# ① Techniques de regularisation d'un algo d'apprentissage d'1 NN

# ② Techniques d'optimisation

↳ 'vanilla' Stochastic Gradient descent (SGD)

↳ algorithmes plus performants.

# (1) Techniques de regularisa°

<span style="color:blue">**Regularisation**</span> : Any modification to the learning algo that its intended to reduce <u>the error</u> but its not <u>training</u> generalization error.

↳ Cost fonction that we minimiz during training

↳ ex : regression case

$$NSE = \frac{1}{N} \sum_{i=1}^{N} \left( f_\theta(x_i) - y_i \right)^2$$

[Test Set] ⟶ erreur de generaliza°
NSE sur le dataset de test.

Overfitting = "Sur apprentissage"

↳ train error small
↳ generaliza° error bigger.

↳ trop de paramètres dans le modèle

↳ dataset d'entrainement ⊃ représente dataset de test toute votre donnée

↳ hyperparamètres de votre algorithme d'entrainement.

↳ EPOCH → 1 mise à jour des paramètres sur tout le dataset d'entrainement.

↳ Finding the trade-off for the errors between:

→ high bias ⇒ underfitting

↘ high variance
　　　　↳ overfitting.

↳ generaliza° error → the smallest.

↳ w/ a small gap between the training error and the generalization error.

Deep Neural Networks → Billions of Parameters! → High Model Capacity

Keep large models but regularized them to avoid overfitting.

Cost function $J(\theta, X, y)$
= loss function

$\llcorner \tilde{J}(\theta, X, y) = J(\theta, X, y) + \boxed{\alpha} \Omega(\theta)$

terme de régularisation
de la loss.  $\alpha > 0$.

$\Omega(\theta)$ penalized the $\boxed{\text{weights}}$

$L^2$ régularisation
$\Omega(\theta) = \frac{1}{2} ||\omega||^2 \longrightarrow$ ridge regression.

# $\ell^1$ parameter regularization

$$\Omega(\theta) = \|w\|_1 = \sum_i |w_i|$$

## norm $L^2$ :

$$\tilde{J}(\theta, X, y) = \frac{\alpha}{2} w^T w$$
$$+ J(\theta, X, y)$$

$\downarrow$

$$\nabla_w \tilde{J}(\theta, X, y) = \alpha w + \nabla_w J(w, X, y)$$

$\alpha$ : norm des weights $w$

$\varepsilon$ : learning rate.

$$W_{k+1} \leftarrow w_k - \varepsilon (\alpha w_k + \nabla_{w_k} J(w_k, X, y))$$

$$w_{k+1} \leftarrow (1 - \varepsilon \alpha) w_k - \varepsilon \nabla_{w_k} J(w_k, X, y)$$

additional term

shrink the weight vector at each
timestep.

$1^\wedge$ regularisation $\rightarrow$ Solution au niveau de la contrainte sur les poids qui est plus 'sparse'

$\hookrightarrow$ w tends to be more equal to 0.

$\downarrow$ mécanisme de selection de variables.

Technique computationally inexpensive but powerful to regularize a large family of models
$\hookrightarrow$ Any kind of Neural Network!

technique ⟷ 'Bagging' Methods
for ensembles of very large
Neural Networks.

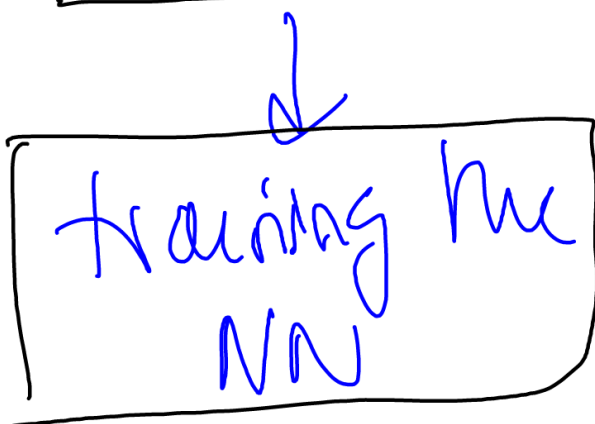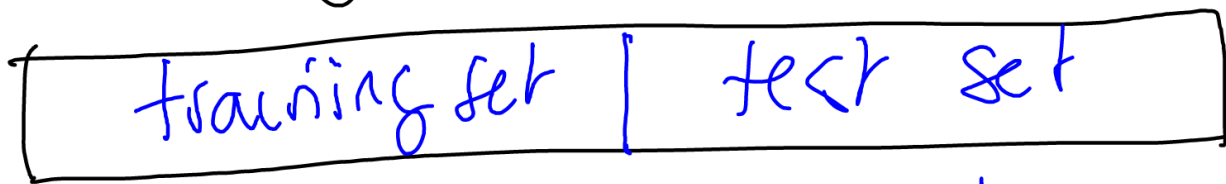[Dropout]

↳ trains sub-Networks
that can be formed by
removing hidden units from
the underlying base Network.

@ Dropout algo

Each time we update the Network's
parameters, we randomly sample
a binary mask to all hidden
units of the
Neural Network
↓
selects randomly a percentage of
all hidden
neurons

# dataset

| training set | test set |
|---|---|

↓        ↓

training the NN     hold-out set to evaluate the capacity of the NN to generalize to **unseen data**

↓

ⓧ parameters $\theta$.

ⓧ hyperparamètrs de votre **NN**
   - ↳ modèle → nbre de couches/nbre de neurones
   - ↳ algorithme d'entrainement
     - @ learning rate.
   - ↳ Multiple training by varying its hyper-parameters.

⟹ DATA LEAKAGE

training set ↓   validation set ↓   test set

↳ utilisé
dans f° coût $J(\theta)$
to find $\theta^{*}$
↳ modèle (NN) spécifique
↳ série d' hyperparamètres
            spécifiques .

evaluation
of ≠
hyperparamètre

Compute généralization error
on the test set
for ( Best hyperparameters
        NN ($\theta^{*}$) )

# II/ Optimization strategies for Deep Neural Networks

## Empirical-risk Minimization

Goal of an algorithm of ML:
reducing the error generalization:

$$J^*(\theta) = E_{(x,y) \sim \boxed{p_{data}}} \mathcal{L}(f(x,\theta),y)$$

"vraie" distribution de la donnée

$\hat{p}_{data}$ : training dataset
$\hookrightarrow$ distribution $\hat{p}_{data}$ pour le training dataset

$$J(\theta) = E_{x,y \sim \hat{p}_{data}} \mathcal{L}(f(x,\theta),y)$$

$\hookrightarrow \mathcal{D} = ((x_i,y_i) \; i \in \{1 \ldots N\})$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(f(x_i, \theta), y_i)$$

↳ In pratice, we minimize the negative log-likelihood:

$$J(\theta) = \mathbb{E}_{x,y \sim \hat{p}_{data}} \log P_{model}(x, y, \theta)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \log P_{model}(x_i, y_i, \theta)$$

→ Mini-batch gradient descent.

Batch $B \subset$ $\underline{\text{Dentraenen}}$

↳ Most optimization algorithms converges faster if they are allowed to rapidly compute approximate

estimation of $\nabla J(\theta)$ ($\boxed{\nabla J(\theta)}$)
than slowly computing exactly
$J(\theta)$ $(\boxed{\nabla J(\theta)})$

training algorithms based on
gradient-based learning for
machine learning problem

| Batch gradient descent | Mini-batch gradient descent | Stochastic online gradient descent |
|---|---|---|
| Compute $\nabla J(\theta)$ on the whole training dataset & then we update $\theta$. | Compute $\nabla J(\theta)$ on the mini-batch & then we update $\theta$ | Compute $\nabla J(\theta)$ per training tuple & then we update $\theta$ |

Most efficient method

(beyond SGD)

SGD Mise à jour de $\theta^k$
de cette manière :

$$\theta_k \leftarrow \theta_{k-1} - \alpha \nabla_{\theta_{k-1}} J(\theta_{k-1})$$

hyperparamètre $\alpha$. } learning rate
(taux d'apprentissage)

↑ big sensivity of the performance
of the learning algorithm to $\alpha$.

ⓐ SGD with decreasing $\alpha_k$.

$$\alpha_k = (1-\epsilon) \alpha_0 + \epsilon \alpha_{\overline{\tau}}$$

$$\sum_{h=1}^{\infty} \alpha_k = \infty \quad \text{serie divergente}$$

$$\sum_{h=1}^{\infty} \alpha_h^2 < \infty \quad \text{serie convergente}$$

$\hookrightarrow$ In theory, convergence guarantees.

**Momentum**

Goal $\rightarrow$ Accelerate learning

$\hookrightarrow$ Faster Convergence

$\rightarrow$ Small but consistent gradients

$\hookrightarrow$ noisy gradients

$\sigma = $ velocity $\rightarrow$ Control direction & speed at which parameters move through the parameter space

Iteration $k$ of the training algo:

$$v_k \leftarrow \boxed{\alpha} v_{k-1} - \boxed{\varepsilon} \nabla_{\theta_{k-1}} J(\theta_{k-1})$$

$$\theta_k \leftarrow \theta_{k-1} + v$$

New hyperparam $\varepsilon$.

## II C/ Algorithms with adaptative learning rates

learning rate $\alpha$ : 1 des hyperparamètres le $\oplus$ durs à tuner.

$\hookrightarrow$ momentum add another hyperparameter.

learning rate $\alpha$ vary depending the $\neq$ parameters. $\theta \in \mathbb{R}^D \longrightarrow$ Billions of parameters.

**Adagrad** optimisation algo w/ adaptative learning rate.

↳ ⊗ parameters w/ the larger partial derivative of the loss

$$\nabla J(\theta) = \begin{bmatrix} \dfrac{\partial J(\theta)}{\partial w_i} \\[2mm] \dfrac{\partial J(\theta)}{\partial b_i} \end{bmatrix}$$

→ rapid decrease in the learning rate.

↳ ⊗ parameters w/ the smallest partial derivatives : small ↓ of the learning rate.

→ RMS prop

↳ Adam

✗~