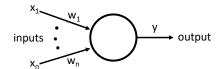
CMSC 373 Artificial Intelligence Fall 2025 13-Backpropagation

Deepak Kumar Bryn Mawr College

1

The Perceptron (1958)

• A single "neuron" (<u>unit</u>) aka Threshold Logic Unit (TLU)



Transfer Function

T is the Threshold value (assume T = 0)

$$I = \sum_{i=1}^{i=n} w_i x_i$$
$$y = \begin{cases} +1, & \text{if } I \ge T \\ -1, & \text{if } I < T \end{cases}$$

2

Learning in Neural Networks (so far)

Epochs
Forward Pass
Hyperparameters
Labelled Dataset
Learning Rule
Model
Parameters

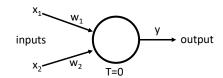
- A labelled dataset e.g., Iris Dataset, MNIST Numbers Dataset
- A model for the network e.g., the Perceptron (single TLU) Parameters refer to the number of weights
- $I = \sum_{i=1}^{l=n} w_i x_i$ $y = \begin{cases} +1, & \text{if } I \ge T \\ 0 & \text{otherwise} \end{cases}$
- A forward pass/prediction algorithm Compute net input. Compute activation.
- A Learning Rule/Algorithm
 Learning occurs by changing weights.
 Specifies change in the weights using the Error/Loss
 In Perceptrons, weights are changed after every prediction.

Perceptron Learning Rule: $\overline{w_{new}} = \overline{w_{old}} + \beta y^* \overline{x}$

• Several **epochs** of training are needed (how many???) This is a **hyperparameter**.

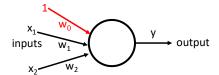
3

Introducing Bias



- Instead of using an arbitrary Threshold value, we can turn it into an input (=1)
- The weight on the bias, \boldsymbol{w}_0 can then be learned using the same algorithm.

$$\overline{\mathbf{w}} = [w_0, w_1, w_2]$$
$$\overline{\mathbf{x}} = [x_0, x_1, x_2]$$



 More often, in other networks, the net input is determined using the following (and no bias is used for output layer):

$$I = \sum_{i=1}^{i=n} w_i x_i + \overline{\boldsymbol{b}}$$

4

4

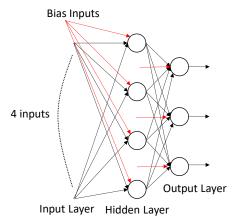
Multi-Layer Perceptron Network

• Example: This could be a network that can recognize all three categories of irises from the Iris dataset.

4 inputs, 3 outputs (Hyperparameters) 4x4 (input to hidden) + 4x3 (hidden to output) weights + 7 bias inputs

#Parameters = 16+12+7 = 35

- Since all units are linear TLUs this network can only learn linear functions.
- · We need to make each unit non-linear.



5

5

Backpropagation Network (Classic Version)

• Net Input

$$I = \sum_{i=1}^{i=n} w_i x_i + \overline{\boldsymbol{b}}$$

• Activation Function (Sigmoid)

$$f(I) = \frac{1}{1 + e^I}$$

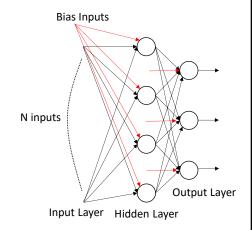
· Learning Rule

$$\Delta w_{ij} = \beta * E * f(I_j)$$

· Error/Loss

$$\mathbf{E}_{j}^{\mathrm{output}} = y_{j}^{\mathrm{desired}} - y_{j}^{\mathrm{actual}}$$

$$\mathbf{E}_{i}^{\text{hidden}} = \frac{\text{df}(\mathbf{I}_{i}^{\text{hidden}})}{\text{dI}} \sum_{j=1}^{n} (w_{ij} \mathbf{E}_{j}^{\text{output}})$$



6

Backpropagation Network (Classic Version)

Net Input

$$I = \sum_{i=1}^{i=n} w_i x_i + \overline{\boldsymbol{b}}$$

Activation Function (Sigmoid)

$$f(I) = \frac{1}{1 + e^{I}}$$

Learning Rule

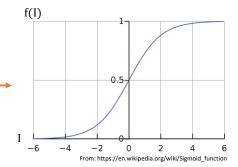
$$w_{ij}^{\text{new}} = w_{ij}^{\text{old}} + \beta * E * f(I_j)$$

where β is the **Learning Constant** (0 < β < 1), E is the error/loss

Error/Loss (E)

$$E_i^{\text{output}} = y_i^{\text{desired}} - y_i^{\text{actual}}$$

$$E_{i}^{hidden} = \frac{df(I_{i}^{hidden})}{dI} \sum_{j=1}^{n} (w_{ij} E_{j}^{output})$$



Sigmoid Function

- Output is always between 0 and 1
- Its derivative (slope) is always positive
- $\frac{\mathrm{df}(I)}{\mathrm{dI}} = \mathrm{f}(I)(1 \mathrm{f}(I))$

Backpropagation Network (Classic Version)

Net Input

$$I = \sum_{i=1}^{i=n} w_i x_i + \overline{\boldsymbol{b}}$$

Activation Function (Sigmoid)

$$f(I) = \frac{1}{1 + e^{I}}$$

Learning Rule¹

$$w_{ij}^{\text{new}} = w_{ij}^{\text{old}} + \beta * \mathbf{E} * \mathbf{f}(I_j)$$

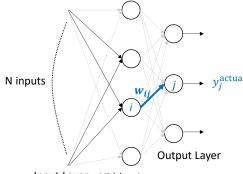
where β is the **Learning Constant** (0 < β < 1), E is the error/loss

Error/Loss (E)

$$E_{j}^{\text{output}} = y_{j}^{\text{desired}} - y_{j}^{\text{actual}}$$

$$E_{j}^{\text{hidden}} = \frac{\text{df}(I_{i}^{\text{hidden}})}{I_{i}^{\text{hidden}}} \sum_{i=1}^{n} (w_{ij} E_{i}^{\text{output}})$$

 $E_{i}^{hidden} = \frac{df(I_{i}^{hidden})}{dI} \sum_{i=1}^{n} (w_{ij} E_{j}^{output})$



Input Layer Hidden Layer

1. Earlier versions used Delta Rule (Adaline): $w_{new} = w_{old} + \frac{\beta_{EX}}{|x|^2}$

8

Backpropagation Network (Classic Version)

Net Input

$$I = \sum_{i=1}^{i=n} w_i x_i + \overline{\boldsymbol{b}}$$

• Activation Function (Sigmoid)

$$f(I) = \frac{1}{1 + e^{I}}$$

· Learning Rule

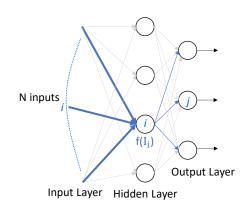
$$w_{ij}^{\text{new}} = w_{ij}^{\text{old}} + \beta * E * f(I_i)$$

where β is the **Learning Constant** (0 < β < 1), E is the error/loss

• Error/Loss (E)

$$E_i^{\text{output}} = y_i^{\text{desired}} - y_i^{\text{actual}}$$

$$E_{i}^{hidden} = \frac{df(I_{i}^{hidden})}{dI} \sum_{j=1}^{n} (w_{ij} E_{j}^{output})$$



9

9

Backpropagation Network (Classic Version)

Net Input

$$I = \sum_{i=1}^{i=n} w_i x_i + \overline{\boldsymbol{b}}$$

• Activation Function (Sigmoid)

$$f(I) = \frac{1}{1 + e^I}$$

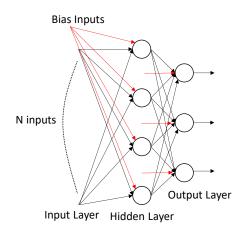
· Learning Rule

$$\Delta w_{ij} = \beta * E * f(I_j)$$

Error/Loss

$$E_{j}^{\text{output}} = y_{j}^{\text{desired}} - y_{j}^{\text{actual}}$$

$$\mathbf{E}_{i}^{\text{hidden}} = \frac{\text{df}(\mathbf{I}_{i}^{\text{hidden}})}{\text{dI}} \sum_{j=1}^{n} (w_{ij} \mathbf{E}_{j}^{\text{output}})$$



10

Backpropagation (Classic) Training Algorithm

```
set #epochs = 0
   total error = 0
   \label{eq:formula} \textbf{for} \ \mathsf{each} \ \mathsf{pattern} \ \mathsf{in} \ \mathsf{training} \ \mathsf{set} \ \textbf{do}
      do a forward pass
         for each unit in the hidden layer do
            compute net input I, and activation, f(I)
save f(I) for backpropagation
         for each unit in the output layer do
            compute net input I, and activation f(I)
            output y = f(I)
      do backward pass
for each unit in the output layer do
            \textit{compute error} = \textit{desired} - \textit{actual output} \, (\texttt{E}_{i}^{\textit{output}} = y_{i}^{\textit{desired}} - y_{i}^{\textit{actual}})
            total error = total error + error
            compute incoming error = weighted sum of output later errors (\sum_{j=1}^n (w_{ij} \mathbf{E}_{\mathbf{i}}^{\mathrm{output}}))
            compute final error = incoming error * f(I) * (1 – I) [derivative)
         for each unit in the output layer do
for each weight from a hidden layer to unit do
               compute weight change \beta * \operatorname{error} * f(I) and update weight
         for each unit in the middle layer do
            for each weight from an input layer unit do compute weight change \beta* final error * f(I) and update weight
      #epochs = #epochs + 1
until total error < maximum acceptable error or #epochs reaches limit
```

set minimum acceptable error and #epochs to train, set eta

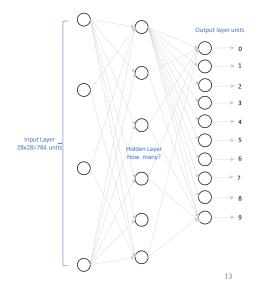
1

11

Example: Recognizing Handwritten Digits

- MNIST Dataset 70,000 images (28x28 pixels), grayscale values (in range 0 (white) to 255 (black).
- Training set: 60,000 images
- Testing set: 10,000 images
- Task: Given an image, classify it as [0,1,...,9]

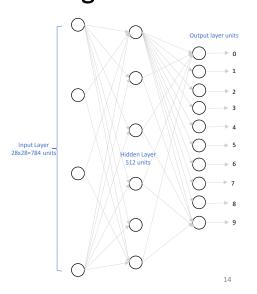
- Assume we will use a 3-layer network: input, hidden, and output
- Input will be 28x28=784 units
- 10 outputs, one for each digit. Say the network is shown a 6, we would then expect the output of 6 to be high (closer to 1.0) compared to others.
- How many hidden units????



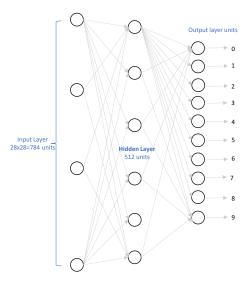
13

MNIST Digit Recognition: The Design

- Assume we will use a 3-layer network: input, hidden, and output
- Input will be 28x28=784 units
- 10 outputs, one for each digit. Say the network is shown a 6, we would then expect the output of 6 to be high (closer to 1.0) compared to others.
- How many hidden units???
 No known science to this.
 Let's say we have 512 units in hidden layer. (Why??)



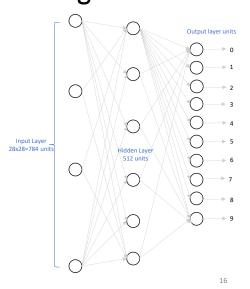
- A 3-layer network: input, hidden, and output layers
- Input will be 28x28=784 units
- 10 outputs, one for each digit. 512 units in hidden layer.
- All units will have a Sigmoid activation function.
- How to determine error/loss at output? (use our formulation?)
- What should be the value of β ?



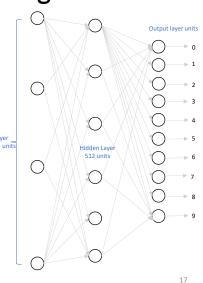
15

MNIST Digit Recognition: The Design

- A 3-layer network: input, hidden, and output layers
- Input will be 28x28=784 units
- 10 outputs, one for each digit. 512 units in hidden layer.
- All units will have a Sigmoid activation function.
- How to determine error/loss at output? (use our formulation?)
- What should be the value of β ?
- How many parameters are there?



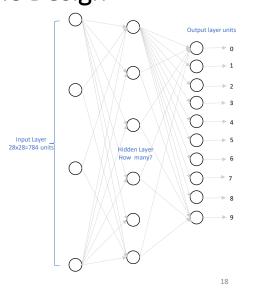
- A 3-layer network: input, hidden, and output layers
- Input will be 28x28=784 units
- 10 outputs, one for each digit. 512 units in hidden layer.
- All units will have a Sigmoid activation function. Input Layer 28x28-784 units
- How to determine error/loss at output? (use our formulation?)
- What should be the value of β ?
- How many parameters are there?
 784x512 + 512 (bias) + 512x10 + 10 (Bias) = 407,050



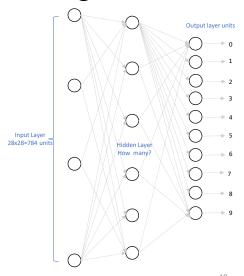
17

MNIST Digit Recognition: The Design

- A 3-layer network: input, hidden, and output layers
- Input will be 28x28=784 units
- 10 outputs, one for each digit. 512 units in hidden layer.
- All units will have a Sigmoid activation function.
- How to determine error/loss at output? (use our formulation?)
- What should be the value of β?
- How many parameters are there? 784x512 + 512 (bias) + 512x10 + 10 (Bias) = 407,050
- What are the hyperparameters?



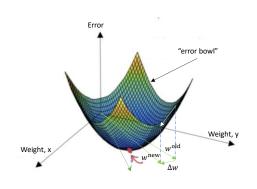
- · A 3-leayer network: input, hidden, and output layers
- Input will be 28x28=784 units
- 10 outputs, one for each digit. 512 units in hidden layer. (Why??)
- All units will have a Sigmoid activation function.
- How to determine error/loss at output? (use our formulation?)
- What should be the value of β ?
- How many parameters are there? 784x512 + 512 (bias) + 512x10 + 10 (Bias) = 407,050
- What are the hyperparameters?
 # Layers, # Units in each layer, Activation Function, β,
 # epochs, Minimum acceptable error, etc.



19

Backpropagation: Gradient Descent

- Learning in a neural network using Backpropagation is essentially a Gradient Descent process.
- Each change in the weights is an attempt to reduce error and descend into the lowest possible position in the "error bowl" (as shown in a 2-D weight vector case)
- In higher dimensional weight vectors (typical ML situations), the error surface can be quite complex.

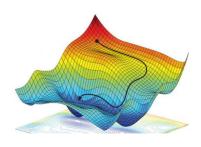


From: https://builtin.com/data-science/gradient-descent

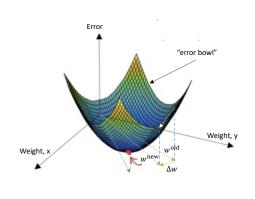
20

Backpropagation: Gradient Descent

 In higher dimensional weight vectors (typical ML situations), the error surface can be quite complex.



From: https://poissonisfish.com/2023/04/11/gradient-descent/



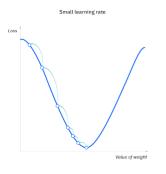
From: https://builtin.com/data-science/gradient-descent

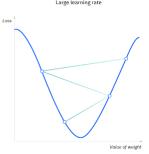
21

21

Backpropagation: Learning rate

- Learning Rate, β (0 < 1)
- The value of β determines how fast or slow the gradient descent takes place.
- Typically, one starts with a higher value (say 0.5 or 0.6) and then decrease it as the learning/epochs progresses. This is called a Learning Rate Schedule.





From: https://www.ibm.com/topics/gradient-descent

22

MNIST Digit Recognition: More Decisions

- In what order do we present the patterns? As they are in the training set? Or, randomly?
- Do we do a backpropagation pass after every input?

Choices:

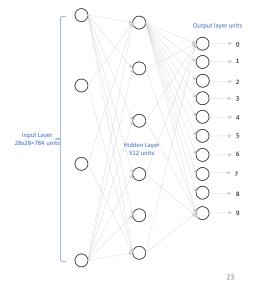
Do a backpropagation pass after every input.

Do the backward pass after all the inputs have been seen, and errors recorded.

Do a small batch of data and the do a backward pass.

- Is there a better way to assess error/loss?
- · Are there any other weight update mechanisms?

Most of the period from 1986 until now has been spent on studying these questions.



23

Decisions/Choices...

 In what order do we present the patterns? As they are in the training set? Or, randomly?
 If patterns are chosen at random (without replacement), we call it Stochastic Gradient Descent (SGD)

Choices:

Do a backpropagation pass after every input. True SGD

Do the backward pass after all the inputs have been seen, and errors recorded. Full batch SGD

Do a small batch of data and then do a backward pass. Mini Batch SGD (each batch is a power of 2)

- Is there a better way to assess error/loss? Loss Functions
- Are there any other weight update mechanisms? Optimizers (also manage Learning rate schedules)

Most of the period from 1986 until now has been spent on studying these questions.

Backpropagation Bias Epochs Forward Pass Full Batch SGD Gradient Descent Hyperparameters Labelled Dataset Learning Rule Loss Function Mini Batch SGD Model Optimizer **Parameters** SGD True SGD

Advances in NN

Better hardware

Laptops became 5000 times faster between 1990 and 2010 Use of GPUs for faster processing (NVIDIA, AMD). Took off in 2011 Google's Tensor Processing Units (TPUs), 2016

- Wide availability of datasets and benchmarks MNIST, ImageNet, etc.
- Better Algorithms

Better activation functions
Better weight initialization schemes
Better optimization schemes (RMSprop, Adam)

• Widely available toolsets for creating and training NNs Theano, TensorFlow, Scikit Learn, PyTorch, Keras

25

25

Backpropagation Network (Updated)

Net Input

$$I = \sum_{i=1}^{i=n} w_i x_i + \overline{\boldsymbol{b}}$$

Activation Functions

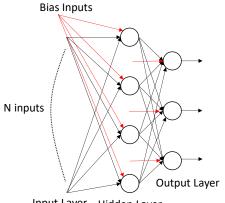
Sigmoid, Relu, Softmax, tanh, exponential, etc.

Loss Functions

Binary Cross Entropy, Categorical Cross Entropy, Poisson, KL Divergence, Mean Squared Error, Mean Absolute Error, Cosine Similarity, etc.

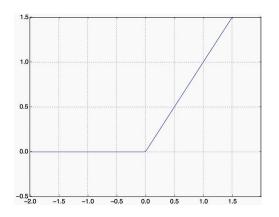
• Optimizer (Learning Rule)

SGD, RMSprop, Adam, Adadelta, Adamax, Namad, etc.

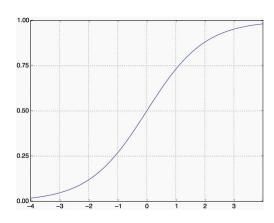


Input Layer Hidden Layer

Popular Activation Functions



Relu – Rectified Linear Unit: $f(I) = \max(0, I)$



Sigmoid: $f(I) = \frac{1}{1+e^{I}}$

27

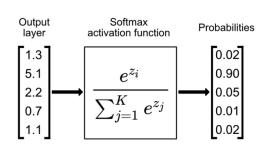
27

Softmax Activation Function

 Transforms a vector of arbitrary numbers into a probability distribution.

Each value is in [0.0..1.0] Summ of all values is 1.0

Useful for multi-class classification problems. E.g., Recognizing handwritten digits (ten possible outcomes [0, ..9])



Vocabulary

 In what order do we present the patterns? As they are in the training set? Or, randomly?
 If patterns are chosen at random (without replacement), we call it Stochastic Gradient Descent (SGD)

Choices:

Do a backpropagation pass after every input. True SGD

Do the backward pass after all the inputs have been seen, and errors recorded. Full batch SGD

Do a small batch of data and then do a backward pass. Mini Batch SGD (each batch is a power of 2)

- Is there a better way to assess error/loss? Loss Functions
- Are there any other weight update mechanisms?
 Optimizers (also manage Learning rate schedules)

Most of the period from 1986 until now has been spent on studying these questions.

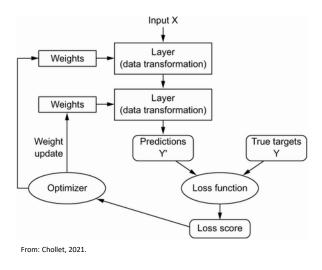
Adam Backpropagation Binary Cross Entropy Categorical Cross Entropy Epochs Exponential Forward Pass Full Batch SGD **Gradient Descent** Hyperparameters Labelled Dataset Learning Rule Loss Function Mean Absolute Error Mean Squared Error Mini Batch SGD Model Optimizer Parameters Relu RMSProp SGD Sigmoid Softmax Tanh True SGD

29

29

The Learning Paradigm Actual Data Test Dataset Finished Trained Learning/Training Model Model Algorithm Outputs Outputs Accuracy Model Design Labelled Training Dataset 30

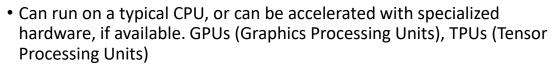
How NN Learning Works



31

Introducing Keras

- Deep Learning API for Python (2016-17)
- Built on top of TensorFlow (2015)



 Makes Neural Network design, implementation, and exploration akin to building with LEGOs!

32

TPU

From: Chollet, 2021.

Keras

GPU

CPU

Typical Keras Workflow

 Acquire, prepare, and load the dataset
 Keras has several predefined datasets available: MNIST Digits, CIFAR10, CIFAR100, IMDB Reviews for
 sentiment classification, Reuters Newswire classification, Fashion MNIST, Boston Housing price regression (see https://keras.io/api/datasets/)

Design and Build the Model

How many layers to use? How many units in each later? What activation function to use? (see https://keras.io/api/layers/activations/)

Compile the Model

Decide which optimizer to use, loss function, accuracy metric https://keras.io/api/optimizers/, https://keras.io/api/optimizers/.

Train/Fit the Model

Provide the training data and its labels, number of epochs to train, batch size

Test/Validate the Model

Use the test data to test how well the trained model performs

33

Over to Colab...

See Lab for Recognizing Handwritten Digits

References

- M. Caudill and C. Butler: Understanding Neural Networks, Volume 1, MIT Press, 1993.
- F. Chollet: Deep Learning with Python, Second Edition, Manning2021.
- A Geron: Hands-on Machine Learning with SciKit-Learn, Keras and TensorFlow, Oreilly, 2019.
- M. Mitchell: *Artificial Intelligence: A Guide For Thinking Humans*, Farrar, Strouss, Giroux, 2019.
- Rumelhart, McClelland, and the PDP Research Group: *Parallel Distributed Processing*, Volumes 1 & 2. MIT Press, 1986.

35