# Deep Learning: Feedforward Neural Nets and Convolutional Neural Nets

Piyush Rai

Machine Learning (CS771A)

Nov 2, 2016

-

**(D)** < ((())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) <

### A Prelude: Linear Models

• Linear models are nice and simple



• Were some of the first models for learning from data (e.g., Perceptron, 1958)

メロト メぼト メヨト メヨト

### A Prelude: Linear Models

• Linear models are nice and simple



- Were some of the first models for learning from data (e.g., Perceptron, 1958)
- But linear models have limitations: Can't learn nonlinear functions



3

A B + A B + A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

∃ ► < ∃ ►</p>

### A Prelude: Linear Models

• Linear models are nice and simple



- Were some of the first models for learning from data (e.g., Perceptron, 1958)
- But linear models have limitations: Can't learn nonlinear functions



• Before kernel methods (e.g., SVMs) were invented, people thought about this a lot and tried to come up with ways to address this

3

• Composed of several Perceptron-like units arranged in multiple layers



• Consists of an input layer, one or more hidden layers, and an output layer

3

イロン 不同と 不同と 不同と

• Composed of several Perceptron-like units arranged in multiple layers



- Consists of an input layer, one or more hidden layers, and an output layer
- Nodes in the hidden layers compute a nonlinear transform of the inputs

-

• Composed of several Perceptron-like units arranged in multiple layers



- Consists of an input layer, one or more hidden layers, and an output layer
- Nodes in the hidden layers compute a nonlinear transform of the inputs
- Also called a Feedforward Neural Network

ъ.

• Composed of several Perceptron-like units arranged in multiple layers



- Consists of an input layer, one or more hidden layers, and an output layer
- Nodes in the hidden layers compute a nonlinear transform of the inputs
- Also called a Feedforward Neural Network
- "Feedforward": no backward connections between layers (no loops)

ъ.

• Composed of several Perceptron-like units arranged in multiple layers



- Consists of an input layer, one or more hidden layers, and an output layer
- Nodes in the hidden layers compute a nonlinear transform of the inputs
- Also called a Feedforward Neural Network
- "Feedforward": no backward connections between layers (no loops)
- Note: All nodes between layers are assumed connected with each other

ъ.

• Composed of several Perceptron-like units arranged in multiple layers



- Consists of an input layer, one or more hidden layers, and an output layer
- Nodes in the hidden layers compute a nonlinear transform of the inputs
- Also called a Feedforward Neural Network
- "Feedforward": no backward connections between layers (no loops)
- Note: All nodes between layers are assumed connected with each other
- Universal Function Approximator (Hornik, 1991): A one hidden layer FFNN with sufficiently large number of hidden nodes can approximate any function

• Composed of several Perceptron-like units arranged in multiple layers



- Consists of an input layer, one or more hidden layers, and an output layer
- Nodes in the hidden layers compute a nonlinear transform of the inputs
- Also called a Feedforward Neural Network
- "Feedforward": no backward connections between layers (no loops)
- Note: All nodes between layers are assumed connected with each other
- Universal Function Approximator (Hornik, 1991): A one hidden layer FFNN with sufficiently large number of hidden nodes can approximate any function
  - Caveat: This results is only in terms of theoretical feasibility. Learning the model can be very difficult in practice (e.g., due to optimization difficulties)

## What do Hidden Layers Learn?

• Hidden layers can automatically extract features from data



◆□▶ ◆□▶ ◆∃▶ ◆∃▶ → ヨ → のへで

# What do Hidden Layers Learn?

• Hidden layers can automatically extract features from data



• The bottom-most hidden layer captures very low level features (e.g., edges). Subsequent hidden layers learn progressively more high-level features (e.g., parts of objects) that are composed of previous laver's features ◆□▶ ◆□▶ ◆∃▶ ◆∃▶ → ヨ → のへで Machine Learning (CS771A)

• Below: FFNN with 4 inputs, one hidden layer with 3 nodes, and 1 output



3

イロン 不同と 不同と 不同と

• Below: FFNN with 4 inputs, one hidden layer with 3 nodes, and 1 output



• Each hidden node computes a nonlinear transformation of its incoming inputs

э.

• Below: FFNN with 4 inputs, one hidden layer with 3 nodes, and 1 output



- Each hidden node computes a nonlinear transformation of its incoming inputs
  - Weighted linear combination followed by a nonlinear "activation function"

э.

• Below: FFNN with 4 inputs, one hidden layer with 3 nodes, and 1 output



- Each hidden node computes a nonlinear transformation of its incoming inputs
  - Weighted linear combination followed by a nonlinear "activation function"
  - Nonlinearity required. Otherwise, the model would reduce to a linear model

ъ.

• Below: FFNN with 4 inputs, one hidden layer with 3 nodes, and 1 output



• Each hidden node computes a nonlinear transformation of its incoming inputs

- Weighted linear combination followed by a nonlinear "activation function"
- Nonlinearity required. Otherwise, the model would reduce to a linear model
- Output y is a weighted comb. of the preceding layer's hidden nodes (followed by another transform if y isn't real valued, e.g., binary/multiclass label)
   Machine Learning (CS771A)



• For an FFNN with D inputs  $\boldsymbol{x} = [x_1, \dots, x_D]$ 

3

イロト 不得下 不良下 不良下



• For an FFNN with *D* inputs  $\mathbf{x} = [x_1, \dots, x_D]$ , a single hidden layer with *K* hidden nodes  $\mathbf{h} = [h_1, \dots, h_K]$ 

3

イロン 不同と 不同と 不同と



• For an FFNN with *D* inputs  $\mathbf{x} = [x_1, \dots, x_D]$ , a single hidden layer with *K* hidden nodes  $\mathbf{h} = [h_1, \dots, h_K]$ , and a scalar-valued output node *y* 

$$y = \mathbf{v}^{\top} \mathbf{h}$$

3

イロン 不同と イヨン イヨン



• For an FFNN with *D* inputs  $\mathbf{x} = [x_1, \dots, x_D]$ , a single hidden layer with *K* hidden nodes  $\mathbf{h} = [h_1, \dots, h_K]$ , and a scalar-valued output node *y* 

$$y = \mathbf{v}^{\top} \mathbf{h} = \mathbf{v}^{\top} f(\mathbf{W}^{\top} \mathbf{x})$$

3

イロン 不同と イヨン イヨン



• For an FFNN with *D* inputs  $\mathbf{x} = [x_1, \dots, x_D]$ , a single hidden layer with *K* hidden nodes  $\mathbf{h} = [h_1, \dots, h_K]$ , and a scalar-valued output node *y* 

$$\mathbf{v} = \mathbf{v}^{\top} \mathbf{h} = \mathbf{v}^{\top} f(\mathbf{W}^{\top} \mathbf{x})$$

where  $\boldsymbol{v} = [v_1 \ v_2 \ \dots \ v_K] \in \mathbb{R}^K$ 

3



• For an FFNN with *D* inputs  $\mathbf{x} = [x_1, \dots, x_D]$ , a single hidden layer with *K* hidden nodes  $\mathbf{h} = [h_1, \dots, h_K]$ , and a scalar-valued output node *y* 

$$y = \mathbf{v}^{\top} \mathbf{h} = \mathbf{v}^{\top} f(\mathbf{W}^{\top} \mathbf{x})$$

where  $\boldsymbol{v} = [v_1 \ v_2 \ \dots \ v_K] \in \mathbb{R}^K$ ,  $\boldsymbol{W} = [\boldsymbol{w}_1 \ \boldsymbol{w}_2 \ \dots \ \boldsymbol{w}_K] \in \mathbb{R}^{D \times K}$ 

3

イロン 不同と 不同と 不同と



• For an FFNN with *D* inputs  $\mathbf{x} = [x_1, \dots, x_D]$ , a single hidden layer with *K* hidden nodes  $\mathbf{h} = [h_1, \dots, h_K]$ , and a scalar-valued output node *y* 

$$y = \mathbf{v}^{\top} \mathbf{h} = \mathbf{v}^{\top} f(\mathbf{W}^{\top} \mathbf{x})$$

where  $\mathbf{v} = [v_1 \ v_2 \ \dots \ v_K] \in \mathbb{R}^K$ ,  $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K] \in \mathbb{R}^{D \times K}$ , f is the nonlinear activation function

イロン 不得 とくほど 不良 とうほう



• For an FFNN with *D* inputs  $\mathbf{x} = [x_1, \dots, x_D]$ , a single hidden layer with *K* hidden nodes  $\mathbf{h} = [h_1, \dots, h_K]$ , and a scalar-valued output node *y* 

$$y = \mathbf{v}^{\top} \mathbf{h} = \mathbf{v}^{\top} f(\mathbf{W}^{\top} \mathbf{x})$$

where  $\mathbf{v} = [v_1 \ v_2 \ \dots \ v_K] \in \mathbb{R}^K$ ,  $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K] \in \mathbb{R}^{D \times K}$ , f is the nonlinear activation function

• Each hidden node's value is computed as:  $h_k = f(\mathbf{w}_k^\top \mathbf{x}) = f(\sum_{d=1}^D w_{dk} \mathbf{x}_d)$ 

Machine Learning (CS771A)

# (Deeper) Feedforward Neural Net



Note: The hidden layer ℓ contains K<sub>ℓ</sub> hidden nodes, W<sup>(1)</sup> is of size D × K<sub>1</sub>, W<sup>(ℓ)</sup> for ℓ ≥ 2 is of size K<sub>ℓ</sub> × K<sub>ℓ+1</sub>, v is of size K<sub>L</sub> × 1

Machine Learning (CS771A)

• Some popular choices for the nonlinear activation function f

- Sigmoid:  $f(x) = \sigma(x) = \frac{1}{1 + \exp(-x)}$  (range between 0-1)
- tanh:  $f(x) = 2\sigma(2x) 1$  (range between -1 and +1)
- Rectified Linear Unit (ReLU):  $f(x) = \max(0, x)$



글 > - - - 글 >

• Some popular choices for the nonlinear activation function f

- Sigmoid:  $f(x) = \sigma(x) = \frac{1}{1 + \exp(-x)}$  (range between 0-1)
- tanh:  $f(x) = 2\sigma(2x) 1$  (range between -1 and +1)
- Rectified Linear Unit (ReLU):  $f(x) = \max(0, x)$



• Sigmoid saturates and can kill gradients. Also not "zero-centered"

글 > - - - 글 >

• Some popular choices for the nonlinear activation function f

- Sigmoid:  $f(x) = \sigma(x) = \frac{1}{1 + \exp(-x)}$  (range between 0-1)
- tanh:  $f(x) = 2\sigma(2x) 1$  (range between -1 and +1)
- Rectified Linear Unit (ReLU):  $f(x) = \max(0, x)$



- Sigmoid saturates and can kill gradients. Also not "zero-centered"
- tanh also saturates but is zero-centered (thus preferred over sigmoid)

• Some popular choices for the nonlinear activation function f

- Sigmoid:  $f(x) = \sigma(x) = \frac{1}{1 + \exp(-x)}$  (range between 0-1)
- tanh:  $f(x) = 2\sigma(2x) 1$  (range between -1 and +1)
- Rectified Linear Unit (ReLU):  $f(x) = \max(0, x)$



- Sigmoid saturates and can kill gradients. Also not "zero-centered"
- tanh also saturates but is zero-centered (thus preferred over sigmoid)
- ReLU is currently the most popular (also cheap to compute)

Machine Learning (CS771A)

• Want to learn the parameters by minimizing some loss function



• Backpropagation (gradient descent + chain rule for derivatives) is commonly used to do this efficiently

Machine Learning (CS771A)

• Consider the feedforward neural net with one hidden layer



• Recall that  $\boldsymbol{h} = [h_1 \ h_2 \ \dots \ h_K] = f(\mathbf{W}^\top \boldsymbol{x})$ 

-

• Consider the feedforward neural net with one hidden layer



- Recall that  $\boldsymbol{h} = [h_1 \ h_2 \ \dots \ h_K] = f(\mathbf{W}^\top \boldsymbol{x})$
- Assuming a regression problem, the optimization problem would be

$$\min_{\mathbf{W},\mathbf{v}} \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \mathbf{v}^{\top} f(\mathbf{W}^{\top} \mathbf{x}_n) \right)^2 = \min_{\mathbf{W},\mathbf{v}} \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\mathbf{w}_k^{\top} \mathbf{x}_n) \right)^2$$

where  $\boldsymbol{w}_k$  is the *k*-th column of the  $D \times K$  matrix  $\boldsymbol{W}$ 

Machine Learning (CS771A)

-

• We can learn the parameters by doing gradient descent (or stochastic gradient descent) on the objective function

$$\mathcal{L} = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right)^2 = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \boldsymbol{v}^{\top} \boldsymbol{h}_n \right)^2$$

-

• We can learn the parameters by doing gradient descent (or stochastic gradient descent) on the objective function

$$\mathcal{L} = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right)^2 = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \boldsymbol{v}^{\top} \boldsymbol{h}_n \right)^2$$

• Gradient w.r.t.  $\boldsymbol{v} = [v_1 \ v_2 \ \dots \ v_K]$  is straightforward

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = -\sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\mathbf{w}_k^{\top} \mathbf{x}_n) \right) \mathbf{h}_n = -\sum_{n=1}^{N} \mathbf{e}_n \mathbf{h}_n$$

<ロ> (四) (四) (三) (三) (三) (三)

• We can learn the parameters by doing gradient descent (or stochastic gradient descent) on the objective function

$$\mathcal{L} = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right)^2 = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \boldsymbol{v}^{\top} \boldsymbol{h}_n \right)^2$$

• Gradient w.r.t.  $\boldsymbol{v} = [v_1 \ v_2 \ \dots \ v_K]$  is straightforward

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{v}} = -\sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right) \boldsymbol{h}_n = -\sum_{n=1}^{N} \boldsymbol{e}_n \boldsymbol{h}_n$$

• Gradient w.r.t. the weights  $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$  is a bit more involved due to the presence of f but can be computed using chain rule

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_k} = \frac{\partial \mathcal{L}}{\partial f_k} \frac{\partial f_k}{\partial \boldsymbol{w}_k} \quad (\text{note: } f_k = f(\boldsymbol{w}_k^\top \boldsymbol{x}))$$

◆□ → ◆□ → ◆三 → ◆三 → ● ● ●

• We can learn the parameters by doing gradient descent (or stochastic gradient descent) on the objective function

$$\mathcal{L} = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right)^2 = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \boldsymbol{v}^{\top} \boldsymbol{h}_n \right)^2$$

• Gradient w.r.t.  $\boldsymbol{v} = [v_1 \ v_2 \ \dots \ v_K]$  is straightforward

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{v}} = -\sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right) \boldsymbol{h}_n = -\sum_{n=1}^{N} \boldsymbol{e}_n \boldsymbol{h}_n$$

• Gradient w.r.t. the weights  $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$  is a bit more involved due to the presence of f but can be computed using chain rule

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_k} = \frac{\partial \mathcal{L}}{\partial f_k} \frac{\partial f_k}{\partial \boldsymbol{w}_k} \quad (\text{note: } f_k = f(\boldsymbol{w}_k^\top \boldsymbol{x}))$$

• We have: 
$$\frac{\partial \mathcal{L}}{\partial f_k} = -\sum_{n=1}^{N} (y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n)) v_k = -\sum_{n=1}^{N} \boldsymbol{e}_n v_k$$

◆□ → ◆□ → ◆三 → ◆三 → ● ● ●

• We can learn the parameters by doing gradient descent (or stochastic gradient descent) on the objective function

$$\mathcal{L} = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right)^2 = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \boldsymbol{v}^{\top} \boldsymbol{h}_n \right)^2$$

• Gradient w.r.t.  $\boldsymbol{v} = [v_1 \ v_2 \ \dots \ v_K]$  is straightforward

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{v}} = -\sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right) \boldsymbol{h}_n = -\sum_{n=1}^{N} \boldsymbol{e}_n \boldsymbol{h}_n$$

• Gradient w.r.t. the weights  $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$  is a bit more involved due to the presence of f but can be computed using chain rule

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_k} = \frac{\partial \mathcal{L}}{\partial f_k} \frac{\partial f_k}{\partial \boldsymbol{w}_k} \quad (\text{note: } f_k = f(\boldsymbol{w}_k^\top \boldsymbol{x}))$$

- We have:  $\frac{\partial \mathcal{L}}{\partial f_k} = -\sum_{n=1}^{N} (y_n \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n)) v_k = -\sum_{n=1}^{N} \boldsymbol{e}_n v_k$
- We have:  $\frac{\partial f_k}{\partial \boldsymbol{w}_k} = \sum_{n=1}^N f'(\boldsymbol{w}_k^\top \boldsymbol{x}_n) \boldsymbol{x}_n$ , where  $f'(\boldsymbol{w}_k^\top \boldsymbol{x}_n)$  is f's derivative at  $\boldsymbol{w}_k^\top \boldsymbol{x}_n$

(ロ) (同) (三) (三) (三) (000)

• We can learn the parameters by doing gradient descent (or stochastic gradient descent) on the objective function

$$\mathcal{L} = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right)^2 = \frac{1}{2} \sum_{n=1}^{N} \left( y_n - \boldsymbol{v}^{\top} \boldsymbol{h}_n \right)^2$$

• Gradient w.r.t.  $\boldsymbol{v} = [v_1 \ v_2 \ \dots \ v_K]$  is straightforward

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{v}} = -\sum_{n=1}^{N} \left( y_n - \sum_{k=1}^{K} v_k f(\boldsymbol{w}_k^{\top} \boldsymbol{x}_n) \right) \boldsymbol{h}_n = -\sum_{n=1}^{N} \boldsymbol{e}_n \boldsymbol{h}_n$$

• Gradient w.r.t. the weights  $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$  is a bit more involved due to the presence of f but can be computed using chain rule

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_k} = \frac{\partial \mathcal{L}}{\partial f_k} \frac{\partial f_k}{\partial \boldsymbol{w}_k} \quad (\text{note: } f_k = f(\boldsymbol{w}_k^\top \boldsymbol{x}))$$

- We have:  $\frac{\partial \mathcal{L}}{\partial f_k} = -\sum_{n=1}^N (y_n \sum_{k=1}^K v_k f(\boldsymbol{w}_k^\top \boldsymbol{x}_n)) v_k = -\sum_{n=1}^N \boldsymbol{e}_n v_k$
- We have:  $\frac{\partial f_k}{\partial \boldsymbol{w}_k} = \sum_{n=1}^N f'(\boldsymbol{w}_k^\top \boldsymbol{x}_n) \boldsymbol{x}_n$ , where  $f'(\boldsymbol{w}_k^\top \boldsymbol{x}_n)$  is f's derivative at  $\boldsymbol{w}_k^\top \boldsymbol{x}_n$
- These calculations can be done efficiently using backpropagation

Machine Learning (CS771A)

・ロト ・ 日 ・ モ ト ・ ヨ ・ うへの

## Backpropagation

• Basically consists of a forward pass and a backward pass



- Forward pass computes the errors  $e_n$  using the current parameters
- Backward pass computes the gradients and updates the parameters, starting from the parameters at the top layer and then moving backwards

## Backpropagation

• Basically consists of a forward pass and a backward pass



- Forward pass computes the errors  $e_n$  using the current parameters
- Backward pass computes the gradients and updates the parameters, starting from the parameters at the top layer and then moving backwards
- Also good at reusing previous computations (updates of parameters at any layer depends on parameters at the layer above)

Machine Learning (CS771A)

### Kernel Methods vs Deep Neural Nets

• Recall the prediction rule for a kernel method (e.g., kernel SVM)

$$y = \sum_{n=1}^{N} \alpha_n k(\boldsymbol{x}_n, \boldsymbol{x})$$

- This is analogous to a single hidden layer NN with fixed/pre-defined hidden nodes  $\{k(\mathbf{x}_n, \mathbf{x})\}_{n=1}^N$ and output layer weights  $\{\alpha_n\}_{n=1}^N$
- The prediction rule for a deep neural network

$$v = \sum_{k=1}^{K} v_k h_k$$

- In this case, the  $h_k$ 's are learned from data (possibly after multiple layers of nonlinear transformations)
- Both kernel methods and deep NNs be seen as using nonlinear basis functions for making predictions. Kernel methods use fixed basis functions (defined by the kernel) whereas NN learns the basis functions adaptively from data

Machine Learning (CS771A)

### Wide vs Deep?

Why might we prefer a deep model over a wide and shallow model?

3

イロト 不得下 不良下 不良下

### Wide vs Deep?

Why might we prefer a deep model over a wide and shallow model? An informal justification:

- Deep "programs" can reuse computational subroutines (and are more compact)



э.

## Wide vs Deep?

Why might we prefer a deep model over a wide and shallow model? An informal justification:

- Deep "programs" can reuse computational subroutines (and are more compact)



Learning Certain functions may require a huge number of units in a shallow model

Machine Learning (CS771A)

• A feedforward neural network with a special structure

3

- A feedforward neural network with a special structure
- Sparse "local" connectivity between layers (except the last output layer). Reduces the number of parameters to be learned



-

**(D)** < ((())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) <

- A feedforward neural network with a special structure
- Sparse "local" connectivity between layers (except the last output layer). Reduces the number of parameters to be learned



• Shared weights (like a "global" filter). Helps capture the local properties of the signal (useful for data such as images or time-series)



**(D)** < ((())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (()) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) < (())) <

• Uses a sequence of 2 operations, convolution and pooling (subsampling), applied repeatedly on the input data



• Uses a sequence of 2 operations, convolution and pooling (subsampling), applied repeatedly on the input data



• Convolution: Extract "local" properties of the signal. Uses a set of "filters" that have to be learned (these are the "weighted" **W** between layers)

A ID > A IP > A

E> < E>

• Uses a sequence of 2 operations, convolution and pooling (subsampling), applied repeatedly on the input data



- Convolution: Extract "local" properties of the signal. Uses a set of "filters" that have to be learned (these are the "weighted" **W** between layers)
- Pooling: Downsamples the outputs to reduce the size of representation

・ 同 ト ・ ヨ ト ・ ヨ ト

• Uses a sequence of 2 operations, convolution and pooling (subsampling), applied repeatedly on the input data



- Convolution: Extract "local" properties of the signal. Uses a set of "filters" that have to be learned (these are the "weighted" **W** between layers)
- Pooling: Downsamples the outputs to reduce the size of representation
- Note: A nonlinearity is also introduced after the convolution layer

Machine Learning (CS771A)

### Convolution

• An operation that captures local (e.g., spatial) properties of a signal



• Mathematically, the operation is defined as

$$h_{ij}^k = f((W^k * \mathbf{X})_{ij} + b_k)$$

where  $W^k$  is a filter, \* is the convolution operator, and f is a nonlinearity

- Usually a number of filters { W<sup>k</sup> }<sup>K</sup><sub>k=1</sub> are applied (each will produce a separate "feature map"). These filters have to be learned
- Size of these filters have to be specified

Machine Learning (CS771A)

3

This operation is used to reduce the size of the representation



э.

• Highly effective in learning good feature representations from data in an "end-to-end" manner

-

- Highly effective in learning good feature representations from data in an "end-to-end" manner
- The objective functions of these models are highly non-convex
  - But lots of recent work on non-convex optimization, so non-convexity doesn't scare us (that much) anymore

-

- Highly effective in learning good feature representations from data in an "end-to-end" manner
- The objective functions of these models are highly non-convex
  - But lots of recent work on non-convex optimization, so non-convexity doesn't scare us (that much) anymore
- Training these models is computationally very expensive
  - But GPUs can help to speed up many of the computations

-

- Highly effective in learning good feature representations from data in an "end-to-end" manner
- The objective functions of these models are highly non-convex
  - But lots of recent work on non-convex optimization, so non-convexity doesn't scare us (that much) anymore
- Training these models is computationally very expensive
  - But GPUs can help to speed up many of the computations
- Training these models can be tricky, especially a proper initialization
  - But now we have several ways to intelligently initialize these models (e.g., unsupervised layer-wise pre-training)

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ ○○

- Highly effective in learning good feature representations from data in an "end-to-end" manner
- The objective functions of these models are highly non-convex
  - But lots of recent work on non-convex optimization, so non-convexity doesn't scare us (that much) anymore
- Training these models is computationally very expensive
  - But GPUs can help to speed up many of the computations
- Training these models can be tricky, especially a proper initialization
  - But now we have several ways to intelligently initialize these models (e.g., unsupervised layer-wise pre-training)
- Deep learning models can also be probabilistic and generative, e.g., deep belief networks (we did not consider these here)

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ ○○