# Neural networks

CISC 5800
Professor Daniel Leeds
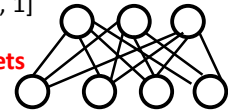
---

## Two breeds of deep networks

Discriminative: $\text{unit}^k(\mathbf{x}) = (\mathbf{w}^k)^T\mathbf{x}+b = [0, 1]$

**Neural networks / Convolutional neural networks**

Generative: $\text{unit}^k(\mathbf{x}) = P(\mathbf{x}; \boldsymbol{\theta}^k) = [0, 1]$

**Bayes Nets / Deep Belief Nets**

2

---

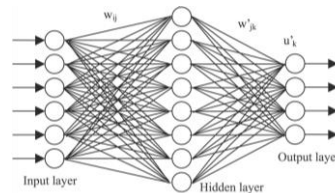## Network architecture

Input layer:
• Compute based on initial features

"Hidden" layers
• Compute based on new features
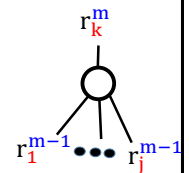
"Output" layer
• Output final class or high-level features

**Each unit takes inputs from past layer, outputs to next layer**
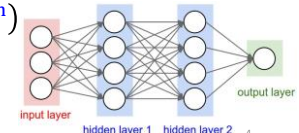
3

---

## Neural network building blocks

$r_k^m$

Individual unit "perceptron":
• Typically logistic function $\text{unit}(\mathbf{x}) = g(h=\mathbf{w}^T\mathbf{x}+b)=\frac{1}{1+e^{-h}}$

$r_1^{m-1} \bullet\bullet\bullet r_j^{m-1}$

Inter-layer computations
• Output $r_{\text{unit\#}}^{\text{level}}$: $r_k^m = g\left(\sum_j w_{k,j}^m r_j^{m-1} + b_k^m\right)$
• Parameters $w_{\text{unit\#,input\#}}^{\text{level}}$ : $w_{k,j}^m$
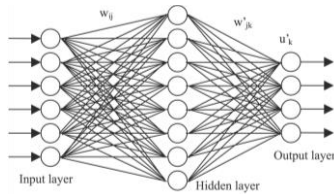
4

1

## Flow of calculation

Calculate output of each unit at layer 1 (based on input)
Calculate output of each unit at layer 2 (based on layer 1)
⋮
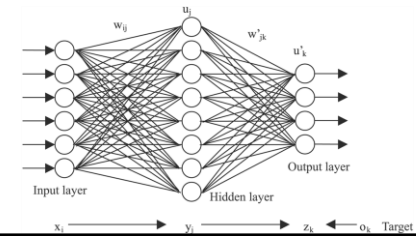Calculate output of each unit at layer out (based on layer out-1)



## Top layer units

$r^{top}_{classY}$    Find the unit with $r^{top}=1$ – that is your class

$r^{top}_{newFeatK}$  Use outputs of all $r^{top}$ for new classifier (e.g., SVM)
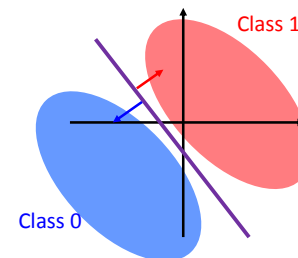


---

Parameters: $w^{m}_{k,i}$ - weights for every unit

Hyper-parameters:
• number of layers
• number of units per layer
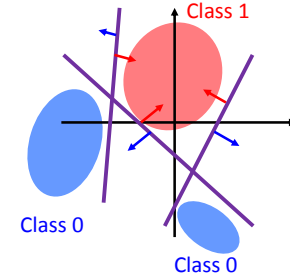• (sigmoid alternatives g(…) with hyper-parameters)
• learning step weight

## Neural Network units dividing feature space

Layer 1 unit          Layer 2 unit
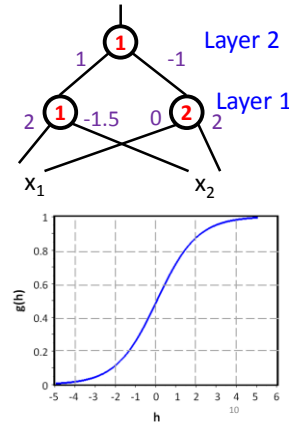


2

3/13/2019

**Simple feedforward practice**

Find $r_1^1$, $r_2^1$, $r_1^2$   Assume b=0

$x_1$=0.1   $x_2$=0.9

$r_1^1$=sigmoid(0.1x2+0.9x-1.5)=
  sigmoid(0.2-1.35)=sigmoid(-1.15)=0.2

$r_2^1$=sigmoid(0.1x0+0.9x2)=
  sigmoid(0+1.8)=sigmoid(1.8)=0.85

$r_1^2$=sigmoid(0.2x1+0.85x-1)=
  sigmoid(0.2-0.85)=sigmoid(-0.65)=0.3



Layer 2

1   **1**   -1

Layer 1

2   **1**  -1.5   0   **2**   2

$x_1$        $x_2$

g(h)   h

---

**Simple feedforward practice**

What is $w_{1,2}^2$?  It is 1
     $w_{2,1}^1$?  It is 0



Layer 2

1   **1**   -1

Layer 1

2   **1**  -1.5   0   **2**   2

$x_1$        $x_2$

g(h)   h

---

**Learning parameters: back-propagation**

Training data: input **x**$^i$ and class **y**$^i$

Compute **x**$^i$'s output for all units, from layer 1 to layer out

Adjust weights in reverse order

   Δw for each unit at layer out (based on y$^i$)
   Δw for each unit at layer out-1 (based on layer out)
      ⋮
   Δw for each unit at layer 1 (based on layer 2)



---

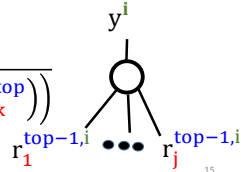Parameters: $w_{k,i}^m$ - weights for every unit

Hyper-parameters:
• number of layers
• number of units per layer
• (sigmoid alternatives g(…) with hyper-parameters)
• $\epsilon$ learning step weight

14

## Learning parameters: back-propagation

First: Change **w** in layer top for each unit

Want to minimize the error, as measured $\sum_i \left(r_k^{top,i} - y_k^i\right)^2$

$$\cdot \, r_k^{top} = g\left(\sum_j w_{k,j}^{top} r_j^{top-1} + b_k^{top}\right)$$

$$= \frac{1}{1 + \exp\left(-\left(\sum_i w_{k,j}^{top} r_j^{top-1} + b_k^{top}\right)\right)}$$

$y^i$

$r_1^{top-1,i} \; \bullet\bullet\bullet \; r_j^{top-1,i}$

15

---

## Δw at each layer

Calculate change to w's at layer top

$$\Delta w_{k,j}^{top} = \epsilon\left(1 - r_k^{top,i}\right)\left(y^i - r_k^{top,i}\right)r_k^{top,i} r_j^{top-1,i}$$

Error correction   input j effect

Define error signal $\delta$: $\delta_k^{top,i} = \left(1 - r_k^{top,i}\right)\left(y^i - r_k^{top,i}\right)r_k^{top,i}$

Send error signal back to layer m-1

16

---

## Simple feedback practice, part 1

Layer 2

Layer 1

Inputs:      $x_1$=0.1        $x_2$=0.9
Outputs:    $r_1^1$=0.3       $r_2^1$=0.8
              $r_1^2$= 0.4

y=1          $\epsilon$=0.1

**Update layer 2 unit:**

$$\Delta w_{k,j}^{top} = \epsilon\left(1 - r_k^{top,i}\right)\left(y^i - r_k^{top,i}\right)r_k^{top,i} r_j^{top-1,i}$$

$x_1$          $x_2$

1          -1
2    -1.5    0    2

---

## Simple feedback practice, part 1

Layer 2

Layer 1

• $x_1$=0.1              $x_2$=0.9
• $r_1^1$=0.2      $r_2^1$=0.8
• $r_1^2$=0.3

• y=1          $\epsilon$=0.1

• $\Delta w_{k,j}^{top} = \epsilon\left(1 - r_k^{top,i}\right)\left(y^i - r_k^{top,i}\right)r_k^{top,i} r_j^{top-1,i}$

$x_1$          $x_2$

1          -1
2    -1.5    0    2

18

## Δw at non-top layer

Top layer error signal $\delta$: $\delta_k^{top,i} = \left(1 - r_k^{top,i}\right)\left(y^i - r_k^{top,i}\right)r_k^{top,i}$

Calculate change to w's at layer m<top

$$\Delta w_{k,j}^m = \epsilon\left(1 - r_k^{m,i}\right)\underbrace{\left(\sum_n w_{n,k}^{m+1,i}\delta_n^{m+1,i}\right)}_{\text{Error correction}}r_k^{m,i}\underbrace{r_j^{m-1,i}}_{\text{input j effect}}$$
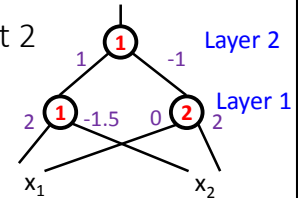
Define error signal $\delta$: $\delta_k^m = \left(1 - r_k^{m,i}\right)\sum_n\left(w_{n,k}^{m+1,i}\delta_n^{m+1,i}\right)r_k^{m,i}$

21

## Simple feedback practice, part 2

Layer 2
Layer 1

Inputs: $x_1$=0.1    $x_2$=0.9
Outputs: $r_1^1$=0.2    $r_2^1$=0.8
$r_1^2$= 0.3

y=1        $\epsilon$=0.1

**Update layer 1, unit 2:**

$$\delta_k^{top,i} = \left(1 - r_k^{top,i}\right)\left(y^i - r_k^{top,i}\right)r_k^{top,i}$$
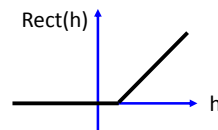$$\Delta w_{k,j}^m = \epsilon\left(1 - r_k^{m,i}\right)\left(\sum_n w_{n,k}^{m+1,i}\delta_n^{m+1,i}\right)r_k^{m,i}r_j^{m-1,i}$$
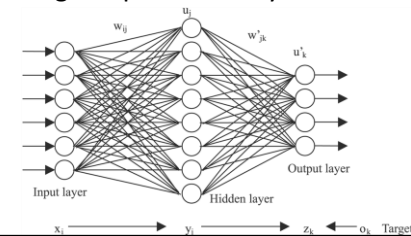


22

## Alternative input transformations

- Traditional: Sum+sigmoid g( $\mathbf{w}^T\mathbf{x}$ + b )
- Straight sum $\mathbf{w}^T\mathbf{x}$ + b
- Sum+rectify
- Max  (no weights!)

Rect(h)

h

24

## Alternative weights

- Traditional: weight inputs from layer m-1
- Competition/Normalization: weight inputs from layer m
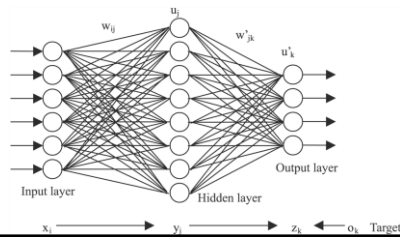- Feedback: weight inputs from layer m+1



25

5

## Top layer units

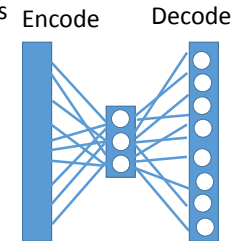$r^{top}_{classY}$     Find the unit with $r^{top}=1$ – that is your class

$r^{top}_{newFeatK}$   Use outputs of all $r^{top}$ for new classifier (e.g., SVM)



26

## Autoencoder

- Input data in large feature space (e.g., 1,000,000 features)

- Intermediate layer has small number of units (e.g., 100 units)

- Output layer has same number of units as input features (e.g., 1,000,000 units)

- Optimize network so output units produce same values as inputs

- Middle units are reduced feature space!

Encode     Decode



27

## Convolutional neural networks

- Wait until end of semester!

28