

Deep Learning Book

Chapter 7

Regularization for Deep Learning

Botian Shi

botianshi@bit.edu.cn

March 7, 2017

You can download the \LaTeX source code of this file from [Here](#).

Generalization and Strategy

- How to make an algorithm that will perform well not just on the training data, but also on new inputs?
- Many strategies designed to reduce the test error, possibly at the expense of increased training error.
- These strategies are known collectively as **regularization**.
- Many regularization algorithm have been developed.
- Developing more effective regularization strategies is one of the major research efforts in the field.
- In this chapter, we describe regularization in more detail, focusing on regularization strategies for deep models or models that may be used as building blocks to form deep models.

Generalization and Strategy

- How to make an algorithm that will perform well not just on the training data, but also on new inputs?
- Many strategies designed to reduce the test error, possibly at the expense of increased training error.
- These strategies are known collectively as **regularization**.
- Many regularization algorithm have been developed.
- Developing more effective regularization strategies is one of the major research efforts in the field.
- In this chapter, we describe regularization in more detail, focusing on regularization strategies for deep models or models that may be used as building blocks to form deep models.

Generalization and Strategy

- How to make an algorithm that will perform well not just on the training data, but also on new inputs?
- Many strategies designed to reduce the test error, possibly at the expense of increased training error.
- These strategies are known collectively as **regularization**.
- Many regularization algorithm have been developed.
- Developing more effective regularization strategies is one of the major research efforts in the field.
- In this chapter, we describe regularization in more detail, focusing on regularization strategies for deep models or models that may be used as building blocks to form deep models.

Generalization and Strategy

- How to make an algorithm that will perform well not just on the training data, but also on new inputs?
- Many strategies designed to reduce the test error, possibly at the expense of increased training error.
- These strategies are known collectively as **regularization**.
- Many regularization algorithm have been developed.
- Developing more effective regularization strategies is one of the major research efforts in the field.
- In this chapter, we describe regularization in more detail, focusing on regularization strategies for deep models or models that may be used as building blocks to form deep models.

Generalization and Strategy

- How to make an algorithm that will perform well not just on the training data, but also on new inputs?
- Many strategies designed to reduce the test error, possibly at the expense of increased training error.
- These strategies are known collectively as **regularization**.
- Many regularization algorithm have been developed.
- Developing more effective regularization strategies is one of the major research efforts in the field.
- In this chapter, we describe regularization in more detail, focusing on regularization strategies for deep models or models that may be used as building blocks to form deep models.

Generalization and Strategy

- There are many regularization strategies.
 1. Put extra constraints on a machine learning model. (Adding restrictions on the parameter values.)
 2. Add extra terms in the objective function that can be thought of as corresponding to a soft constraint on the parameter values.
- If chosen carefully, these extra constraints and penalties can lead to improved performance on the test set.
- Sometimes these constraints and penalties are designed to
 1. **encode** specific kinds of **prior knowledge**.
 2. Express a generic preference for a simpler model class in order to promote generalization.
 3. make an under-determined problem determined. (Provide more information)
- Other forms of regularization, known as ensemble methods, combine multiple hypotheses that explain the training data.

Generalization and Strategy

- There are many regularization strategies.
 1. Put extra constraints on a machine learning model. (Adding restrictions on the parameter values.)
 2. Add extra terms in the objective function that can be thought of as corresponding to a soft constraint on the parameter values.
- If chosen carefully, these extra constraints and penalties can lead to improved performance on the test set.
- Sometimes these constraints and penalties are designed to
 1. **encode** specific kinds of **prior knowledge**.
 2. Express a generic preference for a simpler model class in order to promote generalization.
 3. make an under-determined problem determined. (Provide more information)
- Other forms of regularization, known as ensemble methods, combine multiple hypotheses that explain the training data.

Generalization and Strategy

- There are many regularization strategies.
 1. Put extra constraints on a machine learning model. (Adding restrictions on the parameter values.)
 2. Add extra terms in the objective function that can be thought of as corresponding to a soft constraint on the parameter values.
- If chosen carefully, these extra constraints and penalties can lead to improved performance on the test set.
- Sometimes these constraints and penalties are designed to
 1. encode specific kinds of prior knowledge.
 2. Express a generic preference for a simpler model class in order to promote generalization.
 3. make an under-determined problem determined. (Provide more information)
- Other forms of regularization, known as ensemble methods, combine multiple hypotheses that explain the training data.

Generalization and Strategy

- There are many regularization strategies.
 1. Put extra constraints on a machine learning model. (Adding restrictions on the parameter values.)
 2. Add extra terms in the objective function that can be thought of as corresponding to a soft constraint on the parameter values.
- If chosen carefully, these extra constraints and penalties can lead to improved performance on the test set.
- Sometimes these constraints and penalties are designed to
 1. encode specific kinds of prior knowledge.
 2. Express a generic preference for a simpler model class in order to promote generalization.
 3. make an under-determined problem determined. (Provide more information)
- Other forms of regularization, known as ensemble methods, combine multiple hypotheses that explain the training data.

Generalization and Strategy

- There are many regularization strategies.
 1. Put extra constraints on a machine learning model. (Adding restrictions on the parameter values.)
 2. Add extra terms in the objective function that can be thought of as corresponding to a soft constraint on the parameter values.
- If chosen carefully, these extra constraints and penalties can lead to improved performance on the test set.
- Sometimes these constraints and penalties are designed to
 1. **encode** specific kinds of **prior knowledge**.
 2. Express a generic preference for a simpler model class in order to promote generalization.
 3. make an under-determined problem determined. (Provide more information)
- Other forms of regularization, known as ensemble methods, combine multiple hypotheses that explain the training data.

Generalization and Strategy

- There are many regularization strategies.
 1. Put extra constraints on a machine learning model. (Adding restrictions on the parameter values.)
 2. Add extra terms in the objective function that can be thought of as corresponding to a soft constraint on the parameter values.
- If chosen carefully, these extra constraints and penalties can lead to improved performance on the test set.
- Sometimes these constraints and penalties are designed to
 1. **encode** specific kinds of **prior knowledge**.
 2. Express a generic preference for a simpler model class in order to promote generalization.
 3. make an under-determined problem determined. (Provide more information)
- Other forms of regularization, known as ensemble methods, combine multiple hypotheses that explain the training data.

Generalization and Strategy

- Principle: Trading increased bias for reduced variance.
- An effective regularizer is one that makes a profitable trade, reducing variance significantly while not overly increasing the bias.
- In practice, an overly complex model family does not necessarily include the target function or the true data generating process, or even a close approximation.
- We almost never have access to the true data generating process so we can never know for sure if the model family being estimated includes the generating process or not.

Generalization and Strategy

- Principle: Treading increased bias for reduced variance.
- An effective regularizer is one that makes a profitable trade, **reducing variance** significantly while not overly **increasing the bias**.
- In practice, an **overly complex model family** does not necessarily include the target function or the true data generating process, or even a close approximation.
- We almost never have access to the true data generating process so we can never know for sure if the model family being estimated includes the generating process or not.

Generalization and Strategy

- Principle: Treading increased bias for reduced variance.
- An effective regularizer is one that makes a profitable trade, **reducing variance** significantly while not overly **increasing the bias**.
- In practice, an **overly complex model family** does not necessarily include the target function or the true data generating process, or even a close approximation.
- We almost never have access to the true data generating process so we can never know for sure if the model family being estimated includes the generating process or not.

Generalization and Strategy

- Principle: Treading increased bias for reduced variance.
- An effective regularizer is one that makes a profitable trade, **reducing variance** significantly while not overly **increasing the bias**.
- In practice, an **overly complex model family** does not necessarily include the target function or the true data generating process, or even a close approximation.
- We almost never have access to the true data generating process so we can never know for sure **if the model family being estimated includes the generating process or not**.

Generalization and Strategy

- However, most applications of deep learning algorithms are to domains where the true data generating process is almost certainly outside the model family.
- Deep learning algorithms are typically applied to **extremely complicated domains** such as images, audio sequences and text, for which the true generation process essentially involves **simulating the entire universe**.
- To some extent, we are always trying to fit a square peg(the data generating process) into a round hole (our model family)『持方枘 (ruì) 而欲内圆凿』.
- What this means is that controlling the complexity of the model is not a simple matter of finding the model of the **right size**, with the **right number of parameters**.
- Instead, we might find that the best fitting model is a large model that has been regularized appropriately.
- We now review several strategies for how to create such a large, deep, regularized model.

Generalization and Strategy

- However, most applications of deep learning algorithms are to domains where the true data generating process is almost certainly outside the model family.
- Deep learning algorithms are typically applied to **extremely complicated domains** such as images, audio sequences and text, for which the true generation process essentially involves **simulating the entire universe**.
- To some extent, we are always trying to fit a square peg(the data generating process) into a round hole (our model family)『持方枘 (ruì) 而欲内圆凿』.
- What this means is that controlling the complexity of the model is not a simple matter of finding the model of the **right size**, with the **right number of parameters**.
- Instead, we might find that the best fitting model is a large model that has been regularized appropriately.
- We now review several strategies for how to create such a large, deep, regularized model.

Generalization and Strategy

- However, most applications of deep learning algorithms are to domains where the true data generating process is almost certainly outside the model family.
- Deep learning algorithms are typically applied to **extremely complicated domains** such as images, audio sequences and text, for which the true generation process essentially involves **simulating the entire universe**.
- To some extent, we are always trying to fit a square peg(the data generating process) into a round hole (our model family)
『持方枘 (rui) 而欲内圆凿』.
- What this means is that controlling the complexity of the model is not a simple matter of finding the model of the **right size**, with the **right number of parameters**.
- Instead, we might find that the best fitting model is a large model that has been regularized appropriately.
- We now review several strategies for how to create such a large, deep, regularized model.

Generalization and Strategy

- However, most applications of deep learning algorithms are to domains where the true data generating process is almost certainly outside the model family.
- Deep learning algorithms are typically applied to **extremely complicated domains** such as images, audio sequences and text, for which the true generation process essentially involves **simulating the entire universe**.
- To some extent, we are always trying to fit a square peg(the data generating process) into a round hole (our model family)『持方枘 (rui) 而欲内圆凿』.
- What this means is that controlling the complexity of the model is not a simple matter of finding the model of the **right size**, with the **right number of parameters**.
- Instead, we might find that the best fitting model is a large model that has been regularized appropriately.
- We now review several strategies for how to create such a large, deep, regularized model.

Generalization and Strategy

- However, most applications of deep learning algorithms are to domains where the true data generating process is almost certainly outside the model family.
- Deep learning algorithms are typically applied to **extremely complicated domains** such as images, audio sequences and text, for which the true generation process essentially involves **simulating the entire universe**.
- To some extent, we are always trying to fit a square peg(the data generating process) into a round hole (our model family)『持方枘 (rui) 而欲内圆凿』.
- What this means is that controlling the complexity of the model is not a simple matter of finding the model of the **right size**, with the **right number of parameters**.
- Instead, we might find that the best fitting model is a large model that has been regularized appropriately.
- We now review several strategies for how to create such a large, deep, regularized model.

Generalization and Strategy

- However, most applications of deep learning algorithms are to domains where the true data generating process is almost certainly outside the model family.
- Deep learning algorithms are typically applied to **extremely complicated domains** such as images, audio sequences and text, for which the true generation process essentially involves **simulating the entire universe**.
- To some extent, we are always trying to fit a square peg(the data generating process) into a round hole (our model family)
『持方枘 (rui) 而欲内圆凿』.
- What this means is that controlling the complexity of the model is not a simple matter of finding the model of the **right size**, with the **right number of parameters**.
- Instead, we might find that the best fitting model is a large model that has been regularized appropriately.
- We now review several strategies for how to create such a large, deep, regularized model.

Parameter Norm Penalties

- Regularization has been used for decades prior to the advent of deep learning.
- Linear models allow simple straightforward and effective regularization strategies.
- Most approaches are based on limiting the capacity of models by adding a **parameter norm penalty** $\Omega(\theta)$ to the objective function J :

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.

- Setting α to 0 results in no regularization. Larger values of α correspond to more regularization.
- Optimize both J and norm
- Different Ω has different result.

Parameter Norm Penalties

- Regularization has been used for decades prior to the advent of deep learning.
- Linear models allow simple straightforward and effective regularization strategies.
- Most approaches are based on limiting the capacity of models by adding a **parameter norm penalty** $\Omega(\theta)$ to the objective function J :

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.

- Setting α to 0 results in no regularization. Larger values of α correspond to more regularization.
- Optimize both J and norm
- Different Ω has different result.

Parameter Norm Penalties

- Regularization has been used for decades prior to the advent of deep learning.
- Linear models allow simple straightforward and effective regularization strategies.
- Most approaches are based on limiting the capacity of models by adding a **parameter norm penalty** $\Omega(\theta)$ to the objective function J :

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.

- Setting α to 0 results in no regularization. Larger values of α correspond to more regularization.
- Optimize both J and norm
- Different Ω has different result.

Parameter Norm Penalties

- Regularization has been used for decades prior to the advent of deep learning.
- Linear models allow simple straightforward and effective regularization strategies.
- Most approaches are based on limiting the capacity of models by adding a **parameter norm penalty** $\Omega(\theta)$ to the objective function J :

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.

- Setting α to 0 results in no regularization. Larger values of α correspond to more regularization.
- Optimize both J and norm
- Different Ω has different result.

Parameter Norm Penalties

- Regularization has been used for decades prior to the advent of deep learning.
- Linear models allow simple straightforward and effective regularization strategies.
- Most approaches are based on limiting the capacity of models by adding a **parameter norm penalty** $\Omega(\theta)$ to the objective function J :

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.

- Setting α to 0 results in no regularization. Larger values of α correspond to more regularization.
- Optimize both J and norm
- Different Ω has different result.

Parameter Norm Penalties

- Regularization has been used for decades prior to the advent of deep learning.
- Linear models allow simple straightforward and effective regularization strategies.
- Most approaches are based on limiting the capacity of models by adding a **parameter norm penalty** $\Omega(\theta)$ to the objective function J :

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.

- Setting α to 0 results in no regularization. Larger values of α correspond to more regularization.
- Optimize both J and norm
- Different Ω has different result.

Parameter Norm Penalties

- We penalize **only the weights** of the affine transformation at each layer and leaves the biases unregularized.
- We do not induce too much variance by leaving the biases unregularized.
- Regularizing the bias parameters can introduce a significant amount of under-fitting.
- We therefore use the vector \mathbf{w} to indicate all of the weights that should be affected by a norm penalty, while the vector $\boldsymbol{\theta}$ denotes all of the parameters, including both \mathbf{w} and the unregularized parameters.
- Sometime we use a separate penalty with a different α coefficient for each layer.
- But it can be expensive to search for the correct value of multiple hyper-parameters, it is still reasonable to use the same weight decay at all layers just to reduce the size of search space.

Parameter Norm Penalties

- We penalize **only the weights** of the affine transformation at each layer and leaves the biases unregularized.
- We do not induce too much variance by leaving the biases unregularized.
- Regularizing the bias parameters can introduce a significant amount of under-fitting.
- We therefore use the vector \mathbf{w} to indicate all of the weights that should be affected by a norm penalty, while the vector $\boldsymbol{\theta}$ denotes all of the parameters, including both \mathbf{w} and the unregularized parameters.
- Sometime we use a separate penalty with a different α coefficient for each layer.
- But it can be expensive to search for the correct value of multiple hyper-parameters, it is still reasonable to use the same weight decay at all layers just to reduce the size of search space.

Parameter Norm Penalties

- We penalize **only the weights** of the affine transformation at each layer and leaves the biases unregularized.
- We do not induce too much variance by leaving the biases unregularized.
- Regularizing the bias parameters can introduce a significant amount of under-fitting.
- We therefore use the vector w to indicate all of the weights that should be affected by a norm penalty, while the vector θ denotes all of the parameters, including both w and the unregularized parameters.
- Sometime we use a separate penalty with a different α coefficient for each layer.
- But it can be expensive to search for the correct value of multiple hyper-parameters, it is still reasonable to use the same weight decay at all layers just to reduce the size of search space.

Parameter Norm Penalties

- We penalize **only the weights** of the affine transformation at each layer and leaves the biases unregularized.
- We do not induce too much variance by leaving the biases unregularized.
- Regularizing the bias parameters can introduce a significant amount of under-fitting.
- We therefore use the vector \mathbf{w} to indicate all of the weights that should be affected by a norm penalty, while the vector $\boldsymbol{\theta}$ denotes all of the parameters, including both \mathbf{w} and the unregularized parameters.
- Sometime we use a separate penalty with a different α coefficient for each layer.
- But it can be expensive to search for the correct value of multiple hyper-parameters, it is still reasonable to use the same weight decay at all layers just to reduce the size of search space.

Parameter Norm Penalties

- We penalize **only the weights** of the affine transformation at each layer and leaves the biases unregularized.
- We do not induce too much variance by leaving the biases unregularized.
- Regularizing the bias parameters can introduce a significant amount of under-fitting.
- We therefore use the vector \mathbf{w} to indicate all of the weights that should be affected by a norm penalty, while the vector $\boldsymbol{\theta}$ denotes all of the parameters, including both \mathbf{w} and the unregularized parameters.
- Sometime we use a separate penalty with a different α coefficient for each layer.
- But it can be expensive to search for the correct value of multiple hyper-parameters, it is still reasonable to use the same weight decay at all layers just to reduce the size of search space.

L^2 Parameter Regularization

- The L^2 norm penalty commonly known as *weight decay*.

$$\Omega(\theta) = \frac{1}{2} \|w\|_2^2$$

- This regularization strategy drives the weights closer to the origin. (as well as *ridge regression* or *Tikhonov regularization*)
- We can gain some insight into the behavior of weight decay regularization. (assume no bias for simplification)

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^T w + J(w; X, y)$$

$$\nabla_w \tilde{J}(w; X, y) = \alpha w + \nabla_w J(w; X, y)$$

- The update

$$w \leftarrow w - \epsilon(\alpha w + \nabla_w J(w; X, y))$$

$$w \leftarrow (1 - \epsilon\alpha)w - \epsilon\nabla_w J(w; X, y)$$

- Shrink the weight vector by a constant factor on each step.
- What happens over the entire course of training?

L^2 Parameter Regularization

- The L^2 norm penalty commonly known as *weight decay*.

$$\Omega(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{w}\|_2^2$$

- This regularization strategy drives the weights closer to the origin. (as well as *ridge regression* or *Tikhonov regularization*)
- We can gain some insight into the behavior of weight decay regularization. (assume no bias for simplification)

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

- The update

$$\mathbf{w} \leftarrow \mathbf{w} - \epsilon(\alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}))$$

$$\mathbf{w} \leftarrow (1 - \epsilon\alpha) \mathbf{w} - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

- Shrink the weight vector by a constant factor on each step.
- What happens over the entire course of training?

L^2 Parameter Regularization

- The L^2 norm penalty commonly known as *weight decay*.

$$\Omega(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{w}\|_2^2$$

- This regularization strategy drives the weights closer to the origin. (as well as *ridge regression* or *Tikhonov regularization*)
- We can gain some insight into the behavior of weight decay regularization. (assume no bias for simplification)

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

- The update

$$\mathbf{w} \leftarrow \mathbf{w} - \epsilon(\alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}))$$

$$\mathbf{w} \leftarrow (1 - \epsilon\alpha) \mathbf{w} - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

- Shrink the weight vector by a constant factor on each step.
- What happens over the entire course of training?

L^2 Parameter Regularization

- The L^2 norm penalty commonly known as *weight decay*.

$$\Omega(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{w}\|_2^2$$

- This regularization strategy drives the weights closer to the origin. (as well as *ridge regression* or *Tikhonov regularization*)
- We can gain some insight into the behavior of weight decay regularization. (assume no bias for simplification)

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

- The update

$$\mathbf{w} \leftarrow \mathbf{w} - \epsilon(\alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}))$$

$$\mathbf{w} \leftarrow (1 - \epsilon\alpha) \mathbf{w} - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

- Shrink the weight vector by a constant factor on each step.
- What happens over the entire course of training?

Recall: Quadratic Approximation

- In mathematics, approximation theory is concerned with how functions can best be approximated with simpler functions.
- **local linear approximation** and **taylor expansion**

1. For example, when the independent variable of function $y = x^3$ changes, which is Δx , the variation of y is

$$\Delta y = (x + \Delta x)^3 - x^3 = 3x^2 \Delta x + 3x(\Delta x)^2 + (\Delta x)^3$$

2. When $\Delta x \rightarrow 0$, omit last two terms: $\Delta y = 3x^2 \Delta x$
3. In general:

$$\Delta y = f(x_0 + \Delta x) - f(x_0) \approx f'(x_0) \times \Delta x$$

$$\Delta y = f(x) - f(x_0), \Delta x = x - x_0$$

$$f(x) - f(x_0) = f'(x_0) \times (x - x_0)$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0)$$

4. In order to improve the precision, we can use second-order approximation, which is the second-order Taylor series expansion.

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2$$

Recall: Quadratic Approximation

- In mathematics, approximation theory is concerned with how functions can best be approximated with simpler functions.
- local linear approximation and Taylor expansion

1. For example, when the independent variable of function $y = x^3$ changes, which is Δx , the variation of y is

$$\Delta y = (x + \Delta x)^3 - x^3 = 3x^2 \Delta x + 3x(\Delta x)^2 + (\Delta x)^3$$

2. When $\Delta x \rightarrow 0$, omit last two terms: $\Delta y = 3x^2 \Delta x$
3. In general:

$$\Delta y = f(x_0 + \Delta x) - f(x_0) \approx f'(x_0) \times \Delta x$$

$$\Delta y = f(x) - f(x_0), \Delta x = x - x_0$$

$$f(x) - f(x_0) = f'(x_0) \times (x - x_0)$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0)$$

4. In order to improve the precision, we can use second-order approximation, which is the second-order Taylor series expansion.

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2$$

Recall: Quadratic Approximation

- In mathematics, approximation theory is concerned with how functions can best be approximated with simpler functions.
- **local linear approximation** and **taylor expansion**

1. For example, when the independent variable of function $y = x^3$ changes, which is Δx , the variation of y is

$$\Delta y = (x + \Delta x)^3 - x^3 = 3x^2 \Delta x + 3x(\Delta x)^2 + (\Delta x)^3$$

2. When $\Delta x \rightarrow 0$, omit last two terms: $\Delta y = 3x^2 \Delta x$
3. In general:

$$\Delta y = f(x_0 + \Delta x) - f(x_0) \approx f'(x_0) \times \Delta x$$

$$\Delta y = f(x) - f(x_0), \Delta x = x - x_0$$

$$f(x) - f(x_0) = f'(x_0) \times (x - x_0)$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0)$$

4. In order to improve the precision, we can use second-order approximation, which is the second-order Taylor series expansion.

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2$$

Recall: Quadratic Approximation

- In mathematics, approximation theory is concerned with how functions can best be approximated with simpler functions.
- **local linear approximation** and **taylor expansion**

1. For example, when the independent variable of function $y = x^3$ changes, which is Δx , the variation of y is

$$\Delta y = (x + \Delta x)^3 - x^3 = 3x^2 \Delta x + 3x(\Delta x)^2 + (\Delta x)^3$$

2. When $\Delta x \rightarrow 0$, omit last two terms: $\Delta y = 3x^2 \Delta x$
3. In general:

$$\Delta y = f(x_0 + \Delta x) - f(x_0) \approx f'(x_0) \times \Delta x$$

$$\Delta y = f(x) - f(x_0), \Delta x = x - x_0$$

$$f(x) - f(x_0) = f'(x_0) \times (x - x_0)$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0)$$

4. In order to improve the precision, we can use second-order approximation, which is the second-order Taylor series expansion.

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2$$

Recall: Quadratic Approximation

- In mathematics, approximation theory is concerned with how functions can best be approximated with simpler functions.
- **local linear approximation** and **taylor expansion**

1. For example, when the independent variable of function $y = x^3$ changes, which is Δx , the variation of y is

$$\Delta y = (x + \Delta x)^3 - x^3 = 3x^2 \Delta x + 3x(\Delta x)^2 + (\Delta x)^3$$

2. When $\Delta x \rightarrow 0$, omit last two terms: $\Delta y = 3x^2 \Delta x$

3. In general:

$$\Delta y = f(x_0 + \Delta x) - f(x_0) \approx f'(x_0) \times \Delta x$$

$$\Delta y = f(x) - f(x_0), \Delta x = x - x_0$$

$$f(x) - f(x_0) = f'(x_0) \times (x - x_0)$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0)$$

4. In order to improve the precision, we can use second-order approximation, which is the second-order Taylor series expansion.

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2$$

Recall: Quadratic Approximation

- In mathematics, approximation theory is concerned with how functions can best be approximated with simpler functions.
- **local linear approximation** and **taylor expansion**

1. For example, when the independent variable of function $y = x^3$ changes, which is Δx , the variation of y is

$$\Delta y = (x + \Delta x)^3 - x^3 = 3x^2 \Delta x + 3x(\Delta x)^2 + (\Delta x)^3$$

2. When $\Delta x \rightarrow 0$, omit last two terms: $\Delta y = 3x^2 \Delta x$
3. In general:

$$\Delta y = f(x_0 + \Delta x) - f(x_0) \approx f'(x_0) \times \Delta x$$

$$\Delta y = f(x) - f(x_0), \Delta x = x - x_0$$

$$f(x) - f(x_0) = f'(x_0) \times (x - x_0)$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0)$$

4. In order to improve the precision, we can use second-order approximation, which is the second-order Taylor series expansion.

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2$$

Recall: Quadratic Approximation

- In mathematics, approximation theory is concerned with how functions can best be approximated with simpler functions.
- **local linear approximation** and **taylor expansion**

1. For example, when the independent variable of function $y = x^3$ changes, which is Δx , the variation of y is

$$\Delta y = (x + \Delta x)^3 - x^3 = 3x^2 \Delta x + 3x(\Delta x)^2 + (\Delta x)^3$$

2. When $\Delta x \rightarrow 0$, omit last two terms: $\Delta y = 3x^2 \Delta x$
3. In general:

$$\Delta y = f(x_0 + \Delta x) - f(x_0) \approx f'(x_0) \times \Delta x$$

$$\Delta y = f(x) - f(x_0), \Delta x = x - x_0$$

$$f(x) - f(x_0) = f'(x_0) \times (x - x_0)$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0)$$

4. In order to improve the precision, we can use second-order approximation, which is the second-order Taylor series expansion.

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2$$

L^2 Parameter Regularization

- Let $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$ (unregularized training cost)
- Making a quadratic approximation to the objective function in the neighborhood of the value of the weights. (In DLBook, they used $\hat{J}(\boldsymbol{\theta})$, but here we use $\hat{J}(\mathbf{w})$ to explain easier)

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

- Where \mathbf{H} is the Hessian matrix of J with respect to \mathbf{w} evaluated at \mathbf{w}^* .
- There is no first-order term in this quadratic approximation, because \mathbf{w}^* is defined to be a minimum, where the gradient vanishes.
- The minimum of \hat{J} occurs where its gradient

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

is equal to 0.

L^2 Parameter Regularization

- Let $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$ (unregularized training cost)
- Making a quadratic approximation to the objective function in the neighborhood of the value of the weights. (In DLBook, they used $\hat{J}(\boldsymbol{\theta})$, but here we use $\hat{J}(\mathbf{w})$ to explain easier)

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

- Where \mathbf{H} is the Hessian matrix of J with respect to \mathbf{w} evaluated at \mathbf{w}^* .
- There is no first-order term in this quadratic approximation, because \mathbf{w}^* is defined to be a minimum, where the gradient vanishes.
- The minimum of \hat{J} occurs where its gradient

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

is equal to 0.

L^2 Parameter Regularization

- Let $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$ (unregularized training cost)
- Making a quadratic approximation to the objective function in the neighborhood of the value of the weights. (In DLBook, they used $\hat{J}(\boldsymbol{\theta})$, but here we use $\hat{J}(\mathbf{w})$ to explain easier)

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

- Where \mathbf{H} is the Hessian matrix of J with respect to \mathbf{w} evaluated at \mathbf{w}^* .
- There is no first-order term in this quadratic approximation, because \mathbf{w}^* is defined to be a minimum, where the gradient vanishes.
- The minimum of \hat{J} occurs where its gradient

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

is equal to 0.

L^2 Parameter Regularization

- Let $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$ (unregularized training cost)
- Making a quadratic approximation to the objective function in the neighborhood of the value of the weights. (In DLBook, they used $\hat{J}(\boldsymbol{\theta})$, but here we use $\hat{J}(\mathbf{w})$ to explain easier)

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

- Where \mathbf{H} is the Hessian matrix of J with respect to \mathbf{w} evaluated at \mathbf{w}^* .
- There is no first-order term in this quadratic approximation, because \mathbf{w}^* is defined to be a minimum, where the gradient vanishes.
- The minimum of \hat{J} occurs where its gradient

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

is equal to 0.

L^2 Parameter Regularization

- Let $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$ (unregularized training cost)
- Making a quadratic approximation to the objective function in the neighborhood of the value of the weights. (In DLBook, they used $\hat{J}(\boldsymbol{\theta})$, but here we use $\hat{J}(\mathbf{w})$ to explain easier)

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

- Where \mathbf{H} is the Hessian matrix of J with respect to \mathbf{w} evaluated at \mathbf{w}^* .
- There is no first-order term in this quadratic approximation, because \mathbf{w}^* is defined to be a minimum, where the gradient vanishes.
- The minimum of \hat{J} occurs where its gradient

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

is equal to 0.

L^2 Parameter Regularization

- To study the effect of weight decay, we modify $\nabla_w \hat{J}(w) = H(w - w^*)$ by adding the weight decay gradient.
- We can solve for the minimum of the regularized version of \hat{J} .
- We use the variable \tilde{w} to represent the location of the minimum.

$$\alpha \tilde{w} + H(\tilde{w} - w^*) = 0$$

$$(H + \alpha I) \tilde{w} = H w^*$$

$$\tilde{w} = \frac{H w^*}{(H + \alpha I)}$$

- As α approaches 0, the regularized solution \tilde{w} approaches w^* .
- But what happens as α grows?

L^2 Parameter Regularization

- To study the effect of weight decay, we modify $\nabla_w \hat{J}(w) = H(w - w^*)$ by adding the weight decay gradient.
- We can solve for the minimum of the regularized version of \hat{J} .
- We use the variable \tilde{w} to represent the location of the minimum.

$$\alpha \tilde{w} + H(\tilde{w} - w^*) = 0$$

$$(H + \alpha I) \tilde{w} = H w^*$$

$$\tilde{w} = \frac{H w^*}{(H + \alpha I)}$$

- As α approaches 0, the regularized solution \tilde{w} approaches w^* .
- But what happens as α grows?

L^2 Parameter Regularization

- To study the effect of weight decay, we modify $\nabla_w \hat{J}(w) = H(w - w^*)$ by adding the weight decay gradient.
- We can solve for the minimum of the regularized version of \hat{J} .
- We use the variable \tilde{w} to represent the location of the minimum.

$$\alpha \tilde{w} + H(\tilde{w} - w^*) = 0$$

$$(H + \alpha I) \tilde{w} = H w^*$$

$$\tilde{w} = \frac{H w^*}{(H + \alpha I)}$$

- As α approaches 0, the regularized solution \tilde{w} approaches w^* .
- But what happens as α grows?

L^2 Parameter Regularization

- Because \mathbf{H} is real and symmetric, we can decompose it into a diagonal matrix $\mathbf{\Lambda}$ and an orthonormal basis of eigenvectors, \mathbf{Q} , such that $\mathbf{H} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$.
- Applying the decomposition $\tilde{\mathbf{w}} = (\mathbf{H} + \alpha\mathbf{I})^{-1}\mathbf{H}\mathbf{w}^*$

$$\tilde{\mathbf{w}} = (\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T + \alpha\mathbf{I})^{-1}\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (1)$$

$$= [\mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})\mathbf{Q}^T]^{-1}\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (2)$$

$$= \mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})^{-1}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (3)$$

$$= \mathbf{Q}\frac{\mathbf{\Lambda}}{\mathbf{\Lambda} + \alpha\mathbf{I}}\mathbf{Q}^T\mathbf{w}^* \quad (4)$$

- We see that the effect of weight decay is to rescale \mathbf{w}^* along the axes defined by the eigenvectors of \mathbf{H} .
- Specifically, the component of \mathbf{w}^* that is aligned with the i -th eigenvector of \mathbf{H} is rescaled by a factor of $\frac{\lambda_i}{\lambda_i + \alpha}$

L^2 Parameter Regularization

- Because \mathbf{H} is real and symmetric, we can decompose it into a diagonal matrix $\mathbf{\Lambda}$ and an orthonormal basis of eigenvectors, \mathbf{Q} , such that $\mathbf{H} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$.
- Applying the decomposition $\tilde{\mathbf{w}} = (\mathbf{H} + \alpha\mathbf{I})^{-1}\mathbf{H}\mathbf{w}^*$

$$\tilde{\mathbf{w}} = (\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T + \alpha\mathbf{I})^{-1}\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (1)$$

$$= [\mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})\mathbf{Q}^T]^{-1}\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (2)$$

$$= \mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})^{-1}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (3)$$

$$= \mathbf{Q}\frac{\mathbf{\Lambda}}{\mathbf{\Lambda} + \alpha\mathbf{I}}\mathbf{Q}^T\mathbf{w}^* \quad (4)$$

- We see that the effect of weight decay is to rescale \mathbf{w}^* along the axes defined by the eigenvectors of \mathbf{H} .
- Specifically, the component of \mathbf{w}^* that is aligned with the i -th eigenvector of \mathbf{H} is rescaled by a factor of $\frac{\lambda_i}{\lambda_i + \alpha}$

L^2 Parameter Regularization

- Because \mathbf{H} is real and symmetric, we can decompose it into a diagonal matrix $\mathbf{\Lambda}$ and an orthonormal basis of eigenvectors, \mathbf{Q} , such that $\mathbf{H} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$.
- Applying the decomposition $\tilde{\mathbf{w}} = (\mathbf{H} + \alpha\mathbf{I})^{-1}\mathbf{H}\mathbf{w}^*$

$$\tilde{\mathbf{w}} = (\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T + \alpha\mathbf{I})^{-1}\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (1)$$

$$= [\mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})\mathbf{Q}^T]^{-1}\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (2)$$

$$= \mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})^{-1}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (3)$$

$$= \mathbf{Q}\frac{\mathbf{\Lambda}}{\mathbf{\Lambda} + \alpha\mathbf{I}}\mathbf{Q}^T\mathbf{w}^* \quad (4)$$

- We see that the effect of weight decay is to rescale \mathbf{w}^* along the axes defined by the eigenvectors of \mathbf{H} .
- Specifically, the component of \mathbf{w}^* that is aligned with the i -th eigenvector of \mathbf{H} is rescaled by a factor of $\frac{\lambda_i}{\lambda_i + \alpha}$

L^2 Parameter Regularization

- Because \mathbf{H} is real and symmetric, we can decompose it into a diagonal matrix $\mathbf{\Lambda}$ and an orthonormal basis of eigenvectors, \mathbf{Q} , such that $\mathbf{H} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$.
- Applying the decomposition $\tilde{\mathbf{w}} = (\mathbf{H} + \alpha\mathbf{I})^{-1}\mathbf{H}\mathbf{w}^*$

$$\tilde{\mathbf{w}} = (\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T + \alpha\mathbf{I})^{-1}\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (1)$$

$$= [\mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})\mathbf{Q}^T]^{-1}\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (2)$$

$$= \mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})^{-1}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^* \quad (3)$$

$$= \mathbf{Q}\frac{\mathbf{\Lambda}}{\mathbf{\Lambda} + \alpha\mathbf{I}}\mathbf{Q}^T\mathbf{w}^* \quad (4)$$

- We see that the effect of weight decay is to rescale \mathbf{w}^* along the axes defined by the eigenvectors of \mathbf{H} .
- Specifically, the component of \mathbf{w}^* that is aligned with the i -th eigenvector of \mathbf{H} is rescaled by a factor of $\frac{\lambda_i}{\lambda_i + \alpha}$

References
