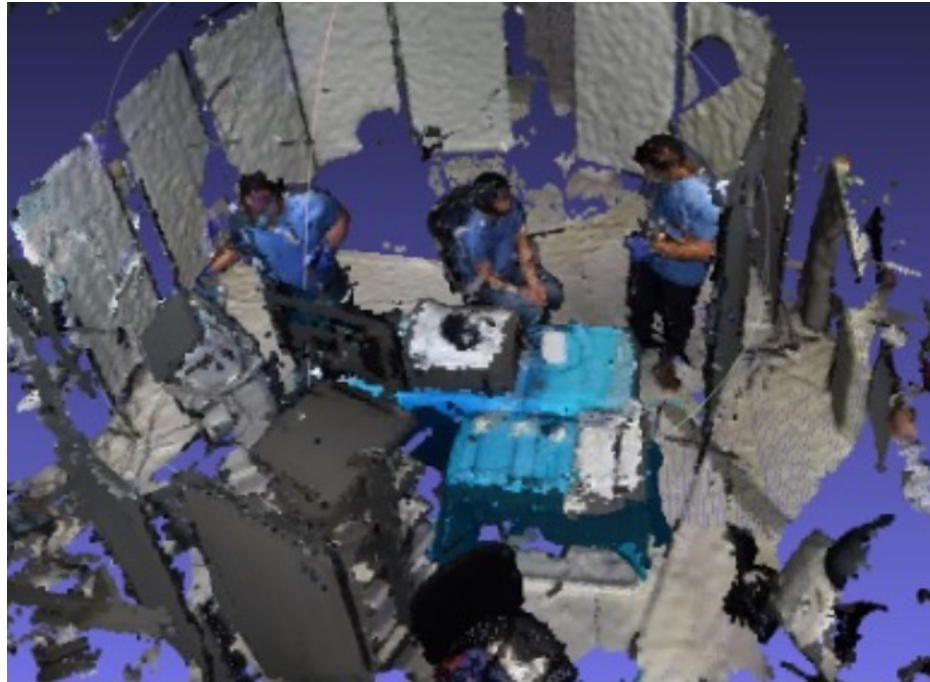# Machine Learning Review

## Tuesday 8h00 – 8h45

Géraldine Conti, Matthew Vowels, Bern Winter School on Machine Learning 2023, Muerren

# Who am I?

- Engineering PhD @ CVSSP, Surrey UK

- Appl. Math / Statistics for Human Sciences PhD UNIL, Switzerland

- Post-doc Sheffield

- Senior Researcher @ The Sense

- Junior Lecturer @ UNIL

- Affiliate member of AI @ Surrey and machine learning CVSSP, University of Surrey

- Research interests:
  - Causal Inference
  - Causal Discovery
  - Semiparametrics
  - Deep Latent Variable Modeling
  - Computer Vision
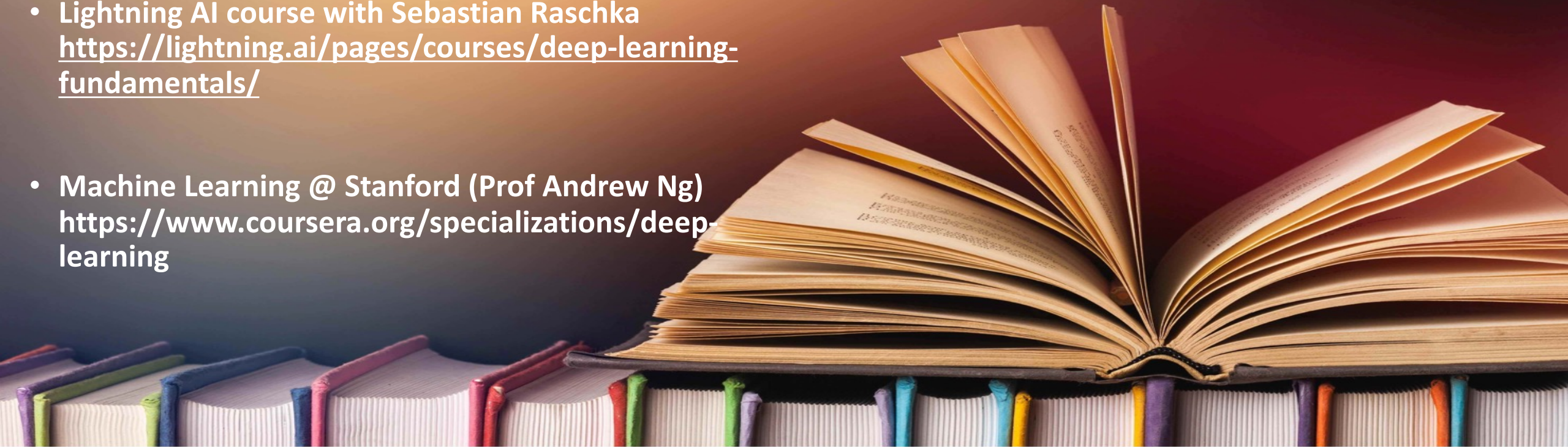  - Multimodal Fusion
  - Algorithmic Finance

Matthew Vowels

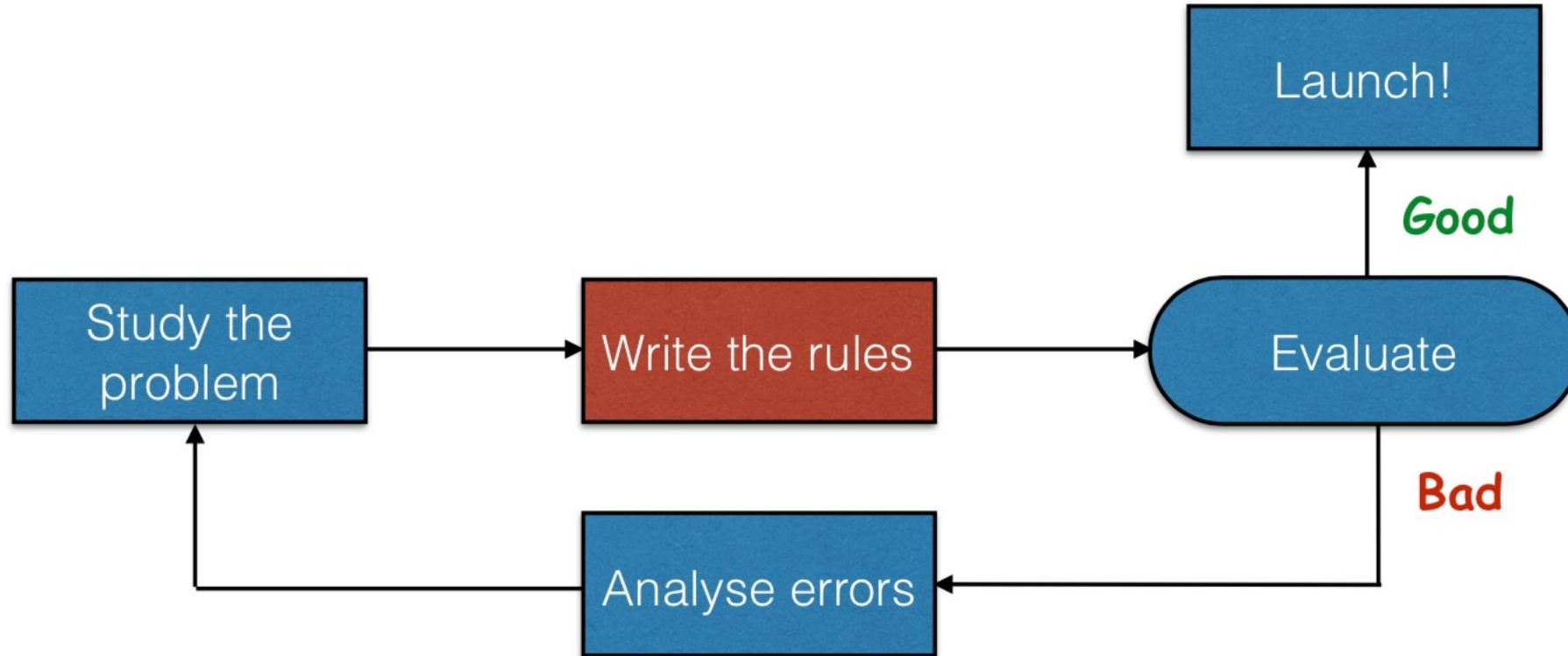# Useful Resources

- **Deep Learning book (Goodfellow, Bengio, Courville)**

- **Lightning AI course with Sebastian Raschka**
**https://lightning.ai/pages/courses/deep-learning-fundamentals/**

- **Machine Learning @ Stanford (Prof Andrew Ng)**
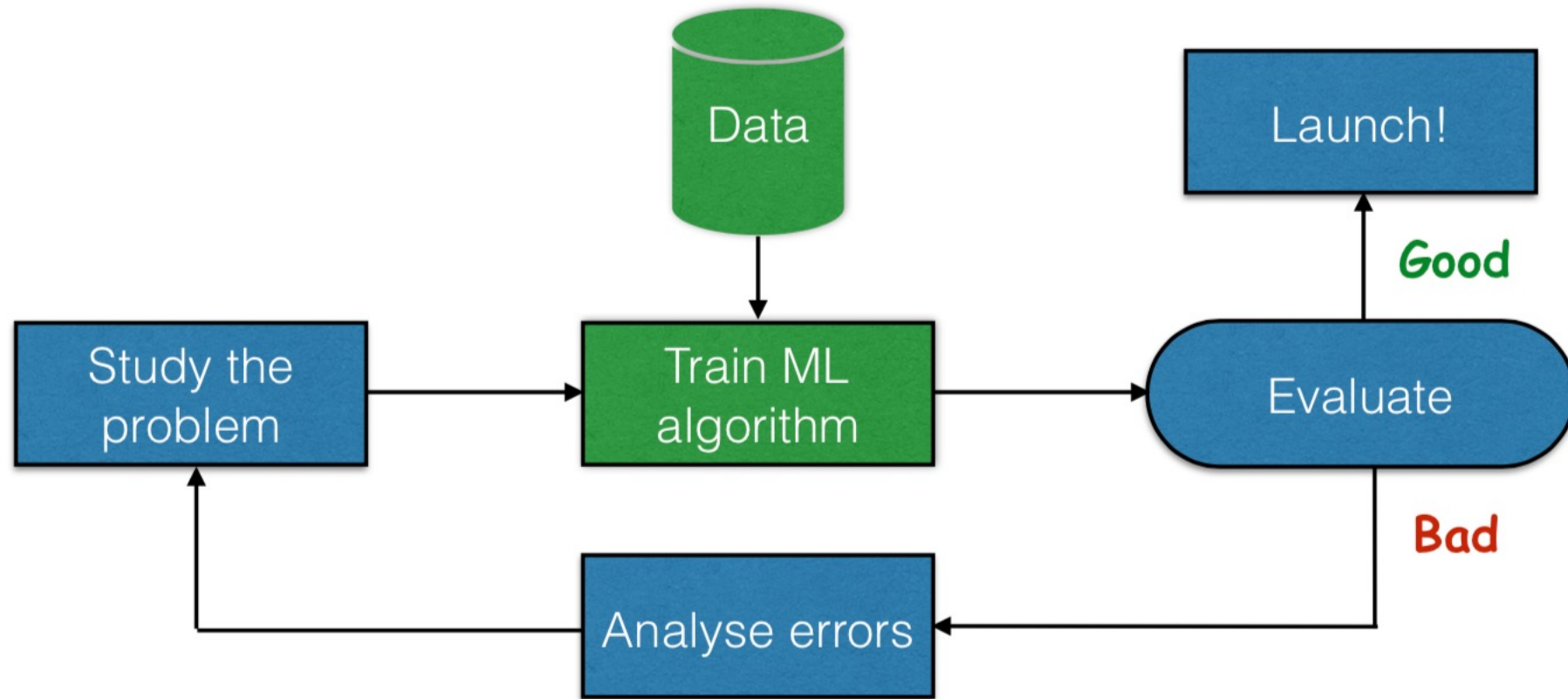**https://www.coursera.org/specializations/deep-learning**
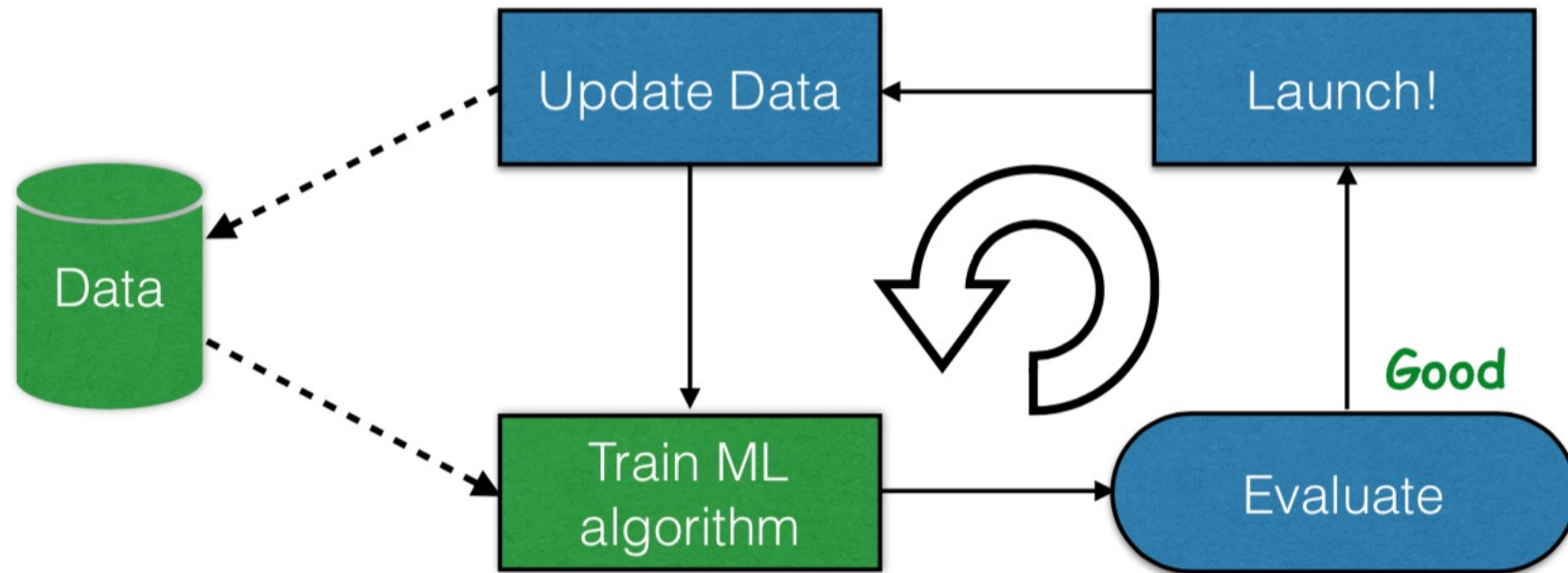
# Traditional approach (Software 1.0)



*List of all the knowledge and formal rules*

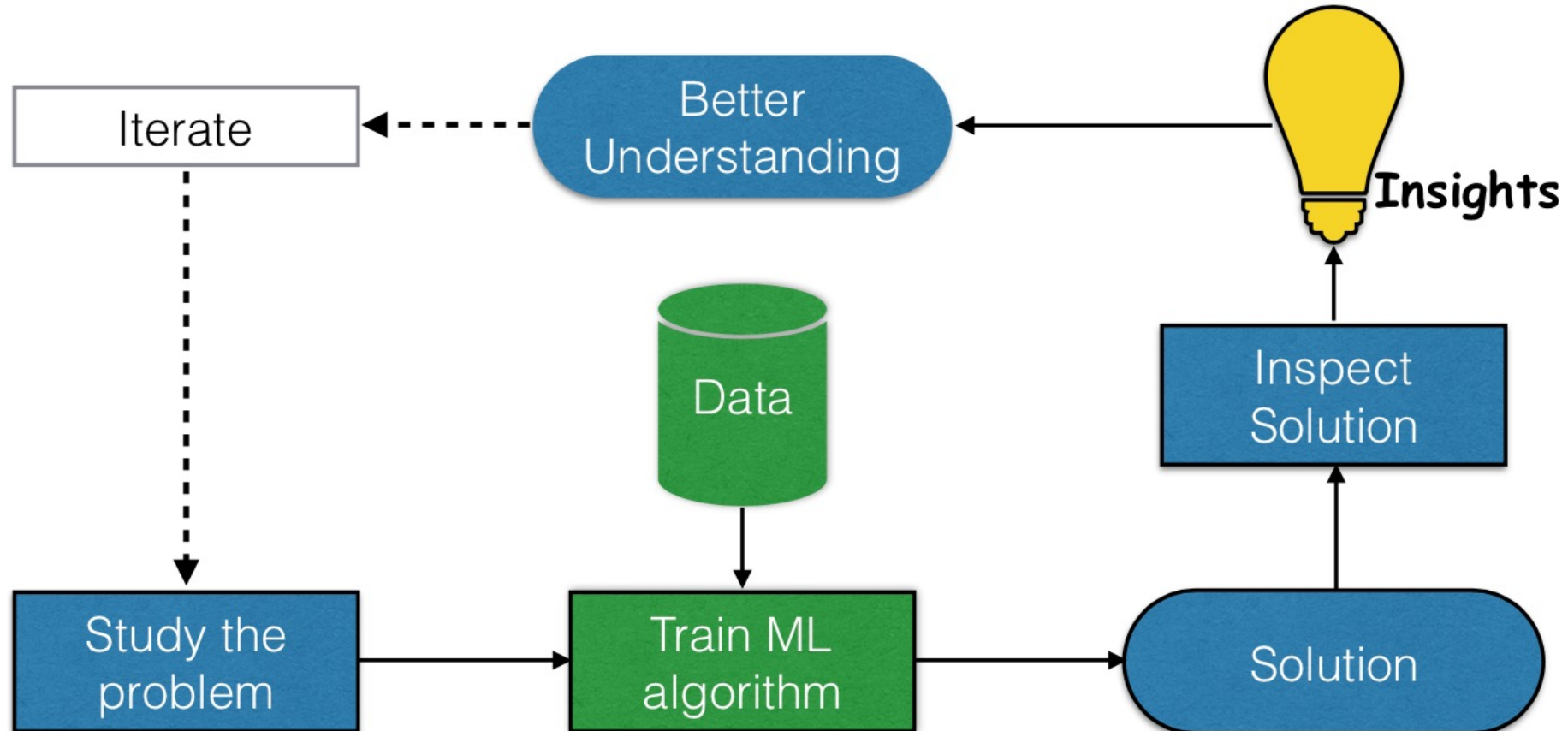# Machine Learning approach (Software 2.0)



*Learning from examples*

# Machine Learning approach (Software 2.0)



*Adapting to change*

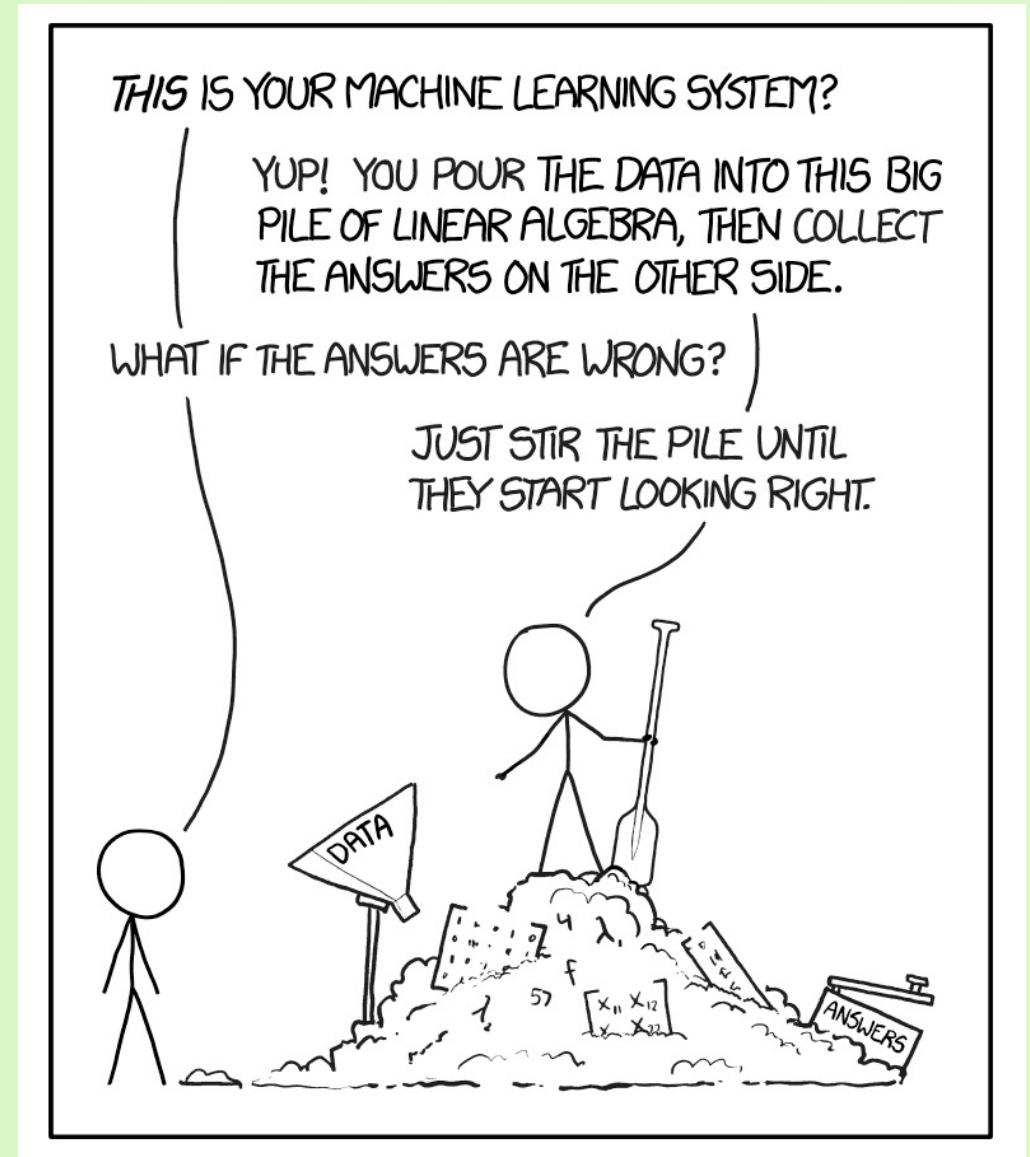# Machine Learning approach (Software 2.0)



*Help Humans learn*

# What is Machine Learning ?

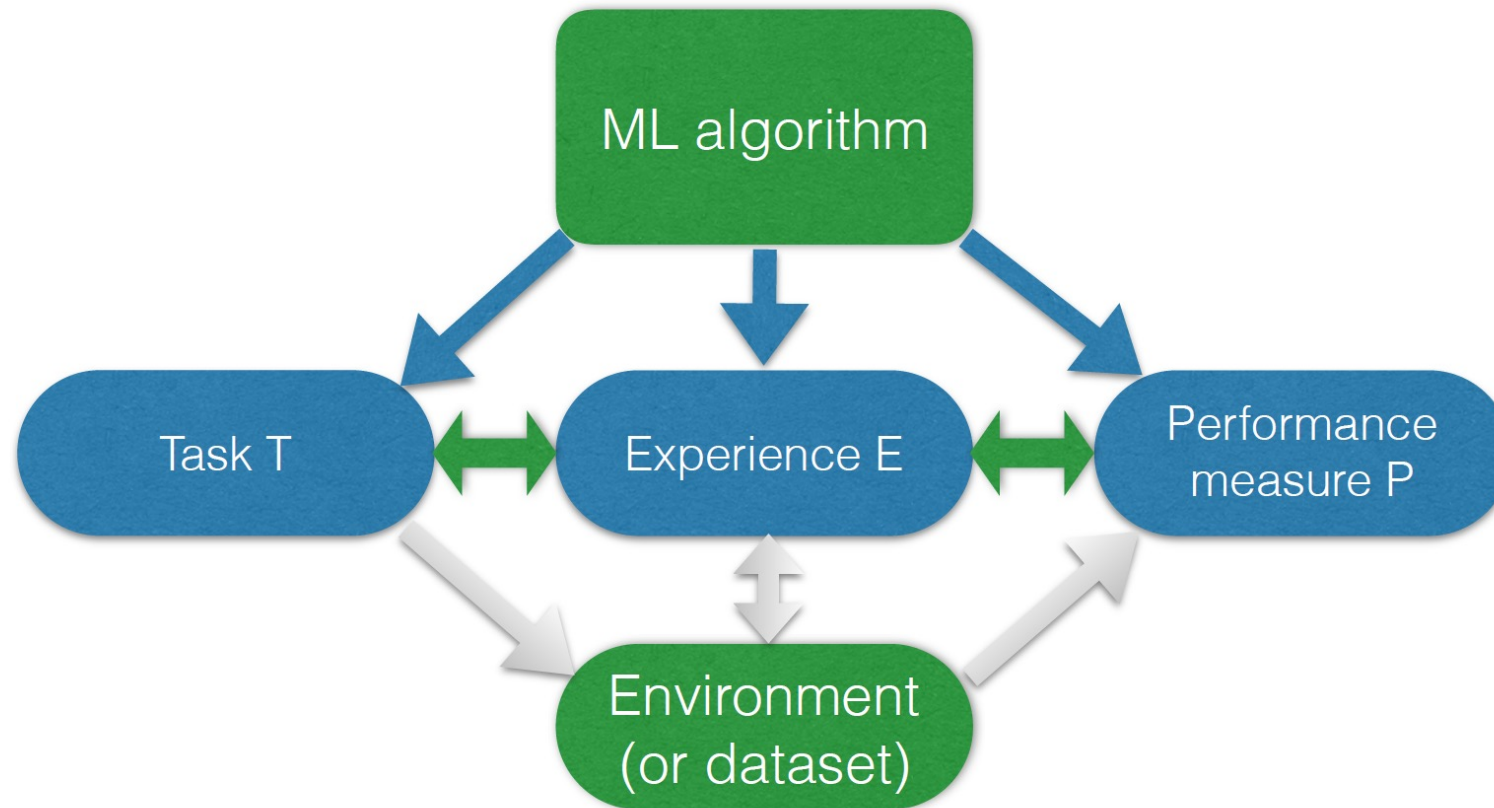*"Can machines do what we (as thinking entities) can do?”*
(Turing)

# Definition… in words

*«  A computer program is said to learn from ==experience E== with respect to some class of ==tasks T== and ==performance measure P==, if its performance at tasks in T, as measured by P, improves with experience E.»*
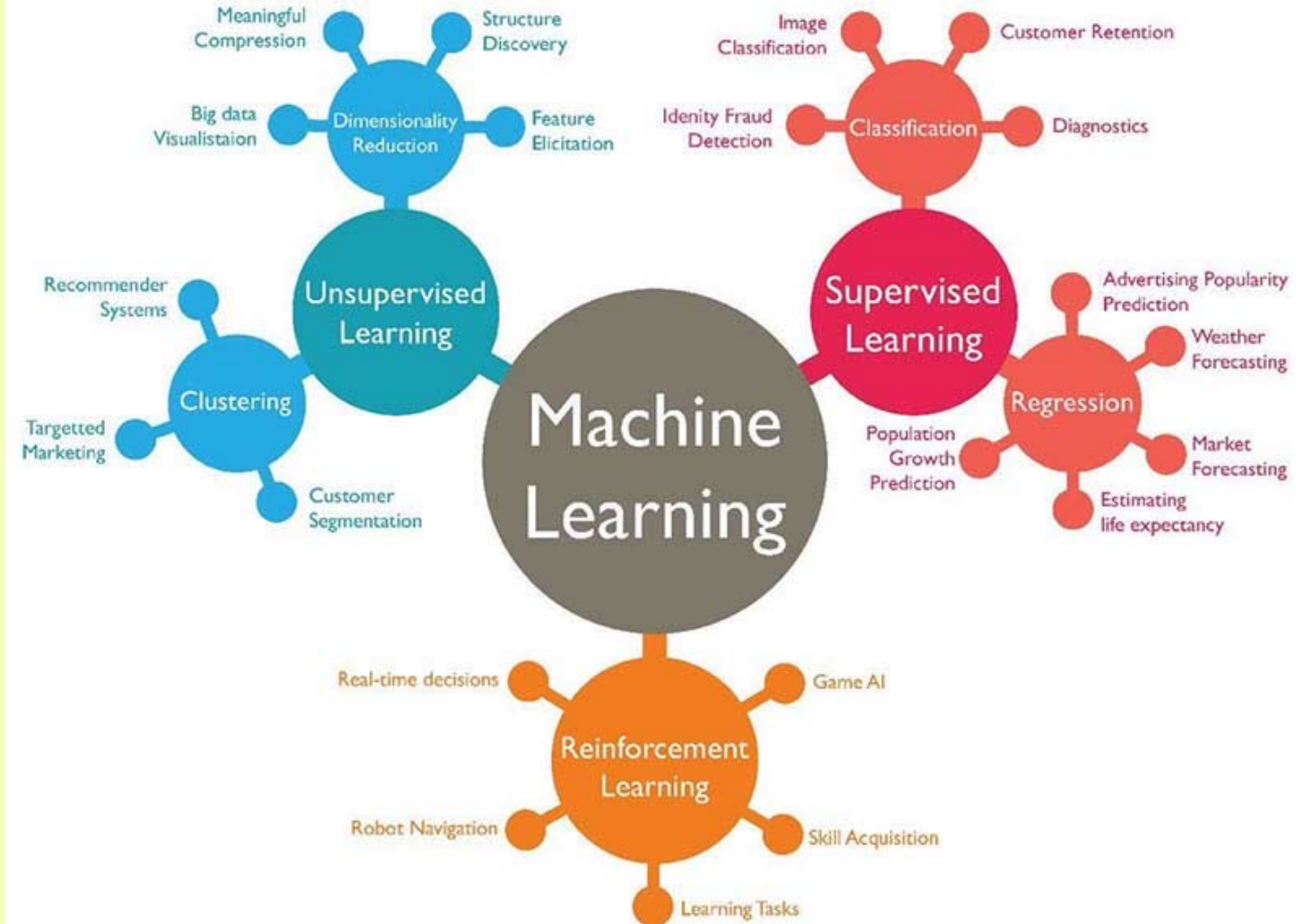
Tom M. Mitchell (1997)

# Definition... schematically

# Experience E

What data to use to solve the task

**Learning Pillars :** How much information is given to the ML algorithm

# Learning Pillars

Supervised
Learning
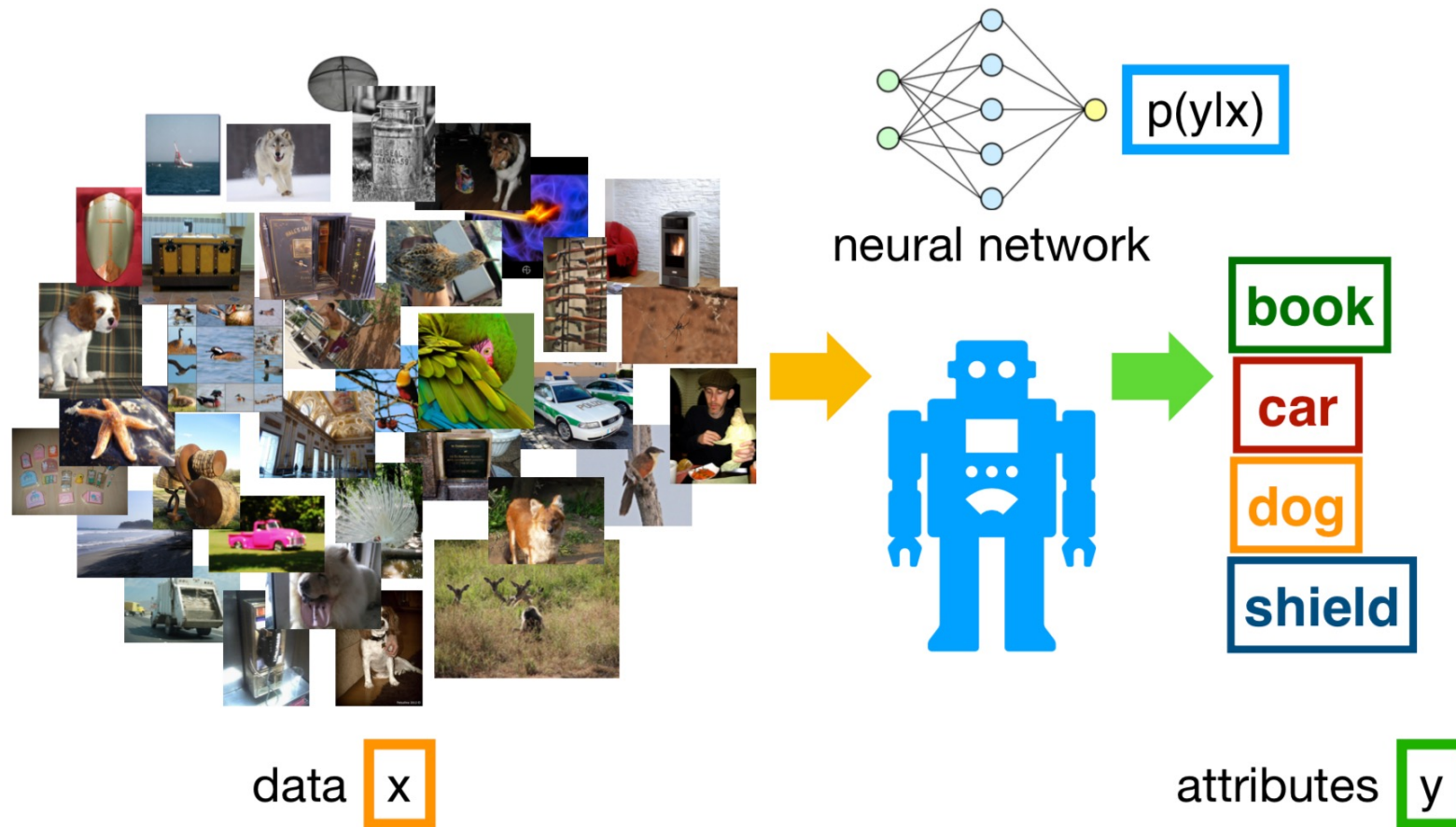
Unsupervised
Learning

Reinforcement
Learning

# Supervised Learning

- Prediction of an output **y** given an input **x**



neural network

$p(y|x)$

data x

attributes y

book
car
dog
shield

# Unsupervised Learning

- Find a *suitable data representation*
  - Preserving all task-relevant information
  - Simpler than the original data and easier to use

Low-dimensional

Sparse

Independent

# Reinforcement Learning



internal state

reward

environment

action

learning rate $\alpha$
inverse temperature $\beta$
discount rate $\gamma$

observation

# Data assumption

- <mark>IID</mark> (independent and identically distributed)

1) Come from the *same distribution*

$$p_{x^{(i)}}(x) = p_{x^{(j)}}(x)$$

2) Are *independent*

$$p\left(x^{(1)}, \ldots, x^{(m)}\right) = \prod_{i=1}^{m} p\left(x^{(i)}\right)$$

# Features

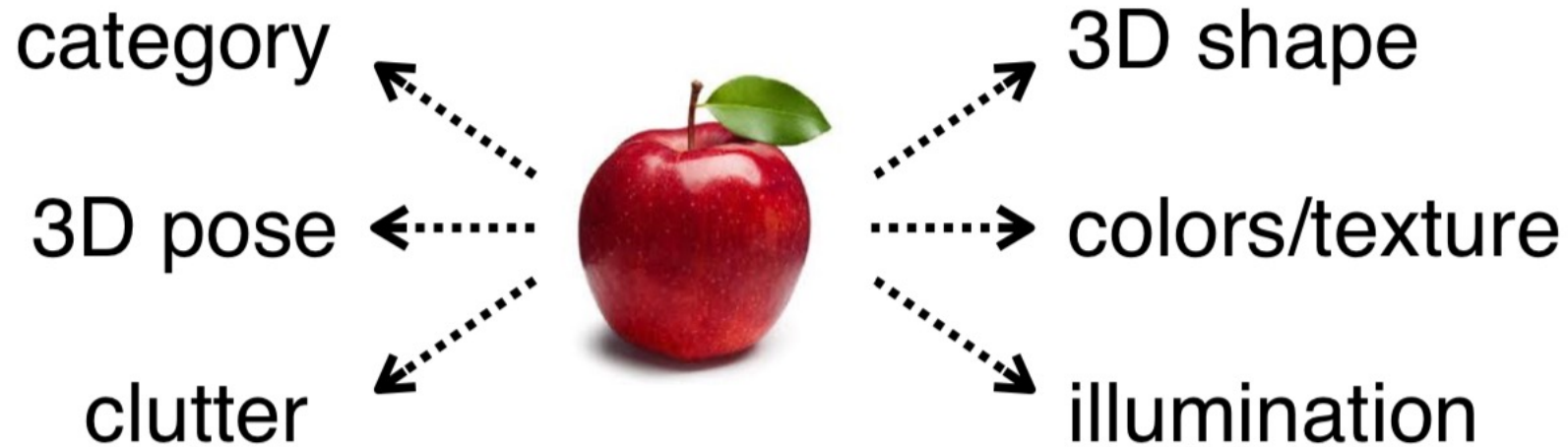- Data often encoded into more focused relevant information (**features** or **internal representation**) to simplify the decision

$$x \rightarrow \phi(x)$$

data ⋰ feature

category ← 🍎 → 3D shape

3D pose ← 🍎 → colors/texture

clutter ← 🍎 → illumination

# Features Example :Image classification

input

feature representation

Learning Algorithm

headlights
tires

pixel 2

pixel 1

× cars
○ non-cars

tires

headlights

# Deep Learning

*« Build a machine that can learn from experience and understand the world as a hierarchy of concepts »*

# Training/Validation/Test sets

- Separate the data into 2(3) sets
  - Training set for training

  - Development / Validation set to find the best parameters

  - (Test set to estimate the performance)

- Separation depends on size of the dataset

- Make sure no algorithmic decisions are being made using data which are also being used to test the algorithm

# Training/Validation/Test sets

# Training/Validation/Test sets

- See/consider also
  - K-fold cross-validation
  - Leave-one-out cross-validation (k=n)
  - Nested cross-validation

# Task T

Find the function f that satisfies f(x) = y using the *training set*

# Problem Types

```
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│   Supervised    │      │  Unsupervised   │      │  Reinforcement  │
│    Learning     │      │    Learning     │      │    Learning     │
└─────────────────┘      └─────────────────┘      └─────────────────┘
```
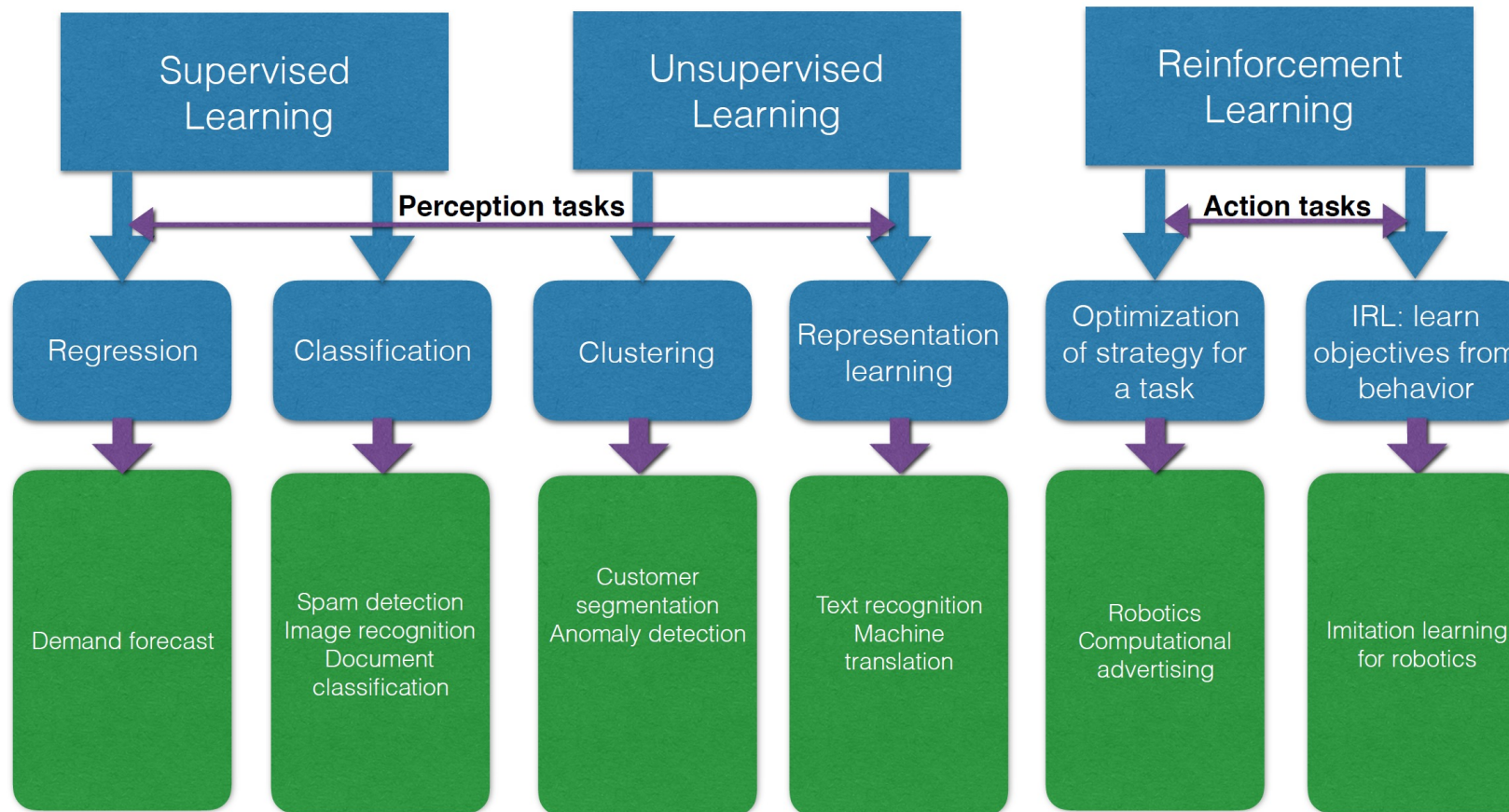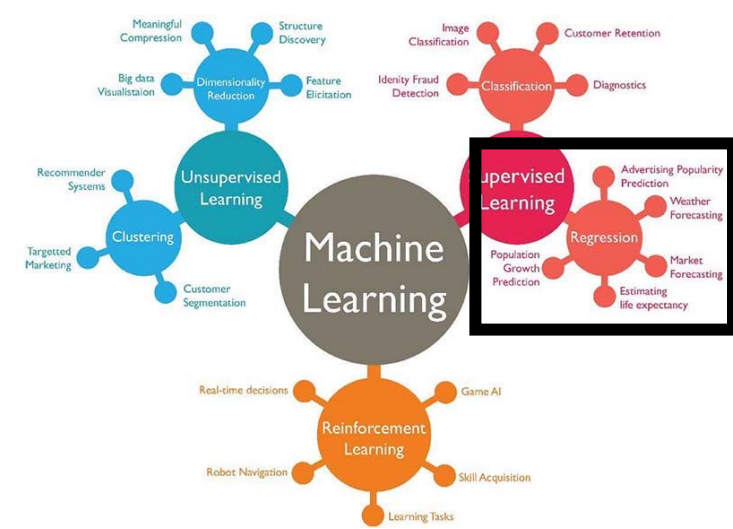
**Perception tasks**                    **Action tasks**

| Regression | Classification | Clustering | Representation learning | Optimization of strategy for a task | IRL: learn objectives from behavior |
|---|---|---|---|---|---|
| Demand forecast | Spam detection Image recognition Document classification | Customer segmentation Anomaly detection | Text recognition Machine translation | Robotics Computational advertising | Imitation learning for robotics |

# Regression

Predict results within *a continuous output*



$82000
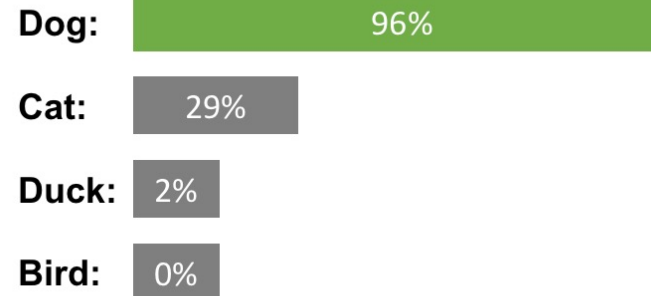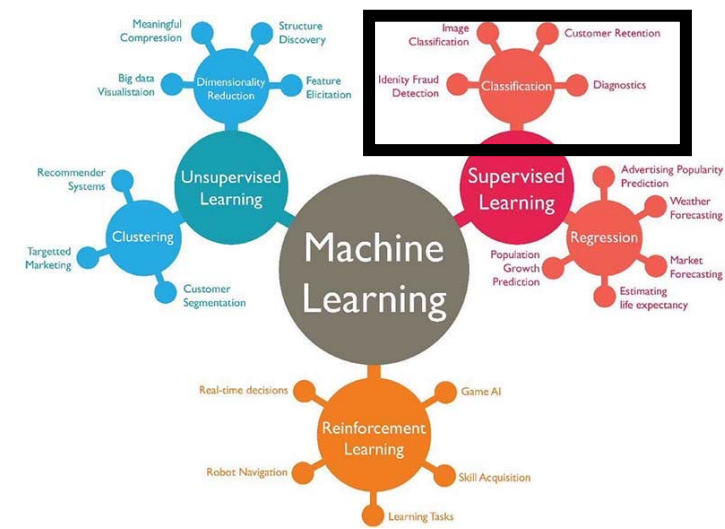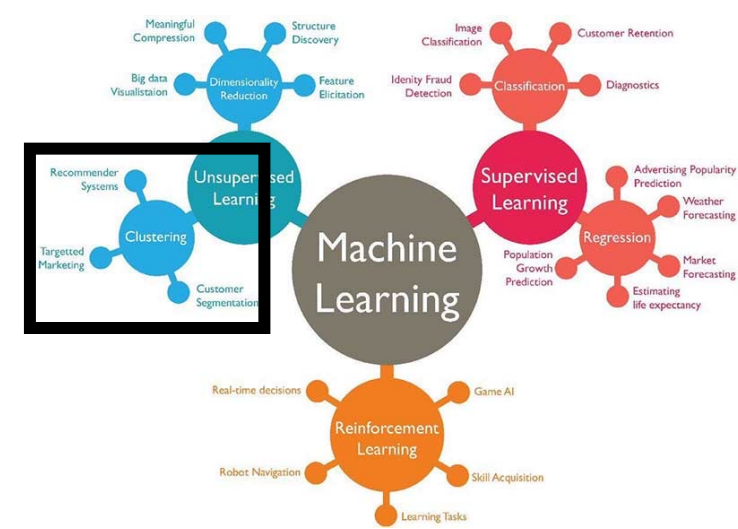
$55500

???

# Classification

- **categorize new inputs** as belonging to one of a set of categories → Predict results within *a discrete output (categories)*



Dog: 96%
Cat: 29%
Duck: 2%
Bird: 0%

Dog: 36%
Cat: 94%
Duck: 2%
Bird: 1%

# Clustering

- create a set of categories, for which individual data instances have a set of common or similar characteristics.

# Data Generation

- generate appropriately novel data



Autoencoder

# Data Generation

- generate appropriately novel data



Autoencoder

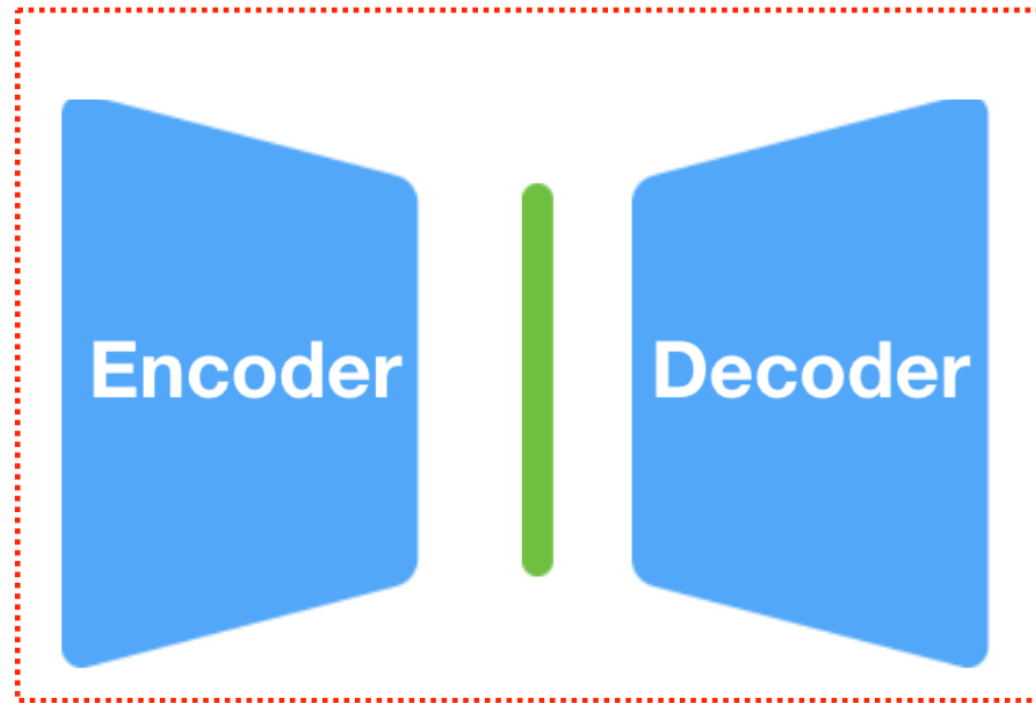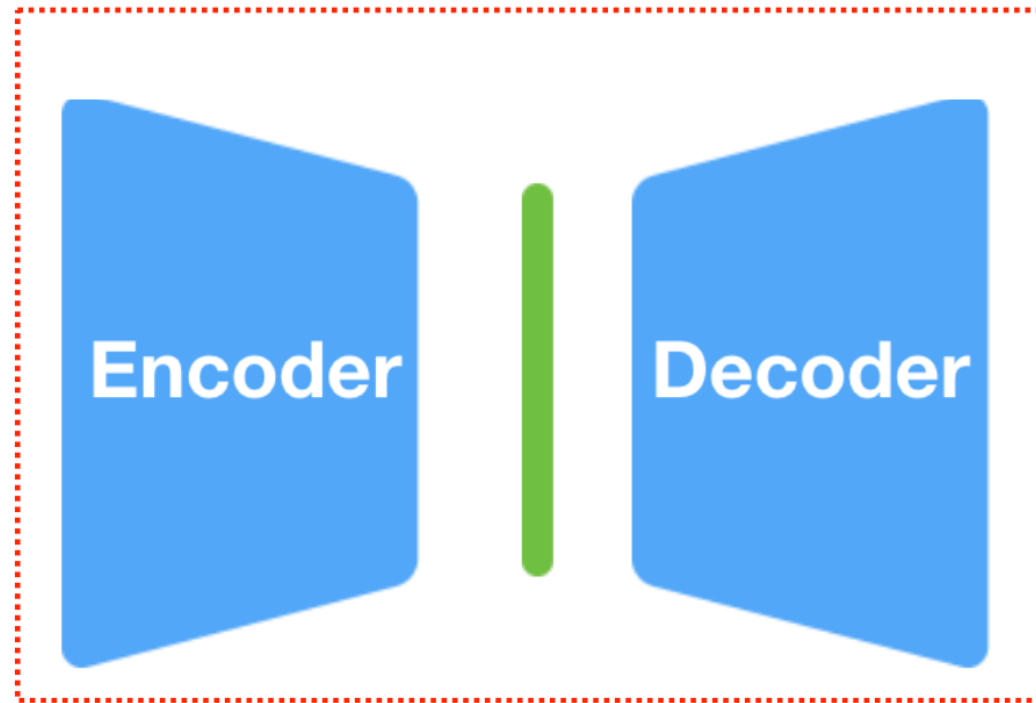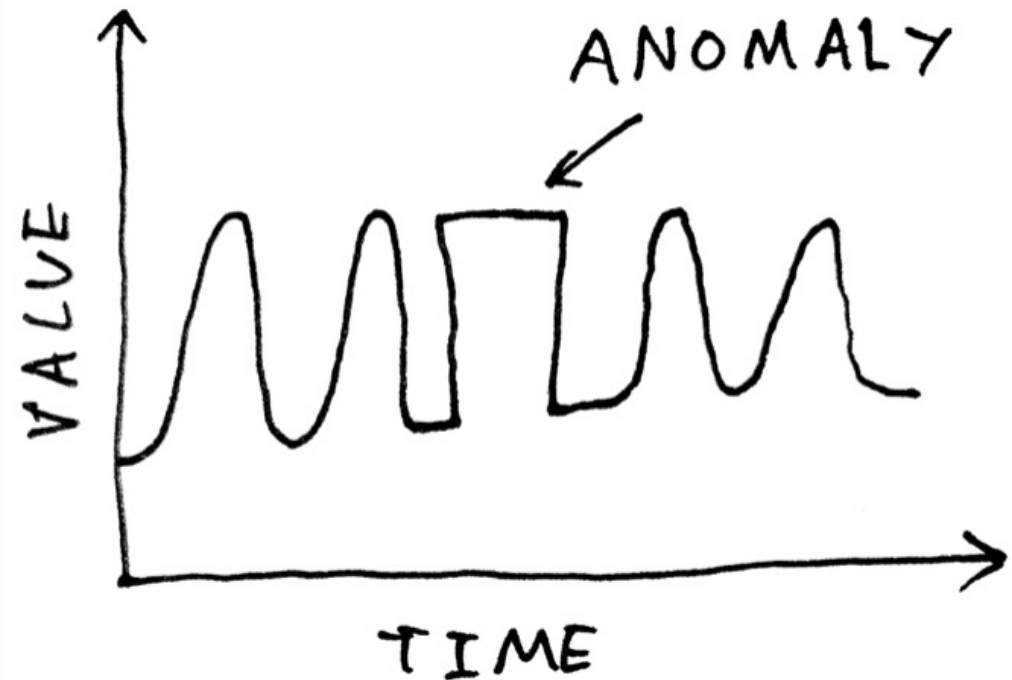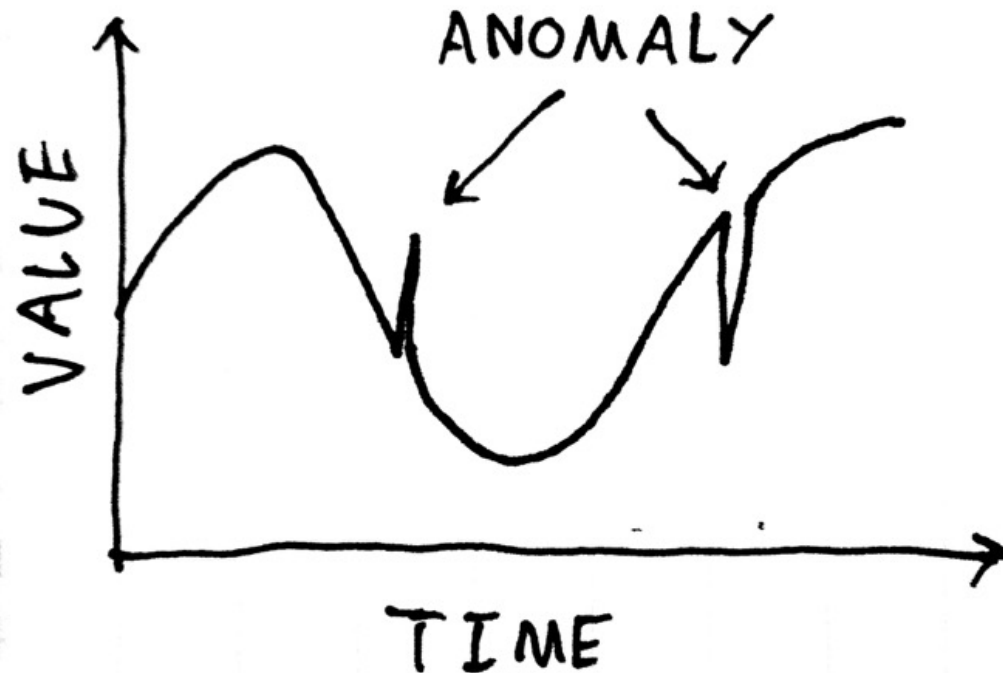# Anomaly Detection

- determine whether specific inputs are
  out of the ordinary.

# Representation Learning



102.4:1 compression
2 principal components

39.4:1 compression
6 principal components

24.4:1 compression
10 principal components

17.7:1 compression
14 principal components

12.5:1 compression
20 principal components

8.4:1 compression
30 principal components

6.3:1 compression
40 principal components

4.2:1 compression
60 principal components

2.8:1 compression
90 principal components

2.1:1 compression
120 principal components

1.7:1 compression
150 principal components

1.4:1 compression
180 principal components

Figure 10: The visual effect of retaining principal components

# Continuous Estimation

- estimate the next numeric value in a sequence (*prediction* for time series data)
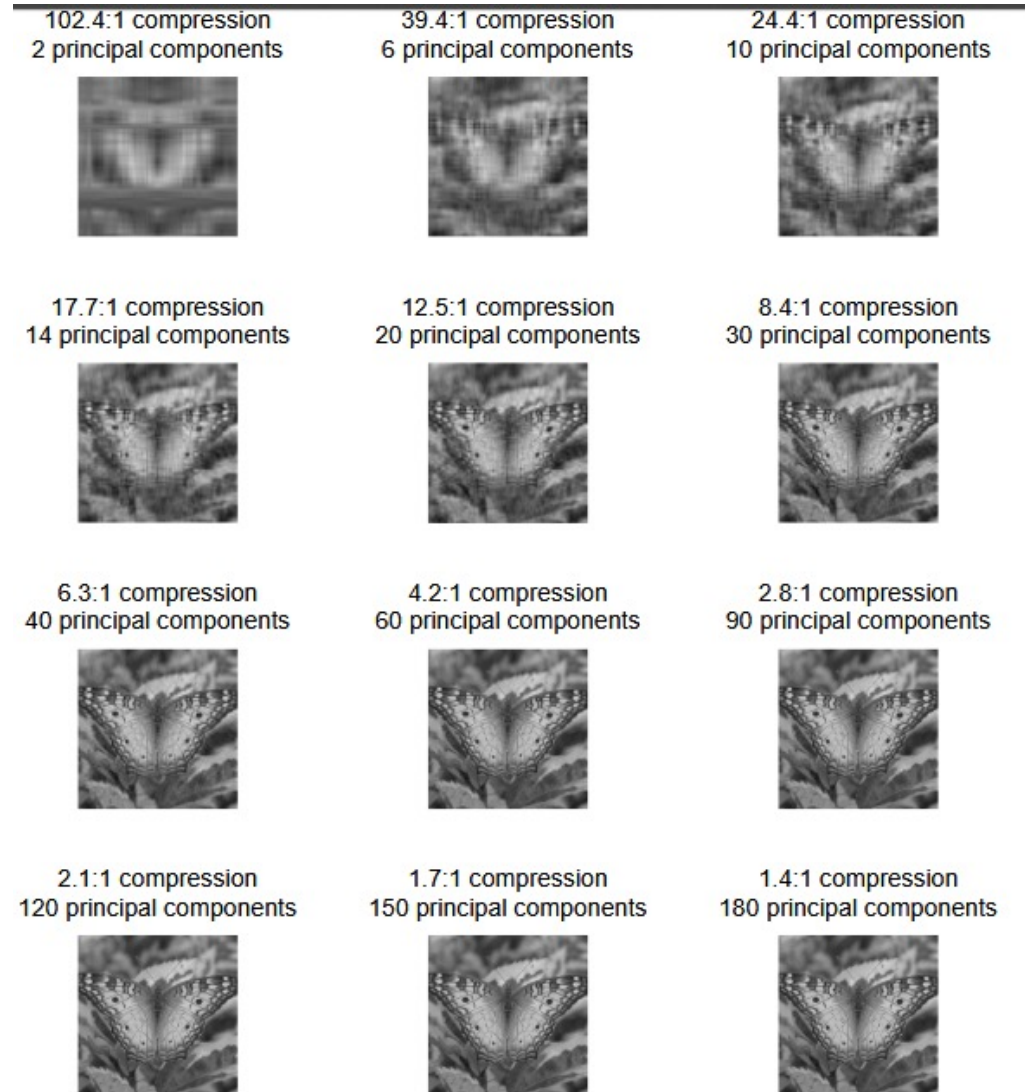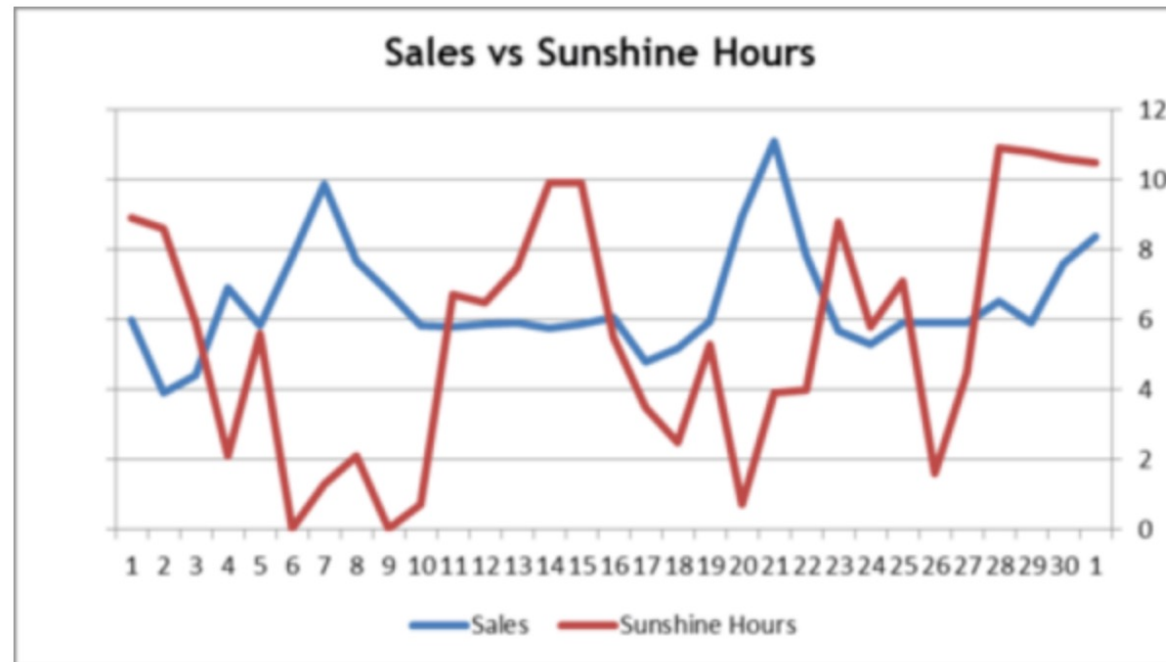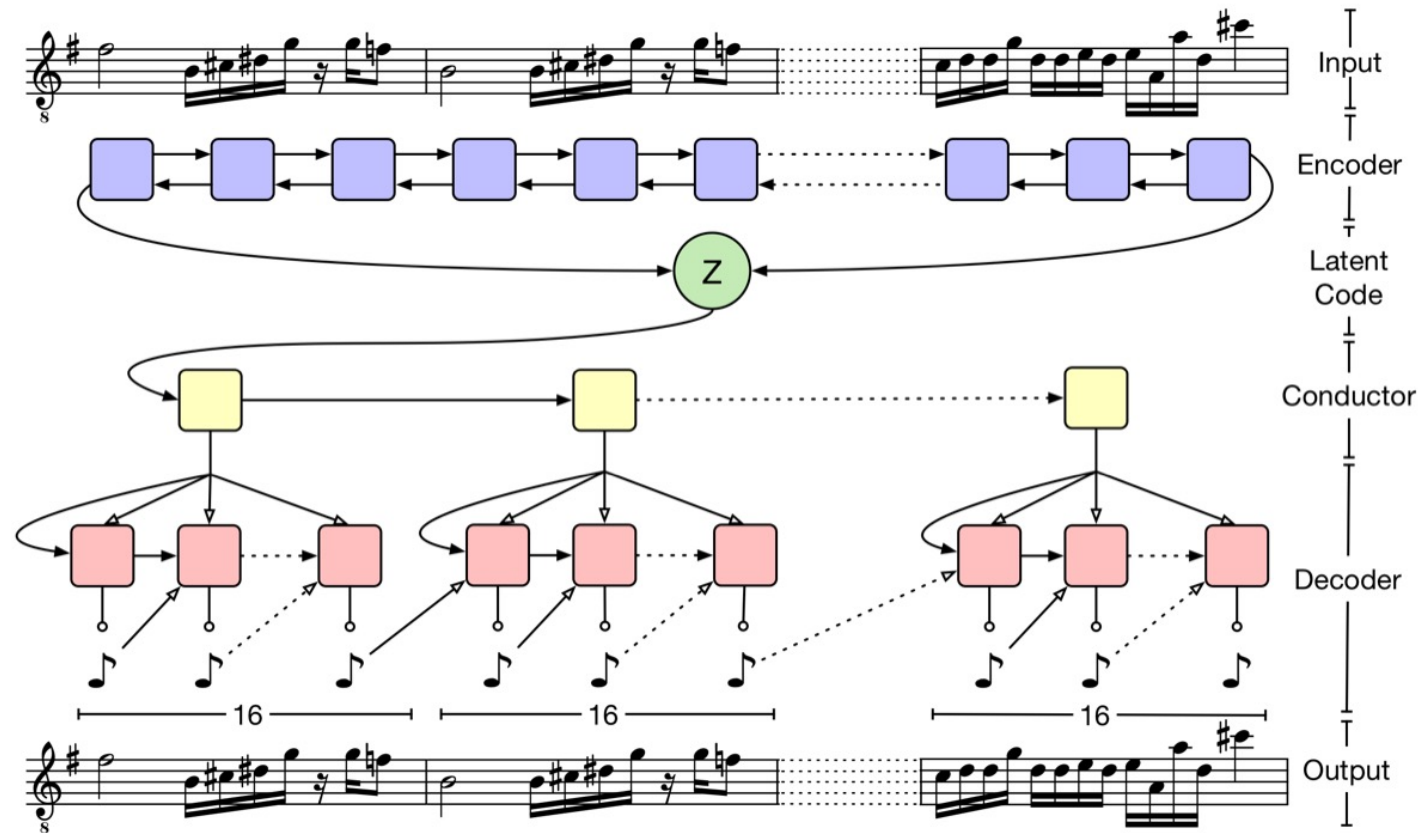
# Data Generation

- generate appropriately <span style="color:blue">novel data</span>

# Ranking

- Used in information retrieval problems
- Used in recommendation systems

# Neural Network (NN)

- Learning algorithm inspired by *how the brain works*

# Deploying a Neural Network

Given a task (in terms of **I/O mappings**), we need :

1) **Network model**

2) **Cost function (/objective/loss function)**

3) **Optimization**



hidden layers

output layer

$h_1$

$h_2$

$x^i$

$z_1$

$z_2$

$z_3$

$q_1$

$q_2$

$q_3$

$f(x^i; \theta)$

network model

cost function

$\text{loss}\left(y^i, f(x^i; \theta)\right)$

$y^i$

training set

optimization

# Network model



sequential processing

$x$  $h_1$  $h_2$  $z_1$  $z_2$  $z_3$  $q_1$  $q_2$  $q_3$  $y$

width

parallel processing

first layer
second layer
output layer

hidden layers

unit (neuron, activation function)

depth

# Activation functions

$$y = f_4(q_1, q_2, q_3)$$

$$y = f_4(f_{3,1}(f_{2,1}(f_{1,1}(x), f_{1,2}(x)), \ldots), \ldots)$$

Hierarchical representation

$$h_1 = f_{1,1}(x)$$
$$h_2 = f_{1,2}(x)$$

$$z_1 = f_{2,1}(h_1, h_2)$$
$$z_2 = f_{2,2}(h_1, h_2)$$
$$z_3 = f_{2,3}(h_1, h_2)$$

$$q_1 = f_{3,1}(z_1, z_2, z_3)$$
$$q_2 = f_{3,2}(z_1, z_2, z_3)$$
$$q_3 = f_{3,3}(z_1, z_2, z_3)$$

Fully connected

$$f_{2,2}(h_1, h_2) = w_1 h_1 + w_2 h_2 + b_{2,2}$$

Weights w and bias b parameters to optimize

# Activation functions

$$y = f_4(q_1, q_2, q_3)$$

$h_1 = f_{1,1}(x)$
$h_2 = f_{1,2}(x)$

$z_1 = f_{2,1}(h_1, h_2)$
$z_2 = f_{2,2}(h_1, h_2)$
$z_3 = f_{2,3}(h_1, h_2)$

$q_1 = f_{3,1}(z_1, z_2, z_3)$
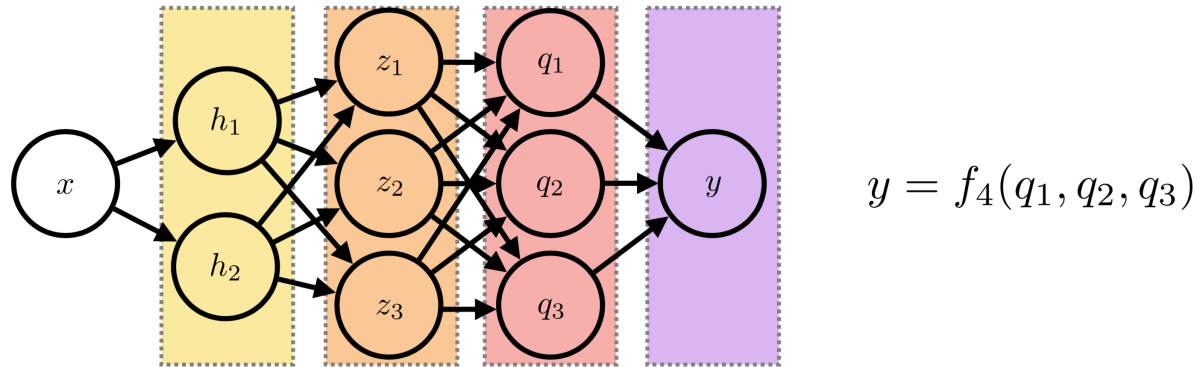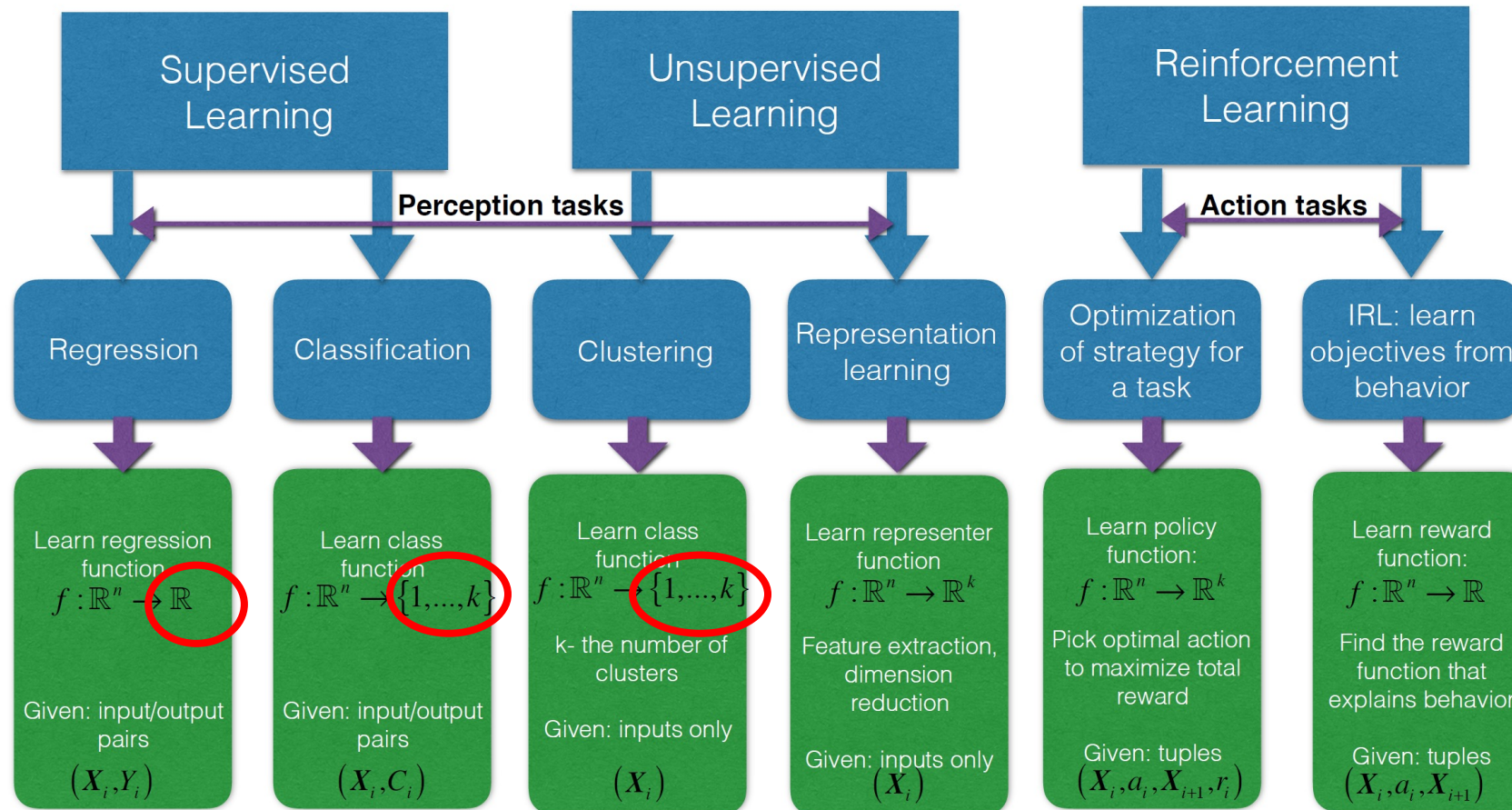$q_2 = f_{3,2}(z_1, z_2, z_3)$
$q_3 = f_{3,3}(z_1, z_2, z_3)$

**(3,1)**    **(3,**$\underline{2}$**) (**$\underline{2}$**,1)**

$$Z = W_Z H$$

Fully connected    $f_{2,2}(h_1, h_2) = w_1 h_1 + w_2 h_2 + b_{2,2}$

# Neural Network Outputs



Supervised Learning

Unsupervised Learning

Reinforcement Learning

**Perception tasks**

**Action tasks**

Regression

Classification

Clustering

Representation learning

Optimization of strategy for a task

IRL: learn objectives from behavior

Learn regression function
$$f : \mathbb{R}^n \to \mathbb{R}$$

Given: input/output pairs
$$(X_i, Y_i)$$

Learn class function
$$f : \mathbb{R}^n \to \{1,...,k\}$$

Given: input/output pairs
$$(X_i, C_i)$$

Learn class function
$$f : \mathbb{R}^n \to \{1,...,k\}$$

k- the number of clusters

Given: inputs only
$$(X_i)$$

Learn representer function
$$f : \mathbb{R}^n \to \mathbb{R}^k$$

Feature extraction, dimension reduction

Given: inputs only
$$(X_i)$$

Learn policy function:
$$f : \mathbb{R}^n \to \mathbb{R}^k$$

Pick optimal action to maximize total reward

Given: tuples
$$(X_i, a_i, X_{i+1}, r_i)$$

Learn reward function:
$$f : \mathbb{R}^n \to \mathbb{R}$$

Find the reward function that explains behavior

Given: tuples
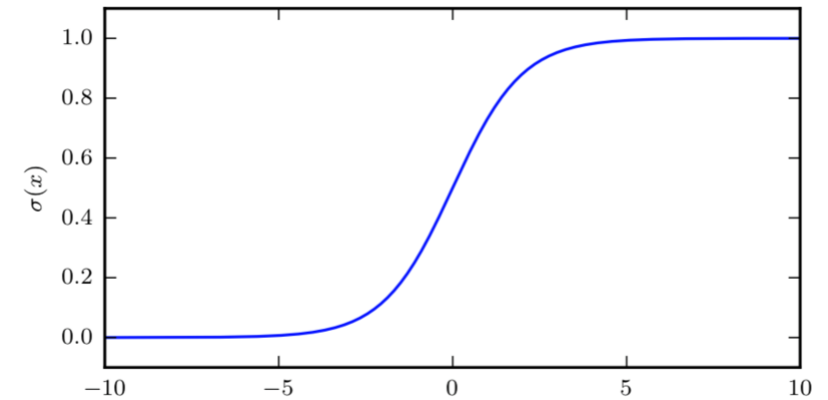$$(X_i, a_i, X_{i+1})$$

# Output layer : activation functions

1) **Classification** : probability vector
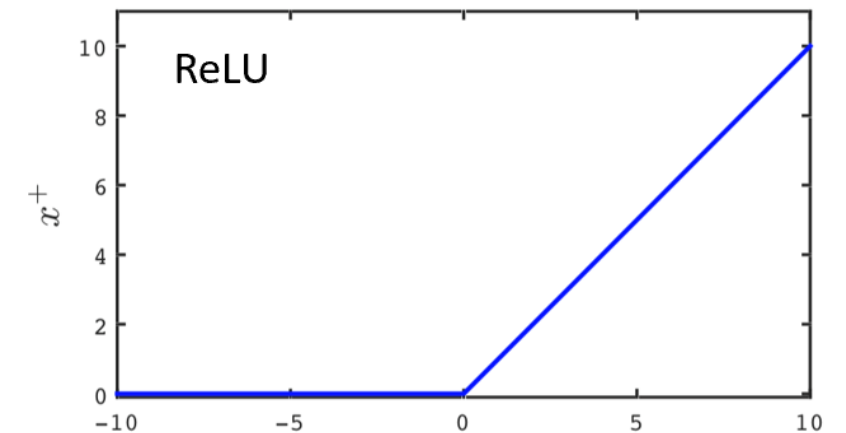   - Sigmoid (binary class)
   - Softmax (multiple class)

$$Z = \sigma(W_z H)$$

**2) Regression** : mean estimate
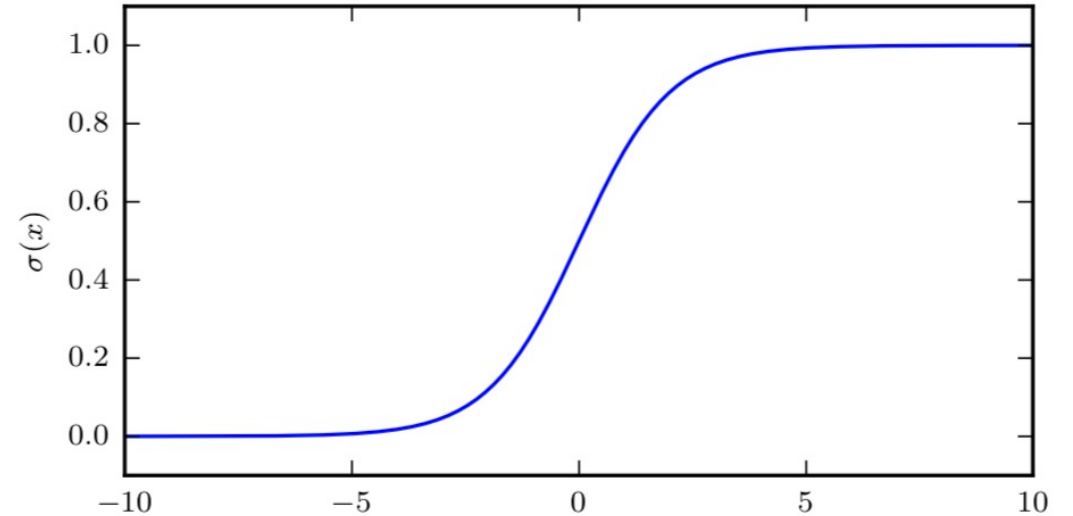   - ReLU
   - Softplus
   - Smoothed max
   - Generalization of ReLU (leaky ReLU,…)

# Sigmoid and softmax

**Sigmoid** (*two-class* classifier) :

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



**Softmax** (*multi-class* classifier) :

$$\text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

# ReLU, softplus and smoothed max
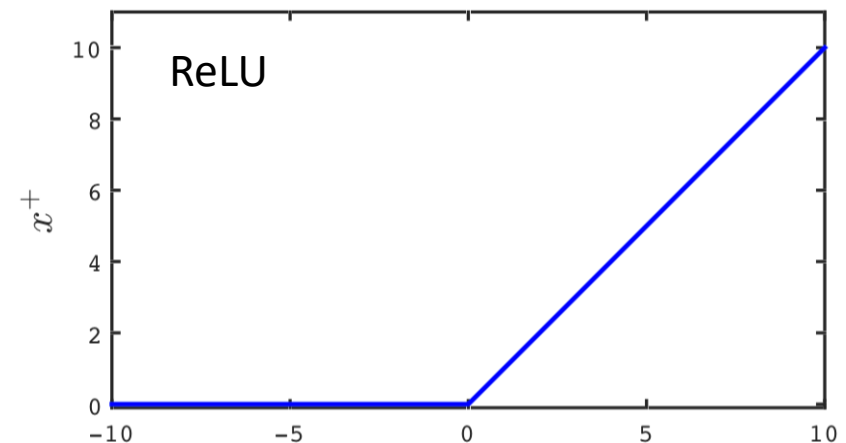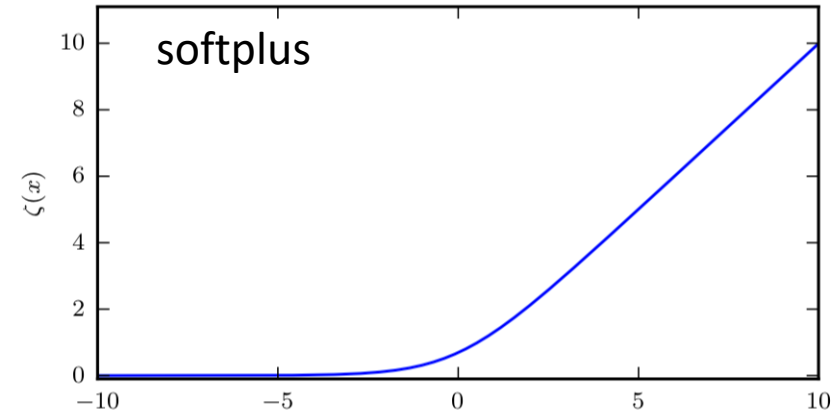
**Softplus** (smooth approx. of ReLU) :

$$\zeta(x) = \log(1 + \exp(x))$$

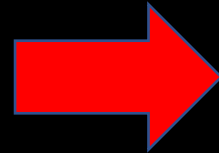**Smoothed max** (*extension* of softplus) :

$$\zeta(x) = \log \sum_j \exp(x_i)$$

**ReLU** (Rectified Linear Unit) :

$$x^+ = \max(0, x)$$



softplus



ReLU

# Let's start playing !

# Tutorial 1: 9h-10h00

## Introduction to Pytorch