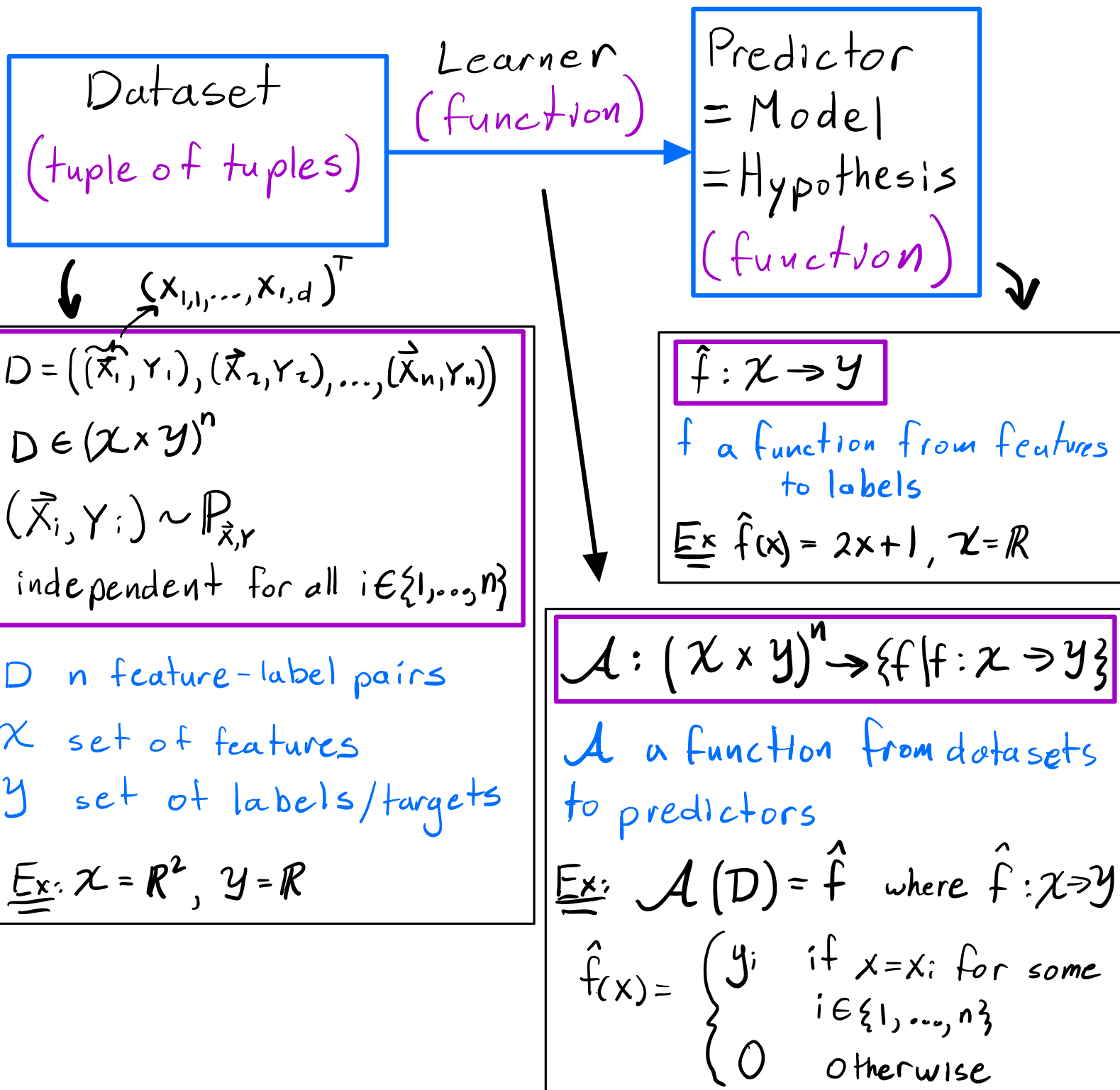


# Supervised Learning: Learning from a randomly sampled batch of labeled data

what does learning well mean?



# Setting

We are given a random dataset

$$D \stackrel{\sim P_D}{=} ((\vec{X}_1, Y_1), \dots, (\vec{X}_n, Y_n)) \in (\mathcal{X} \times \mathcal{Y})^n$$

where  $(\vec{X}_i, Y_i) \sim P_{\vec{X}, Y}$  are independent for all

$\vec{X}_i \in \mathbb{R}^d$  feature vectors  $i \in \{1, \dots, n\}$

$Y_i$  labels, targets

Ex (of features and labels/targets):

$\vec{X}_i \in \mathbb{R}^3$  # of rooms, # of floors, age of a house

$Y_i \in \mathbb{R}$  price

$\vec{X}_i \in \mathbb{R}^2$  amount of chemical 1, amount of chemical 2  
in a wine

$Y_i \in \{0, 1\}$  type of wine

$\vec{X}_i \in \mathbb{R}^{400}$  pixel value of a  $20 \times 20 = 400$  pixel image

$Y_i \in \{\text{cat}, \text{dog}, \text{bird}\}$  type of animal

What is a feature and what is a label is a design choice. Usually a feature is info that is easy to gather. And the label is hard, which is why you want to predict it

## Objective (Informal)

Define a learner  $\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y}\}$

Such that the predictor  $\hat{f}$  is good

where  $\mathcal{A}(D) = \hat{f}$

Ex:  $f(\vec{x})$     suppose  $\vec{x} = 2$  # of rooms

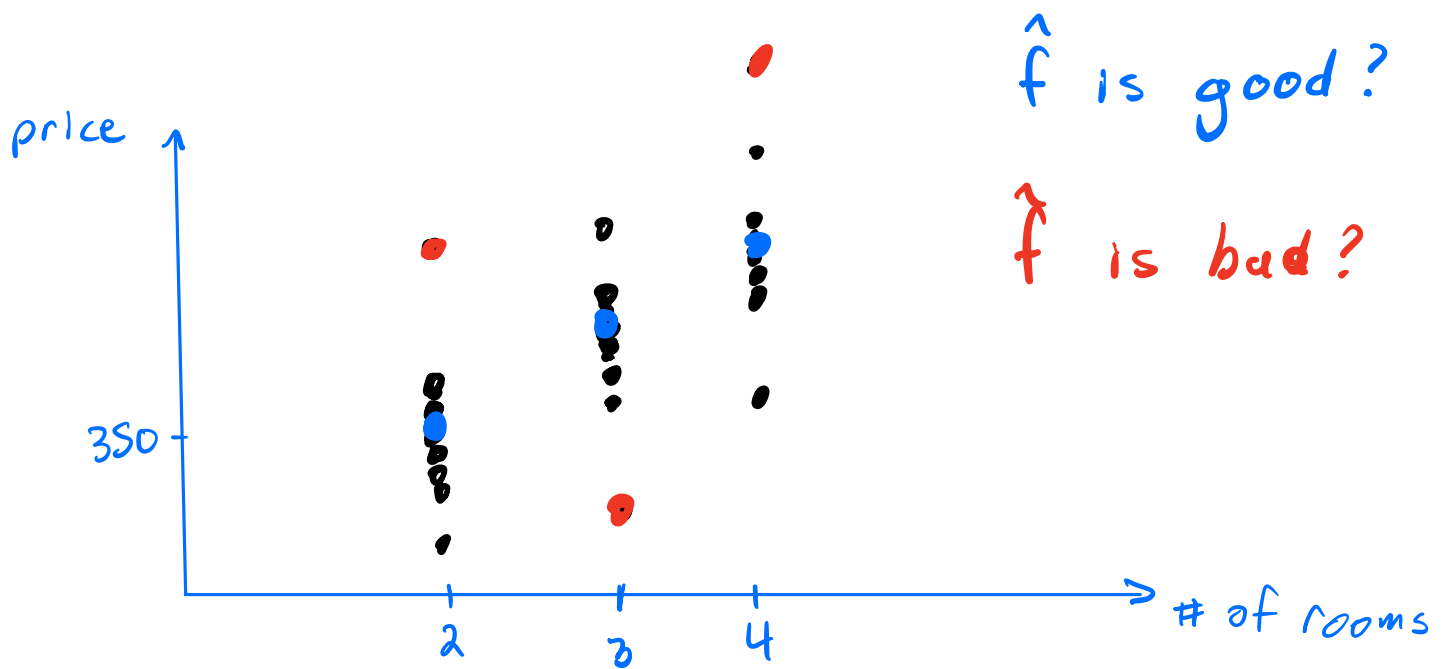
the label is \$300K (predictor  
doesn't know  
this)

$$f(2) = \$300k$$

you get another house with  $\vec{x} = 2$

but the label is \$400k

$$f(2) = \frac{400 + 300}{2} \quad ?$$



We will use a loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$

A good predictor  $\hat{f}$  should have a small loss in expectation

$$L(\hat{f}) = \mathbb{E} [\ell(\hat{f}(\vec{x}), y)] \quad \leftarrow \text{number}$$

$$(\vec{x}, y) \sim P_{\vec{x}, y}$$

"Regression"  $\parallel$   $L(f)$

$p(y|\vec{x}) p(\vec{x})$

$$\begin{aligned} E [\ell(\hat{f}(\vec{x}), Y)] &= \int_{\vec{x}} \int_y \ell(\hat{f}(\vec{x}), y) \underbrace{p(\vec{x}, y)}_{p(y|\vec{x}) p(\vec{x})} dy d\vec{x} \\ &= \int_{\vec{x}} \left( \int_y \ell(\hat{f}(\vec{x}), y) p(y|\vec{x}) dy \right) p(\vec{x}) d\vec{x} \end{aligned}$$

Regression:  $Y \in \mathcal{Y}$  represent something with a notion of order

(Usually  $\mathcal{Y}$  is  $\mathbb{R}$  or some interval)

Ex: house prices, stock prices, energy consumption, weather prediction

We use:

$$\ell(f(\vec{x}), Y) = |f(\vec{x}) - Y| \quad \text{absolute loss}$$

or

$$\ell(f(\vec{x}), Y) = (f(\vec{x}) - Y)^2 \quad \text{squared loss}$$

## Objective (Still Informal)

Define a learner  $\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \Rightarrow \{f \mid f: \mathcal{X} \Rightarrow \mathcal{Y}\}$

Such that the  $L(\hat{f})$  is small

where  $\mathcal{A}(D) = \hat{f}$

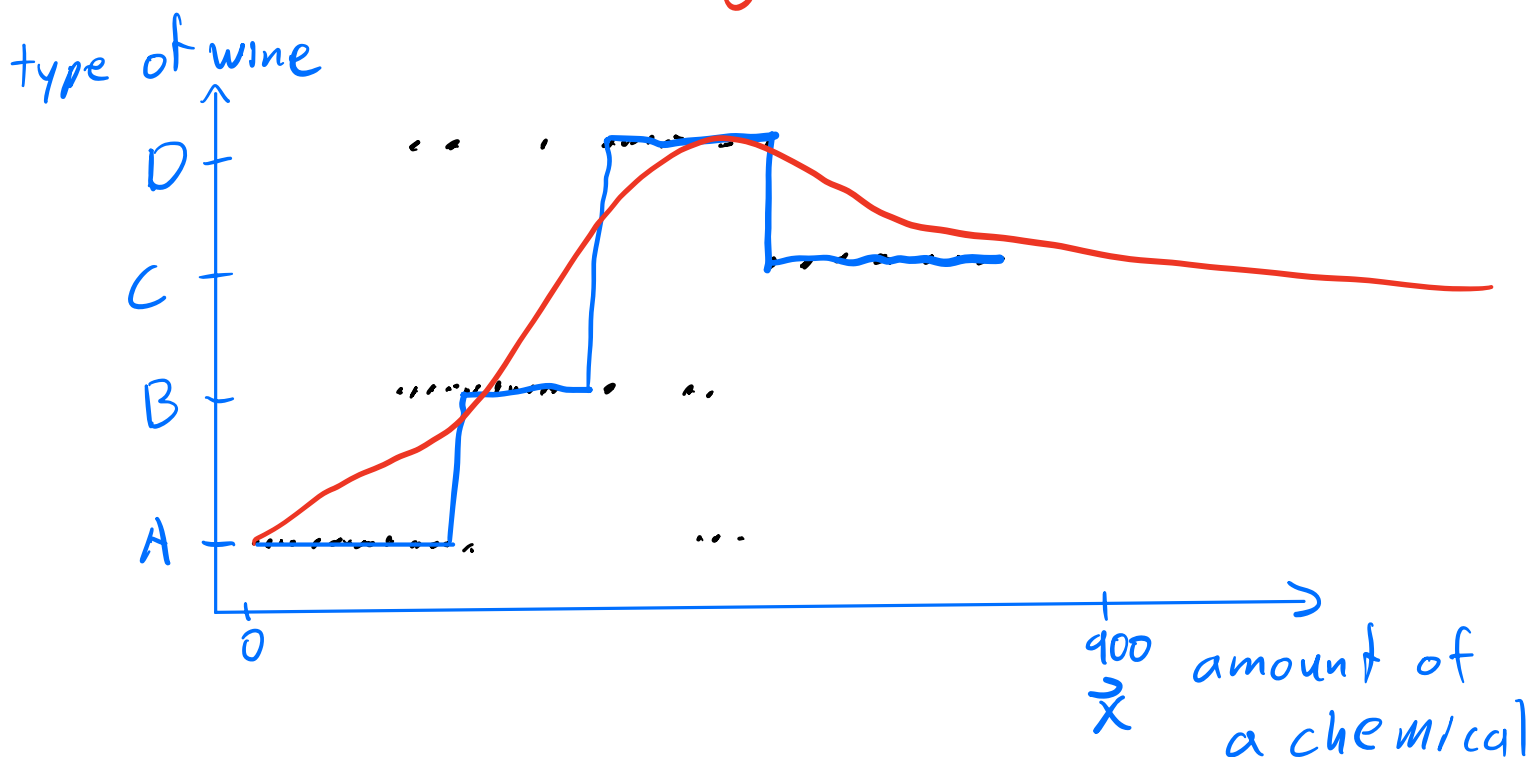
Classification: if  $Y \in \mathcal{Y}$  represents something without order

(Usually  $\mathcal{Y}$  is finite)

Ex(of  $\mathcal{Y}$ ): type wines, type of image, type of email, type of disease

Ex:  $f(\vec{x})$  is a predictor that takes as input the amount of a chemical in a wine and outputs the type of wine

Suppose you got multiple wines, what would a good  $f$  be?



for  $\ell$  we use:

$$\ell(f(\vec{x}), Y) = \begin{cases} 0 & \text{if } f(\vec{x}) = Y \\ 1 & \text{otherwise} \end{cases} \quad \text{0-1 loss}$$

$\mathbb{E}_x$ :  $L(f)$  if we use 0-1 loss  $\mathcal{Y} = \{A, B, C, D\}$

$$\begin{aligned} L(f) &= \mathbb{E}[\ell(f(\vec{x}), Y)] = \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} \ell(f(\vec{x}), y) p(x, y) dx \\ &= \int_{\mathcal{X}} \left( \sum_{y \in \mathcal{Y}} \ell(f(\vec{x}), y) p(y|x) \right) p(x) dx \end{aligned}$$

## Objective (formal)

Define a learner  $\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y}\}$   
such that the  $\mathbb{E}_{\mathcal{D}}[L(\mathcal{A}(\mathcal{D}))]$  is small

ex:  $\mathbb{R}^d$   
↑

For now we will mostly fix a  
dataset  $\mathcal{D}$

$\mathcal{D}$  "random dataset"

## What is the Learner?

We can't explicitly calculate  $L(f)$

for any  $f \in \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y}\}$

because we don't know  $P_{\mathcal{X}, \mathcal{Y}}$

$L(f)$  "risk of  $f$ "

# Defining $A(D)$

## Empirical Risk Minimization (ERM)

Estimation: <sup>function class</sup> use  $D$  to estimate  $L(f)$

for all  $f \in \mathcal{F} \subset \{f | f: \mathcal{X} \Rightarrow \mathcal{Y}\}$

call the estimate  $\hat{L}(f)$

Optimization: pick  $\hat{f}$  to be the  $f \in \mathcal{F}$   
that minimizes  $\hat{L}(f)$



Ex: Let  $\mathcal{F}$  be all linear functions

ERM picks the line that best fits the data

