

An introduction to neural networks: from “classical” to deep learning

Marc Huertas-Company



SOME PRELIMINARY NOTES

I AM NOT A MACHINE LEARNING RESEARCHER

BUT I HAVE BEEN USING MACHINE LEARNING FOR THE
LAST ~15 YEARS FOR MY RESEARCH IN ASTROPHYSICS

THESE LECTURES ARE INTENDED TO PROVIDE A **GENERAL**
INTRODUCTION TO NEURAL NETWORKS

ML@EDE20

CYCLE 1: INTRODUCTION AU DEEP LEARNING:
tutorial session at 2pm today (Monday)

CYCLE 2: 3h LECTURES ON NEURAL NETWORKS AND
DEEP LEARNING 9am (Monday) + Tutorials 2pm (Tuesday)

ALL: Séance Ouverte 5.30 pm

ML@EDE20

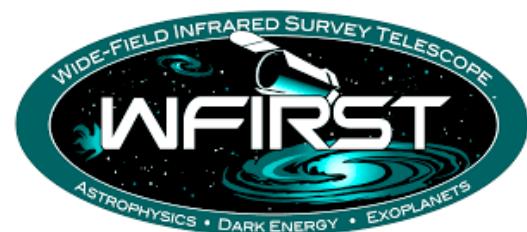
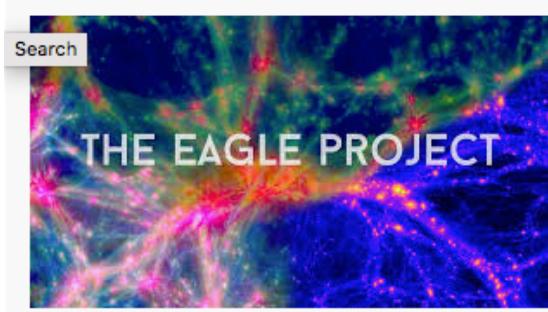
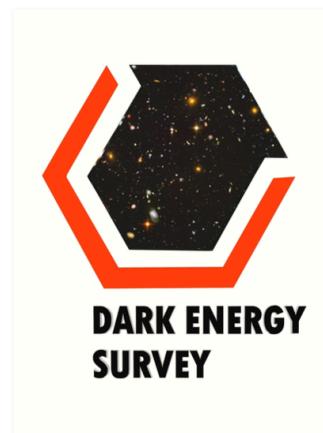
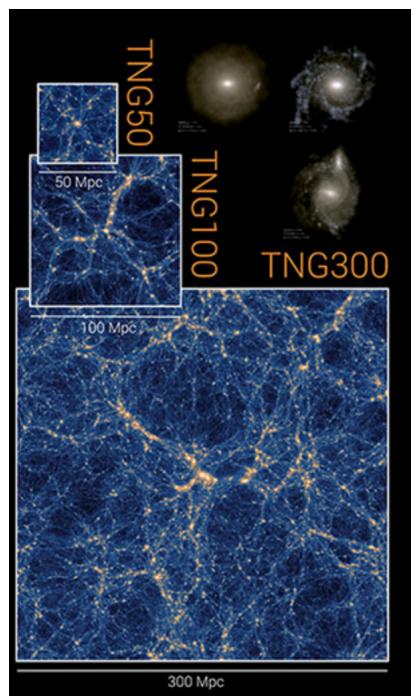
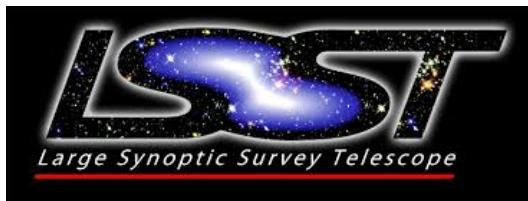
REQUIREMENTS FOR THE TUTORIALS:

- GOOGLE DRIVE ACCOUNT WITH SOME FREE SPACE
- INTERNET CONNECTION ?!

GITHUB REPO WITH MORE COMPLETE SLIDES AND COLAB TUTORIALS:

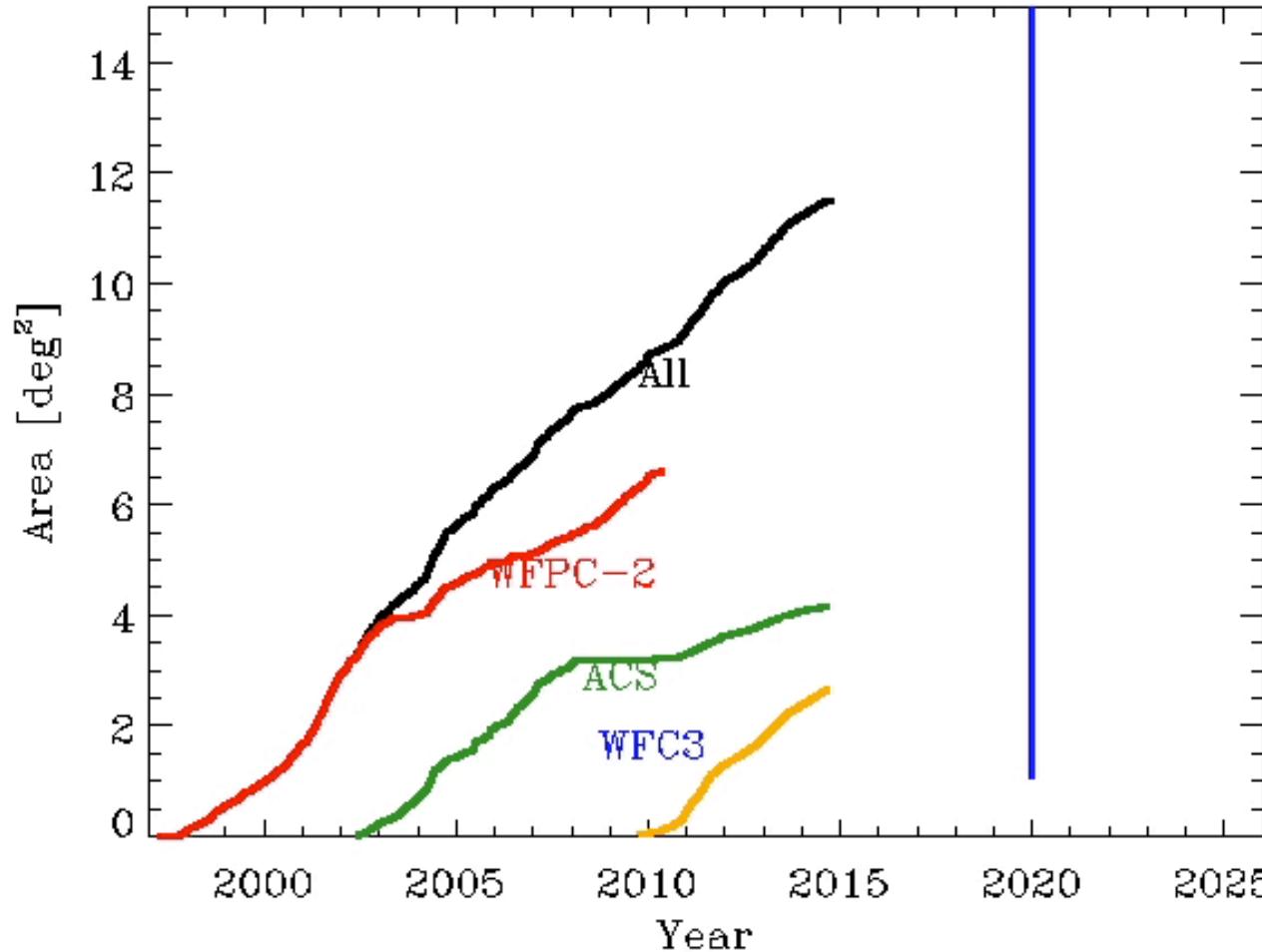
<https://github.com/mhuertascompany/EDE20>

WE DO NOT HAVE THE CAPACITY TO “LOOK” AT FUTURE ASTRONOMICAL (BIG-DATA) SETS

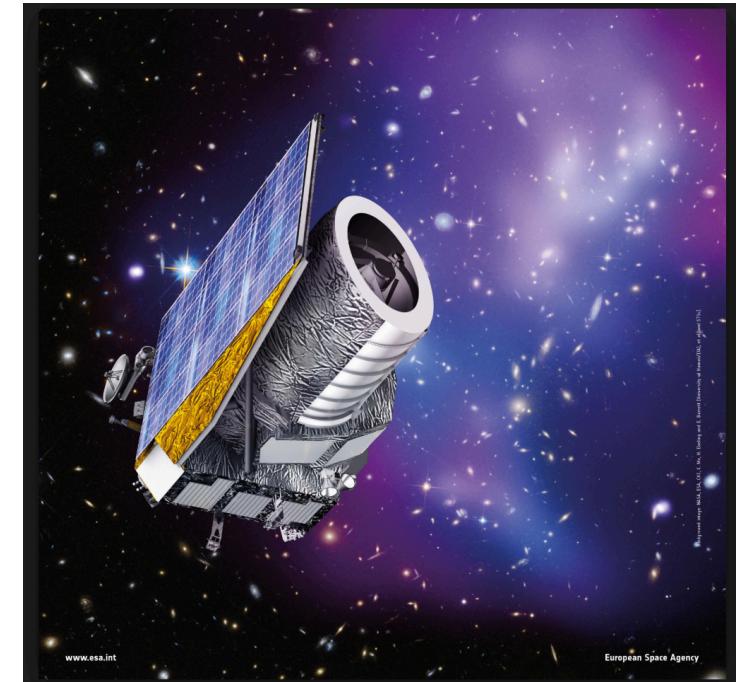


The Horizon Simulation

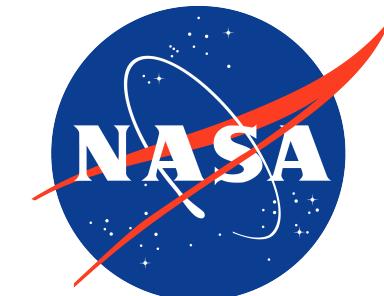
WE WILL SOON ENTER A NEW REGIME



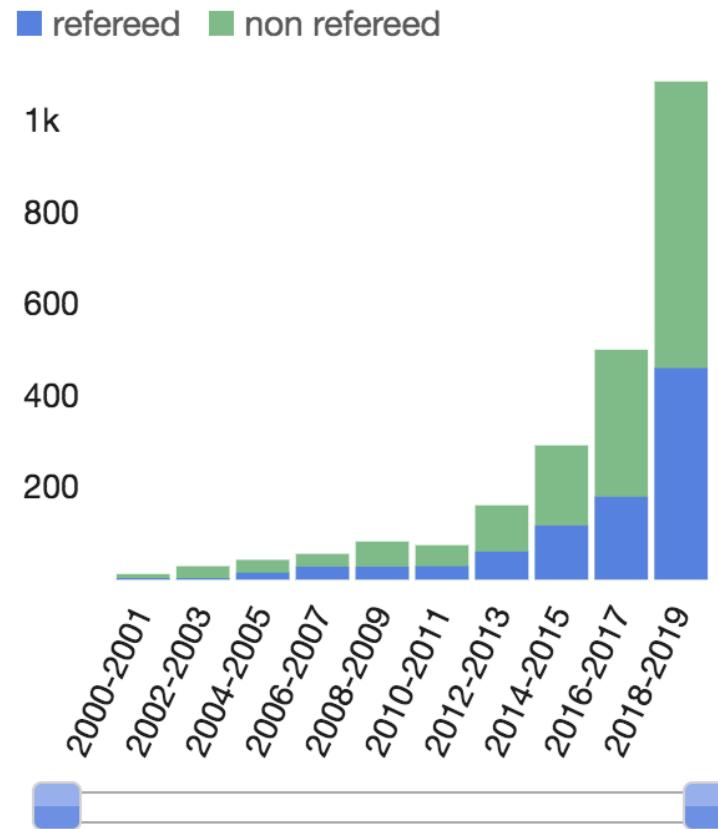
(Thanks to J. Brinchmann)



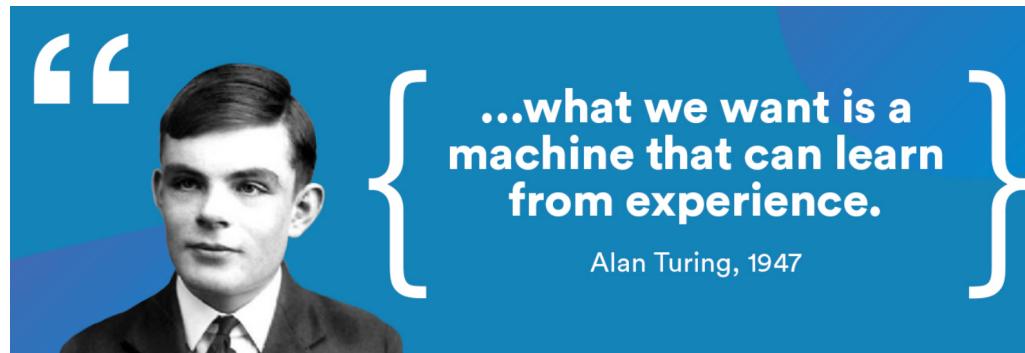
EUCLID space telescope
(2021)



AI TO THE RESCUE?

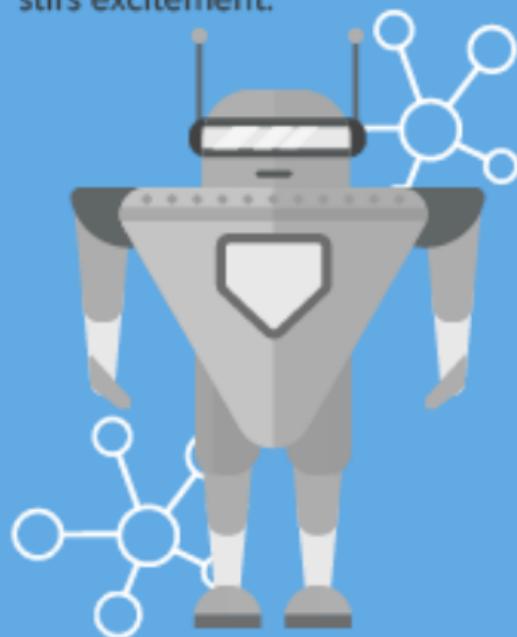


SOURCE: ADS



ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



1950's

1960's

1970's

1980's

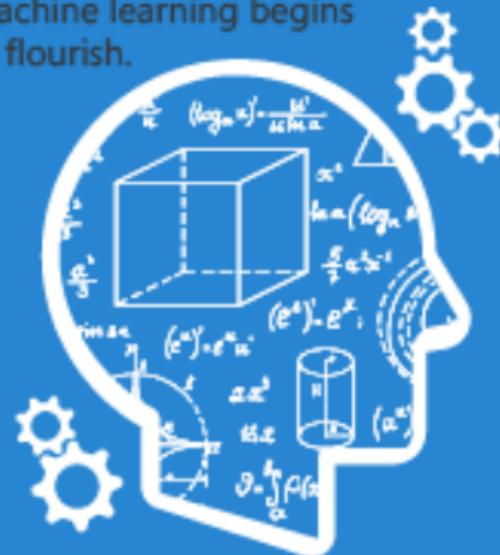
1990's

2000's

2010's

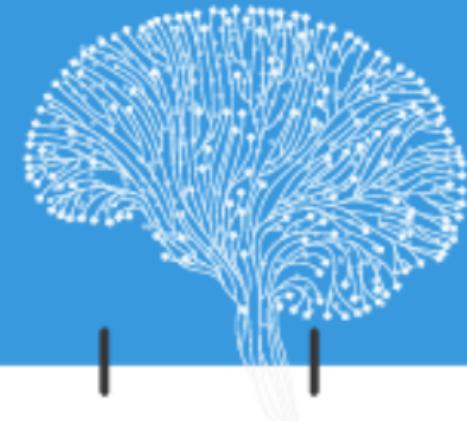
MACHINE LEARNING

Machine learning begins to flourish.

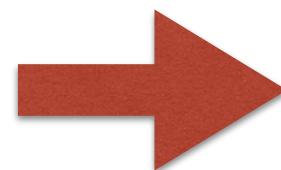


DEEP LEARNING

Deep learning breakthroughs drive AI boom.



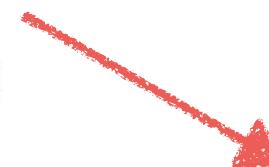
BEFORE 2012....



CAT?

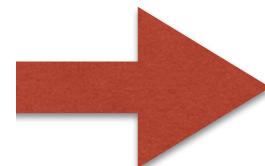


DOG?



**TRIVIAL HUMAN TASKS REMAINED
CHALLENGING FOR COMPUTERS**

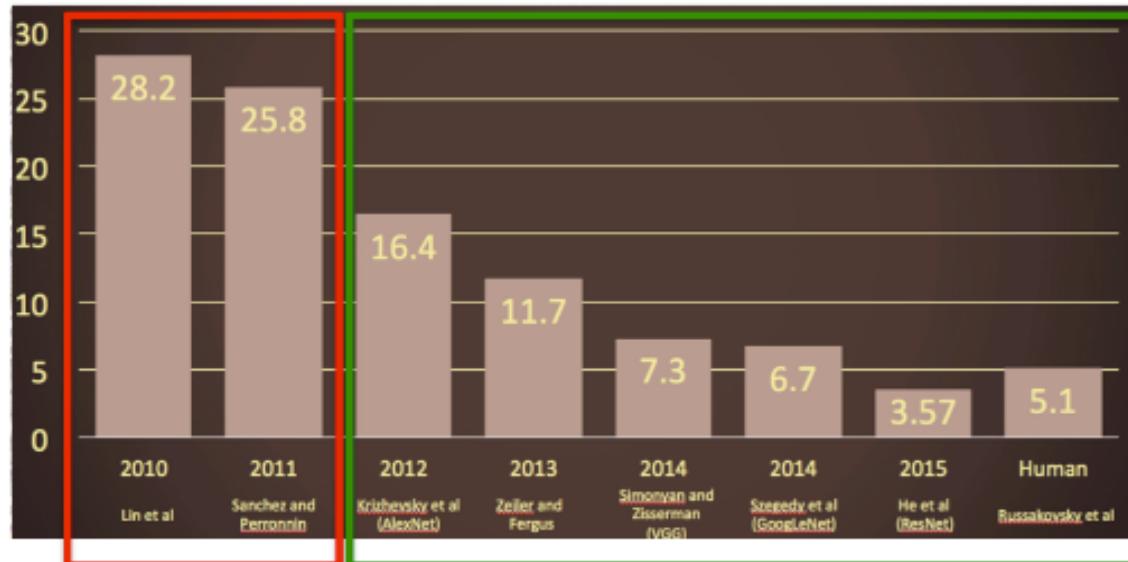
AFTER 2012



IT HAS BECOME TRIVIAL....

THIS IS A CHANGE OF PARADIGM!

Fisher Vectors



CNNs

*ImageNet
top-5 error (%)*



ONE OF THE MAIN REASONS OF THIS BREAKTHROUGH IS THE AVAILABILITY OF VERY LARGE DATASETS TO LEARN



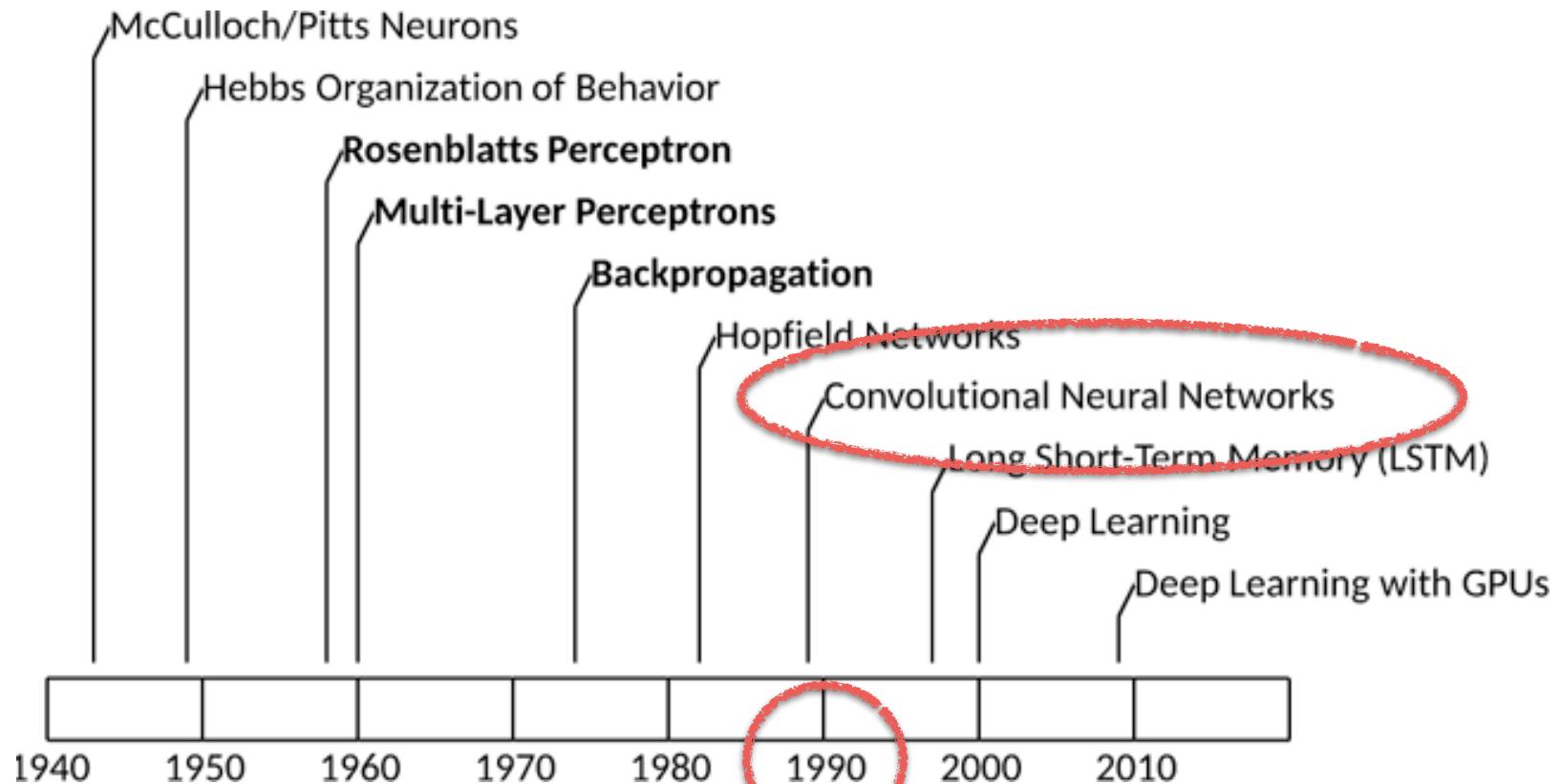
**COMBINED WITH THE TECHNOLOGY TO PROCESS ALL THIS
DATA**



ONE OF THE MAIN REASONS OF THIS BREAKTHROUGH IS THE AVAILABILITY OF VERY LARGE DATASETS TO LEARN



BUILDING BLOCKS OF DEEP LEARNING ARE FROM THE 90s



PROGRAM FOR TODAY

- **PART I: INTRODUCTION**
 - UNSUPERVISED / SUPERVISED
 - GENERAL STEPS TO “TEACH A MACHINE”
 - LOSS FUNCTION, OPTIMIZATION
- **PART II: THE FOUNDATIONS OF NEURAL NETWORKS**
 - PERCEPTRON, NEURON DEFINITION
 - LAYER OF NEURONS, HIDDEN LAYERS
 - ACTIVATION FUNCTIONS
 - OPTIMIZATION [GRADIENT DESCENT, LEARNING RATES]
 - BACKPROPAGATION
 - NEURAL NETWORKS AS STATISTICAL MODELS: LOSS FUNCTIONS

PROGRAM FOR TODAY

- **PART III: THE FOUNDATIONS OF DEEP LEARNING**
 - CONVOLUTIONS AS NEURONS
 - CNNs [POOLING, DROPOUT]
 - VANISHING GRADIENT / BATCH NORMALIZATION
- **PART IV: BEYOND CLASSIFICATION - IMAGE2IMAGE NETWORKS**
 - FCNNs (IMAGE2IMAGE NETWORKS)
 - OBJECT DETECTION
 - INSTANCE AND SEMANTIC SEGMENTATION
 - (A VERY BRIEF INCURSION INTO UNSUPERVISED LEARNING:
AUTOENCODERS)

REFERENCES

SEVERAL SLIDES / INFOS SHOWN HERE ARE INSPIRED/
TAKEN FROM OTHER WORKS / COURSES FOUND ONLINE

- Deep Learning: Do-It-Yourself! [Bursuc, Krzakala, Lelarge]
- DEEPMIND.AI [COURSERA, Ng, Bensouda, Katanforoosh]
- MACHINE LEARNING LECTURES [Keck]
- EPFL DEEP LEARNING COURSE [Fleuret]
- DEEPMIND / UCL Lecture series [DeepMind]

+ many others!

Thanks to all of them!

PART I: AN INTRODUCTION TO “CLASSICAL” SUPERVISED MACHINE LEARNING

WHAT DOES MACHINE LEARNING DO?

the machine is told what to look for

SUPERVISED

the machine is NOT told what to look for

UN-SUPERVISED

TWO VERY BROAD TYPES OF MACHINE LEARNING ALGORITHMS

CAT



CAT

SUPERVISED LEARNING

DOG



HUMAN LABELLING

DOG



the machine is told what to look for

CAT



CAT



DOG



DOG



SUPERVISED LEARNING

the machine is told what to look for

HUMAN LABELLING

TRAINING SET
OF LABELED
EXAMPLES

CAT



CAT



DOG



DOG



SUPERVISED LEARNING



ML



CAT

CAT



CAT



DOG



DOG



SUPERVISED LEARNING



ML



DOG

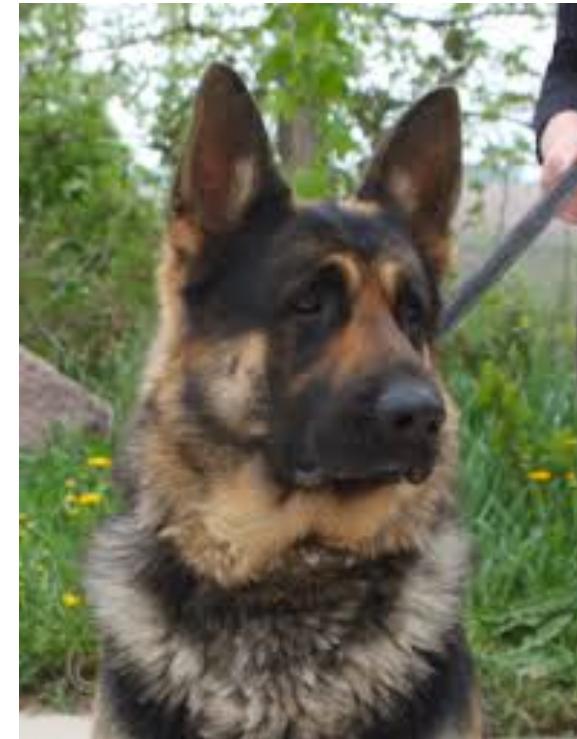
UNSUPERVISED LEARNING



UNSUPERVISED LEARNING



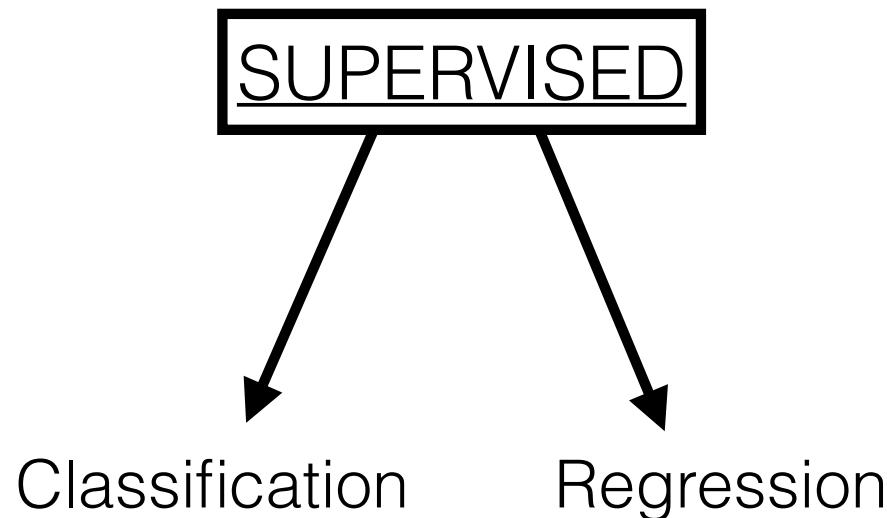
UNSUPERVISED LEARNING



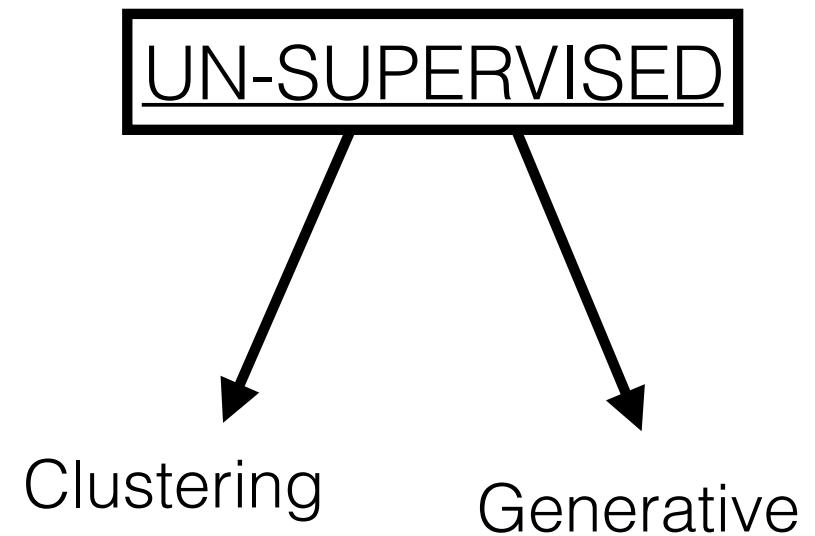
THE DEFINITION OF CLASSES IS SOMETIMES
NOT OBVIOUS

WHAT DOES MACHINE LEARNING DO?

the machine is told what to look for



the machine is NOT told what to look for



WHAT DOES MACHINE LEARNING DO?

the machine is told what to look for

SUPERVISED

Classification

Regression

the machine is NOT told what to look for

UN-SUPERVISED

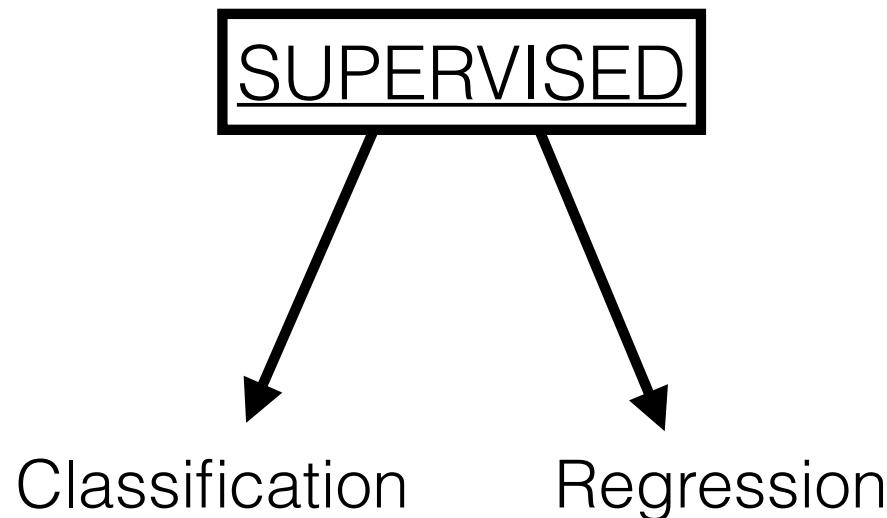
Clustering

Generative

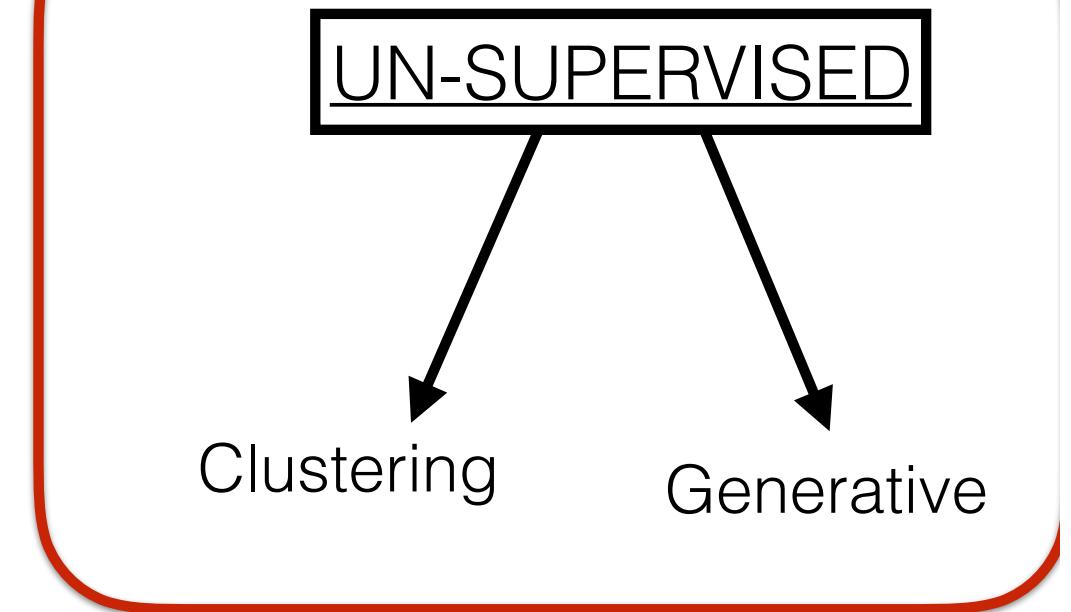
MOSTLY FOCUS ON THIS

WHAT DOES MACHINE LEARNING DO?

the machine is told what to look for



the machine is NOT told what to look for



MAYBE VERY BRIEFLY TODAY
(OR CYCLE3?)

WHAT DOES MACHINE LEARNING DO?

SUPERVISED

Classification

Regression

UN-SUPERVISED

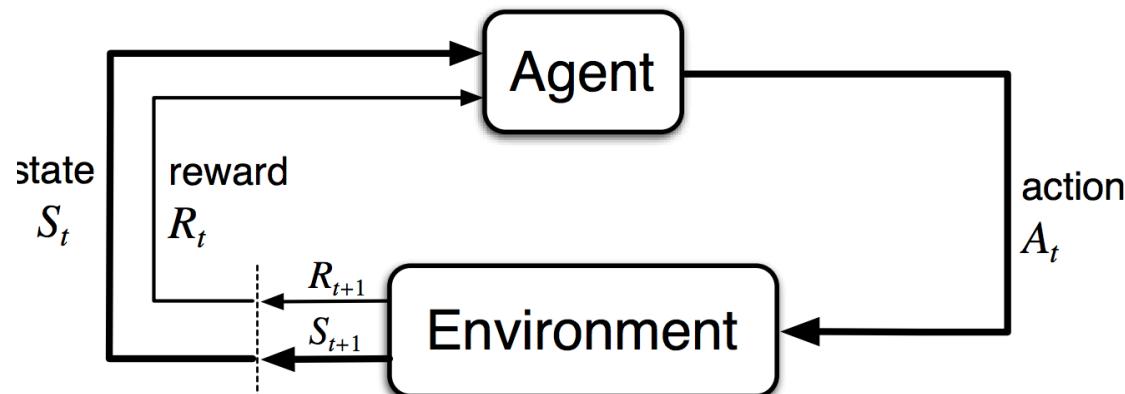
Clustering

Generative
(deep learning)

DEEP LEARNING

ACTUALLY, THERE IS A THIRD AND FORTH TYPE THAT WE WILL NOT HAVE TIME TO COVER

REINFORCEMENT LEARNING:

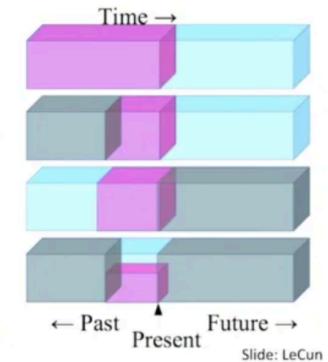


source: Sutton & Barto

**ITERATIVE LEARNING THROUGH TRIAL/ERROR:
USED FOR ALPHAGO FOR EXAMPLE**

SELF-SUPERVISED LEARNING:

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the **occluded** from the **visible**
- ▶ **Pretend there is a part of the input you don't know and predict that.**

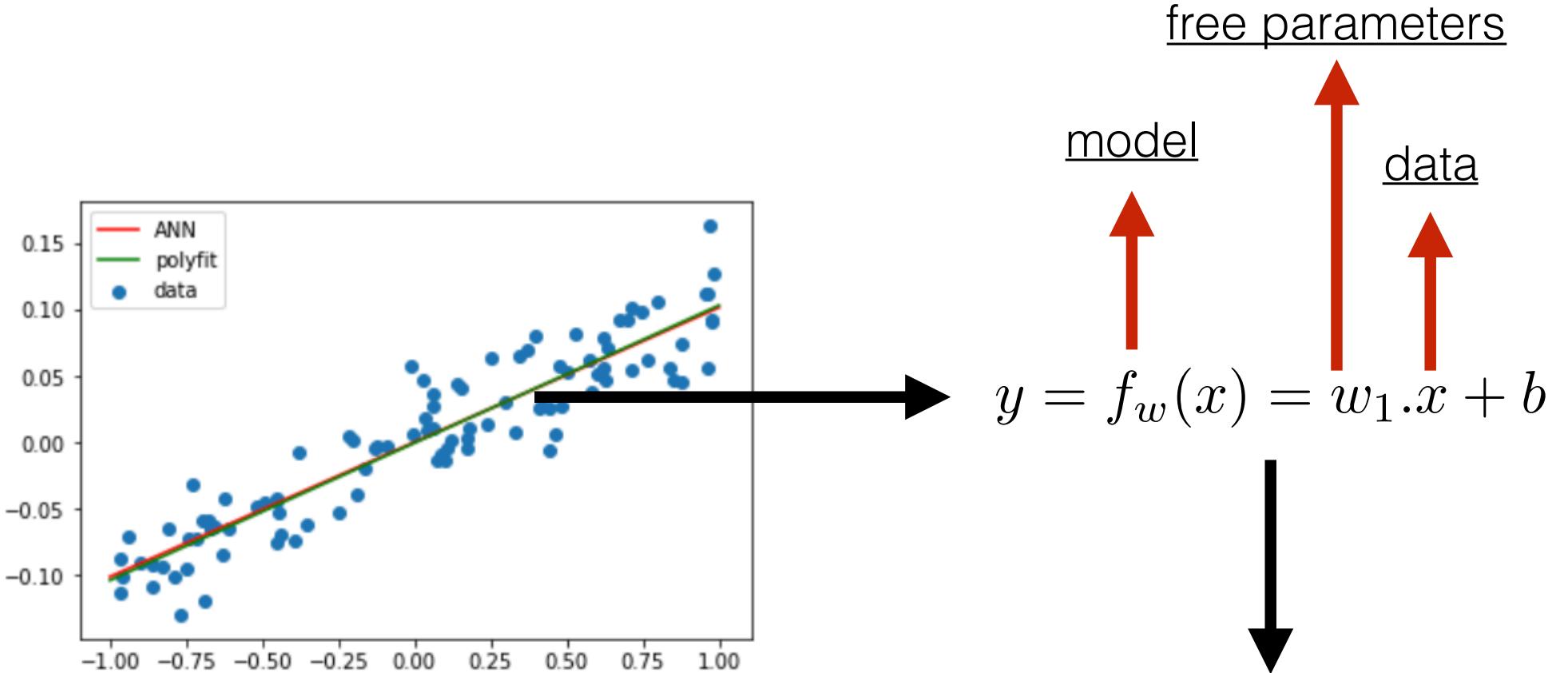


That's also why more knowledge about the structure of the world can be learned through self-supervised learning than from the other two paradigms: the data is unlimited, and amount of feedback provided by each example is huge.

Y. LeCun

a self-supervised learning system attempts to predict parts of its inputs based on the other parts of its inputs

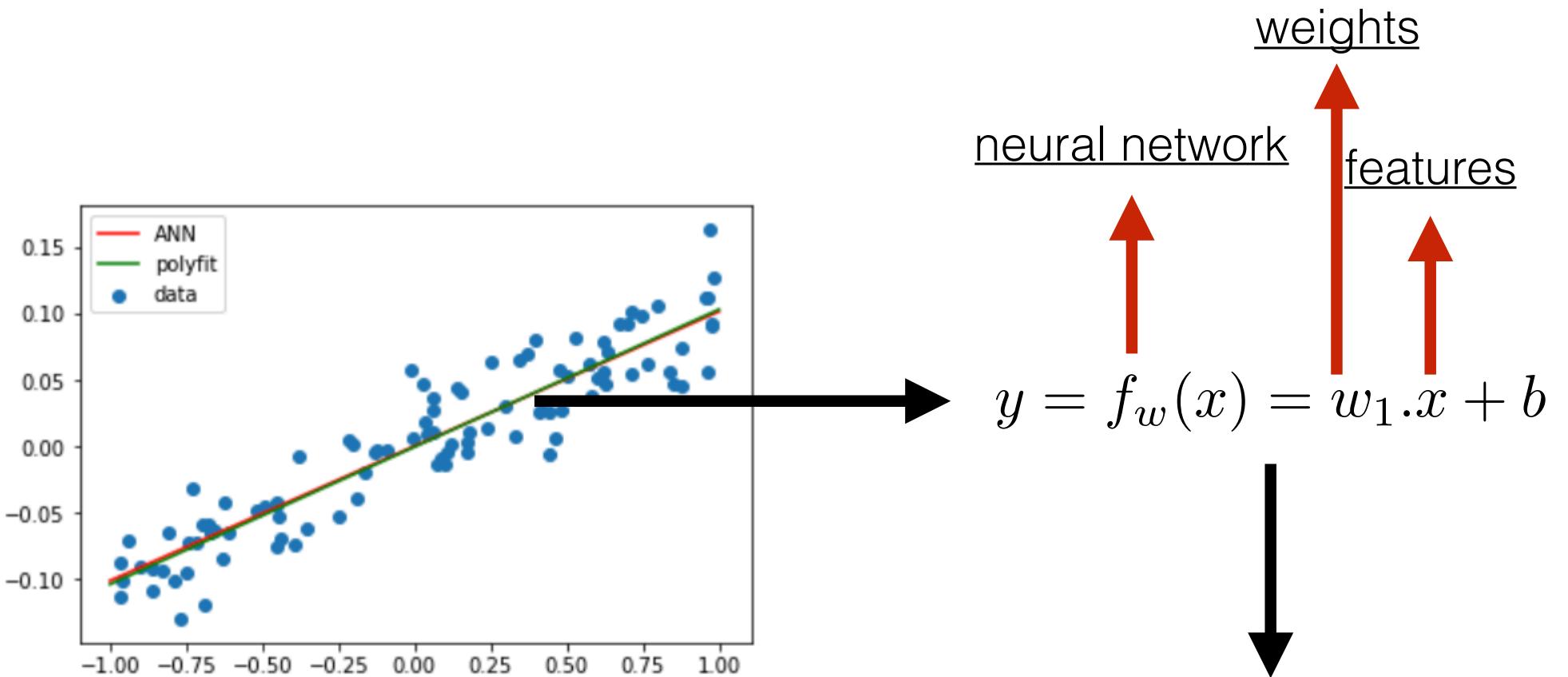
THRE IS NO MAGIC IN MACHINE LEARNING, LINEAR REGRESSION IS MACHINE LEARNING



**we fit a linear model to the data
by minimizing the chi2 error:**

$$\frac{1}{N} \sum (x^2 - f_w(x)^2)$$

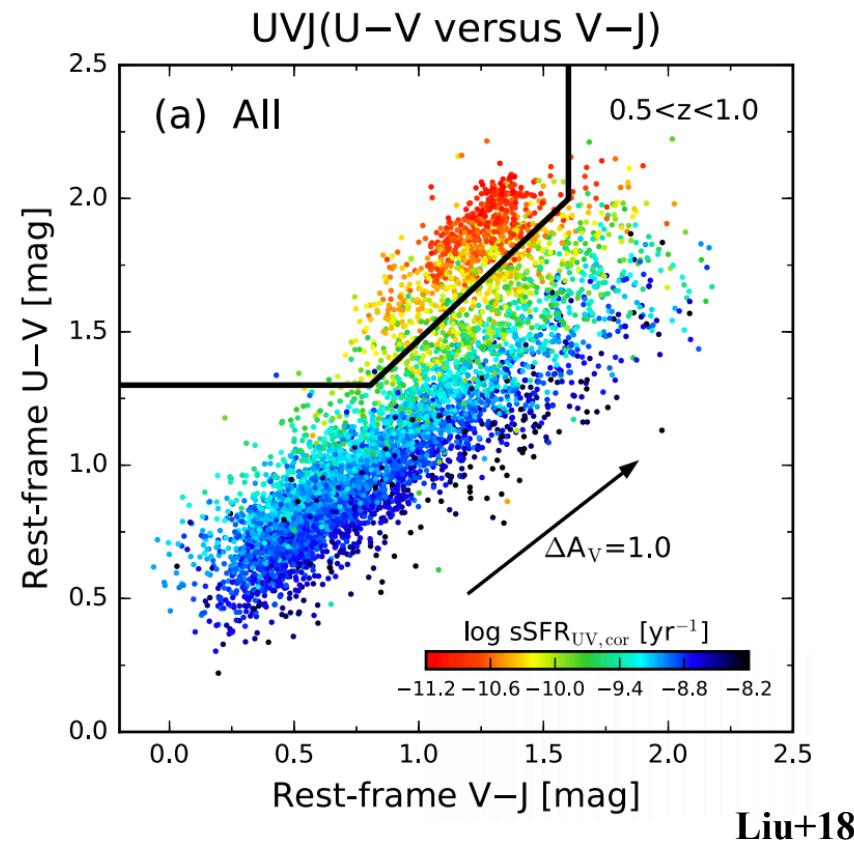
THRE IS NO MAGIC IN MACHINE LEARNING, LINEAR REGRESSION IS MACHINE LEARNING



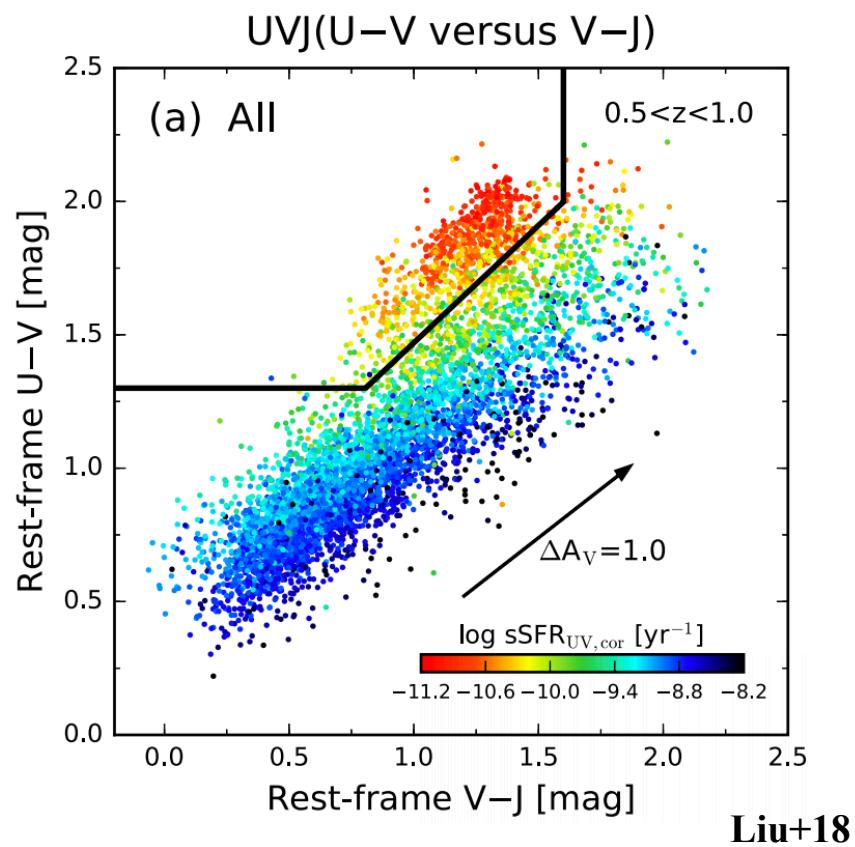
**we train a neural network with a mse
loss function:**

$$\frac{1}{N} \sum (x^2 - f_w(x)^2)$$

THRE IS NO MAGIC IN MACHINE LEARNING, AND IT IS ACTUALLY PRETTY SIMPLE

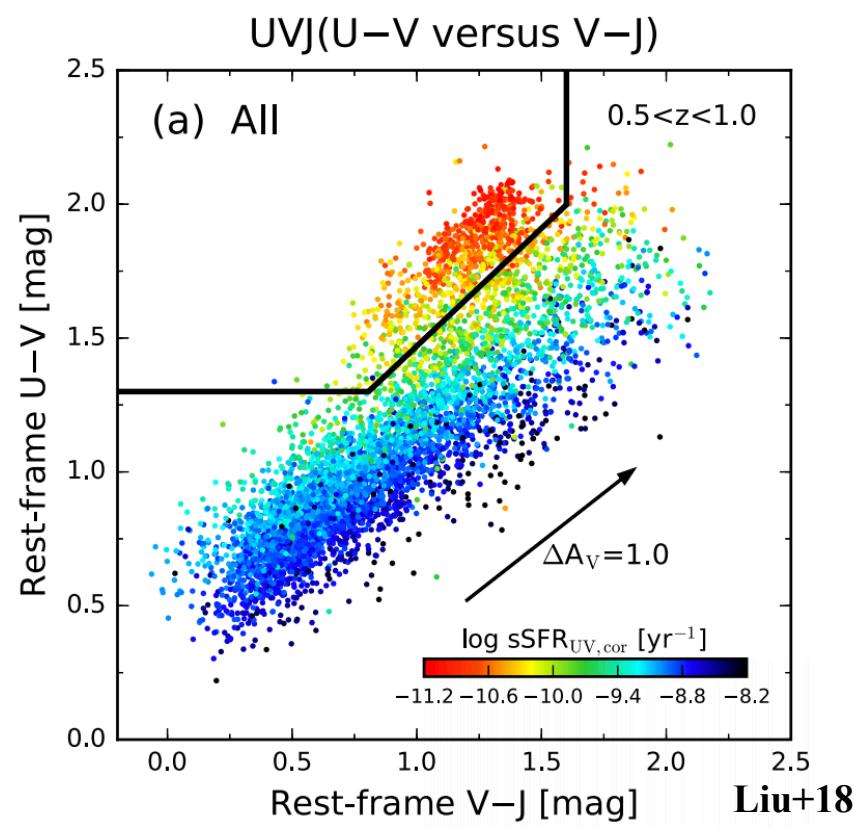


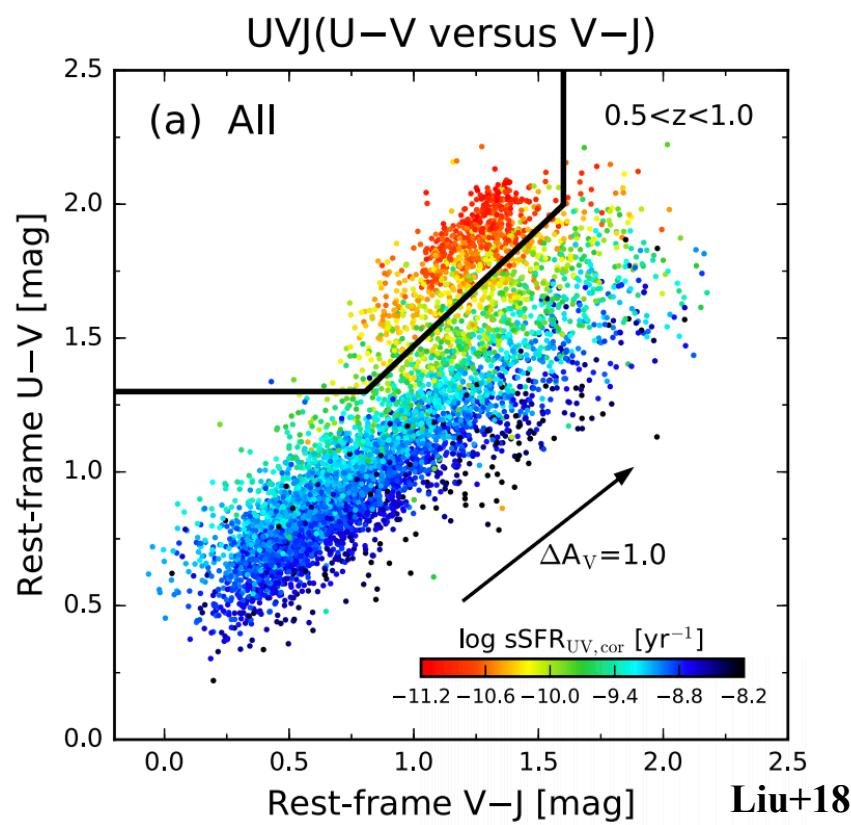
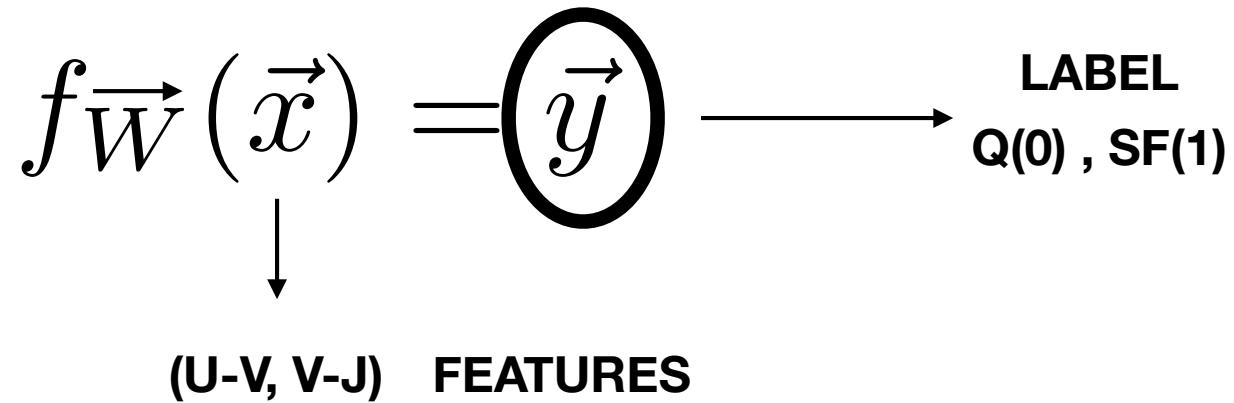
$$f_W(\vec{x}) = \vec{y}$$



$$f_W(\vec{x}) = \vec{y}$$

LABEL
Q , SF





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

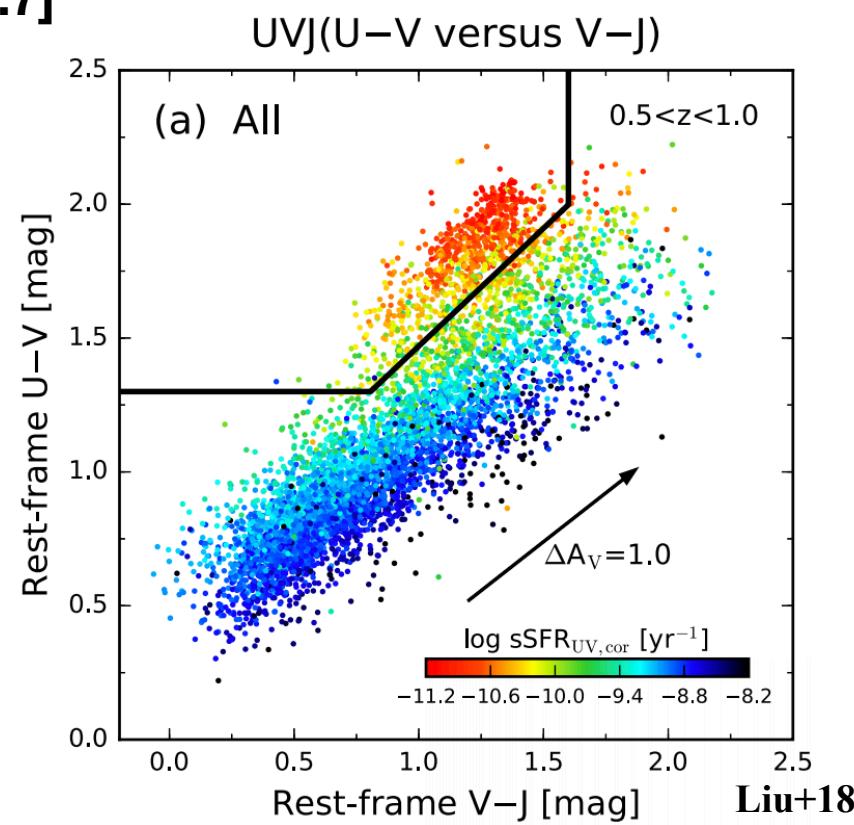
Q(0) , SF(1)

NETWORK FUNCTION

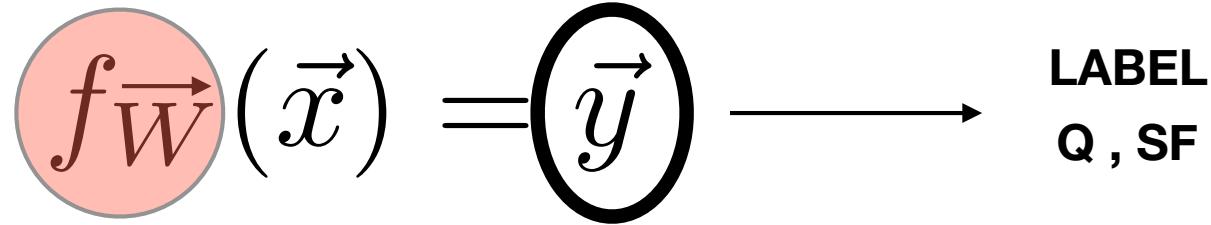
(U-V, V-J) FEATURES

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

WEIGHTS



**“CLASSICAL”
MACHINE LEARNING**

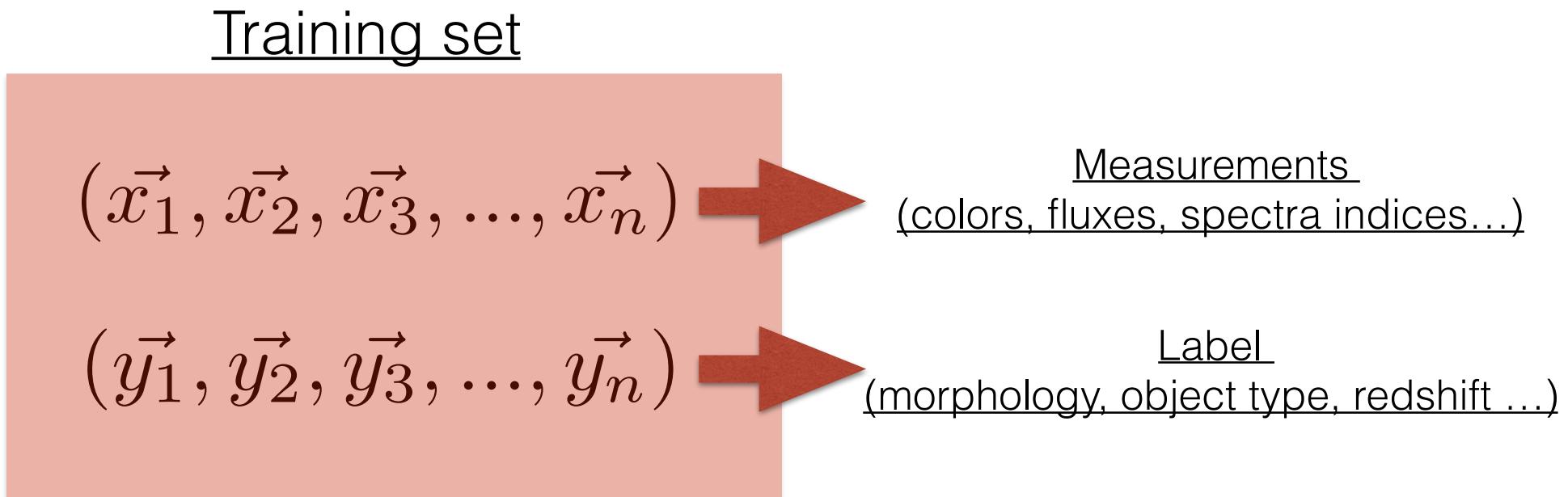


$$\text{sgn}[(u-v)-W_1*(v-j)-W_2]$$

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

SUPERVISED LEARNING

Given a dataset with known labels (measurements) - find a function that can assign (predict) measurements for an unlabeled dataset



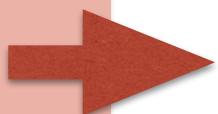
SUPERVISED LEARNING

Given a dataset with known labels (measurements) - find a function that can assign (predict) measurements for an unlabeled dataset

Training set

$$(\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n)$$

$$(\vec{y}_1, \vec{y}_2, \vec{y}_3, \dots, \vec{y}_n)$$

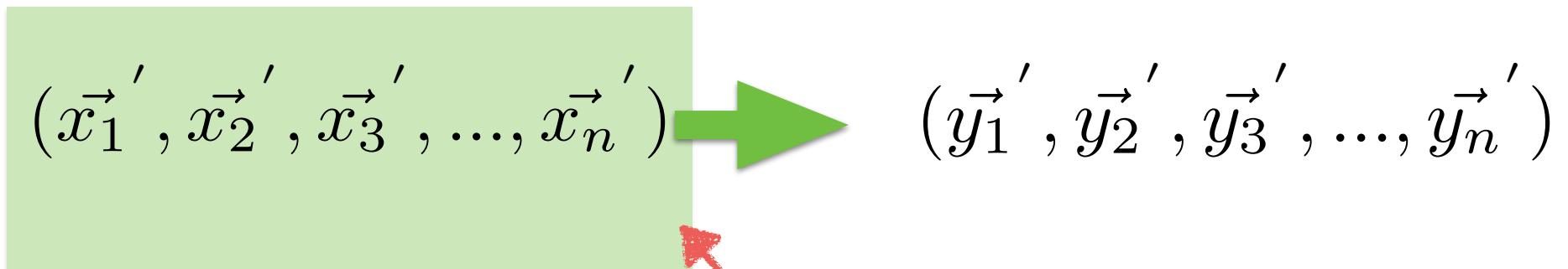


$$f_W(\vec{x}) = \vec{y}$$

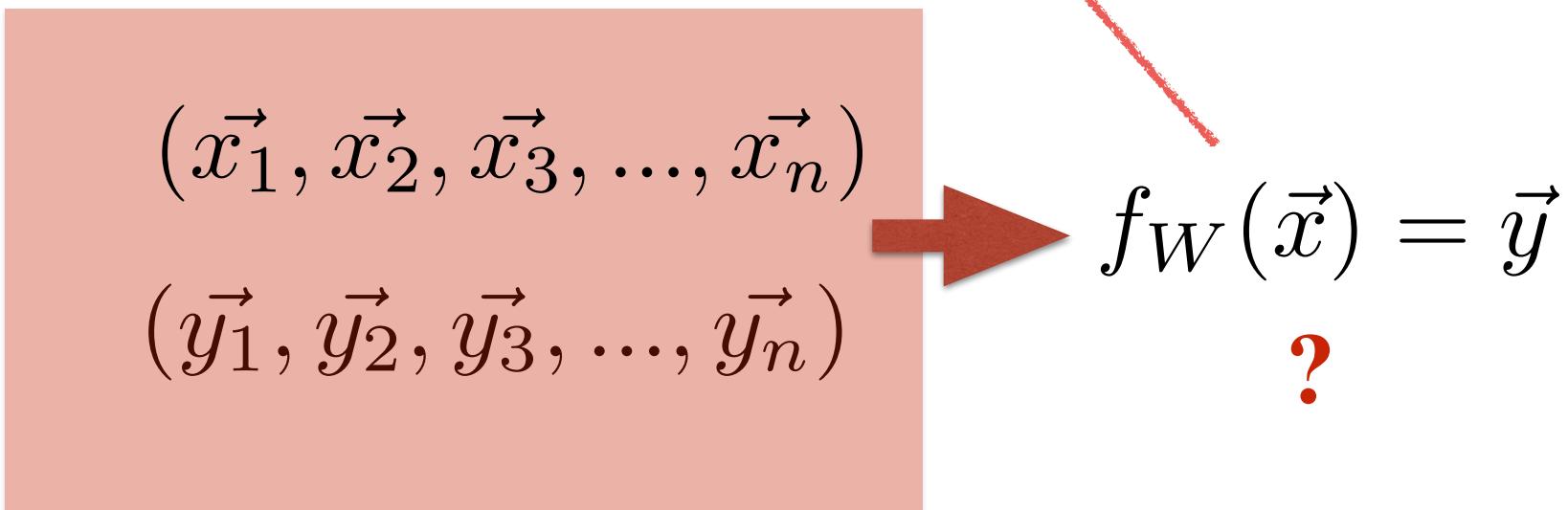
?

SUPERVISED LEARNING

Unlabeled set



Training set



$$(\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n)$$

$$\vec{x} \in \mathbb{R}^d$$

$$(\vec{y}_1, \vec{y}_2, \vec{y}_3, \dots, \vec{y}_n)$$

$$\vec{y} \in \mathbb{R} \quad \vec{y} \in \mathbb{N}$$

GENERAL GOAL: Find a (non-linear) function that outputs the correct class / measurement for a given input object:

$$f_W(\vec{x})$$



Number of parameters - can be large

It is translated into a minimization problem : find \mathbf{W} such as the prediction error is minimal over all unseen vectors

Different “classical” supervised machine learning methods

RANDOM FORESTS

CARTS

decision trees

ARTIFICIAL
NEURAL NETWORKS
(DEEP LEARNING)

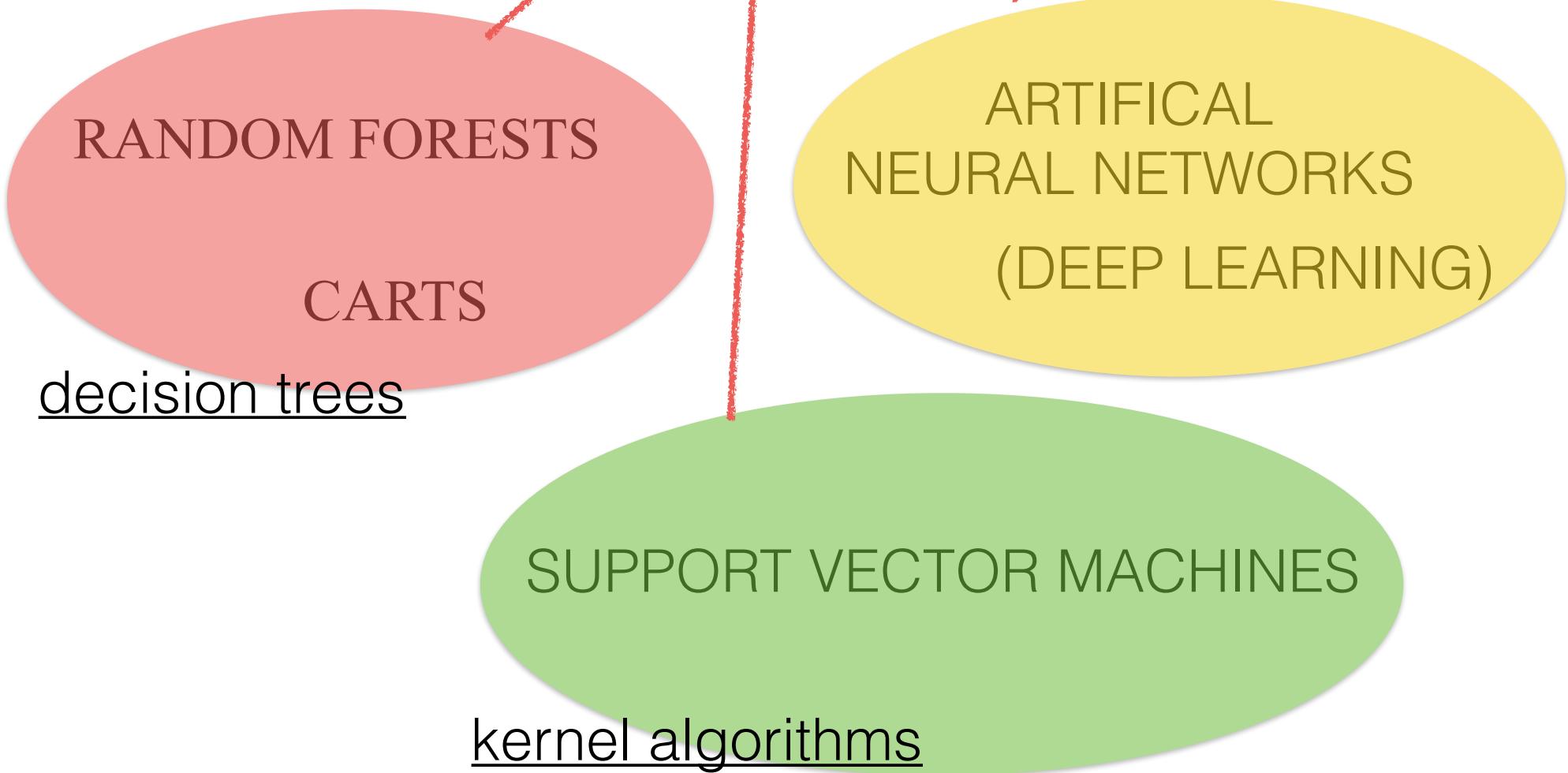
SUPPORT VECTOR MACHINES

kernel algorithms

this is not
classical..

The differences are
in the function
that is used

$$f_W(\vec{x})$$



TO SUMMARIZE: supervised ML **is an optimization problem**

we need therefore 2 key ingredients:

1. A LOSS FUNCTION

**2. A MINIMIZATION OR OPTIMIZATION
ALGORITHM**

TO SUMMARIZE: supervised ML is an optimization problem

we need therefore 2 key ingredients:

1. A LOSS FUNCTION

**2. A MINIMIZATION OR OPTIMIZATION
ALGORITHM**

**THIS IS COMMON TO ALL MACHINE LEARNING
ALGORITHMS**

1. DEFINE A LOSS FUNCTION

$$loss(F_W(\cdot), \vec{x}_i, \vec{y}_i)$$

For example: $(F_W(\vec{x}_i) - \vec{y}_i)^2$ (MSE LOSS FUNCTION)

2. MINIMIZE THE EMPIRICAL RISK WITH OPTIMIZATION

$$\mathcal{R}_{empirical}(W) = \frac{1}{N} \sum_i^N [loss(W, \vec{x}, \vec{y})]$$



MINIMIZE THE RISK

EMPIRICAL RISK?

$$\mathfrak{R}_{\text{empirical}}(W) = \frac{1}{N} \sum_i^N [\text{loss}(W, \vec{x}, \vec{y})]$$

WE ARE MINIMIZING WITH RESPECT TO A FINITE NUMBER OF OBSERVED EXAMPLES

EMPIRICAL RISK?

$$\mathfrak{R}_{\text{empirical}}(W) = \frac{1}{N} \sum_i^N [\text{loss}(W, \vec{x}, \vec{y})]$$

WE ARE MINIMIZING WITH RESPECT TO A FINITE NUMBER OF OBSERVED EXAMPLES

OBSERVED DATASET



EMPIRICAL RISK?

$$\mathfrak{R}_{\text{empirical}}(W) = \frac{1}{N} \sum_i^N [\text{loss}(W, \vec{x}, \vec{y})]$$



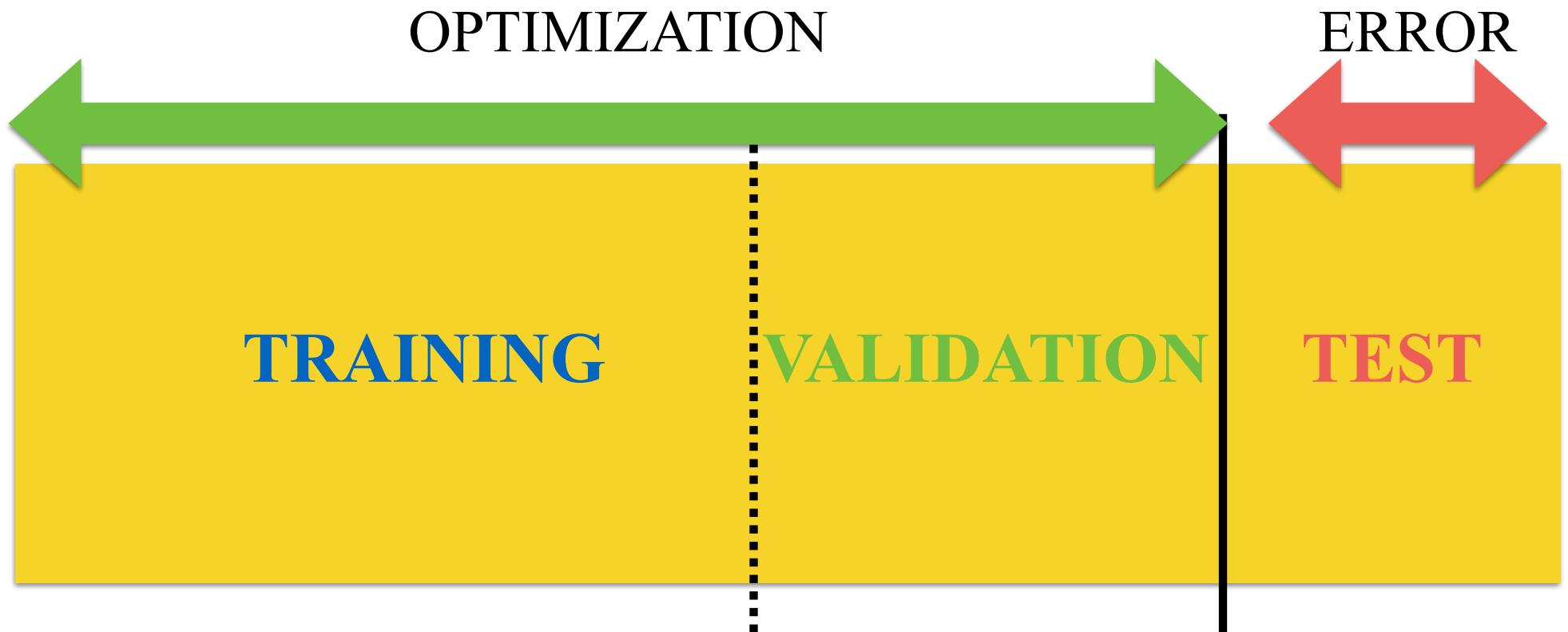
WE ARE MINIMIZING WITH RESPECT TO A FINITE NUMBER OF OBSERVED EXAMPLES

ALL “GALAXIES IN THE UNIVERSE”

OBSERVED DATASET



IN PRACTICE

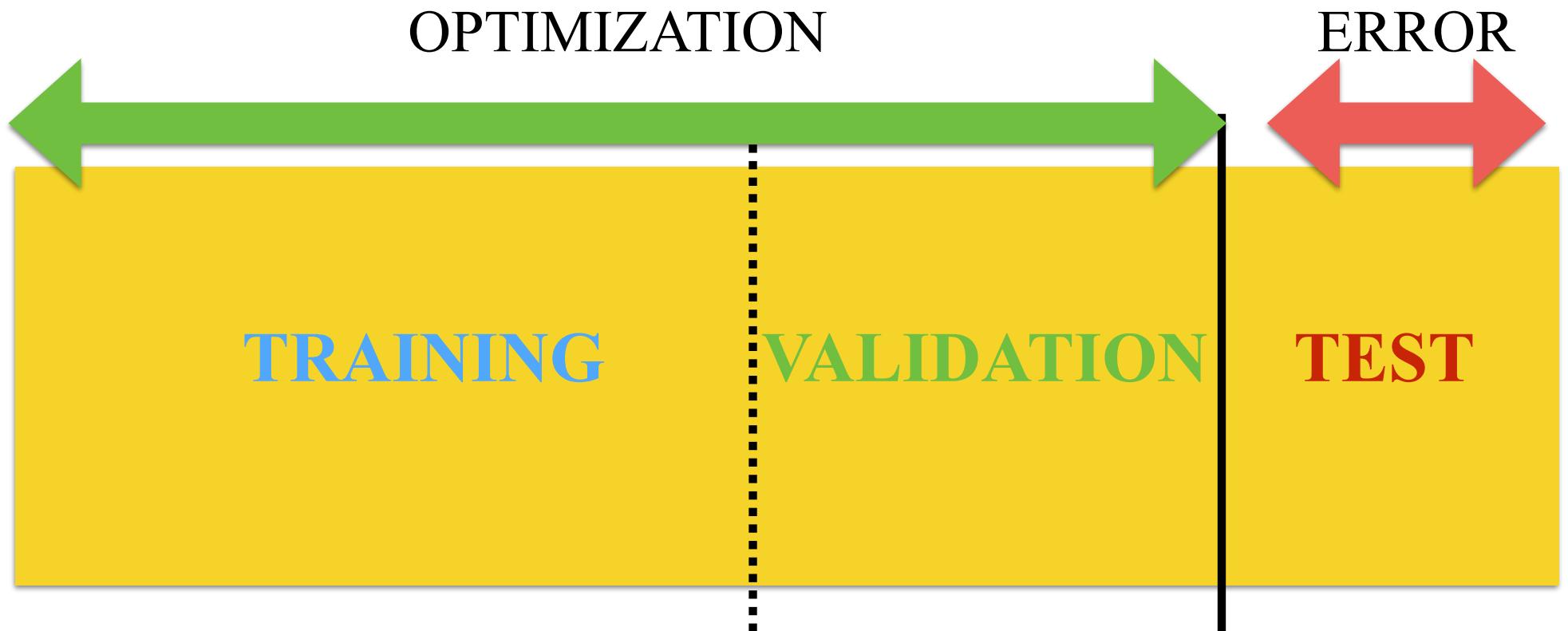


training set: use to train the classifier

validation set: use to monitor performance in real time - check
for overfitting

test set: use to train the classifier

IN PRACTICE



NO CHEATING! NEVER USE TRAINING TO VALIDATE
YOUR ALGORITHM!

The algorithm used to minimize is
called **OPTIMIZATION**

THERE ARE SEVERAL OPTIMIZATION TECHNIQUES

Optimization

THERE ARE SEVERAL OPTIMIZATION TECHNIQUES

THEY DEPEND ON THE MACHINE LEARNING ALGORITHM

Optimization

THERE ARE SEVERAL OPTIMIZATION TECHNIQUES

THEY DEPEND ON THE MACHINE LEARNING ALGORITHM

NEURAL NETWORKS USE THE GRADIENT DESCENT AS WE
WILL SEE LATER

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

weights to be learned



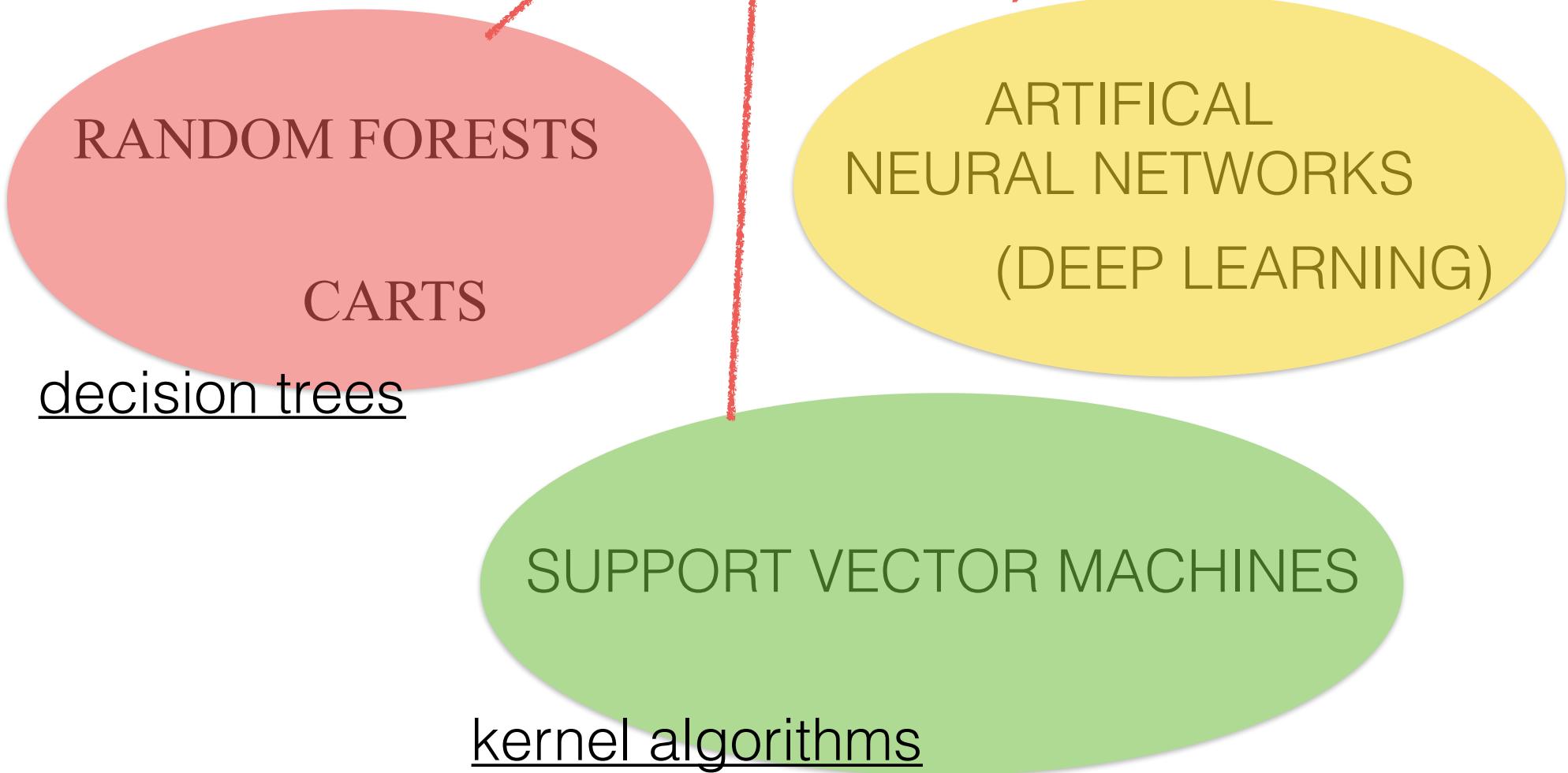
epoch

learning rate
(hyper parameter)

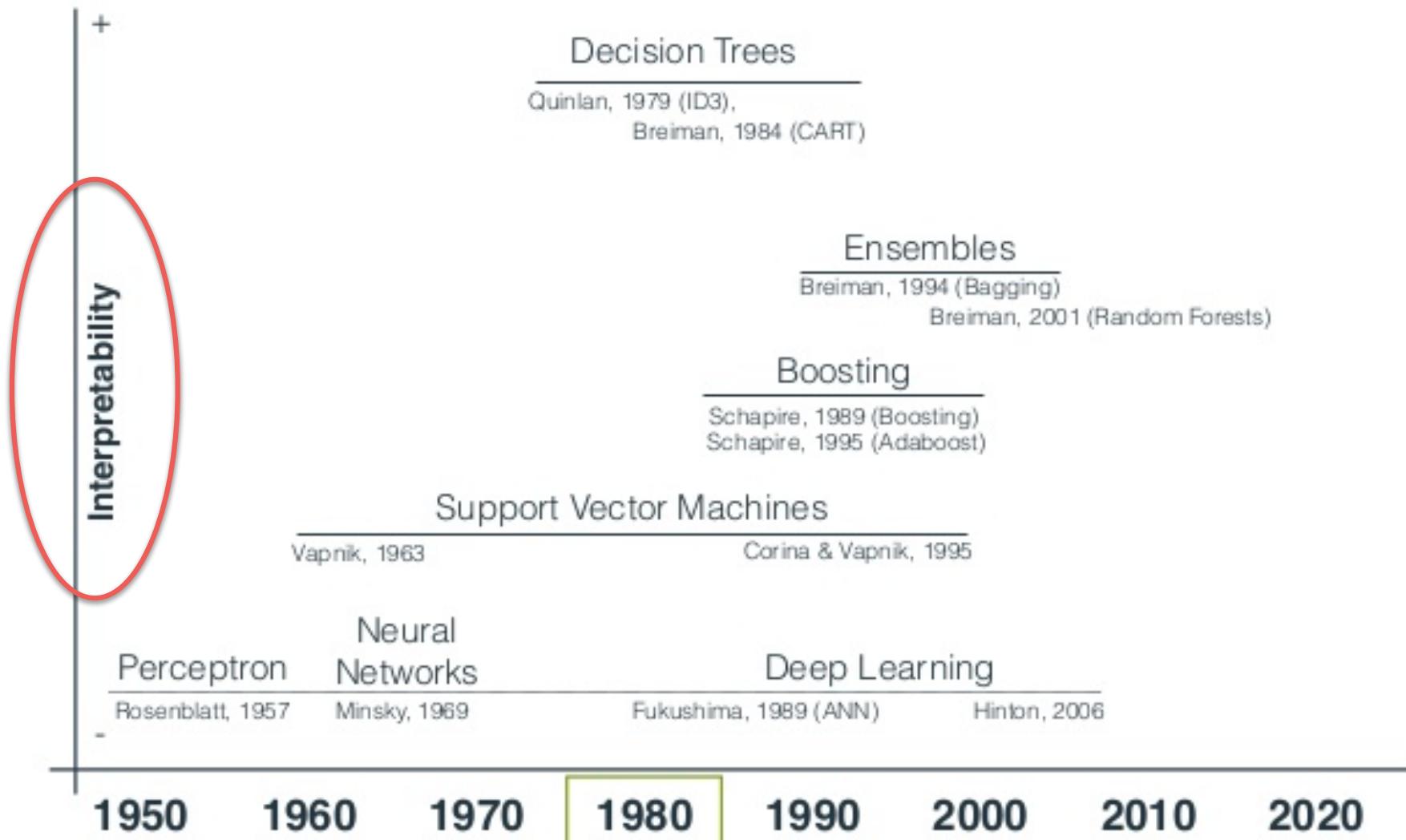


The differences are
in the function
that is used

$$f_W(\vec{x})$$



USING ONE OR ANOTHER METHOD DEPENDS ON YOUR MAIN OBJECTIVES



credit

ALSO INFLUENCED BY “MAINSTREAM” TRENDS

