Introduction to the Theory of Machine Learning The Language of Learning Theory

Rachit Nimavat

IIIT Surat

June 27, 2025

Outline



- 2 The Language of Learning Theory
- The PAC Model

Recognize speech, images, stock market patterns

Why ML?

Machine learning can be used to...

Recognize speech, images, stock market patterns Classify documents, predict protein sequences

Recognize speech, images, stock market patterns Classify documents, predict protein sequences Design course,

Recognize speech, images, stock market patterns

Classify documents, predict protein sequences

Design course, generate course slides, produce video of an instructor teaching those slides, generate assignments,

Recognize speech, images, stock market patterns

Classify documents, predict protein sequences

Design course, generate course slides, produce video of an instructor teaching those slides, generate assignments, solve assignments,

Recognize speech, images, stock market patterns

Classify documents, predict protein sequences

Design course, generate course slides, produce video of an instructor teaching those slides, generate assignments, solve assignments, grade assignments

Recognize speech, images, stock market patterns

Classify documents, predict protein sequences

Design course, generate course slides, produce video of an instructor teaching those slides, generate assignments, solve assignments, grade assignments



Understand Learning!

What kinds of tasks can we hope to learn

Understand Learning!

What kinds of **tasks** can we *hope* to learn What kind of **data** is needed to learn that task

Understand Learning!

What kinds of **tasks** can we *hope* to learn What kind of **data** is needed to learn that task What kind of **guarantees** we can achieve with that data

Understand Learning!

What kinds of **tasks** can we *hope* to learn What kind of **data** is needed to learn that task What kind of **guarantees** we can achieve with that data

Optimize Learning!

How efficient our learning process can be

Understand Learning!

What kinds of **tasks** can we *hope* to learn What kind of **data** is needed to learn that task What kind of **guarantees** we can achieve with that data

Optimize Learning!

How efficient our learning process can be Relating data size, processing time, space, accuracy

Understand Learning!

What kinds of **tasks** can we *hope* to learn What kind of **data** is needed to learn that task What kind of **guarantees** we can achieve with that data

Optimize Learning!

How efficient our learning process can be

Relating data size, processing time, space, accuracy

Improve Learning!

How does biases in data affect learning

Understand Learning!

What kinds of **tasks** can we *hope* to learn What kind of **data** is needed to learn that task What kind of **guarantees** we can achieve with that data

Optimize Learning!

How efficient our learning process can be

Relating data size, processing time, space, accuracy

Improve Learning!

How does biases in data affect learning

Preserve Privacy? Unlabeled data? Interactive environments? ...



Example:

 $\begin{aligned} \mathcal{X} \text{ is the set of all images} \\ \mathcal{Y} &= \{\mathsf{CAT},\mathsf{DOG}\} \text{ is the set of all labels} \\ x \in \mathcal{X} \text{ is a single image} \\ y &= f(x) \text{ is its true label} \end{aligned}$



Example:

 ${\mathcal X}$ is the set of all images

 $\mathcal{Y} = \{\mathsf{CAT}, \mathsf{DOG}\}$ is the set of all labels

 $x \in \mathcal{X}$ is a *single* image

y = f(x) is its true label

Split sampled data into training and test sets



Example:

 $\ensuremath{\mathcal{X}}$ is the set of all images

 $\mathcal{Y} = \{\mathsf{CAT},\mathsf{DOG}\}$ is the set of all labels

 $x \in \mathcal{X}$ is a *single* image

y = f(x) is its true label

The learning phase



Example:

 ${\mathcal X}$ is the set of all images

 $\mathcal{Y} = \{\mathsf{CAT}, \mathsf{DOG}\}$ is the set of all labels

 $x \in \mathcal{X}$ is a *single* image

y = f(x) is its true label

Output of model our guess h for true f

 $\hat{y} = h(x)$ is the models output label



Example:

 ${\mathcal X}$ is the set of all images

 $\mathcal{Y} = \{\mathsf{CAT}, \mathsf{DOG}\}$ is the set of all labels

 $x \in \mathcal{X}$ is a *single* image

y = f(x) is its true label

Evaluating performance

 $\hat{y} = h(x)$ is the models output label Cost: $\mathbf{1}[y \neq \hat{y}]$

_

Physicists	Mathematicians
Start with: Observations / Data	Start with: Axioms / Laws
	λ

Physicists	Mathematicians
Start with: Observations / Data	Start with: Axioms / Laws
	7
~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	
Goal:	Goal:
Find laws that correctly explain the observations/data	Find theorems/observations that follow from the axioms

Physicists	Mathematicians
Start with: Observations / Data	Start with: Axioms / Laws
- <b>X</b> *	λ
Goal:	Goal:
Find laws that correctly explain the observations/data	Find theorems/observations that follow from the axioms

#### Where Machine Learning Fits?

Find rules that make good predictions about unseen data.

The Core Question of Generalization

🗱 Model trained on training data

The Core Question of Generalization

🗱 Model trained on training data

Performs well on test data

#### The Core Question of Generalization

- 🗱 Model trained on training data
- Performs well on test data
- ? Confidence about performance on new, unseen data?

#### The Core Question of Generalization

- 🗱 Model trained on training data
- Performs well on test data
- ? Confidence about performance on new, unseen data?

#### The Fundamental Assumption of ML

All data points are drawn independently from a fixed, but unknown, probability distribution  $\mathcal{D}.$ 

- $\mathcal{X}:$  The set of all data points
- $\mathcal{Y}:$  The set of all labels
- $\mathcal{D}:$  The underlying, but unknown joint probability distribution over  $\mathcal{X}\times\mathcal{Y}.$

- $\mathcal{X}:$  The set of all data points
- $\mathcal{Y}:$  The set of all labels
- $\mathcal{D}:$  The underlying, but unknown joint probability distribution over  $\mathcal{X}\times\mathcal{Y}.$

#### Single Labeled Example

A data point x with (true) label y is drawn from  $\ensuremath{\mathcal{D}}$ 

 $(x, y) \sim \mathcal{D}$ 

- $\mathcal{X} \colon$  The set of all data points
- $\mathcal{Y}:$  The set of all labels
- $\mathcal{D}:$  The underlying, but unknown joint probability distribution over  $\mathcal{X}\times\mathcal{Y}.$

#### Single Labeled Example

A data point x with (true) label y is drawn from  $\ensuremath{\mathcal{D}}$ 

$$(x, y) \sim \mathcal{D}$$

#### Training Set

m labeled examples  $\{(x_1,y_1),\ldots,(x_m,y_m)\}$  are drawn independently and identically from  $\mathcal D$ 

 $S \sim \mathcal{D}^m$ 

- $\mathcal{X} \colon$  The set of all data points
- $\mathcal{Y}:$  The set of all labels
- $\mathcal{D}:$  The underlying, but unknown joint probability distribution over  $\mathcal{X}\times\mathcal{Y}.$

#### Single Labeled Example

A data point x with (true) label y is drawn from  $\ensuremath{\mathcal{D}}$ 

$$(x, y) \sim \mathcal{D}$$

#### Training Set

m labeled examples  $\{(x_1,y_1),\ldots,(x_m,y_m)\}$  are drawn independently and identically from  $\mathcal D$ 

 $S \sim \mathcal{D}^m$ 

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

"The training set S drew m i.i.d examples from a single, unknown, underlying distribution  $\mathcal{D}$ ."

Training Set

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

Training Set

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

The Learning Algorithm

 $\mathcal{A}$  does some computation over *S* to produce a **hypothesis** (prediction rule) *h*.

Training Set

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

#### The Learning Algorithm

 $\mathcal{A}$  does some computation over S to produce a **hypothesis** (prediction rule) h. Under the hood,  $\mathcal{A}$  searches within a **Hypothesis Class** (pre-defined 'universe' of possible hypotheses)  $\mathcal{H}$ .

Training Set

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

#### The Learning Algorithm

 $\mathcal{A}$  does some computation over S to produce a **hypothesis** (prediction rule) h. Under the hood,  $\mathcal{A}$  searches within a **Hypothesis Class** (pre-defined 'universe' of possible hypotheses)  $\mathcal{H}$ .

> ${\cal H}$  is the "library" of models our algorithm can choose from: linear separators, width-100 depth-10 neural networks, 100-line python programs...

Training Set

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

#### The Learning Algorithm

 $\mathcal{A}$  does some computation over S to produce a **hypothesis** (prediction rule) h. Under the hood,  $\mathcal{A}$  searches within a **Hypothesis Class** (pre-defined 'universe' of possible hypotheses)  $\mathcal{H}$ .

#### Some ways in which $\mathcal{A}$ may work:

**SVM**: Out of all linear separators, pick one with maximum margin **Gradient Descent**: Random initialization of weights and back-propagate gradients until we hit a local minima

LLM: Ask ChatGPT to generate a 100 line python program

Training Set

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

#### The Learning Algorithm

 $\mathcal{A}$  does some computation over S to produce a **hypothesis** (prediction rule) h. Under the hood,  $\mathcal{A}$  searches within a **Hypothesis Class** (pre-defined 'universe' of possible hypotheses)  $\mathcal{H}$ .

Some ways in which  $\mathcal{A}$  may **SVM**: Out of all linear se **Gradient Descent**: Ranc until we hit a local minim **LLM**: Ask ChatGPT to generate a 100 line python program

The learning algorithm  $\mathcal{A}$  reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

 $\mathcal{A}(S)$  outputs  $h \in \mathcal{H}$ 

The learning algorithm A reports a hypothesis h when given training data S.

 $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^m$  $\mathcal{A}(S) \text{ outputs } h \in \mathcal{H}$ 

Empirical Error: what we can calculate

$$\operatorname{err}_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i_1}^m \mathbf{1}[h(x_i) \neq y_i]$$

The learning algorithm A reports a hypothesis h when given training data S.

 $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^m$  $\mathcal{A}(S) \text{ outputs } h \in \mathcal{H}$ 

Empirical Error: what we can calculate

$$\operatorname{err}_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i_1}^m \mathbf{1}[h(x_i) \neq y_i]$$

and what  $\mathcal{A}$  typically minimizes:

$$\mathcal{A}(S) \equiv \operatorname*{arg\,min}_{h \in \mathcal{H}} \left( \operatorname{err}_{S}(h) \right)$$

The learning algorithm A reports a hypothesis h when given training data S.

 $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^m$  $\mathcal{A}(S) \text{ outputs } h \in \mathcal{H}$ 

Empirical Error: what we can calculate

$$\operatorname{err}_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i_1}^m \mathbf{1}[h(x_i) \neq y_i]$$

and what  $\mathcal{A}$  typically minimizes:

$$\mathcal{A}(S) \equiv \operatorname*{arg\,min}_{h \in \mathcal{H}} \left( \operatorname{err}_{S}(h) \right)$$

True Error: what we want to minimize

$$\operatorname{err}_{\mathcal{D}}(h) = \operatorname{\mathbf{Pr}}_{x \sim \mathcal{D}}[h(x) \neq f(x)]$$

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^m$$
$$\mathcal{A}(S) \equiv \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} (\operatorname{err}_S(h))$$

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^n$$
$$\mathcal{A}(S) \equiv \arg\min(\operatorname{err}_S(h))$$

 $h \in \mathcal{H}$ 

What is the true error  $err_{\mathcal{D}}(h)$  of the reported hypothesis h?

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^m$$
$$\mathcal{A}(S) \equiv \arg\min(\operatorname{err}_S(h))$$

 $\widetilde{h} \in \mathcal{H}$ 

What is the true error  $err_{\mathcal{D}}(h)$  of the reported hypothesis *h*?

Goal: perform well on fresh data:

 $\mathbf{Pr}_{x\sim\mathcal{D}}\left[h(x)\neq f(x)\right]<\epsilon$ 

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim \mathcal{D}^n$$
$$\mathcal{A}(S) \equiv \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} (\operatorname{err}_S(h))$$

What is the true error  $err_{\mathcal{D}}(h)$  of the reported hypothesis *h*?

Goal: perform well on fresh data:

 $\mathbf{Pr}_{x\sim\mathcal{D}}\left[h(x)\neq f(x)\right]<\epsilon$ 

 $\epsilon$  is a small positive number representing accuracy.  $\epsilon=0.1$  means: hypothesis h is correct on 90% of the data points*.

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

 $\mathcal{A}(S) \equiv \operatorname*{arg\,min}_{h \in \mathcal{H}} \left( \operatorname{err}_{S}(h) \right)$ 

What is the true error  $\operatorname{err}_{\mathcal{D}}(h)$  of the reported hypothesis *h*? Goal: perform well on fresh data:

 $\mathbf{Pr}_{x\sim\mathcal{D}}\left[h(x)\neq f(x)\right]<\epsilon$ 

**Approximately Correct!** 

 $\epsilon$  is a small positive number representing accuracy.  $\epsilon = 0.1$  means: hypothesis *h* is *correct* on 90% of the data points^{*}.

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

 $\mathcal{A}(S) \equiv \operatorname*{arg\,min}_{h \in \mathcal{H}} \left( \operatorname{err}_{S}(h) \right)$ 

What is the true error  $err_{\mathcal{D}}(h)$  of the reported hypothesis *h*?

Goal: perform well on fresh data:

 $\mathbf{Pr}_{x\sim\mathcal{D}}\left[h(x)\neq f(x)\right]<\epsilon$ 

**Approximately Correct!** 

Really??

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

 $\mathcal{A}(S) \equiv \operatorname*{arg\,min}_{h \in \mathcal{H}} \left( \operatorname{err}_{S}(h) \right)$ 

What is the true error  $err_{\mathcal{D}}(h)$  of the reported hypothesis *h*?

Goal: perform well on fresh data:

 $\mathbf{Pr}_{x\sim\mathcal{D}}\left[h(x)\neq f(x)\right]<\epsilon$ 

Approximately Correct!

#### Really??

S is drawn **randomly** from the distribution  $\mathcal{D}$ .

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

 $\mathcal{A}(S) \equiv \operatorname*{arg\,min}_{h \in \mathcal{H}} \left( \operatorname{err}_{S}(h) \right)$ 

What is the true error  $err_{\mathcal{D}}(h)$  of the reported hypothesis *h*?

Goal: perform well on fresh data:

 $\mathbf{Pr}_{x\sim\mathcal{D}}\left[h(x)\neq f(x)\right]<\epsilon$ 

**Approximately Correct!** 

#### Really??

Can only guarantee low error with high probability over the draw of training data S.

 $\mathbf{Pr}_{\mathcal{S}\sim\mathcal{D}^m}\left[\mathrm{err}_{\mathcal{D}}(\mathbf{h})<\epsilon\right]>1-\delta$ 

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

 $\mathcal{A}(S) \equiv \operatorname*{arg\,min}_{h \in \mathcal{H}} \left( \operatorname{err}_{S}(h) \right)$ 

What is the true error  $err_{\mathcal{D}}(h)$  of the reported hypothesis *h*?

Goal: perform well on fresh data:

 $\mathbf{Pr}_{x\sim\mathcal{D}}\left[h(x)\neq f(x)\right]<\epsilon$ 

Approximately Correct!

#### Really??

Can only guarantee low error with high probability over the draw of training data S.

 $\mathbf{Pr}_{\mathcal{S}\sim\mathcal{D}^m}\left[\mathrm{err}_{\mathcal{D}}(h)<\epsilon
ight]>1-\delta$ 

 $\delta$  is a small positive number representing probability of failure.

 $\delta=0.2$  and  $\epsilon=0.1$  means:

 $80 \rm \%$  of the time,  ${\cal A}$  outputs a hypothesis that is correct on  $90 \rm \%$  of the data points.

The learning algorithm A reports a hypothesis h when given training data S.

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \sim \mathcal{D}^m$$

 $\mathcal{A}(S) \equiv \operatorname*{arg\,min}_{h \in \mathcal{H}} \left( \operatorname{err}_{S}(h) \right)$ 

What is the true error  $err_{\mathcal{D}}(h)$  of the reported hypothesis *h*?

Goal: perform well on fresh data:

$$\mathbf{Pr}_{x\sim\mathcal{D}}\left[h(x)\neq f(x)\right]<\epsilon$$

**Approximately Correct!** 

#### Really??

Can only guarantee low error with high probability over the draw of training data S.

$$\Pr_{S \sim \mathcal{D}^m} \left[ \operatorname{err}_{\mathcal{D}}(h) < \epsilon \right] > 1 - \delta$$

"Probably"

 $\delta$  is a small positive number representing probability of failure.

 $\delta = 0.2$  and  $\epsilon = 0.1$  means:

80% of the time,  $\mathcal A$  outputs a hypothesis that is correct on 90% of the data points.

Definition 3.1

A hypothesis class  ${\mathcal H}$  is **PAC-learnable** if

Definition 3.1

A hypothesis class  ${\cal H}$  is **PAC-learnable** if there exists an algorithm  ${\cal A},$  that,

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ ,

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ ,

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data,

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any target  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Are we done?



Leslie Valiant. 1984

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Considerations

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Considerations

1 if there exists an algorithm A, that, ... : What is this algorithm A?

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Considerations

1 if there exists an algorithm A, that, ... : What is this algorithm A?

How to find A? How many samples (*m*) needed for  $(\epsilon, \delta)$ -PAC guarantee? Efficiency of A?

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Considerations

- 1 if there exists an algorithm A, that, ... : What is this algorithm A?
- 2 any target  $f \in \mathcal{H}$ ...: What about target function f?

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Considerations

- 1 if there exists an algorithm A, that, ... : What is this algorithm A?
- 2 any target  $f \in \mathcal{H}$ ...: What about target function f?

What if  $f \notin \mathcal{H}$ ? What if *no true target function exists*? Noise/Errors in labels or features? Errros in training data?

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Considerations

- 1 if there exists an algorithm A, that, ... : What is this algorithm A?
- 2 any target  $f \in \mathcal{H}$ ...: What about target function f?
- 3 any distribution  $\mathcal{D}$ ...: Does such  $\mathcal{D}$  truly exist?

#### Definition 3.1

A hypothesis class  $\mathcal{H}$  is **PAC-learnable** if there exists an algorithm  $\mathcal{A}$ , that, for any desired accuracy  $\epsilon > 0$ , any desired confidence  $\delta < 1$ , any *target*  $f \in \mathcal{H}$ , any distribution  $\mathcal{D}$  over data, with probability at least  $1 - \delta$ , reports a hypothesis  $h \in \mathcal{H}$  with  $\operatorname{err}_{\mathcal{D}}(h) < \epsilon$ .

#### Considerations

- 1 if there exists an algorithm A, that, ... : What is this algorithm A?
- 2 any target  $f \in \mathcal{H}$ ...: What about target function f?
- 3 any distribution  $\mathcal{D}$ ...: Does such  $\mathcal{D}$  truly exist?

Future data may diverge Possibly no *fixed*  $\mathcal{D}$  even for training data

#### Questions?

### Break Time!

# **Questions?**

#### Break Time!

# **Questions?**

Let's take a 10-minute break. We'll resume at 11am.



An excellent introductory course on Theory of ML: TTIC 31250 by Avrim Blum

THE textbook of ML: Understanding Machine Learning: From Theory to Algorithms by Shai Shalev-Shwartz and Shai Ben-David