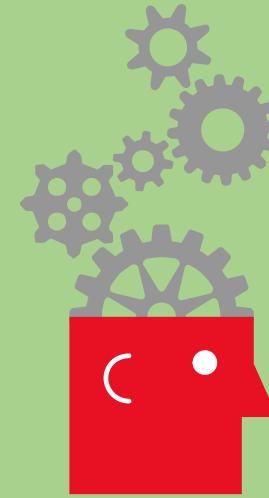


Artificial Intelligence and Machine Learning





1000kg

40mpg



3000kg

20mpg



2000kg



?

inputs

output



1000kg

40mpg



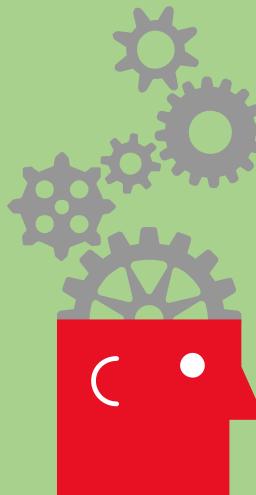
3000kg

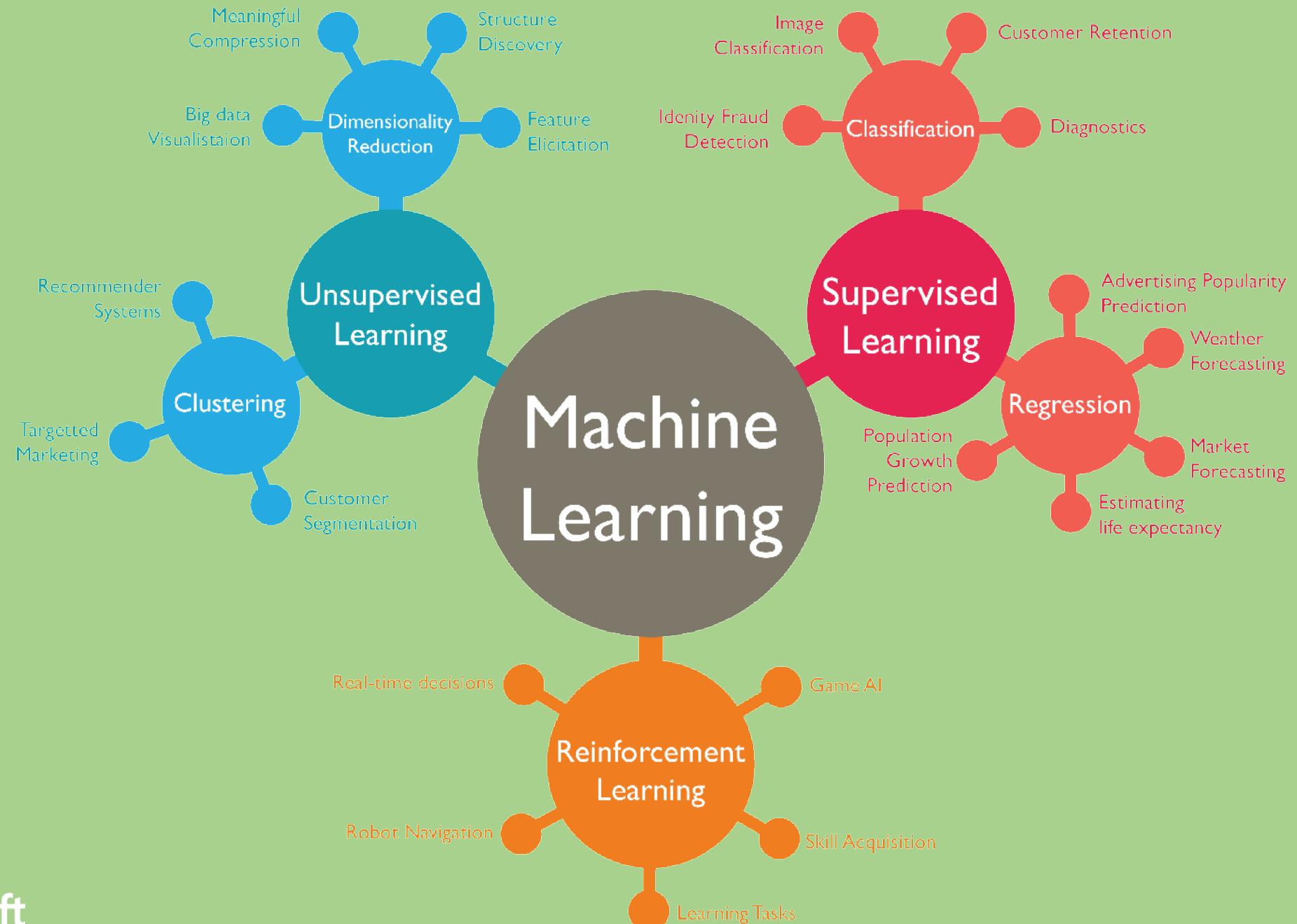
20mpg

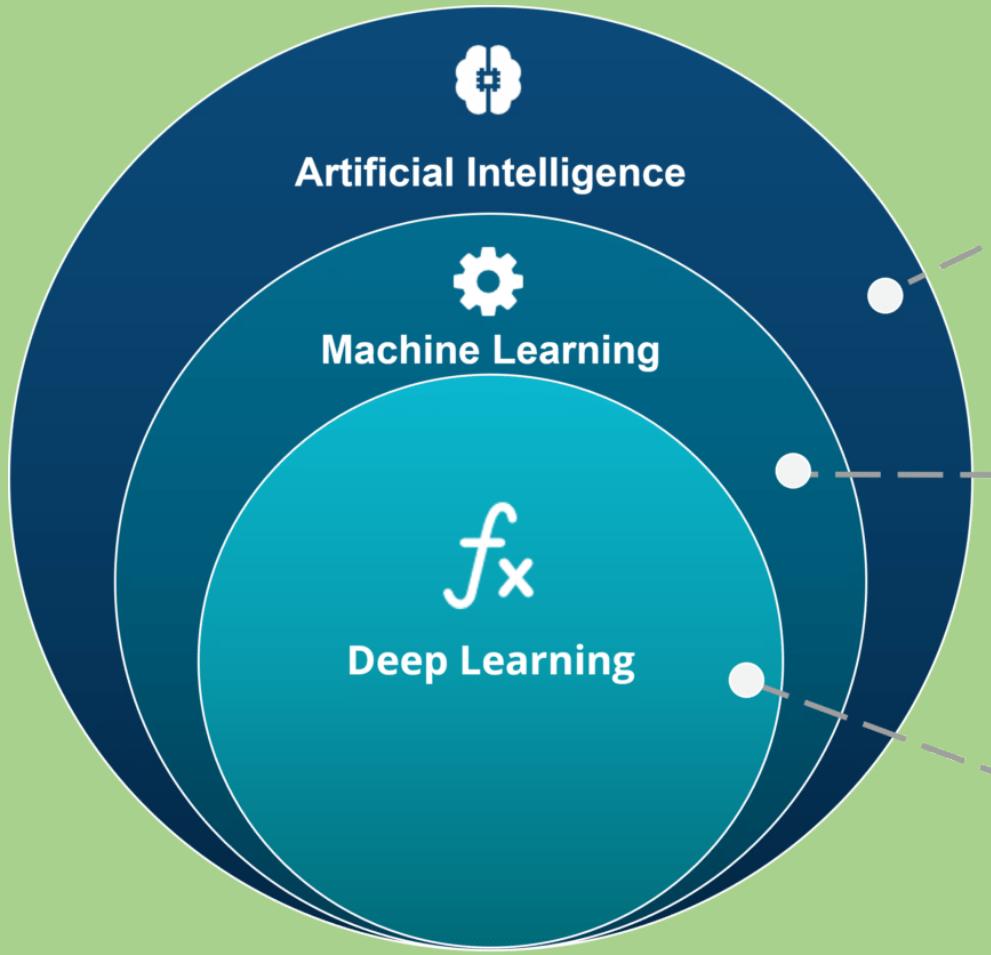


2000kg

30mpg







ARTIFICIAL INTELLIGENCE

A technique which enables machines to mimic human behaviour

MACHINE LEARNING

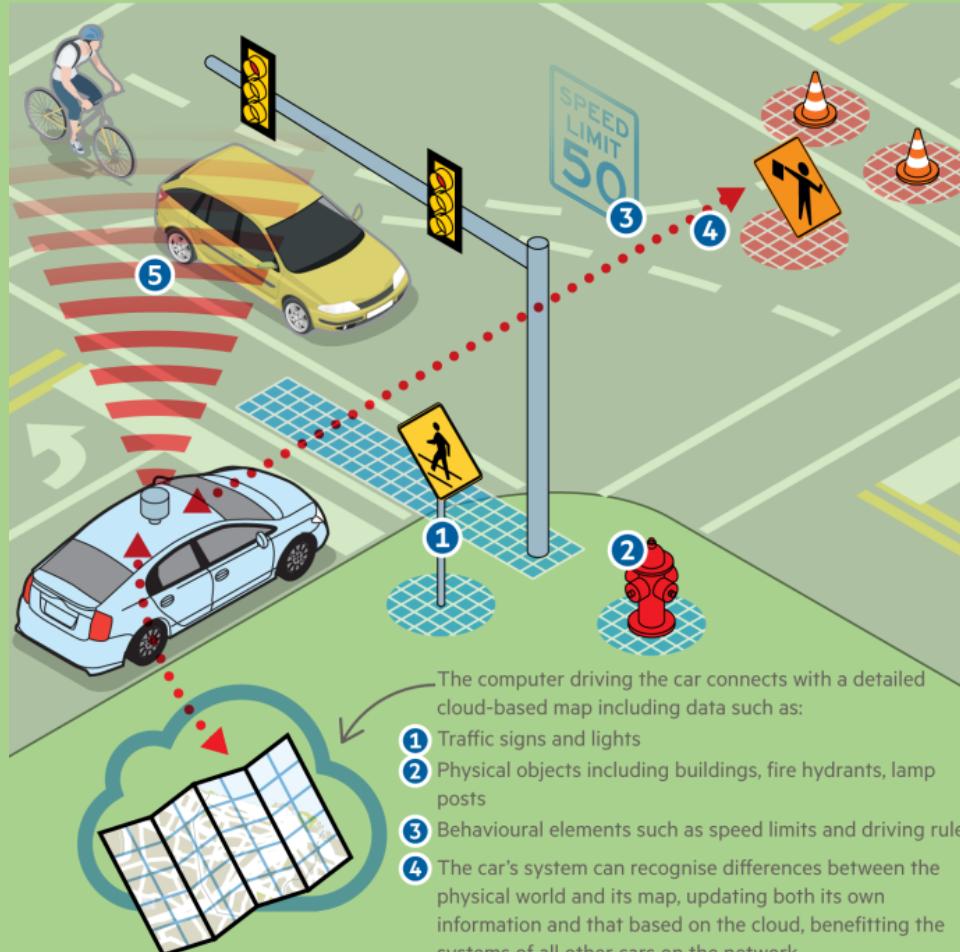
Subset of AI technique which use statistical methods to enable machines to improve with experience

DEEP LEARNING

Subset of ML which make the computation of multi-layer neural network feasible

AI Example

How autonomous cars understand the world around them



Source: FT research Graphic: Ian Bott
© FT



Perception

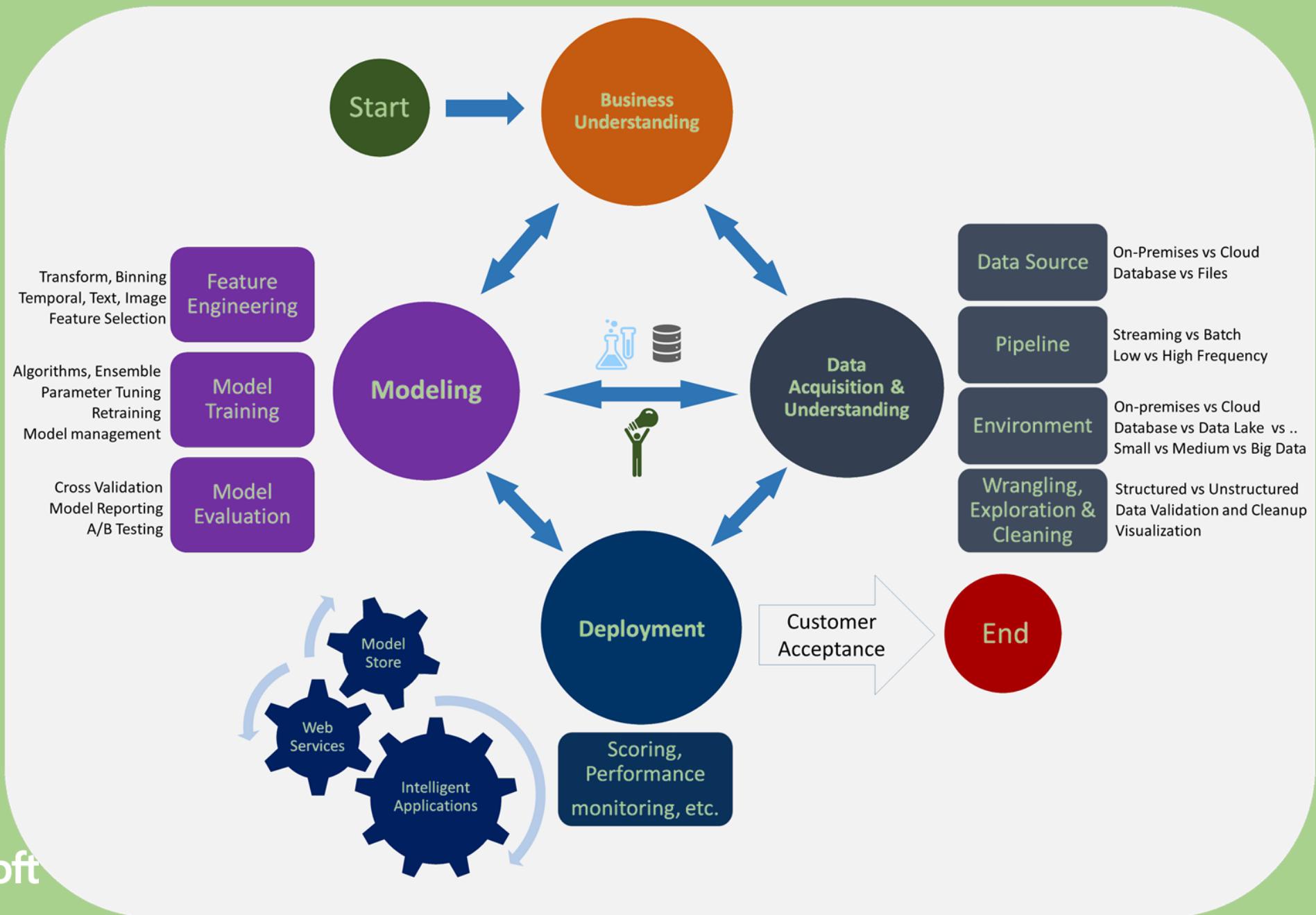


Prediction

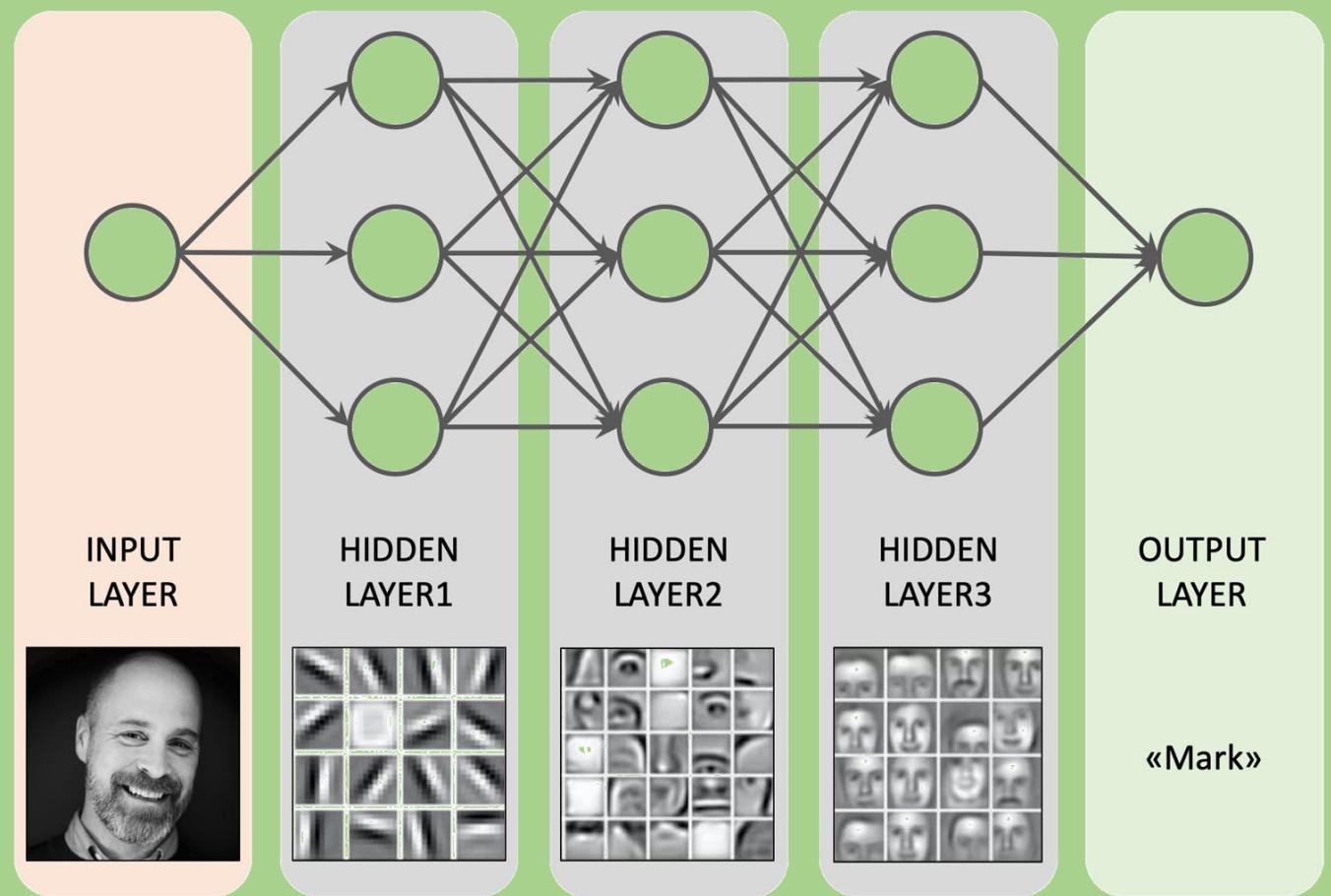
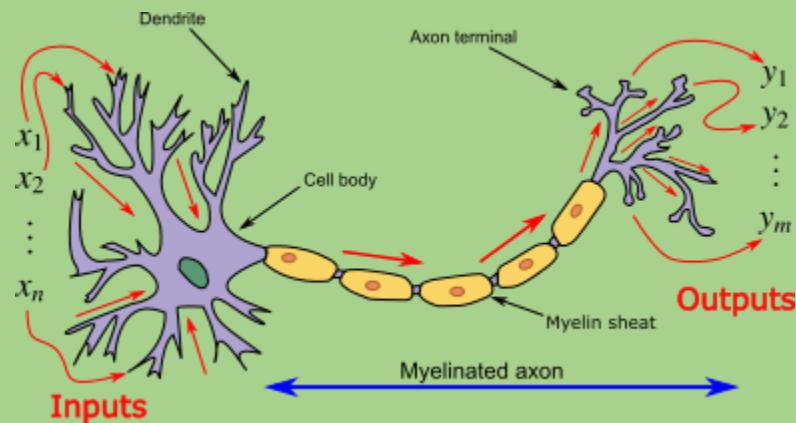


Planning

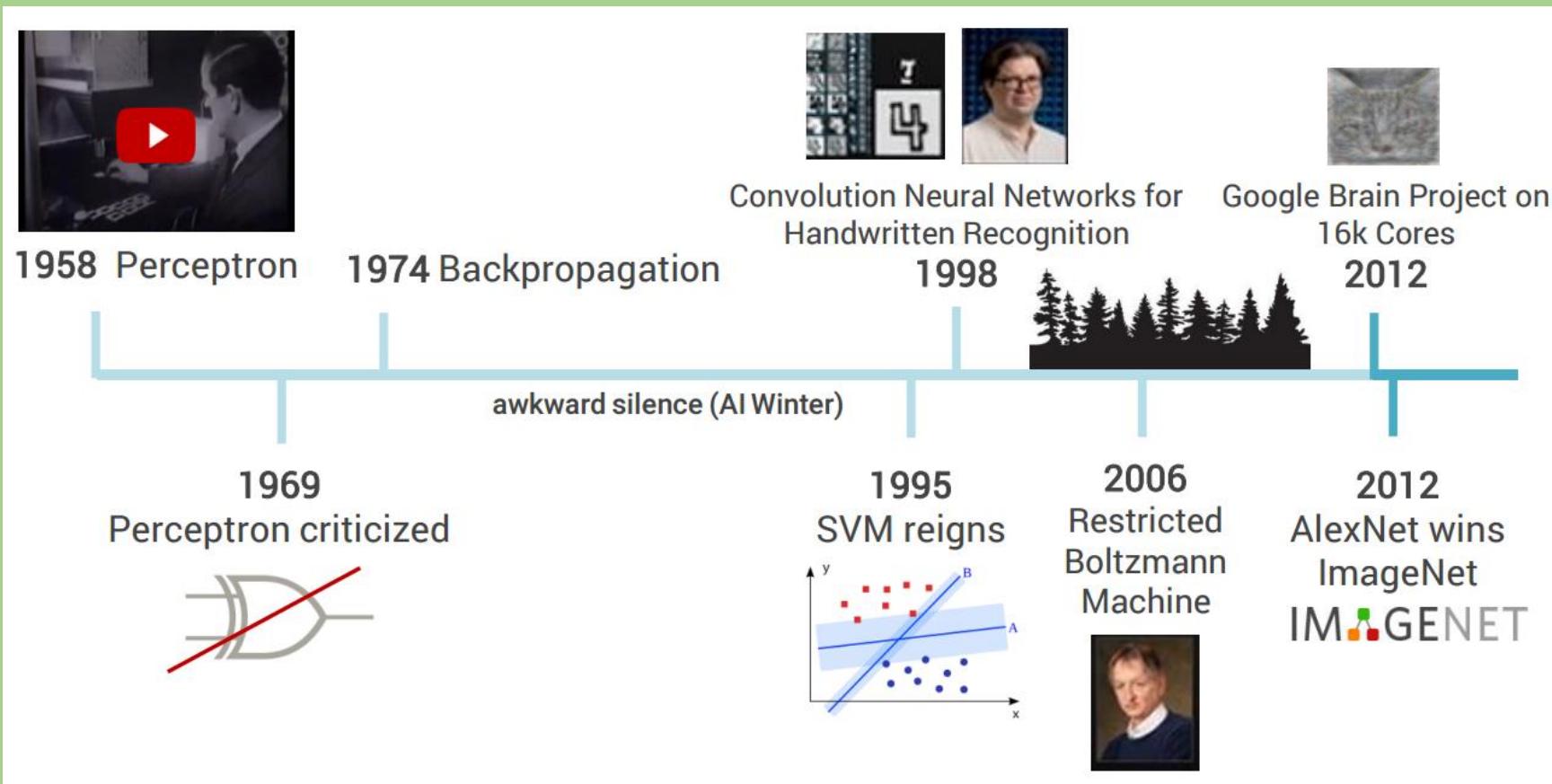
Data Science Lifecycle

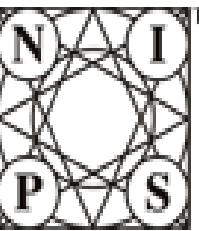


Deep Learning / Neural Networks



Deep Learning





NIPS : Conferences : 2009 : Program

[NIPS Home](#)

[Overview](#)

[Conference Videos](#)

[Workshop Videos](#)

[Program Highlights](#)

[Tutorials](#)

[Conference Sessions](#)

[Workshops](#)

[Publication Models](#)

[Demonstrations](#)

[Mini Symposia](#)

[Accepted Papers](#)

[Dates](#)

[Committees](#)

[Li Deng, Dong Yu, Geoffrey Hinton](#)

[Microsoft Research; Microsoft Research; University of Toronto](#)

[Deep Learning for Speech Recognition and Related Applications](#)

7:30am - 6:30pm Saturday, December 12, 2009

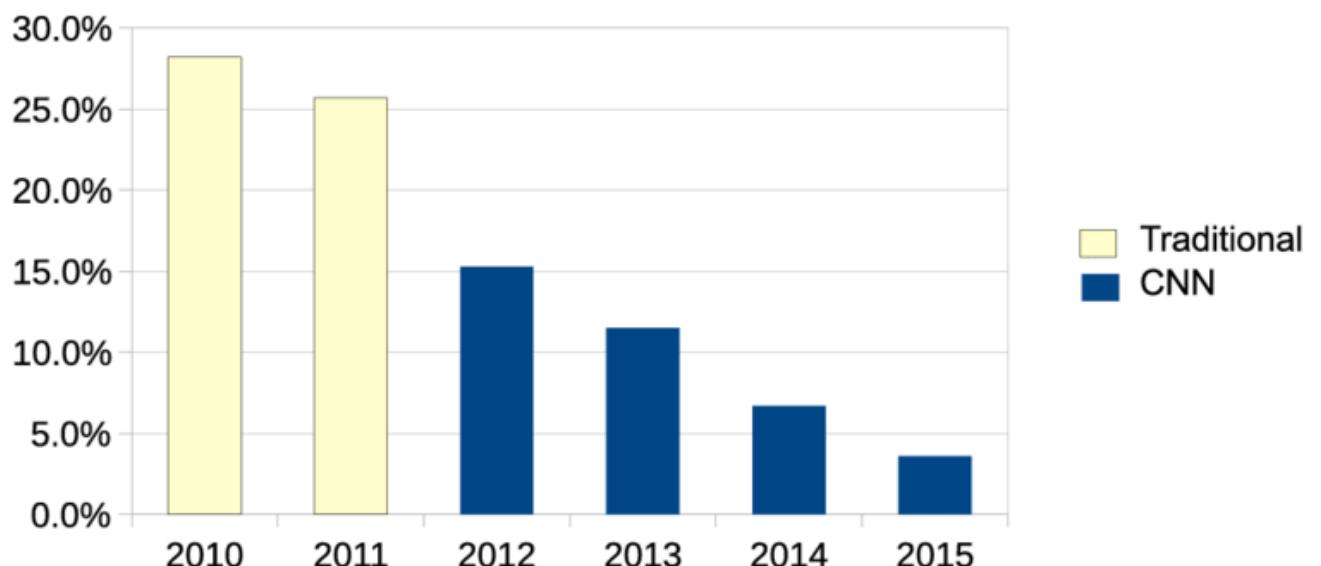
Location: Hilton: Cheakamus

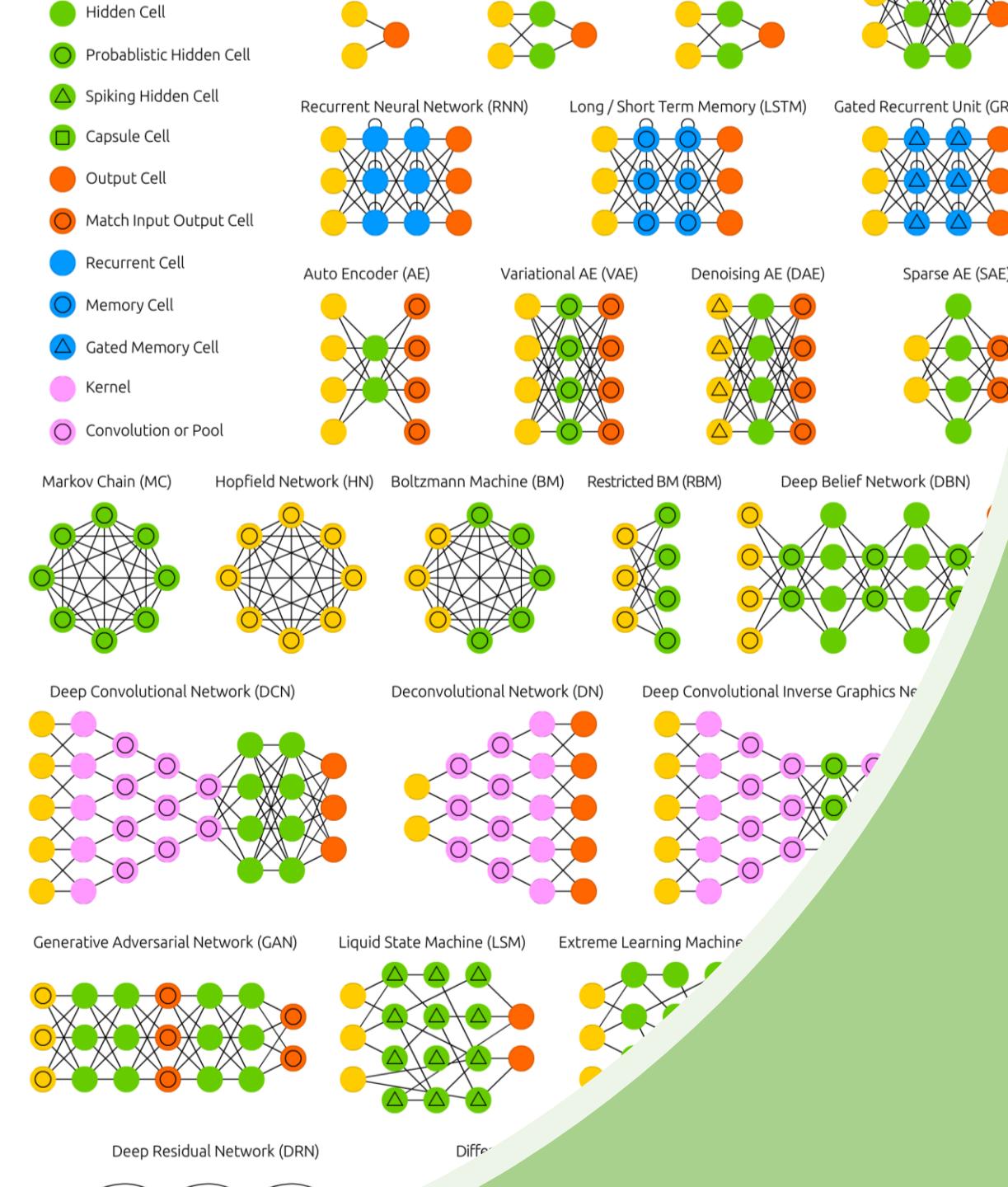
Abstract: Over the past 25 years or so, speech recognition technology has been dominated by a "shallow" architecture -- hidden Markov models (HMMs). Significant technological success has been achieved using complex and carefully engineered variants of HMMs. The next generation of the technology requires solutions to remaining technical challenges under diversified deployment environments. These challenges, not adequately addressed in the past, arise from the many types of variability present in the speech generation process. Overcoming these challenges is likely to require "deep" architectures with efficient learning algorithms. For speech recognition and related sequential pattern recognition applications, some attempts have been made in the past to develop computational architectures that are "deeper" than conventional HMMs, such as

ImageNet Challenge

IMAGENET

- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test.





Neural Networks

Time for some Math

Fish and Chips



x 2



x 1

\$12



x 1

Fish and Chips



x 2

\$?



x 1

\$?



x 1

\$?

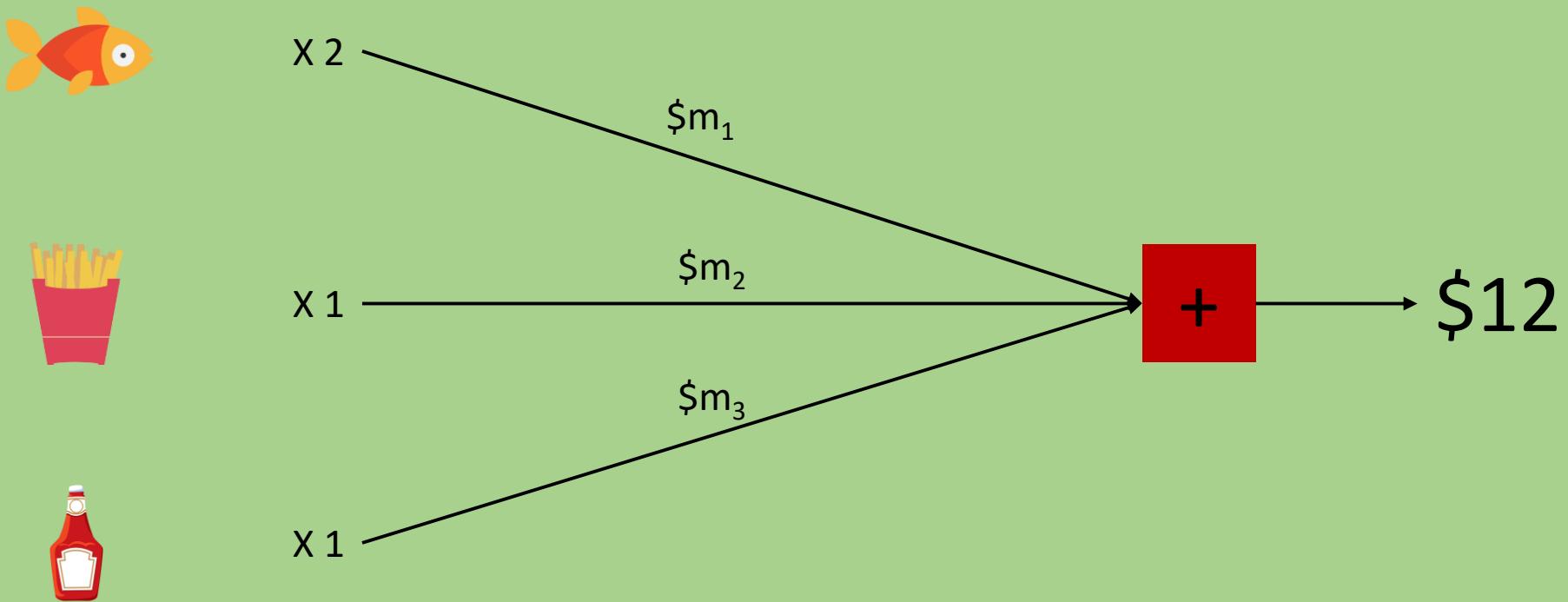
\$12

Dataset

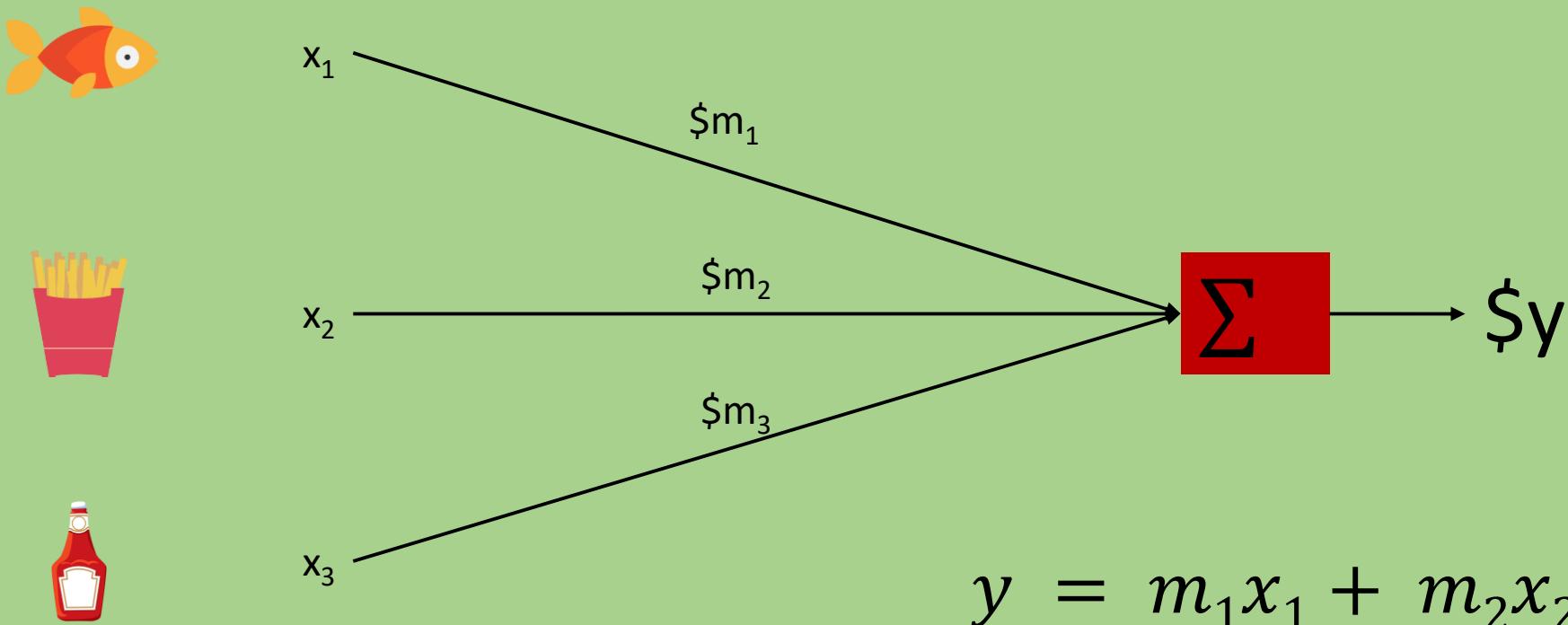


Fish	Chips	Sauce	Total
2	1	1	\$12
3	2	4	\$15
1	1	0	\$7
0	2	2	\$8
1	1	0	\$7

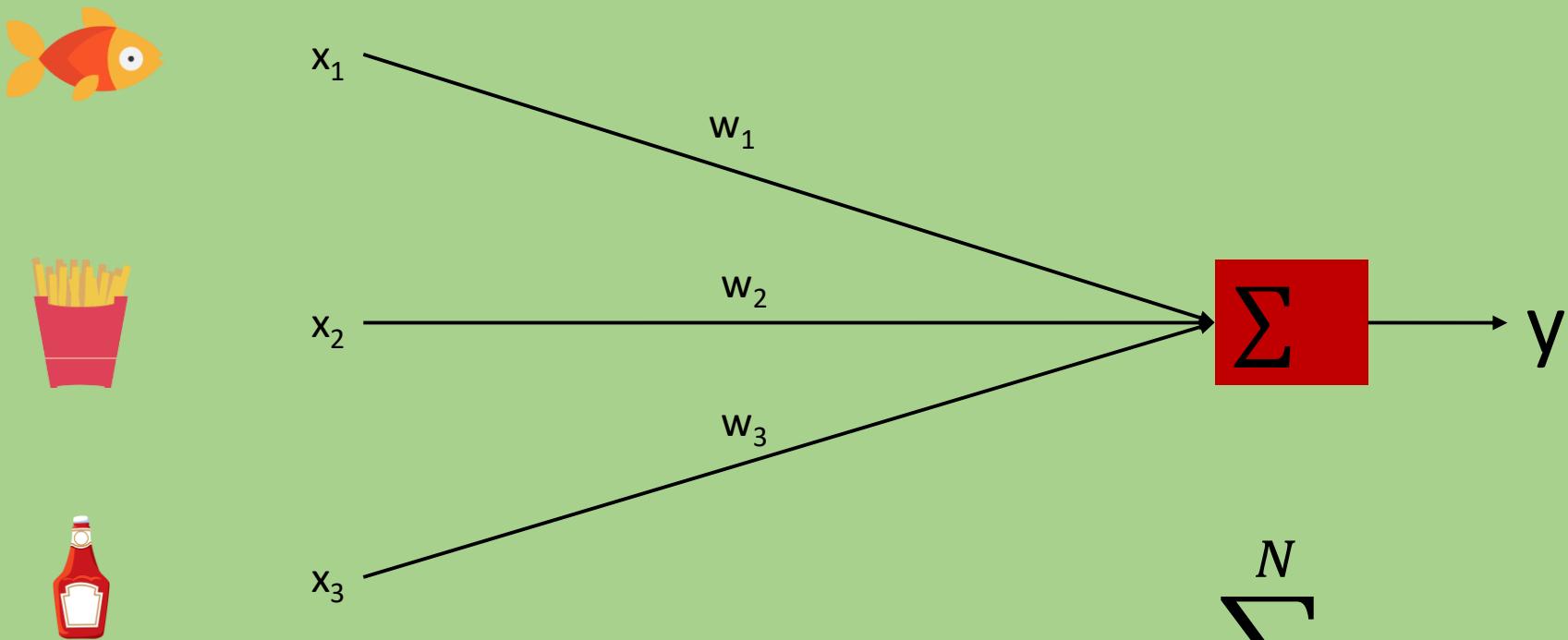
Linear Regression



Linear Regression

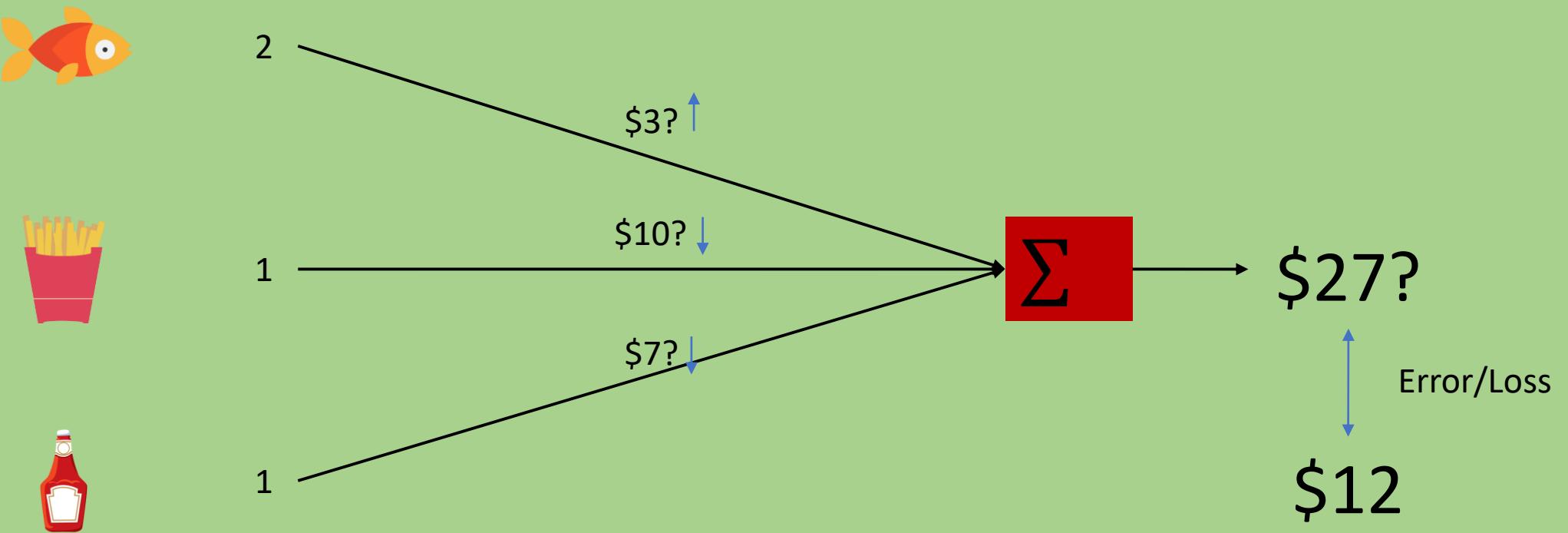


Linear Regression



$$y = \sum_{i=1}^N w_i x_i$$

Fish and Chips



Loss Functions - Regression

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}|$$

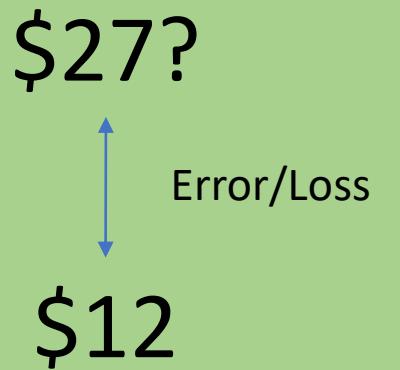
$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2$$

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2}$$

$$R^2 = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2}$$

Where,

\hat{y} – predicted value of y
 \bar{y} – mean value of y



Error Metrics

$$MAE = \frac{1}{n} \sum \left| y - \hat{y} \right|$$

Divide by the total number of data points

Predicted output value

Actual output value

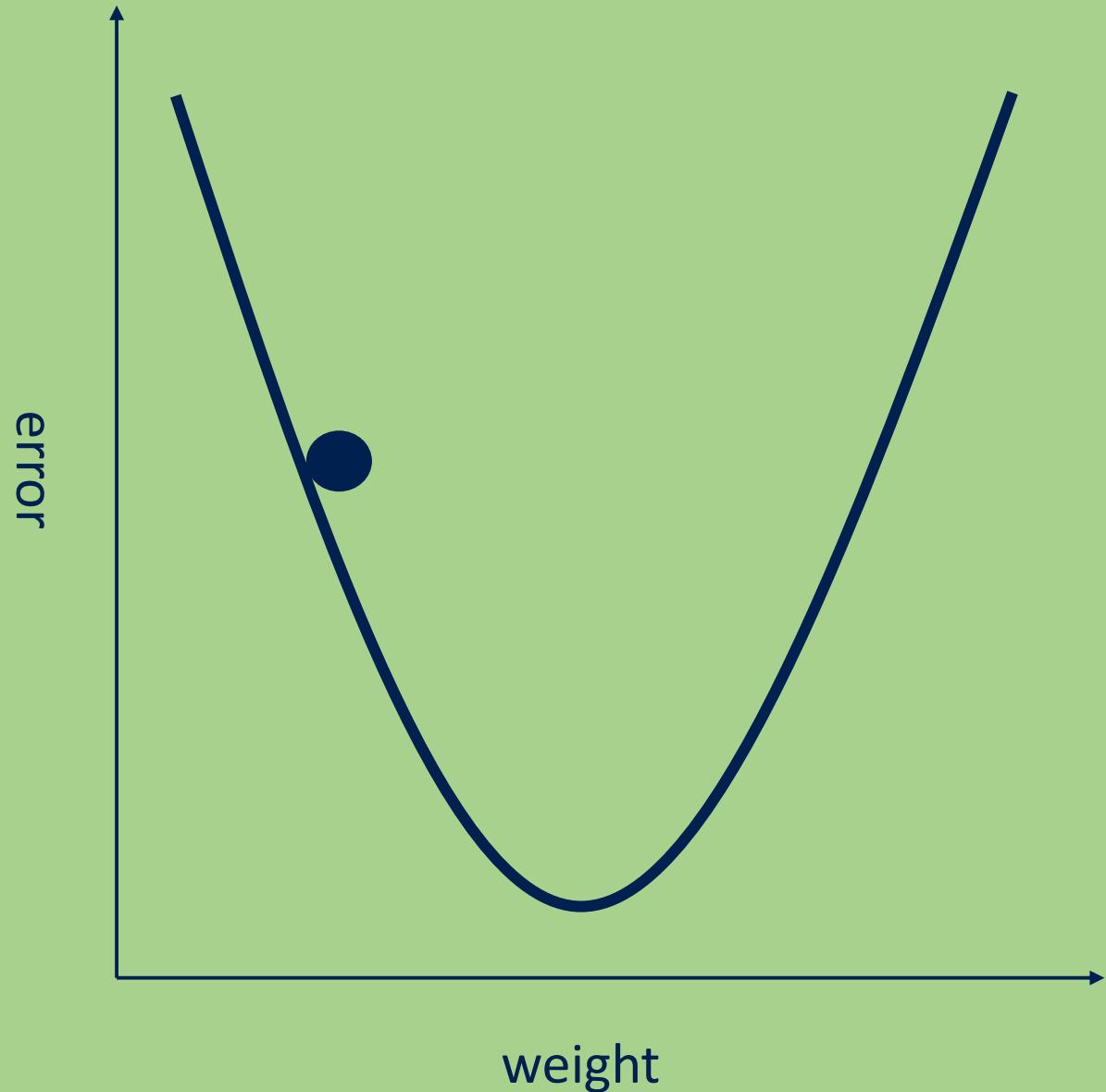
Sum of

The absolute value of the residual

\$27?
↓
Error/Loss
\$12

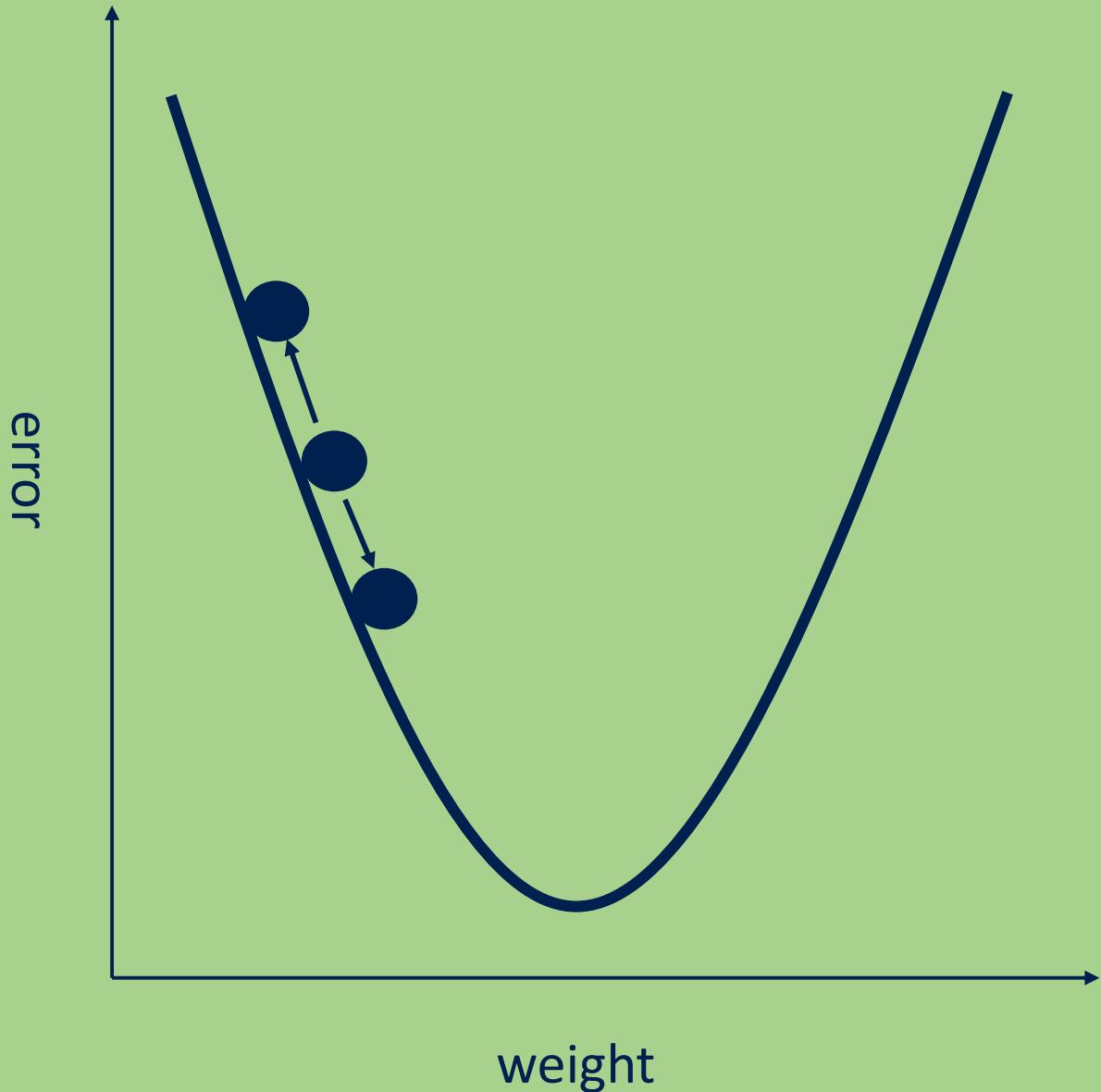
Gradient descent

For each input feature and voting weight, adjust it up and down a bit and see how the error changes.

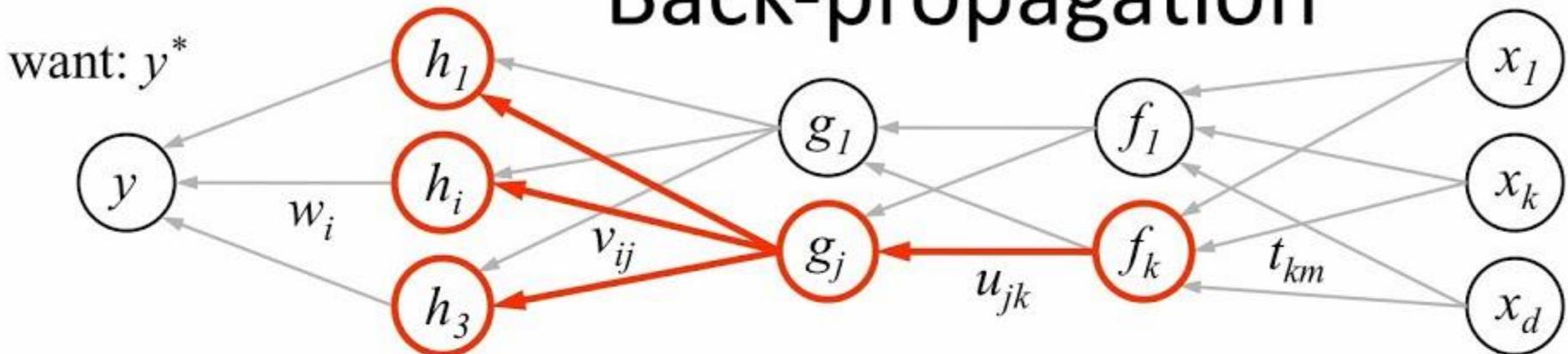


Gradient descent

For each input feature and voting weight, adjust it up and down a bit and see how the error changes.



Back-propagation



1. receive new observation $\mathbf{x} = [x_1 \dots x_d]$ and target y^*
2. **feed forward:** for each unit g_j in each layer $1 \dots L$
compute g_j based on units f_k from previous layer: $g_j = \sigma\left(u_{j0} + \sum_k u_{jk} f_k\right)$
3. get prediction y and error $(y - y^*)$
4. **back-propagate error:** for each unit g_j in each layer $L \dots 1$

(a) compute error on g_j

$$\frac{\partial E}{\partial g_j} = \sum_i \underbrace{\sigma'(h_i)}_{\text{should } g_j \text{ be higher or lower?}} \underbrace{v_{ij}}_{\text{how } h_i \text{ will change as } g_j \text{ changes}} \underbrace{\frac{\partial E}{\partial h_i}}_{\text{was } h_i \text{ too high or too low?}}$$

(b) for each u_{jk} that affects g_j

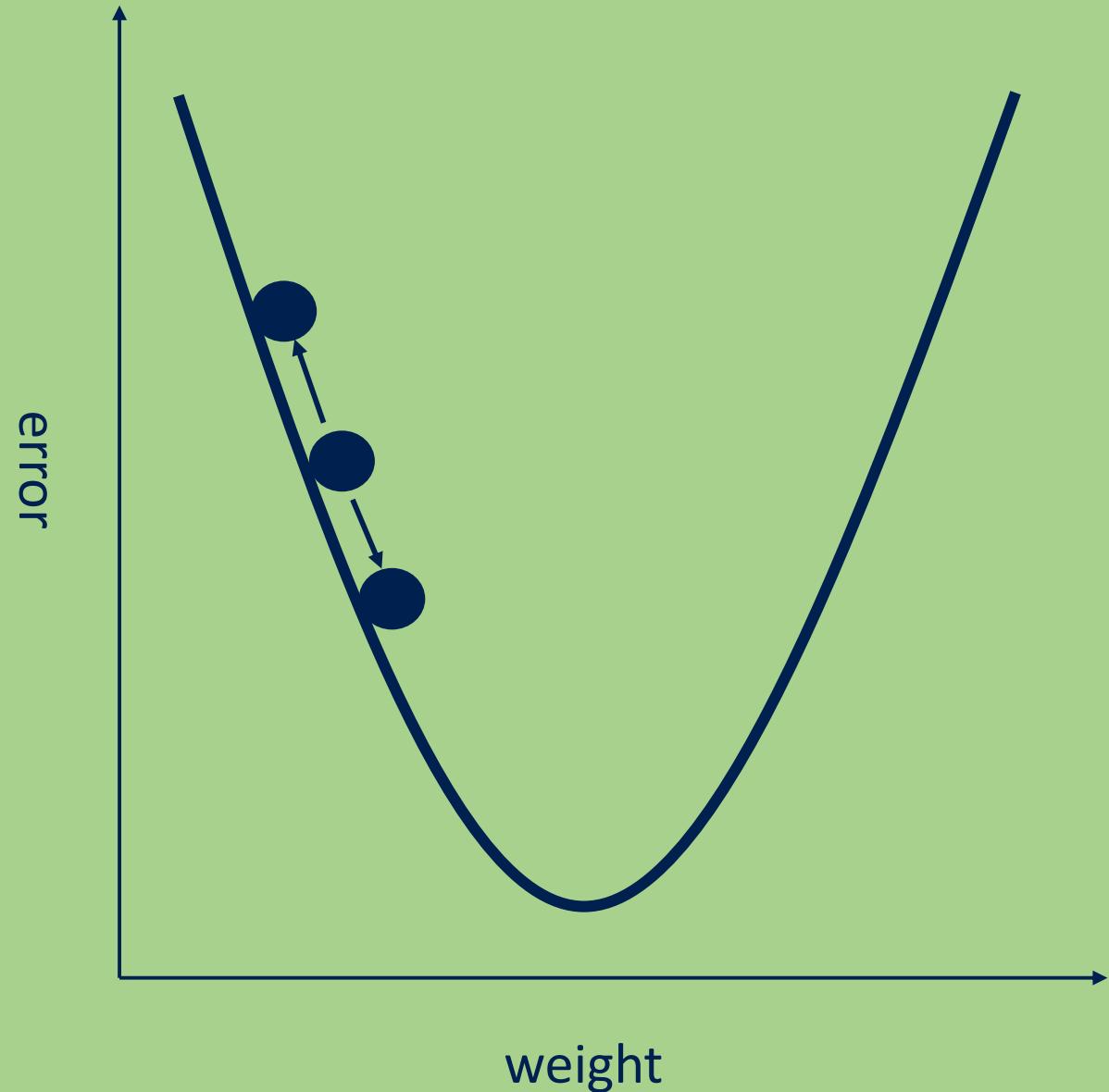
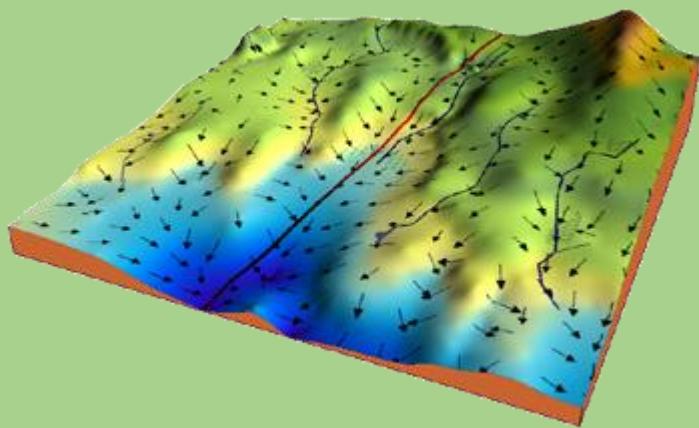
(i) compute error on u_{jk}

$$\frac{\partial E}{\partial u_{jk}} = \underbrace{\frac{\partial E}{\partial g_j}}_{\text{do we want } g_j \text{ to be higher/lower?}} \underbrace{\sigma'(g_j) f_k}_{\text{how } g_j \text{ will change if } u_{jk} \text{ is higher/lower}}$$

(ii) update the weight

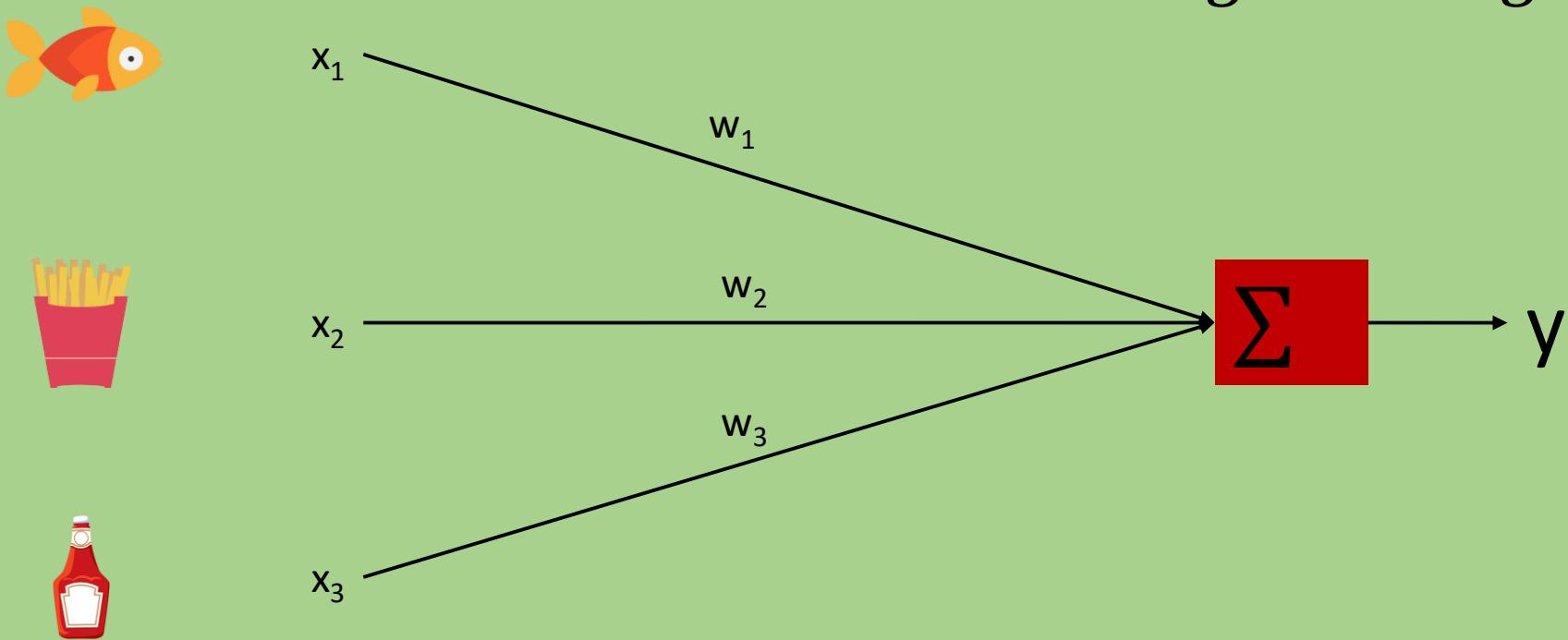
$$u_{jk} \leftarrow u_{jk} - \eta \frac{\partial E}{\partial u_{jk}}$$

Gradient descent



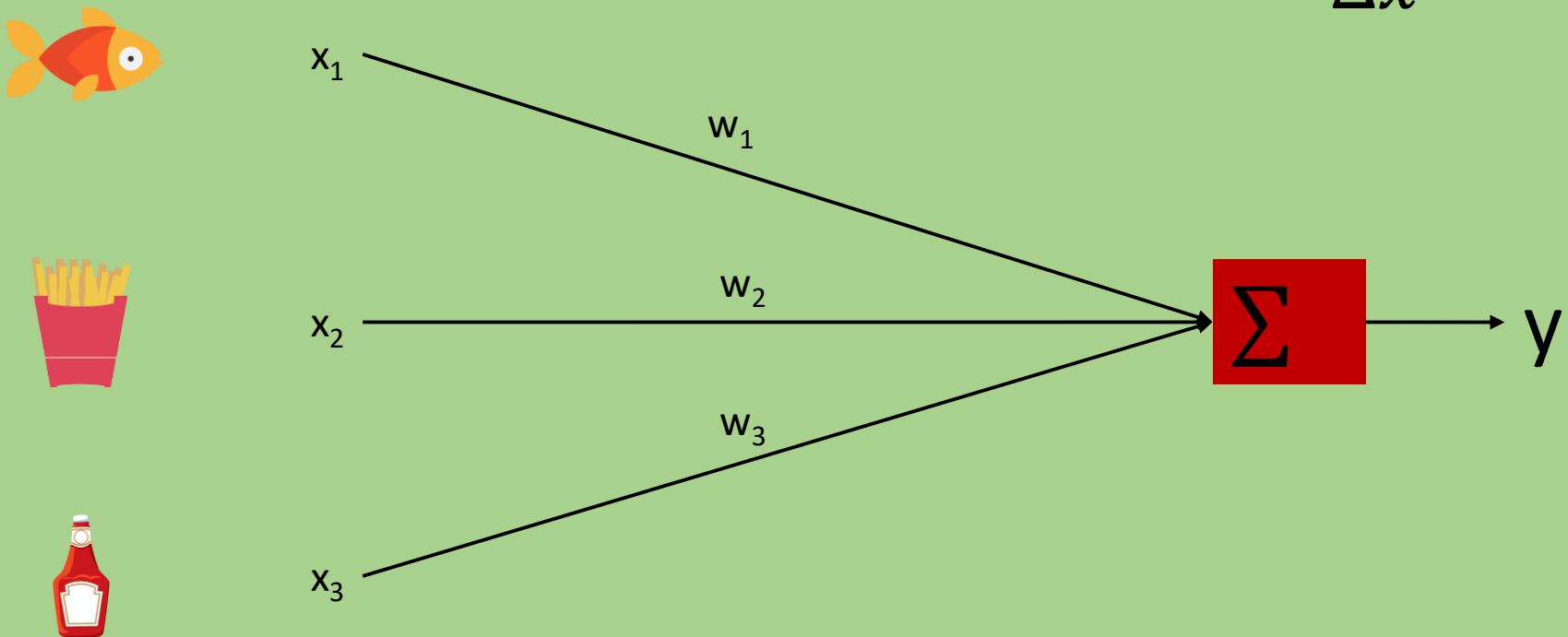
Weight Update

$$w^{new} = w^{old} + learning\ rate * gradient$$



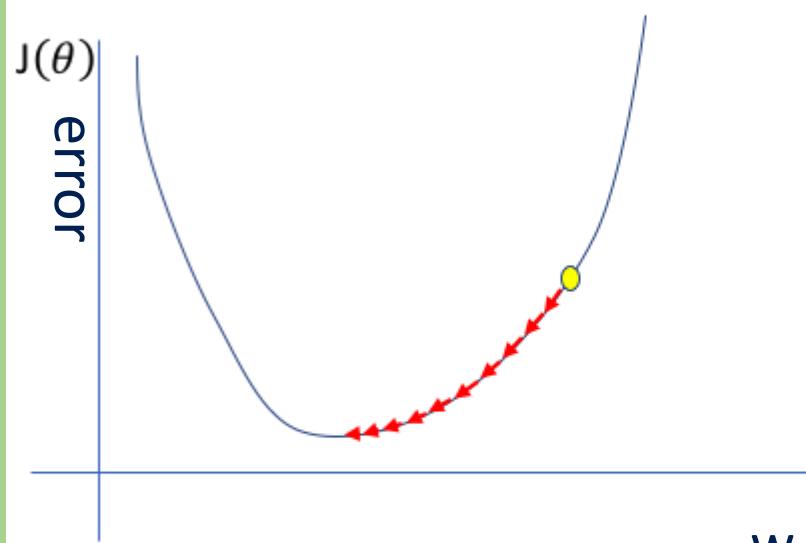
Weight Update

$$w^{new} = w^{old} + \alpha \frac{\Delta y}{\Delta x}$$



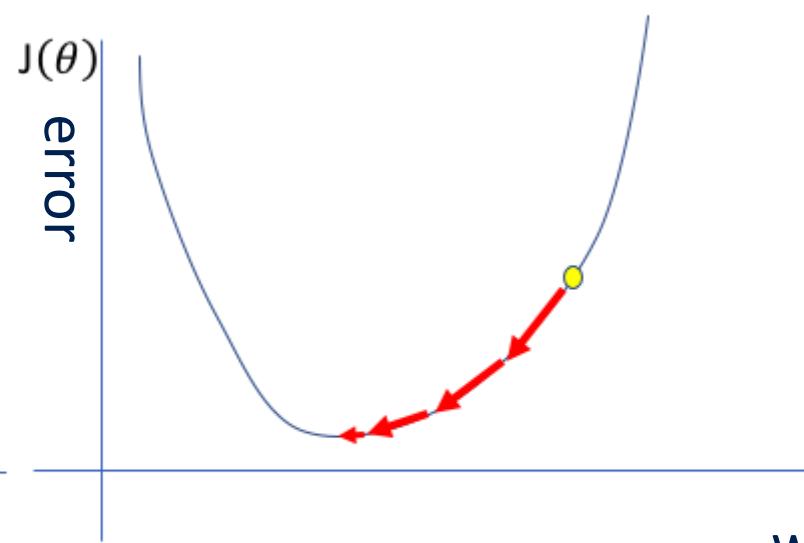
Learning Rate

Too low



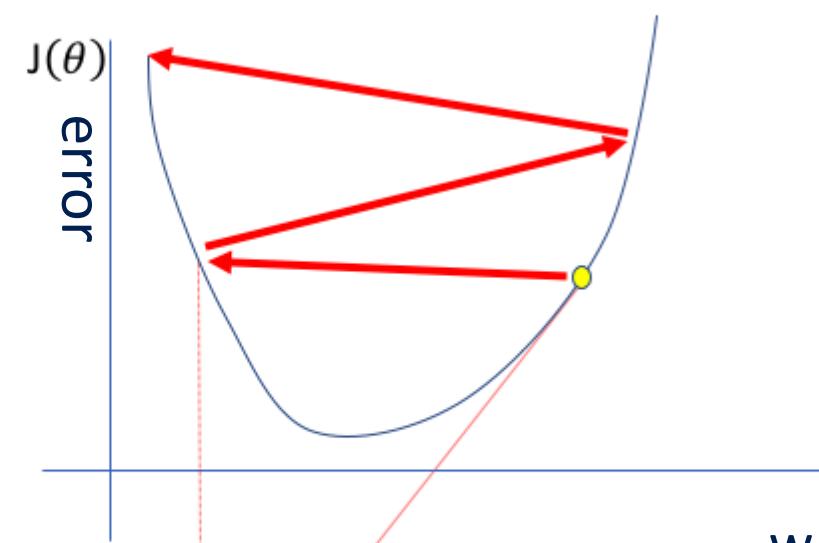
A small learning rate requires many updates before reaching the minimum point

Just right



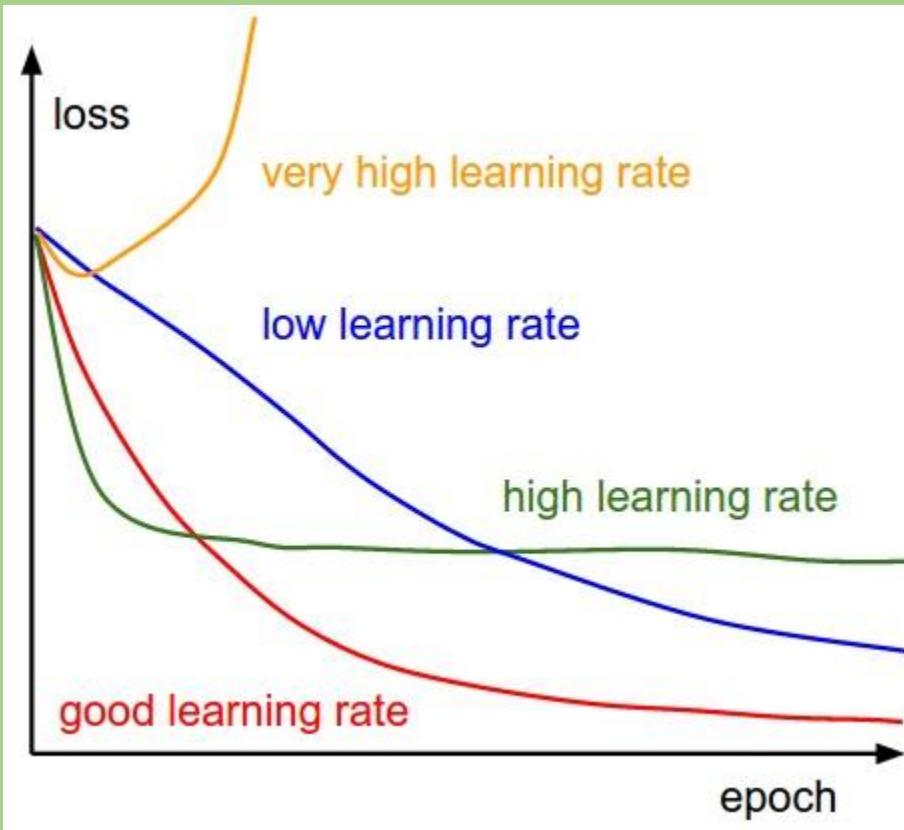
The optimal learning rate swiftly reaches the minimum point

Too high

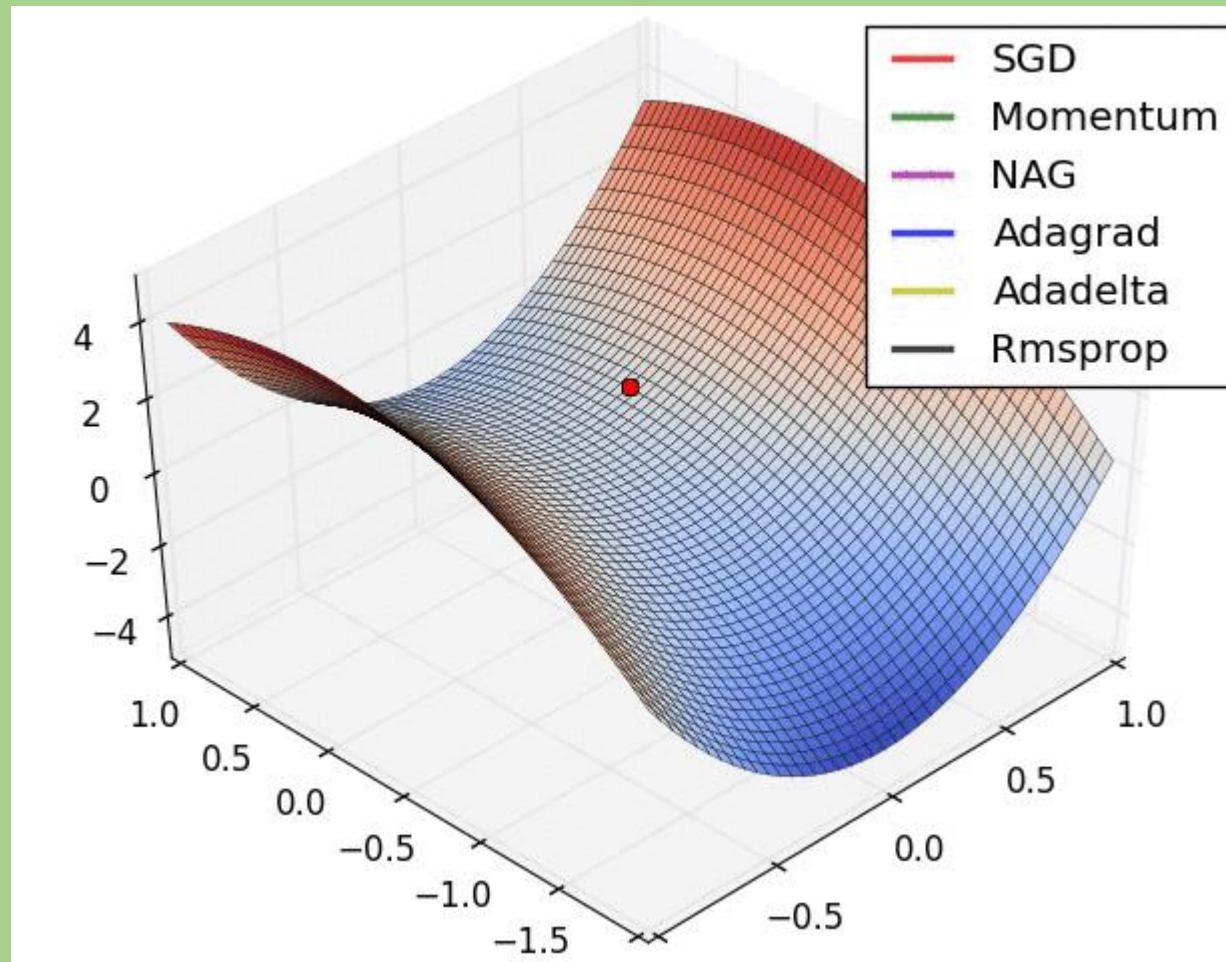


Too large of a learning rate causes drastic updates which lead to divergent behaviors

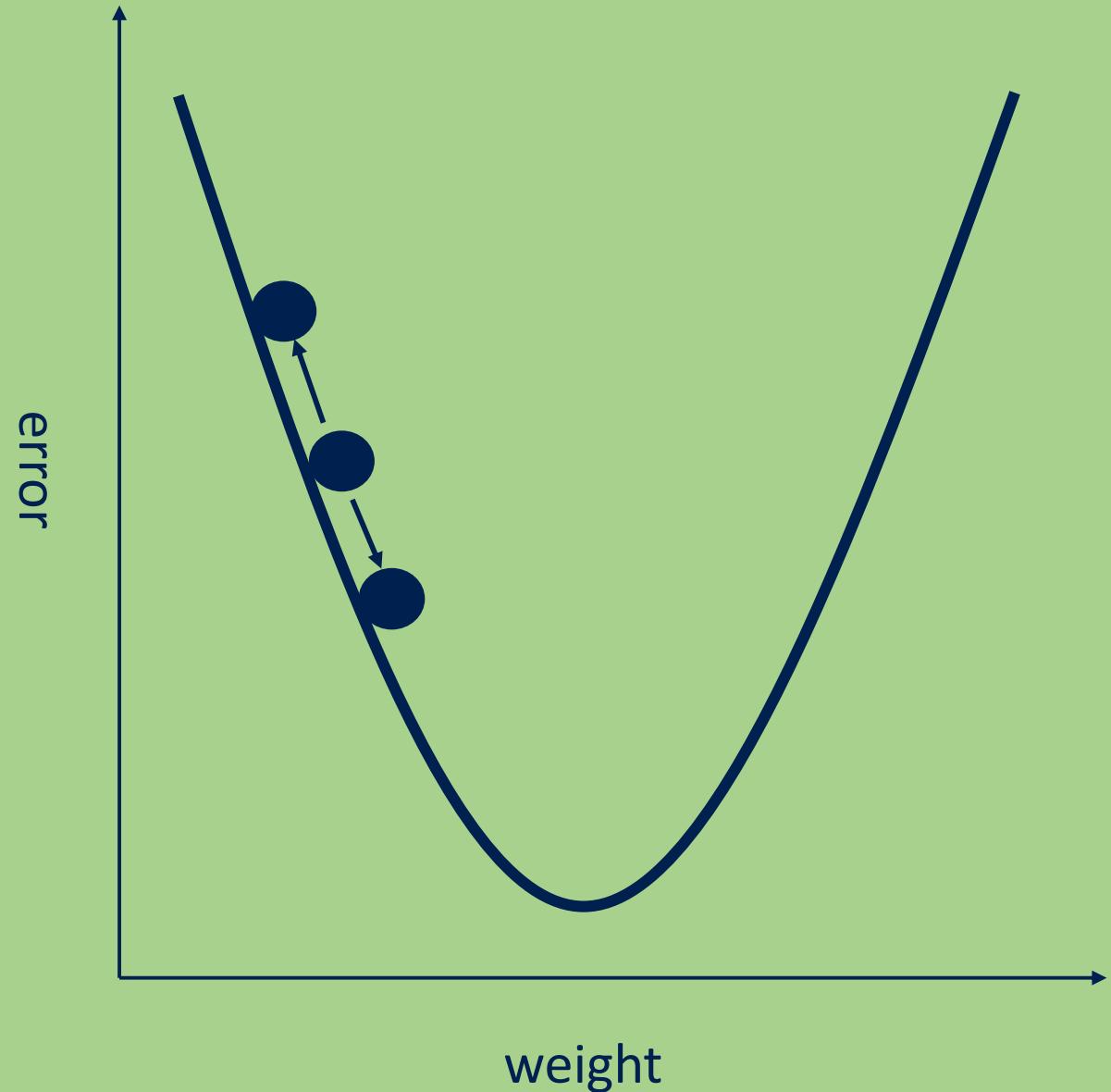
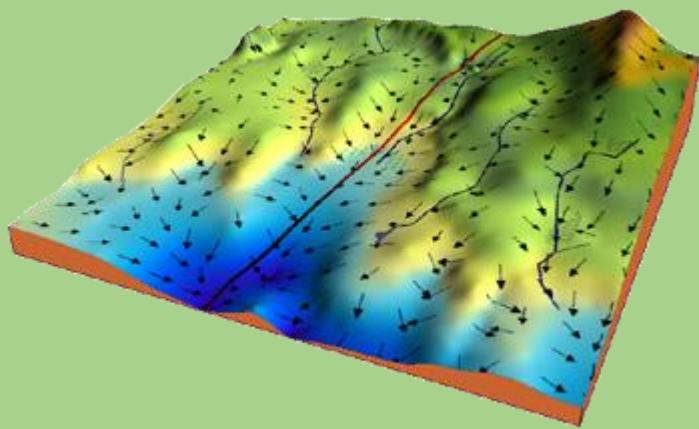
Learning Rate



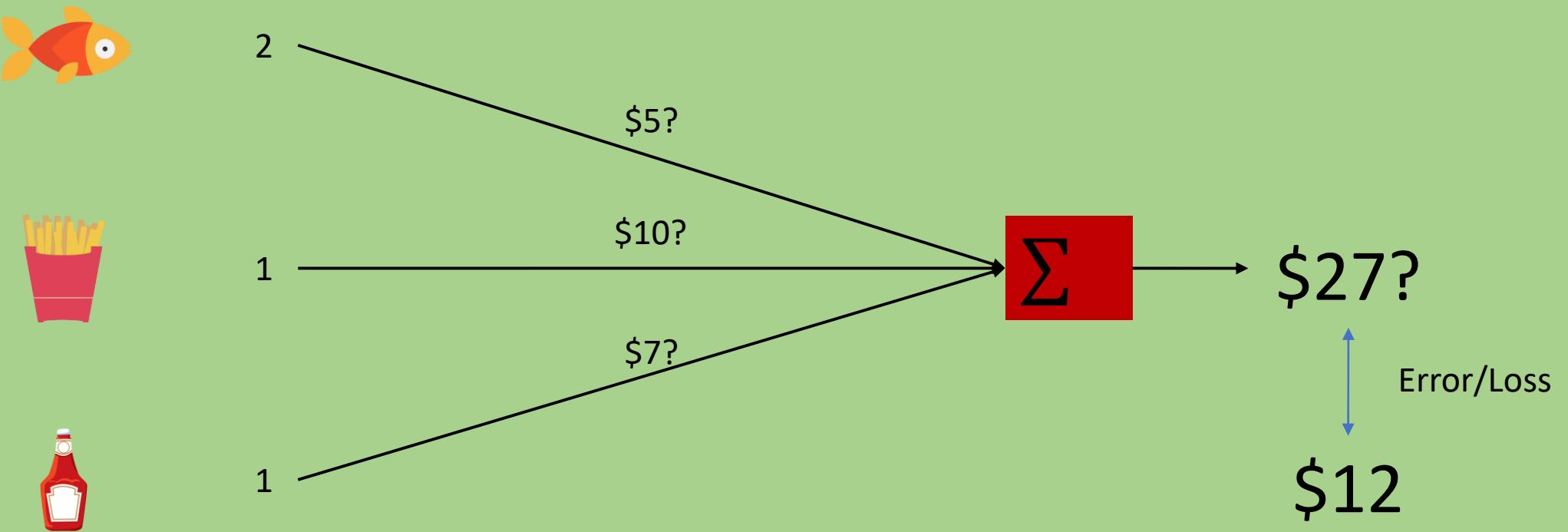
Different Optimizers



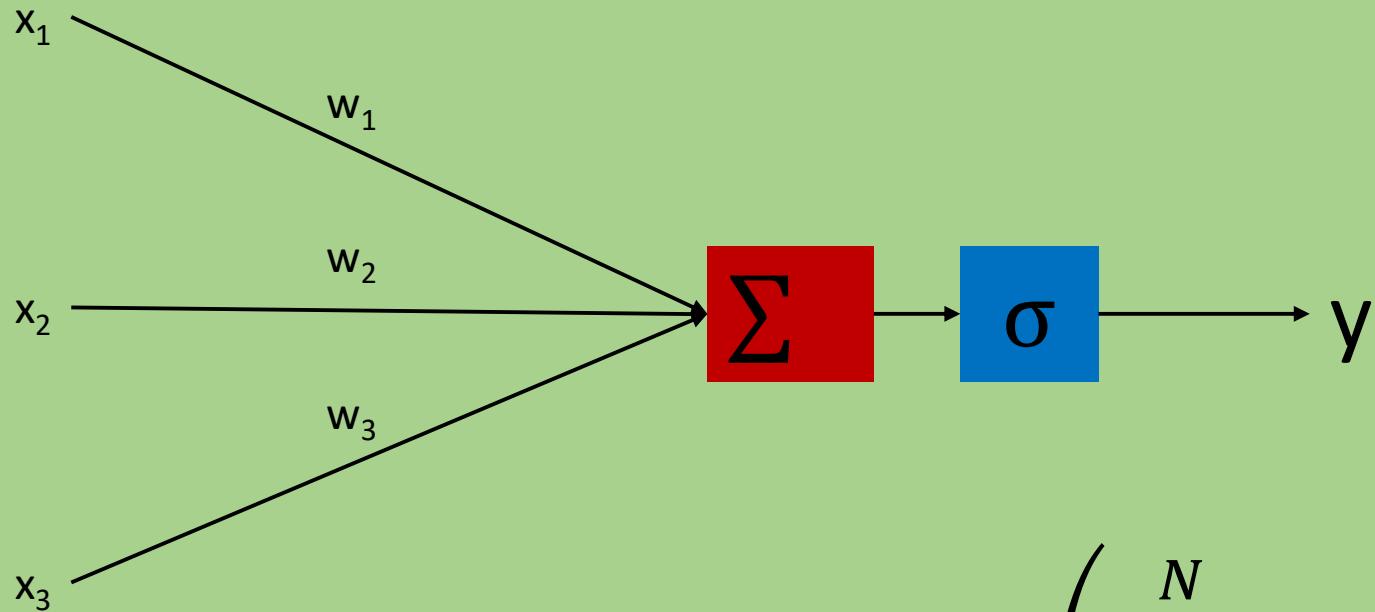
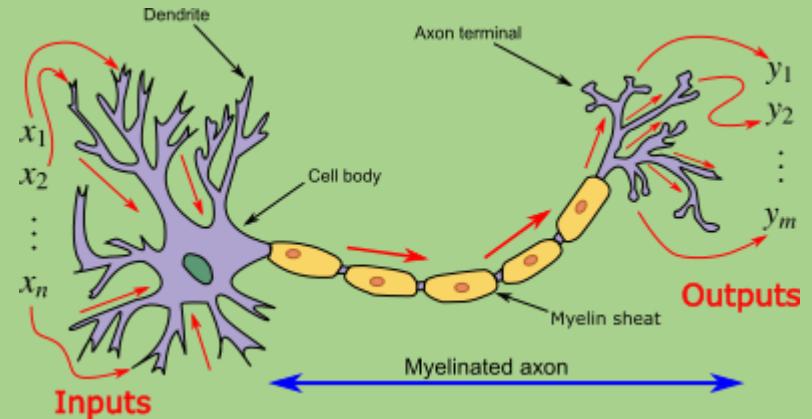
Gradient descent



Fish and Chips

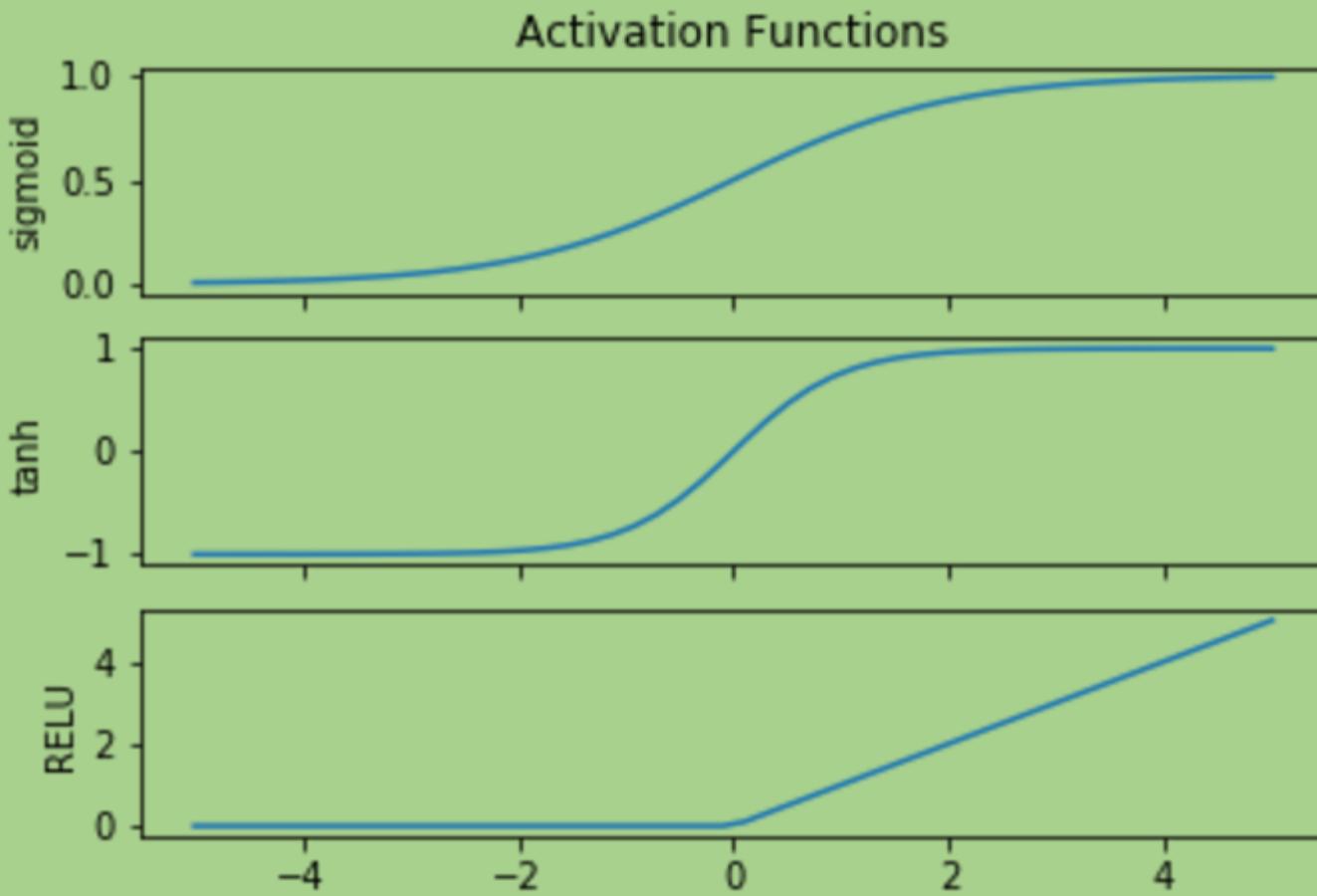


Neuron



$$y = \sigma \left(\sum_{i=1}^N w_i x_i \right)$$

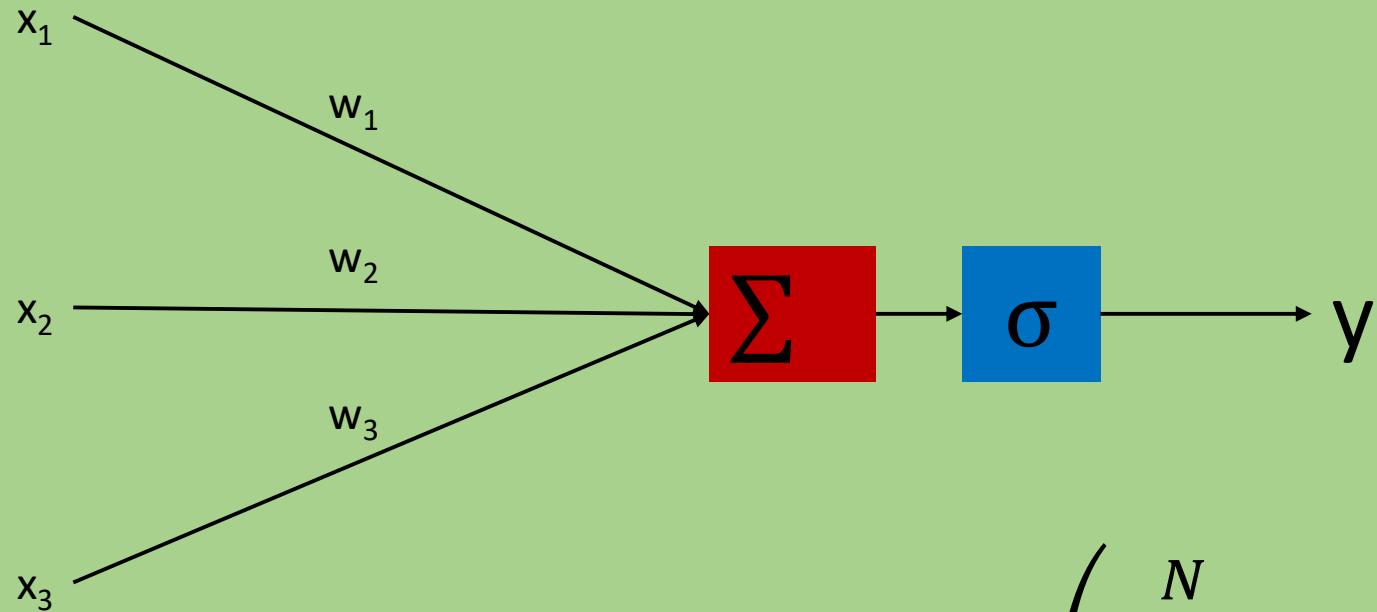
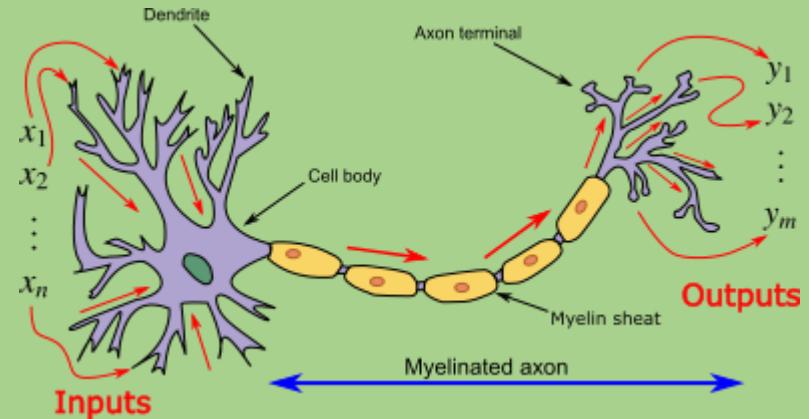
σ



$$x \leftarrow \sum_{i=1}^N m_i x_i$$

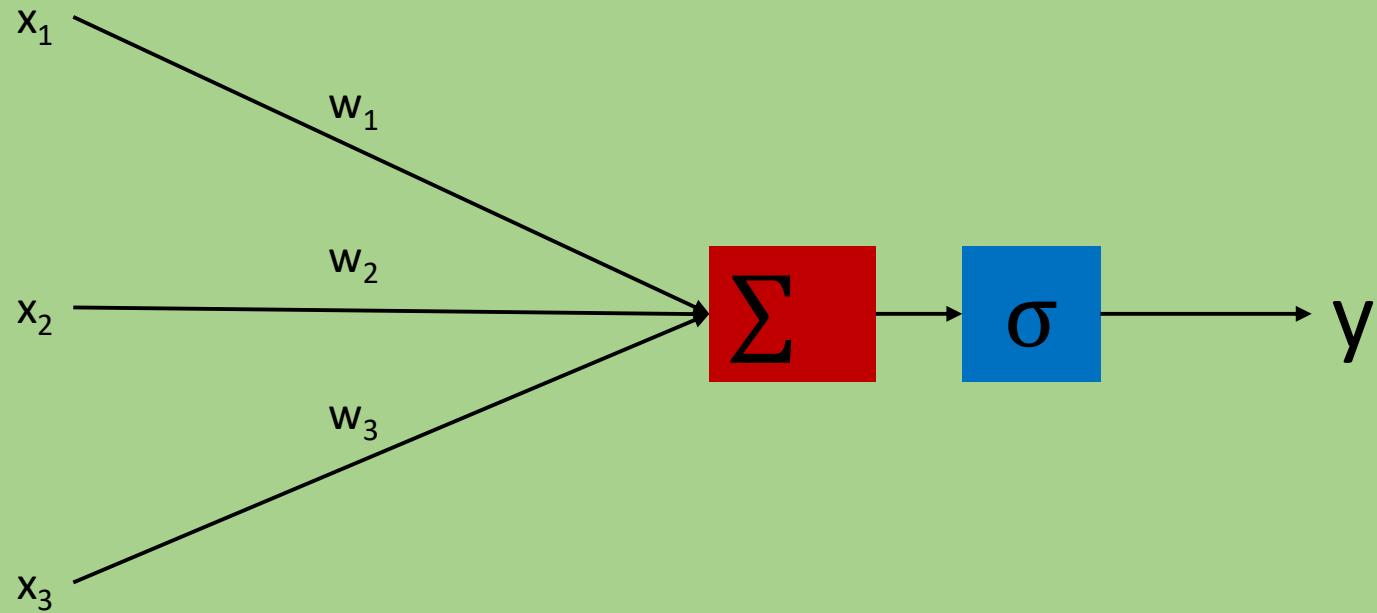
$$y = \sigma \left(\sum_{i=1}^N m_i x_i \right)$$

Neuron

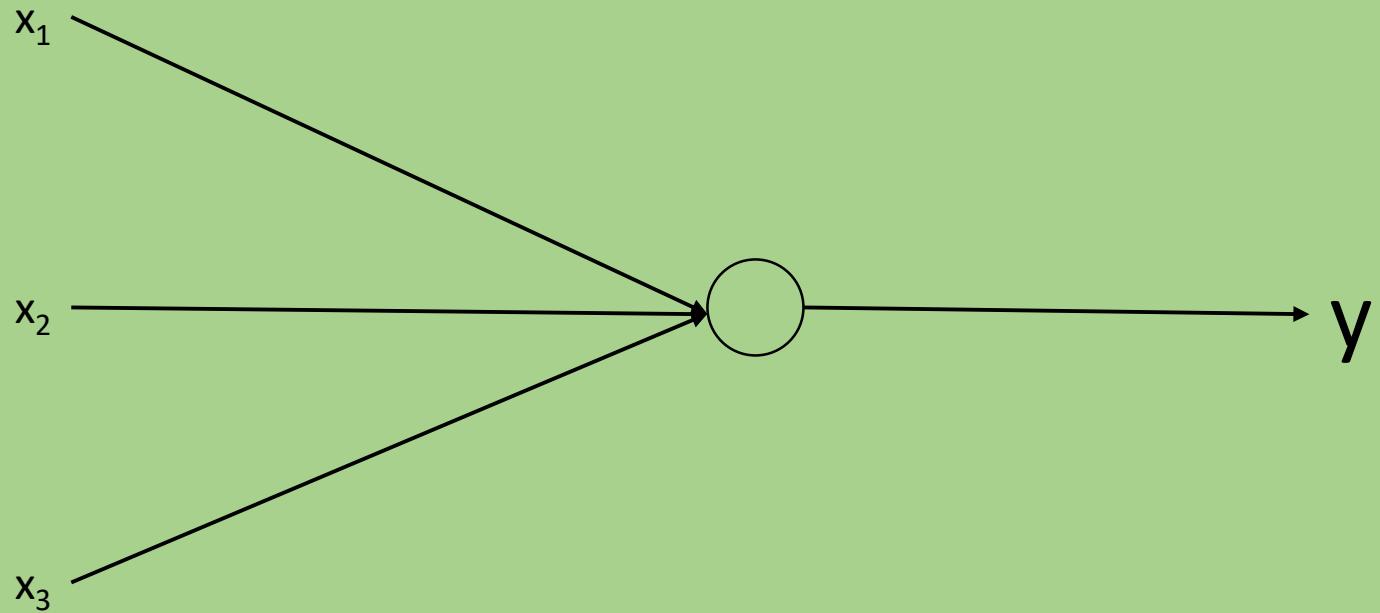


$$y = \sigma \left(\sum_{i=1}^N w_i x_i \right)$$

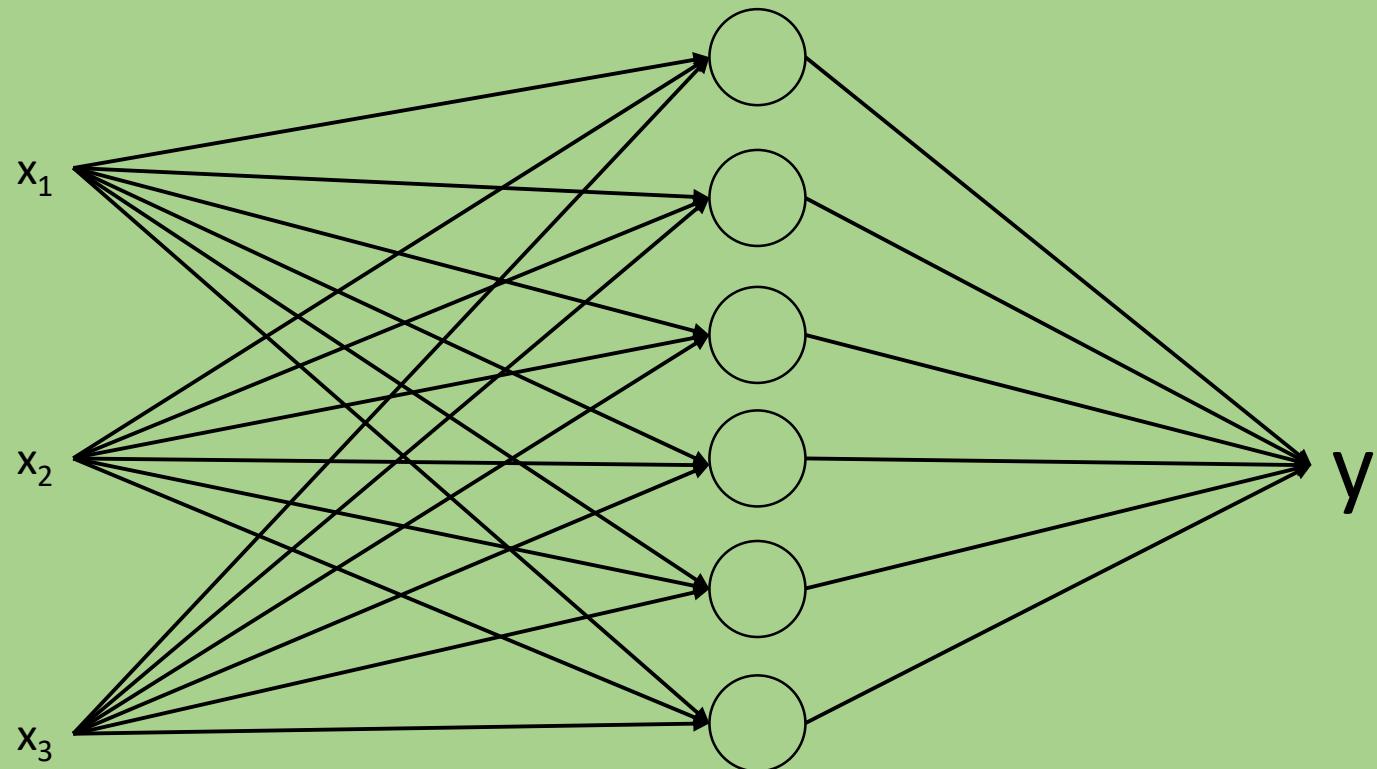
Neuron

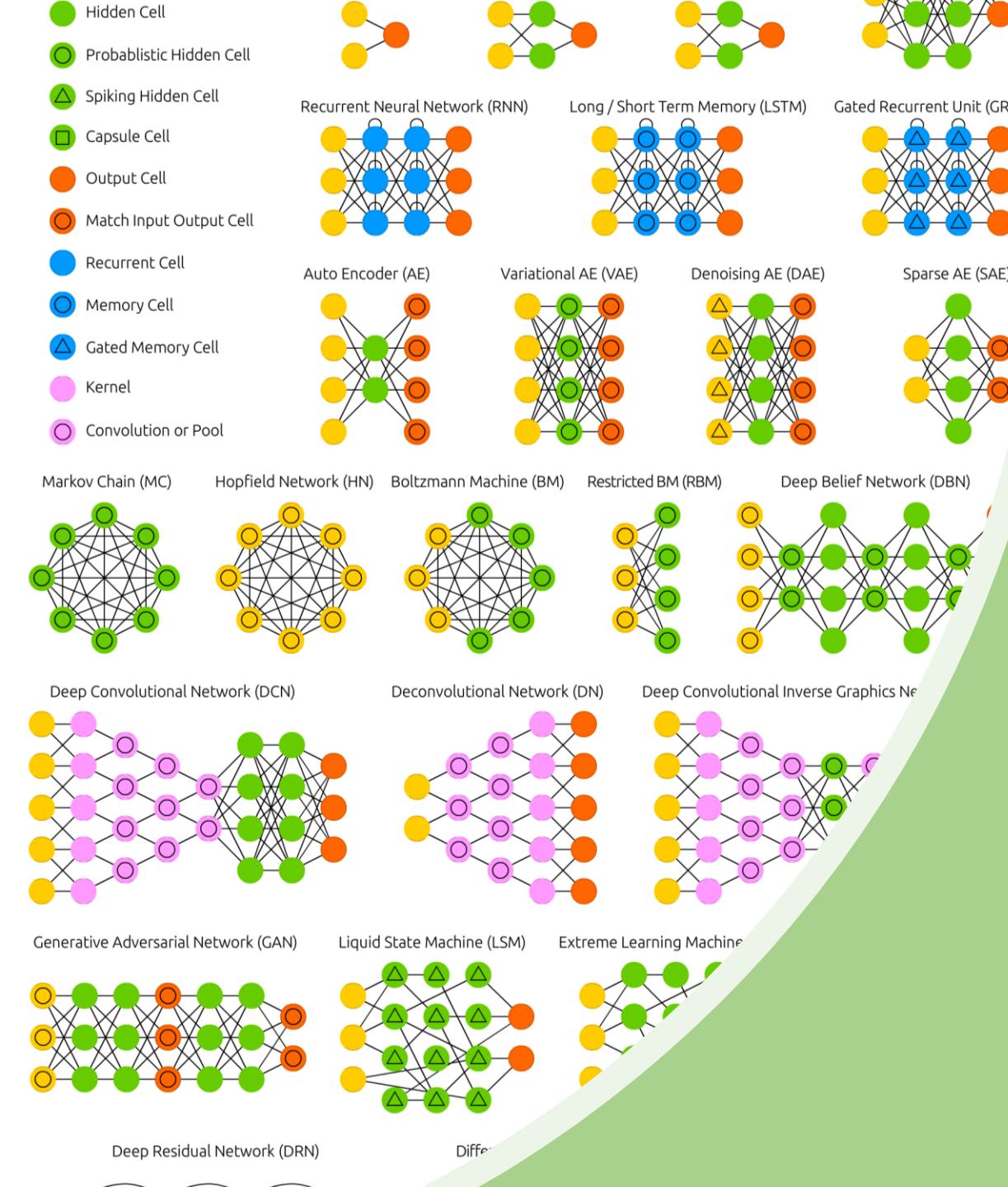


Neuron



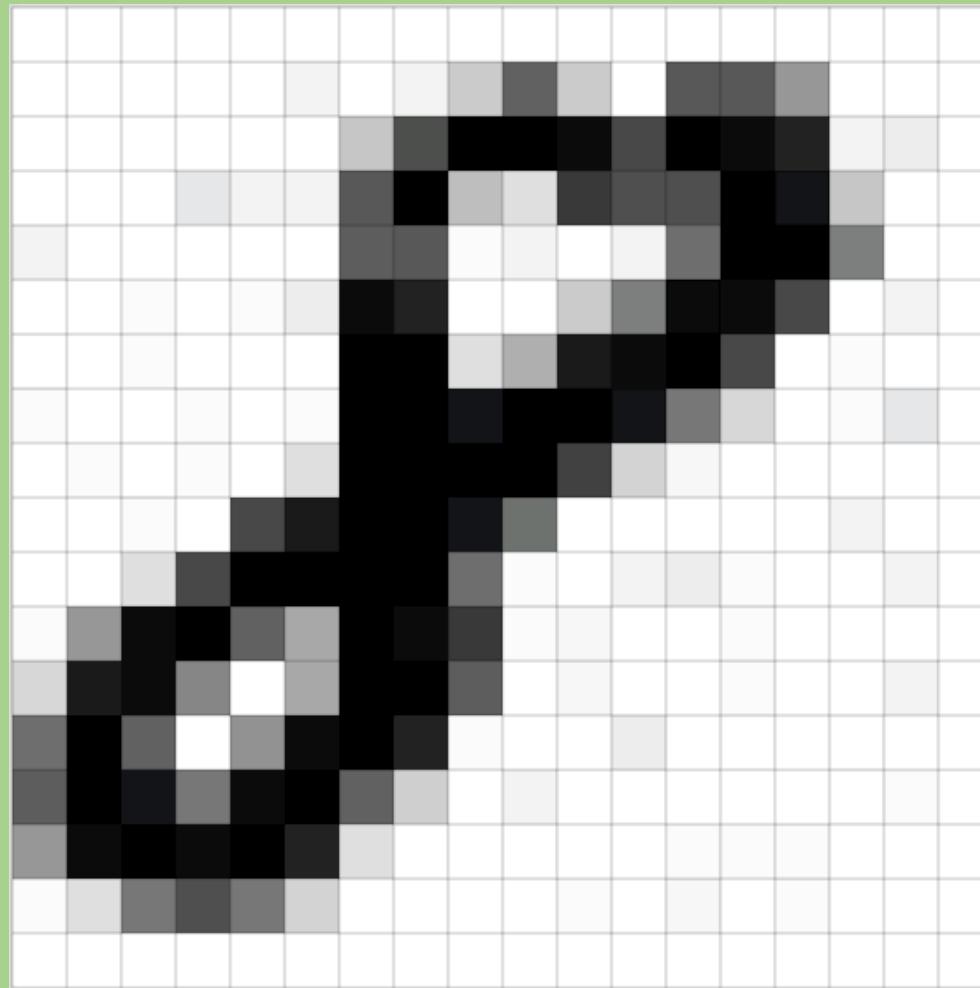
Neural Network





Neural Networks

Can only deal with numbers



The Cat sat on the Mat

One-hot encoding

cat mat on sat the

the =>

0	0	0	0	1
---	---	---	---	---

cat =>

1	0	0	0	0
---	---	---	---	---

sat =>

0	0	0	1	0
---	---	---	---	---

...

...

The Cat sat on the Mat



The Dog sat on my bed



on	mat	sat	the	cat	my	dog	bed
0	1	1	1	2	1	0	0
1	1	0	1	1	0	1	1

TFIDF

- Rather than just counting, we can use the TF-IDF score of a word to rank its importance.

TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

Term frequency
Number of times term t appears in a doc, d

Inverse document frequency
 $\log \frac{1 + n}{1 + df(d, t)}$
 $n \leftarrow$ # of documents
Document frequency of the term t

TFIDF

The Cat sat on the Mat



The Dog sat on my bed



	on	mat	sat	the	cat	my	dog	bed
0	1	1	1	2	1	0	0	0
1	1	0	1	1	0	1	1	1

The Cat sat on the Mat

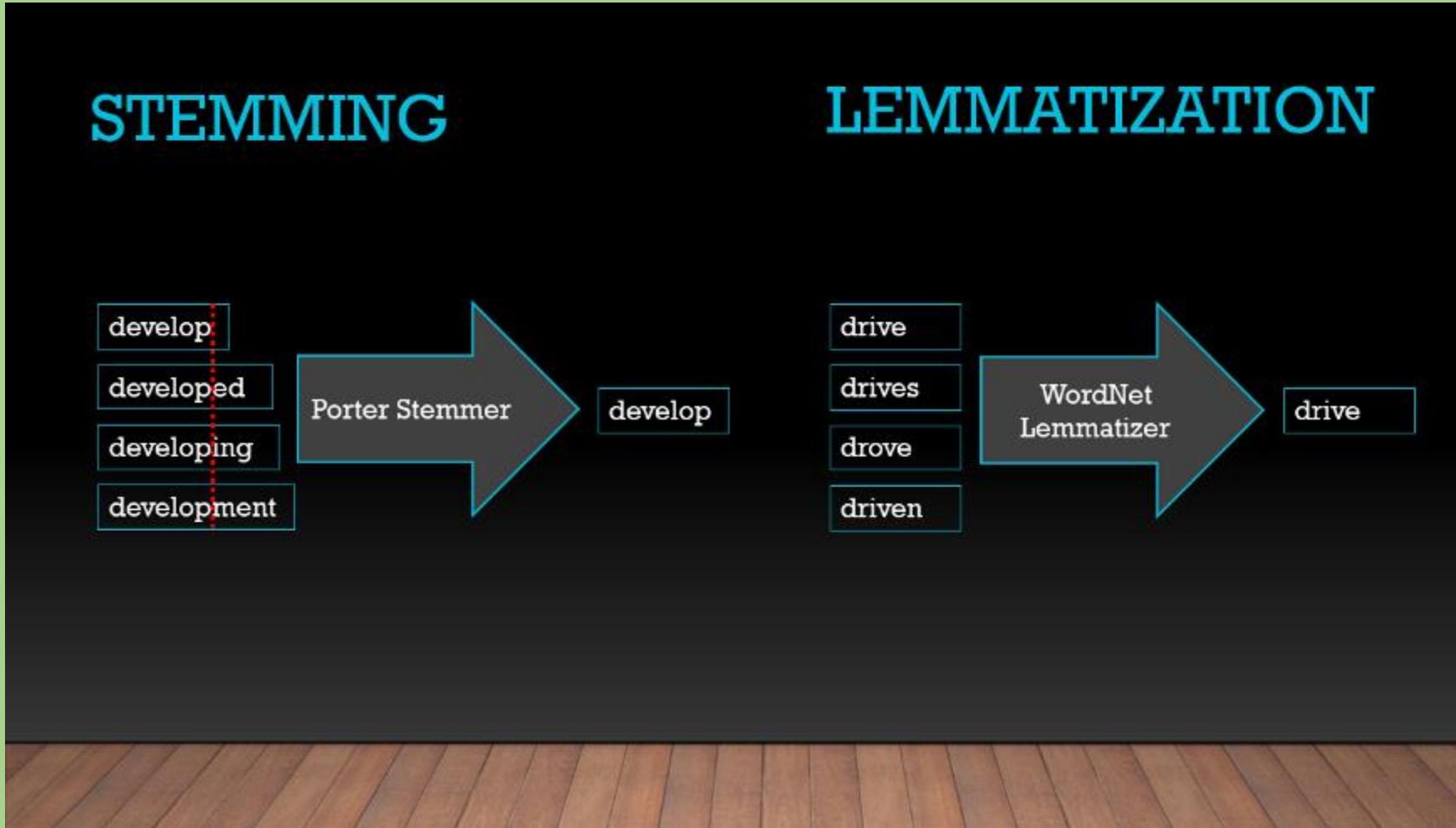


The Dog sat on my bed



	on	mat	sat	the	cat	my	dog	bed
0	0.0	0.115525	0.0	0.0	0.115525	0.000000	0.000000	0.000000
1	0.0	0.000000	0.0	0.0	0.000000	0.115525	0.115525	0.115525

Other NLP techniques



The Cat sat on the Mat

A 4-dimensional embedding

cat =>

1.2	-0.1	4.3	3.2
0.4	2.5	-0.9	0.5
2.1	0.3	0.1	0.4

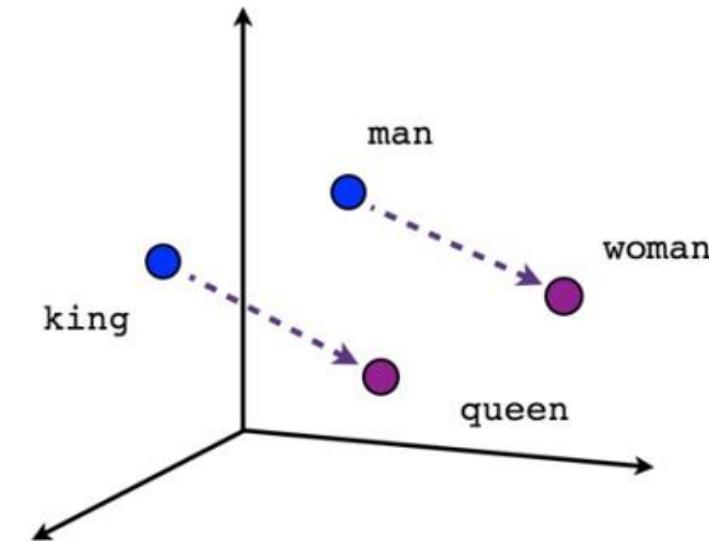
mat =>

on =>

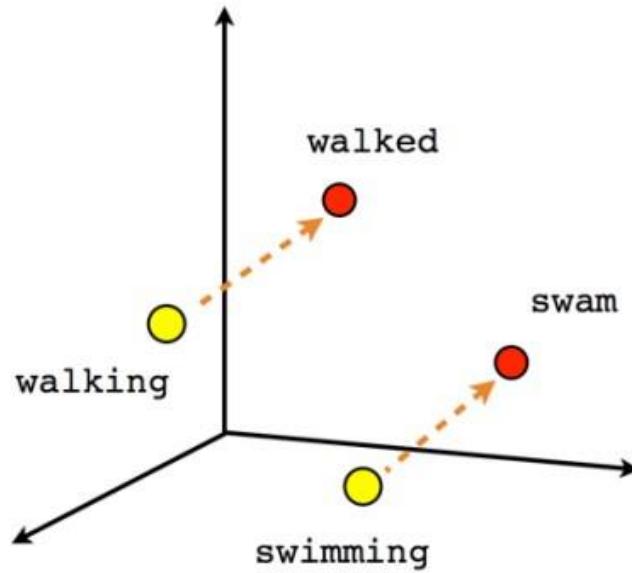
...

...

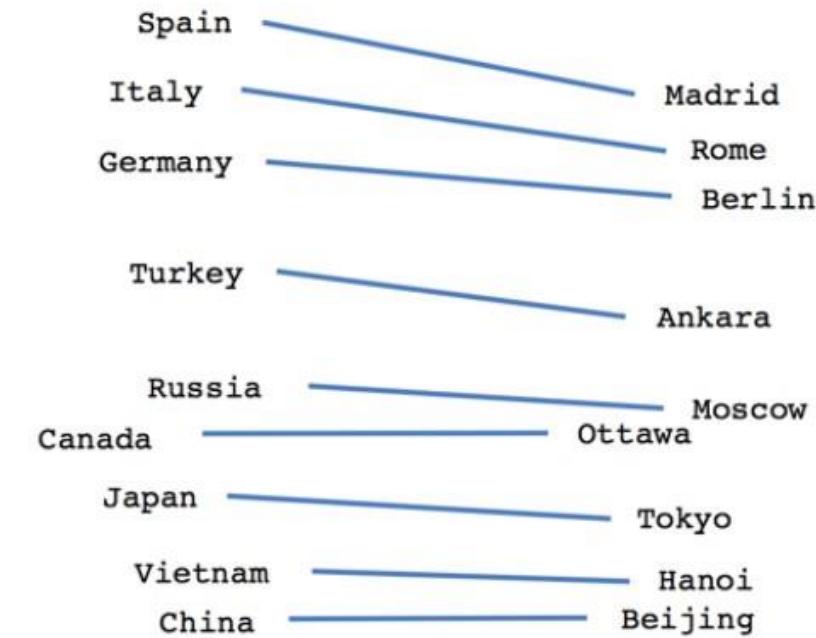
Word Embeddings



Male-Female

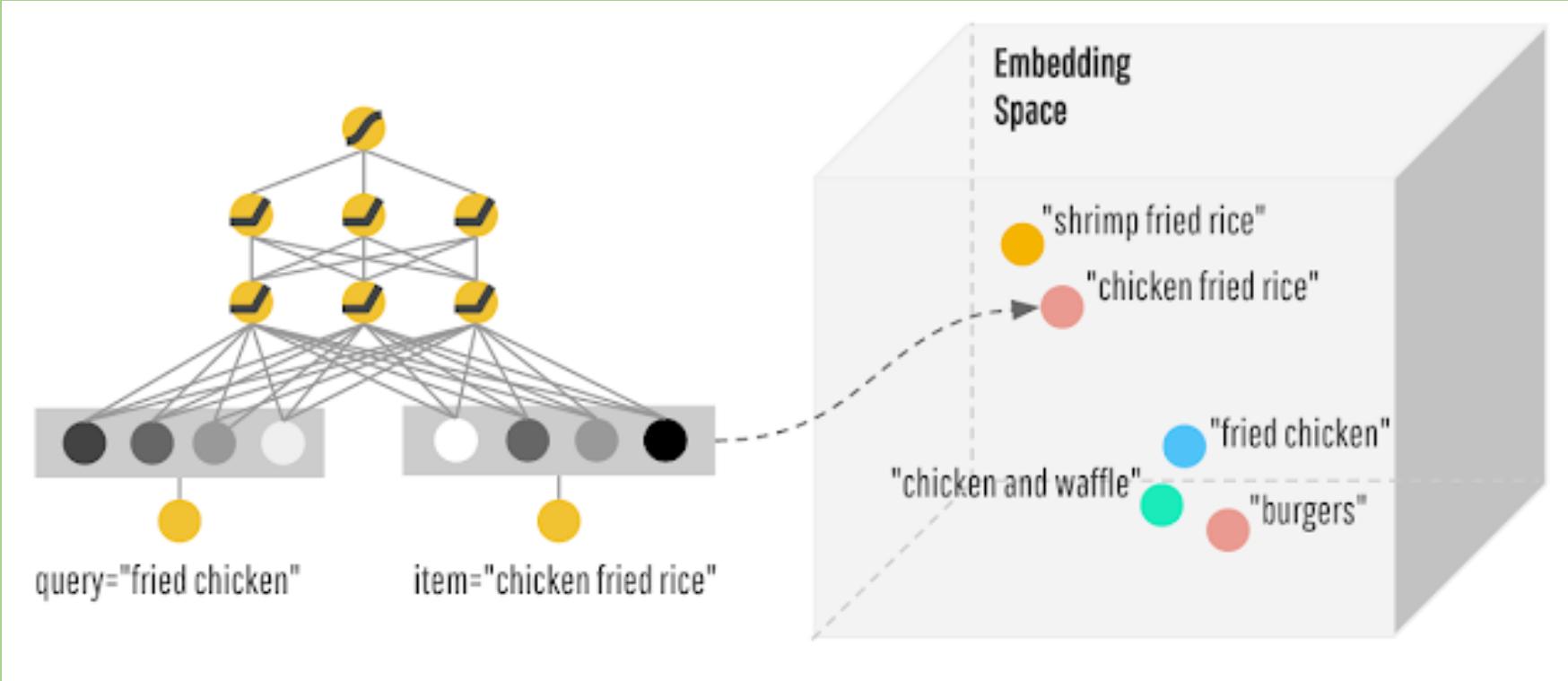


Verb tense

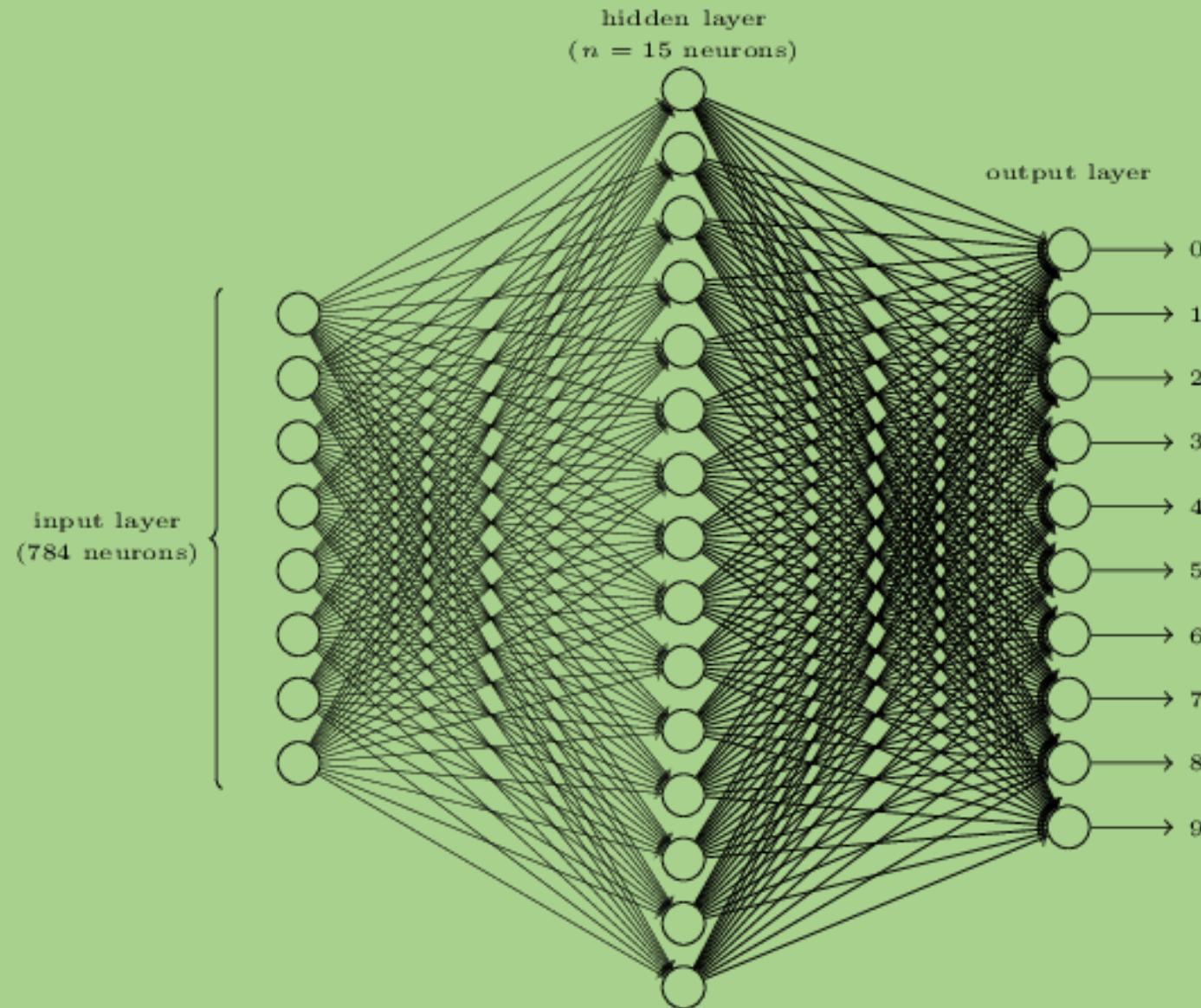


Country-Capital

Word Embeddings



Let's build a ML Model



Training testing and validation

1. Training data

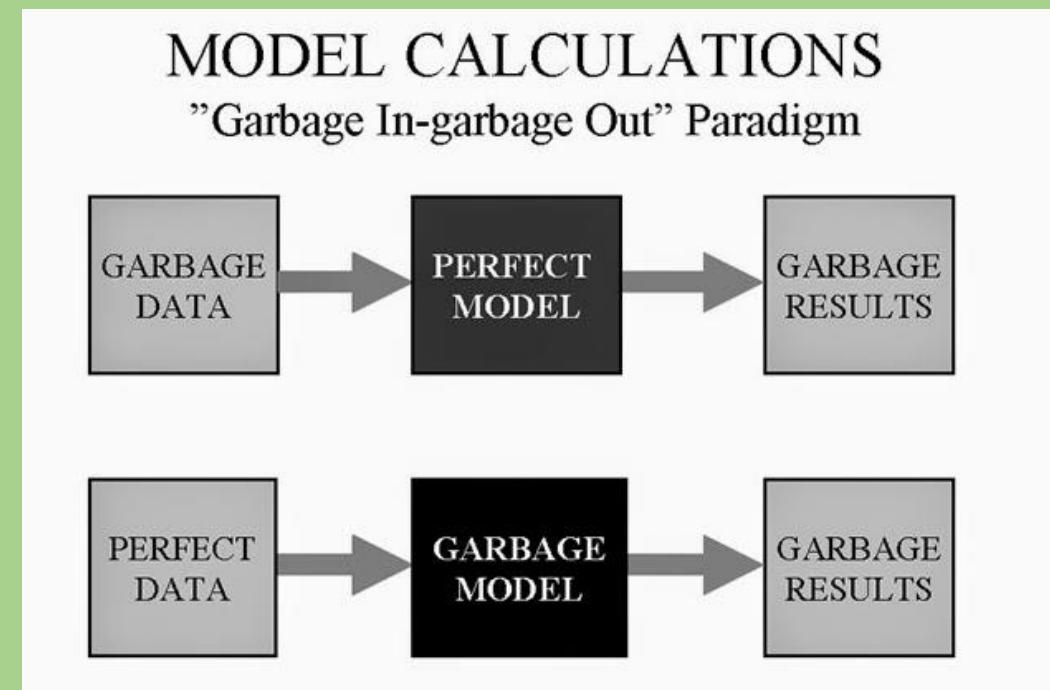
The data your algorithm will 'learn' from.
garbage in = garbage out.

2. Validation data

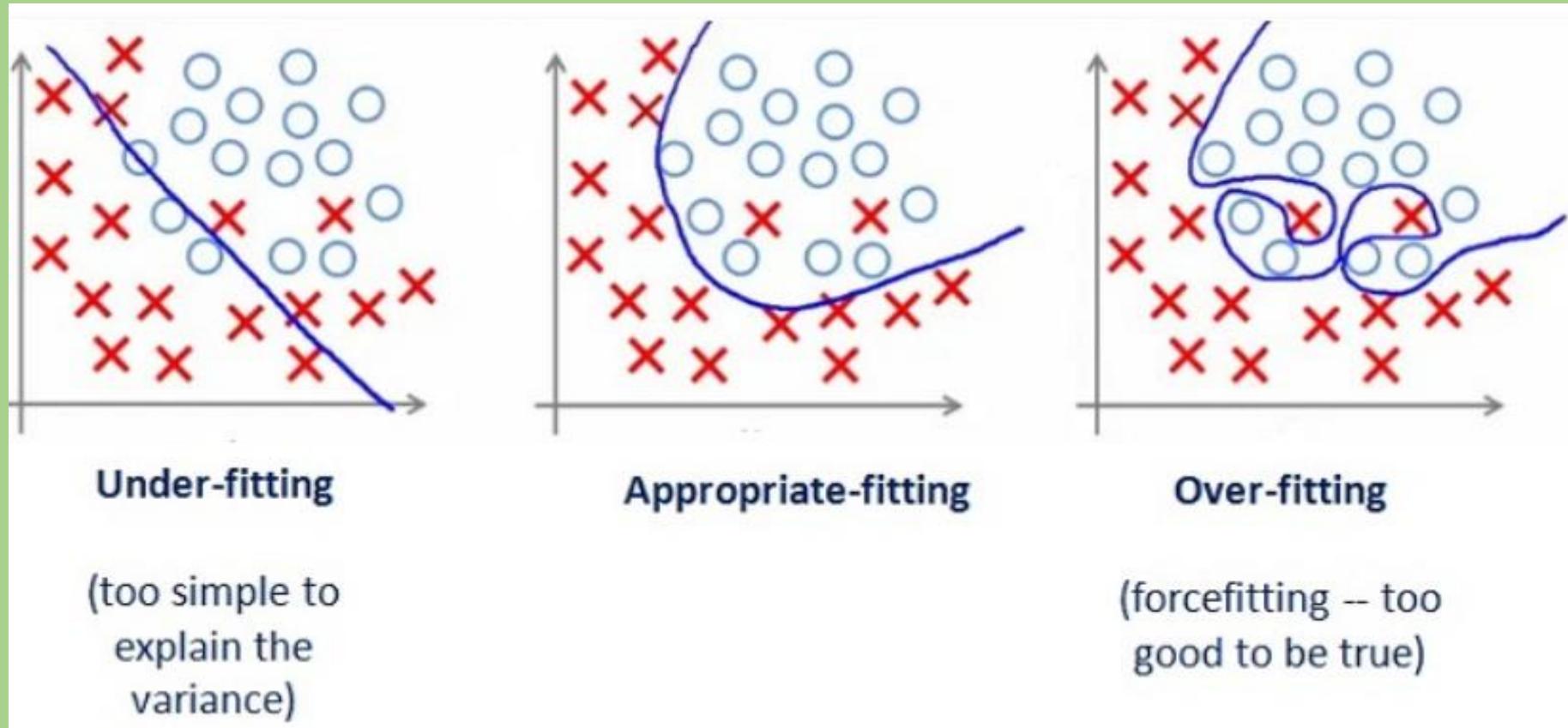
Data used to tune hyperparameters.
Results will be used to tweak the model.

3. Testing data

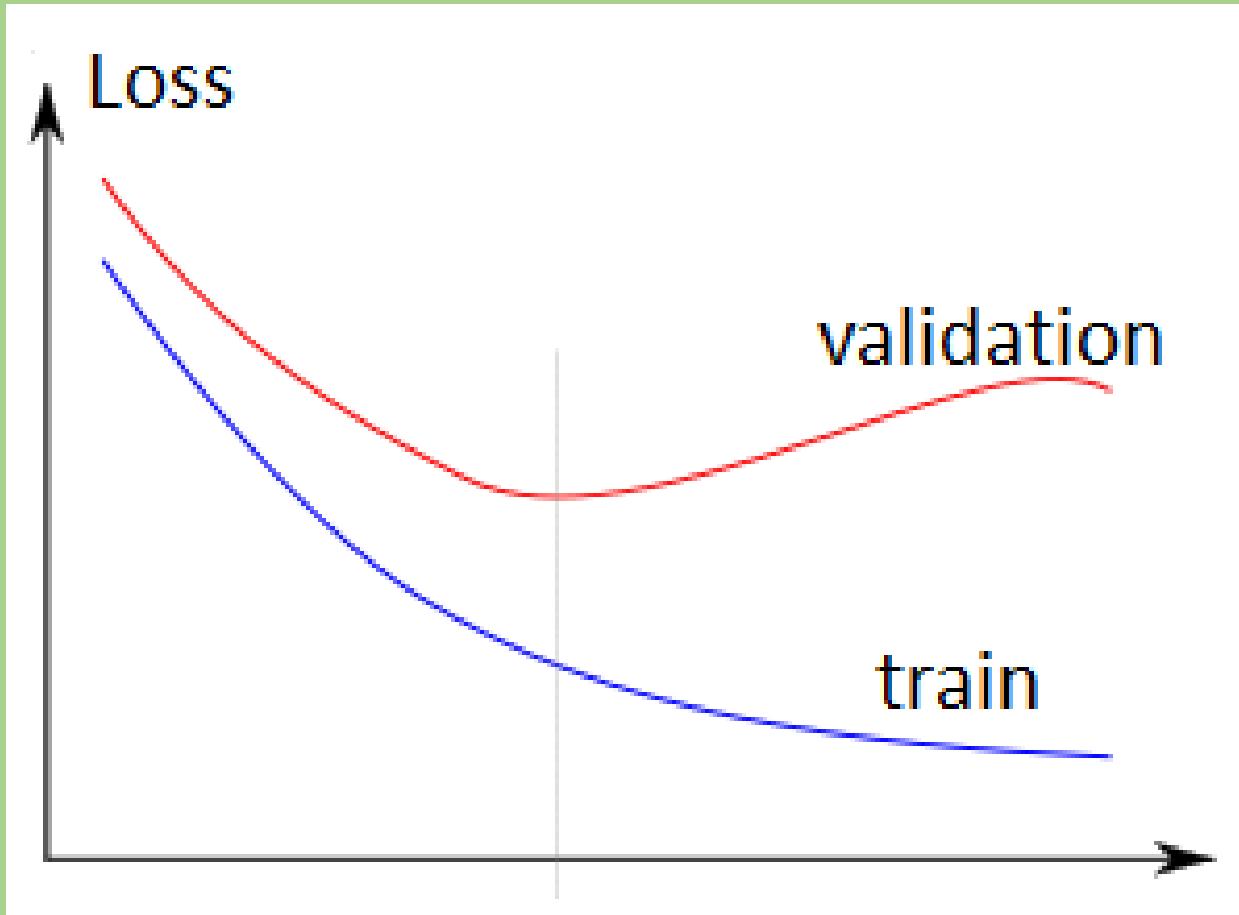
Used to validate final accuracy.
Should not be used for hyperparameter tuning.



Overfitting/Underfitting



Overfitting



Resources

- Youtube
 - <http://3b1b.co/neural-networks>
 - https://www.youtube.com/playlist?list=PLZbbT5o_s2xq7Lwl2y8_QtvuXZedL6tQU
 - https://www.youtube.com/watch?v=JB8T_zN7ZC0
 - ... + a LOT more
- MOOCs
 - fast.ai
 - deeplearning.ai
 - <http://cs231n.stanford.edu>