

Introduction to Machine Learning

Introduction to Machine Learning

Amo G. Tong

1

Lecture 8 Supervised Learning

- Artificial Neural Networks
- Backpropagation Algorithm

- Some materials are courtesy of Vibhava Gogate, Eric Xing and Tom Mitchell.
- All pictures belong to their creators.

Introduction to Machine Learning

Amo G. Tong

2

Supervised Learning

- $\langle x, y \rangle$
- Estimate $\Pr[y|x]$ by Bayesian Theory: $\Pr[y|x] = \frac{\Pr[x|y] \Pr[y]}{\Pr[x]}$
- Naïve Bayes.

- Estimate $\Pr[y|x]$ directed by assuming a certain form
- Logistic regression.

- Assume the form of $y = f(x)$
- Error-driven.

$\Pr[x|y]$ and $\Pr[y]$ are hard to compute.
No good form for $\Pr[y|x]$
No prior knowledge on $f(x)$.

You can try Neural Networks.

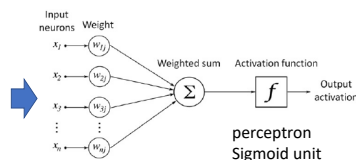
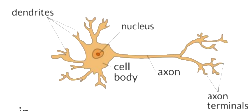
Introduction to Machine Learning

Amo G. Tong

3

Neural Networks

- Input x , output y
- Given x , how to compute y ?
- x and y are generally real vectors.



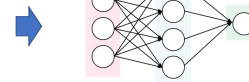
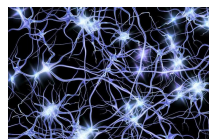
Introduction to Machine Learning

Amo G. Tong

4

Neural Networks

- Input x , output y
- Given x , how to compute y ?
- x and y are generally real vectors.



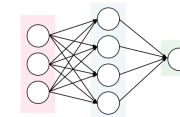
Introduction to Machine Learning

Amo G. Tong

5

Neural Networks

- Input x , output y
- Given x , how to compute y ?
- x and y are generally real vectors.



- Many neuron-like units (perceptron, sigmoid)
- Weighted interconnections between units.
- Parallel distributed process.

Introduction to Machine Learning

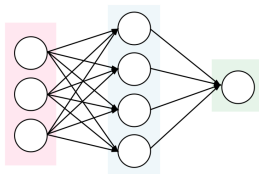
Amo G. Tong

6

Neural Networks

- Input x , output y
- Given x , how to compute y ?
- x and y are generally real vectors.

- Design a neural network:
- What is the function within each unit?
- How are the units connected?



Introduction to Machine Learning

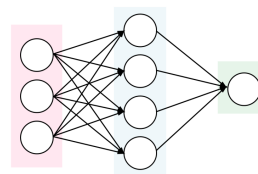
Amo G. Tong

7

Neural Networks

- Input x , output y
- Given x , how to compute y ?
- x and y are generally real vectors.

- Units:
- Input Unit
- Processing Units
 - Receive input from other units
 - Output result to other units
- Edges
 - weights



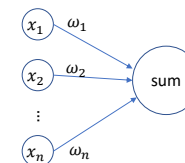
Introduction to Machine Learning

Amo G. Tong

8

Examples

- Linear function $y = \omega_1 x_1 + \dots + \omega_n x_n$
- Input (x_1, \dots, x_n) output y .



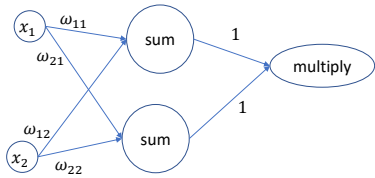
Introduction to Machine Learning

Amo G. Tong

9

Examples

- Quadratic function $y = (\omega_{11}x_1 + \omega_{12}x_2)(\omega_{21}x_1 + \omega_{22}x_2)$
- Input (x_1, \dots, x_n) output y .

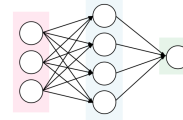


Train a Neural Network

- Define the error E
- Update the weight ω of each edge to minimize E
- $\omega \leftarrow \omega - \eta \frac{\partial E}{\partial \omega}$ (delta rule)

Calculate $\frac{\partial E}{\partial \omega}$

- Batch Mode
 - E : the total error among training data.
 - Consider all the training data each time.
- Incremental Mode
 - E : the error for one training instance.
 - Consider training instances one by one.

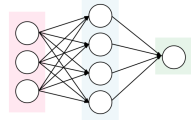


Train a Neural Network

- Define the error E
- Update the weight ω of each edge to minimize E
- $\omega \leftarrow \omega - \eta \frac{\partial E}{\partial \omega}$ (delta rule)

Calculate $\frac{\partial E}{\partial \omega}$

- Incremental Mode
 - E : the error for one training instance.
 - Consider training instances one by one.

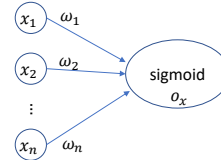


Neural Networks

- One layer of sigmoid with one output
 - Two layers of single functions
 - Two layers of multiple sigmoid units with one output.
 - Two layers of multiple sigmoid units with multiple outputs.
- Goal: how to find the parameters that can minimize the error.

Neural Networks

- One layer of sigmoid with one output

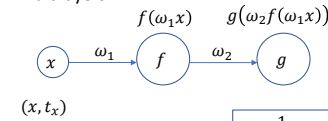


$$o(x_1, \dots, x_n) = \frac{1}{1 + e^{-\omega x}}$$

$$\frac{\partial E_x}{\partial \omega_i} = -(t_x - o_x) \frac{\partial o_x}{\partial \omega_i} = -(t_x - o_x) \cdot x_i \cdot o_x(1 - o_x)$$

Neural Networks

- Multilayers



$$E = \frac{1}{2} |t_x - g|^2$$

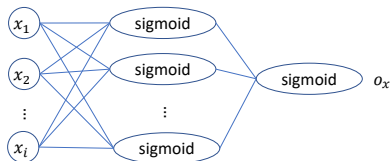
$$\frac{\partial E}{\partial \omega_2} = -(t_x - g) \frac{\partial g}{\partial \omega_2} = \frac{\partial g}{\partial \omega_2} \frac{\partial f}{\partial \omega_1} \frac{\partial \omega_1}{\partial \omega_2} = \frac{\partial g}{\partial \omega_2} \omega_2 \frac{\partial f}{\partial \omega_1} \frac{\partial \omega_1}{\partial \omega_2} = \frac{\partial g}{\partial \omega_2} \omega_2 \frac{\partial f}{\partial \omega_1} x$$

$$\frac{\partial E}{\partial \omega_1} = -(t_x - g) \frac{\partial g}{\partial \omega_1} = \frac{\partial g}{\partial \omega_1} \frac{\partial f}{\partial \omega_1} \frac{\partial \omega_1}{\partial \omega_1} = \frac{\partial g}{\partial \omega_1} \omega_2 \frac{\partial f}{\partial \omega_1} \frac{\partial \omega_1}{\partial \omega_1} = \frac{\partial g}{\partial \omega_1} \omega_2 \frac{\partial f}{\partial \omega_1} x$$

- Update ω_1 or ω_2 first?

Neural Networks

- Two layers of multiple sigmoid units with one output.

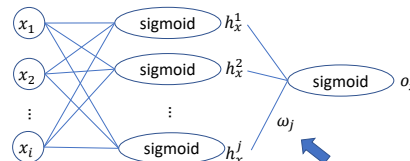


For each x in D
 $E_x = \frac{1}{2} (t_x - o_x)^2$

today's weather x $\xrightarrow{f(x, \omega)}$ tomorrow's temp (normalized) $t_x = f(x)$
 ω : weight on each edge

Neural Networks

- Two layers of multiple sigmoid units with one output.

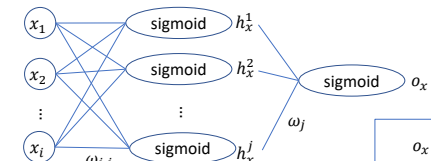


For each x in D
 $E_x = \frac{1}{2} (t_x - o_x)^2$

$$\frac{\partial E_x}{\partial \omega_j} = -(t_x - o_x) \frac{\partial o_x}{\partial \omega_j} = -(t_x - o_x) \cdot h_x^j \cdot o_x(1 - o_x)$$

Neural Networks

- Two layers of multiple sigmoid units with one output.



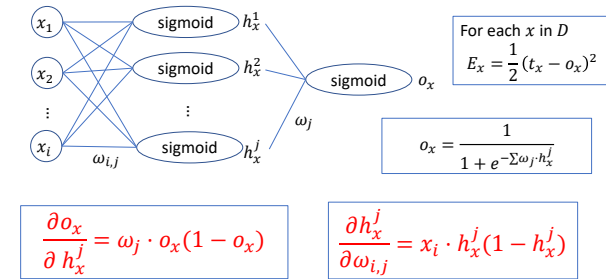
For each x in D
 $E_x = \frac{1}{2} (t_x - o_x)^2$

$$o_x = \frac{1}{1 + e^{-\sum \omega_j h_x^j}}$$

$$\frac{\partial E_x}{\partial \omega_{i,j}} = -(t_x - o_x) \frac{\partial o_x}{\partial \omega_{i,j}} = -(t_x - o_x) \frac{\partial o_x}{\partial h_x^j} \frac{\partial h_x^j}{\partial \omega_{i,j}}$$

Neural Networks

- Two layers of multiple sigmoid units with one output.



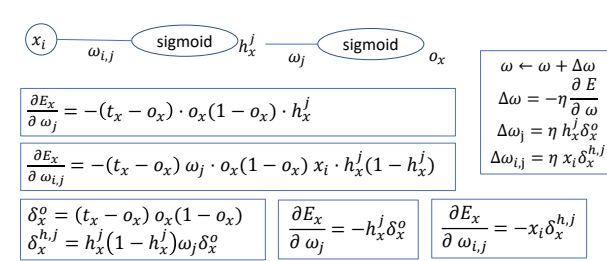
Introduction to Machine Learning

Amo G. Tong

19

Neural Networks

- Two layers of multiple sigmoid units with one output.



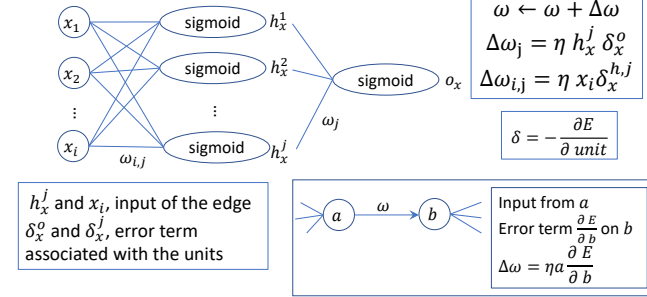
Introduction to Machine Learning

Amo G. Tong

20

Neural Networks

- Two layers of multiple sigmoid units with one output.



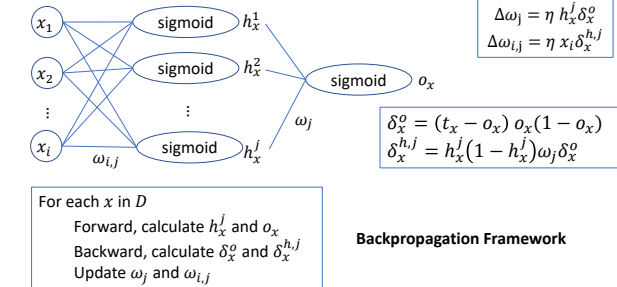
Introduction to Machine Learning

Amo G. Tong

21

Neural Networks

- Two layers of multiple sigmoid units with one output.



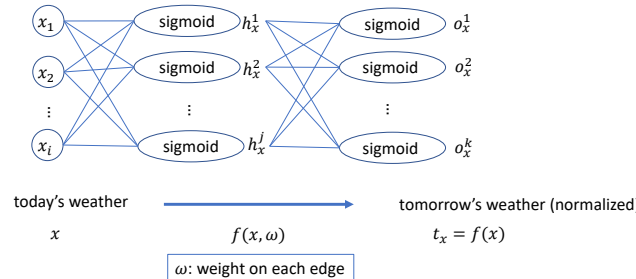
Introduction to Machine Learning

Amo G. Tong

22

Neural Networks

- Two layers of multiple sigmoid units with multiple outputs.



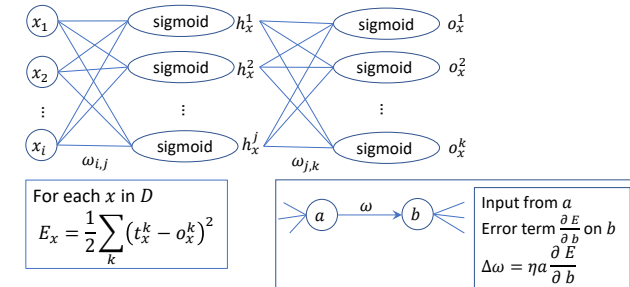
Introduction to Machine Learning

Amo G. Tong

23

Neural Networks

- Two layers of multiple sigmoid units with multiple outputs.



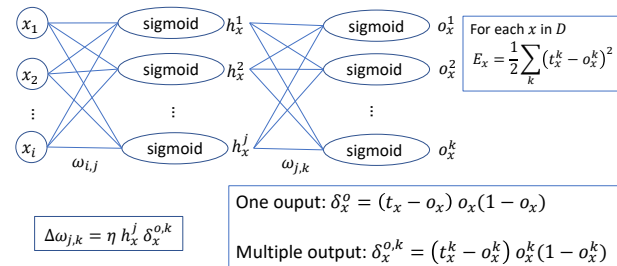
Introduction to Machine Learning

Amo G. Tong

24

Neural Networks

- Two layers of multiple sigmoid units with multiple outputs.



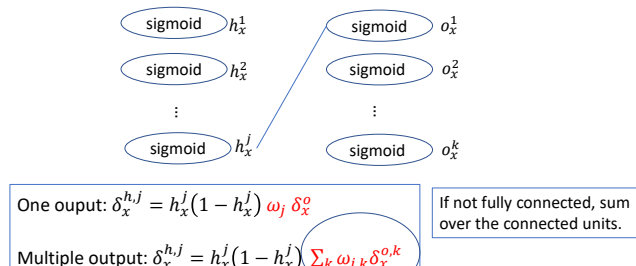
Introduction to Machine Learning

Amo G. Tong

25

Neural Networks

- Two layers of multiple sigmoid units with multiple outputs.



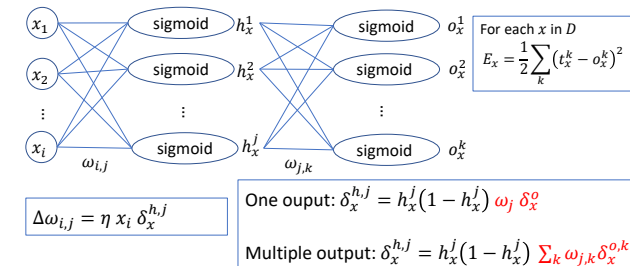
Introduction to Machine Learning

Amo G. Tong

26

Neural Networks

- Two layers of multiple sigmoid units with multiple outputs.



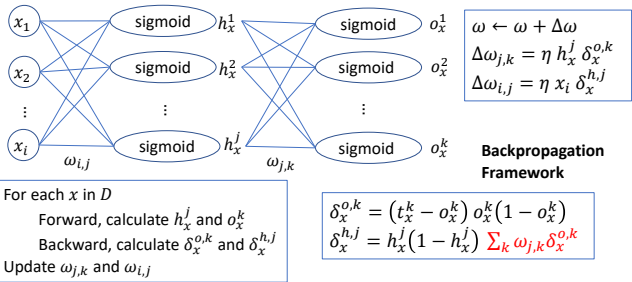
Introduction to Machine Learning

Amo G. Tong

27

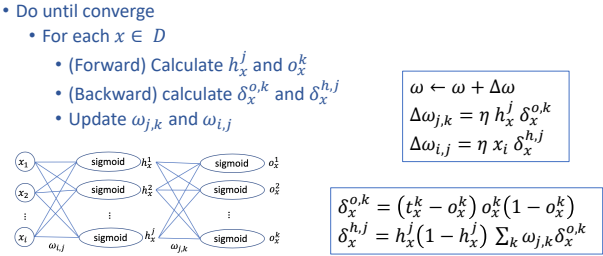
Neural Networks

- Two layers of multiple sigmoid units with multiple outputs.



Neural Networks

- Two layers of multiple sigmoid units with multiple outputs.
- Initialize the weights with random values.
- Do until converge



Neural Networks

- Train a neural network with:
- One sigmoid units.
- Two layers of multiple sigmoid units with one output.
- Two layers of multiple sigmoid units with multiple outputs.

- Backpropagation framework: any units, any acyclic graph.
- Update the weight from right to left.

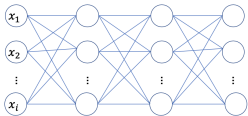


Neural Networks

- Backpropagation framework

- Good news: using the network is fast.
- Bad news: training the network is slow
- More bad news: it may converge to local minima.
- More good news: it performs well in practice.

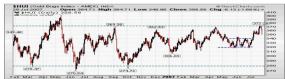
- To avoid local minima
- Add a momentum
 - $\Delta\omega_n^* = \Delta\omega_n + \alpha\Delta\omega_{n-1}^*$
- Train with different initializations.



Neural Networks

- Representational Power of Neural Networks.

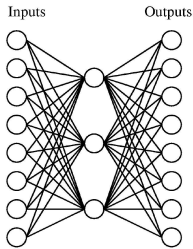
- Boolean functions: every Boolean function can be represented exactly by some network with **two layers** of units.
- Continuous functions: every bounded continuous function can be approximated with arbitrarily small error by a network with **two layers of units**. (sigmoid + linear will do)
- Arbitrary function: any function can be approximated to arbitrary accuracy by a network with **three layers of units**. (sigmoid + sigmoid+ linear will do)



An example (Mitchell)

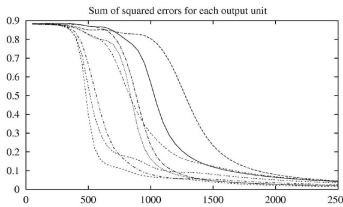
- Learn an identity function with eight training examples.

Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001



An example (Mitchell)

- Learn an identity function with eight training examples.

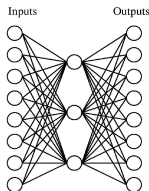


Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001

An example (Mitchell)

- Learn an identity function with eight training examples.

Input	Hidden Values	Output
10000000	→ .89 .04 .08	→ 10000000
01000000	→ .01 .11 .88	→ 01000000
00100000	→ .01 .97 .27	→ 00100000
00010000	→ .99 .97 .71	→ 00010000
00001000	→ .03 .05 .02	→ 00001000
00000100	→ .22 .99 .99	→ 00000100
00000010	→ .80 .01 .98	→ 00000010
00000001	→ .60 .94 .01	→ 00000001

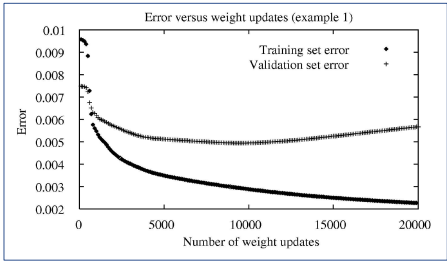


Overfitting

- Overfitting is everywhere...
- Overfitting may happen, when
- The learned model is too complicated
- Fitting data too well when data has noise
- When training set is small
- Avoid large parameters.

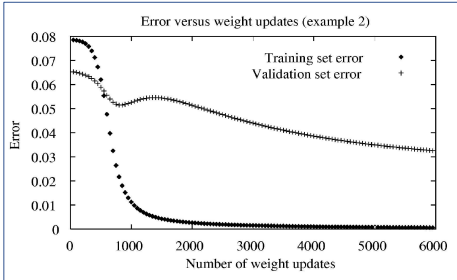
Overfitting

- Overfitting is everywhere...



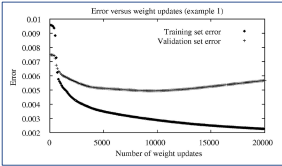
Overfitting

- Overfitting is everywhere...



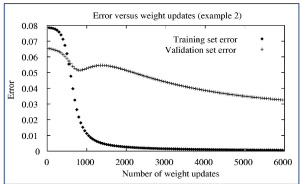
Overfitting

- Overfitting is everywhere...
- Weight Decay:
 - Decrease each weight by some small factor during each iteration?
- Penalize large weights:
 - $Error = Error + \gamma \sum \omega_i^2$
- Using Validation data.



Overfitting

- Overfitting is everywhere...
- Weight Decay:
 - Decrease each weight by some small factor during each iteration?
- Penalize large weights:
 - $Error = Error + \gamma \sum \omega_i^2$
- Using Validation data.



Summary

- What is an artificial neural network?
- Process the input by the **connected units**.
- How to training a neural network?
- Error-driven.
- Backpropagation algorithm.
 - Two layers of sigmoid units.