# Understanding Deep Learning
# (Still) Requires Rethinking Generalization

Alban MARIE [†]

March 10, 2022
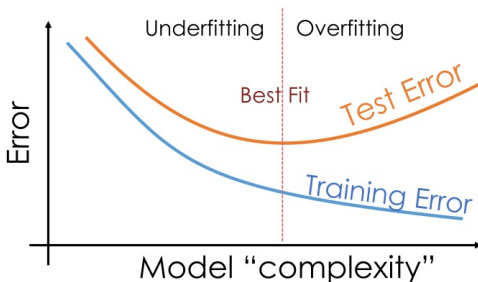
[†]Univ Rennes, INSA Rennes, CNRS, IETR - UMR6164, France
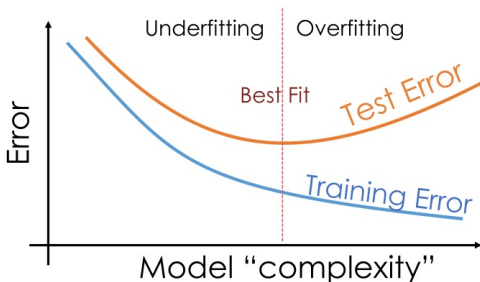
# Why something is wrong with deep learning

## BIAS VS VARIANCE



**Training** vs **Test** Error

## BIAS VS VARIANCE



**Training** vs **Test** Error

Underfitting | Overfitting

Best Fit

Test Error

Training Error

Error

Model "complexity"

➜ Where is deep learning on the x-axis?

**UNDERFITTING VS OVERFITTING**

**Where is deep learning on the x-axis?**

**Authors perform a simple experimental framework to propose an answer to this question.**

Why something is wrong with deep learning
What about regularization?
Bridge the gap with theory
Conclusion

## EXPERIMENTAL FRAMEWORK

➜ True labels

 Truck

 Cat

 Bird house

 Container ship

 Russian airplane
probably

 Dog

Why something is wrong with deep learning
What about regularization?
Bridge the gap with theory
Conclusion

# EXPERIMENTAL FRAMEWORK

➜ Random labels



| | | |
|---|---|---|
| | Truck | Russian airplane probably |
| | Cat | Cat |
| | Bird house | Container ship |
| | Container ship | Truck |
| | Russian airplane probably | Truck |
| | Dog | Bird house |

2

# EXPERIMENTAL FRAMEWORK

➜ Random labels



Russian airplane
probably

Cat

Container ship

Truck

Truck

Bird house

2

## EXPERIMENTAL FRAMEWORK

- → True labels
- → Random labels (previous slide)

## EXPERIMENTAL FRAMEWORK

- ➜ True labels
- ➜ Random labels (previous slide)
- ➜ Partially corrupted labels
  - ➜ Independently for each image and with a probability $p$, draw a random label for this image

Why something is wrong with deep learning    What about regularization?    Bridge the gap with theory    Conclusion

0000●00000    0000    000000    000

## **EXPERIMENTAL FRAMEWORK**

- ➜ True labels
- ➜ Random labels (previous slide)
- ➜ Partially corrupted labels
  - ➜ Independently for each image and with a probability $p$, draw a random label for this image
- ➜ Shuffled pixels
  - ➜ Select one random pixel permutation. Apply this permutation to all images.

Why something is wrong with deep learning    What about regularization?    Bridge the gap with theory    Conclusion

○○○○●○○○○○      ○○○○      ○○○○○○      ○○○

## EXPERIMENTAL FRAMEWORK

→ True labels

→ Random labels (previous slide)

→ Partially corrupted labels

     → Independently for each image and with a probability $p$, draw a random label for this image

→ Shuffled pixels

     → Select one random pixel permutation. Apply this permutation to all images.

→ Random pixels

     → Independently for each image, apply a random permutation.

## **EXPERIMENTAL FRAMEWORK**

→ True labels

→ Random labels (previous slide)

→ Partially corrupted labels

- → Independently for each image and with a probability $p$, draw a random label for this image

→ Shuffled pixels

- → Select one random pixel permutation. Apply this permutation to all images.

→ Random pixels

- → Independently for each image, apply a random permutation.

→ Gaussian pixels

- → Independently for each pixel, draw a random value from gaussian distribution with mean and std from original dataset

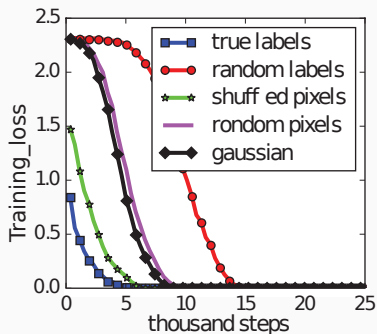**EXPERIMENTAL FRAMEWORK**

**Ok cool but ...**

**... what's the point of these experiments?**

## RESULTS



Training loss of true label experiment
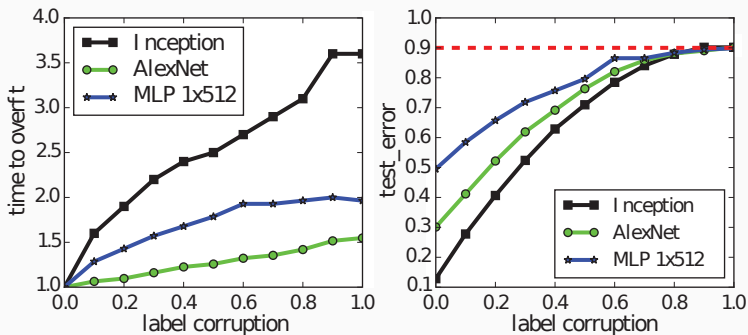decaying with the training steps on CIFAR10

## RESULTS



Training loss of various experiment settings
decaying with the training steps on CIFAR10

# RESULTS



**(left)** Relative convergence time with different label corruption ratio and **(right)** test error (also the generalization error since training error is 0) under different label corruptions.

## HOW BAD IS IT, DOCTOR?

➜ Training loss always converges to 0!

## HOW BAD IS IT, DOCTOR?

- → Training loss always converges to 0!
    - → The effective capacity of Neural Network (NN) is sufficient for memorizing the entire data set.

## HOW BAD IS IT, DOCTOR?

→ Training loss always converges to 0!
  → The effective capacity of Neural Network (NN) is sufficient for memorizing the entire data set.
  → NN are able to capture the remaining signal in the data if any, while at the same time fit the noisy part using brute-force.
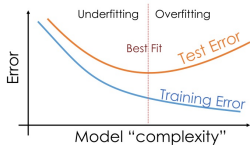
## HOW BAD IS IT, DOCTOR?

→ Training loss always converges to 0!

    → The effective capacity of Neural Network (NN) is sufficient for memorizing the entire data set.

    → NN are able to capture the remaining signal in the data if any, while at the same time fit the noisy part using brute-force.

→ Optimization remains easy, whatever you aim to fit.

## HOW BAD IS IT, DOCTOR?

→ Training loss always converges to 0!
  → The effective capacity of Neural Network (NN) is sufficient for memorizing the entire data set.
  → NN are able to capture the remaining signal in the data if any, while at the same time fit the noisy part using brute-force.
→ Optimization remains easy, whatever you aim to fit.
  → Easy even with random labels (the randomization breaks any relationship between the image and the label).

6

# BIAS VS VARIANCE



**Training** vs **Test** Error

➜ The effective capacity of NN is sufficient for memorizing the entire data set.

## BIAS VS VARIANCE



**Training** vs **Test** Error

→ The effective capacity of NN is sufficient for memorizing the entire data set.

  → Thus, very high overfitting

## BIAS VS VARIANCE



**Training** vs **Test** Error

→ The effective capacity of NN is sufficient for memorizing the entire data set.

    → Thus, very high overfitting

→ At the same time, increasing model complexity allow to reduce test error (and thus generalization error since training error is 0)

## BIAS VS VARIANCE



**Training** vs **Test** Error

→ The effective capacity of NN is sufficient for memorizing the entire data set.

    → Thus, very high overfitting

→ At the same time, increasing model complexity allow to reduce test error (and thus generalization error since training error is 0)

    → Something is wrong

7

# What about regularization?

## REGULARIZERS

**Popular Belief** Neural Network (NN) converges thanks to implicit and explicit regularizers

## REGULARIZERS

**Popular Belief** Neural Network (NN) converges thanks to implicit and explicit regularizers

Explicit regularizers are:

## REGULARIZERS

**Popular Belief** Neural Network (NN) converges thanks to implicit and explicit regularizers

Explicit regularizers are:

➜ Data augmentation

## REGULARIZERS

**Popular Belief**  Neural Network (NN) converges thanks to implicit and explicit regularizers

Explicit regularizers are:

→ Data augmentation

→ Weight decay ($L_2$ norm penalty on weights when too big)

## REGULARIZERS

**Popular Belief**   Neural Network (NN) converges thanks to implicit and explicit regularizers

Explicit regularizers are:

→ Data augmentation

→ Weight decay ($L_2$ norm penalty on weights when too big)

→ Dropout

## REGULARIZERS

**Popular Belief**   Neural Network (NN) converges thanks to implicit and explicit regularizers

Explicit regularizers are:

→ Data augmentation

→ Weight decay ($L_2$ norm penalty on weights when too big)

→ Dropout

Implicit regularizers are:

## REGULARIZERS

**Popular Belief**  Neural Network (NN) converges thanks to implicit and explicit regularizers

Explicit regularizers are:

→ Data augmentation

→ Weight decay ($L_2$ norm penalty on weights when too big)

→ Dropout

Implicit regularizers are:

→ Early stopping (stop training when generalization error is minimal)

## REGULARIZERS

**Popular Belief**   Neural Network (NN) converges thanks to implicit and explicit regularizers

Explicit regularizers are:

→ Data augmentation
→ Weight decay ($L_2$ norm penalty on weights when too big)
→ Dropout

Implicit regularizers are:

→ Early stopping (stop training when generalization error is minimal)
→ SGD

## REGULARIZERS

**Popular Belief** Neural Network (NN) converges thanks to implicit and explicit regularizers
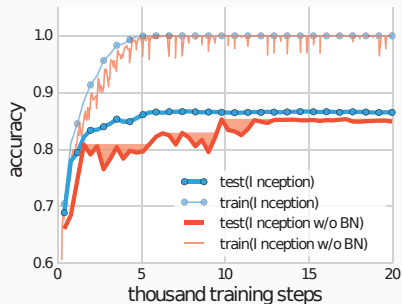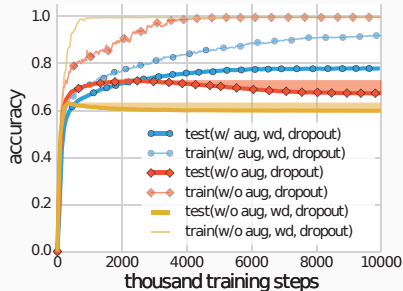
Explicit regularizers are:

→ Data augmentation
→ Weight decay ($L_2$ norm penalty on weights when too big)
→ Dropout

Implicit regularizers are:

→ Early stopping (stop training when generalization error is minimal)
→ SGD
→ NN architecture

8

## REGULARIZERS EXPERIMENTS



Regularizers impact on generalization for **(left)** Imagenet and **(right)** CIFAR10.
**Data augmentation**, **weight decay** and **batch normalization**
are referred as **aug**, **wd** and **BN**, respectively.

## REGULARIZERS CONCLUSIONS

➜ NN training loss converge to $0$, with or without regularizers

## REGULARIZERS CONCLUSIONS

➜ NN training loss converge to $0$, with or without regularizers
  ➜ Training space with regularizers is still huge

## REGULARIZERS CONCLUSIONS

→ NN training loss converge to $0$, with or without regularizers
   → Training space with regularizers is still huge
→ Regularizers improve generalization but...

## REGULARIZERS CONCLUSIONS

→ NN training loss converge to $0$, with or without regularizers

    → Training space with regularizers is still huge

→ Regularizers improve generalization but...

    → it is not necessary for NN to converge

## REGULARIZERS CONCLUSIONS

→ NN training loss converge to $0$, with or without regularizers
  → Training space with regularizers is still huge
→ Regularizers improve generalization but...
  → it is not necessary for NN to converge
  → it is unlikely that regularizers are the fundamental reason for generalization

## REGULARIZERS CONCLUSIONS

➜ NN training loss converge to $0$, with or without regularizers
   ➜ Training space with regularizers is still huge
➜ Regularizers improve generalization but...
   ➜ it is not necessary for NN to converge
   ➜ it is unlikely that regularizers are the fundamental reason for
     generalization
➜ Implicit Regularization with NN architecture is more powerful to
  reduce generalization error

# Bridge the gap with theory

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $f$ be a classification model with weights $\theta$ that aims to predicts labels $y_{i,\ i\in\{1,...N\}}$ based on input features $x_{i,\ i\in\{1,...\ N\}}$.

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $f$ be a classification model with weights $\theta$ that aims to predicts labels $y_{i,\ i\in\{1,...N\}}$ based on input features $x_{i,\ i\in\{1,...\ N\}}$.

It it said that $f$ **shatters** a dataset with $N \in \mathbb{N}$ elements if there exists a configuration for $\theta$ such that model $f$ makes no errors while predicting $y_i$ based on $x_i$ for each $i \in \{1, \ldots\ N\}$.

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $f$ be a classification model with weights $\theta$ that aims to predicts labels $y_{i,\ i\in\{1,...N\}}$ based on input features $x_{i,\ i\in\{1,...\ N\}}$.

It it said that $f$ **shatters** a dataset with $N \in \mathbb{N}$ elements if there exists a configuration for $\theta$ such that model $f$ makes no errors while predicting $y_i$ based on $x_i$ for each $i \in \{1, \dots N\}$.
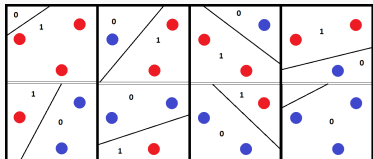
The **VC dimension** of a model $f$ is greater than or equal to $N$ if there exists at least one set of $N$ points where $f$ can shatter all arrangements.

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $f$ be a classification model with weights $\theta$ that aims to predicts labels $y_{i,\ i\in\{1,...N\}}$ based on input features $x_{i,\ i\in\{1,...\ N\}}$.

It it said that $f$ **shatters** a dataset with $N \in \mathbb{N}$ elements if there exists a configuration for $\theta$ such that model $f$ makes no errors while predicting $y_i$ based on $x_i$ for each $i \in \{1, \ldots\ N\}$.

The **VC dimension** of a model $f$ is greater than or equal to $N$ if there exists at least one set of $N$ points where $f$ can shatter all arrangements.
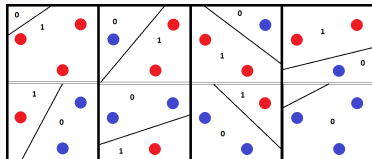


VC dim $\geq 3$

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $f$ be a classification model with weights $\theta$ that aims to predicts labels $y_{i, \, i \in \{1, \dots N\}}$ based on input features $x_{i, \, i \in \{1, \dots N\}}$.

It it said that $f$ **shatters** a dataset with $N \in \mathbb{N}$ elements if there exists a configuration for $\theta$ such that model $f$ makes no errors while predicting $y_i$ based on $x_i$ for each $i \in \{1, \dots N\}$.

The **VC dimension** of a model $f$ is greater than or equal to $N$ if there exists at least one set of $N$ points where $f$ can shatter all arrangements.



There is no set of $4$ points that can be shattered by a line.

VC dim $\geq 3$                    VC dim $< 4$

11

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $N$ be the size of the dataset and $D_{VC}$ the VC dimension of model $f$. With probability $1 - \delta$:

$$err_{test} \leq err_{train} + \sqrt{\frac{1}{N}[D_{VC}(\log(\frac{2N}{D_{VC}}) + 1) - \log(\frac{\delta}{4})]}$$

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $N$ be the size of the dataset and $D_{VC}$ the VC dimension of model $f$. With probability $1 - \delta$:

$$err_{test} \leq err_{train} + \sqrt{\frac{1}{N}[D_{VC}(\log(\frac{2N}{D_{VC}}) + 1) - \log(\frac{\delta}{4})]}$$

✘ For deep learning, $D_{VC} \gg N$ most of the time (complex solutions)

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $N$ be the size of the dataset and $D_{VC}$ the VC dimension of model $f$. With probability $1 - \delta$:

$$err_{test} \leq err_{train} + \sqrt{\frac{1}{N}[D_{VC}(\log(\frac{2N}{D_{VC}}) + 1) - \log(\frac{\delta}{4})]}$$

✘ For deep learning, $D_{VC} \gg N$ most of the time (complex solutions)

→ The effective capacity of NN is sufficient for memorizing the entire data set.

## VC DIMENSION - STATISTICAL LEARNING THEORY

Let $N$ be the size of the dataset and $D_{VC}$ the VC dimension of model $f$. With probability $1 - \delta$:

$$err_{test} \leq err_{train} + \sqrt{\frac{1}{N}[D_{VC}(\log(\frac{2N}{D_{VC}}) + 1) - \log(\frac{\delta}{4})]}$$

✘ For deep learning, $D_{VC} \gg N$ most of the time (complex solutions)

➔ The effective capacity of NN is sufficient for memorizing the entire data set.

   ➔ Thus, $D_{VC} \geq N$

## NN FINITE SAMPLE EXPRESSIVITY

➜ Most work in the litterature try to characterize NN expressivity at the **population level**

## NN FINITE SAMPLE EXPRESSIVITY

→ Most work in the litterature try to characterize NN expressivity at the **population level**

→ **population level**: infinite sample size (dataset with infinite number of elements)

## NN FINITE SAMPLE EXPRESSIVITY

→ Most work in the litterature try to characterize NN expressivity at the **population level**

→ **population level**: infinite sample size (dataset with infinite number of elements)

→ Instead, authors propose to express NN expressivity on a finite sample size $N$

## NN FINITE SAMPLE EXPRESSIVITY

→ Most work in the litterature try to characterize NN expressivity at the **population level**

→ **population level**: infinite sample size (dataset with infinite number of elements)

→ Instead, authors propose to express NN expressivity on a finite sample size $N$
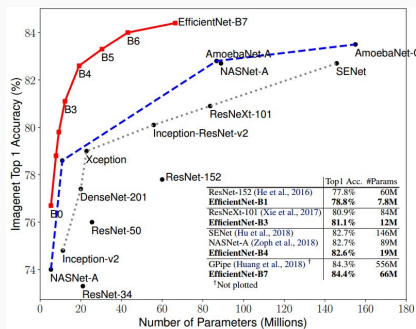
**Theorem 1** *There exists a two-layer neural network with ReLU activations and $2N + d$ weights that can represent any function on a sample of size $N$ in $d$ dimensions.*

## NN FINITE SAMPLE EXPRESSIVITY

| Dataset | $N$ | $d$ | $2N + d$ |
|---|---|---|---|
| MNIST | 70.000 | $28^2 = 784$ | 140.784 |
| CIFAR10 | 50.000 | $3 \times 32^2 = 3.072$ | 103.072 |
| ImageNet | 1.281.165 | $3 \times 224^2 = 150.528$ | 1.431.693 |



Number of parameters of ImageNet state-of-the-art models

14

**NN FINITE SAMPLE EXPRESSIVITY**

**This explains why NN manage to have $0$ training error on random labels**

**They just have way too many parameters!**

## Conclusion

## CONCLUSION

→ Authors propose a simple framework to evaluate NN expressivity

## CONCLUSION

→ Authors propose a simple framework to evaluate NN expressivity

→ Deep learning algorithms are large enough to shatter existing datasets, but still generalize to unseen examples

## CONCLUSION

→ Authors propose a simple framework to evaluate NN expressivity

→ Deep learning algorithms are large enough to shatter existing datasets, but still generalize to unseen examples

→ Regularizers is not the reason why

## CONCLUSION

→ Authors propose a simple framework to evaluate NN expressivity

→ Deep learning algorithms are large enough to shatter existing datasets, but still generalize to unseen examples

→ Regularizers is not the reason why

→ It remains easy to converge on data where generalization is impossible

## CONCLUSION

➜ Authors propose a simple framework to evaluate NN expressivity

➜ Deep learning algorithms are large enough to shatter existing datasets, but still generalize to unseen examples

➜ Regularizers is not the reason why

➜ It remains easy to converge on data where generalization is impossible

➜ Authors show an upper bound on the number of parameter for a NN with ReLU activations to represent any function

## CONCLUSION

→ Authors propose a simple framework to evaluate NN expressivity

→ Deep learning algorithms are large enough to shatter existing datasets, but still generalize to unseen examples

→ Regularizers is not the reason why

→ It remains easy to converge on data where generalization is impossible

→ Authors show an upper bound on the number of parameter for a NN with ReLU activations to represent any function

→ There is more in the paper, but I did not fully understand to present it

**Thank you for listening!**

**Any questions?**