

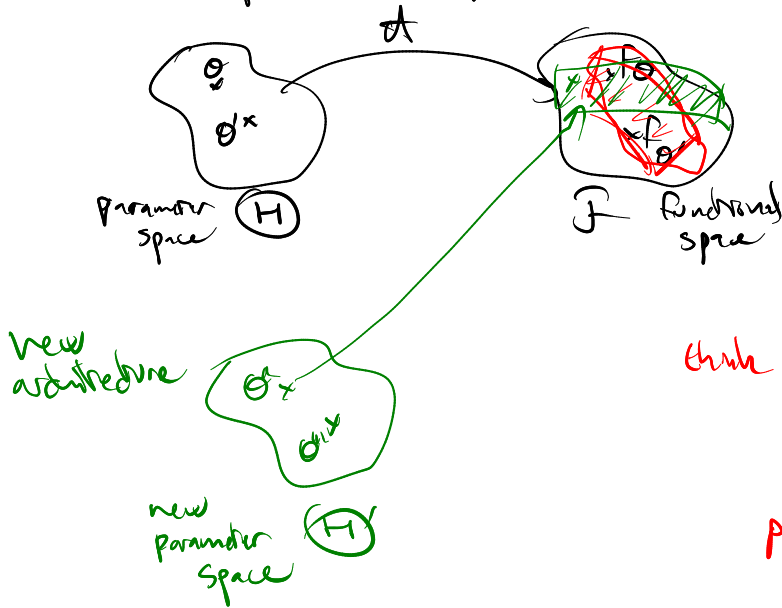
# Chapter 3: Architectures

## II Architectures are priors on function space

### Change of paradigm:

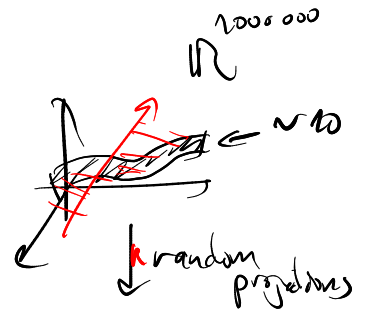
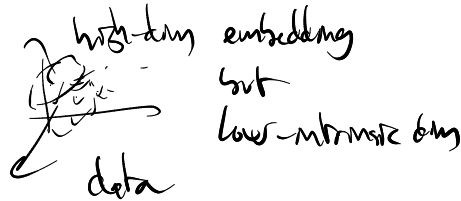
- classical ML: design features by hand
- deep ML: meta-design of features  $\rightarrow$  design the architecture

An architecture = parameterized space of functions



### Architecture bias

$\rightarrow$  random projections:



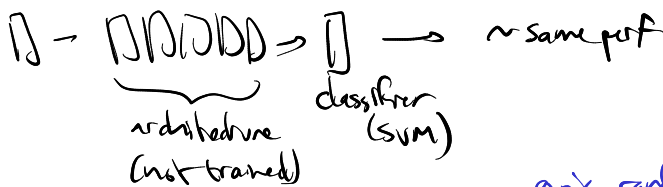
$\rightarrow$  with NN: random non-linear "projections"

all information is there

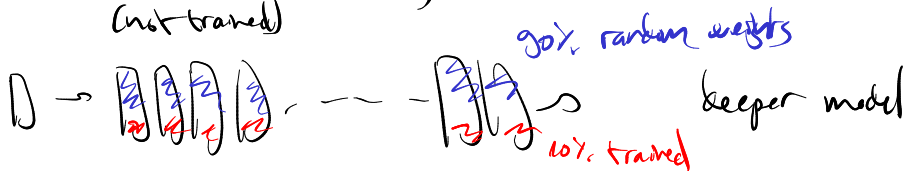


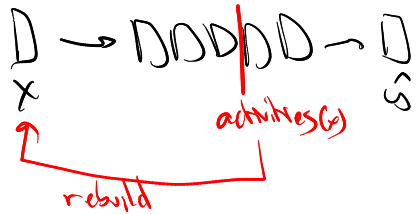
$$\mathbb{D} \xrightarrow{x} \mathbb{D} \xrightarrow{ReLU} \mathbb{D} \xrightarrow{ReLU} \mathbb{D} \xrightarrow{ReLU} \mathbb{D}$$

In some cases: most of the performance is due to the architecture (and to the training) "extreme ML"



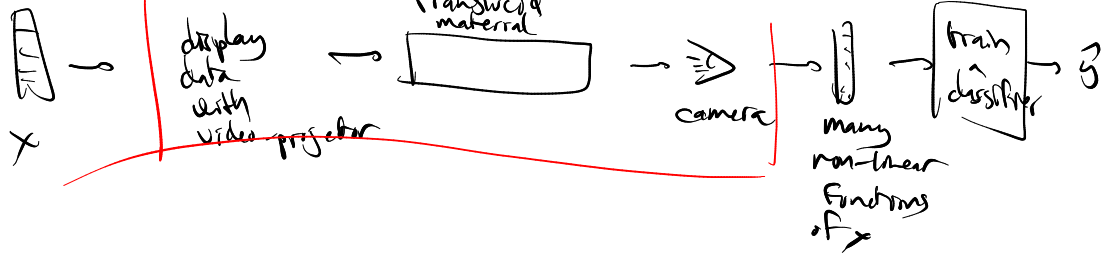
~~VGG~~ AlexNet



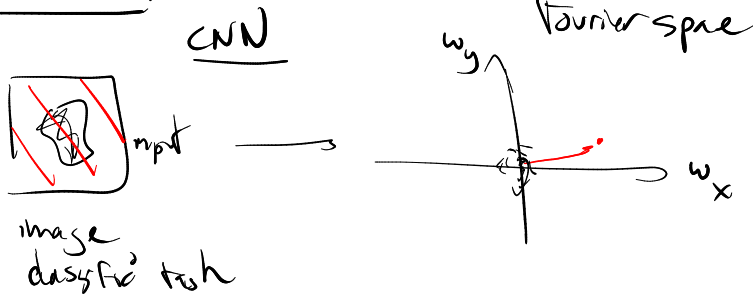


if random weights: good reconstruction of  $x$   
 if trained: not as good (keep only inform. relevant for the task)

Light On:

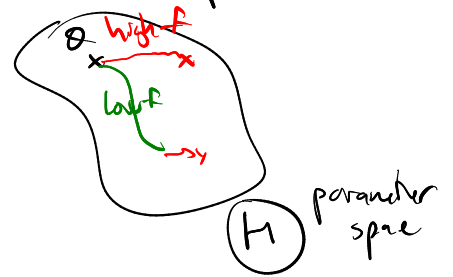
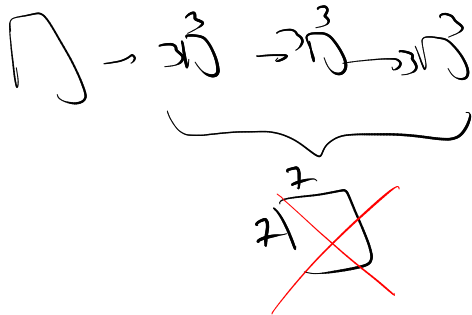


Architecture bias:



train CNN: focus on low frequencies first

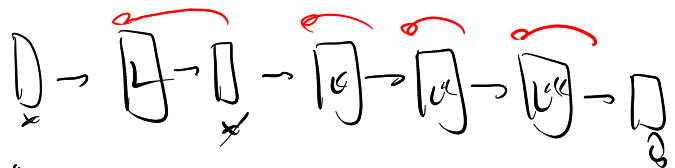
main objects before pixel noise



Random initialization = with which law?

$x \rightarrow$

$$\|x\| = \|x\|_2 = \sqrt{\sum x_i^2}$$



$\|x\| < \|\hat{x}\|$

$\|L(x)\| \approx \gamma \|x\|$

$\gamma = 10$

Recurrent networks

$\gamma > 1 \rightarrow$  exploding

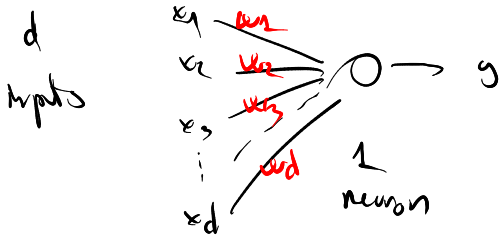
$\gamma < 1 \rightarrow$  vanishing

$\|y\| = \gamma^L \|x\|$

$A(x)$

VF

Xavier Glorot's initialization



$$y = \sum_i w_i x_i \quad (+b)$$

Hypothesis: inputs  $x_i \sim \mathcal{N}(0, 1)$  or mean  $\Rightarrow$   $\sigma = 1$

pick  $w_i \sim \mathcal{L}$

goal: find  $\mathcal{L}$  st.  $\text{var}(y) = 1$

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

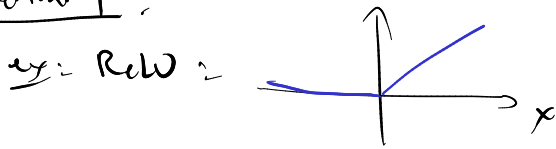
or  $w \sim \mathcal{N}(0, \sigma^2 = 1/d)$

Proof:  $y = \sum_i w_i x_i$

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

variance:  $\mathbb{E}_x \mathbb{E}_w [y^2] = \mathbb{E}_x \mathbb{E}_w \left[ \sum_i w_i^2 x_i^2 + \sum_{i \neq j} w_i w_j x_i x_j \right] = \mathbb{E}_w \left[ \sum_i w_i^2 \mathbb{E}_x [x_i^2] \right]$   
 $= d \mathbb{E}_w [w^2]$   
 $= d \cdot 1$

With activation  $\rho$ :



$$\text{var}(\text{ReLU}(x)) = \frac{1}{2} \text{var}(x)$$

$\Rightarrow$  standard deviation:  $\sqrt{2}$

So correcting factor:  $\times \frac{\sqrt{2}}{\sqrt{d}}$

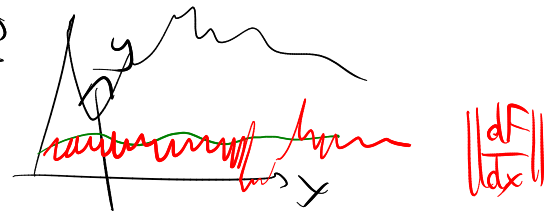
Bias: mit two

Other approaches: Kerning the (d,d)

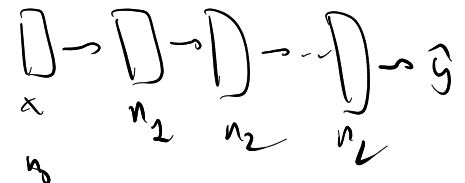
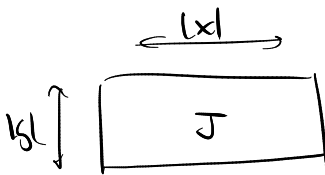


Further: (Boris Hanin, 2018): Jacobian properties

$$D_x \rightarrow \underbrace{DDDD}_{\mathcal{F}} \rightarrow y$$



$$J = \frac{dF}{dx}$$



$$\text{var} \left( \frac{dF}{dx} \right) = 1/n_0$$

$$\alpha \cdot \sum_i 1/n_i$$

$$\leq \text{var} \left( \frac{dF^2}{dx} \right) : \epsilon$$

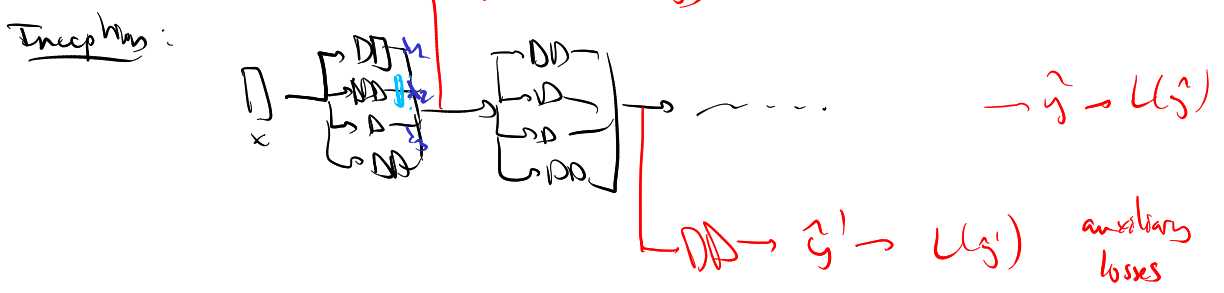
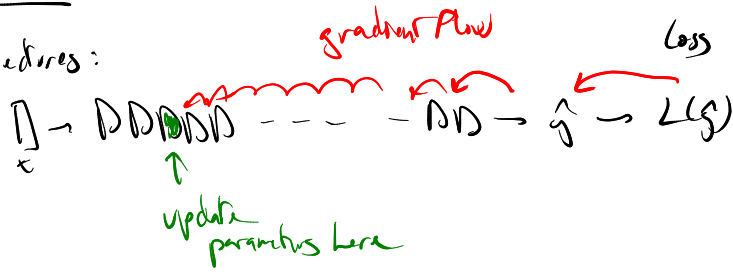
$$\beta \cdot \sum_i 1/n_i$$

Neural Tangent Kernel (NTK) :  $\nabla_{\theta} F \times \nabla_{\theta} F^T$

limit case  $n_l \rightarrow +\infty$

Design easy to train arch.

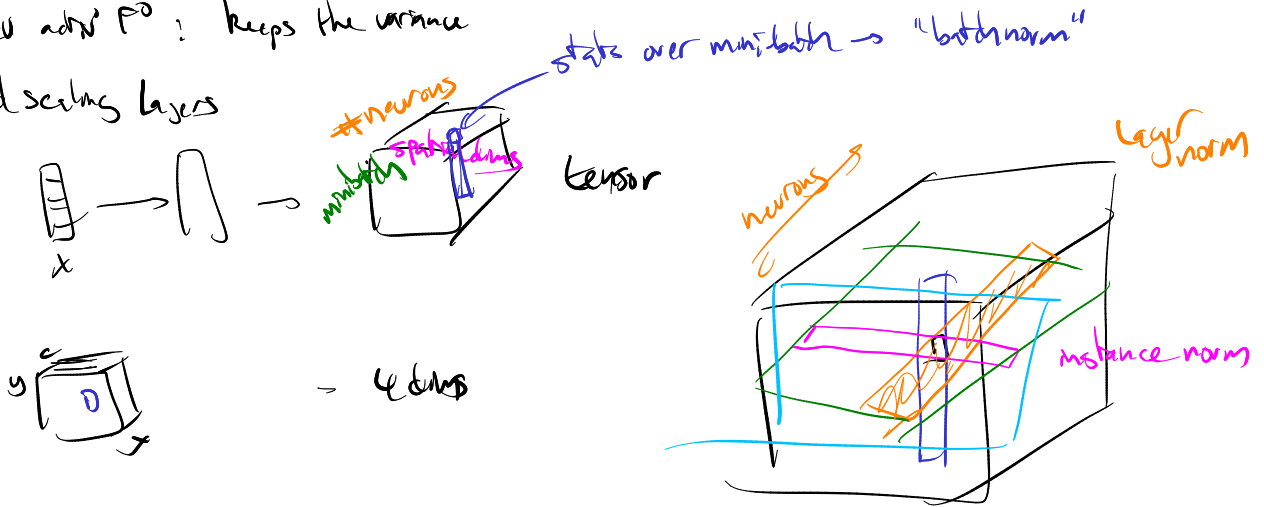
- deeper architectures:



Norms:

- SELU activation: keeps the variance

- add scaling layers



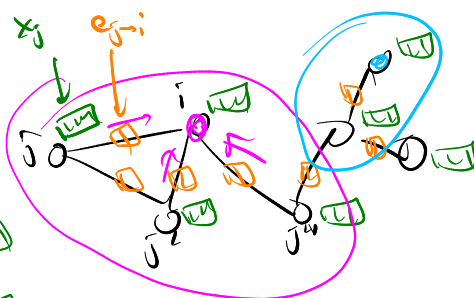
$$x \frac{A_{ij} - \mu}{\sigma} + b$$

parameters

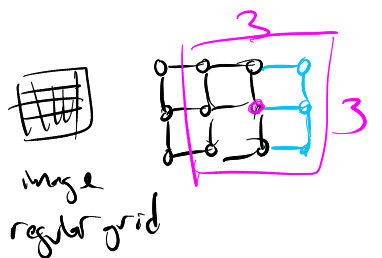
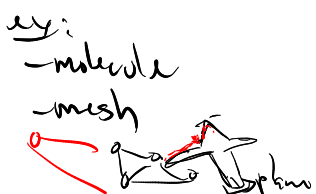
average or sum

Graph-NN

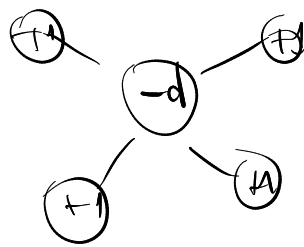
1 input = 1 graph



$$\left( \frac{1}{\sum_j w_j} \sum_j w_j x_j \right) + w_i x_i$$



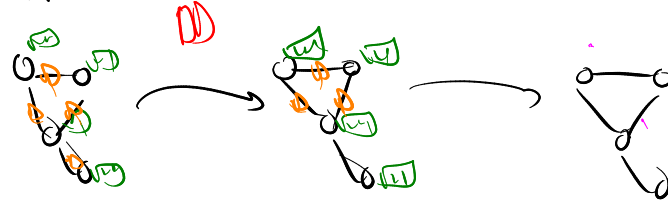
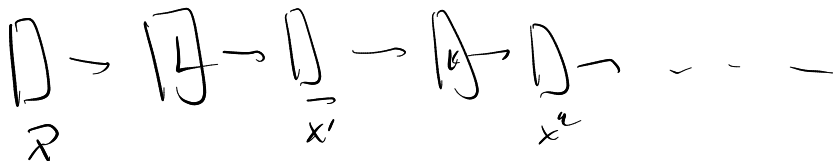
- social graph



$$g(e_{j \rightarrow i}, x_j, x_i)$$

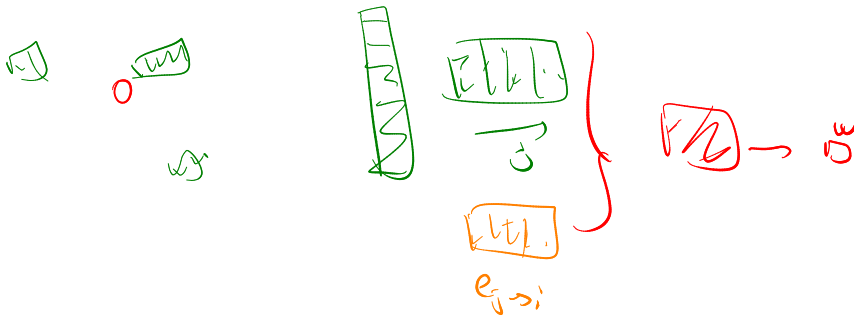
$$w_j = g(e_{j \rightarrow i})$$

graph attention networks  
GAT

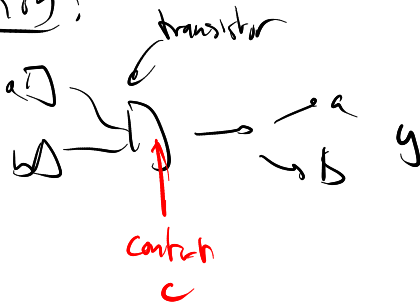


$$e_{j \rightarrow i} = h(e_{j \rightarrow i}, x_j, x_i)$$

average over graph  $\rightarrow$  2 output for full graph

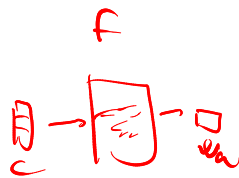


Attention:

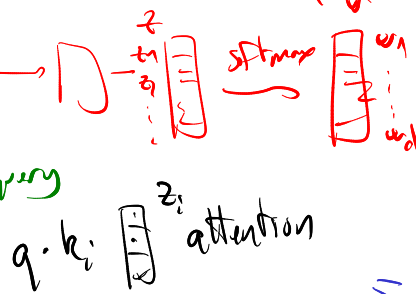
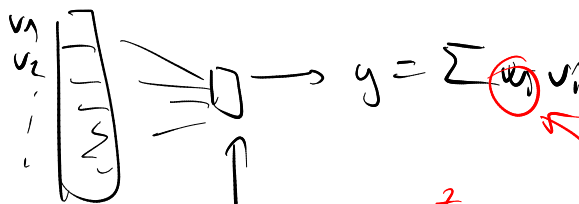


$$g = w_a a + w_b b = F(c)$$

$$\begin{cases} w_a = F(c) \\ w_b = 1 - w_a \end{cases}$$

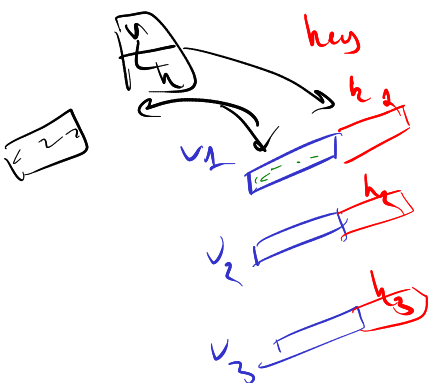


$$w_a + w_b = 1$$

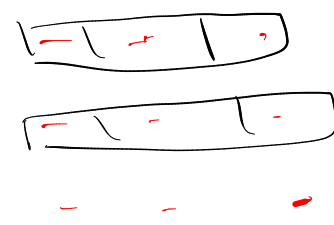
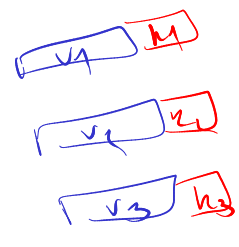


$$\begin{cases} w_i \geq 0 \\ \sum w_i = 1 \end{cases}$$

$$g = \sum_i w_i v_i$$



$h_1$   $h_2$   $h_3$

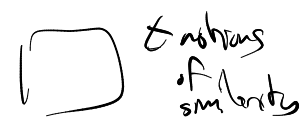


similarity matrix  $\rightarrow v_1^c = \sum_i w_i v_i$   
 $\rightarrow v_2^c = \sum_i w_i v_i$   
 $\rightarrow$

Set of vectors

Transformer

$h_k$   
 $v_k$

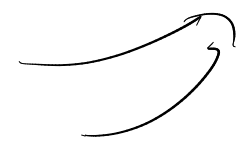


attention of similarity

$h_l$   
 $v_l$



multi-head



concatenate