**LOUISIANA STATE UNIVERSITY**
College of Agriculture
School of Plant, Environmental, and Soil Sciences
HTP in Plant Breeding

# Data mining and modeling
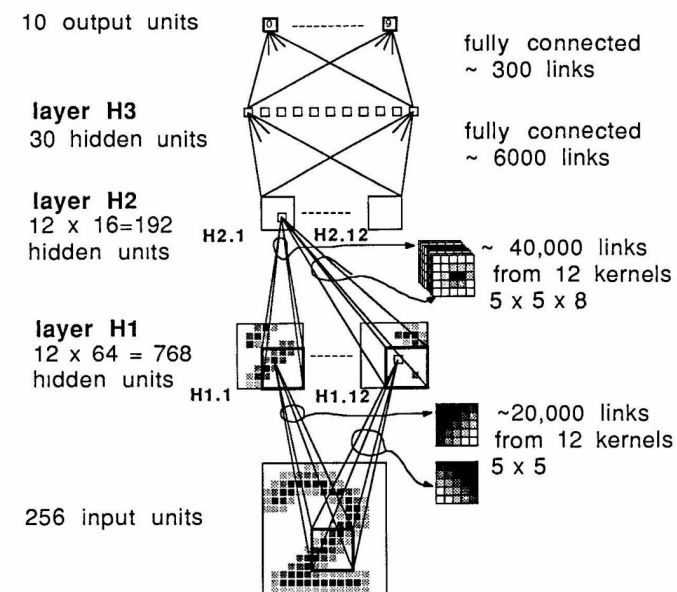
**Prof. Roberto Fritsche-Neto**

**rfneto@agcenter.lsu.edu**
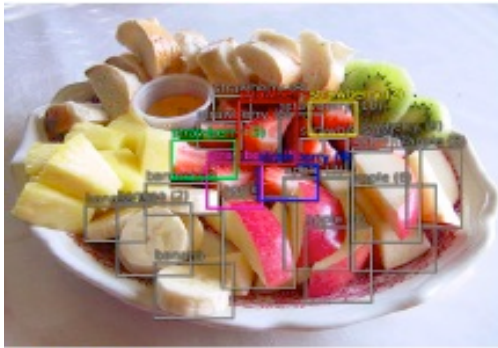
**Baton Rouge, March 20th, 2024**

# LeNet 1989

- **Recognize zip codes in US postal service**



- **Around 5% error rate. Good enough to be useful**

# CNN: Detection



**Groundtruth:**
strawberry
strawberry (2)
strawberry (3)
strawberry (4)
strawberry (5)
strawberry (6)
strawberry (7)
strawberry (8)
strawberry (9)
strawberry (10)
apple
apple (2)
apple (3)

**Groundtruth:**
tv or monitor
tv or monitor (2)
tv or monitor (3)
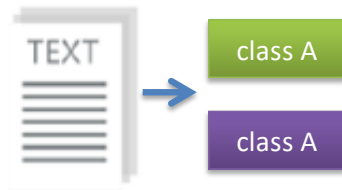person
remote control
remote control (2)

**Sermanet, CVPR 2014**

# Types of Learning
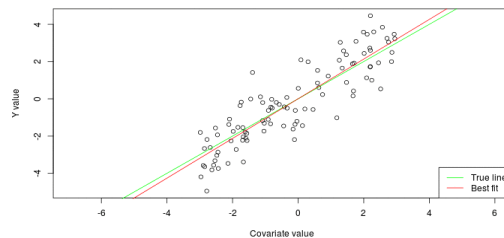
**Supervised:** Learning with a labeled training set
*Example*: seed *classification* based on colors and shape

**Unsupervised:** Discover patterns in unlabeled data
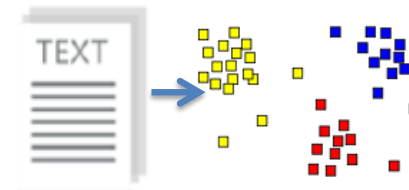*Example*: *cluster* images based on color components



Classification
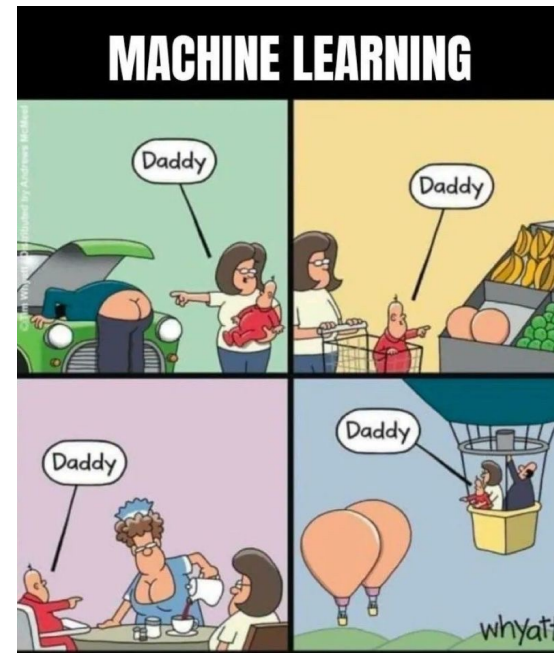
Regression

Clustering

# The Curse of Learning





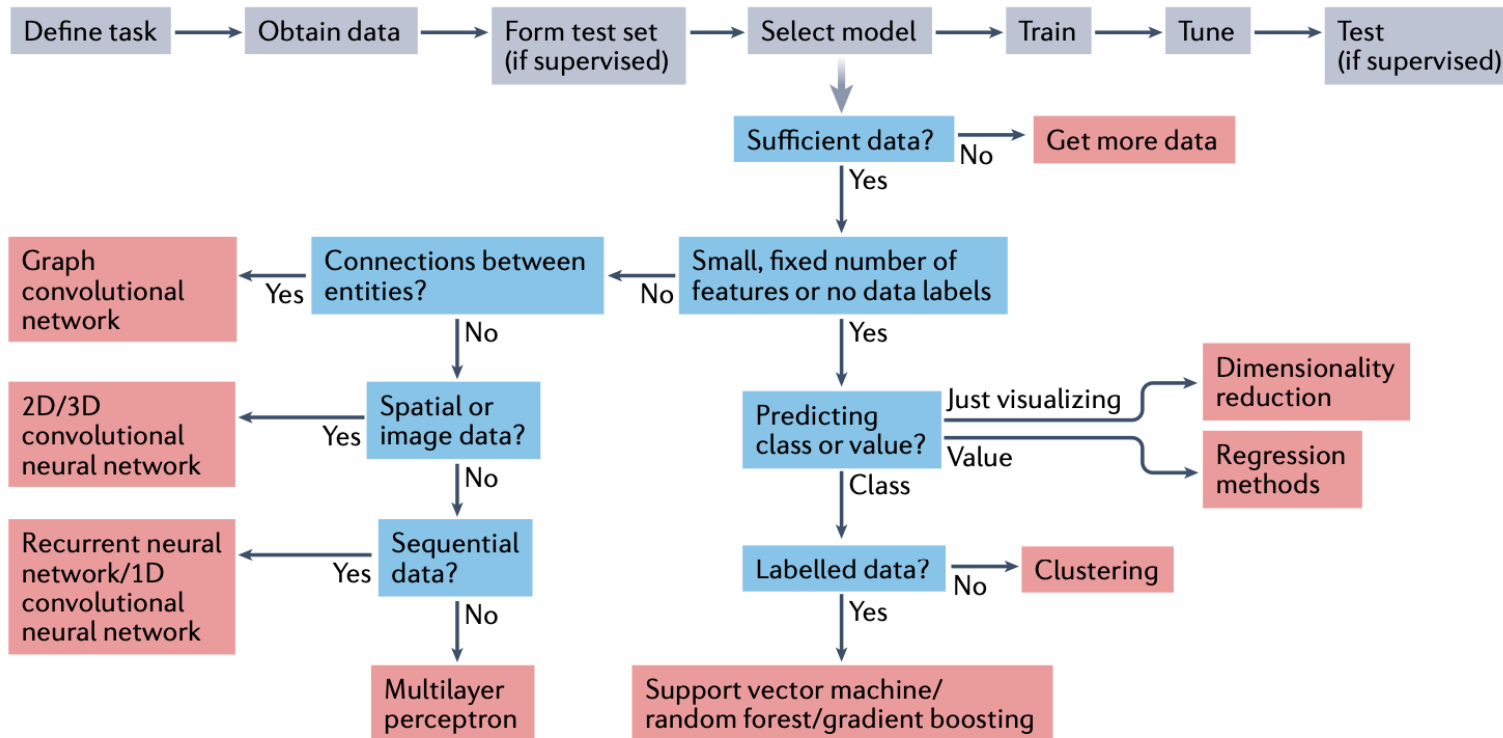**A guide to machine learning for biologists**

Joe G. Greener, Shaun M. Kandathil, Lewis Moffat & David T. Jones ✉

*Nature Reviews Molecular Cell Biology* **23**, 40–55 (2022) | Cite this article

https://www.nature.com/articles/s41580-021-00407-0

# The Way of Learning



Define task → Obtain data → Form test set (if supervised) → Select model → Train → Tune → Test (if supervised)

Sufficient data? — No → Get more data

Yes

Small, fixed number of features or no data labels

Connections between entities? — Yes → Graph convolutional network

No

Spatial or image data? — Yes → 2D/3D convolutional neural network

No

Sequential data? — Yes → Recurrent neural network/1D convolutional neural network

No

Multilayer perceptron

Predicting class or value? — Just visualizing → Dimensionality reduction

Value → Regression methods

Class

Labelled data? — No → Clustering

Yes

Support vector machine/random forest/gradient boosting

# One size does not fit all



**a Regression**

Model
Data point
Weight
Features
Height

**b SVM**

Ordered protein
Margin
Separating hyperplane
Disordered protein

**c Gradient boosting**

Drug property 2
Drug property 1
+ Active
− Inactive

**d PCA**

PC1
PC2
Weight
Height

**e Clustering**

Gene 2 expression
Cell types
Gene 1 expression

# Data structure for supervised learning

- **Classification**

| Var1 | Var2 | Var3 | Var4 | Var5 | Var6 | Var7 | Var8 | Var9 | Var10 | Var-target |
|------|------|------|------|------|------|------|------|------|-------|------------|
| 1.2  | 2.3  | 4.21 | 22   | 1.5  | 0.2  | 0.77 | 2.9  | 3.4  | 0.8   | A          |
| 1.5  | 2.99 | 5.21 | 23   | 41.5 | 1.2  | 1.77 | 2.6  | 1.4  | 3.8   | A          |
| 1.9  | 4.3  | 4.44 | 11   | 8.5  | 0.1  | 9.77 | 6.9  | 0.4  | 6.8   | B          |

- **Regression**

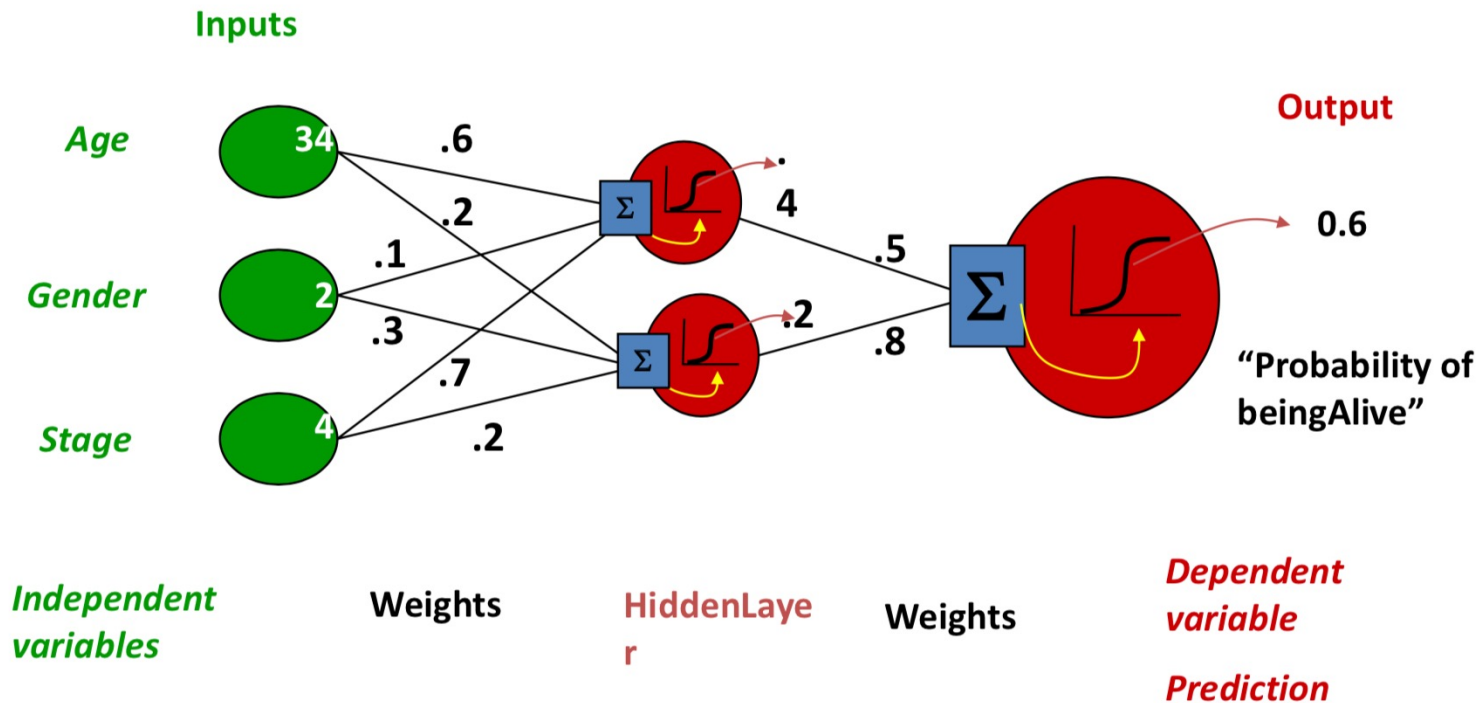| Var1 | Var2 | Var3 | Var4 | Var5 | Var6 | Var7 | Var8 | Var9 | Var10 | Var-target |
|------|------|------|------|------|------|------|------|------|-------|------------|
| 0.2  | 2.3  | 4.23 | 22   | 1.5  | 0.2  | 0.76 | 2.9  | 3.4  | 0.8   | 3.12       |
| 1.5  | 2.99 | 5.21 | 23   | 41.5 | 1.2  | 1.77 | 2.6  | 1.4  | 3.8   | 4.88       |
| 1.9  | 4.3  | 4.44 | 11   | 8.5  | 0.1  | 9.79 | 6.9  | 0.4  | 6.8   | 9.73       |



- **Scale / normalize/standardize the data**
- **Balanced classes**
- **Imputation**
- **Create dummy variables**
- **Replace variables that don't bring any information – zip code -> lat and long**
- **Feature selection - *Reduce dimensionality to help to avoid overfitting***
- *ML models do not work properly with collinearity*

# Artificial Neural Networks (ANN)

**Many layers of staked multiple regressions**
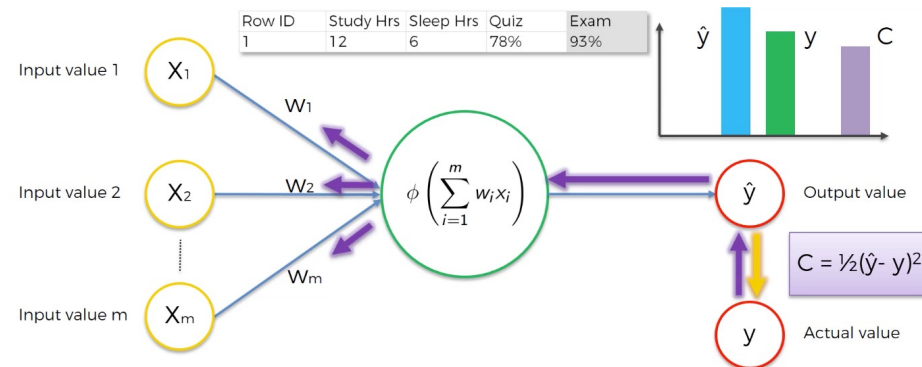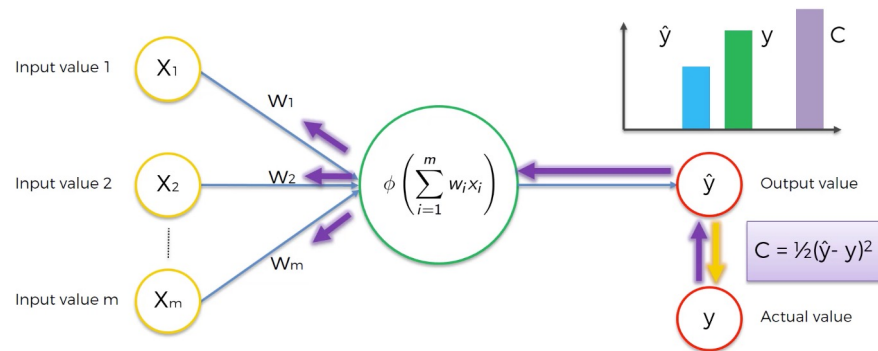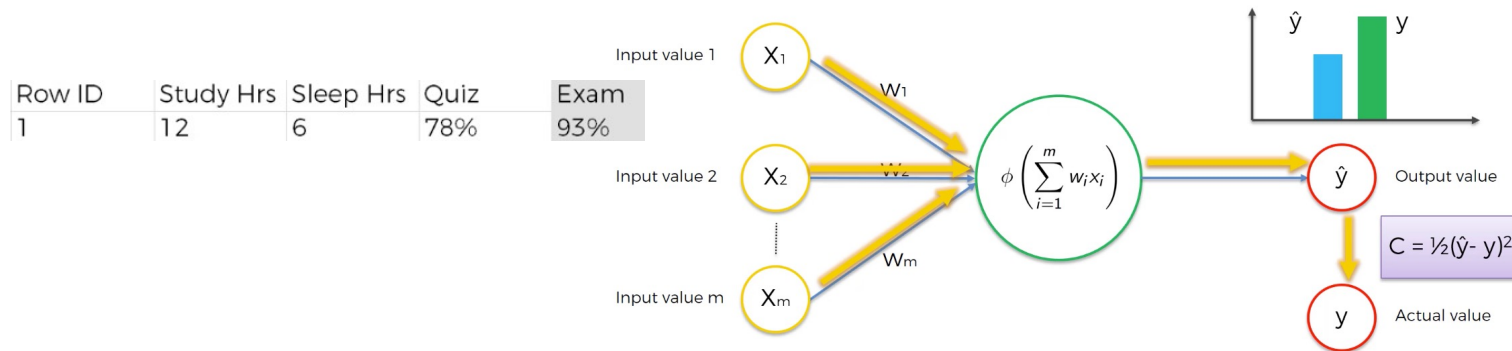**The regression type depends on the activation function**
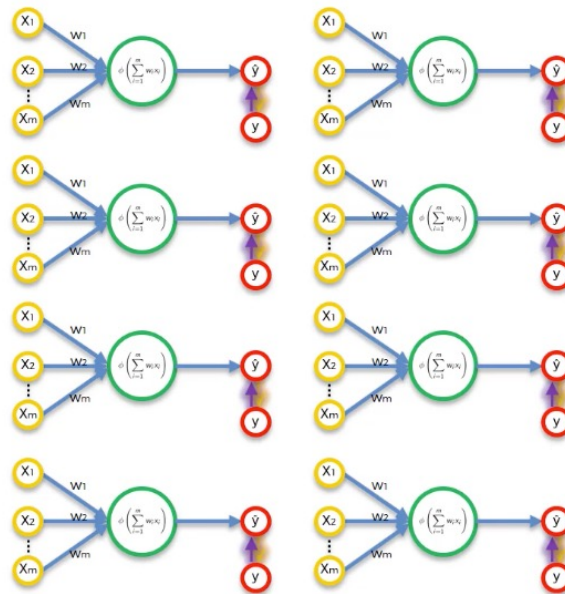
# Algorithm for learning ANN

- **Initialize the weights ($w_0, w_1, \ldots, w_k$)**

- **Adjust the weights in such a way that the output of ANN is consistent with class labels of training examples**

    **Error function:**

$$E = \sum_i \left[ Y_i - f(w_i, X_i) \right]^2$$

- **Find the weights $w_i$'s that minimize the error function**

    **e.g.,** gradient descent, backpropagation algorithm

**Top diagram:**

| Row ID | Study Hrs | Sleep Hrs | Quiz | Exam |
|--------|-----------|-----------|------|------|
| 1 | 12 | 6 | 78% | 93% |

Input value 1 — $X_1$ — $W_1$

Input value 2 — $X_2$ — $W_2$

Input value m — $X_m$ — $W_m$

$\phi\left(\sum_{i=1}^{m} w_i x_i\right)$

$\hat{y}$ — Output value

$\hat{y}$   $y$

$C = \frac{1}{2}(\hat{y} - y)^2$

$y$ — Actual value

**Middle diagram:**

Input value 1 — $X_1$ — $W_1$

Input value 2 — $X_2$ — $W_2$

Input value m — $X_m$ — $W_m$

$\phi\left(\sum_{i=1}^{m} w_i x_i\right)$

$\hat{y}$ — Output value

$\hat{y}$   $y$   $C$

$C = \frac{1}{2}(\hat{y} - y)^2$

$y$ — Actual value

**Bottom diagram:**

| Row ID | Study Hrs | Sleep Hrs | Quiz | Exam |
|--------|-----------|-----------|------|------|
| 1 | 12 | 6 | 78% | 93% |

Input value 1 — $X_1$ — $W_1$

Input value 2 — $X_2$ — $W_2$

Input value m — $X_m$ — $W_m$

$\phi\left(\sum_{i=1}^{m} w_i x_i\right)$

$\hat{y}$ — Output value

$\hat{y}$   $y$   $C$

$C = \frac{1}{2}(\hat{y} - y)^2$

$y$ — Actual value

| Row ID | Study Hrs | Sleep Hrs | Quiz | Exam |
|--------|-----------|-----------|------|------|
| 1 | 12 | 6 | 78% | 93% |
| 2 | 22 | 6.5 | 24% | 68% |
| 3 | 115 | 4 | 100% | 95% |
| 4 | 31 | 9 | 67% | 75% |
| 5 | 0 | 10 | 58% | 51% |
| 6 | 5 | 8 | 78% | 60% |
| 7 | 92 | 6 | 82% | 89% |
| 8 | 57 | 8 | 91% | 97% |

$$C = \sum \tfrac{1}{2}(\hat{y} - y)^2$$

Adjust $w_1$, $w_2$, $w_3$

# Algorithm for learning ANN

**STEP 1:** Randomly initialise the weights to small numbers close to 0 (but not 0).

**STEP 2:** Input the first observation of your dataset in the input layer, each feature in one input node.

**STEP 3:** Forward-Propagation: from left to right, the neurons are activated in a way that the impact of each neuron's activation is limited by the weights. Propagate the activations until getting the predicted result y.

**STEP 4:** Compare the predicted result to the actual result. Measure the generated error.

**STEP 5:** Back-Propagation: from right to left, the error is back-propagated. Update the weights according to how much they are responsible for the error. The learning rate decides by how much we update the weights.

**STEP 6:** Repeat Steps 1 to 5 and update the weights after each observation (Reinforcement Learning). Or: Repeat Steps 1 to 5 but update the weights only after a batch of observations (Batch Learning).
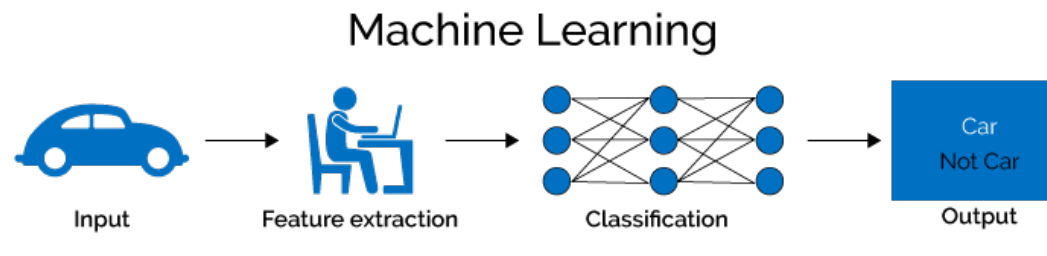
**STEP 7:** When the whole training set passed through the ANN, that makes an epoch. Redo more epochs.

Batch: sample (fold) used for training
Epoch: iteration

# ML *vs*. Deep Learning

- **DL is a machine learning subfield of learning representations of data exceptionally effective at learning patterns**
- **Deep learning algorithms attempt to learn (multiple levels of) representation by using a hierarchy of multiple layers**
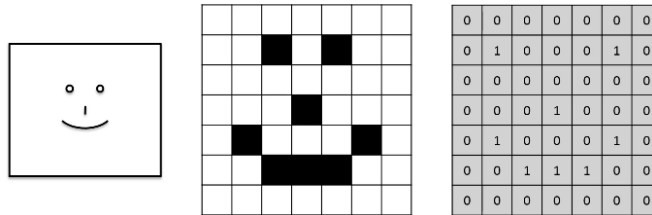- **If you provide the system tons of information, it begins to understand it and respond in useful ways**

**STEP 1:** Convolution

**STEP 2:** Max Pooling

**STEP 3:** Flattening

**STEP 4:** Full Connection

Input Image    Feature Detector    Feature Map

Input Image    Feature Detector    Feature Map

Input Image    Feature Detector    Feature Map

We create many feature maps to obtain our first convolution layer

Feature Maps

Input Image

Convolutional Layer

**STEP 1:** Convolution

**STEP 2:** Max Pooling

**STEP 3:** Flattening

**STEP 4:** Full Connection

Feature Map

Max Pooling

Pooled Feature Map

Max Pooling

Feature Map

Pooled Feature Map

Max Pooling

Feature Map

Pooled Feature Map

Max Pooling

Feature Map

Pooled Feature Map

Input Image

Convolution

Convolutional Layer

Pooling

Pooling Layer

Udemy

**STEP 1:** Convolution

**STEP 2:** Max Pooling

**STEP 3:** Flattening

**STEP 4:** Full Connection



| 1 | 1 | 0 |
|---|---|---|
| 4 | 2 | 1 |
| 0 | 2 | 1 |

Pooled Feature Map

Flattening

| 1 |
|---|
| 1 |
| 0 |
| 4 |
| 2 |
| 1 |
| 0 |
| 2 |
| 1 |

Udemy



Input Image

Convolution

Convolutional Layer

Pooling

Pooling Layer

Flattening

Input layer of a future ANN

Dog

Cat

- AlexNet
- VGG-16
- VGG-19
- GoogLeNet
- ResNet-18
- ResNet-50
- ResNet-101



224 × 224 × 3   224 × 224 × 64

112 × 112 × 128

56 × 56 × 256

28 × 28 × 512

14 × 14 × 512

7 × 7 × 512

1 × 1 × 4096   1 × 1 × 1000

convolution+ReLU
max pooling
fully connected+ReLU
softmax

227 x 227 pixels

**AlexNet**

224 x 224 pixels

**RNCs**

# Learnable parameters



Weights

$$h = \sigma(W_1 x + b_1)$$
$$y = \sigma(W_2 h + b_2)$$

Activation functions

How do we train?

4 + 2 = 6 neurons (not counting inputs)
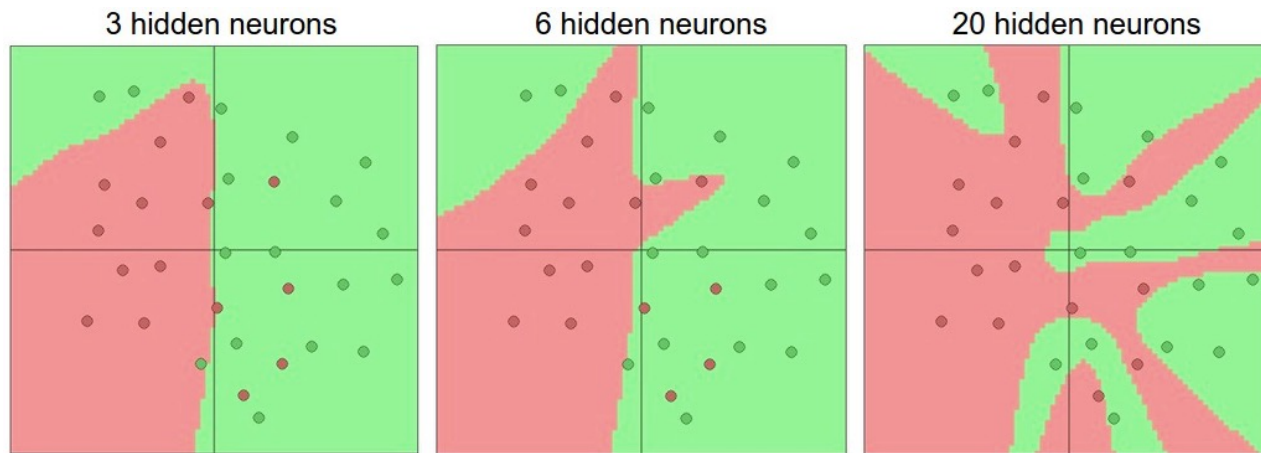[3 x 4] + [4 x 2] = 20 weights
4 + 2 = 6 biases

26 learnable **parameters**
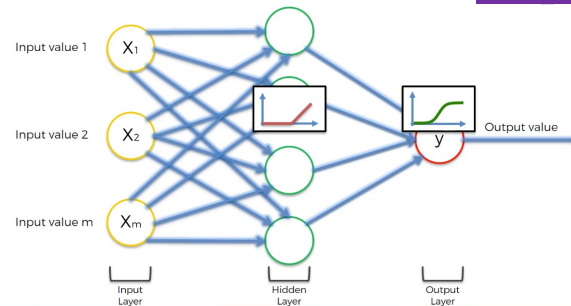
# Number of neurons and layers

- **The number of <span style="color:red">neurons</span> in the <span style="color:red">input layer</span> is the number of <span style="color:red">explanatory variables</span>**
- **Non-linearities needed to learn complex (<span style="color:red">non-linear</span>) representations of data**
- **Otherwise, the ANN would be just a linear function**



3 hidden neurons     6 hidden neurons     20 hidden neurons

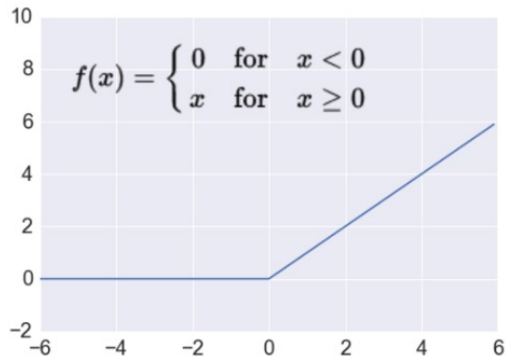- **More <span style="color:red">hidden layers and neurons</span> can approximate more complex functions**

# Activation functions

- **The learned information extracted from the training data is stored and captured by the weight values of the connections between the layers**
- **The final output can be continuous, binary, ordinal, or count, which is controlled for the activation function in the output layer**



| Activation functions | g(z) | String | Use | Data |
|---|---|---|---|---|
| rectifier linear unit | Max(0, Z) | relu | Positive values | CONTINUOS BUT POSITIVE |
| sigmoid | $(1+e^{-z})^{-1}$ | sigmoid | converts independent variables of near infinite range into simple probabilities between 0 and 1 | BINARY |
| Tanh | tanh(z) | tanh | deal more easily with negative numbers | CONTINUOS FROM NEGATIVE TO POSITIVE |
| softmax | $g(z_j) = \frac{\exp(z_j)}{1+\sum_{c=1}^{C} \exp(z_c)}, j=1,..,C$ | softmax | sigmoid activation function that handles multinomial labeling systems, that is, it is appropriate for categorical outcomes | MORE THAN TWO CLASSES |

# Activation: ReLU

$$f(x) = \begin{cases} 0 & \text{for} \quad x < 0 \\ x & \text{for} \quad x \geq 0 \end{cases}$$

*http://adilmoujahid.com/images/activation.png*

Takes a real-valued number and
thresholds it at zero   $f(x) = \max(0, x)$

$$R^n \rightarrow R_+^n$$

**Trains much faster**
- **accelerates the convergence**

**Less expensive operations**
- **compared to sigmoid/tanh (exponentials etc.)**
- **implemented by simply thresholding a matrix at zero**

# Loss or cost functions and output

**Measure the amount of 'disagreement' between the obtained and ideal outputs**

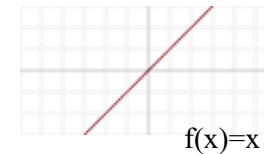|  | **Classification** | **Regression** |
|---|---|---|
| **Training examples** | $R^n$ x {class_1, ..., class_n} (one-hot encoding) | $R^n$ x $R^m$ |
| **Output Layer** | Softmax [map $R^n$ to a probability distribution] $P(y = j \mid \mathbf{x}) = \dfrac{e^{\mathbf{x}^\mathsf{T}\mathbf{w}_j}}{\sum_{k=1}^{K} e^{\mathbf{x}^\mathsf{T}\mathbf{w}_k}}$ | Tanh, linear or Sigmoid f(x)=x |
| **Cost (loss) function** | Cross-entropy | Mean Squared Error |

Classification cost function:

$$J(\theta) = -\frac{1}{n}\sum_{i=1}^{n}\sum_{k=1}^{K} \left[ y_k^{(i)} \log \hat{y}_k^{(i)} + \left(1 - y_k^{(i)}\right) \log\left(1 - \hat{y}_k^{(i)}\right)\right]$$

Regression cost functions:

Mean Squared Error

$$J(\theta) = \frac{1}{n}\sum_{i=1}^{n}\left(y^{(i)} - \hat{y}^{(i)}\right)^2$$
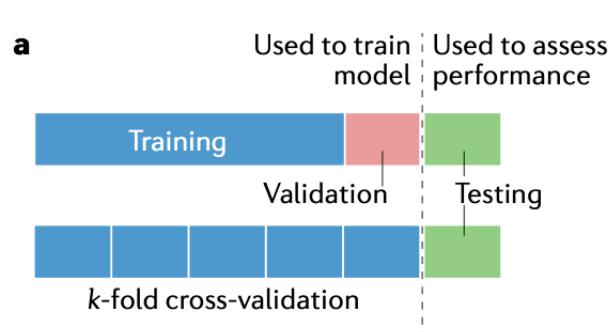
Mean Absolute Error

$$J(\theta) = \frac{1}{n}\sum_{i=1}^{n}\left|y^{(i)} - \hat{y}^{(i)}\right|$$

**Most them weigh more the values near the average, but in breeding, our concern is the extreme values**

# Cross-validation and model selection

- **Randomly select two portions of the data to be used for <span style="color:red">training</span> and <span style="color:red">validation</span>**
- **Then, <span style="color:red">run the model on the test set</span> to see how it performs**



**Confusion Matrix**

Classified As

|  | Blue | Red |
|---|---|---|
| **Blue** | 7 | 1 |
| **Red** | 0 | 5 |

$Accuracy = \dfrac{TP + TN}{TP + FP + TN + FN}$
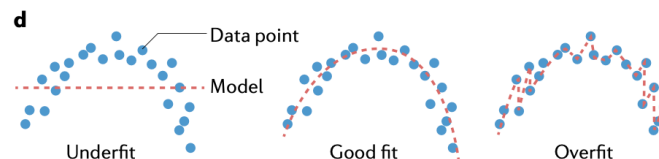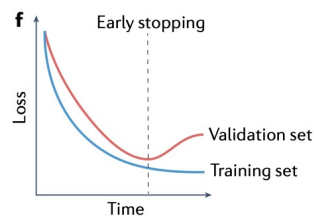
$Sensitivity = \dfrac{TP}{TP + FN}$

$Specificity = \dfrac{TN}{TN + FP}$

$FPR = \dfrac{FP}{FP + TN}$

$F1\ score = 2\ x\ \dfrac{sensitivity\ x\ precision}{sensitivity + precision}$

$AUC = \dfrac{1}{2}\left(\dfrac{TP}{TP + FN} + \dfrac{TN}{TN + FP}\right)$

- **Learn the training data very well, even outliers (<span style="color:red">noise</span>)**
- **Fail to <span style="color:red">generalize</span> to new examples (<span style="color:red">test data</span>)**

# Regularization



**Dropout**
- **Randomly drop units (along with their connections) during training**
- **Each unit retained with fixed probability p, independent of other units**
- **Hyper-parameter p to be chosen (tuned)**

**L2 = weight decay**
- **Regularization term that penalizes big weights, added to the objective**
- **Weight decay value determines how dominant regularization is during gradient computation**
- **Big weight decay coefficient → big penalty for big weights**

$$J_{reg}(\theta) = J(\theta) + \lambda \sum_k \theta_k^2$$

**Early-stopping**
- **Use validation error to decide when to stop training**
- **Stop when monitored quantity has not improved after n subsequent epochs**
- **n is called patience**