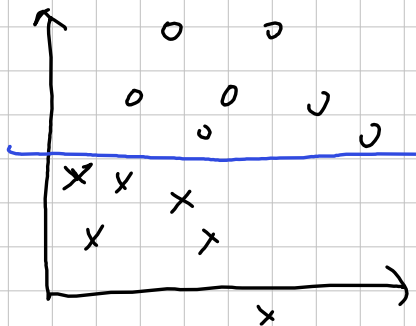
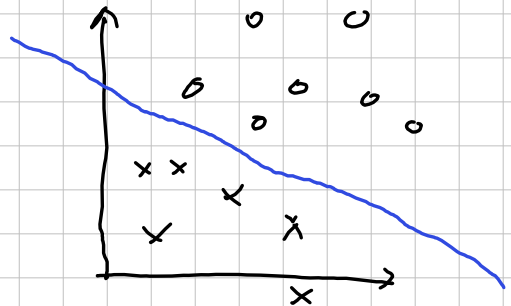


Machine Learning Notes

SVMs (Support Vectors Machine)



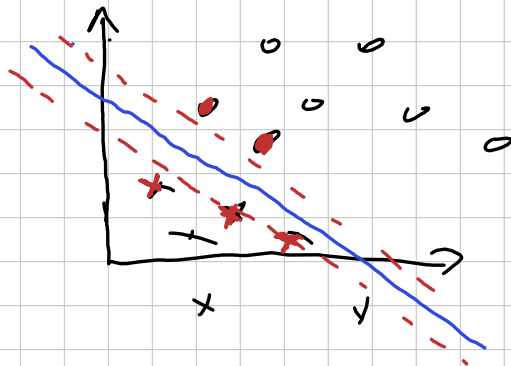
Correct ✓
Best ✗



Correct ✓
Best ✓

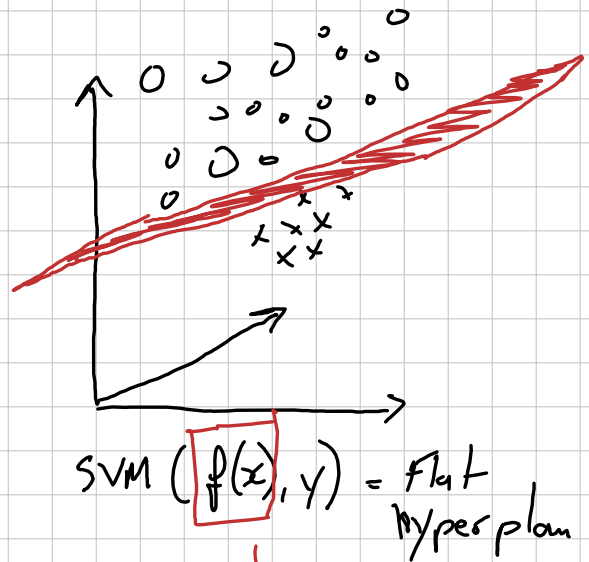
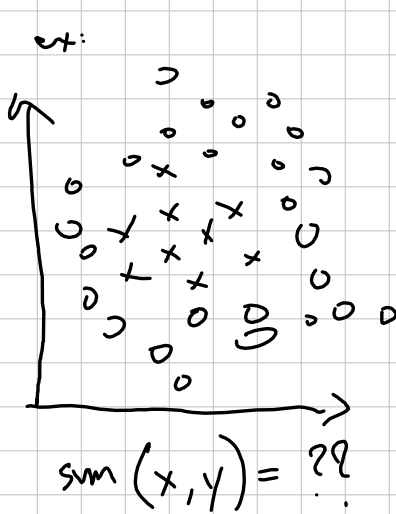
- Support vectors Machine tries to maximise the margin between classes.

↳ The closest points to the hyperplane are the ones responsible of supporting it



SVMs / Kernel Tricks

Limitation: most ^{row} data can not be separated by just a hyperplane



Non Linear Transformation of x ↙

- ↳ problems:
1. Select the correct Transformation
 2. Sophisticated boundaries
 - ↳ increase dimension / output of the transformation
 - ↳ increase computational requirements.

↳ Solution: The Kernel Trick

↙
The idea: instead of caring about $x \rightarrow f(x) \rightarrow x'$

we care about the relation between x, x' comparing them with each other

⇒ in Math it's done using the dot product

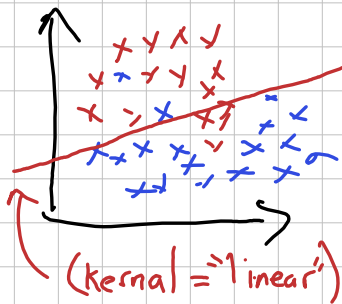
SVMs / Kernel Trick

⇒ in Math it's done using the dot product
also called inner product

$$f(x)^T f(x')$$

⇔ The kernel function $k(x, x')$

Simple Example: The identity Transformation
(we do nothing about x)

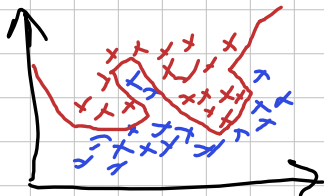


$$\hookrightarrow f(x) = x$$

$$\hookrightarrow k(x, x') = x^T x' \quad (\text{linear Kernel})$$

$$[\text{skl SVC(kernel="linear").fit(x,y)}]$$

Complex Example: The polynomial Transformation



$$\hookrightarrow f(x) = (x_1, x_2, x_1 x_2, x_1^2, x_2^2)$$

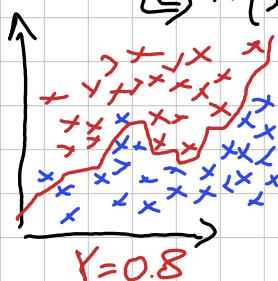
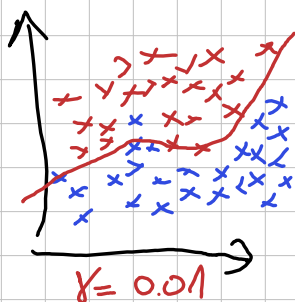
$$\hookrightarrow k(x, x') = (1 + x^T x')^2 \quad (\text{polynomial Kernel})$$

$$[\text{skl SVC(kernel="poly").fit(x,y)}]$$

∞ Complex Example: Radial Basis Function
 $f(x) =$ infinite dimensional!!

$$\hookrightarrow k(x, x') = e^{-\gamma \|x - x'\|^2}$$

Radial Basis Function
Kernel



$$[\text{skl SVC(kernel="rbf", gamma=0.01).fit(x,y)}]$$

Naive Bayes

- ↳ Do a naive assumption(s) about features in the dataset to make it easier to predict/classify

Example:

FACTS

Spann
probability
if:

Probability

• 100% probability = 25
 • No spam probability
 $\Rightarrow \frac{5}{25} = 20\%$
 • Spam = $\frac{20}{25} = 80\%$

• No spam = $\frac{10}{25} = 40\%$
 • Spam = $\frac{15}{25} = 60\%$

• No spam = $\frac{0}{15} = 0\%$
 • Spam = $\frac{15}{15} = 100\%$

based on the dataset
 the model gives 100%
 which is not realistic
 for cases outside the dataset

Selected Dataset
(100 email)

Solution?

- ↳ Use bigger dataset?
- ↳ Guess the number

Naïve assumption that "Buy" and "Cheap" are independent

We can multiply $p_{\text{spam}}(\text{Buy}) \times p_{\text{spam}}(\text{Cheap})$

$$\hookrightarrow \frac{20}{25} \times \frac{15}{25} = \frac{4}{5} \times \frac{3}{5} = \frac{12}{25}$$

↳ Actual number of emails that are spam and contain buy cheap

$$\frac{12}{25} \times 25 = 12 \text{ email}$$

We can multiply $P_{nos}(Buy) \times P_{nos}(Cheap)$

$$\Rightarrow \frac{5}{75} \times \frac{10}{75} = \frac{1}{15} \times \frac{2}{15} = \frac{2}{225}$$

↳ Actual number of pencils that are not spanned and contains buy & cheap

$$\frac{2}{925} \times 75 = \frac{2}{3} \text{ Email}$$

225
→ New Dataset based on Naive assumptions

12

2/3

Spam

No Span

 \llcorner

Spam Probability
"Buy" & "Cheap"

- $\text{span} = \frac{12}{12 + \frac{2}{3}} = \frac{36}{38} = 94,737\%$

$$\therefore \text{No span} = \frac{2/3}{12 + 4/3} = \frac{2}{38} = \boxed{5,263\%}$$

Darstellung

$$P(S|B) = \frac{P(B|S)P(S)}{P(B|S)P(S) + P(B|NS)P(NS)}$$

$$= \frac{\frac{20}{25} \times \frac{25}{100}}{\frac{20}{25} \times \frac{25}{100} + \frac{5}{75} \times \frac{75}{100}} = 80\%$$

NS: Not spam
 B: Bury
 C: Cheap

S: Spam
NS: Not spam
B: Buy
C: Cheap

Naïve Bayes

$$P(B \cap C) = P(B) \times P(C)$$

Naive

$$\Rightarrow P(S|BNC) = \frac{P(B|S)P(C|S)P(S)}{P(B|S)P(C|S)P(S) + P(B|NS)P(C|NS)P(NS)}$$
$$= \frac{20/25 \times 15/25 \times 25/100}{20/25 \times 15/25 \times 25/100 + 5/25 \times 10/25 \times 75/100}$$
$$= \boxed{94,737\%}$$

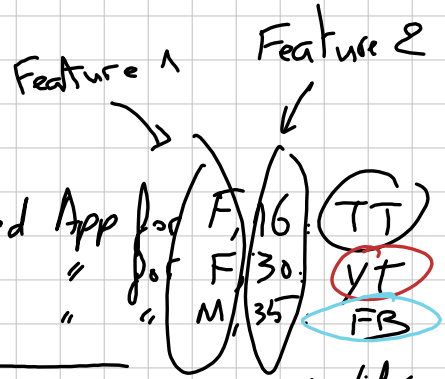
Decision Trees

Main idea:

Example: Existing dataset

Gender	Age	App
F	15	TT
F	25	YT
M	32	FB
F	35	YT
M	12	TT
M	14	TT

Recommended App for



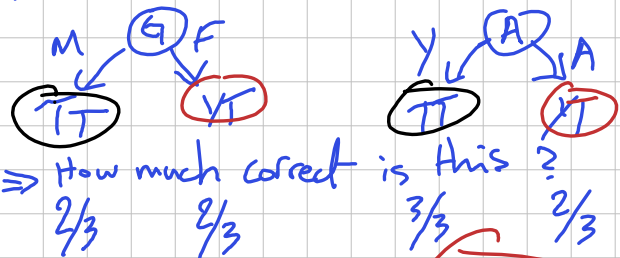
Gender	Age	App
F	Y	TT
F	A	YT
M	A	FB
F	A	YT
M	Y	TT
M	Y	TT

Gender decision stump

Age decision stump



⇒ What would we recommend based on this??



⇒ How much correct is this?

$\frac{2}{3}$

$\frac{2}{3}$

$\frac{3}{3}$

$\frac{2}{3}$

are correct

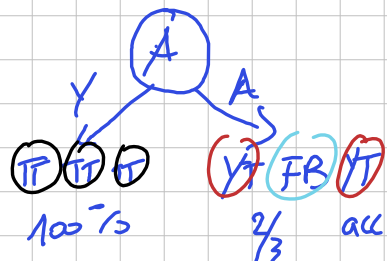
$\frac{4}{6}$

$\frac{5}{6}$

Age decision stump wins!

Accuracy ↑

⇒

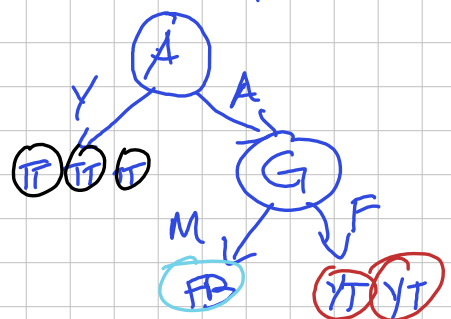


100%

$\frac{2}{3}$

accurate

Could it be improved? ⇌

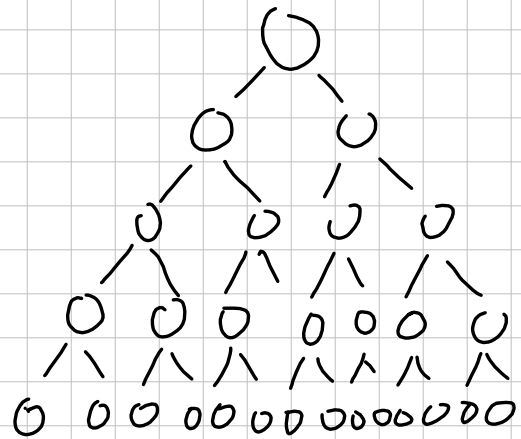


Decision Trees

Example 2: Multiple features

Gender	Age	Country	Likes	History	Subscriptions	App
✓	✓	✓	✓	✓	✓	✓
✓	✓	✓	✓	✓	✓	✓
✓	✓	✓	✓	✓	✓	✓
✓	✓	✓	✓	✓	✓	✓

Deep Tree
= High chance
of overfitting



In example 1 Precision was used to build the decision tree

How to decide what is the best decision stump?

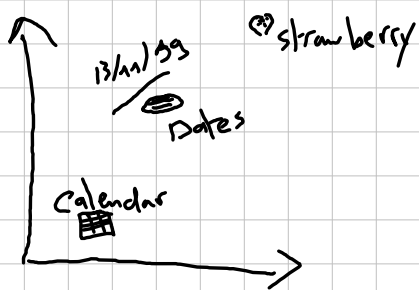
- ↳ Precision
- ↳ Gini Impurity Index
- ↳ Information Gain & Shannon Entropy



Avg Gini

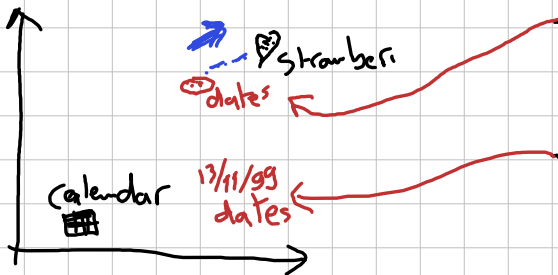
Attention Mechanism

The Challenge



- I bought strawberries & dates
- Metallica announced new dates for their world tour

The solution

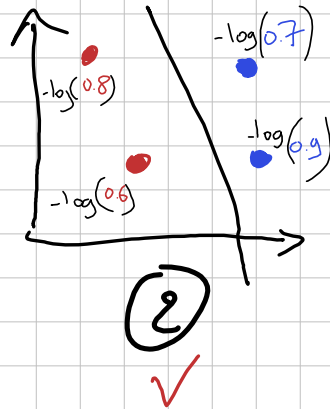
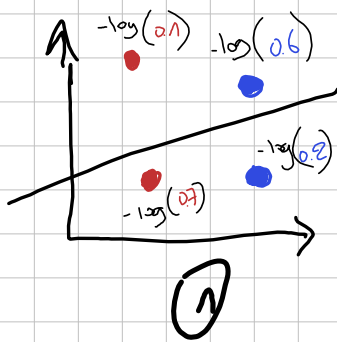


- I bought Strawberries & dates.
- Metallica announced new dates for their world tour.

Neural Networks NN

Principal: minimise error

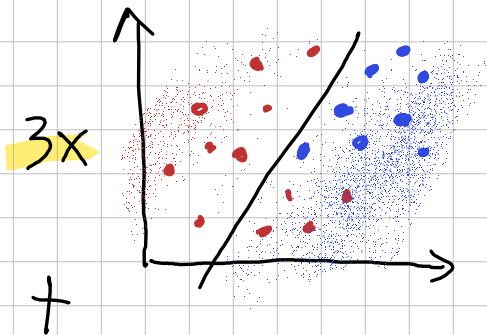
(logistic Regression)



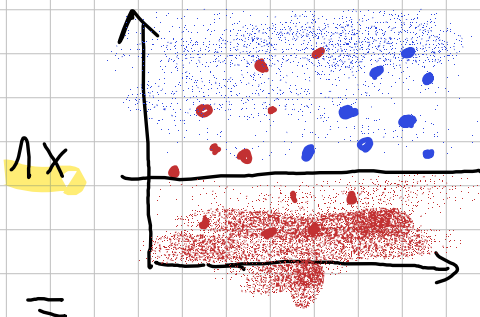
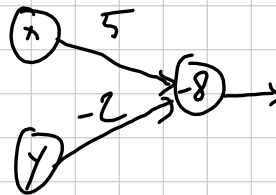
① $-\log(\hat{y}_0) - \log(\hat{y}_4) = 4.8$

② $-\log(\hat{y}_0) - \log(\hat{y}_4) = \underline{\underline{1.2}}$

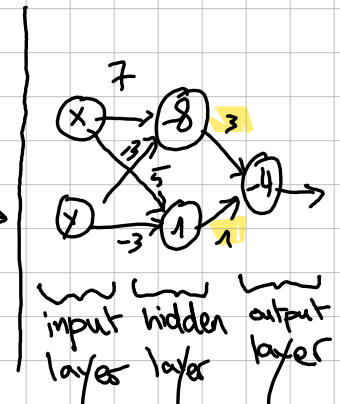
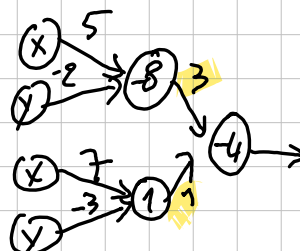
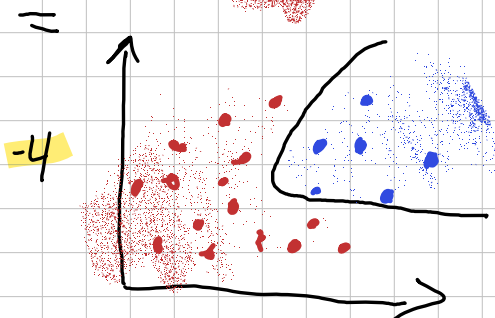
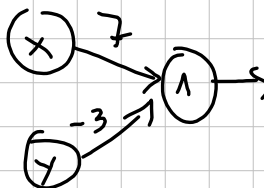
Neurons



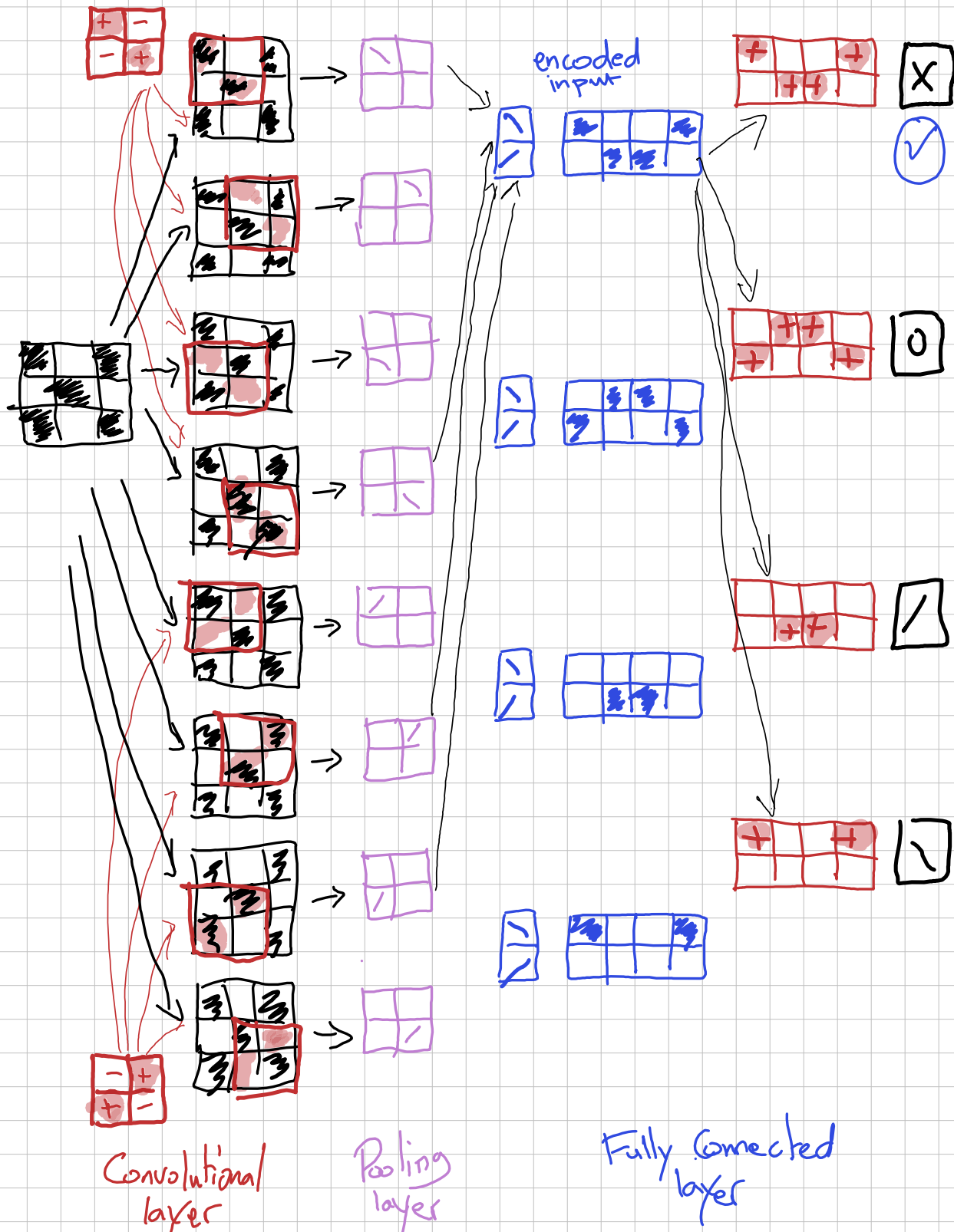
$$5x - 2y = -8$$



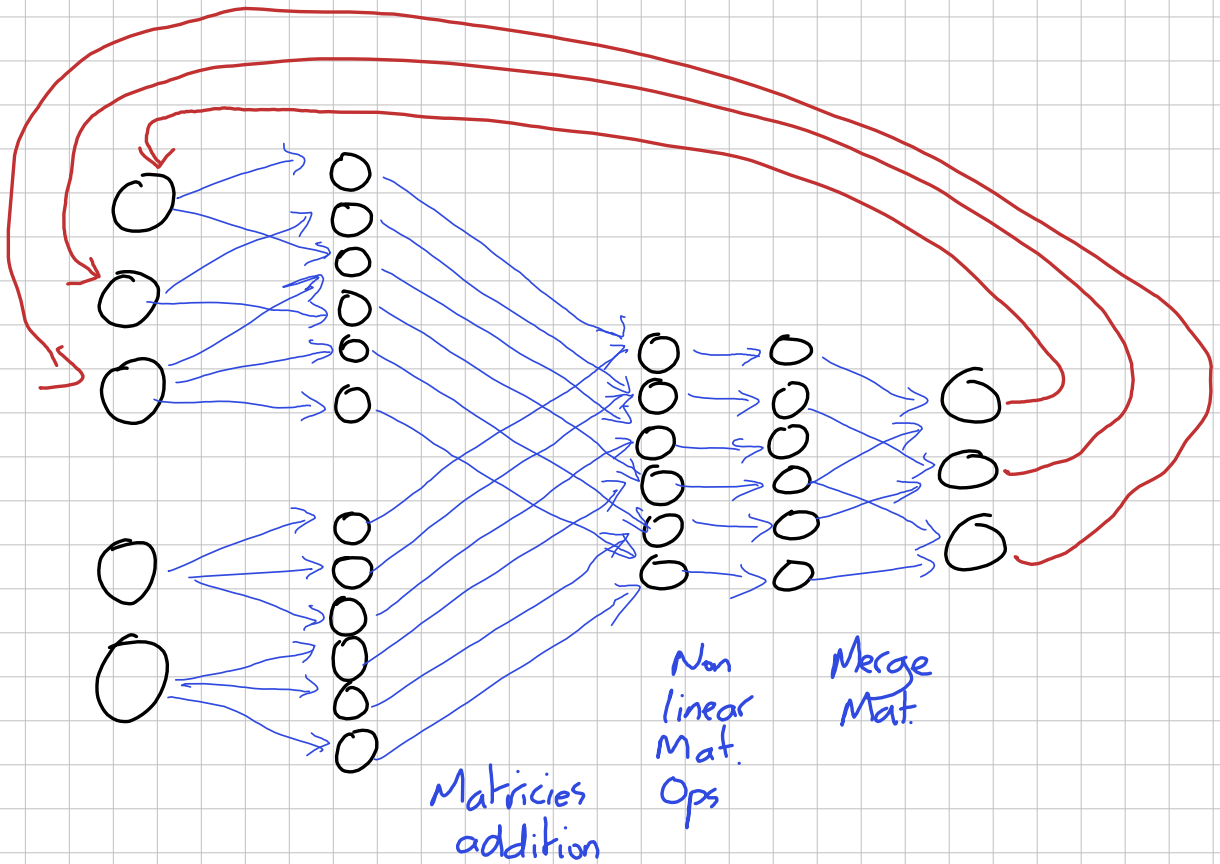
$$7x - 3y = 1$$



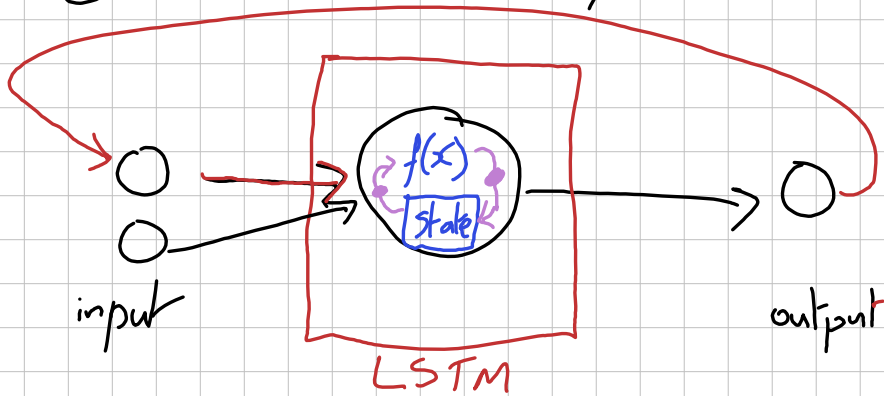
Convolutional Network



Recurrent Neural Networks




Long Short-term Memory Networks

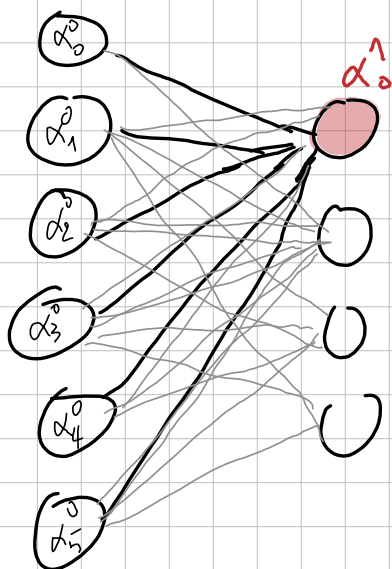


good in learning sequences of data

Activation functions

- Sigmoid (classic)

$$\sigma = \frac{1}{1 + e^{-x}}$$


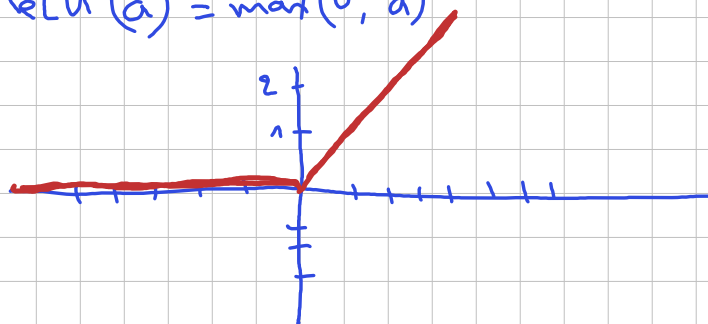


$$x_0^1 = \sigma(w_{0,0}x_0^0 + w_{0,1}x_1^0 + \dots + w_{0,n}x_n^0 + b_0)$$

$$\sigma \left(\begin{bmatrix} w_{0,0} & w_{0,1} & w_{0,2} & \dots & w_{0,n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{n,0} & w_{n,1} & w_{n,2} & \dots & w_{n,n} \end{bmatrix} \begin{bmatrix} x_0^0 \\ x_1^0 \\ x_2^0 \\ \vdots \\ x_n^0 \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix} \right)$$
$$x_0^1 = \sigma(Wa^0 + b)$$

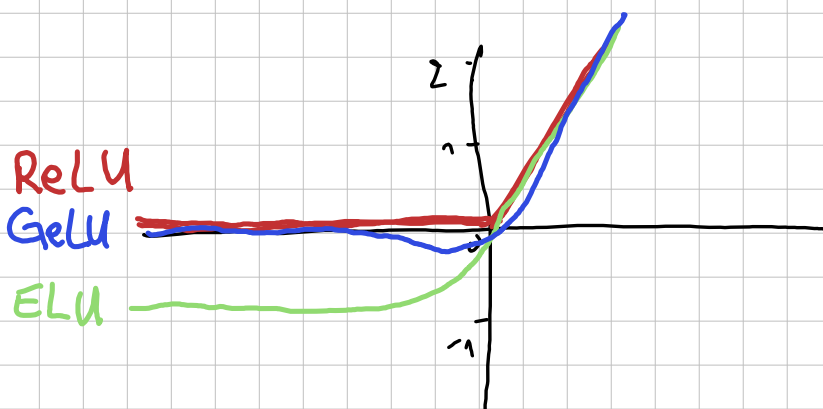
- ReLU (Rectified Linear Unit)

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

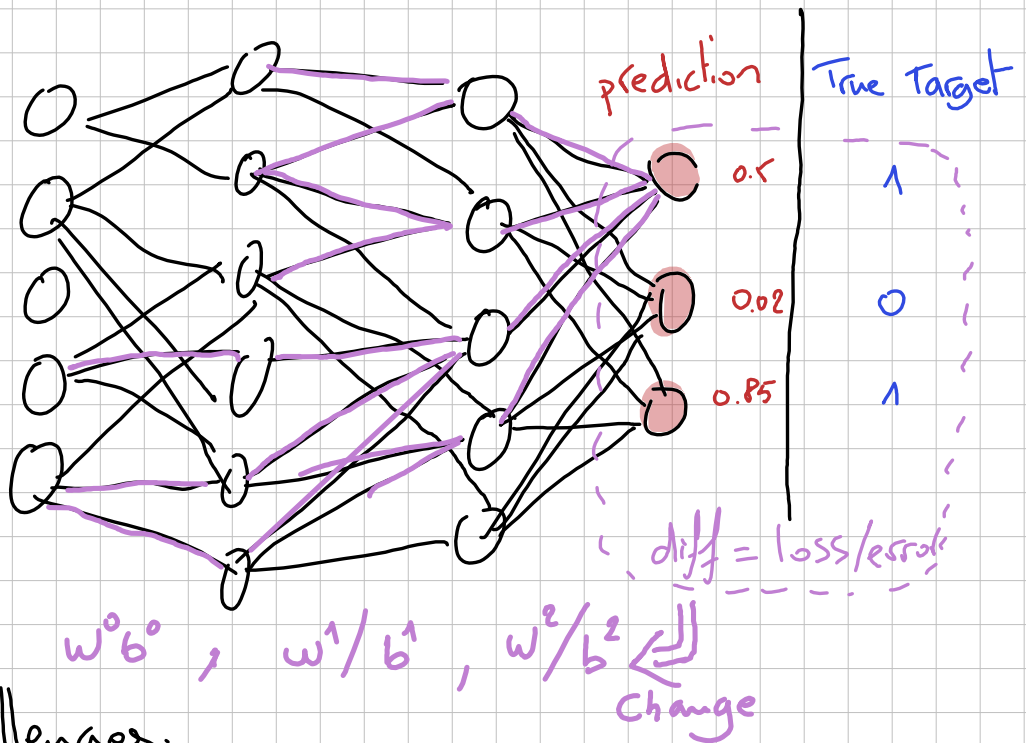


Motivated by biological analogy
↳ ReLU well for very deep N.N.

- GeLU (Gaussian Error Linear Unit)



Backpropagation



Challenges:

- Exploding Gradients
 - ↳ due to accumulation from deep network
 - ↳ Solution: gradient clipping.
- Vanishing Gradients
 - ↳ due to very small weights in earlier layers
 - ↳ solution: residual connections (add input layer to its output) + layer normalisation.
- Space/Time Complexity
 - ↳ Stochastic Gradient (batches of loss function)

Transformer Architecture

