ENCS5341 Machine Learning and Data Science

Neural Networks

STUDENTS-HUB.com

Deep Learning



What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do

STUDENTS-HUB.com



What I think I do

import keras

What I actually do

Historical Trends



STUDENTS-HUB.com

Uploaded By: Jibreel²Bornat

Introduction: Learning Multiple Components



STUDENTS-HUB.com

Uploaded By: Jibreel ³Bornat



IDEA: Use data to optimize features for the given task.

STUDENTS-HUB.com



What we want: Use parameterized function such that a) features are computed efficiently b) features can be trained efficiently

STUDENTS-HUB.com



- Everything becomes adaptive.
- No distiction between feature extractor and classifier.
- Big non-linear system trained from raw pixels to labels.

STUDENTS-HUB.com



Q: How can we build such a highly non-linear system?

A: By combining simple building blocks we can make more and more complex systems.

STUDENTS-HUB.com



NOTE: Each black box can have trainable parameters. Their composition makes a highly non-linear system.

STUDENTS-HUB.com

Thermediate representations/features

NOTE: System produces a hierarchy of features.

STUDENTS-HUB.com



Lee et al. "Convolutional DBN's for scalable unsup. learning..." ICML 2009

Marc'Aurelio Ranzat

STUDENTS-HUB.com





Lee et al. ICML 2009

Marc'Aurelio Ranzat



STUDENTS-HUB.com



Marc'Aurelio Ranzat



Lee et al. ICML 2009

STUDENTS-HUB.com

IDEA # 1

Learn features from data

IDEA # 2

Use differentiable functions that produce features efficiently

IDEA #3

End-to-end learning: no distinction between feature extractor and classifier

IDEA #4

"Deep" architectures: cascade of simpler non-linear modules

STUDENTS-HUB.com

Uploaded By: Jibree¹³Bornat

Perceptrons

Perceptrons

Leading to

- Neural Networks
- aka Multilayer Perceptron Networks
- aka Deep Learning

Neuron inspires Regression

- Graphical Representation of linear regression
- McCullough-Pitts Neuron
- Note: Not a graphical model



STUDENTS-HUB.com

Uploaded By: Jibreel⁵Bornat

Neuron inspires Regression

- Edges multiply the signal x_i by some weight θ_i .
- Nodes sum inputs:

$$f(x,\theta) = \sum_{d=1}^{D} \theta_d x^d + \theta_0$$

• Equivalent to Linear Regression





STUDENTS-HUB.com

Activation function

• How do we construct the neural output

$$f(\vec{x},\theta) = \theta^T \vec{x}$$



Linear Neuron

STUDENTS-HUB.com

Uploaded By: Jibree¹⁷Bornat

Activation function

• Sigmoid function or Squashing function

$$f(\vec{x},\theta) = g(\theta^T \vec{x}) \qquad g(z) = (1 + exp(-z))^{-1}$$



Logistic Neuron

Uploaded By: Jibreef[®]Bornat

7.5 10.0

Sigmoid Activation Function

0.0

Х

2.5

5.0

1.2

1.0

0.8

0.6

0.4

0.2

0.0

-0.2

-10.0 - 7.5 - 5.0 - 2.5

STUDENTS-HUB.com

Perceptron

• Classification squashing function



STUDENTS-HUB.com

The central idea behind Neural Networks

Building a model that is both complex and simple to manipulate by composing multiple functions together.



STUDENTS-HUB.com

Example:

Uploaded By: Jibreel[®]Bornat



We will talk later about the choice of activation function.

STUDENTS-HUB.com

Uploaded By: Jibreef¹Bornat



We will talk later about the choice of the output layer and the loss function. STUDENTS-HUB.com



STUDENTS-HUB.com

Uploaded By: Jibreef³Bornat



STUDENTS-HUB.com

Uploaded By: Jibreef⁴Bornat





STUDENTS-HUB.com

 $h = f(W^T X + b)$

The activation function should:

- Provide non-linearity
- Ensure gradients remain large through hidden unit

Common choices are

- sigmoid, tanh
- ReLU, leaky ReLU, Generalized ReLU
- softplus
- swish

STUDENTS-HUB.com

Uploaded By: Jibreef⁷Bornat

Sigmoid, σ () (aka logistic) and tanh



Derivative is **zero** for much of the domain. This leads to "vanishing gradients" in backpropagation. STUDENTS-HUB.com

Rectified Linear Unit, ReLU(), Exponential ReLU (ELU)



 $y = \max(0, x) + \alpha \min(0, e^{x} - 1)$ where α takes a small value



Two major advantages:

- 1. No vanishing gradient when x > 0
- 2. Provides sparsity (regularization) since y STUDENWSeHUR.com

No vanishing gradients and easy to calculate.

Softplus and Swish



The derivative of the softplus is the sigmoid logistic function, which is a smooth approximation of the derivative of the rectifier. So STUDENTSINGLINE (composite the softplus is continuous. Swish tends to work better than ReLU on deeper models across a number of challenging datasets.

Uploaded By: Jibreel[®]Bornat

Use ReLU or leaky ReLU or ELU to start.

If you are concerned for a few points improvement in performance, experiment with swish and softplus and treat them as hyperparameters. In other words choose the activation that gives you the best validation loss.

STUDENTS-HUB.com

Uploaded By: Jibreel³Bornat

Loss Function

Probabilistic modeling

Likelihood for a given measurement:

 $p(y_i|W;x_i)$

Assume **independency**, likelihood for all measurements:

$$L(W; X, Y) = p(Y|W; X) = \prod_{i} p(y_i|W; x_i)$$

Maximize the likelihood, or equivalently minimizing the -<u>ve</u> log-likelihood:

$$\mathcal{L}(W; X, Y) = -\log L(W; X, Y) = -\sum_{i} \log p(y_i | W; x_i)$$

STUDENTS-HUB.com

Uploaded By: Jibree³²Bornat

Loss Function

Do not need to design separate loss functions if we follow the probabilistic modeling approach, i.e. minimize the -ve likelihood function.

Examples:

• Distribution is **Normal** then -ve log-likelihood is **MSE** :

$$p(y_i|W;x_i) = \frac{1}{\sqrt{\{2\pi^2\sigma\}}} e^{-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2}}$$
$$\mathcal{L}(W;X,Y) = \sum_i (y_i - \hat{y}_i)^2$$

• Distribution is **Bernouli** then -ve log-likelihood is **Binary Cross-Entropy**:

$$p(y_i|W;x_i) = p_i^{y_i}(1-p_i)^{1-y_i}$$

$$\mathcal{L}(W; X, Y) = -\sum_{i} [y_{i} \log p_{i} + (1 - y_{i}) \log(1 - p_{i})]$$

STUDENTS-HUB.com

Uploaded By: Jibreef³Bornat

Loss Function

• For y that follow **Multinouli**, then likelihood is:

$$p(y_i|W; x_i) = \prod_k p(y_i = k)^{\mathbb{I}(y_i = k)}$$

$$Cross-Entropy$$

$$\mathcal{L}(W; X, Y) = -\sum_i \sum_k \mathbb{I}(y_i = k) \log p(y_i = k)$$

where k is the class and $\mathbb{I}(y_i = k)$ is the indicator function.

STUDENTS-HUB.com

Uploaded By: Jibreel⁴Bornat

Output unit for binary classification



STUDENTS-HUB.com

Uploaded By: Jibreel⁵Bornat
Output unit for multi-class classification



STUDENTS-HUB.com

Uploaded By: Jibreel[®]Bornat

SoftMax



STUDENTS-HUB.com

Uploaded By: Jibree³⁷Bornat

Cross-Entropy (Softmax)

Softmax $L_i = -\log(\frac{e^{s_{y_i}}}{\sum_i e^{s_j}})$

Score function $s = f(x_i, W)$ e.g., $f(x_i, W) = W \cdot [x_0, x_1, ..., x_N]^T$

Suppose: 3 training examples and 3 classes



Given a function with weights W, Training pairs $[x_i; y_i]$ (input and labels)



Output unit for regression



STUDENTS-HUB.com

Uploaded By: Jibreel[®]Bornat





STUDENTS-HUB.com

Uploaded By: Jibreef[®]Bornat



Uploaded By: Jibreef¹Bornat







We have seen that neural networks can represent complex functions, but are there limitations on what a neural network can express?

Theorem:

For any continuous function f defined on a bounded domain, we can find a neural network that approximates f with an arbitrary degree of accuracy.

One hidden layer is enough to represent an approximation of any function to an arbitrary degree of accuracy.

So why deeper?

STUDENTS-HUB.com

Uploaded By: Jibreef⁴Bornat

Layers





STUDENTS-HUB.com

×

Uploaded By: Jibreef⁵Bornat

Layers



Uploaded By: Jibreef®ornat

Layers



Why layers?

Representation matters!



Neural networks can **learn useful representations** for the problem. This is another reason why they can be so powerful!

STUDENTS-HUB.com

Uploaded By: Jibreef[®]Bornat







Depth = Repeated Compositions



STUDENTS-HUB.com

Uploaded By: Jibree[§]Bornat

Better Generalization with Depth



STUDENTS-HUB.com

Uploaded By: Jibree⁵¹Bornat

Shallow Nets Overfit More

Depth helps, and it's not just because of more parameters



Learning the coefficients

Start with single neuron



STUDENTS-HUB.com

Uploaded By: Jibree^{§3}Bornat

Start with all randomly selected weights. Most likely it will perform horribly. For example, in our heart data, the model will be giving us the wrong answer.



Start with all randomly selected weights. Most likely it will perform horribly. For example, in our heart data, the model will be giving us the wrong answer.



But what is the idea?

- Loss Function: Takes all of these results and averages them and tells us how bad or good the computer or those weights are.
- Telling the computer how bad or good is, does not help.
- You want to tell it how to change those weights so it gets better.

```
Loss function: \mathcal{L}(w_0, w_1, w_2, w_3, w_4)
```

For now let's only consider a single weight, $\mathcal{L}(w_1)$

STUDENTS-HUB.com

Uploaded By: Jibreef Bornat

Minimizing the Loss function

Ideally we want to know the value of w_1 that gives the minimal $\mathcal{L}(W)$



To find the optimal point of a function $\mathcal{L}(W)$

$$\frac{d\mathcal{L}(W)}{dW} = 0 \qquad \text{solve for } W^* = \operatorname{argmin}_W \mathcal{L}(W)$$

And find the W that satisfies that equation. **Sometimes** there is no explicit solution for that.

STUDENTS-HUB.com

Uploaded By: Jibree[†]Bornat

Estimate of the regression coefficients: gradient descent



A more flexible method is

- Start from a random point
 - 1. Determine which direction to go to reduce the loss (left or right)
 - 2. Compute the slope of the function at this point and step to the right if slope is negative or step to the left if slope is positive
 - 3. Goto to #1

STUDENTS-HUB.com

Uploaded By: Jibreef[®]Bornat

Gradient Descent

- Algorithm for optimization of first order to finding a minimum of a function.
- It is an iterative method.
- *L* is decreasing much faster in the direction of the negative derivative.
- The learning rate is controlled by the magnitude of η .

$$w^{(i+1)} = w^{(i)} - \eta \frac{d\mathcal{L}}{dw}$$



STUDENTS-HUB.com

Uploaded By: Jibree[§]Bornat

Batch and Stochastic Gradient Descent

$$\mathcal{L} = -\sum_{i} [y_i \log p_i + (1 - y_i) \log(1 - p_i)]$$

- Instead of using all the examples for every step, use a subset of them (batch).
- For each iteration *k*, use the following loss function to derive the derivatives:

$$\mathcal{L}^{k} = -\sum_{i \in b^{k}} [y_{i} \log p_{i} + (1 - y_{i}) \log(1 - p_{i})]$$

• which is an **approximation** to the full loss function.

STUDENTS-HUB.com

Uploaded By: Jibreef[®]Bornat

DATA



COMPLETE DATA 2 ONE EPUCH

RESHUFFLE DATA AND RESEAT

STUDENTS-HUB.com

Uploaded By: Jibree¹Bornat

Learning Rate

• Our choice of the learning rate has a significant impact on the performance of gradient descent.



 $L(\omega)$



When η is too small, the algorithm makes very little progress.

When η is too large, the algorithm may overshoot the minimum and has crazy oscillations.

When η is appropriate, the algorithm will find the minimum. The algorithm **converges**!

STUDENTS-HUB.com

Uploaded By: Jibreef²Bornat

How can we tell when gradient descent is converging? We visualize the loss function at each step of gradient descent. This is called the **trace plot**.





While the loss is decreasing throughout training, it does not look like descent hit the bottom.

Loss is mostly oscillating between values rather than converging.

The loss has decreased significantly during training. Towards the end, the loss stabilizes and it can't decrease further.

Uploaded By: Jibreef³Bornat

There are many alternative methods which address how to set or adjust the learning rate, using the derivative or second derivatives and or the momentum.

More on this later.

STUDENTS-HUB.com

Uploaded By: Jibreef⁴Bornat

Multilayer Networks

- Stack Neurons together
- The output from one layer is the input to the next
- Each Layer has its own sets of weights



STUDENTS-HUB.com

Uploaded By: Jibreef⁵Bornat

Linear Regression Neural Networks

• What happens when we arrange linear neurons in a multilayer network?



STUDENTS-HUB.com

Uploaded By: Jibreef Bornat

Linear Regression Neural Networks

- Nothing special happens.
- The product of two linear transformations is itself a linear transformation.



Neural Networks

- We want to introduce non-linearities to the network.
- Non-linearities allow a network to identify complex regions in space



STUDENTS-HUB.com

Uploaded By: Jibreef[®]Bornat

Linear Separability

- More layers can handle more complicated spaces but require more parameters
- Each node splits the feature space with a hyperplane
- A 2-layer network can represent any convex hull.



STUDENTS-HUB.com

Uploaded By: Jibreef⁹Bornat

Feed-Forward Networks

• Predictions are fed forward through the network to classify



STUDENTS-HUB.com

Uploaded By: Jibreel[®]Bornat

Feed-Forward Networks

• Predictions are fed forward through the network to classify



STUDENTS-HUB.com

Uploaded By: Jibreel¹Bornat
• Predictions are fed forward through the network to classify



STUDENTS-HUB.com

Uploaded By: Jibree⁷²Bornat

• Predictions are fed forward through the network to classify



STUDENTS-HUB.com

Uploaded By: Jibreel³Bornat

• Predictions are fed forward through the network to classify



STUDENTS-HUB.com

Uploaded By: Jibreel⁴Bornat

• Predictions are fed forward through the network to classify



STUDENTS-HUB.com

Uploaded By: Jibreel⁵Bornat

- We will do gradient descent on the whole network.
- Training will proceed from the last layer to the first.



STUDENTS-HUB.com

Uploaded By: Jibreel[®]Bornat

• Introduce variables over the neural network

$$\vec{\theta} = \{w_{ij}, w_{jk}, w_{kl}\}$$



Uploaded By: Jibreel⁷Bornat

- Introduce variables over the neural network
- Distinguish the input and output of each node



$$a_{j} = \sum_{i} w_{ij} z_{i} \qquad a_{k} = \sum_{j} w_{jk} z_{j} \qquad a_{l} = \sum_{k} w_{kl} z_{k}$$
$$z_{j} = g(a_{j}) \qquad z_{k} = g(a_{k}) \qquad z_{l} = g(a_{l})$$



• Training: Take the gradient of the last component and iterate backwards





Uploaded By: Jibree⁸¹Bornat

• Optimize last layer weights w_{kl} : $L_n = \frac{1}{2} (y_n - f(x_n))^2$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_{n} \left[\frac{\partial L_n}{\partial a_{l,n}} \right] \left[\frac{\partial a_{l,n}}{\partial w_{kl}} \right] \qquad \text{Calculus chain rule}$$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_{n} \left[\frac{\partial \frac{1}{2} (y_n - g(a_{l,n}))^2}{\partial a_{l,n}} \right] \left[\frac{\partial z_{k,n} w_{kl}}{\partial w_{kl}} \right] = \frac{1}{N} \sum_{n} \left[-(y_n - z_{l,n})g'(a_{l,n}) \right] z_{k,n}$$

$$\frac{z_i}{w_{ij}} \bigvee_{ij} \bigvee_{ij$$

Uploaded By: Jibree^{β2}Bornat

• Optimize last layer weights w_{kl} : $L_n = \frac{1}{2} (y_n - f(x_n))^2$

$$\begin{split} \frac{\partial R}{\partial w_{kl}} &= \frac{1}{N} \sum_{n} \left[\frac{\partial L_n}{\partial a_{l,n}} \right] \left[\frac{\partial a_{l,n}}{\partial w_{kl}} \right] \quad \text{Calculus chain rule} \\ \frac{\partial R}{\partial w_{kl}} &= \frac{1}{N} \sum_{n} \left[\frac{\partial \frac{1}{2} (y_n - g(a_{l,n}))^2}{\partial a_{l,n}} \right] \left[\frac{\partial z_{k,n} w_{kl}}{\partial w_{kl}} \right] = \frac{1}{N} \sum_{n} \left[-(y_n - z_{l,n})g'(a_{l,n}) \right] z_{k,n} \\ &= \frac{1}{N} \sum_{n} \delta_{l,n} z_{k,n} \end{split}$$



Uploaded By: Jibree^{β3}Bornat

• Optimize last hidden weights
$$\mathbf{w}_{jk}$$
: $\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_{n} \delta_{l,n} z_{k,n}$
 $\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_{n} \left[\frac{\partial L_n}{\partial a_{k,n}} \right] \left[\frac{\partial a_{k,n}}{\partial w_{jk}} \right]$
 $\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_{n} \left[\sum_{l} \frac{\partial L_n}{\partial a_{l,n}} \frac{\partial a_{l,n}}{\partial a_{k,n}} \right] \left[\frac{\partial a_{k,n}}{\partial w_{jk}} \right]$ Multivariate chain rule



Uploaded By: Jibree⁸⁴Bornat

• Optimize last hidden weights w_{jk} : $\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum \delta_{l,n} z_{k,n}$





Uploaded By: Jibree[§]Bornat

• Optimize last hidden weights w_{jk} : $\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_{k=1}^{N} \delta_{l,n} z_{k,n}$

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_{n} \left[\sum_{l} \delta_{l} \frac{\partial a_{l,n}}{\partial a_{k,n}} \right] [z_{j,n}]$$

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_{n} \left[\sum_{l} \delta_{l} w_{kl} g'(a_{k,n}) \right] [z_{j,n}] = \frac{1}{N} \sum_{n} [\delta_{k,n}] [z_{j,n}]$$

$$a_l = \sum_k w_{kl} g(a_k)$$



Uploaded By: Jibreel[®]Bornat

• Repeat for all hidden layers:

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_{n} \left[\frac{\partial L_{n}}{\partial a_{l,n}} \right] \left[\frac{\partial a_{l,n}}{\partial w_{kl}} \right] = \frac{1}{N} \sum_{n} \left[-(y_{n} - z_{l,n})g'(a_{l,n}) \right] z_{k,n} = \frac{1}{N} \sum_{n} \delta_{l,n} z_{k,n}$$

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_{n} \left[\frac{\partial L_{n}}{\partial a_{k,n}} \right] \left[\frac{\partial a_{k,n}}{\partial w_{jk}} \right] = \frac{1}{N} \sum_{n} \left[\sum_{l} \delta_{l,n} w_{kl} g'(a_{k,n}) \right] z_{j,n} = \frac{1}{N} \sum_{n} \delta_{k,n} z_{j,n}$$

$$\frac{\partial R}{\partial w_{ij}} = \frac{1}{N} \sum_{n} \left[\frac{\partial L_{n}}{\partial a_{j,n}} \right] \left[\frac{\partial a_{j,n}}{\partial w_{ij}} \right] = \frac{1}{N} \sum_{n} \left[\sum_{k} \delta_{k,n} w_{jk} g'(a_{j,n}) \right] z_{i,n} = \frac{1}{N} \sum_{n} \delta_{j,n} z_{i,n}$$

$$\frac{z_{i}}{w_{ij}} \int w_{ij} \int w_{jk} \int w_{jk} \int w_{kl} \int w_{kl} \int w_{kl} \int f(x,\vec{\theta})$$
STUDENTS-HUB.com

Uploaded By: Jibree^β⁷Bornat

• Now that we have well defined gradients for each parameter, update using Gradient Descent: ∂B



Uploaded By: Jibree⁸Bornat

- Error backpropagation unravels the multivariate chain rule and solves the gradient for each partial component separately.
- The target values for each layer come from the next layer.
- This feeds the errors back along the network.



Uploaded By: Jibree⁶⁹Bornat

Problems with Neural Networks

- Interpretation of Hidden Layers
- Overfitting (Many parameters)
- Require very large annotated datasets

Overfitting

- Augment training data by applying transformations on the training data (mirroring, rotation, affine, ...)
- Norm Penalties: *L*₂ and *L*₁ regularization
- Early Stopping
- Use "dropout" to regularize the weights in the globally connected layers (which contain most of the parameters).
 - Dropout means that half of the hidden units in a layer are randomly removed for each training example.
 - This stops hidden units from relying too much on other hidden units.

STUDENTS-HUB.com

Uploaded By: Jibreel¹Bornat

Early Stopping

STUDENTS-HUB.com

Early stopping: terminate while validation set performance is better. Sometimes is worth waiting a little before stopping. This is called patience.



Patience is defined as the number of epochs to wait before early stop if no progress on the validation set.

The patience is often set somewhere between 10 and 100 (10 or 20 is more common), but it really depends on the dataset and network.

Uploaded By: Jibreel²Bornat

Dropout

• Each time we present a training example, we randomly omit each hidden unit with probability 0.5.



(a) Standard Neural Net STUDENTS-HUB.com



(b) After applying dropout. Uploaded By: Jibreel³Bornat



• At test time, all nodes are used but the weights are scaled.



STUDENTS-HUB.com

Uploaded By: Jibree^βBornat

Optimizers

STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

Local Minima



STUDENTS-HUB.com

Uploaded By: Jibreel[®]Bornat

Critical Points

- Points with zero gradient
- 2nd-derivate (Hessian) determines curvature



STUDENTS-HUB.com

Uploaded By: Jibree^{†7}Bornat

Local Minima

- Old view: local minima is major problem in neural network training
- Recent view:
 - For sufficiently large neural networks, most local minima incur low cost
 - Not important to find true global minimum

Saddle Points



- Recent studies indicate that in high dim, saddle points are more likely than local min
- Gradient can be very small near saddle points

STUDENTS-HUB.com

Uploaded By: Jibree^PBornat

Poor Conditioning

- Poorly conditioned Hessian matrix
 - High curvature: small steps leads to huge increase
- Learning is slow despite strong gradients





STUDENTS-HUB.com

Uploaded By: Jibreel®Bornat

Momentum

• Oscillations because updates do not exploit curvature information



• Average gradient presents faster path to optimal: vertical components cancel out STUDENTS-HUB.com

Momentum

f is the Neural Network

• Old gradient descent:

$$g = \frac{1}{m} \sum_{i} \nabla_{W} L(f(x_i; W), y_i) \qquad \qquad W^* = W - \eta g$$

• New gradient descent with momentum:

$$\nu = \alpha \nu + (1 - \alpha) g$$

$$W^* = W - \eta v$$

STUDENTS-HUB.com

 $\widehat{a} \mid [0,1)$ controls how quickly effect of past gradients decay

Uploaded By: Jibreel⁰²Bornat

Adaptive Learning Rates



- Oscillations along vertical direction
 - Learning must be slower along parameter 2
- Use a different learning rate for each parameter?

STUDENTS-HUB.com

Uploaded By: Jibreel¹³Bornat

AdaGrad

• Accumulate squared gradients:

$$r_i = r_i + g_i^2$$

g is the gradient

• Update each parameter:

$$W_i = W_i - \frac{\epsilon}{\delta + \sqrt{r_i}} g_i$$

Greater progress along gently sloped directions

Inversely proportional to cumulative gradient

 $\delta\,$ is a small number, making sure this does not become too large

STUDENTS-HUB.com

Uploaded By: Jibreel⁴Bornat

Other Adaptive Learning Rate Methods

• RMSProp

• Use exponentially weighted average for gradient accumulation

$$r_i = \Gamma r_i + (1 - \Gamma) g_i^2$$

$$W_i = W_i - \frac{\epsilon}{\delta + \sqrt{r_i}} g_i$$

- Adam
 - RMSProp + Momentum
 - Works well in practice, it is fairly robust to hyper-parameters

STUDENTS-HUB.com

Uploaded By: Jibreel¹⁵Bornat

Other Neural Networks

STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

Convolutional Network

• The network is not fully connected.

STUDENTS-HUB.com

- Different nodes are responsible for different regions of the image.
- This allows for robustness to transformations.



Uploaded By: Jibreel⁷Bornat
A Convolutional Network



STUDENTS-HUB.com

Uploaded By: Jibreel[®]Bornat

A Convolutional Network



STUDENTS-HUB.com

Uploaded By: Jibreel⁹Bornat

Other Neural Networks

- Multiple Outputs
- Skip Layer Network
- Recurrent Neural Networks

Multiple Outputs

- Used for N-way classification.
- Each Node in the output layer corresponds to a different class.
- No guarantee that the sum of the output vector will equal 1.



STUDENTS-HUB.com

Uploaded By: Jibreelf Bornat

Skip Layer Network

STUDENTS-HUB.com

• Input nodes are also sent directly to the output layer.



Recurrent Neural Networks

• Output or hidden layer information is stored in a context or memory layer.



STUDENTS-HUB.com

Uploaded By: Jibreel¹³Bornat

Recurrent Neural Networks

• Vanilla RNN

STUDENTS-HUB.com



• Long-Short term memory (LSTM).



• Gated Recurrent Units (GRU).



Uploaded By: Jibreelf Bornat

Transformer Family



Back to our Object Recognition System

STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

Building an Object Recognition System



- Everything becomes adaptive.
- No distiction between feature extractor and classifier.
- Big non-linear system trained from raw pixels to labels.

STUDENTS-HUB.com

Marc'AUploaded By: Jibreel¹⁷Bornat

Fully Connected Neural Net



Marduploaded By: Jibreef¹⁸Bornat

Locally Connected Neural Net



1M hidden units Filter size: 10x10 100M parameters

Filter/Kernel/Receptive field: input patch which the hidden unit is



Locally Connected Neural Net







Uploaded By: Jibreel¹²¹Bornat



Uploaded By: Jibreef²Bornat



STUDENTS-HUB.com

Uploaded By: Jibreef³Bornat



Uploaded By: JibreeffaBornat



CONVOLUTIONAL LAYER



STUDENTS-HUB.com

Uploaded By: Jibreelf Bornat

"Convolution" Operation



STUDENTS-HUB.com

Uploaded By: Jibreef⁷Bornat

"Convolution" Operation



STUDENTS-HUB.com

Uploaded By: Jibreef®Bornat

Convolutions – what happens at the edges?

If we apply convolutions on a normal image, the result will be down-sampled by an amount depending on the size of the filter.



We can avoid this by padding the edges in different ways.

STUDENTS-HUB.com

Uploaded By: Jibreef⁹Bornat

Padding



Full padding. Introduces zeros such that all pixels are visited the same number of times by the filter. Increases size of output.



Same padding. Ensures that the output has the same size as the input.

Uploaded By: Jibreel³Bornat

STUDENTS-HUB.com



Stride controls how the filter convolves around the input volume.

The formula for calculating the output size is:

$$O = \frac{W - K + 2P}{S} + 1$$

Where O is output dim, W is the input dim, K is the filter size, P is padding and S the stride





STUDENTS-HUB.com

Uploaded By: Jibreel²Bornat



A convolutional layer with 16 3x3 filters that takes 32x32 RGB image as input.



STUDENTS-HUB.com

Uploaded By: Jibreel³Bornat

STUDENTS-HUB.com



Q.: how can we make the detection robust to the exact location of the eye?



STUDENTS-HUB.com





Pooling

POOLING LAYER



STUDENTS-HUB.com Input feature maps

output feature maps

Uploaded By: Jibreeli Bornat

Pooling



STUDENTS-HUB.com

Uploaded By: Jibree¹³⁷Bornat



Example: max pool with 2x2 window and stride 1

STUDENTS-HUB.com

Uploaded By: Jibreel³⁸Bornat

Pooling

Pooling involves selecting:

- A pooling operation, much like a filter to be applied to feature maps: max, mean, median.
- The size of the pooling operation.
- The stride.

The size of the pooling operator must be smaller than the size of the feature map; specifically, it is <u>almost always</u> 2×2 applied with a stride of 2 using max pooling.

Invariant to small, "local transitions"

Face detection: enough to check the presence of eyes, not their precise location

Reduces input size of the final fully connected layers (more later)

No learnable parameters

STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

Pooling: example with stride 2x2

1	1	2	5
5	7	7	8
3	1	1	0
1	2	3	4

max pool with 2x2 window and stride 2x2





mean pool with 2x2 window and stride 2x2

3.5	5.5	
1.75	2	

Uploaded By: Jibreef®Bornat

STUDENTS-HUB.com

Building a CNN

A convolutional neural network is built by stacking layers, typically of 3 types:



One stage



Whole system Input Image + Image +

STUDENTS-HUB.com

Uploaded By: Jibreef¹Bornat

Building a CNN



STUDENTS-HUB.com

Uploaded By: Jibreef²Bornat

Building a CNN



Uploaded By: Jibreef¹³Bornat
Building a CNN



STUDENTS-HUB.com

Uploaded By: Jibreeff4Bornat

A Convolutional Network



STUDENTS-HUB.com

Uploaded By: Jibreeff5Bornat

Example

- Let **C** be a CNN with the following disposition:
 - Input: 32x32x3 images
 - Conv1: 8 3x3 filters, stride 1, padding=same
 - Conv2: 16 5x5 filters, stride 2, padding=same
 - Flatten layer
 - Dense1: 512 nodes
 - Dense2: 4 nodes
- How many parameters does this network have?

 $(8 \times 3 \times 3 \times 3 + 8) + (16 \times 5 \times 5 \times 8 + 16) + (16 \times 16 \times 16 \times 512 + 512) + (512 \times 4 + 4)$

Conv1	Conv2	Dense1	Dense2

STUDENTS-HUB.com

Uploaded By: Jibreeff Bornat

Conv Nets: Training

- In a convolutional layer, we are basically applying multiple filters over the image to extract different features. But most importantly, we are learning those filters!
- Since convolutions and sub-sampling are differentiable, we can use standard backpropagation.
- Algorithm:
 - Given a small mini-batch
 - Forward Propagation
 - Backward Propagation
 - Parameter Update

STUDENTS-HUB.com

Uploaded By: Jibreed¹⁷Bornat

What do CNN layers learn?



STUDENTS-HUB.com

Uploaded By: Jibreeft®Bornat

What do CNN layers learn?



STUDENTS-HUB.com

Uploaded By: Jibreef Bornat

Layers Receptive Field

- The **receptive field** is defined as the **region** in the input space that a particular CNN's feature is looking at (i.e., be affected by).
- The receptive field size is a crucial issue in many visual tasks, as the output must respond to large enough areas in the image to capture information about large objects.
- Let's look at the receptive field in 1D, no padding, stride 1 and kernel 3x1



Layers Receptive Field

In 2D, it works the same way.

The receptive field for layer *L*, *r*₀, can be calculated using the recursive formula:

$$r_0 = \sum_{l=1}^{L} \left((k_l - 1) \prod_{i=1}^{l-1} s_i \right) + 1$$

*k*_lkernel size (positive integer) *s*_l stride (positive integer)



Note: For every max-pooling layer, we multiply the receptive field by the window size (assuming stride is the same as the window size)

STUDENTS-HUB.com

Uploaded By: Jibreef Bornat

Popular Convolutional Neural Networks (CNNs)



STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

Software



The Microsoft Cognitive Toolkit

A free, easy-to-use, open-source, commercial-grade toolkit that trains deep learning algorithms to learn like the human brain.



Uploaded By: Jibreel Bornat

STUDENTS-HUB.com