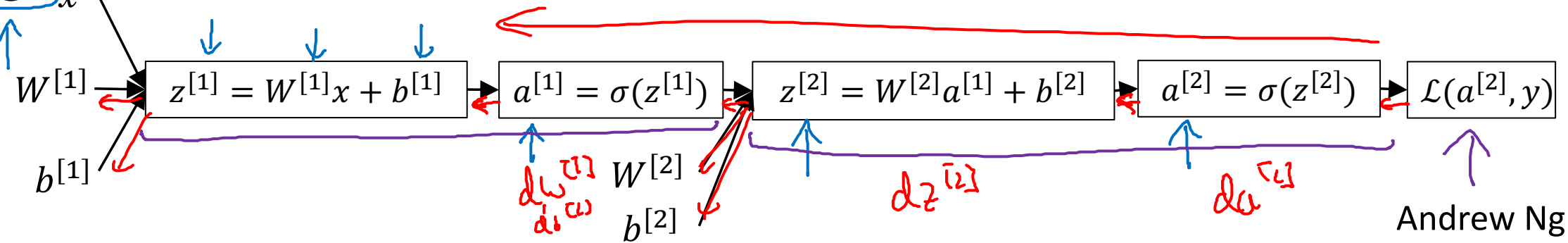
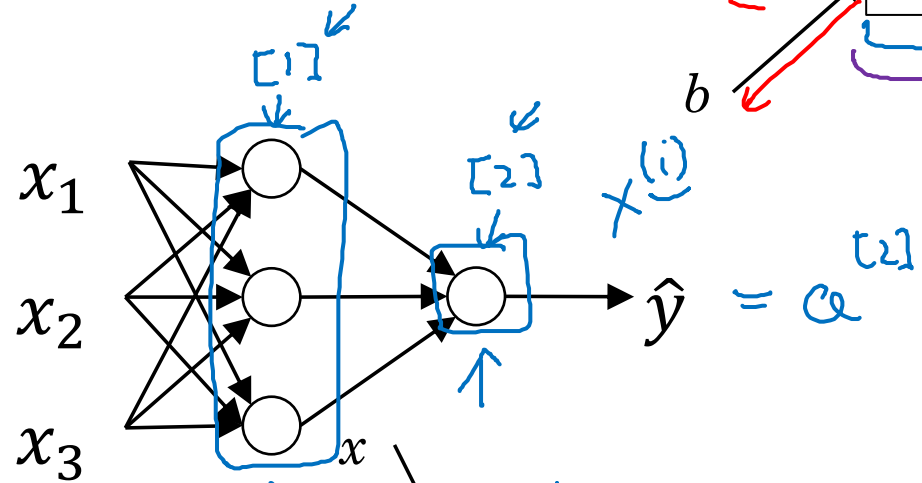
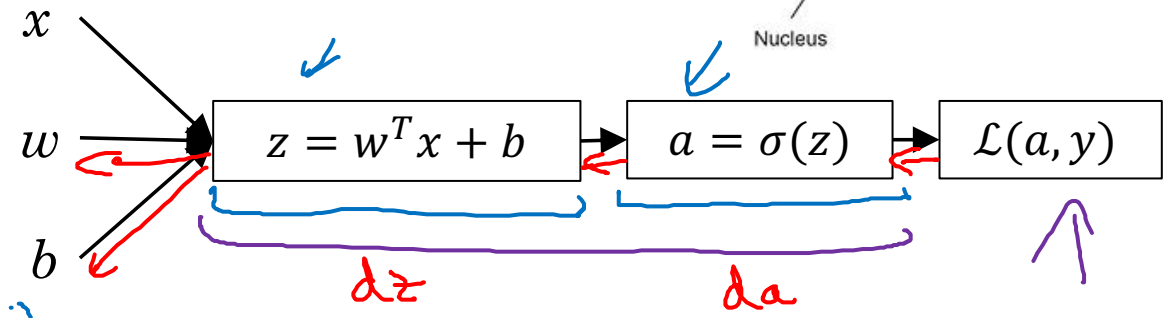
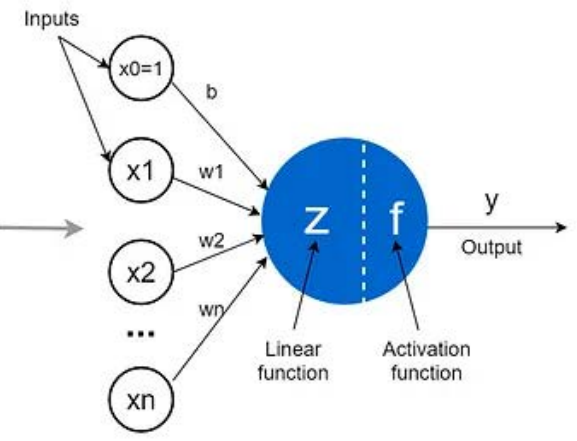
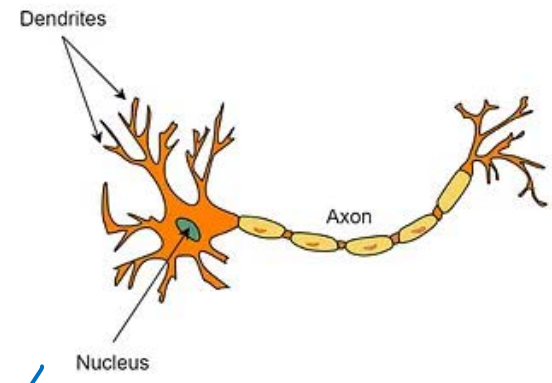
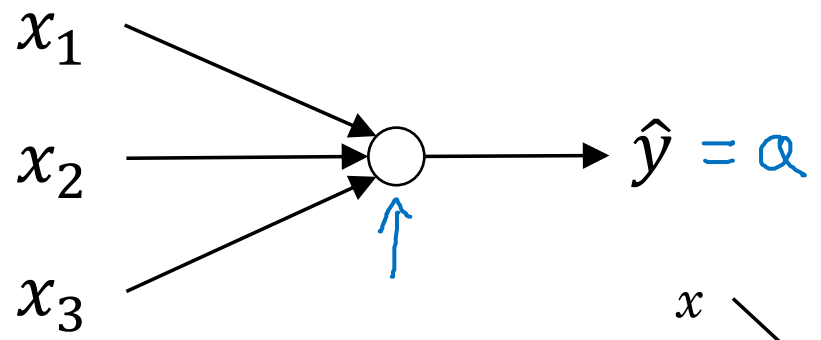
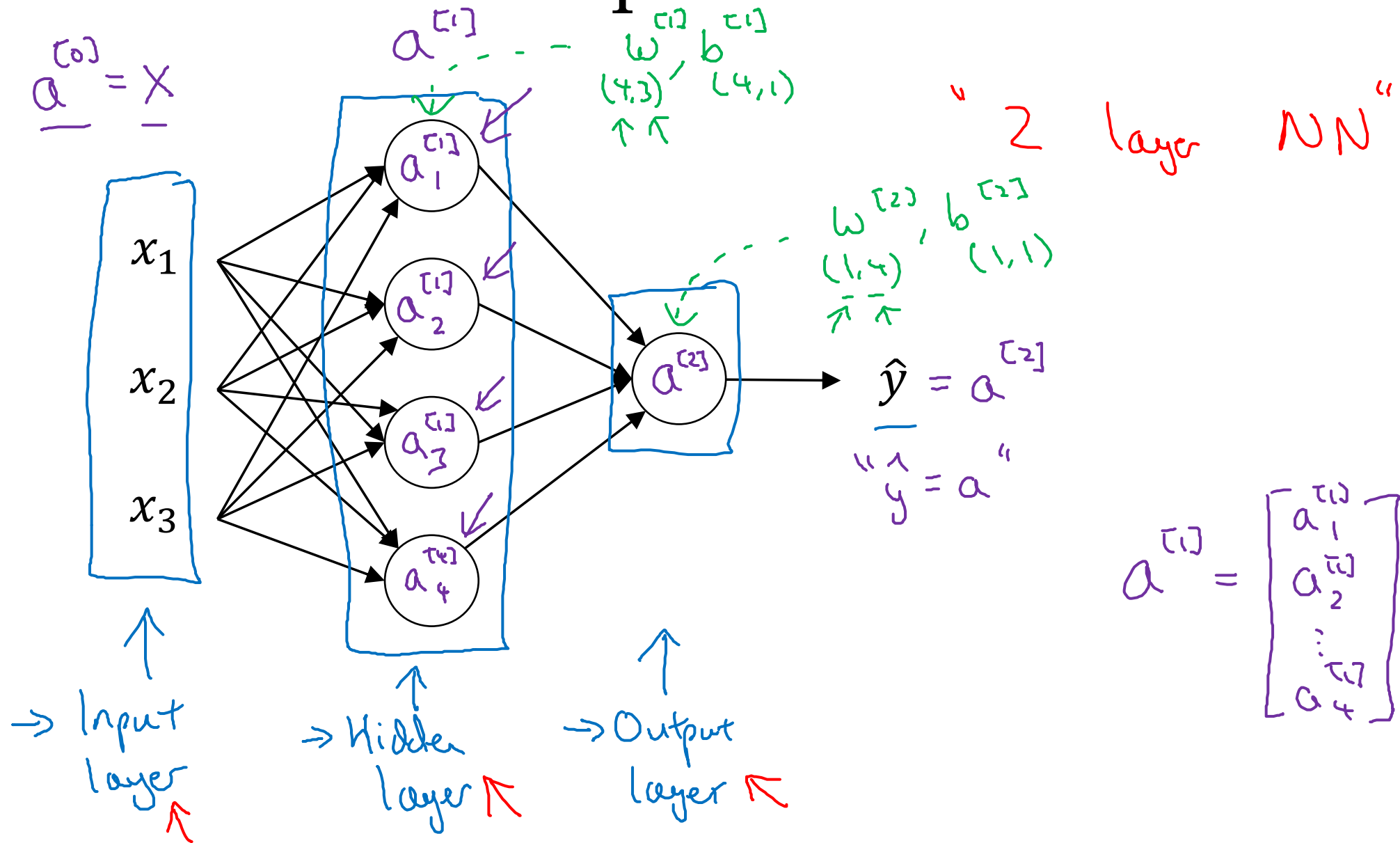


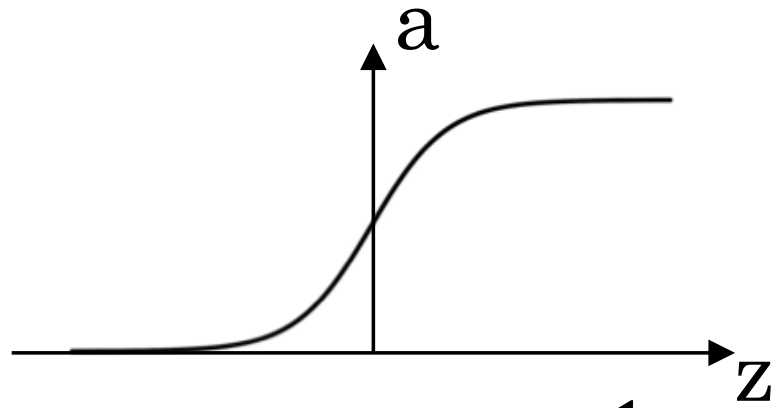
What is a Neural Network?



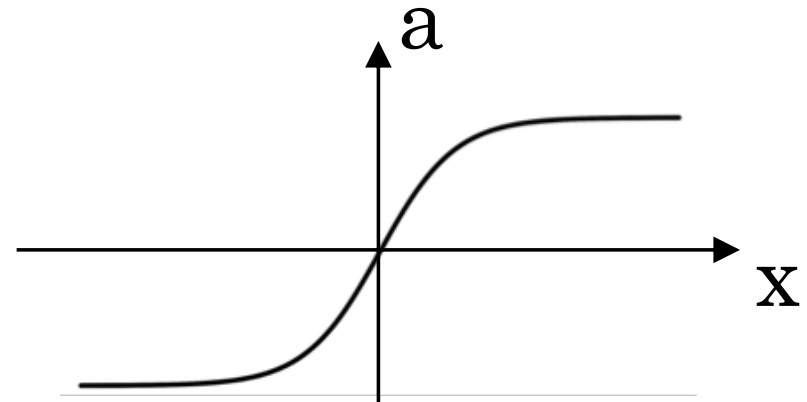
Neural Network Representation



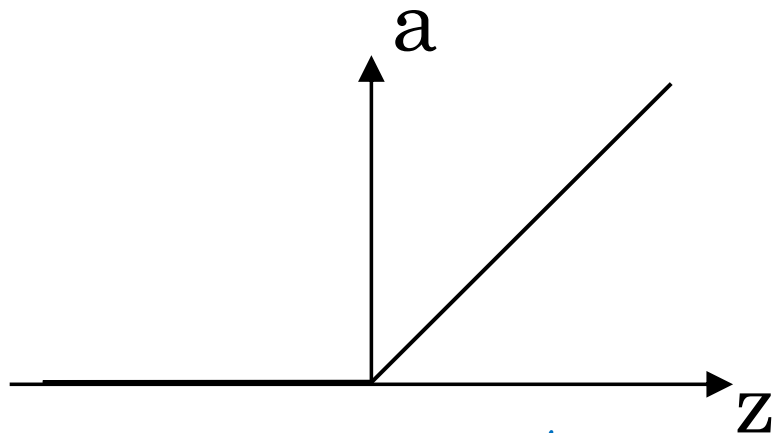
Pros and cons of activation functions



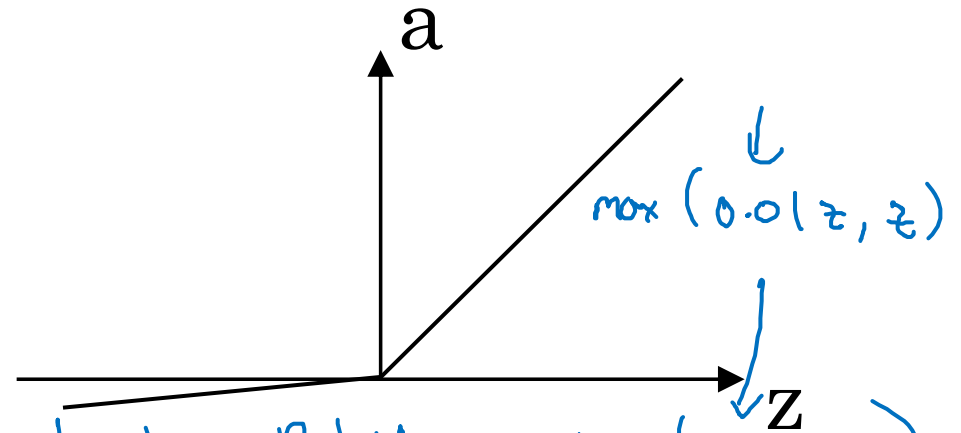
sigmoid: $a = \frac{1}{1 + e^{-z}}$



tanh: $a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$



ReLU $a = \max(0, z)$



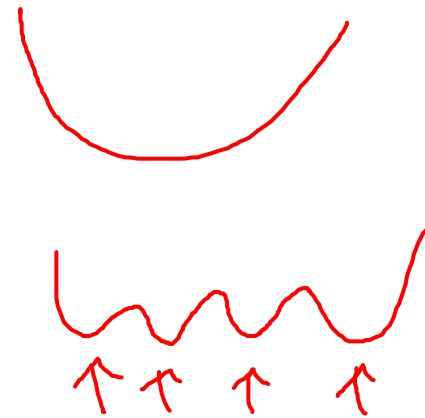
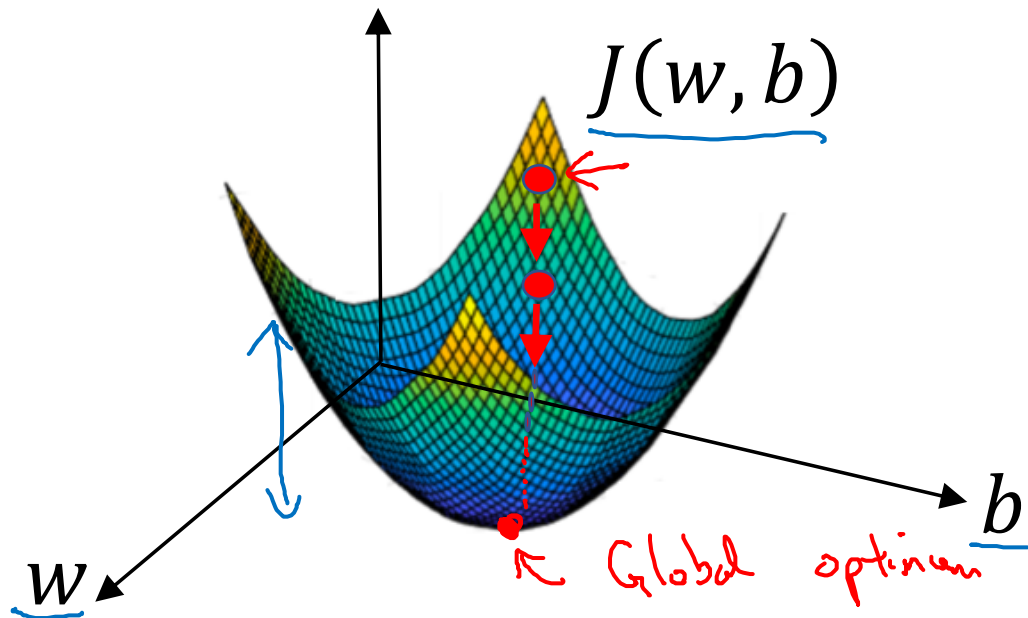
Leaky ReLU $a = \max(0.01z, z)$

Gradient Descent

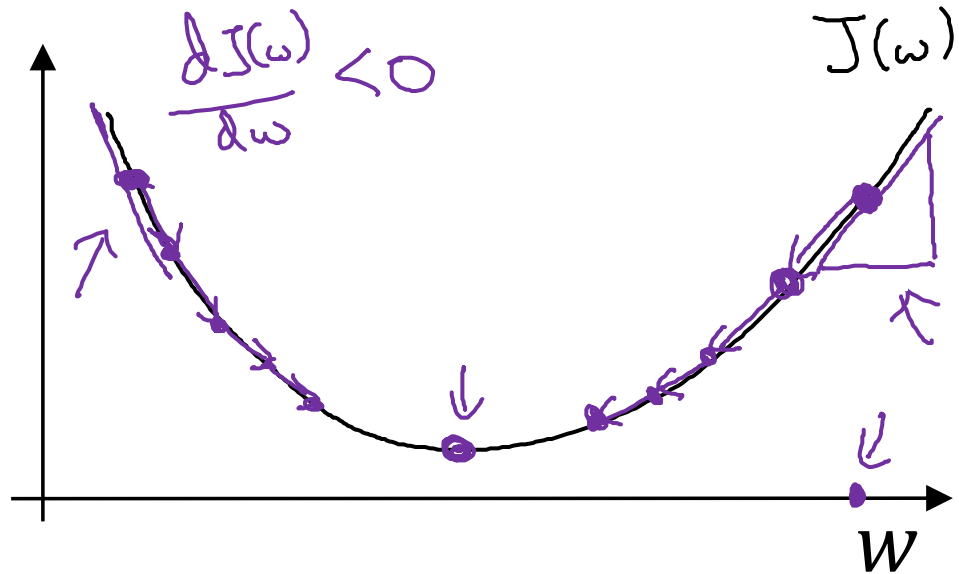
Recap: $\hat{y} = \sigma(w^T x + b)$, $\sigma(z) = \frac{1}{1+e^{-z}}$ ←

$$\underline{J(w, b)} = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\underline{\hat{y}^{(i)}} , \underline{y^{(i)}}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Want to find w, b that minimize $J(w, b)$



Gradient Descent



Repeat {

$$w := w - \alpha \frac{dJ(w)}{dw}$$

learning rate

↑ ↑ "dw"

$$w := w - \alpha \underline{dw}$$

$\frac{dJ(w)}{dw} = ?$

$J(w, b)$

$$w := w - \alpha \frac{dJ(w, b)}{dw}$$

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

$\frac{\partial J(w, b)}{\partial w}$ $\frac{\partial J(w, b)}{\partial b}$

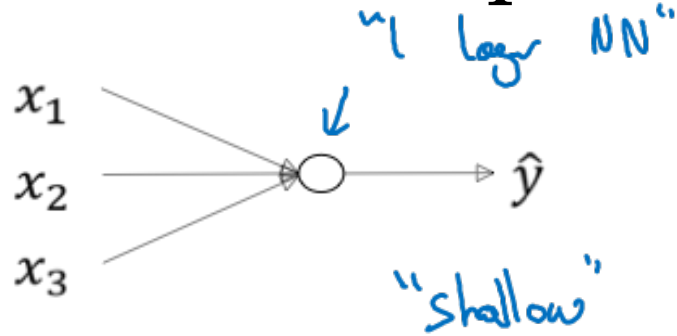
$\frac{\partial}{\partial w}$ $\frac{\partial}{\partial b}$

← "partial derivative" J

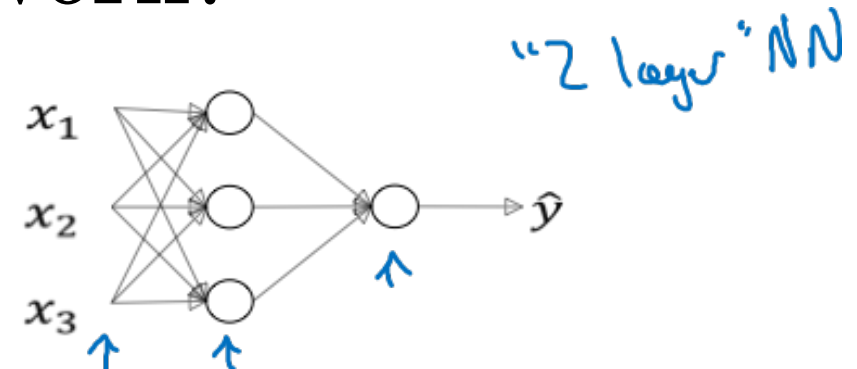
→ dw

→ db

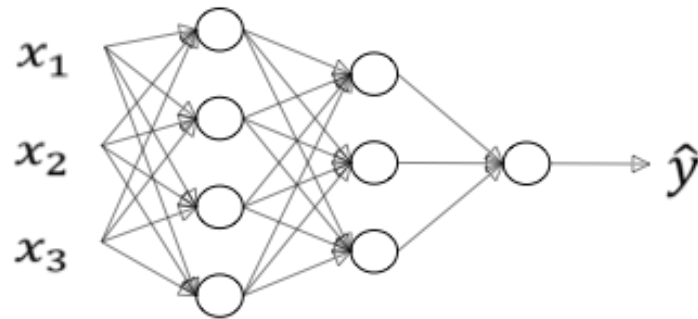
What is a deep neural network?



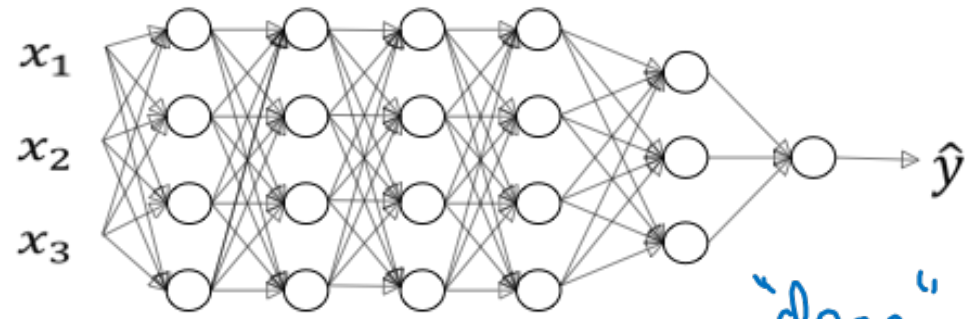
logistic regression



1 hidden layer



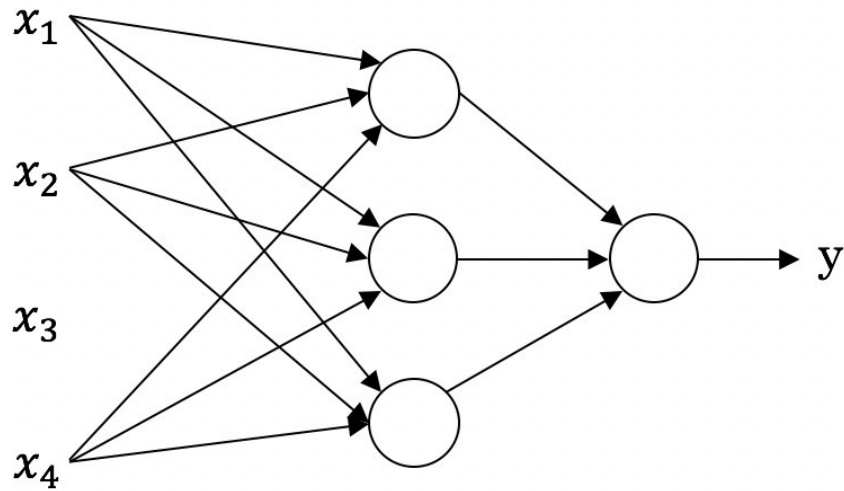
2 hidden layers



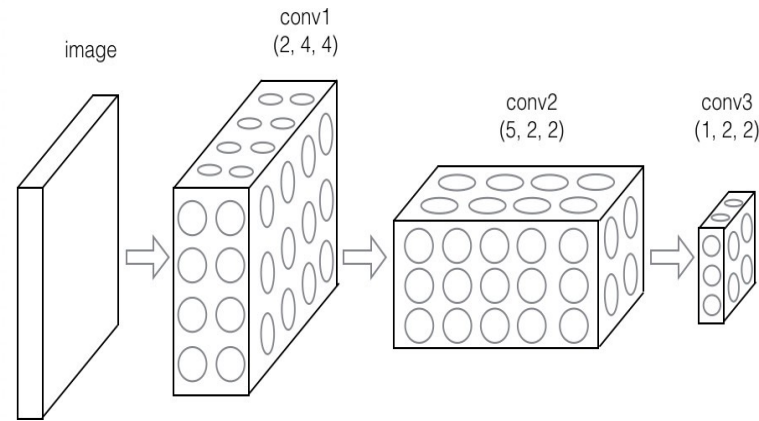
5 hidden layers

"deep"

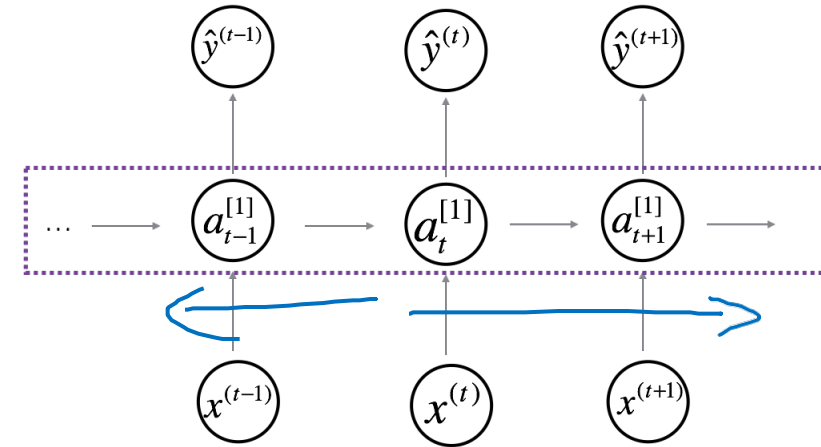
Neural Network examples



Standard NN



Convolutional NN



Recurrent NN

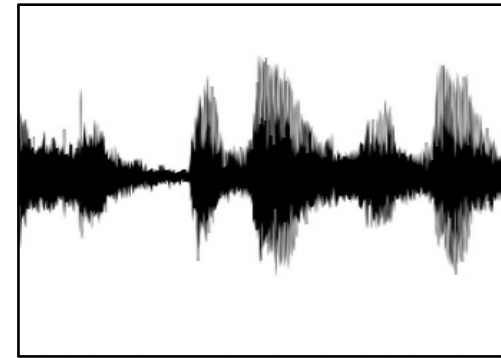
Supervised Learning

Structured Data

Size	#bedrooms	...	Price (1000\$)
2104	3		400
1600	3		330
2400	3		369
...
3000	4		540

User Age	Ad Id	...	Click
41	93242		1
80	93287		0
18	87312		1
...
27	71244		1

Unstructured Data



Audio

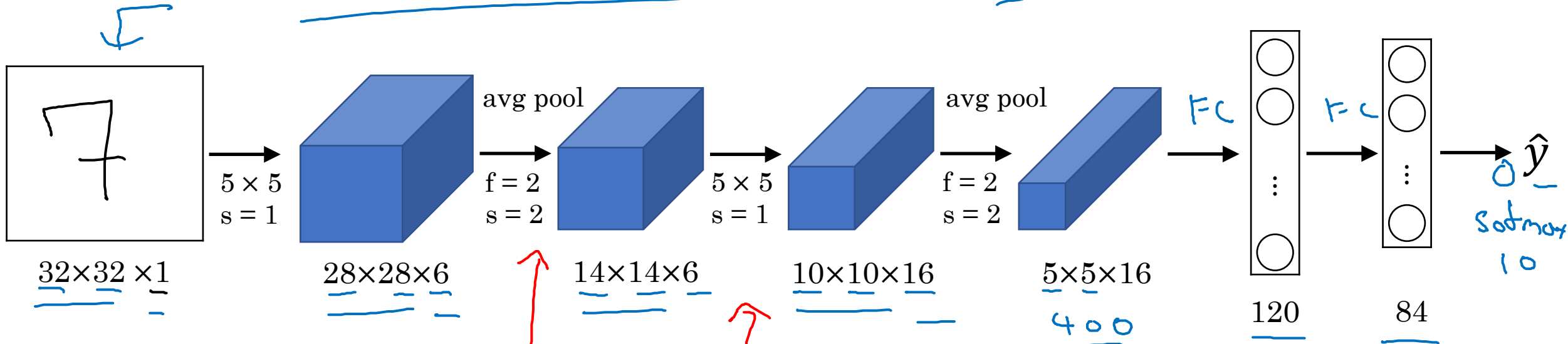


Image

Four scores and seven
years ago...

Text

LeNet - 5



60K parameters.

$n_H, n_w \downarrow$ $n_c \uparrow$

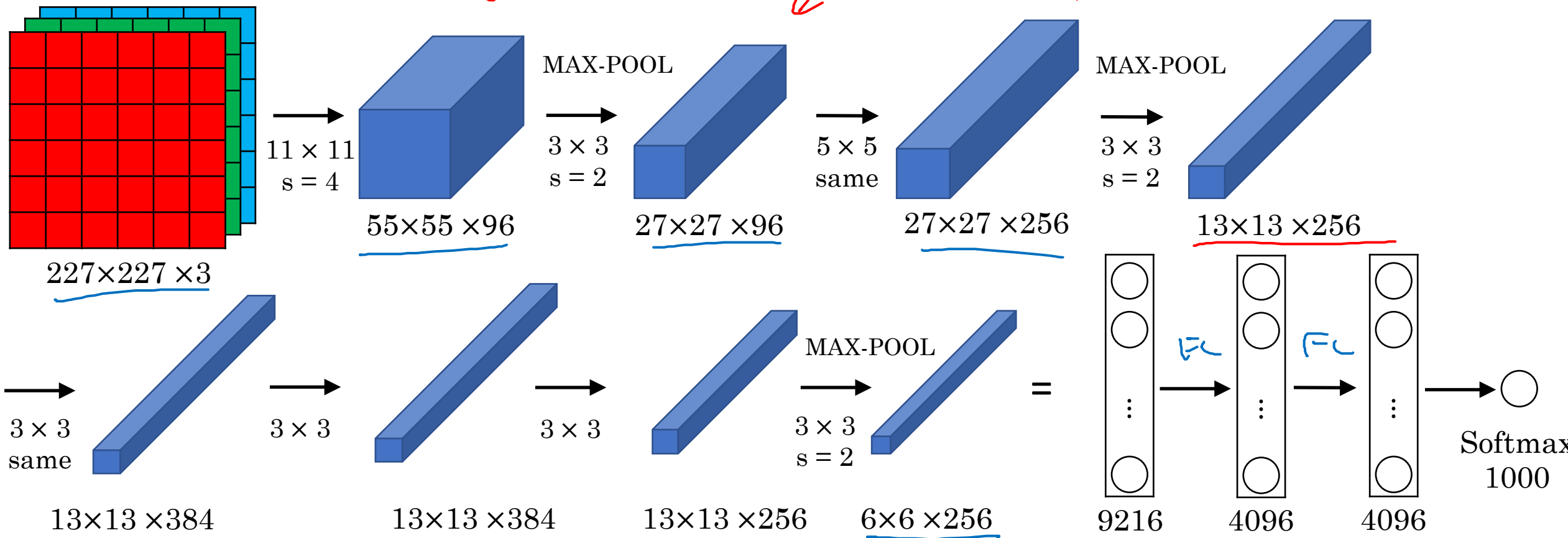
conv pool conv pool fc fc output

non-linearity after pooling $n_H \times n_w \times n_c$ $f \times f \times n_c$

Advanced: sigmoid/tanh ReLU

II, III.

AlexNet

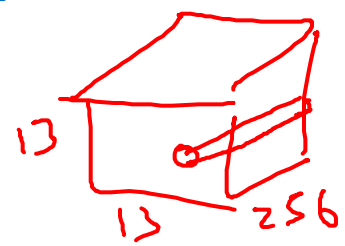


- Similar to LeNet, but much bigger. 9216 160M parameters

- ReLU

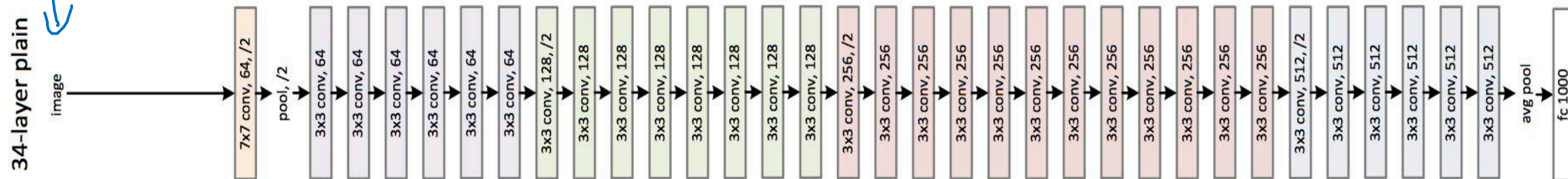
- Multiple GPUs.

- Local Response Normalization (LRN)

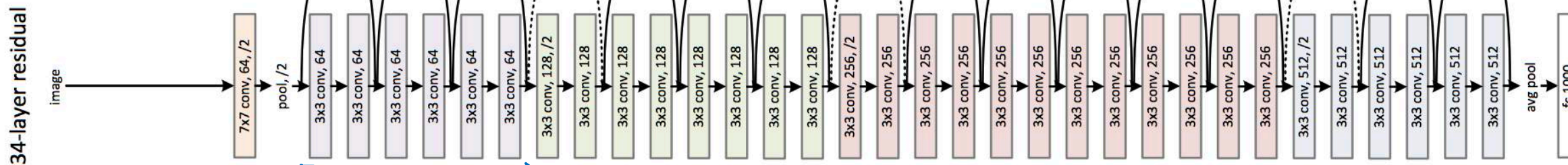


ResNet

Plain



ResNet



$z^{[l+2]} + a^{[l]}$

3x3 same

Pool

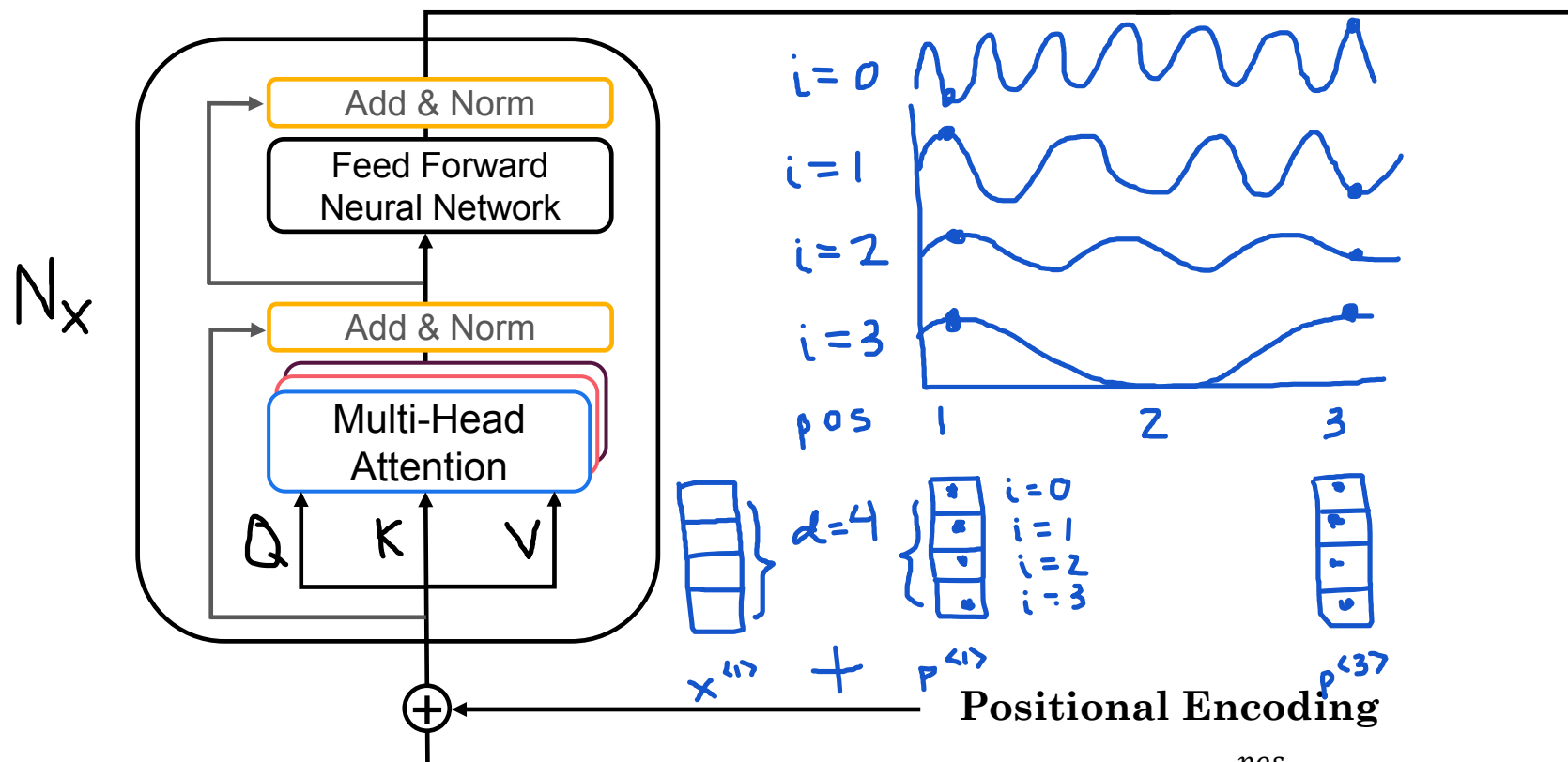
Pool

W_s

Transformer Details

<SOS> Jane visits Africa in September <EOS>

Encoder

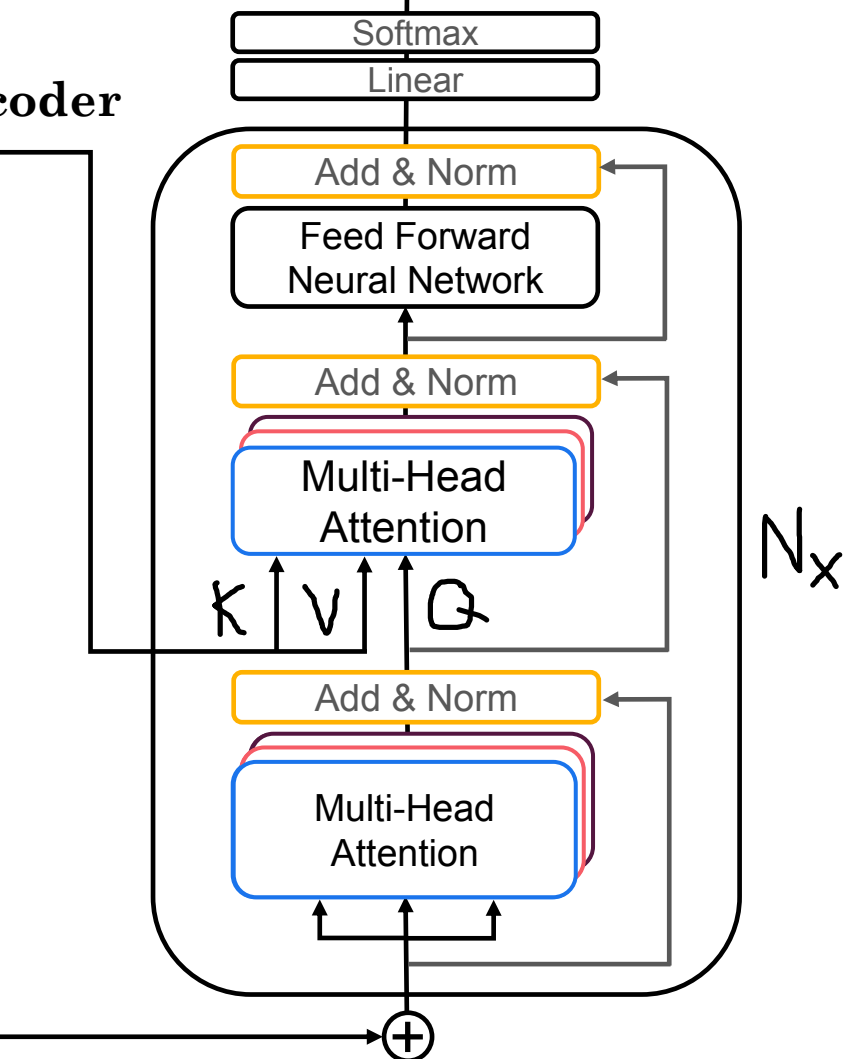


<SOS> $x^{<1>}$ $x^{<2>}$... $x^{<T_x-1>}$ $x^{<T_x>}$ <EOS>
Jane visite l'Afrique en septembre

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{1000^{\frac{2i}{d}}}\right)$$

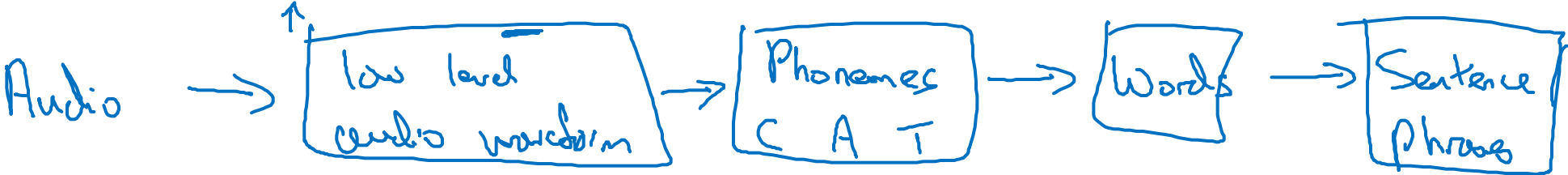
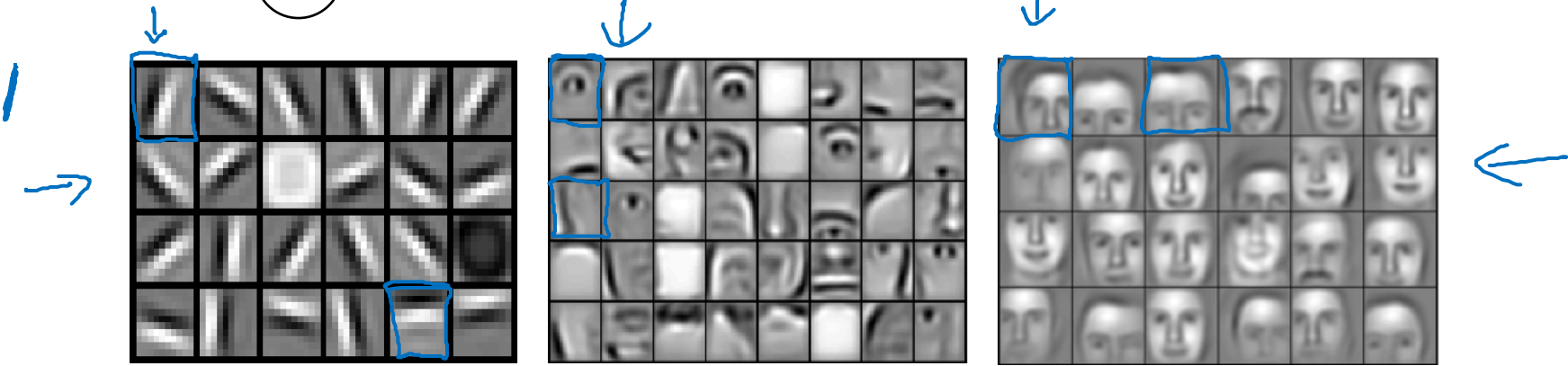
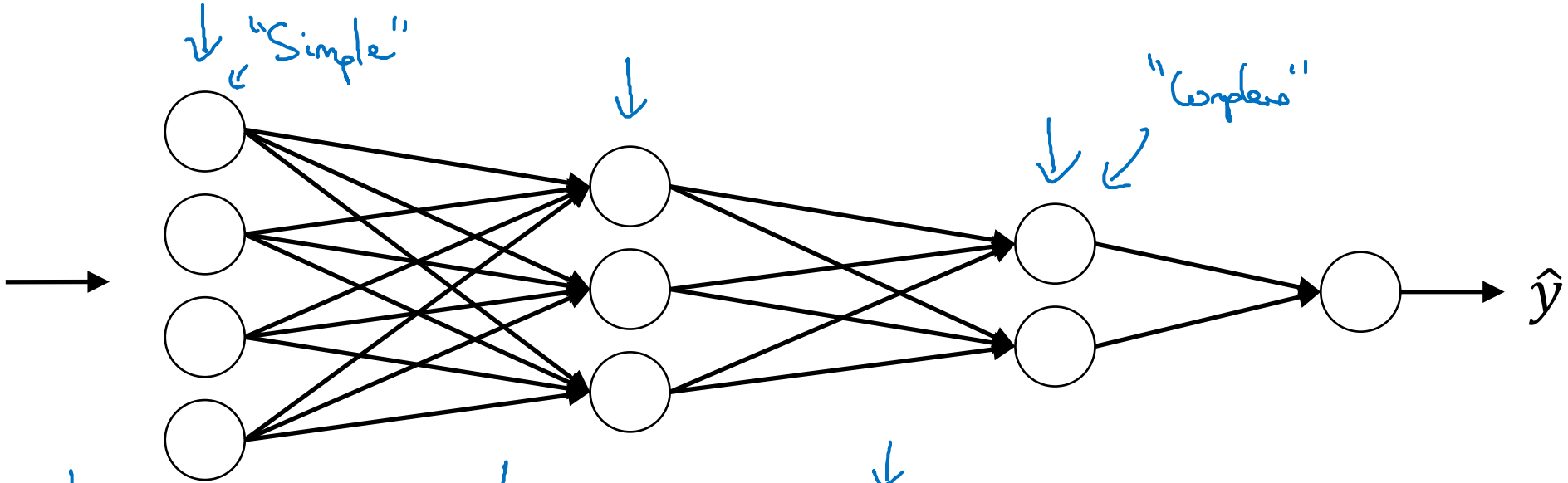
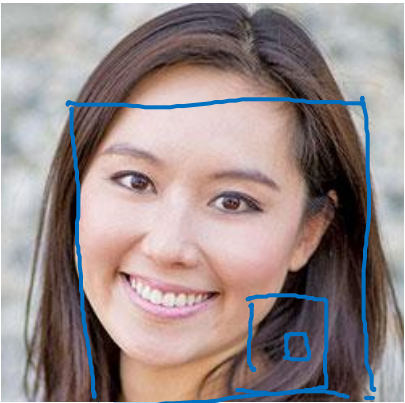
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{1000^{\frac{2i}{d}}}\right)$$

Decoder

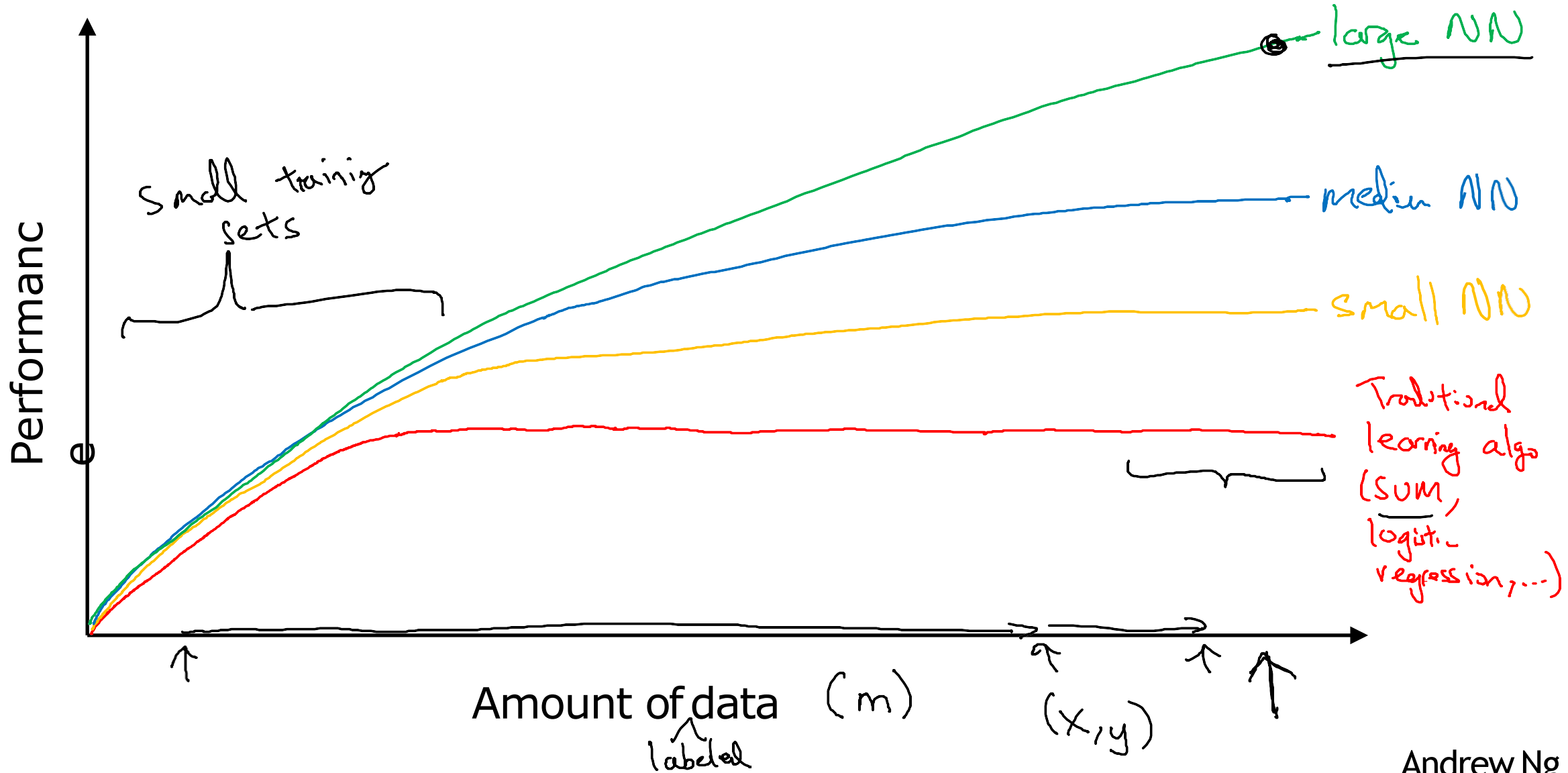


<SOS> $y^{<1>}$ $y^{<2>}$... $y^{<T_y-1>}$ $y^{<T_y>}$
<SOS> Jane visits Africa in September

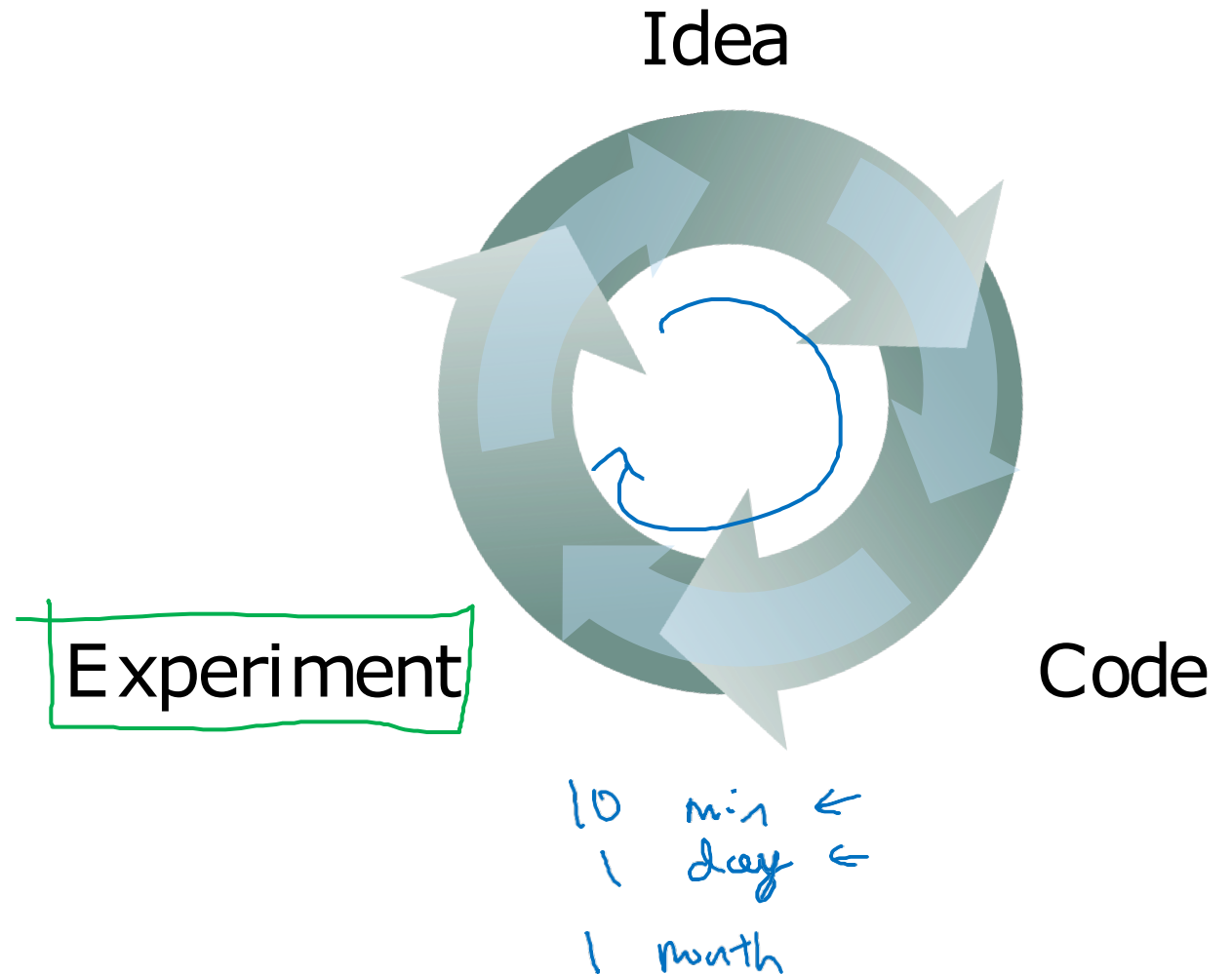
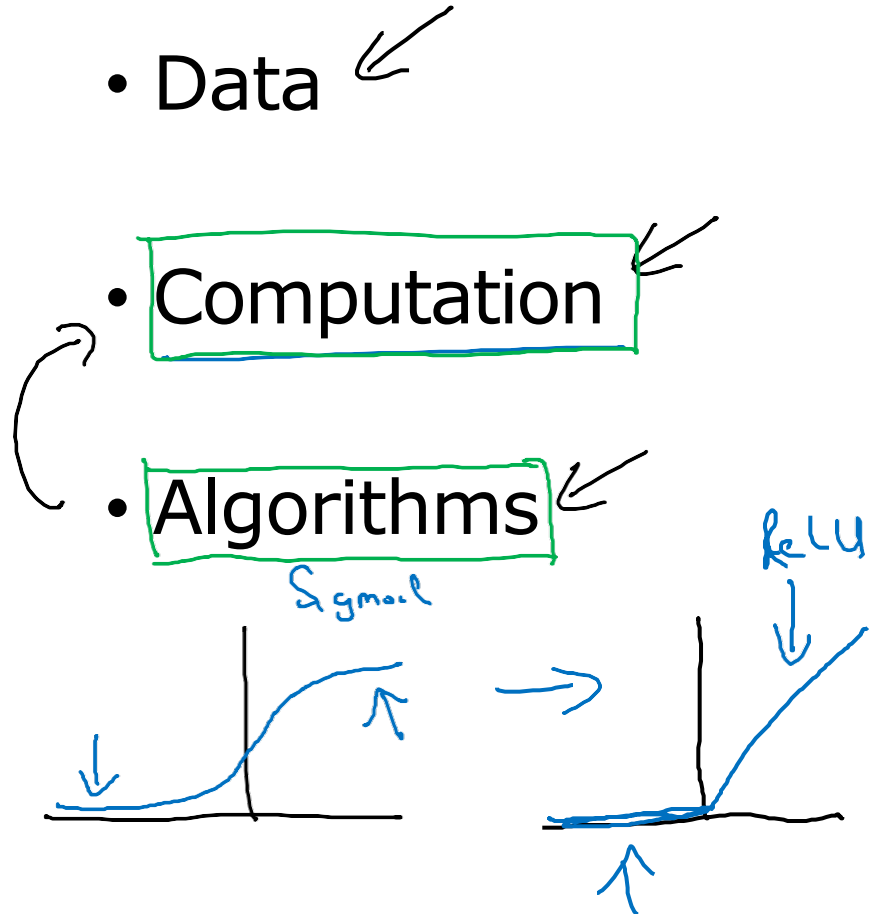
Intuition about deep representation



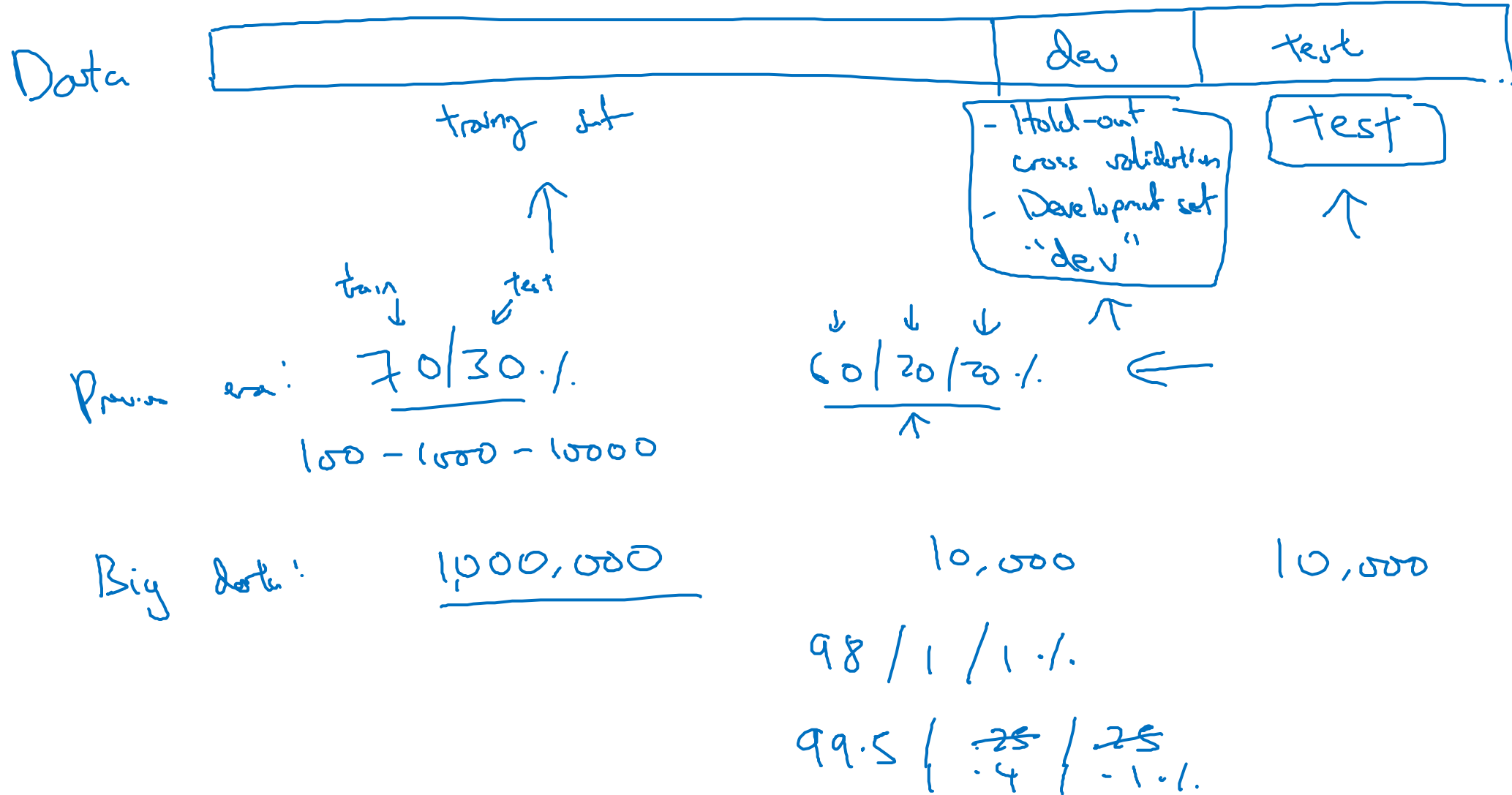
Scale drives deep learning progress



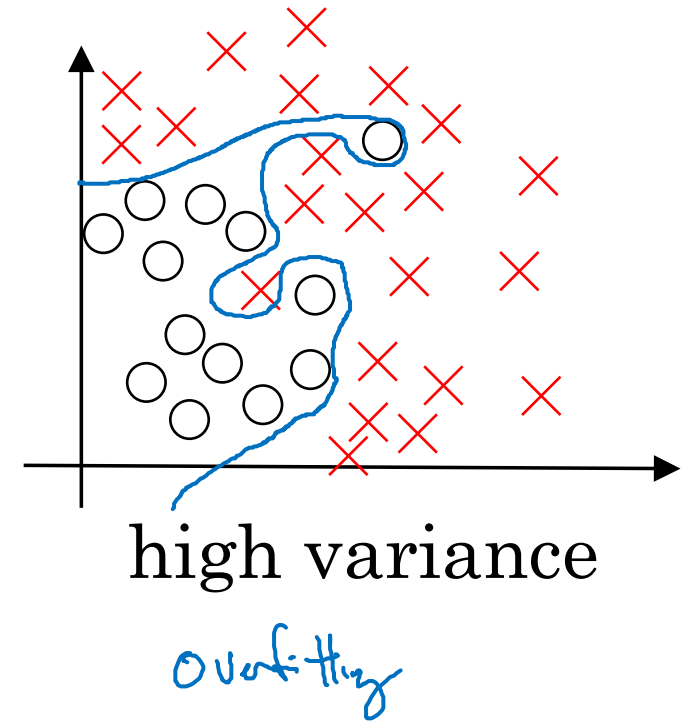
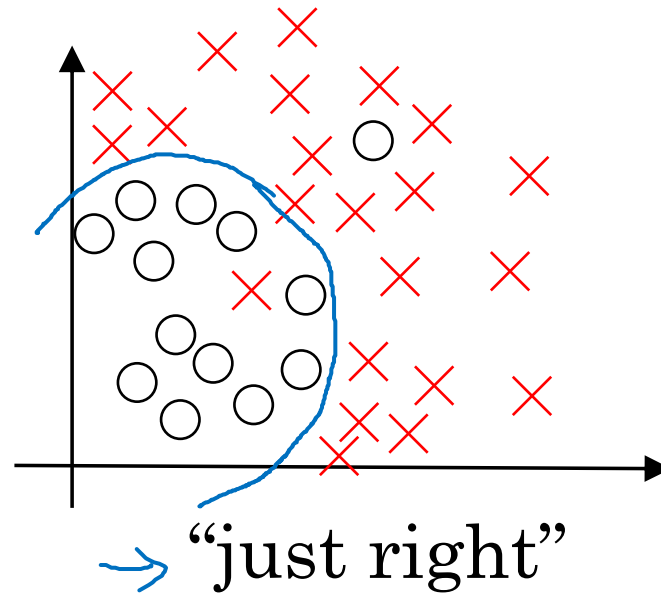
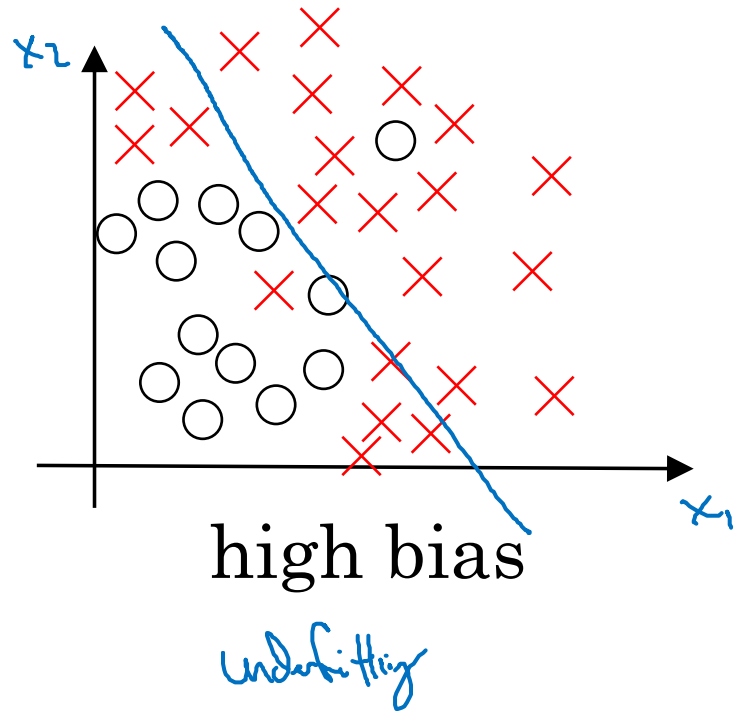
Scale drives deep learning progress



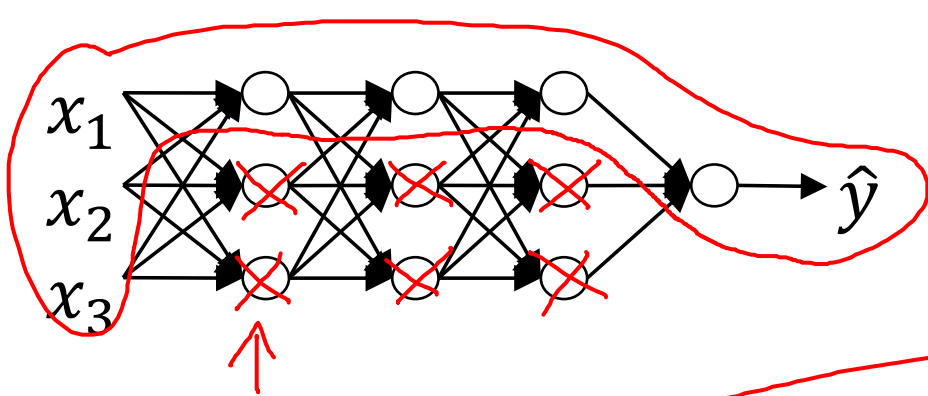
Train/dev/test sets



Bias and Variance

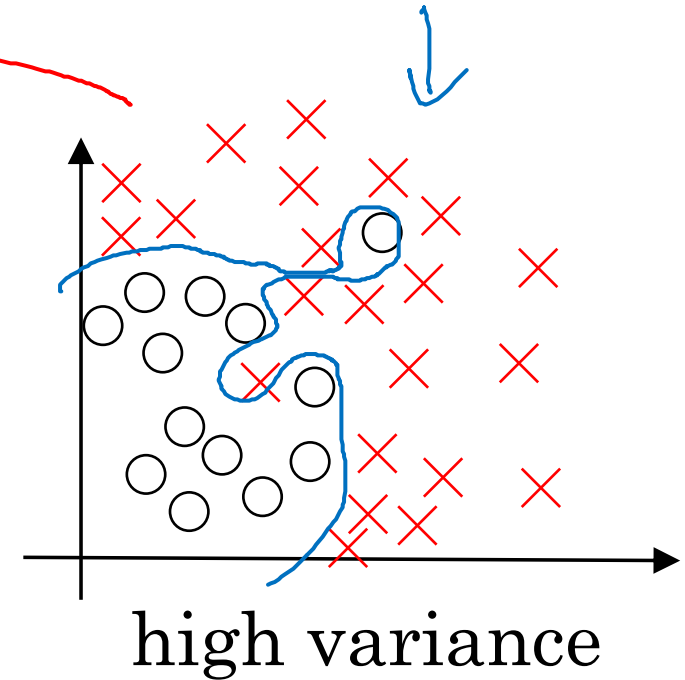
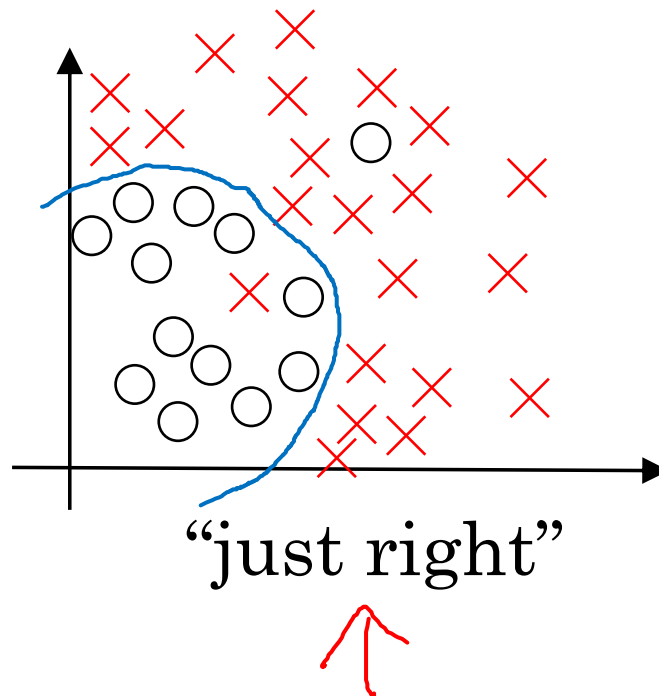
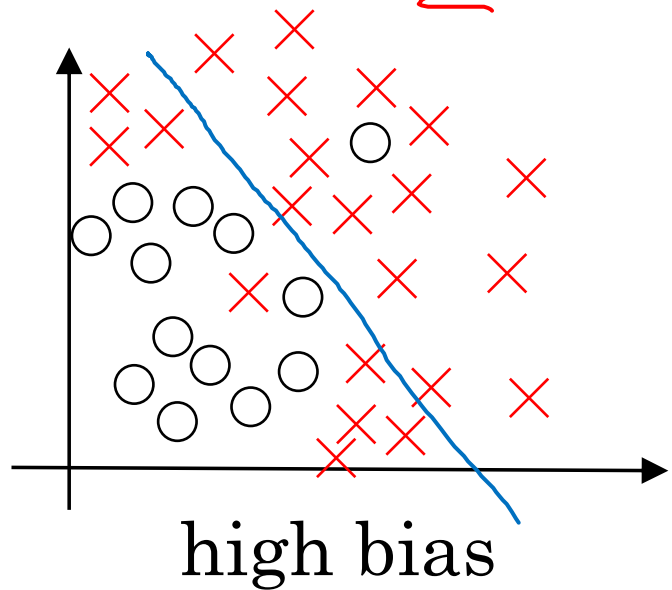


How does regularization prevent overfitting?

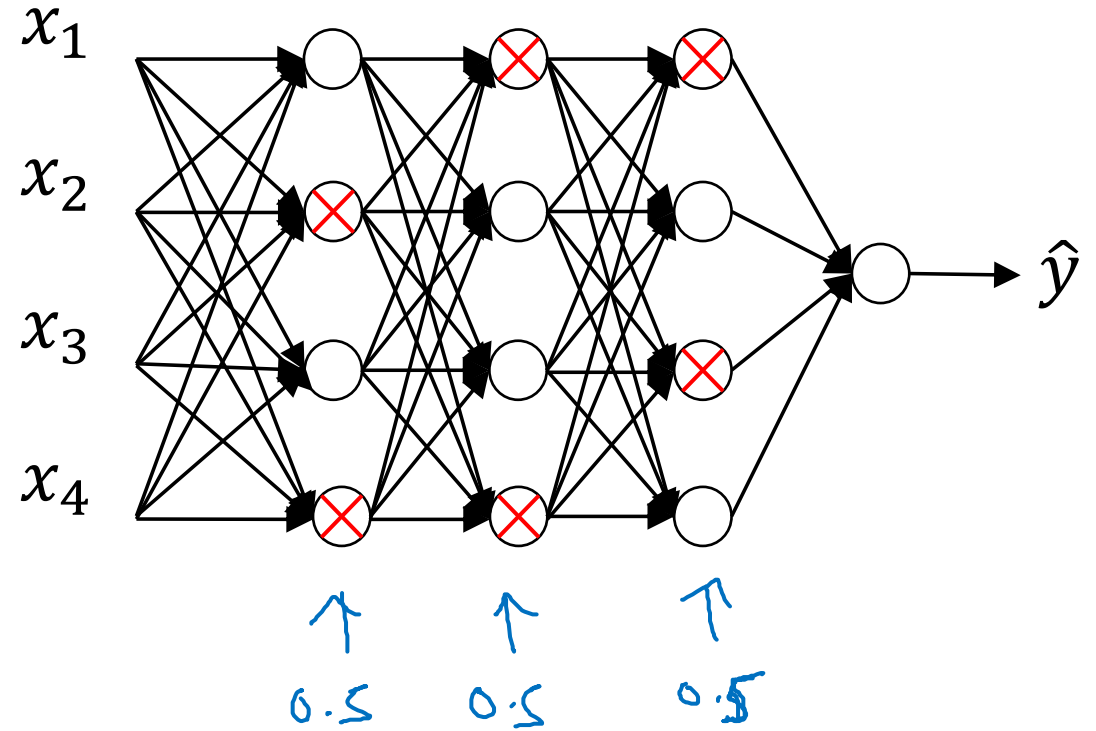
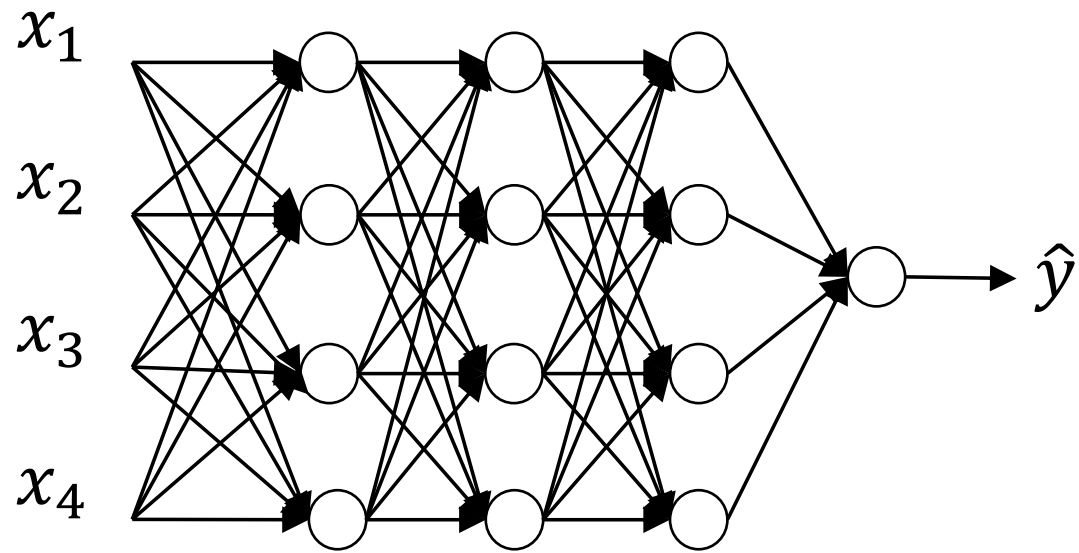


$$J(w^{(l)}, b^{(l)}) = \frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

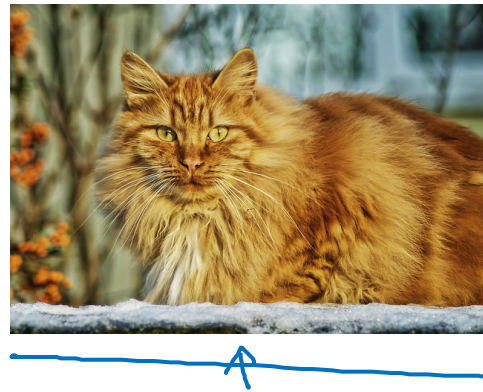
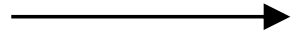
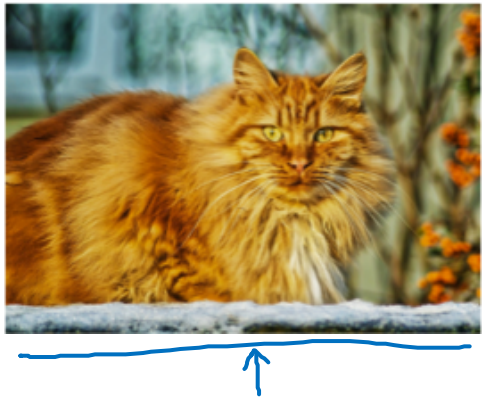
$$w^{(l)} \approx 0$$



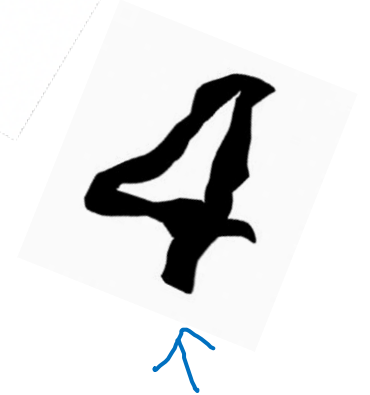
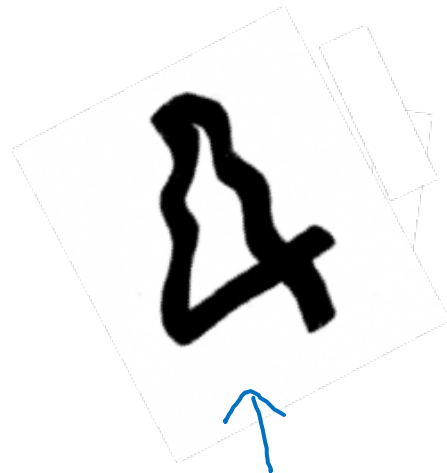
Dropout regularization



Data augmentation



4



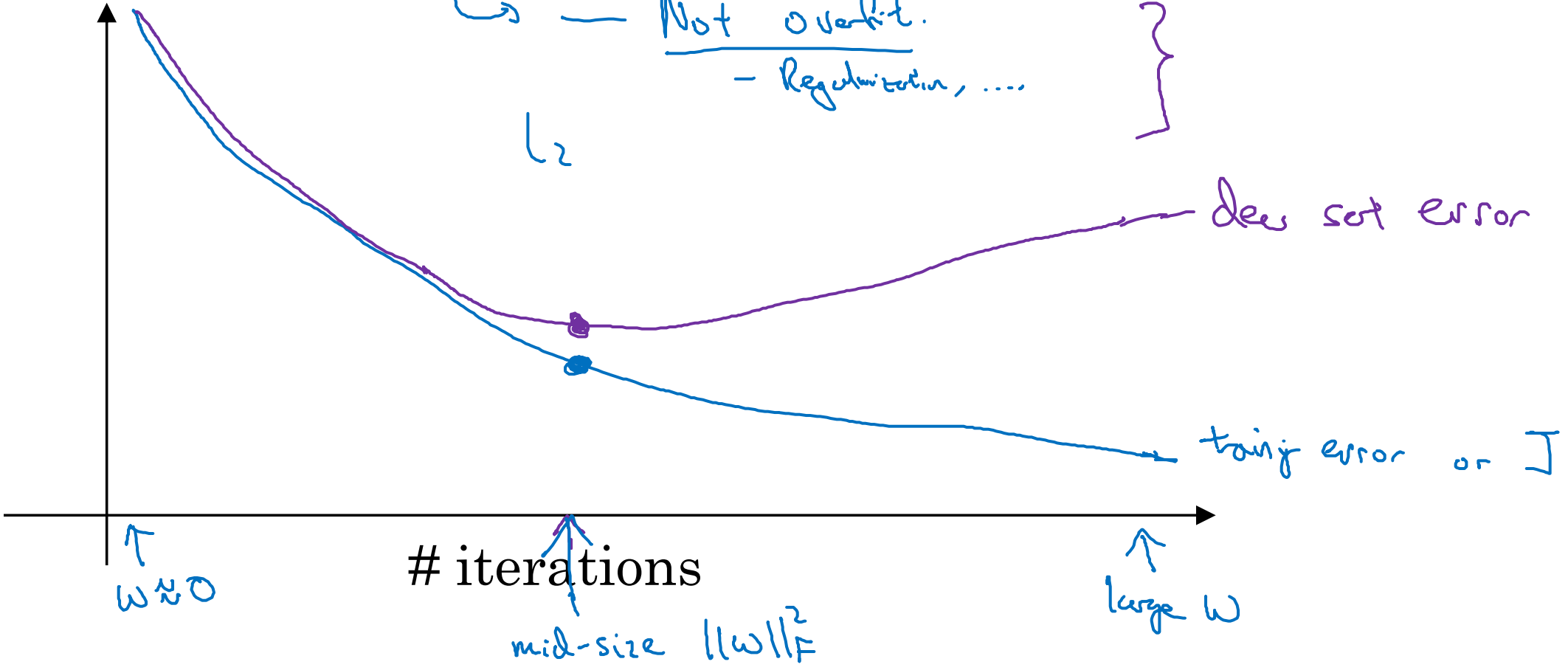
Early stopping

Orthogonalization.

Optimize cost function J
- Gradient, ...

Not overfit.
- Regularization, ...

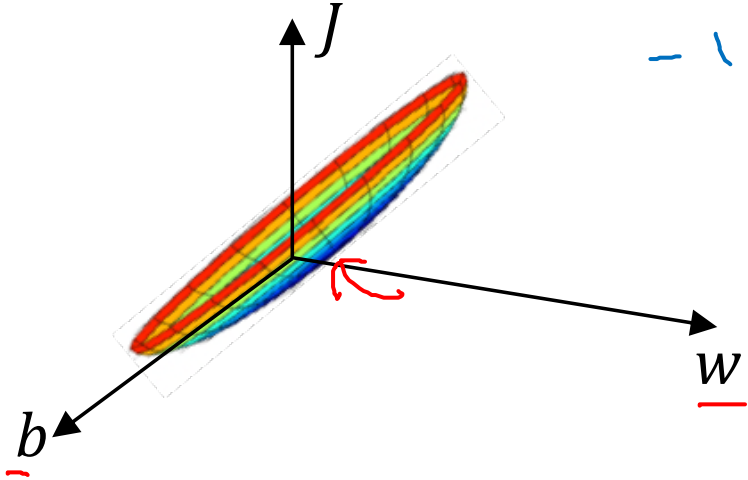
$J(w, b)$



Why normalize inputs?

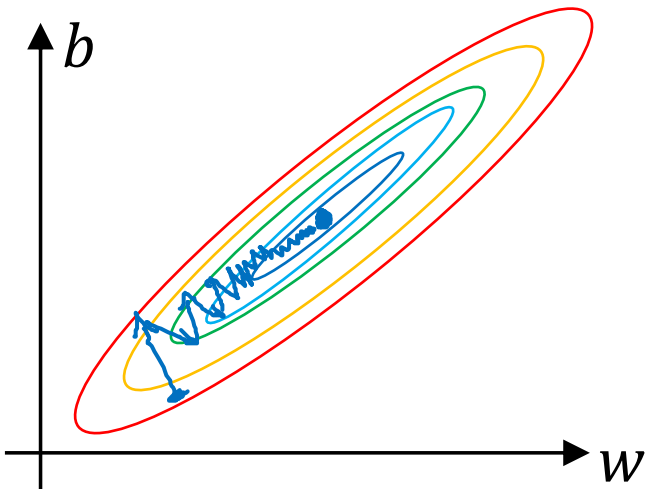
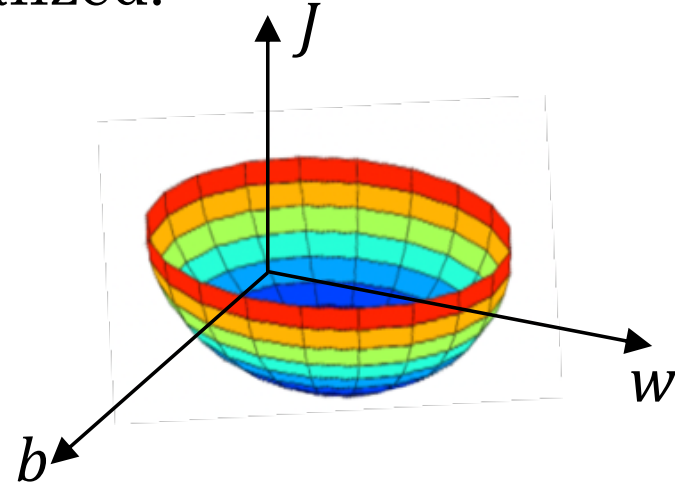
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Unnormalized:

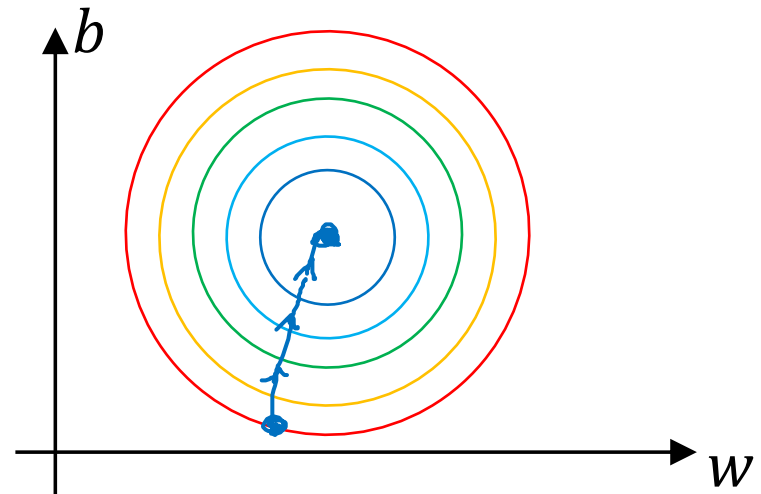


$w_1: x_1: 1 \dots 1000 \leftarrow$
 $w_2: x_2: 0 \dots 1 \leftarrow$
 $-1 \dots 1$

Normalized:



$x_1: 0 \dots 1$
 $x_2: -1 \dots 1$
 $x_3: 1 \dots 2$



Reducing (avoidable) bias and variance

Human-level



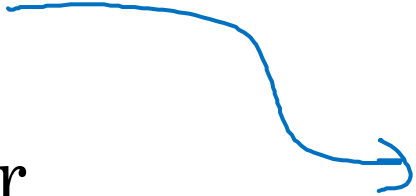
Avoidable bias



Training error



Variance



Dev error

Train bigger model

Train longer/better optimization algorithms
- momentum, RMSprop, Adam

NN architecture/hyperparameters search RNN
CNN

More data

Regularization

- L_2 , dropout, data augmentation

NN architecture/hyperparameters search

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>