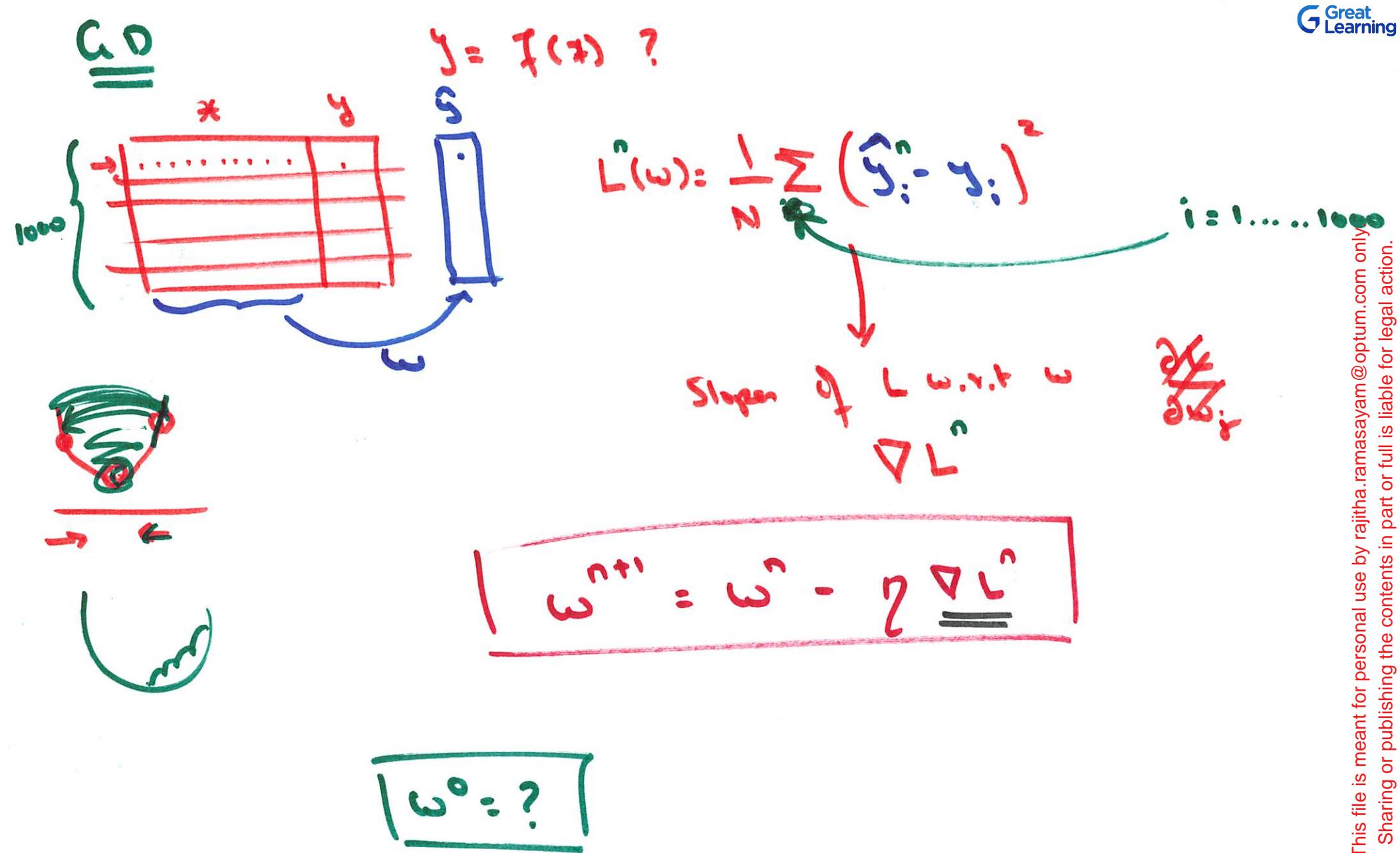


# Agenda

- The Challenges: Over fitting & local optima
- The Training
  - Epochs, Batch Size, Iterations
  - Gradient Descent (GD) Vs Stochastic GD (SGD) Vs Mini-Batch GD
  - SGD with momentum
  - Learning rates and adaptive learning rates
  - Weight Initialization
  - Batch Normalization
- Guarding against over-fitting
  - L1/L2 Regularization
  - Data Augmentation
  - Drop outs
- Neural Network Architectures

weigh decay



SGD

$$(GD) \quad L(\omega) = \frac{1}{N} \sum_{i=1 \dots 1000} (\hat{y}_i - y_i)^2$$

$$(SGD) \quad L(\omega) = (\hat{y}_i - y_i)^2$$

randomly chosen

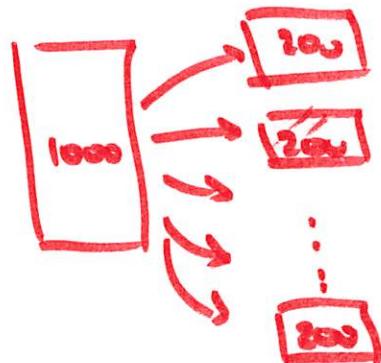
Mini Batch  
SGD

$$N = 1000$$

$$N_b = 200$$

$$\text{iter} = 5$$

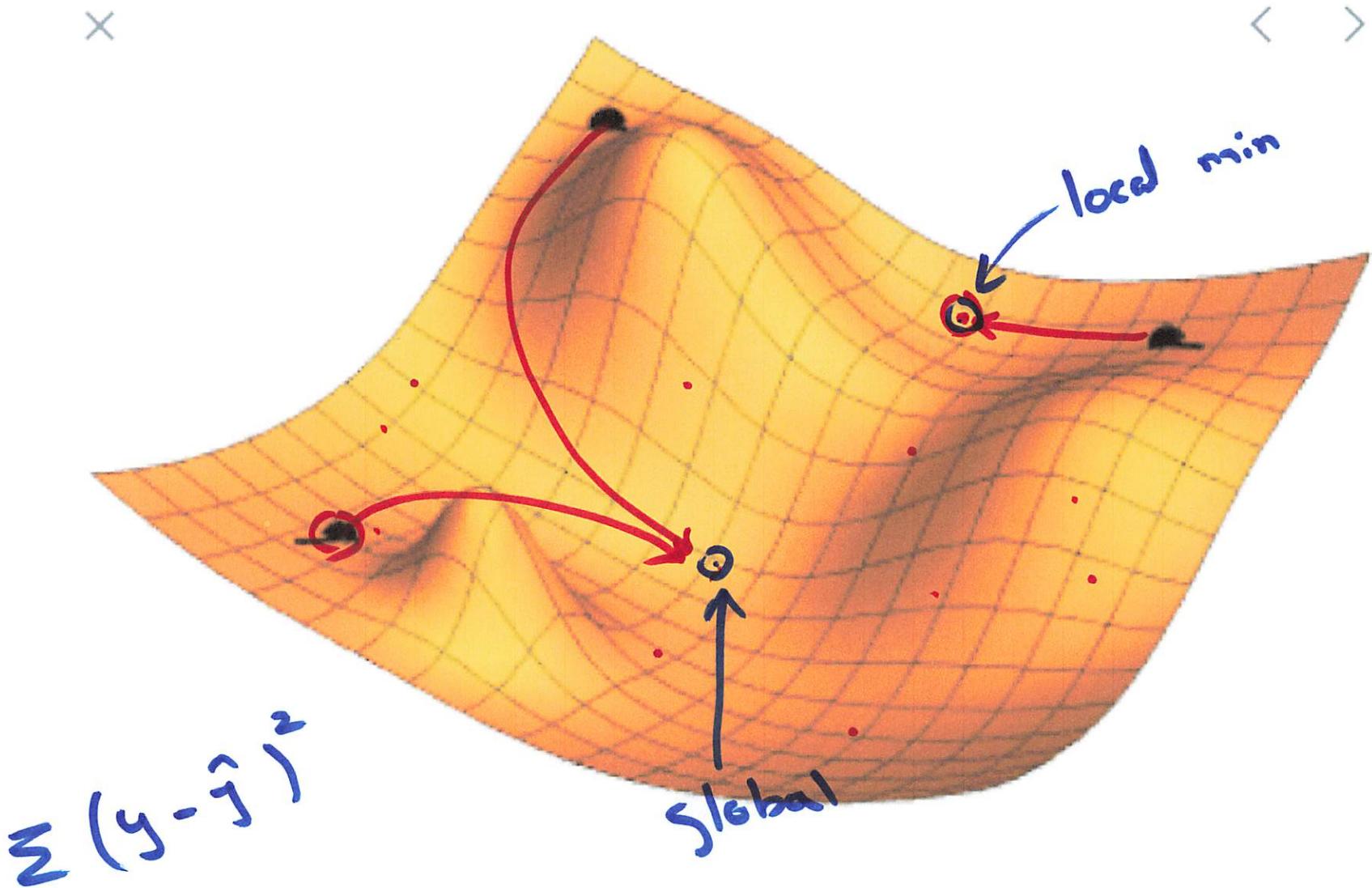
epoch



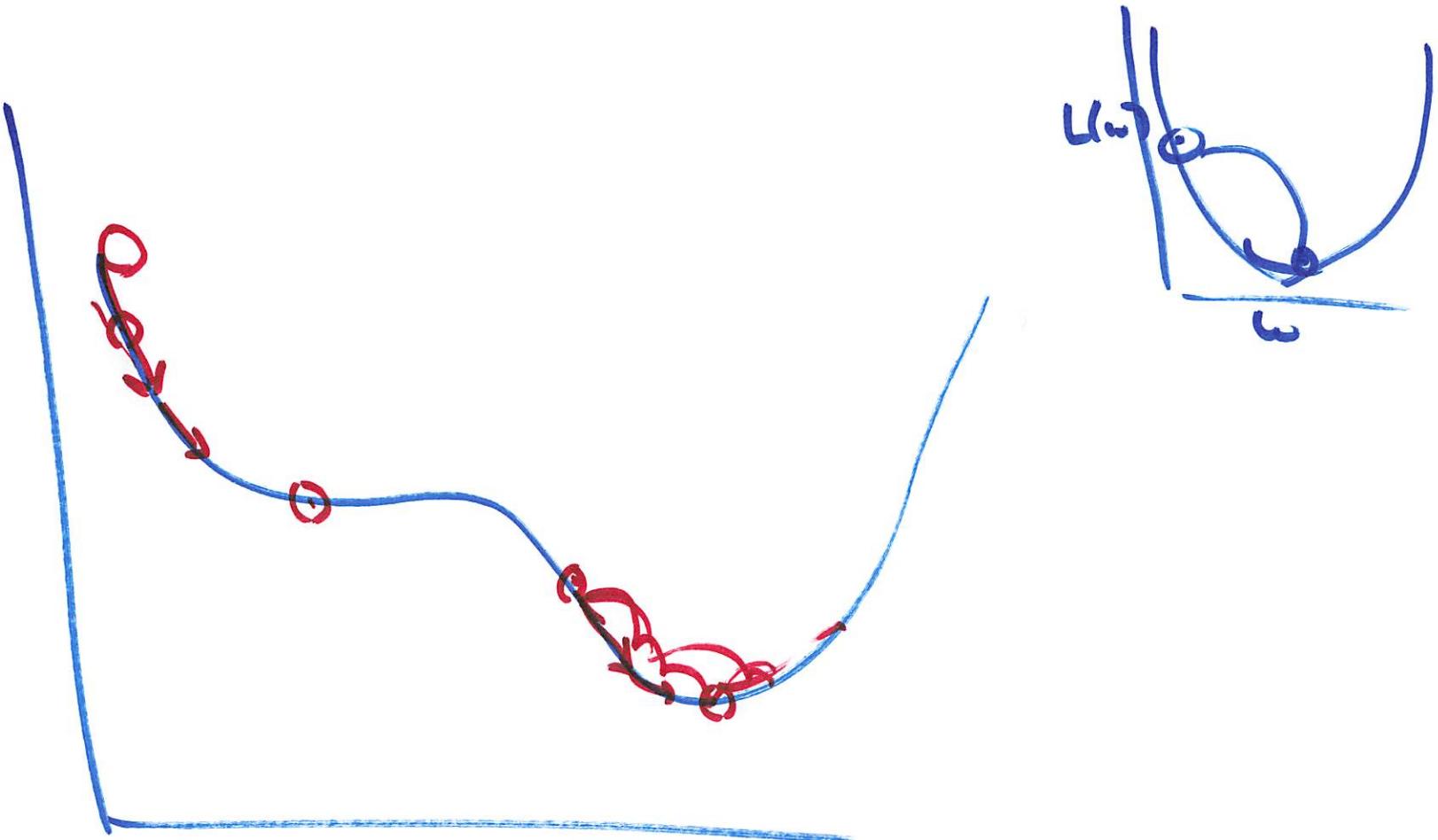
$$L(\omega) = \frac{1}{N_b} \sum_{i \in B} (\hat{y}_i - y_i)^2$$

batch size =  $N_b$

$$\text{iterations} = \frac{N}{N_b}$$



This file is meant for personal use by rajitha.ramasayam@optum.com only.  
Sharing or publishing the contents in part or full is liable for legal action.



$$\hat{\beta}'' = \beta' - \eta \nabla L'$$

SGD with momentum  $\Rightarrow$

$$\hat{\beta}'' = \beta' - \eta (\alpha \nabla L' + (1-\alpha) \nabla L'')$$

# Learning rate

- Choosing the Learning rate ( $\eta$ )
  - Too small, we will need too many iterations for convergence
  - Too large, we may skip the optimal solution
- Adaptive Learning Rate :
  - start with high learning rate and
  - gradually reduce the learning rate with each iteration.
  - Moreover, having different learning rates for different weight updates will help: Adagrad, RMS Prop

Adam

AdaDelta

Adaptive Grad

$$\hat{g}^t = g^t - \frac{\eta}{\sqrt{\delta^t + \epsilon}} \nabla L^t$$

$$\delta^t = \sum_i (\nabla L^t)^i$$

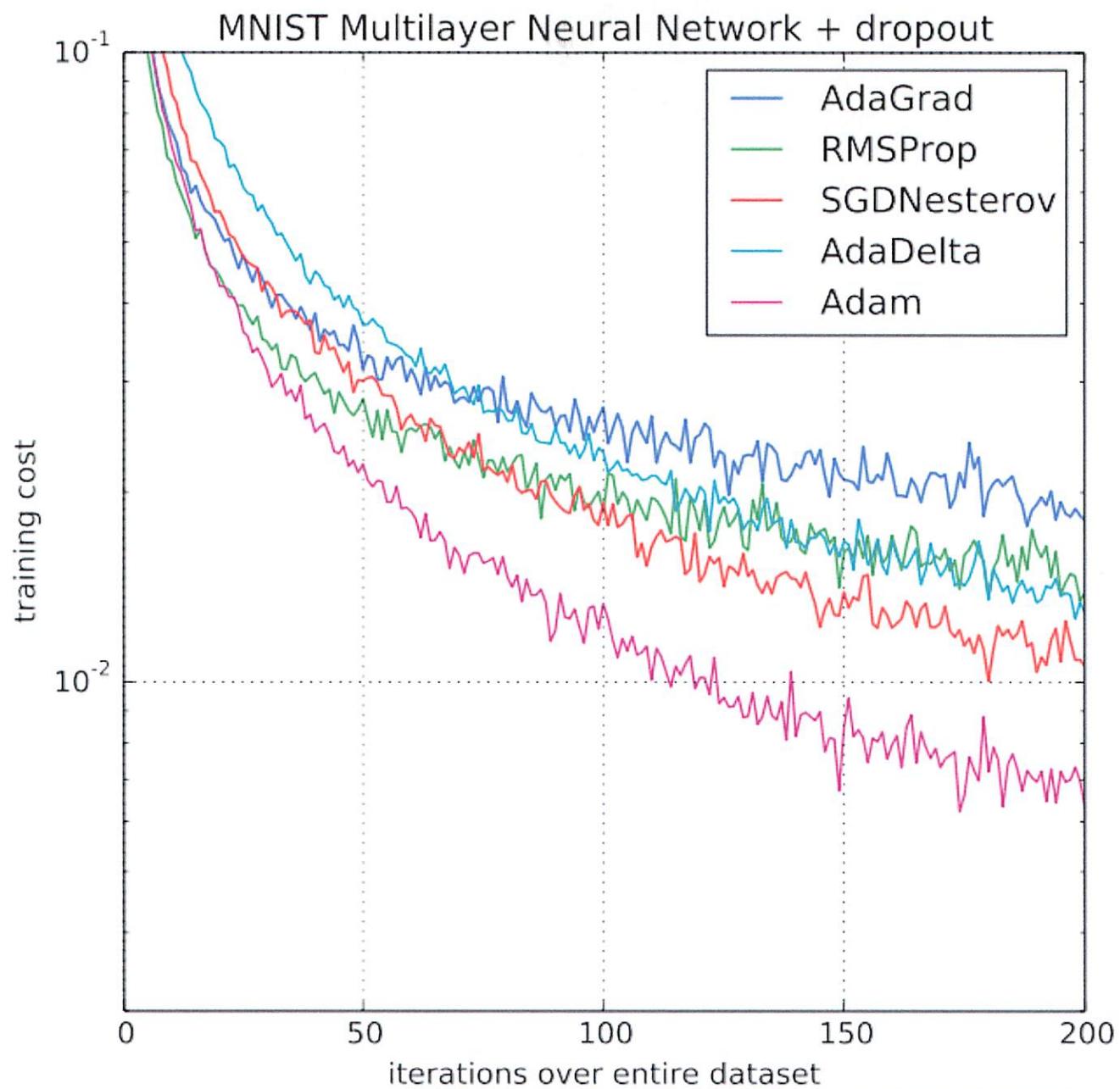
RMS Prop

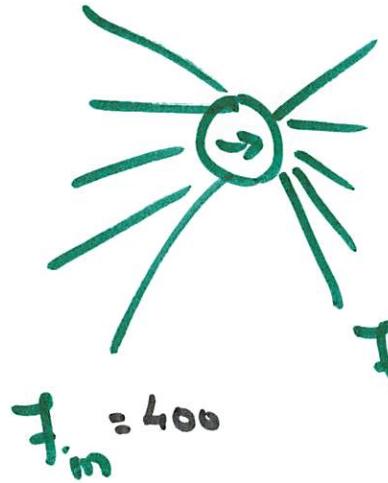
$$\hat{g}^t = g^t - \frac{\eta}{\sqrt{\delta^t + \epsilon}} \nabla L^t$$

$$\delta^t = \alpha \delta^{t-1} + (1-\alpha) \nabla L^t$$

Adam

→ { RMS Prop  
SGD + Momentum



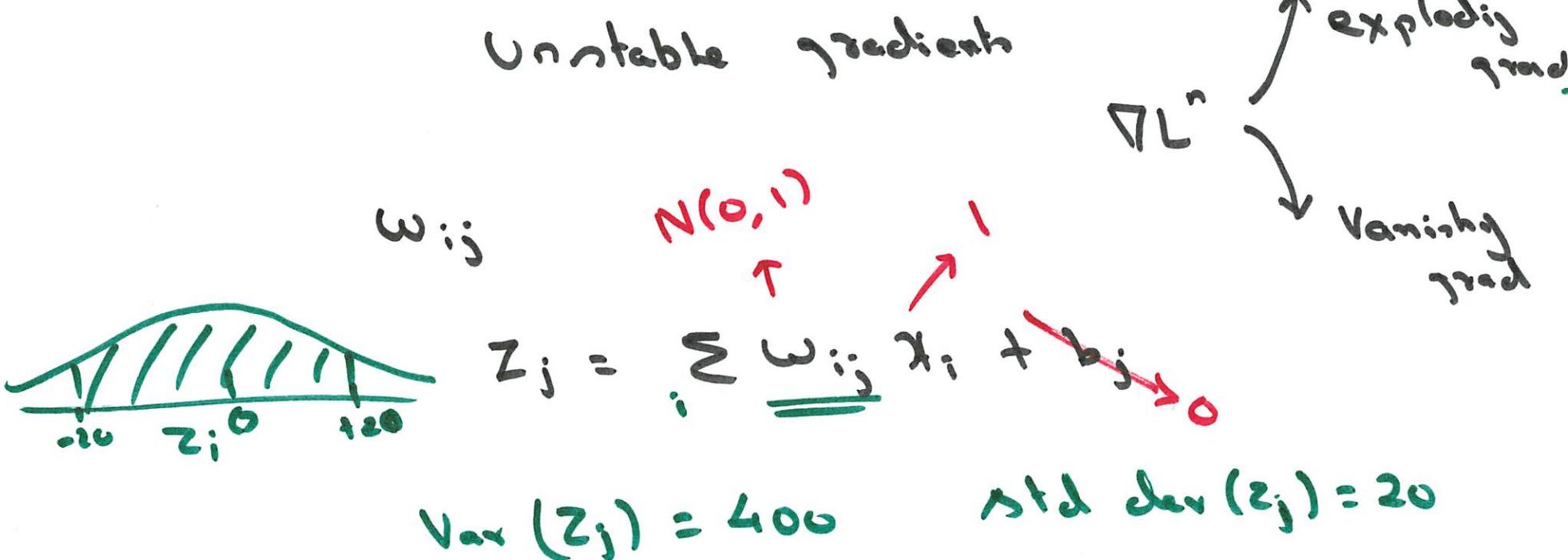


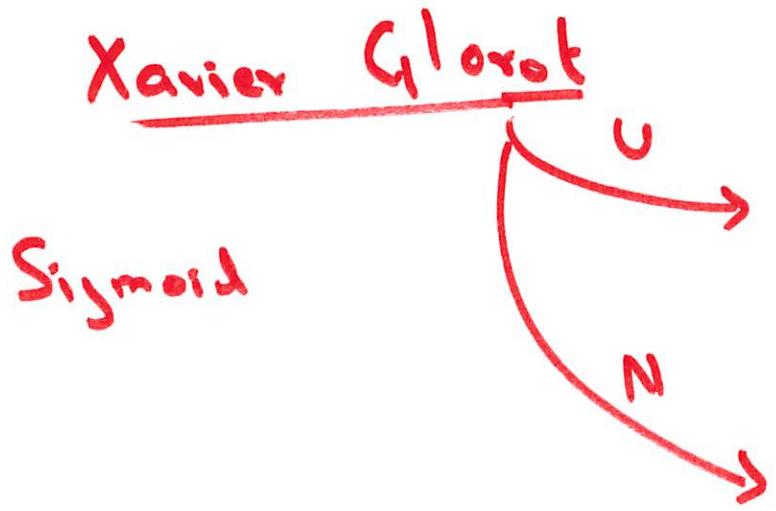
$$a_j = \sigma \left( \sum w_{ij} x_i + b_j \right)$$

$a_j$

$w_{ij} = 0$

$w_{ij} \sim \text{Normal}(0, 1)$





Uniform

$$\left[ -\sqrt{\frac{6}{f_{in} + f_{out}}}, +\sqrt{\frac{6}{f_{in} + f_{out}}} \right]$$

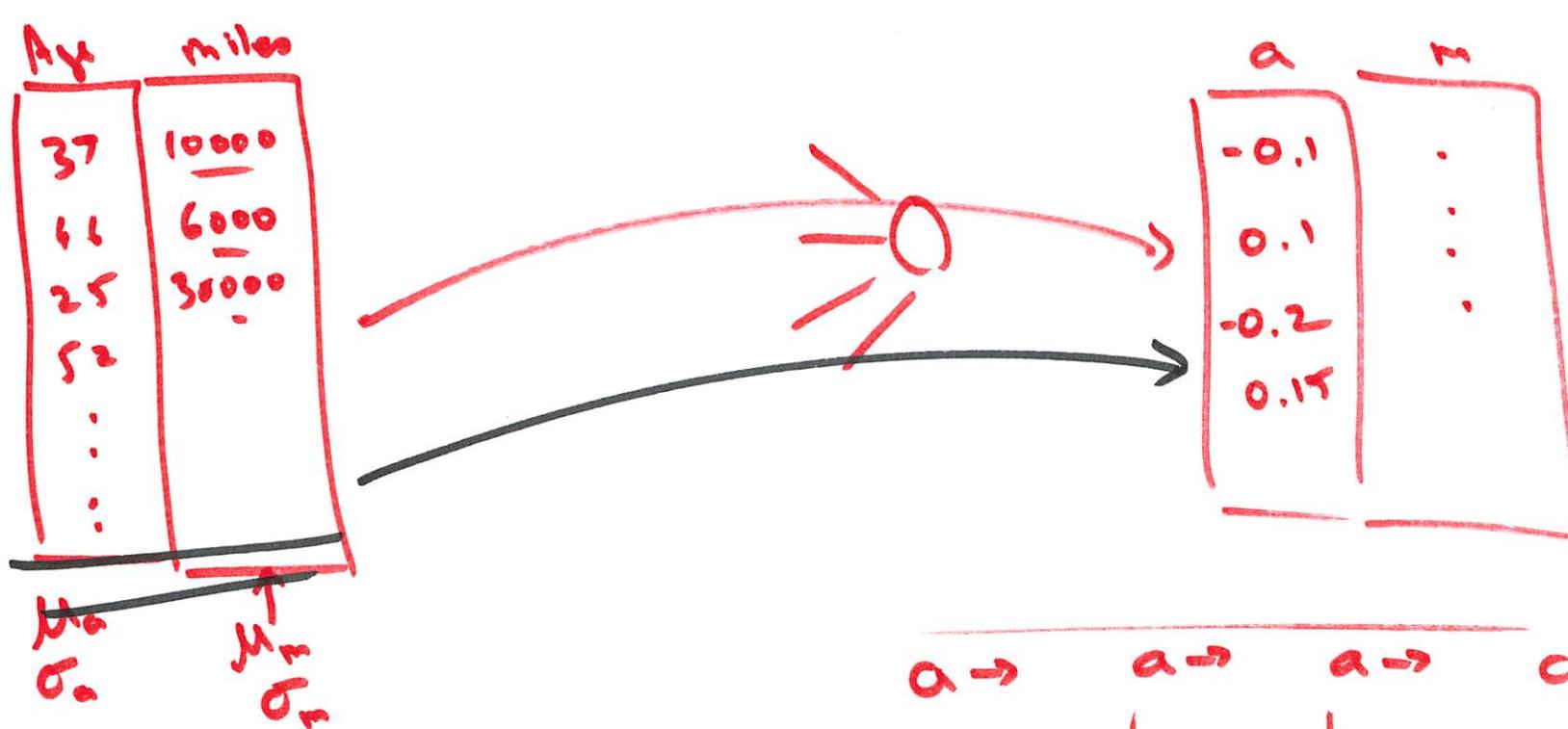
Normal

$$\left( 0, \sqrt{\frac{2}{f_{in} + f_{out}}} \right)$$

ReLU

Normal

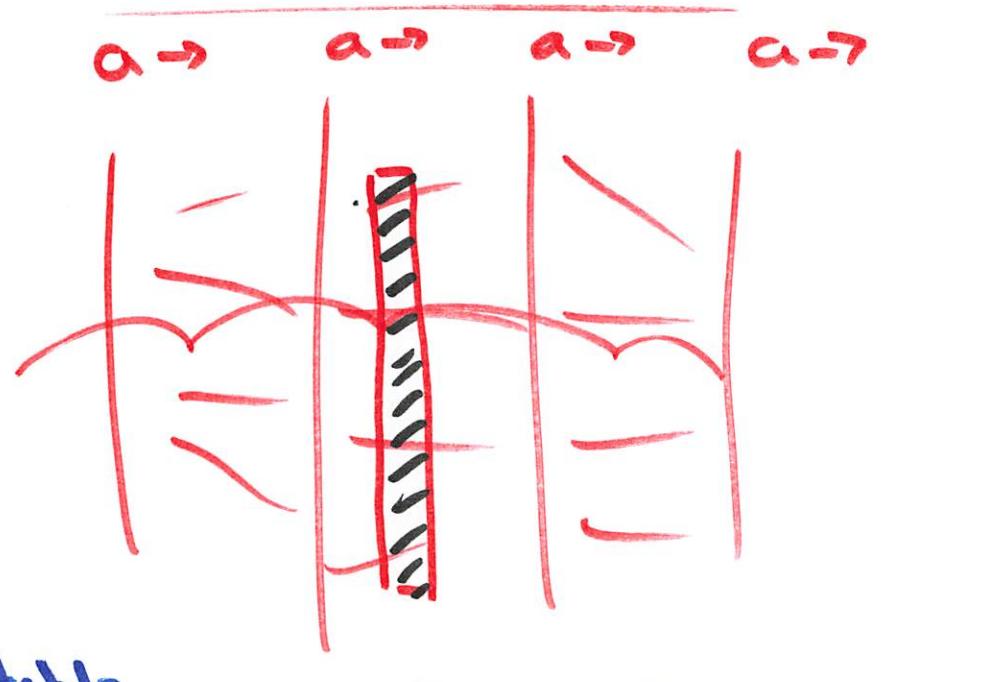
$$\left( 0, \sqrt{\frac{2}{f_{in}}} \right)$$

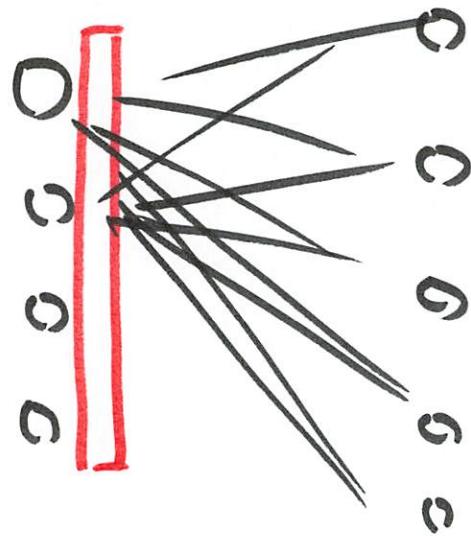


$$\hat{x}_a = \frac{x_a - Ma}{\sigma_a}$$

$$\hat{x}_m = \frac{x_m - Mm}{\sigma_m}$$

long  $\rightarrow$  forward  $\xrightarrow{\text{Stable}}$

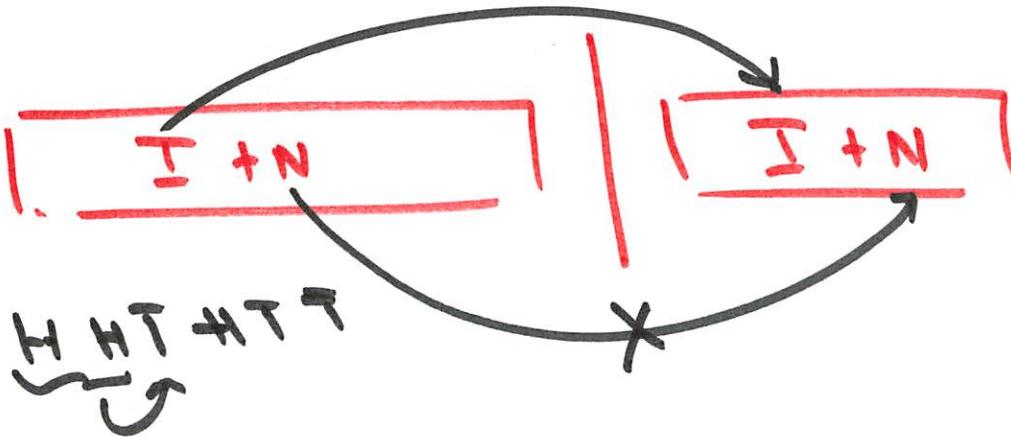




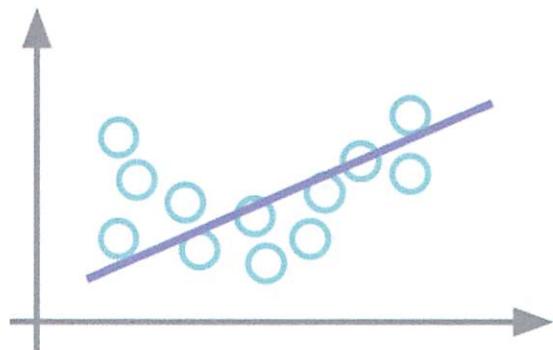
$$\hat{a}_i = \frac{a_i - \mu}{\sigma}$$
$$\hat{a}_i = \gamma \hat{a}_i + \beta$$

mem:  $\mu_p, \sigma_p$

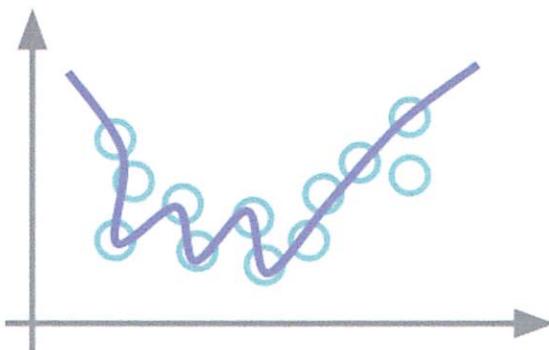
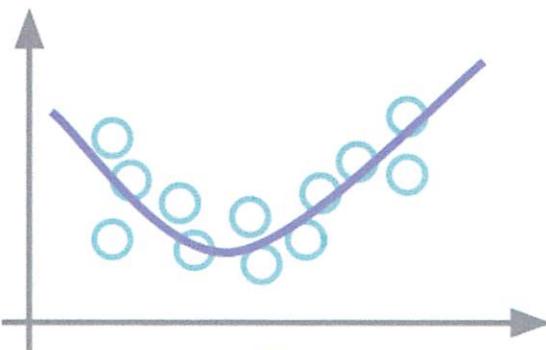
Data =  $I + N$  + Noise



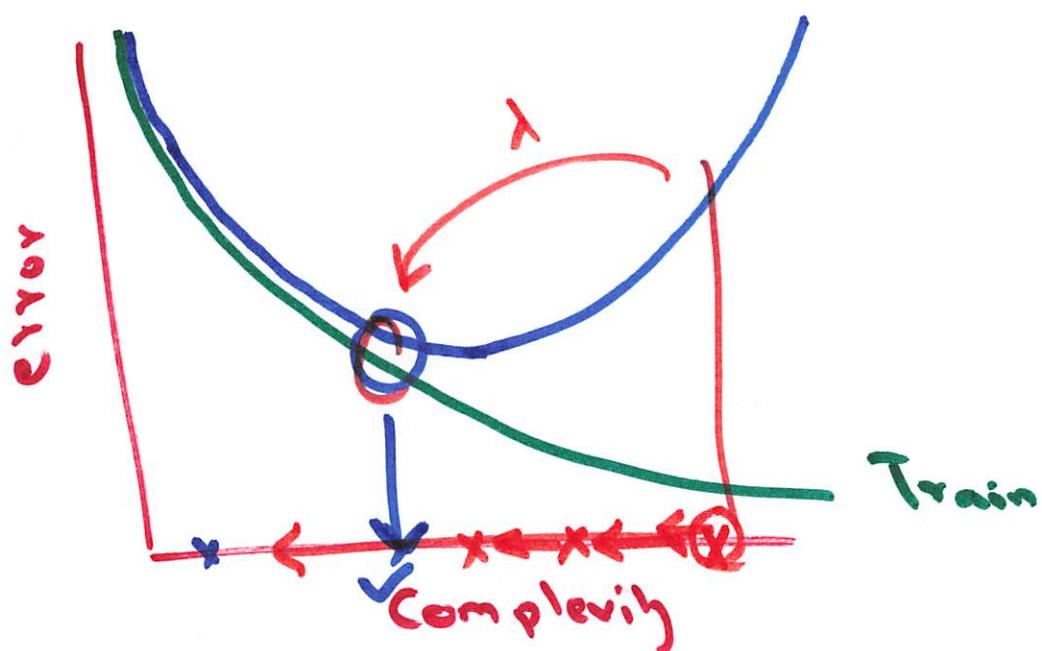
# Under Vs Over-fitting



Underfit



overfit



$L_1 + L_2$  Reg.

$\nabla$

$\min$

$$L(\beta) = \frac{1}{n} \sum (\gamma - \hat{y})^2 + \frac{\lambda}{p} (\text{penalty})$$

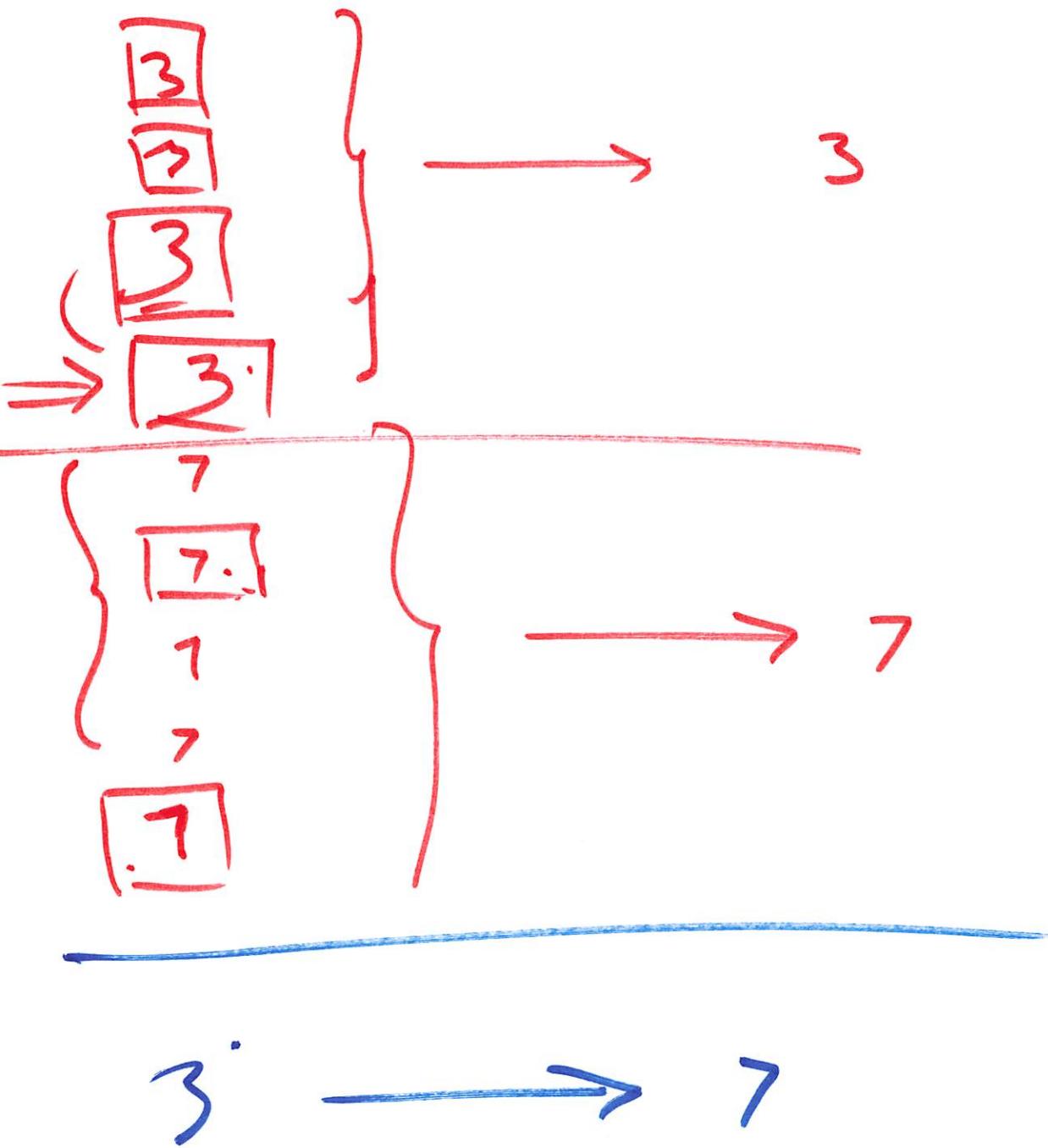
penalty

$$\rightarrow \sum |w_i|$$

$$\rightarrow \sum (\beta_i)^2$$

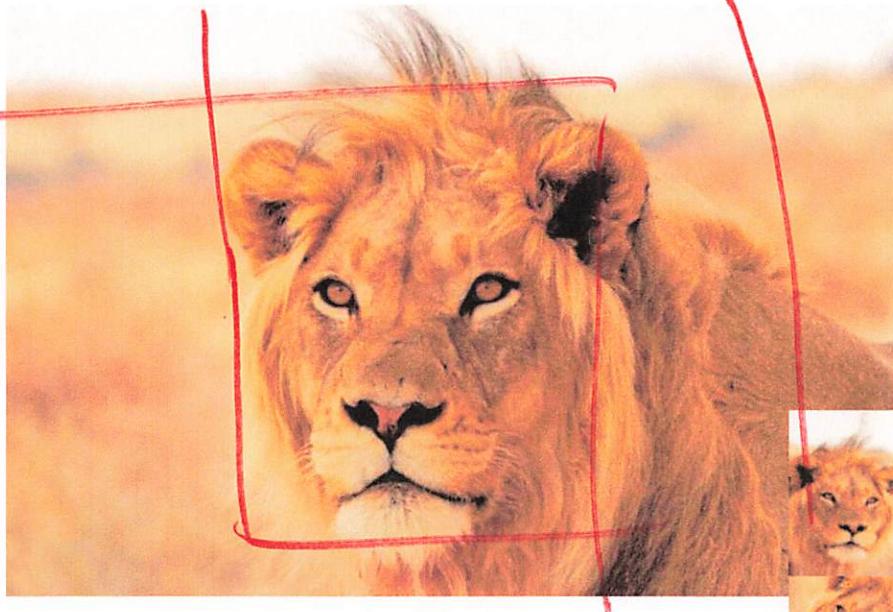
Lasso =

Ridge



This file is meant for personal use by rajitha.ramasayam@optum.com only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Data Augmentation



$\pm 15^\circ$   
Noise  
Shift  
mirror  
Stretch  
Color  
Crop  
GAN

-source: [towardsdatascience.com](https://towardsdatascience.com/machinex-image-data-augmentation-using-keras-2f3a2a2a3e0d), MachineX: Image Data Augmentation Using Keras  
This file is meant for personal use by Rajmata.amasyam@optum.com only.

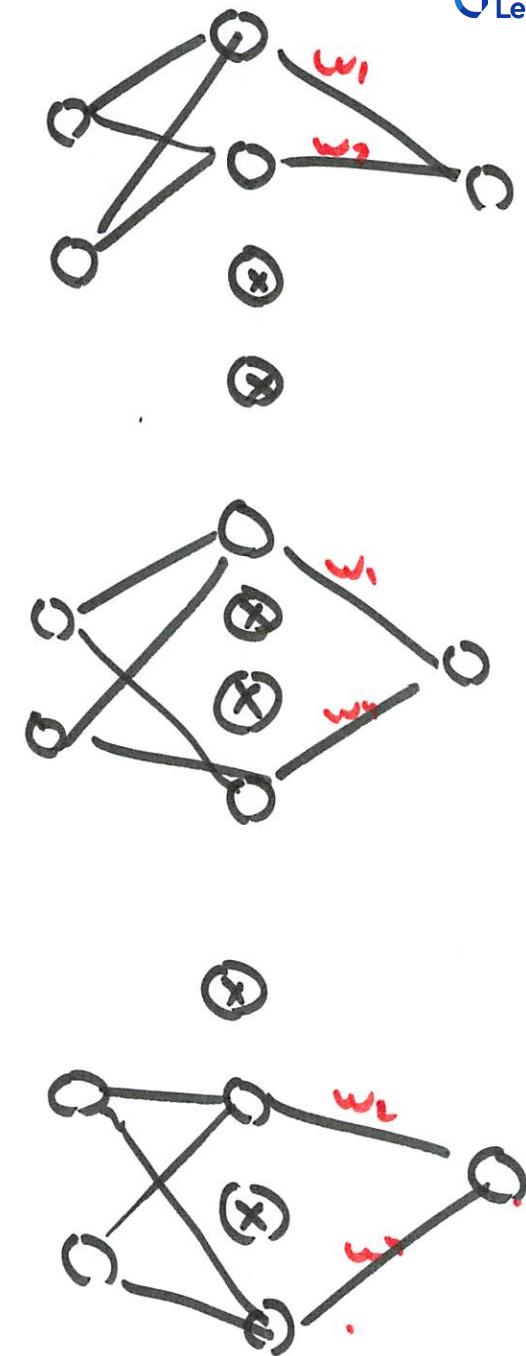
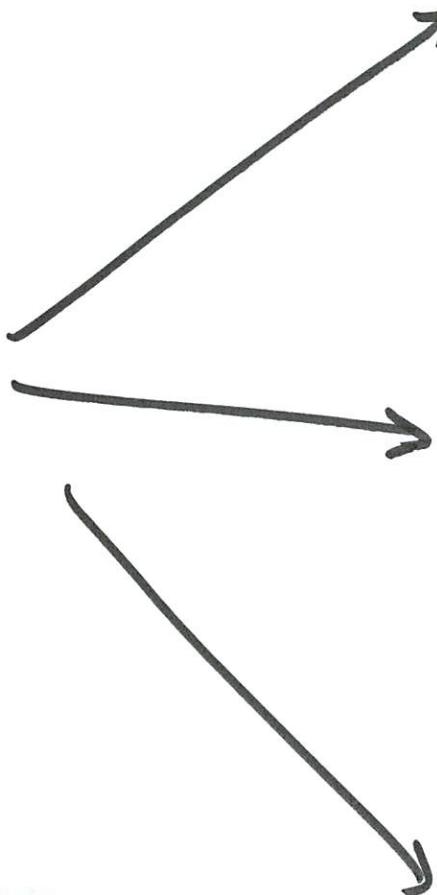
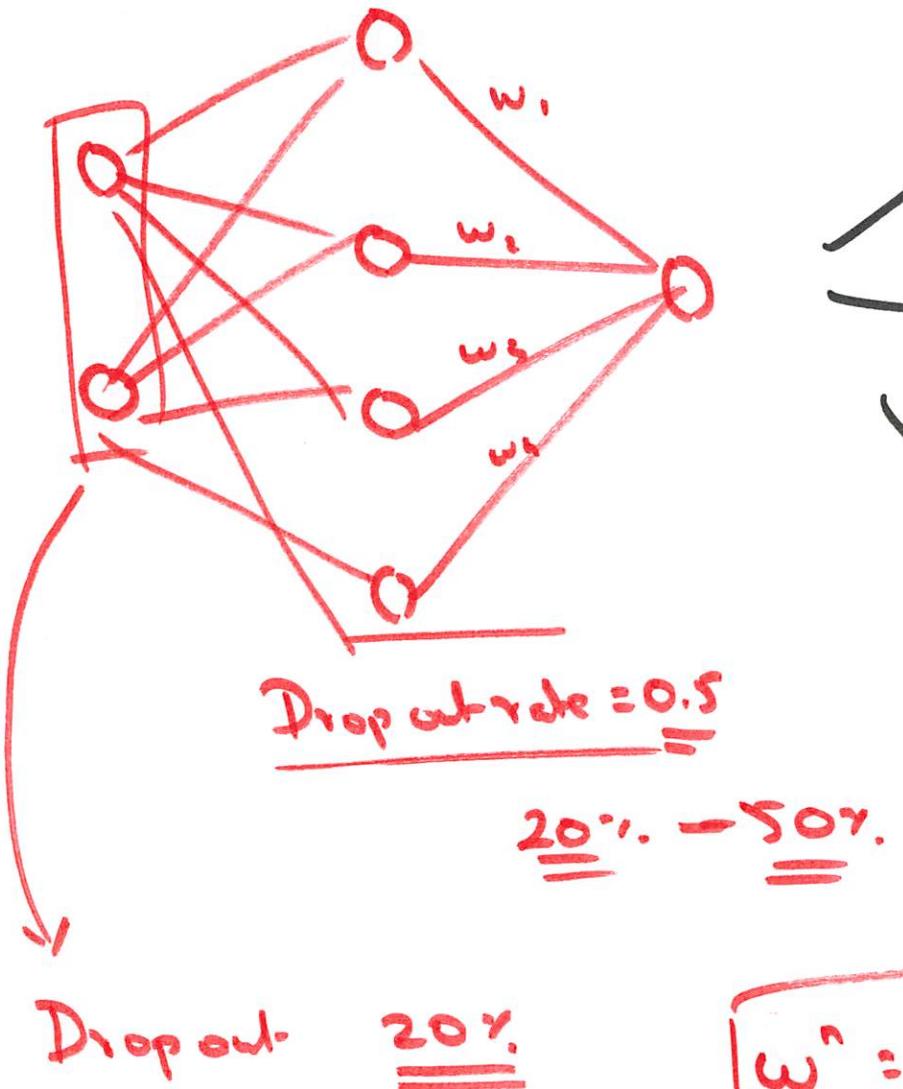
Sharing or publishing the contents in part or full is liable for legal action.

## Co Adaptation



$$-0.1 \alpha_1 + 0.1 \alpha_2$$

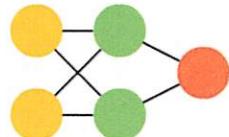
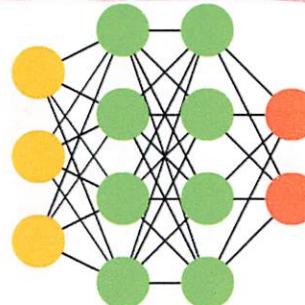
$$\underline{-0.3 \alpha_1 + 0.3 \alpha_2}$$



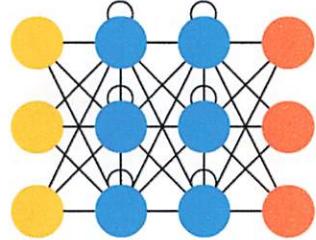
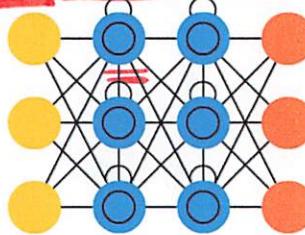
# Types of NN

- Feed Forward
  - MLP
  - DNN
  - CNN
- RNN
- LSTM
- .
- .
- .
- .
- .
- Transf

Feed Forward (FF)

Deep Feed Forward (DFF)

Recurrent Neural Network (RNN)

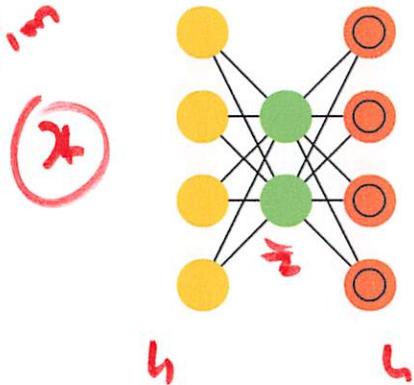
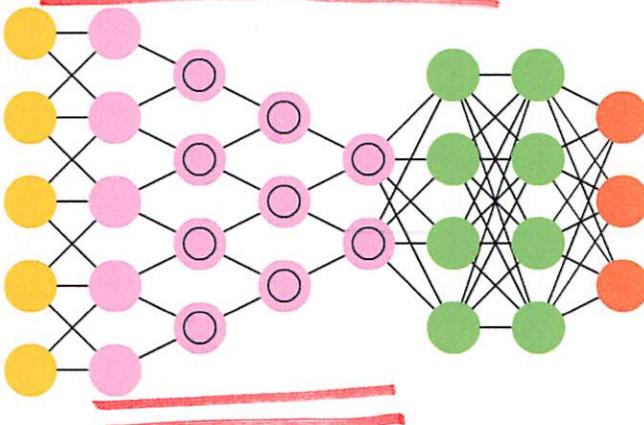
Long / Short Term Memory (LSTM)

17 5  
1

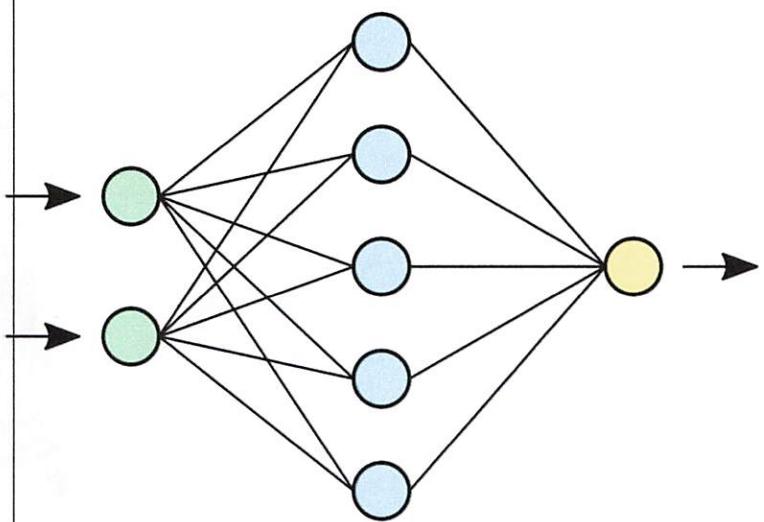
17  
17

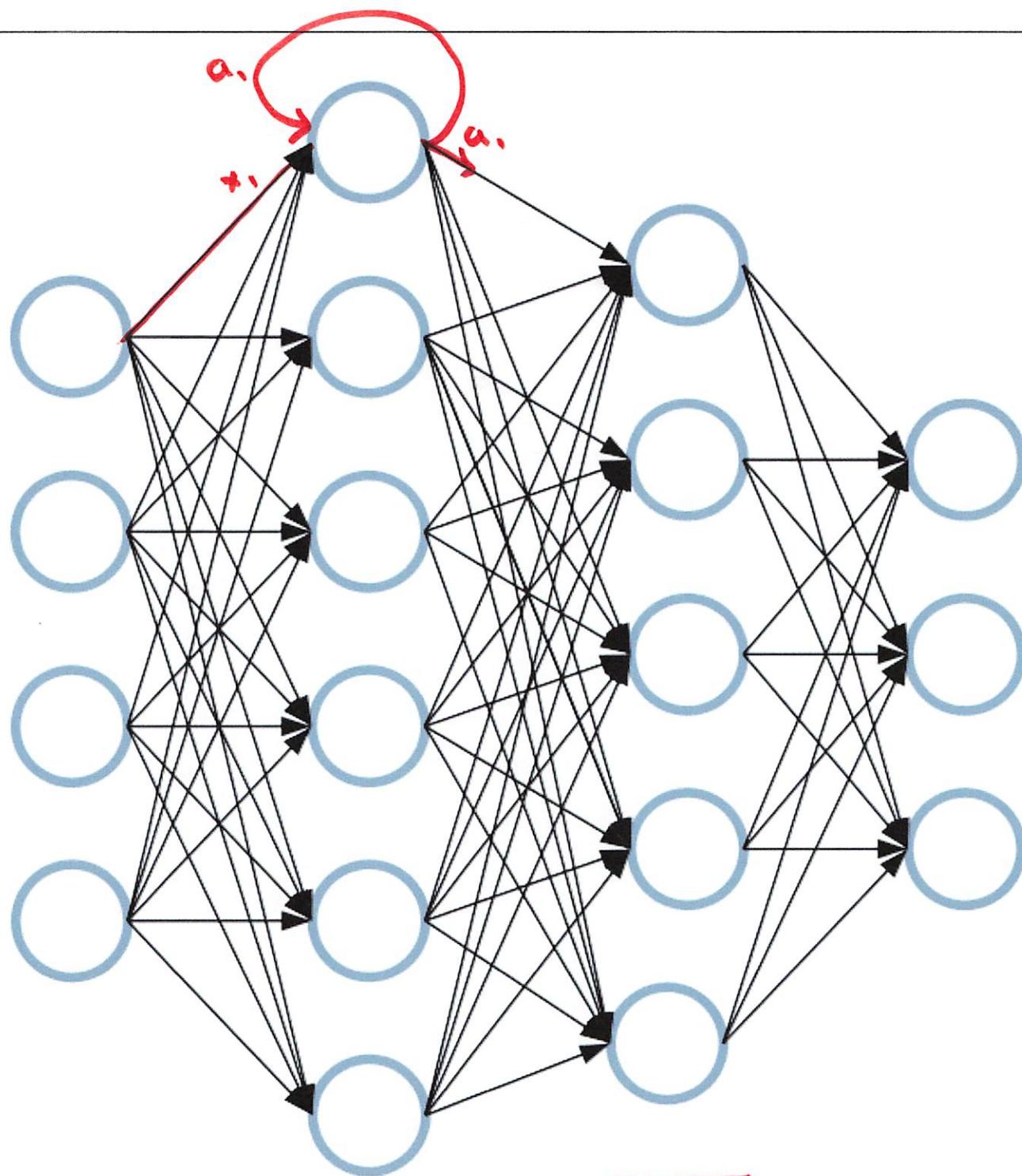
*Compressed  
dim*

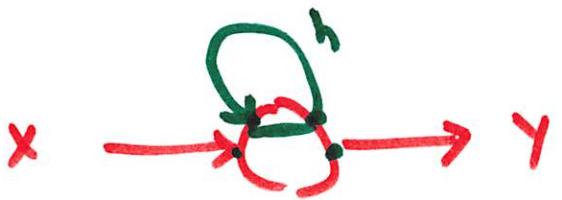
Auto Encoder (AE)

Deep Convolutional Network (DCN)

Feed Forward Net.



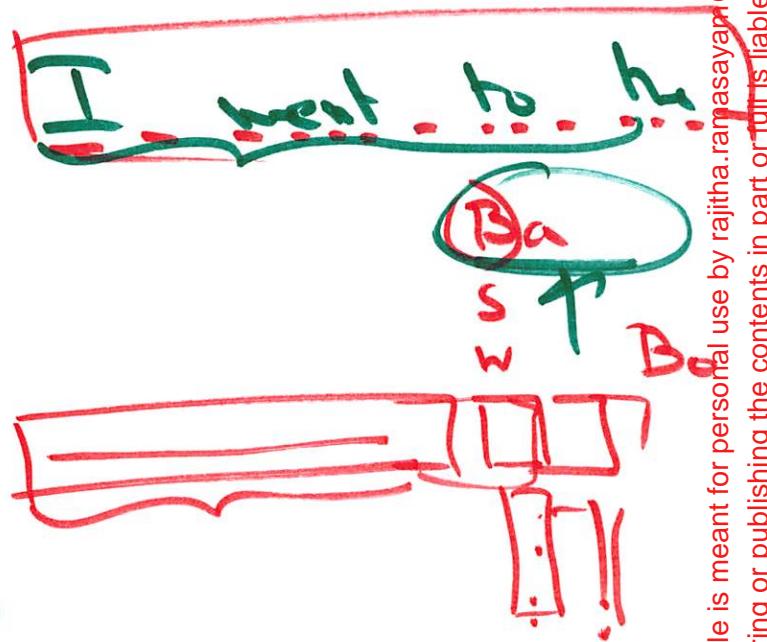
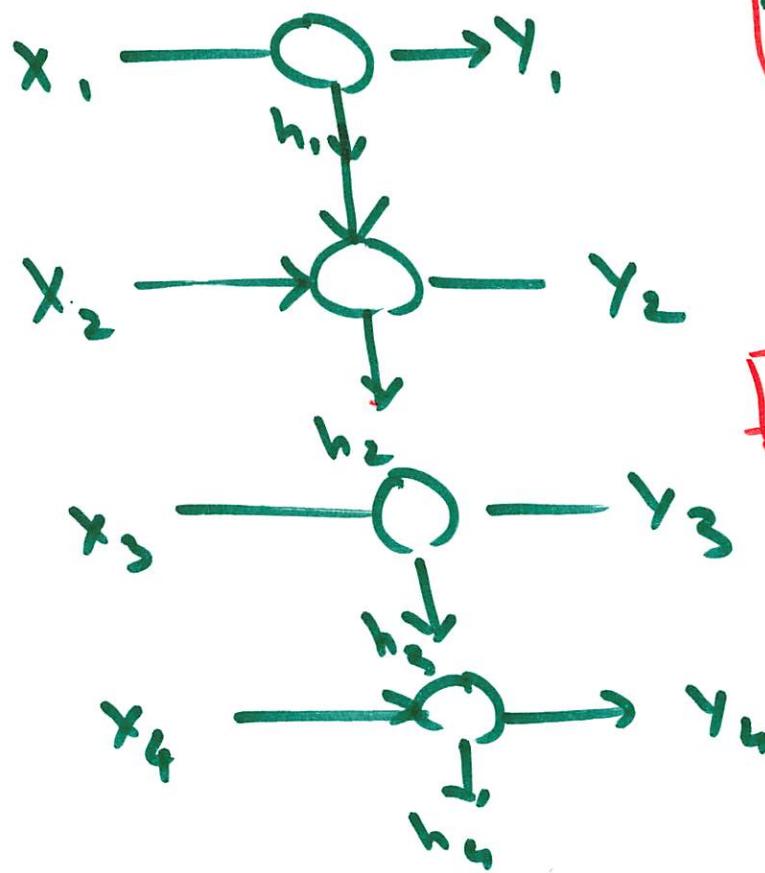




$$y = f(\omega x + b)$$

$$y = f(\omega x + \omega_0 b + b)$$

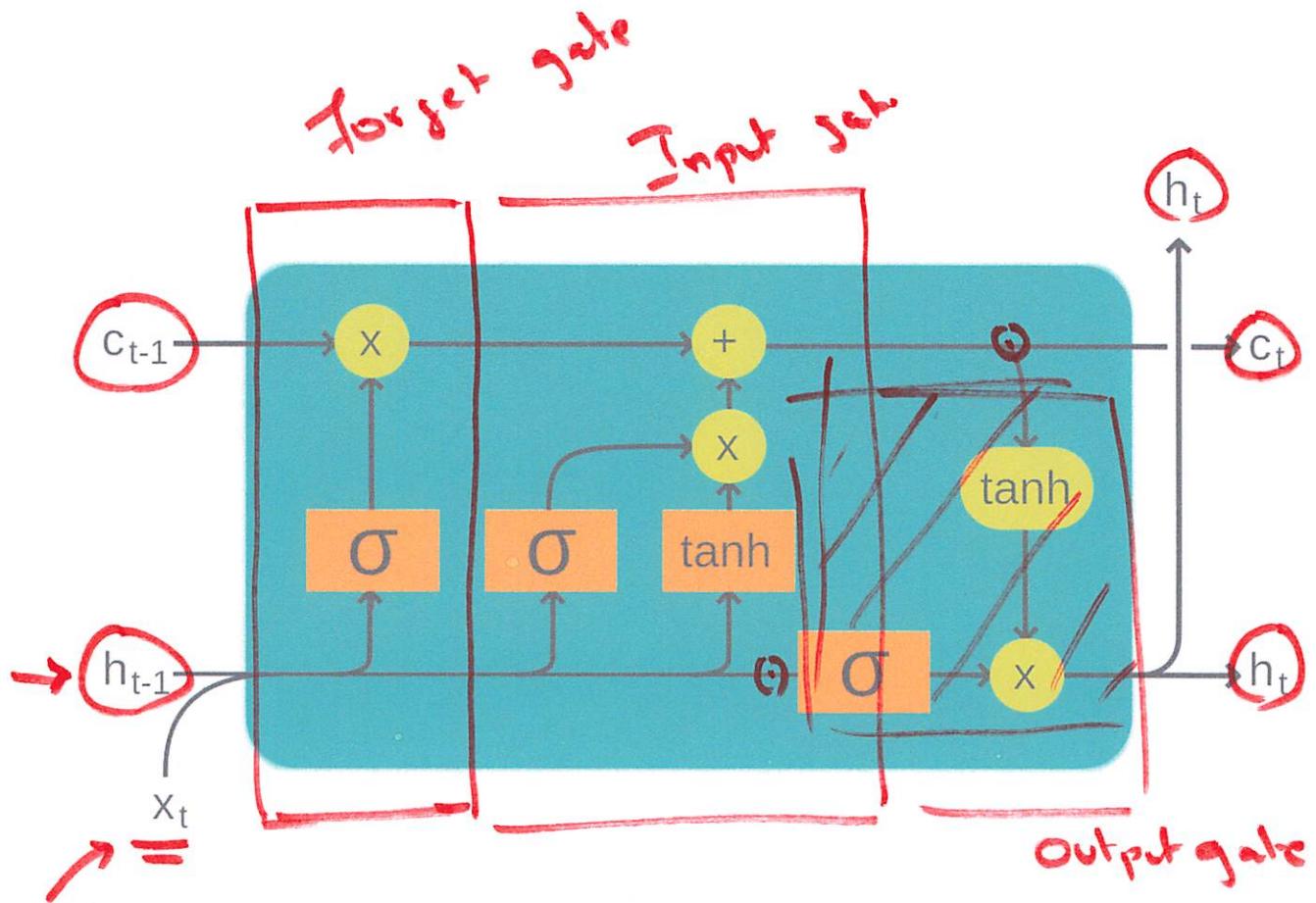
$$f = \sigma(\quad)$$



86

# Text generation using an RNN

The meaning of life is the tradition of the ancient human reproduction: it is less favorable to the good boy for when to remove her bigger. In the show's agreement unanimously resurfaced. The wild pastured with consistent street forests were incorporated by the 15th century BE. In 1996 the primary rapford undergoes an effort that the reserve conditioning, written into Jewish cities, sleepers to incorporate the .St Eurasia that activates the population. Mar??a Nationale, Kelli, Zedlat-Dukastoe, Florendon, Ptu's thought is. To adapt in most parts of North America, the dynamic fairy Dan please believes, the free speech are much related to the



Legend:

Layer	Componentwise	Copy	Concatenate

