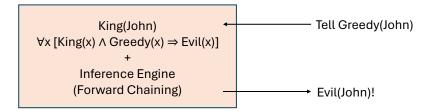
Review for Exam 2

1

Symbolic AI Techniques

Forward Chaining Inference

• Tell-Ask Systems

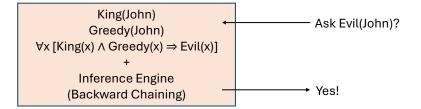


3

3

Backward Chaining Inference

• Tell-Ask Systems



4

Л

PROLOG Knowledgebase

```
isa(X, mammal):-
has(X, hair).
isa(X, bird):-
has(X, feathers),
flies(X).
isa(X, carnivore):-
isa(X, mammal),
has(X, pointed_teeth),
has(X, claws),
has(X, forward_pointing_eyes).
isa(X, ungulate):-
isa(X, mammal),
chews_cud(X);
has(X, hoofs).
isa(X, cheetah):-
isa(X, carnivore),
color(X, tawny),
has(X, dark_spots).
isa(X, giraffe):-
isa(X, ungulate),
has(X, long_legs),
has(X, long_legs),
has(X, long_neck),
color(X, tawny),
has(X, dark_spots).
```

% Facts:
has(stretch, hair).
chews_cud(stretch).
has(stretch, long_legs).
has(stretch, long_neck).
color(stretch, tawny).
has(stretch, dark_spots).
has(swifty, hair).
has(swifty, pointed_teeth).
has(swifty, claws).
has(swifty, forward_pointing_eyes).
color(swifty, tawny).
has(swifty, dark_spots).

5

Vocabulary

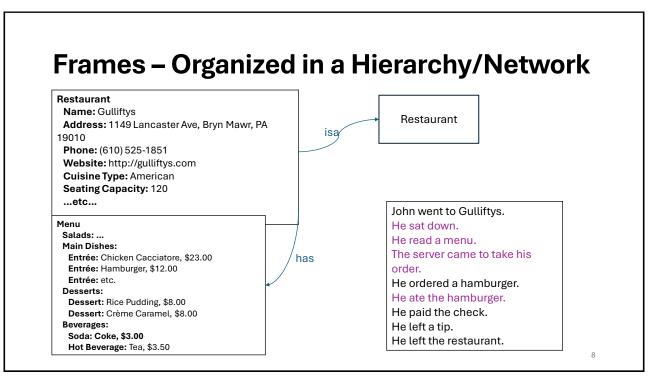
Knowledge Engineering
FOPC
Knowledge Base
Tell-Ask Systems
Forward Chaining Inference
Backward Chaining Inference
Definite Clauses
Logic Programming
PROLOG

Meaning Representation Systems

- Logic
- Semantic Networks
- Frames
- Conceptual Dependency
- many many others...
- Japan's Fifth Generation Project, and several others
- Successes of KR&R
- Limitations of KR&R

7

/



CD – Representations & Inferences

· John went to New York.

ACTOR: John ACTION: PTRANS OBJECT: John DIRECTION TO: New York FROM: unknown

· John bought a book from Mary.

ACTOR: John ACTION: ATRANS OBJECT: money DIRECTION TO: John FROM: Mary

John read a book.

ACTOR: John ACTION: ATTEND OBJECT: eyes DIRECTION TO: book FROM: unknown John drank a glass of milk.

ACTOR: John ACTION: INGEST OBJECT: milk DIRECTION TO: mouth of John FROM: glass

Instrument: ACTOR: John

ACTION: PTRANS
OBJECT: glass containing milk
DIRECTION TO: mouth of John
FROM: table

Instrument:

ACTOR: John
ACTION: MOVE
OBJECT: hand of John
DIRECTION TO: glass
FROM: unknown
Instrument:
ACTOR: John
ACTION: GRASP

OBJECT: glass of milk
DIRECTION TO: hand of John
FROM: unknown

9

9

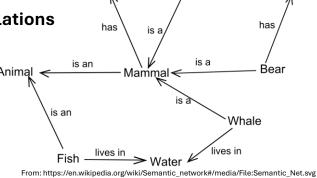
Semantic Networks

• Represents semantic relations between concepts in a network.

It is a **graph vertices** represent **concepts edges** represent **semantic relations**

• Also represented as a triple:

<entity, relation, entity>
<mammal, has, vertebra>
<bear, is a, mammal>
<fish, lives in, water>



Cat

Vertebra

has

10

Fur

Large-Scale AI/KR&R Efforts

Japanese Fifth Generation Project

10-year project started in 1982. Creating supercomputers for future AI development using PROLOG. Develop Knowledge Information Processing Systems. Spent ~\$300 million over twelve years. (IBM's research expenditure in 1982 was \$1.5 billion!)

Limited Results

Developed foundations for concurrent logic programming. Developed Kappa (a parallel DBMS) Developed an automated theorem prover, MGTP Attempt at applications in bioinformatics.

Did not meet commercial success

The arrival of SUN Workstations and x86 based machines far surpassed the abilities. Failed on attempts to develop concurrent KR&R systems.

· Perhaps the ideas were far ahead of its time.

For example, current multi-core processors being used in current AI work.

11

11

CYC - The Ultimate Expert System

Methodology

Developed a representation language, CycL.

Developed a set of representations (ontological engineering)

Developed a massive knowledgebase comprising human knowledge.

Connected CYC's knowledge to Wikidata (Wikipedia) and other large knowledgebases.

An inference engine.

The Seasons of Al

• 1950s - 1966 First Al Summer: Irrational Exuberance

Early successes in game playing, theorem proving, problem solving

1967 – 1977 First Al Winter

No useful deliverables led to loss of research funding and cancellation of AI programs. In UK *The Lighthill Report* (toy AI systems do not scale due to combinatorial explosion).

1978 – 1987 Second Al Summer/Spring

Rise of knowledge-based systems, success of Expert Systems. Boom times.

1988 – 1993 Second Al Winter

Failure of AI Hardware companies (Symbolics, LMI, Lisp Machines) and AI Companies (Teknowledge, Inference Corp. etc.) Commercial deployments of Expert Systems were discontinued.

· 1993 - 2011 Third AI Summer (Mostly academic advances)

Statistical approaches and extensions to logic (Bayesian Nets), Non-Monotonic Reasoning (in Logic), Fuzzy Logic, advances in Machine Learning (Decision Trees, Random Forests, Neural Nets), Cognitive Models, Logic Programming, Case-Based Reasoning, Genetic Algoritms, Agent-based approaches, etc.

· 2011 - Now Third AI Spring

Rise of Deep Learning, Neuro-symbolic AI, ChatGPT and other chatbots, generative AI.

13

13

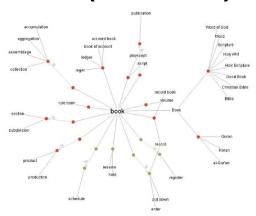
Towards Usable Representations (1985-now)

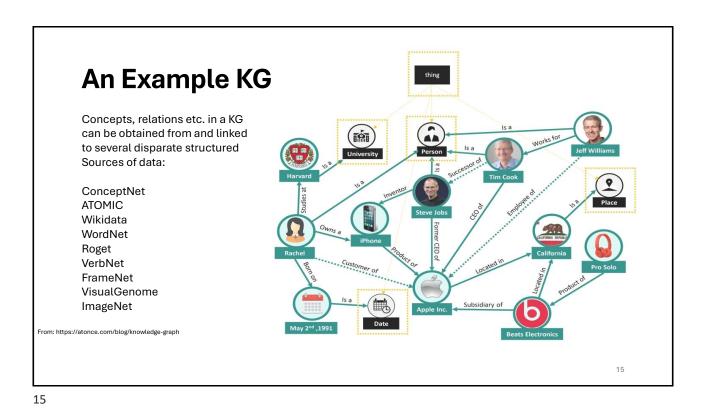
• Wordnet, 1985 (Princeton U.)

A lexical database of semantic relations between words.

Over 150,000 organized in 207,000 word-sense pairs (eat-out, carpool).

Now available for multiple languages.





Commonsense Knowledge

On stage, a woman takes a seat at the piano. She

- 1. sits on a bench as her sister plays with a doll.
- 2. smiles with someone as the music plays.
- 3. is in the crowd, watching the dancers.
- 4. nervously sets her fingers on the keys.

Which one?

Answering the question requires knowledge that humans possess and apply, but machines cannot distill from the communication.

Also, remember Winograd schemas?

16

Review: Meaning Representation Languages

Knowledge Engineering
FOPC
Knowledge Base
Tell-Ask Systems
Forward Chaining
Inference
Backward Chaining
Inference
Definite Clauses
Logic Programming
PROLOG

Frames
Conceptual Dependency
Semantic Networks

5th Gen. Project, SGI, Alvey, ESPRIT
CYC
Wordnet
Knowledge Graphs

Commonsense Knowledge

17

17

Behavioral AI

Subsumption Architecture

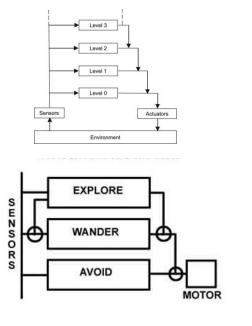
Situatedness Embodiment Intelligence bottom up Emergence

· Layers of control

All layers may have an action to suggest. Only one will be carried out at any time. The action from the "lowest" (highest) module. E.g. AVOID subsumes WANDER. WANDER subsumes EXPLORE.

• Built several robots: Allen, Herbert, Tom & Jerry, Seymour, **Genghis**, Squirt.

Genghis Link: https://www.youtube.com/watch?v=1j6CliOwRng



18

Agent Types

- Table-Driven Agents
 Use a percept-action table in to find next action.
- Simple Reflex Agents
 Based on condition-action rules
- Agents with Memory
 Have internal states that are used to keep track of past states of the world.
- Agents with Goals

 Have state and goal information to take future states into consideration.
- Utility-Based Agents
 Use utility theory to act rationally.
- LLM-Based Agents
 Carry out actions recommended by an LLM.

19

sensors

agent

percepts

actions

effectors

environment

19

Other Successful Approaches

- Rational Agents (Utility based agents)
- Bayesian Inference
- 1997: IBM's Deep Blue beat Garry Kasparov
- 2012: IBM's Watson wins Jeopardy!

Robot Architectures: Vocabulary

Intelligence w/o
Representation
Behavioral AI
Subsumption Architectures
Genghis, Cog
NASA Pathfinder
NASA Spirit & Opportunity
NASA Curiosity, Ingenuity
Boston Dynamics Big Dog
Boston Dynamics Spot, Atlas
Agent-Based AI
Types of Agents
Agentic AI
IBM Watson

21

21

Subsymbolic Al

Symbolic versus Subsymbolic Al

Symbolic Al

Everything is represented using symbols.

A is a block Block(A)

Representation of a state

Expert Systems, Frames, Scripts, Semantic Nets, Knowledge Graphs etc.

Subsymbolic Al

There are NO SYMBOLS.

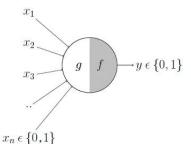
Approaches that employ Neural Networks and other statistical mechanisms

23

23

McCulloch-Pitts Neuron, 1943

- Binary Threshold Units
- Captures the inhibitory and excitatory connections between biological neurons.
- · Limited in what such a model can actually do.
- Missing the learning capability: how to model changes in inhibitory and excitatory connections. This was later included in Hebb's model (1949). Repeated firings can modify the nature of the connections.
- Frank Rosenblatt, 1958 combined these ideas into the model of a Perceptron.



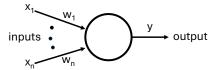
$$g(x_1, x_2, x_3, ..., x_n) = g(\mathbf{x}) = \sum_{i=1}^n x_i$$

$$y = f(g(\mathbf{x})) = 1$$
 if $g(\mathbf{x}) \ge \theta$
= 0 if $g(\mathbf{x}) < \theta$

From: https://towardsdatascience.com/mcculloch-pitts-model-5fdf65ac5dd1

The Perceptron – A Gross Approximation (Rosenblatt, 1958)

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



Transfer Function

T is the Threshold value (assume T = 0)

$$I = \sum_{i=1}^{i=n} w_i x_i$$

$$y = \begin{cases} +1, & \text{if } I \ge T \\ -1, & \text{if } I < T \end{cases}$$

25

25

Perceptron Learning Rule

Changes the weights

$$\overline{\boldsymbol{w}} = [w_1, w_2]$$
 weight vector

$$\overline{x} = [x_1, x_2]$$
 input vector

$$\overline{\boldsymbol{w_{new}}} = \overline{\boldsymbol{w_{old}}} - y^* \, \overline{\boldsymbol{x}}$$
 Training Rule

Perceptron: Vocabulary

Labelled Training Dataset

N samples/patterns/input vector with desired outputs (targets/labels)

Output Error (Loss)

Error = Desired Output - Actual Output

· Learning Rule

Specifies change in the weights using the Error

Prediction/Forward Pass

Application of a pattern to produce output

Epoch

1 pass through the training dataset

27

27

Perceptron Training Algorithm

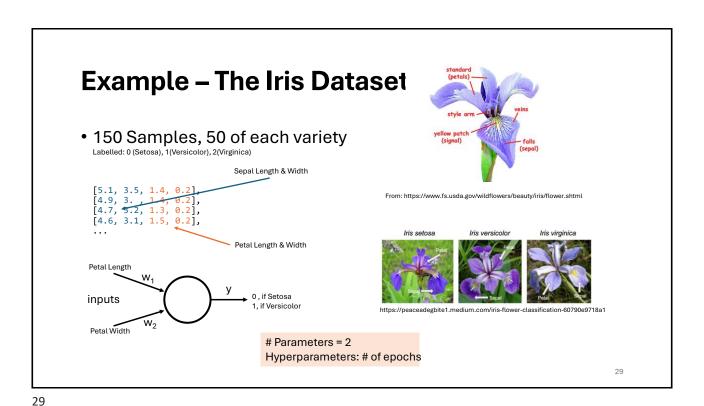
Initialize all weights to random values

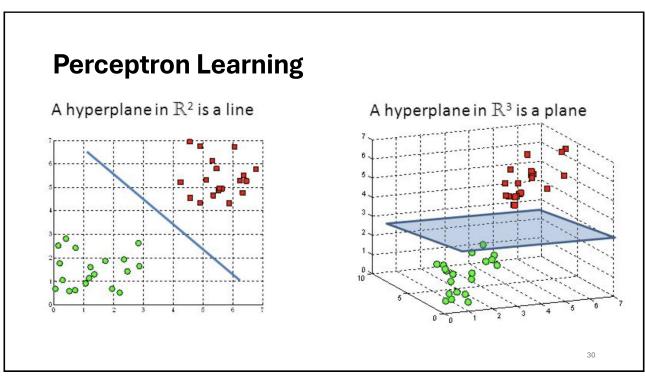
#In what range? Typically [-1.0..1.0]

Set #Epochs to some N

// How to decide what N should be?

Do N times or until all outputs are correct Do for each pattern in the training set apply the pattern to the perceptron change the weight vector as defined





Another Example – MNIST Dataset

- 70,000 images of handwritten digits
- Each image is 28x28 pixels
- Each pixel is in the range [0..255]
 0 = white, 255 = black. Greys in between.
- Training set: 60,000 images
- Testing set: 10,000 images
- Task: Given an image, classify it as [0,1,...,9]

31

31

Limitations of Perceptrons

- · Limited to binary classification tasks only
- · Perceptrons, by Minsky & Papert, 1969

Types of problems Perceptrons could solve were limited to linearly separable problems. Real world problems are not linearly separable.

Perceptron Learning Algorithm would not scale up to tasks requiring large number of weights and thresholds.

For networks with three or more layers there is no obvious way of knowing what the desired output of hidden layers should be.

There is no training procedure possible for networks of three or more layers.

Led to drying up of funding in neural network research in the 1970s.

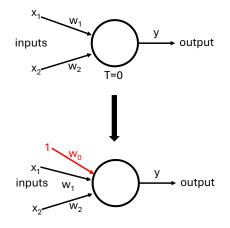
Introducing Bias

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

$$\overline{\mathbf{w}} = [\mathbf{w}_0, \mathbf{w}_1, \mathbf{w}_2]$$
$$\overline{\mathbf{x}} = [\mathbf{1}, \mathbf{x}_1, \mathbf{x}_2]$$

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

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



33

33

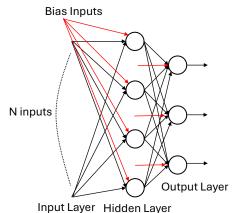
Multi-Layer Perceptron Network

• 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.



Backpropagation Network (Classic Version)

Net Input

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

Activation Function (Sigmoid)

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

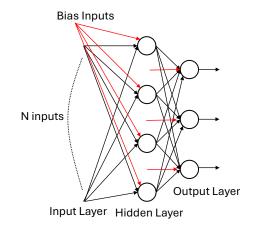
Learning Rule

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

Error/Loss

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

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



35

35

Backpropagation (Classic) Training Algorithm

set minimum ac set #epochs = 0

for each pattern in training set do

do a forward pass for each unit in the hidden layer do compute net input I, and activation, f(I) save f(l) for backpropagation for each unit in the output layer do compute net input I, and activation f(l)

output y = f(I)

do backward pass for each unit in the output layer do

compute error = desired – actual output ($E_i^{\text{output}} = y_j^{\text{desired}} - y_j^{\text{actual}}$) total error = total error + error

for each unit in the **middle layer do**

compute incoming error = weighted sum of output later errors $(\sum_{j=1}^n (w_{ij} E_j^{\text{output}}))$

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 $\beta * \operatorname{error} * f(\mathbf{I})$ and update weight

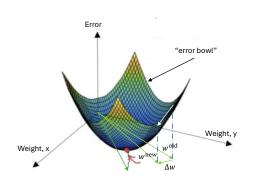
for each unit in the **middle** layer do for each weight from an input layer unit do

compute weight change $\beta*$ final error * f(I) and update weight

until total error < maximum acceptable error or #epochs reaches limit

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.



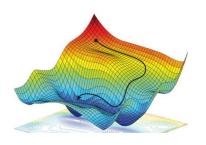
From: https://builtin.com/data-science/gradient-descent

37

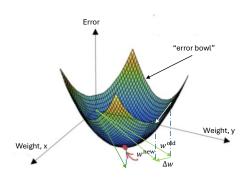
37

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/

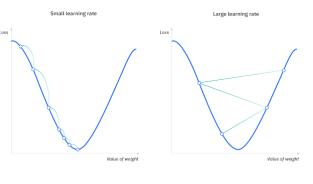


From: https://builtin.com/data-science/gradient-descent

38

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.

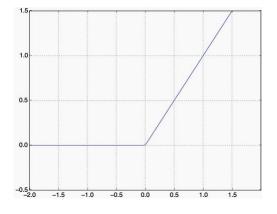


From: https://www.ibm.com/topics/gradient-descent

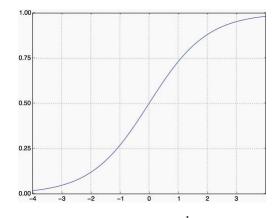
39

39

Popular Activation Functions



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



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

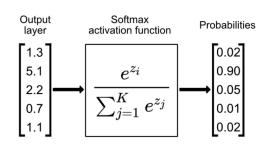
40

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])



41

41

Backpropagation: 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

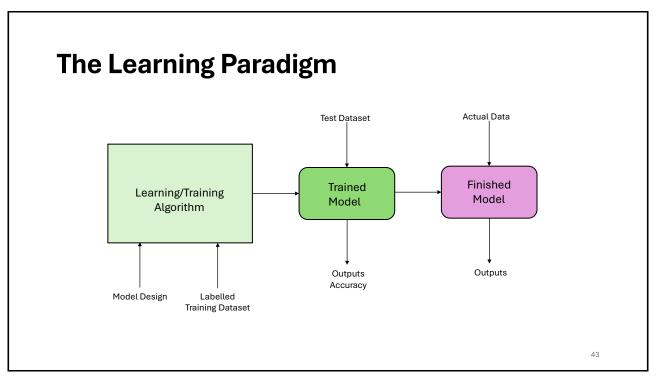
Do the backward pass after all the inputs have been seen, and errors recorded. Full batch ${\bf SGD}$

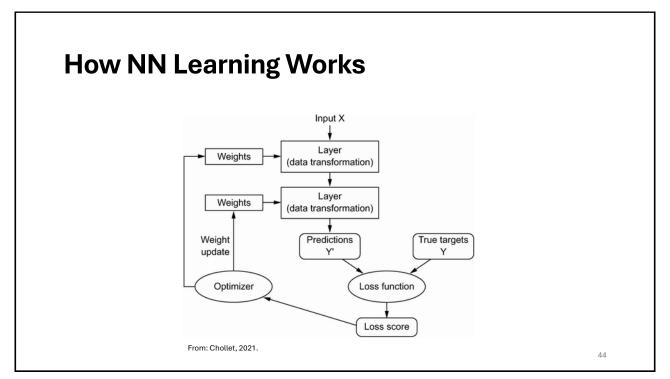
Do a small batch of data and then do a backward pass. Mini Batch SGD (each batch is a power of 2)

- 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.

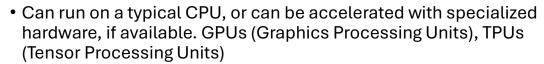
Backpropagation **Binary Cross Entropy** Categorical Cross Entropy Epochs Exponential Forward Pass Full Batch SGD **Gradient Descent** Hyperparameters Labelled Dataset Learning Rule Loss Function Mean Absolute Error Mean Squared Error Mini Batch SGD Optimizer Parameters Relu RMSProp SGD Sigmoid Softmax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





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!

45

From: Chollet, 2021.

Keras

45

Example: Recognizing Handwritten Digits

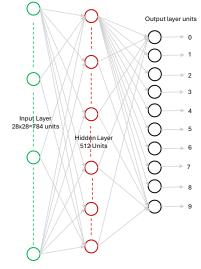
- MNIST Dataset 70,000 images (28x28 pixels), grayscale values (in range 0 (white) to 255 (black).
- Training set: 60,000 images
- Testing set: 10,000 images
- Task: Given an image, classify it as [0,1,...,9]



46

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.
- Parameters 784x512 + 512 (bias) + 512x10 + 10 (Bias) = 407,050
- Hyperparameters
 3 Layers
 784 Units in input layer
 512 units in hidden layer (Sigmoid/Relu)
 10 Units in output layer (Softmax)
 Loss: Mean Squared Error/Categorical Cross-Entropy
 Optimizer: (SGD, RMSProp),
 Learning Rate: β (decided by the optimizer)
 Accuracy Metric: Accuracy (% Correct)
 # epochs (? Start with a max of 10)



47

47

Commonsense Baseline Accuracy

• Before we begin, we should always estimate a baseline accuracy.

That is, given any data set, if we were to randomly assign an output, what would be the accuracy?

For binary classification, the baseline is 50% accuracy.

For MNIST Digits classification, the baseline is 10% accuracy.

 Any neural network model we train should be able to perform far better!

Commonsense Baseline Accuracy

• Before we begin, we should always estimate a baseline accuracy.

That is, given any data set, if we were to randomly assign an output, what would be the accuracy?

For binary classification, the baseline is 50% accuracy.

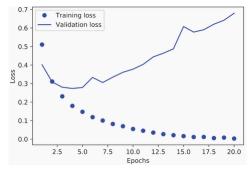
For MNIST Digits classification, the baseline is 10% accuracy.

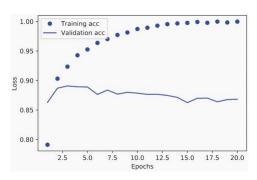
 Any neural network model we train should be able to perform far better!

49

49

Training & Validation Loss and Accuracy



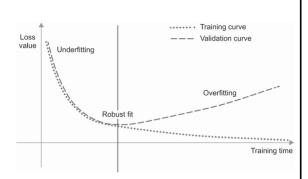


- Notice that the validation loss and validation accuracy both diverge after ~ 4th epoch.
- The model performs better on training data doesn't necessarily do well on the testing data.
- This is called overfitting.

50

Underfitting and Overfitting

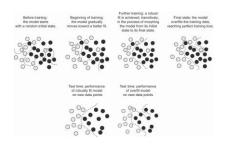
- Initially, as the training proceeds, the lower the loss on training data, the lower the loss on test data. This is underfitting. The network hasn't yet modeled the all the patterns in the training data.
- As training proceeds further, the testing stops improving and starts degrading: This is overfitting. The network is starting to learn patterns specific to the training data.
- Overfitting can occur when the training data is noisy, ambiguous, or involves uncertainty.

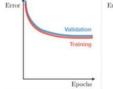


51

51

From a Random Model to Overfitting or Robust Fit





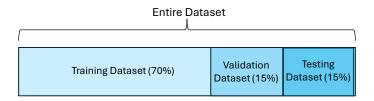




From: Chollet, 2021, and https://medium.com/@datascienceeurope/do-you-know-overfitting-and-underfitting-f27f87ac2f37

52

Training Data, Validation Data, Testing Data



- Training dataset is for use during training
- Validation dataset is to estimate loss/accuracy of the model to tune the hyperparameters
- **Testing dataset** is for evaluating the model after training. (how well does it generalize?)

53

53

Backpropagation: Review

Backpropagation Bias Binary Cross Entropy Categorical Cross Entropy Commonsense Baseline Accuracy Epochs Exponential Forward Pass Full Batch SGD **Gradient Descent** Hyperparameters Keras Labelled Dataset Learning Rule Loss Function Mean Absolute Error Mean Squared Error Mini Batch SGD Model

Optimizer Overfitting Parameters Relu RMSProp Scikit-Learn SGD Sigmoid Softmax Tanh **Testing Data** Training Accuracy Training Data Training Loss True SGD Underfitting Validation Accuracy Validation Data Validation Loss

54