

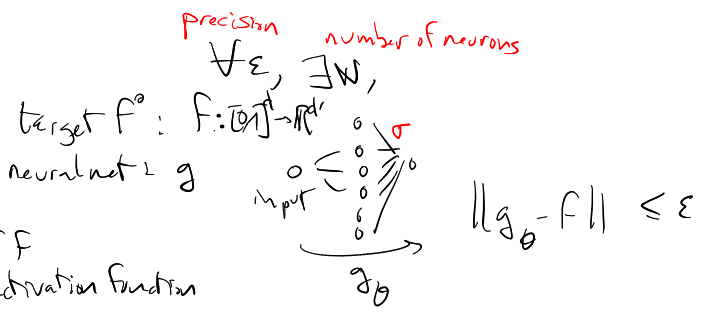
Deep Learning in Practice

Chapter 1: Deep Learning vs. Classical ML & optimisation

I Going Deep or not?

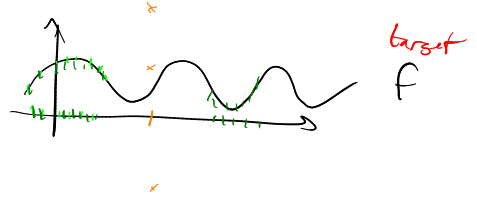
- no guarantee
- Universal approximation theorems

[Cybenko, 1985] [Hornik, 1991]:
 [Sprecher, 1965]
 [Kolmogorov, 1956]: exact decomposition of F
 by tuning the activation function



$$\forall F \in C^0, \exists \sigma, N < +\infty, g_{\theta, \sigma, N} = F$$

Being able to approximate $(\exists \theta, \dots)$ \neq being able to learn & generalize



- Depth simplifies the approximation / estimation task

[Why does deep & deep learning work so well?]

ex: n input variables

input $x \in \mathbb{R}^n$

target: $\prod x_i$

\rightarrow provided σ is smooth, possible to approximate multiplication of 2 variables with 4 neurons

\rightarrow for n variables:



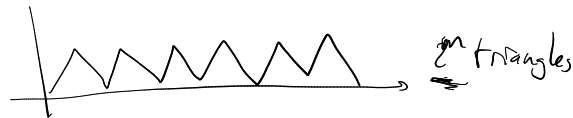
$\sigma(a+b) = \sigma(a) + \epsilon \sigma'(a) + \frac{\epsilon^2}{2} \dots$

\hookrightarrow with only 1 hidden layer: 2^n neurons at least.

\rightarrow solves only the case of binary inputs (0, 1)

- [Telshevsky]

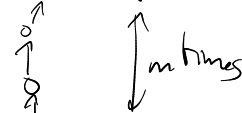
ex: target function:



$$\equiv \left(\text{triangle} \right)^{\circ m}$$

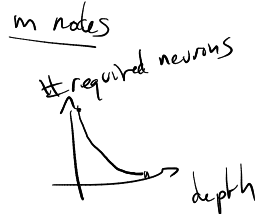


\Rightarrow representable with a thin deep network (depth: m , 1 neuron)

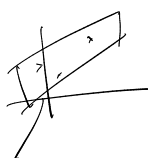


\hookrightarrow with a flat network (2 layers): need $\frac{m^2}{2}$ nodes

\hookrightarrow \sqrt{m} layers $\rightarrow 2^{\sqrt{m}}$ nodes



- [Passis 2016]



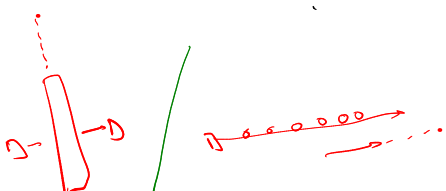
$$F(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

d parameters to tune $\rightarrow \approx d$ samples

$F(\vec{x}) =$ Function of \vec{x} & of d parameters
 \hookrightarrow need $\approx d$ samples

If model not known: $\exp(d)$ samples for a flat network

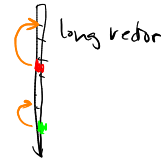
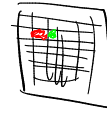
- Deep is sufficient (without width) [Lindesjek, 2018]. ResNet with 1 neuron layers \rightarrow universal approximation property



- Does it work? When?

- big success: computer vision & NLP
 - ↳ convolution: explicit geometry handling
 - translation invariant
 - hierarchical models

task: image classification



≠ random forest: can't incorporate geometry (yet)

↳ few dimensions "independent" in the sense no geometry

- success in Reinforcement Learning

Alpha-Go Alpha-0

- like million training games >> one life ~ humanity

- 4 hours training on ... 5000 TPU > 2 years of 1 TPU ~ 1000s gens on single CPU
↓
quite slow training

↳ computational training cost
NLP: huge networks (BERT, XLnet...) → ≈ 100 000 \$

2) Gap between deep learning & classical ML/optimization

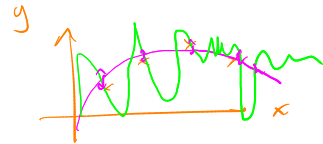
- Reminder: ML setup

- dataset: (x_i, y_i)
- estimate $F: x_i \mapsto y_i$
- loss function Loss: quantifies the quality of the candidate F

$$\hookrightarrow \inf_{F \in \text{predefined parameterized family}} \sum_{\text{examples } i} \text{Loss}(F(x_i), y_i) + \text{Regularizer}(F)$$

$$\| \frac{dF}{dx} \|^2$$

↳ make F smooth \Rightarrow better generalization properties



↳ Minimum Description Length principle (MDL)
↳ Occam's razor

- Now with deep learning:

- millions of parameters
- possible to train networks with 10^6 parameters without overfitting → still, good generaliz^o
- and this with fewer samples than parameters
↳ estimator convergence?
- classical optimization: the fewer parameters, the better
↳ highly non-convex optim^o: thought to be very hard
- add noise in the optim^o process \Rightarrow helps!
- training criterion: cross-entropy
but evaluate on another one: accuracy ← not differentiable
- common recommendation: check that the model is able to overfit
- no regularizer!

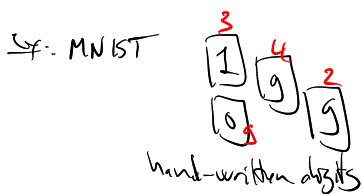
* Closer look at overfitting

- huge models don't necessarily overfit

- [Understanding DL requires rethinking generalizability, 2017]

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

⇒ 2-layer neural network, with ReLU activations and $2n+d$ weights that can associate each sample with a randomly-chosen target



random labels

$2n+d = 2000$ neurons



100% accuracy for these random labels

1000 samples
 $d \geq 1000$

even very small networks are completely able to overfit

⇒ capacity is not the issue (Rademacher, Vapnik-Chervonenkis)

- how to detect overfitting? / how to prevent it?

- possible regularizers: weight decay, drop out, data augmentation, ...

not clear: not sufficient necessary

- early stopping
- SGD (optimizer) regularizes

Alternatives For regularization

What about functional regularizers? i.e. $\int \left\| \frac{dF}{dx} \right\|^2 dx$

[Blaschke's theorem]: first step: check that the norm is 0? i.e. $F=0$ everywhere $[0,1]^d$

ex: all inputs are booleans

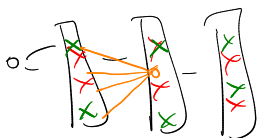
$F: \{0,1,1,0, \dots\}$

$[0,1]^d \rightarrow \mathbb{R} \quad F=0$

↳ NP-hard

↳ stochastic approximation of the norm: feasible

- Drop out training step + pick randomly half of the neurons $\rightarrow 0$



→ improve the performance (→ overfit)

half of inputs (if that neuron) are missing

⇒ robustness to inputs → redundancy
activity a_i (neuron i)

$\frac{dF}{da_i} \approx 0$

→ not far from $\left\| \frac{dF}{dx} \right\|$

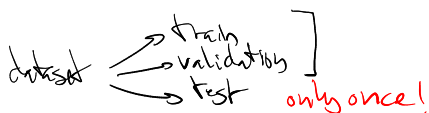
[Drop out as a Bayesian Approximation]

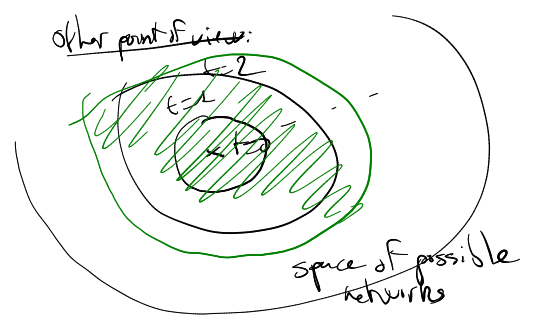
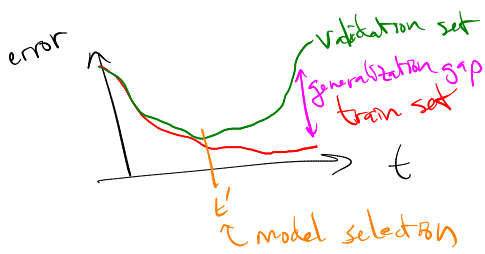
↳ each drop out realization: a \neq sub-network

↳ Family of (sub-)networks

↳ ensemble technique

Early stopping





Optimize noise

[Possio]

- around convergence:

Hessian: $\frac{d^2F}{dx^2}$

↳ if $H > 0$ (definite positive)

[all eigenvalues > 0]: SGD noise will not harm



↳ if not: $Loss = \sum_{\text{training set}} \text{errors}$

⇒ error ≈ stable for points in training set

↳ error ↑ for test points

⇒ SGD overfits after a while

- be sure that $H > 0$:

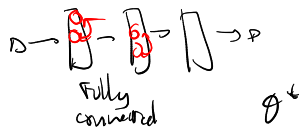
add a weight regularizer: $\sum_p |w_p|^2$

- weight normalization (batch-norm) ⇒ $H \not> 0$ ⇒ still overfit

Optimize landscape

- local minimums:

- global minimum ↔ θ^* particular parameter^s of the network



N! possible θ^* for some function F_{g^k}
From perturbations

↳ bound on the maximal number of local minima - [Possio]

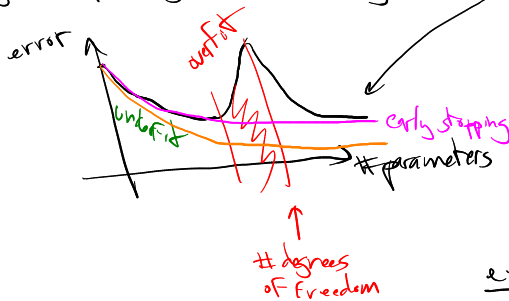
- saddle point:



Overparameterize

[Scaling descrip^o of generaliz^o ---]

"double gradient descent"



new!

- fixed architecture with variable width
- train without early stopping for a fixed # of epochs
- without regularization

→ train 20 models without early stopping
↳ ensemble method on top

ex: linear regression

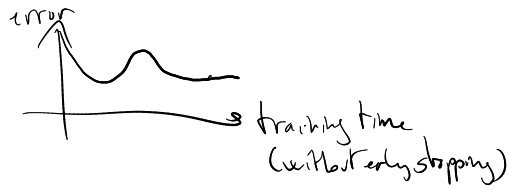
N samples of dim d , desired outputs y_i + noise x_i

N eq^s with N unknown $F(x_i) = y_i$

→ 1 unique solution non stable

[Deep double descent]:

ex: of \uparrow dataset size \Rightarrow \downarrow performance



\Rightarrow take away message: put as many neurons as you can.

H/b: too many neurons lead to overfit? NO!

MDL: high redundancy

\rightarrow counting #parameters: not the right way to estimate model complexity

\hookrightarrow ex: compression techniques \Rightarrow 100 size with similar accuracy
 \hookrightarrow video object detection on smartphone [Welling]