

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help
Convolutional Networks for text classification

<https://powcoder.com>

Add WeChat powcoder

Convolutional Networks

<https://powcoder.com>

- ▶ 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 structure with a pooling function, capturing the local aspects that are most informative of the prediction task at hand.
- ▶ A convolutional network is not fully connected as a feedforward network is.
- ▶ It has been tremendously successful in image procession (or 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://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

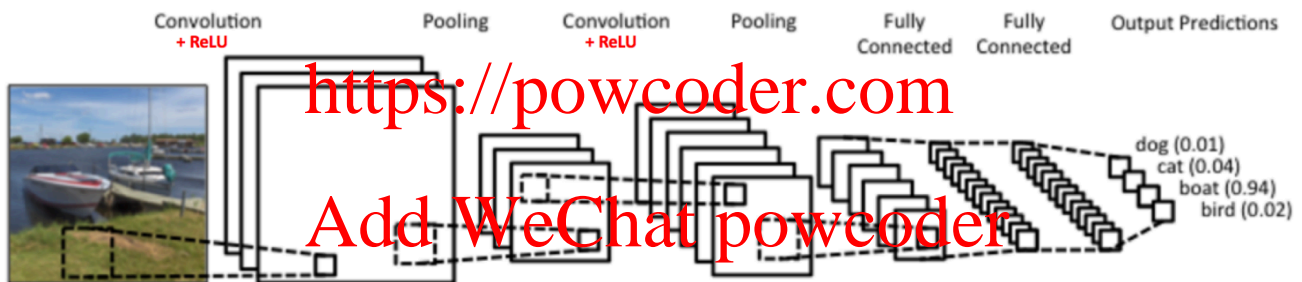


Image pixels

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project PowerUp

<https://powcoder.com>

8 0 121 162 75 235 254 253 126 6 0 10 14 6 3 0 0 0
3 75 242 255 141 66 255 245 189 7 8 0 0 5 0 0 0

[illegible]

Four operations in a convolutional network

<https://powcoder.com>

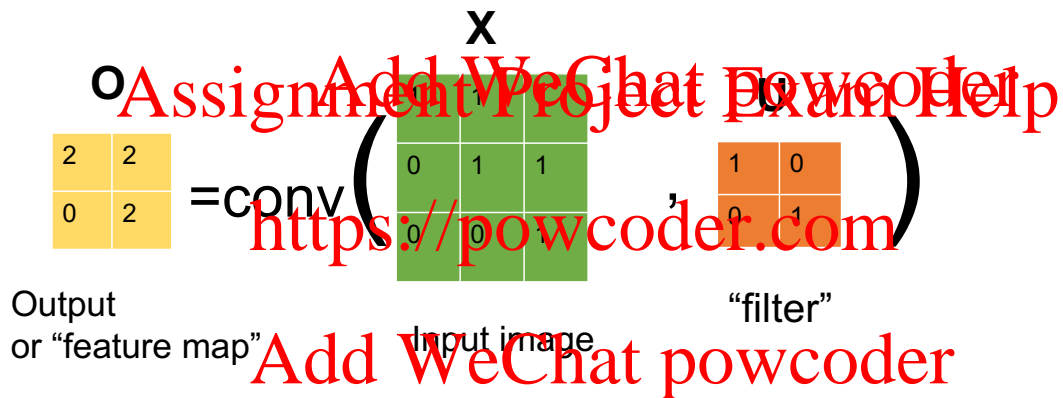
Assignment Project Exam Help

- ▶ Convolution
- ▶ Non-linear activation (ReLU)
- ▶ Pooling or subsampling (Max)
- ▶ Classification with a fully-connected layer

Image convolution

<https://powcoder.com>

Assignment Project Exam Help



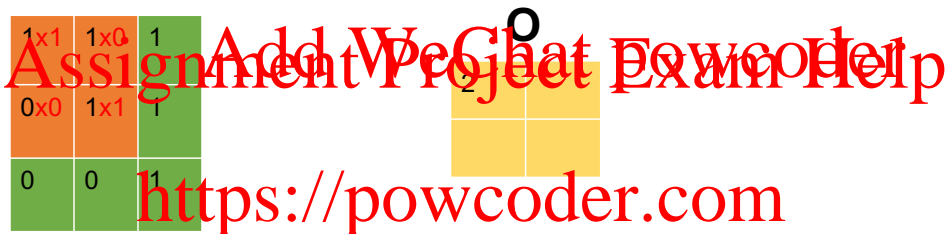
Kernel size = (2,2), strides = 1

Image convolution

<https://powcoder.com>

conv(X, U)

Assignment Project Exam Help



<https://powcoder.com>

Add WeChat powcoder

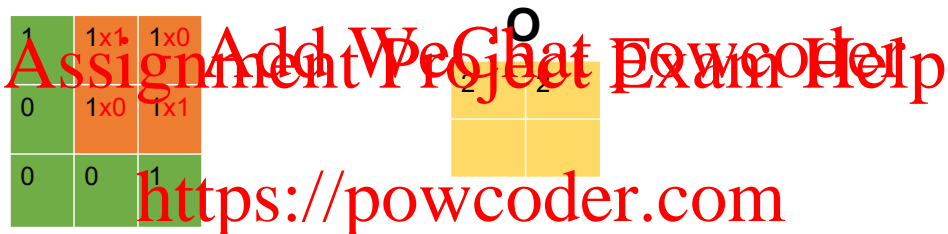
Kernel size = (2,2), strides = 1

Image convolution

<https://powcoder.com>

conv(X, U)

Assignment Project Exam Help



Add WeChat powcoder

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

Image convolution

<https://powcoder.com>

conv(X, U)

1	1	1	0
0x1	1x0	1	2
0x0	0x1	1	2
			0

<https://powcoder.com>

Add WeChat powcoder

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

Image convolution

<https://powcoder.com>

conv(X, U)

1	1	1	0
0	1x1	1x0	2
0	0x0	1x1	2
			0
			2
			0
			2

<https://powcoder.com>

Add WeChat powcoder

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

Forward computation

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

$$O_{11} = x_{11}u_{11} + x_{12}u_{12} + x_{21}u_{21} + x_{22}u_{22}$$

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

$$O_{21} = x_{21}u_{11} + x_{22}u_{12} + x_{31}u_{21} + x_{32}u_{22}$$

$$O_{22} = x_{22}u_{11} + x_{23}u_{12} + x_{32}u_{21} + x_{33}u_{22}$$

<https://powcoder.com>

Add WeChat powcoder

ReLU

<https://powcoder.com>

Assignment Project Exam Help

- ▶ Nonlinear transformation with ReLU

Assignment Project Exam Help

$$\text{Output} = \text{ReLU}(\text{input}) = \max(0, \text{input})$$

- ▶ As we know, ReLU is an element-wise transformation that does not change the dimension of the feature map
- ▶ ReLU replaces all negative pixel values in the featuremap with 0

<https://powcoder.com>

Add WeChat powcoder

Image ReLU

<https://powcoder.com>

Assignment Project Exam Help



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

ReLU activation and Max pooling

<https://powcoder.com>

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

Assignment Project Exam Help

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

<https://powcoder.com>

- ▶ Max pooling does change the dimension of the input. Need to specify the pool size and strides.

Add WeChat powcoder

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

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

Training a CNN

<https://powcoder.com>

Assignment Project Exam Help

- ▶ Loss functions: Cross-entropy loss, squared error loss
- ▶ What are the parameters of a CNN?
 - ▶ The filters (kernels), weight matrices for the feedforward network on top of the convolution and pooling layers, biases
- ▶ Computing the gradient for the convolution layers is different from a feedforward neural network...

Computing the gradient on U

$$\begin{aligned}
 \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}} \\
 &= \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 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 o_{11}} x_{12} + \frac{\partial E}{\partial o_{12}} x_{13} + \frac{\partial E}{\partial o_{21}} x_{22} + \frac{\partial E}{\partial o_{22}} x_{23} \\
 \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}} \\
 &= \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 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_{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}
 \end{aligned}$$

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

Reverse Convolution

<https://powcoder.com>

Assignment Project Exam Help

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

$$\begin{bmatrix} \frac{\partial E}{\partial u_{11}} & \frac{\partial E}{\partial u_{12}} \\ \frac{\partial E}{\partial u_{21}} & \frac{\partial E}{\partial u_{22}} \end{bmatrix} = \text{conv} \left(\begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix}, \begin{bmatrix} \frac{\partial E}{\partial o_{11}} & \frac{\partial E}{\partial o_{12}} \\ \frac{\partial E}{\partial o_{21}} & \frac{\partial E}{\partial o_{22}} \end{bmatrix} \right)$$

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

$$\begin{aligned}\frac{\partial E}{\partial x_{11}} &= \frac{\partial E}{\partial o_{11}} u_{11} + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} 0 \\ \frac{\partial E}{\partial x_{12}} &= \frac{\partial E}{\partial o_{11}} u_{12} + \frac{\partial E}{\partial o_{12}} u_{11} + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} 0 \\ \frac{\partial E}{\partial x_{13}} &= \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} u_{12} + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} 0 \\ \frac{\partial E}{\partial x_{21}} &= \frac{\partial E}{\partial o_{11}} u_{21} + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} u_{11} + \frac{\partial E}{\partial o_{22}} 0 \\ \frac{\partial E}{\partial x_{22}} &= \frac{\partial E}{\partial o_{11}} u_{22} + \frac{\partial E}{\partial o_{12}} u_{21} + \frac{\partial E}{\partial o_{21}} u_{12} + \frac{\partial E}{\partial o_{22}} u_{11} \\ \frac{\partial E}{\partial x_{23}} &= \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} u_{22} + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} u_{12} \\ \frac{\partial E}{\partial x_{31}} &= \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} u_{21} + \frac{\partial E}{\partial o_{22}} 0 \\ \frac{\partial E}{\partial x_{32}} &= \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} u_{22} + \frac{\partial E}{\partial o_{22}} u_{21} \\ \frac{\partial E}{\partial x_{12}} &= \frac{\partial E}{\partial o_{11}} 0 + \frac{\partial E}{\partial o_{12}} 0 + \frac{\partial E}{\partial o_{21}} 0 + \frac{\partial E}{\partial o_{22}} u_{22}\end{aligned}$$

Full convolution

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

$$\begin{bmatrix} \frac{\partial E}{\partial x_{11}} & \frac{\partial E}{\partial x_{12}} & \frac{\partial E}{\partial x_{13}} \\ \frac{\partial E}{\partial x_{21}} & \frac{\partial E}{\partial x_{22}} & \frac{\partial E}{\partial x_{23}} \\ \frac{\partial E}{\partial x_{31}} & \frac{\partial E}{\partial x_{32}} & \frac{\partial E}{\partial x_{33}} \end{bmatrix} = \text{full_conv} \left(\begin{bmatrix} \frac{\partial E}{\partial o_{11}} & \frac{\partial E}{\partial o_{12}} \\ \frac{\partial E}{\partial o_{21}} & \frac{\partial E}{\partial o_{22}} \end{bmatrix}, \begin{bmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \end{bmatrix} \right)$$

Add WeChat powcoder

Gradient on X if it is not the inputs

<https://powcoder.com>

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

Add WeChat powcoder

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

Add WeChat powcoder

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

Why convolutional networks for NLP?

<https://powcoder.com>

Assignment Project Exam Help

- ▶ Even though bag-of-words models are simple and work well in some text classification tasks, they don't account for cases where multiple words combine to create meaning, such as "not interesting".
- ▶ The analogy with image processing is if the pixels are treated as separate features. (The analogy might be going too far).

Alternative text representations

<https://powcoder.com>

- ▶ When using feed-forward networks for text classification, the input text is represented as a sparse vector that encodes bag-of-words features.
- ▶ Alternatively, a text can be represented as a sequence of word tokens $w_1, w_2, w_3, \dots, w_M$. This view is useful for models such as **Convolutional Neural Networks**, or **ConvNets**, which processes text as a sequence.
- ▶ Each word token w_m is represented as a one-hot vector \mathbf{e}_{w_m} , with dimension V . The complete document can be represented by the horizontal concatenation of these one-hot vectors: $\mathbf{W} = [\mathbf{e}_{w_1}, \mathbf{e}_{w_2}, \dots, \mathbf{e}_{w_m}] \in R^{V \times M}$.
- ▶ To show that this is equivalent to the bag-of-words model, we can recover the word count from the matrix-vector product $\mathbf{W}[1, 1, \dots, 1]^T \in R^V$.

Input to a convolutional network in a text classification task

<https://powcoder.com>

Assignment Project Exam Help

When using conv nets for text classification, it is common to first “look up” the pretrained word embedding (e.g., the weight matrix produced by Word2Vec or GLOVE¹) for each word in the sequence:

<https://powcoder.com>

$$\mathbf{x}^{(0)} = \Theta^{(x \rightarrow z)} [\mathbf{e}_{w_1}, \mathbf{e}_{w_2}, \dots, \mathbf{e}_{w_M}]$$

$\mathbf{x}^{(0)} \in \mathbb{R}^{K_e \times M}$

where \mathbf{e}_{w_m} is a column vector of zeros, with a 1 at position w_m , K_e is the size of embeddings

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

“Convolve” the input with a set of filters (kernels)

<https://powcoder.com>

- ▶ A *filter* is a weight matrix of dimension $\mathbf{C}^{(k)} \in \mathbb{R}^{l_e \times h}$ 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 to cover the full width of the image.
- ▶ To merge adjacent words, we convolve $\mathbf{X}^{(0)}$ by sliding a set of filters across it:

<https://powcoder.com>

Add WeChat powcoder

$$\mathbf{X}^{(1)} = f(\mathbf{b} + \mathbf{C} * \mathbf{X}^{(0)})$$

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

Computing the convolution

- ▶ At each position m (the m th word in the sequence), we compute the element-wise product of the k th filter and the sequence of words of window size h (think of it as an ngram of length h) starting at m and take its sum: $\mathbf{C}^{(k)} \odot \mathbf{X}_{m:m+h-1}^{(0)}$
- ▶ The value of the m th position with the k th filter can be computed as:

$$x_m^{(1)} = \sum_{k'=1}^{K_e} \left(h_{k'} \sum_{n=1}^h x_{k',m+n-1}^{(0)} \right)$$

- ▶ 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://powcoder.com>

Assignment Project Exam Help

Conv ($X^{(0)}$, U)

0.4	0.7	0.3	0.3	2.3	3.2	0.7	1	0
0.5	1.3	0.2	1.5	0.8	0.6	1.8	0	1
1.2	2.2	1.1	4.3	0.5	0.3	3.4	0	0

Input text sequence “filter”

=

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

$X^{(1)}$

Padding

<https://powcoder.com>

Assignment Project Exam Help

- ▶ To deal with the beginning and end of the input, the base matrix is often padded with $h - 1$ column vectors of zeros at the beginning and end. this is called **wide convolution**
- ▶ If no padding is applied, then the output of each convolution layer will be $h - 1$ units smaller than the input. This is known as **narrow convolution**.

Pooling

<https://powcoder.com>

Assignment Project Exam Help

- ▶ After D convolutional layers, assuming filters have identical lengths, we have a representation of the document as a matrix $\mathbf{X}^{(D)} \in \mathbb{R}^{n \times m}$.
- ▶ It is very likely that the documents will be of different lengths, so we need to turn them into matrices of the same length before feeding them to a feedward network to perform classification.
- ▶ This can be done by **pooling** across times (over the sequence of words)

<https://powcoder.com>

Add WeChat powcoder

Prediction and training with CNN

<https://powcoder.com>

- ▶ 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 matrices for the feedforward network and the filters $\mathbf{C}^{(k)}$ of the CNN, as well as the biases.
- ▶ The parameters can be updated with backpropagation, which may involve computing the gradient for the max pooling function.

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

Different pooling methods

<https://powcoder.com>

Assignment Project Exam Help

- Max pooling

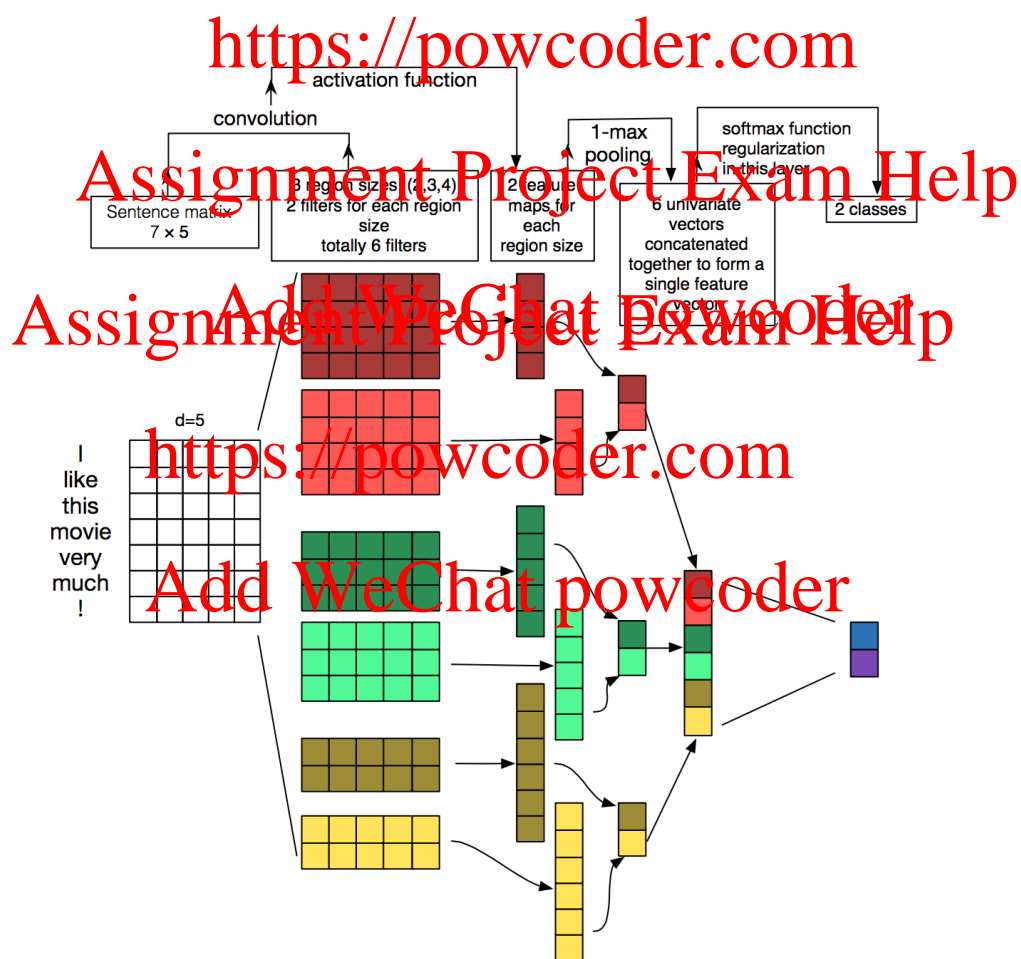
Assignment Project Exam Help
$$z_k = \max \left(x_{k,1}^{(D)}, x_{k,2}^{(D)}, \dots, x_{k,M}^{(D)} \right)$$

- Average pooling

<https://powcoder.com>

Add WeChat powcoder
$$z_k = \frac{1}{M} \sum_{m=1}^M x_{k,m}^{(D)}$$

A graphic representation of a CNN



Using Conv Nets in PyTorch

<https://powcoder.com>

nn.Conv1d Method Assignment Project Exam Help

```
In [19]: from torch import nn
conv1d = nn.Conv1d(3,2,2)
#parameters: input_channel (corresponding to embedding size)
#output channel kernel size
permuted_sequence = embedded_sequence_padded.permute(0,2,1)
conv_out = conv1d(permuted_sequence)
print(conv_out)
relu_out = nn.ReLU()(conv_out)
print(nn.MaxPool1d(2)(relu_out))
#parameters: kernel_size

tensor([[[[-0.2860, -0.4630, -0.0078],
          [-0.9166, -0.8111, -0.4047]],
        [[-0.0078, -0.0493, -0.6538],
          [ 0.4047, -0.5573, -0.5009]]], grad_fn=<SqueezeBackward1>)
tensor([[[[0.0000],
          [0.0000]],
        [[0.0000],
          [0.4047]]], grad_fn=<SqueezeBackward1>)
```

Add WeChat powcoder