

Chapter 1: What is Deep Learning

AI = effort to automate intellectual tasks normally performed by humans

Symbolic AI: a sufficiently large set of explicit rules for manipulating knowledge stored in explicit databases
↳ dominant paradigm 1950 - 1980

Although the Turing Test has sometimes been interpreted as a literal test Turing merely meant it as a conceptual device in a philosophical discussion about the nature of cognition

ML → a very hands-on field driven by empirical findings and deeply reliant on advances in software and hardware

"Deep" → from deep learning means successive layers of representation

Loss function = objective function = cost function

Backpropagation algorithm

There were already 2 AI winters → 1960, 1990

Probabilistic ML → application of the principles of statistic to data analysis

Ex. Naive Bayes (it assumes that all the data points are independent)
Logistic regression

Kernel Methods → classification algorithms best known for SVMs
1) The data is mapped to a new high-dimensional

representation where decision boundaries can be expressed as hyperplanes

2) Find the decision boundary by maximizing the distance between the hyperplane and the closest datapoints from each class

The thing is that it can be computationally expensive to map all the points to a higher dimensional space
⇒ kernel trick allows you to find the distance between two points from the initial space into the target space

Decision Trees → flowchart-like structure that lets you classify stuff
⇒ Random forest.
⇒ Gradient Boosting Machines

Deep learning → one of the reasons why it is so cool is because it automates feature engineering

Best libraries according to Kaggle Teams:

Keras, LightGBM, XGBoost, PyTorch, TensorFlow, Sklearn, FastAI, Caffe

2009 - 2010
↳ better activation functions
↳ better optimization schemes (RMSProp & Adam)

2014 - 2016
↳ improved gradient descent
↳ batch normalization, skip connections

Chapter 2: Mathematical building blocks of neural networks

Softmax classification layer → return an array of 10 probability scores (summing to 1)

Optimizer → the mechanism through which the model will update itself based on the training data it sees, so as to improve its performance

Loss function → how the model will measure its performance on the training data, and how it will be able to steer itself in the right direction

metrics to monitor during training and testing

tensor = container for data

number \rightarrow rank 0 tensor

vector rank-1 tensor

matrix rank-2 tensor

Tensor $\begin{cases} \rightarrow \text{number of axes (rank)} \\ \rightarrow \text{shape: tuple of integers that describes how many dimensions the tensor has along each axis} \\ \rightarrow \text{datatype (or dtype)} \end{cases}$

Deep learning \rightarrow usually process data in batches (chunks)

Common data tensors

1) Vector data
usually rank 2 tensors (samples, features)

2) Time series data
usually rank 3 tensors (samples, timestamp, features)

3) Image tensors (samples, H, W, colors)

(128, 256, 256, 3)
 $\begin{array}{c} \xrightarrow{\text{Samples}} \quad \uparrow \quad \uparrow \\ \text{height} \quad \text{width} \end{array} \quad \brace{3 \text{ colors}}_{\text{RGB}}$

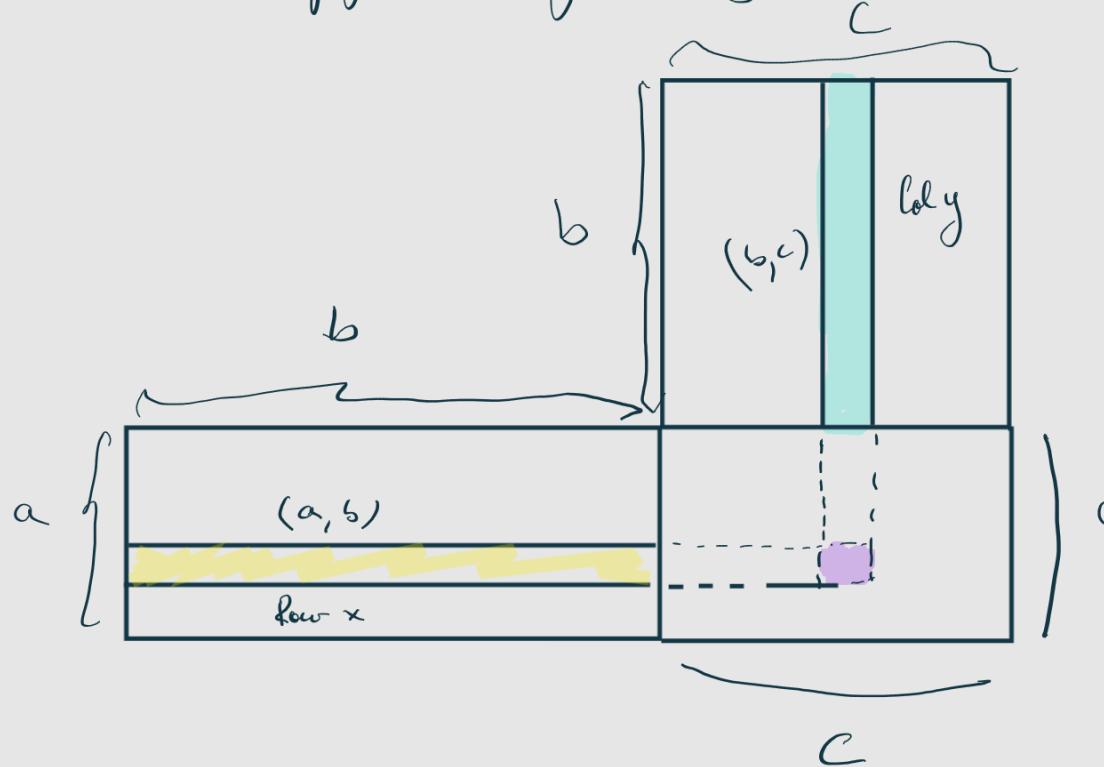
4) Video data
(samples, frames, H, W, colors)

Keras layers Dense (512, activation = "relu")

\Rightarrow identical output = $\text{relu}(\text{dot}(\text{input}, w) + b)$

Element wise operations \rightarrow naive = step by step on each

element
python style using tensors

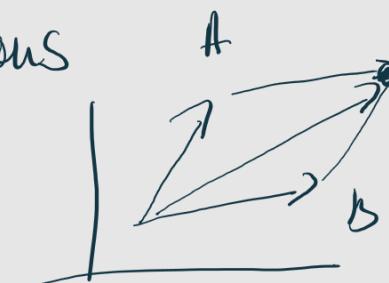


⇒ Reshaping Tensors

Geometric interpretations

1) addition

= translating an object

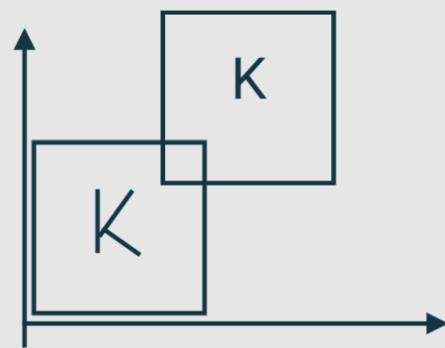


$A + B$

$$\begin{bmatrix} \text{Horizontal} \\ \text{Vertical} \end{bmatrix} + \begin{bmatrix} x \\ y \end{bmatrix}$$

2) Translation

= apply addition to a set of points



3) Rotation

⇒ counter clockwise rotation

by angle θ ⇒ via dot product

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix}$$

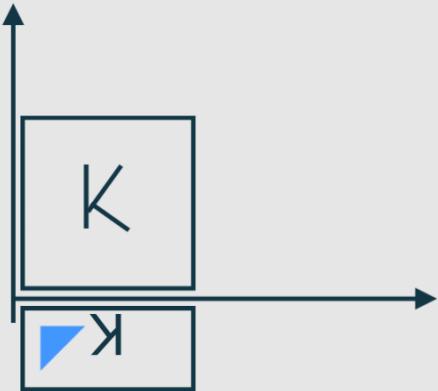
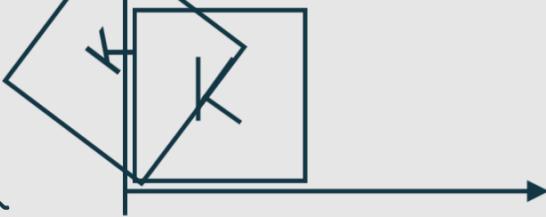
4) Scaling



→ vertical & horizontal

→ with a dot product

with a diagonal matrix



$$\begin{bmatrix} 1 & 0 \\ 0 & 0.5 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix}$$

5) Affine transformations

→ combination of linear transformation and translation

$$y = W \cdot x + b$$

This is what a dense layer is performing. A dense layer without an activation function is an affine transformation.

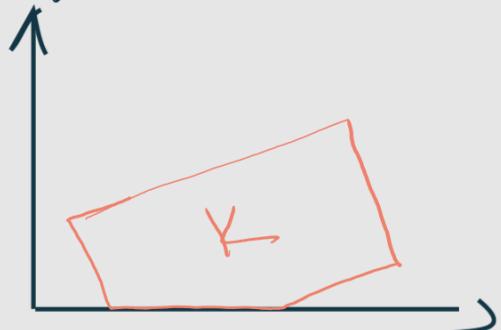
A dense layer with ReLU activation

$$\text{affine}_1(\text{affine}_2(x)) \xrightarrow{\text{Equivalent}} \text{affine}_3(x)$$

where affine_3 can be computed

⇒ no need for multiple layers

affine + ReLU



Gradient based optimisation

→ weight

→ bias

Training loop

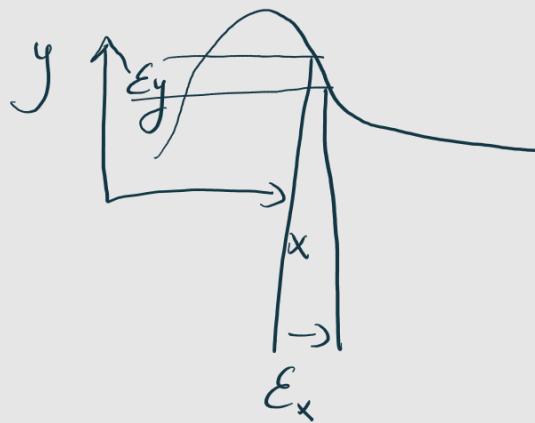
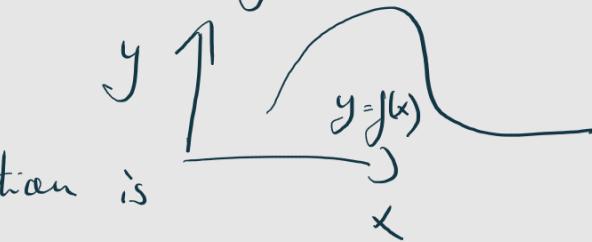
1) Draw a batch of training samples

2) Forward pass ⇒ obtain predictions

- 3) Compute the loss of the model on the batch
- 4) Update the weights of the model in a way that slightly reduces the loss of the batch

$$f(x) = y$$

because the function is continuous



$$f(x + \epsilon_x) = y + a\epsilon_x$$

The slope a = derivative

$$a < 0 \Rightarrow x \uparrow, f(x) \downarrow$$

$a > 0 \Rightarrow x \uparrow, f(x) \uparrow$
magnitude of a \Rightarrow how much

\geq how abrupt is it increases/
the slope decrease the
function

$$y_{-pred} = \text{dot}(w, x)$$

$$\text{loss-value} = \text{loss}(y_{-pred}, y_{-true})$$

\Rightarrow how do we update w to make the loss value smaller?

$$\text{loss-value} = f(w)$$

we can compute partial derivatives with respect to some values

$\text{grad}(\text{loss-value}, w_0)[i, j] \Rightarrow$ indicates the direction and magnitude of change when modifying $w[i, j]$

\Rightarrow you move w in the opposite direction of the gradient

$$w_1 = w_0 - \text{step} * \text{grad}(f(w_0), w_0)$$

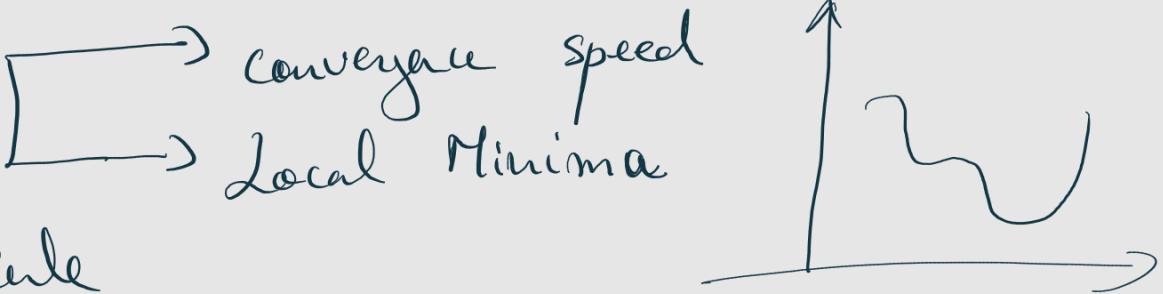
- 1) Draw a batch of training examples
- 2) Obtain predictions
- 3) Compute the loss
- 4) Back propagation (compute the gradient of the loss with regards to the model's parameters)
- 5) Move the parameters a little in the opposite direction for the gradient

True SGD → one sample

mini-batch → with a mini batch

every step on all data \Rightarrow batch gradient descent

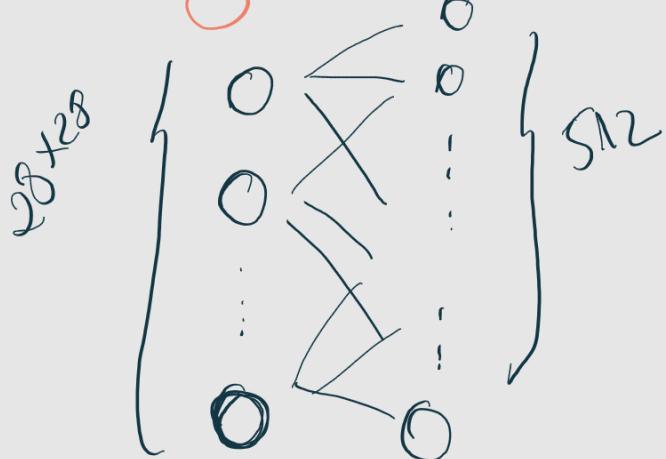
Momentum



The chain Rule

$$\begin{aligned} \text{grad}(y, x) &= \text{grad}(y, x_3) * \text{grad}(x_3, x_2) \\ &\quad * \text{grad}(x_2, x_1) * \text{grad}(x_1, x) \end{aligned}$$

Reimplementing MNIST



Naive Dense

\hookrightarrow 1 layer

output = activation

$$(\text{dot}(W, \text{input}) + b)$$

w

b

$$\begin{array}{c}
 \text{input} \\
 \left(x_1 - \dots - x_n \right) \\
 \underbrace{\hspace{10em}}_{\text{input}}
 \end{array}
 \quad
 \left\{
 \begin{array}{c}
 \left(\begin{array}{c} w_1 - \dots - w_m \\ \vdots \\ \vdots \\ v_{1n} - \dots - v_{mn} \end{array} \right) \\
 \text{output}
 \end{array} \right\}
 \quad
 \left\{
 \begin{array}{c}
 \left(\begin{array}{c} b_0 \\ \vdots \\ b_m \end{array} \right) \\
 \text{output}
 \end{array} \right\}$$

\downarrow

$w.\text{shape}(\text{input}, \text{output})$

