

COMS 4771
Assignment Project Exam Help
Introduction to Machine Learning
<https://powcoder.com>

Add WeChat powcoder

Nakul Verma

Towards formalizing 'learning'

What does it mean to **learn** a concept?

- Gain knowledge or experience of the concept.

Assignment Project Exam Help

<https://powcoder.com>

The basic process of **learning**

- Observe a phenomenon
- Construct a model from observations
- Use that model to make decisions / predictions

Add WeChat powcoder

How can we make this more precise?

A statistical machinery for learning

Phenomenon of interest:

Input space: X Output space: Y

There is an **unknown** distribution \mathcal{D} over $(X \times Y)$

The learner **observes** m examples $(x_1, y_1), \dots, (x_m, y_m)$ drawn from \mathcal{D}

Assignment Project Exam Help

Construct a model:

<https://powcoder.com>

Machine learning

Let \mathcal{F} be a collection of models, where each $f : X \rightarrow Y$ **predicts** y given x

From m observations, **select** a model $f_m \in \mathcal{F}$ which predicts **well**.

$$\text{err}(f) := \mathbb{P}_{(x,y) \sim \mathcal{D}} [f(x) \neq y] \quad (\text{generalization error of } f)$$

We can say that we have **learned** the phenomenon if

$$\text{err}(f_m) - \text{err}(f^*) \leq \epsilon \quad f^* := \operatorname{argmin}_{f \in \mathcal{F}} \text{err}(f)$$

for any tolerance level $\epsilon > 0$ of our choice.

PAC Learning

For all tolerance levels $\epsilon > 0$, and all confidence levels $\delta > 0$, if there exists some model selection algorithm \mathcal{A} that selects $f_m^{\mathcal{A}} \in \mathcal{F}$ from m observations ie, $\mathcal{A} : (x_i, y_i)_{i=1}^m \mapsto f_m^{\mathcal{A}}$, and has the property:

with probability at least $1 - \delta$ over the draw of the sample,

$$\text{err}(f_m^{\mathcal{A}}) - \text{err}(f^*) \leq \epsilon$$

We call

- The model class \mathcal{F} is PAC-learnable.
- If the m is polynomial in $\frac{1}{\epsilon}$ and $\frac{1}{\delta}$, then \mathcal{F} is **efficiently** PAC-learnable

A popular algorithm:

Empirical risk minimizer (ERM) algorithm

$$f_m^{\text{ERM}} := \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{f(x_i) \neq y_i\}$$

PAC learning simple model classes

Theorem (finite size \mathcal{F}):

Pick any tolerance level $\epsilon > 0$, and any confidence level $\delta > 0$

let $(x_1, y_1), \dots, (x_m, y_m)$ be m examples drawn from an unknown \mathcal{D}

if $m \geq C \cdot \frac{1}{\epsilon^2} \ln \frac{|\mathcal{F}|}{\delta}$, then with probability at least $1 - \delta$

$$\text{err}(f_m^{\text{EBM}}) - \text{err}(f^*) \leq \epsilon$$

Add WeChat powcoder

\mathcal{F} is efficiently PAC learnable

Occam's Razor Principle:

All things being equal, usually the simplest explanation of a phenomenon is a good hypothesis.

Simplicity = representational succinctness

Proof sketch

Define:

$$\text{err}(f) := \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbf{1}\{f(x) \neq y\}] \quad \text{err}_m(f) := \frac{1}{m} \sum_{i=1}^m [\mathbf{1}\{f(x_i) \neq y_i\}]$$

(generalization error of f) (sample error of f)

Assignment Project Exam Help

<https://powcoder.com>

We need to analyze:

Add WeChat powcoder

$$\begin{aligned} \text{err}(f_m^{\text{ERM}}) - \text{err}(f^*) &= \text{err}(f_m^{\text{ERM}}) - \text{err}_m(f_m^{\text{ERM}}) \\ &\quad + \text{err}_m(f_m^{\text{ERM}}) - \text{err}_m(f^*) \\ &\quad + \text{err}_m(f^*) - \text{err}(f^*) \\ &\leq 2 \sup_{f \in \mathcal{F}} |\text{err}(f) - \text{err}_m(f)| \end{aligned}$$

≤ 0

Uniform deviations of expectation of a random variable to the sample

Proof sketch

Fix any $f \in \mathcal{F}$ and a sample (x_i, y_i) , define random variable

$$\mathbf{Z}_i^f := \mathbf{1}\{f(x_i) \neq y_i\}$$

$$\mathbb{E}[\mathbf{Z}_1^f] \quad \text{Assignment Project Exam Help} \quad \frac{1}{m} \sum_{i=1}^m [\mathbf{Z}_i^f]$$

(generalization error of f) <https://powcoder.com> (sample error of f)

Add WeChat powcoder

Lemma (Chernoff-Hoeffding bound '63):

Let $\mathbf{Z}_1, \dots, \mathbf{Z}_m$ be m Bernoulli r.v. drawn independently from $\mathbf{B}(\mathbf{p})$.

for any tolerance level $\epsilon > 0$

$$\mathbb{P}_{\mathbf{Z}_i} \left[\left| \frac{1}{m} \sum_{i=1}^m [\mathbf{Z}_i] - \mathbb{E}[\mathbf{Z}_1] \right| > \epsilon \right] \leq 2e^{-2\epsilon^2 m}.$$

A classic result in **concentration of measure**, proof later

Proof sketch

Need to analyze

$$\mathbb{P}_{(x_i, y_i)} \left[\text{exists } f \in \mathcal{F}, \left| \frac{1}{m} \sum_{i=1}^m [\mathbf{Z}_i^f] - \mathbb{E}[\mathbf{Z}_1^f] \right| > \epsilon \right]$$

Assignment Project Exam Help

$$\leq \sum_{f \in \mathcal{F}} \mathbb{P}_{(x_i, y_i)} \left[\left| \frac{1}{m} \sum_{i=1}^m [\mathbf{Z}_i^f] - \mathbb{E}[\mathbf{Z}_1^f] \right| > \epsilon \right]$$

$$\leq \frac{2|\mathcal{F}|}{e^{-2\epsilon^2 m}} \leq \delta$$

Equivalently, by choosing $m \geq \frac{1}{2\epsilon^2} \ln \frac{2|\mathcal{F}|}{\delta}$ with probability at least $1 - \delta$,
for **all** $f \in \mathcal{F}$

$$\left| \frac{1}{m} \sum_{i=1}^m [\mathbf{Z}_i^f] - \mathbb{E}[\mathbf{Z}_1^f] \right| = \left| \text{err}_m(f) - \text{err}(f) \right| \leq \epsilon$$



PAC learning simple model classes

Theorem (Occam's Razor):

Pick any tolerance level $\epsilon > 0$, and any confidence level $\delta > 0$

let $(x_1, y_1), \dots, (x_m, y_m)$ be m examples drawn from an unknown \mathcal{D}

if $m \geq C \cdot \frac{1}{\epsilon^2} \ln \frac{|\mathcal{F}|}{\delta}$, then with probability at least $1 - \delta$

$$\text{err}(f_m^{\text{EBM}}) - \text{err}(f^*) \leq \epsilon$$

Add WeChat powcoder

\mathcal{F} is efficiently PAC learnable

One thing left...

Still need to prove:

Lemma (Chernoff-Hoeffding bound '63):

Let Z_1, \dots, Z_m be m Bernoulli r.v. drawn independently from $\mathbf{B}(p)$.
for any tolerance level $\epsilon > 0$

$$\mathbb{P}_{Z_i} \left[\left| \frac{1}{m} \sum_{i=1}^m [Z_i] - \mathbb{E}[Z_1] \right| > \epsilon \right] \leq 2e^{-2\epsilon^2 m}.$$

How sample average deviates from true average

*as a function of
number of samples (m)*

*Need to analyze: How does the probability measure
concentrates towards a central value (like mean)*

Detour: Concentration of Measure

Let's start with something simple:

Let X be a non-negative random variable.

For a given constant $c > 0$, what is: $\mathbb{P}[X \geq c]$?

Assignment Project Exam Help

<https://powcoder.com>

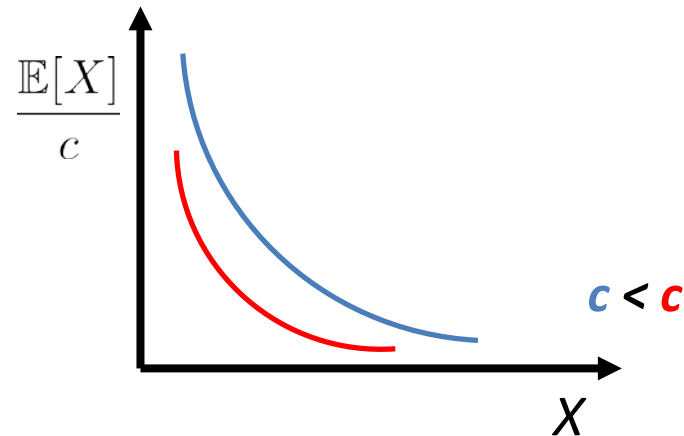
Markov's Inequality

Why?

Add WeChat powcoder

Observation $c \cdot \mathbf{1}[X \geq c] \leq X$

Take expectation on both sides.



Concentration of Measure

Using Markov to bound deviation from mean...

Let X be a random variable (not necessarily non-negative).

Want to examine: $\mathbb{P}[|X - \mathbb{E}X| \geq c]$ for some given constant $c > 0$

Observation: <https://powcoder.com>

$$\mathbb{P}[|X - \mathbb{E}X| \geq c] = \mathbb{P}[(X - \mathbb{E}X)^2 \geq c^2]$$

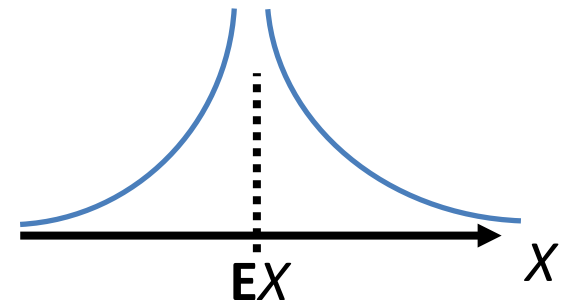
$$\leq \frac{\mathbb{E}(X - \mathbb{E}X)^2}{c^2}$$

by Markov's Inequality

$$= \frac{\text{Var}(X)}{c^2}$$

Chebyshev's Inequality

*True for **all** distributions!*



Concentration of Measure

Sharper estimates using an exponential!

Let X be a random variable (not necessarily non-negative).

For some given constant $c > 0$

Assignment Project Exam Help

Observation:

<https://powcoder.com>

$$\mathbb{P}[X \geq c]$$

Add WeChat powcoder for any $t > 0$

$$\leq \frac{\mathbb{E}[e^{tX}]}{e^{tc}}$$

by Markov's Inequality

*This is called Chernoff's
bounding method*

Concentration of Measure

Now, Given X_1, \dots, X_m i.i.d. random variables (assume $0 \leq X_i \leq 1$)

$$\mathbb{P}\left[\frac{1}{m} \sum_{i=1}^m X_i - \mathbb{E}X_1 \geq c\right] = \mathbb{P}\left[\sum_{i=1}^m X_i - m\mathbb{E}X_1 \geq mc\right] \quad \text{Define } Y_i := X_i - \mathbb{E}X_i$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat [powcoder](#) By Chernoeff's bounding technique

$$= \frac{1}{e^{tmc}} \prod_{i=1}^m \mathbb{E}[e^{tY_i}] \quad Y_i \text{ i.i.d.}$$

$$\mathbb{E}[e^{tY_i}] \leq e^{t^2/8}$$

$$t = 4c$$

$$\leq e^{t^2 m/8 - tmc}$$

$$\leq e^{-2c^2 m}$$

This **implies** the Chernoff-Hoeffding bound!

Back to Learning Theory!

Theorem (Occam's Razor):

Pick any tolerance level $\epsilon > 0$, and any confidence level $\delta > 0$

let $(x_1, y_1), \dots, (x_m, y_m)$ be m examples drawn from an unknown \mathcal{D}

if $m \geq C \cdot \frac{1}{\epsilon^2} \ln \frac{|\mathcal{F}|}{\delta}$, then with probability at least $1 - \delta$

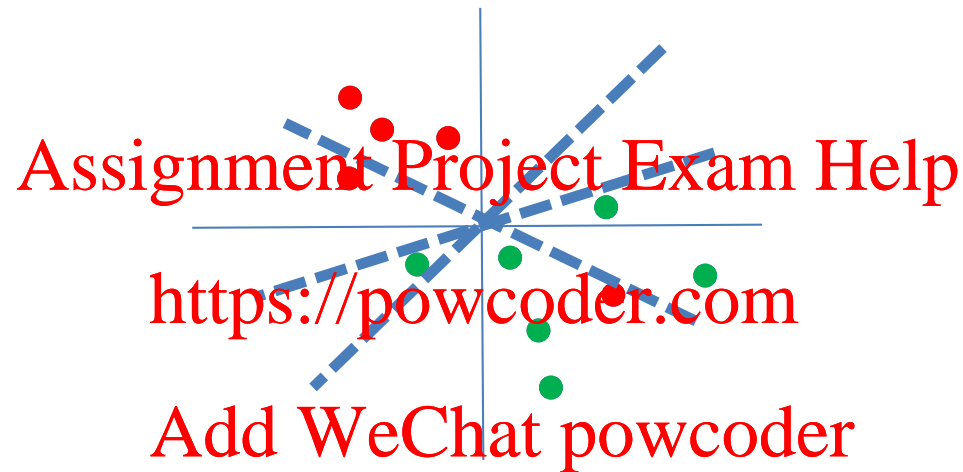
$$\text{err}(f_m^{\text{EBM}}) - \text{err}(f^*) \leq \epsilon$$

Add WeChat powcoder

\mathcal{F} is efficiently PAC learnable

Learning general concepts

Consider linear classification



$$\mathcal{F} = \left\{ \begin{array}{c} \text{dashed lines} \end{array} \right\} \quad |\mathcal{F}| = \infty$$

Occam's Razor bound is ineffective

VC Theory

Need to capture the true richness of \mathcal{F}

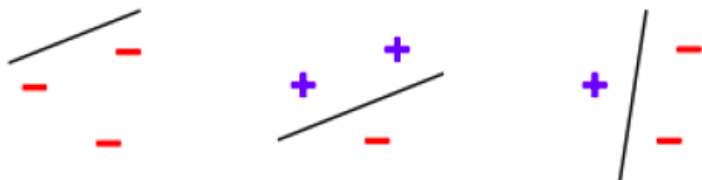
Definition (Vapnik-Chervonenkis or VC dimension):

We say that a model class \mathcal{F} has VC dimension d if d is the largest set of points $x_1, \dots, x_d \subset X$ such that for all possible labelings of x_1, \dots, x_d there exists some $f \in \mathcal{F}$ that achieves that labelling.

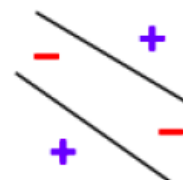
Add WeChat powcoder

Example: \mathcal{F} = linear classifiers in \mathbb{R}^2

linear classifiers can realize all possible labellings of 3 points



linear classifiers CANNOT realize all labellings of 4 points



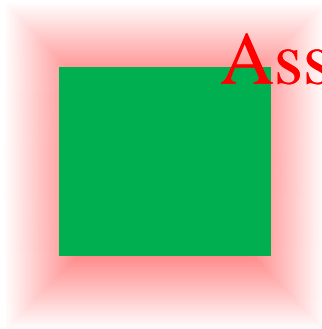
$$VC(\mathcal{F}) = 3$$

VC Dimension

Another example:

\mathcal{F} = Rectangles in \mathbf{R}^2

$$VC(\mathcal{F}) = 4$$



Assignment Project Exam Help

The class of rectangles

<https://powcoder.com> *cannot realize this labelling*

Add WeChat powcoder

VC dimension:

- A **combinatorial concept** to capture the true richness of \mathcal{F}
- Often (but not always!) proportional to the degrees-of-freedom or the number of independent parameters in \mathcal{F}

VC Theorem

Theorem (Vapnik-Chervonenkis '71):

Pick any tolerance level $\epsilon > 0$, and any confidence level $\delta > 0$

let $(x_1, y_1), \dots, (x_m, y_m)$ be m examples drawn from an unknown \mathcal{D}

if $m \geq C \cdot \frac{VC(\mathcal{F}) \ln(1/\delta)}{\epsilon^2}$, then with probability at least $1 - \delta$

$$\text{err}(f_m^{\text{EBM}}) - \text{err}(f^*) \leq \epsilon$$

Add WeChat powcoder

\mathcal{F} is efficiently PAC learnable

VC Theorem \rightarrow Occam's Razor Theorem

Tightness of VC bound

Theorem (VC lower bound):

Let \mathcal{A} be any model selection algorithm that given m samples, returns a model from \mathcal{F} , that is, $\mathcal{A} : (x_i, y_i)_{i=1}^m \mapsto f_m^{\mathcal{A}}$

For all tolerance levels $0 < \epsilon < 1$, and all confidence levels $0 < \delta < 1/4$, there exists a distribution \mathcal{D} such that if $m \leq C \cdot \frac{VC(\mathcal{F})}{\epsilon^2}$

<https://powcoder.com>

$$\mathbb{P}_{(x_i, y_i)} \left[| \text{err}(f_m^{\mathcal{A}}) - \text{err}(f^*) | > \epsilon \right] > \delta$$

Some implications

- VC dimension of a model class fully characterizes its learning ability!

Assignment Project Exam Help

<https://powcoder.com>

- Results are agnostic to the underlying distribution.

Add WeChat powcoder

One algorithm to rule them all?

From our discussion it may seem that ERM algorithm is universally consistent.

This is not the case!

Assignment Project Exam Help

Theorem (no free lunch, Devroye (82))

Pick any sample size m , any algorithm \mathcal{A} and any $\epsilon > 0$

There exists a distribution \mathcal{D} such that

$$\text{err}(f_m^{\mathcal{A}}) > 1/2 - \epsilon$$

while the Bayes optimal error, $\min_f \text{err}(f) = 0$

Further refinements and extensions

- How to do **model class** selection? Structural risk results.

- Dealing with **kernel** – Functional theory

- Incorporating **priors** over the models – PAC-Bayes theory

Add WeChat powcoder

- Is it possible to get **distribution dependent** bound? Rademacher complexity

- How about **regression**? Can derive similar results for nonparametric regression.

What We Learned...

- Formalizing learning
- PAC learnability
- Occam's razor Theorem
- VC dimension and VC theorem
- VC theorem
- No Free-lunch theorem

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Questions?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder