# Slide 1

## ECE408 Lecture 12

## Feed-Forward Networks and Gradient-Based Training

ECE408 / CS483 / CSE 408

Spring 2020

(by Carl Pearson)

1

# Slide 2

## Objective

- To learn the basic approach to feedforward neural networks:
  - neural model
  - common functions
  - training through gradient descent

2

2

# Slide 3

## Let's Look at Classification

In a **classification problem**, we model

- a function mapping an input vector to a set of $C$ categories: $F : \mathbb{R}^N \rightarrow \{1, \ldots, C\}$,
- where the function $F$ **is unknown**.

We **approximate $F$ using a set of functions $f$**

- parametrized by a (large) set of weights, $\Theta$
- that map from a vector of $N$ real values* to an integer value representing a category:
- for category $i$, **prob$(i) = f(x, \Theta)$**

*floating-point values

3

# Slide 4
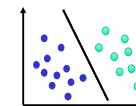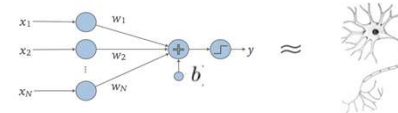
## Perceptron is a Simple Example

- Example: a **perceptron**

$y = \text{sign}(W \bullet x + b)$    $\Theta = \{W, b\}$

The perceptron    The neuron



$\approx$

- Dot product:    $y = W \bullet x$    output    input
- Scalar addition:    $+ b$    bias    weight

4

## One Perceptron is not Enough

**Some functions are non-linear.**

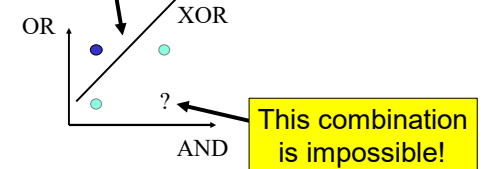**What can we do?**

- ○ FALSE
- ● TRUE



XOR

? = AND + OR

5

---

## Multiple Layers Solve More Problems

**What if input dimensions are AND and OR?**

Now we can divide with one line.

- ○ FALSE
- ● TRUE
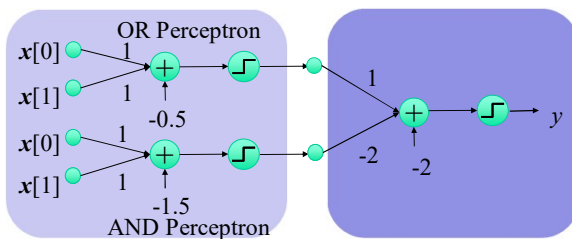


OR

XOR

AND

? ← This combination is impossible!

6

---

| A | B | OR | AND | XOR |
|---|---|----|-----|-----|
| 0 | 0 | -1 | -1 | -1 |
| 0 | 1 | 1 | -1 | 1 |
| 1 | 0 | 1 | -1 | 1 |
| 1 | 1 | 1 | 1 | -1 |

$AND = sign(x[0] + x[1] - 1.5)$

$OR = sign(x[0] + x[1] - 0.5)$

$XOR = sign(2 * OR - AND - 2)$



OR Perceptron

$x[0]$
$x[1]$

1
1
-0.5

$x[0]$
$x[1]$

1
1
-1.5

AND Perceptron

1
-2
-2

y

7

---

## Generalize to Fully-Connected Layer



$x[0]$
$x[1]$
$x[2]$

y

Linear Classifier:
Input vector $x \times$ weight vector $w$ to produce scalar output $y$

$x[0]$
$x[1]$
$x[2]$

$y[0]$
$y[1]$
$y[2]$
$y[3]$
$y[4]$

Fully-connected:
Input vector $x \times$ weight matrix $w$ to produce vector output $y$

8

## Multilayer Terminology

$x[0]$ $x[1]$ $x[2]$ $x[3]$

Weight matrices: Entry $i,j$ is weight between $i^{th}$ input and $j^{th}$ output

Input Layer

$W_1[i, j]$ is [4x4]
$b_1[j]$ is [4x1]

Hidden Layer(s)

"nodes"

$W_2[i, j]$ is [4x3]
$b_2[j]$ is [3x1]

Output Layer $z[0]$ $z[1]$ $z[2]$

Softmax

$k[0]$ $k[1]$ $k[2]$ Probability $k[i]$ that input is class i

9

---

## Example: Digit Recognition

Let's consider an example.

- **handwritten digit recognition**:
- given a $28 \times 28$ **grayscale image**,
- produce a **number from 0 to 9**.

Input dataset

- **60,000** images
- Each labeled by a human with correct answer.

10

---

## How Do We Determine the Weights?

**First layer** of perceptrons

- **784** ($28^2$) inputs, **1024** outputs, **fully connected**
- **[1024 × 784]** weight matrix **W**
- **[1024 x 1]** bias vector **b**

**Use labeled training data to pick weights.**

Idea:

- given enough labeled input data,
- we can **approximate the input-output function**.

11

---

## Forward and Backward Propagation

Forward (**inference**):

- given input $x$ (for example, an image),
- **use parameters $\Theta$** (**W** and $b$ for each layer)
- **to compute probabilities $k[i]$** (ex: for each digit i).

Backward (**training**):

- given input $x$, parameters $\Theta$, and outputs $k[i]$,
- **compute error $E$** based on target label $t$,
- then **adjust $\Theta$** proportional to $E$ to reduce error.

12

# Neural Functions Impact Training

Recall perceptron function: **y = sign (W·x + b)**

**To propagate error** backwards,
- **use chain rule** from calculus.
- **Smooth functions are useful.**

Sign is not a smooth function.

13

# One Choice: Sigmoid/Logistic Function

Until about 2017,
- **sigmoid / logistic function** most popular
$$f(x) = \frac{1}{1+e^{-x}} \quad (\text{f: } \mathbb{R} \to (0,1))$$
for replacing sign.
- Once we have *f(x)*, finding *df/dx* is easy:
$$\frac{df(x)}{dx} = \frac{e^{-x}}{(1+e^{-x})^2} = f(x)\frac{e^{-x}}{(1+e^{-x})} = f(x)(1-f(x))$$

(Our example used this function.)

14

# Today's Choice: ReLU

In 2017, most common choice became
- **rectified linear unit / ReLU / ramp function**
$$f(x) = \max(0, x) \quad (\text{f: } \mathbb{R} \to \mathbb{R}^+)$$
which is much faster (no exponent required).
- A smooth approximation is **softplus/SmoothReLU**
$$f(x) = \ln(1 + e^x) \quad (\text{f: } \mathbb{R} \to \mathbb{R}^+)$$
which is the integral of the logistic function.
- Lots of variations exist. See Wikipedia for an overview and discussion of tradeoffs.

15

# Use Softmax to Produce Probabilities

**How can sigmoid / ReLU produce probabilities?**
They can't.
- Instead, given output vector **Z = (z[0], ..., z[C-1])\***,
- we produce a second vector **K = (k[0], ..., k[C-1])**
- using the **softmax function**
$$k[i] = \frac{e^{z[i]}}{\sum_{j=0}^{C-1} e^{z[j]}}$$

Notice that **the k[i] sum to 1.**

*Remember that we classify into one of C categories.

16

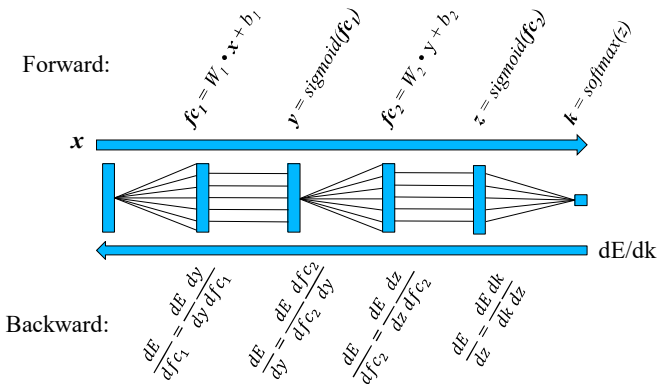## Softmax Derivatives Needed to Train

We also need the **derivatives of softmax**,

$$\frac{dk[i]}{dz[m]} = k[i](\delta_{i,m} - k[m]),$$

where $\delta_{i,m}$ is the Kronecker delta
(1 if i = m, and 0 otherwise).

17

## Forward and Backward Propagation



Forward:

$$fc_1 = W_1 \cdot x + b_1 \quad y = sigmoid(fc_1) \quad fc_2 = W_2 \cdot y + b_2 \quad z = sigmoid(fc_2) \quad k = softmax(z)$$

$x$

dE/dk

Backward:

$$\frac{dE}{dfc_1} = \frac{dE}{dy}\frac{dy}{dfc_1} \quad \frac{dE}{dy} = \frac{dE}{dfc_2}\frac{dfc_2}{dy} \quad \frac{dE}{dfc_2} = \frac{dE}{dz}\frac{dz}{dfc_2} \quad \frac{dE}{dz} = \frac{dE}{dk}\frac{dk}{dz}$$

18

## Choosing an Error Function

Many error functions are possible.

For example, **given label $T$** (digit $T$),

- $E = 1 - k[T]$,
- the **probability of not classifying as $t$**).

**Alternatively**, since our categories are numeric,
we can **penalize quadratically**:

$$E = \sum_{j=0}^{C-1} k[j](j - T)^2$$

Let's **go with the latter**.

19

## Stochastic Gradient Descent

**How do we calculate the weights?**

One common answer: **stochastic gradient descent**.

1. **Calculate**
   - **derivative** of sum **of error $E$**
   - **over all** training **inputs**
   - **for** all network parameters $\Theta$.
2. **Change $\Theta$ slightly** in the opposite direction
   (to decrease error).
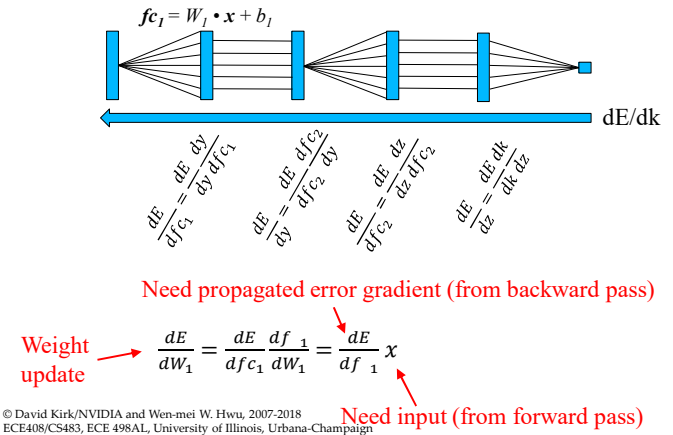3. **Repeat**.

20

## Stochastic Gradient Descent

More precisely,

1. **For every input $X$,**
2. evaluate network to **compute $k[i]$** (forward),
3. then **use $k[i]$ and label $T$** (target digit) **to compute error $E$**.
4. Backpropagate error derivative to **find derivatives for each parameter**.
5. **Adjust $\Theta$ to reduce total $E$: $\Theta_{i+1} = \Theta_i - \varepsilon\Delta\Theta$**

**(Update $\varepsilon$ uses most accurate minima estimation.)**

21

---

## Parameter Updates and Propagation



$$fc_1 = W_1 \cdot x + b_1$$

dE/dk

$$\frac{dE}{dfc_1} = \frac{dE}{dy}\frac{dy}{dfc_1}$$
$$\frac{dE}{dy} = \frac{dE}{dfc_2}\frac{dfc_2}{dy}$$
$$\frac{dE}{dfc_2} = \frac{dE}{dz}\frac{dz}{dfc_2}$$
$$\frac{dE}{dz} = \frac{dE}{dk}\frac{dk}{dz}$$

Need propagated error gradient (from backward pass)

Weight update

$$\frac{dE}{dW_1} = \frac{dE}{dfc_1}\frac{df_1}{dW_1} = \frac{dE}{df_1}x$$

Need input (from forward pass)

22

---

## Example: Gradient Update with One Layer

| | |
|---|---|
| $\Theta_{i+1} = \Theta_i - \varepsilon\Delta\Theta \quad W_{i+1} = W_i - \varepsilon\Delta W$ | Parameter Update |
| $y = W \cdot x + b$ | Network function |
| $\dfrac{dy}{dW} = x$ | Network weight gradient |
| $E = \frac{1}{2}(y - t)^2$ | Error function |
| $\dfrac{dE}{dy} = y - t = Wx + b - t$ | Error function gradient |
| $\Delta W = \dfrac{dE}{dW} = \dfrac{dE}{dy}\dfrac{dy}{dW}$ | Full weight update expression |
| $W_{i+1} = W_i - \varepsilon(Wx+b-t)x$ | Full weight update term |

23

---

## Fully-Connected Gradient Detail

$i^{th}$ entry in $fc_1$     $i^{th}$ row in $W_1$     $j^{th}$ entry in $x_1$



$$
\begin{matrix} fc_1[0] \\ fc_1[1] \\ fc_1[2] \\ \dots \end{matrix}
=
\begin{matrix} W_1[0,:] \\ W_1[1,:] \\ W_1[2,:] \\ \dots \end{matrix}
\quad\quad
\begin{matrix} x_1[0] \\ x_1[1] \\ x_1[2] \\ x_1[3] \\ \dots \end{matrix}
$$

$$fc_1 = \sum_j W_1[i,j]x_1[j]$$

Computed from previous layer

$$\frac{dE}{dW_1[i,j]} = \frac{dE}{dfc_1[i]}\frac{dfc_1[i]}{dW_1[i,j]} = \frac{dE}{dfc_1[i]}x_1[j]$$

Need input to this layer

24

## Batched Stochastic Gradient Descent

- A training *epoch* (a pass through whole training set)
  - *Set ΔΘ = 0*
  - For each labeled image:
    - Read data to initialize input layer
    - Evaluate network to get $y$ (forward)
    - Compare with target label $t$ to get error $E$
    - Backpropagate error derivative to get parameter updates
    - Accumulate parameter updates into $ΔΘ$
  - $Θ_{i+1} = Θ_i - εΔΘ$

Aggregate gradient update most accurately reflects true gradient

25

---

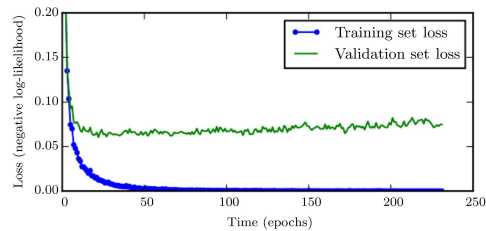## Mini-batch Stochastic Gradient

- For each batch in training set
  - For each labeled image in batch:
    - Read data to initialize input layer
    - Evaluate network to get $y$ (forward)
    - Compare with target label $t$ to get error $E$
    - Backpropagate error derivative to get parameter updates
    - Accumulate parameter updates into $ΔΘ$
  - $Θ_{i+1} = Θ_i - εΔΘ$

Balance between accuracy of gradient estimation and parallelism
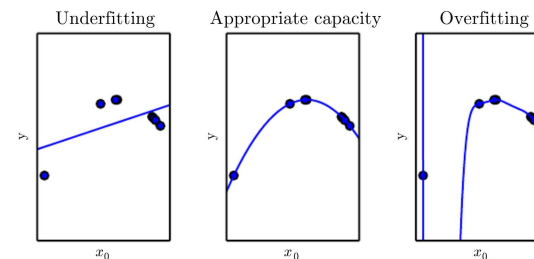
26

---

## When is Training Done?



Split labeled data into *training* and *test* sets.
- Training data to compute parameter updates.
- Test data to check how model generalizes to new inputs (the ultimate goal!)
- The network can become *too good* at classifying training inputs!
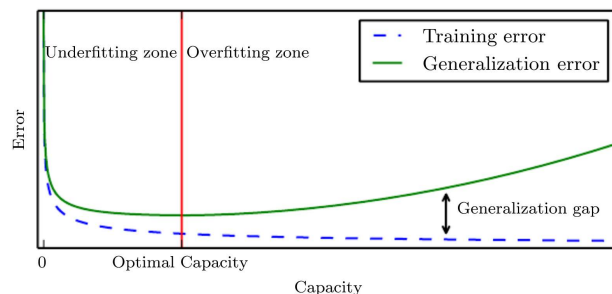
27

---

## How Complicated Should a Network Be?



Intuition: like a polynomial fit. High-order terms improve fit, but add unpredictable swings for inputs outside the training set.

28

## Overtraining Decreases Accuracy



Underfitting zone | Overfitting zone

- - - Training error
—— Generalization error

Error

Generalization gap

0    Optimal Capacity

Capacity

If network works too well for training data,
new inputs cause big unpredictable output changes.

29

---

## No Free Lunch Theorem

- Every classification algorithm has the same error rate when classifying previously unobserved inputs when averaged over all possible input-generating distributions.

- Neural networks must be tuned for specific tasks

30

---

## Summary (1)

- Classification:
  - $f : \mathbb{R}^N \to \{1, \ldots, C\}$
  - $k[i] = f(x, \Theta)$
- Current ML work driven by cheap compute, lots of available data
- Perceptron as a trivial deep network
  - $y = sign(W \bullet x + b)$
- Forward for inference, backward for training

31

---

## Summary (2)

- Chain rule to compute parameter updates
- Stochastic gradient descent for training

32