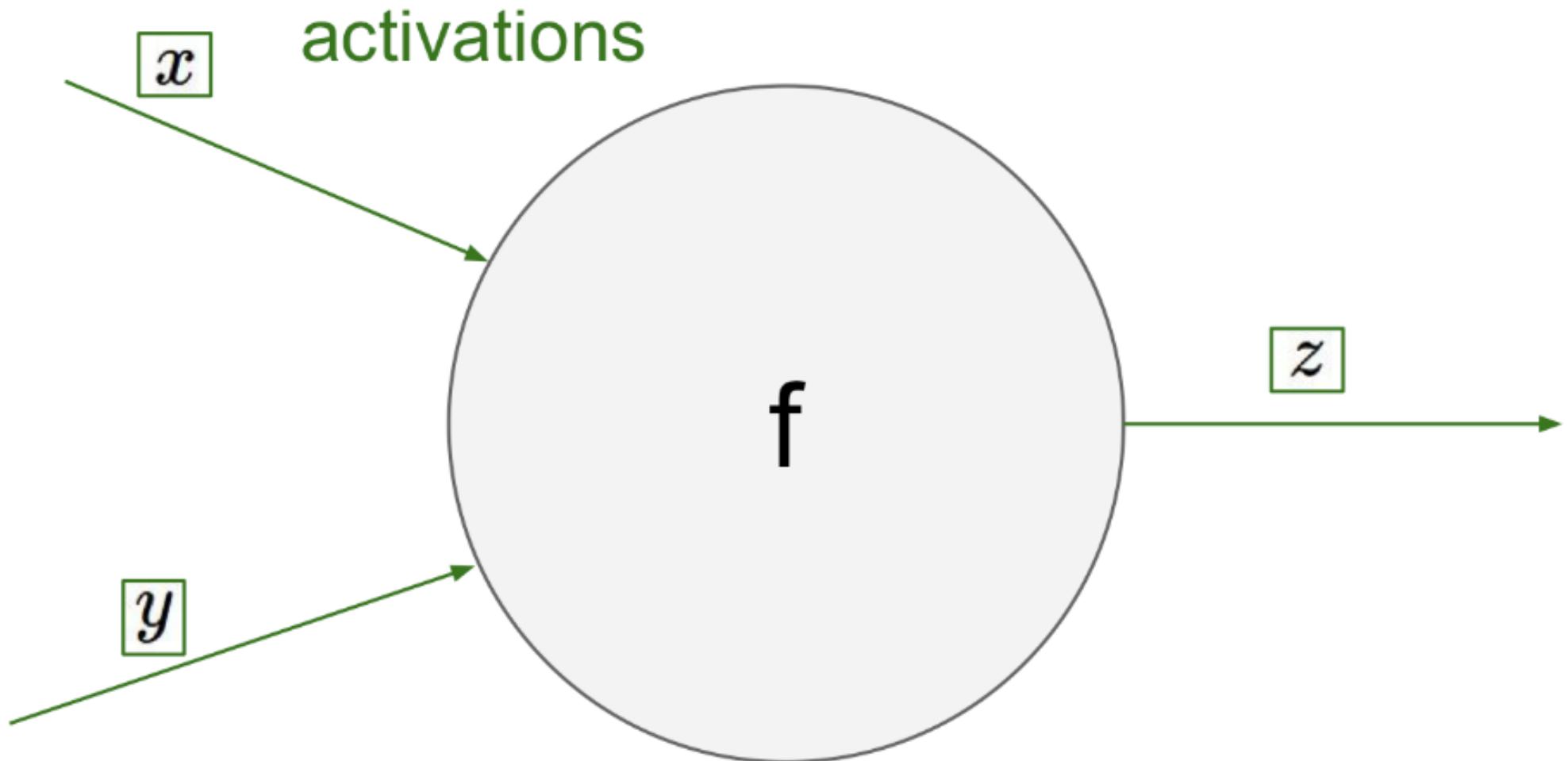


# BACKPROPAGATION

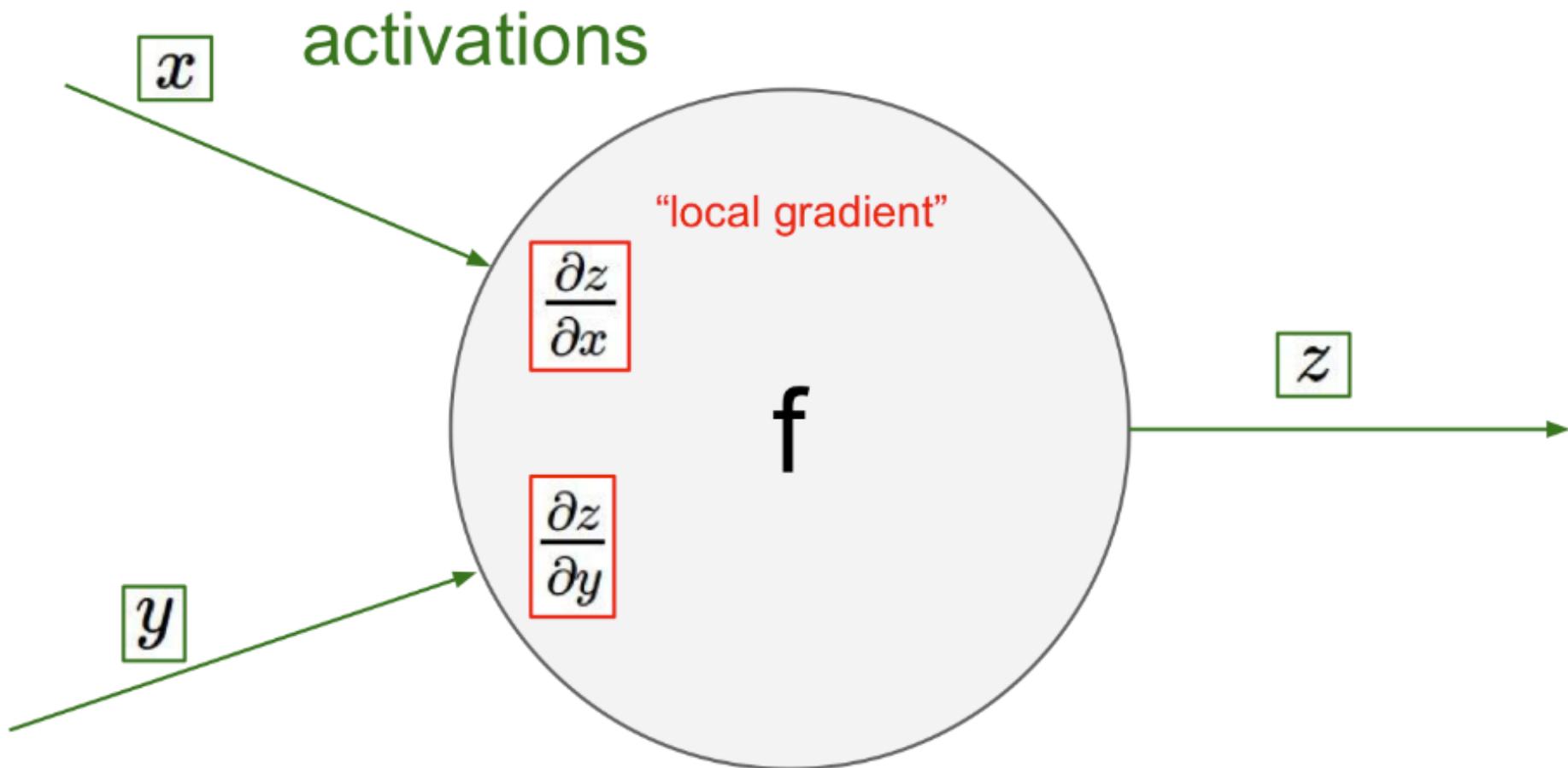
[AT THE NEURON LEVEL]



Credit: A. Karpathy

# BACKPROPAGATION

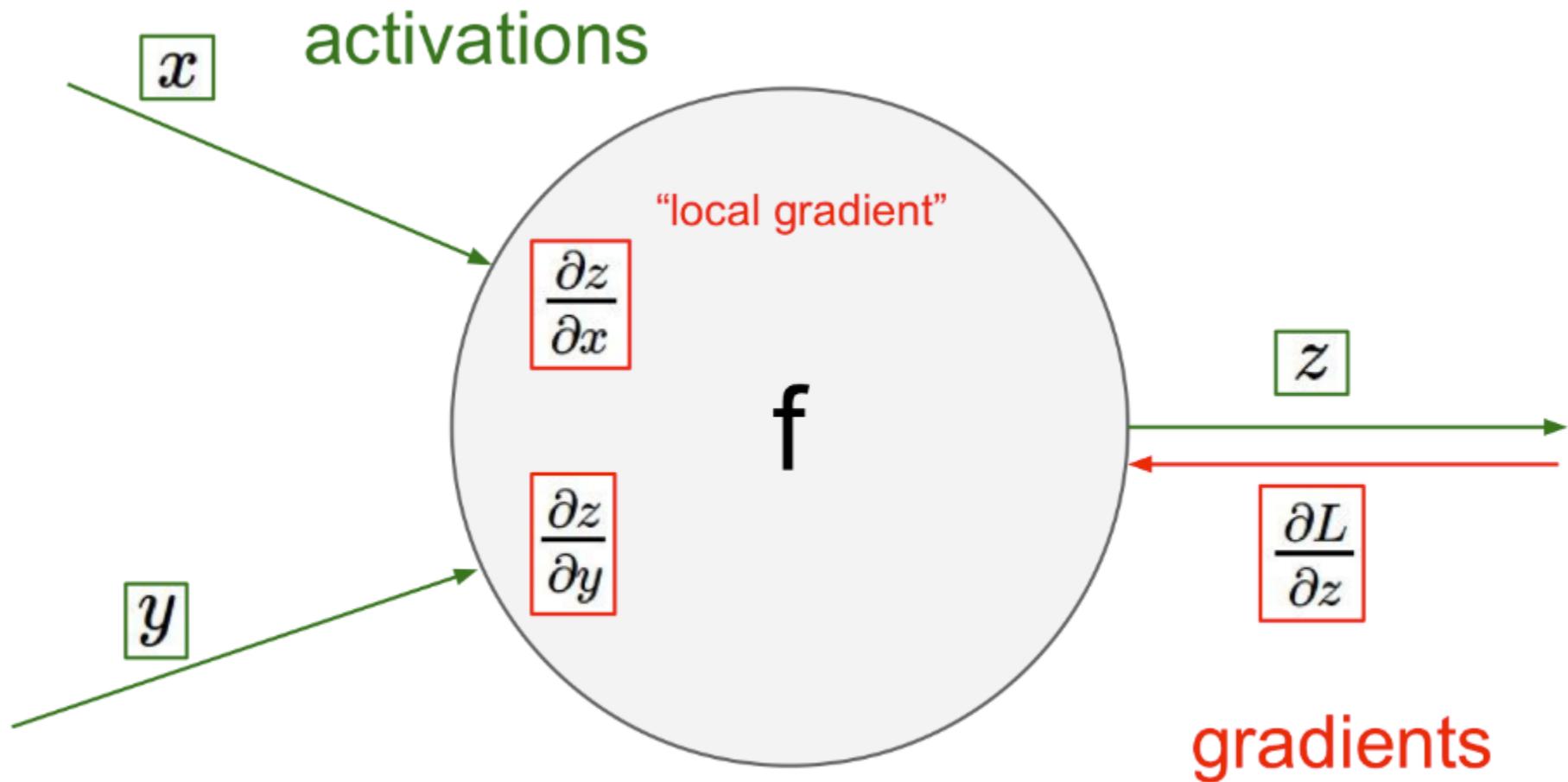
[AT THE NEURON LEVEL]



Credit: A. Karpathy

# BACKPROPAGATION

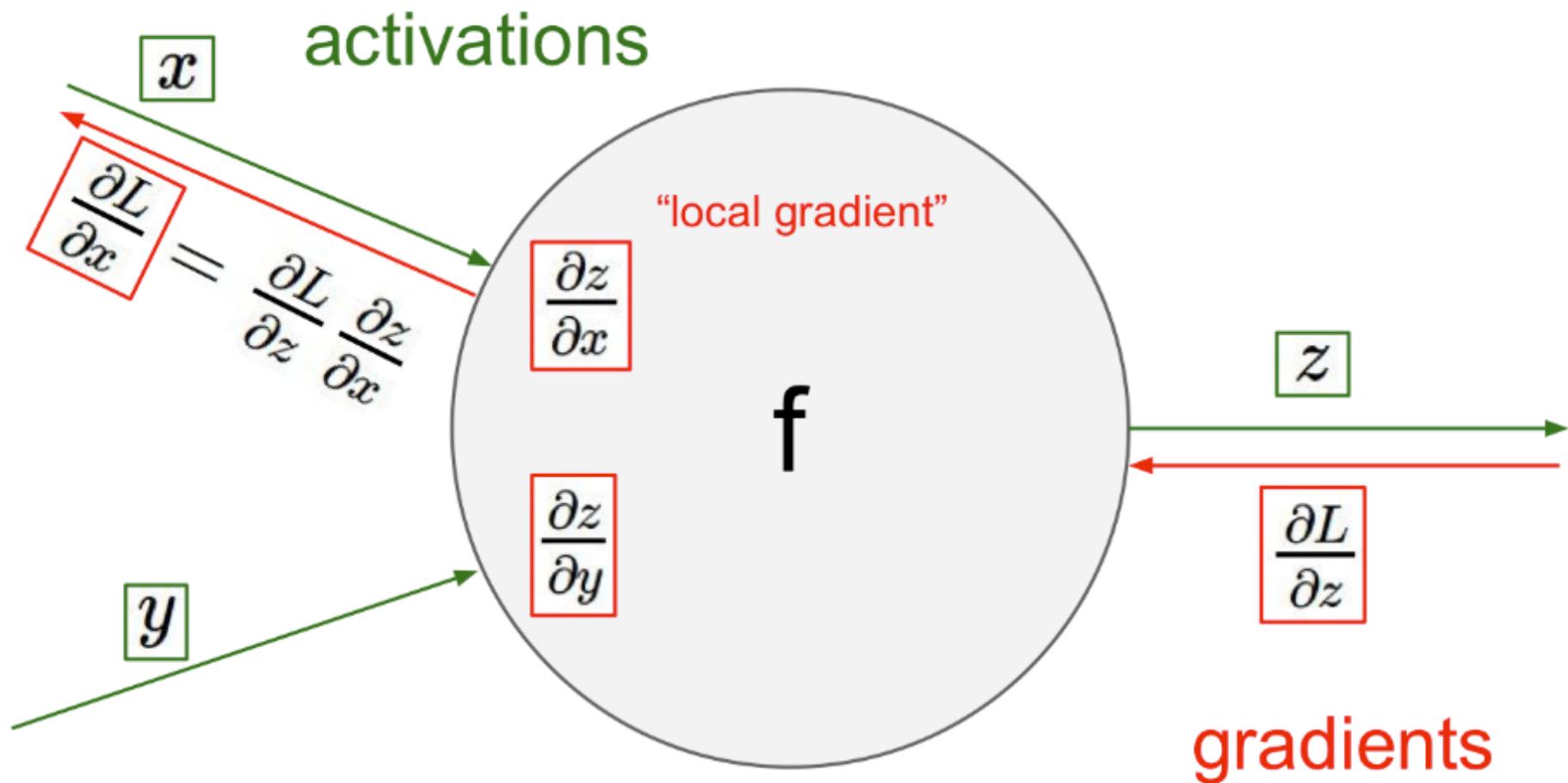
[AT THE NEURON LEVEL]



Credit: A. Karpathy

# BACKPROPAGATION

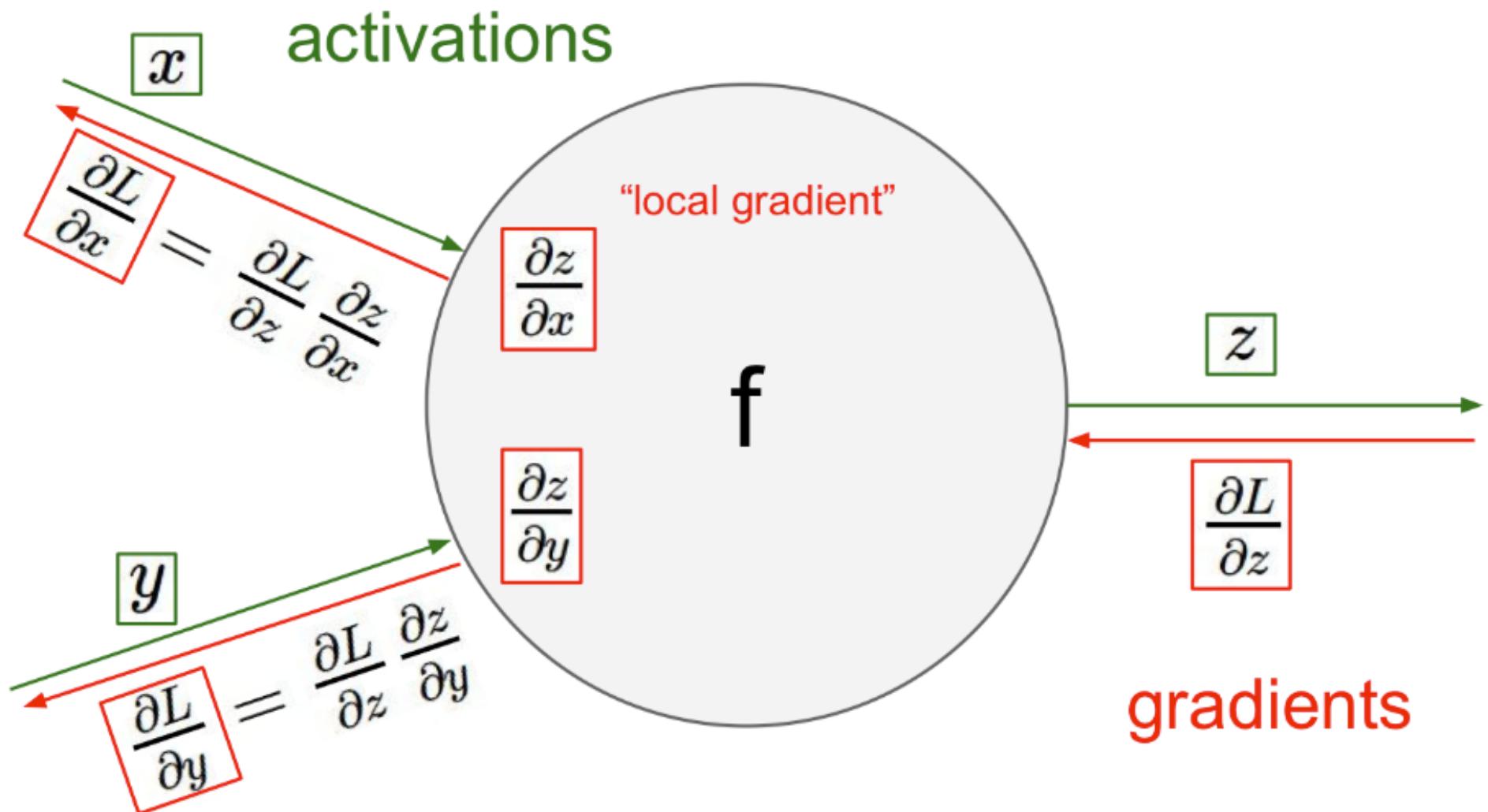
[AT THE NEURON LEVEL]



Credit: A. Karpathy

# BACKPROPAGATION

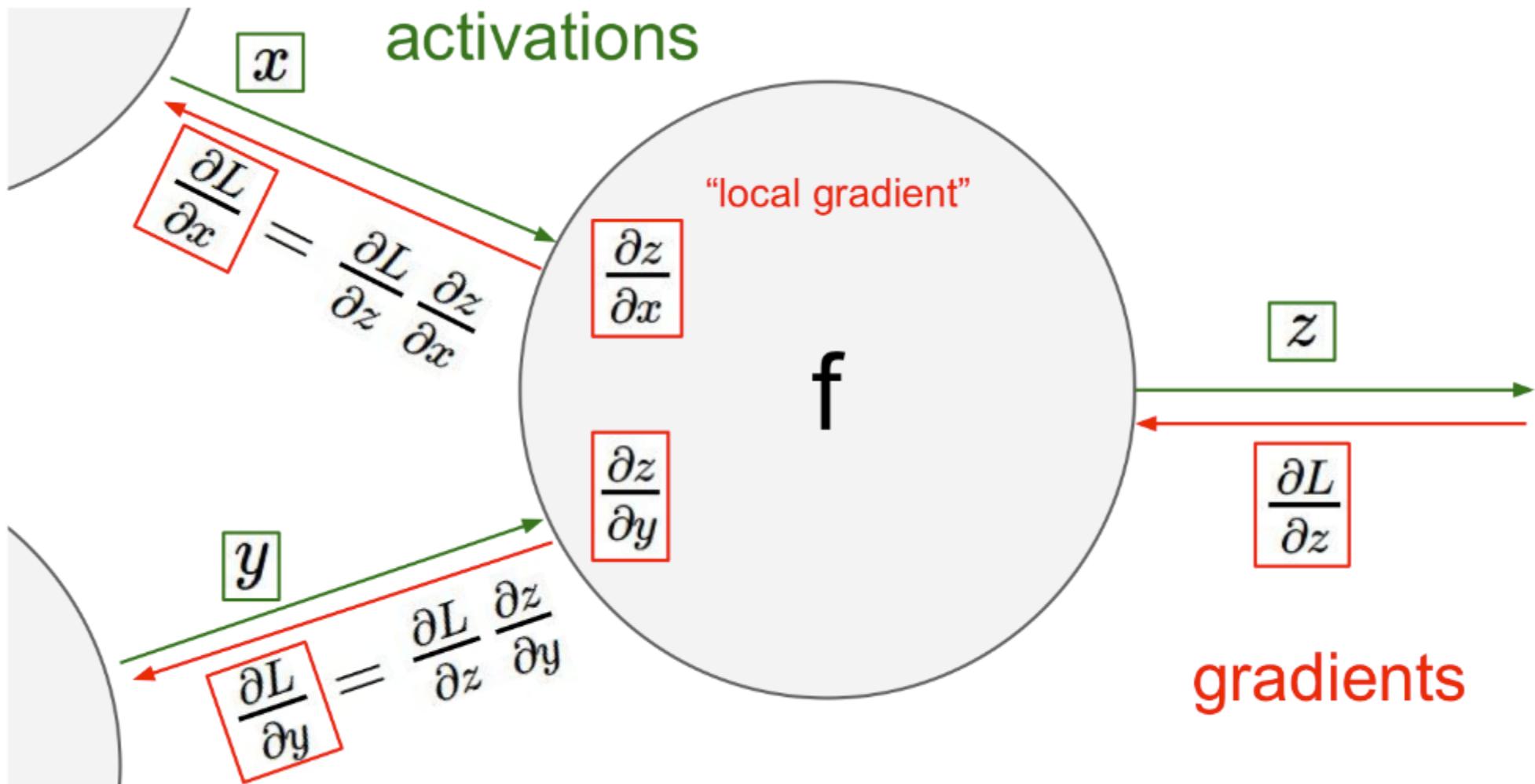
[AT THE NEURON LEVEL]



Credit: A. Karpathy

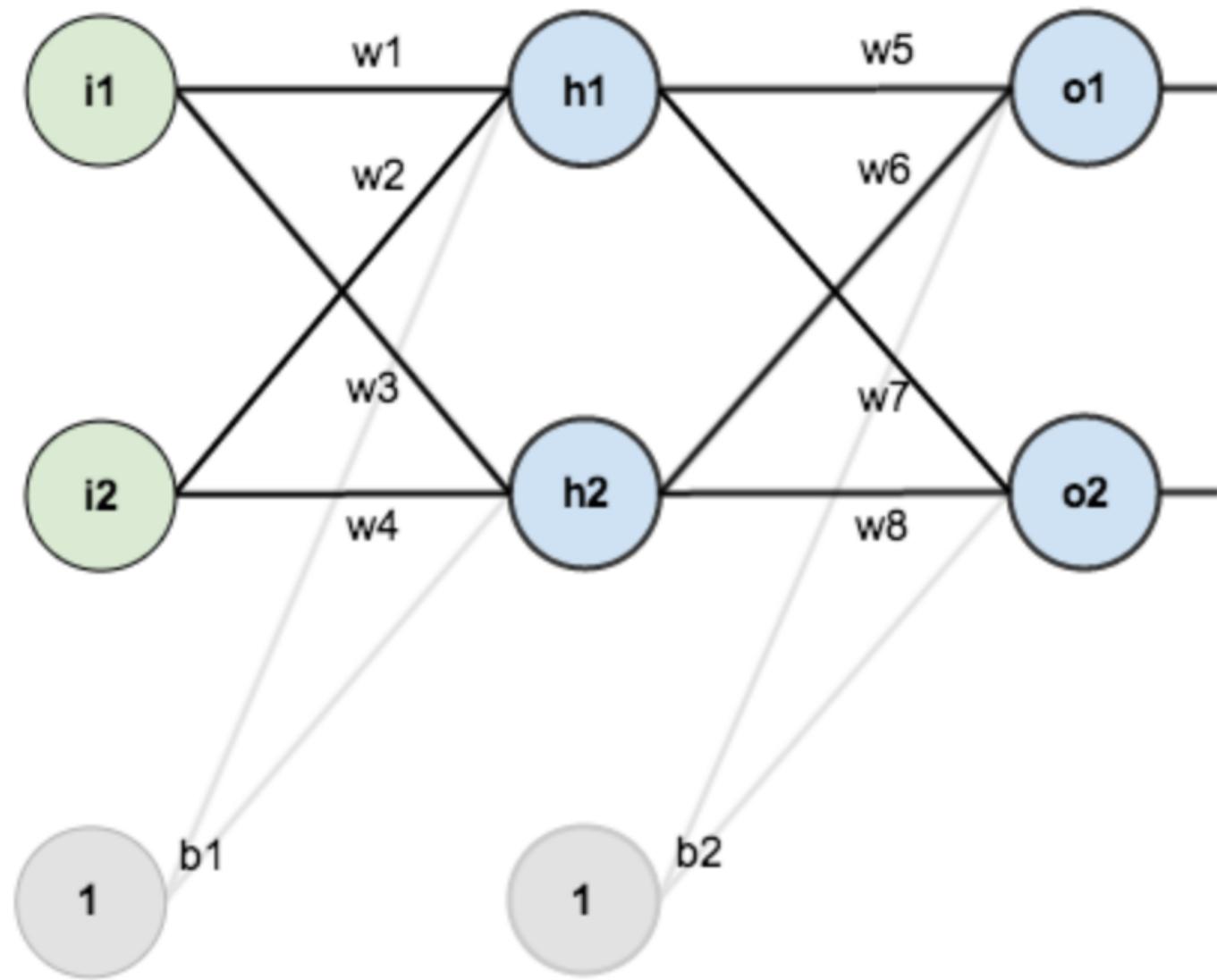
# BACKPROPAGATION

[AT THE NEURON LEVEL]

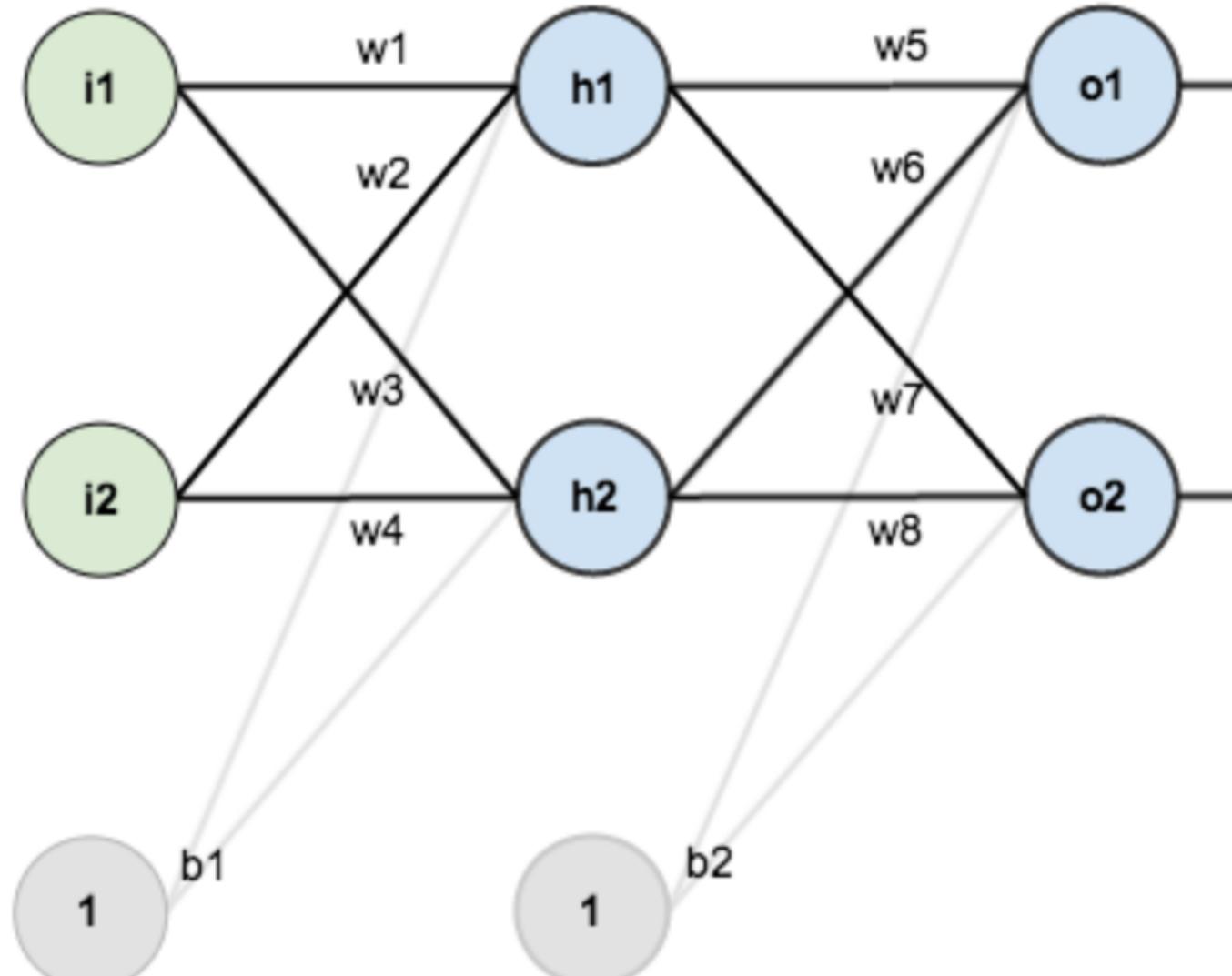


Credit: A. Karpathy

**LET'S FOLLOW A NETWORK  
WHILE IT LEARNS...**

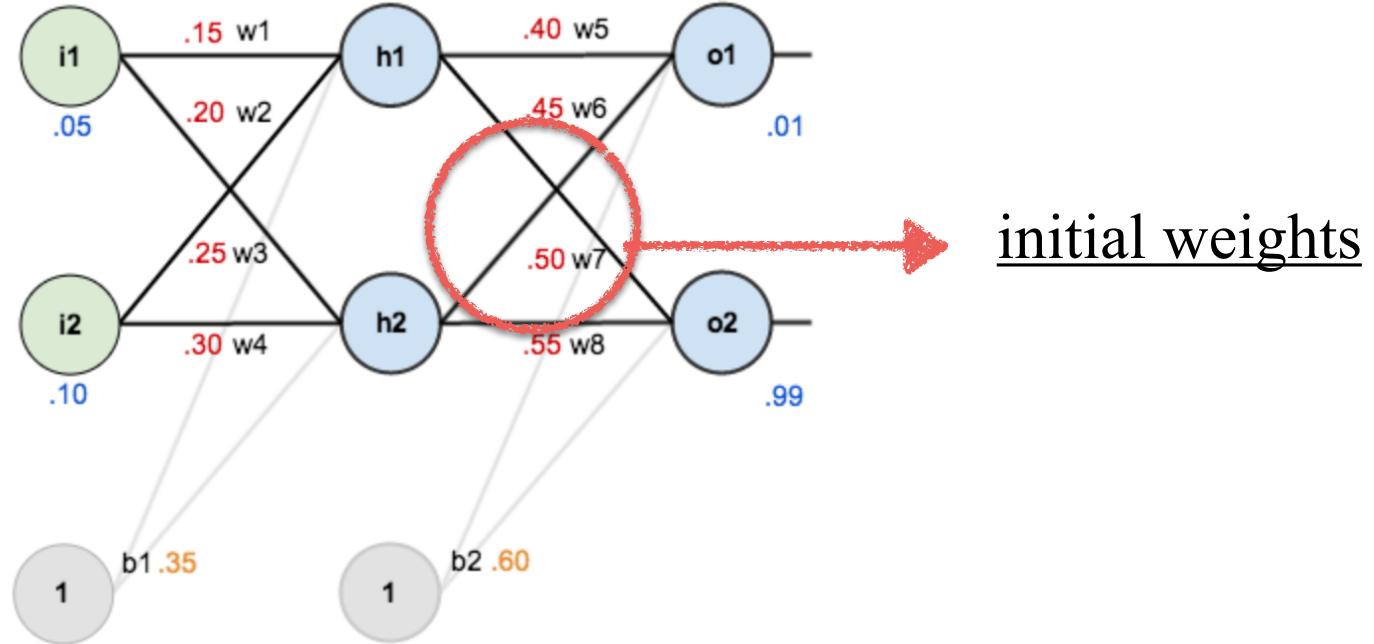


EXAMPLE TAKEN FROM HERE



LET'S ASSUME A VERY SIMPLE TRAINING SET:  
 $X=(0.05, 0.10) \rightarrow Y=(0.01, 0.99)$

EXAMPLE TAKEN FROM HERE

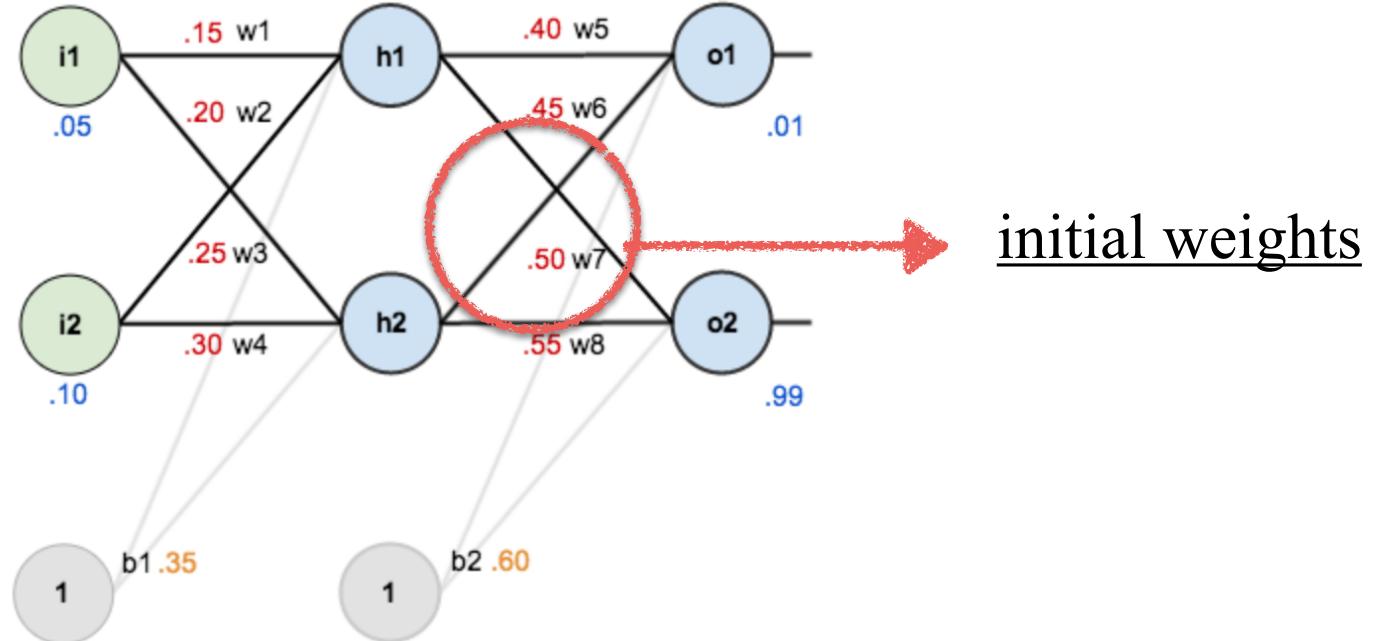


## 1. THE FORWARD PASS

$$in_{h1} = w_1 i_1 + w_2 i_2 + b_1$$

$$in_{h1} = 0.15 \times 0.05 + 0.2 \times 0.1 + 0.35 = 0.3775$$

[with initial weights]



## 1. THE FORWARD PASS

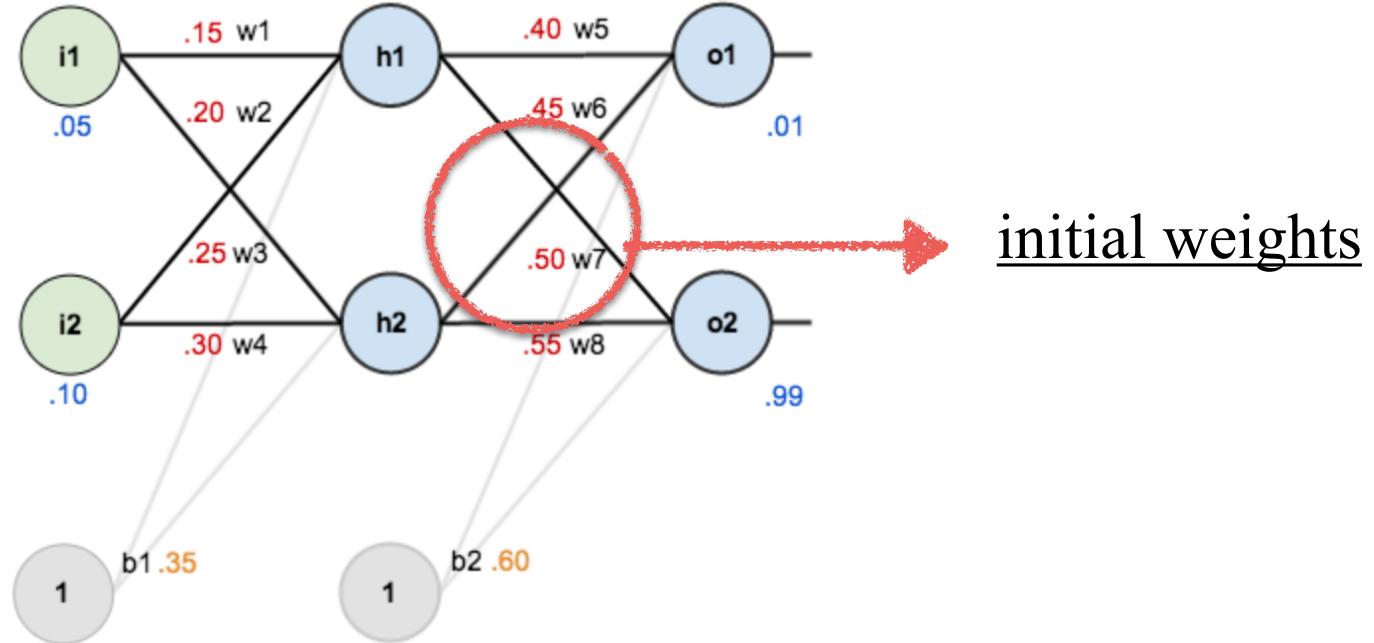
$$in_{h1} = w_1 i_1 + w_2 i_2 + b_1$$

$$in_{h1} = 0.15 \times 0.05 + 0.2 \times 0.1 + 0.35 = 0.3775$$

[with initial weights]

$$out_{h1} = \frac{1}{1 + e^{-in_{h1}}} = 0.5932$$

[after the activation function]



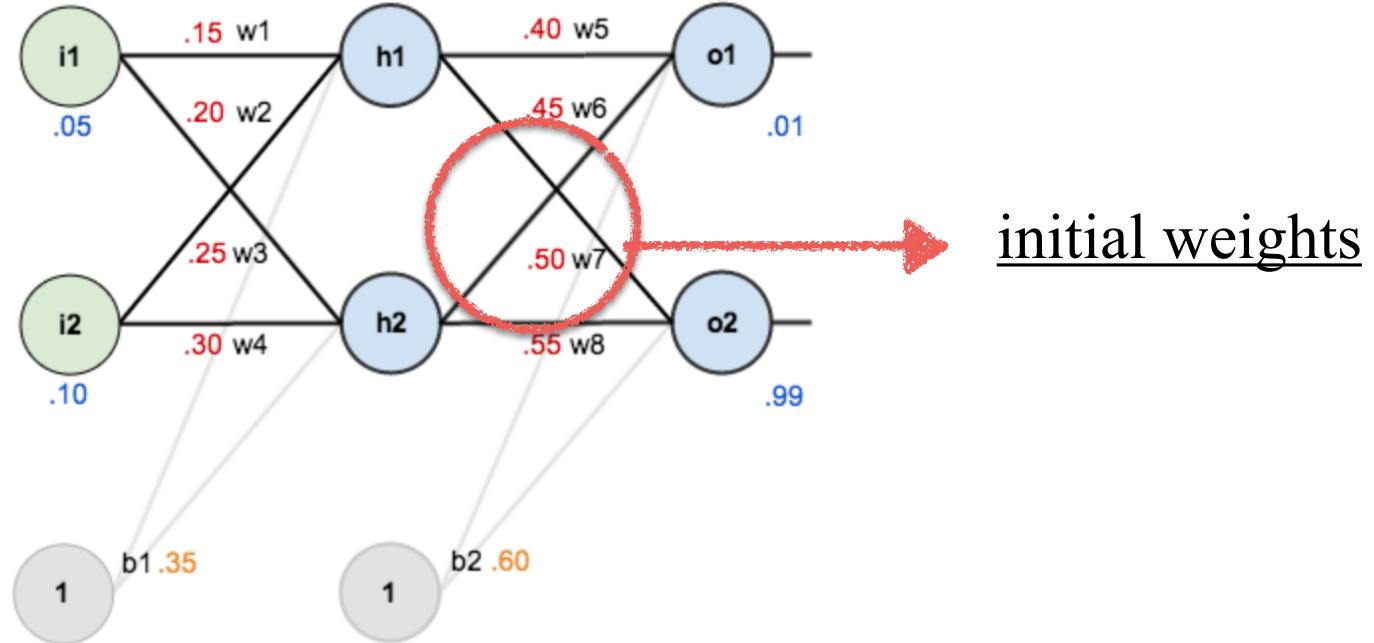
## 1. THE FORWARD PASS

WE CONTINUE TO  $o_1$

$$in_{o1} = w_5 out_{h1} + w_6 out_{h2} + b_2$$

$$in_{o1} = 0.4 \times 0.593 + 0.45 \times 0.596 + 0.6 = 1.105$$

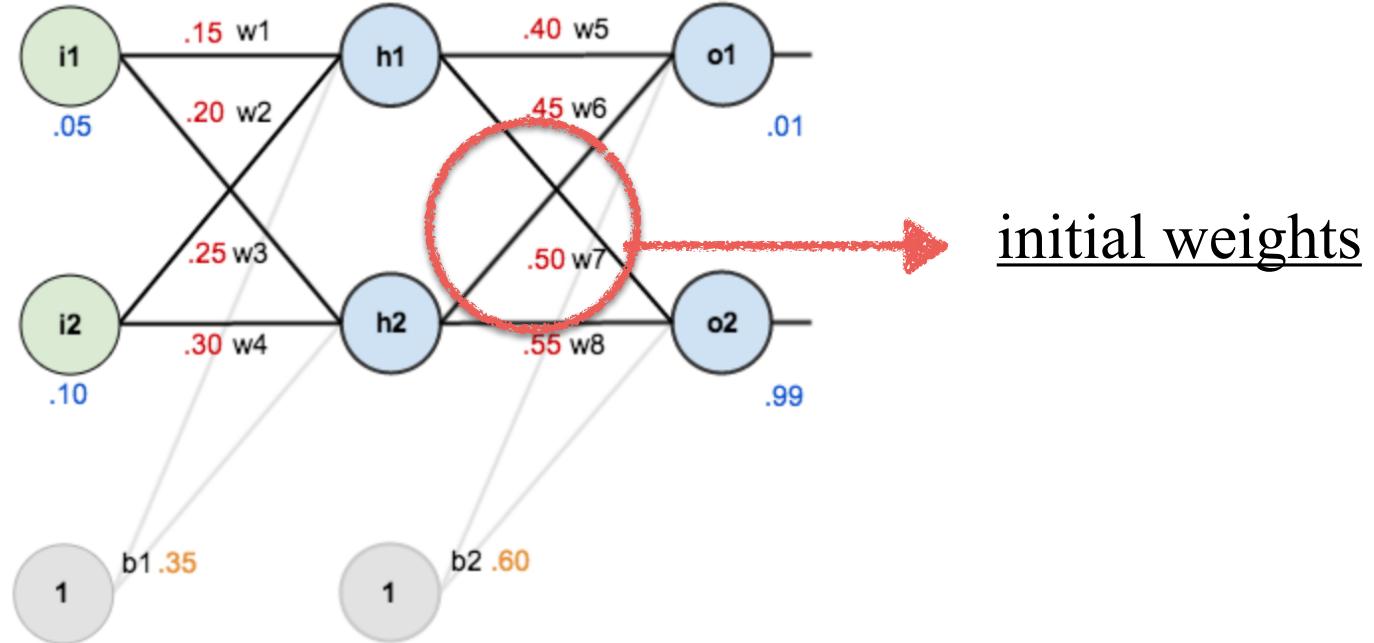
$$out_{o1} = \frac{1}{1 + e^{-1.105}} = 0.751$$



## 1. THE FORWARD PASS

AND THE SAME FOR  $o_2$

$$out_{o2} = 0.7729$$

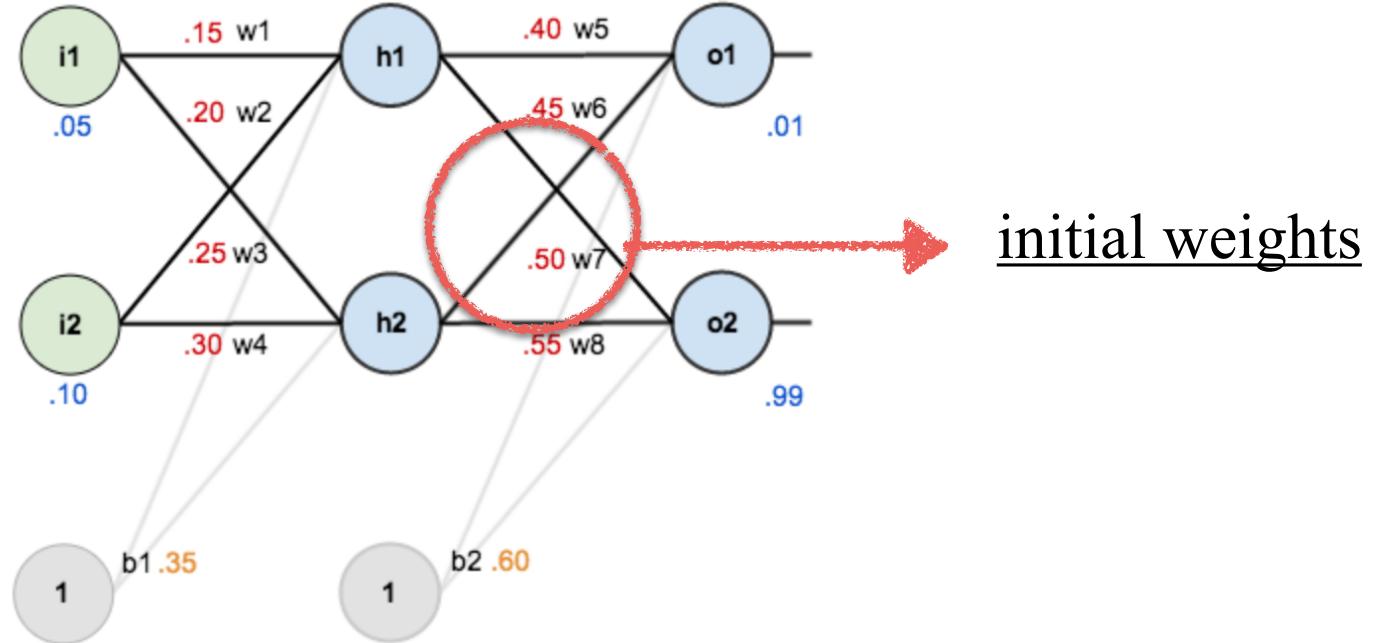


## 2. THE LOSS FUNCTION

$$L_{total} = \sum 0.5(target - output)^2$$

$$L_{o1} = 0.5(target_{o1} - output_{o1})^2 = 0.5 \times (0.01 - 0.751)^2 = 0.274$$

$$L_{o2} = 0.023$$



## 2. THE LOSS FUNCTION

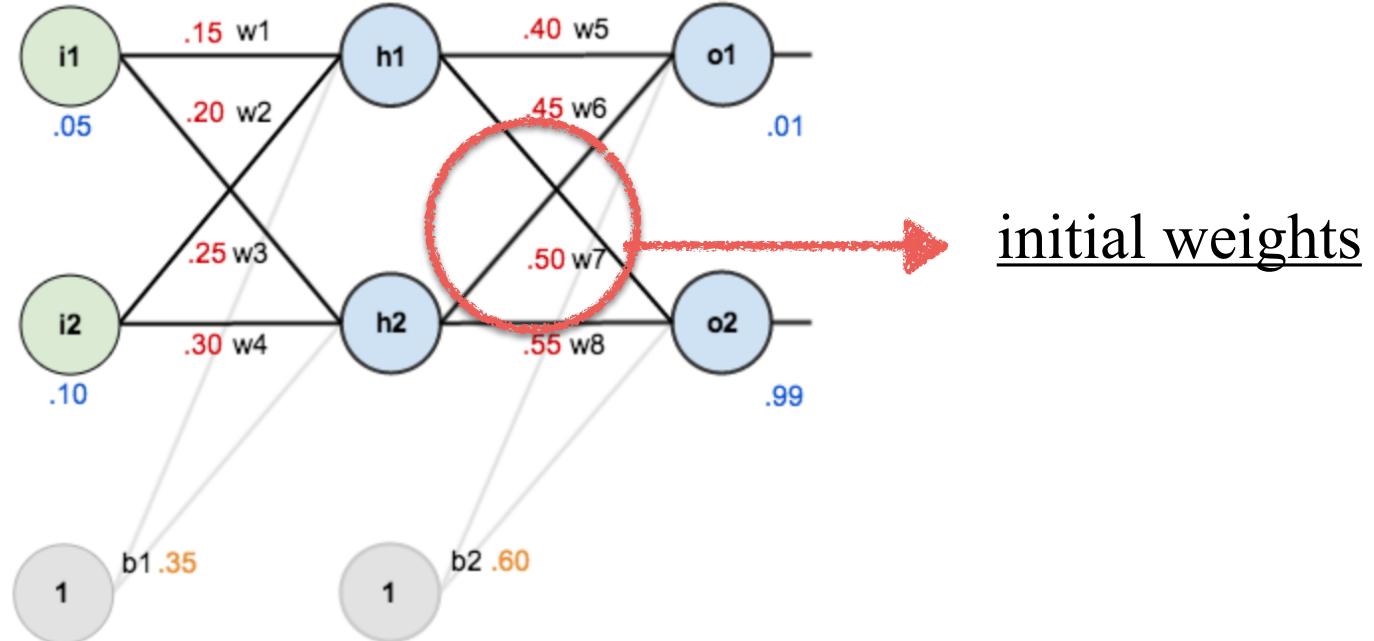
$$L_{total} = \sum 0.5(target - output)^2$$

$$L_{o1} = 0.5(target_{o1} - output_{o1})^2 = 0.5 \times (0.01 - 0.751)^2 = 0.274$$

$$L_{o2} = 0.023$$



$$L_{total} = L_{o1} + L_{o2} = 0.298$$



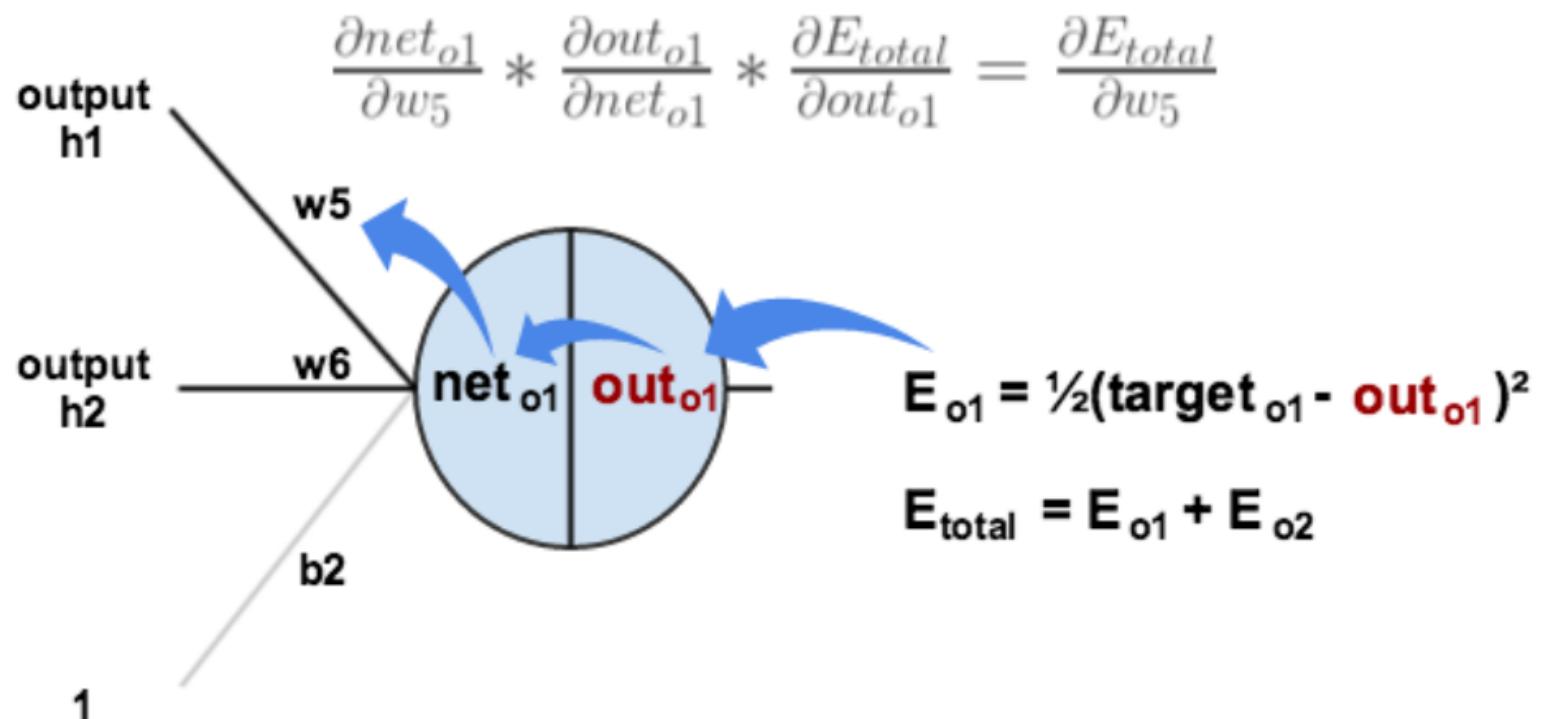
### 3. THE BACKWARD PASS

FOR W5

WE WANT:

$$\frac{\partial L_{total}}{\partial w_5}$$

[gradient of loss function]



### 3. THE BACKWARD PASS

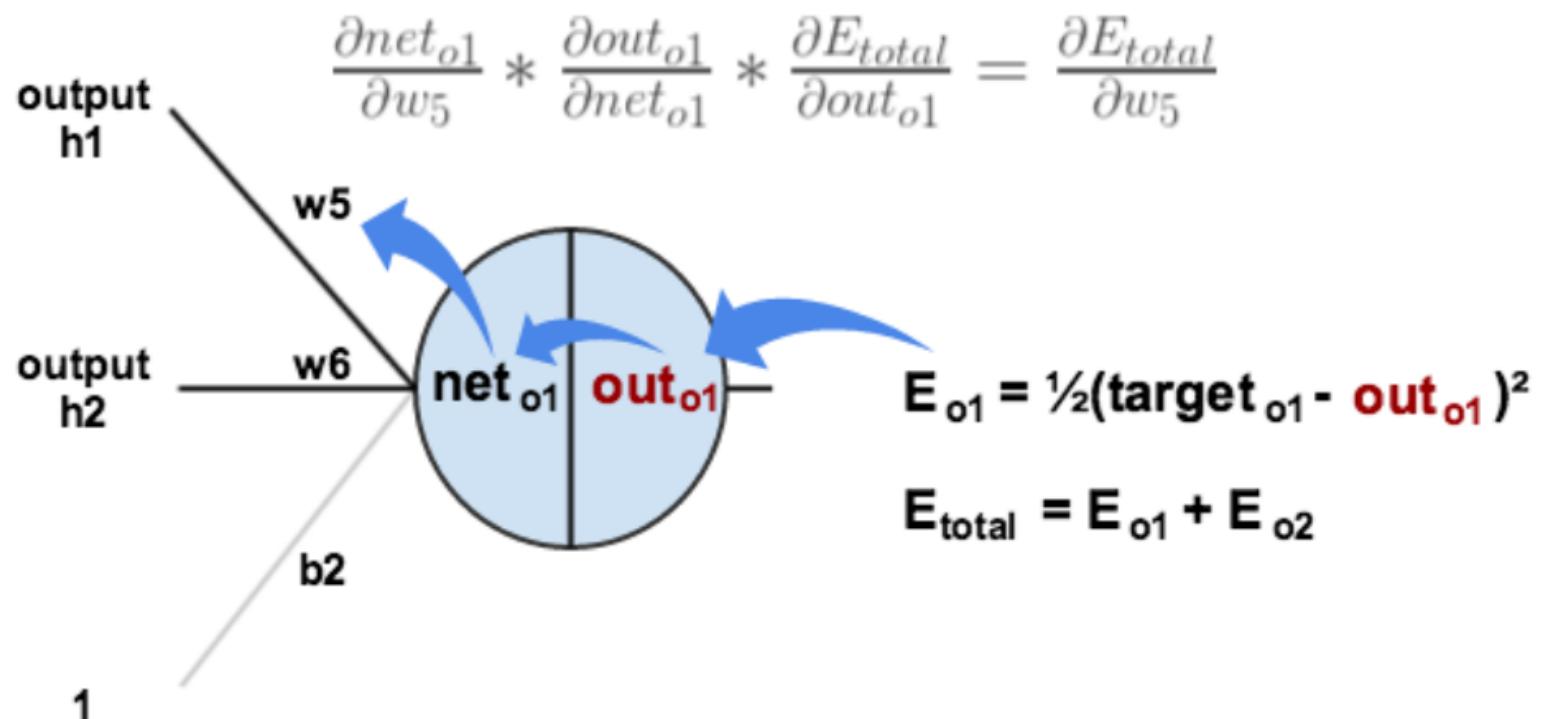
FOR  $w_5$

WE WANT:

$$\frac{\partial L_{total}}{\partial w_5} \quad [\text{gradient of loss function}]$$

WE APPLY THE CHAIN RULE:

$$\frac{\partial L_{total}}{\partial w_5} = \frac{\partial L_{total}}{\partial out_{o1}} \times \frac{\partial out_{o1}}{\partial in_{o1}} \times \frac{\partial in_{o1}}{\partial w_5}$$

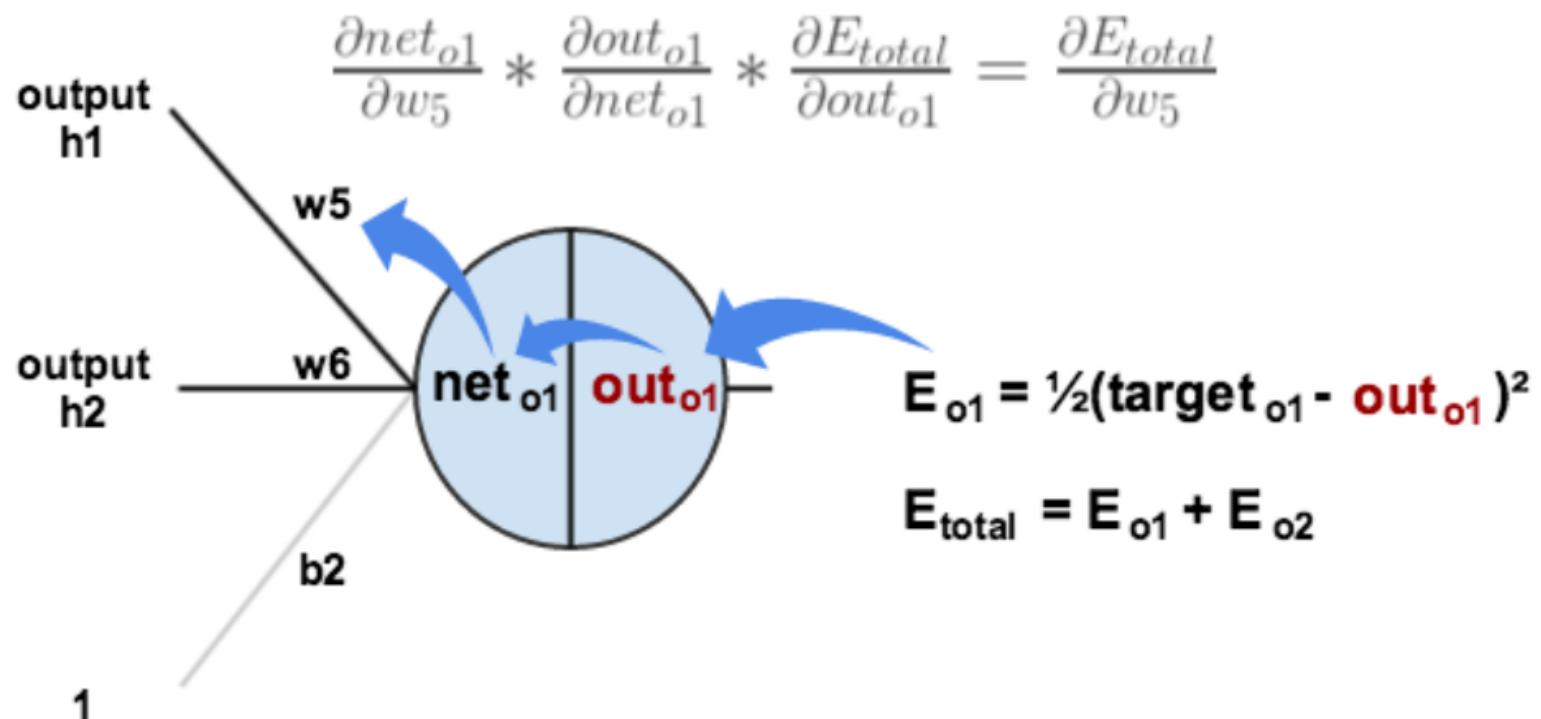


### 3. THE BACKWARD PASS

$$\frac{\partial L_{total}}{\partial w_5} = \frac{\partial L_{total}}{\partial out_{o1}} \times \frac{\partial out_{o1}}{\partial in_{o1}} \times \frac{\partial in_{o1}}{\partial w_5}$$

$$L_{total} = 0.5(target_{o1} - out_{o1})^2 + 0.5(target_{o2} - out_{o2})^2$$

$$\frac{\partial L_{total}}{\partial out_{o1}} = 2 \times 0.5(target_{o1} - out_{o1}) \times (-1) = 0.741$$

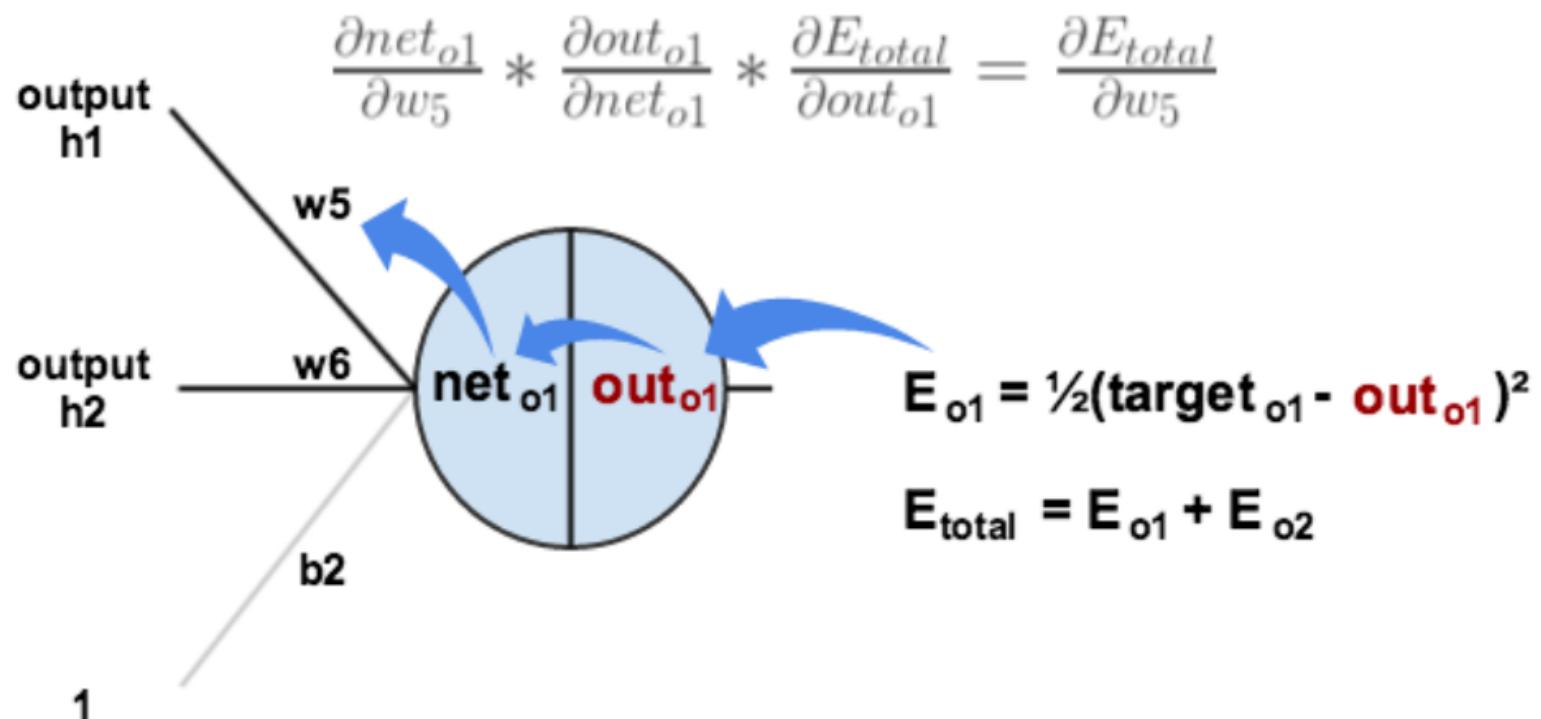


### 3. THE BACKWARD PASS

$$\frac{\partial L_{\text{total}}}{\partial w_5} = \frac{\partial L_{\text{total}}}{\partial out_{o1}} \times \frac{\partial out_{o1}}{\partial in_{o1}} \times \frac{\partial in_{o1}}{\partial w_5}$$

$$out_{o1} = \frac{1}{1 + e^{-in_{o1}}}$$

$$\frac{\partial out_{o1}}{\partial in_{o1}} = out_{o1} \times (1 - out_{o1}) = 0.186$$

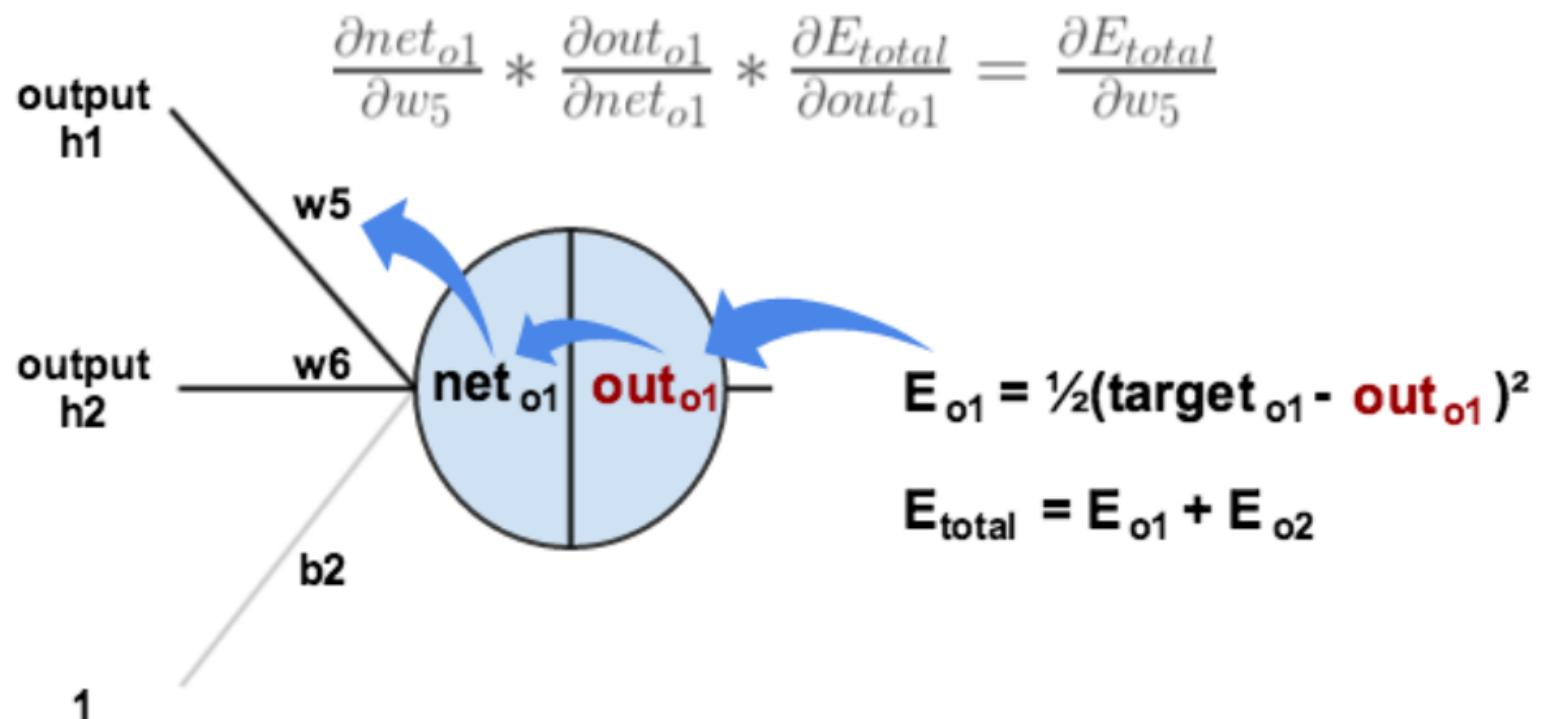


### 3. THE BACKWARD PASS

$$\frac{\partial L_{total}}{\partial w_5} = \frac{\partial L_{total}}{\partial out_{o1}} \times \frac{\partial out_{o1}}{\partial in_{o1}} \times \frac{\partial in_{o1}}{\partial w_5}$$

$$in_{o1} = w_5 \times out_{h1} + w_6 \times out_{h2} + b_2$$

$$\frac{\partial in_{o1}}{\partial w_5} = out_{h1} \times w_5^{1-1} = out_{h1} = 0.593$$

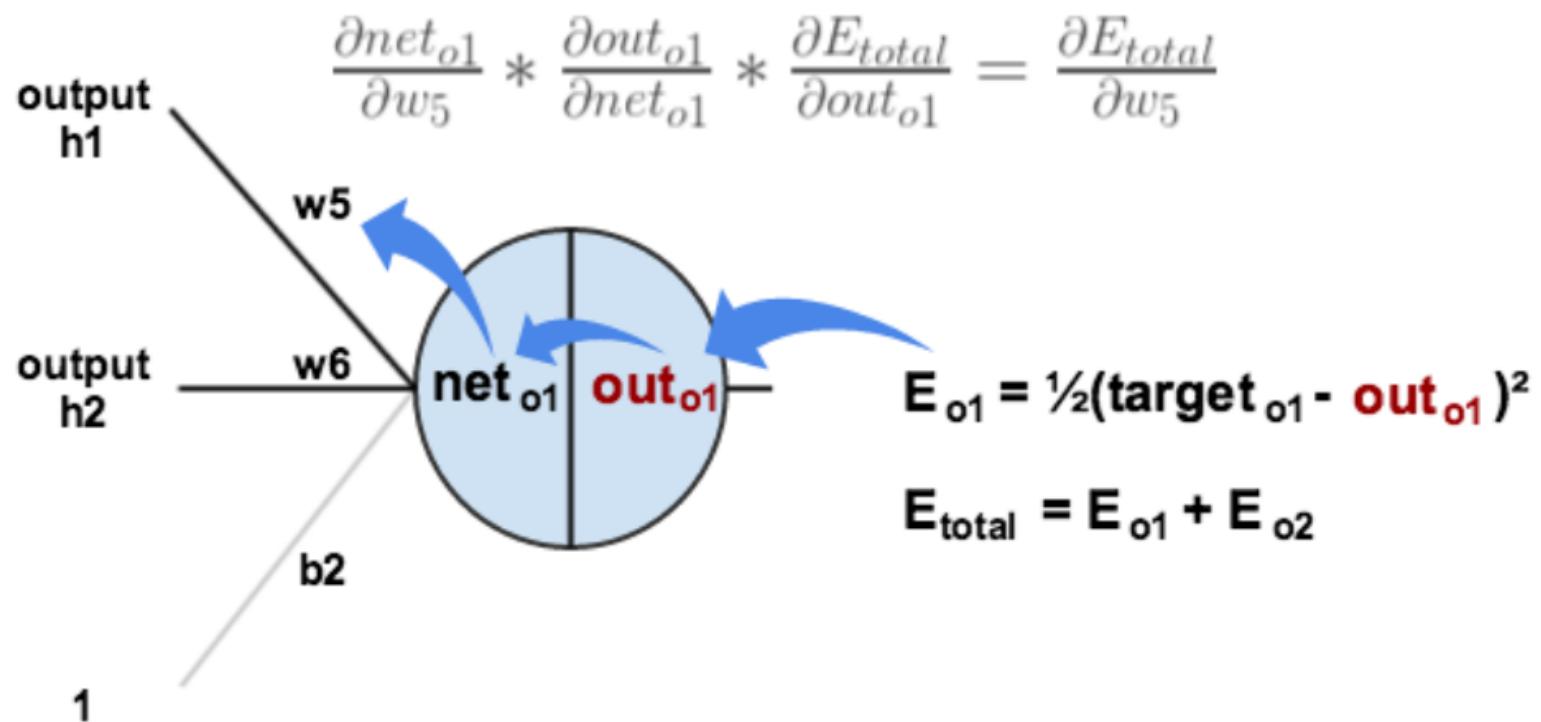


### 3. THE BACKWARD PASS

ALL TOGETHER:

$$\frac{\partial L_{total}}{\partial w_5} = \frac{\partial L_{total}}{\partial out_{o1}} \times \frac{\partial out_{o1}}{\partial in_{o1}} \times \frac{\partial in_{o1}}{\partial w_5}$$

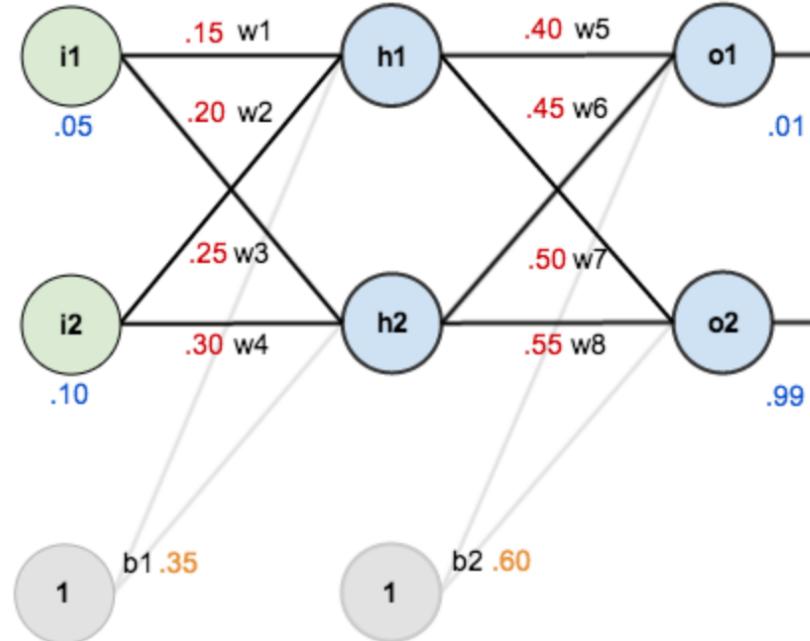
$$\frac{\partial L_{total}}{\partial w_5} = 0.741 \times 0.186 \times 0.593 = 0.082$$



## 4. UPDATE WEIGHTS WITH GRADIENT AND LEARNING RATE

$$w_5^{t+1} = w_5 - \lambda \times \frac{\partial L_{total}}{\partial w_5}$$

$$w_5^{t+1} = 0.4 - 0.5 \times 0.082 = 0.358$$



**THIS IS REPEATED FOR THE OTHER WEIGHTS  
OF THE OUTPUT LAYER**

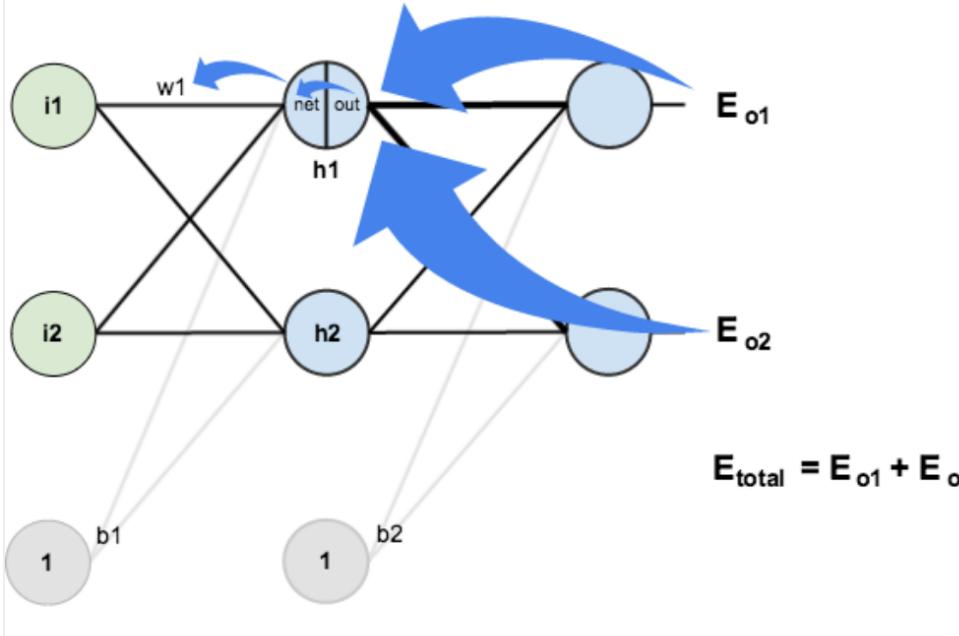
$$w_6^{t+1} = 0.408$$

$$w_7^{t+1} = 0.511$$

$$w_8^{t+1} = 0.561$$

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$



**AND BACK-PROPAGATED TO THE HIDDEN LAYERS**

# VISUALIZE SIMPLE NETWORK LEARNING

IN ORDER TO KNOW MORE  
ABOUT BACKPROPAGATION...

[CLICK HERE](#)

# LET'S TALK LOSS FUNCTIONS....

THE CHOICE OF THE LOSS FUNCTION IS CRITICAL AND WILL DETERMINE THE BEHAVIOR OF THE ALGORITHM  
*(REMEMBER: LOSS IS WHAT IS MINIMIZED BY THE OPTIMIZATION ALGORITHM)*

SOME ALGORITHMS SUCH AS RFs HAVE LESS FLEXIBILITY.  
ONE **ADVANTAGE OF NNs IS THAT THE LOSS FUNCTION CAN BE CHANGED VERY EASILY.**

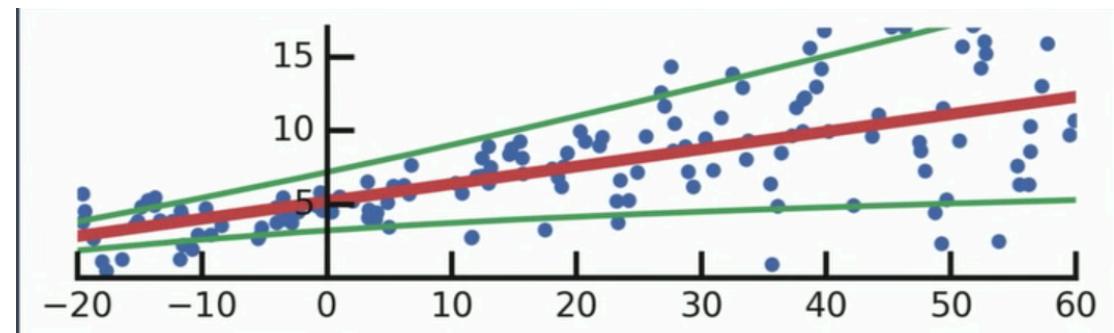
IT IS IMPORTANT TO CHOOSE THE LOSS FUNCTION FOR YOUR PROBLEM

SO, IN A NUTSHELL, WE ARE DOING BAYESIAN VARIATIONAL INFERENCE ASSUMING THEREFORE A PDF FOR OUR OUTPUT VARIABLE (BERNOUILLI FOR CLASSIFICATION, NORMAL FOR REGRESSION)

Tensorflow probability allows now to generalize this to any distribution:

**Guess what! You can.**

```
model = tf.keras.Sequential([
    tf.keras.layers.Dense(hidden_units, ...),
    tf.keras.layers.Dense(1),
    tfp.layers.DistributionLambda(lambda t:
        tfd.Normal(loc=t[...], 0,
                   scale=1)),
```

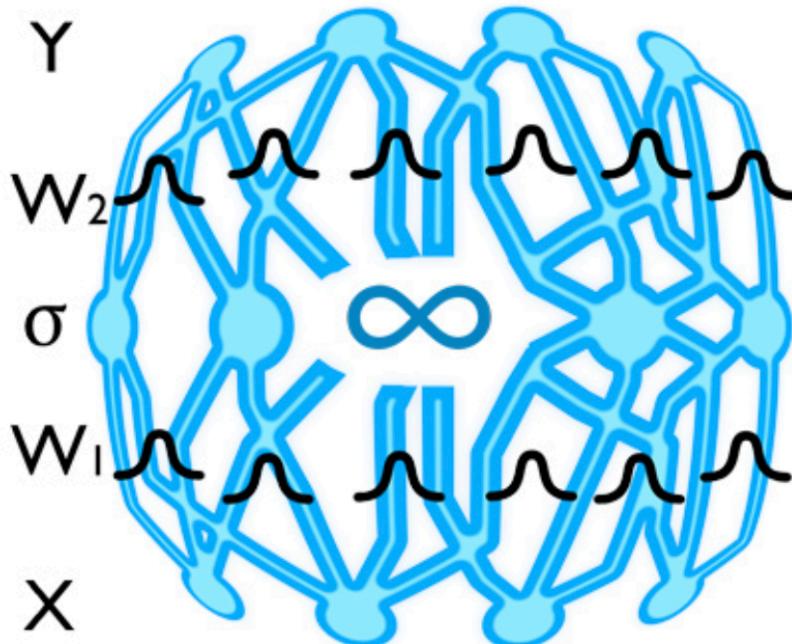


```
model.compile(loss=tf.keras.mean_squared_error)
model.compile(loss=lambda y, rv_y: -rv_y.log_prob(y)) } Our wish.
yhat = model(x_tst) # A Distribution, not a Tensor!
```

**TensorFlow**  
DEV SUMMIT 2019

SEE THIS GREAT TALK: <https://www.youtube.com/watch?v=BrwKURU-wpk>

# BUT REMEMBER, WE ARE CHOOSING ONE MODEL AND ONE VALUE FOR THE WEIGHTS



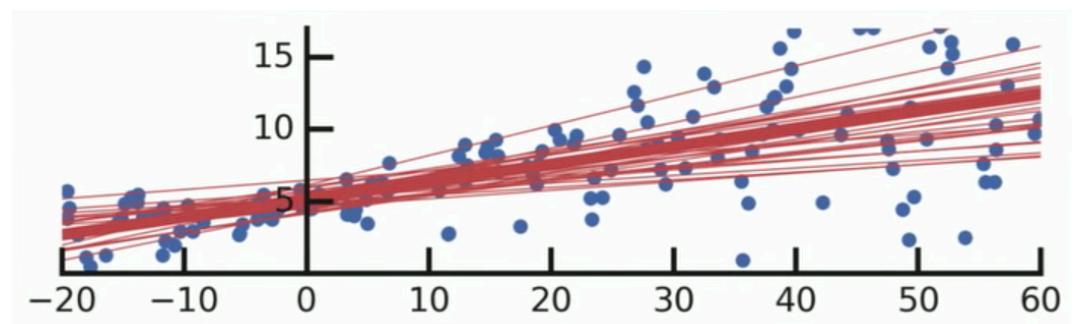
IDEALLY, WE WOULD LIKE TO MARGINALIZE OVER ALL MODELS AND OVER ALL POSSIBLE WEIGHT VALUES

BNNs ADD A PRIOR DISTRIBUTION TO EACH WEIGHT - HARD TO TRAIN

## CAPTURING “MODEL UNCERTAINTY”

```
model = tf.keras.Sequential([
    tfp.layers.DenseVariational(hidden_units, ...),
    tfp.layers.DenseVariational(1),
    tfp.layers.DistributionLambda(lambda t:
        tfd.Normal(loc=t[...], 0),
        scale=1)),
])
```

} Bayesian Weights  
} Linear Regression



# AND YOU STILL ARE LEFT WITH THE CHOICE OF THE MODEL ARCHITECTURE:

## Introduction to the Keras Tuner

[Run in Google Colab](#)[View source on GitHub](#)[Download notebook](#)

### Overview

The Keras Tuner is a library that helps you pick the optimal set of hyperparameters for your TensorFlow program. The process of selecting the right set of hyperparameters for your machine learning (ML) application is called *hyperparameter tuning* or *hypertuning*.

Hyperparameters are the variables that govern the training process and the topology of an ML model. These variables remain constant over the training process and directly impact the performance of your ML program. Hyperparameters are of two types:

1. **Model hyperparameters** which influence model selection such as the number and width of hidden layers
2. **Algorithm hyperparameters** which influence the speed and quality of the learning algorithm such as the learning rate for Stochastic Gradient Descent (SGD) and the number of nearest neighbors for a k Nearest Neighbors (KNN) classifier

In this tutorial, you will use the Keras Tuner to perform hypertuning for an image classification application.

ONE KEY PROBLEM WITH GRADIENT DESCENT IS THAT IT  
EASILY CONVERGES TO LOCAL MINIMA BY FOLLOWING  
THE STEEPEST DESCENT

ONE KEY PROBLEM WITH GRADIENT DESCENT IS THAT IT  
EASILY CONVERGES TO LOCAL MINIMA BY FOLLOWING  
THE STEEPEST DESCENT

THE CHOICES OF THE INITIAL WEIGHTS AND THE  
LEARNING RATES ARE IMPORTANT

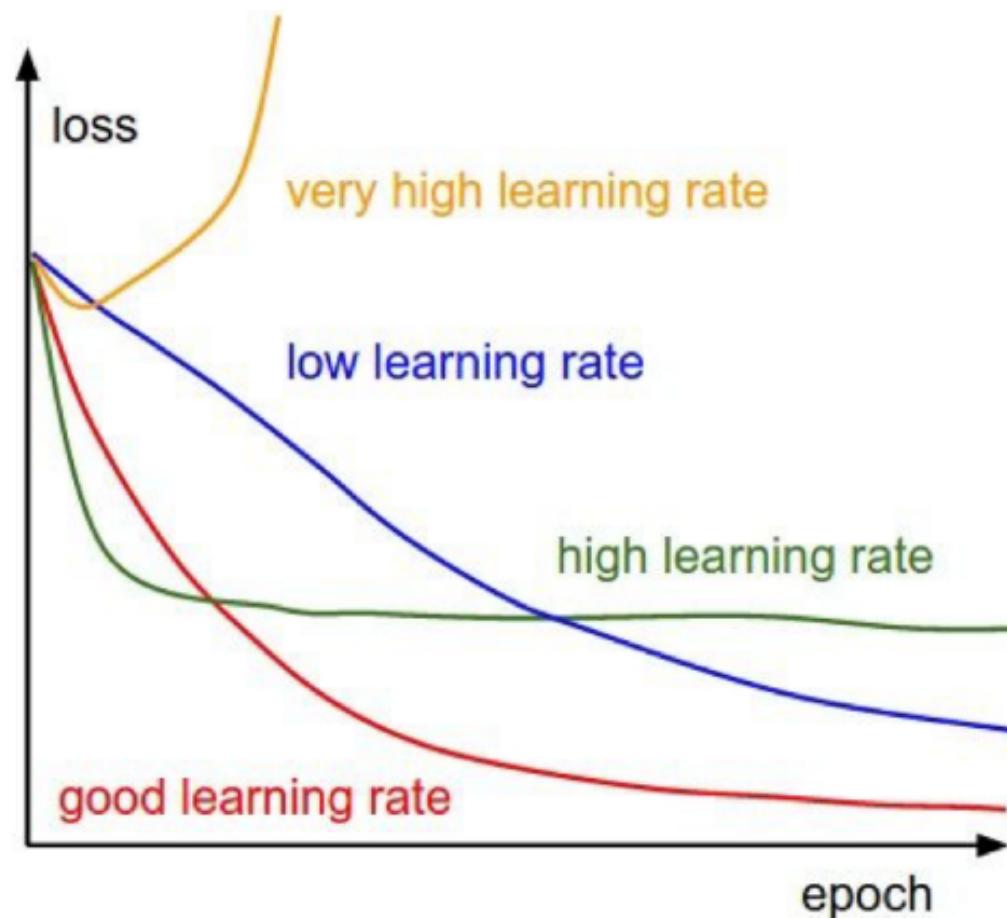
ONE KEY PROBLEM WITH GRADIENT DESCENT IS THAT IT  
EASILY CONVERGES TO LOCAL MINIMA BY FOLLOWING  
THE STEEPEST DESCENT

THE CHOICES OF THE INITIAL WEIGHTS AND THE  
LEARNING RATES ARE IMPORTANT



WE WILL TALK ABOUT  
THIS LATER

# LEARNING RATES



Credit:

# LEARNING RATES

$$W_{t+1} = W_t - \lambda \nabla f(W_t)$$



THERE ARE DIFFERENT WAYS  
TO UPDATE THE LEARNING RATE

Credit:

# LEARNING RATES

$$W_{t+1} = W_t - \lambda \nabla f(W_t)$$



THERE ARE DIFFERENT WAYS  
TO UPDATE THE LEARNING RATE

## ADAGRAD:

THE LEARNING RATE IS SCALED DEPENDING ON THE HISTORY OF PREVIOUS GRADIENTS

$$W_{t+1} = W_t - \frac{\lambda}{\sqrt{G_t + \epsilon}} \nabla f(W_t)$$

G IS A MATRIX CONTAINING ALL PREVIOUS GRADIENTS. WHEN THE GRADIENT BECOMES LARGE THE LEARNING RATE IS DECREASED AND VICE VERSA.

$$G_{t+1} = G_t + (\nabla f)^2$$

Credit:

# LEARNING RATES

$$W_{t+1} = W_t - \lambda \nabla f(W_t)$$



THERE ARE DIFFERENT WAYS  
TO UPDATE THE LEARNING RATE

## RMSPROP:

THE LEARNING RATE IS SCALED DEPENDING ON THE HISTORY OF PREVIOUS GRADIENTS

$$W_{t+1} = W_t - \frac{\lambda}{\sqrt{G_t + \epsilon}} \nabla f(W_t)$$

SAME AS ADAGRAD BUT G IS CALCULATED BY EXPONENTIALLY DECAYING AVERAGE

$$G_{t+1} = \lambda G_t + (1 - \lambda)(\nabla f)^2$$

Credit:

## **ADAM [Adaptive moment estimator]:**

SAME IDEA, USING FIRST AND SECOND ORDER  
MOMENTUMS

$$G_{t+1} = \beta_2 G_t + (1 - \beta_2)(\nabla f)^2 \quad M_{t+1} = \beta_1 M_t + (1 - \beta_1)(\nabla f)$$

$$W_{t+1} = W_t - \frac{\lambda}{\sqrt{\hat{G}_t + \epsilon}} \hat{M}_t$$

with:  $\hat{M}_{t+1} = \frac{M_t}{1 - \beta_1}$        $\hat{G}_{t+1} = \frac{G_t}{1 - \beta_2}$

## **ADAM [Adaptive moment estimator]:**

SAME IDEA, USING FIRST AND SECOND ORDER  
MOMENTUMS

$$G_{t+1} = \beta_2 G_t + (1 - \beta_2)(\nabla f)^2 \quad M_{t+1} = \beta_1 M_t + (1 - \beta_1)(\nabla f)$$

ONLY FOR YOUR  
RECORDS

$$\sqrt{G_t} + \epsilon$$

with:  $\hat{M}_{t+1} = \frac{M_t}{1 - \beta_1}$        $\hat{G}_{t+1} = \frac{G_t}{1 - \beta_2}$

# IN KERAS:

## RMSprop

[source]

```
keras.optimizers.RMSprop(lr=0.001, rho=0.9, epsilon=None, decay=0.0)
```

RMSProp optimizer.

It is recommended to leave the parameters of this optimizer at their default values (except the learning rate, which can be freely tuned).

This optimizer is usually a good choice for recurrent neural networks.

### Arguments

- `lr`: float  $\geq 0$ . Learning rate.
- `rho`: float  $\geq 0$ .
- `epsilon`: float  $\geq 0$ . Fuzz factor. If `None`, defaults to `K.epsilon()`.
- `decay`: float  $\geq 0$ . Learning rate decay over each update.

### References

- [rmsprop](#): Divide the gradient by a running average of its recent magnitude

# IN KERAS:

## Adam

[source]

```
keras.optimizers.Adam(lr=0.001, beta_1=0.9, beta_2=0.999, epsilon=None, decay=0.0, amsgrad=False)
```

Adam optimizer.

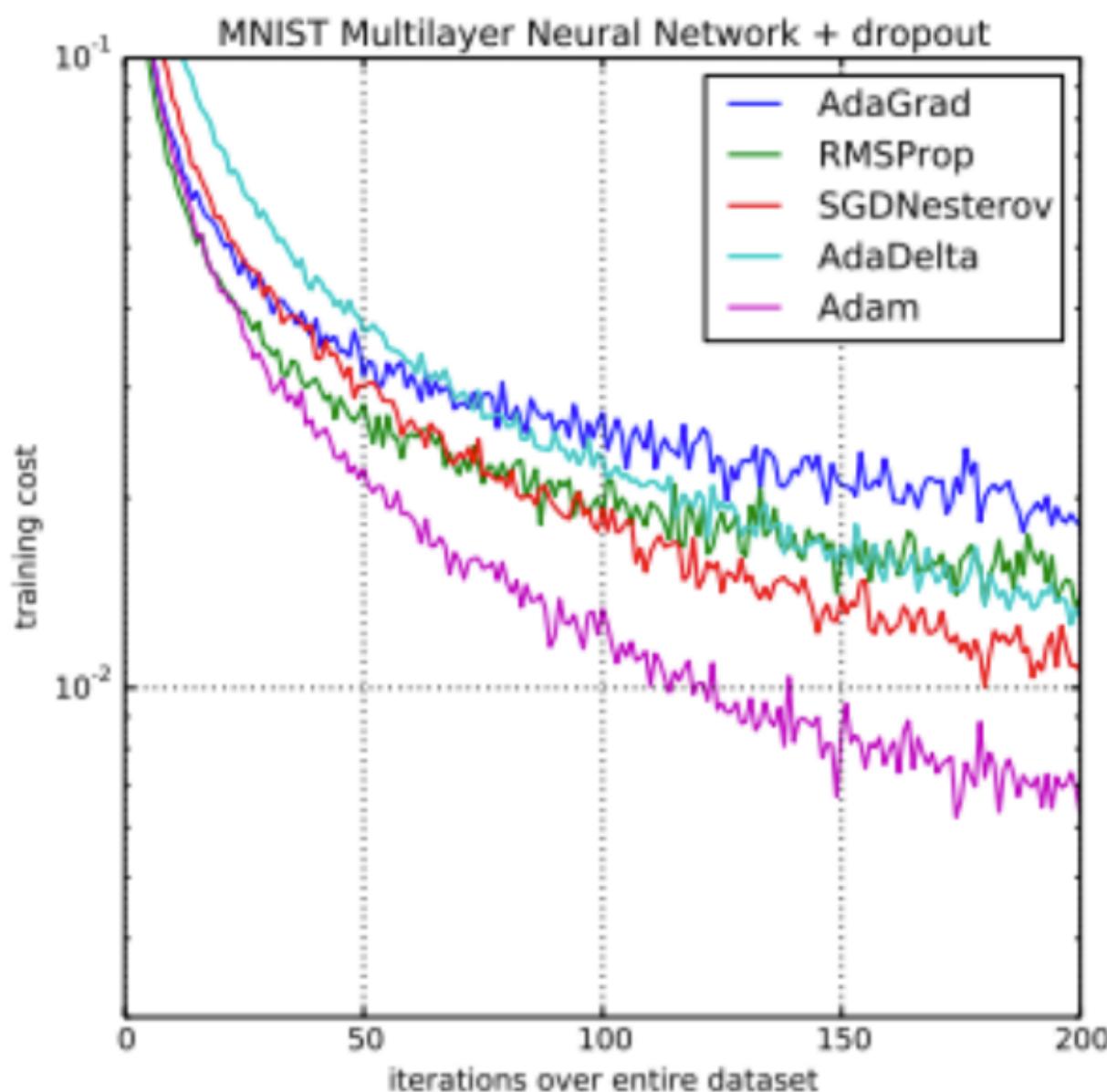
Default parameters follow those provided in the original paper.

### Arguments

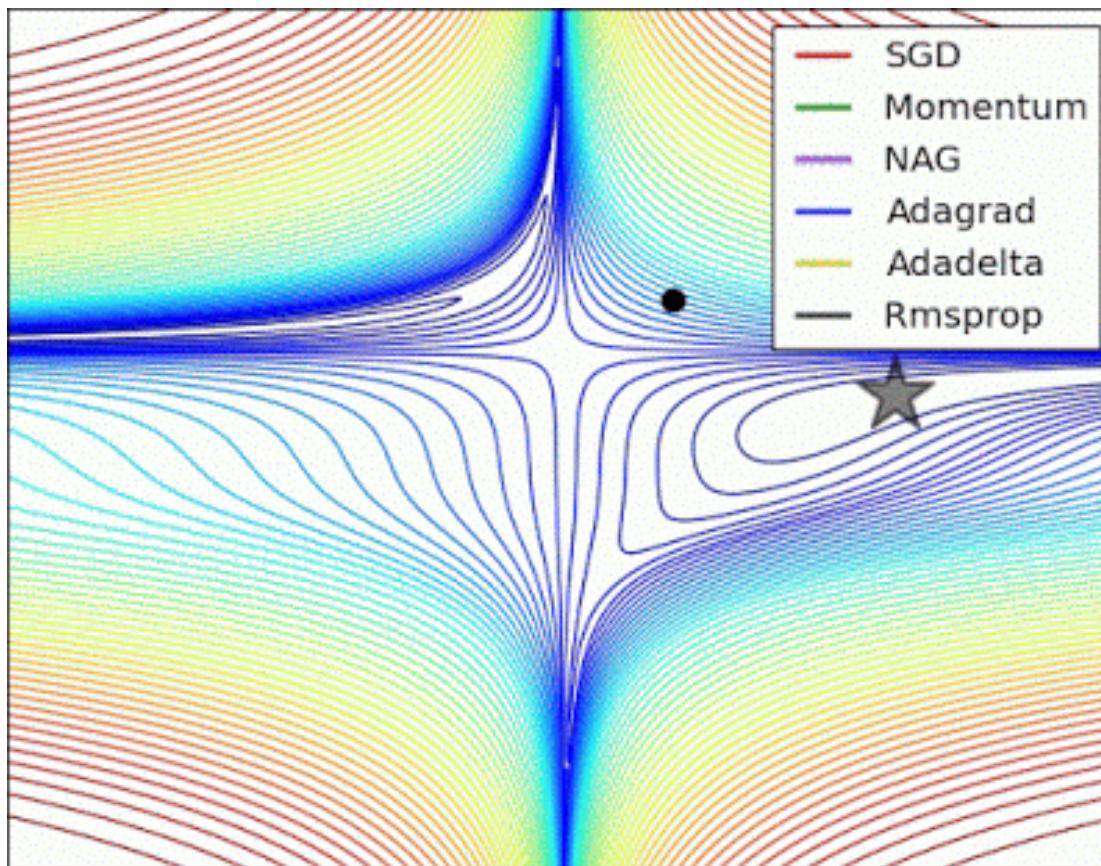
- **lr**: float  $\geq 0$ . Learning rate.
- **beta\_1**: float,  $0 < \beta_1 < 1$ . Generally close to 1.
- **beta\_2**: float,  $0 < \beta_2 < 1$ . Generally close to 1.
- **epsilon**: float  $\geq 0$ . Fuzz factor. If `None`, defaults to `K.epsilon()`.
- **decay**: float  $\geq 0$ . Learning rate decay over each update.
- **amsgrad**: boolean. Whether to apply the AMSGrad variant of this algorithm from the paper "On the Convergence of Adam and Beyond".

### References

- [Adam - A Method for Stochastic Optimization](#)
- [On the Convergence of Adam and Beyond](#)



Credit



Credit

# BATCH GRADIENT DESCENT

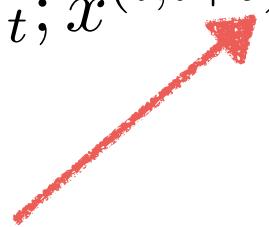
LOCAL MINIMA CAN ALSO BE AVOIDED BY COMPUTING THE  
GRADIENT IN SMALL BATCHES INSTEAD OF OVER THE FULL  
DATASET

# BATCH GRADIENT DESCENT

LOCAL MINIMA CAN ALSO BE AVOIDED BY COMPUTING THE GRADIENT IN SMALL BATCHES INSTEAD OF OVER THE FULL DATASET

## MINI-BATCH GRADIENT DESCENT

$$V_{t+1/num} = W_t - \lambda_t \nabla f(W_t; x^{(i,i+b)}, y^{(i,i+b)})$$



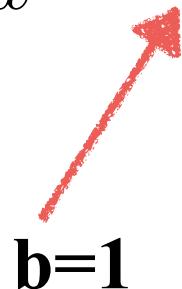
THE GRADIENT IS COMPUTED OVER A BATCH OF SIZE B

# STOCHASTIC GRADIENT DESCENT

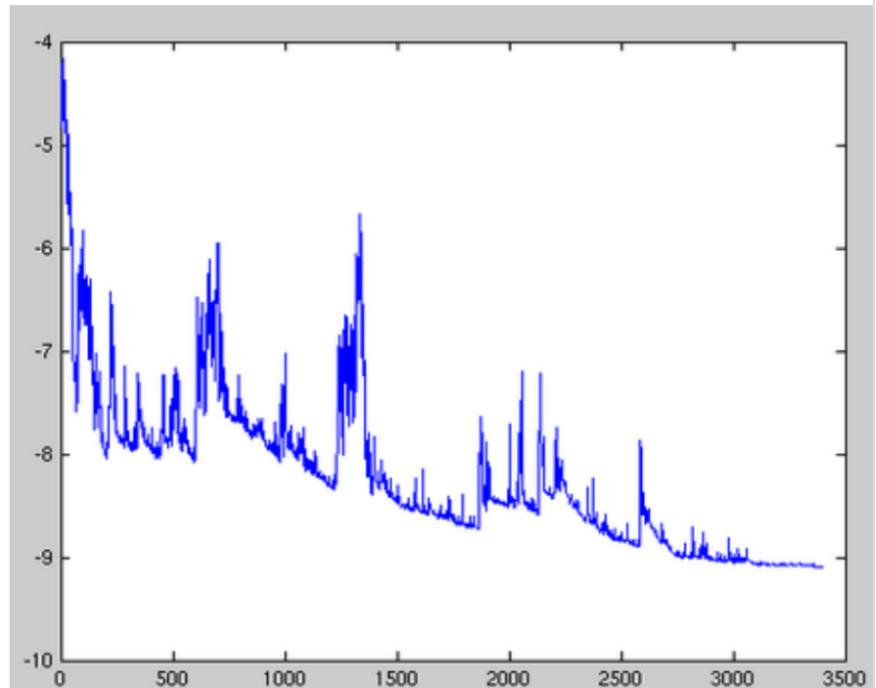
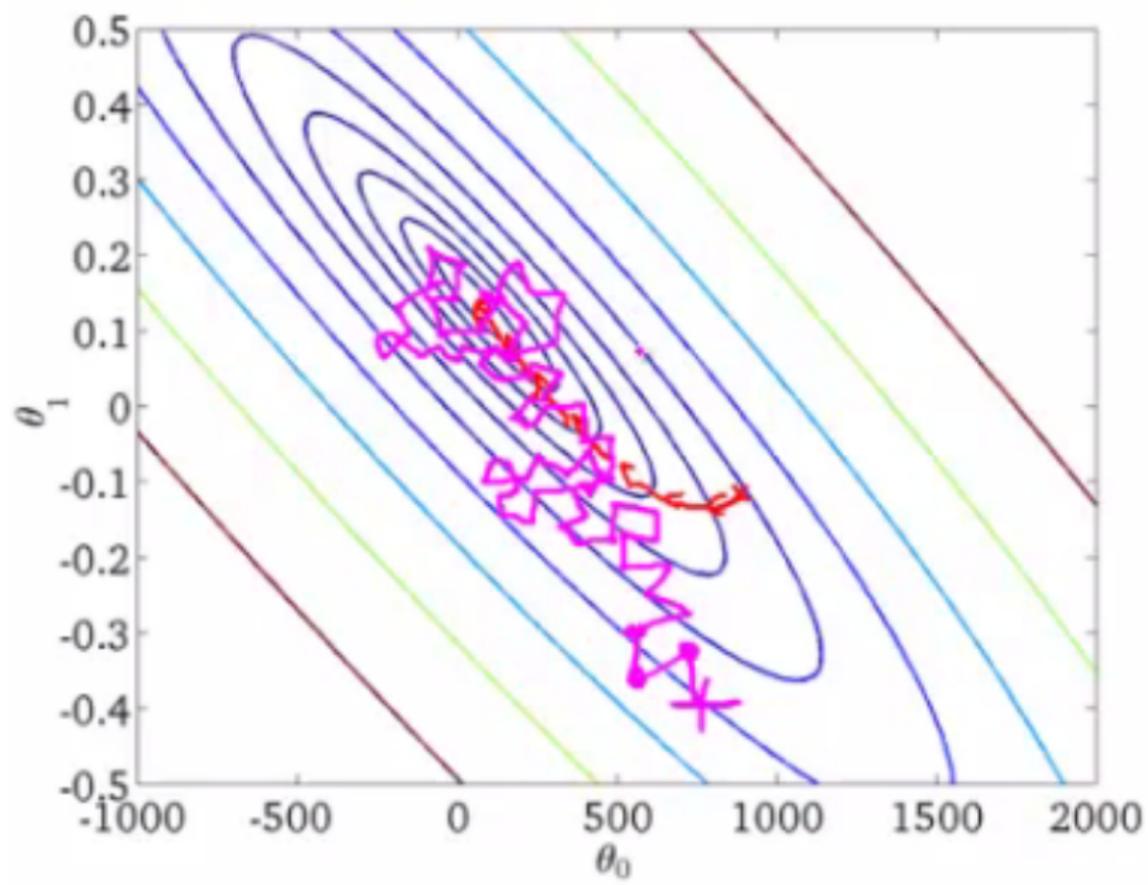
THE EXTREME CASE IS TO COMPUTE THE GRADIENT ON EVERY TRAINING EXAMPLE.

## STOCHASTIC GRADIENT DESCENT

$$V_{t+1/num} = W_t - \lambda_t \nabla f(W_t; x^{(i,i+b)}, y^{(i,i+b)})$$



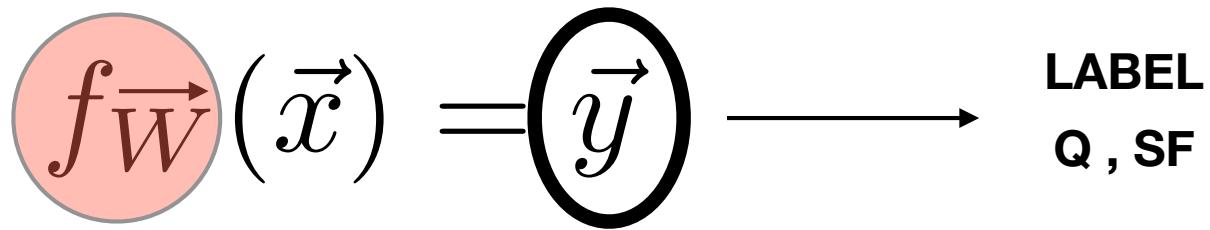
b=1



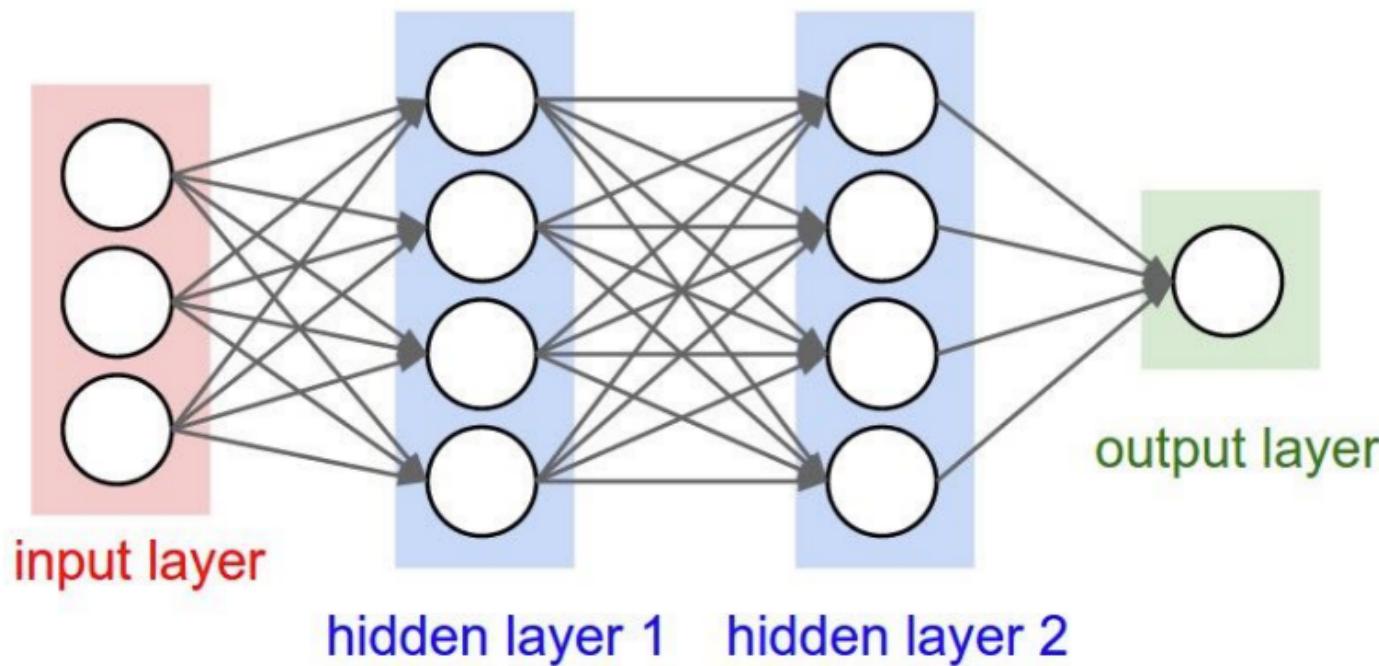
Fluctuations in the total objective function as gradient steps with respect to mini-batches are taken.

Credit

**"CLASSICAL"  
MACHINE LEARNING**



**REPLACE THIS BY A GENERAL  
NON LINEAR FUNCTION WITH SOME PARAMETERS W**



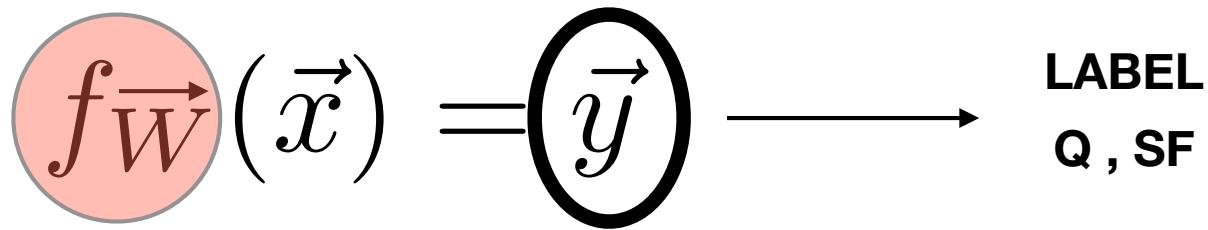
$$p = g_3(W_3g_2(W_2g_1(W_1\vec{x}_0)))$$

NETWORK  
FUNCTION

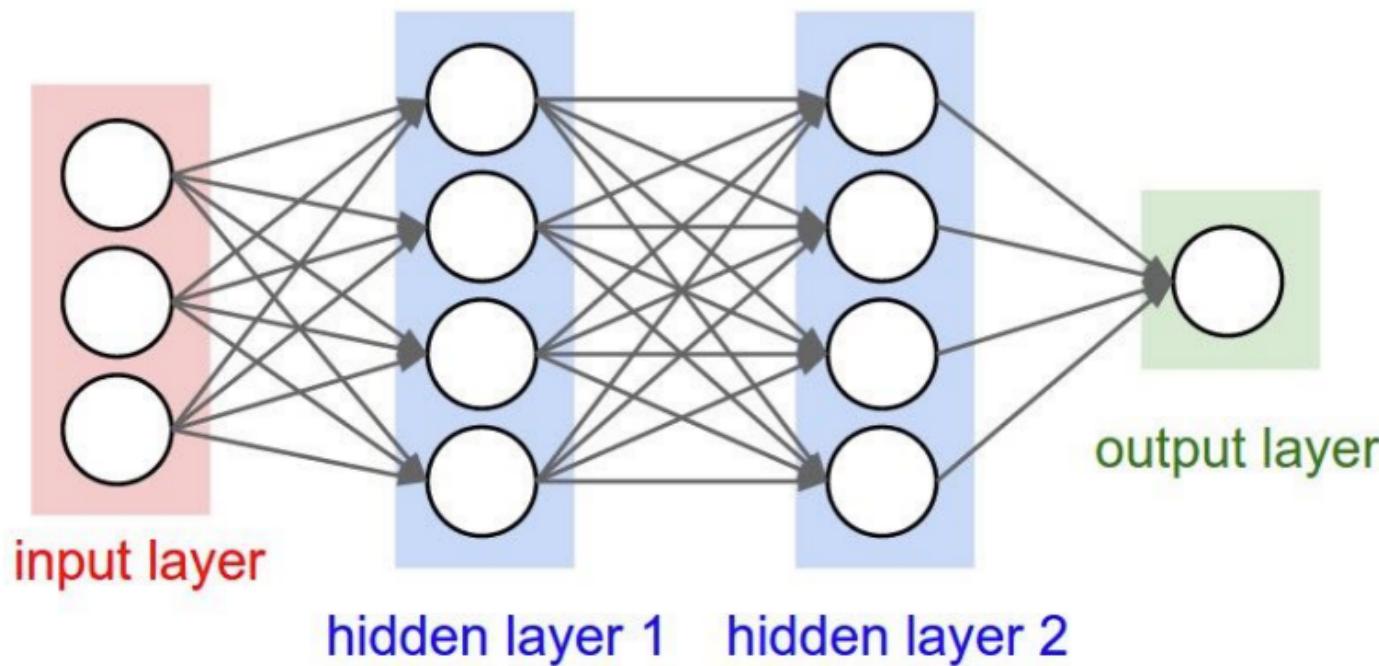
# THE PROBLEMS OF GOING “TOO DEEP”

- DEEP NETWORKS ARE MORE DIFFICULT TO OPTIMIZE
- NEED MORE DATA - MORE SUBJECT TO OVERFITTING
- AND ALSO NEED MORE TIME ...

**"CLASSICAL"  
MACHINE LEARNING**



**REPLACE THIS BY A GENERAL  
NON LINEAR FUNCTION WITH SOME PARAMETERS W**



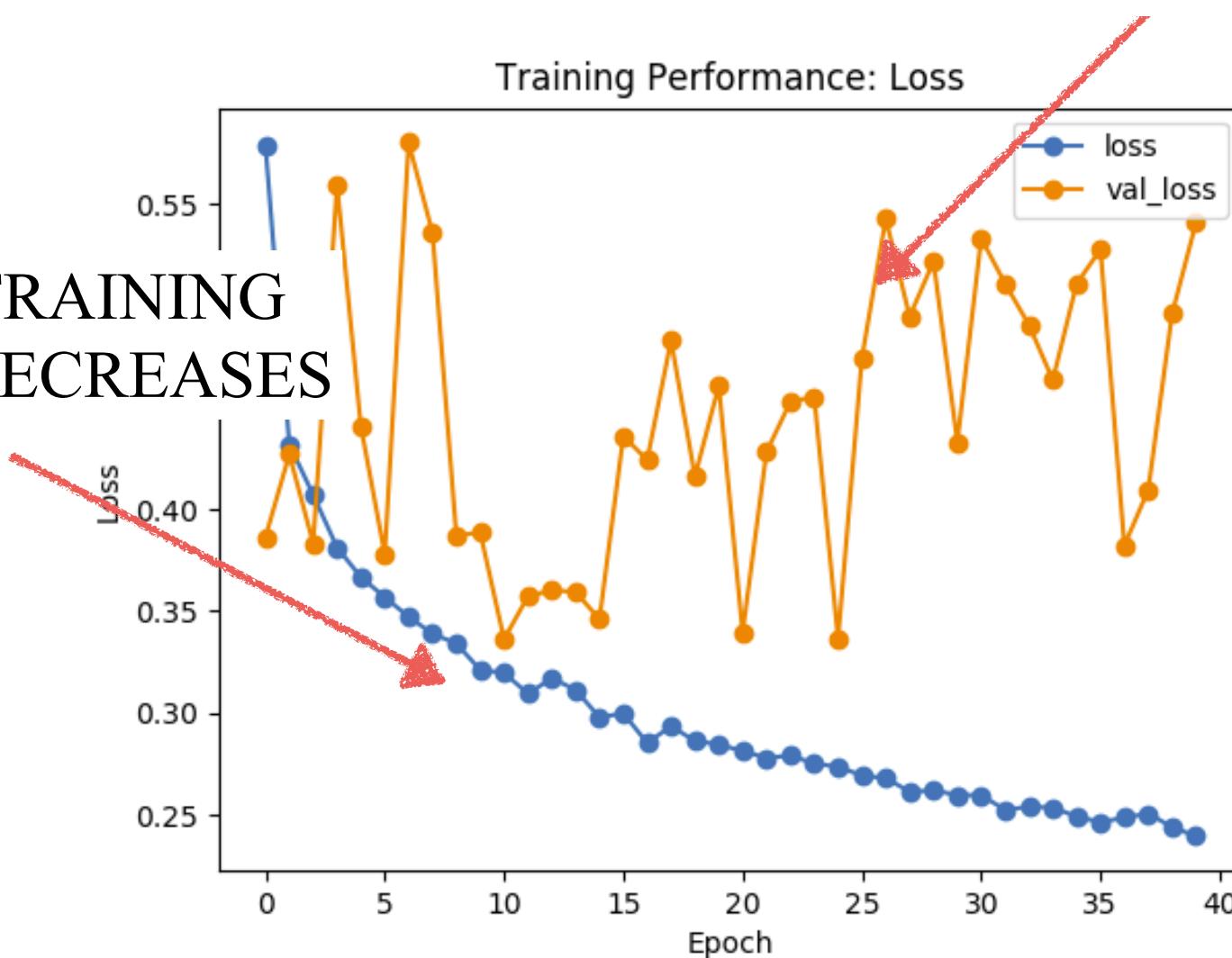
$$p = g_3(W_3g_2(W_2g_1(W_1\vec{x}_0)))$$

← NETWORK  
FUNCTION

# OVER-FITTING

THE TEST STAYS CONSTANT  
OR INCREASES

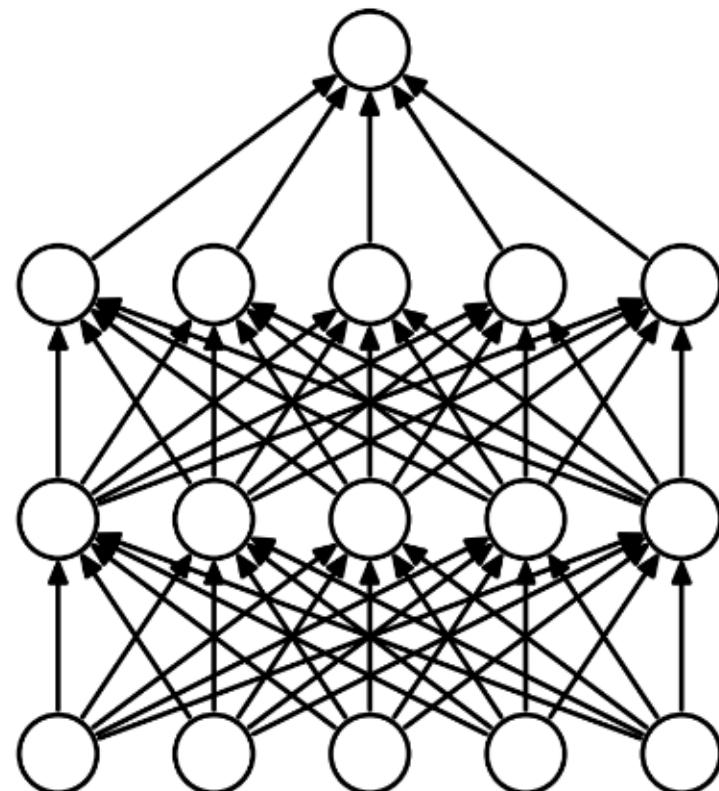
THE TRAINING  
LOSS DECREASES



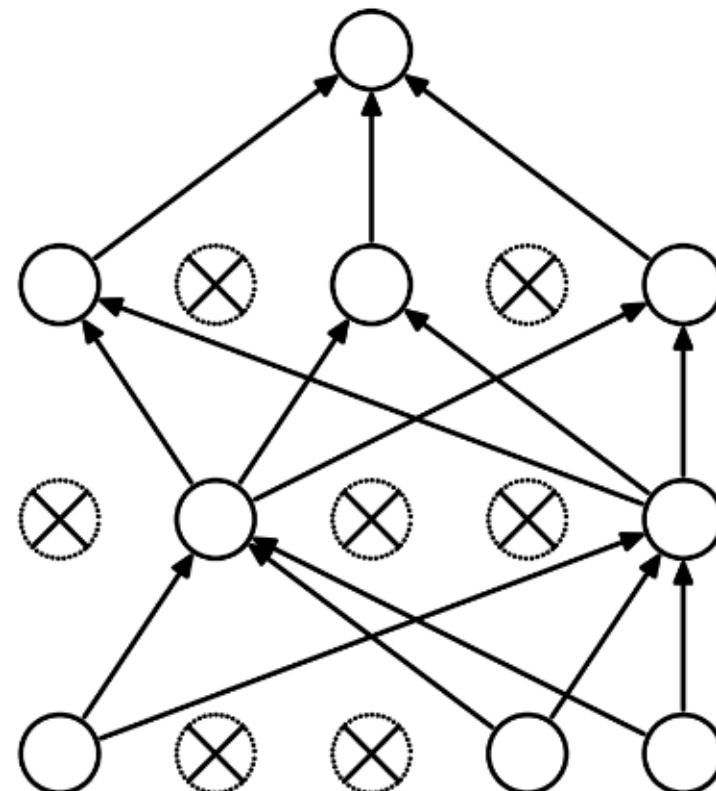
# DROPOUT

[Hinton+12]

- THE IDEA IS TO REMOVE NEURONS RANDOMLY DURING THE TRAINING
- ALL NEURONS ARE PUT BACK DURING THE TEST PHASE



(a) Standard Neural Net



(b) After applying dropout.

# DROPOUT

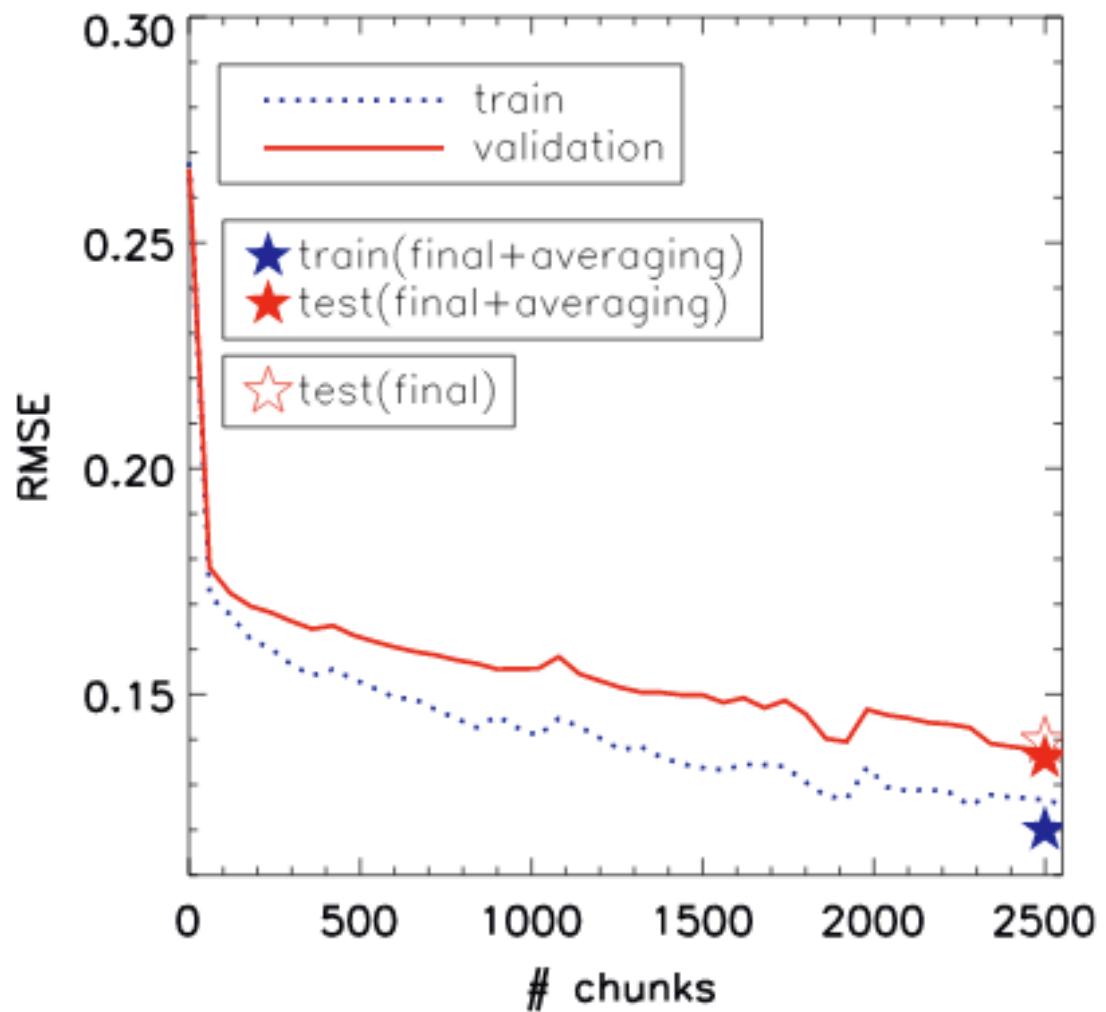
## WHY DOES IT WORK?

1. SINCE NEURONS ARE REMOVED RANDOMLY, IT AVOIDS CO-ADAPTATION AMONG THEMSELVES

2. DIFFERENT SETS OF NEURONS WHICH ARE SWITCHED OFF, REPRESENT A DIFFERENT ARCHITECTURE AND ALL THESE DIFFERENT ARCHITECTURES ARE TRAINED IN PARALLEL. FOR  $N$  NEURONS ATTACHED TO DROPOUT, THE NUMBER OF SUBSET ARCHITECTURES FORMED IS  $2^N$ . SO IT AMOUNTS TO PREDICTION BEING AVERAGED OVER THESE ENSEMBLES OF MODELS.

# DROPOUT

WITH A LITTLE BIT  
OF DROPOUT



**GO HERE:**

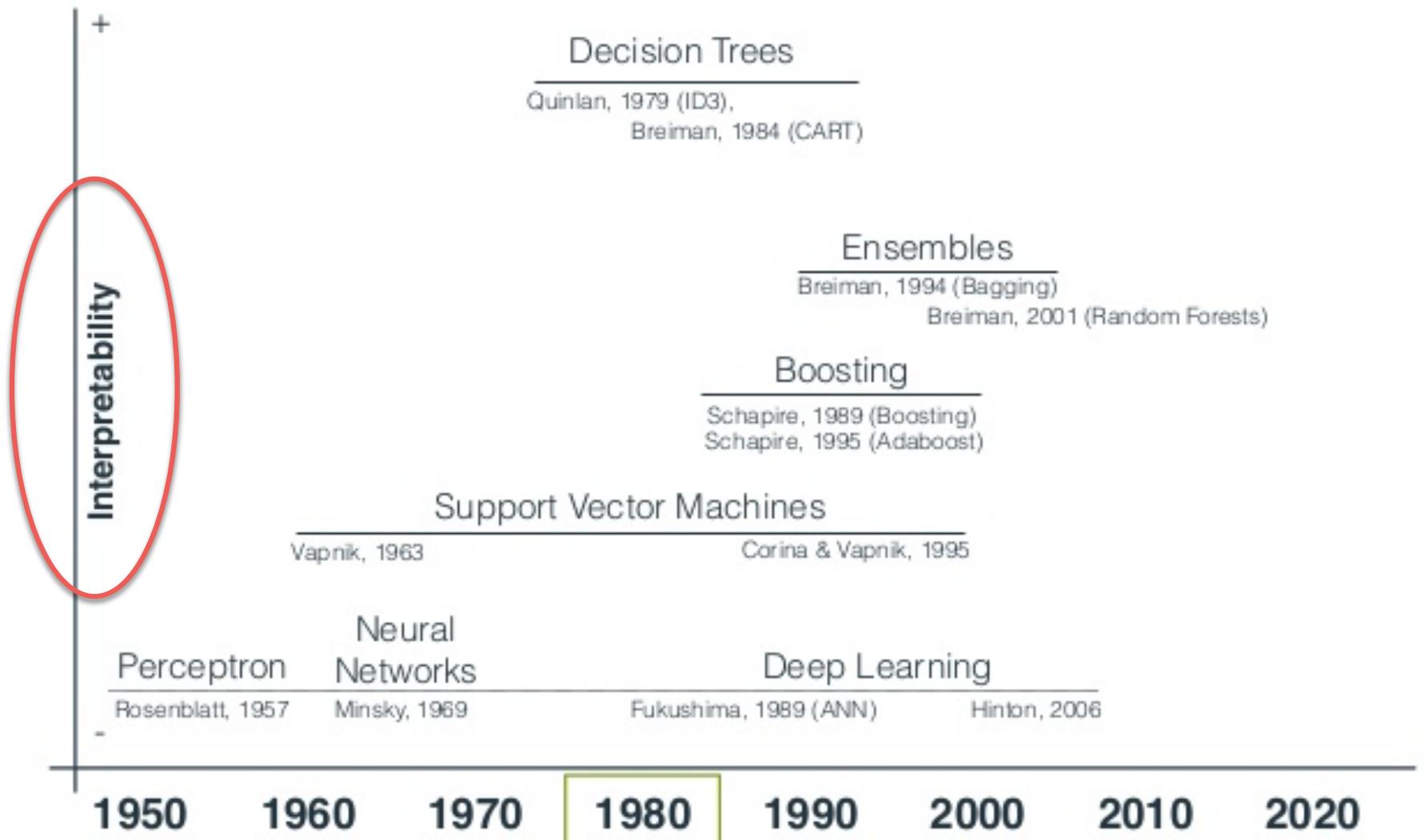
[https://github.com/mhuertascompany/ULL\\_2021/blob/main/tutorials/day2/Galaxy\\_Morphology ANN ULL.ipynb](https://github.com/mhuertascompany/ULL_2021/blob/main/tutorials/day2/Galaxy_Morphology ANN ULL.ipynb)

# HOW TO CHOOSE YOUR CLASSICAL CLASSIFIER?

**NO RULE OF THUMB - REALLY DEPENDS ON APPLICATION**

ML METHOD	++	-	Python
CARTS / RANDOM FOREST	Easy to interpret (“White box”) Litte data preparation Both numerical + categorical	Over-complex trees Unstable Biased trees if some classes dominate	sklearn.ensemble.RandomForestClassifier sklearn.ensemble.RandomForestRegressor
SVM	Easy to interpret + Fast Kernel trick allows no linear problems	not very well suited to multi-class problems	sklearn.svm sklearn.svc
NN	seed of deep-learning very efficient with large amount of data as we will see	more difficult to interpret computing intensive	sklearn.neural_network.MLPClassifier sklearn.neural_network.MLPRegressor

# HOW TO CHOOSE YOUR CLASSICAL CLASSIFIER?



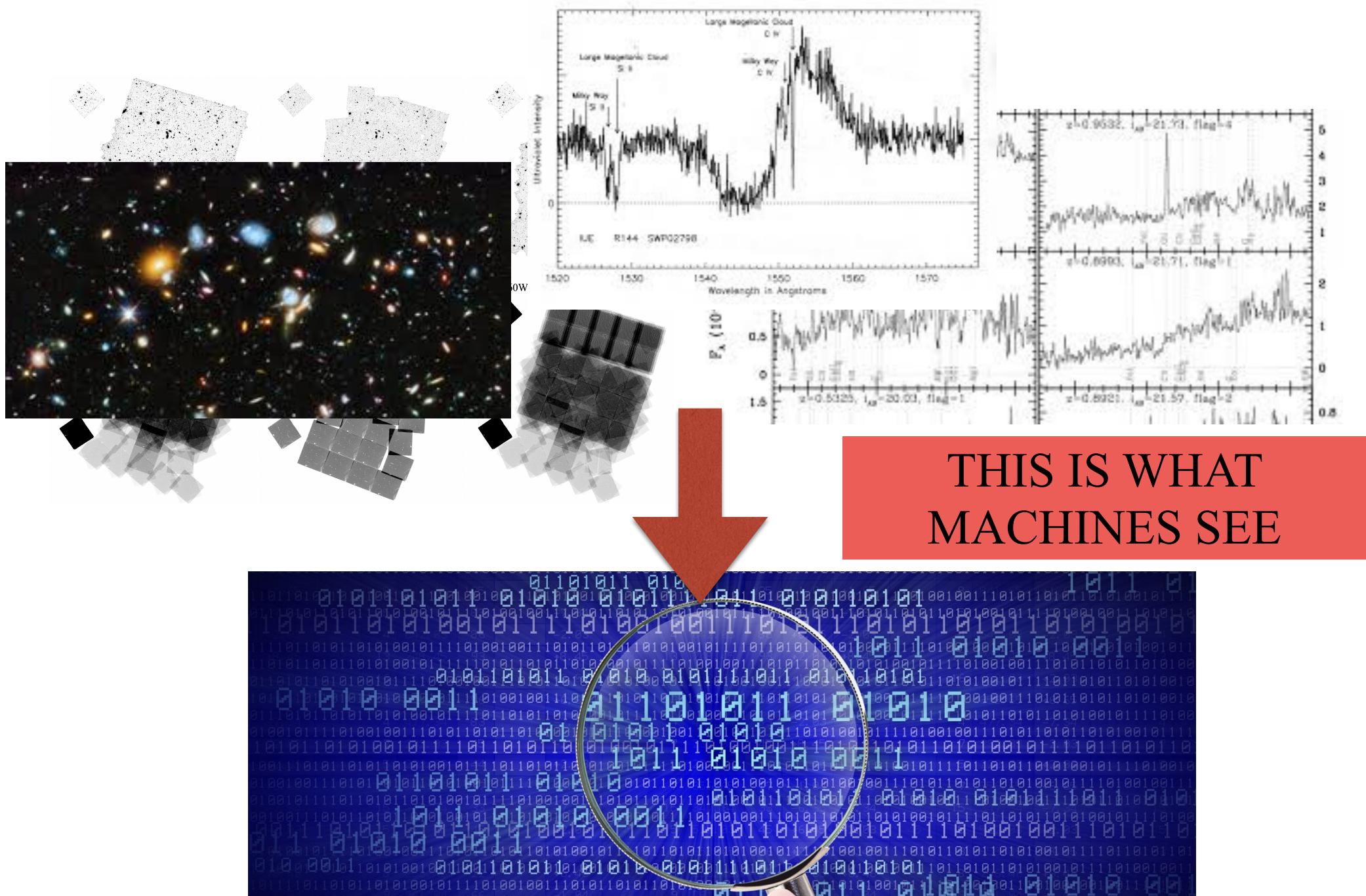
credit

CAN WE GO DEEP NOW?

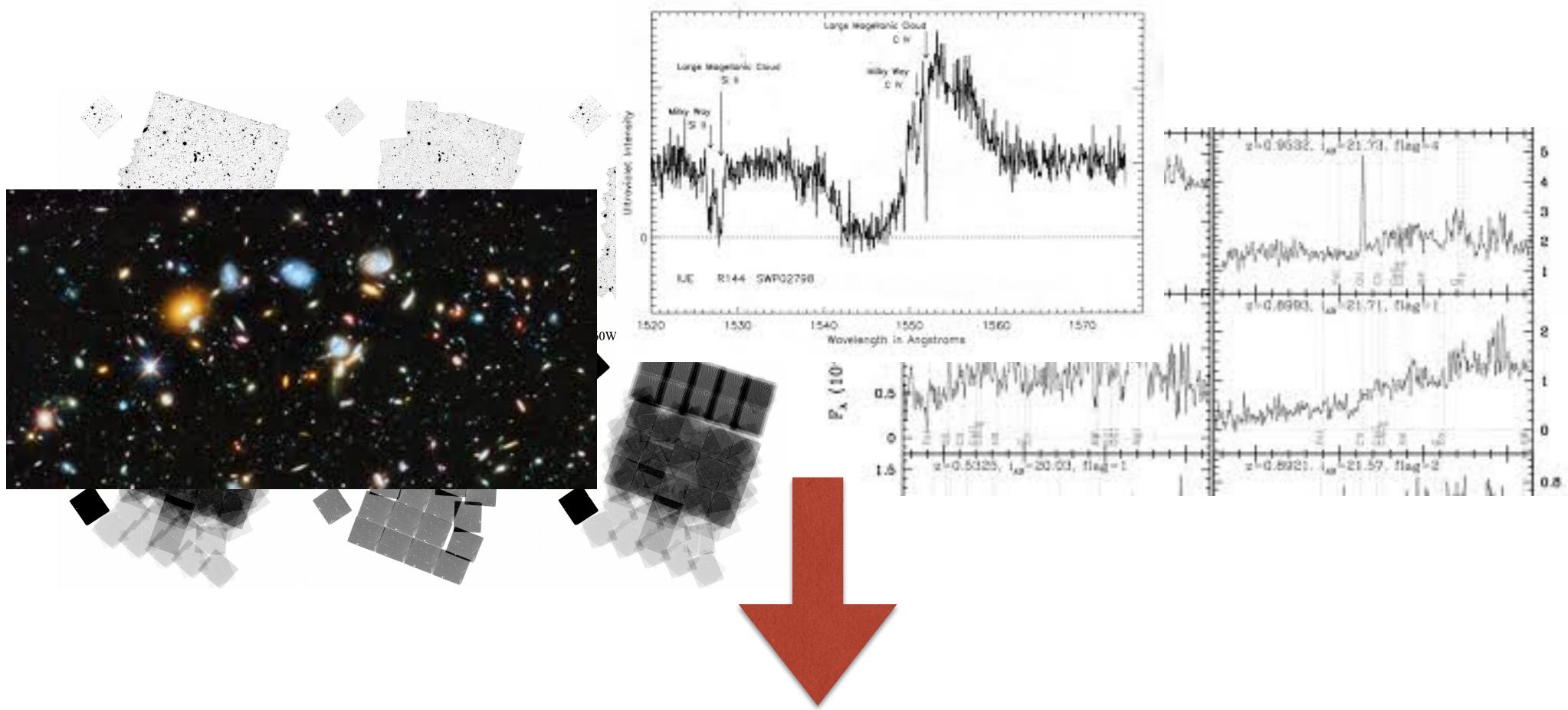
CAN WE GO DEEP NOW?

ALMOST THERE...LET'S THINK FOR A  
MOMENT ABOUT WHAT WE PUT AS  
INPUT...

# What do we put as input?

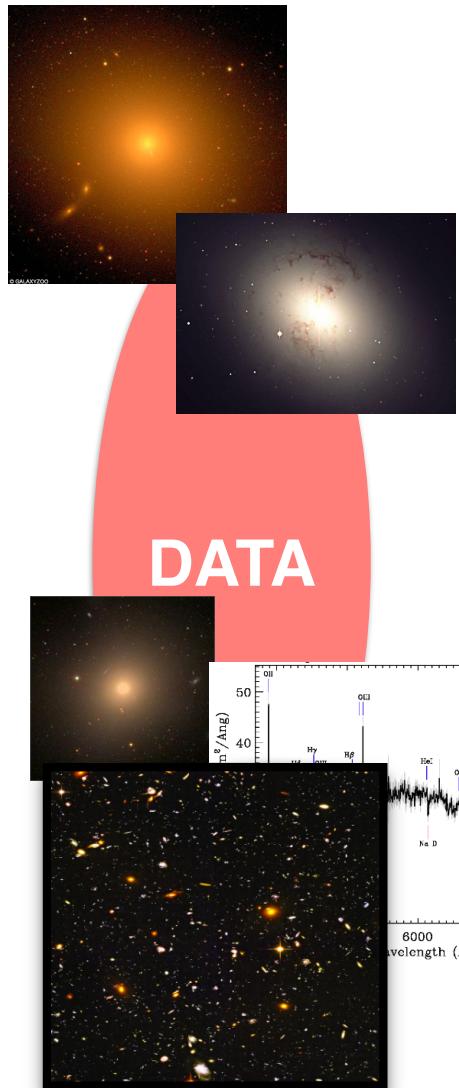


# What do we put as input?



PRE-PROCESS DATA TO EXTRACT MEANINGFUL INFORMATION

THIS IS GENERALLY CALLED **FEATURE EXTRACTION**

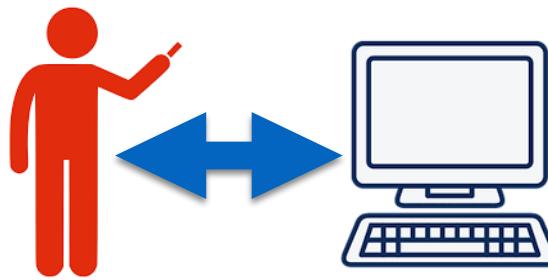
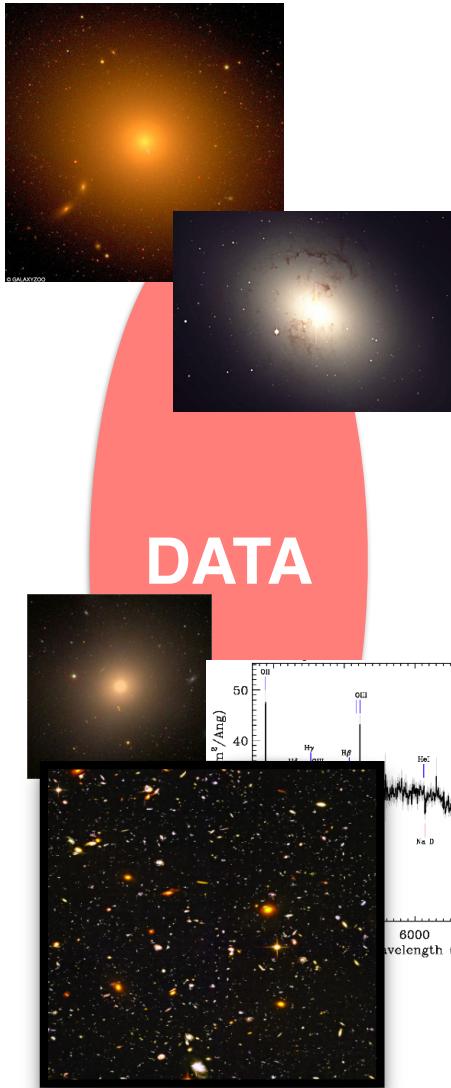


Spiral!

Emission line!

Merger!

Clump!  
AGN!



**Spiral!**

**Emission line!**

**Merger!**

**Clump!**

**AGN!**

$$f_W(\vec{x}) = \vec{y} \longrightarrow \text{LABEL}$$

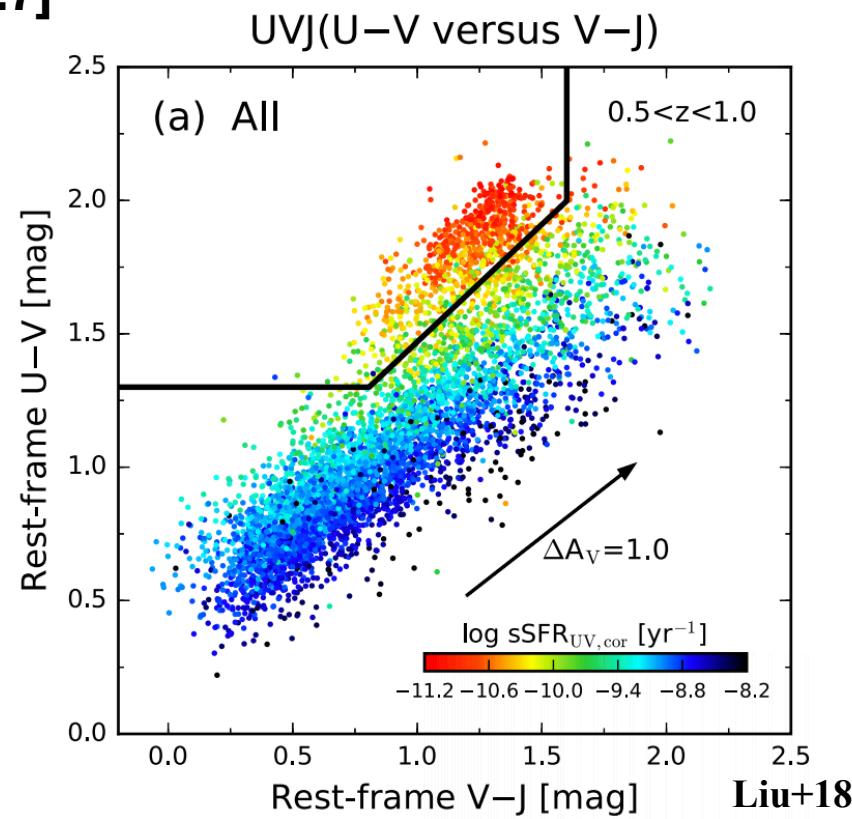
**Q(0) , SF(1)**

**NETWORK FUNCTION**

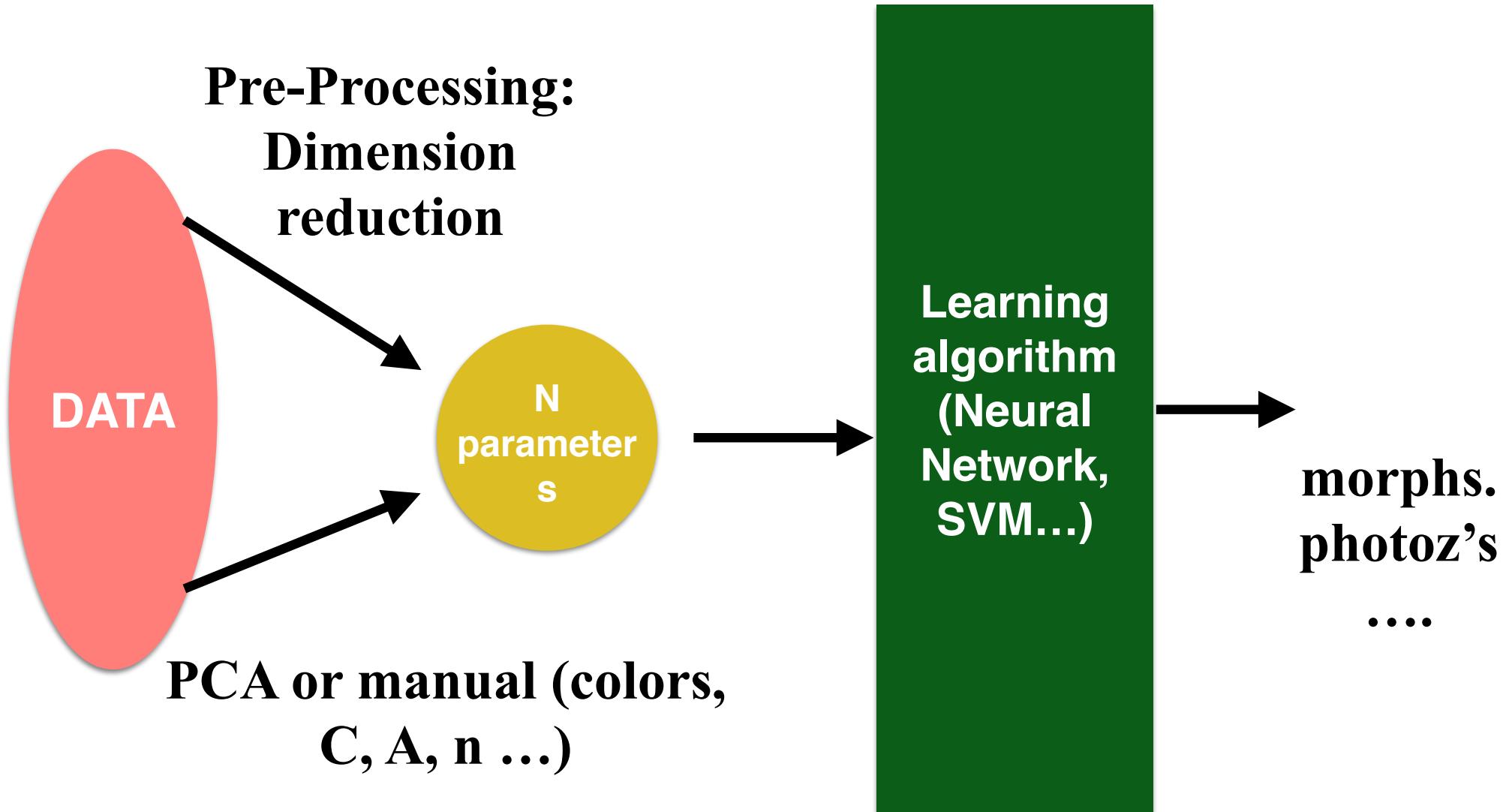
**(U-V, V-J) FEATURES**

$$\text{sgn}[(u-v)-0.8*(v-j)-0.7]$$

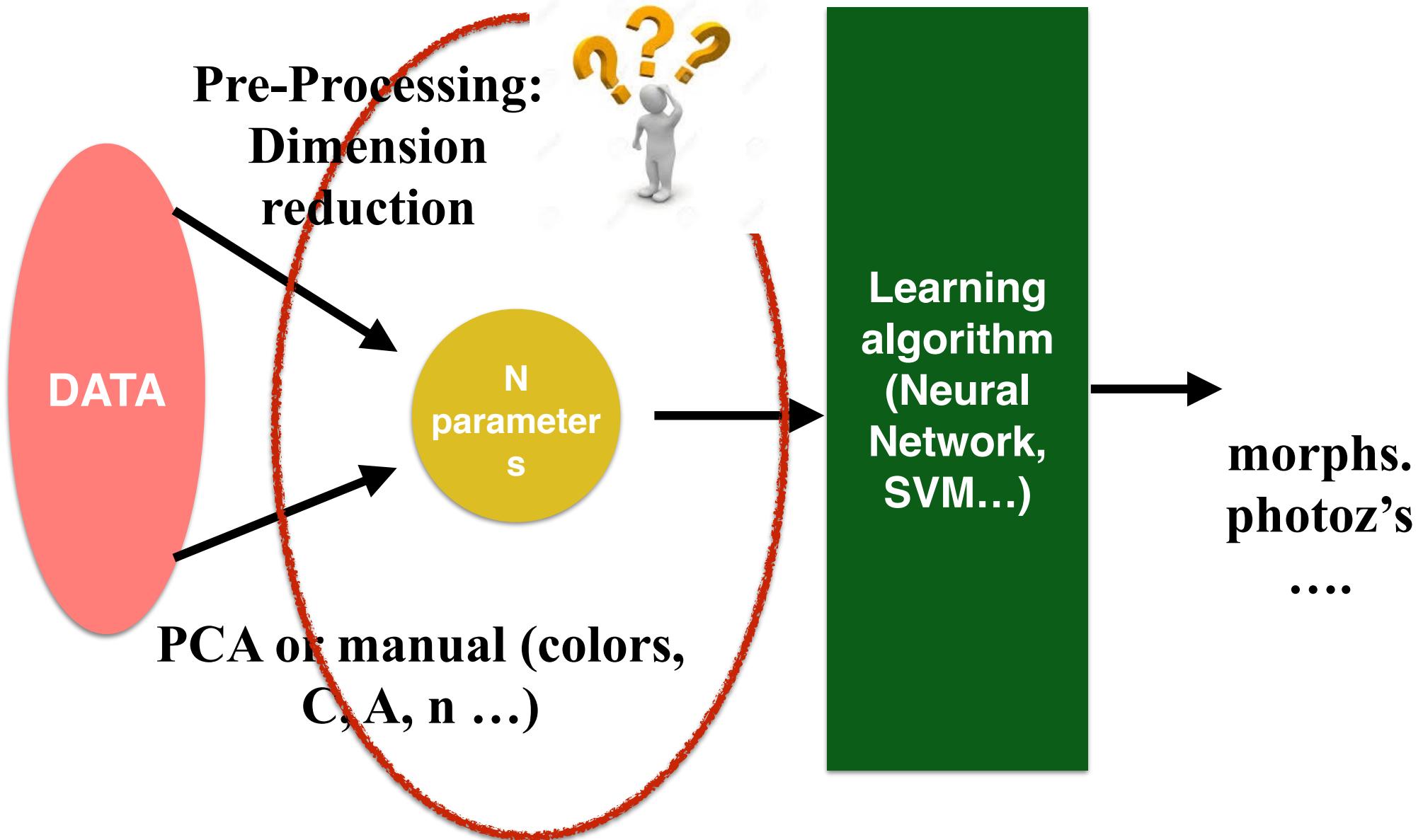
**WEIGHTS**



# THE “CLASSICAL” APPROACH

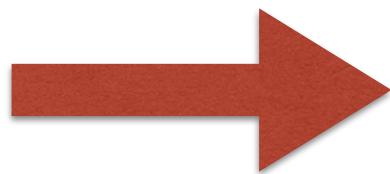


# “CLASSICAL” MACHINE LEARNING



# In Astronomy

- Colors, Fluxes
- Shape indicators
- Line ratios, spectral features
- Stellar Masses, Velocity Dispersions



Requires specialized software before feeding the machine learning algorithm

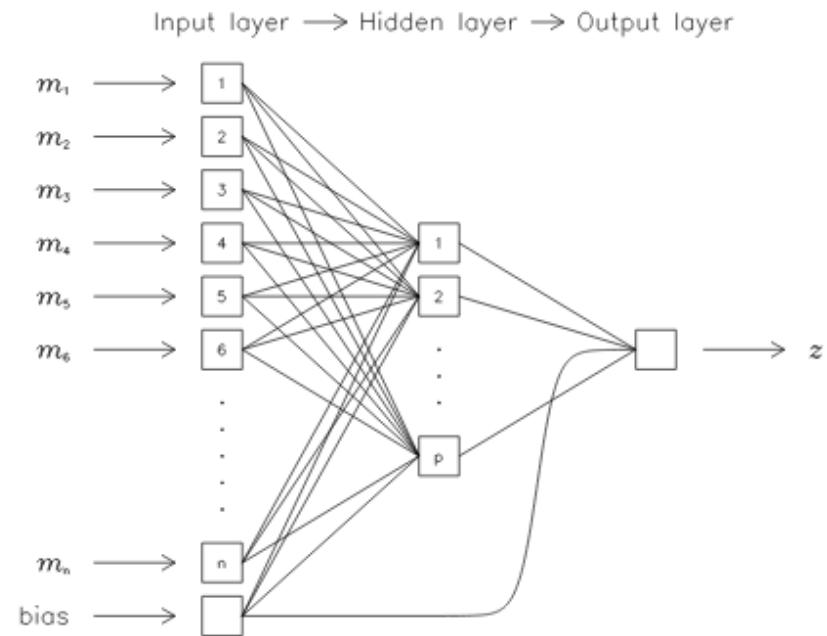
**IT IMPLIES A DIMENSIONALITY REDUCTION!**

# PHOTOMETRIC REDSHIFTS

SDSS

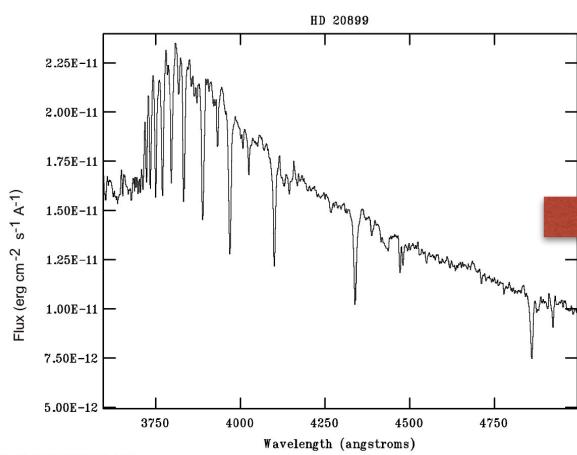


g  
r  
i  
z

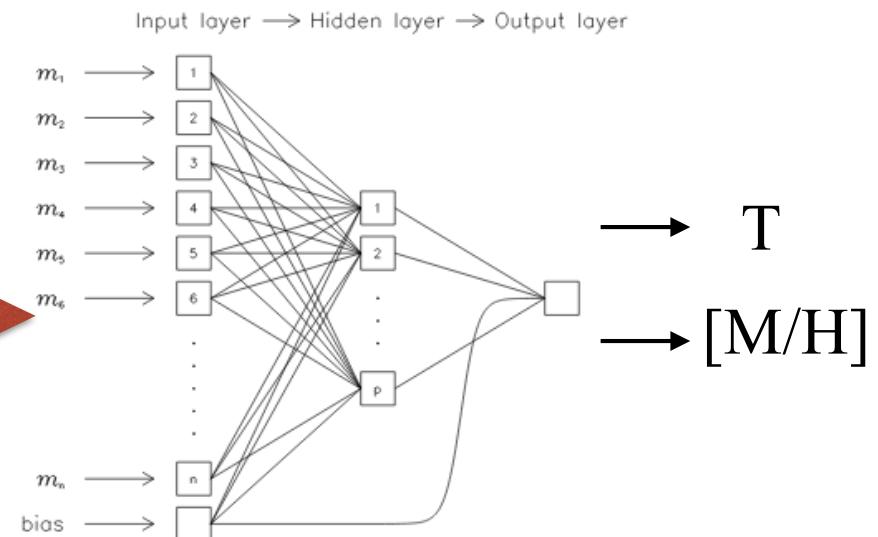


Collister+08

# STELLAR PARAMETERS FROM MEDIUM BAND FILTERS



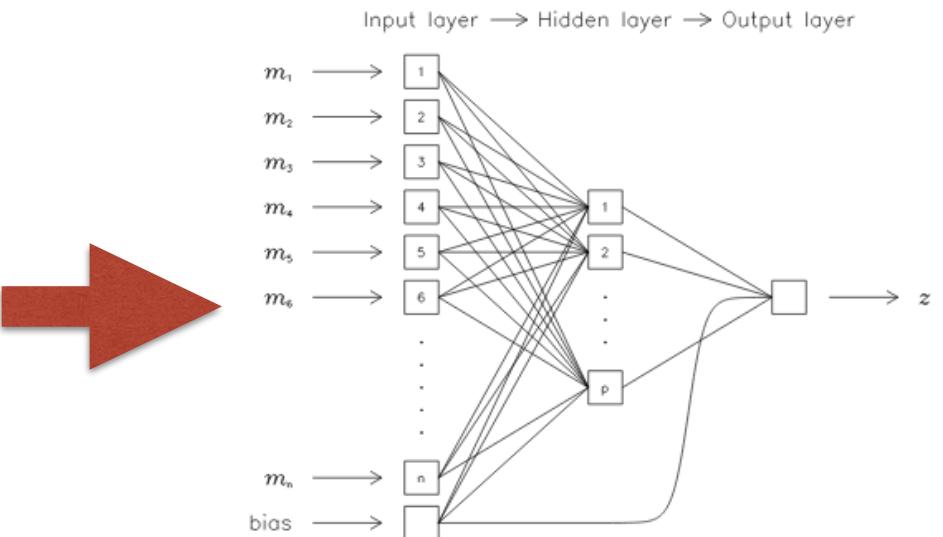
MEDIUM  
BAND  
FLUXES



Bailer-Jones+00

No.	Symbol	Description	Scale <sup>a</sup>
1	CVD	Central velocity dispersion	~1 kpc
2	$M_{\text{bulge}}$	Bulge stellar mass	0.5–4 kpc
3	$R_e$	Bulge effective radius	0.5–4 Kpc
4	B/T	Bulge-to-total stellar mass ratio	0.5–8 kpc
5	$M_*$	Total stellar mass	2–8 kpc
6	$M_{\text{disc}}$	Disc stellar mass	4–10 kpc
7	$M_{\text{halo}}$	Group halo mass	0.1–1 Mpc
8	$\delta_5$	Local density parameter	0.5–3 Mpc

Notes. <sup>a</sup> Approximate  $1\sigma$  range from centre of galaxy. For photometric quantities half-light radii are used.



## HEAVILY PROCESSED DATA

# Other general computer vision features [for images!]

- Pixel Concatenation
- Color histograms
- Texture Features
- Histogram of Gradients
- SIFT

FOR MANY YEARS COMPUTER VISION RESEARCHERS HAVE BEEN TRYING TO FIND THE MOST GENERAL FEATURES

# Other general computer vision features [for images!]

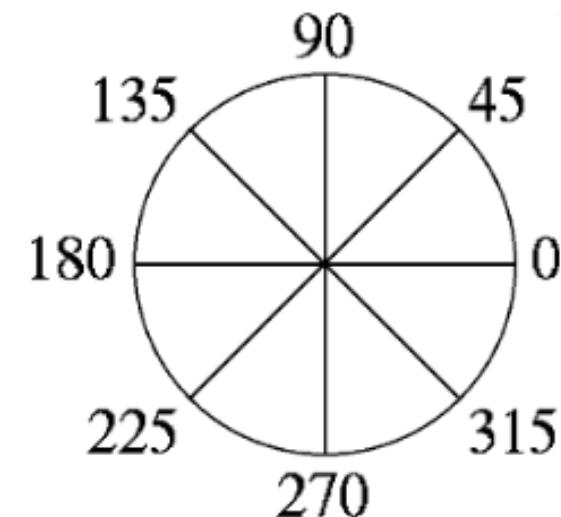
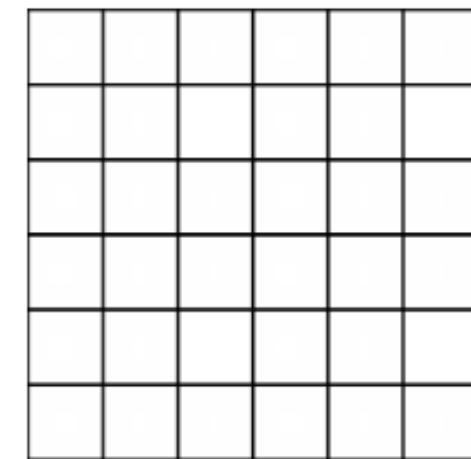
- Pixel Concatenation
- Color histograms
- Texture Features
- Histogram of Gradients
- SIFT

FOR MANY YEARS COMPUTER VISION RESEARCHERS HAVE BEEN TRYING TO FIND THE MOST GENERAL FEATURES

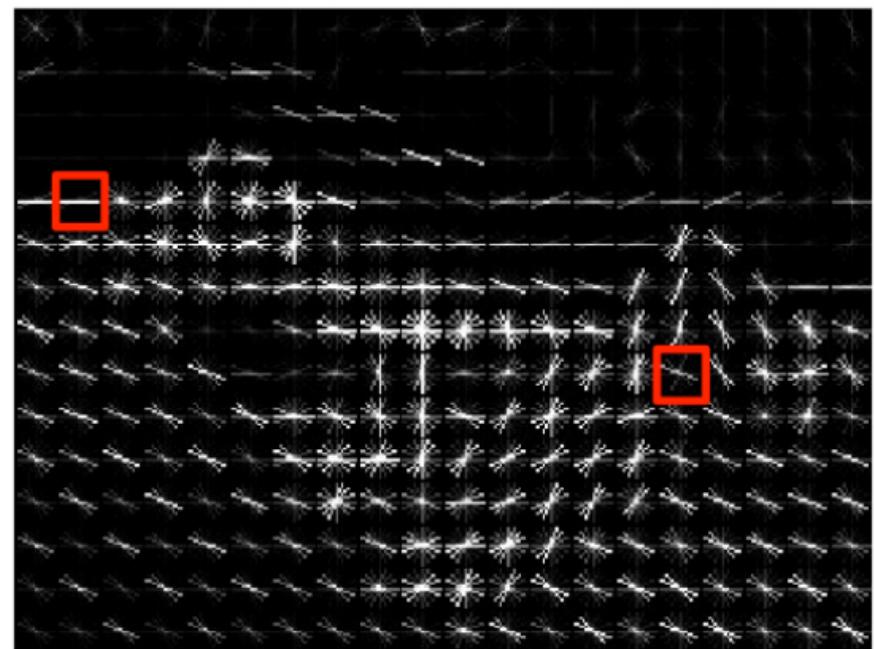
THE BEST CLASSICAL SOLUTION [BEFORE 2012] WHERE BASED ON LOCAL FEATURES

# HISTOGRAM OF ORIENTED GRADIENTS (HoG)

1. DIVIDE IMAGE INTO SMALL SPATIAL REGIONS CALLED CELLS
2. COMPUTE INTENSITY GRADIENTS OVER N DIRECTIONS [TYPICALLY 9 FOR IMAGE ]
3. COMPUTE WEIGHTED 1-D HISTOGRAM OF ALL DIRECTIONS. A CELL IS REDUCED TO N NUMBERS



# HISTOGRAM OF ORIENTED GRADIENTS (HoG)



**EVERYTHING IS IN THE FEATURES....WHAT IF I  
IGNORED SOME IMPORTANT FEATURES?**



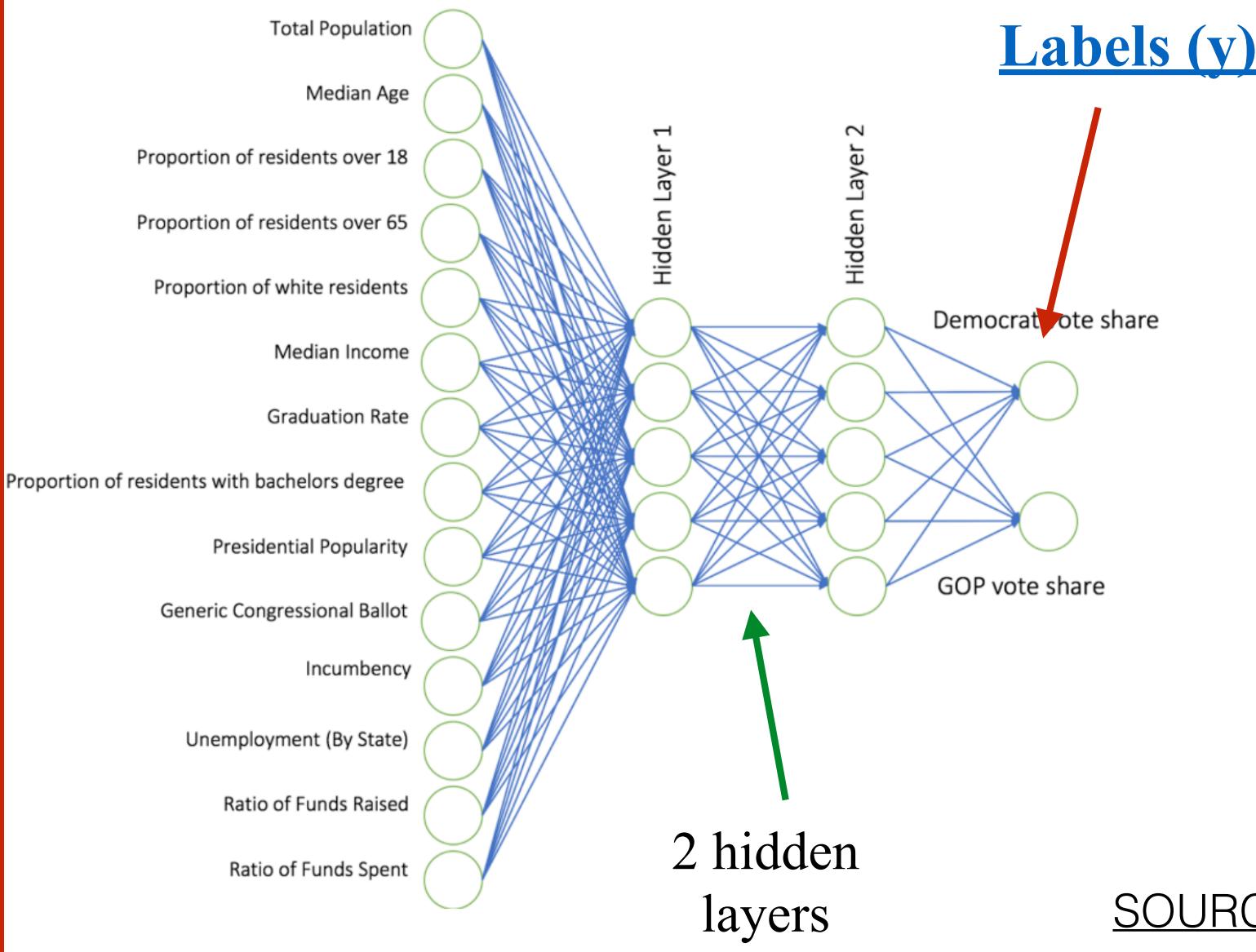
**EVERYTHING IS IN THE FEATURES...WHAT IF I  
IGNORED SOME IMPORTANT FEATURES?**



# NEURAL NETWORK TO PREDICT RESULTS OF MIDTERM ELECTIONS

$$p = g_3(W_3g_2(W_2g_1(W_1\vec{x}_0)))$$

Features  
(x)



SOURCE

# NEURAL NETWORK TO PREDICT RESULTS OF MIDTERM ELECTIONS

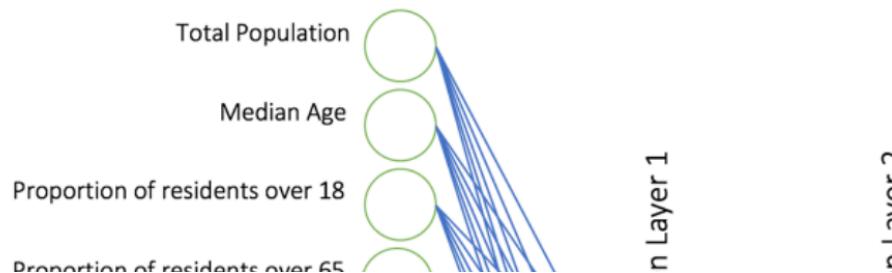
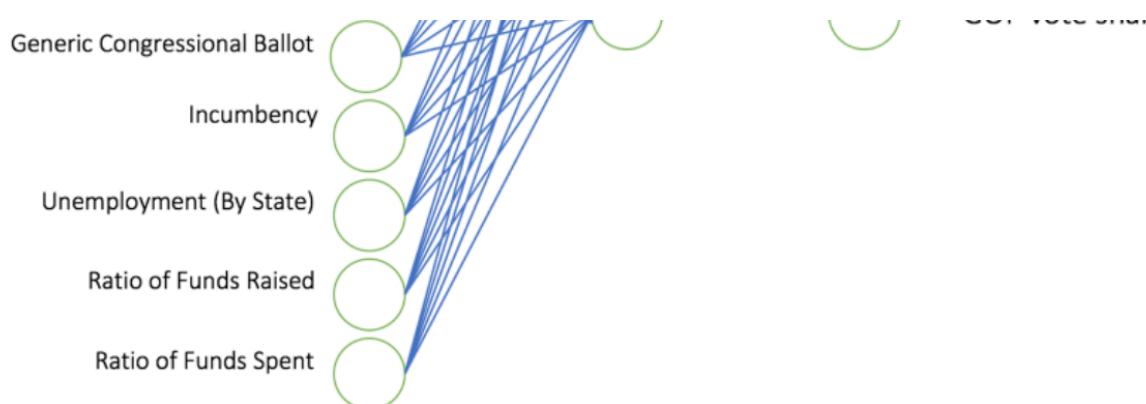


Table 2 – Results of both Models:

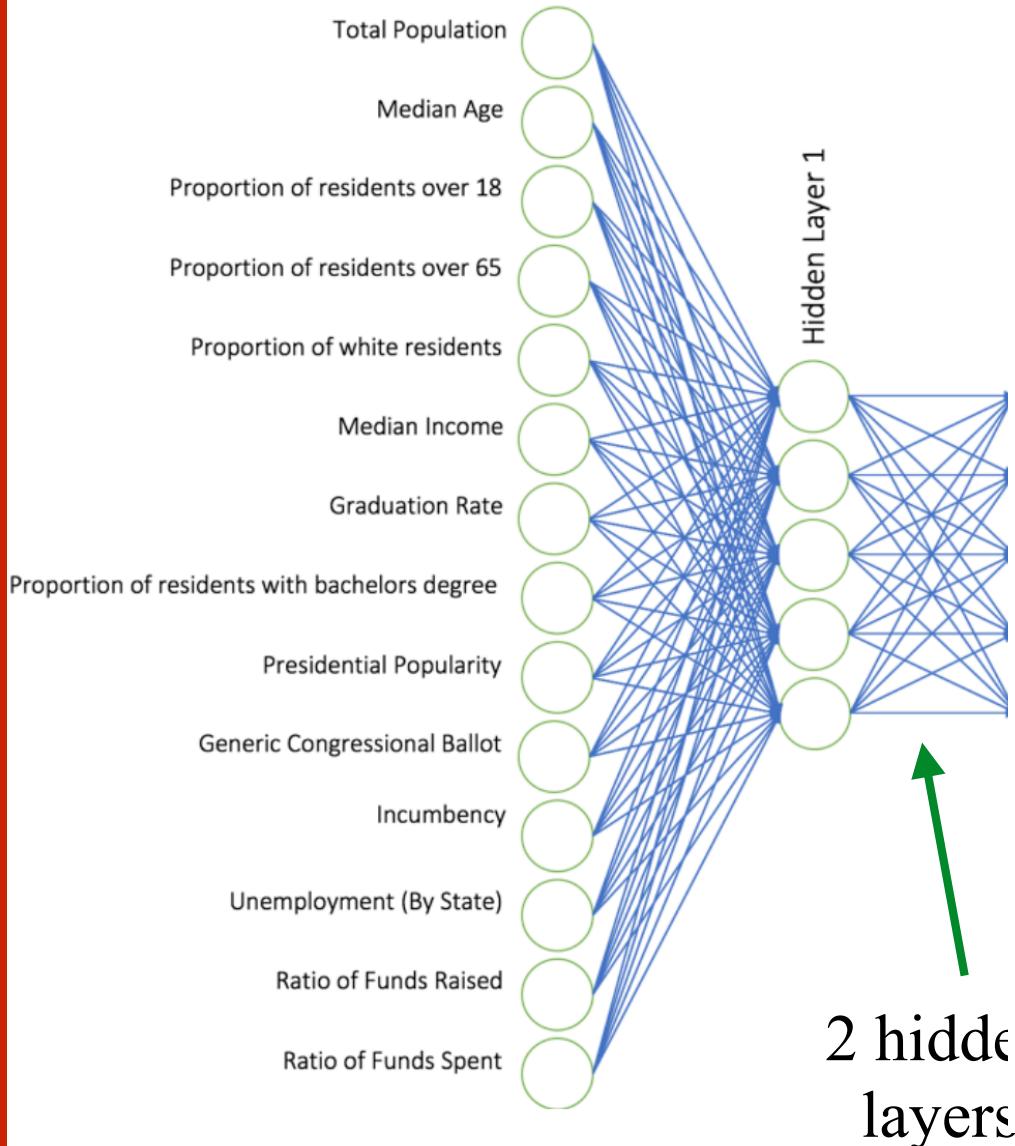
<b><i>Model not including past elections (Model A)</i></b>	<b><i>Model including past elections (Model B)</i></b>
<b>D +3</b>	<b>D +17</b>
<b>Democrat Seats: 219</b>	<b>Democrat Seats: 226</b>
<b>Republican Seats: 216</b>	<b>Republican Seats: 209</b>
<b>33% chance of Republicans keeping house*</b>	<b>0.3% chance of Republicans keeping house*</b>



232+198

# NEURAL NETWORK TO PREDICT RESULTS OF MIDTERM ELECTIONS

## Features (x)



**Bad Weather, Known to Lower Turnout, Will Greet Many Voters**

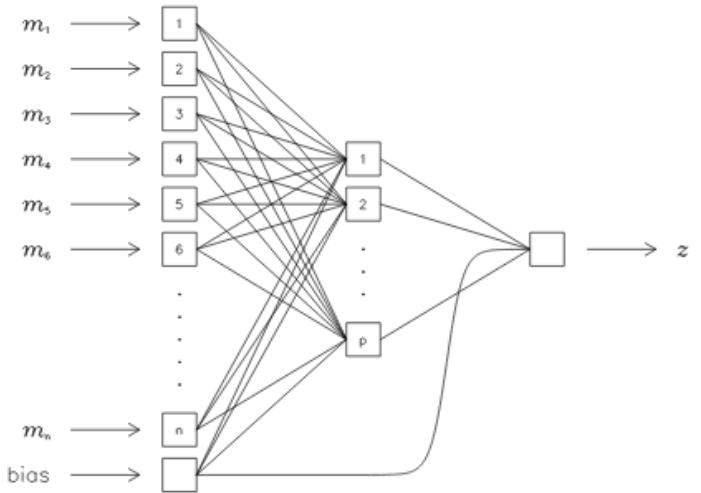
Rain can decrease voter numbers, which studies show tends to help Republicans. “I hope it rains hard tomorrow,” one Republican candidate said.

10h ago

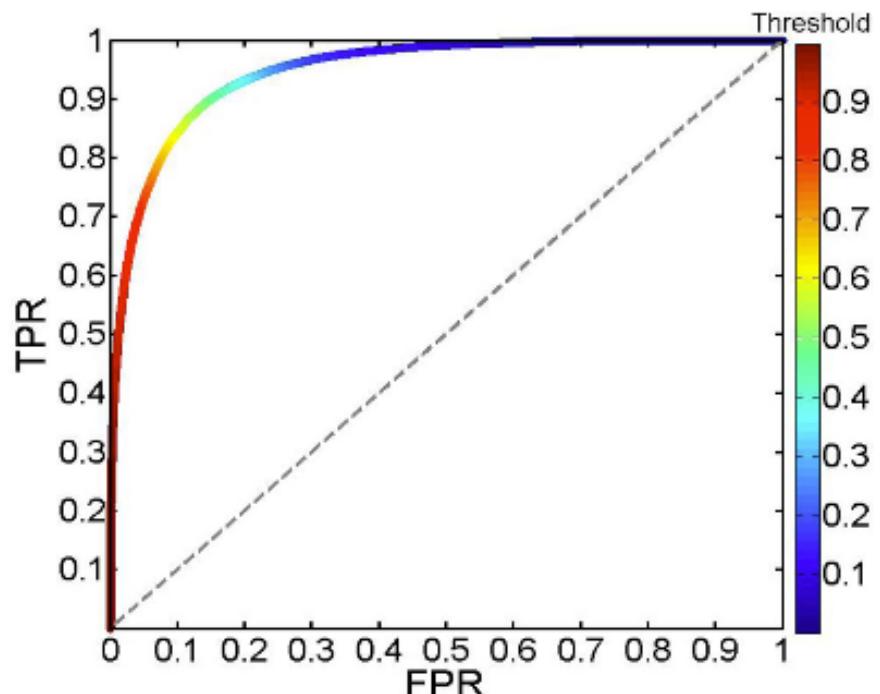
No.	Symbol	Description	Scale <sup>a</sup>
1	CVD	Central velocity dispersion	~1 kpc
2	$M_{\text{bulge}}$	Bulge stellar mass	0.5–4 kpc
3	$R_e$	Bulge effective radius	0.5–4 Kpc
4	B/T	Bulge-to-total stellar mass ratio	0.5–8 kpc
5	$M_*$	Total stellar mass	2–8 kpc
6	$M_{\text{disc}}$	Disc stellar mass	4–10 kpc
7	$M_{\text{halo}}$	Group halo mass	0.1–1 Mpc
8	$\delta_5$	Local density parameter	0.5–3 Mpc

Notes. <sup>a</sup> Approximate  $1\sigma$  range from centre of galaxy. For photometric quantities half-light radii are used.

Input layer → Hidden layer → Output layer

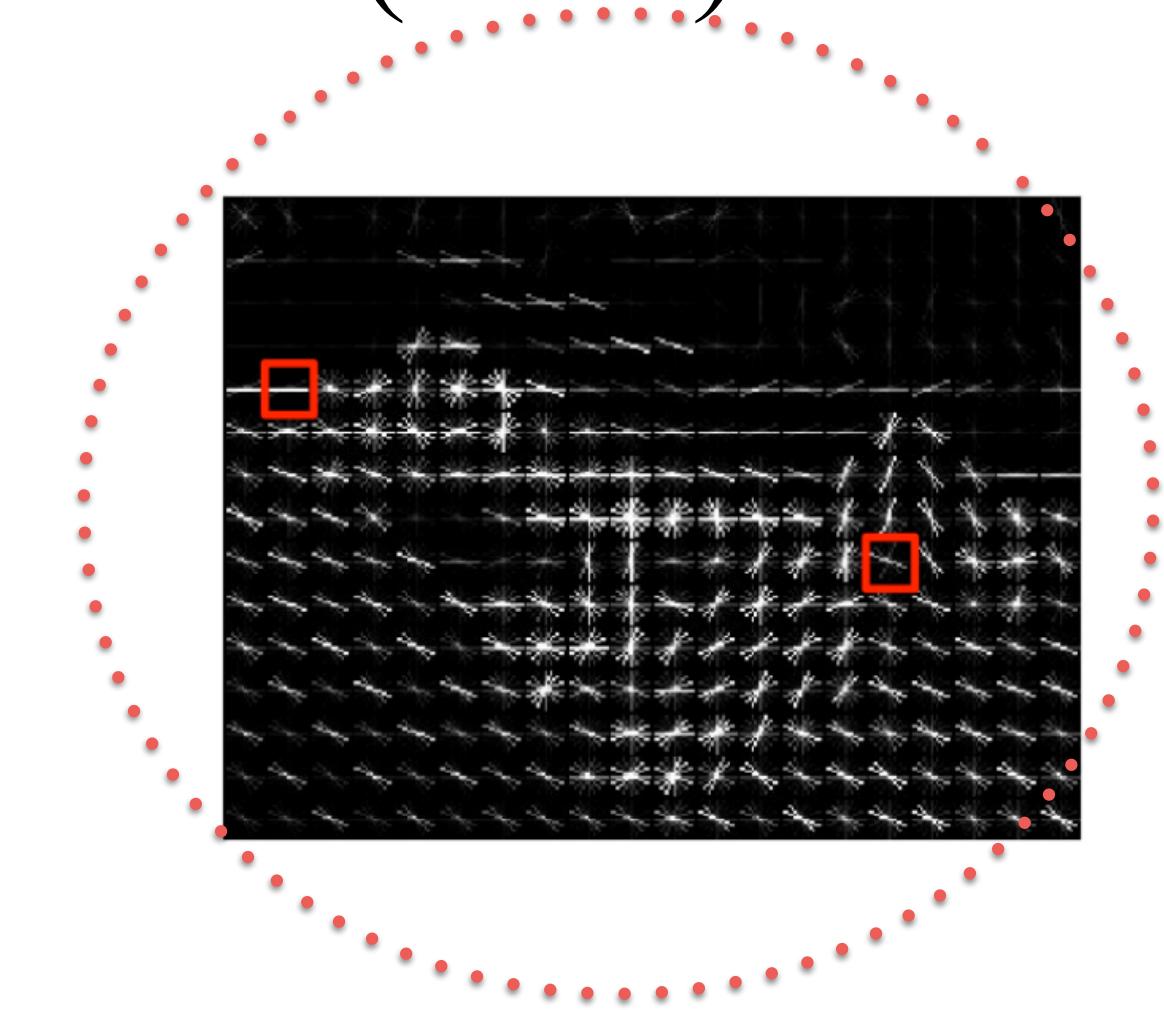


## HEAVILY PROCESSED DATA

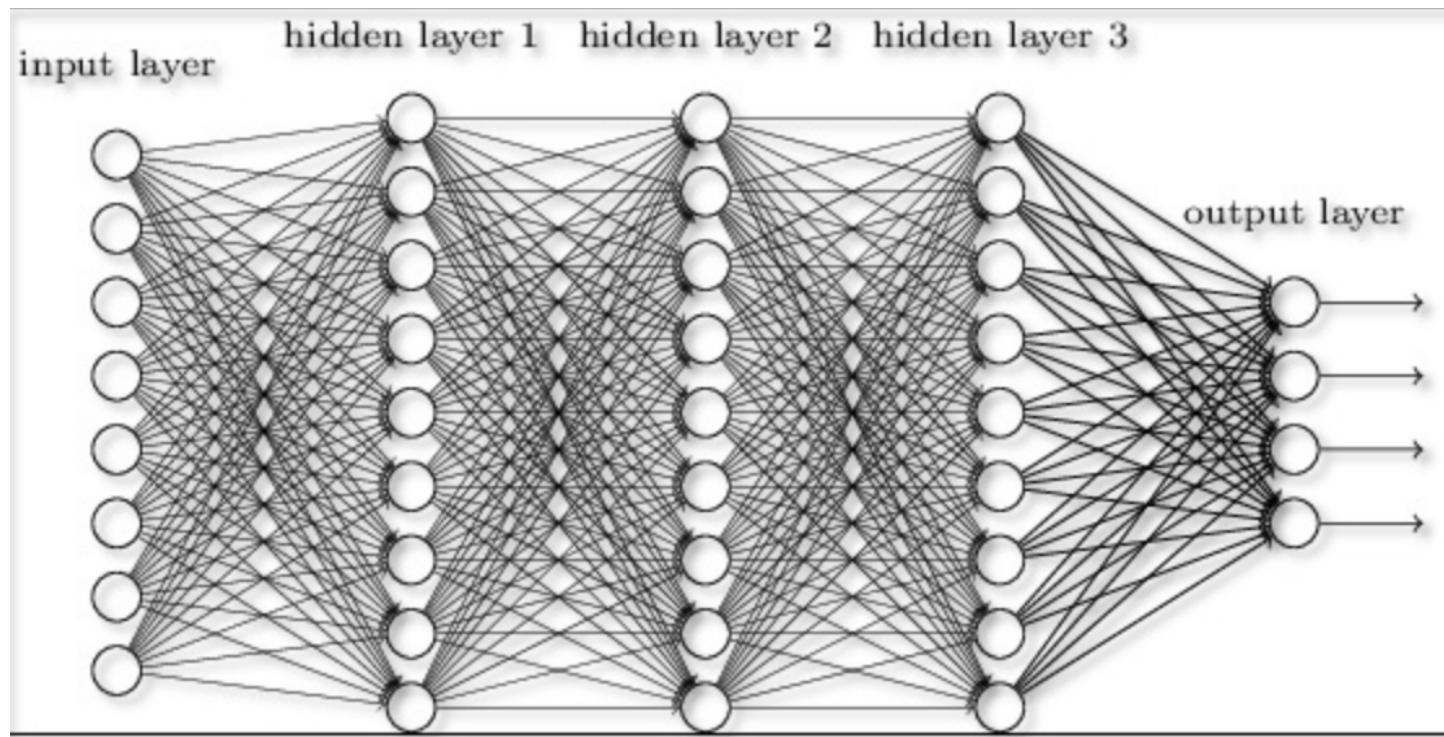


Teimiroonnia+16

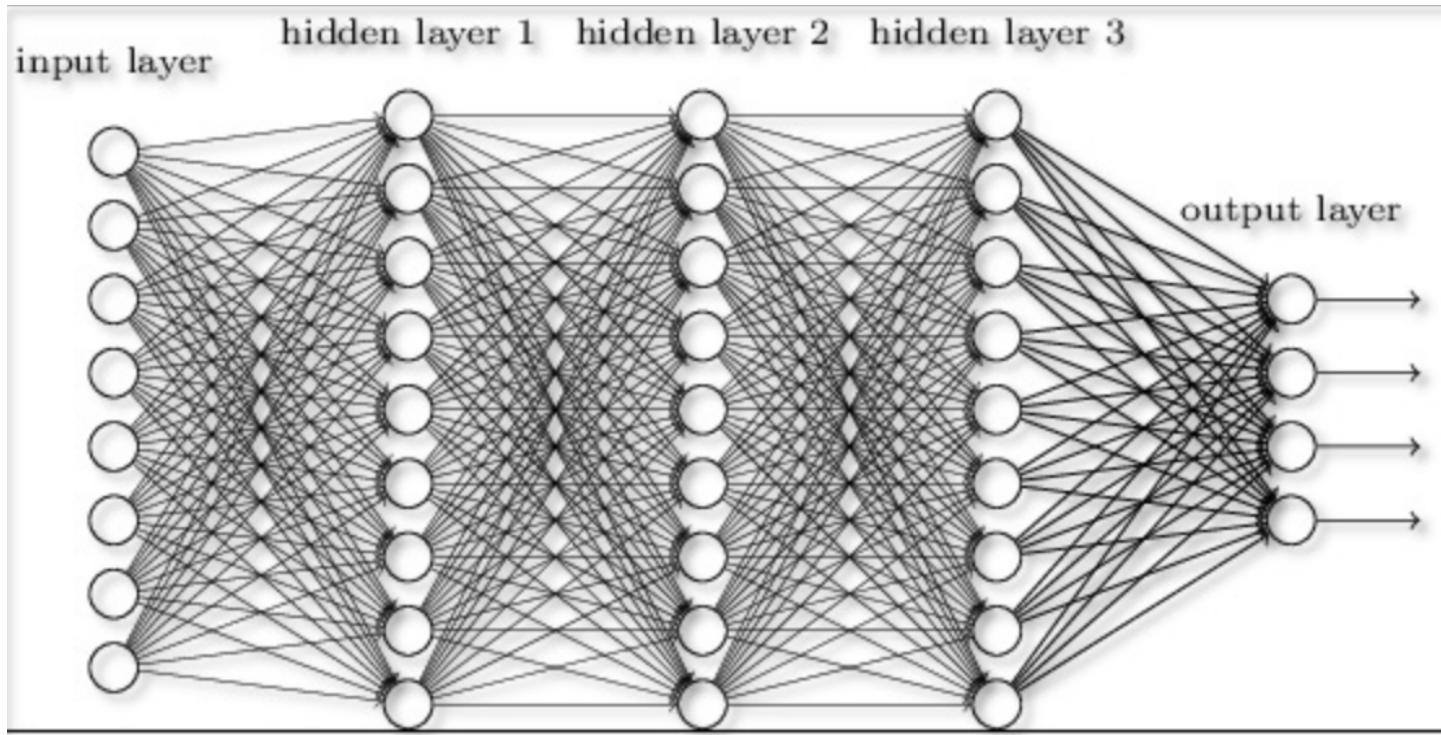
# HISTOGRAM OF ORIENTED GRADIENTS (HoG)



KEEP THIS IMAGE IN MIND FOR LATER...

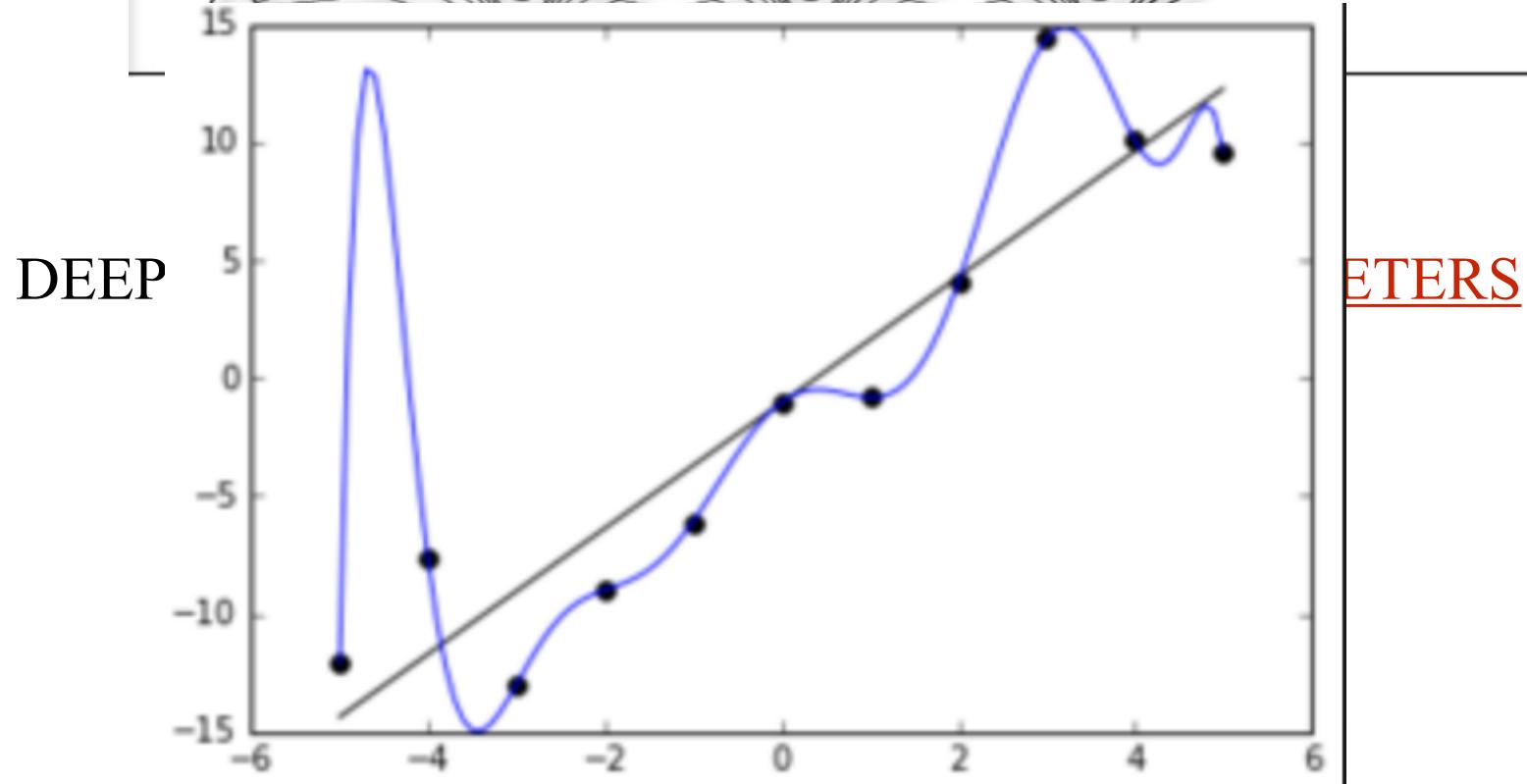
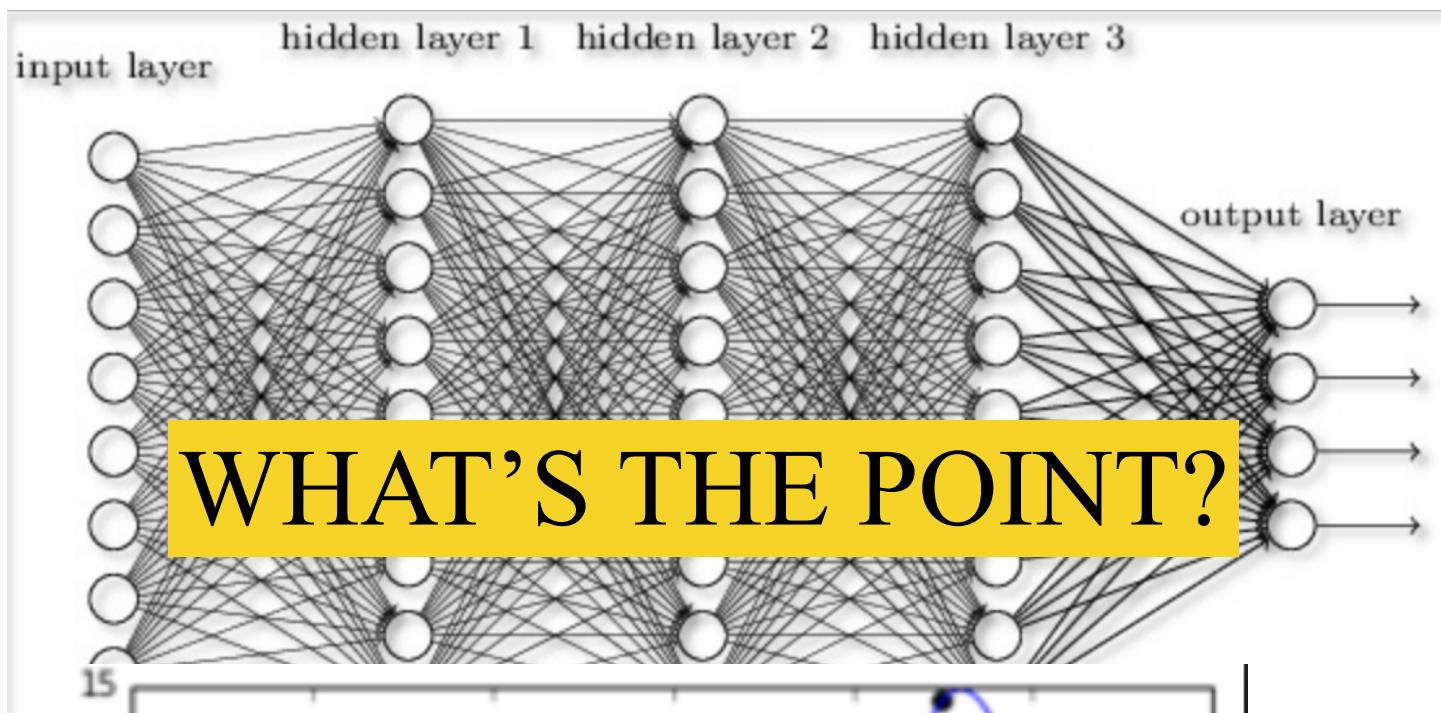


DEEPER → INCREASING NUMBER OF PARAMETERS



DEEPER → INCREASING NUMBER OF PARAMETERS

IF OUR INPUT CONTAINS A HANDFUL  
OF PARAMETERS [colors, sizes, lines etc..]



# WHAT ABOUT USING RAW DATA?

ALL INFORMATION IS IN THE INPUT DATA

WHY REDUCING ?

LET THE NETWORK FIND THE INFO

# WHAT ABOUT USING RAW DATA?

ALL INFORMATION IS IN THE INPUT DATA

WHY REDUCING ?

LET THE NETWORK FIND THE INFO

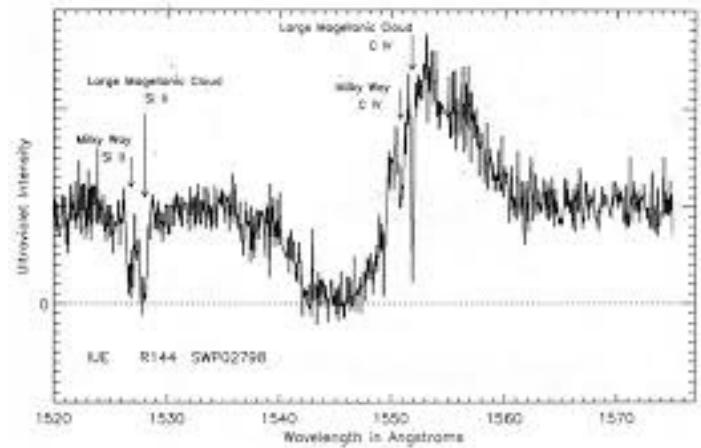
LARGE DIMENSION SIGNALS SUCH AS IMAGES OR SPECTRA WOULD REQUIRE TREMENDOUSLY LARGE MODELS

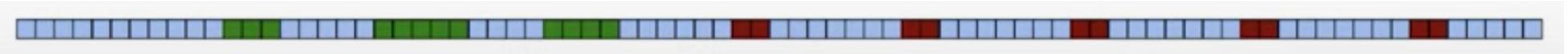
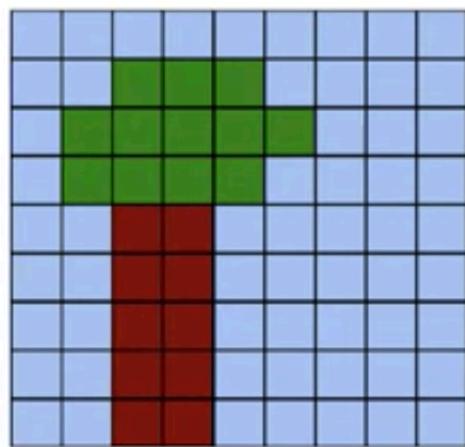
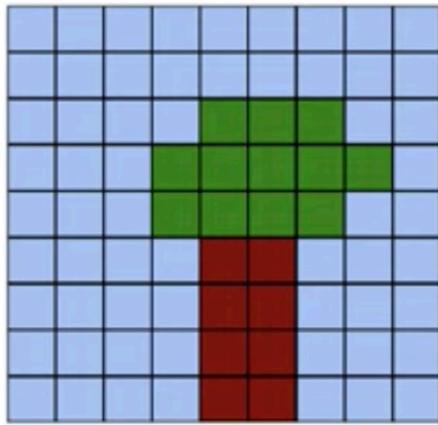
A 512x512 image as input of a fully connected layer producing output of same size:

$$(512 \times 512)^2 = 7e10$$

**BUT**

FEEDING INDIVIDUAL RESOLUTION ELEMENTS IS NOT  
VERY EFFICIENT SINCE IT LOOSES ALL INVARIANCE TO  
TRANSLATION AND IGNORES CORRELATION IN THE DATA  
AT ALL SCALES





(Dielemann@Deepmind)



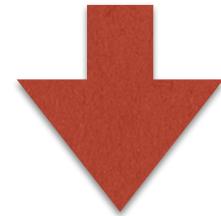
TWO BASIC PROPERTIES OF IMAGING DATA (BUT ALSO  
SPECTROSCOPY IN SOME SENSE) ARE **LOCALITY**  
**TRANSLATION INVARIANCE**

**locality**: nearby pixels are more strongly correlated

**translational invariance**: meaningful patterns can appear anywhere  
in the image

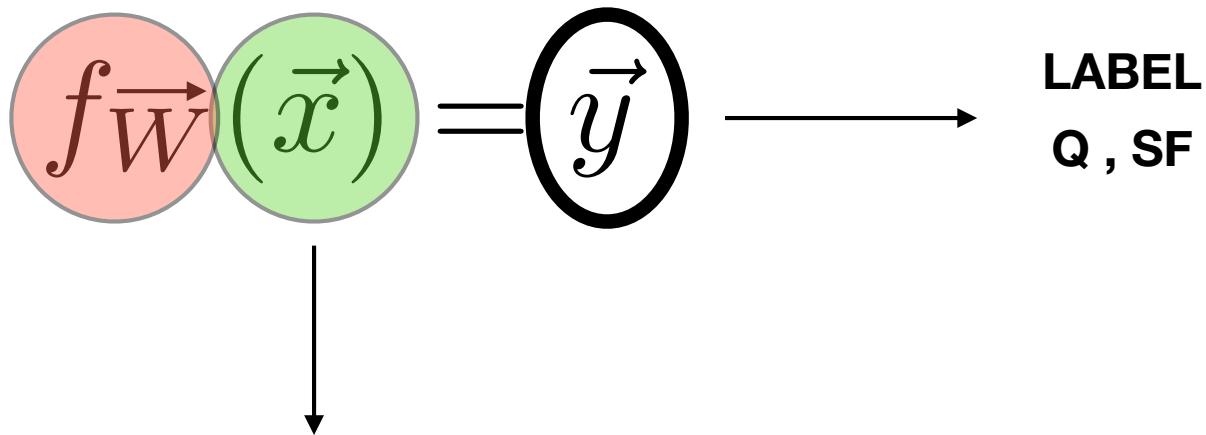
(Dielemann@Deepmind)

FEEDING INDIVIDUAL RESOLUTION ELEMENTS IS NOT  
VERY EFFICIENT SINCE IT LOSES ALL INVARIANCE TO  
TRANSLATION



SO?

## DEEP LEARNING



LET THE MACHINE FIGURE THIS OUT ("unsupervised feature extraction")

LET'S GO A STEP FORWARD INTO LOOSING CONTROL....