

Neural Networks

A lot of hype!

- Very efficient
in practice

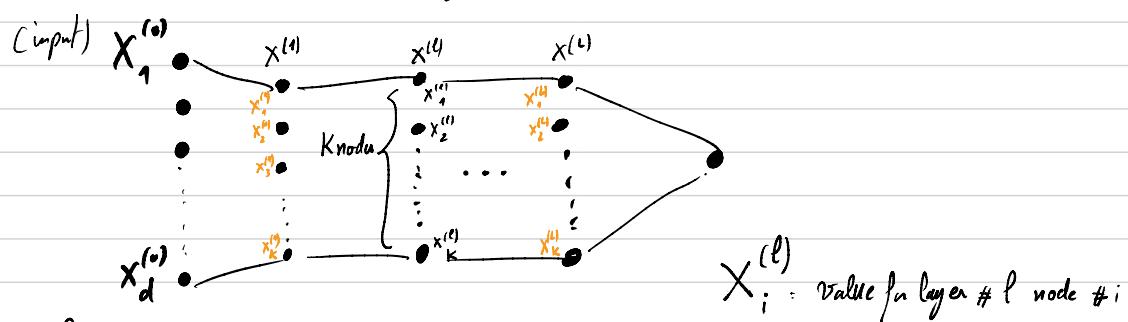
Neural Network?

- So far we have considered supervised learning (X, Y)
→ given a new X predict Y
- Linear methods work well when used with good features
 - Before: domain experts derive efficient features.
 - Now: Neural Net. learning from the data a good representation of them

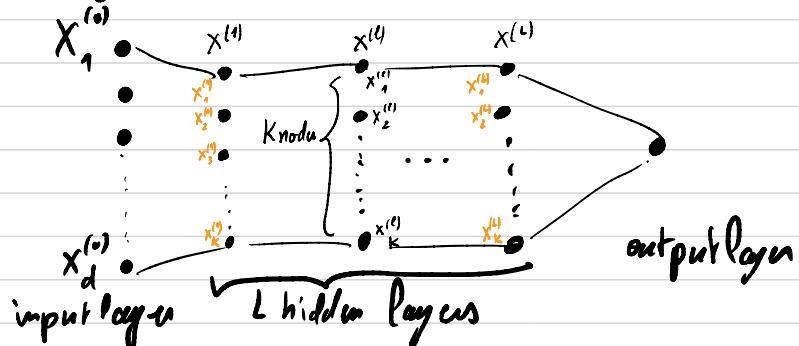
The Basic Structure

$$S_{\text{train}} = \{(x_i, y_i)\}_{i=1}^n, x_i \in \mathbb{R}^d$$

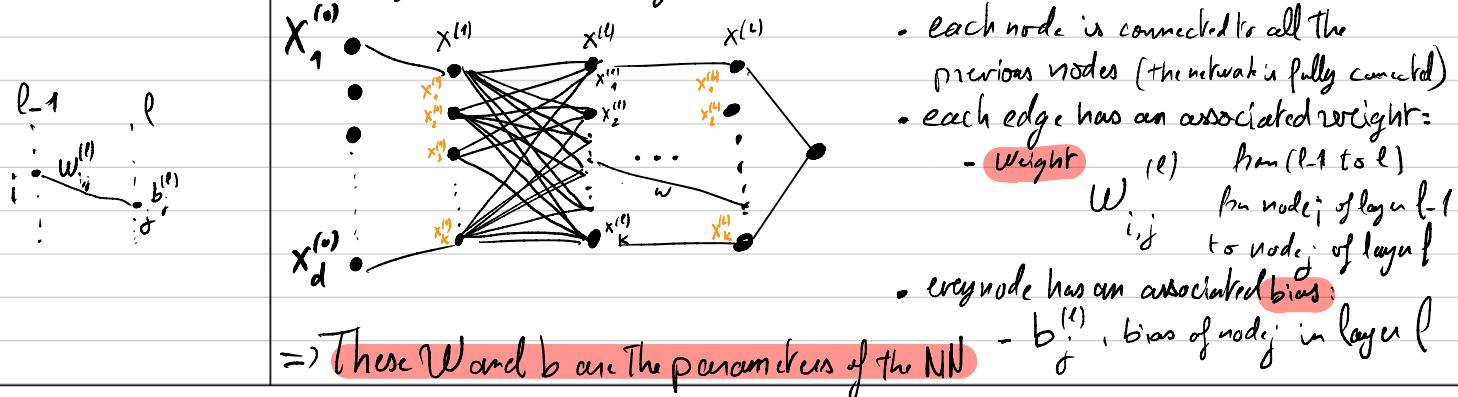
- A N-N with 2 hidden layers of size K and one output layer is:



- each layer has the same # of nodes K
- There is no feedback loop (= **feed forward Network**)



- The first hidden 2 layers represent some transformation of the data.



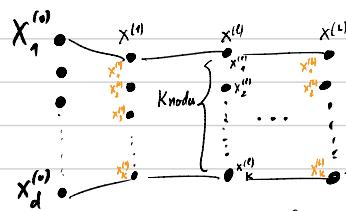
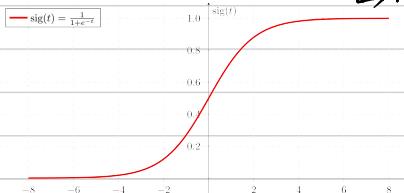
How to compute ϕ at the l^{th} layer:

- add at each layer some Non-linearity

$$x_j^{(l)} = \phi \left(\sum_{i=1}^k x_i^{(l-1)} w_{i,j}^{(l)} + b_j^{(l)} \right) \quad \text{with } \phi: \text{activation function}$$

i.e:

$$\phi(x) = \frac{1}{1+e^{-x}}$$



$$\Leftrightarrow f: \mathbb{R}^d \rightarrow \mathbb{R}^K$$

$$f(x) = \begin{pmatrix} \phi_1(x) \\ \vdots \\ \phi_K(x) \end{pmatrix}$$

The NN is representing a function from \mathbb{R}^d to \mathbb{R}^K
 → This function is characterized by:

- the par. (Weights and bias) you will learn
- the activation function ϕ you will pick

parameters of NN: $O(K^2 L)$

- if we assume that our data comes from some model:

$$y = f(x) + \epsilon$$

- is it possible to approx. f ?
- can we learn the function f ?

Approximation result of Andrew Barron:

let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ and $\tilde{f}(w) = \int_{\mathbb{R}^d} f(x) e^{-i w^T x} dx$ and assume that $\int |w| |\tilde{f}(w)| dw < \infty$

then for all $n \geq 1$, \exists fn s.t. $f_n(x) = \sum_{j=1}^n c_j \phi(x^T w_j + b_j) + C$

s.t.

$$\int_{|x| < r} (f(x) - f_n(x))^2 dx \leq \frac{(RCr)^2}{n}$$

⇒ (a function which is representable by a NN of 1 hidden layer and n nodes with ϕ act. fn.)

Interpretation:

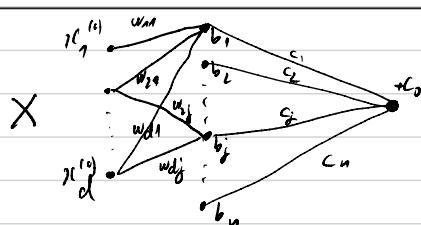
- the more nodes we allow, the smaller the error
- if $C \searrow$, the rate \searrow
- if $r \nearrow$, the rate \nearrow

- result is in ℓ_2 -norm

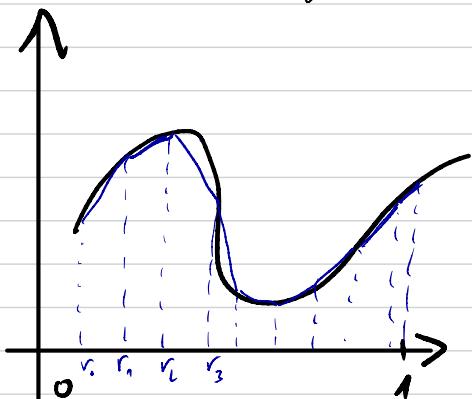
- f - any "sigmoid like" act. fn

"all suff. smooth fn can be approx. by a NN by 1 hidden layer with approx. error going down as $1/n$ "

$$f_n(x) = \sum_{j=1}^n c_j \phi(x^T w_j + b_j) + C$$



Is bound? f continuous in $[0, 1]$



- Polynomials: Stone Weierstrass Theorem
 $\forall \epsilon, \exists P, \sup_{x \in [0,1]} |f(x) - P(x)| < \epsilon$

- Piecewise linear fct

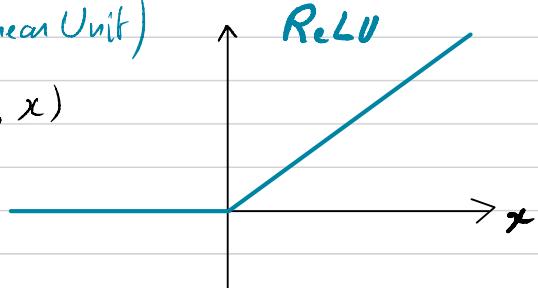
$$\Rightarrow f_n(x) = \sum_{i=1}^n (\alpha_i x + b_i) \mathbb{1}_{\{r_{i-1} \leq x < r_i\}}$$

with $\alpha_i, r_i, b_i = \alpha_{i+1}, r_{i+1}, b_{i+1}$

→ How can we approximate such fct with MN?

ReLU (activation fct) (Rectified Linear Unit)

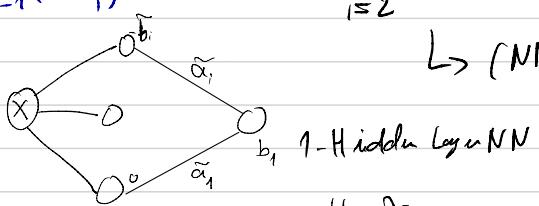
$$\text{ReLU}(x) = (x)_+ = \max(0, x)$$



now:

$$f_n(x) = \sum_{i=1}^n (\alpha_i x + b_i) \mathbb{1}_{\{r_{i-1} \leq x < r_i\}} = \tilde{\alpha}_1 x + \tilde{b}_1 + \sum_{i=2}^n \tilde{\alpha}_i (x - \tilde{b}_i)_+$$

↳ (NN ReLU Act.)

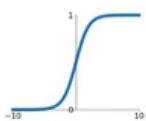


$$\text{with } \tilde{\alpha}_i = \alpha'_i - \alpha'_{i-1}, i \geq 1 \\ \tilde{b}_i = r_{i-1} \text{ and } \alpha'_i = \sum_{j=1}^i \tilde{\alpha}_j$$

Activation Functions

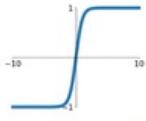
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



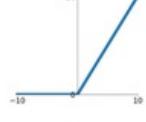
tanh

$$\tanh(x)$$



ReLU

$$\max(0, x)$$



Leaky ReLU
 $\max(0.1x, x)$

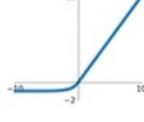


Maxout

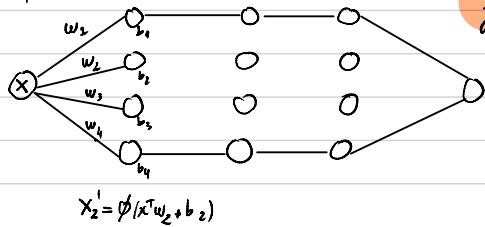
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



w_i vector

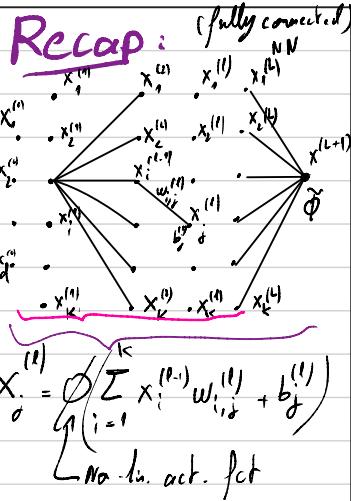


$$f_{w,b}(x) : \mathbb{R}^d \rightarrow \mathbb{R}^k$$

$$\text{for } f, \exists (w, b) \text{ s.t. } \|f - f_{w,b}\|_2 \leq \epsilon$$

$$\|f - f_{w,b}\|_\infty \leq \epsilon$$

Neural Networks training (SGD) & Backpropagation



$$f: \mathbb{R}^d \rightarrow \mathbb{R}^k$$

Repro. power:

- f is smooth
- Bounded domain
- Cond. on the act. fct
- A re. approx. in L_2 , but also worst case approx.

How to Train a NN:

- Training of NN: $S_{\text{Train}} = \{(x_i, y_i)\}_{i=1}^N$ for regression

$$\mathcal{L}_N(f) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2 \quad \text{where } f \text{ is the fct given by the NN}$$

Task: $\min \mathcal{L}_N$ ($O(K^2L)$ par.)

$$w_{ij}^{(1)}$$

$$b_j^{(1)}$$

Rq: \mathcal{L}_N is non convex.

Rq: we can consider some Reg. loss. $\tilde{\mathcal{L}}_N = \mathcal{L}_N + \text{Reg}$

- SGD alg.: Pick (x_i, y_i)
computing gradient of \mathcal{L}_i : $\mathcal{L}_i = \frac{1}{2} (y_i - f(x_i))^2$
(Take small step towards neg. grad. dir.
(stochastic))

• Problem: too many par. $O(K^2L) \Rightarrow$ Sol: Back propagation

NN: composition of functions:

$$x^{(1)} = f_1(x^{(0)}) = \phi \left[(w^{(1)})^T x^{(0)} + b^{(1)} \right]$$

weight matrix: $w^{(1)} \in \mathbb{R}^{d \times K}$ $w_{ij}^{(1)} = w_{ij}^{(1)}$ (weights)

$$x^{(2)} = f_2(x^{(1)}) = \phi \left[(w^{(2)})^T x^{(1)} + b^{(2)} \right]$$

$$w^{(2)} \in \mathbb{R}^{K \times K}$$

$$x^{(l)} = f_l(x^{(l-1)}) = \phi \left[(w^{(l)})^T x^{(l-1)} + b^{(l)} \right]$$

$$y = f_{L+1}(x^{(L)}) = \tilde{\phi} \left[\begin{matrix} w^{(L+1)} \\ \in \mathbb{R}^K \end{matrix} x^{(L)} + b^{(L+1)} \right]$$

$$f = f_{L+1} \circ f_L \circ f_{L-1} \circ \dots \circ f_2 \circ f_1$$

- complexity high
 \rightarrow 1st: Back propagation

Back propagation alg o:

- Forward pass complexity: $O(K^2L)$
- Backward pass

$$\nabla \mathcal{L}_n = \nabla (y_n - f(x_n))^2 \quad \text{Grad: } \frac{\partial \mathcal{L}_n}{\partial w_{ij}^{(l)}}, \frac{\partial \mathcal{L}_n}{\partial b_j^{(l)}}$$

Chain rule: $f(x) = g(h(x)) \rightarrow f'(x) = g'(h(x)) h'(x)$

- We need to compute all values of $X^{(l)}$ (forward pass)

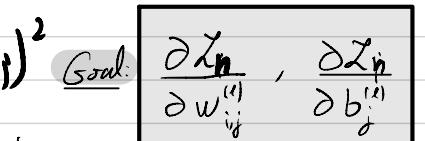
Forward Pass:

$$Z^{(l)} = (w^{(l)})^T X^{(l-1)} + b^{(l)}$$

$$X^{(l)} = \phi(Z^{(l)})$$

$$X^{(0)} = x; \in \mathbb{R}^d$$

Complexity: $O(K^2L)$



Backward Pass:

$$\delta_j^{(l)} = \frac{\partial \mathcal{L}_n}{\partial Z_j^{(l)}} = \sum_{k=1}^K \frac{\partial \mathcal{L}_n}{\partial Z_k^{(l+1)}} \frac{\partial Z_k^{(l+1)}}{\partial Z_j^{(l)}} \quad (l) \quad (l+1)$$

$$\delta_j^{(l)} = \sum_{k=1}^K \frac{\partial \mathcal{L}_n}{\partial Z_j^{(l+1)}} \cdot \delta_k^{(l+1)} \quad j \quad k$$

We have to compute: $\frac{\partial Z_k^{(l+1)}}{\partial Z_j^{(l)}}$

$$\hookrightarrow Z_k^{(l+1)} = \sum_{i=1}^K w_{ik}^{(l+1)} x_i^{(l)} + b_k^{(l+1)} = \sum_{i=1}^K w_{ik}^{(l+1)} \phi(Z_i^{(l)}) + b_k^{(l+1)}$$

$$\Rightarrow \frac{\partial Z_k^{(l+1)}}{\partial Z_j^{(l)}} = \phi'(Z_j^{(l)}) w_{jk}^{(l+1)}$$

$$\Rightarrow \delta_j^{(l)} = \sum_{k=1}^K \phi'(Z_j^{(l)}) w_{jk}^{(l+1)} \cdot \delta_k^{(l+1)} \Rightarrow \boxed{\delta^{(l)} = w^{(l+1)} \delta^{(l+1)} \circ \phi'(Z^{(l)})}$$

$$\delta^{(L+1)} = \frac{\partial}{\partial Z^{(L+1)}} (y_n - \tilde{\phi}(Z^{(L+1)}))^2$$

⊗: pointwise product
(scalar like)

$$\delta^{(L+1)} = 2 \tilde{\phi}'(Z^{(L+1)}) (\tilde{\phi}(Z^{(L+1)}) - y_n) \quad (\text{now we go backward and compute all } \delta^{(l)})$$

Now

$$\bullet Z^{(l)} = (w^{(l)})^T X^{(l-1)} + b^{(l)} ; \Rightarrow Z_j^{(l)} = \sum_{i=1}^K w_{ij}^{(l)} x_i^{(l-1)} + b_j^{(l)}$$

$$\frac{\partial \mathcal{L}_n}{\partial b_j^{(l)}} = \frac{\partial \mathcal{L}_n}{\partial Z_j^{(l)}} \frac{\partial Z_j^{(l)}}{\partial b_j^{(l)}} = \frac{\partial \mathcal{L}_n}{\partial Z_j^{(l)}} ; \quad \frac{\partial \mathcal{L}_n}{\partial w_{ij}^{(l)}} = \frac{\partial \mathcal{L}_n}{\partial Z_j^{(l)}} \cdot \frac{\partial Z_j^{(l)}}{\partial w_{ij}^{(l)}} = x_i^{(l-1)} \cdot \frac{\partial \mathcal{L}_n}{\partial Z_j^{(l)}} = \delta_j^{(l)} \quad \delta_j^{(l)}$$

Summary:

- Forward pass: $X^{(0)} = x_i$ (one sample point)
computing $Z^{(l)} = (w^{(l)})^T X^{(l-1)} + b^{(l)}$
 $X^{(l)} = \phi(Z^{(l)})$
= computing all values

Backward pass:

$$\delta^{(l)} = w^{(l+1)} \delta^{(l+1)} \circ \phi'(Z^{(l)})$$

$$\delta^{(L+1)}$$

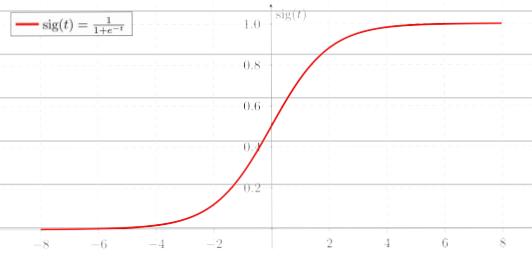
Compute:

$$\frac{\partial \mathcal{L}_n}{\partial w_{ij}^{(l)}} = \delta_j^{(l+1)} x_i^{(l-1)}$$

$$\frac{\partial \mathcal{L}_n}{\partial b_j^{(l)}} = \delta_j^{(l+1)}$$

Activation functions:

Sigmoid



$$\sigma(x) = \frac{1}{1+e^{-x}}$$

⊕ smooth everywhere

⊖ $|\sigma'(x)| \lll 1$
for $|x| \ggg 1$
→ vanishing gradient ≈ 0

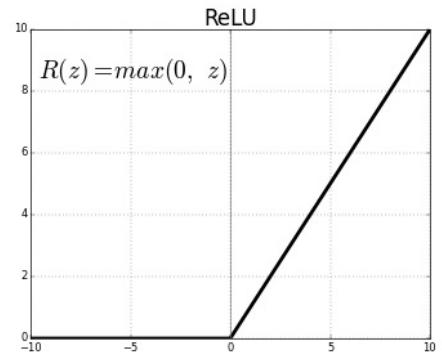
ReLU

$$\text{ReLU}(x) = (x)_+ = \max(0, x)$$

⊕ for $x > 0$, No vanishing gradient

⊖ Not diff. in 0
• The deriv. is 0 for $x < 0$

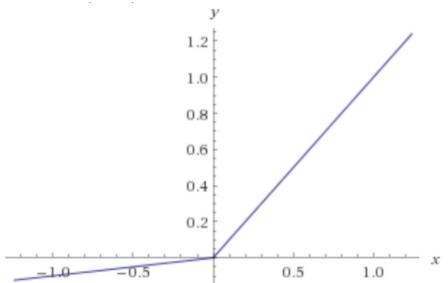
Sol → Leaky ReLU



Leaky ReLU

$$\text{LReLU} = \max(ax, x)$$

$a \ll 1$ small



Max Out

$$f(x) = \max(x^T w_1 + b_1, \dots, x^T w_k + b_k)$$

$$\text{Tanh} \quad \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 2\sigma(x) - 1$$

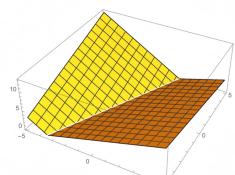
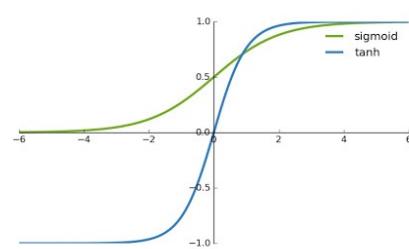


Figure 5: Maxout function with two terms, $\max\{x_1 - 0.5x_2 + 1, -2x_1 + x_2 - 2\}$.

Regularization, Data Augmentation, drop-out & Convolutional Nets

Regularization

$$L(w) + \frac{\lambda}{2} \|w\|_2^2 ; \text{For NN: It is interesting to Regular. the weight but not the bias}$$

- $\nabla_w L + \lambda w \Rightarrow w_{k+1} = w_k - \gamma [\nabla_w L + \lambda w_k]$

$$w_{k+1} = (1 - \gamma \lambda)(w_k) - \gamma \nabla_w L(w_k)$$

\Rightarrow It is called weight decay

- We can also the constrained formulation: $\min_{\|w\|_2 \leq r} L(w)$

$$\Rightarrow \text{can be minimized using Projected GD.}$$

- Both formulations are very close.

(Sometimes equivalent for linear fits)

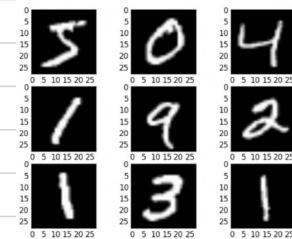
\rightarrow Penalization is simpler to implement

\rightarrow the constraint form is easier to interpret.

\Rightarrow Regularization can be a good idea, but only on the weights for NN.

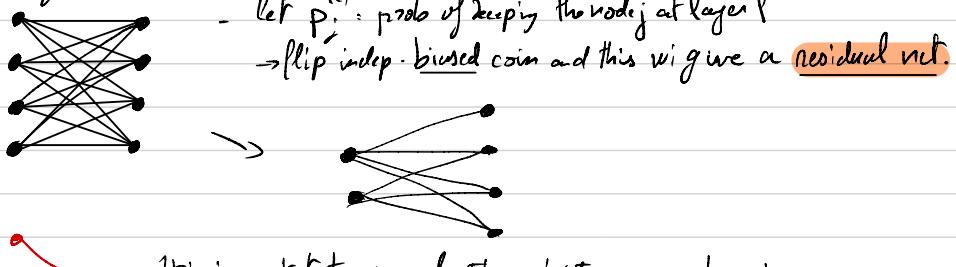
Data augmentation

i.e (MNIST data), $S_{\text{Train}} = \{(x_i, y_i)\}$



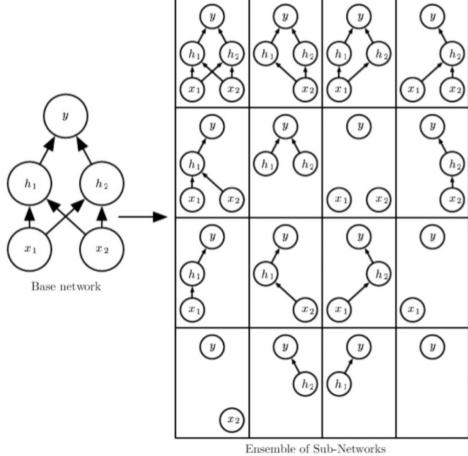
Dropout:

- Experimentally it helps avoid overfitting
- It is also a technique of ensemble averaging
- . Original version on the node



Why?

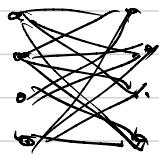
- It gives a more robust way of learning
- It is very close to train w/ many NNs and predict in avg. (ensemble averaging)



(bagging)
 ↳ variance of pred

Convolutional Network

So far, we have considered fully connected NN.



$\mathcal{O}(K^2L)$ parameters

(do we really need all these edges?)

→ Is it possible to get better perf. with a sparse set of edges?

Convolution:

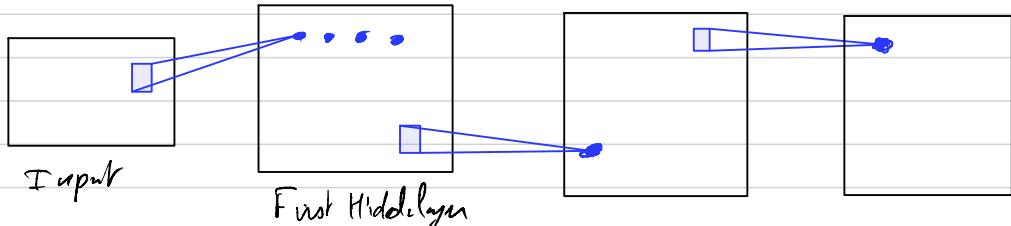
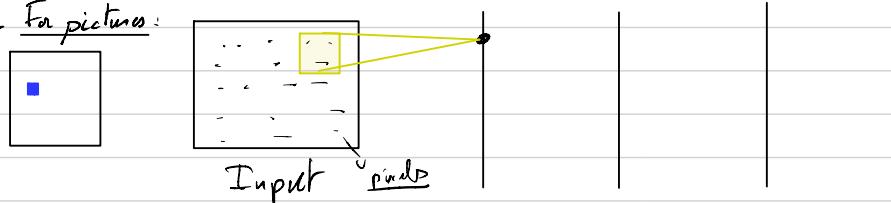
$$h = f * g$$

→ by taking the conv. we are doing local averaging

$$h[n] = \sum_k [f[n-k] \cdot g[n-k]]$$

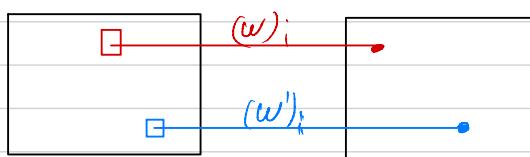
↑ filters, kernel

- For pictures:



\square_p of size p - local averaging window

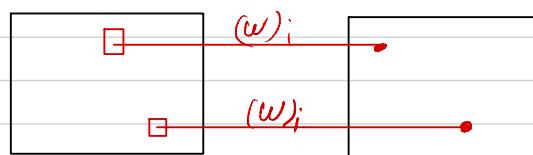
$\mathcal{O}(K^2L)$ parameters to $\mathcal{O}(Kp^2)$ para.



different filters.

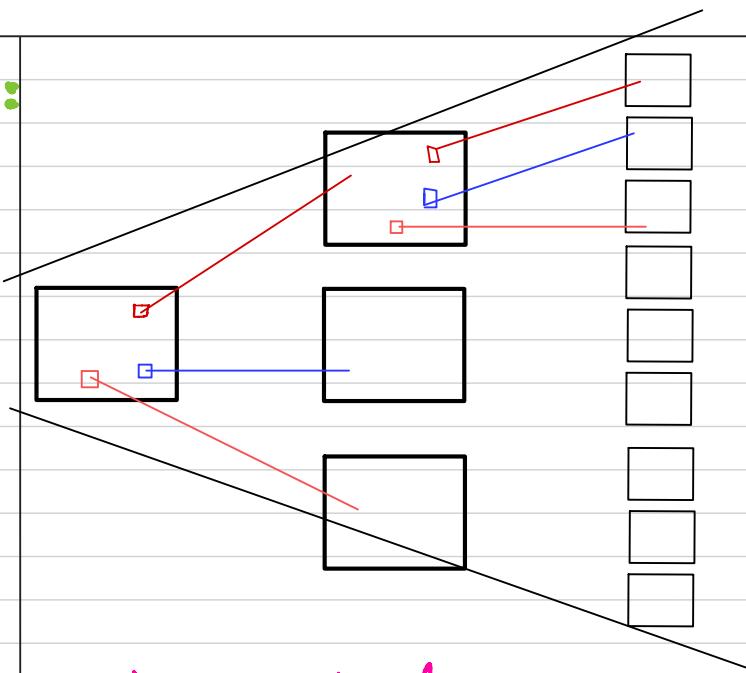
→ often we would use the same filter on all regions.

Weight Sharing:



$\mathcal{O}(PL)$ parameters

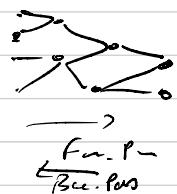
Channels:



Pyramidal form

Back propagation only..

- You can do Back prop. whatever the architecture.



- With weight sharing?



$$f(x, y, z); \nabla f = \begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \\ \frac{\partial f}{\partial z} \end{pmatrix} \in \mathbb{R}^3$$

"weight shring":

$$g(x, y) = f(x, y, x) \text{ How to compute } \nabla g?; \nabla g = \left(\begin{array}{c} \frac{\partial g}{\partial x} \\ \frac{\partial g}{\partial y} \end{array} \right)$$

$$\begin{aligned} \frac{\partial g}{\partial x} &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial x} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial x} + \frac{\partial f}{\partial x_3} \frac{\partial x_3}{\partial x} \\ &= \frac{\partial f}{\partial x_1} + \frac{\partial f}{\partial x_3} \end{aligned}$$

$$\boxed{\nabla g = \left(\begin{array}{c} \frac{\partial f}{\partial x_1} + \frac{\partial f}{\partial x_3} \\ \frac{\partial f}{\partial y} \end{array} \right)}$$

cond: We run Back prop with weight shring,

and then sum the derivatives for the weights you shring.