

Machine Learning

CSE427

Mahbub Majumdar
Typeset by:
Syeda Ramisa Fariha

BRAC University
66 Mohakhali
Dhaka, Bangladesh

May 28, 2018



Superhumungous Thanks

These slides were typeset by Syeda Ramisa Fariha.

Without her tremendous dedication, these slides would not exist.

Table of Contents

The Probability Approximately Correct Learning Model

Agnostic PAC Learnability
Scope of Learning Problems Modeled

The Probability Approximately Correct Learning Model

Topics Covered In Previous Lectures

- Finite *hypothesis classes*

The Probability Approximately Correct Learning Model

Topics Covered In Previous Lectures

- Finite *hypothesis classes*
- Use *ERM* procedure on large enough S

The Probability Approximately Correct Learning Model

Topics Covered In Previous Lectures

- Finite *hypothesis classes*
- Use *ERM* procedure on large enough S
- Then output will be *Probability Approximately Correct*

The Probability Approximately Correct Learning Model

PAC Learnability

Definition: PAC Learnability: \mathcal{H} is PAC Learnable if

- \exists a function $m_{\mathcal{H}}(\epsilon, \delta) : (0, 1) \times (0, 1) \implies \mathbb{N}$, and a learning algorithm such that

The Probability Approximately Correct Learning Model

PAC Learnability

Definition: PAC Learnability: \mathcal{H} is PAC Learnable if

- \exists a function $m_{\mathcal{H}}(\epsilon, \delta) : (0, 1) \times (0, 1) \implies \mathbb{N}$, and a learning algorithm such that
- for all $\epsilon, \delta \in (0, 1)$ and $\forall D$ over X and $\forall f : X \rightarrow \{0, 1\}$, if the *Realizability Assumption* holds, wrt \mathcal{H}, D, f ,

The Probability Approximately Correct Learning Model

PAC Learnability

Definition: PAC Learnability: \mathcal{H} is PAC Learnable if

- \exists a function $m_{\mathcal{H}}(\epsilon, \delta) : (0, 1) \times (0, 1) \implies \mathbb{N}$, and a learning algorithm such that
- for all $\epsilon, \delta \in (0, 1)$ and $\forall D$ over X and $\forall f : X \rightarrow \{0, 1\}$, if the *Realizability Assumption* holds, wrt \mathcal{H}, D, f ,
- then, when running the Learning Algorithm on $m > m_{\mathcal{H}}(\epsilon, \delta)$ iid examples generated by D , and labeled by f ,

The Probability Approximately Correct Learning Model

PAC Learnability

Definition: PAC Learnability: \mathcal{H} is PAC Learnable if

- \exists a function $m_{\mathcal{H}}(\epsilon, \delta) : (0, 1) \times (0, 1) \implies \mathbb{N}$, and a learning algorithm such that
 - for all $\epsilon, \delta \in (0, 1)$ and $\forall D$ over X and $\forall f : X \rightarrow \{0, 1\}$, if the *Realizability Assumption* holds, wrt \mathcal{H}, D, f ,
 - then, when running the Learning Algorithm on $m > m_{\mathcal{H}}(\epsilon, \delta)$ iid examples generated by D , and labeled by f ,
- \Rightarrow the algorithm returns a hypothesis h , such that, with probability of at least $1 - \delta$ (over the choice of examples):

The Probability Approximately Correct Learning Model

PAC Learnability

Definition: PAC Learnability: \mathcal{H} is PAC Learnable if

- \exists a function $m_{\mathcal{H}}(\epsilon, \delta) : (0, 1) \times (0, 1) \implies \mathbb{N}$, and a learning algorithm such that
 - for all $\epsilon, \delta \in (0, 1)$ and $\forall D$ over X and $\forall f : X \rightarrow \{0, 1\}$, if the *Realizability Assumption* holds, wrt \mathcal{H}, D, f ,
 - then, when running the Learning Algorithm on $m > m_{\mathcal{H}}(\epsilon, \delta)$ iid examples generated by D , and labeled by f ,
- \Rightarrow the algorithm returns a hypothesis h , such that, with probability of at least $1 - \delta$ (over the choice of examples):
- \Rightarrow then

$$L_{D,f}(h) \leq \epsilon.$$

The Probability Approximately Correct Learning Model

PAC Learnability

Note:

ϵ : difference between output classifier and optimal classifier

δ : how likely h is inaccurate

- Might accidentally sample the same point over and over,
 $S = \{\text{single point}\}$

The Probability Approximately Correct Learning Model

PAC Learnability

Note:

ϵ : difference between output classifier and optimal classifier

δ : how likely h is inaccurate

- Might accidentally sample the same point over and over,
 $S = \{\text{single point}\}$
- *Nonzero* ϵ enables forgiveness. The learner is allowed to make minor errors

The Probability Approximately Correct Learning Model

Sample Complexity

- if \mathcal{H} is *PAC learnable*
 - \Rightarrow many $m_{\mathcal{H}}$ satisfy requirements of *PAC learnability*
 - \Rightarrow define *sample complexity* as the minimal $m_{\mathcal{H}}(\epsilon, \delta)$

The Probability Approximately Correct Learning Model

Sample Complexity

- if \mathcal{H} is *PAC learnable*
 - \Rightarrow many $m_{\mathcal{H}}$ satisfy requirements of *PAC learnability*
 - \Rightarrow define *sample complexity* as the minimal $m_{\mathcal{H}}(\epsilon, \delta)$
- Recall from last lecture: Every finite hypothesis class is PAC learnable with sample complexity $m_{\mathcal{H}} \leq \frac{1}{\epsilon} \log\left(\frac{|\mathcal{H}|}{\delta}\right)$

The Probability Approximately Correct Learning Model

Sample Complexity

- if \mathcal{H} is *PAC learnable*
 - \Rightarrow many $m_{\mathcal{H}}$ satisfy requirements of *PAC learnability*
 - \Rightarrow define *sample complexity* as the minimal $m_{\mathcal{H}}(\epsilon, \delta)$
- Recall from last lecture: Every finite hypothesis class is PAC learnable with sample complexity $m_{\mathcal{H}} \leq \frac{1}{\epsilon} \log\left(\frac{|\mathcal{H}|}{\delta}\right)$
- Because we have chosen the the sample complexity to be the *minimal* $m_{\mathcal{H}}$, it will not be bigger than $\frac{1}{\epsilon} \log\left(\frac{|\mathcal{H}|}{\delta}\right)$. That's why we have $m_{\mathcal{H}} \leq \frac{1}{\epsilon} \log\left(\frac{|\mathcal{H}|}{\delta}\right)$

The Probability Approximately Correct Learning Model

Sample Complexity

- if \mathcal{H} is *PAC learnable*
 \Rightarrow many $m_{\mathcal{H}}$ satisfy requirements of *PAC learnability*
 \Rightarrow define *sample complexity* as the minimal $m_{\mathcal{H}}(\epsilon, \delta)$
- Recall from last lecture: Every finite hypothesis class is PAC learnable with sample complexity $m_{\mathcal{H}} \leq \frac{1}{\epsilon} \log\left(\frac{|\mathcal{H}|}{\delta}\right)$
- Because we have chosen the the sample complexity to be the *minimal* $m_{\mathcal{H}}$, it will not be bigger than $\frac{1}{\epsilon} \log\left(\frac{|\mathcal{H}|}{\delta}\right)$. That's why we have $m_{\mathcal{H}} \leq \frac{1}{\epsilon} \log\left(\frac{|\mathcal{H}|}{\delta}\right)$
- There are infinite classes that are learnable also, *VC Dimension* determines learnability.

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- If we remove *Realizability*, this becomes an *Agnostic PAC Model*

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- If we remove *Realizability*, this becomes an *Agnostic PAC Model*
- Realizability requires $\exists h^*$ such that,

$$\mathbb{P}_{x \sim D} \left(h^*(x) = f(x) \right) = 1$$

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- If we remove *Realizability*, this becomes an *Agnostic PAC Model*
- Realizability requires $\exists h^*$ such that,

$$\mathbb{P}_{x \sim D} \left(h^*(x) = f(x) \right) = 1$$

- In real life the output labels will not be fully determined by the *features* we measure

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- If we remove *Realizability*, this becomes an *Agnostic PAC Model*
- Realizability requires $\exists h^*$ such that,

$$\mathbb{P}_{x \sim D} \left(h^*(x) = f(x) \right) = 1$$

- In real life the output labels will not be fully determined by the *features* we measure
- For example, player quality not determined by only 2 features such as, $x_i = \begin{pmatrix} \text{speed}_i \\ \text{stamina}_i \end{pmatrix}$

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- In the absence of realizability, often label the distribution $D(x)$ as $D(x, y)$. (Notation used by UML.)

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- In the absence of realizability, often label the distribution $D(x)$ as $D(x, y)$. (Notation used by UML.)
- $D(x, y)$ = Joint Probability Distribution over X and Y .

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- In the absence of realizability, often label the distribution $D(x)$ as $D(x, y)$. (Notation used by UML.)
- $D(x, y)$ = Joint Probability Distribution over X and Y .
- $D_x \equiv D(x) = \sum_{y_i} D(x, y)$ is the *marginal distribution* over unlabeled points.

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- In the absence of realizability, often label the distribution $D(x)$ as $D(x, y)$. (Notation used by UML.)
- $D(x, y)$ = Joint Probability Distribution over X and Y .
- $D_x \equiv D(x) = \sum_{y_i} D(x, y)$ is the *marginal distribution* over unlabeled points.
- For example, $D_x \equiv$ probability that:

$$a \leq \text{feature}_1 \leq b$$

$$c \leq \text{feature}_2 \leq d$$

$$\vdots$$

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- $D(y | x)$ = Conditional Probability of getting label y given x

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- $D(y \mid x)$ = Conditional Probability of getting label y given x
- Notation used in "*Understanding Machine Learning*"

$$D(y \mid x) \equiv D\left((x, y) \mid x\right)$$

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- $D(y \mid x)$ = Conditional Probability of getting label y given x
- Notation used in "*Understanding Machine Learning*"

$$D(y \mid x) \equiv D\left((x, y) \mid x\right)$$

- By using a probability distribution for y ,

The Probability Approximately Correct Learning Model

Removing Some Assumptions

- $D(y \mid x)$ = Conditional Probability of getting label y given x
- Notation used in "*Understanding Machine Learning*"

$$D(y \mid x) \equiv D\left((x, y) \mid x\right)$$

- By using a probability distribution for y ,
- Allows two different mangoes with identical x 's to have different y 's because the feature set doesn't fully parameterize the set of mangoes.

The Probability Approximately Correct Learning Model

Empirical And True Error Revised

- We want to measure how likely h will mislabel points.

The Probability Approximately Correct Learning Model

Empirical And True Error Revised

- We want to measure how likely h will mislabel points.
- True error of Prediction Rule h

$$L_D(h) \equiv \mathbb{P}_{x \sim D}(h(x) \neq y) \equiv D(x, y \mid h(x) \neq y)$$

The Probability Approximately Correct Learning Model

Empirical And True Error Revised

- We want to measure how likely h will mislabel points.
- True error of Prediction Rule h

$$L_D(h) \equiv \mathbb{P}_{x \sim D} (h(x) \neq y) \equiv D(x, y \mid h(x) \neq y)$$

- Note, because we are not assuming realizability, the true error is labeled $L_D(h)$ not $L_{D,f}(h)$.

The Probability Approximately Correct Learning Model

Empirical And True Error Revised

- We want to measure how likely h will mislabel points.
- True error of Prediction Rule h

$$L_D(h) \equiv \mathbb{P}_{x \sim D} (h(x) \neq y) \equiv D(x, y \mid h(x) \neq y)$$

- Note, because we are not assuming realizability, the true error is labeled $L_D(h)$ not $L_{D,f}(h)$.
- Empirical Risk

$$L_S(h) \equiv \frac{|\{i \in [m] \mid h(x_i) \neq y_i\}|}{m}$$

The Probability Approximately Correct Learning Model

Empirical And True Error Revised

- We want to measure how likely h will mislabel points.
- True error of Prediction Rule h

$$L_D(h) \equiv \mathbb{P}_{x \sim D} (h(x) \neq y) \equiv D(x, y \mid h(x) \neq y)$$

- Note, because we are not assuming realizability, the true error is labeled $L_D(h)$ not $L_{D,f}(h)$.
- Empirical Risk

$$L_S(h) \equiv \frac{|\{i \in [m] \mid h(x_i) \neq y_i\}|}{m}$$

- Goal: Want an h minimizing $L_D(h)$ that is PAC

The Probability Approximately Correct Learning Model

Bayes' Optimal Predictor

- For a given D over $X \times \{0, 1\}$, the best labeling function $f : X \rightarrow \{0, 1\}$ is

$$f_D(x) = \begin{cases} 1 & \text{if } \mathbb{P}(y = 1 \mid x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

The Probability Approximately Correct Learning Model

Bayes' Optimal Predictor

- For a given D over $X \times \{0, 1\}$, the best labeling function $f : X \rightarrow \{0, 1\}$ is

$$f_D(x) = \begin{cases} 1 & \text{if } \mathbb{P}(y = 1 \mid x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

- Note, $f_D(x)$ is not the *true* labeling function f since we don't know what the true f is

The Probability Approximately Correct Learning Model

Bayes' Optimal Predictor

- For a given D over $X \times \{0, 1\}$, the best labeling function $f : X \rightarrow \{0, 1\}$ is

$$f_D(x) = \begin{cases} 1 & \text{if } \mathbb{P}(y = 1 \mid x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

- Note, $f_D(x)$ is not the *true* labeling function f since we don't know what the true f is
- We can show that $\forall D$, the *Bayes' Optimal Predictor* f_D , is optimal.

The Probability Approximately Correct Learning Model

Bayes' Optimal Predictor

- For a given D over $X \times \{0, 1\}$, the best labeling function $f : X \rightarrow \{0, 1\}$ is

$$f_D(x) = \begin{cases} 1 & \text{if } \mathbb{P}(y = 1 \mid x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

- Note, $f_D(x)$ is not the *true* labeling function f since we don't know what the true f is
- We can show that $\forall D$, the *Bayes' Optimal Predictor* f_D , is optimal.
- No other classifier $g : X \rightarrow \{0, 1\}$ has lower error $L_D(f_D) \leq L_D(g)$

The Probability Approximately Correct Learning Model

Bayes' Optimal Predictor

- Don't know D , thus can't utilize f_D

The Probability Approximately Correct Learning Model

Bayes' Optimal Predictor

- Don't know D , thus can't utilize f_D
- No algorithm can be guaranteed to find a predictor as good as f_D

The Probability Approximately Correct Learning Model

Bayes' Optimal Predictor

- Don't know D , thus can't utilize f_D
- No algorithm can be guaranteed to find a predictor as good as f_D
- Seek a *predictor* as close to f_D as possible

The Probability Approximately Correct Learning Model

Bayes' Optimal Predictor

- Don't know D , thus can't utilize f_D
- No algorithm can be guaranteed to find a predictor as good as f_D
- Seek a *predictor* as close to f_D as possible
- Ability to do this depends on the hypothesis class of h

What is Agnostic PAC Learnability?

Definition: Agnostic PAC Learnability

A hypothesis class h is *agnostic PAC learnable* if

- $\exists m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and \exists a learning algorithm such that

What is Agnostic PAC Learnability?

Definition: Agnostic PAC Learnability

A hypothesis class h is *agnostic PAC learnable* if

- $\exists m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and \exists a learning algorithm such that
- $\forall \epsilon, \delta \in (0, 1)$. and $\forall D$ over $X \times Y$, when running the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ iid samples generated by D

What is Agnostic PAC Learnability?

Definition: Agnostic PAC Learnability

A hypothesis class h is *agnostic PAC learnable* if

- $\exists m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and \exists a learning algorithm such that
- $\forall \epsilon, \delta \in (0, 1)$. and $\forall D$ over $X \times Y$, when running the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ iid samples generated by D
- the algorithm return a hypothesis h ,

What is Agnostic PAC Learnability?

Definition: Agnostic PAC Learnability

A hypothesis class h is *agnostic PAC learnable* if

- $\exists m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and \exists a learning algorithm such that
- $\forall \epsilon, \delta \in (0, 1)$. and $\forall D$ over $X \times Y$, when running the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ iid samples generated by D
- the algorithm return a hypothesis h ,
- such that with probability of at least $1 - \delta$ over the training samples

$$L_D(h) \leq \min_{h' \in \mathcal{H}} L_D(h') + \epsilon$$

Agnostic PAC Learnability

- If realizability holds, agnostic PAC learnability provides the same guarantee as PAC learning

Agnostic PAC Learnability

- If realizability holds, agnostic PAC learnability provides the same guarantee as PAC learning
- ⇒ Agnostic PAC Learning generalizes PAC learning

Agnostic PAC Learnability

- If realizability holds, agnostic PAC learnability provides the same guarantee as PAC learning
- ⇒ Agnostic PAC Learning generalizes PAC learning
- When *realizability* doesn't hold,

Agnostic PAC Learnability

- If realizability holds, agnostic PAC learnability provides the same guarantee as PAC learning
- ⇒ Agnostic PAC Learning generalizes PAC learning
- When *realizability* doesn't hold,
- ⇒ Can't guarantee arbitrarily small error

Agnostic PAC Learnability

- If realizability holds, agnostic PAC learnability provides the same guarantee as PAC learning
- ⇒ Agnostic PAC Learning generalizes PAC learning
- When *realizability* doesn't hold,
- ⇒ Can't guarantee arbitrarily small error
- But using *agnostic PAC learning*

Agnostic PAC Learnability

- If realizability holds, agnostic PAC learnability provides the same guarantee as PAC learning
- ⇒ Agnostic PAC Learning generalizes PAC learning
- When *realizability* doesn't hold,
- ⇒ Can't guarantee arbitrarily small error
- But using *agnostic PAC learning*
- ⇒ Can get smallest error predictor in the class \mathcal{H}

Scope of Learning Problems Modeled

- *Multiclass Classification:*

$h : \text{Documents} \rightarrow \begin{matrix} \nearrow \text{sports} \\ \text{news} \\ \searrow \text{entertainment} \end{matrix}$

error: Probability of misclassifying a document.

Scope of Learning Problems Modeled

- *Multiclass Classification:*

$$h : \text{Documents} \rightarrow \begin{matrix} \nearrow \text{sports} \\ \text{news} \\ \searrow \text{entertainment} \end{matrix}$$

error: Probability of misclassifying a document.

- **Regression:** Looking for simple patterns

$$h : \text{ultrasound measurements} \rightarrow \text{baby's weight}$$

error = Expected difference between true labels and predicted labels

$$L_D(h) \equiv \mathbb{E}_{(x,y) \sim D} [(h(x) - y)^2]$$

Scope of Learning Problems Modeled

Generalized Measure of Success

- Given \mathcal{H} , domain $Z = X \times Y$,

Scope of Learning Problems Modeled

Generalized Measure of Success

- Given \mathcal{H} , domain $Z = X \times Y$,
- Let $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$, where

ℓ = loss function

\mathbb{R}_+ = set of non negative real numbers

Scope of Learning Problems Modeled

Generalized Measure of Success

- Given \mathcal{H} , domain $Z = X \times Y$,

- Let $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$, where

ℓ = loss function

\mathbb{R}_+ = set of non negative real numbers

- In unsupervised problems, Z is not $X \times Y$, since there is "no" Y

Scope of Learning Problems Modeled

Generalized Measure of Success

- Given \mathcal{H} , domain $Z = X \times Y$,

- Let $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$, where

ℓ = loss function

\mathbb{R}_+ = set of non negative real numbers

- In unsupervised problems, Z is not $X \times Y$, since there is "no" Y
- Definition:** Risk function of $h \in \mathcal{H}$ wrt D over Z

$$L_D(h) \equiv \mathbb{E}_{Z \sim D} [\ell(h, z)]$$

This is the expected loss of h sampled from Z using D .

Scope of Learning Problems Modeled

Definition: Empirical Risk

Given $S = \{z_1, z_2, \dots, z_m\} \in Z^m$

$$L_S(h) \equiv \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$$

Loss Functions

- **0-1 Loss:** $Z \in X \times Y$

$$\ell_{01}(h, (x, y)) \equiv \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{if } h(x) \neq y \end{cases}$$

Loss Functions

- **0-1 Loss:** $Z \in X \times Y$

$$\ell_{01}(h, (x, y)) \equiv \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{if } h(x) \neq y \end{cases}$$

⇒ Binary/multiclass classification problem

Scope of Learning Problems Modeled

Loss Functions

- **0-1 Loss:** $Z \in X \times Y$

$$\ell_{01}(h, (x, y)) \equiv \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{if } h(x) \neq y \end{cases}$$

⇒ Binary/multiclass classification problem

⇒ For random variable $\alpha \in \{0, 1\}$

$$\mathbb{E}_{\alpha \sim D} [\alpha] = \mathbb{P}_{\alpha \sim D} (\alpha = 1)$$

Scope of Learning Problems Modeled

Loss Functions

⇒ Then the different definitions of $L_D(h)$ coincide

$$\text{a) } L_D(h) \equiv \mathbb{E}_{z \sim D} (\ell(h, z))$$

$$\text{b) } L_D(h) \equiv \mathbb{P}_{z \sim D} (h(x) \neq y)$$

Loss Functions

⇒ Then the different definitions of $L_D(h)$ coincide

$$\text{a) } L_D(h) \equiv \mathbb{E}_{Z \sim D} (\ell(h, z))$$

$$\text{b) } L_D(h) \equiv \mathbb{P}_{z \sim D} (h(x) \neq y)$$

- **Square Loss:** $Z \in X \times Y$

$$\ell_{01}(h(x, y)) \equiv (h(x) - y)^2$$

Scope of Learning Problems Modeled

Agnostic PAC Learnability For General Loss Functions

Definition: Same as before with \mathcal{H}, Z, D, ℓ where $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$ and

$$L_D(h) = \mathbb{E}_{Z \sim D}[\ell(h, z)]$$

Sidenote: Proper vs Representational Independent Learning

- $\mathcal{H} \in \mathcal{H}'$

Scope of Learning Problems Modeled

Agnostic PAC Learnability For General Loss Functions

Definition: Same as before with \mathcal{H}, Z, D, ℓ where $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$ and

$$L_D(h) = \mathbb{E}_{Z \sim D}[\ell(h, z)]$$

Sidenote: Proper vs Representational Independent Learning

- $\mathcal{H} \in \mathcal{H}'$
- Extend *Loss Function* to $\mathcal{H}' \times Z \rightarrow \mathbb{R}_+$

Scope of Learning Problems Modeled

Agnostic PAC Learnability For General Loss Functions

Definition: Same as before with \mathcal{H}, Z, D, ℓ where $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$ and

$$L_D(h) = \mathbb{E}_{Z \sim D}[\ell(h, z)]$$

Sidenote: Proper vs Representational Independent Learning

- $\mathcal{H} \in \mathcal{H}'$
- Extend *Loss Function* to $\mathcal{H}' \times Z \rightarrow \mathbb{R}_+$
- Return algorithm from \mathcal{H}' instead of \mathcal{H} (*This is representational independent learning*) as long as it satisfies

$$L_D(h' \in \mathcal{H}') \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$$

Scope of Learning Problems Modeled

Agnostic PAC Learnability For General Loss Functions

Definition: Same as before with \mathcal{H}, Z, D, ℓ where $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$ and

$$L_D(h) = \mathbb{E}_{Z \sim D}[\ell(h, z)]$$

Sidenote: Proper vs Representational Independent Learning

- $\mathcal{H} \in \mathcal{H}'$
- Extend *Loss Function* to $\mathcal{H}' \times Z \rightarrow \mathbb{R}_+$
- Return algorithm from \mathcal{H}' instead of \mathcal{H} (*This is representational independent learning*) as long as it satisfies

$$L_D(h' \in \mathcal{H}') \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$$

- Proper learning is when the algorithm outputs an h from \mathcal{H} not from \mathcal{H}'