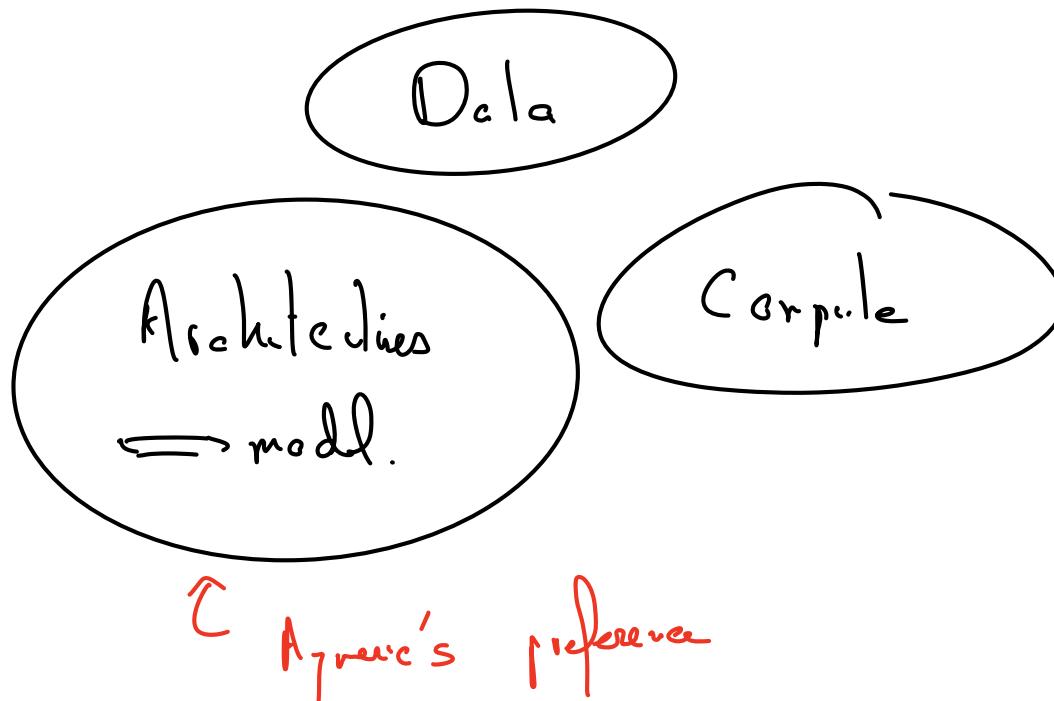


Deep Learning - Part 1

Aymeric DIEULEVEUT

DS4M - April 2025

1 Goals

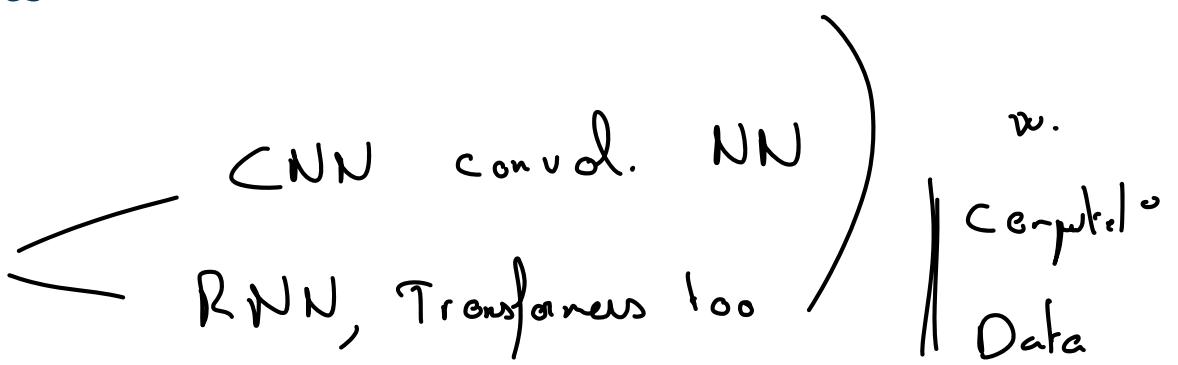


2 The first Neural Network: the Perceptron

- Logistic regression
- Discovering Fully connected NN through the NNPG.
- Formalization, advanced topics
- Further discussion

1/3 Polarized but drove the progress

3. Other Architectures



4. Apply and discuss

Outline

1 Goals

2 The first Neural Network: the Perceptron

- Logistic regression
- Discovering Fully connected NN through the NNPG.
- Formalization, advanced topics
- Further discussion

Deep Learning in one slide + Goals of the afternoon

- How does it work:

- 1 | ▶ Automatically learn representations of observations
- 1 | ▶ Learn highly non-linear models.

♡ not crucial aspect

understood

- What does it require:

- 2 | ▶ Large datasets with structure ✓
- 3 | ▶ Computational power ✓

+ architecture

- Why now:

- ▶ Combination of the 3 points above
- ▶ investment !

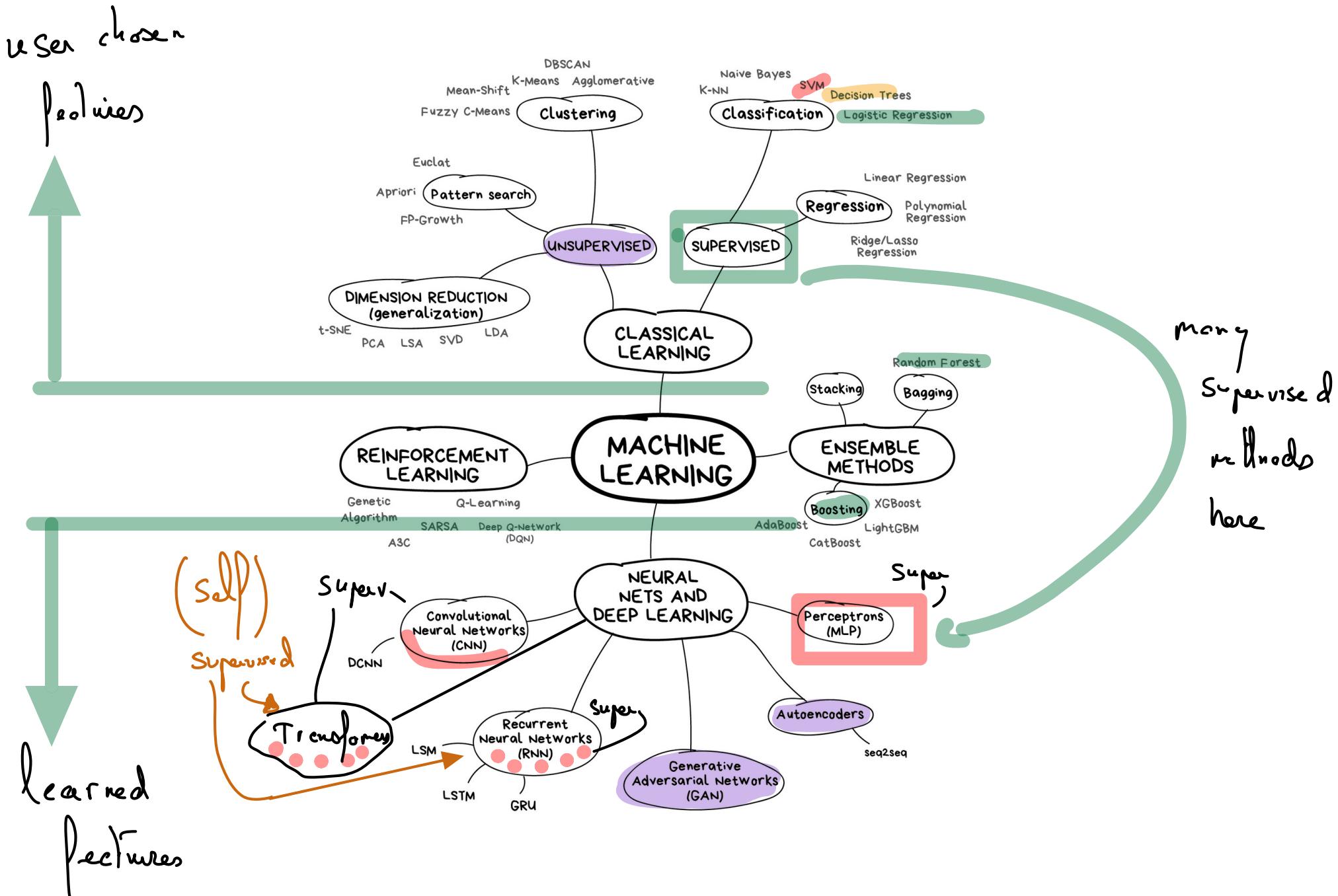
- Some Applications

- ▶ Image classification; object / face recognition
- ▶ Self driving cars
- ▶ Automatic Translation, Information extraction
- ▶ Caption Generation
- ▶ Ads, recommendation systems, etc.

Goals: Understanding core concepts of Deep Learning.

- When and how it works on paper (data, models, architecture)
- How to implement a simple neural network with Python.
- Overview of some of the main applications and challenges.

Machine Learning

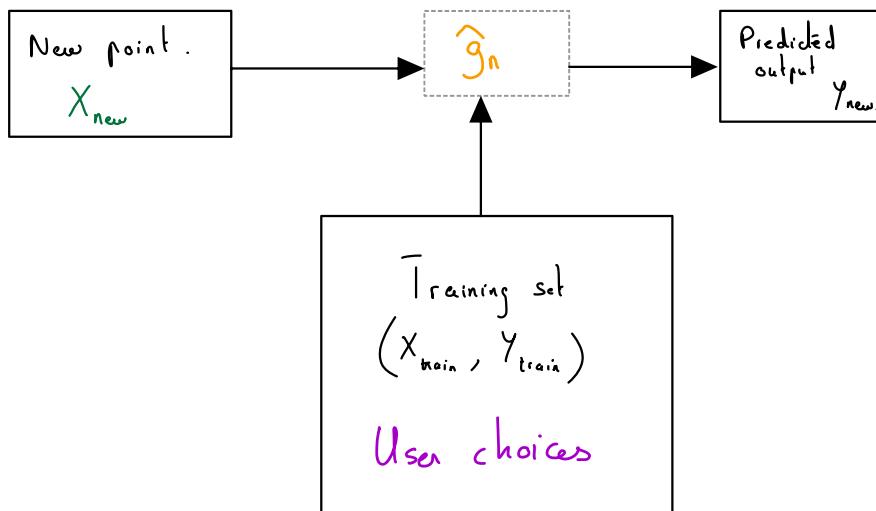


Summary of previous lectures on Machine Learning

Important points

Inductive reasoning : Data \rightarrow prediction rule
Generalization
Overfitting

① General framework: loss, risk, generalization risk vs training risk



Loss function, Generalization Risk and data-dependent predictors

Loss Function

- Loss function: $\ell(y, g(x))$ quantifies the quality of the prediction $g(x)$ of y , e.g.:
 - ▶ 0-1 Loss (classification): $\ell(y, g(x)) = 1_{y \neq g(x)}$, $y \in \mathcal{Y} = \{-1, 1\}$
 - ▶ Quadratic Loss (regression): $\ell(y, g(x)) = \|y - g(x)\|^2$, $y \in \mathbb{R}^d$

Risk of a Decision Rule = average loss

- Risk: $R(g) = \mathbb{E}[\ell(Y, g(x))] = \int \ell(y, g(x)) d\mathbb{P}(x, y)$, e.g.:
 - ▶ 0-1 Risk (classification): $R(g) = \mathbb{P}(Y \neq g(x))$
 - ▶ Quadratic Risk (regression): $R(g) = \mathbb{E}[\|Y - g(x)\|^2]$
- ↪ R is also referred to as *Generalization risk*, or *True risk*.

⚠ The distribution \mathbb{P} is unknown.

Data-dependent predictors

- Training set: $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ (i.i.d. $\sim \mathbb{P}$)
- Goal: build a **data dependent predictor / classifier** \hat{g}_n based on the training data
- That has a **minimal generalization risk** $R(\hat{g}_n)$.

$R(\hat{g}_n)$ = average loss on a "new point" $(X, Y) \sim \mathbb{P}$, independent from D_n

↪ Our goal is to predict well on inputs/outputs pairs that we have not seen in the dataset, i.e. *generalize well*.

THM 1.

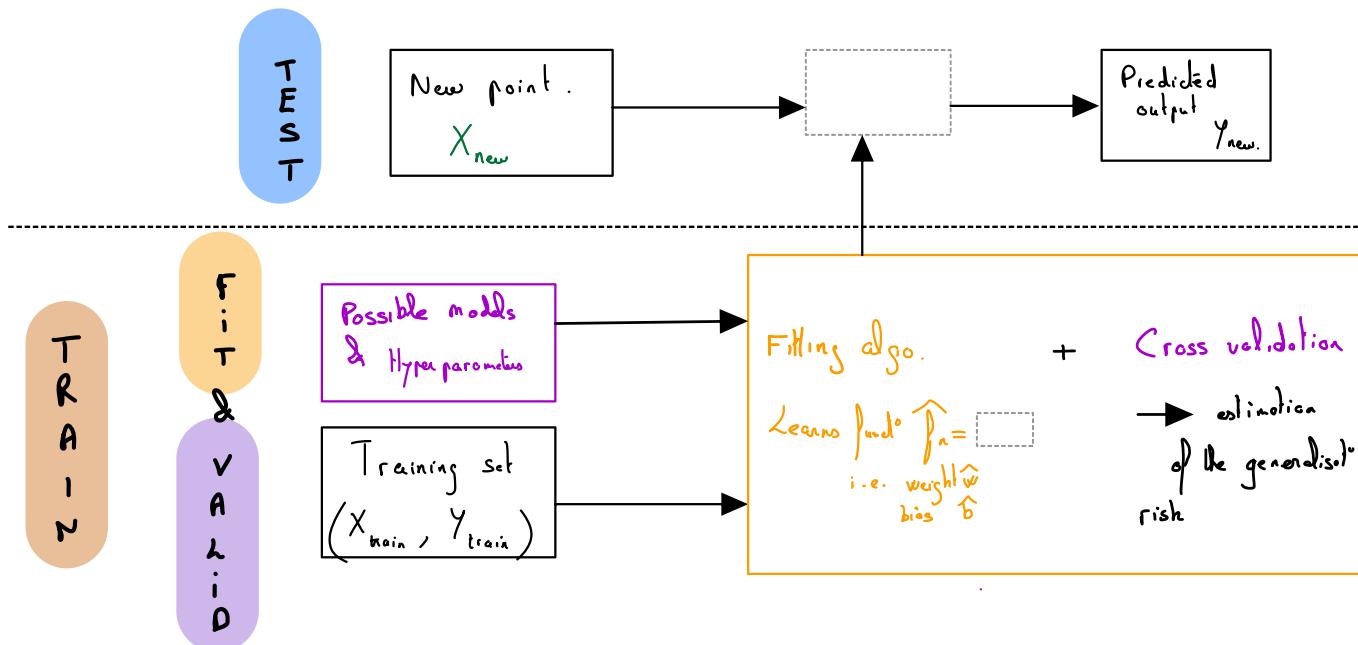
Summary of previous lectures on Machine Learning

Key concept : Cross validation and its lesson to be

① General framework: loss, risk, generalization risk vs training risk

② SML Pipeline

- ▶ Train (Fit + validation) / Test / Deployment
- ▶ Evaluation of validation scores through cross validation to pick the best method/hyperparams



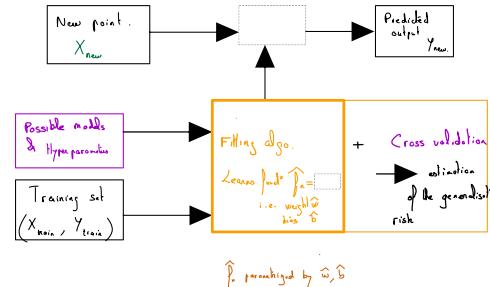
Summary of previous lectures on Machine Learning

① General framework: loss, risk, generalization risk vs training risk

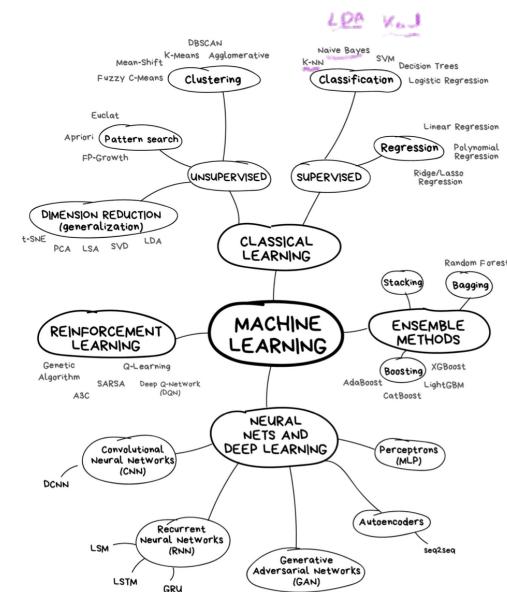
② SML Pipeline

- ▶ Train (Fit + validation) / Test / Deployment
- ▶ Evaluation of validation scores through cross validation to pick the best method/hyperparams

③ Methods: Linear, Logistic regression, Decision Trees, Random Forests, Boosting



Fitting algo.
Learn f s.t. $\hat{Y}_n = \boxed{\quad}$
i.e. weight \hat{w} bias \hat{b}



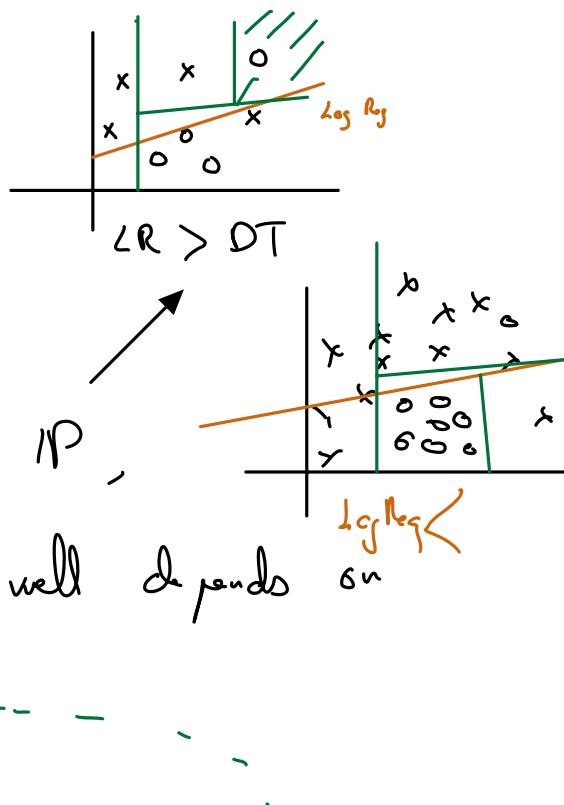
Summary of previous lectures on Machine Learning

- ① General framework: loss, risk, generalization risk vs training risk
- ② SML Pipeline
 - ▶ Train (Fit + validation) / Test / Deployment
 - ▶ Evaluation of validation scores through cross validation to pick the best method/hyperparams
- ③ Methods: Linear, Logistic regression, Decision Trees, Random Forests, Boosting
- ④ Opening the black box: for each method
 - ▶ How to predict (\leftrightarrow interpretability)
 - ▶ How fit is made (\leftrightarrow speed, complexity) [Computer learning parameters]
 - ▶ Important Hyper-parameters [User chosen]
 - ▶ Drawbacks and Strengths
 - ▶ Underfitting, overfitting risks (\leftrightarrow Performance)
 - regularization techniques (avoiding overfitting by modifying the method)

$\hat{Y}_n =$	Lin. Reg	Polynomial Linear reg	Logistic Reg	Dec Trees	RF/ Boosting
How to predict	y_{pred} $w_1 X_1 + w_2 X_2 + \dots + w_d X_d$	Linear model. It's lin reg on poly. $T \mapsto 1, T, T^2, T^3$	Linear model for classif $\text{Sign}(w^T X)$	Yes	
How to fit	incremental improvement SGD (or closed formula)		II	one split by or split	not II
Key hyperparameters	None yes if regularized (sparsity) LASSO		Same as LR	depth ...	# trees subsampling, diversity
⊕	simplicity / speed stability / (interpre)	more powerful (non-linear rel)	Same as LR	speed non linear interpret	more stable less overfitting!
⊖	Accuracy (linear relat)		S = LR	overfitting less interpret	slower.
Overfitting / Underfitting → regularization techniq	better: often underfit possibly overfit if degree too high	make higher degree $80 \rightarrow 6400$ → overfitting	S = LR	yes	not too much
Approach	ERM.	ERM	ERM	not ERM	not ERM

e.g. House Price predict dataset, with $d=80$
 polynomial features of degree 2 create more than $80^2 = 6400$ features
 leads to Boosting.

Which is the best model?



→ no complete answer

- * even for a give data desribl^o IP,
the ability of a model to perform well depends on
the nbr of obsrvd^o.
- * The best model depends on IP (! and n)
- * The best model depends on Business perspective
(loss is , goal, $\frac{S}{P} / \frac{I}{C} / A$)
Speed
Accuracy
Precision
+

Outline

1 Goals

2 The first Neural Network: the Perceptron

- Logistic regression
- Discovering Fully connected NN through the NNPG.
- Formalization, advanced topics
- Further discussion

Outline

1 Goals

2 The first Neural Network: the Perceptron

- Logistic regression
- Discovering Fully connected NN through the NNPG.
- Formalization, advanced topics
- Further discussion

Supervised ML problem:

- ① Data $D_n = (X_i, Y_i)_{i=1,\dots,n}, X_i \in \mathbb{R}^d, Y_i \in \{-1, 1\}$
- ② Goal: Predict a the label for a new point.

Empirical risk minimization (Discriminative approach)

- No statistical model!
- ~~choose~~ Define a loss function ℓ
- Choose a function class \mathcal{F}



- Goal:

$$\min_{f \in \mathcal{F}} \mathbb{E}[\ell(f(X), Y)]$$

- ERM

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)]$$

We minimize the empirical risk.

\hookrightarrow we have an opt prob →

we use an opt algo to find a soln = SGD = "IJ"

loss

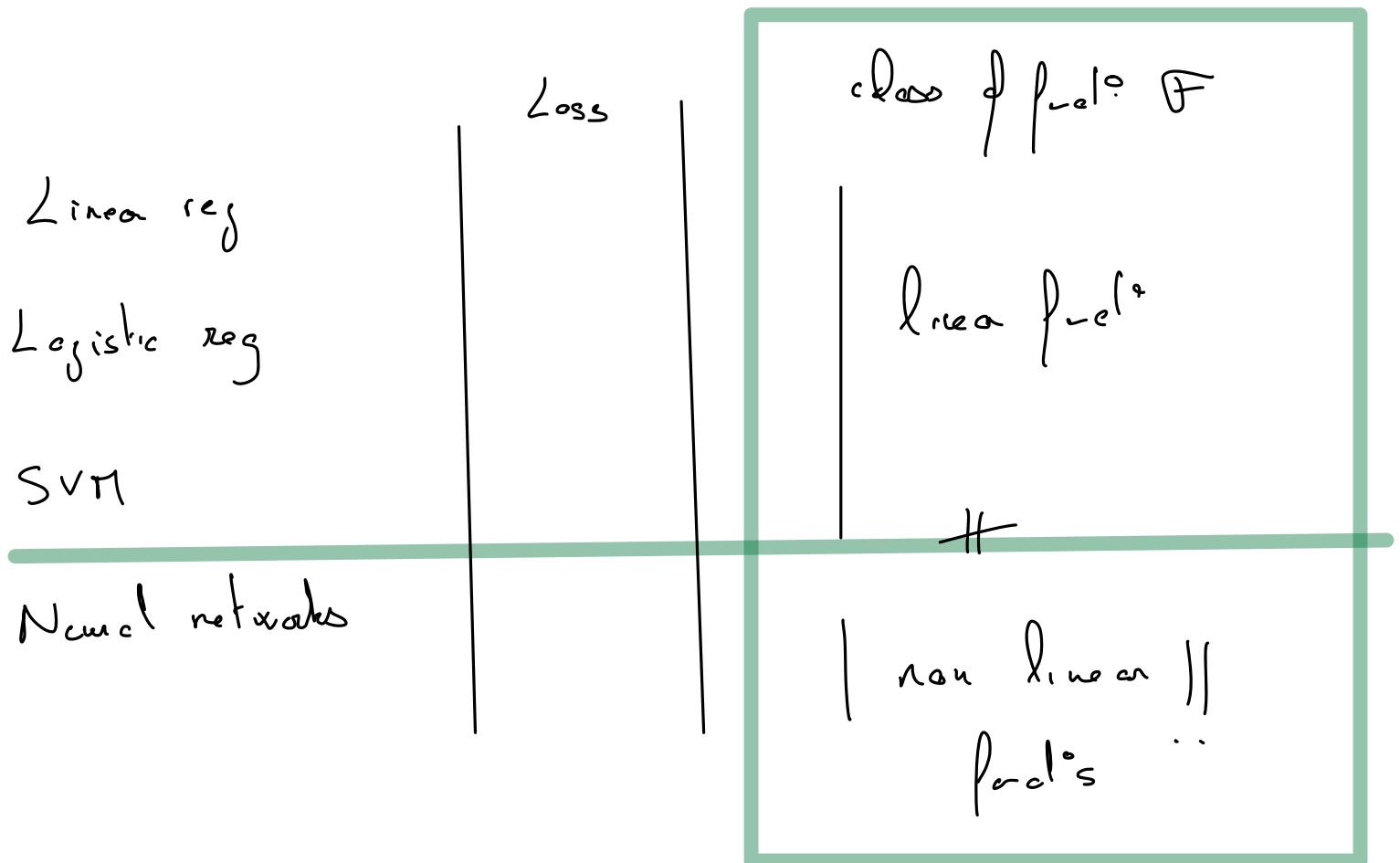
$$\ell(Y_i, f_\omega(X_i))$$

$$\frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_\omega(X_i))$$

average loss on fit

empirical risk

def.

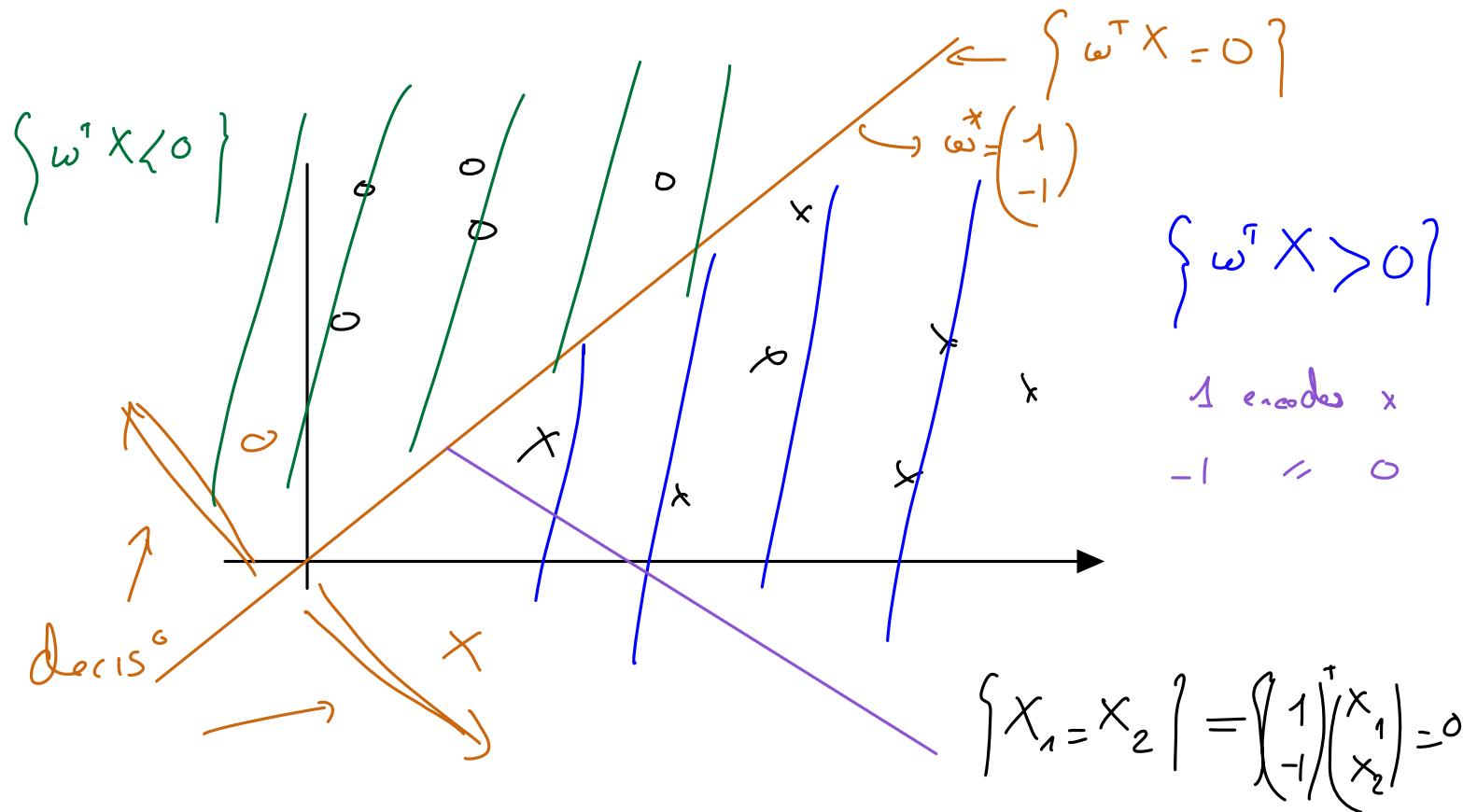


Linear models

Class of *soft* linear separators:

$$\mathcal{F}_{lin} = \{g_w : X \mapsto w^\top X \quad , w \in \mathbb{R}^d\}$$

Binary prediction from a soft predictor: $X \mapsto \text{sign}(w^\top X)$



Linear models

Class of *soft* linear separators:

$$\mathcal{F}_{lin} = \{g_w : X \mapsto w^T X \quad , w \in \mathbb{R}^d\}$$

01 loss \leftrightarrow accuracy

Binary prediction from a soft predictor: $X \mapsto \text{sign}(w^T X)$

truth
my predict

(Pb 1)

01 loss:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{Y_i \neq \text{sign}(w^T X_i)} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{-Y_i w^T X_i > 0}$$

Logistic regression loss

$$\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-Y_i w^T X_i))$$

l_{logistic} loss

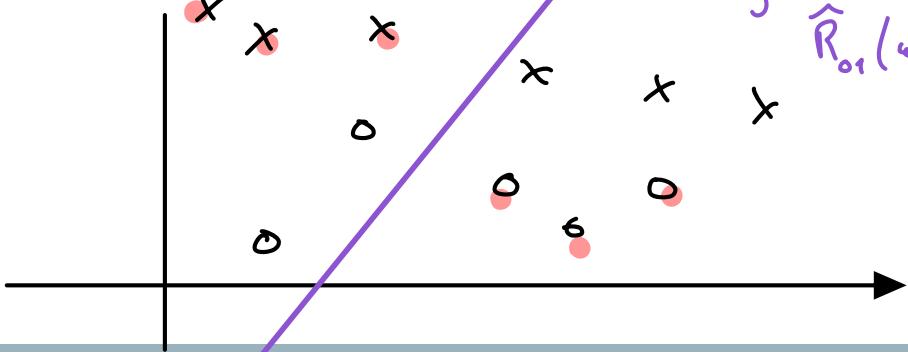
why not solve ERN for
01 loss

$$\min_w \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i \neq \text{sign}(w^T X_i)\}}$$

I could solve that, this would
be a great method

this is my initial
guess. $w^{(0)}$

$$\hat{R}_{01}(w^{(0)}) = \frac{6}{11}$$



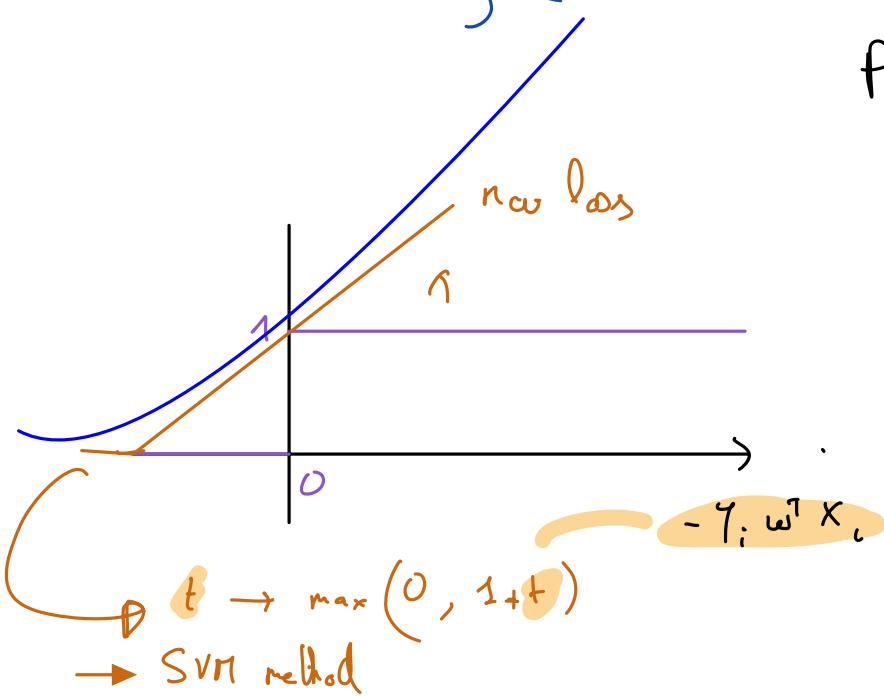
01 loss:

$$\text{Loss of } 1 \text{ if } -y_i \omega^\top x_i > 0$$

$$t = \log(1 + e^+)$$

01 Loss | is not almost everywhere
non differentiable!
non cvx

Pb1 is IMPOSSIBLE
to solve algorithmically
(or exponential complexity
in d).



Logistic regression summary

- * I want to use linear func
- * The 0-1 loss is $\left\{ \begin{array}{l} \text{bad} \\ \text{impossible} \end{array} \right\}$ to optimize
- * Change the 0-1 loss to logistic loss
→ get logistic regress!

Outline

1 Goals

2 The first Neural Network: the Perceptron

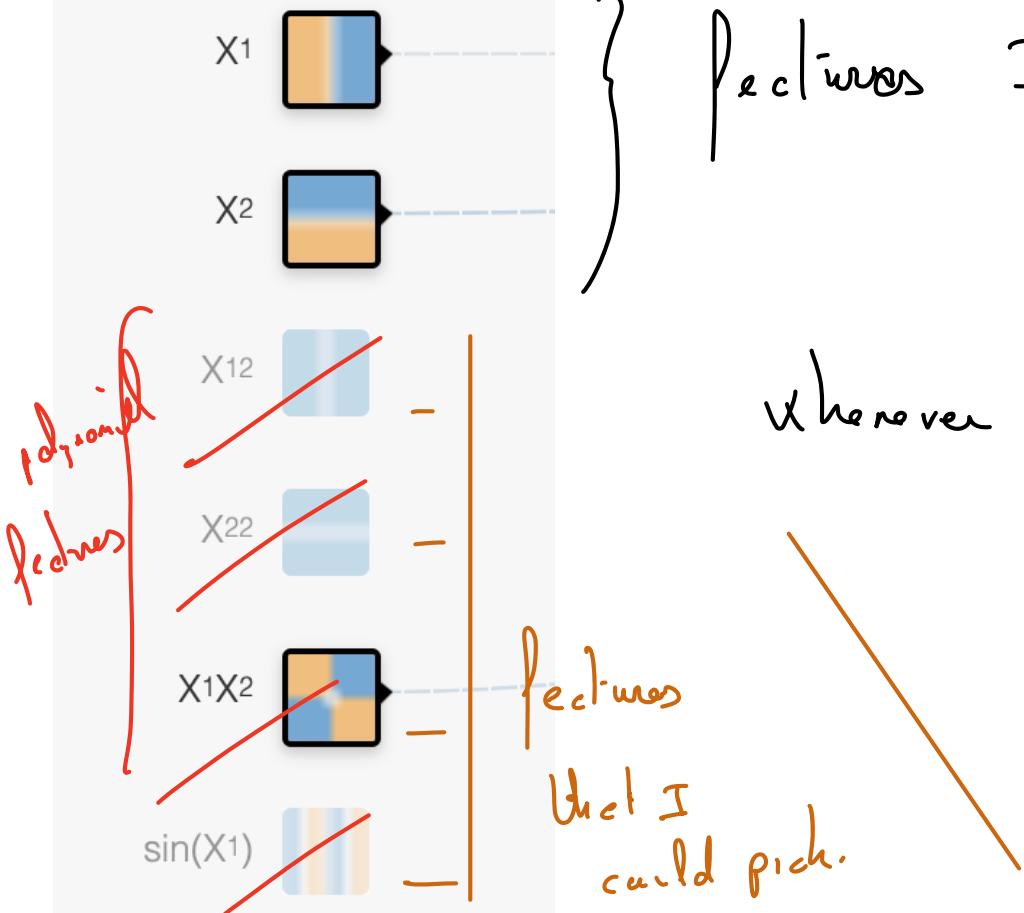
- Logistic regression
- Discovering Fully connected NN through the NNPG.
- Formalization, advanced topics
- Further discussion

Let's NN playground: [Link](#)

1. Understand what the features are. Understand what colors represent.

FEATURES

Which properties do you want to feed in?



lectures I use .

$x^1 \rightarrow$ the x axis
value

$x^2 \rightarrow$ the y axis
value .

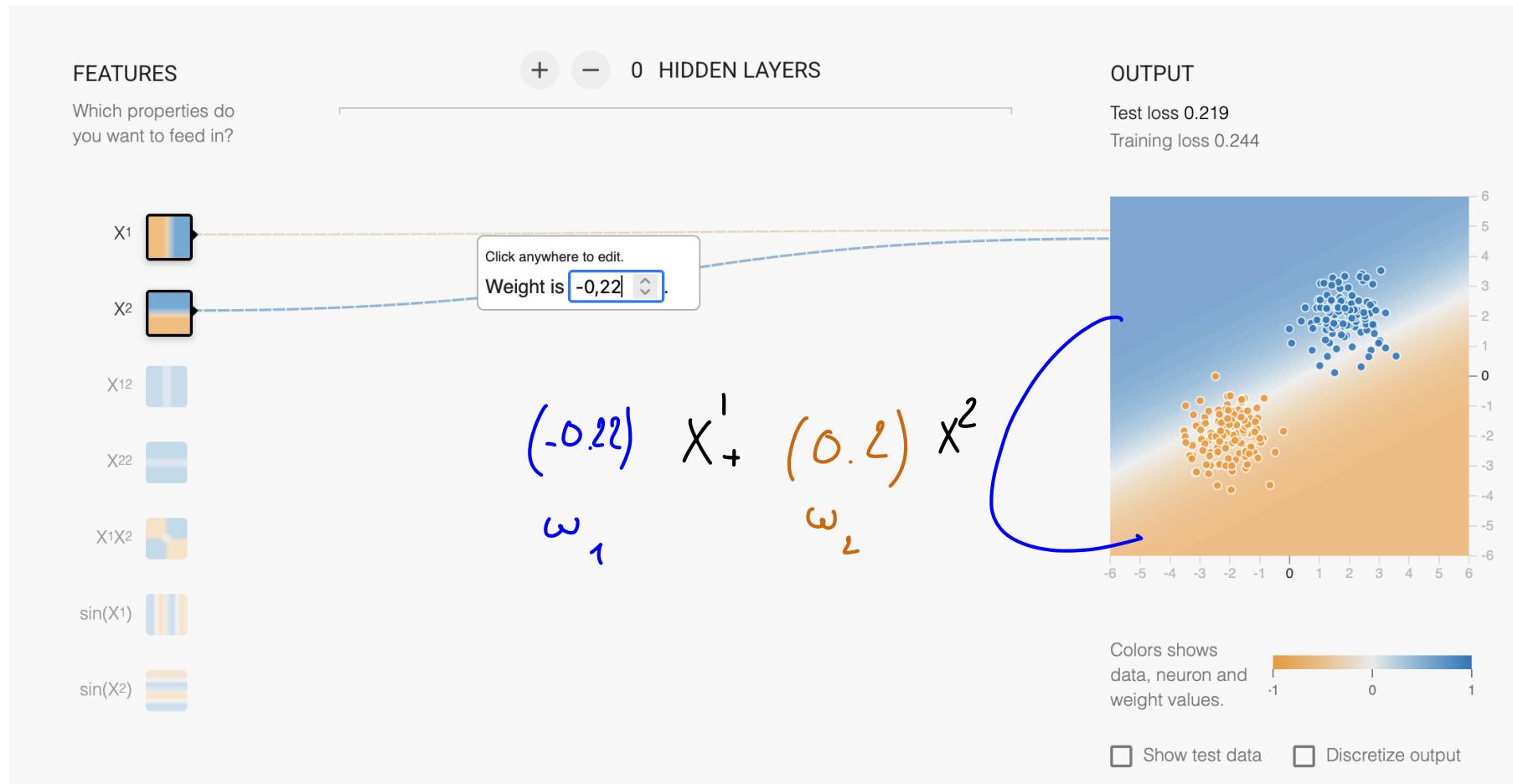
Whenever you have over a square-

I am restricting myself to $x^1 x^2$

that I
could pick.

Let's NN playground: [Link](#)

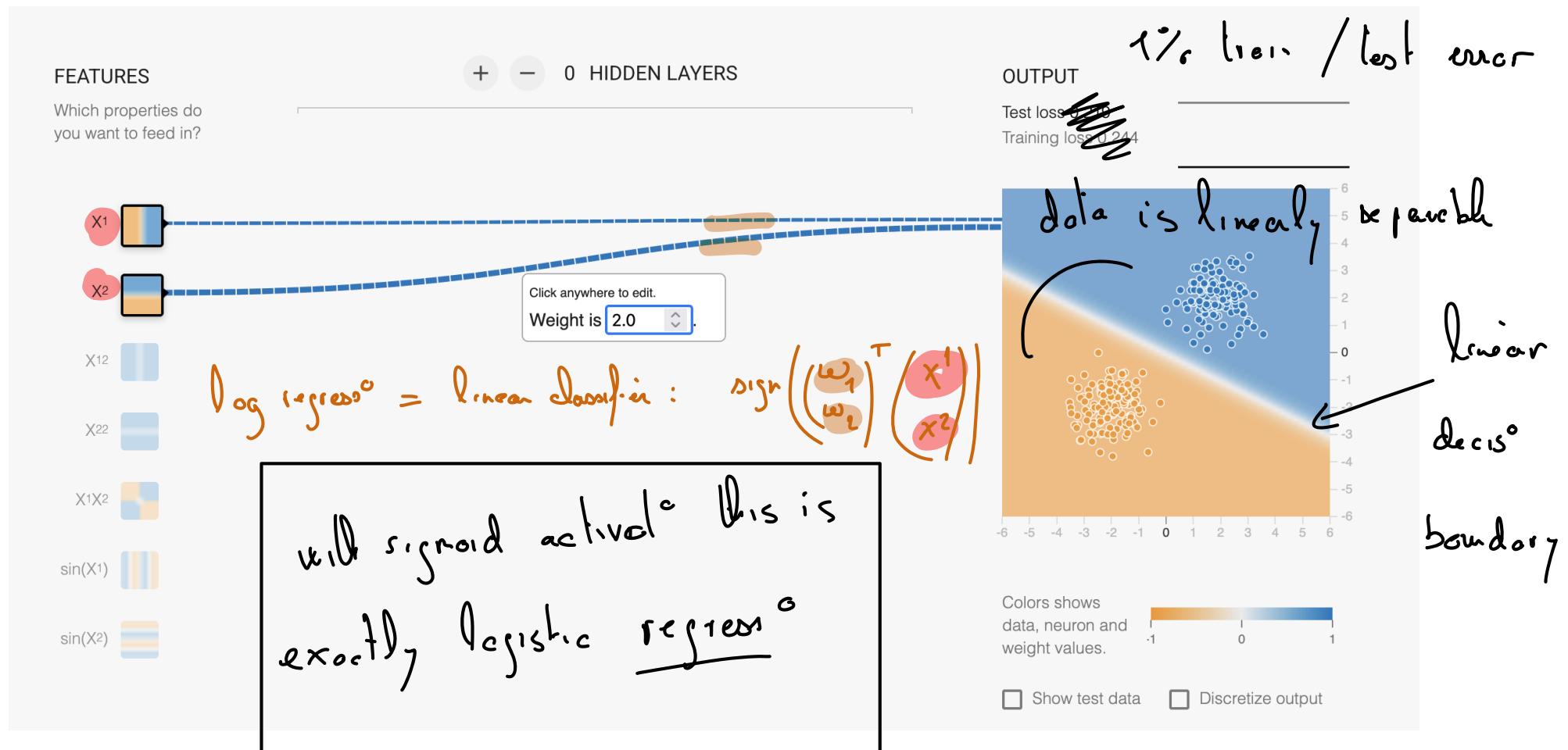
2. Consider the simplest dataset, an 0 hidden layer. Manually pick weights to obtain a good classification.



What model does this correspond to? How many weights does it have?

Let's NN playground: [Link](#)

2. Consider the simplest dataset, an 0 hidden layer. Manually pick weights to obtain a good classification.

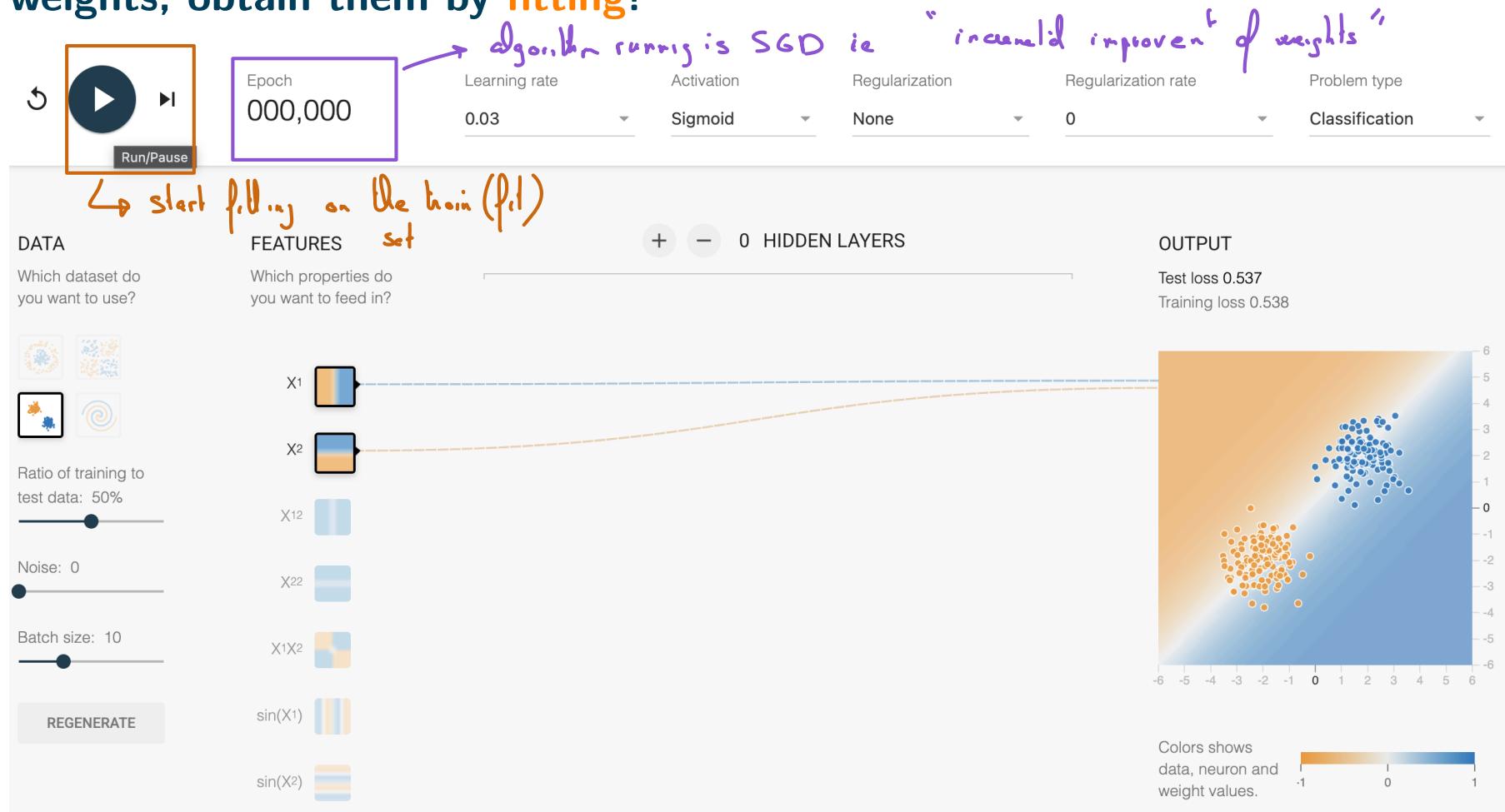


What model does this correspond to? How many weights does it have?

↪(2)

Let's NN playground: [Link](#)

3. Consider the simplest dataset, an 0 hidden layer. Instead of manually tuning the weights, obtain them by **fitting!**



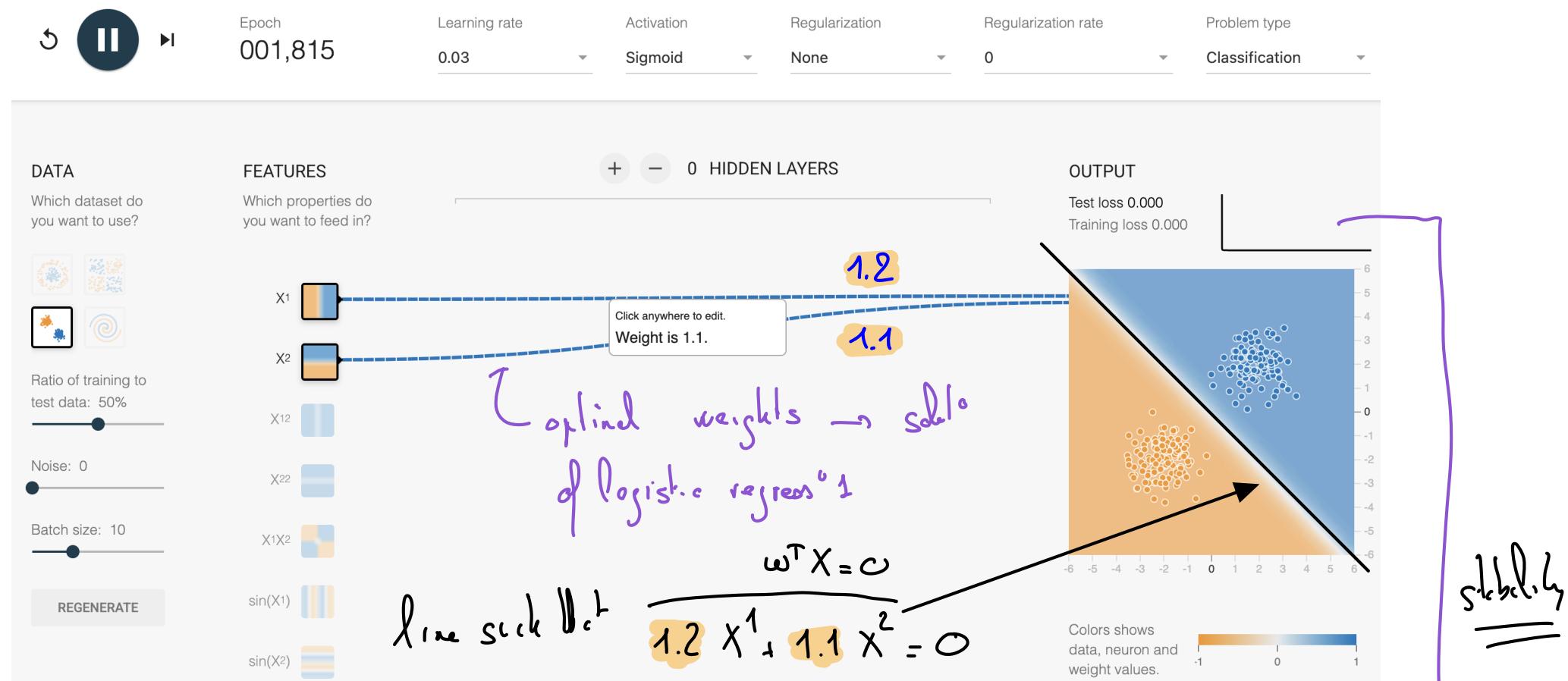
What algorithm are we running? → SGD

What can you say about the loss? →

Rerun the algorithm? What happens?

Let's NN playground: [Link](#)

3. Consider the simplest dataset, an 0 hidden layer. Instead of manually tuning the weights, obtain them by **fitting!**



What algorithm are we running?

What can you say about the loss?

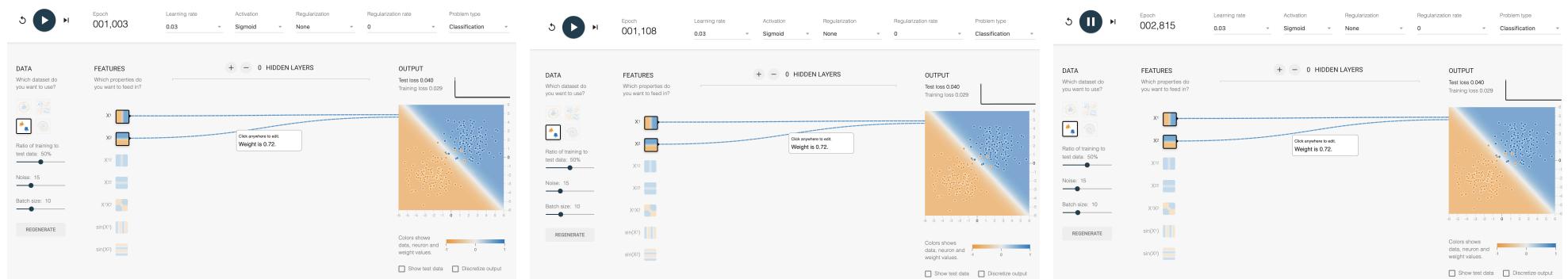
Rerun the algorithm? What happens?

Loss, train and test!

→ almost the same boundary - almost the same weights

Let's NN playground: [Link](#)

4. Increase the data noise to 15. Is the data linearly separable then? Run 3 times for 1000 epochs, compare the weights obtained



What do you conclude / confirm on logistic regression?

Let's NN playground: [Link](#)

5. Change the dataset to . Run logistic regression. What do you observe?



What do you suggest?

Let's NN playground: [Link](#)

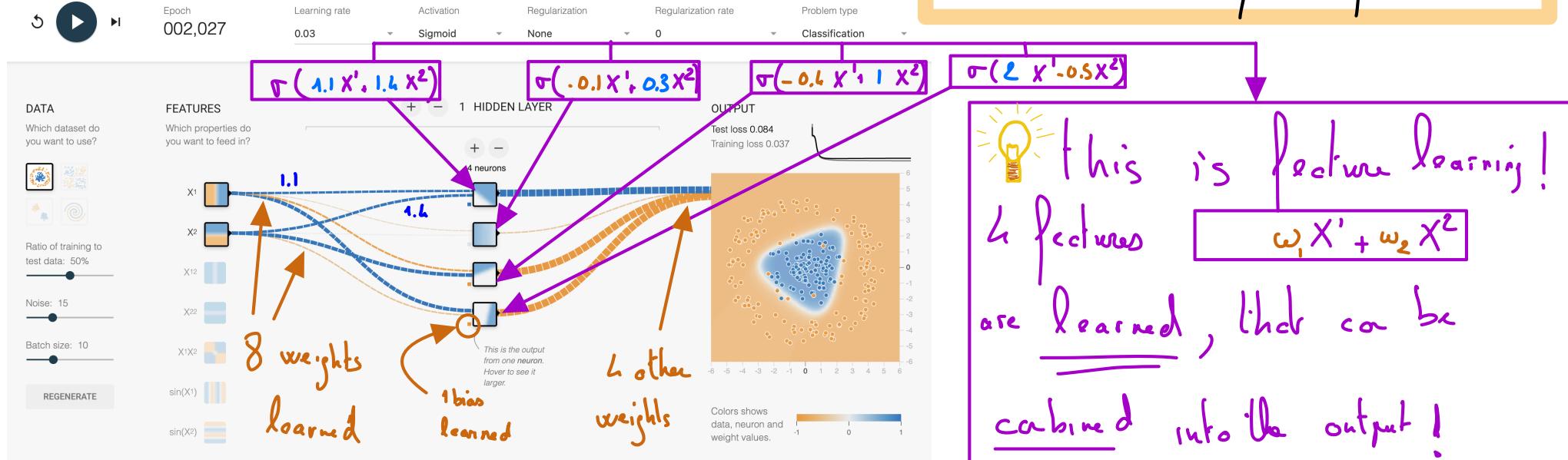
6. Add a single hidden layer, first with a single neuron, then with 4 neurons.
Keep only features X^1 and X^2 .

↳ Single neuron : not enough

Let's NN playground: Link

6. Add a single hidden layer, first with a single neuron, then with 4 neurons. Keep only features X^1 and X^2 .

Epoch 002,027 Learning rate 0.03 Activation Sigmoid Regularization None Regularization rate 0 Problem type Classification



This is your first NN
 this is feature learning!
 $w_1 X^1 + w_2 X^2$

DATA
Which dataset do you want to use?

 Ratio of training to test data: 50%
 Noise: 15
 Batch size: 10
 REGENERATE

FEATURES
Which properties do you want to feed in?
 X1, X2, X12, X22, X1X2, sin(X1), sin(X2)

1 HIDDEN LAYER
 $\sigma(1.1X^1 + 1.4X^2)$, $\sigma(-0.1X^1 + 0.3X^2)$, $\sigma(-0.6X^1 + 1X^2)$, $\sigma(2X^1 - 0.5X^2)$

OUTPUT
 Test loss 0.084
 Training loss 0.037

8 weights learned
 1 bias learned
 6 other weights

This is the output from one neuron. Hover to see it larger.

Colors shows data, neuron and weight values.

Show test data Discretize output

What do you observe? → 11 values! → Loss gets to 5.8% on the test
 What have we learned? → 16 weights! 16 : 12 weights + 4 biases
 How does that illustrate feature learning?
 What happens if you restart learning? → completed, unstable / completed, ≠ after restart.

Let's NN playground: [Link](#)

6. Switch to the more complex spiral dataset  Increase the number of neurons. Fit the model.

Keep only features X^1 and X^2 .

Let's NN playground: [Link](#)

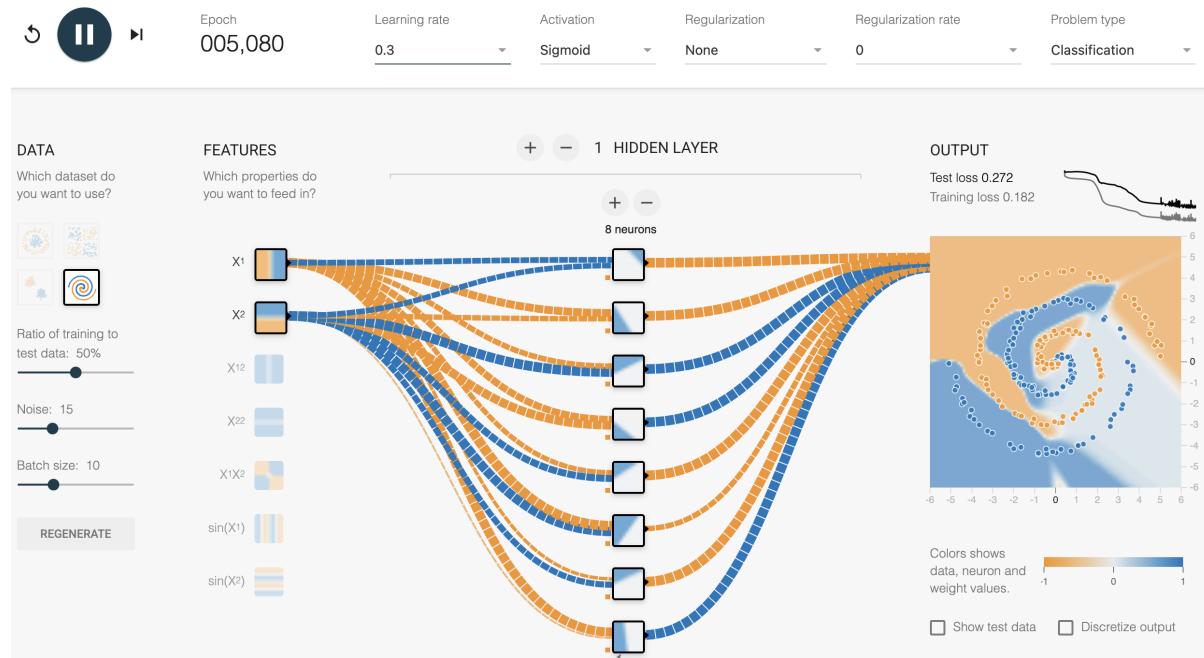
6. Switch to the more complex spiral dataset  Increase the number of neurons. Fit the model.

Keep only features X^1 and X^2 .



Let's NN playground: [Link](#)

6. Switch to the more complex spiral dataset . Increase the number of neurons. Fit the model.
Keep only features X^1 and X^2 .



OUTPUT
Test loss 0.261
Training loss 0.173

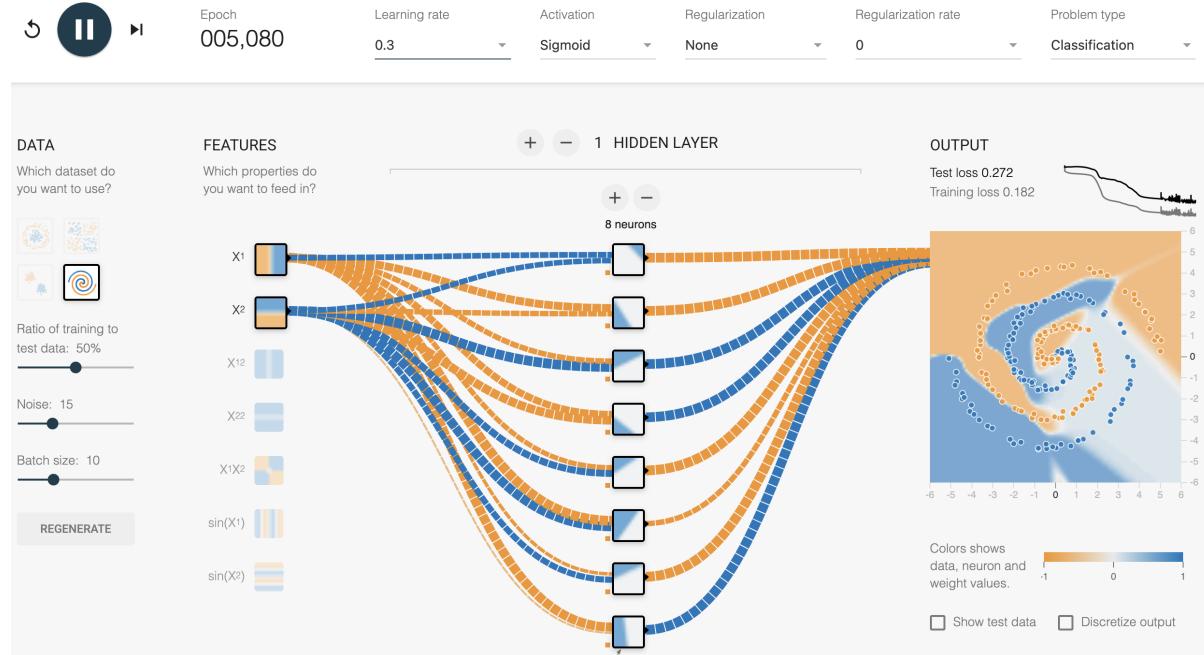


What can you say about the algorithm dynamic ?

Let's NN playground: [Link](#)

6. Switch to the more complex spiral dataset  Increase the number of neurons. Fit the model.

Keep only features X^1 and X^2 .



What do you suggest next?



increase the number of layers

Let's NN playground: [Link](#)

7. Explore configurations : what do you obtain?

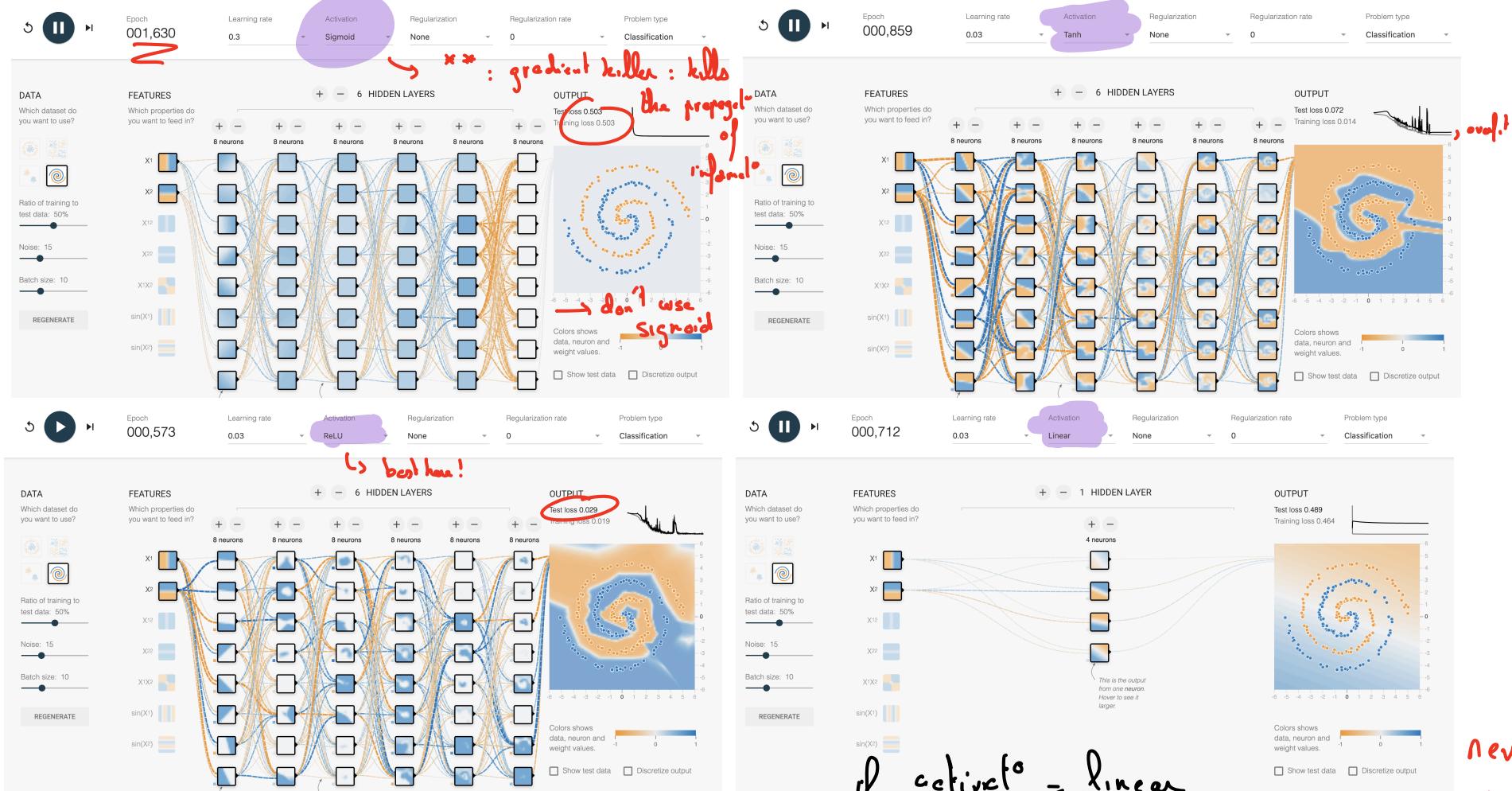
Let's NN playground: [Link](#)

7. Explore configurations : what do you obtain? What is the role of the activation function? What happens for a deep architecture for Linear / Sigmoid / Relu / Tanh?

Let's NN playground: Link

hyperparam

7. Explore configurations : what do you obtain? What is the role of the activation function? What happens for a deep architecture for Linear / Sigmoid / Relu / Tanh?



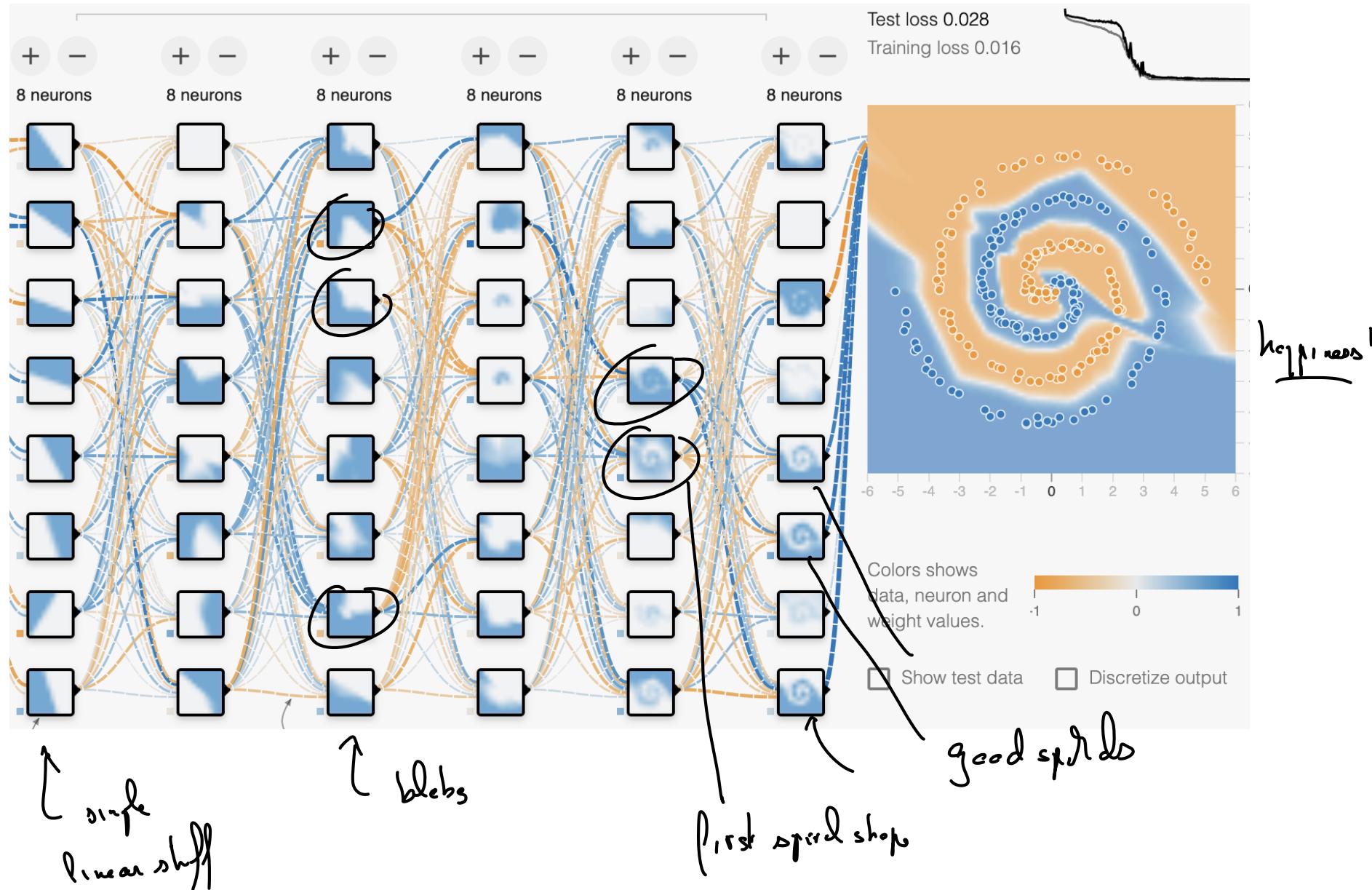
Activation = Linearity breaker

Let's NN playground: [Link](#)

8. How do the features learned evolve depending on the layer?

Let's NN playground: [Link](#)

8. How do the features learned evolve depending on the layer?



9. What are the hyper-parameters of the SGD algorithm? Play with them.

step size : ~~million dollar quest~~ - very hard to pick
many strategies / lots of engineering.

Advanced

- ① What happens if you initialize all weights to 0?

Conclusion on NNPG

- ① “Architecture”:
- * neuron
 - * hidden layer
 - * activation

learned

- * weights
- * biases

...

- * SGD looks like

- ② What we have seen :
- feature learning ❤
 - Instability - randomness in the process
 - Flat moments then progress
 - Impact of depth.

- ③ Protip :
- do it again in one year !!
- add it in your calendar

Wooclap time!



1 Allez sur wooclap.com

2 Entrez le code d'événement dans le bandeau supérieur

Code d'événement
BLCOFS



1 Envoyez @BLCOFS au 06 44 60 96 62

Désactiver les réponses par SMS

2 Vous pouvez participer

Outline

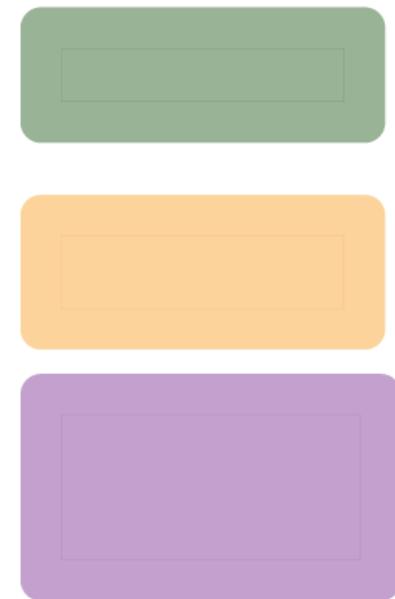
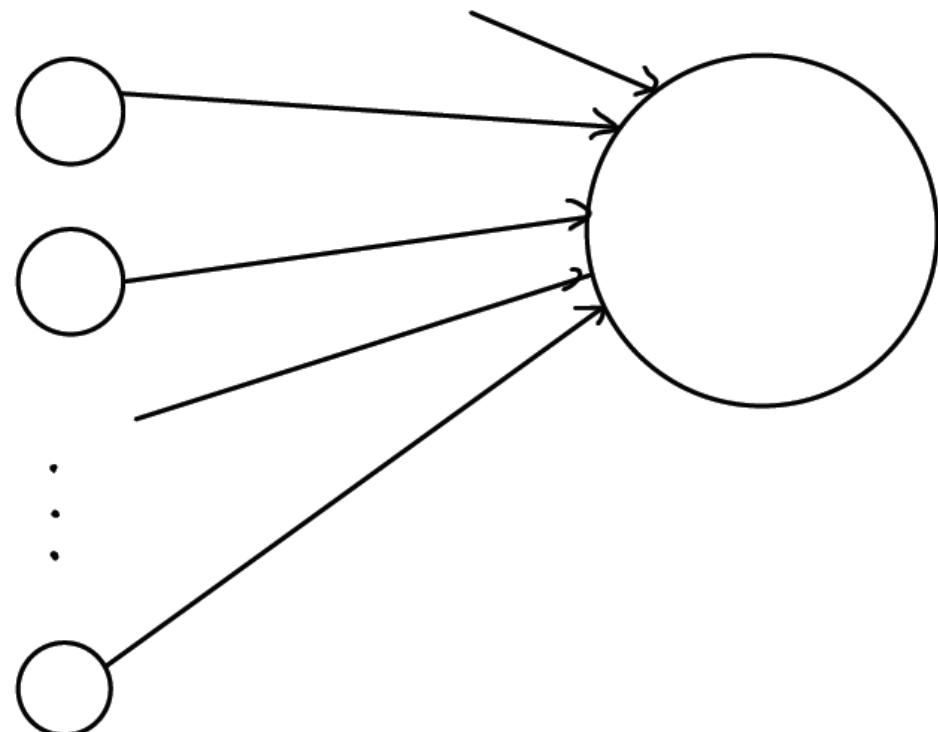
1 Goals

2 The first Neural Network: the Perceptron

- Logistic regression
- Discovering Fully connected NN through the NNPG.
- Formalization, advanced topics
- Further discussion

From logistic regression to Multi Layer Perceptron 3

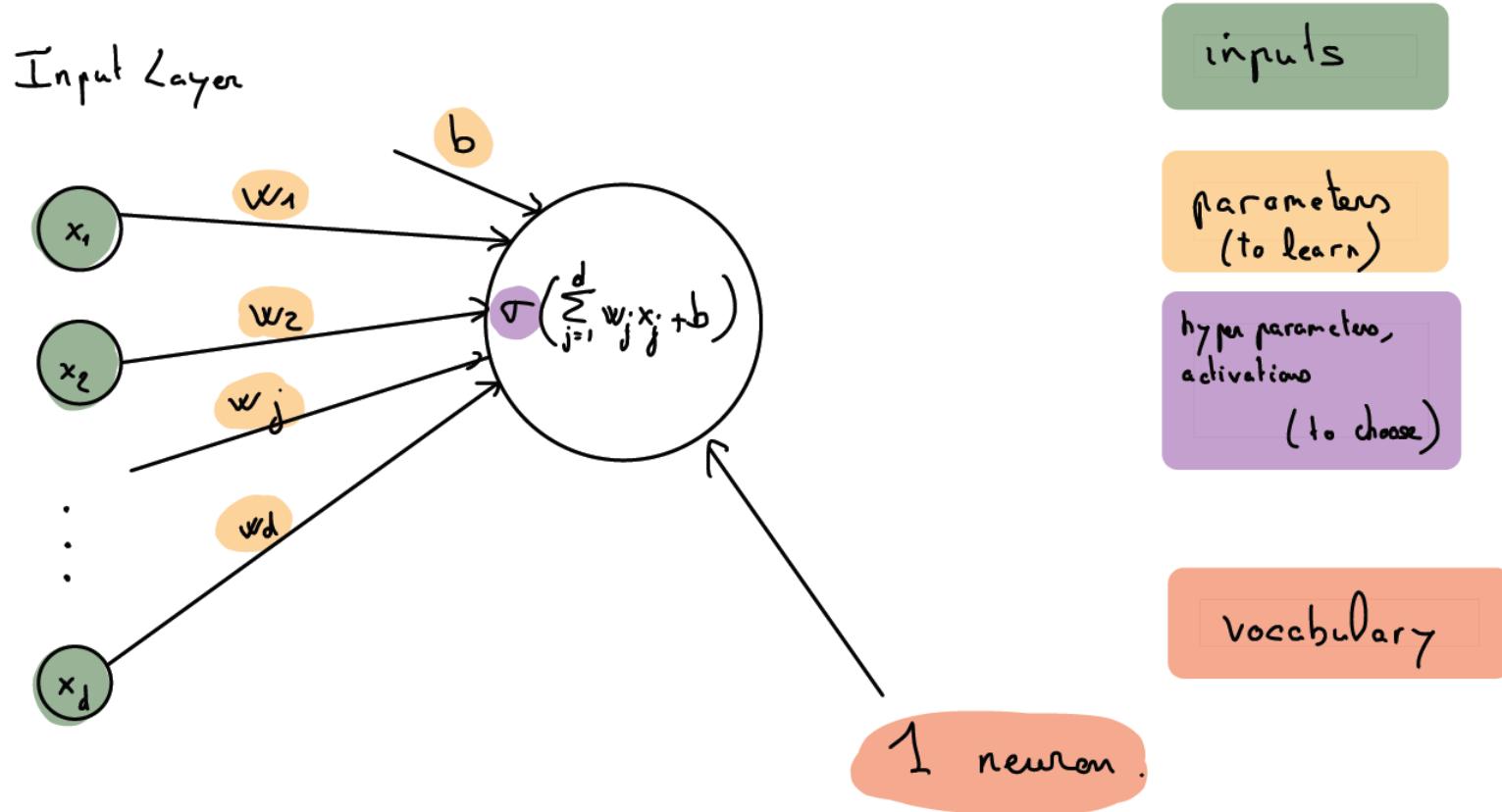
- Class of functions for logistic regression: $\mathcal{F} = \{X \mapsto w^T X, w \in \mathbb{R}^d\}$.
- Prediction for a new point is $1_{\sigma(X_i) > 1/2}$
- Can be seen as predicting a probability $\sigma(X_i)$ that the output is 1



This is a neural network, with one neuron !

From logistic regression to Multi Layer Perceptron 3

- Class of functions for logistic regression: $\mathcal{F} = \{X \mapsto w^T X, w \in \mathbb{R}^d\}$.
- Prediction for a new point is $1_{\sigma(X_i) > 1/2}$
- Can be seen as predicting a probability $\sigma(X_i)$ that the output is 1



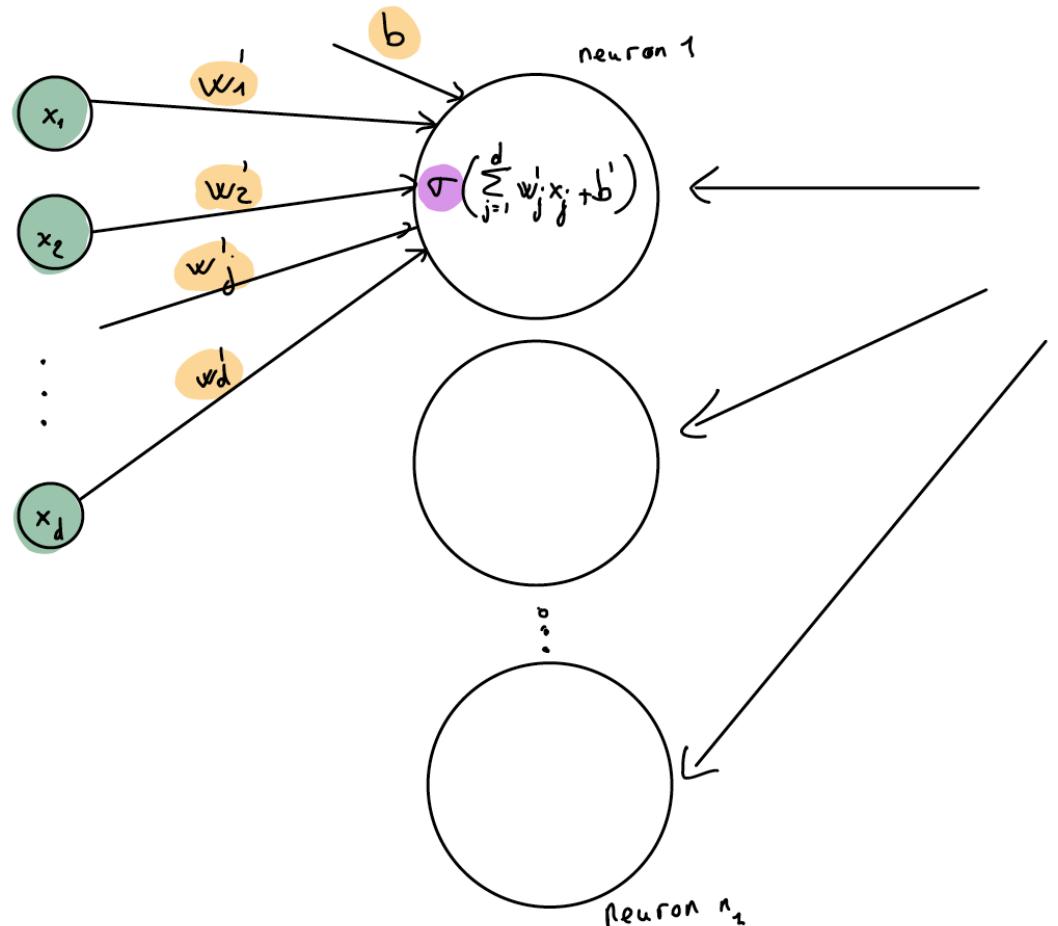
This is a neural network, with one neuron !

= This is logistic regress°

From logistic regression to Multi Layer Perceptron 4

Extend to

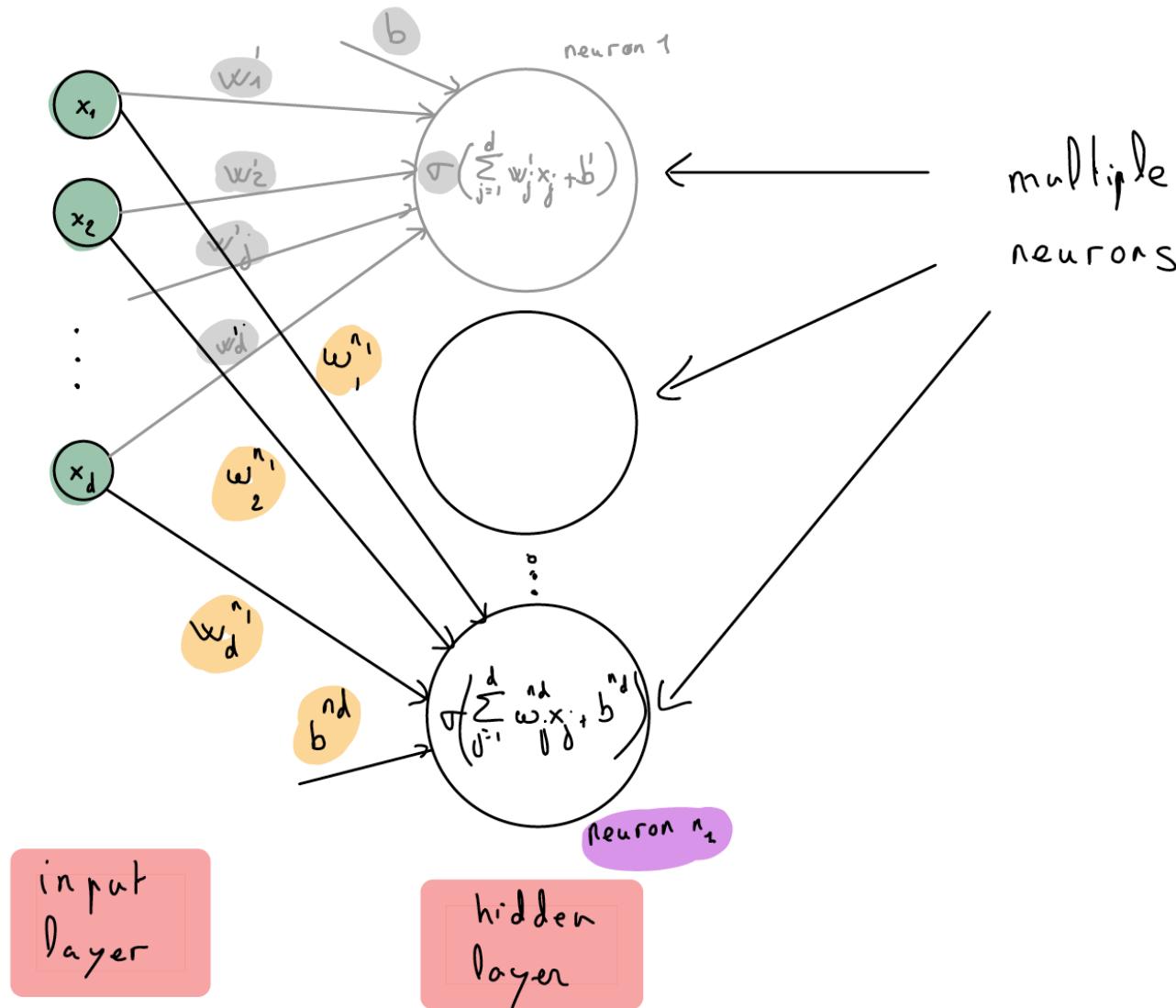
- Multiple Neurons
- Multiple Layers



From logistic regression to Multi Layer Perceptron 4

Extend to

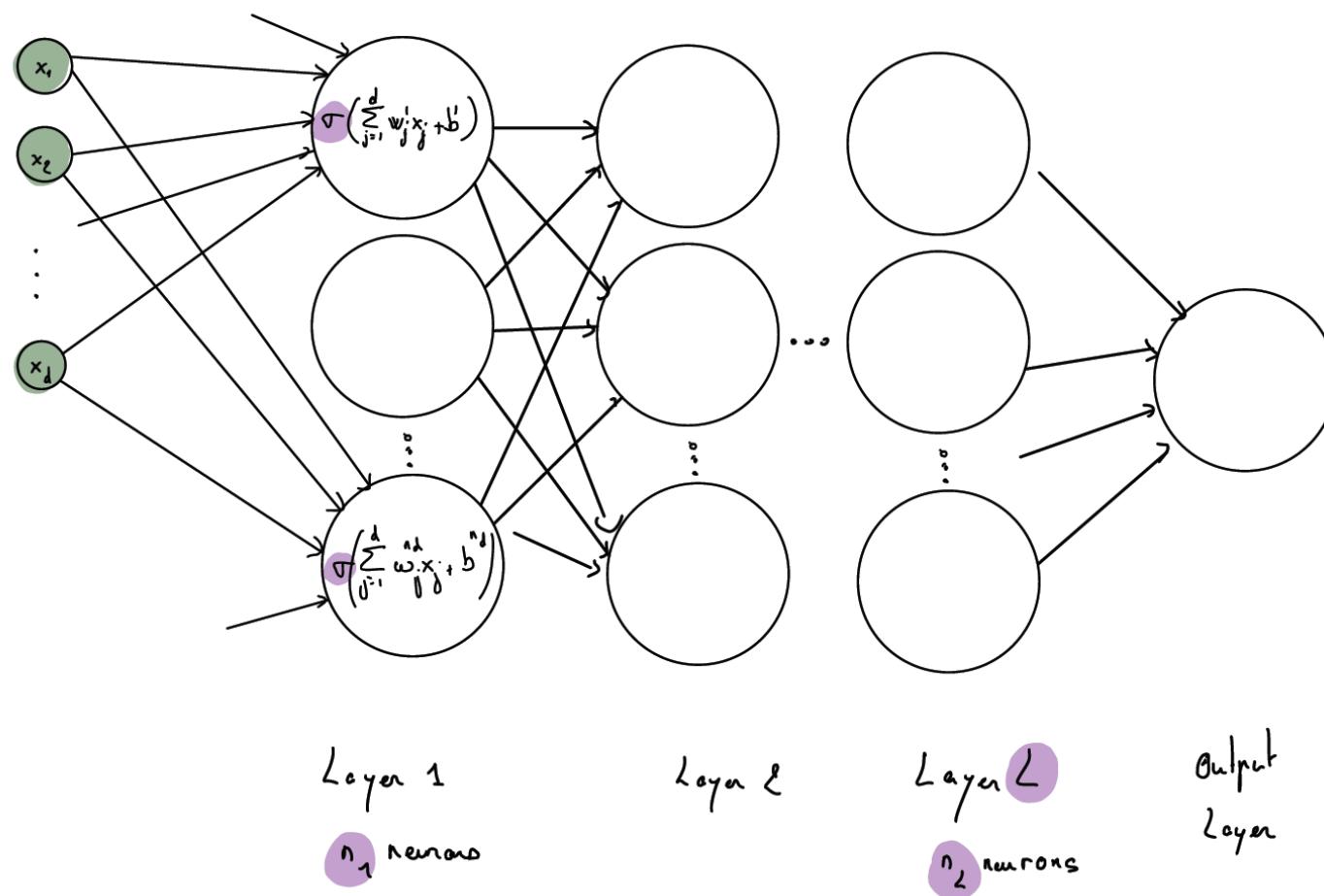
- Multiple Neurons
- Multiple Layers



From logistic regression to Multi Layer Perceptron 5

Extend to

- Multiple Neurons
- Multiple Layers



From logistic regression to Multi Layer Perceptron 5

Neural network:

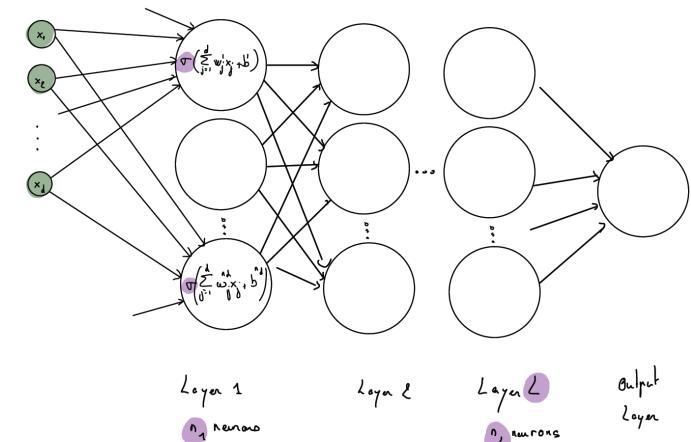
- Goal $\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)]$

- Class of non linear functions.

$$\mathcal{F}_{\text{MLP}} = \{f(X; W, b) = \sigma(b_n + W_n \sigma(\dots(b_1 + W_1 x))),$$

$$W \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n\}$$

math. portfolio of the class of func^s we've looked at!



Summary: NN generalize regression by looking for candidates in a much larger class of functions than linear ones.

What needs to be learned

W

What needs to be chosen

archt. Layers

Vocabulary

Neuron, SGD,

b

Neuron

act. rel.

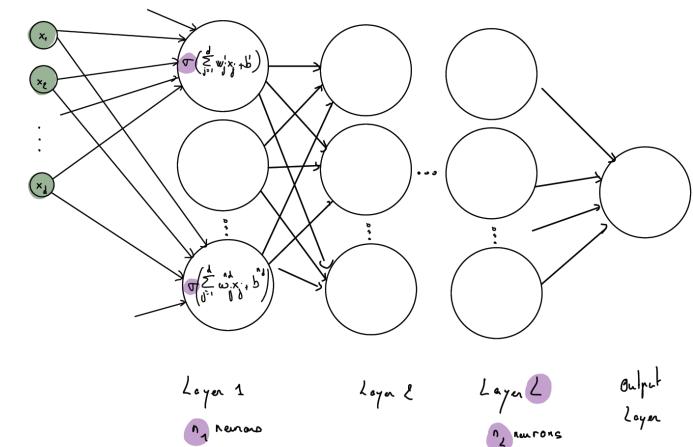
etc

From logistic regression to Multi Layer Perceptron 5

Neural network:

- Goal $\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)]$
- Class of non linear functions.
$$\mathcal{F}_{\text{MLP}} = \{f(X; W, b) = \sigma(b_n + W_n \sigma(\dots(b_1 + W_1 x))),$$

$$W \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n\}$$



Summary: NN generalize regression by looking for candidates in a much larger class of functions than linear ones.

What needs to be learned

Parameters W, b

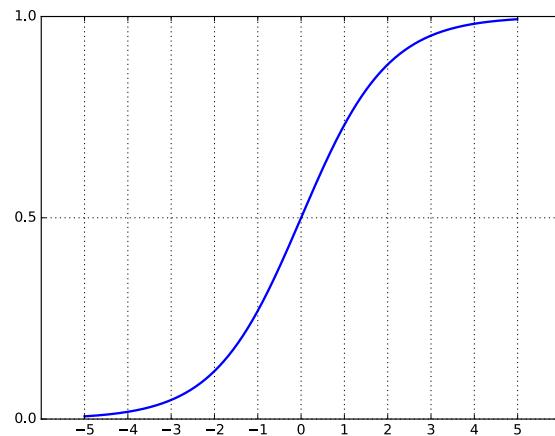
What needs to be chosen

- Activation σ
- Number of layers L
- Number of neurons per hidden layer
 $n_i, i = 1, \dots, L$.

Vocabulary

- Neuron
- Activation
- Hidden Layer
- Width, Depth
- Fully connected layer

2 possible activations: Sigmoid and Rectified Linear Unit



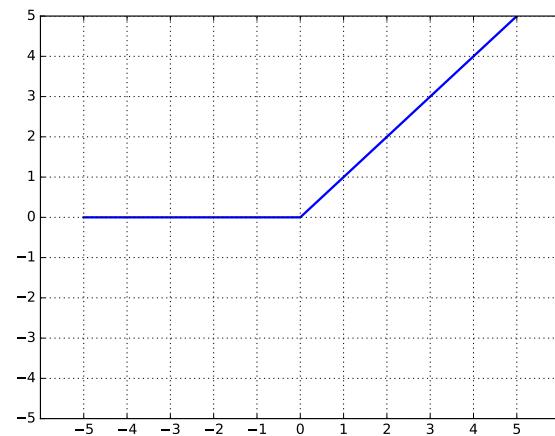
Sigmoid function

- $x \mapsto \frac{\exp(x)}{1+\exp(x)}$

- Problems:

- ① Saturated function: gradient killer \rightarrow need for rescaling data

not to be used



Rectified Linear Unit (ReLU)

- $x \mapsto \max(0, x)$

- Pros & cons:

- ① Not a saturated function. Kills negative values.
- ② Empirically, convergence is faster than sigmoid/tanh.
- ③ Plus: biologically plausible

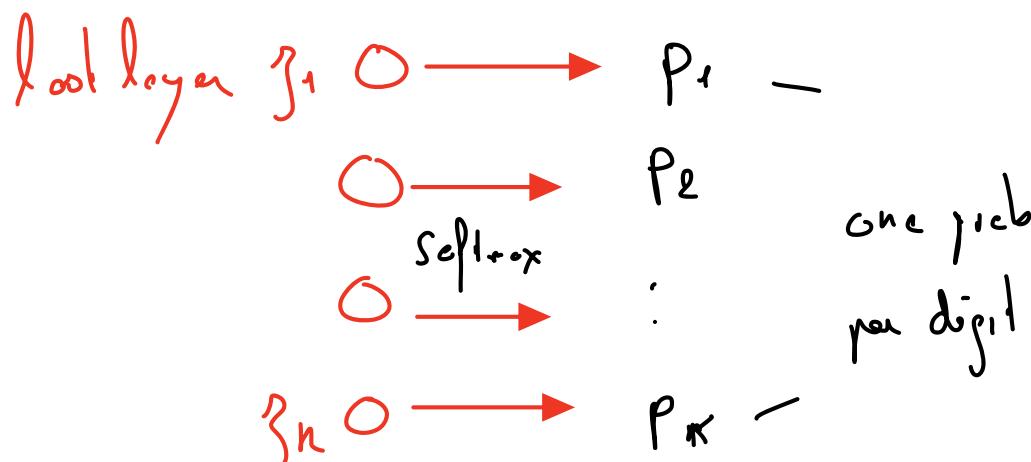
How to deal with multi-class output: softmax activation

When we do multi-class classification, i.e., $Y_i \in 1, \dots, K$, we output K different values:

- Each of them corresponds to the probability of belonging to the class j , $1 \leq j \leq K$
- To obtain probabilities (adding up to 1):
 - ① We have K neurons on the last layer
 - ② And use a **Softmax** activation to renormalize.

Softmax output unit, used to predict $\{1, \dots, K\}$:

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}$$



Summary: if you
have K classes ($K = K$
for d.g.s)

change the last layer activation
the loss

Up to that point, we have seen :

- What a Neural Network was
- Why it is expected to learn better.

Next Question: how to implement it ?

Python - Keras

1. Hardware:

- CPU
- GPU
- TPU

2. Software (Python packages):

- Pytorch/Tensorflow
- → Keras : TF high level API. Ideal for applications.



Keras - A few lines !¹

What do we need to specify?

- * architecture L , $n_1 \dots n_L$
Layers neurons or layers $1 \dots L$
- * activation
- * regularization

What do we need to do next?

* properties
of learning
algo

```
# import tensorflow and keras (tf.keras not "vanilla" Keras)
import tensorflow as tf
from tensorflow import keras
# get data
(train_images, train_labels), (test_images,
                               test_labels) = keras.datasets.mnist.load_data()

# setup model
model = keras.Sequential([
    keras.layers.Flatten(input_shape=(28, 28)), not h.l?
    keras.layers.Dense(128, activation=tf.nn.relu), h.l.
    keras.layers.Dense(64, activation= tf.nn.relu),
    keras.layers.Dense(10, activation=tf.nn.softmax)]) output layer

model.compile(optimizer=tf.optimizers.Adam(),
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

# train models
model.fit(train_images, train_labels, epochs=5)

print('\n Evaluation')
# evaluate
test_loss, test_acc = model.evaluate(test_images, test_labels)
print('test accuracy:', test_acc)
# make predictions
predictions = model.predict(test_images)
```

Wooclap!



- how many layers?
- what is the final test accuracy ?
- what happens if you remove the hidden layers?
- how many parameters and why?
- how many points in the train set?
- what does 1865 and 313 correspond to?

```
# import tensorflow and keras (tf.keras not "vanilla" Keras)
import tensorflow as tf
from tensorflow import keras
# get data
(train_images, train_labels), (test_images,
                               test_labels) = keras.datasets.mnist.load_data()
# setup model
model = keras.Sequential([
    keras.layers.Flatten(input_shape=(28, 28)),
    keras.layers.Dense(128, activation=tf.nn.relu),
    keras.layers.Dense(64, activation= tf.nn.relu),
    keras.layers.Dense(10, activation=tf.nn.softmax)
])

model.compile(optimizer=tf.optimizers.Adam(),
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])
# train models
model.fit(train_images, train_labels, epochs=5)

print('\n Evaluation')
# evaluate
test_loss, test_acc = model.evaluate(test_images, test_labels)
print('test accuracy:', test_acc)
# make predictions
predictions = model.predict(test_images)
```



What can you say about this code bit

- ① The flatten layer is a dense layer
- ⑤ I use SGD for the algorithm
- ② I have 2 hidden layers
- ⑥ AdamOptimizer is EveOptimizer's boyfriend
- ③ I have one hidden layer
- ⑦ I run 5 epochs of the stochastic algorithm
- ④ I should use sigmoid activation on the last layer
- ⑧ I use relu activation on all but the last layer

Wooclap

```
# import tensorflow and keras (tf.keras not "vanilla" Keras)
import tensorflow as tf
from tensorflow import keras
# get data
(train_images, train_labels), (test_images,
                               test_labels) = keras.datasets.mnist.load_data()

# setup model
model = keras.Sequential([
    keras.layers.Flatten(input_shape=(28, 28)),
    keras.layers.Dense(128, activation=tf.nn.relu),
    keras.layers.Dense(64, activation=tf.nn.relu),
    keras.layers.Dense(10, activation=tf.nn.softmax)
])

model.compile(optimizer=tf.optimizers.Adam(),
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

# train models
model.fit(train_images, train_labels, epochs=5)

print('\n Evaluation')
# evaluate
test_loss, test_acc = model.evaluate(test_images, test_labels)
print('test accuracy:', test_acc)
# make predictions
predictions = model.predict(test_images)

Epoch 1/5
1875/1875 6s 3ms/step - accuracy: 0.8124 - loss: 4.6972
Epoch 2/5
1875/1875 4s 2ms/step - accuracy: 0.9239 - loss: 0.3186 ONCE
Epoch 3/5
1875/1875 5s 2ms/step - accuracy: 0.9464 - loss: 0.2017
Epoch 4/5
1875/1875 5s 2ms/step - accuracy: 0.9548 - loss: 0.1581
Epoch 5/5
1875/1875 4s 2ms/step - accuracy: 0.9610 - loss: 0.1388

Evaluation → evaluate last error.
313/313 1s 3ms/step - accuracy: 0.9592 - loss: 0.1602
test accuracy: 0.9621999859809875
313/313 1s 2ms/step
```

→ evaluate last error.

your choice
on how long
algo runs for
one epoch
look at all fit data

steps of algo / epoch

at each step I look at
a fraction (a batch) of data
by default, 32 (in tf. fit)

5 epochs | 1875 steps | 60k
x 32 batchsize images in fit dataset

We now see the output. What can we say



- ① We overfit
- ② There are 313 points in the test set
- ③ There are 10 000 points in the test set
- ④ 1875 is about 6 times 313
- ⑤ the loss is the accuracy
- ⑥ the loss decreases along the trajectory

$$32 \times 1875 = 60.000$$

$$313 \times 32 =$$

$$\boxed{10016}$$
 total dataset
$$10000$$

Wooclap

model.summary()

Model: "sequential"

Layer (type)	Output Shape	Param #
flatten (Flatten)	(None, 784)	0
dense (Dense)	(None, 128)	100480
dense_1 (Dense)	(None, 64)	8256
dense_2 (Dense)	(None, 10)	650

Total params: 109,386
Trainable params: 109,386
Non-trainable params: 0

Optimizer parameters $218 \cdot 774$ $\downarrow_{1,2}$ $\curvearrowright 2x$
Adam total per $3 \times 109,386$ $\downarrow 12$

SGD : 2 !

Explain the number of parameters on each layer

model.summary()

Model: "sequential"

$$784 \times 128 = 100,280$$

Layer (type)	Output Shape	Param #
flatten (Flatten)	(None, 784)	0
dense (Dense)	(None, 128)	$784 \times 128 + 128$ biases
dense_1 (Dense)	(None, 64)	$128 \times 64 = 8192$
dense_2 (Dense)	(None, 10)	$64 \times 10 = 640$

Total params: 109,386
Trainable params: 109,386
Non-trainable params: 0

$$128 \times 64 + 1 \times 64$$

one "vector of size $109k$ " for momentum
one "vector of size $109k$ " for step size

$$128 \times 64$$

Outline

1 Goals

2 The first Neural Network: the Perceptron

- Logistic regression
- Discovering Fully connected NN through the NNPG.
- Formalization, advanced topics
- Further discussion

Intuition

Up to that point, we have seen :

- What a Neural Network was.
- How to implement it in Python.

Next Question: why and when does it work?

Intuition

Up to that point, we have seen :

- What a Neural Network was.
- How to implement it in Python.

Next Question: why and when does it work?

- ➊ Approximation theorem
- ➋ Optimization

Approximation Theorem

Goal

$$\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)].$$

As said before, a MLP can output non-linear functions... But how powerful is it?

Approximation Theorem

Goal

$$\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)].$$

As said before, a MLP can output non-linear functions... But how powerful is it?

Continuous Neural Networks with one single hidden layer and any bounded and non constant activation function can approximate any function in L^p , provided a sufficient number of hidden units.

- ↪ Very powerful in terms of approximations.
- Not very surprising: non linear function with millions of parameters!

Optimization - Fitting the network

Goal

$$\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)].$$

How to minimize a function?

Gradient Methods. = cheapest way to update (learn) the weights iteratively to minimize the loss.

Optimization - Fitting the network

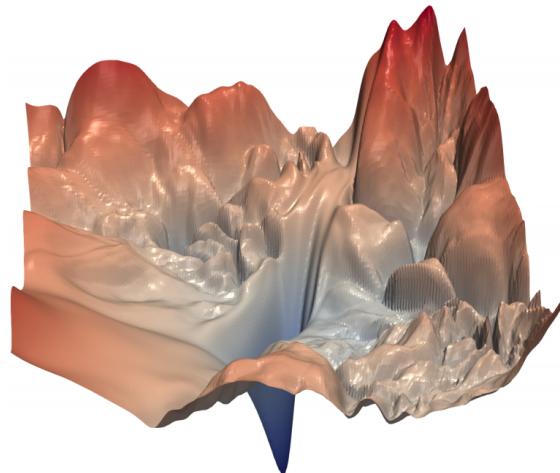
Goal

$$\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)].$$

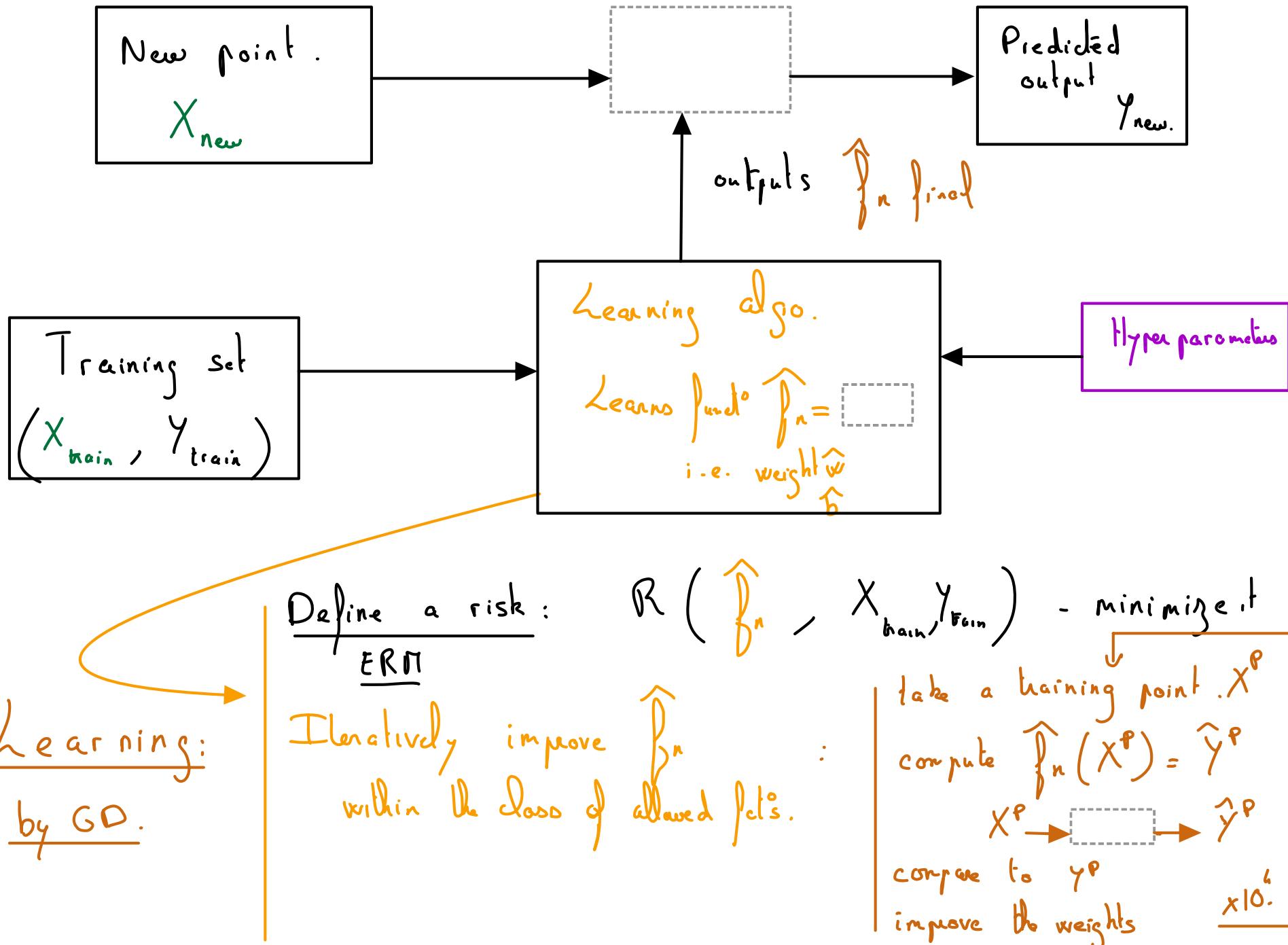
Gradient Methods. = cheapest way to update (learn) the weights iteratively to minimize the loss.

Two “miracles”:

- ① Autodifferentiation: even though the function is very complex and high dimensional, it is possible to compute gradients!
- ② Optimization seems to work ok, though the function is non convex - we do not end up in too bad local minima. (specialized optim algos: Adam, AdaDelta, etc.)



These 2 miracles make it possible to learn a good regressor/classifier with NN and to benefit from the powerful approximation properties



Conclusion

Neural networks :

- ① learn very complex non linear functions in high dimension
- ② can approximate nearly any function
- ③ can be optimized even though highly non convex.

Are thus expected to work:

- When depth or width increases.
- Thus with very large datasets
- Requiring a lot of computational power.

⇒ **Explains why they were so successful since 2010**

Conclusion

Neural networks :

- ① learn very complex non linear functions in high dimension
- ② can approximate nearly any function
- ③ can be optimized even though highly non convex.

Are thus expected to work:

- When depth or width increases.
- Thus with very large datasets
- Requiring a lot of computational power.

⇒ **Explains why they were so successful since 2010**

What my layer wants me to say: we haven't talked about many points !!

- Initialization
- Regularization

Conclusion

Neural networks :

- ① learn very complex non linear functions in high dimension
- ② can approximate nearly any function
- ③ can be optimized even though highly non convex.

Are thus expected to work:

- When depth or width increases.
- Thus with very large datasets
- Requiring a lot of computational power.

⇒ **Explains why they were so successful since 2010**

What my layer wants me to say: we haven't talked about many points !!

- Initialization
- Regularization

Next:

- Brief history
- How to use the structure?

A brief history of NN

- 1943: Neural networks
- 1957: Perceptron
- 1974-86: Backpropagation, RBM, RNN
- 1989-98: CNN, MNIST
- 2009: ImageNet
- 2012: AlexNet,
- 2014: GANss
- 2016: AlphaGo
- 2018: BERT

The Computational power made the major change (+ investment and creativity).

Directions

- When does it work?

Large or huge datasets. Structured tasks.

↳ **Images and text.**

⚠ do not try to apply DL everywhere

Each application requires a special architecture:

- Images : Convolutional Neural Network (CNN).
- Text : Recurrent Neural Networks (RNN) (2012-2018), Transformers (2018-2025).

2027

2025