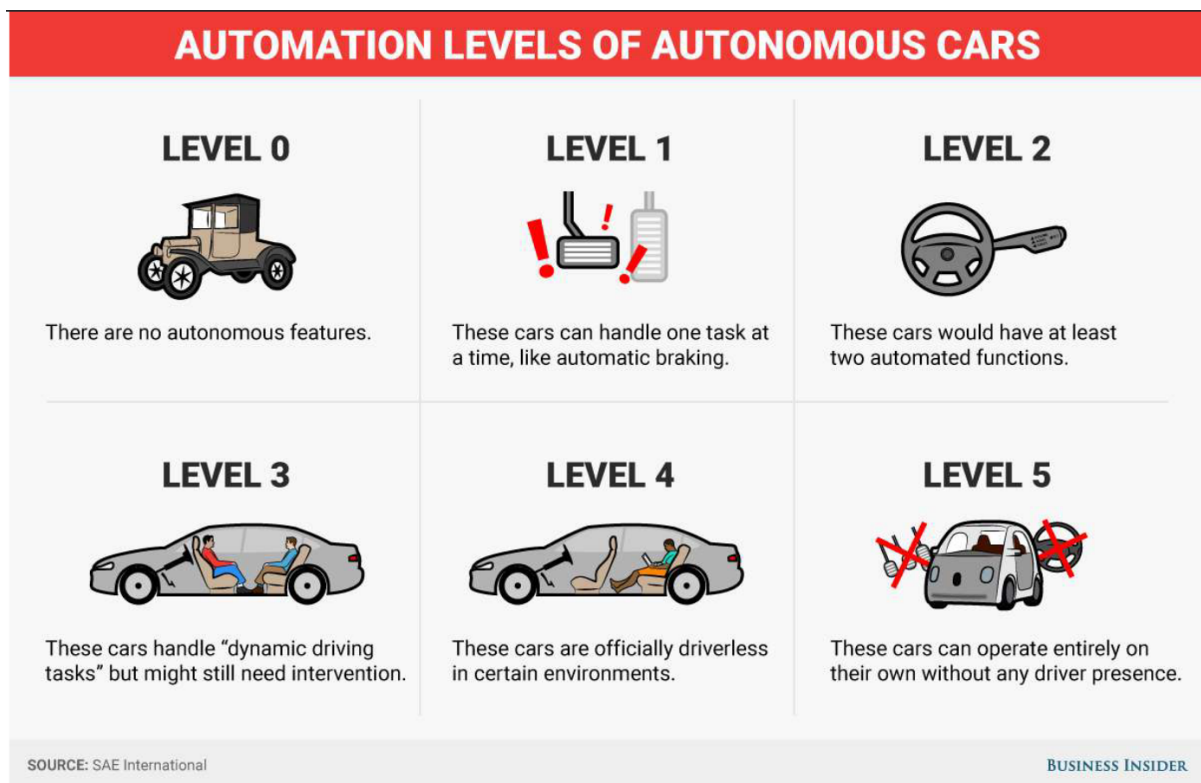


# Lesson 5\_Ethical\_Knob

## A Genetic Approach to the Ethical Knob

When we talk about Ethical Knob we are talking about Autonomous Vehicles (AVs), which is a particular area of research. We are still in semi-autonomous devices era.

Autonomous Driving is classified according to the amount of human driver intervention:



AVs is a very interesting challenge because with the increasing of level of autonomy the amount of collected data grows very steeply because of the amount of sensors. We estimate about 4.4 GB/s Data Logging for full Autonomous Driving

## CAR AUTOMATION SENSORS & DATA VOLUMES

Sensor type	Quantity	Data generated
Radar	4–6	0.1–15 Mbit/s
LIDAR	1–5	20–100 Mbit/s
Camera	6–12	500–3,500 Mbit/s
Ultrasonic	8–16	<0.01 Mbit/s
Vehicle motion, GNSS, IMU	-	<0.1 Mbit/s

### TOTAL ESTIMATED BANDWIDTH

**3 Gbit/s (~1.4TB/h) to 40 Gbit/s (~19 TB/h)**

Of course, AVs can fail, in particular when they have to deal with the uncertainty of the real-world.

**At the very beginning**, one of the most famous experiment was the **Moral Machine** that collected people's decision when they faced a moral dilemma. The idea was to collect different answers in different scenarios to describe the behavior of people in some specific situations. The knob expresses directly the ethical attitude of the AV passengers: the value passengers attribute to their life relative to the value of the lives of third parties (altruistic, selfish, etc.)

**In the new proposal**, the position of the knob no longer indicates the passengers' moral attitude, but it indicates the AV's assessment of the relative importance of the lives of passenger(s) and third parties. This assessment is made by looking at how many pedestrians are present, how much these pedestrians may be hurt by the vehicle, which are the passengers, which is the moral attitude selected by the subject, how much the passengers may be hurt if the AV swerves etc. Thus, in this new setting the vehicle is able to perceive the environment and compute an assessment of risk for who is inside and outside the vehicle.

### HOW TO DO THAT?

#### Combination of AI techniques:

**Neural networks** to compute the right action to take based on the given scenario

**Genetic Algorithm** to find an (almost) optimal configuration of neural networks

Thus, the NNs are not trained, but it is the GA who searches for the (almost) optimal weights configuration.

### Genetic Algorithms

- Inspired by Charles Darwin's theory of natural evolution.
- The fittest individuals are selected for reproduction in order to produce offspring of the next generation
- Heuristic Search in the solution space
- Mostly used in optimization tasks

Nelle slides è riportato l'esempio delle n-queens per spiegare che cos'è un GA, ma direi che lo sappiamo bene tutti e 4 cos'è e non insulto la nostra intelligenza riscrivendolo -.-'

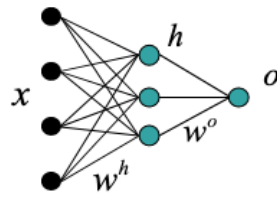
Spiega poi cos'è una rete neurale e anche qui mi sembra inutile riscriverlo :)

### Why do we use GA for the ethical knob instead of backpropagation and gradient descent?

Non l'ha spiegato, ha saltato tutte le slides delle reti neurali ed è passato direttamente a 'Simulation'.

### Simulation

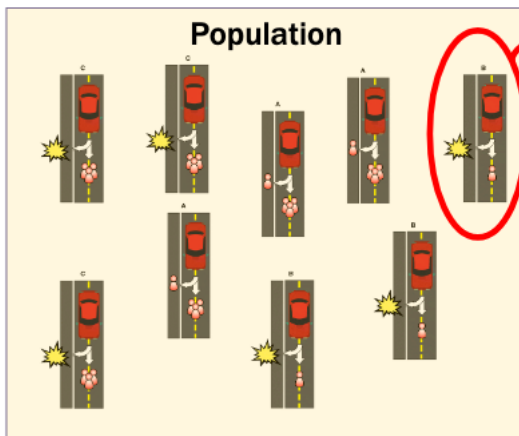
An individual in the simulation corresponds to an AV, so each neural network represents an AV that at the beginning behave randomly.



We represent an AV using a NN. The NN:

- Analyzes the scenario
- Outputs the level of the knob

The knob value is used to take an action



Any scenario has:

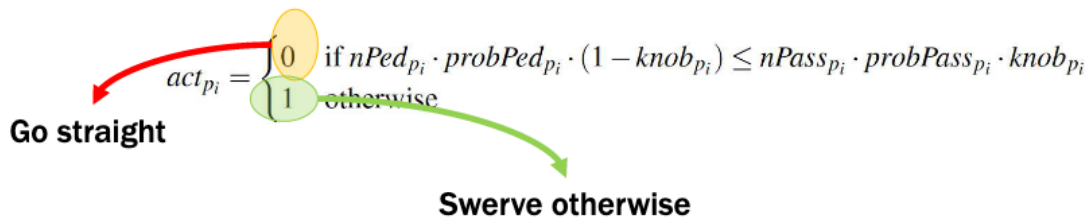
- Altruism level
- Number of passengers
- Prob. of harming passengers
- Number of pedestrians
- Prob. of harming pedestrians

Le formule sotto non sono da imparare, ma da capire come concetto generale.

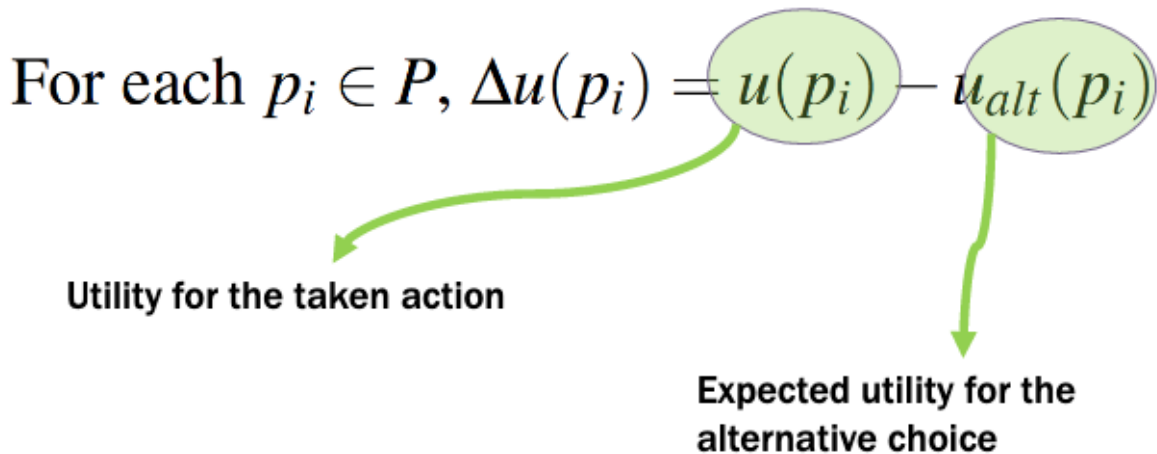
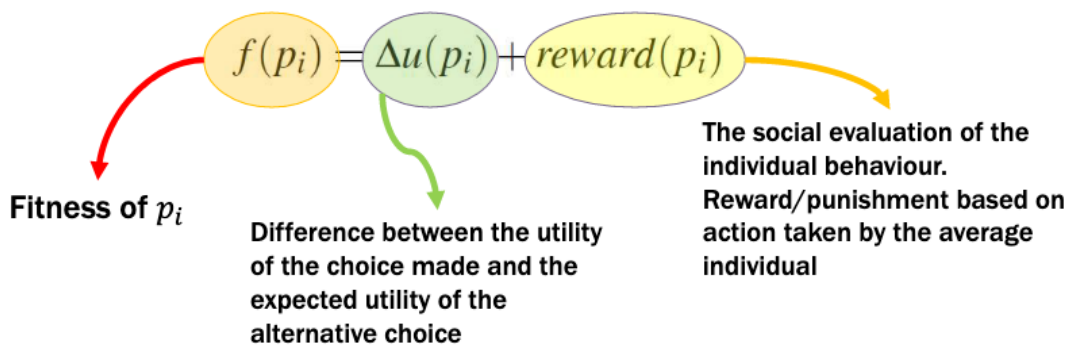
**The notation:**

- $nPed_{p_i}$ : number of pedestrians
- $nPass_{p_i}$ : number of passengers
- $a_{p_i}$ : intrinsic level of altruism for passengers in  $p_i$
- $s_{p_i}$ : intrinsic level of selfishness for passengers in  $p_i$
- $prodPed_{p_i}$ : probability of injuring pedestrians when the AV goes straight
- $prodPass_{p_i}$ : probability of injuring passengers when the AV swerves

The action is taken based on the assessment computed by the NN. The idea is pondering which action minimize harm with respect to relative importance of lives:



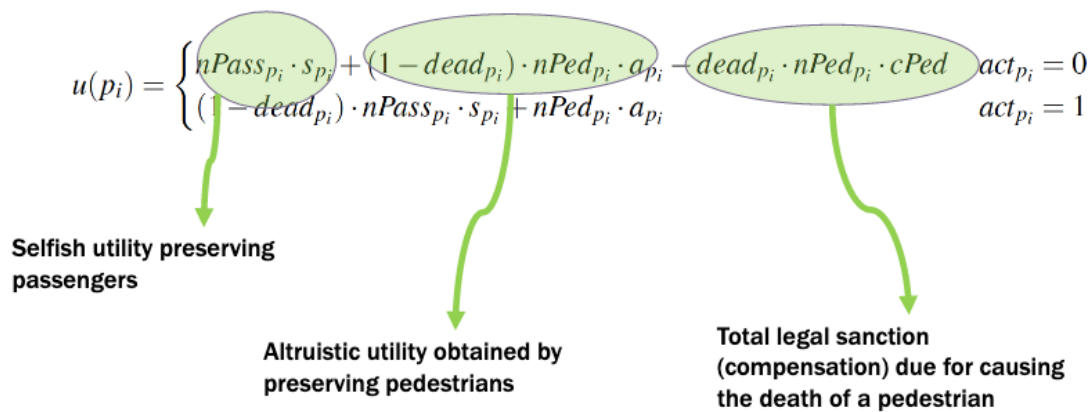
Individual is evaluated using the following fitness function:



Depending on the taken action, the utility is computed based on the response of the scenario:

$$u(p_i) = \begin{cases} nPass_{p_i} \cdot s_{p_i} + (1 - dead_{p_i}) \cdot nPed_{p_i} \cdot a_{p_i} - dead_{p_i} \cdot nPed_{p_i} \cdot cPed & act_{p_i} = 0 \\ (1 - dead_{p_i}) \cdot nPass_{p_i} \cdot s_{p_i} + nPed_{p_i} \cdot a_{p_i} & act_{p_i} = 1 \end{cases}$$

where  $dead_{p_i}$  is 0 if people survived, 1 otherwise.



The second component is computed based on the alternative action:

$$u_{alt}(p_i) = \begin{cases} nPass_{p_i} \cdot s_{p_i} \cdot (1 - probPass_{p_i}) + nPed_{p_i} \cdot a_{p_i} & act_{p_i} = 0 \\ nPass_{p_i} \cdot s_{p_i} + nPed_{p_i} \cdot a_{p_i} \cdot (1 - probPed_{p_i}) + \\ -nPed_{p_i} \cdot cPed \cdot probPed_{p_i} & act_{p_i} = 1 \end{cases}$$

Notice that single components are weighted using the likelihood of harming pedestrian/passengers in this case.

The reward depends on whether the AV's behaviour differs from the average behaviour of the community:

If the average individual would go straight and the AV turns, then the action is rewarded (having done an action that is meritorious, since it minimizes the risk of losses more than the average)

On the other hand, if the average individual would turn and the AV goes straight, then it is punished.

Loreggia si è fermato qui, non è andato oltre nelle slides....