

Introduction to Neural Networks and Deep Learning

Regularization

Andres Mendez-Vazquez

August 24, 2020

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Outline

1 Bias-Variance Dilemma

● Introduction

- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Introduction

What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

Introduction

What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

Going back to the probability models

We might think in the machine to be learned as a function $g(x|\mathcal{D})$...

- Something as curve fitting...

Introduction

What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

Going back to the probability models

We might think in the machine to be learned as a function $g(x|\mathcal{D})$...

- Something as curve fitting...

Uniform data set

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\} \quad (1)$$

Remark: Where the $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$!!!

Introduction

What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

Going back to the probability models

We might think in the machine to be learned as a function $g(\mathbf{x}|\mathcal{D})$

- Something as curve fitting...

Uniform dataset

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\} \quad (1)$$

Remark: Where the $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$!!!

Introduction

What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

Going back to the probability models

We might think in the machine to be learned as a function $g(\mathbf{x}|\mathcal{D})$

- Something as curve fitting...

Definition: Dataset

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\} \quad (1)$$

Remark: Where the $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$!!!

Introduction

What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

Going back to the probability models

We might think in the machine to be learned as a function $g(\mathbf{x}|\mathcal{D})$

- Something as curve fitting...

Under a data set

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\} \quad (1)$$

Remark: Where the $\mathbf{x}_i \sim p(\mathbf{x}|\theta)$!!!

Introduction

What did we see until now?

The design of learning machines from two main points:

- Statistical Point of View
- Linear Algebra and Optimization Point of View

Going back to the probability models

We might think in the machine to be learned as a function $g(\mathbf{x}|\mathcal{D})$

- Something as curve fitting...

Under a data set

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\} \quad (1)$$

Remark: Where the $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$!!!

Thus, we have that

Two main functions

- A function $g(x|\mathcal{D})$ obtained using some algorithm!!!
- $E[y|x]$ the optimal regression...

Thus, we have that

Two main functions

- A function $g(\mathbf{x}|\mathcal{D})$ obtained using some algorithm!!!
- $E[y|\mathbf{x}]$ the optimal regression...

important

The key factor here is the dependence of the approximation on \mathcal{D} .

Thus, we have that

Two main functions

- A function $g(\mathbf{x}|\mathcal{D})$ obtained using some algorithm!!!
- $E[y|\mathbf{x}]$ the optimal regression...

Important

The key factor here is the dependence of the approximation on \mathcal{D} .

The approximation may be very good for a specific training data set but very bad for another.

- This is the reason of studying fusion of information at decision level...

Thus, we have that

Two main functions

- A function $g(x|\mathcal{D})$ obtained using some algorithm!!!
- $E[y|x]$ the optimal regression...

Important

The key factor here is the dependence of the approximation on \mathcal{D} .

Why?

The approximation may be very good for a specific training data set but very bad for another.

• This is the reason of studying fusion of information at decision level...

Thus, we have that

Two main functions

- A function $g(x|\mathcal{D})$ obtained using some algorithm!!!
- $E[y|x]$ the optimal regression...

Important

The key factor here is the dependence of the approximation on \mathcal{D} .

Why?

The approximation may be very good for a specific training data set but very bad for another.

- This is the reason of studying fusion of information at decision level...

Outline

1 Bias-Variance Dilemma

- Introduction
- **Measuring the difference between optimal and learned**
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

How do we measure the difference? [1]

We have that

$$\text{Var}(X) = E((X - \mu)^2)$$

How do we measure the difference? [1]

We have that

$$\text{Var}(X) = E((X - \mu)^2)$$

We can do that for our data

$$\text{Var}_{\mathcal{D}}(g(\mathbf{x}|\mathcal{D})) = E_{\mathcal{D}}\left(\left(g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}]\right)^2\right)$$

How do we measure the difference? [1]

We have that

$$\text{Var}(X) = E((X - \mu)^2)$$

We can do that for our data

$$\text{Var}_{\mathcal{D}}(g(\mathbf{x}|\mathcal{D})) = E_{\mathcal{D}}\left(\left(g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}]\right)^2\right)$$

Now, if we add and subtract

$$E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] \tag{2}$$

Remark: The expected output of the machine $g(\mathbf{x}|\mathcal{D})$

How do we measure the difference? [1]

We have that

$$\text{Var}(X) = E((X - \mu)^2)$$

We can do that for our data

$$\text{Var}_{\mathcal{D}}(g(\mathbf{x}|\mathcal{D})) = E_{\mathcal{D}}\left(\left(g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}]\right)^2\right)$$

Now, if we add and subtract

$$E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] \tag{2}$$

Remark: The expected output of the machine $g(\mathbf{x}|\mathcal{D})$

Thus, we have that

Or Original variance

$$\begin{aligned} \text{Var}_{\mathcal{D}}(g(\mathbf{x}|\mathcal{D})) &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}])^2) \\ &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] + E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2) \\ &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})])^2 + \dots \\ &\quad \dots 2((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})]) (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}]) + \dots \\ &\quad \dots (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2) \end{aligned}$$

Thus, we have that

Or Original variance

$$\begin{aligned} \text{Var}_{\mathcal{D}}(g(\mathbf{x}|\mathcal{D})) &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}])^2) \\ &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] + E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2) \\ &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})])^2 + \dots \\ &\quad \dots + 2((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})]) (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}]) + \dots \\ &\quad \dots + (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2) \end{aligned}$$

Finally

$$E_{\mathcal{D}}(((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})])) (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])) =? \quad (3)$$

Thus, we have that

Or Original variance

$$\begin{aligned} \text{Var}_{\mathcal{D}}(g(\mathbf{x}|\mathcal{D})) &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}])^2) \\ &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] + E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2) \\ &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})])^2 + \dots \\ &\quad \dots 2((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})]))(E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}]) + \dots \\ &\quad \dots (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2) \end{aligned}$$

Finally

$$E_{\mathcal{D}}(((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})]))(E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])) =? \quad (3)$$

Thus, we have that

Or Original variance

$$\begin{aligned} \text{Var}_{\mathcal{D}}(g(\mathbf{x}|\mathcal{D})) &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}])^2) \\ &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] + E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2) \\ &= E_{\mathcal{D}}((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})])^2 + \dots \\ &\quad \dots 2((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})]) (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}]) + \dots \\ &\quad \dots (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2) \end{aligned}$$

Finally

$$E_{\mathcal{D}}(((g(\mathbf{x}|\mathcal{D}) - E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})])) (E_{\mathcal{D}}[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])) =? \quad (3)$$

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- **The Bias-Variance**
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

We have the Bias-Variance

Our Final Equation

$$E_D \left((g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}])^2 \right) = \underbrace{E_D \left((g(\mathbf{x}|\mathcal{D}) - E_D[g(\mathbf{x}|\mathcal{D})])^2 \right)}_{\text{VARIANCE}} + \underbrace{(E_D[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2}_{\text{BIAS}}$$

We have the Bias-Variance

Our Final Equation

$$E_D \left((g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}])^2 \right) = \underbrace{E_D \left((g(\mathbf{x}|\mathcal{D}) - E_D[g(\mathbf{x}|\mathcal{D})])^2 \right)}_{\text{VARIANCE}} + \underbrace{(E_D[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2}_{\text{BIAS}}$$

Where the variance

It represents the measure of the error between our machine $g(\mathbf{x}|\mathcal{D})$ and the expected output of the machine under $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$.

We have the Bias-Variance

Our Final Equation

$$E_D \left((g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}])^2 \right) = \underbrace{E_D \left((g(\mathbf{x}|\mathcal{D}) - E_D[g(\mathbf{x}|\mathcal{D})])^2 \right)}_{\text{VARIANCE}} + \underbrace{(E_D[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2}_{\text{BIAS}}$$

Where the variance

It represents the measure of the error between our machine $g(\mathbf{x}|\mathcal{D})$ and the expected output of the machine under $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$.

Where the bias

It represents the quadratic error between the expected output of the machine under $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$ and the expected output of the optimal regression.

We have the Bias-Variance

Our Final Equation

$$E_D \left((g(\mathbf{x}|\mathcal{D}) - E[y|\mathbf{x}])^2 \right) = \underbrace{E_D \left((g(\mathbf{x}|\mathcal{D}) - E_D[g(\mathbf{x}|\mathcal{D})])^2 \right)}_{\text{VARIANCE}} + \underbrace{(E_D[g(\mathbf{x}|\mathcal{D})] - E[y|\mathbf{x}])^2}_{\text{BIAS}}$$

Where the variance

It represents the measure of the error between our machine $g(\mathbf{x}|\mathcal{D})$ and the expected output of the machine under $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$.

Where the bias

It represents the quadratic error between the expected output of the machine under $\mathbf{x}_i \sim p(\mathbf{x}|\Theta)$ and the expected output of the optimal regression.

Remarks

We have then

Even if the estimator is unbiased, it can still result in a large mean square error due to a large variance term.

Remarks

We have then

Even if the estimator is unbiased, it can still result in a large mean square error due to a large variance term.

The situation is more dire in a finite set of data \mathcal{D}

We have then a trade-off:

- Increasing the bias decreases the variance and vice versa.
- This is known as the bias-variance dilemma.

Remarks

We have then

Even if the estimator is unbiased, it can still result in a large mean square error due to a large variance term.

The situation is more dire in a finite set of data \mathcal{D}

We have then a trade-off:

- 1 Increasing the bias decreases the variance and vice versa.

⦿ This is known as the bias-variance dilemma.

Remarks

We have then

Even if the estimator is unbiased, it can still result in a large mean square error due to a large variance term.

The situation is more dire in a finite set of data \mathcal{D}

We have then a trade-off:

- 1 Increasing the bias decreases the variance and vice versa.
- 2 This is known as the **bias–variance dilemma**.

Similar to...

Curve Fitting

If, for example, the adopted model is complex (many parameters involved) with respect to the number N , the model will fit the idiosyncrasies of the specific data set.

Similar to...

Curve Fitting

If, for example, the adopted model is complex (many parameters involved) with respect to the number N , the model will fit the idiosyncrasies of the specific data set.

Thus

Thus, it will result in low bias but will yield high variance, as we change from one data set to another data set.

Similar to...

Curve Fitting

If, for example, the adopted model is complex (many parameters involved) with respect to the number N , the model will fit the idiosyncrasies of the specific data set.

Thus

Thus, it will result in low bias but will yield high variance, as we change from one data set to another data set.

Furthermore

If N grows we can have a more complex model to be fitted which reduces bias and ensures low variance.

• However, N is always finite!!!

Similar to...

Curve Fitting

If, for example, the adopted model is complex (many parameters involved) with respect to the number N , the model will fit the idiosyncrasies of the specific data set.

Thus

Thus, it will result in low bias but will yield high variance, as we change from one data set to another data set.

Furthermore

If N grows we can have a more complex model to be fitted which reduces bias and ensures low variance.

- However, N is always finite!!!

Thus

You always need to compromise

However, you always have some a priori knowledge about the data

Allowing you to impose restrictions

Lowering the bias and the variance

Nevertheless

We have the following example to grasp better the bothersome bias-variance dilemma.

Thus

You always need to compromise

However, you always have some a priori knowledge about the data

Allowing you to impose restrictions

Lowering the bias and the variance

Nonetheless

We have the following example to grasp better the bothersome bias-variance dilemma.

Thus

You always need to compromise

However, you always have some a priori knowledge about the data

Allowing you to impose restrictions

Lowering the bias and the variance

Nevertheless

We have the following example to grasp better the bothersome **bias–variance dilemma**.

For this

Assume

The data is generated by the following function

$$y = f(x) + \epsilon,$$

$$\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

We know that

The optimum regressor is $E[y|x] = f(x)$

Furthermore

Assume that the randomness in the different training sets, \mathcal{D} , is due to the y_i 's (Affected by noise), while the respective points, x_i , are fixed.

For this

Assume

The data is generated by the following function

$$y = f(x) + \epsilon,$$

$$\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

We know that

The optimum regressor is $E[y|x] = f(x)$

Assumptions

Assume that the randomness in the different training sets, \mathcal{D} , is due to the y_i 's (Affected by noise), while the respective points, x_i , are fixed.

For this

Assume

The data is generated by the following function

$$y = f(x) + \epsilon,$$
$$\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

We know that

The optimum regressor is $E[y|x] = f(x)$

Furthermore

Assume that the randomness in the different training sets, \mathcal{D} , is due to the y_i 's (Affected by noise), while the respective points, x_i , are fixed.

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Sampling the Space [2]

Imagine that $\mathcal{D} \subset [x_1, x_2]$ in which x lies

For example, you can choose $x_i = x_1 + \frac{x_2 - x_1}{N-1} (i - 1)$ with $i = 1, 2, \dots, N$

Case 1

Choose the estimate of $f(x)$, $g(x|\mathcal{D})$, to be independent of \mathcal{D}

For example, $g(x) = w_1x + w_0$

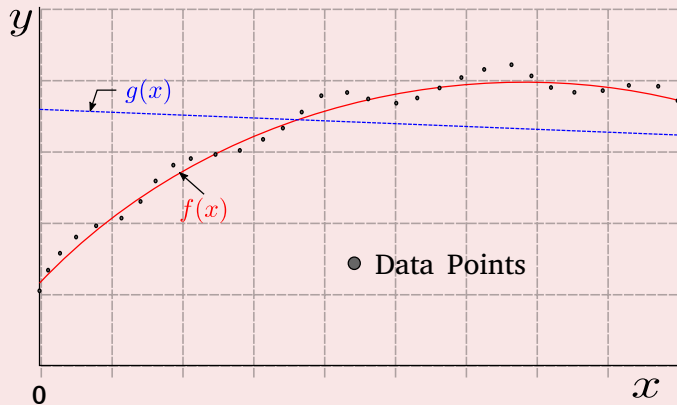
For example, the points are spread around $f(x)$ (σ)

Case 1

Choose the estimate of $f(x)$, $g(x|\mathcal{D})$, to be independent of \mathcal{D}

For example, $g(x) = w_1x + w_0$

For example, the points are spread around $(x, f(x))$



Case 1

Since $g(x)$ is fixed

$$E_{\mathcal{D}} [g(x|\mathcal{D})] = g(x|\mathcal{D}) \equiv g(x) \quad (4)$$

With

$$\text{Var}_{\mathcal{D}} [g(x|\mathcal{D})] = 0 \quad (5)$$

On the other hand

Because $g(x)$ was chosen arbitrarily the expected bias must be large.

$$\underbrace{(E_{\mathcal{D}} [g(x|\mathcal{D})] - E[y|x])^2}_{\text{BIAS}} \quad (6)$$

Case 1

Since $g(x)$ is fixed

$$E_{\mathcal{D}} [g(x|\mathcal{D})] = g(x|\mathcal{D}) \equiv g(x) \quad (4)$$

With

$$\text{Var}_{\mathcal{D}} [g(x|\mathcal{D})] = 0 \quad (5)$$

On the other hand:

Because $g(x)$ was chosen arbitrarily the expected bias must be large.

$$\underbrace{(E_{\mathcal{D}} [g(x|\mathcal{D})] - E[y|x])^2}_{\text{BIAS}} \quad (6)$$

Case 1

Since $g(x)$ is fixed

$$E_{\mathcal{D}} [g(x|\mathcal{D})] = g(x|\mathcal{D}) \equiv g(x) \quad (4)$$

With

$$\text{Var}_{\mathcal{D}} [g(x|\mathcal{D})] = 0 \quad (5)$$

On the other hand

Because $g(x)$ was chosen arbitrarily the expected bias must be large.

$$\underbrace{(E_{\mathcal{D}} [g(x|\mathcal{D})] - E[y|\mathbf{x}])^2}_{\text{BIAS}} \quad (6)$$

Case 2

In the other hand

Now, $g_1(x)$ corresponds to a polynomial of high degree so it can pass through each training point in \mathcal{D} .

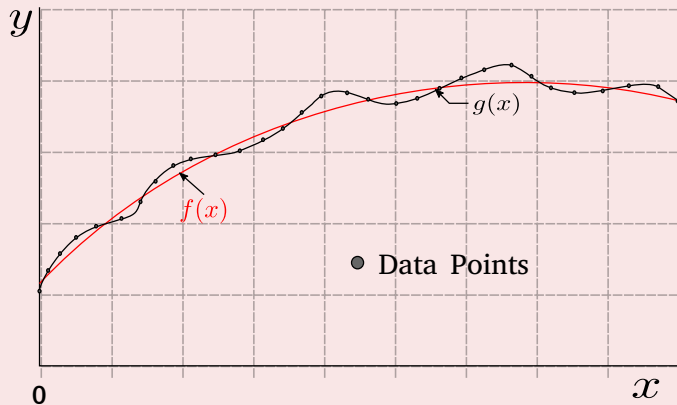
Example of $g_1(x)$

Case 2

In the other hand

Now, $g_1(x)$ corresponds to a polynomial of high degree so it can pass through each training point in \mathcal{D} .

Example of $g_1(x)$



Case 2

Due to the zero mean of the noise source

$$E_D [g_1(\mathbf{x}|\mathcal{D})] = f(x) = E[y|x] \text{ for any } x = x_i \quad (7)$$

Remark: At the training points the bias is zero.

However the variance increases

$$\begin{aligned} E_D \left[(g_1(\mathbf{x}|\mathcal{D}) - E_D [g_1(\mathbf{x}|\mathcal{D})])^2 \right] &= E_D \left[(f(x) + \epsilon - f(x))^2 \right] \\ &= \sigma_\epsilon^2, \text{ for } x = x_i, i = 1, 2, \dots, N \end{aligned}$$

In other words

The bias becomes zero (or approximately zero) but the variance is now equal to the variance of the noise source.

Case 2

Due to the zero mean of the noise source

$$E_D [g_1 (\mathbf{x}|\mathcal{D})] = f (x) = E [y|x] \text{ for any } x = x_i \quad (7)$$

Remark: At the training points the bias is zero.

However the variance increases

$$\begin{aligned} E_D \left[(g_1 (\mathbf{x}|\mathcal{D}) - E_D [g_1 (\mathbf{x}|\mathcal{D})])^2 \right] &= E_D \left[(f (x) + \epsilon - f (x))^2 \right] \\ &= \sigma_\epsilon^2, \text{ for } x = x_i, i = 1, 2, \dots, N \end{aligned}$$

In other words

The bias becomes zero (or approximately zero) but the variance is now equal to the variance of the noise source.

Case 2

Due to the zero mean of the noise source

$$E_D [g_1 (\mathbf{x}|\mathcal{D})] = f (x) = E [y|x] \text{ for any } x = x_i \quad (7)$$

Remark: At the training points the bias is zero.

However the variance increases

$$\begin{aligned} E_D \left[(g_1 (\mathbf{x}|\mathcal{D}) - E_D [g_1 (\mathbf{x}|\mathcal{D})])^2 \right] &= E_D \left[(f (x) + \epsilon - f (x))^2 \right] \\ &= \sigma_\epsilon^2, \text{ for } x = x_i, i = 1, 2, \dots, N \end{aligned}$$

In other words

The bias becomes zero (or approximately zero) but the variance is now equal to the variance of the noise source.

Observations

First

Everything that has been said so far applies to both the regression and the classification tasks.

However:

Mean squared error is not the best way to measure the power of a classifier.

Think about:

A classifier that sends everything far away of the hyperplane!!! Away from the values $+ - 1$!!!

Observations

First

Everything that has been said so far applies to both the regression and the classification tasks.

However

Mean squared error is not the best way to measure the power of a classifier.

Think about

A classifier that sends everything far away of the hyperplane!!! Away from the values $+ - 1$!!!

Observations

First

Everything that has been said so far applies to both the regression and the classification tasks.

However

Mean squared error is not the best way to measure the power of a classifier.

Think about

A classifier that sends everything far away of the hyperplane!!! Away from the values $+ - 1!!!$

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- **Intuition from Overfitting**
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

The house example (From Andrew Ng Course)

Imagine the following data set



Now assume that we use a regressor

For the fitting

$$\frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2$$

We can then run one of our machines to see what minimize better the previous equation.

Question: Did you notice that I did not impose any structure to $h_w(x)$?

Now assume that we use a regressor

For the fitting

$$\frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2$$

We can then run one of our machine to see what minimize better the previous equation

Question: Did you notice that I did not impose any structure to $h_w(x)$?

Then, First fitting

What about using $h_1(x) = \theta_0 + \theta_1x + \theta_2x^2$?



Second fitting

What about using $h_2(x) = \theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3 + \theta_4x^4 + \theta_5x^5$?



Therefore, we have a problem

We get weird over fitting effects!!!

What do we do? What about minimizing the influence of $\theta_3, \theta_4, \theta_5$?

How do we do that?

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2$$

What about integrating those values to the cost function? Ideas

Therefore, we have a problem

We get weird over fitting effects!!!

What do we do? What about minimizing the influence of $\theta_3, \theta_4, \theta_5$?

How do we do that?

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2$$

What about integrating those values to the cost function? Ideas

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- **The Idea of Regularization**
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

We have

Regularization intuition is as follow

Small values for parameters $\theta_0, \theta_1, \theta_2, \dots, \theta_n$

It implies

- "Simpler" function
- Less prone to overfitting

We have

Regularization intuition is as follow

Small values for parameters $\theta_0, \theta_1, \theta_2, \dots, \theta_n$

It implies

- 1 "Simpler" function
- 2 Less prone to overfitting

We can do the previous idea for the other parameters

We can do the same for the other parameters

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2 + \sum_{i=1}^d \lambda_i \theta_i^2 \quad (8)$$

However, handling such many parameters can be so difficult

Combinatorial problem in reality!!!

We can do the previous idea for the other parameters

We can do the same for the other parameters

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2 + \sum_{i=1}^d \lambda_i \theta_i^2 \quad (8)$$

However handling such many parameters can be so difficult

Combinatorial problem in reality!!!

Better, we can

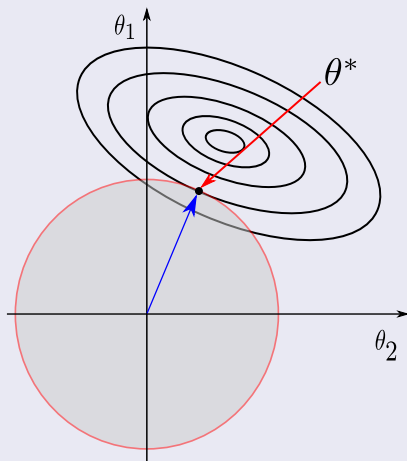
We better use the following

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2 + \lambda \sum_{i=1}^d \theta_i^2 \quad (9)$$

Graphically

Geometrically Equivalent to send our function to something quadratic

$$\frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2 + \lambda \sum_{i=1}^d \theta_i^2$$



An interesting Observation, when using linear estimators

The function $\sum_{i=1}^N (\theta^T \mathbf{x}_i - y_i)^2$

- It is a convex function...

And also $\|\theta\|_2^2$

- It is also a convex function...

Therefore the final Lagrangian is a Convex function

- Here, Regularization basically remove dimensions that could not be useful in the minimization of the linear estimator.

An interesting Observation, when using linear estimators

The function $\sum_{i=1}^N (\theta^T \mathbf{x}_i - y_i)^2$

- It is a convex function...

And also $\sum_{i=1}^d \theta_i$

- It is also a convex function...

Therefore the final loss function is a Convex function

- Here, Regularization basically remove dimensions that could not be useful in the minimization of the linear estimator.

An interesting Observation, when using linear estimators

The function $\sum_{i=1}^N (\theta^T \mathbf{x}_i - y_i)^2$

- It is a convex function...

And also $\sum_{i=1}^d \theta_i$

- It is also a convex function...

Therefore the final Lagrangian is a Convex function

- Here, Regularization basically remove dimensions that could not be useful in the minimization of the linear estimator.

However

The game changes a lot

- When the estimator is a complex non-convex function

Implications

- Deep Learners

However

The game changes a lot

- When the estimator is a complex non-convex function

In our case

- Deep Learners

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- **Ridge Regression**
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Ridge Regression

Equation

$$\hat{\theta} = \arg \min_w \left\{ \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^d x_{ij} \theta_j \right)^2 + \lambda \sum_{j=1}^d \theta_j^2 \right\}$$

Here

- $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage

The Larger $\lambda \geq 0$

- The coefficients are shrunk toward zero (and each other).

Ridge Regression

Equation

$$\hat{\theta} = \arg \min_w \left\{ \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^d x_{ij} \theta_j \right)^2 + \lambda \sum_{j=1}^d \theta_j^2 \right\}$$

Here

- $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage

The Larger $\lambda \geq 0$

- The coefficients are shrunk toward zero (and each other).

This is also can be written

Optimization Solution

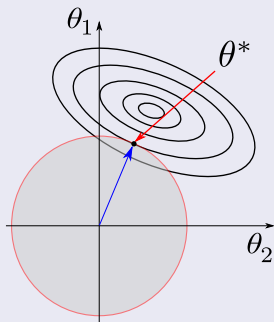
$$\arg \min_{\theta} \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^d x_{ij} \theta_j \right)^2$$

subject to $\sum_{j=1}^d \theta_j^2 < t$

Graphically

Geometrically Equivalent to

$$\begin{aligned} \arg \min \quad & \sum_{i=1}^N (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2 \\ \text{subject to} \quad & \sum_{i=1}^{d+1} \theta_i^2 < t \end{aligned}$$



Outline

1

Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2

The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- **The LASSO**
- Generalization
- What can be done?

3

Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\theta}^{garrote} = \arg \min_{\theta} \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^d x_{ij} \theta_j \right)^2 + N\lambda \sum_{j=1}^d \theta_j$$

s.t. $\theta_j > 0 \forall j$

This is quite desirable

However, Tibshirani realized that you could get a more flexible model by using the absolute value at the constraint!!!

Robert Tibshirani proposed the use of the L_1 norm

$$\|\theta\|_1 = \sum_{i=1}^d |\theta_i|$$

Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\theta}^{garrote} = \arg \min_{\theta} \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^d x_{ij} \theta_j \right)^2 + N\lambda \sum_{j=1}^d \theta_j$$

s.t. $\theta_j > 0 \forall j$

This is quite derivable

However, Tibshirani realized that you could get a more flexible model by using the absolute value at the constraint!!!

Robert Tibshirani proposed the use of the L_1 norm

$$\|\theta\|_1 = \sum_{i=1}^d |\theta_i|$$

Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\boldsymbol{\theta}}^{garrote} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^d x_{ij} \theta_j \right)^2 + N\lambda \sum_{j=1}^d \theta_j$$

s.t. $\theta_j > 0 \forall j$

This is quite derivable

However, Tibshirani realized that you could get a more flexible model by using the absolute value at the constraint!!!

Robert Tibshirani proposed the use of the L_1 norm

$$\|\boldsymbol{\theta}\|_1 = \sum_{i=1}^d |\theta_i|$$

The Final Optimization Problem

LASSO

$$\hat{\boldsymbol{\theta}}^{LASSO} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^d x_{ij} \theta_j \right)^2$$
$$\text{s.t. } \sum_{i=1}^d |\theta_i| \leq t$$

This is not derivable

More advanced methods are necessary to solve this problem!!!

The Final Optimization Problem

LASSO

$$\hat{\boldsymbol{\theta}}^{LASSO} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^d x_{ij} \theta_j \right)^2$$
$$\text{s.t. } \sum_{i=1}^d |\theta_i| \leq t$$

This is not derivable

More advanced methods are necessary to solve this problem!!!

The Lagrangian Version

The Lagrangian

$$\hat{\boldsymbol{\theta}}^{LASSO} = \arg \min_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^N (y_i - \mathbf{x}^T \boldsymbol{\theta})^2 + \lambda \sum_{i=1}^d |\theta_i| \right\}$$

However

You have other regularizations as $\|\boldsymbol{\theta}\|_2 = \sqrt{\sum_{i=1}^d |\theta_i|^2}$

The Lagrangian Version

The Lagrangian

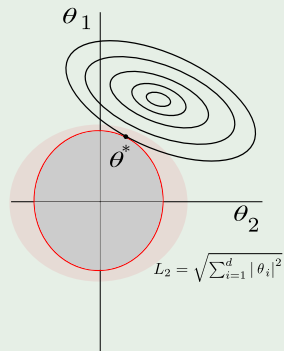
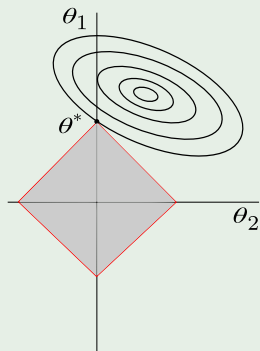
$$\hat{\boldsymbol{\theta}}^{LASSO} = \arg \min_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^N (y_i - \mathbf{x}^T \boldsymbol{\theta})^2 + \lambda \sum_{i=1}^d |\theta_i| \right\}$$

However

You have other regularizations as $\|\boldsymbol{\theta}\|_2 = \sqrt{\sum_{i=1}^d |\theta_i|^2}$

Graphically

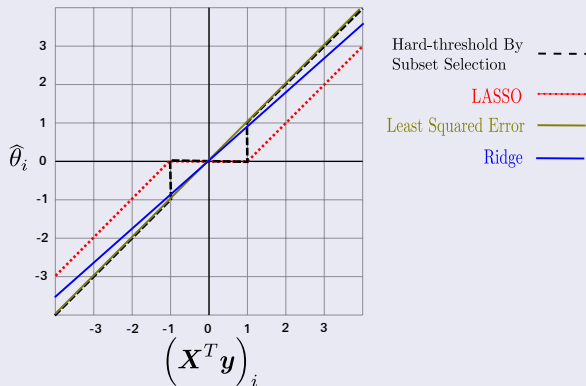
Yes the circle defined as $\|\theta\|_2 = \sqrt{\sum_{i=1}^d |\theta_i|^2}$



For Example

In the Case of \mathbf{X} is a Orthogonal Matrix, we have

$$\hat{\theta}_i = \text{sgn} \left(X^T y \right)_i \left(\left(X^T y \right)_i - \sigma^2 \alpha \right)_+$$



The seminal paper by Robert Tibshirani

An initial study of this regularization can be seen in

“Regression Shrinkage and Selection via the LASSO” by Robert Tibshirani
- 1996

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- **Generalization**
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

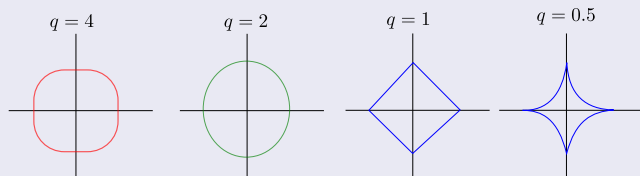
Furthermore

We can generalize ridge regression and the lasso, and view them as Bayes estimates

$$\hat{\boldsymbol{\theta}}^{LASSO} = \arg \min_{\boldsymbol{w}} \left\{ \sum_{i=1}^N (y_i - L(\boldsymbol{x}_i, \boldsymbol{\theta}))^2 + \lambda \sum_{i=1}^d |\theta_i|^q \right\} \text{ with } q \geq 0$$

For Example

We have when $d = 2$

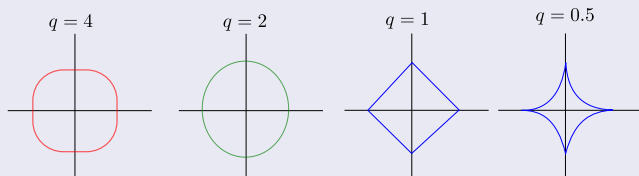


Here when $d = 1$

- You are having a derivable Lagrangian, but you lose the LASSO properties

For Example

We have when $d = 2$



Here, when $q > 1$

- You are having a derivable Lagrangian, but you lose the LASSO properties

Therefore

Zou and Hastie (2005) introduced the elastic-net penalty [3]

$$\lambda \sum_{i=1}^d \left\{ \alpha \theta_i^2 + (1 - \alpha) |\theta_i| \right\}$$

This is basically

- A Compromise Between the Ridge and LASSO.

Therefore

Zou and Hastie (2005) introduced the elastic-net penalty [3]

$$\lambda \sum_{i=1}^d \left\{ \alpha \theta_i^2 + (1 - \alpha) |\theta_i| \right\}$$

This is Basically

- A Compromise Between the Ridge and LASSO.

Outline

1

Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2

The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- **What can be done?**

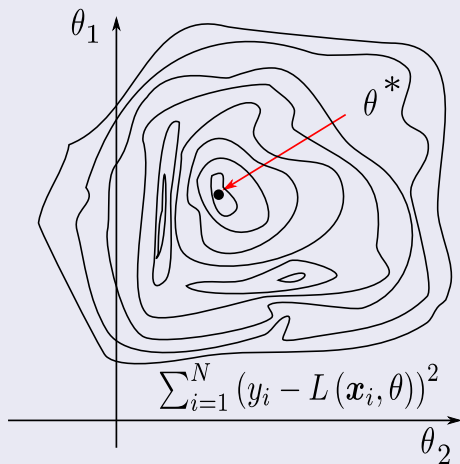
3

Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

What can be done?

Remember that our optimization Landscape is highly variable



Over-fitting?

Basically (Intuition)

$$(y_i - L(\mathbf{x}_i, \theta))^2 = 0 \text{ for } i \in \textit{Training}$$

$$(y_j - L(\mathbf{x}_j, \theta))^2 \gg 0 \text{ for } i \in \textit{Validation}$$

At the other side, you have BIAS \rightarrow Simplification

- Then, Regularization is an operator moving the model toward a bias

Over-fitting?

Basically (Intuition)

$$(y_i - L(\mathbf{x}_i, \theta))^2 = 0 \text{ for } i \in \textit{Training}$$

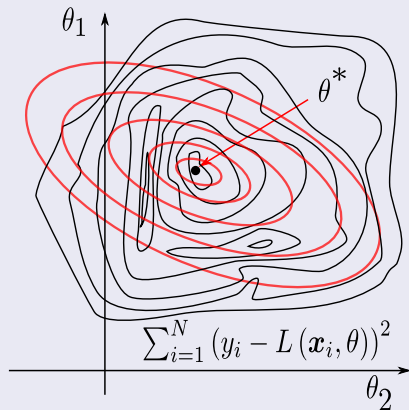
$$(y_j - L(\mathbf{x}_i, \theta))^2 \gg 0 \text{ for } i \in \textit{Validation}$$

A the other side, you have BIAS==Simplification

- Then, Regularization is an operator moving the model toward a bias

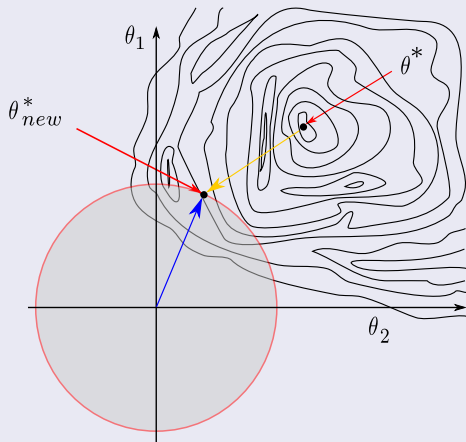
However, we do not want too much simplification

Look at this, the worst case Bias toward Red



Basically this simplification is due to the constrained optimization landscape

Basically our constraint is too Euclidean for Optimization Landscape



Well-Posed Problem

Definition by Hadamard (Circa 1902)

- Models of physical phenomena should have the following properties
 - 1 A solution exists,
 - 2 The solution is unique,
 - 3 The solution's behavior changes continuously with the initial conditions.

Any other problem that fails in any of these conditions

- It is considered an Ill-Posed Problem.

Well-Posed Problem

Definition by Hadamard (Circa 1902)

- Models of physical phenomena should have the following properties
 - 1 A solution exists,
 - 2 The solution is unique,
 - 3 The solution's behavior changes continuously with the initial conditions.

Any other problem that fails in any of these conditions

- It is considered an Ill-Posed Problem.

It seems to be that

The Deep Learners are highly ill-posed problems

- Ridge and LASSO have two possible effects

Too much simplification

- The Deep Learners losses power of representation.
 - ▶ Weights are eliminated

The constraints forces them to

- They are forced to live in a too smooth optimization landscape

It seems to be that

The Deep Learners are highly ill-posed problems

- Ridge and LASSO have two possible effects

Too much simplification

- The Deep Learners losses power of representation.
 - ▶ Weights are eliminated

The constraints forces them to

- They are forced to live in a too smooth optimization landscape

It seems to be that

The Deep Learners are highly ill-posed problems

- Ridge and LASSO have two possible effects

Too much simplification

- The Deep Learners losses power of representation.
 - ▶ Weights are eliminated

The constraints forces the θ 's

- They are forced to live in a too smooth optimization landscape

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- **Gaussian Noise on Hidden Units for Regularization**
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

DeVris and Taylor [5]

For many years

- Dataset augmentation has been a standard regularization technique used to reduce overfitting while training supervised learning models

For Example, LeCun et al. [8] when training the LeNet5

- They applied a series of transformations to the input images in order to improve the robustness of the model.

Unfortunately

- Dataset augmentation is not as straightforward to apply in all domains as it is for images.

DeVris and Taylor [5]

For many years

- Dataset augmentation has been a standard regularization technique used to reduce overfitting while training supervised learning models

For Example, LeCun et al. [4] when training the LeNet5

- They applied a series of transformations to the input images in order to improve the robustness of the model.

Unfortunately

- Dataset augmentation is not as straightforward to apply in all domains as it is for images.

DeVris and Taylor [5]

For many years

- Dataset augmentation has been a standard regularization technique used to reduce overfitting while training supervised learning models

For Example, LeCun et al. [4] when training the LeNet5

- They applied a series of transformations to the input images in order to improve the robustness of the model.

Unfortunately

- Dataset augmentation is not as straightforward to apply in all domains as it is for images.

For Example

In voice detection, adding

- 1 Gaussian noise to the input,
- 2 Shifting the pitch of the audio signal,
- 3 Time stretching,
- 4 Varying the loudness of the audio signal,
- 5 Applying random frequency filters,
- 6 Interpolating between samples in input space.

For Example

In voice detection, adding

- 1 Gaussian noise to the input,
- 2 Shifting the pitch of the audio signal,
- 3 Time stretching,
- 4 Varying the loudness of the audio signal,
- 5 Applying random frequency filters,
- 6 Interpolating between samples in input space.

Actually, only the following techniques worked out

- Pitch shifting and random frequency filtering

For Example

In voice detection, adding

- 1 Gaussian noise to the input,
- 2 Shifting the pitch of the audio signal,
- 3 Time stretching,
- 4 Varying the loudness of the audio signal,
- 5 Applying random frequency filters,
- 6 Interpolating between samples in input space.

Actually, only the following techniques worked out

- Pitch shifting and random frequency filtering

For Example

In voice detection, adding

- 1 Gaussian noise to the input,
 - 2 Shifting the pitch of the audio signal,
 - 3 Time stretching,
 - 4 Varying the loudness of the audio signal,
- 5 Applying random frequency filters,
 - 6 Interpolating between samples in input space.

Actually, only the following techniques worked out

- Pitch shifting and random frequency filtering

For Example

In voice detection, adding

- 1 Gaussian noise to the input,
 - 2 Shifting the pitch of the audio signal,
 - 3 Time stretching,
 - 4 Varying the loudness of the audio signal,
 - 5 Applying random frequency filters,
- 6 Interpolating between samples in input space.

Actually, only the following techniques worked out:

- Pitch shifting and random frequency filtering

For Example

In voice detection, adding

- 1 Gaussian noise to the input,
- 2 Shifting the pitch of the audio signal,
- 3 Time stretching,
- 4 Varying the loudness of the audio signal,
- 5 Applying random frequency filters,
- 6 Interpolating between samples in input space.

Actually, only the following techniques worked out:

- Pitch shifting and random frequency filtering

For Example

In voice detection, adding

- 1 Gaussian noise to the input,
- 2 Shifting the pitch of the audio signal,
- 3 Time stretching,
- 4 Varying the loudness of the audio signal,
- 5 Applying random frequency filters,
- 6 Interpolating between samples in input space.

Actually, only the following techniques worked out

- Pitch shifting and random frequency filtering

DeVris and Taylor [5]

They did something different

- First learning a data representation
- Then applying transformations to samples mapped to that representation.

They hypothesized

- Due to manifold unfolding in feature space, simple transformations applied to encoded rather than raw inputs
 - ▶ They will result in more plausible synthetic data.

DeVris and Taylor [5]

They did something different

- First learning a data representation
- Then applying transformations to samples mapped to that representation.

They hypothesized

- Due to manifold unfolding in feature space, simple transformations applied to **encoded** rather than raw inputs
 - ▶ They will result in more plausible synthetic data.

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

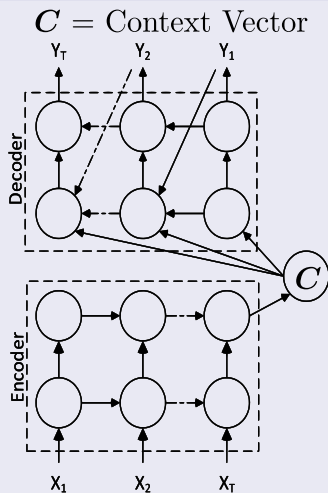
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Decoder/Encoder Part

We have a Decoder and Encoder Architecture



Basically

They used a context C to pass information between the encoder and decoder

- Here is where the authors performed the augmentation

Explicitly

- At the context, something like the embeddings at document level.

Basically

They used a context C to pass information between the encoder and decoder

- Here is where the authors performed the augmentation

Basically

- At the context, something like the embeddings at document level.

Here

We have a K-coding symbol set

- The Encoder and Decoder are based in a novel hidden unit.

We have the following configuration per element j :

$$r_j = \sigma \left([W_r x]_j + [U_r h_{t-1}]_j \right) \leftarrow \text{Reset Gate}$$

- σ a sigmoid function

The Update gate

$$z_j = \sigma \left([W_z x]_j + [U_z h_{t-1}]_j \right)$$

Here

We have a K-coding symbol set

- The Encoder and Decoder are based in a novel hidden unit.

We have the following configuration per row element j

$$r_j = \sigma \left([\mathbf{W}_r \mathbf{x}]_j + [\mathbf{U}_r \mathbf{h}_{t-1}]_j \right) \leftarrow \text{Reset Gate}$$

- σ a sigmoid function

The Update gate

$$z_j = \sigma \left([\mathbf{W}_z \mathbf{x}]_j + [\mathbf{U}_z \mathbf{h}_{t-1}]_j \right)$$

Here

We have a K-coding symbol set

- The Encoder and Decoder are based in a novel hidden unit.

We have the following configuration per row element j

$$r_j = \sigma \left([\mathbf{W}_r \mathbf{x}]_j + [\mathbf{U}_r \mathbf{h}_{t-1}]_j \right) \leftarrow \text{Reset Gate}$$

- σ a sigmoid function

The Update gate

$$z_j = \sigma \left([\mathbf{W}_z \mathbf{x}]_j + [\mathbf{U}_z \mathbf{h}_{t-1}]_j \right)$$

Where

The Activation Gate update

$$h_j^t = z_j h_j^{t-1} + (1 - z_j) \tilde{h}_j^t$$

- Where $\tilde{h}_j^t = \phi \left([\mathbf{W}\mathbf{x}]_j + [\mathbf{U}(\mathbf{r} \odot \mathbf{h}_{t-1})]_j \right)$

In the formulation

- When the reset gate is close to 0, the hidden state is forced to ignore the previous hidden state!!!

Where

The Activation Gate update

$$h_j^t = z_j h_j^{t-1} + (1 - z_j) \tilde{h}_j^t$$

- Where $\tilde{h}_j^t = \phi \left([\mathbf{W}\mathbf{x}]_j + [\mathbf{U}(\mathbf{r} \odot \mathbf{h}_{t-1})]_j \right)$

In this formulation

- When the reset gate is close to 0, the hidden state is forced to ignore the previous hidden state!!!

Finally, at output

We have a probability of producing a symbol of a set of at the Decoder

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{c}) = \frac{\exp(W_o \mathbf{h}_t + U_o y_{t-1} + \mathbf{c}_{t-1})}{\sum_{j=1}^K \exp(W_j \mathbf{h}_t + U_o y_{t-1} + \mathbf{c}_{t-1})}$$

Then, at the Encoder

- The encoder learns to predict the next symbol x_t based in the previous $x_{t-1}, x_{t-2}, \dots, x_1$ by using the maximization

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N p(y_n | x_n)$$

Finally, at output

We have a probability of producing a symbol of a set of at the Decoder

$$p(y_t | y_{t-1}, \dots, y_1, \mathbf{c}) = \frac{\exp(W_o \mathbf{h}_t + U_o y_{t-1} + \mathbf{c}_{t-1})}{\sum_{j=1}^K \exp(W_j \mathbf{h}_t + U_o y_{t-1} + \mathbf{c}_{t-1})}$$

Then, at the Encoder

- The encoder learns to predict the next symbol x_t based in the previous $x_{t-1}, x_{t-2}, \dots, x_1$ by using the maximization

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n)$$

Here, the Noise

Generate noise by drawing from

- A Gaussian distribution with **zero mean** and **per-element standard deviation** calculated across all context vectors in the dataset

$$c'_i = c_i + \gamma X, \quad X \sim N(0, \sigma_i^2)$$

We can generate this using a more direct approach

- For each sample in the dataset, we find its K nearest neighbors in feature space, then

$$c' = (c_k - c_j)\lambda + c_j$$

- $\lambda = 0.5$

Here, the Noise

Generate noise by drawing from

- A Gaussian distribution with **zero mean** and **per-element standard deviation** calculated across all context vectors in the dataset

$$c'_i = c_i + \gamma X, \quad X \sim N(0, \sigma_i^2)$$

We can generate this using a more direct approach

- For each sample in the dataset, we find its K nearest neighbors in feature space, then

$$c' = (c_k - c_j) \lambda + c_j$$

- $\lambda = 0.5$

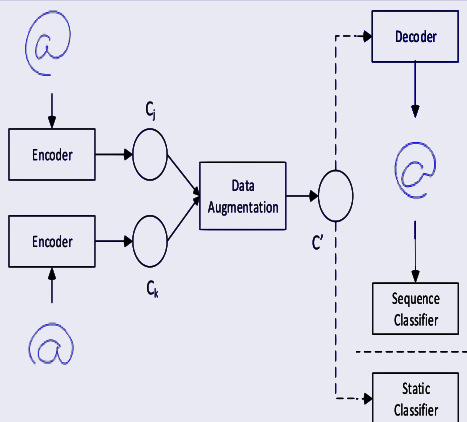
Then

Once this new augmented context vectors with noise are ready

- As input for a learning task,
- They can be decoded to generate new sequences

Finally, we have

The following architecture where two symbols are encoded



Results

Not so much improvement

Image Size	Description	Test Error	Test Error (Reconstructions of original data)
32 × 32	Original dataset	8.59 ± 0.24	-
24 × 24	Center crop	11.28 ± 0.25	18.54 ± 0.38
24 × 24	Center crop + extrapolation	13.90 ± 0.22	17.69 ± 0.39
24 × 24	Simple data augmentation	7.33 ± 0.17	13.60 ± 0.17
24 × 24	Simple data augmentation + extrapolation	8.80 ± 0.24	12.00 ± 0.23

Why is this happening?

It is the same problem at the exit point

- We are regularizing at the encoded input space... but the architecture is still there...

Therefore

- It is necessary to do something quite different...

Why is this happening?

It is the same problem at the exit point

- We are regularizing at the encoded input space... but the architecture is still there...

Therefore

- It is necessary to do something quite different...

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- **Dropout as Regularization**
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Outline

- 1 Bias-Variance Dilemma
 - Introduction
 - Measuring the difference between optimal and learned
 - The Bias-Variance
 - "Extreme" Example
- 2 The Problem with Overfitting
 - Intuition from Overfitting
 - The Idea of Regularization
 - Ridge Regression
 - The LASSO
 - Generalization
 - What can be done?
- 3 Methods of Regularization for Deep Networks
 - Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
 - **Dropout as Regularization**
 - **Introduction**
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
 - Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
 - Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
 - For More in Normalization

Regularization in Deep Forward

In Layers of a Deep Forward

- We want to find an estimation x_t^r to an input at $x_0 \in \mathbb{R}^d$ in layer t satisfying

$$\sigma(A_t^r x_t) = y_{t+1}$$

Regularization in Deep Forward

In Layers of a Deep Forward

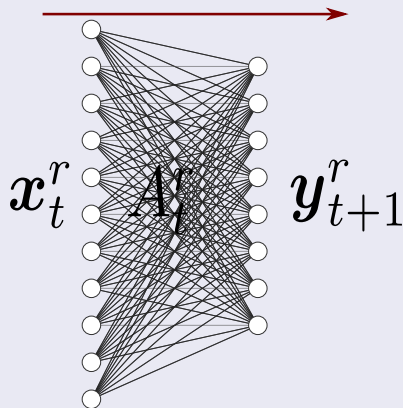
- We want to find an estimation \mathbf{x}_t^r to an input at $\mathbf{x}_0 \in \mathbb{R}^d$ in layer t satisfying

$$\sigma(A_t^r \mathbf{x}_t) = \mathbf{y}_{t+1}$$

We can see this

A flow of information

FORWARD FLOW OF INFORMATION



In all such situations

The vector x_t is generated by y_{t+1} using back-propagation

$$A_t^r = A_t^{r-1} - \eta \frac{\partial L(A_T^{r-1}, \dots, A_0^{r-1}, x_0)}{\partial A_t^{r-1}}$$

It is usually a meaningless, bad approximation

- to x^* optimal at layer t for all possible inputs x'_0 's.

In all such situations

The vector x_t is generated by y_{t+1} using back-propagation

$$A_t^r = A_t^{r-1} - \eta \frac{\partial L(A_T^{r-1}, \dots, A_0^{r-1}, x_0)}{\partial A_t^{r-1}}$$

It is usually a meaningless bad approximation

- to x^* optimal at layer t for all possible inputs x'_0 's.

Then

We can see the Deep Forward Network as

$$y_T = \sigma (A_T \sigma (A_{T-1} \sigma (A_{T-2} (\dots \sigma (A_0 x_0))))))$$

File

- The σ is applied to the generated vectors point wise...

Then

We can see the Deep Forward Network as

$$y_T = \sigma (A_T \sigma (A_{T-1} \sigma (A_{T-2} (\dots \sigma (A_0 x_0))))))$$

Here

- The σ is applied to the generated vectors point wise...

The Jacobian of the Gradient Descent

Here, we assume a Least Squared Error cost function

$$\frac{\partial L(A_T^{r-1}, \dots, A_0^{r-1}, x_0^i)}{\partial A_t^{r-1}} = -(z^i - y_T) \times \sigma'(A_{T-1}^r \mathbf{x}_{T-1}) \times \frac{\partial A_{T-1}^r \mathbf{x}_{T-1}}{\partial \mathbf{x}_{T-1}} \times \dots \times \sigma'(A_t^r \mathbf{x}_t) \times \frac{\partial A_t^r \mathbf{x}_t}{\partial \mathbf{x}_t}$$

Where

$$\sigma'(A_k^r \mathbf{x}_k) = \begin{pmatrix} \sigma'(a_{1k}^r x_k) & 0 & \dots & 0 \\ 0 & \sigma'(a_{2k}^r x_k) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma'(a_{Mk}^r x_k) \end{pmatrix}$$

The Jacobian of the Gradient Descent

Here, we assume a Least Squared Error cost function

$$\frac{\partial L(A_T^{r-1}, \dots, A_0^{r-1}, x_0^i)}{\partial A_t^{r-1}} = -(z^i - y_T) \times \sigma'(A_{T-1}^r \mathbf{x}_{T-1}) \times \frac{\partial A_{T-1}^r \mathbf{x}_{T-1}}{\partial \mathbf{x}_{T-1}} \times \dots \times \sigma'(A_t^r \mathbf{x}_t) \times \frac{\partial A_t^r \mathbf{x}_t}{\partial \mathbf{x}_t}$$

Where

$$\sigma'(A_k^r \mathbf{x}_k) = \begin{pmatrix} \sigma'(a_{1k}^r \mathbf{x}_k) & 0 & \dots & 0 \\ 0 & \sigma'(a_{2k}^r \mathbf{x}_k) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma'(a_{Mk}^r \mathbf{x}_k) \end{pmatrix}$$

What will happen in the following situation?

Imagine that A'_k 's are diagonal matrix

$$A_k^r = \begin{pmatrix} a_{1k} & 0 & \cdots & 0 \\ 0 & a_{2k} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{Mk} \end{pmatrix}$$

Therefore, we have

$$\sigma'(A_k^r x_k) = \begin{pmatrix} \sigma'(a_{1k}^r x_{1k}) & 0 & \cdots & 0 \\ 0 & \sigma'(a_{2k}^r x_{2k}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(a_{Mk}^r x_{Mk}) \end{pmatrix}$$

What will happen in the following situation?

Imagine that A_k^r 's are diagonal matrix

$$A_k^r = \begin{pmatrix} a_{1k} & 0 & \cdots & 0 \\ 0 & a_{2k} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{Mk} \end{pmatrix}$$

Therefore, we have

$$\sigma'(A_k^r \mathbf{x}_k) = \begin{pmatrix} \sigma'(a_{1k}^r x_{1k}) & 0 & \cdots & 0 \\ 0 & \sigma'(a_{2k}^r x_{2k}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(a_{Mk}^r x_{2k}) \end{pmatrix}$$

Then, we have that

First

$$\sigma' \left(A_{T-1}^r \mathbf{x}_{T-1} \right) \times \frac{\partial A_{T-1}^r \mathbf{x}_{T-1}}{\partial \mathbf{x}_{T-1}} \times \dots \times \sigma' \left(A_t^r \mathbf{x}_t \right) \times \frac{\partial A_t^r \mathbf{x}_t}{\partial \mathbf{x}_t} = *$$

Then, we have that

$$* = \begin{pmatrix} \prod_{k=T-1}^t \sigma' \left(a_{1k}^r x_{1k} \right) a_{1k} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \prod_{k=T-1}^t \sigma' \left(a_{Mk}^r x_{2k} \right) a_{2k} \end{pmatrix}$$

Then, we have that

First

$$\sigma' (A_{T-1}^r \mathbf{x}_{T-1}) \times \frac{\partial A_{T-1}^r \mathbf{x}_{T-1}}{\partial \mathbf{x}_{T-1}} \times \dots \times \sigma' (A_t^r \mathbf{x}_t) \times \frac{\partial A_t^r \mathbf{x}_t}{\partial \mathbf{x}_t} = *$$

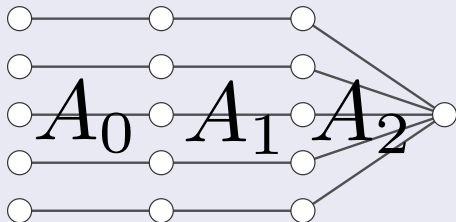
Then, we have that

$$* = \begin{pmatrix} \prod_{k=T-1}^t \sigma' (a_{1k}^r x_{1k}) a_{1k} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \prod_{k=T-1}^t \sigma' (a_{Mk}^r x_{2k}) a_{2k} \end{pmatrix}$$

Actually

Choosing Matrices in such way

- It is like a heavy simplification of the Deep Forward Network



Something happens with the LASSO and Ridge

At the top of the Optimization Cost Function

- We do not know how such shallow regularization can affect the Neural Network

So heavy regularization

- It can not be a so good idea...

We need a new way of doing stuff

- For example, we could do the following...

Something happens with the LASSO and Ridge

At the top of the Optimization Cost Function

- We do not know how such shallow regularization can affect the Neural Network

So heavy regularization

- It can not be a so good idea...

We need a new way of doing stuff

- For example, we could do the following...

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- **Dropout as Regularization**
 - Introduction
 - **Dropout Process**
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
 - Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
 - Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
 - For More in Normalization

Dropout

It was introduced by Hinton and Google [6]

- To avoid the problem of over-fitting

You can see it as a regularization

- From [7] "Dropout training as adaptive regularization" by Wager et al.

Dropout

It was introduced by Hinton and Google [6]

- To avoid the problem of over-fitting

You can see it as a regularization

- From [7] “Dropout training as adaptive regularization” by Wager et al.

He comments that with unlimited computations

- “the best way to “regularize” a fixed-sized model is to average the predictions of all possible settings of the parameters”

Something like Boosting [4]

- By Using simpler and smaller models

He comments that with unlimited computations

- “the best way to “regularize” a fixed-sized model is to average the predictions of all possible settings of the parameters”

Something like Boosting [1]

- By Using simpler and smaller models

Problem

We have Deep Architectures with thousands of parameters and hyperparameters

- Therefore, we have a problem!!! We need to solve this in some way!!!

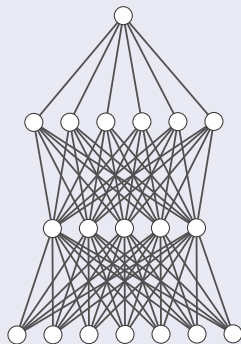
What if we fix our architecture

Problem

We have Deep Architectures with thousands of parameters and hyperparameters

- Therefore, we have a problem!!! We need to solve this in some way!!!

What if we fix our architecture



How it works?

You have forward layers

$$z_i^{l+1} = W_i^{l+1} \mathbf{x}^l + b_i^{l+1}$$
$$x_i^{l+1} = \sigma(z_i^{l+1})$$

With dropout, the test-forward operation becomes

$$r_j^l \sim \text{Bernoulli}(p)$$
$$\tilde{\mathbf{x}}^l = \mathbf{r}^l \odot \mathbf{x}^l$$
$$z_i^{l+1} = W_i^{l+1} \tilde{\mathbf{x}}^l + b_i^{l+1}$$
$$x_i^{l+1} = \sigma(z_i^{l+1})$$

How it works?

You have forward layers

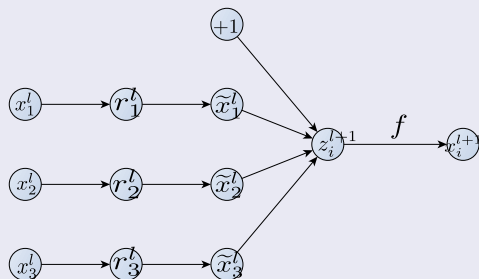
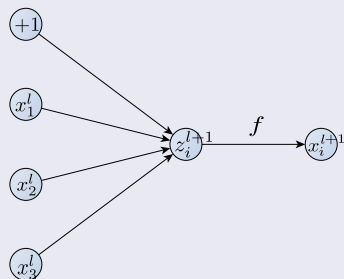
$$z_i^{l+1} = W_i^{l+1} \mathbf{x}^l + b_i^{l+1}$$
$$x_i^{l+1} = \sigma(z_i^{l+1})$$

With dropout, the feed-forward operation becomes

$$r_j^l \sim \text{Bernoulli}(p)$$
$$\tilde{\mathbf{x}}^l = \mathbf{r}^l \odot \mathbf{x}^l$$
$$z_i^{l+1} = W_i^{l+1} \tilde{\mathbf{x}}^l + b_i^{l+1}$$
$$x_i^{l+1} = \sigma(z_i^{l+1})$$

The Network

It looks like a series of gates



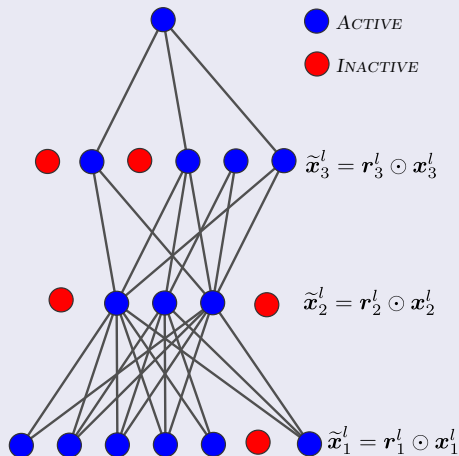
Therefore

We have that sampling is done in a Bernoulli to generate the r^l , a vector of Bernoulli random variables

- Then, the layers are thinned by the wise multiplication with the nodes at each layer

Then, we erase randomly connections through the network

We generate sparser version with input layer such that $p_{1j}^1 \rightarrow 1.0$



Then assuming a Multilayer Perceptron

We have the following Architecture without bias to simplify with a single output

$$\min \frac{1}{N} \sum_{i=1}^N (z_i - t_i)^2$$
$$z_i = \sigma_1 (W_{oh} \mathbf{y}_i)$$
$$\mathbf{y}_i = \sigma_2 (W_{hi} \mathbf{x}_i)$$

Then, we get the following network after the sampling

$$L(W_{oh}, W_{hl}) = (t - z)^2$$
$$z = \sigma_1 \left(W_{oh} \left(r^2 \odot \mathbf{y} \right) \right)$$
$$\mathbf{y} = \sigma_2 \left(W_{hl} \left(r^1 \odot \mathbf{x} \right) \right)$$

Then assuming a Multilayer Perceptron

We have the following Architecture without bias to simplify with a single output

$$\min \frac{1}{N} \sum_{i=1}^N (z_i - t_i)^2$$
$$z_i = \sigma_1 (W_{oh} \mathbf{y}_i)$$
$$\mathbf{y}_i = \sigma_2 (W_{hi} \mathbf{x}_i)$$

Then, we get the following network after the sampling

$$L(W_{oh}, W_{hI}) = (t - z)^2$$
$$z = \sigma_1 \left(W_{oh} \left(\mathbf{r}^2 \odot \mathbf{y} \right) \right)$$
$$\mathbf{y} = \sigma_2 \left(W_{hI} \left(\mathbf{r}^1 \odot \mathbf{x} \right) \right)$$

Then, we have that

The Backpropagation at hidden weights

$$\frac{\partial L}{\partial W_{oh}} = -2(t - z) \times \frac{\partial \sigma'_1(\text{net}_{oh})}{\partial \text{net}_{oh}} \times (\mathbf{r}^2 \odot \mathbf{y})$$

Basically,

$$(W_{oh}^{t+1})_j = \begin{cases} (W_{oh}^t)_j + \eta 2(t - z) \times \frac{\partial \sigma'_1(\text{net}_{oh})}{\partial \text{net}_{oh}} (\mathbf{y})_j & \text{if } r_j = 1 \\ (W_{oh}^t)_j & \text{if } r_j = 0 \end{cases}$$

Then, we have that

The Backpropagation at hidden weights

$$\frac{\partial L}{\partial W_{oh}} = -2(t - z) \times \frac{\partial \sigma'_1(net_{oh})}{\partial net_{oh}} \times (\mathbf{r}^2 \odot \mathbf{y})$$

Basically

$$(W_{oh}^{t+1})_j = \begin{cases} (W_{oh}^t)_j + \eta 2(t - z) \times \frac{\partial \sigma'_1(net_{oh})}{\partial net_{oh}} (\mathbf{y})_j & \text{if } r_j = 1 \\ (W_{oh}^t)_j & \text{if } r_j = 0 \end{cases}$$

However, At Testing

There are a exponential number of possible sparse networks

- A neural net with n units, can be seen as a collection of 2^n possible thinned neural networks.

Assuming

- These networks all share weights so that the total number of parameters is still $O(n^2)$ given that you this many connections

$$\frac{n(n-1)}{2} = O(n^2)$$

Problem: we cannot average such amount of different networks

- We average over the different passes to obtain a p for each node in the network
 - Meaning the probability of being active in the network.

$$p_{ik} = \frac{\text{\#of subnets wehre node } ik \text{ was active}}{\text{\#Of total subnets}}$$

However, At Testing

There are a exponential number of possible sparse networks

- A neural net with n units, can be seen as a collection of 2^n possible thinned neural networks.

Assuming

- These networks all share weights so that the total number of parameters is still $O(n^2)$ given that you this many connections

$$\frac{n(n-1)}{2} = O(n^2)$$

Problem: we cannot average over all possible subnetworks

- We average over the different passes to obtain a p for each node in the network
 - Meaning the probability of being active in the network.

$$p_{ik} = \frac{\text{\#of subnets wehre node } ik \text{ was active}}{\text{\#Of total subnets}}$$

However, At Testing

There are a exponential number of possible sparse networks

- A neural net with n units, can be seen as a collection of 2^n possible thinned neural networks.

Assuming

- These networks all share weights so that the total number of parameters is still $O(n^2)$ given that you this many connections

$$\frac{n(n-1)}{2} = O(n^2)$$

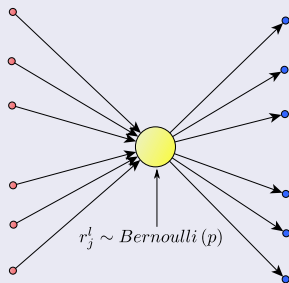
Problem, we cannot average such amount of sub-networks

- We average over the different passes to obtain a p for each node in the network
 - ▶ Meaning the probability of being active in the network.

$$p_{ik} = \frac{\text{\#of subnets wehre node } ik \text{ was active}}{\text{\#Of total subnets}}$$

Then, we have

At Training



The mixture of the models

We know that

$$E(w_{ik}) = \sum_{m=1}^M w_{ik}^m p(w_{ik}^m | \text{BackProp}_M, \mathbf{X})$$

Clearly, we need to get $p(w_{ik}^m | \text{BackProp}_M, \mathbf{X})$

- A simple solution, we can use

$$p_{ik} = \frac{\text{\#of subnets where node } ik \text{ was active}}{\text{\#Of total subnets}}$$

The mixture of the models

We know that

$$E(w_{ik}) = \sum_{m=1}^M w_{ik}^m p(w_{ik}^m | \text{BackProp}_M, \mathbf{X})$$

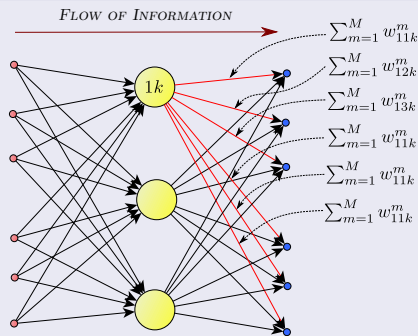
Clearly, we need to get $p(w_{ik}^m | \text{BackProp}_M, \mathbf{X})$

- A simple solution, we can use

$$p_{ik} = \frac{\text{\#of subnets wehre node } ik \text{ was active}}{\text{\#Of total subnets}}$$

Therefore, Using the fact that Forward has a Flow of Information

Add flow of information between all the different generated trained networks



Mathematically

We have the following ideas

- Each node has associated matrices for exit weights

$$W_{out} = \begin{pmatrix} \sum_{i=1}^m w_{i1k}^m \\ \sum_{i=1}^m w_{i2k}^m \\ \vdots \\ \sum_{i=1}^m w_{iJk}^m \end{pmatrix}$$

Then use the probability p to get the new final weights

$$P_{ik} W_{out} = \begin{pmatrix} \sum_{i=1}^m w_{i1k}^m P_{ik} \\ \sum_{i=1}^m w_{i2k}^m P_{ik} \\ \vdots \\ \sum_{i=1}^m w_{iJk}^m P_{ik} \end{pmatrix}$$

Mathematically

We have the following ideas

- Each node has associated matrices for exit weights

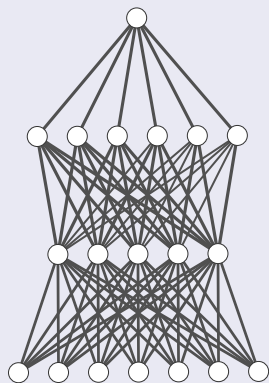
$$W_{out} = \begin{pmatrix} \sum_{i=1}^m w_{i1k}^m \\ \sum_{i=1}^m w_{i2k}^m \\ \vdots \\ \sum_{i=1}^m w_{iJk}^m \end{pmatrix}$$

Then use the probability p to get the new final weights

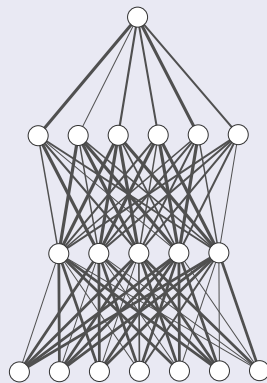
$$p_{ik} W_{out} = \begin{pmatrix} \sum_{i=1}^m w_{i1k}^m p_{ik} \\ \sum_{i=1}^m w_{i2k}^m p_{ik} \\ \vdots \\ \sum_{i=1}^m w_{iJk}^m p_{ik} \end{pmatrix}$$

Then

We have the following structure where thinner lines represent smaller weights



The Original Structure



At Testing

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- **Dropout as Regularization**
 - Introduction
 - Dropout Process
 - **Dropout as Bagging/Bootstrap Aggregation**
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Why dropout?

Srivastava et al. [6]

- A motivation for dropout comes from the theory of evolution!!!
 - ▶ Yes a original network and after a mutated one!!!

The most accepted interpretation of dropout

- It is implicitly bagging at test time a large number of neural networks which share parameters.

Why dropout?

Srivastava et al. [6]

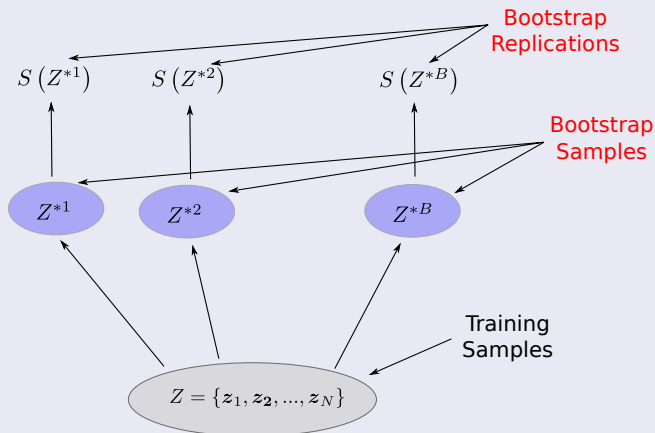
- A motivation for dropout comes from the theory of evolution!!!
 - ▶ Yes a original network and after a mutated one!!!

The most accepted interpretation of dropout

- It is implicitly bagging at test time a large number of neural networks which share parameters.

Bagging/Bootstrap Aggregation

Schematic of the Bootstrap Aggregation process [1]



Thus

Use each of them to train a copy $y_b(\mathbf{x})$ of a predictive regression model to predict a single continuous variable

$$y_{com}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B y_b(\mathbf{x})$$

Results

We have that

Method	CIFAR-10 Error	CIFAR-100 Error
CNN+max pooling (hand tuned)	15.60%	43.48%
CNN+stochastic pooling (Zeiler and Fergus, 2013)	15.13%	42.51%
CNN+max pooling (Snoek et al., 2012)	14.98%	-
CNN+max pooling + dropout fully connected layers	14.32%	41.26%
CNN+max pooling + dropout in all layers	12.61%	37.20%
CNN+maxout (Goodfellow et al., 2013)	11.68%	38.57%

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- **Dropout as Regularization**
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - **Beyond an Empirical Probabilities, LASSO and Data Flow**
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Given the previous ideas

Why not to use the Data Flow for Sparsity?

- Basically, we can assume that a pattern exist in the data you are looking at
 - ▶ The shifts on the weights are not so great...

is to train because it does not represents the real distribution (BackProp, X)

- Actually, you should use the min-batch values, x_t and y_{t+1} , to generate the real distribution

Given the previous ideas

Why not to use the Data Flow for Sparsity?

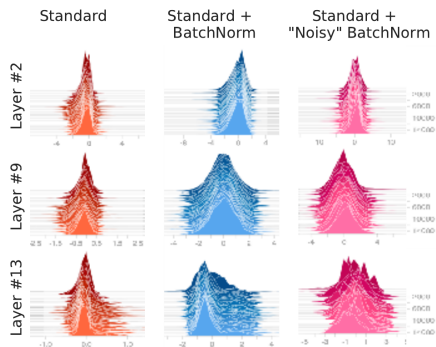
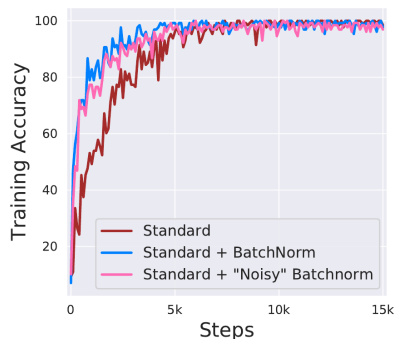
- Basically, we can assume that a pattern exist in the data you are looking at
 - ▶ The shifts on the weights are not so great...

p_{ik} is too broad because it does not represent the real
 $p(w_{ik}^m | \text{BackProp}_M, \mathbf{X})$

- Actually, you should use the min-batch values, \mathbf{x}_t and \mathbf{y}_{t+1} , to generate the real distribution

Based in the paper

"How does batch normalization help optimization?", in Advances in Neural Information Processing Systems (2018), pp. 2483--2493.



Then, we can use a Gaussian Distribution to model this

Actually, the paper is telling us that, given the noise that is injected at each time step t

$$\mu^t \sim U(-n_\mu, n_\mu)$$

$$\sigma^t \sim U(1, n)$$

Something Notable

Properties

Then, we can use a Gaussian Distribution to model this

Actually, the paper is telling us that, given the noise that is injected at each time step t

$$\mu^t \sim U(-n_\mu, n_\mu)$$

$$\sigma^t \sim U(1, n)$$

Something Notable

Properties

Then, we can use a Gaussian Distribution to model this

Actually, the paper is telling us that, given the noise that is injected at each time step t

$$\mu^t \sim U(-n_\mu, n_\mu)$$

$$\sigma^t \sim U(1, n)$$

Something Notable

Properties

Why not use for the Data for enforcing Sparsity?

We have

$$p(\mathbf{y}^{l+1} | \mathbf{x}^l, W) = \mathcal{N}(\sigma(W\mathbf{x}^l), \sigma^2 I)$$

$$p(\sigma^2) \propto \text{"constant"}$$

$$p(W^l | \tau) = \prod_{i=1}^d \mathcal{N}(w_j^l | 0, \tau_j^l) = \mathcal{N}(W^l | 0, (\Upsilon(\tau))^{-1})$$

$$p(\tau | \gamma) = \left(\frac{\gamma}{2}\right)^d \prod_{i=1}^d \exp\left\{-\frac{\gamma}{2} \tau_i\right\}$$

- With $\Upsilon(\tau) = \text{diag}(\tau_1^{-1}, \dots, \tau_d^{-1})$ is the diagonal matrix with the inverse variances of all the w_i 's.

Why not use for the Data for enforcing Sparsity?

We have

$$p(\mathbf{y}^{l+1} | \mathbf{x}^l, W) = \mathcal{N}(\sigma(W\mathbf{x}^l), \sigma^2 I)$$

$$p(\sigma^2) \propto \text{"constant"}$$

$$p(W^l | \tau) = \prod_{i=1}^d \mathcal{N}(w_j^l | 0, \tau_j^l) = \mathcal{N}(W^l | 0, (\Upsilon(\tau))^{-1})$$

$$p(\tau | \gamma) = \left(\frac{\gamma}{2}\right)^d \prod_{i=1}^d \exp\left\{-\frac{\gamma}{2} \tau_i\right\}$$

- With $\Upsilon(\tau) = \text{diag}(\tau_1^{-1}, \dots, \tau_d^{-1})$ is the diagonal matrix with the inverse variances of all the w_i 's.

Why not use for the Data for enforcing Sparsity?

We have

$$p(\mathbf{y}^{l+1} | \mathbf{x}^l, W) = \mathcal{N}(\sigma(W\mathbf{x}^l), \sigma^2 I)$$

$$p(\sigma^2) \propto \text{"constant"}$$

$$p(W^l | \tau) = \prod_{i=1}^d \mathcal{N}(w_j^l | 0, \tau_j^l) = \mathcal{N}(W^l | 0, (\Upsilon(\boldsymbol{\tau}))^{-1})$$

$$p(\boldsymbol{\tau} | \gamma) = \left(\frac{\gamma}{2}\right)^d \prod_{i=1}^d \exp\left\{-\frac{\gamma}{2} \tau_i\right\}$$

- With $\Upsilon(\boldsymbol{\tau}) = \text{diag}(\tau_1^{-1}, \dots, \tau_d^{-1})$ is the diagonal matrix with the inverse variances of all the w_i 's.

Why not use for the Data for enforcing Sparsity?

We have

$$p(\mathbf{y}^{l+1} | \mathbf{x}^l, W) = \mathcal{N}(\sigma(W\mathbf{x}^l), \sigma^2 I)$$

$$p(\sigma^2) \propto \text{"constant"}$$

$$p(W^l | \tau) = \prod_{i=1}^d \mathcal{N}(w_j^l | 0, \tau_j^l) = \mathcal{N}(W^l | 0, (\Upsilon(\boldsymbol{\tau}))^{-1})$$

$$p(\boldsymbol{\tau} | \gamma) = \left(\frac{\gamma}{2}\right)^d \prod_{i=1}^d \exp\left\{-\frac{\gamma}{2}\tau_i\right\}$$

- With $\Upsilon(\boldsymbol{\tau}) = \text{diag}(\tau_1^{-1}, \dots, \tau_d^{-1})$ is the diagonal matrix with the inverse variances of all the w_i 's.

How do we build such distribution

Given that each w_i has a zero-mean Gaussian prior

$$p(w_i|\tau_i) = \mathcal{N}(w_i|0, \tau_i) \quad (10)$$

Where τ_i has the following exponential hyp-prior

$$p(\tau_i|\gamma) = \frac{\gamma}{2} \exp\left\{-\frac{\gamma}{2}\tau_i\right\} \text{ for } \tau_i \geq 0 \quad (11)$$

Then, we have

$$w_i \sim p(w_i|\gamma) = \int_0^\infty p(w_i|\tau_i) p(\tau_i|\gamma) d\tau_i = \frac{\sqrt{\gamma}}{2} \exp\{-\sqrt{\gamma}|w_i|\} \quad (12)$$

How do we build such distribution

Given that each w_i has a zero-mean Gaussian prior

$$p(w_i|\tau_i) = \mathcal{N}(w_i|0, \tau_i) \quad (10)$$

Where τ_i has the following exponential hyper-prior

$$p(\tau_i|\gamma) = \frac{\gamma}{2} \exp\left\{-\frac{\gamma}{2}\tau_i\right\} \text{ for } \tau_i \geq 0 \quad (11)$$

Then, we have

$$w_i \sim p(w_i|\gamma) = \int_0^\infty p(w_i|\tau_i) p(\tau_i|\gamma) d\tau_i = \frac{\sqrt{\gamma}}{2} \exp\{-\sqrt{\gamma}|w_i|\} \quad (12)$$

How do we build such distribution

Given that each w_i has a zero-mean Gaussian prior

$$p(w_i|\tau_i) = \mathcal{N}(w_i|0, \tau_i) \quad (10)$$

Where τ_i has the following exponential hyper-prior

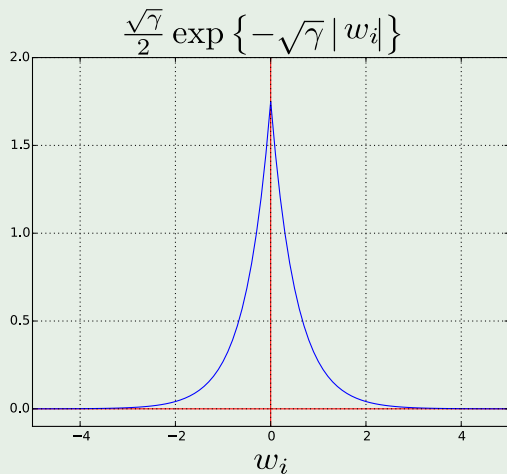
$$p(\tau_i|\gamma) = \frac{\gamma}{2} \exp\left\{-\frac{\gamma}{2}\tau_i\right\} \text{ for } \tau_i \geq 0 \quad (11)$$

Then, we have

$$w_i \sim p(w_i|\gamma) = \int_0^\infty p(w_i|\tau_i) p(\tau_i|\gamma) d\tau_i = \frac{\sqrt{\gamma}}{2} \exp\{-\sqrt{\gamma}|w_i|\} \quad (12)$$

Example

The double exponential



Then using the Monte Carlo Method

We have

$$E [W^t | f (W_b^{tl} \mathbf{x}_b), \sigma^2 I] = \frac{p(\sigma^2)}{B} \sum_{b=1}^B \mathcal{N} (f (W_b^{tl} \mathbf{x}_b), \sigma^2 I) p (W_b^{tl} | \tau_i) p (\tau_i | \gamma)$$

Then, we use the mini batch per epoch to decide if we drop a weight

- Basically, the previous

Then using the Monte Carlo Method

We have

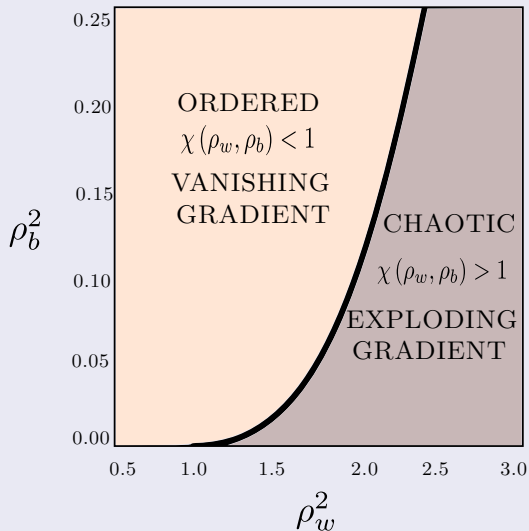
$$E [W^t | f(W_b^{tl} \mathbf{x}_b), \sigma^2 I] = \frac{p(\sigma^2)}{B} \sum_{b=1}^B \mathcal{N}(f(W_b^{tl} \mathbf{x}_b), \sigma^2 I) p(W_b^{tl} | \tau_i) p(\tau_i | \gamma)$$

Then, we use the mini batch per epoch to decide if we drop a weight

- Basically, the previous

We are using the following idea

Basically, we are using the fact that



Thus, we have that

The layer output can be bounded by

$$\mathcal{N}\left(f\left(W_b^{tl}\mathbf{x}_b\right), \sigma^2 I\right)$$

The other part of the equation is the sparsity part:

$$p\left(W_b^{tl}|\tau_i\right)p\left(\tau_i|\gamma\right)$$

Thus, we have that

The layer output can be bounded by

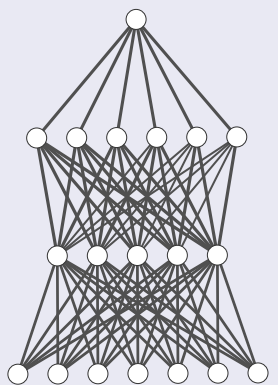
$$\mathcal{N}\left(f\left(W_b^{tl}\mathbf{x}_b\right), \sigma^2 I\right)$$

The other part of the equation is the sparsity part

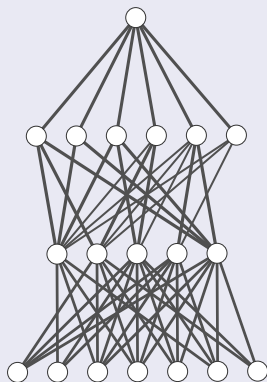
$$p\left(W_b^{tl}|\tau_i\right)p\left(\tau_i|\gamma\right)$$

As the process progress

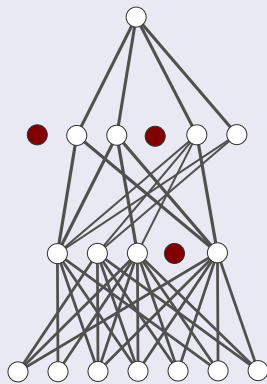
Once the weights fall below certain level we shutdown the weight



The Original Structure



After Some Epochs



More Epochs

Outline

- 1 Bias-Variance Dilemma
 - Introduction
 - Measuring the difference between optimal and learned
 - The Bias-Variance
 - "Extreme" Example
- 2 The Problem with Overfitting
 - Intuition from Overfitting
 - The Idea of Regularization
 - Ridge Regression
 - The LASSO
 - Generalization
 - What can be done?
- 3 Methods of Regularization for Deep Networks
 - Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
 - Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
 - **Random dropout probability**
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
 - Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
 - For More in Normalization

The main goal when using dropout

- It is to regularize the neural network we are training

These random modifications of the network's structure

- They are believed to avoid co-adaptation of neurons by making it impossible for two subsequent neurons to rely solely on each other [6]

The main goal when using dropout

- It is to regularize the neural network we are training

Those random modifications of the network's structure

- They are believed to avoid co-adaptation of neurons by making it impossible for two subsequent neurons to rely solely on each other [6]

Therefore

We have a function that projects from a dimensional space to another

$$h(\mathbf{x}) = W\mathbf{x} + \mathbf{b}$$

Then, given the noisy version of an activation function where $W \rightarrow \tilde{W}$,

$$\tilde{f}(h) = M \odot \text{rect}(h) \quad (\text{Training})$$

- Where $f(h) = \text{rect}(h)$ (Testing)

Actually, Srivastava et al. [6]

- He mentions to use

$$P_{ijk} = \frac{\text{\#of subnets where node } ijk \text{ was active}}{\text{\#Of total subnets}}$$

Therefore

We have a function that projects from a dimensional space to another

$$h(\mathbf{x}) = W\mathbf{x} + \mathbf{b}$$

Then, given the noisy version of an activation function where $M \sim \mathcal{B}(p_h)$

$$\tilde{f}(h) = M \odot \text{rect}(h) \text{ (Training)}$$

- Where $f(h) = \text{rect}(h)$ (Testing)

Actually, Svestava et al. [6]

- He mentions to use

$$p_{ijk} = \frac{\text{\#of subnets where node } ijk \text{ was active}}{\text{\#Of total subnets}}$$

Therefore

We have a function that projects from a dimensional space to another

$$h(\mathbf{x}) = W\mathbf{x} + \mathbf{b}$$

Then, given the noisy version of an activation function where $M \sim \mathcal{B}(p_h)$

$$\tilde{f}(h) = M \odot \text{rect}(h) \text{ (Training)}$$

- Where $f(h) = \text{rect}(h)$ (Testing)

Actually Srivastava et al. [6]

- He mentions to use

$$p_{ijk} = \frac{\text{\#of subnets wehre node } ijk \text{ was active}}{\text{\#Of total subnets}}$$

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- **Random dropout probability**
 - **Projecting Noise into Input Space**
 - Augmenting by Noise
 - Co-adaptation/Overfitting
 - Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
 - For More in Normalization

Data Augmentation

In many previous works [5, 4]

- It has been shown that augmenting data by using domain specific transformations helps in learning better models

Therefore, the main idea:

- It is to map input data to output labels

One way to learn such a mapping function:

- It is to augment the data using noise:
 - ▶ Hypothesis!!! Noise based regularization techniques seems to be increasing training data coverage as augmentation

Data Augmentation

In many previous works [5, 4]

- It has been shown that augmenting data by using domain specific transformations helps in learning better models

Therefore, the main idea

- It is to map input data to output labels

One way to think about data augmentation

- It is to augment the data using noise:
 - ▶ Hypothesis!!! Noise based regularization techniques seems to be increasing training data coverage as augmentation

Data Augmentation

In many previous works [5, 4]

- It has been shown that augmenting data by using domain specific transformations helps in learning better models

Therefore, the main idea

- It is to map input data to output labels

One way to learn such a mapping function

- It is to augment the data using noise:
 - ▶ Hypothesis!!! Noise based regularization techniques seems to be increasing training data coverage as augmentation

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- **Random dropout probability**
 - Projecting Noise into Input Space
 - **Augmenting by Noise**
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Augmenting by Noise [8]

We assume that for a given $\tilde{f}(h)$, there is an optimal \mathbf{x}^*

$$(f \circ h)(\mathbf{x}^*) = \text{rect}(h(\mathbf{x}^*)) \approx M \odot \text{rect}(h) = (\tilde{f} \circ h)(\mathbf{x}^*)$$

This \mathbf{x}^* can be found by minimizing by stochastic gradient descent

$$L(\mathbf{x}, \mathbf{x}^*) = [(f \circ h)(\mathbf{x}^*) - (\tilde{f} \circ h)(\mathbf{x}^*)]^2$$

Augmenting by Noise [8]

We assume that for a given $\tilde{f}(h)$, there is an optimal \mathbf{x}^*

$$(f \circ h)(\mathbf{x}^*) = \text{rect}(h(\mathbf{x}^*)) \approx M \odot \text{rect}(h) = (\tilde{f} \circ h)(\mathbf{x}^*)$$

This \mathbf{x}^* can be found by minimizing by stochastic gradient descent

$$L(\mathbf{x}, \mathbf{x}^*) = \left[(f \circ h)(\mathbf{x}^*) - (\tilde{f} \circ h)(\mathbf{x}^*) \right]^2$$

Extending to n layers

For this, we define

$$\begin{aligned}\tilde{g}^{(i)}(\mathbf{x}) &= \left[\tilde{f}^{(i)} \circ h^{(i)} \circ \dots \circ \tilde{f}^{(1)} \circ h^{(1)} \right] (\mathbf{x}) \\ g^{(i)}(\mathbf{x}^*) &= \left[f^{(i)} \circ h^{(i)} \circ \dots \circ f^{(1)} \circ h^{(1)} \right] (\mathbf{x}^*)\end{aligned}$$

Then, it is possible to compute the back-propagation projection corresponding to all hidden layer activations at once

$$L(\mathbf{x}, \mathbf{x}^{(1)*}, \dots, \mathbf{x}^{(n)*}) = \sum_{i=1}^n \lambda_i \left[g^{(i)}(\mathbf{x}^{(i)*}) - \tilde{g}^{(i)}(\mathbf{x}) \right]^2$$

Extending to n layers

For this, we define

$$\begin{aligned}\tilde{g}^{(i)}(\mathbf{x}) &= \left[\tilde{f}^{(i)} \circ h^{(i)} \circ \dots \circ \tilde{f}^{(1)} \circ h^{(1)} \right] (\mathbf{x}) \\ g^{(i)}(\mathbf{x}^*) &= \left[f^{(i)} \circ h^{(i)} \circ \dots \circ f^{(1)} \circ h^{(1)} \right] (\mathbf{x}^*)\end{aligned}$$

Then, it is possible to compute the back propagation projection corresponding to all hidden layer activations at once

$$L(\mathbf{x}, \mathbf{x}^{(1)*}, \dots, \mathbf{x}^{(n)*}) = \sum_{i=1}^n \lambda_i \left[g^{(i)}(\mathbf{x}^{(i)*}) - \tilde{g}^{(i)}(\mathbf{x}) \right]^2$$

However

Small Problem

- It is possible to show by contradiction that one is unlikely to find a single $\mathbf{x}^* = \mathbf{x}^{(1)*} = \dots = \mathbf{x}^{(n)*}$
 - ▶ Such that you can significantly reduce L

Proof of the unlikeness of $\mathbf{x}^* = \mathbf{x}^{(1)*} = \dots = \mathbf{x}^{(n)*}$

By the associative property of function composition

$$g^{(i)}(\mathbf{x}^*) = (f^{(i)} \circ h^{(i)})(g^{(i-1)}(\mathbf{x}^*))$$

Suppose there exist $i < n$ and $\mathbf{x} \in \mathbb{R}^n$ such that

$$\begin{aligned} (f^{(i)} \circ h^{(i)})(g^{(i-1)}(\mathbf{x}^*)) &= (\bar{f}^{(i)} \circ h^{(i)})(\bar{g}^{(i-1)}(\mathbf{x})) \\ (f^{(i-1)} \circ h^{(i-1)})(g^{(i-2)}(\mathbf{x}^*)) &= (\bar{f}^{(i-1)} \circ h^{(i-1)})(\bar{g}^{(i-2)}(\mathbf{x})) \end{aligned}$$

Proof of the unlikeness of $\mathbf{x}^* = \mathbf{x}^{(1)*} = \dots = \mathbf{x}^{(n)*}$

By the associative property of function composition

$$g^{(i)}(\mathbf{x}^*) = (f^{(i)} \circ h^{(i)})(g^{(i-1)}(\mathbf{x}^*))$$

Suppose there exist $\mathbf{x}^* = \mathbf{x}^{(1)*} = \dots = \mathbf{x}^{(n)*}$ an such that

$$\begin{aligned}(f^{(i)} \circ h^{(i)})(g^{(i-1)}(\mathbf{x}^*)) &= (\tilde{f}^{(i)} \circ h^{(i)})(\tilde{g}^{(i-1)}(\mathbf{x})) \\ (f^{(i-1)} \circ h^{(i-1)})(g^{(i-2)}(\mathbf{x}^*)) &= (\tilde{f}^{(i-1)} \circ h^{(i-1)})(\tilde{g}^{(i-2)}(\mathbf{x}))\end{aligned}$$

Then

Based on the previous equations

$$g^{(i-1)}(\mathbf{x}^*) = \tilde{g}^{(i-1)}(\mathbf{x})$$

Then, we get

$$\left(f^{(i)} \circ h^{(i)}\right)\left(g^{(i-1)}\left(\mathbf{x}^*\right)\right)=\left(\tilde{f}^{(i)} \circ h^{(i)}\right)\left(\tilde{g}^{(i-1)}(\mathbf{x})\right)$$

Finally

$$\text{rect}\left(h^{(i)}\left(g^{(i-1)}\left(\mathbf{x}^*\right)\right)\right)=M^{(i)} \odot \text{rect}\left(h^{(i)}\left(\tilde{g}^{(i-1)}(\mathbf{x})\right)\right)$$

Then

Based on the previous equations

$$g^{(i-1)}(\mathbf{x}^*) = \tilde{g}^{(i-1)}(\mathbf{x})$$

Then, we get

$$(f^{(i)} \circ h^{(i)})(g^{(i-1)}(\mathbf{x}^*)) = (\tilde{f}^{(i)} \circ h^{(i)})(\tilde{g}^{(i-1)}(\mathbf{x}))$$

Finally

$$\text{rect}(h^{(i)}(g^{(i-1)}(\mathbf{x}^*))) = M^{(i)} \odot \text{rect}(h^{(i)}(\tilde{g}^{(i-1)}(\mathbf{x})))$$

Then

Based on the previous equations

$$g^{(i-1)}(\mathbf{x}^*) = \tilde{g}^{(i-1)}(\mathbf{x})$$

Then, we get

$$(f^{(i)} \circ h^{(i)}) (g^{(i-1)}(\mathbf{x}^*)) = (\tilde{f}^{(i)} \circ h^{(i)}) (\tilde{g}^{(i-1)}(\mathbf{x}))$$

Finally

$$\text{rect} \left(h^{(i)} \left(g^{(i-1)}(\mathbf{x}^*) \right) \right) = M^{(i)} \odot \text{rect} \left(h^{(i)} \left(\tilde{g}^{(i-1)}(\mathbf{x}) \right) \right)$$

Therefore

This is only true if $M^{(i)} = 1$

- When $\text{rect}_j \left(h^{(i)} \left(g^{(i-1)} \left(\mathbf{x}^* \right) \right) \right) > 0$

This only happens with a probability of

- Where:
 - ▶ $p_{(i)}$ is the Bernoulli success probability.
 - ▶ $d_{(i)}$ is the number of hidden units.
 - ▶ $s_{(i)}$ is the mean sparsity level at i (Mean percentage of active hidden units).

Therefore

This is only true if $M^{(i)} = 1$

- When $\text{rect}_j \left(h^{(i)} \left(g^{(i-1)} (\mathbf{x}^*) \right) \right) > 0$

This only happens with a probability $p_{(i)}^{d_{(i)} s_{(i)}}$

- Where:
 - ▶ $p_{(i)}$ is the Bernoulli success probability.
 - ▶ $d_{(i)}$ is the number of hidden units.
 - ▶ $s_{(i)}$ is the mean sparsity level at i (Mean percentage of active hidden units).

Which is quite low!!!

This probability is very low for standard hyper-parameters values

- With $p_{(i)} = 0.5$, $d_{(i)} = 1000$ and $s_{(i)} = 0.15$

$$p_{(i)}^{d_{(i)} s_{(i)}} = 10^{-47}$$

However

Fortunately

- It is easy to find a different x^* for each hidden layer

by providing multiple inputs

$$\left(x, x^{(1)*}, x^{(2)*}, \dots, x^{(n)*}\right)$$

However

- This raises the question whether we can train the network deterministically on the $x^{(i)*}$ instead of using dropout

However

Fortunately

- It is easy to find a different \mathbf{x}^* for each hidden layer

by providing multiple inputs

$$\left(\mathbf{x}, \mathbf{x}^{(1)*}, \mathbf{x}^{(2)*}, \dots, \mathbf{x}^{(n)*} \right)$$

However

- This raises the question whether we can train the network deterministically on the $\mathbf{x}^{(i)*}$ instead of using dropout

However

Fortunately

- It is easy to find a different \mathbf{x}^* for each hidden layer

by providing multiple inputs

$$\left(\mathbf{x}, \mathbf{x}^{(1)*}, \mathbf{x}^{(2)*}, \dots, \mathbf{x}^{(n)*} \right)$$

However

- This raises the question whether we can train the network deterministically on the $\mathbf{x}^{(i)*}$ instead of using dropout

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- **Random dropout probability**
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - **Co-adaptation/Overfitting**
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Co-adaptation/Overfitting

Definition

- Co-adaptation is the accumulation of interacting genes in the gene pool of a population by selection.
 - ▶ Selection pressures on one of the genes will affect its interacting proteins, after which compensatory changes occur.

In Neural Networks

- In neural network, co-adaptation means that some neurons are highly dependent on others:
 - ▶ Getting into over-fitting!!!

Co-adaptation/Overfitting

Definition

- Co-adaptation is the accumulation of interacting genes in the gene pool of a population by selection.
 - ▶ Selection pressures on one of the genes will affect its interacting proteins, after which compensatory changes occur.

In Neural Networks

- In neural network, co-adaptation means that some neurons are highly dependent on others:
 - ▶ Getting into over-fitting!!!

Question

We have that

- Question: Can we train the network deterministically on $\mathbf{x}^{(i)*}$?

Question

We have that

- Question: Can we train the network deterministically on $\mathbf{x}^{(i)*}$?

This is not trivial given that

- Dropout is not effectively applied to every layer at the same time when using

$$\left(\mathbf{x}, \mathbf{x}^{(1)*}, \mathbf{x}^{(2)*}, \dots, \mathbf{x}^{(n)*}\right)$$

- The gradients of the linear projections will differ greatly, different from dropout!!!

Question

We have that

- Question: Can we train the network deterministically on $\mathbf{x}^{(i)*}$?

This is not trivial given that

- Dropout is not effectively applied to every layer at the same time when using

$$\left(\mathbf{x}, \mathbf{x}^{(1)*}, \mathbf{x}^{(2)*}, \dots, \mathbf{x}^{(n)*}\right)$$

- The gradients of the linear projections will differ greatly, different from dropout!!!

Therefore

We can then

- Modifying the probability distribution is the most straightforward way to improve the set of transformations.

For example:

- A simple way to vary the transformation magnitude randomly is to replace p_{hi} by a random variable!!!

Therefore

We can then

- Modifying the probability distribution is the most straightforward way to improve the set of transformations.

For example

- A simple way to vary the transformation magnitude randomly is to replace p_{hij} by a random variable!!!

Therefore

Define

$$M_{hij} \sim \mathcal{B}(\rho_h) \text{ (Bernoulli)}$$
$$\rho_h \sim U(0, p_h) \text{ (Uniform)}$$

- where h defines the layer, i the sample, and j the layer's neuron.

Here, the authors use the same ρ for all the layers of the neurons, then

$$\tilde{f}(h) = \frac{1}{1-\rho} M \odot \text{rect}(h)$$

Therefore

Define

$$M_{hij} \sim \mathcal{B}(\rho_h) \text{ (Bernoulli)}$$
$$\rho_h \sim U(0, p_h) \text{ (Uniform)}$$

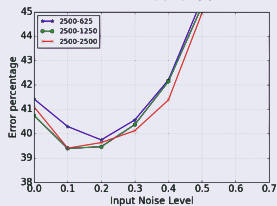
- where h defines the layer, i the sample, and j the layer's neuron.

Here, the authors use the same ρ for all the layers of the neurons, then

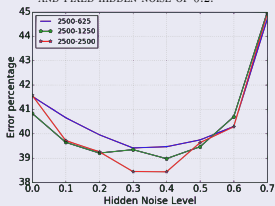
$$\tilde{f}(h) = \frac{1}{1-\rho} M \odot \text{rect}(h)$$

Something Notable

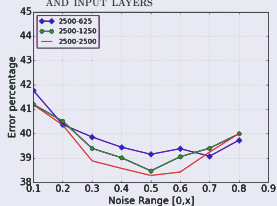
USING DROPOUT WITH VARYING INPUT NOISE
AND FIXED HIDDEN NOISE OF 0.5.



USING DROPOUT WITH VARYING INPUT NOISE
AND FIXED HIDDEN NOISE OF 0.2.



USING RANDOM-DROPOUT WITH VARYING
NOISE RANGE $[0, x]$ USED AT HIDDEN
AND INPUT LAYERS



Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- **Batch normalization**
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- For More in Normalization

Here, the people at Google [9] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer’s input

What claim

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

In Neural Networks, they define this

- Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training.

Here, the people at Google [9] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer’s input

They claim

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

In Neural Networks they define this

- Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training.

Here, the people at Google [9] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer’s input

They claim

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

In Neural Networks, they define this

- Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training.

Transformation

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

- ① $\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$
- ② $\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$
- ③ $\hat{x} = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$
- ④ $y_i = \gamma^{(k)} \hat{x}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Transformation

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

① $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$

② $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu_{\mathcal{B}})^2$

③ $\hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$

④ $y_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Transformation

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

① $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$

② $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu_{\mathcal{B}})^2$

③ $\hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$

④ $y_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Transformation

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

- 1 $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$
- 2 $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu_{\mathcal{B}})^2$
- 3 $\hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$

$$y_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$$

Transformation

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

- 1 $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$
- 2 $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu_{\mathcal{B}})^2$
- 3 $\hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$
- 4 $\mathbf{y}_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Outline

- 1 Bias-Variance Dilemma
 - Introduction
 - Measuring the difference between optimal and learned
 - The Bias-Variance
 - "Extreme" Example
- 2 The Problem with Overfitting
 - Intuition from Overfitting
 - The Idea of Regularization
 - Ridge Regression
 - The LASSO
 - Generalization
 - What can be done?
- 3 Methods of Regularization for Deep Networks
 - Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
 - Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
 - Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
 - **Batch normalization**
 - **Improving the Google Layer Normalization**
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
 - For More in Normalization

Remember

Using Min-Batch inputs, we have

$$\mu_B = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

And Variance

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu_B)^2$$

Remember

Using Min-Batch inputs, we have

$$\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

And Variance

$$\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu_{\mathcal{B}})^2$$

Therefore, Ba et al. [10]

We get the mean over the output of the layer l with H number of hidden units

$$\mu^l = \frac{1}{H} \sum_{i=1}^H y_i^l$$

- Basically, do the forward process then add over the output $y_i^l = w_i^{lT} h^l$ where $h_i^{l+1} = f(y_i^l + b_i^l)$

Then the standard deviation layer:

$$\sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^l - \mu^l)^2}$$

Therefore, Ba et al. [10]

We get the mean over the output of the layer l with H number of hidden units

$$\mu^l = \frac{1}{H} \sum_{i=1}^H y_i^l$$

- Basically, do the forward process then add over the output $y_i^l = w_i^{lT} h^l$ where $h_i^{l+1} = f(y_i^l + b_i^l)$

Then the standard deviation layer l

$$\sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^l - \mu^l)^2}$$

We have that

- All the hidden units in a layer share the same normalization terms μ and σ
 - ▶ but different training cases have different normalization terms.

Layer normalization does not impose any constraint

- On the size of a mini-batch and it can be used in the pure on-line regime with batch size 1

Remarks

We have that

- All the hidden units in a layer share the same normalization terms μ and σ
 - ▶ but different training cases have different normalization terms.

Layer normalization does not impose any constraint

- On the size of a mini-batch and it can be used in the pure on-line regime with batch size 1.

Outline

- 1 Bias-Variance Dilemma
 - Introduction
 - Measuring the difference between optimal and learned
 - The Bias-Variance
 - "Extreme" Example
- 2 The Problem with Overfitting
 - Intuition from Overfitting
 - The Idea of Regularization
 - Ridge Regression
 - The LASSO
 - Generalization
 - What can be done?
- 3 Methods of Regularization for Deep Networks
 - Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
 - Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
 - Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
 - **Batch normalization**
 - Improving the Google Layer Normalization
 - **Layer Normalization in RNN**
 - Invariance Under Weights and Data Transformations
 - For More in Normalization

The Flow of Information through time

First, the new \mathbf{h}^t with a gain vector \mathbf{g}

$$\mathbf{h}^t = f \left[\frac{\mathbf{g}}{\sigma^t} \odot (\mathbf{y}^t - \mu^t) + b \right]$$

The Temporal Layer Mean Normalization

$$\mu^t = \frac{1}{H} \sum_{i=1}^H y_i^t$$

The Temporal Layer STD Normalization

$$\sigma^t = \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^t - \mu^t)^2}$$

The Flow of Information through time

First, the new h^t with a gain vector g

$$h^t = f \left[\frac{g}{\sigma^t} \odot (y^t - \mu^t) + b \right]$$

The Temporal Layer Mean Normalization

$$\mu^t = \frac{1}{H} \sum_{i=1}^H y_i^t$$

The Temporal Layer STD Normalization

$$\sigma^t = \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^t - \mu^t)^2}$$

The Flow of Information through time

First, the new \mathbf{h}^t with a gain vector \mathbf{g}

$$\mathbf{h}^t = f \left[\frac{\mathbf{g}}{\sigma^t} \odot (\mathbf{y}^t - \mu^t) + b \right]$$

The Temporal Layer Mean Normalization

$$\mu^t = \frac{1}{H} \sum_{i=1}^H y_i^t$$

The Temporal Layer STD Normalization

$$\sigma^t = \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^t - \mu^t)^2}$$

Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- **Batch normalization**
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - **Invariance Under Weights and Data Transformations**
- For More in Normalization

Weight re-scaling and re-centering

Observe that under batch normalization and weight normalization

- Any re-scaling to the incoming weights w_i of a single neuron has no effect on the normalized summed inputs to a neuron.

Meaning

- If the weight vector is scaled by δ_i the two scalars μ and σ will also be scaled by δ

Properties

- The batch and weight normalization are invariant to the re-scaling of the weights.

Weight re-scaling and re-centering

Observe that under batch normalization and weight normalization

- Any re-scaling to the incoming weights w_i of a single neuron has no effect on the normalized summed inputs to a neuron.

Meaning

- If the weight vector is scaled by δ_i the two scalars μ and σ will also be scaled by δ

Properties

- The batch and weight normalization are invariant to the re-scaling of the weights.

Weight re-scaling and re-centering

Observe that under batch normalization and weight normalization

- Any re-scaling to the incoming weights w_i of a single neuron has no effect on the normalized summed inputs to a neuron.

Meaning

- If the weight vector is scaled by δ_i the two scalars μ and σ will also be scaled by δ

Properties

- The batch and weight normalization are invariant to the re-scaling of the weights.

In the other hand

Layer normalization

- It is not invariant to the individual scaling of the single weight vectors.

In the other hand

Layer normalization

- It is not invariant to the individual scaling of the single weight vectors.

However

- Layer normalization is invariant to scaling of the entire weight matrix.
- Also it is invariant to a shift to all of the incoming weights in the weight matrix.

In the other hand

Layer normalization

- It is not invariant to the individual scaling of the single weight vectors.

However

- Layer normalization is invariant to scaling of the entire weight matrix.
- Also it is invariant to a shift to all of the incoming weights in the weight matrix.

How?

Imagine the following

- Let there be two sets of model parameters θ , θ' with weigh matrices

$$W' = \delta W + 1\gamma^T$$

We have

Given that $y_i^l = w_i^{lT} \mathbf{x}^l$

$$y_i^l = (\delta W + 1\gamma^T)_i \mathbf{x}^l$$

Then, we have

$$\mu^l = \frac{\delta}{H} \sum_{i=1}^H W_i \mathbf{x}^l + \frac{1}{H} \sum_{i=1}^H (1\gamma^T)_i \mathbf{x}^l = \delta\mu + (1\gamma^T)_i \mathbf{x}^l$$

We have

Given that $y_i^l = w_i^{lT} \mathbf{x}^l$

$$y_i^l = (\delta W + 1\gamma^T)_i \mathbf{x}^l$$

Then, we have

$$\mu^l = \frac{\delta}{H} \sum_{i=1}^H W_i \mathbf{x}^l + \frac{1}{H} \sum_{i=1}^H (1\gamma^T)_i \mathbf{x}^l = \delta\mu + (1\gamma^T)_i \mathbf{x}^l$$

Standard Deviation

$$\sigma' = \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^l - \mu')^2} = \delta \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^l - \mu)^2}$$

Finally, Under Layer Normalization, we have the same output

$$\begin{aligned} h' &= f \left[\frac{g}{\sigma'} (W'x - \mu') + b \right] \\ &= f \left[\frac{g}{\sigma'} \left([\delta W + 1\gamma^T] x - \mu' \right) + b \right] \\ &= f \left[\frac{g}{\sigma} (Wx - \mu) + b \right] = h \end{aligned}$$

Now

Standard Deviation

$$\sigma' = \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^l - \mu')^2} = \delta \sqrt{\frac{1}{H} \sum_{i=1}^H (y_i^l - \mu)^2}$$

Finally, Under Layer Normalization, we have the same output

$$\begin{aligned} \mathbf{h}' &= f \left[\frac{\mathbf{g}}{\sigma'} (W' \mathbf{x} - \mu') + \mathbf{b} \right] \\ &= f \left[\frac{\mathbf{g}}{\sigma'} \left([\delta W + 1\gamma^T] \mathbf{x} - \mu' \right) + \mathbf{b} \right] \\ &= f \left[\frac{\mathbf{g}}{\sigma} (W \mathbf{x} - \mu) + \mathbf{b} \right] = \mathbf{h} \end{aligned}$$

Something Notable

- if normalization is only applied to the input before the weights, the model will not be invariant to re-scaling and re-centering of the weights.

Data re-scaling and re-centering

We can show

- All the normalization methods are invariant to re-scaling the dataset

Layer normalization is invariant to re-scaling of individual training cases

$$h'_i = f \left[\frac{g_i}{\sigma'} \left(w_i^T x' - \mu' \right) + b_i \right] = f \left[\frac{g_i}{\delta \sigma} \left(\delta w_i^T x - \delta \mu \right) + b_i \right] = h_i$$

Data re-scaling and re-centering

We can show

- All the normalization methods are invariant to re-scaling the dataset

Layer normalization is invariant to re-scaling of individual training cases

$$h'_i = f \left[\frac{g_i}{\sigma'} \left(w_i^T \mathbf{x}' - \mu' \right) + b_i \right] = f \left[\frac{g_i}{\delta \sigma} \left(\delta w_i^T \mathbf{x} - \delta \mu \right) + b_i \right] = h_i$$

Additionally

Layer Normalization has a relation with the Fisher Information Matrix

$$F(\theta) = E_{\mathbf{x} \sim P(\mathbf{x}), y \sim P(y|\mathbf{x})} \left[\frac{\partial \log P(y|\mathbf{x})}{\partial \theta} \left(\frac{\partial \log P(y|\mathbf{x})}{\partial \theta} \right)^T \right]$$

Basically, we can write the generalized linear model as

$$\log P(y|\mathbf{x}, w, b) = \frac{(a+b)y - \eta(a+b)}{\Phi} + c(y, \Phi)$$

$$E[y|\mathbf{x}] = f(a+b) = f(w^T \mathbf{x} + b)$$

$$\text{Var}[y|\mathbf{x}] = \Phi f'(a+b)$$

Additionally

Layer Normalization has a relation with the Fisher Information Matrix

$$F(\theta) = E_{\mathbf{x} \sim P(\mathbf{x}), y \sim P(y|\mathbf{x})} \left[\frac{\partial \log P(y|\mathbf{x})}{\partial \theta} \left(\frac{\partial \log P(y|\mathbf{x})}{\partial \theta} \right)^T \right]$$

Basically, we can write the generalized linear model as

$$\log P(y|\mathbf{x}, w, b) = \frac{(a + b)y - \eta(a + b)}{\Phi} + c(y, \Phi)$$

$$E[y|\mathbf{x}] = f(a + b) = f(w^T \mathbf{x} + b)$$

$$\text{Var}[y|\mathbf{x}] = \Phi f'(a + b)$$

The curvature of a Riemannian manifold

It is entirely captured by its Riemannian metric

$$ds^2 \approx \frac{1}{2} \delta^T F(\theta) \delta$$

- where, δ is a small change to the parameters.

Then, under Layer Normalization, we have

$$F(\theta) = \frac{1}{\sigma^2} E_{x \sim \mathcal{P}(x)} \begin{bmatrix} \text{Cov}(y_1, y_2 | x) \frac{(\alpha_1 - \mu)^2}{\sigma^2} & \dots & \text{Cov}(y_1, y_H | x) \frac{(\alpha_1 - \mu)(\alpha_H - \mu)}{\sigma^2} \\ \vdots & \ddots & \vdots \\ \text{Cov}(y_H, y_1 | x) \frac{(\alpha_1 - \mu)(\alpha_H - \mu)}{\sigma^2} & \dots & \text{Cov}(y_H, y_H | x) \frac{(\alpha_H - \mu)^2}{\sigma^2} \end{bmatrix}$$

The curvature of a Riemannian manifold

It is entirely captured by its Riemannian metric

$$ds^2 \approx \frac{1}{2} \delta^T F(\theta) \delta$$

- where, δ is a small change to the parameters.

Then, under Layer Normalization, we have

$$F(\theta) = \frac{1}{\Phi^2} E_{x \sim P(x)} \begin{bmatrix} \text{Cov}(y_1, y_2 | \mathbf{x}) \frac{(a_1 - \mu)^2}{\sigma^2} & \cdots & \text{Cov}(y_1, y_H | \mathbf{x}) \frac{(a_1 - \mu)(a_H - \mu)}{\sigma^2} \\ \vdots & \ddots & \vdots \\ \text{Cov}(y_H, y_1 | \mathbf{x}) \frac{(a_1 - \mu)(a_H - \mu)}{\sigma^2} & \cdots & \text{Cov}(y_H, y_H | \mathbf{x}) \frac{(a_H - \mu)^2}{\sigma^2} \end{bmatrix}$$

Where

We have that $a_i = w_i^T \mathbf{x}$

- We project the gradient updates to the gain parameter δ_{gi} of the i^{th} neuron to its weight vector as

$$\frac{\delta_{gi} \delta_{gj}}{2\Phi^2} E_{x \sim P(\mathbf{x})} \left[Cov(y_i, y_j | \mathbf{x}) \frac{(a_1 - \mu)(a_H - \mu)}{\sigma^2} \right]$$

Essentially

- We have that the normalization layer is more robust to the scaling of the input and parameters

Where

We have that $a_i = w_i^T \mathbf{x}$

- We project the gradient updates to the gain parameter δ_{gi} of the i^{th} neuron to its weight vector as

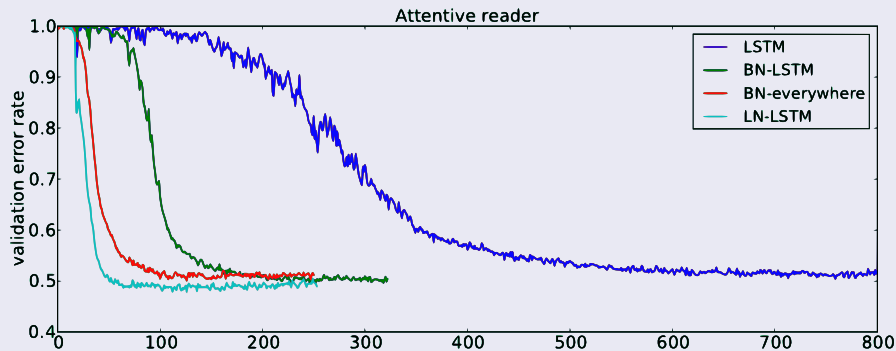
$$\frac{\delta_{gi} \delta_{gj}}{2\Phi^2} E_{x \sim P(\mathbf{x})} \left[\text{Cov}(y_i, y_j | \mathbf{x}) \frac{(a_1 - \mu)(a_H - \mu)}{\sigma^2} \right]$$

Basically

- We have that the normalization layer is more robust to the scaling of the input and parameters

Results

In a LSTM



Outline

1 Bias-Variance Dilemma

- Introduction
- Measuring the difference between optimal and learned
- The Bias-Variance
- "Extreme" Example

2 The Problem with Overfitting

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- The LASSO
- Generalization
- What can be done?

3 Methods of Regularization for Deep Networks

- Gaussian Noise on Hidden Units for Regularization
 - Application into a Decoder/Encoder
- Dropout as Regularization
 - Introduction
 - Dropout Process
 - Dropout as Bagging/Bootstrap Aggregation
 - Beyond an Empirical Probabilities, LASSO and Data Flow
- Random dropout probability
 - Projecting Noise into Input Space
 - Augmenting by Noise
 - Co-adaptation/Overfitting
- Batch normalization
 - Improving the Google Layer Normalization
 - Layer Normalization in RNN
 - Invariance Under Weights and Data Transformations
- **For More in Normalization**

We have the following paper

Please Take a Look

- Kukačka, J., Golkov, V., & Cremers, D. (2017). Regularization for deep learning: A taxonomy. arXiv preprint arXiv:1710.10686.

Conclusions

There is still a lot to understand on the Deep Learning Architectures

- The Last 10 years have shown us a lot on the need of regularization...

Therefore

- When connecting with the paper
 - ▶ "How Does Batch Normalization Help Optimization?" by Santurkar, Tsipras, Ilyas and Madry

We have this if we were able to connect these normalizations

- With the building of the Jacobian on the Gradient Descent, we could improve
 - ▶ The speed of optimization + The regularization properties of such Gradient Descent

Conclusions

There is still a lot to understand on the Deep Learning Architectures

- The Last 10 years have shown us a lot on the need of regularization...

Therefore

- When connecting with the paper
 - ▶ “How Does Batch Normalization Help Optimization?” by Santurkar, Tsipras, Ilyas and Madry

Well, even if we were able to compute the Hessian-normalization

- With the building of the Jacobian on the Gradient Descent, we could improve
 - ▶ The speed of optimization + The regularization properties of such Gradient Descent

Conclusions

There is still a lot to understand on the Deep Learning Architectures






- The Last 10 years have shown us a lot on the need of regularization...






Therefore

- When connecting with the paper
 - ▶ “How Does Batch Normalization Help Optimization?” by Santurkar, Tsipras, Ilyas and Madry

We have the if we were able to connect these normalizations

- With the building of the Jacobian on the Gradient Descent, we could improve
 - ▶ The speed of optimization + The regularization properties of such Gradient Descent

-  T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. Springer Series in Statistics, Springer New York, 2009.
-  S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*. Academic Press, 1st ed., 2015.
-  H. Zou, T. Hastie, and R. Tibshirani, “Sparse principal component analysis,” *Journal of computational and graphical statistics*, vol. 15, no. 2, pp. 265–286, 2006.
-  Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, *et al.*, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
-  T. DeVries and G. W. Taylor, “Dataset augmentation in feature space,” *arXiv preprint arXiv:1702.05538*, 2017.

-  N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
-  S. Wager, S. Wang, and P. S. Liang, “Dropout training as adaptive regularization,” in *Advances in neural information processing systems*, pp. 351–359, 2013.
-  X. Bouthillier, K. Konda, P. Vincent, and R. Memisevic, “Dropout as data augmentation,” *arXiv preprint arXiv:1506.08700*, 2015.
-  S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” *arXiv preprint arXiv:1502.03167*, 2015.
-  J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.