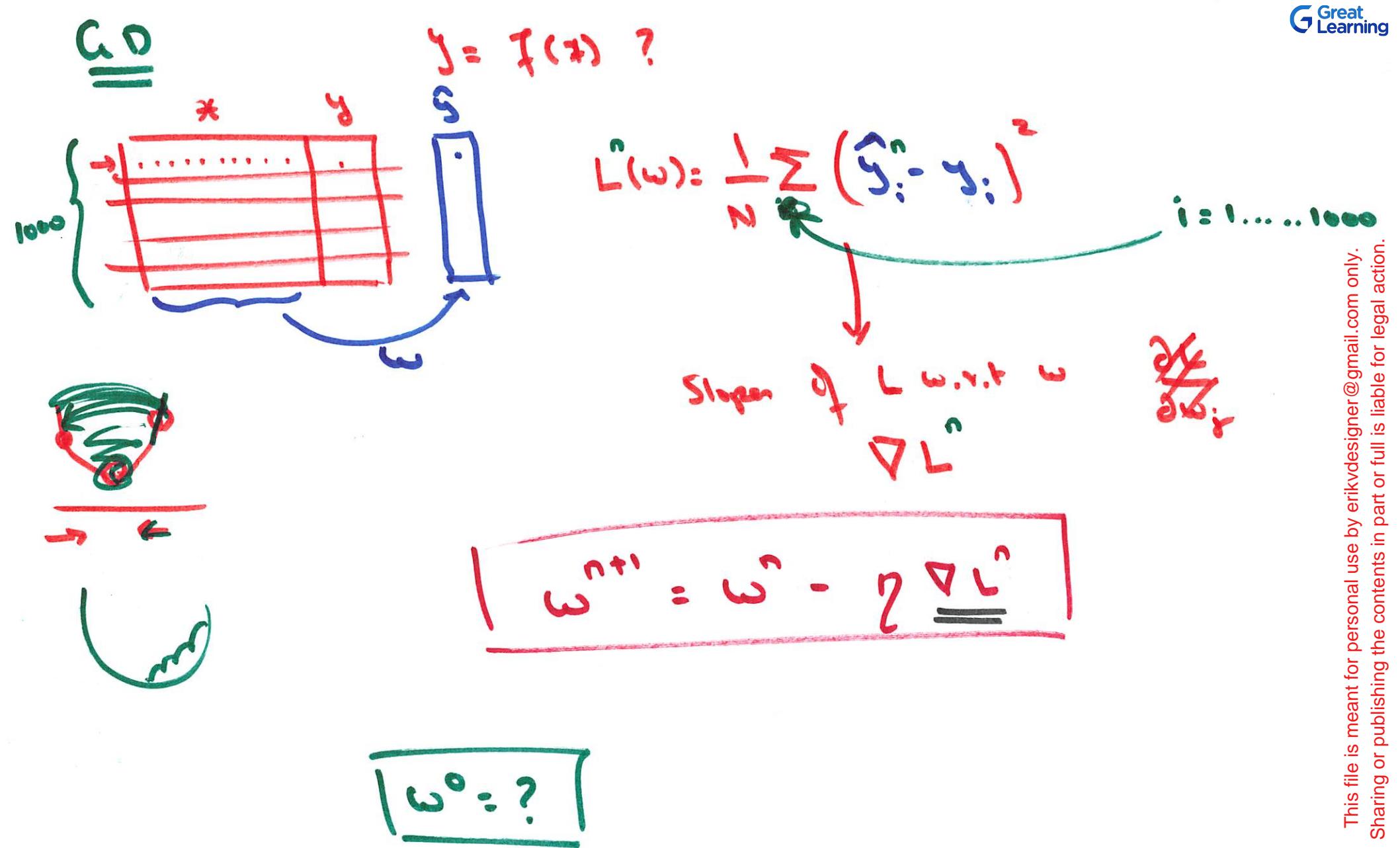


Agenda

- 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

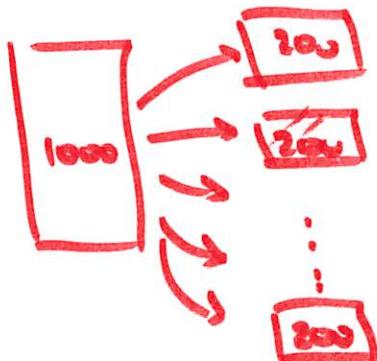
$$N = 1000$$

$$N_b = 200$$

$$iter = 5$$

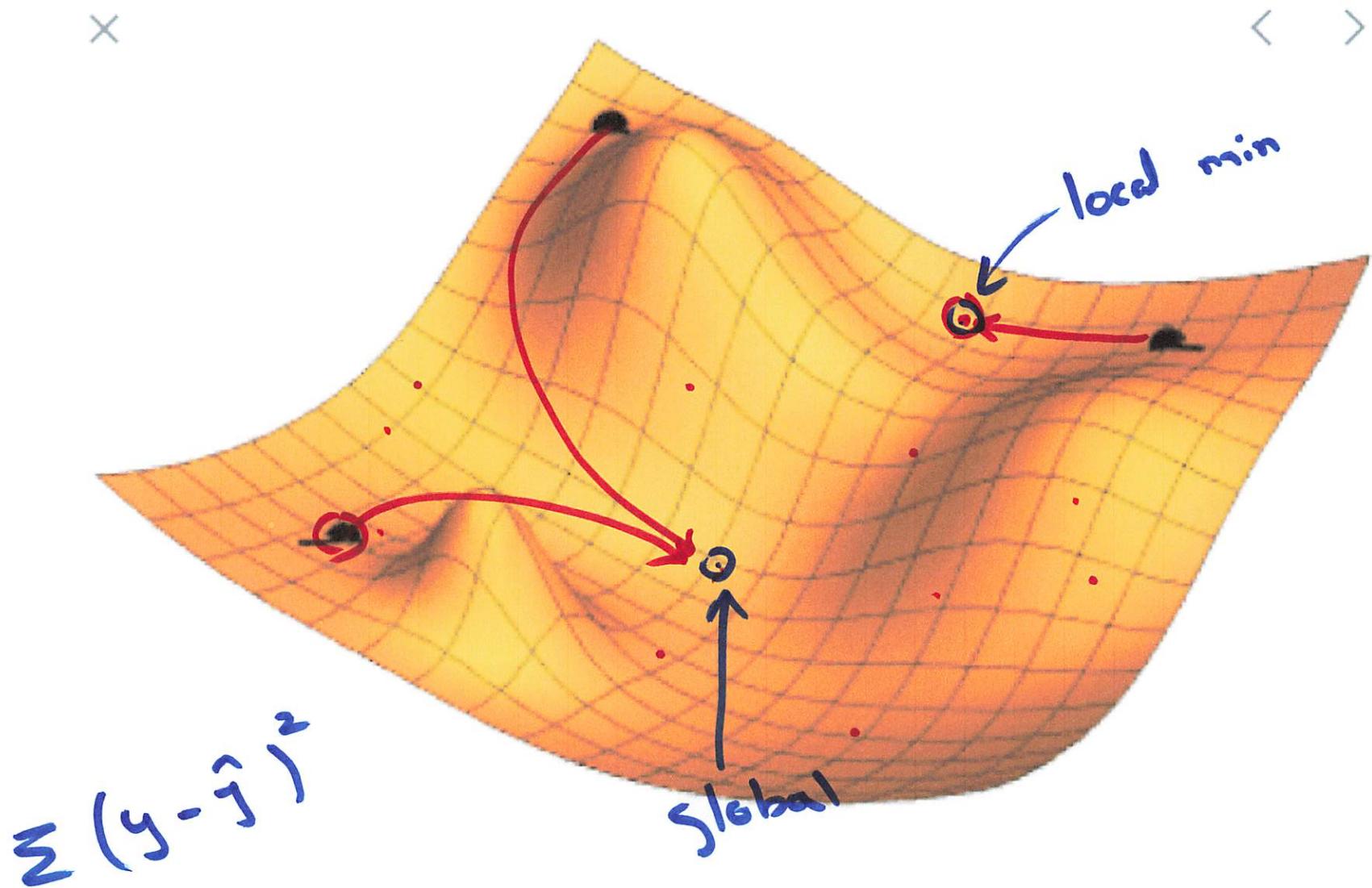
epoch

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

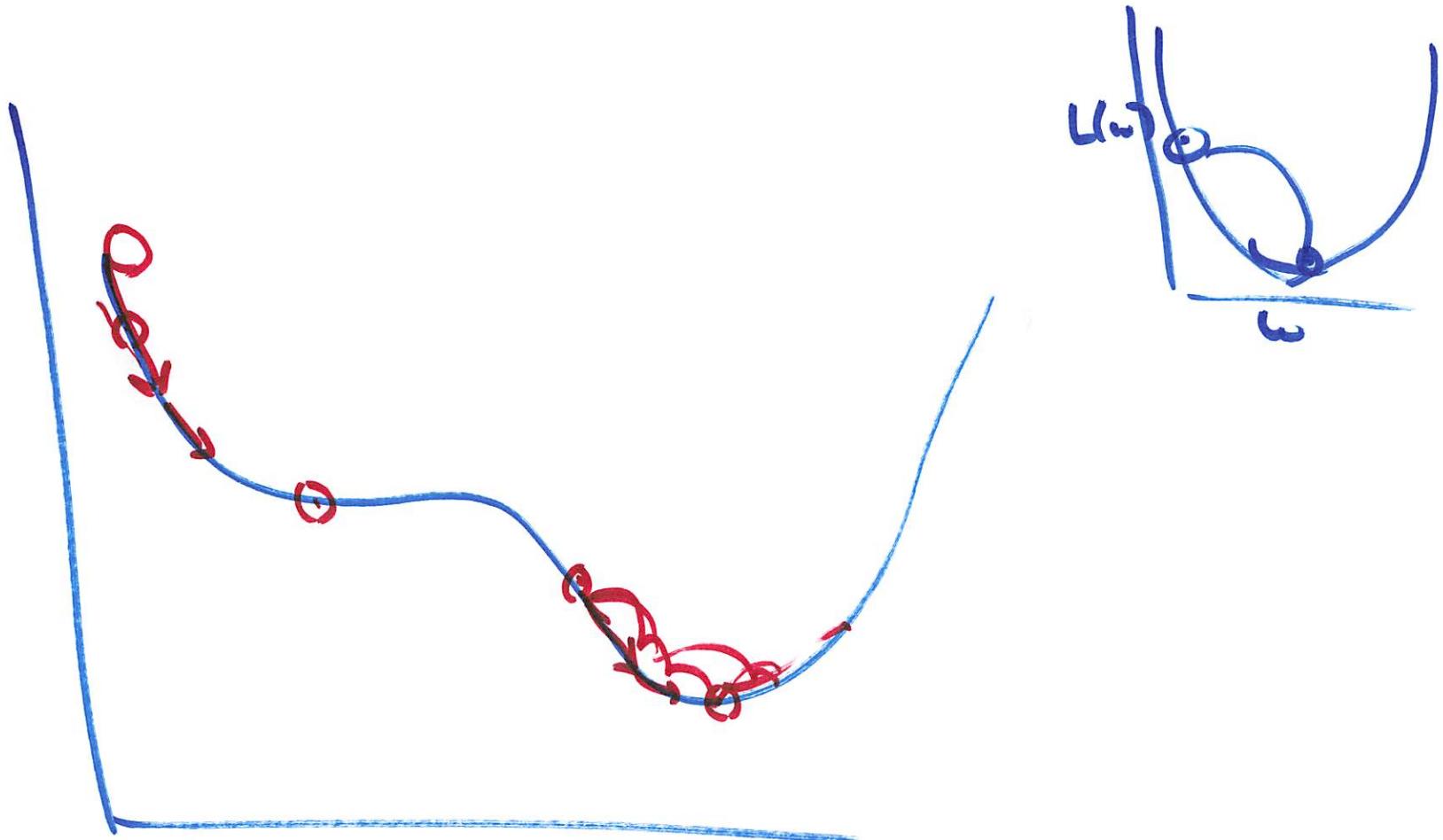


$$L(\omega) = \frac{1}{N_b} \sum_{i \in B} (\hat{y}_i - y_i)^2$$

batch size = N_b



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



$$E'' = \omega - 2 \Delta'$$

SGD with momentum \rightarrow

$$\omega^{n+1} = \omega^n - \eta \left(\alpha \nabla L^n + (1-\alpha) \nabla L^{n-1} \right)$$

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

$$\hat{g}^t = g^t - \frac{2}{\sqrt{\delta^t + \epsilon}} \nabla L^t$$

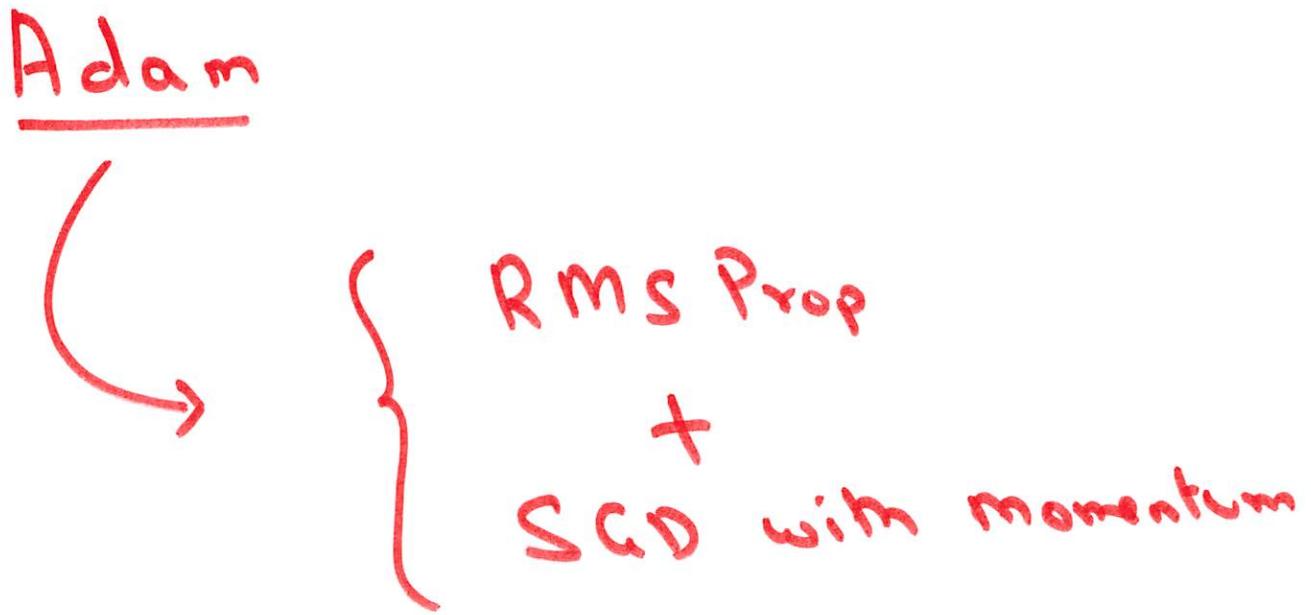
$$\delta^t = \sum_i (\nabla L^t)^i$$

Res Prop

$$\hat{g}^t = g^t - \frac{2}{\sqrt{\delta^t + \epsilon}} \nabla L^t$$

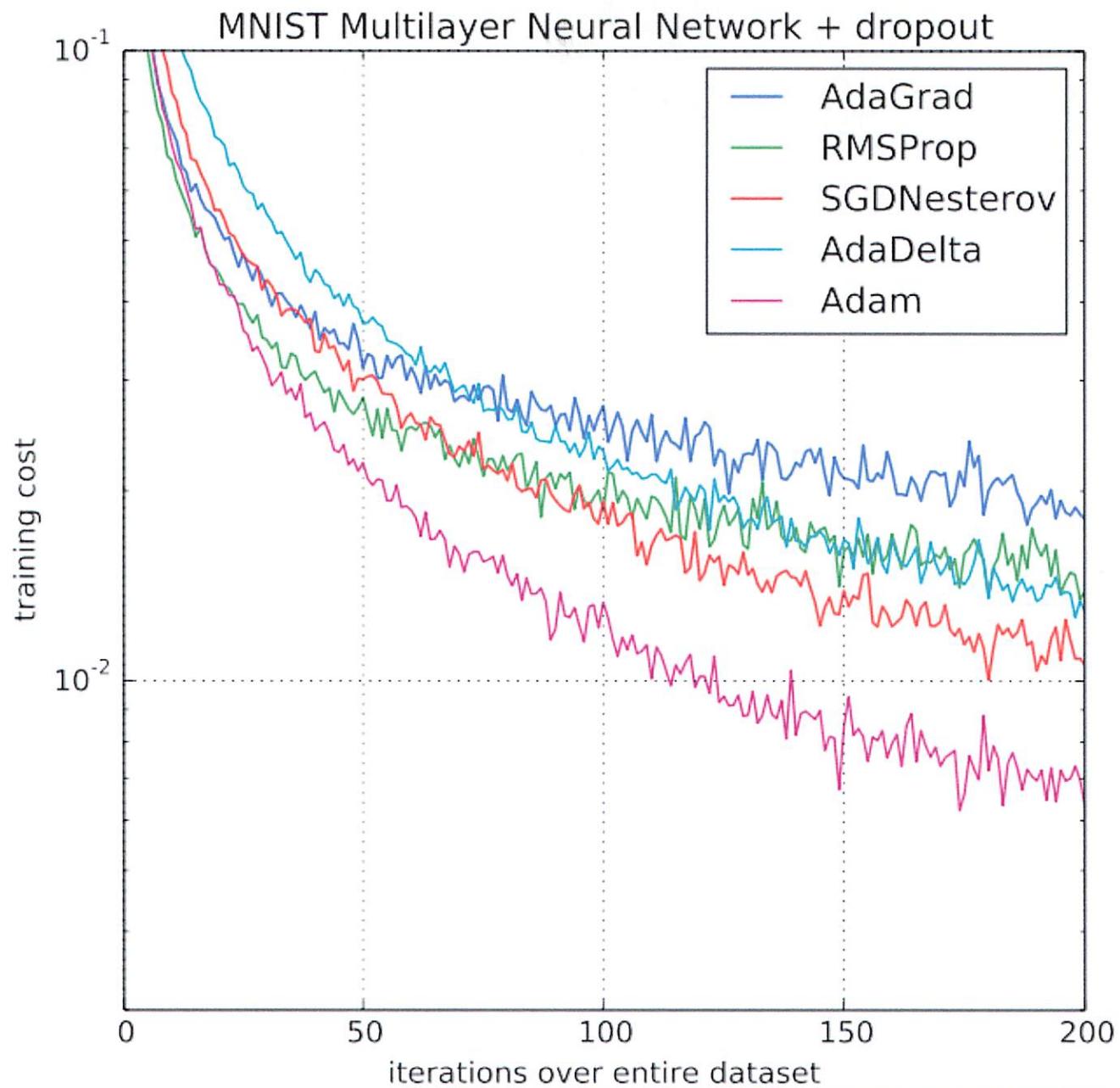
$$\delta^t = \alpha \delta^{t-1} + (1-\alpha) \nabla L^t$$

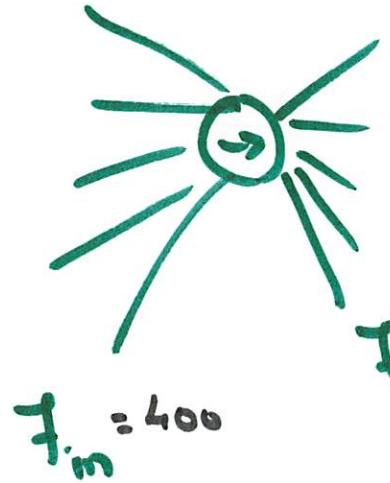
Adam



RMS Prop

+ SGD with Momentum



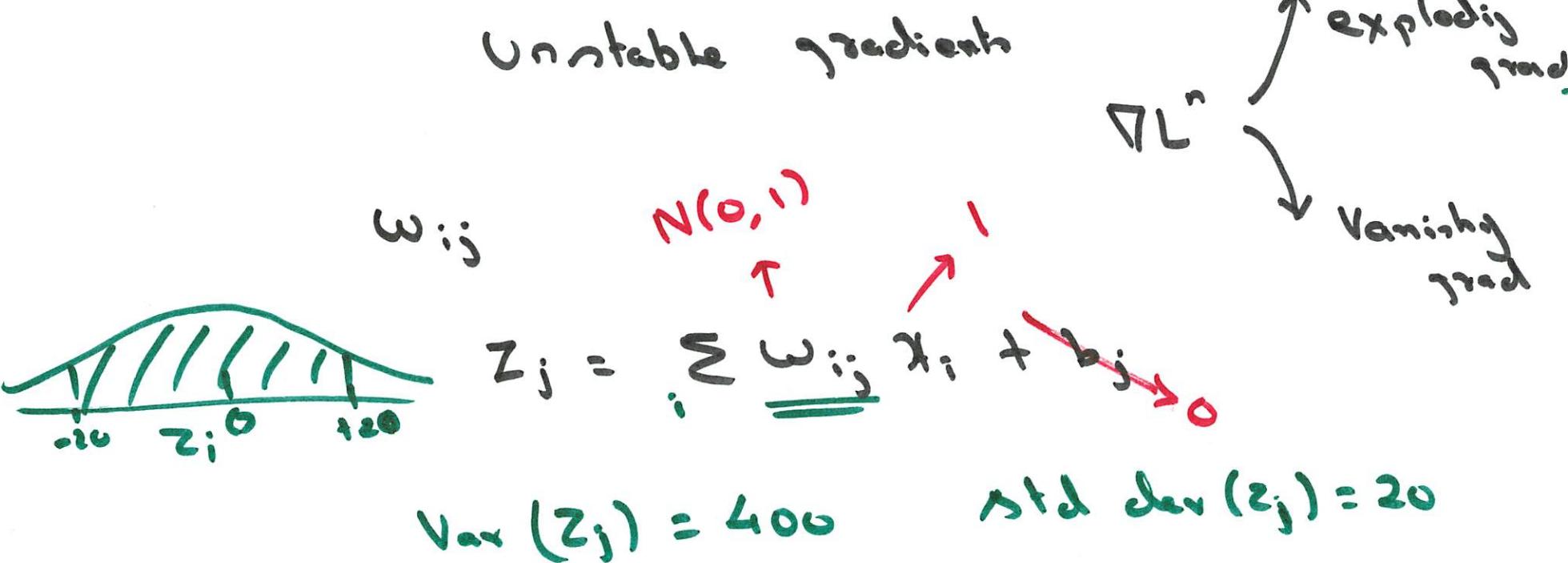


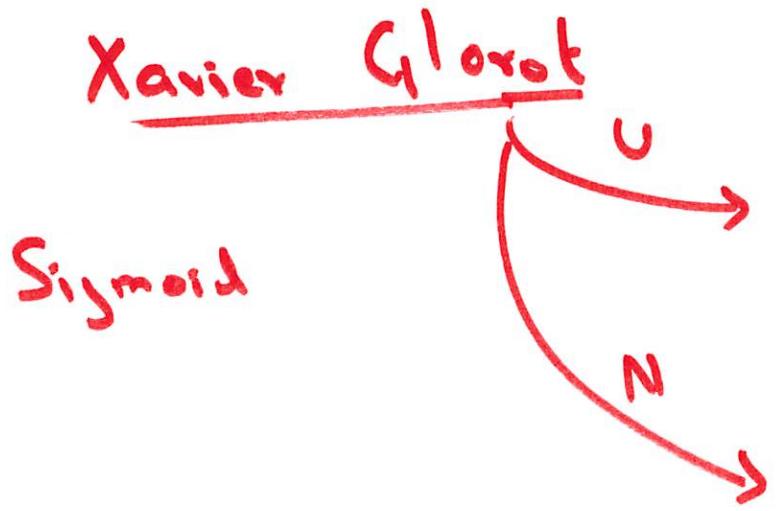
$$a_j = \sigma(z_j) = \sigma(\sum w_{ij} x_i + b_j)$$

Graph of a sigmoid function $\sigma(z)$ with a red circle highlighting the point where the function is zero.

$$w_{ij} = 0$$

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



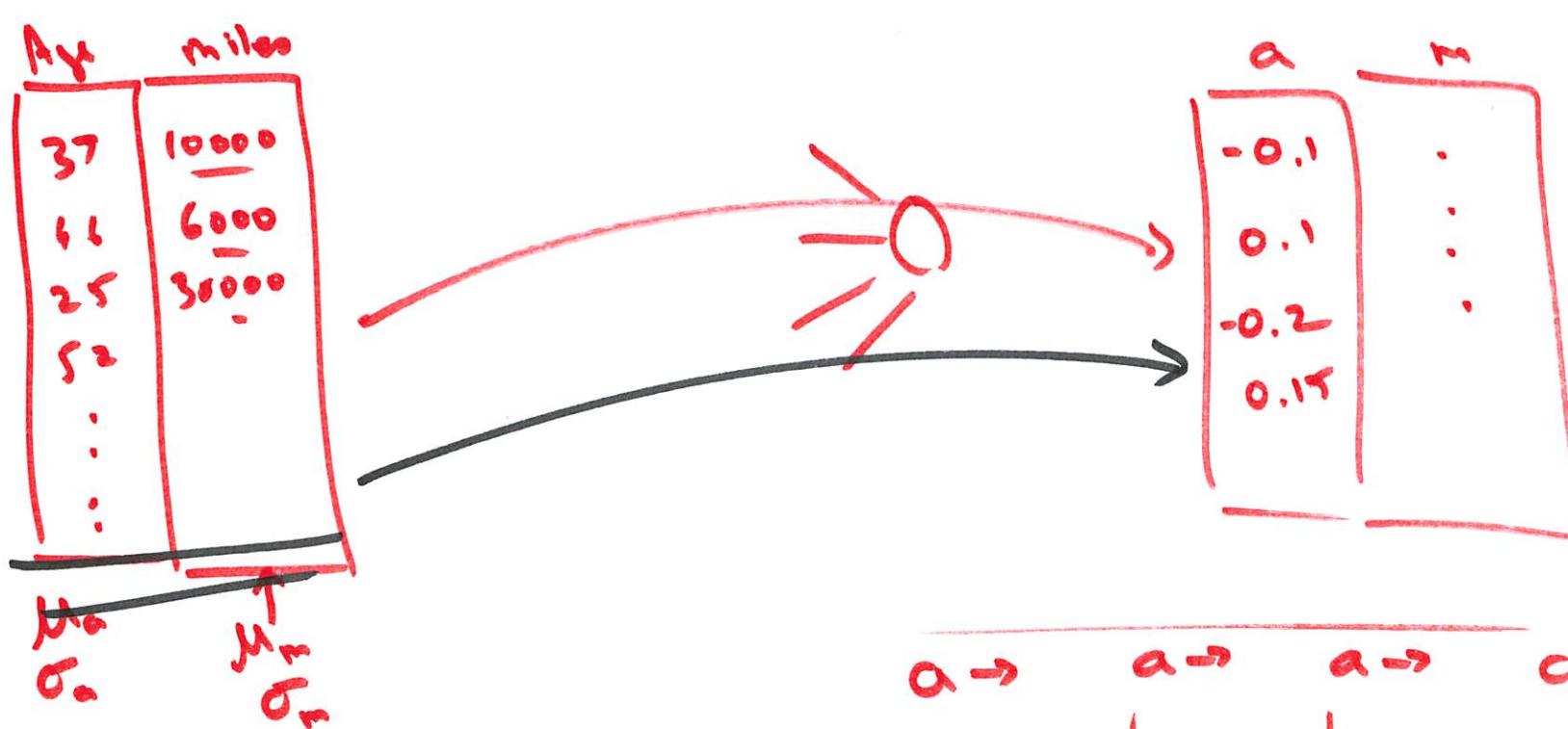


Uniform $\left[-\sqrt{\frac{6}{f_{in} + f_{out}}}, +\sqrt{\frac{6}{f_{in} + f_{out}}} \right]$

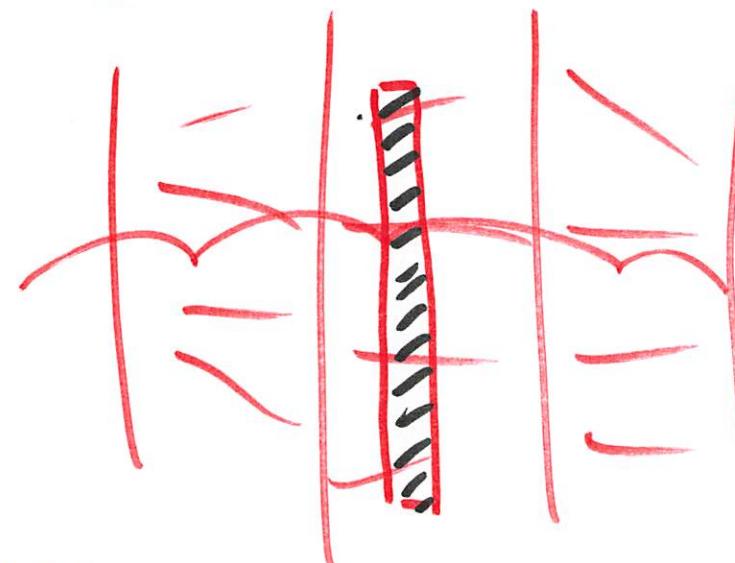
Normal $\left(0, \sqrt{\frac{2}{f_{in} + f_{out}}} \right)$

ReLU \xrightarrow{z}

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

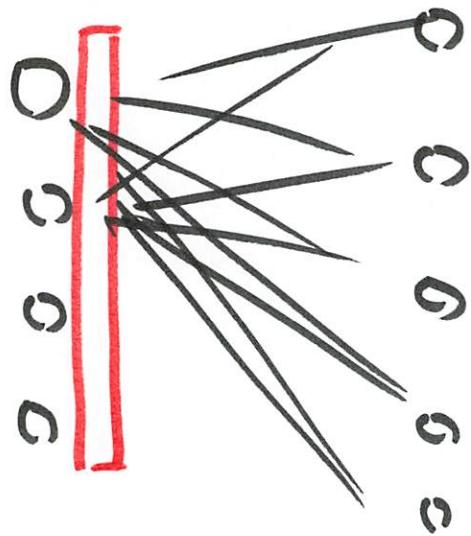


$$a \rightarrow a \rightarrow a \rightarrow a \rightarrow$$



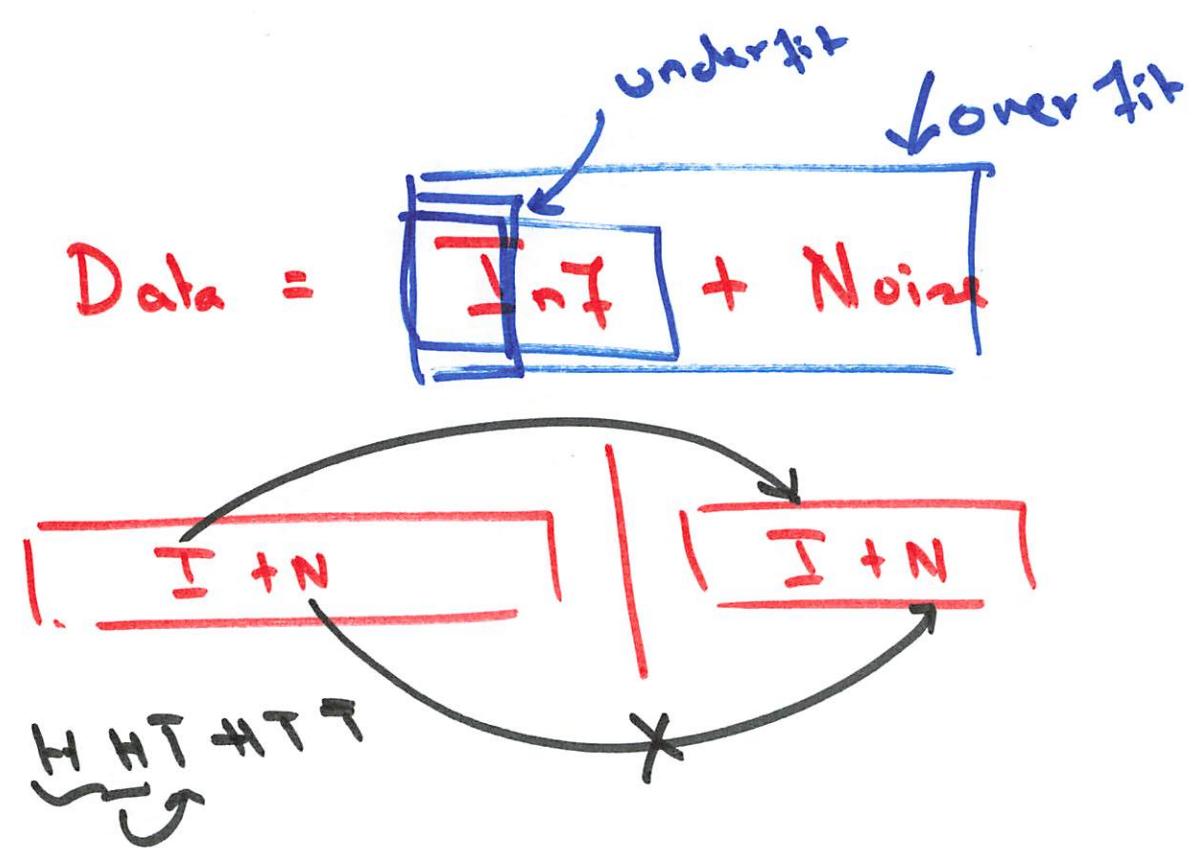
long → forward → stable

Covariate Shifts

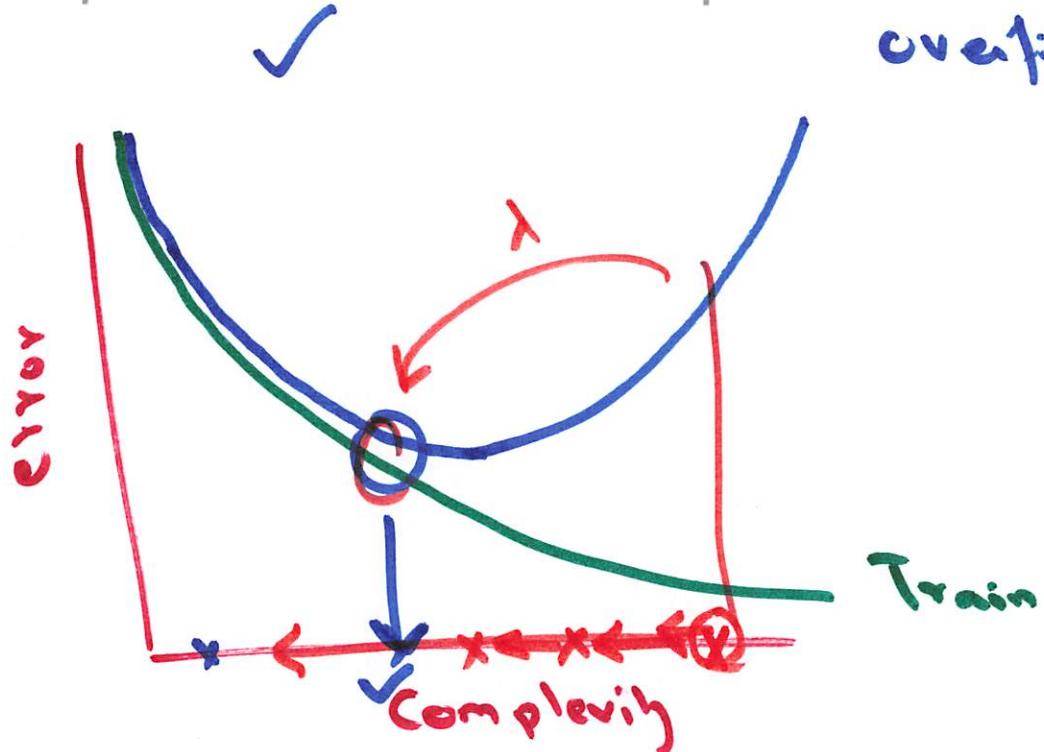
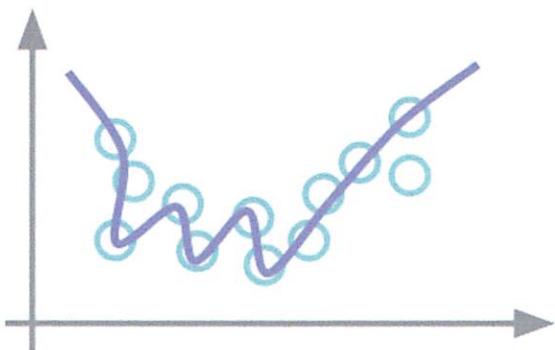
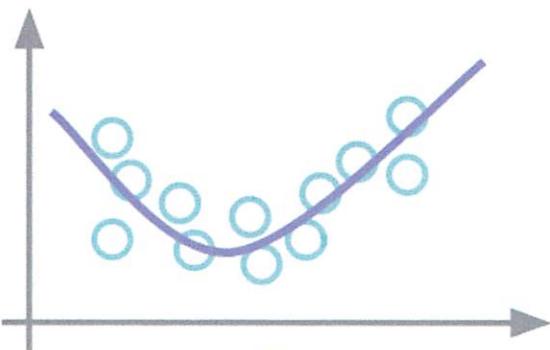
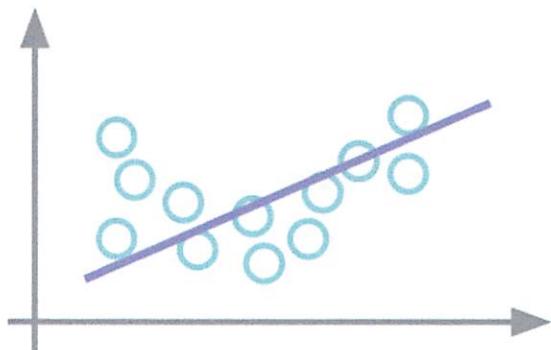


$$a_i \rightarrow \frac{\mu}{\sigma}$$
$$\hat{a}_i = \frac{a_i - \mu}{\sigma}$$
$$\hat{a}_i = \gamma \hat{a}_i + \beta$$

mem: μ_p, σ_p



Under Vs Over-fitting



$L_1 + L_2$ Reg.

∇

\min

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

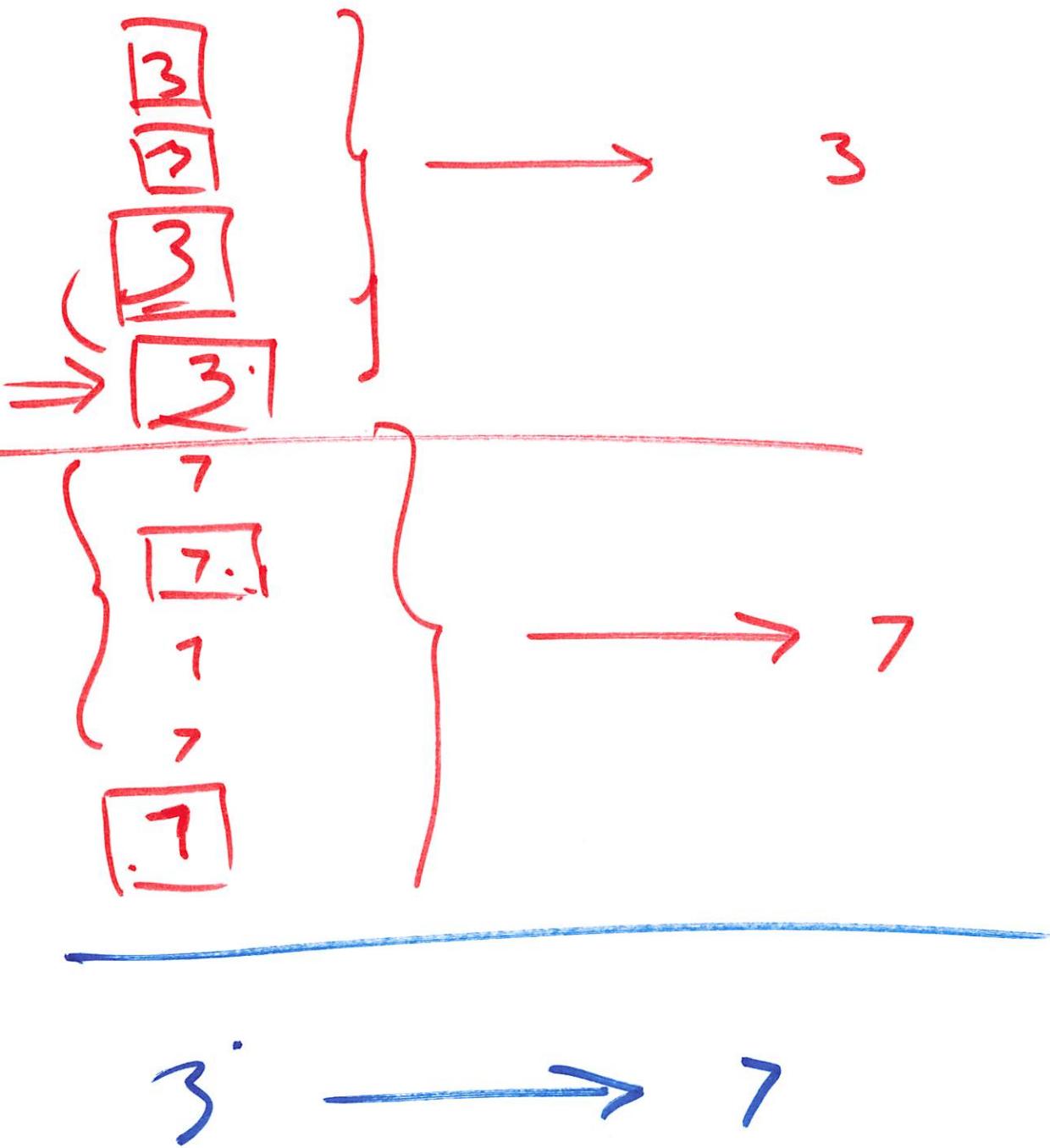
$$\hat{y} = \sum \beta_i x_i + b$$

penalty

$$\rightarrow \sum |w_i|$$

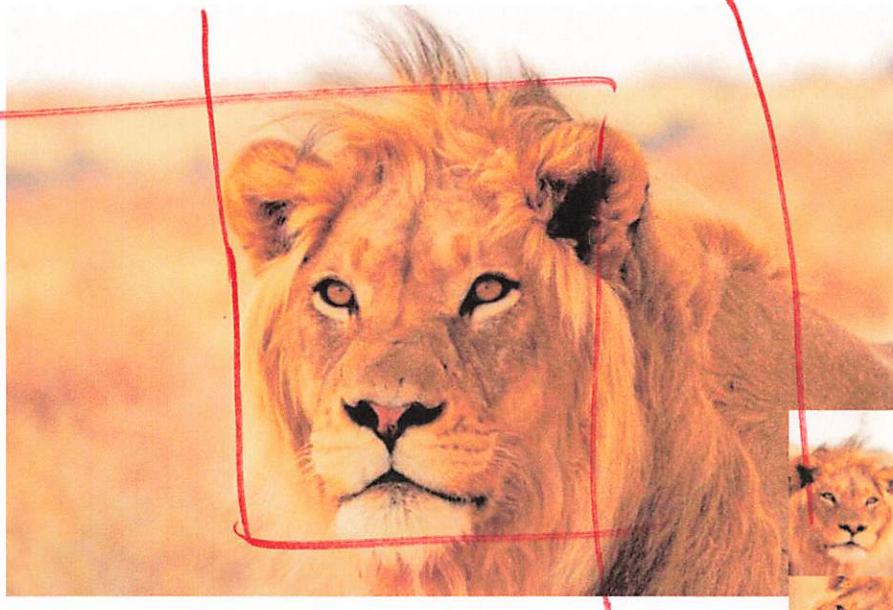
$$\rightarrow \sum (\beta_i)^2$$

Lasso =
L2 ridge



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

Data Augmentation



± 15°
Noise
Shift
mirror
Stretch
Color
Crop
GAN

-source: [towardsdatascience.com](https://towardsdatascience.com/machinex-image-data-augmentation-using-keras-1000f3a2a23), MachineX: Image Data Augmentation Using Keras

Proprietary content. © Great Learning. All Rights Reserved. Unauthorised use or distribution is illegal.

This file is meant for personal use by shrikadeshwar@gmail.com only.

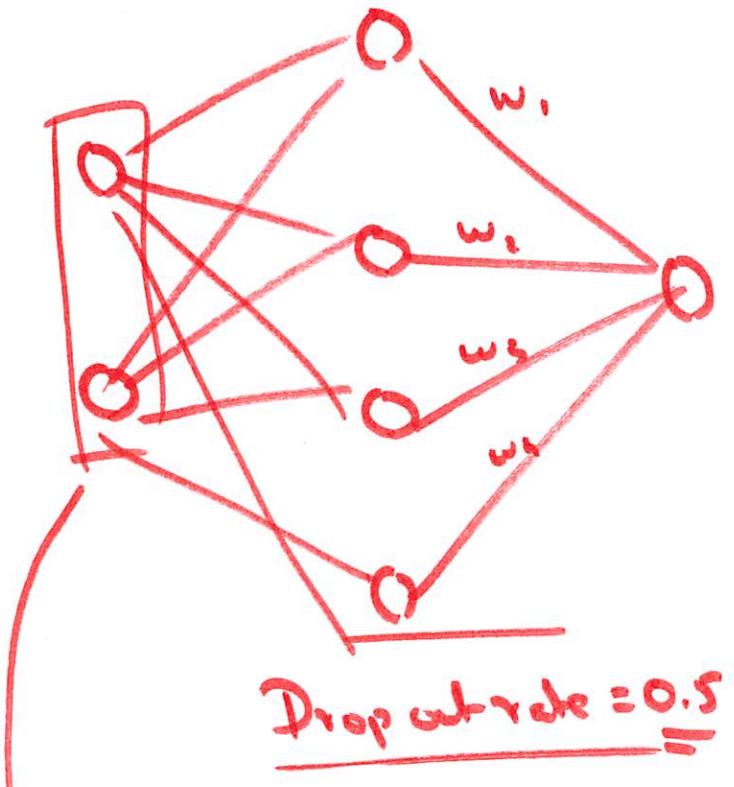
Sharing or publishing the contents in part or full is liable for legal action.

Co Adaptation



$$-0.1a_1 + 0.1a_2$$

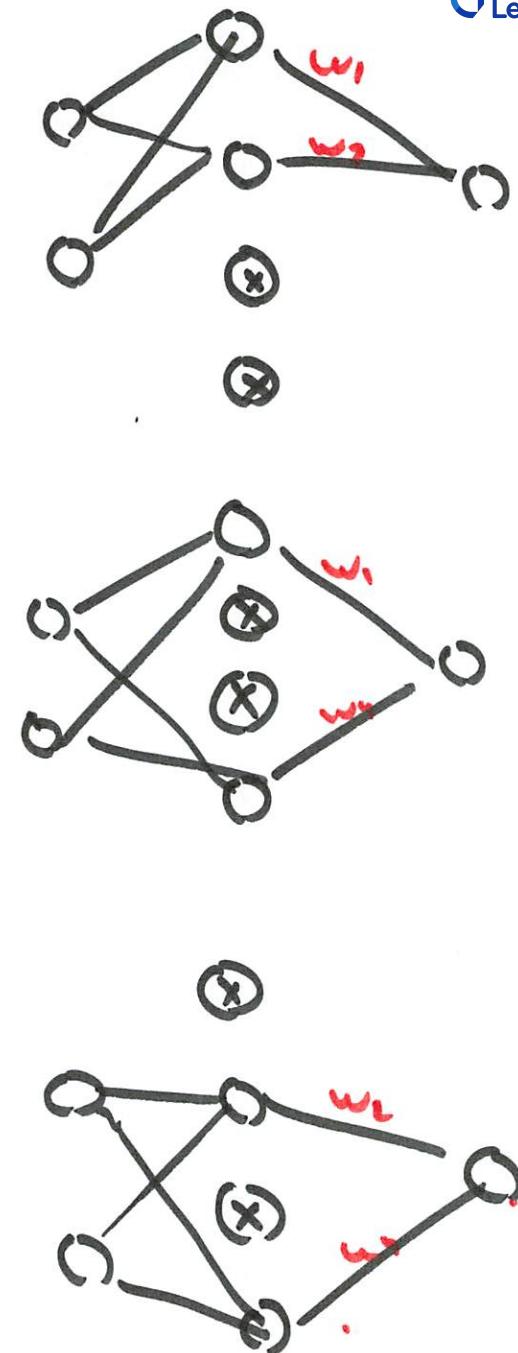
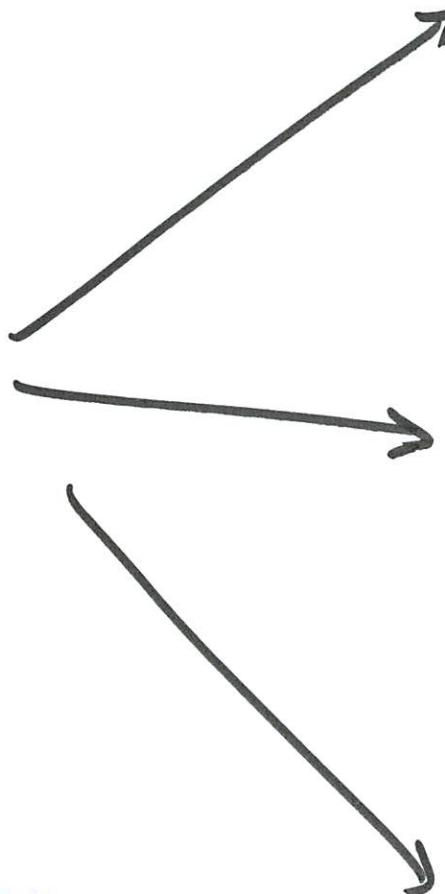
$$\underline{-0.3a_1 + 0.3a_2}$$



$$20\% - 50\% =$$

Drop out $\underline{\underline{20\%}}$

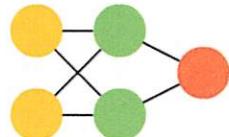
$$3' = \underline{\underline{w^0 \times 0.5}}$$



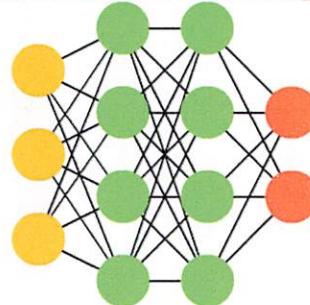
Types of NN

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

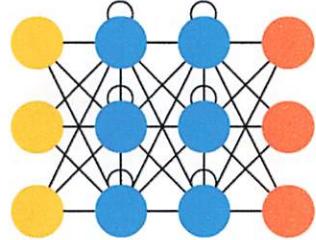
Feed Forward (FF)



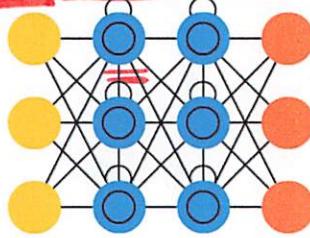
Deep Feed Forward (DFF)



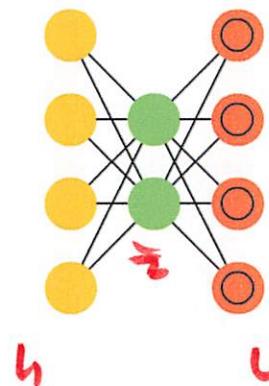
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Compressed
dim
(x)



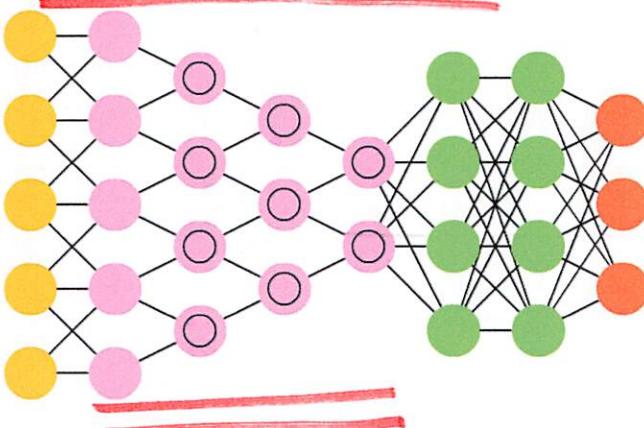
$\rightarrow |11011|$

This file is meant for personal use by envydesigner@gmail.com only.

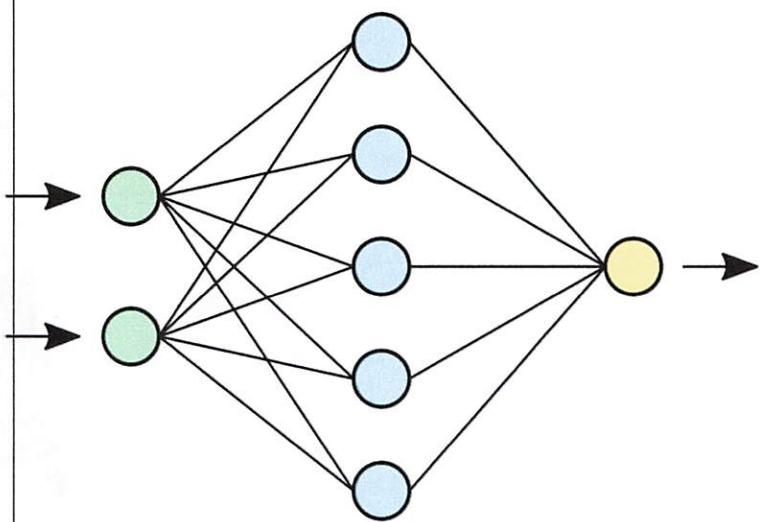
Sharing or publishing the contents in part or full is liable for legal action.

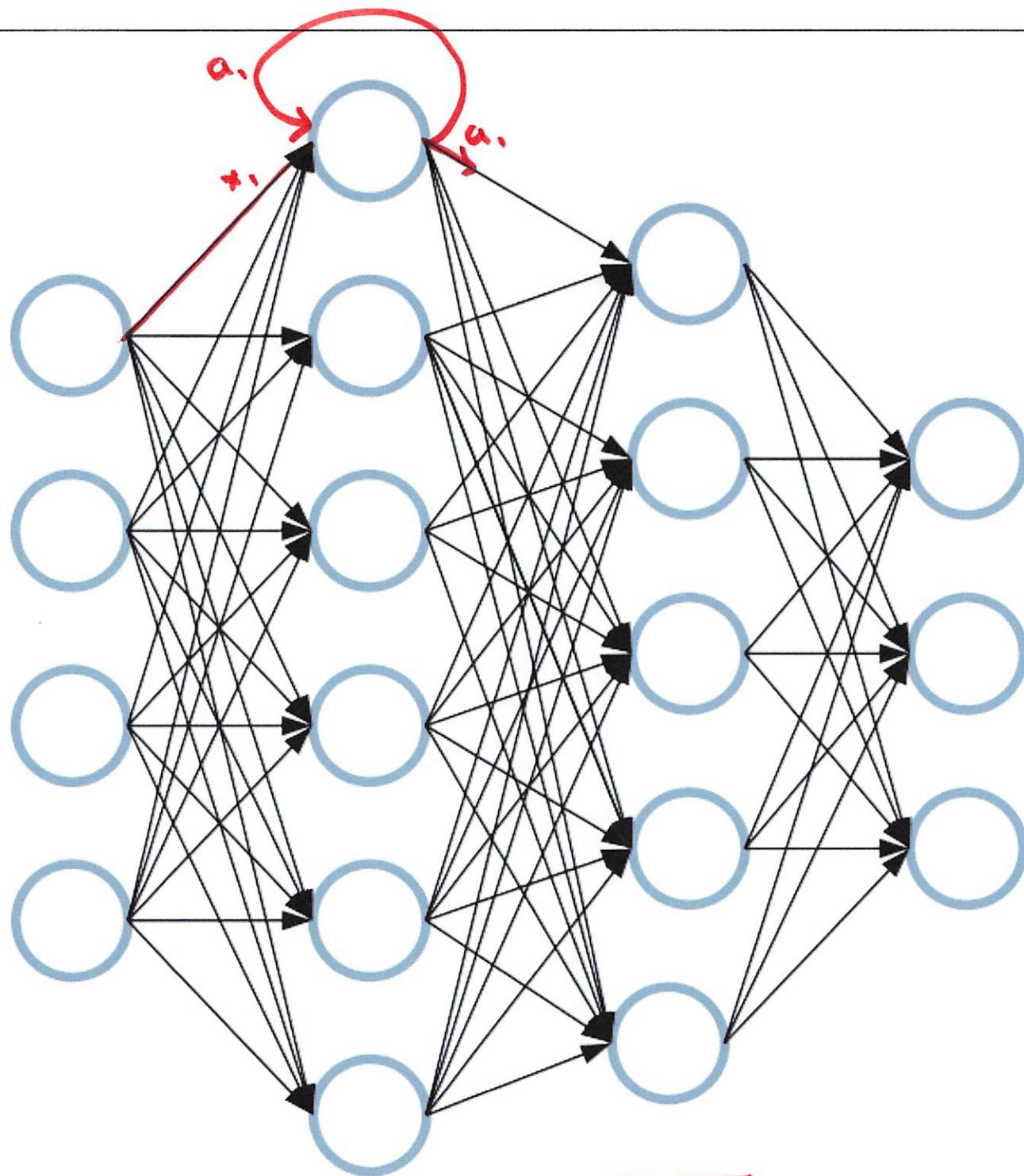
17 05
1
17
05 17

Deep Convolutional Network (DCN)



Feed Forward Net.





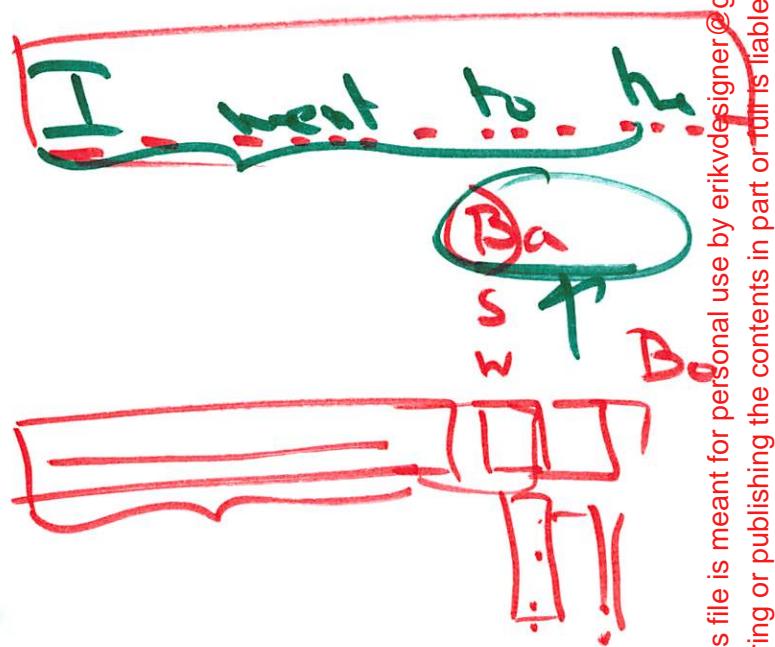
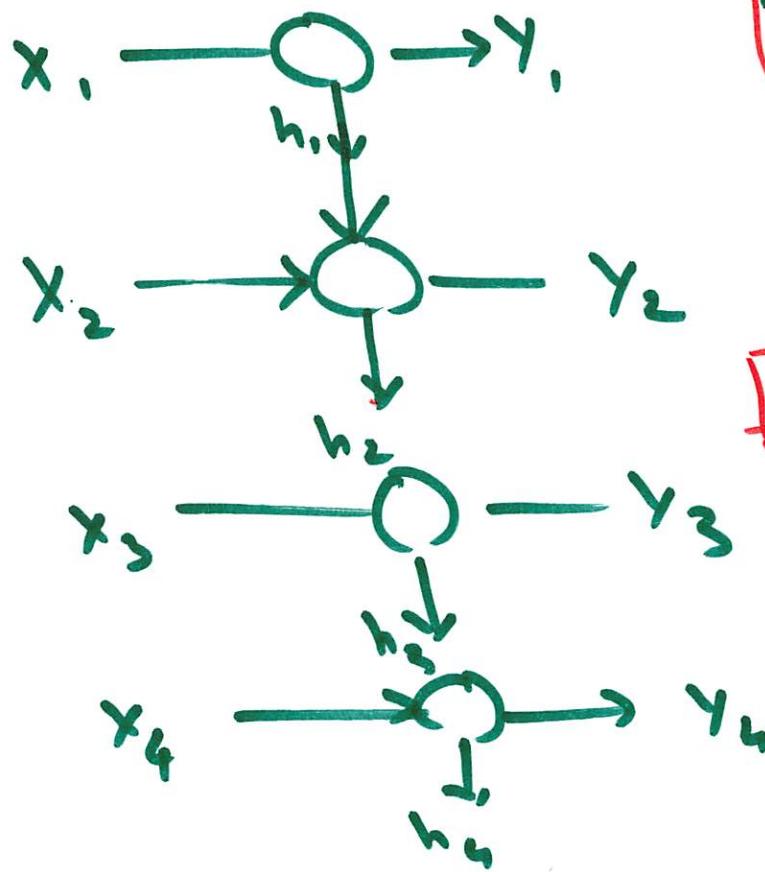
This file is meant for personal use by erikvdesigner@gmail.com only.
Sharing or publishing the contents in part or full is liable for legal action.



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

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

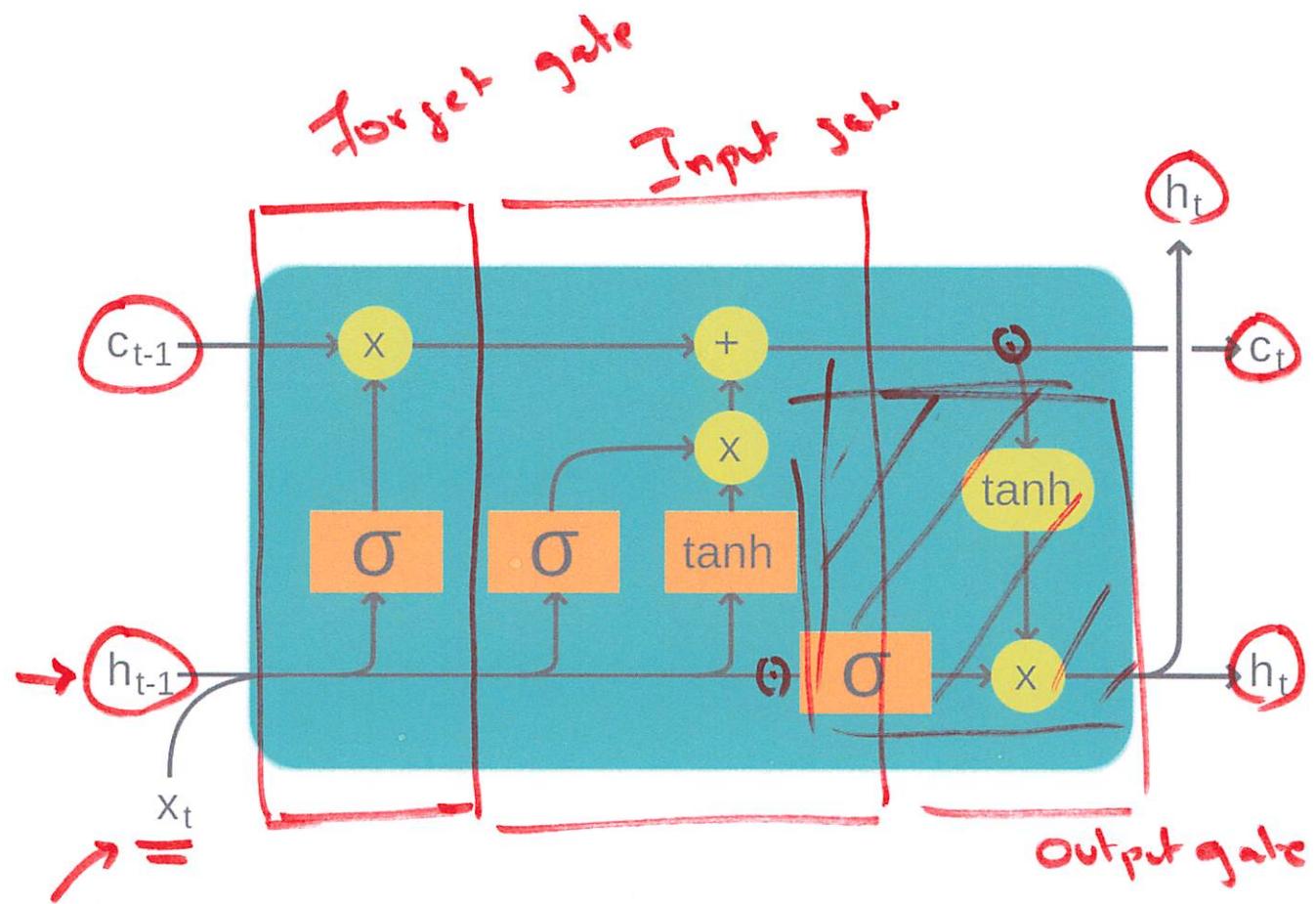
$$h = f(\omega x + \omega_0 b + b)$$



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

