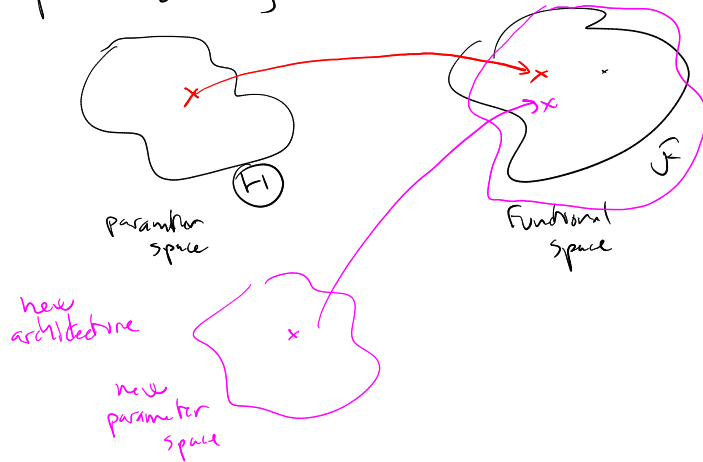


Architectures

I Architectures as priors on function space

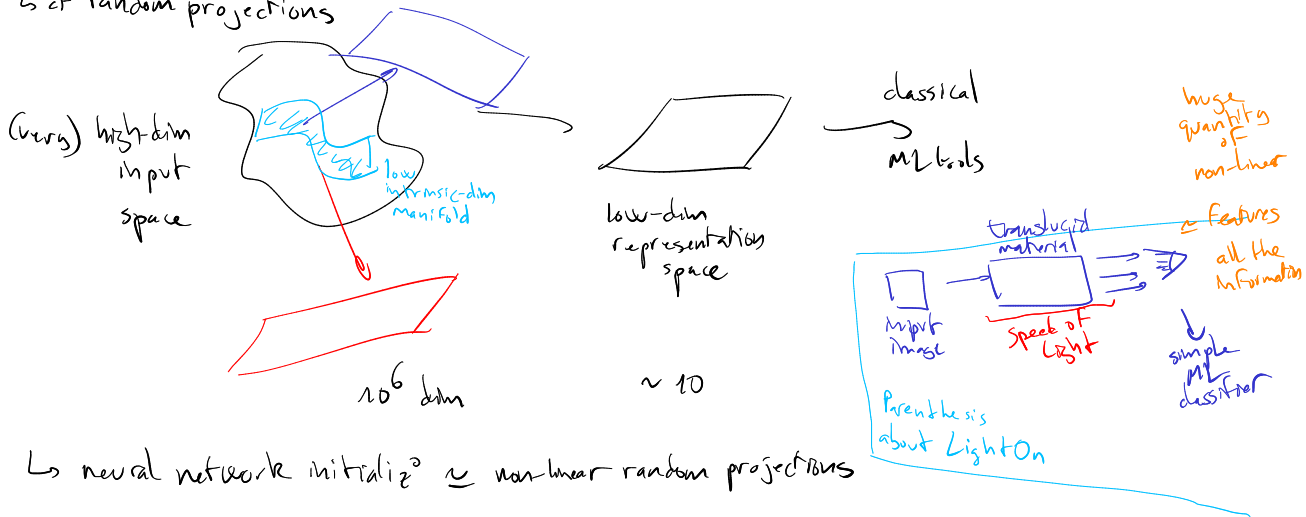
Change of Paradigm

- classical ML: design features by hand
- vs deep L: meta-design of features → design architectures able to produce the features we think about
- an architecture = a parameterized family



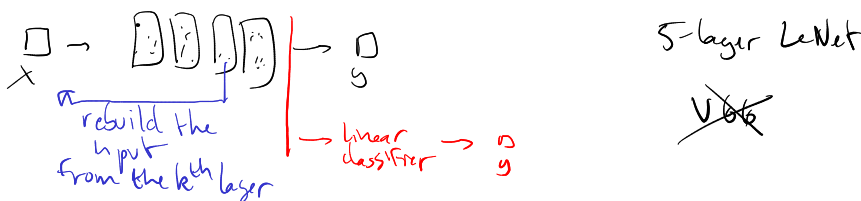
Architecture: prior on the function

- prior:
 - as a constraint: what is expressible with this architecture?
 - most (reasonable) networks have already a huge expressive power
 - probabilities: easy to reach?
- with random weights: reasonable function or not?
 - ↳ good architecture: grabs good features "naturally"
 - ↳ of random projections



↳ neural network initializ^o ≈ non-linear random projections

* In some cases: most of the performance is due to the architecture, not to the training



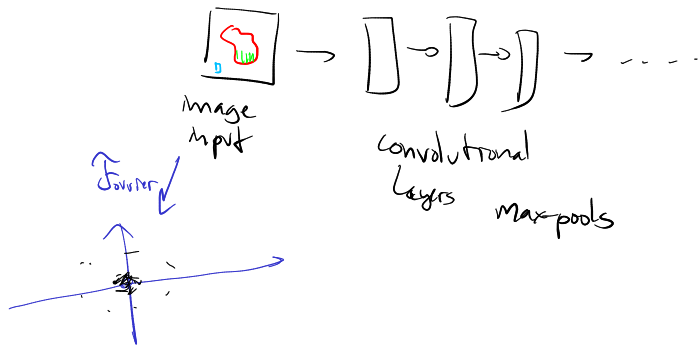
↳ Extreme Machine Learning: learn only the last layer ~~VGG~~

↳ quality of layers with random weights

↳ quantify the information flow: information theory → "information bottleneck"

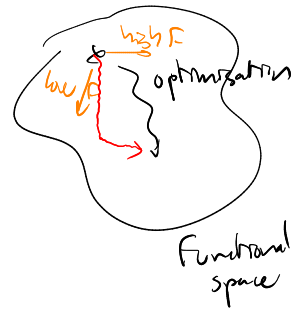
Bias of the architecture

→ move to Fourier space (CNN)



show that lower frequencies are learned first

→ look at the "big picture" before the details



→ stacking fully connected layers

Linear:



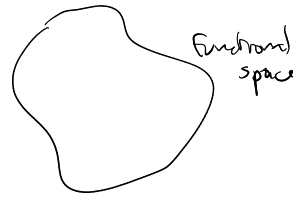
$$y = M_3 M_2 M_1 x$$

M'

prob \sim law;
 $\Rightarrow M_{ij} \sim \mathcal{N}(0, 1)$
 \Rightarrow proba over M

distribution of M, M'
 \neq of M

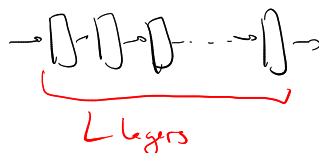
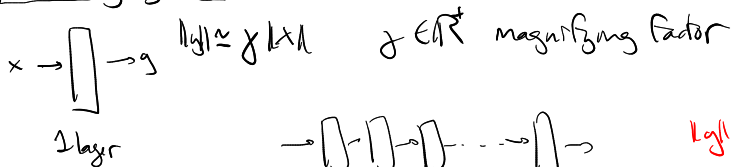
expressive power = the same



Random initialization:

→ random, according to which law?
 ↳ choose a law yielding good properties

* exploding/vanishing gradient



$\|y\| = \gamma^L \|x\|$
 $\hookrightarrow \gamma < 1: y \approx 0$ (vanish) \Rightarrow wish for good start
 $\hookrightarrow \gamma > 1: y \approx \infty$ (explode)

Same for backpropagation

$$\| \frac{dL}{dW} \| \approx (\gamma')^L \left\| \frac{dL}{dx} \right\|$$

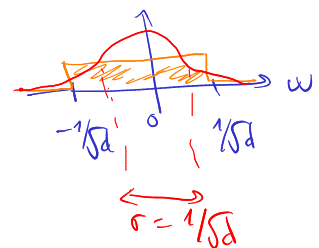
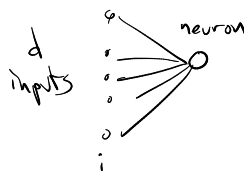
prevent optimize

→ wish for good optimization properties

Xavier Glorot initialization

$$w \sim \mathcal{U} \left[-\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}} \right]$$

or $\sim \mathcal{N} \left(0, \sigma^2 = \frac{1}{d} \right)$



justification: $\begin{cases} w_i \text{ iid of mean } 0, \text{ standard deviation } 1/\sqrt{d} \\ x_i \text{ iid of mean } 0, \text{ variance } 1 \end{cases}$

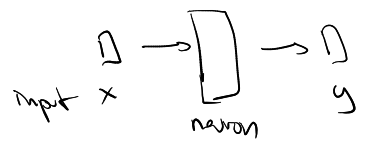
neuron computes: $y = \sum_i w_i x_i$

→ statistical properties?

$$\hookrightarrow \text{mean: } \mathbb{E}_x \mathbb{E}_w [y] = \mathbb{E}_w \left[\sum_i w_i \mathbb{E}_x [x_i] \right] = 0$$

↳ variance: $E[y^2] = E[E[\sum_i w_i x_i]^2] = \sum_i E[x_i^2] E[w_i^2] = 1$

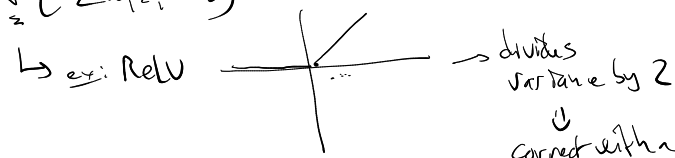
over
↓
1
1/d



$x_i \sim \text{law mean } 0$ $y \sim \text{law mean } < 1$

$\Rightarrow \gamma = 1 \rightarrow$ no exploding/vanishing issues (reasonable outputs)

* non-linear activation: $y = \sigma(\sum w_i x_i + b)$



- bias b? → initialize to 0

- other initialization laws (similar properties): He et al

[Boris Hanin, 2018] Jacobian properties

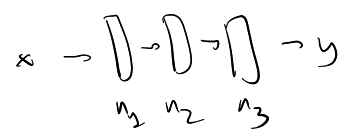
→ at initialization,

$\frac{df}{dx}$

Jacobian matrix: dim outputs × dim inputs



- n_l : number of neurons in layer l
- $\text{var}(\frac{df}{dx})$ is fixed = $2/n_l$
- $\text{var}(\left(\frac{df}{dx}\right)^2) = E\left[\left(\frac{df}{dx}\right)^4\right]$
- $\beta e^{\sum 1/n_l} \leq \dots \leq \alpha e^{\sum 1/n_l}$



⇒ better to choose equal widths than one thin layer

⇒ put sufficiently many neurons in each layer

Designing architectures easy to train

→ "deep" networks: not that deep in terms of optim^o

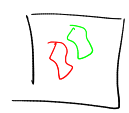


often: small number of layers → next part to cross

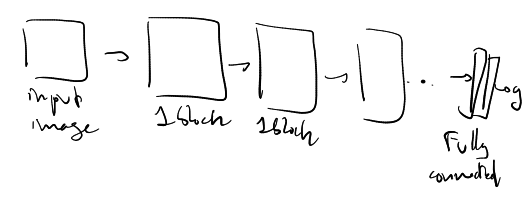
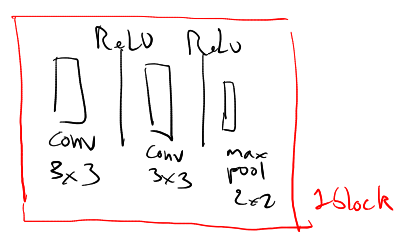
→ need "depth"?

II Architecture zoo

For vision problems: CNN → share local filters ⇒ invariance to translation
 ↳ hierarchical model ⇒ much greater generaliz^o power

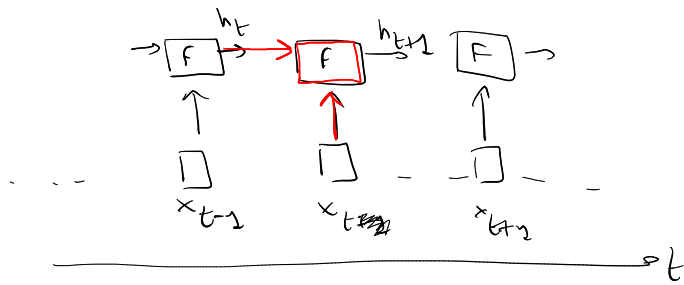


typical block:



Recurrent networks:

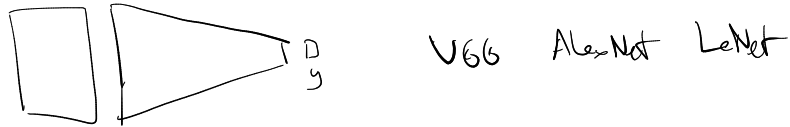
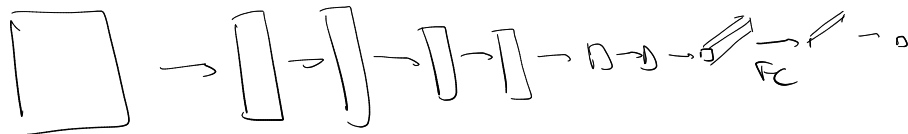
Basic RNN: $h_{t+1} = F(h_t, x_t)$



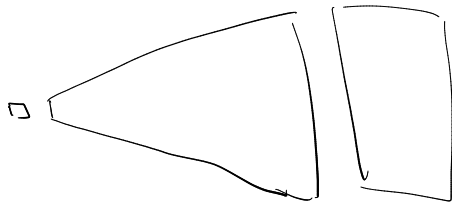
- memory
 - leaky: $h_{t+1} = h_t + F(h_t, x_t)$
 - gates
- ⇒ LSTM
GRU
Minimal RNN

Scale & resolution

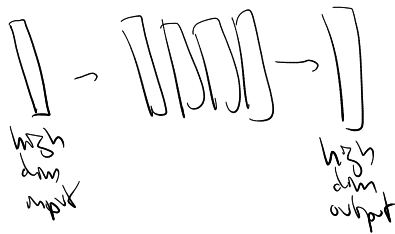
- classification tasks in computer vision: pyramidal approach



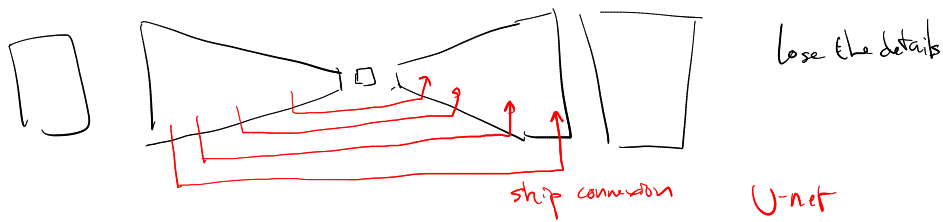
- generation:



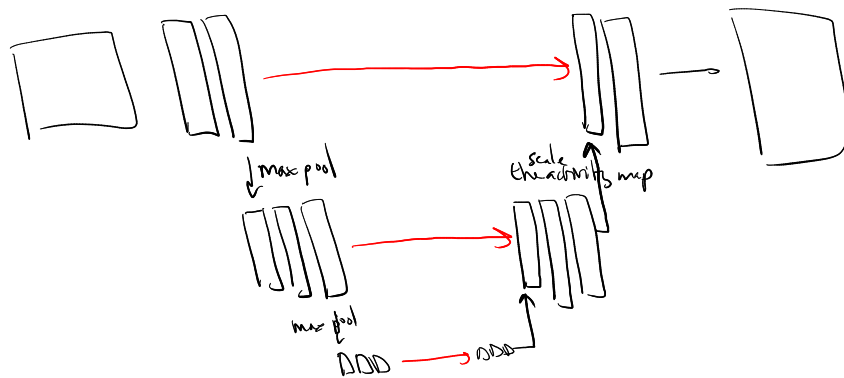
- regression / high input/output dim



ex: transformed image (of same resolution)



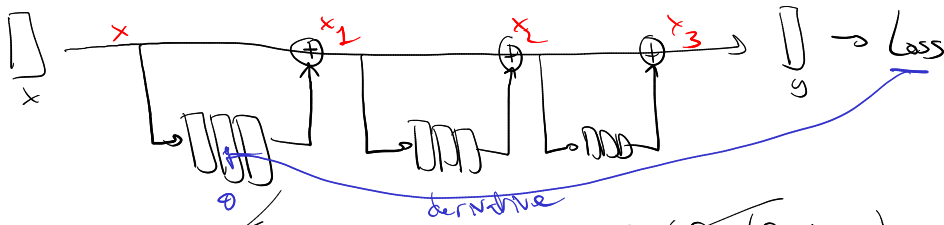
U-net



ex: segmentation

Depth & mixing blocks

ResNet:



$$x' = F(\theta)$$

$$x = x + F(\theta)$$

$$y = F_3(F_2(F_1(x)))$$

$$y = x + F_1(x) + F_2(x + F_1(x)) + F_3(x + F_1(x) + F_2(x))$$

$$\frac{\partial L(y)}{\partial \theta} = \frac{\partial L}{\partial F_3} \times \frac{\partial F_3}{\partial F_2} \times \frac{\partial F_2}{\partial F_1} \times \frac{\partial F_1}{\partial \theta}$$

$$= \frac{\partial L}{\partial y} \times \left(\underbrace{\frac{\partial F_1}{\partial \theta}}_{\text{direct Feedback}} + \frac{\partial F_2}{\partial F_2} \frac{\partial F_2}{\partial \theta} + \frac{\partial F_3}{\partial F_2} \left(\frac{\partial F_2}{\partial F_2} \frac{\partial F_1}{\partial \theta} + \frac{\partial F_1}{\partial \theta} \right) \right)$$

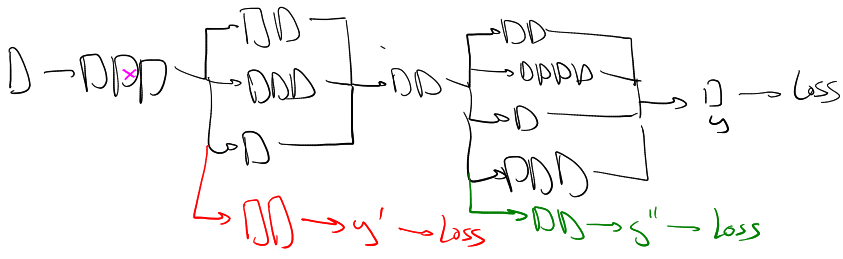
Highway networks

→ same idea + attention

Orthogonal matrices

→ train a FeedForward network: 10^4 layers

Inception



Auxiliary losses

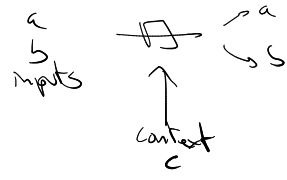
Start training: $Loss(y) + Loss(y') + Loss(y'')$

↓

$Loss(y)$ / / when optimizer goes fine

Attention

basic ex: transistor:



$$y = w_a a + w_b b$$

Linear combination
w's: fixed

$$y = \underbrace{w_a(c)}_{\alpha} a + \underbrace{w_b(c)}_{(1-\alpha)} b$$

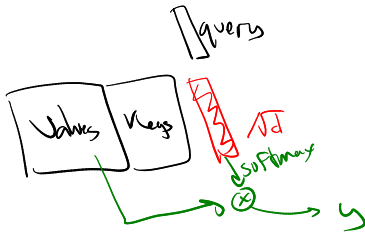
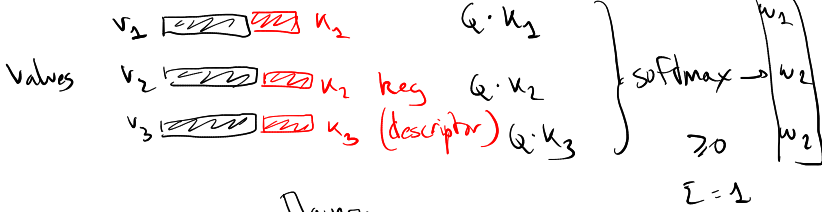
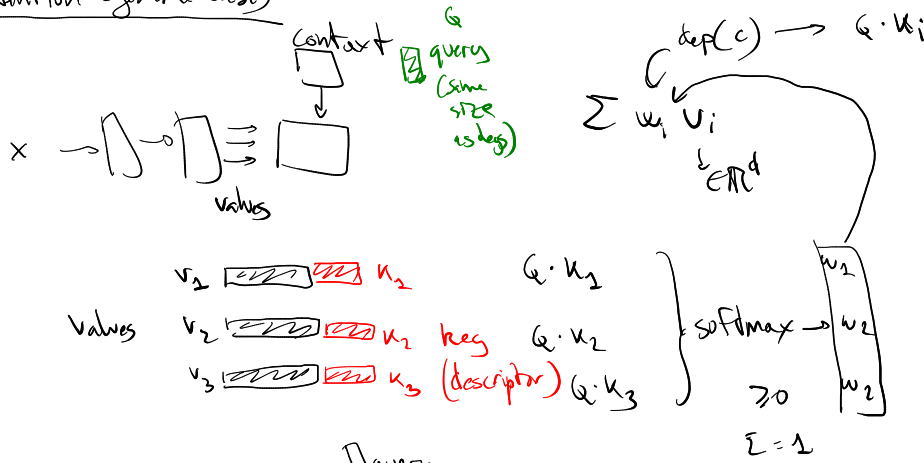
interpolation

$$y = \alpha a + (1-\alpha) b$$

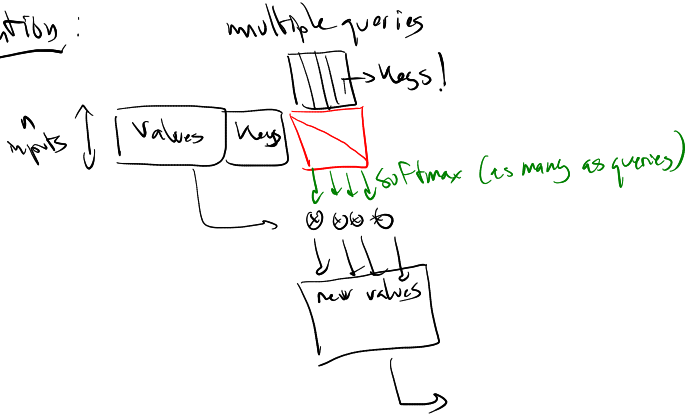
$$\alpha = F(c) \quad \alpha \in [0, 1]$$

$$= \text{softmax}(\dots) \text{ sigmoid}$$

Attention (general case)

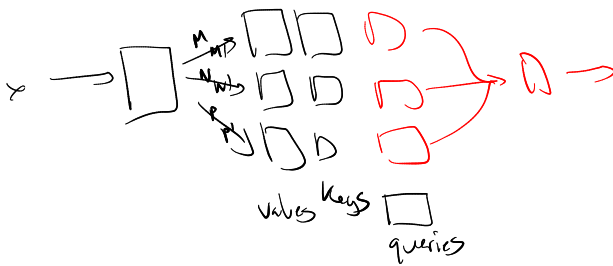


Self-attention:



Multihhead:

multiple attention blocks in parallel



Multihhead

$$V = M \times N \times P$$

$$K = M' \times N' \times P'$$

Q Fixed (parameters)
 $M'' \times$ or P'' (other things)

Graph-NN (GNN)

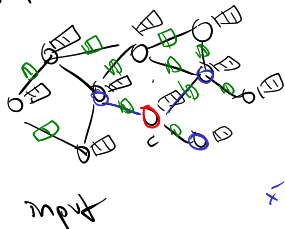
- analogy of conv for graph
- computer vision



$$\sum w_i x_i$$

x_i patch window size

- any graph!



needs to adapt to the shape of the subgraph

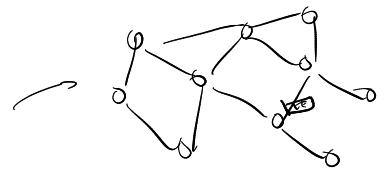
Layer

$$z_c = 2a_c + (0.3) \sum_{\text{neighbors of } c} a_n$$

parameters

neighbors of c

or average



same graph output

for each node: $F(a_c, \{a_{\text{neighbors}}\})$
 ↳ function needs to be generic to adapt to varying # of neighbors

using edge information:

$$F(a_c, \{(a_i, e_{i \rightarrow c})\}_{i \in \text{neighbors}(c)})$$

$$g_c = 2a_c + \sum_{i \in \mathcal{N}(c)} w_i(e_{i \rightarrow c}) a_i$$

↳ GAT: graph attention network

↳ towards attention mechanism
 ↳ M edge parameters

⇒ very flexible architecture

{ set of nodes with activities
 set of edges with activities

F

{ new values for nodes
 new activities for edges
 same graph

→ keeps the same graph

↳ graph classification?

input: graph

→

decrease regularly the size of the graph

or FC

output: GPR?