# 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.

- 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.

3

## 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 constrains 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 constrains 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 constrains 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 constrains 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.

4

## Generalization and Strategy

- There are many regularization strategies.
    1. Put extra constrains 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 constrains 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.

- 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

- 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 (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**.
- Insteamd, 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.

6

## 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**.
- Insteamd, 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.

6

## 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**.
- Insteamd, 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.

6

## 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**.
- Insteamd, 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.

6

## 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**.
- Insteamd, 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.

6

## 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**.
- Insteamd, 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 $\alpha$ to 0 results in no regularization. Larger values of $\alpha$ correspond to more regularization.
- Optimize both $J$ and norm
- Different $\Omega$ has different result.

7

## 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 $\alpha$ to 0 results in no regularization. Larger values of $\alpha$ correspond to more regularization.
- Optimize both $J$ and norm
- Different $\Omega$ 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 $J$:

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

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

- Setting $\alpha$ to 0 results in no regularization. Larger values of $\alpha$ correspond to more regularization.
- Optimize both $J$ and norm
- Different $\Omega$ 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 $J$:

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

  where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.
- Setting $\alpha$ to 0 results in no regularization. Larger values of $\alpha$ correspond to more regularization.
- Optimize both $J$ and norm
- Different $\Omega$ 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 $J$:

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

  where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.
- Setting $\alpha$ to 0 results in no regularization. Larger values of $\alpha$ correspond to more regularization.
- Optimize both $J$ and norm
- Different $\Omega$ 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 $J$:

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

  where $\alpha \in [0, +\infty)$ weights the relative contribution of the norm penalty term.
- Setting $\alpha$ to 0 results in no regularization. Larger values of $\alpha$ correspond to more regularization.
- Optimize both $J$ and norm
- Different $\Omega$ 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 $w$ to indicate all of the weights that should be affected by a norm penalty, while the vector $\theta$ denotes all of the parameters, including both $w$ and the unregularized parameters.
- Sometime we use a separate penalty with a different $\alpha$ 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 $\theta$ denotes all of the parameters, including both $w$ and the unregularized parameters.
- Sometime we use a separate penalty with a different $\alpha$ 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 $\theta$ denotes all of the parameters, including both $w$ and the unregularized parameters.
- Sometime we use a separate penalty with a different $\alpha$ 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 $\theta$ denotes all of the parameters, including both $w$ and the unregularized parameters.
- Sometime we use a separate penalty with a different $\alpha$ 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 $\alpha$ 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*.

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

- 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.*

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

- 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?

9

# $L^2$ Parameter Regularization

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

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

- 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?

9

# $L^2$ Parameter Regularization

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

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

- 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?

9

## 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 $\Delta x$, 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 $\Delta x \to 0$, 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 $\Delta x$, 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 $\Delta x \to 0$, 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 $\Delta x$, 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 $\Delta x \to 0$, 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 $\Delta x$, 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 $\Delta x \to 0$, 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 $\Delta x$, 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 $\Delta x \to 0$, 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.

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

10

## 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 $\Delta x$, 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 $\Delta x \to 0$, 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 $\Delta x$, 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 $\Delta x \to 0$, 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 $w^* = \arg\min_w J(w)$ (unregularized training cost)
- Making a quadratic approximation to the objective function in the neighborhood of the value of the weights. (In DLBook, they used $\hat{J}(\theta)$, but here we use $\hat{J}(w)$ to explain easier)

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

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

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

is equal to 0.

## $L^2$ Parameter Regularization

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

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

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

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

is equal to 0.

## $L^2$ Parameter Regularization

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

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

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

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

is equal to 0.

## $L^2$ Parameter Regularization

- Let $w^* = \arg\min_w J(w)$ (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 $\hat{J}(w)$ to explain easier)

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

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

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

is equal to 0.

## $L^2$ Parameter Regularization

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

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

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

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

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.

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

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

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

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

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

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

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

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

12

## $L^2$ 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.

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

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

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

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

# $L^2$ Parameter Regularization

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

$$
\begin{align}
\tilde{w} &= (Q\Lambda Q^T + \alpha I)^{-1} Q\Lambda Q^T w^* \tag{1}\\
&= \left[ Q(\Lambda + \alpha I)Q^T \right]^{-1} Q\Lambda Q^T w^* \tag{2}\\
&= Q(\Lambda + \alpha I)^{-1} \Lambda Q^T w^* \tag{3}\\
&= Q\frac{\Lambda}{\Lambda + \alpha I} Q^T w^* \tag{4}
\end{align}
$$

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

# $L^2$ Parameter Regularization

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

$$\begin{align}
\tilde{w} &= (Q\mathbf{\Lambda}Q^T + \alpha\mathbf{I})^{-1}Q\mathbf{\Lambda}Q^Tw^* \tag{1} \\
&= \left[Q(\mathbf{\Lambda} + \alpha\mathbf{I})Q^T\right]^{-1}Q\mathbf{\Lambda}Q^Tw^* \tag{2} \\
&= Q(\mathbf{\Lambda} + \alpha\mathbf{I})^{-1}\mathbf{\Lambda}Q^Tw^* \tag{3} \\
&= Q\frac{\mathbf{\Lambda}}{\mathbf{\Lambda} + \alpha\mathbf{I}}Q^Tw^* \tag{4}
\end{align}$$

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

13

## $L^2$ Parameter Regularization

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

$$
\begin{align}
\tilde{w} &= (Q\Lambda Q^T + \alpha I)^{-1}Q\Lambda Q^T w^* \tag{1}\\
&= \left[Q(\Lambda + \alpha I)Q^T\right]^{-1} Q\Lambda Q^T w^* \tag{2}\\
&= Q(\Lambda + \alpha I)^{-1}\Lambda Q^T w^* \tag{3}\\
&= Q\frac{\Lambda}{\Lambda + \alpha I}Q^T w^* \tag{4}
\end{align}
$$

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

# $L^2$ Parameter Regularization

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

$$
\begin{aligned}
\tilde{w} &= (Q\Lambda Q^T + \alpha I)^{-1}Q\Lambda Q^T w^* & (1) \\
&= \left[Q(\Lambda + \alpha I)Q^T\right]^{-1} Q\Lambda Q^T w^* & (2) \\
&= Q(\Lambda + \alpha I)^{-1}\Lambda Q^T w^* & (3) \\
&= Q\frac{\Lambda}{\Lambda + \alpha I}Q^T w^* & (4)
\end{aligned}
$$

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

13

# $L^2$ Parameter Regularization

This effect is illustrated in figure:

Fig. 1: An illustration of the effect of $L^2$ (or weight decay) regularization on the value of the optimal **w**



- The solid ellipses represent contours of equal value of the unregularized objective.
- The dotted circles represent contours of equal value of the $L^2$ regularizer.
- At the point $\tilde{w}$, these competing objectives reach an equilibrium.

# $L^2$ Parameter Regularization

This effect is illustrated in figure:

Fig. 1: An illustration of the effect of $L^2$ (or weight decay) regularization on the value of the optimal *w*



- The solid ellipses represent contours of equal value of the **unregularized objective**.
- The dotted circles represent contours of equal value of the $L^2$ regularizer.
- At the point $\tilde{w}$, these competing objectives reach an equilibrium.

14

# $L^2$ Parameter Regularization

This effect is illustrated in figure:

Fig. 1: An illustration of the effect of $L^2$ (or weight decay) regularization on the value of the optimal **w**
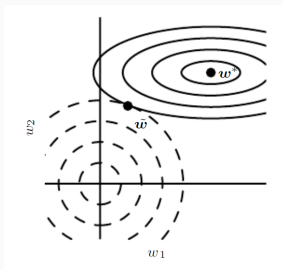


- The solid ellipses represent contours of equal value of the unregularized objective.
- The dotted circles represent contours of equal value of the $L^2$ regularizer.
- At the point $\tilde{w}$, these competing objectives reach an equilibrium.

# $L^2$ Parameter Regularization

This effect is illustrated in figure:

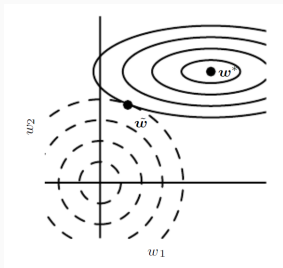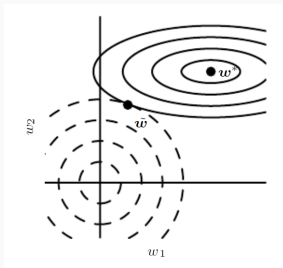Fig. 1: An illustration of the effect of $L^2$ (or weight decay) regularization on the value of the optimal **w**



- The solid ellipses represent contours of equal value of the unregularized objective.
- The dotted circles represent contours of equal value of the $L^2$ regularizer.
- At the point $\tilde{w}$, these competing objectives reach an equilibrium.

# $L^2$ Parameter Regularization

- How do these effects relate to machine learning in particular?
- We can find out by studying linear regression, the cost function is the sum of squared errors:

$$(Xw - y)^T(Xw - y)$$

- Add $L^2$ regularization, the objective function changes to:

$$(Xw - y)^T(Xw - y) + \frac{1}{2}\alpha w^T w$$

- This changes the normal equations for the solution from:

$$w = (X^TX)^{-1}X^Ty \text{ to } w = (X^TX + \alpha I)^{-1}X^Ty$$

- The new matrix has the addition of $\alpha$ to the diagonal.
- Diagonal correspond to the variance of each input feature.
- We can see that $L^2$ regularization causes the learning algorithm to "perceive" the input $X$ as having higher variance, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

# $L^2$ Parameter Regularization

- How do these effects relate to machine learning in particular?
- We can find out by studying linear regression, the cost function is the sum of squared errors:

$$(Xw - y)^T(Xw - y)$$

- Add $L^2$ regularization, the objective function changes to:

$$(Xw - y)^T(Xw - y) + \frac{1}{2}\alpha w^T w$$

- This changes the normal equations for the solution from:

$$w = (X^T X)^{-1} X^T y \text{ to } w = (X^T X + \alpha I)^{-1} X^T y$$

- The new matrix has the addition of $\alpha$ to the diagonal.
- Diagonal correspond to the variance of each input feature.
- We can see that $L^2$ regularization causes the learning algorithm to "perceive" the input $X$ as having higher variance, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

## $L^2$ Parameter Regularization

- How do these effects relate to machine learning in particular?
- We can find out by studying linear regression, the cost function is the sum of squared errors:

$$(Xw - y)^T(Xw - y)$$

- Add $L^2$ regularization, the objective function changes to:

$$(Xw - y)^T(Xw - y) + \frac{1}{2}\alpha w^T w$$

- This changes the normal equations for the solution from:

$$w = (X^TX)^{-1}X^Ty \text{ to } w = (X^TX + \alpha I)^{-1}X^Ty$$

- The new matrix has the addition of $\alpha$ to the diagonal.
- Diagonal correspond to the variance of each input feature.
- We can see that $L^2$ regularization causes the learning algorithm to "perceive" the input $X$ as having higher variance, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

15

## $L^2$ Parameter Regularization

- How do these effects relate to machine learning in particular?
- We can find out by studying linear regression, the cost function is the sum of squared errors:

$$(Xw - y)^T(Xw - y)$$

- Add $L^2$ regularization, the objective function changes to:

$$(Xw - y)^T(Xw - y) + \frac{1}{2}\alpha w^T w$$

- This changes the normal equations for the solution from:

$$w = (X^T X)^{-1} X^T y \text{ to } w = (X^T X + \alpha I)^{-1} X^T y$$

- The new matrix has the addition of $\alpha$ to the diagonal.
- Diagonal correspond to the variance of each input feature.
- We can see that $L^2$ regularization causes the learning algorithm to "perceive" the input $X$ as having higher variance, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

15

# $L^2$ Parameter Regularization

- How do these effects relate to machine learning in particular?
- We can find out by studying linear regression, the cost function is the sum of squared errors:

$$(Xw - y)^T(Xw - y)$$

- Add $L^2$ regularization, the objective function changes to:

$$(Xw - y)^T(Xw - y) + \frac{1}{2}\alpha w^T w$$

- This changes the normal equations for the solution from:

$$w = (X^TX)^{-1}X^Ty \text{ to } w = (X^TX + \alpha I)^{-1}X^Ty$$

- The new matrix has the addition of $\alpha$ to the diagonal.
- Diagonal correspond to the variance of each input feature.
- We can see that $L^2$ regularization causes the learning algorithm to "perceive" the input $X$ as having higher variance, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

## $L^2$ Parameter Regularization

- How do these effects relate to machine learning in particular?
- We can find out by studying linear regression, the cost function is the sum of squared errors:

$$(Xw - y)^T(Xw - y)$$

- Add $L^2$ regularization, the objective function changes to:

$$(Xw - y)^T(Xw - y) + \frac{1}{2}\alpha w^T w$$

- This changes the normal equations for the solution from:

$$w = (X^T X)^{-1} X^T y \text{ to } w = (X^T X + \alpha I)^{-1} X^T y$$

- The new matrix has the addition of $\alpha$ to the diagonal.
- Diagonal correspond to the variance of each input feature.
- We can see that $L^2$ regularization causes the learning algorithm to "perceive" the input $X$ as having higher variance, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

15

# $L^2$ Parameter Regularization

- How do these effects relate to machine learning in particular?
- We can find out by studying linear regression, the cost function is the sum of squared errors:

$$(Xw - y)^T(Xw - y)$$

- Add $L^2$ regularization, the objective function changes to:

$$(Xw - y)^T(Xw - y) + \frac{1}{2}\alpha w^T w$$

- This changes the normal equations for the solution from:

$$w = (X^TX)^{-1}X^Ty \text{ to } w = (X^TX + \alpha I)^{-1}X^Ty$$

- The new matrix has the addition of $\alpha$ to the diagonal.
- Diagonal correspond to the variance of each input feature.
- We can see that $L^2$ regularization causes the learning algorithm to "perceive" the input $X$ as having higher variance, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

15

## $L^1$ Regularization

- $L^1$ regularization on the model parameter $w$ is defined as:

$$\Omega(\boldsymbol{\theta}) = \|\boldsymbol{w}\|_1 = \sum_i |w_i|$$

- We will now discuss the effect of $L^1$ regularization on the simple linear regression model, with no bias parameters, that we studied in our analysis of $L^2$ regularization.

- In particular, we are interested in delineating the differences between $L^1$ and $L^2$ forms of regularization.

- $L^1$ regularization on the model parameter $w$ is defined as:

$$\Omega(\boldsymbol{\theta}) = \|\boldsymbol{w}\|_1 = \sum_i |w_i|$$

- We will now discuss the effect of $L^1$ regularization on the simple linear regression model, with no bias parameters, that we studied in our analysis of $L^2$ regularization.

- In particular, we are interested in delineating the differences between $L^1$ and $L^2$ forms of regularization.

## $L^1$ Regularization

- As with $L^2$ weight decay, $L^1$ weight decay controls the strength of the regularization by scaling the penalty $\Omega$ using a positive hyperparameter $\alpha$.

- Thus, the regularized objective function $\tilde{J}(w; X, y)$ is given by

$$\tilde{J}(w; X, y) = \alpha \|w\|_1 + J(w; X, y)$$

with the corresponding gradient:

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

where $\text{sign}(w)$ is simply the sign of $w$ applied element-wise.

# $L^1$ Regularization

- As with $L^2$ weight decay, $L^1$ weight decay controls the strength of the regularization by scaling the penalty $\Omega$ using a positive hyperparameter $\alpha$.

- Thus, the regularized objective function $\tilde{J}(w; X, y)$ is given by

$$\tilde{J}(w; X, y) = \alpha \|w\|_1 + J(w; X, y)$$

with the corresponding gradient:

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

where $\operatorname{sign}(w)$ is simply the sign of $w$ applied element-wise.

# $L^1$ Regularization

- As with $L^2$ weight decay, $L^1$ weight decay controls the strength of the regularization by scaling the penalty $\Omega$ using a positive hyperparameter $\alpha$.

- Thus, the regularized objective function $\tilde{J}(w; X, y)$ is given by

$$\tilde{J}(w; X, y) = \alpha\|w\|_1 + J(w; X, y)$$

  with the corresponding gradient:

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

  where $\text{sign}(w)$ is simply the sign of $w$ applied element-wise.

# $L^1$ Regularization

- As with $L^2$ weight decay, $L^1$ weight decay controls the strength of the regularization by scaling the penalty $\Omega$ using a positive hyperparameter $\alpha$.

- Thus, the regularized objective function $\tilde{J}(w; X, y)$ is given by

$$\tilde{J}(w; X, y) = \alpha \|w\|_1 + J(w; X, y)$$

with the corresponding gradient:

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

where $\text{sign}(w)$ is simply the sign of $w$ applied element-wise.

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

- From this equation, we can see that the effect of $L^1$ regularization is quite different from that of $L^2$ regularization.

- We can see that the regularization contribution to the gradient no longer scales linearly with each $w_i$; instead it is a constant factor with a sign equal to $\text{sign}(w_i)$.

- One consequence of this form of the gradient is that we will not necessarily see clean algebraic solutions to quadratic approximations of $J(X, y; w)$ as we did for $L^2$ regularization.

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

- From this equation, we can see that the effect of $L^1$ regularization is quite different from that of $L^2$ regularization.
- We can see that the regularization contribution to the gradient no longer scales linearly with each $w_i$; instead it is a constant factor with a sign equal to $\text{sign}(w_i)$.
- One consequence of this form of the gradient is that we will not necessarily see clean algebraic solutions to quadratic approximations of $J(X, y; w)$ as we did for $L^2$ regularization.

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

- From this equation, we can see that the effect of $L^1$ regularization is quite different from that of $L^2$ regularization.
- We can see that the regularization contribution to the gradient no longer scales linearly with each $w_i$; instead it is a constant factor with a sign equal to $\mathrm{sign}(w_i)$.
- One consequence of this form of the gradient is that we will not necessarily see clean algebraic solutions to quadratic approximations of $J(X, y; w)$ as we did for $L^2$ regularization.

## $L^1$ Regularization

- Out simple linear model has a quadratic cost function that we can represent via its Taylor series.
- Alternately, we could imaging that this is a truncated Taylor series approximating the cost function of a more sophisticated model.
- The gradient in this setting is given by

$$\nabla_w \tilde{J}(w) = H(w - w^*)$$

- Because the $L^1$ penalty does not admit clean algebraic expressions in the case of a full general Hessian, we will also make the further simplifying assumption that the Hessian is a diagonal, $H = \mathrm{diag}([H_{1,1}, \ldots, H_{n,n}])$, where each $H_{i,i} > 0$.
- This assumption holds if the data for the linear regression problem has been preprocessed to remove all correlation between the input features, which may be accomplished using PCA.

## $L^1$ Regularization

- Out simple linear model has a quadratic cost function that we can represent via its Taylor series.
- Alternately, we could imaging that this is a truncated Taylor series approximating the cost function of a more sophisticated model.
- The gradient in this setting is given by

$$\nabla_w \tilde{J}(w) = H(w - w^*)$$

- Because the $L^1$ penalty does not admit clean algebraic expressions in the case of a full general Hessian, we will also make the further simplifying assumption that the Hessian is a diagonal, $H = \mathrm{diag}([H_{1,1}, \dots, H_{n,n}])$, where each $H_{i,i} > 0$.
- This assumption holds if the data for the linear regression problem has been preprocessed to remove all correlation between the input features, which may be accomplished using PCA.

# $L^1$ Regularization

- Out simple linear model has a quadratic cost function that we can represent via its Taylor series.
- Alternately, we could imaging that this is a truncated Taylor series approximating the cost function of a more sophisticated model.
- The gradient in this setting is given by

$$\nabla_w \tilde{J}(w) = H(w - w^*)$$

- Because the $L^1$ penalty does not admit clean algebraic expressions in the case of a full general Hessian, we will also make the further simplifying assumption that the Hessian is a diagonal, $H = \text{diag}([H_{1,1}, \ldots, H_{n,n}])$, where each $H_{i,i} > 0$.
- This assumption holds if the data for the linear regression problem has been preprocessed to remove all correlation between the input features, which may be accomplished using PCA.

# $L^1$ Regularization

- Out simple linear model has a quadratic cost function that we can represent via its Taylor series.
- Alternately, we could imaging that this is a truncated Taylor series approximating the cost function of a more sophisticated model.
- The gradient in this setting is given by

$$\nabla_w \tilde{J}(w) = \mathsf{H}(w - w^*)$$

- Because the $L^1$ penalty does not admit clean algebraic expressions in the case of a full general Hessian, we will also make the further simplifying assumption that the Hessian is a diagonal, $\mathsf{H} = \mathrm{diag}([H_{1,1}, \ldots, H_{n,n}])$, where each $H_{i,i} > 0$.
- This assumption holds if the data for the linear regression problem has been preprocessed to remove all correlation between the input features, which may be accomplished using PCA.

# $L^1$ Regularization

- Out simple linear model has a quadratic cost function that we can represent via its Taylor series.
- Alternately, we could imaging that this is a truncated Taylor series approximating the cost function of a more sophisticated model.
- The gradient in this setting is given by

$$\nabla_w \tilde{J}(w) = \mathsf{H}(w - w^*)$$

- Because the $L^1$ penalty does not admit clean algebraic expressions in the case of a full general Hessian, we will also make the further simplifying assumption that the Hessian is a diagonal, $\mathsf{H} = \mathrm{diag}([H_{1,1}, \ldots, H_{n,n}])$, where each $H_{i,i} > 0$.
- This assumption holds if the data for the linear regression problem has been preprocessed to remove all correlation between the input features, which may be accomplished using PCA.

# $L^1$ Regularization

- Our quadratic approximation of the $L^1$ regularized objective function decomposes into a sum over the parameters:

$$\tilde{J}(w; X, y) = J(w^*; X, y) + \sum_i \left[ \frac{1}{2} H_{i,i}(w_i - w_i^*)^2 + \alpha |w_i| \right]$$

- The problem of minimizing this approximate cost function has an analytical solution (for each dimension $i$), with the following form:

$$w_i = \text{sign}(w_i^*) max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$

## $L^1$ Regularization

- Our quadratic approximation of the $L^1$ regularized objective function decomposes into a sum over the parameters:

$$\tilde{J}(w; X, y) = J(w^*; X, y) + \sum_i \left[ \frac{1}{2} H_{i,i}(w_i - w_i^*)^2 + \alpha |w_i| \right]$$

- The problem of minimizing this approximate cost function has an analytical solution (for each dimension $i$), with the following form:

$$w_i = \text{sign}(w_i^*) max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$

# $L^1$ Regularization

$$w_i = \text{sign}(w_i^*) max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$

- Consider the situation where $w_i^* > 0$ for all $i$. There are two possible outcomes:

  1. The case where $w_i^* \leq \frac{\alpha}{H_{i,i}}$. Here the optimal value of $w_i$ under the regularized objective is simply $w_i = 0$. This occurs because the contribution of $J(w; X, y)$ to the regularized objective $\tilde{J}(w; X, y)$ is overwhelmed–in direction $i$–by the $L^1$ regularization which pushes the value of $w_i$ to zero.

  2. The case where $w_i^* > \frac{\alpha}{H_{i,i}}$. In this case, the regularization does not move the optimal value of $w_i$ to zero but instead it just shifts it in that direction by a distance equal to $\frac{\alpha}{H_{i,i}}$.

- A similar process happens when $w_i^* < 0$, but with the $L^1$ penalty making $w_i$ less negative by $\frac{\alpha}{H_{i,i}}$, or 0.

## $L^1$ Regularization

$$w_i = \text{sign}(w_i^*) max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$

- Consider the situation where $w_i^* > 0$ for all $i$. There are two possible outcomes:
  1. The case where $w_i^* \leq \frac{\alpha}{H_{i,i}}$. Here the optimal value of $w_i$ under the regularized objective is simply $w_i = 0$. This occurs because the contribution of $J(w; X, y)$ to the regularized objective $\tilde{J}(w; X, y)$ is overwhelmed–in direction $i$–by the $L^1$ regularization which pushes the value of $w_i$ to zero.
  2. The case where $w_i^* > \frac{\alpha}{H_{i,i}}$. In this case, the regularization does not move the optimal value of $w_i$ to zero but instead it just shifts it in that direction by a distance equal to $\frac{\alpha}{H_{i,i}}$.

- A similar process happens when $w_i^* < 0$, but with the $L^1$ penalty making $w_i$ less negative by $\frac{\alpha}{H_{i,i}}$, or 0.

# $L^1$ Regularization

$$w_i = \text{sign}(w_i^*) max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$

- Consider the situation where $w_i^* > 0$ for all $i$. There are two possible outcomes:
    1. The case where $w_i^* \leq \frac{\alpha}{H_{i,i}}$. Here the optimal value of $w_i$ under the regularized objective is simply $w_i = 0$. This occurs because the contribution of $J(w; X, y)$ to the regularized objective $\tilde{J}(w; X, y)$ is overwhelmed–in direction $i$–by the $L^1$ regularization which pushes the value of $w_i$ to zero.
    2. The case where $w_i^* > \frac{\alpha}{H_{i,i}}$. In this case, the regularization does not move the optimal value of $w_i$ to zero but instead it just shifts it in that direction by a distance equal to $\frac{\alpha}{H_{i,i}}$.
- A similar process happens when $w_i^* < 0$, but with the $L^1$ penalty making $w_i$ less negative by $\frac{\alpha}{H_{i,i}}$, or 0.

# $L^1$ Regularization

$$w_i = \text{sign}(w_i^*)max\left\{|w_i^*| - \frac{\alpha}{H_{i,i}}, 0\right\}$$

· Consider the situation where $w_i^* > 0$ for all $i$. There are two
possible outcomes:
   1. The case where $w_i^* \leq \frac{\alpha}{H_{i,i}}$. Here the optimal value of $w_i$ under the
      regularized objective is simply $w_i = 0$. This occurs because the
      contribution of $J(w; X, y)$ to the regularized objective $\tilde{J}(w; X, y)$ is
      overwhelmed–in direction $i$–by the $L^1$ regularization which pushes
      the value of $w_i$ to zero.
   2. The case where $w_i^* > \frac{\alpha}{H_{i,i}}$. In this case, the regularization does not
      move the optimal value of $w_i$ to zero but instead it just shifts it in
      that direction by a distance equal to $\frac{\alpha}{H_{i,i}}$.

· A similar process happens when $w_i^* < 0$, but with the $L^1$ penalty
making $w_i$ less negative by $\frac{\alpha}{H_{i,i}}$, or 0.

# $L^1$ Regularization

$$w_i = \text{sign}(w_i^*) max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$

- Consider the situation where $w_i^* > 0$ for all $i$. There are two possible outcomes:
  1. The case where $w_i^* \leq \frac{\alpha}{H_{i,i}}$. Here the optimal value of $w_i$ under the regularized objective is simply $w_i = 0$. This occurs because the contribution of $J(w; X, y)$ to the regularized objective $\tilde{J}(w; X, y)$ is overwhelmed–in direction $i$–by the $L^1$ regularization which pushes the value of $w_i$ to zero.
  2. The case where $w_i^* > \frac{\alpha}{H_{i,i}}$. In this case, the regularization does not move the optimal value of $w_i$ to zero but instead it just shifts it in that direction by a distance equal to $\frac{\alpha}{H_{i,i}}$.

- A similar process happens when $w_i^* < 0$, but with the $L^1$ penalty making $w_i$ less negative by $\frac{\alpha}{H_{i,i}}$, or 0.

# $L^1$ Regularization

- In comparison to $L^2$ regularization, $L^1$ regularization results in a solution that is more *sparse*.
- Sparsity in this context refers to the fact that some parameters have an optimal value of zero.
- The sparsity property induced by $L^1$ regularization has been used extensively as a *feature selection* mechanism.
- Feature selection simplifies a machine learning problem by choosing which subset of the available features should be used.
- In particular, the well known LASSO ([Tibshirani, 1996]) (least absolute shrinkage and selection operator) model integrates an $L^1$ penalty with a linear model and a least squares cost function.

# $L^1$ Regularization

- In comparison to $L^2$ regularization, $L^1$ regularization results in a solution that is more *sparse*.
- Sparsity in this context refers to the fact that some parameters have an optimal value of zero.
- The sparsity property induced by $L^1$ regularization has been used extensively as a *feature selection* mechanism.
- Feature selection simplifies a machine learning problem by choosing which subset of the available features should be used.
- In particular, the well known LASSO ([Tibshirani, 1996]) (least absolute shrinkage and selection operator) model integrates an $L^1$ penalty with a linear model and a least squares cost function.

22

# $L^1$ Regularization

- In comparison to $L^2$ regularization, $L^1$ regularization results in a solution that is more *sparse*.
- Sparsity in this context refers to the fact that some parameters have an optimal value of zero.
- The sparsity property induced by $L^1$ regularization has been used extensively as a *feature selection* mechanism.
- Feature selection simplifies a machine learning problem by choosing which subset of the available features should be used.
- In particular, the well known LASSO ([Tibshirani, 1996]) (least absolute shrinkage and selection operator) model integrates an $L^1$ penalty with a linear model and a least squares cost function.

22

## $L^1$ Regularization

- In comparison to $L^2$ regularization, $L^1$ regularization results in a solution that is more *sparse*.
- Sparsity in this context refers to the fact that some parameters have an optimal value of zero.
- The sparsity property induced by $L^1$ regularization has been used extensively as a *feature selection* mechanism.
- Feature selection simplifies a machine learning problem by choosing which subset of the available features should be used.
- In particular, the well known LASSO ([Tibshirani, 1996]) (least absolute shrinkage and selection operator) model integrates an $L^1$ penalty with a linear model and a least squares cost function.
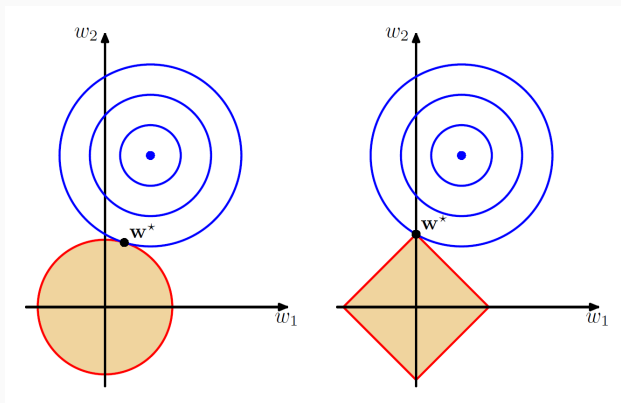
# $L^1$ Regularization

- In comparison to $L^2$ regularization, $L^1$ regularization results in a solution that is more *sparse*.
- Sparsity in this context refers to the fact that some parameters have an optimal value of zero.
- The sparsity property induced by $L^1$ regularization has been used extensively as a *feature selection* mechanism.
- Feature selection simplifies a machine learning problem by choosing which subset of the available features should be used.
- In particular, the well known LASSO ([Tibshirani, 1996]) (least absolute shrinkage and selection operator) model integrates an $L^1$ penalty with a linear model and a least squares cost function.

Fig. 2: Plot of the contours of the unregularized error function (blue) along with the constraint region for the quadratic regularizer on the left and the lasso regularizer on the right.

- Consider the cost function regularized by a parameter norm penalty:

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

- If we want to constrain $\Omega(\boldsymbol{\theta})$ to be less than some constant $k$, we could construct a generalized Language function

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha}; X, y) = J(\boldsymbol{\theta}; X, y) + \alpha(\Omega(\boldsymbol{\theta}) - k)$$

- The solution to the constrained problem is given by

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha})$$

- Consider the cost function regularized by a parameter norm penalty:

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

- If we want to constrain $\Omega(\boldsymbol{\theta})$ to be less than some constant $k$, we could construct a generalized Language function

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha}; X, y) = J(\boldsymbol{\theta}; X, y) + \alpha(\Omega(\theta) - k)$$

- The solution to the constrained problem is given by

$$\boldsymbol{\theta}^* = \arg\min_{\theta} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha})$$

# Norm Penalties as Constrained Optimization

- Consider the cost function regularized by a parameter norm penalty:

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

- If we want to constrain $\Omega(\boldsymbol{\theta})$ to be less than some constant $k$, we could construct a generalized Language function

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha}; X, y) = J(\boldsymbol{\theta}; X, y) + \alpha(\Omega(\theta) - k)$$

- The solution to the constrained problem is given by

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha})$$

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha})$$

- Solving this problem requires modifying both $\boldsymbol{\theta}$ and $\alpha$.

- Many different procedures are possible–some may use gradient descent, while others may use analytical solutions for where the gradient is zero–but in all procedures $\alpha$ must increase whenever $\Omega(\boldsymbol{\theta}) > k$ and decrease whenever $\Omega(\boldsymbol{\theta}) < k$.

- All positive $\alpha$ encourage $\Omega(\boldsymbol{\theta})$ to shrink.

- The optimal value $a^*$ will encourage $\Omega(\theta)$ to shrink, but not so strongly to make $\Omega(\boldsymbol{\theta})$ become less than $k$.

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha})$$

- Solving this problem requires modifying both $\boldsymbol{\theta}$ and $\alpha$.
- Many different procedures are possible–some may use gradient descent, while others may use analytical solutions for where the gradient is zero–but in all procedures $\alpha$ must increase whenever $\Omega(\boldsymbol{\theta}) > k$ and decrease whenever $\Omega(\boldsymbol{\theta}) < k$.
- All positive $\alpha$ encourage $\Omega(\boldsymbol{\theta})$ to shrink.
- The optimal value $a^*$ will encourage $\Omega(\boldsymbol{\theta})$ to shrink, but not so strongly to make $\Omega(\boldsymbol{\theta})$ become less than $k$.

## Norm Penalties as Constrained Optimization

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha})$$

- Solving this problem requires modifying both $\boldsymbol{\theta}$ and $\alpha$.
- Many different procedures are possible–some may use gradient descent, while others may use analytical solutions for where the gradient is zero–but in all procedures $\alpha$ must increase whenever $\Omega(\boldsymbol{\theta}) > k$ and decrease whenever $\Omega(\boldsymbol{\theta}) < k$.
- All positive $\alpha$ encourage $\Omega(\boldsymbol{\theta})$ to shrink.
- The optimal value $a^*$ will encourage $\Omega(\theta)$ to shrink, but not so strongly to make $\Omega(\boldsymbol{\theta})$ become less than $k$.

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha})$$

- Solving this problem requires modifying both $\boldsymbol{\theta}$ and $\alpha$.
- Many different procedures are possible–some may use gradient descent, while others may use analytical solutions for where the gradient is zero–but in all procedures $\alpha$ must increase whenever $\Omega(\boldsymbol{\theta}) > k$ and decrease whenever $\Omega(\boldsymbol{\theta}) < k$.
- All positive $\alpha$ encourage $\Omega(\boldsymbol{\theta})$ to shrink.
- The optimal value $a^*$ will encourage $\Omega(\theta)$ to shrink, but not so strongly to make $\Omega(\boldsymbol{\theta})$ become less than $k$.

## Regularization and Under-Constrained Problems

- In some cases, regularization is necessary.
- Many linear models in machine learning, including linear regression and PCA, depend on inverting the matrix $X^T X$.
- This is not possible whenever $X^T X$ is singular.
- This matrix can be singular whenever the data generating distribution truly has no variance in some direction, or when no variance in observed in some direction because there are fewer examples (rows of $X$) than input features (columns of $X$).
- In this case, many forms of regularization correspond to inverting $X^T X + \alpha I$ instead. This regularized matrix is guaranteed to be invertible.

## Regularization and Under-Constrained Problems

- In some cases, regularization is necessary.
- Many linear models in machine learning, including linear regression and PCA, depend on inverting the matrix $X^TX$.
- This is not possible whenever $X^TX$ is singular.
- This matrix can be singular whenever the data generating distribution truly has no variance in some direction, or when no variance in observed in some direction because there are fewer examples (rows of $X$) than input features (columns of $X$).
- In this case, many forms of regularization correspond to inverting $X^TX + \alpha I$ instead. This regularized matrix is guaranteed to be invertible.

# Regularization and Under-Constrained Problems

- In some cases, regularization is necessary.
- Many linear models in machine learning, including linear regression and PCA, depend on inverting the matrix $X^TX$.
- This is not possible whenever $X^TX$ is singular.
- This matrix can be singular whenever the data generating distribution truly has no variance in some direction, or when no variance in **observed** in some direction because there are fewer examples (rows of $X$) than input features (columns of $X$).
- In this case, many forms of regularization correspond to inverting $X^TX + \alpha I$ instead. This regularized matrix is guaranteed to be invertible.

## Regularization and Under-Constrained Problems

- In some cases, regularization is necessary.
- Many linear models in machine learning, including linear regression and PCA, depend on inverting the matrix $X^TX$.
- This is not possible whenever $X^TX$ is singular.
- This matrix can be singular whenever the data generating distribution truly has no variance in some direction, or when no variance in **observed** in some direction because there are fewer examples (rows of $X$) than input features (columns of $X$).
- In this case, many forms of regularization correspond to inverting $X^TX + \alpha I$ instead. This regularized matrix is guaranteed to be invertible.

## Regularization and Under-Constrained Problems

- In some cases, regularization is necessary.
- Many linear models in machine learning, including linear regression and PCA, depend on inverting the matrix $X^T X$.
- This is not possible whenever $X^T X$ is singular.
- This matrix can be singular whenever the data generating distribution truly has no variance in some direction, or when no variance in **observed** in some direction because there are fewer examples (rows of $X$) than input features (columns of $X$).
- In this case, many forms of regularization correspond to inverting $X^T X + \alpha I$ instead. This regularized matrix is guaranteed to be invertible.

## Regularization and Under-Constrained Problems

- These linear problems have closed form solutions when the relevant matrix is invertible.

- It is also possible for a problem with no closed form solution to be underdetermined.

- An example is logistic regression applied to a problem where the classes are linearly separable.

- If a weight vector $w$ is able to achieve perfect classification, then $2w$ will also achieve perfect classification and higher likelihood.

- An iterative optimization procedure like SGD will continually increase the magnitude of $w$ and, in theory, will never halt.

## Regularization and Under-Constrained Problems

- These linear problems have closed form solutions when the relevant matrix is invertible.
- It is also possible for a problem with no closed form solution to be underdetermined.
- An example is logistic regression applied to a problem where the classes are linearly separable.
- If a weight vector $w$ is able to achieve perfect classification, then $2w$ will also achieve perfect classification and higher likelihood.
- An iterative optimization procedure like SGD will continually increase the magnitude of $w$ and, in theory, will never halt.

## Regularization and Under-Constrained Problems

- These linear problems have closed form solutions when the relevant matrix is invertible.
- It is also possible for a problem with no closed form solution to be underdetermined.
- An example is logistic regression applied to a problem where the classes are linearly separable.
- If a weight vector $w$ is able to achieve perfect classification, then $2w$ will also achieve perfect classification and higher likelihood.
- An iterative optimization procedure like SGD will continually increase the magnitude of $w$ and, in theory, will never halt.

## Regularization and Under-Constrained Problems

- These linear problems have closed form solutions when the relevant matrix is invertible.
- It is also possible for a problem with no closed form solution to be underdetermined.
- An example is logistic regression applied to a problem where the classes are linearly separable.
- If a weight vector $w$ is able to achieve perfect classification, then $2w$ will also achieve perfect classification and higher likelihood.
- An iterative optimization procedure like SGD will continually increase the magnitude of $w$ and, in theory, will never halt.

## Regularization and Under-Constrained Problems

- These linear problems have closed form solutions when the relevant matrix is invertible.
- It is also possible for a problem with no closed form solution to be underdetermined.
- An example is logistic regression applied to a problem where the classes are linearly separable.
- If a weight vector $w$ is able to achieve perfect classification, then $2w$ will also achieve perfect classification and higher likelihood.
- An iterative optimization procedure like SGD will continually increase the magnitude of $w$ and, in theory, will never halt.

## Regularization and Under-Constrained Problems

- We can solve underdetermined linear equations using the Moore-Penrose pseudoinverse. Recall that one definition of the pseudoinverse $X^+$ of a matrix $X$ is

$$X^+ = \lim_{a \to 0} (X^T X + \alpha I)^{-1} X^T$$

- We can now recognize this equation as performing linear regression with weight decay.

- We can interpret the pseudoinverse as stabilizing underdetermined problems using regularization.

## Regularization and Under-Constrained Problems

- We can solve underdetermined linear equations using the Moore-Penrose pseudoinverse. Recall that one definition of the pseudoinverse $X^+$ of a matrix $X$ is

$$X^+ = \lim_{a \to 0}(X^T X + \alpha I)^{-1}X^T$$

- We can now recognize this equation as performing linear regression with weight decay.

- We can interpret the pseudoinverse as stabilizing underdetermined problems using regularization.

## Regularization and Under-Constrained Problems

- We can solve underdetermined linear equations using the Moore-Penrose pseudoinverse. Recall that one definition of the pseudoinverse $X^+$ of a matrix $X$ is

$$X^+ = \lim_{a \to 0}(X^TX + \alpha I)^{-1}X^T$$

- We can now recognize this equation as performing linear regression with weight decay.

- We can interpret the pseudoinverse as stabilizing underdetermined problems using regularization.

## Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on more data.

- In practice, it is limited.

- Create fake data and add it to the training set.

- This approach is easiest for classification.

- A classifier needs to take a complicated, high dimensional input $x$ and summarize it with a single category identity $y$.

- This means that the main task facing a classifier is to be invariant to a wide variety of transformations.

- We can generate new $(x, y)$ pairs easily just by transforming the $x$ inputs in our training set.

## Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on more data.

- In practice, it is limited.

- Create fake data and add it to the training set.

- This approach is easiest for classification.

- A classifier needs to take a complicated, high dimensional input $x$ and summarize it with a single category identity $y$.

- This means that the main task facing a classifier is to be invariant to a wide variety of transformations.

- We can generate new ($x, y$) pairs easily just by transforming the $x$ inputs in our training set.

## Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on more data.
- In practice, it is limited.
- Create fake data and add it to the training set.
- This approach is easiest for classification.
- A classifier needs to take a complicated, high dimensional input *x* and summarize it with a single category identity *y*.
- This means that the main task facing a classifier is to be invariant to a wide variety of transformations.
- We can generate new (*x*, *y*) pairs easily just by transforming the *x* inputs in our training set.

## Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on more data.
- In practice, it is limited.
- Create fake data and add it to the training set.
- This approach is easiest for classification.
- A classifier needs to take a complicated, high dimensional input $x$ and summarize it with a single category identity $y$.
- This means that the main task facing a classifier is to be invariant to a wide variety of transformations.
- We can generate new $(x, y)$ pairs easily just by transforming the $x$ inputs in our training set.

## Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on more data.
- In practice, it is limited.
- Create fake data and add it to the training set.
- This approach is easiest for classification.
- A classifier needs to take a complicated, high dimensional input $x$ and summarize it with a single category identity $y$.
- This means that the main task facing a classifier is to be invariant to a wide variety of transformations.
- We can generate new $(x, y)$ pairs easily just by transforming the $x$ inputs in our training set.

## Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on more data.
- In practice, it is limited.
- Create fake data and add it to the training set.
- This approach is easiest for classification.
- A classifier needs to take a complicated, high dimensional input $x$ and summarize it with a single category identity $y$.
- This means that the main task facing a classifier is to be invariant to a wide variety of transformations.
- We can generate new $(x, y)$ pairs easily just by transforming the $x$ inputs in our training set.

## Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on more data.
- In practice, it is limited.
- Create fake data and add it to the training set.
- This approach is easiest for classification.
- A classifier needs to take a complicated, high dimensional input *x* and summarize it with a single category identity *y*.
- This means that the main task facing a classifier is to be invariant to a wide variety of transformations.
- We can generate new (*x*, *y*) pairs easily just by transforming the *x* inputs in our training set.

## Dataset Augmentation

- Dataset augmentation has been a particularly effective technique for a specific classification problem: object recognition.

- Images are high dimensional and include an enormous variety of factors of variation, many of which can be easily simulated.

- One must be careful not to apply transformations that would change the correct class. (e.g. '6' and '9', 'b' and 'd').

## Dataset Augmentation

- Dataset augmentation has been a particularly effective technique for a specific classification problem: object recognition.
- Images are high dimensional and include an enormous variety of factors of variation, many of which can be easily simulated.
- One must be careful not to apply transformations that would change the correct class. (e.g. '6' and '9', 'b' and 'd').

## Dataset Augmentation

- Dataset augmentation has been a particularly effective technique for a specific classification problem: object recognition.
- Images are high dimensional and include an enormous variety of factors of variation, many of which can be easily simulated.
- One must be careful not to apply transformations that would change the correct class. (e.g. '6' and '9', 'b' and 'd').

## Dataset Augmentation

- Dataset augmentation is effective for speech recognition task as well (Jaitly and Hinton [2013]).
- Inject noise in the input to a neural network can also be seen as a form of data augmentation (Sietsma and Dow [1991]).
- For many classification and even some regression tasks, the task should still be possible to solve even if small random noise is added to the input.
- One way to improve the robustness of neural networks is simply to train them with random noise applied to their inputs.
- Input noise injection is part of some unsupervised learning algorithms such as the denoising autoencoder (Vincent et al. [2008]).
- Dropout, a powerful regularization strategy can be seen as a process of constructing new inputs by **multiplying** by noise.

## Dataset Augmentation

- Dataset augmentation is effective for speech recognition task as well (Jaitly and Hinton [2013]).
- Inject noise in the input to a neural network can also be seen as a form of data augmentation (Sietsma and Dow [1991]).
- For many classification and even some regression tasks, the task should still be possible to solve even if small random noise is added to the input.
- One way to improve the robustness of neural networks is simply to train them with random noise applied to their inputs.
- Input noise injection is part of some unsupervised learning algorithms such as the denoising autoencoder (Vincent et al. [2008]).
- Dropout, a powerful regularization strategy can be seen as a process of constructing new inputs by **multiplying** by noise.

## Dataset Augmentation

- Dataset augmentation is effective for speech recognition task as well (Jaitly and Hinton [2013]).
- Inject noise in the input to a neural network can also be seen as a form of data augmentation (Sietsma and Dow [1991]).
- For many classification and even some regression tasks, the task should still be possible to solve even if small random noise is added to the input.
- One way to improve the robustness of neural networks is simply to train them with random noise applied to their inputs.
- Input noise injection is part of some unsupervised learning algorithms such as the denoising autoencoder (Vincent et al. [2008]).
- Dropout, a powerful regularization strategy can be seen as a process of constructing new inputs by **multiplying** by noise.

## Dataset Augmentation

- Dataset augmentation is effective for speech recognition task as well (Jaitly and Hinton [2013]).
- Inject noise in the input to a neural network can also be seen as a form of data augmentation (Sietsma and Dow [1991]).
- For many classification and even some regression tasks, the task should still be possible to solve even if small random noise is added to the input.
- One way to improve the robustness of neural networks is simply to train them with random noise applied to their inputs.
- Input noise injection is part of some unsupervised learning algorithms such as the denoising autoencoder (Vincent et al. [2008]).
- Dropout, a powerful regularization strategy can be seen as a process of constructing new inputs by **multiplying** by noise.

## Dataset Augmentation

- Dataset augmentation is effective for speech recognition task as well (Jaitly and Hinton [2013]).
- Inject noise in the input to a neural network can also be seen as a form of data augmentation (Sietsma and Dow [1991]).
- For many classification and even some regression tasks, the task should still be possible to solve even if small random noise is added to the input.
- One way to improve the robustness of neural networks is simply to train them with random noise applied to their inputs.
- Input noise injection is part of some unsupervised learning algorithms such as the denoising autoencoder (Vincent et al. [2008]).
- Dropout, a powerful regularization strategy can be seen as a process of constructing new inputs by **multiplying** by noise.

## Dataset Augmentation

- Dataset augmentation is effective for speech recognition task as well (Jaitly and Hinton [2013]).
- Inject noise in the input to a neural network can also be seen as a form of data augmentation (Sietsma and Dow [1991]).
- For many classification and even some regression tasks, the task should still be possible to solve even if small random noise is added to the input.
- One way to improve the robustness of neural networks is simply to train them with random noise applied to their inputs.
- Input noise injection is part of some unsupervised learning algorithms such as the denoising autoencoder (Vincent et al. [2008]).
- Dropout, a powerful regularization strategy can be seen as a process of constructing new inputs by **multiplying** by noise.

## Dataset Augmentation

- When comparing machine learning benchmark results, it is important to take the effect of dataset augmentation into account.

- Often, hand-designed dataset augmentation schemes can dramatically reduce the generalization error.

- When comparing machine learning algorithm A and machine learning algorithm B, it is necessary to make sure that both algorithms were evaluated using the same hand-designed dataset augmentation schemes.

## Dataset Augmentation

- When comparing machine learning benchmark results, it is important to take the effect of dataset augmentation into account.
- Often, hand-designed dataset augmentation schemes can dramatically reduce the generalization error.
- When comparing machine learning algorithm A and machine learning algorithm B, it is necessary to make sure that both algorithms were evaluated using the same hand-designed dataset augmentation schemes.

# Dataset Augmentation

- When comparing machine learning benchmark results, it is important to take the effect of dataset augmentation into account.
- Often, hand-designed dataset augmentation schemes can dramatically reduce the generalization error.
- When comparing machine learning algorithm A and machine learning algorithm B, it is necessary to make sure that both algorithms were evaluated using the same hand-designed dataset augmentation schemes.

## Noise Robustness

- For some models, the addition of noise with infinitesimal variance at the input of the model is equivalent to imposing a penalty on the norm of the weights (Bishop [1995b,a]).

- Noise injection can be much more powerful than simply shrinking the parameters, especially when the noise is added to the hidden units.

- Noise applied to the hidden units is such an important topic; the dropout algorithm describe later.

- Another way that noise can be added into the weights.

- This technique has been used primarily in the context of recurrent neural networks (Jim et al. [1996], Graves [2011]).

- This can also be interpreted as equivalent (under some assumptions) to a more traditional form of regularization.

## Noise Robustness

- For some models, the addition of noise with infinitesimal variance at the input of the model is equivalent to imposing a penalty on the norm of the weights (Bishop [1995b,a]).

- Noise injection can be much more powerful than simply shrinking the parameters, especially when the noise is added to the hidden units.

- Noise applied to the hidden units is such an important topic; the dropout algorithm describe later.

- Another way that noise can be added into the weights.

- This technique has been used primarily in the context of recurrent neural networks (Jim et al. [1996], Graves [2011]).

- This can also be interpreted as equivalent (under some assumptions) to a more traditional form of regularization.

## Noise Robustness

- For some models, the addition of noise with infinitesimal variance at the input of the model is equivalent to imposing a penalty on the norm of the weights (Bishop [1995b,a]).

- Noise injection can be much more powerful than simply shrinking the parameters, especially when the noise is added to the hidden units.

- Noise applied to the hidden units is such an important topic; the dropout algorithm describe later.

- Another way that noise can be added into the weights.

- This technique has been used primarily in the context of recurrent neural networks (Jim et al. [1996], Graves [2011]).

- This can also be interpreted as equivalent (under some assumptions) to a more traditional form of regularization.

## Noise Robustness

- For some models, the addition of noise with infinitesimal variance at the input of the model is equivalent to imposing a penalty on the norm of the weights (Bishop [1995b,a]).
- Noise injection can be much more powerful than simply shrinking the parameters, especially when the noise is added to the hidden units.
- Noise applied to the hidden units is such an important topic; the dropout algorithm describe later.
- Another way that noise can be added into the weights.
- This technique has been used primarily in the context of recurrent neural networks (Jim et al. [1996], Graves [2011]).
- This can also be interpreted as equivalent (under some assumptions) to a more traditional form of regularization.

## Noise Robustness

- For some models, the addition of noise with infinitesimal variance at the input of the model is equivalent to imposing a penalty on the norm of the weights (Bishop [1995b,a]).
- Noise injection can be much more powerful than simply shrinking the parameters, especially when the noise is added to the hidden units.
- Noise applied to the hidden units is such an important topic; the dropout algorithm describe later.
- Another way that noise can be added into the weights.
- This technique has been used primarily in the context of recurrent neural networks (Jim et al. [1996], Graves [2011]).
- This can also be interpreted as equivalent (under some assumptions) to a more traditional form of regularization.

## Noise Robustness

- For some models, the addition of noise with infinitesimal variance at the input of the model is equivalent to imposing a penalty on the norm of the weights (Bishop [1995b,a]).
- Noise injection can be much more powerful than simply shrinking the parameters, especially when the noise is added to the hidden units.
- Noise applied to the hidden units is such an important topic; the dropout algorithm describe later.
- Another way that noise can be added into the weights.
- This technique has been used primarily in the context of recurrent neural networks (Jim et al. [1996], Graves [2011]).
- This can also be interpreted as equivalent (under some assumptions) to a more traditional form of regularization.

## Noise Robustness

- We study the regression setting, where we wish to train a function $\tilde{y}(x)$ that maps a set of features $x$ to a scalar using the least-squares cost function between the model predictions $\tilde{y}(x)$ and the true values $y$:

$$J = \mathbb{E}_{p(x,y)} \left[ (\hat{y}(x) - y)^2 \right]$$

- The training set with $m$ examples: $(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})$.
- We now assume that with each input presentation we also include a random perturbation $\epsilon_W \, \mathcal{N}(\epsilon; 0, \eta I)$ of the network weights.
- We denote the perturbed model as $\hat{y}_{\epsilon_W}(x)$. The objective function thus becomes:

$$\begin{aligned} \tilde{J}_W &= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ (\hat{y}_{\epsilon_W}(x) - y)^2 \right] \\ &= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ \hat{y}_{\epsilon_W}^2 (x - 2y\hat{y}_{\epsilon_W} + y^2) \right] \end{aligned}$$

## Noise Robustness

- We study the regression setting, where we wish to train a function $\tilde{y}(x)$ that maps a set of features $x$ to a scalar using the least-squares cost function between the model predictions $\tilde{y}(x)$ and the true values $y$:

$$J = \mathbb{E}_{p(x,y)}\left[(\hat{y}(x) - y)^2\right]$$

- The training set with $m$ examples: $(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})$.
- We now assume that with each input presentation we also include a random perturbation $\epsilon_W \, \mathcal{N}(\epsilon; 0, \eta I)$ of the network weights.
- We denote the perturbed model as $\hat{y}_{\epsilon_W}(x)$. The objective function thus becomes:

$$\begin{aligned}
\tilde{J}_W &= \mathbb{E}_{p(x,y,\epsilon_W)}\left[(\hat{y}_{\epsilon_W}(x) - y)^2\right] \\
&= \mathbb{E}_{p(x,y,\epsilon_W)}\left[\hat{y}_{\epsilon_W}^2(x - 2y\hat{y}_{\epsilon_W} + y^2)\right]
\end{aligned}$$

## Noise Robustness

- We study the regression setting, where we wish to train a function $\tilde{y}(x)$ that maps a set of features $x$ to a scalar using the least-squares cost function between the model predictions $\tilde{y}(x)$ and the true values $y$:

$$J = \mathbb{E}_{p(x,y)} \left[ (\hat{y}(x) - y)^2 \right]$$

- The training set with $m$ examples: $(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})$.
- We now assume that with each input presentation we also include a random perturbation $\epsilon_W \, \mathcal{N}(\epsilon; 0, \eta I)$ of the network weights.
- We denote the perturbed model as $\hat{y}_{\epsilon_W}(x)$. The objective function thus becomes:

$$\begin{aligned} \tilde{J}_W &= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ (\hat{y}_{\epsilon_W}(x) - y)^2 \right] \\ &= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ \hat{y}_{\epsilon_W}^2 (x - 2y\hat{y}_{\epsilon_W} + y^2) \right] \end{aligned}$$

## Noise Robustness

- We study the regression setting, where we wish to train a function $\tilde{y}(x)$ that maps a set of features $x$ to a scalar using the least-squares cost function between the model predictions $\tilde{y}(x)$ and the true values $y$:

$$J = \mathbb{E}_{p(x,y)} \left[ (\hat{y}(x) - y)^2 \right]$$

- The training set with $m$ examples: $(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})$.
- We now assume that with each input presentation we also include a random perturbation $\epsilon_W \, \mathcal{N}(\epsilon; 0, \eta I)$ of the network weights.
- We denote the perturbed model as $\hat{y}_{\epsilon_W}(x)$. The objective function thus becomes:

$$\begin{aligned} \tilde{J}_W &= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ (\hat{y}_{\epsilon_W}(x) - y)^2 \right] \\ &= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ \hat{y}_{\epsilon_W}^2 (x - 2y\hat{y}_{\epsilon_W} + y^2) \right] \end{aligned}$$

## Noise Robustness

- We study the regression setting, where we wish to train a function $\tilde{y}(x)$ that maps a set of features $x$ to a scalar using the least-squares cost function between the model predictions $\tilde{y}(x)$ and the true values $y$:

$$J = \mathbb{E}_{p(x,y)} \left[ (\hat{y}(x) - y)^2 \right]$$

- The training set with $m$ examples: $(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})$.
- We now assume that with each input presentation we also include a random perturbation $\epsilon_W \, \mathcal{N}(\epsilon; 0, \eta I)$ of the network weights.
- We denote the perturbed model as $\hat{y}_{\epsilon_W}(x)$. The objective function thus becomes:

$$\begin{aligned}
\tilde{J}_W &= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ (\hat{y}_{\epsilon_W}(x) - y)^2 \right] \\
&= \mathbb{E}_{p(x,y,\epsilon_W)} \left[ \hat{y}_{\epsilon_W}^2 (x - 2y\hat{y}_{\epsilon_W} + y^2) \right]
\end{aligned}$$

## Noise Robustness

$$\tilde{J}_W = \mathbb{E}_{p(x,y,\epsilon_W)} \left[ \hat{y}_{\epsilon_W}^2 (x - 2y\hat{y}_{\epsilon_W} + y^2) \right]$$

- For small $\eta$, the minimization of $J$ with added weight noise (with covariance $\eta I$) is equivalent to minimization of $J$ with an additional regularization.

- This form of regularization encourages the parameters to go to regions of parameter space where small pertubations of the weights have a relatively small influence on the output.

- In other words, it pushes the model weights, finding points that are not merely minima, but minima surrounded by flat regions (Hochreiter et al. [1995]).

## Noise Robustness

$$\tilde{J}_W = \mathbb{E}_{p(x,y,\epsilon_W)} \left[ \hat{y}_{\epsilon_W}^2 (x - 2y\hat{y}_{\epsilon_W} + y^2) \right]$$

- For small $\eta$, the minimization of $J$ with added weight noise (with covariance $\eta I$) is equivalent to minimization of $J$ with an additional regularization.
- This form of regularization encourages the parameters to go to regions of parameter space where small pertubations of the weights have a relatively small influence on the output.
- In other words, it pushes the model weights, finding points that are not merely minima, but minima surrounded by flat regions (Hochreiter et al. [1995]).

$$\tilde{J}_W = \mathbb{E}_{p(x,y,\epsilon_W)} \left[ \hat{y}^2_{\epsilon_W}(x - 2y\hat{y}_{\epsilon_W} + y^2) \right]$$

- For small $\eta$, the minimization of $J$ with added weight noise (with covariance $\eta I$) is equivalent to minimization of $J$ with an additional regularization.
- This form of regularization encourages the parameters to go to regions of parameter space where small perturbations of the weights have a relatively small influence on the output.
- In other words, it pushes the model weights, finding points that are not merely minima, but minima surrounded by flat regions (Hochreiter et al. [1995]).

- Most datasets have some amount of mistakes in the *y* labels.

- It can be harmful to maximize $\log p(y|x)$ when *y* is a mistake.

- One way to prevent this is to explicitly model the noise on the labels.

- For example, we can assume that for some small constant $\epsilon$, the training set label *y* is correct with probability $1 - \epsilon$, and otherwise any of the other possible labels might be correct.

- Most datasets have some amount of mistakes in the $y$ labels.
- It can be harmful to maximize $\log p(y|x)$ when $y$ is a mistake.
- One way to prevent this is to explicitly model the noise on the labels.
- For example, we can assume that for some small constant $\epsilon$, the training set label $y$ is correct with probability $1 - \epsilon$, and otherwise any of the other possible labels might be correct.

- Most datasets have some amount of mistakes in the $y$ labels.
- It can be harmful to maximize $\log p(y|x)$ when $y$ is a mistake.
- One way to prevent this is to explicitly model the noise on the labels.
- For example, we can assume that for some small constant $\epsilon$, the training set label $y$ is correct with probability $1 - \epsilon$, and otherwise any of the other possible labels might be correct.

- Most datasets have some amount of mistakes in the $y$ labels.
- It can be harmful to maximize $\log p(y|\boldsymbol{x})$ when $y$ is a mistake.
- One way to prevent this is to explicitly model the noise on the labels.
- For example, we can assume that for some small constant $\epsilon$, the training set label $y$ is correct with probability $1 - \epsilon$, and otherwise any of the other possible labels might be correct.

## Semi-Supervised Learning

- In the paradigm of semi-supervised learning, both unlabeled examples from $P(x)$ and labeled examples from $P(x, y)$ are used to estimate $P(y|x)$ or predict $y$ from $x$.

- One can construct models in which a generative model of either ($x$) or $P(x, y)$ shares parameters with a discriminative model of $P(y|x)$.

- The generative criterion then express a particular form of prior belief about the solution to the supervised learning problem, namely that the structure of $P(x)$ is connected to the structure of $P(y|x)$ in a way that is captured by the shared parameterization.

- By controlling how much of the generative criterion is included in the total criterion, one can find a better trade-off than with a purely generative or purely discriminative training criterion.

- Hinton and Salakhutdinov [2008] describe a method for learning the kernel function of a kernel machine used for regression, in which the usage of unlabeled examples for modeling $P(x)$ improves $P(y|x)$ quite significantly.

37

## Semi-Supervised Learning

- In the paradigm of semi-supervised learning, both unlabeled examples from $P(x)$ and labeled examples from $P(x, y)$ are used to estimate $P(y|x)$ or predict $y$ from $x$.
- One can construct models in which a generative model of either $(x)$ or $P(x, y)$ shares parameters with a discriminative model of $P(y|x)$.
- The generative criterion then express a particular form of prior belief about the solution to the supervised learning problem, namely that the structure of $P(x)$ is connected to the structure of $P(y|x)$ in a way that is captured by the shared parameterization.
- By controlling how much of the generative criterion is included in the total criterion, one can find a better trade-off than with a purely generative or purely discriminative training criterion.
- Hinton and Salakhutdinov [2008] describe a method for learning the kernel function of a kernel machine used for regression, in which the usage of unlabeled examples for modeling $P(x)$ improves $P(y|x)$ quite significantly.

37

## Semi-Supervised Learning

- In the paradigm of semi-supervised learning, both unlabeled examples from $P(x)$ and labeled examples from $P(x, y)$ are used to estimate $P(y|x)$ or predict $y$ from $x$.
- One can construct models in which a generative model of either ($x$) or $P(x, y)$ shares parameters with a discriminative model of $P(y|x)$.
- The generative criterion then express a particular form of prior belief about the solution to the supervised learning problem, namely that the structure of $P(x)$ is connected to the structure of $P(y|x)$ in a way that is captured by the shared parameterization.
- By controlling how much of the generative criterion is included in the total criterion, one can find a better trade-off than with a purely generative or purely discriminative training criterion.
- Hinton and Salakhutdinov [2008] describe a method for learning the kernel function of a kernel machine used for regression, in which the usage of unlabeled examples for modeling $P(x)$ improves $P(y|x)$ quite significantly.

## Semi-Supervised Learning

- In the paradigm of semi-supervised learning, both unlabeled examples from $P(x)$ and labeled examples from $P(x, y)$ are used to estimate $P(y|x)$ or predict $y$ from $x$.
- One can construct models in which a generative model of either $(x)$ or $P(x, y)$ shares parameters with a discriminative model of $P(y|x)$.
- The generative criterion then express a particular form of prior belief about the solution to the supervised learning problem, namely that the structure of $P(x)$ is connected to the structure of $P(y|x)$ in a way that is captured by the shared parameterization.
- By controlling how much of the generative criterion is included in the total criterion, one can find a better trade-off than with a purely generative or purely discriminative training criterion.
- Hinton and Salakhutdinov [2008] describe a method for learning the kernel function of a kernel machine used for regression, in which the usage of unlabeled examples for modeling $P(x)$ improves $P(y|x)$ quite significantly.
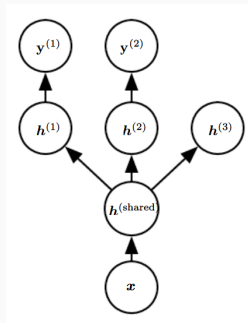
## Semi-Supervised Learning

- In the paradigm of semi-supervised learning, both unlabeled examples from $P(x)$ and labeled examples from $P(x, y)$ are used to estimate $P(y|x)$ or predict $y$ from $x$.
- One can construct models in which a generative model of either ($x$) or $P(x, y)$ shares parameters with a discriminative model of $P(y|x)$.
- The generative criterion then express a particular form of prior belief about the solution to the supervised learning problem, namely that the structure of $P(x)$ is connected to the structure of $P(y|x)$ in a way that is captured by the shared parameterization.
- By controlling how much of the generative criterion is included in the total criterion, one can find a better trade-off than with a purely generative or purely discriminative training criterion.
- Hinton and Salakhutdinov [2008] describe a method for learning the kernel function of a kernel machine used for regression, in which the usage of unlabeled examples for modeling $P(x)$ improves $P(y|x)$ quite significantly.

## Multi-Task Learning

- Multi-task learning is a way to improve generalization by pooling the examples arising out of several tasks.
- In the same way that additional training examples put more pressure on the parameters of the model towards values that generalize well, when part of a model is shared across tasks, that part of the model is more constrained towards good values, often yielding better generalization.

# Multi-Task Learning

- Here is a very common form of multi-task learning.
- Different supervised tasks (predicting $y^{(i)}$ given $x$) share the same input $x$, as well as some intermediate-level representation $h^{(\text{shared})}$ capturing a common pool of factors.

- The model has two kinds of parts:
    1. Task-specific parameters (which only benefit from the examples of their task to achieve good generalization). These are the upper layers.
    2. Generic parameters, shared across all the tasks (which benefit from the pooled data of all the tasks). These are the lower years.



- The factors that explain the variations are shared across two or more tasks.

## References

Chris M Bishop. Training with noise is equivalent to tikhonov regularization. *Neural computation*, 7(1):108–116, 1995a.

Christopher M Bishop. Regularization and complexity control in feed-forward networks. 1995b.

Alex Graves. Practical variational inference for neural networks. In *Advances in Neural Information Processing Systems*, pages 2348–2356, 2011.

Geoffrey E Hinton and Ruslan R Salakhutdinov. Using deep belief nets to learn covariance kernels for gaussian processes. In *Advances in neural information processing systems*, pages 1249–1256, 2008.

Sepp Hochreiter, Jürgen Schmidhuber, et al. Simplifying neural nets by discovering flat minima. *Advances in Neural Information Processing Systems*, pages 529–536, 1995.

Navdeep Jaitly and Geoffrey E Hinton. Vocal tract length perturbation (vtlp) improves speech recognition. In *Proc. ICML Workshop on Deep Learning for Audio, Speech and Language*, 2013.

Kam-Chuen Jim, C Lee Giles, and Bill G Horne. An analysis of noise in recurrent neural networks: Convergence and generalization. *IEEE Transactions on neural networks*, 7(6):1424–1438, 1996.

Jocelyn Sietsma and Robert JF Dow. Creating artificial neural networks that generalize. *Neural networks*, 4(1):67–79, 1991.

Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.

Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103. ACM, 2008.