

parameters VS Hyperparameters

parameters:

$w^{[3]}, b^{[3]}, \dots, w^{[2]}, b^{[2]}$

Hyperparameters

Optimization algorithm:

Learning rate α / Learning rate decay

iterations

mini batch size

momentum

regularization parameters λ , β , ρ (dropout)

networks:

hidden layers L

hidden units in each layer

choice of activation functions

the variance of $w^{[2]}$ when initializing them

Regularization techniques

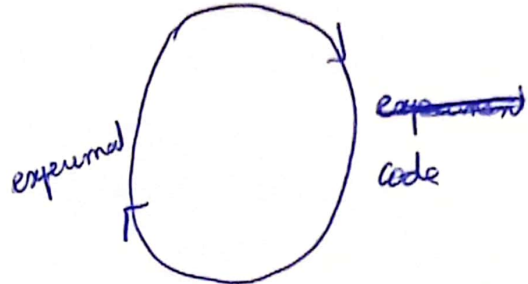
- normal regularization
- dropout
- data augmentation
- early stopping
- * batch normalization
- * bagging

~~Hyperparameters~~

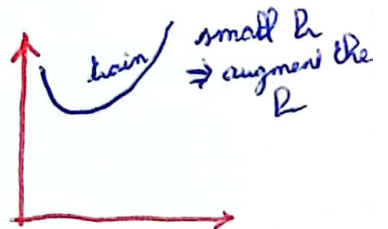
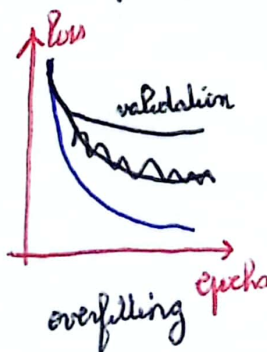
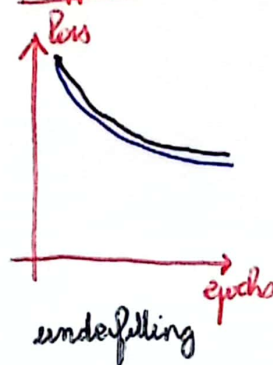
We plot cost function of # iterations
changing a value of a hyperparameter

\Rightarrow applied ML is a highly iterative process

Idea

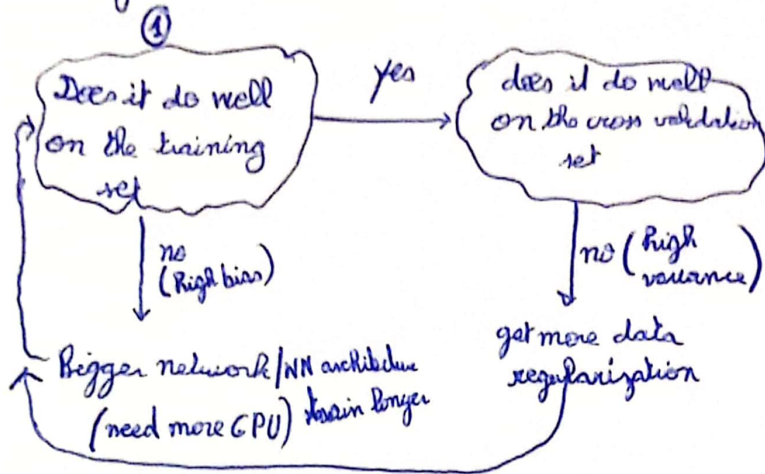


overfitting / underfitting with loss curves
d'apprentissage



bias / variance and neural networks

* Large neural networks are low bias machines



* a large neural network will usually do as well or better than a small one so long as regularization is chosen correctly

Apple "Bias Variance tradeoff"

we don't have as many tools to reduce bias without affecting variance (and vice versa) in ML algorithms

but in modern big data era:

- 1/ training bigger network almost reduces bias without hurting your variance as you regularize well
- 2/ getting more data always reduces variance without hurting bias

⇒ we don't have necessarily to give attention to the tradeoff

Regularized MNIST model

l1 = Dense(25, activation = 'relu',

kernel_regularizer = L2(0.01))

1

to have an idea about bias and variance of our model

* train set error

* dev set error

example 1

train set error : 1%
test set error : 11%

⇒ high variance

example 2

train set error : 15%

test set error : 16%

we suppose human error in this field ≈ 1%

⇒ train error >> human error
⇒ under-fitting (high bias)

example 3

train set error : 15%

test set error : 30%

human error : 1%

⇒ high bias

high variance (test error >> train error)

example 4

train error : 0.5%

test error : 1%

applied

Regularization

Logistic Regression

$$\min_{w,b} J(w,b)$$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

$$\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$$

L2 regularization

regularization parameter

$$L1: \text{regularization} : \frac{\lambda}{m} \|w\|_1 = \frac{\lambda}{m} \sum_{i=1}^n |w_i|$$

$\Rightarrow w$ will be sparse (a lot of zeros)

⚠ L2: regularization is much more used
 \downarrow
it is a weight decay

NN

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]})$$

$$= \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_2^2$$

$$\|w^{[l]}\|_2^2 = \sum_{i=1}^n \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2$$

norme de Euclidean

$$w^{[l]} : (n^{[l]}, n^{[l-1]})$$

why regularization prevents overfitting

setting $\|w^{[l]}\| \approx 0$ et donc $w^{[l]} \approx 0$

\Rightarrow ~~reducing~~ the impact of some neurons

\Rightarrow simpler network

if λ is very big \Rightarrow very very simple NN

\Rightarrow (underfitting)

if λ is small \Rightarrow overfitting

\Rightarrow we need to change λ to get the best parameter

dropout regularization

with dropout

supervised learning

during training: a percentage of ~~neurons~~ were set to 0

during testing: all neurons are used

⇒ the model during testing is more robust and can lead to higher testing accuracies.

① set a certain probability of eliminating a node in a neural network



we get a smaller NN

but we need to fix this probability

How to implement dropout with "inverted dropout"

illustrate with layer 3; keep_prob = 0.8

$d3 = \text{np.random.rand}(a3.\text{shape}[0], a3.\text{shape}[1]) < \text{keep_prob}$

$a3 = \text{np.multiply}(a3, d3)$

$a3 /= \text{keep_prob}$

if we have a total of 50 units

10 were shut off

$$z[4] = w[4] a[3] + b[4]$$

↑
20% of these elements

were = 0



% 0.8

if we have a lot of neurons in a layer

⇒ more risk of overfitting

⇒ we choose a high keep_prob (0.5)

else

we chose a high keep_prob (0.8, 0.9, 1.0)

if we have only one neuron in a layer: chose 1.0

other regularization methods

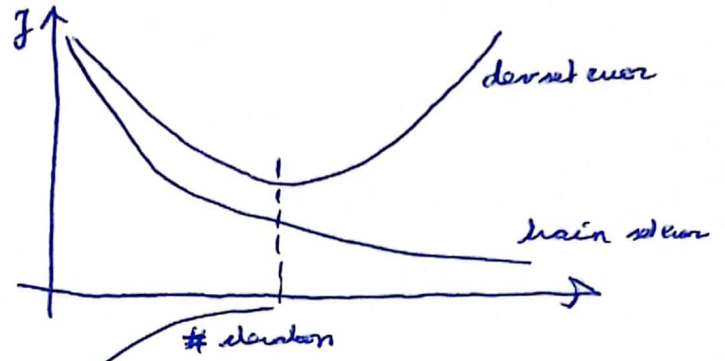
1/ data augmentation

add extra fake training examples

applying rotations, zooming

adding distortions

2/ early stopping



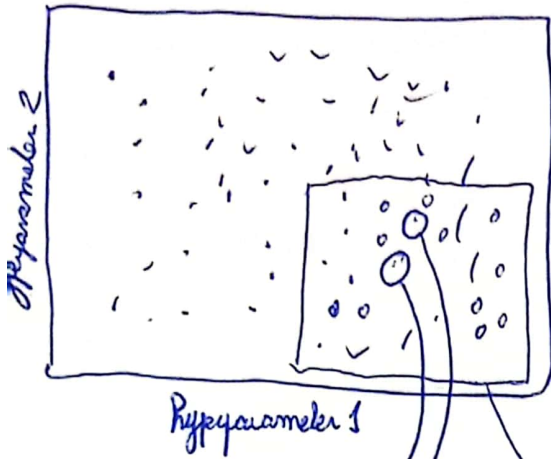
we stop the optimization algorithm here

Hyperparameter tuning

1/ use a grid for the Hyperparameters

2/ but in deep learning try random values instead of using a grid.

3/ coarse to fine



suppose we find the best results are when using
these set of hyperparameters

⇒ we zoom over this zone and generate more points
in it

appropriate scale for Hyperparameters

α :

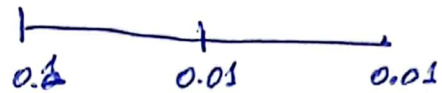
$r = -1 \times \text{np.random.randn}() \rightarrow [-1, 0]$

$\alpha = 10^r$

β : exponentially weighted averages Hyperparameters

$[0, 0, \dots, 0, 999]$

$1 - \beta = 0.1 \dots \dots 0.001$



$r = -$