

Séminaire FONISEN

A quick introduction to Surrogate modeling with **SMT**

joseph.morlier@isae-superaero.fr

**P. Saves and R. Lafage and
N. Bartoli , T. Lefevbre,
R. Charayron ...**



About Me?

<https://ica.cnrs.fr/en/author/jmorlier/>

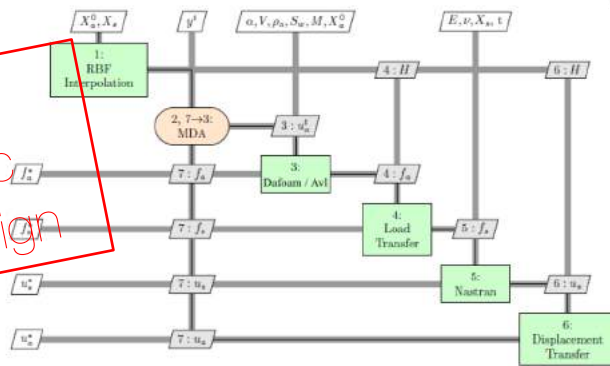
- 6 PhDs, 3 MsCs



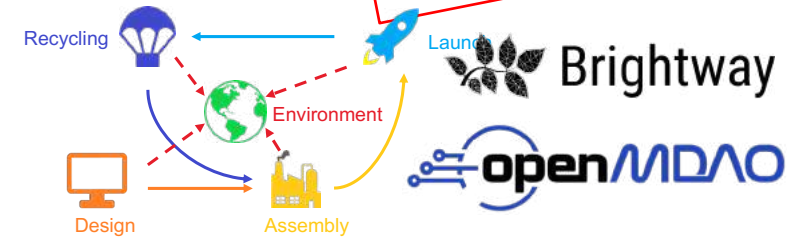
Digital fabrication

EcoOptimization

MDO for Aeroelastic Wing design



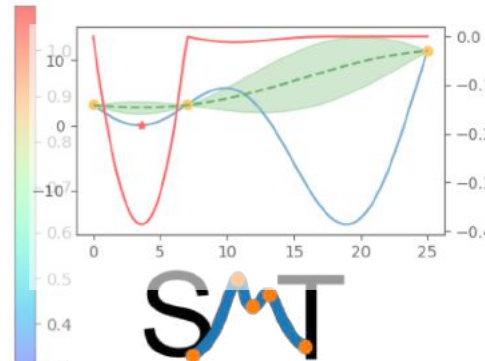
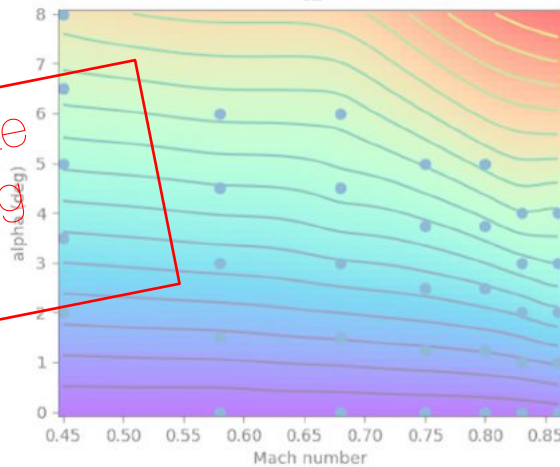
<https://github.com/SMTorg/SMT>



LCA4MDAO (needed LCA database ecoinvent)

<https://github.com/mid2SUPAERO/LCA4MDAO>

Surrogate modeling AI4E



INRIA Seminar 7/12/23



?

Can we Learn from

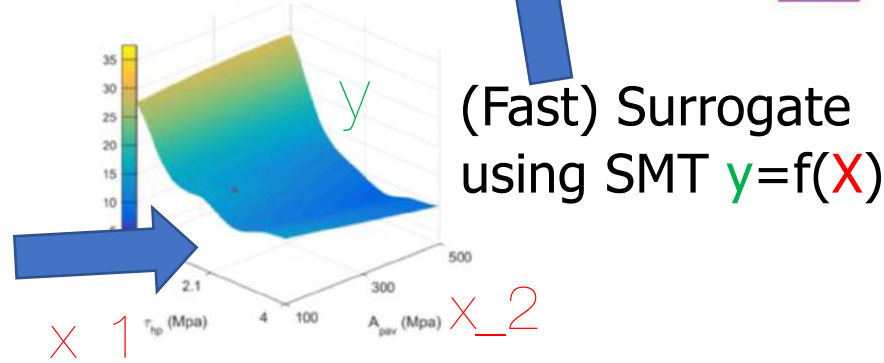
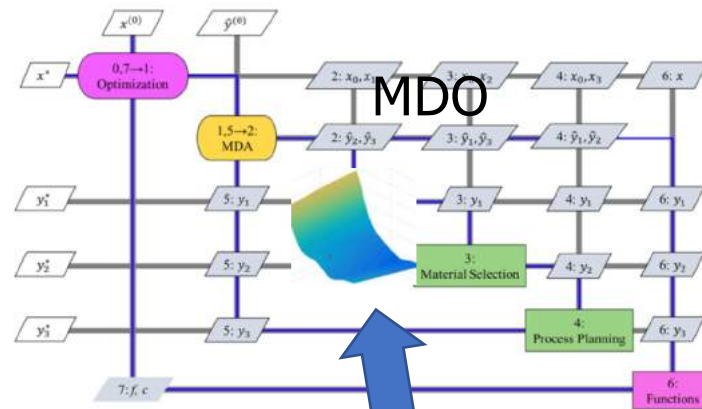
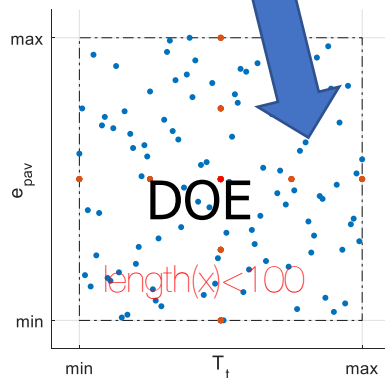
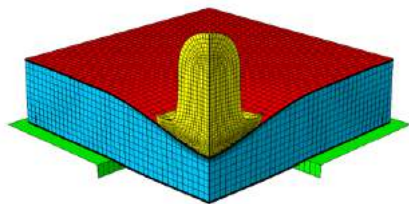
data $y=f(x)$ to predict
engineering (costly) QoI

BTW, How to start with an open source Python toolbox?

What is Surrogate Modeling ?

A. Forrester, A. Sobester, and A. Keane. "Engineering Design via Surrogate Modelling: A Practical Guide". Coll. John Wiley & Sons (2008).

(Expensive) simulation



A surrogate model of a function is an **approximation of Expensive Computer simulation**: It's a supervised learning process in AI.

As the surrogate is less costly to evaluate it can be **used as a "fast" code in a Multidisciplinary Design Optimization loop**. [or do Uncertainty Quantification or do Bayesian Optimization etc...]

Au programme

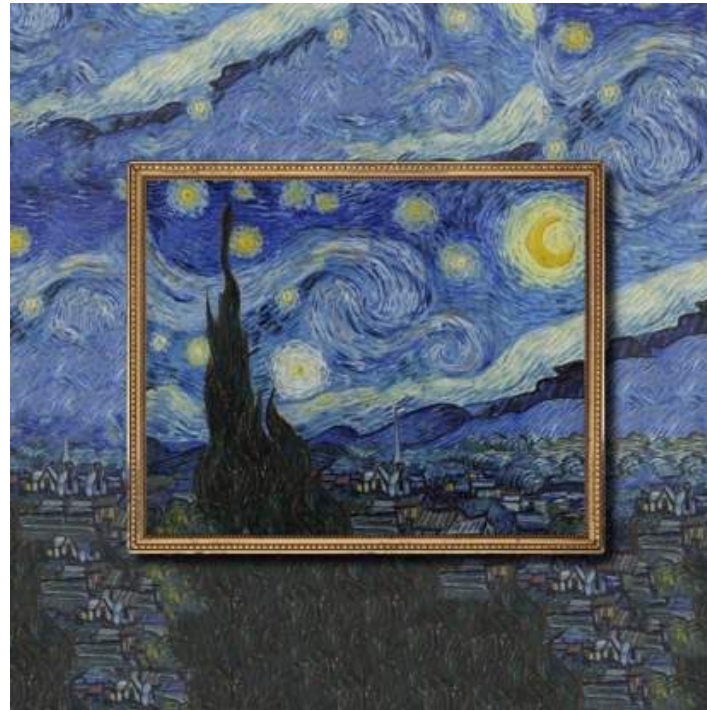
- **Part1: GP aka Kriging**
- **Part2: About SMT**
- **Part3: Engineering applications**

Part 1 : GP aka Kriging

Mathematical foundations

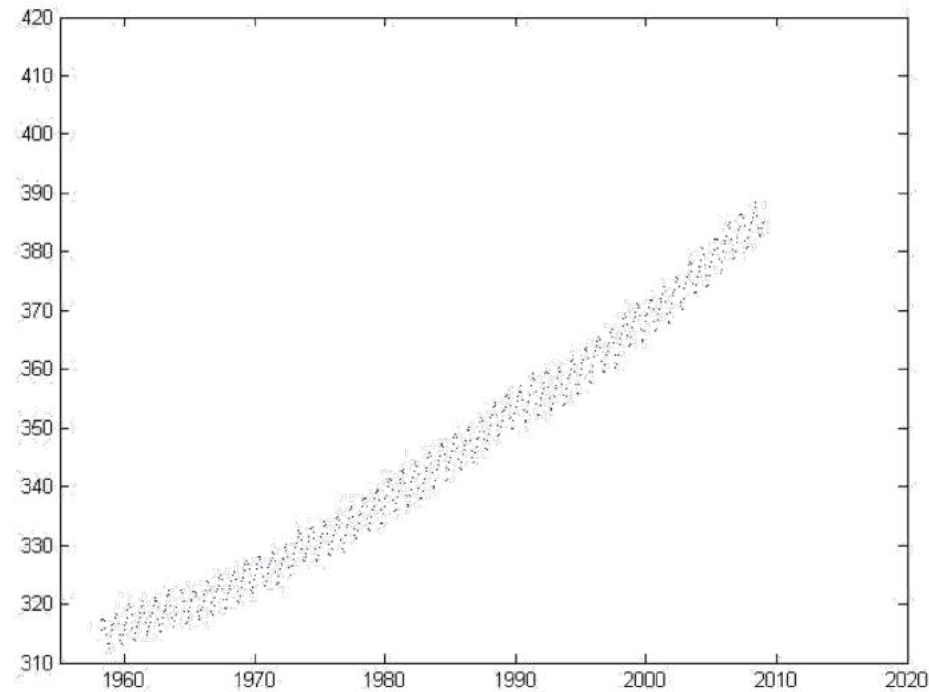
ML vs Engineering

Kriging (Pioneer)	Gaussian Processes (link with AI)
Developed by Daniel Krige – 1951; formalized by Georges Mathéron in the 60's (Mines Paris)	Neural network with infinite neurons tend to Gaussian Process 1994
<p>Krige, D. G., 1951, A statistical approach to some basic mine valuation problems on the Witwatersrand: J. Chem. Metal. Min. Soc. South Africa, v. 52, p. 119-139.</p> <p>Matheron, G., 1963b, Principles of geostatistics: Economic Geol., v. 58, p. 1246-1266.</p>	<p>Neal, R. Priors for infinite networks. Tech. rep., University of Toronto, 1994.</p> <p>— Williams, C. K. I., and Rasmussen, C. E. Gaussian processes for regression. <i>Advances in Neural Information Processing Systems 8</i> (1996), 514-520.</p>



<http://extrapolated-art.com>

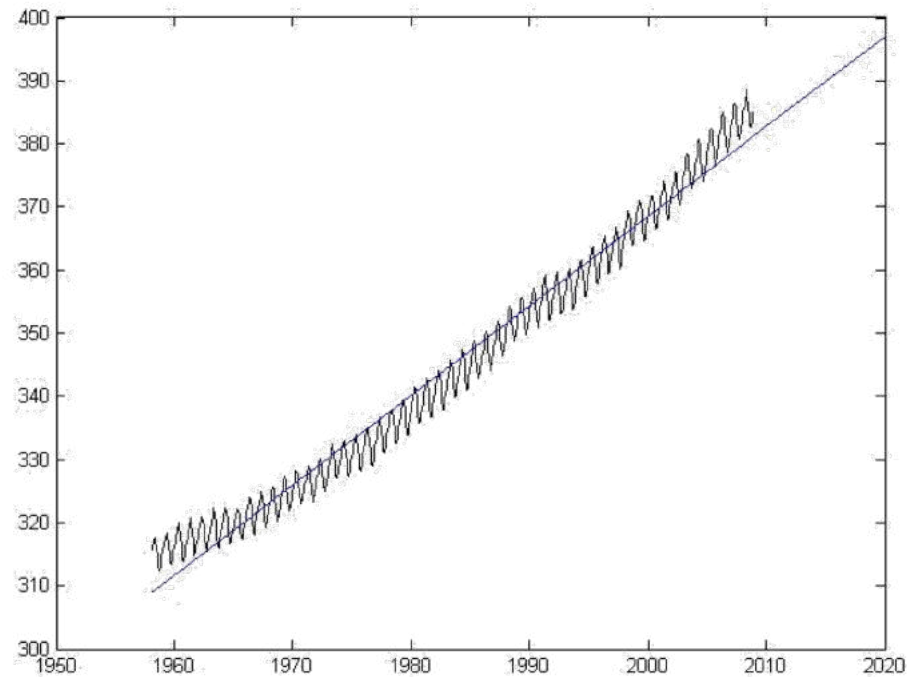
Limit of linear models for prediction



Month-wise data of Co2 concentration in atmosphere at Hawaii

Image Source: <http://mlg.eng.cam.ac.uk/teaching/4f13/1314/>

Example – Linear Regression



Should we choose a **polynomial**?

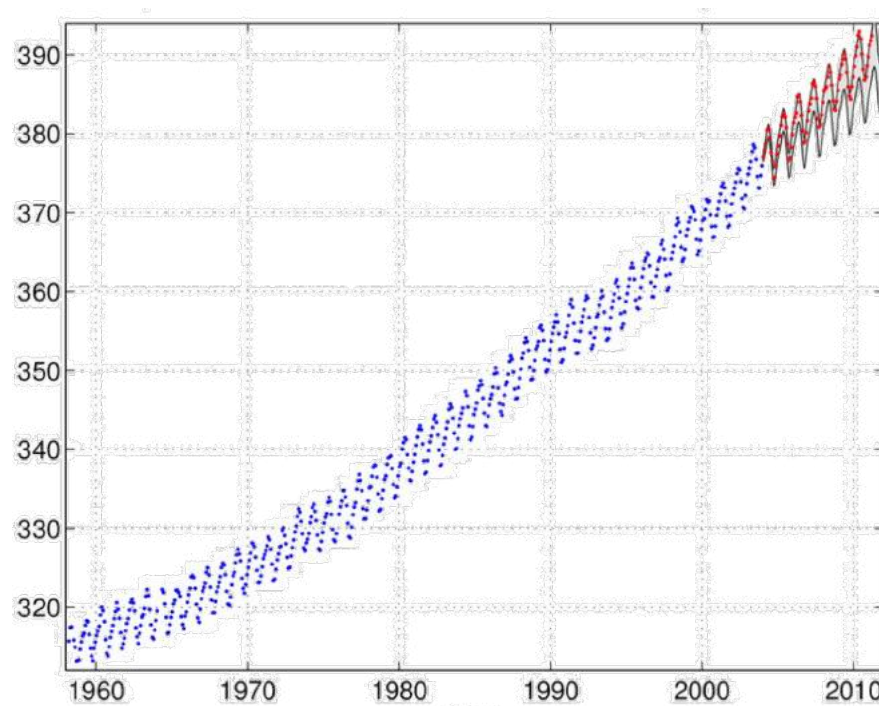
What **degree** of polynomial should we choose? (overfitting)

For a given degree, what **parameters** of polynomial should we choose

Image Source: <http://mlg.eng.cam.ac.uk/teaching/4f13/1314/>

Example – GP

https://scikit-learn.org/stable/auto_examples/gaussian_process/plot_gpr_co2.html



Predicted variance after year 2005 in grey, real data-points in red

Image Source: <http://mlg.eng.cam.ac.uk/teaching/4f13/1314/>

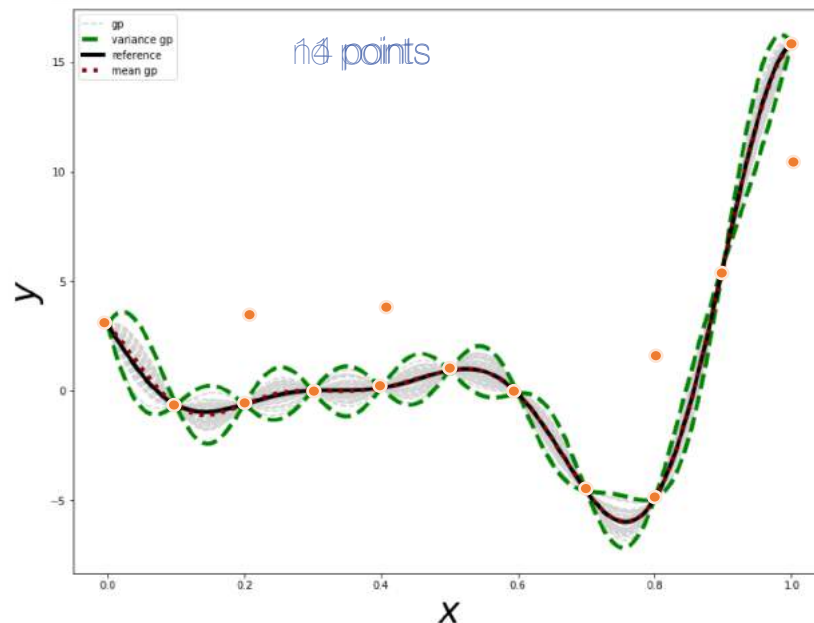
Main steps to build a surrogate model

1. Black box to replace by a surrogate model

$$y(x^{(i)}) = f(x^{(i)}) \text{ with } x^{(i)} \in \mathbb{R}^d$$

where f is the model (expensive, free-derivative, noise-free) such as CFD code or FEM code

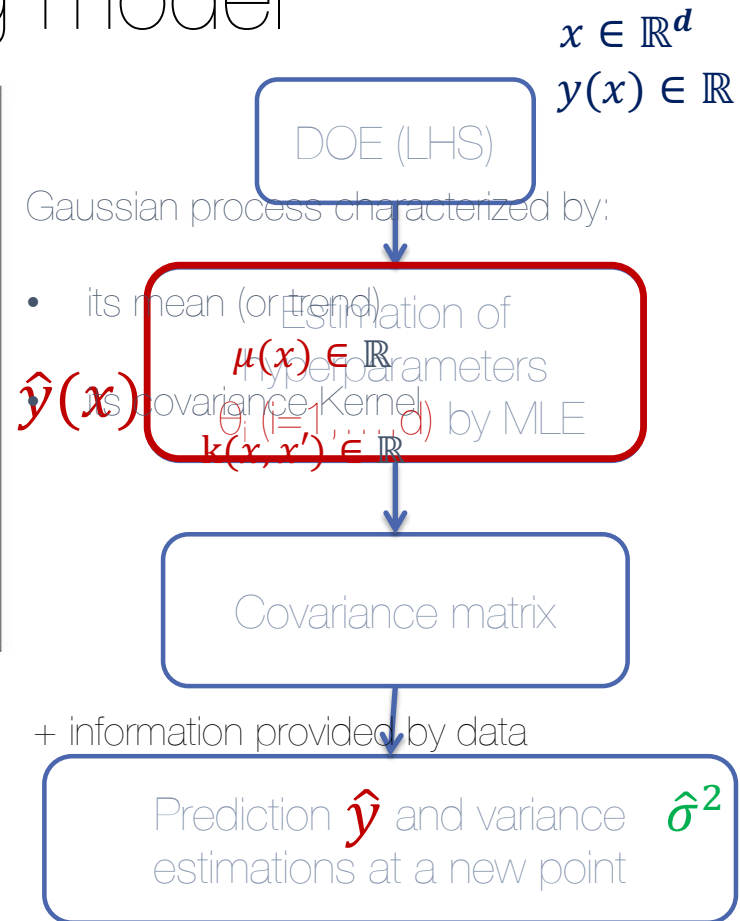
Gaussian process or Kriging model



- Hyperparameters tuning
- Number of hyperparameters increases with the dimension d (number of design variables)
- Curse of dimensionality

$$f(\mathbf{x}) \Rightarrow Y(\mathbf{x}) = \mathcal{N}(\hat{\mathbf{y}}(\mathbf{x}), s^2(\mathbf{x}))$$

D. G. Krige, A statistical approach to some basic mine valuation problems on the Witwatersrand, 1951, Journal of the Southern African Institute of Mining and Metallurgy.



Matrix view of Gaussian Process

1/ Get your inputs/outputs data

2/ You want to predict at x^*

x
 $y(x)$

x^*

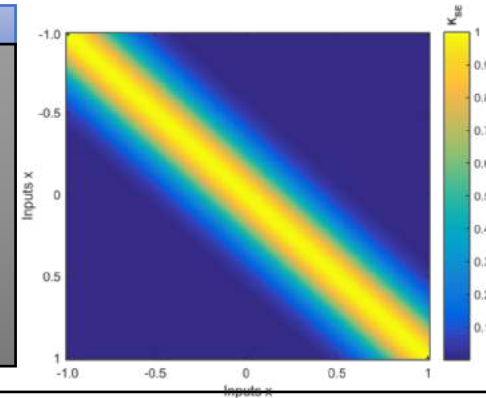
HYP1 Choose a Kernel here: SE

$$k(x, x') = \theta_1^2 \exp\left(-\frac{(x - x')^2}{2\theta_2^2}\right)$$

$=$

3bis/ Construct K_{xx} and do Hyperparameters tuning (2)

x^T
 $[K_{xx}]$



$$m(y_*) = [K_{x_*x_s}] [K_{xx}]^{-1} y(x)$$

4/ compute mean

$$m(x_*) = K_* [K_{xx}]^{-1} y$$

$$\text{cov}(y_*) = [K_{x_sx_s}] - [K_{x_sx_*}] [K_{xx}]^{-1} [K_{xx}]^{-1} [K_{xx}]^{-1} [K_{x_*x_s}]$$

4/ compute variance of estimate

$$\text{var}(x_*, x'_*) = K_{**} - K_*^T [K_{xx}]^{-1} K_*$$

posterior_mean = covXXs @ np.linalg.inv(covXX_noisy) @ y

posterior_cov = covXsXs - covXXs @ np.linalg.inv(covXX_noisy) @

Need deeper understanding?

The Art of Gaussian Processes: Classical and Contemporary

César Lincoln C. Mattos¹ Felipe Tobar²

¹Department of Computer Science
Federal University of Ceará
Fortaleza, Ceará, Brazil
cesarlincoln@dc.ufc.br

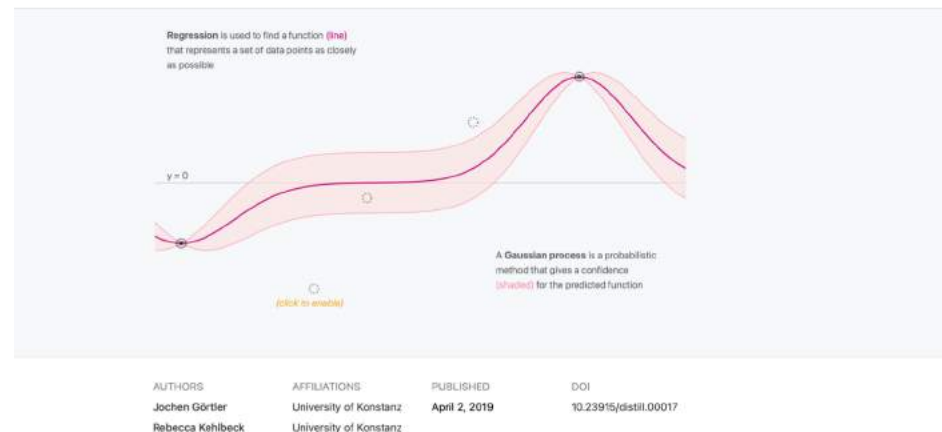
²Initiative for Data & Artificial Intelligence
University of Chile
Santiago, Chile
ftobar@dim.uchile.cl

2021

<https://neurips.cc/virtual/2021/tutorial/21890>

A Visual Exploration of Gaussian Processes

How to turn a collection of small building blocks into a versatile tool for solving regression problems.



<https://distill.pub/2019/visual-exploration-gaussian-processes/>

Part2: About SMT

FAQ

X_0 in practice

Since 2017

<https://smt.readthedocs.io/en/latest>

<https://github.com/SMTorg/smt>



SMT 2.2.1 documentation » SMT: Surrogate Modeling Toolbox

next | index





Table of Contents

- SMT: Surrogate Modeling Toolbox
- Cite us
- Focus on derivatives
- Documentation contents
 - Indices and tables

Next topic

Getting started

This Page

Show Source

Quick search

Go

SMT: Surrogate Modeling Toolbox

The surrogate modeling toolbox (SMT) is an open-source Python package consisting of libraries of surrogate modeling methods (e.g., radial basis functions, kriging), sampling methods, and benchmarking problems. SMT is designed to make it easy for developers to implement new surrogate models in a well-tested and well-document platform, and for users to have a library of surrogate modeling methods with which to use and compare methods.

The code is available open-source on [GitHub](#).

Cite us

To cite SMT 2.0: P. Saves and R. Lafage and N. Bartoli and Y. Diouane and J. H. Bussemaker and T. Lefebvre and J. T. Hwang and J. Morlier and J. R. R. A. Martins.

[SMT 2.0: A Surrogate Modeling Toolbox with a focus on Hierarchical and Mixed Variables Gaussian Processes, ArXiv 2023.](#)

```
nd N. Bartoli and Y. Diouane and J. H. Bussemaker and T. Lefebvre and J. T. Hwang and J. Morlier and J. R. R. A. Martins},  
deling Toolbox with a focus on Hierarchical and Mixed Variables Gaussian Processes},  
xiv.2305.13998},
```

To cite SMT legacy: M. A. Bouhlef and J. T. Hwang and N. Bartoli and R. Lafage and J. Morlier and J. R. R. A. Martins.

[A Python surrogate modeling framework with derivatives, Advances in Engineering Software, 2019.](#)

```
@article{SMT2019,  
  Author = {Mohamed Amine Bouhlef and John T. Hwang and Nathalie Bartoli and Rémi Lafage and Joseph Morlier and Joao  
  Journal = {Advances in Engineering Software},  
  Title = {A Python surrogate modeling framework with derivatives},  
  pages = {102662},  
  issn = {0965-9978},  
  doi = {https://doi.org/10.1016/j.advengsoft.2019.03.005},  
  Year = {2019}}
```

Y. latest

Required packages

SMT depends on the following modules: numpy, scipy, scikit-learn, pyDOE2 and Cython.

Installation

If you want to install the latest release

```
pip install smt
```

or else if you want to install from the current master branch

```
pip install git+https://github.com/SMTorg/smt.git@master
```

Usage

For examples demonstrating how to use SMT, you can take a look at the [tutorial notebook](#) or go to the 'smt/examples' folder.



The very First question:

Q1 : Do you want to create a surrogate model or to do Global Optimization?

The Second question:

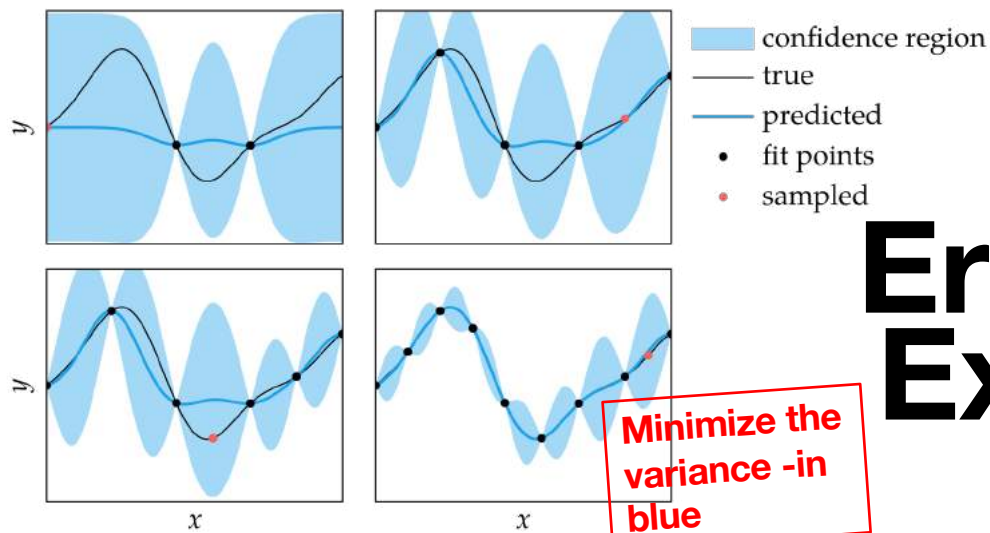
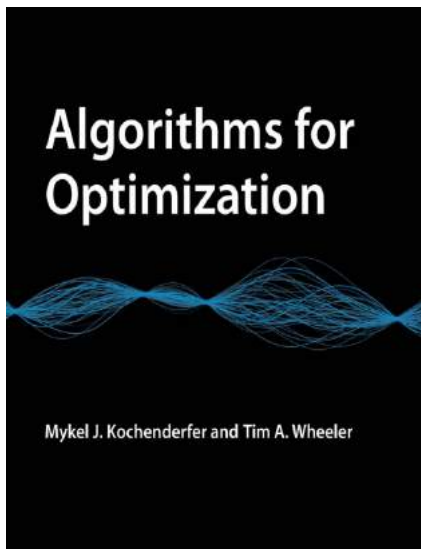
Q2: Are your data cheap or expensive to evaluate?



For this:
you need to understand one key concept

Given a surrogate model (GP) with both prediction and confidence interval, an optimization procedure must balance the search for the expected optimal point and decreasing uncertainty

*In other words, the algorithm must **balance exploitation with exploration to find the global optimum***



I want the most precise surrogate

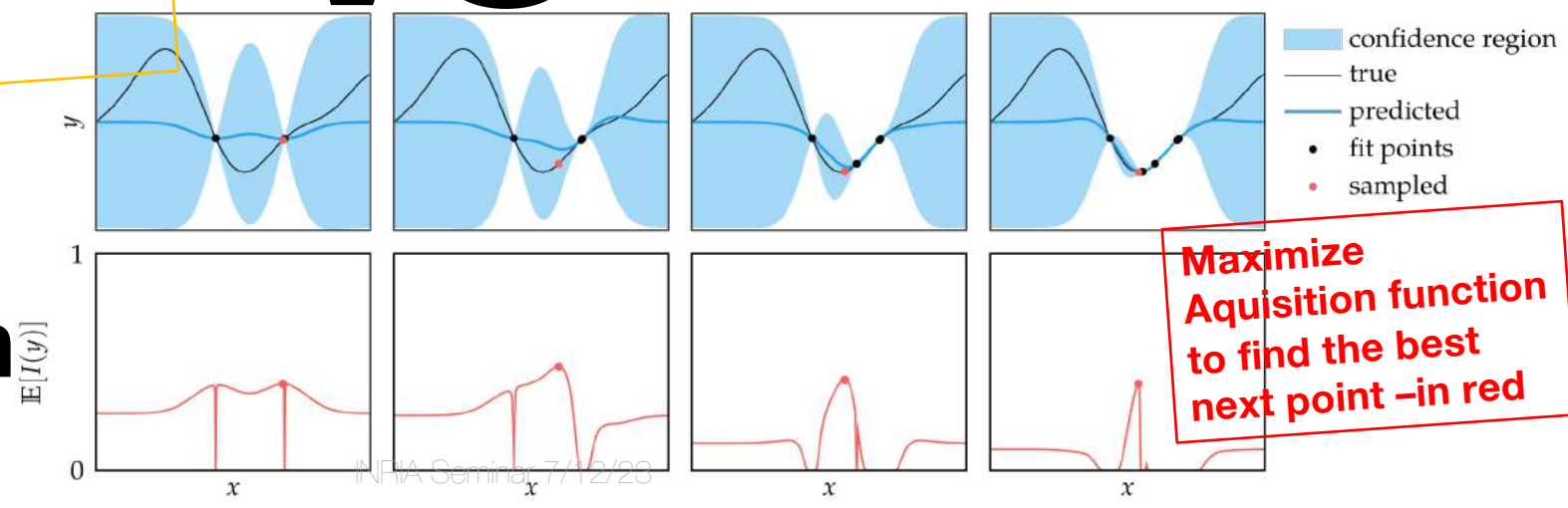
Error-Based Exploration

Minimize the variance -in blue

VS

I want a global optimum at fixed budget

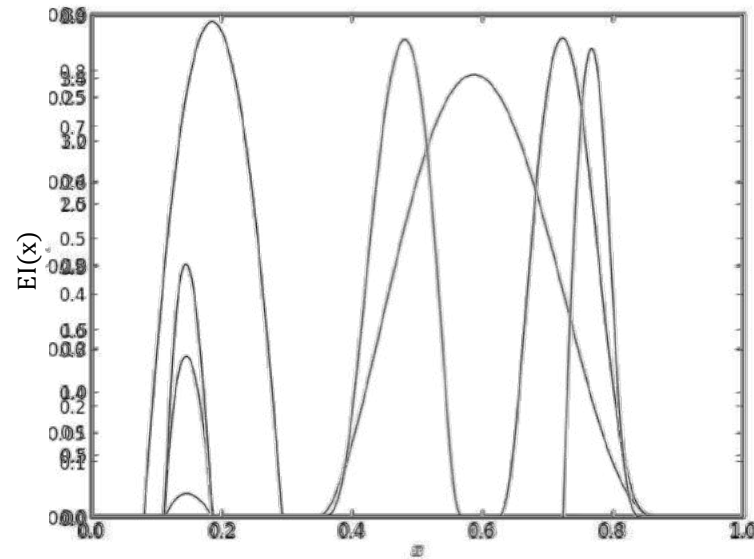
Efficient Global Optimization



Efficient Global Optimization: Illustration

$$\begin{cases} \min (6x - 2)^2 \sin(12x - 4) \\ s.t. \\ 0 \leq x \leq 1 \end{cases}$$

— EI function



N. Bartoli, Optimisation adaptative basée sur les métamodèles, HdR, Université Toulouse III.

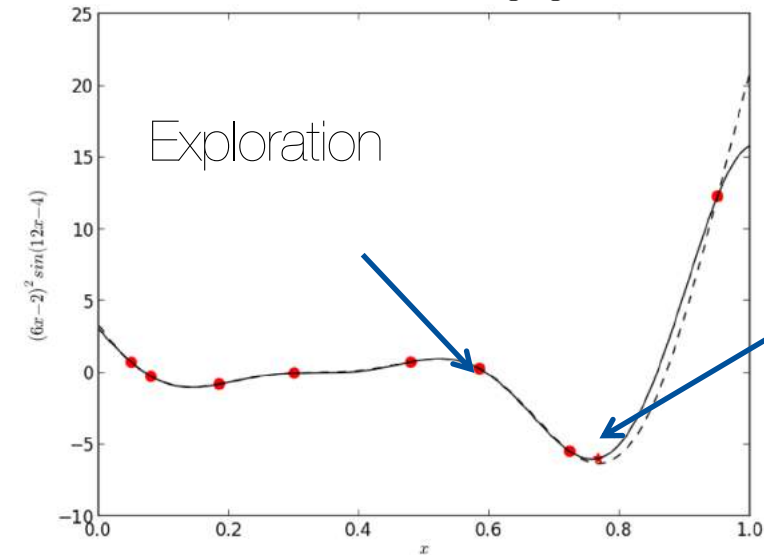
- 4 points for the initial DOE
- 6 iterations

★ Enrichment Samples

• Training Samples

— True function

--- Kriging function



Exploitation



How to use SMT?

Surrogate?

expensive data?
the initial DOE is
small (LHS)

Error-based exploration

if data are not really
smooth, use MOE

cheap data?
the initial DOE is
large

Explainability?

Use first a quadratic
models with interactions!

use MOE=1 pick
the best

Unconstrained
Optimization?

Gradient?

Use GEK with EGO

Derivative free?

Use SE or Matern kernel
with EGO

Large input space
d?

Use KPLS with
EGO

**Rule of thumb:
use initial DOE
= $4 \cdot d$**



Next question:

Can you handle different kinds of variable types?



What kind of
variable types
are available in
SMT 2.0 ?

Continuous

Discrete

Categorical

Hierarchical

Mixed

<https://automl.github.io/ConfigSpace/main/index.html>

Mixed-type:

One-hot encoding:

Red: [1,0,0]

Green: [0,1,0]

Blue: [0,0,1]

Categorical, continuous, discrete, binary variables

Categorical: Red, Green, Blue

Continuous: [0,1]

Discrete: 1,2,3,4,5

Binary: Yes, No

Tuning hyperparameters for deep neural network

$$= f(\underbrace{x_1, x_2}_{\substack{\text{continuous variables} \\ \text{learning rate } \in [1e^{-6}, 1e^{-1}] \quad \text{weight decay } \in [1e^{-6}, 1e^{-1}]}}, \underbrace{x_3, x_4}_{\substack{\text{categorical variables} \\ \text{optimiser type } \in \{SGD, Adam, \dots\} \quad \text{activation type } \in \{tanh, sigmoid, \dots\}}})$$



Multiple categorical - each categorical has multiple options

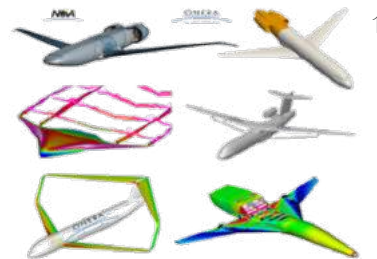
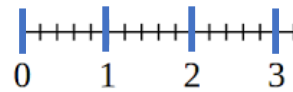
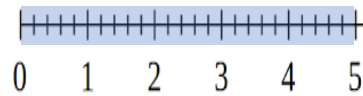


Models to handle mixed variables

Hybrid variables

Variables types:

- Continuous (x) Ex: wing length
- Integer (z) Ex: winglet number
- Categorical (u) Ex: Plane shape / material properties



Categorical variables: n variables, n=2

u1= shape

u2= color

Levels: L_i levels for i in $1, \dots, n$, $L_1=3$, $L_2=2$

Levels(u1)= square, circle, rhombus

Levels(u2)= blue, red

Categories: $\prod_{i=1}^n L_i$, $2*3=6$

- Blue square
- Blue circle
- Blue rhombus
- Red square
- Red circle
- Red rhombus

6 possibilities



Mixed variables Kriging: $x \in [a, b]$ $x \in \{0, 1, \dots, N\}$ $x \in \{blue, red, green\}$

$$k(x, x', \Theta) = k_{con}(x_{con}, x'_{con}, \theta_{con}) k_{int}(x_{int}, x'_{int}, \theta_{int}) \underbrace{k_{cat}(x_{cat}, x'_{cat}, \theta_{cat})}$$

$$k_{cat}^{CR}(x_{cat}, x'_{cat}, \theta_{cat}) = \prod_{i=1}^{d_{cat}} k_{cont,i}(e_{x_{cat,i}}, e_{x'_{cat,i}}, \theta_{cat,i})$$

SO RELAXED DIMENSION IS now REALLY

... **BIGGER**

Encoding **x-mixed** → continuous space (bigger) x_{encoded} → Optimize x^*_{encoded}
→ Decoding $x^*_{\text{engineering}}$

SoA- Continuous relaxation

Example with 1 categorical variable and two levels

- Red color
- Blue color

→ One-hot encoding: Categorical variable replaced by two continuous variables denoted by X_1 and X_2

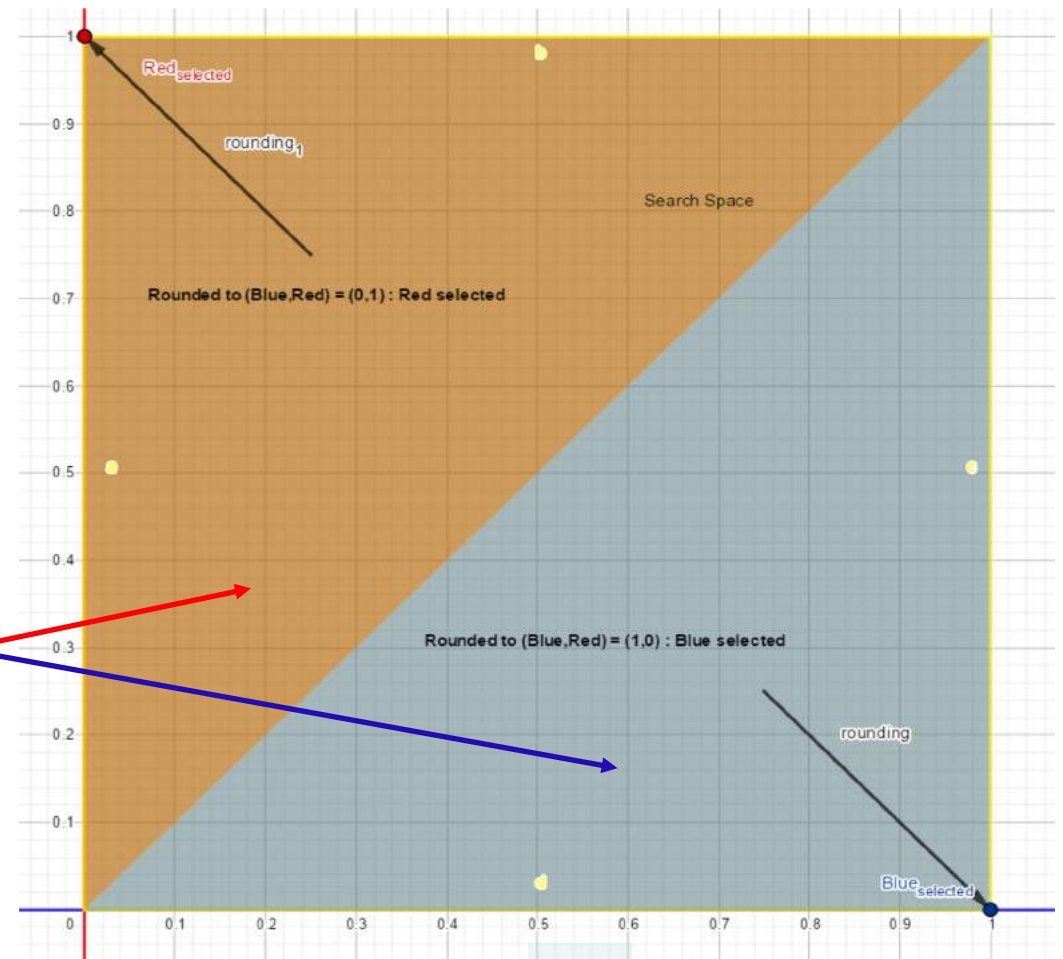
- If $X_1 > X_2 \Rightarrow e_{c_1^b} = (1., 0.) \Rightarrow$ Blue color
- If $X_1 < X_2 \Rightarrow e_{c_1^r} = (0., 1.) \Rightarrow$ Red color

n relaxed dimension

$$x^r, x^s \in \mathbb{R}^n$$

A continuous kernel

$$k(x^r, x^s, \theta^{cont}) = \prod_{j=1}^n \exp \left(-(x_j^r - x_j^s) \theta_j^{cont} (x_j^r - x_j^s) \right)$$



Mixed GP

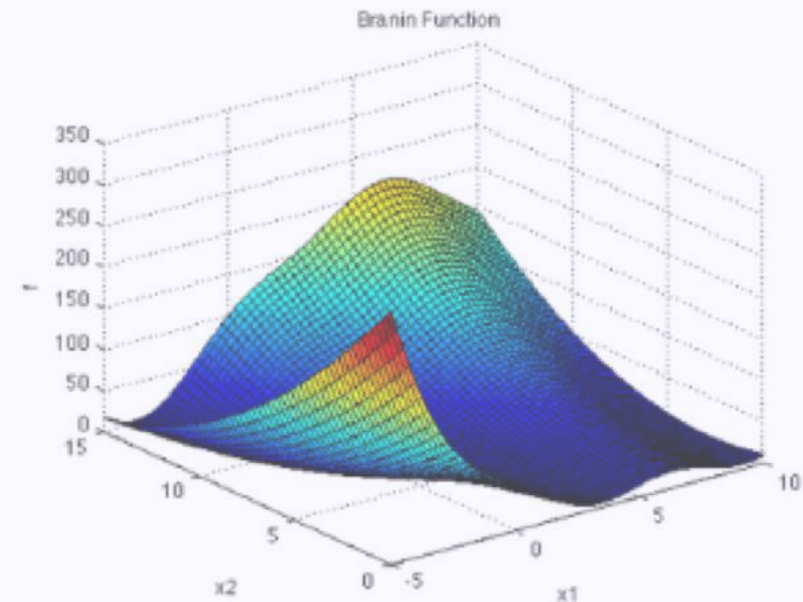
Validation problem $n = 2$

- **Variable types: continuous and categorical with 10 levels. $n' = 11$**

$$f(x, z) = \begin{cases} \cos(3.6\pi(x-2)) + x - 1 & \text{if } z = 1, \\ 2 \cos(1.1\pi \exp(x)) - \frac{x}{2} + 2 & \text{if } z = 2, \\ \cos(2\pi x) + \frac{1}{2}x & \text{if } z = 3, \\ x \left(\cos(3.4\pi(x-1)) - \frac{x-1}{2} \right) & \text{if } z = 4, \\ -\frac{x^2}{2} & \text{if } z = 5, \\ 2 \cos\left(\frac{\pi}{4} \exp(-x^4)\right)^2 - \frac{x}{2} + 1 & \text{if } z = 6, \\ x \cos(3.4\pi x) - \frac{x}{2} + 1 & \text{if } z = 7, \\ x \left(-\cos\left(7\frac{\pi}{2}x\right) - \frac{x}{2} \right) + 2 & \text{if } z = 8, \\ -\frac{x^5}{2} + 1 & \text{if } z = 9, \\ -\cos\left(5\frac{\pi}{2}x\right)^2 \sqrt{x} - \frac{\ln(x+0.5)}{2} - 1.3 & \text{if } z = 10. \end{cases}$$

Branin 2D with

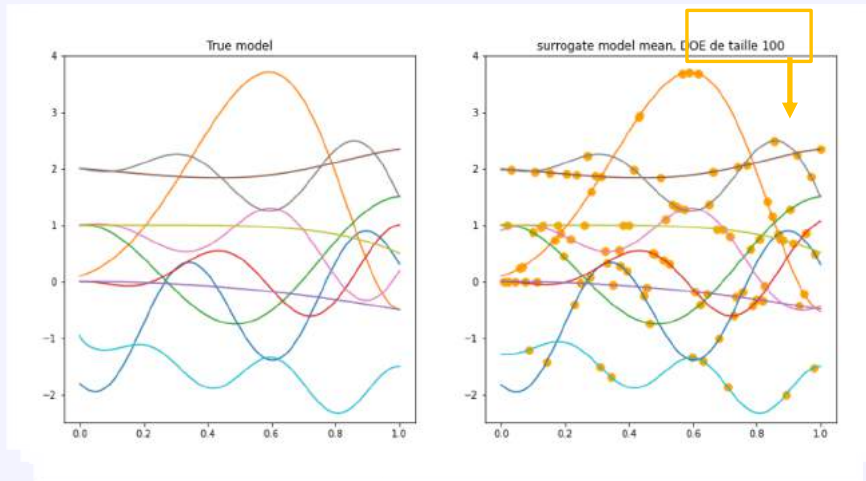
\mathbf{X}_1 as an integer in $[-5, 10]$ and \mathbf{X}_2 continuous in $[0, 15]$



Mixed GP

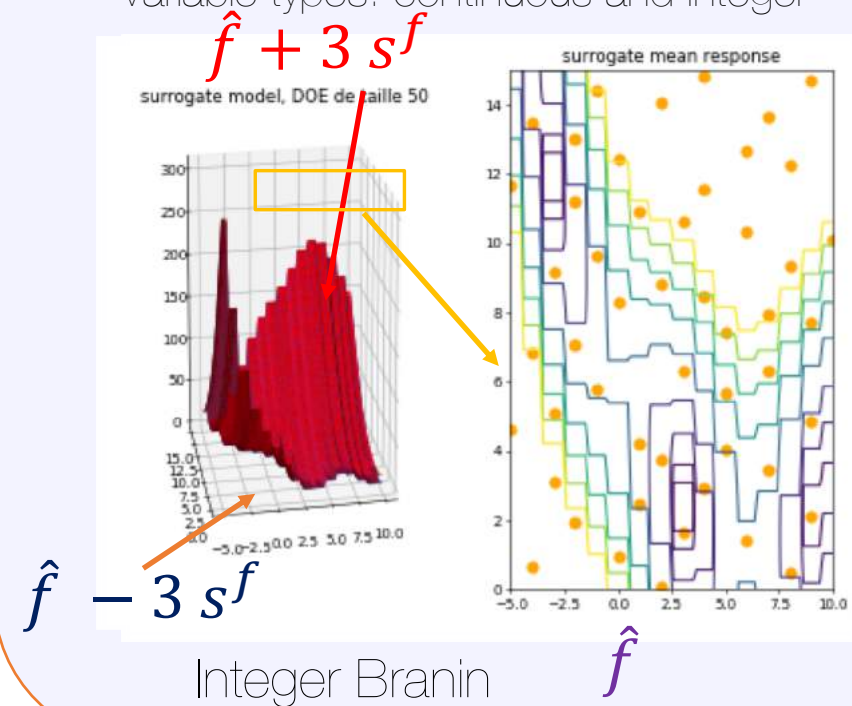
Validation problem $n = 2$

- Variable types: continuous and categorical with 10 levels. $n' = 11$



Validation problem $n = 2$

- Variable types: continuous and integer



Part3: Engineering applications


Most of them are opensource

Example: Accelerating Material discovery

Predict y the crystal structure type (monoclinic, orthorhombic, triclinic)

from x Lithium-ion physical and chemical compound information

i.e. learn from learning database $y=f(x)$

Materials Id	Formula	Spacegroup	Formation Ene...
The unique ID of the material as stated on materialsproject.org	Chemical formula of the material	Spacegroup	Formation Energy
339 unique values	LiFeSiO4 12% LiCoSiO4 9% Other (268) 79%	P1 21% P21/c 9% Other (235) 69%	







The function inputs are:

5 continuous variables [Formation Energy (eV) E Above Hull (eV) Band Gap (eV) Nsites Density (gm/cc) Volume]

1 ordinal variable [Nsites]

4 categorical variables [Materials Id Formula Spacegroup Has Bandstructure] with **respectively: 339, 114 and 44 levels.**

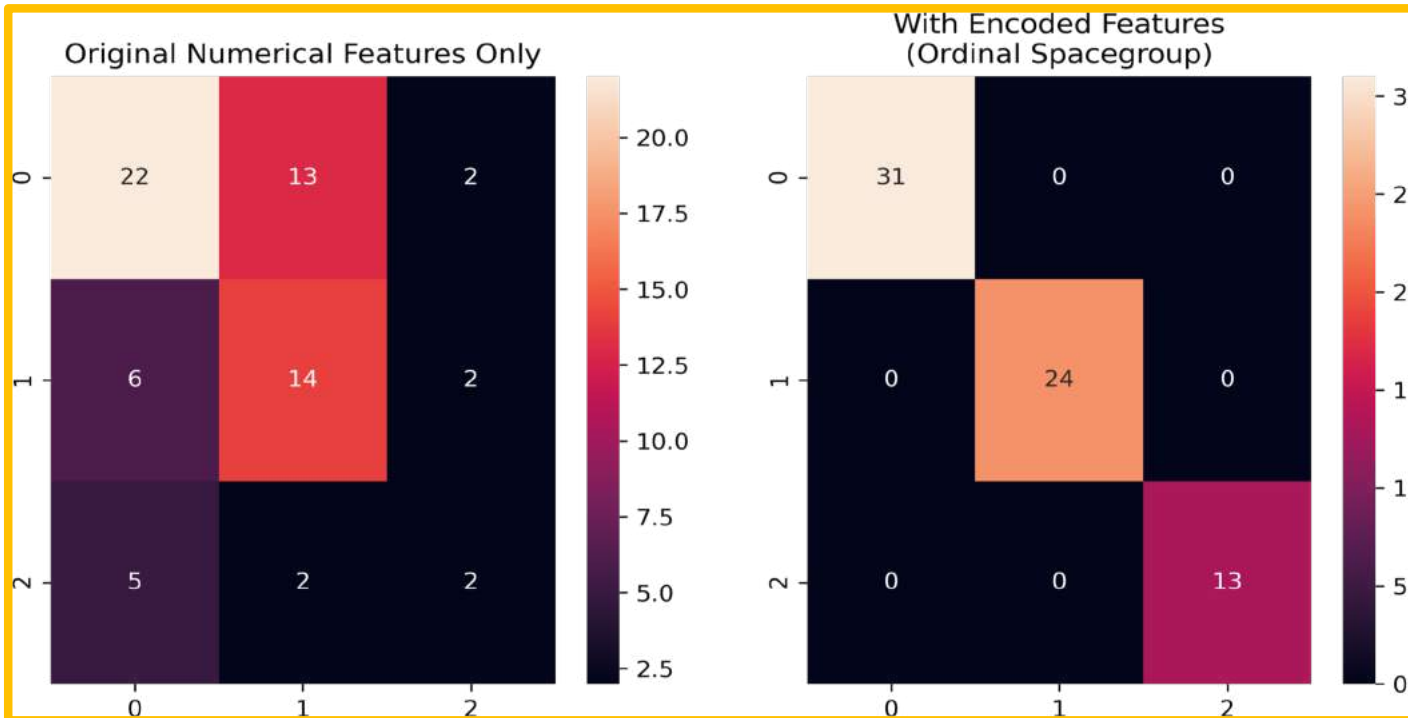
The relaxed dimension by **one-hot-encoding** is **503**

# E Above Hull (...)	# Band Gap (eV)	# Nsites	# Density (gm/cc)	# Volume	Has Bandstruc...
Energy if decomposition of material into most stable ones	Band Gap (in eV)	Number of atoms in the unit cell of the crystal	The density of bulk crystalline materials	The unit cell volume of the material	Boolean variable for bandstructure
					
0 0.19	0 3.82	10 132	2.2 4.2	123 1.52k	true 274 81% false 65 19%

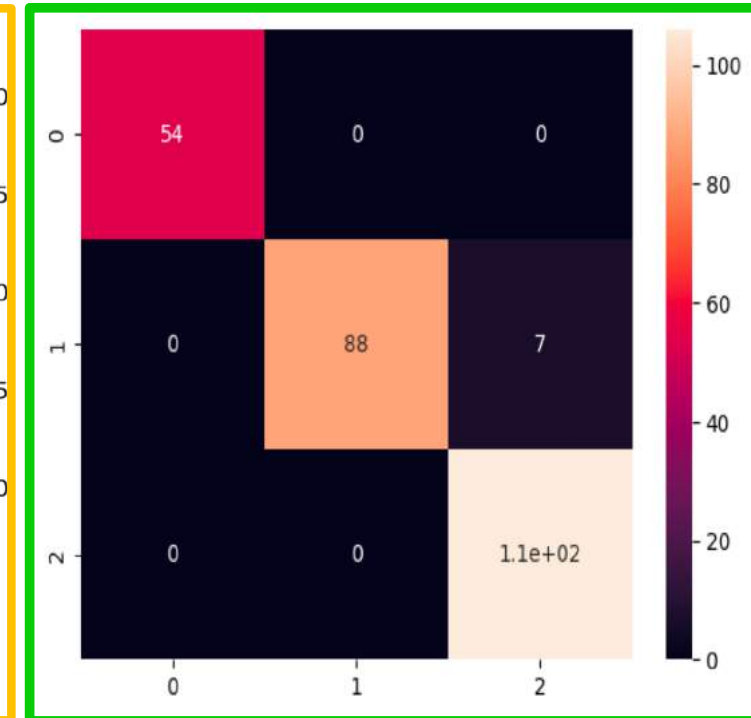
... REALLY
HIGH DIMENSION

Agrawal, D., Crystal System Properties for Li-ion batteries, Properties of Li-ion silicate to predict the crystal system class of the battery, Kaggle, March 2020.

Results: Accelerating Material discovery



scikit-learn with 80/20
dtree w/wo specific features (*sf*)



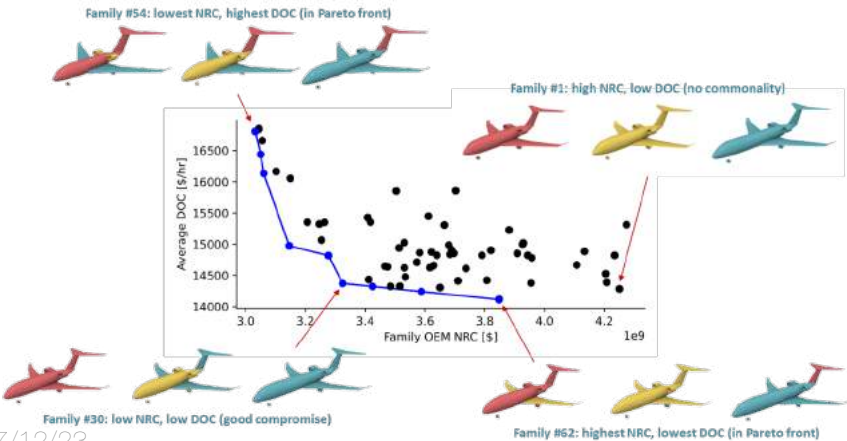
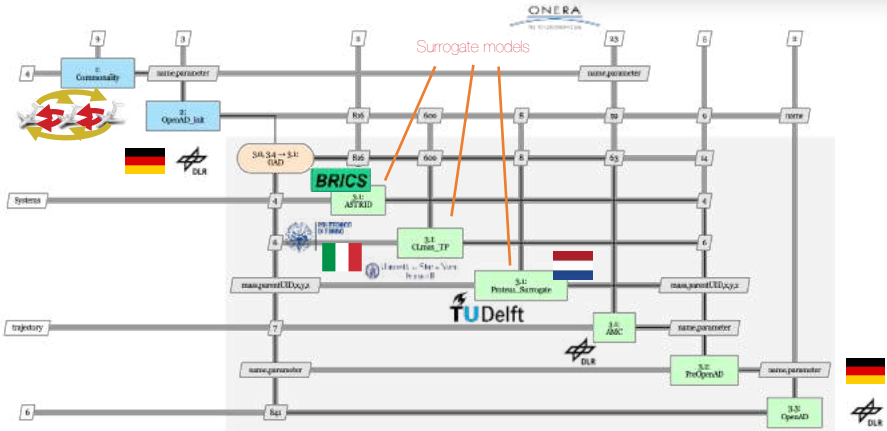
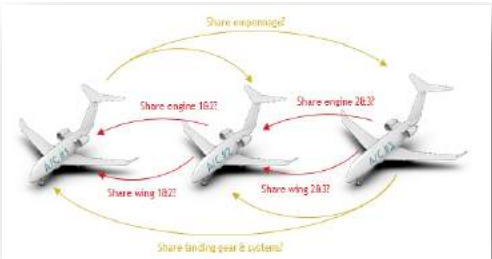
SMT with 10/90 (!!) wo *sf*

Example in BO with hierarchical variables

Select the best family with commonality choices (wing, engine, empennage)

Objectives	2 obj: Min (Direct Operating Costs, OEM Non – Recurring Costs)
Design variables	10 categorical: commonality with 2 levels 9 continuous: Leading edge sweep, rear spar location, Wing t/c for each family member
Constraints	2 ineq. : Balanced Field Length, Landing Field Length

1. Handling of hierarchical variables using imputation method (mean value used when the var. is inactive)
2. DOE and offline optimization to fine tune the optimization parameter

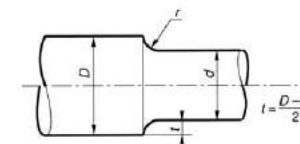


Conclusions

and more

Surrogate is the new abacus

Coefficient de concentration de contrainte : K_t .



$$\sigma_{\text{nomiale}} = \frac{N}{S} \text{ d'où } \sigma_{\text{max}} = K_t \sigma_{\text{nomiale}}$$

Condition de résistance : $\sigma_{\text{max}} < R_{pe}$

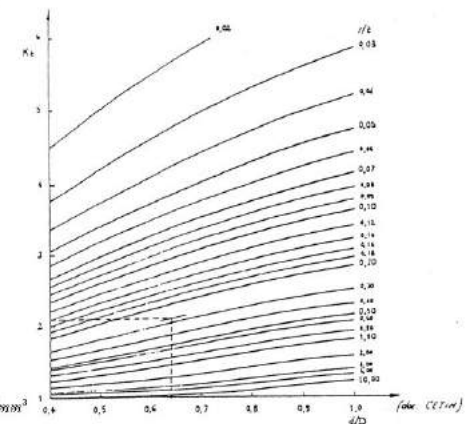
Exemple : $D=100, d=64, r=5$
 $N = 5000 \text{ daN}$

$$\left. \begin{aligned} \frac{d}{D} &= \frac{64}{100} = 0,64 \\ \frac{r}{D-d} &= \frac{5}{100-64} = 0,278 \end{aligned} \right\} K_t = 2,1$$

$$\sigma_{\text{nomiale}} = \frac{4 \times 5000}{\pi \times 64^2} = 1,55 \text{ daN/mm}^2$$

$$\sigma_{\text{max}} = K_t \times \sigma_{\text{nomiale}} = 2,1 \times 1,55 = 3,26 \text{ daN/mm}^2$$

Arbre épaulé en traction



Conclusions

- SMT 1.0: A Python surrogate modeling framework **with derivatives** (2019)
- SMT 2.0: A Surrogate Modeling Toolbox with a focus on Hierarchical and **Mixed Variables** Gaussian Processes (2023)
- Explainability, ...

<https://doi.org/10.1016/j.advengsoft.2019.03.005>

<https://doi.org/10.1016/j.advengsoft.2023.103571>



Documentation:

<https://smt.readthedocs.io/en/latest/>

Code:

<https://github.com/SMTorg/SMT>

How to start:

<https://github.com/SMTorg/smt/tree/master/tutorial#readme>



Mixed Optimization in 200 dimensions?

- High dimension causes problem for optimization.
 - Statistical challenge: the search space **grows exponentially**
 - Computational challenge: global optimizers **fail to return an optimum** within limited time and resource



--> refactoring & opensourcing our Constrained Bayesian Optimization code based on SMT called SEGOMOE

GPyOpt

Tune your algorithms and your design wetlab experiment



botorch

Bayesian optimization in PyTorch (by pytorch)



Ax

Adaptive Experimentation Platform (by facebook)



https://scikit-learn.org/stable/modules/gaussian_process.html

https://scikit-optimize.github.io/stable/auto_examples/bayesian-optimization.html

SEGOMOE is available Through a
webservice

<https://github.com/whatsopt/wopsego>

joseph.morlier@isae-supaero.fr
nathalie.bartoli@onera.fr

MERCI

ONERA: P. Saves and R. Lafage and N. Bartoli , T. Lefevbre, R. Charayron...

ISAE-SUPAERO: R. Rufato Careira, P. Saves, R. Charayron

Polytechnique Montréal: Y. Diouane

DLR: J. H. Bussemaker

UCSD: J. T. Hwang

UoM: J. R. R. A. Martins

SMT 1.0 coding due to Mohamed Amine Bouhlel (INTEL)