

Statistical and computational foundations of machine learning

- What is machine learning?
 - Study of computer algorithms that improve through experience and use of data.
-

- As scientific discipline machine learning is part of
 - 1) Statistics
 - 2) Optimization (Operations research)
 - 3) Computer science

It is part of AI.

These days it is almost synonymous with AI.

Applications

- Security applications
 - Spam filtering (emails, discussion forums)
 - Credit card fraud detection
 - Network intrusion detection
- Content optimization
 - Web search
 - Recommendation systems (Netflix, Amazon, ...)

- online advertising
 - Computer vision
 - OCR
 - Face detection in cameras
 - Face recognition
 - Object recognition
 - Natural language processing
 - Speech-to-text
 - Translation
 - Chatbots
 - Reinforcement learning
 - Robotics
 - Trading
-

Various types of ML

- Supervised vs. Unsupervised

- Are there labels that need be predicted?

- No obvious labels

(clustering, topic modelling)

- Active vs. Passive

(How are labels collected?)

- Passive: For random data examples?

- Active: For carefully and adaptively chosen examples

- Online vs. Batch

- Learning and prediction

- is interleaved
 - Learn, then predict
 - Adversarial vs. IID vs. Helpful
 - How are examples constructed
-

We will focus on supervised learning, passive learning.

Mostly in batch setting with some excursions to online. Mostly i.i.d. with some excursions to adversarial.

We will almost exclusively look at binary classification.

Papayas

- Tropical fruit.
 - Can observe color and softness
 - Predict tasty/non-tasty
-

Statistical model

- Domain X
(Domain = non-empty set)
(Feature space)

• Papaya's $X = \mathbb{R}^2$

• Emails $X = \{a, b, \dots, z\}^*$
⋮

• Unlabeled example
 $x \in X$

• Set of possible labels

$$Y = \{0, 1\} \text{ or } Y = \{+1, -1\}$$

(We are dealing with binary classification only.)

• Labeled example

$$(x, y) \in X \times Y$$

... ..

- Unlabeled sample

$$S = (x_1, x_2, \dots, x_m) \in X^*$$

- Labeled sample

$$S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)) \\ \in (X \times Y)^*$$

- Classifier

$$h: X \rightarrow Y$$

- Predictor $h \in \mathcal{Y}^X$
- Hypothesis
- Concept
- Model
- Set $h \leftrightarrow h^{-1}(1)$

• Learning algorithm

= mapping from labeled samples to classifiers

$$A: \underbrace{(X \times Y)^*}_{\text{set of all possible}} \rightarrow \underbrace{Y^X}_{\text{set of all possible}}$$

set of all possible

set of all possible

labeled samples
Learner / Training algorithm, classifiers

- Mistake/error of a classifier h on labeled example (x, y)

$$\text{iff } h(x) \neq y$$

- Probability distribution D over $X \times Y$

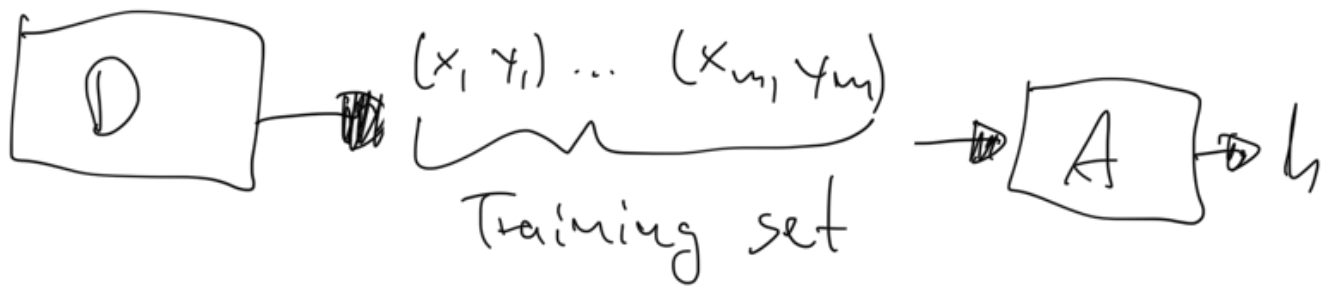
- D implicitly induces a distribution over labeled samples of size m :

$$\mathcal{T}_m$$

↓
 $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$
is an i.i.d. sequence
sampled from D .

- $(x_i, y_i) \sim D$
- $(x_1, y_1), \dots, (x_m, y_m)$ are independent
 - Beware x_i and y_i might be dependent!
- IID = independently identically distributed

System diagram



D models both the distribution of unlabeled examples and the "correct" labels.

Papaya example

- Collect bunch of papayas at random
 - Document their features
 - Taste them all
 - Produce labeled sample
-

Email spam example

Learning from

- Collect bunch of emails at random
 - Manually determine for each if it's spam or not
 - The resulting collection is a labeled sample
-

Success criterion

- Generalization error of a predictor $h: X \rightarrow Y$
$$\text{err}_D(h) = \Pr[h(x) \neq y]$$

where $(x, y) \sim D$.

$$\text{err}_D(h) = \Pr[h(x) \neq y]$$

More generally we can look at

$$L_D(h) = \mathbb{E} \left[\underbrace{l(h(x), y)}_{\text{loss of } h \text{ on } (x, y)} \right]$$

$L_D(h) / \text{err}_D(h)$ is generalization error, expected loss risk.

Note for fixed h
generalization error is a fixed number. It is not random.

The problem however is

that \mathcal{D} is often unknown
and thus $\text{err}_{\mathcal{D}}(h)$ is
unknown as well.

Optimal Bayes classifier

Suppose you know \mathcal{D} .

The best classifier you
can construct is called
"Bayes optimal classifier".

It is defined as

$$h^*(x) = \begin{cases} 1 & \text{if } \Pr[y=1 | x] \geq 1/2 \\ 0 & \text{if } \Pr[y=1 | x] < 1/2 \end{cases}$$

• When $\Pr[Y=1|x] = 1/2$ we can predict 0 or 1. It does not matter.

• If labels are $\{0, 1\}$ we can also express it as

$$h^*(x) = \begin{cases} 1 & \text{if } \mathbb{E}[Y|x] \geq 1/2 \\ 0 & \text{if } \mathbb{E}[Y|x] < 1/2 \end{cases}$$

• If labels are $\{+1, -1\}$ we can express it as

$$h^*(x) = \begin{cases} +1 & \text{if } \mathbb{E}[Y|x] \geq 0 \end{cases}$$

$$L = -1 \quad \text{if } \mathbb{E}[\gamma|x] < 0$$

$$= \text{sign}(\mathbb{E}[\gamma|x])$$

where

$$\text{sign}(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

(This version of sign function is very common in ML. In some other fields

$$\curvearrowright +1 \quad \text{if } z > 0$$

$$\text{sign}(z) = \begin{cases} 0 & \text{if } z=0 \\ -1 & \text{if } z < 0 \end{cases}$$

• The generalization error of a Bayes optimal classifier is the lowest among all possible functions:

Consider any $h: X \rightarrow Y$

Then

$$\text{err}_D(h) = \Pr[h(x) \neq y]$$

$$= \mathbb{E}[\mathbb{1}[h(x) \neq y]]$$

$$= \mathbb{E} \left[\mathbb{E} \left[\mathbb{1} [h(x) \neq y] \mid x \right] \right]$$

$$\geq \mathbb{E} \left[\mathbb{1} [h^*(x) \neq y] \mid y \right]$$

$$\Pr[h(x) \neq y \mid x] \geq \Pr[h^*(x) \neq y \mid x]$$

If $h(x) = h^*(x)$ then both sides are equal. Suppose $h(x) \neq h^*(x)$

$$\Pr[h(x) \neq y \mid x] = \begin{cases} \Pr[y=0 \mid x] & \text{if } h(x)=1 \\ \Pr[y=1 \mid x] & \text{if } h(x)=0 \end{cases}$$

$$\Pr[h^*(x) \neq y \mid x] = \begin{cases} \Pr[y=0 \mid x] & \text{if } h^*(x)=1 \\ \Pr[y=1 \mid x] & \text{if } h^*(x)=0 \end{cases}$$

~ Pr[y=1|x] 1/2

Case 1: $h(x)=0, h^*(x)=1$

$$h^*(x)=1, \Pr[h^*(x) \neq y | x] = \Pr[y=0 | x]$$

$$h(x)=0, \Pr[h(x) \neq y | x] = \Pr[y=1 | x]$$

$$h^*(x)=1 \Rightarrow \Pr[y=1 | x] \geq 1/2$$

$$\Pr[y=0 | x] \leq 1/2$$

Case 2: $h(x)=1, h^*(x)=0$

Symmetric.

Note that

$$\dots (h^*) \leq 1/2$$

$$\text{err}_D(h^*) - \text{err}_D(h)$$

- Note that $\text{err}_D(h^*)$ can be positive, and even it can be pretty high.

This is the fundamental limit on generalization of any classifier.

ERM

- ERM = empirical risk minimization
- Also called SEM = sample error minimization
- Pick a class predictors
 $H \subseteq \mathcal{Y}^X$

" = /

- Given labeled sample

$$S = \{(x_1, y_1), (x_2, y_2) \dots (x_m, y_m)\}$$

define empirical error
or sample error or
training error

$$\widehat{\text{err}}_S(h) = \frac{|\{i : 1 \leq i \leq m, h(x_i) \neq y_i\}|}{m}$$

fraction of mistakes
h makes on S.

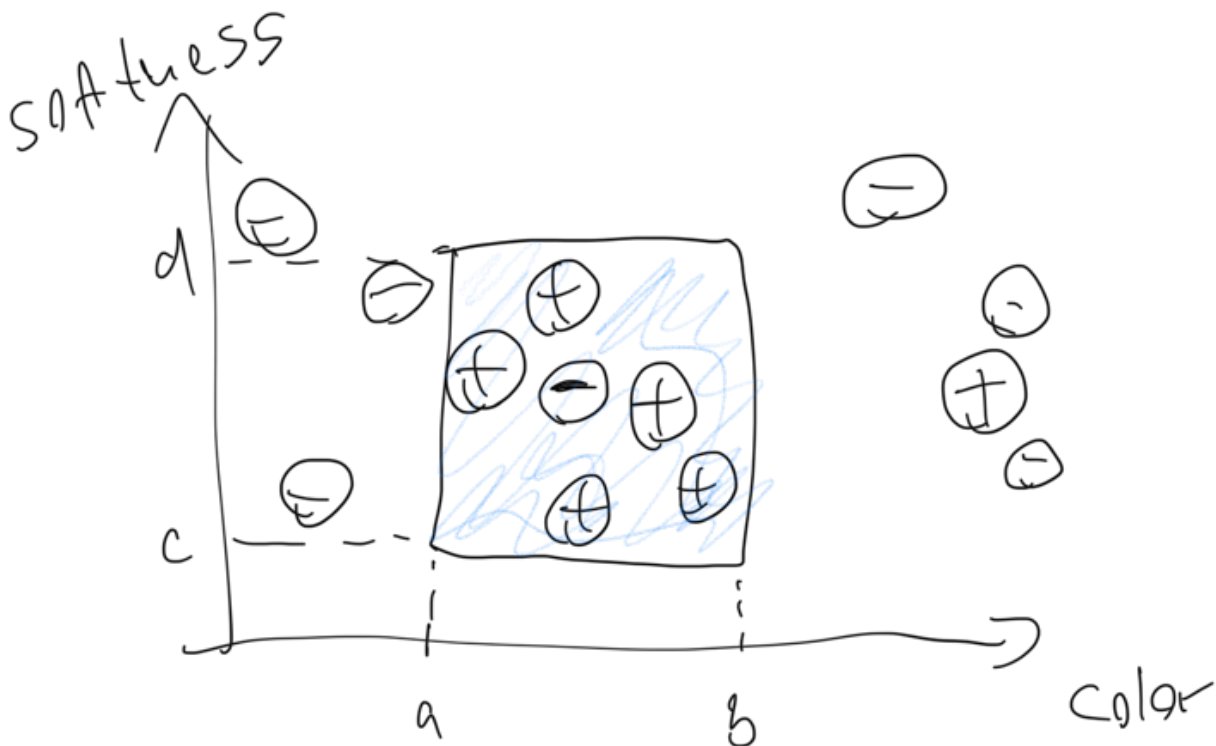
- ERM algorithm:

Find $h \in H$ with

the smallest possible empirical error.

$$\text{ERM}(S) = \underset{h \in H}{\text{argmin}} \widehat{\text{err}}_S(h)$$

Papaya example + Rectangles



$$h_{a,b,c,d}(x) = \begin{cases} 1 & \text{if } a \leq x_1 \leq b, c \leq x_2 \leq d \\ 0 & \text{otherwise} \end{cases}$$

ERM can go horribly

Imagine that U is
uniform distribution over $[0, 1]$.

Define \mathcal{D} over $[0, 1] \times \{0, 1\}$
as follows $(x, y) \sim \mathcal{D}$

$$x \sim U$$

$$y = 1$$

$\text{ERM}(S) = h$ where

$$h(x) = \begin{cases} 1 & \text{if } (x, 1) \in S \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{\text{err}}_S(h) = 0$$

+ + + + + Sample

$$\text{err}_D(h) = 1 \quad \begin{array}{c} \uparrow \uparrow \uparrow \\ 000 \end{array}$$

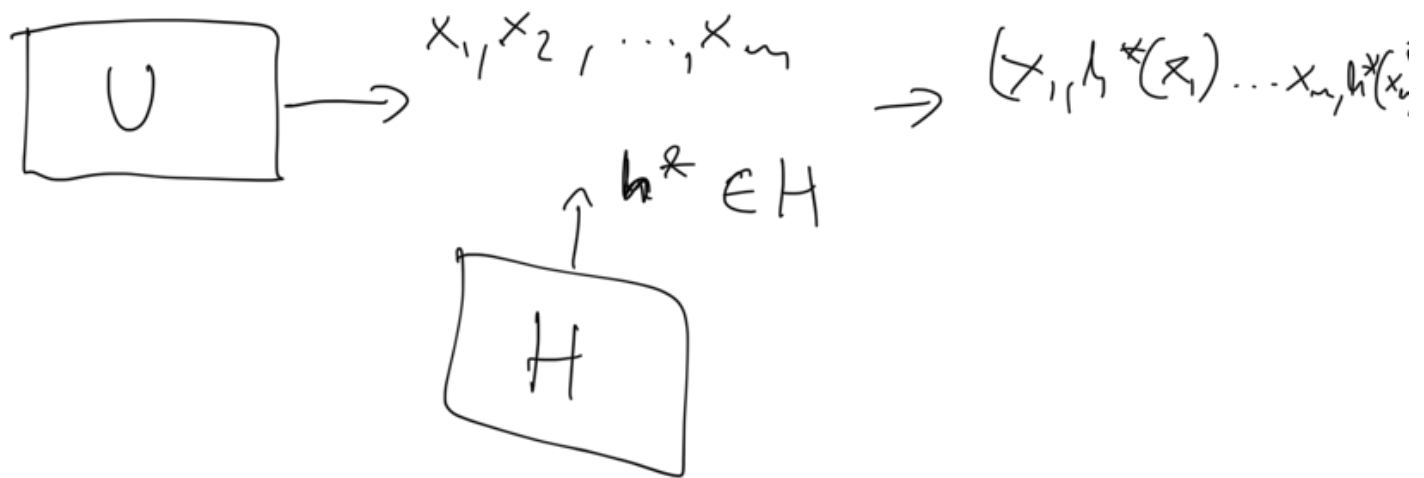
• When $\widehat{\text{err}}_S(h)$ is small
and $\text{err}_D(h)$ is big
we call it overfitting.

• Overfitting is the opposite
of generalization

Sometimes ERM goes right

• Assume H is finite

- Assume Bayes optimal classifier belongs to H and $\text{err}_D(h) = 0$.



Theorem: Suppose $H \subseteq Y^X$ is finite.

Let \mathcal{D} be a distribution ^{over $X \times Y$} such that for some $h^* \in H$, $\text{err}_D(h^*) = 0$.

Let S be a labeled sample from \mathcal{D} of size $m \geq \frac{\ln H + \ln 1/\delta}{\epsilon}$

then with probability at least $1 - \delta$, $\text{err}_D(\text{ERM}(S)) < \epsilon$.

Proof :

$$\text{err}_{D_{H^*}}(h) = \Pr[h^*(x) \neq h(x)]$$

$$\text{Let } H_{\text{BAD}} = \{h \in H : \text{err}_{D_{H^*}}(h) \geq \epsilon\}$$

$$1) \text{err}_D(\text{ERM}(S)) \geq \epsilon$$

\Leftrightarrow

$$2) \text{ERM}(S) \in H_{\text{BAD}}$$

\Rightarrow

$$3) \exists h \in H_{\text{BAD}} \text{ s.t. } \widehat{\text{err}}_S(h) = 0$$

$$1) \forall h \in H_{\text{BAD}} : \text{err}_S(h) > 0$$

Complementary events

$$2) \text{ERM}(S) \in H_{\text{BAD}}$$

$$3) \text{err}_D(\text{ERM}(S)) < \epsilon$$

• Consider $h \in H_{\text{BAD}}$

• Consider $\widehat{\text{err}}_S(h)$.

• What is the probability that

$$\widehat{\text{err}}_S(h) = 0 \quad ?$$

$$\Pr[\widehat{\text{err}}_S(h) = 0] \leq (1 - \epsilon)^m$$

It suffice to show that

$$\Pr[\exists h \in H_{\text{BAD}} \text{ s.t. } \text{err}_D(h) \geq \epsilon] \leq \delta$$

$$\Pr \left[\exists h \in H_{\text{BAD}} \text{ s.t. } \widehat{\text{err}}_S(h) \geq \epsilon \right]$$

$$= \Pr \left[\bigcup_{h \in H_{\text{BAD}}} \{S : \widehat{\text{err}}_S(h) \geq \epsilon\} \right]$$

$$\leq \sum_{h \in H_{\text{BAD}}} \Pr \left[\{S : \widehat{\text{err}}_S(h) \geq \epsilon\} \right] \quad (\text{Union Bound})$$

$$= \sum_{h \in H_{\text{BAD}}} \Pr \left[\widehat{\text{err}}_S(h) \geq \epsilon \right]$$

$$\leq |H_{\text{BAD}}| \cdot (1 - \epsilon)^m$$

$$\leq |H| \cdot (1 - \epsilon)^m$$

$$< |H| \cdot e^{-\epsilon m} \quad (1 + x < e^{\pm x})$$

$\vdash \text{||||}$

$\leq \delta$

|||||

