

Learning From Data

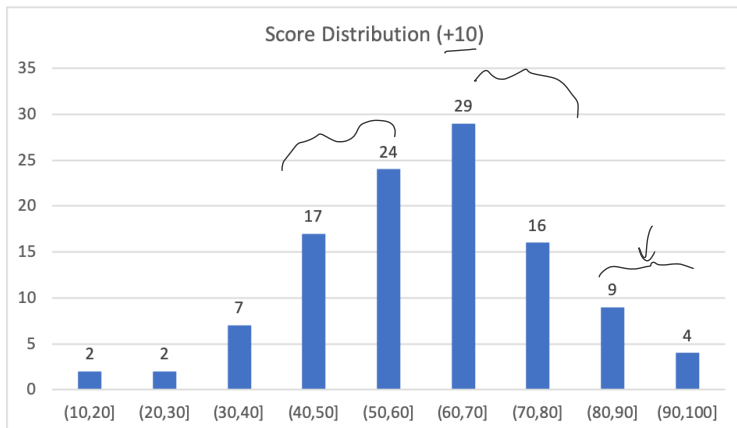
Lecture 7: Learning Theory

Yang Li yangli@sz.tsinghua.edu.cn

TBSI

November 18, 2022

Midterm Results

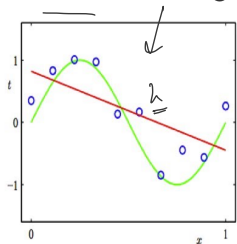


	max	mean	median
curved score	100	60	61

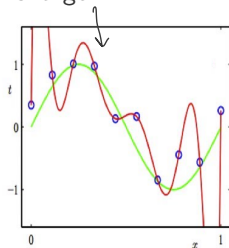
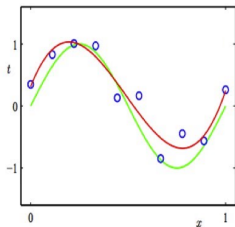
Review

Overfit & Underfit

- **Underfit** Both training error and testing error are large
- **Overfit** Training error is small, testing error is large



underfit

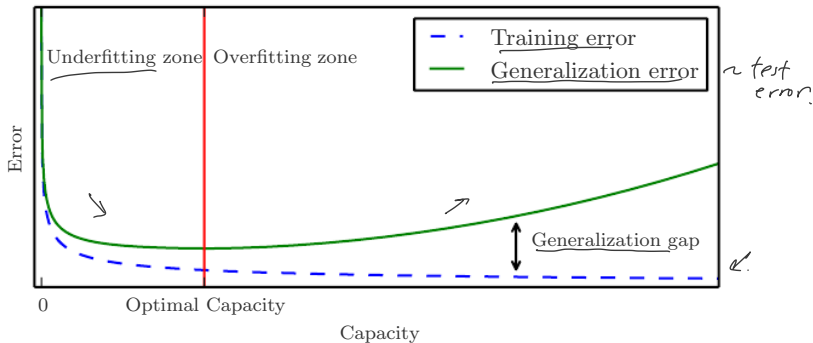


overfit

Model capacity: the ability to fit a wide variety of functions

Model Capacity

Changing a model's **capacity** controls whether it is more likely to overfit or underfit



How to formalize this idea?

Bias and Variance

Suppose data is generated by the following model:

$$\underline{y} = h(\underline{x}) + \underline{\epsilon} \quad \sim \underline{P_{XY}}$$

with $\mathbb{E}[\epsilon] = 0$, $\text{Var}(\epsilon) = \sigma^2$

- ▶ $h(x)$: true hypothesis function, unknown
- ▶ $\hat{h}_D(x)$: estimated hypothesis function based on training data $\underline{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ sampled from $\underline{P_{XY}}$
- ▶ **Model bias:** $\text{Bias}(\hat{h}_D(x)) = \mathbb{E}_D[\hat{h}_D(x) - h(x)]$ *Expected estimation error of the model over all choices of training data D*
- ▶ **Model variance:** $\text{Var}(\hat{h}_D(x)) = \mathbb{E}_D[\hat{h}_D(x)^2] - \mathbb{E}_D[\hat{h}_D(x)]^2$ *Variance of the model over all choices of D*

Bias - Variance Tradeoff

If we measure generalization error by MSE $y = \underbrace{h(x)}_{\text{fit}} + \underbrace{\varepsilon}_{\text{error}} \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$

$\underbrace{MSE}_{\text{fit error}} = \mathbb{E}[(\hat{h}_D(x) - y)^2] = \underbrace{\text{Bias}(\hat{h}_D(x))^2}_{\text{fit error}} + \underbrace{\text{Var}(\hat{h}_D(x))}_{\text{fit error}} + \underbrace{\sigma^2}_{\text{fit error}},$

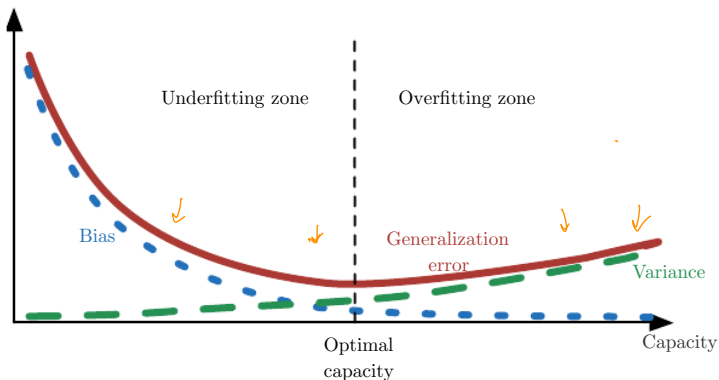
- ▶ σ^2 represents irreducible error (*caused by noisy data*)
- ▶ in practice, increasing capacity tends to increase variance and decrease bias.

Bias - Variance Tradeoff

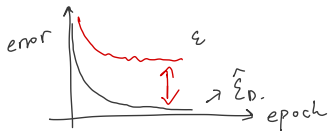
If we measure generalization error by MSE

$$MSE = \mathbb{E}[(\hat{h}_D(x) - y)^2] = \text{Bias}(\hat{h}_D(x))^2 + \text{Var}(\hat{h}_D(x)) + \sigma^2,$$

- ▶ σ^2 represents irreducible error (*caused by noisy data*)
- ▶ in practice, increasing capacity tends to increase variance and decrease bias.



Exercise:

overfitting

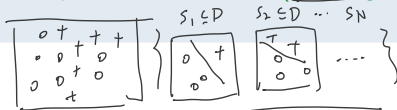
When the training error is much smaller than the testing error in a regression problem, what should be done? Select all that apply. }

- ▶ A) Add more training data. ✓
- ▶ B) Reduce model complexity. ✓
- ▶ C) Add more features. → reduce Bias ✗

- ▶ D) Apply random transformation to the training data (data augmentation). ✓ → more training data
→ equivalent to adding implicit regularization



Train multiple models on random subsets of the training data; Make prediction by averaging of the output of each model. (bagging a.k.a. bootstrap aggregation) ✓



noisy training data D

$$\text{Var}(h_i) = \sigma^2$$

$$\text{Var}\left(\sum_{i=1}^N h_i\right) = \frac{\sigma^2}{N} \text{ reduces variance}$$

Today's Lecture

- ▶ How to measure model capacity?
- ▶ Can we find a theoretical guarantee for model generalization?

A brief introduction to learning theory

- ▶ Empirical risk ^{minimization} estimation ←
- ▶ Generalization bound for finite and infinite hypothesis space
↙

Final project information.

Learning Theory

Empirical Risk Estimation

Uniform Convergence and Sample Complexity

Infinite H

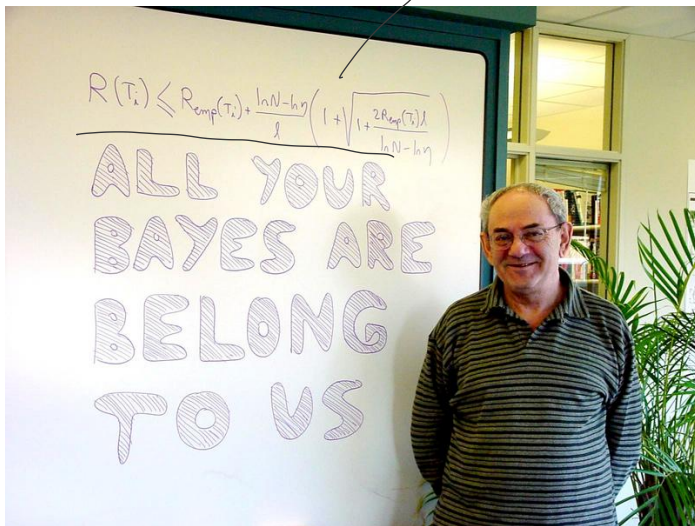
Introduction to Learning Theory

- ▶ Empirical risk estimation
- ▶ Learning bounds
 - ▶ Finite Hypothesis Class
 - ▶ Infinite Hypothesis Class

Learning theory

How to quantify generalization error?

uniform
convergence



Prof. Vladimir Vapnik in front of his famous theorem

Empirical risk

Simplified assumption: $y \in (0, 1)$

- ▶ Training set: $\underline{S} = (x^{(i)}, y^{(i)}); i = 1, \dots, m$ with $(x^{(i)}, y^{(i)}) \sim \underline{D}$
- ▶ For hypothesis h , the **training error** or **empirical risk/error** in learning theory is defined as

$$\checkmark \quad \hat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^m \underbrace{1\{h(x^{(i)}) \neq y^{(i)}\}}_{\text{0-1 error}}$$

- ▶ The **generalization error** is

$$\underline{\epsilon}(h) = P_{(x,y) \sim \underline{D}} \underbrace{(h(x) \neq y)}_{\mathbb{P}_{x,y}}$$

- ▶ **PAC assumption**: assume that training data and test data (for evaluating generalization error) were drawn from the same distribution \underline{D}

Hypothesis Class and ERM

$$\mathcal{H} = \{h_{\theta}(x) = 1\{\theta^T x \geq 0\} \mid \theta \in \mathbb{R}^n\}$$

Hypothesis class

The **hypothesis class** \mathcal{H} used by a learning algorithm is the set of all classifiers considered by it.

e.g. Linear classification considers $h_{\theta}(x) = 1\{\theta^T x \geq 0\} = \begin{cases} 1 & \theta^T x \geq 0. \\ 0 & \text{o.w.} \end{cases}$

Empirical Risk Minimization (ERM): the "simplest" learning algorithm: pick the best hypothesis \hat{h} from hypothesis class \mathcal{H}

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \underbrace{\hat{\epsilon}(h)}_{\text{empirical risk}} \quad \hat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^m 1\{L(x_i) \neq y_i\}$$

How to measure the generalization error of empirical risk minimization over \mathcal{H} ?

- ▶ Case of finite \mathcal{H}
- ▶ Case of infinite \mathcal{H}

Case of Finite \mathcal{H}

Goal: give guarantee on generalization error $\underline{\epsilon}(h) = \mathbb{E}_{x,y \sim \mathcal{D}} \mathbb{1}\{h(x) \neq y\}$

- ▶ Show $\hat{\underline{\epsilon}}(h)$ (training error) is a good estimate of $\underline{\epsilon}(h)$
- ▶ Derive an upper bound on $\underline{\epsilon}(h)$

For any $\underline{h}_i \in \underline{\mathcal{H}}$, the event of \underline{h}_i miss-classification given sample $(x, y) \sim \underline{\mathcal{D}}$:

$$\underline{Z} = \mathbb{1}\{\underline{h}_i(x) \neq y\} \quad \underline{Z} \in \{0, 1\}.$$

$\underline{Z}_j = \mathbb{1}\{\underline{h}_i(x^{(j)}) \neq y^{(j)}\}$: event of \underline{h}_i miss-classifying sample $x^{(j)}$

Case of Finite \mathcal{H}

Goal: give guarantee on generalization error $\epsilon(h)$

- ▶ Show $\hat{\epsilon}(h)$ (training error) is a good estimate of $\epsilon(h)$
- ▶ Derive an upper bound on $\epsilon(h)$

For any $h_i \in \mathcal{H}$, the event of h_i miss-classification given sample $(x, y) \sim \mathcal{D}$:

$$Z = 1\{h_i(x) \neq y\}$$

$Z_j = 1\{h_i(x^{(j)}) \neq y^{(j)}\}$: event of h_i miss-classifying sample $x^{(j)}$

Training error of $h_i \in \mathcal{H}$ is:

$$\hat{\epsilon}(h_i) = \frac{1}{m} \sum_{j=1}^m \underbrace{1\{h_i(x^{(j)}) \neq y^{(j)}\}}_{Z_j}$$

$$\hat{\epsilon}(h_i) = \frac{1}{m} \sum_{j=1}^m Z_j$$

Preliminaries

Here we make use of two famous inequalities:

Lemma 1 (Union Bound)

Let A_1, A_2, \dots, A_k be k different events, then

$$P(A_1 \cup \dots \cup A_k) \leq P(A_1) + \dots + P(A_k)$$

Probability of any one of k events happening is less the sums of their probabilities.

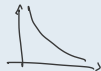
Preliminaries

Lemma 2 (Hoeffding Inequality, Chernoff bound)

Let Z_1, \dots, Z_m be m i.i.d. random variables drawn from a Bernoulli(ϕ) distribution. i.e. $P(Z_i = 1) = \phi$, $P(Z_i = 0) = 1 - \phi$. Let $\hat{\phi} = \frac{1}{m} \sum_{i=1}^m Z_i$ be the sample mean of RVs.

For any $\gamma > 0$,

$$P(|\phi - \hat{\phi}| > \gamma) \leq 2 \exp(-2\gamma^2 m)$$



The probability of $\hat{\phi}$ having large estimation error is small when m is large!

$\underbrace{\hat{\phi}}_{\text{sample mean}}$

Case of Finite \mathcal{H}

Training error of $h_i \in \mathcal{H}$ is:

$$\hat{\epsilon}(h_i) = \frac{1}{m} \sum_{j=1}^m Z_j$$

where $Z_j \sim \text{Bernoulli}(\epsilon(h_i))$

Case of Finite \mathcal{H}

Training error of $\underline{h}_i \in \mathcal{H}$ is:

$$\hat{\epsilon}(h_i) = \frac{1}{m} \sum_{j=1}^m Z_j$$

Hoeffding.

Given γ ,

$$P(|\hat{\phi} - \phi| > \gamma) \leq 2e^{-2\gamma^2 m}$$

where $Z_j \sim \text{Bernoulli}(\epsilon(h_i))$ Given γ .

By Hoeffding inequality, P for every $h_i \in \mathcal{H}$

$$P(|\epsilon(h_i) - \hat{\epsilon}(h_i)| > \gamma) \leq 2e^{-2\gamma^2 m} \quad (1)$$

Let A_i be the event (R.V) that $|\epsilon(h_i) - \hat{\epsilon}(h_i)| > \gamma$, $i=1, \dots, k$

Then $P_r(\exists h \in \mathcal{H}) |\epsilon(h) - \hat{\epsilon}(h)| > \gamma = P(A_1 \cup A_2 \dots \cup A_k) \leftarrow$

there exists $\xrightarrow{\text{By the Union bound,}} \leq \sum_{i=1}^k P(A_i) = \sum_{i=1}^k P(|\epsilon(h_i) - \hat{\epsilon}(h_i)| > \gamma)$.

Then, by negation

By (1),

$$\leq \sum_{i=1}^k 2e^{-2\gamma^2 m} = 2ke^{-2\gamma^2 m}$$

$$P_r(\forall h \in \mathcal{H}) |\epsilon(h) - \hat{\epsilon}(h)| \leq \gamma \geq 1 - 2ke^{-2\gamma^2 m}$$

Case of Finite \mathcal{H}

Training error of $\underline{h}_i \in \mathcal{H}$ is:

$$\hat{\epsilon}(h_i) = \frac{1}{m} \sum_{j=1}^m Z_j$$

where $Z_j \sim \text{Bernoulli}(\epsilon(h_i))$

By Hoeffding inequality,

$$P(\underbrace{|\epsilon(h_i) - \hat{\epsilon}(h_i)|}_{A_i} > \gamma) \leq 2e^{-2\gamma^2 m}$$

By Union bound,

$$P(\forall h \in \mathcal{H}. \underbrace{|\epsilon(h) - \hat{\epsilon}(h)|}_{\leq \gamma}) \geq \underbrace{1 - 2ke^{-2\gamma^2 m}}_{\delta.}$$

Uniform Convergence Results

Proposition.

Given γ, m

$$P(\forall h \in \mathcal{H} \mid |\hat{\epsilon}(h) - \epsilon(h)| \leq \gamma) \geq 1 - 2k e^{-2\gamma^2 m}$$

Corollary 3

Given γ and $\delta > 0$, If

$$m \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta} \quad \text{where } k = |\mathcal{H}|$$

Then with probability at least $1 - \delta$, we have $|\hat{\epsilon}(h) - \epsilon(h)| \leq \gamma$ for all h .
 m is called the algorithm's **sample complexity**.

$$\text{Let } \delta = 2k e^{-2\gamma^2 m}$$

$$\log \delta = \log(2k) + (-2\gamma^2 m)$$

$$m = \frac{\log \delta - \log 2k}{-2\gamma^2} = \frac{\log 2k - \log \delta}{2\gamma^2} = \frac{1}{2\gamma^2} \log \left(\frac{2k}{\delta} \right) \leftarrow \begin{array}{l} \text{minimum} \\ \text{sample size} \\ \text{for} \end{array}$$

$$P(\forall h \in \mathcal{H} \mid |\hat{\epsilon}(h) - \epsilon(h)| \leq \gamma) \geq 1 - \delta.$$

Uniform Convergence Results

Corollary 3

Given γ and $\delta > 0$, If

$$\underline{m} \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta}$$

Then with probability at least $1 - \delta$, we have $|\epsilon(h) - \hat{\epsilon}(h)| \leq \gamma$ for all h .
 m is called the algorithm's **sample complexity**.

Remarks

- ▶ Lower bound on \underline{m} tell us how many training examples we need to make generalization guarantee.
- ▶ # of training examples needed is logarithm in k

Uniform Convergence Results

Proposition.

Given γ, m

$$P(\forall h \in \mathcal{H} \mid |\hat{\epsilon}(h) - \epsilon(h)| \leq \gamma) \geq 1 - 2ke^{-2\gamma^2 m}$$

} γ
 δ
 m

Corollary 3. For $|\hat{\epsilon}(h) - \epsilon(h)| \leq \gamma$ to hold for all $h \in \mathcal{H}$, $m \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta}$.

Corollary 4

With probability $1 - \delta$, for all $h \in \mathcal{H}$, given m samples

$$|\hat{\epsilon}(h) - \epsilon(h)| \leq \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}}$$

By corollary 3, solve for γ : $2\gamma^2 = \frac{1}{m} \log \frac{2k}{\delta}$.

$$\gamma = \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}}$$

Then $|\hat{\epsilon}(h) - \epsilon(h)| \leq \gamma = \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}}$ for all $h \in \mathcal{H}$.

Uniform Convergence Results

Corollary 4

With probability $1 - \delta$, for all $h \in \mathcal{H}$,

$$|\hat{\epsilon}(h) - \epsilon(h)| \leq \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}}$$

What is the convergence result when we pick $\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{\epsilon}(h)$

Uniform Convergence Theorem for Finite \mathcal{H}

Using previous corollaries, we can bound $\epsilon(\hat{h})$:

Theorem 5 (Uniform convergence)

Let $|\mathcal{H}| = k$, and m, δ be fixed. With probability at least $1 - \delta$, we have

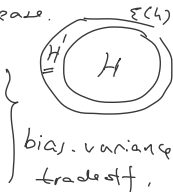
$$\epsilon(\hat{h}) \leq \left(\min_{h \in \mathcal{H}'} \epsilon(h) \right) + 2 \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}} \quad \leftarrow \text{variance}$$

1), Choose a larger $\mathcal{H}' \supseteq \mathcal{H}$, $\min_{h \in \mathcal{H}} \epsilon(h)$ will decrease.

i.e. $\min_{h \in \mathcal{H}'} \epsilon(h) \leq \min_{h \in \mathcal{H}} \epsilon(h) \rightarrow$ smaller bias.

2) When \mathcal{H} is larger, $K = |\mathcal{H}|$ increases.

then $2 \sqrt{\frac{1}{2m} \log \frac{K}{\delta}}$ increases



training risk/error $\hat{\mathcal{E}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{h(x_i) \neq y_i\}$,

testing/generalization risk: $\mathcal{E}(h) = \mathbb{E}_{x \sim \mathcal{D}} \mathbb{1}\{h(x) \neq y\}$.

$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{\mathcal{E}}(h) \leftarrow$ empirical estimator (using ERM)

$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathcal{E}(h) \leftarrow$ true hypothesis

With probability at least $1 - \delta$, sample size $m \geq \frac{1}{2\gamma^2} \log \frac{2^k}{\delta}$, we have

$$|\mathcal{E}(h) - \hat{\mathcal{E}}(h)| \leq \gamma \text{ for all } h \in \mathcal{H}. \quad (\text{by corollary 3})$$

Then $|\mathcal{E}(\hat{h}) - \hat{\mathcal{E}}(\hat{h})| \leq \gamma$ since $\hat{h} \in \mathcal{H}$

$$\mathcal{E}(\hat{h}) - \hat{\mathcal{E}}(\hat{h}) \leq \gamma.$$

$$\mathcal{E}(\hat{h}) \leq \gamma + \hat{\mathcal{E}}(\hat{h}) \leq \gamma + \hat{\mathcal{E}}(h^*). \quad (1)$$

$$\text{since } \hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{\mathcal{E}}(h), \hat{\mathcal{E}}(\hat{h}) \leq \hat{\mathcal{E}}(h^*)$$

Similarly, $|\mathcal{E}(h^*) - \hat{\mathcal{E}}(h^*)| \leq \gamma$

$$\mathcal{E}(h^*) - \hat{\mathcal{E}}(h^*) \geq -\gamma.$$

$$\hat{\mathcal{E}}(h^*) \leq \mathcal{E}(h^*) + \gamma,$$

By (1), $\mathcal{E}(\hat{h}) \leq \gamma + \hat{\mathcal{E}}(h^*) \leq \gamma + \mathcal{E}(h^*) + \gamma.$

$$\leq \mathcal{E}(h^*) + 2\gamma \leq \mathcal{E}(h^*) + 2\sqrt{\frac{1}{2m} \log \frac{2^k}{\delta}}.$$

$$\underbrace{\min_{h \in \mathcal{H}} \mathcal{E}(h)}$$

$$\mathcal{E}(\hat{h}) \leq \mathcal{E}(h^*) + 2\sqrt{\frac{1}{2m} \log \frac{2^k}{\delta}}.$$

Infinite hypothesis class: Challenges

Can we apply the same theorem to infinite \mathcal{H} ?

Example

- ▶ Suppose \mathcal{H} is parameterized by d real numbers. e.g.
 $\theta = [\theta_1, \theta_2, \dots, \theta_d] \in \mathbb{R}^d$ in linear regression with $d - 1$ unknowns.

Infinite hypothesis class: Challenges

Can we apply the same theorem to infinite \mathcal{H} ?

Example

$$\theta_d \in \mathbb{R}.$$

- ▶ Suppose \mathcal{H} is parameterized by d real numbers. e.g.
 $\theta = [\theta_1, \theta_2, \dots, \theta_d] \in \mathbb{R}^d$ in linear regression with $d - 1$ unknowns.
- ▶ In a 64-bit floating point representation, size of hypothesis class:

$$|\mathcal{H}| = 2^{64d} \leftarrow \begin{array}{l} \# \text{ of parameter.} \\ \uparrow \\ \text{one real \#} \end{array}$$

Infinite hypothesis class: Challenges

Can we apply the same theorem to infinite \mathcal{H} ?

Example

- ▶ Suppose \mathcal{H} is parameterized by d real numbers. e.g.
 $\theta = [\theta_1, \theta_2, \dots, \theta_d] \in \mathbb{R}^d$ in linear regression with $d - 1$ unknowns.
- ▶ In a 64-bit floating point representation, size of hypothesis class:
 $|\mathcal{H}| = 2^{64d}$
- ▶ How many samples do we need to guarantee $\epsilon(\hat{h}) \leq \epsilon(h^*) + 2\gamma$ to hold with probability at least $1 - \delta$?

$$\underline{m} \geq O\left(\frac{1}{\gamma^2} \log \frac{2^{64d}}{\delta}\right) = O\left(\frac{d}{\gamma^2} \log \frac{1}{\delta}\right) = \underline{O_{\gamma, \delta}(d)}$$

Infinite hypothesis class: Challenges

Can we apply the same theorem to infinite \mathcal{H} ?

Example

- ▶ Suppose \mathcal{H} is parameterized by d real numbers. e.g.
 $\theta = [\theta_1, \theta_2, \dots, \theta_d] \in \mathbb{R}^d$ in linear regression with $d - 1$ unknowns.
- ▶ In a 64-bit floating point representation, size of hypothesis class:
 $|\mathcal{H}| = 2^{64d}$
- ▶ How many samples do we need to guarantee $\epsilon(\hat{h}) \leq \epsilon(h^*) + 2\gamma$ to hold with probability at least $1 - \delta$?

$$m \geq O\left(\frac{1}{\gamma^2} \log \frac{2^{64d}}{\delta}\right) = O\left(\frac{d}{\gamma^2} \log \frac{1}{\delta}\right) = O_{\gamma, \delta}(d)$$

Infinite hypothesis class: Challenges

Can we apply the same theorem to infinite \mathcal{H} ?

Example

- ▶ Suppose \mathcal{H} is parameterized by d real numbers. e.g. $\theta = [\theta_1, \theta_2, \dots, \theta_d] \in \mathbb{R}^d$ in linear regression with $d - 1$ unknowns.
- ▶ In a 64-bit floating point representation, size of hypothesis class: $|\mathcal{H}| = 2^{64d}$
- ▶ How many samples do we need to guarantee $\epsilon(\hat{h}) \leq \epsilon(h^*) + 2\gamma$ to hold with probability at least $1 - \delta$?

$$m \geq O\left(\frac{1}{\gamma^2} \log \frac{2^{64d}}{\delta}\right) = O\left(\frac{d}{\gamma^2} \log \frac{1}{\delta}\right) = O_{\gamma, \delta}(d)$$

To learn well, the number of samples has to be linear in d

Infinite hypothesis class: Challenges

Size of \mathcal{H} depends on the choice of parameterization

Example

$2n + 2$ parameters:

$$h_{u,v} = \mathbf{1}\{ \overbrace{(u_0^2 - v_0^2)}^{\theta_0} + \overbrace{(u_1^2 - v_1^2)}^{\theta_1} x_1 + \dots + \overbrace{(u_n^2 - v_n^2)}^{\theta_n} x_n \geq 0 \}$$

is equivalent the hypothesis with $n + 1$ parameters:

$$h_{\theta}(x) = \mathbf{1}\{\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0\}$$

Infinite hypothesis class: Challenges

Size of \mathcal{H} depends on the choice of parameterization

Example

$2n + 2$ parameters:

$$h_{u,v} = \mathbf{1}\{(u_0^2 - v_0^2) + (u_1^2 - v_1^2)x_1 + \dots + (u_n^2 - v_n^2)x_n \geq 0\}$$

is equivalent the hypothesis with $n + 1$ parameters:

$$h_\theta(x) = \mathbf{1}\{\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0\}$$

We need a complexity measure of a hypothesis class invariant to parameterization choice

Infinite hypothesis class: Vapnik-Chervonenkis theory

A computational learning theory developed during 1960-1990 explaining the learning process from a statistical point of view.



Alexey Chervonenkis (1938-2014), Russian mathematician



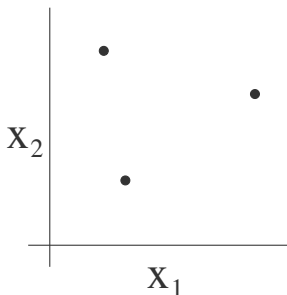
Vladimir Vapnik (Facebook AI Research, Vencore Labs)
Most known for his contribution in statistical learning theory

Shattering a point set

- Given d points $x^{(i)} \in \mathcal{X}$, $i = 1, \dots, d$, \mathcal{H} shatters S if \mathcal{H} can realize any labeling on S .

① \mathcal{H} realize a labeling on S :

Example: $S = \{x^{(1)}, x^{(2)}, x^{(3)}\}$ where $x^{(i)} \in \mathbb{R}^2$.



$\exists h \in \mathcal{H}$, $h(x)$ makes no mistakes in predicting the labeling

② \mathcal{H} can realize any labeling on S .

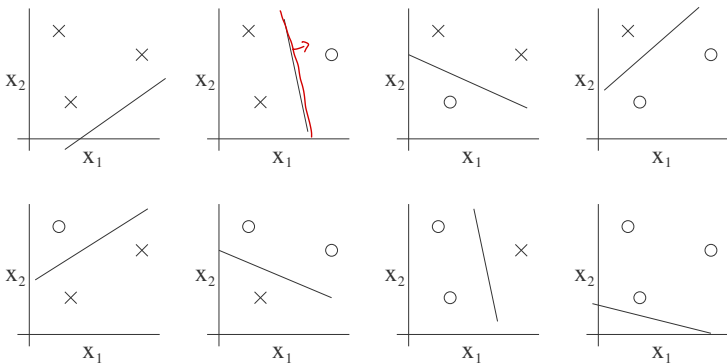
Suppose $y^{(i)} \in \{0, 1\}$, how many possible labelings does S have? 2^3

Shattering a point set

- ▶ Example: Let $\mathcal{H}_{LTF,2}$ ^{linear threshold function in \mathbb{R}^2 .} be the linear threshold function in \mathbb{R}^2 (e.g. in the perceptron algorithm)

$$h(x) = \begin{cases} 1 & \underline{w_1 x_1 + w_2 x_2} \geq b \\ 0 & \text{otherwise} \end{cases}$$

$\mathcal{H}_{LTF,2}$ can shatter S
 $|S| = 3$.



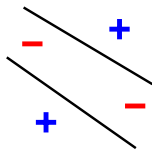
$\mathcal{H}_{LTF,2}$ shatters $S = \{x^{(1)}, x^{(2)}, x^{(3)}\}$

VC Dimension

The **Vapnik-Chervonenkis** dimension of \mathcal{H} , or $VC(\mathcal{H})$, is the cardinality of the largest set shattered by \mathcal{H} .

▶ Example: $VC(H_{LTF,2}) = 3$

$$VC(H_{LTF,2}) \geq 3,$$



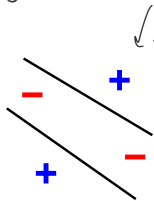
$$\rightarrow VC(H_{LTF,2}) < 4.$$

\mathcal{H}_{LTF} can not shatter 4 points: for any 4 points, label points on the diagonal as '+'. (See Radon's theorem)

VC Dimension

The **Vapnik-Chervonenkis** dimension of \mathcal{H} , or $VC(\mathcal{H})$, is the cardinality of the largest set shattered by \mathcal{H} .

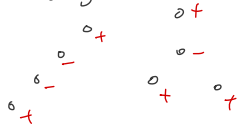
▶ Example: $VC(\mathcal{H}_{LTF,2}) = 3$



To show $VC(\mathcal{H}_{LTF,2}) < 4$,

Given any set S , with $|S| = 4$.

Construct a labeling by assigning opposite labels on the diagonal.



\mathcal{H}_{LTF} can not shatter 4 points: for any 4 points, label points on the diagonal as '+'. (See Radon's theorem)

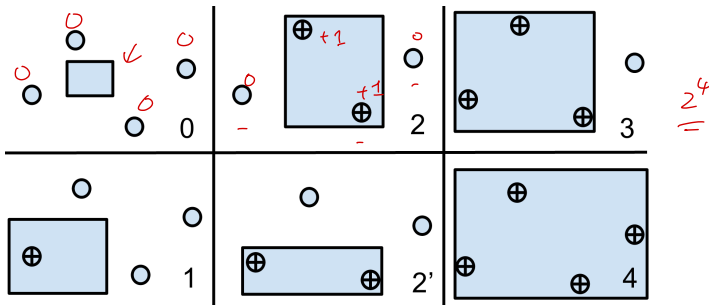
- ▶ To show $VC(\mathcal{H}) \geq d$, it's sufficient to find one set of d points shattered by \mathcal{H} ✓
- ▶ To show $VC(\mathcal{H}) < d$, need to prove \mathcal{H} doesn't shatter any set of d points ✓

VC Dimension

$$h_{\text{AAR}} = \begin{cases} 1. & x \in \text{Rectangle.} \\ 0. & \text{otherwise} \end{cases}$$

► Example: $VC(\text{AxisAlignedRectangles}) = 4$

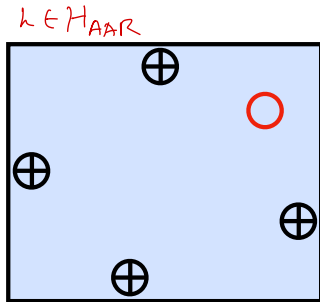
$$|S| = 4.$$



Axis-aligned rectangles can shatter 4 points. $VC(\text{AxisAlignedRectangles}) \geq 4$

VC Dimension

- ▶ Example: $VC(\text{AxisAlignedRectangles}) = 4$



← For all $|S| = 5$,
we can find such
labeling that H_{AAR}
can not realize!

For any 5 points, label topmost, bottommost, leftmost and rightmost points as “+”.

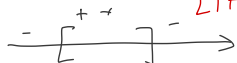
$VC(\text{AxisAlignedRectangles}) < 5$

Discussion on VC Dimension

More VC results of common \mathcal{H} :

- ▶ $VC(\text{Positive Half-Lines}) = 1, \mathcal{X} = \mathbb{R}$

LTF in \mathbb{R} .



- ▶ $VC(\text{Intervals}) = 2, \mathcal{X} = \mathbb{R}$

- ▶ $VC(\text{LTF in } \mathbb{R}^n) = n + 1, \mathcal{X} = \mathbb{R}^n \leftarrow \text{prove this at home!}$

$$\textcircled{1} VC(\text{PHL}) \leq 1.$$



$$\textcircled{2} VC(\text{PHL}) < 2.$$

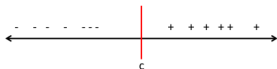


Let $c = \frac{x_2 - x_1}{2}$ assuming $x_2 > x_1$.
 label $\min(x_1, x_2)$ as +1.
 and the other as 0.

Discussion on VC Dimension

More VC results of common \mathcal{H} :

- ▶ $VC(\text{PositiveHalf-Lines}) = 1, \mathcal{X} = \mathbb{R}$



- ▶ $VC(\text{Intervals}) = 2, \mathcal{X} = \mathbb{R}$
- ▶ $VC(\text{LTF in } \mathbb{R}^n) = n + 1, \mathcal{X} = \mathbb{R}^n \leftarrow \text{prove this at home!}$

Proposition 1

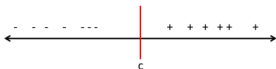
If \mathcal{H} is finite, VC dimension is related to the cardinality of \mathcal{H} :

$$VC(\mathcal{H}) \leq \log |\mathcal{H}|$$

Discussion on VC Dimension

More VC results of common \mathcal{H} :

- ▶ $VC(\text{PositiveHalf-Lines}) = 1, \mathcal{X} = \mathbb{R}$



- ▶ $VC(\text{Intervals}) = 2, \mathcal{X} = \mathbb{R}$
- ▶ $VC(\text{LTF in } \mathbb{R}^n) = n + 1, \mathcal{X} = \mathbb{R}^n \leftarrow \text{prove this at home!}$

Proposition 1

If \mathcal{H} is finite, VC dimension is related to the cardinality of \mathcal{H} :

$$\underline{VC(\mathcal{H})} \leq \log \underline{|\mathcal{H}|}$$

Let $VC(\mathcal{H}) = d$.
 \mathcal{H} has to shatter d points
 $\geq 2^d$ labelings. $|\mathcal{H}| \geq 2^d$.

Proof. Let $d = VC|\mathcal{H}|$. There must exist a shattered set of size d on which \mathcal{H} realizes all possible labelings. Every labeling must have a corresponding hypothesis, then $|\mathcal{H}| \geq 2^d$



Learning bound for infinite \mathcal{H}

Theorem 6

Given \mathcal{H} , let $\underline{d} = \underline{VC}(\mathcal{H})$.

- With probability at least $1 - \delta$, we have that for all h generation $n \geq \frac{1}{\epsilon}$.

replaced $k = |\mathcal{H}|$
by $VC(\mathcal{H}) \leq \log |\mathcal{H}|$

$$|\epsilon(h) - \hat{\epsilon}(h)| \leq O \left(\sqrt{\frac{\overbrace{d}^{VC(\mathcal{H})}}{m} \log \frac{m}{d} + \frac{1}{m} \log \frac{1}{\delta}} \right)$$

\downarrow training error \uparrow sample size

Learning bound for infinite \mathcal{H}

Theorem 6

Given \mathcal{H} , let $d = VC(\mathcal{H})$.

- ▶ With probability at least $1 - \delta$, we have that for all h

$$|\epsilon(h) - \hat{\epsilon}(h)| \leq O\left(\sqrt{\frac{d}{m} \log \frac{m}{d} + \frac{1}{m} \log \frac{1}{\delta}}\right)$$

- ▶ Thus, with probability at least $1 - \delta$, we also have

$$\epsilon(\hat{h}) \leq \epsilon(h^*) + O\left(\sqrt{\frac{d}{m} \log \frac{m}{d} + \frac{1}{m} \log \frac{1}{\delta}}\right)$$

$h = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathcal{E}(h)$

testing error \sim learned from ERM

variance

Learning bound for infinite \mathcal{H}

Corollary 7

For $|\epsilon(h) - \hat{\epsilon}(h)| \leq \gamma$ to hold for all $h \in \mathcal{H}$ with probability at least $1 - \delta$, it suffices that $\underline{m} = \underline{O}_{\gamma, \delta}(d)$.

Learning bound for infinite \mathcal{H}

Corollary 7

For $|\epsilon(h) - \hat{\epsilon}(h)| \leq \gamma$ to hold for all $h \in \mathcal{H}$ with probability at least $1 - \delta$, it suffices that $m = O_{\gamma, \delta}(d)$.

Remarks

- ▶ Sample complexity using \mathcal{H} is linear in $\underline{VC}(\mathcal{H})$
- ▶ For “most”^a hypothesis classes, the VC dimension is linear in terms of parameters
- ▶ For algorithms minimizing training error, # training examples needed is roughly linear in number of parameters in \mathcal{H} .

^aNot always true for deep neural networks

VC Dimension of Deep Neural Networks

Theorem 8 (Cover, 1968; Baum and Haussler, 1989)

Let \mathcal{N} be an arbitrary feedforward neural net with w weights that consists of linear threshold activations, then $VC(\mathcal{N}) = O(w \log w)$.

VC Dimension of Deep Neural Networks

Theorem 8 (Cover, 1968; Baum and Haussler, 1989)

Let \mathcal{N} be an arbitrary feedforward neural net with w weights that consists of linear threshold activations, then $VC(\mathcal{N}) = O(w \log w)$.

Recent progress

- ▶ For feed-forward neural networks with piecewise-linear activation functions (e.g. ReLU), let w be the number of parameters and l be the number of layers, $VC(\mathcal{N}) = O(w/l \log(w))$ [Bartlett et. al., 2017]

Bartlett and W. Maass (2003) Vapnik-Chervonenkis Dimension of Neural Nets

Bartlett et. al., (2017) Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks.

VC Dimension of Deep Neural Networks

Theorem 8 (Cover, 1968; Baum and Haussler, 1989)

Let \mathcal{N} be an arbitrary feedforward neural net with w weights that consists of linear threshold activations, then $VC(\mathcal{N}) = O(w \log w)$.

Recent progress

- ▶ For feed-forward neural networks with piecewise-linear activation functions (e.g. ReLU), let w be the number of parameters and l be the number of layers, $VC(\mathcal{N}) = O(wl \log(w))$ [Bartlett et. al., 2017]
- ▶ *Among all networks with the same size (number of weights), more layers have larger VC dimension*, thus more training samples are needed to learn a deeper network

Bartlett and W. Maass (2003) Vapnik-Chervonenkis Dimension of Neural Nets

Bartlett et. al., (2017) Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks.

Final Project Information

See <http://yangli-feasibility.com/home/classes/lfd2022fall/project.html>