

CSE 541: Interactive Learning

Kevin Jamieson





CSE 541, Spring 2025 Interactive Learning

Lecture: Wednesday, Friday 10:00-11:20, [ECE 003](#)

Instructor: [Professor Kevin Jamieson](#)

Contact: cse541-staff@cs.washington.edu

TA office hours:

- Zhihan Xiong: Thursday 4:00-5:00, remote

Instructor office hours:

- Kevin Jamieson: Tuesday 11:00-12:00, CSE2 340

Grading and Evaluation

There will be 3 homeworks (each worth 20%) and a project to be completed in the last few weeks of the class (details forthcoming).

We will cover selected topics from [SzepesvariLattimore]:

- (Non)-stochastic Online learning
- (Non)-stochastic Multi-armed Bandits
- (Non)-stochastic Linear Bandits and experimental design
- (Non)-stochastic Contextual bandits (model-free and model-based)

Prerequisites: The course will make frequent references to introductory concepts of machine learning (e.g., CSE 446/546) but it is not a prerequisite. However, fluency in basic concepts from linear algebra, statistics, and calculus will be assumed (see HW0). Some review materials:

- [Linear Algebra Review](#) by Zico Kolter and Chuong Do.
- [Linear Algebra](#), David Cherney, Tom Denton, Rohit Thomas and Andrew Waldron. Introductory linear algebra text.
- [Probability Review](#) by Arian Maleki and Tom Do. Also see Chapter 5 of [SzepesvariLattimore] below.

The course will be analysis heavy, with a focus on methods that work well in practice. You are strongly encouraged to complete the [self-test](#) of fundamental prerequisites on your own (not to be turned in or graded). You should be able to complete most of these in your head or with minimal computation.

Class materials

The course will pull from textbooks and course notes.

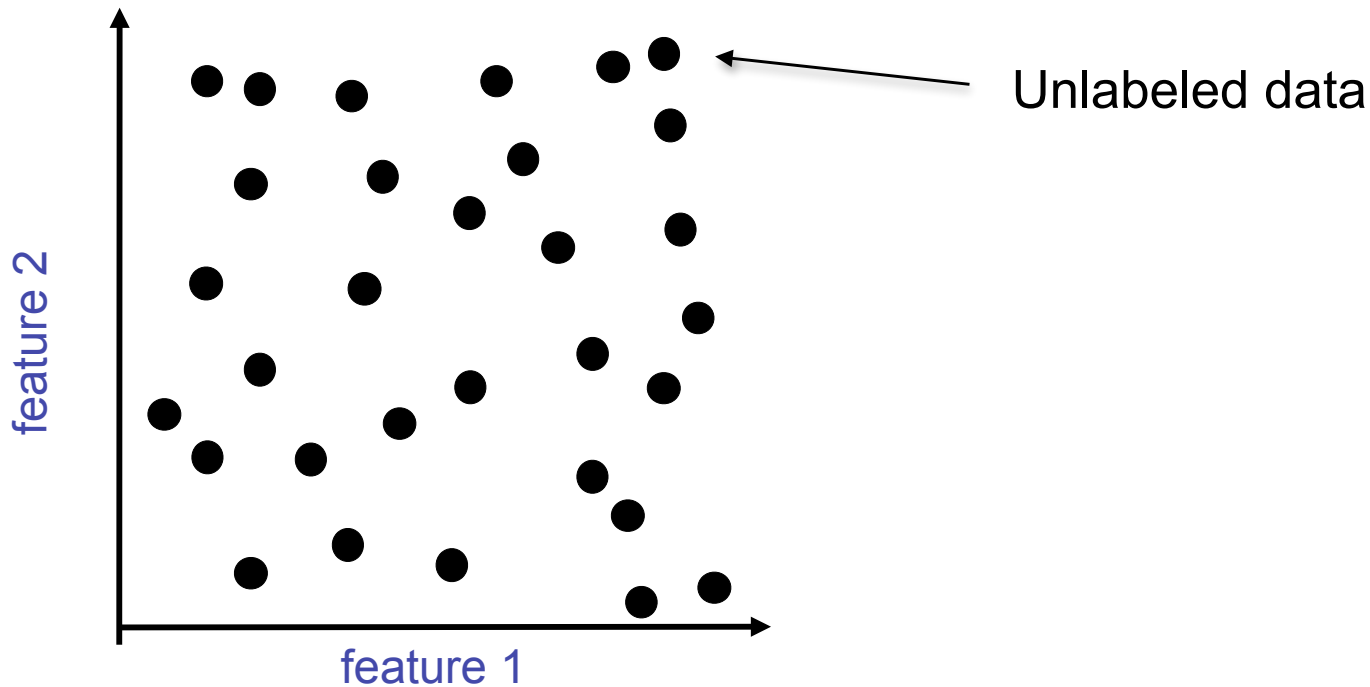
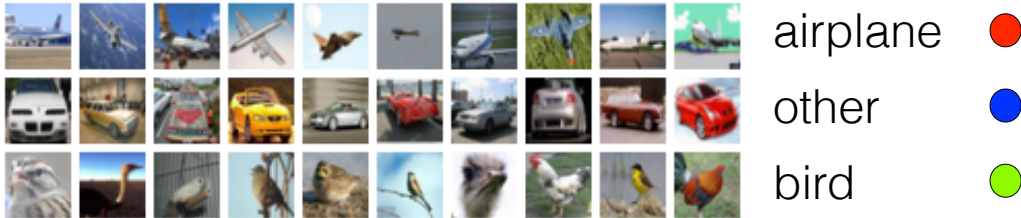
- [SzepesvariLattimore] [Bandit Algorithms course notes](#) Csaba Szepesvari and Tor Lattimore

Assignments

- Homework 0: (Self-examination, Not due but recommend you complete within the first week) [PDF](#)

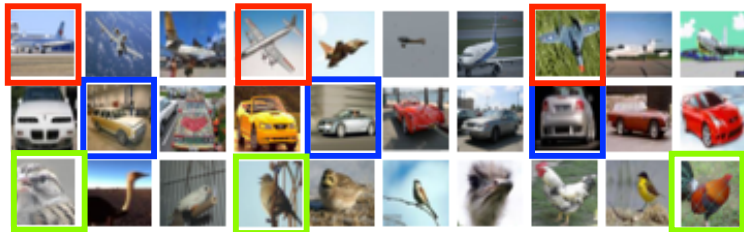
Standard Machine Learning Paradigm

- **Data:** past observations
- **Hypotheses/Models:** devised to capture the patterns in data
- **Prediction:** apply model to forecast future observations

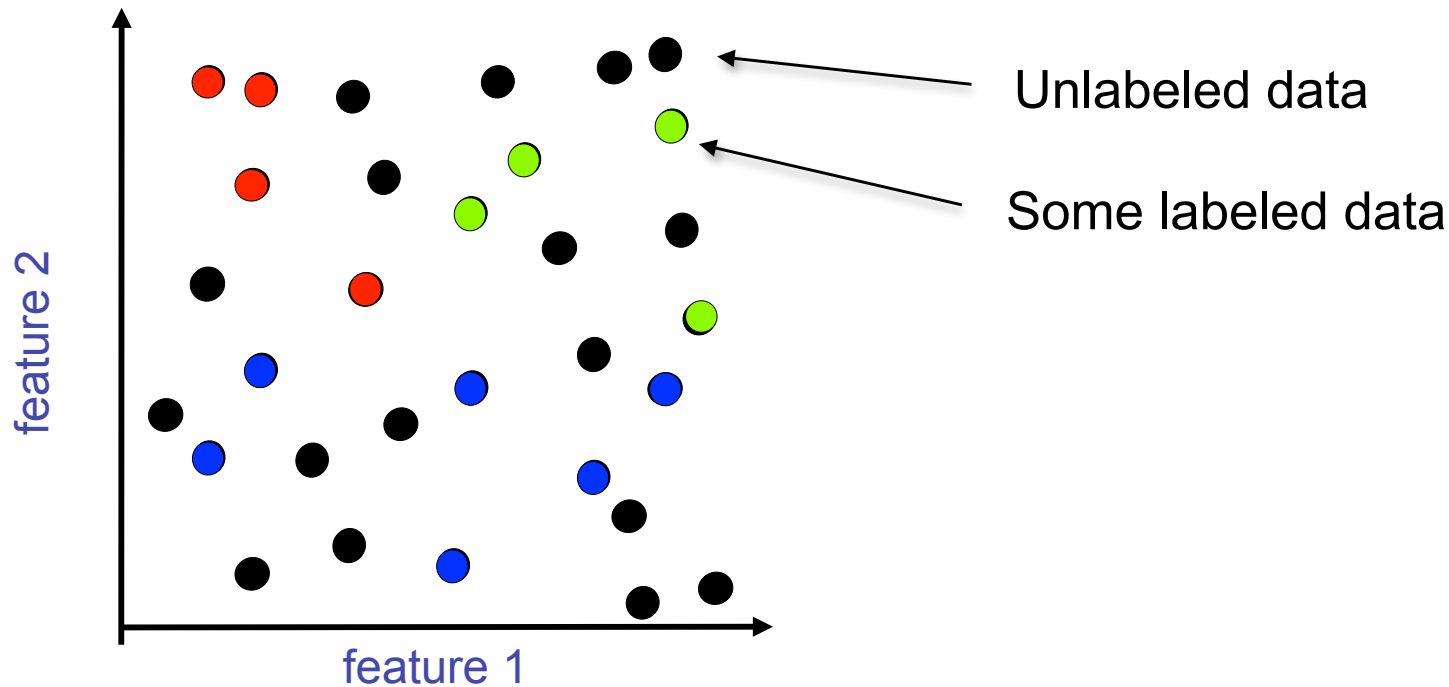


Standard Machine Learning Paradigm

- **Data:** past observations
- **Hypotheses/Models:** devised to capture the patterns in data
- **Prediction:** apply model to forecast future observations

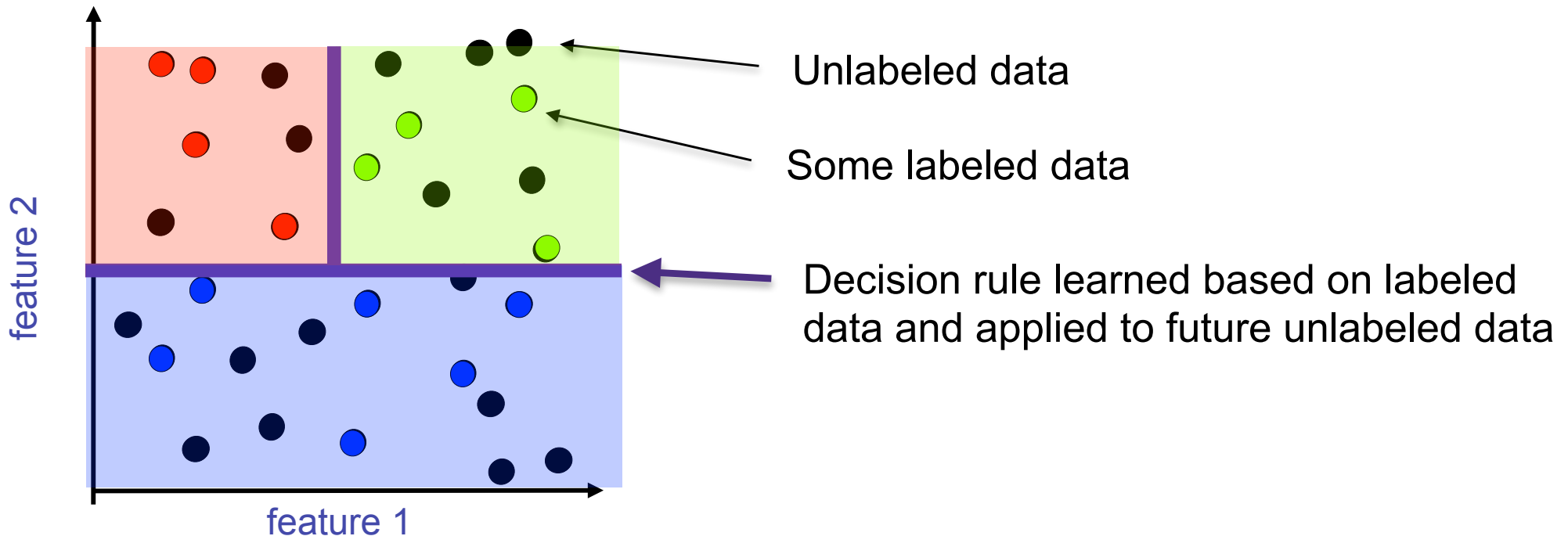
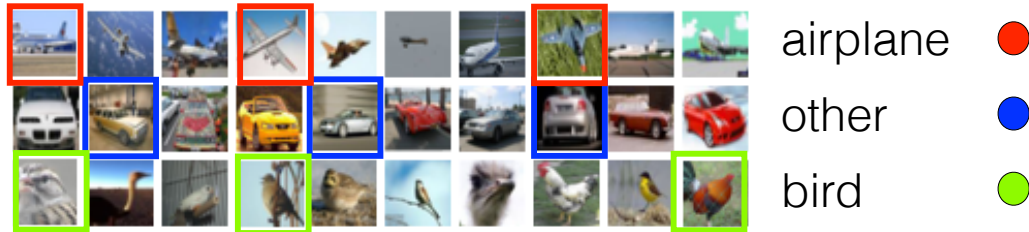


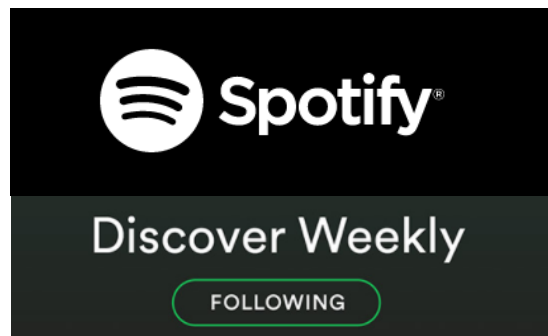
airplane ●
other ●
bird ●



Standard Machine Learning Paradigm

- **Data:** past observations
- **Hypotheses/Models:** devised to capture the patterns in data
- **Prediction:** apply model to forecast future observations





You may also like...



Do these applications actually fall into the standard machine learning paradigm?

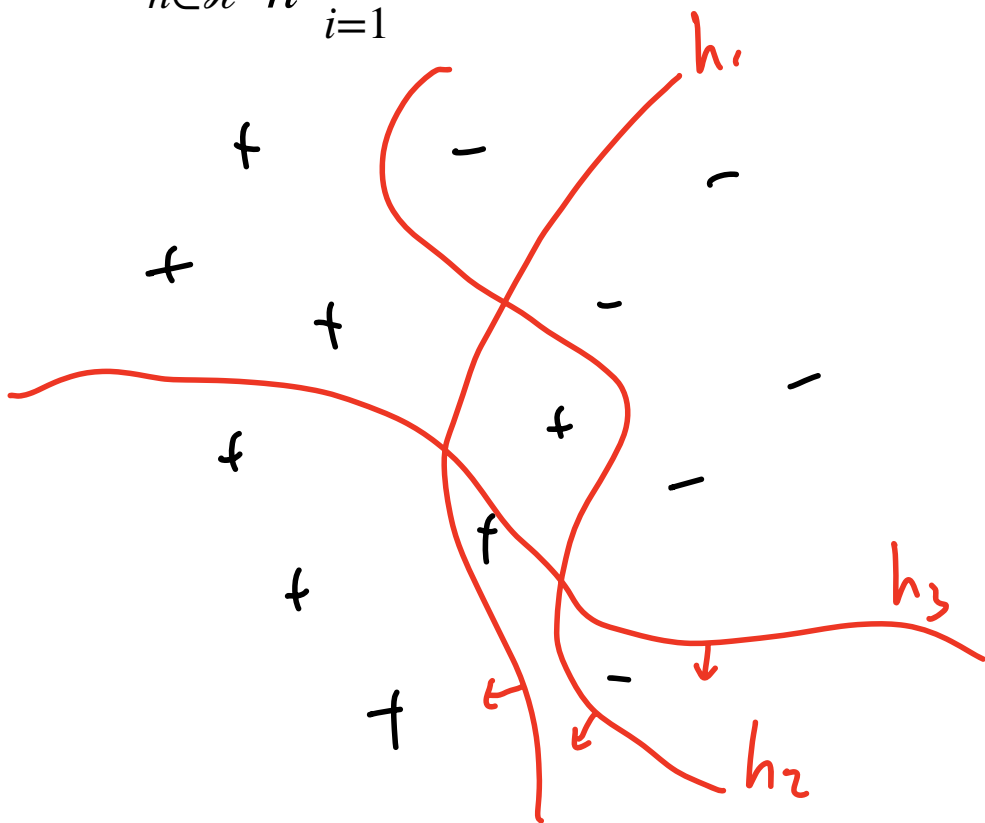
Generalization Bounds

Realizable case

Fix a finite hypothesis class $\mathcal{H} = \{h_1, h_2, \dots\}$ where $h(x) \in \{-1, 1\}$.

You are given a data set $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i = h_*(x_i)$ for some $h_* \in \mathcal{H}$

Let $\hat{h} \in \arg \min_{h \in \mathcal{H}} \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}}_{\tilde{R}_n(h) \text{ "empirical risk"}}$ how "good" is \hat{h} ?



$(x, y) \stackrel{iid}{\sim} \nu$

$$\mathbb{P}(\hat{h}(x) \neq y) = \underbrace{R(\hat{h})}_{\text{"True risk"}}$$

Realizable case

$$R(h) = \mathbb{E}[\hat{R}_n(h)]$$

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i \in \{-1, 1\}$. For any $h \in \mathcal{H}$ define $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$ and $R(h) = \mathbb{P}(h(X) \neq Y)$ where $(X, Y) \sim \nu$. Assume there exists an $h_* \in \mathcal{H}$ such that $R(h_*) = 0$. If $\hat{h} \in \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$ then with probability at least $1 - \delta$ we have

$$R(\hat{h}) \leq \frac{\log(|\mathcal{H}|/\delta)}{n}$$

$$\mathbb{P}(\hat{R}_n(\hat{h}) > \frac{\log|\mathcal{H}|/\delta}{n}) \leq \delta$$

where $(X, Y) \sim \nu$.

$$\mathbb{P}(R(\hat{h}) > \varepsilon) \leq \delta.$$

Probably Approximately Correct (PAC)

$$\mathbb{E}[R(\hat{h})]$$

Realizable case - Proof

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \\ \leq P(A) + P(B)$$

$$P(R(\hat{h}) > \varepsilon) = P(R(\hat{h}) > \varepsilon \text{ and } \bigcap_{i=1}^n \{\hat{h}(x_i) = y_i\})$$

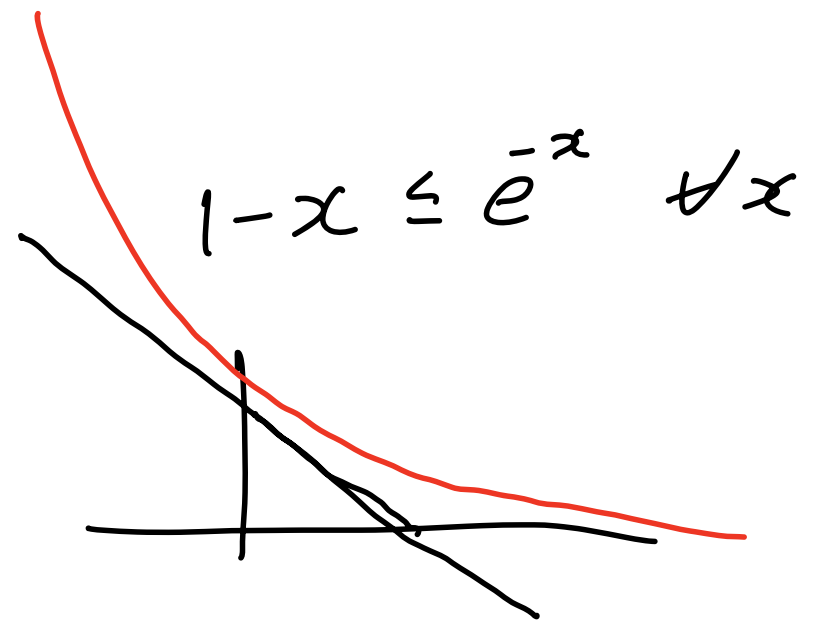
$$\leq P\left(\bigcup_{h \in \mathcal{H}} \left\{ R(h) > \varepsilon \text{ and } \bigcap_{i=1}^n \{h(x_i) = y_i\} \right\}\right)$$

$$\leq \sum_{h \in \mathcal{H}} P(R(h) > \varepsilon \text{ and } \bigcap_{i=1}^n \{h(x_i) = y_i\})$$

$$\leq \sum_{h \in \mathcal{H}} (1 - \varepsilon)^n$$

$$\leq |\mathcal{H}| (1 - \varepsilon)^n$$

$$\leq |\mathcal{H}| e^{-n\varepsilon}$$



Realizable case - Proof

Union bound: $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B) \leq \mathbb{P}(A) + \mathbb{P}(B)$

$$\begin{aligned}\mathbb{P}(R(\hat{h}) \geq \epsilon) &= \mathbb{P}(R(\hat{h}) \geq \epsilon) \\ &= \mathbb{P}(R(\hat{h}) \geq \epsilon \text{ and } \cap_{i=1}^n \{\hat{h}(x_i) = y_i\}) \\ &\leq \mathbb{P}\left(\bigcup_{h \in \mathcal{H}} \left\{R(h) \geq \epsilon \text{ and } \cap_{i=1}^n \{h(x_i) = y_i\}\right\}\right) \\ &\leq \sum_{h \in \mathcal{H}} \mathbb{P}(R(h) \geq \epsilon \text{ and } \cap_{i=1}^n \{h(x_i) = y_i\}) \\ &\leq \sum_{h \in \mathcal{H}} (1 - \epsilon)^n \\ &\leq |\mathcal{H}| \exp(-n\epsilon) \qquad \exp(-x) \geq (1 - x) \quad \forall x\end{aligned}$$

Realizable case

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i \in \{-1, 1\}$. For any $h \in \mathcal{H}$ define $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$ and $R(h) = \mathbb{P}(h(X) \neq Y)$ where $(X, Y) \sim \nu$. Assume there exists an $h_* \in \mathcal{H}$ such that $R(h_*) = 0$. If $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$ then with probability at least $1 - \delta$ we have

$$R(\hat{h}) \leq \frac{\log(|\mathcal{H}|/\delta)}{n}$$

where $(X, Y) \sim \nu$.

for $z \geq 0$

$$\mathbb{E}[z] = \int_0^{\infty} \mathbb{P}(z > \epsilon) d\epsilon$$

$$\mathbb{P}\left(R(\hat{h}) > \underbrace{\frac{\log(|\mathcal{H}|/\delta)}{n}}_{=\epsilon}\right) \leq \delta$$

$$\Rightarrow \delta = |\mathcal{H}| e^{-n\epsilon}$$

Corollary Under the conditions of the theorem (i.e., there exists an $h_* \in \mathcal{H}$ such that $R(h_*) = 0$, $(x_i, y_i) \stackrel{iid}{\sim} \nu$, and $\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$) we have $\mathbb{E}[R(\hat{h})] \leq \int_{\epsilon=0}^{\infty} \underbrace{\mathbb{P}(R(\hat{h}) \geq \epsilon)}_{\leq \delta} d\epsilon \leq \frac{2 \log(|\mathcal{H}|)}{n}$

$$\leq \min\{1, |\mathcal{H}| e^{-n\epsilon}\}$$

Agnostic (Non-realizable) case

$$h_* = \operatorname{argmin}_{h \in \mathcal{H}} R(h)$$

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i \in \{-1, 1\}$. For any $h \in \mathcal{H}$ define $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$ and $R(h) = \mathbb{P}(h(X) \neq Y)$ where $(X, Y) \sim \nu$. If $\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_n(h)$ then with probability at least $1 - \delta$ we have

$$\mathbb{E}[\hat{R}_n(h)] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{P}(h(x_i) \neq y_i) = R(h)$$

$$\max_{h \in \mathcal{H}} R(\hat{h}) - R(h) = \underbrace{R(\hat{h}) - R(h_*)}_{\text{"Excess Risk"}} \leq \sqrt{\frac{2 \log(|\mathcal{H}|/\delta)}{n}}$$

$$\mathbb{E}[R(\hat{h})] \leq \mathbb{E}[R(\hat{h})]$$

$$\mathbb{E}[\hat{R}_n(h)] = R(h)$$

≤ 0

$$R(\hat{h}) - R(h_*) = R(\hat{h}) - \hat{R}_n(\hat{h}) + \hat{R}_n(\hat{h}) - \hat{R}_n(h_*) + \hat{R}_n(h_*) - R(h_*)$$

$$\leq \left(\max_{h \in \mathcal{H}} R(h) - \hat{R}_n(h) \right) + \hat{R}_n(h_*) - R(h_*)$$

$$\leq \sqrt{\frac{\log 2|\mathcal{H}|/\delta}{2n}} + \sqrt{\frac{\log 2|\mathcal{H}|/\delta}{2n}} =$$

$$\mathbb{P}\left(\bigcup_{h \in \mathcal{H}} (R(h) - \hat{R}_n(h) > \epsilon)\right) \leq |\mathcal{H}| e^{-2n\epsilon^2} = \delta$$

Agnostic (Non-realizable) case - Proof

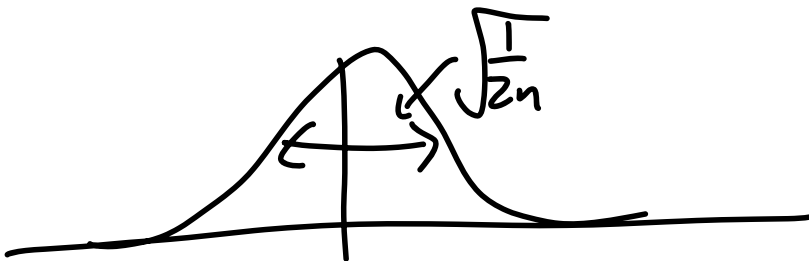
Lemma (Hoeffding's inequality): Let $Z_1, \dots, Z_n \stackrel{iid}{\sim} \nu$ where $\mathbb{E}[Z_i] = \mu$ and $Z_i \in [a, b]$ almost surely. Then

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n Z_i \geq \mu + \epsilon \right) \leq \exp \left(\frac{-2n\epsilon^2}{|b-a|^2} \right) \doteq e^{-2n\epsilon^2} = \delta$$

$$\hat{R}_n(h) - R(h) = \frac{1}{n} \sum_{i=1}^n \underbrace{(\mathbb{1}\{h(x_i) \neq y_i\} - R(h))}_{:= Z_i \in [-R(h), 1-R(h)]} \approx \mathcal{N} \left(0, \frac{1}{4n} \right)$$

$$\hat{R}_n(h) - R(h) \xrightarrow{n \rightarrow \infty} \mathcal{N} \left(0, \frac{\text{Var}(Z_1)}{n} \right)$$

$$\begin{aligned} \text{Var}(Z_0) &= R(h)(1-R(h)) \\ &\leq \frac{1}{4} \end{aligned}$$



$$\hat{R}_n(h) - R(h) \leq \sqrt{\frac{\log(1/\delta)}{2n}}$$

Agnostic (Non-realizable) case - Proof

Agnostic (Non-realizable) case

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i \in \{-1, 1\}$. For any $h \in \mathcal{H}$ define $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$ and $R(h) = \mathbb{P}(h(X) \neq Y)$ where $(X, Y) \sim \nu$. If $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$ then with probability at least $1 - \delta$ we have

$$R(\hat{h}) - R(h_*) \leq \sqrt{\frac{2 \log(|\mathcal{H}|/\delta)}{n}}.$$

Corollary Under the conditions of the theorem (i.e., $(x_i, y_i) \stackrel{iid}{\sim} \nu$, and $\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$) and $|\mathcal{H}| \geq n$, we have $\mathbb{E}[R(\hat{h})] - R(h_*) \leq \sqrt{\frac{8 \log(|\mathcal{H}|)}{n}}$

Agnostic (Non-realizable) case - Interpolation

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i \in \{-1, 1\}$. For any $h \in \mathcal{H}$ define $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$ and $R(h) = \mathbb{P}(h(X) \neq Y)$ where $(X, Y) \sim \nu$. If $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$ then with probability at least $1 - \delta$ we have

$$R(\hat{h}) - R(h_*) \leq \sqrt{\frac{2R(h_*) \log(2|\mathcal{H}|/\delta)}{n}} + \frac{\log(2|\mathcal{H}|/\delta)}{n}.$$

Proof: Use Bernstein's inequality instead of Hoeffding. ■

Infinite classes

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i \in \{-1, 1\}$. For any $h \in \mathcal{H}$ define $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$ and $R(h) = \mathbb{P}(h(X) \neq Y)$ where $(X, Y) \sim \nu$. If $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$ then with probability at least $1 - \delta$ we have

$$R(\hat{h}) - R(h_*) \leq \sqrt{\frac{2R(h_*) \log(2|\mathcal{H}|/\delta)}{n}} + \frac{\log(2|\mathcal{H}|/\delta)}{n}.$$

What if $|\mathcal{H}|$ is *infinite* such as the space of all hyperplane classifiers?

Infinite classes

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i \in \{-1, 1\}$. For any $h \in \mathcal{H}$ define $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$ and $R(h) = \mathbb{P}(h(X) \neq Y)$ where $(X, Y) \sim \nu$. If $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$ then with probability at least $1 - \delta$ we have

$$R(\hat{h}) - R(h_*) \leq \sqrt{\frac{2R(h_*) \log(2|\mathcal{H}|/\delta)}{n}} + \frac{\log(2|\mathcal{H}|/\delta)}{n}.$$

What if $|\mathcal{H}|$ is *infinite* such as the space of all hyperplane classifiers?

Lots of tools to address this:

- minimum description length
- VC-dimension and Rademacher complexity
- Covering number / log-entropy bounds

Online Learning

Realizable case

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n) \stackrel{iid}{\sim} \nu$ where $y_i \in \{-1, 1\}$. For any $h \in \mathcal{H}$ define $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$ and $R(h) = \mathbb{P}(h(X) \neq Y)$ where $(X, Y) \sim \nu$. Assume there exists an $h_* \in \mathcal{H}$ such that $R(h_*) = 0$. If $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_n(h)$ then with probability at least $1 - \delta$ we have

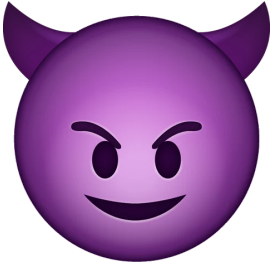
$$R(\hat{h}) \leq \frac{\log(|\mathcal{H}|/\delta)}{n}$$

where $(X, Y) \sim \nu$.

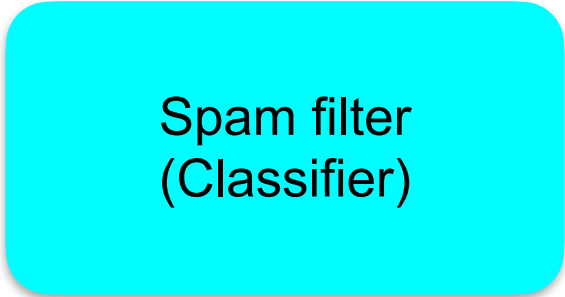
All the guarantees of the previous section (and the entirety of this class so far) has relied critically on (x, y) being drawn **IID**. Can we say anything if (x, y) are chosen **adversarially**?

Online learning

Spammer



x_t

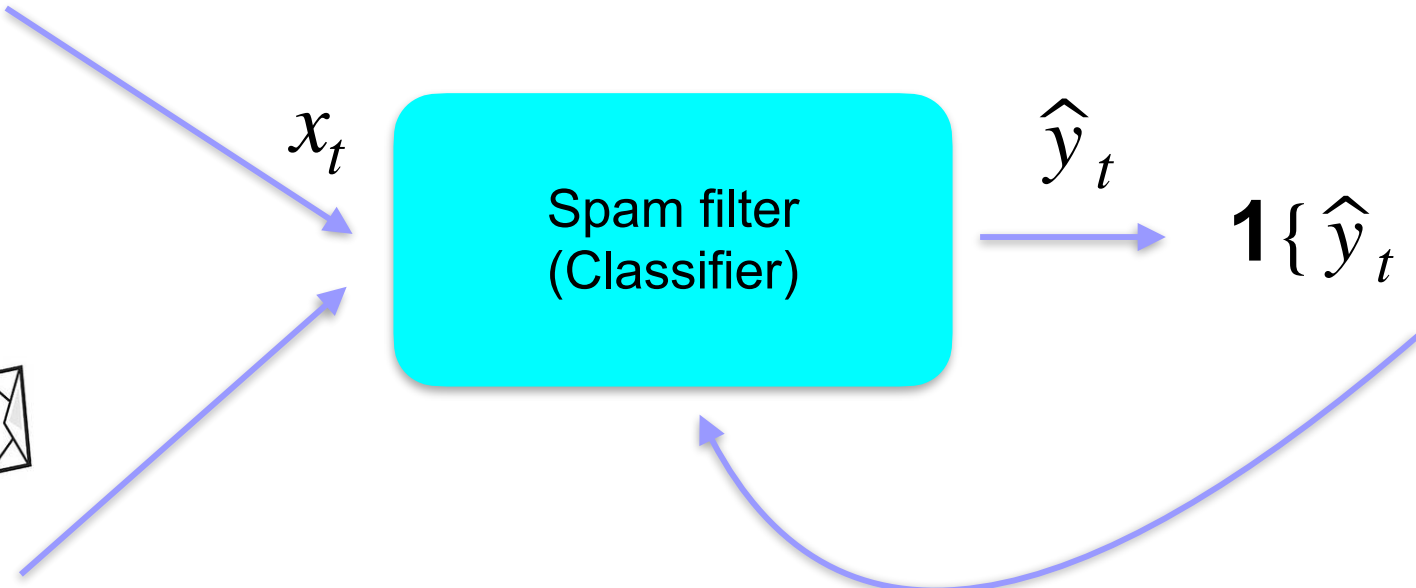


\hat{y}_t

$\mathbf{1}\{\hat{y}_t \neq y_t\}$



Real mail



Online learning

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

Settings of interest:

IID $(x_t, y_t) \sim \nu$

Adversarial (x_t, y_t) arbitrary

Online learning - Realizable IID

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

IID $(x_t, y_t) \sim \nu \quad y_t = h_*(x_t)$

We know learning theory! Choose $h_t \in \arg \min_{h \in \mathcal{H}} \sum_{s=1}^{t-1} \mathbf{1}\{h(x_s) \neq y_s\}$

Online learning - IID

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

IID $(x_t, y_t) \sim \nu \quad y_t = h_*(x_t)$

Corollary Under the conditions of the theorem (i.e., there exists an $h_* \in \mathcal{H}$ such that $R(h_*) = 0$, $(x_i, y_i) \stackrel{iid}{\sim} \nu$, and $\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$) we have $\mathbb{E}[R(\hat{h})] \leq \int_{\epsilon=0}^d \mathbb{P}(R(\hat{h}) \geq \epsilon) \leq \frac{2 \log(|\mathcal{H}|)}{n}$

Online learning - IID

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

IID $(x_t, y_t) \sim \nu \quad y_t = h_*(x_t)$

Corollary Under the conditions of the theorem (i.e., there exists an $h_* \in \mathcal{H}$ such that $R(h_*) = 0$, $(x_i, y_i) \stackrel{iid}{\sim} \nu$, and $\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$) we have $\mathbb{E}[R(\hat{h})] \leq \int_{\epsilon=0}^d \mathbb{P}(R(\hat{h}) \geq \epsilon) \leq \frac{2 \log(|\mathcal{H}|)}{n}$

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} \right] &\leq 1 + \sum_{t=2}^T \mathbb{E}[\mathbb{P}(h_t(x_t) \neq y_t)] \\ &\leq 1 + \sum_{t=2}^T \mathbb{E}[R(h_t)] \leq 1 + \sum_{t=2}^T \frac{2 \log(|\mathcal{H}|)}{t-1} \leq 2 + 2 \log(|\mathcal{H}|) \log(T) \end{aligned}$$

of mistakes grows only logarithmically!

Online learning - Adversarial

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

Adversarial (x_t, y_t) arbitrary $y_t = h_*(x_t)$

Online learning - Adversarial

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

Adversarial (x_t, y_t) arbitrary $y_t = h_*(x_t)$

We know learning theory! Choose $h_t \in \arg \min_{h \in \mathcal{H}} \sum_{s=1}^{t-1} \mathbf{1}\{h(x_s) \neq y_s\}$?

Online learning - Adversarial

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

Adversarial (x_t, y_t) arbitrary $y_t = h_*(x_t)$

We know learning theory! Choose $h_t \in \arg \min_{h \in \mathcal{H}} \sum_{s=1}^{t-1} \mathbf{1}\{h(x_s) \neq y_s\}$?

Claim There exists a sequence $\{(x_t, y_t)\}_{t=1}^T$ and $\hat{h}_t \in \arg \min_{h \in \mathcal{H}} \sum_{s=1}^{t-1} \mathbf{1}\{h(x_s) \neq y_s\}$ such that the strategy makes $\min\{|\mathcal{H}|, T\}$ mistakes.

Hint: many classifiers achieve minimum, assume adversary knows your tie-breaking strategy

Online learning - Adversarial

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

Adversarial (x_t, y_t) arbitrary $y_t = h_*(x_t)$

Halving Algorithm

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

Initialize: $V_1 = \mathcal{H}$

for $t = 1, 2, \dots$

x_t arrives

Player picks a $h_t \in V_t : \sum_{h \in V_t} \mathbf{1}\{h(x_t) = h_t(x_t)\} > \sum_{h \in V_t} \mathbf{1}\{h(x_t) = -h_t(x_t)\}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$

Online learning - Adversarial

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

Adversarial (x_t, y_t) arbitrary $y_t = h_*(x_t)$

Halving Algorithm

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

Initialize: $V_1 = \mathcal{H}$

for $t = 1, 2, \dots$

x_t arrives

Player picks a $h_t \in V_t : \sum_{h \in V_t} \mathbf{1}\{h(x_t) = h_t(x_t)\} > \sum_{h \in V_t} \mathbf{1}\{h(x_t) = -h_t(x_t)\}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$

Either the algorithm doesn't make mistake,
or *at least half* of hypotheses are discarded

Online learning - Adversarial

Input: \mathcal{H} with $|\mathcal{H}| < \infty$

for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Goal:

Minimize mistakes

$$\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\}$$

Adversarial (x_t, y_t) arbitrary $y_t = h_*(x_t)$

Theorem: Fix a finite hypothesis class \mathcal{H} so that $|\mathcal{H}| < \infty$ and for all $h \in \mathcal{H}$ we have $h(x) \in \{-1, 1\}$. Let $(x_1, y_1), \dots, (x_n, y_n)$ where x_t is arbitrary and $y_t = h_*(x_t)$ for some $h_* \in \mathcal{H}$. Then if h_t is recommended by the Halving algorithm, we have that $\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} \leq \log_2(|\mathcal{H}|)$

Online learning

Assuming that your data is IID is a **very** strong assumption that is almost never true in practice. Online learning is a different paradigm that makes no assumptions but still yields meaningful guarantees.

Assuming there exists a perfect classifier h_* :

- When x_t is drawn IID, empirical risk minimization results in only a number of mistakes that grows like $\log(T)\log(H)$
- When x_t is chosen **adversarially** empirical risk minimization can do arbitrarily badly. But there exist smarter approaches (like Halving algorithm) that make only $\log(H)$ mistakes

Questions?

Online learning in non-separable case

Online learning

Goal: Minimize regret wrt best

$$\max_{h \in \mathcal{H}} \sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\}$$

Input: \mathcal{H} with $|\mathcal{H}| < \infty$
for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Settings of interest:

IID $(x_t, y_t) \sim \nu$

Adversarial (x_t, y_t) arbitrary

Online learning

Goal: Minimize regret wrt best

$$\max_{h \in \mathcal{H}} \sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\}$$

Input: \mathcal{H} with $|\mathcal{H}| < \infty$
for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Settings of interest:

IID $(x_t, y_t) \sim \nu$

Choose $h_t \in \arg \min_{h \in \mathcal{H}} \sum_{s=1}^{t-1} \mathbf{1}\{h(x_s) \neq y_s\}$

Corollary Under the conditions of the theorem (i.e., $(x_i, y_i) \stackrel{iid}{\sim} \nu$, and $\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(x_i) \neq y_i\}$) and $|\mathcal{H}| \geq n$, we have $\mathbb{E}[R(\hat{h})] - R(h_*) \leq \sqrt{\frac{8 \log(|\mathcal{H}|)}{n}}$

$$\implies \max_{h \in \mathcal{H}} \mathbb{E} \left[\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\} \right] \leq \sqrt{8T \log(|\mathcal{H}|)}$$

Online learning

Goal: Minimize regret wrt best

$$\max_{h \in \mathcal{H}} \sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\}$$

Input: \mathcal{H} with $|\mathcal{H}| < \infty$
for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Settings of interest:

IID $(x_t, y_t) \sim \nu$

Adversarial (x_t, y_t) arbitrary

Theorem: If $z_t \in [0, 1]^d \forall t$, and I_t, p_t are chosen by exponential weights then
 $\max_{i \in [d]} \mathbb{E} \left[\sum_{t=1}^T \langle I_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle \right] = \max_{i \in [d]} \sum_{t=1}^T \langle p_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle \leq \sqrt{T \log(d)/2}$

$$\implies \max_{h \in \mathcal{H}} \mathbb{E} \left[\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\} \right] \leq \sqrt{T \log(|\mathcal{H}|)/2}$$

Online learning

Goal: Minimize regret wrt best

$$\max_{h \in \mathcal{H}} \sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\}$$

Input: \mathcal{H} with $|\mathcal{H}| < \infty$
for $t = 1, 2, \dots$

x_t arrives

Player picks $h_t \in \mathcal{H}$

y_t is revealed

Player receives loss $\ell(h_t, (x_t, y_t)) = \mathbf{1}\{h_t(x_t) \neq y_t\}$

Settings of interest:

IID $(x_t, y_t) \sim \nu$

$$\implies \max_{h \in \mathcal{H}} \mathbb{E} \left[\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\} \right] \leq \sqrt{8T \log(|\mathcal{H}|)}$$

Adversarial (x_t, y_t) arbitrary

$$\implies \max_{h \in \mathcal{H}} \mathbb{E} \left[\sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\} \right] \leq \sqrt{T \log(|\mathcal{H}|)/2}$$

Online learning

Assuming that your data is IID is a **very** strong assumption that is almost never true in practice. Online learning is a different paradigm that makes no assumptions but still yields meaningful guarantees.

Questions?

Exponential weights

Expert prediction

Suppose $b_t \in [0,1]^d$ is a vector of d experts predictions of tomorrow's temperature.

$t=1$ $t=2$ $t=3$ $t=4$ $t=5$...

Expert 1

Expert 2

Expert 3

Expert prediction

Suppose $b_t \in [0,1]^d$ is a vector of d experts predictions of tomorrow's temperature.

$t=1$ $t=2$ $t=3$ $t=4$ $t=5$...

Expert 1

Expert 2

Expert 3

$$z_t(i) = |b_t(i) - y_t|$$

i th expert's prediction

True temperature

Input: d experts

for $t = 1, 2, \dots$

Player picks $p_t \in \Delta_d$ and plays $I_t \sim p_t$

Adversary simultaneously reveals expert losses $z_t \in [0, 1]^d$

Player pays loss $\langle p_t, z_t \rangle = \mathbb{E}[z_t(I_t)]$

Expert prediction

Suppose $b_t \in [0,1]^d$ is a vector of d experts predictions of tomorrow's temperature.

$t=1$ $t=2$ $t=3$ $t=4$ $t=5$...

Expert 1

Expert 2

Expert 3

$$z_t(i) = |b_t(i) - y_t|$$

i th expert's prediction

True temperature

Input: d experts

for $t = 1, 2, \dots$

Player picks $p_t \in \Delta_d$ and plays $I_t \sim p_t$

Adversary simultaneously reveals expert losses $z_t \in [0, 1]^d$

Player pays loss $\langle p_t, z_t \rangle = \mathbb{E}[z_t(I_t)]$

Goal: Minimize
regret wrt best

$$\max_{i \in [d]} \sum_{t=1}^T \langle p_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle$$

Expert prediction

Goal: Minimize
regret wrt best

$$\max_{i \in [d]} \sum_{t=1}^T \langle p_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle$$

Input: d experts

for $t = 1, 2, \dots$

Player picks $p_t \in \Delta_d$ and plays $I_t \sim p_t$

Adversary simultaneously reveals expert losses $z_t \in [0, 1]^d$

Player pays loss $\langle p_t, z_t \rangle = \mathbb{E}[z_t(I_t)]$

Exponential weights algorithm

Input: d experts, $\eta > 0$

Initialize: $w_1 \in [1, \dots, 1]^T \in \mathbb{R}^d$

for $t = 1, 2, \dots$

Player plays $I_t \sim p_t$ where $p_t(i) = w_t(i) / \sum_{j=1}^d w_t(j)$

Adversary simultaneously reveals expert losses $z_t \in [0, 1]^d$

Player pays loss $\langle p_t, z_t \rangle = \mathbb{E}[z_t(I_t)]$

Player updates weights $w_{t+1}(i) = w_t(i) \exp(-\eta z_t(i))$

Expert prediction

Goal: Minimize regret wrt best

$$\max_{i \in [d]} \sum_{t=1}^T \langle p_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle$$

Exponential weights algorithm

Input: d experts, $\eta > 0$

Initialize: $w_1 \in [1, \dots, 1]^T \in \mathbb{R}^d$

for $t = 1, 2, \dots$

Player plays $I_t \sim p_t$ where $p_t(i) = w_t(i) / \sum_{j=1}^d w_t(j)$

Adversary simultaneously reveals expert losses $z_t \in [0, 1]^d$

Player pays loss $\langle p_t, z_t \rangle = \mathbb{E}[z_t(I_t)]$

Player updates weights $w_{t+1}(i) = w_t(i) \exp(-\eta z_t(i))$

Theorem: If $z_t \in [0, 1]^d \forall t$, and I_t, p_t are chosen by exponential weights then
$$\max_{i \in [d]} \mathbb{E} \left[\sum_{t=1}^T \langle I_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle \right] = \max_{i \in [d]} \sum_{t=1}^T \langle p_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle \leq \frac{\log(d)}{\eta} + \frac{T\eta}{8}$$

Choosing $\eta = \sqrt{\frac{8 \log(d)}{T}}$ gives regret bound of $\sqrt{T \log(d)/2}$

Expert prediction

Goal: Minimize
regret wrt best

$$\max_{i \in [d]} \sum_{t=1}^T \langle p_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle$$

Exponential weights algorithm, proof: Let $W_t = \sum_{i=1}^d w_t(i)$ so that

Expert prediction

Goal: Minimize regret wrt best

$$\max_{i \in [d]} \sum_{t=1}^T \langle p_t, z_t \rangle - \langle \mathbf{e}_i, z_t \rangle$$

Exponential weights algorithm, proof:

Let $W_t = \sum_{i=1}^d w_t(i)$ so that

$$\begin{aligned} \log \frac{W_{T+1}}{W_1} &= \sum_{t=1}^T \log \frac{W_{t+1}}{W_t} \\ &= \sum_{t=1}^T \log \left(\sum_{i=1}^d \frac{w_{t+1}(i)}{W_t} \right) \\ &= \sum_{t=1}^T \log \left(\sum_{i=1}^d \frac{w_t(i) \exp(-\eta z_t(i))}{W_t} \right) \\ &= \sum_{t=1}^T \log \left(\sum_{i=1}^d p_t(i) \exp(-\eta z_t(i)) \right) \\ &= \sum_{t=1}^T \log \left(\exp(-\eta \mathbb{E}[z_t(I_t)]) \sum_{i=1}^d p_t(i) \exp(-\eta(z_t(i) - \mathbb{E}[z_t(I_t)])) \right) \\ &= \sum_{t=1}^T -\eta \mathbb{E}[z_t(I_t)] + \log \left(\mathbb{E}[\exp(-\eta(z_t(I_t) - \mathbb{E}[z_t(I_t)]))] \right) \\ &\leq \sum_{t=1}^T -\eta \mathbb{E}[z_t(I_t)] + \eta^2/8 \end{aligned}$$

$$\begin{aligned} \log \frac{W_{T+1}}{W_1} &\geq \log \frac{w_{T+1}(i)}{W_1} \\ &= -\log(d) + \log \left(\prod_{t=1}^T \exp(-\eta z_t(i)) \right) \\ &= -\log(d) - \sum_{t=1}^T \eta z_t(i) \end{aligned}$$

$$\implies \sum_{t=1}^T \eta \mathbb{E}[z_t(I_t)] - \sum_{t=1}^T \eta z_t(i) \leq \log(d) + \eta^2 T/8$$

Online Convex Optimization

Convex surrogate loss functions

Previous section for the **adversarial** case suggested using multiplicative weights over the $|H|$ hypotheses, which is completely intractable in practice.

And in the **stochastic** case we used $h_t \in \arg \min_{h \in \mathcal{H}} \sum_{s=1}^{t-1} \mathbf{1}\{h(x_s) \neq y_s\}$ which is also intractable to compute!

So it seems we have no practical algorithm! Solution: relax the objective.

Convex surrogate loss functions

Previous section for the **adversarial** case suggested using multiplicative weights over the $|\mathcal{H}|$ hypotheses, which is completely intractable in practice.

And in the **stochastic** case we used $h_t \in \arg \min_{h \in \mathcal{H}} \sum_{s=1}^{t-1} \mathbf{1}\{h(x_s) \neq y_s\}$ which is also intractable to compute!

So it seems we have no practical algorithm! Solution: relax the objective.

Instead of $\max_{h \in \mathcal{H}} \sum_{t=1}^T \mathbf{1}\{h_t(x_t) \neq y_t\} - \mathbf{1}\{h(x_t) \neq y_t\}$

We use $\max_{h \in \mathcal{H}} \sum_{t=1}^T \ell(h_t, (x_t, y_t)) - \ell(h, (x_t, y_t))$ with \mathcal{H} convex

Example: Linear classification takes $\mathcal{H} \subset \mathbb{R}^d$ and $\ell(h, (x_t, y_t)) = \log(1 + \exp(-y_t h^\top x_t))$

Convex surrogate loss functions

Goal: $\max_{h \in \mathcal{H}} \sum_{t=1}^T \ell(h_t, (x_t, y_t)) - \ell(h, (x_t, y_t))$ with \mathcal{H} convex

Online gradient descent

Input: $\mathcal{H} \subset \mathbb{R}^d$, convex loss function ℓ , step size $\eta > 0$

Initialize: Choose any $h_1 \in \mathcal{H}$

for $t = 1, 2, \dots$

Player plays $h_t \in \mathcal{H}$

Adversary simultaneously reveals (x_t, y_t)

Player pays loss $\ell_t(h_t) := \ell(h_t, (x_t, y_t))$

Player updates $w_{t+1} = \Pi_{\mathcal{H}}(w_t - \eta \nabla_h \ell_t(h_t))$

Theorem Online gradient descent satisfies for any $h_* \in \mathcal{H}$

$$\sum_{t=1}^T \ell(h_t, (x_t, y_t)) - \ell(h_*, (x_t, y_t)) \leq \frac{\|h_*\|_2^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla_h \ell_t(h_t)\|_2^2$$

if $\max_{h \in \mathcal{H}} \|h\|_2 \leq R$ and $\ell(\cdot)$ is G -Lipschitz then $\text{regret} \leq RB\sqrt{T}$

Proof

Theorem Online gradient descent satisfies for any $h_* \in \mathcal{H}$

$$\sum_{t=1}^T \ell(h_t, (x_t, y_t)) - \ell(h_*, (x_t, y_t)) \leq \frac{\|h_*\|_2^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla_h \ell_t(h_t)\|_2^2$$

Questions?
