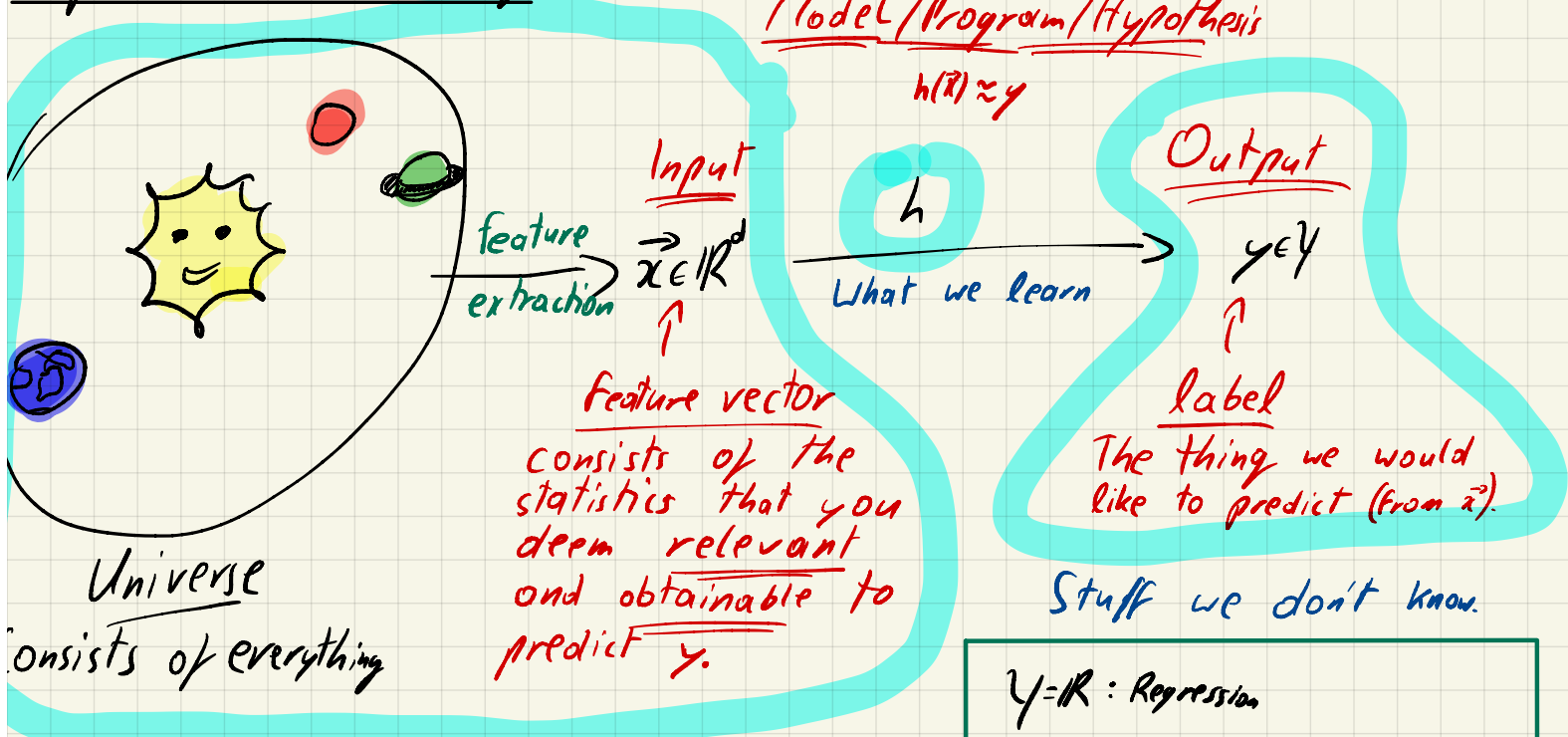


ML Setup

Supervised Learning:



- $\mathcal{Y} = \mathbb{R}$: Regression
- $\mathcal{Y} = \{0, 1\}$: Binary Classification
- $\mathcal{Y} = \{1, 2, \dots, K\}$: Multiclass classification

Stuff we do know

- Examples:
- Predict if the Coca-Cola stock will go up tomorrow.
 - Predict if an email is spam or not.
 - " - if a photo contains a human face.
 - " - what a user said to a home assistant device (e.g. Alexa)

Goal: Learn h from available data.

Ingredients:

<u>labeled Data: D</u>	<u>Hypothesis Class H:</u>	<u>loss function L:</u>	<u>Algorithm A:</u>	<u>Data Scientist:</u>
$D = \{(x_1, y_1), \dots, (x_n, y_n)\} \sim P^n$ where we know $(x_i, y_i) \sim P$ i.i.d. \leftarrow unknown distribution	a set of possible functions $h: X \rightarrow Y$	$L: X \rightarrow \mathbb{R}_0^+$ to tell us how bad any $h \in H$ is	To pick a good $h \in H$ for D	To extract features, to choose pick A

example:

'stop'	'yield'	convolutional neural networks	Cross-entropy loss	Stochastic Gradient Descent	
--------	---------	-------------------------------	--------------------	-----------------------------	--

ML Stages and Concepts:

Learning / Training

Use A to find $h \in \mathcal{H}$ with low loss, $\ell(h)$, on training data \mathcal{D} .

Inference / Testing

For some testing data \tilde{x} , not in the training data, predict the label $y = h(\tilde{x})$

Train and Test data must be drawn i.i.d. from the same distribution \mathcal{P}

Training Data:

Data used to find $h \in \mathcal{H}$

gives rise to the training loss:

$$\frac{1}{|\mathcal{D}|} \sum_{(x,y)} \ell(h(x), y)$$

Testing Data:

Data used to evaluate h .

approximates:

Generalization loss:

$$E_{(x,y) \sim \mathcal{P}} [\ell(x, y)]$$

WLLN

Typically you split your data $\approx 80/20$ into train/test.

(Often people split into train/validation/test. Why?)
80% 10% 10%

Typical ways to split your data: - uniformly at random

Standard rule: Simulate the test case

Never predict the past from the future!

- by time (eg. Jan, Feb, Mar | Apr)
- by patient / instance

Assumptions:

No free lunch Theorem: You must make assumptions in order to learn.

\Rightarrow there is no single ML algorithm that works for all settings.

Example assumptions: - data is locally smooth
- "-" consists of natural images
- \mathcal{P} is a mixture of Gaussians
- Features are independent given the label

