The background features a series of concentric circles in light gray, some solid and some dashed, creating a ripple effect. A large, solid green oval is positioned in the center, containing the main text. A dark gray, curved shape is visible on the left side, partially overlapping the green oval.

How do I teach my second  
grade kid what AI is?

Cupid Chan



# Disclaimer

- This presentation session was prepared or accomplished by Cupid Chan in his personal capacity. The opinions expressed in this talk are the author's own and do not reflect the view of any United States Government Agency or Department.



# Cupid Chan



- CTO of Index Analytics
- Board of Directors and TSC, Linux Foundation ODPI
- Organizer of a Big Data In Action Meetup in Washington/Baltimore area with 1600 members

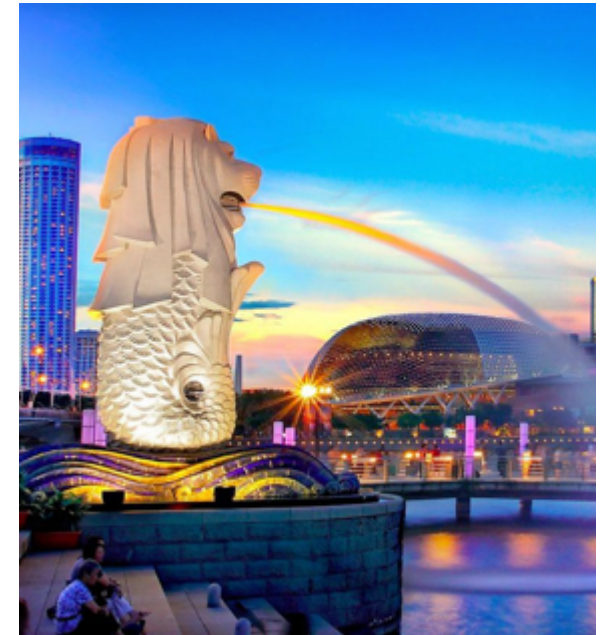


[www.linkedin.com/in/cupidchan/](http://www.linkedin.com/in/cupidchan/)



[@cupidckchan](https://twitter.com/cupidckchan)





Let's play a game

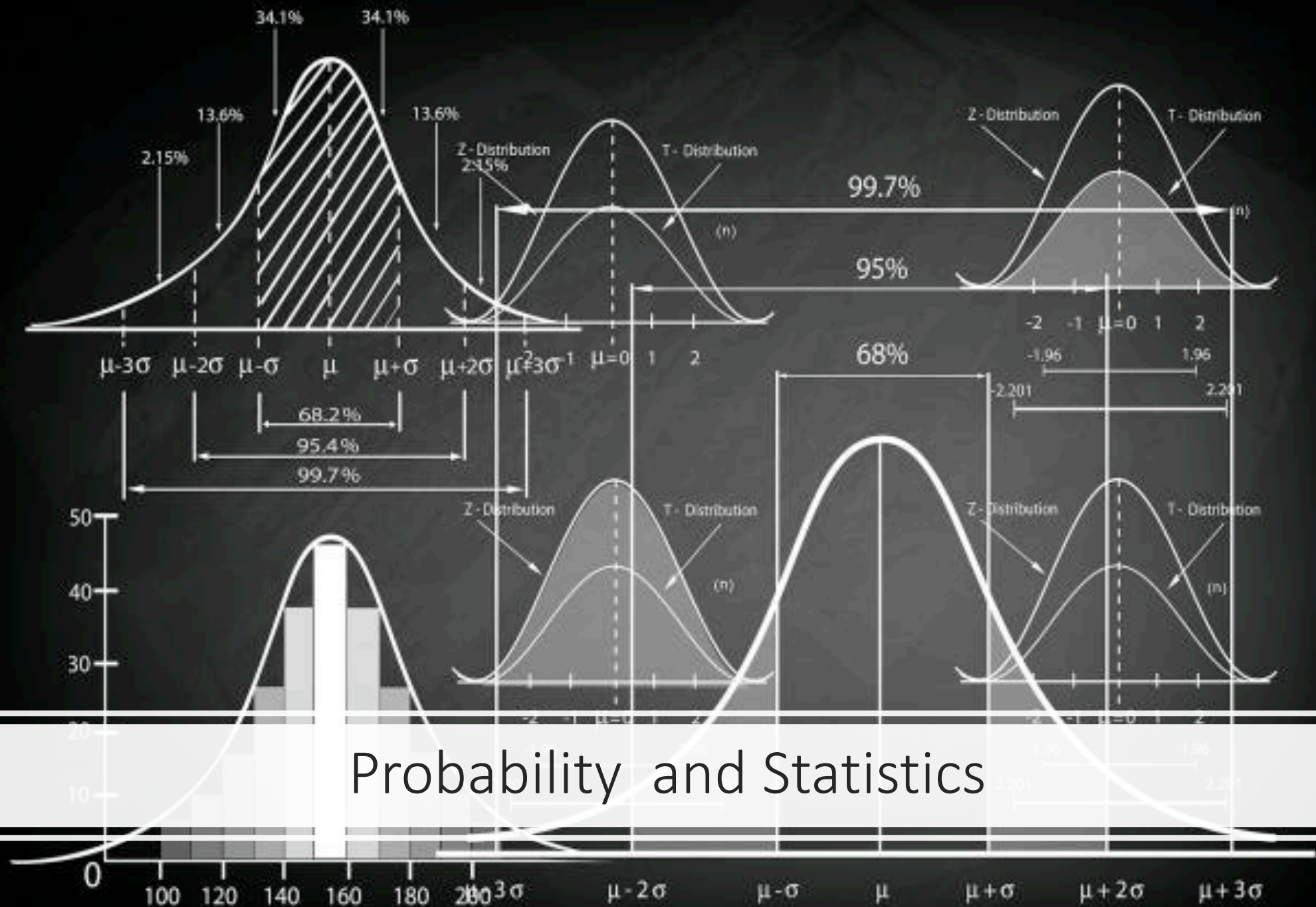
---

Probability to guess the 3 cards correctly

1

---

132,600



Probability and Statistics



6)  $x_1, x_2 \in \mathbb{R}$ ,  $y_1, y_2 \in \mathbb{R}$ , for all scalars  $a \in \mathbb{R}$ , and all vectors  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^n$ .  
 7) For any  $x, y \in \mathbb{R}^n$ ,  $\langle x, y \rangle = \langle y, x \rangle$  (for scalar 1).  
 8) For any two scalars  $a, b \in \mathbb{R}$  and any  $x \in \mathbb{R}^n$ ,  $\langle ax + by, x \rangle = a\langle x, x \rangle + b\langle y, x \rangle$ .  
 9)  $\langle ax + by, x \rangle = a\langle x, x \rangle + b\langle y, x \rangle$ . 10)  $\langle x, x \rangle \geq 0$  and  $\langle x, x \rangle = 0$  if and only if  $x = 0$ .

**Linear Independence:** Consider  $n$  vectors  $\{x_1, x_2, \dots, x_n\}$ . If there exists  $n$  scalars  $a_1, a_2, \dots, a_n$ , at least one of which is nonzero, such that  $a_1x_1 + a_2x_2 + \dots + a_nx_n = 0$ , then the  $\{x_i\}$  are linearly dependent.

#### Spanning a Space:

Let  $L$  be a linear vector space and let  $\{x_1, x_2, \dots, x_n\}$  be a subset of vectors in  $L$ . This subset spans  $L$  if and only if for every vector  $x \in L$  there exist scalars  $a_1, a_2, \dots, a_n$  such that  $x = a_1x_1 + a_2x_2 + \dots + a_nx_n$ .

**Inner Product:**  $\langle x, y \rangle$  for any scalar function of  $x$  and  $y$ .

1.  $\langle x, y \rangle = \langle y, x \rangle$  2.  $\langle ax + by, z \rangle = a\langle x, z \rangle + b\langle y, z \rangle$

3.  $\langle x, x \rangle \geq 0$ , where equality holds iff  $x$  is the zero vector.

**Norm:** A scalar function  $\|x\|$  is called a norm if it satisfies:

1.  $\|x\| \geq 0$  2.  $\|x\| = 0$  if and only if  $x = 0$ .

3.  $\|ax\| = |a|\|x\|$  4.  $\|x + y\| \leq \|x\| + \|y\|$

**Angle:** The angle  $\theta$  bet. 2 vectors  $x$  and  $y$  is defined by  $\cos \theta = \frac{\langle x, y \rangle}{\|x\|\|y\|}$

**Orthogonality:** 2 vectors  $x, y \in \mathbb{R}^n$  are said to be orthogonal if  $\langle x, y \rangle = 0$ .

#### Gram Schmidt Orthogonalization:

Assume that we have  $n$  independent vectors  $y_1, y_2, \dots, y_n$ . From these vectors we will obtain  $n$  orthogonal vectors  $v_1, v_2, \dots, v_n$ .

$$v_1 = y_1 \quad v_k = y_k - \sum_{i=1}^{k-1} \frac{\langle y_k, v_i \rangle}{\langle v_i, v_i \rangle} v_i$$

where  $\frac{\langle y_k, v_i \rangle}{\langle v_i, v_i \rangle} v_i$  is the projection of  $y_k$  on  $v_i$

#### Vector Expansions:

$$x = \sum_{i=1}^n x_i v_i = x_1 v_1 + x_2 v_2 + \dots + x_n v_n$$

for orthogonal vectors,  $x_i = \frac{\langle x, v_i \rangle}{\langle v_i, v_i \rangle}$

#### Reciprocal Basis Vectors:

$$(v_i, v_j) = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \quad x_j = \langle v_j, x \rangle$$

To compute the reciprocal basis vectors: set  $B = [v_1 \ v_2 \ \dots \ v_n]$

$B^{-1} = [v_1^* \ v_2^* \ \dots \ v_n^*]$  In matrix form:  $x^* = B^{-1} x^T$

#### Transformations:

A transformation consists of three parts:

domain:  $X = \{x_i\}$ , range:  $Y = \{y_i\}$ , and a rule relating each  $x_i \in X$  to an element  $y_i \in Y$ .

**Linear Transformations:** transformation  $A$  is linear if:

1. for all  $x_1, x_2 \in X$ ,  $A(x_1 + x_2) = A(x_1) + A(x_2)$

2. for all  $x \in X$ ,  $a \in \mathbb{R}$ ,  $A(ax) = aA(x)$

#### Matrix Representations:

Let  $\{v_1, v_2, \dots, v_n\}$  be a basis for vector space  $X$  and let  $\{u_1, u_2, \dots, u_m\}$  be a basis for vector space  $Y$ . Let  $A$  be a linear transformation with domain  $X$  and range  $Y$ .  $A(x) = y$

The coefficients of the matrix representation are obtained from

$$A(v_i) = \sum_{j=1}^m a_{ji} u_j$$

**Change of Basis:**  $B_1 = [t_1 \ t_2 \ \dots \ t_n]$ ,  $B_2 = [w_1 \ w_2 \ \dots \ w_n]$

$$A' = [B_2^{-1} A B_1]$$

**Eigenvalues & Eigenvectors:**  $Ax = \lambda x$ ,  $[\lambda I - A] = 0$

**Diagonalization:**  $B = [x_1 \ x_2 \ \dots \ x_n]$

Single layer perceptrons can classify linearly separable vectors.

#### Perceptron Learning Rule

$$W^{new} = W^{old} + e p^T, b^{new} = b^{old} + e, \text{ where } e = t - a$$

**Hebb's Postulate:** When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased.

**Linear Associator:**  $a = \text{purelin}(Wp)$

**The Hebb Rule:** Supervised Form:  $W_{ij}^{new} = W_{ij}^{old} + t_i p_j$

$$W = t_1 P_1^T + t_2 P_2^T + \dots + t_q P_q^T$$

$$W = [t_1 \ t_2 \ \dots \ t_q] \begin{bmatrix} P_1^T \\ P_2^T \\ \vdots \\ P_q^T \end{bmatrix} = TP^T$$

**Pseudoinverse Rule:**  $W = TP^+$

When the number,  $R$ , of rows of  $P$  is greater than the number of columns,  $Q$ , of  $P$  and the columns of  $P$  are independent, then the pseudoinverse can be computed by  $P^+ = (P^T P)^{-1} P^T$

**Variations of Hebbian Learning:**

**Filtered Learning:**  $W^{new} = (1 - \gamma)W^{old} + \alpha t_i p_i^T$

**Delta Rule (Ch10):**  $W^{new} = W^{old} + \alpha(t_i - a_i)p_i^T$

**Unsupervised Hebb (Ch15):**  $W^{new} = W^{old} + \alpha a_i p_i^T$

**Taylor:**  $F(x) = F(x') + \nabla F(x')^T (x - x') + \frac{1}{2} (x - x')^T \nabla^2 F(x') (x - x') + \dots$

**Grad:**  $\nabla F(x) = \left[ \frac{\partial}{\partial x_1} F(x) \quad \frac{\partial}{\partial x_2} F(x) \quad \dots \quad \frac{\partial}{\partial x_n} F(x) \right]^T$

**Hessian:**  $\nabla^2 F(x) = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} F(x) & \frac{\partial^2}{\partial x_1 \partial x_2} F(x) & \dots & \frac{\partial^2}{\partial x_1 \partial x_n} F(x) \\ \frac{\partial^2}{\partial x_2 \partial x_1} F(x) & \frac{\partial^2}{\partial x_2^2} F(x) & \dots & \frac{\partial^2}{\partial x_2 \partial x_n} F(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} F(x) & \frac{\partial^2}{\partial x_n \partial x_2} F(x) & \dots & \frac{\partial^2}{\partial x_n^2} F(x) \end{bmatrix}$

**Directional Derivatives:**

**1<sup>st</sup> Dir. Der.:**  $\frac{p^T \nabla F(x)}{\|p\|}$ , **2<sup>nd</sup> Dir. Der.:**  $\frac{p^T \nabla^2 F(x) p}{\|p\|^3}$

**Minima:**

**Strong Minimum:** if a scalar  $\delta > 0$  exists, such that  $F(x) < F(x + \Delta x)$  for all  $\Delta x$  such that  $\delta > \|\Delta x\| > 0$ .

**Global Minimum:** if  $F(x) < F(x + \Delta x)$  for all  $\Delta x \neq 0$

**Weak Minimum:** if it is not a strong minimum, and a scalar  $\delta > 0$  exists, such that  $F(x) \leq F(x + \Delta x)$  for all  $\Delta x$  such that  $\delta > \|\Delta x\| > 0$ .

**Necessary Conditions for Optimality:**

**1<sup>st</sup>-Order Condition:**  $\nabla F(x)|_{x=x^*} = 0$  (Stationary Points)

**2<sup>nd</sup>-Order Condition:**  $\nabla^2 F(x)|_{x=x^*} \geq 0$  (Positive Semi-definite Hessian Matrix).

**Quadratic fn.:**  $F(x) = \frac{1}{2} x^T A x + d^T x + c$

**Stable Learning Rule:** ( $\hat{u}_k = \hat{u}$ , constant)  $\hat{u} < \frac{1}{\lambda_{max}}$

$\{\lambda_1 \ \lambda_2 \ \dots \ \lambda_n\}$  Eigenvalues of Hessian matrix  $A$   
**Learning Rate to Minimize Along the Line:**

$$x_{k+1} = x_k + \alpha_k p_k \Rightarrow \hat{n}_k = -\frac{g_k^T p_k}{p_k^T A p_k} \quad (\text{For quadratic fn.})$$

**After Minimization Along the Line:**

$$x_{k+1} = x_k + \alpha_k p_k \Rightarrow g_{k+1}^T p_k = 0$$

**ADALINE:**  $a = \text{purelin}(Wp + b)$

**Mean Square Error:** (for ADALINE it is a quadratic fn.)

$$F(x) = E[e^2] = E[(t - a)^2] = E[(t - x^T z)^2]$$

$$F(x) = c - 2x^T h + x^T R x,$$

$$c = E[t^2], \quad h = E[tx] \quad \text{and} \quad R = E[xx^T] \Rightarrow A = 2R, \quad d = -2h$$

Unique minimum, if it exists, is  $x^* = R^{-1}h$ ,

$$\text{where } x = \begin{bmatrix} w \\ b \end{bmatrix} \text{ and } z = \begin{bmatrix} p \\ 1 \end{bmatrix}$$

**LMS Algorithm:**  $W(k+1) = W(k) + 2\alpha e(k) p^T(k)$

$$b(k+1) = b(k) + 2\alpha e(k)$$

**Convergence Point:**  $x^* = R^{-1}h$

**Stable Learning Rate:**  $0 < \alpha < 1/\lambda_{max}$  where

$\lambda_{max}$  is the maximum eigenvalue of  $R$

**Adaptive Filter ADALINE:**

$$a(k) = \text{purelin}(Wp(k) + b) = \sum_{i=1}^p w_i y(k-i+1) + b$$

**Backpropagation Algorithm:**

**Performance Index:**

$$\text{Mean Square error: } F(x) = E[e^2] = E[(t - a)^2]$$

**Approximate Performance Index:** (single sample)

$$\hat{F}(x) = e^T(k) e(k) = (t(k) - a(k))^2$$

$$Sensitivity: s^m = \frac{\partial \hat{F}}{\partial a^m} = \left[ \frac{\partial \hat{F}}{\partial a_1^m} \quad \frac{\partial \hat{F}}{\partial a_2^m} \quad \dots \quad \frac{\partial \hat{F}}{\partial a_n^m} \right]^T$$

**Forward Propagation:**  $a^0 = p$ ,

$$a^{m+1} = f^{m+1}(W^{m+1} a^m + b^{m+1}) \text{ for } m = 0, 1, \dots, M-1$$

$$a = a^M$$

**Backward Propagation:**  $s^M = -2\hat{F}'(a^M)(t - a)$ ,

$$s^m = \hat{F}'^m(a^m)(W^{m+1})^T s^{m+1} \text{ for } m = M-1, \dots, 1, \text{ where}$$

$$\hat{F}'^m(a^m) = \text{diag}([f'^m(n_1^m) \quad f'^m(n_2^m) \quad \dots \quad f'^m(n_n^m)])$$

$$f'^m(n_j^m) = \frac{\partial f^m(n_j^m)}{\partial n_j^m}$$

**Weight Update (Approximate Steepest Descent):**

$$W^m(k+1) = W^m(k) - \alpha s^m (a^{m-1})^T$$

$$b^m(k+1) = b^m(k) - \alpha s^m$$

**hardlim:**  $a = \begin{cases} 1 & a < 0 \\ 0 & a \geq 0 \end{cases}$ , **hardlims:**  $a = \begin{cases} -1 & a < 0 \\ 1 & a \geq 0 \end{cases}$ , **purelin:**  $a = x$ , **logsig:**  $a = \frac{1}{1 + e^{-x}}$ , **tanstg:**  $a = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ , **poslin:**  $a = \begin{cases} 0 & a < 0 \\ m & a \geq 0 \end{cases}$

**Backpropagation with Momentum (MOBP):**

$$\Delta W^m(k) = \gamma \Delta W^m(k-1) - (1 - \gamma) \alpha s^m (a^{m-1})^T$$

$$\Delta b^m(k) = \gamma \Delta b^m(k-1) - (1 - \gamma) \alpha s^m$$

**Variable Learning Rate Backpropagation (VLBP)**

1. If the squared error (over the entire training set) increases by more than some set percentage,  $\zeta$  (typically one to five percent) after a weight update, then the weight update is discarded, the learning rate is multiplied by some factor  $\rho < 1$ , and the momentum coefficient  $\gamma$  (if it is used) is set to zero.

2. If the squared error decreases after a weight update, then the weight update is accepted and the learning rate is multiplied by some factor  $\eta > 1$ . If  $\gamma$  has been previously set to zero, it is reset to its original value.

3. If the squared error increases by less than  $\zeta$ , then the weight update is accepted but the learning rate and the momentum coefficient are unchanged.

**Association:**  $a = \text{hardlim}(W^T p + b)$

An association is a link between the inputs and outputs of a network, so that when a stimulus  $A$  is presented to the network, it will output a response  $B$ .

**Associative Learning Rules:**

**Unsupervised Hebb Rule:**

$$W(q) = W(q-1) + \alpha a(q) p^T(q)$$

**Hebb with Decay:**

$$W(q) = (1 - \gamma)W(q-1) + \alpha a(q) p^T(q)$$

**Instar:**  $a = \text{hardlim}(Wp + b)$ ,  $a = \text{hardlim}(w^T p + b)$

The instar is activated for  $-1 w^T p = \|w\| \|p\| \cos \theta \geq -b$

where  $\theta$  is the angle between  $p$  and  $w$ .

**Instar Rule:**

$$w(q) = w(q-1) + \alpha a(q)(p(q) - w(q-1))$$

$$w(q) = (1 - \alpha) w(q-1) + \alpha p(q), \text{ if } (a_i(q) = 1)$$

**Kohonen Rule:**

$$w_j(q) = w_j(q-1) + \alpha (a(q) - w_j(q-1)) p_j(q)$$

**Outstar Rule:**  $a = \text{satslin}(Wp)$

$$w_j(q) = w_j(q-1) + \alpha (a(q) - w_j(q-1)) p_j(q)$$

**Competitive Layer:**  $a = \text{compet}(Wp) = \text{compet}(u)$

**Competitive Learning with the Kohonen Rule:**

$$w_j(q) = w_j(q-1) + \alpha (p(q) - w_j(q-1))$$

$$= (1 - \alpha) w_j(q-1) + \alpha p(q)$$

$$w_j(q) = w_j(q-1), \quad i \neq i^* \text{ where } i^* \text{ is the winning neuron.}$$

**Self-Organizing with the Kohonen Rule:**

$$w_j(q) = w_j(q-1) + \alpha (p(q) - w_j(q-1))$$

$$= (1 - \alpha) w_j(q-1) + \alpha p(q), \quad i \in N_c(d)$$

$$N_c(d) = \{j, d_{ij} \leq d\}$$

**LVO Network:**  $(w_i^k = 1) \Rightarrow$  subcell  $i$  is a part of class  $k$

$$n_i^k = -\|w_i^k - p\|, a^k = \text{compet}(n^k), a^2 = W^2 a^1$$

**LVO Network Learning with the Kohonen Rule:**

$$w_i^k(q) = w_i^k(q-1) + \alpha (p(q) - w_i^k(q-1))$$

$$\text{if } a_i^k = t_k = 1$$

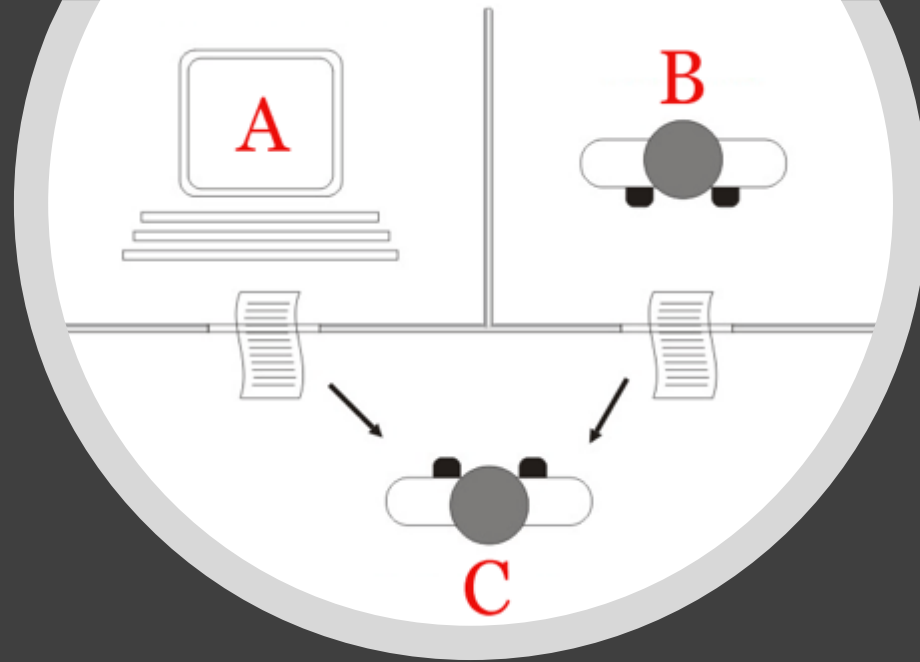
$$w_i^k(q) = w_i^k(q-1) - \alpha (p(q) - w_i^k(q-1))$$

$$\text{if } a_i^k = 1 \neq t_k = 0$$





“Dad, what are you doing at work?”



Computer can help you to do something for you and people will think you are doing that.

# Based on Turing Test...







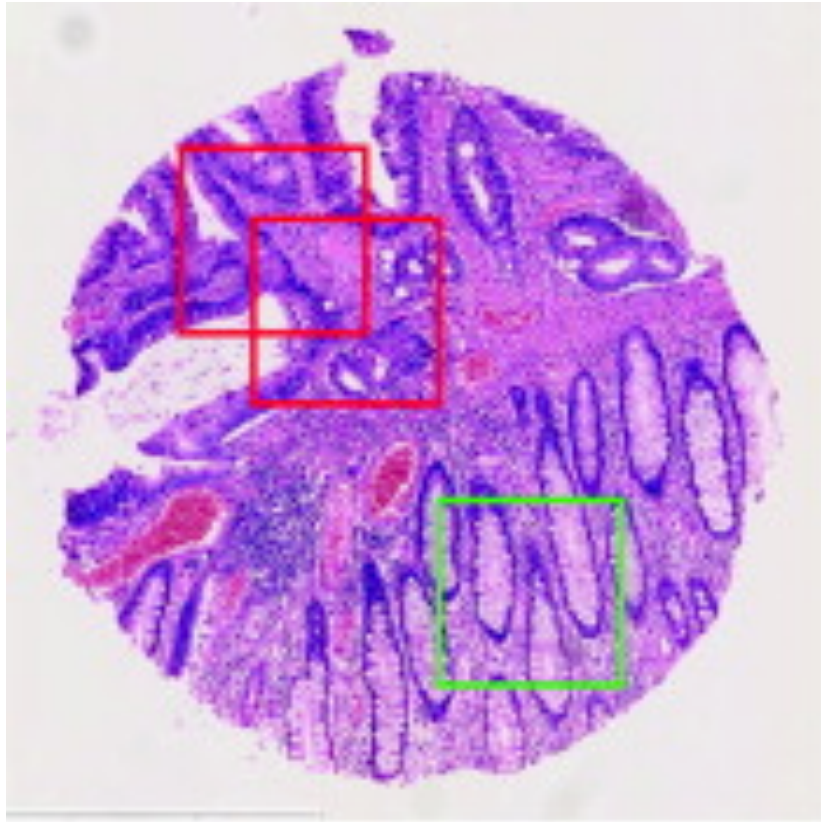
# How did you teach your younger sister?

---

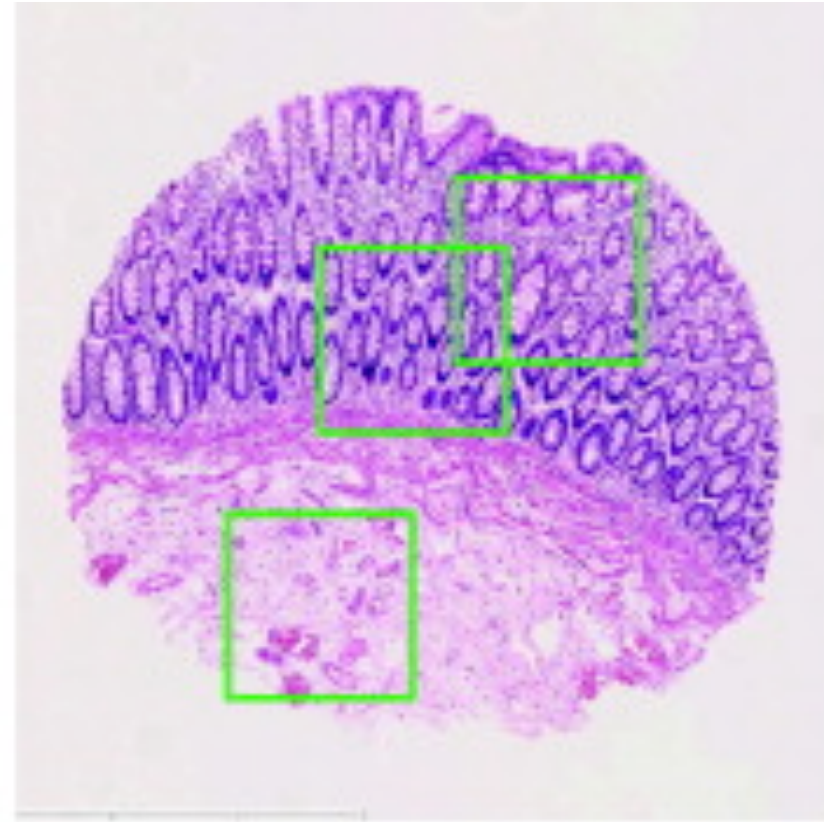
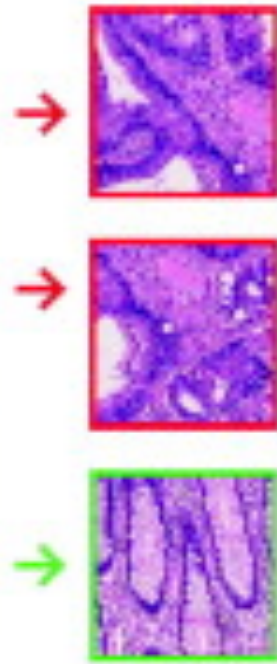
Supervised Learning - Maps an input to an output based on example input-output pairs



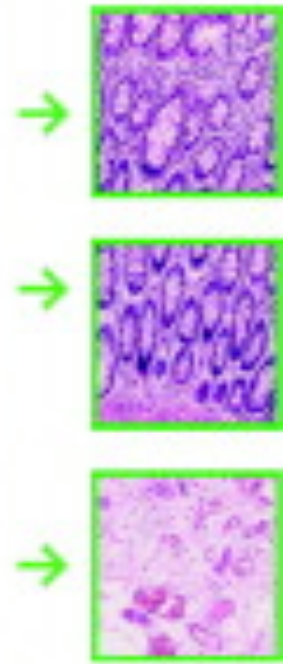
# Supervised Learning Example – Detecting Cancer



(a) cancer image



(b) non-cancer image



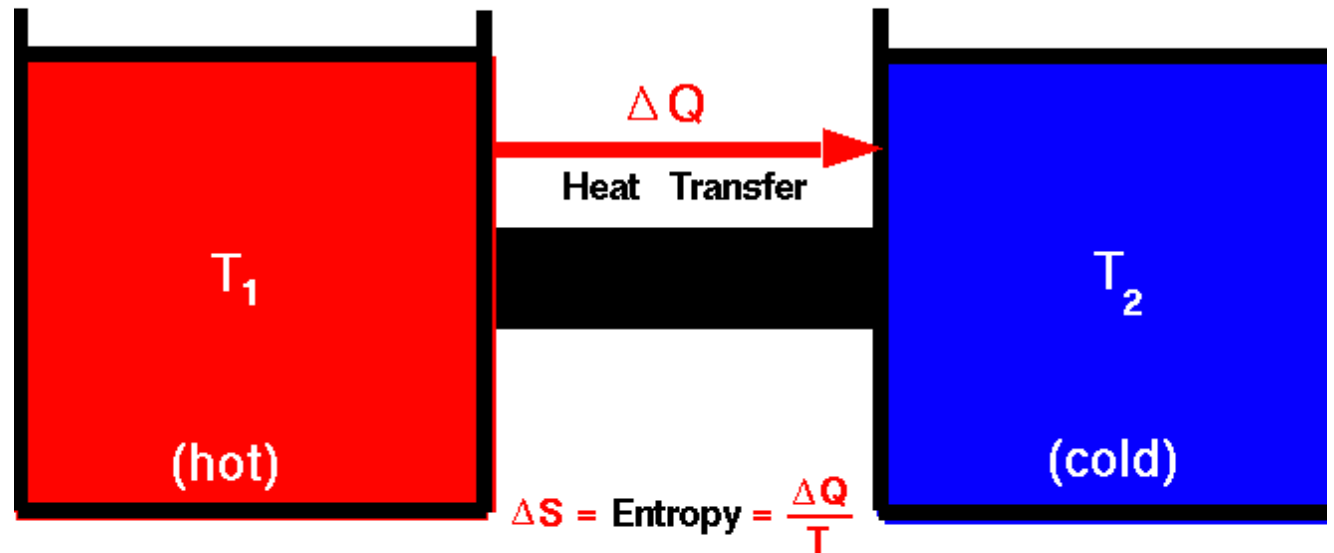
## Law of Entropy

- The total entropy of an isolated system can never decrease over time



## Second Law of Thermodynamics

Glenn  
Research  
Center



There exists a useful thermodynamic variable called entropy (S). A natural process that starts in one equilibrium state and ends in another will go in the direction that causes the entropy of the system plus the environment to increase for an irreversible process and to remain constant for a reversible process.

$$S_f = S_i \text{ (reversible)}$$

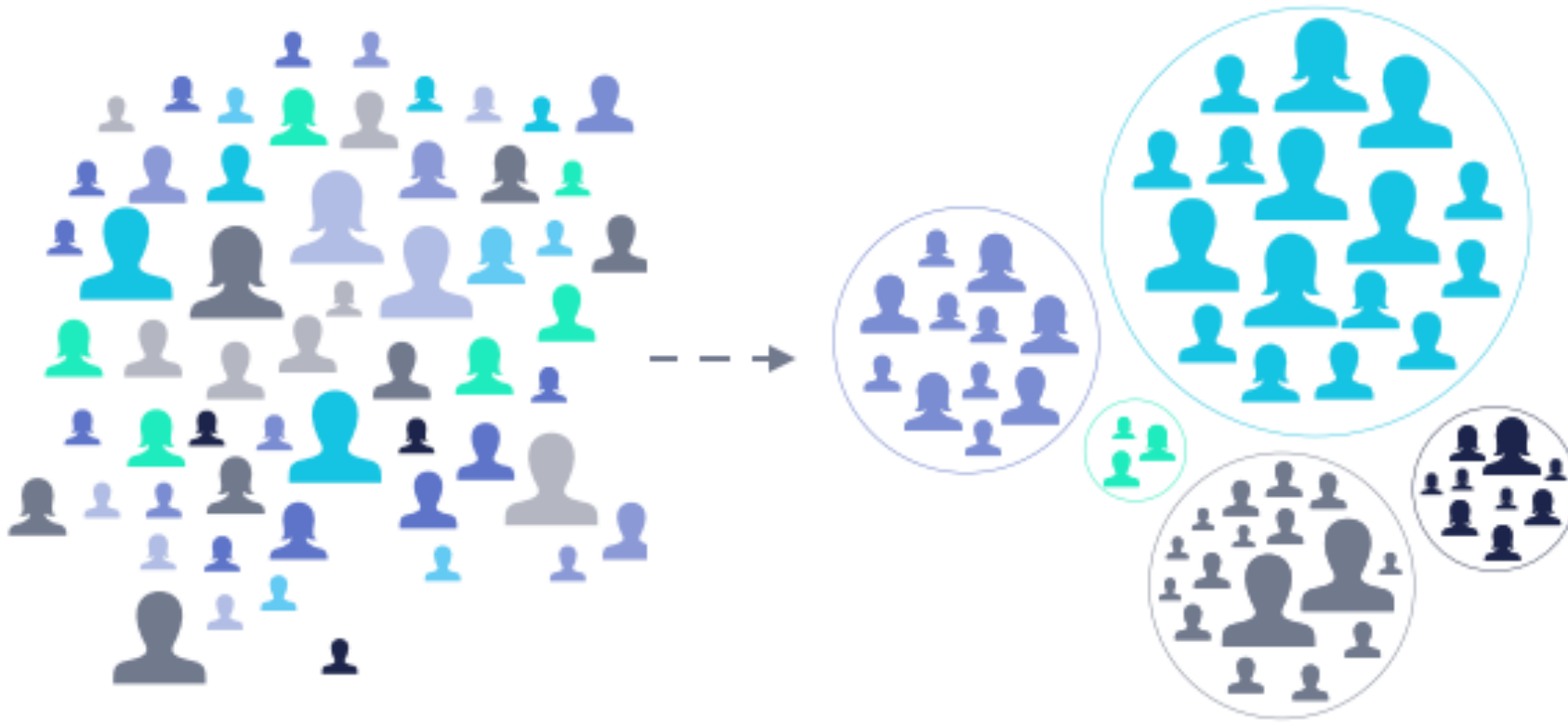
$$S_f > S_i \text{ (irreversible)}$$



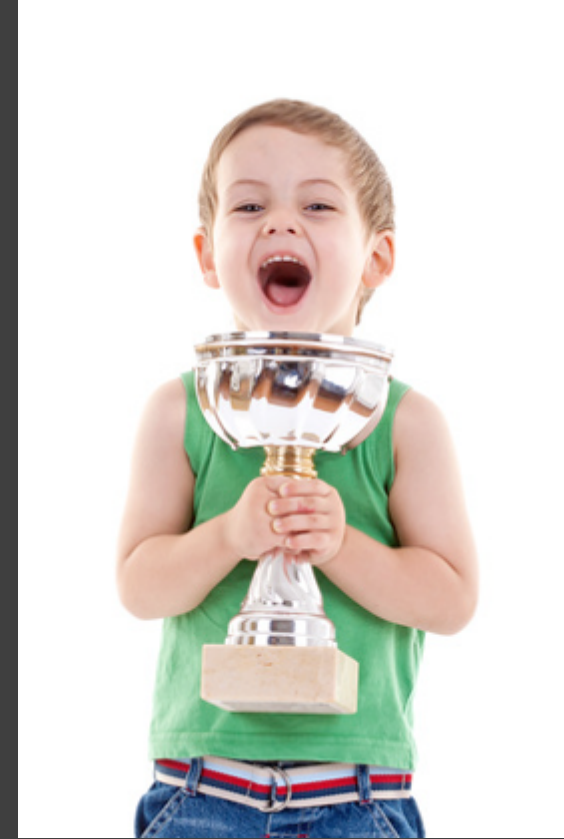
A photograph of a young child sitting amidst a vast, chaotic sea of toys and cardboard boxes. The child is looking towards the camera. The room is filled with various items, including a large white teddy bear, a blue toy car, a yellow toy car, and numerous other colorful toys and boxes. The scene is cluttered and disorganized, visually representing the concept of 'unlabeled data' mentioned in the text.

# Unsupervised Learning – Uncovering hidden patterns from unlabeled data

# Unsupervised Learning Example – Customer Segmentation

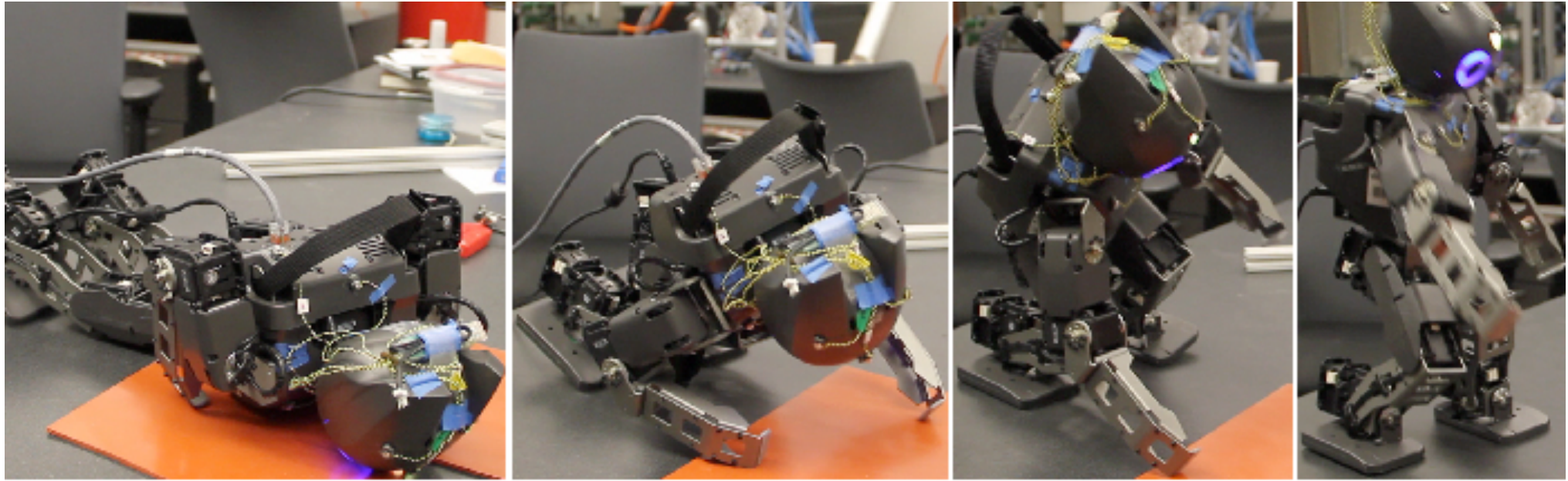






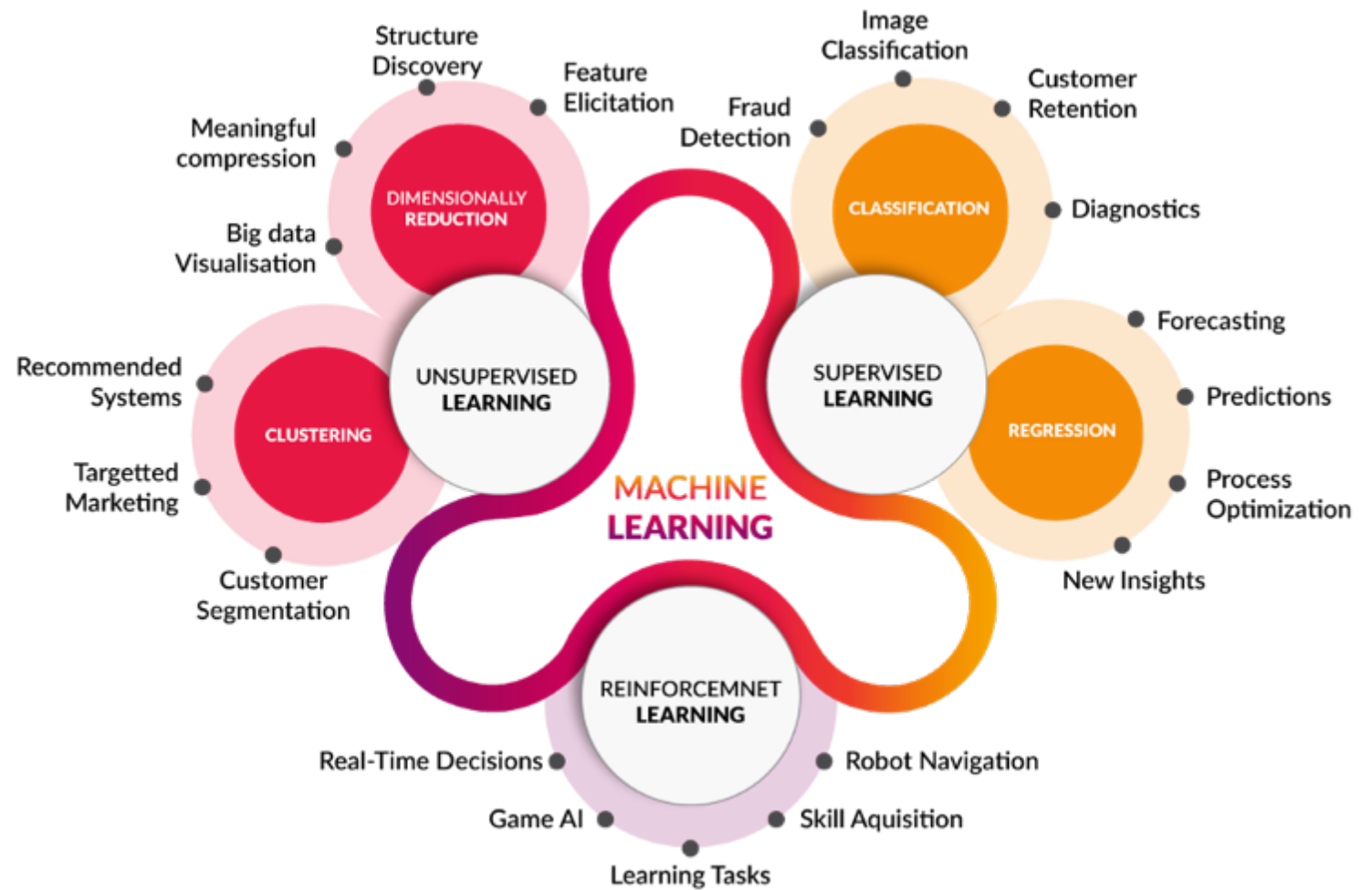
# Reinforcement Learning

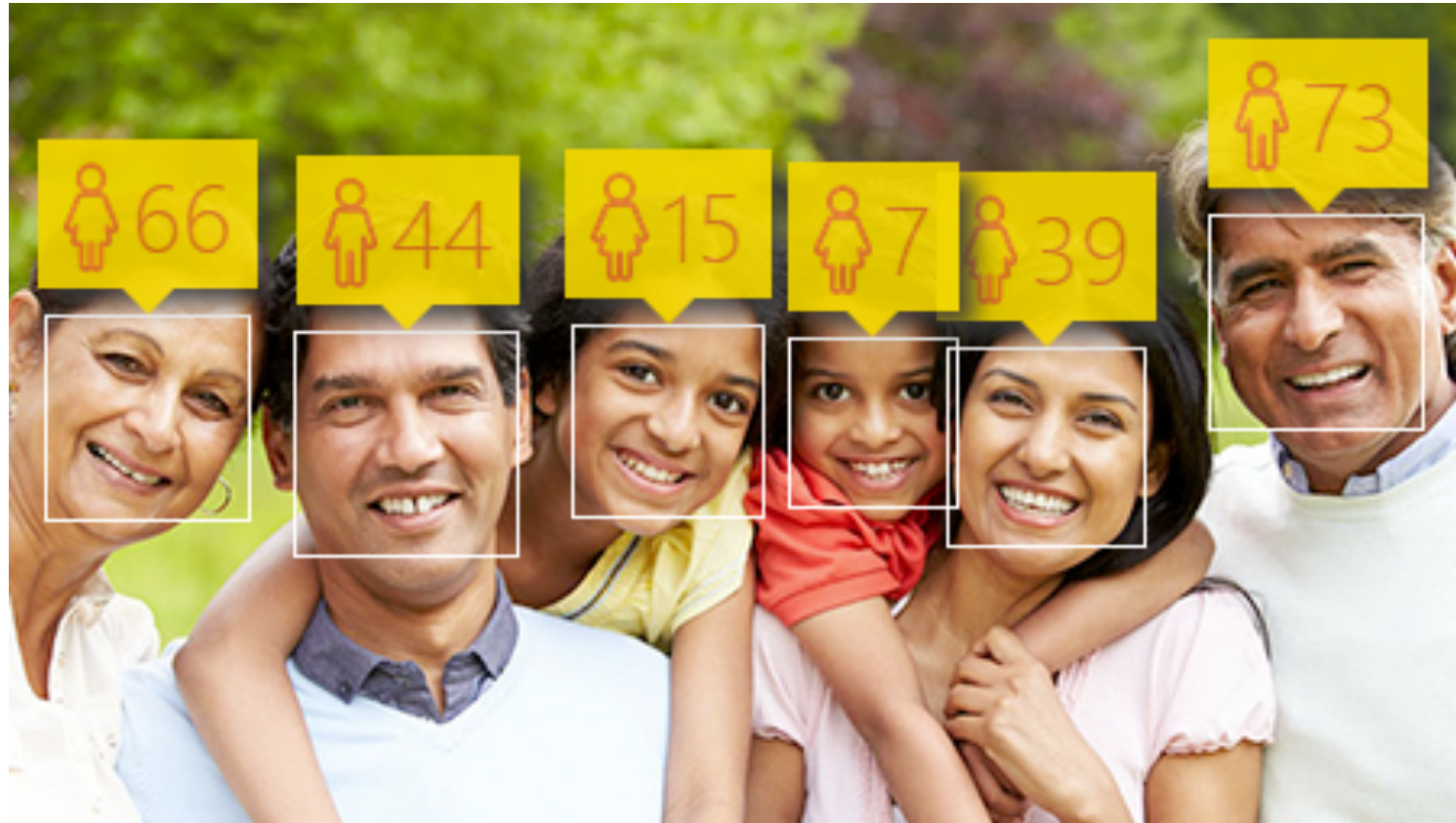
Take actions in an environment so as to maximize some notion of cumulative reward



Reinforcement Learning Example - Robotics

# Major Types of AI



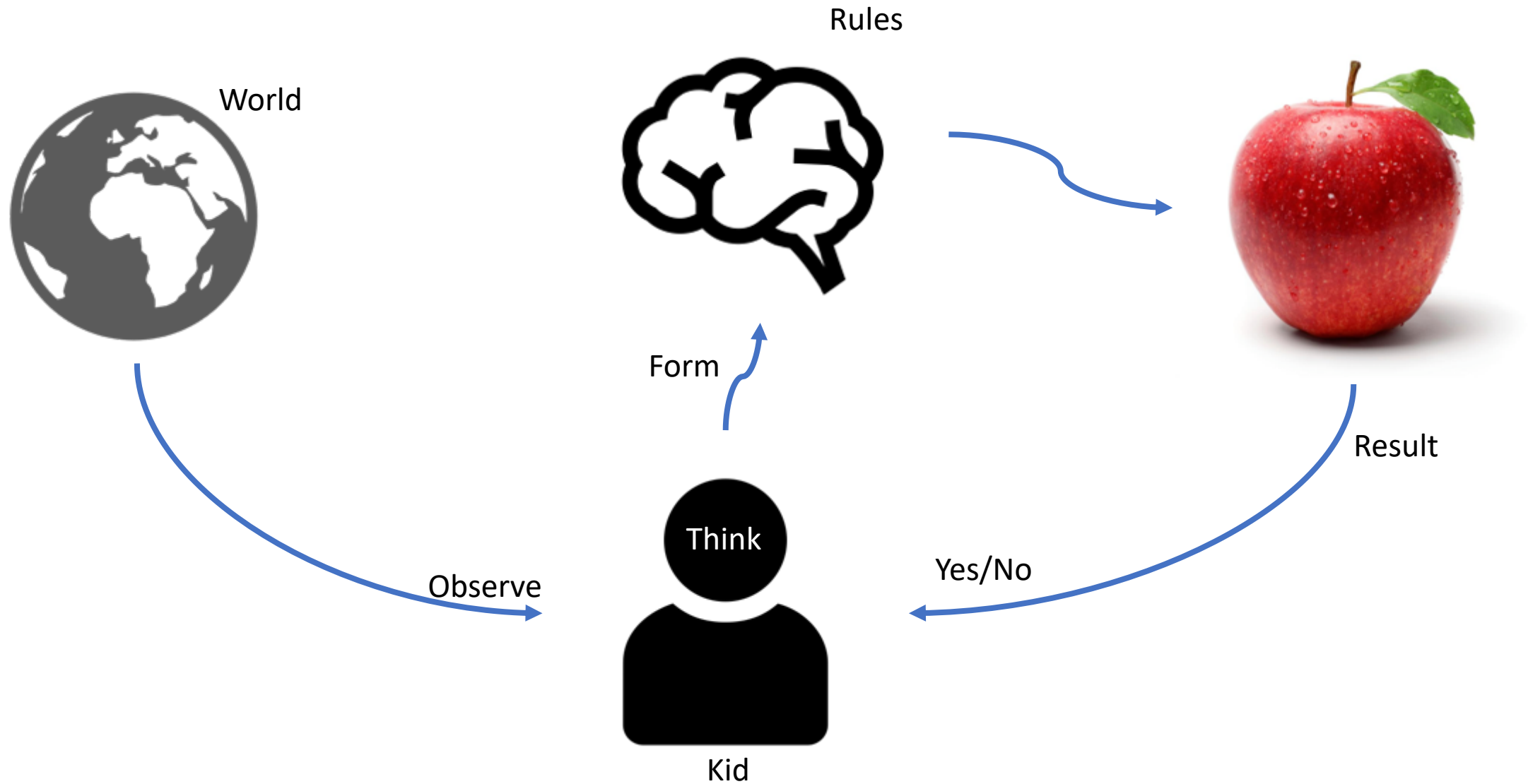


**Regression** – a prediction method with a real number as an output  
**Classification** – a prediction method with a predefined category

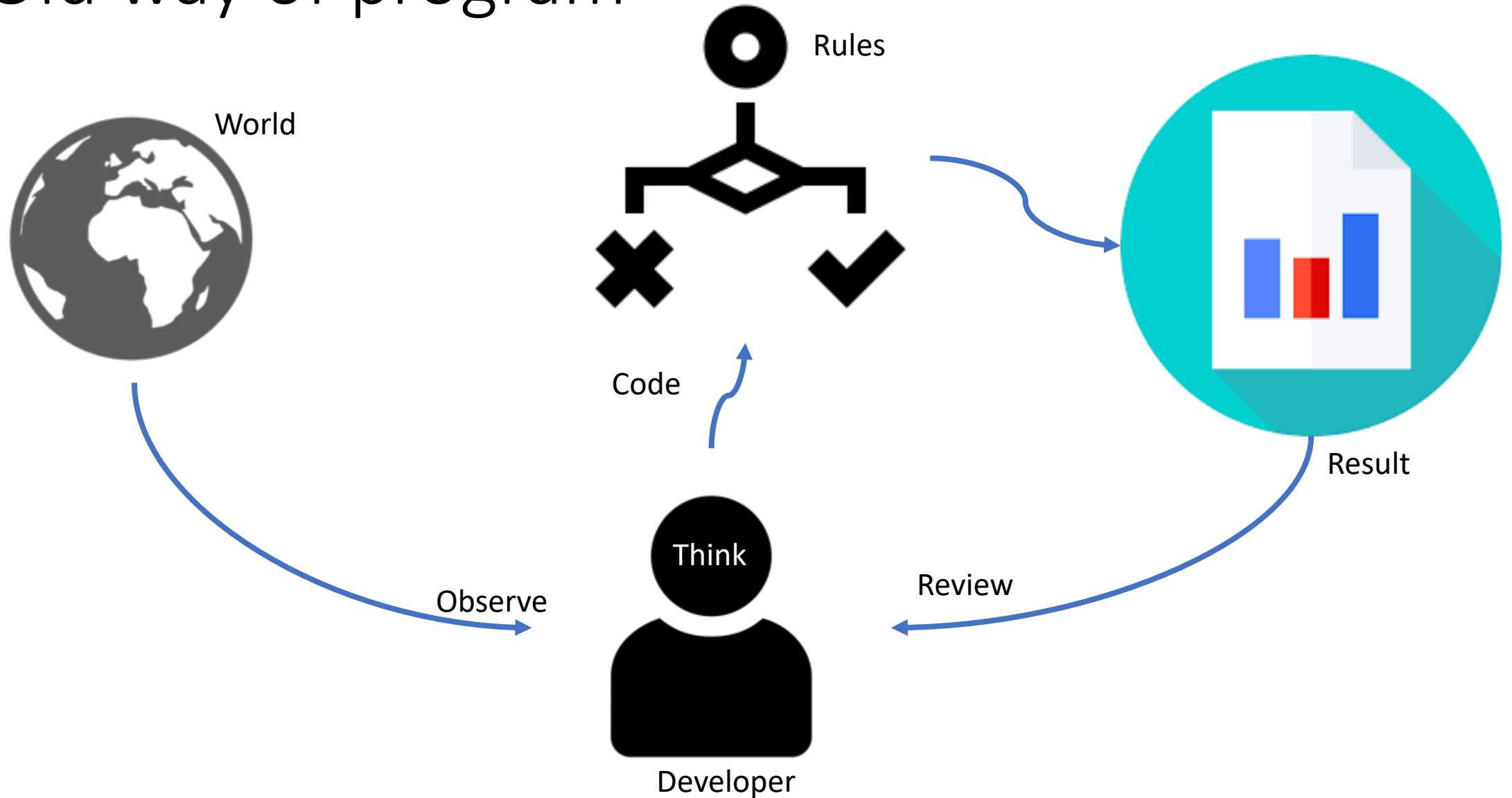
---



# How a kid learn?

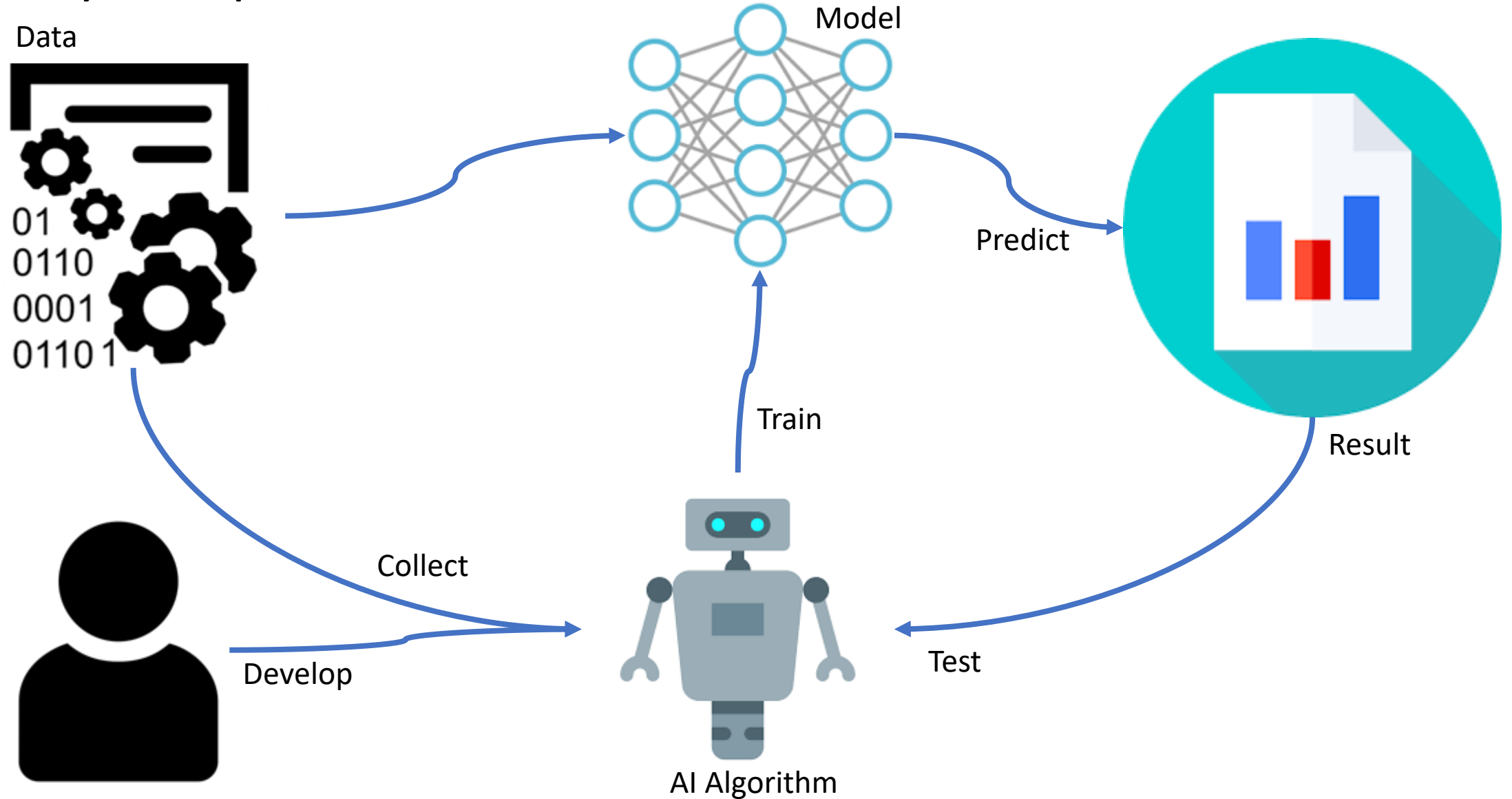


# Old way of program





# Key Steps of AI



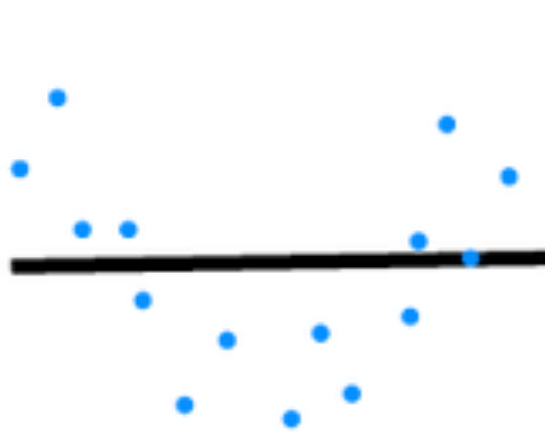


Give a Man a Fish, and You Feed Him for a Day...  
Teach a Man To Fish, and You Feed Him for a Lifetime – **Proverb**





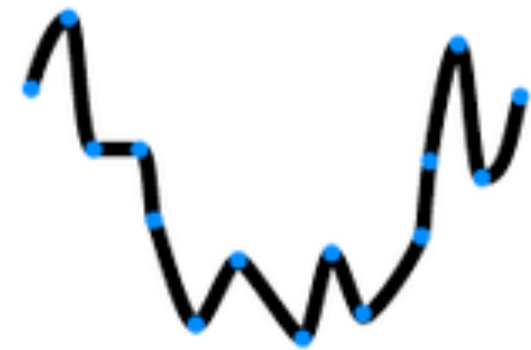
Give a Computer a Code, and You Help It Solving for one Situation...  
Teach a Computer To Code, and You Help It Solving for a Collection of Situations – **Cupid Proverb**



Underfitting



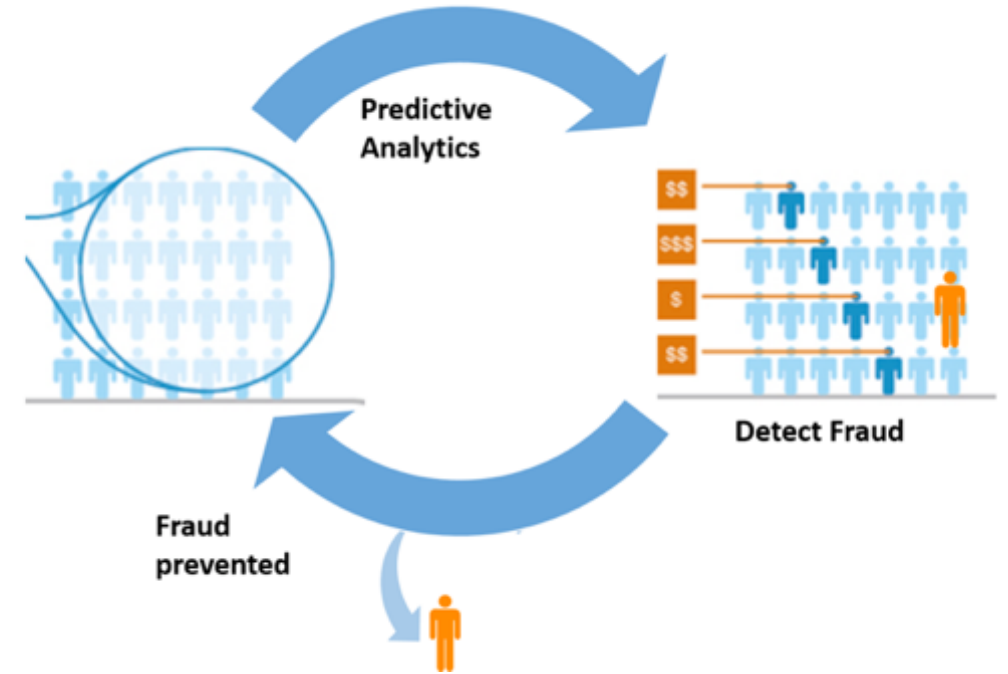
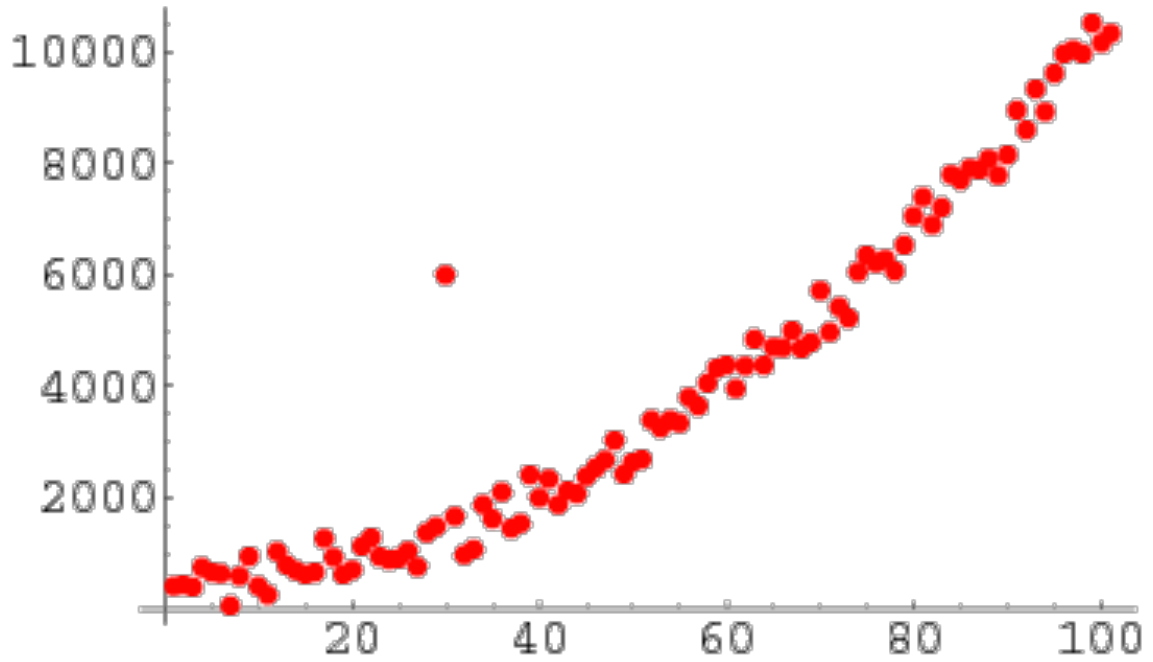
Desired



Overfitting

## Memorize VS Generalize

- What is  $2 + 1$ ?
- How about  $1 + 2$ ?



# Fraud Detection



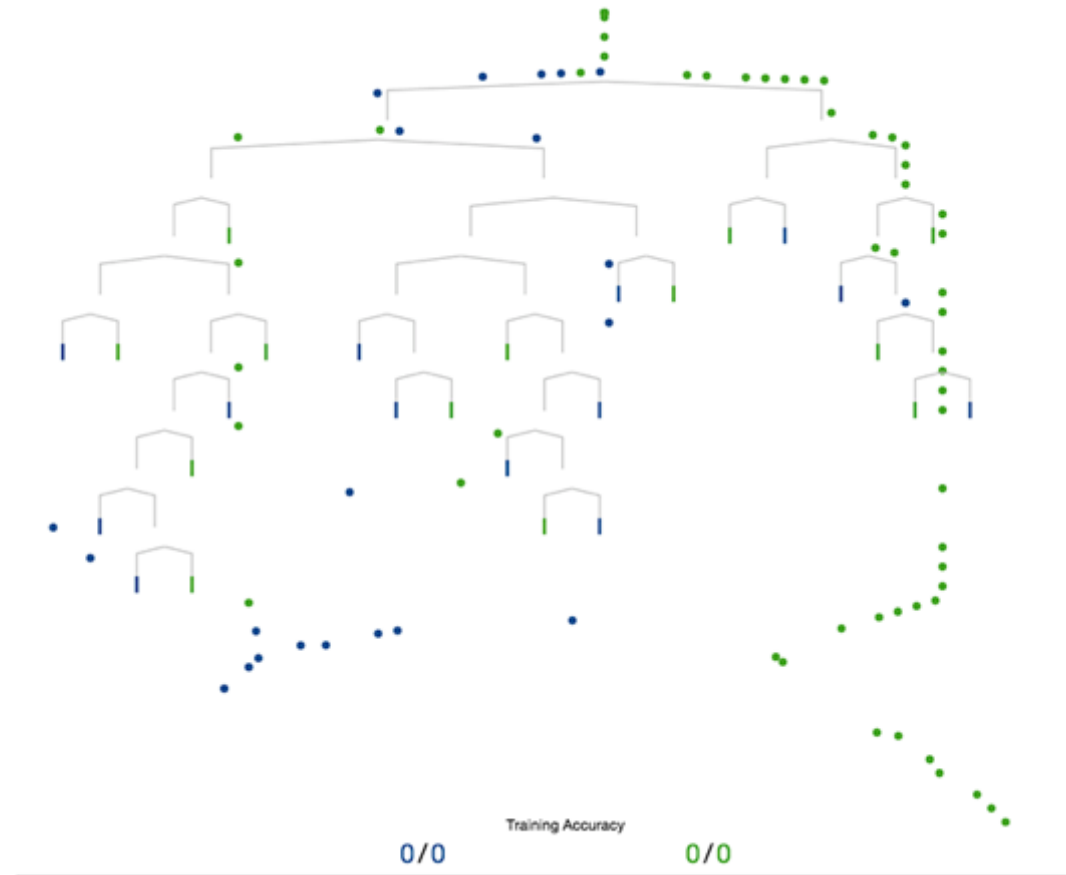
McDonald "Guess Who?"



# Decision Tree

## Making predictions

The newly-trained decision tree model determines whether a home is in San Francisco or New York by running each data point through the branches.





# Main Algorithms of AI







“AI can help me to do homework and the teacher will not find out I didn’t do that?”



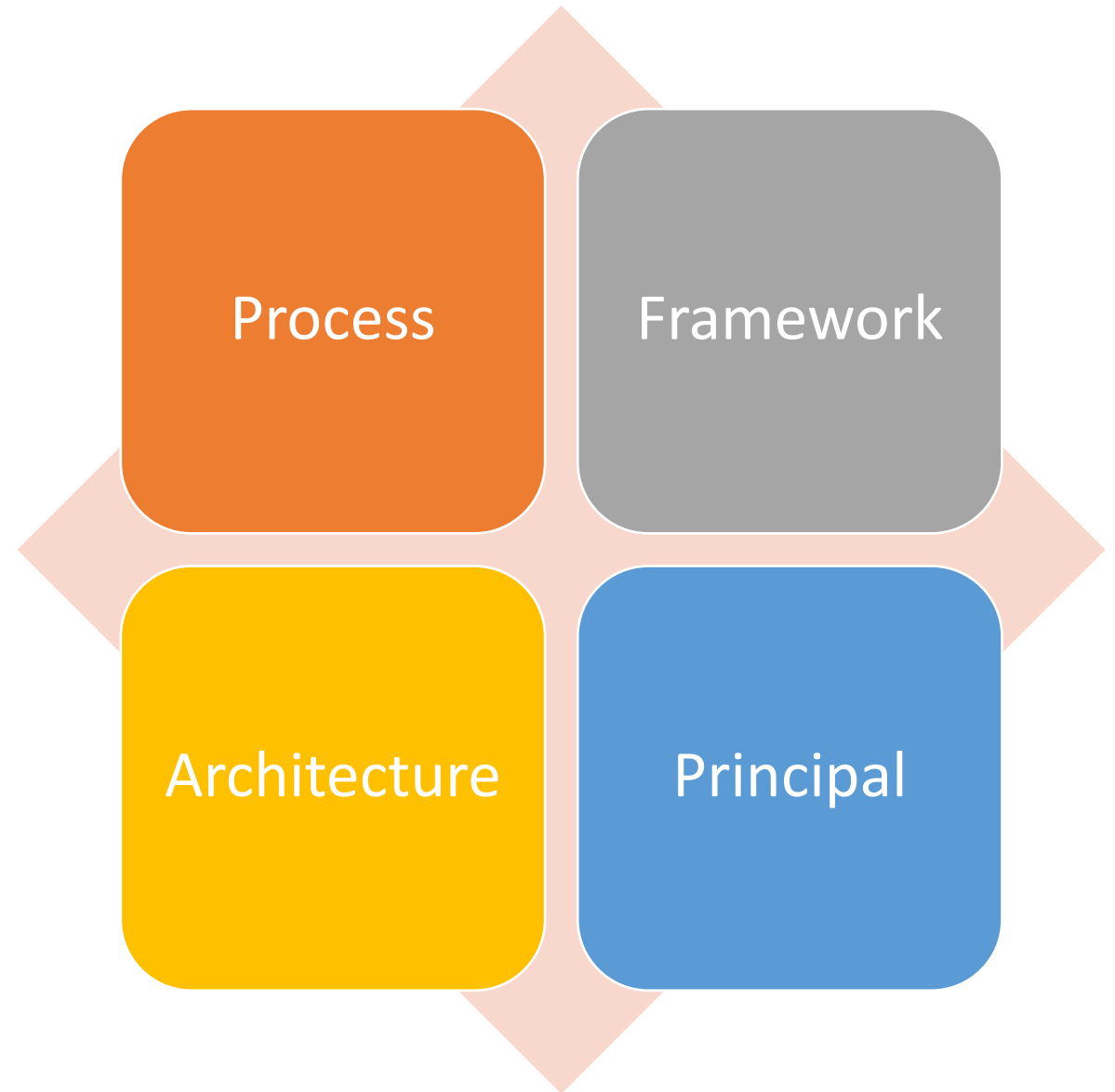
There is no Right or Wrong, only Right or Left. But no matter which direction you pick, be persistent and you will cross the finish line of success via either route



# BI is dead

Unless you incorporate AI and other approaches into your BI Strategy

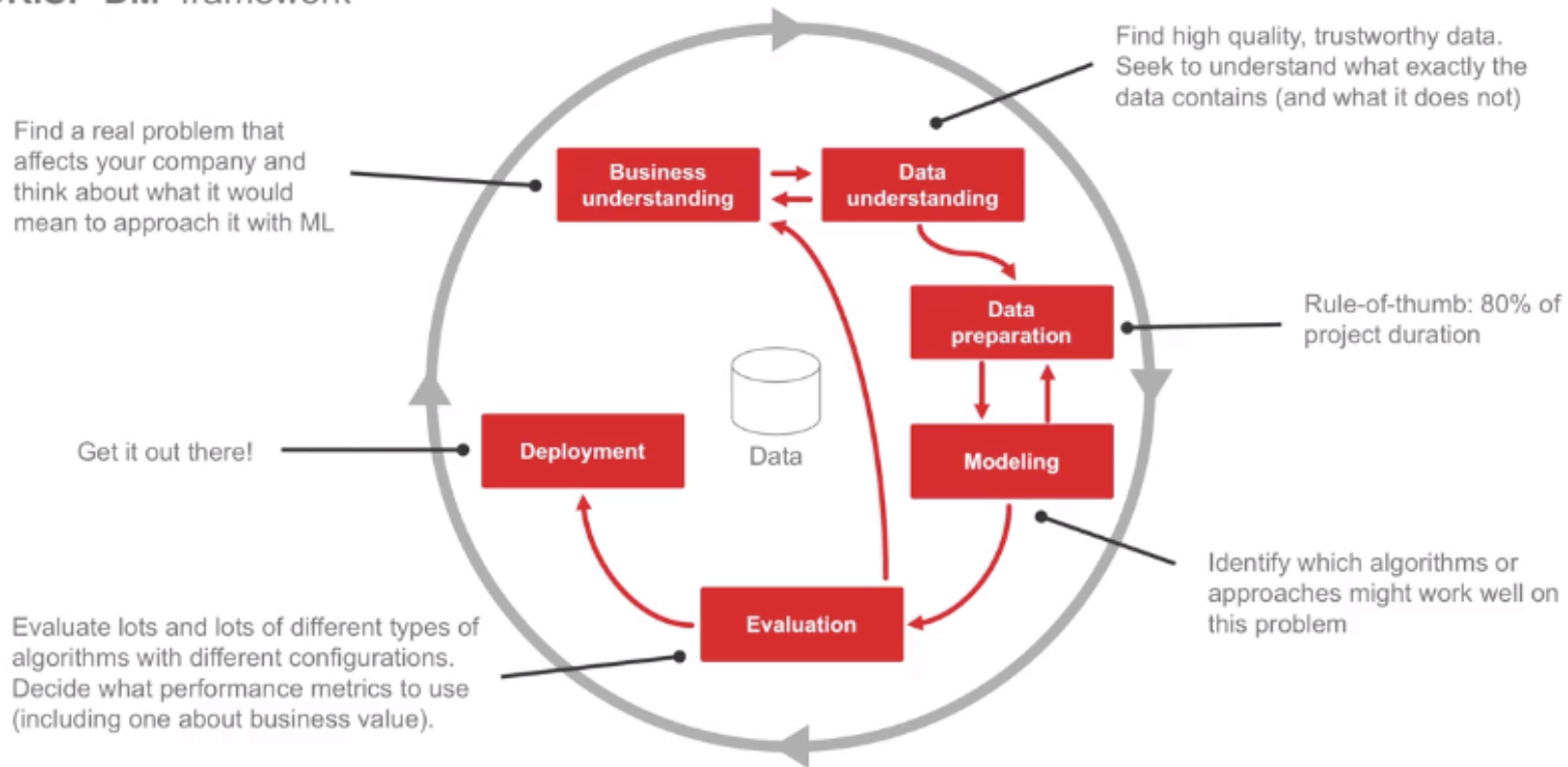
# Strategies to resurrect BI



# Process: Cross-industry standard process for data mining

## How its done in the real world

## CRISP-DM<sup>1</sup> framework





Data Consumer



Business Analyst / Business Owners

Service



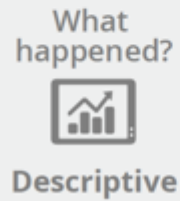
Data Science Platform



API

Common

Process



Descriptive



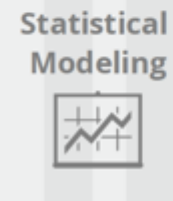
Diagnostic



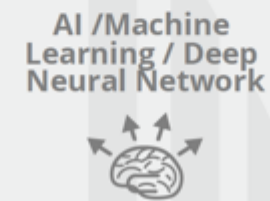
Predictive



Prescriptive



Statistical Modeling



AI / Machine Learning / Deep Neural Network



Metadata



Governance



Security



Operation & Orchestration

Store



Landing Zone



Exploration Zone



Enterprise Zone



Olap Cube



Data Mart



Operational Data Store

Ingest



Batch Processing (SFTP, files)



Stream Processing (API call to data sources, streaming system or message subscription)

Data Origin



Social



Apps



System Logs



IoT



INDEX | ANALYTICS ARCHITECTURE

# Principal

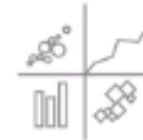
---



**Single comprehensive platform** to centralize business data used for machine learning



**One version of the truth** helps protect against data scientists training models on the wrong data



**Enterprise grade self-service** that empowers business users and data scientists to work together



**Open to the ecosystem** of machine learning technology and services



**APIs** to expose your governed business data to cutting edge open-source and vendor ML products used by your data scientists today



**Enable the workforce** to take action on predictive workflows wherever they are

# Framework: S.O.W.E

---







Very Good Story for a presentation Cupid  
But I don't believe you really taught your kid AI as you said

# Cupid Chan

- CTO of Index Analytics
- Board of Directors and TSC, Linux Foundation ODPI
- Organizer of a Big Data In Action Meetup in Washington/Baltimore area with 1600 members



[www.linkedin.com/in/cupidchan/](http://www.linkedin.com/in/cupidchan/)



[@cupidckchan](https://twitter.com/cupidckchan)

