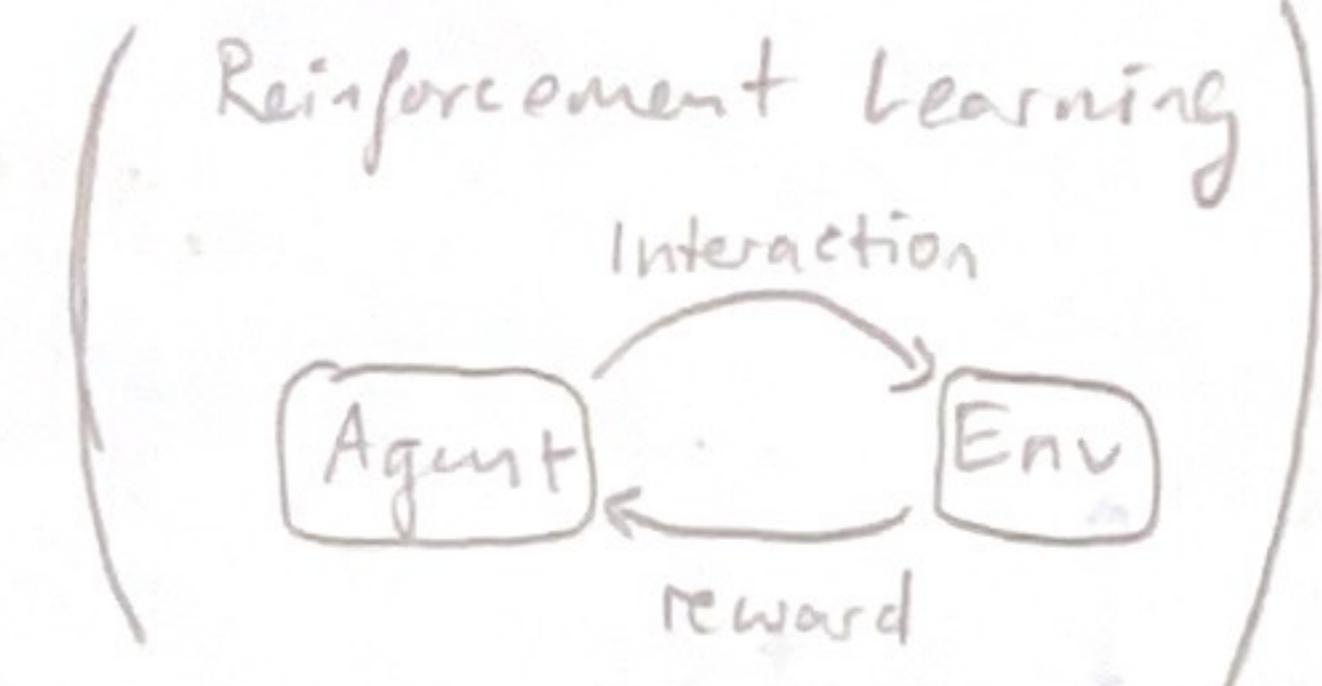


## 12DL 2: ML Basics

1

- ML: Unsupervised vs. Supervised Learning:

- \* No labels
- \* Labels
- \* Find structure of the data
- \* Clustering



- Properties of k-NN:

- Hyperparameters
  - Distance metric: L<sub>1</sub>-Norm, L<sub>2</sub>-Norm
  - No. of neighbors k
- Values are Problem dependent

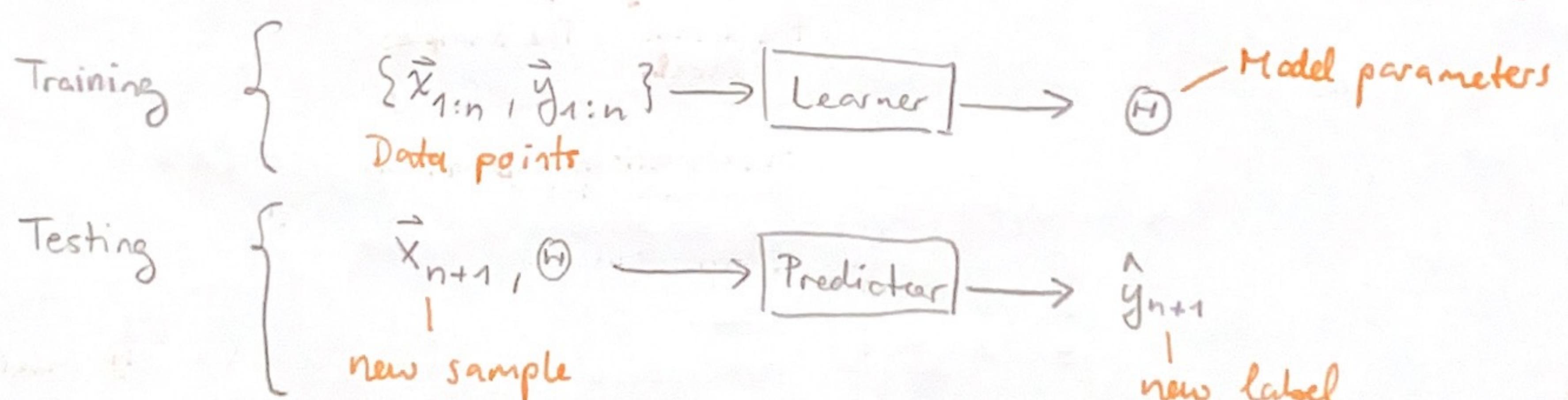
Training data	Validation d.	Test d.
---------------	---------------	---------

Find hyperparameters

- What is Cross Validation?

- Split the ~~data~~ training data into N folds
- ⇒ Train on N-1 folds and validate on 1 fold

- Linear Regression: Find a linear model that maps inputs x to target y.



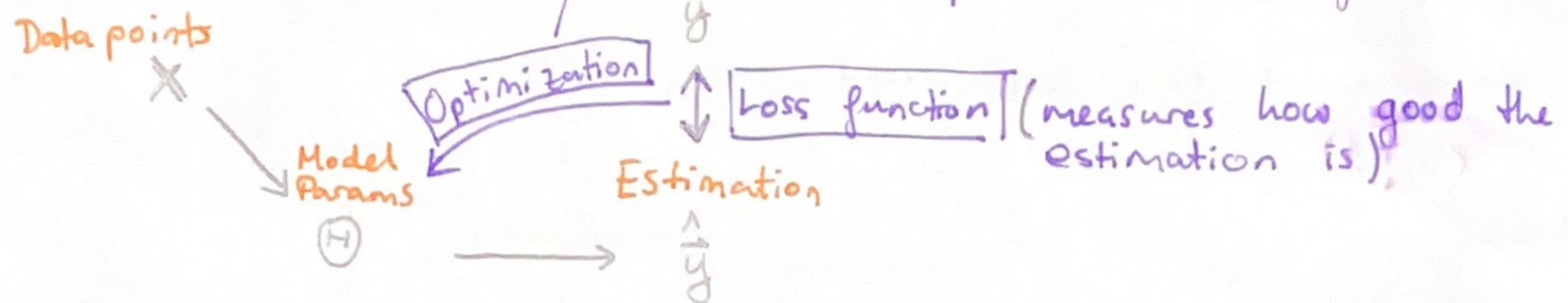
Linear model:  $\hat{y}_i = \Theta_0 + \sum_{j=1}^d x_{ij} \Theta_j$  — weights (for each feature)

Prediction for ith sample | Bias | Input sample

$$\Rightarrow \hat{y} = X\vec{\Theta} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ 1 & x_{21} & \dots & x_{2d} \\ \vdots & \ddots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{bmatrix} \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \vdots \\ \Theta_d \end{bmatrix}$$

Input sample with d features

How to obtain the model?



Lin. Regression loss function:  
(Objective function/  
Cost function)

$$J(\vec{\theta}) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Optimization : (Least squares estimate)

$$\begin{aligned} \min_{\vec{\theta}} J(\vec{\theta}) &= \frac{1}{n} \sum_{i=1}^n (\vec{x}_i^\top \vec{\theta} - y_i)^2 \\ &= (\vec{X}\vec{\theta} - \vec{y})^\top (\vec{X}\vec{\theta} - \vec{y}) \\ \Rightarrow \frac{\partial J(\vec{\theta})}{\partial \vec{\theta}} &= 2\vec{X}^\top \vec{\theta} - 2\vec{X}^\top \vec{y} = 0 \\ \Rightarrow \vec{\theta} &= (\vec{X}^\top \vec{X})^{-1} \vec{X}^\top \vec{y} \end{aligned}$$

- What's the maximum likelihood estimate? (MLE)

$$\vec{\theta}_{ML} = \underset{\vec{\theta}}{\operatorname{argmax}} \ p_{\text{Model}}(\vec{y} | \vec{X}, \vec{\theta})$$

↓  
 the actual labels      The samples  
 observations from  $p_{\text{data}}(\vec{y} | \vec{X})$

"Find the parameters that max. the probability of that the model's predictions match the actual observations"

(Or "Probabilities of seeing outcomes that I actually see")

Data is i.i.d.  $\Rightarrow$  independence  $\Rightarrow p(a, b) = p(a)p(b)$

$$= \underset{\vec{\theta}}{\operatorname{argmax}} \prod_{i=1}^n p_{\text{Model}}(y_i | \vec{x}_i, \vec{\theta})$$

" $\vec{\theta}$  maximizes [...] as well as  $\log(\cdot)$ , since  $\log$  is a monotonic function"

\* Gaussian Assumption

Assume  $y_i \sim N(\vec{x}_i^\top \vec{\theta}, \sigma^2)$   
Mean

$$\Rightarrow p(y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - \mu)^2}$$

$$= \underset{\vec{\theta}}{\operatorname{argmax}} \left( \sum_{i=1}^n -\frac{1}{2} \log(2\pi\sigma^2) + \sum_{i=1}^n -\frac{1}{2\sigma^2} (y_i - \vec{x}_i^\top \vec{\theta}) \right)$$

$$= \underset{\vec{\theta}}{\operatorname{argmax}} \left( -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\vec{y} - \vec{X}\vec{\theta})^\top (\vec{y} - \vec{X}\vec{\theta}) \right)$$

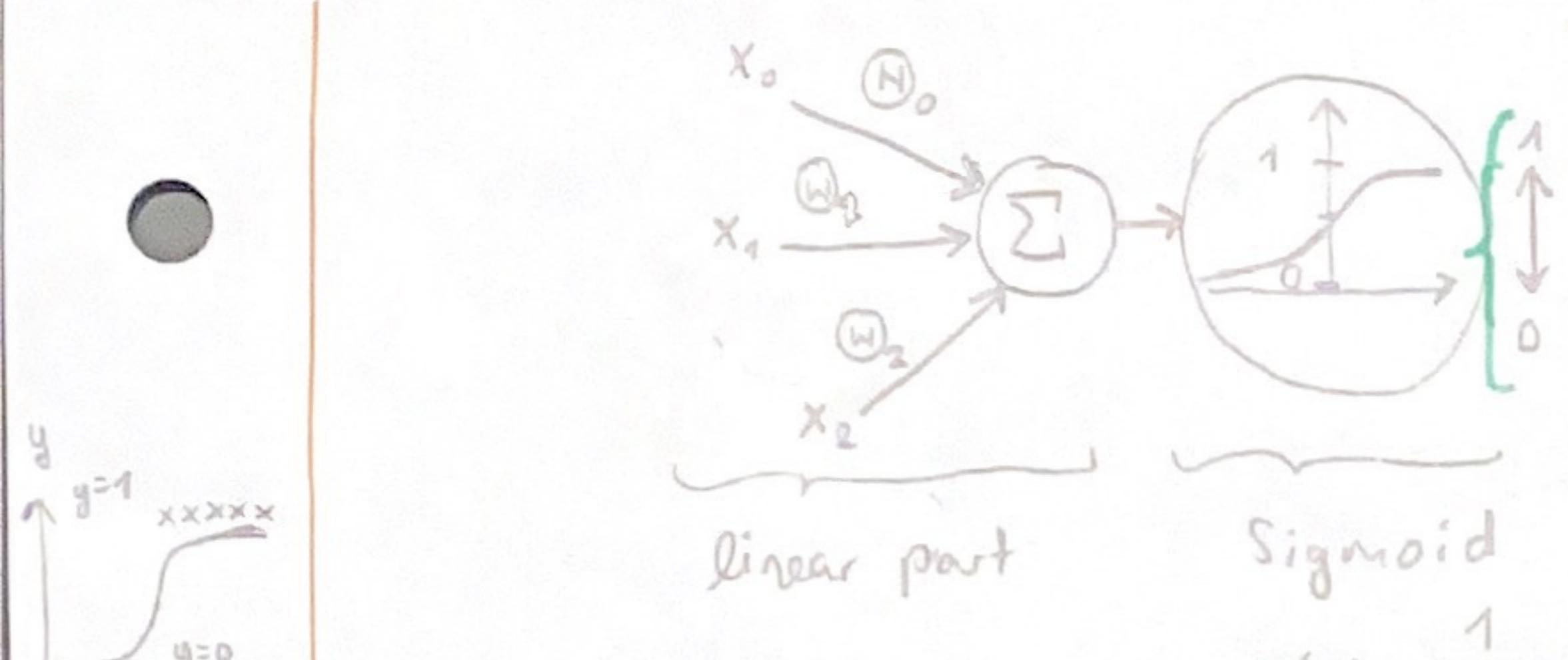
$$\frac{\partial J(\vec{\theta})}{\partial \vec{\theta}}$$

$$\Rightarrow \vec{\theta} = (\vec{X}^\top \vec{X})^{-1} \vec{X}^\top \vec{y}$$

⇒ MLE corresponds to the Least Squares Estimate (same solution)

- What is Logistic Regression? → for binary predictions

→ Taking the output from a lin. regression and feeding it into a sigmoid function



⇒ The label could be 0 or 1.

(sigmoid)

⇒ The model's prediction:

$$\hat{y}_i = \sigma(x_i \cdot \vec{\theta})$$

model for the

Is a Probability of a binary output:

Model of Coins:

$$p(z|\phi) = \phi^z (1-\phi)^{1-z} =$$

$$= \begin{cases} \phi, & z=1 \\ 1-\phi, & z=0 \end{cases}$$

The prediction of  
the sigmoid

$$\hat{y} = p(y=1 | X, \vec{\theta}) = \prod_{i=1}^n p(y_i=1 | x_i, \vec{\theta})$$

Bernoulli Trial label (discrete)

$$= \prod_{i=1}^n \hat{y}_i^{y_i} (1-\hat{y}_i)^{1-y_i}$$

prediction (continuous)

$$= p(\vec{y} | X, \vec{\theta})$$

MLE: Find the  $\vec{\theta}$  that max. the likelihood of the observations given  $\vec{\theta}$

$$\Rightarrow \vec{\theta}_{ML} = \underset{\vec{\theta}}{\operatorname{argmax}} \log p(\vec{y} | X, \vec{\theta})$$

$$= \underset{\vec{\theta}}{\operatorname{argmax}} \sum_{i=1}^n \log (\hat{y}_i^{y_i} (1-\hat{y}_i)^{1-y_i})$$

$$= \underset{\vec{\theta}}{\operatorname{argmax}} \sum_{i=1}^n (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i))$$

$$\Rightarrow \text{Loss } \boxed{L(\hat{y}_i, y_i) = -[y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)]}$$

e.g.:

$$y_i = 1 \rightarrow L = \log(\hat{y}_i) \rightarrow \hat{y}_i \text{ should be large}$$

binary-cross-entropy (BCE)

$$y_i = 0 \rightarrow L = \log(1-\hat{y}_i) \rightarrow \hat{y}_i \text{ should be small}$$

Minus in  $L$

→ Optimization:

$$\left[ \text{Cost function: } C(\vec{\theta}) = +\frac{1}{n} \sum_{i=1}^n L(\hat{y}_i, y_i) \right]$$

⇒ Minimize via Gradient descent, since there is no closed form solution.

## I2DL 03: Introduction to Neural Networks

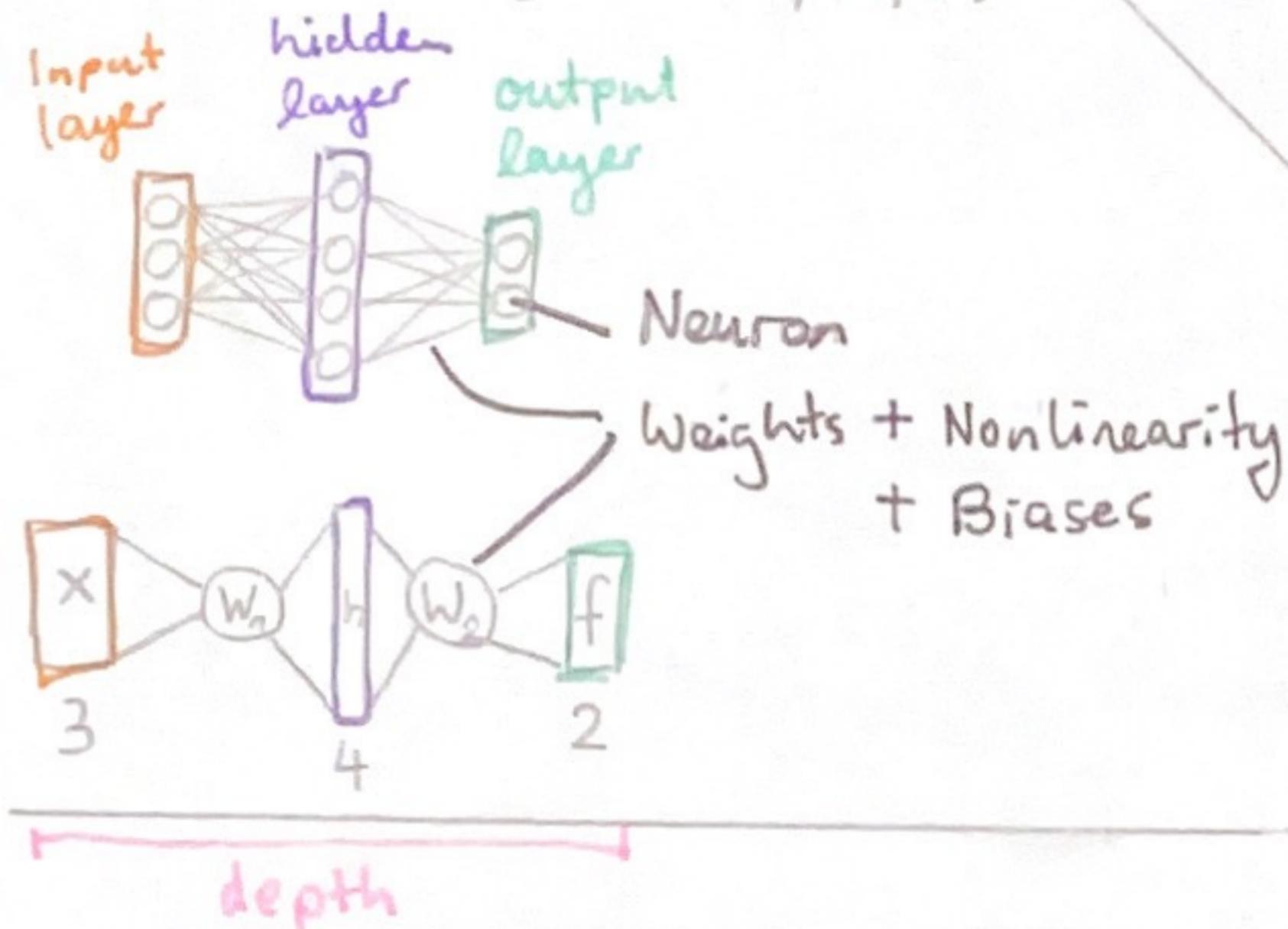
(3)

↳ Can be written as non-linear score function

$$f = \dots (\underbrace{\max(0, \underbrace{w_1 \cdot x}_{\text{1. layer}})}_{\text{weight of 1. layer}})$$

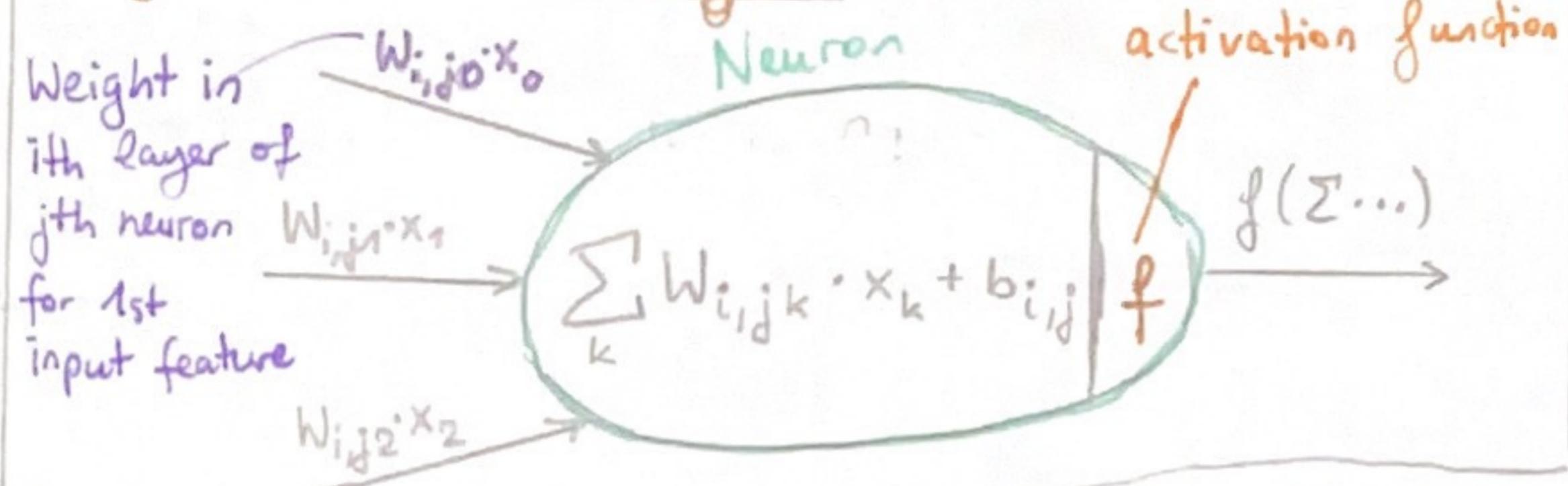
Example: 2-layer NN

$$f = w_2 \max(0, w_1 \cdot x)$$



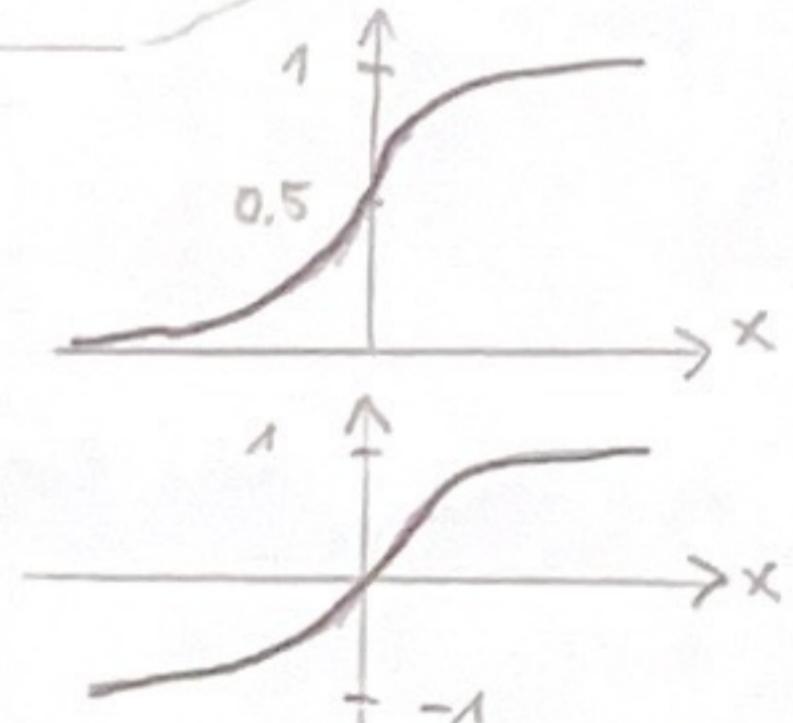
→ ANNs are inspired by the brain

Single Neuron in a Layer:



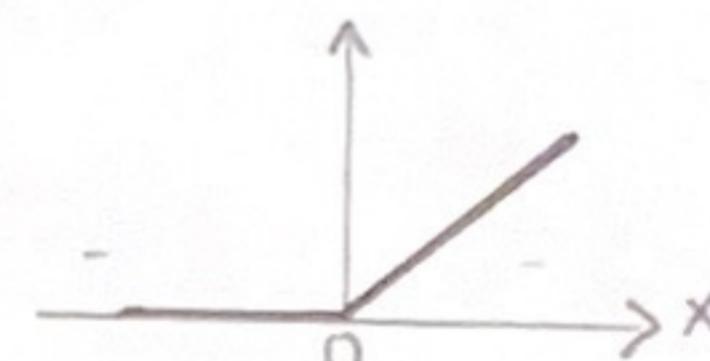
- Activation functions: \* Sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$

\* tanh:  $\tanh(x)$

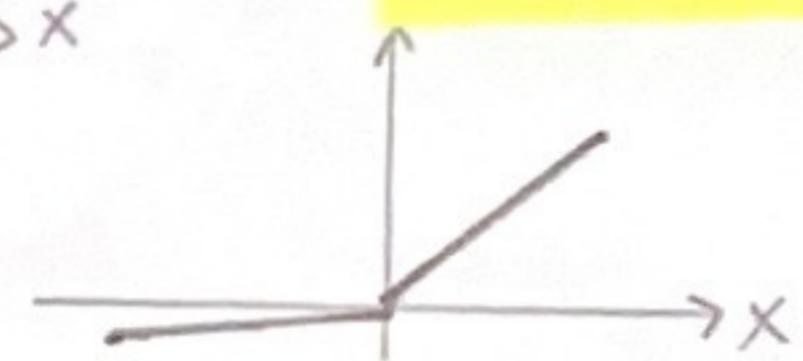


Should be steeper than Sigmoid

\* ReLU:  $\max(0, x)$



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



\* Parametric ReLU:  $\max(\alpha x, x)$  (learnable parameter)

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

\* ELU:  $f(x) = \begin{cases} x, & x > 0 \\ x(e^x - 1), & x \leq 0 \end{cases}$

- Why is a NN structured into layers?

- 1) Usually all neurons per layer are the same w/ very simple math
- 2) Structure allows efficient computation
- 3) Inspired by neurons in the brain

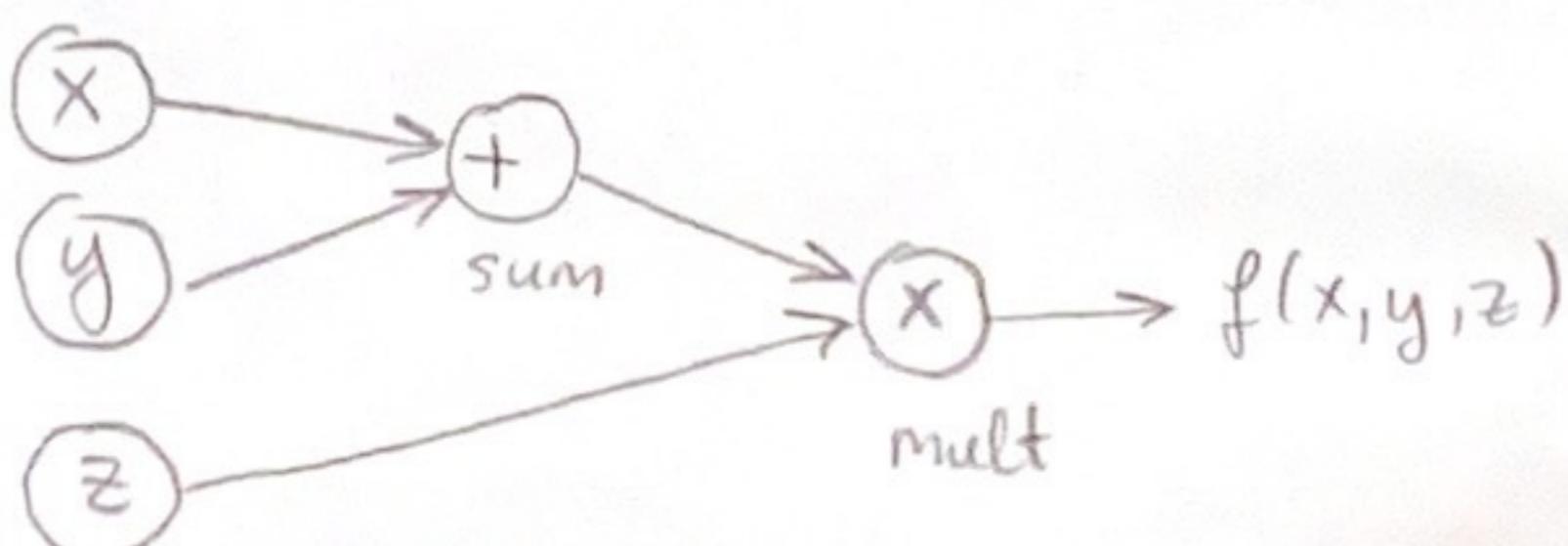
- What's a compute graph/computational graph?

→ Directional graph

→ Matrix operations represented as compute nodes

→ Variables or operators (+, -, \*, /, log(), exp()) as vertex nodes

e.g.:  $f(x, y, z) = (x + y)z$



- Loss functions: A function to measure the network performance

\* Small loss = good predictions

\* depends heavily on the problem Ground truth

L1-Loss:

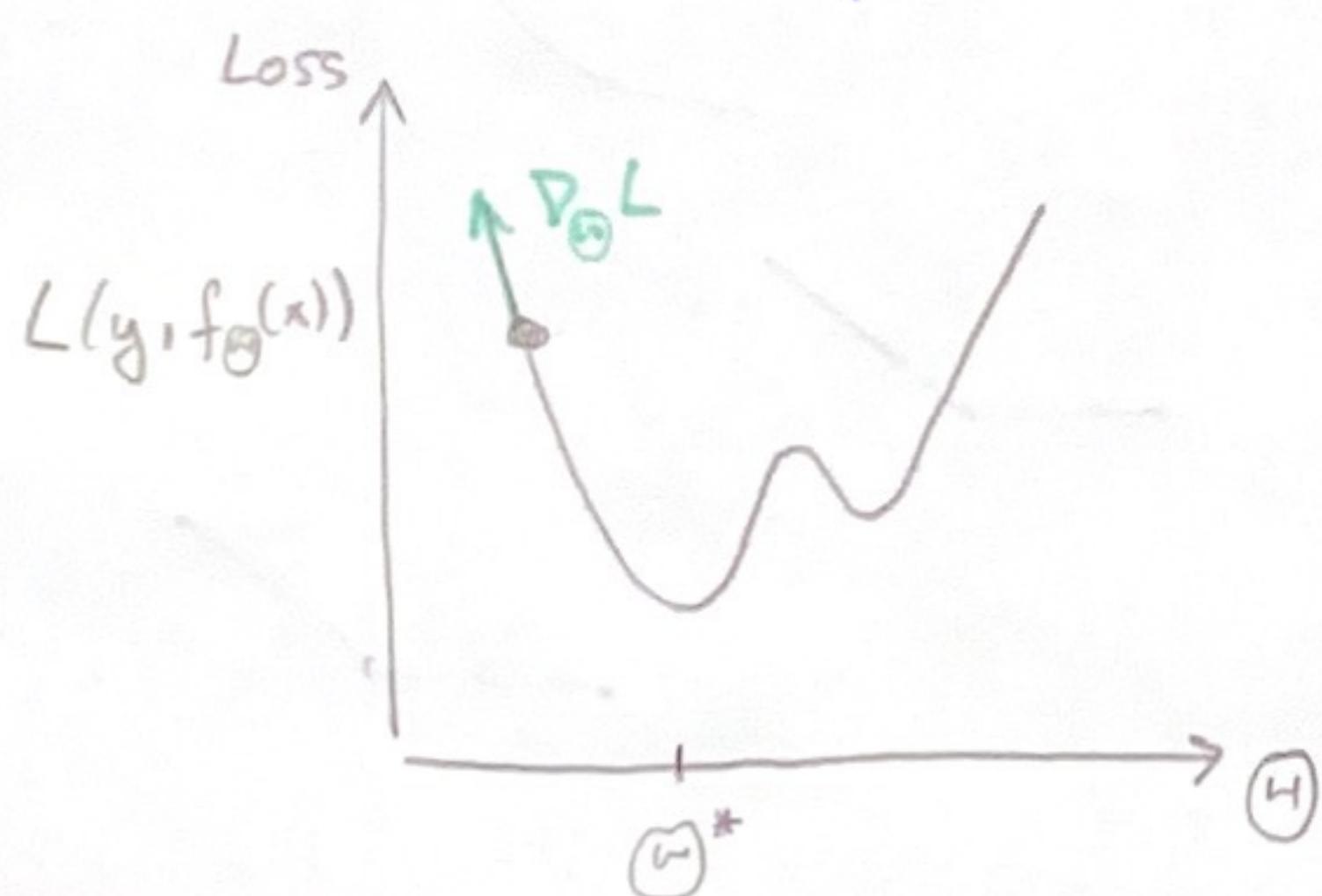
$$L(y, \hat{y}; \Theta) = \frac{1}{n} \sum_i^n \|y_i - \hat{y}_i\|_1$$

$$\text{MSE-Loss: } L(y, \hat{y}; \Theta) = \frac{1}{n} \sum_i^n \|y_i - \hat{y}_i\|_2^2$$

$$\text{Binary-Cross-Entropy: } L(y, \hat{y}; \Theta) = -\frac{1}{n} \sum_{i=1}^n [y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)]$$

$$\text{Cross-Entropy: } L(y, \hat{y}; \Theta) = -\sum_{i=1}^n \sum_{k=1}^K (y_{ik} \log \hat{y}_{ik})$$

- How is the NN improved?



Goal: Find  $\Theta^* = \operatorname{argmin} L(y, f_\Theta(x))$

$\Rightarrow$  Iterate:  $\Theta = \Theta - \alpha \nabla_\Theta L(y, f_\Theta(x))$

step size / learning rate

Since the gradient  
points in the direction  
of the greatest  
ascent of L

- Why is gradient descent used?

→ NNs can be formulated as compute graphs  
and GD is easy to compute on a compute graph.

## I2DL 04: Optimization & Backpropagation

(4)

- What's the idea of Backprop?

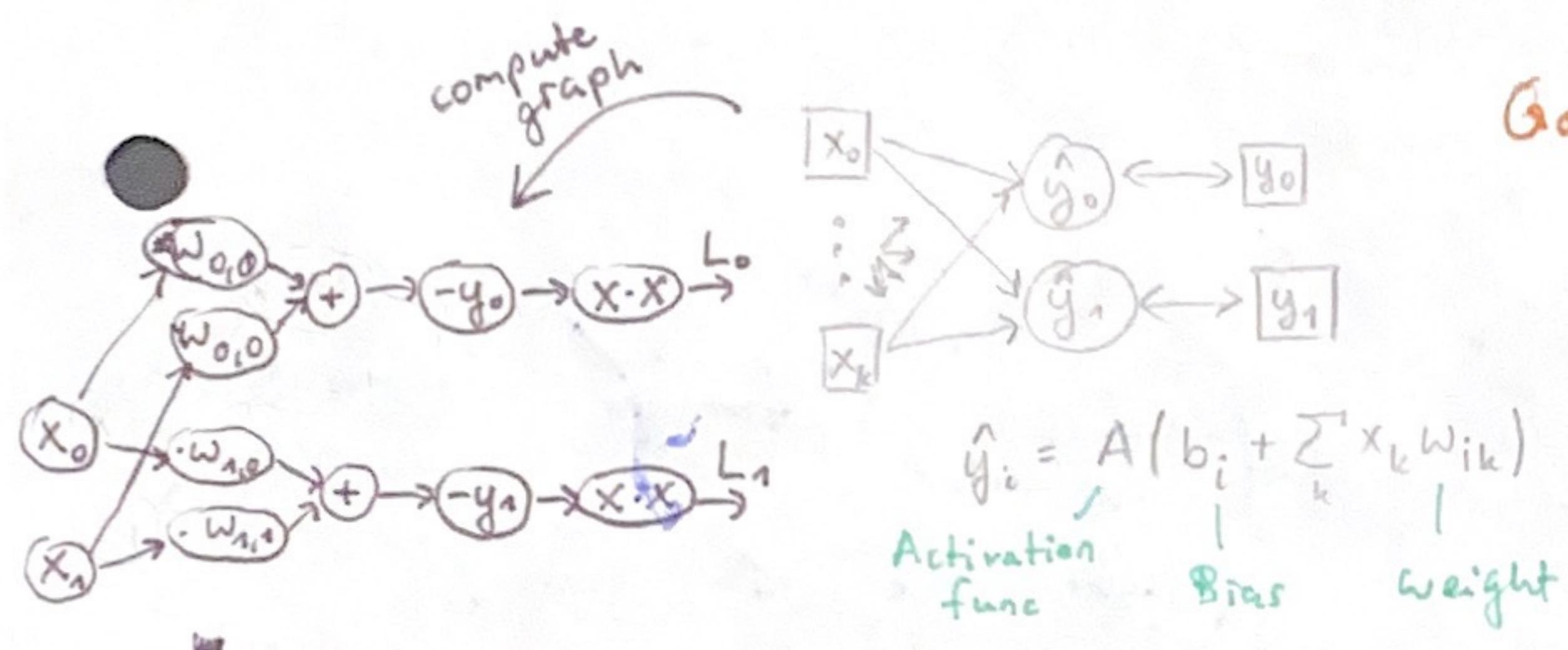
→ Break down the gradient computation  $\nabla_{\theta} L(\theta)$

e.g. for  $f(x,y,z) = (x+y)z$

- \*  $X_k$ : Input variables
- \*  $w_{l,m,n}$ : weights  
layer / weight in neuron  
neuron in layer
- \*  $\hat{y}_i$ : prediction
- \*  $y_i$ : ground truth targets
- \*  $L$ : loss

$$\begin{aligned} \frac{\partial d}{\partial x} &= \frac{\partial d}{\partial y} = 1 \\ &= \frac{\partial}{\partial d} (dz) = z \\ \frac{\partial f}{\partial x} &= \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial x} \\ \left( \frac{\partial f}{\partial y} = \dots \right) & \\ \text{Apply chain rule} \end{aligned}$$

NN example:

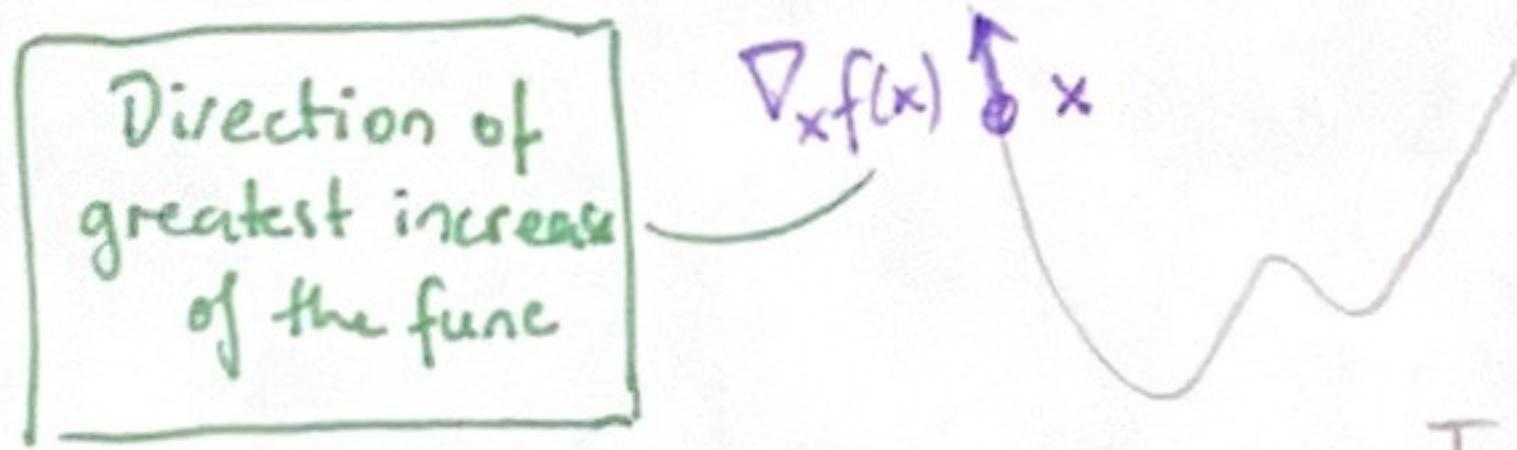


Goal: Compute  $\frac{\partial L}{\partial w_k}$  (Gradient w.r.t. all weights and biases)

$$\text{with } L = \sum_i L_i = \sum_i (\hat{y}_i - y_i)^2$$

→ Use chain rule:  $\frac{\partial L_i}{\partial w_{ik}} = \frac{\partial L_i}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w_{ik}}$

- What's Gradient descent?



$$\Rightarrow \text{Update rule: } x' = x - \alpha \nabla_x f(x)$$

/ learning rate

$$\text{For a NN: } \vec{W}' = \vec{W} - \alpha \nabla_{\vec{W}} f_{\{\vec{x}, \vec{y}\}}(\vec{W})$$

$$\text{with } \nabla_{\vec{W}} f_{\{\vec{x}, \vec{y}\}}(\vec{W}) = \begin{pmatrix} \frac{\partial f}{\partial w_{0,0,0}} \\ \vdots \\ \frac{\partial f}{\partial w_{l,m,n}} \end{pmatrix}$$

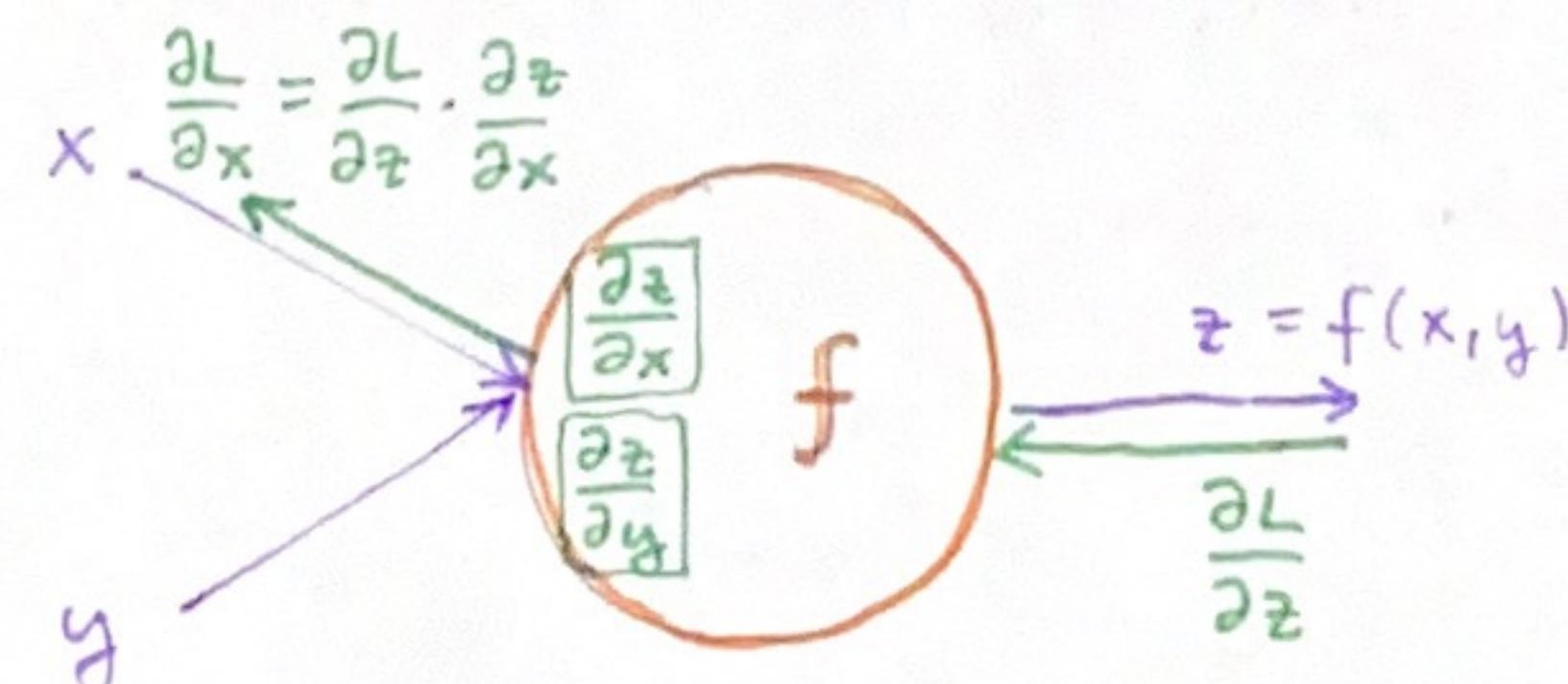
for a given training pair

- Which functions does each neuron implement?

→ Forward pass:  $\Theta$

→ Backward pass:  $\Theta$

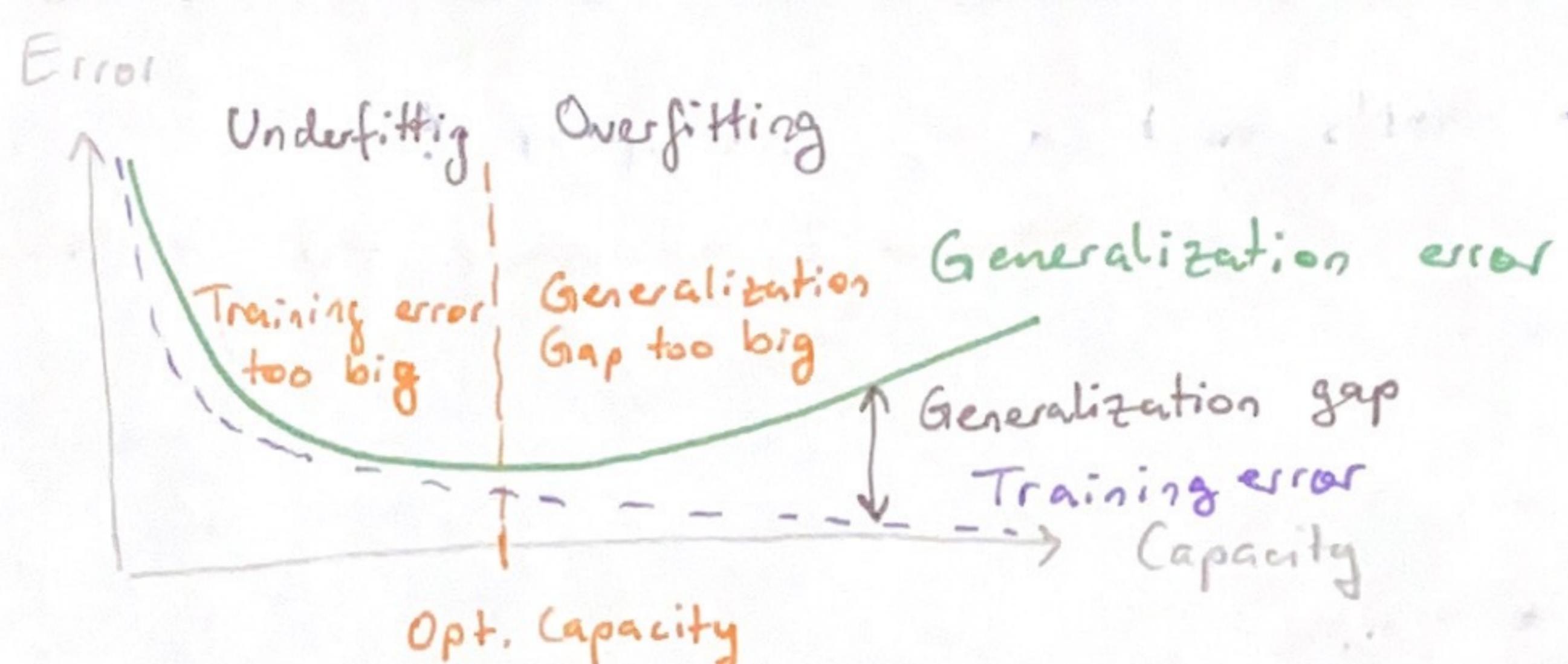
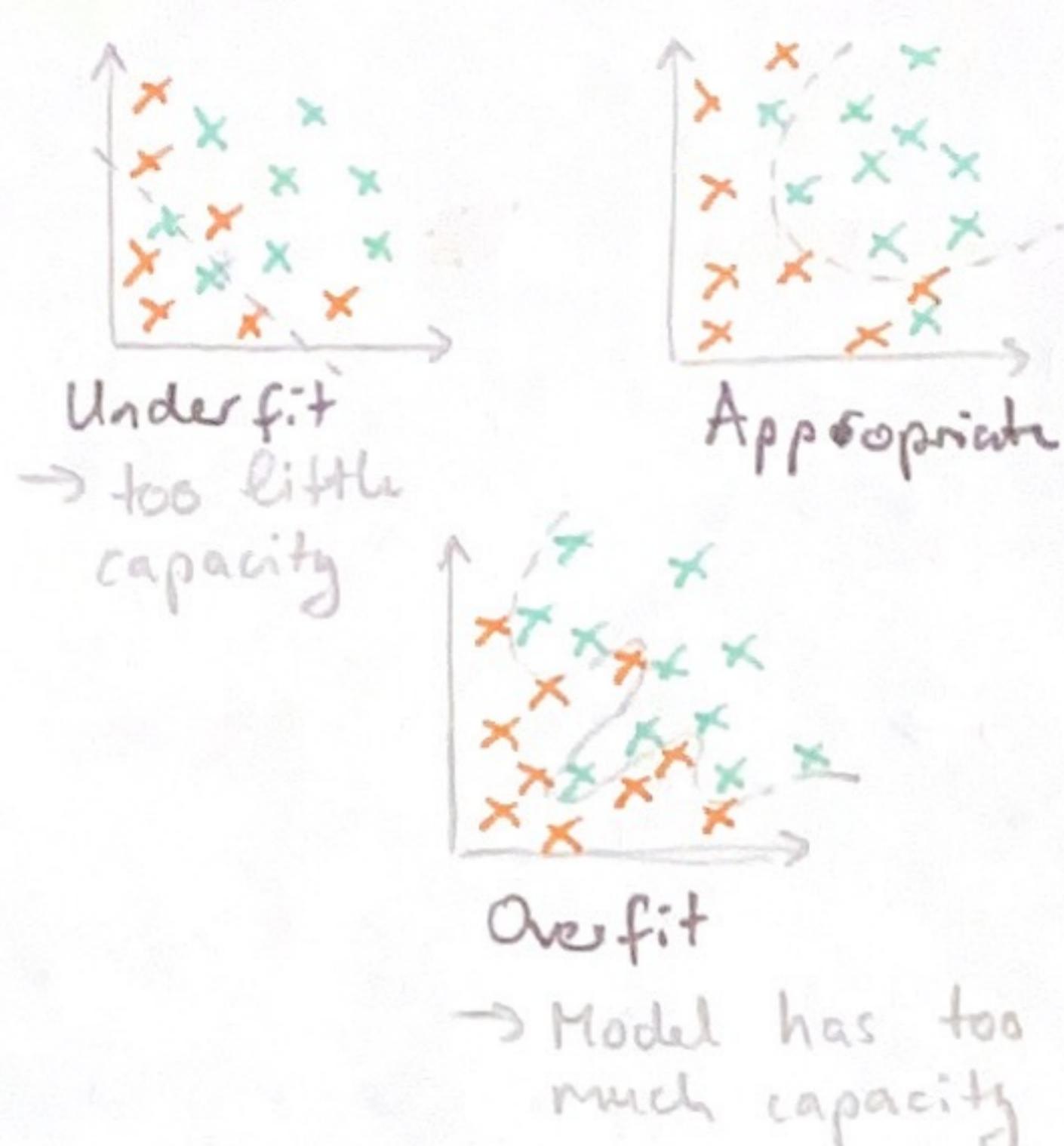
Flow of the forward pass & gradients



Unknown parameter pass  
backward

# neurons · # input channels + # bias

## - What's over- & underfitting?



## - How to prevent overfitting? → Regularization

=D Add regularization term to loss function

a.k.a. weight decay

$$\Rightarrow L(\vec{y}, \vec{\hat{y}}, \vec{\theta}) = \sum_i (x_i \theta_{ji} - y_i)^2 + \lambda R(\vec{\theta})$$

Hyper-parameter      | Regularization term

L2 Regularization:

$$R(\vec{\theta}) = \sum_i \theta_i^2$$

*favours smaller weights  
→ enforces that weights have similar values*

L1 Regularization:

$$R(\vec{\theta}) = \sum_i |\theta_i| \rightarrow \text{considers all features}$$

(only in regularization i.e. if the L-norm is taken of the weights  $\|w\|$ )

enforces sparsity of weights (a few large ones)

→ considers key features

→ Large regularization smoothes the decision / prediction

Problem: GD possibly only optimizes for the regulariz. term

Aim of regularization:

Training error ↑

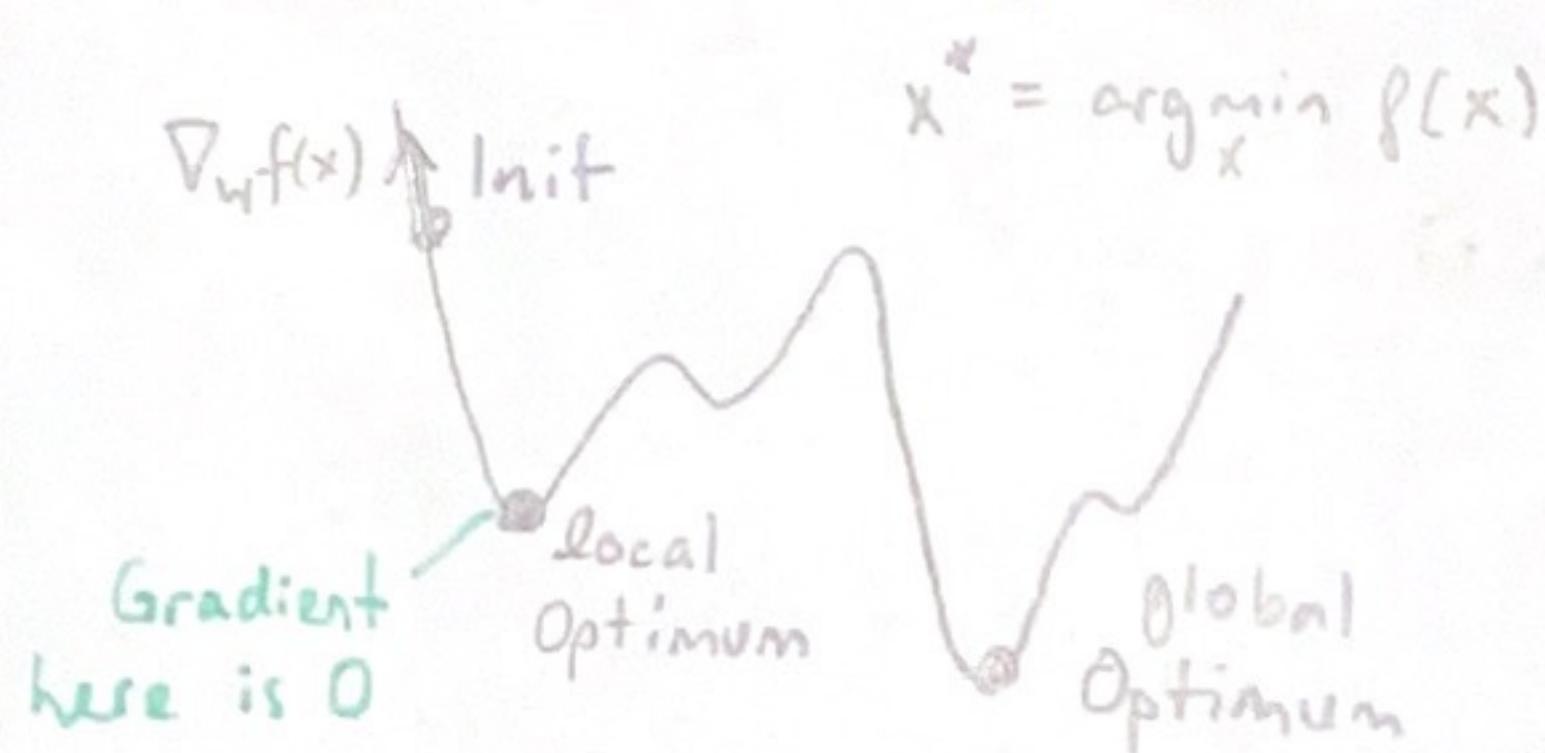
Validation error ↓

Regularization = Any technique that lowers the validation error (variance) at the cost of increasing the training error (bias)

## 12DL 05: Scaling Optimization

(5)

- How do we know if the local minimum found with G.D. is a global minimum?



→ Use Convexity

↪ If  $f$  is convex, then all local minima are global minima

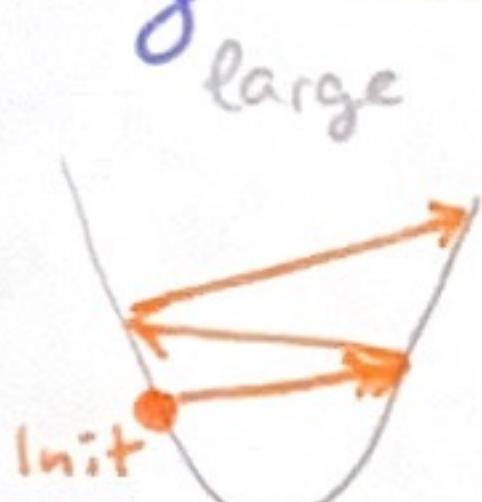


→ In general the Loss function of a NN is not convex

→ No (practical) method to know this

- Small vs large learning rate ?

in 1-D



- ⊖ Can overshoot
- ⊕ good progress



- ⊖ slow
- ⊕ does not overshoot

\* Also the learning rate may cause slow convergence during plateaus but still overshoot during steep drops of the Loss curve

In 2-D:



"Contour plot of Loss"

- What is the equation for G.D. for a single training sample and a batch?

Single Sample:

$$\vec{\theta}^{k+1} = \vec{\theta}^k - \alpha \nabla_{\vec{\theta}} L_i(\vec{\theta}^k, x_i, y_i)$$

↓                      ↓                      ↓  
 Weights & Biases    Training sample    Loss function  
 at training step                     Gradient w.r.t.  $\vec{\theta}$  (Computed via Backprop)  
 k and k+1              Learning rate

Batch:

$$\vec{\theta}^{k+1} = \vec{\theta}^k - \alpha \nabla_{\vec{\theta}} L(\vec{\theta}^k, x_{\{1..n\}}, y_{\{1..n\}})$$

$$\text{avg. Gradient} = \frac{1}{n} \sum_{i=1}^n \nabla_{\vec{\theta}} L_i(\vec{\theta}^k, x_i, y_i)$$

$$\text{"Slang"} = \nabla L = \sum_{i=1}^n \nabla_{\vec{\theta}} L_i$$

omitting  $\frac{1}{n}$  rescales the learning rate

- Can one compute the optimal learning rate?

→ Optimize for  $\alpha$ :

- 1)
- 2)

$$\alpha^* = \underset{x}{\operatorname{argmin}} L(\vec{\theta}^k - \alpha \nabla_{\vec{\theta}} L)$$

$$\vec{\theta}^{k+1} = \vec{\theta}^k - \alpha^* \nabla_{\vec{\theta}} L$$

} In practice, this is too computationally expensive

- How to tackle the computational complexity of naive G.D.?
  - has  $\mathcal{O}(n)$  per training step
- Use **Stochastic Gradient Descent (SGD)** Gradient for  $n$  samples
- + can help escape local minima and saddle points

$$\frac{1}{n} \sum_{i=1}^n L_i(\vec{\theta}, x_i, y_i) = E_{i \sim [1..n]}[L_i(\vec{\theta}, x_i, y_i)] \approx \frac{1}{|S|} \sum_{j \in S} L_j(\vec{\theta}, x_j, y_j)$$

express total loss  
over training data  
as expectation

with  $S \subseteq \{1, \dots, n\}$

→ Minibatches w/ length  $m$   
and  $m \ll n$

- What's a good minibatch size?

Typically power of 2:

8	16	32	64	128
small				large

- ⊕ noisy updates
- ⊕ fast
- ⊕ less noisy
- ⊖ slow

→ limited by GPU memory

- Iteration: Single gradient update step with one minibatch
- Epoch: Complete pass through training set

- ① When does SGD converge?

Robbins & Monro Condition

$$\vec{\theta}^{k+1} = \vec{\theta}^k - \alpha_k H(\vec{\theta}^k, X)$$

estimates  $\nabla F(\vec{\theta}^k)$

converges to a local minimum if:

1)  $x_n \geq 0 \quad \forall n \geq 0$

2)  $\sum_{n=1}^{\infty} x_n = \infty$

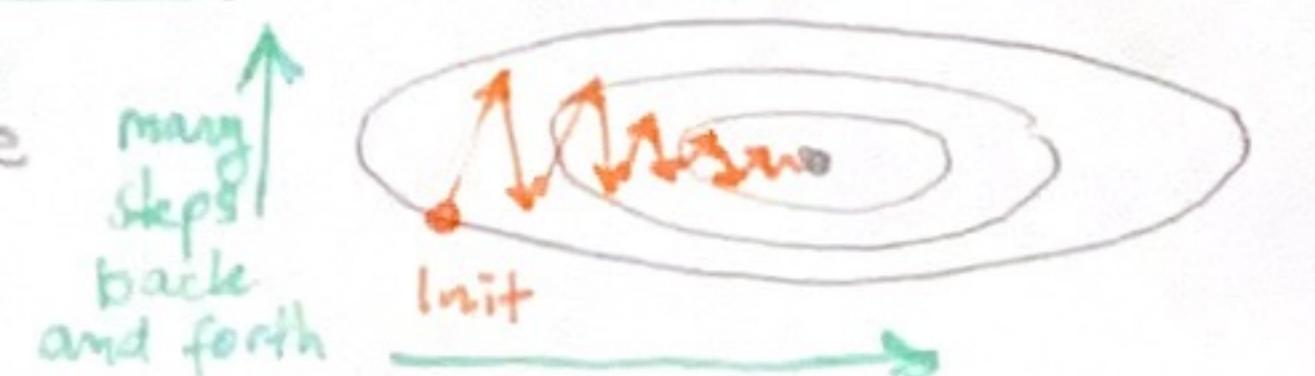
3)  $\sum_{n=1}^{\infty} \alpha_n^2 < \infty$

If also 4)  $F(\vec{\theta})$  is convex  $\rightarrow$  global minimum

- What are problems of SGD?

→ Gradient is scaled equally across dimensions

→ Difficult to find a good learning rate



→ Solution to equally scaled gradients?

→ Gradient Descent with Momentum

(6)

① Gradient Descent w/ Momentum:

$$\vec{v}^{k+1} = \beta \vec{v}^k - \alpha \nabla_{\theta} L(\vec{\theta}^k)$$

↑ velocity

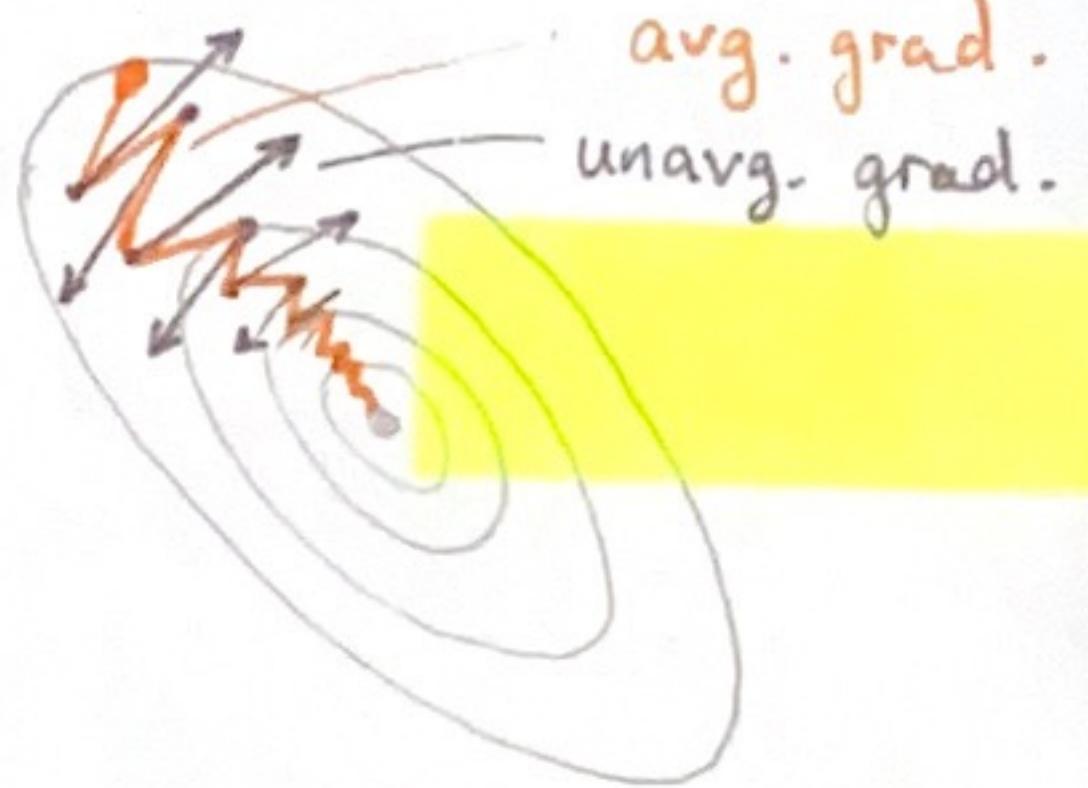
accumulation rate  
(e.g. = 0.9)

$$\vec{\theta}^{k+1} = \vec{\theta}^k + \vec{v}^{k+1}$$

- \* If gradient steps in a dim. always go in the same direction → larger step
- \* If steps go back & forth → smaller steps

=> Exponentially weighted avg. of the gradient

because previous gradients "decay" exponentially with  $\beta$   
(i.e. naive G.D. for  $\beta=0$ )



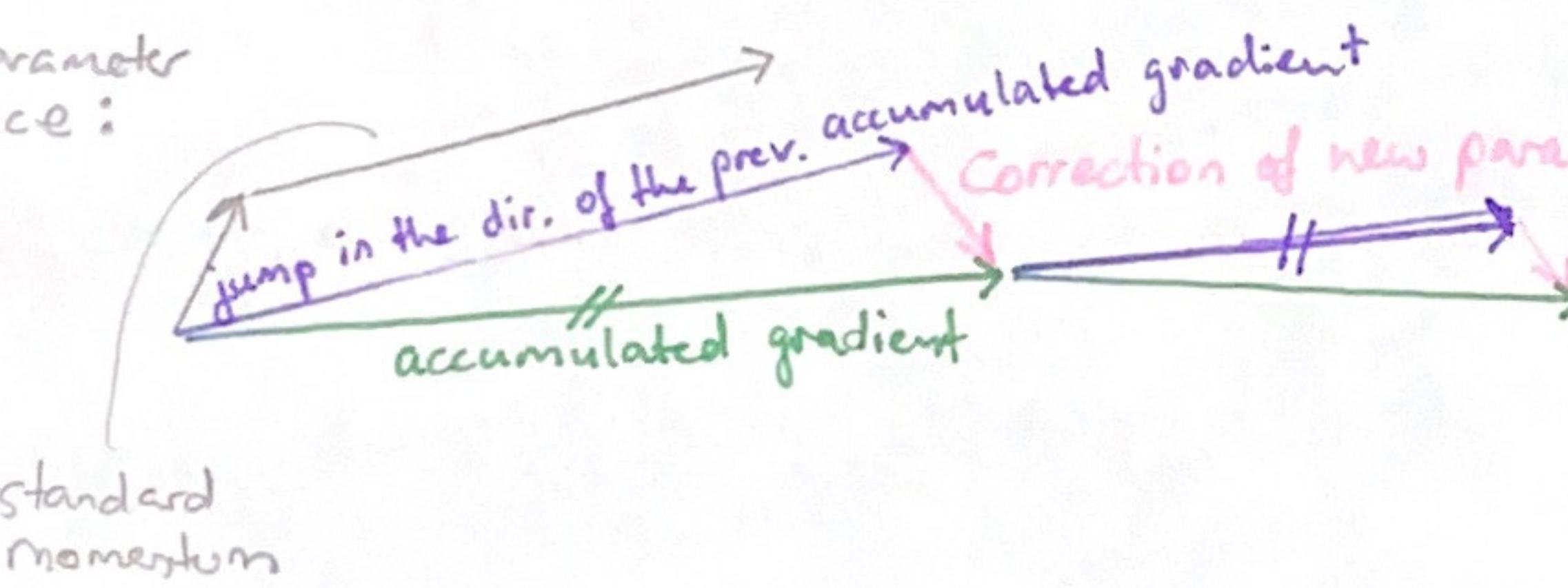
→ it can overcome local minima + accelerate optimization

- What's the Nesterov Momentum? ("Look-ahead momentum")

$$\boxed{\vec{v}^{k+1}} = \boxed{\beta \vec{v}^k} - \boxed{\alpha \nabla_{\theta} L(\vec{\theta}^{k+1})} \quad \text{with } \vec{\theta}^{k+1} = \vec{\theta}^k + \vec{v}^{k+1}$$

↳ Motivation: Computing the gradient at the next look-ahead-step gives us a better gradient

In parameter space:



Root Mean Squared Propagation

- What's RMS Prop?

→ If gradients in a dim. have a high variance, the  $\alpha$  are damped

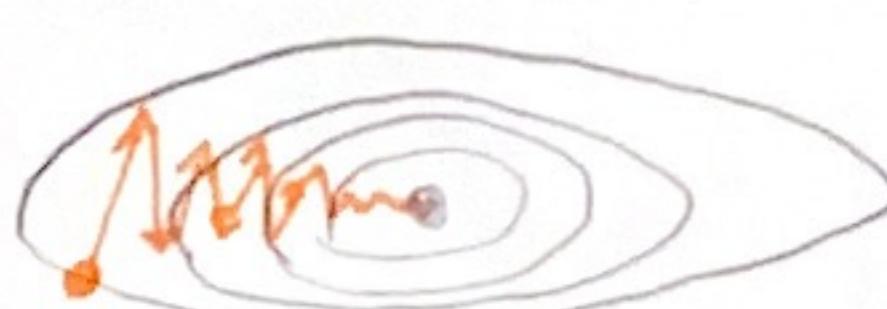
$$\vec{s}^{k+1} = \beta \cdot \vec{s}^k + (1-\beta) [\nabla_{\theta} L \circ \nabla_{\theta} L]$$

$$\vec{\theta}^{k+1} = \vec{\theta}^k - \alpha \cdot$$

$$\frac{\nabla_{\theta} L}{\sqrt{\vec{s}^{k+1}} + \epsilon}$$

Element-wise operations  
e.g.  $10^{-8}$

divide learning rate by exp. decaying avg. of squared gradients



Divide by squared gradients:

• In Y-direction: divide by large number

• In X-direction: divide by small number

⇒ allows for larger learning rates

## - What's Adaptive Momentum estimation? : (ADAM)

→ Combines Momentum and RMSProp

$$1. \text{ Momentum: } \hat{m}^{k+1} = \beta_1 \cdot \hat{m}^k + (1 - \beta_1) \nabla_{\theta} L(\vec{\theta}^k)$$

$$2. \text{ Momentum: } \hat{v}^{k+1} = \beta_2 \cdot \hat{v}^k + (1 - \beta_2) [\nabla_{\theta} L(\vec{\theta}^k) \circ \nabla_{\theta} L(\vec{\theta}^k)]$$

$$\vec{\theta}^{k+1} = \vec{\theta}^k - \alpha \cdot \frac{\hat{m}^{k+1}}{\sqrt{\hat{v}^{k+1}} + \epsilon}$$

↳ mean of gradients (exp. decaying)  
↳ Variance of gradients (exp. decaying)  
↳ Needs tuning

$$\text{At the start: } \hat{v}^0 = \hat{m}^0 = \vec{\theta}$$

For stable training, we need "bias correction" (explained on StackOverflow)

Actual update rule in ADAM

$$\Rightarrow \hat{m}^{k+1} = \frac{\hat{m}^k}{1 - \beta_1} \quad \text{and} \quad \hat{v}^{k+1} = \frac{\hat{v}^k}{1 - \beta_2} \text{ — in the } (k+1)\text{th step}$$

$$\Rightarrow \vec{\theta}^{k+1} = \vec{\theta}^k - \alpha \cdot \frac{\hat{m}^{k+1}}{\sqrt{\hat{v}^{k+1}} + \epsilon} \text{ — to the power of } k+1$$

## - What's Newton's Method?

→ Approx. loss func. by second-order derivative (taylor expansion)

$$L(\vec{\theta}) \approx L(\vec{\theta}_0) + (\vec{\theta} - \vec{\theta}_0)^T \nabla_{\vec{\theta}} L(\vec{\theta}_0) + \frac{1}{2} (\vec{\theta} - \vec{\theta}_0)^T H(\vec{\theta} - \vec{\theta}_0)$$

Diff. and set equal to 0

$$\Rightarrow \vec{\theta}^* = \vec{\theta}_0 - H^{-1} \nabla_{\vec{\theta}} L(\vec{\theta})$$

Hessian replaces the learning rate!

→ Scales each gradient dim. by the curvature

Example:



→ Not used in DL, because of high complexity (inverting is  $\mathcal{O}(k^3)$ )

$\Rightarrow$  Quasi-Newton methods approximate  $H^{-1}$

↳ BFGS & L-BFGS ( $\mathcal{O}(n^2)$  &  $\mathcal{O}(n)$ )

→ Gauss-Newton

→ Levenberg

→ Levenberg-Marguardt

## - Which optimizer to use:

Standard: Adam

Fallback: SGD with momentum

→ Newton, L-BFGS, ... don't work on minibatches

- Which (general) optimization algorithms exist?

7

• for lin. Systems ( $Ax = b$ )

- LU, QR, Cholesky, Jacobi, Gauss-Seidel, CG, PCG

• Non-lin. (gradient-based)

- Newton, Gauss-Newton, (L) BFGS

← second order derivative

- GD, SGD

← first order derivative

• Others

- Genetic algorithms, ...

### Universal Approximation Theorem:

Any function on a compact set can be approximated within an arbitrary error  $\epsilon > 0$ , using a single-layer neural network (1 hidden layer + sigmoid)

↳ (but the layer may be extremely large)

↳ inefficient use of comp. power

## I2DL 06 : Training NNs

(8)

- Learning rate decay: → decrease step size during training to converge to optimum

different Options

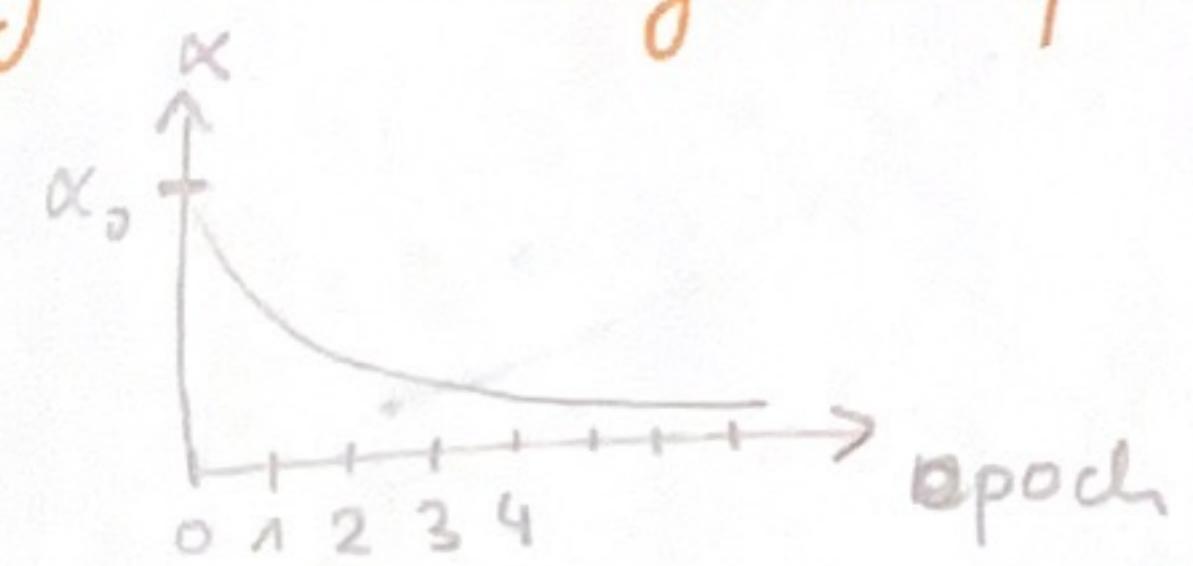
Step decay

$$\alpha = \alpha_0 - t \cdot \alpha \quad (\text{apply only every } n \text{ steps})$$

exp. decay

$$\alpha = \frac{\alpha_0}{t + epoch} \quad t < 1$$

$$\alpha = \frac{t}{epoch} \cdot \alpha_0$$



- Training Schedule (specify  $\alpha$  for entire training process)
  - every  $n$ -epochs

- Basic training recipe:

Given: - Dataset  $\{x_i, y_i\}$  with usually  $\dim(x) \gg \dim(y)$   
- Network  $f$  with parameters  $\vec{w}, \vec{b}$

Method: Use SGD to optimize parameters

Learning means generalization to unknown dataset.

Procedure

- Train on training set → Training error is avg minibatch error
- Optimize hyperparameters & check generalization on validation set → Evaluate loss on subset of validation set every  $n$  iterations
- After developing, test performance on test set

① Error:

mislabeled data → Ground truth error  
e.g. 1%

Bias  
(underfitting)  
Deviation

Training set error

e.g. 5 %

too high?

No

- Bigger model
- Train longer
- Model architecture

Variance  
(overfitting)  
Deviation

Val / test set error

e.g. 8 %

too high?

No

- More data
- Regularization
- New model architecture

test error too high?  
No

val error still too high?  
No

Yes

More val set data  
start over

Yes

Yes

- Data augmentation
- Train - Val data mismatch

More val set data  
start over

ML Recipe START

## - Training curves :



## - Hyperparameter tuning methods :

- Manual search
- Grid Search (Select points from parameter spaces within a range at fixed distances between them)
- Random Search (Select random points)

## - Various Tips for DL :

1. Start with a single sample to check if the model works and overfits (train accuracy should be 100%)
  - Then increase to a few samples → Then move from overfitting to more samples (see generalization)
2. Find learning rate so that loss drops exponentially within 100 iterations
3. Use coarse grid search and then refine grid around best hyperparameters (if the compute budget allows)
4. Get precise timings (find performance bottleneck)
  - An iteration should not exceed 500 ms
    - Dataloading
    - Backprop
  - How long until convergence
5. Start with the simplest network possible
  - Rule of thumb : # of layers divided by 5
6. Make one change at a time.

## I2DL 07: Training NNs 2

9

- Loss functions:

•  $L^1$ -Loss :

$$\sum_{i=1}^n |y_i - f(\vec{x}_i)|$$

+ Robust to outliers  
- costly to optimize  
→ optimum is median

•  $L^2$ -Loss :

$$\sum_{i=1}^n (y_i - f(\vec{x}_i))^2$$

+ prone to outliers  
- compute-efficient optimization  
→ optimum is mean

• Cross-Entropy Loss :

$$-\frac{1}{N} \sum_{i=1}^n \sum_{k=1}^K [y_{ik} \cdot \log \left( \frac{e^{s_{ik}}}{\sum_{k'=1}^K e^{s_{ik'}}} \right)]$$

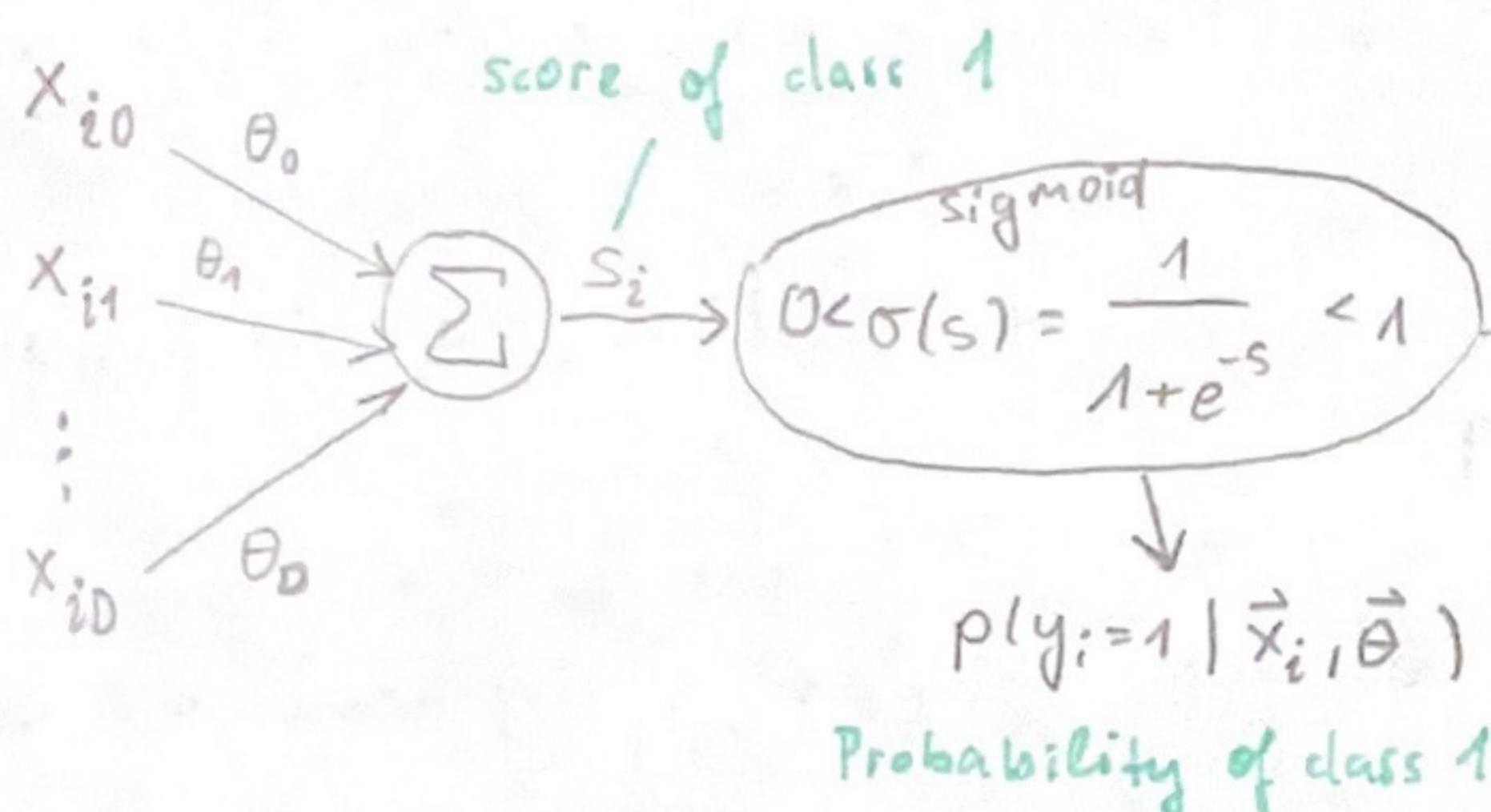
} general  
derivation

used in  
the lecture

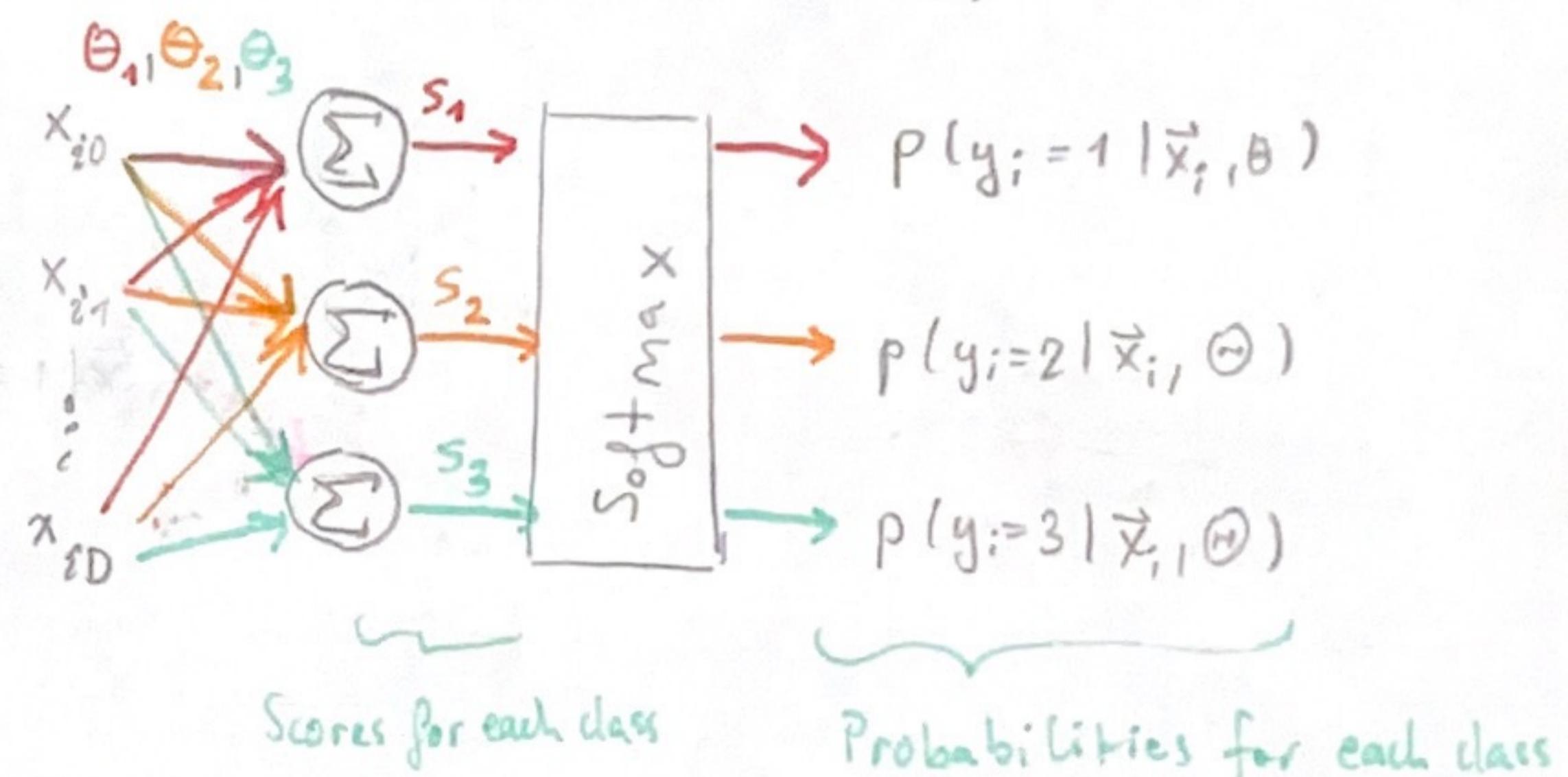
or  $L_i = -\log \left( \frac{e^{s_{y_i}}}{\sum_{k=1}^K e^{s_k}} \right)$  for every class  $y_i$   
and use only the  $L_i$  from the correctly  
classified class.

Network Outputs:

• Binary classification ( $k=2$ )



• Multiclass Classification ( $k>2$ )



• Hinge Loss :

$$L_i = \sum_{k \neq y_i} \max(0, s_k - s_{y_i} + 1)$$

score of the "label class"

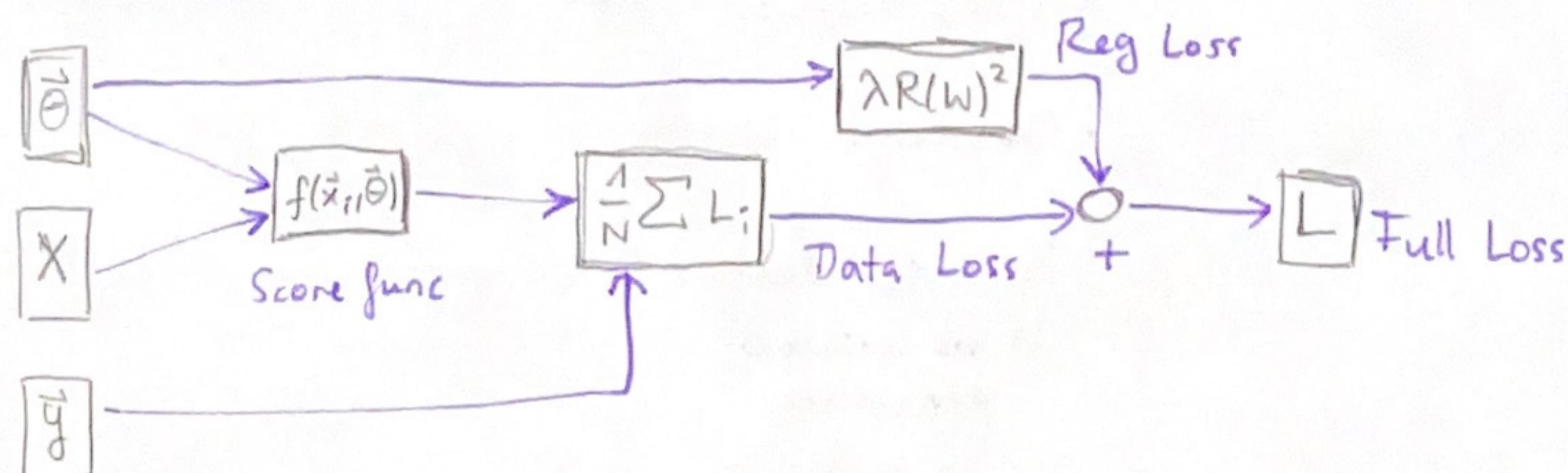
Once the NN is very sure about its decision,  
the loss becomes 0 and learning stops  
⇒ "The hinge loss saturates"  
↳ stops optimizing (as opposed to CE)

• Regularization-Loss :

$$\frac{1}{N} \sum_{i=1}^N L_i + \lambda R^2(W)$$

Data Loss

Reg Loss



- Activation functions:

- Sigmoid:
  - $\sigma(s) = \frac{1}{1 + e^{-s}}$
  - not really used
- 1st problem
  - In the backward pass:  $\frac{\partial L}{\partial s} = \frac{\partial \sigma}{\partial s} \cdot \frac{\partial L}{\partial \sigma} \approx 0$
  - ⇒ Saturated neurons kill the gradient flow
  - ⇒ Sigmoids have a small active region / not ideal
  - \*  $\sigma(s) > 0 \wedge s$  (sigmoid is always positive)
  - $\Rightarrow \frac{\partial L}{\partial w_1} = \left[ \begin{array}{c|c} \frac{\partial L}{\partial y} & \frac{\partial y}{\partial s} \\ \hline \frac{\partial y}{\partial s} & \frac{\partial s}{\partial w_1} \end{array} \right] \cdot \left[ \begin{array}{c|c} \frac{\partial s}{\partial w_1} \\ \hline \end{array} \right]$
  - $\frac{\partial L}{\partial w_2} = \left[ \begin{array}{c|c} \frac{\partial L}{\partial y} & \frac{\partial y}{\partial s} \\ \hline \frac{\partial y}{\partial s} & \frac{\partial s}{\partial w_2} \end{array} \right] \cdot \left[ \begin{array}{c|c} \frac{\partial s}{\partial w_2} \\ \hline \end{array} \right]$
  - either  $> 0$  or  $< 0$
  - or sigmoid input from previous layer ( $> 0$ )
  - $> 0$  for all positive data.
- 2nd problem
  - and assume all positive data  $x_1, x_2 > 0$
  - allowed gradient update directions
  - All
    - ⇒ Gradients are simultaneously positive or neg.
    - \* The hypothetical optimal parameters are only reached via a zigzag path
    - also
      - ⇒ For this reason we use zero-centered data! (and a zero-centered activation)
- Tanh:
  - used in RNNs
  - still saturates  $\times$  (→ kills gradient flow)
  - is zero-centered ✓
- ReLU:
  - standard choice
  - dead ReLU → zero grad
  - large & consistent gradients
  - Initialising ReLUs with slightly pos. bias makes it likely for them to stay active
- Rectified Lin. unit
  - second choice to ReLU
  - Generalization of ReLUs ✓
  - Lin. regimes (nice grads) ✓
  - Does not die ✓
  - Does not saturate ✓
  - Increases the # of params  $\times$
  - Leaky ReLU:  $f(x) = \max(0.01x, x)$
  - Parametric ReLU:  $\sigma(x) = \max(\alpha x, x)$
  - Maxout Units:  $\text{Maxout} = \max(\vec{w}_1^T \vec{x} + b_1, \vec{w}_2^T \vec{x} + b_2)$
  - ↳ has trainable parameters
  - can approximate convex functions (e.g. Absolut val)
  - example with two functions

- Weight Initialization:

- Why not all zero? → All hidden units compute the same function
  - same gradients for all neurons
  - "no symmetry-breaking" / weights updated equally

→ Why not with small random numbers?

- Cause a small output of each layer when using tanh as an activation func

$$f_l(\sum_i w_i^l x_i^l + b^l) \approx 0 = x_i^{l+1}$$

⇒ The input to the next layer is small, but even if the grad. of the activation func. is  $> 0$  → the grad. of the affine layer are the inputs!

$$\Rightarrow \frac{\partial L}{\partial w_i^{l+1}} = \frac{\partial L}{\partial f_{l+1}} \cdot \frac{\partial f_{l+1}}{\partial w_i^{l+1}} = \frac{\partial L}{\partial f_{l+1}} \cdot x_i^{l+1} \approx 0$$

⇒ vanishing gradients in the affine layers caused by small output

→ Why not with large random numbers?

- When using tanh, the output is saturated to -1 and +1.

⇒ Gradient of the activation func  $\approx 0$

⇒ Vanishing gradient, caused by saturated activation function.

\* Xavier Initialization:

- ⇒ The layer output does not collapse to 0 or saturate.
- ⇒ The gradient flow is not killed

→ Use Gaussian with zero mean and  $\sqrt{\frac{1}{n}}$  # of input neurons

$$\text{Var}(s) = \text{Var}\left(\sum_{i=1}^n w_i x_i\right) = \sum_{i=1}^n \text{Var}(w_i x_i)$$

$$= \sum_{i=1}^n E[w_i^2] \text{Var}(x_i) + E[x_i^2]^2 \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i)$$

↓  
0, zero-mean      ↓  
0, zero-mean

i.i.d.

$$= n (\text{Var}(w) \text{Var}(x))$$

When using ReLU: since half of the inputs are set to zero

$$\text{Var}(w) = \frac{2}{n}$$

Xavier-haff initialization

⇒ We want the Variance of the output equal the var. of the input

$$\text{Var}(s) = \text{Var}(x) \Rightarrow$$

$$\text{Var}(w) = \frac{1}{n} (\text{wtanh})$$

## I2DL 08: Training NNs

(11)

- Data Augmentation: "is regularisation"

→ A classifier should be invariant to many transformations (pose, appearance, illumination of obj)

⇒ Synthesize data simulating such plausible transformations.

- Methods {
- \* Random brightness and contrast changes
  - \* Random crops of larger image (→ fixed crops in testing)
  - \* Flipping

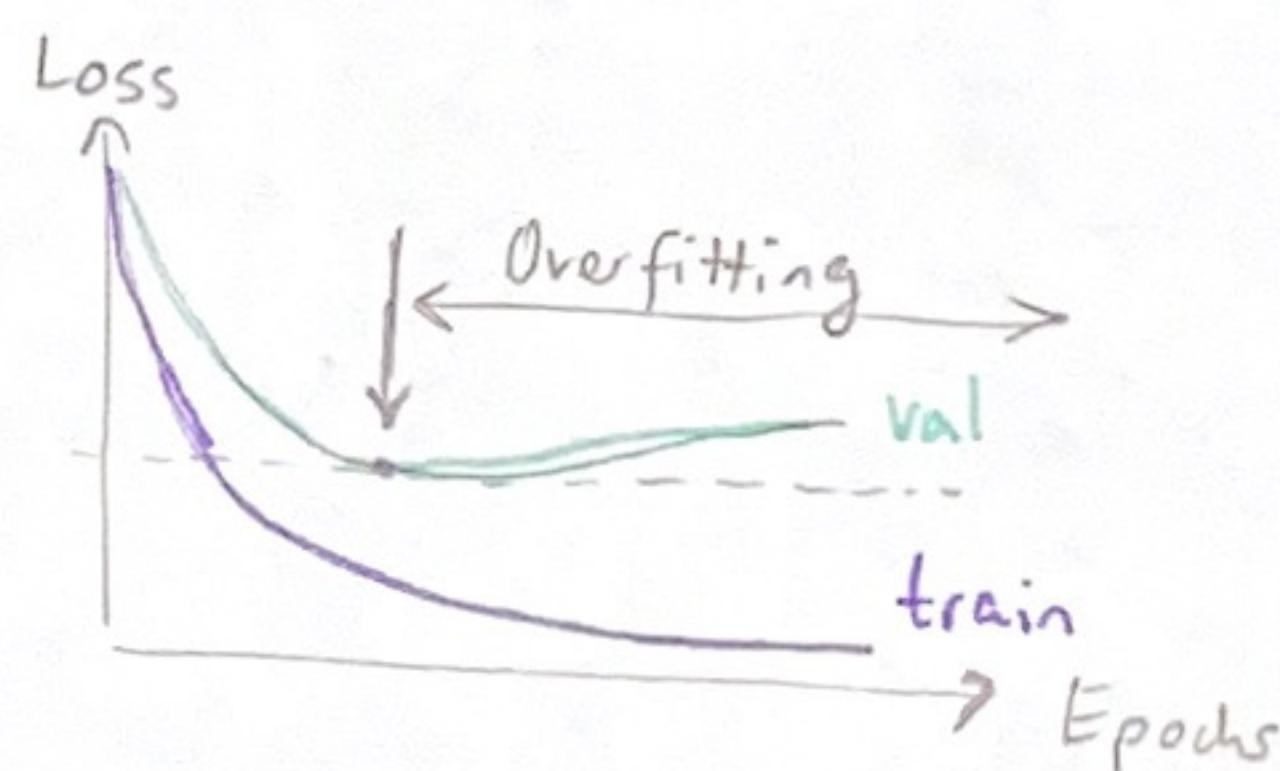
→ Use the same data augmentation when comparing networks

- Advanced Regularization:

\* L<sub>2</sub>-Regularization:  $\Theta_{k+1} = \Theta_k - \varepsilon (\nabla_{\Theta}(\Theta_k, x, y) + \lambda \Theta_k)$

↳ penalizes large weights  
→ Improves generalization

\* Early stopping



\* Bagging & Ensemble methods

Use a different dataset for each model  
(is an ensemble method)

train multiple models and average their results  
→ change optimizer/architecture/...

model ensemble

→ different set of neurons w/ each batch (diff. data), but shared params.

→ If the errors are uncorrelated, the expected combined error will decrease lin. w/ the ensemble size.

\* Dropout → Disable a set of random neurons (e.g. 50%)

↳ In the forward pass: fewer neurons impact the decision } half the capacity  
(→ fewer information → harder)

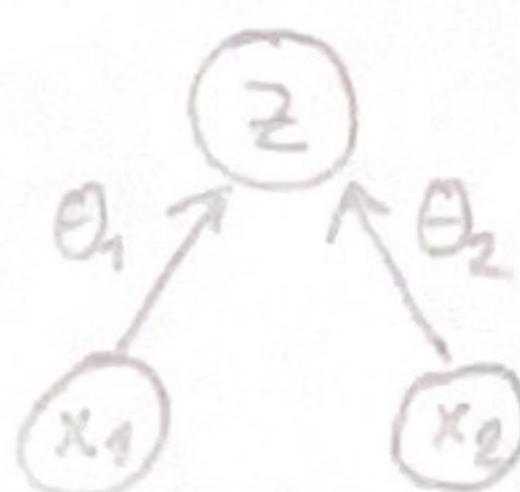
→ In the backward pass: Gradients are only applied to active neurons

reduces co-adaptation between neurons / scores based on more features

→ creates redundant representations (e.g. neurons for detecting tails are not active, so I create a new one)

- \* At test time the conditions in the NN during training time must be "simulated"
- Weight scaling inference rule

During testing:



$$\begin{aligned} E[z] &= \frac{1}{4} (\theta_1 \cdot 0 + \theta_2 \cdot 0) + \frac{1}{4} (\theta_1 x_1 + \theta_2 \cdot 0) \\ &\quad + \frac{1}{4} (\theta_1 \cdot 0 + \theta_2 x_2) + \frac{1}{4} (\theta_1 x_1 + \theta_2 x_2) \\ &= \left(\frac{1}{2}\right) \theta_1 x_1 + \theta_2 x_2 \end{aligned}$$

Dropout probability  $p = 0.5$

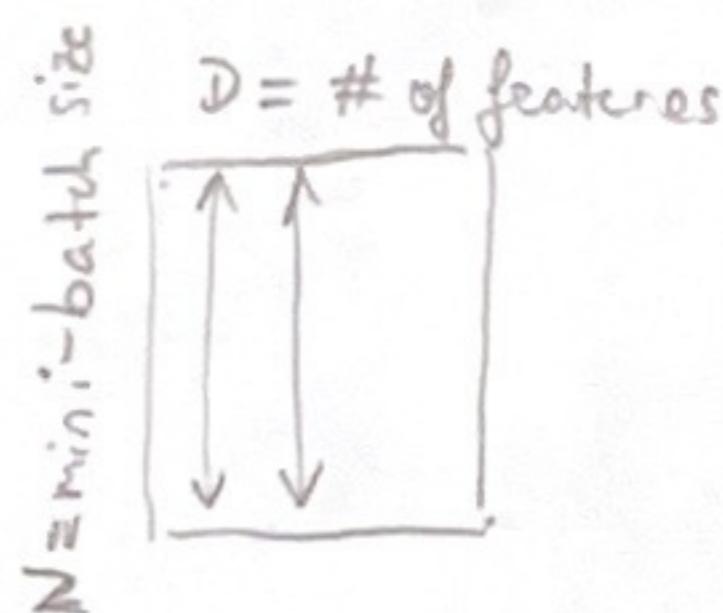
⇒ At test time the weights need to be adapted to match the conditions at training time

→ Dropout improves accuracy but reduces model capacity

## \* Batch Normalization:

→ Scale mini-batches to unit gaussian in each feature/dim.

$$\begin{aligned} E[\vec{x}] &= \frac{1}{N} \sum_{i=1}^N x_i \\ \text{Var}[\vec{x}] &= \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 \end{aligned}$$



$$\hat{\vec{x}}^{(k)} = \frac{\vec{x}^{(k)} - E[\vec{x}^{(k)}]}{\sqrt{\text{Var}[\vec{x}^{(k)}]}}$$

\* BN layer is applied between affine layer and activation

↳ 1. Normalize — " —

2. Allow netw. to change the range  
→ possibility to undo normalization

$$\hat{y}^{(k)} = g^{(k)} \hat{\vec{x}}^{(k)} + \beta^{(k)}$$

optimized in backprop

\* Empirical results show, it's ok to treat dims separately

\* One can set all biases before BN to zero, as they're wasted

⇒ At test time use a mean and variance obtained from an exponentially weighted running average during training

### Benefits:

- More stable gradients w/ deeper nets
- larger range of good hyperparams

### Drawbacks:

- cannot be used w/ small batches → Group norm instead of Batch norm

since we often don't test a mini-batch, but single samples

## 12DL 08: Convolutional neural networks

(12)

- \* FC layers with image data have too many weights ("brute force approach")
  - share the weights for different parts of the image ("use structured approach")

- What are convolutions? "Applying a filter to a function"

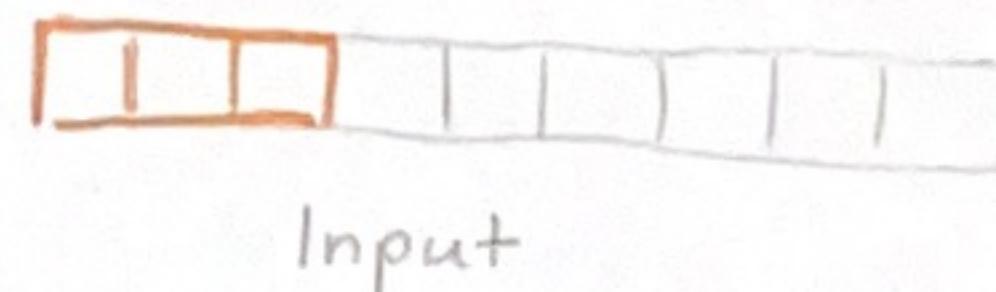
Continuous:  $f * g = \int_{-\infty}^{\infty} f(\tau)g(t-\tau) d\tau :=$  "the area underneath the product"

e.g.

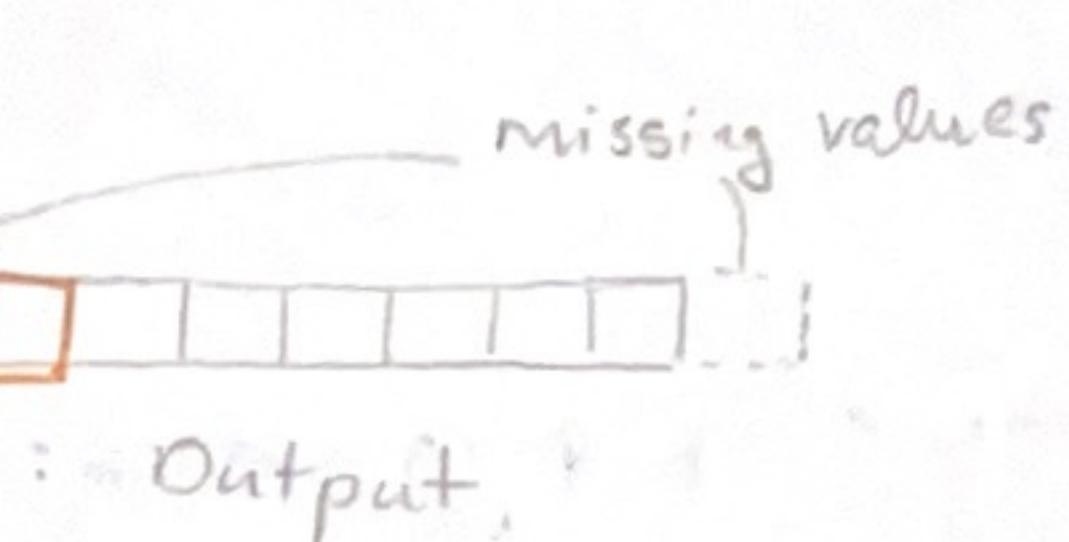
(2)  
in conv layers  
• flip kernel  
• slide kernel / shift kernel w/ t  
① multiply flipped and shifted kernel w/ func

Discrete:

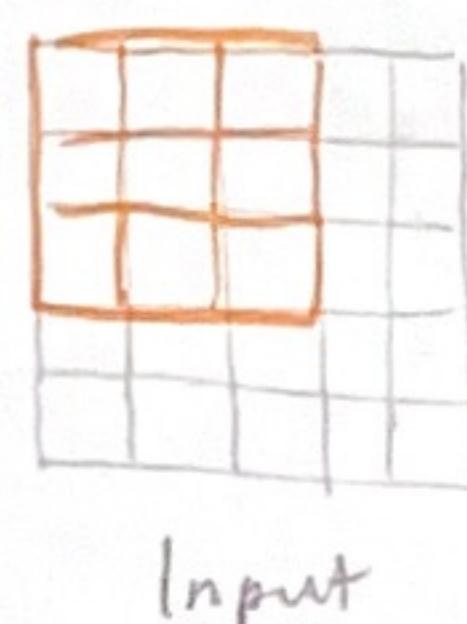
1D:  
(w/o padding)



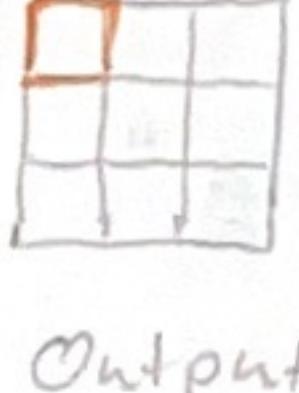
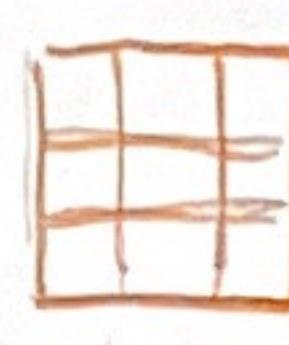
Convolve →



2D:  
(w/o padding)

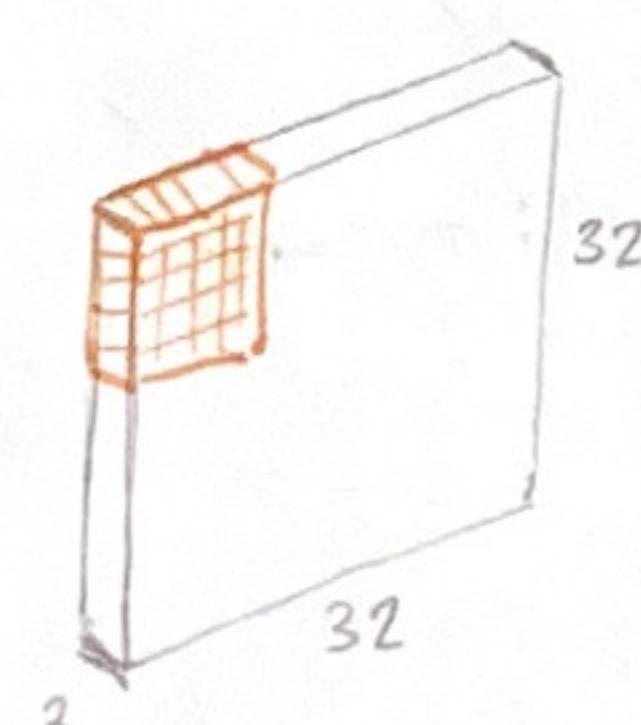


Convolve →

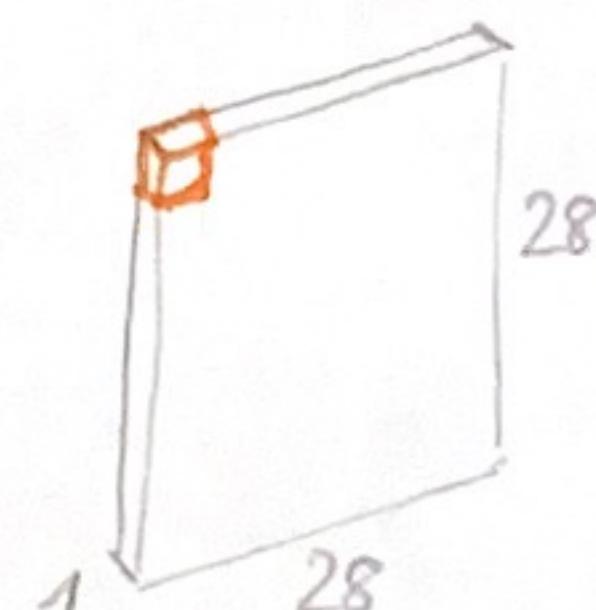
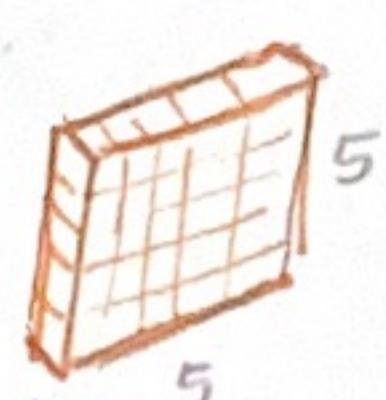


Kernel

RGB images  
3D:  
(w/o padding)



Convolve →



One channel / ~~volume~~  
One Activation map / Feature map

"depth" dim of  
input and filter must match  
(the number of channels)

Using multiple  
filters:

(a convolutional  
layer)

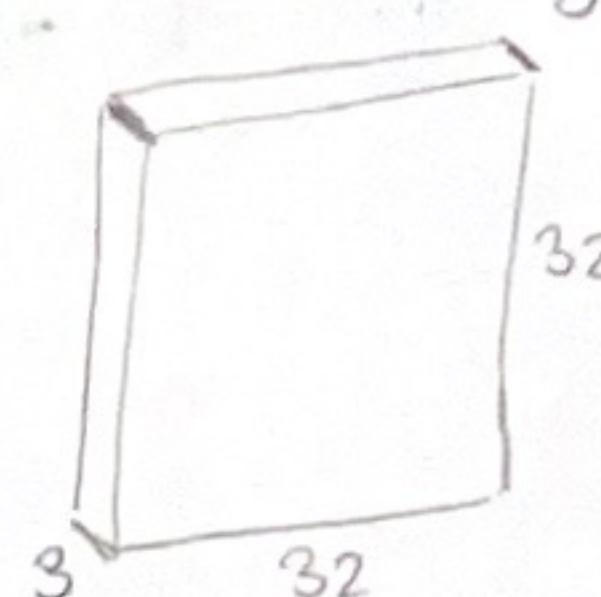
→ Layer is defined by

1. kernel height/width

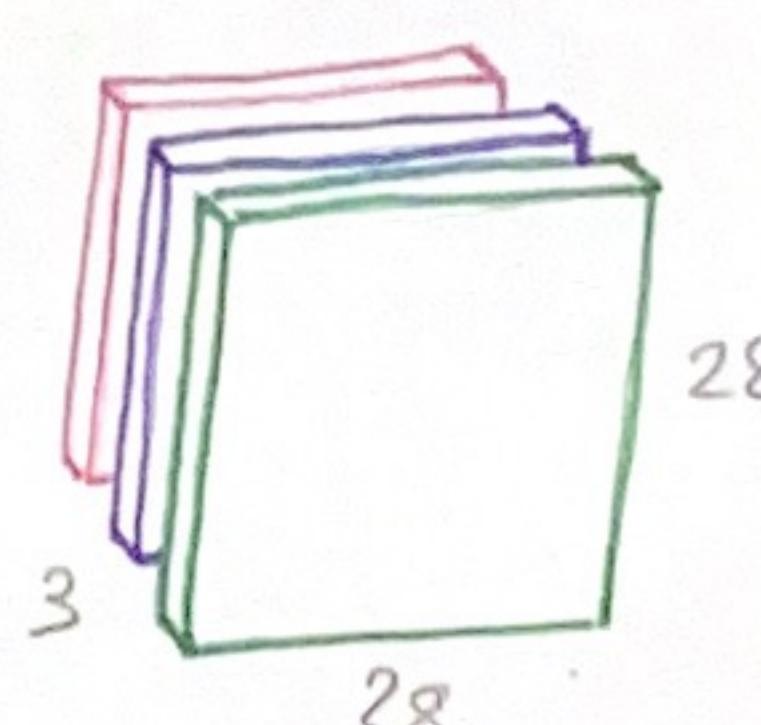
2. number of different kernels

⇒ Each filter captures a different characteristic

32x32x3 image



Convolve  
w/ three  
filters



Activation maps

$28 \times 28 \times 3$

output channels

## - Conv. layer dimensions:

↓  
does "Feature extraction"  
(computes feature in a given region)

Input:  $N \times N$

Filter:  $F \times F \Rightarrow$  Output:  $\left(\frac{N-F}{S} + 1\right) \times \left(\frac{N-F}{S} + 1\right)$

Stride:  $S$

Fractions are illegal

→ certain strides not possible

## - Padding:

- To avoid feature maps shrinking too quickly
- To use corner pixels multiple times

$$\text{output size} = \left(\left\lfloor \frac{N+2P-F}{S} \right\rfloor + 1\right) \times \left(\left\lfloor \frac{N+2P-F}{S} \right\rfloor + 1\right)$$

↓  
floor operation

e.g. zero padding

→ Valid convolution: no padding

→ Same convolution: output = input size  $\Rightarrow P = \frac{F-1}{2}$

## - Number of parameters:

Each filter has  $\underbrace{F \times F \times \text{in-channels}}_{\text{weights}} + \underbrace{1}_{\text{bias}}$  parameters

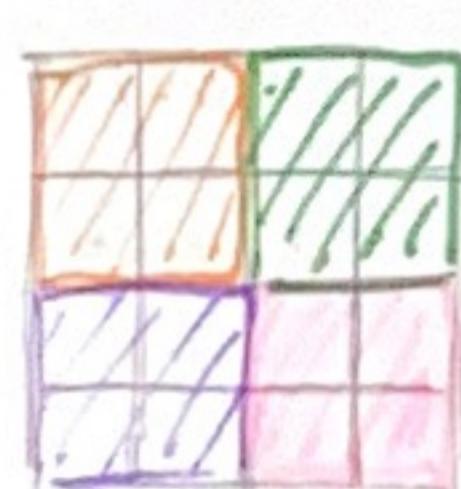
## - Pooling: "Is downsampling"

↳ does "Feature selection" (picks the strongest activation in a region)

Input: size  $W_{in} \times H_{in} \times D_{in}$

Output:  $W_{out} = \frac{W_{in}-F}{S} + 1$  spatial filter extent  
 $H_{out} = \frac{H_{in}-F}{S} + 1$  stride } 2 hyperparameters

$$D_{out} = D_{in}$$



$$F=2$$

$$S=2$$

Max-Pooling

Arg-Pooling  
etc.



typically used deeper in the network

## - Final FC layer:

→ Typically one or two FC layers to make final decision

## - Receptive field: "the input pixels that are connected to a single output pixel"

→ Spatial extent of the connectivity of a convolutional filter

⇒ "Information from larger parts of an image are condensed into a single pixel"

## 12DL 10: CNN-2

(13)

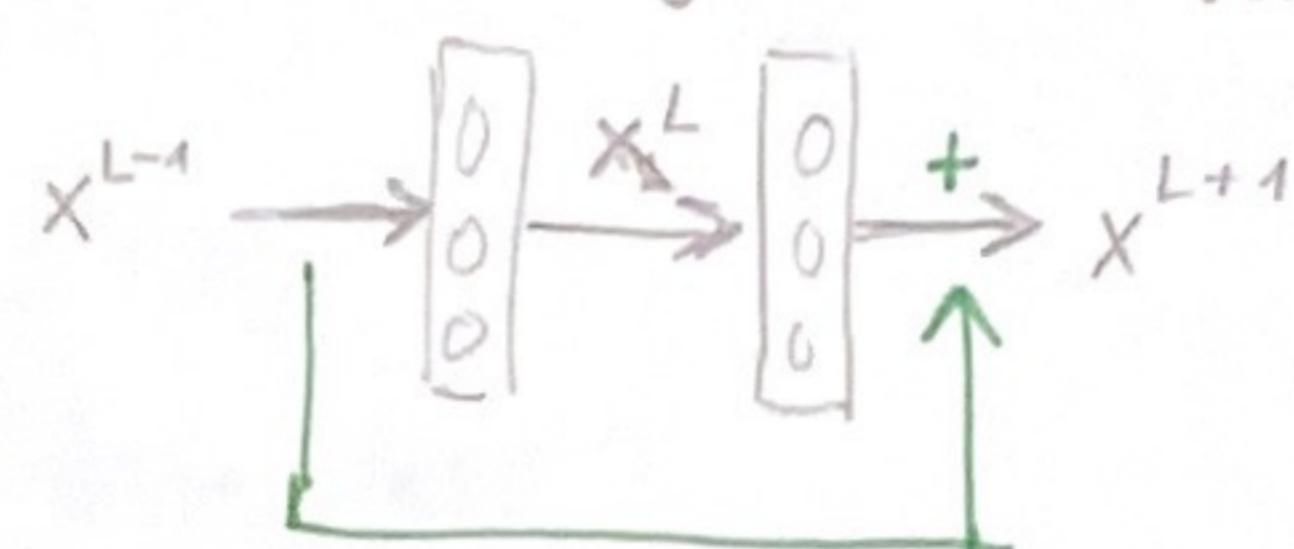
- Common performance metrics in image detection:
  - \* Top-1 score: If the predicted class is the target label
  - \* Top-5 score: If the top 5 predicted classes contain the target
  - \* Top-5 error: percentage of samples for which the top 5 predicted classes did not contain the target

### Skip connections:

\* Problem of vanishing & exploding gradients in deep nets  
(with normal nets, the performance degrades at some depth)

→ Residual Block

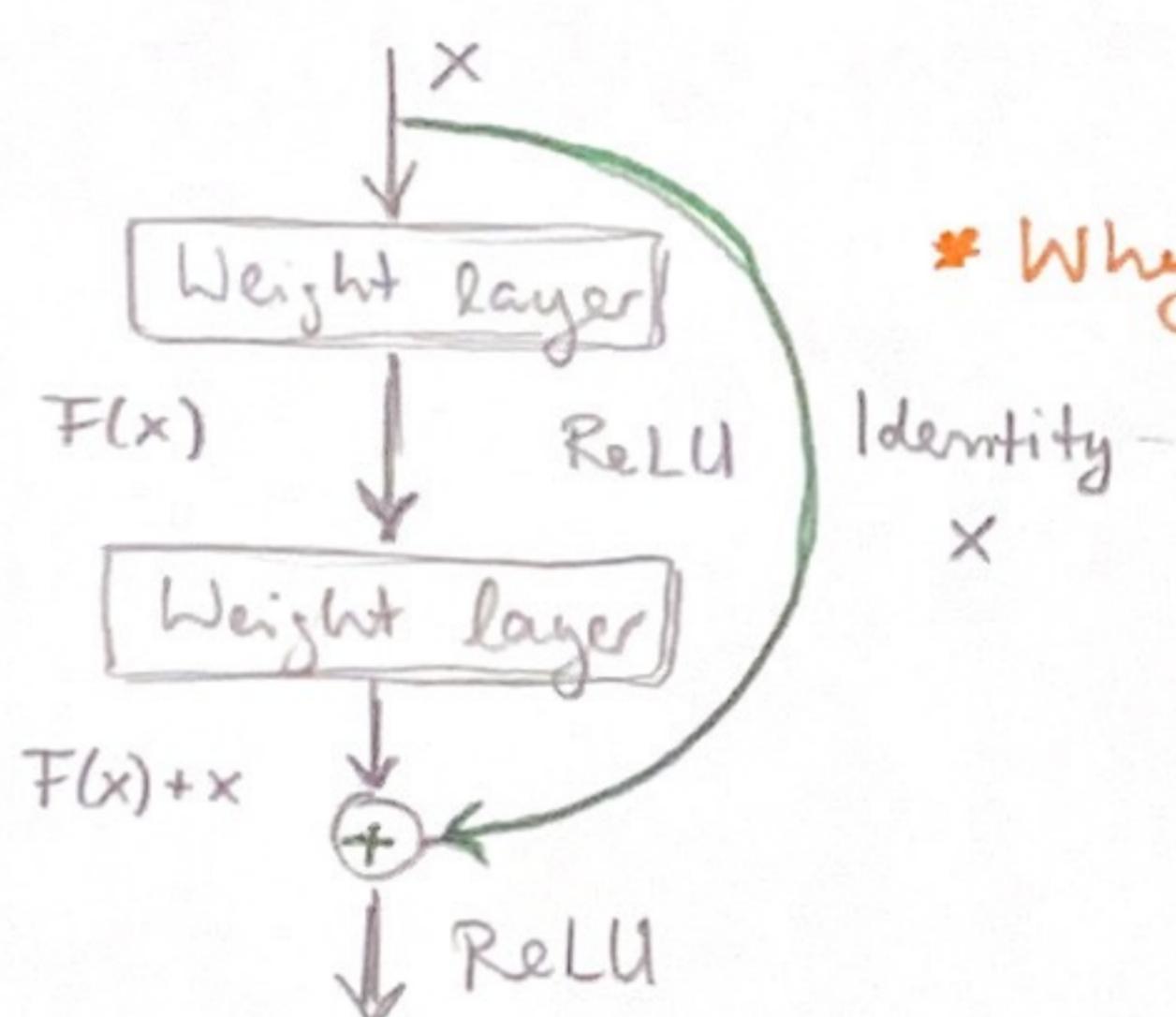
(used in ResNet)



$$\Rightarrow x^{L+1} = f(W^{L+1}x^L + b^{L+1} + x^{L-1})$$

\* The convolution within the block must be SAME, such that the size of  $x^{L-1}$  is the size of  $x^{L+1}$

→ Otherwise convert the dims w/ a matrix of learned weights or zero padding



\* Why do ResNets work?

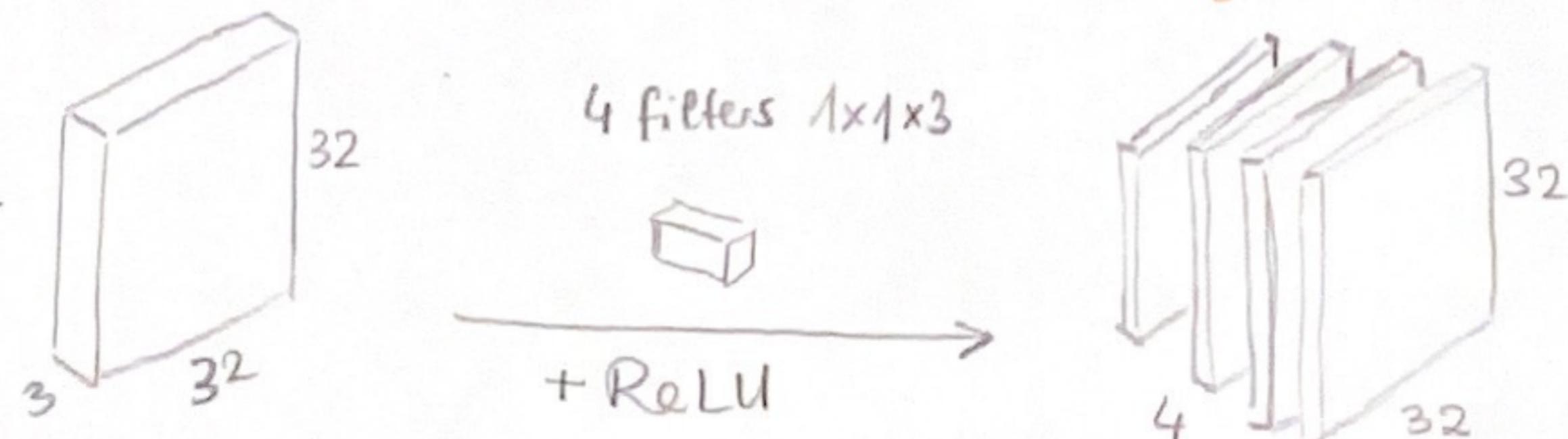
- If the layers in the ResBlock hurt performance, the network can "turn them off" by learning to produce the identity (if all weights & biases are 0)

$\Rightarrow$  Guaranteed to not hurt performance  
(can only improve)

$\Rightarrow$  "Revolution of depth"

### $1 \times 1$ Convolutions:

→  $1 \times 1$  Kernel: Keeps the dims & scales the input

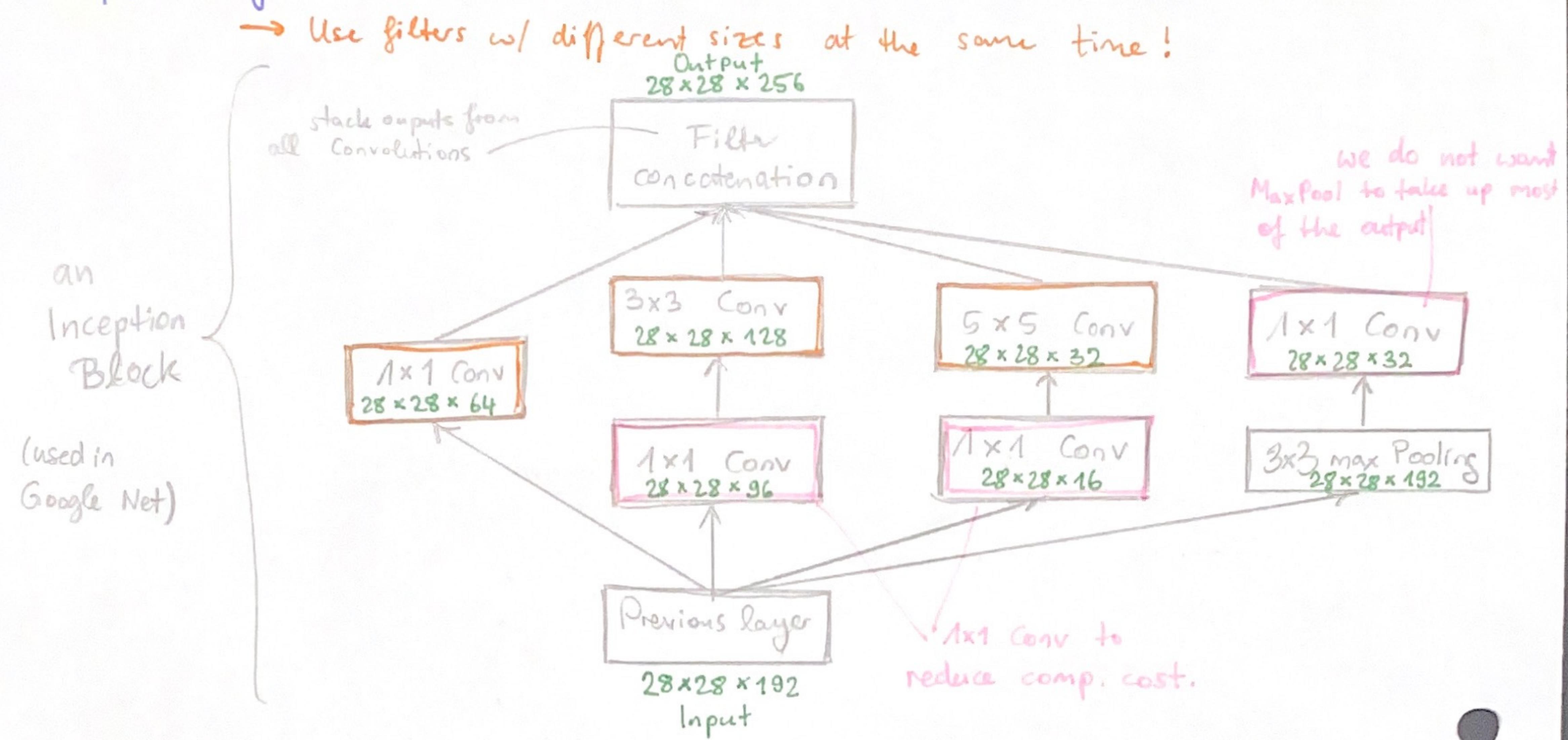


→ Used to shrink the num of channels but keep spatial resolution  
+ Adds non-lin.

→ Subsequent computations are cheaper

More complex func can be learnt

## - Inception layer:

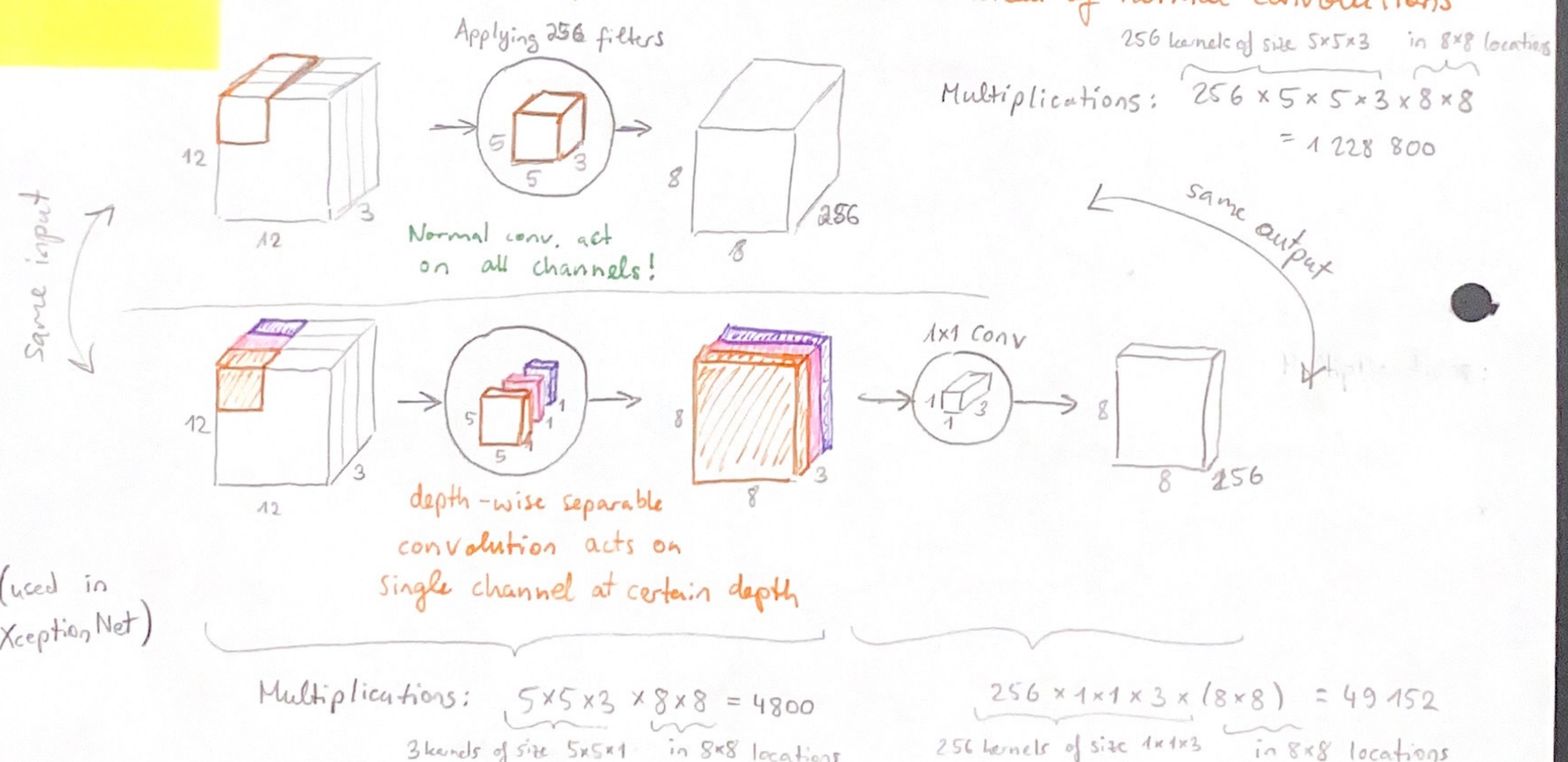


\* Extra max pool layers between Inception Blocks to reduce dimensionality

→ to avoid choosing a kernel size + reduces computational cost

- Xception: "Extreme version of an inception."

→ Uses depth wise separable convolutions instead of normal convolutions



(used in Xception Net)

- Fully convolutional ~~layer~~ network: → Replace lin. layers with equivalent conv. layers

↳ Continuously reduces the width and height

→ In semantic segmentation we want to output the input size

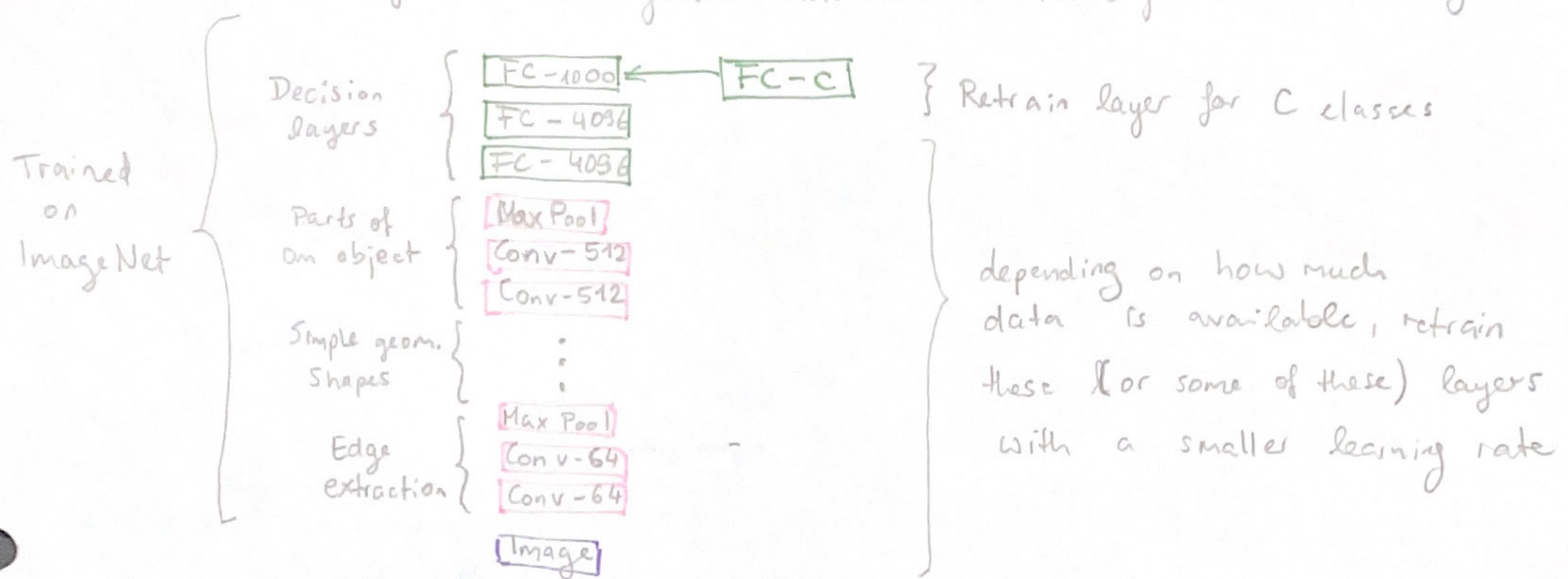
↳ Upsampling: 1) Interpolation (few artefacts)

2) Transposed convolution (unpooling + learned filter)  
→ artefacts

# I2DL 11 : RNNs

(14)

- Transfer learning:
  - Use pretrained model learned for another setting
    - ↳ Finetune it on small dataset
    - (e.g. low-level feature extraction is similar for all RGB images)

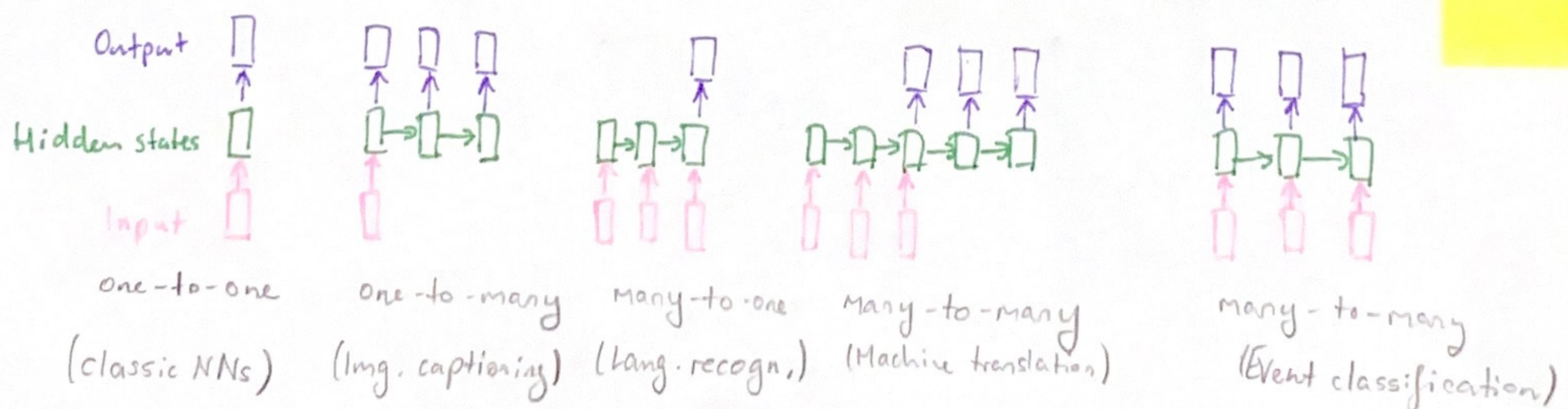


- When to use transfer learning:
  - 1) If both tasks have the same input
  - 2) When we have more data for task 1 than for task 2
  - 3) When low-level features from T1 could be useful for T2

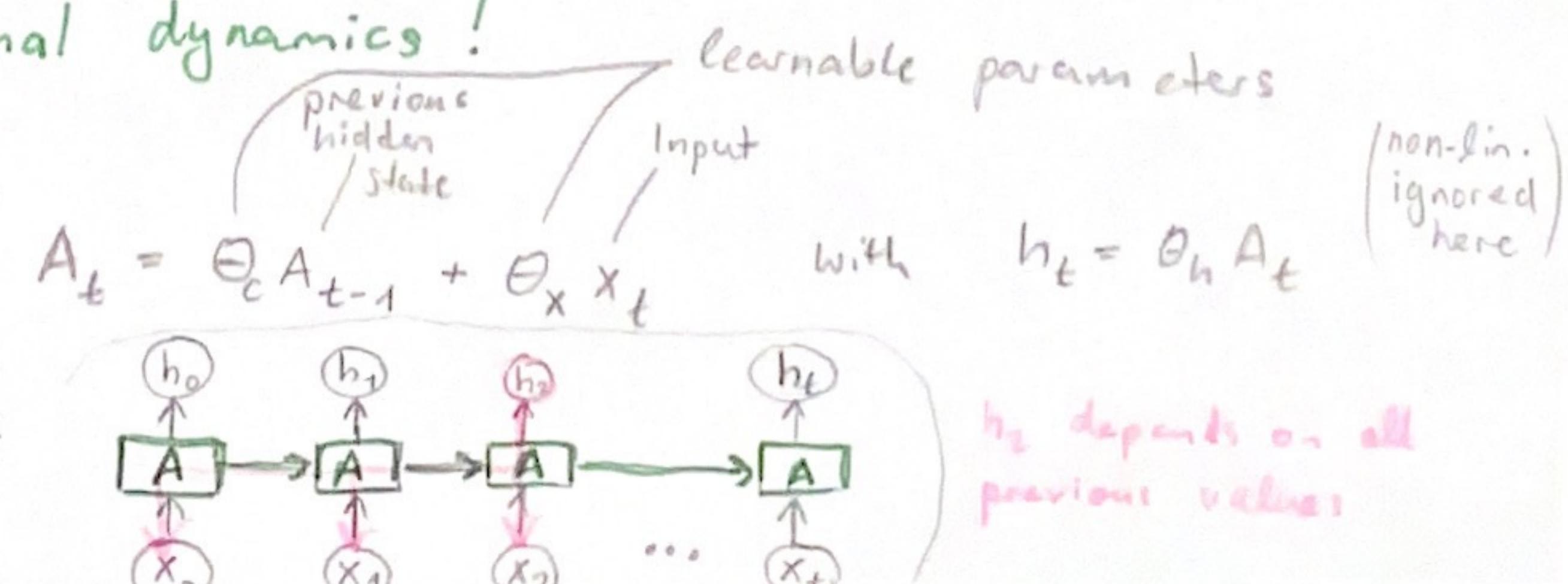
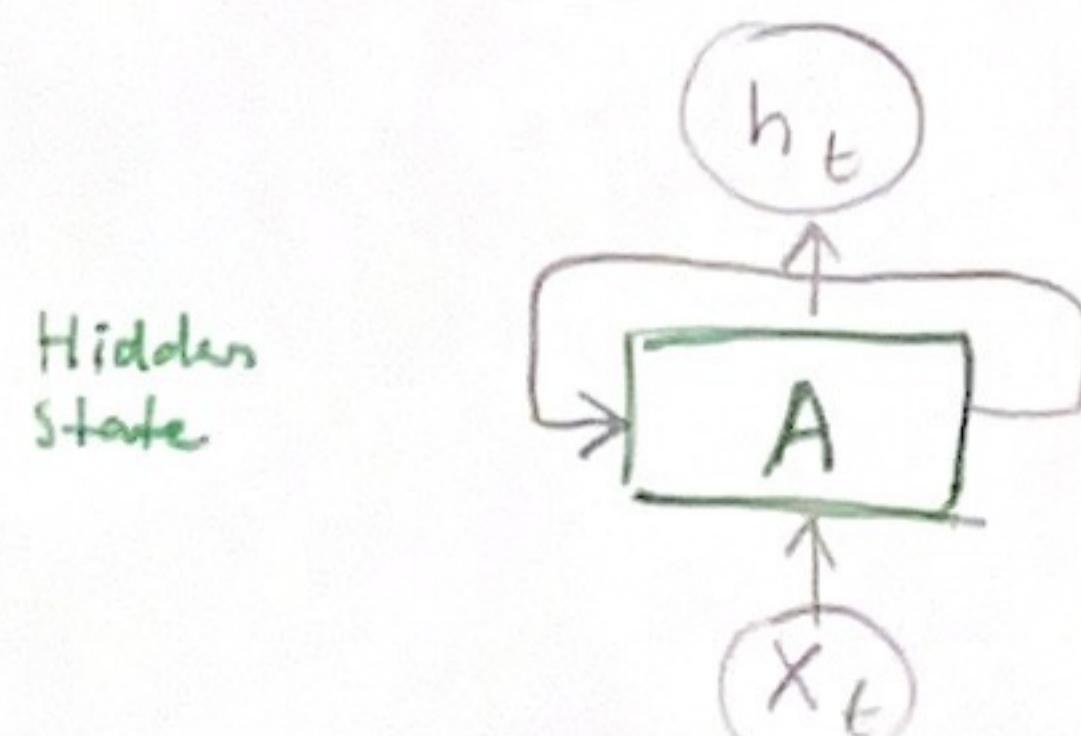
## - Recurrent Neural Networks: RNNs

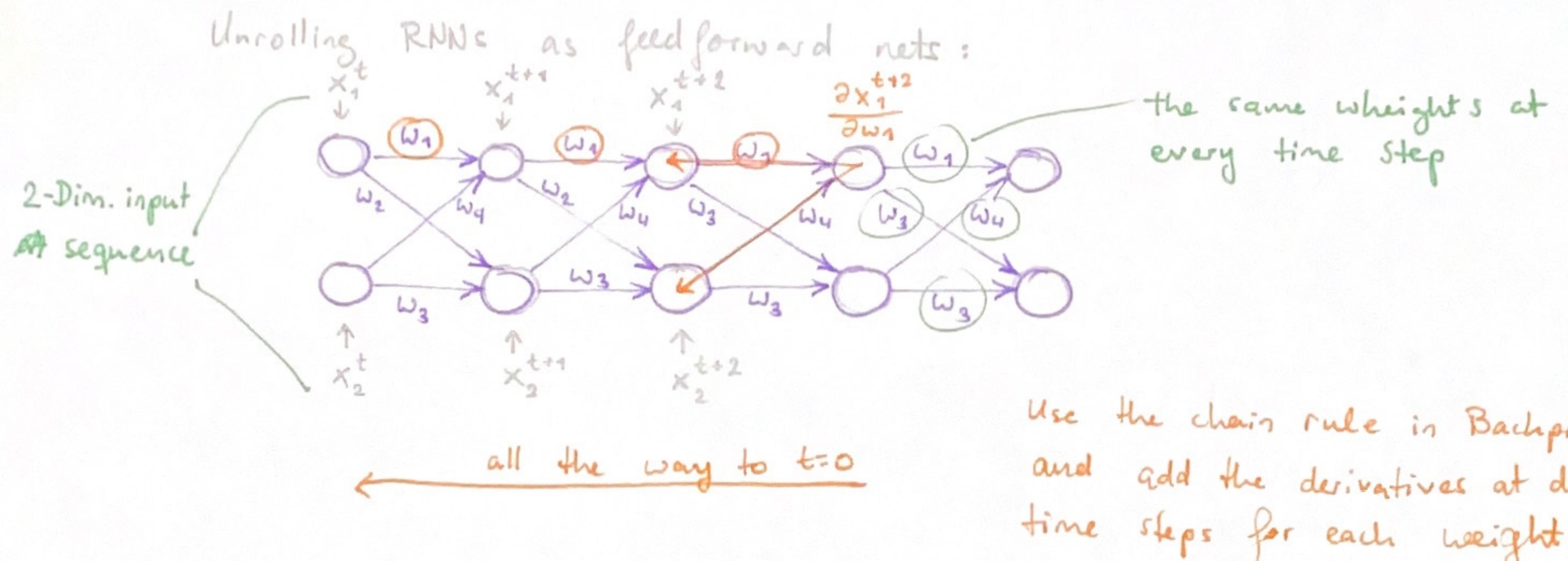
→ Process sequence data (Input/output can be sequences)

Input-Output-mappings:



→ Hidden states have internal dynamics!





\* It is a major challenge to keep enough information in the hidden states for long-term-dependencies (when an output depends on input long before it)

$$A_t = \Theta_c A_{t-1} + \Theta_x x_t \approx \Theta_c^t A_0 \quad (\text{same weights are multiplied over and over again})$$

Eigen-decomposition

$$\Theta = Q \Lambda Q^T$$

$$\Rightarrow A_t = Q \Lambda^t Q^T A_0$$

largest eigenvalues have magnitude

small weights  $\rightarrow$  vanishing grad.  
large weights  $\rightarrow$  exploding grad.

$< 1 \Rightarrow$  vanishing grad  
 $> 1 \Rightarrow$  expl. grad.

grad. clipping

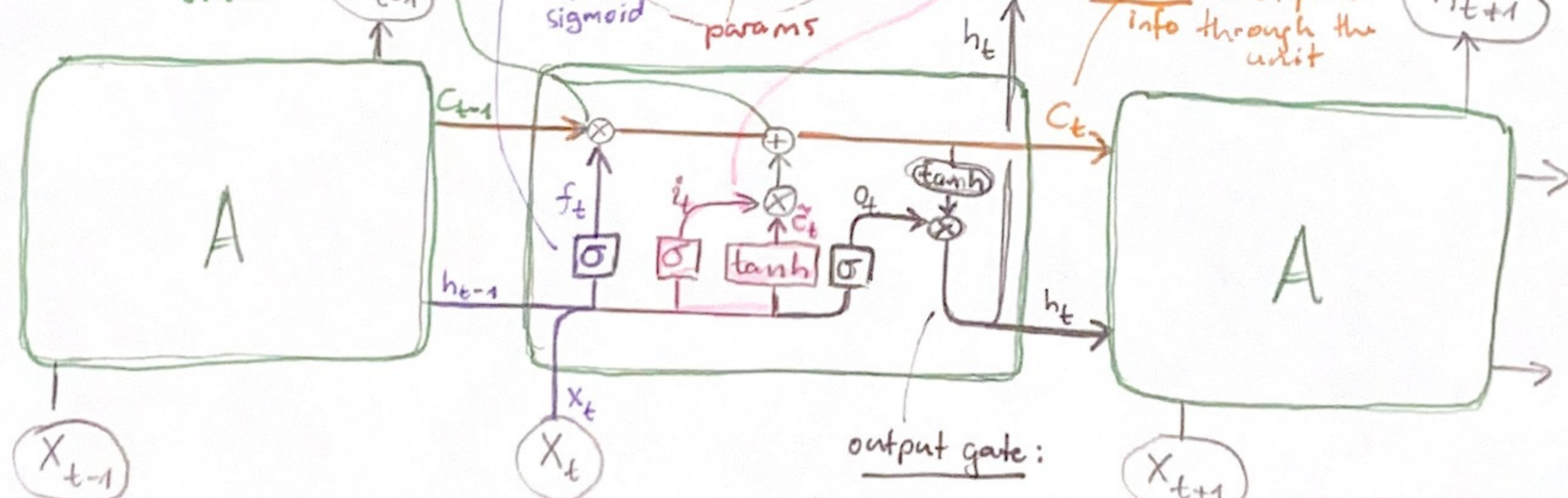
$\Rightarrow$  Set weight / eigenvalues to 1

## - Long-Short-Term-Memory: (LSTM)

element-wise operations:

$$C_t = f_t \odot C_{t-1} + i_t \odot g_t$$

previous state current state



All inputs into a cell must have the same dims

The orange line is a highway for the gradient in back prop.

$f_t$  is 1 for important info

forget gate deals with past information (only here over and over multiplication w/ weights)

only here vanishing grad. from activation

$\Rightarrow$  No problems w/ vanishing gradients

## I2DL 12: Advanced deep learning topics

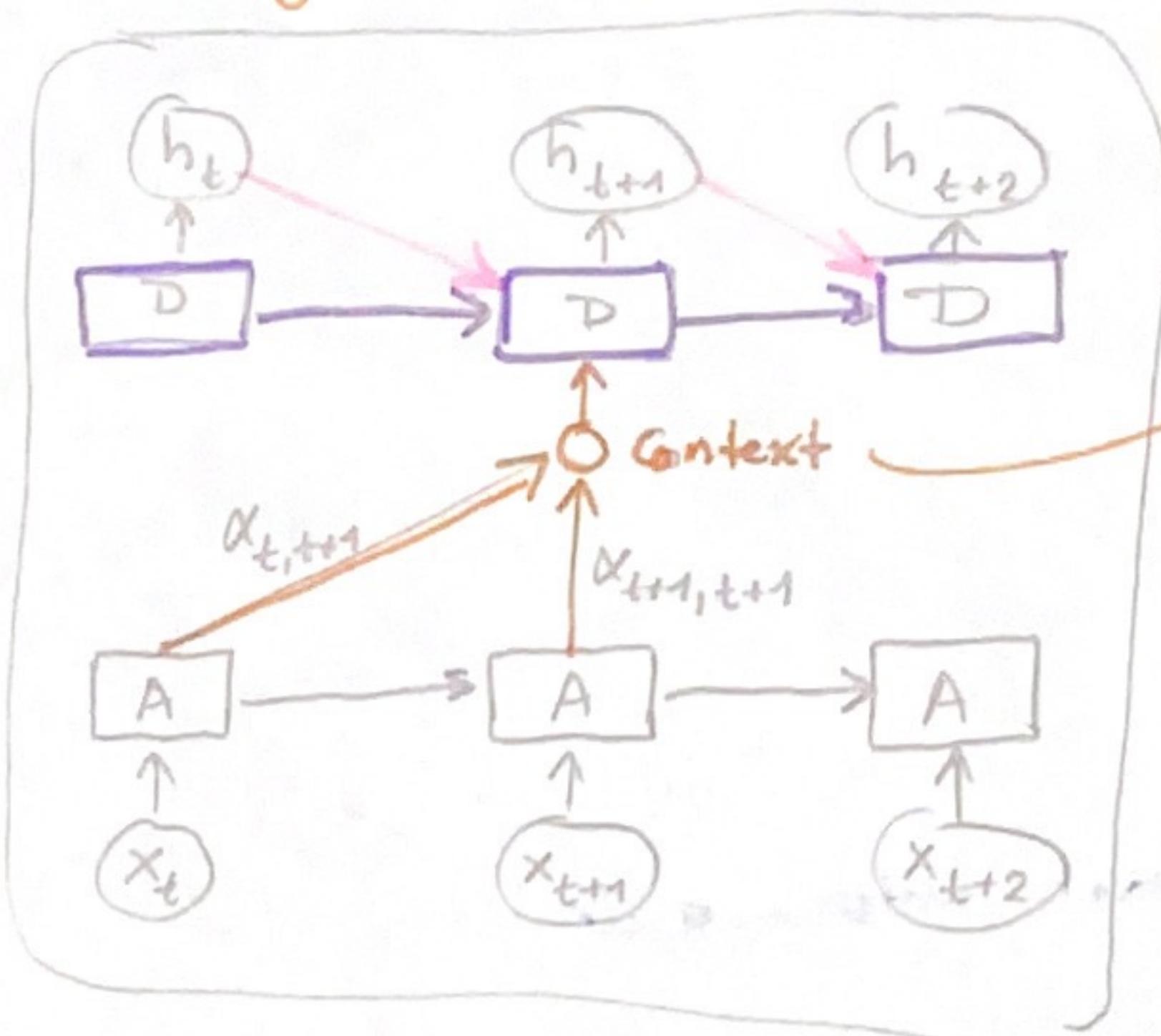
- Attention: → wants to tackle the problem w/ long term dependencies (similar to LSTM)

Architecture:

\* Input to decoders:

- prev. decoder hidden state
- prev. decoder output
- attention/context

=D The decoder then processes the information.



relationship of different time steps w.r.t. the current output at  $t+1$

- Transformers: → get rid of the RNN and only use attention

RNN: No connections between the value at a time step w/ the val. at a later time step

Transformers: All values in the sequence are connected to all other values!

- Graph neural networks:

Node = a concept, Edge = connection between concepts

- \* General idea:
- Input is a graph w/ optional node & edge feature vectors
  - Info is propagated across the graph for several iterations (each iteration is a "hidden layer")
  - Output is a graph w/ updated, context-aware nodes and edge feature vectors.

\* During propagation, we have

- "node to edge" update : 
$$h_{(i,j)}^{(l)} = \mathcal{N}_e([h_i^{(l-1)}, h_{(i,j)}^{(l-1)}, h_j^{(l-1)}])$$
- "edge to node" update :

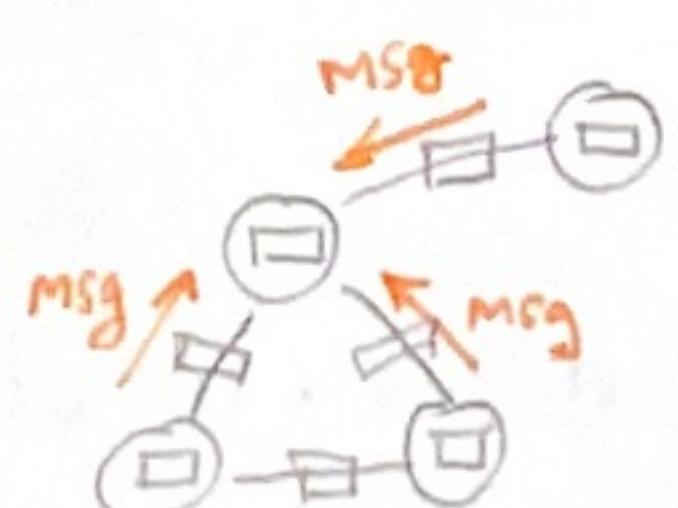
→ After an edge update, each edge embedding has info about its incident nodes.

→ Use edge embedding to update nodes.

chosen msg for node i

$$m_i^{(l)} = \underbrace{\oplus}_{\text{order-invariant operation (sum/max/mean)}} \left( \{ h_{(i,j)}^{(l)} \}_{j \in N_i} \right)$$

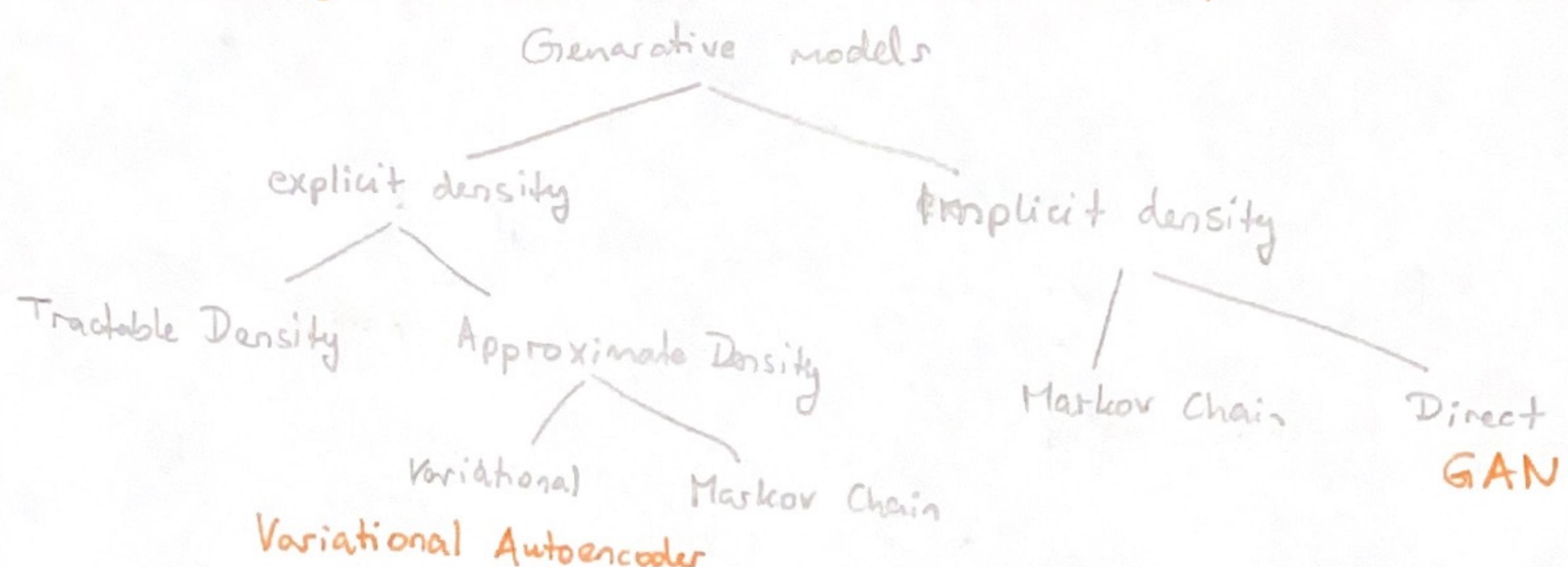
neighbors of node i



$$h_i^{(l)} = \mathcal{N}_v \left( [m_i^{(l)}, h_i^{(l-1)}] \right)$$

new embedding for node i learnable func with weights shared for all nodes previous embedding for node i

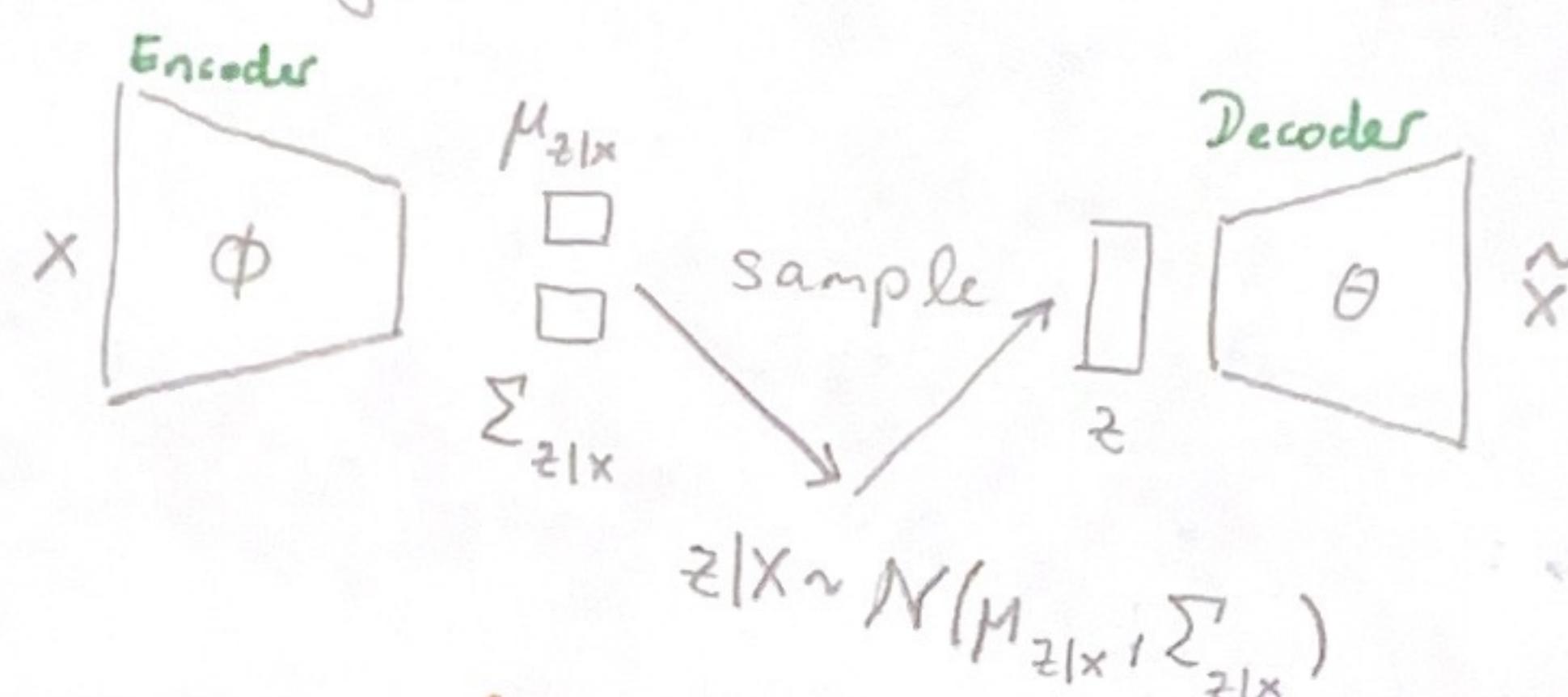
- Generative models:  
→ Given training data generate "realistic" models from the same distribution



→ Variational autoencoder:

Goal: Sample from latent distr. to generate new outputs.

⇒ Latent space is Gaussian



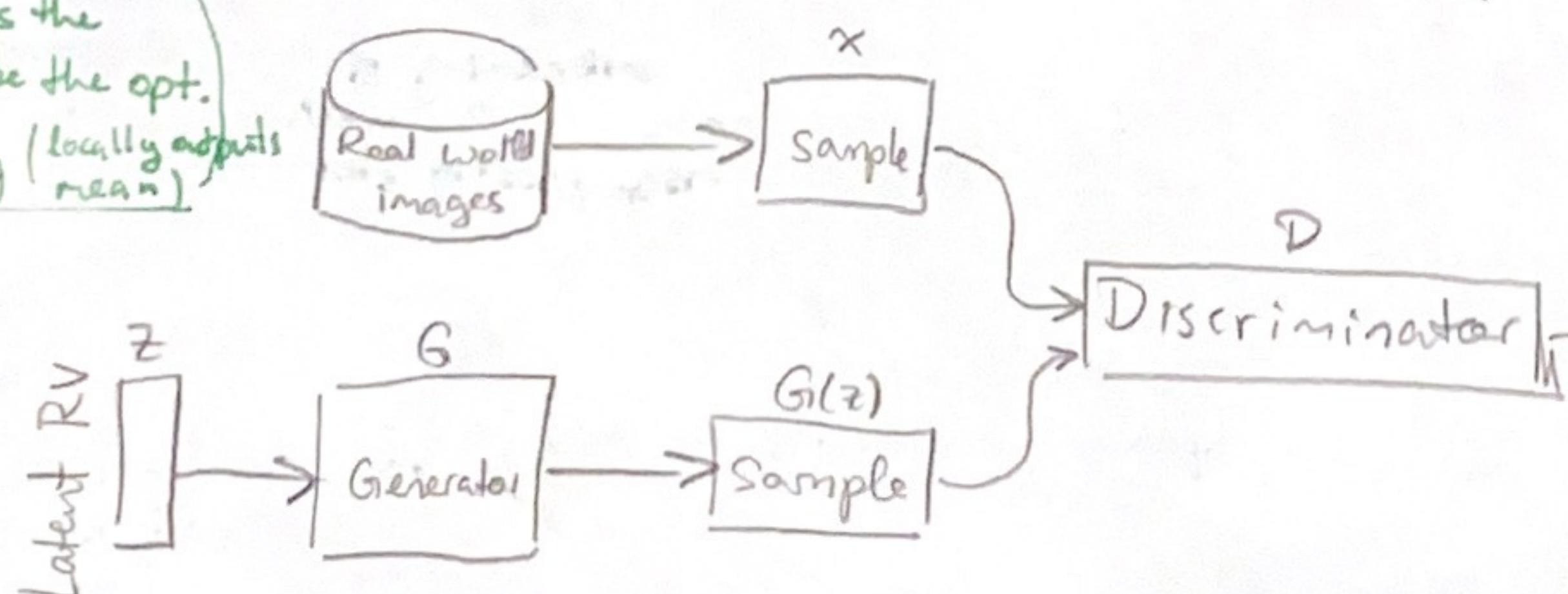
→ Prob. distr. in latent space → sample from it

interpretable (sample specific head pose / degree of smile)

- Generative Adversarial Networks: (GANs)

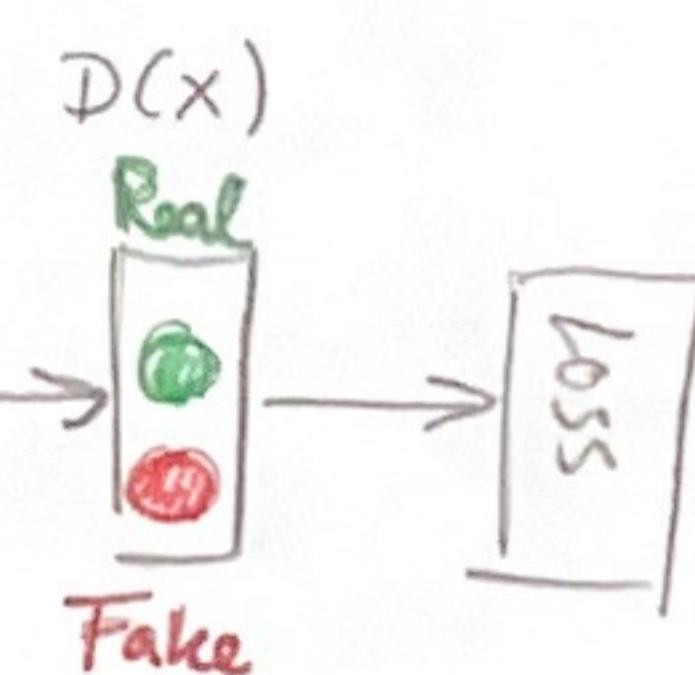
The L2-loss in Autoencoders distributes the error equally among the image  
→ for constant imgs the mean would be the opt.  
⇒ Result is blurry (locally outputs mean)

→ Learn the loss function



On the real data:  $D(x)$  tries to be near 1

On the fake data:  $D$  tries to make  $D(G(z)) \approx 0$   
 $G$  tries to make  $D(G(z)) \approx 1$



$D(x)$

$D(G(z))$

Components of RL agents

\* Policy: Mapping from state to action  $\pi(s) = a$

\* Value-/Q-func: How good is a state/state-action-pair  
→ future reward

- Reinforcement Learning:

→ Given sequential data, learn by interaction with the environment

\* data is non i.i.d. (time matters / it is sequential)

\* Actions affect the env (change future input)

\* No supervisor (target is approximated by the reward signal)

\* Markov assumption: "The state only depends on the previous state"

