

## Linear models: Recap

<https://eduassistpro.github.io/>

Linear models:

- ▶ Perceptron [Assignment Project Exam Help](#)

[Assignment Help](#) [WeChat](#) [edu\\_assist\\_pro](#)

- ▶ Naïve Bay

<https://eduassistpro.github.io/>

$$\log P(y|x; \theta) = \log P(x|y; \phi) + \log B(x) + \theta \cdot f(x, y)$$

[Add WeChat](#) [edu\\_assist\\_pro](#)

- ▶ Logistic Regression

$$\log P(y|x; \theta) = \theta \cdot f(x, y) - \log \sum_{y' \in \mathcal{Y}} \exp \theta \cdot f(x, y')$$

## Features and weights in linear models: Recap

- ▶ Feature representation:  
 $\theta = [\underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=1}; \underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=2}; \dots; \underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=K}]$

~~Assignment Project Exam Help  
Assignment Add WeChat edu\_assist\_pro~~  
 $\theta = [\underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=1}; \underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=2}; \dots; \underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=K}]$   
<https://eduassistpro.github.io/>  
~~Add WeChat edu\_assist\_pro~~

- ▶ Weights:  $\theta$

$$\theta = [\underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=1}; \underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=2}; \dots; \underbrace{\theta_1; \theta_2; \dots; \theta_V}_{y=K}]$$

## Rearranging the features and weights

<https://eduassistpro.github.io/>

- ▶ Represent the features  $\mathbf{x}$  as a *column* vector of length  $V$ , and represent the weights as a  $\Theta$  as  $K$

~~Assignment Project Exam Help  
Assignment Help WeChat edu\_assist\_pro~~

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_V \end{bmatrix} \quad \Theta = \begin{bmatrix} \theta_{1,1} & \theta_{1,2} & \cdots & \theta_{1,V} \\ \theta_{2,1} & \theta_{2,2} & \cdots & \theta_{2,V} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{K,1} & \theta_{K,2} & \cdots & \theta_{K,V} \end{bmatrix}$$

<https://eduassistpro.github.io/>  
[Add WeChat](#) [edu\\_assist\\_pro](#)

- ▶ What is  $\Theta\mathbf{x}$ ?

Scores for each class

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ Verify that  $\psi_1, \psi_2, \dots, \psi_K$  corresponds to the class [Assignment Help](#) [Add WeChat](#) [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

$\Psi = \Theta x =$   
[Add WeChat](#) [edu\\_assist\\_pro](#)

Implementation in Pytorch

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

## Digression: Matrix multiplication

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- Matrix with  $m$  rows and  $n$  columns:

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

where  $C_{ij} = \sum_{k=1}^n A_{ik} B_{kj}$

- Example:

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

$$\begin{bmatrix} 2 & 3 \\ 1 & 2 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 1 & 4 \end{bmatrix}$$

Digression: 3-D matrix multiplication

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

Tensor shape: (batch-size, sentence-length, embedding size)

## SoftMax

- ▶ SoftMax, also known as <https://eduassistpro.github.io/>

$$\text{Assignment Project Exam Help}$$
$$\text{SoftMax}_i(\psi) = \frac{\exp(\psi_i)}{K}$$

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)  
for  $i = 1, 2, \dots, K$

- ▶ Applying SoftMax distribution <https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$$\text{SoftMax}(\Psi) = \begin{bmatrix} P(y=1) \\ P(y=2) \\ \dots \\ P(y=K) \end{bmatrix}$$

- ▶ Verify this is exactly logistic regression

Logistic regression as a neural network

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$$\mathbf{y} = \text{SoftMax}(\Theta \mathbf{x})$$

$$V = 5 \ K = 3$$

## Going deep

- ▶ There is no reason. <https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$$\mathbf{z} = \sigma(\Theta_1 \mathbf{x})$$

$$\mathbf{y} = \text{SoftMax}(\Theta_2 \mathbf{z})$$

## Going even deeper

- ▶ There is no reason w  
<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$$\mathbf{z}_1 = \sigma(\Theta_1 \mathbf{x})$$

$$\mathbf{z}_2 = \sigma(\Theta_2 \mathbf{z}_1)$$

$$\mathbf{y} = \text{SoftMax}(\Theta_3 \mathbf{z}_2)$$

- ▶ But why?

## Non-linear classification

Linear models like Logistic regression work well in a high-dimensional vector space, and they are expressive enough and work well for many NLP problems, why do we need more complex non-linear models?

- ▶ There have been many advances in deep learning, a family of nonlinear models, through multiple <https://eduassistpro.github.io/>
- ▶ Deep learning facilitates the incorporation of word embeddings, which are dense vectors that can be learned from massive amounts of text. It has evolved from early static embeddings (e.g., Word2vec, Glove) to recent dynamic embeddings (ELMO, BERT, XLNet).
- ▶ Rapid advances in specialized hardware called graphic processing units (GPUs). Many deep learning models can be implemented efficiently on GPUs.

## Feedforward Neural networks: an intuitive justification

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ In image classification, instead of using the input (pixels) to predict the image type directly, you can imagine that you can predict the shapes of parts of an image such as hand, ear.
- ▶ In text processing, say we want to classify movie reviews (or movie descriptions) into a label set of {Good, Bad, ...}. Instead of predicting these labels directly, we first predict a set of components such as the story, acting, soundtrack, cinematography, etc. from raw input (words in the text).

Face Recognition

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

## Feedforward neural networks

<https://eduassistpro.github.io/>

Formally, this is what we do:

### Assignment Project Exam Help

- ▶ **Use the text  $x$  to predict the feature  $z$ .** Specifically, train a logistic regression classifier to compute  $P(z=k)$  for each  $k \in \{1, 2, \dots, K_z\}$ .
- ▶ **Use the features  $z$  to predict the class  $y$ .** Train a logistic regression classifier to compute  $P(y=k|z)$  for each  $k \in \{1, 2, \dots, K_y\}$ .

Caveat: it's easy to demonstrate what this is what's actually going on in image processing, but it's hard to show this is what's actually going on in language processing. Interpretability is a major issue in neural models for language processing.

## The hidden layer: computing the composite features

<https://eduassistpro.github.io/>

- If we assume each  $z_k$  is binary, that is,  $z_k \in \{0, 1\}$ , then  $P(z_k|x)$  can be modeled with binary logistic regression:

$$P(z_k = 1|x; \Theta^{(x \rightarrow z)}) = \sigma(\theta_k^T x) = \frac{1}{1 + e^{-\theta_k^T x}}$$

- The weight matrix  $\Theta^{(x \rightarrow z)}$  is produced by

stacking (not concatenating, as in linear algebra) vectors for each  $z_k$ ,

$$\Theta^{(x \rightarrow z)} = [\theta_1^{x \rightarrow z}, \theta_2^{x \rightarrow z}, \dots, \theta_{k_z}^{x \rightarrow z}]$$

- We assume an offset/bias term is included in  $x$  and its parameter is included in each  $\theta_k^{x \rightarrow z}$

Notations:  $\Theta^{(x \rightarrow z)} \in \mathbb{R}^{k_z \times V}$  is a real number matrix with a dimension of  $k_z$  rows and  $V$  columns

## Activation functions

- ▶ Sigmoid: The range is  $(0, 1)$ .  
<https://eduassistpro.github.io/>

$$\sigma(x) = \frac{e^x}{1 + e^{-x}}$$

**Assignment Project Exam Help**

Question: what's the value of the sigmoid function at  $x = 0$ ?

[Assignment Help WeChat edu\\_assist\\_pro](#)

- ▶ Tanh: The range of the tanh activation function is  $(-1, 1)$ .

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

**Add WeChat edu\_assist\_pro**

Question: what's the value of the tanh function at  $x = 0$ ?

- ▶ ReLU: The **rectified linear unit (ReLU)** is zero for negative inputs, and linear for positive inputs.

$$ReLU(x) = \max(x, 0) = \begin{cases} 0 & x < 0 \\ x & \text{otherwise} \end{cases}$$

Sigmoid and tanh are sometimes described as **squashing functions**.

## Activation functions in Pytorch

<https://eduassistpro.github.io/>

```
from torch import nn  
import torch  
  
input = torch.FloatTensor(4)  
sigmoid = nn.Si  
output = sigmoid(input)  
  
tanh = nn.Tanh()  
output = tanh(input)  
  
relu = nn.ReLU()  
output = relu(input)
```

## The output layer

<https://eduassistpro.github.io/>

- ▶ The output layer is computed by the multiclass logistic regression probability

$$P(y=j|z; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \frac{e^{(\mathbf{z} \cdot \mathbf{b}_j')}}{\sum_{i=1}^{K_y} e^{(\mathbf{z} \cdot \mathbf{b}'_i)}}$$

<https://eduassistpro.github.io/>

- ▶ The weight matrix  $\Theta^{(z \rightarrow y)} \in \mathbb{R}^{y \times z}$  is constructed by stacking weight vectors for each class.

$$\Theta^{(z \rightarrow y)} = [\theta_1^{z \rightarrow y}, \theta_2^{z \rightarrow y}, \dots, \theta_{K_y}^z]$$

- ▶ The vector of probabilities over each possible value of  $y$  is denoted:

$$P(\mathbf{y}|z; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \text{SoftMax}(\Theta^{(z \rightarrow y)} \mathbf{z} + \mathbf{b})$$

## The negative loglikelihood or cross-entropy loss

In a multi-class classification, it measures the negative conditional likelihood of the true label  $y^{(i)}$  given a probabilistic distribution over possible labels. It works well together with negative loglikelihood (just like logistic regression)

Add WeChat <https://eduassistpro.github.io/>  
Assignment Project Exam Help  
 $\hat{y}_j \triangleq \Pr_{\theta}^{(N)}(y^{(i)} = j | x^{(i)})$

or cross entropy loss  
<https://eduassistpro.github.io/>

Add WeChat <https://eduassistpro.github.io/>  
 $\hat{y}_j \triangleq \Pr(y = j)$

$$-\mathcal{L} = - \sum_{i=1}^N \mathbf{e}_{y^{(i)}} \cdot \log \tilde{\mathbf{y}}$$

where  $\mathbf{e}_{y^{(i)}}$  is a **one-hot vector** of zeros with a value of one at the position  $y^{(i)}$

## Alternative loss functions

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ There are alternatives to SoftMax and cross-entropy loss, just as there are alternatives in linear models.
- ▶ Pairing an affine transformation (remember h) with a margin loss:

<https://eduassistpro.github.io/>

$\Psi(y; \mathbf{x}^{(i)}, \Theta)$  Add WeChat edu\_assist\_pro

$$\ell_{\text{MARGIN}}(\Theta; \mathbf{x}^{(i)}, y^{(i)}) = \max_{y \neq y^{(i)}} \left( 1 + \Psi(y; \mathbf{x}^{(i)}, \Theta) - \Psi(y^{(i)}; \mathbf{x}^{(i)}, \Theta) \right)$$

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Training a feedforward network with the  
AssignmentWeChat edu\_assist\_pro

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

## A feedforward network with a cross-entropy loss

<https://eduassistpro.github.io/>

The following sums up a feed-forward neural network with one hidden layer, complete with a cross-entropy loss th parameter estimation  
[Assignment Project Exam Help](#)  
[Assignment Help WeChat edu\\_assist\\_pro](#)

z <https://eduassistpro.github.io/>

$$\begin{aligned}\tilde{\mathbf{y}} &\leftarrow \text{SoftMax}(\Theta^{z \rightarrow y} z) \\ \ell^{(i)} &\leftarrow -\mathbf{e}_{y^{(i)}} \cdot \log \tilde{\mathbf{y}}\end{aligned}$$

where  $f$  is an elementwise activation function (e.g.,  $\sigma$  or ReLU),  $\ell^{(i)}$  is the loss at instance  $i$

## Updating the parameters of a feedforward network

<https://eduassistpro.github.io/>

As usual,  $\Theta^{x \rightarrow z}$ ,  $\Theta^{z \rightarrow y}$ , and  $b$  can be estimated by online gradient based optimization such as stochastic gradient descent:

[Assignment Project Exam Help](#)  
Drop  $\eta$  here

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

where  $\eta^{(t)}$  is the learning rate at iteration  $t$ ,  $\ell^{(i)}$  is the loss on instance (or minibatch)  $i$ ,  $\theta_n^{x \rightarrow z}$  is the  $n$ th column of the matrix  $\Theta^{x \rightarrow z}$ , and  $\theta_k^{z \rightarrow y}$  is the  $k$ th column of the matrix  $\Theta^{z \rightarrow y}$ .

Compute the gradient of the cross-entropy loss on hidden layer weights  $\Theta^{z \rightarrow}$

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

<https://eduassistpro.github.io/>

<https://eduassistpro.github.io/>

$$\begin{aligned}\frac{\partial \ell^{(i)}}{\partial \theta_{k,j}^{z \rightarrow y}} &= -\frac{\partial}{\partial \theta_{k,j}^{z \rightarrow y}} \left( \theta_{y^{(i)}}^{(z \rightarrow y)} \right)_{y \in \mathcal{Y}} \\ &= \left( P(y = j | \mathbf{z}; \Theta^{(z \rightarrow y)}, \mathbf{b}) - \delta(j = y^{(i)}) \right) z_k\end{aligned}$$

where  $\delta(j = y^{(i)})$  returns 1 if  $j = y^{(i)}$  and 0 otherwise,  $z_k$  is the  $k$ th element in  $\mathbf{z}$ .

Applying the chain rule to compute the derivatives

<https://eduassistpro.github.io/>

We use the chain rule to break down the gradient of the loss on the hidden layer weights:

~~Assignment Project Exam Help~~ Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

where

<https://eduassistpro.github.io/>

$$\ell^{(i)} = -\mathbf{e}_{y^{(i)}} \log \hat{\mathbf{y}} = -\log \left( \frac{\sum_k e^{o_k}}{\sum_k e^{o_k}} \right)$$
$$o_j = \theta_j^{(z \rightarrow y)} z$$

Note:  $o_i$  is the logit that corresponds to the true label  $y^{(i)}$

Derivative of the cross-entropy loss with respect to the logits

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>

$$\begin{aligned} \frac{\partial \ell^{(i)}}{\partial o_j} &= \frac{\partial}{\partial o_j} \left( -\sum_k \log \frac{e^{o_k}}{\sum_l e^{o_l}} \right) \\ &= \frac{e^{o_j}}{\sum_k e^{o_k}} - \frac{\partial}{\partial o_j} \sum_k \log \frac{e^{o_k}}{\sum_l e^{o_l}} \\ &= \tilde{y}_j - \delta(j = y^{(i)}) \end{aligned}$$

## Gradient on the hidden layer weights $\Theta^{z \rightarrow y}$

One more step to compute  
the weights:

<https://eduassistpro.github.io/>

$$\frac{\partial o_j}{\partial \theta_{k,j}^{(z \rightarrow y)}} = \frac{\partial}{\partial \theta_{k,j}^{(z \rightarrow y)}} \theta_{j,k}^{(z \rightarrow y)} \cdot z = \frac{\partial}{\partial \theta_{j,k}^{(z \rightarrow y)}} \theta_{j,k}^{(z \rightarrow y)} z_k = z_k$$

[Assignment](#) [Project](#) [Exam](#) [Help](#)  
[Add](#) [WeChat](#) [edu\\_assist\\_pro](#)

$$\frac{\partial \ell^{(i)}}{\partial \theta_{k,j}^{(z \rightarrow y)}} = \frac{\partial \ell^{(i)}}{\partial o_j} \frac{\partial o_j}{\partial \theta_{k,j}^{(z \rightarrow y)}} = (\tilde{y}_j - \delta(j = y^{(i)})) z_k$$

<https://eduassistpro.github.io/>

Similarly, we can also compute the derivative w  
hidden units: [Add](#) [WeChat](#) [edu\\_assist\\_pro](#)

$$\frac{\partial o_j}{\partial z_k} = \theta_{k,j}^{(z \rightarrow y)}$$

$$\frac{\partial \ell^{(i)}}{\partial z_k} = \frac{\partial \ell^{(i)}}{\partial o_j} \frac{\partial o_j}{\partial z_k} = (\tilde{y}_j - \delta(j = y^{(i)})) \theta_{k,j}^{(z \rightarrow y)}$$

This value will be useful for computing the gradient of the loss  
function from the inputs to the hidden layer.

Compute the gradient on the input weights  $\Theta^{(x \rightarrow z)}$

<https://eduassistpro.github.io/>

Apply the old chain rule in differentiation:

$$\frac{\partial \ell^{(l)}}{\partial \theta_{jk}^{(x \rightarrow z)}} = \frac{\partial \ell^{(l)}}{\partial z_k} \frac{\partial z_k}{\partial \theta_{jk}^{(x \rightarrow z)}}$$

~~Assignment Project Exam Help  
Add WeChat edu\_assist\_pro~~  
 $(x \rightarrow z)$

<https://eduassistpro.github.io/>

Add WeChat  $\frac{\partial \ell^{(l)}}{\partial z_k} f'$  (~~Assignment Project Exam Help~~  
 $\theta_{jk}^{(x \rightarrow z)} \cdot x$ )

where  $f'(\theta_k^{(x \rightarrow z)} \cdot x)$  is the derivative of the activation function  $f$ , applied at the input  $\theta_k^{(x \rightarrow z)} \cdot x$ . Depending on what the actual activation is, the derivative will also be different.

## Derivative of the sigmoid activation function

<https://eduassistpro.github.io/>

$$\frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{x \rightarrow z}} = \frac{\partial \ell^{(i)}}{\partial z_k} \sigma \left( \theta_k^{(x \rightarrow z)} \cdot \mathbf{x} \right) \left( 1 - \sigma \left( \theta_k^{(x \rightarrow z)} \cdot \mathbf{x} \right) \right) x_n$$

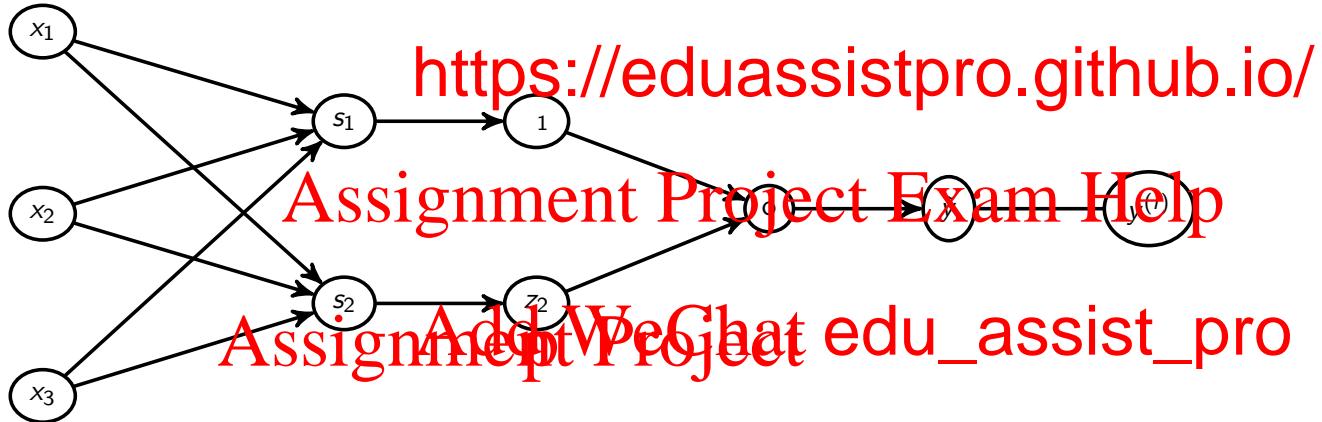
~~Assignment Project Exam Help  
Assignment Help WeChat edu\_assist\_pro~~

$$= \frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{x \rightarrow z}}$$

<https://eduassistpro.github.io/>

- ▶ If the negative log-likelihood  $\ell^{(i)}$  is much on  $z_k$ , then  $\frac{\partial \ell^{(i)}}{\partial z_k} \approx 0$ . In this case it doesn't matter if  $z_k$  is computed, and so  $\frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{x \rightarrow z}} \approx 0$
- ▶ If  $x_n = 0$ , then it does not matter how we set the weights  $\frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{x \rightarrow z}}$ , so  $\frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{x \rightarrow z}} = 0$

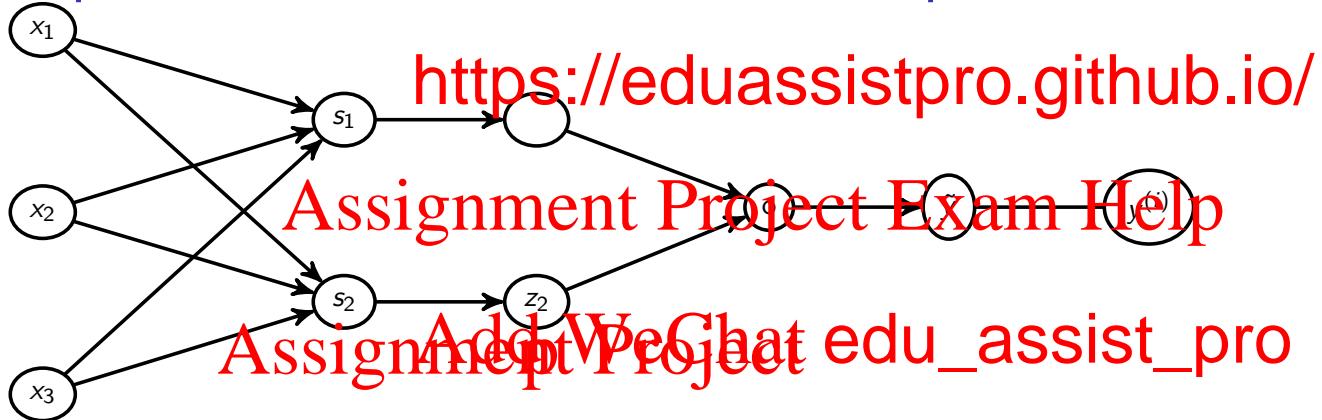
## A simple neural network



$$\Theta^{(x \rightarrow s)} = \begin{bmatrix} s_1 & s_2 \\ \theta_{11}^{(x \rightarrow s)} & \theta_{21}^{(x \rightarrow s)} & \theta_{12}^{(x \rightarrow s)} & \theta_{22}^{(x \rightarrow s)} & \theta_{13}^{(x \rightarrow s)} & \theta_{23}^{(x \rightarrow s)} \end{bmatrix}$$

$$\Theta^{(z \rightarrow o)} = \begin{bmatrix} z_1 & z_2 \\ \theta_{11}^{(z \rightarrow o)} & \theta_{12}^{(z \rightarrow o)} \end{bmatrix}$$

## A simple neural network: forward computation



$s_1 = \sigma(x_1)$

$s_2 = \sigma(x_2)$

$z_1 = \sigma(s_1)$

$z_2 = \sigma(s_2)$

$o = \theta_{11}^{(z \rightarrow o)} z_1 + \theta_{12}^{(z \rightarrow o)} z_2$

$\tilde{y} = \sigma(o)$

Note: When making predictions with the sigmoid function, choose the label 1 if  $\tilde{y} > 0.5$ . Otherwise choose 0

## A simple neural network with cross-entropy sigmoid loss

<https://eduassistpro.github.io/>

- ▶ Cross-entropy sigmoid loss:

$$\mathcal{L}_{H(p,q)} = - \sum_i p_i \log q_i = -y \lvert g(1 - \tilde{y}) \rvert$$

- ▶  $y \in \{0, 1\}$  i

- ▶  $\tilde{y}$  is the predicted value

- ▶ If the true label is 1, loss is  $-\log(\tilde{y})$ . Loss is bigger if  $\tilde{y}$  is smaller

- ▶ If the true label is 0, loss is  $-\log(1 - \tilde{y})$ . Loss is bigger if  $\tilde{y}$  is bigger

A simple neural network: Compute the derivatives

<https://eduassistpro.github.io/>

$\frac{\partial \mathcal{L}}{\partial o} = \tilde{y} - y$  Assignment Project Exam Help

$\frac{\partial o}{\partial z_1} = \theta_{11}^{(z \rightarrow o)}$   $\frac{\partial o}{\partial z_2} = \theta_{12}^{(z \rightarrow o)}$   
Assignment Help Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$\frac{\partial \theta_{11}^{(z)}}{\partial s_1}$  <https://eduassistpro.github.io/>

$\frac{\partial z_1}{\partial s_1} = z_1(1 - z_1)$   $\frac{\partial z_2}{\partial s_2} =$   
Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$\frac{\partial s_1}{\partial \theta_{11}^{(x \rightarrow s)}} = x_1$   $\frac{\partial s_1}{\partial \theta_{12}^{(x \rightarrow s)}} = x_2$   $\frac{1}{\partial \theta_{13}^{(x \rightarrow s)}} = x_3$

$\frac{\partial s_2}{\partial \theta_{21}^{(x \rightarrow s)}} = x_1$   $\frac{\partial s_2}{\partial \theta_{22}^{(x \rightarrow s)}} = x_2$   $\frac{\partial s_2}{\partial \theta_{23}^{(x \rightarrow s)}} = x_3$

Compute the derivatives on the weights  $\Theta^{z \rightarrow o}$

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Apply the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \theta_{12}^{(z \rightarrow o)}} = \frac{\partial \mathcal{L}}{\partial \theta} \frac{\partial \theta}{\partial \theta_{12}^{(z \rightarrow o)}} = (\tilde{y} - y) z_1$$

<https://eduassistpro.github.io/>

$$\frac{\partial \theta_{12}^{(z \rightarrow o)}}{\partial \theta_{12}^{(z \rightarrow o)}} = \frac{\partial \theta}{\partial \theta_{12}^{(z \rightarrow o)}} \frac{\partial \theta}{\partial \theta_{12}^{(z \rightarrow o)}} = 1$$

Compute derivatives with respect to the weights  $\Theta^{x \rightarrow s}$

$$\frac{\partial \mathcal{L}}{\partial \theta_{11}^{(x \rightarrow s)}} = \frac{\partial \mathcal{L}}{\partial o} \frac{\partial o}{\partial z_1} \frac{1}{\partial s_1} \frac{1}{\partial \theta_{11}^{(x \rightarrow s)}} = (\tilde{y} - y) \theta_{11}^{(z \rightarrow o)} z_1 (1 - z_1) x_1$$

Assignment Project Exam Help  
Assignment Help WeChat edu\_assist\_pro

$$\frac{\partial \mathcal{L}}{\partial \theta_{12}^{(x \rightarrow s)}} = \frac{\partial \mathcal{L}}{\partial o} \frac{\partial o}{\partial z_1} \frac{\partial z_1}{\partial s_1} \frac{\partial s_1}{\partial \theta_{12}^{(x \rightarrow s)}} = (y - \tilde{y}) \theta_{12}^{(z \rightarrow o)} z_1^2 (1 - z_1) x_2$$

Assignment Help WeChat edu\_assist\_pro

$$\frac{\partial \mathcal{L}}{\partial \theta_{13}^{(x \rightarrow s)}} = \frac{\partial \mathcal{L}}{\partial o} \frac{\partial o}{\partial z_1} \frac{\partial z_1}{\partial s_1} \frac{\partial s_1}{\partial \theta_{13}^{(x \rightarrow s)}} = (\tilde{y} - y) \theta_{13}^{(z \rightarrow o)} z_1 (1 - z_1) x_3$$

https://eduassistpro.github.io/

$$\frac{\partial \mathcal{L}}{\partial \theta_{21}^{(x \rightarrow s)}} = \frac{\partial \mathcal{L}}{\partial o} \frac{\partial o}{\partial z_2} \frac{\partial z_2}{\partial s_2} \frac{\partial s_2}{\partial \theta_{21}^{(x \rightarrow s)}} = (\tilde{y} - y) \theta_{21}^{(z \rightarrow o)} z_2 (1 - z_2) x_1$$

Add WeChat edu\_assist\_pro

$$\frac{\partial \mathcal{L}}{\partial \theta_{22}^{(x \rightarrow s)}} = \frac{\partial \mathcal{L}}{\partial o} \frac{\partial o}{\partial z_2} \frac{\partial z_2}{\partial s_2} \frac{\partial s_2}{\partial \theta_{22}^{(x \rightarrow s)}} = (y - \tilde{y}) \theta_{22}^{(z \rightarrow o)} z_2^2 (1 - z_2) x_2$$
$$\frac{\partial \mathcal{L}}{\partial \theta_{23}^{(x \rightarrow s)}} = \frac{\partial \mathcal{L}}{\partial o} \frac{\partial o}{\partial z_2} \frac{\partial z_2}{\partial s_2} \frac{\partial s_2}{\partial \theta_{23}^{(x \rightarrow s)}} = (\tilde{y} - y) \theta_{23}^{(z \rightarrow o)} z_2 (1 - z_2) x_3$$

## Vectorizing forward computation

<https://eduassistpro.github.io/>

$$\mathbf{s} = \Theta^{(x \rightarrow s)} \mathbf{x} = \begin{bmatrix} \theta_{11}^{(x \rightarrow s)} & \theta_{12}^{(x \rightarrow s)} & \theta_{13}^{(x \rightarrow s)} \\ \theta_{21}^{(x \rightarrow s)} & \theta_{22}^{(x \rightarrow s)} & \theta_{23}^{(x \rightarrow s)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

~~Assignment Project Exam Help  
Assignment Add WeChat edu\_assist\_pro~~

$$\theta_{11}^{(x \rightarrow s)} x_1 + \theta_{12}^{(x \rightarrow s)} x_2 + \theta_{13}^{(x \rightarrow s)} x_3 = \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$$

<https://eduassistpro.github.io/>

$$\mathbf{z} = \sigma \left( \begin{bmatrix} s_1 \\ s_2 \end{bmatrix} \right) = \begin{bmatrix} \sigma(s_1) \\ \sigma(s_2) \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

~~Add WeChat edu\_assist\_pro~~

$$\mathbf{o} = \Theta^{(z \rightarrow o)} \mathbf{z} = \begin{bmatrix} \theta_{11}^{(z \rightarrow o)} & \theta_{12}^{(z \rightarrow o)} \end{bmatrix} \times \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} \theta_{11}^{(z \rightarrow o)} z_1 + \theta_{12}^{(z \rightarrow o)} z_2 \end{bmatrix}$$

$$\tilde{y} = \sigma(\mathbf{o})$$

Note: Summing over inputs in forward computation

# Backpropagation

<https://eduassistpro.github.io/>

- ▶ When applying the chain rule, local derivatives are frequently reused. For example,  $\frac{\partial \mathcal{L}}{\partial \theta}$  is used to compute the gradient on both  $\Theta^{z \rightarrow o}$  and  $\Theta^{x \rightarrow s}$ . It would be useful to reuse them.
- ▶ We can only compute the derivatives of the parameters if we have all the necessary derivatives from the previous layer. So careful sequencing of differentiation is required.
- ▶ This combination of *sequencing* and *differentiation* is called **backpropagation**.
- ▶ In order to make this process manageable, backpropagation is typically done in vectorial form (tensors)

## Elements in a neural network

<https://eduassistpro.github.io/>

- ▶ Variables includes input  $x$ , hidden units  $z$ , outputs  $y$ , and the loss function.
- ▶ Inputs are not computed from other nodes i  
example, inputs to a feedforward neural net  
feature vector extracted from a training (o
- ▶ Backpropagation is used to compute gradients with respect to all variables.
- ▶ Parameters include weights and bias. They are initialized and then updated during descent.
- ▶ Loss is not used to compute any other nodes in the graph, and is usually computed with the predicted label  $\hat{y}$  and the true label  $y(i)$ .

Backpropagation: caching “error signals”

<https://eduassistpro.github.io/>

Let  $D_o$  represent the gradient of the loss function with respect to  $o$  and  $D_s$  be the gradient of the loss function with respect to  $s$ .

Assuming a cross-entropy loss and a sigmoid (which is generalized to SoftMax) function, the error signal a  $o$  is:

$$D_o = [\tilde{y} - y^{(i)}]$$

This error signal can be calculated with respect to  $\Theta^{(z \rightarrow o)}$

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

$$\begin{aligned}\nabla_{\Theta^{(z \rightarrow o)}} &= D_o z = [\tilde{y} - y] [z_1 \quad z_2] \\ &= [(\tilde{y} - y)z_1 \quad (\tilde{y} - y)z_2]\end{aligned}$$

## Backpropagation: Computing error signals

<https://eduassistpro.github.io/>

Using  $D_o$ , we can also compute the error signals at the hidden layer:

$$\begin{aligned} D_s &= (\Theta^{(z \rightarrow o)} \\ &= \left( \begin{bmatrix} \frac{\partial o}{\partial z_1} & \frac{\partial o}{\partial z_2} & \dots & \frac{\partial o}{\partial s_1} & \dots \end{bmatrix} \right. \\ &= \left( \begin{bmatrix} \theta_{11}^{(z \rightarrow o)} & \theta_{12}^{(z \rightarrow o)} & \dots & \theta_{n1}^{(z \rightarrow o)} & \dots \end{bmatrix} \times \begin{bmatrix} \tilde{y} - y \\ z_1(1 - z_1) \\ z_2(1 - z_2) \\ \vdots \\ z_n(1 - z_n) \end{bmatrix} \right) \\ &= \left[ (\tilde{y} - y)\theta_{11}^{(z \rightarrow o)}z_1(1 - z_1) \quad (\tilde{y} - y)\theta_{12}^{(z \rightarrow o)}z_2(1 - z_2) \quad \dots \right] \end{aligned}$$

Backpropagation: caching “error signals”

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

Using error signals  $D_s$ , we can now compute th

$\Theta^{(x \rightarrow s)}$ :

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

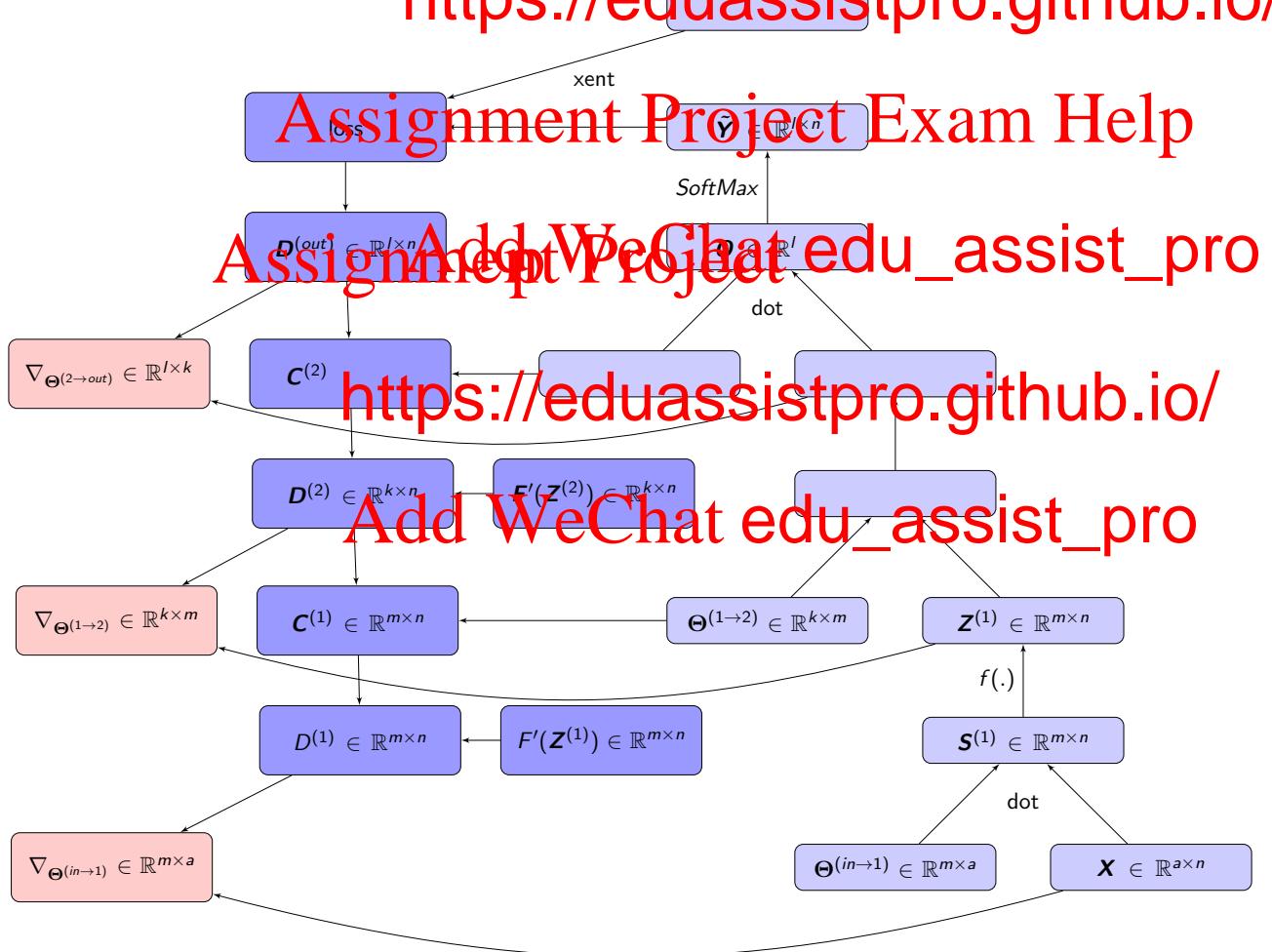
$$\nabla_{\Theta^{(x \rightarrow s)}} = D_s^\top X$$

$$= \begin{bmatrix} (\tilde{y} - y)\theta_{11}^{(z \rightarrow o)} z_1 \\ (\tilde{y} - y)\theta_{12}^{(z \rightarrow o)} z_2(1 - z_2) \end{bmatrix}$$

$$= \begin{bmatrix} (\tilde{y} - y)\theta_{11}^{(z \rightarrow o)} z_1(1 - z_1)x_1 & (\tilde{y} - y)\theta_{11}^{(z \rightarrow o)} z_1(1 - z_1)x_3 \\ (\tilde{y} - y)\theta_{12}^{(z \rightarrow o)} z_2(1 - z_2)x_1 & (\tilde{y} - y)\theta_{12}^{(z \rightarrow o)} z_2(1 - z_2)x_3 \end{bmatrix}$$

# Computation graph of a three-layer feedforward neural network

<https://eduassistpro.github.io/>



## Forward computation in matrix form

<https://eduassistpro.github.io/>

$$\mathbf{S}^{(1)} = \Theta^{in \rightarrow 1} \mathbf{x}$$

$$\mathbf{Z}^{(1)} = f_1(\mathbf{S}^{(1)})$$

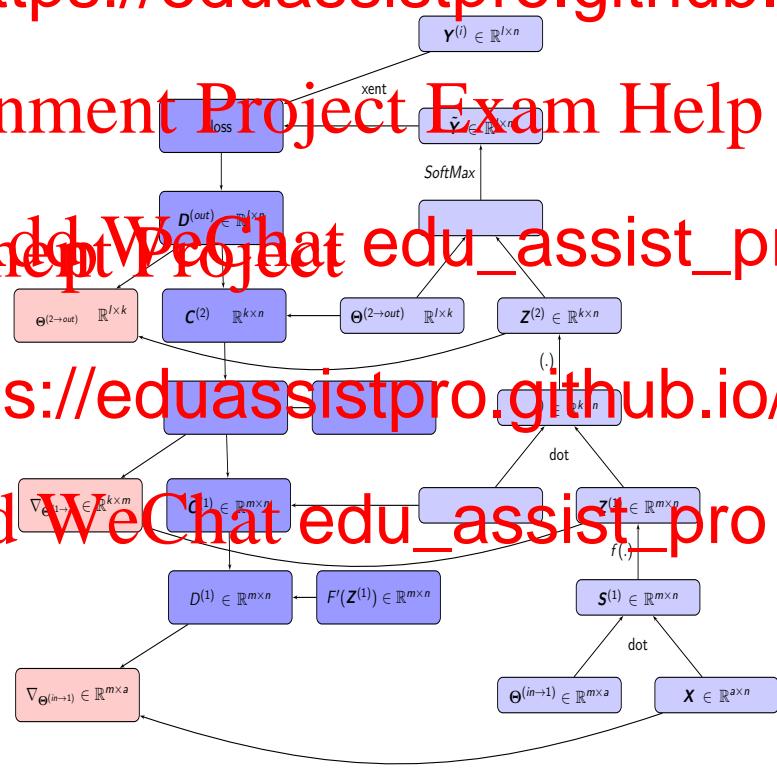
$$\mathbf{S}^{(2)} = \Theta^{1 \rightarrow 2} \mathbf{Z}^{(1)}$$

$$\mathbf{Z}^{(2)} = f_2(\mathbf{S}^{(2)})$$

$$\mathbf{O} = \Theta^{2 \rightarrow o} \mathbf{Z}^{(2)}$$

$$\tilde{\mathbf{y}}_{[j]} = \text{SoftMax}(\mathbf{O}^\top \mathbf{z}_{[j]})$$

Note SoftMax  
applies to a vector.



Matrix multiplication in forward computation sums over inputs, hidden units. We assume the input and hidden units in batches

## Backpropagation: Computing error signals

<https://eduassistpro.github.io/>

$$\mathbf{D}^{out} = \tilde{\mathbf{Y}} - \mathbf{Y}^{(i)}$$

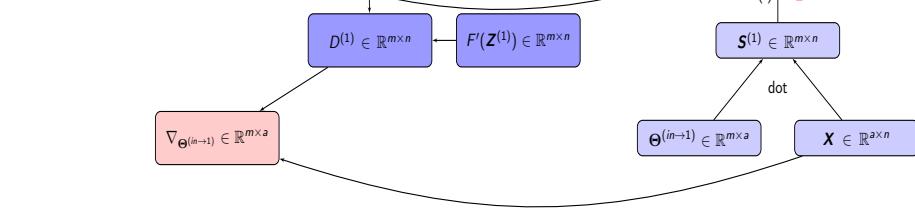
$$\mathbf{C}^{(2)} = (\Theta^{2 \rightarrow o})^\top \mathbf{D}^{out}$$

$$\mathbf{D}^{(2)} = F'(\mathbf{Z}^{(2)}) \odot \mathbf{C}^{(2)}$$

$$\mathbf{C}^{(1)} = (\Theta^{1 \rightarrow 2})^\top \mathbf{D}^{(2)}$$

$$\mathbf{D}^{(1)} = F'(\mathbf{Z}^{(1)}) \odot$$

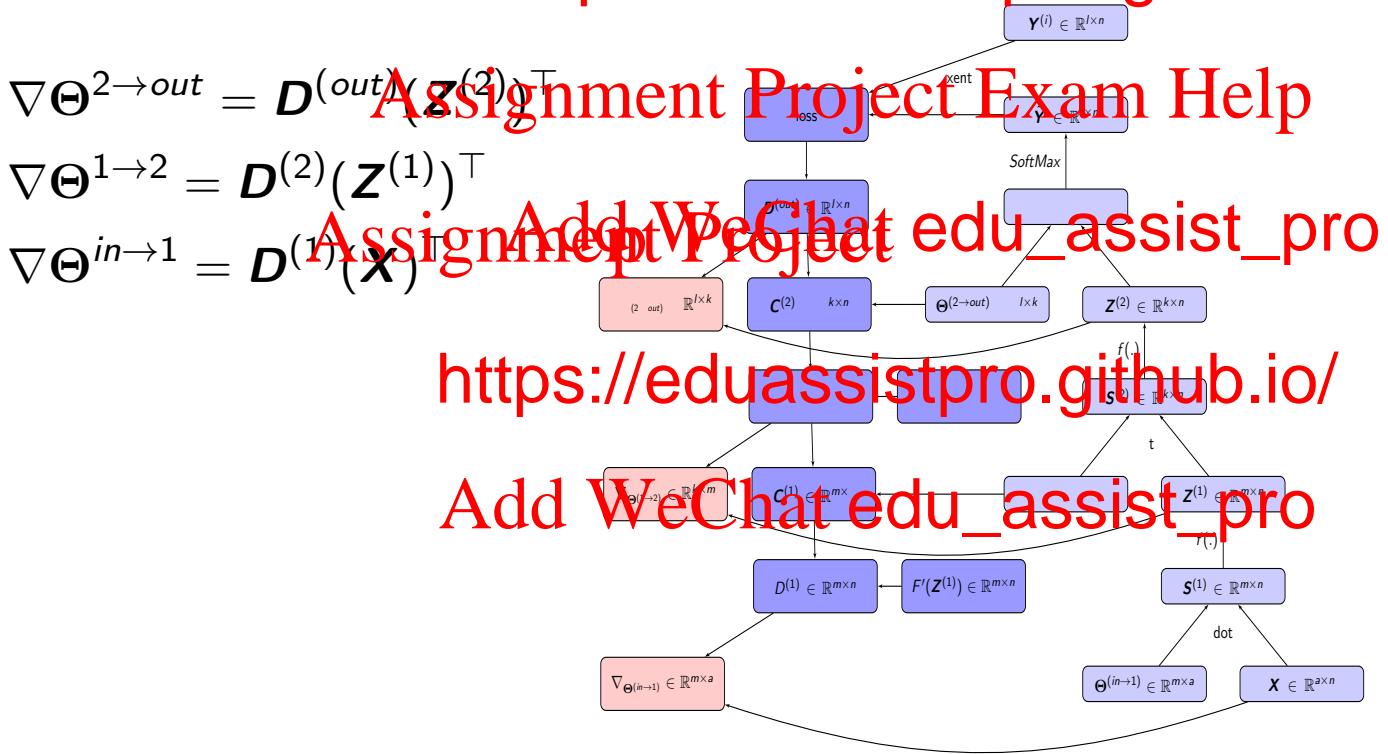
Add WeChat edu\_assist\_pro



Matrix multiplication in error signal computation involves summing over outputs and hidden units.

# Backpropagation: Compute the gradient with error signals

<https://eduassistpro.github.io/>



Matrix multiplication when computing gradient sums over instances in mini-batch.

Notes in backpropagation

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat edu\_assist\_pro

- ▶ Where does caching take place?
- ▶ Why is sequence important?
- ▶ What computation patterns can you observe?

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Pay attention to what you sum over

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ In forward computation, you sum over all input vectors for each output.
- ▶ In backward computation, you sum over all outputs for each input vector.
- ▶ When you compute the gradient for each data point, (and average), make sure you line up the columns of the first rows of the second matrix. That's what you sum over.

## Automatic differentiation software

<https://eduassistpro.github.io/>

- ▶ Backpropagation is mechanical. There is no reason for everyone to repeat the same work
- ▶ Currently many automatic differentiation exist, e.g., Tor
- ▶ Using these libraries, the computation, and the gradient can be computed automatically. These libraries have had development in deep learning consider
- ▶ Libraries that support dynamic computation graphs are better suited for many NLP problems.

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat edu\_assist\_pro  
Bells an

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

## Tricks in training neural networks

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

There are various tricks that people use when training neural networks: <https://eduassistpro.github.io/>

- ▶ Regularization
- ▶ Dropout: Adjusting the learning rate
- ▶ Optimization methods: Adjusting the learning rate
- ▶ Initialization: Using particular forms of initialization

## Regularization

<https://eduassistpro.github.io/>

Neural networks can be regularized in a similar way as linear models. Neural networks can also with **Frobenius norm**, which is a trivial extension to L2 norm for matrices. In fact, in machine learning it is just referred to as L2 regularization.

$$\mathcal{L} = \sum_{i=1}^N \ell^{(i)} + \lambda_{z \rightarrow y} \|\Theta^{(z \rightarrow y)}\|_F^2$$

where  $\|\Theta\|_F^2 = \sum_{i,j} \theta_{i,j}^2$  is the squared **Frobenius norm**, which generalizes the  $L_2$  norm to matrices. The bias parameters  $b$  are not regularized, as they do not contribute to the classifier to the inputs.

## L2 regularization

- ▶ Compute the gradient: <https://eduassistpro.github.io/>

Assignment Project Exam Help

$$\frac{\partial \mathcal{L}}{\partial \theta} = \sum_{i=1}^N \frac{\partial \mathcal{L}_i}{\partial \theta} + \lambda \theta$$

Assignment Help WeChat: [edu\\_assist\\_pro](https://edu_assist_pro)

- ▶ Update the weight: <https://eduassistpro.github.io/>

Add WeChat: [edu\\_assist\\_pro](https://edu_assist_pro)

$$\theta = \theta - \eta \left( \sum_{i=1}^N \frac{\partial \mathcal{L}_i}{\partial \theta} + \lambda \theta \right)$$

- ▶ “Weight decay factor”:  $\lambda$  is a tunable hyper parameter that pulls a weight back when it has become too big
- ▶ Question: Does it matter which layer  $\theta$  is from when computing the regularization term?

## L1 regularization

- ▶ L1 regularization  
<https://eduassistpro.github.io/>

**Assignment Project Exam Help**  
$$\mathcal{L} = \sum_{i=1}^N \ell^{(i)} + \lambda_{z \rightarrow y} \|\Theta^{(z \rightarrow y)}\|_1$$
  
**Add WeChat edu\_assist\_pro**

- ▶ Compute the g

<https://eduassistpro.github.io/>

**Add WeChat edu\_assist\_pro**

- ▶ update the weights

$$\theta = \theta - \eta \left( \sum_{i=1}^N \frac{\partial \ell^{(i)}}{\partial \theta} + \lambda \operatorname{sign}(\theta) \right)$$

## Comparison of L1 and L2

<https://eduassistpro.github.io/>

- ▶ In L1 regularization, the weights shrink by a constant amount toward 0. In L2 regularization, the weights shrink by an amount which is proportional to  $w$ .
- ▶ When a particular weight has a large absolute value, L1 regularization tends to make it small, while L2 regularization tends to make it even larger.
- ▶ The net result is that L1 regularization tends to reduce the weight of the network in a relatively small number of high-importance connections, while the other weights are driven toward zero. So L1 regularization effectively does *feature selection*.

## Dropout

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ Randomly drops a certain percentage of the input over-reliance feature co-adaptation working together goal is to avoid overfitting

## Dropout

<https://eduassistpro.github.io/>

- ▶ Dropout can be ach

Assignment Project Exam Help

$z^{(1)} = g(\Theta^{(1)}x + b^{(1)})$

$m^1 \sim Bernoulli$

$\tilde{z}^{(1)} = m^1 \odot z^{(1)}$

<https://eduassistpro.github.io/>

$m^2 \sim Bernoulli$

Add WeChat edu\_assist\_pro

$\tilde{z}^{(2)} = m^2 \odot z^{(2)}$

$y = \Theta^{(3)} \tilde{z}^{(2)}$

where  $m^1$  and  $m^2$  are mask vectors. The values of the elements in these vectors are either 1 or 0, drawn from a Bernouli distribution with parameter  $r$  (usually  $r = 0.5$ )

## Optimization methods

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ Moment <https://eduassistpro.github.io/> Add WeChat edu\_assist\_pro
- ▶ Adgrad <https://eduassistpro.github.io/> Add WeChat edu\_assist\_pro
- ▶ Root Mean Sq <https://eduassistpro.github.io/> Add WeChat edu\_assist\_pro
- ▶ Adam <https://eduassistpro.github.io/> Add WeChat edu\_assist\_pro

## Momentum

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

At each timestep  $t$ , compute  $\nabla_{\Theta}$  and  $\nabla_{\Theta}$ , and then compute the momentum as follows:

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>

$$\Theta = \Theta - \frac{\nabla_{\Theta}}{t}$$

### Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

The momentum term increases for dimensions pointing in the same directions and reduces updates for dimensions whose gradient change directions.

Adgrad

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

Weight and bias update for Adgrad at each time step  $t$ :

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)  
 $\frac{1}{2}$

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

e.g.,  $\epsilon = 10^{-8}$

## Root Mean Square Prop (RMSProp)

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

Weight update for RMSprop at each time step  $t$ :

~~Assignment Help WeChat edu\_assist\_pro~~  
 $\frac{\partial}{\partial w} \text{MSE}(w) = \frac{2}{n} \sum_{i=1}^n (y_i - w)^2$

<https://eduassistpro.github.io/>

~~Add WeChat edu\_assist\_pro~~

e.g.  $\beta = 0.9, \eta = 0.001, \epsilon = 10^{-8}$

## Adaptive Moment Estimation (Adam)

<https://eduassistpro.github.io/>

Weight update at time st

Assignment Project Exam Help

$$\begin{aligned} \mathbf{V}_{\nabla \Theta} &= \beta_1 \mathbf{V}_{\nabla \Theta} + (1 - \beta_1) \nabla \Theta \\ \mathbf{S}_{\nabla \Theta} &= \beta_2 \mathbf{S}_{\nabla \Theta} + (1 - \beta_2) \nabla^2 \Theta \end{aligned}$$

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

$$\Theta = \Theta - \eta \frac{\mathbf{V}_{\nabla \Theta}^{corrected}}{\sqrt{\mathbf{S}_{\nabla \Theta}^{corrected}} + \epsilon}$$

Adam combines Momentum and RMSProp

## Define a neural net

```
from torch import nn
class Net(nn.Module):
    ''' subclass from nn.Module is important to insp
    def __init__(self, in_dim=25, out_dim=3, batch_
        super(Net, self).__init__()
        self.in_dim = in_dim
        self.out_dim = out_dim
        self.softmax = nn.Softmax(dim=1) #softmax o
    Add WeChat edu_assist_pro
    def forward(self, input_
        logit = self.linear(input_matrix)
        return logit #return raw score, not normalized
    Add WeChat edu_assist_pro
    def xentropy_loss(self, input_matrix, target_label_
        loss = nn.CrossEntropyLoss()
        logits = self.forward(input_matrix)
        return loss(logits, target_label_vec)
```

## Use optimizers in Pytorch

```
import torch.optim
net = Net(input_dim, output_dim)
optimizer = optim.Adam(net.parameters(), lr=lr_rate)
for epoch in range(epochs):
    total_nll = 0
    for batch in dataloader:
        optimizer.zero_grad() #ze
        vecto
        feat_v
        label_vec = map(item
        feat_list = list(feat_
        label_list = list(lab_
        x = torch.Tensor(feat_list)
        y = torch.LongTensor(label_list)
        loss = net.cross_entropy_loss(x,y)
        total_nll += loss
        loss.backward()
        optimizer.step()
    torch.save(net.state_dict(), net_path)
```

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat edu\_assist\_pro  
Sparse and Dense embeddings as input to ne

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

## Input to feedforward neural networks

- ▶ Assuming a bag-of words representation,  $x$  is the count of each word ( $x$ )
- ▶ To compute the hidden unit  $z_k$ :

**Assignment Project Exam Help**

$$z_k = \sum_{j=1}^V \theta_j^x$$

**Assignment Project Exam Help**

- ▶ The connections between the input layer and the hidden layer form a vector embedding of word  $j$ .
- ▶ If there is a lot of training data, word embeddings can be learned within the same network as the classification.
- ▶ Word embeddings can also be learned separately from unlabeled data, using techniques such as Word2Vec and GLOVE.
- ▶ The latest trend is to learn *contextualized* word embeddings which are computed dynamically for each classification instance (e.g., ELMO, BERT). This requires more advanced architectures (Transformers) that we will talk about later in the course.

## One-hot encodings for features

<https://eduassistpro.github.io/>

A *one-hot* encoding is one in which each dimension corresponds to a unique feature, and the resulting feature vector of a c instance can be thought of as a set of indicator features in which a single dimension has a value of one while all others have a value of zero.

<https://eduassistpro.github.io/>

### Example

When considering a bag-of-words representation of 40000 words. A short document of 20 words will be represented with a very sparse 40000-dimensional vector in which **at most** 20 dimensions have non-zero values

Combination of sparse vectors

<https://eduassistpro.github.io/>

Assignment Project Exam Help

~~Assignment Help Project~~  $\begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$  ~~WeChat~~  $\text{edu\_assist\_pro}$   
+

<https://eduassistpro.github.io/>

~~Add WeChat~~  $\text{edu\_assist\_pro}$

=

$[0 \ 1 \ 0 \ 1 \ 0 \ 0 \ 1 \ 0]$

## Sparse vs Dense representations

- ▶ Sparse representation
  - ▶ Each feature is a sparse vector in which one dimension is 1 and the rest are 0s (thus, “one-hot”)
  - ▶ Dimensionality of one-hot vector is same as number of distinct features
  - ▶ Features can be completely independent as it is to “ $w_1$ ”
  - ▶ Feature  $w_d$  is obtained by summation.
- ▶ Dense representation
  - ▶ Each feature is a  $d$ -dimensional vector, with a  $d$  that is generally much shorter than that of a one-hot vector.
  - ▶ Model training will cause similar features to have similar vectors - information is shared between similar features.
  - ▶ Features can be combined via summation (or averaging), concatenation, or some combination of the two.
    - ▶ Concatenation if we care about relative position.

## Using dense vectors in a classifier

<https://eduassistpro.github.io/>

Each core feature is embedded into a  $d$  dimensional space (typically 50-300), and represented as a vector in that space.

- ▶ Extract a set of core linguistic features that are relevant for predicting the output class
- ▶ For each feature vector  $v_i$ . <https://eduassistpro.github.io/>
- ▶ Combine the vectors (either by concatenation, or a combination of both) into an input vector
  - ▶ Note: concatenation doesn't work for variable-length vectors such as document classification
- ▶ Feed  $x$  into a nonlinear classifier (feed-forward neural network).

## Relationship between one-hot and dense vectors

<https://eduassistpro.github.io/>

- ▶ Dense representations are pre-trained word embeddings.
- ▶ One-hot and dense representations may not be as different as one might think.
- ▶ In fact, using sparse, one-hot vectors as input to a neural network to learn [e.g., feature] based on training data.
- ▶ With task-specific word embedding, the dimensionality is typically smaller, but the training objective for the embedding and the task objective are one and same.
- ▶ With pre-trained word embeddings, the training data is easy to come by (just unannotated text), but the embedding object and task objective may diverge.

Two approaches of getting dense word vectors

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ Count based methods, known in NLP as Distr Semantico Models (DPM) or vector Semantico
- ▶ Predictive models for word representations
  - ▶ Distributed word representations, product of neural language models and later became its own

## Distributional semantics

<https://eduassistpro.github.io/>

- ▶ Based on the well-known observation of Z. Harris: Words are similar if they occur in the same context (Harris, 1954)
- ▶ Further summarized it as a slogan: “You shall know a word by the company it keeps.” (J. R. Firth, 1957)
- ▶ A long history of word meaning  
represents a context word it can occur with
- ▶ Each word is represented as a sparse vector in a high-dimensional space
- ▶ Then word distances and similarities can be computed with such a matrix

## Steps for building a distributional semantic model

<https://eduassistpro.github.io/>

- ▶ Preprocess a (large) corpus: tokenization at a minimum, possibly lemmatization, POS tagging, or s
- ▶ Define the “context” for a target term (word).  
The context can be a window centered on the ta  
terms that are s  
(subject-of, o)  
<https://eduassistpro.github.io/>
- ▶ Compute a term-context matrix where  
to a term and each column corresponds to a c  
the target term.
- ▶ Each target term is then represented with a high-dimensional  
vector of context terms.

## Mathematical processing for building a DSM

<https://eduassistpro.github.io/>

- ▶ Weight the term-metrics such as Positive Pointwise Mutual Information (PPMI) to correct frequency bias

$$\text{Assignment Project Exam Help}$$
$$PPMI(x, y) = \max(0, \frac{\log \frac{f(x, y)}{f(x)f(y)}}{-y})$$

- ▶ Its dimension techniques such as singular value decomposition (SVD)

<https://eduassistpro.github.io/>

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$$

$$\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{U} \in \mathbb{R}^{m \times k}, \Sigma \in \mathbb{R}^{k \times k}, \mathbf{V} \in \mathbb{R}^{n \times k}, n \gg k$$

- ▶ This will result in a matrix that has much lower dimension but retains most of the information of the original matrix.

## Predictive methods

- ▶ Learns word embeddings using various language model objectives.
- ▶ Decide on the context window
- ▶ Define the objective function that is used to predict context words based on the target word or previous word based on context.
- ▶ Train the neural network
- ▶ The resulting weight matrix will serve as the representation for the target word
- ▶ “Don’t count, predict!” (Baroni et al, 2014) conducted systematic studies and found predict-based word embeddings outperform count-based embeddings.
- ▶ One of popular early word embeddings are Word2vec embeddings.

## Word2vec

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ Word2vec is a software package that consists of two main models: CBOW (Continuous Bag of Words)
- ▶ It popularized the use of distributed representations to neural networks, inspired many followers.  
<https://eduassistpro.github.io/> (LNet)
- ▶ It has its roots in language modeling (the use of context to predict the target word), but gives up getting good language models and focus instead on getting good word embeddings.

## Understanding word2vec: A simple CBOW model with only one context word in

- ▶ Input  $x \in \mathbb{R}^V$  and  $x_{k'} = 0$  for  $k' \neq k$ .  $\Theta \in \mathbb{R}^{N \times V}$  is the weight matrix from the input layer to the hidden layer. Each column of  $\Theta$  is an  $N$ -dimensional vector representation of the input word  $w_k$ .

- ▶  $\Theta' \in \mathbb{R}^{V \times N}$  is the weight matrix from the hidden layer to the output layer and  $\mathbf{u}_{w_j}$  is the  $j$ -th row of  $\Theta'$ . The score  $o_j$  for each target word  $w_j$  and context  $w$  is computed as:
$$o_j = \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i}$$
- ▶ Finally we use softmax to obtain a posterior distribution

$$p(w_j|w_i) = y_j = \frac{\exp(o_j)}{\sum_{j'=1}^V \exp(o_{j'})}$$

where  $y_j$  is the output of the  $j$ -th unit in the output layer

Computing the hidden layer is just embedding lookup

Hidden layer computation  
<https://eduassistpro.github.io/>

$$v_{w_i} = z = \Theta x =$$

Assignment Project Exam Help

Assignment Project Exam Help WeChat edu\_assist\_pro

$$\begin{bmatrix} 0.1 & 0.3 & 0.5 \\ 0.2 & 0.5 & 0.8 \\ 0.2 & 0.5 & 0.8 \end{bmatrix} \times \begin{bmatrix} 0.6 \\ 0.1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.8 \\ 0.8 \end{bmatrix}$$

Add WeChat edu\_assist\_pro

Note there is no activation at the hidden layer (or there is a linear activation function), so this is a “degenerate neural network”.

Computing the output layer

<https://eduassistpro.github.io/>

**Assignment Project Exam Help**

$$\Theta' = \begin{bmatrix} 0.3 & 0.4 & 0.6 & 0.95 \\ 0.7 & 0.1 & 0.6 & 0.91 \\ 0.5 & 0.2 & 0.7 & 0.97 \\ 0.4 & 0.6 & 0.3 & 0.82 \\ 0.2 & 0.4 & 0.6 & 1.18 \\ 0.3 & 0.2 & 0.6 & 0.31 \\ 0.3 & 0.4 & 0.6 & 1.06 \\ 0.3 & 0.5 & 0.1 & 0.39 \\ 0.3 & 0.6 & 0.5 & 0.95 \\ 0.3 & 0.5 & 0.1 & 0.63 \end{bmatrix}$$

**Add WeChat edu\_assist\_pro**

$$o = \Theta'z$$

<https://eduassistpro.github.io/>

Each row of  $\Theta'$  correspond to vector for a target word  $w_j$ .

Taking the softmax over the output

<https://eduassistpro.github.io/>

Assignment Project Exam Help  
[0.11039215  
0.  
AssignmentProjectExamHelp  
AddWeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

[0.11039215  
0.08016116]

The output  $\mathbf{y}$  is a probabilistic distribution over the entire vocabulary.

Input vector and output vector

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

Since there is no activation function at the hidden layer, the output is really just the dot product of the vector of the in context word and the vector of the output target word.

<https://eduassistpro.github.io/>

$$p(w_j|w_i) = y_j = \frac{(\mathbf{u}_{w_j}^\top \mathbf{v}_{w_i})}{\sum_{j=1}^V \exp(p_j)}$$

where  $\mathbf{v}_{w_i}$  from  $\Theta$  is the **input vector** for word  $w_i$  and  $\mathbf{u}_{w_j}$  from  $\Theta'$  is the **output vector** for word  $w_j$

## Computing the gradient on the hidden-output weights

- ▶ Use the familiar cost function:  
<https://eduassistpro.github.io/>

$$\mathcal{L} = -\sum_j t_j \log y_j - \log y_{j^*}$$

where  $j^*$  is the index of the target word

- ▶ Given  $y_j$  is the output and  $t_j$  is the target, compute the gradient on the output is:  
<https://eduassistpro.github.io/>

$$\frac{\partial \mathcal{L}}{\partial o_j} = y_j - t_j$$

$$\frac{\partial \mathcal{L}}{\partial \theta'_{ji}} = \frac{\partial \mathcal{L}}{\partial o_j} \frac{\partial o_j}{\partial \theta'_{ji}} = (y_j - t_j) z_i$$

- ▶ Update the hidden→output weights

$$\theta'_{ji} = \theta'_{ji} - \eta(y_j - t_j) z_i$$

## Updating input→hidden weights

- ▶ Compute the error

<https://eduassistpro.github.io/>

$$\frac{\partial \mathcal{L}}{\partial z_i} = \sum_{j=1}^V \frac{\partial \mathcal{L}}{\partial o_j} \frac{\partial o_j}{\partial z_i} = \sum_{j=1}^V (y_j - t_j) \theta'_{ji}$$

- ▶ Since

<https://eduassistpro.github.io/>

<https://eduassistpro.github.io/>

The derivative of  $\mathcal{L}$  on the input

[Add WeChat edu\\_assist\\_pro](https://eduassistpro.github.io/)

$$\frac{\partial \mathcal{L}}{\partial \theta_{ik}} = \frac{\partial \mathcal{L}}{\partial z_i} \frac{z_i}{\theta_{ik}} = \sum_{j=1}^V (y_j - t_j) \theta'_{ji} x_k$$

- ▶ Update the input→hidden weights

$$\theta_{ki} = \theta_{ki} - \eta \sum_{j=1}^V (y_j - t_j) \theta'_{ij} x_k$$

Gradient computation in matrix form

<https://eduassistpro.github.io/>

Assignment Project Exam Help

$$D_o = Y - \left( \begin{array}{c|c} \begin{matrix} 0.11039215 \\ 0.10404302 \\ 0.11262222 \end{matrix} & \begin{matrix} 39215 \\ 06362 \\ 37778 \end{matrix} \\ \hline \text{Add WeChat } \text{edu\_assist\_pro} & \begin{matrix} 0.09693485 \\ 0.13893957 \\ 0.13893957 \end{matrix} \end{array} \right)$$
$$= \left( \begin{array}{c|c} \begin{matrix} 0.05820895 \\ 0.12322834 \\ 0.063057 \\ 0.11039215 \\ 0.08016116 \end{matrix} & \begin{matrix} 5820895 \\ 2312834 \\ 0 \\ 0 \\ 0 \end{matrix} \\ \hline \text{Add WeChat } \text{edu\_assist\_pro} & \begin{matrix} 0.063057 \\ 0.11039215 \\ 0.08016116 \end{matrix} \end{array} \right)$$

Computing the errors at the hidden layer

<https://eduassistpro.github.io/>

$$\mathbf{D}_z = \mathbf{D}_o^\top \Theta' =$$

**Assignment Project Exam Help**

$$\begin{bmatrix} 0.110 & 0.106 & -0.887 & 0.097 & 0.139 & 0.0 \\ 0.3 & 0.4 & 0.6 & 0.7 & 0.1 & 0.6 \\ 0.5 & 0.2 & 0.7 & 0.2 & 0.6 & 0.3 \\ 0.2 & 0.6 & 0.3 & 0.6 & 0.5 & 0.6 \\ 0.3 & 0.1 & 0.1 & 0.3 & 0.1 & 0.1 \\ 0.2 & 0.4 & 0.8 & 0.3 & 0.2 & 0.1 \\ 0.3 & 0.4 & 0.6 & 0.3 & 0.5 & 0.1 \end{bmatrix} \times \begin{bmatrix} 1.10 & 0.080 \\ 0.115 & 0.157 \end{bmatrix}$$

**Assignment Project Exam Help  
Add WeChat edu\_assist\_pro**

Computing the updates to  $\Theta'$

<https://eduassistpro.github.io/>

Assignment Project Exam Help

$$\nabla_{\Theta'} = \mathbf{D}_{\Theta} \mathbf{z}^T =$$

$$\begin{bmatrix} 0.410 \\ 0.106 \\ -0.887 \\ 0.097 \\ 0.139 \\ 0.058 \\ 0.123 \\ 0.063 \\ 0.110 \\ 0.080 \end{bmatrix} \times \begin{bmatrix} 0.5 & 0.8 & 0.8 \end{bmatrix} = \begin{bmatrix} 0.10 \\ 0.053 \\ 10 \\ 0.085 \\ 0.085 \\ 0.0710 \\ 0.0718 \\ .111 \\ .0467 \\ .099 \end{bmatrix}$$

<https://eduassistpro.github.io/>

Computing the update to  $\Theta$

$$\nabla_{\Theta} = \mathbf{x}^T D = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} -0.11538456 & 0 & 388611 \\ 0.15687943 & 0 & -0.19388611 \\ 0. & 0. & 0. \\ 0. & 0. & 0. \\ 0. & 0. & 0. \\ 0. & 0. & 0. \\ 0. & 0. & 0. \end{bmatrix}$$

<https://eduassistpro.github.io/>

$$= \begin{bmatrix} 0. & 0. & 0. \\ 0. & 0. & 0. \\ -0.11538456 & 0.15687943 & -0.19388611 \\ 0. & 0. & 0. \\ 0. & 0. & 0. \\ 0. & 0. & 0. \\ 0. & 0. & 0. \\ 0. & 0. & 0. \end{bmatrix}$$

## CBOW for multiple context words

<https://eduassistpro.github.io/>

$$\mathbf{z} = \frac{1}{M} \Theta(\mathbf{x}_1 + \mathbf{x}_2 + \cdots + \mathbf{x}_M)$$

Assignment Project Exam Help

$$= \frac{1}{M} (\mathbf{v}_{w_1} + \mathbf{v}_{w_2} + \cdots + \mathbf{v}_{w_M})$$

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

where  $M$  is the number of words in the context,  $w_1, w_2, \dots, w_M$  are the words in the context. The loss function is

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$$E = -\log p(w_j | w_1, w_2, \dots, w_M)$$

$$= -o_{j^*} + \log \sum_{j'=1}^V \exp(o_{j'})$$

$$= -\mathbf{u}_{w_j}^\top \mathbf{z} + \log \sum_{j'=1}^V \exp(\mathbf{u}_{w_{j'}}^\top \mathbf{z})$$

Computing the hidden layer for multiple context words

<https://eduassistpro.github.io/>

$z = \Theta x = \begin{matrix} \text{Assignment} & \text{Project} & \text{Exam} & \text{Help} \end{matrix}$

$$\begin{matrix} \text{Assignment} & \text{Project} & \text{Exam} & \text{Help} \\ \text{Add} & \text{We} & \text{Chat} & \text{edu\_assist\_pro} \end{matrix} \times \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1.5 \\ 3.0 \\ 3.0 \end{bmatrix}$$

$\begin{bmatrix} 0.1 & 0.3 & 0.5 & 0.6 \\ 0.2 & 0.5 & 0.8 & 0.7 \\ 0.2 & 0.5 & 0.8 & 0.7 \end{bmatrix} \times \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1.5 \\ 3.0 \\ 3.0 \end{bmatrix}$

During backprop, update vectors for four words instead of just one.

Skip-gram: model

<https://eduassistpro.github.io/>

Assignment Project Exam Help

<https://eduassistpro.github.io/>

where  $w_{c,j}$  is the output layer,  $w_{O,c}$  is the actual input word,  $y_{c,j}$  is the output of the  $c$ -th panel of the output layer,  $o_{c,j}$  is the net of the  $j$ -th unit on the  $c$ -th panel of the output layer.

$$o_{c,j} = o_j = \mathbf{u}_{w_j} \cdot \mathbf{z}, \text{ for } c = 1, 2, \dots, C$$

## Skip-gram: loss function

<https://eduassistpro.github.io/>

Assignment Project Exam Help  
 $\mathcal{L} = -\log p(w_{O,1}, w_{O,2}, \dots, w_{O,C} | w_I)$

Assignment Project Exam Help  
Add WeChat  $\frac{\sum_{c=1}^C \exp(o_c)}{V}$  edu\_assist\_pro

<https://eduassistpro.github.io/>  
 $= -\log \frac{o_{c,j_c^*} + C}{V} \quad , \quad c=1, \dots, C$   
Add WeChat edu\_assist\_pro

where  $j_c^*$  is the index of the actual  $c$ -th output context word.

Combine the loss of  $C$  context words with multiplication. Note:  $o_j$  is the same for all  $C$  panels

## Skip-gram: updating the weights

- ▶ We take the derivative of  $\mathcal{L}$  with respect to every unit on every position  $j$ , and obtain

Assignment Project Exam Help

$$e_{c,j} = \frac{\partial \mathcal{L}}{\partial o_{c,j}} = y_{c,j} - t_{c,j}$$

Assignment Project Exam Help Add WeChat edu\_assist\_pro  
which is the prediction error of the unit.

- ▶ We define a loss function  $E_j$  as the sum of the prediction errors:

Add WeChat edu\_assist\_pro

$$\frac{\partial \mathcal{L}}{\partial \theta'_{ji}} = \sum_{c=1}^C \frac{\partial \mathcal{L}}{\partial o_{c,j}} \cdot \frac{\partial o_{c,j}}{\partial \theta'_{ji}} = E_j \cdot z_i$$

- ▶ Updating the hidden→output weight matrix:

$$\theta'_{ji} = \theta'_{ji} - \eta \cdot E_j \cdot z_i$$

- ▶ No change in how the input→hidden weights are updated.

## Optimizing computational efficiency

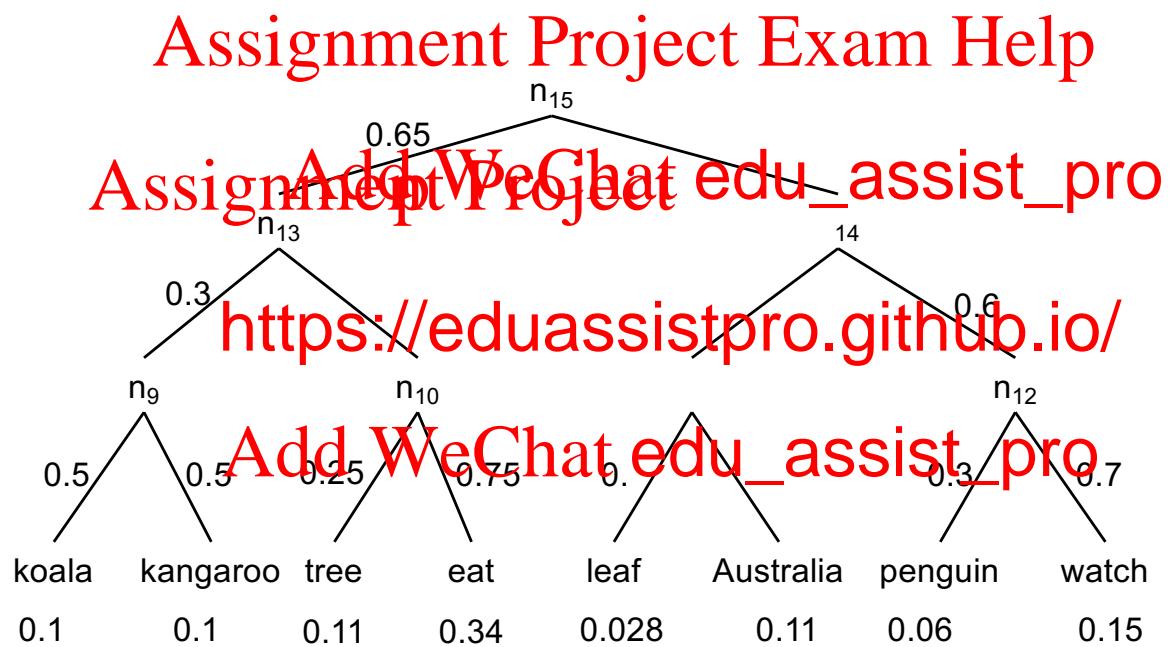
<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ Computing softmax at the output layer is expensive. It involves iterative operations over the entire vocabulary.
- ▶ Two methods for optimizing computational efficiency
  - ▶ **Hierarchical softmax:** Instead of computing the probability distribution over the entire vocabulary from  $|V|$  to  $\log |V|$ .
  - ▶ **Negative sampling:** Instead of updating all the words in the vocabulary, only sample a small number of words that are not actual context words in the training corpus.

## Hierarchical softmax

<https://eduassistpro.github.io/>



Computing the probabilities of the leaf nodes

<https://eduassistpro.github.io/>

Assignment Project Exam Help

$P(\text{"Kangaroo"}|z) = P_n(\text{Left}|z) \times P_n(\text{Right}|z)$

<https://eduassistpro.github.io/>

$P_n(\text{Right}|z) = 1$

[Add WeChat edu\\_assist\\_pro](https://eduassistpro.github.io/)

where  $\gamma_n$  is a vector from a set of new parameters that replace  $\Theta$

# Huffman Tree Building

<https://eduassistpro.github.io/>

A simple algorithm:

## Assignment Project Exam Help

- ▶ Prepare a collection of  $n$  initial Huffman trees as a single leaf node. Put them in a queue organized by weight (frequency).
- ▶ Remove the first two trees from the queue. Add WeChat edu\_assist\_pro
  - ▶ These two trees as children, and whose weight is the sum of the two children trees. Put this new tree in the queue.
- ▶ Repeat steps 2-3 until all of the partial Huffman trees have been combined into one.

## Negative sampling

<https://eduassistpro.github.io/>

- ▶ Computing softmax over the vocabulary is expensive. Another alternative is to approximate softmax by only updating a small sample of (context) words at a time.
- ▶ Given a pair of words  $(w, c)$ , let  $P$  be the probability of  $t$  and  $P(D = 1|w, c)$  come from the corpus.
- ▶ This probability can be modeled as a sigmoid.

$$P(D = 1|w, c) = \sigma(\mathbf{u}_w^\top \mathbf{v}_c) = \frac{1}{1 + e^{\mathbf{u}_w^\top \mathbf{v}_c}}$$

## New learning objective for negative sampling

- We need a new objective function to minimize the following loss function:

$$\mathcal{L} = - \sum_{w_j \in D} \log \sigma(o_{w_j}) - \log \sigma(-o_{w_j})$$

where  $D$  is a set of correct context - target word pairs, and  $D'$  is a set of incorrect word pairs.

- Note that we use multiple positive context words in the original skip-gram algorithm, there will be only one positive target word.
- The derivative of the loss function with respect to the output word will be:

$$\frac{\partial \mathcal{L}}{\partial o_{w_j}} = \sigma(o_{w_j}) - t_{w_j}$$

where  $t_{w_j} = 1$  if  $w_j \in D$  and  $t_{w_j} = 0$  if  $w_j \in D'$

Updates to the hidden output weights

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ Compute the gradient on the output weight  
~~Assignment Help WeChat edu\_assist\_pro~~

<https://eduassistpro.github.io/>

- ▶ When updating the output weights, only for words in the positive sample and negative samples should be updated:  
~~Add WeChat edu\_assist\_pro~~

Updates to the input hidden weights

<https://eduassistpro.github.io/>

- ▶ Computing the derivative of the loss function with respect to the hidden layer

~~Assignment Project Exam Help~~  
~~Assignment Add WeChat edu\_assist\_pro~~  
 $\frac{\partial L}{\partial w} = \sigma(o_w - t_w) w_j$

<https://eduassistpro.github.io/>

- ▶ In the CBOW algorithm, the weights for a words will be updated. In the Skip-gram model for the target word will be updated.

$$\mathbf{v}_{w_i} = \mathbf{v}_{w_i} - \eta(\sigma(o_{w_j}) - t_{w_j}) \mathbf{u}_{w_j} x_i$$

## How to pick the negative samples?

- ▶ If we just randomly pick any given word  $w_i$  getting picked is:

**Assignment Project Exam Help**  
 $p(w_i) = \frac{\text{freq}(w_i)}{\sum_{j=0}^V \text{freq}(w_j)}$

More frequent words may not be ideal

- ▶ Adjust the formula to give less frequent words a chance to get picked:

$$p(w_i) = \frac{\text{freq}(w_i)^{\frac{3}{4}}}{\sum_{j=0}^V \text{freq}(w_j)^{\frac{3}{4}}}$$

- ▶ Generate a sequence of words using the adjusted probability, and randomly pick  $n_{D'}$  words

Use of embeddings: word and short document similarity

<https://eduassistpro.github.io/>

- ▶ Word embeddings can be used to compute word similarity with cosine similarity

~~Assignment Project Exam Help  
Assignment Help WeChat Add edu\_assist\_pro~~

- ▶ How accurate word embeddings can also be used to compute document similarity
- ▶ They can also be used to compute the similarity between documents

$$sim_{doc}(D_1, D_2) = \sum_{i=1}^m \sum_{j=1}^n \cos(\mathbf{w}_i^1, \mathbf{w}_j^2)$$

## Use of embeddings: word analogy

<https://eduassistpro.github.io/>

- ▶ What's even more impressive is that they can be used to compute word analogy

$$\text{analog}(m : w \rightarrow k : ?) = \underset{\mathbf{v} \in V}{\operatorname{argmax}} \frac{\cos(\mathbf{v}, \mathbf{k}) - \cos(\mathbf{v}, \mathbf{m})}{\cos(\mathbf{v}, \mathbf{m}) + \epsilon}$$

$$\text{analog}(m : w \rightarrow k : ?) = \underset{\mathbf{v} \in V \setminus \{m, k, w\}}{\operatorname{argmax}} \frac{\cos(\mathbf{v}, \mathbf{k}) - \cos(\mathbf{v}, \mathbf{m})}{\cos(\mathbf{v}, \mathbf{m}) + \epsilon}$$

$$\text{analog}(m : w \rightarrow k : ?) = \underset{\mathbf{v} \in V \setminus \{m, k, w\}}{\operatorname{argmax}} \frac{\cos(\mathbf{v}, \mathbf{k}) \cos(\mathbf{v}, \mathbf{w})}{\cos(\mathbf{v}, \mathbf{m}) + \epsilon}$$

Word analogy

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help Project WeChat Add edu\_assist\_pro

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

## Use of word embeddings

<https://eduassistpro.github.io/>

- ▶ Computing word similarities is not a “real” problem in the eyes of many
- ▶ The most important use word embeddings is predict the outcome of tasks that have real-world applications
- ▶ Many follow-up embeddings:
  - ▶ word2vec: <http://vectors.nlpl.eu>
  - ▶ fasttext: <https://fasttext.cc/docs.html>
  - ▶ GLOVE: <https://nlp.stanford.edu>
- ▶ Contextualized word embeddings:
  - ▶ <https://allennlp.org/elmo>
  - ▶ BERT: <https://github.com/google-research/bert>
  - ▶ Roberta: [https://pytorch.org/hub/pytorch\\_fairseq\\_roberta](https://pytorch.org/hub/pytorch_fairseq_roberta)

Embeddings in Pytorch

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

## Commonly used neural architectures

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ One important aspect of neural network mo  
out the right representation for a problem
- ▶ Commonly u  
  - ▶ Variant <https://eduassistpro.github.io/> which have  
been inst
  - ▶ Convolutional Networks (CNN)  
in image processing and some NLP pro  
classification)

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)  
Convolut ation

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

## Convolutional Networks

<https://eduassistpro.github.io/>

- ▶ A convolutional network is designed to identify indicative local indicators in a large structure, and combine them to produce a fixed size vector representation of the st pooling function, capturing the local aspect informative o
- ▶ A convolutional feedforward network is.
- ▶ It has been tremendously successful in im computer vision), where the input is the raw pixels of an image
- ▶ In NLP, it has been shown to be effective in sentence classification, etc.

Why it has been so effective in image processing

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

Image pixels

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

Four operations in a convolutional network

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ Convolution <https://eduassistpro.github.io/> Add WeChat edu\_assist\_pro
- ▶ Non-linear ac <https://eduassistpro.github.io/>
- ▶ Pooling or sub <https://eduassistpro.github.io/>
- ▶ Classification with a fully connected layer <https://eduassistpro.github.io/> Add WeChat edu\_assist\_pro

## Image convolution

<https://eduassistpro.github.io/>

Assignment Project Exam Help

$$\begin{matrix} & & \mathbf{x} \\ \circ & \text{Assign} & \text{Add WeChat edu_assist_pro} \\ \begin{matrix} 2 & 2 \\ 0 & 2 \end{matrix} & = c & \begin{matrix} \text{Input image} \\ \text{Output} \\ \text{or "feature map"} \end{matrix} \end{matrix}$$

The diagram illustrates the process of image convolution. An input image  $x$  is shown as a 4x4 grid of green squares. A kernel, represented by a 2x2 matrix with values [2, 2; 0, 2], is applied to the input. The result is an output feature map  $c$ , which is a 2x2 grid of orange squares. The output value in the bottom-right square of the feature map is highlighted in red.

Kernel size = (2,2), strides = 1

## Image convolution

<https://eduassistpro.github.io/>

Assignment Project Exam Help  
 $\text{conv}(x, U)$



Add WeChat edu\_assist\_pro  
<https://eduassistpro.github.io/>

Kernel size = (2,2), strides = 1

## Image convolution

<https://eduassistpro.github.io/>

Assignment Project Exam Help  
 $\text{conv}(x, U)$



Add WeChat edu\_assist\_pro  
<https://eduassistpro.github.io/>

Kernel size = (2,2), strides = 1

## Image convolution

<https://eduassistpro.github.io/>

Assignment Project Exam Help  
 $\text{conv}(x, U)$



Add WeChat edu\_assist\_pro  
<https://eduassistpro.github.io/>

Kernel size = (2,2), strides = 1

## Image convolution

<https://eduassistpro.github.io/>

Assignment Project Exam Help  
 $\text{conv}(x, U)$



Add WeChat edu\_assist\_pro  
<https://eduassistpro.github.io/>

Kernel size = (2,2), strides = 1

Forward computation

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

$o_{11} = x_{11}u_{11} + x_{12}u_{12} + x_{13}u_{13}$   
 $o_{12} = x_{21}u_{11} + x_{22}u_{12} + x_{23}u_{13}$   
 $o_{21} = x_{11}u_{21} + x_{12}u_{22} + x_{13}u_{23}$   
 $o_{22} = x_{21}u_{21} + x_{22}u_{22} + x_{23}u_{23}$

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

## ReLU

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ Nonlinear transformation with ReLU  
[Assignment Help WeChat edu\\_assist\\_pro](https://eduassistpro.github.io/)  
 $x(0, input)$
- ▶ As we know, ReLU does not change the dimension of the feature vector  $x$ .
- ▶ ReLU replaces all negative pixel values in the image  $h$  by 0.

Image ReLU

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

## ReLU activation and Max pooling

<https://eduassistpro.github.io/>

- ▶ ReLU activation is a component-wise function and does not change the dimension of the input

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

- ▶ Max pooling does change the dimension  
specify the pool size and strides

$$[2] = \text{Max} \left( \begin{bmatrix} 2 & 2 \\ 0 & 2 \end{bmatrix} \right)$$

pool size = (2, 2), strides = 2

## Training a CNN

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ Loss functions: Cross entropy loss, square error
- ▶ What are the parameters of a CNN?
  - ▶ The filter network
- ▶ Computing the gradient for the convolutional layer from a feedforward neural network...

## Computing the gradient on $U$

$$\frac{\partial E}{\partial u_{11}} = \frac{\partial E}{\partial o_{11}} \frac{\partial o_{11}}{\partial u_{11}} + \frac{\partial E}{\partial o_{12}} \frac{\partial o_{12}}{\partial u_{11}} + \frac{\partial E}{\partial o_{21}} \frac{\partial o_{21}}{\partial u_{11}} + \frac{\partial E}{\partial o_{22}} \frac{\partial o_{22}}{\partial u_{11}}$$

**Assignment Project Exam Help**  
<https://eduassistpro.github.io/>  
**Add WeChat** [edu\\_assist\\_pro](https://edu_assist_pro)

$$\frac{\partial E}{\partial u_{12}} = \frac{\partial E}{\partial o_{11}} \frac{\partial o_{11}}{\partial u_{12}} + \frac{\partial E}{\partial o_{12}} \frac{\partial o_{12}}{\partial u_{12}} + \frac{\partial E}{\partial o_{21}} \frac{\partial o_{21}}{\partial u_{12}} + \frac{\partial E}{\partial o_{22}} \frac{\partial o_{22}}{\partial u_{12}}$$

$$\frac{\partial E}{\partial u_{21}} = \frac{\partial E}{\partial o_{11}} \frac{\partial o_{11}}{\partial u_{21}} + \frac{\partial E}{\partial o_{12}} \frac{\partial o_{12}}{\partial u_{21}} + \frac{\partial E}{\partial o_{21}} \frac{\partial o_{21}}{\partial u_{21}} + \frac{\partial E}{\partial o_{22}} \frac{\partial o_{22}}{\partial u_{21}}$$

**Assignment Project Exam Help**  
<https://eduassistpro.github.io/>  
**Add WeChat** [edu\\_assist\\_pro](https://edu_assist_pro)

$$\frac{\partial E}{\partial u_{22}} = \frac{\partial E}{\partial o_{11}} \frac{\partial o_{11}}{\partial u_{22}} + \frac{\partial E}{\partial o_{12}} \frac{\partial o_{12}}{\partial u_{22}} + \frac{\partial E}{\partial o_{21}} \frac{\partial o_{21}}{\partial u_{22}} + \frac{\partial E}{\partial o_{22}} \frac{\partial o_{22}}{\partial u_{22}}$$

$$\frac{\partial E}{\partial o_{11}} x_{11} + \frac{\partial E}{\partial o_{12}} x_{12} + \frac{\partial E}{\partial o_{21}} x_{21} + \frac{\partial E}{\partial o_{22}} x_{22}$$

$$\frac{\partial E}{\partial o_{11}} x_{21} + \frac{\partial E}{\partial o_{12}} x_{22} + \frac{\partial E}{\partial o_{21}} x_{31} + \frac{\partial E}{\partial o_{22}} x_{32}$$

$$\frac{\partial E}{\partial o_{11}} \frac{\partial o_{11}}{\partial u_{22}} + \frac{\partial E}{\partial o_{12}} \frac{\partial o_{12}}{\partial u_{22}} + \frac{\partial E}{\partial o_{21}} \frac{\partial o_{21}}{\partial u_{22}} + \frac{\partial E}{\partial o_{22}} \frac{\partial o_{22}}{\partial u_{22}}$$

$$\frac{\partial E}{\partial o_{22}} x_{11} + \frac{\partial E}{\partial o_{12}} x_{23} + \frac{\partial E}{\partial o_{21}} x_{32} + \frac{\partial E}{\partial o_{22}} x_{33}$$

Summing up errors from all outputs that the filter component has contributed to.

## Reverse Convolution

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

The computation of the gradient on the filter can be via  
a reverse convolution:

$$\left[ \begin{array}{cc} \frac{\partial E}{\partial u_{11}} & \frac{\partial E}{\partial u_{12}} \\ \frac{\partial E}{\partial u_{21}} & \frac{\partial E}{\partial u_{22}} \end{array} \right] \xrightarrow[Add\ WeChat\ edu\_assist\_pro]{21} \left( \begin{array}{cc} \frac{\partial E}{\partial o_{11}} & \frac{\partial E}{\partial o_{12}} \\ \frac{\partial E}{\partial o_{21}} & \frac{\partial E}{\partial o_{22}} \end{array} \right)$$

Computing the gradient on  $X$  (if this is not the input layer)

$$\frac{\partial E}{\partial x_{11}} = \frac{\partial E}{\partial o_{11}} u_{11} + \frac{\partial E}{\partial o_{12}} u_{12} + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} 0$$

$$\frac{\partial E}{\partial x_{12}} = \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} u_{12} + \frac{\partial E}{\partial o_{21}} u_{11} + \frac{\partial E}{\partial o_{22}} 0$$

$$\frac{\partial E}{\partial x_{13}} = \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} u_{12} + \frac{\partial E}{\partial o_{22}} u_{11}$$

$$\frac{\partial E}{\partial x_{21}} = \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} 0$$

$$\frac{\partial E}{\partial x_{22}} = \frac{\partial E}{\partial o_{11}} u_{22} + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} u_{11}$$

$$\frac{\partial E}{\partial x_{23}} = \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} u_{22} + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} u_{22}$$

$$\frac{\partial E}{\partial x_{31}} = \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} u_{21} + \frac{\partial E}{\partial o_{22}} 0$$

$$\frac{\partial E}{\partial x_{32}} = \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} u_{22} + \frac{\partial E}{\partial o_{22}} u_{21}$$

$$\frac{\partial E}{\partial x_{12}} = \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} u_{22}$$

Full convolution

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)  
Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

Gradient on  $X$  if it is not the inputs

$u_{22}$	$u_{21}$				$u_{22}$	$u_{21}$
$u_{12}$	$\delta o_{11}u_{11}$	$\delta o_{12}u_{12}$	$\delta o_{21}u_{21}$	$\delta o_{22}u_{22}$	$\delta o_{11}$	$\delta o_{12}u_{12}$
	$\delta o_{21}$	$\delta o_{22}$	$\delta o_{21}u_{12}$	$\delta o_{22}$		$u_{11}$

<https://eduassistpro.github.io/>  
 Assignment Project Exam Help  
 Add WeChat edu\_assist\_pro

$u_{22}$	$\delta o_{11}u_{21}$	$\delta o_{12}u_{22}$			$x_{12}$	$x_{21}$
$u_{12}$	$\delta o_{21}u_{11}$	$\delta o_{22}$	$\delta o_{21}u_{12}$	$\delta o_{22}$	$x_{22}$	$u_{11}$

<https://eduassistpro.github.io/>  
 Add WeChat edu\_assist\_pro

	$\delta o_{11}$	$\delta o_{12}$			$\delta o_{11}$	$\delta o_{12}$	
$u_{22}$	$\delta o_{21}u_{21}$	$\delta o_{22}$	$\delta o_{21}u_{22}$	$\delta o_{22}u_{21}$	$\delta o_{21}$	$\delta o_{22}u_{22}$	$u_{21}$
$u_{12}$	$u_{11}$		$u_{12}$	$u_{11}$	$u_{12}$	$u_{11}$	

## Sample code of 2D convolution with Keras

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

```
model = Sequential()
model.add(Conv2D(32, kernel_size=(1, 1),
                 activation='relu',
                 input_shape=inp))
model.add(MaxPool2D(pool_size=(2, 2)))
model.add(Conv2D(64, (3, 3), activation='relu'))
model.add(MaxPool2D(pool_size=(2, 2)))
model.add(Flatten())
model.add(Dense(1000, activation='relu'))
model.add(Dense(num_classes, activation='softmax'))
```

Tutorials on how to train CNNs can be found on the Tensorflow website:

<https://www.tensorflow.org/tutorials/images/cnn>

Why convolutional networks for NLP?

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ Even though bag-of-word models are simple, they can perform well on some text classification tasks, they don't account for the context where multiple words appear together, "not interesting".  
<https://eduassistpro.github.io/>
- ▶ The analogy with image processing is if the words are features, then each word is a separate feature. (The analogy might be clearer if we add WeChat to the analogy.)  
<https://eduassistpro.github.io/>

Input to a convolutional network in a text classification task

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

The input to a convolutional network can be pretrained word embeddings (e.g., the weight matrix produced by GLOVE<sup>1</sup>) and the input sentence:

<https://eduassistpro.github.io/>  
 $\mathbf{x}^{(0)} = \Theta [e_{w_1}, e_{w_2}, \dots, e_{w_m}, \mathbf{e}_{w_m}]^T \in \mathbb{R}^{K_e \times M}$

where  $e_{w_m}$  is a column vector of zeros, with a 1 at position  $w_m$ ,  $K_e$  is the size of embeddings

---

<sup>1</sup><https://nlp.stanford.edu/projects/glove>

## Alternative text representations

<https://eduassistpro.github.io/>

- ▶ Alternatively, a text can be represented as a sequence of word tokens  $w_1, w_2, w_3, \dots, w_M$ . This is useful for models such as Convolutional Neural Networks (cNNs), which processes text as a sequence.
- ▶ Each word token  $w$  has a vector representation  $e_{w_m}$ , with dimension  $V$ , represented by the horizontal concatenation of vectors:  $\mathbf{W} = [e_{w_1}, e_{w_2}, \dots, e_{w_M}]$ .
- ▶ To show that this is equivalent to the bag-of-words model, we can recover the word count from the matrix-vector product  $\mathbf{W}[1, 1, \dots, 1]^\top \in R^V$ .

“Convolve” the input with a set of filters

<https://eduassistpro.github.io/>

- ▶ A *filter* is a weight matrix of dimension  $\mathbf{C}^{(k)} \in \mathbb{R}^{K_e \times h}$  where  $\mathbf{C}^{(k)}$  is the  $k$ th filter. Note the first dimension of the filter is the same as the size of the embedding.
  - ▶ Unlike image processing, the filter doesn't have full width
- ▶ To merge adjacent words by sliding a set of filters across it:

<https://eduassistpro.github.io/>  
 $\mathbf{x}^{(1)} = f(\mathbf{b})$

where  $f$  is an activation function (e.g., tanh, ReLU),  $b$  is a vector of bias terms, and  $*$  is the convolution operator.

## Computing the convolution

- ▶ At each position  $\text{https://eduassistpro.github.io/}$  compute the element-wise product of the  $k$ th filter and the sequence of words of window size  $h$  (think of it as an n-gram of length  $h$ ) starting at  $m$  and take it  $(k) \odot \mathbf{x}_{m:m+h-1}^{(0)}$
- ▶ The value of the  $m$ th position in the output  $\text{Add WeChat edu_assist_pro}$  computed as:

$$\text{https://eduassistpro.github.io/}$$
$$x_m^{(1)} = f \left( \sum_{k'=1}^{K_f} h_k \odot \sum_{n=1}^{K_e} \mathbf{x}_{m+n-1}^{(0)} \right)$$

- ▶ When we finish the convolution step, if we have  $K_f$  filters of dimension  $\mathbb{R}^{K_e \times h}$ , then  $\mathbf{X}^{(1)} \in \mathbb{R}^{K_f \times M-h+1}$
- ▶ In practice, filters of different sizes are often used to capture ngrams of different lengths, so  $\mathbf{X}^{(1)}$  will be  $K_f$  vectors of variable lengths, and we can write the size of each vector of  $h_k$

Convolution step when processing text

<https://eduassistpro.github.io/>

Assignment Project Exam Help  
 $\mathbf{X}^{(0)}$

Conv(  
Assignment Help Project Exam Help  
https://eduassistpro.github.io/  
Add WeChat edu\_assist\_pro  
Input text sequence  
))

= 

2.9	3.1	3.4	?	?	?
-----	-----	-----	---	---	---

$\mathbf{X}^{(1)}$

## Padding

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ To deal with the beginning and end of the input, the matrix is often padded with  $h - 1$  columns at the beginning and  $c - 1$  columns at the end. This is known as **narrow convolution**.
- ▶ If no padding is applied, the output layer will be  $h - 1$  units smaller than the input. This is known as **narrow convolution**.

## Pooling

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ After  $D$  convolutional layers, assuming filters have identical lengths, we have a representation of the document  $\mathbf{x}^{(D)} \in \mathbb{R}^{1 \times 2M}$ .
- ▶ It is very likely that so we need to trim them before feeding them to a feedforward network for classification.
- ▶ This can be done by **pooling** across times (over the sequence of words)

## Prediction and training with CNN

<https://eduassistpro.github.io/>

- ▶ The CNN needs to be fed into a feedforward network to make a prediction  $\hat{y}$  and compute the loss  $\ell^{(i)}$  in training.
- ▶ Parameters of a CNN includes the weight matrix of the feedforward network and the filters, as well as the biases.
- ▶ The parameter <https://eduassistpro.github.io/>, which may involve computing the gradient for the function. [Add WeChat edu\\_assist\\_pro](#)

$$\frac{\partial z_k}{\partial x_{k,m}^{(D)}} = \begin{cases} 1, & x_{k,m}^{(D)} = \max(x_{k,1}^{(D)}, x_{k,2}^{(D)}, \dots, x_{k,M}^{(D)}) \\ 0, & \text{Otherwise} \end{cases}$$

## Different pooling methods

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ Max pooling

Add WeChat edu\_assist\_pro  
$$z_k = \max_{k,1}^{(D), k,M} x_{k,1}^{(D)}, x_{k,2}^{(D)}$$

- ▶ Average pool

Add WeChat edu\_assist\_pro

$$z_k = \frac{1}{M} \sum_{m=1}^{k,m} x_m$$

A graphic representation of a CNN

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

Figure 1: Caption

## Sample code of convolution with Keras

<https://eduassistpro.github.io/>

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, \\\
    Embedding, Flatten, Conv1D, GlobalMaxPooling1D

model = Sequential()
model.add(Embedding(input_dim=n_input_words,
                    g_dim,
                    n))
model.add(Conv1D(num_filters, kernel_size=ker_size, \\
activation='relu'))
model.add(GlobalMaxPooling1D())
model.add(Dense(50, activation='relu'))
model.add(Dense(3, activation='softmax'))
model.compile(optimizer='adam',
              loss='categorical_crossentropy',
              metrics=['accuracy'])
```

# An Interim Summary of Supervised Learning Methods

<https://eduassistpro.github.io/>

- ▶ In all the linear and no  
far, we assume we have labeled training data where we can  
perform supervised learning:
  - ▶ a **training set** where you get observa  
bels  $y$ ;
  - ▶ a **test set** where you only get observat  
[Add WeChat edu\\_assist\\_pro](https://eduassistpro.github.io/)
- ▶ A summary of t  
discussed so fa  
<https://eduassistpro.github.io/>
  - ▶ Linear m  
Support Vector Machines
  - ▶ Non-linear models: feed-forward n  
erceptional Networks
  - ▶ Sparse vs dense feature representations as input to classifiers
- ▶ Given sufficient amounts of high-quality data, supervised  
learning methods tend to produce more accurate classifiers  
than alternative learning paradigms

NLP problems that can be formulated as simple text classifications

<https://eduassistpro.github.io/>

- ▶ An NLP problem can be formulated as a simple text classification if there is no inter-dependency between the classification instances
  - ▶ Word sense
  - ▶ Sentiment
  - ▶ Genre classification
  - ▶ Others
- ▶ NLP problems that cannot be formulated as simple text classifications (or you can, but the results won't be optimal)
  - ▶ Sequence labeling problems such as POS tagging, Named Entity Recognition
  - ▶ Structured prediction problems such as syntactic parsing

## Beyond Supervised Learning

<https://eduassistpro.github.io/>

There are other learning scenarios where labeled training sets are available to various degree or not available at all

- ▶ When there is no labeled target data, we'll have to **unsupervised learning**, variants of the EM algorithm
- ▶ When there is a small amount of labeled data to try **semi-supervised learning**
- ▶ When there is a lot of labeled data in one domain but there is only a small amount of labeled data in the target domain, we might try **domain adaptation**

## K-Means clustering algorithm

<https://eduassistpro.github.io/>

---

K-means clustering algorithm

## Assignment Project Exam Help

---

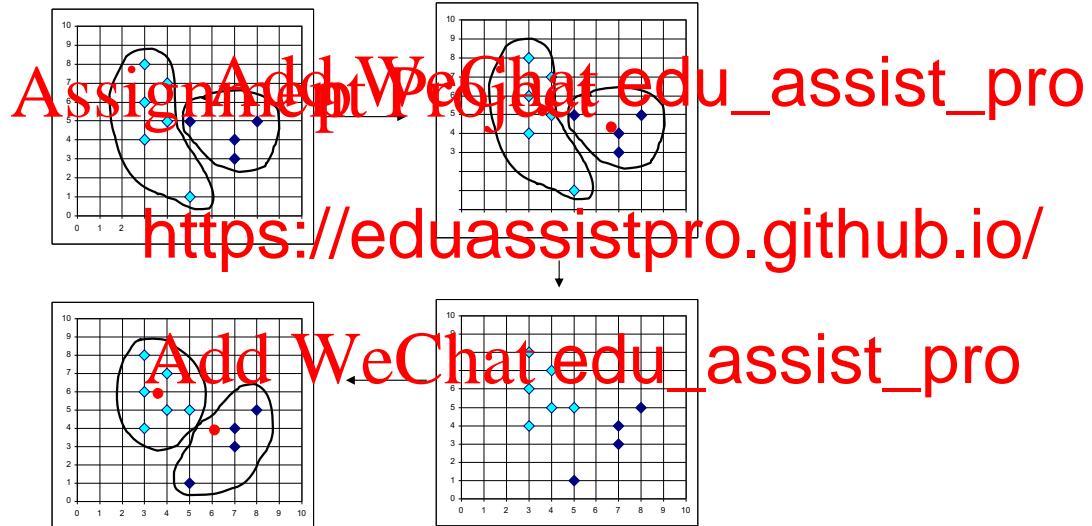
```
procedure K-MEANS( $x_{1:N}$ ,  $K$ )
    for  $i \in 1 \dots N$  do
         $z^{(i)} \leftarrow \text{RANDOMINT}(1, K)$   $\triangleright$  Assign  $x_i$  to cluster  $k$ 
    repeat
        for  $k$  https://eduassistpro.github.io/ do
             $v_k \leftarrow \frac{1}{\sum_{i=1}^N \delta(z^{(i)}=k)} \sum_i^N x^{(i)}$   $\triangleright$  Compute cluster centers
        Add WeChat edu_assist_pro
        for  $i \in 1 \dots N$  do
             $z^{(i)} \leftarrow \operatorname{argmin}_K \|x^{(i)} - v_k\|^2$ 
    until Converged
    return  $z$ 
```

---

## K-Means training

<https://eduassistpro.github.io/>

Assignment Project Exam Help



- ▶ K-means clustering is non-parametric and has no parameters to update
- ▶ The number of clusters need to be pre-specified before the training process starts

## Semi-supervised learning

<https://eduassistpro.github.io/>

- ▶ Initialize parameters with supervised learning and then apply unsupervised learning (such as the EM algorithm)
- ▶ Multi-view learning: Co-training
  - ▶ divide features into multiple views, and train each view
  - ▶ Each class is trained on all instances, using only the features available in its view. Predictions are then used as ground truth for classifiers associated with the other views.
- ▶ Named entity example: named entity view and local context view
- ▶ Word sense disambiguation: local context view and global context view

## Domain adaptation

Supervised domain adaptation (Daumé)

<https://eduassistpro.github.io/>

- ▶ Creates copies of each feature: one for each domain and one for the cross-domain setting

$f(x, y, d) = \{((boring, NEG, Movie):1,$

$(f1, EG, *):1,$   
 $(t, day-old, NEG, *):1,$

$\dots\}$  Add WeChat edu\_assist\_pro

where  $d$  is the domain.

- ▶ Let the learning algorithm allocate weights between domain specific features and cross-domain features: for words that facilitate prediction in both domains, the learner will use cross-domain features. For words that are only relevant to a particular domain, domain-specific features will be used.

## Other learning paradigms

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ Active learning: A learning that is often used to reduce the number of instances that have to be annotated to produce the same level of accuracy.  
<https://eduassistpro.github.io/>
- ▶ Distant supervision: generate some external resource such as a dictionary. For example, generate named entity annotation with  
<https://eduassistpro.github.io/>
- ▶ Multitask learning: The learning induces a representation that can be used to solve multiple tasks (learning POS tagging with syntactic parsing)