# Gems of TCS

## PAC Learning

Sasha Golovnev

November 29, 2021

# CLASSES OF LEARNING PROBLEMS

- Classification

# Classes of Learning Problems

- Classification

- Ranking

# CLASSES OF LEARNING PROBLEMS

- Classification

- Ranking

- Regression

# Classes of Learning Problems

- Classification

- Ranking

- Regression

- Clustering

# CLASSES OF LEARNING PROBLEMS

- Classification

- Ranking

- Regression

- Clustering

- . . .

- Given a set of <u>labeled</u> emails

- Given a set of <u>labeled</u> emails

- Build a classifier that predicts
  spam/non-spam labels for incoming emails

# SETUP

- Partition labeled data into three sets:
    - training sample
    - validation sample
    - test sample

# Setup

- Partition labeled data into three sets:
  - training sample
  - validation sample
  - test sample
- Identify relevant features

# Setup

- Partition labeled data into three sets:
  - training sample
  - validation sample
  - test sample
- Identify relevant features
- Train on training sample

# Setup

- Partition labeled data into three sets:
    - training sample
    - validation sample
    - test sample
- Identify relevant features
- Train on training sample
- Tune parameters using validation sample

# Setup

- Partition labeled data into three sets:
  - training sample
  - validation sample
  - test sample
- Identify relevant features
- Train on training sample
- Tune parameters using validation sample
- Evaluate using test sample

# What can be learned?

- What can be learned?

- What cannot be learned?

- How many samples do we need to learn?

# WHAT CAN BE LEARNED?

- What can be learned?

- What cannot be learned?

- How many samples do we need to learn?

- Framework of PAC learning (L. Valiant, 1984)

- *X*—set of all possible instances/examples

- $X$—set of all possible instances/examples

- $\mathcal{D}$—target distribution over $X$

# Definitions

- $X$—set of all possible instances/examples

- $\mathcal{D}$—target distribution over $X$

- $c$—target concept

# DEFINITIONS

- $X$—set of all possible instances/examples

- $\mathcal{D}$—target distribution over $X$

- $c$—target concept

- $C$—concept class

# Definitions

- $X$—set of all possible instances/examples

- $\mathcal{D}$—target distribution over $X$

- $c$—target concept

- $C$—concept class

- Goal: given training set, select $h$ that approximates $c$ well

### Generalization Error

For hypothesis $h$, target concept $c$, and target distribution $D$:

$$R(h) = \Pr_{x \sim D}[h(x) \neq c(x)].$$

# ERRORS

## Generalization Error

For hypothesis $h$, target concept $c$, and target distribution $D$:

$$R(h) = \Pr_{x \sim D}[h(x) \neq c(x)].$$

## Empirical Error

For hypothesis $h$, target concept $c$, and sample $S = (x_1, \ldots, x_m)$:

$$\widehat{R}(h) = \frac{|\{x_i : h(x_i) \neq c(x_i)\}|}{m}.$$

$$\mathbb{E}_S[\widehat{R}(h)] = R(h)\,.$$

## PAC (Probably Approximately Correct)

Concept class $C$ is PAC-learnable if there exists learning algorithm s.t.

- for all $c \in C, \varepsilon > 0, \delta > 0$, all distributions $D$:

$$\Pr_{S \sim D^m}[R(h_S) \leq \varepsilon] \geq 1 - \delta \,,$$

## PAC (Probably Approximately Correct)

Concept class $C$ is PAC-learnable if there exists learning algorithm s.t.

- for all $c \in C, \varepsilon > 0, \delta > 0$, all distributions $D$:

$$\Pr_{S \sim D^m} [R(h_S) \leq \varepsilon] \geq 1 - \delta \,,$$

- for random samples of size

$$m \leq \text{poly}(1/\varepsilon, 1/\delta, n) \,.$$

- Probably: confidence $1 - \delta$
- Approximately correct: accuracy $1 - \varepsilon$

# EXAMPLE