# Challenges for Real Applications

Data Science and Engineering (I)
Master's Degree in Computer Engineering

UNIVERSIDAD DE MÁLAGA

E.T.S. INGENIERÍA INFORMÁTICA
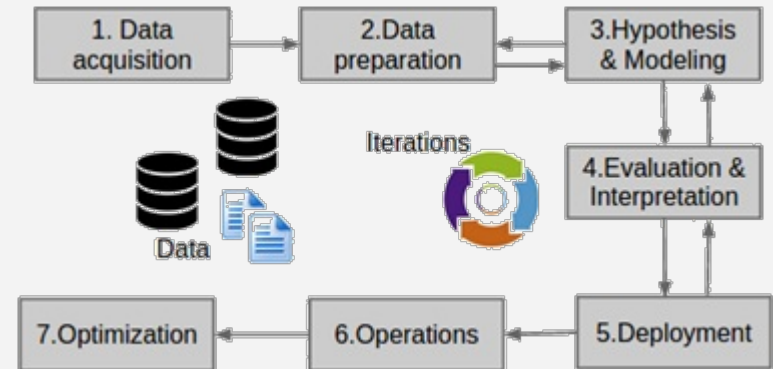UNIVERSIDAD DE MÁLAGA

# Challenges for real products

1. **Managing** a data science and engineering project

2. Search, **optimization**, and learning

3. **Needs of real projects**: scalability, dynamism, robustness, multiple objectives, restrictions, and self-control

4. **Examples** of real products and services

5. **Modern techniques** for real applications

# Management of a DS project (I)

- The typical data science project is an **engineering procedure**: start, steps, end
- Full of **informed decisions** on whether to continue based on pre-defined criteria
- **Goal**: optimize **resource** utilization, get high-quality **results** and maximize **benefits**
- **Money** is an issue, but **realistic** hypotheses and ideas are a must
- **The data science life-cycle**:
  1. Data acquisition
  2. Data preparation
  3. Hypothesis and modeling
  4. Evaluation & Interpretation
  5. Deployment
  6. Operations
  7. Optimization

# Management of a DS project (II)

1. **Data acquisition** – acquiring data from **internal** and **external** sources

2. **Data preparation** ("data wrangling") - involves cleaning the data and reshaping it into a readily **usable**

3. **Hypothesis and modeling** – **applying ML techniques** to all data (MS: model selection). MS involves to identify **training/test** sets

4. **Evaluation and interpretation** – comparing model performances

- *Steps 2-3-4 are repeated; as the understanding of data and business becomes clearer*

5. **Deployment** – the project is run in a production environment. It could include fast-tweaks *after* deployment, based on the **continuous deployment** model.

6. **Operations** (maintenance) – This phase could also follow a **DevOps** model which gels well with the continuous deployment model, given the rapid time-to-market requirements in big data projects. Steps 5 and 6 are mixed usually (like agile software in software engineering).

7. **Optimization** – This could be triggered by failing performance, or due to the need to add new data sources and retraining the model, …
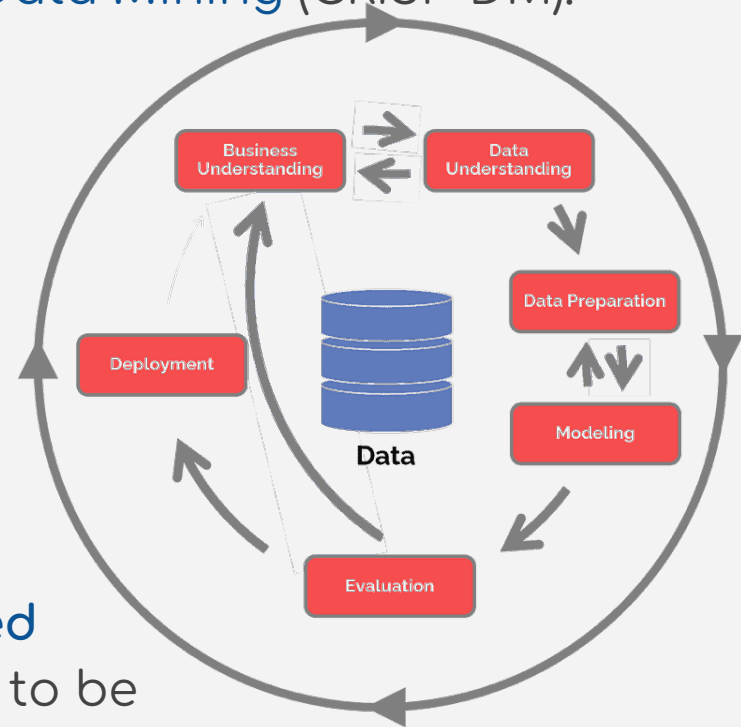
# Examples

In which step is the activity done? (1. Data acquisition, 2. Data preparation, 3. Hypothesis and modeling, 4. Evaluation & Interpretation, 5. Deployment, 6. Operations, and 7. Optimization)

A.- Removing outliers
B.- Calculating vehicle speed from its positions
C.- Running the complete system on a Docker infrastructure for testing
D.- Splitting the dataset into train and test set
E.- Deploy improved version of the model
F.- Rebooting the complete system after an unrecoverable failure
G.- Applying cross-validation
H.- Tuning the model parameters
I.- Examining which models can be applied to the data
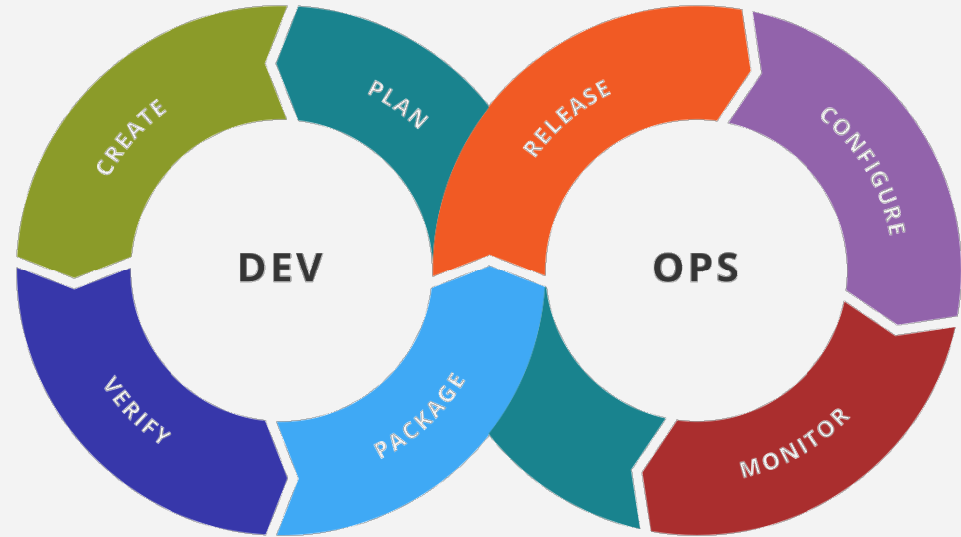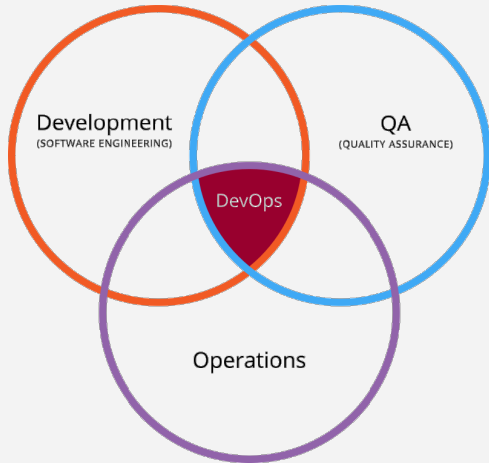J.- Obtaining values from car's sensors (OBD-II)

Management

Challenges for real Applications
Data Science and Engineering I (MUII)

- **CRoss Industry Standard Process for Data Mining** (CRISP-DM):
  1. Business understanding
  2. Data understanding
  3. Data preparation
  4. Modeling
  5. Evaluation
  6. Deployment
- **Library of assets** (expertise/maturity):
  1. Library of business use case
  2. Data requirements
  3. Minimum data quality requirements
  4. ...
- **Data scientists are likely to have limited business domain expertise**. They need to be paired with business people and those with expertise in understanding the data.

Business Understanding
Data Understanding
Data Preparation
Modeling
Evaluation
Deployment
Data

# You should now also on (II) DevOps

- Set of practices to **reduce the time** between committing a change to a system and the change being placed into normal production, while ensuring **quality**
- It uses different sets of tools (toolchains) rather than a single one
- **Steps** – Coding + Building + Testing + Packaging + Releasing + Configuring + Monitoring
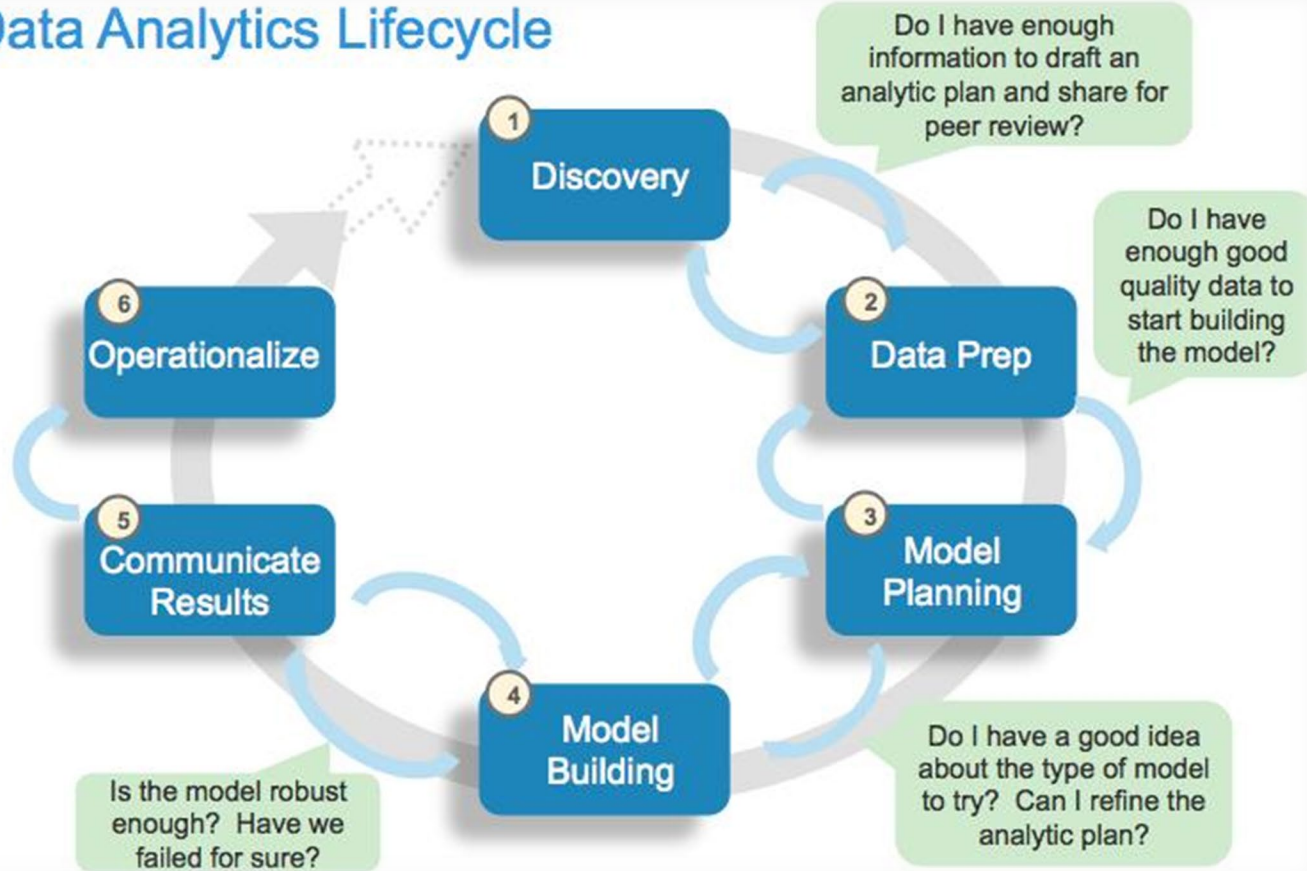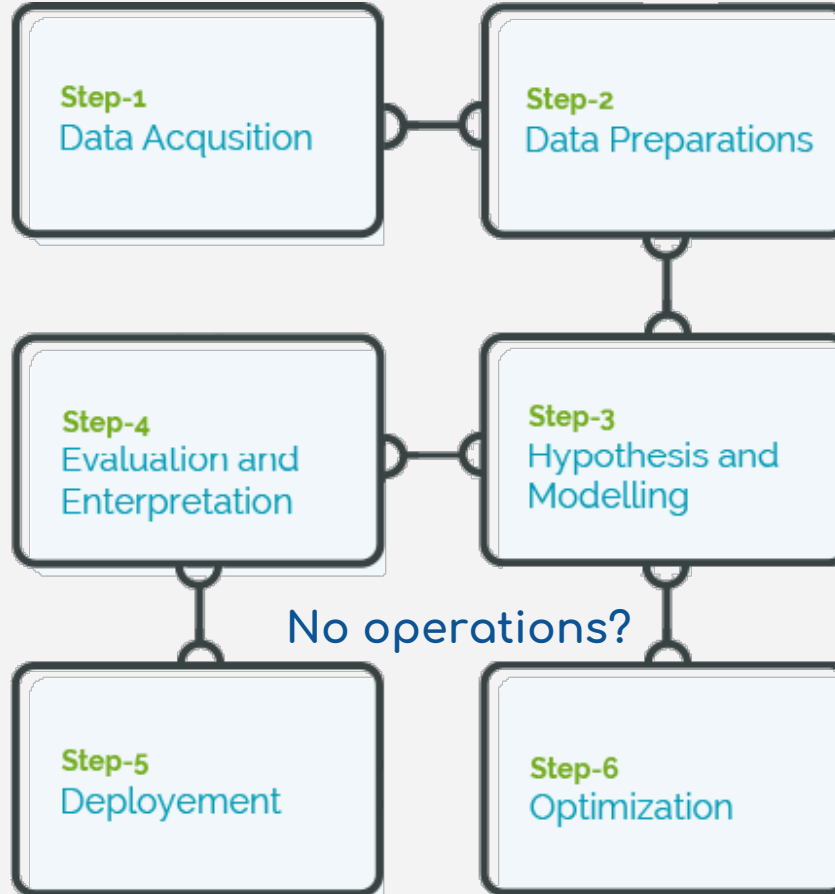
**DATA SCIENCE LIFECYCLE**

sudeep.co

01 BUSINESS UNDERSTANDING — Ask relevant questions and define objectives for the problem that needs to be tackled.

02 DATA MINING — Gather and scrape the data necessary for the project.

03 DATA CLEANING — Fix the inconsistencies within the data and handle the missing values.

04 DATA EXPLORATION — Form hypotheses about your defined problem by visually analyzing the data.

05 FEATURE ENGINEERING — Select important features and construct more meaningful ones using the raw data that you have.

06 PREDICTIVE MODELING — Train machine learning models, evaluate their performance, and use them to make predictions.

07 DATA VISUALIZATION — Communicate the findings with key stakeholders using plots and interactive visualizations.

Too new! Standards just born…

Business Discovery · Data Discovery · Data Preparation · Exploratory Analysis · Feature Engineering · Predictive Modeling · Production Deployment

Management

Challenges for real Applications
Data Science and Engineering I (MUII)

# Other models having six steps (II)

**Step-1**
Data Acqusition

**Step-2**
Data Preparations

**Step-4**
Evaluation and Enterpretation

**Step-3**
Hypothesis and Modelling

No operations?

**Step-5**
Deployement

**Step-6**
Optimization

It always starts with some data ...

| Data Preparation | Model Training | Model Optimization | Model Evaluation | Deployment |
|---|---|---|---|---|
| Data Manipulation | Model Training | Parameter Tuning | Performance Measures | Files & DBs |
| Data Blending | Bag of Models | Parameter Optimization | Accuracy | Dashboards |
| Missing Values Handling | Model Selection | Regularization | ROC Curve | REST API |
| Feature Generation | Ensemble Models | Model Size | Cross-Validation | SQL Code Export |
| Dimensionality Reduction | Own Ensemble Model | No Iterations | ... | Reporting |
| Feature Selection | External Models | ... | | ... |
| Outlier Removal | Import Existing Models | | | |
| Normalization | Model Factory | | | |
| Partitioning | ... | | | |
| ... | | | | |

Data Science Process

| OBTAIN | SCRUB | EXPLORE | MODEL | INTERPRET |
|--------|-------|---------|-------|-----------|
| O | S | E | M | N |
| Gather data from relevant sources | Clean data to formats that machine understands | Find significant patterns and trends using statistical methods | Construct models to predict and forecast | Put the results into good use |

# Models having a larger number of steps!

# ... OK; let's go wild!

# Optimization problems are everywhere!

- Logistics, transportation, supply change management
- Manufacturing, production lines
- Timetabling
- Cutting & packing
- Computer networks and telecommunications
- Health
- Videogames
- Software (SBSE)

... even in **data science** (data fitting, NN training, feature selection...)
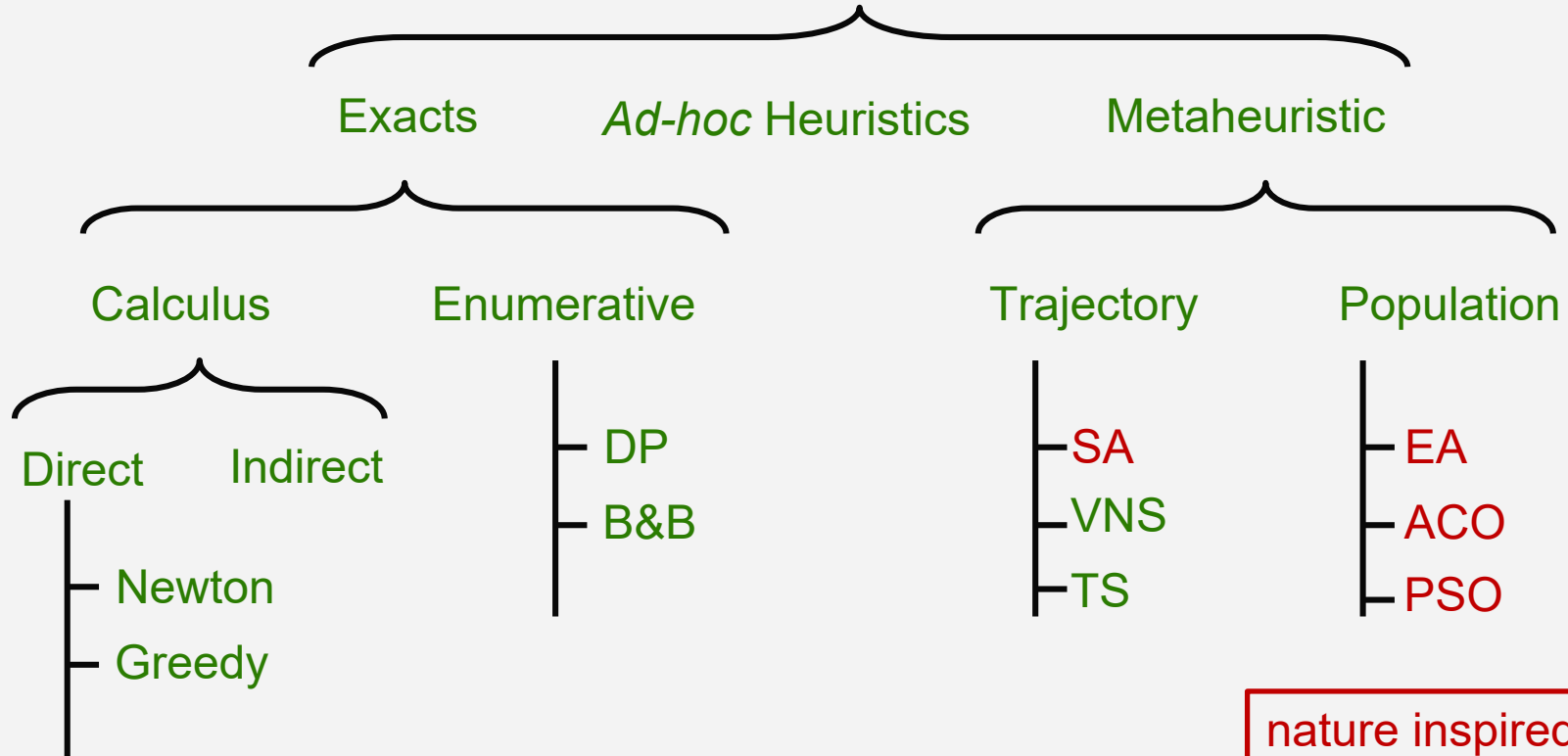
# Optimization problems

- Most general form:

$$min_{x \in X_{ad}} f(x)$$

Terminology:
- $f: X_{ad} \to R$: **fitness function**, objective function, usually real-valued
- $min \leftrightarrow max$ by replacement $f \leftrightarrow -f$
- $x$: **control** or **optimization parameters**
  - integer/discrete, continuous, or mixed-integer problems
- $X$: usually vector space or unbound set
- $X_{ad} \subset X$: **admissible** or **feasible set**
  - $X_{ad} = X$: **unconstrained problem**
  - $X_{ad} \neq X$: **constrained problem**

# A taxonomy of modern AI techniques

Optimization Algorithms

Exacts          *Ad-hoc* Heuristics          Metaheuristic

Calculus          Enumerative          Trajectory          Population

Direct     Indirect          DP          SA          EA

Newton          B&B          VNS          ACO

Greedy          TS          PSO

nature inspired in red

# Evolutionary Algorithm

- Based on the ideas of Darwinian Evolution theory

```
t := 0
initialize(P(t))
evaluate(P(t))
while  not  end condition  do
        P'(t)  := selection(P(t))
        P'(t)  := recombination(P'(t))
        P'(t)  := mutation(P'(t))
        evaluate(P'(t))
        P(t+1)  := replacement(P(t), P'(t))
        t  := t+1
end  while
```



In order to use one EA several steps of instantiation are needed:
- Problem: genotype (encoding) and fitness function
- Operators and their parameters
- Stopping criterion

# Evolutionary Algorithm: example (I)

- 0-1 Knapsack problem

?

$4  12 kg

$2  2 kg

15 kg

$2  1 kg

$1  1 kg

$10  4 kg

Maximize $\sum_{i=1}^{n} v_i x_i$

Subject to $\sum_{i=1}^{n} w_i x_i \leq W$ and $x_i \in \{0,1\}$

# Evolutionary Algorithm: example (II)

- **Genotype:** bit string

| 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

- **Fitness function**: $\sum_{i=1}^{n} v_i x_i$

- **Stop condition**: 100000 evaluations

- **Population**:
  - Size: 100
  - Random generated

| 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 |
| 2 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 |
| ... | | | | | | | | | | | | | | | |
| 99 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 |

- **Selection**: Random
- **Replacement**: Worst

- **Crossover / recombination:** SPX (simple point crossover)

| 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 |

| 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 |

| 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 |

| 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 |

- **Mutation:** bit-flip

| 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 |

| 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 |

**Unimodal**

Rastrigin function

**Multimodal**

**Deceptive**

**Population generation**

**Selection**

**Crossover**

**Population generation**

**Selection**

**Crossover**

**Mutation**

**Replacement**

## Current solution

## Neighbours

# Particle Swarm Optimization (I)

- Particle Swarm Optimization (PSO) is a population based metaheuristic inspired in the social behavior of **birds** within a flock
- It was initially designed for **continuous optimization** problems, but can be used in discrete ones also
- In PSO, each potential solution is called a **particle** and the population of particles is called a swarm
- In this algorithm, each **particle position** $\rho_i$ is updated each generation k by means of this equation ($v^i$ is its **velocity**):

$$p_i^{k+1} \leftarrow p_i^k + v_i^{k+1}$$

- The velocity of the particle is given by the expression:

$$v_i^{k+1} \leftarrow w \cdot v_i^k + c_1 \cdot r \cdot g_{best}^k + c_2 \cdot r \cdot p_i^{best}$$



**Algorithm 1** Pseudocode of PSO

```
 1:  initializeSwarm()
 2:  locateLeader(b)
 3:  while !stopCondition() or g < maxGenerations do
 4:      for each particle x_g^i do
 5:          updateVelocity(v_g^i) // Equation 2
 6:          updatePosition(p_g^i) // Equation 1
 7:          evaluate(p_g^i)
 8:          update(bp_g^i)
 9:      end for
10:      updateLeader(b_g)
11:  end while
```

# Simulated Annealing

- It is based on **annealing** in metallurgy
- SA is a **hill-climbing** method
- It **accepts worse solutions** to avoid getting stuck in local optima, according a criterion (y is new solution and x the old one):

$$rand(0,1) \leq \min\left(1, e^{\frac{f(x)-f(y)}{T}}\right)$$



**Algorithm 1** Simulated annealing algorithm

```
1:  procedure SA(f, N, Ω, x⁰, T₀)
2:      k ← 0
3:      x_min ← xᵏ
4:      f_min ← f(x_min)
5:      T_k ← T₀
6:      while stopping criterion is not satisfied do
7:          zᵏ ← rand (N(xᵏ, T_k))
8:          yᵏ ← xᵏ + zᵏ
9:          if rand(0,1) ≤ min{1, exp{(f(xᵏ) − f(yᵏ))/T_k}} then
10:             xᵏ⁺¹ ← yᵏ
11:         else
12:             xᵏ⁺¹ ← xᵏ
13:         if f(xᵏ⁺¹) < f_min then
14:             x_min ← xᵏ⁺¹
15:             f_min ← f(x_min)
16:         k ← k + 1
17:         T_{k+1} ← temperature is updated
18:     return x_min
```

- VNS is a stochastic algorithm with a **set of neighbourhood structures** are defined,
- Each iteration: **shaking**, **local search** and **move**
- VNS **explores** a set of **neighbourhoods** to get different local optima and **escape from local optima**

```
Procedure Algorithm of variable neighborhood search:
begin
    find the best solution found  x;
    k := 1
    while (k ≤ k_max) do
        randomly generate a new solution  y ∈ N_k(x);
        if (f(x) > f(y)) then
            x := y;
            k := 1;
        else
            k := k + 1;
        endif
    endwhile
end
```

# Real problems pose real challenges

## Reality is challenging:

- **Large scale**, every is really big
- **Time** consuming and real time
- **Dynamic**, everything changes in time
- **Uncertainty** in all tasks and phases
- **Complex** relations, interdependences
- Several **goals** at the same time
- **Human** preferences and interfaces
- Lots of **restrictions** (legal, technical…)
- **Mobile** plus **desktop** applications



WELCOME TO THE REAL WORLD

**Scalability** is the property of a system to handle a growing amount of work by adding resources to the system

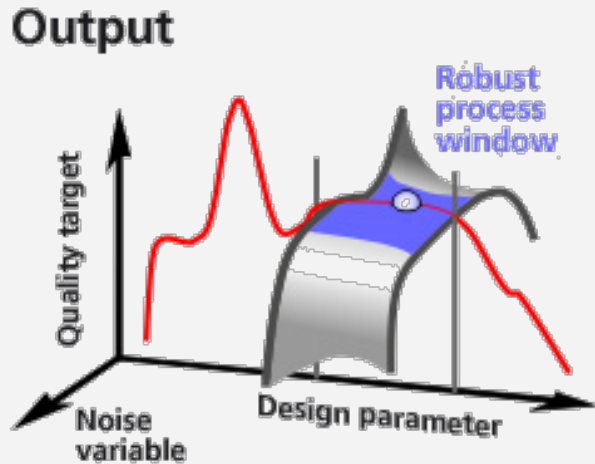**Dynamism:** the problem conditions change over the time in an unpredictable way

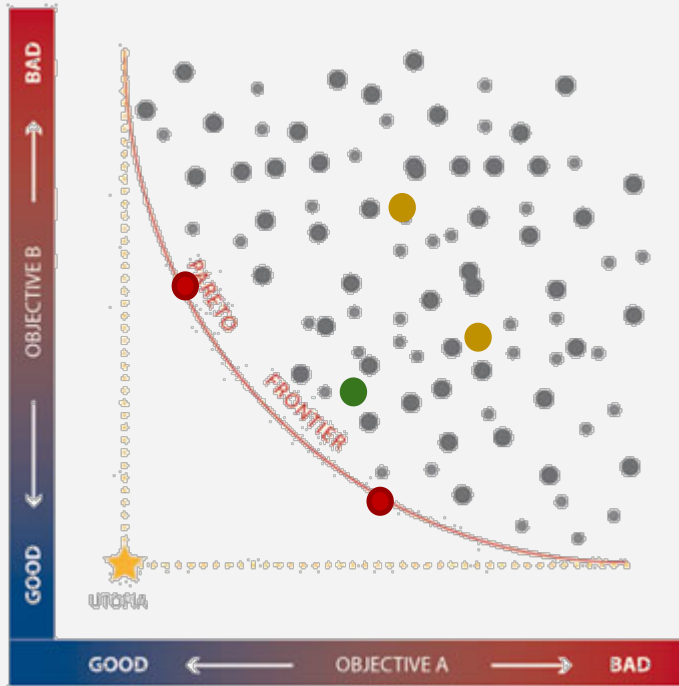**Uncertainty:** the problem involves imperfect or unkown information

# Robustness

**Robustness:** the performance is stable after adding some noise to the environment

**Green** dominates **yellow**. **Red** are non-dominated.
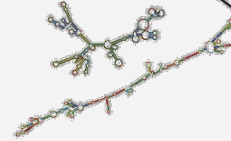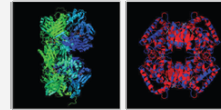
# Constraints

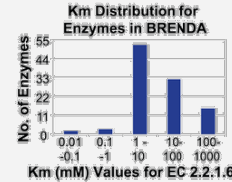## … many types of constrains …



Project Constraints

### Types of Data

A. Gene Expression

B. Protein Expression

C. Metabolite Concentration

D. Kinetic Parameters

Km Distribution for Enzymes in BRENDA

Km (mM) Values for EC 2.2.1.6

### Types of Constraints

Flux Capacity- Boolean (On/Off)

$$y_j \cdot v_j^{min} \leq v_j \leq y_j \cdot v_j^{max}$$
$$\text{where} \quad y_j = \{0,1\}$$

Flux Capacity- Continuous

$$p_j \cdot v_j^{min} \leq v_j \leq p_j \cdot v_j^{max}$$
$$\text{where} \quad p_j = [0,1]$$

Thermodynamic Constraints

if $v_j \geq 0$ then $\Delta G_j \leq 0$
if $v_j \leq 0$ then $\Delta G_j \geq 0$

$$\Delta G_j = \Delta G_j^\circ + RT \sum_i S_{i,j} \ln C_i$$

Molecular Crowding Constraints

$$\sum_j w_j \cdot v_j \leq 1$$

Kinetic Constraints

$$v_j = \frac{k_{cat,j} \cdot C_i}{k_{m,j} + C_i} \quad \text{(Biochemical)}$$

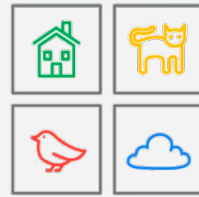$$v_j = k_j \prod_i C_i^{S_{ij}} \quad \text{(Mass Action)}$$

# AutoML
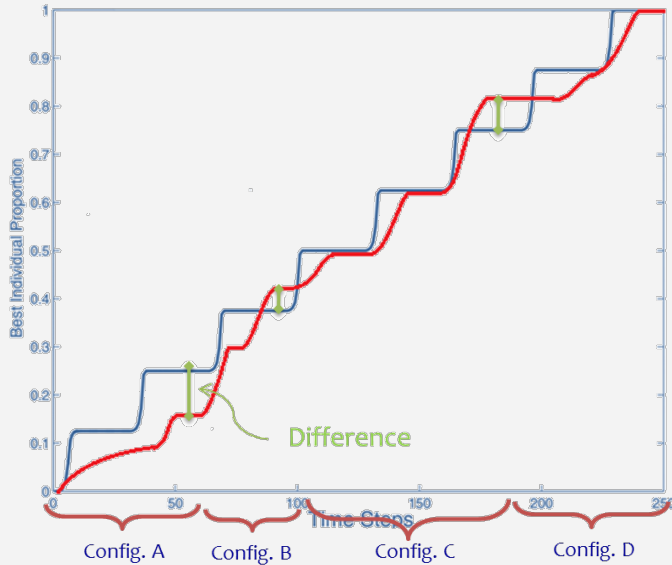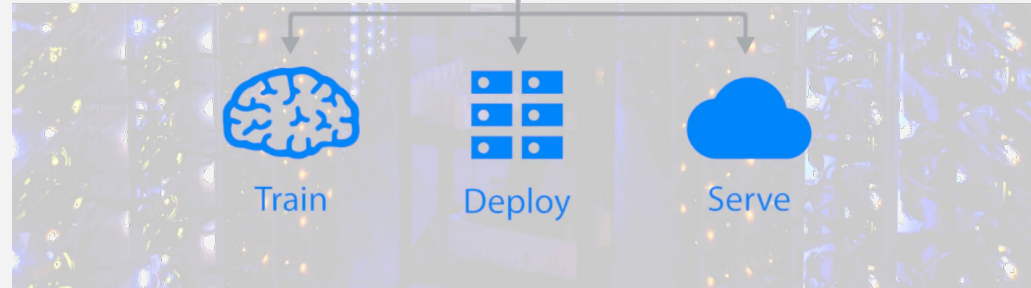


Photo Dataset → Cloud AutoML Vision → Rest API

Generate predictions with a REST API

Train    Deploy    Serve

Difference

Best Individual Proportion

Time Steps

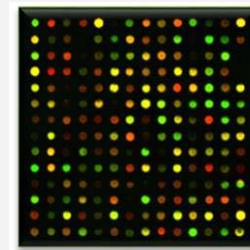Config. A    Config. B    Config. C    Config. D

# An example (I)

**Problem:**
- **Gene selection** and cancer classification of DNA
- Microarray, feature selection

**Objectives:**
- Maximize accuracy of prediction
- Minimize the number of selected genes
- Maximize sensibility and specificity ( ROC factors )

**Phases:**
- Feature selection
- Training
- Validation
- Fitness calculation



Solution (S). Provided by
Metaheuristic (PSO, GA)

Microarray
Initial Dataset

1    N

1 0 1 1 ... 0 0 1

Subset

Trainig
set

Test set

Subset Evaluation

SVM – Classification
+
10-Fold Cross
Validation

Fitness(S)=Accuracy &
Subset Size

**Fitness:**

- **Monobjective:** aggregative: (alpha*100/accuracy + beta * #features )
- **Multiobjective:**
    - 2 (accuracy, #features) or 3 (sensibility, specificity, #features)
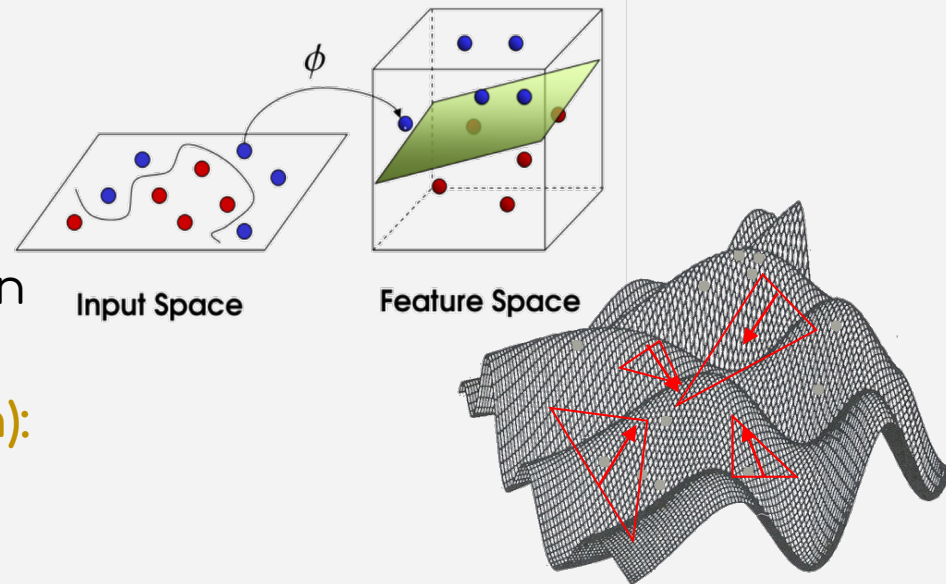
**Classification:**

- SVM
- K-means

**Validation:**

- Leave one out cross-validation
- 10-fold cross-validation

**Algorithms (for feature selection):**

- PSO variant
- GA variant



Input Space          Feature Space

Example

Challenges for real Applications
Data Science and Engineering I (MUII)

## Instances:

- Large scale datasets of well-known cancer DNA Microarrays: Leukemia, Colon, Prostate, Lung, Ovarian, Breast (e.g. breast 24481 genes and 97 patient samples)

Results: comparison against other techniques (S.O.T.A.)

| Dataset | GPSO | GA | Huerta et al. | Juliusdotir et al. | Deb et al. | Guyon et al. | Yu et al. | Liu et al. | Shen et al. |
|---------|------|-----|-----|------|------|------|------|------|------|
| *Leukemia* | 97.38(3) | 97.27(4) | 100(25) | - | 100(4) | **100(2)** | 87.44(4) | - | - |
| *Breast* | 86.35(4) | **95.86(4)** | - | - | - | - | 79.38(67) | - | - |
| *Colon* | **100(2)** | 100(3) | 99.41(10) | 94.12(37) | 97(7) | 98(4) | 93.55(4) | 85.48(-) | 94(4) |
| *Lung* | 99.00(4) | **99.49(4)** | - | - | - | - | 98.34(6) | - | - |
| *Ovarian* | **99.44(4)** | 98.83(4) | - | - | - | - | - | 99.21(75) | - |
| *Prostate* | **98.66(4)** | 98.65(4) | - | 88.88(20) | - | - | - | - | - |

## Leukemia Gene Subset:

**PSO**: K01383,  U03056, J04130 vs **GA**: L40379, S85963, U83192, Z49099

# Advanced Tools

Data Science and Engineering (I)
Master's Degree in Computer Engineering

UNIVERSIDAD
DE MÁLAGA

E.T.S. INGENIERÍA
INFORMÁTICA
UNIVERSIDAD DE MÁLAGA

# Advanced techniques and technologies

1. **Complex** problems need advanced tools

2. **Measuring** efficacy and efficiency

3. **Parallel hardware**, or how new technology helps

4. **Algorithm hybridization**, or how new techniques can help

5. **Practical Examples**

# Real problems need more than you expect / know

- **Graduate** students know some tools to deal with engineering apps
- Most graduate programs offer a **small sample** of algorithms and technologies
- Graduate students then only know **very basic concepts**
- **Real** problems seldom admit the **constraints** of basic tools
- A complex real application needs **advanced** algorithms and technologies
- **Research** in algorithms, software, AI, and new technologies is full of them
- Just **few techniques** that can be used as described in books
- To work well, they need to be improved...

## How do we improve them?

# Important questions

## What is the computational complexity of your algorithm?

Measure it as $O(n)$ $O(\log n)$ $O(n \cdot \log n)$ $O(n^2)$ $O(n^3)$ ... $O(2^n)$ ... $O(n!)$ ... $O(n^n)$ ...

## How much simple is a technique? Occam's razor principle applies

If more complex than one with similar behaviour, then not interesting

## How measure complexity: computational, software, understanding...?

In terms of input, branches, length of description, time to learn it, ...

## What is defining the limits of a technique or a technology?

Its complexity, but also its accuracy in solving a problem, its robustness...

## How a technique could be improved? And a technology?

New design (operations, concepts), new implementation, latest hardware, ...

## Similarities to other existing tools? Do they inspire to improve?

See the basics of the tool, similar structure, know on cross-fertilization

## Can we quantify all the decisions? Identify all the needed steps?

Data driven decision making, always measure ... scientific method

# Always measure: efficacy measures

## MSE – Mean Squared Error:
- Risk measure (quality of estimator)

$$MSE = \frac{1}{n}\sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2$$

## Logarithmic Loss:
- Penalising the false classification

$$LogLoss = \frac{-1}{N}\sum_{i=1}^{N}\sum_{j=1}^{M} y_{ij} * \log(p_{ij})$$

## Accuracy

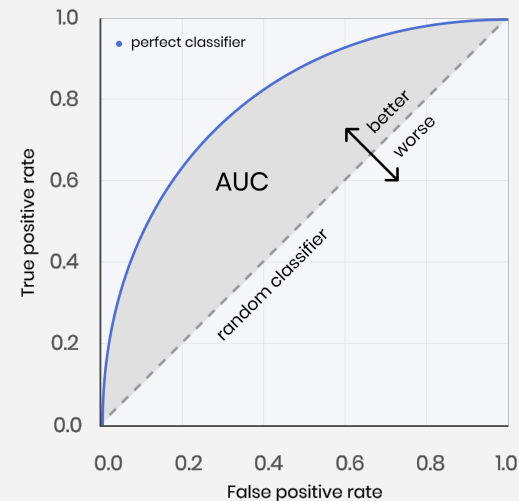$$Accuracy = \frac{\#\ correct\ predictions}{\#\ input\ samples}$$

## Confusion Matrix



## Area under de curve (AUC)
- How well the test separates the group being tested into 2 classes?
- TPR = TP/(TP + FN)      FPR = FP/(FP + TN)

## Wall Clock Time

- $T = t_{end} - t_{start}$



User time
CPU time
Communication time
...

Battery consumption (phone)
Kwh consumption (data center)



## Speedup

$$S_m = \frac{\bar{T}_1}{\bar{T}_m}$$

Memory usage
Size of data files

# Statistical analysis mandatory

- Even very advanced algorithms reach a **maximum efficiency**. This happens in large problem instances, or when using simulators, or in real time scenarios, or in web services for clients, …
- Advances in **parallel hardware** like clusters, multicores, GPUs, cloud, etc. allow to make more than one step per unit time in the used techniques
- Sometimes you are not only looking for **reduced times**, but for **new types** of techniques that search for different solutions at the same time collaborating
- Sometimes you have a multicore or a lab full of cores: **explode them** !!!
- Thus, you can make new techniques and also run them faster, **both** !!!

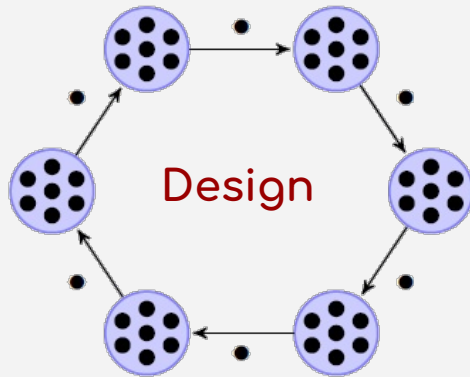**Problems not solved before, now become solvable by using parallel tools**

## Parallelism and Metaheuristics:

The increasing availability of new kinds of CPUs and the parallel nature of metaheuristics have allowed the fast development of parallel metaheuristics
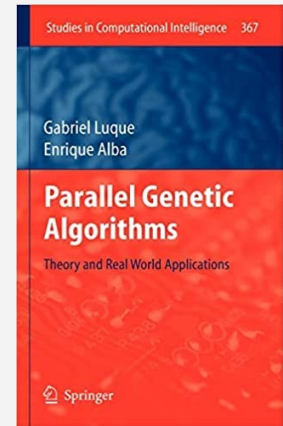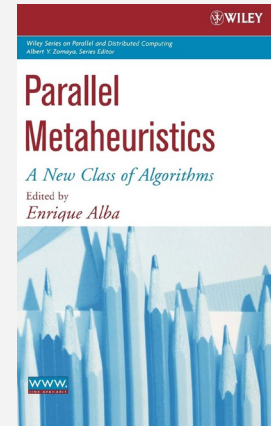
## Advantages:

- Allow to tackle more complex problems/instances
- Allow to reduce the execution time
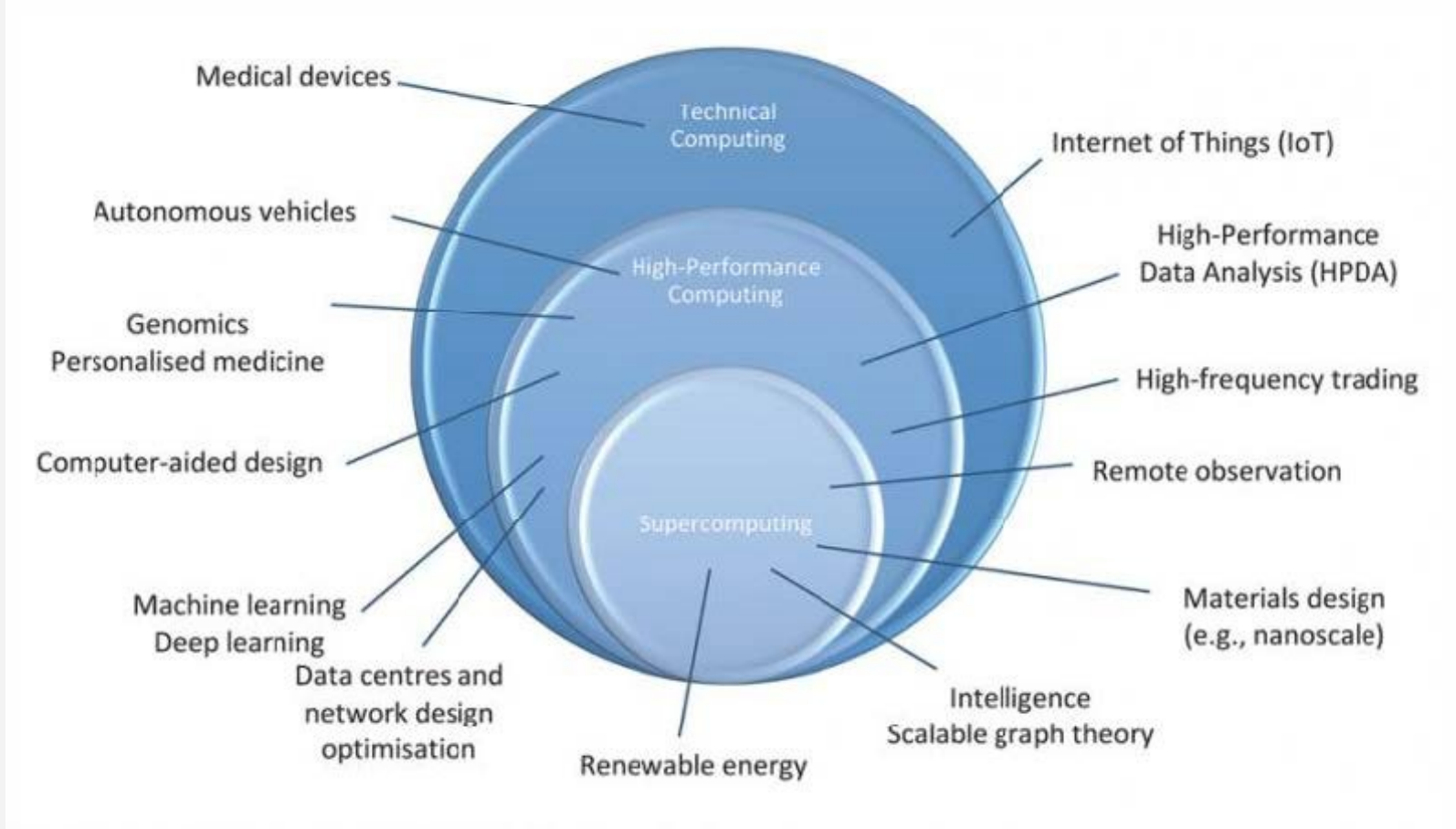- Allow to improve quality of the found solutions

Design

Implementations

S.O.L.

Challenges for real Applications
Data Science and Engineering I (MUII)

# Hardware is important

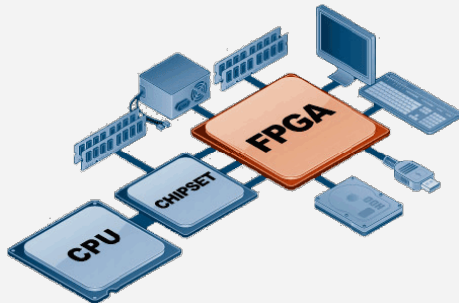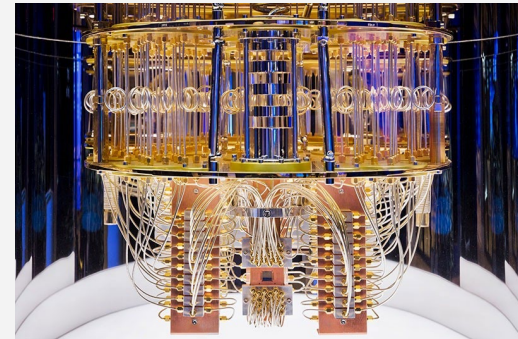## Grid and cloud computing

## Cluster computing

## GPU

## FPGA

## Manycores

## Quantum Computers

# Design vs implementation, not the same

- Node in a distributed EA

```
t := 0
initialize(P(t))
evaluate(P(t))
while  not end condition do
        P'(t) := selection(P(t))
        P'(t) := recombination(P'(t))
        P'(t) := mutation(P'(t))
        evaluate(P'(t))
        P(t+1) := replacement(P(t), P'(t))
        <<Communication  with  neighbours  >>
        t := t+1
end while
```
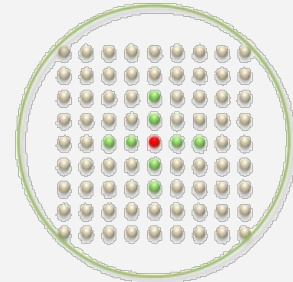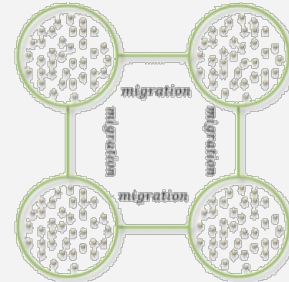


Master

Individuals

Fitness

Slaves

migration

migration

migration

migration

# Hybridization: a good way to build techniques

**Hybridization** is the inclusion of problem-dependent information in the algorithm, but also combining fields, operations, data, technologies, frameworks ...

## Types
- **Weak**
- **Strong**

# Hybridizing ML with metaheuristics

Machine Learning (ML) applications

Enhancing Metaheuristics

Improving ML techniques

Metaheuristics applications

**Specifically-located hybridizations**

- Parameter fine-tuning
- Initialization
- Evaluation
- Population management
- Operators
- Local search

**Global hybridizations**

- Reduction of search space
- Algorithm selection
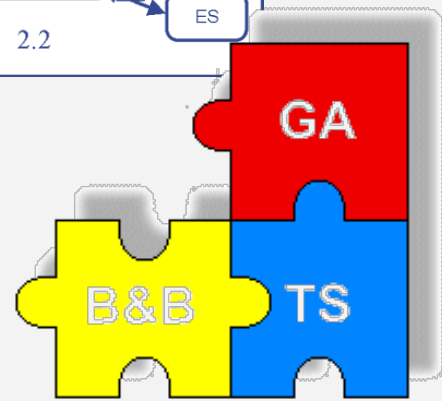- Hyperheuristics
- Cooperative strategies
- New types of metaheuristics

- Classification
- Regression
- Clustering
- Rule mining

**Optimization + Learning**

**Hybrids = ML + MetaH**

**Tuning ML tools**

**Surrogate systems**

**Collaboration processes**

# Academic problem domains

## Mathematical Optimization:

- Rastrigin, Rosenbrock, Mishra's Bird...

## Combinatorial optimization:

- Routes, scheduling, graphs...

## Domain dependent benchmarks:

- Multiobjective
- Temporal series
- Data mining
- Neuronal network training

Unscheduled tasks

Scheduled tasks

TIME

## Know on standard benchmarking!!!

# Sectoral domains: Telecoms in this case

## Radio Network Design



## GSM Frequency Assignment



## MANETs



## Sensor Network Layout



## Location Area in 4G/5G



## VANETS

# Every single domain is a good target

Designing Quantum Circuits

Data Based (Data Mining, Query Optimization)

Dynamic Optimization Problems (DOPs)

Tasks Scheduling in Operating Systems

Genomics (Fragment assembly, protein structure)

Games

# Parallelism, Hybridization, and a real application

**MICROARRAY**
(Scanned Data)

QUANTIFICATION DATA

**GENE EXPRESSION DATA MATRIX**

Samples

Gene expression levels

Genes

**Gene subset selection**
A solution (particle) looking for maximizing the classification with a small number of genes

Supervised classification
(SVM, KNN, HNN, ...)

A particle is a solution codifying a gene subset

**PARTICLE SWARM OPTIMIZATION**
• Optimization algorithm modeling bird flocking self-organizing behavior
• Embedding feature (gene) selection mechanism

New particle position:
a new subset of genes has been produced by PSO

Best-known position of current particle

Global best particle in the swarm

For each subGPSO *Si* in *{1,...,m}*
<do in parallel>
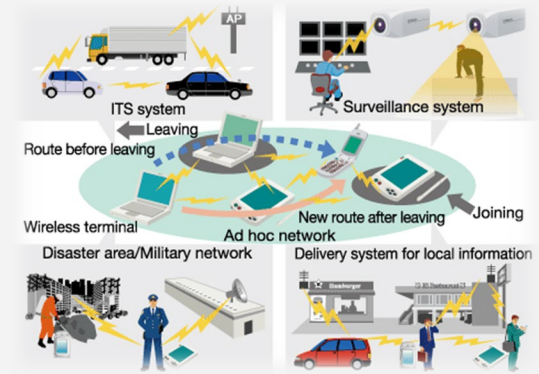
... *Sm*

*Si+1*

*Si*

Microarray dataset

Data preparation

*Initialize(Si)*

For each particle *j* of *Si*

Gene selection *(j)* and training SVM

Training dataset

*10 fold cross-validation*
Compute *fitness(j)*

*Three parent mask-based crossover (3PMBCX) and mutation*

Validation dataset

Migration?

*Yes*

Migrate *j* to *Si+1 from Si-1*

Test dataset

Compute new swarm *S'i*

Stop?

*No*

Best particle test evaluation

# Real Use Cases

Data Science and Engineering (I)
Master's Degree in Computer Engineering

UNIVERSIDAD
DE MÁLAGA

E.T.S. INGENIERÍA
INFORMÁTICA
UNIVERSIDAD DE MÁLAGA

# Learn more in going for real applications

1. Feature selection

2. Neuroevolution

3. Surrogate models

4. Real applications

# Feature selection

- The computational ability of machine learning models **depends** a lot on the **feature set**.
- Retaining the significant features vastly **improves** the learning **time**, and also improves **accuracy**.
- In **feature selection**, we find the **optimal feature subset** that contributes most to our predicted variable.
- **Advantages**:
  - **Improve generalization** of models by reducing overfitting of data.
  - **Remove** unnecessary/redundant data.
  - Curtail the **Curse Of Dimensionality**
  - **Optimize** training time

Feature Selection



Full Feature Set

Identify Useful Features

Selected Feature Set

# Feature selection as optimization problem

- **Solution encoding:**
  - The solution can be implemented as a bit string.
  - The **solution's length** is taken as the **number of features** in the dataset.
  - **0/1** indicates the **presence/absence** of the ith **feature** in the solution.

| 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

- **Fitness function:**
  - **Number of features** selected (nof)
  - **Model accuracy** (acc)
- **Approaches:**
  - Multiobjetive
  - Aggregate function:

$$max\,f = \alpha \cdot \frac{MaxFeatures - nof}{MaxFeatures} + (1 - \alpha) \cdot acc$$

I. Forward-propagate Input Signal

$$a_k = g_k(b_k + \Sigma_j g_j(b_j + \Sigma_i a_i w_{ij})w_{jk})$$

II. Back-propagate Error Signals

$$E'(a_k, t_k)$$

$$\delta_k = g_k'(z_k)E'(a_k, t_k)$$
$$\delta_j = g_j'(z_j)\Sigma_k \delta_k w_{jk}$$

III. Calculate Parameter Gradients

$$E'(a_k, t_k)$$

$$\partial E/\partial w_{ij} = a_i \delta_j$$
$$\partial E/\partial w_{jk} = a_j \delta_k$$

IV. Update Parameters

$$w_{ij} = w_{ij} - \eta(\partial E/\partial w_{ij})$$
$$w_{jk} = w_{jk} - \eta(\partial E/\partial w_{jk})$$
for learning rate $\eta$

S.O.L.

**Challenges for real Applications**
Data Science and Engineering I (MUII)

$$y = \sigma\left(w_0 + \sum_i w_i h_i\right) \quad h_i = \sigma\left(v_{i0} + \sum_j v_{ij} g_j\right) \quad g_j = \sigma\left(u_{j0} + \sum_k u_{jk} f_k\right) \quad f_k = \sigma\left(t_{k0} + \sum_m t_{km} x_m\right)$$

$$E = \tfrac{1}{2}\left(y - y^*\right)^2$$

$$\frac{\partial E}{\partial h_i} = \left(y - y^*\right) y(1 - y) w_i$$

$$\frac{\partial E}{\partial g_j} = \left(y - y^*\right) y(1 - y) \sum_i w_i h_i (1 - h_i) v_{ij}$$

- Prone to local optima
- Not appropriate for complex NN
- Oscillations
- Depends on structural functions
- Need unfolding in deep learning

**INPUT ENCODING**

(a-n)...(m-n)(a-o)...(m-o)(a-p)..(m-p)(n-q)(o-q)(p-q)...(n-z)(o-z)(p-z)

NEURAL NETWORK TRAINING

**Input Layer**

a  b  c  •••  m

**Hidden Layer**

n  o  p

**Output Layer**

q  r  s  •••  z

**INPUT-OUTPUT ENCODING**

(a-n)...(m-n)(n-q)...(n-z)(a-o)..(m-o)(o-q)...(o-z)(a-p)...(m-p)(p-q)...(p-z)

- **Solution encoding:**
  - **List** of **float** numbers
  - **Mapping** between **poisition** in list and **weight** in NN

- **Fitness function:**
  - Model accuracy

# Evolving NN structures!!!



(i) Mixed-coding scheme

1: represents active neuron

0: represents inactive neuron

(ii) The configurations of three-layer feedforward ANN

**genome**

gene: #1 | 000001001

gene: #2 | 001100001

gene: #3 | 000100100

gene: #4 | 000000011

**000000000**

↳ *Trainable Flag*

↳ *Activation Function*
• 00: tanh
• 01: sigmoid
• ...

↳ *Unit Size*
• 0000:   1
• 0001: 16
• 0010: 32
• ...

↳ *Layer Type*
• 00: Dense
• 01: LSTM
• ...

**Genotype**

**Translation
Faculty**

**Phenotype**

# Surrogate models

- A lot of engineering **problems** require experiments and/or **simulations** to evaluate design objective and constraint functions as a function of design variables.
- A single **simulation** can take many minutes, hours, or even days to complete, thus rendering them infeasible in practice.

**Surrogate models** are a statistical model to accurately approximate the simulation output.

This **trained model** can be deployed to replace the original computer simulation.



Design parameters

$X_1$

$X_2$

$\vdots$

$X_N$

- *Sensitivity analysis*
- *Optimization*
- *Risk analysis*

Simulations

$y = f(X)$

$(x^1, f(x^1))$
$(x^2, f(x^2))$
$\vdots$

Training

Surrogate model

$y \approx \hat{f}(X)$

Output

$y_1$

$y_2$

$\vdots$

$y_N$

Expensive, since it involves many simulation runs

Cheap, since training and employing a surrogate model is not expensive

# Surrogate models



**Sampling:**
- Random
- Latin hypercube

**Construct model:**
- Model selection
- Tuning parameters
- Feature selection?

**Utilization:**
- Only surrogate model
- Mixed model: surrogate + simulation
- New samples?

**Introduction**
- Problem
- Approaches
- Proposal

**Traffic Control Signals**
- Concepts
- Current systems
- Proposed automatic system

**Solving**
- System modeling
- Realistic data
- Surrogate model

# Problem

- **City evolution:**
  - Nowadays, cities are growing in the number of inhabitants, many of whom are arriving at the city for the first time
  - By 2050 the human population will reach 9 billion with **75%** of the world's inhabitants living in towns and cities

- As consequence, the **number of vehicles** in streets is continuously **increasing**, affecting all aspects of daily life:
  - Traffic jams
  - Pollution
  - Security
  - Stress
  - Economic losses
  - …

# Potential Solutions



- "Classic" solutions:
  - Infrastructures
  - Promote the use of **public transportation** or green vehicles (bikes)
  - Promote the use of **car-sharing** (VAO lanes)
  - **Limiting car** access to city centers
- "Intelligent" solutions:
  - Provide **real-time and accurate data** to citizens to make **informed decisions** (traffic intensity, free park slots, …)
  - Automatic assistance tools: **adaptive and/or customized routes**
  - Better **tuning** of **existing elements**: routes and frequencies of public transportation, traffic light timing…

# Proposal

## Automatic Traffic Control Signals

- Reduce the traffic jams
- Minimize the waiting times in red lights
- Faster routes
- Reduce the gas emissions

# Traffic Control Signals

- First, we need to **study** the elements, constraints, and regulations in **the problem domain**

- Multiple sources of information:
  - **International regulations**. (I.e., U.S. Transport Department):
    - Manual on Uniform Traffic Control Devices (862 pages)
    - Traffic Signal Timing Manual (274 pages)
  - **National regulations**. (I.e., DGT):
    - Regulación semafórica (32 pages)
    - Cruces semafóricos y sincronismo (32 pages)
  - **Specialized personnel** (city traffic managers)
  - **Scientific literature**



- Information filtering

# Traffic Control Signals

- **Important concepts**:
  - Intersections
  - Cycle
  - Phases
  - Traffic light schedule or plan





Intersection

# Traffic Control Signals

- More information:

  - The **duration of phases and cycles** can be modified

  - The **phases CANNOT** be modified

  - Recommended **duration of a cycle**: 60-120 seconds

  - **Yellow phases** (before red light): 4 seconds

  - **Minimum duration** of some phases (i.e., red phases at crosswalks should allow to safely cross the road => minimum duration = 1 m/s * 4 m/lane * #lanes)

  - Promote **green waves** in important avenues

  - **Traffic-dependent planes** (time, weekday, season,... )

  - ...



Green wave

# How are traffic lights configured?

- **Location/type:**
  - Some locations are mandatory according to regulations
  - Other locations are recommended but not mandatory



- **Phases of traffic lights:**
  - Regulations establish a procedure to set the phases

- **Duration of phases/cycles:**
  - It is defined by city traffic managers according to some constraints
  - Usually, it is manually done in each intersection
  - Based on experience and accumulated knowledge
  - There are dynamic systems (they react to current traffic). Problem: Quick changes that don't improve traffic

CRUCE 1.01

c/Vicente Espinel

Avda. Juan Sebastián Elcano

1004
1005

Bus

c/Pez de Plata 1008

fase 1 y 2

fase 4

fase 3

- 3 traffic flows
- 6 traffic lights:
  - 4 for vehicles (1 for bus)
  - 2 for citizens

# How are traffic lights configured?

- **...ns** that are chosen



**Cruce:** 01010 **PLAN:** Actual: 003

**Descripción:** Juan Sebastián Elcano - c/ Vicente Espinel

**Comentario:**

06-jun-2016 14:31:43

| SELECCIÓN HORARIA - VERANO 2015 | | | |
|---|---|---|---|
| **SUB 1 LABORABLE** | | | |
| Hora | Plan | Ciclo | Sentido |
| 0:00 | 12 | 80 | Oeste |
| 2:00 | 13 | 70 | Oeste |
| 5:00 | 12 | 80 | Oeste |
| 6:30 | 20 | 115 | Oeste-Simultáneo (118 a 115) |
| 11:00 | 5 | 110 | Oeste |
| 16:00 | 2 | 115 | Oeste |
| 20:30 | 5 | 110 | Oeste |
| 22:30 | 7 | 100 | Oeste |
| **SUB 1 SABADO** | | | |
| Hora | Plan | Ciclo | Sentido |
| 0:00 | 7 | 100 | Oeste |
| 3:00 | 12 | 80 | Oeste |
| 10:00 | 5 | 110 | Oeste |
| 14:00 | 7 | 100 | Oeste |
| 17:30 | 5 | 110 | Oeste |
| 22:30 | 7 | 100 | Oeste |
| **SUB 1 DOMINGO** | | | |
| Hora | Plan | Ciclo | Sentido |
| 0:00 | 7 | 100 | Oeste |
| 3:00 | 12 | 80 | Oeste |
| 10:00 | 11 | 90 | Oeste |
| 11:30 | 10 | 95 | Oeste |
| 16:00 | 7 | 100 | Oeste |
| 17:00 | 20 | 115 | Oeste |
| 22:30 | 7 | 100 | Oeste |

# Proposed system

- **Automatically generated cycle and phase times**:

  - Those times must respect the constraints

- Simultaneously consider **all traffic lights in the city** or the area defined by the traffic control center

- Obtain **different plans** (offline) according to traffic intensity

- The final goal is to obtain **more fluid traffic** that reduces the pollution
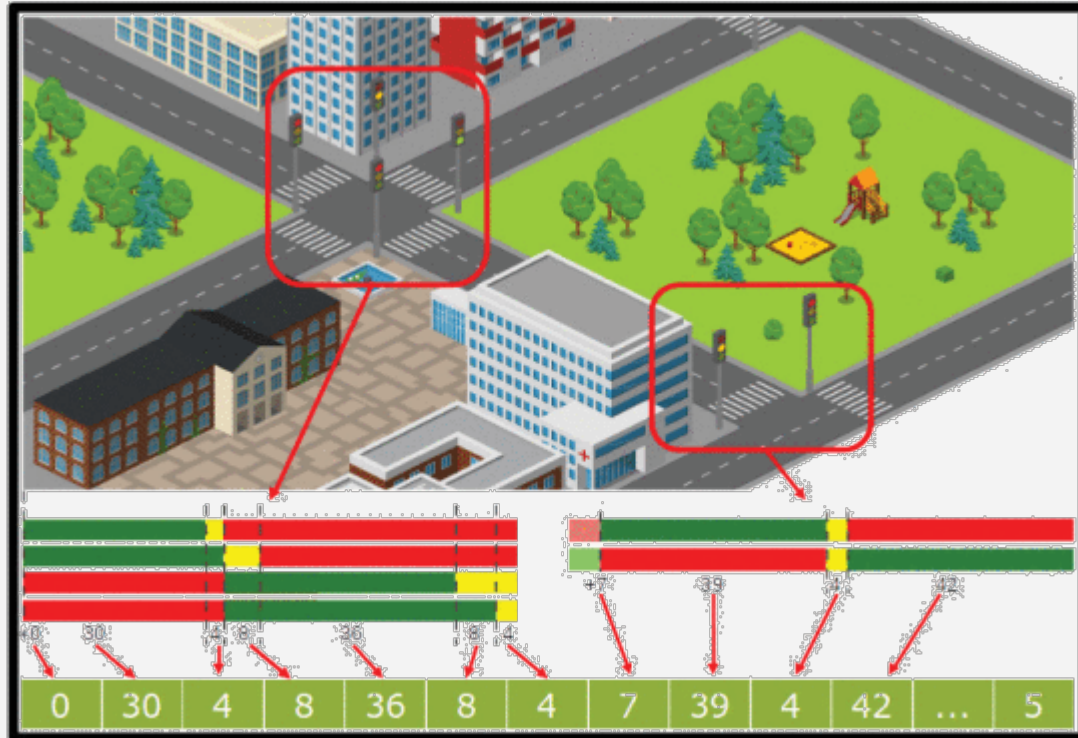
# Modeling the problem

- Given the following **input data**:
  - **Intersections** to improve (location, phases)
  - And the **traffic flows**

- **Objective**: to find the configuration (duration for the phases) that outperforms the rest of the existing configurations

- **Questions**:
  - What is computationally a solution (**representation**)?
  - When is a configuration **better** than another one (from a numerical point of view)?
  - Where do we get the **information** from?
  - How do we **find the best**?

Example

Challenges for real Applications
Data Science and Engineering I (MUII)

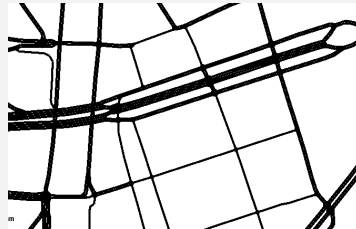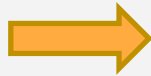# Solution encoding

- List of numbers (duration of phases):

# Solution encoding

- **Restrictions** not met by default:

  - Yellow phases to 4
  - Cycles > 60 and < 120
  - Green Wave Promotion

- Other **alternative representations**:

  - Each intersection: cycle time + percentage that each phase occupies
  - Reduce the number of traffic lights:
    - Cluster the intersections into groups
    - Only one in the group is optimized
    - The rest are small variants of the optimized
  - Others?

# Solution fitness

- To compare solutions, we need quantitative values (fitness)
- The traffic system is very complex

  - No (realistic) mathematical models

  - Utilization of simulators
- SUMO (Simulator of Urban Mobility)
  - Input: Roadmap, traffic flows, traffic light plans
  - Output: statistics of the simulation

| 20 | 8 | 47 | 23 | 30 | 10 | 60 | 4 | … | 29 |



SUMO

Vehicles reaching destination: 30
Average trip time: 120
CO2: 543452
Average waiting time: 10
…

# Solution fitness

- **Statistics about simulation**:
  - Number of vehicles that reach their destination during a given simulation time
  - Average trip time
  - Emissions (CO2, CO, NOx, PMx, HC...)
  - W...

- **Fitness**
  - O...
  - C...
  - Multi-objective approach
  - Derivate values (green waves?)

- **Additional challenge**:
  - SUMO only calculates statistics for vehicles that reach their destination

$$f_{obj} = \frac{T_{trip} + Tsw + VNRTsim}{V_R^2 + P}$$

Trip duration

Waiting times

Vehicles not reaching destination

Vehicles reaching destination

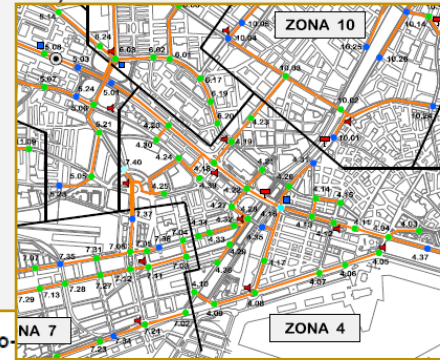Green phases duration

# Realistic data

- **Input for our system:**

  ○ City map

  ○ Location of traffic lights

  ○ Vehicle flow according to several factors (time, weekday, or season)

  ○ Constraints in phase duration

- **Source of the data:**

  ○ **Maps**: OpenStreetMap

  ○ **Traffic lights**: traffic control center of the city

  ○ **Routes**: Mobility department

# Realistic data

- ● Challenges: Maps and traffic lights

  - ○ Incomplete maps: missing road, road directions, traffic lights, …

  - ○ Errors in conversion (OSM => SUMO)

  - ○ Manual correction of the maps. Labour intensive process

- ● Challenges: Routes

  - ○ Not enough details

  - ○ Non-automatable format

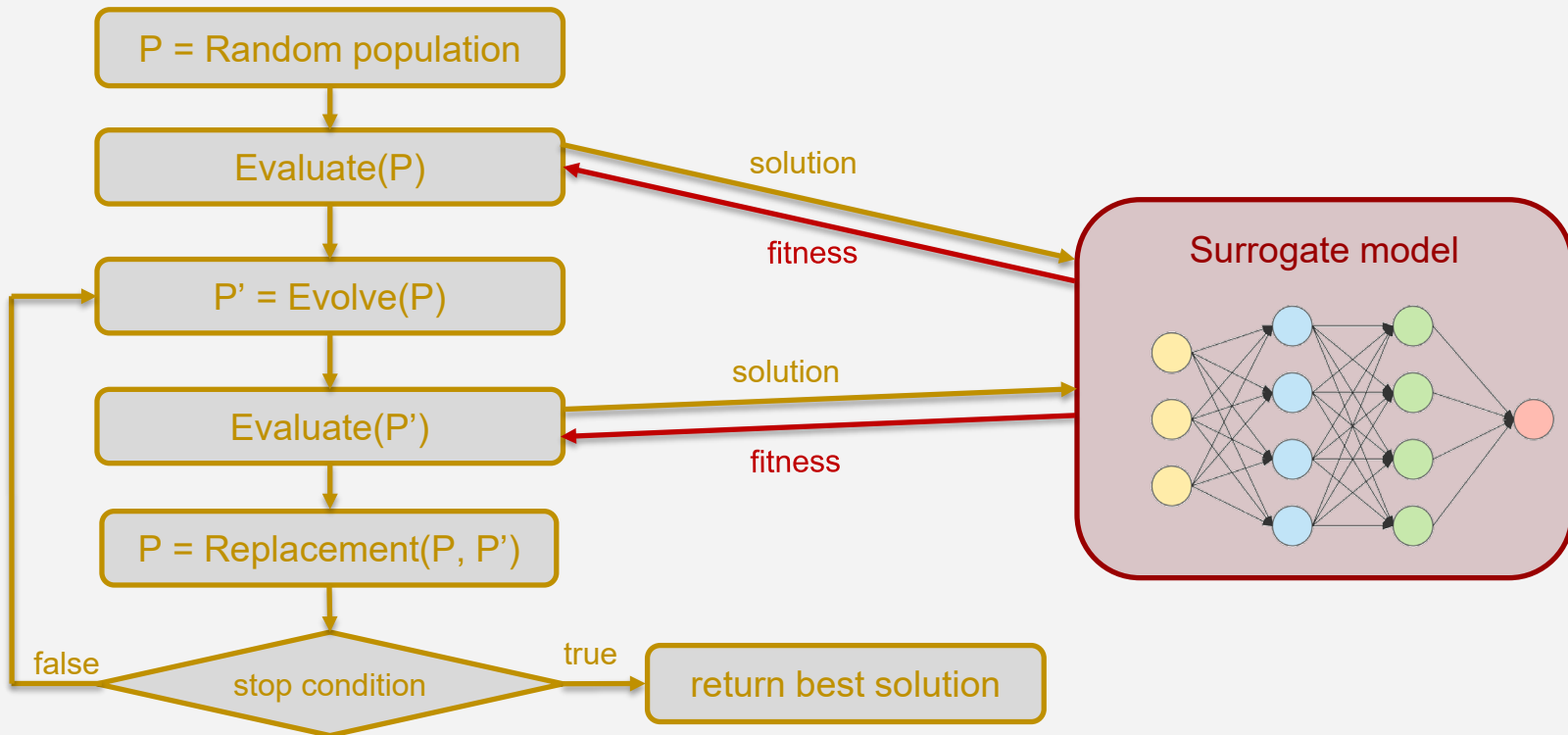  - ○ Conversion of traffic intensity to traffic flows

**Intensidad de vehículos Mayo**

| PM | Ubicación | I.M.D.L | I.M.D.S. | I.M.D.D. | I.M.H.P.L. | H.P.L.M. | I.M.H.P.S. | H.P.S.M. | I.M.H.P.D. | H.P.D.M. |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Avda. Juan Sebastián Elcano - Este | 2.004 | 1.636 | 1.388 | 169 | 12:00 | 149 | 12:00 | 138 | 13:00 |
| 2 | Avda. Juan Sebastián Elcano - Oeste | 20.383 | 17.890 | 15.186 | 1.465 | 8:00 | 1.278 | 21:00 | 1.236 | 21:00 |
| 3 | Bolivia - Este | 16.775 | 15.174 | 12.698 | 1.413 | 14:00 | 1.250 | 14:00 | 1.091 | 13:00 |
| 4 | P.M. Pablo Ruiz Picasso - Este | 19.376 | 15.574 | 12.631 | 1.908 | 14:00 | 1.337 | 14:00 | 1.044 | 13:00 |
| 5 | P.M. Pablo Ruiz Picasso - Oeste | 28.829 | 23.675 | 20.368 | 2.192 | 8:00 | 1.632 | 21:00 | 1.582 | 21:00 |
| 6 | Pso. Reding - Este | 8.528 | 6.767 | 5.296 | 687 | 14:00 | 554 | 14:00 | 418 | 13:00 |
| 7 | Pso. Reding - Oeste | 6.987 | 5.930 | 4.655 | 546 | 9:00 | 445 | 12:00 | 342 | 21:00 |
| 8 | Victoria - Sur * | 5.874 | 5.255 | 4.285 | 419 | 8:00 | 443 | 22:00 | 319 | 22:00 |
| 9 | Victoria - Norte * | 6.474 | 5.888 | 4.923 | 483 | 14:00 | 454 | 22:00 | 360 | 22:00 |
| 10 | Túnel Alcazaba - Este | 15.870 | 14.744 | 12.561 | 973 | 8:00 | 849 | 14:00 | 701 | 1:00 |

Management of DS proyects



Real-world Application: TLSP

Challenges in real-world problems

- Scalability
- Dynamism
- Uncertainty
- Robutness
- Multiple Objectives
- Contraints

Approaches

Efficiency: Parallelism

Efficacy: Hybridization

Surrogate models

Neuro-evolution

Feature Selection