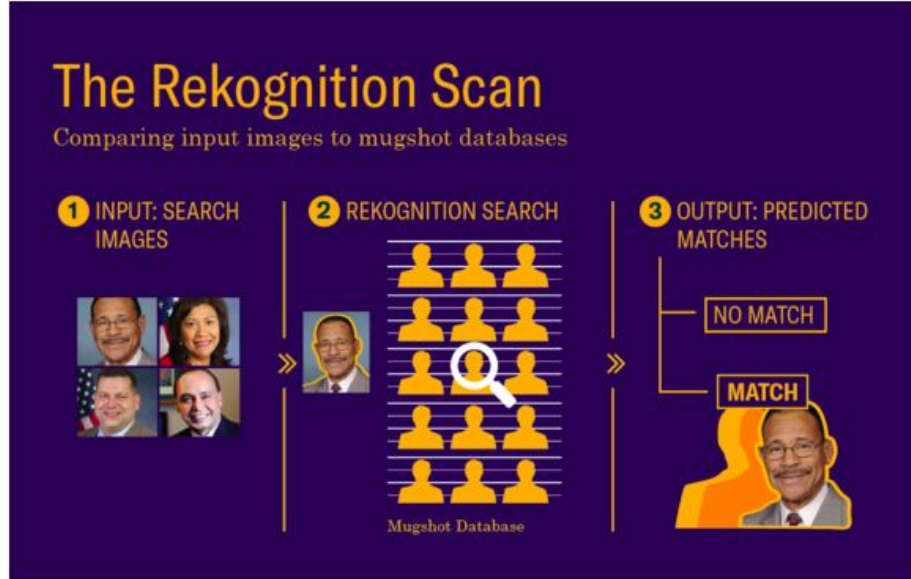
# Amazon ditched AI recruiting tool that favored men for technical jobs

### Specialists had been building computer programs since 2014 to review résumés in an effort to automate the search process
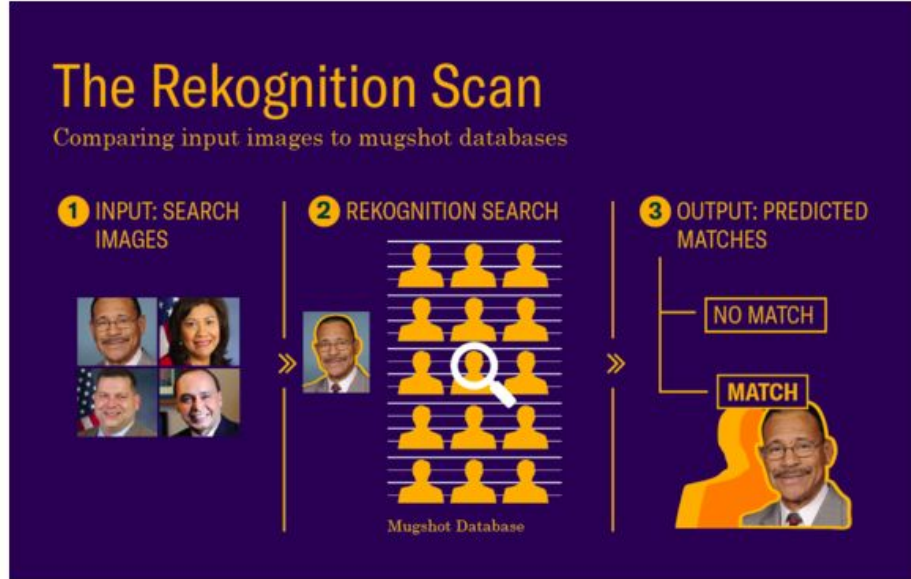
Before we get started, a word of caution: machine learning is more than just dumping a lot of data into a magic black box to make predictions.
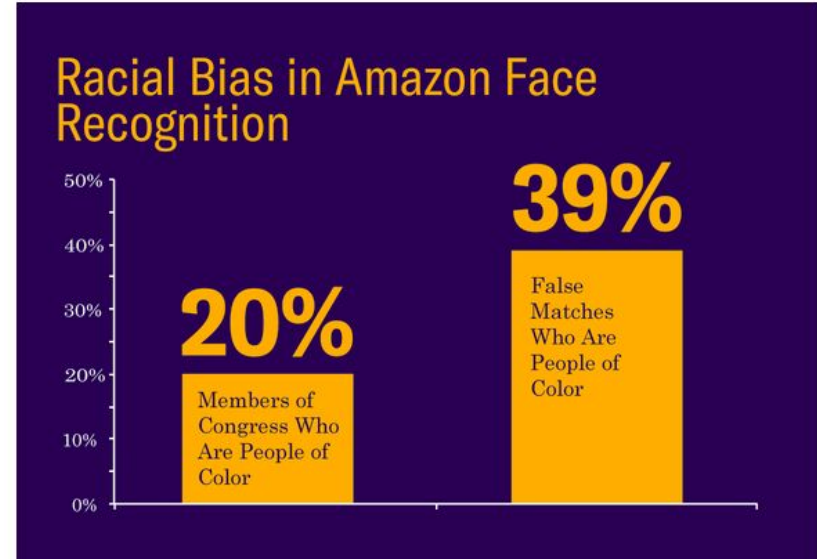


Rep. Sanford Bishop (D-Ga.) was falsely identified by Amazon Rekognition as someone who had been arrested for a crime.

https://www.aclu.org/blog/privacy-technology/surveillance-technologies/amazons-face-recognition-falsely-matched-28

Before we get started, a word of caution: machine learning is more than just dumping a lot of data into a magic black box to make predictions.
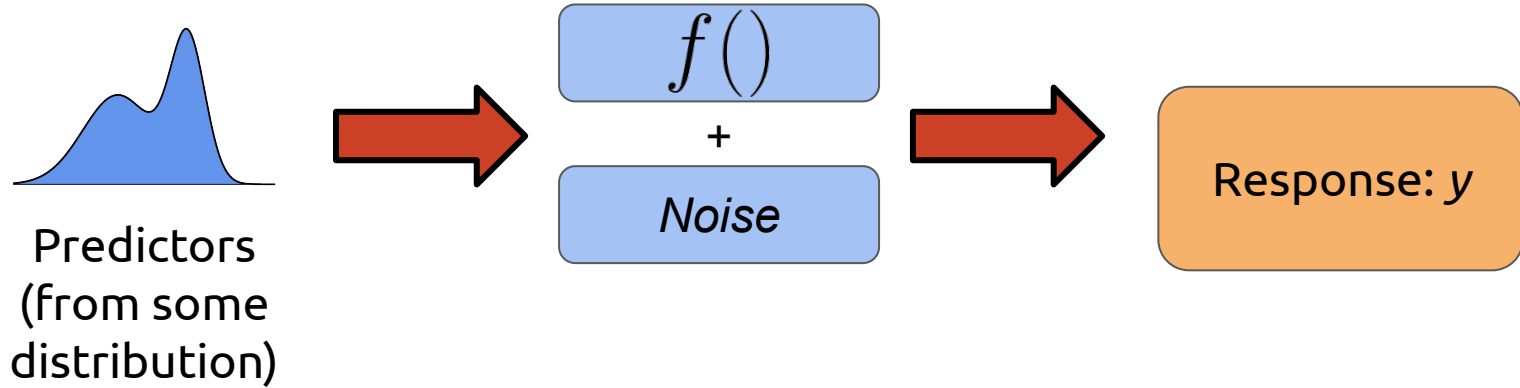


Rep. Sanford Bishop (D-Ga.) was falsely identified by Amazon Rekognition as someone who had been arrested for a crime.
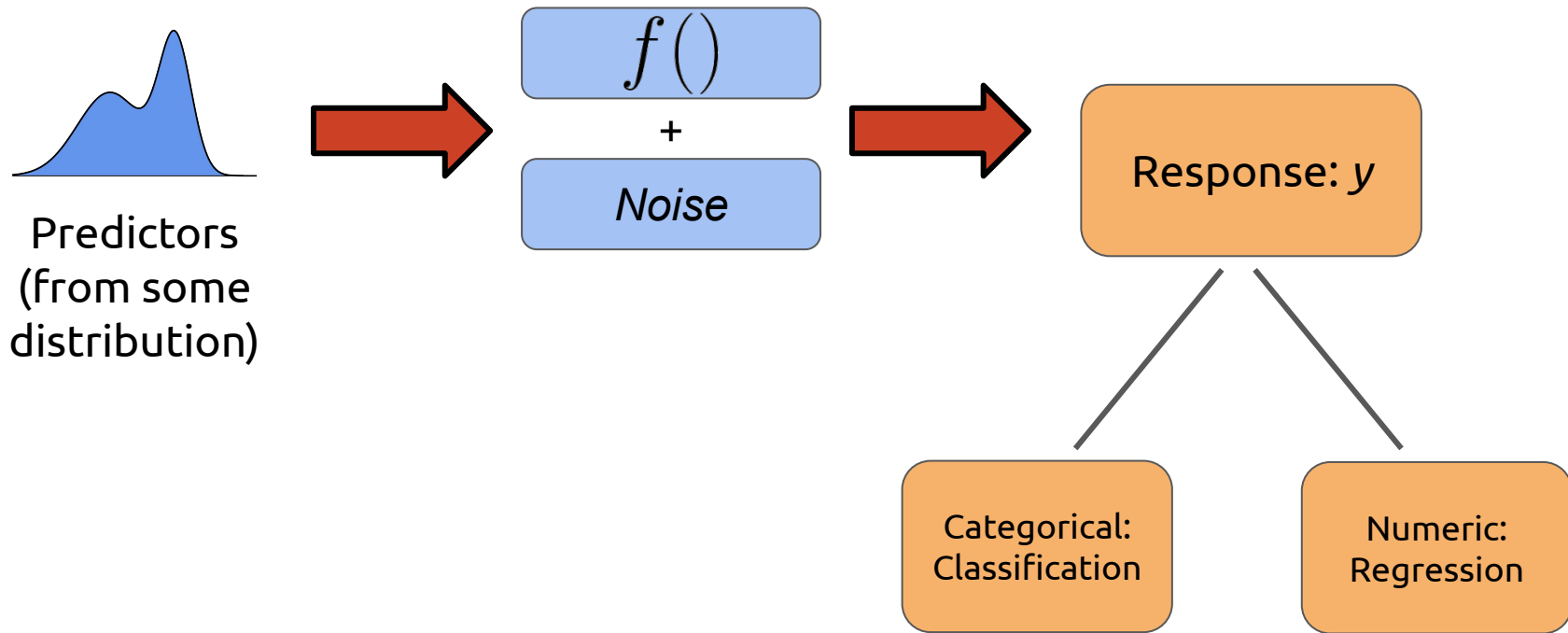
People of color were disproportionately falsely matched in our test.

https://www.aclu.org/blog/privacy-technology/surveillance-technologies/amazons-face-recognition-falsely-matched-28

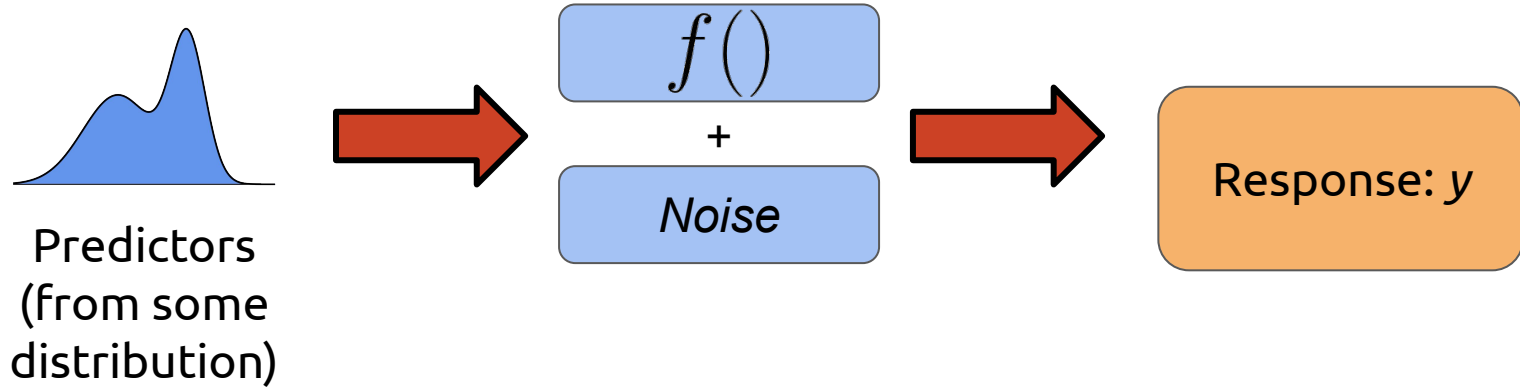# Supervised Learning - Setup



Predictors (from some distribution) → $f() + Noise$ → Response: $y$

# Supervised Learning - Setup



Predictors
(from some
distribution)

$f()$

+

*Noise*

Response: *y*

Categorical:
Classification

Numeric:
Regression

# Supervised Learning - Setup



Predictors
(from some
distribution)

$f()$

+

*Noise*
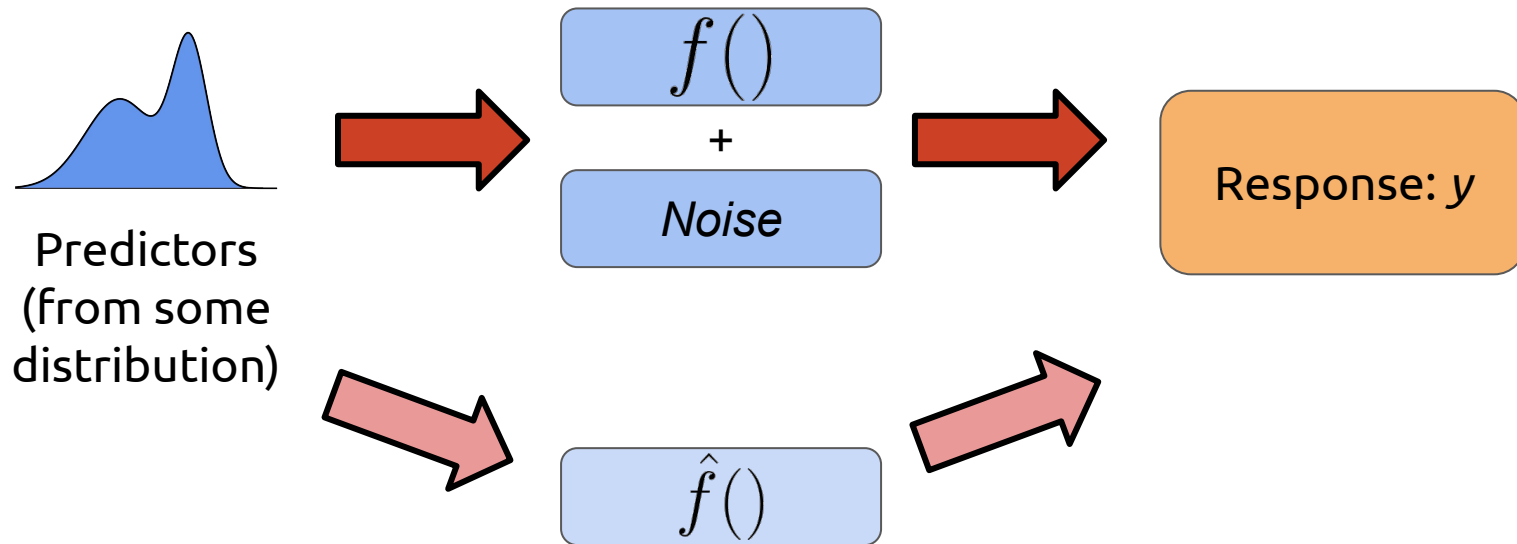
Response: *y*
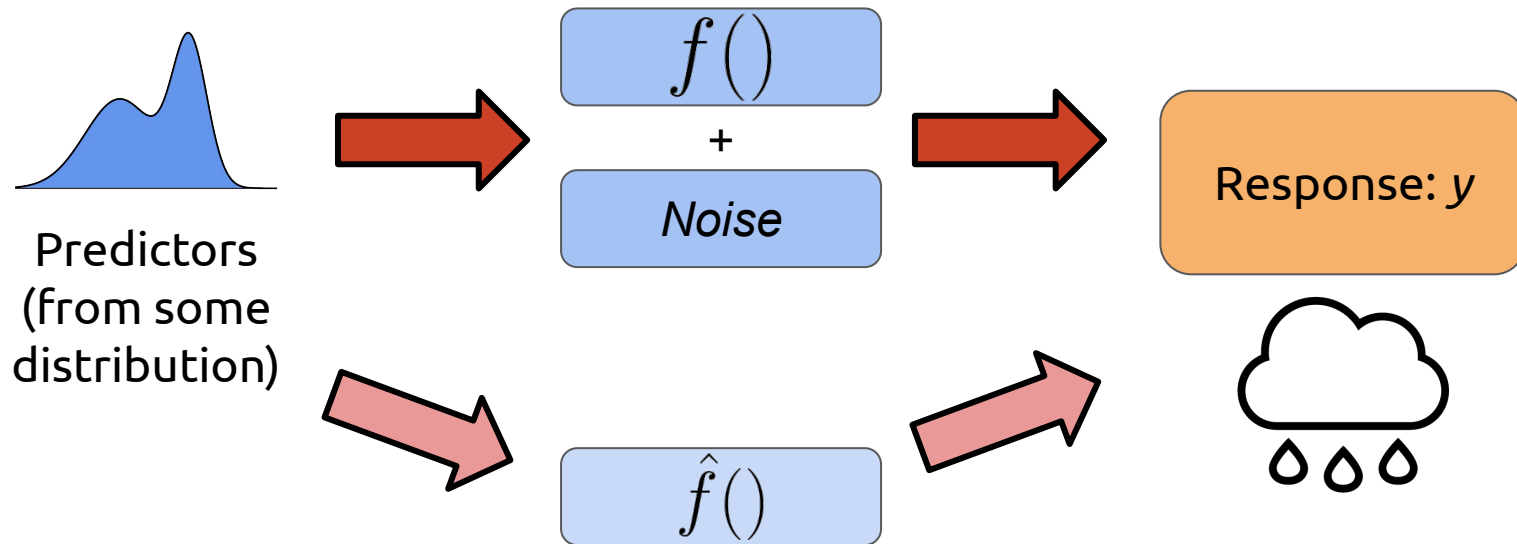
# Supervised Learning - Goals
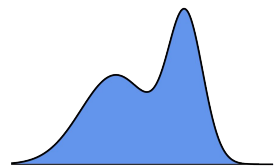
# Supervised Learning - Goals



**Goal:** Choose a function so that the our predictions are close (on average) to the true values.

# Supervised Learning - Grossly Oversimplified

# Supervised Learning - Grossly Oversimplified



Predictors
(from some
distribution)

$$\begin{cases} \partial_t \rho_t + \mathrm{div}\big((\nabla P_t^* - x)^\perp \rho_t\big) = 0 \\ \det(D^2 P_t^*) = \rho_t \end{cases}$$

$$\begin{cases} \partial_t \nabla P_t + (u_t \cdot \nabla)\nabla P_t + (\nabla P_t - x)^\perp = 0 \\ \mathrm{div}\, u_t = 0 \end{cases}$$

Response: $y$

$\hat{f}()$

# Supervised Learning - Grossly Oversimplified

# Example - Weather Prediction



$$f()$$

$$+$$

*Noise*

Response: *y*

Solar Radiation

# Example - Weather Prediction

$$f()$$

$+$

*Noise*

Response: *y*

air temperature
precipitation rate
visibility
wind speed
wind direction
dew point
temperature
air pressure
relative humidity

Solar Radiation

# Example - Weather Prediction



air temperature
precipitation rate
visibility
wind speed
wind direction
dew point
temperature
air pressure
relative humidity

$$f()$$

$+$

*Noise*

$$\hat{f}()$$

Response: *y*

Solar Radiation

# Example - Readmission or Death of Stroke Patients



$f()$

$+$

*Noise*

Response: *y*

90 day readmission or mortality

https://www.mdpi.com/2076-3417/10/18/6337/pdf

# Example - Readmission or Death of Stroke Patients



demographic data
initial vital signs
laboratory results
past medical history
comorbidities
treatment-seeking behavior
pre-stroke functional status

$$f() + \textit{Noise}$$

Response: $y$

90 day readmission or mortality

https://www.mdpi.com/2076-3417/10/18/6337/pdf

# Example - Readmission or Death of Stroke Patients

demographic data
initial vital signs
laboratory results
past medical history
comorbidities
treatment-seeking behavior
pre-stroke functional status

$f()$

$+$

*Noise*

$\hat{f}()$

Response: *y*

90 day
readmission or
mortality

https://www.mdpi.com/2076-3417/10/18/6337/pdf

# Supervised Learning - How

# Supervised Learning - Goals

To measure how "good" our model is, we need some way to measure "error" (eg. mean squared error).

Our goal is to minimize the expected loss over *new* data.

**Important:** We are not trying to minimize loss over the observed data (which is often very easy to do), but to minimize the *generalization error* - the performance on unseen data.

# Measuring Generalization Data - How

If we only care about how well our model performs on unseen data, how do we measure that?

# Measuring Generalization Data - How

If we only care about how well our model performs on unseen data, how do we measure that?

We can't - it's unseen!

# Measuring Generalization Data - How

If we only care about how well our model performs on unseen data, how do we measure that?

We can't - it's unseen!

But, we can *estimate* it.

# Measuring Generalization Error - How

The most simple way to estimate generalization error is through employing a train/test split.

**Full Dataset**

# Measuring Generalization Error - How

The most simple way to estimate generalization error is through employing a train/test split.

# Measuring Generalization Error - How

The most simple way to estimate generalization error is through employing a train/test split.

| Training Data | Test Data |
|---|---|

Build a model on this

# Measuring Generalization Error - How

The most simple way to estimate generalization error is through employing a train/test split.



**Training Data**    **Test Data**

Build a model on this

Estimate generalization error by seeing how it performs on this.