



## Contents

- The goal of deep learning
- Why nonlinear? How?
- General framework of feedforward nets
- Learning tasks
  - Regression
  - Classification
- Gradient Descent, mini batch, SGD
  - Backpropagation
- Regularization
  - Early stopping
  - Dropout

3

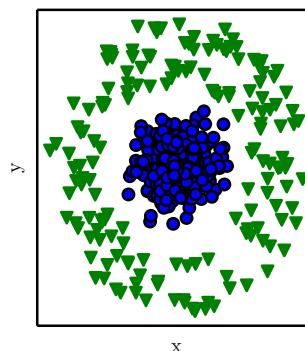
## References

- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press.
- Lectures: [http://www.deeplearningbook.org/lecture\\_slides.html](http://www.deeplearningbook.org/lecture_slides.html)
- <https://www.tensorflow.org/tutorials/word2vec>

2

## Representations Matter

Cartesian coordinates



Polar coordinates

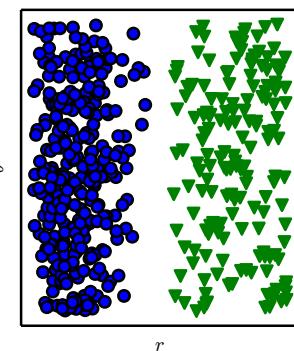


Figure 1.1

4

## Learning Multiple Components

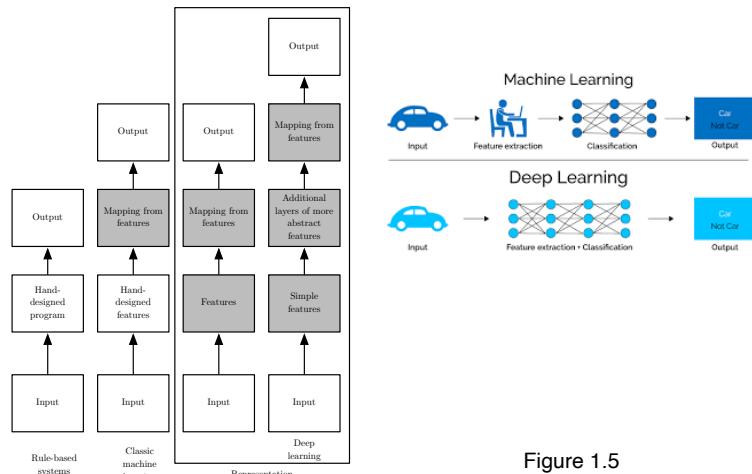


Figure 1.5

5

## The goal of Deep Learning

Look for the (nonlinear) function

$$f : \mathbb{R}^m \longrightarrow \mathbb{R}^n, f(\mathbf{x}, \theta) = f_\theta(\mathbf{x})$$

With the greatest **consistency** with the data

$$D = \{(\mathbf{x}, \mathbf{y}) : \mathbf{x} \sim \mathbf{y}\}$$

What is meaning of **consistency** in this context?

- Ease estimations  $\mathbf{y} \approx \Psi(f(\mathbf{x}, \theta))$
- Make minimum a loss function (*inconsistency with the data*)

6

## The goal of Deep Learning

Finally, the goal of Deep Learning is to solve an optimization equation like this

$$\theta \leftarrow \operatorname{argmin}_\theta \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(y, f_\theta(\mathbf{x}))$$

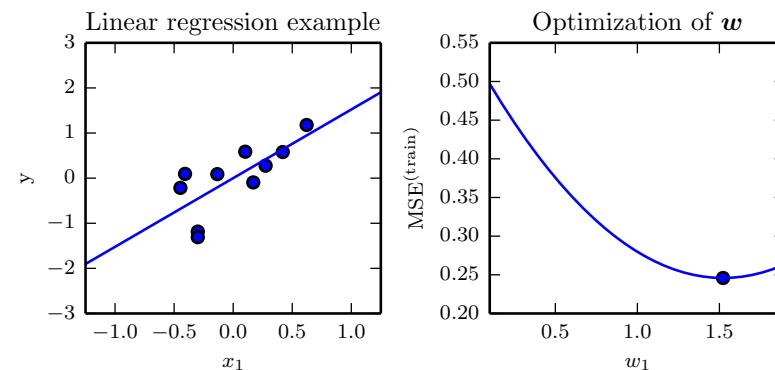
Example: linear regression ( $n = 1$ ) with a quadratic loss

$$f_w(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$$

$$\text{loss}(\mathbf{y}, f_w(\mathbf{x})) = \frac{1}{2}(\mathbf{y} - f_w(\mathbf{x}))^2$$

7

## Linear Regression



$$y \approx w_1 x_1$$

Figure 5.1

8

# Objetivo del Deep Learning

Given

$$D = \{(x, y) : x \rightsquigarrow y\}$$

Given

Design and find a meta parameter  $\Theta$

$$f : \mathbb{R}^m \longrightarrow \mathbb{R}^n, f(x, \theta) = f_\theta(x)$$

Design

Such that minimizes the mean of the loss function

$$\theta \leftarrow \operatorname{argmin}_\theta \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(y, f_\theta(x))$$

Program

that yields to

$$y \approx \Psi(f(x, \theta))$$

Design

9

# Representations Matter

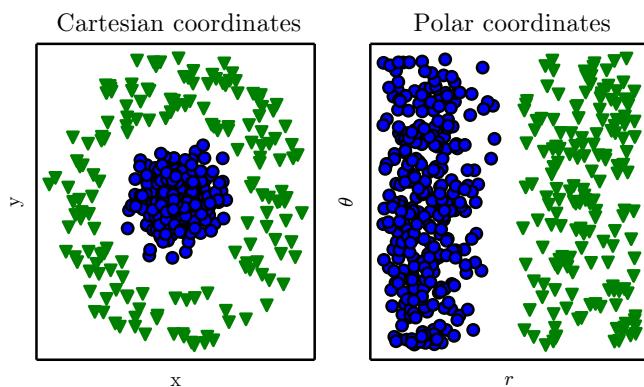


Figure 1.1

11

# Why nonlinear?

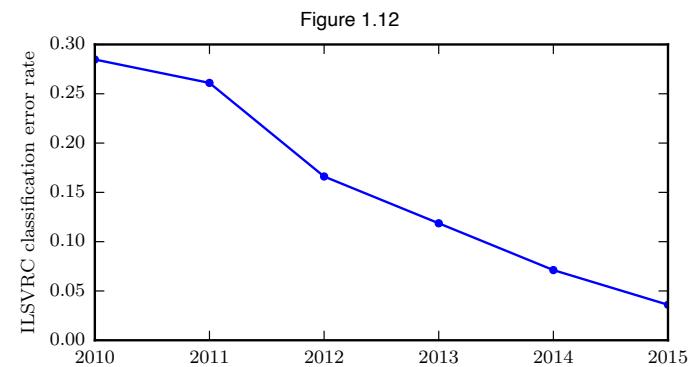
When the function  $f$  is nonlinear, the optimization problem is not convex

The SGD (all versions) only guarantee the convergence (find a solution) whenever the optimization problem is convex

Why nonlinear? **Because it makes up for the result**

10

# Solving Object Recognition

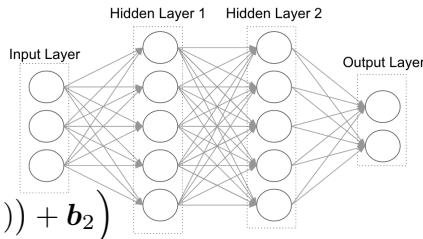


The **ImageNet** project is a large visual database designed for use in visual object recognition software research. As of 2016, over ten million URLs of images have been hand-annotated by ImageNet to indicate what objects are pictured. Since 2010, the ImageNet project runs an annual software contest, the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), where software programs compete to correctly classify and detect objects and scenes.

12

## How to achieve nonlinearity?

$$f : \mathbb{R}^m \longrightarrow \mathbb{R}^n, f_{\theta}(\mathbf{x}) \quad \theta \leftarrow \operatorname{argmin}_{\theta} \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(y, f_{\theta}(\mathbf{x}))$$



$$\mathbf{h}_1 = g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{h}_2 = g(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2)$$

$$f(\mathbf{x}) = g\left(\mathbf{W}_2(g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)) + \mathbf{b}_2\right)$$

$$g(z) = \max\{\mathbf{0}, z\} \quad \text{Rectifier Linear Unit (Function) (ReLU)}$$

13

## How to achieve nonlinearity?

$$f : \mathbb{R}^m \longrightarrow \mathbb{R}^n, f_{\theta}(\mathbf{x}) \quad \theta \leftarrow \operatorname{argmin}_{\theta} \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(y, f_{\theta}(\mathbf{x}))$$

Optimization requires derivatives

14

## Gradient Descent

To solve

$$\theta \leftarrow \operatorname{argmin}_{\theta} \mathbf{J}(\theta)$$

GD is the algorithm



```

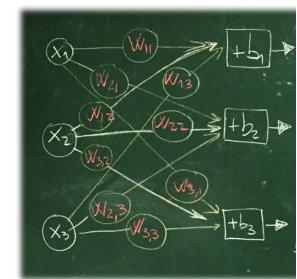
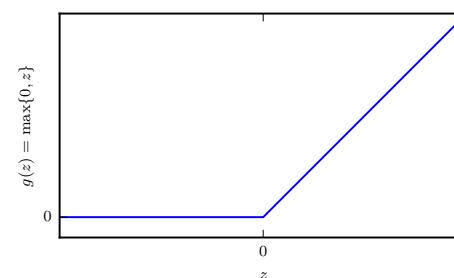
 $\theta \leftarrow \theta_0$ 
repeat
 $\theta \leftarrow \theta - \gamma * \frac{\partial J}{\partial \theta}$ 
until no more improvement can be reached
return  $\theta$ 

```

15

## Rectifier Linear Activation (ReLU)

$$\mathbf{h}_1 = g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$



*In practice one can safely disregard the non-differentiability of the hidden unit activation functions. It is very unlikely that the underlying value was 0*

16

## Old fashion activation functions

- Instead of ReLU, functions such as sigmoid or hyperbolic tangent (similar) were used.
- Nowadays, are not used because they have almost zero slopes for most of the values. They saturate and are not useful.

17

## General Framework

19

## Universal Approximator Theorem

- One hidden layer is enough to *represent* (not *learn*) an approximation of any function to an arbitrary degree of accuracy
- So why deeper?
  - Shallow net may need (exponentially) more width
  - Shallow net may overfit more

18

## General Framework

$$D = \{(\mathbf{x}, \mathbf{y}) : \mathbf{x} \rightsquigarrow \mathbf{y}\}$$

We learn a distribution

$$\Pr(\mathbf{y}|\mathbf{x}, \theta)$$

That involves the function

$$f : \mathbb{R}^m \longrightarrow \mathbb{R}^n, f_\theta(\mathbf{x})$$

This requires maximum likelihood

$$\operatorname{argmax}_\theta \prod_D \Pr(\mathbf{y}|\mathbf{x}, \theta)$$

20

# General Framework

This means to maximize -log

$$\operatorname{argmin}_{\theta} - \log \left( \prod_D \Pr(y|x, \theta) \right) = \operatorname{argmin}_{\theta} - \sum_D \log \Pr(y|x, \theta)$$

In other words, using as loss functions

$$\text{Cross entropy} \quad - \log \Pr(y|x, \theta)$$

21

# Learning Tasks

- **Feedforward** (feedforward networks, feedforward neural networks, or multilayer perceptrons (MLPs)) fully connected
- **Convolutional** networks (convolutional neural networks, CNNs), are a specialized kind of neural network for processing data that has a known grid-like topology. Examples include time-series data, which can be thought of as a 1-D grid taking samples at regular time intervals, and image data, which can be thought of as a 2-D grid of pixels
- **Recurrent** neural networks (RNN) are a family of neural networks for processing sequential data. They extend feedforward neural networks to include feedback connections.

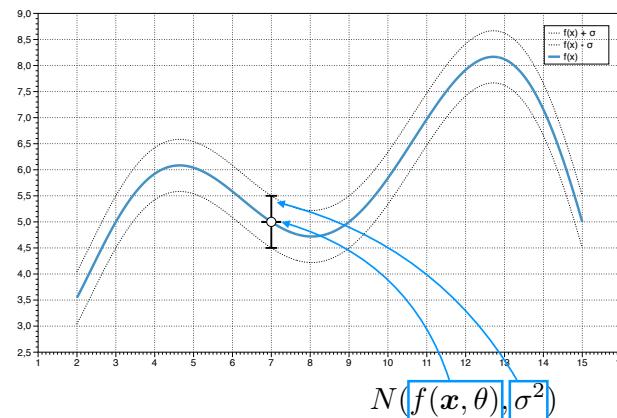
22

# Learning Tasks

- General setting for some learning tasks:
  - (Multi)-regression
  - classification
- Regularization

23

# Regression



$\sigma$  is constant and it is estimated using the sample deviation

24

# Regression

$$y \in \mathbb{R}; D = \{(\mathbf{x}, y), \dots\}$$

$$\Pr(y|\mathbf{x}, \theta) = N(y; f(\mathbf{x}, \theta), \sigma^2)$$

Probability distribution involving  $f$

$$N(y; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^2\right)$$

25

# Regression

$$y \in \mathbb{R}; D = \{(\mathbf{x}, y), \dots\}$$

$$\Pr(y|\mathbf{x}, \theta) = N(y; f(\mathbf{x}, \theta), \sigma^2)$$

Probability distribution involving  $f$

$$N(y; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^2\right)$$

$$-\log \Pr(y|\mathbf{x}, \theta) = \cancel{\log \sigma} + \frac{1}{2} \cancel{\log(2\pi)} + \frac{1}{2\sigma^2} (y - f(\mathbf{x}, \theta))^2$$

constants

$$\arg \min_{\theta} -\log \prod_D \Pr(y|\mathbf{x}, \theta) = \arg \min_{\theta} \sum_D (y - f(\mathbf{x}, \theta))^2$$

Loss function

26

# Binary classification

$$y \in \{0, 1\}; D = \{(\mathbf{x}, y), \dots\}$$

We search for a probability distribution proportional to

$$\Pr(y|\mathbf{x}, \theta) \propto \exp(yf(\mathbf{x}, \theta))$$

This is NOT a probability distribution

Thus,

$$\Pr(y|\mathbf{x}, \theta) = \frac{\exp(yf(\mathbf{x}, \theta))}{\sum_{y'} \exp(y'f(\mathbf{x}, \theta))} = \frac{\exp(yf(\mathbf{x}, \theta))}{1 + \exp(f(\mathbf{x}, \theta))}$$

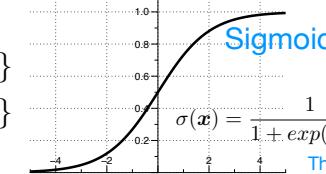
Softmax

27

# Binary classification

$$y \in \{0, 1\}$$

$$D = \{(\mathbf{x}, y), \dots\}$$

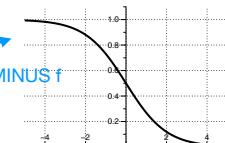


$$\Pr(y=1|\mathbf{x}, \theta) = \frac{\exp(f(\mathbf{x}, \theta))}{1 + \exp(f(\mathbf{x}, \theta))} = \frac{1}{\exp(f(\mathbf{x}, \theta)) + 1} = \sigma(f(\mathbf{x}, \theta))$$

The probability increase with  $f$  value

$$\Pr(y=0|\mathbf{x}, \theta) = \frac{1}{1 + \exp(f(\mathbf{x}, \theta))} = \sigma(-f(\mathbf{x}, \theta))$$

The probability increase with MINUS  $f$  value



$$\Pr(y|\mathbf{x}, \theta) = \sigma((2y - 1)f(\mathbf{x}, \theta))$$

28

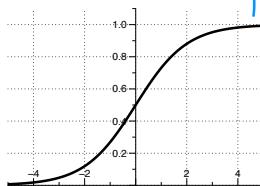
# Binary classification

$$y \in \{0, 1\}; D = \{(\mathbf{x}, y), \dots\}$$

$$\Pr(y|\mathbf{x}, \theta) = \sigma((2y - 1)f(\mathbf{x}, \theta))$$

This probability  $> 0.5$  when are equal the signs of  $(2y - 1)$  and  $f$ : predictions are right whenever

- Right •  $y=1 \Rightarrow (2y-1)=+1$  and  $f$  are +
- $y=0 \Rightarrow (2y-1)=-1$  and  $f$  are -



Erroneous whenever

Error  $(2y - 1)f(\mathbf{x}, \theta) < 0$

29

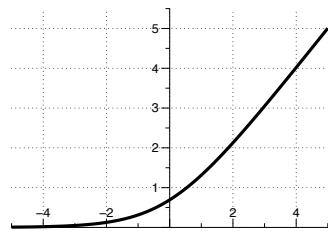
# Binary classification

$$y \in \{0, 1\}; D = \{(\mathbf{x}, y), \dots\}$$

$$\Pr(y|\mathbf{x}, \theta) = \sigma((2y - 1)f(\mathbf{x}, \theta))$$

The loss function is

$$\text{softplus}((1 - 2y)f(\mathbf{x}, \theta))$$



The classification error appears whenever the argument of softplus is positive

Error  $(1 - 2y)f(\mathbf{x}, \theta) > 0$

Error  $\rightarrow$  a lot of slope

**NO** error  $\rightarrow$  slope almost zero (saturates)

31

# Binary classification

$$y \in \{0, 1\}; D = \{(\mathbf{x}, y), \dots\}$$

$$\Pr(y|\mathbf{x}, \theta) = \sigma((2y - 1)f(\mathbf{x}, \theta))$$

The loss function is

$$-\log \Pr(y|\mathbf{x}, \theta) = -\log \sigma((2y - 1)f(\mathbf{x}, \theta))$$

Given that

$$-\log(\sigma(z)) = -\log \frac{1}{1 + e^{-z}} = -\log(1) + \log(1 + e^{-z}) = \text{softplus}(-z)$$

**Softplus**

Thus, we have

$$-\log \Pr(y|\mathbf{x}, \theta) = \text{softplus}((1 - 2y)f(\mathbf{x}, \theta))$$

Loss function

30

# Multiclass classification

$$y \in \{1, 2, \dots, n\}; D = \{(\mathbf{x}, y), \dots\}$$

$$f : \mathbb{R}^m \longrightarrow \mathbb{R}^n, f(\mathbf{x}, \theta) = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n)$$

$$\Pr(y = i|\mathbf{x}, \theta) \propto \exp(\mathbf{z}_i)$$

"Distribution" involving f

$$\Pr(y = i|\mathbf{x}, \theta) = \frac{\exp(\mathbf{z}_i)}{\sum_j \exp(\mathbf{z}_j)} = \text{softmax}(\mathbf{z})_i$$

**Softmax**

The goal is to minimize

$$-\log \text{softmax}(\mathbf{z})_i = \log \sum_j \exp(\mathbf{z}_j) - \mathbf{z}_i$$

Loss function

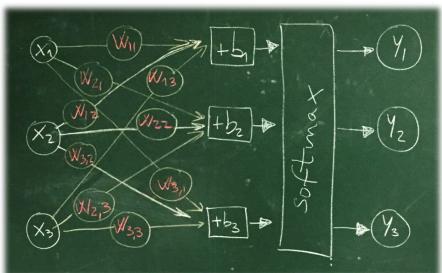
32

# Multiclass classification

$$y \in \{1, 2, \dots, n\}; D = \{(\mathbf{x}, y), \dots\}$$

$$f : \mathbb{R}^m \longrightarrow \mathbb{R}^n, f(\mathbf{x}, \theta) = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n)$$

$$\begin{aligned} \Pr(y = i|\mathbf{x}, \theta) &= \frac{\exp(\mathbf{z}_i)}{\sum_j \exp(\mathbf{z}_j)} \\ &= \text{softmax}(\mathbf{z})_i \end{aligned}$$



33

# Gradient

34

# Gradient

The **gradient** of a differentiable function

$$f : \mathbb{R}^n \longrightarrow \mathbb{R}$$

of several variables  $\mathbf{x} = (x_1, \dots, x_n)$

is the vector whose value at a point  $\mathbf{a}$  is the **vector** whose components are the **partial derivatives** of  $f$  at  $\mathbf{a}$

$$\nabla_{\mathbf{x}} f(\mathbf{a}) = \left( \frac{\partial f}{\partial x_1}(\mathbf{a}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{a}) \right)$$

35

# Gradient Descent

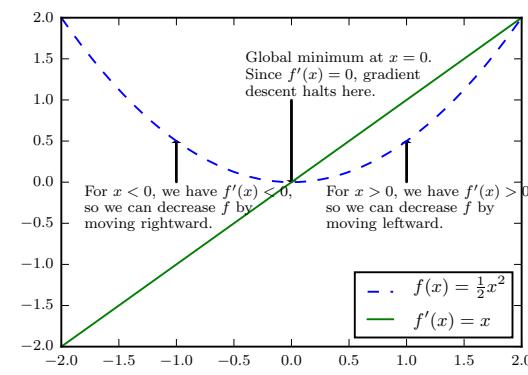


Figure 4.1

36

# Gradient Descent

The Taylor polynomial of order 1 states that for values  $\mathbf{f}(\mathbf{x})$  when  $\mathbf{x}$  is close to  $\mathbf{a}$  can be approximated using the derivative as follows

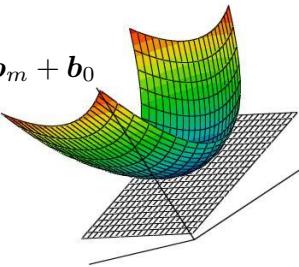
$$f(\mathbf{x}) = f(\mathbf{a}) + (\mathbf{x} - \mathbf{a})^T \nabla_x f(\mathbf{a})$$

That is to say,

$$f(\mathbf{x}) = f(\mathbf{a}) + \mathbf{x}_1 \mathbf{b}_1 + \dots + \mathbf{x}_m \mathbf{b}_m + b_0$$

$$(\mathbf{b}_1, \dots, \mathbf{b}_m) = \nabla_x f(\mathbf{a})$$

$$\mathbf{b}_0 = -\mathbf{a}^T \nabla_x f(\mathbf{a})$$



37

# Gradient Descent

$$f(\mathbf{x}) = f(\mathbf{a}) + (\mathbf{x} - \mathbf{a})^T \nabla_x f(\mathbf{a})$$

If we state that  $\mathbf{x}$  is  $\mathbf{a}$  modified by a small increase ( $\gamma$ ) in the address of a unit vector  $\mathbf{u}$  ( $\|\mathbf{u}\| = 1$ )

$$\mathbf{x} = \mathbf{a} + \gamma \mathbf{u}$$

We have that

$$f(\mathbf{a} + \gamma \mathbf{u}) = f(\mathbf{a}) + \gamma \mathbf{u}^T \nabla_x f(\mathbf{a})$$

38

# Gradient Descent

To find a value  $\mathbf{x}$  such that  $\mathbf{f}(\mathbf{x}) < \mathbf{f}(\mathbf{a})$ , we must find a direction ( $\mathbf{u}$ ) that

$$\begin{aligned} \operatorname{argmin}_{\mathbf{u}} \mathbf{u}^T \nabla_x f(\mathbf{a}) \\ = \operatorname{argmin}_{\mathbf{u}} \|\mathbf{u}\| \|\nabla_x f(\mathbf{a})\| \cos(\mathbf{u}, \nabla_x f(\mathbf{a})) \\ = \operatorname{argmin}_{\mathbf{u}} \cos(\mathbf{u}, \nabla_x f(\mathbf{a})) \end{aligned}$$

Therefore  $\mathbf{u}$  must be parallel to the opposite of the gradient ( $\cos = -1$ )

$$\mathbf{u} \parallel -\nabla_x f(\mathbf{a})$$

39

# Gradient Descent

Therefore (if  $\gamma$  is small)

$$\begin{aligned} f(\mathbf{a} - \gamma \nabla_x f(\mathbf{a})) &= f(\mathbf{a}) - \gamma \nabla_x f(\mathbf{a})^T \nabla_x f(\mathbf{a}) \\ &= f(\mathbf{a}) - \gamma \|\nabla_x f(\mathbf{a})\|^2 \\ &< f(\mathbf{a}) \end{aligned}$$

After the assignment

$$\mathbf{a} \leftarrow \mathbf{a} - \gamma \nabla_x f(\mathbf{a})$$

We have a point with lower f-value than  $f(\mathbf{a})$ , whenever  $\gamma$  is small

40

# Gradient Descent

Tying up loose ends. The goal was to solve

$$\theta \leftarrow \operatorname{argmin}_{\theta} \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(y, f_{\theta}(x))$$

We look for the  $\theta$  value that minimizes the mean of the loss function over the dataset D

$$\begin{aligned} \mathbf{J}(\theta) &= \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(y, f_{\theta}(x)) \\ &= \mathbb{E}_{(x,y) \sim p(D)} \text{loss}(y, f_{\theta}(x)) \end{aligned}$$

41

# Mini Batch and SGD

GD has the problem to compute

$$\frac{\partial \mathbf{J}(\theta)}{\partial \theta} = \frac{1}{|D|} \sum_{(x,y) \in D} \frac{\partial}{\partial \theta} \text{loss}(y, f_{\theta}(x))$$

Thus, to make GD useful, we estimate this average with the mean of a subset  $D'$  of D

$$\frac{\partial \mathbf{J}(\theta)}{\partial \theta} \cong \frac{1}{|D'|} \sum_{(x,y) \in D'} \frac{\partial}{\partial \theta} \text{loss}(y, f_{\theta}(x))$$

This is called the Stochastic Gradient Descent with mini batch. Whenever  $D'$  has only one element, we have SGD

43

# Gradient Descent

To solve

$$\theta \leftarrow \operatorname{argmin}_{\theta} \mathbf{J}(\theta)$$

GD is the algorithm

```
θ ← θ₀
repeat
    θ ← θ - γ * ∂J / ∂θ
until no more improvement can be reached
return θ
```

42

# Final remarks

44

# Back-Propagation

- Back-propagation is “just the chain rule” of calculus

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}. \quad (6.44)$$

$$\nabla_{\mathbf{x}} z = \left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z, \quad (6.46)$$

- But it's a particular implementation of the chain rule
  - Uses dynamic programming (table filling)
  - Avoids recomputing repeated subexpressions
  - Speed vs memory tradeoff

45

# Regularization

“Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.”

46

# Regularization

- Parameter norm penalization
- Early stopping
- Dropout

47

# Norm Penalties

Squared L2: Encourages small weights, equivalent to MAP Bayesian estimation with Gaussian prior

$$\begin{aligned} \mathbf{J}(\theta) &+ \frac{\nu}{2} \|\theta\|^2 \\ &= \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(\mathbf{y}, f(\mathbf{x}, \theta)) + \frac{\nu}{2} \|\theta\|^2 \end{aligned}$$

Also known as ridge regression or Tikhonov regularization

Regression and classification

48

# Norm Penalties

L1: Encourages sparsity, equivalent to MAP Bayesian estimation with Laplace prior. It has been used extensively as a feature selection mechanism. LASSO.

$$\begin{aligned} \mathbf{J}(\theta) + \nu \|\theta\|_1 \\ = \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(\mathbf{y}, f(\mathbf{x}, \theta)) + \nu \|\theta\|_1 \end{aligned}$$

(sub-gradient)

$$\nabla_{\theta} \mathbf{J}(\theta) + \nu \text{ sign}(\theta)$$

49

# Early stopping

Early stopping: terminate when validation set performance is better

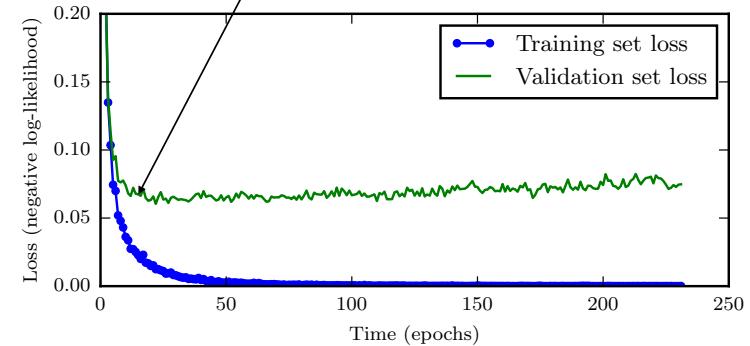


Figure 7.3

50

# Early stopping

- It is probably the most commonly used form of regularization in deep learning.
- Requires a validation set
- After initial training with early stopping has completed, a second training with all data (train + validation) is performed with new starting values

51

# Dropout

- Dropout trains an *ensemble*: *all sub-nets* which can be formed by removing nodes that are not output
- There is an exponential number of subnets. Then dropout uses a sampling
- Typically, inputs are included with probability 0.8 and hidden nodes with 0.5

52

# Dropout

Figure 7.6

