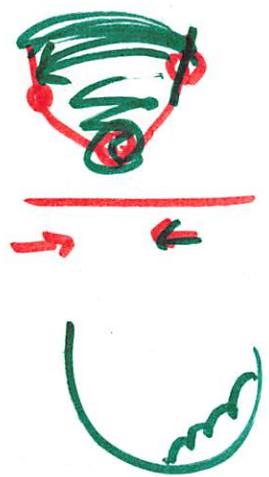
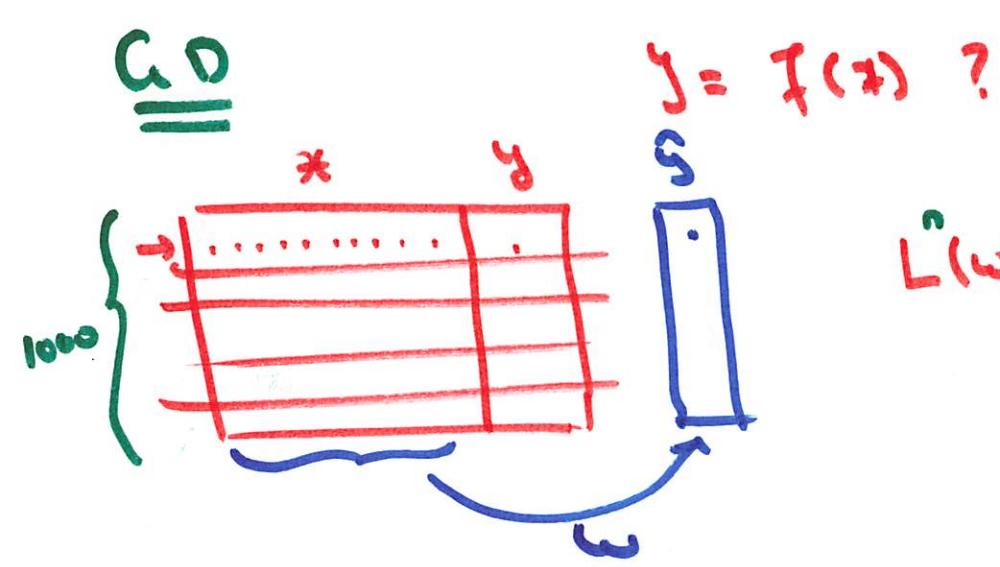


Agenda

- 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



$$L(w) = \frac{1}{n} \sum_i (\hat{y}_i - y_i)^2$$

$i = 1, \dots, 1000$

Slope of L w.r.t w

$$\nabla L$$



$$w^{n+1} = w^n - \eta \nabla L$$

$$w^0 = ?$$

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

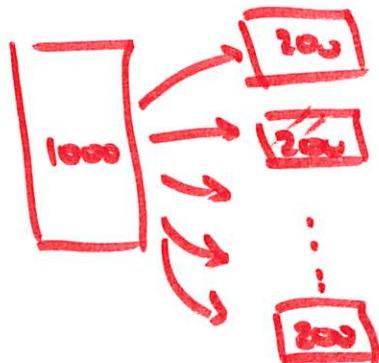
Mini Batch
SGD

$$N = 1000$$

$$N_b = 200$$

$$\text{iter} = 5$$

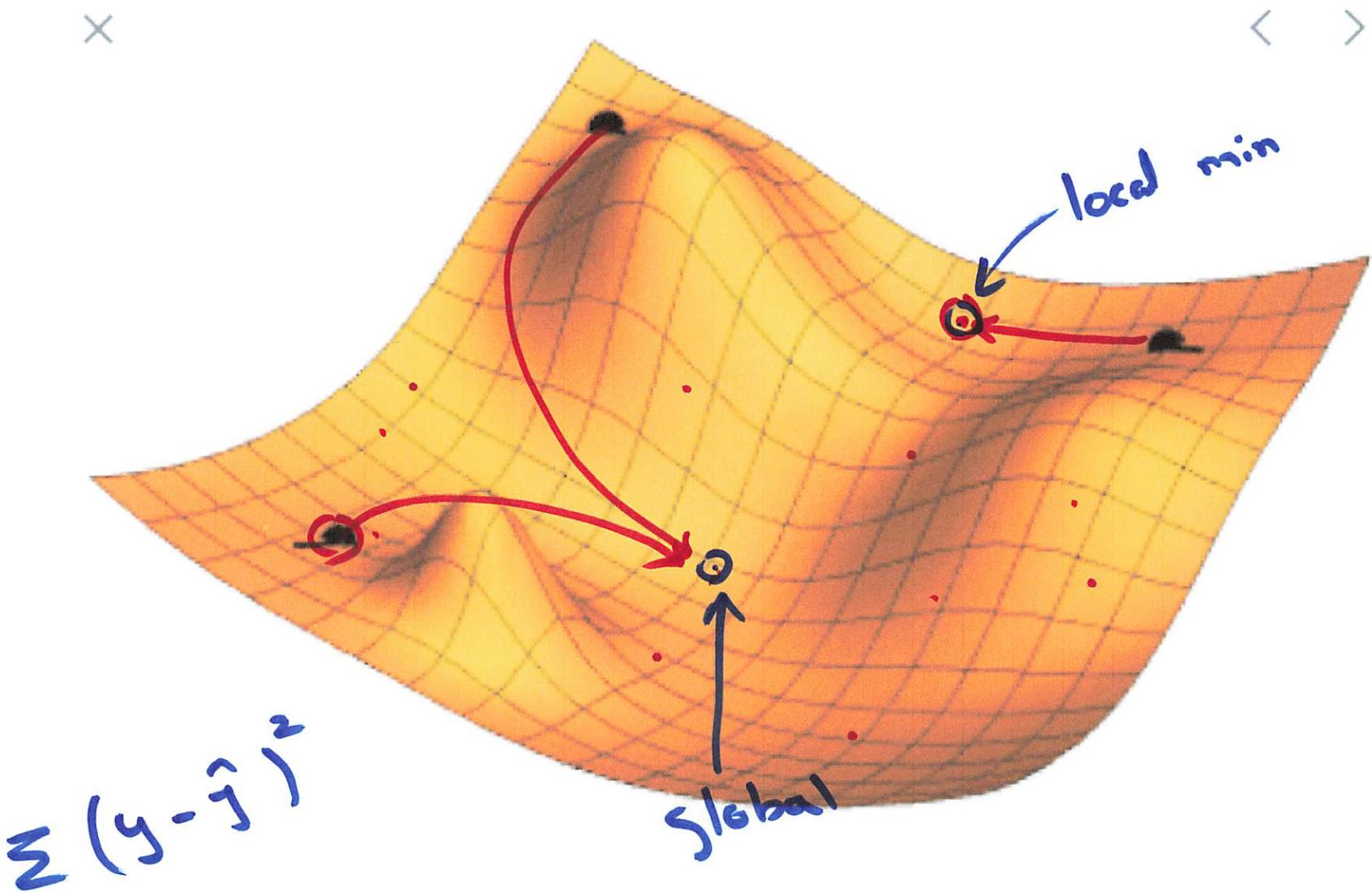
epoch



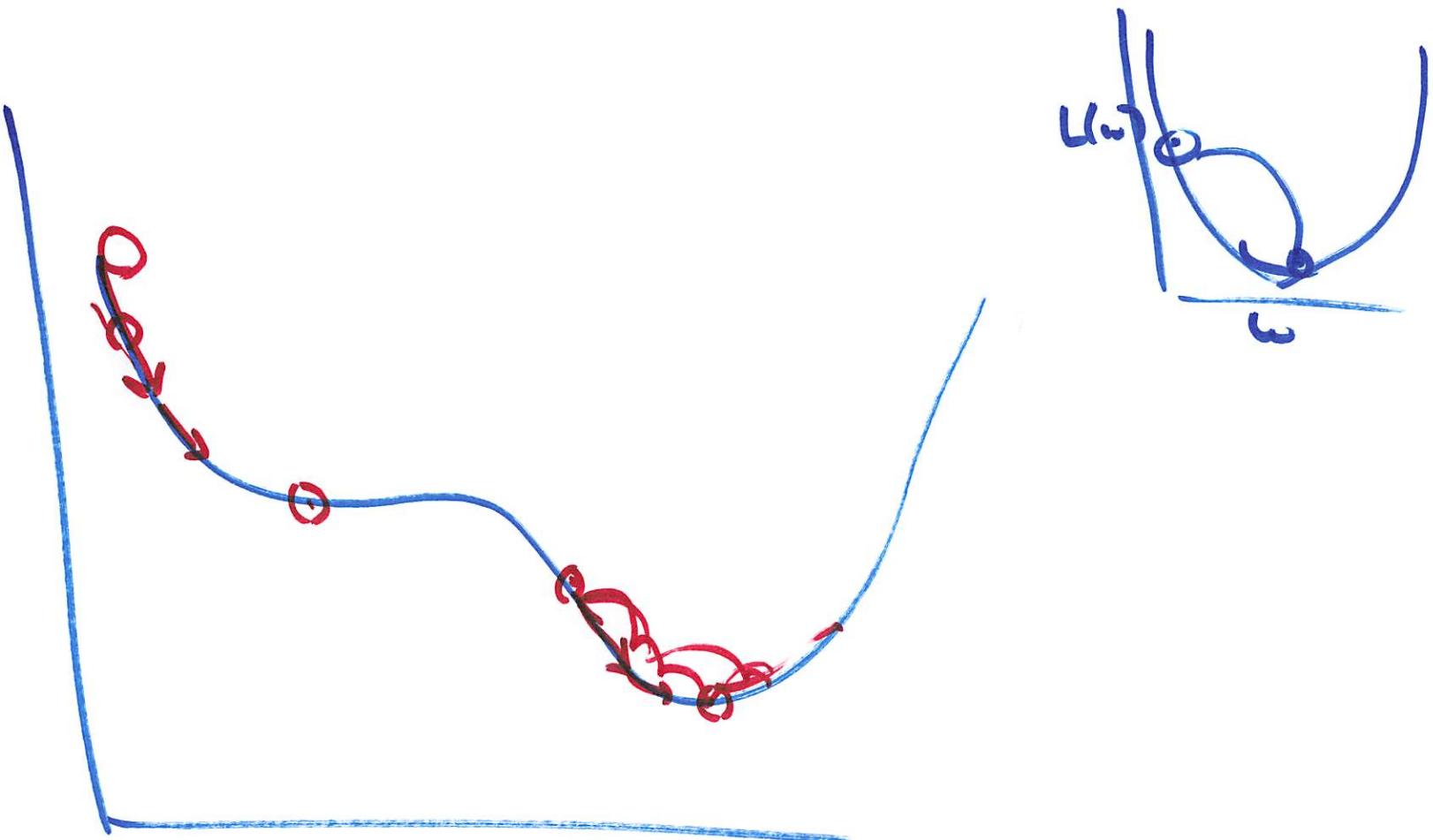
$$L(\omega) = \frac{1}{N_b} \sum_{i \in B} (\hat{y}_i - y_i)^2$$

batch size = N_b

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



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



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

SGD with momentum \Rightarrow

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

Learning rate

- Choosing the Learning rate (η)
 - 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

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

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

RMS Prop

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

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

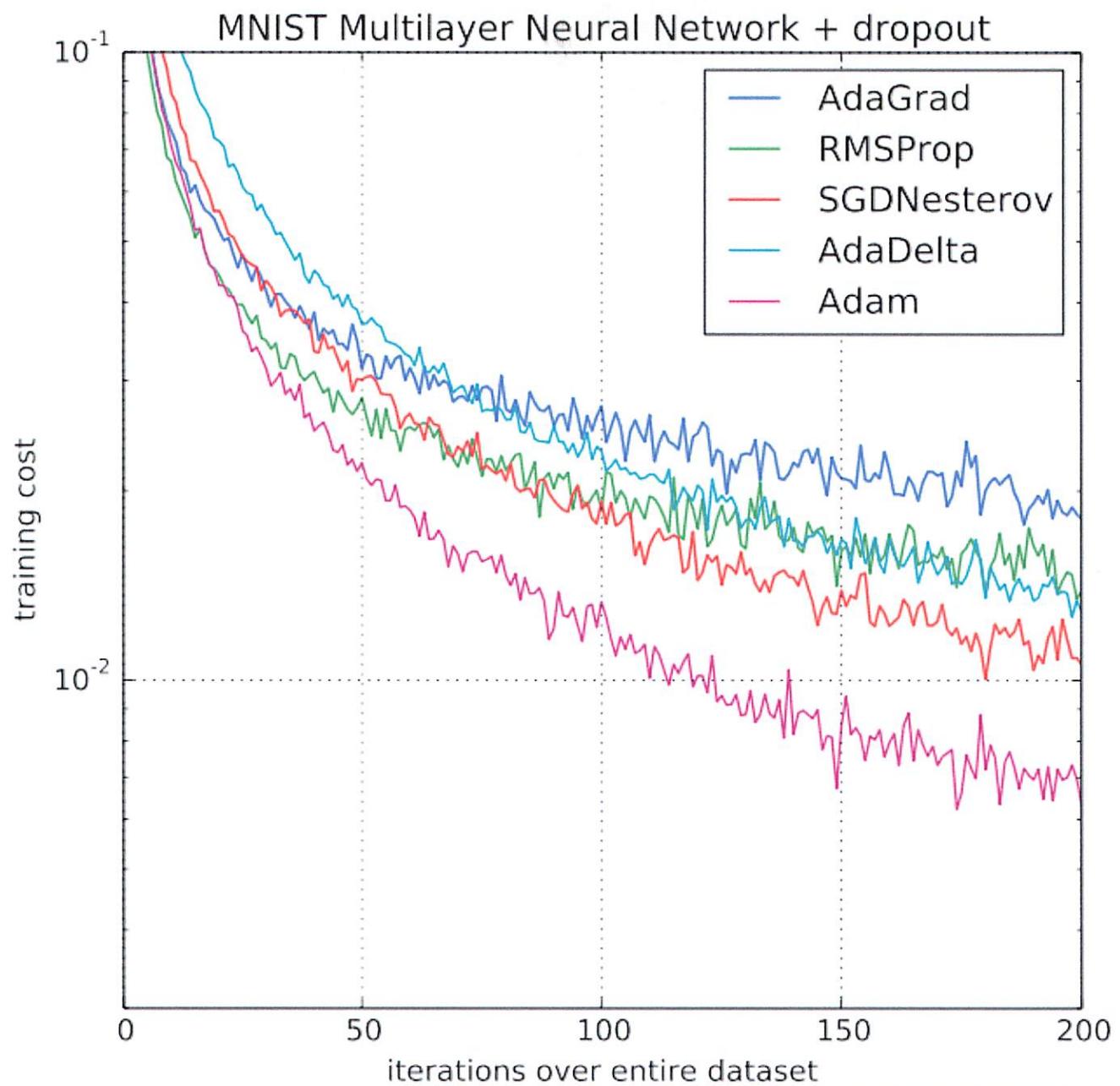
Adam

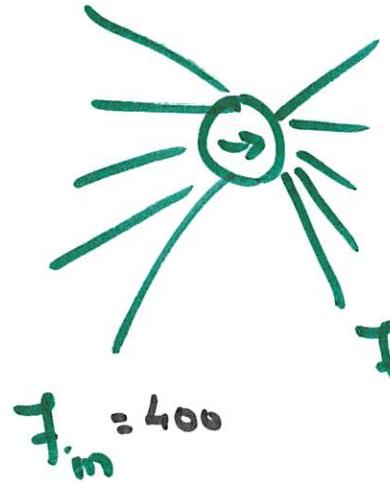


RMS Prop

+

SGD with Momentum

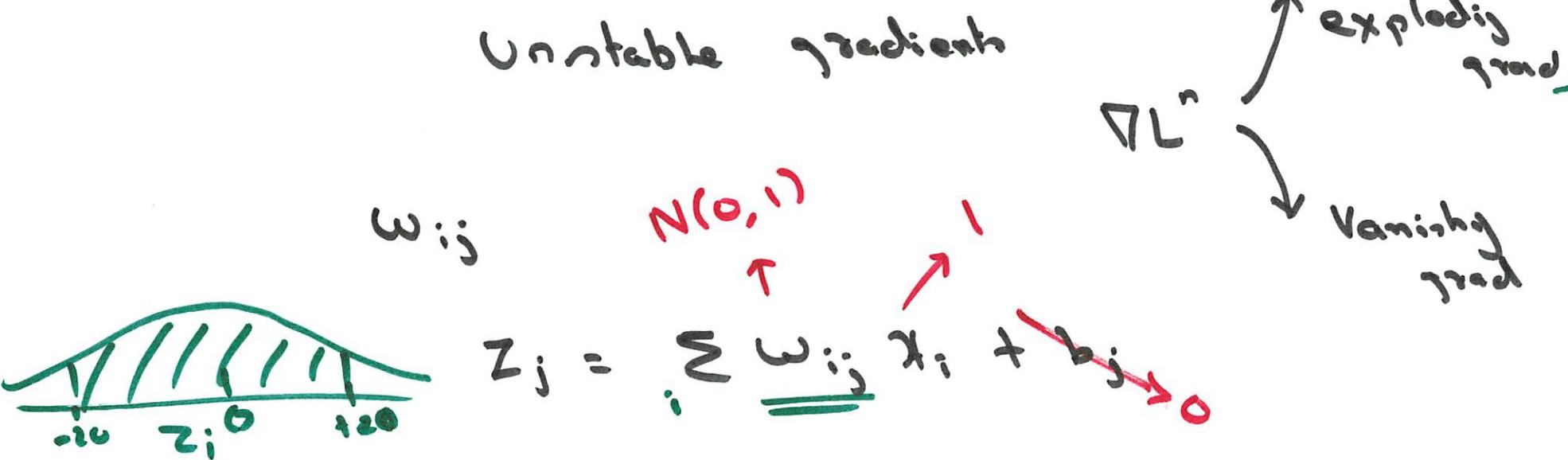




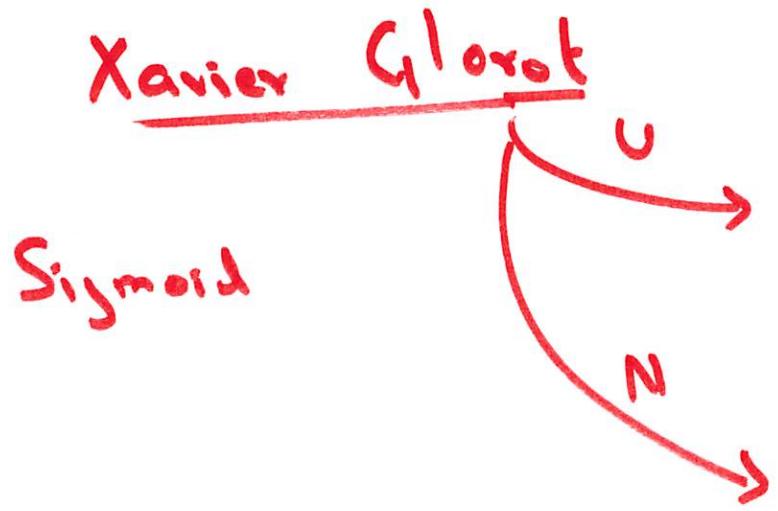
$$a_j = \sigma(\sum w_{ij} x_i + b_j)$$

$w_{ij} = 0$

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



$\text{Var}(z_j) = 400$ $\text{std dev}(z_j) = 20$



Uniform

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

Normal

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

A diagram showing the Xavier Glorot initialization rule for a layer with uniform or normal activation. It lists two ranges: a uniform range from $-\sqrt{\frac{6}{f_{in} + f_{out}}}$ to $+\sqrt{\frac{6}{f_{in} + f_{out}}}$, and a normal range centered at 0 with a standard deviation of $\sqrt{\frac{2}{f_{in} + f_{out}}}$.

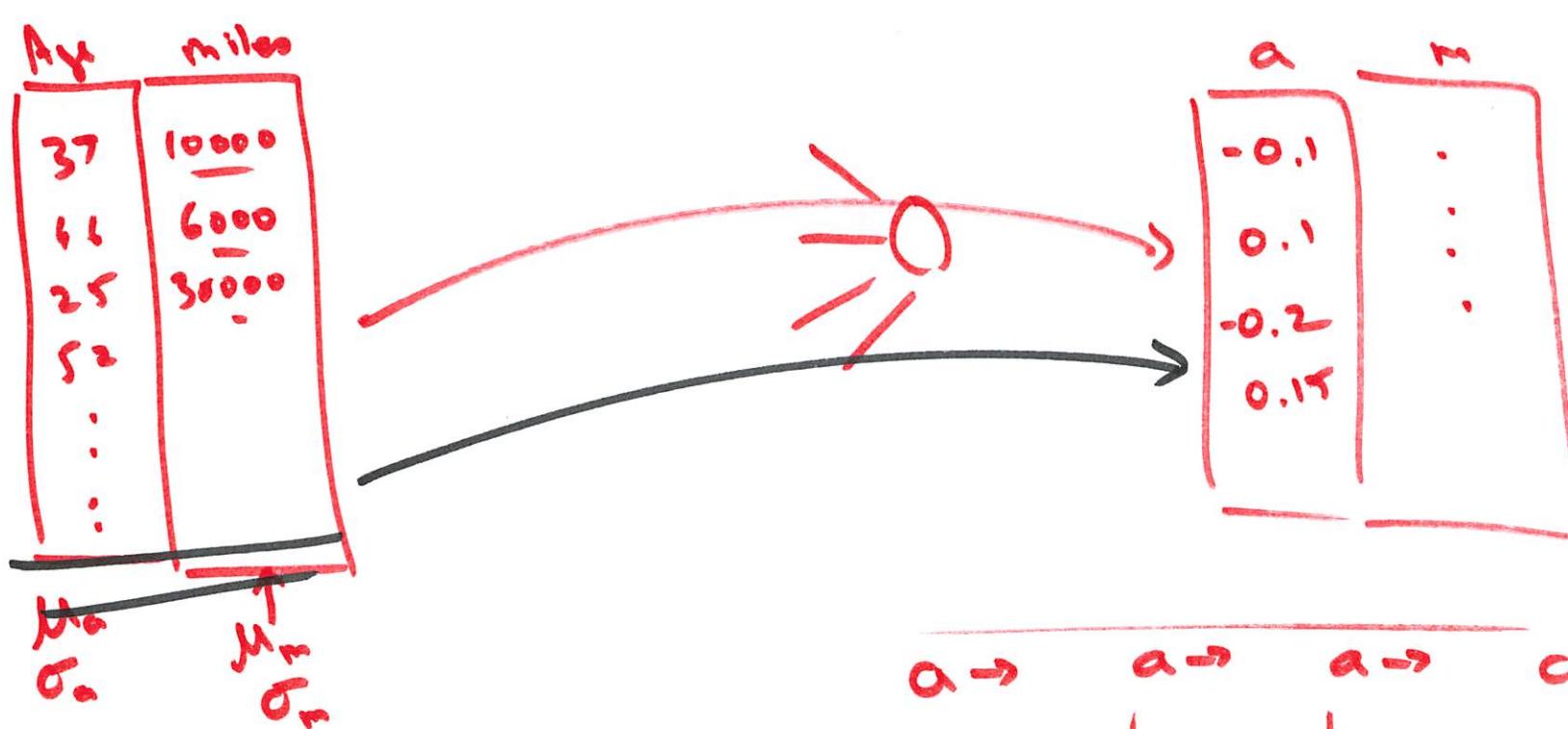
ReLU

$z \rightarrow c$

Normal

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

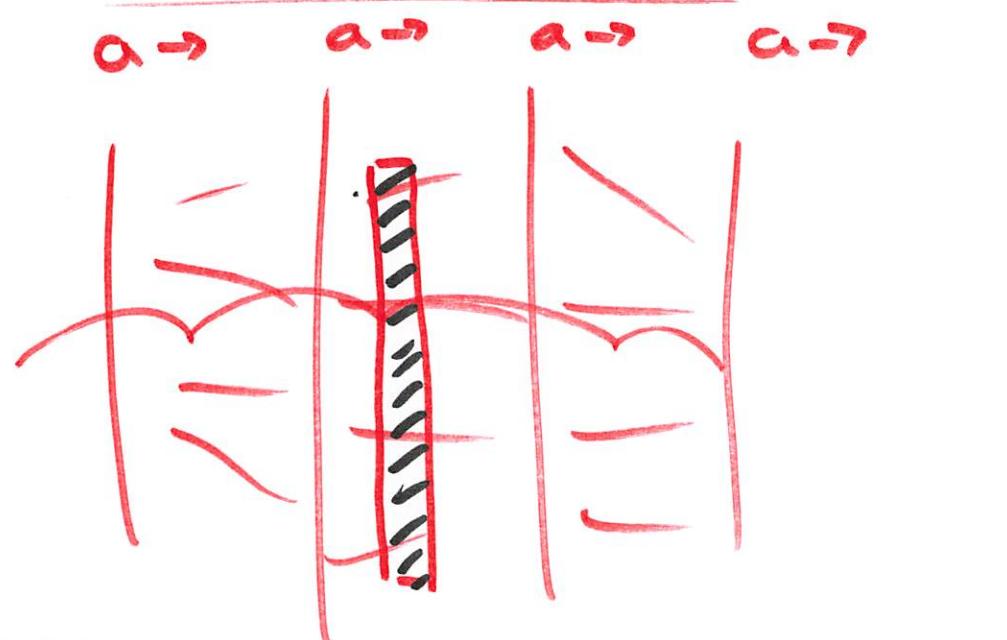
A diagram showing the Xavier Glorot initialization rule for a layer with ReLU activation. A red arrow labeled z points from the input to a point labeled c on a red curve representing the ReLU function. It specifies a normal distribution for the weights with a mean of 0 and a standard deviation of $\sqrt{\frac{2}{f_{in}}}$.



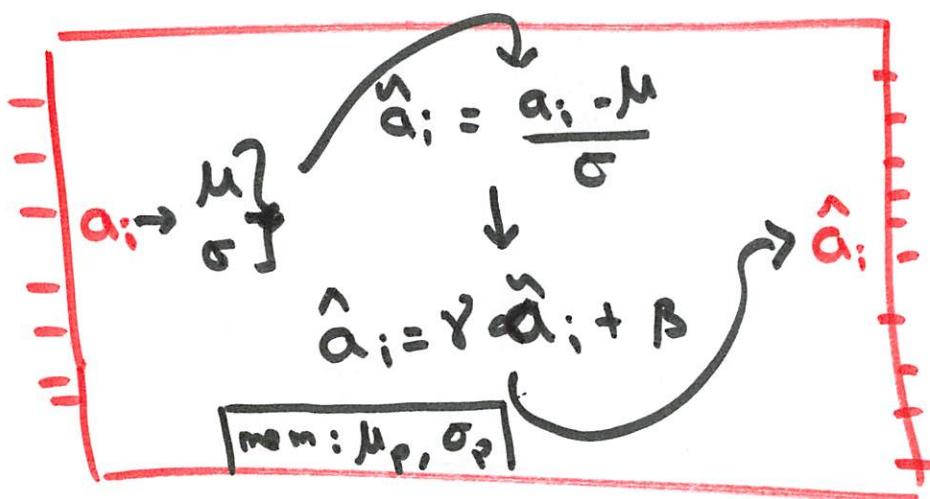
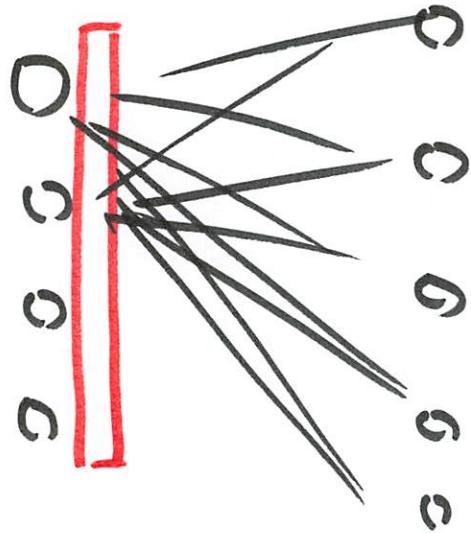
$$\hat{x}_a = \frac{x_a - \bar{X}_a}{\sigma_a}$$

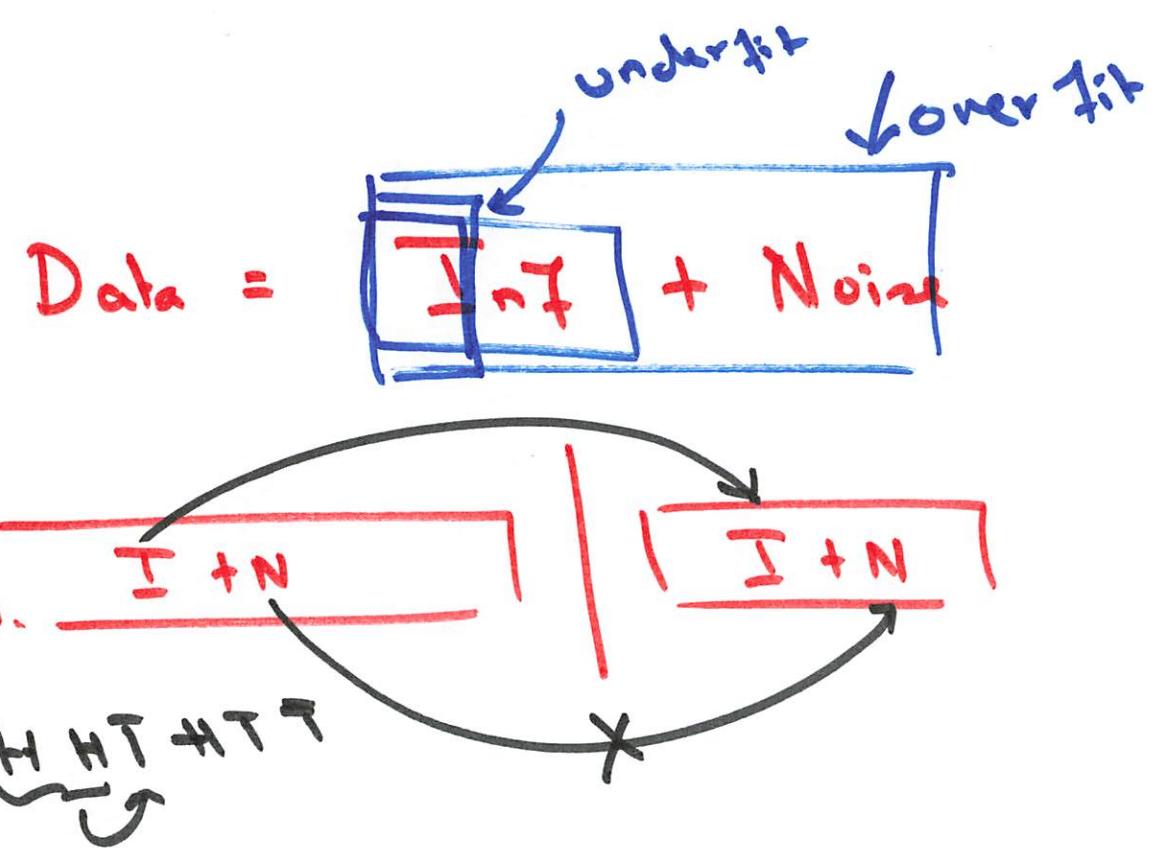
$$\hat{x}_m = \frac{x_m - \bar{X}_m}{\sigma_m}$$

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

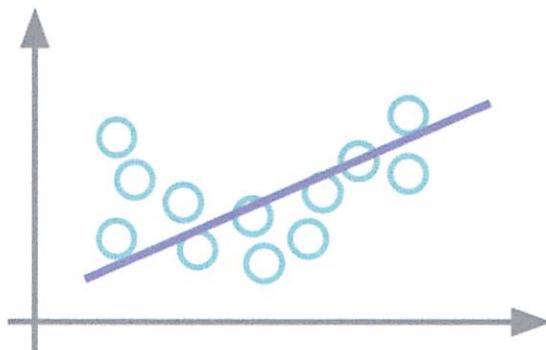


Covariate Shifts

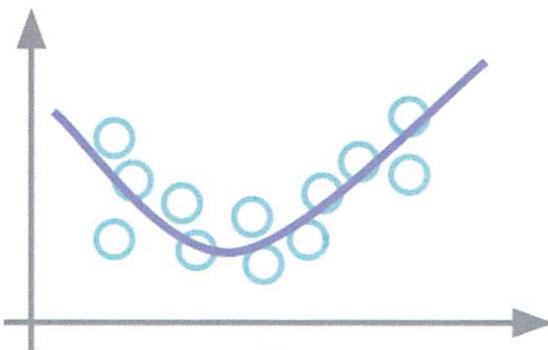




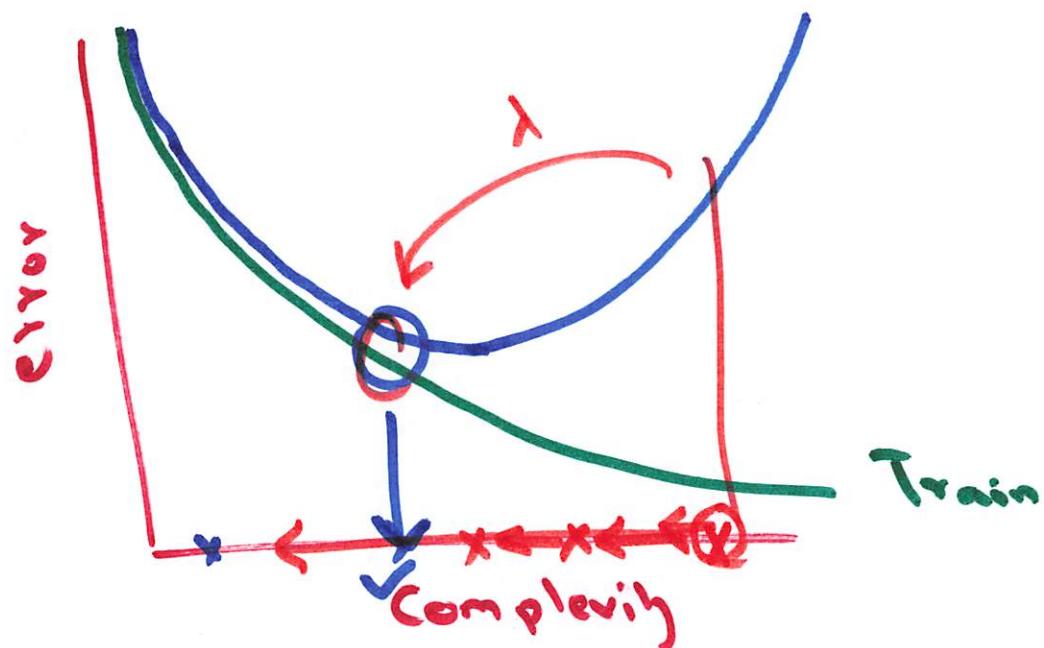
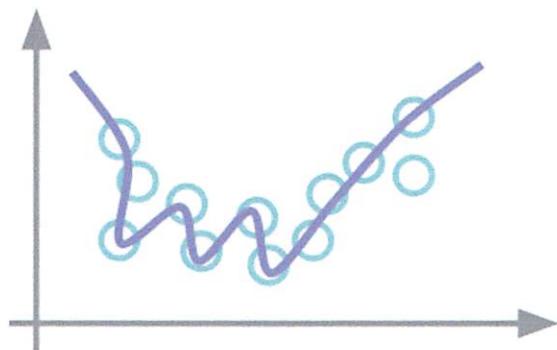
Under Vs Over-fitting



Underfit



overfit



$L_1 + L_2$ Reg.

∇

min

$$L(\beta) = \frac{1}{N} \sum (\gamma - g)^2 + \frac{\lambda}{\rho} (\text{penalty})$$

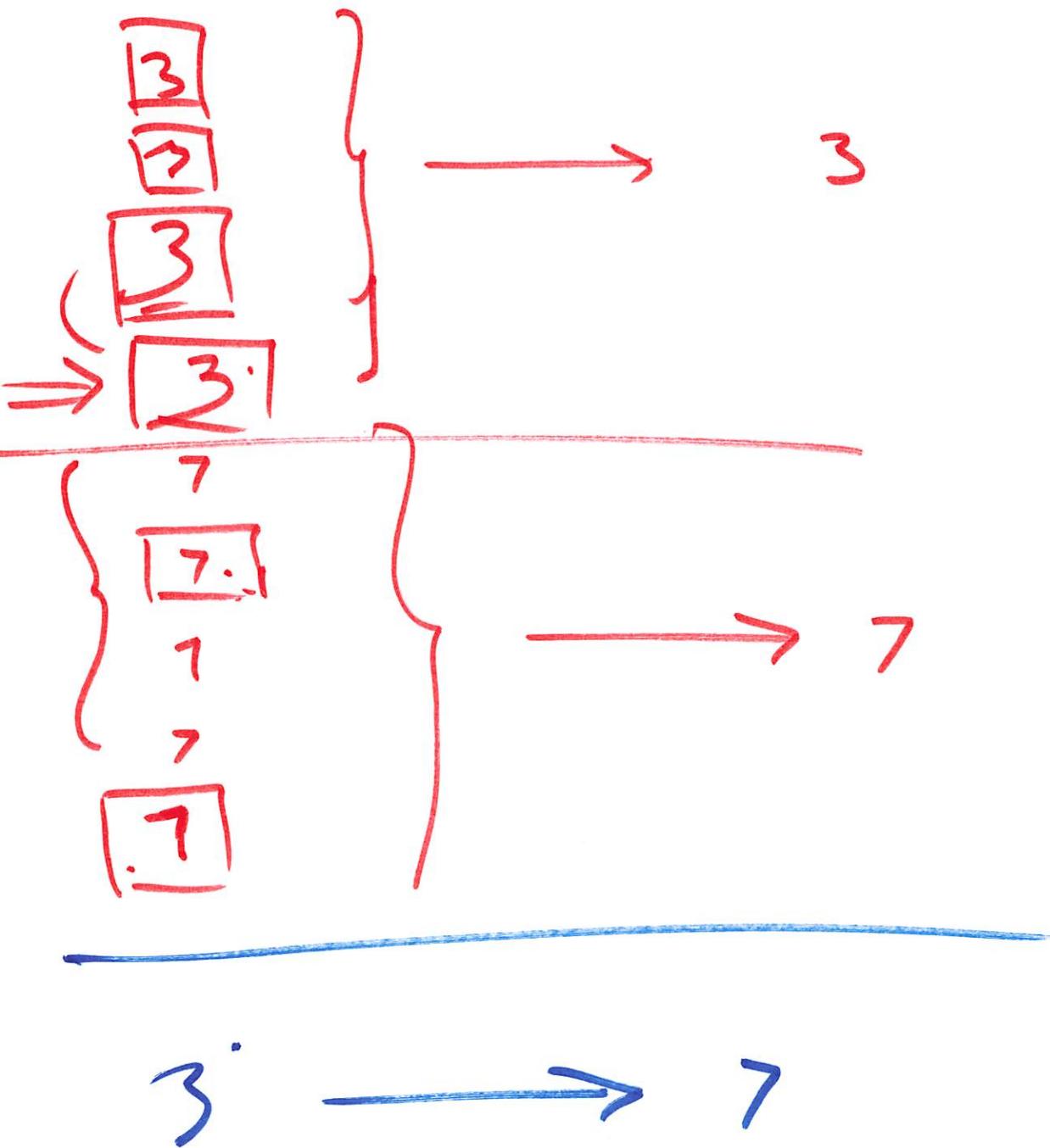
$$g = \sum \beta_i x_i + b$$

penalty

$$\rightarrow \sum |w_i|$$

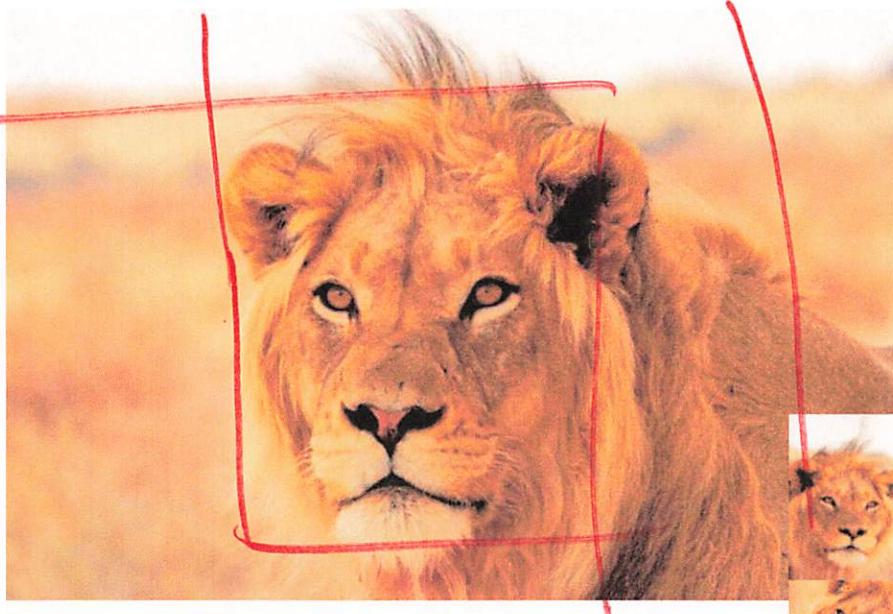
$$\rightarrow \sum (w_i)^2$$

L_1
Lasso =
 L_2
ridge

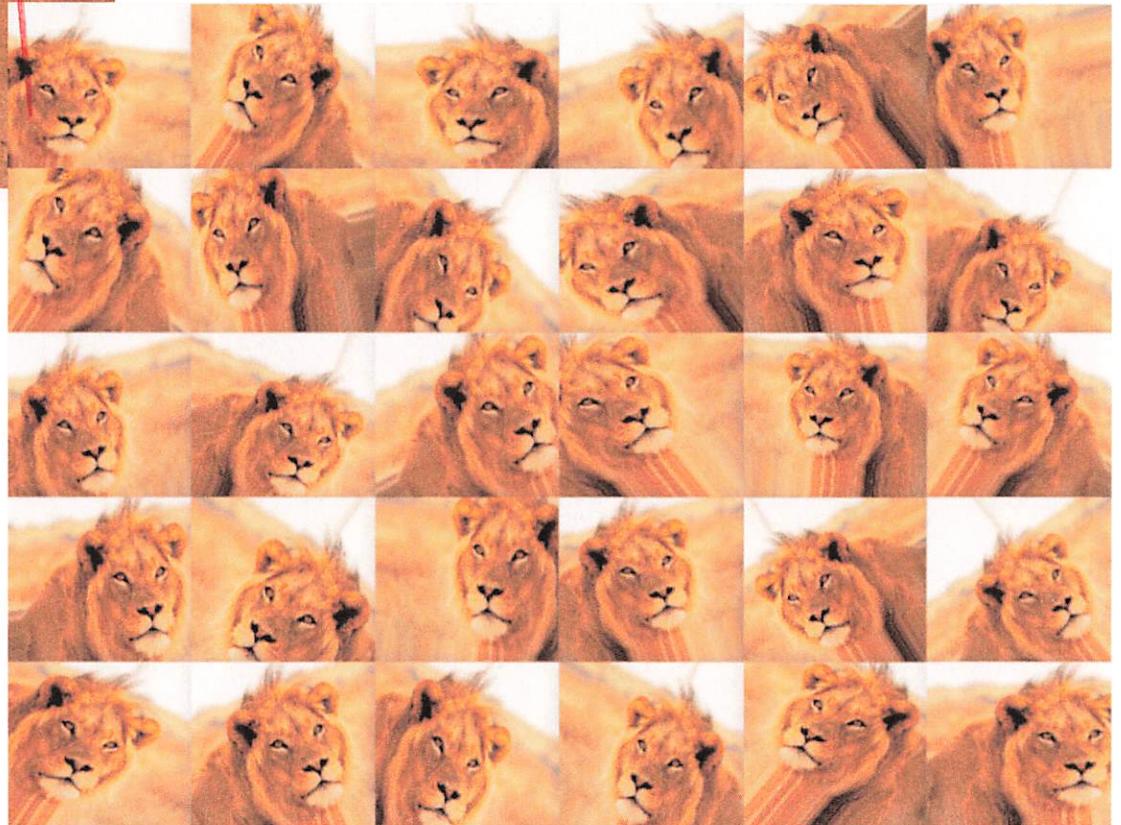


This file is meant for personal use by bpidugu@hotmail.com only.
Sharing or publishing the contents in part or full is liable for legal action.

Data Augmentation



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

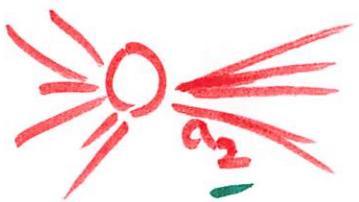
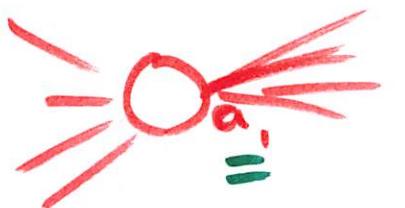


-source: [towardsdatascience.com](https://towardsdatascience.com/machinex-image-data-augmentation-using-keras-23a2a2a2e3d0), MachineX: Image Data Augmentation Using Keras
Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution

This file is meant for personal use by sriyugu@hotmail.com only.

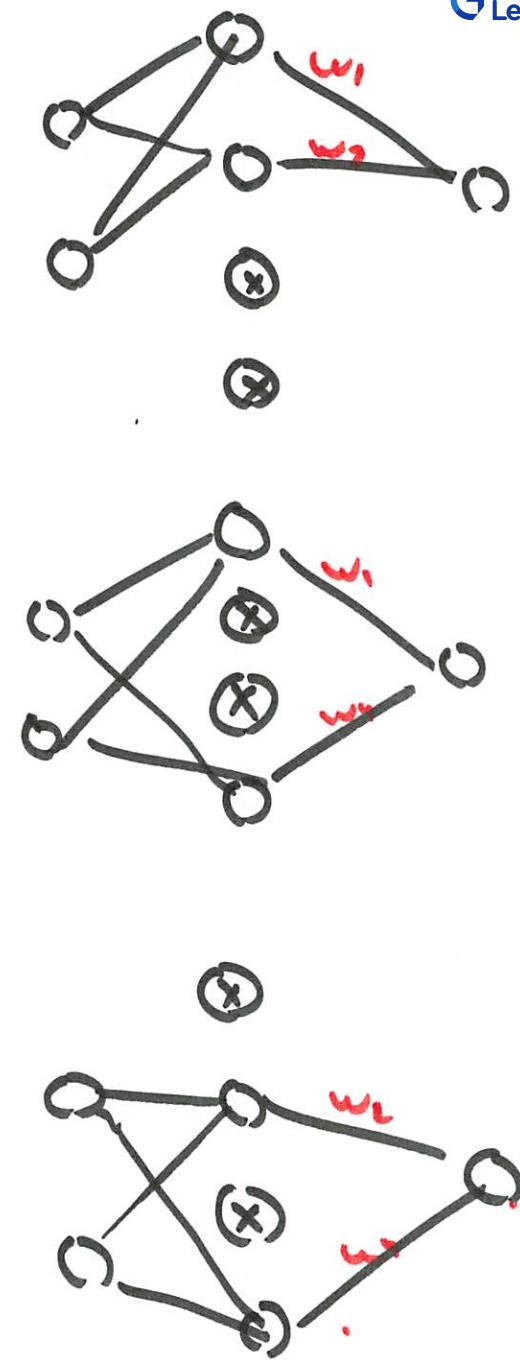
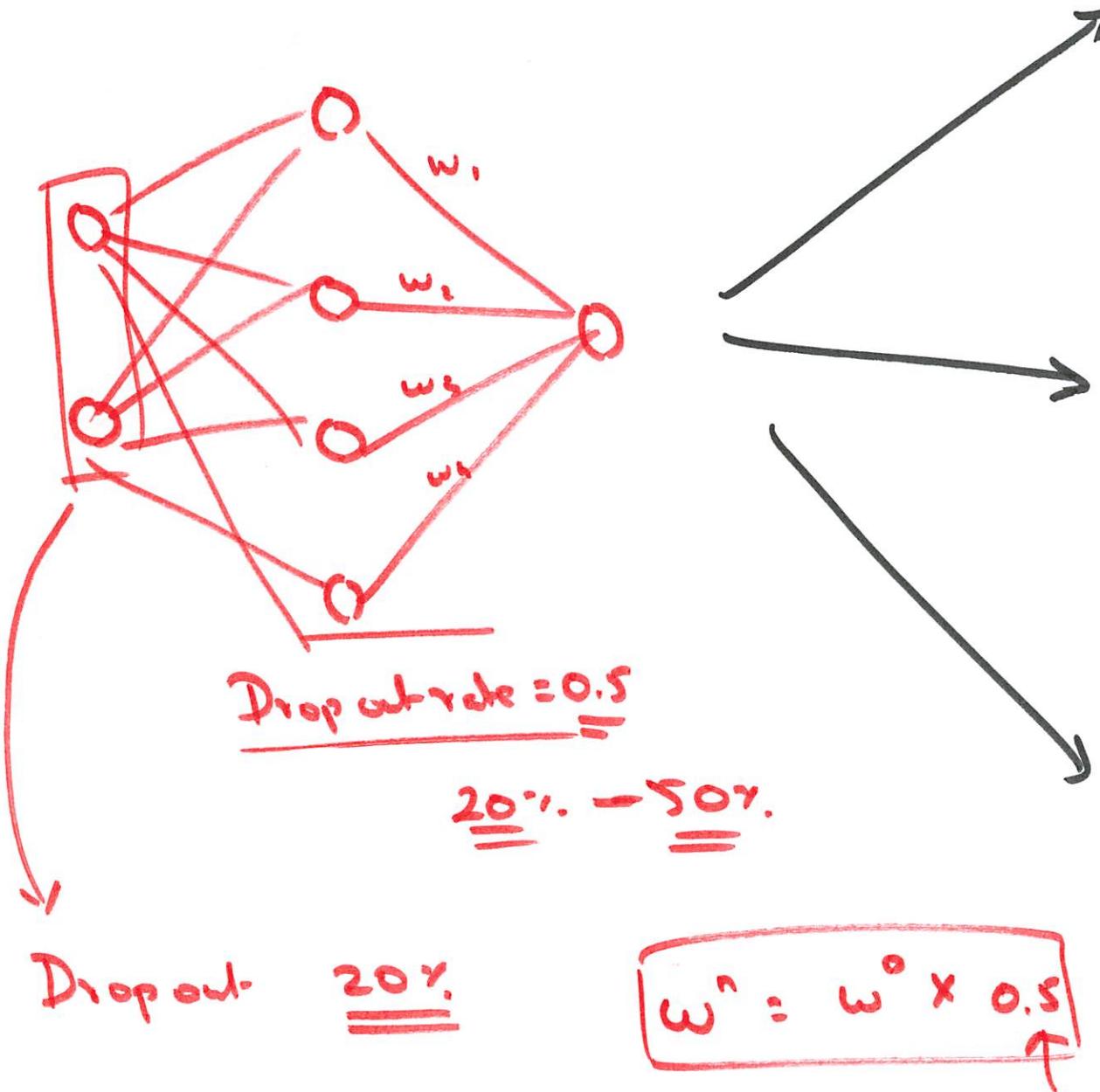
Sharing or publishing the contents in part or full is liable for legal action.

Co Adaptation



$$\overbrace{-0.1a_1 + 0.1a_2}$$

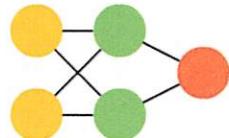
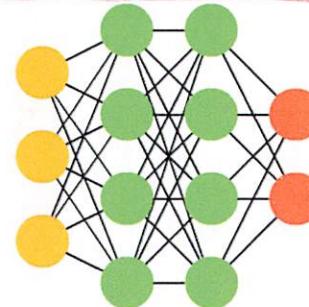
$$\overbrace{-0.3a_1 + 0.3a_2}$$



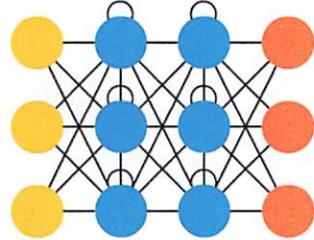
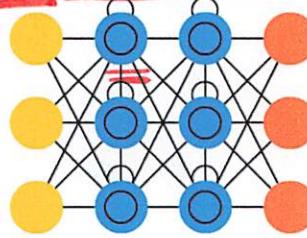
Types of NN

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

Feed Forward (FF)

Deep Feed Forward (DFF)

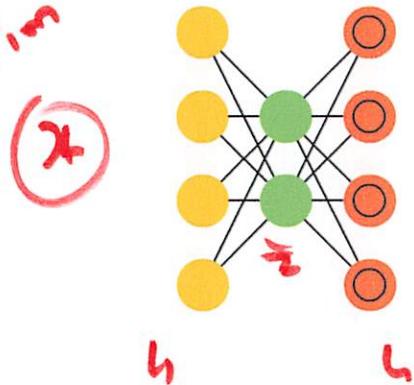
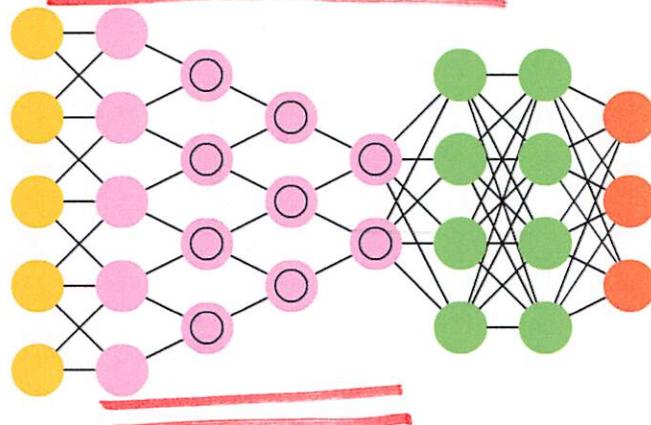
Recurrent Neural Network (RNN)

Long / Short Term Memory (LSTM)

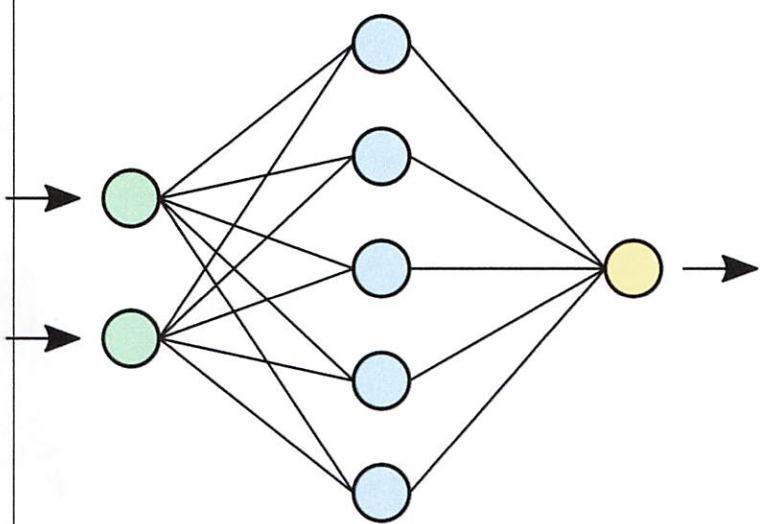
17
1
17
17

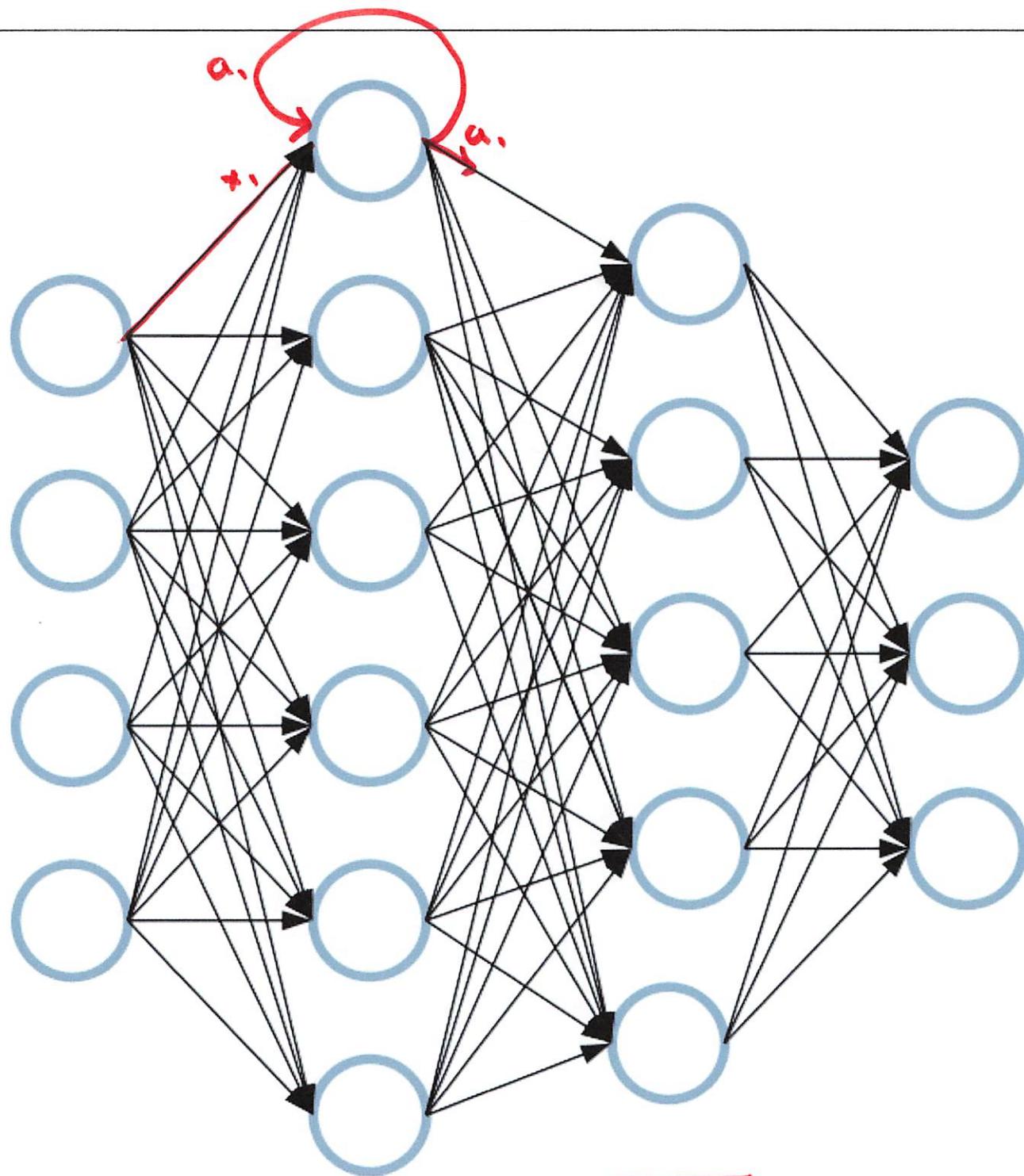
*Compressed
dim*

Auto Encoder (AE)

Deep Convolutional Network (DCN)

Feed Forward Net.





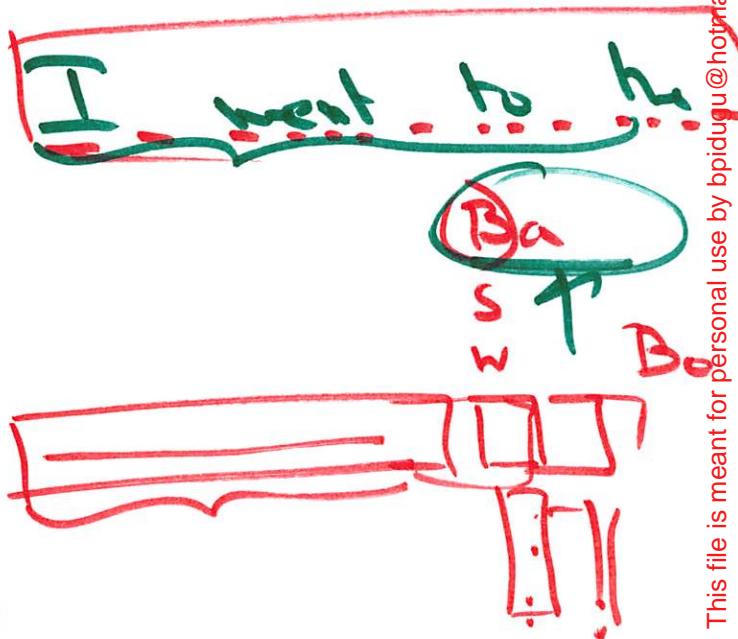
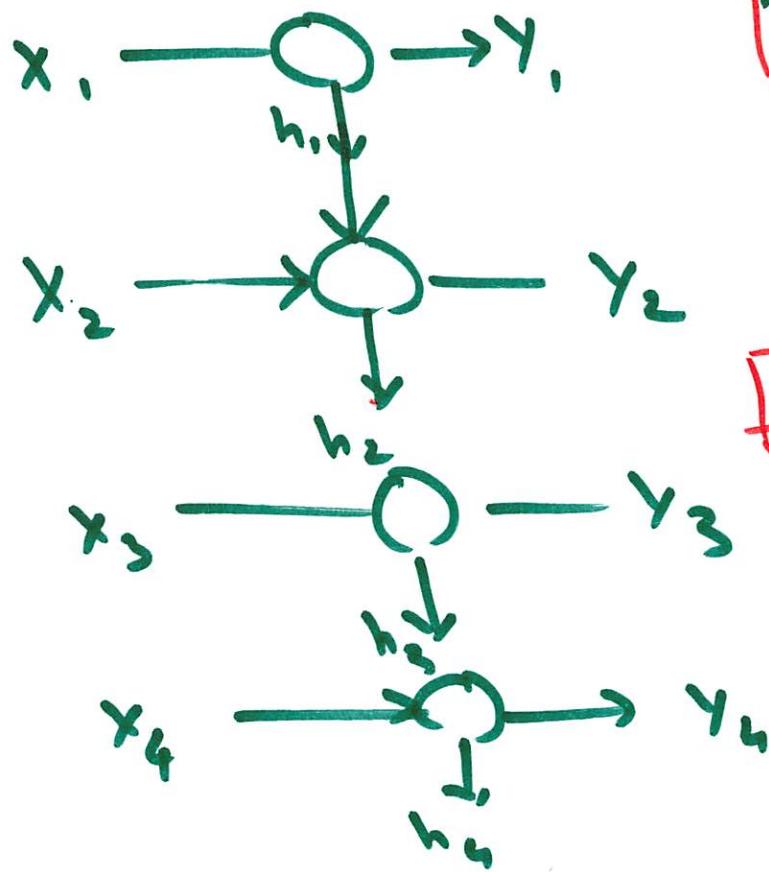
This file is meant for personal use by bpidugu@hotmail.com only.
Sharing or publishing the contents in part or full is liable for legal action.



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

$$y = f(\underline{\omega} \underline{x} + \underline{w}, \underline{b})$$

$$\underline{w} = f(\underline{\omega})$$

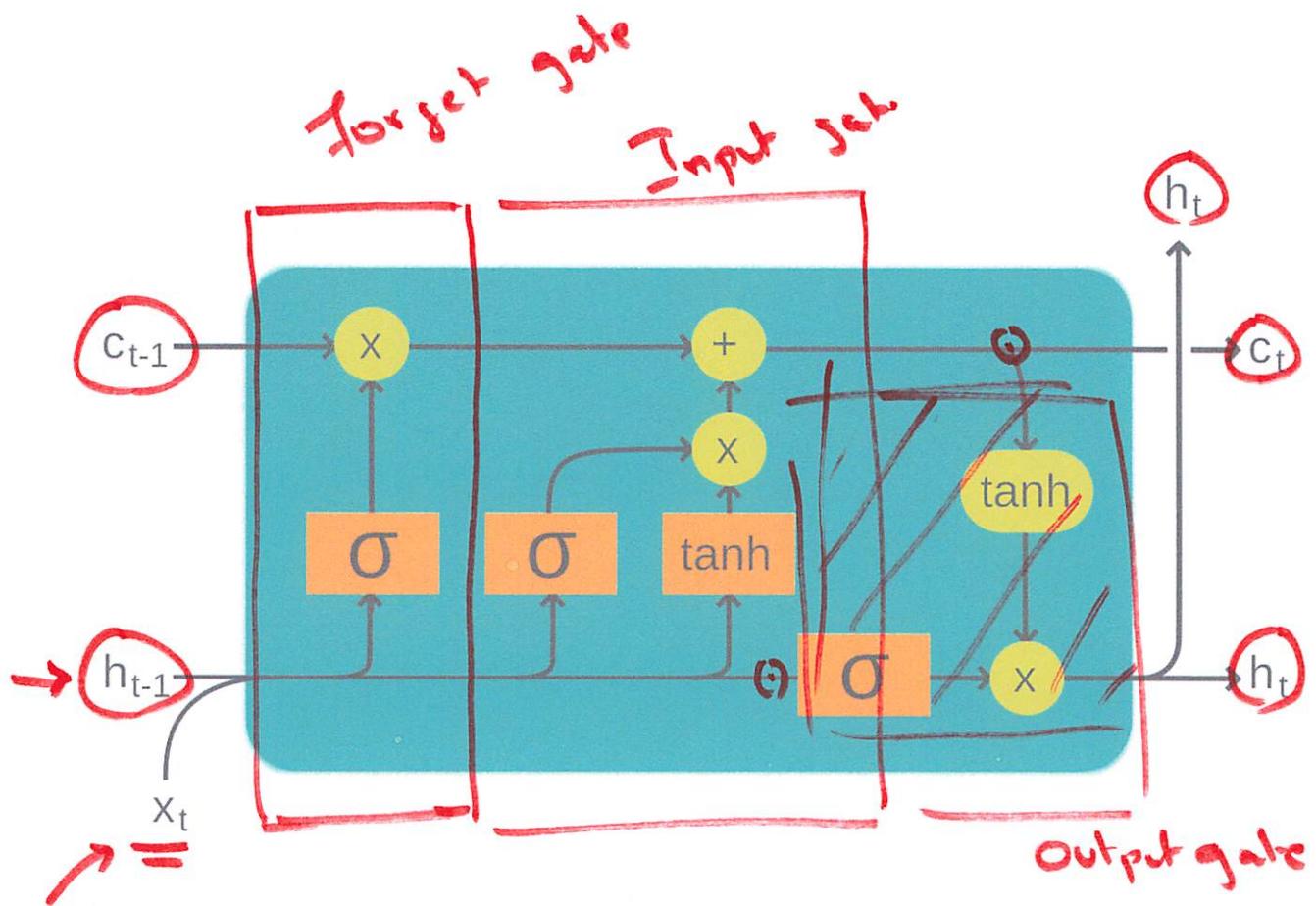


86
==

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:



