

Introduction

phenomenon: state-of-art models are over-parameterized.

data interpolation, memorization

shows in

contradiction: classical theory of generalization vs. modern ML practice
extracting patterns
memorize everything needed.
trade-offs between
performs good on training sets
less error on training sets
with optimal generalization errors
and
generalization abilities.

explanation: the insufficient data hurdles the accuracy of model.
long-tailed datasets.

The aggregate of small subsets outweigh the larger ones.

"Fragmented" dataset.

For small subsets, it's better to just remember them.

a little bit like "rules" and "exceptions" in life.

question: if memorizing these "exceptions" is a must to reach close-to-optimal generalization error?

- ① the algorithm can't decide the borderline between memorizing and not-memorizing
- ② the long-tailed dataset → even the smallest subsets can count when aggregating together

∴ ①, ②

∴ Need to memorize

results.

$$\underline{\text{err}}(\pi, F, A) \geq \underline{\text{opt}}(\pi, F) + \tau_i \cdot \underline{E[\text{err}_{ns}(A, I)]}$$

expected generalization error $\min_{\alpha} \text{achievable error}$ num of examples appear only once and mislabeled.
↓
mistakes on "memory"

$$\tau_i = \frac{E_{\alpha \sim \bar{\pi}^N} [\alpha^2 \cdot (1-\alpha)^{n-1}]}{E_{\alpha \sim \bar{\pi}^N} [\alpha \cdot (1-\alpha)^{n-1}]}$$

Memorization. in Unstructured Classification

$$[n] = \{1, \dots, n\} \quad n \in \{1, 2, 3, \dots\}$$

$I(O)$: under the condition O , whether the event will happen or not can be expressed by '0' for happen and '1' for not happen.

dataset $S = (x_1, y_1) \dots (x_n, y_n)$ x_s = all points in S .

Probability distribution of $F(x)$:

value of $F(x)$ F is a function on X .
 $D_{x \sim D}[F(x) | x \in O]$ under the condition O
 $x \in X$ for all x subjective to distribution D .
for all the x in X which subjective to distribution D and under condition O .
the distribution of $F(x)$, for F is a function on X .

$TV(D_1, D_2)$: variation distance between D_1, D_2

function $h: X \rightarrow Y$ prediction function

dataset $S: (x_1, y_1), \dots, (x_n, y_n)$ x for points y for labels.

$\text{err}_P(h) := E_{(x,y) \sim P} [h(x) \neq y]$ the expectation of the prediction mistakes of given x .

For an algorithm A , the expected generalization error:

$\text{err}_P(A, S) := \underbrace{E_{h \sim A(S)}}_{\text{h being the output of the algorithm.}} [\text{err}_P(h)]$ the expectation of $\text{err}_P(h)$

For an algorithm A on any sample S in P , the expected generalization error.

$\text{err}_P(A) := E_{S \sim P^n} [\text{err}_P(A, S)]$ the expectation of $\text{err}_P(A)$

$\text{err}_P(h)$ is the indicator of prediction action.

$E \downarrow$ E for all predict. action. for one $h(x)$

$\text{err}_P(A, S)$ is the indicator of the algorithm.

$E \downarrow$ E for all the possible $h(x)$

$\text{err}_P(A)$ is the indicator of all of the data points.

$E \downarrow$ E for all the possible samples S on P .

Problem setup.

$$|X| = N \quad |Y| = m$$

data domain label domain.

Lemma 2.1. $E[D(x) | U = V] = \frac{\sum_{l=1}^n \pi_l l}{\sum_{l=1}^n \pi_l}$ (Each point in V) $\sim D$ $\rightarrow l: \text{in set } V, x \text{ shows } l \text{ times}$

$$D \sim D_\pi^X, U \sim D^n \quad E[D(x) | U = V] = \frac{\sum_{l=1}^n \pi_l l \cdot (1-\pi_l)^{n-l}}{\sum_{l=1}^n \pi_l l}$$

Distribution over $D(x)$ (x_1, x_2, \dots, x_n) the marginal distribution

in condition $U = V$, the expectation of $D(x)$ can be represented as ..

generalization error on a randomly chosen problem:

$$\overline{\text{err}}(\pi, F, A) := \underbrace{E_{D \sim D^\pi} [\text{err}_{D, f}(A)]}_{f \sim F}$$

for the whole data domain and all the prediction functions.

dataset $S \in (X \times Y)^n$, $l \in [n]$

$X_{S=l}$: set of points shows l times in S .

$h: X \rightarrow Y$

$$\text{err}_S(h, l) := |\{x \in X_{S=l} \mid h(x) \neq y\}|$$

the number of x for which $h(x) \neq y$ and x belongs to $X_{S=l}$

$$\text{err}_S(A, l) := E_{h \sim A(S)} [\text{err}_S(h, l)] \quad \text{for all of the prediction function } h \text{ for Algorithm A.}$$

For all the algorithms and datasets :

$$(D, f, S): D \sim D_\pi^X \quad f \sim F \quad S \sim (D, f)^n$$

↑
the data domain ↑
all the labeling
functions

G is the distribution of (D, f, S)
the datasets on domain D and
using labeling function f

$$\text{For } Z \in (X \times Y)^n \quad G(Z) := \underbrace{D}_{(D, f, S) \sim G} [(D, f) \mid S = Z]$$

$Z \in (X \times Y)^n$
the f only define the
labeling function
 S is still a dataset.

the expected error of A conditioned on dataset $= Z$

$$\overline{\text{err}}(\pi, F, A | Z) := \underbrace{E_{(D, f) \sim G(Z)} [\text{err}_{D, f}(A, Z)]}_{(D, f) \sim G(Z)}$$

for all the datasets and labeling functions.

$$\min[\overline{\text{err}}(\pi, F, A | Z)] = \overline{\text{err}}(\pi, F, A' | Z) = \text{opt}(\pi, F | Z)$$

for all the algorithm A s.

$$\text{Theorem 2.3. } \overline{\text{err}}(\pi, F, A | Z) \geq \underline{\text{opt}}(\pi, F | Z) + \sum_{l \in [n]} \tau_l \cdot \text{err}_{\pi}(A, l)$$

The not-fitting error of prior distribution π
labeling function $f \in F$

$$\min_{\text{not-fitting error}}$$

$$\sum_{l=1,2,3 \dots n} \tau_l$$

not-fitting error
For algorithm A and x showing l times

Algorithm A

dataset $Z = (D, f)$

All of the mistakes \geq min mistakes + number of all mistakes on points showing 1, 2, 3 ... n times.

$$\tau_l := \frac{E_{\alpha \sim \bar{\pi}^N} [\alpha^{l+1} \cdot (1-\alpha)^{n-l}]}{E_{\alpha \sim \bar{\pi}^N} [\alpha^l \cdot (1-\alpha)^{n-l}]}$$

The number of singleton is determined by "How long the tail is".

the fraction of examples in frequency range $[\beta_1, \beta_2]$:

$$\text{Weight}(\bar{\pi}^N, [\alpha, \beta]) := E_{D \sim D_{\bar{\pi}}^X} \left[\sum_{x \in X} D(x) \cdot I(D(x) \in [\beta_1, \beta_2]) \right]$$

↑ distribution {0,1} indicator
 ↓ in range → 1
 not in range → 0

Whether the data is counted depends on the data in the range or not.

$$= N \cdot E_{\alpha \sim \bar{\pi}^N} [\alpha \cdot I(\alpha \in [\beta_1, \beta_2])]$$

Get the data out of $D(x)$

divide the $D(x)$ into data (N) and distribution E .

The expected number of singleton points:

$$\text{single}(\bar{\pi}^N) := E_{D \sim P, V \sim D^n} [|X_{V=1}|] \text{ the expectation for the number of all the "single" points}$$

$$= E_{D \sim D_x^\pi} \left[\sum_{x \in X} \Pr_{V \sim D^n} [x \in X_{V=1}] \right]$$

the add results of all the prob of all of the single points.

$$= E_{D \sim D_x^\pi} \left[\sum_{x \in X} n \cdot \underbrace{D(x)}_{\text{distribution of } x} \underbrace{(1 - D(x))^{n-1}}_{\text{distribution of "not } x"} \right]$$

$$= \sum_{x \in X} n E_{D \sim D_x^\pi} [D(x) (1 - D(x))^{n-1}]$$

$$= n N \cdot E_{\alpha \sim \bar{\pi}^N} [\alpha (1 - \alpha)^{n-1}]$$

$$\sum_{x \in X} E[D(x)] = N \cdot E[\alpha]$$

for $\alpha < \frac{1}{3}$, $(1 - \alpha)^{n-1} \geq \frac{1}{3}$

$$\text{single}(\bar{\pi}^N) \geq n N \cdot E_{\alpha \sim \bar{\pi}^N} [\underbrace{\alpha (1 - \alpha)^{n-1}}_{\geq \frac{1}{3}} \cdot \underbrace{1(\alpha \in [0, \frac{1}{n}])}_{\text{add range for } E}]$$

$$\frac{n}{3} N E_{\alpha \sim \bar{\pi}^N} [\alpha \cdot 1(\alpha \in [0, \frac{1}{n}])]$$

$$\text{weight}(\bar{\pi}^N, [\alpha, \beta]) = N \cdot E_{\alpha \sim \bar{\pi}^N} [\alpha \cdot 1(\alpha \in [\beta_1, \beta_2])]$$

$$\therefore \frac{n}{3} N E_{\alpha \sim \bar{\pi}^N} [\alpha \cdot 1(\alpha \in [0, \frac{1}{n}])] = \frac{n}{3} \text{weight}(\bar{\pi}^N, [0, \frac{1}{n}])$$

$$\text{single}(\bar{\pi}^N) \geq \frac{n}{3} \text{weight}(\bar{\pi}^N, [0, \frac{1}{n}])$$

the cost of not fitting any of the singleton examples is at least
the weight of frequency of $\frac{1}{n}$

$$\text{Loostab}(P, A) := \frac{1}{n} \sum_{i=1, 2, 3, \dots, n} E_{S \sim P^n} \left[\left| \Pr_{h \sim A(S)} [h(x_i) = y_i] - \Pr_{h \sim A(S \setminus i)} [h(x_i) = y_i] \right| \right]$$

$i=1, 2, 3, \dots, n$. the probability of
"correct"
on dataset S .
using algorithm A

the probability of
"correct"
on dataset S . without point i
using algorithm A

$$\text{mem}(A, S, i) := \max \left\{ 0, \underbrace{\Pr_{h \sim A(S)} [h(x_i) = y_i] - \Pr_{h \sim A(S \setminus i)} [h(x_i) = y_i]}_{\begin{array}{l} \text{the probability of} \\ \text{"correct"} \\ \text{on dataset } S \\ \text{using algorithm } A \end{array}} \right\}$$

the probability of
 "correct"
 on dataset S .
 using algorithm A

$$\underbrace{\Pr_{h \sim A(S \setminus i)} [h(x_i) = y_i]}_{\begin{array}{l} \text{the probability of} \\ \text{"correct"} \\ \text{on dataset } S \text{ without point } i \\ \text{using algorithm } A \end{array}}$$

$$\min(\text{mem}(A, S, i)) = 0$$

mem evaluate the increase of prediction correctness

"How large can the correctness increase by 'knowing' point x_i ."

γ -memorization bounded: $x_i \in X_{S=1}$

for singletons

$$\text{mem}(A, S, i) \leq \gamma$$

"The maximum enhancement on correctness
by memorizing points in $X_{S=1}$."

$$\frac{1}{n} E_{S \sim P^n} \left[\sum_{i \in [n]} \text{mem}(A, S, i) \right] = \frac{1}{n} \sum_{i \in [n]} E_{S \sim P^n} \left[\max \left\{ 0, \underbrace{\Pr_{h \sim A(S)} [h(x_i) = y_i] - \Pr_{h \sim A(S \setminus i)} [h(x_i) = y_i]}_{\begin{array}{l} \text{the probability of} \\ \text{"correct"} \\ \text{on dataset } S \\ \text{using algorithm } A \end{array}} \right\} \right]$$

$$\text{Loostab}(P, A) := \frac{1}{n} \sum_{i \in [n]} E_{S \sim P^n} \left[\left| \Pr_{h \sim A(S)} [h(x_i) = y_i] - \Pr_{h \sim A(S \setminus i)} [h(x_i) = y_i] \right| \right]$$

$$e = \Pr_{h \sim A(S)} [h(x_i) = y_i] - \Pr_{h \sim A(S \setminus i)} [h(x_i) = y_i]$$

$$\max \{0, e\} \in [0, e]$$

$$|e| \in [-e, e]$$

$$\text{for } e \geq 0 \quad \max \{0, e\} = |e|$$

$$\text{for } e < 0 \quad \max \{0, e\} = 0. \quad |e| = -e \quad \Rightarrow \frac{1}{n} E_{S \sim P^n} \left[\sum_{i \in [n]} \text{mem}(A, S, i) \right] \leq \text{Loostab}(P, A)$$

$$0 < -e$$

Lemma 4.2.

$$E[\text{errn}_S(A, l)] \geq E \left[\sum_{\substack{i \in [n], \\ h \sim A(S \setminus \{x_i\})}} \Pr_{x_i \in X_{S \setminus \{i\}}} [h(x_i) \neq y_i] - \text{mem}(A, S, i) \right]$$

$f \sim F$ $f \sim F$ $\Pr_{\substack{i \in [n], \\ h \sim A(S \setminus \{x_i\})}} [h(x_i) \neq y_i]$ $\text{mem}(A, S, i)$
 $S \sim (D, f)^n$ $S \sim (D, f)^n$ The probability of error The correctness increasement
expectation of error Algorithm A for singeltons.
made on points only showing 1 time: on dataset S removing x_i
for Algorithm A.

$$\text{errn}_S(h, l) := |\{x \in X_{S \setminus \{l\}} \mid h(x) \neq y\}|$$

$$\text{errn}_S(A, l) := E_{h \sim A(S)} [\text{errn}_S(h, l)] =$$

$$E \left[\text{errn}_S(A, l) \right] = E_{f \sim F} \left[E_{h \sim A(S)} \left[|\{x \in X_{S \setminus \{l\}} \mid h(x) \neq y\}| \right] \right]$$

$f \sim F$ $E_{h \sim A(S)}$ $|\{x \in X_{S \setminus \{l\}} \mid h(x) \neq y\}|$
 $S \sim (D, f)^n$ E_f The number of error on singeltons for every algorithm on S .
every labeling function f and every dataset S .