## I Going deep learning or not?
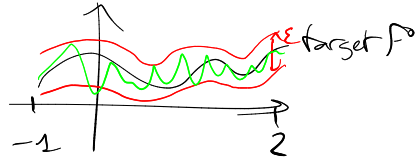
→ no guarantee to obtain a good solution

### Universal approximation theorems    (expressive power)

[Cybenko, 1989 ; Hornik, 1991] with just one hidden layer, one can approximate any $C^0$ function arbitrarily well (on a compact set)

[Sprecher 1965]

[Kolmogorov 1956] $\forall C^0 f, \exists N < +\infty,$

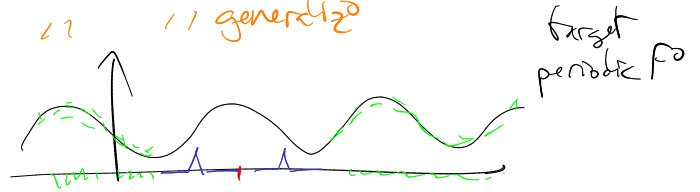$\exists \sigma \atop R \to R'$     $f \circleddash$ Network $_{N,\sigma}$

↳ hashing functions

⚠ existence theorems ⟹ doesn't provide the solution

⟹ doesn't tell whether the solution is easy to find (anything about the optim°)

⟹   "   "   "   " generaliz°

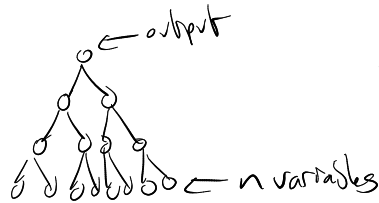target periodic $f°$

## Depth simplifies the approximation (estimation) task

• [Lin&co, 2017] multiplication of n variables

→ multiplicat° $\underline{\approx}$ 4 neurons of 2 variables

→ " of n variables :

← output

← n variables

binary tree ↓↑ $\log n$

#nodes = 4n

→ flat network: need $2^n$ nodes to be correct on binary inputs ⟹ learn by heart

exponential in input dimension
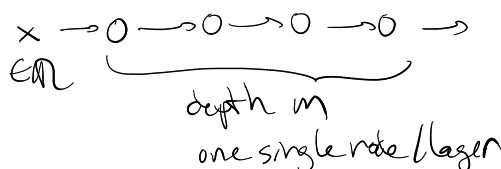
• [Telgarsky 2015]  target funct° :

$2^m$ triangles

$∘m$

1 / 2 / 3 / $2^m$

### Flat network

2 layers → require $2^{m/2}$ nodes

$\sqrt{m}$ layers → $2^{\sqrt{m}}$ nodes

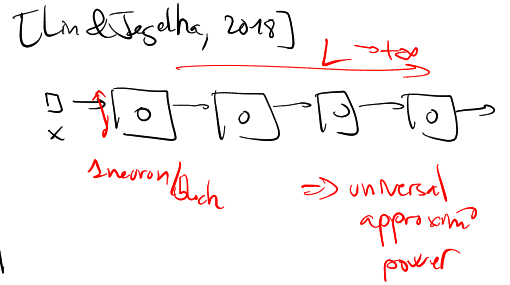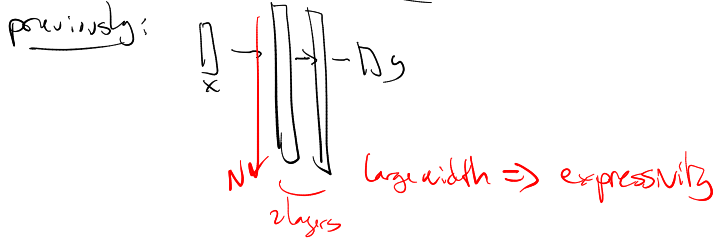x → ○ — ○ — ○ — ○ →

depth m
one single node / layer

thin deep network

[Mhaskar Liao & Poggio] about architecture suitability for $f^0$ estimation
2016

Theorem: if target $f^0$ = computational tree of input variables

then the # of samples required to train a

similar-shaped network : $O(d)$

vs $exp d$

target $f^0$

← inputs
$\in \mathbb{R}^d$

new network

• Depth is sufficient (without width)

previously:

x

— $N_y$

$N \downarrow$

2 layers

Large width ⇒ expressivity

[Lin & Jegelka, 2018]

$L$ → steps

x

1 neuron / block

⇒ universal
approxim$^o$
power

Does it work? When?

– computer vision & NLP (natural language processing)

↳ big success → designed to handle image properties
CNN

↳ invariant to translations
↳ hierarchical models
↳ handle pixel location (geometry)

vs random forests (e.g.)

spatially
structured
data

Convolution ⇒ any place

Ex: 1D temporal signal

– 3D videos

✓

Geometry is lost → no exploitat$^o$

– on the opposite:
if no data structure

ex: medical
observations;
temperature
coughing;

permut

random forests or SVM
might the job
will

+ need to learn the appearance of
objects at any possible location

⇒ if small data:

Neural Network → bad performance

⇒ #samples ⇒ huge

– big success in Reinforcement Learning?
α–Go, StarCraft, …
α-Go/α-0 : game of go : – 44 millions training games ( ≫ one life)
⪅ total Humanity?

– trained for only 4 hours on 5000 TPU
≃ GPU graphic card
very
good

= ≫2 years of 1 TPU

≃ 1000s years of 1 CPU

– computational power: NLP → huge networks
BERT, XLnet, GPT-3 … $10^{to}$ parameters
↓
state of the art
for
text translation,
etc.

⇒ re-used
through transfer learning

⇒ very long time to train
on huge GPU clusters

⇒ $ electricity

⇒ environmental impact

# Gap between classical ML & DL

Bayesian view $\Rightarrow$ information theory

## Reminder: classical ML

- samples $(x_i, y_i)$
- estimate $F : x_i \mapsto \sim y_i$
- quantifies "goodness" with "loss fo" =criterion

$exp(-\cdots)$

$$\Rightarrow \inf_{F \in \mathcal{F}} \sum_{examples} Loss(F(x_i), y_i) + Regularizer(F)$$

predefined parametrized family

e.g.
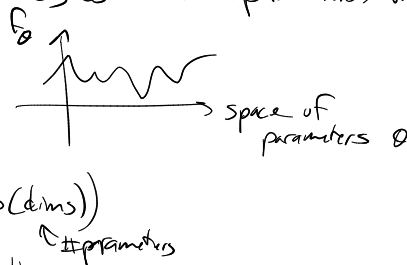$$\int_{\Omega} \left\| \frac{\partial F}{\partial x} \right\|^2 dx$$

- w/thout regularizer: **overfit**
- with " $\Rightarrow$ hope for good generalization

$\Rightarrow$ MDL: minimum description length paradigm

Occam's razor $\rightarrow$ prefer simpler models

$\hookrightarrow$ prefer models with fewer parameters

## Potential Issues with DL

- $10^6$ parameters ... Occam's razor? MDL?
- models : are able to overfit (easily)
- possible to train huge models without overfitting (still good generaliz°) **+ no regularizer!**
- " " " " " " with fewer samples than parameters

$\hookleftarrow$ gap between train error & test error

$\hookrightarrow$ estimator (of parameters) convergence?

- highly non-convex optimis°
  in a high-dimensional space
  $\Rightarrow$ Supposed to be very hard!
  (difficulty $= \sim exp(dims)$)

$f_\theta$

space of parameters $\theta$

$\hookleftarrow$ #parameters

- add noise to optimis° process $\Rightarrow$ works better
- train to optimize a criterion : cross-entropy
  but evaluate : with accuracy $\hookrightarrow$ not differential

- common recommandation : new task $\rightarrow$ new architecture $\rightarrow$ check able to overfit !
  (small part of data)
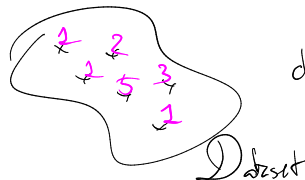  $\Rightarrow$ means enough expressive power

## A closer look at overfitting

[Zhang & co, 2017]

- huge models can do the job (might not overfit)
- " " can completely overfit also (able to)

theorem : $n$ samples $\in \mathbb{R}^d$

dataset   input dim

$\exists$ 2-layer network with $2n+d$ weights that can
with ReLU   represent any function
activat° f°   on such a dataset


Dataset

classific° task

random labels $\rightarrow$ perfect fit
$\Downarrow$
overfit

$\Rightarrow$ capacity of networks is not the issue
{ Rademacher complexity
{ Vapnik-Chervonenkis

# Palliatives for regularization

→ what about adding a functional regulariser?  $\mathscr{Y} = \int_{x\in\Omega} \|\frac{\partial f}{\partial x}\|^2 dx$

$\Rightarrow$ norm of a function

→ First: checking whether $f'' = 0$?

$\quad\hookrightarrow$ NP-hard  (ex: inputs = binary $\Rightarrow$ SAT pb)

$\Rightarrow$ stochastic approximations of norms = tractable

$$\mathscr{Y} = \sum_i \|\frac{\partial f}{\partial x}\|^2_{(x_i)}$$

## Dropout

each neuron will be replaced by 0
with probability $1/2$

important information to the next layer → duplicated

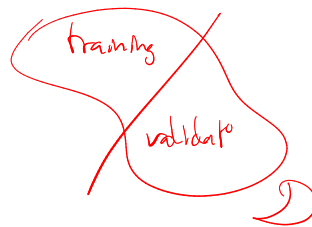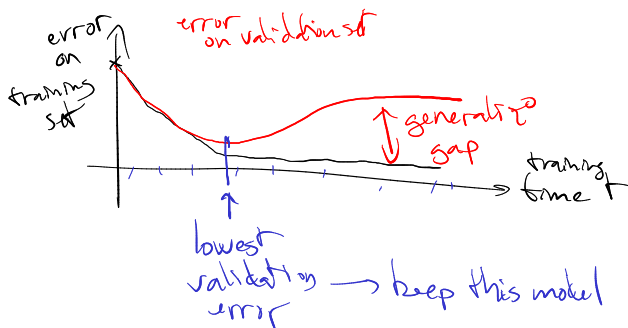$\Rightarrow$ many ≠ ways to express similar things

$\Rightarrow$ <u>robust</u>

* $\frac{\partial F}{\partial a_i} \approx 0 \Rightarrow$ make f smooth $\Rightarrow$ regularizer

* Bayesian point of view → ensemble method → robust

\* at test time,
use all neurons
with activities $/2$

$a_i$  activ° $= \sigma(\sum_i w_{i,i} x_i + S)$
<u>pre-activation</u>
activity $\Sigma$ activation

## Early stopping

error on training set

error on validation set

↑ generaliz° gap

training time t

lowest validation error → keep this model

training / validation

radius: maximum training time
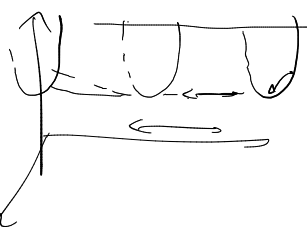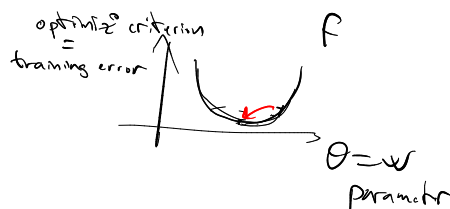↓
maximum irregularity

## Optimiz° noise acts as a regularizer
(due to SGD, or other)

[Poggio] : around convergence,
to a local minimum

Hessian matrix $\frac{d^2 F}{dw^2}$

→ if $>0$ : then no pb
$\hookrightarrow$ if not : pb  degenerated

optimiz° criterion = training error

$\theta = w$  parameter

$>0$ directions;
due to training sample

$= 0$ directions;
include test sample errors

$\Rightarrow$ make the Hessian $>0$ !

how? → <u>add a weight regularizer</u> : criterion $+ \lambda \sum_i \|w_i\|^2$
"weight decay"  works!

how not? → batch norm → improve optim° but doesn't make H $>0$

# Optimiz° Landscape

## Local minima & saddle points



Fully-connected

$N_1$ neurons

$x \to$

$\to \hat{s}$

permutation of neurons in one layer
$\Rightarrow$ get the same function

1 local minimum
$\Rightarrow$ duplicates of this minimum

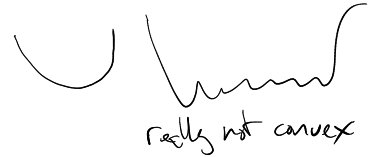$N_1! \times N_2! \times N_3! \cdots$

really not convex

- [Poggio] bound on the number of minima
  if activat°func° $\sigma$ = polynomial
  $\mathbb{R} \to \mathbb{R}$   ex: $\sigma(x) = x^3$

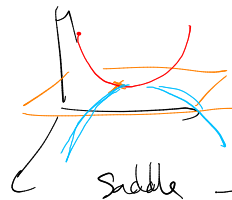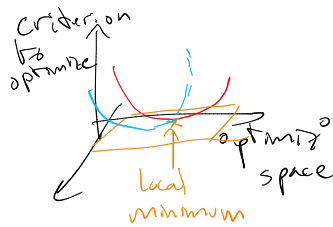  then (output)$(x)$ = polynomial$(x)$
  $\hookrightarrow$ loss$(x)$
  $\hookrightarrow$ Bezout's theorem $\Rightarrow$ bound depending on the degree

- So many parameters : local minimum = very strong notion   (huge neighborhood in optimiz° space)
  $\Rightarrow$ local minima = very good

- many local minima $\Rightarrow$ even more : saddle points



criterion to optimize

local minimum

optimiz° space

Saddle point $\to$ Issue : slow down the optimiz° process

critical points

[Dauphin, 2014]
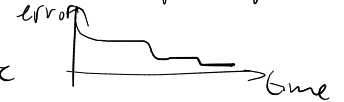Hessian : second derivative
$\hookrightarrow$ symmetric matrices

. standard law on random symmetric matrices :

$$P^T D P$$

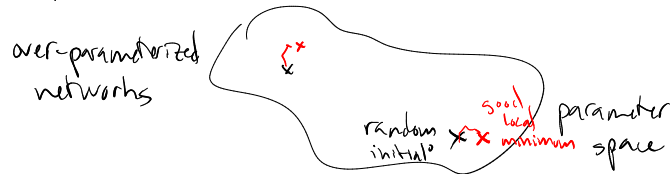$\uparrow$ pick a random diagonal matrix

$\uparrow$ pick an random orthogonal matrix

$\Rightarrow$ chances to have all eigenvalues of the same sign $\sim 2^{\#parameters}$

error

time

## Lots of works on convergence :
- always under strong hypotheses
  { specific to a particular architecture (ex: 2 layers)



over-parametrized networks
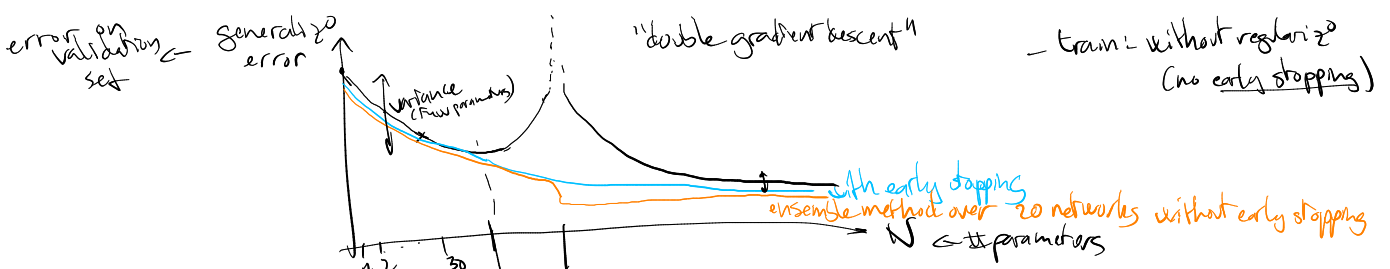
random initial°

good local minimum

parameter space

$\Rightarrow$ Francis Bach
$\hookrightarrow$ "lazy training regime"

## Over-parametrization helps

[Geiger et al, 2019]



$x$

$\to \hat{s} \to L$

fixed architecture

$N$ #neurons/layer

error on validation set — generaliz° error

"double gradient descent"

— train: without regulariz° (no early stopping)



variance (few parameters)

with early stopping

ensemble method over 20 networks without early stopping

N ⟵ #parameters

1 2 — 30

← classical training regime

critical N explosion

over-parametrized regime →

perfect number of degrees of freedom to fit the data

e.g: linear regression (without regulariz°)

linear system $\begin{cases} \alpha_1 x_1 + \alpha_2 x_2 \cdots = b_2 \\ \vdots \\ \end{cases}$

N variables ⟹ exactly one solution
N equations

⟶ data: a bit noisy ⟹ overfit exactly the noise



random init ×
complete overfit
Parameter space

→ prove theoretically: ⟵ without early stopping

denote by $\begin{cases} F_N & \text{a trained network with N neurons/layer} \\ \overline{F_N} & \text{: average of trained networks with all possible random init°} \end{cases}$

$\overline{F_N} \to F_\infty$   $\| F_N - \overline{F_N} \| = O\left(\frac{1}{\sqrt{N}}\right)$   or α?

⟹ the training of neural networks gets robust to init° for large N (#neurons)

⟹ don't worry about large # of neurons causing overfit / optimiz° issues
   ⟹ put as many neurons as you can

Further: [Nakhiran & co, 2019]

valid° error



training time

#data

Large architecture without regulariz°

valid° error

more data



#Neurons

More: [Choromanska & al, 2015]
— optimiz° a network ⟶ Hamiltonian of spherical spin-glass model
   ⟶ results on statistics over critical points

training loss



global optimum

most in a narrow band

local minima

⟶ #minima outside that band ∝ $e^{-\text{Network size}}$

⟶ gets harder to find ⟶ but: global minimum = prone to overfit

parameter value

About Minimum Description Length paradigm

— MDL is lost? #parameters = huge, but:
   — no precision needed when encoding parameters
   — high redundancy: different parts of the network may compute similar things
      ⟶ NN compression $\begin{cases} \text{prune} \\ \text{quantize weights} \\ \text{tensor factorizat°} \end{cases}$ compression factor ≥ 100 ⟶ [Lavros & al, 2017]
      ⟹ run online video object detection on a smartphone



fully connected

weights ≈ low rank

double precision: 8 bytes
⟵ float → 4 bytes

parameter values

8 possible values ⟶ 3 bit (< ½ byte)