

# APPLYING MACHINE LEARNING TO STRING THEORY

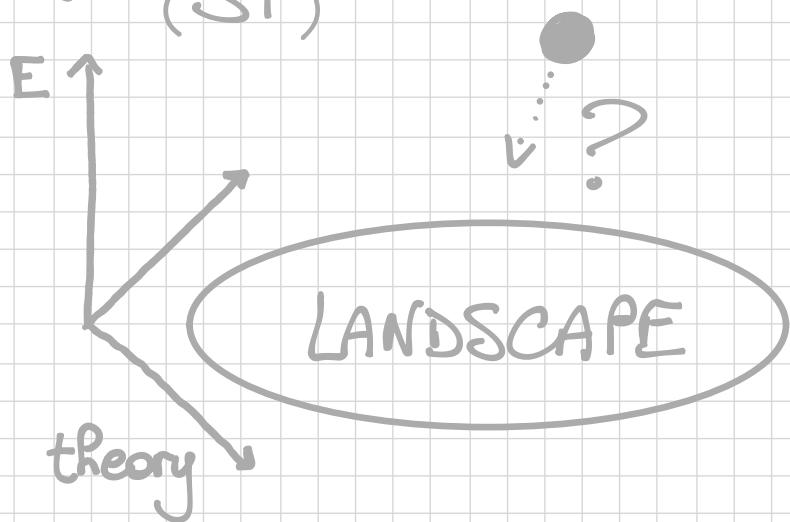
## 1. MOTIVATIONS

Complex analytical problem  $\xleftarrow{\text{"more feasible"}}$  optimisation problem  $\Rightarrow$  Insights ?

### 1.1 MACHINE LEARNING FOR PHYSICS

#### PATTERN EXPLORATION

e.g. 1) String Theory:  $\mathbb{R}^{1,9} \rightarrow \mathbb{R}^{1,3} \times X_6$   
(ST)



How many possibilities?

F-theory:  $10^{272\,000}$

[Taylor, Wang, 15.11.2019]

How difficult?

NP-Complete [Denef, Douglas,  
hep-th/0602072]  
NP-Hard  
undecidable [Halverson, Ruehle,  
1809.08279]

Non deterministic Searches?

[Abel, Rizos, 1404.7359]

[Ruehle, 1706.07024]

[He, 1706.02714]

[Altman et al., 1811.06490]

---

[Ruehle, Phys. Rep. 839 (2020)]

e.g. 2) algebraic geometry and ST:

TOPOLOGICAL  
SPACE  $\xrightarrow{?}$  TOPOLOGICAL  
INVARIANTS

[2007.15706]

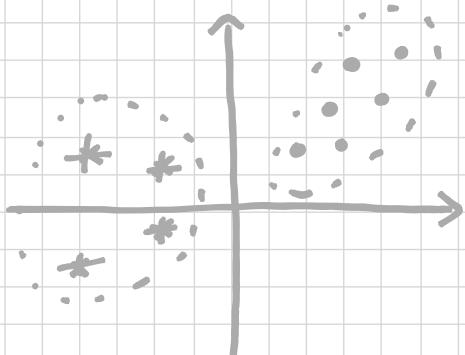
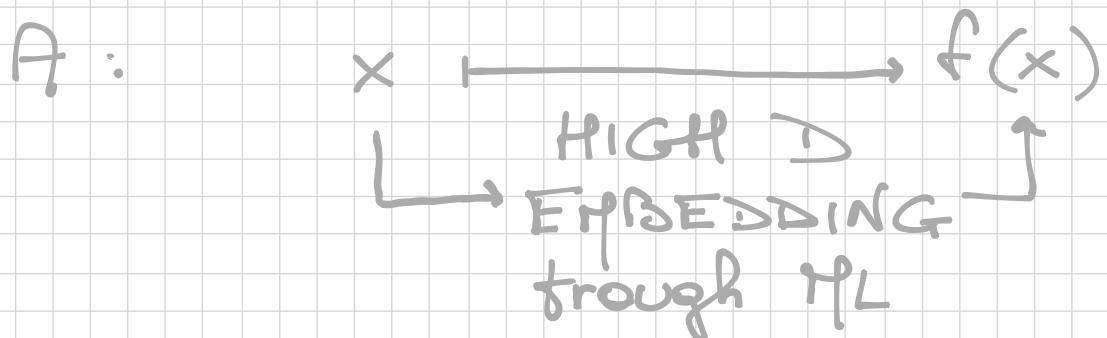
[2108.02221]

### e.g. 3) Symmetry detection

$$S : A \rightarrow A \\ x \mapsto S(x)$$

$$f : A \rightarrow B \\ z \mapsto f(z)$$

Q: how to detect  $x \in A \mid f(S(x)) = f(x)$ ?



$\Rightarrow$  find clusters  
of symmetric objects  
[Krippendorff,  
Syraeri, 2003. 13679]

## 1.2 PHYSICS FOR MACHINE LEARNING

EXPLANATION and  
BEHAVIOURAL ANALYSIS

e.g. 1.1) NN - QFT duality

Let  $f_{\theta, N} : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$

in the limit  $N \rightarrow \infty$ :

$$P[f] \sim e^{-S_{GP}}$$

where

$$S_{GP} = \frac{1}{2} \int d^{d_{in}}x \int d^{d_{in}}x' f(x) \sum(x, x') f(x')$$

is a GAUSSIAN PROCESS:

$$G^{(k)}(x_1, \dots, x_k) = \frac{\int \partial f f(x_1) \dots f(x_k) e^{-S_{GP}}}{\int \partial f e^{-S_{GP}}}$$

KEY POINT: propagator  $K(x, x')$

$$\int d^{d_{in}}x' \underbrace{K(x, x') \sum(x', x'')}_{\text{COVARIANCE}} = \delta^{d_{in}}(x - x'')$$

strongly dep.  
of the output ] on type of NN

What for  $N < \infty$ ?

$$S_{GP} \rightarrow S = S_{GP} + \text{INTERACTIONS}$$

[Halverson, Maiti, Stoner, 2008.08601]

[Maiti, Stoner, Halverson, 2106.00694]

[Di Luò, Halverson, 2112.00723]

$$\text{e.g. 2) } S = S_{\text{eff}} + \boxed{\Delta S}$$

↓

coarse graining

↓

RG

Non pert. approaches show links with NN behaviour.

[Lahoche, Samary, Tamaazousti, 2002.1054]

[Erbin, Lahoche, Samary, 2108.0403]

### 1.3 MACHINE LEARNING FOR "ANALYTICAL METHODS"

e.g. 1) study F-theory compactif.

ML { INPUT : reflexive polytopes  
OUTPUT : Fano toric 3-folds

↳ conjecture on  $E_6$  compactif.

↳ analytical proof

[Cerifio et al., 1707.00655]

e.g. 2) triangulations of reflexive polytopes

ML { INPUT: 4D reflexive polytopes  
OUTPUT: triangulations

↳ estimate no. of triang. in the KS dataset → proof

[Altman et al., 1811.06490]

e.g. 3) analytic rules from 4D F-theory

ML { INPUT: divisors of complex 3-folds base  
OUTPUT: non-Higgsable gauge groups

↳ analytic rules → proof

[Wang, Zhang, 1804.07296]

e.g. 4) generation of superpotentials

ML { INPUT: Known Superp.  
OUTPUT: unknown Superp.

[Erbin, Krippendorf, 1809.02612]

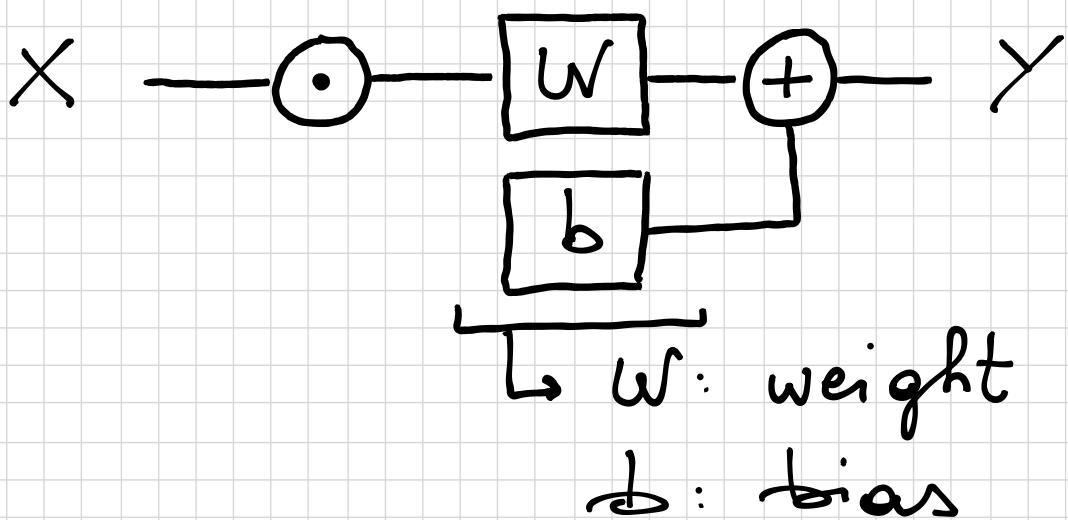
## 2. NEURAL NETWORKS

- impact
- building blocks
- how to use them in physics

### 2.1 PERCEPTRON

First conceived by Rosenblatt, 1958  
(Cornell Aeronautical Lab.):

- i) pattern recognition
- ii) first HARDWARE  
(software implementation came later)
- iii) provide linear classification



In formulae:

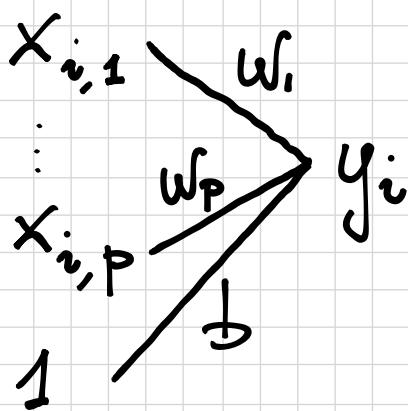
Let  $x_i \in \mathbb{R}^p$ ,  $i = 1, \dots, N$  and

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_N^T \end{pmatrix} \in \mathbb{R}^{N \times p}$$

then

$$P(X) = Xw + bI_N \in \mathbb{R}^N$$

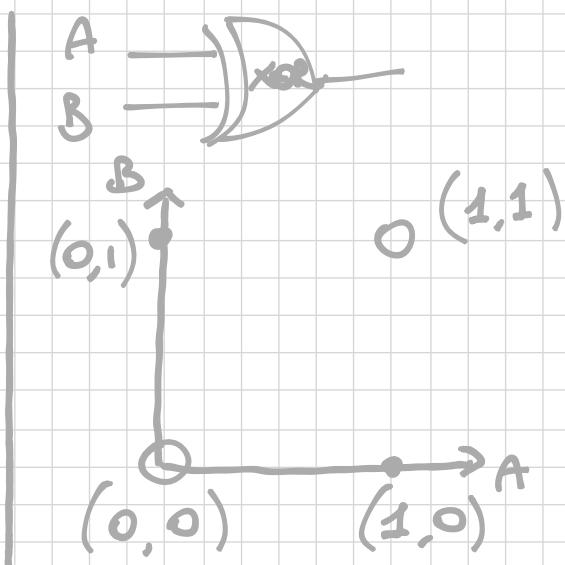
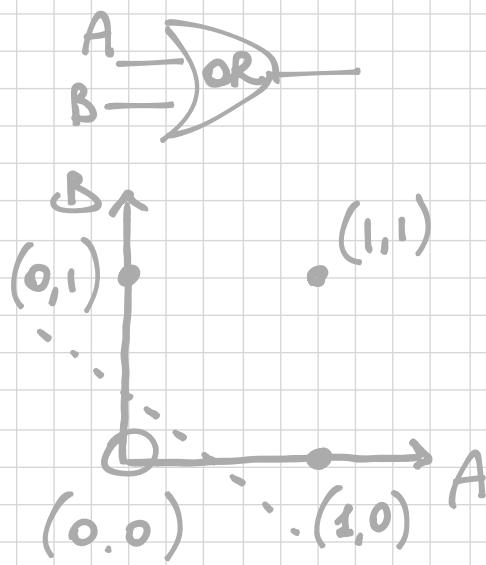
that is



N.B.: LINEAR ALGORITHM

[Historically, it lead to abandon NN research for ~20 years]

e.g.: LOGIC GATES

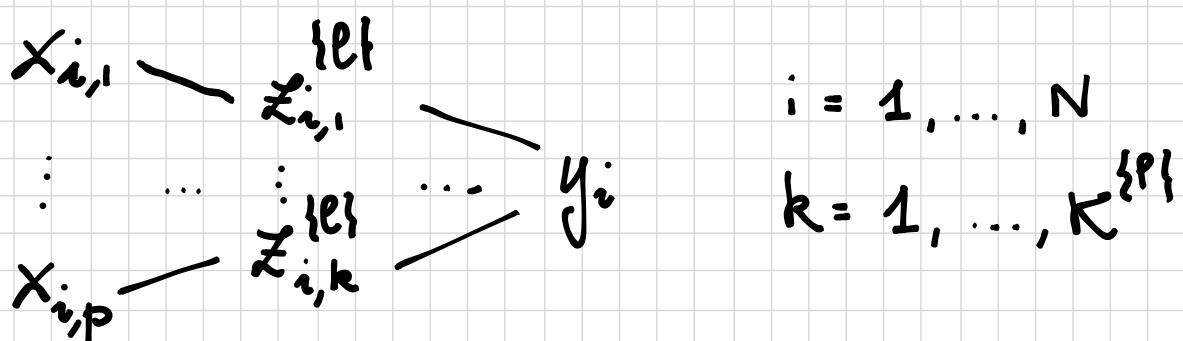


## 2.2 MULTI-LAYERED PERCEPTRON AND FULLY CONNECTED NETWORKS

New function

$$Y = \prod_{l=1}^L \left\{ (\omega^{lp}, b^{lp}) \right\} (x)$$

where :



(the bias term is omitted for simplicity)

In formulae,

$$\text{def: } \omega^{lp} \in \mathbb{R}^{K^{lp-1} \times K^{lp}}$$

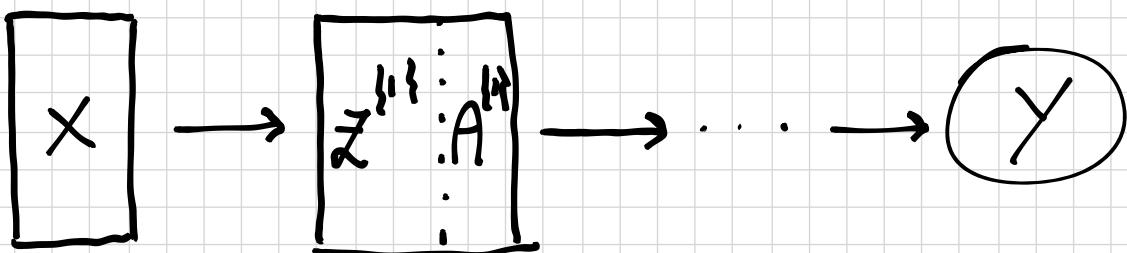
$$\begin{aligned} \text{then: } Y &= ((X \omega^{(1)}) \omega^{(2)}) \omega^{(3)} \dots \\ &= X \underbrace{\tilde{\omega}}_{\text{LINEAR}} \end{aligned}$$

New ingredient:

$$Y = \text{MLP}_{\{(w^l, b^l, f^l)\}_{l \in \{1, \dots, L\}}} (X)$$

for each layer :

$$\begin{cases} z^l = A^{l-1} w^l \\ A^l = f^l(z^l) \\ A^0 = X \end{cases}$$

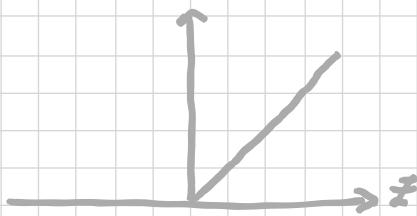


"ACTIVATION FUNCTIONS"

- non linear
- differentiable

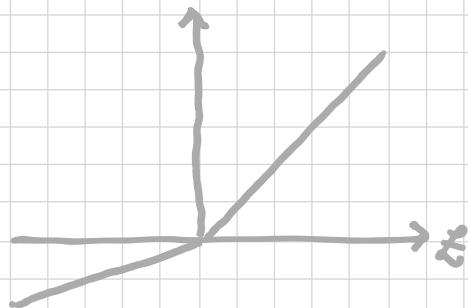
e.g. 1) RECTIFIED LINEAR UNIT

$$\text{ReLU}(z) = \max(0, z)$$

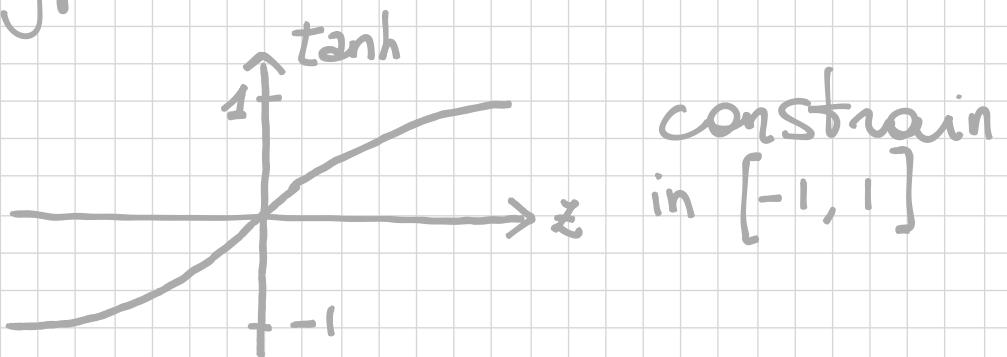


e.g. 2) "LEAKY"-RELU

$$\text{LeakyReLU}_\alpha(z) = \max(\alpha z, z), \alpha > 0$$

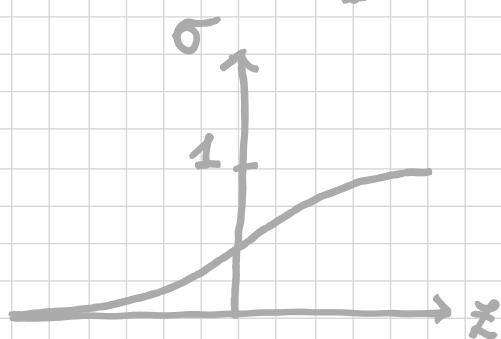


e.g. 3) Hyp. Tan.



e.g. 4) Sigmoid (CLASSIFICATION)

$$\sigma(z) = \frac{1}{1 + e^{-z}}, z \in \mathbb{R}^N$$



model probabilities

let  $p \in [0, 1)$  be the prob. of an event:

$$\text{logit}(p) = \ln\left(\frac{p}{1-p}\right) = \sigma^{-1}(p).$$

e.g. 4) Softmax

$$S(\mathbf{z})_{i,k} = \frac{e^{z_{i,k}}}{\sum_{k=1}^K e^{z_{i,k}}}$$

where  $\mathbf{z} \in \mathbb{R}^{N \times K}$ ,  $k = 1, \dots, K$ .

Used in MULTICLASS CLASS. :

- one class

$$y_i \in \{0, 1\}$$

- multiclass

$$y_i \in \{0, 1, \dots, K\}$$

Introduce ONE-HOT ENCODING

$$y_i \rightarrow \tilde{y}_{i,k} = \begin{cases} 1 & \text{if } y_i = k \\ 0 & \text{otherwise} \end{cases}$$

finally  $\tilde{\mathbf{y}}_i \in \mathbb{R}^K$ :

$$\mathbf{Y} = \begin{pmatrix} 0 \\ 3 \\ 1 \end{pmatrix} \rightarrow \tilde{\mathbf{Y}} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

## 2.4 Loss Functions

So far : FORWARD PASS

[ Input  $\rightarrow$  NN  $\rightarrow$  Output ]

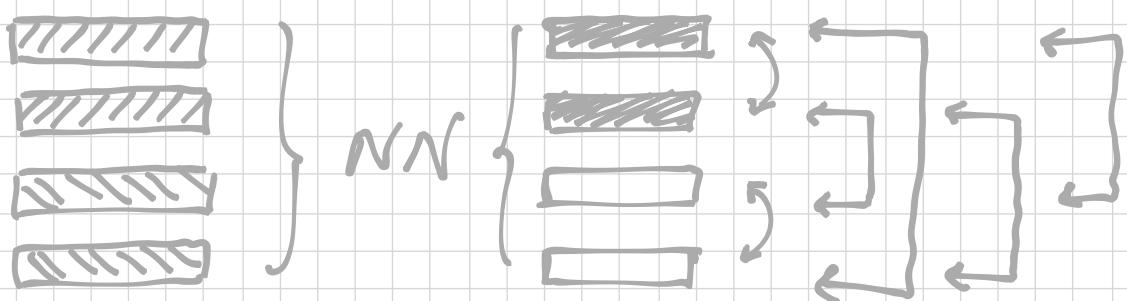
How to guide the "learning"?

CASE 1 : you do not know the true result (the "ground truth")

$\rightarrow$  UNSUPERVISED LEARNING

(tricky for neural networks)

e.g.



[Chen et al., 2002.05#09]

[Kochla et al., 2004.11362]

"CONTRASTIVE LEARNING"

CASE 2 : you know the answer  
 $\Rightarrow$  SUPERVISED LEARNING

- define a "distance" from the ground truth:

$$L(Y, \hat{Y}_\theta) = L(\theta)$$

where  $\hat{Y}_\theta = f_\theta(x) \leftarrow \text{NN}_\theta(x)$ .

- criterion based on LIKELIHOOD :

$$P(\hat{Y}_\theta | \theta) := \mathcal{L}(\theta | \hat{Y}_\theta)$$

$\Rightarrow$  joint prob. of the outp. given the model param. SEEN AS A FUNCTION OF THE PARAMS! (NOT a PDF!)

$\Rightarrow$  "which params are more likely to generate that output?"

$$\hat{\theta} := \arg \max_{\theta} \mathcal{L}(\theta | \hat{Y}_\theta)$$

MAXIMUM LIKELIHOOD ESTIMATE

or

$$\hat{\theta} := \arg \min_{\theta} L(\theta), \quad L(\theta) := -\log \mathcal{L}(\theta | \hat{Y}_\theta)$$

# a) CLASSIFICATION

$$y_i \in \mathbb{R}^K \quad (\text{OH enc.}) \\ i = 1, \dots, N$$

then

$$\mathcal{L}(\theta | \hat{y}_{\theta, i}) = \prod_{k=1}^K \hat{y}_{\theta, i, k}^{y_{i,k}}$$

(only the observation for which  
 $y_{i,k} = 1$  appears)

e.g. BINARY CLASS.  $K=2$

$$\begin{cases} y_0 + y_1 = 1 \\ \hat{y}_0 + \hat{y}_1 = 1 \end{cases} \rightarrow \mathcal{L}(\theta | \hat{y}_\theta) = \hat{y}_0^{y_0} (1 - \hat{y}_0)^{1-y_0}$$

then the NLL or "Loss function":

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \left( -\log \mathcal{L}(\theta | \hat{y}_{\theta, i}) \right)$$

$$= -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{i,k} \log \hat{y}_{\theta, i, k}$$

$$= -\mathbb{E}_y [\log \hat{\gamma}_\theta] := \text{CROSS-ENTROPY}$$

## b) REGRESSION

$$y_i \in \mathbb{R} \quad i = 1, \dots, N$$

s.t.  $y_i = \hat{y}_{\theta,i} + \varepsilon_i$

where  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

"no depend. on the  
indep. var."

Then

$$\mathcal{L}(\theta | \hat{y}_{\theta,i}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \hat{y}_{\theta,i})^2}{2\sigma^2}}$$

$$\Rightarrow L(\theta) = -\frac{1}{N} \sum_{i=1}^N \left\{ -\frac{(y_i - \hat{y}_{\theta,i})^2}{2\sigma^2} - \frac{1}{2} \log 2\pi\sigma^2 \right\}$$

$$\Rightarrow \hat{\theta} = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_{\theta,i})^2$$

"MEAN SQUARED ERROR"

Other more general examples:

i) MEAN ABSOLUTE ERROR

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_{\theta,i}|$$

ii) HUBER LOSS

$$L_\delta(\theta) = H_\delta(\theta) = \frac{1}{N} \sum_{i=1}^N \begin{cases} \frac{1}{2}(y_i - \hat{y}_{\theta,i})^2 & \text{if } |y_i - \hat{y}_{\theta,i}| < \delta \\ \delta(|y_i - \hat{y}_{\theta,i}| - \frac{\delta}{2}) & \text{otherwise} \end{cases}$$

iii) KULLBACK-LEIBLER DIV.

- let  $P(x)$  be the cumulative distrib of an input  $x$
- $Q$ : how similar are two distrib.  
 $P(z|x)$  and  $Q(z)$  ?

$$\text{EVIDENCE} := \log P(x)$$

$$= \log \int dz P(x|z)P(z)$$

$$= \log \int dz \frac{Q(z)}{P(z)} P(x,z)$$

$$= \log \mathbb{E}_{z \sim Q} \left[ \frac{P(x,z)}{Q(z)} \right]$$

$$\text{Jensen's inequality} \geq \mathbb{E}_{z \sim Q} \left[ \log \frac{P(x, z)}{Q(z)} \right]$$

$\therefore$  EVIDENCE LOWER Bound

In other words (let  $\theta$  be params):

$$\text{ELBO} := \mathbb{E}_{z \sim Q} [\log L(\theta | x, z)] - \mathbb{E}_{z \sim Q} [Q(z)]$$

Now introduce the KL div.:

$$\begin{aligned} \text{KL}(Q(z) \| P(z|x)) &:= \\ &:= \mathbb{E}_{z \sim Q} \left[ \log \frac{Q(z)}{P(z|x)} \right] \end{aligned}$$

$$= \mathbb{E}_{z \sim Q} [\log Q(z)] - \mathbb{E}_{z \sim Q} [\log \frac{P(x, z)}{P(x)}]$$

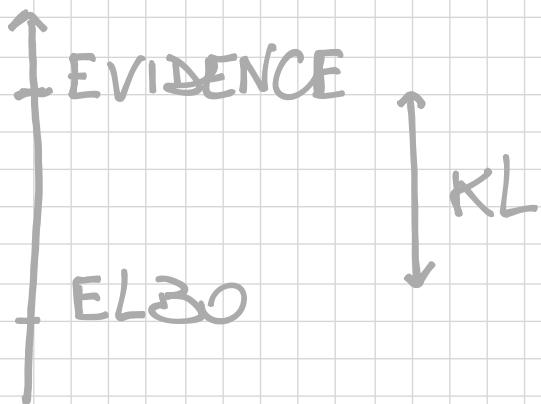
With params  $\theta$ :

$$\text{KL}(Q(z) \| P(z|x)) = \log L(\theta | x)$$

$$- \left[ \mathbb{E}_{z \sim Q} [\log L(\theta | x, z)] - \mathbb{E}_{z \sim Q} [\log Q(z)] \right]$$

$$\therefore \text{EVIDENCE} - \text{ELBO}$$

Visually



at the heart of (Bayes) VARIATIONAL INFERENCE:

: if  $Q(z)$  is impossible to compute  
we use  $P(z|X)$  to approximate:

$$\hat{Q} = \arg \max_Q \text{ELBO}(P, Q)$$

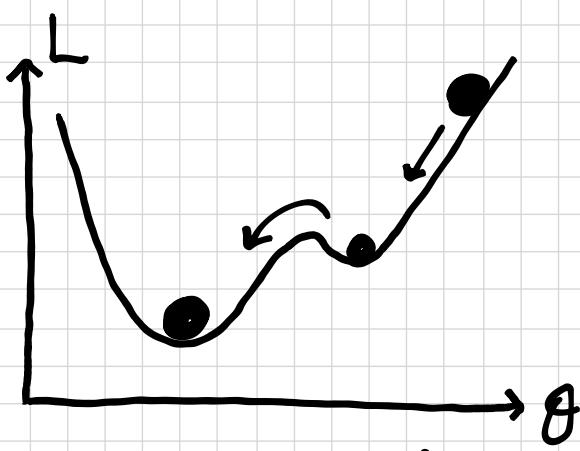
$$\Rightarrow \hat{Q} = \arg \min_Q \text{KL}(Q \| P).$$

## 2.5 Backpropagation

So far

FORWARD PASS → LOSS FUNCTION

Where is training?



$$\hat{\theta} := \underset{\theta}{\operatorname{argmin}} L(\theta)$$

"GRADIENT DESCENT"

i)  $\dot{\theta} = -\nabla_{\theta} L(\theta)$

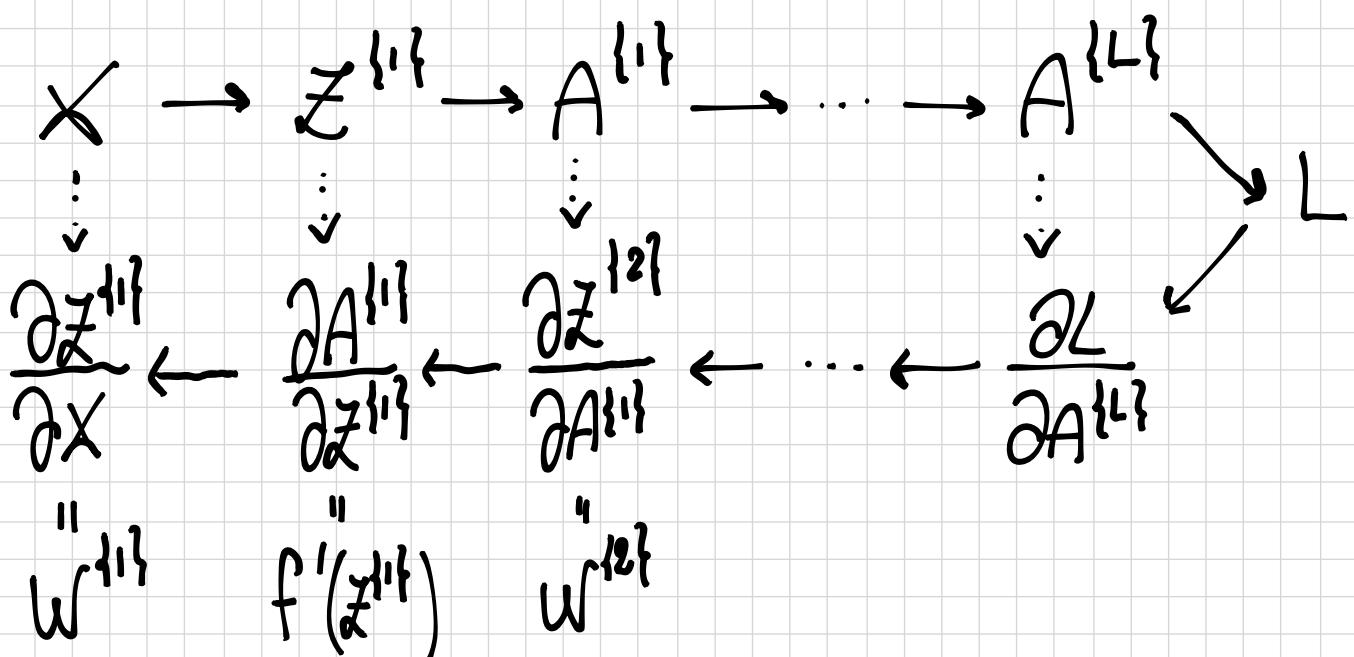
ii) let  $\alpha > 0$  "LEARNING RATE"

iii) update  $\theta \leftarrow \theta + \alpha \dot{\theta}$

How to do this efficiently?

BACKPROPAGATION

[Rumelhart, Hinton, Williams,  
Nature 323 (1986)]



Then let

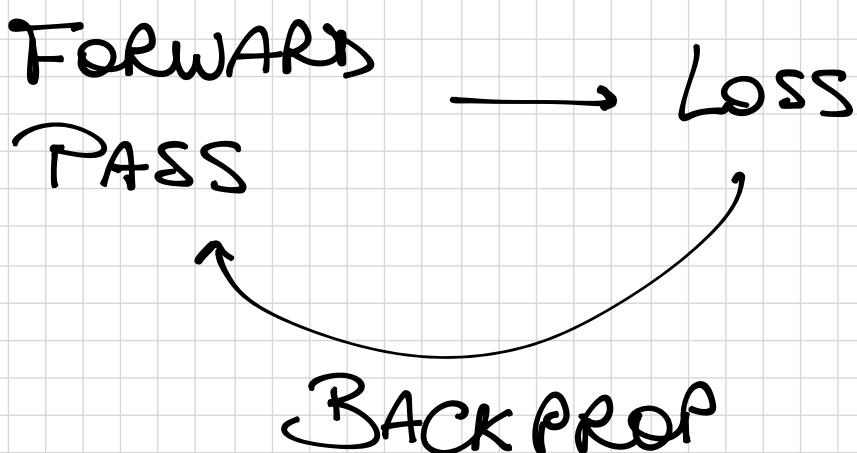
$$g^{\{l\}} = f'(z^{\{l\}}) w^{\{l+1\}} \dots f'(z^{\{L\}}) \nabla_A L$$

and

$$\begin{aligned} \frac{\partial L}{\partial w^{\{l\}}} &= \frac{\partial z^{\{p\}}}{\partial w^{\{l\}}} \cdot \frac{\partial L}{\partial z^{\{l\}}} \\ &= A^{\{p-1\}} \cdot g^{\{l\}} \end{aligned}$$

N.B.: frameworks such as PyTorch or Tensorflow are auto-grad. Systems!

Thus we have

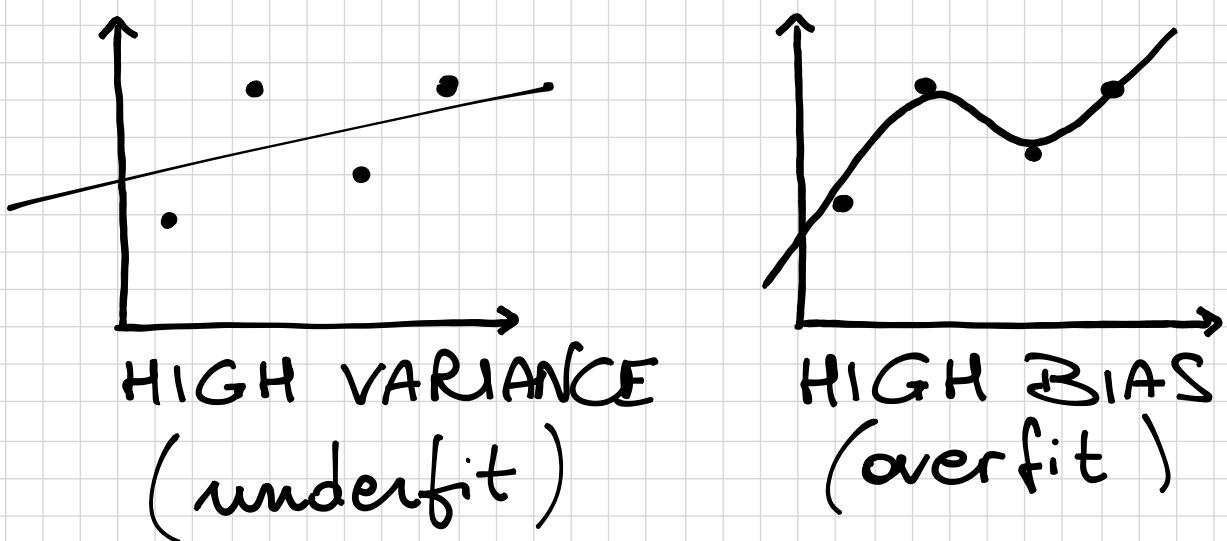


N.B.: there are good optimisations of GD, such as ADAM

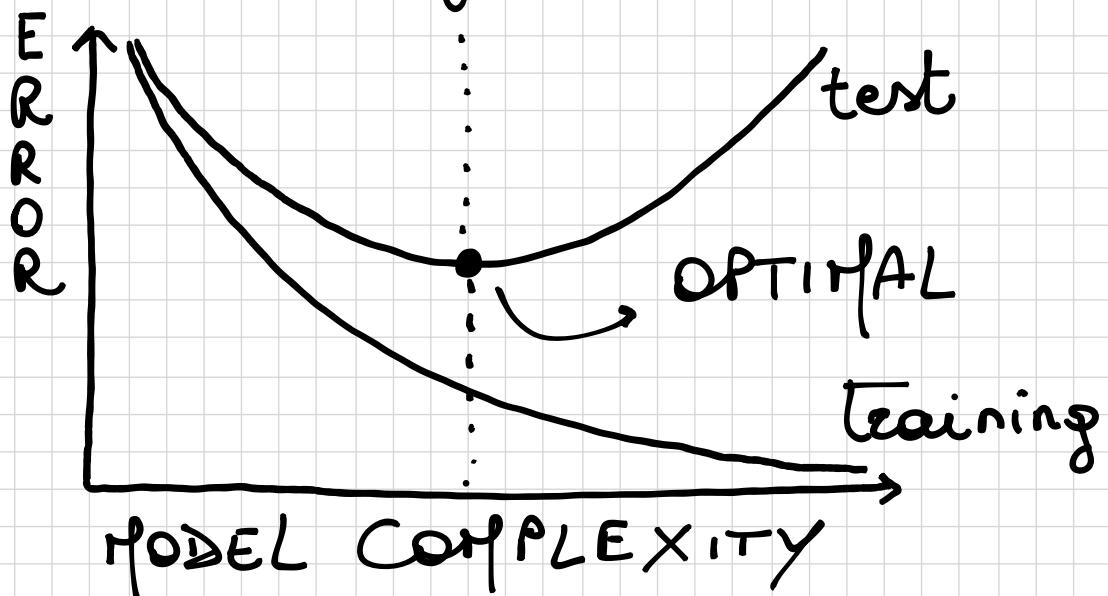
[Kingma, Ba, 1412.6980]

## 2.7 Variance - Bias Tradeoff

GD may induce issues:

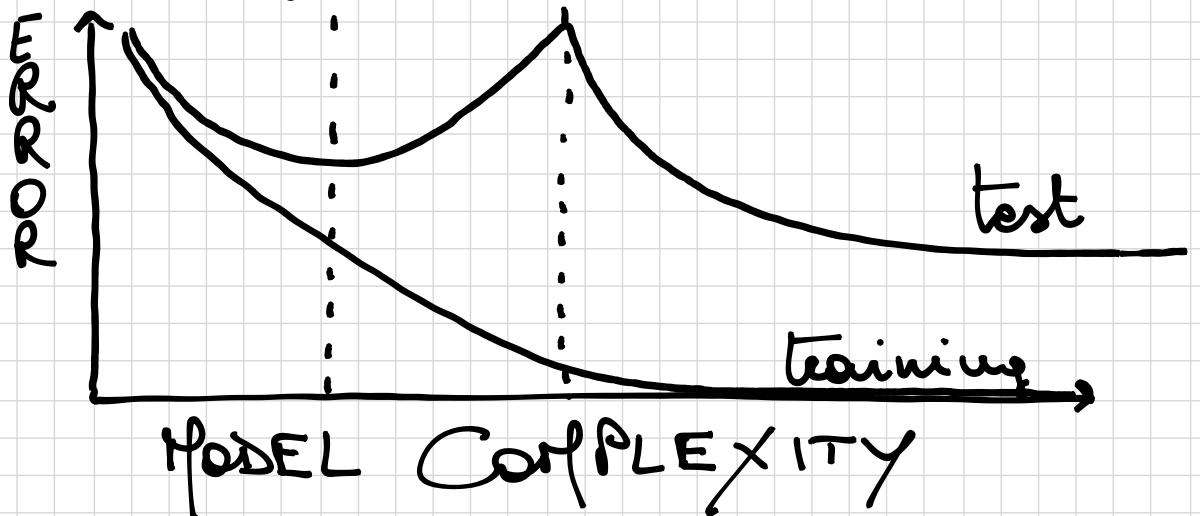


Classical argument:



⇒ must find the optimal point

# Modern argument



← CLASSICAL REGIM X "OVERPARAM"

[Belkin et al., 1812.11118]

NB for physicists:

- high complexity = hard interpr.
- less params = easier train.
- geometrical tasks are "easier": do we need all those parameters?

NB: the training error can be made arbitrarily small: NNs are UNIVERSAL APPROXIM.

[Hornik, Stinchcombe, White,  
Neural Networks 2 (1989)]

### 3. Modern Applications

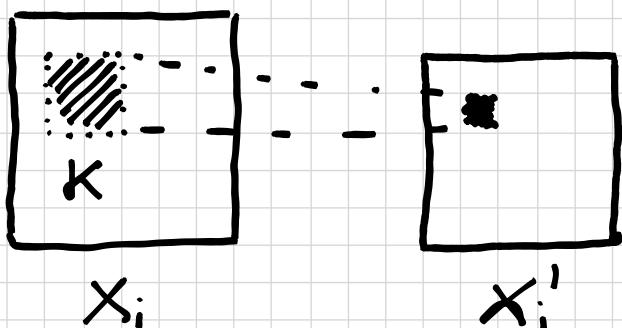
Requirements for training:

- i) differentiable forward pass
- ii) Convex loss function  
(at least locally)

Are there more interesting functions,  
other than MLPs?

#### 3.1 Convolutional Neural Networks

Convolution filter:



$$x_i \in \mathbb{R}^{N \times M}$$
$$K \in \mathbb{R}^{R \times S}$$

$$i = 1, \dots, N^p$$

then

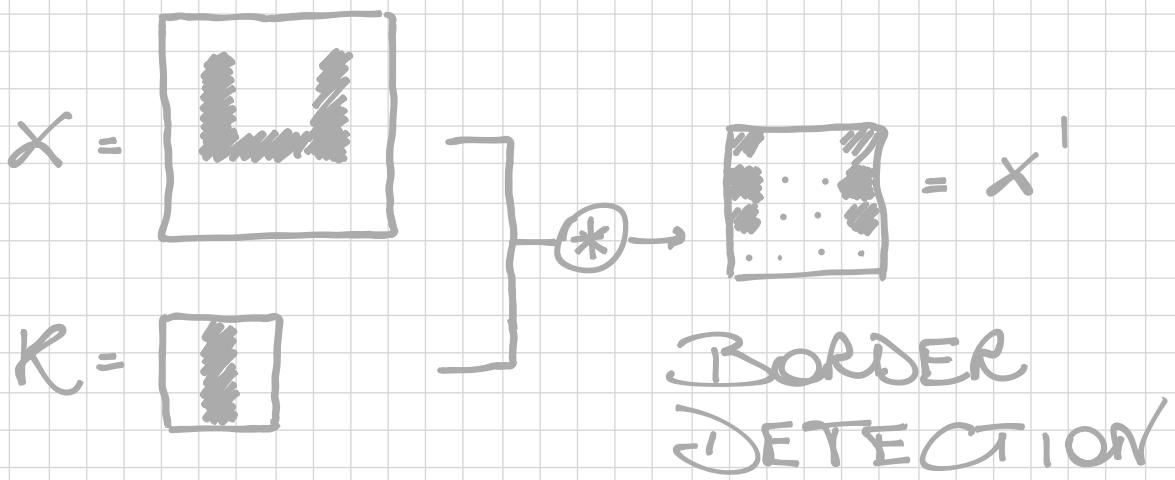
$$x'_{i,uv} = \sum_{r=0}^{R-1} \sum_{s=0}^{S-1} x_{i,u+r, u+s} K_{r,s}$$

thus  $x'_i \in \mathbb{R}^{N-R+1, M-S+1}$

Why useful?

- local features
- shared parameters
- neighbouring properties

e.g.



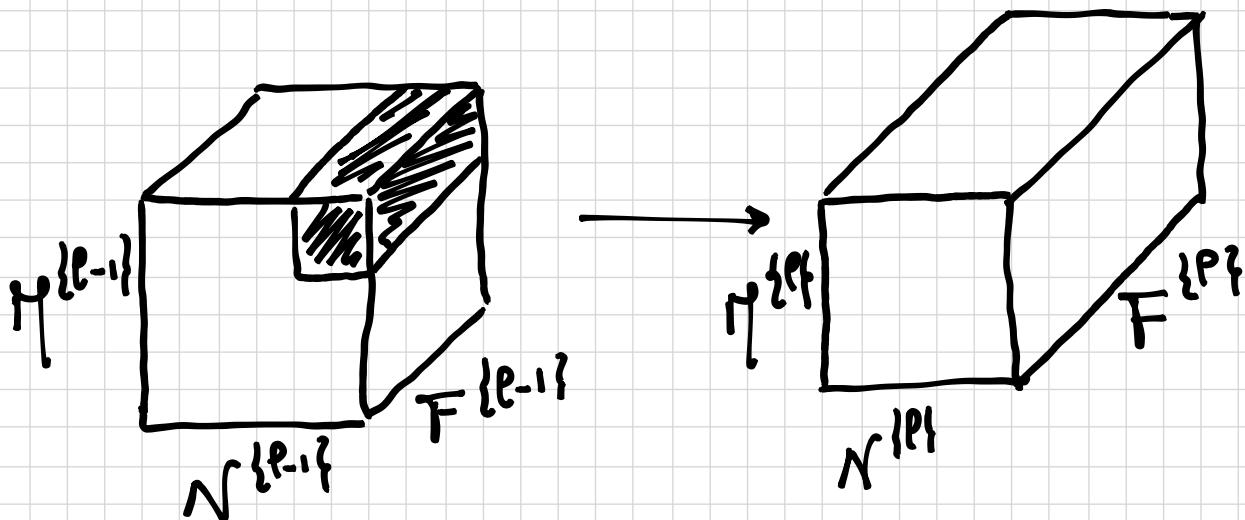
As it is a DIFFERENTIABLE OPERATION, we can learn the filters (more than one at a time).

For each layer  $l \in \{1, \dots, L\}$ :

$$a_i^{\{l\}} \in \mathbb{R}^{N^{\{l\}} \times H^{\{l\}} \times F^{\{l\}}}$$
$$K^{\{l\}} \in \mathbb{R}^{F'^{\{l\}} \times F^{\{l\}} \times R^{\{l\}} \times S^{\{l\}}}$$

we have:

$$\left\{ \begin{array}{l} z_{i,uv} = \sum_{r=0}^{P-1} \sum_{s=0}^{P-1} \sum_{f=0}^{P-1} a_{i,u+r,v+s} K_{r,s} \\ q_{i,uv} = f(z_{i,uv}) \end{array} \right.$$

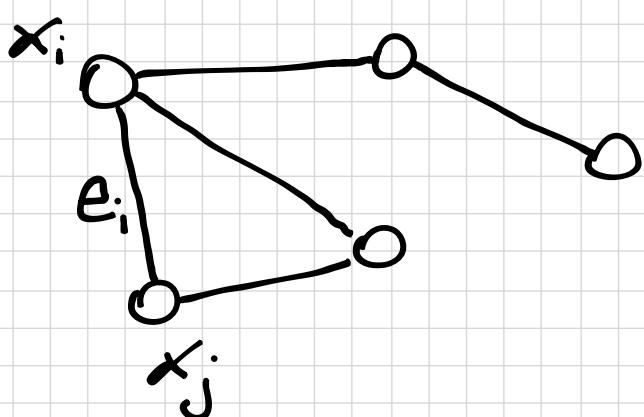


[LeCun et al., *Neur. Comp.* 1  
(1989)]

### 3.2 Graph Neural Network

Graph  $G = (V, E)$

vertices      edges



# How to perform graph Computations?

## 1) EDGE UPDATE

$$e'_k = \phi^e(e_k, x_k^{(r)}, x_k^{(s)}, u)$$

"compute forces"

## 2) EDGE AGGREGATION

$$\bar{e}'_i = \rho^{e \rightarrow v}(E_{\{i\}})$$

"Compute resulting forces"

## 3) NODE UPDATE

$$x'_i = \phi^v(\bar{e}'_i, x_i, u)$$

"update pos. and vel."

## 4) NODE AGGREGATION

$$\bar{e} = \rho^{e \rightarrow u}(E')$$

"Sum ALL forces"

## 5) GLOBAL UPDATE

$$u' = \phi^u(\bar{e}', \bar{x}', u)$$

"kinetic energy"

## 5) GLOBAL AGGREGATION

$$\bar{x}' = \rho^{v \rightarrow u}(v') \text{ "total en."}$$

How to use NNs?

$$\phi_{\{e, v, u\}} \longrightarrow NN_{\theta}^{\{e, v, u\}}$$

and for instance

$$\rho_{(i)}^{e \rightarrow v} = \sum_{\{k | r_k = i\}} e'_{k(i)}$$

$$\rho^{v \rightarrow u} = \sum_i x'_i$$

$$\rho^{e \rightarrow u} = \sum_i e'_{ik}$$

This is encoded in MESSAGE

PARSING

[Güneş et al., 1704.01212]

$$x'_i = \gamma \left( x_i, \bigcup_{j \in N(i)} \phi(x_i, x_j, e_{j \rightarrow i}) \right)$$

update ↲

message  
↳ aggregation

## e.g.: GRAPH CONVOLUTION

let  $A = (a_{ij})$  s.t.

$$a_{ij} = \begin{cases} 1 & \text{if } x_i \text{ conn. to } x_j \\ 0 & \text{otherwise} \end{cases}$$

then

$$\mathcal{D} = (d_{ij}) \Leftarrow d_{ii} = \sum_j a_{ij}$$

$$\hat{A} = A + \mathbb{I}$$

then:

$$x' = f\left[\left(\mathcal{D}^{-\frac{1}{2}} \hat{A} \mathcal{D}^{-\frac{1}{2}}\right) x w + b\right]$$

- only local features
- aggregation of neighbours.

# 4. APPLICATIONS TO STRING THEORY

## 4.1 COMPUTER VISION FOR COMPLETE INTERSECTIONS CALABI-YAU MANIFOLDS

DEF: CY : complex manifold  $X$   
s.t.  $\dim_{\mathbb{C}} X = N$   
and  $H^0(X) \subseteq SU(N)$   
(or  $C_1(X) = 0$ )

$\Rightarrow$  ST: no known metrics for  
COMPACT CY

A simple way to build CY is  
through COMPLETE INTERSECTIONS:

let  $A = \mathbb{P}^{n_1} \times \dots \times \mathbb{P}^{n_m}$

det

$$\begin{cases} P_{a_k}(z) = \prod_{I_{1k} \dots I_{ak}} z^{I_{1k}} \dots z^{I_{ak}} = 0 \\ P_{a_k}(\lambda z) = \lambda^{a_k} P_{a_k}(z) \end{cases}$$

Topologically:

$$X = \left[ \begin{array}{c|cc} P^{n_1} & a'_1 & \dots & a'_k \\ \vdots & \vdots & & \vdots \\ P^{n_m} & a''_1 & \dots & a''_k \end{array} \right] \in N^{m \times k}$$

s.t.  $\dim_{\mathbb{C}} X = \sum_{r=1}^m n_r - k = N$

$$n_r + 1 - \sum_{\alpha=1}^k a_\alpha^r = 0 \quad \forall r = 1, \dots, m$$

Why?

$X \rightarrow$  predict topological  
invariants

$\Rightarrow$  direct connection to  
EFT in ST.

Particularly:

$$h^{(r,s)} = H_d^{(r,s)}(x, c)$$

e.g.:  $N=3 \Rightarrow h^{(1,1)}$  and  $h^{(2,1)}$   
count the no. of  
chiral mult. in HE  
or hyper- and vector  
mult. in Type II

[He, 1706.02414]

[He, Lukas, 2009.02544]

## Datasets

$N=3$  [Candela, Dale, Lutken,  
Nucl. Phys. B 298 (1988)]

[Green, Habsch, Lutken,  
Class. Quant. Grav. 6 (1989)]

[Anderson et al., 1708.07907]

$N=4$  [Gray, Haupt, Lukas,  
1303.1832 and 1405.2073]

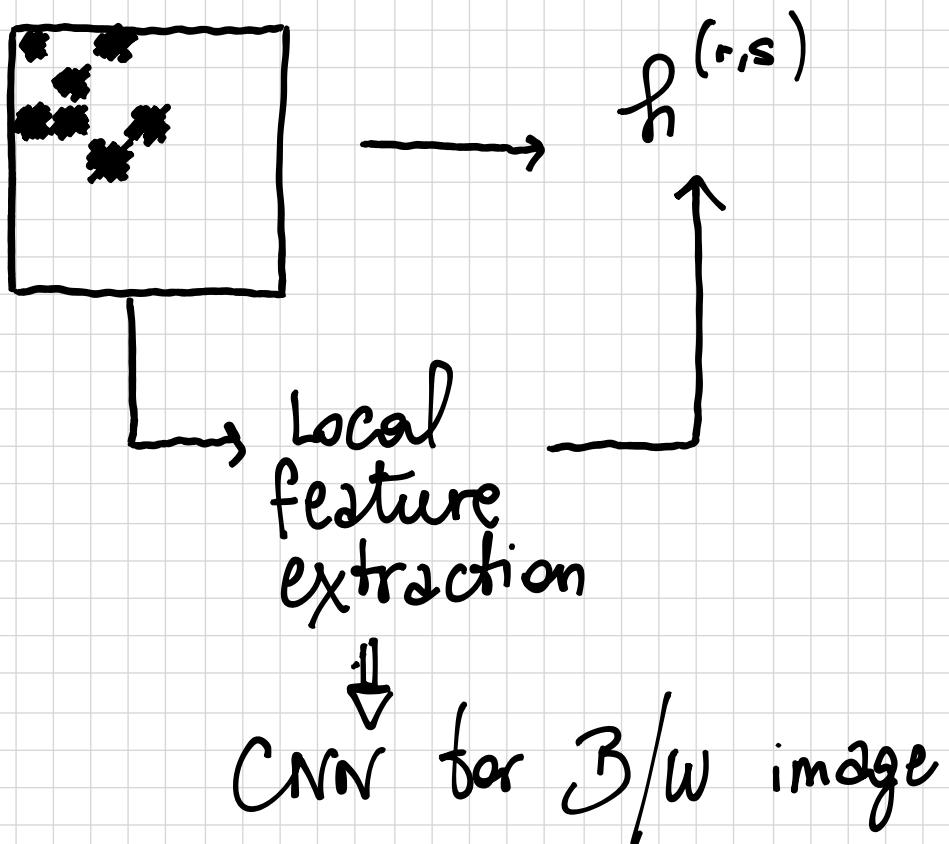
# Attempts (through NN's)

$$f: \mathbb{N}^{m,k} \longrightarrow \mathbb{N}$$
$$x \longmapsto f(x) = h^{(r,s)}$$

replace  $f$  by  $NN_\theta$ :

FULLY CONNECTED  $\Rightarrow \sim 75\%$  acc. on  $h^{(1,1)}$   
[Bull et al., ]

However:



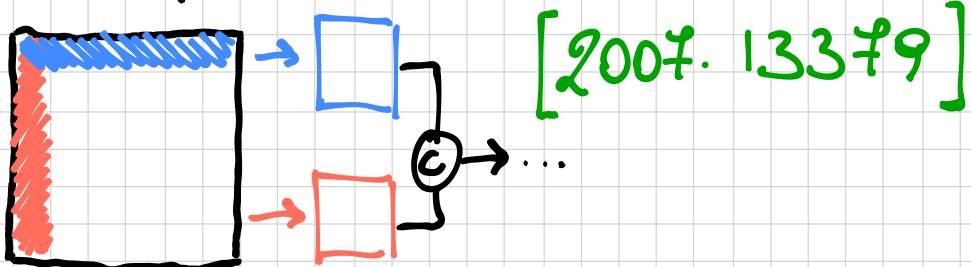
How to build such arch. ?

## 1) DATA ANALYSIS [2007. 1SF06]

- i) "good info" in small portions of  $X$
- ii) strong influence of the no. of proj. spaces and the dim. of the spaces  
(notice: rows and columns of  $X$ )
- iii) engineered var. are strongly corr.  $\Rightarrow$  no need for all of them

## 2) INCEPTION MODULE

- 1d kernels [Szegedy et al., 1409. 4842]
- maximal size



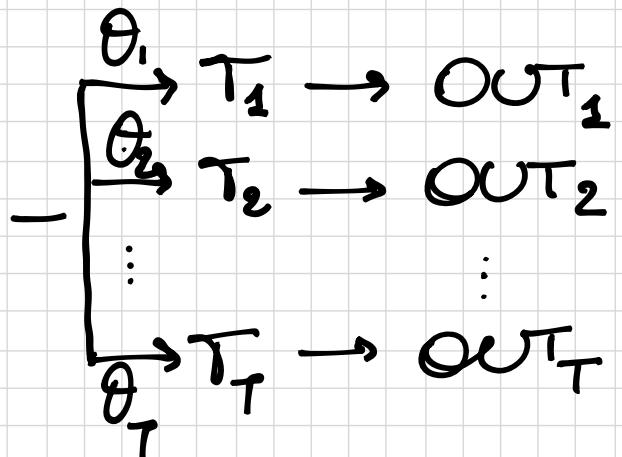
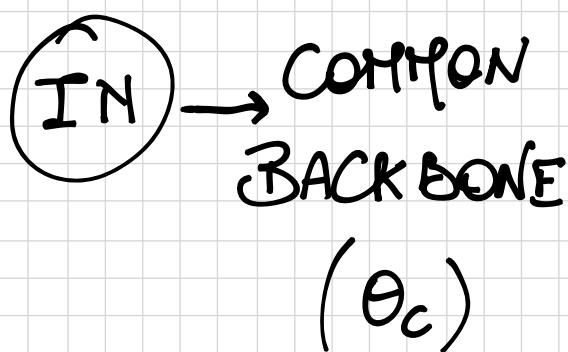
- let AI decide everything
- no prior math knowledge

$$N=3 \Rightarrow \begin{aligned} h^{(1,1)} &\sim 100\% \\ (h^{(2,1)} &\sim 50\%) \end{aligned}$$

- Deep multi-task learning:

[2108.02221]

$$L(\theta) = \sum_{t=1}^T L_t(\theta), \quad \theta = \{\theta_c, \theta_1, \dots\}$$



- only 1 model for 921497 conf. matr.
- highlight dep.
- no overfit + outliers

$$N=4 \Rightarrow (h^{(1,1)}, h^{(2,1)}, h^{(3,1)}, h^{(2,2)}) \sim (100\%, 100\%, 96\%, 83\%)$$

using 80 : 10 : 10