

Large-Scale Machine Learning

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

Machine Learning

Outline

① When ML Meets Big Data

② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Outline

1 When ML Meets Big Data

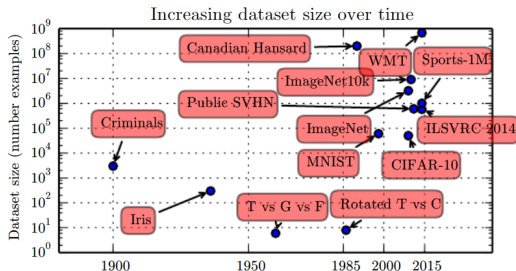
2 Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

3 Learning Theory Revisited

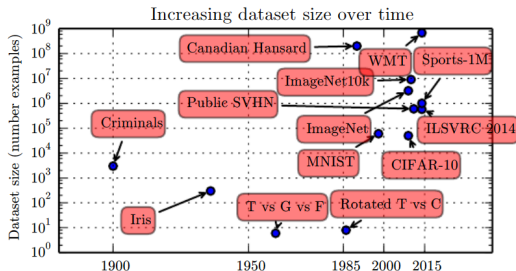
- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

The Big Data Era



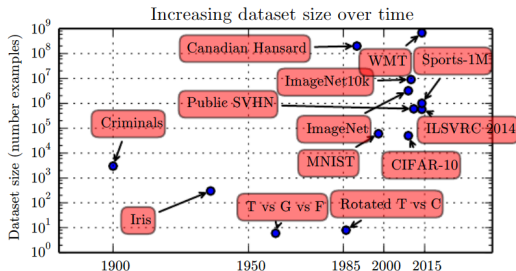
- Today, more and more of our activities are recorded by ubiquitous computing devices

The Big Data Era



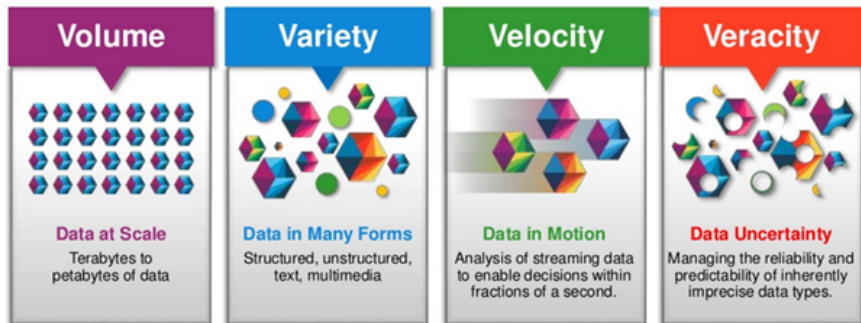
- Today, more and more of our activities are recorded by ubiquitous computing devices
- Networked computers make it easy to centralize these records and curate them into a **big** dataset

The Big Data Era



- Today, more and more of our activities are recorded by ubiquitous computing devices
- Networked computers make it easy to centralize these records and curate them into a **big** dataset
- **Large-scale machine learning** techniques solve problems by leveraging the posteriori knowledge learned from the big data

Characteristics of Big Data



Challenges of Large-Scale ML

- Variety and veracity
 - Feature engineering gets *even harder*

Challenges of Large-Scale ML

- Variety and veracity
 - Feature engineering gets *even harder*
 - Multi-task/transfer learning



A group of young people playing a game of Frisbee

Challenges of Large-Scale ML

- Variety and veracity
 - Feature engineering gets *even harder*
 - Multi-task/transfer learning
- Volume
 - Large D : curse of dimensionality
 - Large N : training efficiency



A group of young people playing a game of Frisbee

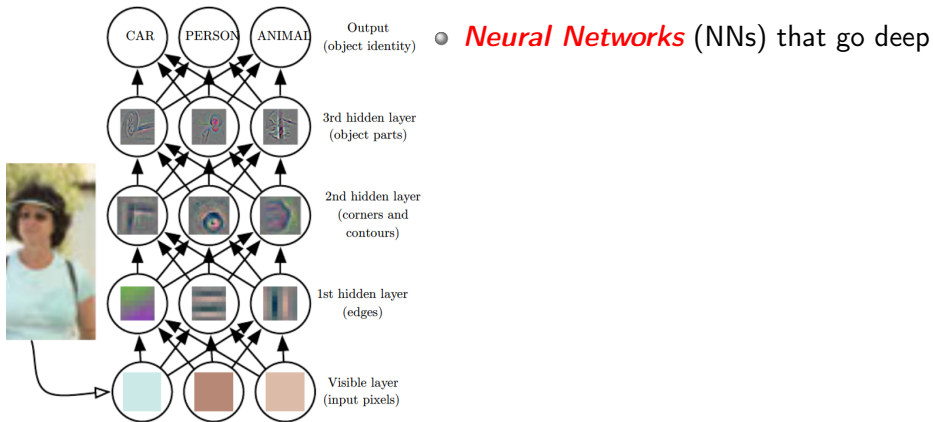
Challenges of Large-Scale ML

- Variety and veracity
 - Feature engineering gets *even harder*
 - Multi-task/transfer learning
- Volume
 - Large D : curse of dimensionality
 - Large N : training efficiency
- Velocity
 - Online learning

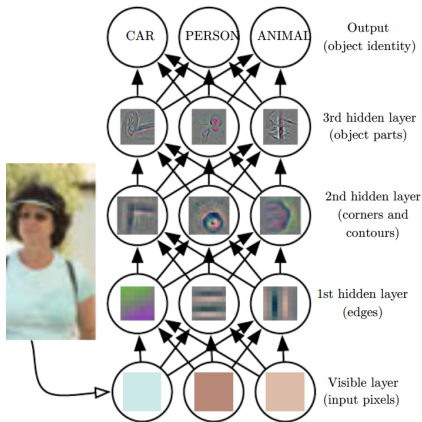


A group of young people playing a game of Frisbee

Advantages of Deep Learning

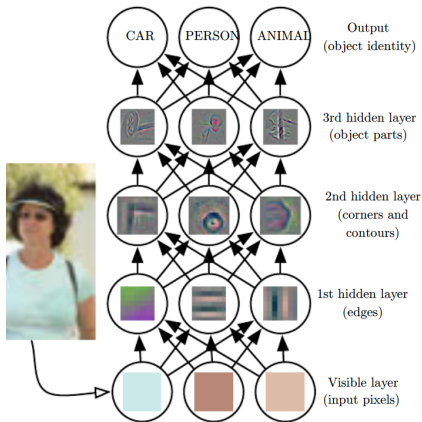


Advantages of Deep Learning



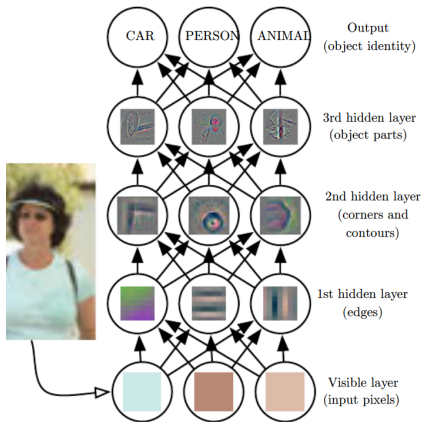
- **Neural Networks** (NNs) that go deep
- Automatic feature engineering
 - A kind of **representation learning**

Advantages of Deep Learning



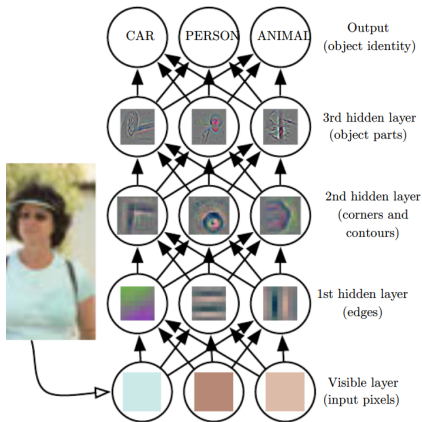
- **Neural Networks** (NNs) that go deep
- Automatic feature engineering
 - A kind of **representation learning**
- Exponential gain of expressiveness
 - Counters the curse of dimensionality

Advantages of Deep Learning



- **Neural Networks** (NNs) that go deep
- Automatic feature engineering
 - A kind of **representation learning**
- Exponential gain of expressiveness
 - Counters the curse of dimensionality
- Memory and GPU friendliness
 - SGD
 - GPU-based parallelism

Advantages of Deep Learning



- **Neural Networks** (NNs) that go deep
- Automatic feature engineering
 - A kind of **representation learning**
- Exponential gain of expressiveness
 - Counters the curse of dimensionality
- Memory and GPU friendliness
 - SGD
 - GPU-based parallelism
- Supporting online & transfer learning

Is Deep Learning a Panacea?

- I have big data, so I have to use deep learning

Is Deep Learning a Panacea?

- I have big data, so I have to use deep learning
- *Wrong!* No free lunch theorem: there is no single ML algorithm that outperforms others in every domain

Is Deep Learning a Panacea?

- I have big data, so I have to use deep learning
- **Wrong!** No free lunch theorem: there is no single ML algorithm that outperforms others in every domain
- Deep learning is more useful when the function f to learn is **complex** (nonlinear to the input dimension) and has **composed patterns**
 - E.g., image recognition, natural language processing

Is Deep Learning a Panacea?

- I have big data, so I have to use deep learning
- **Wrong!** No free lunch theorem: there is no single ML algorithm that outperforms others in every domain
- Deep learning is more useful when the function f to learn is **complex** (nonlinear to the input dimension) and has **composed patterns**
 - E.g., image recognition, natural language processing
- For simple (linear) f , there are specialized large-scale ML techniques (e.g., LIBLINEAR [6]) that are much more efficient

Outline

① When ML Meets Big Data

② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Outline

① When ML Meets Big Data

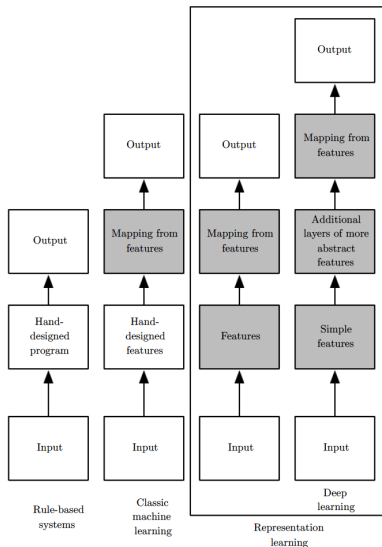
② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

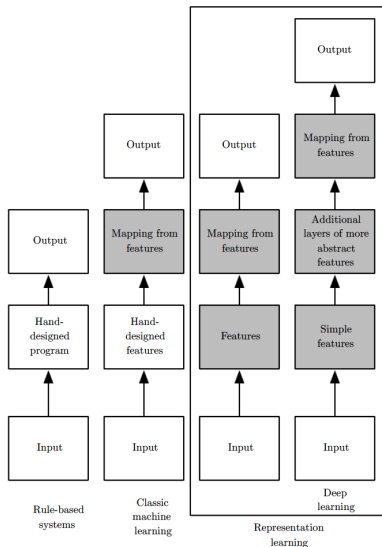
- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Representation Learning



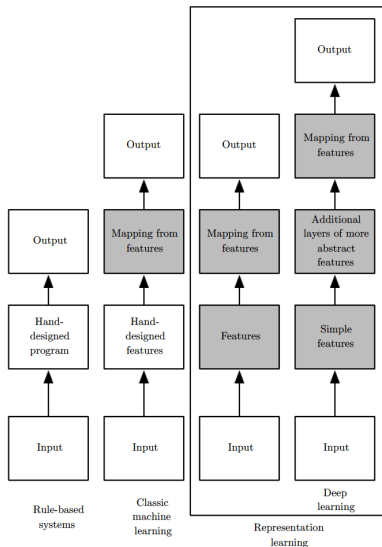
- Gray boxes are learned automatically

Representation Learning



- Gray boxes are learned automatically
- Deep learning maps the **most abstract** (deepest) features to the output
 - Usually, a simple linear function suffices

Representation Learning

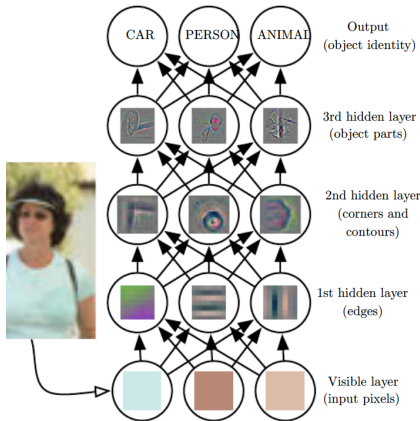


- Gray boxes are learned automatically
- Deep learning maps the **most abstract** (deepest) features to the output
 - Usually, a simple linear function suffices
- In deep learning, features/presentations are **distributed**

Distributed Representations of Data

- In deep learning, we assume that x 's were generated by *compositions of factors*, potentially *at multiple levels* in a hierarchy

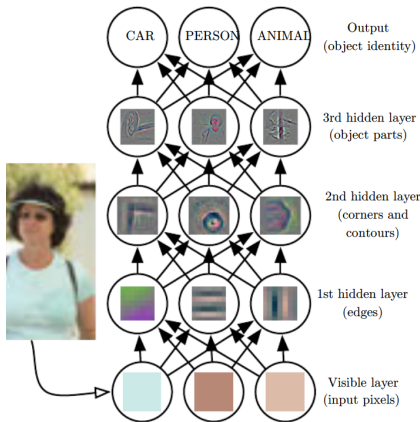
- E.g., layer 3: face = 0.3 [corner] + 0.7 [circle] + 0 [curve]



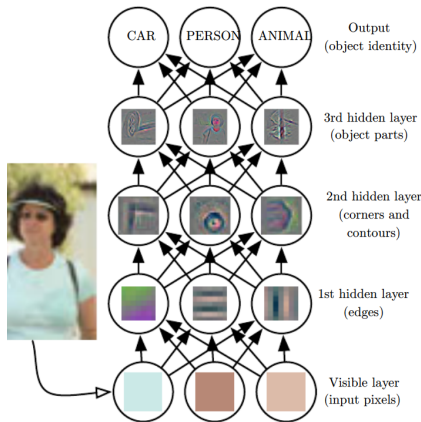
Distributed Representations of Data

- In deep learning, we assume that x 's were generated by *compositions of factors*, potentially *at multiple levels* in a hierarchy

- E.g., layer 3: face = 0.3 [corner] + 0.7 [circle] + 0 [curve]
- [.] a predefined non-linear function
- Weights (arrows) learned from training examples

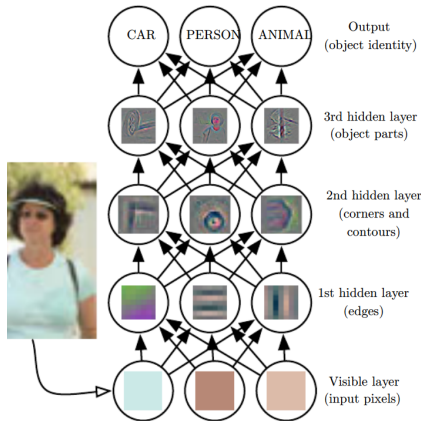


Distributed Representations of Data



- In deep learning, we assume that x 's were generated by *compositions of factors*, potentially *at multiple levels* in a hierarchy
 - E.g., layer 3: $\text{face} = 0.3 [\text{corner}] + 0.7 [\text{circle}] + 0 [\text{curve}]$
 - $[\cdot]$ a predefined non-linear function
 - **Weights** (arrows) learned from training examples
- Given x , factors at the same level output *a layer of features* of x
 - Layer 2: 1, 2, 0.5 for $[\text{corner}]$, $[\text{circle}]$, and $[\text{curve}]$ respectively

Distributed Representations of Data



- In deep learning, we assume that x 's were generated by *compositions of factors*, potentially *at multiple levels* in a hierarchy
 - E.g., layer 3: $\text{face} = 0.3 [\text{corner}] + 0.7 [\text{circle}] + 0 [\text{curve}]$
 - $[\cdot]$ a predefined non-linear function
 - **Weights** (arrows) learned from training examples
- Given x , factors at the same level output *a layer of features* of x
 - Layer 2: 1, 2, 0.5 for $[\text{corner}]$, $[\text{circle}]$, and $[\text{curve}]$ respectively
- To be fed into the factors in the next (deeper) level
 - $\text{Face} = 0.3 * 1 + 0.7 * 2$

Outline

① When ML Meets Big Data

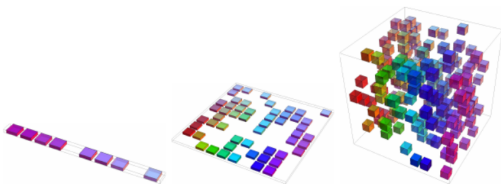
② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

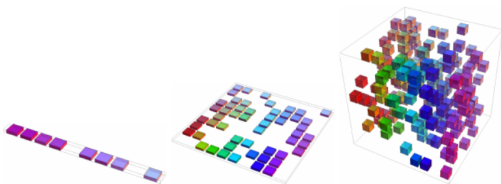
Curse of Dimensionality



- Most classic nonlinear ML models find θ by assuming function smoothness:

$$\text{if } \mathbf{x} \sim \mathbf{x}^{(i)} \in \mathbb{X}, \text{ then } f(\mathbf{x}; \mathbf{w}) \sim f(\mathbf{x}^{(i)}; \mathbf{w})$$

Curse of Dimensionality



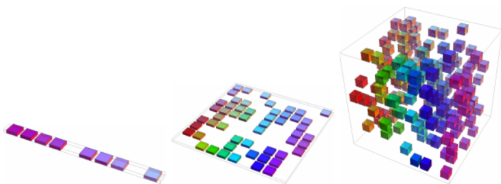
- Most classic nonlinear ML models find θ by assuming function smoothness:

$$\text{if } \mathbf{x} \sim \mathbf{x}^{(i)} \in \mathbb{X}, \text{ then } f(\mathbf{x}; \mathbf{w}) \sim f(\mathbf{x}^{(i)}; \mathbf{w})$$

- E.g., the non-parametric methods predict the label \hat{y} of \mathbf{x} by simply interpolating the labels of examples $\mathbf{x}^{(i)}$'s **close to \mathbf{x}** :

$$\hat{y} = \sum_i \alpha_i y^{(i)} k(\mathbf{x}^{(i)}, \mathbf{x}) + b, \text{ where } k(\mathbf{x}^{(i)}, \mathbf{x}) = \exp(-\gamma \|\mathbf{x}^{(i)} - \mathbf{x}\|^2)$$

Curse of Dimensionality



- Most classic nonlinear ML models find θ by assuming function smoothness:

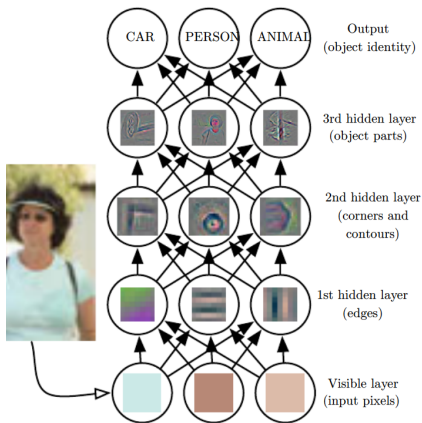
$$\text{if } \mathbf{x} \sim \mathbf{x}^{(i)} \in \mathbb{X}, \text{ then } f(\mathbf{x}; \mathbf{w}) \sim f(\mathbf{x}^{(i)}; \mathbf{w})$$

- E.g., the non-parametric methods predict the label \hat{y} of \mathbf{x} by simply interpolating the labels of examples $\mathbf{x}^{(i)}$'s *close to \mathbf{x}* :

$$\hat{y} = \sum_i \alpha_i y^{(i)} k(\mathbf{x}^{(i)}, \mathbf{x}) + b, \text{ where } k(\mathbf{x}^{(i)}, \mathbf{x}) = \exp(-\gamma \|\mathbf{x}^{(i)} - \mathbf{x}\|^2)$$

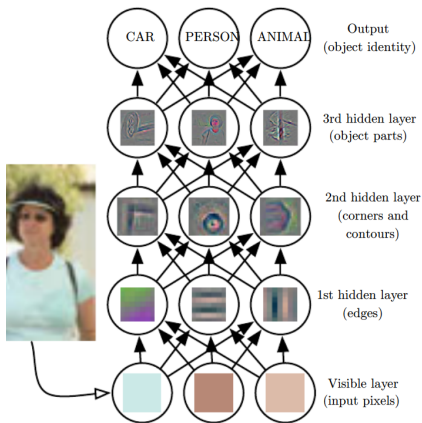
- Suppose f is smooth within a bin, we need *exponentially more examples* to get a good interpolation as D increases

Exponential Gains from Depth I



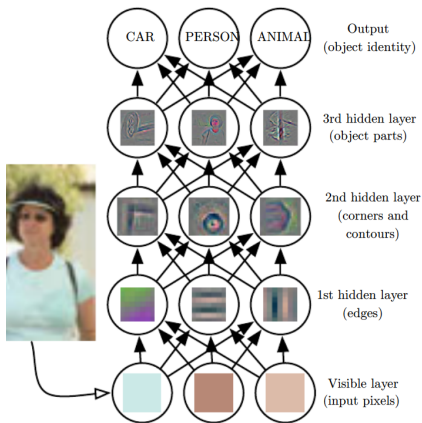
- Functions representable with a deep rectifier NN require an exponential number of hidden units in a shallow NN [12]

Exponential Gains from Depth I



- Functions representable with a deep rectifier NN require an exponential number of hidden units in a shallow NN [12]
- In deep learning, a deep factor is defined by “reusing” the shallow ones
 - Face = 0.3 [corner] + 0.7 [circle]

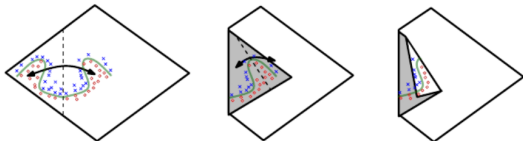
Exponential Gains from Depth I



- Functions representable with a deep rectifier NN require an exponential number of hidden units in a shallow NN [12]
- In deep learning, a deep factor is defined by “reusing” the shallow ones
 - Face = 0.3 [corner] + 0.7 [circle]
- With a shallow structure, a deep factor needs to be replaced by *exponentially many* factors
 - Face = 0.3 [0.5 [vertical] + 0.5 [horizontal]] + 0.7 [...]

Exponential Gains from Depth II

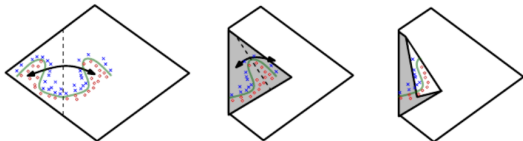
- Another example: an NN with absolute value rectification units



- Each hidden unit specifies where to fold the input space in order to create mirror responses (on both sides of the absolute value)
- A single fold in a deep layer creates an exponentially large number of piecewise linear regions in input space
 - No need to see examples in each linear regions in input space

Exponential Gains from Depth II

- Another example: an NN with absolute value rectification units



- Each hidden unit specifies where to fold the input space in order to create mirror responses (on both sides of the absolute value)
- A single fold in a deep layer creates an exponentially large number of piecewise linear regions in input space
 - No need to see examples in each linear regions in input space
- This exponential gain counters the exponential challenges posed by the curse of dimensionality

Outline

① When ML Meets Big Data

② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Stochastic Gradient Descent

Gradient Descent (GD)

$\mathbf{w}^{(0)} \leftarrow$ a random vector;

Repeat until convergence {

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} C_N(\mathbf{w}^{(t)}; \mathbb{X});$$

}

Stochastic Gradient Descent

Gradient Descent (GD)

$\mathbf{w}^{(0)} \leftarrow$ a random vector;

Repeat until convergence {

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} C_N(\mathbf{w}^{(t)}; \mathbb{X});$$

}

- Needs to scan the entire dataset to descent (many I/Os)

Stochastic Gradient Descent

Gradient Descent (GD)

$\mathbf{w}^{(0)} \leftarrow$ a random vector;
Repeat until convergence {
 $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} C_N(\mathbf{w}^{(t)}; \mathbb{X});$
}

- Needs to scan the entire dataset to descent (many I/Os)

(Mini-Batched) Stochastic Gradient Descent (SGD)

$\mathbf{w}^{(0)} \leftarrow$ a random vector;
Repeat until convergence {
 Randomly partition the training set \mathbb{X} into *minibatches* $\{\mathbb{X}^{(j)}\}_j$;
 $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} C(\mathbf{w}^{(t)}; \mathbb{X}^{(j)});$
}

Stochastic Gradient Descent

Gradient Descent (GD)

$\mathbf{w}^{(0)} \leftarrow$ a random vector;
Repeat until convergence {
 $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} C_N(\mathbf{w}^{(t)}; \mathbb{X});$
}

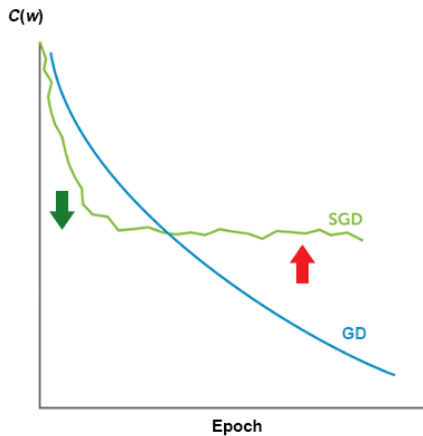
- Needs to scan the entire dataset to descent (many I/Os)

(Mini-Batched) Stochastic Gradient Descent (SGD)

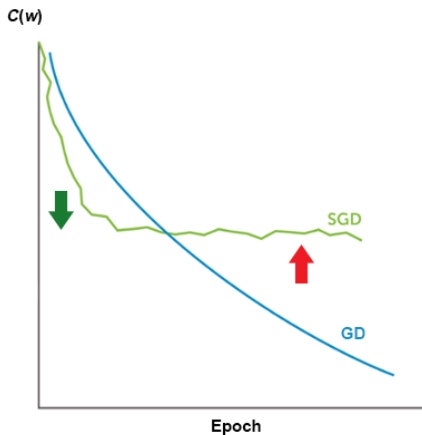
$\mathbf{w}^{(0)} \leftarrow$ a random vector;
Repeat until convergence {
 Randomly partition the training set \mathbb{X} into **minibatches** $\{\mathbb{X}^{(j)}\}_j$;
 $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} C(\mathbf{w}^{(t)}; \mathbb{X}^{(j)});$
}

- **No I/O** if the next mini-batch can be prefetched

GD vs. SGD



GD vs. SGD



- Is SGD really a better algorithm?

Yes, If You Have Big Data



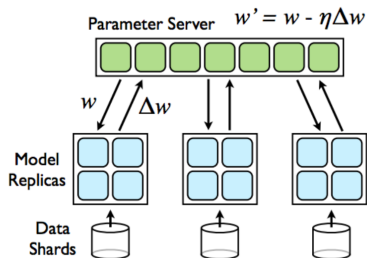
- Performance is limited by *training time*

Asymptotic Analysis [3]

	GD	SGD
Time per iteration	N	1
#Iterations to opt. error ρ	$\log \frac{1}{\rho}$	$\frac{1}{\rho}$
Time to opt. error ρ	$N \log \frac{1}{\rho}$	$\frac{1}{\rho}$
Time to excess error ε	$\frac{1}{\varepsilon^{1/\alpha}} \log \frac{1}{\varepsilon}$, where $\alpha \in [\frac{1}{2}, 1]$	$\frac{1}{\varepsilon}$

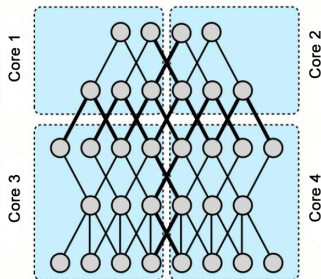
Parallelizing SGD

Data Parallelism



Every core/GPU trains the full model given partitioned data.

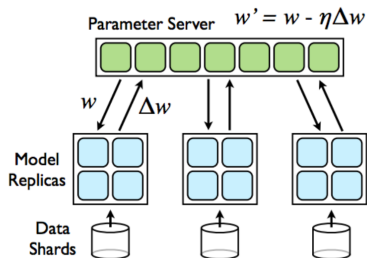
Model Parallelism



Every core/GPU train a partitioned model given full data.

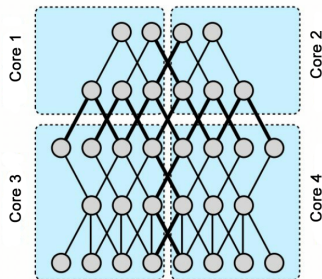
Parallelizing SGD

Data Parallelism



Every core/GPU trains the full model given partitioned data.

Model Parallelism



Every core/GPU train a partitioned model given full data.

- The effectiveness depends on applications and available hardware
 - E.g., CPU/GPU speed, communication latency, bandwidth, etc.

Outline

① When ML Meets Big Data

② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Online Learning

- So far, we assume that the training data \mathbb{X} comes at once
- What if data come sequentially?

Online Learning

- So far, we assume that the training data \mathbb{X} comes at once
- What if data come sequentially?
- *Online learning*: to update model when new data arrive

Online Learning

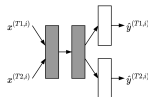
- So far, we assume that the training data \mathbb{X} comes at once
- What if data come sequentially?
- *Online learning*: to update model when new data arrive
- This is already supported by SGD

Muti-Task and Transfer Learning

- *Multi-task learning*: to learning a single model for multiple tasks
- *Transfer learning*: to reuse the knowledge learned from one task to help another

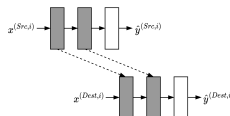
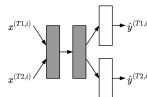
Muti-Task and Transfer Learning

- **Multi-task learning**: to learning a single model for multiple tasks
 - Via shared layers
- **Transfer learning**: to reuse the knowledge learned from one task to help another



Muti-Task and Transfer Learning

- **Multi-task learning**: to learning a single model for multiple tasks
 - Via shared layers
- **Transfer learning**: to reuse the knowledge learned from one task to help another
 - Via pretrained layers (whose weights may be further updated when a smaller learning rate)



Outline

① When ML Meets Big Data

② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Learning Theory

- How to learn a function f_N from N examples \mathbb{X} that is close to the true function f^* ?

Learning Theory

- How to learn a function f_N from N examples \mathbb{X} that is close to the true function f^* ?
- **Empirical risk:** $C_N[f_N] = \frac{1}{N} \sum_{i=1}^N \text{loss}(f_N(\mathbf{x}^{(i)}), y^{(i)})$
- **Expected risk:** $C[f_N] = \int \text{loss}(f(\mathbf{x}), y) dP(\mathbf{x}, y)$

Learning Theory

- How to learn a function f_N from N examples \mathbb{X} that is close to the true function f^* ?
- **Empirical risk:** $C_N[f_N] = \frac{1}{N} \sum_{i=1}^N \text{loss}(f_N(\mathbf{x}^{(i)}), y^{(i)})$
- **Expected risk:** $C[f_N] = \int \text{loss}(f(\mathbf{x}), y) dP(\mathbf{x}, y)$
- Let $f^* = \arg \min_f C[f]$ be the true function (our goal)

Learning Theory

- How to learn a function f_N from N examples \mathbb{X} that is close to the true function f^* ?
- **Empirical risk:** $C_N[f_N] = \frac{1}{N} \sum_{i=1}^N \text{loss}(f_N(\mathbf{x}^{(i)}), y^{(i)})$
- **Expected risk:** $C[f_N] = \int \text{loss}(f(\mathbf{x}), y) dP(\mathbf{x}, y)$
- Let $f^* = \arg \min_f C[f]$ be the true function (our goal)
- Since we are seeking a function in a model (hypothesis space) \mathbb{F} , this is what can have at best: $f_{\mathbb{F}}^* = \arg \min_{f \in \mathbb{F}} C[f]$

Learning Theory

- How to learn a function f_N from N examples \mathbb{X} that is close to the true function f^* ?
- **Empirical risk:** $C_N[f_N] = \frac{1}{N} \sum_{i=1}^N \text{loss}(f_N(\mathbf{x}^{(i)}), y^{(i)})$
- **Expected risk:** $C[f_N] = \int \text{loss}(f(\mathbf{x}), y) dP(\mathbf{x}, y)$
- Let $f^* = \arg \min_f C[f]$ be the true function (our goal)
- Since we are seeking a function in a model (hypothesis space) \mathbb{F} , this is what can have at best: $f_{\mathbb{F}}^* = \arg \min_{f \in \mathbb{F}} C[f]$
- But we only minimize errors on limited examples in our objective, so we only have $f_N = \arg \min_{f \in \mathbb{F}} C_N[f]$

Learning Theory

- How to learn a function f_N from N examples \mathbb{X} that is close to the true function f^* ?
- **Empirical risk**: $C_N[f_N] = \frac{1}{N} \sum_{i=1}^N \text{loss}(f_N(\mathbf{x}^{(i)}), y^{(i)})$
- **Expected risk**: $C[f_N] = \int \text{loss}(f(\mathbf{x}), y) dP(\mathbf{x}, y)$
- Let $f^* = \arg \min_f C[f]$ be the true function (our goal)
- Since we are seeking a function in a model (hypothesis space) \mathbb{F} , this is what can have at best: $f_{\mathbb{F}}^* = \arg \min_{f \in \mathbb{F}} C[f]$
- But we only minimize errors on limited examples in our objective, so we only have $f_N = \arg \min_{f \in \mathbb{F}} C_N[f]$
- The **excess error** $\mathcal{E} = C[f_N] - C[f^*]$:

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}}$$

Excess Error

- Wait, we may not have enough training time, so we stop iterations early and have \tilde{f}_N , where $C_N[\tilde{f}_N] \leq C_N[f_N] + \rho$

Excess Error

- Wait, we may not have enough training time, so we stop iterations early and have \tilde{f}_N , where $C_N[\tilde{f}_N] \leq C_N[f_N] + \rho$
- The excess error becomes $\mathcal{E} = C[\tilde{f}_N] - C[f^*]$:

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

Excess Error

- Wait, we may not have enough training time, so we stop iterations early and have \tilde{f}_N , where $C_N[\tilde{f}_N] \leq C_N[f_N] + \rho$
- The excess error becomes $\mathcal{E} = C[\tilde{f}_N] - C[f^*]$:

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- **Approximation error \mathcal{E}_{app}** : reduced by choosing a larger model

Excess Error

- Wait, we may not have enough training time, so we stop iterations early and have \tilde{f}_N , where $C_N[\tilde{f}_N] \leq C_N[f_N] + \rho$
- The excess error becomes $\mathcal{E} = C[\tilde{f}_N] - C[f^*]$:

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- **Approximation error \mathcal{E}_{app}** : reduced by choosing a larger model
- **Estimation error \mathcal{E}_{est}** : reduced by
 - ① Increasing N , or
 - ② Choosing smaller model [4, 11, 14]

Excess Error

- Wait, we may not have enough training time, so we stop iterations early and have \tilde{f}_N , where $C_N[\tilde{f}_N] \leq C_N[f_N] + \rho$
- The excess error becomes $\mathcal{E} = C[\tilde{f}_N] - C[f^*]$:

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- **Approximation error \mathcal{E}_{app}** : reduced by choosing a larger model
- **Estimation error \mathcal{E}_{est}** : reduced by
 - ① Increasing N , or
 - ② Choosing smaller model [4, 11, 14]
- **Optimization error \mathcal{E}_{opt}** : reduced by
 - ① Running optimization alg. longer (with smaller ρ)
 - ② Choosing more efficient optimization alg.

Minimizing Excess Error

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- Small-scale ML tasks:

Minimizing Excess Error

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- Small-scale ML tasks:
 - Mainly constrained by N

Minimizing Excess Error

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- Small-scale ML tasks:
 - Mainly constrained by N
 - Computing time is not an issue, so \mathcal{E}_{opt} can be insignificant by choosing small ρ

Minimizing Excess Error

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- Small-scale ML tasks:
 - Mainly constrained by N
 - Computing time is not an issue, so \mathcal{E}_{opt} can be insignificant by choosing small ρ
 - Size of hypothesis is important to balance the trade-off between \mathcal{E}_{app} and \mathcal{E}_{est}

Minimizing Excess Error

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- Small-scale ML tasks:
 - Mainly constrained by N
 - Computing time is not an issue, so \mathcal{E}_{opt} can be insignificant by choosing small ρ
 - Size of hypothesis is important to balance the trade-off between \mathcal{E}_{app} and \mathcal{E}_{est}
- Large-scale ML tasks:

Minimizing Excess Error

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- Small-scale ML tasks:
 - Mainly constrained by N
 - Computing time is not an issue, so \mathcal{E}_{opt} can be insignificant by choosing small ρ
 - Size of hypothesis is important to balance the trade-off between \mathcal{E}_{app} and \mathcal{E}_{est}
- Large-scale ML tasks:
 - Mainly constrained by time (significant \mathcal{E}_{opt}), so **SGD** is preferred

Minimizing Excess Error

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

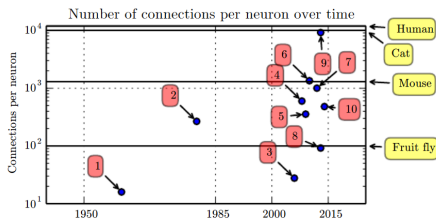
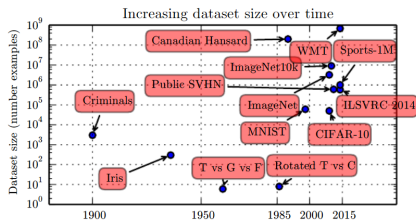
- Small-scale ML tasks:
 - Mainly constrained by N
 - Computing time is not an issue, so \mathcal{E}_{opt} can be insignificant by choosing small ρ
 - Size of hypothesis is important to balance the trade-off between \mathcal{E}_{app} and \mathcal{E}_{est}
- Large-scale ML tasks:
 - Mainly constrained by time (significant \mathcal{E}_{opt}), so **SGD** is preferred
 - N is large, so \mathcal{E}_{est} can be reduced

Minimizing Excess Error

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}} + \underbrace{C[\tilde{f}_N] - C[f_N]}_{\mathcal{E}_{\text{opt}}}$$

- Small-scale ML tasks:
 - Mainly constrained by N
 - Computing time is not an issue, so \mathcal{E}_{opt} can be insignificant by choosing small ρ
 - Size of hypothesis is important to balance the trade-off between \mathcal{E}_{app} and \mathcal{E}_{est}
- Large-scale ML tasks:
 - Mainly constrained by time (significant \mathcal{E}_{opt}), so **SGD** is preferred
 - N is large, so \mathcal{E}_{est} can be reduced
 - **Large model** is preferred to reduce \mathcal{E}_{app}

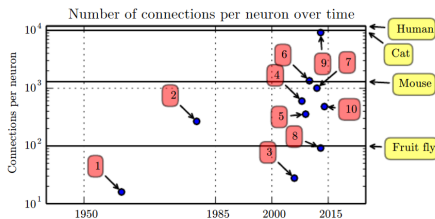
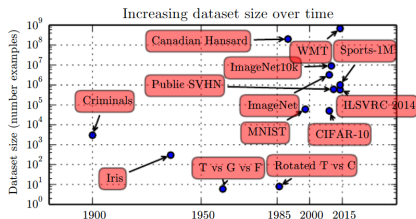
Big Data + Big Models



9. COTS HPC unsupervised convolutional network [5]

10. GoogleLeNet [13]

Big Data + Big Models



9. COTS HPC unsupervised convolutional network [5]

10. GoogleLeNet [13]

- With domain-specific architecture such as *convolutional NNs* (CNNs) and *recurrent NNs* (RNNs)

Outline

① When ML Meets Big Data

② Advantages of Deep Learning

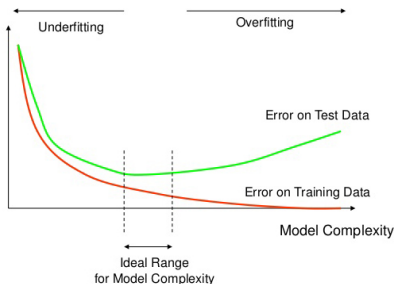
- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

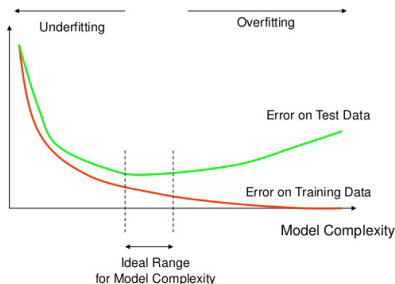
Over-Parametrized NNs

- Let $D^{(l)}$ be the output dimension (“width”) of a layer $f^{(l)}(\cdot; \theta^{(l)})$ of an NN
 - Examples $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{D^{(0)}} \times \mathbb{R}^{D^{(L)}}$
 - $D = D^{(1)} + \dots + D^{(L)}$ the total number of neurons
- From the statistical learning theory point of view, the larger the D , the worse the generalizability



Over-Parametrized NNs

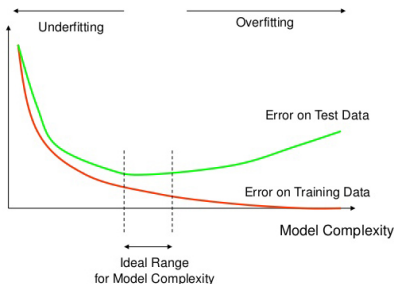
- Let $D^{(l)}$ be the output dimension (“width”) of a layer $f^{(l)}(\cdot; \theta^{(l)})$ of an NN
 - Examples $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{D^{(0)}} \times \mathbb{R}^{D^{(L)}}$
 - $D = D^{(1)} + \dots + D^{(L)}$ the total number of neurons
- From the statistical learning theory point of view, the larger the D , the worse the generalizability



- However, as D grows, the generalizability actually *increases* [19]; i.e., over-parametrization leads to better performance

Over-Parametrized NNs

- Let $D^{(l)}$ be the output dimension (“width”) of a layer $f^{(l)}(\cdot; \theta^{(l)})$ of an NN
 - Examples $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{D^{(0)}} \times \mathbb{R}^{D^{(L)}}$
 - $D = D^{(1)} + \dots + D^{(L)}$ the total number of neurons
- From the statistical learning theory point of view, the larger the D , the worse the generalizability



- However, as D grows, the generalizability actually *increases* [19]; i.e., over-parametrization leads to better performance
- Why such a paradox?

Wide-and-Deep NNs as Gaussian Processes

- Recent studies [9, 8, 10] show that *a wide NN of any depth can be approximated by a Gaussian process (GP)*
 - Either before, during, or after training
- Recall that a GP is a non-parametric model whose complexity depends only on the size of training set $|\mathbb{X}|$ and the hyperparameters of kernel function $k(\cdot, \cdot)$:

$$\begin{bmatrix} \mathbf{y}_N \\ \mathbf{y}_M \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_N \\ \mathbf{m}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N} & \mathbf{K}_{N,M} \\ \mathbf{K}_{M,N} & \mathbf{K}_{M,M} \end{bmatrix} \right)$$

with Bayesian inference for test points \mathbb{X}' :

$$\mathbf{P}(\mathbf{y}_M | \mathbb{X}', \mathbb{X}) = \mathcal{N}(\mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{y}_N, \mathbf{K}_{M,M} - \mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{K}_{N,M})$$

Wide-and-Deep NNs as Gaussian Processes

- Recent studies [9, 8, 10] show that *a wide NN of any depth can be approximated by a Gaussian process (GP)*
 - Either before, during, or after training
- Recall that a GP is a non-parametric model whose complexity depends only on the size of training set $|\mathbb{X}|$ and the hyperparameters of kernel function $k(\cdot, \cdot)$:

$$\begin{bmatrix} \mathbf{y}_N \\ \mathbf{y}_M \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_N \\ \mathbf{m}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N} & \mathbf{K}_{N,M} \\ \mathbf{K}_{M,N} & \mathbf{K}_{M,M} \end{bmatrix} \right)$$

with Bayesian inference for test points \mathbb{X}' :

$$\mathbf{P}(\mathbf{y}_M | \mathbb{X}', \mathbb{X}) = \mathcal{N}(\mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{y}_N, \mathbf{K}_{M,M} - \mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{K}_{N,M})$$

- Therefore, wide-and-deep NNs do not overfit as one may expect
 - The D , once becoming large, does *not* reflect true model complexity

Outline

① When ML Meets Big Data

② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Example: NN for Regression

- For simplicity, we consider an L -layer NN $f(\cdot; \theta)$ for the regression problem:

$$f(\mathbf{x}; \theta) = \mathbf{a}^{(l)} = \phi^{(l)}(\mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}), \text{ for } l = 1, \dots, L,$$

where

- the activation functions $\phi^{(1)}(\cdot) = \dots = \phi^{(L-1)}(\cdot) \equiv \phi(\cdot)$ and $\phi^{(L-1)}(\cdot)$ is an identity function
- $\mathbf{a}^{(0)} = \mathbf{x}$ and $\hat{y} = a^{(L)} = z^{(L)} \in \mathbb{R}$ the mean of a Gaussian
- $\theta^{(l)} = \text{vec}(\mathbf{W}^{(l)}, \mathbf{b}^{(l)})$ and $\theta = \text{vec}(\theta^{(1)}, \dots, \theta^{(L)})$

Example: NN for Regression

- For simplicity, we consider an L -layer NN $f(\cdot; \theta)$ for the regression problem:

$$f(\mathbf{x}; \theta) = \mathbf{a}^{(l)} = \phi^{(l)}(\mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}), \text{ for } l = 1, \dots, L,$$

where

- the activation functions $\phi^{(1)}(\cdot) = \dots = \phi^{(L-1)}(\cdot) \equiv \phi(\cdot)$ and $\phi^{(L-1)}(\cdot)$ is an identity function
- $\mathbf{a}^{(0)} = \mathbf{x}$ and $\hat{y} = a^{(L)} = z^{(L)} \in \mathbb{R}$ the mean of a Gaussian
- $\theta^{(l)} = \text{vec}(\mathbf{W}^{(l)}, \mathbf{b}^{(l)})$ and $\theta = \text{vec}(\theta^{(1)}, \dots, \theta^{(L)})$
- Let $\hat{\mathbf{y}} = [f(\mathbf{x}^{(1)}; \theta), \dots, f(\mathbf{x}^{(N)}; \theta)]^\top \in \mathbb{R}^N$ be the predictions for the points $\mathbf{X} \in \mathbb{R}^{N \times D^{(0)}}$ in training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i = \{\mathbf{X}, \mathbf{y}\}$
- Maximum-likelihood estimation:

$$\arg \max_{\theta} \mathbb{P}(\mathbb{X} | \theta) = \arg \min_{\theta} C(\hat{\mathbf{y}}, \mathbf{y}) = \arg \min_{\theta} \frac{1}{2} \|\hat{\mathbf{y}} - \mathbf{y}\|^2$$

Weight Initialization and Normalization

$$\mathbf{a}^{(l)} = \phi^{(l)}(\mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)})$$

- Common initialization: $W_{i,j}^{(l)} \sim \mathcal{N}(0, \sigma_w^2)$ and $b_i^{(l)} \sim \mathcal{N}(0, \sigma_b^2)$

Weight Initialization and Normalization

$$\mathbf{a}^{(l)} = \phi^{(l)}(\mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)})$$

- Common initialization: $W_{ij}^{(l)} \sim \mathcal{N}(0, \sigma_w^2)$ and $b_i^{(l)} \sim \mathcal{N}(0, \sigma_b^2)$
- To normalize the forward and backward gradient signals w.r.t. layer width $D^{(l)}$, we can define an equivalent NN:

$$\mathbf{a}^{(l)} = \phi^{(l)}(\mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}),$$

where $W_{ij}^{(l)} = \frac{\sigma_w}{\sqrt{D^{(l-1)}}} \omega_{ij}^{(l)}$, $b_i^{(l)} = \sigma_b \beta_i^{(l)}$, and $\omega_{ij}^{(l)}, \beta_i^{(l)} \sim \mathcal{N}(0, 1)$

Distribution of \hat{y}

- Given an x , what is the distribution of its prediction \hat{y} ?

Distribution of \hat{y}

- Given an \mathbf{x} , what is the distribution of its prediction \hat{y} ?
- Recall that

$$\hat{y} = z^{(L)} = \mathbf{w}^{(L)\top} \mathbf{a}^{(L-1)} + b^{(L)} = \frac{\sigma_w}{\sqrt{D^{(l-1)}}} \sum_j \omega_j^{(L)} \phi(z_j^{(L-1)}) + \sigma_b \beta^{(L)}$$

- Since $\omega_j^{(l)}$'s and $\beta^{(l)}$ are Gaussian random variables with zero means, their sum \hat{y} is also a zero-mean Gaussian

Distribution of $\hat{\mathbf{y}}$

- Given an \mathbf{x} , what is the distribution of its prediction \hat{y} ?
- Recall that

$$\hat{y} = z^{(L)} = \mathbf{w}^{(L)\top} \mathbf{a}^{(L-1)} + b^{(L)} = \frac{\sigma_w}{\sqrt{D^{(l-1)}}} \sum_j \omega_j^{(L)} \phi(z_j^{(L-1)}) + \sigma_b \beta^{(L)}$$

- Since $\omega_j^{(l)}$'s and $\beta^{(l)}$ are Gaussian random variables with zero means, their sum \hat{y} is also a zero-mean Gaussian
- Now consider the predictions $\hat{\mathbf{y}} = [\hat{y}(\mathbf{x}^{(1)}), \dots, \hat{y}(\mathbf{x}^{(N)})]^\top \in \mathbb{R}^N$ for N points, we have

$$\begin{bmatrix} \hat{y}(\mathbf{x}^{(1)}) \\ \vdots \\ \hat{y}(\mathbf{x}^{(N)}) \end{bmatrix} = \frac{\sigma_w}{\sqrt{D^{(l-1)}}} \sum_j \omega_{j,i}^{(l)} \begin{bmatrix} \phi(z_j^{(l-1)}(\mathbf{x}^{(1)})) \\ \vdots \\ \phi(z_j^{(l-1)}(\mathbf{x}^{(N)})) \end{bmatrix} + \sigma_b \beta_i^{(l)} \mathbf{1}_N$$

- As $D^{(L-1)} \rightarrow \infty$, by multidimensional Central Limit Theorem, $\hat{\mathbf{y}}$ is a multivariate Gaussian with mean $\mathbf{0}_N$ and covariance Σ

Wide-and-Deep NN as a Gaussian Process

- The covariance Σ completely describes the behavior of our NN $\hat{y}(\cdot) = f(\cdot)$ over N points
- Furthermore, we will show that Σ can be describe by a **deterministic** kernel function $k^{(L)}(\cdot, \cdot)$ independent of a particular initialization such that

$$\Sigma = \begin{bmatrix} k^{(L)}(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k^{(L)}(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & \ddots & \vdots \\ k^{(L)}(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \dots & k^{(L)}(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix} \equiv \mathbf{K}_{N,N}^{(L)}$$

- This implies that the NN is in correspondent with a GP called **NN-GP**:

$$\begin{bmatrix} \hat{\mathbf{y}}_N \\ \hat{\mathbf{y}}_M \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0}_N \\ \mathbf{0}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N}^{(L)} & \mathbf{K}_{N,M}^{(L)} \\ \mathbf{K}_{M,N}^{(L)} & \mathbf{K}_{M,M}^{(L)} \end{bmatrix} \right)$$

Wide-and-Deep NN as a Gaussian Process

- The covariance Σ completely describes the behavior of our NN $\hat{y}(\cdot) = f(\cdot)$ over N points
- Furthermore, we will show that Σ can be describe by a **deterministic** kernel function $k^{(L)}(\cdot, \cdot)$ independent of a particular initialization such that

$$\Sigma = \begin{bmatrix} k^{(L)}(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \cdots & k^{(L)}(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & \ddots & \vdots \\ k^{(L)}(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \cdots & k^{(L)}(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix} \equiv \mathbf{K}_{N,N}^{(L)}$$

- This implies that the NN is in correspondent with a GP called **NN-GP**:

$$\begin{bmatrix} \hat{\mathbf{y}}_N \\ \hat{\mathbf{y}}_M \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0}_N \\ \mathbf{0}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N}^{(L)} & \mathbf{K}_{N,M}^{(L)} \\ \mathbf{K}_{M,N}^{(L)} & \mathbf{K}_{M,M}^{(L)} \end{bmatrix} \right)$$

- What's the $k^{(L)}(\cdot, \cdot)$?

Deriving $k^{(1)}(\cdot, \cdot)$

- We use induction to show that $z_i^{(1)}(\cdot), z_i^{(2)}(\cdot), \dots, z_i^{(L)}(\cdot) = \hat{y}(\cdot)$ are GPs, which are governed by kernels $k^{(1)}(\cdot, \cdot), \dots, k^{(L)}(\cdot, \cdot)$ independent with i , respectively

Deriving $k^{(1)}(\cdot, \cdot)$

- We use induction to show that $z_i^{(1)}(\cdot), z_i^{(2)}(\cdot), \dots, z_i^{(L)}(\cdot) = \hat{y}(\cdot)$ are GPs, which are governed by kernels $k^{(1)}(\cdot, \cdot), \dots, k^{(L)}(\cdot, \cdot)$ independent with i , respectively
- Consider $z_i^{(1)}(\mathbf{x}) = \frac{\sigma_w}{\sqrt{D^{(0)}}} \Sigma_j \omega_{j,i}^{(l)} \mathbf{x} + \sigma_b \beta_i^{(l)}$ a zero-mean Gaussian
- As $D^{(0)} \rightarrow \infty$, we have $[z_i^{(1)}(\mathbf{x}^{(1)}), \dots, z_i^{(1)}(\mathbf{x}^{(N)})]^\top \sim N(\mathbf{0}_N, \mathbf{K}_{N,N}^{(1)})$ by multidimensional Central Limit Theorem, where

$$\begin{aligned} k^{(1)}(\mathbf{x}, \mathbf{x}') &= \text{Cov}[z_i^{(1)}(\mathbf{x}), z_i^{(1)}(\mathbf{x}')] = \mathbb{E}_{\omega_{:,i}^{(l)}, \beta_i^{(l)}} [z_i^{(1)}(\mathbf{x}) z_i^{(1)}(\mathbf{x}')] \\ &= \frac{\sigma_w^2}{D^{(0)}} \mathbb{E} \left[\Sigma_{j,k} \omega_{j,i}^{(l)} \omega_{k,i}^{(l)} x_j x'_k \right] + \frac{\sigma_w \sigma_b}{\sqrt{D^{(0)}}} \mathbb{E} \left[\beta_i^{(l)} \Sigma_j \omega_{j,i}^{(l)} x_j \right] \\ &\quad + \frac{\sigma_w \sigma_b}{\sqrt{D^{(0)}}} \mathbb{E} \left[\beta_i^{(l)} \Sigma_j \omega_{j,i}^{(l)} x'_j \right] + \sigma_b^2 \mathbb{E} \left[\beta_i^{(l)} \beta_i^{(l)} \right] \\ &= \frac{\sigma_w^2}{D^{(0)}} \Sigma_j \mathbb{E} \left[\omega_{j,i}^{(l)} \omega_{j,i}^{(l)} \right] x_j x'_j + \sigma_b^2 \mathbb{E} \left[\beta_i^{(l)} \beta_i^{(l)} \right] \\ &= \frac{\sigma_w^2}{D^{(0)}} \mathbf{x}^\top \mathbf{x}' + \sigma_b^2, \end{aligned}$$

is independent with i

- Note that $z_i^{(1)}(\cdot)$ and $z_j^{(1)}(\cdot)$ are independent with each other, $\forall i \neq j$

Deriving $k^{(l)}(\cdot, \cdot)$

- Given that $D^{(0)} \rightarrow \infty, \dots, D^{(l-2)} \rightarrow \infty$ and
 - $[z_i^{(l-1)}(\mathbf{x}^{(1)}), \dots, z_i^{(l-1)}(\mathbf{x}^{(N)})]^\top \sim N(\mathbf{0}_N, \mathbf{K}_{N,N}^{(l-1)})$
 - $z_i^{(l-1)}(\cdot)$ and $z_j^{(l-1)}(\cdot)$ are independent with each other, $\forall i \neq j$

Deriving $k^{(l)}(\cdot, \cdot)$

- Given that $D^{(0)} \rightarrow \infty, \dots, D^{(l-2)} \rightarrow \infty$ and
 - $[z_i^{(l-1)}(\mathbf{x}^{(1)}), \dots, z_i^{(l-1)}(\mathbf{x}^{(N)})]^\top \sim N(\mathbf{0}_N, \mathbf{K}_{N,N}^{(l-1)})$
 - $z_i^{(l-1)}(\cdot)$ and $z_j^{(l-1)}(\cdot)$ are independent with each other, $\forall i \neq j$
- Consider $z_i^{(l)}(\mathbf{x}) = \frac{\sigma_w}{\sqrt{D^{(l-1)}}} \sum_j \omega_{j,i}^{(l)} \phi(z_j^{(l-1)}(\mathbf{x})) + \sigma_b \beta_i^{(l)}$ a zero-mean Gaussian
- As $D^{(l-1)} \rightarrow \infty$, we have $[z_i^{(l)}(\mathbf{x}^{(1)}), \dots, z_i^{(l)}(\mathbf{x}^{(N)})]^\top \sim N(\mathbf{0}_N, \mathbf{K}_{N,N}^{(l)})$ by multidimensional Central Limit Theorem, where

$$\begin{aligned}
 k^{(l)}(\mathbf{x}, \mathbf{x}') &= \text{Cov}[z_i^{(l)}(\mathbf{x}), z_i^{(l)}(\mathbf{x}')] = \mathbb{E}_{\omega_{:,i}^{(l)}, \beta_i^{(l)}, z^{(l-1)}(\mathbf{x})} [z_i^{(l)}(\mathbf{x}) z_i^{(l)}(\mathbf{x}')] \\
 &= \frac{\sigma_w^2}{D^{(l-1)}} \mathbb{E} \left[\sum_{j,k} \omega_{j,i}^{(l)} \omega_{k,i}^{(l)} \phi(z_j^{(l-1)}(\mathbf{x})) \phi(z_k^{(l-1)}(\mathbf{x}')) \right] + \sigma_b^2 \mathbb{E} \left[\beta_i^{(l)} \beta_i^{(l)} \right] \\
 &\quad + \frac{\sigma_w \sigma_b}{\sqrt{D^{(l-1)}}} \left(\mathbb{E} \left[\beta_i^{(l)} \sum_j \omega_{j,i}^{(l)} \phi(z_j^{(l-1)}(\mathbf{x})) \right] + \mathbb{E} \left[\beta_i^{(l)} \sum_j \omega_{j,i}^{(l)} \phi(z_j^{(l-1)}(\mathbf{x}')) \right] \right) \\
 &= \frac{\sigma_w^2}{D^{(l-1)}} \sum_j \mathbb{E} \left[\omega_{j,i}^{(l)} \omega_{j,i}^{(l)} \right] \mathbb{E} \left[\phi(z_j^{(l-1)}(\mathbf{x})) \phi(z_j^{(l-1)}(\mathbf{x}')) \right] + \sigma_b^2 \mathbb{E} \left[\beta_i^{(l)} \beta_i^{(l)} \right] \\
 &= \sigma_w^2 \mathbb{E}_{(z_i^{(l-1)}(\mathbf{x}), z_i^{(l-1)}(\mathbf{x}')) \sim \mathcal{N}(\mathbf{0}_2, \mathbf{K}_{2,2}^{(l-1)})} \left[\phi(z_i^{(l-1)}(\mathbf{x})) \phi(z_i^{(l-1)}(\mathbf{x}')) \right] + \sigma_b^2,
 \end{aligned}$$

where

$$\mathbf{K}_{2,2}^{(l-1)} = \begin{bmatrix} k^{(l-1)}(\mathbf{x}, \mathbf{x}) & k^{(l-1)}(\mathbf{x}, \mathbf{x}') \\ k^{(l-1)}(\mathbf{x}, \mathbf{x}') & k^{(l-1)}(\mathbf{x}', \mathbf{x}') \end{bmatrix}$$

Evaluating $\mathbf{K}^{(l)}$

- For certain activation functions $\phi(\cdot)$, such as tanh and ReLU, $k^{(l)}(\mathbf{x}, \mathbf{x}')$ has a closed form [9]
- For other $\phi(\cdot)$'s, Markov Chain Monte Carlo (MCMC) sampling is required to devaluate $k^{(l)}(\mathbf{x}, \mathbf{x}')$

Outline

① When ML Meets Big Data

② Advantages of Deep Learning

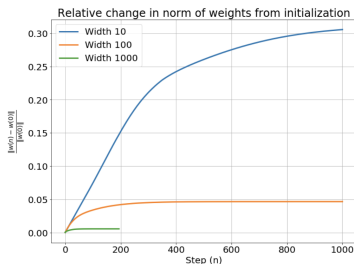
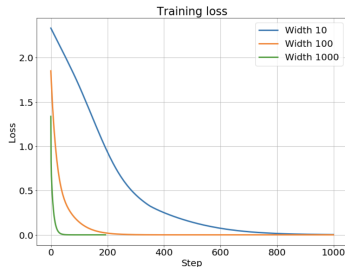
- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Weight Dynamics

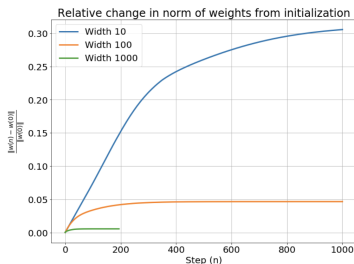
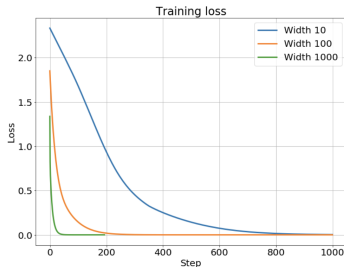
- Observation: the weights of a wide NN do not change much during gradient descent



- Why?

Weight Dynamics

- Observation: the weights of a wide NN do not change much during gradient descent



- Why? A small change in a large number of neurons is enough to significantly change the output
- This allows us to approximate an NN $f(\cdot; \theta)$ *w.r.t. weights* using the first-order Taylor expansion

Linearization of $f(\cdot; \theta)$

- Let $\theta^{(t)}$ be the parameters of the NN at the t -th step of gradient descent
 - $\hat{\mathbf{y}}^{(t)} = [f(\mathbf{x}^{(1)}; \theta^{(t)}), \dots, f(\mathbf{x}^{(N)}; \theta^{(t)})]^\top$ be the predictions over training points
- Since $\theta^{(t)}$ is close to $\theta^{(0)}$ at any time t , we can approximate $f(\cdot; \theta^{(t)})$ using the first-order Taylor expansion **w.r.t. $\theta^{(t)}$** around $\theta^{(0)}$:

$$f(\mathbf{x}, \theta^{(t)}) \approx \bar{f}(\mathbf{x}, \theta^{(t)}) = f(\mathbf{x}, \theta^{(0)}) + \nabla_{\theta} f(\mathbf{x}, \theta^{(0)})^\top (\theta^{(t)} - \theta^{(0)})$$

- \bar{f} is still **non-linear in terms of \mathbf{x}**
- Let $\bar{\mathbf{y}}^{(t)} = [\bar{f}(\mathbf{x}^{(1)}; \theta^{(t)}), \dots, \bar{f}(\mathbf{x}^{(N)}; \theta^{(t)})]^\top$ be the predictions of \bar{f} at time t

Weight and Prediction Dynamics

$$f(\mathbf{x}, \boldsymbol{\theta}^{(t)}) \approx \bar{f}(\mathbf{x}, \boldsymbol{\theta}^{(t)}) = f(\mathbf{x}, \boldsymbol{\theta}^{(0)}) + \nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}^{(0)})^\top (\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^{(0)})$$

- Gradient descent with learning rate η makes the following changes:

$$\begin{aligned}\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)} &\approx -\eta \nabla_{\boldsymbol{\theta}} C(\bar{\mathbf{y}}^{(t)}, \mathbf{y}) \\ &= -\eta \nabla_{\boldsymbol{\theta}} \bar{\mathbf{y}}^{(t)} \nabla_{\bar{\mathbf{y}}^{(t)}} C(\bar{\mathbf{y}}^{(t)}, \mathbf{y}) \\ &= -\eta \nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)} \nabla_{\bar{\mathbf{y}}^{(t)}} C(\bar{\mathbf{y}}^{(t)}, \mathbf{y})\end{aligned}$$

and

$$\begin{aligned}\bar{\mathbf{y}}^{(t+1)} - \bar{\mathbf{y}}^{(t)} &= \nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)\top} (\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}) \\ &\approx -\eta \underbrace{\nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)\top}}_{N \times D} \underbrace{\nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)}}_{D \times N} \nabla_{\bar{\mathbf{y}}^{(t)}} C(\bar{\mathbf{y}}^{(t)}, \mathbf{y}),\end{aligned}$$

where $\mathbf{T}_{N,N}^{(0)} \equiv \nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)\top} \nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)} \in \mathbb{R}^{N \times N}$ is called the **Neural Tangent Kernel (NTK)** matrix

Prediction Dynamics in Regression

- In regression where $C(\bar{\mathbf{y}}^{(0)}, \mathbf{y}) = \frac{1}{2} \|\bar{\mathbf{y}}^{(0)} - \mathbf{y}\|^2$, we have

$$\bar{\mathbf{y}}^{(t+1)} - \bar{\mathbf{y}}^{(t)} \approx -\eta \mathbf{T}_{N,N}^{(0)} \nabla_{\bar{\mathbf{y}}^{(t)}} C(\bar{\mathbf{y}}^{(t)}, \mathbf{y}) = -\eta \mathbf{T}_{N,N}^{(0)} (\bar{\mathbf{y}}^{(t)} - \mathbf{y})$$

- With a sufficiently small learning rate η , we can think t as continuous time and each GD step as Δt , where

$$\lim_{\Delta t \rightarrow 0} \frac{\bar{\mathbf{y}}^{(t+\Delta t)} - \bar{\mathbf{y}}^{(t)}}{\Delta t} = \frac{\partial \bar{\mathbf{y}}^{(t)}}{\partial t} \approx -\eta \mathbf{T}_{N,N}^{(0)} (\bar{\mathbf{y}}^{(t)} - \mathbf{y})$$

- Letting $\mathbf{u}^{(t)} = \bar{\mathbf{y}}^{(t)} - \mathbf{y}$, we have an ordinary differential equation:

$$\begin{aligned} \frac{\partial \bar{\mathbf{y}}^{(t)}}{\partial t} &\approx -\eta \mathbf{T}_{N,N}^{(0)} (\bar{\mathbf{y}}^{(t)} - \mathbf{y}) \\ \Rightarrow \frac{\partial \mathbf{u}^{(t)}}{\partial t} &\approx -\eta \mathbf{T}_{N,N}^{(0)} \mathbf{u}^{(t)} \end{aligned}$$

Prediction Dynamics in Regression

- In regression where $C(\bar{\mathbf{y}}^{(0)}, \mathbf{y}) = \frac{1}{2} \|\bar{\mathbf{y}}^{(0)} - \mathbf{y}\|^2$, we have

$$\bar{\mathbf{y}}^{(t+1)} - \bar{\mathbf{y}}^{(t)} \approx -\eta \mathbf{T}_{N,N}^{(0)} \nabla_{\bar{\mathbf{y}}^{(t)}} C(\bar{\mathbf{y}}^{(t)}, \mathbf{y}) = -\eta \mathbf{T}_{N,N}^{(0)} (\bar{\mathbf{y}}^{(t)} - \mathbf{y})$$

- With a sufficiently small learning rate η , we can think t as continuous time and each GD step as Δt , where

$$\lim_{\Delta t \rightarrow 0} \frac{\bar{\mathbf{y}}^{(t+\Delta t)} - \bar{\mathbf{y}}^{(t)}}{\Delta t} = \frac{\partial \bar{\mathbf{y}}^{(t)}}{\partial t} \approx -\eta \mathbf{T}_{N,N}^{(0)} (\bar{\mathbf{y}}^{(t)} - \mathbf{y})$$

- Letting $\mathbf{u}^{(t)} = \bar{\mathbf{y}}^{(t)} - \mathbf{y}$, we have an ordinary differential equation:

$$\begin{aligned} \frac{\partial \bar{\mathbf{y}}^{(t)}}{\partial t} &\approx -\eta \mathbf{T}_{N,N}^{(0)} (\bar{\mathbf{y}}^{(t)} - \mathbf{y}) \\ \Rightarrow \frac{\partial \mathbf{u}^{(t)}}{\partial t} &\approx -\eta \mathbf{T}_{N,N}^{(0)} \mathbf{u}^{(t)} \end{aligned}$$

- Therefore, $\mathbf{u}^{(t)} = e^{-\eta \mathbf{T}_{N,N}^{(0)} t} \mathbf{u}^{(0)}$
 - Recall that $e^{\mathbf{A}t} = \frac{1}{0!} \mathbf{I} + \frac{t}{1!} \mathbf{A} + \frac{t^2}{2!} \mathbf{A}^2 + \dots$ for a symmetric \mathbf{A}
 - So, $\frac{\partial e^{\mathbf{A}t}}{\partial t} = \frac{1}{0!} \mathbf{A} + \frac{t}{1!} \mathbf{A}^2 + \dots = (\frac{1}{0!} \mathbf{I} + \frac{t}{1!} \mathbf{A} + \dots) \mathbf{A} = \mathbf{A} e^{\mathbf{A}t}$
- This implies that

$$\bar{\mathbf{y}}^{(t)} = e^{-\eta \mathbf{T}_{N,N}^{(0)} t} \bar{\mathbf{y}}^{(0)} + (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N}^{(0)} t}) \mathbf{y} = e^{-\eta \mathbf{T}_{N,N}^{(0)} t} \hat{\mathbf{y}}^{(0)} + (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N}^{(0)} t}) \mathbf{y}$$

Weight Dynamics in Regression

- By definition of $\bar{\mathbf{y}}^{(t)}$, we also have $\bar{\mathbf{y}}^{(t)} = \hat{\mathbf{y}}^{(0)} + \nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)\top} (\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^{(0)})$
- Solving $\boldsymbol{\theta}^{(t)}$ in

$$e^{-\eta \mathbf{T}_{N,N}^{(0)} t} \hat{\mathbf{y}}^{(0)} + (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N}^{(0)} t}) \mathbf{y} = \hat{\mathbf{y}}^{(0)} + \nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)\top} (\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^{(0)}),$$

we have

$$\boldsymbol{\theta}^{(t)} = \boldsymbol{\theta}^{(0)} - \nabla_{\boldsymbol{\theta}} \hat{\mathbf{y}}^{(0)\top} \mathbf{T}_{N,N}^{(0)-1} (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N}^{(0)} t}) (\hat{\mathbf{y}}^{(0)} - \mathbf{y})$$

Predictions of Trained NN

- Substituting $\theta^{(t)}$ in $\bar{y}^{(t)} = \hat{y}^{(0)} + \nabla_{\theta} \hat{y}^{(0)\top} (\theta^{(t)} - \theta^{(0)})$, we have that:
- For an arbitrary (training or test) point \mathbf{x}' , the prediction of trained NN is

$$f(\mathbf{x}', \theta^{(t)}) \approx \bar{f}(\mathbf{x}'; \theta^{(t)}) = \mathbf{p}^\top \begin{bmatrix} \hat{\mathbf{y}}^{(0)} \\ \hat{\mathbf{y}}'^{(0)} \end{bmatrix} + q,$$

where

$$\begin{aligned} \mathbf{p} &= [-\mathbf{T}_{1',N}^{(0)} \mathbf{T}_{N,N}^{(0)-1} (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N}^{(0)}}), 1]^\top \in \mathbb{R}^{N+1}, \\ q &= \mathbf{T}_{1',N}^{(0)} \mathbf{T}_{N,N}^{(0)-1} (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N}^{(0)}}) \mathbf{y} \end{aligned}$$

- $\mathbf{T}_{N,N}^{(0)} = \nabla_{\theta} \hat{\mathbf{y}}^{(0)\top} \nabla_{\theta} \hat{\mathbf{y}}^{(0)} \in \mathbb{R}^{N \times N}$ is the NTK matrix for \mathbf{X}_N
- $\mathbf{T}_{1',N}^{(0)} = \nabla_{\theta} \hat{\mathbf{y}}'^{(0)\top} \nabla_{\theta} \hat{\mathbf{y}}^{(0)} \in \mathbb{R}^{1 \times N}$ is the NTK matrix between \mathbf{x}' and \mathbf{X}_N

Predictions of Trained NN

- Substituting $\theta^{(t)}$ in $\bar{y}^{(t)} = \hat{y}^{(0)} + \nabla_{\theta} \hat{y}^{(0)\top} (\theta^{(t)} - \theta^{(0)})$, we have that:
- For an arbitrary (training or test) point \mathbf{x}' , the prediction of trained NN is

$$f(\mathbf{x}', \theta^{(t)}) \approx \bar{f}(\mathbf{x}'; \theta^{(t)}) = \mathbf{p}^\top \begin{bmatrix} \hat{\mathbf{y}}^{(0)} \\ \hat{\mathbf{y}}'^{(0)} \end{bmatrix} + q,$$

where

$$\mathbf{p} = [-\mathbf{T}_{1',N}^{(0)} \mathbf{T}_{N,N}^{(0)-1} (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N}^{(0)}}), 1]^\top \in \mathbb{R}^{N+1},$$
$$q = \mathbf{T}_{1',N}^{(0)} \mathbf{T}_{N,N}^{(0)-1} (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N}^{(0)}}) \mathbf{y}$$

- $\mathbf{T}_{N,N}^{(0)} = \nabla_{\theta} \hat{\mathbf{y}}^{(0)\top} \nabla_{\theta} \hat{\mathbf{y}}^{(0)} \in \mathbb{R}^{N \times N}$ is the NTK matrix for \mathbf{X}_N
- $\mathbf{T}_{1',N}^{(0)} = \nabla_{\theta} \hat{\mathbf{y}}'^{(0)\top} \nabla_{\theta} \hat{\mathbf{y}}^{(0)} \in \mathbb{R}^{1 \times N}$ is the NTK matrix between \mathbf{x}' and \mathbf{X}_N
- **No actual training** needed!

Gradient Descent as an Affine Transformation

Theorem (NTK in infinite width)

*As the NN's width goes to infinity, $\mathbf{T}_{N,N}^{(0)}$ and $\mathbf{T}_{1',N}^{(0)}$ converges to $\mathbf{T}_{N,N}$ and $\mathbf{T}_{1',N}$, which can be described by a **deterministic** kernel function independent of a particular initialization [8, 10].*

- $\mathbf{T}_{N,N}$ and $\mathbf{T}_{1',N}$ have closed forms for certain activation functions $\phi(\cdot)$'s, including erf and ReLU

Gradient Descent as an Affine Transformation

Theorem (NTK in infinite width)

As the NN's width goes to infinity, $\mathbf{T}_{N,N}^{(0)}$ and $\mathbf{T}_{1',N}^{(0)}$ converges to $\mathbf{T}_{N,N}$ and $\mathbf{T}_{1',N}$, which can be described by a **deterministic** kernel function independent of a particular initialization [8, 10].

- $\mathbf{T}_{N,N}$ and $\mathbf{T}_{1',N}$ have closed forms for certain activation functions $\phi(\cdot)$'s, including erf and ReLU
- $f(\mathbf{x}', \boldsymbol{\theta}^{(t)}) \approx \mathbf{p}^\top \begin{bmatrix} \hat{\mathbf{y}}^{(0)} \\ \hat{\mathbf{y}}'^{(0)} \end{bmatrix} + q$ is an **affine transformation** of a random vector $\begin{bmatrix} \hat{\mathbf{y}}^{(0)} \\ \hat{\mathbf{y}}'^{(0)} \end{bmatrix}$

Outline

① When ML Meets Big Data

② Advantages of Deep Learning

- Representation Learning
- Exponential Gain of Expressiveness
- Memory and GPU Friendliness
- Online & Transfer Learning

③ Learning Theory Revisited

- Generalizability and Over-Parametrization
- Wide-and-Deep NN is a Gaussian Process before Training*
- Gradient Descent is an Affine Transformation
- Wide-and-Deep NN is a Gaussian Process after Training

Wide-and-Deep NN as a Gaussian Process I

- As $D \rightarrow \infty$, randomly initialized NN has a corresponding NN-GP:

$$\begin{bmatrix} \hat{\mathbf{y}}_N \\ \hat{\mathbf{y}}_M \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0}_N \\ \mathbf{0}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N}^{(L)} & \mathbf{K}_{N,M}^{(L)} \\ \mathbf{K}_{M,N}^{(L)} & \mathbf{K}_{M,M}^{(L)} \end{bmatrix} \right)$$

- As $D \rightarrow \infty$, GD-based training is an affine transformation:

$$f(\mathbf{x}', \boldsymbol{\theta}^{(t)}) \approx \bar{f}(\mathbf{x}'; \boldsymbol{\theta}^{(t)}) = \mathbf{p}^\top \begin{bmatrix} \hat{\mathbf{y}}^{(0)} \\ \hat{\mathbf{y}}'^{(0)} \end{bmatrix} + q$$

where

- $\mathbf{p} = [-\mathbf{T}_{1',N} \mathbf{T}_{N,N}^{-1} (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N} t}), 1]^\top \in \mathbb{R}^{N+1}$
- $q = \mathbf{T}_{1',N} \mathbf{T}_{N,N}^{-1} (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N} t}) \mathbf{y}$
- $\mathbf{T}_{N,N}$ and $\mathbf{T}_{1',N}$ the NTK matrices

Wide-and-Deep NN as a Gaussian Process II

- Therefore, as $D \rightarrow \infty$, the trained NN is still in correspondent with a GP, called **NTK-GP**, whose predictions for M test points are

$$\begin{bmatrix} \hat{\mathbf{y}}_N \\ \hat{\mathbf{y}}_M \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{A}\mathbf{y} \\ \mathbf{B}\mathbf{y} \end{bmatrix}, \mathbf{C}^\top \begin{bmatrix} \mathbf{K}_{N,N}^{(L)} & \mathbf{K}_{N,M}^{(L)} \\ \mathbf{K}_{M,N}^{(L)} & \mathbf{K}_{M,M}^{(L)} \end{bmatrix} \mathbf{C} \right),$$

where

- $\mathbf{A} = (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N} t}) \in \mathbb{R}^{N \times N}$
- $\mathbf{B} = \mathbf{T}_{M,N} \mathbf{T}_{N,N}^{-1} (\mathbf{I} - e^{-\eta \mathbf{T}_{N,N} t}) \in \mathbb{R}^{M \times N}$
- $\mathbf{C} = \begin{bmatrix} -\mathbf{A} & \mathbf{1}_{N,M} \\ -\mathbf{B} & \mathbf{1}_{M,M} \end{bmatrix} \in \mathbb{R}^{(N+M) \times (N+M)}$

Mean Predictions of NTK-GP

- Prior (unconditioned) mean predictions for training set:

$$\hat{\mathbf{y}}_N = \mathbf{A}\mathbf{y} = (\mathbf{I} - e^{-\eta T_{N,N}t})\mathbf{y}$$

- As $t \rightarrow \infty$, **the $\hat{\mathbf{y}}_N$ always approaches true labels \mathbf{y}**
- This explains why the SGD-based training of large NNs seldom encounters significant obstacles such as local minima [7]

Mean Predictions of NTK-GP

- Prior (unconditioned) mean predictions for training set:

$$\hat{\mathbf{y}}_N = \mathbf{A}\mathbf{y} = (\mathbf{I} - e^{-\eta T_{N,N}t})\mathbf{y}$$

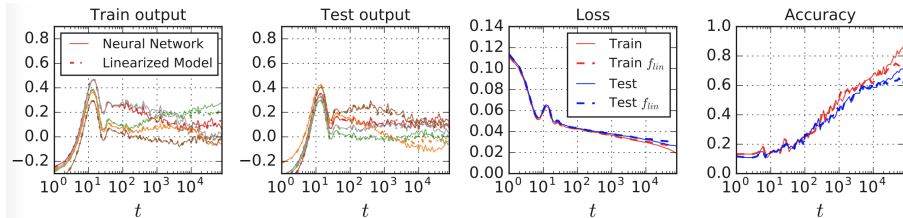
- As $t \rightarrow \infty$, **the $\hat{\mathbf{y}}_N$ always approaches true labels \mathbf{y}**
- This explains why the SGD-based training of large NNs seldom encounters significant obstacles such as local minima [7]
- Prior mean predictions for test set:

$$\hat{\mathbf{y}}_M = \mathbf{B}\mathbf{y} = \mathbf{T}_{M,N}\mathbf{T}_{N,N}^{-1}(\mathbf{I} - e^{-\eta T_{N,N}t})\mathbf{y}$$

- As $t \rightarrow \infty$, we have $\hat{\mathbf{y}}_M = \mathbf{T}_{M,N}\mathbf{T}_{N,N}^{-1}\mathbf{y}$
- **Weight hyperparameters are important** because they determines $\mathbf{T}_{M,N}\mathbf{T}_{N,N}^{-1}$

Analytic vs. Real Predictions

- Wide residual network [18] trained by SGD with momentum on MSE loss on CIFAR-10
 - First two panes shows the output dynamics for a randomly selected subset of train and test points



Remarks

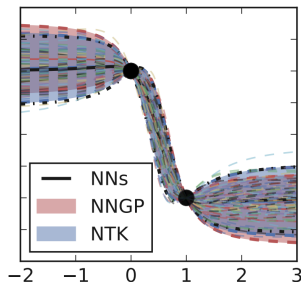
- Wide-and-deep NNs can be approximated by a class of GPs
 - Either before, during, or after training
- Therefore, complexity of wide-and-deep NNs grows with N , not $|\theta|$

Remarks

- Wide-and-deep NNs can be approximated by a class of GPs
 - Either before, during, or after training
- Therefore, complexity of wide-and-deep NNs grows with N , not $|\theta|$
- Applicable to other architectures incl. CNN [2, 15], RNN [16, 1], and any architecture [17]

Limitations

- Approximation holds only when the NNs have:
 - Infinite width
 - Small learning rate: $\eta \leq O(1/\lambda_{\max})$ where λ_{\max} is the max eigenvalue of $\mathbf{T}_{N,N}$
 - Proper initialization (to be discussed next)
- The prior NTK-GP inference is **inconsistent** with the Bayesian inference of NN-GP
 - SGP introduces bias



Reference I

- [1] Sina Alemohammad, Zichao Wang, Randall Balestriero, and Richard Baraniuk.
The recurrent neural tangent kernel.
arXiv preprint arXiv:2006.10246, 2020.
- [2] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang.
On exact computation with an infinitely wide neural net.
In *Advances in Neural Information Processing Systems*, pages 8141–8150, 2019.
- [3] Léon Bottou.
Large-scale machine learning with stochastic gradient descent.
In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010.

Reference II

- [4] Olivier Bousquet.
Concentration inequalities and empirical processes theory applied to the analysis of learning algorithms.
Ph.D. thesis, Ecole Polytechnique, Palaiseau, France, 2002.
- [5] Adam Coates, Brody Huval, Tao Wang, David Wu, Bryan Catanzaro, and Ng Andrew.
Deep learning with cots hpc systems.
In *Proceedings of The 30th International Conference on Machine Learning*, pages 1337–1345, 2013.
- [6] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin.
Liblinear: A library for large linear classification.
Journal of machine learning research, 9(Aug):1871–1874, 2008.

Reference III

- [7] Ian J Goodfellow, Oriol Vinyals, and Andrew M Saxe.
Qualitatively characterizing neural network optimization problems.
arXiv preprint arXiv:1412.6544, 2014.
- [8] Arthur Jacot, Franck Gabriel, and Clément Hongler.
Neural tangent kernel: Convergence and generalization in neural networks.
In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [9] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein.
Deep neural networks as gaussian processes.
arXiv preprint arXiv:1711.00165, 2017.

Reference IV

- [10] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington.
Wide neural networks of any depth evolve as linear models under gradient descent.
In *Advances in neural information processing systems*, pages 8572–8583, 2019.
- [11] Pascal Massart.
Some applications of concentration inequalities to statistics.
In *Annales de la Faculté des sciences de Toulouse: Mathématiques*, volume 9, pages 245–303, 2000.
- [12] Guido F Montufar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio.
On the number of linear regions of deep neural networks.
In *Advances in neural information processing systems*, pages 2924–2932, 2014.

Reference V

- [13] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich.
Going deeper with convolutions.
In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 1–9, 2015.
- [14] Vladimir N Vapnik and A Ya Chervonenkis.
On the uniform convergence of relative frequencies of events to their probabilities.
In Measures of Complexity, pages 11–30. Springer, 2015.
- [15] Greg Yang.
Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation.
arXiv preprint arXiv:1902.04760, 2019.

Reference VI

- [16] Greg Yang.
Wide feedforward or recurrent neural networks of any architecture are gaussian processes.
In Advances in Neural Information Processing Systems, pages 9951–9960, 2019.
- [17] Greg Yang.
Tensor programs ii: Neural tangent kernel for any architecture.
arXiv preprint arXiv:2006.14548, 2020.
- [18] Sergey Zagoruyko and Nikos Komodakis.
Wide residual networks.
arXiv preprint arXiv:1605.07146, 2016.
- [19] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals.
Understanding deep learning requires rethinking generalization.
arXiv preprint arXiv:1611.03530, 2016.

Reference VII