



# Hands-on Introduction to Deep Learning

---

## Methodology



**IDRIS**

- 1 What **work** can be done to **improve** the **data** used for training?
- 2 How can a model be **evaluated**?
- 3 Is it possible to make the training more **robust**?
- 4 Can we **benefit** from an **already** trained model?
- 5 Bonus: Any good **practices**? Good architectures?



Learning from exercise with a teacher to guide us



Applying what we learn to the real world

## CIFAR-10

AllConv



SHIP  
CAR (99.7%)

NiN



HORSE  
FROG (99.9%)

VGG16



DEER  
AIRPLANE (85.3%)

## ImageNet

BVLC AlexNet

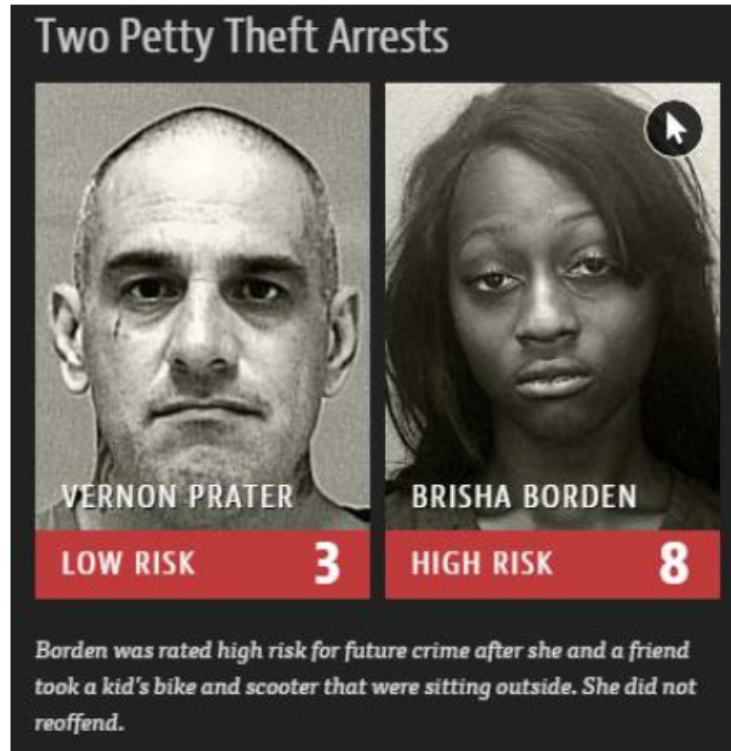


Cup (16.48%)  
Soup Bowl (16.74%)



Bassinet (16.59%)  
Paper Towel (16.21%)

## COMPAS



Source : *Propublica*

## ALGORITHME DE RECRUTEMENT

### Quand le logiciel de recrutement d'Amazon discrimine les femmes

En 2014, le géant du e-commerce a voulu confier ses candidatures à un algorithme, mais celui-ci a commencé à écarter les profils féminins.

Source : *Les Echos*

**Amazon a dû désactiver une IA qui discriminait les candidatures de femmes à l'embauche**

Source : *Numerama*

**Generalization Lack Issue : Discrimination**

# Dataset is the guilty ?!

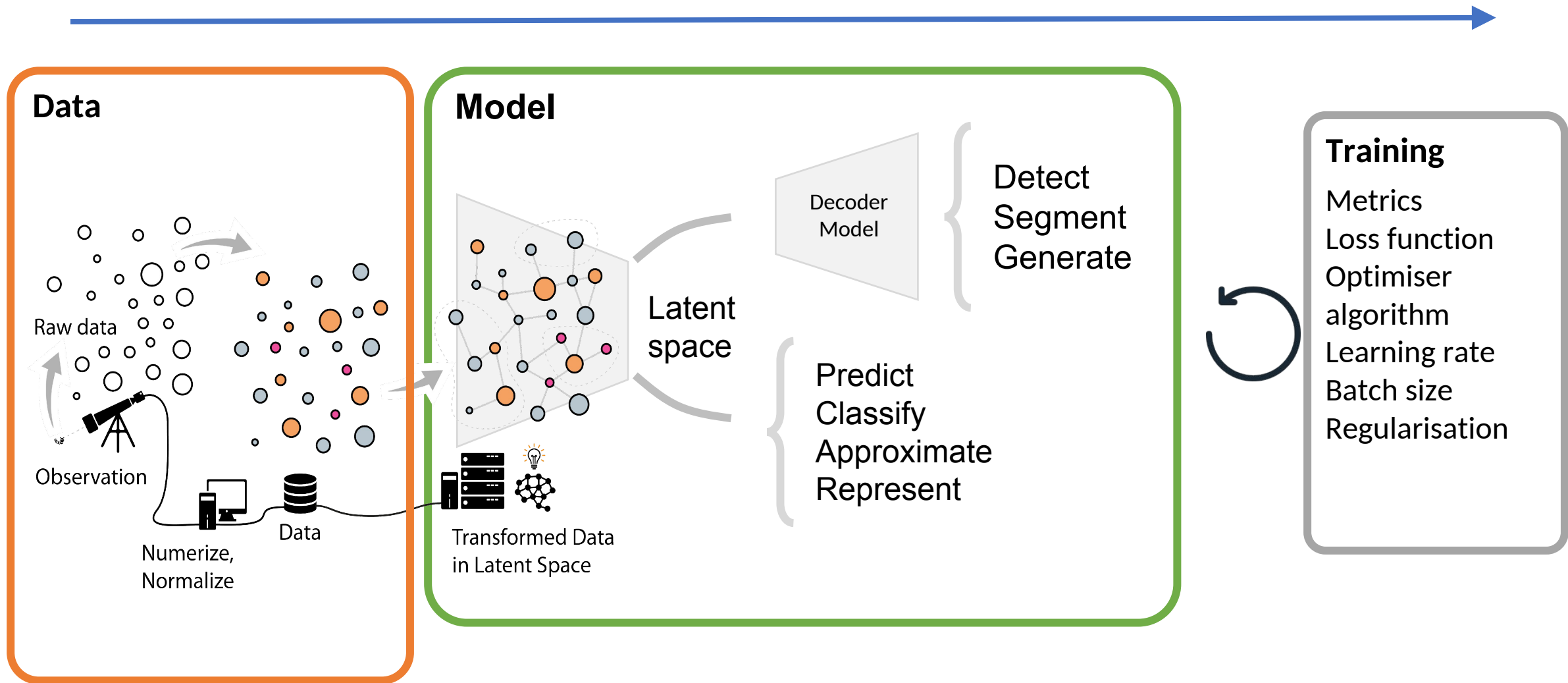


Augment It, Transform It !!!

# Super Regularizator !!



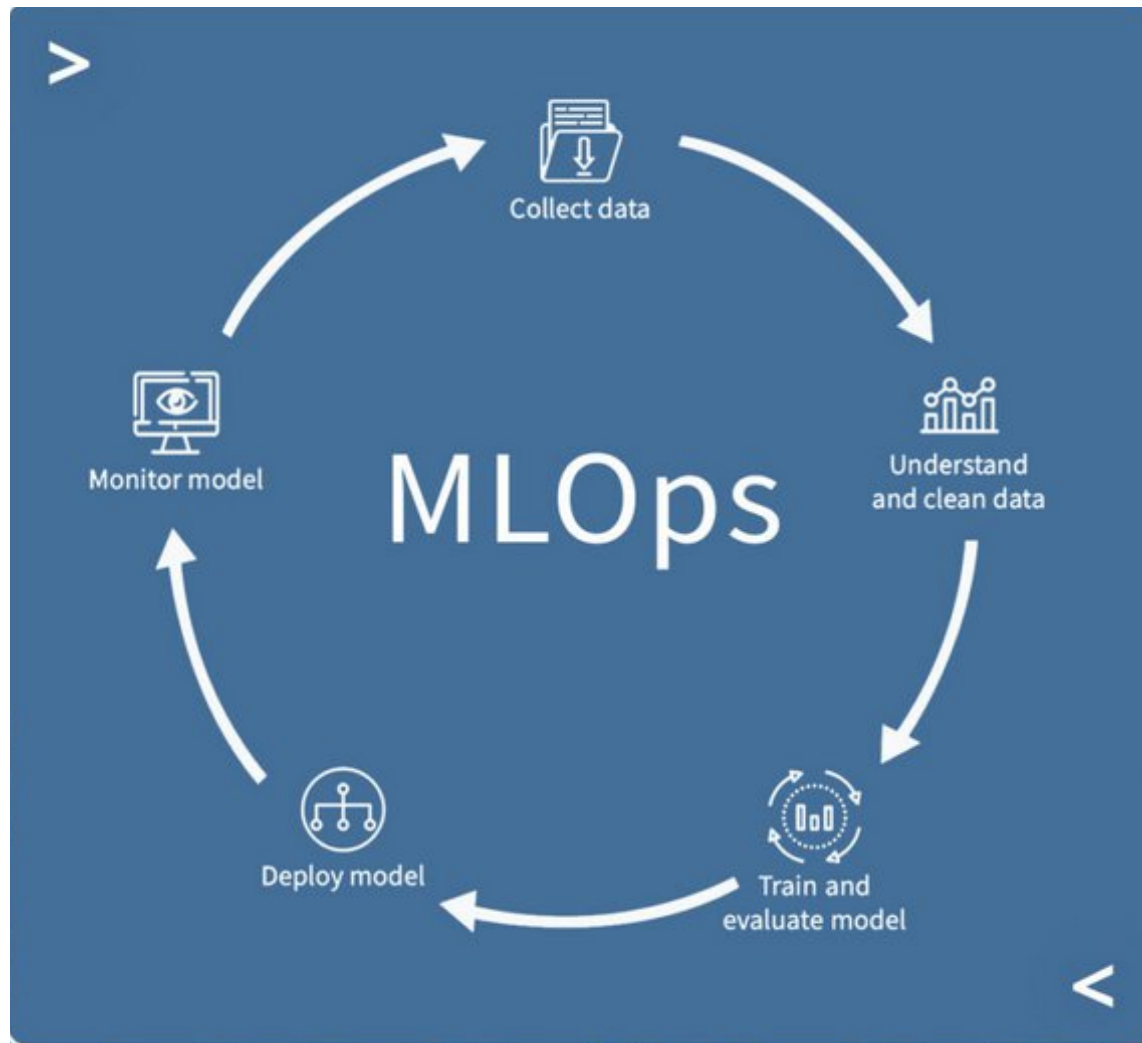
Use his Super Powers !!



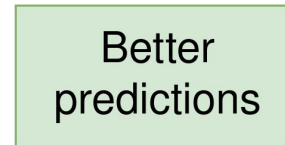
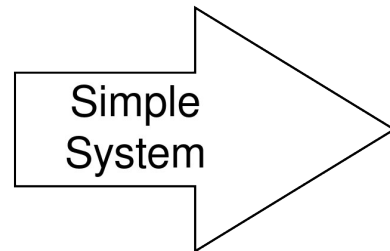
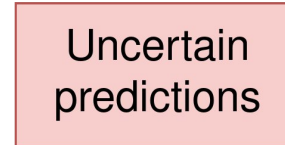
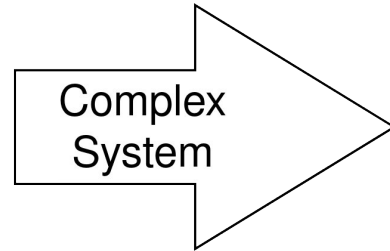
# Deep Learning Pipeline Principle

- 1 Quel **travail faire** pour **améliorer** les **données** utilisées pour **l'entraînement** ?
- 2 Comment **évaluer** un **modèle** ?
- 3 Est-il possible de rendre **l'entraînement** plus **robuste** ?
- 4 Peut-on **profiter** d'un modèle **déjà entraîné** ?
- 5 Bonus : Quelques **bonnes pratiques** ?





50 to 80% time spent on data



## Diabetes risk prediction system

All the features



Selection

Selected features



**All the features**



**Extraction**

**Extracted features**



### Image



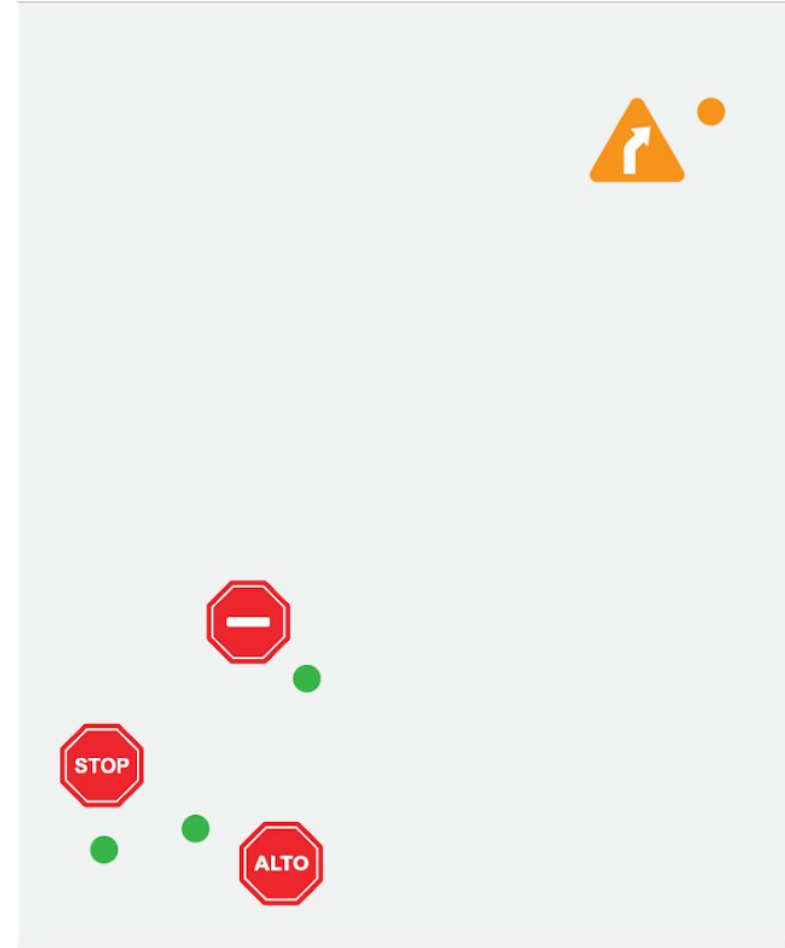
### Image Embedding

Vector representation of the image that captures some of the semantic meaning

|   | sign | stop | english | symbol | octagon |
|---|------|------|---------|--------|---------|
| → | 1    | 1    | 1       | -1     | 1       |
| → | 1    | -1   | 0.5     | 0.6    | -1      |
| → | 1    | -1   | -0.3    | -1     | 1       |
| → | 1    | 0.8  | -0.1    | 0.8    | 1       |

### Dimensionality Reduction Visualization

Lower dimensional representation of the vector.  
Places similar objects closer to each other.



## Features extraction - Embedding example

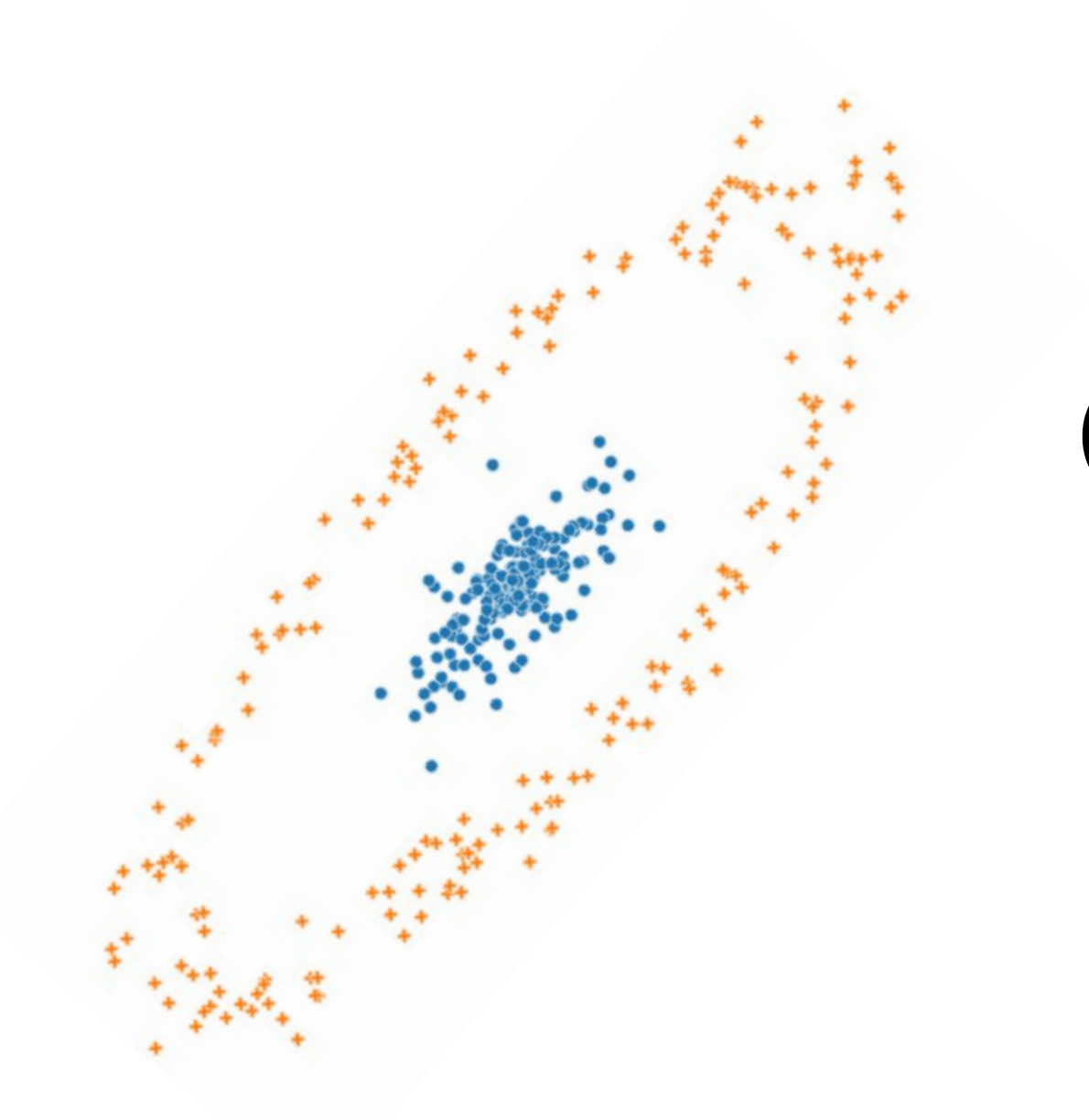
All the features



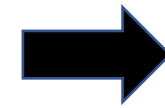
Transformation

Transformed features





$(x, y)$

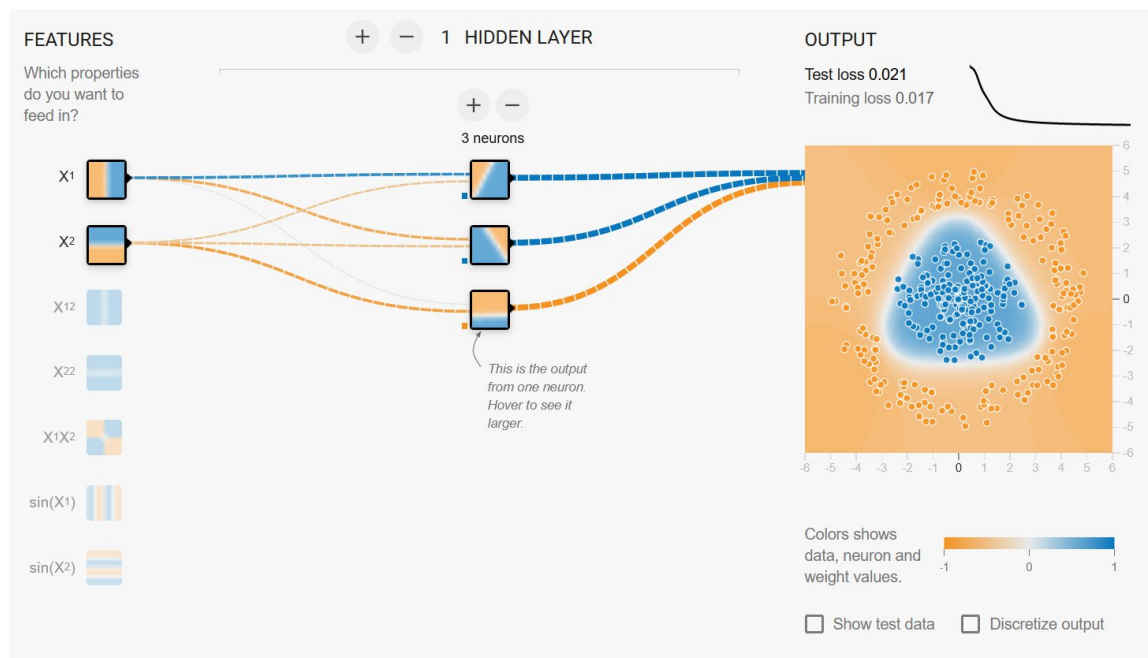


**Complex relation  
between  $x$  and  $y$**

$(r, \theta)$

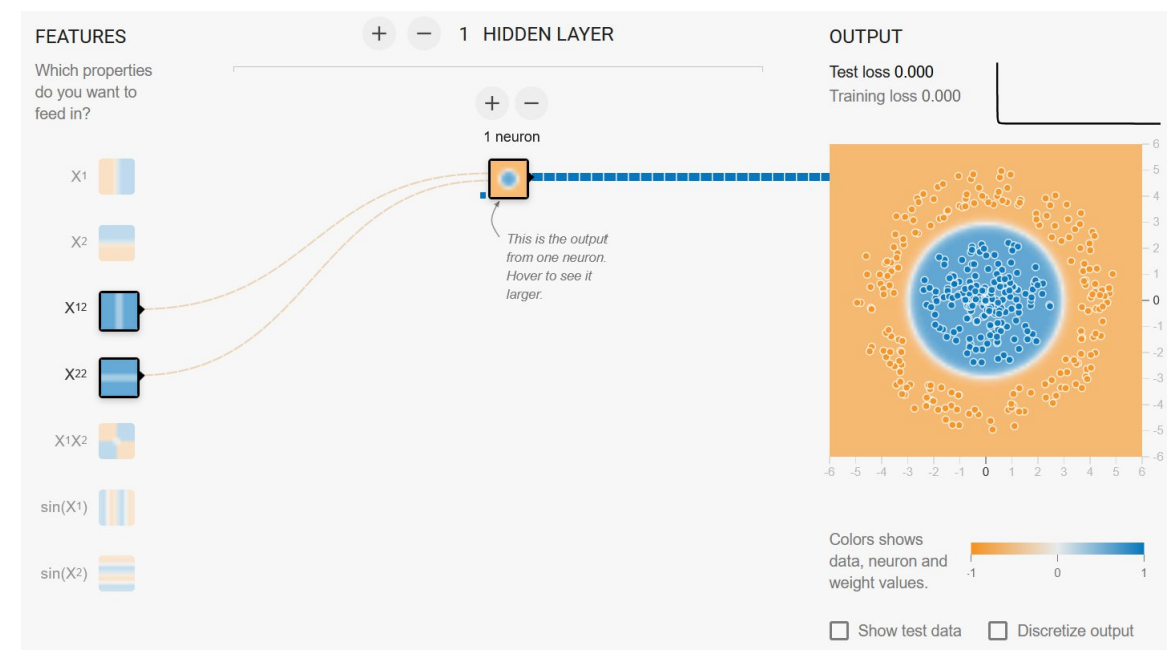


**Simple relation  
with  $r$  and  $\theta$**



Données de position (x, y)  
Réseau à 3 neurones

Données transformées  
Réseau à 1 neurone

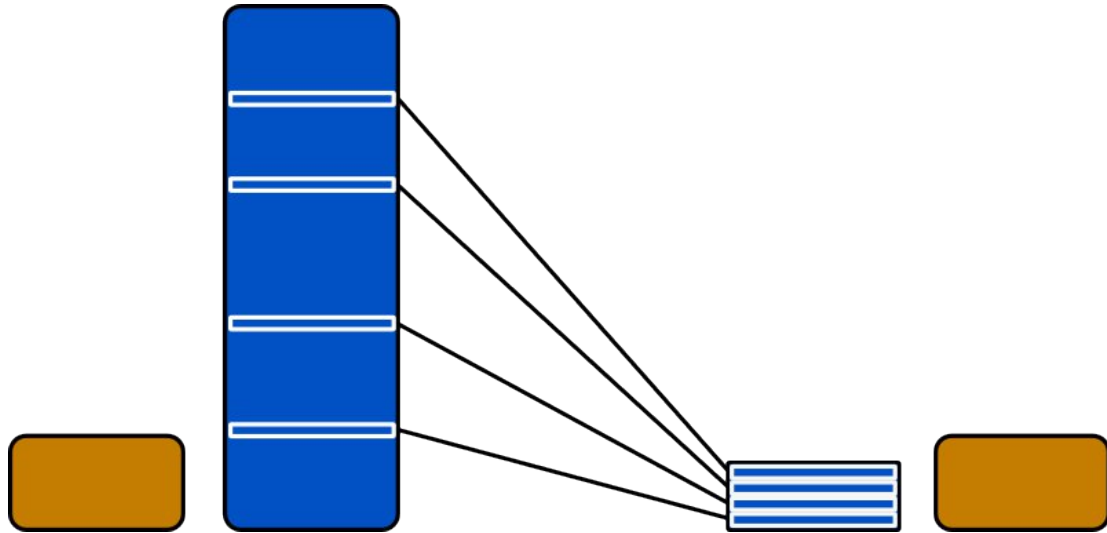




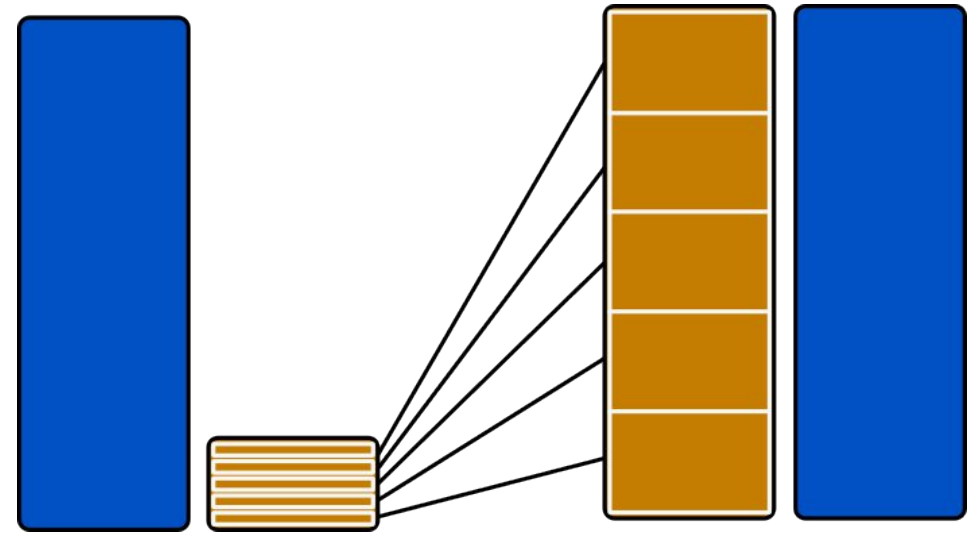


## Balance the classes - Example

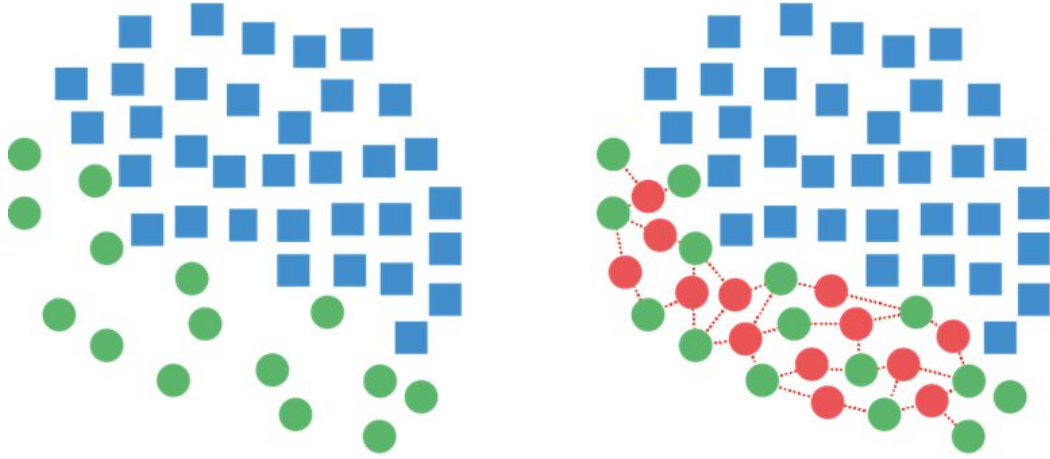
Undersampling



Oversampling



# Synthetic Minority Oversampling Technique



Data Augmentation



Original Image



De-texturized



De-colored



Edge Enhanced

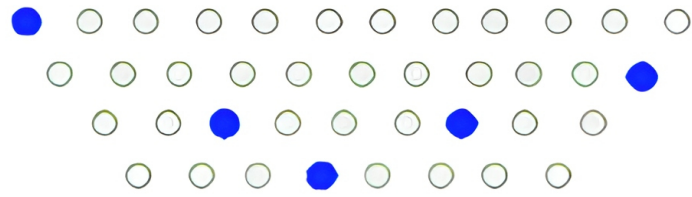


Salient Edge Map

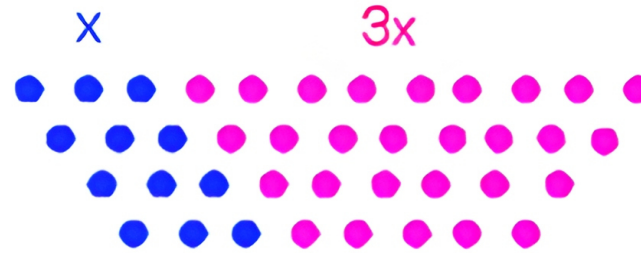


Flip/Rotate

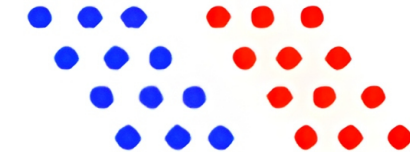
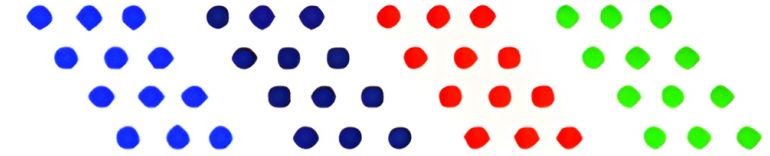
```
# Define the transformations
transform = transforms.Compose([
    transforms.RandomHorizontalFlip(),
    transforms.RandomRotation(10),
    transforms.RandomResizedCrop(224),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
])
```



Random Sampling



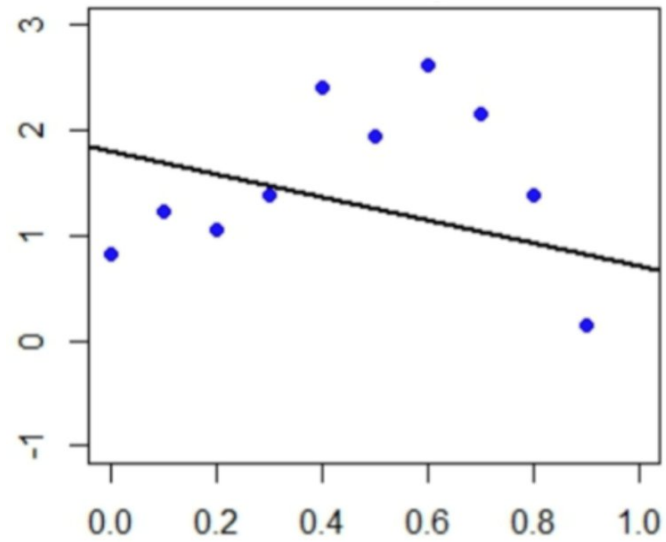
Stratified Sampling



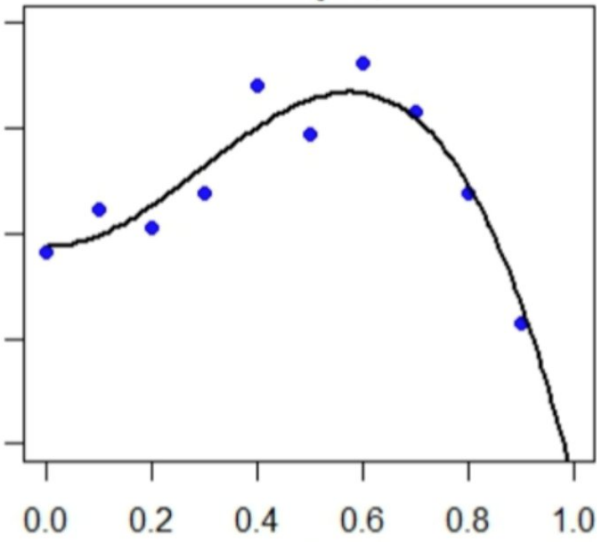
Cluster Sampling

- 1 Quel **travail** faire pour **améliorer** les **données** utilisées pour **l'entraînement**
- 2 **Comment évaluer un modèle ?**
- 3 Est-il possible de rendre **l'entraînement** plus **robuste** ?
- 4 Peut-on **profiter** d'un modèle **déjà entraîné** ?
- 5 Bonus : Quelques **bonnes pratiques** ?

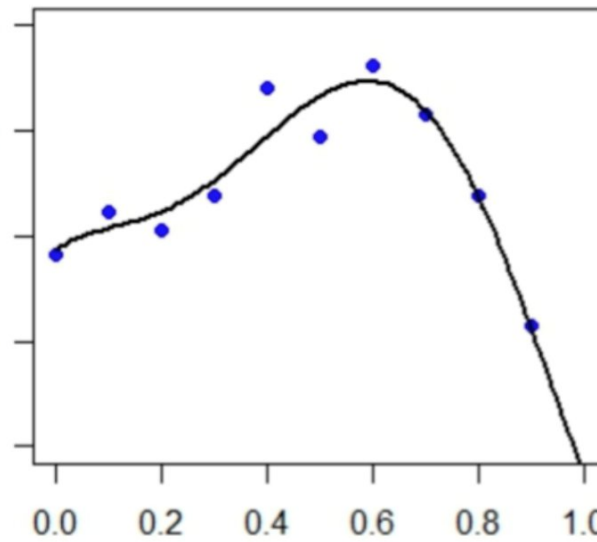
$$\hat{y}_i = a_0 + a_1x_i$$



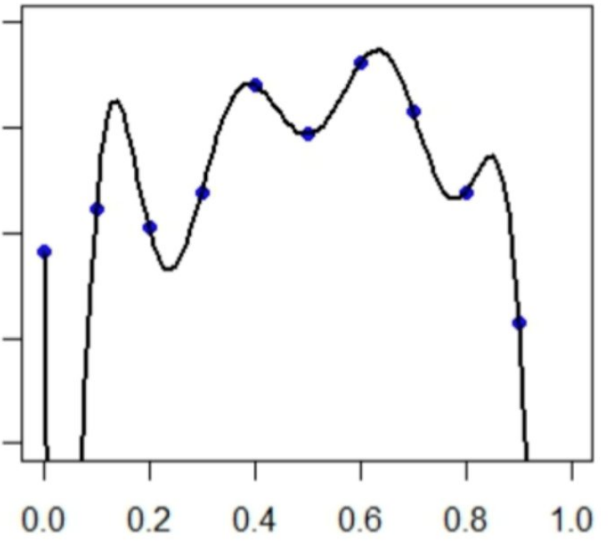
$$\hat{y}_i = a_0 + a_1x_i + a_2x_i^2 + a_3x_i^3$$



$$\hat{y}_i = a_0 + a_1x_i + \dots + a_5x_i^5$$



$$\hat{y}_i = a_0 + a_1x_i + \dots + a_{10}x_i^{10}$$







**Bias**



**Noise**

Noise or bias ?





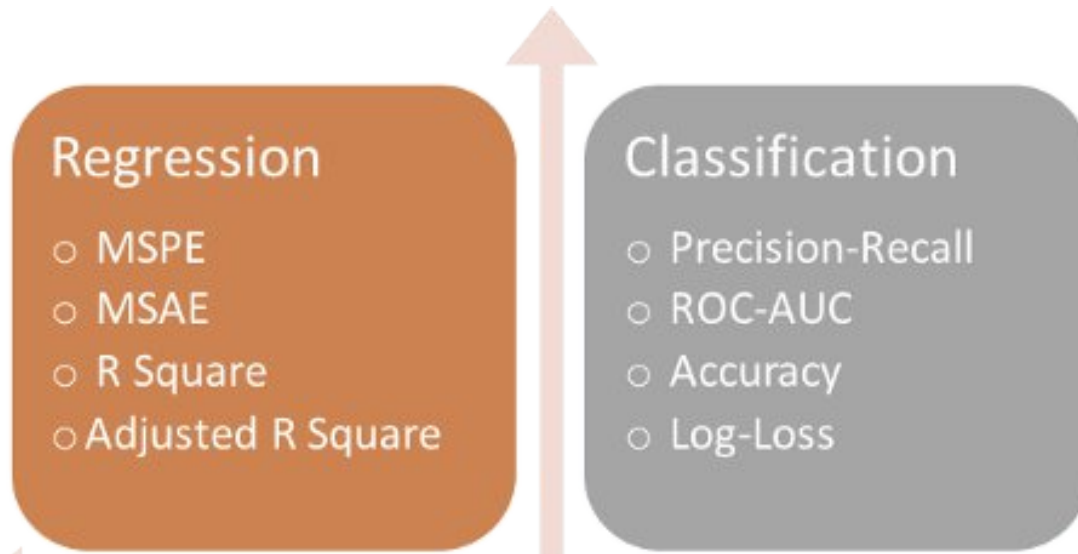
**Bias**



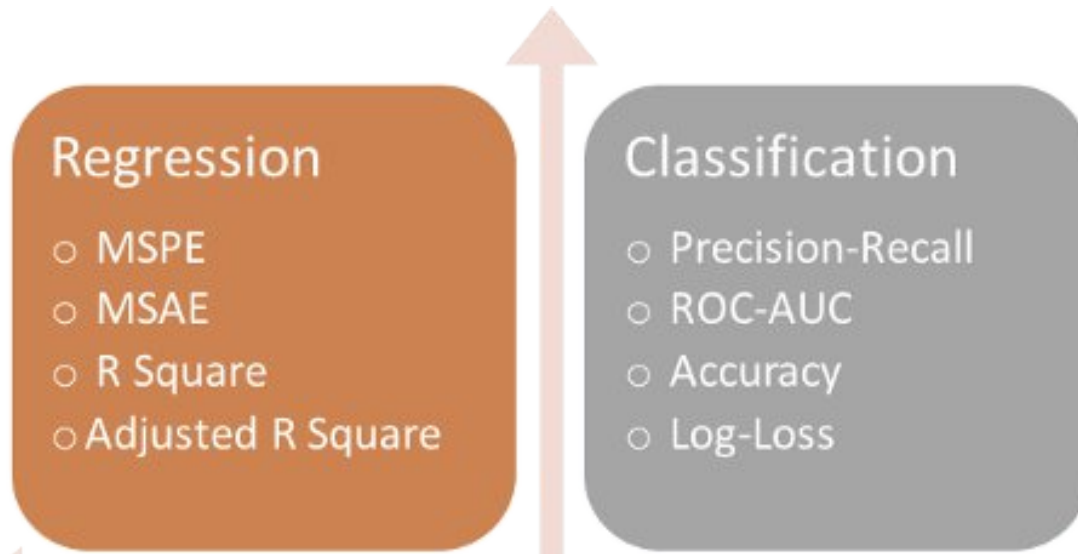
**Noise**

- My model says it's a cat no matter the image I give him
- There is no correlation between my model predictions and the label

**Noise or bias ?**



- **What is a metric?**
- **Is it different from a loss function?**



- **What is a metric?**
- **Is it different from a loss function?**
  - **Differentiability**
  - **Training vs evaluation**
  - **Number of samples**
  - **Interpretability/Meaning**



New system for illness detection - the accuracy is not a good metric for this case



All negative : accuracy = 99% ✓ ?

|                 |          | True Class |          |
|-----------------|----------|------------|----------|
|                 |          | Positive   | Negative |
| Predicted Class | Positive | TP         | FP       |
|                 | Negative | FN         | TN       |

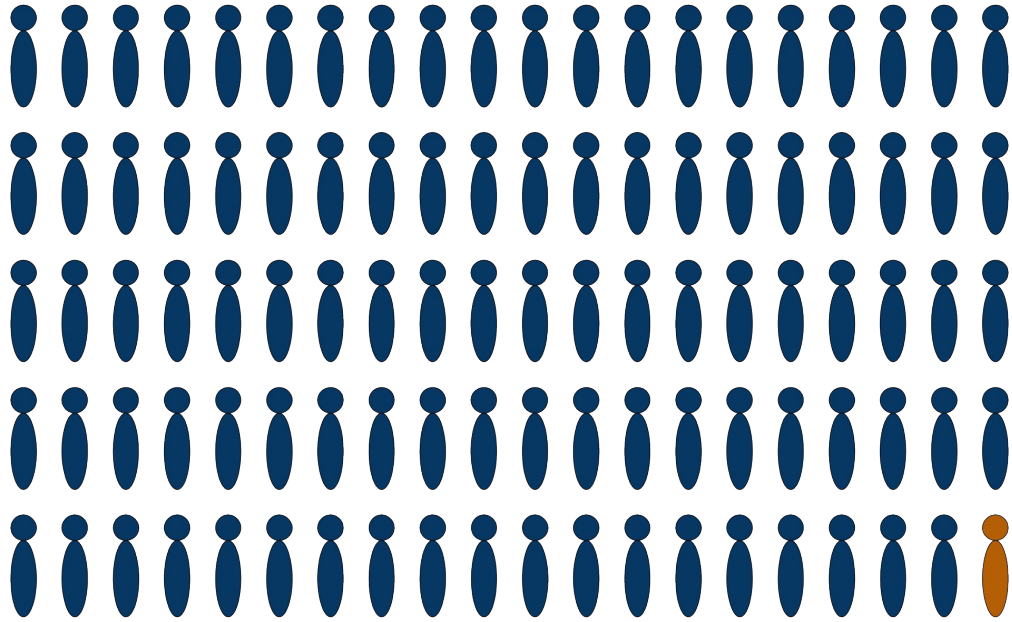
$$\text{precision} = \frac{TP}{TP + FP}$$

Above all positive prediction, how many are positive data

---


$$\text{recall} = \frac{TP}{TP + FN}$$

Above all positive data, how many have been predicted positive



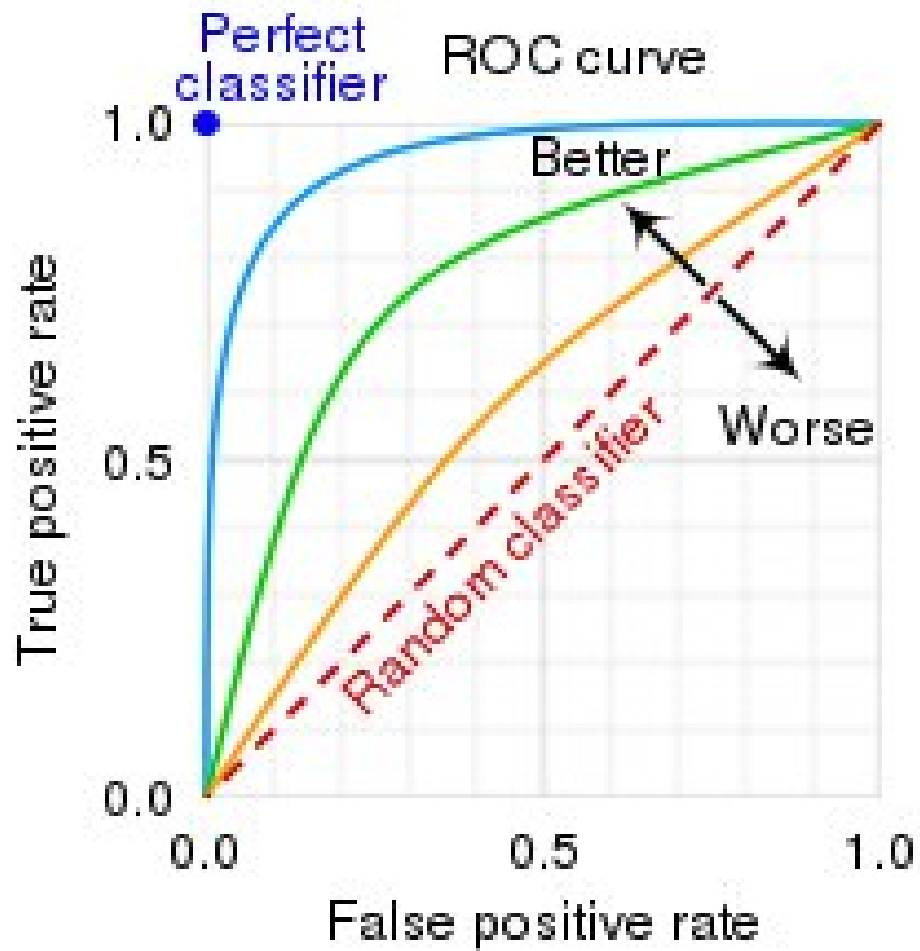
|                 |   | True class |    |
|-----------------|---|------------|----|
|                 |   | 0          | 1  |
| Predicted class | 0 | 0          | 0  |
|                 | 1 | 1          | 99 |

All negative :  
accuracy = 99%



precision = nan  
recall = 0





$$TPR = \frac{TP}{TP + FN}$$

Above all positive data,  
how many have been  
predicted positive

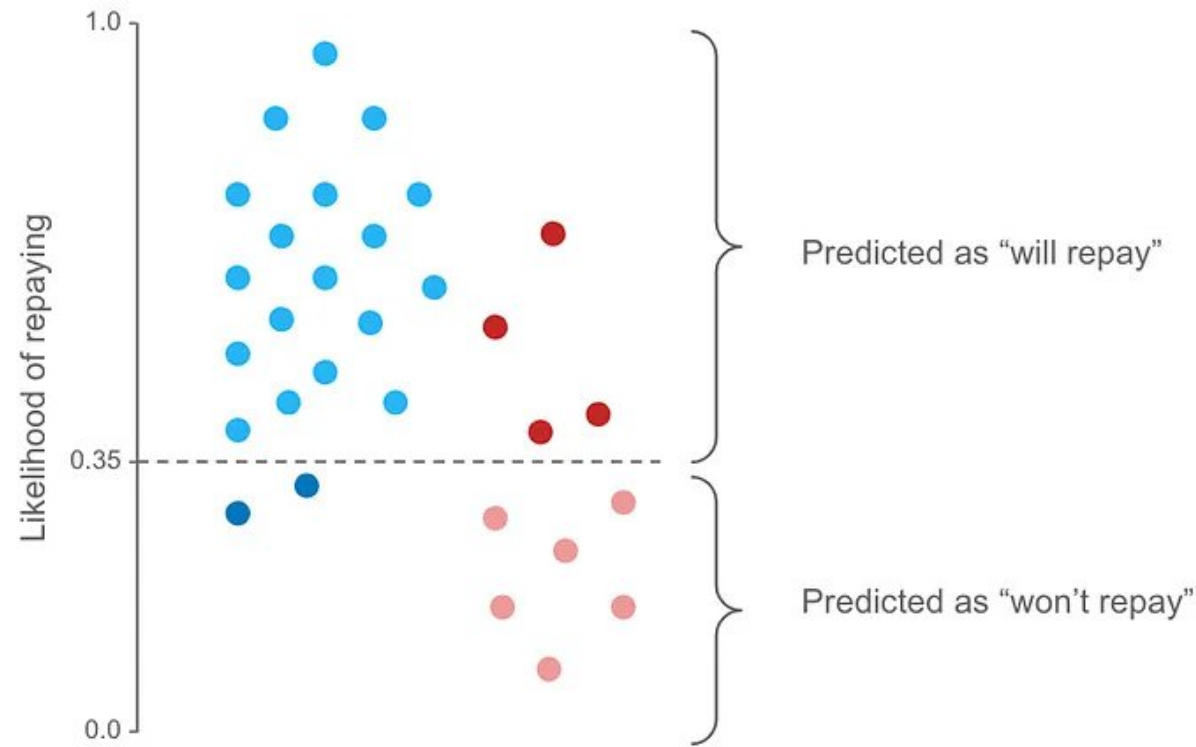
$$FPR = \frac{FP}{FP + TN}$$

Above all negative data,  
how many have been  
predicted positive

Variation of the acceptance threshold of a class  
to obtain the curve

NB : ROC = Receiver Operating Characteristics

## Metrics - ROC curve



Actual positives: *users who repaid the loan*

● Predicted as "will repay"

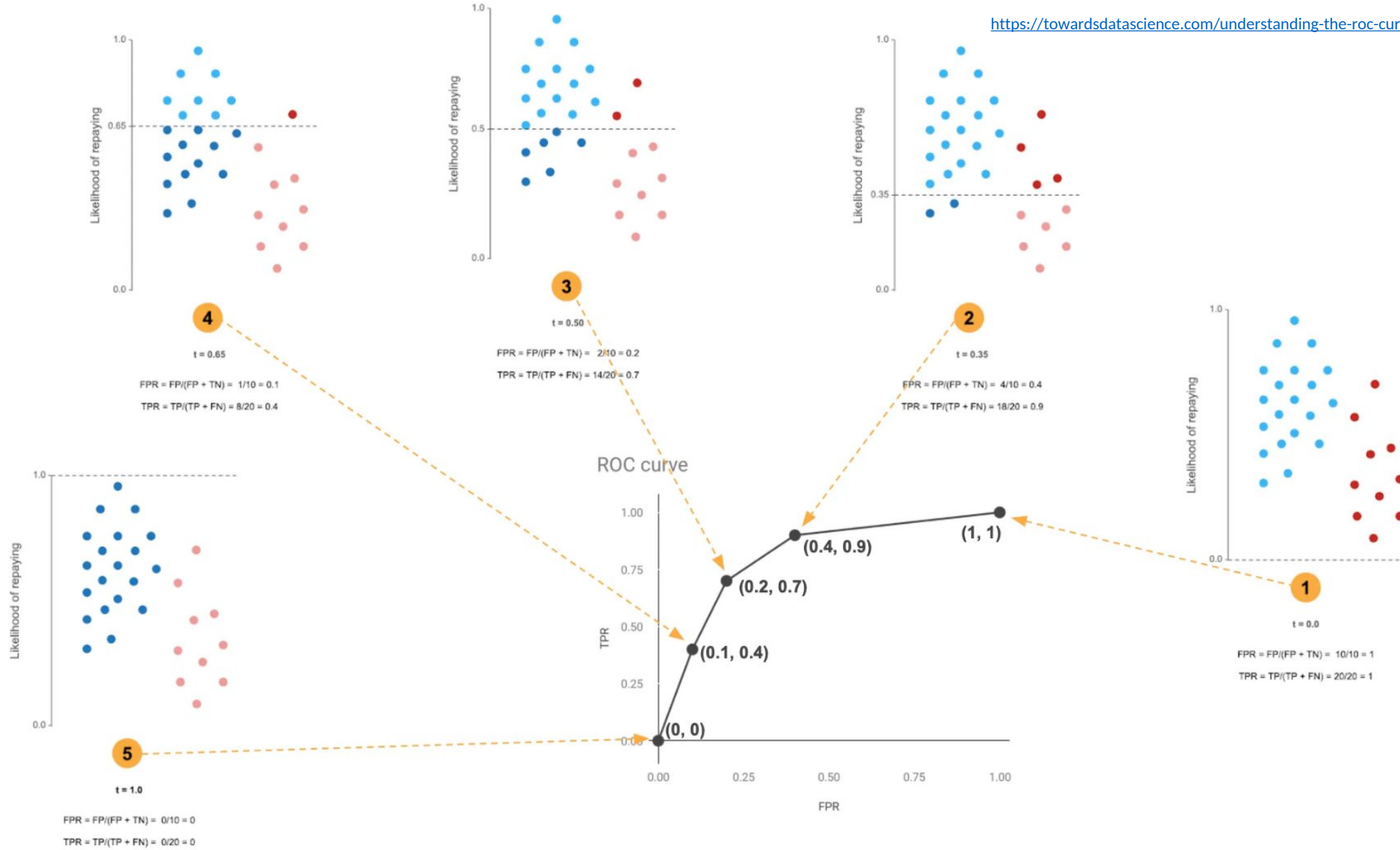
● Predicted as "won't repay"

Actual negatives: *users who didn't repaid the loan*

● Predicted as "won't repay"

● Predicted as "will repay"

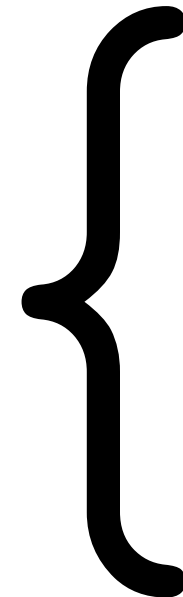




**Perplexity :**  
Is model surprised to see this text?

$$\text{PPL}(X) = \exp \left\{ -\frac{1}{t} \sum_i^t \log p_{\theta}(x_i | x_{<i}) \right\}$$

This person is innocent, she needs a lawyer



|             |           |
|-------------|-----------|
| “ssistance” | - 1.5%    |
| “ lawyer”   | - 0.5%    |
| “ car”      | - 0.05%   |
| ...         |           |
| “ sea”      | - 0.0001% |

# Perplexity is not All you Need !!

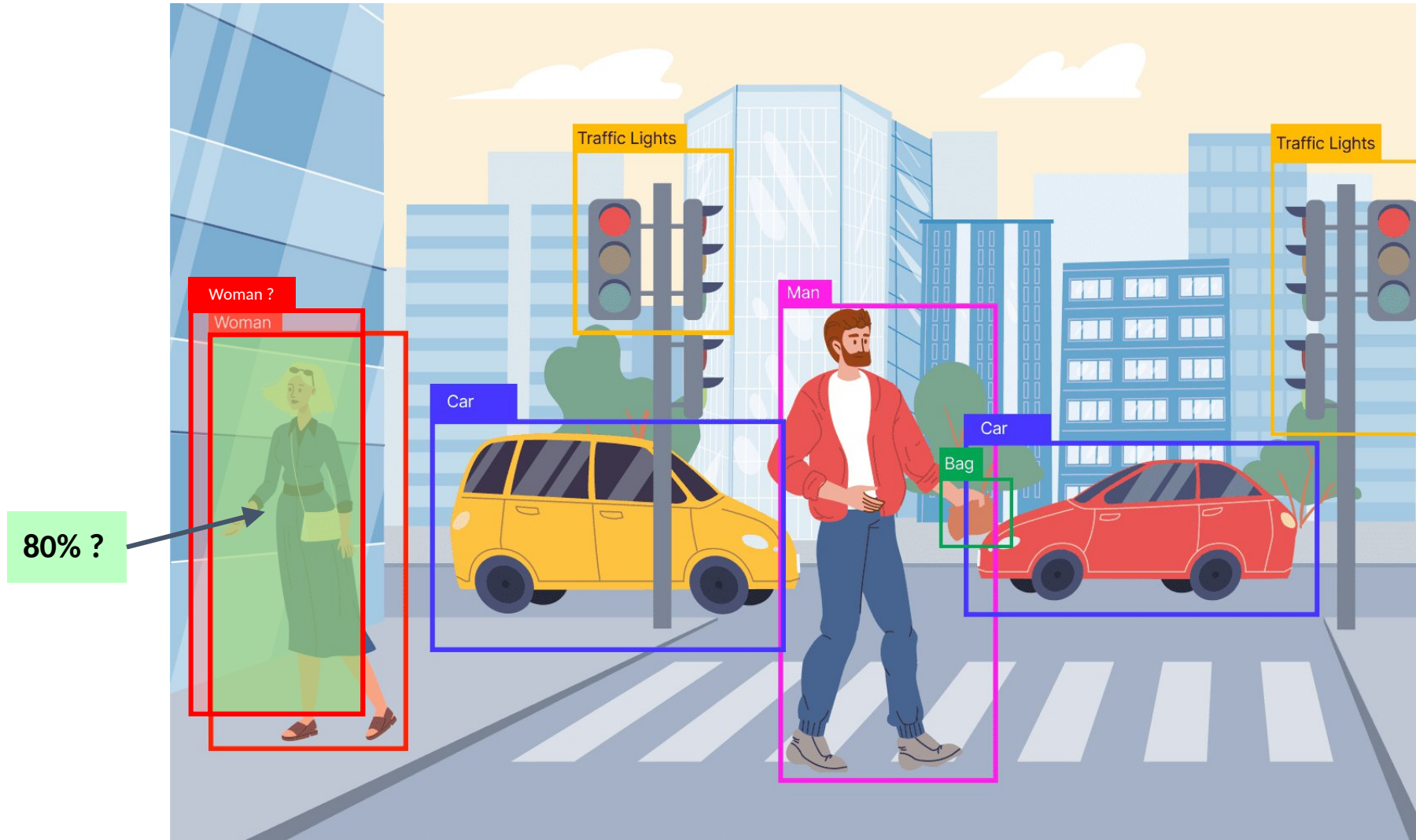
| T ▲ | Model ▲  | Average 📈 ▼ | ARC ▲ | HellaSwag ▲ | MMLU ▲ | TruthfulQA ▲ | Winogrande ▲ | GSM8K ▲ | Architecture ▲     | #Params (B) ▲ |
|-----|--|-------------|-------|-------------|--------|--------------|--------------|---------|--------------------|---------------|
| ■   | <a href="#">moreh/MoMo-70B-lora-1.8.6-DPO</a>          | 77.29       | 70.14 | 86.03       | 77.4   | 69           | 84.37        | 76.8    | LlamaForCausalLM   | 72.29         |
| ■   | <a href="#">moreh/MoMo-70B-lora-1.8.4-DPO</a>          | 76.23       | 69.62 | 85.35       | 77.33  | 64.64        | 84.14        | 76.27   | LlamaForCausalLM   | 72.29         |
| ◆   | <a href="#">TomGrc/FusionNet_7Bx2_MoE_14B</a>          | 75.91       | 73.55 | 88.84       | 64.68  | 69.6         | 88.16        | 70.66   | MixtralForCausalLM | 12.88         |
| ◆   | <a href="#">Weyaxi/Helion-4x34B</a>                    | 75.48       | 69.71 | 85.28       | 77.33  | 63.91        | 84.37        | 72.25   | MixtralForCausalLM | 113.66        |
| ◆   | <a href="#">one-man-army/UNA-34Beagles-32K-bf16-v1</a> | 75.41       | 73.55 | 85.93       | 76.45  | 73.55        | 82.95        | 60.05   | LlamaForCausalLM   | 34.39         |
| ◆   | <a href="#">Weyaxi/Cosmosis-3x34B</a>                  | 75.39       | 69.71 | 85.18       | 77.25  | 63.82        | 84.14        | 72.25   | MixtralForCausalLM | 87.24         |
| ◆   | <a href="#">Weyaxi/Bagel-Hermes-2x34b</a>              | 75.1        | 69.8  | 85.26       | 77.24  | 64.82        | 84.77        | 68.69   | MixtralForCausalLM | 60.81         |
| ○   | <a href="#">jondurbin/bagel-dpo-34b-v0.2</a>           | 74.69       | 71.93 | 85.25       | 76.58  | 70.05        | 83.35        | 60.96   | LlamaForCausalLM   | 34.39         |
| ◆   | <a href="#">jondurbin/nontoxic-bagel-34b-v0.2</a>      | 74.69       | 72.44 | 85.64       | 76.41  | 72.7         | 82.48        | 58.45   | LlamaForCausalLM   | 34.39         |
| ◆   | <a href="#">moreh/MoMo-70B-LoRA-V1.4</a>               | 74.67       | 69.2  | 85.07       | 77.12  | 62.66        | 83.74        | 70.2    | LlamaForCausalLM   | 72.29         |
| ■   | <a href="#">udkai/Turdus</a>                           | 74.66       | 73.38 | 88.56       | 64.52  | 67.11        | 86.66        | 67.7    | MistralForCausalLM | 7.24          |

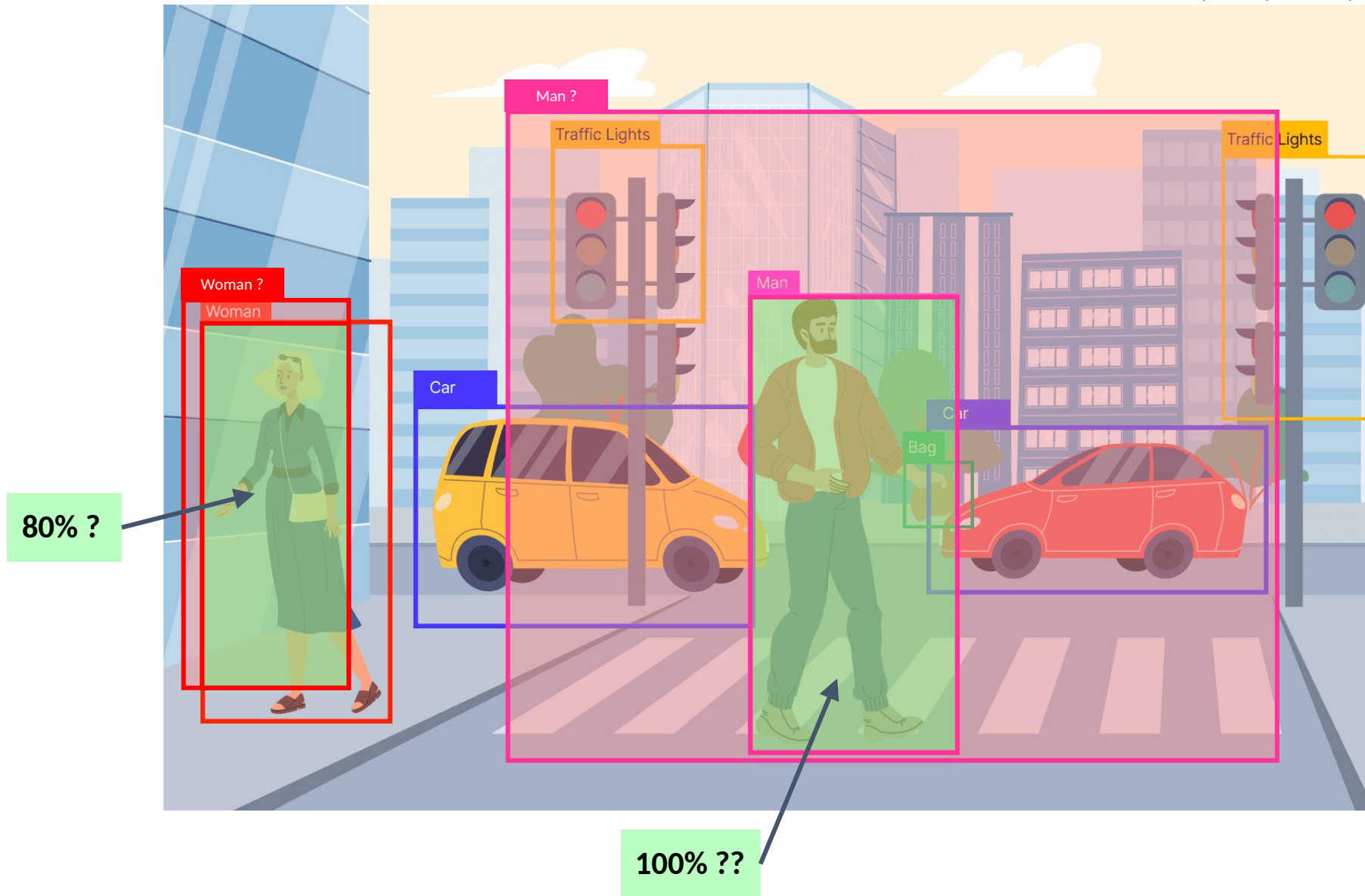
## Language models evaluation benchmark

The evaluation is complicated, there are many different tasks

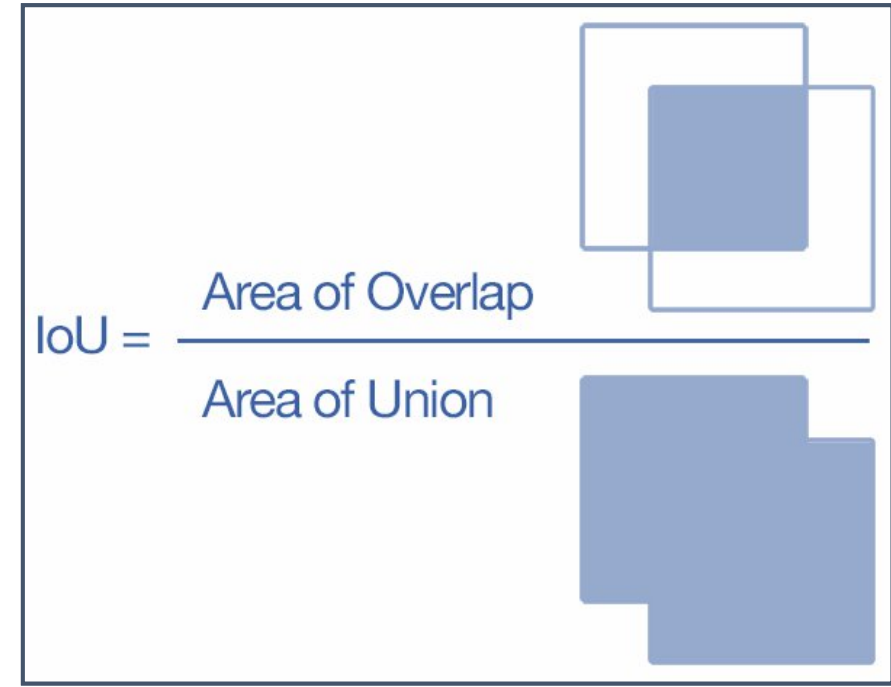
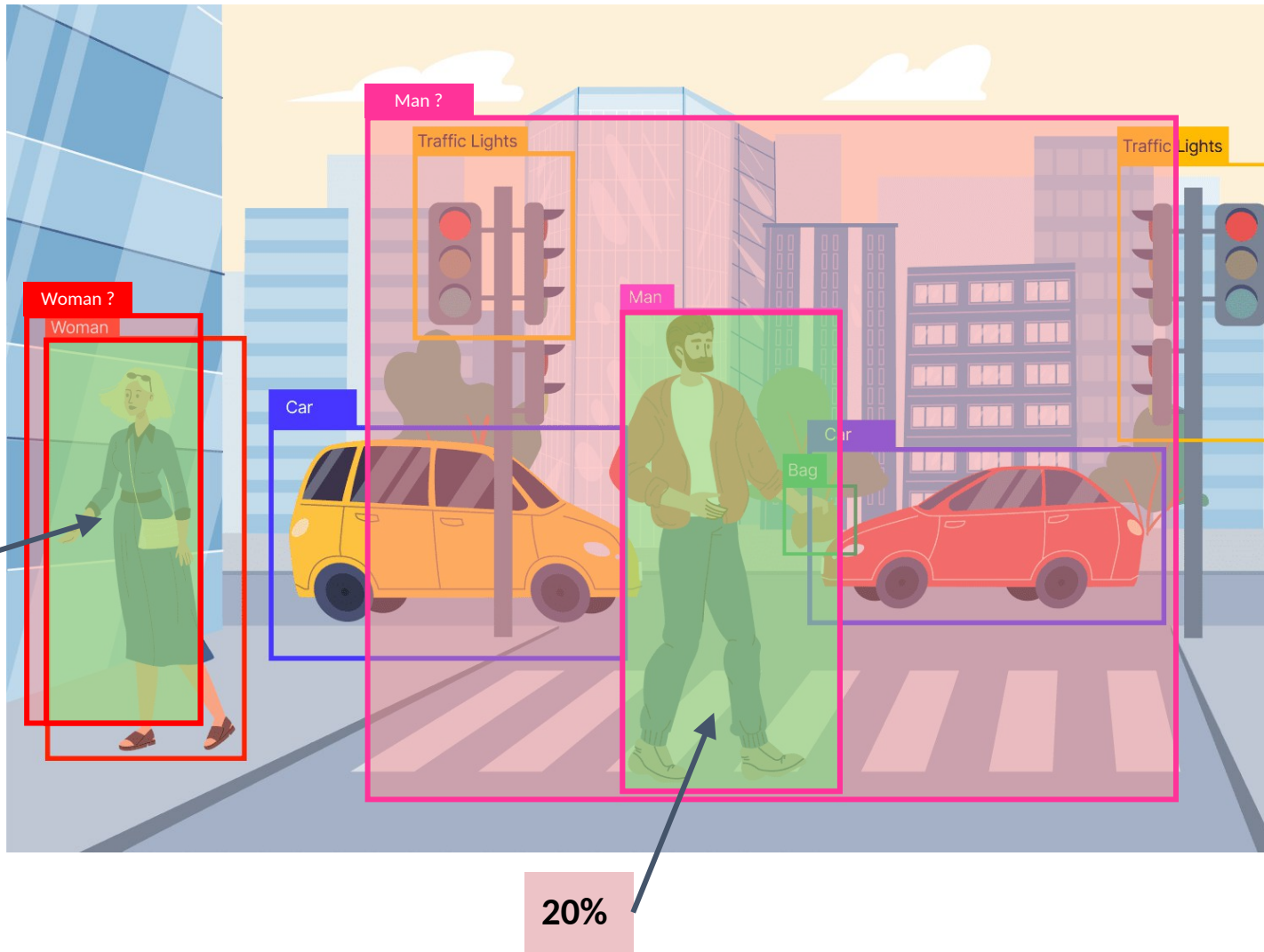
It can depends on the targeted application

It is hard to be fair/objective without knowing the training data





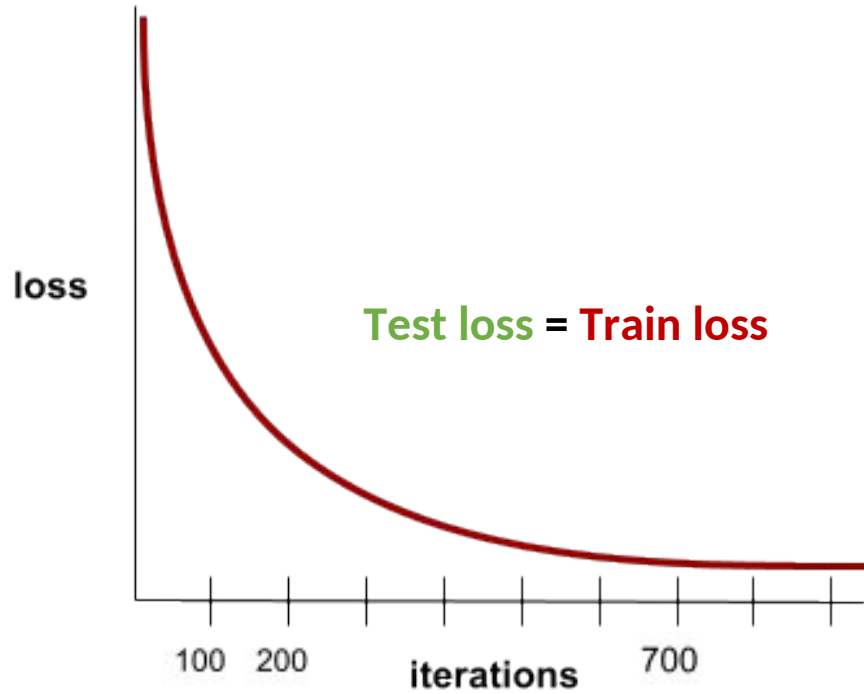
## Metrics - Intersection over Union



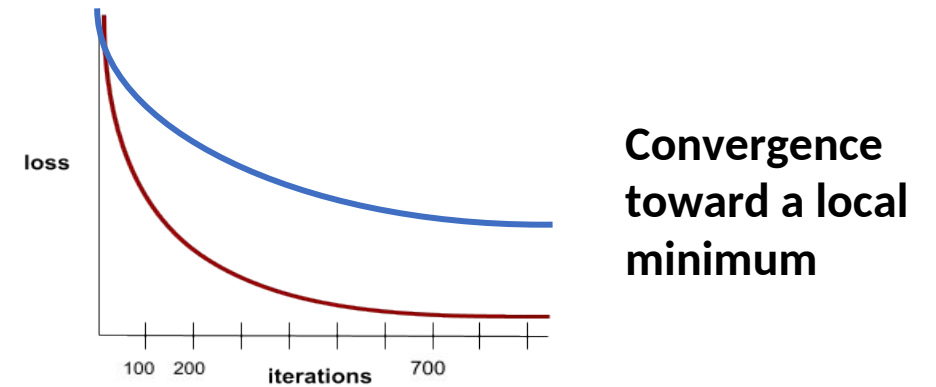
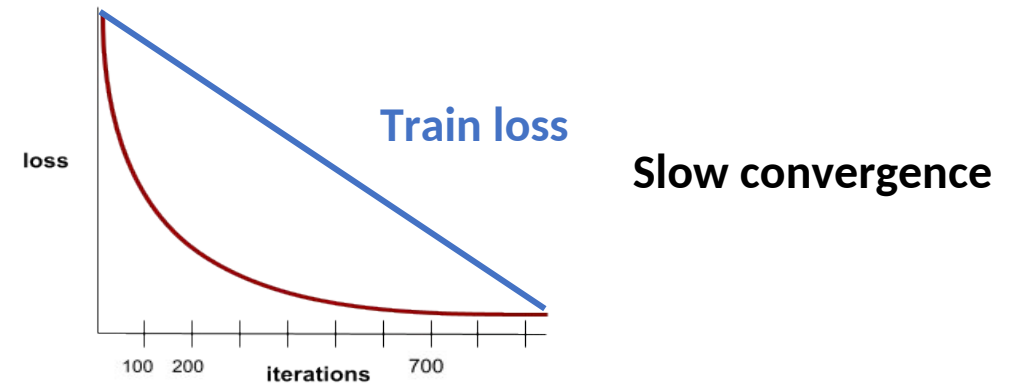
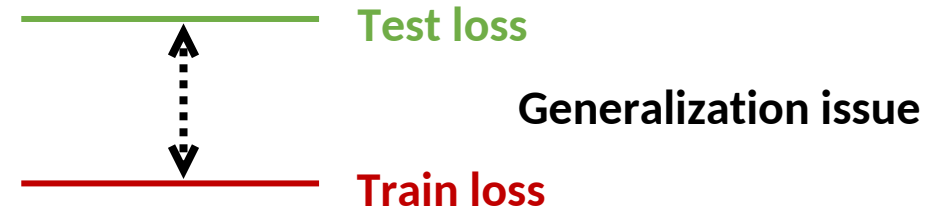
**Intersection over Union**

- 1 Quel **travail** faire pour **améliorer** les **données** utilisées pour **l'entraînement**
- 2 Comment **évaluer** un **modèle** ?
- 3 Est-il possible de rendre **l'entraînement** plus **robuste** ?
- 4 Peut-on **profiter** d'un modèle **déjà entraîné** ?
- 5 Bonus : Quelques **bonnes pratiques** ?

## Expectation

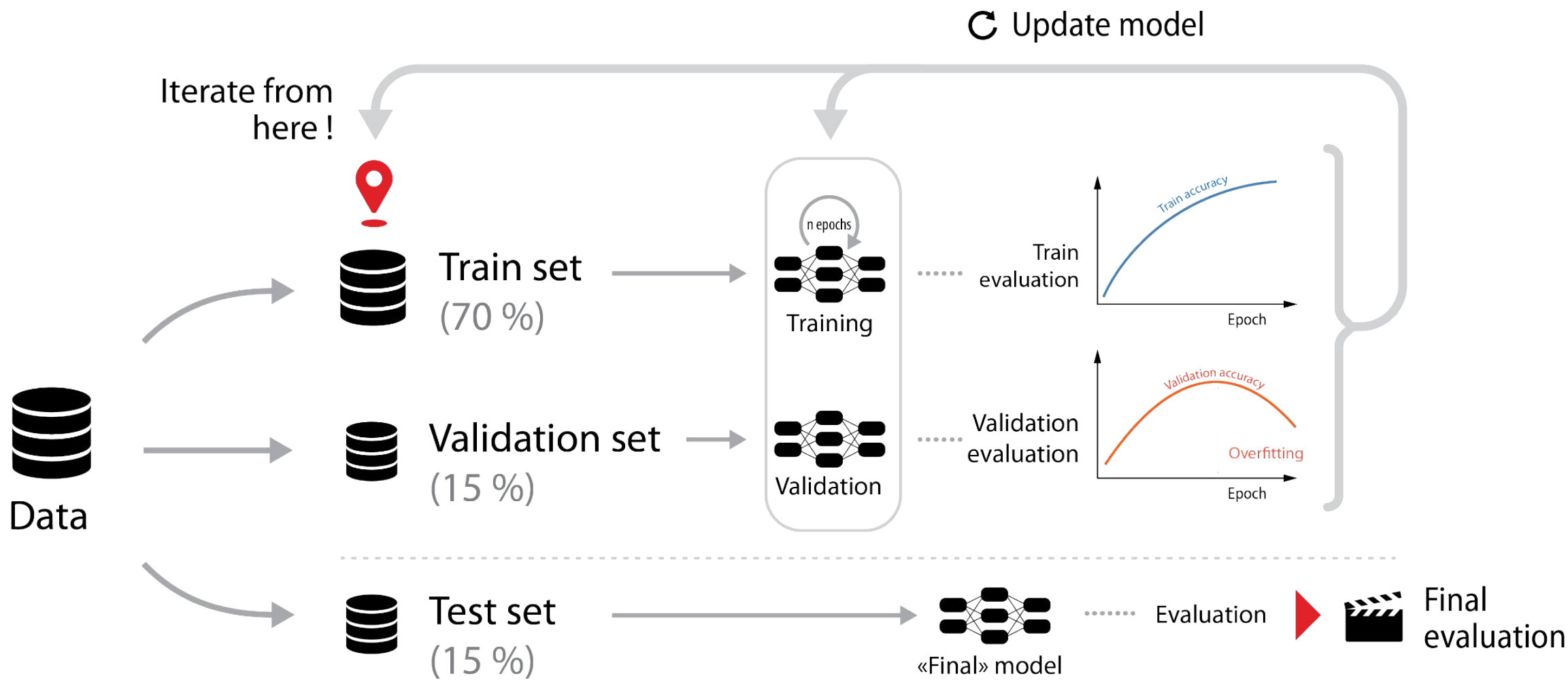


## Reality



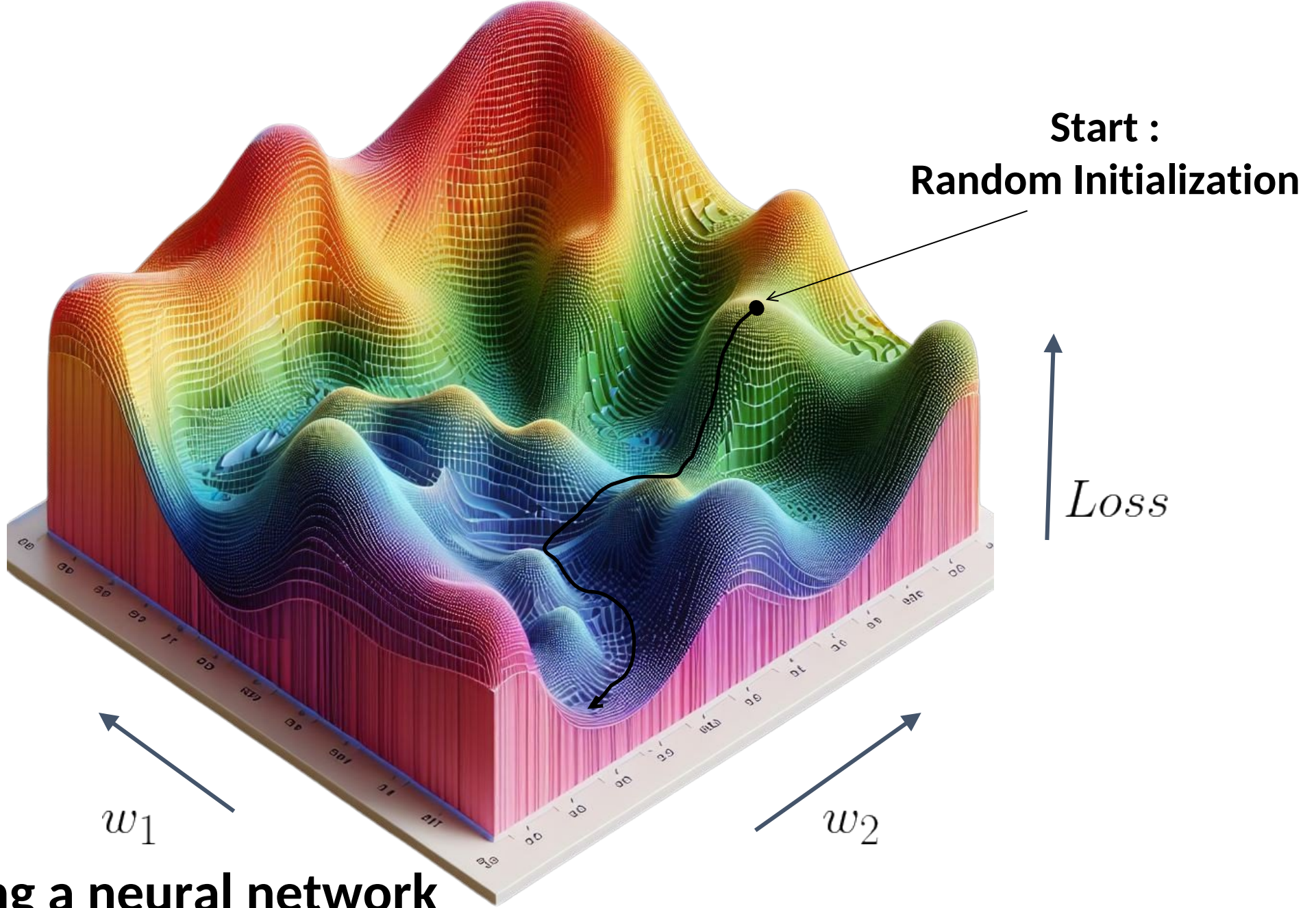
What happens during training ?





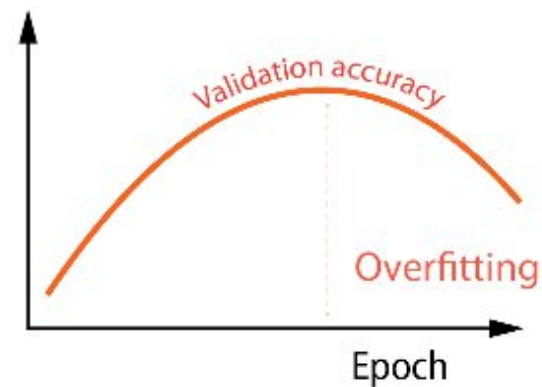
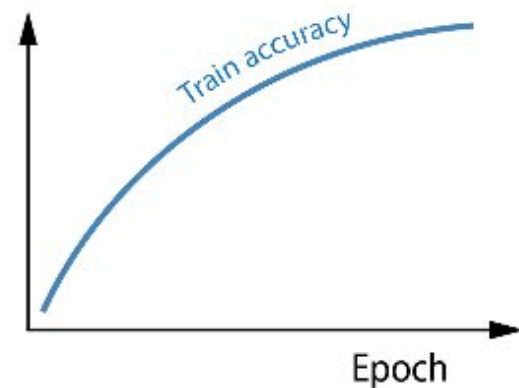
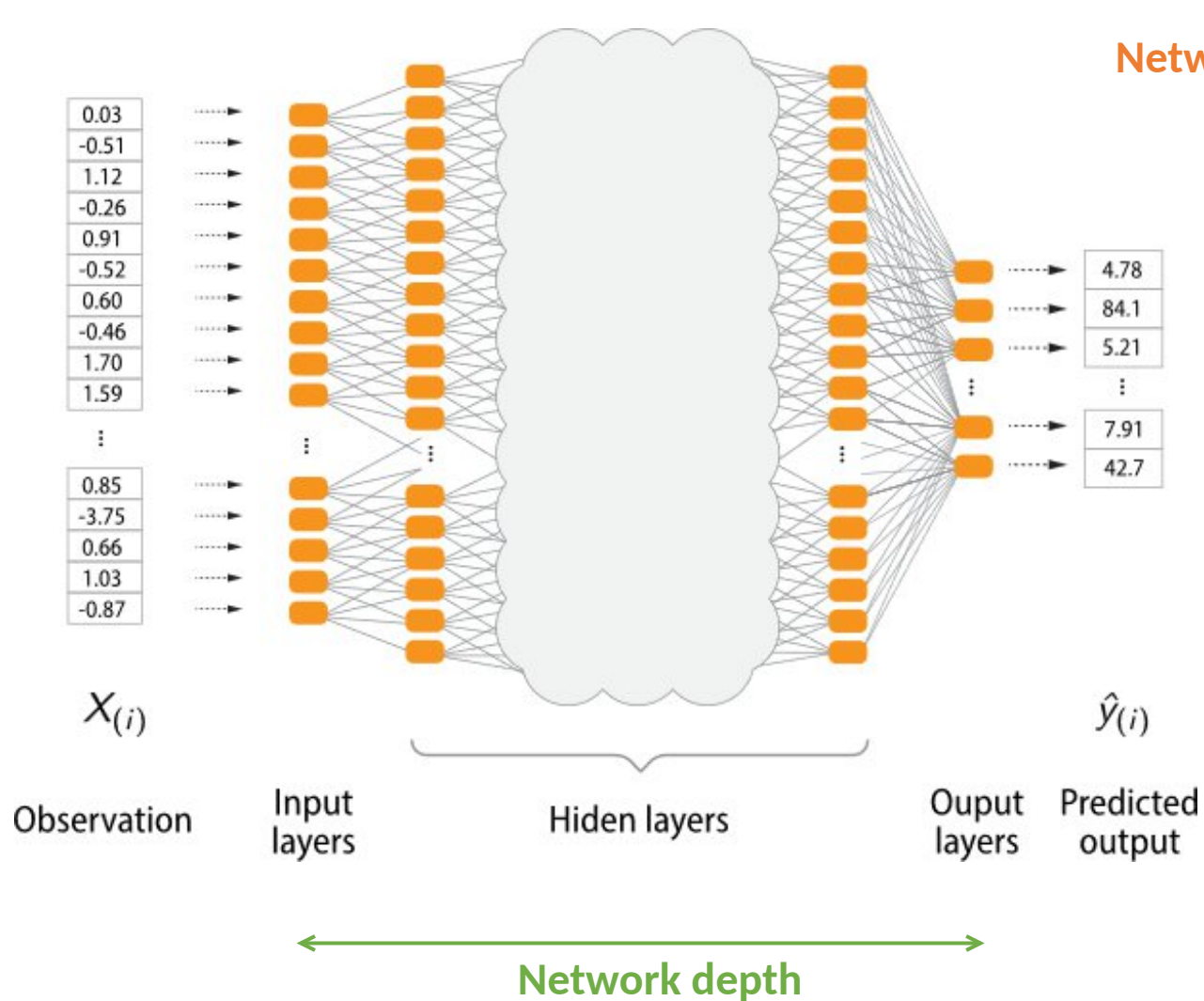
# Data splitting

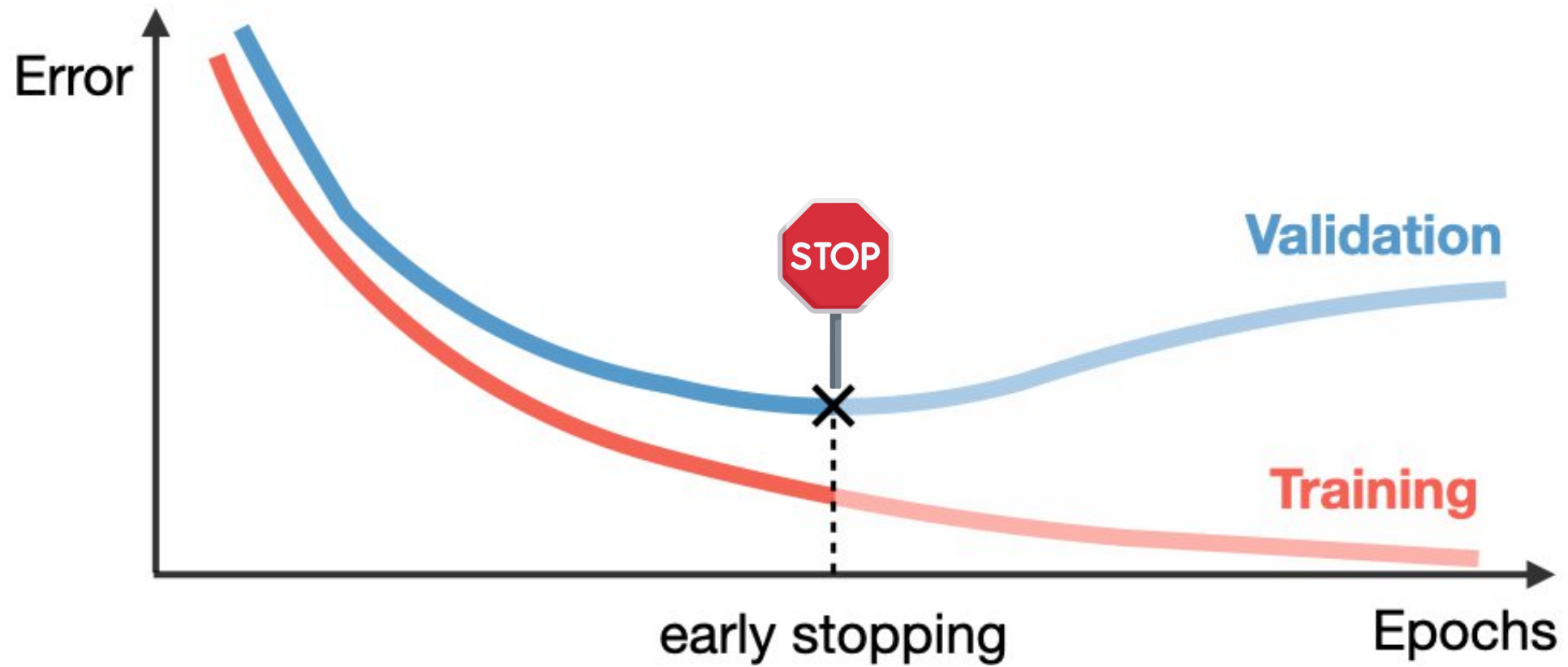
# Gradient Descent



## Training a neural network







## Basic way : Early stopping

## Hyper-parameters

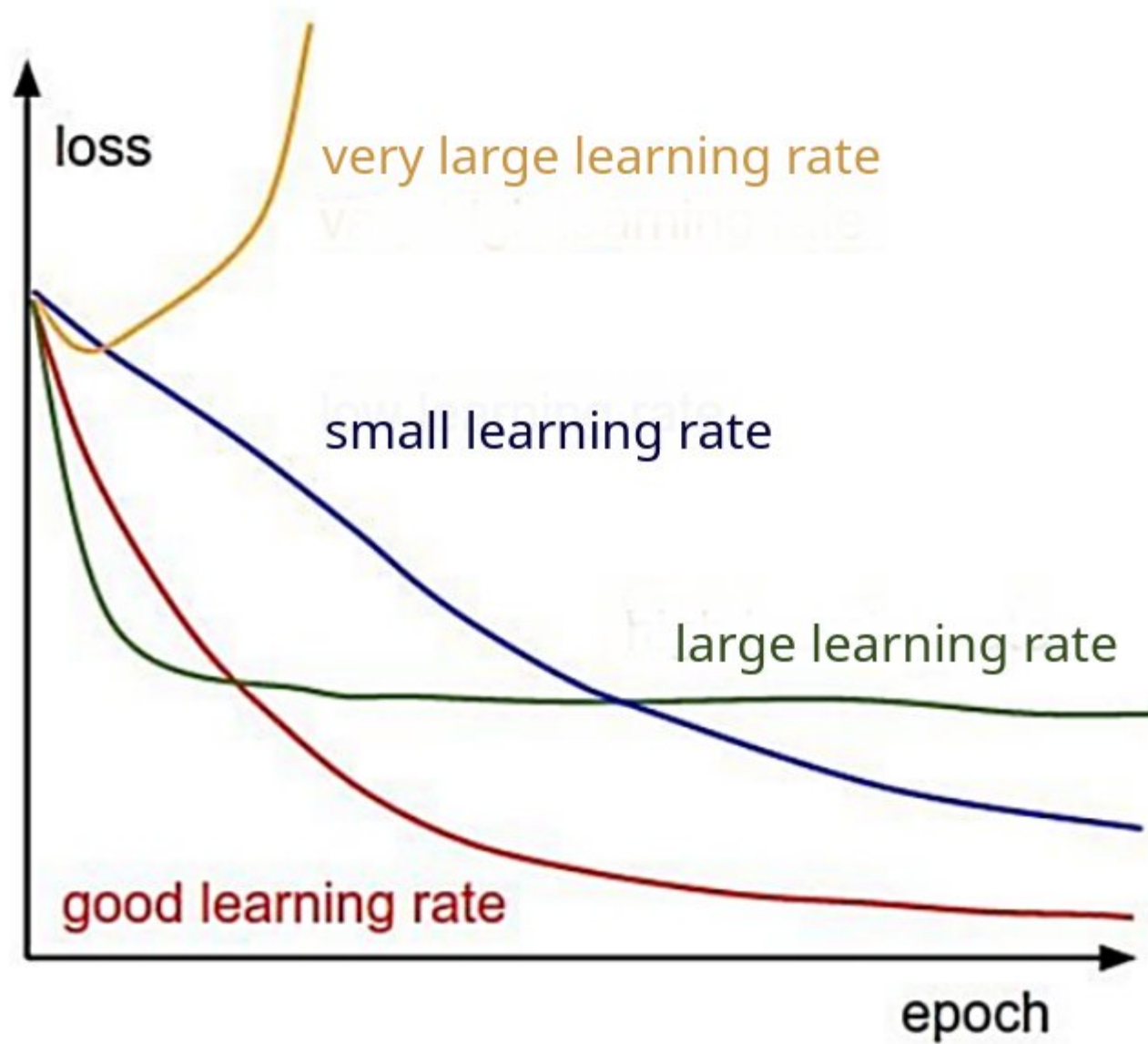
- Learning rate
- Regularization
- Optimizer
- Model architecture
- Batch size
- ...

$$\Theta_{t+1} = \Theta_t - \eta \nabla_{\Theta} [\mathcal{L}(\hat{y}_i, y_i)]$$

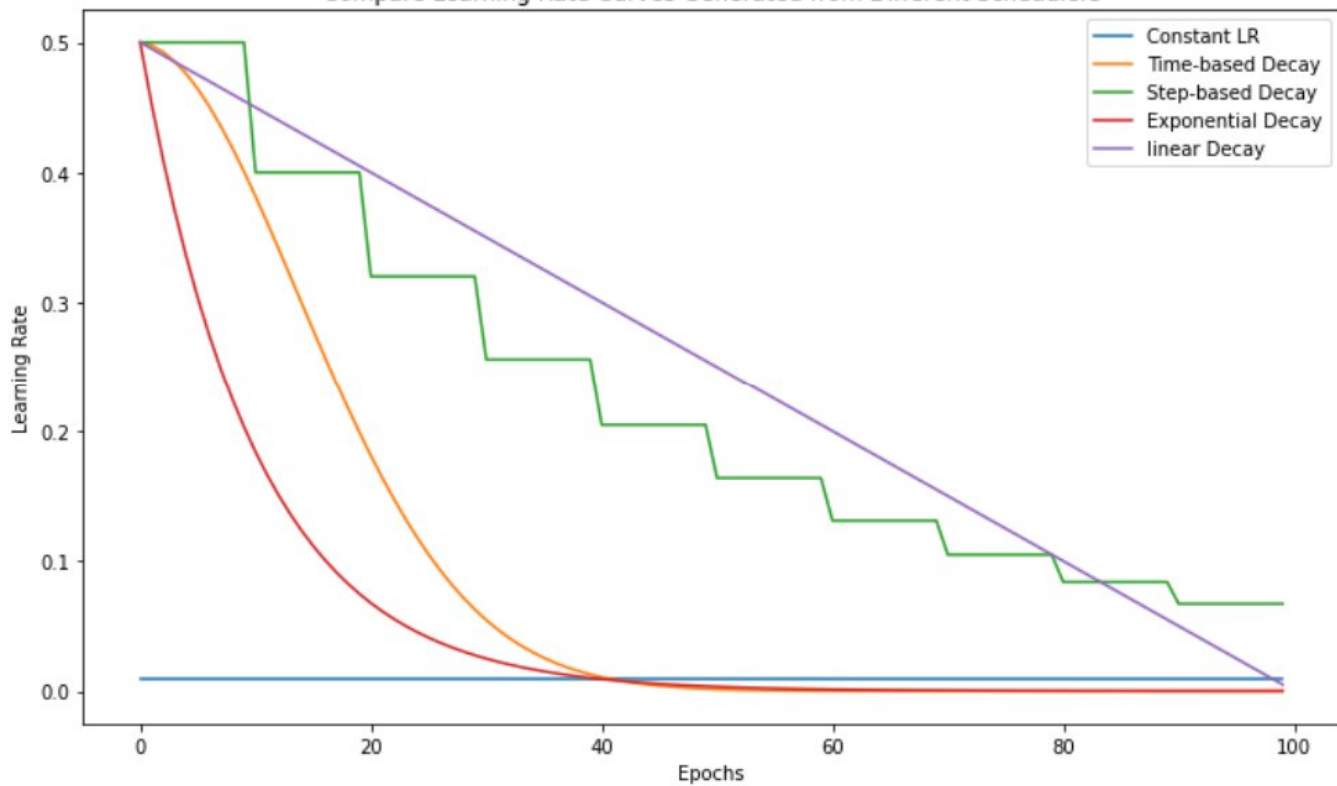
Updated weights = Weights before update - Learning rate \* Gradient [ Cost function ( Prediction, Label ) ]

Weight update equation

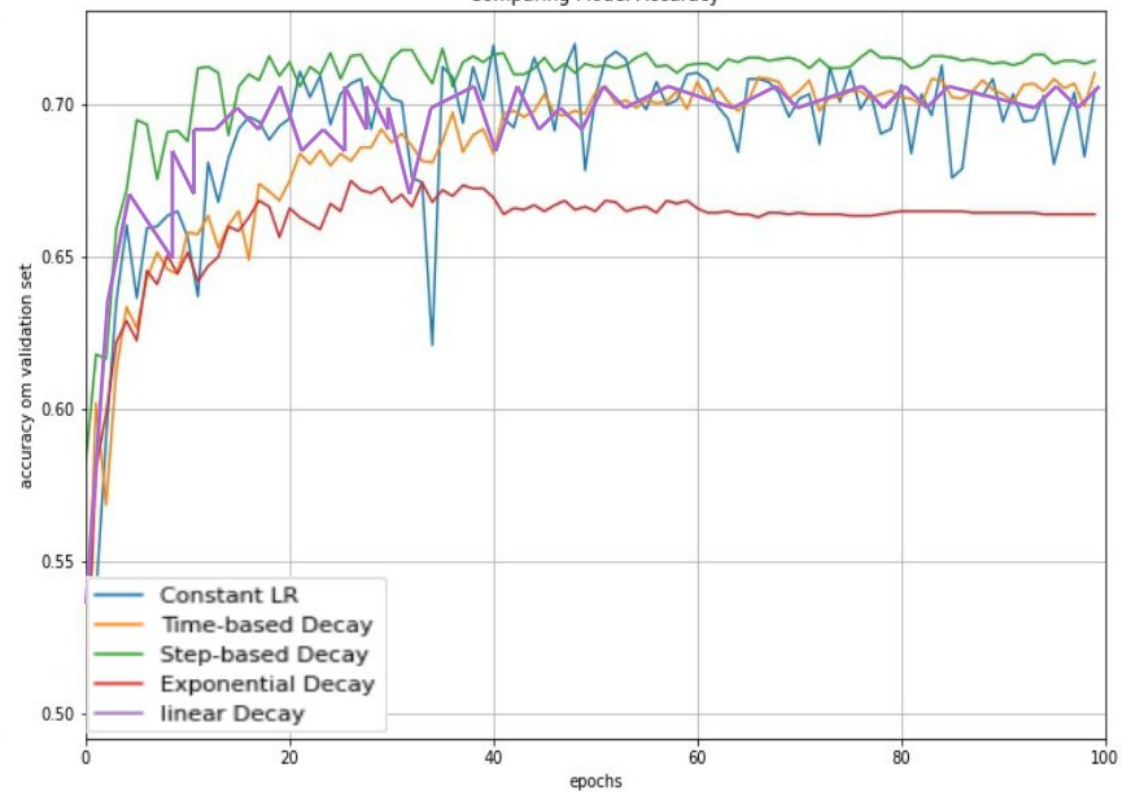
Regularization - Loss function and weight update



Compare Learning Rate Curves Generated from Different Schedulers

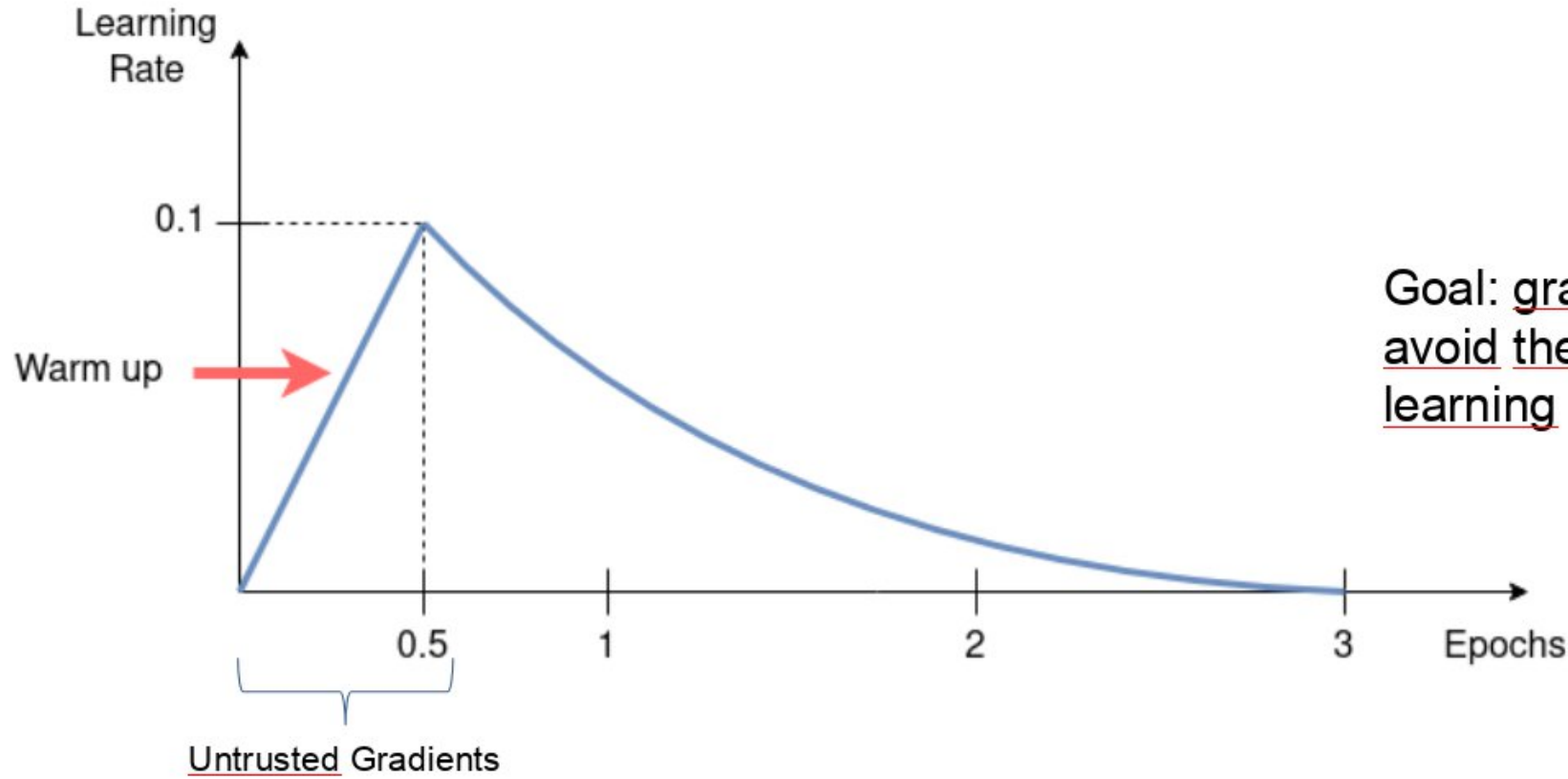


Comparing Model Accuracy



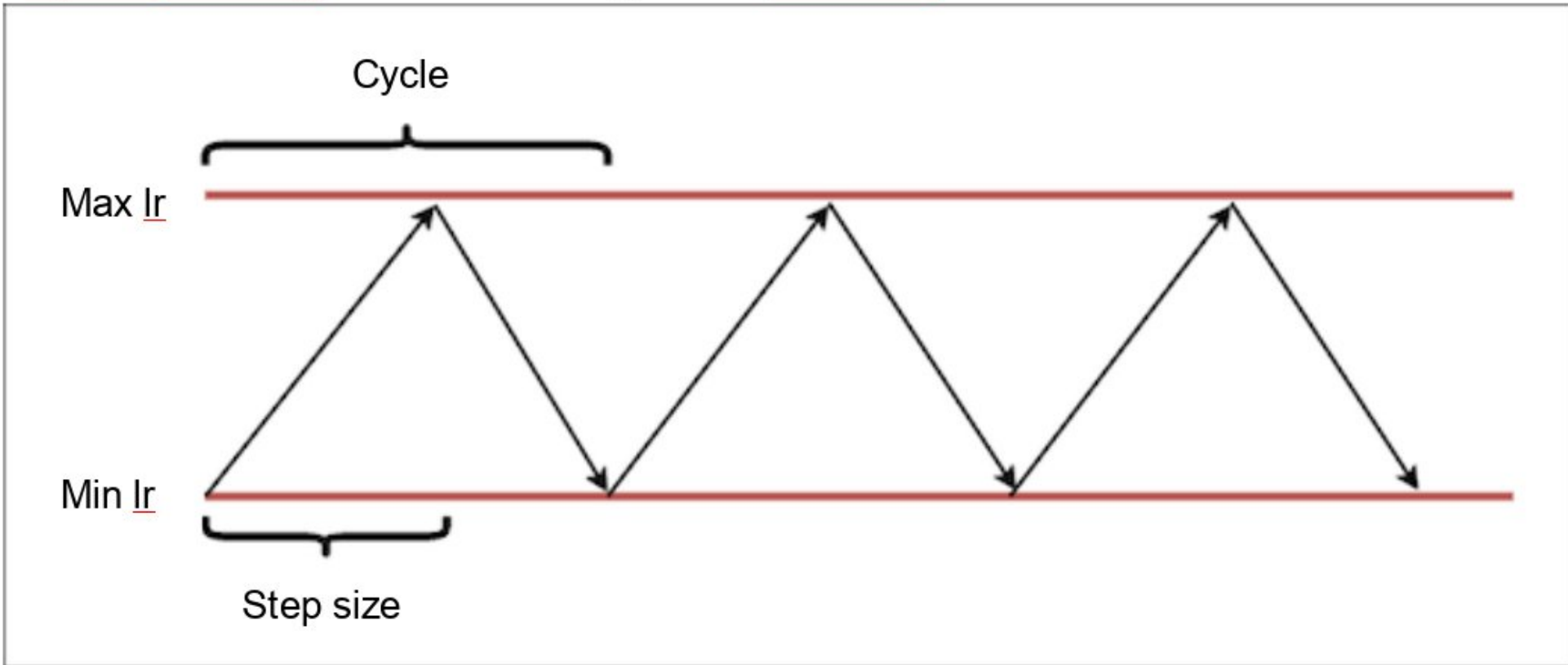


**Problems** : The first iterations have too much effect on the model (significant losses, high gradients, bias, etc.)  
A high learning rate can cause strong instability or divergence



Goal: gradually increase the learning rate to avoid the risk of divergence at the start of learning

**Learning rate scheduler : warmup**

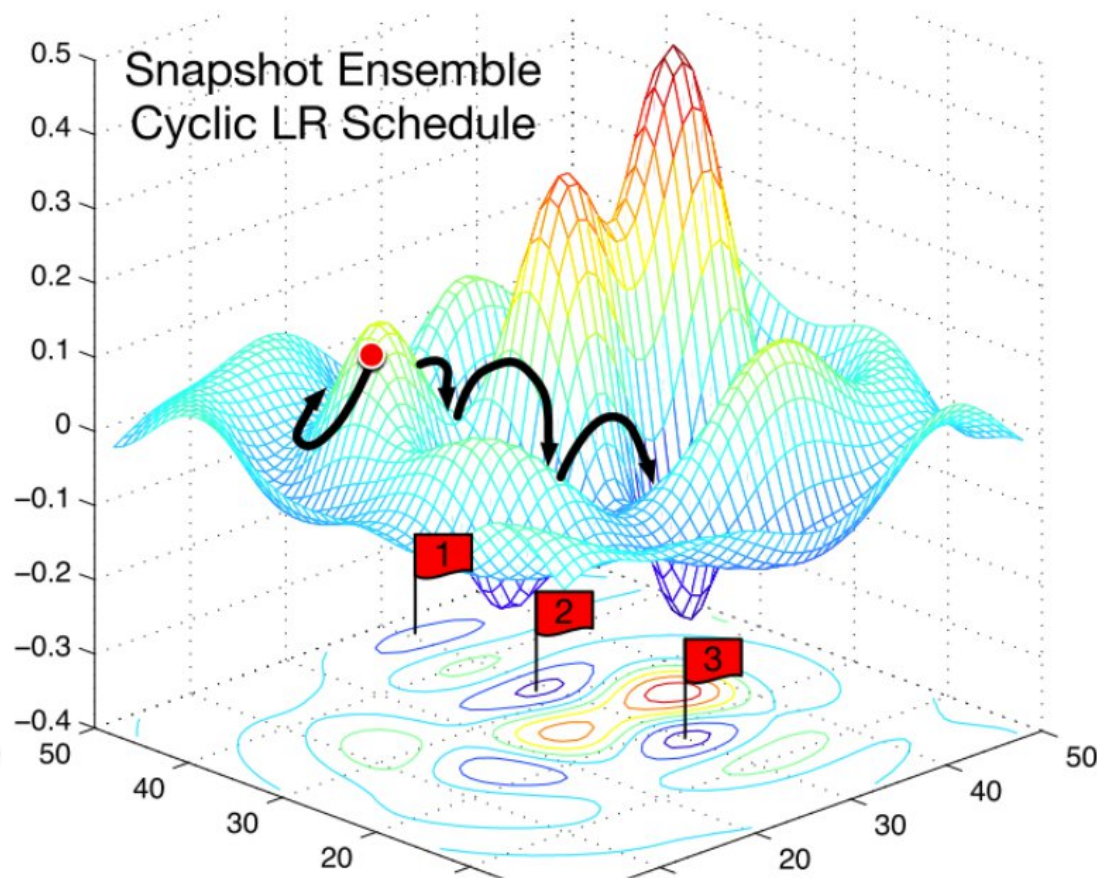
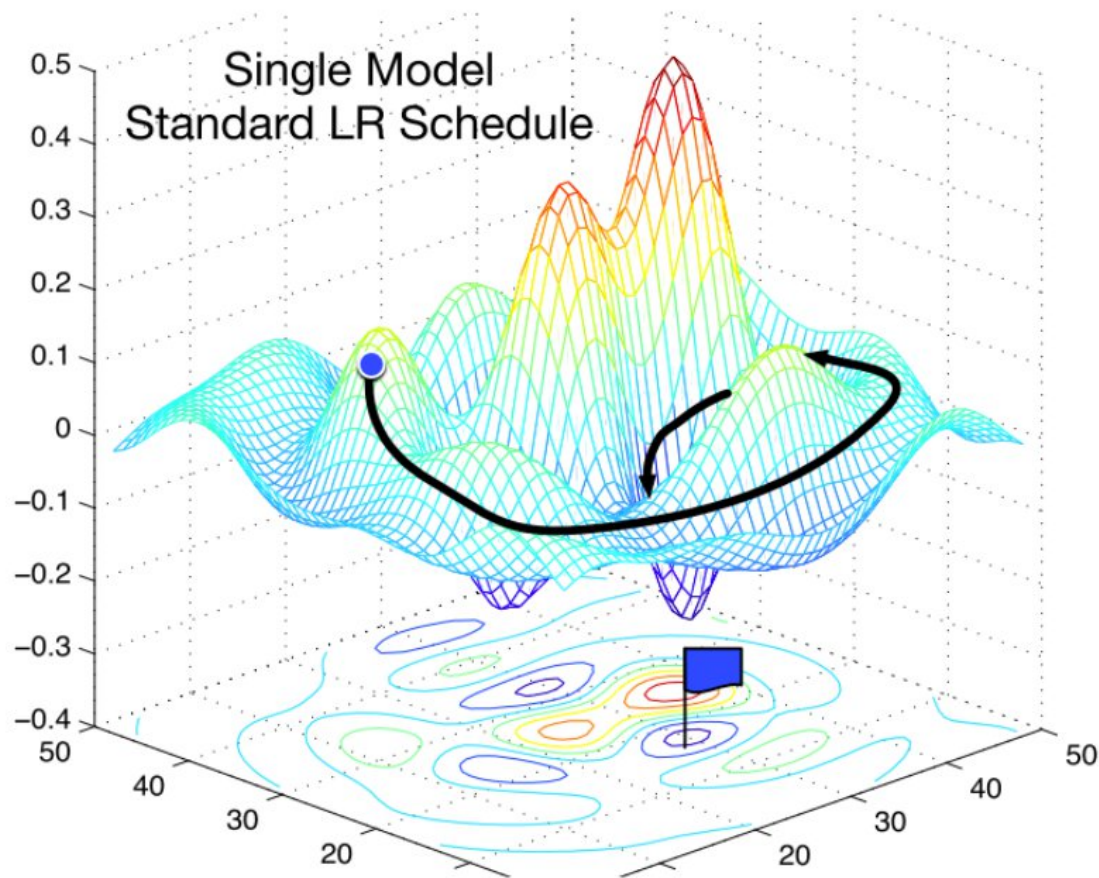


Paramètres :

- Step\_size =  $x * \text{epoch}$  ( $2 \leq x \leq 10$ )
- Base\_lr -> min convergence value
- max\_lr -> max convergence value

Succession of warmups and learning rate decays

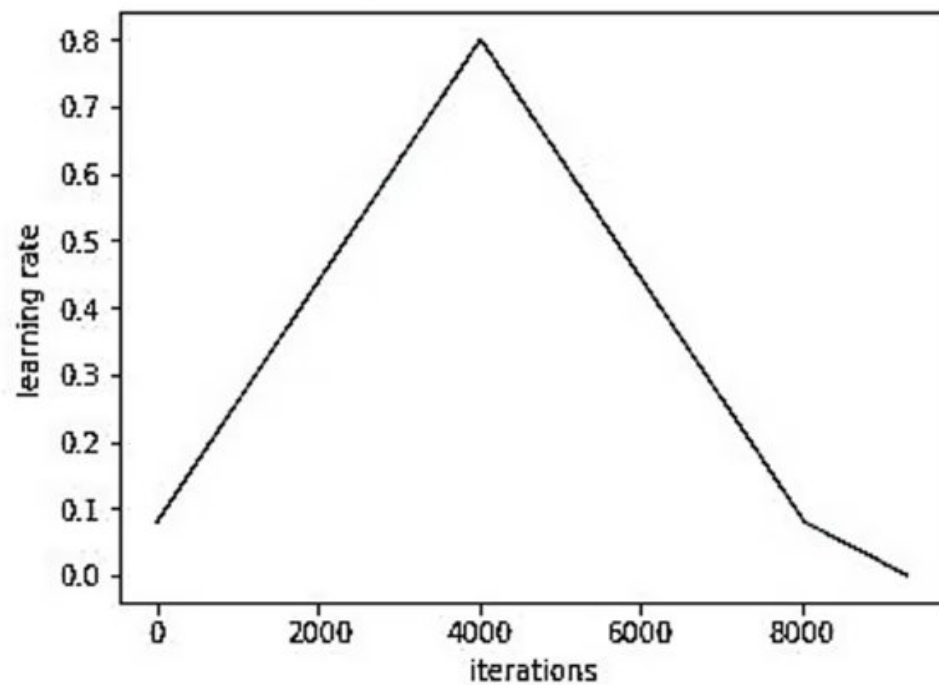
## Cyclic Learning Rate Scheduler



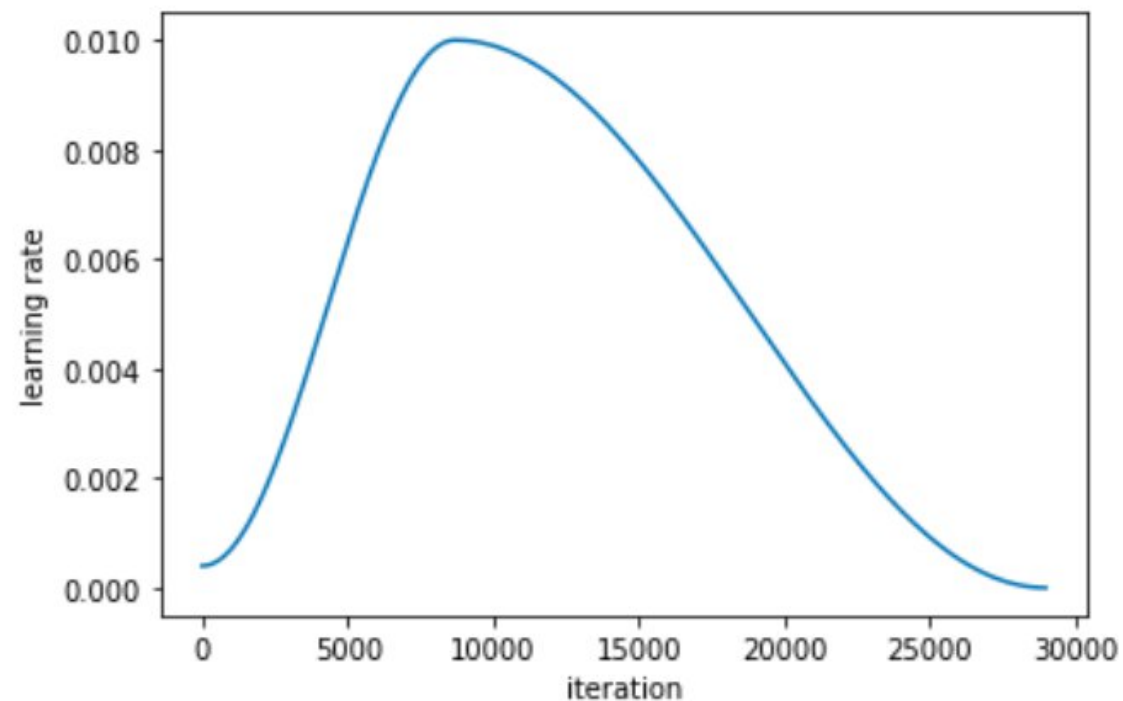
SNAPSHOT ENSEMBLES: TRAIN 1, GET M FOR FREE  
*Gao Huang, Yixuan Li, Geoff Pleiss*

A disciplined approach to neural network hyper-parameters - [Leslie N. Smith](#)

Proposition initiale

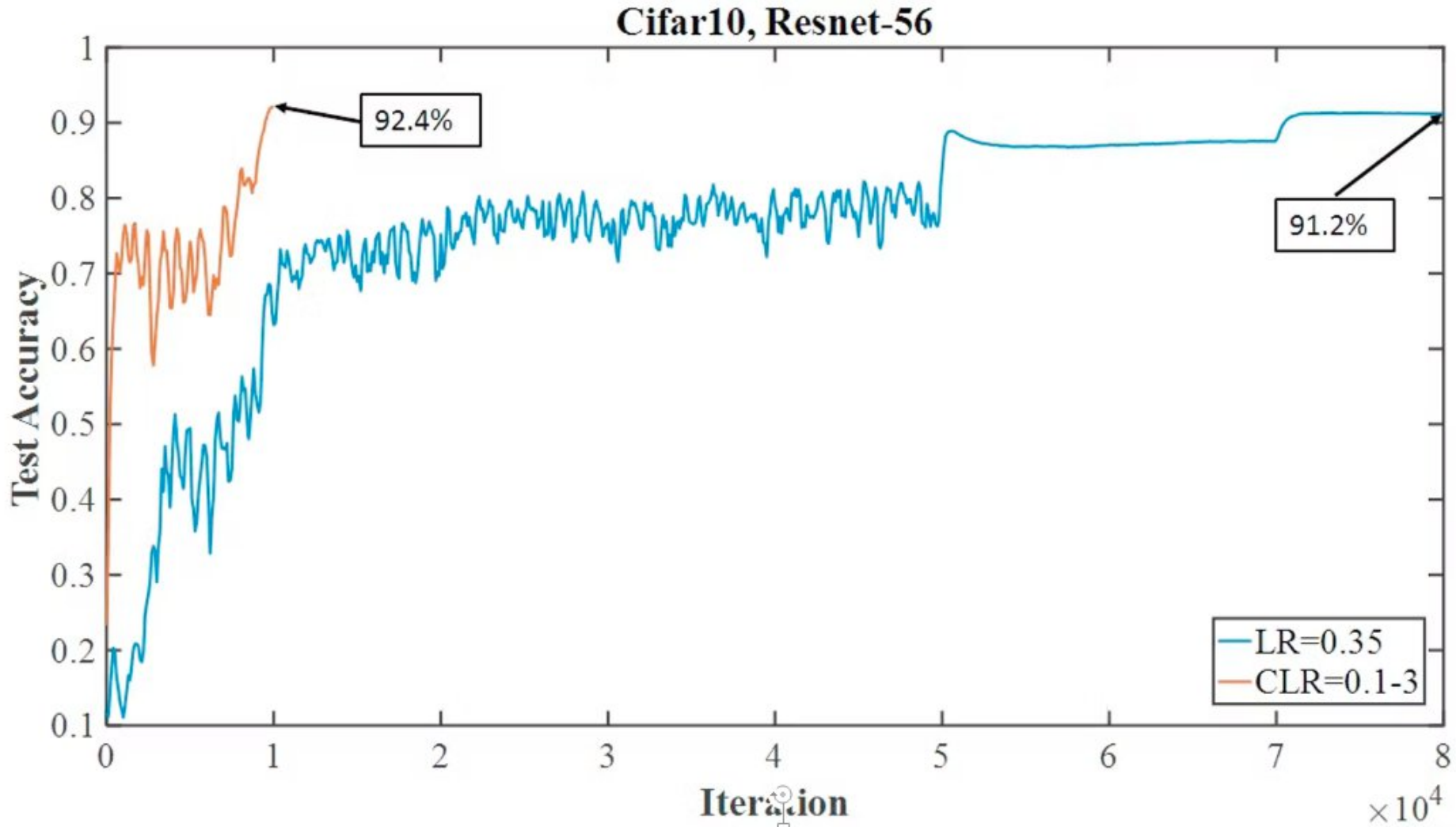


cosine annealing : Recommandation par FastAI



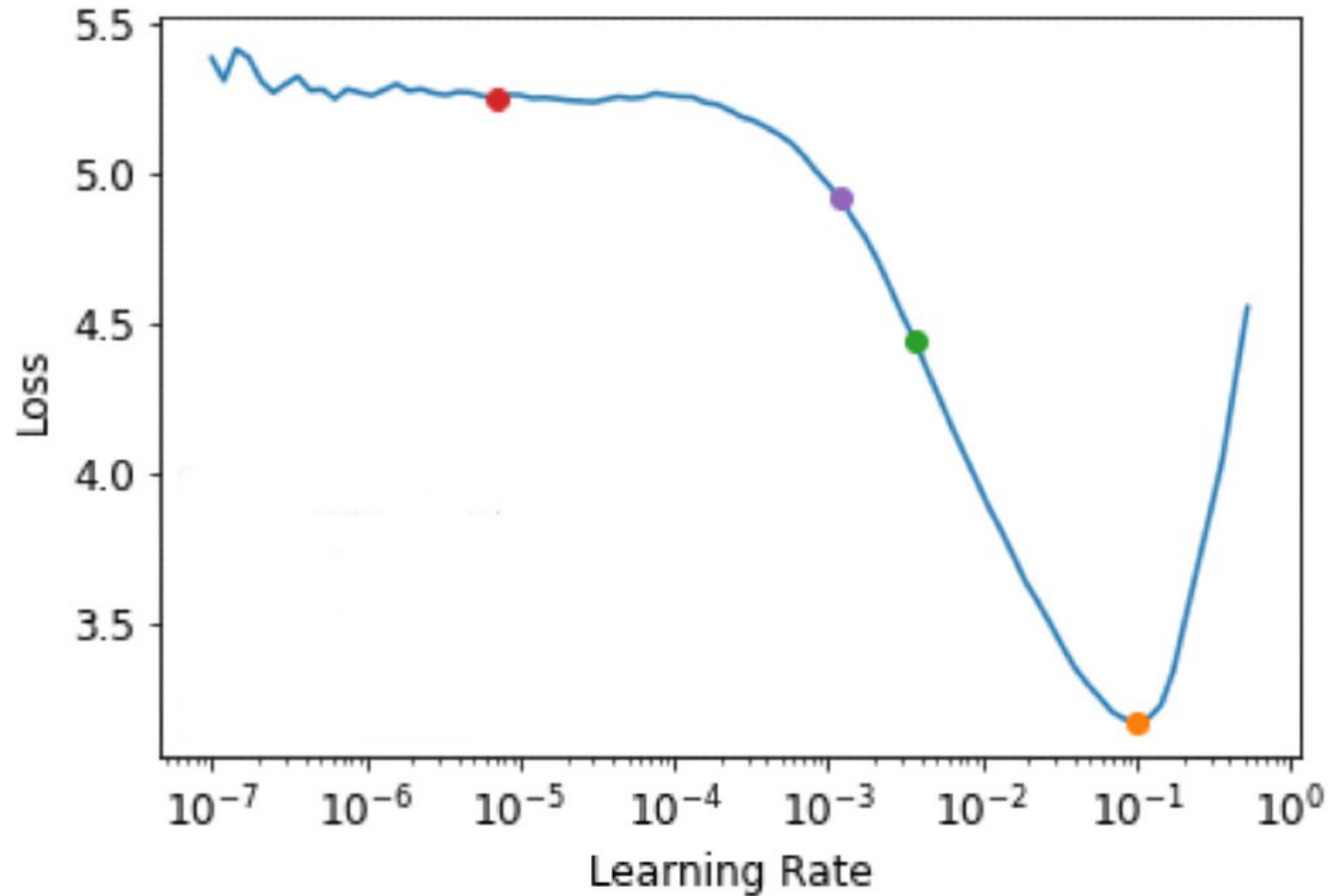
**One Cycle is enough**

Faster convergence for equivalent final precision



**Cyclic learning rate scheduler**

Goal: Find the **optimal learning rate** values for your model, particularly for **the maximum value** of a *cyclic scheduler*



Each **scheduler** has *its own settings*

```
import torch.optim as opt
```

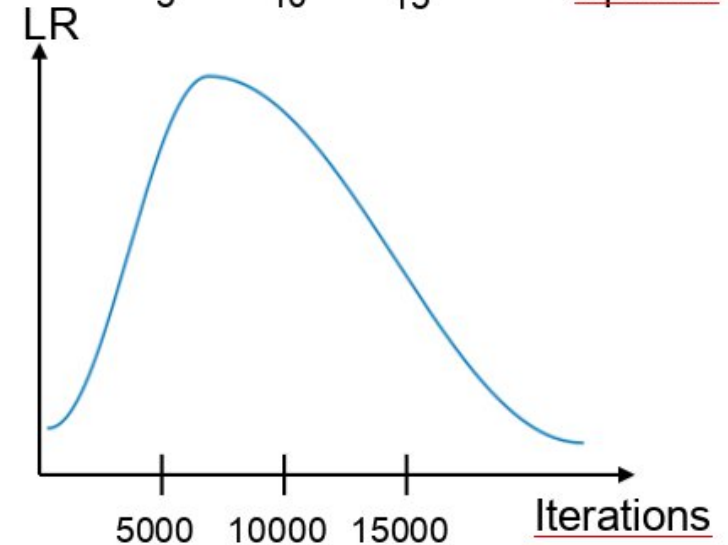
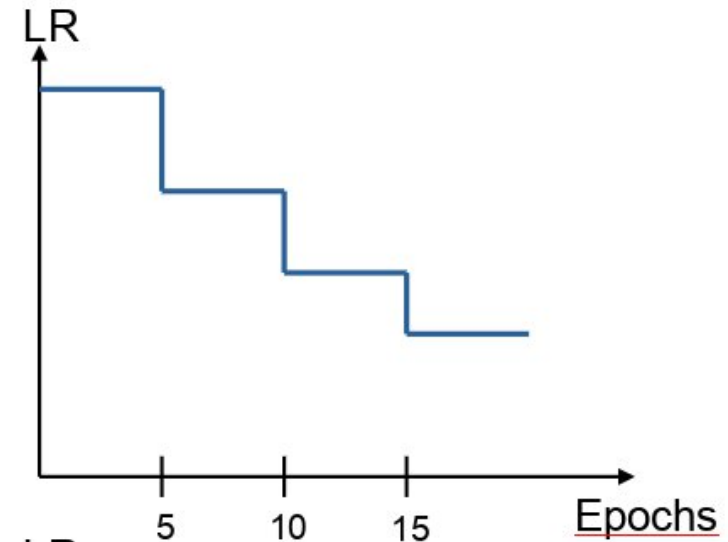
```
scheduler = opt.lr_scheduler.StepLR(optimizer, step_size=5, gamma=0.1)
```

```
for epoch in range(100):  
    train(...)  
    validate(...)  
    scheduler.step()
```

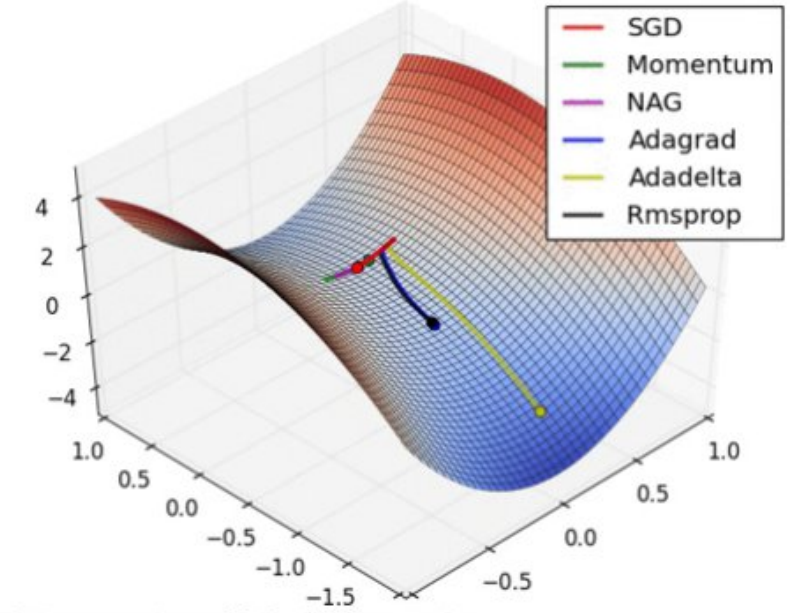
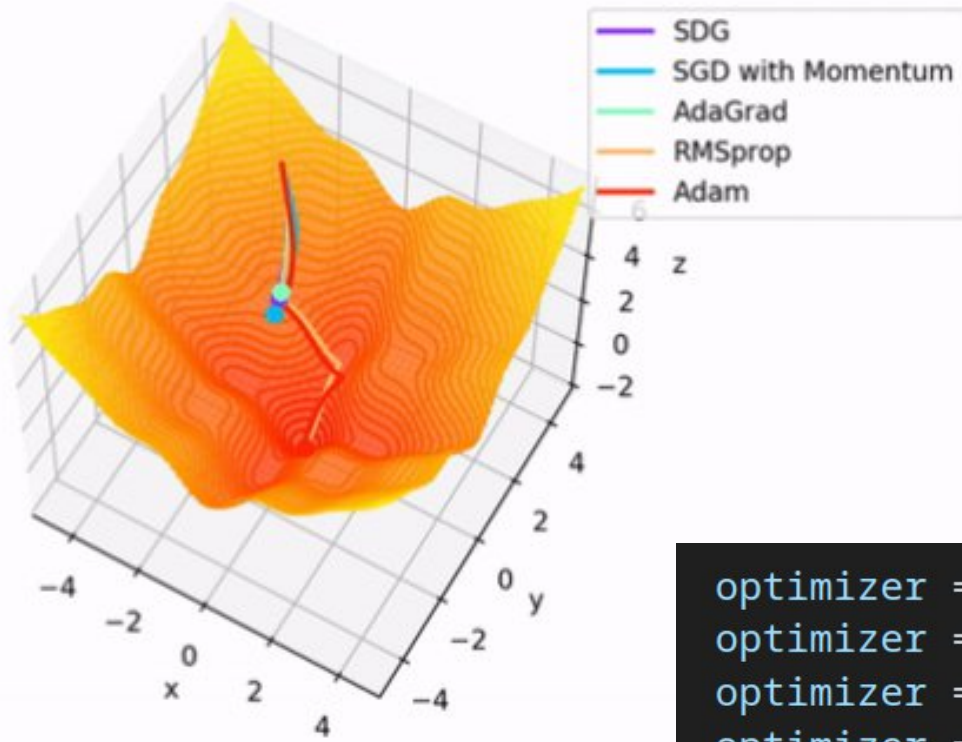
```
import torch.optim as opt
```

```
scheduler = opt.lr_scheduler.CyclicLR(optimizer, base_lr=0.01, max_lr=0.1)
```

```
for epoch in range(10):  
    for batch in data_loader:  
        train_batch(...)  
        scheduler.step()
```



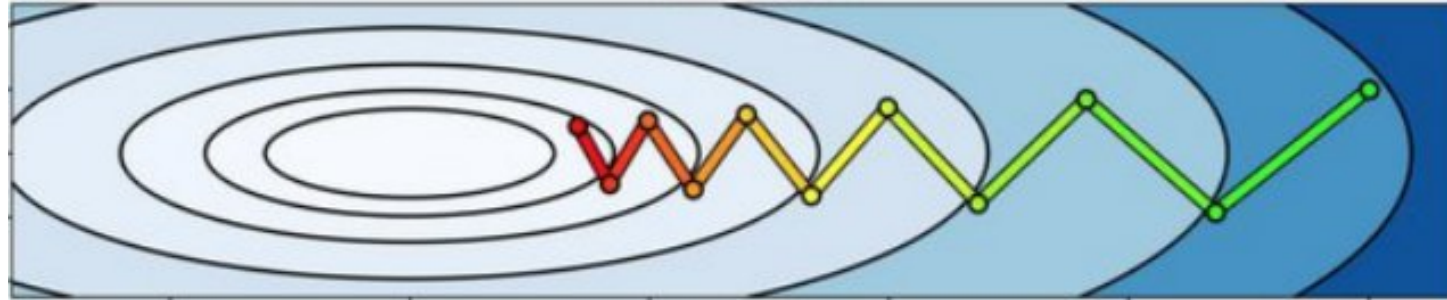
## Learning rate scheduler implementation



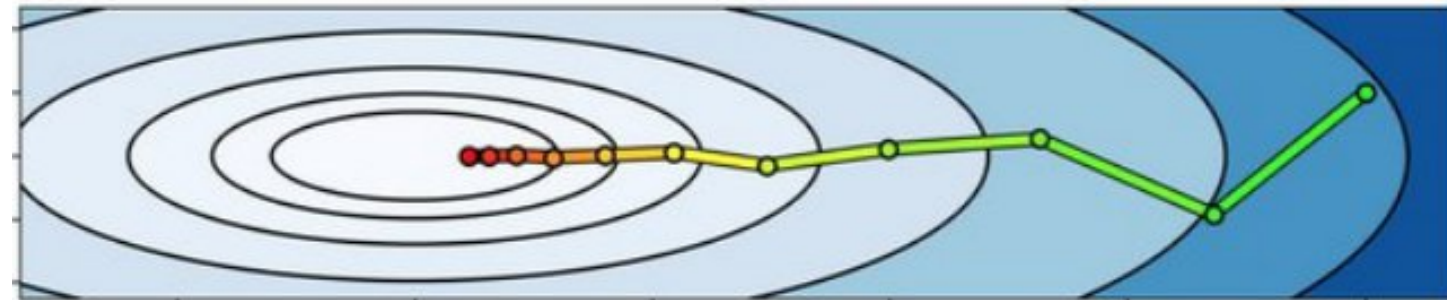
```
optimizer = torch.optim.Adadelta
optimizer = torch.optim.Adagrad
optimizer = torch.optim.Adam
optimizer = torch.optim.AdamW
optimizer = torch.optim.Adamax
optimizer = torch.optim.ASGD
optimizer = torch.optim.LBFGS
optimizer = torch.optim.NAdam
optimizer = torch.optim.RAdam
optimizer = torch.optim.RMSprop
optimizer = torch.optim.Rprop
optimizer = torch.optim.SGD
```



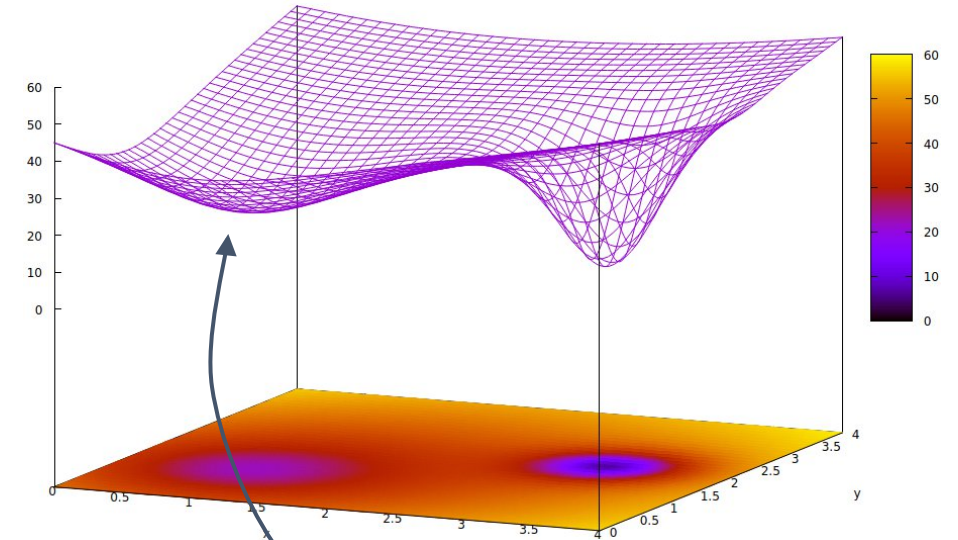
## Convergence



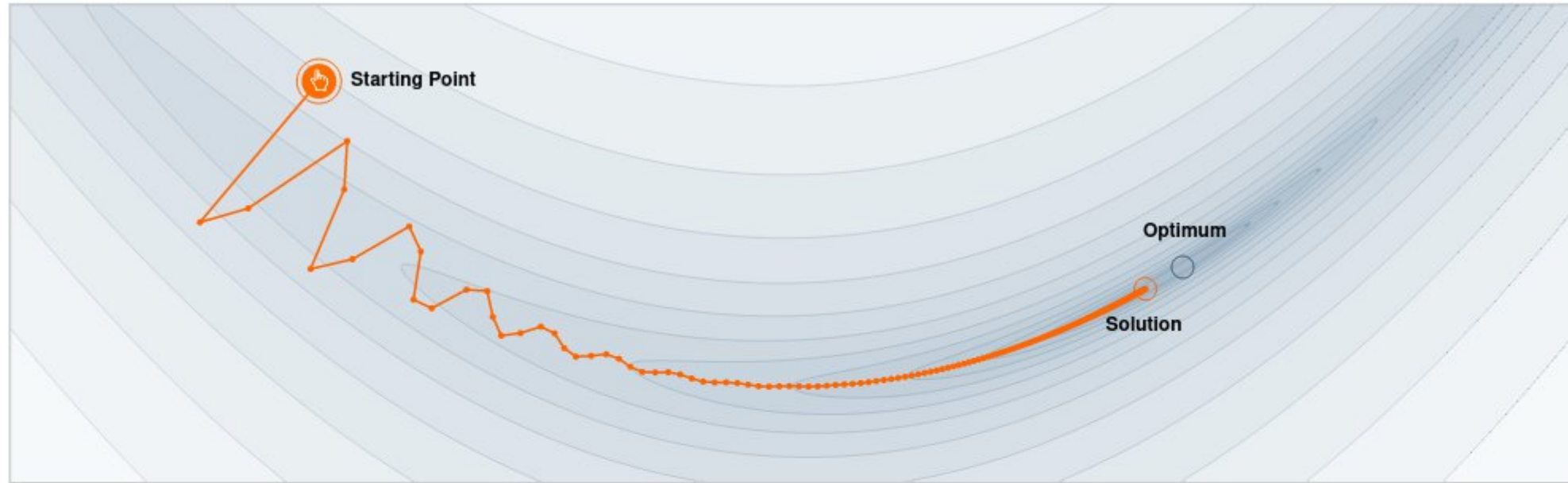
Without momentum



With momentum



Local minimum



Step-size  $\alpha = 0.02$



Momentum  $\beta = 0.99$



We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

GABRIEL GOH  
UC Davis

April. 4  
2017

Citation:  
Goh, 2017

## Why momentum works ?

$$\Theta_{t+1} = \Theta_t - \eta \nabla_{\Theta} [\mathcal{L}(\hat{y}_i, y_i) + \lambda R(\Theta_t)]$$

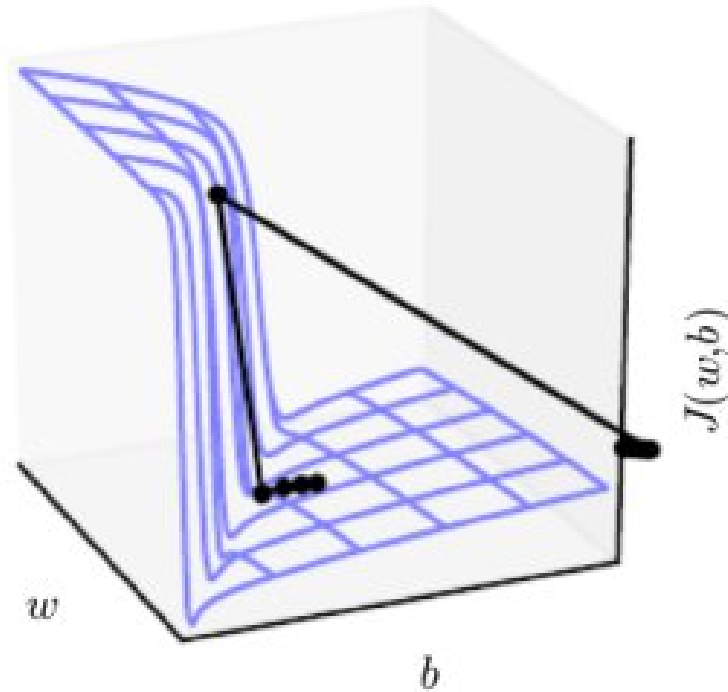
Updated weights = Weights before update - Learning rate \* Gradient [ Cost function ( Prediction, Label ) + Regularization rate \* Regularization function ( Weights before update ) ]

Weight update equation

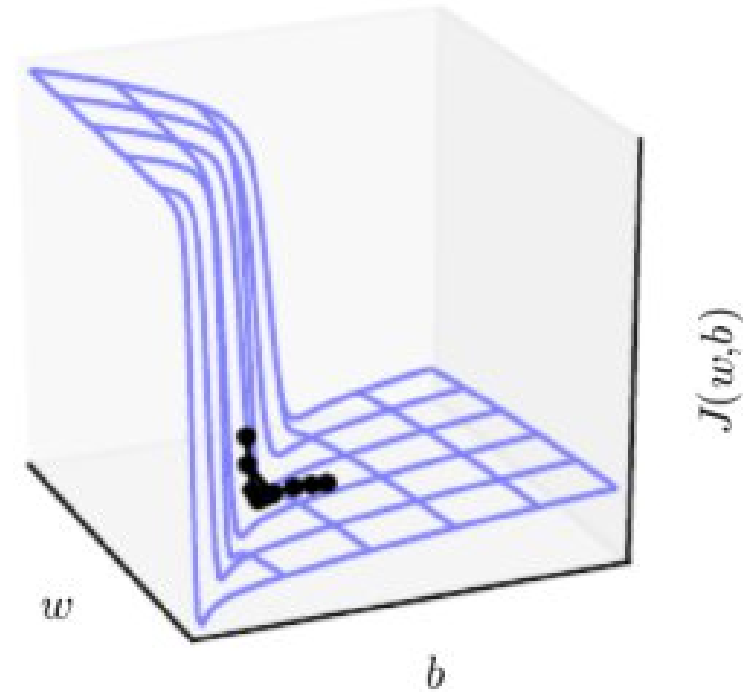
- L1 Regularization
- L2 Regularization
- Max norm Regularization
- Regularization with the cost function
- Dropout

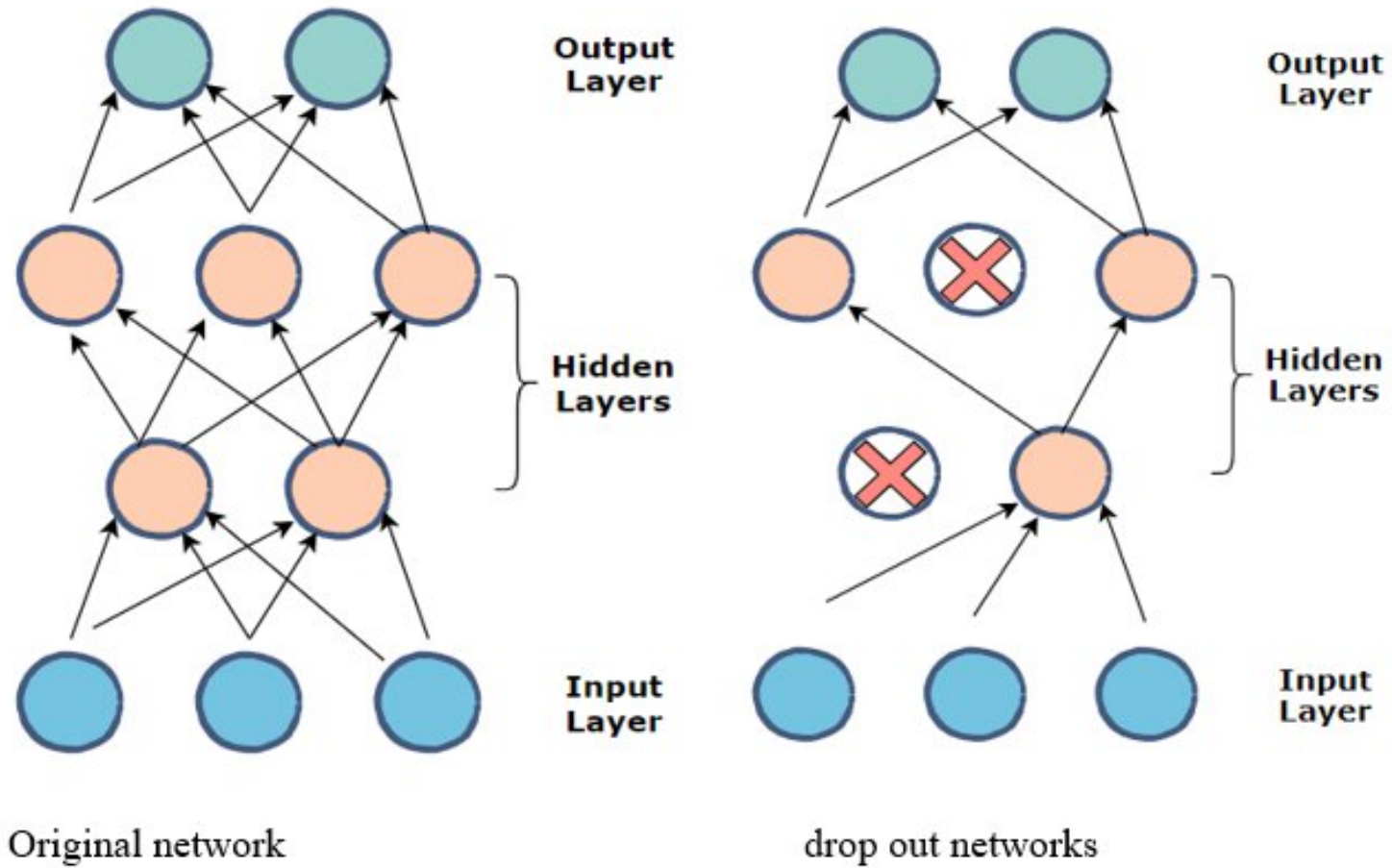
| L1 : LASSO | L2 : Ridge |
|------------|------------|
| $ \Theta $ | $\Theta^2$ |

Without clipping



With clipping

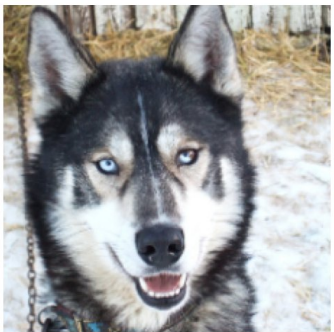




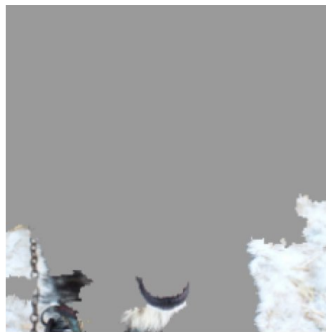
## Regularization : Dropout



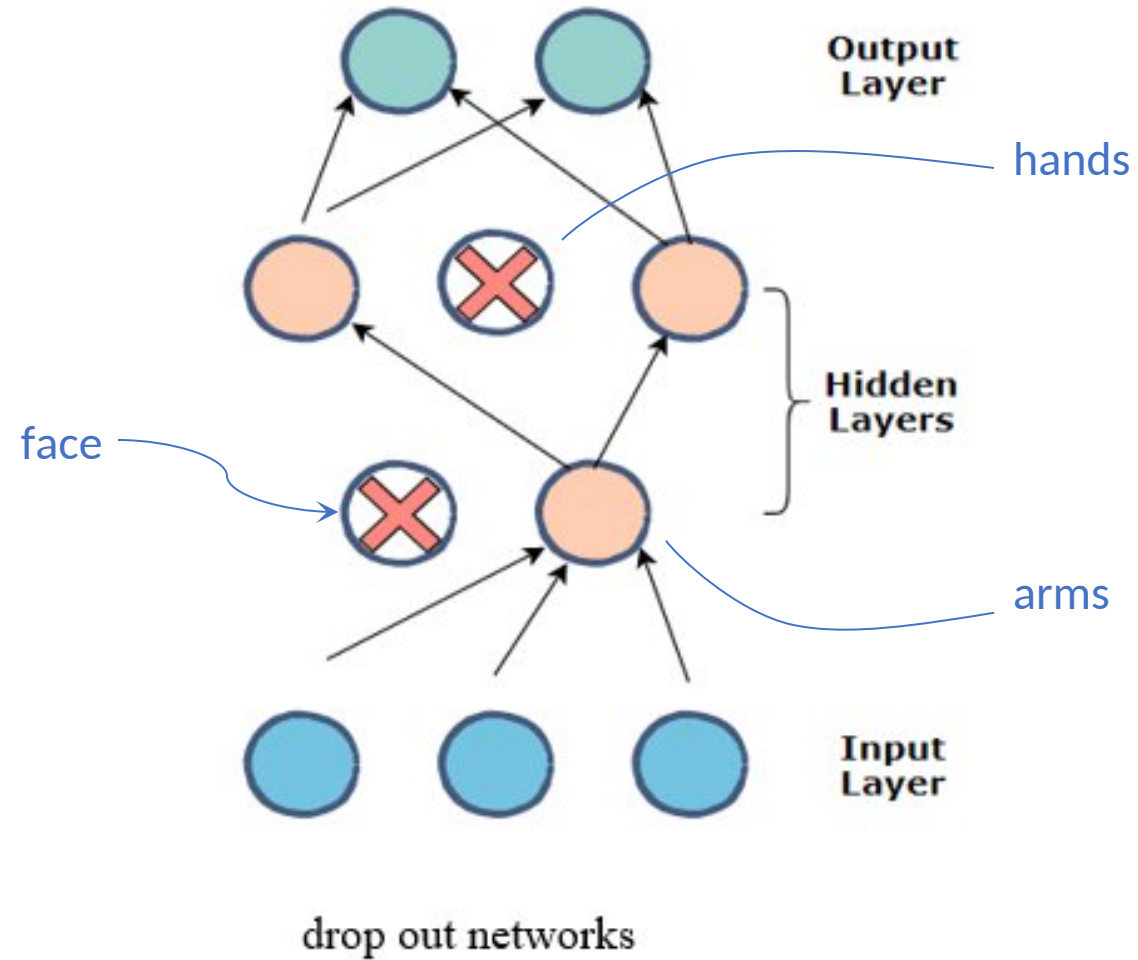
Human = face ? ❌



(a) Husky classified as wolf

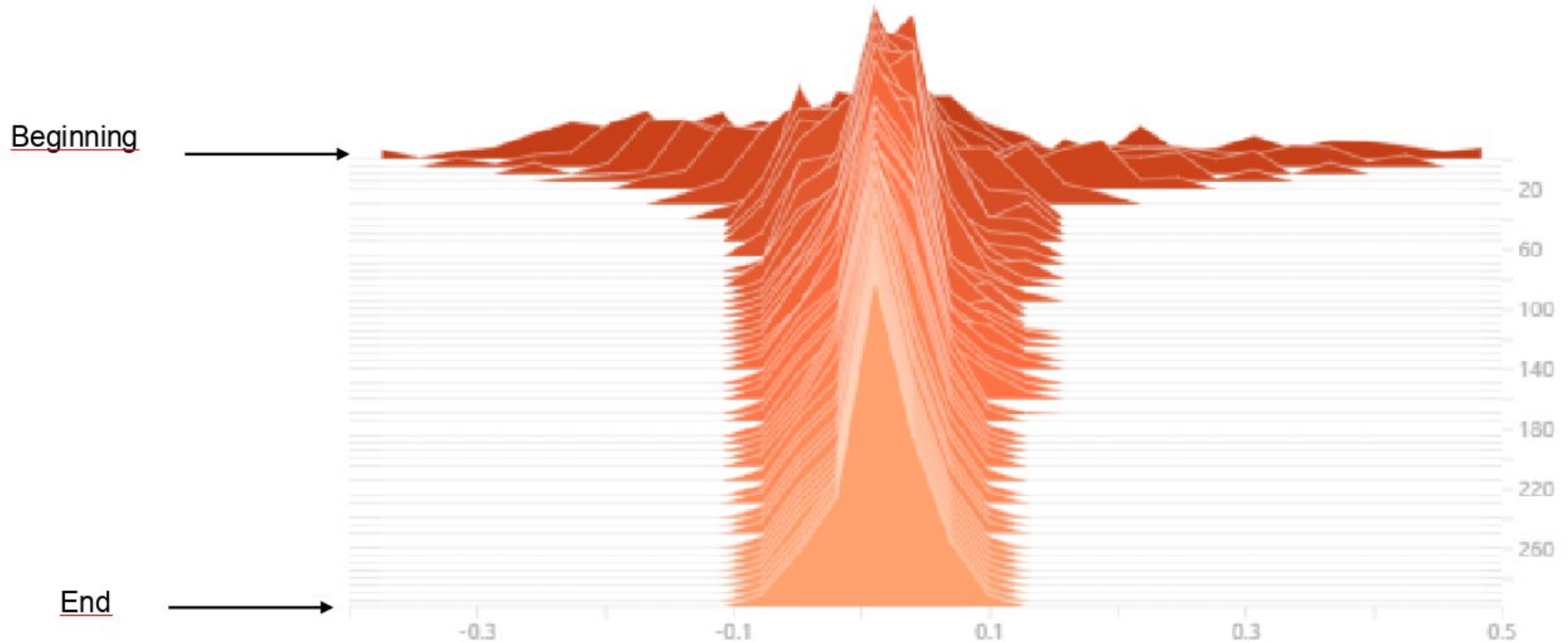


(b) Explanation



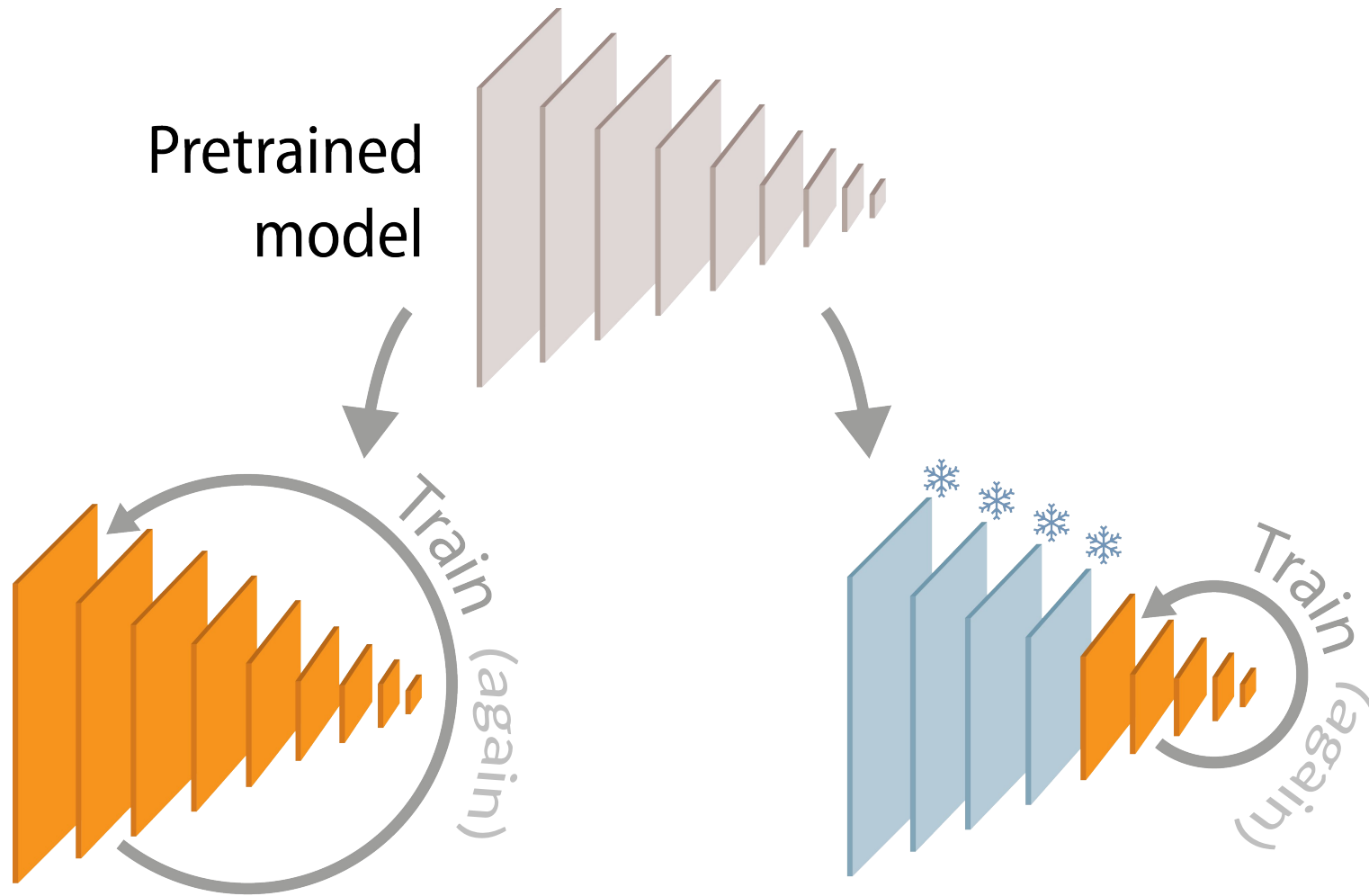
A neural network that **converges and generalizes correctly**\* generally has weights that **tend to 0**.  
\*(neither underfitting nor overfitting)

## Distribution of weights during learning:



- 1 Quel **travail** faire pour **améliorer** les **données** utilisées pour **l'entraînement** ?
- 2 Comment **évaluer** un **modèle** ?
- 3 Est-il possible de rendre **l'entraînement** plus **robuste** ?
- 4 **Peut-on profiter** d'un **modèle déjà entraîné** ?
- 5 Bonus : Quelques **bonnes pratiques** ?



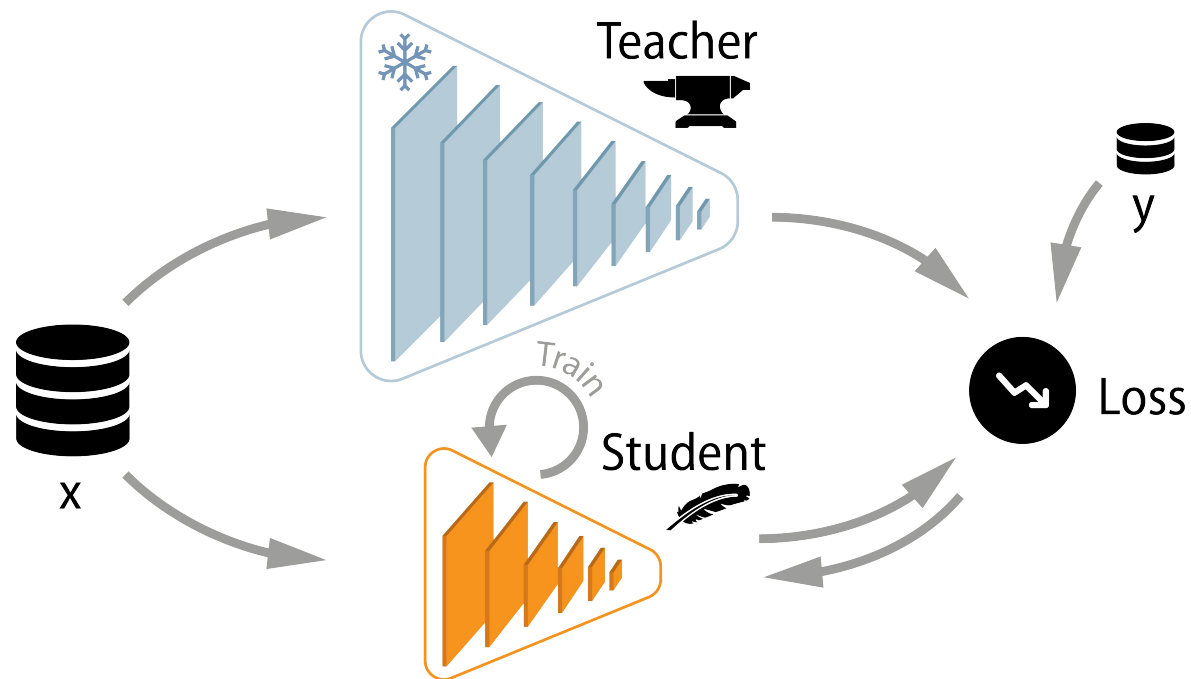


**Hugging Face**

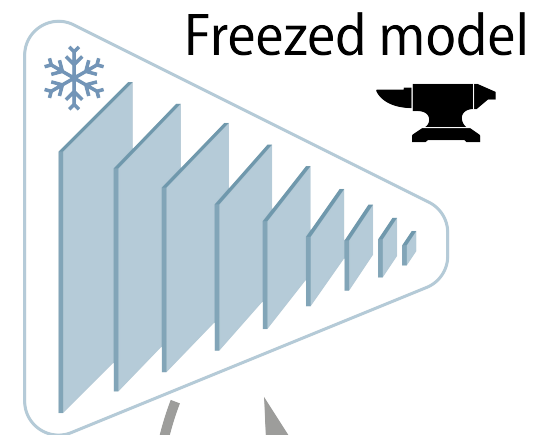
**Transformers**

**Timm**





**Knowledge distillation**



Trained adapters

**Adapters**

 **Hugging Face**

**PEFT**  
Parameter-Efficient  
Fine-Tuning

- 1 Quel **travail** faire pour **améliorer** les **données** utilisées pour **l'entraînement** ?
- 2 Comment **évaluer** un **modèle** ?
- 3 Est-il possible de rendre **l'entraînement** plus **robuste** ?
- 4 Peut-on **profiter** d'un modèle **déjà entraîné** ?
- 5 **Bonus : Quelques bonnes pratiques ?**

## Hyper-parameters

- Learning rate
- Regularization
- Optimizer
- Model architecture
- Batch size
- ...

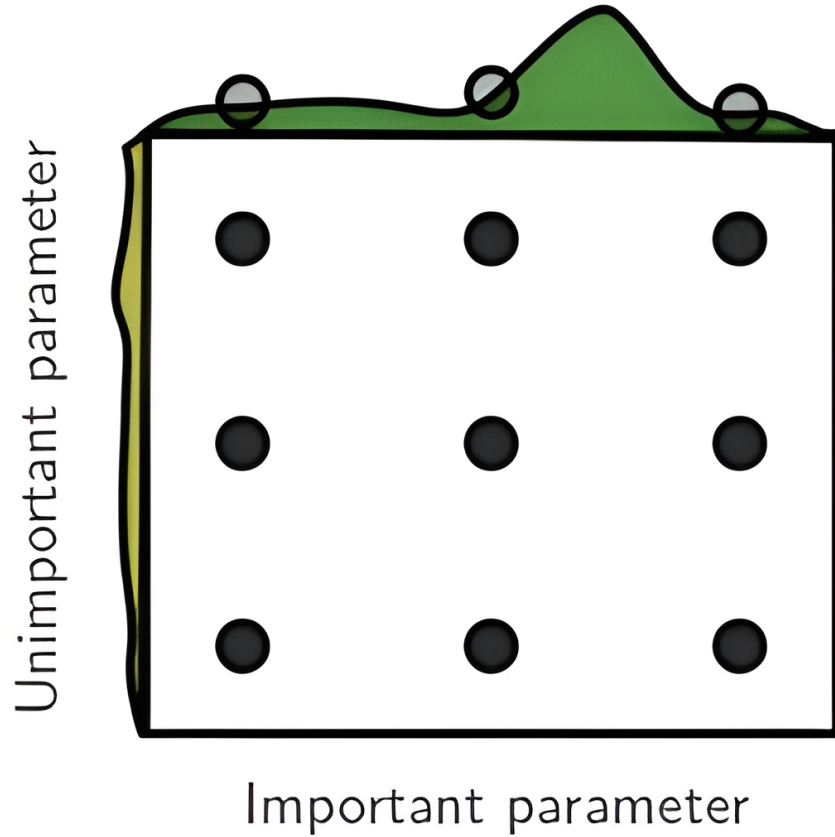
## Methods

- Manual research
- Grid search
- Random research
- Gradient
- Evolutionary algorithms
- ...

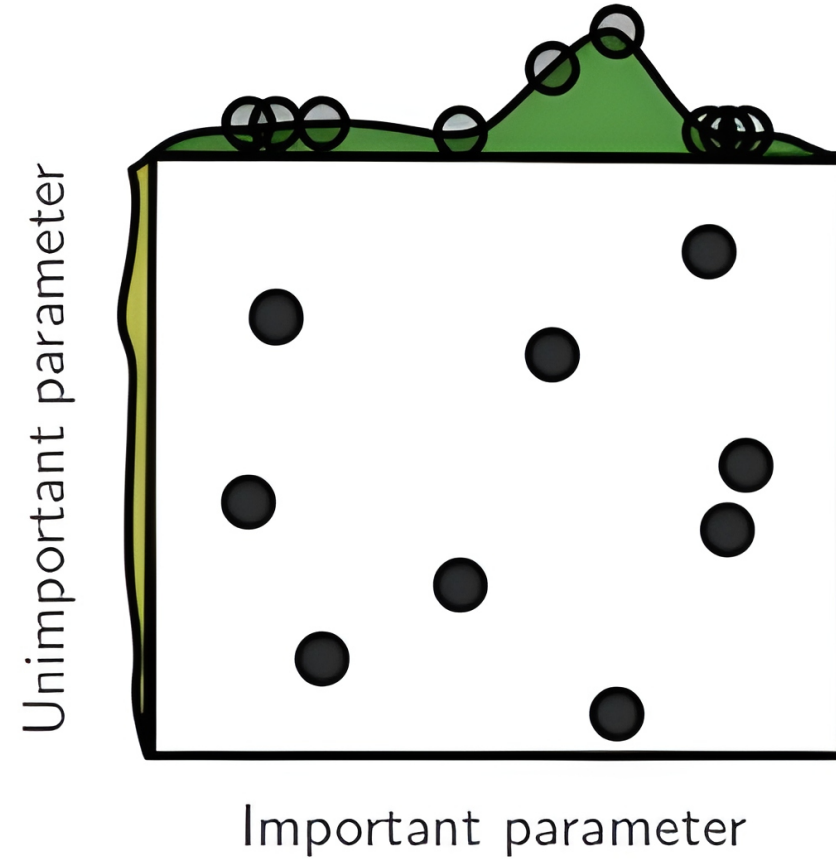
**HPO (Hyperparameter Optimization)**

**Find the good hyper-parameters**

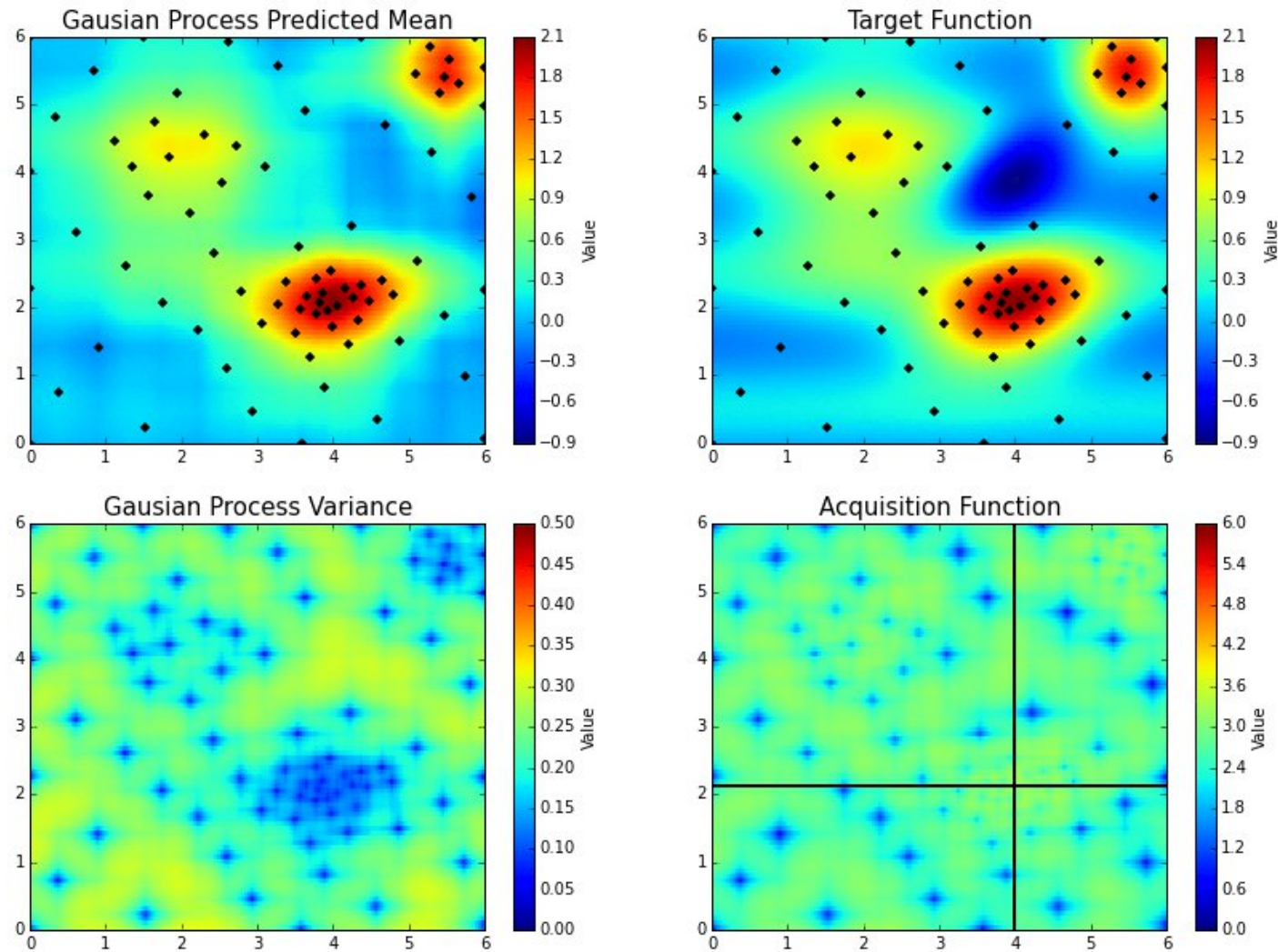
## Grid Layout



## Random Layout



# Bayesian Optimization in Action



Time created: All time
 State: Active

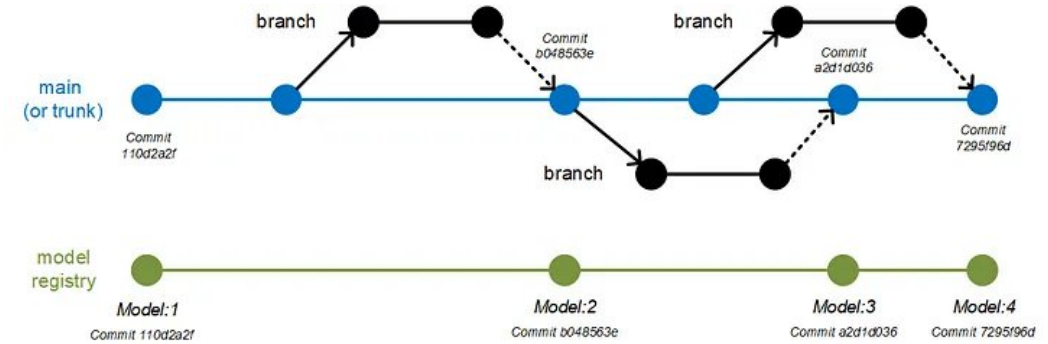
|                          |                          |                     |           | Metrics  |               | Parameters    |            |        |               |          |                |            |
|--------------------------|--------------------------|---------------------|-----------|----------|---------------|---------------|------------|--------|---------------|----------|----------------|------------|
| <input type="checkbox"/> | <input type="checkbox"/> | Run Name            | Created   | Duration | test_accuracy | training_loss | batch_size | epochs | learning_rate | momentum | weight_clampir | world_size |
| <input type="checkbox"/> | <input type="checkbox"/> | adorable-flea-594   | 1 day ago | 1.5min   | 58.73         | 0.935         | 32         | 10     | 0.01          | 0.9      | False          | 2          |
| <input type="checkbox"/> | <input type="checkbox"/> | righteous-calf-205  | 1 day ago | 1.5min   | 61.84         | 0.738         | 32         | 10     | 0.01          | 0.9      | True           | 2          |
| <input type="checkbox"/> | <input type="checkbox"/> | loud-dog-130        | 1 day ago | 2.8min   | 55.87         | 1.303         | 32         | 10     | 0.01          | 0.9      | True           | 1          |
| <input type="checkbox"/> | <input type="checkbox"/> | industrious-yak-917 | 1 day ago | 2.9min   | 57.57         | 0.915         | 32         | 10     | 0.01          | 0.9      | False          | 1          |



# Collaboration & Reproducibility

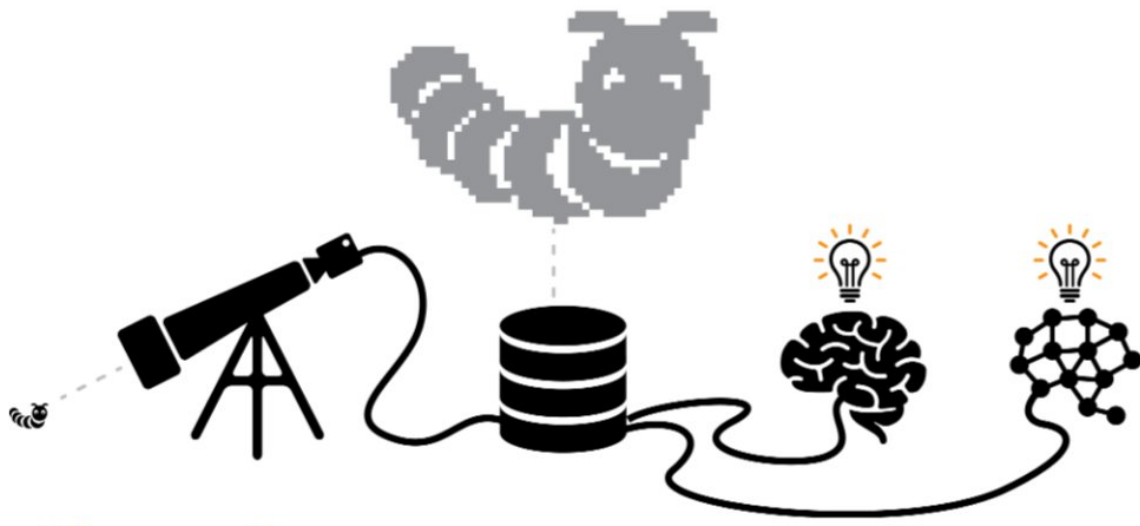
Gitlab Org > Issue Boards

The screenshot shows a GitLab Issue Board with three columns: Open, Deliverable, and Closed. The top navigation bar includes a dropdown for 'Development', a search filter, and buttons for 'Show labels', 'Group by', 'Edit board', and 'Create list'. The 'Open' column (56 issues) includes issues like 'Milestones swimlanes', 'Assign issue to epic', 'Create group', 'Remove issue from board', and 'Add lists for assignees and milestones'. The 'Deliverable' column (32 issues) includes 'Update issue due date from sidebar', 'Update issue's labels from sidebar', 'Update issue labels', 'Drag and drop issue between epics', and 'Paginate issues in Swimlanes'. The 'Closed' column (59 issues) includes 'Persist collapsed state of Swimlanes', 'Remove list from board', 'Remove issue from Swimlane', 'Expand diff to entire file', and 'Laboriosam commodi ab in eum qui suscipit necessitatibus modi fuga'.

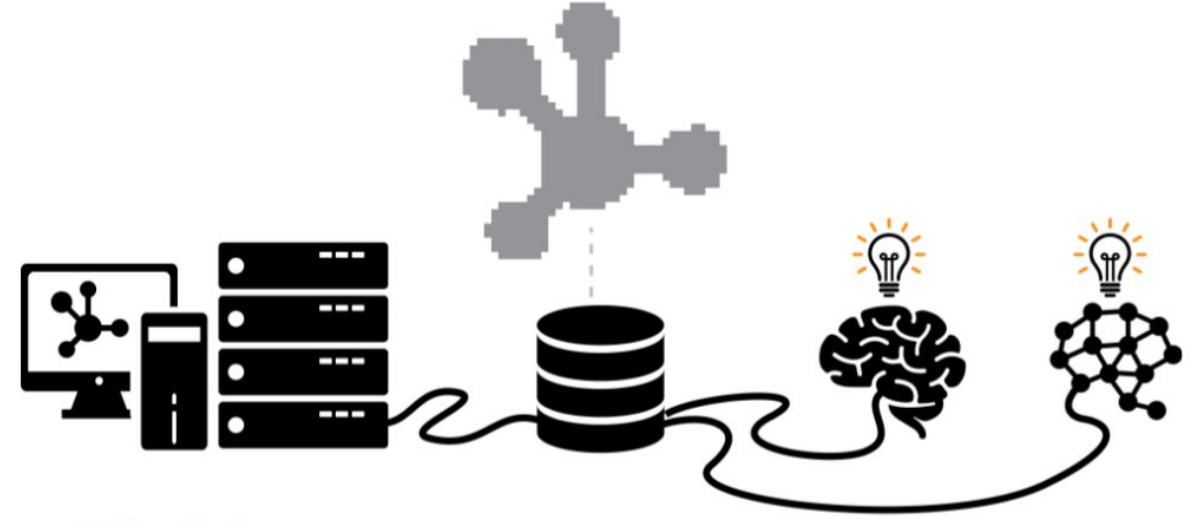


# A structured project





Observations



Générées

- |                 |             |          |                |                |            |
|-----------------|-------------|----------|----------------|----------------|------------|
|                 |             |          |                |                |            |
| Characteristics | Composition | Text     | Images         | Audio          | Trajectory |
|                 |             |          |                |                |            |
| Atom            | Molecule    | Material | Transformation | Social network | ???        |

1

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix}$$
$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$
$$\begin{pmatrix} \begin{pmatrix} 1 & 2 \end{pmatrix} & \begin{pmatrix} 3 & 4 \end{pmatrix} \\ \begin{pmatrix} 5 & 6 \end{pmatrix} & \begin{pmatrix} 7 & 8 \end{pmatrix} \end{pmatrix}$$

Scalar

Vector

Matrix

Tensor

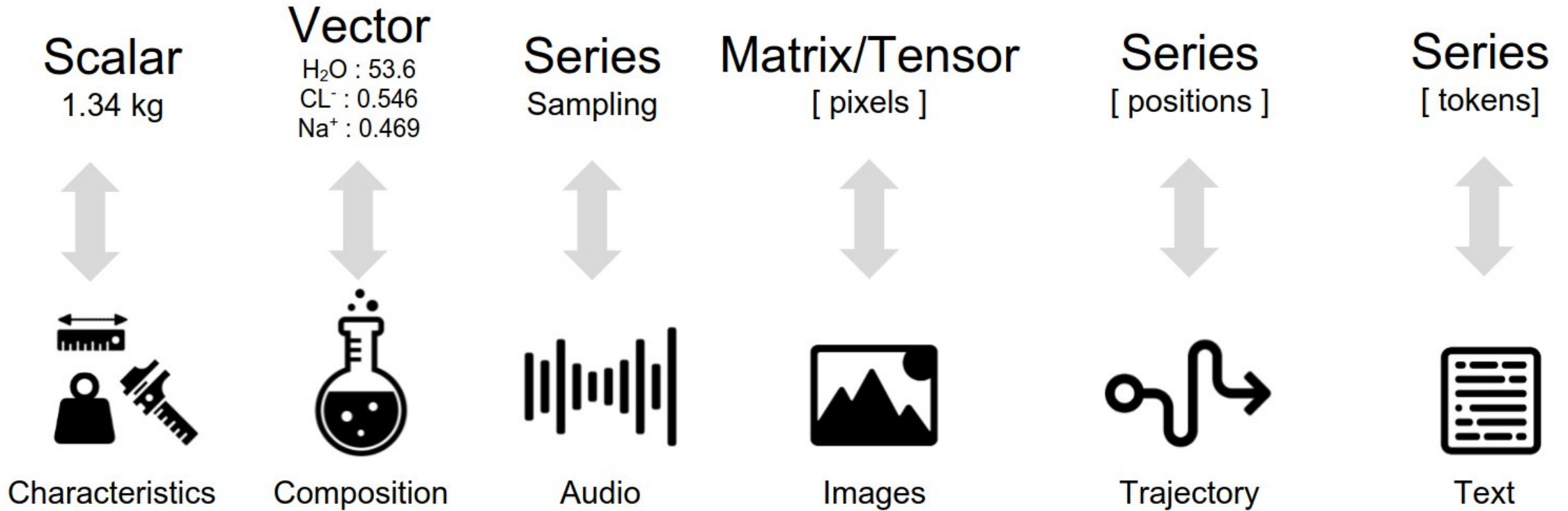
(Series of scalars)

(table of scalars)  
(series of series)

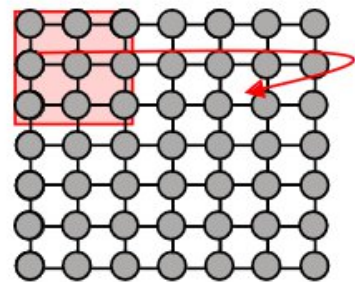
(series of series of series of...)



Some descriptors are relatively simple and intuitive

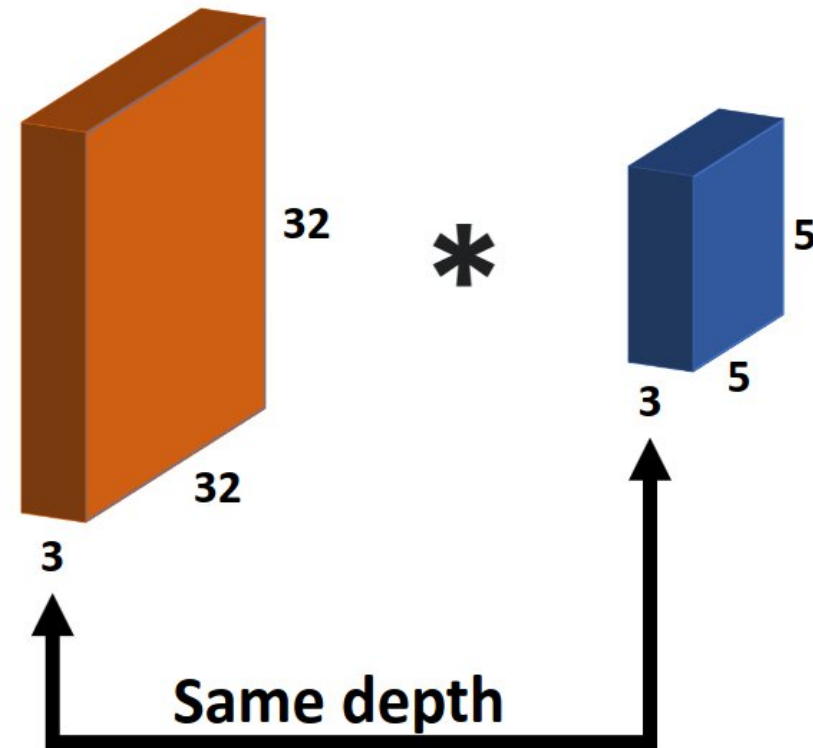
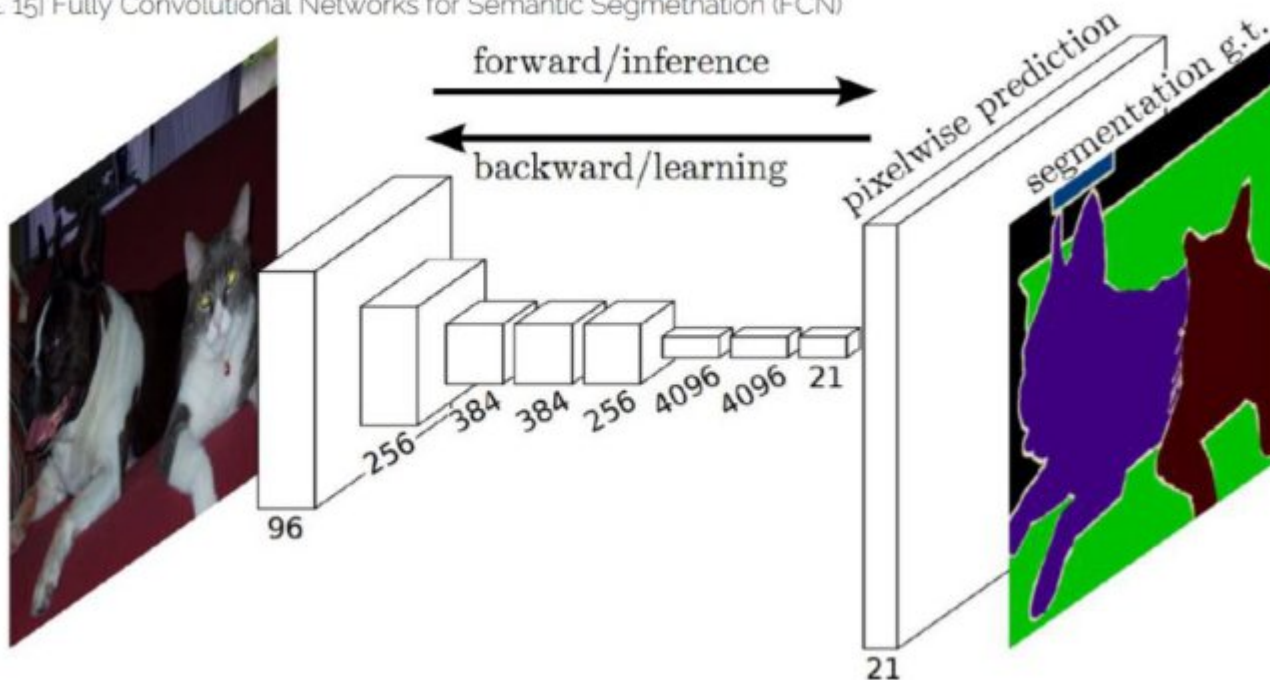


**Data ?**



Images 2D / 3D

[Long et al. 15] Fully Convolutional Networks for Semantic Segmentation (FCN)



# Convolutional Neural Network

Reverse diffusion



$x_0$



Markov  
Chain

$x_1$



$x_t$



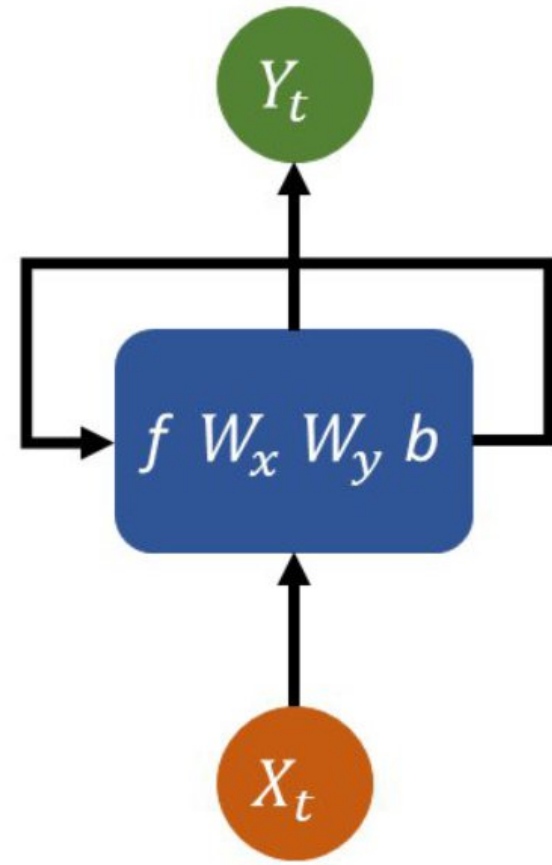
$x_T$



Forward diffusion



|         | day 1 | day 2 | day 3 |
|---------|-------|-------|-------|
| asset 1 | 9.77  | 79.94 | 64.13 |
| asset 2 | 47.66 | 74.07 | 70.90 |
| asset 3 | 94.25 | 76.34 | 99.95 |
| asset 4 | 41.19 | 9.99  | 89.50 |
| asset 5 | 65.44 | 63.79 | 67.14 |



$$Y_t = f(W_x \cdot X_t + W_y Y_{t-1} + b)$$

The <sup>Focus</sup> → The big red dog  
 big → The big red dog  
 red → The big red dog  
 dog → The big red dog

