What Neural Networks Memorize and Why

Discovering the Long Tail via Influence Estimation

Chiyuan Zhang, 2020.06.10 Joint work with Vitaly Feldman

Agenda

01 Motivation

02 Formal Definitions

03 Efficient Estimators

04 Case Studies

05 Theoretical Characterization

⁰⁶ Further Investigations

Motivation

Memorization is ubiquitous in deep learning

Memorization puzzles our intuition on generalization

Memorization is harmful in practice (security / privacy concerns)

Memorization is ubiquitous in deep learning





fig source: https://arxiv.org/abs/1611.03530

ImageNet Trainset: 1.2M

fig source: https://arxiv.org/abs/1905.11946

Memorization puzzles our intuition on generalization



fig source: <u>https://arxiv.org/abs/2002.03206</u>

Memorization is harmful in practice

For example, users may find that the input "my social-security number is..." gets auto-completed to an obvious secret (such as a valid-looking SSN not their own), or find that other inputs are auto-completed to text with oddly-specific details. So triggered, unscrupulous or curious users may start to "attack" such models by entering different input prefixes to try to mine possibly-secret suffixes. Therefore, for generative text models, assessing and reducing the chances that secrets may be disclosed in this manner is a key practical concern.

"Knowledge" vs "Facts"

fig source: https://arxiv.org/abs/1802.08232

Definition 2.4 (Differential Privacy). A randomized algorithm \mathcal{M} with domain $\mathbb{N}^{|\mathcal{X}|}$ is (ε, δ) -differentially private if for all $\mathcal{S} \subseteq \text{Range}(\mathcal{M})$ and for all $x, y \in \mathbb{N}^{|\mathcal{X}|}$ such that $||x - y||_1 \leq 1$:

 $\Pr[\mathcal{M}(x) \in \mathcal{S}] \le \exp(\varepsilon) \Pr[\mathcal{M}(y) \in \mathcal{S}] + \delta,$

fig source: https://www.nowpublishers.com/article/Details/TCS-042

02

A Formal Framework for Studying Memorization in Deep Learning

High level intuition why memorization could be useful

Formal definition of "memorization" and "useful"

Could a memorized example be useful at all?

а



Real world data follows power law

A plot of the rank versus frequency for the first 10 million words in 30 Wikipedias (dumps from October 2015) in a log-log scale.

For the second s

Categorization is a crude approximation for human's convenience, the hierarchy continues inside each category.

С

fig source: https://animalhybrids.files.wordpress.com/2012/05/dog-breeds.jpg

"Learning"

... a continuous spectrum ...

"Memorization"

A Formal Definition of Memorization

- A: a learning algorithm (including architecture, optimizer, hyper-parameters)
- S: a training set
- i: (the index of) a training example

$$\operatorname{mem}(\mathcal{A}, S, i) := \Pr_{h \leftarrow \mathcal{A}(S)}[h(x_i) = y_i] - \Pr_{h \sim \mathcal{A}(S \setminus i)}[h(x_i) = y_i]$$

If a model is able to correctly classify an example i with high probability only when this example is included in the training set, then this example is considered to be **memorized**.

Quantifying the Utility of a Memorized Example

- A: a learning algorithm (including architecture, optimizer, hyper-parameters)
- S: a training set
- i: (the index of) a training example
- j: (the index of) a test example

$$\operatorname{infl}(\mathcal{A}, S, i, j) := \Pr_{h \leftarrow \mathcal{A}(S)}[h(x'_j) = y'_j] - \Pr_{h \sim \mathcal{A}(S \setminus i)}[h(x'_j) = y'_j]$$

If a model is able to correctly classify an example j with high probability only when another example i is included in the training set, then i is considered **crucially important** for j.

Note $mem(\mathcal{A}, S, i) = infl(\mathcal{A}, S, i, (x_i, y_i))$

03

Efficient Estimators for Influence and Memorization

With the previous definitions, directly estimating the influence of n training examples within standard deviation σ requires training n/ σ^2 models (e.g. 5M for n=50,000, σ =0.1).

Subsampled Influence Estimation

the influence estimation in multiple examples.

$$\begin{split} \inf & \operatorname{Imfl}_m(\mathcal{A}, S, i, z) := \mathop{\mathbf{E}}_{I \sim P([n] \setminus \{i\}, m-1)} \left[\operatorname{imfl}(\mathcal{A}, S_{I \cup \{i\}}, i, z) \right] \\ & \text{If we sample many different subsets, for each example i,} \\ & \text{there are subsets containing i and subsets not containing i.} \\ & \text{Therefore, training on each subset sample can be reused in} \\ \end{split}$$

Sub-sampled training sets

Subsampled Influence Estimation

Algorithm 1 Memorization and influence value estimators

Require: Training dataset: $S = ((x_1, y_1), \dots, (x_n, y_n))$, testing dataset $S_{test} = ((x'_1, y'_1), \dots, (x'_{n'}, y'_{n'}))$, learning algorithm \mathcal{A} , subset size m, number of trials t.

- 1: Sample t random subsets of [n] of size $m: I_1, I_2, \ldots, I_t$.
- 2: **for** k = 1 to t **do**
- 3: Train model h_k by running \mathcal{A} on S_{T_k} .
- 4: **for** i = 1 to n **do**
- 5: $\widetilde{\mathrm{mem}}_m(\mathcal{A}, S, i) := \mathbf{Pr}_{k \sim [t]}[h_k(x_i) = y_i \mid i \in I_k] \mathbf{Pr}_{k \sim [t]}[h_k(x_i) = y_i \mid i \notin I_k].$
- 6: **for** j = 1 to n' **do**
- 7: $\widetilde{\mathsf{infl}}_m(\mathcal{A}, S, i, j) := \mathbf{Pr}_{k \sim [t]}[h_k(x'_j) = y'_j \mid i \in I_k] \mathbf{Pr}_{k \sim [t]}[h_k(x'_j) = y'_j \mid i \notin I_k].$
- 8: **return** $\widetilde{\text{mem}}_m(\mathcal{A}, S, i)$ for all $i \in [n]$; $\inf \mathbb{1}_m(\mathcal{A}, S, i, j)$ for all $i \in [n], j \in [n']$.

Lemma 2.1. There exists an algorithm that for every dataset $S \in (X \times Y)^n$, learning algorithm \mathcal{A} , $m \in [n]$ and integer t, runs \mathcal{A} t times and outputs estimates $(\mu_i)_{i \in [n]}$ such that for every $i \in [n]$ and $p = \min(m/n, 1 - m/n)$,

$$\mathbf{E}\left[(\mathtt{infl}_m(\mathcal{A}, S, i, z) - \mu_i)^2\right] \le \frac{1}{pt} + \frac{1}{(1-p)t} + \frac{e^{-pt/16}}{2},$$

where the expectation is with respect to the randomness of A and the randomness of the estimation algorithm.

Illustration of the Estimation Procedure



Random 70% subset

Training Pipeline

data-augmentation, regularization, stable initialization, SoTA activation function, fancy Ir schedule, momentum, preconditioning, label smoothing, loss tempering, unsupervised aux loss

Repeat 2,000 Times

2,000 x 1,281,167 prediction correctness for each example after training



i-th c<mark>o</mark>lumn

2,000 x 1,281,167 binary mask of which train examples are excluded



50,000 x 2,000 prediction correctness for each test examples

2,000 x 1,281,167 binary mask of which train examples are excluded



Case Studies MNIST, CIFAR100 & ImageNet

Which examples are "memorized"?

Are the "memorized" examples useful?

Is the utility of "memorized" examples consistent with our intuition?

Visualization of Memorized Examples

ImageNet "bobsled"



CIFAR-100 "bee"



ImageNet "black swan"



MNIST 2, 3, 5, 6



Are Memorized Examples Useful?



Effect on the test set accuracy of removing examples with memorization value estimate above a given threshold and the same number of randomly chosen examples. Fraction of the training set remaining after the removal is in the bottom plots. Shaded area in the accuracy represents one standard deviation on 100 (CIFAR-100, MNIST) and 5 (ImageNet) trials.

Why the Memorized Examples Are Useful?

Finding high influence train-test pairs

Criterion: a pair (i,j) is selected if mem(i) ≥ 0.25 and infl(i,j) ≥ 0.15

	MNIST	CIFAR-100	ImageNet
# High influence pairs	35	1015	1641
# Unique test examples selected	33	888	1462
# Test examples influenced by a single training example	31	774	1298

Removing the 964 unique training examples in these pairs on CIFAR-100 reduces the test accuracy by 2.46 \pm 0.36%, which is comparable to the effect of removing 11,000 random examples. Essentially all of that effect on the accuracy comes from the drop in accuracy on the test examples in the high-influence pairs from 72.1 \pm 1.3% to 45.4 \pm 1.4%.

The benefits of memorized examples are fully captured by their high influence on individual test examples.

Histogram of Influences in High-influence Pairs



MNIST is a very easy dataset with low diversity in the input examples. There are very few training examples that need to be memorized, or could induce high influence on a test example.

Visualization of High-influence Pairs on ImageNet



	M=0.8970	I=0.7029	I=0.0356	I=0.0317	I=0.0272	I=0.0246
toy poodle	Ŕ		P			
ļ	M=0.9663	I=0.7066	I=0.0599	I=0.0280	I=0.0233	I=0.0226
chain					0	
	M=0.9392	I=0.7181	I=0.0360	I=0.0167	I=0.0137	I=0.0104
sea lion					1	A A
	M=0.9066	I=0.7182	I=0.0100	I=0.0081	I=0.0058	I=0.0044
giant panda				P		
	M=0.8241	I=0.7483	I=0.0634	I=0.0368	I=0.0183	I=0.0145
kit fox		6.9		À	-C* -Exc	Contraction of the second seco
	M=0.9598	I=0.7946	I=0.0288	I=0.0284	I=0.0269	I=0.0233
rain barrel						
	M=0.5604	I=0.8451	I=0.0319	I=0.0301	I=0.0298	I=0.0265
pot						

Visualization of High-influence Pairs on ImageNet



bassoon	M=0.3721	I=0.2604	I=0.0469	I=0.0383	I=0.0265	I=0.0211
/	M=0.5733	I=0.2610	I=0.1281	I=0.0905	I=0.0637	I=0.0540
wing			tool .	Lon Lon		
	M 0 5061		I=0.0413	I=0.0371	I=0.0233	1-0.0230
silky terrier 1		1=0.2621				
	M=0.9504	I=0.2625	I=0.0312	I=0.0238	I=0.0197	I=0.0139
koala	A A					Ś
	M=0.3792	I=0.2626	I=0.0364	I=0.0219	I=0.0063	I=0.0053
black widow				K		1 Jac
	M=0.7842	I=0.2626	I=0.0594	I=0.0513	I=0.0399	1-0.0310
mailbag						
	M=0.9742	I=0.2629	I=0.0415	I=0.0386	I=0.0369	I=0.0325
fountain			CREV GAL		J.	THE SECOND

Visualization of High-influence Pairs on ImageNet





Visualization of High Influence Pairs on CIFAR-100

0.075

0.030

0.025

0.032

0.054



Top-ranking pairs on CIFAR-100 are near duplicates.



tt:5434 infl=0.031

tt-5353 infl=0.024

tt-1501 infl=0.028

tt-7571 infl=0.059

tt-5979 infl=0.026

9

tt-9968 infl=0.023

-

tt-5586 infl=0.064

tt-4043 infl=0.022

High Influence Pairs on MNIST

Less interesting but still show some visual correlation, and some ambiguous / mislabeled examples.



Theoretical Characterization

Construction of a simplified learning model with long tail data distribution to demonstrate that optimal performance cannot be achieved without memorization.

The cost of not fitting: a discrete case

To capture the main phenomenon we are interested in, we start by considering a simple and general prediction problem in which the domain does not have any underlying structure (such as the notion of distance). The domains X and Y are discrete, |X| = N and |Y| = m (for concreteness one can think of X = [N] and Y = [m]).

Theorem 2.3. Let π be a frequency prior with a corresponding marginal frequency distribution $\overline{\pi}^N$, and \mathcal{F} be a distribution over Y^X . Then for every learning algorithm \mathcal{A} and every dataset $Z \in (X \times Y)^n$:

where

theorem from: https://arxiv.org/abs/1906.05271

The cost of not fitting: a discrete case

Theorem 2.3. Let π be a frequency prior with a corresponding marginal frequency distribution $\overline{\pi}^N$, and \mathcal{F} be a distribution over Y^X . Then for every learning algorithm \mathcal{A} and every dataset $Z \in (X \times Y)^n$:



where

- tau_1 can be lower bounded, with characterization of the tail distribution.
- The discrete setup can be extended to continuous case of mixture models.

06

Further Investigations

Are influence and memorization consistent across architectures? In which layer does memorization happen?

Consistency of Estimation Across Architectures



Comparing two estimations:

For each threshold artheta

- → Select the subset of examples using ϑ by each of the estimation
- → Compare the two selected subsets by Jaccard similarity (Intersection-over-Union)
- → Also compare by the average difference of the estimations on the union of the two selected subsets

Some variations are observed, but overall the estimations are quite consistent across different architectures.

In Which Layer Does Memorization Happen?



- The body is considered to be computing feature representations for the inputs -- this feature is commonly finetuned or even directly reused in downstream tasks.
- The head is usually a densely connected linear layer, which could have potentially representation power to memorize the label mapping.

The "body": computing the representations The "head": computing the classification

Influence Estimation via "Transfer Learning"



Random 70% subset

The feature representation (trained on the full training set) is fixed and re-used, only the head is trained (from fresh random init each run).

Training Pipeline

data-augmentation, regularization, stable initialization, SoTA activation function, fancy Ir schedule, momentum, preconditioning, label smoothing, loss tempering, unsupervised aux loss

Repeat 2,000 Times

Influence Estimation via "Transfer Learning"

If memorization indeed happens in the classifier layer, our "transfer-learning" style experiments that trains only linear classifier on pre-trained representations could extract the influence estimation at a much lower cost than training many full ResNet models.

75.9%

ResNet50 trained on CIFAR-100 training set

72.3 ± 0.3%

4k **ResNet50** trained on 70% of CIFAR-100 training set

75.8 ± 0.1%

4k Linear Models trained on 70% of CIFAR-100 training set with pre-trained features

Influence Estimation via "Transfer Learning"

Number of training examples with memorization estimates \geq 0.25:

18,099

4k ResNet50 trained on 70% of CIFAR-100 training set 38

4k Linear Models trained on 70% of CIFAR-100 training set with pre-trained features



Conclusion

- A formal definition of memorization is given, efficiently estimated and shown to be consistent with high level intuition.
- Memorized examples are useful for test performance.
- The benefits of memorized examples are explained by strong influence on individual test examples.
- The observations verified the theoretical model of learning with long-tailed distribution [Fel19] in real data.