

# Introduction

---

Introduction to Deep Learning

# Learning outcomes

After this lesson you will be able to:

- Understand ML foundations
- Understand DL foundations
- Formulate problems from a ML/DL perspective

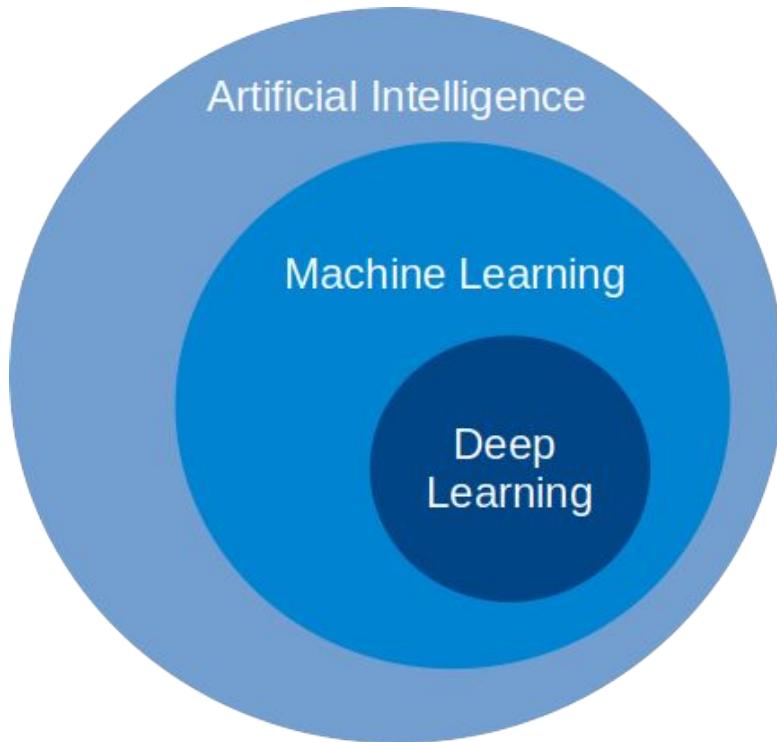
# Introduction



*“Machine Learning is a lot like teenage sex. Everybody talks about it. Only some really know how to do it. Everyone thinks everyone else is doing it. So, everyone claims they’re doing it.”*

— <https://medium.com/@peterx/dcdc08f17d17>

# Where we are



## Artificial Intelligence

Try to create machines that can simulate human behavior

## Machine Learning

Automatically learn from data without explicit programming

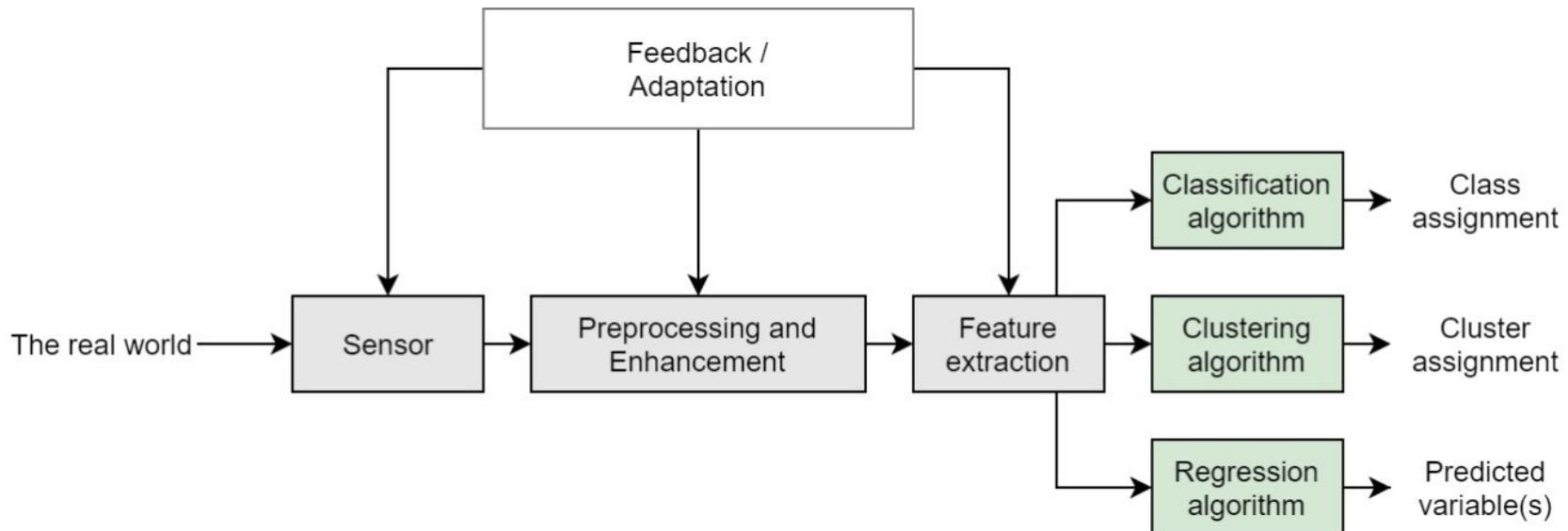
## Deep Learning

Machine learning with neural networks

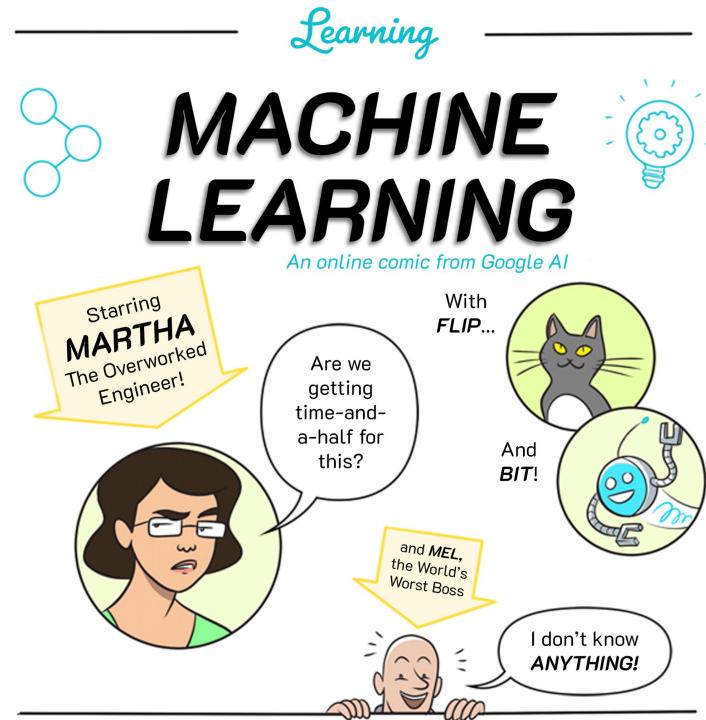
# Machine Learning

Pattern recognition and inference

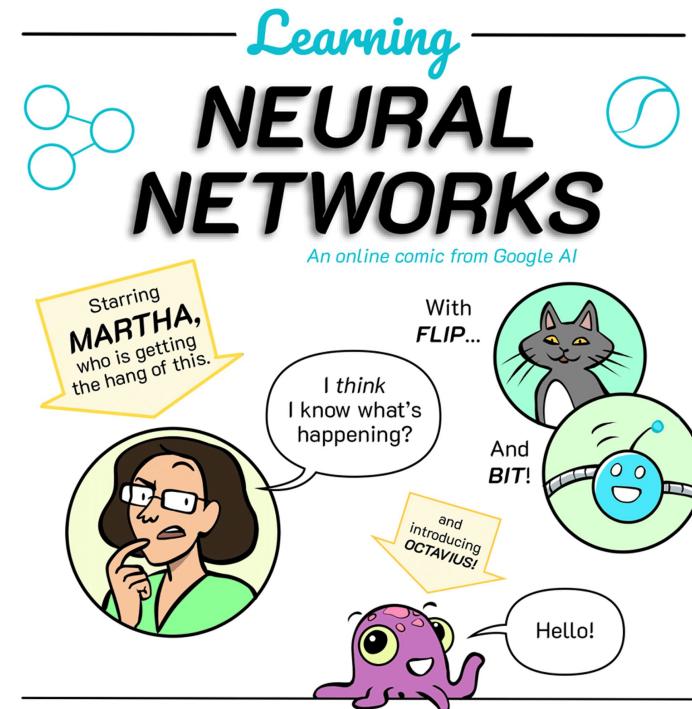
Often guided by examples



# ML preliminaries



<https://cloud.google.com/products/ai/ml-comic-1/>



<https://cloud.google.com/products/ai/ml-comic-2>

# Why is it called DL?

Input layer

One or more hidden layers

Output layer

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG  
PILE OF LINEAR ALGEBRA, THEN COLLECT  
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL  
THEY START LOOKING RIGHT.



<https://xkcd.com/1838/>

# Why DL now?



More hardware

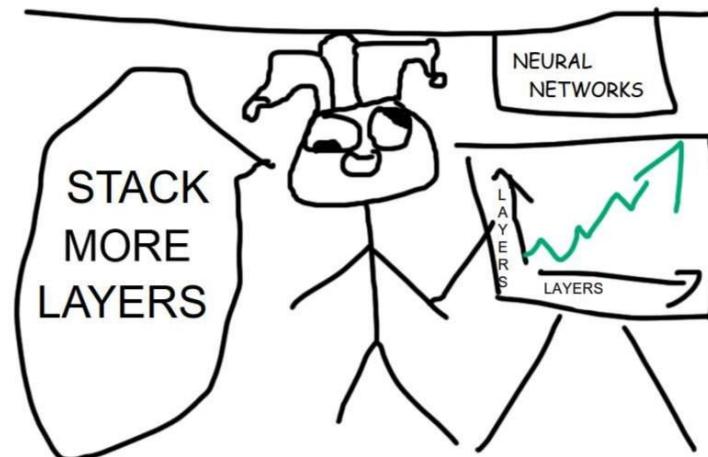
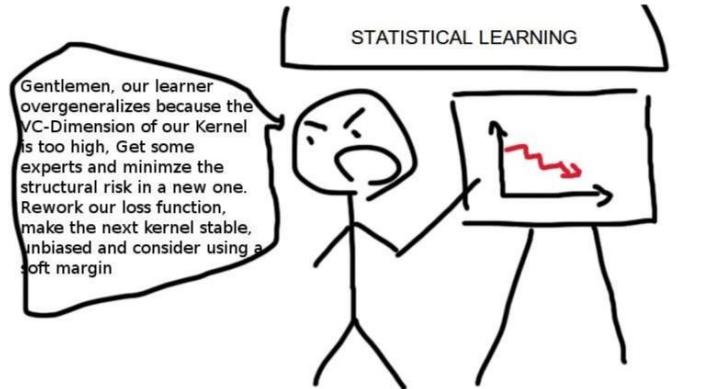


More software

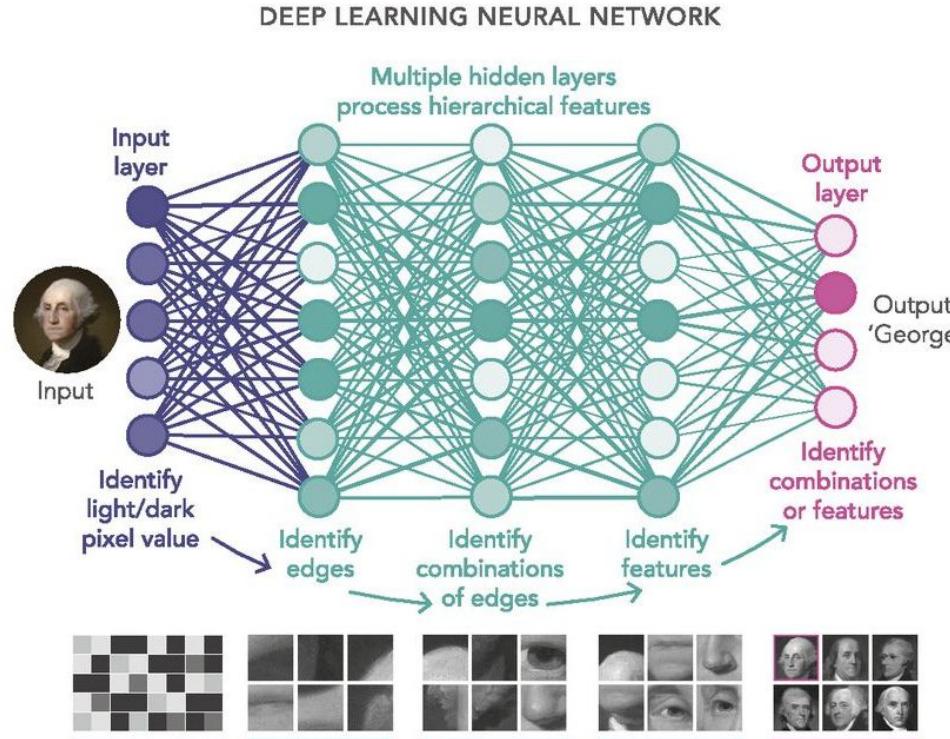


More data

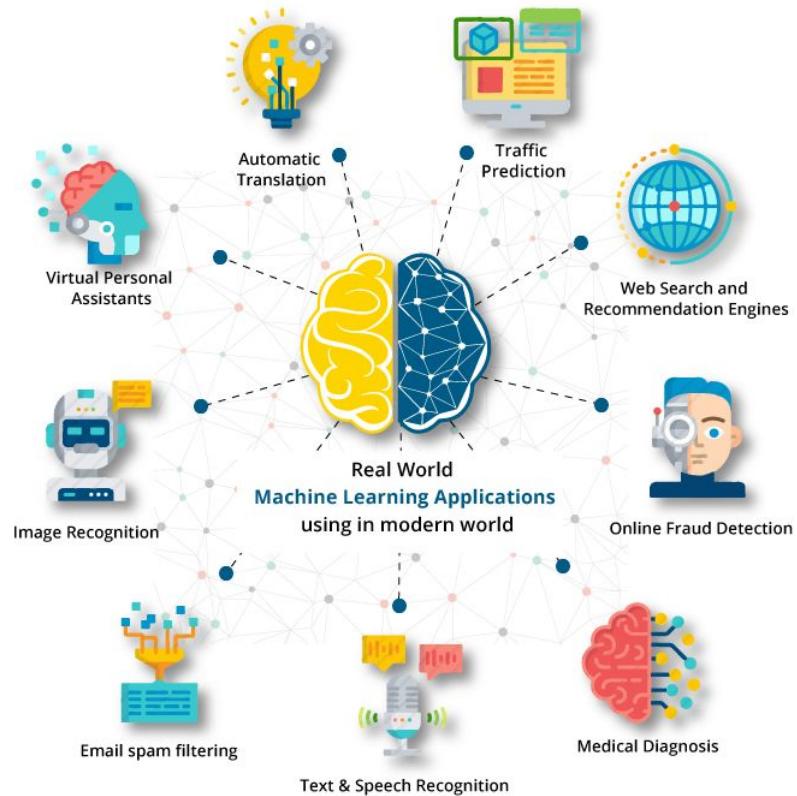
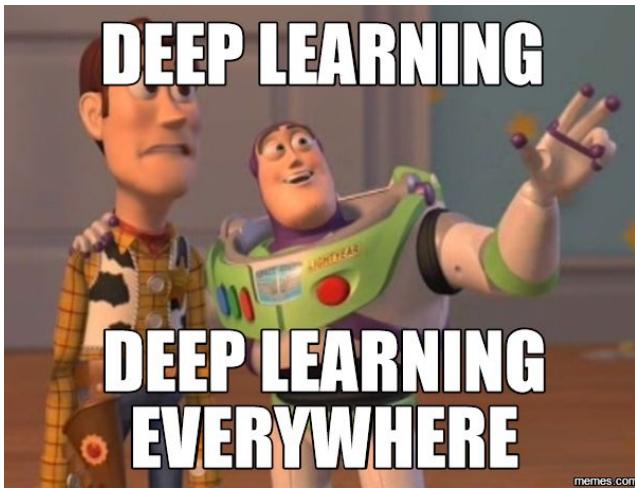
# Why DL now?



# Example



# Applications



# Learning paradigms

## Supervised Learning

**Examples:**  
Classification  
Regression

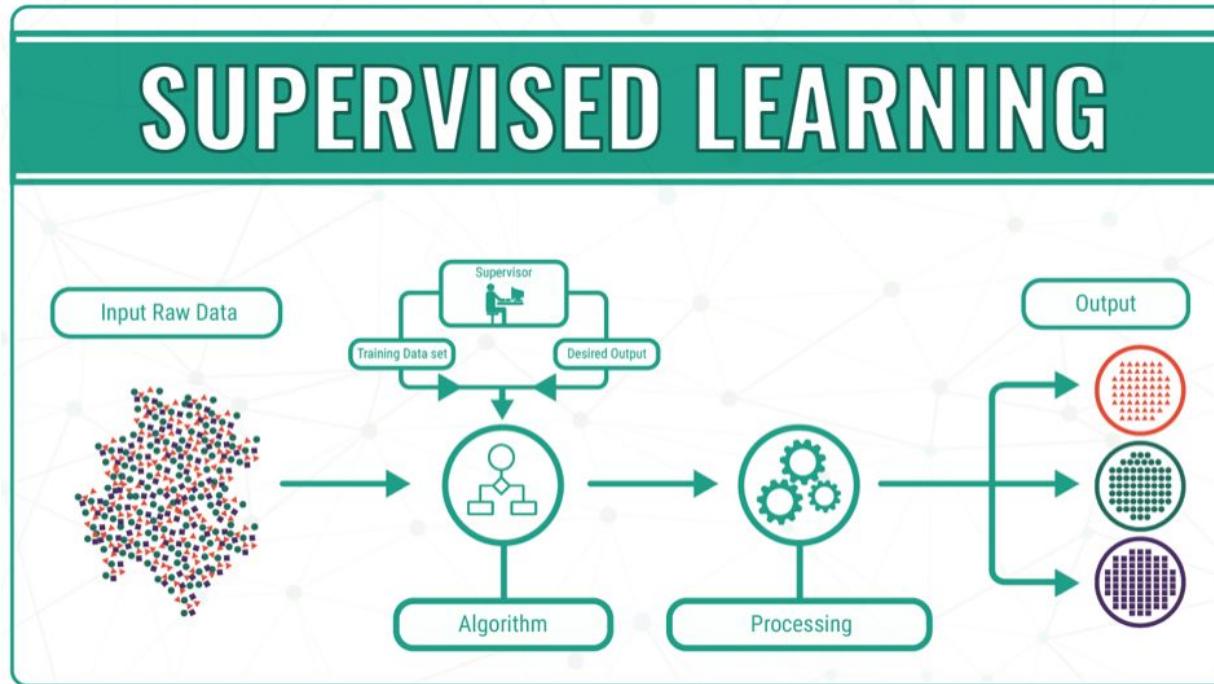
## Unsupervised Learning

**Examples:**  
Autoencoders  
Clustering

## Reinforcement Learning

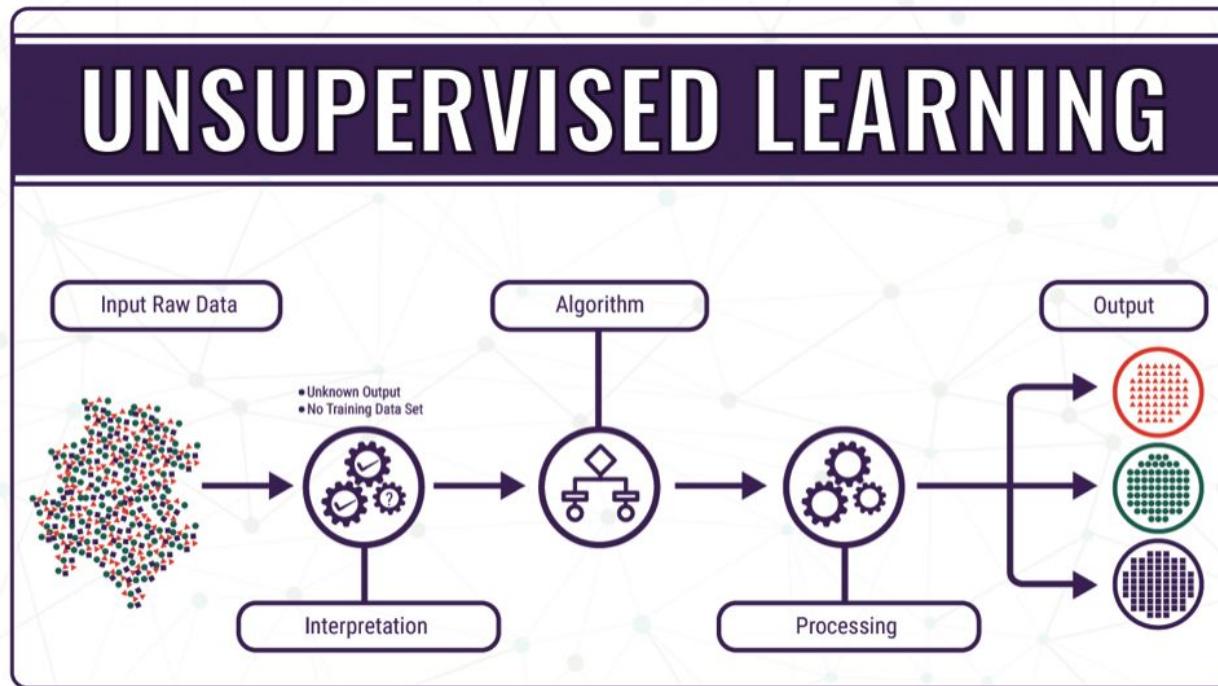
**Examples:**  
Navigation  
Planning

# Learning paradigms



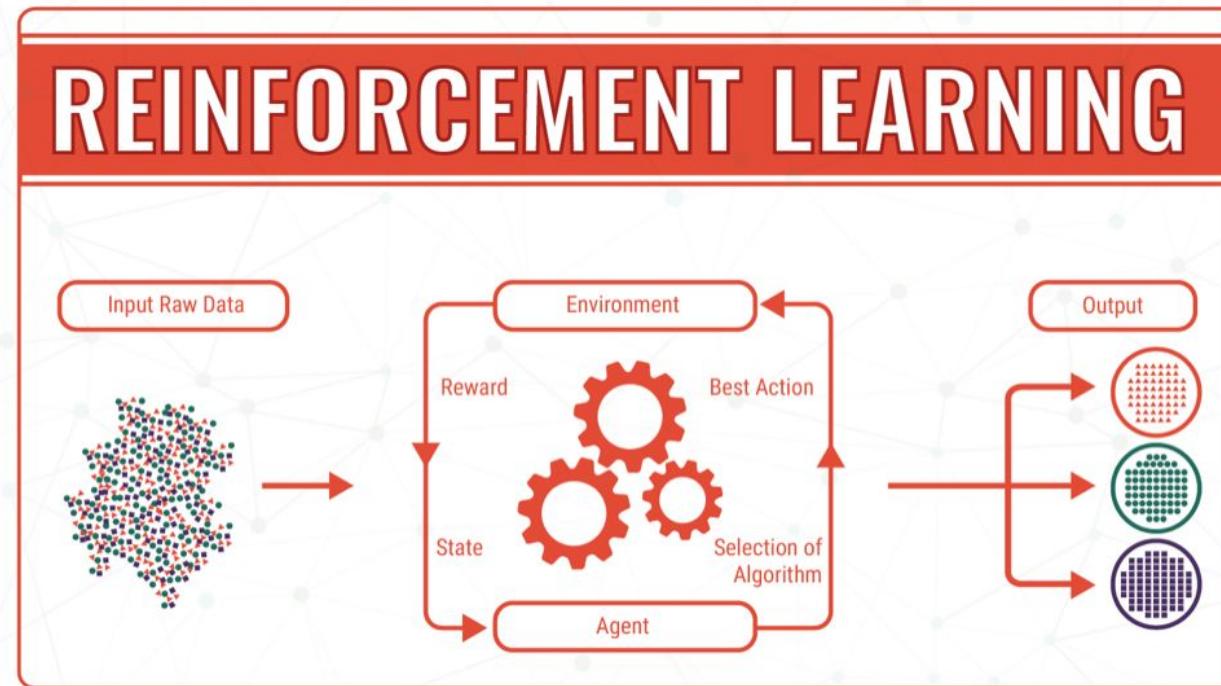
<https://www.linkedin.com/pulse/machine-learning-explained-understanding-supervised-ronald-van-loon>

# Learning paradigms



<https://www.linkedin.com/pulse/machine-learning-explained-understanding-supervised-ronald-van-loon>

# Learning paradigms



# Modeling paradigms

## Discriminative modeling

Learn class boundaries

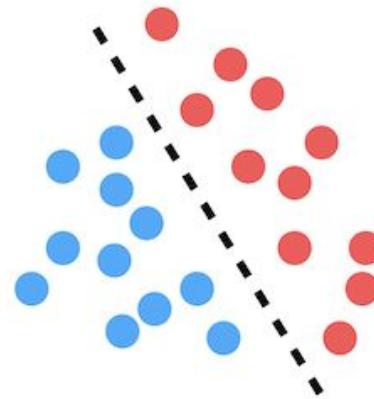
Supervised learning

## Generative modeling

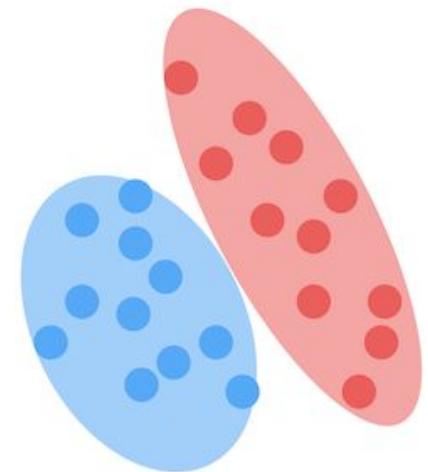
Learn class distributions

Unsupervised learning

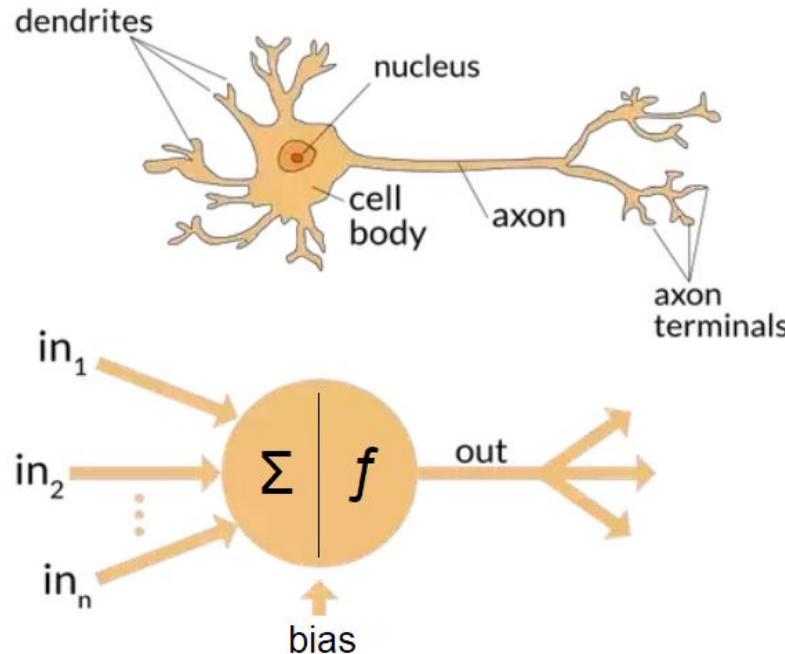
## Discriminative



## Generative

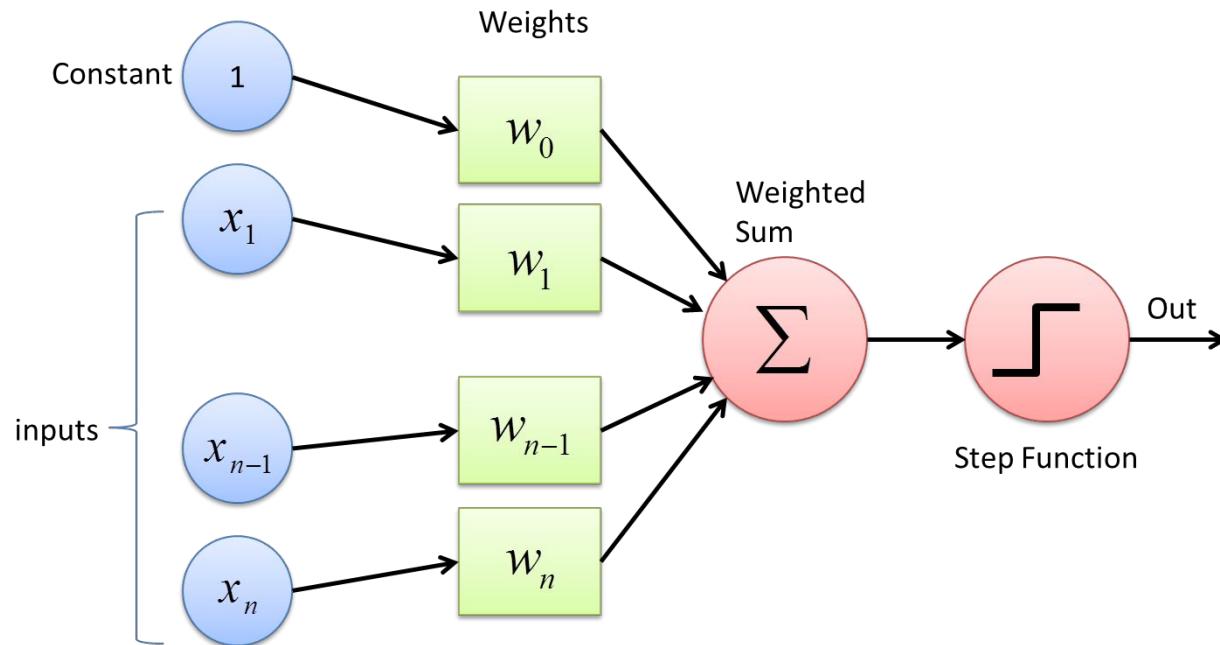


# Perceptron algorithm



<https://towardsdatascience.com/a8b46db828b7>

# Perceptron algorithm



# Perceptron algorithm

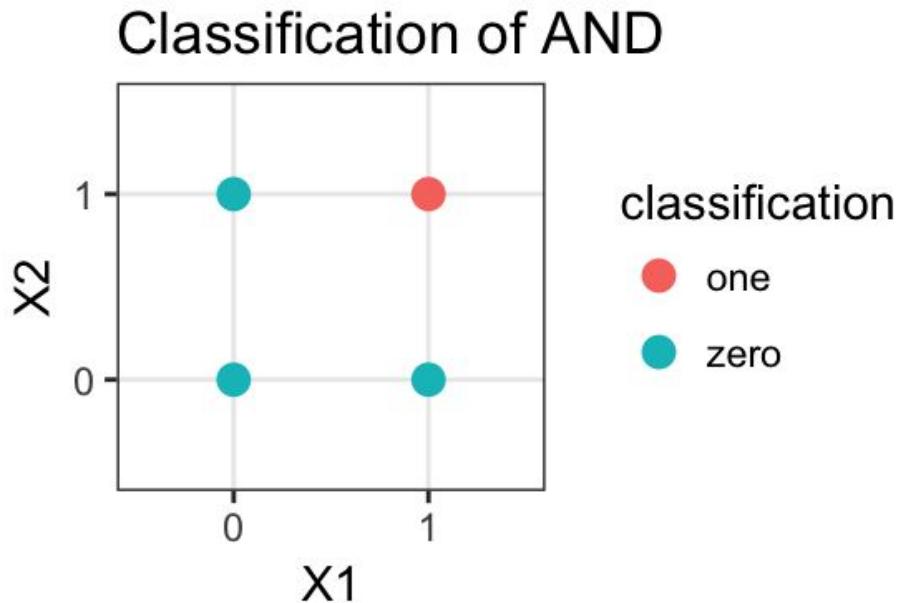
$$\hat{y} = f(z) = \begin{cases} 1 & w_0 + \sum_{i=1}^n x_i w_i > \theta \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{y} = f(z) = \begin{cases} 1 & \sum_{i=0}^n x_i w_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

# Perceptron algorithm

1. Initialize weights  $\mathbf{w}$
2. For  $N$  iterations (or until convergence) do:
  - 2.1. Process each instance  $f(z) = f(\mathbf{w}^T \mathbf{x})$
  - 2.2. Update weights:  $w_i += [y - f(z)] x_i$
3. Return weights

# Perceptron algorithm



$x_1$	$x_2$	$y$
1	1	1
1	0	0
0	1	0
0	0	0

# Perceptron algorithm

Let  $y = f(b, x_1, x_2)$  and  $\mathbf{w} = (0,0,0)$ .

Train  $\mathbf{x} = (1,1,1)$ ;  $y = 1$

$$z = \mathbf{x} \cdot \mathbf{w}^T = 1 \cdot 0 + 1 \cdot 0 + 1 \cdot 0 = 0 \rightarrow f(z) = 0$$

$$\mathbf{w} = (0,0,0) + (1 - 0) \cdot (1,1,1) = (1,1,1)$$

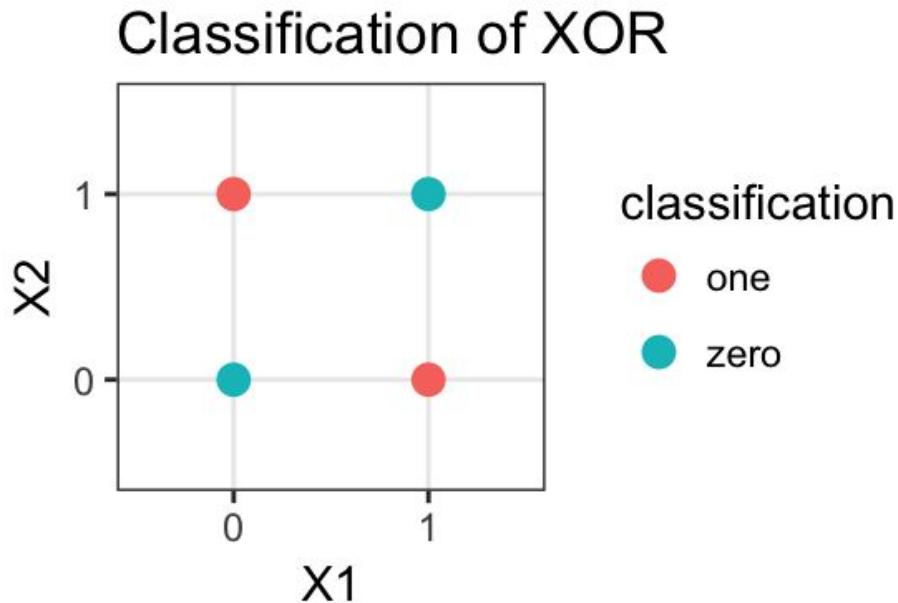
Train  $\mathbf{x} = (1,1,0)$ ;  $y = 0$

$$z = \mathbf{x} \cdot \mathbf{w}^T = 1 \cdot 1 + 1 \cdot 1 + 0 \cdot 1 > 0 \rightarrow f(z) = 1$$

$$\mathbf{w} = (1,1,1) + (0 - 1) \cdot (1,1,0) = (0,0,1)$$

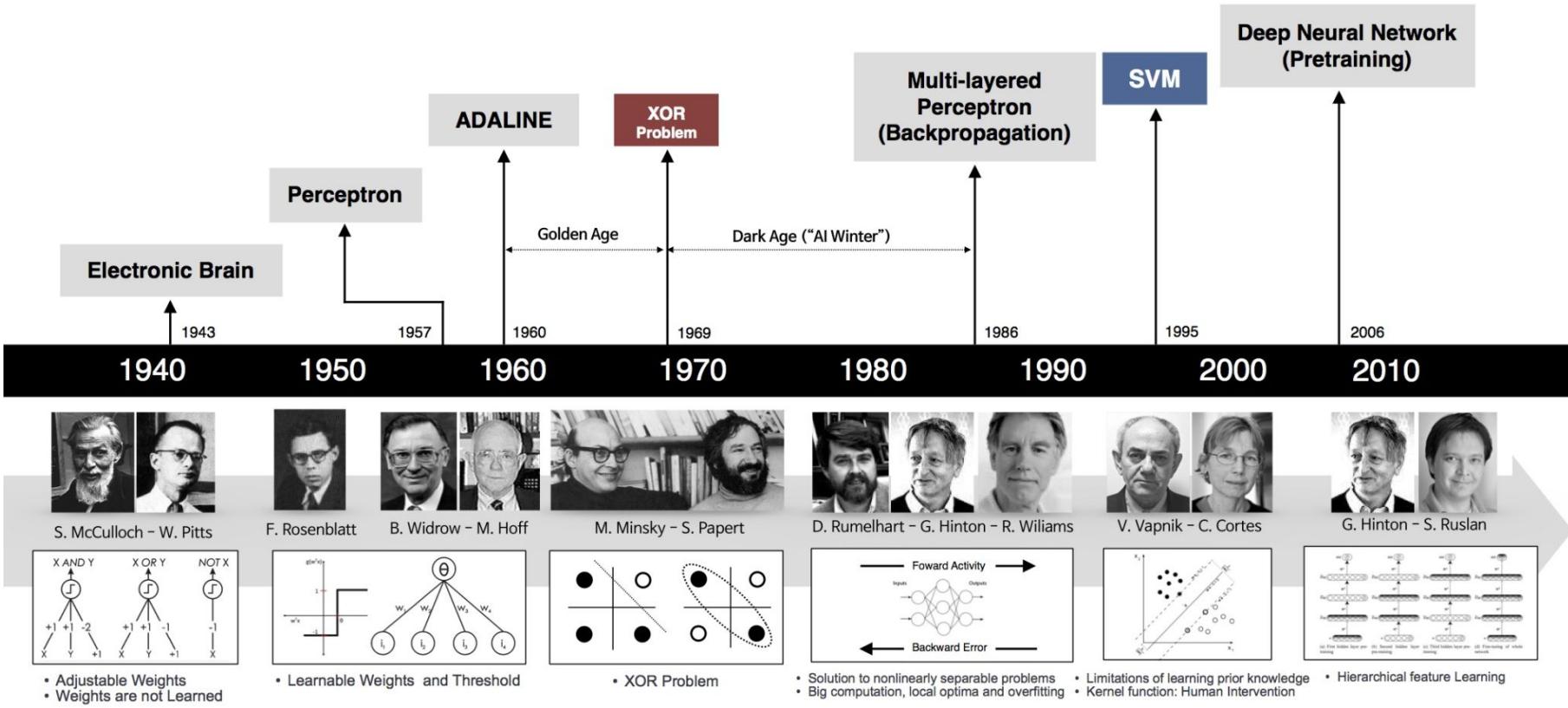
$x_1$	$x_2$	$y$
1	1	1
1	0	0
0	1	0
0	0	0

# Perceptron algorithm

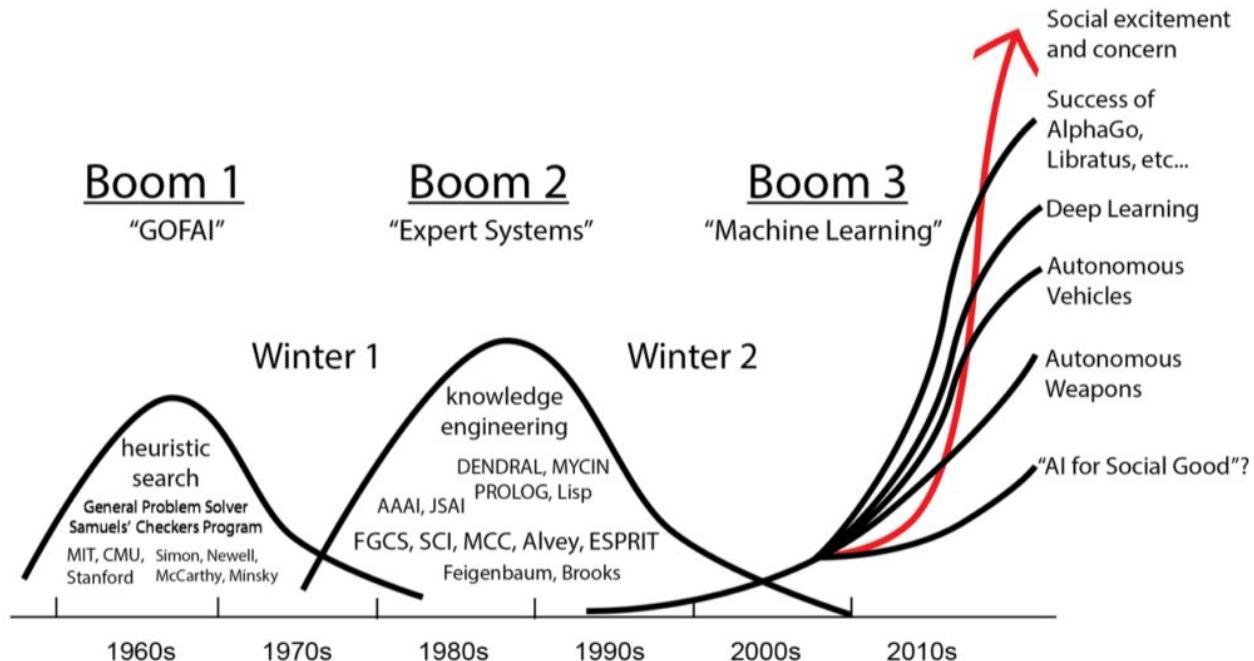


$x_1$	$x_2$	$y$
1	1	0
1	0	1
0	1	1
0	0	0

# Timeline

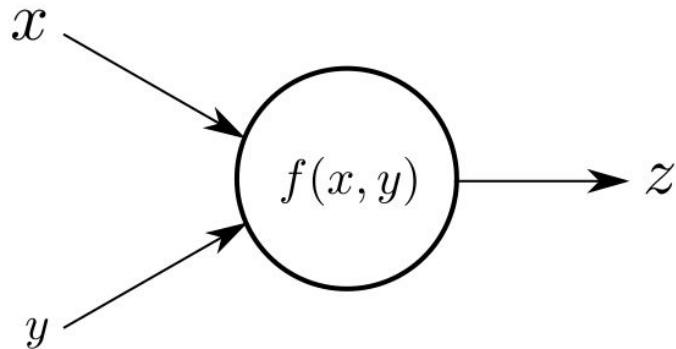


# AI winters

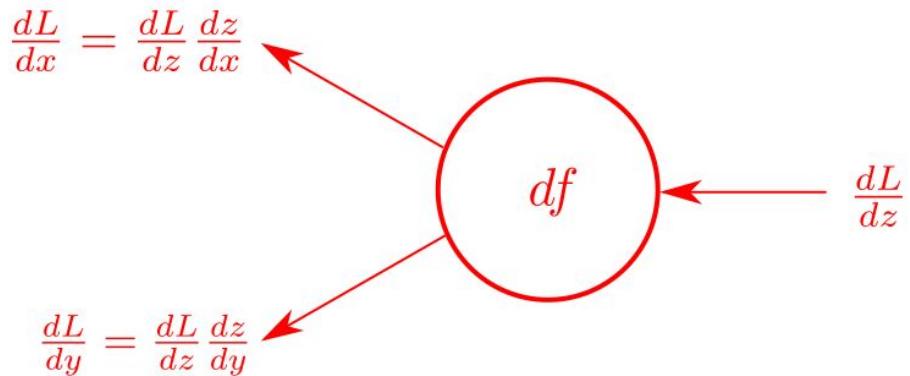


# Backprop

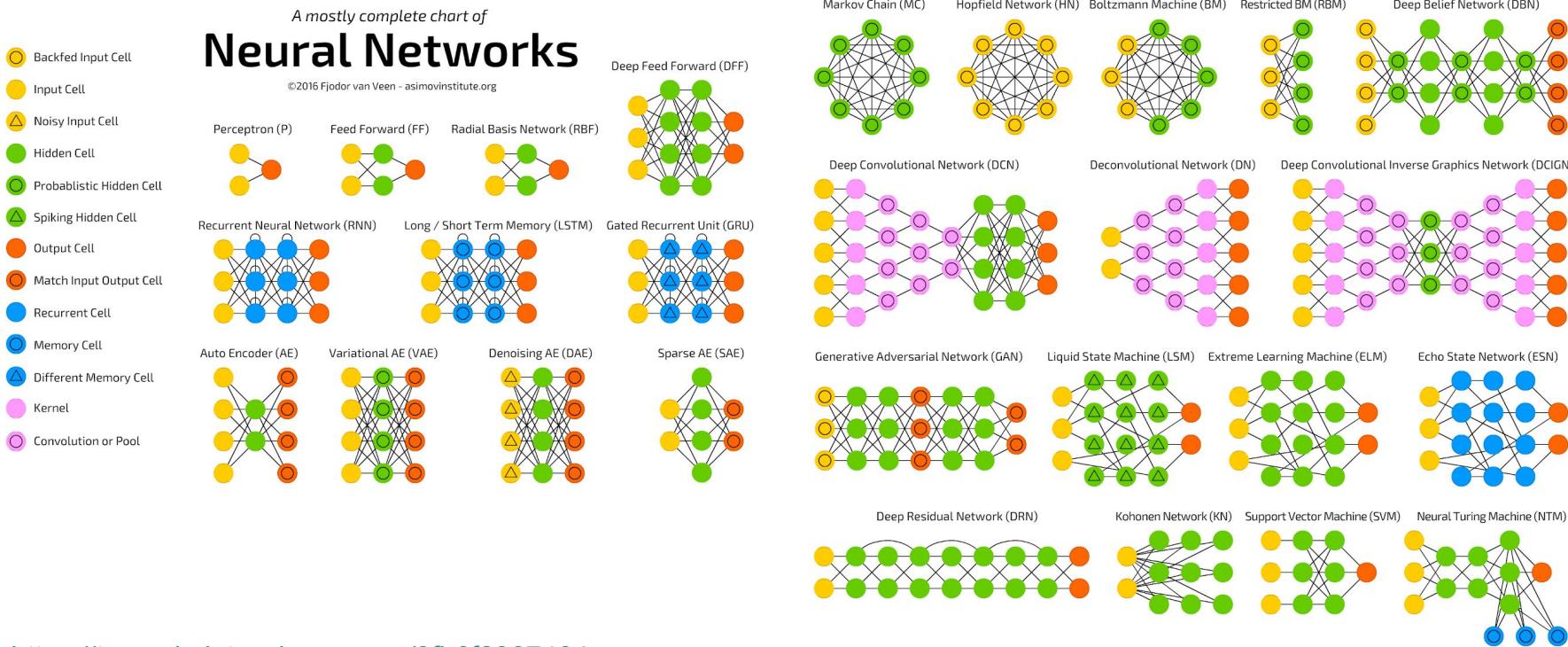
Forwardpass



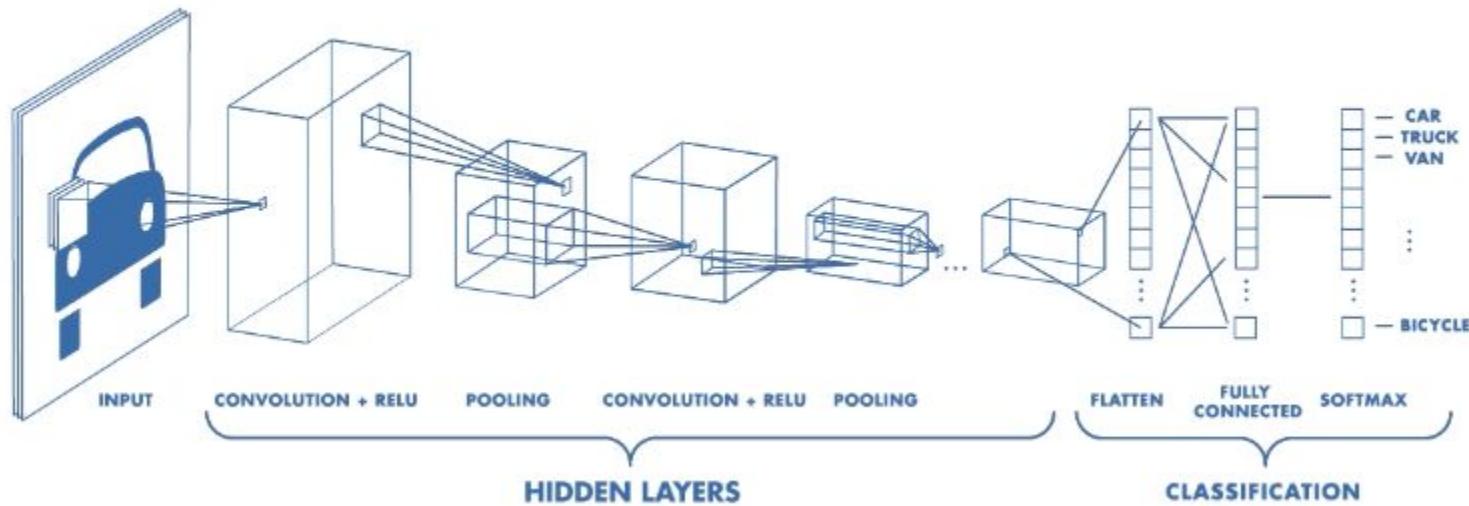
Backwardpass



# Some architectures



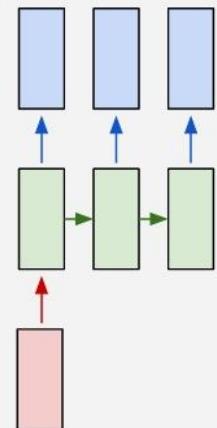
# Convolutional Neural Nets



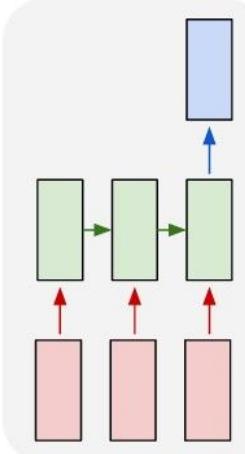
<https://towardsdatascience.com/3bd2b1164a53>

# Recurrent Neural Nets

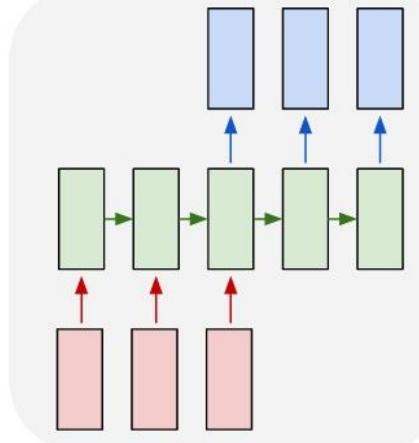
one to many



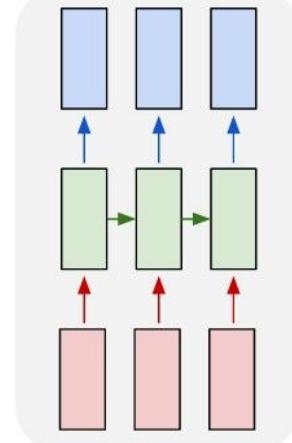
many to one



many to many



many to many



# Building blocks



Layers

Activation functions

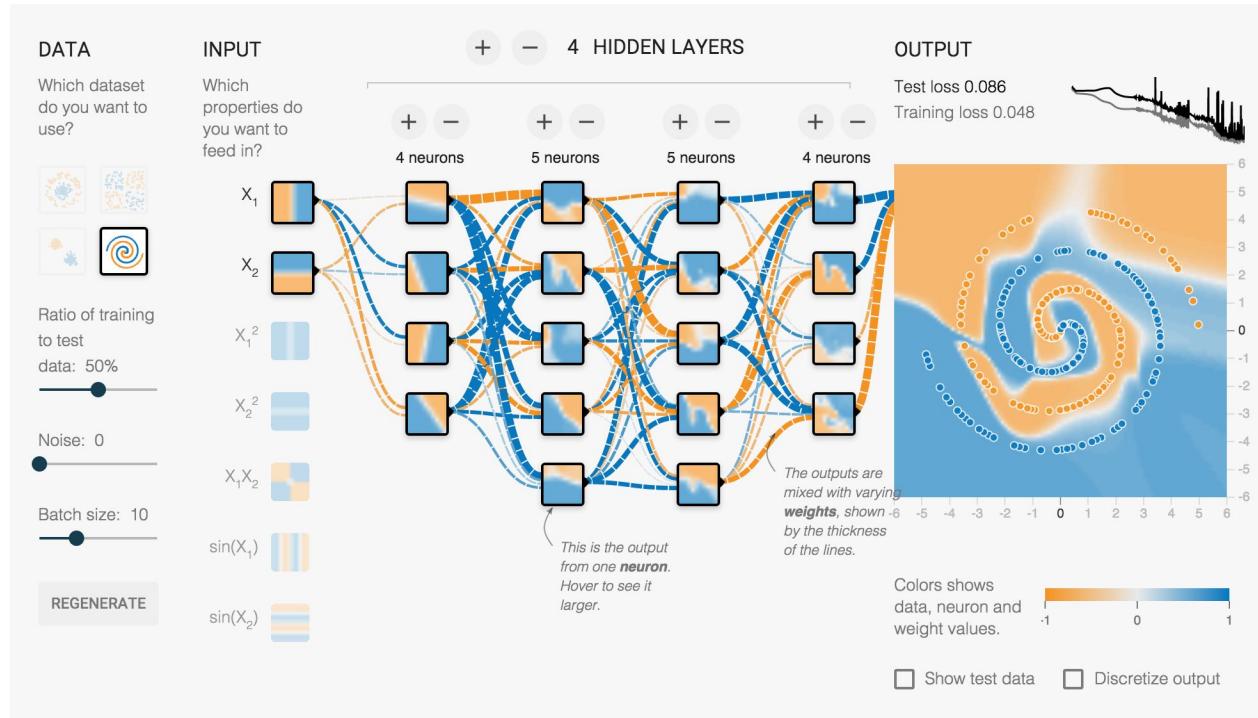
Loss functions

Optimizers

Regularizers

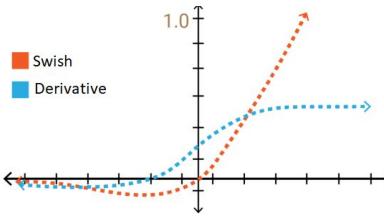
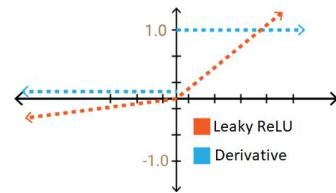
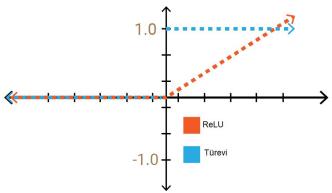
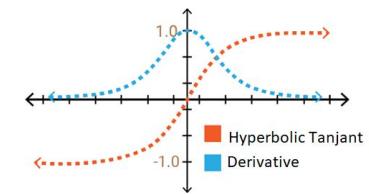
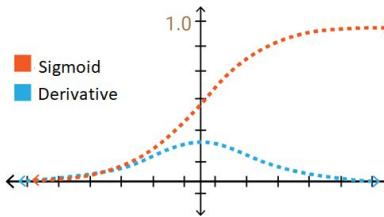
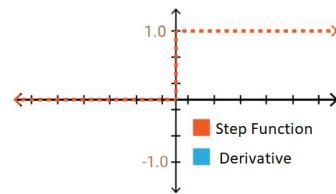
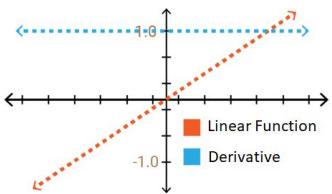
Evaluation metrics

# Layers



<https://blogs.scientificamerican.com/sa-visual/unveiling-the-hidden-layers-of-deep-learning/>

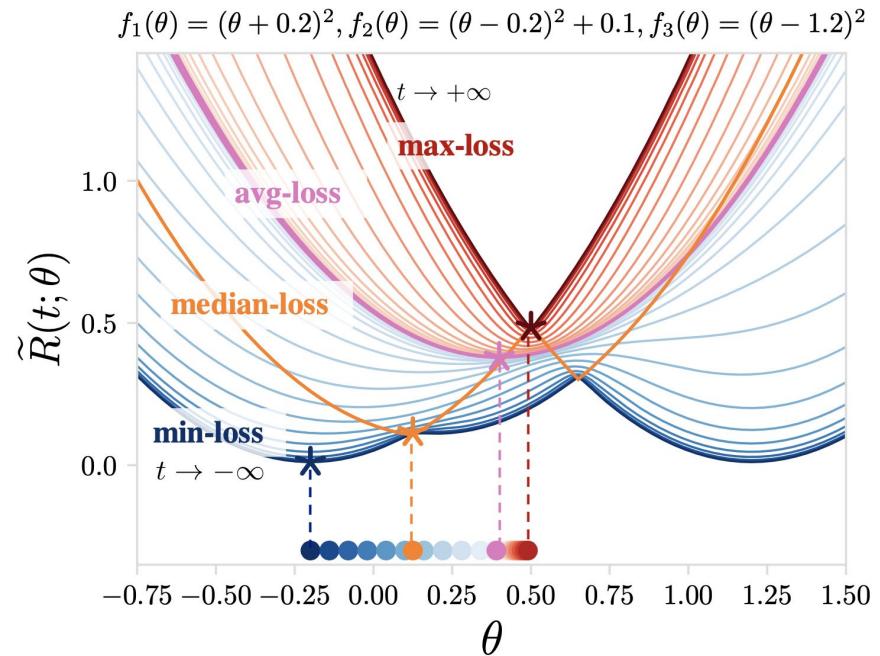
# Activation functions



# Loss functions

**Goal:** Minimize empirical loss

$$J(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\hat{y}_i, y_i)$$



<https://pythonawesome.com/tilted-empirical-risk-minimization-in-python/>

# Loss functions for classification tasks

**Binary Crossentropy loss** for  $C = 2$  classes:

$$\mathcal{L}(\hat{y}_i, y_i) = -\frac{1}{N} \sum_i^N y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)$$

**Categorical Crossentropy loss** for  $C > 2$  classes:

$$\mathcal{L}(\hat{y}_i, y_i) = -\frac{1}{N} \sum_i^N \sum_c^C \mathbb{1}_{y_i \in C_c} \log \hat{y}_i$$

Note:  $\hat{y}_i = p(y_i)$

# Loss functions for regression tasks

**Mean Squared Error loss:**

$$\mathcal{L}(\hat{y}_i, y_i) = (y_i - \hat{y}_i)^2$$

*Notes:*

$$\text{RMSE} = \sqrt{\text{MSE}}$$

$$\text{MAE} \leq \text{RMSE}$$

**Mean Absolute Error loss:**

$$\mathcal{L}(\hat{y}_i, y_i) = |y_i - \hat{y}_i|$$

$$\text{RMSE} \leq \text{MAE} \cdot \sqrt{n}$$

# Loss functions for regression tasks

CASE 1: Evenly distributed errors

ID	Error	Error	Error^2
1	2	2	4
2	2	2	4
3	2	2	4
4	2	2	4
5	2	2	4
6	2	2	4
7	2	2	4
8	2	2	4
9	2	2	4
10	2	2	4

CASE 2: Small variance in errors

ID	Error	Error	Error^2
1	1	1	1
2	1	1	1
3	1	1	1
4	1	1	1
5	1	1	1
6	3	3	9
7	3	3	9
8	3	3	9
9	3	3	9
10	3	3	9

CASE 3: Large error outlier

ID	Error	Error	Error^2
1	0	0	0
2	0	0	0
3	0	0	0
4	0	0	0
5	0	0	0
6	0	0	0
7	0	0	0
8	0	0	0
9	0	0	0
10	20	20	400

MAE	RMSE
2.000	2.000

MAE	RMSE
2.000	2.236

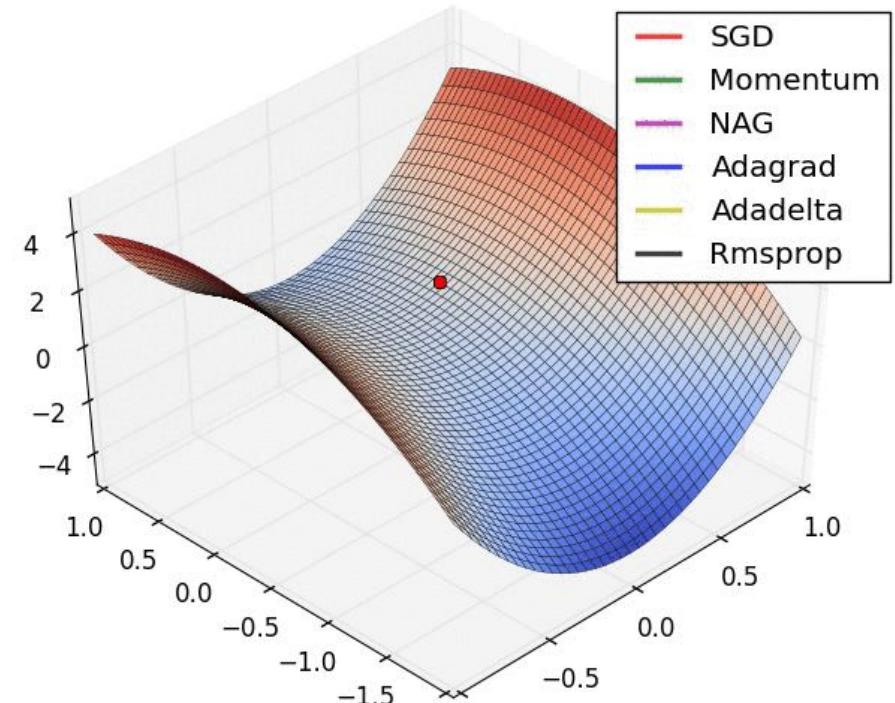
MAE	RMSE
2.000	6.325

# Optimizers

Mostly based on Gradient Descent:

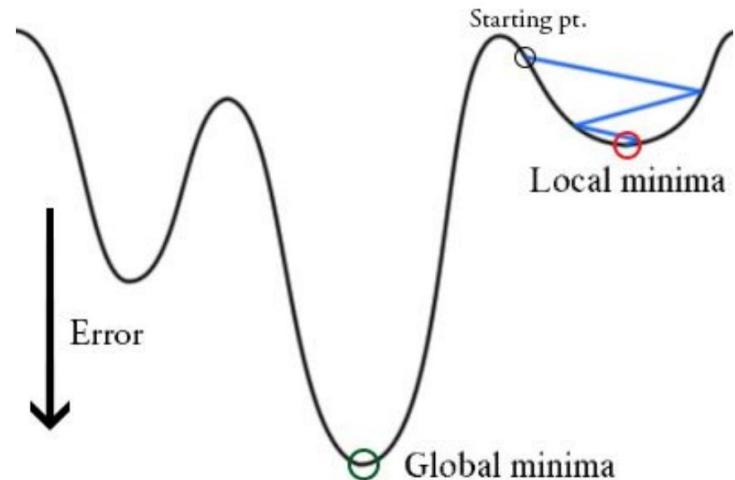
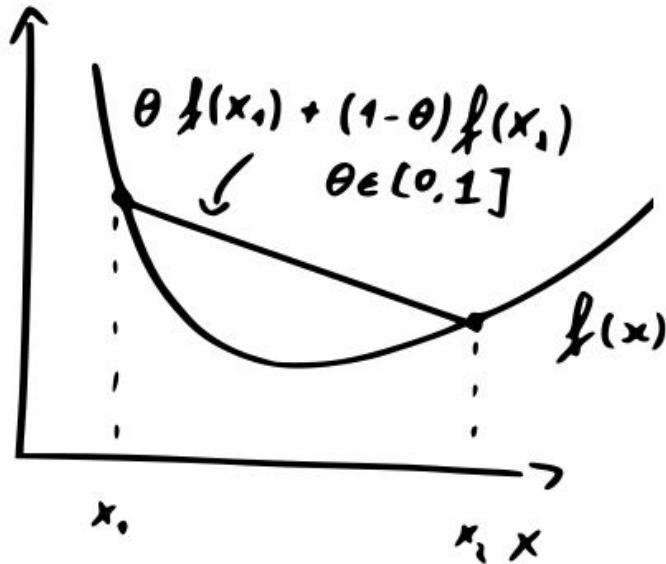
$$\theta = \theta - \eta \cdot \nabla J(\mathbf{W}; \theta)$$

- SGD
- RMSProp
- Adam
- AdaGrad
- AdaDelta



<https://ruder.io/optimizing-gradient-descent/>

# Non-convex optimization



<https://michielsstock.github.io/posts/2018/2018-03-07-ConvexSummary/>

[https://www.cs.ubc.ca/labs/lci/mlrg/slides/non\\_convex\\_optimization.pdf](https://www.cs.ubc.ca/labs/lci/mlrg/slides/non_convex_optimization.pdf)

# Regularizers

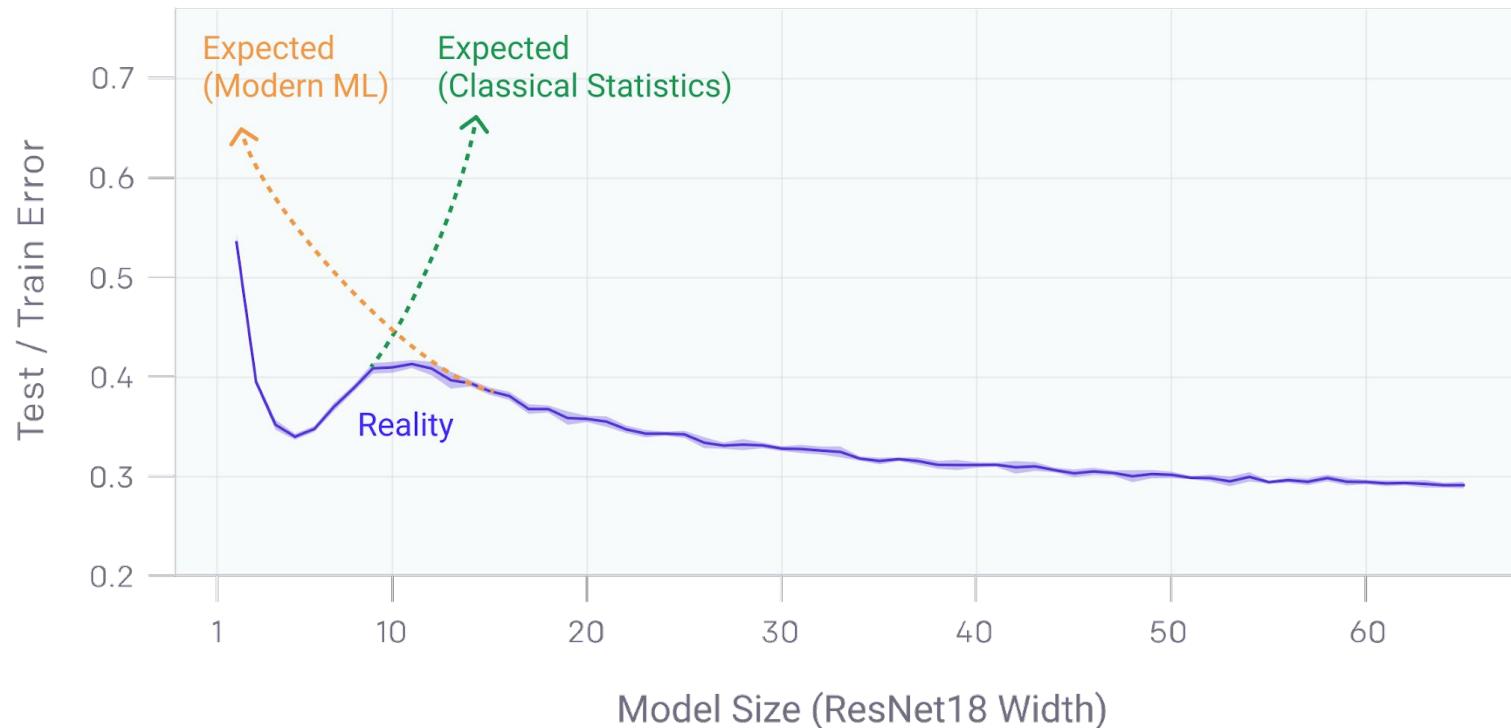


**Goal:** Improve model generalization

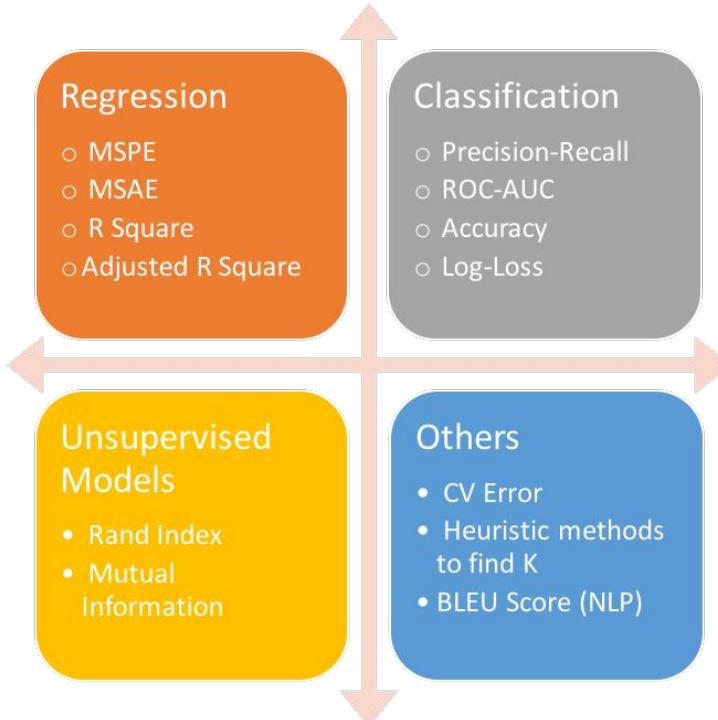
- Dropout\*
- BatchNorm
- Early stopping\*
- L1 & L2 penalty
- Adaptive learning rate\*
- Data augmentation

\* Very very useful, really

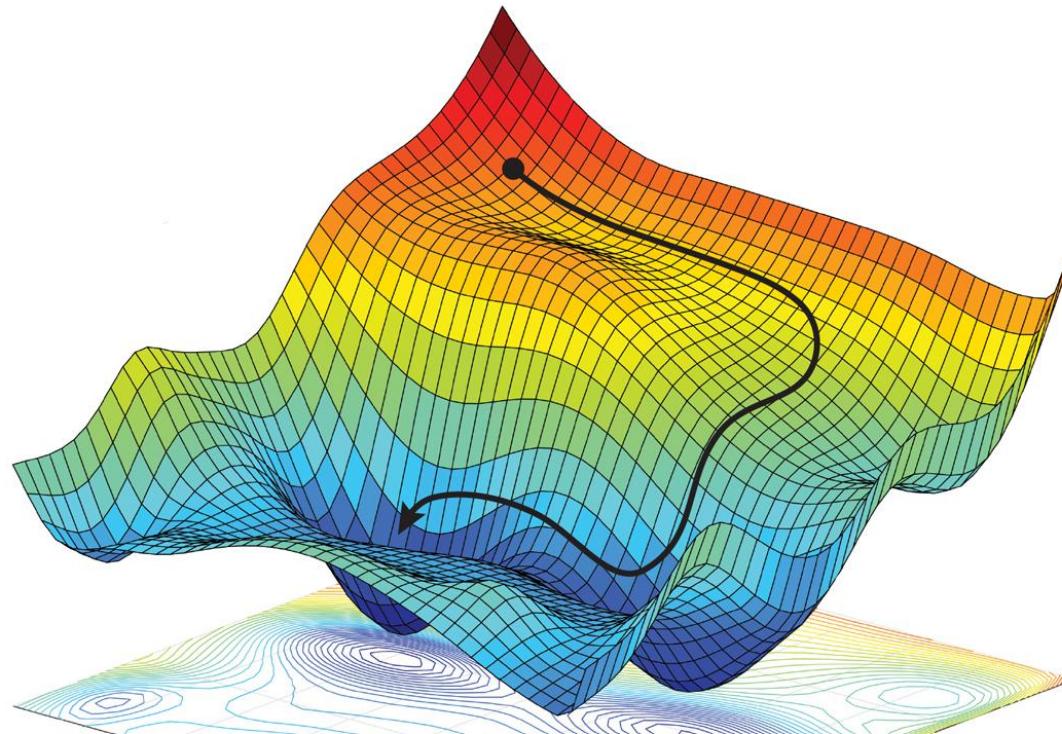
# Note: Double Descent



# Evaluation metrics



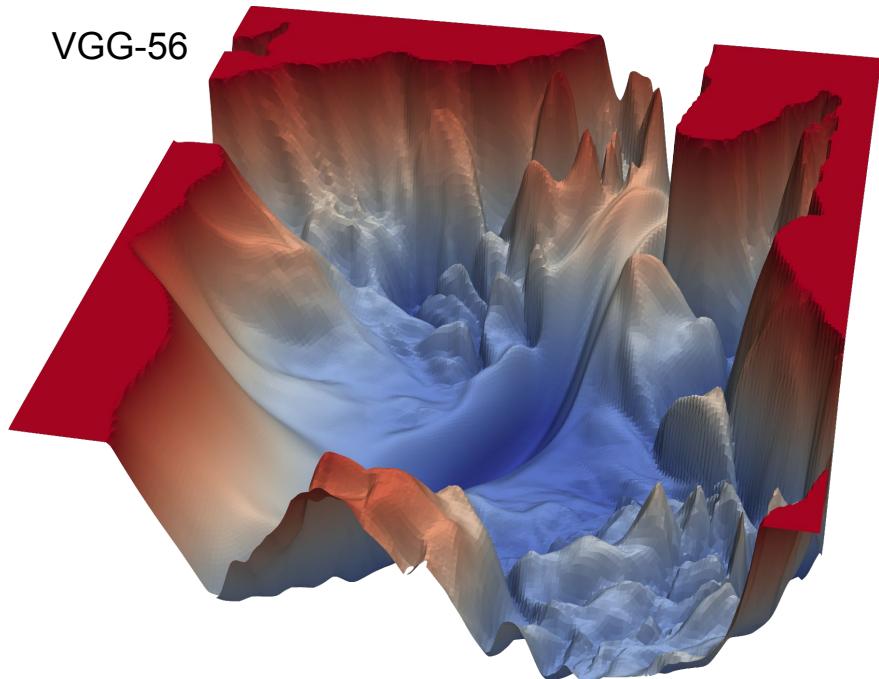
# DL in practice



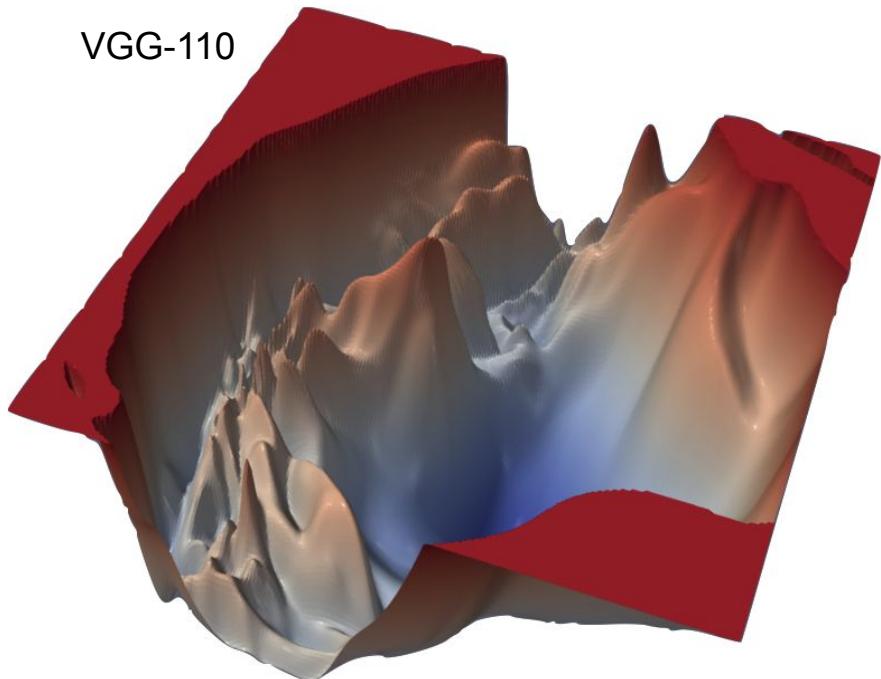
<https://www.science.org/content/article/ai-researchers-allege-machine-learning-alchemy>

# Loss landscapes

VGG-56



VGG-110



# Generalization

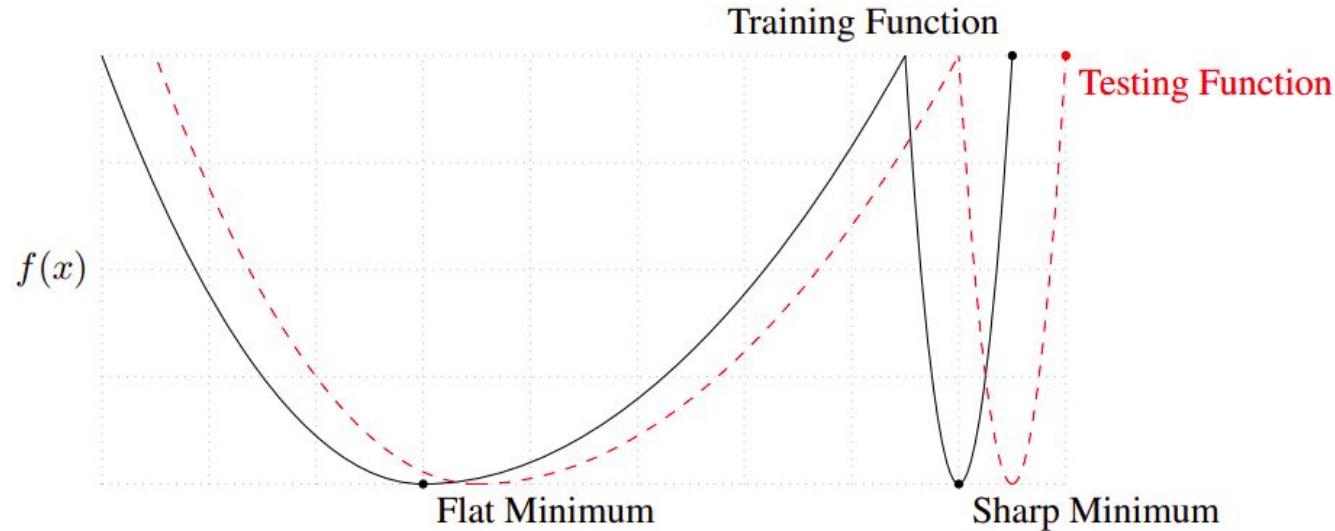
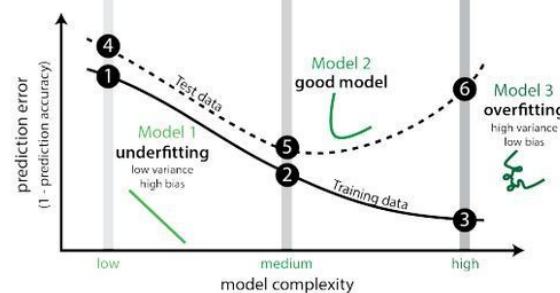
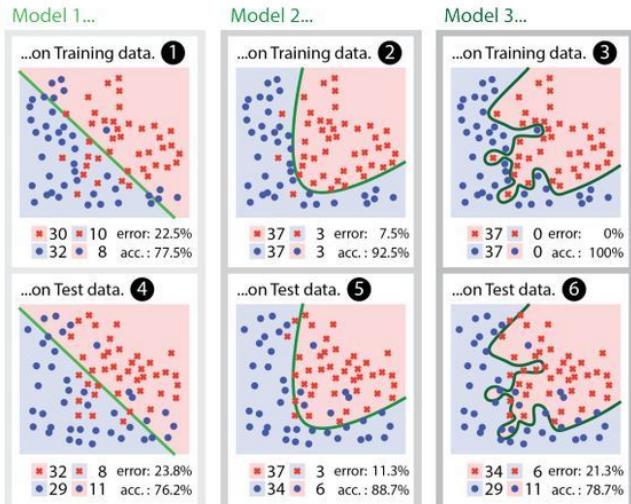
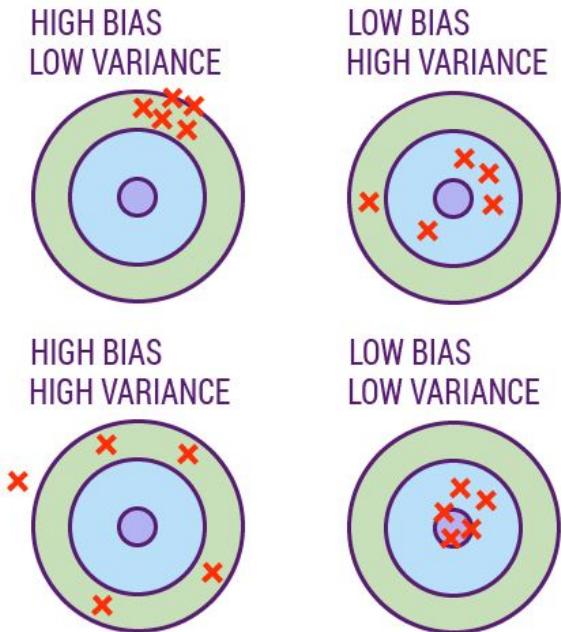
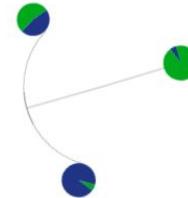
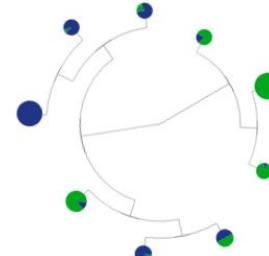
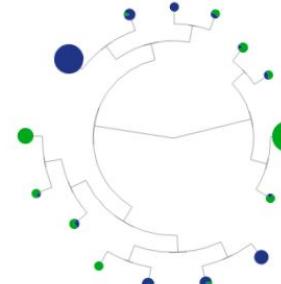
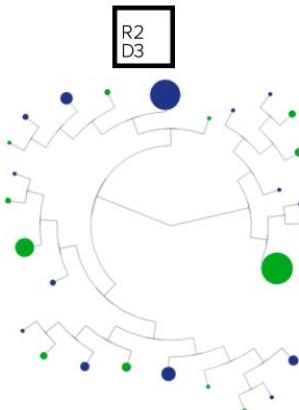


Figure 1: A Conceptual Sketch of Flat and Sharp Minima. The Y-axis indicates value of the loss function and the X-axis the variables (parameters)

# Bias-Variance tradeoff



# Bias-Variance tradeoff interactive demo

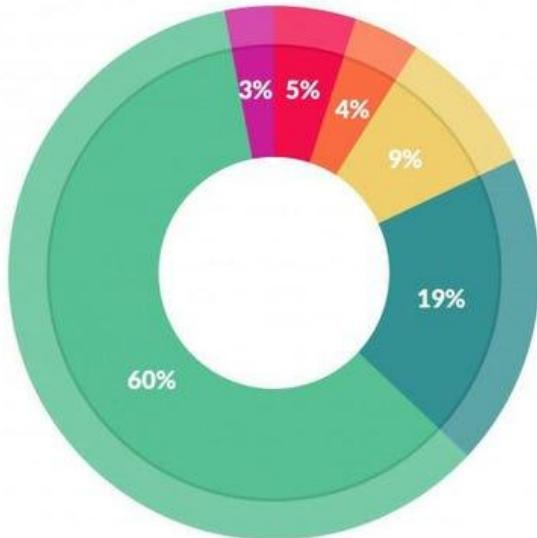


A VISUAL INTRODUCTION  
TO MACHINE LEARNING—PART II

## Model Tuning and the Bias-Variance Tradeoff

<http://www.r2d3.us/visual-intro-to-machine-learning-part-2/>

# Data science

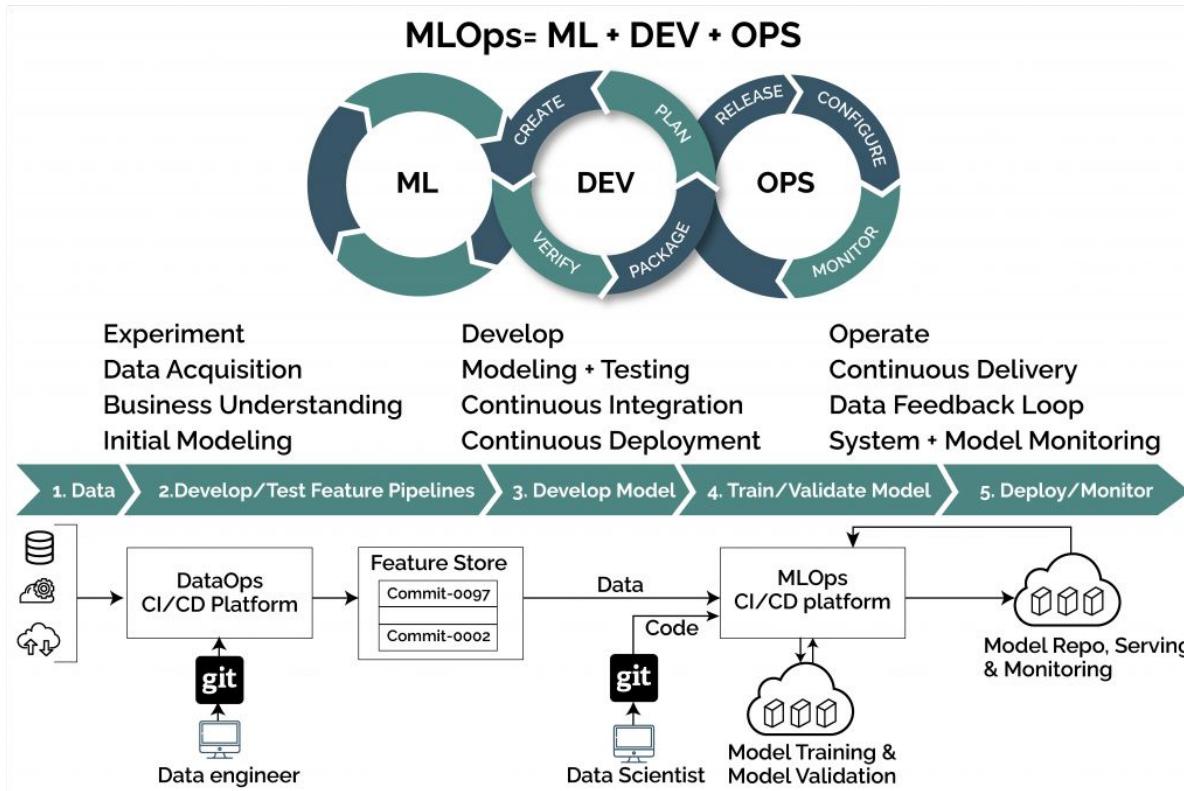


## What data scientists spend the most time doing

- *Building training sets: 3%*
- *Cleaning and organizing data: 60%*
- *Collecting data sets; 19%*
- *Mining data for patterns: 9%*
- *Refining algorithms: 4%*
- *Other: 5%*

<https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/>

# DL in practice



# Limitations of DL

Lack of reasoning

Lack of generalization → DL excels at *automating* tasks

Lack of prior knowledge → solvable with transfer learning

Lack of interpretability (black-box models) → XAI working on this

Lots of training data often needed → Recent research on *Efficient DL*

Models prone to adversarial noise → Hot topic!

DL is a “too-empirical” field

# Recommended readings

1. Zachary C. Lipton, Jacob Steinhardt. [Troubling Trends in Machine Learning Scholarship](#). arXiv:1807.03341, 2018.
2. D. Sculley, G. Holt, D. Golovin, E. Davydov, T. Phillips, D. Ebner, V. Chaudhary, M. Young, J.-F. Crespo, D. Dennison. [Hidden Technical Debt in Machine Learning Systems](#). Proc. NeurIPS, 2015.

# Some tricks and tips



<http://cs231n.github.io/neural-networks-3/>

<http://karpathy.github.io/2019/04/25/recipe/>

<https://towardsdatascience.com/3ca24c31ddc8>

<https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-deep-learning-tips-and-tricks>

# Bibliography

1. F. Chollet. **Deep Learning with Python**, 2021
2. M. Nielsen. **Neural Networks and Deep Learning**, 2015
3. I. Goodfellow, Y. Bengio, A. Courville. **Deep Learning**, 2016
4. A. Gibson, J. Patterson. **Deep Learning: A Practitioners Approach**, 2017
5. R. O. Duda, P. E. Hart, D. G. Stork. **Pattern Classification**, 2000

