

# **Hyperparameter Tuning Cookbook**

**A guide for scikit-learn, PyTorch, river, and spotPython**

Thomas Bartz-Beielstein

Jan 9, 2024

# Table of contents

<b>Preface</b>	<b>12</b>
Book Structure . . . . .	12
Software Used in this Book . . . . .	13
Citation . . . . .	13
<b>I Optimization</b>	<b>15</b>
<b>1 Introduction: Optimization</b>	<b>16</b>
1.1 Optimization, Simulation, and Surrogate Modeling . . . . .	16
1.2 Surrogates . . . . .	16
1.2.1 Costs of Simulation . . . . .	17
1.2.2 Mathematical Models and Meta-Models . . . . .	17
1.2.3 Surrogates = Trained Meta-models . . . . .	17
1.2.4 Computer Experiments . . . . .	17
1.2.5 Limits of Mathematical Modeling . . . . .	18
1.2.6 Example: Why Computer Simulations are Necessary . . . . .	18
1.2.7 Simulation Requirements . . . . .	18
1.3 Jupyter Notebook . . . . .	19
<b>2 Aircraft Wing Weight Example</b>	<b>20</b>
2.1 AWWE Equation . . . . .	20
2.2 AWWE Parameters and Equations (Part 1) . . . . .	20
2.3 Goals: Understanding and Optimization . . . . .	21
2.4 Properties of the Python “Solver” . . . . .	22
2.5 Plot 1: Load Factor ( $N_z$ ) and Aspect Ratio ( $A$ ) . . . . .	23
2.6 Plot 2: Taper Ratio and Fuel Weight . . . . .	25
2.7 The Big Picture: Combining all Variables . . . . .	26
2.8 AWWE Landscape . . . . .	30
2.9 Summary of the First Experiments . . . . .	31
2.10 Exercise . . . . .	31
2.10.1 Adding Paint Weight . . . . .	31
2.11 Jupyter Notebook . . . . .	32

<b>3 Introduction to <code>scipy.optimize</code></b>	<b>33</b>
3.1 Derivative-free Optimization Algorithms . . . . .	34
3.1.1 Nelder-Mead Simplex Algorithm . . . . .	34
3.1.2 Powell's Method . . . . .	35
3.2 Gradient-based optimization algorithms . . . . .	36
3.2.1 An Introductory Example: Broyden-Fletcher-Goldfarb-Shanno Algorithm (BFGS) . . . . .	36
3.2.2 Background and Basics for Gradient-based Optimization . . . . .	38
3.2.3 Gradient . . . . .	38
3.2.4 Jacobian Matrix . . . . .	38
3.2.5 Hessian Matrix . . . . .	39
3.2.6 Gradient for Optimization . . . . .	40
3.2.7 Newton Method . . . . .	42
3.2.8 BFGS-Algorithm . . . . .	45
3.2.9 Procedure: . . . . .	45
3.2.10 Visualization BFGS for Rosenbrock . . . . .	46
3.3 Gradient- and Hessian-based optimization algorithms . . . . .	47
3.3.1 Newton-Conjugate-Gradient Algorithm . . . . .	47
3.3.2 Trust-Region Newton-Conjugate-Gradient Algorithm . . . . .	47
3.3.3 Trust-Region Truncated Generalized Lanczos / Conjugate Gradient Algorithm . . . . .	47
3.4 Global Optimization . . . . .	47
3.4.1 Dual Annealing Optimization . . . . .	51
3.4.2 Differential Evolution . . . . .	51
3.4.3 DIRECT . . . . .	51
3.4.4 SHGO . . . . .	52
3.4.5 Basin-hopping . . . . .	52
3.5 Jupyter Notebook . . . . .	52
<b>4 Sequential Parameter Optimization: Using <code>scipy Optimizers</code></b>	<b>53</b>
4.1 The Objective Function Branin . . . . .	53
4.2 The Optimizer . . . . .	54
4.2.1 TensorBoard . . . . .	55
4.3 Print the Results . . . . .	56
4.4 Show the Progress . . . . .	57
4.5 Exercises . . . . .	58
4.5.1 <code>dual_annealing</code> . . . . .	58
4.5.2 <code>direct</code> . . . . .	59
4.5.3 <code>shgo</code> . . . . .	61
4.5.4 <code>basinhopping</code> . . . . .	63
4.5.5 Performance Comparison . . . . .	65
4.6 Jupyter Notebook . . . . .	66

<b>II Numerical Methods</b>	<b>67</b>
<b>5 Introduction: Numerical Methods</b>	<b>68</b>
5.1 Response Surface Methods: What is RSM? . . . . .	68
5.1.1 Visualization: Problems in Practice . . . . .	71
5.1.2 RSM: Strategies . . . . .	71
5.1.3 RSM: Noise in the Empirical Model . . . . .	72
5.1.4 RSM: Natural and Coded Variables . . . . .	72
5.1.5 RSM Low-order Polynomials . . . . .	73
5.2 First-Order Models (Main Effects Model) . . . . .	73
5.2.1 First-Order Model Properties . . . . .	74
5.2.2 First-order Model with Interactions in python . . . . .	75
5.2.3 Observations: First-Order Model with Interactions . . . . .	76
5.3 Second-Order Models . . . . .	76
5.3.1 Second-Order Models: Properties . . . . .	77
5.3.2 Example: Stationary Ridge . . . . .	77
5.3.3 Observations: Second-Order Model (Ridge) . . . . .	78
5.3.4 Example: Rising Ridge . . . . .	79
5.3.5 Summary: Rising Ridge . . . . .	80
5.3.6 Falling Ridge . . . . .	80
5.3.7 Saddle Point . . . . .	80
5.3.8 Interpretation: Saddle Points . . . . .	81
5.3.9 Summary: Ridge Analysis . . . . .	81
5.4 General RSM Models . . . . .	82
5.4.1 Ordinary Least Squares . . . . .	82
5.5 Designs . . . . .	82
5.5.1 Different Designs . . . . .	82
5.6 RSM Experimentation . . . . .	83
5.6.1 First Step . . . . .	83
5.6.2 Second Step . . . . .	83
5.6.3 Third Step . . . . .	83
5.7 RSM: Review and General Considerations . . . . .	83
5.7.1 Historical Considerations about RSM . . . . .	84
5.7.2 Status Quo . . . . .	84
5.7.3 The Role of Statistics . . . . .	85
5.7.4 New RSM is needed: DACE . . . . .	85
5.8 Exercises . . . . .	86
5.9 Jupyter Notebook . . . . .	86
<b>6 Kriging (Gaussian Process Regression)</b>	<b>87</b>
6.1 DACE and RSM . . . . .	87
6.2 Background: Expectation, Mean, Standard Deviation . . . . .	88
6.2.1 Sample Mean . . . . .	88

6.2.2	Variance and Standard Deviation . . . . .	88
6.2.3	Standard Deviation . . . . .	89
6.2.4	Calculation of the Standard Deviation with Python . . . . .	89
6.2.5	The Empirical Standard Deviation . . . . .	89
6.2.6	The Argument “axis” . . . . .	90
6.3	Data Types and Precision in Python . . . . .	90
6.4	Distributions and Random Numbers in Python . . . . .	92
6.4.1	The Uniform Distribution . . . . .	93
6.4.2	The Normal Distribution . . . . .	94
6.4.3	Visualization of the Standard Deviation . . . . .	96
6.4.4	Standardization of Random Variables . . . . .	96
6.4.5	Realizations of a Normal Distribution . . . . .	96
6.4.6	The Multivariate Normal Distribution . . . . .	97
6.4.7	The Bivariate Normal Distribution with Mean Zero and Zero Covariances $\sigma_{12} = \sigma_{21} = 0$ . . . . .	101
6.4.8	The Bivariate Normal Distribution with Mean Zero and Negative Covariances $\sigma_{12} = \sigma_{21} = -4$ . . . . .	101
6.5	Cholesky Decomposition and Positive Definite Matrices . . . . .	101
6.6	Maximum Likelihood Estimation: Multivariate Normal Distribution . . . . .	104
6.7	Introduction to Gaussian Processes . . . . .	104
6.7.1	Gaussian Process Prior . . . . .	105
6.7.2	Covariance Function . . . . .	105
6.7.3	Construction of the Covariance Matrix . . . . .	107
6.7.4	Generation of Random Samples and Plotting the Realizations of the Random Function . . . . .	109
6.7.5	Properties of the 1d Example . . . . .	111
6.8	Kriging: Modeling Basics . . . . .	114
6.8.1	The Kriging Idea in a Nutshell . . . . .	114
6.8.2	The Kriging Basis Function . . . . .	114
6.8.3	The Correlation Coefficient . . . . .	115
6.8.4	Covariance Matrix and Correlation Matrix . . . . .	116
6.8.5	The Kriging Model . . . . .	117
6.8.6	Correlations . . . . .	118
6.8.7	The Condition Number . . . . .	122
6.8.8	MLE to estimate $\theta$ and $p$ . . . . .	123
6.8.9	Tuning $\theta$ and $p$ . . . . .	124
6.9	Kriging Prediction . . . . .	124
6.9.1	The Augmented Correlation Matrix . . . . .	124
6.9.2	Properties of the Predictor . . . . .	125
6.10	Kriging Example: Sinusoid Function . . . . .	125
6.10.1	Calculating the Correlation Matrix $\Psi$ . . . . .	125
6.10.2	Computing the $\psi$ Vector . . . . .	126
6.10.3	Predicting at New Locations . . . . .	127

6.10.4	Visualization . . . . .	127
6.11	Cholesky Example With Two Points . . . . .	128
6.11.1	Cholesky Decomposition . . . . .	128
6.11.2	Computation of the Inverse Matrix . . . . .	129
6.12	Jupyter Notebook . . . . .	130
<b>7</b>	<b>Introduction to spotPython</b>	<b>131</b>
7.1	Example: Spot and the Sphere Function . . . . .	131
7.1.1	The Objective Function: Sphere . . . . .	132
7.1.2	The Spot Method as an Optimization Algorithm Using a Surrogate Model	133
7.2	Spot Parameters: <code>fun_evals</code> , <code>init_size</code> and <code>show_models</code> . . . . .	135
7.3	Print the Results . . . . .	136
7.4	Show the Progress . . . . .	136
7.5	Visualizing the Optimization and Hyperparameter Tuning Process with TensorBoard . . . . .	137
7.6	Jupyter Notebook . . . . .	140
<b>8</b>	<b>Multi-dimensional Functions</b>	<b>141</b>
8.1	Example: Spot and the 3-dim Sphere Function . . . . .	141
8.1.1	The Objective Function: 3-dim Sphere . . . . .	141
8.1.2	Results . . . . .	142
8.1.3	A Contour Plot . . . . .	143
8.1.4	TensorBoard . . . . .	145
8.2	Conclusion . . . . .	146
8.3	Exercises . . . . .	146
8.3.1	The Three Dimensional <code>fun_cubed</code> . . . . .	146
8.3.2	The Ten Dimensional <code>fun_wing_wt</code> . . . . .	147
8.3.3	The Three Dimensional <code>fun_runge</code> . . . . .	147
8.3.4	The Three Dimensional <code>fun_linear</code> . . . . .	147
8.4	Jupyter Notebook . . . . .	148
<b>9</b>	<b>Isotropic and Anisotropic Kriging</b>	<b>149</b>
9.1	Example: Isotropic Spot Surrogate and the 2-dim Sphere Function . . . . .	149
9.1.1	The Objective Function: 2-dim Sphere . . . . .	149
9.1.2	Results . . . . .	150
9.2	Example With Anisotropic Kriging . . . . .	150
9.2.1	Taking a Look at the <code>theta</code> Values . . . . .	153
9.3	Exercises . . . . .	154
9.3.1	<code>fun_branin</code> . . . . .	154
9.3.2	<code>fun_sin_cos</code> . . . . .	155
9.3.3	<code>fun_runge</code> . . . . .	155
9.3.4	<code>fun_wingwt</code> . . . . .	155
9.4	Jupyter Notebook . . . . .	156

<b>10 Using sklearn Surrogates in spotPython</b>	<b>157</b>
10.1 Example: Branin Function with <code>spotPython</code> 's Internal Kriging Surrogate . . . . .	157
10.1.1 The Objective Function Branin . . . . .	157
10.1.2 Running the surrogate model based optimizer <code>Spot</code> : . . . . .	158
10.1.3 TensorBoard . . . . .	159
10.1.4 Print the Results . . . . .	160
10.1.5 Show the Progress and the Surrogate . . . . .	160
10.2 Example: Using Surrogates From <code>scikit-learn</code> . . . . .	161
10.2.1 <code>GaussianProcessRegressor</code> as a Surrogate . . . . .	162
10.3 Example: One-dimensional Sphere Function With <code>spotPython</code> 's Kriging . . . . .	163
10.3.1 Results . . . . .	169
10.4 Example: <code>Sklearn</code> Model <code>GaussianProcess</code> . . . . .	170
10.5 Exercises . . . . .	176
10.5.1 <code>DecisionTreeRegressor</code> . . . . .	176
10.5.2 <code>RandomForestRegressor</code> . . . . .	176
10.5.3 <code>linear_model.LinearRegression</code> . . . . .	176
10.5.4 <code>linear_model.Ridge</code> . . . . .	177
10.6 Exercise 2 . . . . .	177
<b>11 Sequential Parameter Optimization: Gaussian Process Models</b>	<b>178</b>
11.1 Gaussian Processes Regression: Basic Introductory <code>scikit-learn</code> Example . . . . .	178
11.1.1 Train and Test Data . . . . .	178
11.1.2 Building the Surrogate With <code>Sklearn</code> . . . . .	179
11.1.3 Plotting the <code>SklearnModel</code> . . . . .	179
11.1.4 The <code>spotPython</code> Version . . . . .	180
11.1.5 Visualizing the Differences Between the <code>spotPython</code> and the <code>sklearn</code> Model Fits . . . . .	181
11.2 Exercises . . . . .	182
11.2.1 Schonlau Example Function . . . . .	182
11.2.2 Forrester Example Function . . . . .	182
11.2.3 <code>fun_runge</code> Function (1-dim) . . . . .	183
11.2.4 <code>fun_cubed</code> (1-dim) . . . . .	184
11.2.5 The Effect of Noise . . . . .	184
<b>12 Expected Improvement</b>	<b>185</b>
12.1 Example: <code>Spot</code> and the 1-dim Sphere Function . . . . .	185
12.1.1 The Objective Function: 1-dim Sphere . . . . .	185
12.1.2 Results . . . . .	186
12.2 Same, but with EI as <code>infill_criterion</code> . . . . .	187
12.3 Non-isotropic Kriging . . . . .	189
12.4 Using <code>sklearn</code> Surrogates . . . . .	192
12.4.1 The <code>spot</code> Loop . . . . .	192
12.4.2 <code>spot</code> : The Initial Model . . . . .	193

12.4.3	Init: Build Initial Design . . . . .	194
12.4.4	Evaluate . . . . .	197
12.4.5	Build Surrogate . . . . .	197
12.4.6	A Simple Predictor . . . . .	197
12.5	Gaussian Processes regression: basic introductory example . . . . .	197
12.6	The Surrogate: Using scikit-learn models . . . . .	200
12.7	Additional Examples . . . . .	203
12.7.1	Optimize on Surrogate . . . . .	205
12.7.2	Evaluate on Real Objective . . . . .	205
12.7.3	Impute / Infill new Points . . . . .	205
12.8	Tests . . . . .	205
12.9	EI: The Famous Schonlau Example . . . . .	207
12.10	EI: The Forrester Example . . . . .	209
12.11	Noise . . . . .	212
12.12	Cubic Function . . . . .	215
12.13	Factors . . . . .	221
<b>13</b>	<b>Handling Noise</b>	<b>223</b>
13.1	Example: Spot and the Noisy Sphere Function . . . . .	223
13.1.1	The Objective Function: Noisy Sphere . . . . .	223
13.1.2	Reproducibility: Noise Generation and Seed Handling . . . . .	225
13.2	spotPython’s Noise Handling Approaches . . . . .	227
13.3	Print the Results . . . . .	233
13.4	Noise and Surrogates: The Nugget Effect . . . . .	233
13.4.1	The Noisy Sphere . . . . .	233
13.5	Exercises . . . . .	236
13.5.1	Noisy <code>fun_cubed</code> . . . . .	236
13.5.2	<code>fun_runge</code> . . . . .	237
13.5.3	<code>fun_forrester</code> . . . . .	237
13.5.4	<code>fun_xsin</code> . . . . .	237
<b>14</b>	<b>Optimal Computational Budget Allocation in Spot</b>	<b>238</b>
14.1	Example: Spot, OCBA, and the Noisy Sphere Function . . . . .	238
14.1.1	The Objective Function: Noisy Sphere . . . . .	238
14.2	Print the Results . . . . .	245
14.3	Noise and Surrogates: The Nugget Effect . . . . .	245
14.3.1	The Noisy Sphere . . . . .	245
14.4	Exercises . . . . .	248
14.4.1	Noisy <code>fun_cubed</code> . . . . .	248
14.4.2	<code>fun_runge</code> . . . . .	249
14.4.3	<code>fun_forrester</code> . . . . .	249
14.4.4	<code>fun_xsin</code> . . . . .	249

<b>III Hyperparameter Tuning with PyTorch Lightning</b>	<b>250</b>
<b>15 HPT PyTorch Lightning: Diabetes</b>	<b>251</b>
15.1 Step 1: Setup . . . . .	251
15.2 Step 2: Initialization of the <code>fun_control</code> Dictionary . . . . .	252
15.3 Step 3: Loading the Diabetes Data Set . . . . .	253
15.4 Step 4: Preprocessing . . . . .	254
15.5 Step 5: Select the Core Model ( <code>algorithm</code> ) and <code>core_model_hyper_dict</code> . . . . .	254
15.6 Step 6: Modify <code>hyper_dict</code> Hyperparameters for the Selected Algorithm aka <code>core_model</code> . . . . .	255
15.7 Step 7: Data Splitting, the Objective (Loss) Function and the Metric . . . . .	256
15.7.1 Evaluation . . . . .	256
15.7.2 Loss Function . . . . .	257
15.7.3 Metric . . . . .	257
15.8 Step 8: Calling the SPOT Function . . . . .	257
15.8.1 Preparing the SPOT Call . . . . .	257
15.8.2 The Objective Function <code>fun</code> . . . . .	258
15.8.3 Showing the <code>fun_control</code> Dictionary . . . . .	258
15.8.4 Starting the Hyperparameter Tuning . . . . .	262
15.9 Step 9: Tensorboard . . . . .	265
15.10 Step 10: Results . . . . .	265
15.10.1 Get the Tuned Architecture . . . . .	266
15.10.2 Parallel Coordinates Plot . . . . .	273
15.10.3 Cross Validation With Lightning . . . . .	274
15.10.4 Plot all Combinations of Hyperparameters . . . . .	275
15.10.5 Visualizing the Activation Distribution (Under Development) . . . . .	275
<b>16 HPT PyTorch Lightning: Diabetes Using a Recurrent Neural Network</b>	<b>277</b>
16.1 Step 1: Setup . . . . .	277
16.2 Step 2: Initialization of the <code>fun_control</code> Dictionary . . . . .	278
16.3 Step 3: Loading the Diabetes Data Set . . . . .	279
16.4 Step 4: Preprocessing . . . . .	280
16.5 Step 5: Select the Core Model ( <code>algorithm</code> ) and <code>core_model_hyper_dict</code> . . . . .	281
16.6 Step 6: Modify <code>hyper_dict</code> Hyperparameters for the Selected Algorithm aka <code>core_model</code> . . . . .	282
16.7 Step 7: Data Splitting, the Objective (Loss) Function and the Metric . . . . .	283
16.7.1 Evaluation . . . . .	283
16.7.2 Loss Function . . . . .	283
16.7.3 Metric . . . . .	284
16.8 Step 8: Calling the SPOT Function . . . . .	284
16.8.1 Preparing the SPOT Call . . . . .	284
16.8.2 The Objective Function <code>fun</code> . . . . .	284
16.8.3 Showing the <code>fun_control</code> Dictionary . . . . .	285

16.8.4 Starting the Hyperparameter Tuning . . . . .	288
16.9 Step 9: Tensorboard . . . . .	293
16.10 Step 10: Results . . . . .	293
16.10.1 Get the Tuned Architecture . . . . .	295
16.10.2 Parallel Coordinates Plot . . . . .	297
16.10.3 Cross Validation With Lightning . . . . .	297
16.10.4 Plot all Combinations of Hyperparameters . . . . .	298
16.10.5 Visualizing the Activation Distribution (Under Development) . . . . .	299

## **Appendices** **301**

<b>A Introduction to Jupyter Notebook</b>	<b>301</b>
A.1 Different Notebook cells . . . . .	301
A.1.1 Code cells . . . . .	301
A.1.2 Markdown cells . . . . .	301
A.1.3 Raw cells . . . . .	302
A.2 Install Packages . . . . .	302
A.3 Load Packages . . . . .	303
A.4 Functions in Python . . . . .	303
A.5 List of Useful Jupyter Notebook Shortcuts . . . . .	304
<b>B Git Introduction</b>	<b>306</b>
B.1 Learning Objectives . . . . .	306
B.2 Basics of Git . . . . .	306
B.2.1 Initializing a Repository: <code>git init</code> . . . . .	306
B.2.2 Ignoring Files: <code>.gitignore</code> . . . . .	307
B.2.3 Adding Changes to the Staging Area: <code>git add</code> . . . . .	307
B.2.4 Transferring Changes to Memory: <code>git commit</code> . . . . .	308
B.2.5 Check the Status of Your Repository: <code>git status</code> . . . . .	309
B.2.6 Review Your Repository's History: <code>git log</code> . . . . .	310
B.3 Branches (Timelines) . . . . .	310
B.3.1 Creating an Alternative Timeline: <code>git branch</code> . . . . .	310
B.3.2 The Pointer to the Current Branch: <code>HEAD</code> . . . . .	311
B.3.3 Switching to an Alternative Timeline: <code>git switch</code> . . . . .	311
B.3.4 Switching to an Alternative Timeline and Making Changes: <code>git checkout</code> . . . . .	311
B.3.5 The Difference Between <code>checkout</code> and <code>switch</code> . . . . .	312
B.4 Merging Branches and Resolving Conflicts . . . . .	314
B.4.1 <code>git merge</code> : Merging Two Timelines . . . . .	314
B.4.2 Resolving Conflicts When Merging . . . . .	315
B.4.3 <code>git revert</code> : Undoing Something . . . . .	316
B.5 Downloading from GitLab . . . . .	318

B.6	Advanced . . . . .	319
B.6.1	git rebase: Moving the Base of a Branch . . . . .	319
B.7	Exercises . . . . .	321
B.7.1	Create project folder . . . . .	322
B.8	Initialize repo . . . . .	322
B.8.1	Do not upload / ignore certain file types . . . . .	322
B.8.2	Create file and stage it . . . . .	322
B.8.3	Create another file and check status . . . . .	322
B.8.4	Commit changes . . . . .	322
B.8.5	Create a new branch and switch to it . . . . .	323
B.8.6	Commit changes in the new branch . . . . .	323
B.8.7	Merge branch into main . . . . .	323
B.8.8	Resolve merge conflict . . . . .	323
<b>C</b>	<b>Python Introduction</b>	<b>324</b>
C.1	Recommendations . . . . .	324
<b>D</b>	<b>Documentation of the Sequential Parameter Optimization</b>	<b>325</b>
D.1	Example: spot . . . . .	325
D.1.1	The Objective Function . . . . .	325
D.1.2	External Parameters . . . . .	327
D.2	The fun_control Dictionary . . . . .	330
D.3	The design_control Dictionary . . . . .	330
D.4	The surrogate_control Dictionary . . . . .	331
D.5	The optimizer_control Dictionary . . . . .	331
D.6	Run . . . . .	332
D.7	Print the Results . . . . .	332
D.8	Show the Progress . . . . .	332
D.9	Visualize the Surrogate . . . . .	333
D.10	Run With a Specific Start Design . . . . .	333
D.11	Init: Build Initial Design . . . . .	334
D.12	Replicability . . . . .	335
D.13	Surrogates . . . . .	336
D.13.1	A Simple Predictor . . . . .	336
D.14	Demo/Test: Objective Function Fails . . . . .	337
D.15	PyTorch: Detailed Description of the Data Splitting . . . . .	339
D.15.1	Description of the "train_hold_out" Setting . . . . .	339
<b>R</b>	<b>References</b>	<b>350</b>

# Preface

This document provides a comprehensive guide to hyperparameter tuning using spotPython for scikit-learn, scipy-optimize, River, and PyTorch. The first part introduces fundamental ideas from optimization. The second part discusses numerical issues and introduces spotPython’s surrogate model-based optimization process. The thirs part focuses on hyperparameter tuning. Several case studies are presented, including hyperparameter tuning for sklearn models such as Support Vector Classification, Random Forests, Gradient Boosting (XGB), and K-nearest neighbors (KNN), as well as a Hoeffding Adaptive Tree Regressor from river. The integration of spotPython into the PyTorch and PyTorch Lightning training workflow is also discussed. With a hands-on approach and step-by-step explanations, this cookbook serves as a practical starting point for anyone interested in hyperparameter tuning with Python. Highlights include the interplay between Tensorboard, PyTorch Lightning, spotPython, spotRiver, and River. This publication is under development, with updates available on the corresponding webpage.

! Important: This book is still under development.

The most recent version of this book is available at <https://sequential-parameter-optimization.github.io/Hyperparameter-Tuning-Cookbook/>

## Book Structure

This document is structured in three parts. The first part presents an introduction to optimization. The second part describes numerical methods, and the third part presents hyperparameter tuning.

💡 Hyperparameter Tuning Reference

- The open access book Bartz et al. (2022) provides a comprehensive overview of hyperparameter tuning. It can be downloaded from <https://link.springer.com/book/10.1007/978-981-19-5170-1>.

### Note

The `.ipynb` notebook (Bartz-Beielstein 2023) is updated regularly and reflects updates and changes in the `spotPython` package. It can be downloaded from [https://github.com/sequential-parameter-optimization/spotPython/blob/main/notebooks/14\\_spot\\_ray\\_hpt\\_torch\\_cifar10.ipynb](https://github.com/sequential-parameter-optimization/spotPython/blob/main/notebooks/14_spot_ray_hpt_torch_cifar10.ipynb).

## Software Used in this Book

`scikit-learn` is a Python module for machine learning built on top of SciPy and is distributed under the 3-Clause BSD license. The project was started in 2007 by David Cournapeau as a Google Summer of Code project, and since then many volunteers have contributed.

`PyTorch` is an optimized tensor library for deep learning using GPUs and CPUs. `Lightning` is a lightweight PyTorch wrapper for high-performance AI research. It allows you to decouple the research from the engineering.

`River` is a Python library for online machine learning. It is designed to be used in real-world environments, where not all data is available at once, but streaming in.

`spotPython` (“Sequential Parameter Optimization Toolbox in Python”) is the Python version of the well-known hyperparameter tuner SPOT, which has been developed in the R programming environment for statistical analysis for over a decade. The related open-access book is available here: [Hyperparameter Tuning for Machine and Deep Learning with R—A Practical Guide](#).

`spotRiver` provides an interface between `spotPython` and `River`.

## Citation

If this document has been useful to you and you wish to cite it in a scientific publication, please refer to the following paper, which can be found on arXiv: <https://arxiv.org/abs/2307.10262>.

```
@ARTICLE{bart23iArXiv,
    author = {{Bartz-Beielstein}, Thomas},
    title = "{Hyperparameter Tuning Cookbook:
        A guide for scikit-learn, PyTorch, river, and spotPython}",
    journal = {arXiv e-prints},
    keywords = {Computer Science - Machine Learning,
        Computer Science - Artificial Intelligence, 90C26, I.2.6, G.1.6},
    year = 2023,
```

```
month = jul,
    eid = {arXiv:2307.10262},
    pages = {arXiv:2307.10262},
    doi = {10.48550/arXiv.2307.10262},
archivePrefix = {arXiv},
    eprint = {2307.10262},
primaryClass = {cs.LG},
    adsurl = {https://ui.adsabs.harvard.edu/abs/2023arXiv230710262B},
    adsnote = {Provided by the SAO/NASA Astrophysics Data System}
}
```

# **Part I**

# **Optimization**

# 1 Introduction: Optimization

## 1.1 Optimization, Simulation, and Surrogate Modeling

- We will consider the interplay between
  - mathematical models,
  - numerical approximation,
  - simulation,
  - computer experiments, and
  - field data
- Experimental design will play a key role in our developments, but not in the classical regression and response surface methodology sense
- Challenging real-data/real-simulation examples benefiting from modern surrogate modeling methodology
- We will consider the classical, response surface methodology (RSM) approach, and then move on to more modern approaches
- All approaches are based on surrogates

## 1.2 Surrogates

- Gathering data is **expensive**, and sometimes getting exactly the data you want is impossible or unethical
- **Surrogate:** substitute for the real thing
- In statistics, draws from predictive equations derived from a fitted model can act as a surrogate for the data-generating mechanism
- Benefits of the surrogate approach:
  - Surrogate could represent a cheaper way to explore relationships, and entertain “what ifs?”
  - Surrogates favor faithful yet pragmatic reproduction of dynamics:
    - \* interpretation,
    - \* establishing causality, or
    - \* identification
  - Many numerical simulators are **deterministic**, whereas field observations are noisy or have measurement error

### 1.2.1 Costs of Simulation

- Computer simulations are generally cheaper (but not always!) than physical observation
- Some computer simulations can be just as expensive as field experimentation, but computer modeling is regarded as easier because:
  - the experimental apparatus is better understood
  - more aspects may be controlled.

### 1.2.2 Mathematical Models and Meta-Models

- Use of mathematical models leveraging numerical solvers has been commonplace for some time
- Mathematical models became more complex, requiring more resources to simulate/solve numerically
- Practitioners increasingly relied on **meta-models** built off of limited simulation campaigns

### 1.2.3 Surrogates = Trained Meta-models

- Data collected via expensive computer evaluations tuned flexible functional forms that could be used in lieu of further simulation to
  - save money or computational resources;
  - cope with an inability to perform future runs (expired licenses, off-line or over-impacted supercomputers)
- Trained meta-models became known as **surrogates**

### 1.2.4 Computer Experiments

- **Computer experiment:** design, running, and fitting meta-models.
  - Like an ordinary statistical experiment, except the data are generated by computer codes rather than physical or field observations, or surveys
- **Surrogate modeling** is statistical modeling of computer experiments

### 1.2.5 Limits of Mathematical Modeling

- Mathematical biologists, economists and others had reached the limit of equilibrium-based mathematical modeling with cute closed-form solutions
- **Stochastic simulations replace deterministic solvers** based on FEM, Navier–Stokes or Euler methods
- Agent-based simulation models are used to explore predator-prey (Lotka–Volterra) dynamics, spread of disease, management of inventory or patients in health insurance markets
- Consequence: the distinction between surrogate and statistical model is all but gone

### 1.2.6 Example: Why Computer Simulations are Necessary

- You can't seed a real community with Ebola and watch what happens
- If there's (real) field data, say on a historical epidemic, further experimentation may be almost entirely limited to the mathematical and computer modeling side
- Classical statistical methods offer little guidance

### 1.2.7 Simulation Requirements

- Simulation should
  - enable rich **diagnostics** to help criticize that models
  - **understanding** its sensitivity to inputs and other configurations
  - providing the ability to **optimize** and
  - refine both **automatically** and with expert intervention
- And it has to do all that while remaining **computationally tractable**
- One perspective is so-called **response surface methods** (RSMs):
  - a poster child from industrial statistics' heyday, well before information technology became a dominant industry

#### ! Goals

- How to choose models and optimizers for solving real-world problems
- How to use simulation to understand and improve processes

## 1.3 Jupyter Notebook

### Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

# 2 Aircraft Wing Weight Example

## 2.1 AWWE Equation

- Example from Forrester et al.
- Understand the **weight** of an unpainted light aircraft wing as a function of nine design and operational parameters:

$$W = 0.036 S_W^{0.758} \times W_{fw}^{0.0035} \left( \frac{A}{\cos^2 \Lambda} \right)^{0.6} \times q^{0.006} \times \lambda^{0.04} \\ \times \left( \frac{100 R_{tc}}{\cos \Lambda} \right)^{-0.3} \times (N_z W_{dg})^{0.49}$$

## 2.2 AWWE Parameters and Equations (Part 1)

Table 2.1: Aircraft Wing Weight Parameters

Symbol	Parameter	Baseline	Minimum	Maximum
$S_W$	Wing area ( $ft^2$ )	174	150	200
$W_{fw}$	Weight of fuel in wing (lb)	252	220	300
$A$	Aspect ratio	7.52	6	10
$\Lambda$	Quarter-chord sweep (deg)	0	-10	10
$q$	Dynamic pressure at cruise ( $lb/ft^2$ )	34	16	45
$\lambda$	Taper ratio	0.672	0.5	1
$R_{tc}$	Aerofoil thickness to chord ratio	0.12	0.08	0.18
$N_z$	Ultimate load factor	3.8	2.5	6
$W_{dg}$	Flight design gross weight (lb)	2000	1700	2500
$W_p$	paint weight ( $lb/ft^2$ )	0.064	0.025	0.08

The study begins with a baseline Cessna C172 Skyhawk Aircraft as its reference point. It aims to investigate the impact of wing area and fuel weight on the overall weight of the aircraft. Two crucial parameters in this analysis are the aspect ratio ( $A$ ), defined as the ratio of the

wing's length to the average chord (thickness of the airfoil), and the taper ratio ( $\lambda$ ), which represents the ratio of the maximum to the minimum thickness of the airfoil or the maximum to minimum chord.

It's important to note that the equation used in this context is not a computer simulation but will be treated as one for the purpose of illustration. This approach involves employing a true mathematical equation, even if it's considered unknown, as a useful tool for generating realistic settings to test the methodology. The functional form of this equation was derived by "calibrating" known physical relationships to curves obtained from existing aircraft data, as referenced in Raymer 2012. Essentially, it acts as a surrogate for actual measurements of aircraft weight.

Examining the mathematical properties of the AWWE (Aircraft Weight With Wing Area and Fuel Weight Equation), it is evident that the response is highly nonlinear concerning its inputs. While it's common to apply the logarithm to simplify equations with complex exponents, even when modeling the logarithm, which transforms powers into slope coefficients and products into sums, the response remains nonlinear due to the presence of trigonometric terms. Given the combination of nonlinearity and high input dimension, simple linear and quadratic response surface approximations are likely to be inadequate for this analysis.

## 2.3 Goals: Understanding and Optimization

The primary goals of this study revolve around understanding and optimization:

1. **Understanding:** One of the straightforward objectives is to gain a deep understanding of the input-output relationships in this context. Given the global perspective implied by this setting, it becomes evident that a more sophisticated model is almost necessary. At this stage, let's focus on this specific scenario to establish a clear understanding.
2. **Optimization:** Another application of this analysis could be optimization. There may be an interest in minimizing the weight of the aircraft, but it's likely that there will be constraints in place. For example, the presence of wings with a nonzero area is essential for the aircraft to be capable of flying. In situations involving (constrained) optimization, a global perspective and, consequently, the use of flexible modeling are vital.

The provided Python code serves as a genuine computer implementation that "solves" a mathematical model. It accepts arguments encoded in the unit cube, with defaults used to represent baseline settings, as detailed in the table labeled as Table 2.1. To map values from the interval  $[a, b]$  to the interval  $[0, 1]$ , the following formula can be employed:

$$y = f(x) = \frac{x - a}{b - a}.$$

To reverse this mapping and obtain the original values, the formula

$$g(y) = a + (b - a)y$$

can be used.

```
import numpy as np

def wingwt(Sw=0.48, Wfw=0.4, A=0.38, L=0.5, q=0.62, l=0.344, Rtc=0.4, Nz=0.37, Wdg=0.38):
    # put coded inputs back on natural scale
    Sw = Sw * (200 - 150) + 150
    Wfw = Wfw * (300 - 220) + 220
    A = A * (10 - 6) + 6
    L = (L * (10 - (-10)) - 10) * np.pi/180
    q = q * (45 - 16) + 16
    l = l * (1 - 0.5) + 0.5
    Rtc = Rtc * (0.18 - 0.08) + 0.08
    Nz = Nz * (6 - 2.5) + 2.5
    Wdg = Wdg*(2500 - 1700) + 1700
    # calculation on natural scale
    W = 0.036 * Sw**0.758 * Wfw**0.0035 * (A/np.cos(L)**2)**0.6 * q**0.006
    W = W * l**0.04 * (100*Rtc/np.cos(L))**(-0.3) * (Nz*Wdg)**(0.49)
    return(W)
```

## 2.4 Properties of the Python “Solver”

The compute time required by the “wingwt” solver is extremely short and can be considered trivial in terms of computational resources. The approximation error is exceptionally small, effectively approaching machine precision, which indicates the high accuracy of the solver’s results.

To simulate time-consuming evaluations, a deliberate delay is introduced by incorporating a `sleep(3600)` command, which effectively synthesizes a one-hour execution time for a particular evaluation.

Moving on to the AWWE visualization, plotting in two dimensions is considerably simpler than dealing with nine dimensions. To aid in creating visual representations, the code provided below establishes a grid within the unit square to facilitate the generation of sliced visuals. This involves generating a “meshgrid” as outlined in the code.

```

import numpy as np
x = np.linspace(0, 1, 3)
y = np.linspace(0, 1, 3)
X, Y = np.meshgrid(x, y)
zp = zip(np.ravel(X), np.ravel(Y))
list(zp)

```

```

[(0.0, 0.0),
 (0.5, 0.0),
 (1.0, 0.0),
 (0.0, 0.5),
 (0.5, 0.5),
 (1.0, 0.5),
 (0.0, 1.0),
 (0.5, 1.0),
 (1.0, 1.0)]

```

The coding used to transform inputs from natural units is largely a matter of taste, so long as it's easy to undo for reporting back on original scales

```

%matplotlib inline
import matplotlib.pyplot as plt
# plt.style.use('seaborn-white')
import numpy as np
x = np.linspace(0, 1, 100)
y = np.linspace(0, 1, 100)
X, Y = np.meshgrid(x, y)

```

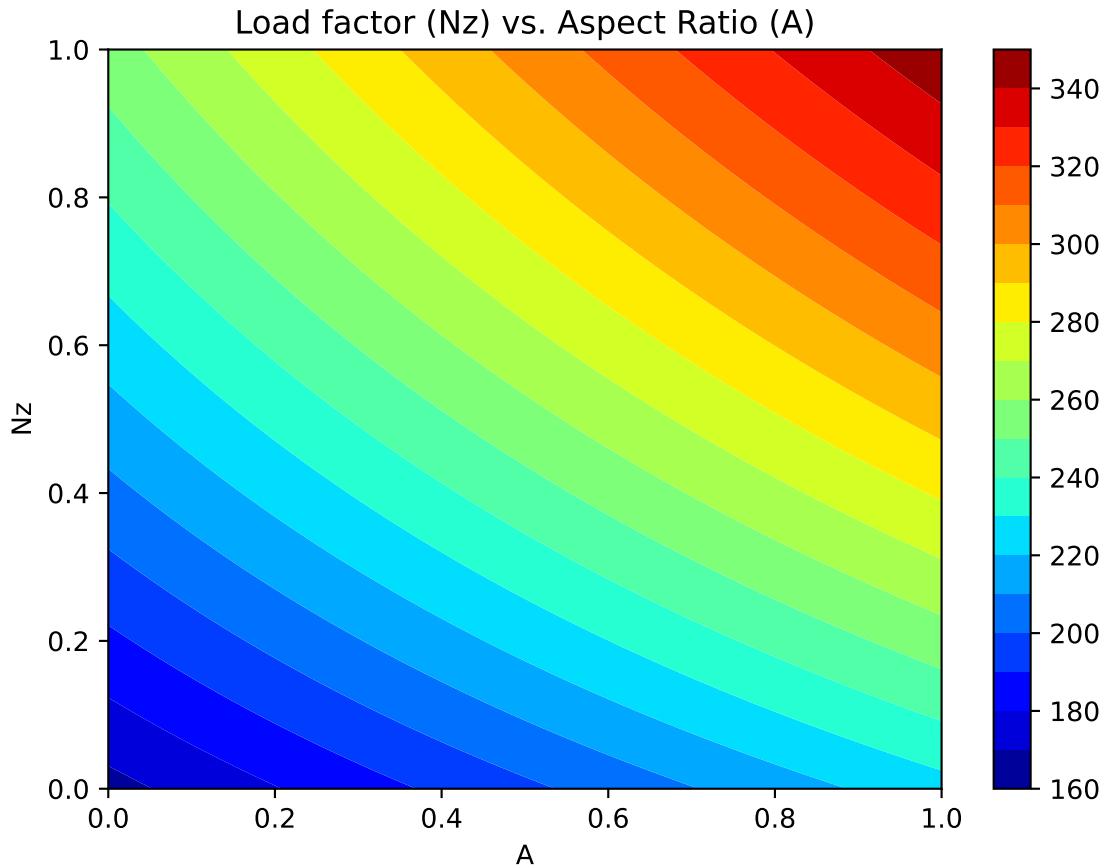
## 2.5 Plot 1: Load Factor ( $N_z$ ) and Aspect Ratio ( $A$ )

We will vary  $N_z$  and  $A$ , with other inputs fixed at their baseline values.

```

z = wingwt(A = X, Nz = Y)
fig = plt.figure(figsize=(7., 5.))
plt.contourf(X, Y, z, 20, cmap='jet')
plt.xlabel("A")
plt.ylabel("Nz")
plt.title("Load factor (Nz) vs. Aspect Ratio (A)")
plt.colorbar()
plt.show()

```



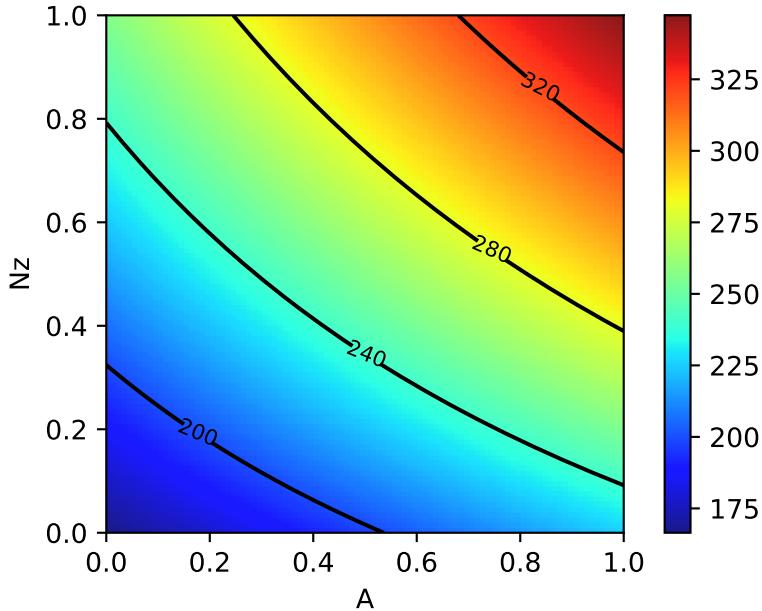
Contour plots can be refined, e.g., by adding explicit contour lines as shown in the following figure.

```

contours = plt.contour(X, Y, z, 4, colors='black')
plt.clabel(contours, inline=True, fontsize=8)
plt.xlabel("A")
plt.ylabel("Nz")

plt.imshow(z, extent=[0, 1, 0, 1], origin='lower',
           cmap='jet', alpha=0.9)
plt.colorbar()

```



The interpretation of the AWWE plot can be summarized as follows:

- The figure displays the weight response as a function of two variables,  $N_z$  and  $A$ , using an image-contour plot.
- The slight curvature observed in the contours suggests an interaction between these two variables.
- Notably, the range of outputs depicted in the figure, spanning from approximately 160 to 320, nearly encompasses the entire range of outputs observed from various input settings within the full 9-dimensional input space.
- The plot indicates that aircraft wings tend to be heavier when the aspect ratios ( $A$ ) are high.
- This observation aligns with the idea that wings are designed to withstand and accommodate high gravitational forces ( $g$ -forces, large  $N_z$ ), and there may be a compounding effect where larger values of  $N_z$  contribute to increased wing weight.
- It's plausible that this phenomenon is related to the design considerations of fighter jets, which cannot have the efficient and lightweight glider-like wings typically found in other types of aircraft.

## 2.6 Plot 2: Taper Ratio and Fuel Weight

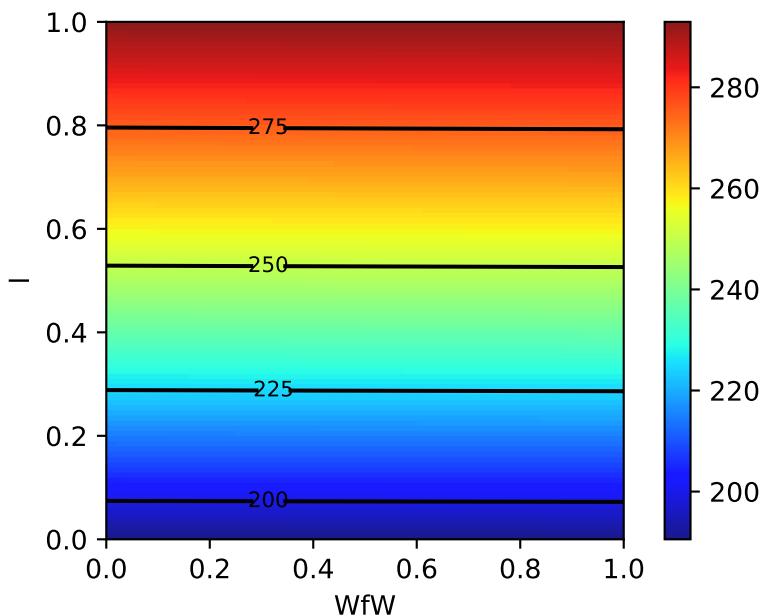
- The same experiment for two other inputs, e.g., taper ratio  $\lambda$  and fuel weight  $W_{fw}$

```

z = wingwt(Wfw = X, Nz = Y)
contours = plt.contour(X, Y, z, 4, colors='black')
plt.clabel(contours, inline=True, fontsize=8)
plt.xlabel("Wfw")
plt.ylabel("l")

plt.imshow(z, extent=[0, 1, 0, 1], origin='lower',
           cmap='jet', alpha=0.9)
plt.colorbar();

```



- Interpretation of Taper Ratio ( $l$ ) and Fuel Weight ( $W_{fw}$ )
  - Apparently, neither input has much effect on wing weight:
    - \* with  $\lambda$  having a marginally greater effect, covering less than 4 percent of the span of weights observed in the  $A \times N_z$  plane
  - There's no interaction evident in  $\lambda \times W_{fw}$

## 2.7 The Big Picture: Combining all Variables

```

pl = ["Sw", "Wfw", "A", "L", "q", "l", "Rtc", "Nz", "Wdg"]

```

```

import math

Z = []
Zlab = []
l = len(pl)
# lc = math.comb(l,2)
for i in range(l):
    for j in range(i+1, l):
        # for j in range(l):
        # print(pl[i], pl[j])
        d = {pl[i]: X, pl[j]: Y}
        Z.append(wingwt(**d))
        Zlab.append([pl[i],pl[j]])

```

Now we can generate all 36 combinations, e.g., our first example is combination  $p = 19$ .

```

p = 19
Zlab[p]

```

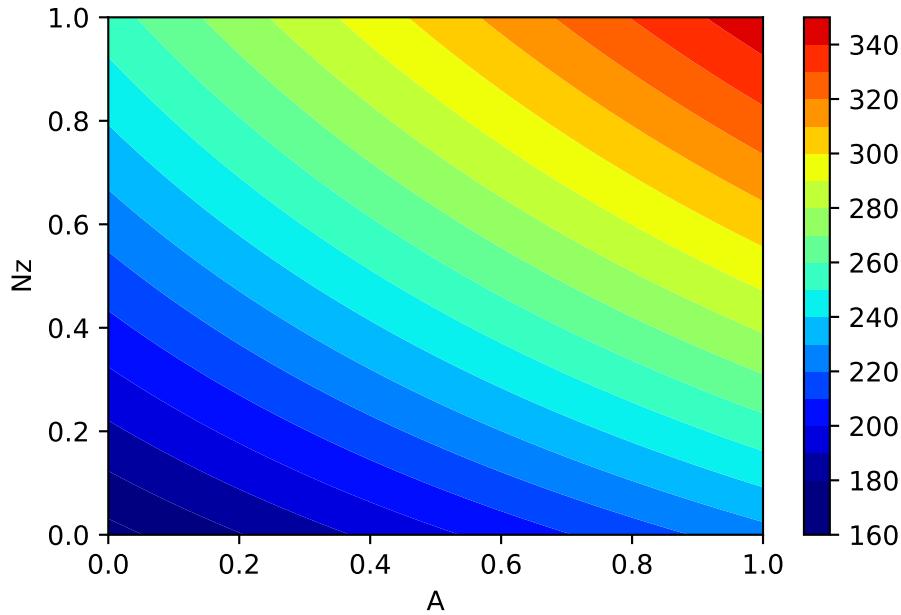
`['A', 'Nz']`

To help interpret outputs from experiments such as this one—to level the playing field when comparing outputs from other pairs of inputs—code below sets up a color palette that can be re-used from one experiment to the next. We use the arguments `vmin=180` and `vmax =360` to implement comparability

```

plt.contourf(X, Y, Z[p], 20, cmap='jet', vmin=180, vmax=360)
plt.xlabel(Zlab[p][0])
plt.ylabel(Zlab[p][1])
plt.colorbar()

```

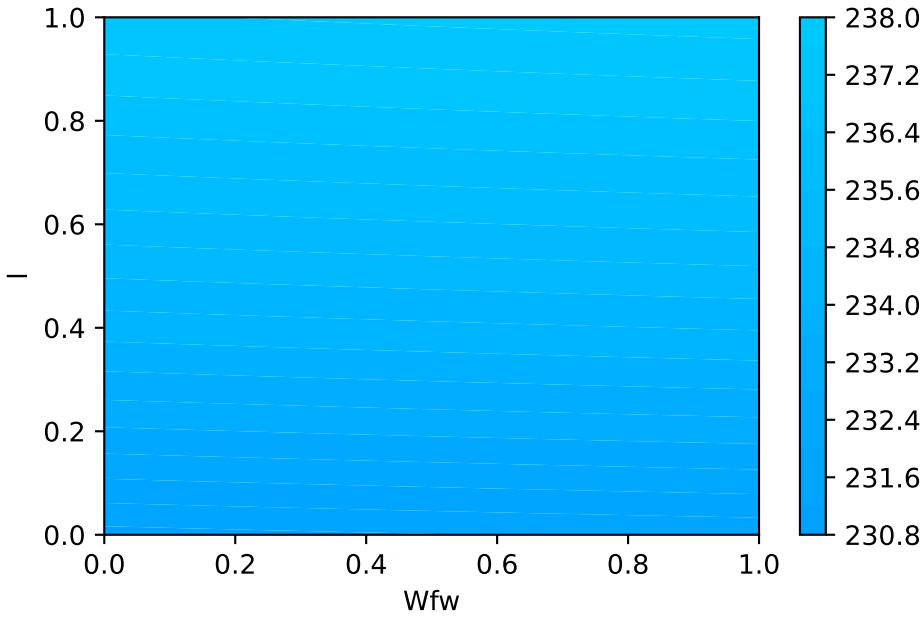


- Let's plot the second example, taper ratio  $\lambda$  and fuel weight  $W_{fw}$
- This is combination 11:

```
p = 11
Zlab[p]
```

```
['Wfw', '1']
```

```
plt.contourf(X, Y, Z[p], 20, cmap='jet', vmin=180, vmax=360)
plt.xlabel(Zlab[p][0])
plt.ylabel(Zlab[p][1])
plt.colorbar()
```



- Using a global colormap indicates that these variables have minor effects on the wing weight.
- Important factors can be detected by visual inspection
- Plotting the Big Picture: we can plot all 36 combinations in one figure.

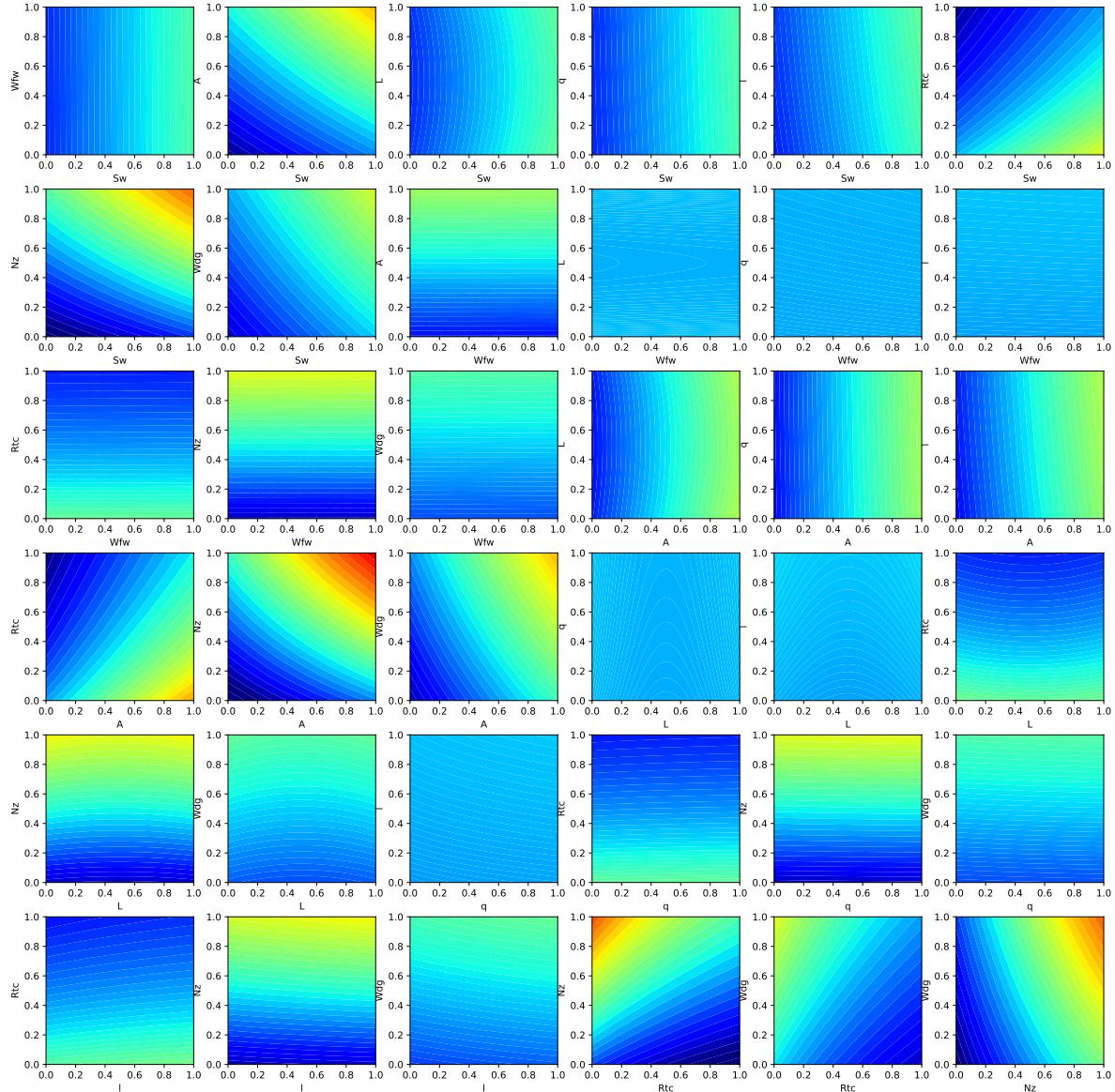
```

import matplotlib.pyplot as plt
from mpl_toolkits.axes_grid1 import ImageGrid
import numpy as np

fig = plt.figure(figsize=(20., 20.))
grid = ImageGrid(fig, 111, # similar to subplot(111)
                 nrows_ncols=(6,6), # creates 2x2 grid of axes
                 axes_pad=0.5, # pad between axes in inch.
                 share_all=True,
                 label_mode="0",
                 )
i = 0
for ax, im in zip(grid, Z):
    # Iterating over the grid returns the Axes.
    ax.set_xlabel(Zlab[i][0])
    ax.set_ylabel(Zlab[i][1])
    # ax.set_title(Zlab[i][1] + " vs. " + Zlab[i][0])
    ax.contourf(X, Y, im, 30, cmap = "jet", vmin = 180, vmax = 360)
    i = i + 1

```

```
plt.show()
```



## 2.8 AWWE Landscape

- Our Observations

1. The load factor  $N_z$ , which determines the magnitude of the maximum aerodynamic load on the wing, is very active and involved in interactions with other variables.
    - Classic example: the interaction of  $N_z$  with the aspect ratio  $A$  indicates a heavy wing for high aspect ratios and large  $g$ -forces
    - This is the reason why highly manoeuvrable fighter jets cannot have very efficient, glider wings)
  2. Aspect ratio  $A$  and airfoil thickness to chord ratio  $R_{tc}$  have nonlinear interactions.
  3. Most important variables:
    - Ultimate load factor  $N_z$ , wing area  $S_w$ , and flight design gross weight  $W_{dg}$ .
  4. Little impact: dynamic pressure  $q$ , taper ratio  $l$ , and quarter-chord sweep  $L$ .
- Expert Knowledge
    - Aircraft designers know that the overall weight of the aircraft and the wing area must be kept to a minimum
    - the latter usually dictated by constraints such as required stall speed, landing distance, turn rate, etc.

## 2.9 Summary of the First Experiments

- First, we considered two pairs of inputs, out of 36 total pairs
- Then, the “Big Picture”:
  - For each pair we evaluated `wingwt` 10,000 times
- Doing the same for all pairs would require 360K evaluations:
  - not a reasonable number with a real computer simulation that takes any non-trivial amount of time to evaluate
  - Only 1s per evaluation: > 100 hours
- Many solvers take minutes/hours/days to execute a single run
- And: three-way interactions?
- Consequence: a different strategy is needed

## 2.10 Exercise

### 2.10.1 Adding Paint Weight

- Paint weight is not considered.

- Add Paint Weight  $W_p$  to formula (the updated formula is shown below) and update the functions and plots in the notebook.

$$W = 0.036 S_W^{0.758} \times W_{fw}^{0.0035} \times \left( \frac{A}{\cos^2 \Lambda} \right)^{0.6} \times q^{0.006} \times \lambda^{0.04} \\ \times \left( \frac{100 R_{tc}}{\cos \Lambda} \right)^{-0.3} \times (N_z W_{dg})^{0.49} + S_w W_p$$

## 2.11 Jupyter Notebook

**i** Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

## 3 Introduction to `scipy.optimize`

[SciPy](#) provides algorithms for optimization, integration, interpolation, eigenvalue problems, algebraic equations, differential equations, statistics and many other classes of problems. SciPy is a collection of mathematical algorithms and convenience functions built on NumPy. It adds significant power to Python by providing the user with high-level commands and classes for manipulating and visualizing data.

[SciPy optimize](#) provides functions for minimizing (or maximizing) objective functions, possibly subject to constraints. It includes solvers for nonlinear problems (with support for both local and global optimization algorithms), linear programming, constrained and nonlinear least-squares, root finding, and curve fitting.

In this notebook, we will learn how to use the `scipy.optimize` module to solve optimization problems. See: <https://docs.scipy.org/doc/scipy/tutorial/optimize.html>

### Note

- This content is based on information from the `scipy.optimize` package.
- The `scipy.optimize` package provides several commonly used optimization algorithms. A detailed listing is available in `scipy.optimize` (can also be found by `help(scipy.optimize)`).

Common functions and objects, shared across different SciPy optimize solvers, are shown in Table 3.1.

Table 3.1: Common functions and objects, shared across different SciPy optimize solvers

Function or Object	Description
<code>show_options([solver, method, disp])</code>	Show documentation for additional options of optimization solvers.
<code>OptimizeResult</code>	Represents the optimization result.
<code>OptimizeWarning</code>	Warning issued by solvers.

We will introduce unconstrained minimization of multivariate scalar functions in this chapter. The `minimize` function provides a common interface to unconstrained and constrained minimization algorithms for multivariate scalar functions in `scipy.optimize`. To demonstrate

the minimization function, consider the problem of minimizing the Rosenbrock function of  $N$  variables:

$$f(J) = \sum_{i=1}^{N-1} 100(x_{i+1} - x_i^2)^2 + (1 - x_i)^2$$

The minimum value of this function is 0, which is achieved when ( $x_i = 1$ ).

Note that the Rosenbrock function and its derivatives are included in `scipy.optimize`. The implementations shown in the following sections provide examples of how to define an objective function as well as its Jacobian and Hessian functions. Objective functions in `scipy.optimize` expect a numpy array as their first parameter, which is to be optimized and must return a float value. The exact calling signature must be `f(x, *args)`, where `x` represents a numpy array, and `args` is a tuple of additional arguments supplied to the objective function.

## 3.1 Derivative-free Optimization Algorithms

Section 3.1.1 and Section 3.1.2 present two approaches that do not need gradient information to find the minimum. They use function evaluations to find the minimum.

### 3.1.1 Nelder-Mead Simplex Algorithm

`method='Nelder-Mead'`: In the example below, the `minimize` routine is used with the *Nelder-Mead* simplex algorithm (selected through the `method` parameter):

```
import numpy as np
from scipy.optimize import minimize

def rosen(x):
    """The Rosenbrock function"""
    return sum(100.0 * (x[1:] - x[:-1]**2.0)**2.0 + (1 - x[:-1])**2.0)

x0 = np.array([1.3, 0.7, 0.8, 1.9, 1.2])
res = minimize(rosen, x0, method='nelder-mead',
               options={'xtol': 1e-8, 'disp': True})

print(res.x)
```

```

Optimization terminated successfully.
    Current function value: 0.000000
    Iterations: 339
    Function evaluations: 571
[1. 1. 1. 1. 1.]

```

The simplex algorithm is probably the simplest way to minimize a well-behaved function. It requires only function evaluations and is a good choice for simple minimization problems. However, because it does not use any gradient evaluations, it may take longer to find the minimum.

### 3.1.2 Powell's Method

Another optimization algorithm that needs only function calls to find the minimum is *Powell's* method, which can be selected by setting the `method` parameter to '`'powell'`' in the `minimize` function.

To demonstrate how to supply additional arguments to an objective function, let's consider minimizing the Rosenbrock function with an additional scaling factor  $a$  and an offset  $b$ :

$$f(J, a, b) = \sum_{i=1}^{N-1} a(x_{i+1} - x_i^2)^2 + (1 - x_i)^2 + b$$

You can achieve this using the `minimize` routine with the example parameters  $a = 0.5$  and  $b = 1$ :

```

def rosen_with_args(x, a, b):
    """The Rosenbrock function with additional arguments"""
    return sum(a * (x[1:] - x[:-1]**2.0)**2.0 + (1 - x[:-1])**2.0) + b

x0 = np.array([1.3, 0.7, 0.8, 1.9, 1.2])
res = minimize(rosen_with_args, x0, method='nelder-mead',
               args=(0.5, 1.), options={'xtol': 1e-8, 'disp': True})

print(res.x)

```

```

Optimization terminated successfully.
    Current function value: 1.000000
    Iterations: 319
    Function evaluations: 525
[1.          1.          1.          1.          0.99999999]

```

As an alternative to using the `args` parameter of `minimize`, you can wrap the objective function in a new function that accepts only `x`. This approach is also useful when it is necessary to pass additional parameters to the objective function as keyword arguments.

```
def rosen_with_args(x, a, *, b): # b is a keyword-only argument
    return sum(a * (x[1:] - x[:-1]**2.0)**2.0 + (1 - x[:-1])**2.0) + b

def wrapped_rosen_without_args(x):
    return rosen_with_args(x, 0.5, b=1.) # pass in `a` and `b`

x0 = np.array([1.3, 0.7, 0.8, 1.9, 1.2])
res = minimize(wrapped_rosen_without_args, x0, method='nelder-mead',
               options={'xtol': 1e-8,})

print(res.x)
```

```
[1.          1.          1.          1.          0.99999999]
```

Another alternative is to use `functools.partial`.

```
from functools import partial

partial_rosen = partial(rosen_with_args, a=0.5, b=1.)
res = minimize(partial_rosen, x0, method='nelder-mead',
               options={'xtol': 1e-8,})

print(res.x)
```

```
[1.          1.          1.          1.          0.99999999]
```

## 3.2 Gradient-based optimization algorithms

### 3.2.1 An Introductory Example: Broyden-Fletcher-Goldfarb-Shanno Algorithm (BFGS)

This section introduces an optimization algorithm that uses gradient information to find the minimum. The Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm (selected by setting `method='BFGS'`) is an optimization algorithm that aims to converge quickly to the solution. This algorithm uses the gradient of the objective function. If the gradient is not provided by the user, it is estimated using first-differences. The BFGS method typically requires fewer function calls compared to the simplex algorithm, even when the gradient needs to be estimated.

## Example: BFGS

To demonstrate the BFGS algorithm, let's use the Rosenbrock function again. The gradient of the Rosenbrock function is a vector described by the following mathematical expression:

$$\frac{\partial f}{\partial x_j} = \sum_{i=1}^N 200(x_i - x_{i-1}^2)(\delta_{i,j} - 2x_{i-1}\delta_{i-1,j}) - 2(1 - x_{i-1})\delta_{i-1,j} \quad (3.1)$$

$$= 200(x_j - x_{j-1}^2) - 400x_j(x_{j+1} - x_j^2) - 2(1 - x_j) \quad (3.2)$$

This expression is valid for interior derivatives, but special cases are:

$$\frac{\partial f}{\partial x_0} = -400x_0(x_1 - x_0^2) - 2(1 - x_0)$$

$$\frac{\partial f}{\partial x_{N-1}} = 200(x_{N-1} - x_{N-2}^2)$$

Here's a Python function that computes this gradient:

```
def rosen_der(x):
    xm = x[1:-1]
    xm_m1 = x[:-2]
    xm_p1 = x[2:]
    der = np.zeros_like(x)
    der[1:-1] = 200*(xm-xm_m1**2) - 400*(xm_p1 - xm**2)*xm - 2*(1-xm)
    der[0] = -400*x[0]*(x[1]-x[0]**2) - 2*(1-x[0])
    der[-1] = 200*(x[-1]-x[-2]**2)
    return der
```

You can specify this gradient information in the minimize function using the jac parameter as illustrated below:

```
res = minimize(rosen, x0, method='BFGS', jac=rosen_der,
               options={'disp': True})
```

```
print(res.x)
```

```
Optimization terminated successfully.
      Current function value: 0.000000
      Iterations: 25
      Function evaluations: 30
      Gradient evaluations: 30
[1.00000004 1.0000001  1.00000021  1.00000044  1.00000092]
```

### 3.2.2 Background and Basics for Gradient-based Optimization

#### 3.2.3 Gradient

The gradient  $\nabla f(J)$  for a scalar function  $f(J)$  with  $n$  different variables is defined by its partial derivatives:

$$\nabla f(J) = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

#### 3.2.4 Jacobian Matrix

The Jacobian matrix  $J(J)$  for a vector-valued function  $F(J) = [f_1(J), f_2(J), \dots, f_m(J)]$  is defined as:

$$J(J) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

It consists of the first order partial derivatives and gives therefore an overview about the gradients of a vector valued function.

##### Example: Jacobian matrix

Consider a vector-valued function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^3$  defined as follows:

$$f(J) = \begin{bmatrix} x_1^2 + 2x_2 \\ 3x_1 - \sin(x_2) \\ e^{x_1+x_2} \end{bmatrix}$$

Let's compute the partial derivatives and construct the Jacobian matrix:

$$\frac{\partial f_1}{\partial x_1} = 2x_1, \quad \frac{\partial f_1}{\partial x_2} = 2$$

$$\frac{\partial f_2}{\partial x_1} = 3, \quad \frac{\partial f_2}{\partial x_2} = -\cos(x_2)$$

$$\frac{\partial f_3}{\partial x_1} = e^{x_1+x_2}, \quad \frac{\partial f_3}{\partial x_2} = e^{x_1+x_2}$$

So, the Jacobian matrix is:

$$J(J) = \begin{bmatrix} 2x_1 & 2 \\ 3 & -\cos(x_2) \\ e^{x_1+x_2} & e^{x_1+x_2} \end{bmatrix}$$

This Jacobian matrix provides information about how small changes in the input variables  $x_1$  and  $x_2$  affect the corresponding changes in each component of the output vector.

### 3.2.5 Hessian Matrix

The Hessian matrix  $H(J)$  for a scalar function  $f(J)$  is defined as:

$$H(J) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

So, the Hessian matrix consists of the second order derivatives of the function. It provides information about the local curvature of the function with respect to changes in the input variables.

#### **i** Example: Hessian matrix

Consider a scalar-valued function:

$$f(J) = x_1^2 + 2x_2^2 + \sin(x_1 x_2)$$

The Hessian matrix of this scalar-valued function is the matrix of its second-order partial derivatives with respect to the input variables:

$$H(J) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}$$

Let's compute the second-order partial derivatives and construct the Hessian matrix:

$$\frac{\partial^2 f}{\partial x_1^2} = 2 + \cos(x_1 x_2) x_2^2 \quad (3.3)$$

$$\frac{\partial^2 f}{\partial x_1 \partial x_2} = 2x_1 x_2 \cos(x_1 x_2) - \sin(x_1 x_2) \quad (3.4)$$

$$\frac{\partial^2 f}{\partial x_2 \partial x_1} = 2x_1 x_2 \cos(x_1 x_2) - \sin(x_1 x_2) \quad (3.5)$$

$$\frac{\partial^2 f}{\partial x_2^2} = 4x_2^2 + \cos(x_1 x_2) x_1^2 \quad (3.6)$$

So, the Hessian matrix is:

$$H(J) = \begin{bmatrix} 2 + \cos(x_1 x_2) x_2^2 & 2x_1 x_2 \cos(x_1 x_2) - \sin(x_1 x_2) \\ 2x_1 x_2 \cos(x_1 x_2) - \sin(x_1 x_2) & 4x_2^2 + \cos(x_1 x_2) x_1^2 \end{bmatrix}$$

### 3.2.6 Gradient for Optimization

In optimization, the goal is to find the minimum or maximum of a function. Gradient-based optimization methods utilize information about the gradient (or derivative) of the function to guide the search for the optimal solution. This is particularly useful when dealing with complex, high-dimensional functions where an exhaustive search is impractical.

The gradient descent method can be divided in the following steps:

- **Initialize:** start with an initial guess for the parameters of the function to be optimized.
- **Compute Gradient:** Calculate the gradient (partial derivatives) of the function with respect to each parameter at the current point. The gradient indicates the direction of the steepest increase in the function.
- **Update Parameters:** Adjust the parameters in the opposite direction of the gradient, scaled by a learning rate. This step aims to move towards the minimum of the function:

- $x_{k+1} = x_k - \alpha \times \nabla f(x_k)$
- $x_x$  is current parameter vector or point in the parameter space.
- $\alpha$  is the learning rate, a positive scalar that determines the step size in each iteration.
- $\nabla f(x)$  is the gradient of the objective function.

- **Iterate:** Repeat the above steps until convergence or a predefined number of iterations. Convergence is typically determined when the change in the function value or parameters becomes negligible.

#### i Example: Gradient Descent

Let's consider a simple quadratic function as an example:

$$f(x) = x^2 + 4x + y^2 + 2y + 4.$$

We'll use gradient descent to find the minimum of this function.

```

import numpy as np
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D

# Define the quadratic function
def quadratic_function(x, y):
    return x**2 + 4*x + y**2 + 2*y + 4

# Define the gradient of the quadratic function
def gradient_quadratic_function(x, y):
    grad_x = 2*x + 4
    grad_y = 2*y + 2
    return np.array([grad_x, grad_y])

# Gradient Descent for optimization in 2D
def gradient_descent(initial_point, learning_rate, num_iterations):
    points = [np.array(initial_point)]

    for _ in range(num_iterations):
        current_point = points[-1]
        gradient = gradient_quadratic_function(*current_point)
        new_point = current_point - learning_rate * gradient

        points.append(new_point)

    return points

# Visualization of optimization process with 3D surface and consistent arrow sizes
def plot_optimization_process_3d_consistent_arrows(points):
    fig = plt.figure(figsize=(10, 8))
    ax = fig.add_subplot(111, projection='3d')

    x_vals = np.linspace(-10, 2, 100)
    y_vals = np.linspace(-10, 2, 100)
    X, Y = np.meshgrid(x_vals, y_vals)
    Z = quadratic_function(X, Y)

    ax.plot_surface(X, Y, Z, cmap='viridis', alpha=0.6)
    ax.scatter(*zip(*points), [quadratic_function(*p) for p in points], c='red', label='Optimal Path')

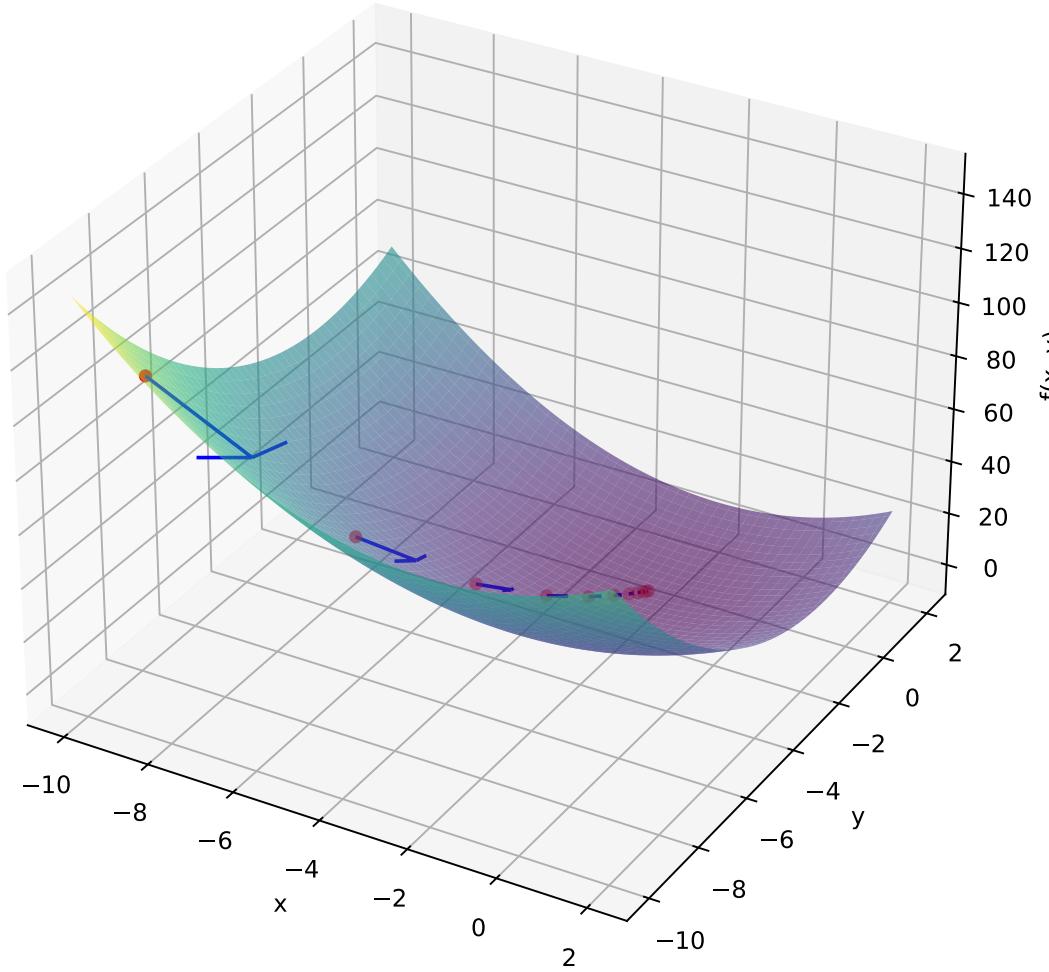
    for i in range(len(points) - 1):
        x, y = points[i]
        dx, dy = points[i + 1] - points[i]
        dz = quadratic_function(*(points[i + 1])) - quadratic_function(*points[i])
        gradient_length = 0.5
        41
        ax.quiver(x, y, quadratic_function(*points[i]), dx, dy, dz, color='blue', length=gradient_length)

    ax.set_title('Gradient-Based Optimization with 2D Quadratic Function')
    ax.set_xlabel('x')
    ax.set_ylabel('y')
    ax.set_zlabel('f(x, y)')
    ax.legend()

```

## Gradient-Based Optimization with 2D Quadratic Function

● Optimization Trajectory



### 3.2.7 Newton Method

**Initialization:** Start with an initial guess for the optimal solution:  $x_0$ .

**Iteration:** Repeat the following three steps until convergence or a predefined stopping criterion is met:

- 1) Calculate the gradient ( $\nabla$ ) and the Hessian matrix ( $\nabla^2$ ) of the objective function at the

current point:

$$\nabla f(x_k) \quad \text{and} \quad \nabla^2 f(x_k)$$

2) Update the current solution using the Newton-Raphson update formula

$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k),$$

where

- $\nabla f(x_k)$  is the gradient (first derivative) of the objective function with respect to the variable  $x$ , evaluated at the current solution  $x_k$ .
- $\nabla^2 f(x_k)$ : The Hessian matrix (second derivative) of the objective function with respect to  $x$ , evaluated at the current solution  $x_k$ .
- $x_k$ : The current solution or point in the optimization process.
- $[\nabla^2 f(x_k)]^{-1}$ : The inverse of the Hessian matrix at the current point, representing the approximation of the curvature of the objective function.
- $x_{k+1}$ : The updated solution or point after applying the Newton-Raphson update.

3) Check for convergence.

#### i Example: Newton Method

We want to optimize the Rosenbrock function and use the Hessian and the Jacobian (which is equal to the gradient vector for scalar objective function) to the `minimize` function.

```

def rosenbrock(x):
    return 100 * (x[1] - x[0]**2)**2 + (1 - x[0])**2

def rosenbrock_gradient(x):
    dfdx0 = -400 * x[0] * (x[1] - x[0]**2) - 2 * (1 - x[0])
    dfdx1 = 200 * (x[1] - x[0]**2)
    return np.array([dfdx0, dfdx1])

def rosenbrock_hessian(x):
    d2fdx0 = 1200 * x[0]**2 - 400 * x[1] + 2
    d2fdx1 = -400 * x[0]
    return np.array([[d2fdx0, d2fdx1], [d2fdx1, 200]])

def classical_newton_optimization_2d(initial_guess, tol=1e-6, max_iter=100):
    x = initial_guess.copy()

    for i in range(max_iter):
        gradient = rosenbrock_gradient(x)
        hessian = rosenbrock_hessian(x)

        # Solve the linear system H * d = -g for d
        d = np.linalg.solve(hessian, -gradient)

        # Update x
        x += d

        # Check for convergence
        if np.linalg.norm(gradient, ord=np.inf) < tol:
            break

    return x

# Initial guess
initial_guess_2d = np.array([0.0, 0.0])

# Run classical Newton optimization for the 2D Rosenbrock function
result_2d = classical_newton_optimization_2d(initial_guess_2d)

# Print the result
print("Optimal solution:", result_2d)
print("Objective value:", rosenbrock(result_2d))

```

Optimal solution: [1. 1.]  
 Objective value: 0.0

### 3.2.8 BFGS-Algorithm

BFGS is an optimization algorithm designed for unconstrained optimization problems. It belongs to the class of quasi-Newton methods and is known for its efficiency in finding the minimum of a smooth, unconstrained objective function.

### 3.2.9 Procedure:

#### 1. Initialization:

- Start with an initial guess for the parameters of the objective function.
- Initialize an approximation of the Hessian matrix (inverse) denoted by  $H$ .

#### 2. Iterative Update:

- At each iteration, compute the gradient vector at the current point.
- Update the parameters using the BFGS update formula, which involves the inverse Hessian matrix approximation, the gradient, and the difference in parameter vectors between successive iterations:

$$x_{k+1} = x_k - H_k^{-1} \nabla f(x_k).$$

- Update the inverse Hessian approximation using the BFGS update formula for the inverse Hessian.

$$H_{k+1} = H_k + \frac{\Delta x_k \Delta x_k^T}{\Delta x_k^T \Delta g_k} - \frac{H_k g_k g_k^T H_k}{g_k^T H_k g_k},$$

where:

- $x_k$  and  $x_{k+1}$  are the parameter vectors at the current and updated iterations, respectively.
- $\nabla f(x_k)$  is the gradient vector at the current iteration.
- $\Delta x_k = x_{k+1} - x_k$  is the change in parameter vectors.
- $\Delta g_k = \nabla f(x_{k+1}) - \nabla f(x_k)$  is the change in gradient vectors.

#### 3. Convergence:

- Repeat the iterative update until the optimization converges. Convergence is typically determined by reaching a sufficiently low gradient or parameter change.

### **i** Example: BFGS for Rosenbrock

```
import numpy as np
from scipy.optimize import minimize

# Define the 2D Rosenbrock function
def rosenbrock(x):
    return (1 - x[0])**2 + 100 * (x[1] - x[0]**2)**2

# Initial guess
initial_guess = np.array([0.0, 0.0])

# Minimize the Rosenbrock function using BFGS
minimize(rosenbrock, initial_guess, method='BFGS')
```

```
message: Optimization terminated successfully.
success: True
status: 0
    fun: 2.843987518235081e-11
    x: [ 1.000e+00  1.000e+00]
    nit: 19
    jac: [ 3.987e-06 -2.844e-06]
    hess_inv: [[ 4.948e-01  9.896e-01]
                [ 9.896e-01  1.984e+00]]
    nfev: 72
    njev: 24
```

### **3.2.10 Visualization BFGS for Rosenbrock**

A visualization of the BFGS search process on Rosenbrock's function can be found here: <https://upload.wikimedia.org/wikipedia/de/f/ff/Rosenbrock-bfgs-animation.gif>

### **i** Tasks

- In which situations is it possible to use algorithms like BFGS, but not the classical Newton method?
- Investigate the Newton-CG method
- Use an objective function of your choice and apply Newton-CG
- Compare the Newton-CG method with the BFGS. What are the similarities and differences between the two algorithms?

### 3.3 Gradient- and Hessian-based optimization algorithms

Section 3.3.1 presents an optimization algorithm that uses gradient and Hessian information to find the minimum. Section 3.3.2 presents an optimization algorithm that uses gradient and Hessian information to find the minimum. Section 3.3.3 presents an optimization algorithm that uses gradient and Hessian information to find the minimum.

The methods Newton-CG, trust-ncg and trust-krylov are suitable for dealing with large-scale problems (problems with thousands of variables). That is because the conjugate gradient algorithm approximately solve the trust-region subproblem (or invert the Hessian) by iterations without the explicit Hessian factorization. Since only the product of the Hessian with an arbitrary vector is needed, the algorithm is specially suited for dealing with sparse Hessians, allowing low storage requirements and significant time savings for those sparse problems.

#### 3.3.1 Newton-Conjugate-Gradient Algorithm

Newton-Conjugate Gradient algorithm is a modified Newton's method and uses a conjugate gradient algorithm to (approximately) invert the local Hessian.

#### 3.3.2 Trust-Region Newton-Conjugate-Gradient Algorithm

#### 3.3.3 Trust-Region Truncated Generalized Lanczos / Conjugate Gradient Algorithm

### 3.4 Global Optimization

Global optimization aims to find the global minimum of a function within given bounds, in the presence of potentially many local minima. Typically, global minimizers efficiently search the parameter space, while using a local minimizer (e.g., minimize) under the hood. SciPy contains a number of good global optimizers. Here, we'll use those on the same objective function, namely the (aptly named) eggholder function:

```
def eggholder(x):
    return -(x[1] + 47) * np.sin(np.sqrt(abs(x[0]/2 + (x[1] + 47))))
    -x[0] * np.sin(np.sqrt(abs(x[0] - (x[1] + 47)))))

bounds = [(-512, 512), (-512, 512)]
```

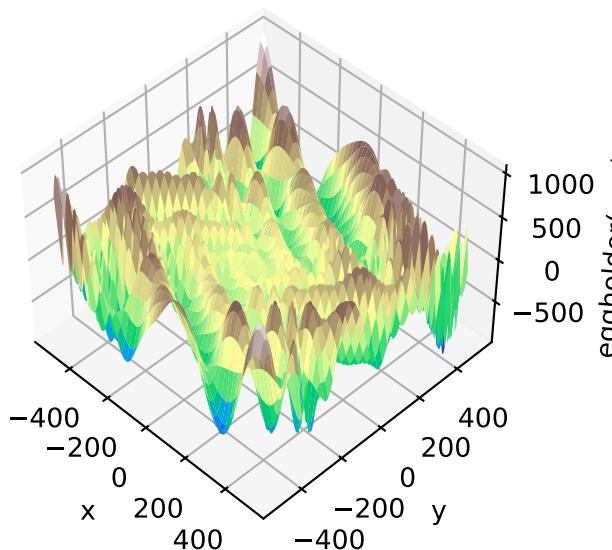
```

import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D

x = np.arange(-512, 513)
y = np.arange(-512, 513)
xgrid, ygrid = np.meshgrid(x, y)
xy = np.stack([xgrid, ygrid])

fig = plt.figure()
ax = fig.add_subplot(111, projection='3d')
ax.view_init(45, -45)
ax.plot_surface(xgrid, ygrid, eggholder(xy), cmap='terrain')
ax.set_xlabel('x')
ax.set_ylabel('y')
ax.set_zlabel('eggholder(x, y)')
plt.show()

```



We now use the global optimizers to obtain the minimum and the function value at the minimum. We'll store the results in a dictionary so we can compare different optimization results later.

```

from scipy import optimize
results = dict()
results['shgo'] = optimize.shgo(eggholder, bounds)
results['shgo']

```

```

message: Optimization terminated successfully.
success: True
    fun: -935.3379515605789
    funl: [-9.353e+02]
        x: [ 4.395e+02  4.540e+02]
        xl: [[ 4.395e+02  4.540e+02]]
    nit: 1
    nfev: 45
    nlfev: 40
    nljev: 10
    nlhev: 0

results['DA'] = optimize.dual_annealing(eggholder, bounds)
results['DA']

```

```

message: ['Maximum number of iteration reached']
success: True
status: 0
    fun: -959.6406627208282
        x: [ 5.120e+02  4.042e+02]
    nit: 1000
    nfev: 4091
    njev: 30
    nhev: 0

```

All optimizers return an `OptimizeResult`, which in addition to the solution contains information on the number of function evaluations, whether the optimization was successful, and more. For brevity, we won't show the full output of the other optimizers:

```

results['DE'] = optimize.differential_evolution(eggholder, bounds)
results['DE']

message: Optimization terminated successfully.
success: True
    fun: -956.9182316230113
        x: [ 4.824e+02  4.329e+02]
    nit: 42
    nfev: 1326
    jac: [ 1.137e-04 -1.137e-04]

```

`shgo` has a second method, which returns all local minima rather than only what it thinks is the global minimum:

```

results['shgo_sobol'] = optimize.shgo(eggholder, bounds, n=200, iters=5,
                                      sampling_method='sobol')
results['shgo_sobol']

```

```

message: Optimization terminated successfully.
success: True
  fun: -959.640662720831
  funl: [-9.596e+02 -9.353e+02 ... -6.591e+01 -6.387e+01]
    x: [ 5.120e+02  4.042e+02]
    xl: [[ 5.120e+02  4.042e+02]
          [ 4.395e+02  4.540e+02]
          ...
          [ 3.165e+01 -8.523e+01]
          [ 5.865e+01 -5.441e+01]]
  nit: 5
  nfev: 3529
  nlfev: 2327
  nljev: 634
  nlhev: 0

```

We'll now plot all found minima on a heatmap of the function:

```

fig = plt.figure()
ax = fig.add_subplot(111)
im = ax.imshow(eggholder(xy), interpolation='bilinear', origin='lower',
               cmap='gray')
ax.set_xlabel('x')
ax.set_ylabel('y')

def plot_point(res, marker='o', color=None):
    ax.plot(512+res.x[0], 512+res.x[1], marker=marker, color=color, ms=10)

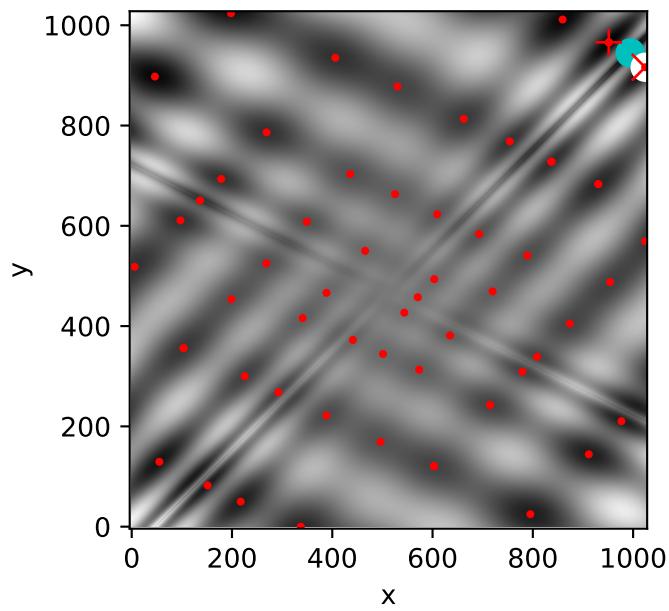
plot_point(results['DE'], color='c') # differential_evolution - cyan
plot_point(results['DA'], color='w') # dual_annealing. - white

# SHGO produces multiple minima, plot them all (with a smaller marker size)
plot_point(results['shgo'], color='r', marker='+')
plot_point(results['shgo_sobol'], color='r', marker='x')
for i in range(results['shgo_sobol'].xl.shape[0]):
    ax.plot(512 + results['shgo_sobol'].xl[i, 0],
            512 + results['shgo_sobol'].xl[i, 1],

```

```
'ro', ms=2)

ax.set_xlim([-4, 514*2])
ax.set_ylim([-4, 514*2])
plt.show()
```



### 3.4.1 Dual Annealing Optimization

This function implements the Dual Annealing optimization.

### 3.4.2 Differential Evolution

Differential Evolution is an algorithm used for finding the global minimum of multivariate functions. It is stochastic in nature (does not use gradient methods), and can search large areas of candidate space, but often requires larger numbers of function evaluations than conventional gradient based techniques.

### 3.4.3 DIRECT

DIviding RECTangles (DIRECT) is a deterministic global optimization algorithm capable of minimizing a black box function with its variables subject to lower and upper bound constraints

by sampling potential solutions in the search space

#### 3.4.4 SHGO

SHGO stands for “simplicial homology global optimization”. It is considered appropriate for solving general purpose NLP and blackbox optimization problems to global optimality (low-dimensional problems).

#### 3.4.5 Basin-hopping

Basin-hopping is a two-phase method that combines a global stepping algorithm with local minimization at each step. Designed to mimic the natural process of energy minimization of clusters of atoms, it works well for similar problems with “funnel-like, but rugged” energy landscapes

### 3.5 Jupyter Notebook

#### i Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

## 4 Sequential Parameter Optimization: Using `scipy` Optimizers

As a default optimizer, `spotPython` uses `differential_evolution` from the `scipy.optimize` package. Alternatively, any other optimizer from the `scipy.optimize` package can be used. This chapter describes how different optimizers from the `scipy.optimize` package can be used on the surrogate. The optimization algorithms are available from <https://docs.scipy.org/doc/scipy/reference/optimize.html>

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from scipy.optimize import shgo
from scipy.optimize import direct
from scipy.optimize import differential_evolution
from scipy.optimize import dual_annealing
from scipy.optimize import basinhopping
from spotPython.utils.init import fun_control_init, design_control_init, optimizer_control_i
```

### 4.1 The Objective Function Branin

The `spotPython` package provides several classes of objective functions. We will use an analytical objective function, i.e., a function that can be described by a (closed) formula. Here we will use the Branin function. The 2-dim Branin function is

```
$$y = a * (x2 - b * x1**2 + c * x1 - r) ** 2 + s * (1 - t) * \cos(x1) + s,$$
where values of a, b, c, r, s and t are:
$a = 1$, $b = 5.1 / (4\pi^2)$, $c = 5 / \pi$, $r = 6$, $s = 10$ and $t = 1 / (8\pi)$.
```

- It has three global minima:

$$f(x) = 0.397887 \text{ at } (-\pi, 12.275), (\pi, 2.275), \text{ and } (9.42478, 2.475).$$

- Input Domain: This function is usually evaluated on the square  $x_1$  in  $[-5, 10]$   $\times$   $x_2$  in  $[0, 15]$ .

```

from spotPython.fun.objectivefunctions import analytical
lower = np.array([-5,-0])
upper = np.array([10,15])
fun = analytical(seed=123).fun_branin

```

## 4.2 The Optimizer

Differential Evolution (DE) from the `scikit.optimize` package, see [https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.differential\\_evolution.html#scipy.optimize.differential\\_evolution](https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.differential_evolution.html#scipy.optimize.differential_evolution) is the default optimizer for the search on the surrogate. Other optimiers that are available in `spotPython`, see <https://docs.scipy.org/doc/scipy/reference/optimize.html#global-optimization>.

- `dual_annealing`
- `direct`
- `shgo`
- `basinhopping`

These optimizers can be selected as follows:

```
surrogate_control = "model_optimizer": differential_evolution
```

As noted above, we will use `differential_evolution`. The optimizer can use 1000 evaluations. This value will be passed to the `differential_evolution` method, which has the argument `maxiter` (int). It defines the maximum number of generations over which the entire differential evolution population is evolved, see [https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.differential\\_evolution.html#scipy.optimize.differential\\_evolution](https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.differential_evolution.html#scipy.optimize.differential_evolution)

### TensorBoard

Similar to the one-dimensional case, which is discussed in Section 7.5, we can use TensorBoard to monitor the progress of the optimization. We will use a similar code, only the prefix is different:

```

fun_control=fun_control_init(
    lower = lower,
    upper = upper,
    fun_evals = 20,
    PREFIX = "04_DE_"
)
surrogate_control=surrogate_control_init(
    n_theta=len(lower))

```

```
Created spot_tensorboard_path: runs/spot_logs/04_DE_p040025_2024-01-09_20-00-19 for Summary
```

```
spot_de = spot.Spot(fun=fun,
                     fun_control=fun_control,
                     surrogate_control=surrogate_control)
spot_de.run()
```

```
spotPython tuning: 3.8004580634289518 [#####----] 55.00%
spotPython tuning: 3.8004580634289518 [#####----] 60.00%
spotPython tuning: 3.158983526047736 [#####----] 65.00%
spotPython tuning: 3.1338444083542836 [#####----] 70.00%
spotPython tuning: 2.917653611971259 [#####----] 75.00%
spotPython tuning: 0.40458354637529226 [#####----] 80.00%
spotPython tuning: 0.40458354637529226 [#####----] 85.00%
spotPython tuning: 0.3987528104146545 [#####----] 90.00%
spotPython tuning: 0.3987528104146545 [#####----] 95.00%
spotPython tuning: 0.3987528104146545 [#####----] 100.00% Done...
```

```
<spotPython.spot.spot.Spot at 0x2d353b050>
```

#### 4.2.1 TensorBoard

If the `prefix` argument in `fun_control_init()` is not `None` (as above, where the `prefix` was set to `04_DE_`) , we can start TensorBoard in the background with the following command:

```
tensorboard --logdir=".runs"
```

We can access the TensorBoard web server with the following URL:

```
http://localhost:6006/
```

The TensorBoard plot illustrates how `spotPython` can be used as a microscope for the internal mechanisms of the surrogate-based optimization process. Here, one important parameter, the learning rate  $\theta$  of the Kriging surrogate is plotted against the number of optimization steps.

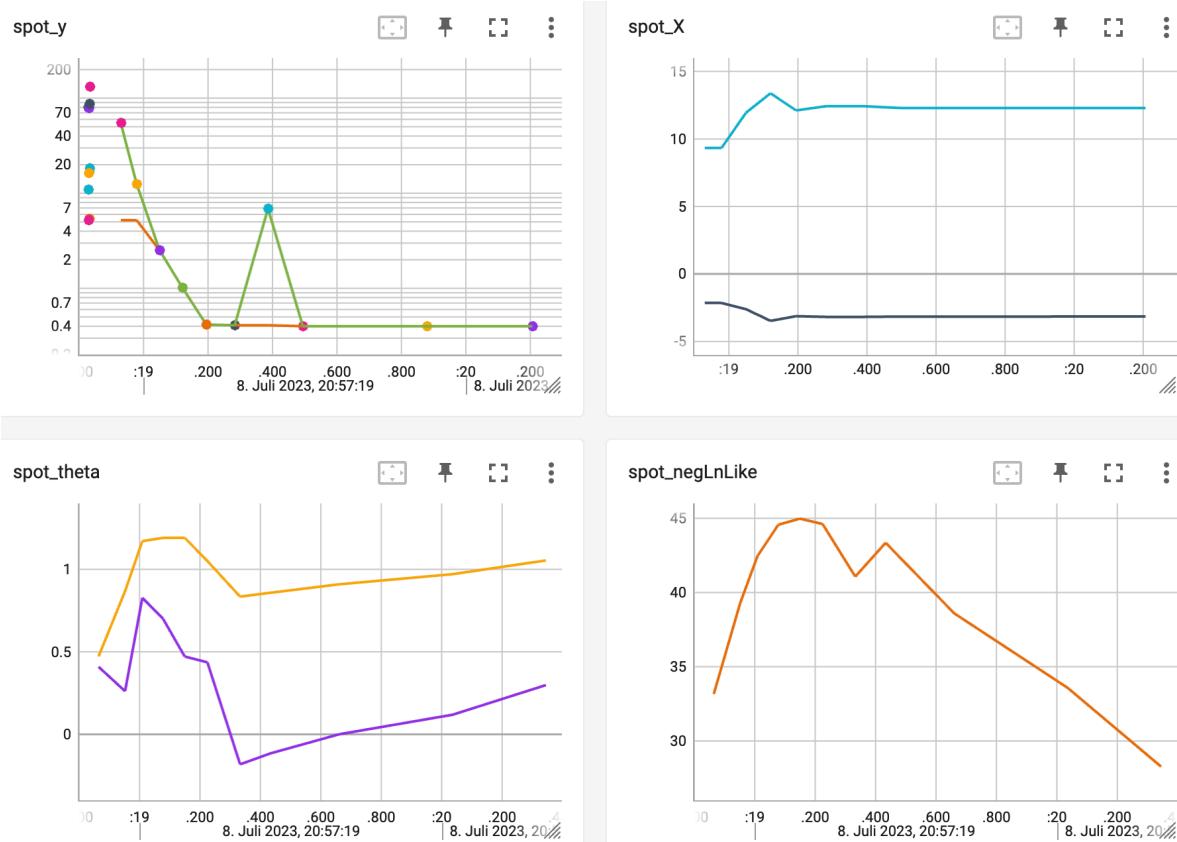


Figure 4.1: TensorBoard visualization of the spotPython optimization process and the surrogate model.

### 4.3 Print the Results

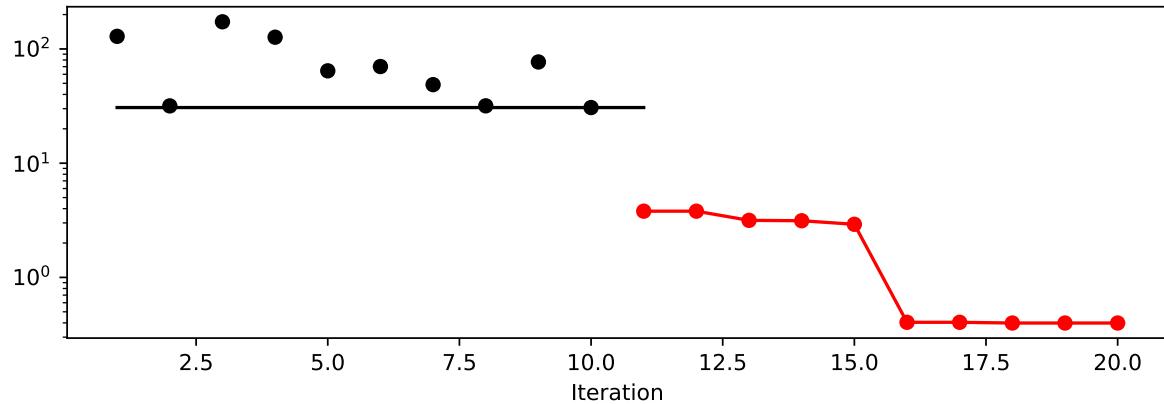
```
spot_de.print_results()
```

```
min y: 0.3987528104146545
x0: 3.14748607975711
x1: 2.2968413897617554
```

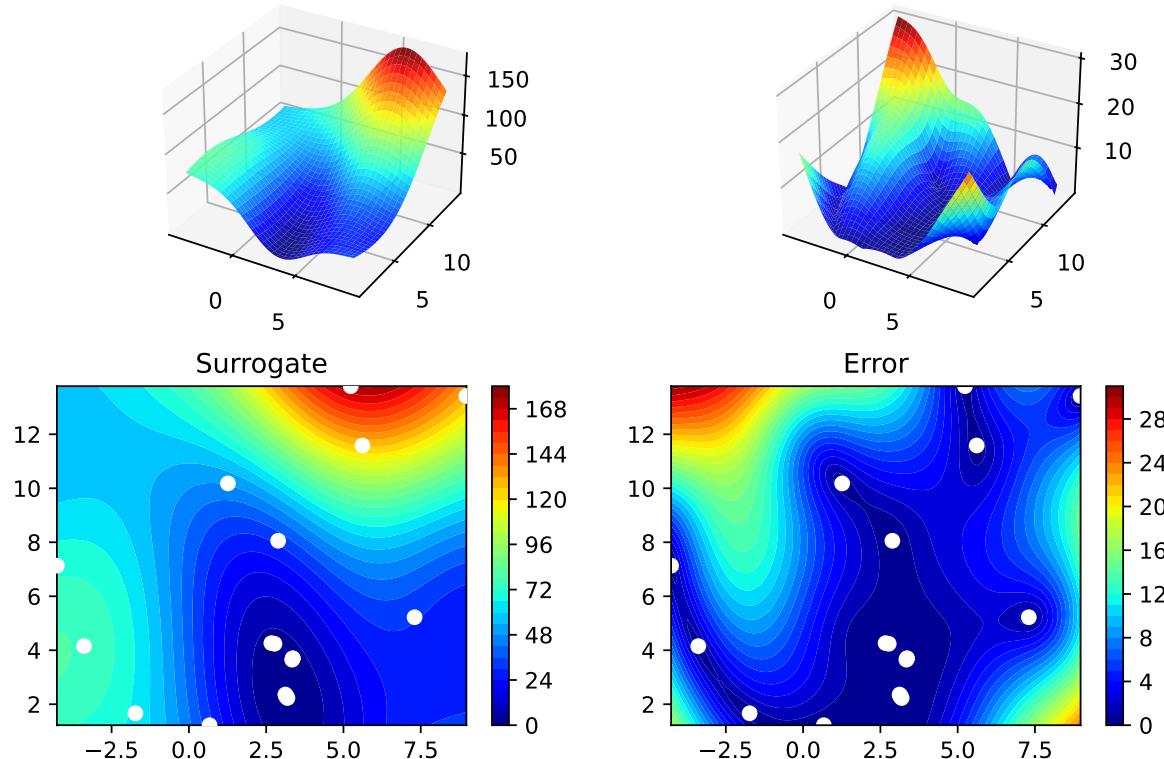
```
[['x0', 3.14748607975711], ['x1', 2.2968413897617554]]
```

## 4.4 Show the Progress

```
spot_de.plot_progress(log_y=True)
```



```
spot_de.surrogate.plot()
```



## 4.5 Exercises

### 4.5.1 dual\_annealing

- Describe the optimization algorithm, see [scipy.optimize.dual\\_annealing](#).
- Use the algorithm as an optimizer on the surrogate.

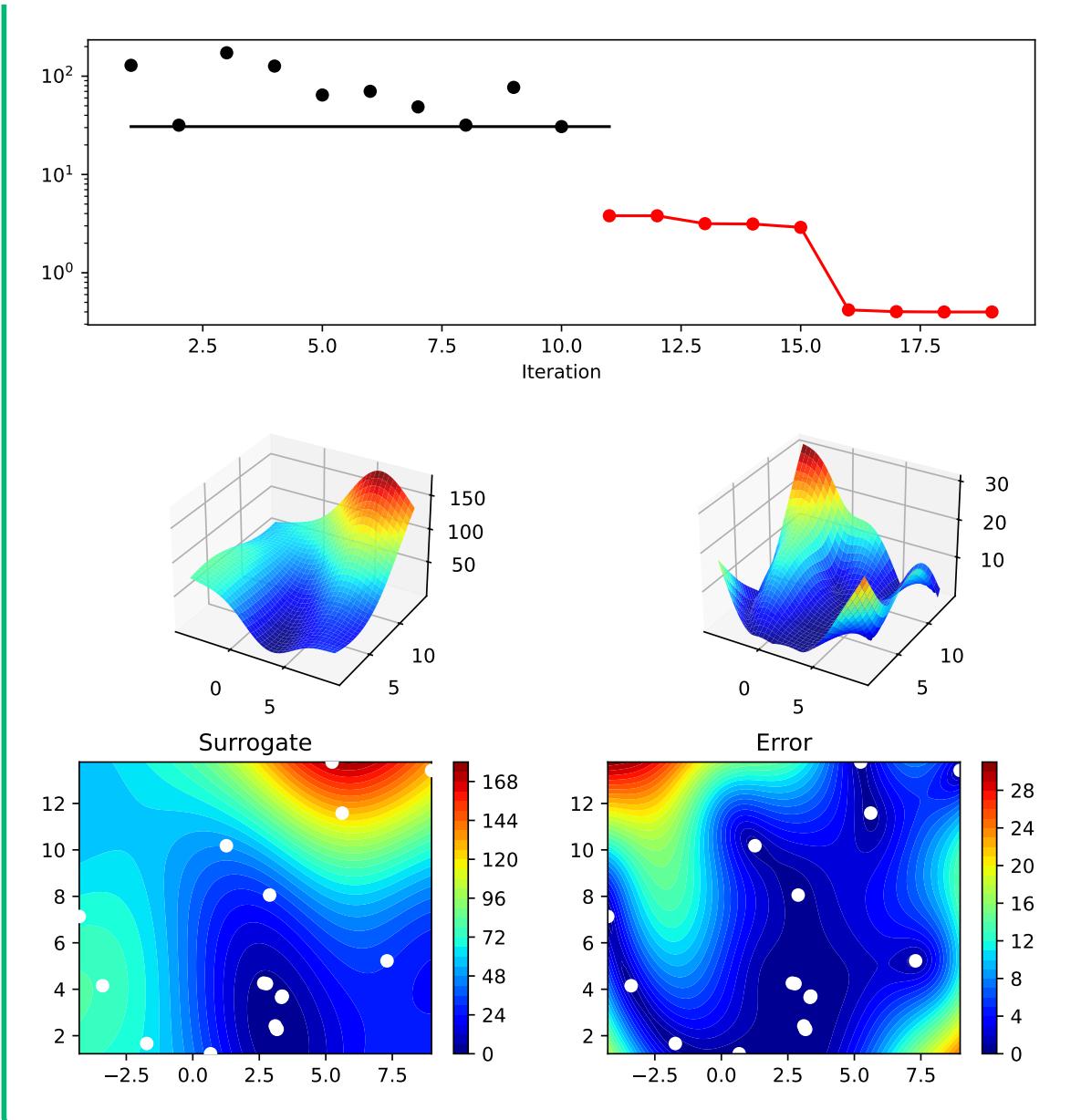
 Tip: Selecting the Optimizer for the Surrogate

We can run spotPython with the `dual_annealing` optimizer as follows:

```
spot_da = spot.Spot(fun=fun,
                     fun_control=fun_control,
                     optimizer=dual_annealing,
                     surrogate_control=surrogate_control)

spot_da.run()
spot_da.print_results()
spot_da.plot_progress(log_y=True)
spot_da.surrogate.plot()

spotPython tuning: 3.8004506180745494 [#####----] 55.00%
spotPython tuning: 3.8004506180745494 [#####----] 60.00%
spotPython tuning: 3.159026259576751 [#####----] 65.00%
spotPython tuning: 3.1343634581305215 [#####---] 70.00%
spotPython tuning: 2.8965876127956935 [#####---] 75.00%
spotPython tuning: 0.41905758842574414 [#####---] 80.00%
spotPython tuning: 0.4020256285473973 [#####---] 85.00%
spotPython tuning: 0.39921903284476734 [#####---] 90.00%
spotPython tuning: 0.39921903284476734 [#####---] 95.00%
min y: 0.39921903284476734
x0: 3.1508062869062012
x1: 2.2982248293554903
```



#### 4.5.2 direct

- Describe the optimization algorithm
- Use the algorithm as an optimizer on the surrogate

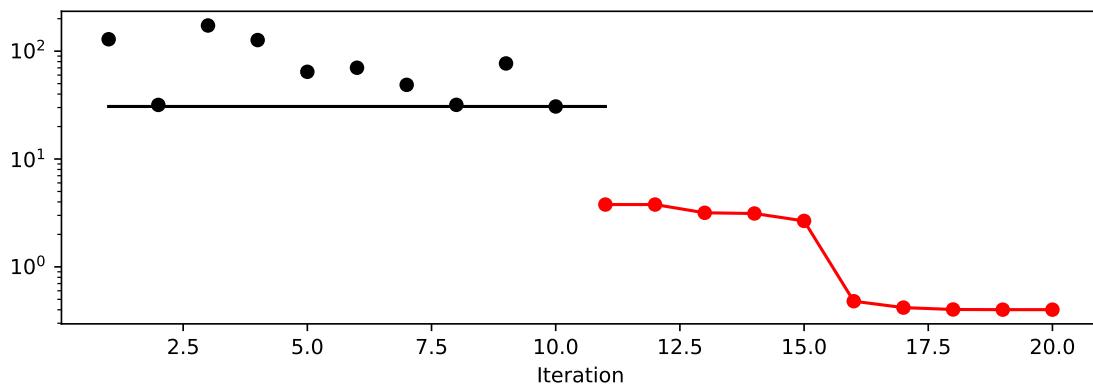
### Tip: Selecting the Optimizer for the Surrogate

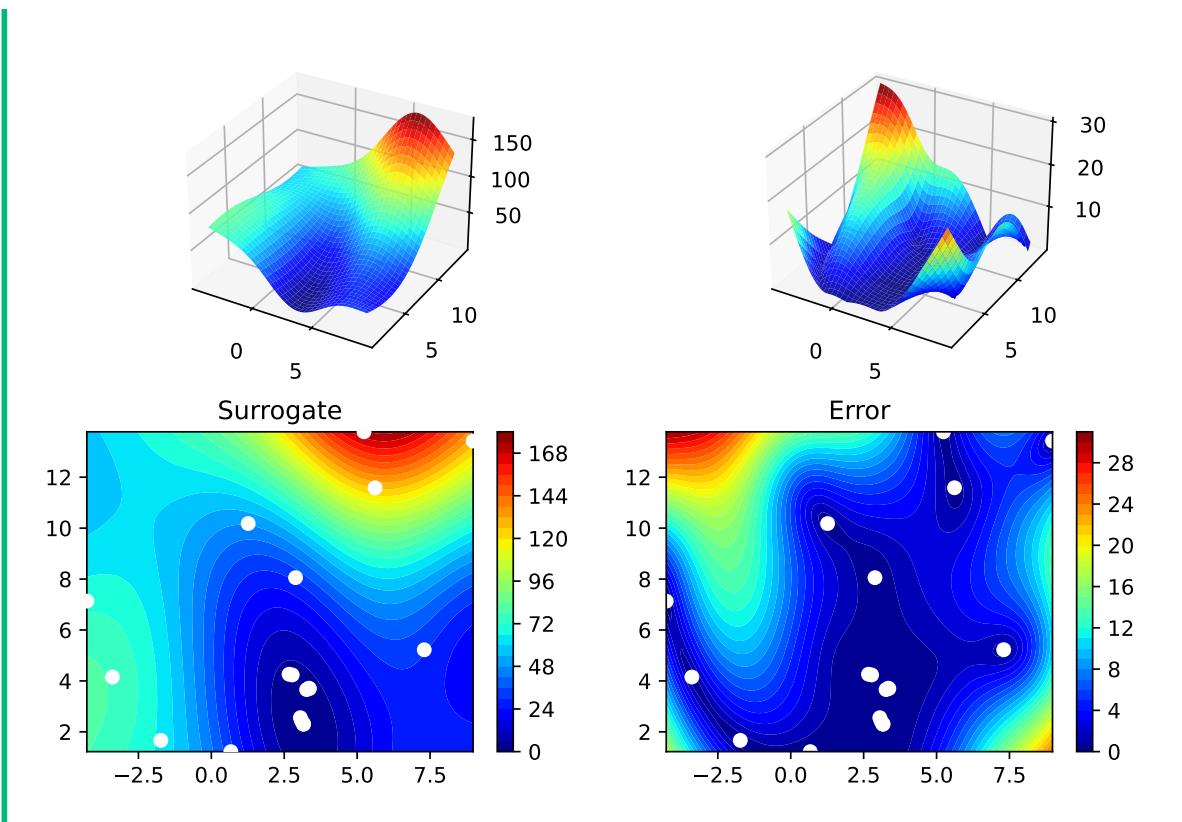
We can run spotPython with the `direct` optimizer as follows:

```
spot_di = spot.Spot(fun=fun,
                     fun_control=fun_control,
                     optimizer=direct,
                     surrogate_control=surrogate_control)
spot_di.run()
spot_di.print_results()
spot_di.plot_progress(log_y=True)
spot_di.surrogate.plot()
```

```
spotPython tuning: 3.78192024900577 [#####----] 55.00%
spotPython tuning: 3.78192024900577 [#####----] 60.00%
spotPython tuning: 3.1707843299428866 [#####----] 65.00%
spotPython tuning: 3.1253295886690413 [#####---] 70.00%
spotPython tuning: 2.6673899789334117 [#####---] 75.00%
spotPython tuning: 0.48037397889434175 [#####---] 80.00%
spotPython tuning: 0.41903636460779303 [#####---] 85.00%
spotPython tuning: 0.4025655271878108 [#####---] 90.00%
spotPython tuning: 0.4010733027914206 [#####---] 95.00%
spotPython tuning: 0.4010733027914206 [#####---] 100.00% Done...

min y: 0.4010733027914206
x0: 3.1561499771376322
x1: 2.3102423411065387
```





#### 4.5.3 shgo

- Describe the optimization algorithm
- Use the algorithm as an optimizer on the surrogate

Tip: Selecting the Optimizer for the Surrogate

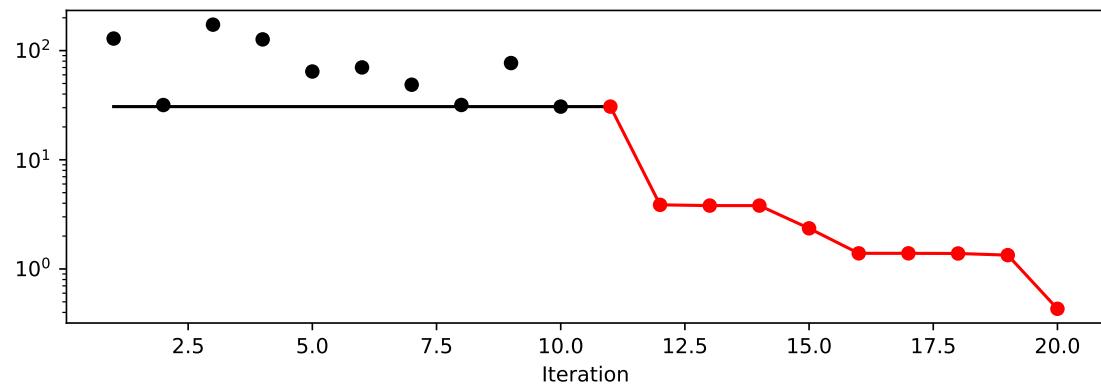
We can run spotPython with the `direct` optimizer as follows:

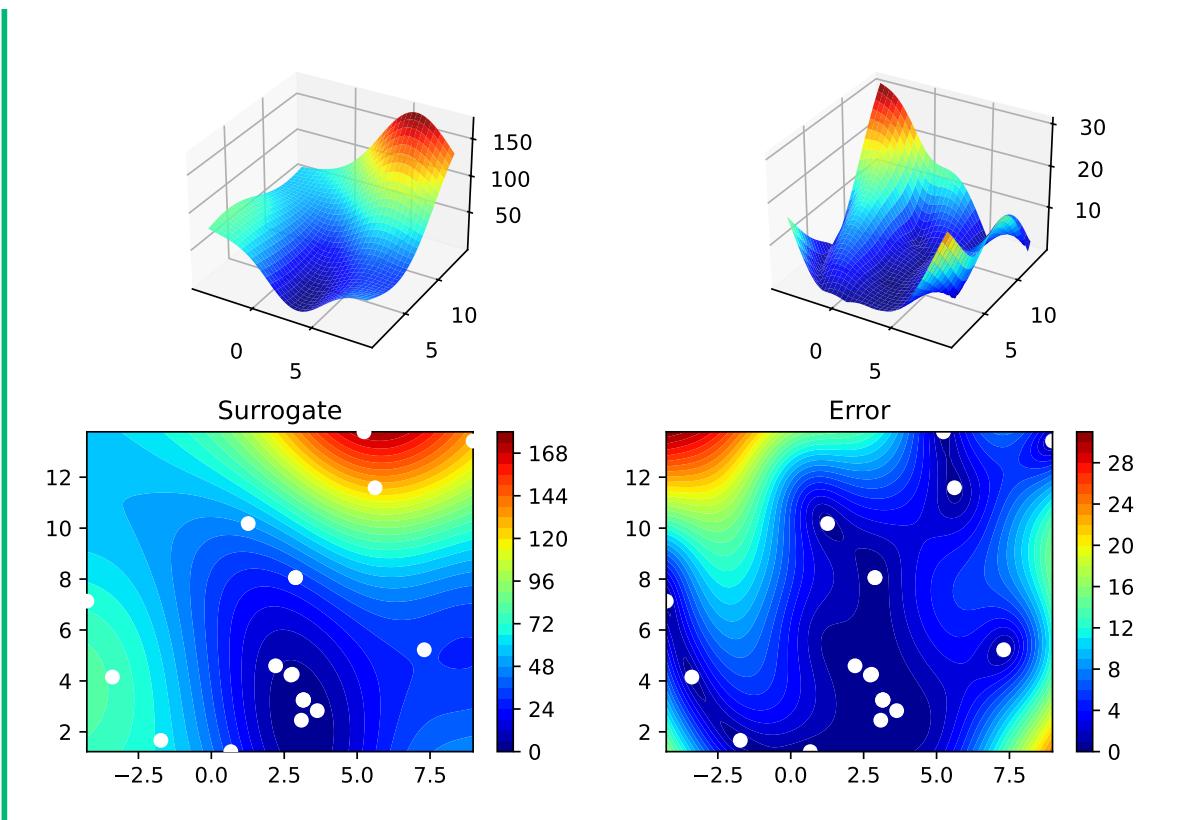
```
spot_sh = spot.Spot(fun=fun,
                     fun_control=fun_control,
                     optimizer=shgo,
                     surrogate_control=surrogate_control)

spot_sh.run()
spot_sh.print_results()
spot_sh.plot_progress(log_y=True)
spot_sh.surrogate.plot()
```

```
spotPython tuning: 30.69410528614059 [#####----] 55.00%
spotPython tuning: 3.8670090115148232 [#####----] 60.00%
spotPython tuning: 3.8060289138706764 [#####----] 65.00%
spotPython tuning: 3.8060289138706764 [#####----] 70.00%
spotPython tuning: 2.355799878130849 [#####---] 75.00%
spotPython tuning: 1.3905725351665694 [#####---] 80.00%
spotPython tuning: 1.3905725351665694 [#####---] 85.00%
spotPython tuning: 1.3864169812254215 [#####---] 90.00%
spotPython tuning: 1.3386472543572232 [#####---] 95.00%
spotPython tuning: 0.4315864524692383 [#####---] 100.00% Done...
```

```
min y: 0.4315864524692383
x0: 3.087070911243333
x1: 2.457298105703394
```





#### 4.5.4 basinhopping

- Describe the optimization algorithm
- Use the algorithm as an optimizer on the surrogate

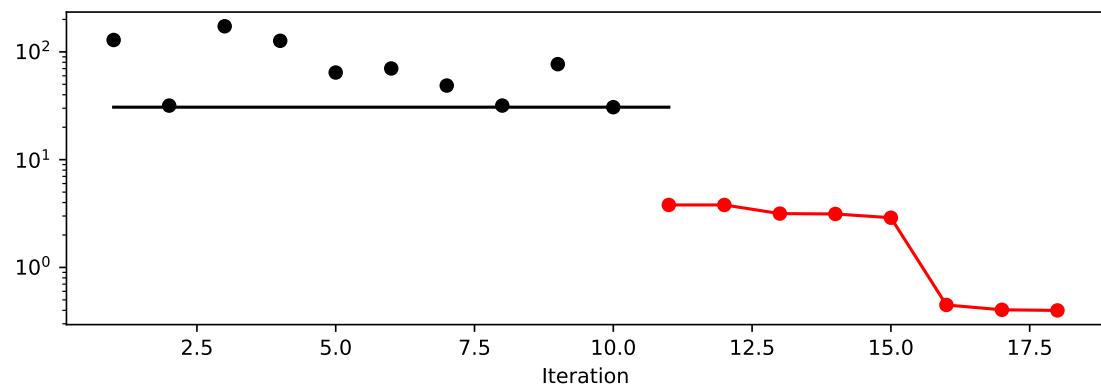
Tip: Selecting the Optimizer for the Surrogate

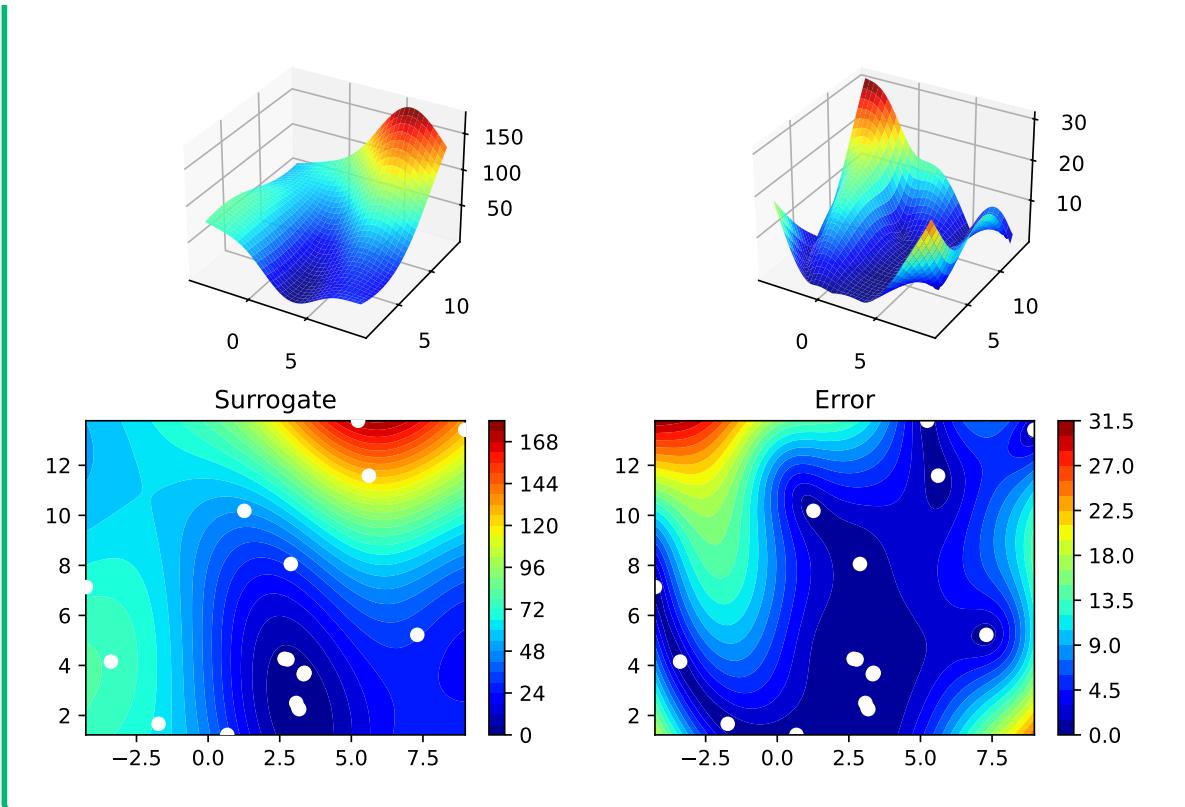
We can run spotPython with the `direct` optimizer as follows:

```
spot_bh = spot.Spot(fun=fun,
                     fun_control=fun_control,
                     optimizer=basinhopping,
                     surrogate_control=surrogate_control)

spot_bh.run()
spot_bh.print_results()
spot_bh.plot_progress(log_y=True)
spot_bh.surrogate.plot()
```

```
spotPython tuning: 3.800453417609317 [#####----] 55.00%
spotPython tuning: 3.800453417609317 [#####----] 60.00%
spotPython tuning: 3.1590142664203835 [#####----] 65.00%
spotPython tuning: 3.1341409786773404 [#####----] 70.00%
spotPython tuning: 2.8923708554157326 [#####---] 75.00%
spotPython tuning: 0.4484530687273516 [#####---] 80.00%
spotPython tuning: 0.404489521852879 [#####---] 85.00%
spotPython tuning: 0.39929273099035534 [#####---] 90.00%
min y: 0.39929273099035534
x0: 3.152060522318641
x1: 2.2965036205182314
```





#### 4.5.5 Performance Comparison

Compare the performance and run time of the 5 different optimizers:

- differential\_evolution
- dual\_annealing
- direct
- shgo
- basinhopping.

The Branin function has three global minima:

- $f(x) = 0.397887$  at
  - $(-\pi, 12.275)$ ,
  - $(\pi, 2.275)$ , and
  - $(9.42478, 2.475)$ .
- Which optima are found by the optimizers?

- Does the `seed` argument in `fun = analytical(seed=123).fun_branin` change this behavior?

## 4.6 Jupyter Notebook

### Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

## **Part II**

# **Numerical Methods**

# 5 Introduction: Numerical Methods

This part deals with numerical implementations of optimization methods. The goal is to understand the implementation of optimization methods and to solve real-world problems numerically and efficiently. We will focus on the implementation of surrogate models, because they are the most efficient way to solve real-world problems.

Starting point is the well-established response surface methodology. It will be extended to the design and analysis of computer experiments (DACE). The DACE methodology is a modern extension of the response surface methodology. It is based on the use of surrogate models, which are used to replace the real-world problem with a simpler problem. The simpler problem is then solved numerically. The solution of the simpler problem is then used to solve the real-world problem.

## ! Numerical methods: Goals

- Understand implementation of optimization methods
- Solve real-world problems numerically and efficiently

## 5.1 Response Surface Methods: What is RSM?

Response Surface Methods (RSM) refer to a collection of statistical and mathematical tools that are valuable for developing, improving, and optimizing processes. The overarching theme of RSM involves studying how input variables that control a product or process can potentially influence a response that measures performance or quality characteristics.

The advantages of RSM include a rich literature, well-established methods often used in manufacturing, the importance of careful experimental design combined with a well-understood model, and the potential to add significant value to scientific inquiry, process refinement, optimization, and more. However, there are also drawbacks to RSM, such as the use of simple and crude surrogates, the hands-on nature of the methods, and the limitation of local methods.

RSM is related to various fields, including Design of Experiments (DoE), quality management, reliability, and productivity. Its applications are widespread in industry and manufacturing, focusing on designing, developing, and formulating new products and improving existing ones, as well as from laboratory research. RSM is commonly applied in domains such as materials science, manufacturing, applied chemistry, climate science, and many others.

An example of RSM involves studying the relationship between a response variable, such as yield ( $y$ ) in a chemical process, and two process variables: reaction time ( $\xi_1$ ) and reaction temperature ( $\xi_2$ ). The provided code illustrates this scenario, following a variation of the so-called “banana function.”

In the context of visualization, RSM offers the choice between 3D plots and contour plots. In a 3D plot, the independent variables  $\xi_1$  and  $\xi_2$  are represented, with  $y$  as the dependent variable.

```
import numpy as np
import matplotlib.pyplot as plt

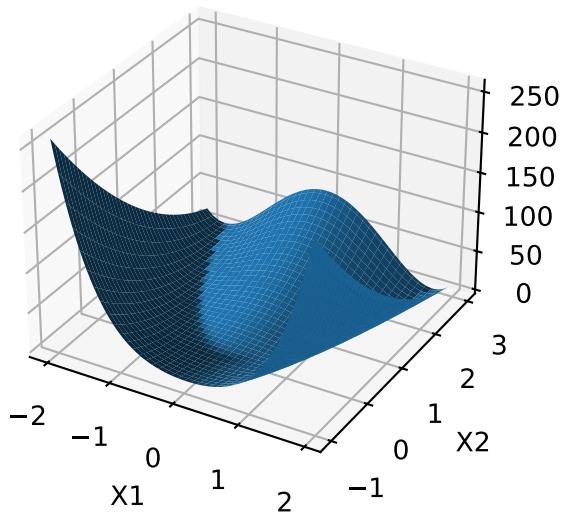
def fun_rosen(x1, x2):
    b = 10
    return (x1-1)**2 + b*(x2-x1**2)**2

fig = plt.figure()
ax = fig.add_subplot(111, projection='3d')
x = np.arange(-2.0, 2.0, 0.05)
y = np.arange(-1.0, 3.0, 0.05)
X, Y = np.meshgrid(x, y)
zs = np.array(fun_rosen(np.ravel(X), np.ravel(Y)))
Z = zs.reshape(X.shape)

ax.plot_surface(X, Y, Z)

ax.set_xlabel('X1')
ax.set_ylabel('X2')
ax.set_zlabel('Y')

plt.show()
```



- contour plot example:
  - $x_1$  and  $x_2$  are the independent variables
  - $y$  is the dependent variable

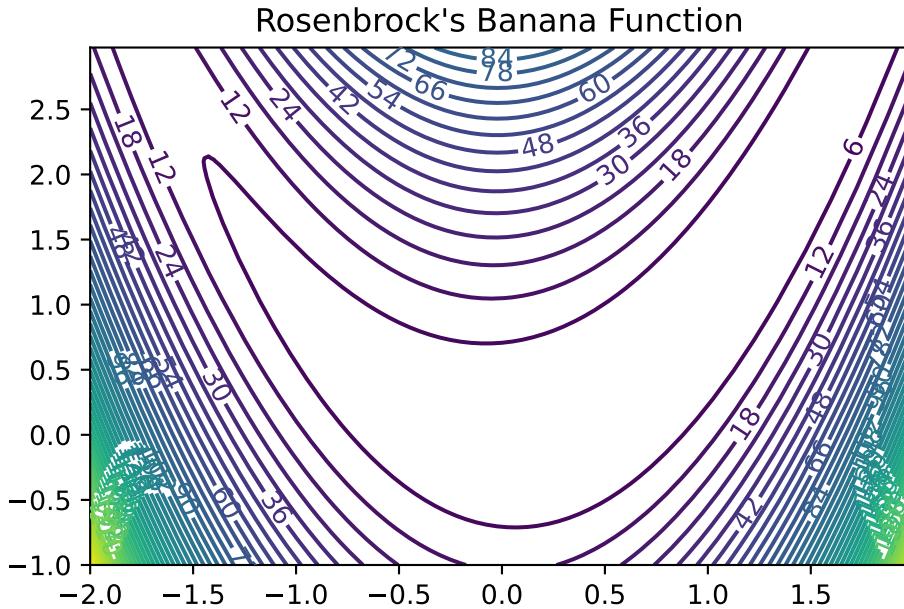
```

import numpy as np
import matplotlib.cm as cm
import matplotlib.pyplot as plt

delta = 0.025
x1 = np.arange(-2.0, 2.0, delta)
x2 = np.arange(-1.0, 3.0, delta)
X1, X2 = np.meshgrid(x1, x2)
Y = fun_rosen(X1, X2)
fig, ax = plt.subplots()
CS = ax.contour(X1, X2, Y , 50)
ax.clabel(CS, inline=True, fontsize=10)
ax.set_title("Rosenbrock's Banana Function")

```

Text(0.5, 1.0, "Rosenbrock's Banana Function")



- Visual inspection: yield is optimized near  $(\xi_1, \xi_2)$

### 5.1.1 Visualization: Problems in Practice

- True response surface is unknown in practice
- When yield evaluation is not as simple as a toy banana function, but a process requiring care to monitor, reconfigure and run, it's far too expensive to observe over a dense grid
- And, measuring yield may be a noisy/inexact process
- That's where stats (RSM) comes in

### 5.1.2 RSM: Strategies

- RSMs consist of experimental strategies for
- **exploring** the space of the process (i.e., independent/input) variables (above  $\xi_1$  and  $\xi_2$ )
- empirical statistical **modeling** targeted toward development of an appropriate approximating relationship between the response (yield) and process variables local to a study region of interest
- **optimization** methods for sequential refinement in search of the levels or values of process variables that produce desirable responses (e.g., that maximize yield or explain variation)
- RSM used for fitting an Empirical Model

- True response surface driven by an unknown physical mechanism
- Observations corrupted by noise
- Helpful: fit an empirical model to output collected under different process configurations
- Consider response  $Y$  that depends on controllable input variables  $\xi_1, \xi_2, \dots, \xi_m$
- RSM: Equations of the Empirical Model
  - $Y = f(\xi_1, \xi_2, \dots, \xi_m) + \epsilon$
  - $\mathbb{E}\{Y\} = \eta = f(\xi_1, \xi_2, \dots, \xi_m)$
  - $\epsilon$  is treated as zero mean idiosyncratic noise possibly representing
    - \* inherent variation, or
    - \* the effect of other systems or
    - \* variables not under our purview at this time

### 5.1.3 RSM: Noise in the Empirical Model

- Typical simplifying assumption:  $\epsilon \sim N(0, \sigma^2)$
- We seek estimates for  $f$  and  $\sigma^2$  from noisy observations  $Y$  at inputs  $\xi$

### 5.1.4 RSM: Natural and Coded Variables

- Inputs  $\xi_1, \xi_2, \dots, \xi_m$  called **natural variables**:
  - expressed in natural units of measurement, e.g., degrees Celsius, pounds per square inch (psi), etc.
- Transformed to **coded variables**  $x_1, x_2, \dots, x_m$ :
  - to mitigate hassles and confusion that can arise when working with a multitude of scales of measurement
- Typical **Transformations** offering dimensionless inputs  $x_1, x_2, \dots, x_m$ 
  - in the unit cube, or
  - scaled to have a mean of zero and standard deviation of one, are common choices.
- Empirical model becomes  $\eta = f(x_1, x_2, \dots, x_m)$

### 5.1.5 RSM Low-order Polynomials

- Low-order polynomial make the following simplifying Assumptions
  - Learning about  $f$  is lots easier if we make some simplifying approximations
  - Appealing to **Taylor's theorem**, a low-order polynomial in a small, localized region of the input ( $x$ ) space is one way forward
  - Classical RSM:
    - \* disciplined application of **local analysis** and
    - \* **sequential refinement** of locality through conservative extrapolation
  - Inherently a **hands-on process**

## 5.2 First-Order Models (Main Effects Model)

- **First-order model** (sometimes called main effects model) useful in parts of the input space where it's believed that there's little curvature in  $f$ :

$$\eta = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

- For example:

$$\eta = 50 + 8x_1 + 3x_2$$

- In practice, such a surface would be obtained by fitting a model to the outcome of a designed experiment
- First-Order Model in python Evaluated on a Grid
- Evaluate model on a grid in a double-unit square centered at the origin
- Coded units are chosen arbitrarily, although one can imagine deploying this approximating function nearby  $x^{(0)} = (0, 0)$

```
def fun_1(x1,x2):  
    return 50 + 8*x1 + 3*x2
```

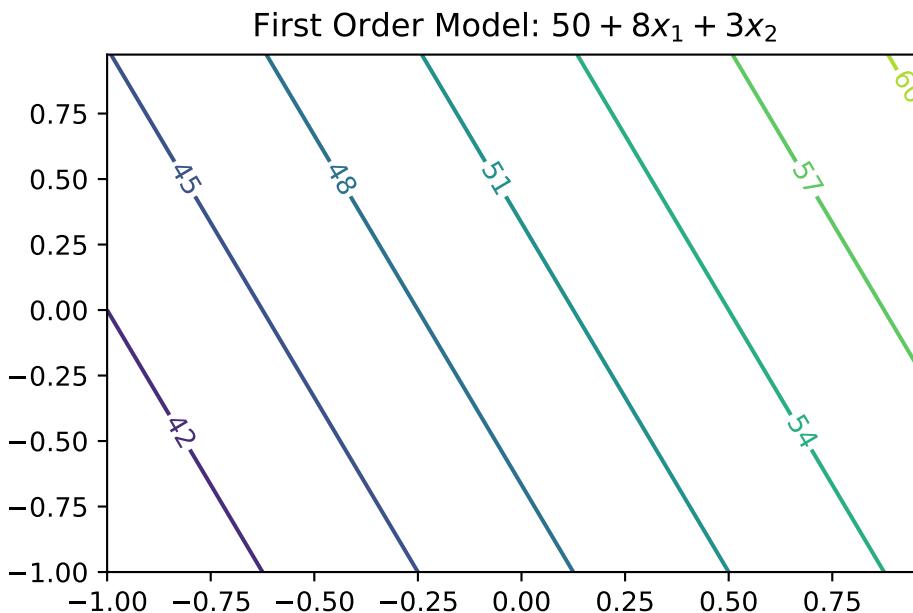
```
import numpy as np  
import matplotlib.cm as cm  
import matplotlib.pyplot as plt  
  
delta = 0.025  
x1 = np.arange(-1.0, 1.0, delta)  
x2 = np.arange(-1.0, 1.0, delta)  
X1, X2 = np.meshgrid(x1, x2)  
Y = fun_1(X1,X2)  
fig, ax = plt.subplots()
```

```

CS = ax.contour(X1, X2, Y)
ax.clabel(CS, inline=True, fontsize=10)
ax.set_title('First Order Model: $50 + 8x_1 + 3x_2$')

```

```
Text(0.5, 1.0, 'First Order Model: $50 + 8x_1 + 3x_2$')
```



### 5.2.1 First-Order Model Properties

- First-order model in 2d traces out a **plane** in  $y \times (x_1, x_2)$  space
- Only be appropriate for the most trivial of response surfaces, even when applied in a highly localized part of the input space
- Adding **curvature** is key to most applications:
  - First-order model with **interactions** induces limited degree of curvature via different rates of change of  $y$  as  $x_1$  is varied for fixed  $x_2$ , and vice versa:

$$\eta = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_{12}$$

- For example  $\eta = 50 + 8x_1 + 3x_2 - 4x_1 x_2$

### 5.2.2 First-order Model with Interactions in python

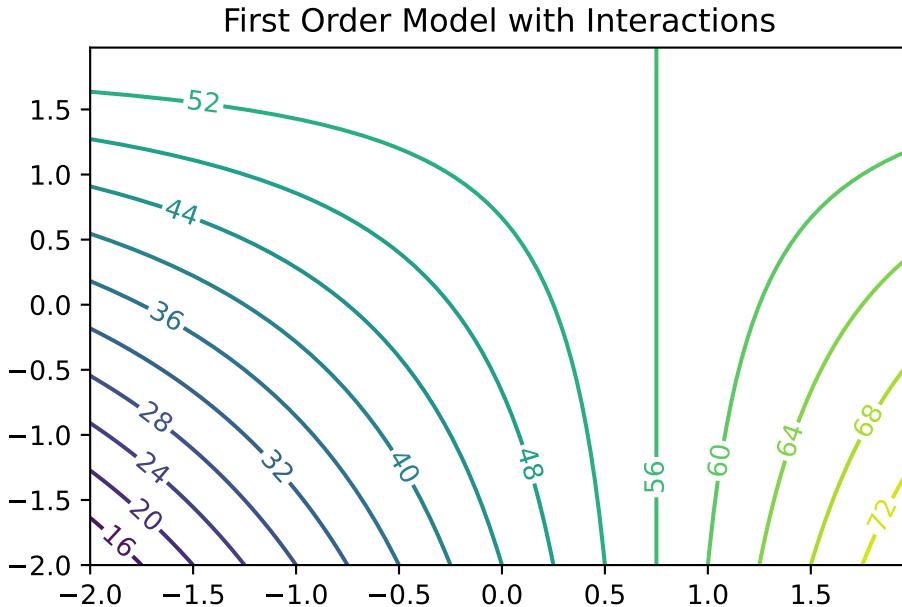
- Code below facilitates evaluations for pairs  $(x_1, x_2)$
- Responses may be observed over a mesh in the same double-unit square

```
def fun_11(x1,x2):
    return 50 + 8 * x1 + 3 * x2 - 4 * x1 * x2

import numpy as np
import matplotlib.cm as cm
import matplotlib.pyplot as plt

delta = 0.025
x1 = np.arange(-2.0, 2.0, delta)
x2 = np.arange(-2.0, 2.0, delta)
X1, X2 = np.meshgrid(x1, x2)
Y = fun_11(X1,X2)
fig, ax = plt.subplots()
CS = ax.contour(X1, X2, Y, 20)
ax.clabel(CS, inline=True, fontsize=10)
ax.set_title('First Order Model with Interactions')
```

```
Text(0.5, 1.0, 'First Order Model with Interactions')
```



### 5.2.3 Observations: First-Order Model with Interactions

- Mean response  $\eta$  is increasing marginally in both  $x_1$  and  $x_2$ , or conditional on a fixed value of the other until  $x_1$  is 0.75
- Rate of increase slows as both coordinates grow simultaneously since the coefficient in front of the interaction term  $x_1x_2$  is negative
- Compared to the first-order model (without interactions): surface is far more useful locally
- Least squares regressions often flag up significant interactions when fit to data collected on a design far from local optima

## 5.3 Second-Order Models

- Second-order model may be appropriate near local optima where  $f$  would have substantial curvature:

$$\eta = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{11} x_1^2 + \beta_{22} x_2^2 + \beta_{12} x_1 x_2$$

- For example

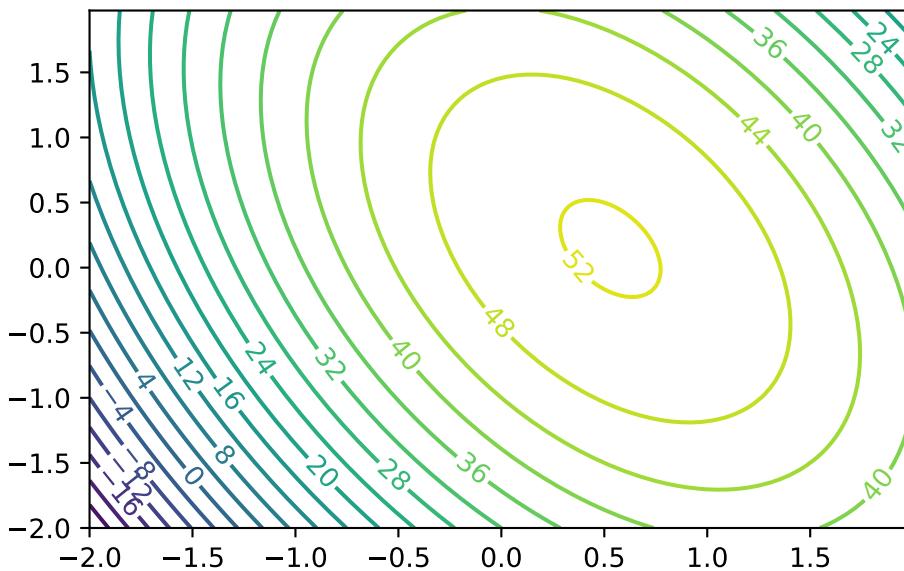
$$\eta = 50 + 8x_1 + 3x_2 - 7x_1^2 - 3x_2^2 - 4x_1 x_2$$

- Implementation of the Second-Order Model as `fun_2()`.

```
def fun_2(x1,x2):  
    return 50 + 8 * x1 + 3 * x2 - 7 * x1**2 - 3*x2**2 - 4 * x1 * x2
```

```
import numpy as np  
import matplotlib.cm as cm  
import matplotlib.pyplot as plt  
  
delta = 0.025  
x1 = np.arange(-2.0, 2.0, delta)  
x2 = np.arange(-2.0, 2.0, delta)  
X1, X2 = np.meshgrid(x1, x2)  
Y = fun_2(X1,X2)  
fig, ax = plt.subplots()  
CS = ax.contour(X1, X2, Y, 20)  
ax.clabel(CS, inline=True, fontsize=10)  
ax.set_title('Second Order Model with Interactions. Maximum near about $(0.6,0.2)$')  
  
Text(0.5, 1.0, 'Second Order Model with Interactions. Maximum near about $(0.6,0.2)$')
```

Second Order Model with Interactions. Maximum near about (0.6, 0.2)



### 5.3.1 Second-Order Models: Properties

- Not all second-order models would have a single stationary point (in RSM jargon called “a simple maximum”)
- In “yield maximizing” setting we’re presuming response surface is **concave** down from a global viewpoint
  - even though local dynamics may be more nuanced
- Exact criteria depend upon the eigenvalues of a certain matrix built from those coefficients
- Box and Draper (2007) provide a diagram categorizing all of the kinds of second-order surfaces in RSM analysis, where finding local maxima is the goal

### 5.3.2 Example: Stationary Ridge

- Example set of coefficients describing what’s called a **stationary ridge** is provided by the code below

```
def fun_ridge(x1, x2):
    return 80 + 4*x1 + 8*x2 - 3*x1**2 - 12*x2**2 - 12*x1*x2
```

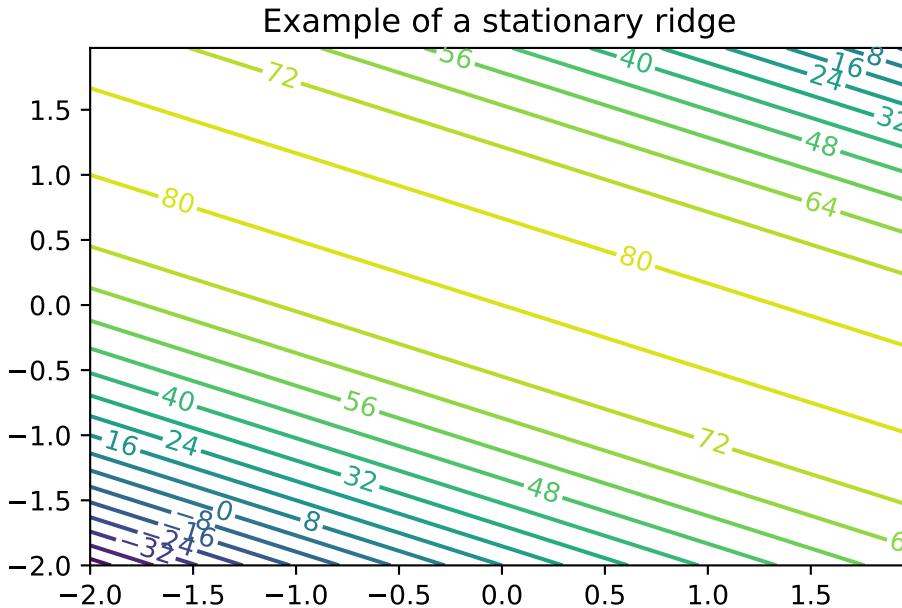
```

import numpy as np
import matplotlib.cm as cm
import matplotlib.pyplot as plt

delta = 0.025
x1 = np.arange(-2.0, 2.0, delta)
x2 = np.arange(-2.0, 2.0, delta)
X1, X2 = np.meshgrid(x1, x2)
Y = fun_ridge(X1,X2)
fig, ax = plt.subplots()
CS = ax.contour(X1, X2, Y, 20)
ax.clabel(CS, inline=True, fontsize=10)
ax.set_title('Example of a stationary ridge')

```

Text(0.5, 1.0, 'Example of a stationary ridge')



### 5.3.3 Observations: Second-Order Model (Ridge)

- **Ridge:** a whole line of stationary points corresponding to maxima
- Situation means that the practitioner has some flexibility when it comes to optimizing:
  - can choose the precise setting of  $(x_1, x_2)$  either arbitrarily or (more commonly) by consulting some tertiary criteria

### 5.3.4 Example: Rising Ridge

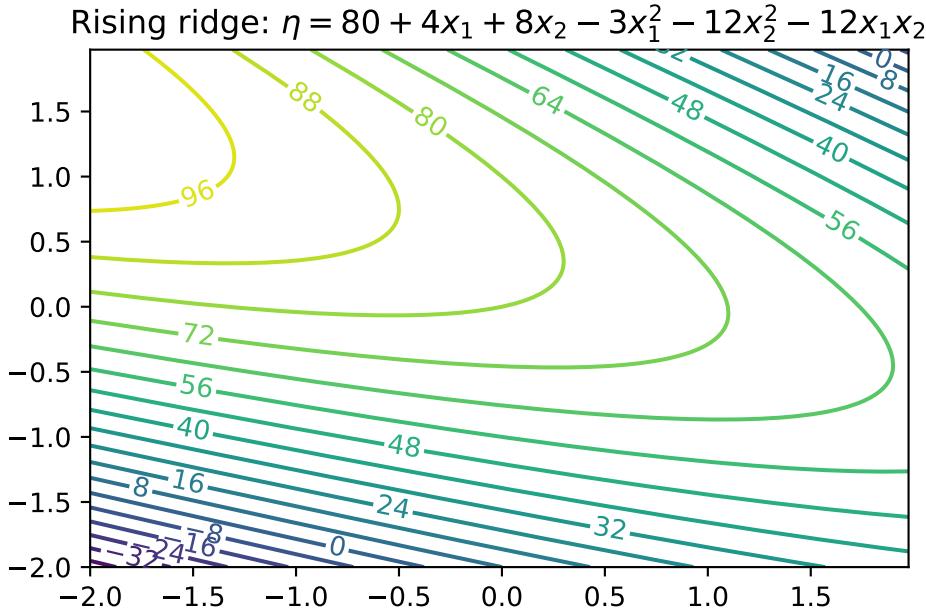
- An example of a rising ridge is implemented by the code below.

```
def fun_ridge_rise(x1, x2):
    return 80 - 4*x1 + 12*x2 - 3*x1**2 - 12*x2**2 - 12*x1*x2

import numpy as np
import matplotlib.cm as cm
import matplotlib.pyplot as plt

delta = 0.025
x1 = np.arange(-2.0, 2.0, delta)
x2 = np.arange(-2.0, 2.0, delta)
X1, X2 = np.meshgrid(x1, x2)
Y = fun_ridge_rise(X1, X2)
fig, ax = plt.subplots()
CS = ax.contour(X1, X2, Y, 20)
ax.clabel(CS, inline=True, fontsize=10)
ax.set_title('Rising ridge: $\eta = 80 + 4x_1 + 8x_2 - 3x_1^2 - 12x_2^2 - 12x_1x_2$')

Text(0.5, 1.0, 'Rising ridge: $\eta = 80 + 4x_1 + 8x_2 - 3x_1^2 - 12x_2^2 - 12x_1x_2$')
```



### 5.3.5 Summary: Rising Ridge

- The stationary point is remote to the study region
- Continuum of (local) stationary points along any line going through the 2d space, excepting one that lies directly on the ridge
- Although estimated response will increase while moving along the axis of symmetry toward its stationary point, this situation indicates
  - either a poor fit by the approximating second-order function, or
  - that the study region is not yet precisely in the vicinity of a local optima—often both.

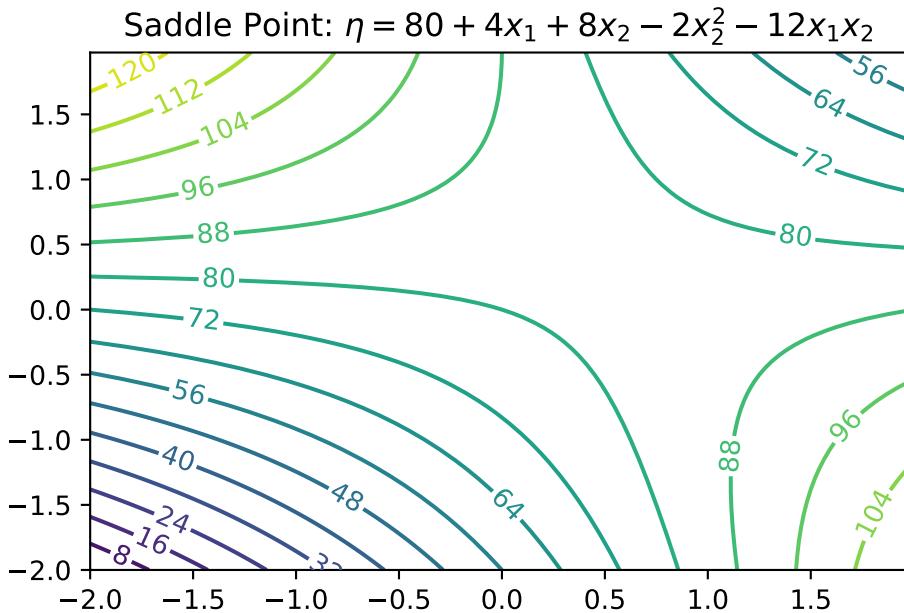
### 5.3.6 Falling Ridge

- Inversion of a rising ridge is a falling ridge
- Similarly indicating one is far from local optima, except that the response decreases as you move toward the stationary point
- Finding a falling ridge system can be a back-to-the-drawing-board affair.

### 5.3.7 Saddle Point

- Finally, we can get what's called a saddle or minimax system.

```
def fun_saddle(x1, x2):  
    return 80 + 4*x1 + 8*x2 - 2*x2**2 - 12*x1*x2  
  
import numpy as np  
import matplotlib.cm as cm  
import matplotlib.pyplot as plt  
  
delta = 0.025  
x1 = np.arange(-2.0, 2.0, delta)  
x2 = np.arange(-2.0, 2.0, delta)  
X1, X2 = np.meshgrid(x1, x2)  
Y = fun_saddle(X1,X2)  
fig, ax = plt.subplots()  
CS = ax.contour(X1, X2, Y, 20)  
ax.clabel(CS, inline=True, fontsize=10)  
ax.set_title('Saddle Point: $\eta = 80 + 4x_1 + 8x_2 - 2x_2^2 - 12x_1x_2$')  
  
Text(0.5, 1.0, 'Saddle Point: $\eta = 80 + 4x_1 + 8x_2 - 2x_2^2 - 12x_1x_2$')
```



### 5.3.8 Interpretation: Saddle Points

- Likely further data collection, and/or outside expertise, is needed before determining a course of action in this situation

### 5.3.9 Summary: Ridge Analysis

- Finding a simple maximum, or stationary ridge, represents ideals in the spectrum of second-order approximating functions
- But getting there can be a bit of a slog
- Using models fitted from data means uncertainty due to noise, and therefore uncertainty in the type of fitted second-order model
- A ridge analysis attempts to offer a principled approach to navigating uncertainties when one is seeking local maxima
- The two-dimensional setting exemplified above is convenient for visualization, but rare in practice
- Complications compound when studying the effect of more than two process variables

## 5.4 General RSM Models

- General **first-order model** on  $m$  process variables  $x_1, x_2, \dots, x_m$  is

$$\eta = \beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m$$

- General **second-order model** on  $m$  process variables

$$\eta = \beta_0 + \sum_{j=1}^m + \sum_{j=1}^m x_j^2 + \sum_{j=2}^m \sum_{k=1}^j \beta_{kj} x_k x_j.$$

### 5.4.1 Ordinary Least Squares

- Inference from data is carried out by **ordinary least squares** (OLS)
- For an excellent review including R examples, see Sheather (2009)
- OLS and maximum likelihood estimators (MLEs) are in the typical Gaussian linear modeling setup basically equivalent

## 5.5 Designs

- Important: Organize the data collection phase of a response surface study carefully
- **Design:** choice of  $x$ 's where we plan to observe  $y$ 's, for the purpose of approximating  $f$
- Analyses and designs need to be carefully matched
- When using a first-order model, some designs are preferred over others
- When using a second-order model to capture curvature, a different sort of design is appropriate
- Design choices often contain features enabling modeling assumptions to be challenged
  - e.g., to check if initial impressions are supported by the data ultimately collected

### 5.5.1 Different Designs

- **Screening designs:** determine which variables matter so that subsequent experiments may be smaller and/or more focused
- Then there are designs tailored to the form of model (first- or second-order, say) in the screened variables
- And then there are more designs still

## 5.6 RSM Experimentation

### 5.6.1 First Step

- RSM-based experimentation begins with a **first-order model**, possibly with interactions
- Presumption: current process operating **far from optimal** conditions
- Collect data and apply **method of steepest ascent** (gradient) on fitted surfaces to move to the optimum

### 5.6.2 Second Step

- Eventually, if all goes well after several such carefully iterated refinements, **second-order models** are used on appropriate designs in order to zero-in on ideal operating conditions
- Careful analysis of the fitted surface:
  - Ridge analysis with further refinement using gradients of, and
  - standard errors associated with, the fitted surfaces, and so on

### 5.6.3 Third Step

- Once the practitioner is satisfied with the full arc of
  - design(s),
  - fit(s), and
  - decision(s):
- A small experiment called **confirmation test** may be performed to check if the predicted optimal settings are realizable in practice

## 5.7 RSM: Review and General Considerations

- First Glimpse, RSM seems sensible, and pretty straightforward as quantitative statistics-based analysis goes
- But: RSM can get complicated, especially when input dimensions are not very low
- Design considerations are particularly nuanced, since the goal is to obtain reliable estimates of main effects, interaction, and curvature while minimizing sampling effort/expense
- RSM Downside: Inefficiency

- Despite intuitive appeal, several RSM downsides become apparent upon reflection
  - Problems in practice
  - Stepwise nature of sequential decision making is inefficient:
    - \* Not obvious how to re-use or update analysis from earlier phases, or couple with data from other sources/related experiments
- RSM Downside: Locality
  - In addition to being local in experiment-time (stepwise approach), it's local in experiment-space
  - Balance between
    - \* exploration (maybe we're barking up the wrong tree) and
    - \* exploitation (let's make things a little better) is modest at best
- RSM Downside: Expert Knowledge
  - Interjection of expert knowledge is limited to hunches about relevant variables (i.e., the screening phase), where to initialize search, how to design the experiments
  - Yet at the same time classical RSMs rely heavily on constant examination throughout stages of modeling and design and on the instincts of seasoned practitioners
- RSM Downside: Replicability
  - Parallel analyses, conducted according to the same best intentions, rarely lead to the same designs, model fits and so on
  - Sometimes that means they lead to different conclusions, which can be cause for concern

### **5.7.1 Historical Considerations about RSM**

- In spite of those criticisms, however, there was historically little impetus to revise the status quo
- Classical RSM was comfortable in its skin, consistently led to improvements or compelling evidence that none can reasonably be expected
- But then in the late 20th century came an explosive expansion in computational capability, and with it a means of addressing many of those downsides

### **5.7.2 Status Quo**

- Nowadays, field experiments and statistical models, designs and optimizations are coupled with mathematical models
- Simple equations are not regarded as sufficient to describe real-world systems anymore

- Physicists figured that out fifty years ago; industrial engineers followed, biologists, social scientists, climate scientists and weather forecasters, etc.
- Systems of equations are required, solved over meshes (e.g., finite elements), or stochastically interacting agents
- Goals for those simulation experiments are as diverse as their underlying dynamics
- Optimization of systems is common, e.g., to identify worst-case scenarios

### 5.7.3 The Role of Statistics

- Solving systems of equations, or interacting agents, requires computing
- Statistics involved at various stages:
  - choosing the mathematical model
  - solving by stochastic simulation (Monte Carlo)
  - designing the computer experiment
  - smoothing over idiosyncrasies or noise
  - finding optimal conditions, or
  - calibrating mathematical/computer models to data from field experiments

### 5.7.4 New RSM is needed: DACE

- Classical RSMs are not well-suited to any of those tasks, because
  - they lack the fidelity required to model these data
  - their intended application is too local
  - they're also too hands-on.
- Once computers are involved, a natural inclination is to automate—to remove humans from the loop and set the computer running on the analysis in order to maximize computing throughput, or minimize idle time
- **Design and Analysis of Computer Experiments** as a modern extension of RSM
- Experimentation is changing due to advances in machine learning
- **Gaussian process (GP) regression** is the canonical surrogate model
- Origins in geostatistics (gold mining)
- Wide applicability in contexts where prediction is king
- Machine learners exposed GPs as powerful predictors for all sorts of tasks:
  - from regression to classification,
  - active learning/sequential design,
  - reinforcement learning and optimization,
  - latent variable modeling, and so on

## 5.8 Exercises

1. Generate 3d Plots for the Contour Plots in this notebook.
2. Write a `plot_3d` function, that takes the objective function `fun` as an argument.
  - It should provide the following interface: `plot_3d(fun)`.
3. Write a `plot_contour` function, that takes the objective function `fun` as an argument:
  - It should provide the following interface: `plot_contour(fun)`.
4. Consider further arguments that might be useful for both function, e.g., ranges, size, etc.

## 5.9 Jupyter Notebook

 Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

# 6 Kriging (Gaussian Process Regression)

## 6.1 DACE and RSM

Mathematical models implemented in computer codes are used to circumvent the need for expensive field data collection. These models are particularly useful when dealing with highly nonlinear response surfaces, high signal-to-noise ratios (which often involve deterministic evaluations), and a global scope. As a result, a new approach is required in comparison to Response Surface Methodology (RSM).

With the improvement in computing power and simulation fidelity, researchers gain higher confidence and a better understanding of the dynamics in physical, biological, and social systems. However, the expansion of configuration spaces and increasing input dimensions necessitates more extensive designs. High-performance computing (HPC) allows for thousands of runs, whereas previously only tens were possible. This shift towards larger models and training data presents new computational challenges.

Research questions for DACE (Design and Analysis of Computer Experiments) include how to design computer experiments that make efficient use of computation and how to meta-model computer codes to save on simulation effort. The choice of surrogate model for computer codes significantly impacts the optimal experiment design, and the preferred model-design pairs can vary depending on the specific goal.

The combination of computer simulation, design, and modeling with field data from similar real-world experiments introduces a new category of computer model tuning problems. The ultimate goal is to automate these processes to the greatest extent possible, allowing for the deployment of HPC with minimal human intervention.

One of the remaining differences between RSM and DACE lies in how they handle noise. DACE employs replication, a technique that would not be used in a deterministic setting, to separate signal from noise. Traditional RSM is best suited for situations where a substantial proportion of the variability in the data is due to noise, and where the acquisition of data values can be severely limited. Consequently, RSM is better suited for a different class of problems, aligning with its intended purposes.

Two very good texts on computer experiments and surrogate modeling are Santner, Williams, and Notz (2003) and Forrester, Sóbester, and Keane (2008). The former is the canonical reference in the statistics literature and the latter is perhaps more popular in engineering.

## 6.2 Background: Expectation, Mean, Standard Deviation

The distribution of a random vector is characterized by some indexes. One of them is the expected value, which is defined as

$$E[X] = \sum_{x \in D_X} xp_X(x) \quad \text{if } X \text{ is discrete}$$

$$E[X] = \int_{x \in D_X} xf_X(x)dx \quad \text{if } X \text{ is continuous.}$$

The mean,  $\mu$ , of a probability distribution is a measure of its central tendency or location. That is,  $E(X)$  is defined as the average of all possible values of  $X$ , weighted by their probabilities.

### i Example: Expectation

Let  $X$  denote the number produced by rolling a fair die. Then

$$E(X) = 1 \times 1/6 + 2 \times 1/6 + 3 \times 1/6 + 4 \times 1/6 + 5 \times 1/6 + 6 \times 1/6 = 3.5.$$

### 6.2.1 Sample Mean

The sample mean is an important estimate of the population mean. The sample mean of a sample  $\{x_i\}$  ( $i = 1, 2, \dots, n$ ) is defined as

$$\bar{x} = \frac{1}{n} \sum_i x_i.$$

### 6.2.2 Variance and Standard Deviation

If we are trying to predict the value of a random variable  $X$  by its mean  $\mu = E(X)$ , the error will be  $X - \mu$ . In many situations it is useful to have an idea how large this deviation or error is. Since  $E(X - \mu) = E(X) - \mu = 0$ , it is necessary to use the absolute value or the square of  $(X - \mu)$ . The squared error is the first choice, because the derivatives are easier to calculate. These considerations motivate the definition of the variance:

The variance of a random variable  $X$  is the mean squared deviation of  $X$  from its expected value  $\mu = E(X)$ .

$$Var(X) = E[(X - \mu)^2]. \tag{6.1}$$

### 6.2.3 Standard Deviation

Taking the square root of the variance to get back to the same scale of units as  $X$  gives the standard deviation. The standard deviation of  $X$  is the square root of the variance of  $X$ .

$$sd(X) = \sqrt{Var(X)}. \quad (6.2)$$

### 6.2.4 Calculation of the Standard Deviation with Python

The function `numpy.std` returns the standard deviation, a measure of the spread of a distribution, of the array elements. The argument `ddof` specifies the Delta Degrees of Freedom. The divisor used in calculations is  $N - ddof$ , where  $N$  represents the number of elements. By default `ddof` is zero, i.e., `std` uses the formula

$$\sqrt{\frac{1}{N} \sum_i (x_i - \bar{x})^2} \quad \text{with} \quad \bar{x} = \sum_{i=1}^N x_i / N. \quad (6.3)$$

#### i Example: Standard Deviation with Python

Consider the array `[1, 2, 3]`: Since  $\bar{x} = 2$ , the following value is computed:

$$\sqrt{1/3 \times ((1-2)^2 + (2-2)^2 + (3-2)^2)} = \sqrt{2/3}.$$

```
import numpy as np
a = np.array([[1, 2, 3]])
np.std(a)
```

0.816496580927726

### 6.2.5 The Empirical Standard Deviation

The empirical standard deviation (which uses  $N-1$ ),  $\sqrt{1/2 \times ((1-2)^2 + (2-2)^2 + (3-2)^2)} = \sqrt{2/2}$ , can be calculated as follows:

```
np.std(a, ddof=1)
```

1.0

### 6.2.6 The Argument “axis”

**i** Axes along which the standard deviation is computed

- When you compute `np.std` with `axis=0`, it calculates the standard deviation along the vertical axis, meaning it computes the standard deviation for each column of the array.
- On the other hand, when you compute `np.std` with `axis=1`, it calculates the standard deviation along the horizontal axis, meaning it computes the standard deviation for each row of the array.
- If the `axis` parameter is not specified, `np.std` computes the standard deviation of the flattened array.

```
A = np.array([[1, 2], [3, 4]])  
A
```

```
array([[1, 2],  
       [3, 4]])
```

```
np.std(A)
```

```
1.118033988749895
```

```
np.std(A, axis=0)
```

```
array([1., 1.])
```

```
np.std(A, axis=1)
```

```
array([0.5, 0.5])
```

## 6.3 Data Types and Precision in Python

We consider single versus double precision in Python. In single precision, `std()` can be inaccurate:

```
a = np.zeros((2, 4*4), dtype=np.float32)
a[0, :] = 1.0
a[1, :] = 0.1
a
```

```
array([[1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ,
       1. , 1. , 1. ],
       [0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1,
        0.1, 0.1, 0.1]], dtype=float32)
```

```
np.std(a, axis=0)
```

```
array([0.45, 0.45, 0.45, 0.45, 0.45, 0.45, 0.45, 0.45, 0.45, 0.45,
       0.45, 0.45, 0.45, 0.45], dtype=float32)
```

```
np.std(a, axis=1)
```

```
array([0., 0.], dtype=float32)
```

```
abs(0.45 - np.std(a))
```

```
1.7881393421514957e-08
```

### i Float data types

- float32 and float64 are data types in numpy that specify the precision of floating point numbers.
- float32 is a single-precision floating point number that occupies 32 bits of memory. It has a precision of about 7 decimal digits.
- float64 is a double-precision floating point number that occupies 64 bits of memory. It has a precision of about 15 decimal digits.
- The main difference between float32 and float64 is the precision and memory usage. float64 provides a higher precision but uses more memory, while float32 uses less memory but has a lower precision.

Computing the standard deviation in float64 is more accurate (result may vary), see <https://numpy.org/devdocs/reference/generated/numpy.std.html>.

```
abs(0.45 - np.std(a, dtype=np.float64))
```

7.450580707946131e-10

**i** Example: 32 versus 64 bit

```
import numpy as np

# Define a number
num = 0.123456789123456789

# Convert to float32 and float64
num_float32 = np.float32(num)
num_float64 = np.float64(num)

# Print the number in both formats
print("float32: ", num_float32)
print("float64: ", num_float64)

float32: 0.12345679
float64: 0.12345678912345678
```

The float32 data type in numpy represents a single-precision floating point number. It uses 32 bits of memory, which gives it a precision of about 7 decimal digits. On the other hand, float64 represents a double-precision floating point number. It uses 64 bits of memory, which gives it a precision of about 15 decimal digits.

The reason float32 shows fewer digits is because it has less precision due to using less memory. The bits of memory are used to store the sign, exponent, and fraction parts of the floating point number, and with fewer bits, you can represent fewer digits accurately.

## 6.4 Distributions and Random Numbers in Python

Results from computers are deterministic, so it sounds like a contradiction in terms to generate random numbers on a computer. Standard computers generate pseudo-random numbers, i.e., numbers that behave as if they were drawn randomly.

## Deterministic Random Numbers

- Idea: Generate deterministically numbers that **look** (behave) as if they were drawn randomly.

### 6.4.1 The Uniform Distribution

The probability density function of the uniform distribution is defined as:

$$f_X(x) = \frac{1}{b-a} \quad \text{for } x \in [a, b].$$

Generate 10 random numbers from a uniform distribution between  $a = 0$  and  $b = 1$ :

```
import numpy as np
# Initialize the random number generator
rng = np.random.default_rng(seed=123456789)
n = 10
x = rng.uniform(low=0.0, high=1.0, size=n)
x

array([0.02771274, 0.90670006, 0.88139355, 0.62489728, 0.79071481,
       0.82590801, 0.84170584, 0.47172795, 0.95722878, 0.94659153])
```

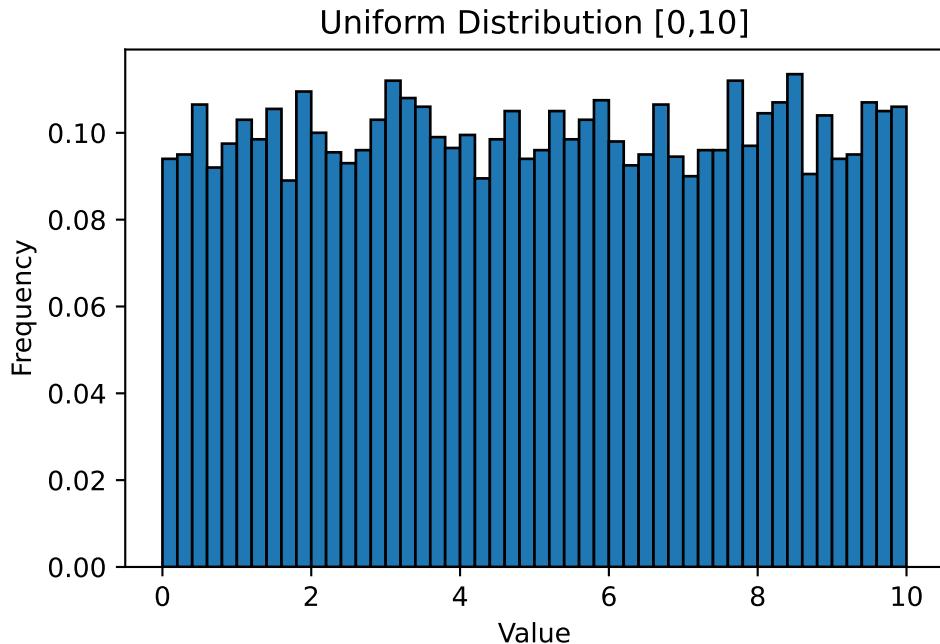
Generate 10,000 random numbers from a uniform distribution between 0 and 10 and plot a histogram of the numbers:

```
import numpy as np
import matplotlib.pyplot as plt

# Initialize the random number generator
rng = np.random.default_rng(seed=123456789)

# Generate random numbers from a uniform distribution
x = rng.uniform(low=0, high=10, size=10000)

# Plot a histogram of the numbers
plt.hist(x, bins=50, density=True, edgecolor='black')
plt.title('Uniform Distribution [0,10]')
plt.xlabel('Value')
plt.ylabel('Frequency')
plt.show()
```



#### 6.4.2 The Normal Distribution

The probability density function of the normal distribution is defined as:

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right), \quad (6.4)$$

where:  $\mu$  is the mean;  $\sigma$  is the standard deviation.

To generate ten random numbers from a normal distribution, the following command can be used.

```
# generate 10 random numbers between from a normal distribution
import numpy as np
rng = np.random.default_rng()
n = 10
mu, sigma = 2, 0.1
x = rng.normal(mu, sigma, n)
x
```

```
array([2.16246891, 2.05965816, 2.11034775, 1.93658891, 2.10379593,
       1.95185597, 1.96748382, 2.05627427, 1.85056219, 2.00557113])
```

Verify the mean:

```
abs(mu - np.mean(x))
```

0.02046070417147572

Note: To verify the standard deviation, we use `ddof = 1` (empirical standard deviation):

```
abs(sigma - np.std(x, ddof=1))
```

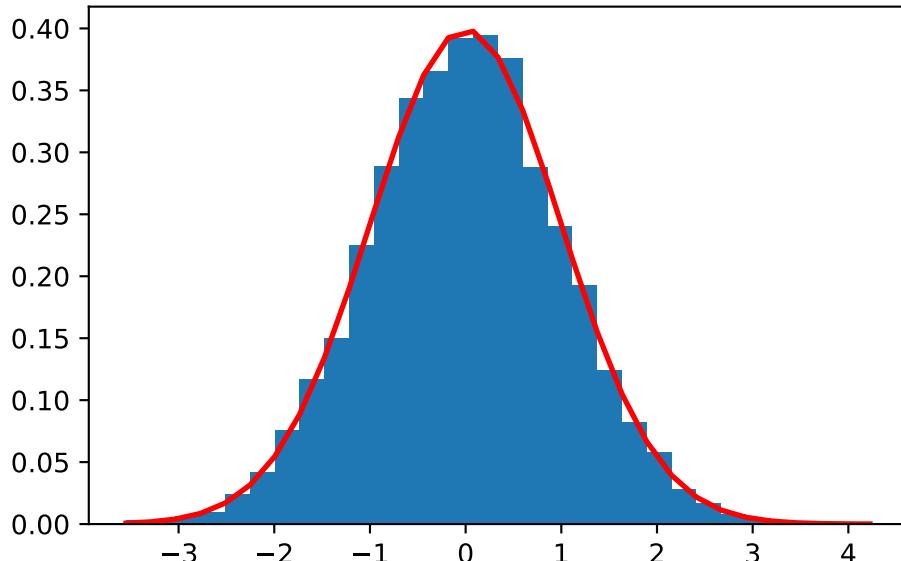
0.004754897102072941

A normally distributed random variable is a random variable whose associated probability distribution is the normal (or Gaussian) distribution. The normal distribution is a continuous probability distribution characterized by a symmetric bell-shaped curve.

The distribution is defined by two parameters: the mean  $\mu$  and the standard deviation  $\sigma$ . The mean indicates the center of the distribution, while the standard deviation measures the spread or dispersion of the distribution.

This distribution is widely used in statistics and the natural and social sciences as a simple model for random variables with unknown distributions.

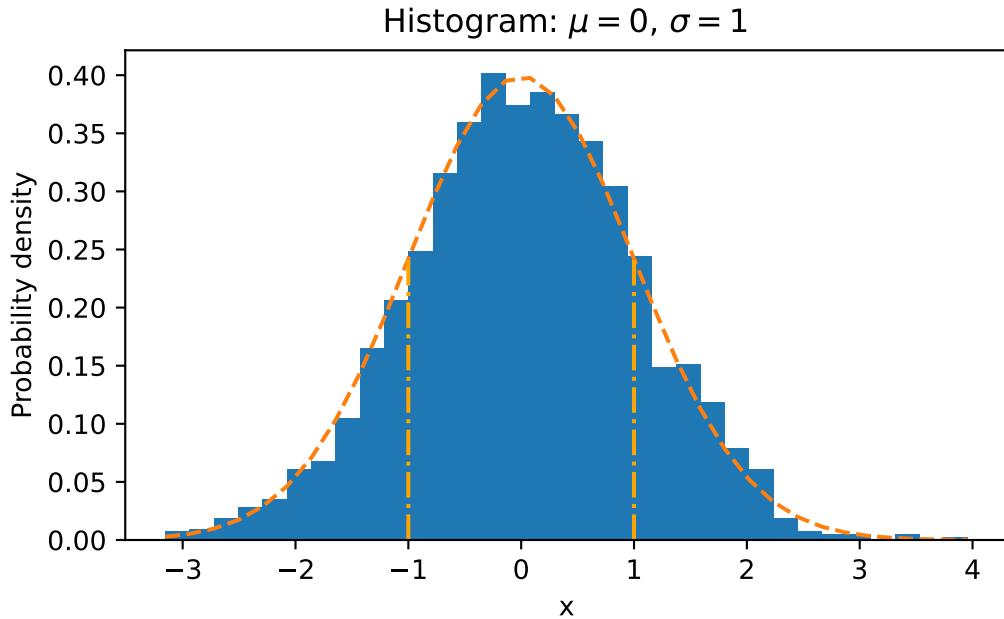
```
plot_normal_distribution(mu=0, sigma=1, num_samples=10000)
```



### 6.4.3 Visualization of the Standard Deviation

The standard deviation of normal distributed can be visualized in terms of the histogram of  $X$ :

- about 68% of the values will lie in the interval within one standard deviation of the mean
- 95% lie within two standard deviation of the mean
- and 99.9% lie within 3 standard deviations of the mean.



### 6.4.4 Standardization of Random Variables

To compare statistical properties of random variables which use different units, it is a common practice to transform these random variables into standardized variables. If a random variable  $X$  has expectation  $E(X) = \mu$  and standard deviation  $sd(X) = \sigma > 0$ , the random variable

$$X^* = (X - \mu)/\sigma$$

is called  $X$  in standard units. It has  $E(X^*) = 0$  and  $sd(X^*) = 1$ .

### 6.4.5 Realizations of a Normal Distribution

Realizations of a normal distribution refers to the actual values that you get when you draw samples from a normal distribution. Each sample drawn from the distribution is a realization of that distribution.

For example, if you have a normal distribution with a mean of 0 and a standard deviation of 1, each number you draw from that distribution is a realization.

Here's a Python example:

```
import numpy as np

# Define the parameters of the normal distribution
mu = 0
sigma = 1

# Draw 10 samples (realizations) from the normal distribution
realizations = np.random.normal(mu, sigma, 10)

print(realizations)
```

```
[ 0.48951662  0.23879586 -0.44811181 -0.610795   -2.02994507  0.60794659
 -0.35410888  0.15258149  0.50127485 -0.78640277]
```

In this code, `np.random.normal` generates 10 realizations of a normal distribution with a mean of 0 and a standard deviation of 1. The `realizations` array contains the actual values drawn from the distribution.

#### 6.4.6 The Multivariate Normal Distribution

The multivariate normal, multinormal, or Gaussian distribution serves as a generalization of the one-dimensional normal distribution to higher dimensions. We will consider  $k$ -dimensional random vectors  $X = (X_1, X_2, \dots, X_k)$ . When drawing samples from this distribution, it results in a set of values represented as  $\{x_1, x_2, \dots, x_k\}$ . To fully define this distribution, it is necessary to specify its mean  $\mu$  and covariance matrix  $\Sigma$ . These parameters are analogous to the mean, which represents the central location, and the variance (squared standard deviation) of the one-dimensional normal distribution introduced in Equation 6.4.

In the context of the multivariate normal distribution, the mean takes the form of a coordinate within an  $k$ -dimensional space. This coordinate represents the location where samples are most likely to be generated, akin to the peak of the bell curve in a one-dimensional or univariate normal distribution.

##### Covariance of two random variables

For two random variables  $X$  and  $Y$ , the covariance is defined as the expected value (or

mean) of the product of their deviations from their individual expected values:

$$\text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

The covariance within the multivariate normal distribution denotes the extent to which two variables vary together. The elements of the covariance matrix, such as  $\Sigma_{ij}$ , represent the covariances between the variables  $x_i$  and  $x_j$ . These covariances describe how the different variables in the distribution are related to each other in terms of their variability. The probability density function (PDF) of the multivariate normal distribution is defined as:

$$f_X(x) = \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right),$$

where:  $\mu$  is the  $k \times 1$  mean vector;  $\Sigma$  is the  $k \times k$  covariance matrix. The covariance matrix  $\Sigma$  is assumed to be positive definite, so that its determinant is strictly positive. For discrete random variables, covariance can be written as:

$$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - E(X))(y_i - E(Y)).$$

Figure 6.1 shows draws from a bivariate normal distribution with  $\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  and  $\Sigma = \begin{pmatrix} 9 & 4 \\ 4 & 9 \end{pmatrix}$ .

```
import numpy as np
rng = np.random.default_rng()
import matplotlib.pyplot as plt
mean = [0, 0]
cov = [[9, 4], [4, 9]] # diagonal covariance
x, y = rng.multivariate_normal(mean, cov, 1000).T
# Create a scatter plot of the numbers
plt.scatter(x, y, s=2)
plt.axis('equal')
plt.grid()
plt.title(f"Bivariate Normal. Mean zero and positive covariance: {cov}")
plt.show()
```

Bivariate Normal. Mean zero and positive covariance:  $[[9, 4], [4, 9]]$

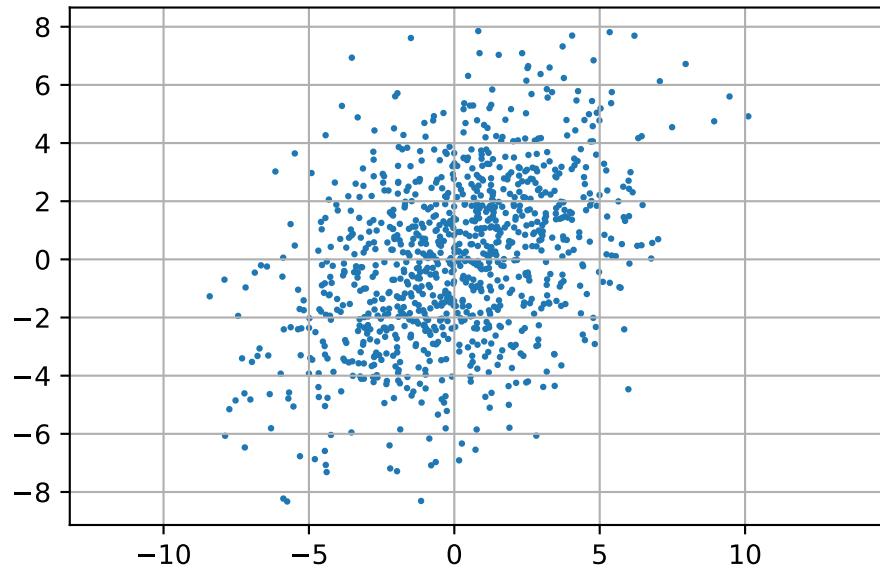


Figure 6.1: Bivariate Normal. Mean zero and covariance  $\Sigma = \begin{pmatrix} 9 & 4 \\ 4 & 9 \end{pmatrix}$

The covariance matrix of a bivariate normal distribution determines the shape, orientation, and spread of the distribution in the two-dimensional space.

The diagonal elements of the covariance matrix ( $\sigma_1^2, \sigma_2^2$ ) are the variances of the individual variables. They determine the spread of the distribution along each axis. A larger variance corresponds to a greater spread along that axis.

The off-diagonal elements of the covariance matrix ( $\sigma_{12}, \sigma_{21}$ ) are the covariances between the variables. They determine the orientation and shape of the distribution. If the covariance is positive, the distribution is stretched along the line  $y = x$ , indicating that the variables tend to increase together. If the covariance is negative, the distribution is stretched along the line  $y = -x$ , indicating that one variable tends to decrease as the other increases. If the covariance is zero, the variables are uncorrelated and the distribution is axis-aligned.

In Figure 6.1, the variances are identical and the variables are correlated (covariance is 4), so the distribution is stretched along the line  $y = x$ .

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import multivariate_normal

# Parameters
```

```

mu = np.array([0, 0])
cov = np.array([[9, 4], [4, 9]])

# Create grid and multivariate normal
x = np.linspace(-10,10,100)
y = np.linspace(-10,10,100)
X, Y = np.meshgrid(x,y)
pos = np.empty(X.shape + (2,))
pos[:, :, 0] = X; pos[:, :, 1] = Y
rv = multivariate_normal(mu, cov)

fig = plt.figure()
ax = plt.axes(projection='3d')
surf=ax.plot_surface(X, Y, rv.pdf(pos),cmap='viridis', linewidth=0)
ax.set_xlabel('X axis')
ax.set_ylabel('Y axis')
ax.set_zlabel('Z axis')
ax.set_title('Bivariate Normal Distribution')
fig.colorbar(surf, shrink=0.5, aspect=10)
plt.show()

```

Bivariate Normal Distribution

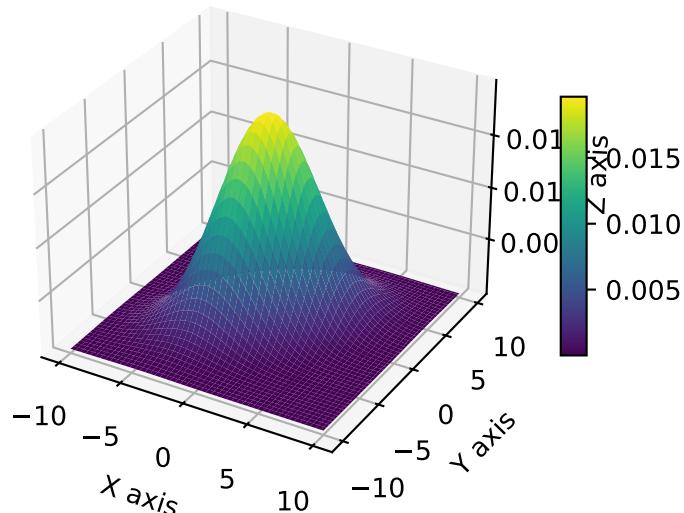


Figure 6.2: Bivariate Normal. Mean zero and covariance  $\Sigma = \begin{pmatrix} 9 & 4 \\ 4 & 9 \end{pmatrix}$

#### 6.4.7 The Bivariate Normal Distribution with Mean Zero and Zero Covariances

$$\sigma_{12} = \sigma_{21} = 0$$

$$\Sigma = \begin{pmatrix} 9 & 0 \\ 0 & 9 \end{pmatrix}$$

Bivariate Normal. Mean zero and covariance: [[9, 0], [0, 9]]

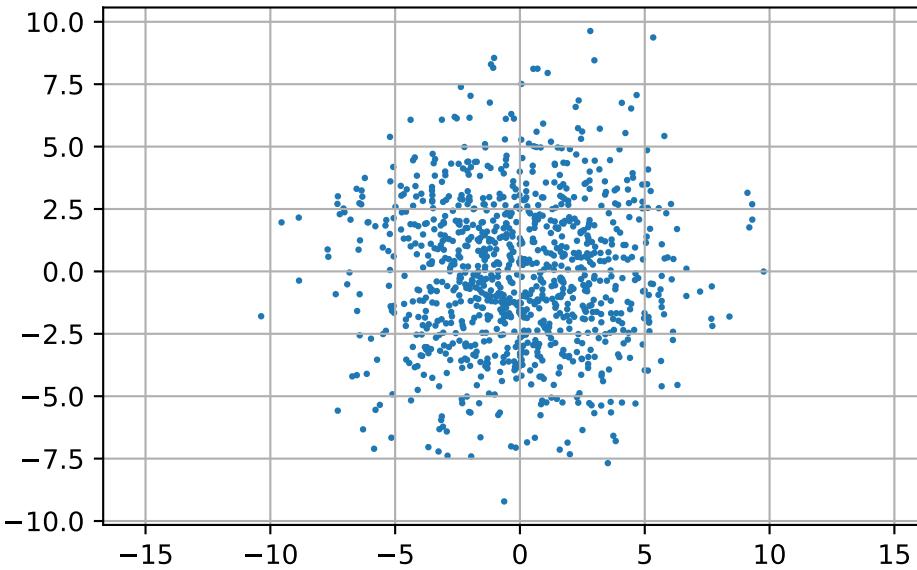


Figure 6.3: Bivariate Normal. Mean zero and covariance  $\Sigma = \begin{pmatrix} 9 & 0 \\ 0 & 9 \end{pmatrix}$

#### 6.4.8 The Bivariate Normal Distribution with Mean Zero and Negative Covariances $\sigma_{12} = \sigma_{21} = -4$

$$\Sigma = \begin{pmatrix} 9 & -4 \\ -4 & 9 \end{pmatrix}$$

### 6.5 Cholesky Decomposition and Positive Definite Matrices

The covariance matrix must be positive definite for a multivariate normal distribution for a couple of reasons:

- Semidefinite vs Definite: A covariance matrix is always symmetric and positive semidefinite. However, for a multivariate normal distribution, it must be positive definite, not

Bivariate Normal. Mean zero and covariance:  $[[9, -4], [-4, 9]]$

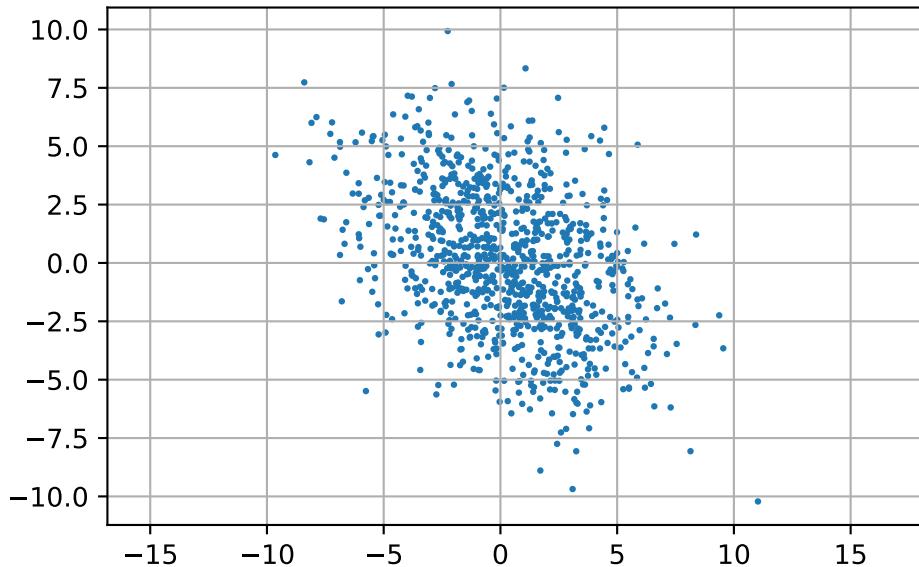


Figure 6.4: Bivariate Normal. Mean zero and covariance  $\Sigma = \begin{pmatrix} 9 & -4 \\ -4 & 9 \end{pmatrix}$

just semidefinite. This is because a positive semidefinite matrix can have zero eigenvalues, which would imply that some dimensions in the distribution have zero variance, collapsing the distribution in those dimensions. A positive definite matrix has all positive eigenvalues, ensuring that the distribution has positive variance in all dimensions.

- Invertibility: The multivariate normal distribution's probability density function involves the inverse of the covariance matrix. If the covariance matrix is not positive definite, it may not be invertible, and the density function would be undefined.

In summary, the covariance matrix being positive definite ensures that the multivariate normal distribution is well-defined and has positive variance in all dimensions.

```
import numpy as np

def is_positive_definite(matrix):
    return np.all(np.linalg.eigvals(matrix) > 0)

matrix = np.array([[9, 4], [4, 9]])
print(is_positive_definite(matrix)) # Outputs: True
```

True

More efficient (and check if symmetric) is based on Cholesky decomposition.

```
import numpy as np

def is_pd(K):
    try:
        np.linalg.cholesky(K)
        return True
    except np.linalg.linalg.LinAlgError as err:
        if 'Matrix is not positive definite' in err.message:
            return False
        else:
            raise
matrix = np.array([[9, 4], [4, 9]])
print(is_pd(matrix)) # Outputs: True
```

True

**i** Example: Cholesky decomposition.

`linalg.cholesky` computes the Cholesky decomposition of a matrix, i.e., it computes a lower triangular matrix  $L$  such that  $LL^T = A$ . If the matrix is not positive definite, an error (`LinAlgError`) is raised.

```
import numpy as np

# Define a Hermitian, positive-definite matrix
A = np.array([[9, 4], [4, 9]])

# Compute the Cholesky decomposition
L = np.linalg.cholesky(A)

print("L = \n", L)
print("L*LT = \n", np.dot(L, L.T))

L =
[[3.          0.         ]
 [1.33333333 2.68741925]]
L*LT =
[[9. 4.]
 [4. 9.]]
```

## 6.6 Maximum Likelihood Estimation: Multivariate Normal Distribution

Consider the first  $n$  terms of an identically and independently distributed (i.i.d.) sequence  $X^{(j)}$  of  $k$ -dimensional multivariate normal random vectors, i.e.,  $X^{(j)} \sim N(\mu, \Sigma)$ ,  $j = 1, 2, \dots$ . The joint probability density function of the  $j$ -th term of the sequence is

$$f_X(x_j) = \frac{1}{\sqrt{(2\pi)^k \det(\Sigma)}} \exp\left(-\frac{1}{2}(x_j - \mu)^T \Sigma^{-1} (x_j - \mu)\right),$$

where:  $\mu$  is the  $k \times 1$  mean vector;  $\Sigma$  is the  $k \times k$  covariance matrix. The covariance matrix  $\Sigma$  is assumed to be positive definite, so that its determinant is strictly positive. We use  $x_1, \dots, x_n$ , i.e., the realizations of the first  $n$  random vectors in the sequence, to estimate the two unknown parameters  $\mu$  and  $\Sigma$ .

The likelihood function is defined as the joint probability density function of the observed data, viewed as a function of the unknown parameters. Since the terms in the sequence are independent, their joint density is equal to the product of their marginal densities. As a consequence, the likelihood function can be written as the product of the individual densities:

$$\begin{aligned} L(\mu, \Sigma) &= \prod_{j=1}^n f_X(x_j) = \prod_{j=1}^n \frac{1}{\sqrt{(2\pi)^k \det(\Sigma)}} \exp\left(-\frac{1}{2}(x_j - \mu)^T \Sigma^{-1} (x_j - \mu)\right) \\ &= \frac{1}{(2\pi)^{nk/2} \det(\Sigma)^{n/2}} \exp\left(-\frac{1}{2} \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)\right). \end{aligned}$$

The log-likelihood function is

$$\ell(\mu, \Sigma) = -\frac{nk}{2} \ln(2\pi) - \frac{n}{2} \ln(\det(\Sigma)) - \frac{1}{2} \sum_{j=1}^n (x_j - \mu)^T \Sigma^{-1} (x_j - \mu).$$

The likelihood function is well-defined only if  $\det(\Sigma) > 0$ .

## 6.7 Introduction to Gaussian Processes

The concept of GP (Gaussian Process) regression can be understood as a simple extension of linear modeling. It is worth noting that this approach goes by various names and acronyms, including “kriging,” a term derived from geostatistics, as introduced by Matheron in 1963. Additionally, it is referred to as Gaussian spatial modeling or a Gaussian stochastic process, and machine learning (ML) researchers often use the term Gaussian process regression (GPR). In all of these instances, the central focus is on regression. This involves training on both inputs

and outputs, with the ultimate objective of making predictions and quantifying uncertainty (referred to as uncertainty quantification or UQ).

However, it's important to emphasize that GPs are not a universal solution for every problem. Specialized tools may outperform GPs in specific, non-generic contexts, and GPs have their own set of limitations that need to be considered.

### 6.7.1 Gaussian Process Prior

In the context of GP, any finite collection of realizations, which is represented by  $n$  observations, is modeled as having a multivariate normal (MVN) distribution. The characteristics of these realizations can be fully described by two key parameters:

1. Their mean, denoted as an  $n$ -vector  $\mu$ .
2. The covariance matrix, denoted as an  $n \times n$  matrix  $\Sigma$ . This covariance matrix encapsulates the relationships and variability between the individual realizations within the collection.

### 6.7.2 Covariance Function

The covariance function is defined by inverse exponentiated squared Euclidean distance:

$$\Sigma(\vec{x}, \vec{x}') = \exp\{-||\vec{x} - \vec{x}'||^2\},$$

where  $\vec{x}$  and  $\vec{x}'$  are two points in the  $k$ -dimensional input space and  $\|\cdot\|$  denotes the Euclidean distance, i.e.,

$$||\vec{x} - \vec{x}'||^2 = \sum_{i=1}^k (x_i - x'_i)^2.$$

An 1-d example is shown in Figure 6.5.

```
visualize_inverse_exp_squared_distance(5, 0.0, [0.5, 1, 2.0])
```

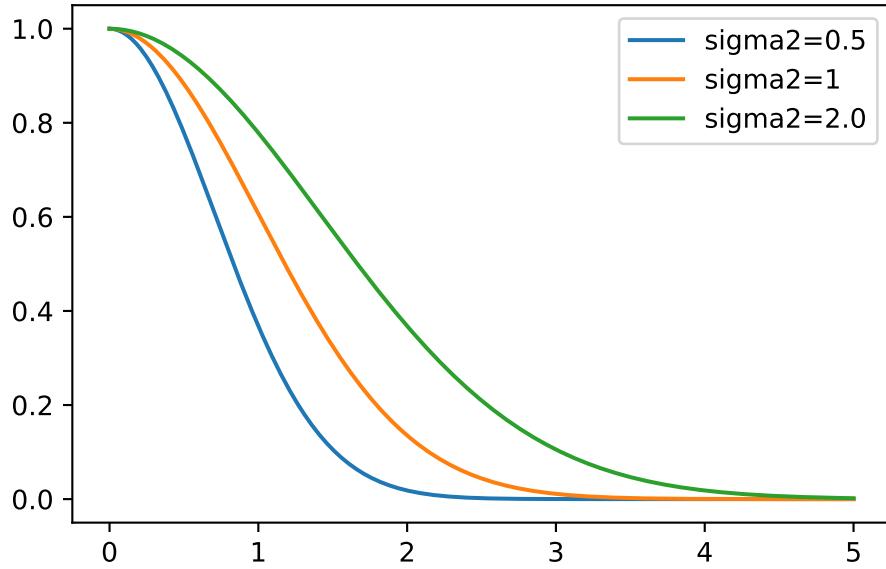


Figure 6.5: One-dim inverse exponentiated squared Euclidean distance

The covariance function is also referred to as the kernel function. The *Gaussian* kernel uses an additional parameter,  $\sigma^2$ , to control the rate of decay. This parameter is referred to as the length scale or the characteristic length scale. The covariance function is then defined as

$$\Sigma(\vec{x}, \vec{x}') = \exp\{-||\vec{x} - \vec{x}'||^2/(2\sigma^2)\}. \quad (6.5)$$

The covariance decays exponentially fast as  $\vec{x}$  and  $\vec{x}'$  become farther apart. Observe that

$$\Sigma(\vec{x}, \vec{x}) = 1$$

and

$$\Sigma(\vec{x}, \vec{x}') < 1$$

for  $\vec{x} \neq \vec{x}'$ . The function  $\Sigma(\vec{x}, \vec{x}')$  must be positive definite.

### i Positive Definiteness

Positive definiteness in the context of the covariance matrix  $\Sigma_n$  is a fundamental requirement. It is determined by evaluating  $\Sigma(x_i, x_j)$  at pairs of  $n$   $\vec{x}$ -values, denoted as  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ . The condition for positive definiteness is that for all  $\vec{x}$  vectors that are not equal to zero, the expression  $\vec{x}^\top \Sigma_n \vec{x}$  must be greater than zero. This property is

essential when intending to use  $\Sigma_n$  as a covariance matrix in multivariate normal (MVN) analysis. It is analogous to the requirement in univariate Gaussian distributions where the variance parameter,  $\sigma^2$ , must be positive.

Gaussian Processes (GPs) can be effectively utilized to generate random data that follows a smooth functional relationship. The process involves the following steps:

1. Select a set of  $\vec{x}$ -values, denoted as  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$ .
2. Define the covariance matrix  $\Sigma_n$  by evaluating  $\Sigma_n^{ij} = \Sigma(\vec{x}_i, \vec{x}_j)$  for  $i, j = 1, 2, \dots, n$ .
3. Generate an  $n$ -variate realization  $Y$  that follows a multivariate normal distribution with a mean of zero and a covariance matrix  $\Sigma_n$ , expressed as  $Y \sim \mathcal{N}_n(0, \Sigma_n)$ .
4. Visualize the result by plotting it in the  $x$ - $y$  plane.

### 6.7.3 Construction of the Covariance Matrix

Here is an one-dimensional example. The process begins by creating an input grid using  $\vec{x}$ -values. This grid consists of 100 elements, providing the basis for further analysis and visualization.

```
import numpy as np
n = 100
X = np.linspace(0, 10, n, endpoint=False).reshape(-1,1)
```

In the context of this discussion, the construction of the covariance matrix, denoted as  $\Sigma_n$ , relies on the concept of inverse exponentiated squared Euclidean distances. However, it's important to note that a modification is introduced later in the process. Specifically, the diagonal of the covariance matrix is augmented with a small value, represented as "eps" or  $\epsilon$ .

The reason for this augmentation is that while inverse exponentiated distances theoretically ensure the covariance matrix's positive definiteness, in practical applications, the matrix can sometimes become numerically ill-conditioned. By adding a small value to the diagonal, such as  $\epsilon$ , this ill-conditioning issue is mitigated. In this context,  $\epsilon$  is often referred to as "jitter."

```
import numpy as np
from numpy import array, zeros, power, ones, exp, multiply, eye, linspace, mat, spacing, sqrt
from numpy.linalg import cholesky, solve
from numpy.random import multivariate_normal
def build_Sigma(X, sigma2):
    n = X.shape[0]
    k = X.shape[1]
    D = zeros((k, n, n))
```

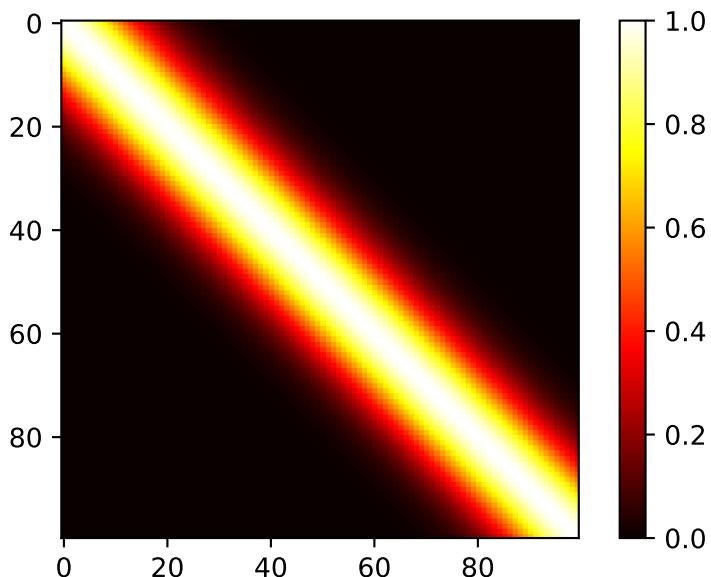
```

for l in range(k):
    for i in range(n):
        for j in range(i, n):
            D[l, i, j] = 1/(2*sigma2[l])*(X[i,l] - X[j,l])**2
D = sum(D)
D = D + D.T
return exp(-D)

```

```
sigma2 = np.array([1.0])
Sigma = build_Sigma(X, sigma2)
np.round(Sigma[:3,:], 3)
```

```
import matplotlib.pyplot as plt
plt.imshow(Sigma, cmap='hot', interpolation='nearest')
plt.colorbar()
plt.show()
```



#### 6.7.4 Generation of Random Samples and Plotting the Realizations of the Random Function

In the context of the multivariate normal distribution, the next step is to utilize the previously constructed covariance matrix denoted as `Sigma`. It is used as an essential component in generating random samples from the multivariate normal distribution.

The function `multivariate_normal` is employed for this purpose. It serves as a random number generator specifically designed for the multivariate normal distribution. In this case, the mean of the distribution is set equal to `mean`, and the covariance matrix is provided as `Psi`.

The argument `size` specifies the number of realizations, which, in this specific scenario, is set to one.

By default, the mean vector is initialized to zero. To match the number of samples, which is equivalent to the number of rows in the `X` and `Sigma` matrices, the argument `zeros(n)` is used, where `n` represents the number of samples (here taken from the size of the matrix, e.g.,: `Sigma.shape[0]`).

```
rng = np.random.default_rng(seed=12345)
```

```
Y = rng.multivariate_normal(zeros(Sigma.shape[0]), Sigma, size = 1, check_valid="raise").reshape(1, -1)
```

```
(100, 1)
```

Now we can plot the results, i.e., a finite realization of the random function  $Y()$  under a GP prior with a particular covariance structure. We will plot those `X` and `Y` pairs as connected points on an  $x$ - $y$  plane.

```
import matplotlib.pyplot as plt
plt.plot(X, Y)
plt.title("Realization of Random Functions under a GP prior.\n sigma2: {}".format(sigma2[0]))
plt.show()
```

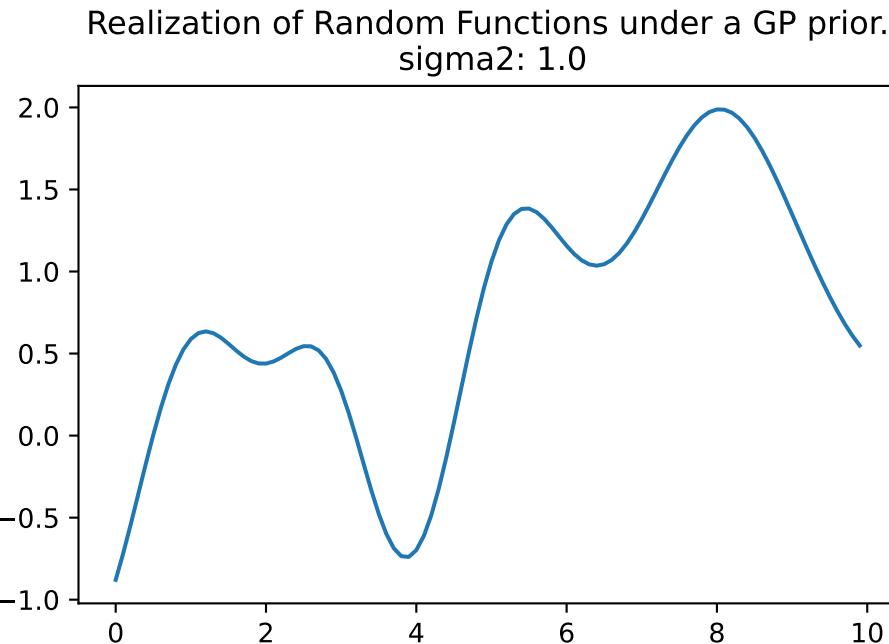


Figure 6.6: Realization of one random function under a GP prior.  $\sigma^2: 1.0$

```

rng = np.random.default_rng(seed=12345)
Y = rng.multivariate_normal(zeros(Sigma.shape[0]), Sigma, size = 3, check_valid="raise")
plt.plot(X, Y.T)
plt.title("Realization of Three Random Functions under a GP prior.\n sigma2: {}".format(sigma2))
plt.show()

```

Realization of Three Random Functions under a GP prior.  
 $\sigma^2: 1.0$

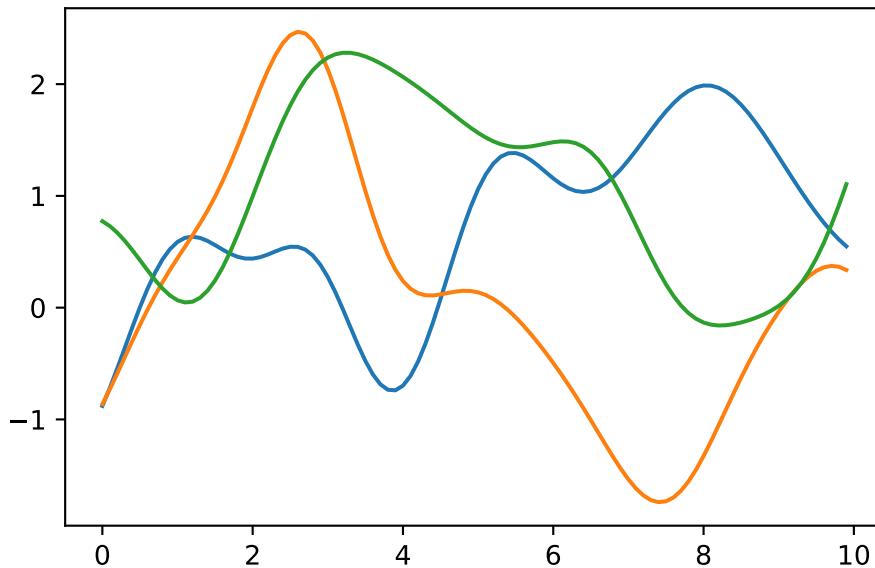


Figure 6.7: Realization of three random functions under a GP prior.  $\sigma^2: 1.0$

### 6.7.5 Properties of the 1d Example

#### 6.7.5.1 Several Bumps:

In this analysis, we observe several bumps in the  $x$ -range of  $[0, 10]$ . These bumps in the function occur because shorter distances exhibit high correlation, while longer distances tend to be essentially uncorrelated. This leads to variations in the function's behavior:

- When  $x$  and  $x'$  are one  $\sigma$  unit apart, the correlation is  $\exp(-\sigma^2/(2\sigma^2)) = \exp(-1/2) \approx 0.61$ , i.e., a relative high correlation.
- $2\sigma$  apart means correlation  $\exp(-2^2/2) \approx 0.14$ , i.e., only small correlation.
- $4\sigma$  apart means correlation  $\exp(-4^2/2) \approx 0.0003$ , i.e., nearly no correlation—variables are considered independent for almost all practical application.

### 6.7.5.2 Smoothness:

The function plotted in Figure 6.6 represents only a finite realization, which means that we have data for a limited number of pairs, specifically 100 points. These points appear smooth in a tactile sense because they are closely spaced, and the plot function connects the dots with lines to create the appearance of smoothness. The complete surface, which can be conceptually extended to an infinite realization over a compact domain, is exceptionally smooth in a calculus sense due to the covariance function's property of being infinitely differentiable.

### 6.7.5.3 Scale of Two:

Regarding the scale of the  $Y$  values, they have a range of approximately  $[-2, 2]$ , with a 95% probability of falling within this range. In standard statistical terms, 95% of the data points typically fall within two standard deviations of the mean, which is a common measure of the spread or range of data.

```
import numpy as np
from numpy import array, zeros, power, ones, exp, multiply, eye, linspace, mat, spacing, sqrt
from numpy.random import multivariate_normal

def build_Sigma(X, sigma2):
    n = X.shape[0]
    k = X.shape[1]
    D = zeros((k, n, n))
    for l in range(k):
        for i in range(n):
            for j in range(i, n):
                D[l, i, j] = 1/(2*sigma2[l])*(X[i,l] - X[j,l])**2
    D = sum(D)
    D = D + D.T
    return exp(-D)

def plot_mvnb( a=0, b=10, sigma2=1.0, size=1, n=100, show=True):
    X = np.linspace(a, b, n, endpoint=False).reshape(-1,1)
    sigma2 = np.array([sigma2])
    Sigma = build_Sigma(X, sigma2)
    rng = np.random.default_rng(seed=12345)
    Y = rng.multivariate_normal(zeros(Sigma.shape[0]), Sigma, size = size, check_valid="raise")
    plt.plot(X, Y.T)
    plt.title("Realization of Random Functions under a GP prior.\n sigma2: {}".format(sigma2))
    if show:
        plt.show()
```

```
plot_mvnr(a=0, b=10, sigma2=10.0, size=3, n=250)
```

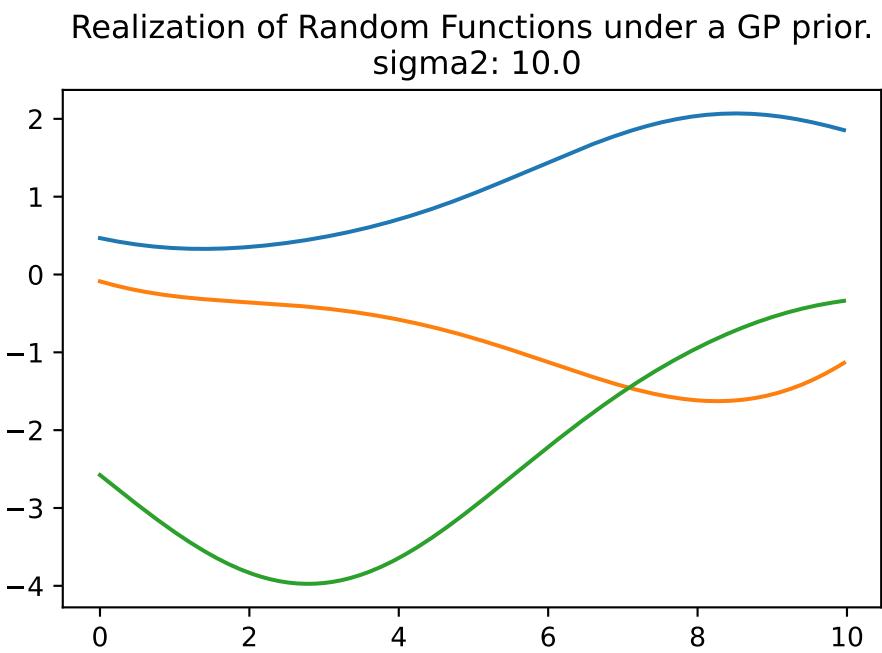


Figure 6.8: Realization of Random Functions under a GP prior. sigma2: 10

```
plot_mvnr(a=0, b=10, sigma2=0.1, size=3, n=250)
```

Realization of Random Functions under a GP prior.  
sigma2: 0.1

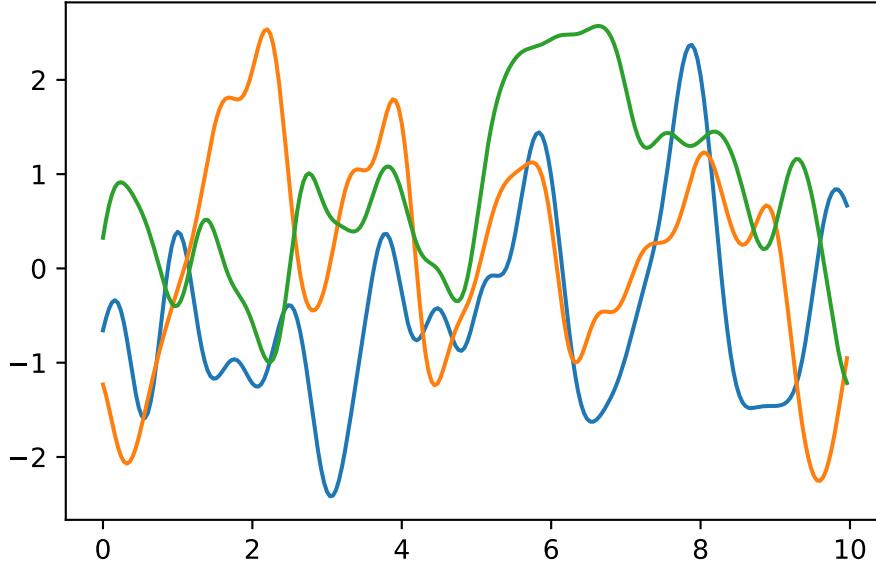


Figure 6.9: Realization of Random Functions under a GP prior. sigma2: 0.1

## 6.8 Kriging: Modeling Basics

### 6.8.1 The Kriging Idea in a Nutshell

We consider observed data of an unknown function  $f$  at  $n$  points  $x_1, \dots, x_n$ , see Figure 6.10. These measurements are considered as realizations of MVN random variables  $Y_1, \dots, Y_n$  with mean  $\mu$  and covariance matrix  $\Sigma_n$  as shown in Figure 6.7, Figure 6.8 or Figure 6.9. In Kriging, a more general covariance matrix (or equivalently, a correlation matrix  $\Psi$ ) is used, see Equation 6.6. Using a maximum likelihood approach, we can estimate the unknown parameters  $\mu$  and  $\Sigma_n$  from the data so that the likelihood function is maximized.

### 6.8.2 The Kriging Basis Function

$k$ -dimensional basis functions of the form

$$\psi(\vec{x}^{(i)}, \vec{x}^{(j)}) = \exp\left(-\sum_{l=1}^k \theta_l |x_l^{(i)} - x_l^{(j)}|^{p_l}\right) \quad (6.6)$$

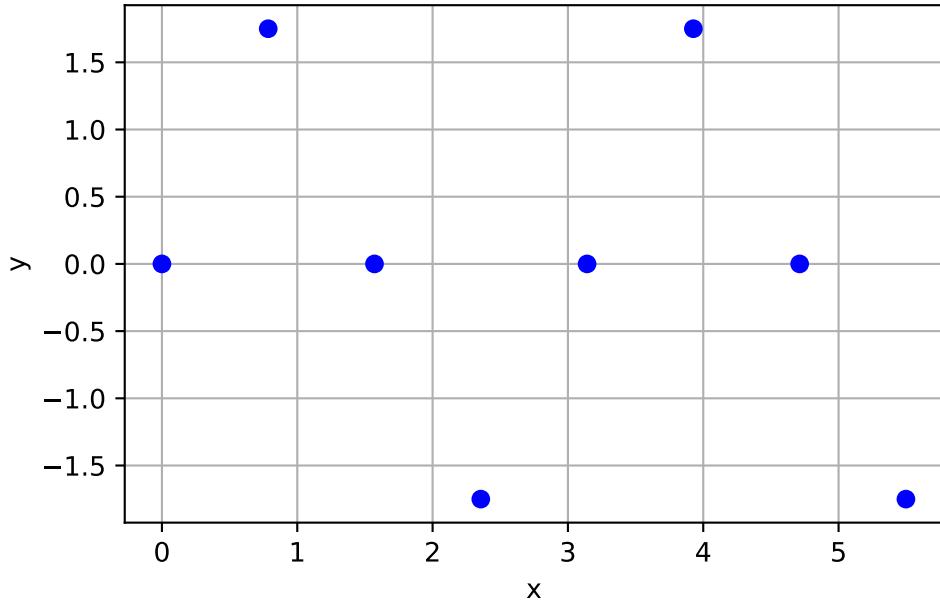


Figure 6.10: Eight measurements of an unknown function

are used in a method known as Kriging. Note,  $\vec{x}^{(i)}$  denotes the  $k$ -dim vector  $\vec{x}^{(i)} = (x_1^{(i)}, \dots, x_k^{(i)})^T$ .

The Kriging basis function is related to the 1-dim Gaussian basis function (Equation 6.5), which is defined as

$$\Sigma(\vec{x}^{(i)}, \vec{x}^{(j)}) = \exp\{-||\vec{x}^{(i)} - \vec{x}^{(j)}||^2/(2\sigma^2)\}. \quad (6.7)$$

There are some differences between Gaussian basis functions and Kriging basis functions:

- Where the Gaussian basis function has  $1/(2\sigma^2)$ , the Kriging basis has a vector  $\theta = [\theta_1, \theta_2, \dots, \theta_k]^T$ .
- The  $\theta$  vector allows the width of the basis function to vary from dimension to dimension.
- In the Gaussian basis function, the exponent is fixed at 2, Kriging allows this exponent  $p_l$  to vary (typically from 1 to 2).

### 6.8.3 The Correlation Coefficient

In a bivariate normal distribution, the covariance matrix and the correlation coefficient are closely related. The covariance matrix  $\Sigma$  for a bivariate normal distribution is a  $2 \times 2$  matrix that looks like this:

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix},$$

where  $\sigma_1^2$  and  $\sigma_2^2$  are the variances of  $X_1$  and  $X_2$ , and  $\sigma_{12} = \sigma_{21}$  is the covariance between  $X_1$  and  $X_2$ .

The correlation coefficient, often denoted as  $\rho$ , is a normalized measure of the linear relationship between two variables. It is calculated from the covariance and the standard deviations  $\sigma_1$  and  $\sigma_2$  (or the square roots of the variances) of  $X_1$  and  $X_2$  as follows:

$$\rho = \sigma_{12}/(\sqrt{\sigma_1^2} \times \sqrt{\sigma_2^2}) = \sigma_{12}/(\sigma_1 \times \sigma_2).$$

So we can express the correlation coefficient  $\rho$  in terms of the elements of the covariance matrix  $\Sigma$ . It can be interpreted as follows: The correlation coefficient ranges from -1 to 1. A value of 1 means that  $X_1$  and  $X_2$  are perfectly positively correlated, a value of -1 means they are perfectly negatively correlated, and a value of 0 means they are uncorrelated. This gives the same information as the covariance, but on a standardized scale that does not depend on the units of  $X_1$  and  $X_2$ .

#### 6.8.4 Covariance Matrix and Correlation Matrix

##### Covariance and Correlation (taken from @Forr08a)

Covariance is a measure of the correlation between two or more sets of random variables.

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$$

From the covariance, we can derive the correlation

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sigma_X\sigma_Y}. \quad (6.8)$$

For a vector of random variables

$$Y = ((Y^{(1)}, \dots, Y^{(n)}))^T$$

the covariance matrix is a matrix of covariances between the random variables

$$\Sigma = \text{Cov}(Y, Y) = \begin{pmatrix} \text{Cov}(Y^{(1)}, Y^{(1)}) & \dots & \text{Cov}(Y^{(1)}, Y^{(n)}) \\ \vdots & \ddots & \vdots \\ \text{Cov}(Y^{(n)}, Y^{(1)}) & \dots & \text{Cov}(Y^{(n)}, Y^{(n)}) \end{pmatrix},$$

and from Equation 6.8

$$\text{Cov}(Y, Y) = \sigma_Y^2 \text{Cor}(Y, Y).$$

You can compute the correlation matrix  $\Psi$  from a covariance matrix  $\Sigma$  in Python using the numpy library. The correlation matrix is computed by dividing each element of the covariance matrix by the product of the standard deviations of the corresponding variables.

The function `covariance_to_correlation` first computes the standard deviations of the variables with `np.sqrt(np.diag(cov))`. It then computes the correlation matrix by dividing each element of the covariance matrix by the product of the standard deviations of the corresponding variables with `cov / np.outer(std_devs, std_devs)`.

```
import numpy as np

def covariance_to_correlation(cov):
    # Compute standard deviations
    std_devs = np.sqrt(np.diag(cov))

    # Compute correlation matrix
    corr = cov / np.outer(std_devs, std_devs)

    return corr

cov = np.array([[9, -4], [-4, 9]])
print(covariance_to_correlation(cov))
```

```
[[ 1.          -0.44444444]
 [-0.44444444  1.          ]]
```

### 6.8.5 The Kriging Model

Consider sample data  $\vec{X}$  and  $\vec{y}$  from  $n$  locations that are available in matrix form:  $\vec{X}$  is a  $(n \times k)$  matrix, where  $k$  denotes the problem dimension and  $\vec{y}$  is a  $(n \times 1)$  vector.

The observed responses  $\vec{y}$  are considered as if they are from a stochastic process, which will be denoted as

$$\begin{pmatrix} \vec{Y}(\vec{x}^{(1)}) \\ \vdots \\ \vec{Y}(\vec{x}^{(n)}) \end{pmatrix}.$$

The set of random vectors (also referred to as a *random field*) has a mean of  $\vec{\mu}$ , which is a  $(n \times 1)$  vector.

### 6.8.6 Correlations

The random vectors are correlated with each other using the basis function expression from Equation 6.6:

$$\text{cor}(\vec{Y}(\vec{x}^{(i)}), \vec{Y}(\vec{x}^{(l)})) = \exp \left\{ - \sum_{j=1}^k \theta_j |x_j^{(i)} - x_j^{(l)}|^{p_j} \right\}.$$

The  $(n \times n)$  correlation matrix of the observed sample data is

$$\vec{\Psi} = \begin{pmatrix} \text{cor}(\vec{Y}(\vec{x}^{(i)}), \vec{Y}(\vec{x}^{(l)})) & \dots & \text{cor}(\vec{Y}(\vec{x}^{(i)}), \vec{Y}(\vec{x}^{(l)})) \\ \vdots & \ddots & \vdots \\ \text{cor}(\vec{Y}(\vec{x}^{(i)}), \vec{Y}(\vec{x}^{(l)})) & \dots & \text{cor}(\vec{Y}(\vec{x}^{(i)}), \vec{Y}(\vec{x}^{(l)})) \end{pmatrix}.$$

Note: correlations depend on the absolute distances between sample points  $|x_j^{(n)} - x_j^{(l)}|$  and the parameters  $p_j$  and  $\theta_j$ .

Correlation is intuitive, because when two points move close together, then  $|x_l^{(i)} - x_l| \rightarrow 0$  and  $\exp(-|x_l^{(i)} - x_l|) \rightarrow 1$ , points show very close correlation and  $Y(x_l^{(i)}) = Y(x_l)$ .

$\theta$  can be seen as a width parameter:

- low  $\theta_j$  means that all points will have a high correlation, with  $Y(x_j)$  being similar across the sample.
- high  $\theta_j$  means that there is a significant difference between the  $Y(x_j)$ 's.
- $\theta_j$  is a measure of how active the function we are approximating is.
- High  $\theta_j$  indicate important parameters, see Figure 6.11.

```
visualize_inverse_exp_squared_distance(5, 0, theta_values=[0.5, 1, 2.0])
```

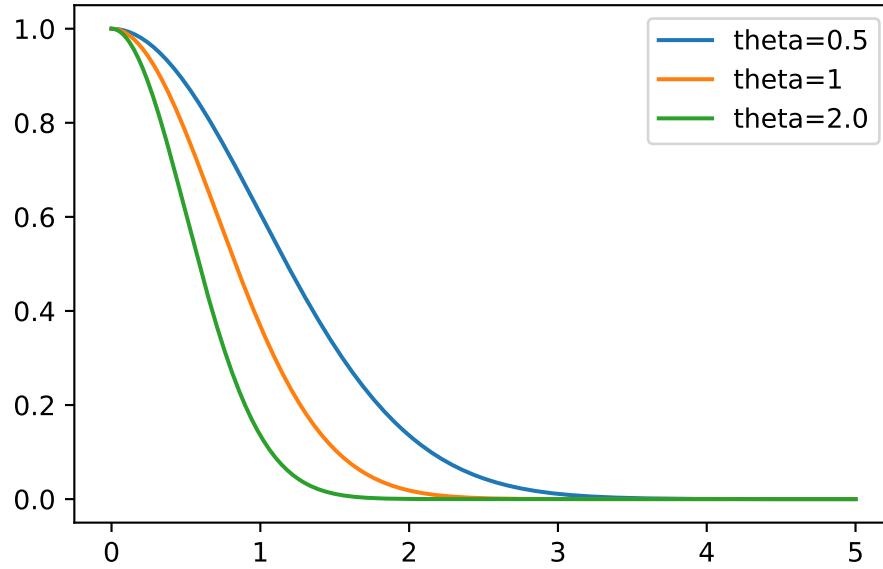


Figure 6.11: Theta set to 1/2, 1, and 2

**i** Example: The Correlation Matrix (Detailed Computation)

Let  $n = 4$  and  $k = 3$ . The sample plan is represented by the following matrix  $X$ :

$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \\ x_{41} & x_{42} & x_{43} \end{pmatrix}$$

To compute the elements of the matrix  $\Psi$ , the following  $k$  (one for each of the  $k$  dimensions)  $(n, n)$ -matrices have to be computed:

$$D_1 = \begin{pmatrix} x_{11} - x_{11} & x_{11} - x_{21} & x_{11} - x_{31} & x_{11} - x_{41} \\ x_{21} - x_{11} & x_{21} - x_{21} & x_{21} - x_{31} & x_{21} - x_{41} \\ x_{31} - x_{11} & x_{31} - x_{21} & x_{31} - x_{31} & x_{31} - x_{41} \\ x_{41} - x_{11} & x_{41} - x_{21} & x_{41} - x_{31} & x_{41} - x_{41} \end{pmatrix}$$

$$D_2 = \begin{pmatrix} x_{12} - x_{12} & x_{12} - x_{22} & x_{12} - x_{32} & x_{12} - x_{42} \\ x_{22} - x_{12} & x_{22} - x_{22} & x_{22} - x_{32} & x_{22} - x_{42} \\ x_{32} - x_{12} & x_{32} - x_{22} & x_{32} - x_{32} & x_{32} - x_{42} \\ x_{42} - x_{12} & x_{42} - x_{22} & x_{42} - x_{32} & x_{42} - x_{42} \end{pmatrix}$$

$$D_3 = \begin{pmatrix} x_{13} - x_{13} & x_{13} - x_{23} & x_{13} - x_{33} & x_{13} - x_{43} \\ x_{23} - x_{13} & x_{23} - x_{23} & x_{23} - x_{33} & x_{23} - x_{43} \\ x_{33} - x_{13} & x_{33} - x_{23} & x_{33} - x_{33} & x_{33} - x_{43} \\ x_{43} - x_{13} & x_{43} - x_{23} & x_{43} - x_{33} & x_{43} - x_{43} \end{pmatrix}$$

Since the matrices are symmetric and the main diagonals are zero, it is sufficient to compute the following matrices:

$$D_1 = \begin{pmatrix} 0 & x_{11} - x_{21} & x_{11} - x_{31} & x_{11} - x_{41} \\ 0 & 0 & x_{21} - x_{31} & x_{21} - x_{41} \\ 0 & 0 & 0 & x_{31} - x_{41} \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$D_2 = \begin{pmatrix} 0 & x_{12} - x_{22} & x_{12} - x_{32} & x_{12} - x_{42} \\ 0 & 0 & x_{22} - x_{32} & x_{22} - x_{42} \\ 0 & 0 & 0 & x_{32} - x_{42} \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$D_3 = \begin{pmatrix} 0 & x_{13} - x_{23} & x_{13} - x_{33} & x_{13} - x_{43} \\ 0 & 0 & x_{23} - x_{33} & x_{23} - x_{43} \\ 0 & 0 & 0 & x_{33} - x_{43} \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

We will consider  $p_l = 2$ . The differences will be squared and multiplied by  $\theta_i$ , i.e.:

$$D_1 = \theta_1 \begin{pmatrix} 0 & (x_{11} - x_{21})^2 & (x_{11} - x_{31})^2 & (x_{11} - x_{41})^2 \\ 0 & 0 & (x_{21} - x_{31})^2 & (x_{21} - x_{41})^2 \\ 0 & 0 & 0 & (x_{31} - x_{41})^2 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$D_2 = \theta_2 \begin{pmatrix} 0 & (x_{12} - x_{22})^2 & (x_{12} - x_{32})^2 & (x_{12} - x_{42})^2 \\ 0 & 0 & (x_{22} - x_{32})^2 & (x_{22} - x_{42})^2 \\ 0 & 0 & 0 & (x_{32} - x_{42})^2 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$D_3 = \theta_3 \begin{pmatrix} 0 & (x_{13} - x_{23})^2 & (x_{13} - x_{33})^2 & (x_{13} - x_{43})^2 \\ 0 & 0 & (x_{23} - x_{33})^2 & (x_{23} - x_{43})^2 \\ 0 & 0 & 0 & (x_{33} - x_{43})^2 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

The sum of the three matrices  $D = D_1 + D_2 + D_3$  will be calculated next:

$$\begin{pmatrix} 0 & \theta_1(x_{11} - x_{21})^2 + \theta_2(x_{12} - x_{22})^2 + \theta_3(x_{13} - x_{23})^2 & \theta_1(x_{11} - x_{31})^2 + \theta_2(x_{12} - x_{32})^2 + \theta_3(x_{13} - x_{33})^2 & \theta_1(x_{11} - x_{41})^2 + \theta_2(x_{12} - x_{42})^2 + \theta_3(x_{13} - x_{43})^2 \\ 0 & 0 & \theta_1(x_{21} - x_{31})^2 + \theta_2(x_{22} - x_{32})^2 + \theta_3(x_{23} - x_{33})^2 & \theta_1(x_{21} - x_{41})^2 + \theta_2(x_{22} - x_{42})^2 + \theta_3(x_{23} - x_{43})^2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

Finally,

$$\Psi = \exp(-D)$$

is computed.

Next, we will demonstrate how this computation can be implemented in Python.

```
from numpy import (array, zeros, power, ones, exp, multiply,
                   eye, linspace, mat, spacing, sqrt, arange,
                   append, ravel)
from numpy.linalg import cholesky, solve
theta = np.array([1,2,3])
X = np.array([[1,0,0], [0,1,0], [100, 100, 100], [101, 100, 100]])
X

array([[ 1,    0,    0],
       [ 0,    1,    0],
       [100, 100, 100],
       [101, 100, 100]])

def build_Psi(X, theta):
    n = X.shape[0]
    k = X.shape[1]
    D = zeros((k, n, n))
    for l in range(k):
        for i in range(n):
            for j in range(i, n):
                D[l, i, j] = theta[l]*(X[i,l] - X[j,l])**2
    D = sum(D)
    D = D + D.T
    return exp(-D)

Psi = build_Psi(X, theta)
Psi

array([[1.          , 0.04978707, 0.          , 0.          ],
       [0.04978707, 1.          , 0.          , 0.          ],
       [0.          , 0.          , 1.          , 0.36787944],
       [0.          , 0.          , 0.36787944, 1.        ]])
```

### Example: The Correlation Matrix (Using Existing Functions)

The same result as computed in the previous example can be obtained with existing python functions, e.g., from the package `scipy`.

```
from scipy.spatial.distance import squareform
from scipy.spatial.distance import pdist

def build_Psi(X, theta, eps=sqrt(spacing(1))):
    return exp(- squareform(pdist(X,
                                    metric='sqeuclidean',
                                    out=None,
                                    w=theta))) + multiply(eye(X.shape[0]),
                                                          eps)

Psi = build_Psi(X, theta, eps=.0)
Psi

array([[1.          , 0.04978707, 0.          , 0.          ],
       [0.04978707, 1.          , 0.          , 0.          ],
       [0.          , 0.          , 1.          , 0.36787944],
       [0.          , 0.          , 0.36787944, 1.        ]])
```

#### 6.8.7 The Condition Number

A small value, `eps`, can be passed to the function `build_Psi` to improve the condition number. For example, `eps=sqrt(spacing(1))` can be used. The numpy function `spacing()` returns the distance between a number and its nearest adjacent number.

The condition number of a matrix is a measure of its sensitivity to small changes in its elements. It is used to estimate how much the output of a function will change if the input is slightly altered.

A matrix with a low condition number is well-conditioned, which means its behavior is relatively stable, while a matrix with a high condition number is ill-conditioned, meaning its behavior is unstable with respect to numerical precision.

```
import numpy as np

# Define a well-conditioned matrix (low condition number)
A = np.array([[1, 0.1], [0.1, 1]])
print("Condition number of A: ", np.linalg.cond(A))
```

```
# Define an ill-conditioned matrix (high condition number)
B = np.array([[1, 0.9999999], [0.9999999, 1]])
print("Condition number of B: ", np.linalg.cond(B))
```

```
Condition number of A: 1.222222222222225
Condition number of B: 200000000.53159264
```

```
np.linalg.cond(Psi)
```

```
2.163953413738652
```

### 6.8.8 MLE to estimate $\theta$ and $p$

We know what the correlations mean, but how do we estimate the values of  $\theta_j$  and where does our observed data  $y$  come in? To estimate the values of  $\vec{\theta}$  and  $\vec{p}$ , they are chosen to maximize the likelihood of  $\vec{y}$ , which can be expressed in terms of the sample data

$$L(\vec{Y}(\vec{x}^{(1)}), \dots, \vec{Y}(\vec{x}^{(n)}) | \mu, \sigma) = \frac{1}{(2\pi\sigma)^{n/2} |\vec{\Psi}|^{1/2}} \exp \left\{ \frac{-(\vec{y} - \vec{1}\mu)^T \vec{\Psi}^{-1} (\vec{y} - \vec{1}\mu)}{2\sigma^2} \right\},$$

and formulated as the log-likelihood:

$$\ln(L) = -\frac{n}{2} \ln(2\pi\sigma) - \frac{1}{2} \ln |\vec{\Psi}| \frac{-(\vec{y} - \vec{1}\mu)^T \vec{\Psi}^{-1} (\vec{y} - \vec{1}\mu)}{2\sigma^2}.$$

Optimization of the log-likelihood by taking derivatives with respect to  $\mu$  and  $\sigma$  results in

$$\hat{\mu} = \frac{\vec{1}^T \vec{\Psi}^{-1} \vec{y}^T}{\vec{1}^T \vec{\Psi}^{-1} \vec{1}^T}$$

and

$$\hat{\sigma} = \frac{(\vec{y} - \vec{1}\mu)^T \vec{\Psi}^{-1} (\vec{y} - \vec{1}\mu)}{n}.$$

Combining the equations leads to the concentrated log-likelihood:

$$\ln(L) = -\frac{n}{2} \ln(\hat{\sigma}) - \frac{1}{2} \ln |\vec{\Psi}|. \quad (6.9)$$

**i** Note: The Concentrated Log-Likelihood

- The first term in Equation 6.9 requires information about the measured point (observations)  $y_i$ .
- To maximize  $\ln(L)$ , optimal values of  $\vec{\theta}$  and  $\vec{p}$  are determined numerically, because the equation is not differentiable.

### 6.8.9 Tuning $\theta$ and $p$

Optimizers such as Nelder-Mead, Conjugate Gradient, or Simulated Annealing can be used to determine optimal values for  $\theta$  and  $p$ . After the optimization, the correlation matrix  $\Psi$  is build with the optimized  $\theta$  and  $p$  values. This is best (most likely) Kriging model for the given data  $y$ .

## 6.9 Kriging Prediction

### 6.9.1 The Augmented Correlation Matrix

We will use the Kriging correlation  $\Psi$  to predict new values based on the observed data. The matrix algebra involved for calculating the likelihood is the most computationally intensive part of the Kriging process. Care must be taken that the computer code is as efficient as possible.

Basic elements of the Kriging based surrogate optimization such as interpolation, expected improvement, and regression are presented. The presentation follows the approach described in Forrester, Sóbester, and Keane (2008) and Bartz et al. (2022).

Main idea for prediction is that the new  $\vec{Y}(\vec{x})$  should be consistent with the old sample data  $X$ . For a new prediction  $\hat{y}$  at  $\vec{x}$ , the value of  $\hat{y}$  is chosen so that it maximizes the likelihood of the sample data  $\vec{X}$  and the prediction, given the (optimized) correlation parameter  $\vec{\theta}$  and  $\vec{p}$  from above. The observed data  $\vec{y}$  is augmented with the new prediction  $\hat{y}$  which results in the augmented vector  $\vec{\tilde{y}} = (\vec{y}^T, \hat{y})^T$ . A vector of correlations between the observed data and the new prediction is defined as

$$\vec{\psi} = \begin{pmatrix} \text{cor}(\vec{Y}(\vec{x}^{(1)}), \vec{Y}(\vec{x})) \\ \vdots \\ \text{cor}(\vec{Y}(\vec{x}^{(n)}), \vec{Y}(\vec{x})) \end{pmatrix} = \begin{pmatrix} \vec{\psi}^{(1)} \\ \vdots \\ \vec{\psi}^{(n)} \end{pmatrix}.$$

The augmented correlation matrix is constructed as

$$\tilde{\Psi} = \begin{pmatrix} \vec{\Psi} & \vec{\psi} \\ \vec{\psi}^T & 1 \end{pmatrix}.$$

The log-likelihood of the augmented data is

$$\ln(L) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\hat{\sigma}^2) - \frac{1}{2} \ln |\vec{\Psi}| - \frac{(\vec{y} - \vec{1}\hat{\mu})^T \vec{\Psi}^{-1} (\vec{y} - \vec{1}\hat{\mu})}{2\hat{\sigma}^2}.$$

The MLE for  $\hat{y}$  can be calculated as

$$\hat{y}(\vec{x}) = \hat{\mu} + \vec{\psi}^T \vec{\Psi}^{-1} (\vec{y} - \vec{1}\hat{\mu}). \quad (6.10)$$

### 6.9.2 Properties of the Predictor

Equation 6.10 reveals two important properties of the Kriging predictor:

1. Basis functions: The basis function impacts the vector  $\vec{\psi}$ , which contains the  $n$  correlations between the new point  $\vec{x}$  and the observed locations. Values from the  $n$  basis functions are added to a mean base term  $\mu$  with weightings  $\vec{w} = \vec{\Psi}^{(-1)} (\vec{y} - \vec{1}\hat{\mu})$ .
2. Interpolation: The predictions interpolate the sample data. When calculating the prediction at the  $i$ th sample point,  $\vec{x}^{(i)}$ , the  $i$ th column of  $\vec{\Psi}^{-1}$  is  $\vec{\psi}$ , and  $\vec{\psi}^T \vec{\Psi}^{-1}$  is the  $i$ th unit vector. Hence,  $\hat{y}(\vec{x}^{(i)}) = y^{(i)}$ .

## 6.10 Kriging Example: Sinusoid Function

Toy example in 1d where the response is a simple sinusoid measured at eight equally spaced  $x$ -locations in the span of a single period of oscillation.

### 6.10.1 Calculating the Correlation Matrix $\Psi$

The correlation matrix  $\Psi$  is based on the pairwise squared distances between the input locations. Here we will use  $n = 8$  sample locations and  $\theta$  is set to 1.0.

```
n = 8
X = np.linspace(0, 2*np.pi, n, endpoint=False).reshape(-1,1)
# theta should be an array (of one value, for the moment, will be changed later)
theta = np.array([1.0])
Psi = build_Psi(X, theta)
```

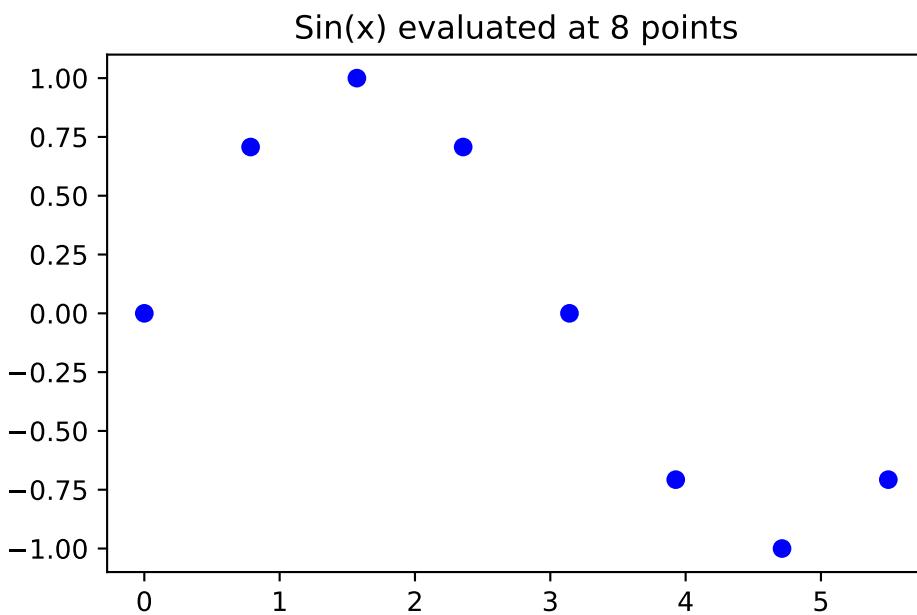
Evaluate at sample points

```

y = np.sin(X)

import matplotlib.pyplot as plt
plt.plot(X, y, "bo")
plt.title(f"Sin(x) evaluated at {n} points")
plt.show()

```



### 6.10.2 Computing the $\psi$ Vector

Distances between testing locations  $x$  and training data locations  $X$ .

```

from scipy.spatial.distance import cdist

def build_psi(X, x, theta, eps=sqrt(spacing(1))):
    n = X.shape[0]
    k = X.shape[1]
    m = x.shape[0]
    psi = zeros((n, m))
    theta = theta * ones(k)
    D = zeros((n, m))
    D = cdist(x.reshape(-1, k),
              X.reshape(-1, k),

```

```

        metric='sqeuclidean',
        out=None,
        w=theta)
print(D.shape)
psi = exp(-D)
# return psi transpose to be consistent with the literature
return(psi.T)

```

### 6.10.3 Predicting at New Locations

We would like to predict at  $m = 100$  new locations in the interval  $[0, 2\pi]$ . The new locations are stored in the variable  $\mathbf{x}$ .

```

m = 100
x = np.linspace(0, 2*np.pi, m, endpoint=False).reshape(-1,1)
psi = build_psi(X, x, theta)

```

(100, 8)

Computation of the predictive equations.

```

U = cholesky(Psi).T
one = np.ones(n).reshape(-1,1)
mu = (one.T.dot(solve(U, solve(U.T, y)))) / one.T.dot(solve(U, solve(U.T, one)))
f = mu * ones(m).reshape(-1,1) + psi.T.dot(solve(U, solve(U.T, y - one * mu)))

```

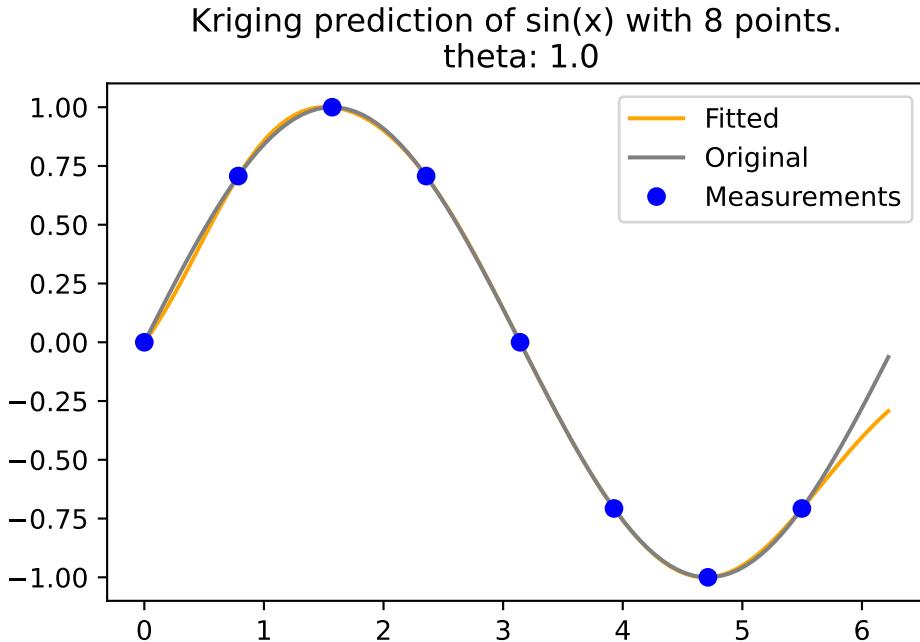
To compute  $f$ , Equation 6.10 is used.

### 6.10.4 Visualization

```

import matplotlib.pyplot as plt
plt.plot(x, f, color = "orange", label="Fitted")
plt.plot(x, np.sin(x), color = "grey", label="Original")
plt.plot(X, y, "bo", label="Measurements")
plt.title("Kriging prediction of sin(x) with {} points.\n theta: {}".format(n, theta[0]))
plt.legend(loc='upper right')
plt.show()

```



## 6.11 Cholesky Example With Two Points

### 6.11.1 Cholesky Decomposition

We consider  $k = 1$  and  $n = 2$  sample points. The sample points are located at  $x_1 = 1$  and  $x_2 = 5$ . The response values are  $y_1 = 2$  and  $y_2 = 10$ . The correlation parameter is  $\theta = 1$  and  $p$  is set to 1. Using Equation 6.6, we can compute the correlation matrix  $\Psi$ :

$$\Psi = \begin{pmatrix} 1 & e^{-1} \\ e^{-1} & 1 \end{pmatrix}.$$

To determine MLE as in Equation 6.10, we need to compute  $\Psi^{-1}$ :

$$\Psi^{-1} = \frac{e}{e^2 - 1} \begin{pmatrix} e & -1 \\ -1 & e \end{pmatrix}.$$

Cholesky-decomposition of  $\Psi$  is recommended to compute  $\Psi^{-1}$ . Cholesky decomposition is a decomposition of a positive definite symmetric matrix into the product of a lower triangular matrix  $L$ , a diagonal matrix  $D$  and the transpose of  $L$ , which is denoted as  $L^T$ . Consider the following example:

$$\begin{aligned}
LDL^T &= \begin{pmatrix} 1 & 0 \\ l_{21} & 1 \end{pmatrix} \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix} \begin{pmatrix} 1 & l_{21} \\ 0 & 1 \end{pmatrix} = \\
&\begin{pmatrix} d_{11} & 0 \\ d_{11}l_{21} & d_{22} \end{pmatrix} \begin{pmatrix} 1 & l_{21} \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} d_{11} & d_{11}l_{21} \\ d_{11}l_{21} & d_{11}l_{21}^2 + d_{22} \end{pmatrix}. \tag{6.11}
\end{aligned}$$

Using Equation 6.11, we can compute the Cholesky decomposition of  $\Psi$ :

1.  $d_{11} = 1$ ,
2.  $l_{21}d_{11} = e^{-1} \Rightarrow l_{21} = e^{-1}$ , and
3.  $d_{11}l_{21}^2 + d_{22} = 1 \Rightarrow d_{22} = 1 - e^{-2}$ .

The Cholesky decomposition of  $\Psi$  is

$$\Psi = \begin{pmatrix} 1 & 0 \\ e^{-1} & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 1 - e^{-2} \end{pmatrix} \begin{pmatrix} 1 & e^{-1} \\ 0 & 1 \end{pmatrix} = LDL^T$$

Some programs use  $U$  instead of  $L$ . The Cholesky decomposition of  $\Psi$  is

$$\Psi = LDL^T = U^T DU.$$

Using

$$\sqrt{D} = \begin{pmatrix} 1 & 0 \\ 0 & \sqrt{1 - e^{-2}} \end{pmatrix},$$

we can write the Cholesky decomposition of  $\Psi$  without a diagonal matrix  $D$  as

$$\Psi = \begin{pmatrix} 1 & 0 \\ e^{-1} & \sqrt{1 - e^{-2}} \end{pmatrix} \begin{pmatrix} 1 & e^{-1} \\ 0 & \sqrt{1 - e^{-2}} \end{pmatrix} = U^T U.$$

### 6.11.2 Computation of the Inverse Matrix

To compute the inverse of a matrix using the Cholesky decomposition, you can follow these steps:

1. Decompose the matrix  $A$  into  $L$  and  $L^T$ , where  $L$  is a lower triangular matrix and  $L^T$  is the transpose of  $L$ .
2. Compute  $L^{-1}$ , the inverse of  $L$ .
3. The inverse of  $A$  is then  $(L^{-1})^T L^{-1}$ .

Please note that this method only applies to symmetric, positive-definite matrices.

The inverse of the matrix  $\Psi$  from above is:

$$\Psi^{-1} = \frac{e}{e^2 - 1} \begin{pmatrix} e & -1 \\ -1 & e \end{pmatrix}.$$

Here's an example of how to compute the inverse of a matrix using Cholesky decomposition in Python:

```
import numpy as np
from scipy.linalg import cholesky, inv
E = np.exp(1)

# Psi is a symmetric, positive-definite matrix
Psi = np.array([[1, 1/E], [1/E, 1]])
L = cholesky(Psi, lower=True)
L_inv = inv(L)
# The inverse of A is (L^-1)^T * L^-1
Psi_inv = np.dot(L_inv.T, L_inv)

print("Psi:\n", Psi)
print("Psi Inverse:\n", Psi_inv)
```

```
Psi:
[[1.          0.36787944]
 [0.36787944 1.          ]]
Psi Inverse:
[[ 1.15651764 -0.42545906]
 [-0.42545906  1.15651764]]
```

## 6.12 Jupyter Notebook

### Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

# 7 Introduction to spotPython

Surrogate model based optimization methods are common approaches in simulation and optimization. SPOT was developed because there is a great need for sound statistical analysis of simulation and optimization algorithms. SPOT includes methods for tuning based on classical regression and analysis of variance techniques. It presents tree-based models such as classification and regression trees and random forests as well as Bayesian optimization (Gaussian process models, also known as Kriging). Combinations of different meta-modeling approaches are possible. SPOT comes with a sophisticated surrogate model based optimization method, that can handle discrete and continuous inputs. Furthermore, any model implemented in `scikit-learn` can be used out-of-the-box as a surrogate in `spotPython`.

SPOT implements key techniques such as exploratory fitness landscape analysis and sensitivity analysis. It can be used to understand the performance of various algorithms, while simultaneously giving insights into their algorithmic behavior.

The `spot` loop consists of the following steps:

1. Init: Build initial design  $X$
2. Evaluate initial design on real objective  $f$ :  $y = f(X)$
3. Build surrogate:  $S = S(X, y)$
4. Optimize on surrogate:  $X_0 = \text{optimize}(S)$
5. Evaluate on real objective:  $y_0 = f(X_0)$
6. Impute (Infill) new points:  $X = X \cup X_0$ ,  $y = y \cup y_0$ .
7. Goto 3.

Central Idea: Evaluation of the surrogate model  $S$  is much cheaper (or / and much faster) than running the real-world experiment  $f$ . We start with a small example.

## 7.1 Example: Spot and the Sphere Function

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.utils.init import fun_control_init, design_control_init
from spotPython.hyperparameters.values import set_control_key_value
```

```
from spotPython.spot import spot
import matplotlib.pyplot as plt
```

### 7.1.1 The Objective Function: Sphere

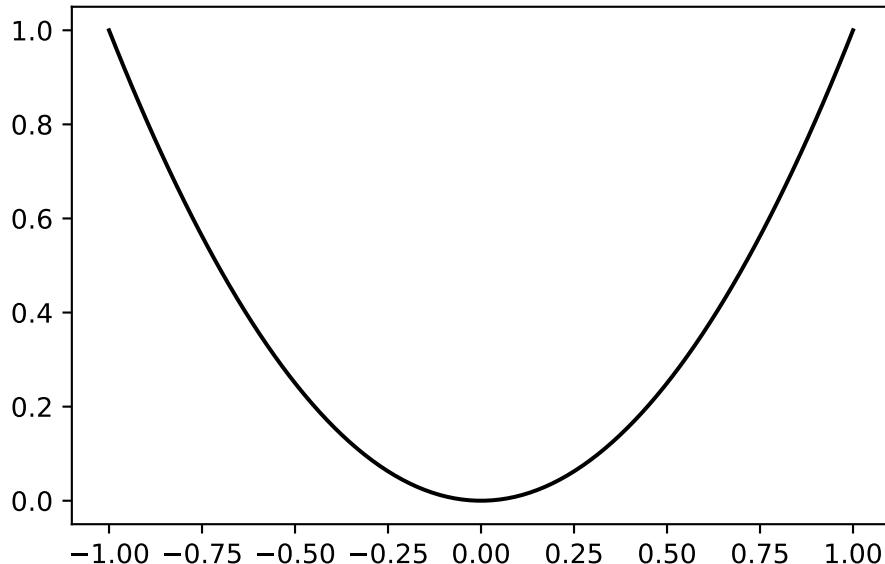
The `spotPython` package provides several classes of objective functions. We will use an analytical objective function, i.e., a function that can be described by a (closed) formula:

$$f(x) = x^2$$

```
fun = analytical().fun_sphere
```

We can apply the function `fun` to input values and plot the result:

```
x = np.linspace(-1,1,100).reshape(-1,1)
y = fun(x)
plt.figure()
plt.plot(x, y, "k")
plt.show()
```



### 7.1.2 The Spot Method as an Optimization Algorithm Using a Surrogate Model

We initialize the `fun_control` dictionary. The `fun_control` dictionary contains the parameters for the objective function. The `fun_control` dictionary is passed to the `Spot` method.

```
fun_control=fun_control_init(lower = np.array([-1]),
                             upper = np.array([1]))
spot_0 = spot.Spot(fun=fun,
                   fun_control=fun_control)
spot_0.run()
```

```
spotPython tuning: 1.1986325668379847e-08 [#####---] 73.33%
spotPython tuning: 1.1986325668379847e-08 [#####---] 80.00%
spotPython tuning: 1.1986325668379847e-08 [#####---] 86.67%
spotPython tuning: 1.1385311249270152e-08 [#####---] 93.33%
spotPython tuning: 5.189351944607845e-10 [#####---] 100.00% Done...
```

```
<spotPython.spot.spot.Spot at 0x2c3ed7b50>
```

Print the results:

```
spot_0.print_results()
```

```
min y: 5.189351944607845e-10
x0: 2.2780149131662516e-05

[['x0', 2.2780149131662516e-05]]
```

Plot the search progress:

```
spot_0.plot_progress(log_y=True)
```

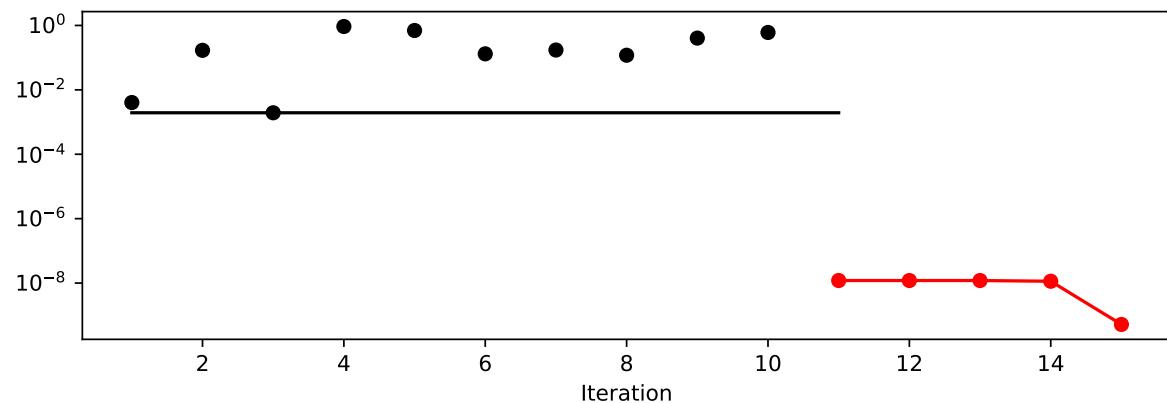


Figure 7.1: Visualization of the search progress of the `Spot` method. The black elements (points and line) represent the initial design, before the surrogate is build. The red elements represent the search on the surrogate.

```
spot_0.plot_model()
```

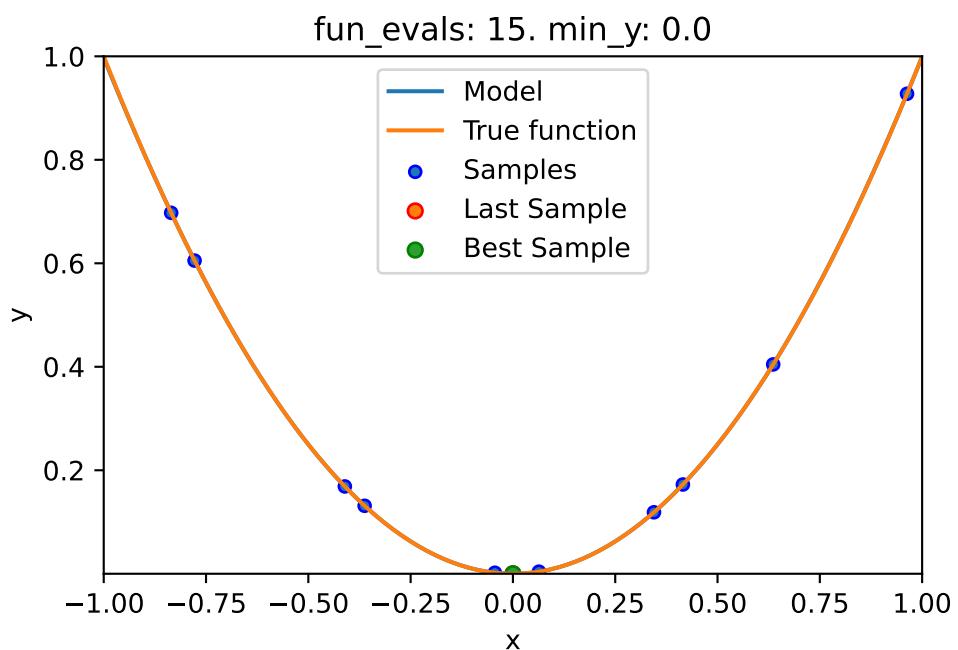


Figure 7.2: Visualization of the model and the underlying objective function values.

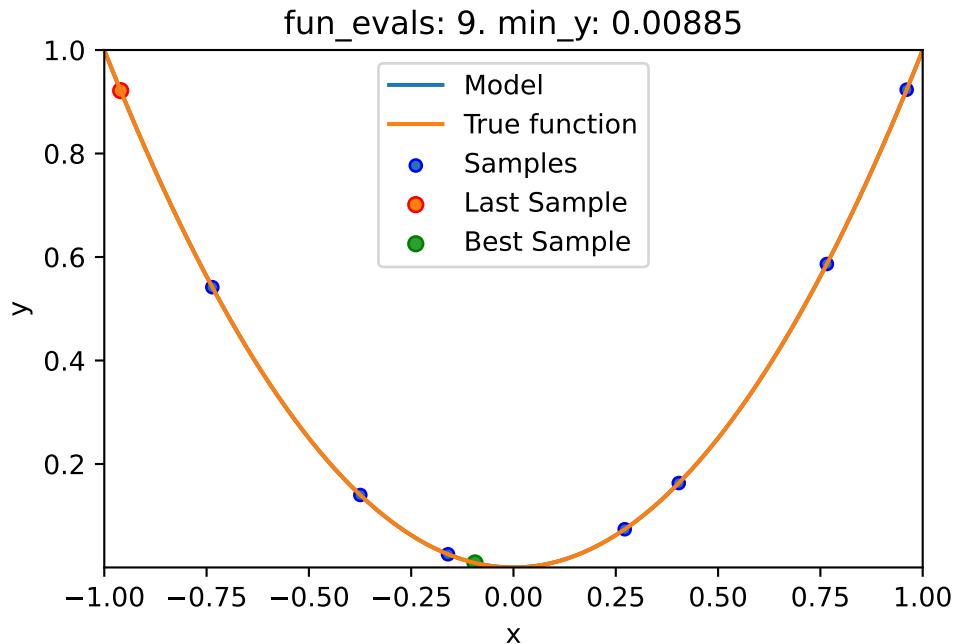
## 7.2 Spot Parameters: `fun_evals`, `init_size` and `show_models`

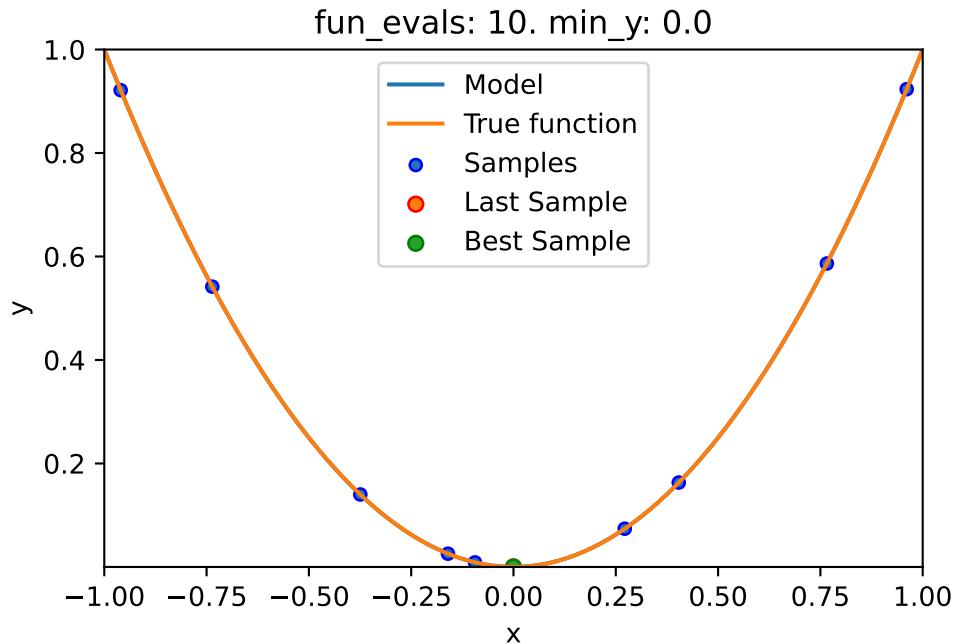
We will modify three parameters:

1. The number of function evaluations (`fun_evals`) will be set to 10 (instead of 15, which is the default value) in the `fun_control` dictionary.
2. The parameter `show_models`, which visualizes the search process for each single iteration for 1-dim functions, in the `fun_control` dictionary.
3. The size of the initial design (`init_size`) in the `design_control` dictionary.

The full list of the Spot parameters is shown in code reference on GitHub, see [Spot](#).

```
fun_control=fun_control_init(lower = np.array([-1]),
                             upper = np.array([1]),
                             fun_evals = 10,
                             show_models = True)
design_control = design_control_init(init_size=9)
spot_1 = spot.Spot(fun=fun,
                    fun_control=fun_control,
                    design_control=design_control)
spot_1.run()
```





```
spotPython tuning: 1.1624309649016718e-08 [#####] 100.00% Done...
```

### 7.3 Print the Results

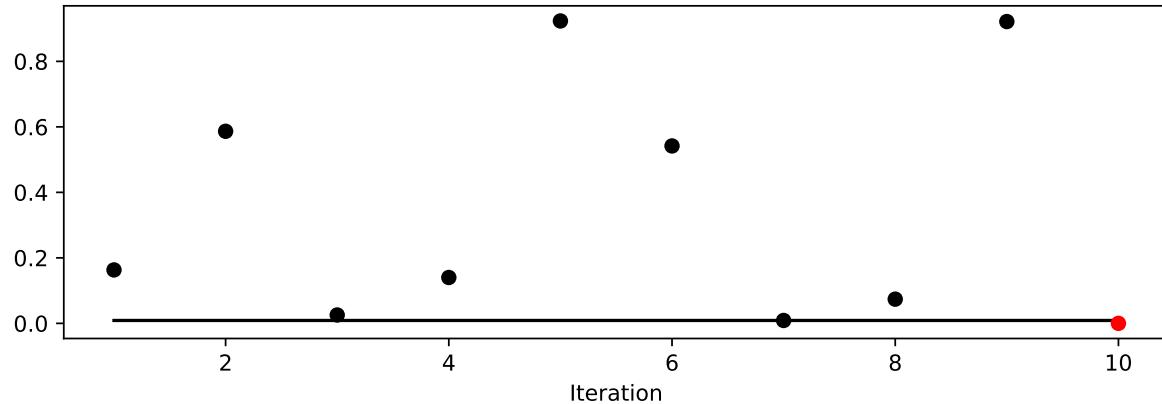
```
spot_1.print_results()
```

```
min y: 1.1624309649016718e-08
x0: -0.00010781609179068177
```

```
[['x0', -0.00010781609179068177]]
```

### 7.4 Show the Progress

```
spot_1.plot_progress()
```



## 7.5 Visualizing the Optimization and Hyperparameter Tuning Process with TensorBoard

`spotPython` supports the visualization of the hyperparameter tuning process with TensorBoard. The following example shows how to use TensorBoard with `spotPython`.

First, we define an “PREFIX” to identify the hyperparameter tuning process. The PREFIX is used to create a directory for the TensorBoard files.

```
fun_control = fun_control_init(
    PREFIX = "01",
    lower = np.array([-1]),
    upper = np.array([2]))
design_control = design_control_init(init_size=5)
```

```
Created spot_tensorboard_path: runs/spot_logs/01_p040025_2024-01-09_20-11-31 for SummaryWriter
```

Since the `spot_tensorboard_path` is defined, `spotPython` will log the optimization process in the TensorBoard files.

```
spot_tuner = spot.Spot(fun=fun,
                      fun_control=fun_control,
                      design_control=design_control)
spot_tuner.run()
spot_tuner.print_results()
```

```

spotPython tuning: 2.7704798574631272e-05 [#####-----] 40.00%
spotPython tuning: 7.761255945313609e-07 [#####-----] 46.67%
spotPython tuning: 7.729540372746933e-07 [#####-----] 53.33%
spotPython tuning: 3.677408011584826e-07 [#####----] 60.00%
spotPython tuning: 1.1649997893485877e-09 [#####----] 66.67%
spotPython tuning: 1.1649997893485877e-09 [#####---] 73.33%
spotPython tuning: 1.1649997893485877e-09 [#####--] 80.00%
spotPython tuning: 1.1649997893485877e-09 [#####-] 86.67%
spotPython tuning: 1.1649997893485877e-09 [#####-] 93.33%
spotPython tuning: 1.1649997893485877e-09 [#####-] 100.00% Done...

```

```

min y: 1.1649997893485877e-09
x0: -3.4132093245926007e-05

```

```
[['x0', -3.4132093245926007e-05]]
```

Now we can start TensorBoard in the background. The TensorBoard process will read the TensorBoard files and visualize the hyperparameter tuning process. From the terminal, we can start TensorBoard with the following command:

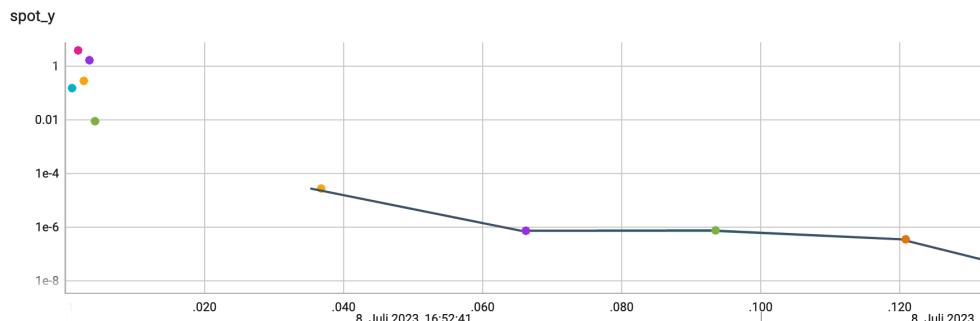
```
tensorboard --logdir=".runs"
```

`logdir` is the directory where the TensorBoard files are stored. In our case, the TensorBoard files are stored in the directory `./runs`.

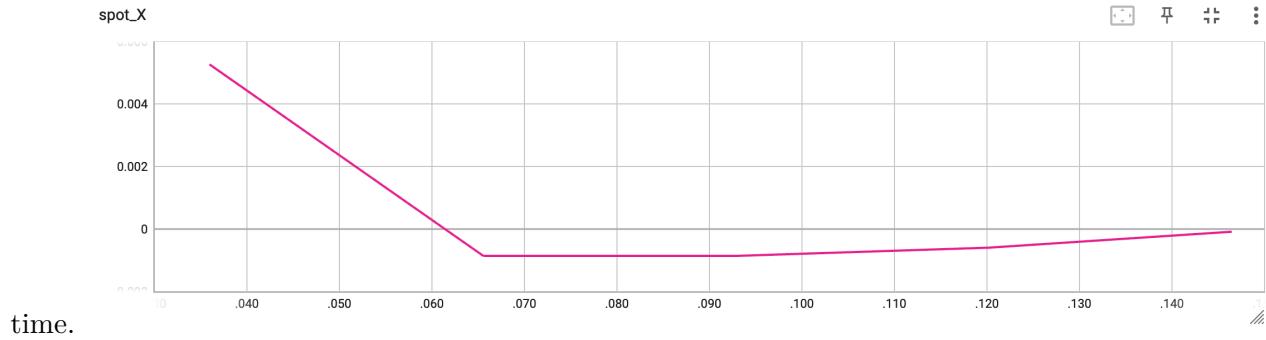
TensorBoard will start a web server on port 6006. We can access the TensorBoard web server with the following URL:

```
http://localhost:6006/
```

The first TensorBoard visualization shows the objective function values plotted against the wall time. The wall time is the time that has passed since the start of the hyperparameter tuning process. The five initial design points are shown in the upper left region of the plot. The line visualizes the optimization process.



The second TensorBoard visualization shows the input values, i.e.,  $x_0$ , plotted against the wall time.



The third TensorBoard plot illustrates how `spotPython` can be used as a microscope for the internal mechanisms of the surrogate-based optimization process. Here, one important parameter, the learning rate  $\theta$  of the Kriging surrogate is plotted against the number of optimization steps.

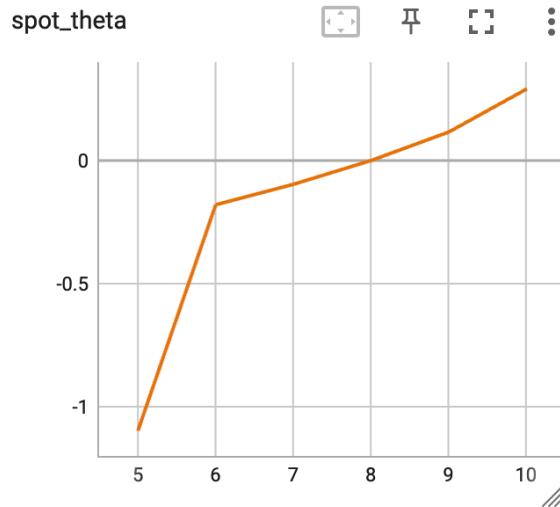


Figure 7.3: TensorBoard visualization of the `spotPython` process.

## 7.6 Jupyter Notebook

**i** Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

# 8 Multi-dimensional Functions

This chapter illustrates how high-dimensional functions can be optimized and analyzed.

## 8.1 Example: Spot and the 3-dim Sphere Function

```
import numpy as np
from spotPython.fun.objectivefunctions import analytical
from spotPython.utils.init import fun_control_init, surrogate_control_init
from spotPython.spot import spot
```

### 8.1.1 The Objective Function: 3-dim Sphere

The `spotPython` package provides several classes of objective functions. We will use an analytical objective function, i.e., a function that can be described by a (closed) formula:

$$f(x) = \sum_i^k x_i^2$$

. It is available as `fun_sphere` in the `analytical` class.

```
fun = analytical().fun_sphere
```

Here we will use problem dimension  $k = 3$ .

- The size of the `lower` bound array determines the problem dimension.
- Here we will use `-1.0 * np.ones(3)`, i.e., a three-dim function.
- We will use three different `theta` values (one for each dimension), i.e., we set `n_theta=3` in the `surrogate_control`.

Similar to the one-dimensional case, which was introduced in Section Section 7.5, we can use TensorBoard to monitor the progress of the optimization. We will use the same code, only the prefix is different:

```

fun_control = fun_control_init(
    PREFIX="01",
    lower = -1.0*np.ones(3),
    upper = np.ones(3),
    var_name=["Pressure", "Temp", "Lambda"],
    show_progress=True)
surrogate_control = surrogate_control_init(n_theta=3)
spot_3 = spot.Spot(fun=fun,
                    fun_control=fun_control,
                    surrogate_control=surrogate_control)
spot_3.run()

```

```

Created spot_tensorboard_path: runs/spot_logs/01_p040025_2024-01-09_20-11-46 for SummaryWriter
spotPython tuning: 0.03443414446852215 [#####---] 73.33%
spotPython tuning: 0.031346790207570745 [#####---] 80.00%
spotPython tuning: 0.0009630779801151435 [#####---] 86.67%
spotPython tuning: 8.57685255248343e-05 [#####---] 93.33%
spotPython tuning: 7.337173529617274e-05 [#####] 100.00% Done...

```

```
<spotPython.spot.spot.Spot at 0x2cd92f8d0>
```

### Note

Now we can start TensorBoard in the background with the following command:

```
tensorboard --logdir=".runs"
```

and can access the TensorBoard web server with the following URL:

```
http://localhost:6006/
```

## 8.1.2 Results

```
spot_3.print_results()
```

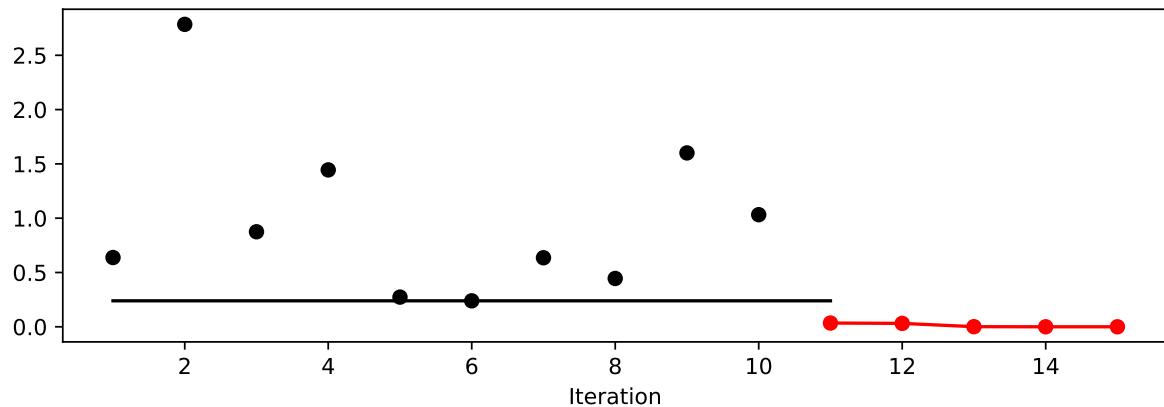
```

min y: 7.337173529617274e-05
Pressure: 0.005530391439119171
Temp: 0.001931974999553984
Lambda: 0.00624931823700702

```

```
[['Pressure', 0.005530391439119171],  
 ['Temp', 0.001931974999553984],  
 ['Lambda', 0.00624931823700702]]
```

```
spot_3.plot_progress()
```



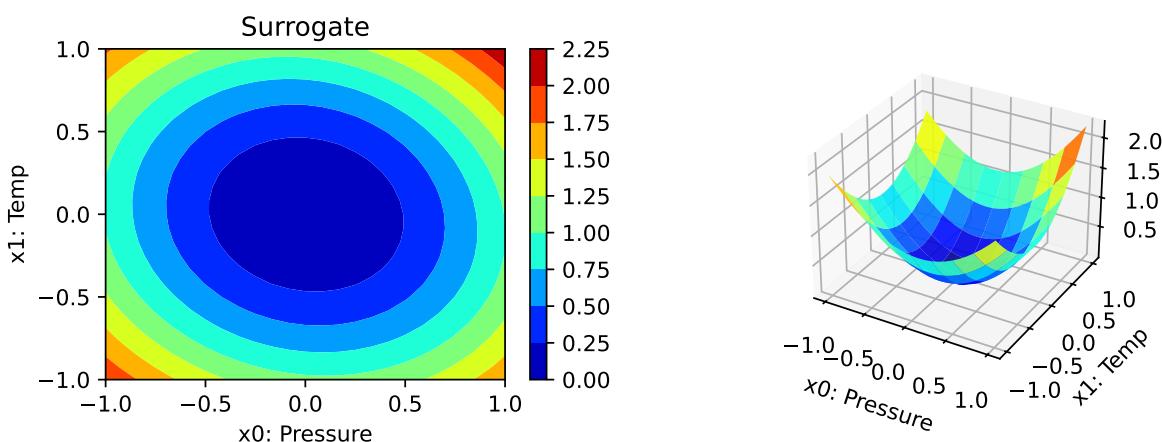
### 8.1.3 A Contour Plot

We can select two dimensions, say  $i = 0$  and  $j = 1$ , and generate a contour plot as follows.

**i** Note:

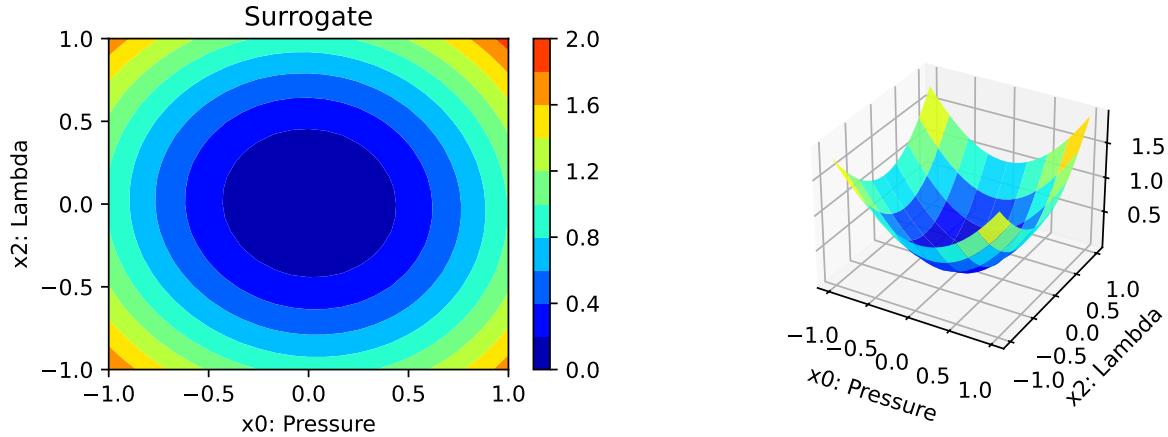
We have specified identical `min_z` and `max_z` values to generate comparable plots.

```
spot_3.plot_contour(i=0, j=1, min_z=0, max_z=2.25)
```



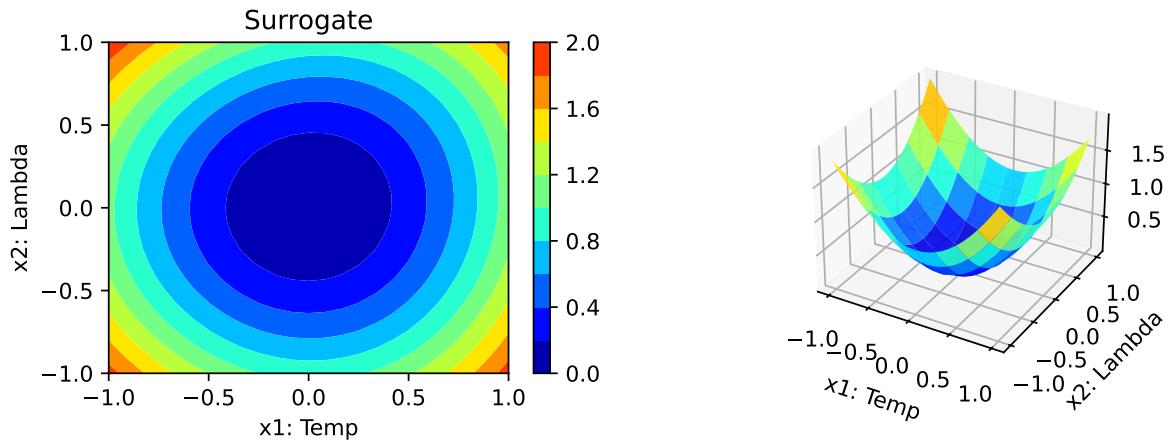
- In a similar manner, we can plot dimension  $i = 0$  and  $j = 2$ :

```
spot_3.plot_contour(i=0, j=2, min_z=0, max_z=2.25)
```



- The final combination is  $i = 1$  and  $j = 2$ :

```
spot_3.plot_contour(i=1, j=2, min_z=0, max_z=2.25)
```



- The three plots look very similar, because the `fun_sphere` is symmetric.
- This can also be seen from the variable importance:

```
spot_3.print_importance()
```

```

Pressure: 95.28112765785582
Temp: 100.0
Lambda: 86.8087715378692

```

```

[['Pressure', 95.28112765785582],
 ['Temp', 100.0],
 ['Lambda', 86.8087715378692]]

```

### 8.1.4 TensorBoard

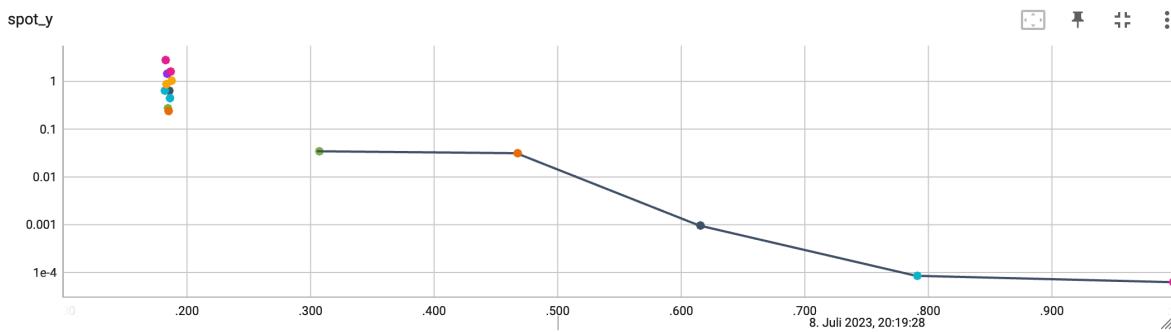
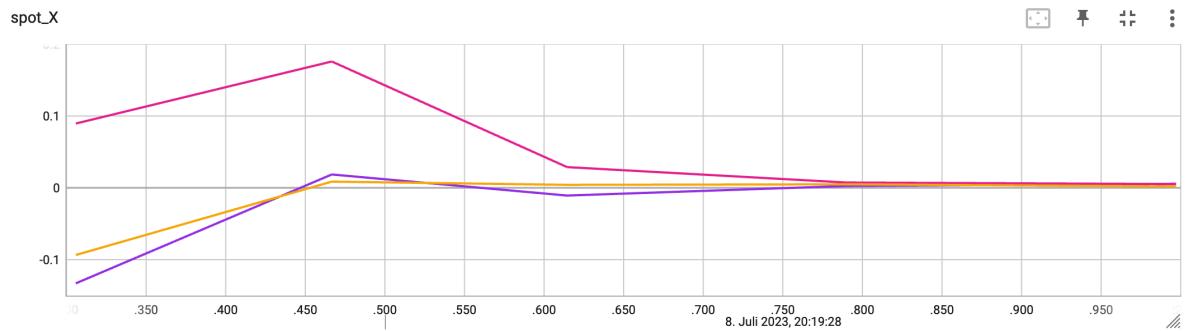


Figure 8.1: TensorBoard visualization of the `spotPython` process. Objective function values plotted against wall time.

The second TensorBoard visualization shows the input values, i.e.,  $x_0, \dots, x_2$ , plotted against



the wall time.

The third TensorBoard plot illustrates how `spotPython` can be used as a microscope for the internal mechanisms of the surrogate-based optimization process. Here, one important parameter, the learning rate  $\theta$  of the Kriging surrogate is plotted against the number of optimization steps.

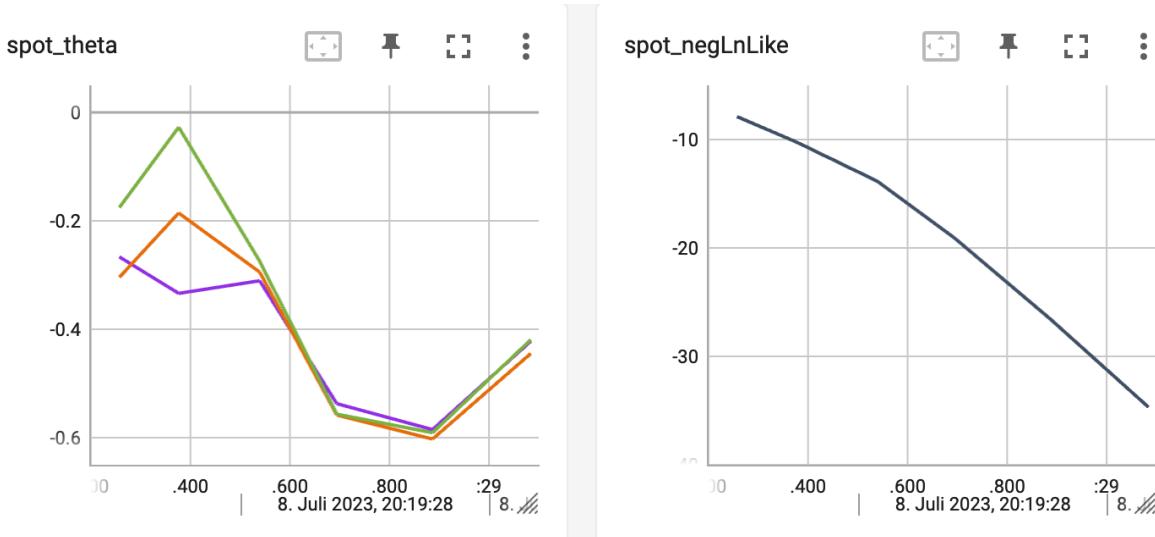


Figure 8.2: TensorBoard visualization of the spotPython surrogate model.

## 8.2 Conclusion

Based on this quick analysis, we can conclude that all three dimensions are equally important (as expected, because the analytical function is known).

## 8.3 Exercises

- Important:
  - Results from these exercises should be added to this document, i.e., you should submit an updated version of this notebook.
  - Please combine your results using this notebook.
  - Only one notebook from each group!
  - Presentation is based on this notebook. No additional slides are required!
  - spotPython version 0.16.11 (or greater) is required

### 8.3.1 The Three Dimensional fun\_cubed

- The input dimension is 3. The search range is  $-1 \leq x \leq 1$  for all dimensions.
- Generate contour plots
- Calculate the variable importance.
- Discuss the variable importance:

- Are all variables equally important?
- If not:
  - \* Which is the most important variable?
  - \* Which is the least important variable?

### **8.3.2 The Ten Dimensional `fun_wing_wt`**

- The input dimension is 10. The search range is  $0 \leq x \leq 1$  for all dimensions.
- Calculate the variable importance.
- Discuss the variable importance:
  - Are all variables equally important?
  - If not:
    - \* Which is the most important variable?
    - \* Which is the least important variable?
  - Generate contour plots for the three most important variables. Do they confirm your selection?

### **8.3.3 The Three Dimensional `fun_runge`**

- The input dimension is 3. The search range is  $-5 \leq x \leq 5$  for all dimensions.
- Generate contour plots
- Calculate the variable importance.
- Discuss the variable importance:
  - Are all variables equally important?
  - If not:
    - \* Which is the most important variable?
    - \* Which is the least important variable?

### **8.3.4 The Three Dimensional `fun_linear`**

- The input dimension is 3. The search range is  $-5 \leq x \leq 5$  for all dimensions.
- Generate contour plots
- Calculate the variable importance.
- Discuss the variable importance:
  - Are all variables equally important?
  - If not:
    - \* Which is the most important variable?
    - \* Which is the least important variable?

## 8.4 Jupyter Notebook

**i** Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

# 9 Isotropic and Anisotropic Kriging

This chapter illustrates the difference between isotropic and anisotropic Kriging models. The difference is illustrated with the help of the `spotPython` package. Isotropic Kriging models use the same `theta` value for every dimension. Anisotropic Kriging models use different `theta` values for each dimension.

## 9.1 Example: Isotropic Spot Surrogate and the 2-dim Sphere Function

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from spotPython.utils.init import fun_control_init, surrogate_control_init
```

### 9.1.1 The Objective Function: 2-dim Sphere

- The `spotPython` package provides several classes of objective functions.
- We will use an analytical objective function, i.e., a function that can be described by a (closed) formula:

$$f(x, y) = x^2 + y^2$$

```
fun = analytical().fun_sphere
fun_control = fun_control_init(lower = np.array([-1, -1]),
                                upper = np.array([1, 1]))
```

- The size of the `lower` bound vector determines the problem dimension.
- Here we will use `np.array([-1, -1])`, i.e., a two-dim function.

```
spot_2 = spot.Spot(fun=fun,
                    fun_control=fun_control)

spot_2.run()
```

```

spotPython tuning: 1.961941240809186e-05 [#####---] 73.33%
spotPython tuning: 1.961941240809186e-05 [#####---] 80.00%
spotPython tuning: 1.961941240809186e-05 [#####---] 86.67%
spotPython tuning: 1.961941240809186e-05 [#####---] 93.33%
spotPython tuning: 1.961941240809186e-05 [#####---] 100.00% Done...

```

```
<spotPython.spot.spot.Spot at 0x2ccc4a250>
```

### 9.1.2 Results

```
spot_2.print_results()
```

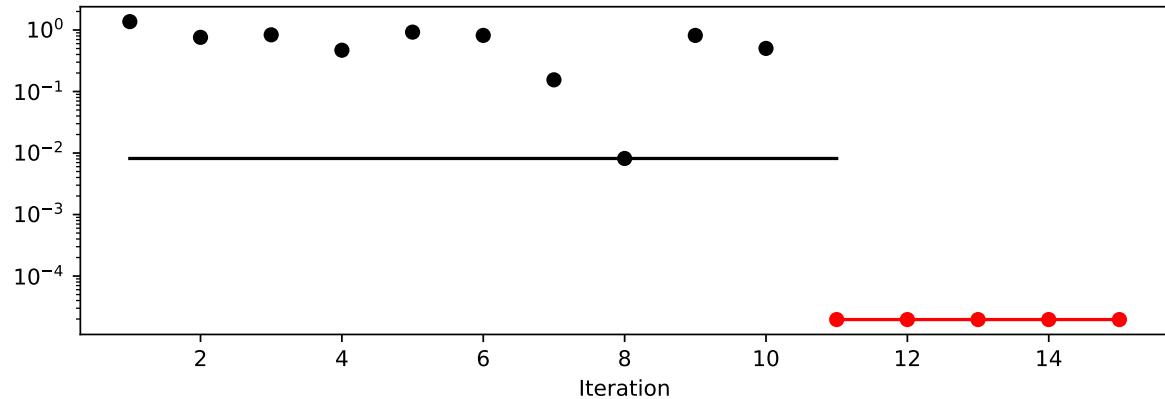
```

min y: 1.961941240809186e-05
x0: 0.0016552607778336345
x1: 0.004108469808268944

```

```
[['x0', 0.0016552607778336345], ['x1', 0.004108469808268944]]
```

```
spot_2.plot_progress(log_y=True)
```



## 9.2 Example With Anisotropic Kriging

- The default parameter setting of `spotPython`'s Kriging surrogate uses the same `theta` value for every dimension.
- This is referred to as “using an isotropic kernel”.

- If different `theta` values are used for each dimension, then an anisotropic kernel is used
- To enable anisotropic models in `spotPython`, the number of `theta` values should be larger than one.
- We can use `surrogate_control={"n_theta": 2}` to enable this behavior (2 is the problem dimension).

```
:::{.callout-note} ##### TensorBoard
```

Similar to the one-dimensional case, which was introduced in Section Section 7.5, we can use TensorBoard to monitor the progress of the optimization. We will use the same code, only the prefix is different:

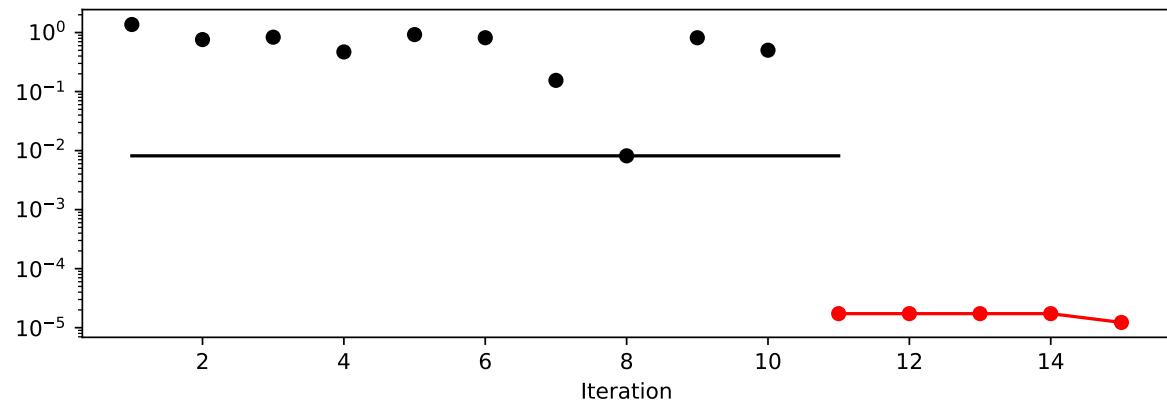
```
from spotPython.utils.init import fun_control_init
PREFIX = "03"
fun_control = fun_control_init(
    PREFIX=PREFIX,
    lower = np.array([-1, -1]),
    upper = np.array([1, 1]))
surrogate_control = surrogate_control_init(n_theta=2)
spot_2_anisotropic = spot.Spot(fun=fun,
    fun_control=fun_control,
    surrogate_control=surrogate_control)
spot_2_anisotropic.run()
```

```
Created spot_tensorboard_path: runs/spot_logs/03_p040025_2024-01-09_20-12-06 for SummaryWriter
spotPython tuning: 1.7290078116309417e-05 [#####---] 73.33%
spotPython tuning: 1.7290078116309417e-05 [#####---] 80.00%
spotPython tuning: 1.7290078116309417e-05 [#####---] 86.67%
spotPython tuning: 1.7290078116309417e-05 [#####---] 93.33%
spotPython tuning: 1.2267580433005057e-05 [#####---] 100.00% Done...
```

```
<spotPython.spot.spot.Spot at 0x2c2709e10>
```

- The search progress of the optimization with the anisotropic model can be visualized:

```
spot_2_anisotropic.plot_progress(log_y=True)
```

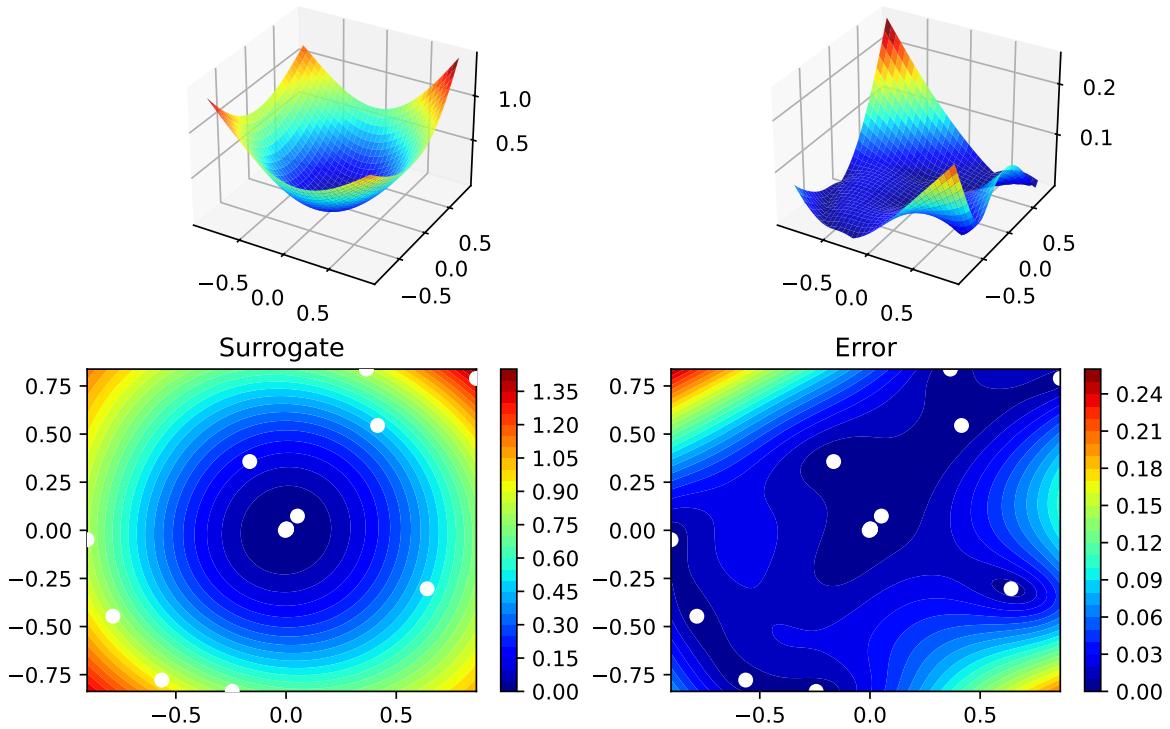


```
spot_2_anisotropic.print_results()
```

```
min y: 1.2267580433005057e-05
x0: -0.003457387366511919
x1: -0.000560404345886819
```

```
[['x0', -0.003457387366511919], ['x1', -0.000560404345886819]]
```

```
spot_2_anisotropic.surrogate.plot()
```



## 9.2.1 Taking a Look at the theta Values

### 9.2.1.1 theta Values from the spot Model

- We can check, whether one or several `theta` values were used.
- The `theta` values from the surrogate can be printed as follows:

```
spot_2_anisotropic.surrogate.theta
```

```
array([-0.30709835, -0.13263039])
```

- Since the surrogate from the isotropic setting was stored as `spot_2`, we can also take a look at the `theta` value from this model:

```
spot_2.surrogate.theta
```

```
array([-0.15830282])
```

### 9.2.1.2 TensorBoard

Now we can start TensorBoard in the background with the following command:

```
tensorboard --logdir=". ./runs"
```

We can access the TensorBoard web server with the following URL:

```
http://localhost:6006/
```

The TensorBoard plot illustrates how `spotPython` can be used as a microscope for the internal mechanisms of the surrogate-based optimization process. Here, one important parameter, the learning rate  $\theta$  of the Kriging surrogate is plotted against the number of optimization steps.

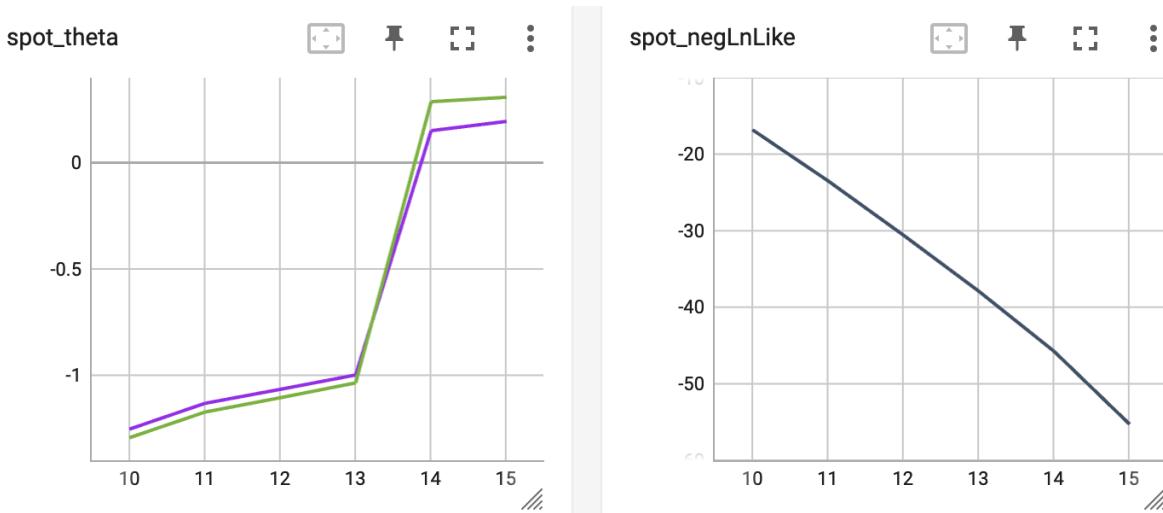


Figure 9.1: TensorBoard visualization of the `spotPython` surrogate model.

## 9.3 Exercises

### 9.3.1 fun\_branin

- Describe the function.
  - The input dimension is 2. The search range is  $-5 \leq x_1 \leq 10$  and  $0 \leq x_2 \leq 15$ .
- Compare the results from `spotPython` run a) with isotropic and b) anisotropic surrogate models.

- Modify the termination criterion: instead of the number of evaluations (which is specified via `fun_evals`), the time should be used as the termination criterion. This can be done as follows (`max_time=1` specifies a run time of one minute):

```
fun_evals=inf,
max_time=1,
```

### 9.3.2 fun\_sin\_cos

- Describe the function.
  - The input dimension is 2. The search range is  $-2\pi \leq x_1 \leq 2\pi$  and  $-2\pi \leq x_2 \leq 2\pi$ .
- Compare the results from `spotPython` run a) with isotropic and b) anisotropic surrogate models.
- Modify the termination criterion (`max_time` instead of `fun_evals`) as described for `fun_branin`.

### 9.3.3 fun\_runge

- Describe the function.
  - The input dimension is 2. The search range is  $-5 \leq x_1 \leq 5$  and  $-5 \leq x_2 \leq 5$ .
- Compare the results from `spotPython` run a) with isotropic and b) anisotropic surrogate models.
- Modify the termination criterion (`max_time` instead of `fun_evals`) as described for `fun_branin`.

### 9.3.4 fun\_wingwt

- Describe the function.
  - The input dimension is 10. The search ranges are between 0 and 1 (values are mapped internally to their natural bounds).
- Compare the results from `spotPython` run a) with isotropic and b) anisotropic surrogate models.
- Modify the termination criterion (`max_time` instead of `fun_evals`) as described for `fun_branin`.

## 9.4 Jupyter Notebook

**i** Note

- The Jupyter-Notebook of this lecture is available on GitHub in the [Hyperparameter-Tuning-Cookbook Repository](#)

# 10 Using sklearn Surrogates in spotPython

Besides the internal kriging surrogate, which is used as a default by `spotPython`, any surrogate model from `scikit-learn` can be used as a surrogate in `spotPython`. This chapter explains how to use `scikit-learn` surrogates in `spotPython`.

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
```

## 10.1 Example: Branin Function with spotPython's Internal Kriging Surrogate

### 10.1.1 The Objective Function Branin

- The `spotPython` package provides several classes of objective functions.
- We will use an analytical objective function, i.e., a function that can be described by a (closed) formula.
- Here we will use the Branin function:

```
y = a * (x2 - b * x1**2 + c * x1 - r) ** 2 + s * (1 - t) * np.cos(x1) + s,
where values of a, b, c, r, s and t are: a = 1, b = 5.1 / (4*pi**2),
c = 5 / pi, r = 6, s = 10 and t = 1 / (8*pi).
```

- It has three global minima:

```
f(x) = 0.397887 at (-pi, 12.275), (pi, 2.275), and (9.42478, 2.475).
```

```
from spotPython.fun.objectivefunctions import analytical
fun = analytical().fun_branin
```

## TensorBoard

Similar to the one-dimensional case, which was introduced in Section [Section 7.5](#), we can use TensorBoard to monitor the progress of the optimization. We will use the same code, only the prefix is different:

```
from spotPython.utils.init import fun_control_init, design_control_init
PREFIX = "04"
fun_control = fun_control_init(
    PREFIX=PREFIX,
    lower = np.array([-5,-0]),
    upper = np.array([10,15]),
    fun_evals=20,
    max_time=inf)

design_control = design_control_init(
    init_size=10)
```

Created spot\_tensorboard\_path: runs/spot\_logs/04\_p040025\_2024-01-09\_20-13-12 for SummaryWriter

### 10.1.2 Running the surrogate model based optimizer Spot:

```
spot_2 = spot.Spot(fun=fun,
                    fun_control=fun_control,
                    design_control=design_control)
```

```
spot_2.run()
```

```
spotPython tuning: 3.1468324340130396 [#####----] 55.00%
spotPython tuning: 3.1468324340130396 [#####----] 60.00%
spotPython tuning: 3.1468324340130396 [#####----] 65.00%
spotPython tuning: 3.1468324340130396 [#####---] 70.00%
spotPython tuning: 1.1487178651597567 [#####---] 75.00%
spotPython tuning: 1.023873960670663 [#####---] 80.00%
spotPython tuning: 0.4157039136851921 [#####---] 85.00%
spotPython tuning: 0.40177808730611275 [#####---] 90.00%
spotPython tuning: 0.39909168407470474 [#####---] 95.00%
spotPython tuning: 0.39909168407470474 [#####---] 100.00% Done...
```

```
<spotPython.spot.spot.Spot at 0x2d4150f10>
```

### 10.1.3 TensorBoard

Now we can start TensorBoard in the background with the following command:

```
tensorboard --logdir=". ./runs"
```

We can access the TensorBoard web server with the following URL:

```
http://localhost:6006/
```

The TensorBoard plot illustrates how `spotPython` can be used as a microscope for the internal mechanisms of the surrogate-based optimization process. Here, one important parameter, the learning rate  $\theta$  of the Kriging surrogate is plotted against the number of optimization steps.

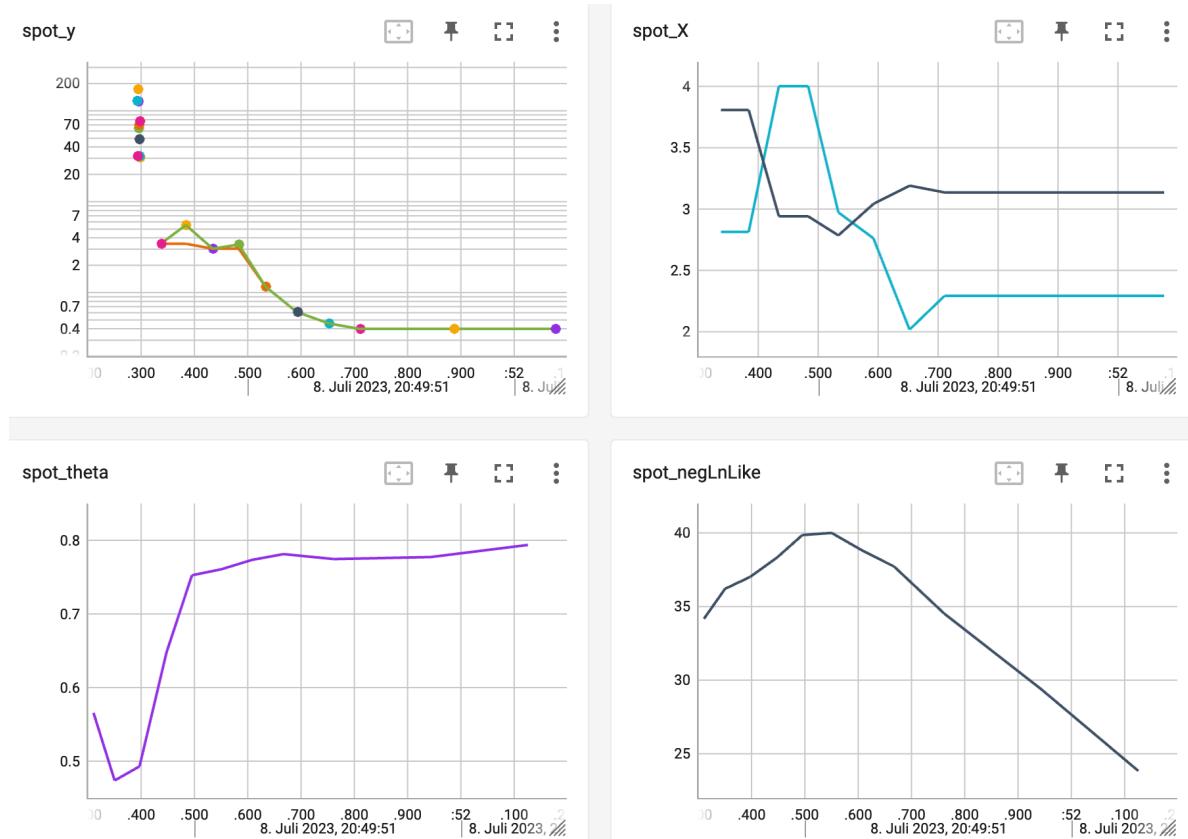


Figure 10.1: TensorBoard visualization of the `spotPython` optimization process and the surrogate model.

#### 10.1.4 Print the Results

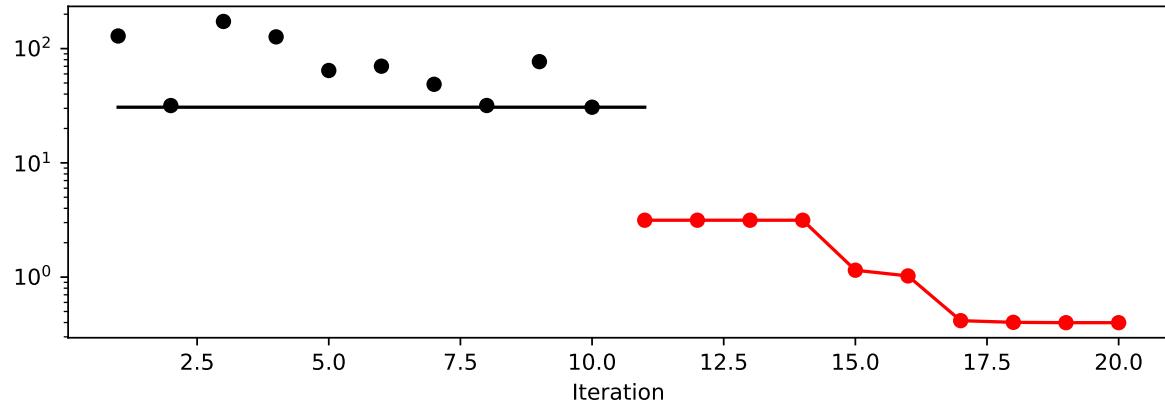
```
spot_2.print_results()
```

```
min y: 0.39909168407470474
x0: 3.1523096108334228
x1: 2.292209322210569

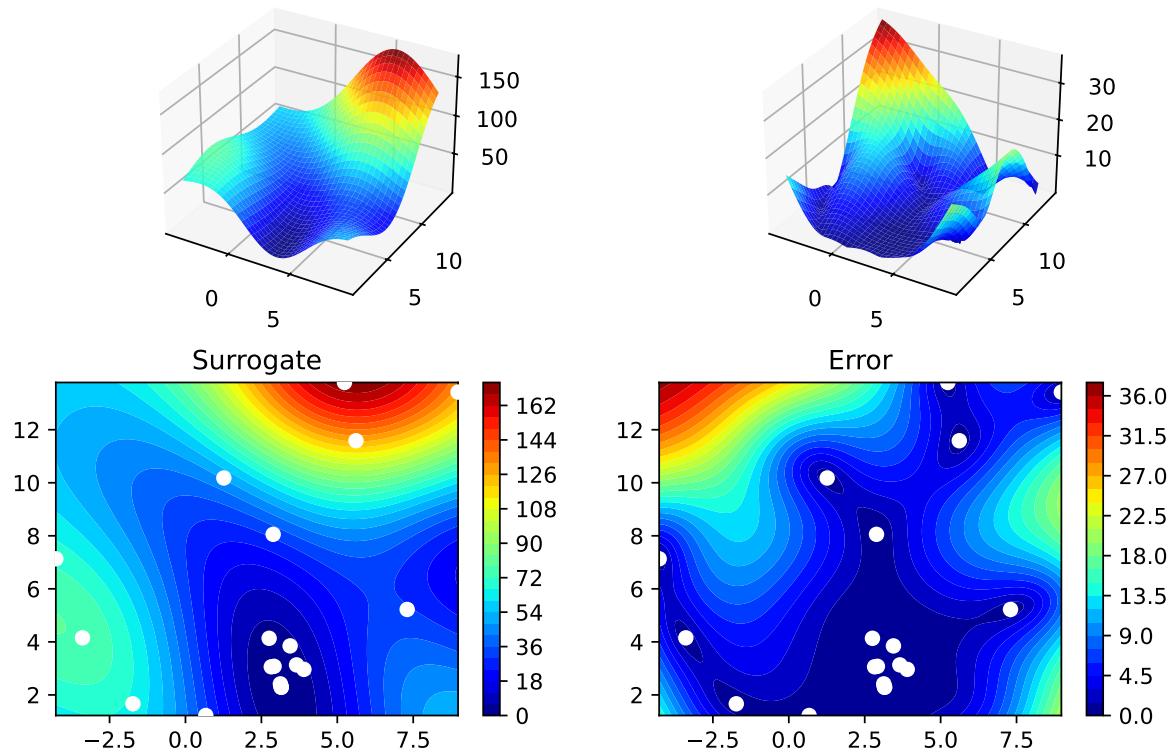
[['x0', 3.1523096108334228], ['x1', 2.292209322210569]]
```

#### 10.1.5 Show the Progress and the Surrogate

```
spot_2.plot_progress(log_y=True)
```



```
spot_2.surrogate.plot()
```



## 10.2 Example: Using Surrogates From scikit-learn

- Default is the `spotPython` (i.e., the internal) kriging surrogate.
- It can be called explicitly and passed to `Spot`.

```
from spotPython.build.kriging import Kriging
S_0 = Kriging(name='kriging', seed=123)
```

- Alternatively, models from `scikit-learn` can be selected, e.g., Gaussian Process, RBFs, Regression Trees, etc.

```
# Needed for the sklearn surrogates:
from sklearn.gaussian_process import GaussianProcessRegressor
from sklearn.gaussian_process.kernels import RBF
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn import linear_model
```

```
from sklearn import tree
import pandas as pd
```

- Here are some additional models that might be useful later:

```
S_Tree = DecisionTreeRegressor(random_state=0)
S_LM = linear_model.LinearRegression()
S_Ridge = linear_model.Ridge()
S_RF = RandomForestRegressor(max_depth=2, random_state=0)
```

### 10.2.1 GaussianProcessRegressor as a Surrogate

- To use a Gaussian Process model from `sklearn`, that is similar to `spotPython`'s Kriging, we can proceed as follows:

```
kernel = 1 * RBF(length_scale=1.0, length_scale_bounds=(1e-2, 1e2))
S_GP = GaussianProcessRegressor(kernel=kernel, n_restarts_optimizer=9)
```

- The scikit-learn GP model `S_GP` is selected for `Spot` as follows:

```
surrogate = S_GP
```

- We can check the kind of surrogate model with the command `isinstance`:

```
isinstance(S_GP, GaussianProcessRegressor)
```

True

```
isinstance(S_0, Kriging)
```

True

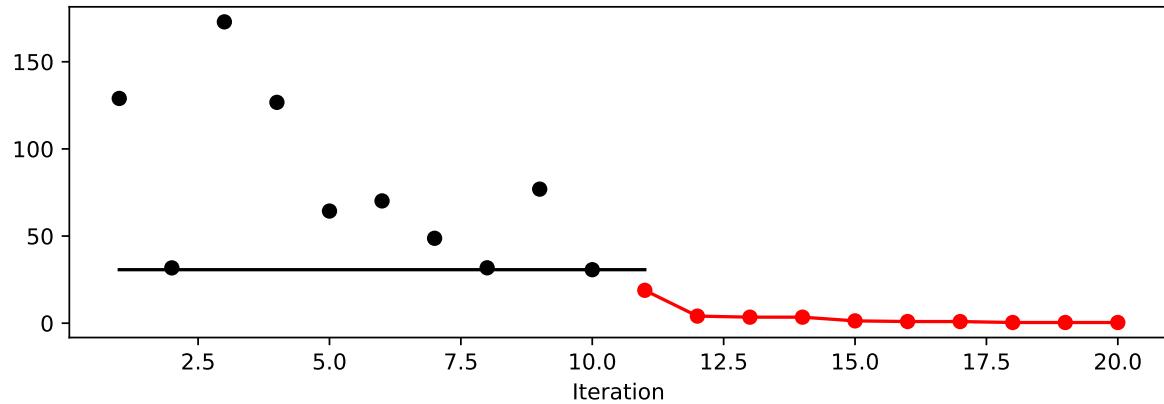
- Similar to the `Spot` run with the internal Kriging model, we can call the run with the `scikit-learn` surrogate:

```
fun = analytical(seed=123).fun_branin
spot_2_GP = spot.Spot(fun=fun,
                      fun_control=fun_control,
                      design_control=design_control,
                      surrogate = S_GP)
spot_2_GP.run()
```

```
spotPython tuning: 18.865121449825782 [#####----] 55.00%
spotPython tuning: 4.06700305855078 [#####----] 60.00%
spotPython tuning: 3.461906927549384 [#####----] 65.00%
spotPython tuning: 3.461906927549384 [#####---] 70.00%
spotPython tuning: 1.3280944252046556 [#####---] 75.00%
spotPython tuning: 0.9548334920645392 [#####---] 80.00%
spotPython tuning: 0.9344485781421579 [#####---] 85.00%
spotPython tuning: 0.39916716809341857 [#####---] 90.00%
spotPython tuning: 0.3982254000779708 [#####---] 95.00%
spotPython tuning: 0.3982254000779708 [#####---] 100.00% Done...
```

```
<spotPython.spot.spot.Spot at 0x2d9a82f10>
```

```
spot_2_GP.plot_progress()
```



```
spot_2_GP.print_results()
```

```
min y: 0.3982254000779708
x0: 3.1499822680266343
x1: 2.268811272474469

[['x0', 3.1499822680266343], ['x1', 2.268811272474469]]
```

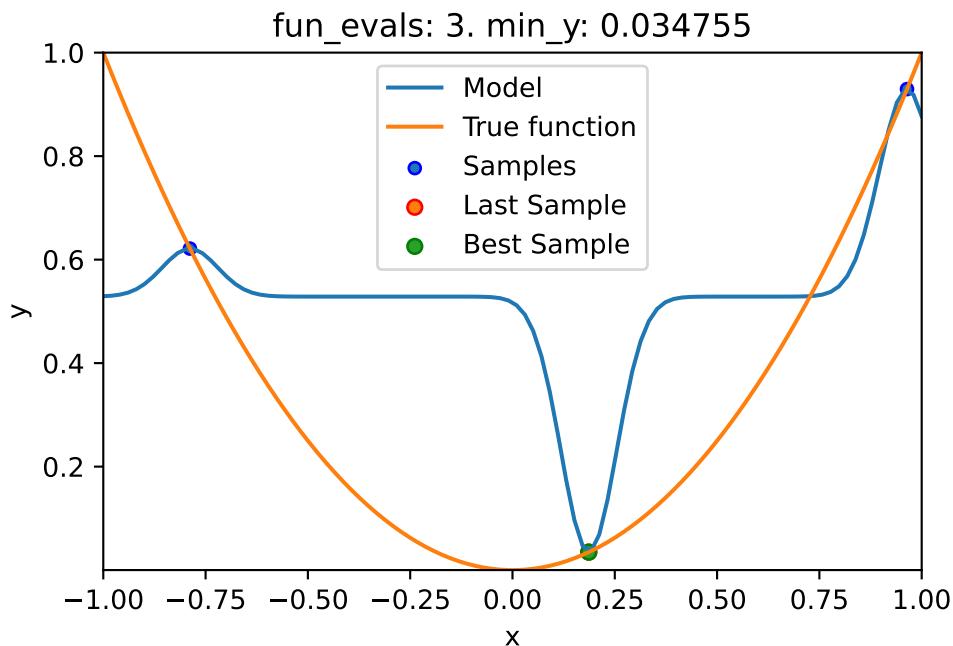
### 10.3 Example: One-dimensional Sphere Function With spotPython's Kriging

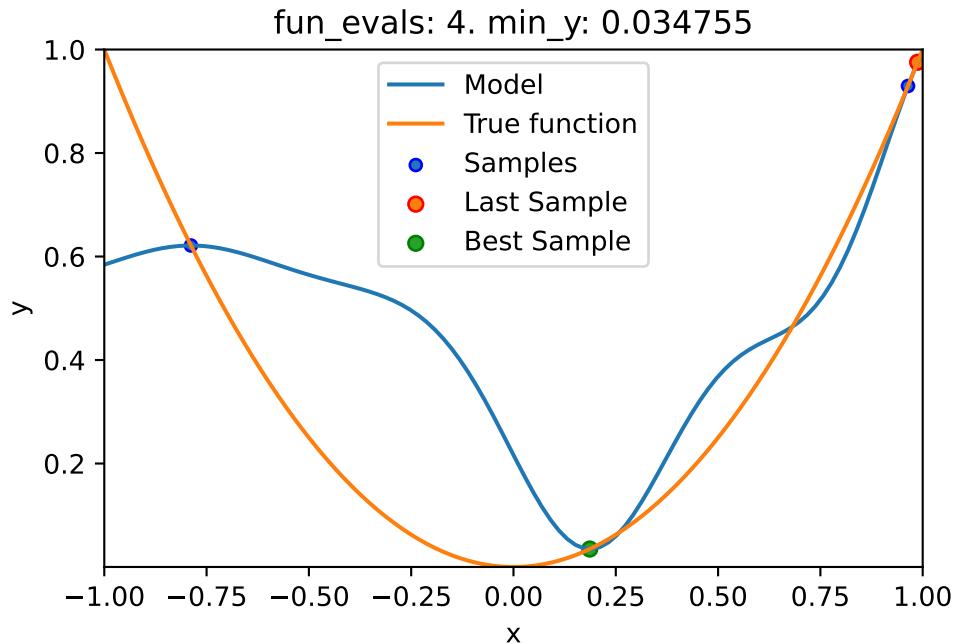
- In this example, we will use an one-dimensional function, which allows us to visualize the optimization process.

– `show_models= True` is added to the argument list.

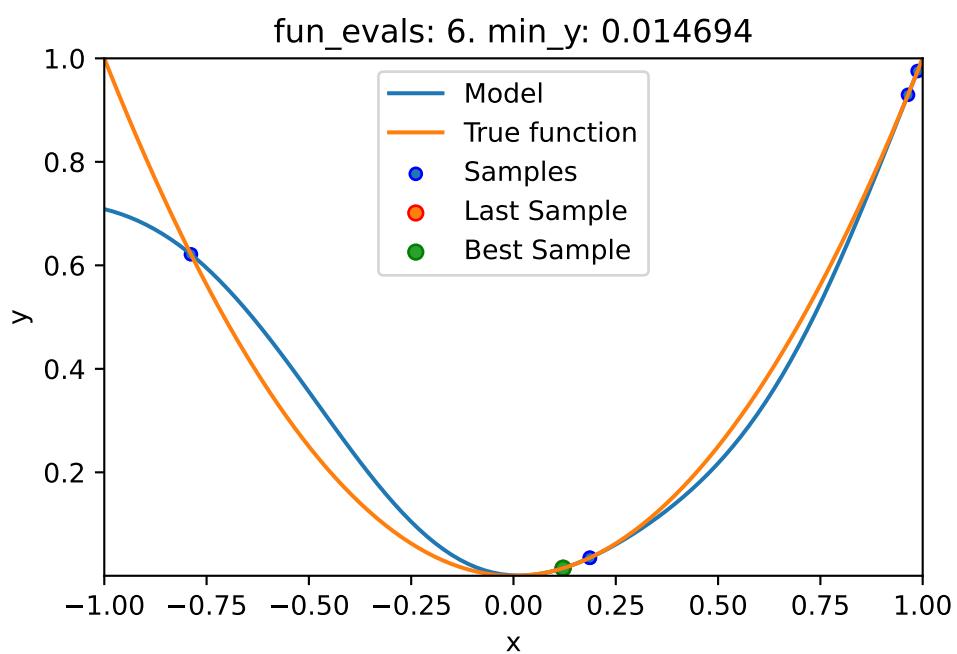
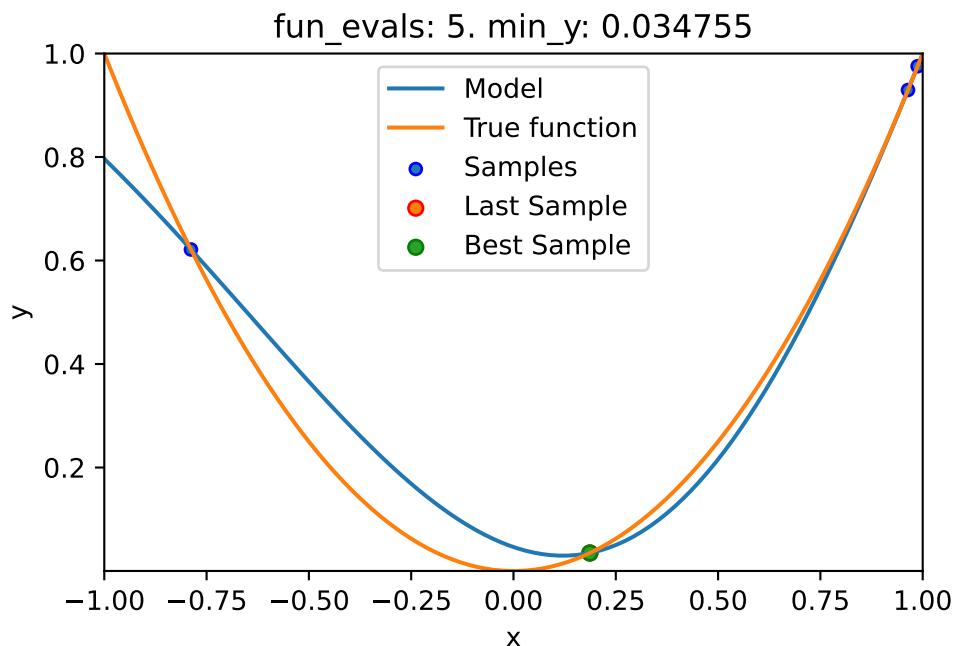
```
from spotPython.fun.objectivefunctions import analytical
fun_control = fun_control_init(
    lower = np.array([-1]),
    upper = np.array([1]),
    fun_evals=10,
    max_time=inf,
    show_models= True,
    tolerance_x = np.sqrt(np.spacing(1)))
fun = analytical(seed=123).fun_sphere
design_control = design_control_init(
    init_size=3)
```

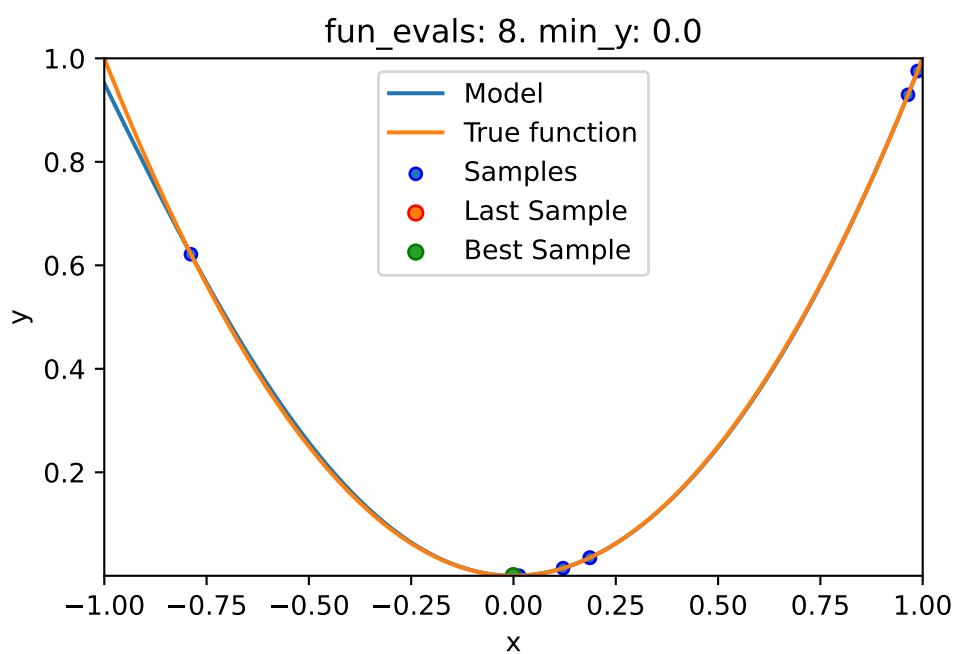
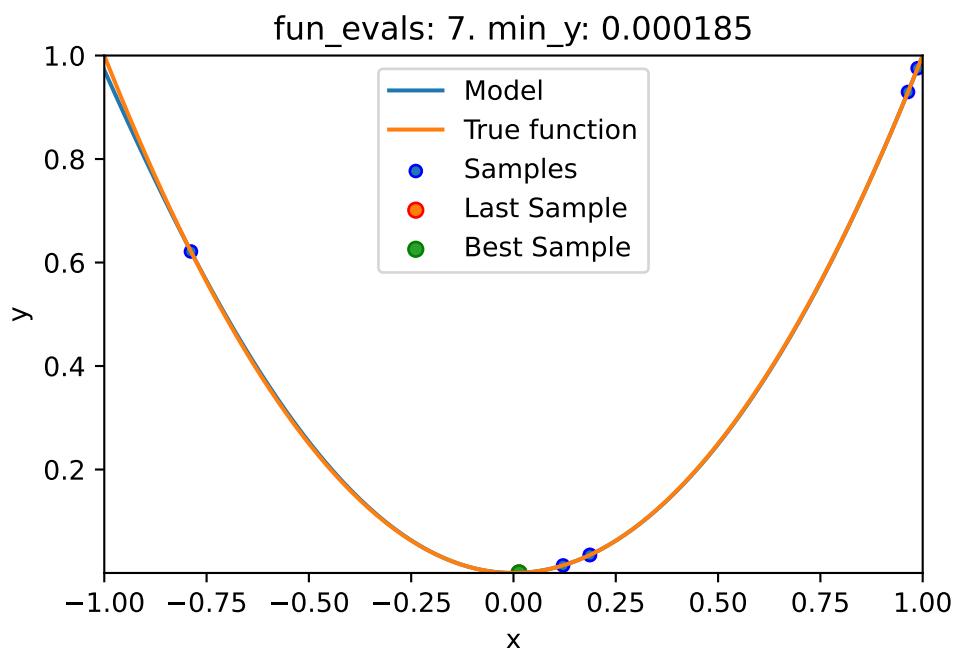
```
spot_1 = spot.Spot(fun=fun,
                    fun_control=fun_control,
                    design_control=design_control)
spot_1.run()
```

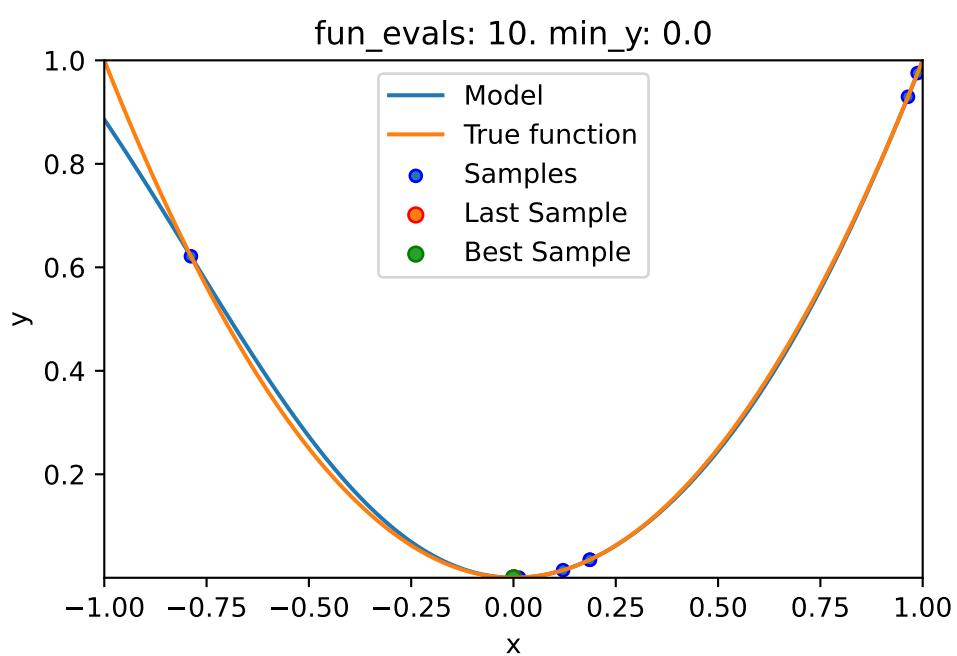
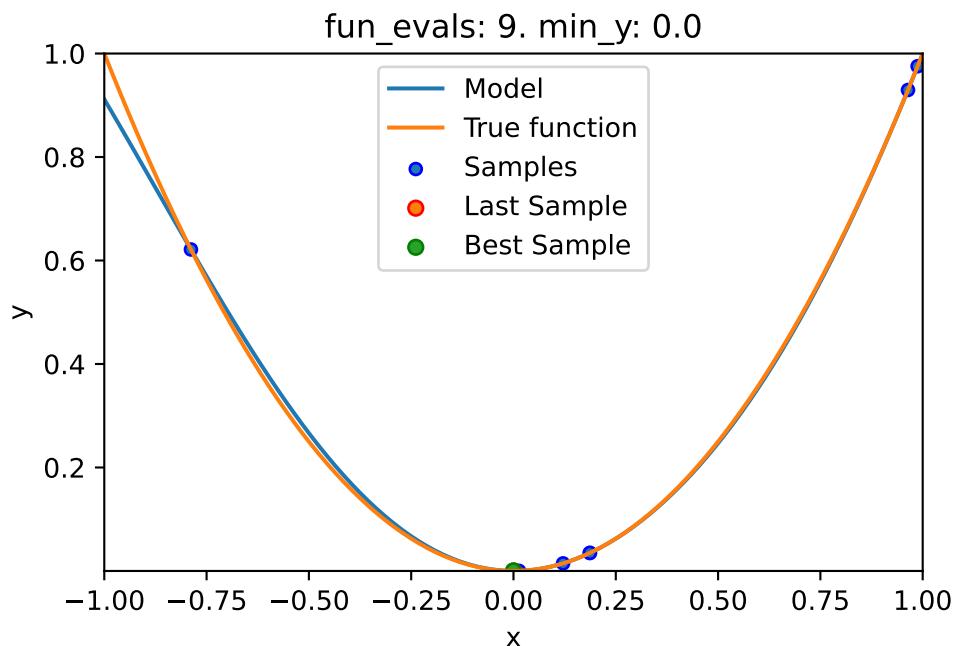




```
spotPython tuning: 0.03475493366922229 [#####-----] 40.00%
spotPython tuning: 0.03475493366922229 [#####-----] 50.00%
spotPython tuning: 0.014693634061986652 [#####----] 60.00%
spotPython tuning: 0.0001850752755504842 [#####---] 70.00%
spotPython tuning: 4.643556417950029e-08 [#####--] 80.00%
spotPython tuning: 4.643556417950029e-08 [#####-] 90.00%
spotPython tuning: 4.643556417950029e-08 [#####] 100.00% Done...
```







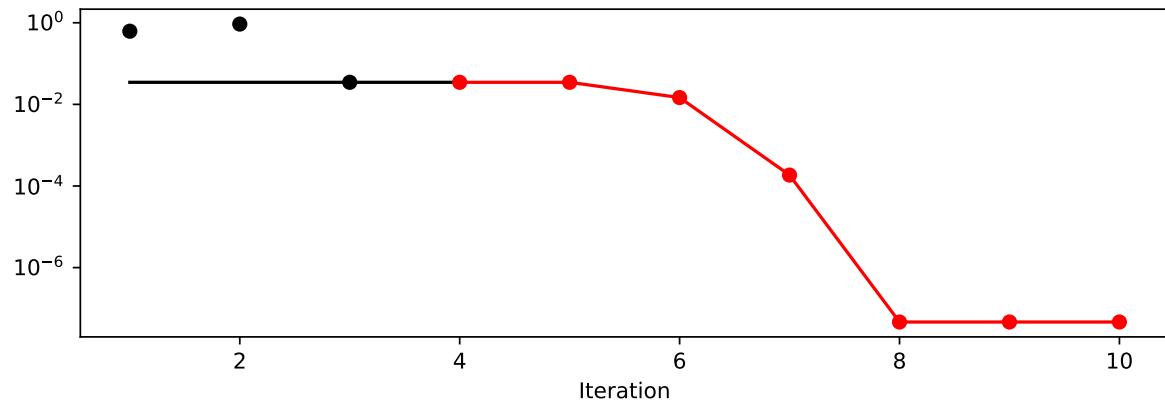
### 10.3.1 Results

```
spot_1.print_results()
```

```
min y: 4.643556417950029e-08
x0: -0.0002154891277524235
```

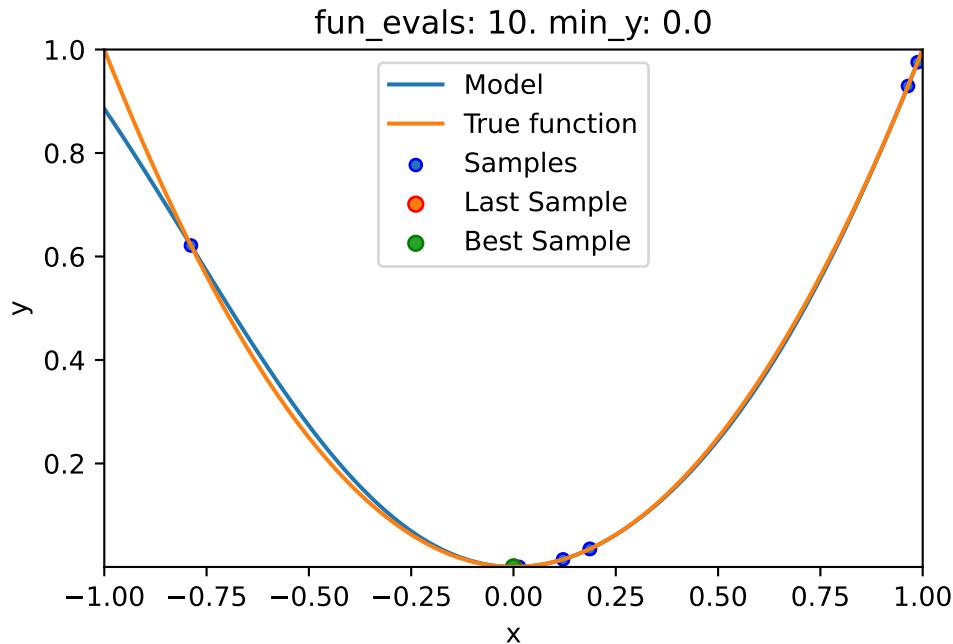
```
[['x0', -0.0002154891277524235]]
```

```
spot_1.plot_progress(log_y=True)
```



- The method `plot_model` plots the final surrogate:

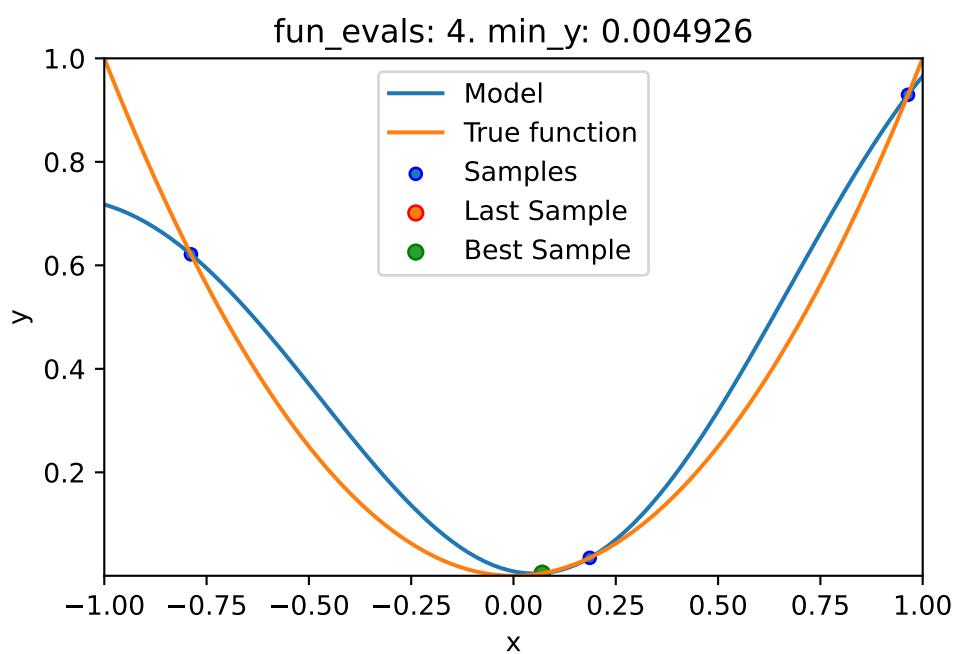
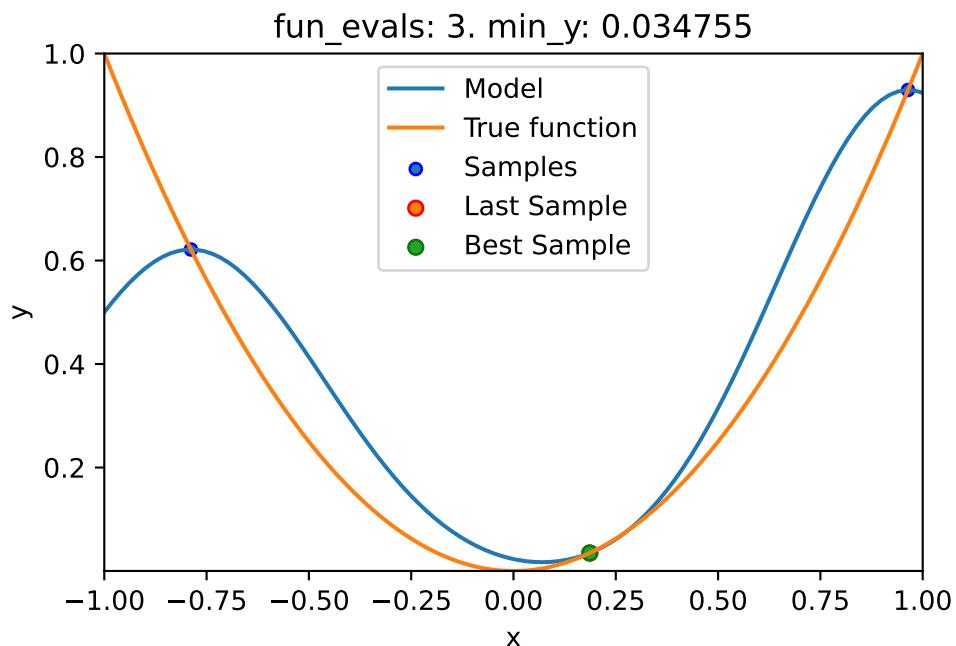
```
spot_1.plot_model()
```



## 10.4 Example: Sklearn Model GaussianProcess

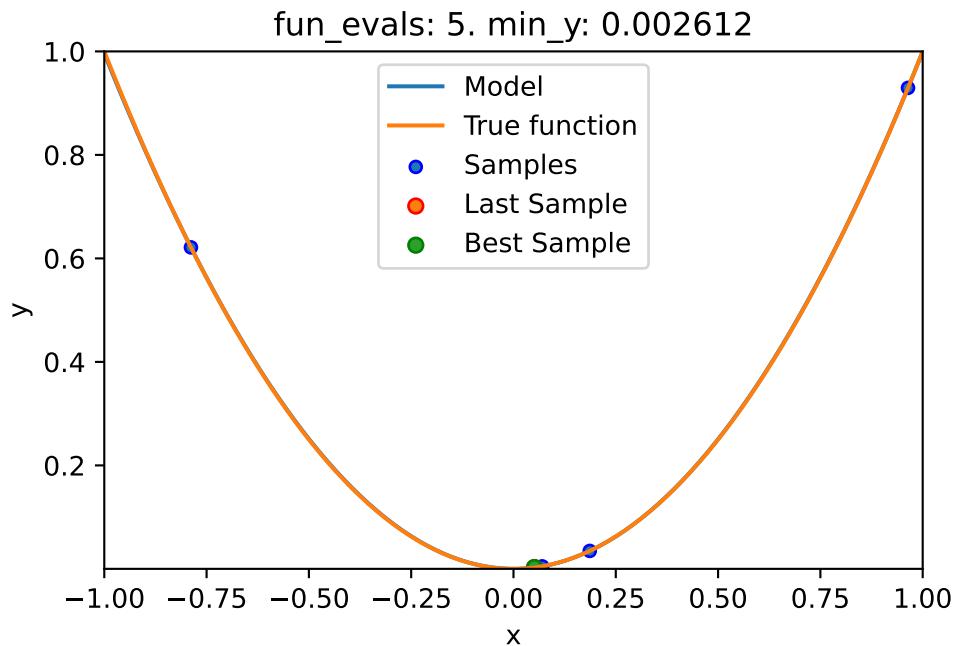
- This example visualizes the search process on the `GaussianProcessRegression` surrogate from `sklearn`.
- Therefore `surrogate = S_GP` is added to the argument list.

```
fun = analytical(seed=123).fun_sphere
spot_1_GP = spot.Spot(fun=fun,
                      fun_control=fun_control,
                      design_control=design_control,
                      surrogate = S_GP)
spot_1_GP.run()
```

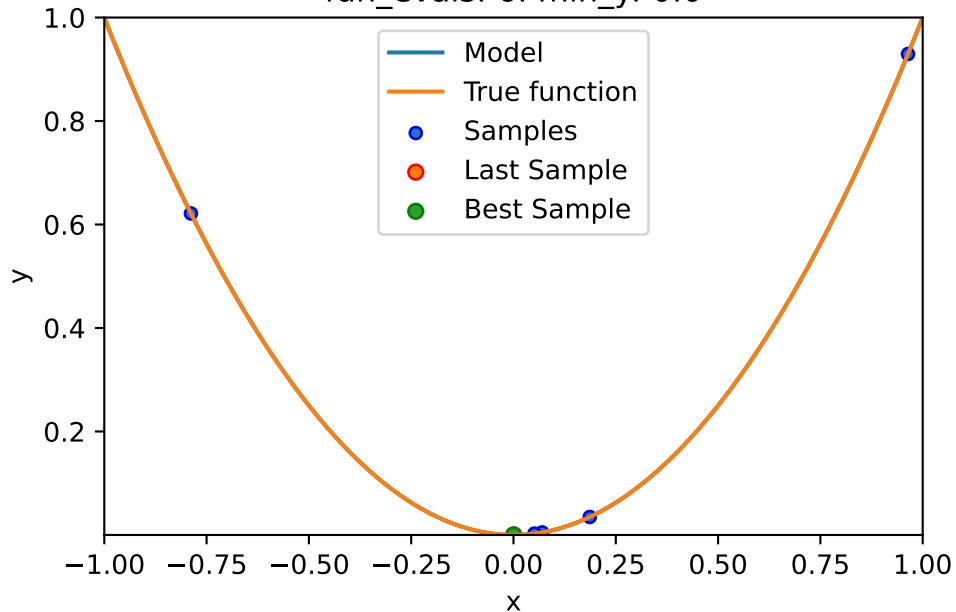


```
spotPython tuning: 0.004925671374769521 [#####-----] 40.00%
spotPython tuning: 0.002612062924748803 [#####-----] 50.00%
spotPython tuning: 3.6666409852957783e-07 [#####-----] 60.00%
```

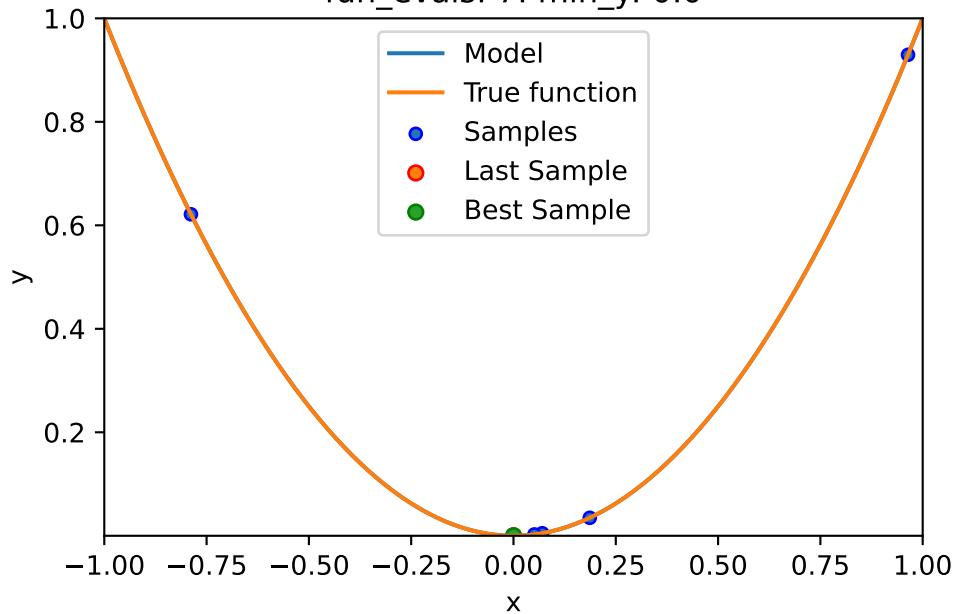
```
spotPython tuning: 4.638244203084832e-08 [#####---] 70.00%
spotPython tuning: 3.2711094860544125e-09 [#####--] 80.00%
spotPython tuning: 2.2493573831304313e-10 [#####--] 90.00%
spotPython tuning: 2.2493573831304313e-10 [#######] 100.00% Done...
```



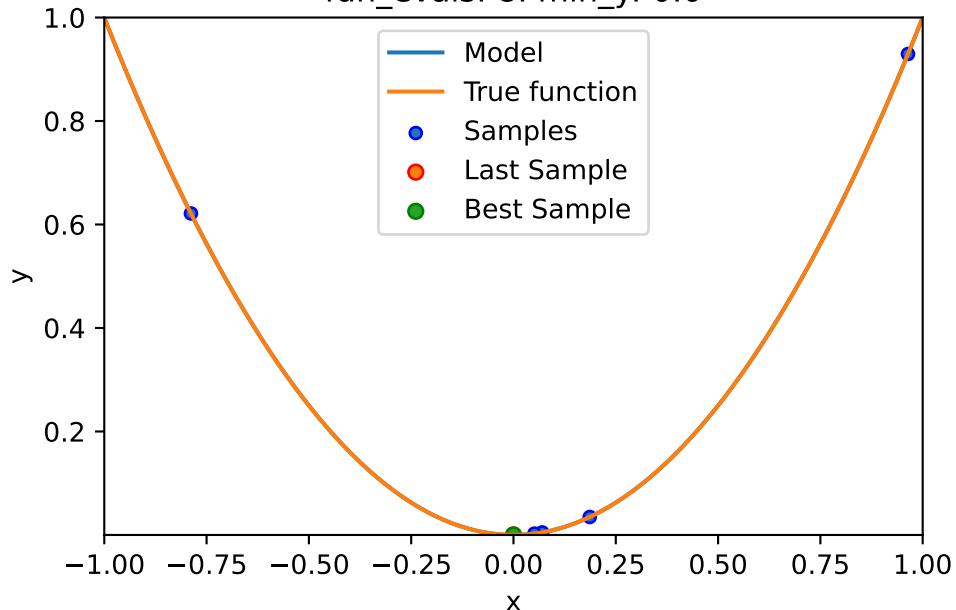
fun\_evals: 6. min\_y: 0.0



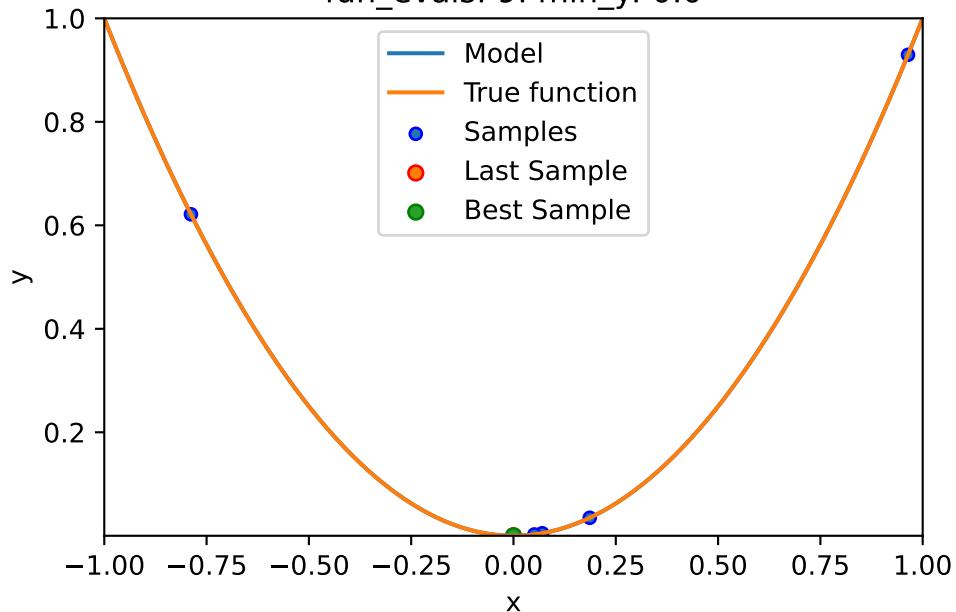
fun\_evals: 7. min\_y: 0.0

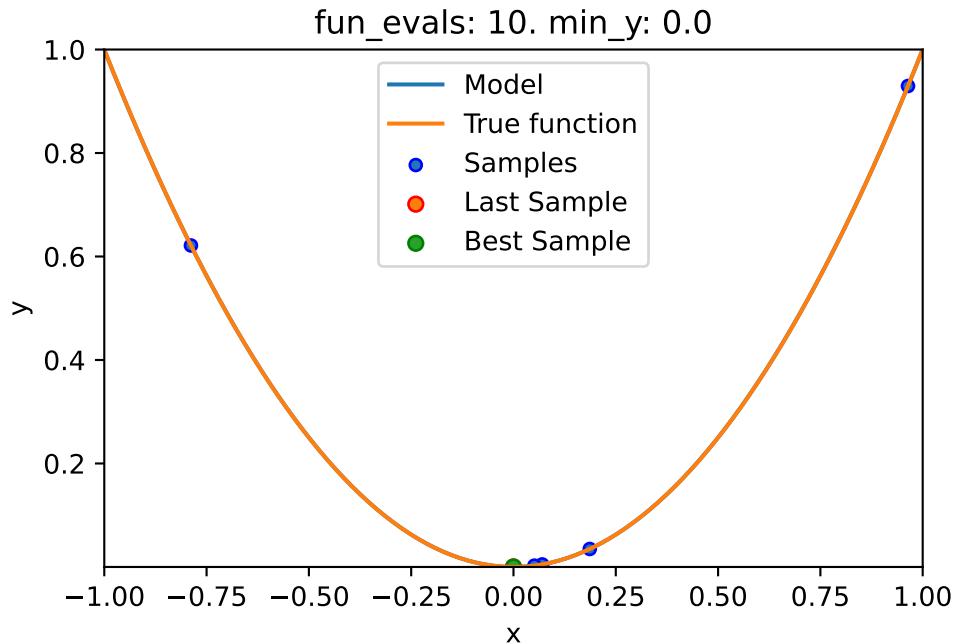


fun\_evals: 8. min\_y: 0.0



fun\_evals: 9. min\_y: 0.0



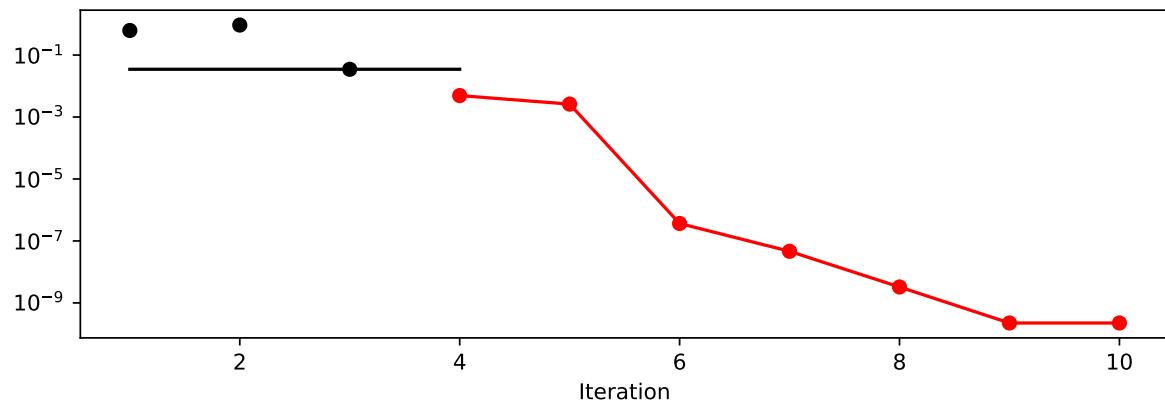


```
spot_1_GP.print_results()
```

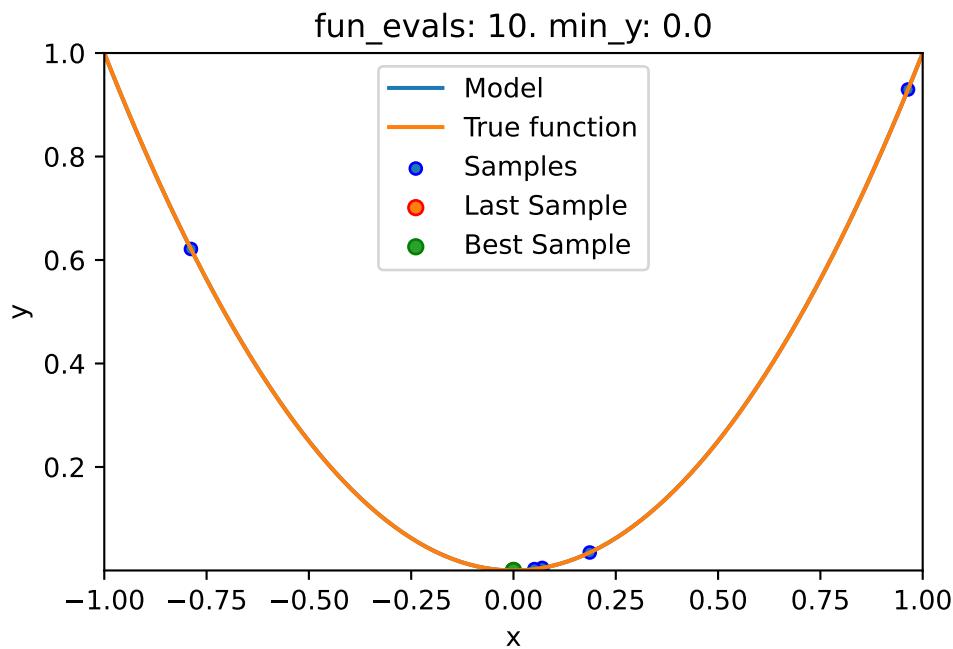
```
min y: 2.2493573831304313e-10
x0: 1.499785779079943e-05
```

```
[['x0', 1.499785779079943e-05]]
```

```
spot_1_GP.plot_progress(log_y=True)
```



```
spot_1_GP.plot_model()
```



## 10.5 Exercises

### 10.5.1 DecisionTreeRegressor

- Describe the surrogate model.
- Use the surrogate as the model for optimization.

### 10.5.2 RandomForestRegressor

- Describe the surrogate model.
- Use the surrogate as the model for optimization.

### 10.5.3 linear\_model.LinearRegression

- Describe the surrogate model.
- Use the surrogate as the model for optimization.

#### **10.5.4 `linear_model.Ridge`**

- Describe the surrogate model.
- Use the surrogate as the model for optimization.

### **10.6 Exercise 2**

- Compare the performance of the five different surrogates on both objective functions:
  - `spotPython`'s internal Kriging
  - `DecisionTreeRegressor`
  - `RandomForestRegressor`
  - `linear_model.LinearRegression`
  - `linear_model.Ridge`

# 11 Sequential Parameter Optimization: Gaussian Process Models

This chapter analyzes differences between the Kriging implementation in `spotPython` and the `GaussianProcessRegressor` in `scikit-learn`.

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.design.spacefilling import spacefilling
from spotPython.spot import spot
from spotPython.build.kriging import Kriging
from scipy.optimize import shgo
from scipy.optimize import direct
from scipy.optimize import differential_evolution
import matplotlib.pyplot as plt
import math as m
from sklearn.gaussian_process import GaussianProcessRegressor
from sklearn.gaussian_process.kernels import RBF
```

## 11.1 Gaussian Processes Regression: Basic Introductory `scikit-learn` Example

- This is the example from `scikit-learn`: [https://scikit-learn.org/stable/auto\\_examples/gaussian\\_process/pl](https://scikit-learn.org/stable/auto_examples/gaussian_process/pl)
- After fitting our model, we see that the hyperparameters of the kernel have been optimized.
- Now, we will use our kernel to compute the mean prediction of the full dataset and plot the 95% confidence interval.

### 11.1.1 Train and Test Data

```

X = np.linspace(start=0, stop=10, num=1_000).reshape(-1, 1)
y = np.squeeze(X * np.sin(X))
rng = np.random.RandomState(1)
training_indices = rng.choice(np.arange(y.size), size=6, replace=False)
X_train, y_train = X[training_indices], y[training_indices]

```

### 11.1.2 Building the Surrogate With Sklearn

- The model building with `sklearn` consists of three steps:
  1. Instantiating the model, then
  2. fitting the model (using `fit`), and
  3. making predictions (using `predict`)

```

kernel = 1 * RBF(length_scale=1.0, length_scale_bounds=(1e-2, 1e2))
gaussian_process = GaussianProcessRegressor(kernel=kernel, n_restarts_optimizer=9)
gaussian_process.fit(X_train, y_train)
mean_prediction, std_prediction = gaussian_process.predict(X, return_std=True)

```

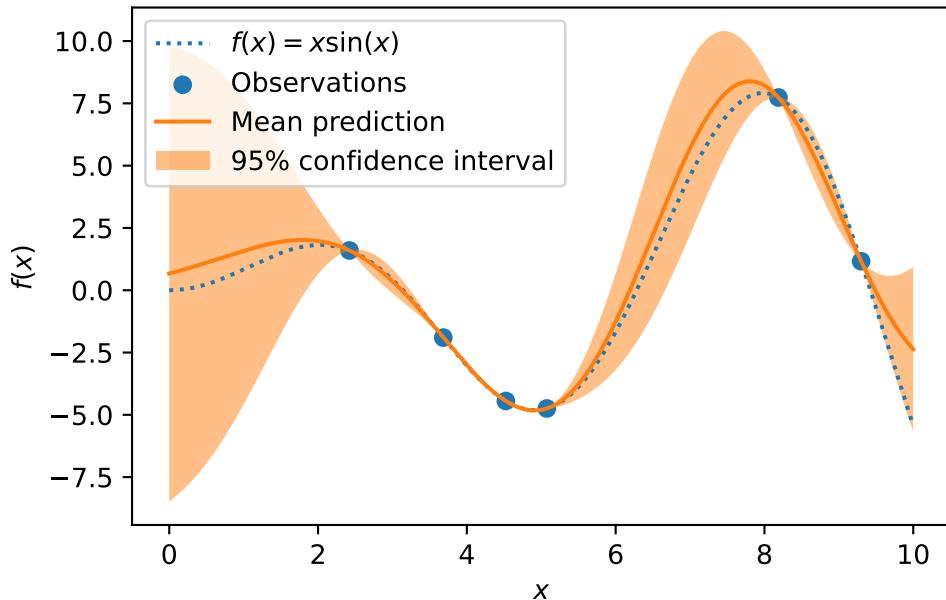
### 11.1.3 Plotting the SklearnModel

```

plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, mean_prediction, label="Mean prediction")
plt.fill_between(
    X.ravel(),
    mean_prediction - 1.96 * std_prediction,
    mean_prediction + 1.96 * std_prediction,
    alpha=0.5,
    label=r"95% confidence interval",
)
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("sk-learn Version: Gaussian process regression on noise-free dataset")

```

## sk-learn Version: Gaussian process regression on noise-free dataset



### 11.1.4 The spotPython Version

- The spotPython version is very similar:
  1. Instantiating the model, then
  2. fitting the model and
  3. making predictions (using predict).

```
S = Kriging(name='kriging', seed=123, log_level=50, cod_type="norm")
S.fit(X_train, y_train)
S_mean_prediction, S_std_prediction, S_ei = S.predict(X, return_val="all")
```

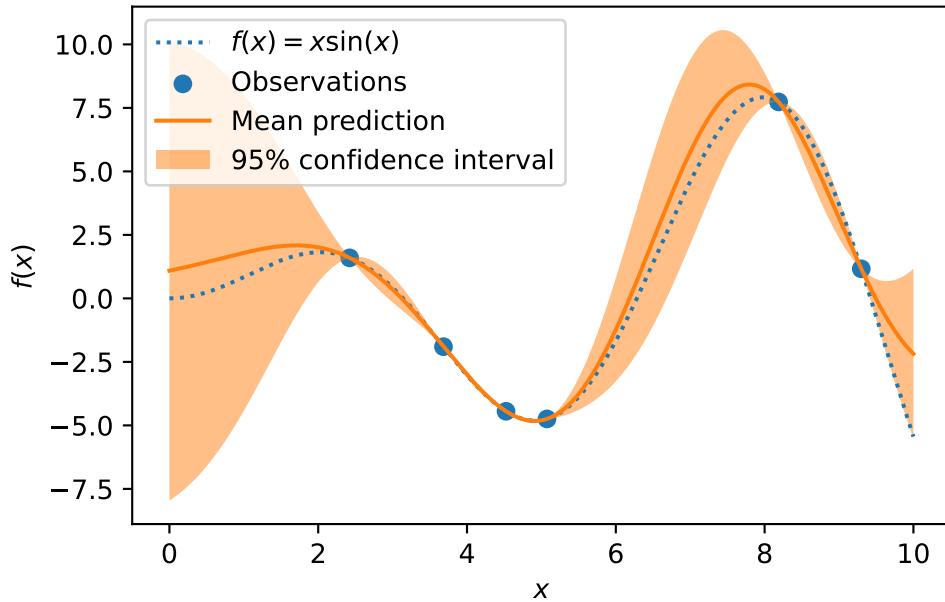
```
plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, S_mean_prediction, label="Mean prediction")
plt.fill_between(
    X.ravel(),
    S_mean_prediction - 1.96 * S_std_prediction,
    S_mean_prediction + 1.96 * S_std_prediction,
    alpha=0.5,
    label=r"95% confidence interval",
)
```

```

plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("spotPython Version: Gaussian process regression on noise-free dataset")

```

spotPython Version: Gaussian process regression on noise-free dataset

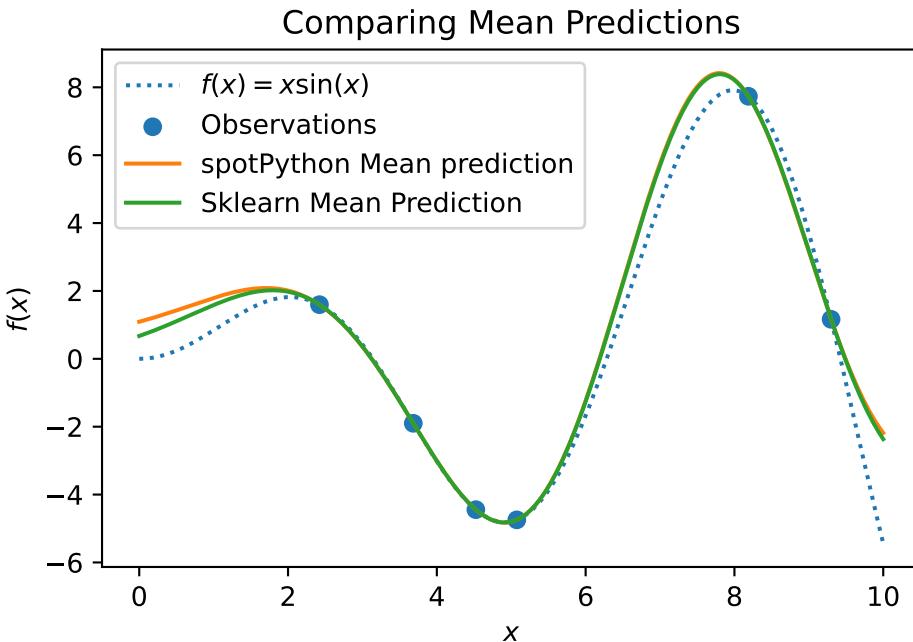


### 11.1.5 Visualizing the Differences Between the spotPython and the sklearn Model Fits

```

plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, S_mean_prediction, label="spotPython Mean prediction")
plt.plot(X, mean_prediction, label="Sklearn Mean Prediction")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Comparing Mean Predictions")

```



## 11.2 Exercises

### 11.2.1 Schonlau Example Function

- The Schonlau Example Function is based on sample points only (there is no analytical function description available):

```
X = np.linspace(start=0, stop=13, num=1_000).reshape(-1, 1)
X_train = np.array([1., 2., 3., 4., 12.]).reshape(-1,1)
y_train = np.array([0., -1.75, -2, -0.5, 5.])
```

- Describe the function.
- Compare the two models that were build using the `spotPython` and the `sklearn` surrogate.
- Note: Since there is no analytical function available, you might be interested in adding some points and describe the effects.

### 11.2.2 Forrester Example Function

- The Forrester Example Function is defined as follows:

$f(x) = (6x - 2)^2 \sin(12x - 4)$  for  $x$  in  $[0, 1]$ .

- Data points are generated as follows:

```
from spotPython.utils import fun_control_init
X = np.linspace(start=-0.5, stop=1.5, num=1_000).reshape(-1, 1)
X_train = np.array([0.0, 0.175, 0.225, 0.3, 0.35, 0.375, 0.5, 1]).reshape(-1, 1)
fun = analytical().fun_forrester
fun_control = fun_control_init(sigma = 0.1)
y = fun(X, fun_control=fun_control)
y_train = fun(X_train, fun_control=fun_control)
```

- Describe the function.
- Compare the two models that were build using the `spotPython` and the `sklearn` surrogate.
- Note: Modify the noise level ("sigma"), e.g., use a value of 0.2, and compare the two models.

```
fun_control = fun_control_init(sigma = 0.2)
```

### 11.2.3 `fun_runge` Function (1-dim)

- The Runge function is defined as follows:

$f(x) = 1 / (1 + \sum(x_i))^2$

- Data points are generated as follows:

```
gen = spacefilling(1)
rng = np.random.RandomState(1)
lower = np.array([-10])
upper = np.array([10])
fun = analytical().fun_runge
fun_control = fun_control_init(sigma = 0.025)
X_train = gen.scipy_lhd(10, lower=lower, upper = upper).reshape(-1, 1)
y_train = fun(X, fun_control=fun_control)
X = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
y = fun(X, fun_control=fun_control)
```

- Describe the function.
- Compare the two models that were build using the `spotPython` and the `sklearn` surrogate.

- Note: Modify the noise level ("sigma"), e.g., use a value of 0.05, and compare the two models.

```
fun_control = fun_control_init(sigma = 0.5)
```

#### 11.2.4 fun\_cubed (1-dim)

- The Cubed function is defined as follows:

```
np.sum(X[i]** 3)
```

- Data points are generated as follows:

```
gen = spacefilling(1)
rng = np.random.RandomState(1)
fun_control = fun_control_init(sigma = 0.025,
                                lower = np.array([-10]),
                                upper = np.array([10]))
fun = analytical().fun_cubed
X_train = gen.scipy_lhd(10, lower=lower, upper = upper).reshape(-1,1)
y_train = fun(X, fun_control=fun_control)
X = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
y = fun(X, fun_control=fun_control)
```

- Describe the function.
- Compare the two models that were build using the `spotPython` and the `sklearn` surrogate.
- Note: Modify the noise level ("sigma"), e.g., use a value of 0.05, and compare the two models.

```
fun_control = fun_control_init(sigma = 0.025)
```

#### 11.2.5 The Effect of Noise

How does the behavior of the `spotPython` fit changes when the argument `noise` is set to True, i.e.,

```
S = Kriging(name='kriging', seed=123, n_theta=1, noise=True)
```

is used?

# 12 Expected Improvement

This chapter describes, analyzes, and compares different infill criterion. An infill criterion defines how the next point  $x_{n+1}$  is selected from the surrogate model  $S$ . Expected improvement is a popular infill criterion in Bayesian optimization.

## 12.1 Example: Spot and the 1-dim Sphere Function

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from spotPython.utils.init import fun_control_init, surrogate_control_init, design_control_init
import matplotlib.pyplot as plt
```

### 12.1.1 The Objective Function: 1-dim Sphere

- The `spotPython` package provides several classes of objective functions.
- We will use an analytical objective function, i.e., a function that can be described by a (closed) formula:

$$f(x) = x^2$$

```
fun = analytical().fun_sphere
```

- The size of the `lower` bound vector determines the problem dimension.
- Here we will use `np.array([-1])`, i.e., a one-dim function.

#### i TensorBoard

Similar to the one-dimensional case, which was introduced in Section Section 7.5, we can use TensorBoard to monitor the progress of the optimization. We will use the same code, only the prefix is different:

```

from spotPython.utils.init import fun_control_init
PREFIX = "07_Y"
fun_control = fun_control_init(
    PREFIX=PREFIX,
    fun_evals = 25,
    lower = np.array([-1]),
    upper = np.array([1]),
    tolerance_x = np.sqrt(np.spacing(1)),)
design_control = design_control_init(init_size=10)

```

Created spot\_tensorboard\_path: runs/spot\_logs/07\_Y\_p040025\_2024-01-09\_20-17-02 for SummaryWriter

```

spot_1 = spot.Spot(
    fun=fun,
    fun_control=fun_control,
    design_control=design_control)
spot_1.run()

```

```

spotPython tuning: 1.1986325668379847e-08 [#####-----] 44.00%
spotPython tuning: 1.1986325668379847e-08 [#####-----] 48.00%
spotPython tuning: 1.1986325668379847e-08 [#####-----] 52.00%
spotPython tuning: 1.1385311249270152e-08 [#####----] 56.00%
spotPython tuning: 5.189351944607845e-10 [#####----] 60.00%
spotPython tuning: 5.189351944607845e-10 [#####----] 64.00%
spotPython tuning: 5.189351944607845e-10 [#####---] 68.00%
spotPython tuning: 5.189351944607845e-10 [#####---] 72.00%
spotPython tuning: 5.189351944607845e-10 [#####---] 76.00%
spotPython tuning: 5.189351944607845e-10 [#####---] 80.00%
spotPython tuning: 5.189351944607845e-10 [#####---] 84.00%
spotPython tuning: 5.189351944607845e-10 [#####---] 88.00%
spotPython tuning: 4.9233547910699585e-11 [#####---] 92.00%
spotPython tuning: 4.9233547910699585e-11 [#####---] 96.00%
spotPython tuning: 4.9233547910699585e-11 [#####---] 100.00% Done...

```

<spotPython.spot.spot.Spot at 0x2ce04cad0>

### 12.1.2 Results

```
spot_1.print_results()
```

```
min y: 4.9233547910699585e-11
x0: -7.016662163072951e-06
```

```
[['x0', -7.016662163072951e-06]]
```

```
spot_1.plot_progress(log_y=True)
```

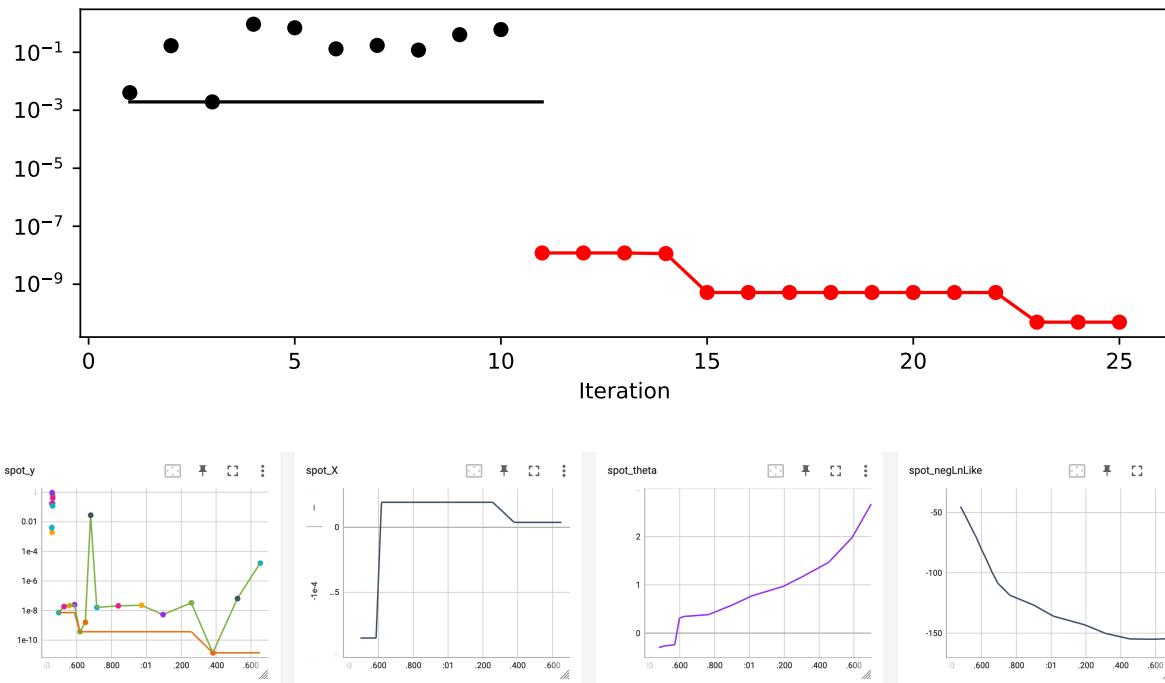


Figure 12.1: TensorBoard visualization of the spotPython optimization process and the surrogate model.

## 12.2 Same, but with EI as infill\_criterion

```
PREFIX = "07_EI_ISO"
fun_control = fun_control_init(
    PREFIX=PREFIX,
```

```
lower = np.array([-1]),  
upper = np.array([1]),  
fun_evals = 25,  
tolerance_x = np.sqrt(np.spacing(1)),  
infill_criterion = "ei")
```

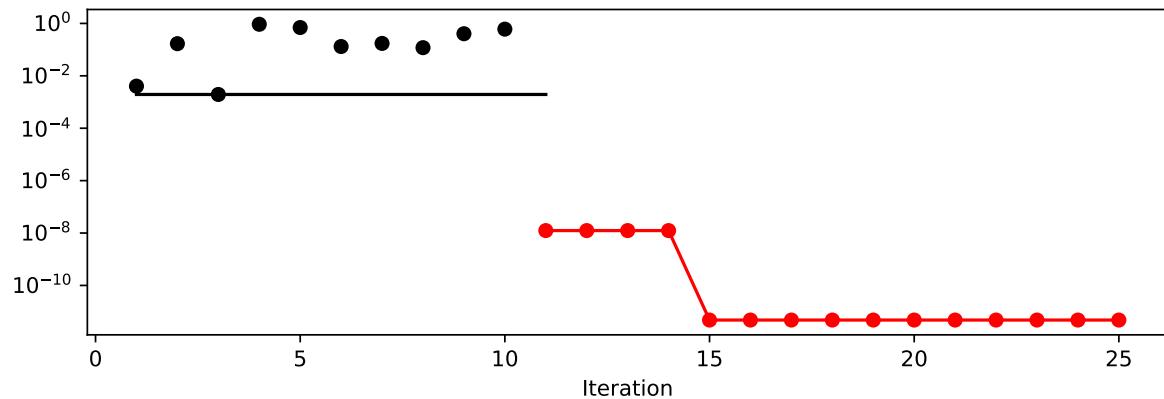
Created spot\_tensorboard\_path: runs/spot\_logs/07\_EI\_IS0\_p040025\_2024-01-09\_20-17-15 for Summary

```
spot_1_ei = spot.Spot(fun=fun,  
                      fun_control=fun_control)  
spot_1_ei.run()
```

```
spotPython tuning: 1.2401208585321248e-08 [#####-----] 44.00%  
spotPython tuning: 1.2401208585321248e-08 [#####-----] 48.00%  
spotPython tuning: 1.2401208585321248e-08 [#####-----] 52.00%  
spotPython tuning: 1.2401208585321248e-08 [#####-----] 56.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 60.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 64.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 68.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 72.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 76.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 80.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 84.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 88.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 92.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 96.00%  
spotPython tuning: 4.773576053400193e-12 [#####-----] 100.00% Done...
```

```
<spotPython.spot.spot.Spot at 0x2cead8c90>
```

```
spot_1_ei.plot_progress(log_y=True)
```



```
spot_1_ei.print_results()
```

```
min y: 4.773576053400193e-12
x0: -2.184851494587262e-06
```

```
[['x0', -2.184851494587262e-06]]
```

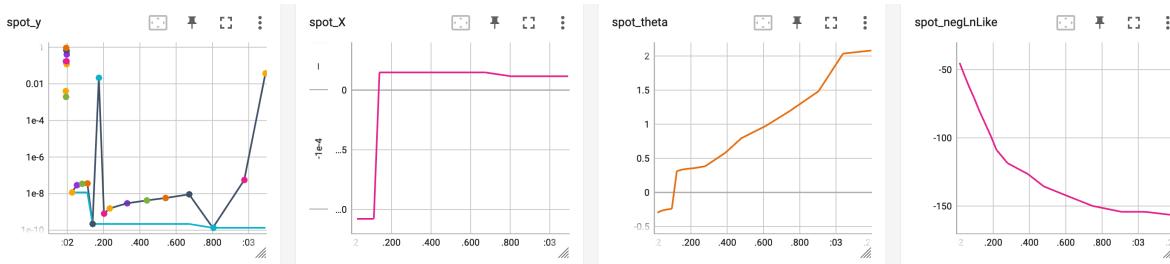


Figure 12.2: TensorBoard visualization of the spotPython optimization process and the surrogate model. Expected improvement, isotropic Kriging.

## 12.3 Non-isotropic Kriging

```
PREFIX = "07_EI_NONISO"
fun_control = fun_control_init(
    PREFIX=PREFIX,
    lower = np.array([-1, -1]),
    upper = np.array([1, 1]),
```

```
    fun_evals = 25,
    tolerance_x = np.sqrt(np.spacing(1)),
    infill_criterion = "ei")
surrogate_control = surrogate_control_init(
    n_theta=2,
    noise=False,
)
```

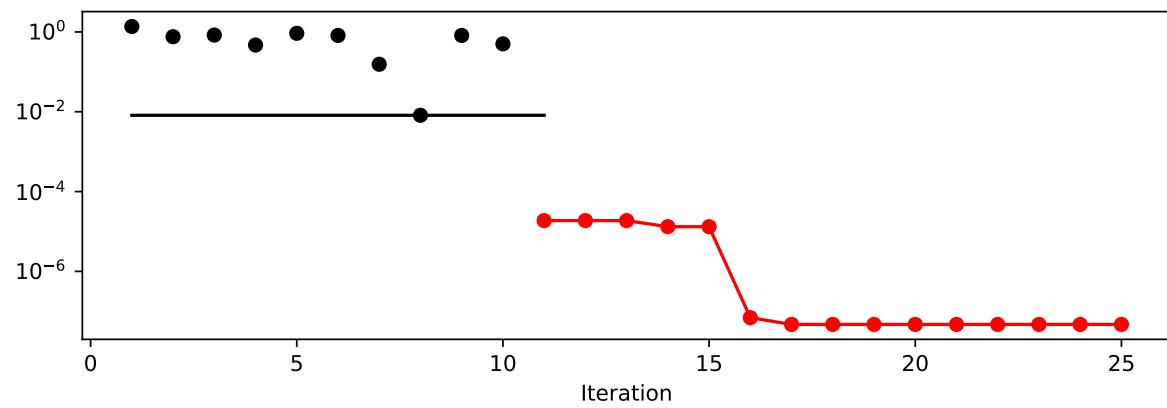
```
Created spot_tensorboard_path: runs/spot_logs/07_EI_NONISO_p040025_2024-01-09_20-17-27 for S
```

```
spot_2_ei_noniso = spot.Spot(fun=fun,
                               fun_control=fun_control,
                               surrogate_control=surrogate_control)
spot_2_ei_noniso.run()
```

```
spotPython tuning: 1.8691113149381826e-05 [#####-----] 44.00%
spotPython tuning: 1.8691113149381826e-05 [#####-----] 48.00%
spotPython tuning: 1.8691113149381826e-05 [#####-----] 52.00%
spotPython tuning: 1.3171284020530374e-05 [#####----] 56.00%
spotPython tuning: 1.3171284020530374e-05 [#####----] 60.00%
spotPython tuning: 6.912401156134061e-08 [#####----] 64.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 68.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 72.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 76.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 80.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 84.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 88.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 92.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 96.00%
spotPython tuning: 4.6808832368953396e-08 [#####---] 100.00% Done...
```

```
<spotPython.spot.spot.Spot at 0x2ce24ed90>
```

```
spot_2_ei_noniso.plot_progress(log_y=True)
```



```
spot_2_ei_noniso.print_results()
```

```
min y: 4.6808832368953396e-08
x0: -0.00020937062789341742
x1: 5.452313769831611e-05
```

```
[['x0', -0.00020937062789341742], ['x1', 5.452313769831611e-05]]
```

```
spot_2_ei_noniso.surrogate.plot()
```

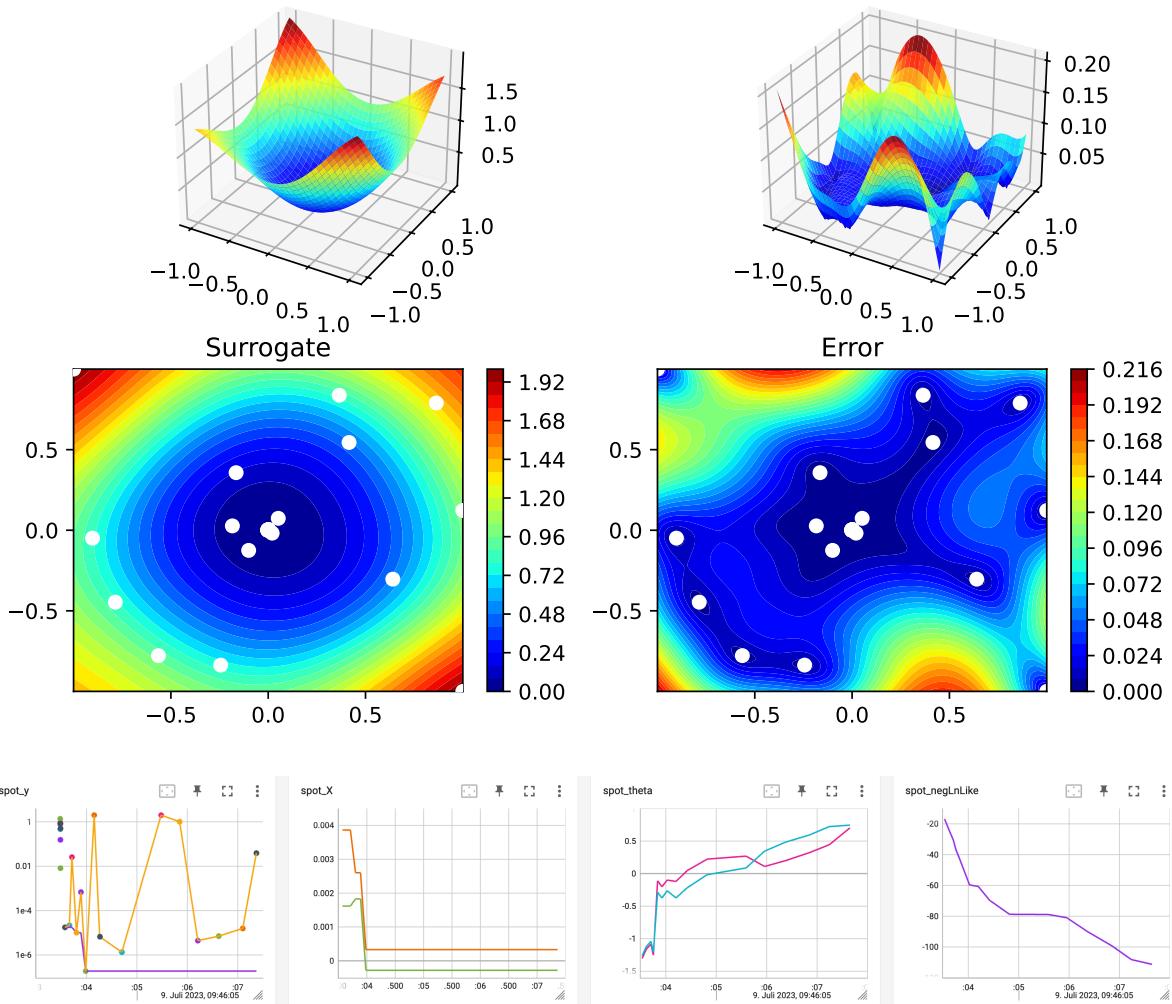


Figure 12.3: TensorBoard visualization of the spotPython optimization process and the surrogate model. Expected improvement, isotropic Kriging.

## 12.4 Using sklearn Surrogates

### 12.4.1 The spot Loop

The `spot` loop consists of the following steps:

1. Init: Build initial design  $X$
2. Evaluate initial design on real objective  $f$ :  $y = f(X)$
3. Build surrogate:  $S = S(X, y)$

4. Optimize on surrogate:  $X_0 = \text{optimize}(S)$
5. Evaluate on real objective:  $y_0 = f(X_0)$
6. Impute (Infill) new points:  $X = X \cup X_0$ ,  $y = y \cup y_0$ .
7. Got 3.

The spot loop is implemented in R as follows:

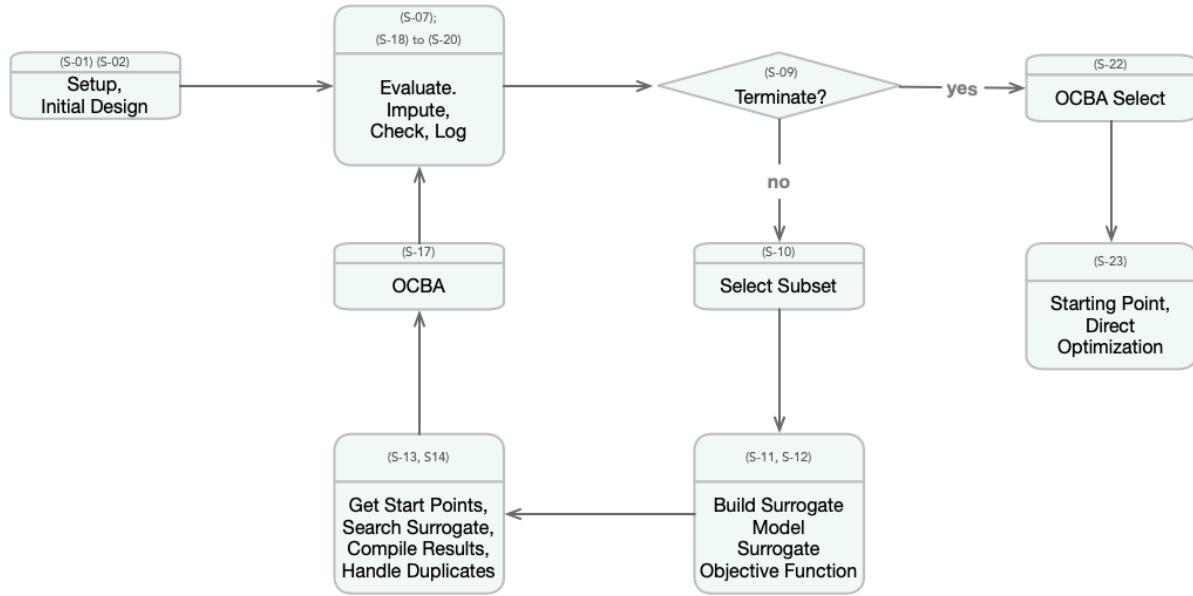


Figure 12.4: Visual representation of the model based search with SPOT. Taken from: Bartz-Beielstein, T., and Zaefferer, M. Hyperparameter tuning approaches. In Hyperparameter Tuning for Machine and Deep Learning with R - A Practical Guide, E. Bartz, T. Bartz-Beielstein, M. Zaefferer, and O. Mersmann, Eds. Springer, 2022, ch. 4, pp. 67–114.

## 12.4.2 spot: The Initial Model

### 12.4.2.1 Example: Modifying the initial design size

This is the “Example: Modifying the initial design size” from Chapter 4.5.1 in [bart21i].

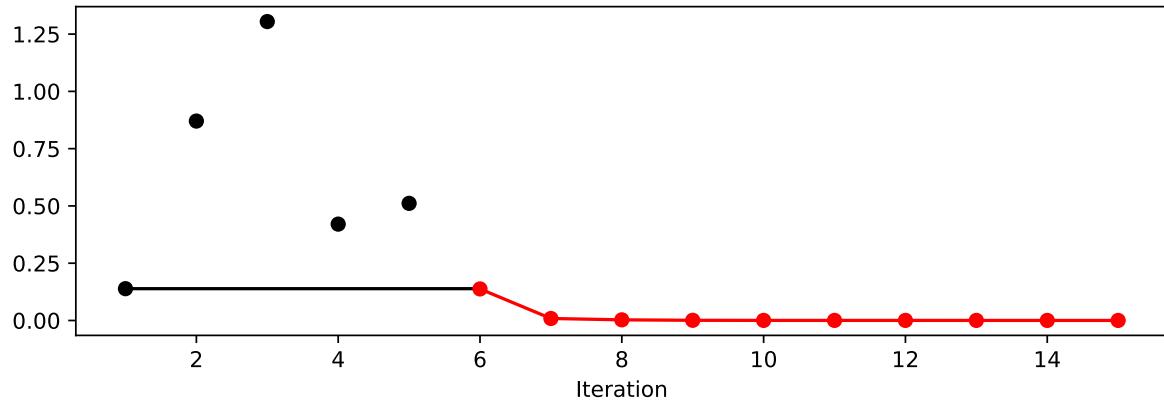
```

spot_ei = spot.Spot(fun=fun,
                     fun_control=fun_control_init(
                     lower = np.array([-1,-1]),
                     upper= np.array([1,1])),
                     design_control = design_control_init(init_size=5))
spot_ei.run()
  
```

```
spotPython tuning: 0.13771718056980337 [#####-----] 40.00%
spotPython tuning: 0.00873860889878173 [#####-----] 46.67%
spotPython tuning: 0.0028322099789017184 [#####-----] 53.33%
spotPython tuning: 0.0008112636408310097 [#####----] 60.00%
spotPython tuning: 0.00036558411628539707 [#####----] 66.67%
spotPython tuning: 0.00035830687535911014 [#####----] 73.33%
spotPython tuning: 0.00035830687535911014 [#####---] 80.00%
spotPython tuning: 0.00032833591068864664 [#####---] 86.67%
spotPython tuning: 0.0002757670364359785 [#####---] 93.33%
spotPython tuning: 0.0001601071101026825 [#####---] 100.00% Done...
```

```
<spotPython.spot.spot at 0x2ceb2a650>
```

```
spot_ei.plot_progress()
```



```
np.min(spot_1.y), np.min(spot_ei.y)
```

```
(4.9233547910699585e-11, 0.0001601071101026825)
```

### 12.4.3 Init: Build Initial Design

```
from spotPython.design.spacefilling import spacefilling
from spotPython.build.kriging import Kriging
from spotPython.fun.objectivefunctions import analytical
gen = spacefilling(2)
rng = np.random.RandomState(1)
```

```

lower = np.array([-5, -0])
upper = np.array([10,15])
fun = analytical().fun_branin

X = gen.scipy_lhd(10, lower=lower, upper = upper)
print(X)
y = fun(X, fun_control=fun_control)
print(y)

```

```

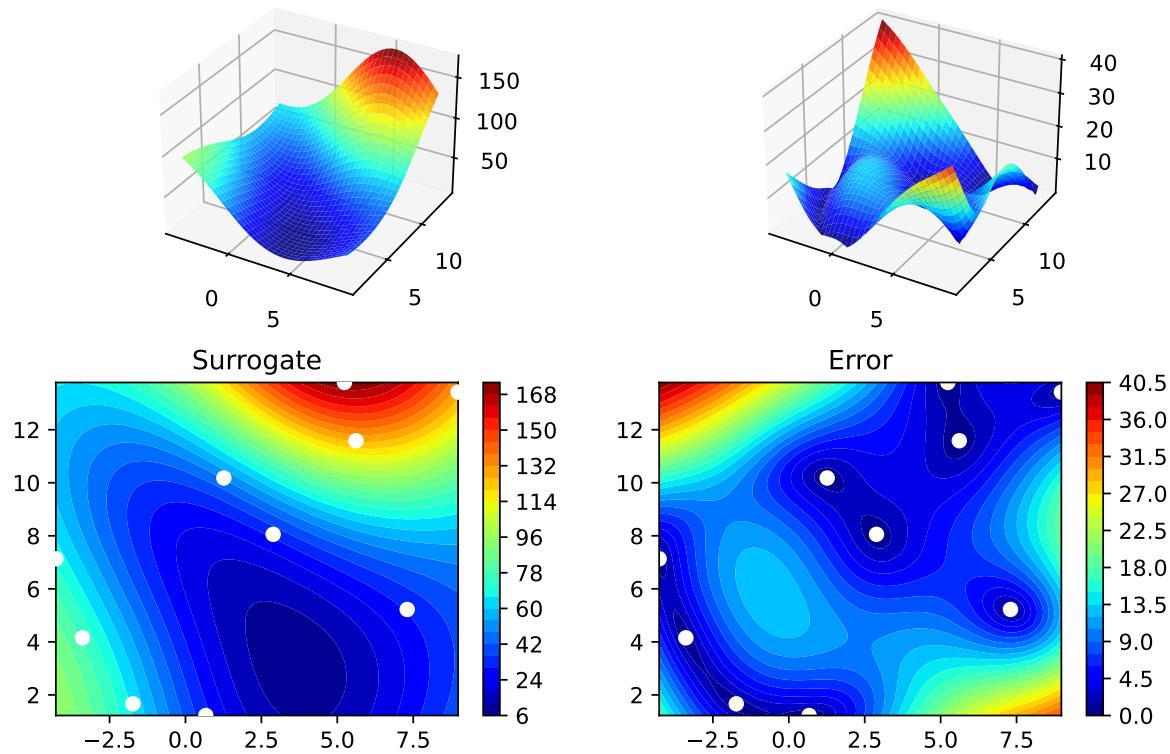
[[ 8.97647221 13.41926847]
 [ 0.66946019  1.22344228]
 [ 5.23614115 13.78185824]
 [ 5.6149825  11.5851384 ]
 [-1.72963184  1.66516096]
 [-4.26945568  7.1325531 ]
 [ 1.26363761 10.17935555]
 [ 2.88779942  8.05508969]
 [-3.39111089  4.15213772]
 [ 7.30131231  5.22275244]]
[128.95676449 31.73474356 172.89678121 126.71295908 64.34349975
 70.16178611 48.71407916 31.77322887 76.91788181 30.69410529]

```

```

S = Kriging(name='kriging', seed=123)
S.fit(X, y)
S.plot()

```



```
gen = spacefilling(2, seed=123)
X0 = gen.scipy_lhd(3)
gen = spacefilling(2, seed=345)
X1 = gen.scipy_lhd(3)
X2 = gen.scipy_lhd(3)
gen = spacefilling(2, seed=123)
X3 = gen.scipy_lhd(3)
X0, X1, X2, X3
```

```
(array([[0.77254938, 0.31539299],
       [0.59321338, 0.93854273],
       [0.27469803, 0.3959685 ]]),
 array([[0.78373509, 0.86811887],
       [0.06692621, 0.6058029 ],
       [0.41374778, 0.00525456]]),
 array([[0.121357 , 0.69043832],
       [0.41906219, 0.32838498],
       [0.86742658, 0.52910374]]),
 array([[0.77254938, 0.31539299],
```

```
[0.59321338, 0.93854273],  
[0.27469803, 0.3959685 ])))
```

#### 12.4.4 Evaluate

#### 12.4.5 Build Surrogate

#### 12.4.6 A Simple Predictor

The code below shows how to use a simple model for prediction.

- Assume that only two (very costly) measurements are available:
  1.  $f(0) = 0.5$
  2.  $f(2) = 2.5$
- We are interested in the value at  $x_0 = 1$ , i.e.,  $f(x_0 = 1)$ , but cannot run an additional, third experiment.

```
from sklearn import linear_model  
X = np.array([[0], [2]])  
y = np.array([0.5, 2.5])  
S_lm = linear_model.LinearRegression()  
S_lm = S_lm.fit(X, y)  
X0 = np.array([[1]])  
y0 = S_lm.predict(X0)  
print(y0)
```

[1.5]

- Central Idea:
  - Evaluation of the surrogate model  $S_{lm}$  is much cheaper (or / and much faster) than running the real-world experiment  $f$ .

### 12.5 Gaussian Processes regression: basic introductory example

This example was taken from [scikit-learn](#). After fitting our model, we see that the hyperparameters of the kernel have been optimized. Now, we will use our kernel to compute the mean prediction of the full dataset and plot the 95% confidence interval.

```

import numpy as np
import matplotlib.pyplot as plt
import math as m
from sklearn.gaussian_process import GaussianProcessRegressor
from sklearn.gaussian_process.kernels import RBF

X = np.linspace(start=0, stop=10, num=1_000).reshape(-1, 1)
y = np.squeeze(X * np.sin(X))
rng = np.random.RandomState(1)
training_indices = rng.choice(np.arange(y.size), size=6, replace=False)
X_train, y_train = X[training_indices], y[training_indices]

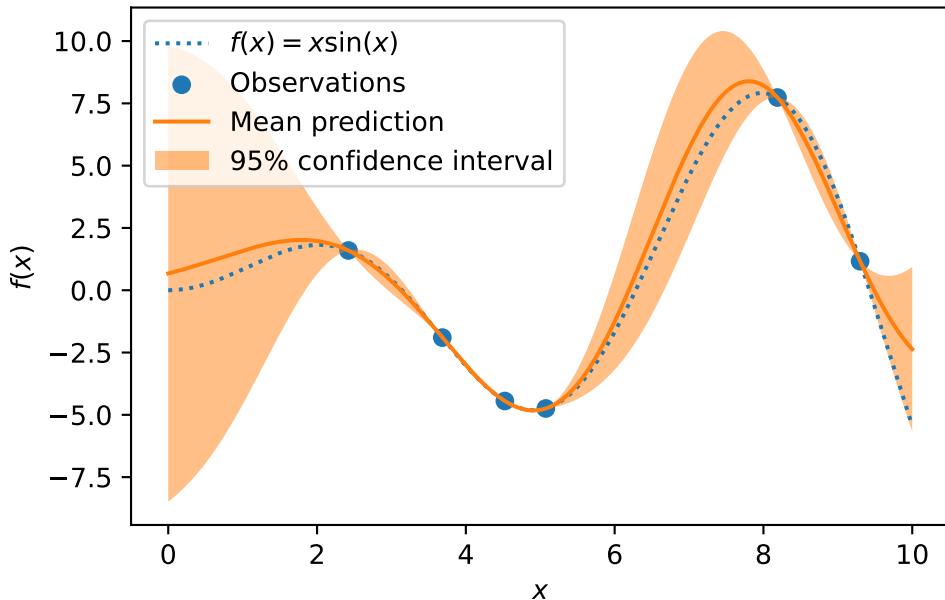
kernel = 1 * RBF(length_scale=1.0, length_scale_bounds=(1e-2, 1e2))
gaussian_process = GaussianProcessRegressor(kernel=kernel, n_restarts_optimizer=9)
gaussian_process.fit(X_train, y_train)
gaussian_process.kernel_

mean_prediction, std_prediction = gaussian_process.predict(X, return_std=True)

plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, mean_prediction, label="Mean prediction")
plt.fill_between(
    X.ravel(),
    mean_prediction - 1.96 * std_prediction,
    mean_prediction + 1.96 * std_prediction,
    alpha=0.5,
    label=r"95% confidence interval",
)
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("sk-learn Version: Gaussian process regression on noise-free dataset")

```

## sk-learn Version: Gaussian process regression on noise-free dataset



```

from spotPython.build.kriging import Kriging
import numpy as np
import matplotlib.pyplot as plt
rng = np.random.RandomState(1)
X = np.linspace(start=0, stop=10, num=1_000).reshape(-1, 1)
y = np.squeeze(X * np.sin(X))
training_indices = rng.choice(np.arange(y.size), size=6, replace=False)
X_train, y_train = X[training_indices], y[training_indices]

S = Kriging(name='kriging', seed=123, log_level=50, cod_type="norm")
S.fit(X_train, y_train)

mean_prediction, std_prediction, ei = S.predict(X, return_val="all")

std_prediction

plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, mean_prediction, label="Mean prediction")
plt.fill_between(
    X.ravel(),
    mean_prediction - 1.96 * std_prediction,
    mean_prediction + 1.96 * std_prediction,
    color="orange",
    alpha=0.5)

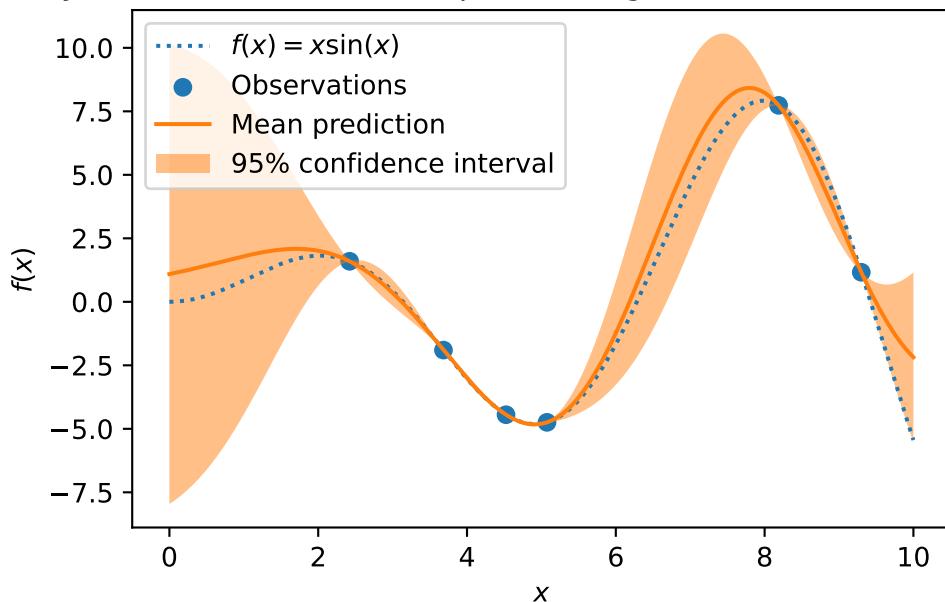
```

```

        mean_prediction + 1.96 * std_prediction,
        alpha=0.5,
        label=r"95% confidence interval",
    )
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("spotPython Version: Gaussian process regression on noise-free dataset")

```

spotPython Version: Gaussian process regression on noise-free dataset



## 12.6 The Surrogate: Using scikit-learn models

Default is the internal `kriging` surrogate.

```
S_0 = Kriging(name='kriging', seed=123)
```

Models from `scikit-learn` can be selected, e.g., Gaussian Process:

```
# Needed for the sklearn surrogates:
from sklearn.gaussian_process import GaussianProcessRegressor
from sklearn.gaussian_process.kernels import RBF
```

```
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn import linear_model
from sklearn import tree
import pandas as pd
```

```
kernel = 1 * RBF(length_scale=1.0, length_scale_bounds=(1e-2, 1e2))
S_GP = GaussianProcessRegressor(kernel=kernel, n_restarts_optimizer=9)
```

- and many more:

```
S_Tree = DecisionTreeRegressor(random_state=0)
S_LM = linear_model.LinearRegression()
S_Ridge = linear_model.Ridge()
S_RF = RandomForestRegressor(max_depth=2, random_state=0)
```

- The scikit-learn GP model S\_GP is selected.

```
S = S_GP
```

```
isinstance(S, GaussianProcessRegressor)
```

```
True
```

```
from spotPython.fun.objectivefunctions import analytical
fun = analytical().fun_branin
fun_control = fun_control_init(
    lower = np.array([-5,-0]),
    upper = np.array([10,15]),
    fun_evals = 15)
design_control = design_control_init(init_size=5)
spot_GP = spot.Spot(fun=fun,
                     fun_control=fun_control,
                     surrogate=S,
                     design_control=design_control)
spot_GP.run()
```

```
spotPython tuning: 24.51465459019188 [#####-----] 40.00%
spotPython tuning: 11.003077541587748 [#####-----] 46.67%
spotPython tuning: 11.003077541587748 [#####-----] 53.33%
```

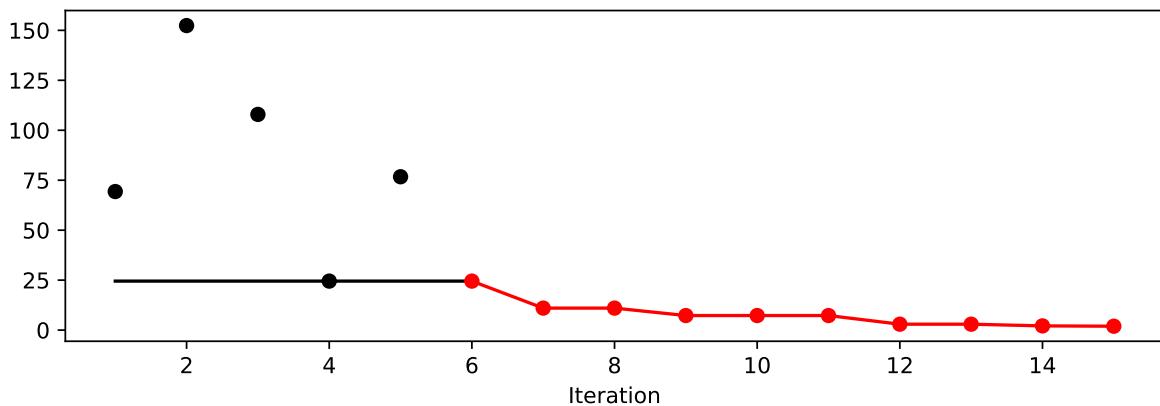
```
spotPython tuning: 7.281227279299504 [#####----] 60.00%
spotPython tuning: 7.281227279299504 [#####---] 66.67%
spotPython tuning: 7.281227279299504 [#####--] 73.33%
spotPython tuning: 2.9519489314482 [#####---] 80.00%
spotPython tuning: 2.9519489314482 [#####--] 86.67%
spotPython tuning: 2.104972804244822 [#####--] 93.33%
spotPython tuning: 1.9431600962086772 [#####--] 100.00% Done...
```

```
<spotPython.spot.spot at 0x2d7b66350>
```

```
spot_GP.y
```

```
array([ 69.32459936, 152.38491454, 107.92560483, 24.51465459,
       76.73500031, 86.30425303, 11.00307754, 16.11742138,
       7.28122728, 21.82317903, 10.96088904, 2.95194893,
       3.02910742, 2.1049728 , 1.9431601 ])
```

```
spot_GP.plot_progress()
```



```
spot_GP.print_results()
```

```
min y: 1.9431600962086772
x0: 10.0
x1: 2.9985482809555464
```

```
[['x0', 10.0], ['x1', 2.9985482809555464]]
```

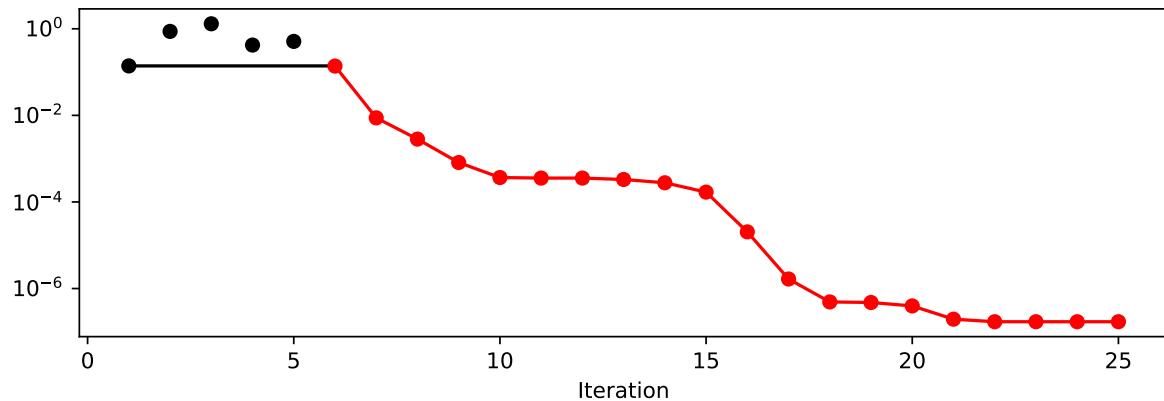
## 12.7 Additional Examples

```
# Needed for the sklearn surrogates:  
from sklearn.gaussian_process import GaussianProcessRegressor  
from sklearn.gaussian_process.kernels import RBF  
from sklearn.tree import DecisionTreeRegressor  
from sklearn.ensemble import RandomForestRegressor  
from sklearn import linear_model  
from sklearn import tree  
import pandas as pd  
  
kernel = 1 * RBF(length_scale=1.0, length_scale_bounds=(1e-2, 1e2))  
S_GP = GaussianProcessRegressor(kernel=kernel, n_restarts_optimizer=9)  
  
from spotPython.build.kriging import Kriging  
import numpy as np  
import spotPython  
from spotPython.fun.objectivefunctions import analytical  
from spotPython.spot import spot  
  
S_K = Kriging(name='kriging',  
               seed=123,  
               log_level=50,  
               infill_criterion = "y",  
               n_theta=1,  
               noise=False,  
               cod_type="norm")  
fun = analytical().fun_sphere  
  
fun_control = fun_control_init(  
    lower = np.array([-1,-1]),  
    upper = np.array([1,1]),  
    fun_evals = 25)  
  
spot_S_K = spot.Spot(fun=fun,  
                      fun_control=fun_control,  
                      surrogate=S_K,  
                      design_control=design_control,  
                      surrogate_control=surrogate_control)  
spot_S_K.run()
```

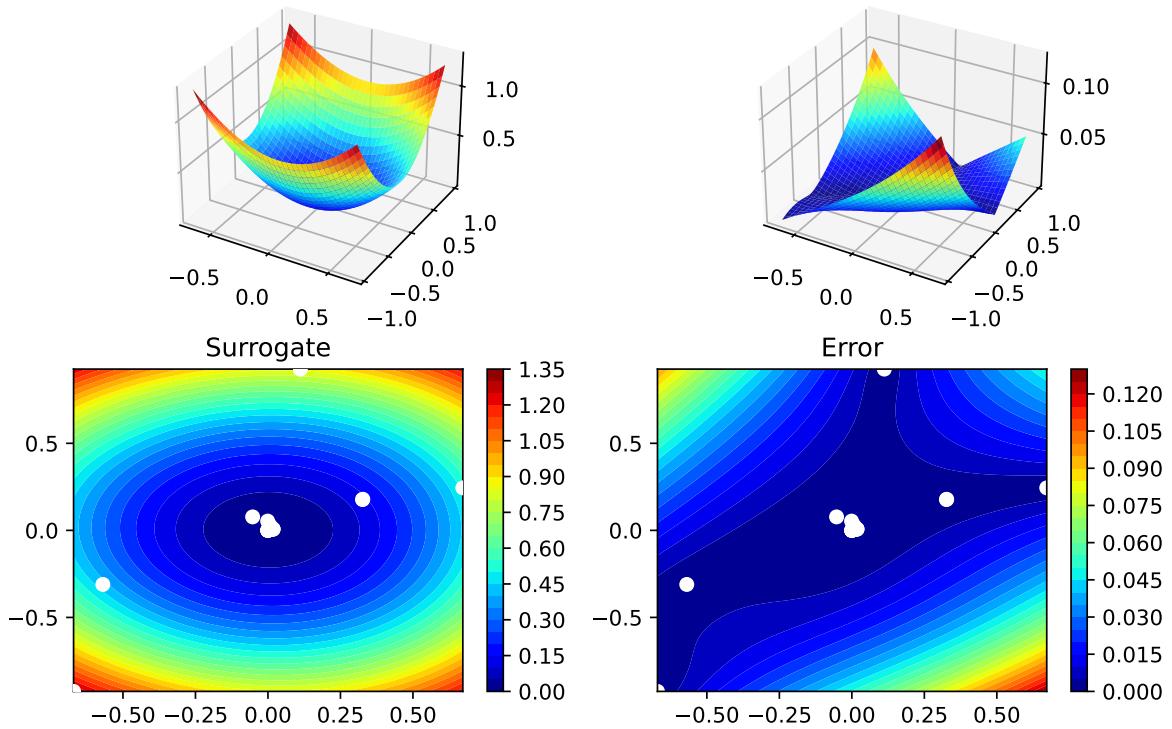
```
spotPython tuning: 0.13771718778810743 [##-----] 24.00%
spotPython tuning: 0.008768000187888899 [###-----] 28.00%
spotPython tuning: 0.0028300907437246053 [###-----] 32.00%
spotPython tuning: 0.0008148020998531609 [####-----] 36.00%
spotPython tuning: 0.00036681248440550095 [####-----] 40.00%
spotPython tuning: 0.00035607605553701025 [####-----] 44.00%
spotPython tuning: 0.00035607605553701025 [#####-----] 48.00%
spotPython tuning: 0.00033033596693814263 [#####-----] 52.00%
spotPython tuning: 0.0002774179969789593 [#####-----] 56.00%
spotPython tuning: 0.00016886412273302311 [#####-----] 60.00%
spotPython tuning: 2.0349536932144563e-05 [#####-----] 64.00%
spotPython tuning: 1.6621220007683266e-06 [#####-----] 68.00%
spotPython tuning: 4.905822935561126e-07 [#####-----] 72.00%
spotPython tuning: 4.7634545282279014e-07 [#####----] 76.00%
spotPython tuning: 3.966290585455581e-07 [#####----] 80.00%
spotPython tuning: 1.9602185212475464e-07 [#####----] 84.00%
spotPython tuning: 1.7115221726800905e-07 [#####----] 88.00%
spotPython tuning: 1.7115221726800905e-07 [#####----] 92.00%
spotPython tuning: 1.7115221726800905e-07 [#####----] 96.00%
spotPython tuning: 1.7115221726800905e-07 [#####----] 100.00% Done...
```

```
<spotPython.spot.spot at 0x2d80250d0>
```

```
spot_S_K.plot_progress(log_y=True)
```



```
spot_S_K.surrogate.plot()
```



```
spot_S_K.print_results()
```

```
min y: 1.7115221726800905e-07
x0: 0.0003105897139994429
x1: 0.0002732878460995902
```

```
[['x0', 0.0003105897139994429], ['x1', 0.0002732878460995902]]
```

### 12.7.1 Optimize on Surrogate

### 12.7.2 Evaluate on Real Objective

### 12.7.3 Impute / Infill new Points

## 12.8 Tests

```

import numpy as np
from spotPython.spot import spot
from spotPython.fun.objectivefunctions import analytical

fun_sphere = analytical().fun_sphere

fun_control = fun_control_init(
    lower=np.array([-1, -1]),
    upper=np.array([1, 1]),
    n_points = 2)
spot_1 = spot.Spot(
    fun=fun_sphere,
    fun_control=fun_control,
)

# (S-2) Initial Design:
spot_1.X = spot_1.design.scipy_lhd(
    spot_1.design_control["init_size"], lower=spot_1.lower, upper=spot_1.upper
)
print(spot_1.X)

# (S-3): Eval initial design:
spot_1.y = spot_1.fun(spot_1.X)
print(spot_1.y)

spot_1.fit_surrogate()
X0 = spot_1.suggest_new_X()
print(X0)
assert X0.size == spot_1.n_points * spot_1.k

```

```

[[ 0.86352963  0.7892358 ]
 [-0.24407197 -0.83687436]
 [ 0.36481882  0.8375811 ]
 [ 0.415331     0.54468512]
 [-0.56395091 -0.77797854]
 [-0.90259409 -0.04899292]
 [-0.16484832  0.35724741]
 [ 0.05170659  0.07401196]
 [-0.78548145 -0.44638164]
 [ 0.64017497 -0.30363301]]
[1.36857656  0.75992983  0.83463487  0.46918172  0.92329124  0.8170764

```

```

0.15480068 0.00815134 0.81623768 0.502017 ]
[[0.00165526 0.00410847]
 [0.00165685 0.0039177 ]]

```

## 12.9 EI: The Famous Schonlau Example

```

X_train0 = np.array([1, 2, 3, 4, 12]).reshape(-1,1)
X_train = np.linspace(start=0, stop=10, num=5).reshape(-1, 1)

from spotPython.build.kriging import Kriging
import numpy as np
import matplotlib.pyplot as plt

X_train = np.array([1., 2., 3., 4., 12.]).reshape(-1,1)
y_train = np.array([0., -1.75, -2, -0.5, 5.])

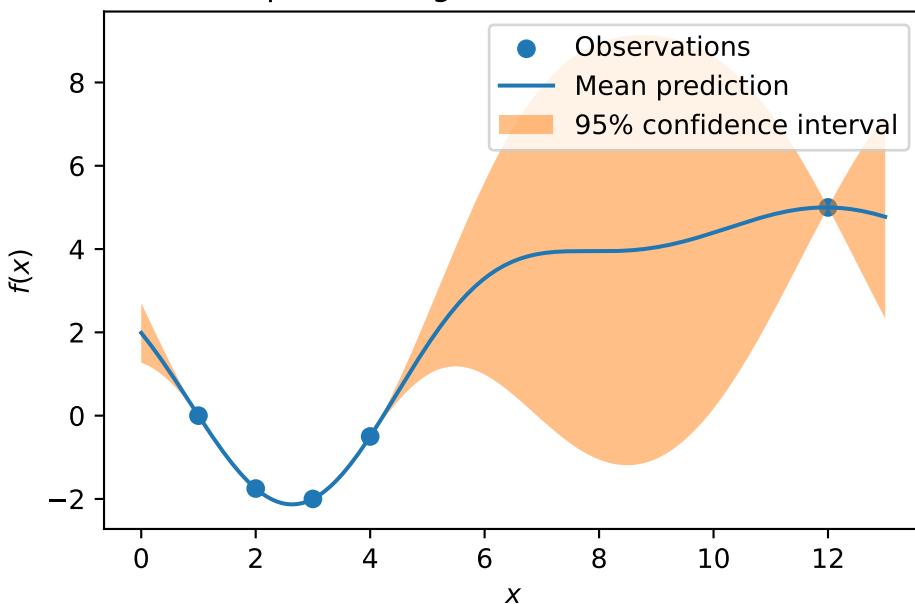
S = Kriging(name='kriging', seed=123, log_level=50, n_theta=1, noise=False, cod_type="norm")
S.fit(X_train, y_train)

X = np.linspace(start=0, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X, return_val="all")

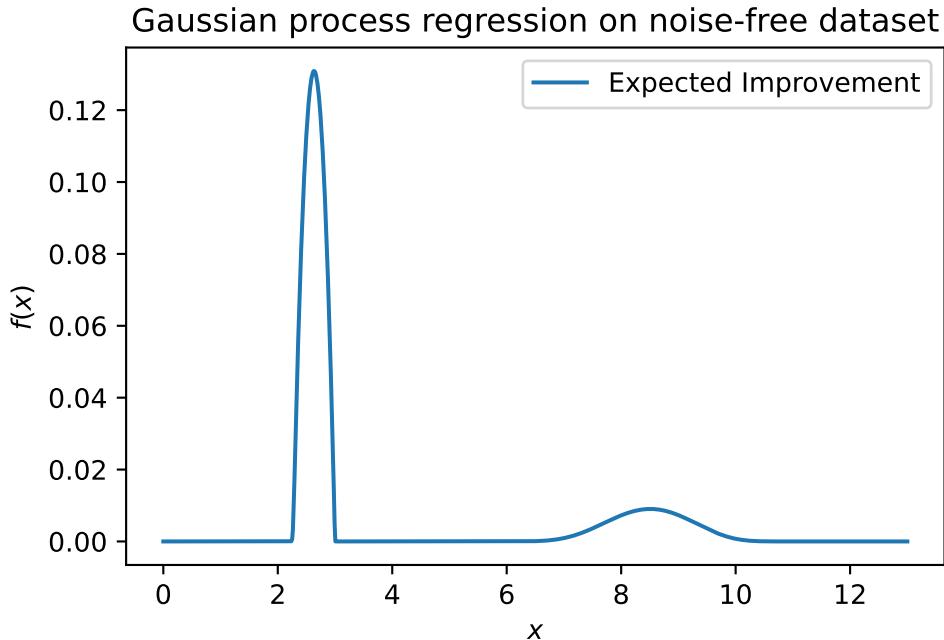
plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, mean_prediction, label="Mean prediction")
if True:
    plt.fill_between(
        X.ravel(),
        mean_prediction - 2 * std_prediction,
        mean_prediction + 2 * std_prediction,
        alpha=0.5,
        label=r"95% confidence interval",
    )
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Gaussian process regression on noise-free dataset")

```

### Gaussian process regression on noise-free dataset



```
#plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
# plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, -ei, label="Expected Improvement")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Gaussian process regression on noise-free dataset")
```



```
S.log
```

```
{'negLnLike': array([1.20788205]),
 'theta': array([-0.9900252]),
 'p': [],
 'Lambda': []}
```

## 12.10 EI: The Forrester Example

```
from spotPython.build.kriging import Kriging
import numpy as np
import matplotlib.pyplot as plt
import spotPython
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot

# exact x locations are unknown:
X_train = np.array([0.0, 0.175, 0.225, 0.3, 0.35, 0.375, 0.5, 1]).reshape(-1,1)

fun = analytical().fun_forrester
```

```

fun_control = fun_control_init(
    PREFIX="07_EI_FORRESTER",
    sigma=1.0,
    seed=123)
y_train = fun(X_train, fun_control=fun_control)

S = Kriging(name='kriging', seed=123, log_level=50, n_theta=1, noise=False, cod_type="norm")
S.fit(X_train, y_train)

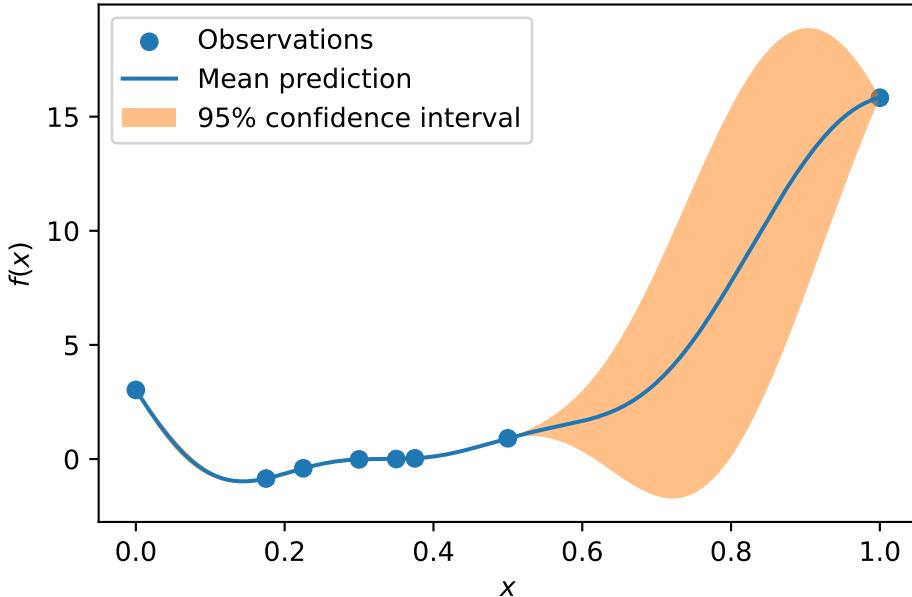
X = np.linspace(start=0, stop=1, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X, return_val="all")

plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, mean_prediction, label="Mean prediction")
if True:
    plt.fill_between(
        X.ravel(),
        mean_prediction - 2 * std_prediction,
        mean_prediction + 2 * std_prediction,
        alpha=0.5,
        label=r"95% confidence interval",
    )
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Gaussian process regression on noise-free dataset")

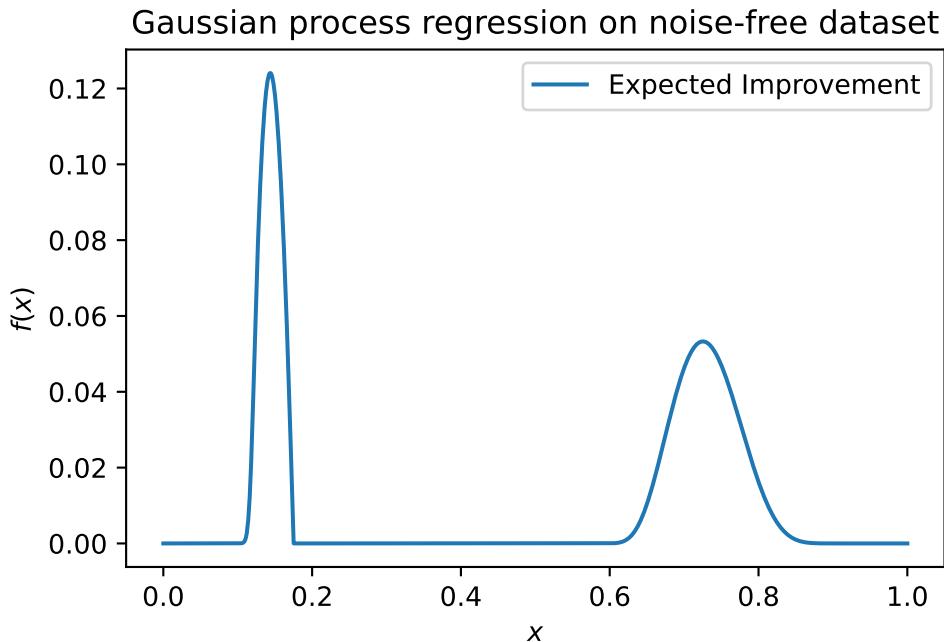
```

Created spot\_tensorboard\_path: runs/spot\_logs/07\_EI\_FORRESTER\_p040025\_2024-01-09\_20-19-28 for

### Gaussian process regression on noise-free dataset



```
#plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
# plt.scatter(X_train, y_train, label="Observations")
plt.plot(X, -ei, label="Expected Improvement")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Gaussian process regression on noise-free dataset")
```



## 12.11 Noise

```

import numpy as np
import spotPython
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from spotPython.design.spacefilling import spacefilling
from spotPython.build.kriging import Kriging
import matplotlib.pyplot as plt

gen = spacefilling(1)
rng = np.random.RandomState(1)
lower = np.array([-10])
upper = np.array([10])
fun = analytical().fun_sphere
fun_control = fun_control_init(
    PREFIX="07_Y",
    sigma=2.0,
    seed=123,)
X = gen.scipy_lhd(10, lower=lower, upper = upper)

```

```

print(X)
y = fun(X, fun_control=fun_control)
print(y)
y.shape
X_train = X.reshape(-1,1)
y_train = y

S = Kriging(name='kriging',
             seed=123,
             log_level=50,
             n_theta=1,
             noise=False)
S.fit(X_train, y_train)

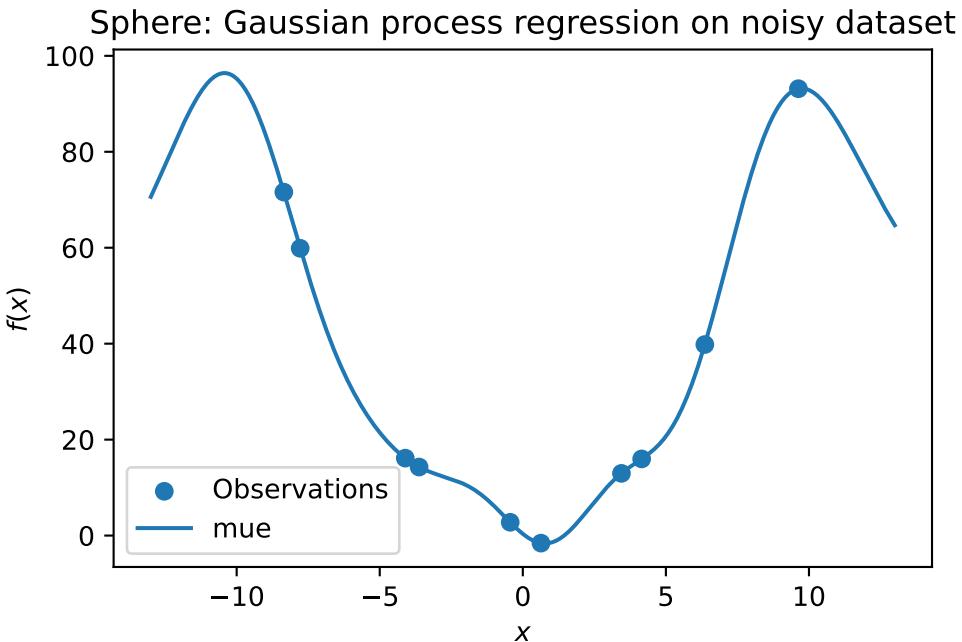
X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X_axis, return_val="all")

#plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
plt.scatter(X_train, y_train, label="Observations")
#plt.plot(X, ei, label="Expected Improvement")
plt.plot(X_axis, mean_prediction, label="mue")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Sphere: Gaussian process regression on noisy dataset")

```

```

Created spot_tensorboard_path: runs/spot_logs/07_Y_p040025_2024-01-09_20-19-28 for SummaryWriter
[[ 0.63529627]
 [-4.10764204]
 [-0.44071975]
 [ 9.63125638]
 [-8.3518118 ]
 [-3.62418901]
 [ 4.15331   ]
 [ 3.4468512 ]
 [ 6.36049088]
 [-7.77978539]
 [-1.57464135 16.13714981  2.77008442  93.14904827  71.59322218  14.28895359
 15.9770567 12.96468767 39.82265329 59.88028242]
```



```
S.log
```

```
{
  'negLnLike': array([26.18505386]),
  'theta': array([-1.10547474]),
  'p': [],
  'Lambda': []
}

S = Kriging(name='kriging',
             seed=123,
             log_level=50,
             n_theta=1,
             noise=True)
S.fit(X_train, y_train)

X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X_axis, return_val="all")

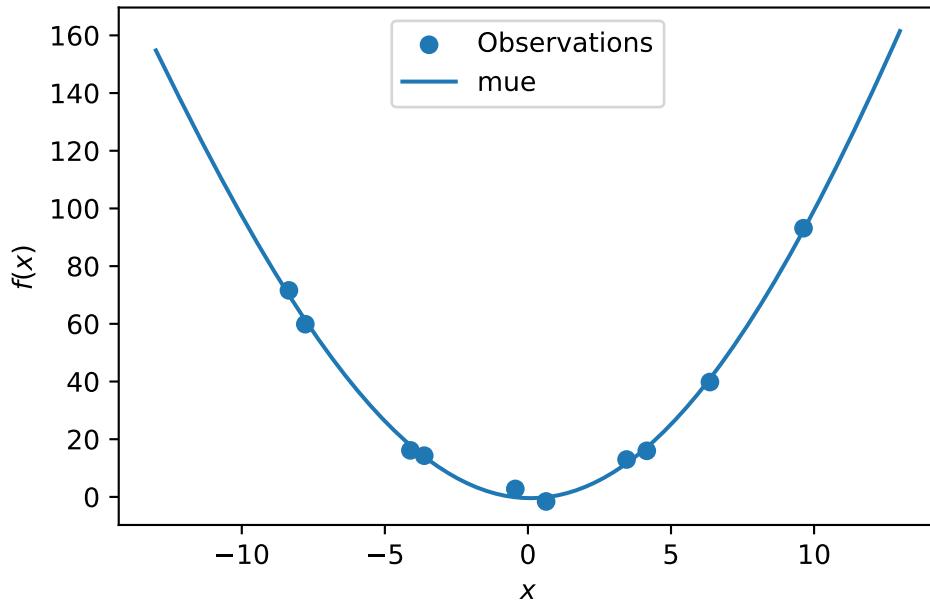
# plt.plot(X, y, label=r"$f(x) = x \sin(x)$", linestyle="dotted")
plt.scatter(X_train, y_train, label="Observations")
# plt.plot(X, ei, label="Expected Improvement")
plt.plot(X_axis, mean_prediction, label="mle")
plt.legend()
```

```

plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Sphere: Gaussian process regression with nugget on noisy dataset")

```

Sphere: Gaussian process regression with nugget on noisy dataset



S.log

```

{'negLnLike': array([21.82059174]),
 'theta': array([-2.96946062]),
 'p': [],
 'Lambda': array([4.28985898e-05])}

```

## 12.12 Cubic Function

```

import numpy as np
import spotPython
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from spotPython.design.spacefilling import spacefilling
from spotPython.build.kriging import Kriging

```

```

import matplotlib.pyplot as plt

gen = spacefilling(1)
rng = np.random.RandomState(1)
lower = np.array([-10])
upper = np.array([10])
fun = analytical().fun_cubed
fun_control = fun_control_init(
    PREFIX="07_Y",
    sigma=10.0,
    seed=123,)

X = gen.scipy_lhd(10, lower=lower, upper = upper)
print(X)
y = fun(X, fun_control=fun_control)
print(y)
y.shape
X_train = X.reshape(-1,1)
y_train = y

S = Kriging(name='kriging', seed=123, log_level=50, n_theta=1, noise=False)
S.fit(X_train, y_train)

X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X_axis, return_val="all")

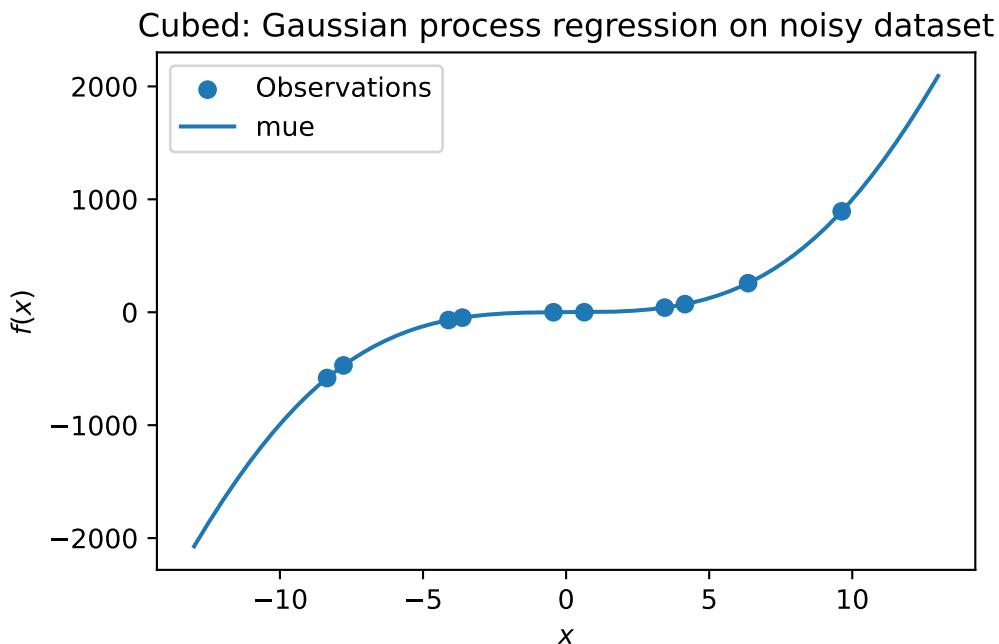
plt.scatter(X_train, y_train, label="Observations")
#plt.plot(X, ei, label="Expected Improvement")
plt.plot(X_axis, mean_prediction, label="mue")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Cubed: Gaussian process regression on noisy dataset")

```

```

Created spot_tensorboard_path: runs/spot_logs/07_Y_p040025_2024-01-09_20-19-28 for SummaryWriter
[[ 0.63529627]
 [-4.10764204]
 [-0.44071975]
 [ 9.63125638]
 [-8.3518118 ]
 [-3.62418901]
 [ 4.15331   ]
```

```
[ 3.4468512 ]
[ 6.36049088]
[-7.77978539]
[ 2.56406437e-01 -6.93071067e+01 -8.56027124e-02  8.93405931e+02
-5.82561927e+02 -4.76028022e+01  7.16445311e+01  4.09512920e+01
 2.57319028e+02 -4.70871982e+02]
```

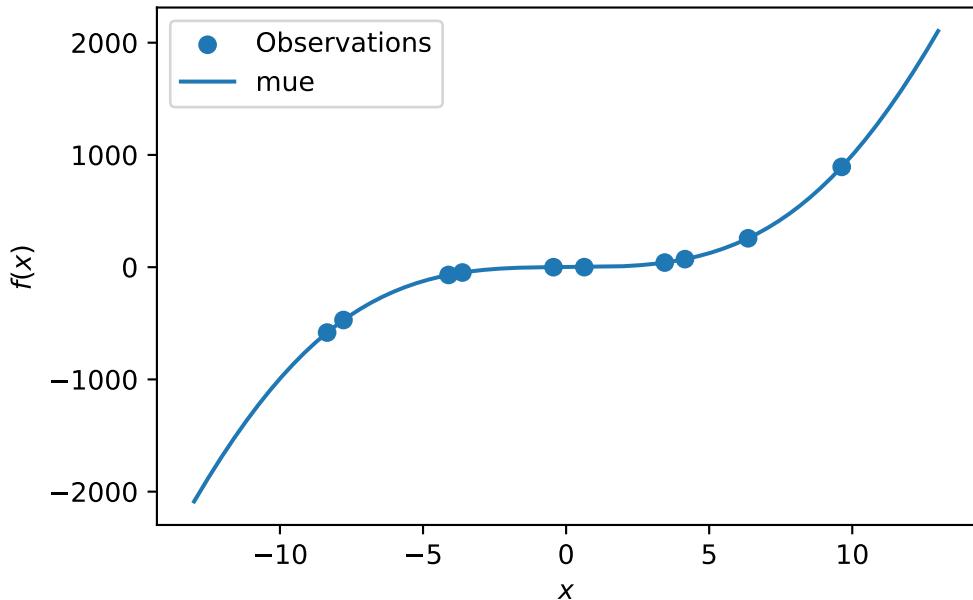


```
S = Kriging(name='kriging', seed=123, log_level=0, n_theta=1, noise=True)
S.fit(X_train, y_train)

X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X_axis, return_val="all")

plt.scatter(X_train, y_train, label="Observations")
# plt.plot(X, ei, label="Expected Improvement")
plt.plot(X_axis, mean_prediction, label="mue")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Cubed: Gaussian process with nugget regression on noisy dataset")
```

Cubed: Gaussian process with nugget regression on noisy dataset



```

import numpy as np
import spotPython
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from spotPython.design.spacefilling import spacefilling
from spotPython.build.kriging import Kriging
import matplotlib.pyplot as plt

gen = spacefilling(1)
rng = np.random.RandomState(1)
lower = np.array([-10])
upper = np.array([10])
fun = analytical().fun_runge
fun_control = fun_control_init(
    PREFIX="07_Y",
    sigma=0.25,
    seed=123,)

X = gen.scipy_lhd(10, lower=lower, upper = upper)
print(X)
y = fun(X, fun_control=fun_control)
print(y)
y.shape

```

```

X_train = X.reshape(-1,1)
y_train = y

S = Kriging(name='kriging', seed=123, log_level=50, n_theta=1, noise=False)
S.fit(X_train, y_train)

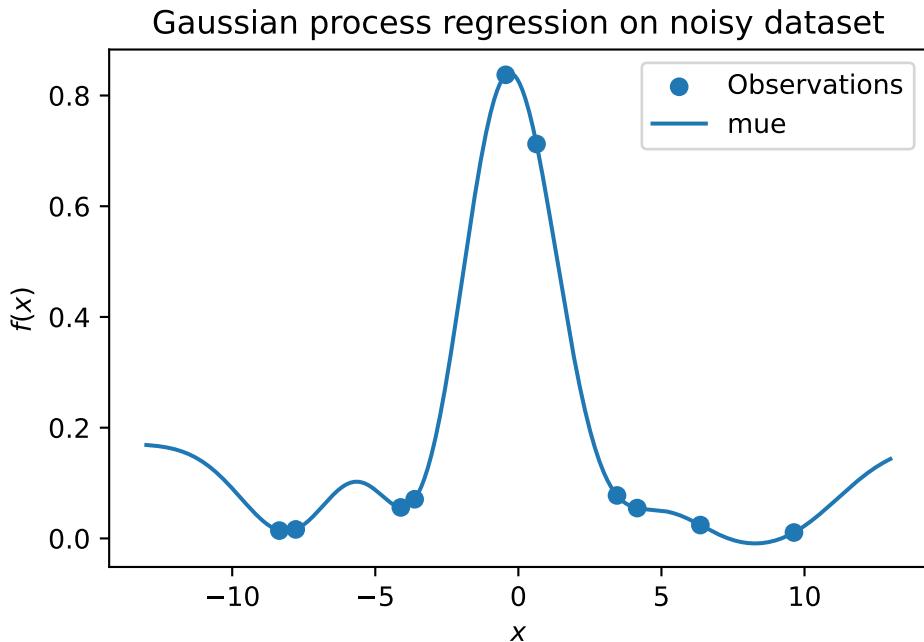
X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X_axis, return_val="all")

plt.scatter(X_train, y_train, label="Observations")
#plt.plot(X, ei, label="Expected Improvement")
plt.plot(X_axis, mean_prediction, label="mue")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Gaussian process regression on noisy dataset")

```

```

Created spot_tensorboard_path: runs/spot_logs/07_Y_p040025_2024-01-09_20-19-29 for SummaryWriter
[[ 0.63529627]
 [-4.10764204]
 [-0.44071975]
 [ 9.63125638]
 [-8.3518118 ]
 [-3.62418901]
 [ 4.15331   ]
 [ 3.4468512 ]
 [ 6.36049088]
 [-7.77978539]]
[0.712453  0.05595118 0.83735691 0.0106654  0.01413372 0.07074765
 0.05479457 0.07763503 0.02412205 0.01625354]
```

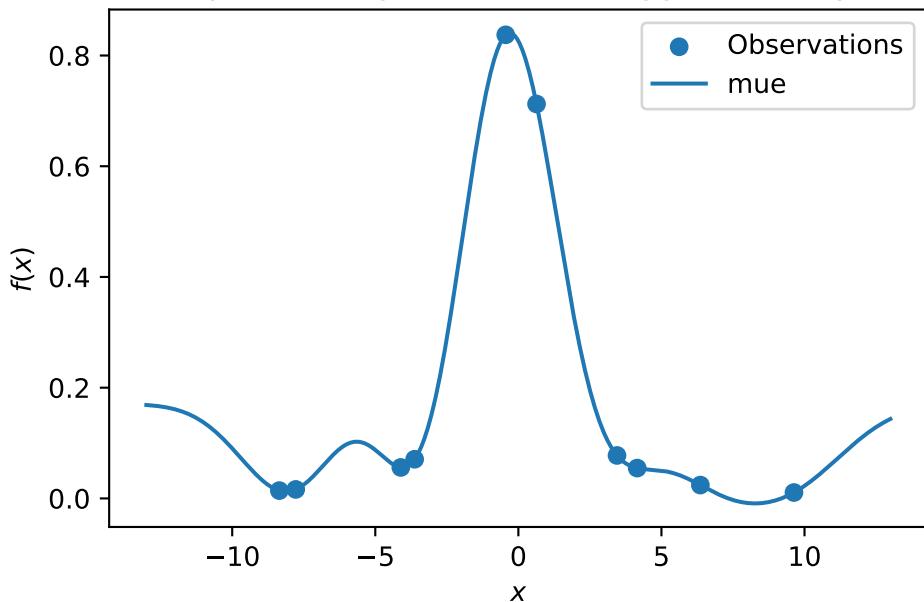


```
S = Kriging(name='kriging',
            seed=123,
            log_level=50,
            n_theta=1,
            noise=True)
S.fit(X_train, y_train)

X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X_axis, return_val="all")

plt.scatter(X_train, y_train, label="Observations")
# plt.plot(X, ei, label="Expected Improvement")
plt.plot(X_axis, mean_prediction, label="mle")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Gaussian process regression with nugget on noisy dataset")
```

Gaussian process regression with nugget on noisy dataset



## 12.13 Factors

```
["num"] * 3
```

```
['num', 'num', 'num']
```

```
from spotPython.design.spacefilling import spacefilling
from spotPython.build.kriging import Kriging
from spotPython.fun.objectivefunctions import analytical
import numpy as np
```

```
gen = spacefilling(2)
n = 30
rng = np.random.RandomState(1)
lower = np.array([-5,-0])
upper = np.array([10,15])
fun = analytical().fun_branin_factor
#fun = analytical(sigma=0).fun_sphere
```

```

X0 = gen.scipy_lhd(n, lower=lower, upper = upper)
X1 = np.random.randint(low=1, high=3, size=(n,))
X = np.c_[X0, X1]
y = fun(X)
S = Kriging(name='kriging', seed=123, log_level=50, n_theta=3, noise=False, var_type=["num"])
S.fit(X, y)
Sf = Kriging(name='kriging', seed=123, log_level=50, n_theta=3, noise=False, var_type=["num"])
Sf.fit(X, y)
n = 50
X0 = gen.scipy_lhd(n, lower=lower, upper = upper)
X1 = np.random.randint(low=1, high=3, size=(n,))
X = np.c_[X0, X1]
y = fun(X)
s=np.sum(np.abs(S.predict(X)[0] - y))
sf=np.sum(np.abs(Sf.predict(X)[0] - y))
sf - s

```

-55.49075685088792

```
# vars(S)
```

```
# vars(Sf)
```

# 13 Handling Noise

This chapter demonstrates how noisy functions can be handled by Spot and how noise can be simulated, i.e., added to the objective function.

## 13.1 Example: Spot and the Noisy Sphere Function

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
import matplotlib.pyplot as plt
from spotPython.utils.init import fun_control_init
from spotPython.utils.file import get_spot_tensorboard_path
from spotPython.utils.init import fun_control_init, design_control_init, surrogate_control_i
PREFIX = "08"
```

### 13.1.1 The Objective Function: Noisy Sphere

The `spotPython` package provides several classes of objective functions, which return a one-dimensional output  $y = f(x)$  for a given input  $x$  (independent variable). Several objective functions allow one- or multidimensional input, some also combinations of real-valued and categorial input values.

An objective function is considered as “analytical” if it can be described by a closed mathematical formula, e.g.,

$$f(x, y) = x^2 + y^2.$$

To simulate measurement errors, adding artificial noise to the function value  $y$  is a common practice, e.g.,:

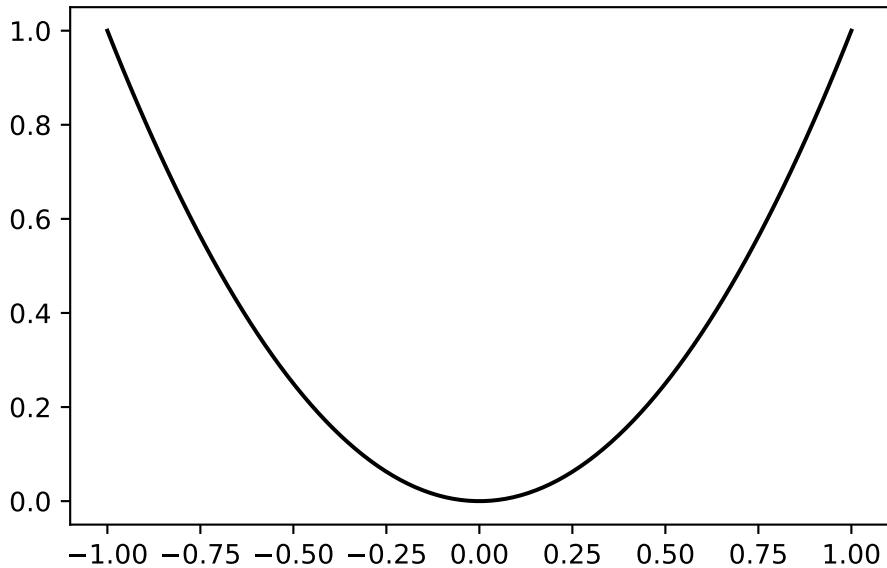
$$f(x, y) = x^2 + y^2 + \epsilon.$$

Usually, noise is assumed to be normally distributed with mean  $\mu = 0$  and standard deviation  $\sigma$ . spotPython uses numpy's `scale` parameter, which specifies the standard deviation (spread or "width") of the distribution is used. This must be a non-negative value, see <https://numpy.org/doc/stable/reference/random/generated/numpy.random.normal.html>.

#### Example: The sphere function without noise

The default setting does not use any noise.

```
from spotPython.fun.objectivefunctions import analytical
fun = analytical().fun_sphere
x = np.linspace(-1,1,100).reshape(-1,1)
y = fun(x)
plt.figure()
plt.plot(x,y, "k")
plt.show()
```



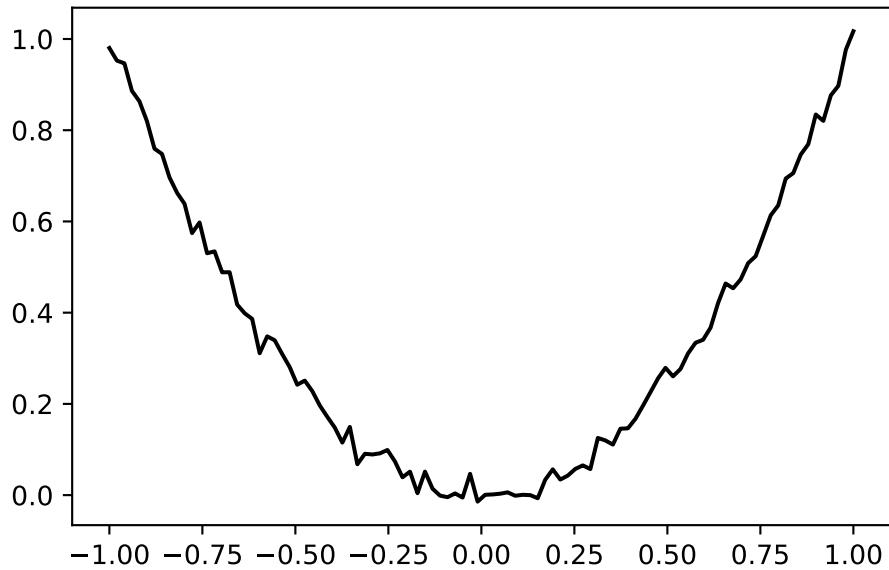
#### Example: The sphere function with noise

Noise can be added to the sphere function as follows:

```

from spotPython.fun.objectivefunctions import analytical
fun = analytical(seed=123, sigma=0.02).fun_sphere
x = np.linspace(-1,1,100).reshape(-1,1)
y = fun(x)
plt.figure()
plt.plot(x,y, "k")
plt.show()

```



### 13.1.2 Reproducibility: Noise Generation and Seed Handling

spotPython provides two mechanisms for generating random noise:

1. The seed is initialized once, i.e., when the objective function is instantiated. This can be done using the following call: `fun = analytical(sigma=0.02, seed=123).fun_sphere`.
2. The seed is set every time the objective function is called. This can be done using the following call: `y = fun(x, sigma=0.02, seed=123)`.

These two different ways lead to different results as explained in the following tables:

### **i** Example: Noise added to the sphere function

Since `sigma` is set to 0.02, noise is added to the function:

```
from spotPython.fun.objectivefunctions import analytical
fun = analytical(sigma=0.02, seed=123).fun_sphere
x = np.array([1]).reshape(-1,1)
for i in range(3):
    print(f"{i}: {fun(x)}")
```

```
0: [0.98021757]
1: [0.99264427]
2: [1.02575851]
```

The seed is set once. Every call to `fun()` results in a different value. The whole experiment can be repeated, the initial seed is used to generate the same sequence as shown below:

### **i** Example: Noise added to the sphere function

Since `sigma` is set to 0.02, noise is added to the function:

```
from spotPython.fun.objectivefunctions import analytical
fun = analytical(sigma=0.02, seed=123).fun_sphere
x = np.array([1]).reshape(-1,1)
for i in range(3):
    print(f"{i}: {fun(x)}")
```

```
0: [0.98021757]
1: [0.99264427]
2: [1.02575851]
```

If `spotPython` is used as a hyperparameter tuner, it is important that only one realization of the noise function is optimized. This behaviour can be accomplished by passing the same seed via the dictionary `fun_control` to every call of the objective function `fun` as shown below:

### **i** Example: The same noise added to the sphere function

Since `sigma` is set to 0.02, noise is added to the function:

```

from spotPython.fun.objectivefunctions import analytical
fun = analytical().fun_sphere
fun_control = fun_control_init(
    PREFIX=PREFIX,
    sigma=0.02)
y = fun(x, fun_control=fun_control)
x = np.array([1]).reshape(-1,1)
for i in range(3):
    print(f"{i}: {fun(x)}")

```

```

Created spot_tensorboard_path: runs/spot_logs/08_p040025_2024-01-09_20-20-12 for SummaryWriter
0: [0.98021757]
1: [0.98021757]
2: [0.98021757]

```

## 13.2 spotPython's Noise Handling Approaches

The following setting will be used for the next steps:

```

fun = analytical().fun_sphere
fun_control = fun_control_init(
    PREFIX=PREFIX,
    sigma=0.02,
)

```

```

Created spot_tensorboard_path: runs/spot_logs/08_p040025_2024-01-09_20-20-12 for SummaryWriter

```

spotPython is adopted as follows to cope with noisy functions:

1. `fun_repeats` is set to a value larger than 1 (here: 2)
2. `noise` is set to `true`. Therefore, a nugget (`Lambda`) term is added to the correlation matrix
3. `init size` (of the `design_control` dictionary) is set to a value larger than 1 (here: 3)

```

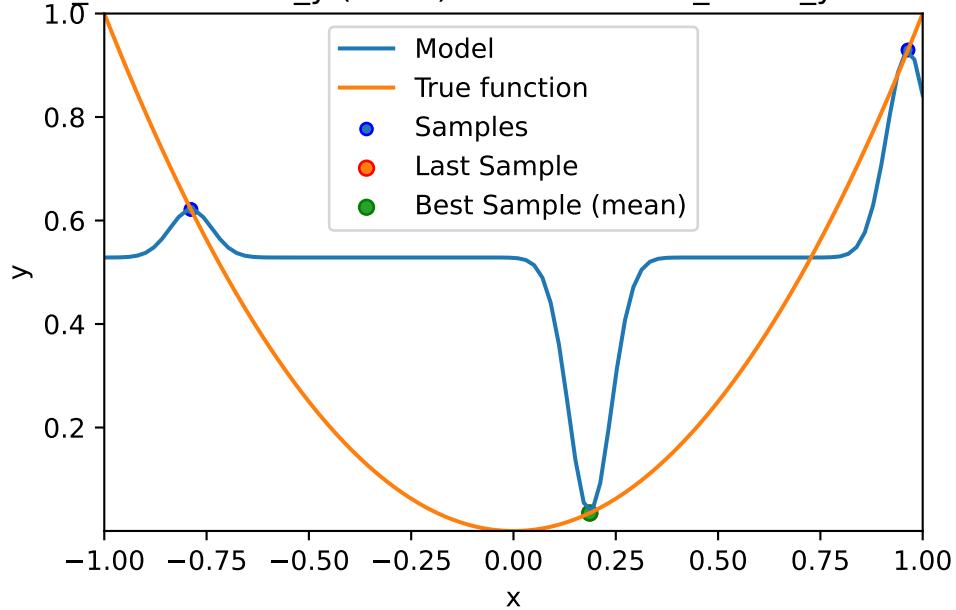
spot_1_noisy = spot.Spot(fun=fun,
                        fun_control=fun_control_init(
                            lower = np.array([-1]),
                            upper = np.array([1]),
                            fun_evals = 20,
)

```

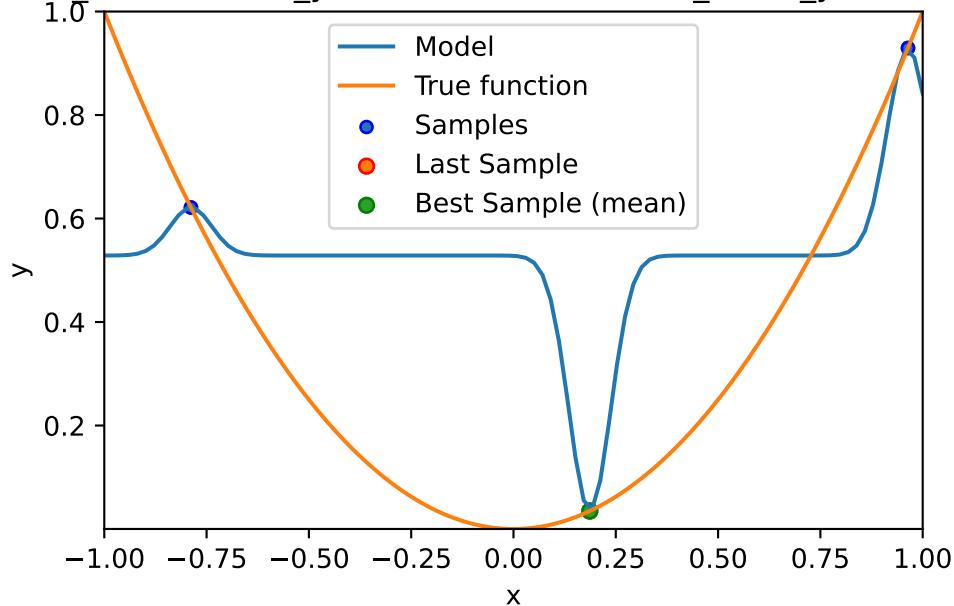
```
        fun_repeats = 2,  
        noise = True,  
        show_models=True),  
    design_control=design_control_init(init_size=3, repeats=2),  
    surrogate_control=surrogate_control_init(noise=True))
```

```
spot_1_noisy.run()
```

fun\_evals: 6. min\_y (noise): 0.034755 min\_mean\_y: 0.034755

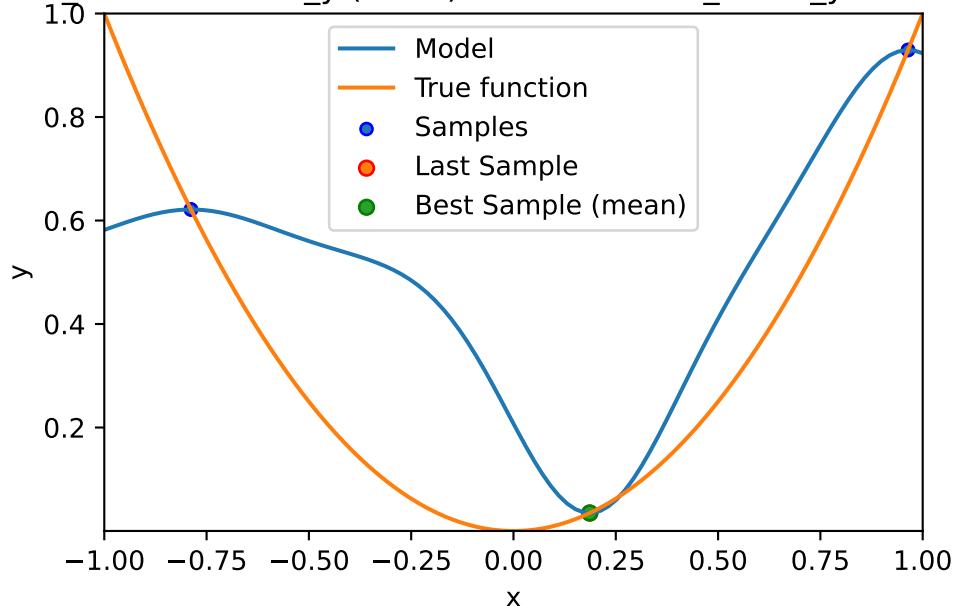


fun\_evals: 8. min\_y (noise): 0.034755 min\_mean\_y: 0.034755

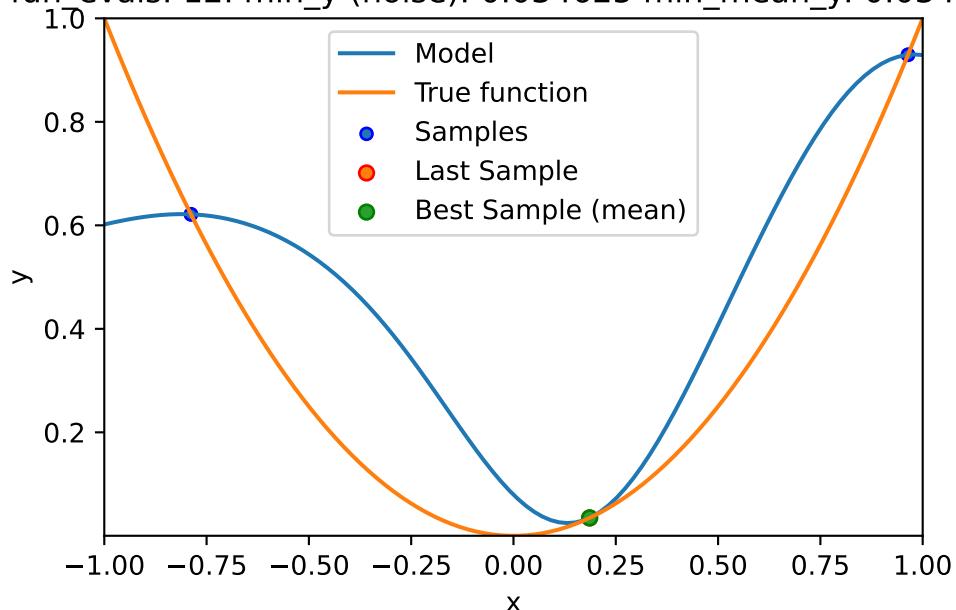


```
spotPython tuning: 0.034754931797000475 [#####-----] 40.00%
spotPython tuning: 0.0347549291052682 [#####-----] 50.00%
spotPython tuning: 0.03462521823117037 [#####----] 60.00%
spotPython tuning: 0.017469986696189304 [#####----] 70.00%
spotPython tuning: 0.004000038051356374 [#####----] 80.00%
spotPython tuning: 0.00022058859554626256 [#####----] 90.00%
spotPython tuning: 9.544809881586792e-06 [#####----] 100.00% Done...
```

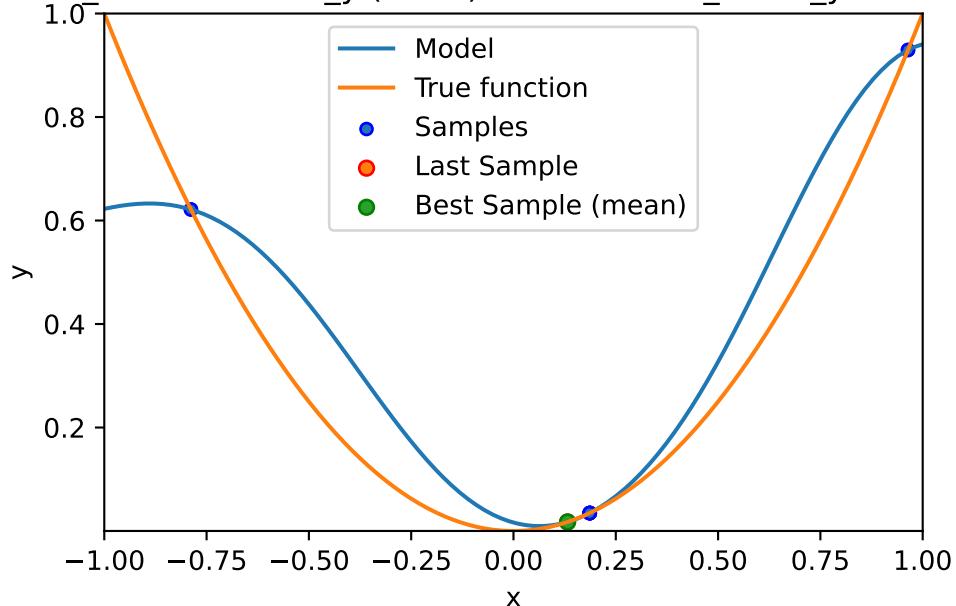
fun\_evals: 10. min\_y (noise): 0.034755 min\_mean\_y: 0.034755



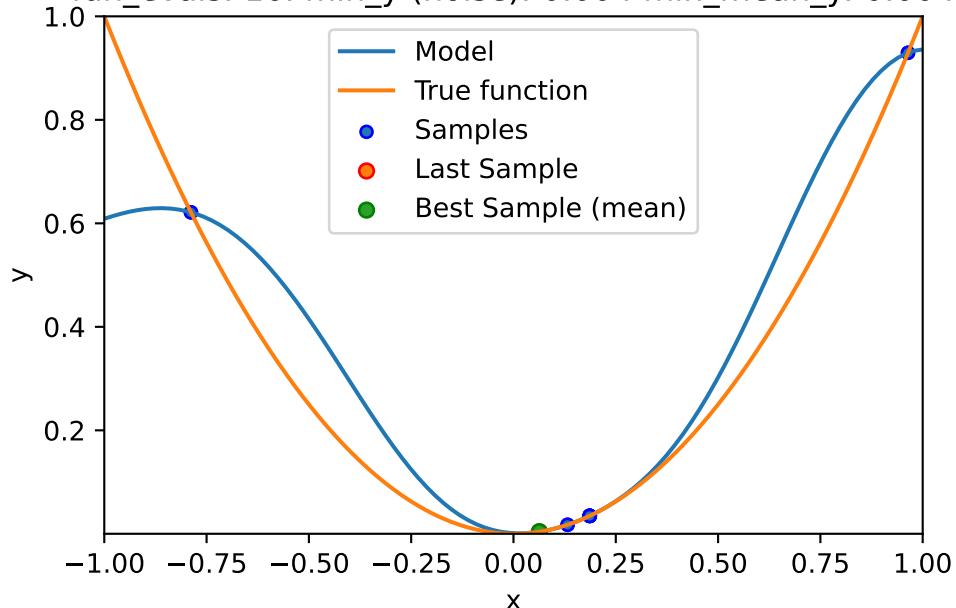
fun\_evals: 12. min\_y (noise): 0.034625 min\_mean\_y: 0.034625



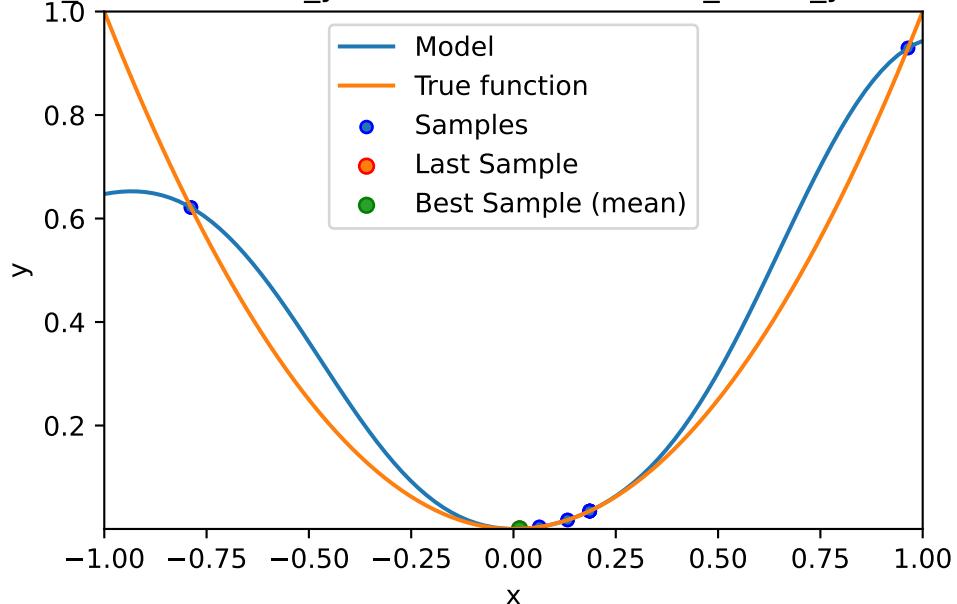
fun\_evals: 14. min\_y (noise): 0.01747 min\_mean\_y: 0.01747



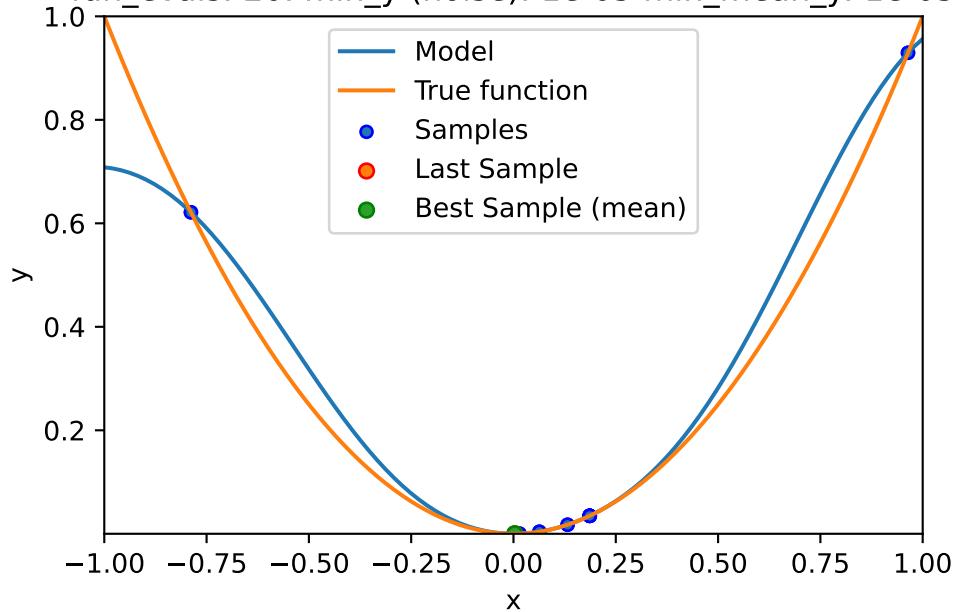
fun\_evals: 16. min\_y (noise): 0.004 min\_mean\_y: 0.004



fun\_evals: 18. min\_y (noise): 0.000221 min\_mean\_y: 0.000221



fun\_evals: 20. min\_y (noise): 1e-05 min\_mean\_y: 1e-05



### 13.3 Print the Results

```
spot_1_noisy.print_results()
```

```
min y: 9.544809881586792e-06
x0: 0.003089467572509346
min mean y: 9.544809881586792e-06
x0: 0.003089467572509346
```

```
[['x0', 0.003089467572509346], ['x0', 0.003089467572509346]]
```

```
spot_1_noisy.plot_progress(log_y=False,
    filename="./figures/" + PREFIX + "_progress.png")
```

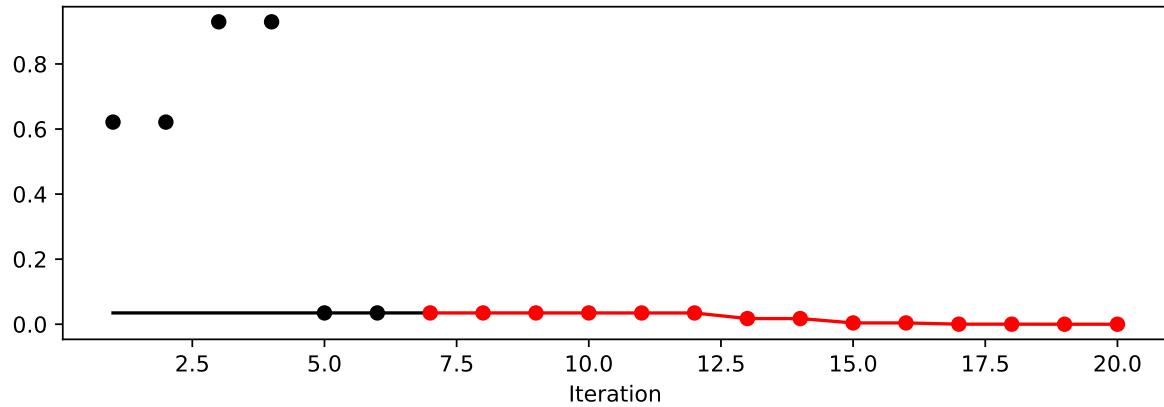


Figure 13.1: Progress plot. *Black* dots denote results from the initial design. *Red* dots illustrate the improvement found by the surrogate model based optimization.

### 13.4 Noise and Surrogates: The Nugget Effect

#### 13.4.1 The Noisy Sphere

##### 13.4.1.1 The Data

- We prepare some data first:

```

import numpy as np
import spotPython
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from spotPython.design.spacefilling import spacefilling
from spotPython.build.kriging import Kriging
import matplotlib.pyplot as plt

gen = spacefilling(1)
rng = np.random.RandomState(1)
lower = np.array([-10])
upper = np.array([10])
fun = analytical().fun_sphere
fun_control = fun_control_init(
    PREFIX=PREFIX,
    sigma=4)
X = gen.scipy_lhd(10, lower=lower, upper = upper)
y = fun(X, fun_control=fun_control)
X_train = X.reshape(-1,1)
y_train = y

```

Created spot\_tensorboard\_path: runs/spot\_logs/08\_p040025\_2024-01-09\_20-20-30 for SummaryWriter

- A surrogate without nugget is fitted to these data:

```

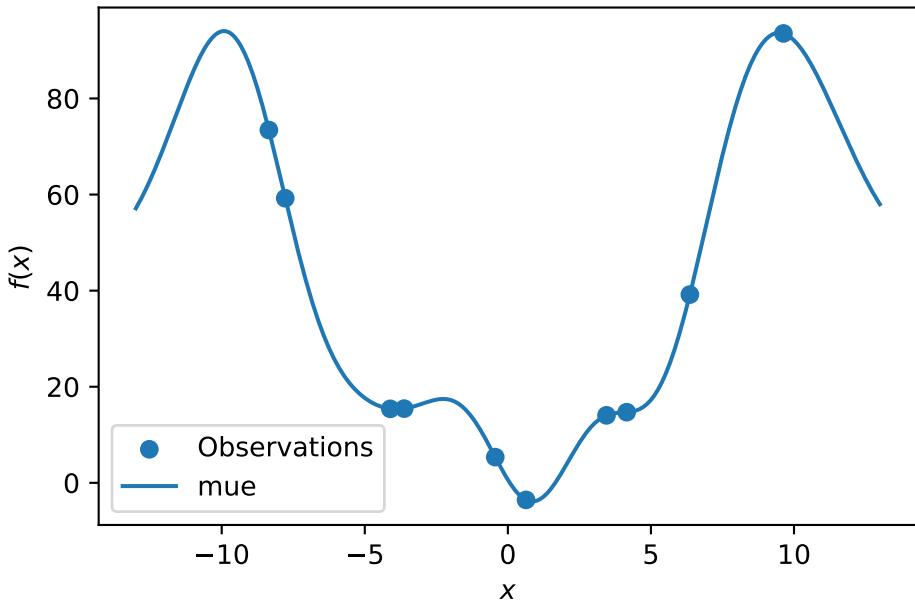
S = Kriging(name='kriging',
            n_theta=1,
            noise=False)
S.fit(X_train, y_train)

X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X_axis, return_val="all")

plt.scatter(X_train, y_train, label="Observations")
plt.plot(X_axis, mean_prediction, label="mu")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Sphere: Gaussian process regression on noisy dataset")

```

Sphere: Gaussian process regression on noisy dataset



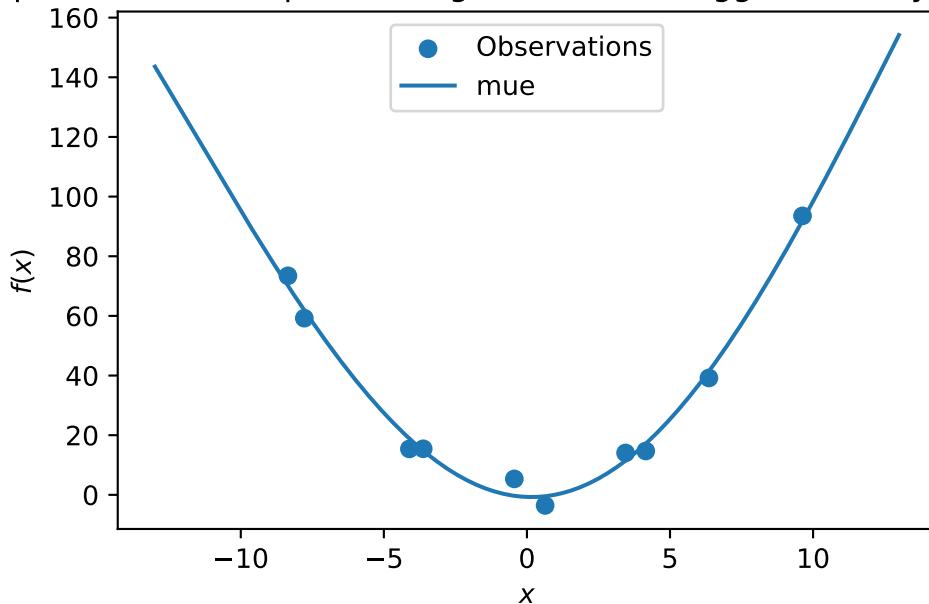
- In comparison to the surrogate without nugget, we fit a surrogate with nugget to the data:

```

S_nug = Kriging(name='kriging',
                  n_theta=1,
                  noise=True)
S_nug.fit(X_train, y_train)
X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S_nug.predict(X_axis, return_val="all")
plt.scatter(X_train, y_train, label="Observations")
plt.plot(X_axis, mean_prediction, label="mle")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Sphere: Gaussian process regression with nugget on noisy dataset")

```

Sphere: Gaussian process regression with nugget on noisy dataset



- The value of the nugget term can be extracted from the model as follows:

```
S.Lambda
```

```
S_nug.Lambda
```

```
0.00055921881757264
```

- We see:
  - the first model `S` has no nugget,
  - whereas the second model has a nugget value (`Lambda`) larger than zero.

## 13.5 Exercises

### 13.5.1 Noisy fun\_cubed

- Analyse the effect of noise on the `fun_cubed` function with the following settings:

```
fun = analytical().fun_cubed
fun_control = fun_control_init(
    sigma=10)
lower = np.array([-10])
upper = np.array([10])
```

### 13.5.2 fun\_runge

- Analyse the effect of noise on the `fun_runge` function with the following settings:

```
lower = np.array([-10])
upper = np.array([10])
fun = analytical().fun_runge
fun_control = fun_control_init(
    sigma=0.25)
```

### 13.5.3 fun\_forrester

- Analyse the effect of noise on the `fun_forrester` function with the following settings:

```
lower = np.array([0])
upper = np.array([1])
fun = analytical().fun_forrester
fun_control = fun_control_init(
    sigma=5)
```

### 13.5.4 fun\_xsin

- Analyse the effect of noise on the `fun_xsin` function with the following settings:

```
lower = np.array([-1.])
upper = np.array([1.])
fun = analytical().fun_xsin
fun_control = fun_control_init(
    sigma=0.5)
```

# 14 Optimal Computational Budget Allocation in Spot

This chapter demonstrates how noisy functions can be handled with Optimal Computational Budget Allocation (OCBA) by Spot.

## 14.1 Example: Spot, OCBA, and the Noisy Sphere Function

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
import matplotlib.pyplot as plt
from spotPython.utils.init import fun_control_init
from spotPython.utils.file import get_spot_tensorboard_path
from spotPython.utils.init import fun_control_init, design_control_init, surrogate_control_i
PREFIX = "09"
```

### 14.1.1 The Objective Function: Noisy Sphere

The `spotPython` package provides several classes of objective functions. We will use an analytical objective function with noise, i.e., a function that can be described by a (closed) formula:

$$f(x) = x^2 + \epsilon$$

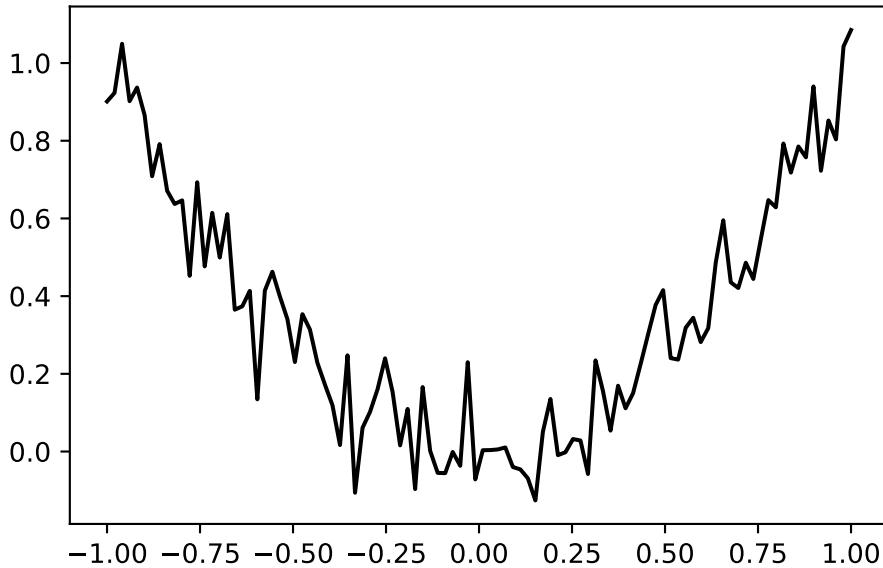
Since `sigma` is set to 0.1, noise is added to the function:

```
fun = analytical().fun_sphere
fun_control = fun_control_init(
    PREFIX=PREFIX,
    sigma=0.1)
```

```
Created spot_tensorboard_path: runs/spot_logs/09_p040025_2024-01-09_20-21-14 for SummaryWriter
```

A plot illustrates the noise:

```
x = np.linspace(-1,1,100).reshape(-1,1)
y = fun(x, fun_control=fun_control)
plt.figure()
plt.plot(x,y, "k")
plt.show()
```



Spot is adopted as follows to cope with noisy functions:

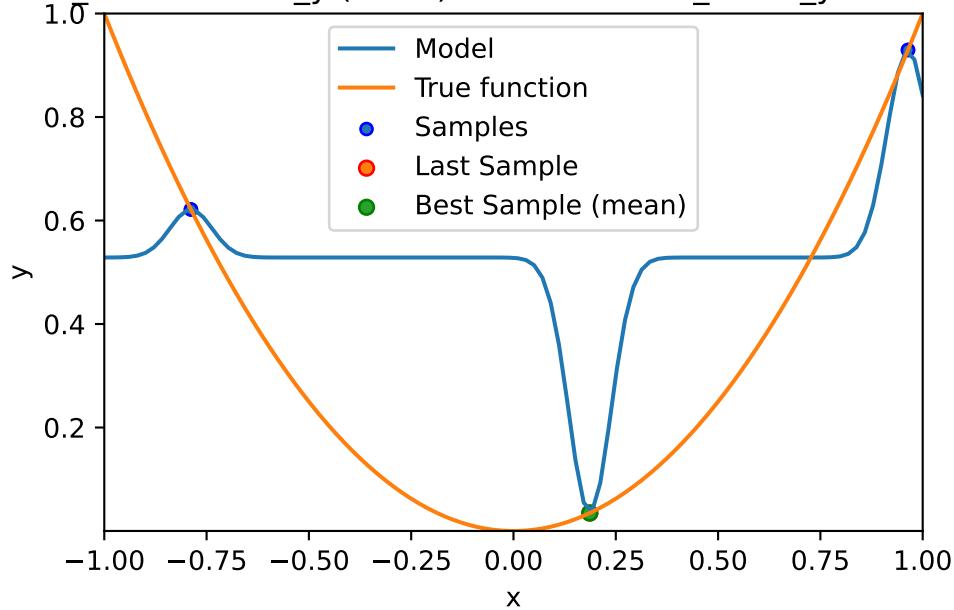
1. `fun_repeats` is set to a value larger than 1 (here: 2)
2. `noise` is set to `true`. Therefore, a nugget (`Lambda`) term is added to the correlation matrix
3. `init_size` (of the `design_control` dictionary) is set to a value larger than 1 (here: 2)

```
spot_1_noisy = spot.Spot(fun=fun,
                          fun_control=fun_control_init(
                            lower = np.array([-1]),
                            upper = np.array([1]),
                            fun_evals = 20,
                            fun_repeats = 2,
                            infill_criterion="ei",
                            noise = True,
```

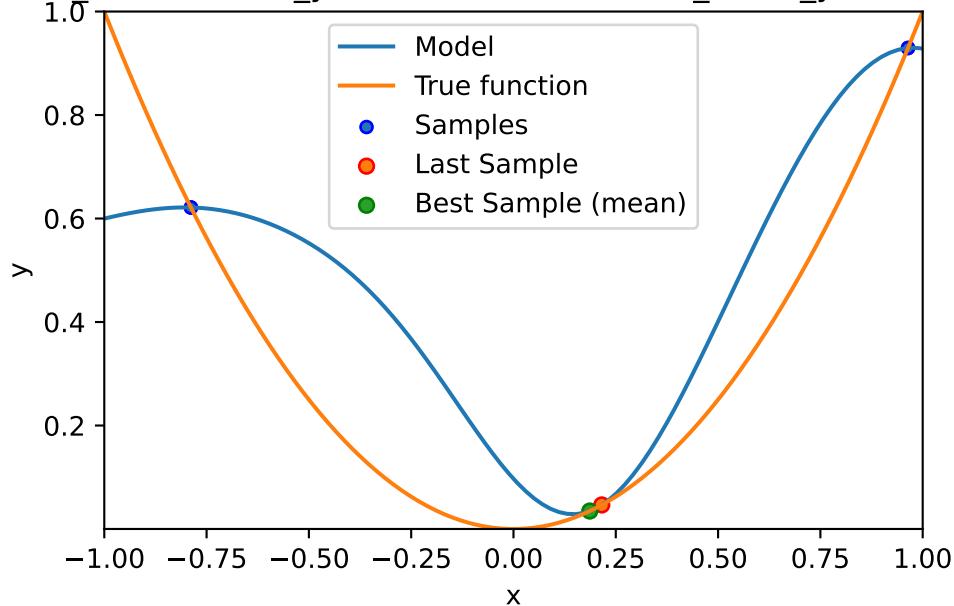
```
tolerance_x=0.0,  
ocba_delta = 1,  
show_models=True),  
design_control=design_control_init(init_size=3, repeats=2),  
surrogate_control=surrogate_control_init(noise=True))
```

```
spot_1_noisy.run()
```

fun\_evals: 6. min\_y (noise): 0.034755 min\_mean\_y: 0.034755

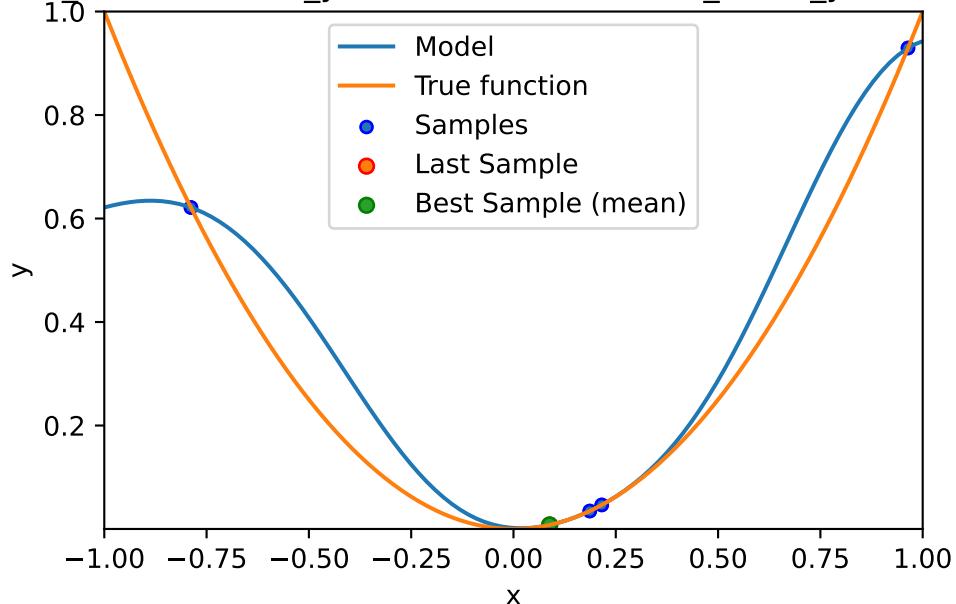


fun\_evals: 8. min\_y (noise): 0.034755 min\_mean\_y: 0.034755

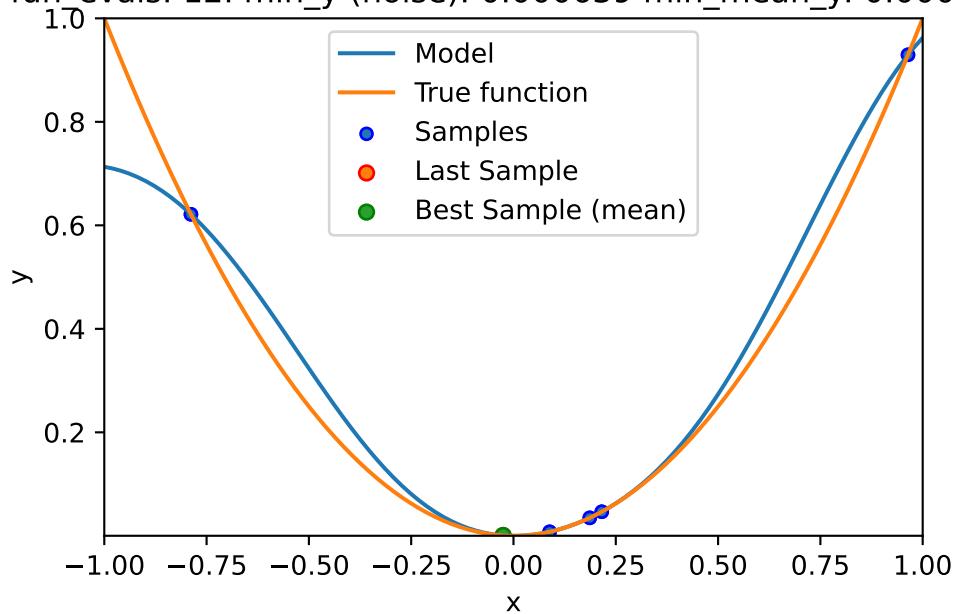


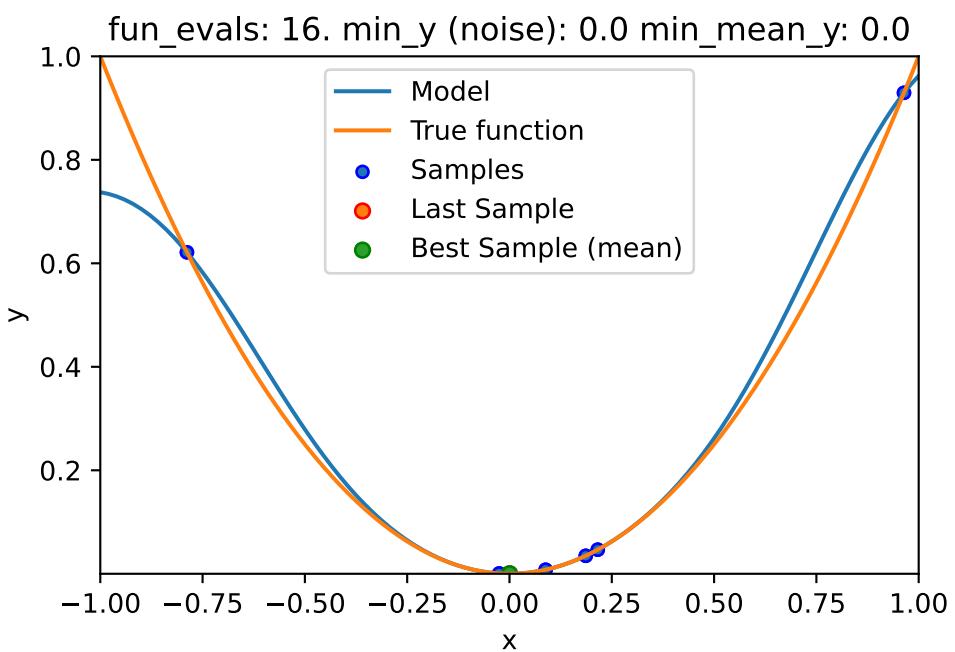
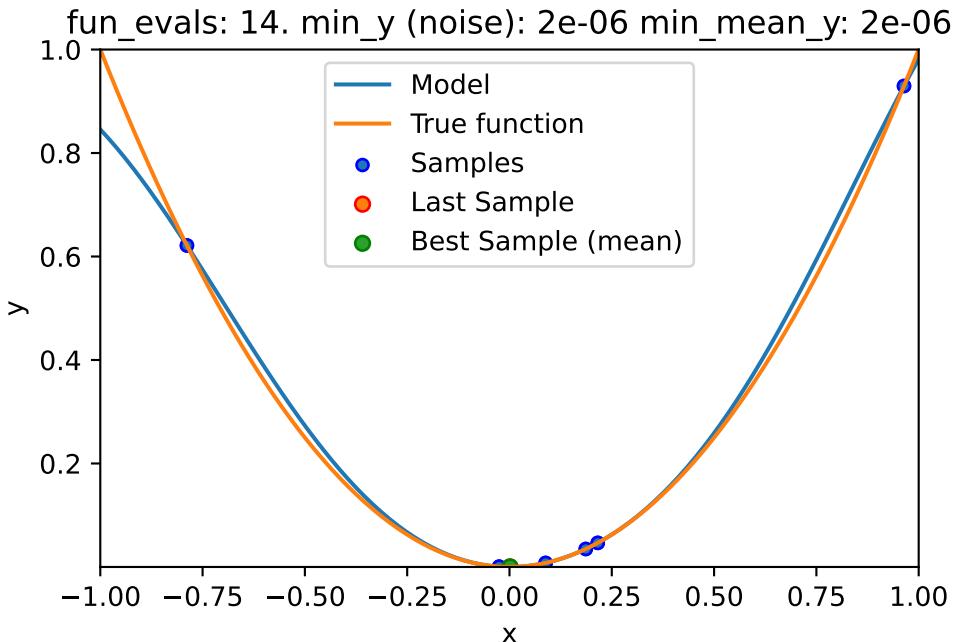
```
spotPython tuning: 0.03475493366922229 [#####-----] 40.00%
spotPython tuning: 0.00779944837528439 [######-----] 50.00%
spotPython tuning: 0.0006387949380800254 [#######----] 60.00%
spotPython tuning: 1.69966282482316e-06 [########---] 70.00%
spotPython tuning: 3.8093880816701285e-08 [#######---] 80.00%
spotPython tuning: 4.094750152333457e-09 [########---] 90.00%
spotPython tuning: 2.5113140448254767e-09 [#######---] 100.00% Done...
```

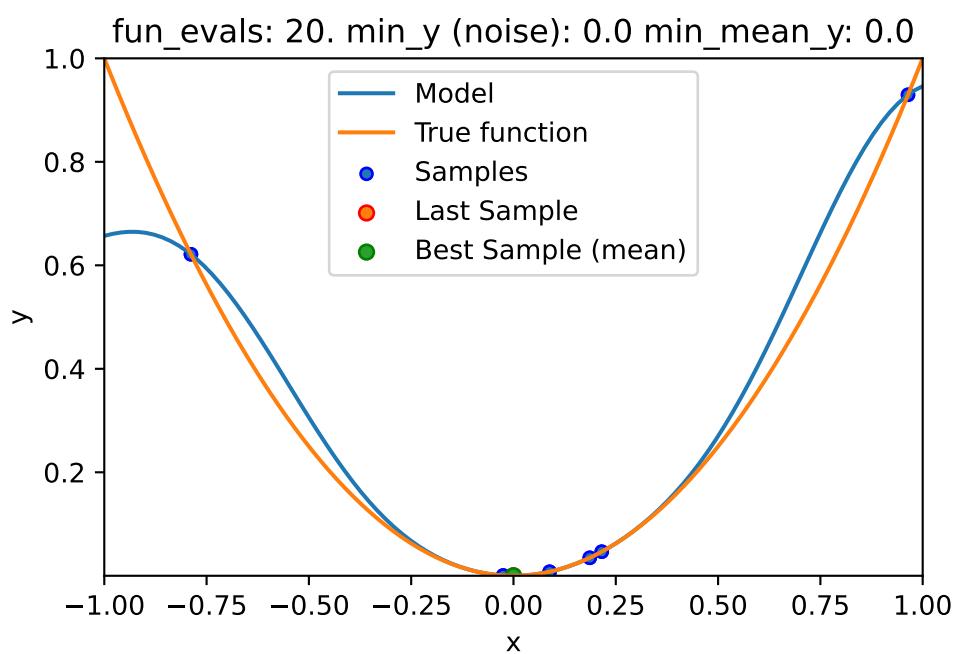
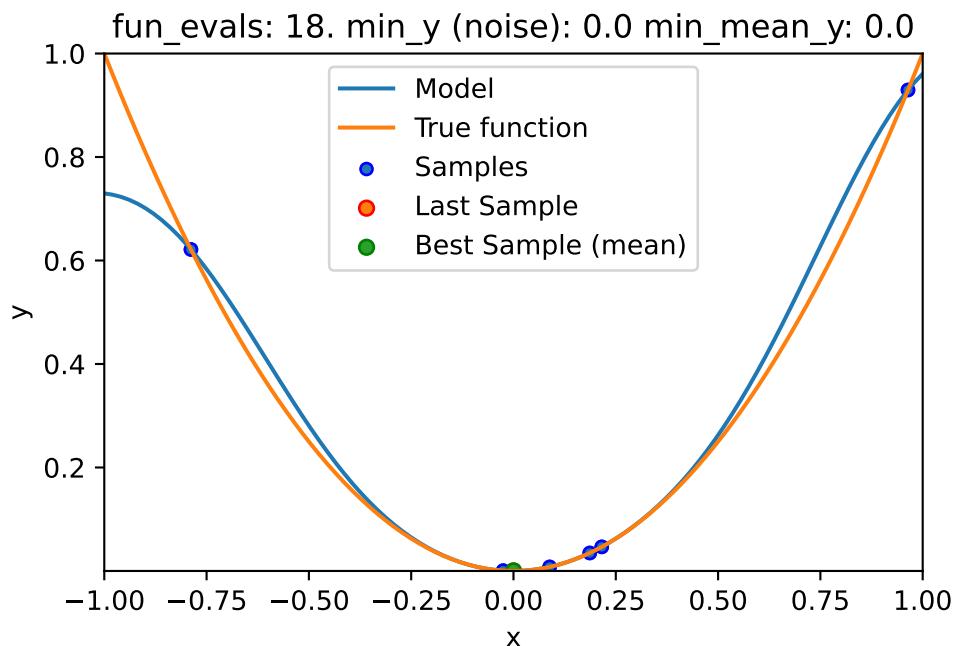
fun\_evals: 10. min\_y (noise): 0.007799 min\_mean\_y: 0.007799



fun\_evals: 12. min\_y (noise): 0.000639 min\_mean\_y: 0.000639







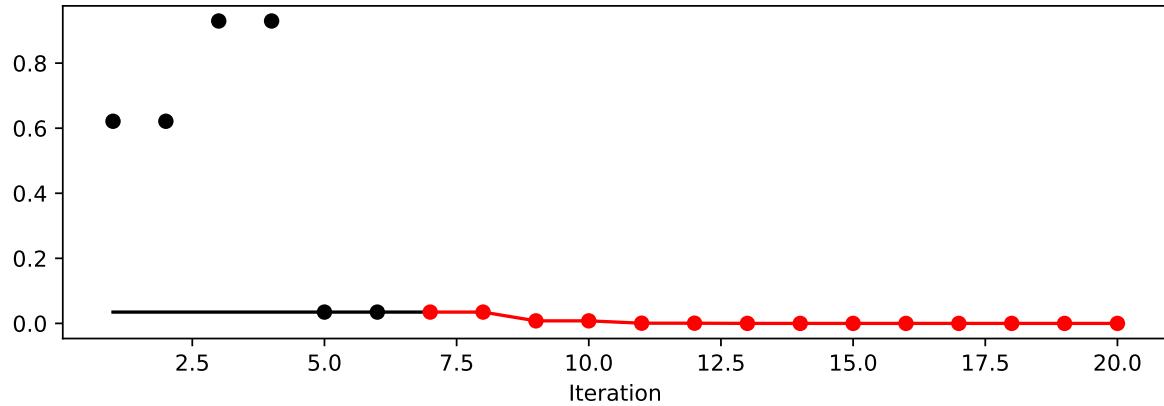
## 14.2 Print the Results

```
spot_1_noisy.print_results()
```

```
min y: 2.5113140448254767e-09
x0: 5.011301272948451e-05
min mean y: 2.5113140448254767e-09
x0: 5.011301272948451e-05

[['x0', 5.011301272948451e-05], ['x0', 5.011301272948451e-05]]
```

```
spot_1_noisy.plot_progress(log_y=False)
```



## 14.3 Noise and Surrogates: The Nugget Effect

### 14.3.1 The Noisy Sphere

#### 14.3.1.1 The Data

We prepare some data first:

```

import numpy as np
import spotPython
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from spotPython.design.spacefilling import spacefilling
from spotPython.build.kriging import Kriging
import matplotlib.pyplot as plt

gen = spacefilling(1)
rng = np.random.RandomState(1)
lower = np.array([-10])
upper = np.array([10])
fun = analytical().fun_sphere
fun_control = fun_control_init(
    sigma=2,
    seed=125)
X = gen.scipy_lhd(10, lower=lower, upper = upper)
y = fun(X, fun_control=fun_control)
X_train = X.reshape(-1,1)
y_train = y

```

A surrogate without nugget is fitted to these data:

```

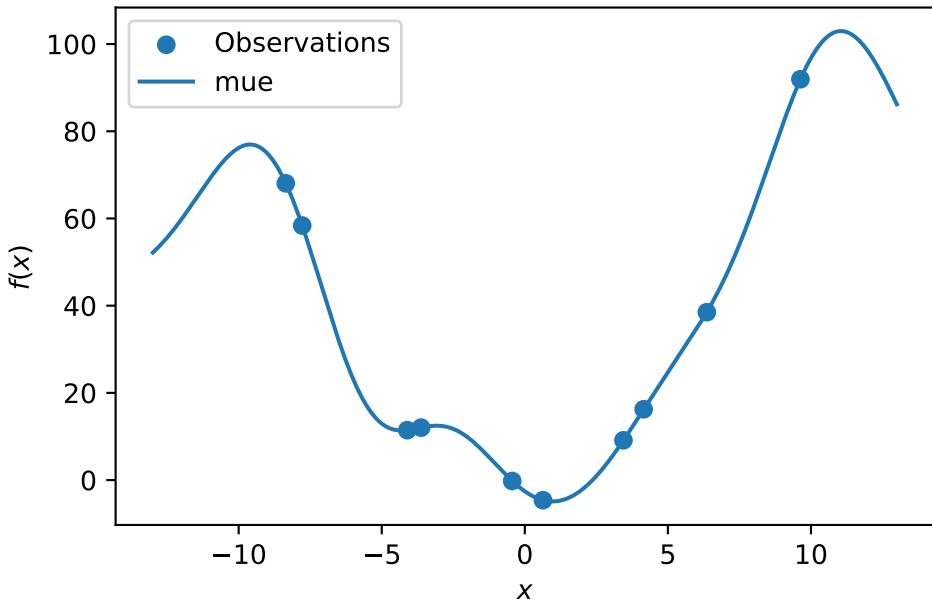
S = Kriging(name='kriging',
            seed=123,
            log_level=50,
            n_theta=1,
            noise=False)
S.fit(X_train, y_train)

X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S.predict(X_axis, return_val="all")

plt.scatter(X_train, y_train, label="Observations")
plt.plot(X_axis, mean_prediction, label="mu")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Sphere: Gaussian process regression on noisy dataset")

```

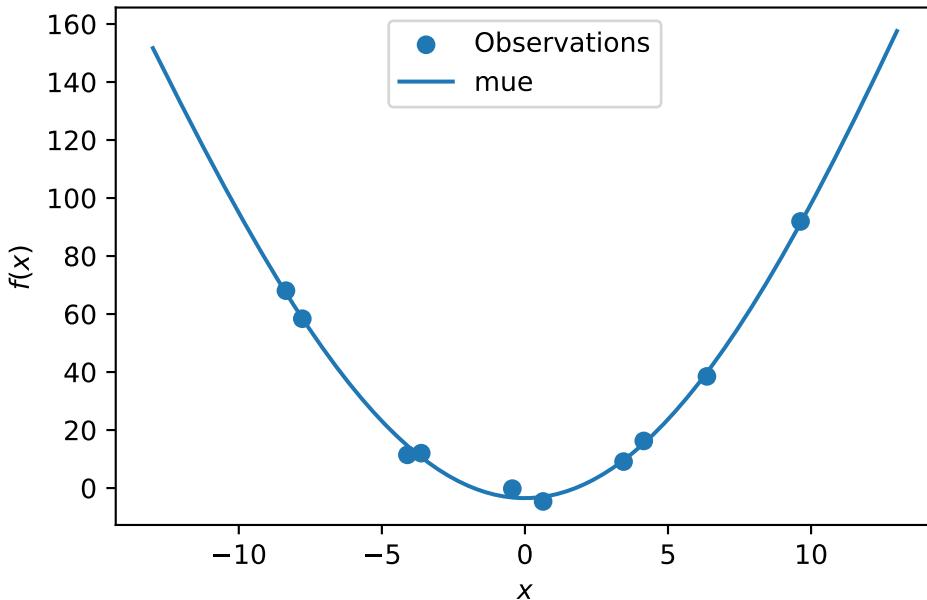
### Sphere: Gaussian process regression on noisy dataset



In comparison to the surrogate without nugget, we fit a surrogate with nugget to the data:

```
S_nug = Kriging(name='kriging',
                  seed=123,
                  log_level=50,
                  n_theta=1,
                  noise=True)
S_nug.fit(X_train, y_train)
X_axis = np.linspace(start=-13, stop=13, num=1000).reshape(-1, 1)
mean_prediction, std_prediction, ei = S_nug.predict(X_axis, return_val="all")
plt.scatter(X_train, y_train, label="Observations")
plt.plot(X_axis, mean_prediction, label="mle")
plt.legend()
plt.xlabel("$x$")
plt.ylabel("$f(x)$")
_ = plt.title("Sphere: Gaussian process regression with nugget on noisy dataset")
```

Sphere: Gaussian process regression with nugget on noisy dataset



The value of the nugget term can be extracted from the model as follows:

```
S.Lambda
```

```
S_nug.Lambda
```

```
8.374496269458742e-05
```

We see:

- the first model `S` has no nugget,
- whereas the second model has a nugget value (`Lambda`) larger than zero.

## 14.4 Exercises

### 14.4.1 Noisy fun\_cubed

Analyse the effect of noise on the `fun_cubed` function with the following settings:

```
fun = analytical().fun_cubed
fun_control = fun_control_init(
    sigma=10,
    seed=123)
lower = np.array([-10])
upper = np.array([10])
```

#### 14.4.2 fun\_runge

Analyse the effect of noise on the `fun_runge` function with the following settings:

```
lower = np.array([-10])
upper = np.array([10])
fun = analytical().fun_runge
fun_control = fun_control_init(
    sigma=0.25,
    seed=123)
```

#### 14.4.3 fun\_forrester

Analyse the effect of noise on the `fun_forrester` function with the following settings:

```
lower = np.array([0])
upper = np.array([1])
fun = analytical().fun_forrester
fun_control = {"sigma": 5,
               "seed": 123}
```

#### 14.4.4 fun\_xsin

Analyse the effect of noise on the `fun_xsin` function with the following settings:

```
lower = np.array([-1.])
upper = np.array([1.])
fun = analytical().fun_xsin
fun_control = fun_control_init(
    sigma=0.5,
    seed=123)
```

## **Part III**

# **Hyperparameter Tuning with PyTorch Lightning**

# 15 HPT PyTorch Lightning: Diabetes

In this tutorial, we will show how `spotPython` can be integrated into the PyTorch Lightning training workflow for a regression task.

This chapter describes the hyperparameter tuning of a PyTorch Lightning network on the Diabetes data set. This is a PyTorch Dataset for regression. A toy data set from scikit-learn. Ten baseline variables, age, sex, body mass index, average blood pressure, and six blood serum measurements were obtained for each of  $n = 442$  diabetes patients, as well as the response of interest, a quantitative measure of disease progression one year after baseline.

## 15.1 Step 1: Setup

- Before we consider the detailed experimental setup, we select the parameters that affect run time, initial design size, etc.
- The parameter `MAX_TIME` specifies the maximum run time in seconds.
- The parameter `INIT_SIZE` specifies the initial design size.
- The parameter `WORKERS` specifies the number of workers.
- The prefix `PREFIX` is used for the experiment name and the name of the log file.
- The parameter `DEVICE` specifies the device to use for training.

```
from spotPython.utils.device import getDevice
from math import inf

MAX_TIME = 1
FUN_EVALS = inf
INIT_SIZE = 5
WORKERS = 0
PREFIX="031"
DEVICE = getDevice()
DEVICES = 1
TEST_SIZE = 0.1
```



Caution: Run time and initial design size should be increased for real experiments

- `MAX_TIME` is set to one minute for demonstration purposes. For real experiments, this should be increased to at least 1 hour.
- `INIT_SIZE` is set to 5 for demonstration purposes. For real experiments, this should be increased to at least 10.
- `WORKERS` is set to 0 for demonstration purposes. For real experiments, this should be increased. See the warnings that are printed when the number of workers is set to 0.



Note: Device selection

- Although there are no `.cuda()` or `.to(device)` calls required, because Lightning does these for you, see [LIGHTNINGMODULE](#), we would like to know which device is used. Therefore, we imitate the `LightningModule` behaviour which selects the highest device.
- The method `spotPython.utils.device.getDevice()` returns the device that is used by Lightning.

## 15.2 Step 2: Initialization of the `fun_control` Dictionary

`spotPython` uses a Python dictionary for storing the information required for the hyperparameter tuning process.

```
from spotPython.utils.init import fun_control_init
import numpy as np
fun_control = fun_control_init(
    _L_in=10,
    _L_out=1,
    PREFIX=PREFIX,
    TENSORBOARD_CLEAN=True,
    device=DEVICE,
    enable_progress_bar=False,
    fun_evals=FUN_EVALS,
    log_level=10,
    max_time=MAX_TIME,
    num_workers=WORKERS,
    show_progress=True,
    test_size=0.1,
    tolerance_x=np.sqrt(np.spacing(1)),
)
```

```
Moving TENSORBOARD_PATH: runs/ to TENSORBOARD_PATH_OLD: runs_OLD/runs_2024_01_09_20_26_16  
Created spot_tensorboard_path: runs/spot_logs/031_p040025_2024-01-09_20-26-16 for SummaryWriter
```

### 15.3 Step 3: Loading the Diabetes Data Set

```
from spotPython.hyperparameters.values import set_control_key_value  
from spotPython.data.diabetes import Diabetes  
dataset = Diabetes()  
set_control_key_value(control_dict=fun_control,  
                      key="data_set",  
                      value=dataset,  
                      replace=True)  
print(len(dataset))
```

442

#### Note: Data Set and Data Loader

- As shown below, a DataLoader from `torch.utils.data` can be used to check the data.

```
# Set batch size for DataLoader  
batch_size = 5  
# Create DataLoader  
from torch.utils.data import DataLoader  
dataloader = DataLoader(dataset, batch_size=batch_size, shuffle=False)  
  
# Iterate over the data in the DataLoader  
for batch in dataloader:  
    inputs, targets = batch  
    print(f"Batch Size: {inputs.size(0)}")  
    print(f"Inputs Shape: {inputs.shape}")  
    print(f"Targets Shape: {targets.shape}")  
    print("-----")  
    print(f"Inputs: {inputs}")  
    print(f"Targets: {targets}")  
    break
```

```
Batch Size: 5  
Inputs Shape: torch.Size([5, 10])
```

```

Targets Shape: torch.Size([5])
-----
Inputs: tensor([[ 0.0381,  0.0507,  0.0617,  0.0219, -0.0442, -0.0348, -0.0434, -0.0026,
                 0.0199, -0.0176],
               [-0.0019, -0.0446, -0.0515, -0.0263, -0.0084, -0.0192,  0.0744, -0.0395,
                -0.0683, -0.0922],
               [ 0.0853,  0.0507,  0.0445, -0.0057, -0.0456, -0.0342, -0.0324, -0.0026,
                 0.0029, -0.0259],
               [-0.0891, -0.0446, -0.0116, -0.0367,  0.0122,  0.0250, -0.0360,  0.0343,
                0.0227, -0.0094],
               [ 0.0054, -0.0446, -0.0364,  0.0219,  0.0039,  0.0156,  0.0081, -0.0026,
                -0.0320, -0.0466]])
Targets: tensor([151.,  75., 141., 206., 135.])

```

## 15.4 Step 4: Preprocessing

Preprocessing is handled by Lightning and PyTorch. It is described in the [LIGHTNING-DATAMODULE](#) documentation. Here you can find information about the `transforms` methods.

## 15.5 Step 5: Select the Core Model (algorithm) and core\_model\_hyper\_dict

spotPython includes the `NetLightRegression` class [\[SOURCE\]](#) for configurable neural networks. The class is imported here. It inherits from the class `Lightning.LightningModule`, which is the base class for all models in `Lightning`. `Lightning.LightningModule` is a subclass of `torch.nn.Module` and provides additional functionality for the training and testing of neural networks. The class `Lightning.LightningModule` is described in the [Lightning documentation](#).

- Here we simply add the NN Model to the `fun_control` dictionary by calling the function `add_core_model_to_fun_control`:

```

from spotPython.light.regression.netlightregression import NetLightRegression
from spotPython.hyperdict.light_hyper_dict import LightHyperDict
from spotPython.hyperparameters.values import add_core_model_to_fun_control
add_core_model_to_fun_control(fun_control=fun_control,
                             core_model=NetLightRegression,
                             hyper_dict=LightHyperDict)

```

The hyperparameters of the model are specified in the `core_model_hyper_dict` dictionary [SOURCE].

## 15.6 Step 6: Modify `hyper_dict` Hyperparameters for the Selected Algorithm aka `core_model`

`spotPython` provides functions for modifying the hyperparameters, their bounds and factors as well as for activating and de-activating hyperparameters without re-compilation of the Python source code.



Caution: Small number of epochs for demonstration purposes

- `epochs` and `patience` are set to small values for demonstration purposes. These values are too small for a real application.
- More resonable values are, e.g.:
  - `set_control_hyperparameter_value(fun_control, "epochs", [7, 9])`  
and
  - `set_control_hyperparameter_value(fun_control, "patience", [2, 7])`

```
from spotPython.hyperparameters.values import set_control_hyperparameter_value

set_control_hyperparameter_value(fun_control, "l1", [4, 6])
set_control_hyperparameter_value(fun_control, "epochs", [9, 10])
set_control_hyperparameter_value(fun_control, "batch_size", [4, 5])
set_control_hyperparameter_value(fun_control, "optimizer", [
    "Adadelta",
    "Adagrad",
    "Adam",
    "AdamW",
    "Adamax",
    "NAdam",
    "RAdam",
    "RMSprop",
    "Rprop"
])
set_control_hyperparameter_value(fun_control, "dropout_prob", [0.01, 0.1])
set_control_hyperparameter_value(fun_control, "lr_mult", [0.5, 5.0])
set_control_hyperparameter_value(fun_control, "patience", [5, 7])
set_control_hyperparameter_value(fun_control, "act_fn", [
```

```

    "Sigmoid",
    "ReLU",
    "LeakyReLU",
    "Swish"
]
)

```

Now, the dictionary `fun_control` contains all information needed for the hyperparameter tuning. Before the hyperparameter tuning is started, it is recommended to take a look at the experimental design. The method `gen_design_table` [SOURCE] generates a design table as follows:

```
from spotPython.utils.eda import gen_design_table
print(gen_design_table(fun_control))
```

name	type	default	lower	upper	transform
l1	int	3	4	6	transform_power_2_int
epochs	int	4	9	10	transform_power_2_int
batch_size	int	4	4	5	transform_power_2_int
act_fn	factor	ReLU	0	3	None
optimizer	factor	SGD	0	8	None
dropout_prob	float	0.01	0.01	0.1	None
lr_mult	float	1.0	0.5	5	None
patience	int	2	5	7	transform_power_2_int
initialization	factor	Default	0	2	None

This allows to check if all information is available and if the information is correct.

**i** Note: Hyperparameters of the Tuned Model and the `fun_control` Dictionary

The updated `fun_control` dictionary can be shown with the command `fun_control["core_model_hyper_dict"]`.

## 15.7 Step 7: Data Splitting, the Objective (Loss) Function and the Metric

### 15.7.1 Evaluation

The evaluation procedure requires the specification of two elements:

1. the way how the data is split into a train and a test set
2. the loss function (and a metric).

 Caution: Data Splitting in Lightning

The data splitting is handled by **Lightning**.

### 15.7.2 Loss Function

The loss function is specified in the configurable network class [\[SOURCE\]](#). We will use MSE.

### 15.7.3 Metric

- Similar to the loss function, the metric is specified in the configurable network class [\[SOURCE\]](#).

 Caution: Loss Function and Metric in Lightning

- The loss function and the metric are not hyperparameters that can be tuned with `spotPython`.
- They are handled by **Lightning**.

## 15.8 Step 8: Calling the SPOT Function

### 15.8.1 Preparing the SPOT Call

```
from spotPython.utils.init import design_control_init, surrogate_control_init
design_control = design_control_init(init_size=INIT_SIZE)

surrogate_control = surrogate_control_init(noise=True,
                                            n_theta=2)
```

 Note: Modifying Values in the Control Dictionaries

- The values in the control dictionaries can be modified with the function `set_control_key_value` [\[SOURCE\]](#), for example:

```
set_control_key_value(control_dict=surrogate_control,
                      key="noise",
                      value=True,
                      replace=True)
set_control_key_value(control_dict=surrogate_control,
                      key="n_theta",
                      value=2,
                      replace=True)
```

### 15.8.2 The Objective Function fun

The objective function `fun` from the class `HyperLight` [SOURCE] is selected next. It implements an interface from PyTorch's training, validation, and testing methods to `spotPython`.

```
from spotPython.fun.hyperlight import HyperLight
fun = HyperLight(log_level=50).fun
```

### 15.8.3 Showing the fun\_control Dictionary

```
import pprint
pprint.pprint(fun_control)
```

```
{'CHECKPOINT_PATH': 'runs/saved_models/',
'DATASET_PATH': 'data/',
'RESULTS_PATH': 'results/',
'TENSORBOARD_PATH': 'runs/',
'_L_in': 10,
'_L_out': 1,
'accelerator': 'auto',
'core_model': <class 'spotPython.light.regression.netlightregression.NetLightRegression'>,
'core_model_hyper_dict': {'act_fn': {'class_name': 'spotPython.torch.activation',
                                      'core_model_parameter_type': 'instance()',
                                      'default': 'ReLU',
                                      'levels': ['Sigmoid',
                                                 'ReLU',
                                                 'LeakyReLU',
                                                 'Swish'],
                                      'lower': 0,
```

```

        'transform': 'None',
        'type': 'factor',
        'upper': 3},
    'batch_size': {'default': 4,
                   'lower': 4,
                   'transform': 'transform_power_2_int',
                   'type': 'int',
                   'upper': 5},
    'dropout_prob': {'default': 0.01,
                      'lower': 0.01,
                      'transform': 'None',
                      'type': 'float',
                      'upper': 0.1},
    'epochs': {'default': 4,
                'lower': 9,
                'transform': 'transform_power_2_int',
                'type': 'int',
                'upper': 10},
    'initialization': {'core_model_parameter_type': 'str',
                        'default': 'Default',
                        'levels': ['Default',
                                   'Kaiming',
                                   'Xavier'],
                        'lower': 0,
                        'transform': 'None',
                        'type': 'factor',
                        'upper': 2},
    'l1': {'default': 3,
            'lower': 4,
            'transform': 'transform_power_2_int',
            'type': 'int',
            'upper': 6},
    'lr_mult': {'default': 1.0,
                 'lower': 0.5,
                 'transform': 'None',
                 'type': 'float',
                 'upper': 5.0},
    'optimizer': {'class_name': 'torch.optim',
                  'core_model_parameter_type': 'str',
                  'default': 'SGD',
                  'levels': ['Adadelta',
                             'Adagrad',
                             'Adam'],

```

```

        'AdamW',
        'Adamax',
        'NAdam',
        'RAdam',
        'RMSprop',
        'Rprop'],
    'lower': 0,
    'transform': 'None',
    'type': 'factor',
    'upper': 8},
    'patience': {'default': 2,
                  'lower': 5,
                  'transform': 'transform_power_2_int',
                  'type': 'int',
                  'upper': 7}},

'counter': 0,
'data': None,
'data_dir': './data',
'data_module': None,
'data_set': <spotPython.data.diabetes.Diabetes object at 0x2d6a00410>,
'design': None,
'device': 'mps',
'devices': 1,
'enable_progress_bar': False,
'eval': None,
'fun_evals': inf,
'fun_repeats': 1,
'infill_criterion': 'y',
'k_folds': 3,
'log_level': 10,
'loss_function': None,
'lower': array([3. , 4. , 1. , 0. , 0. , 0. , 0.1, 2. , 0. ]),
'max_time': 1,
'metric_params': {},
'metric_river': None,
'metric_sklearn': None,
'metric_torch': None,
'model_dict': {},
'n_points': 1,
'n_samples': None,
'noise': False,
'num_workers': 0,
'ocba_delta': 0,

```

```
'optimizer': None,
'path': None,
'prep_model': None,
'save_model': False,
'seed': 123,
'show_batch_interval': 1000000,
'show_models': False,
'show_progress': True,
'shuffle': None,
'sigma': 0.0,
'spot_tensorboard_path': 'runs/spot_logs/031_p040025_2024-01-09_20-26-16',
'spot_writer': <torch.utils.tensorboard.writer.SummaryWriter object at 0x2d1442a10>,
'target_column': None,
'task': None,
'test': None,
'test_seed': 1234,
'test_size': 0.1,
'tolerance_x': 1.4901161193847656e-08,
'train': None,
'upper': array([ 8. ,  9. ,  4. ,  5. , 11. ,  0.25, 10. ,  6. ,  2. ]),
'ver_name': ['l1',
             'epochs',
             'batch_size',
             'act_fn',
             'optimizer',
             'dropout_prob',
             'lr_mult',
             'patience',
             'initialization'],
'ver_type': ['int',
             'int',
             'int',
             'factor',
             'factor',
             'float',
             'float',
             'int',
             'factor'],
'verbosity': 0,
'weights': 1.0}
```

#### 15.8.4 Starting the Hyperparameter Tuning

The `spotPython` hyperparameter tuning is started by calling the `Spot` function [SOURCE].

```
from spotPython.spot import spot
spot_tuner = spot.Spot(fun=fun,
                       fun_control=fun_control,
                       design_control=design_control,
                       surrogate_control=surrogate_control)
spot_tuner.run()
```

```
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 32
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 32
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 3704.84814453125, 'hp_metric': 3704.84814453125}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 32
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 32
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 3881.30419921875, 'hp_metric': 3881.30419921875}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 16
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 16
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 2666.55078125, 'hp_metric': 2666.55078125}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 16
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 16
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 5838.70849609375, 'hp_metric': 5838.70849609375}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 32
```

```

LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 32
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 3217.487060546875, 'hp_metric': 3217.487060546875}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 16
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 16
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 3510.2412109375, 'hp_metric': 3510.2412109375}
spotPython tuning: 2666.55078125 [#####----] 36.22%
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 16
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 16
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 4158.83154296875, 'hp_metric': 4158.83154296875}
spotPython tuning: 2666.55078125 [#####----] 77.03%
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 16
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 16
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 4188.34912109375, 'hp_metric': 4188.34912109375}
spotPython tuning: 2666.55078125 [#####----] 100.00% Done...

```

Validate metric	DataLoader 0
hp_metric	3704.84814453125
val_loss	3704.84814453125

Validate metric	DataLoader 0
hp_metric	3881.30419921875
val_loss	3881.30419921875

```
Validate metric           DataLoader 0
```

hp_metric	2666.55078125
val_loss	2666.55078125

```
Validate metric           DataLoader 0
```

hp_metric	5838.70849609375
val_loss	5838.70849609375

```
Validate metric           DataLoader 0
```

hp_metric	3217.487060546875
val_loss	3217.487060546875

```
Validate metric           DataLoader 0
```

hp_metric	3510.2412109375
val_loss	3510.2412109375

```
Validate metric           DataLoader 0
```

hp_metric	4158.83154296875
val_loss	4158.83154296875

```
Validate metric           DataLoader 0
```

hp_metric	4188.34912109375
val_loss	4188.34912109375

```
<spotPython.spot.spot.Spot at 0x2d7222790>
```

## 15.9 Step 9: Tensorboard

The textual output shown in the console (or code cell) can be visualized with Tensorboard.

```
tensorboard --logdir="runs/"
```

Further information can be found in the [PyTorch Lightning documentation](#) for Tensorboard.

## 15.10 Step 10: Results

After the hyperparameter tuning run is finished, the results can be analyzed.

```
spot_tuner.plot_progress(log_y=False,  
                         filename=".//figures/" + PREFIX + "_progress.png")
```

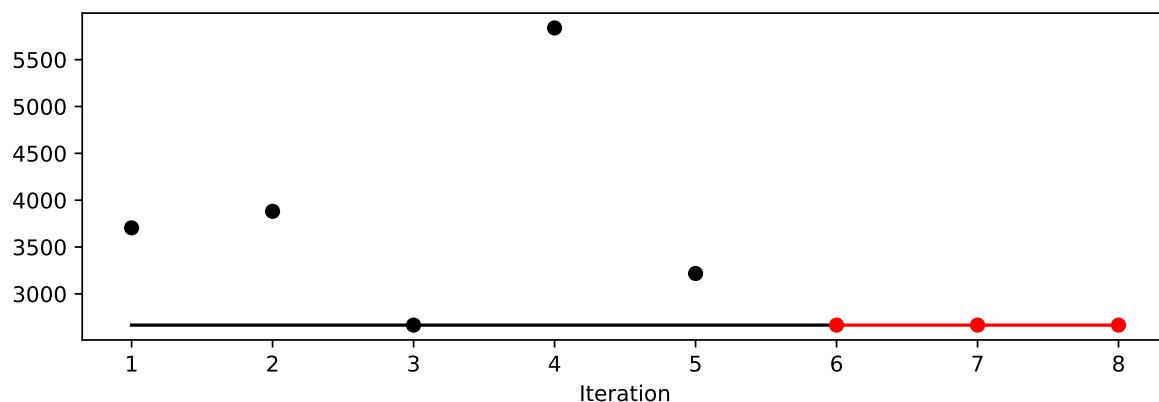


Figure 15.1: Progress plot. *Black* dots denote results from the initial design. *Red* dots illustrate the improvement found by the surrogate model based optimization.

```
from spotPython.utils.eda import gen_design_table  
print(gen_design_table(fun_control=fun_control, spot=spot_tuner))
```

name	type	default	lower	upper	tuned	transform
l1	int	3	4.0	6.0	5.0	transform_l
epochs	int	4	9.0	10.0	9.0	transform_l
batch_size	int	4	4.0	5.0	4.0	transform_l

act_fn	factor	ReLU	0.0	3.0	Swish	None
optimizer	factor	SGD	0.0	8.0	RMSprop	None
dropout_prob	float	0.01	0.01	0.1	0.08951095646892425	None
lr_mult	float	1.0	0.5	5.0	3.9340958113143683	None
patience	int	2	5.0	7.0	7.0	transform_
initialization	factor	Default	0.0	2.0	Xavier	None

```
spot_tuner.plot_importance(threshold=0.025,
    filename="./figures/" + PREFIX + "_importance.png")
```

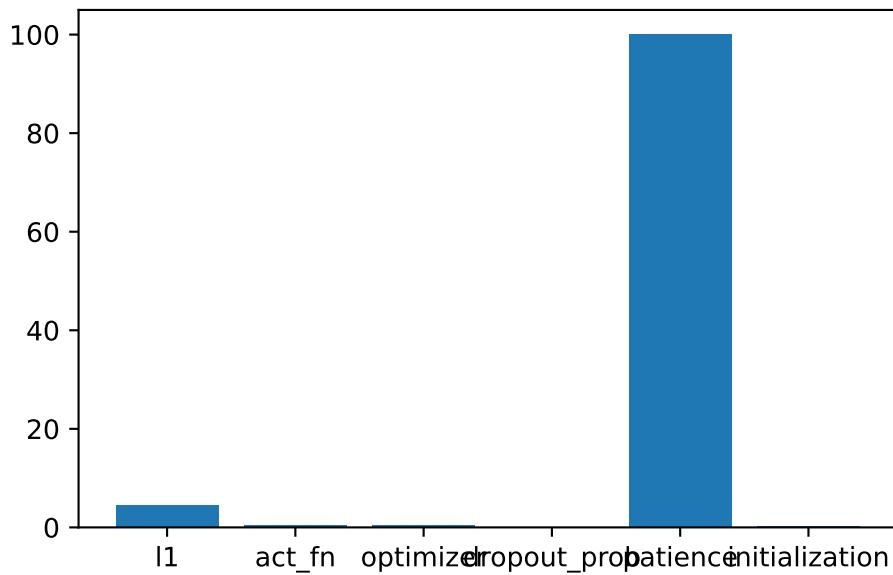


Figure 15.2: Variable importance plot, threshold 0.025.

### 15.10.1 Get the Tuned Architecture

```
from spotPython.hyperparameters.values import get_tuned_architecture
config = get_tuned_architecture(spot_tuner, fun_control)
print(config)
```

```
{'l1': 32, 'epochs': 512, 'batch_size': 16, 'act_fn': Swish(), 'optimizer': 'RMSprop', 'dropo
```

- Test on the full data set

```
from spotPython.light.testmodel import test_model
test_model(config, fun_control)
```

```
LightDataModule: train_dataloader(). Training set size: 71
LightDataModule: train_dataloader(). batch_size: 16
LightDataModule: train_dataloader(). num_workers: 0
LightDataModule: test_dataloader(). Training set size: 266
LightDataModule: test_dataloader(). batch_size: 16
LightDataModule: test_dataloader(). num_workers: 0
test_model result: {'val_loss': 5688.5185546875, 'hp_metric': 5688.5185546875}
```

Test metric	DataLoader 0
hp_metric	5688.5185546875
val_loss	5688.5185546875

(5688.5185546875, 5688.5185546875)

```
from spotPython.light.loadmodel import load_light_from_checkpoint
model_loaded = load_light_from_checkpoint(config, fun_control)
```

```
config: {'l1': 32, 'epochs': 512, 'batch_size': 16, 'act_fn': Swish(), 'optimizer': 'RMSprop'}
Loading model with 32_512_16_Swish_RMSprop_0.0895_3.9341_128_Xavier_TEST from runs/saved_models
Model: NetLightRegression(
(layers): Sequential(
(0): Linear(in_features=10, out_features=32, bias=True)
(1): Swish()
(2): Dropout(p=0.08951095646892425, inplace=False)
(3): Linear(in_features=32, out_features=16, bias=True)
(4): Swish()
(5): Dropout(p=0.08951095646892425, inplace=False)
(6): Linear(in_features=16, out_features=16, bias=True)
(7): Swish()
(8): Dropout(p=0.08951095646892425, inplace=False)
(9): Linear(in_features=16, out_features=8, bias=True)
(10): Swish()
(11): Dropout(p=0.08951095646892425, inplace=False)
```

```

(12): Linear(in_features=8, out_features=1, bias=True)
)
)

filename = "./figures/" + PREFIX
spot_tuner.plot_important_hyperparameter_contour(filename=filename)

```

l1: 4.540266379877663  
 act\_fn: 0.3898508891447861  
 optimizer: 0.45864634455539977  
 dropout\_prob: 0.07111682213320393  
 patience: 100.0  
 initialization: 0.1790210228848177

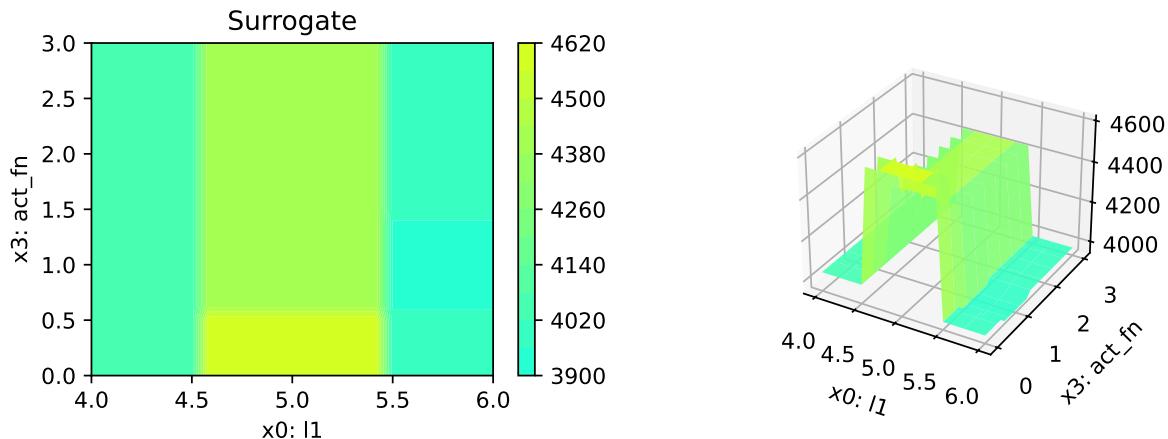
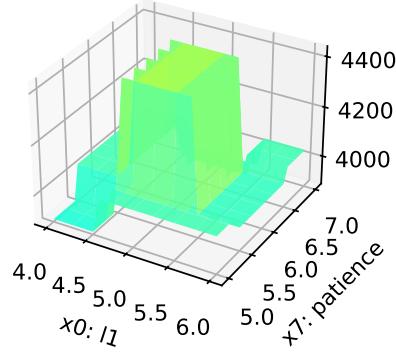
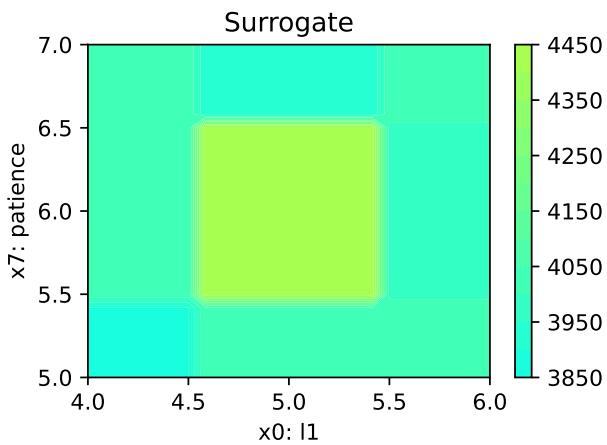
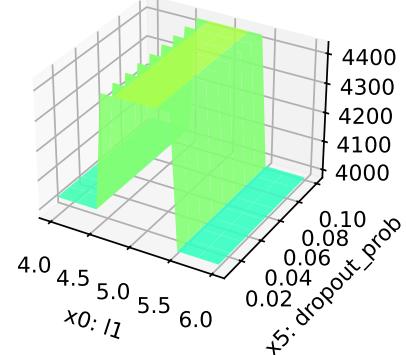
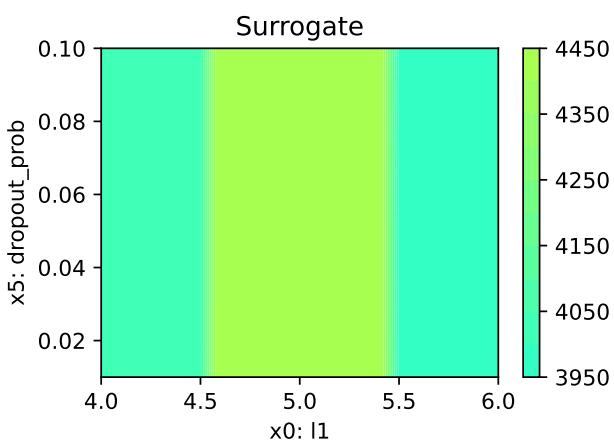
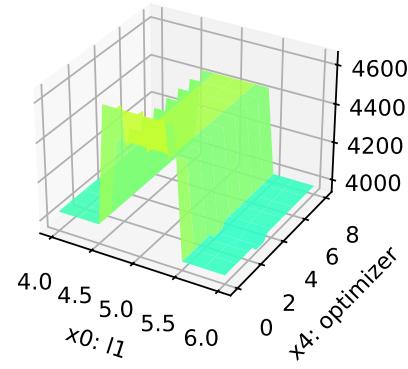
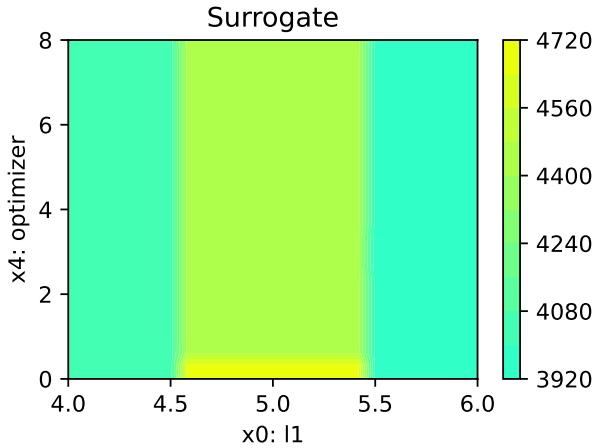
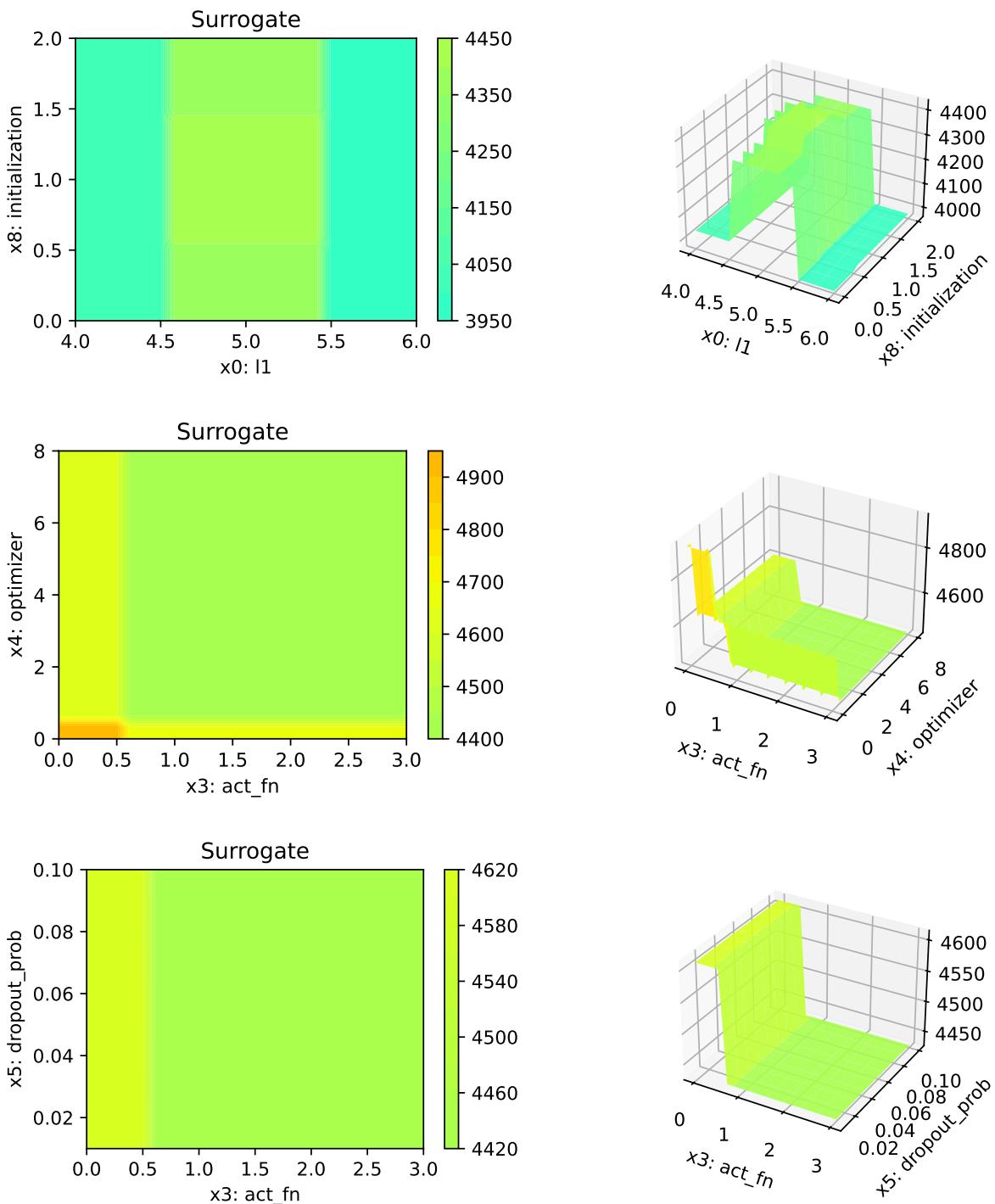
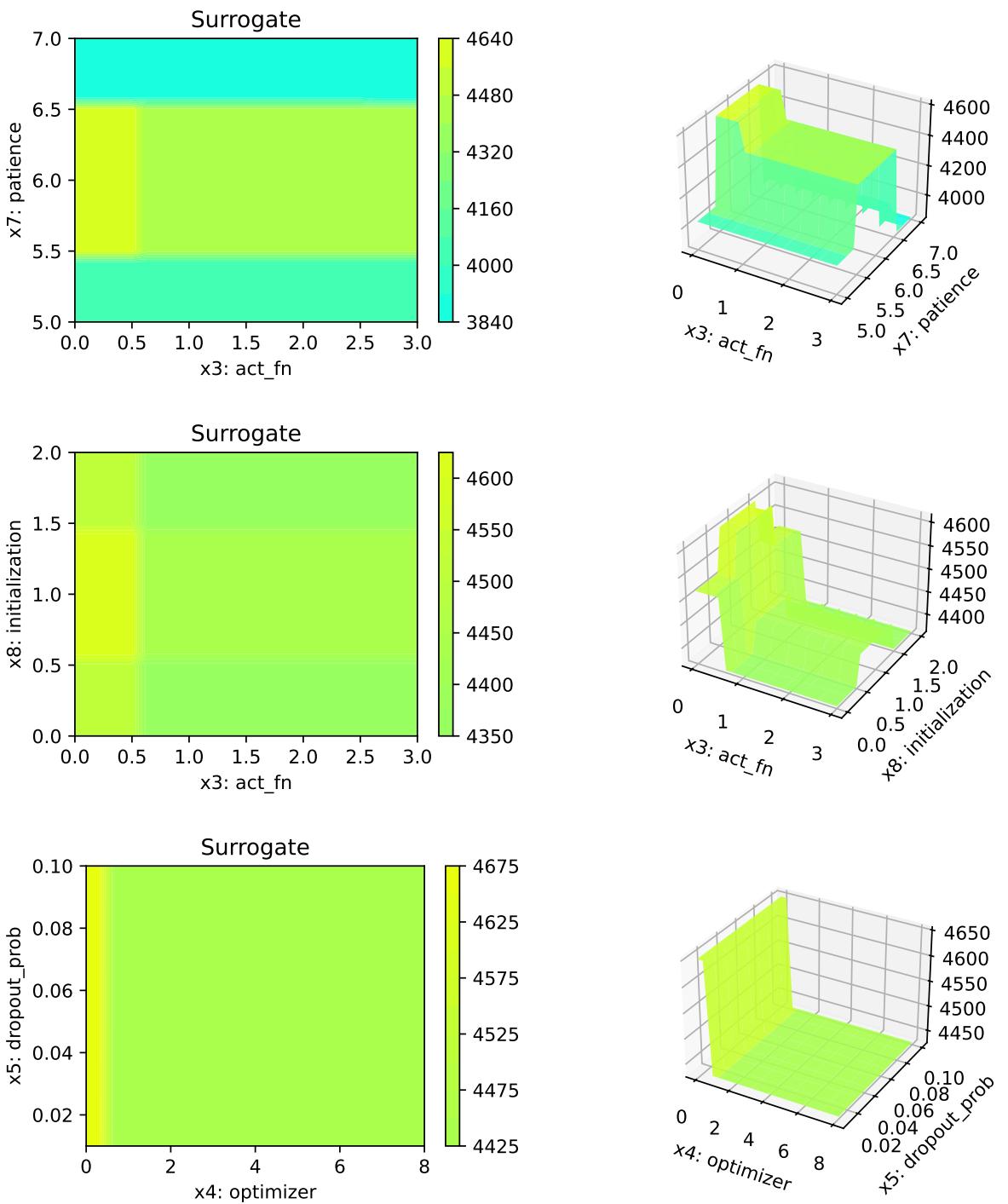
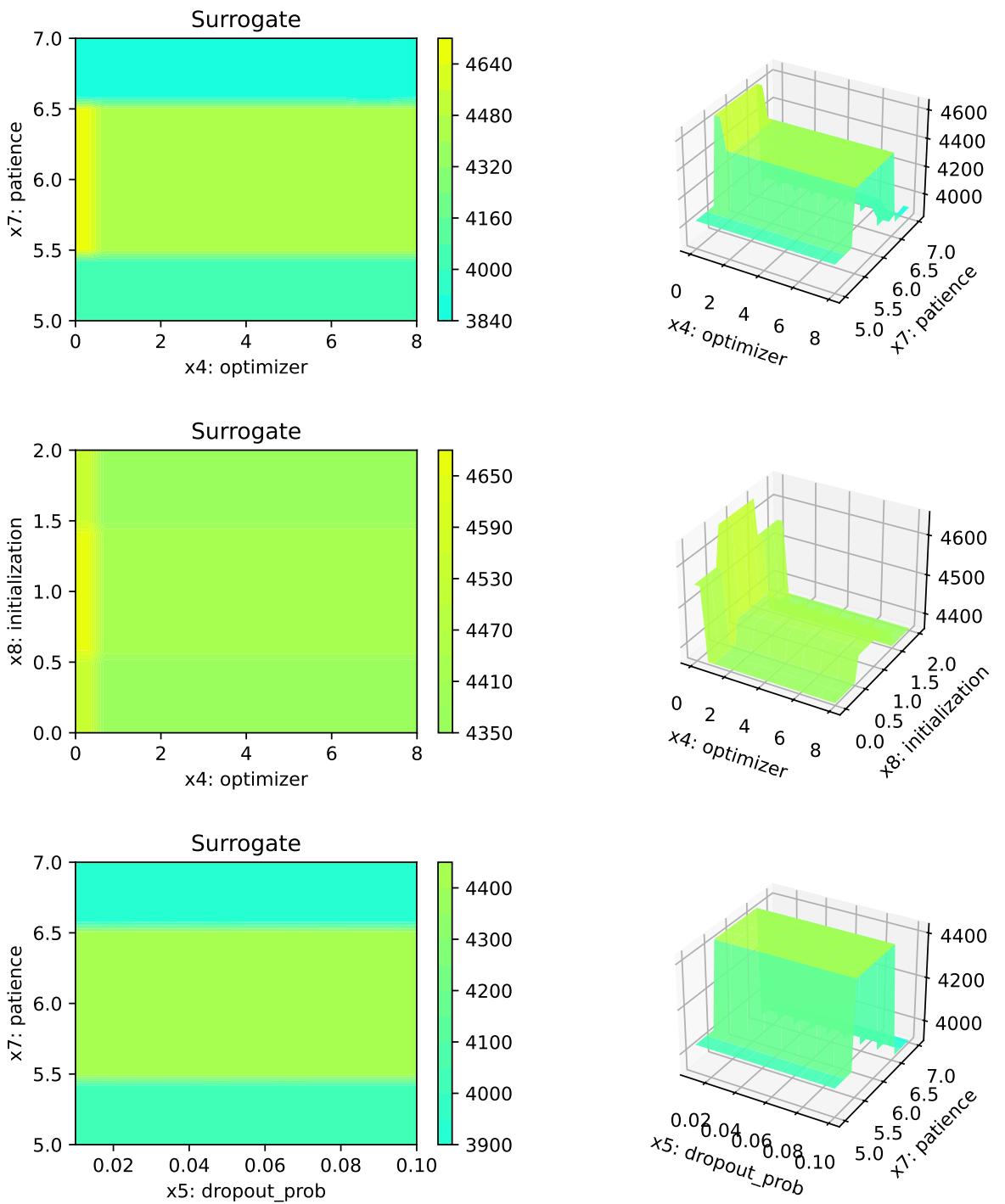


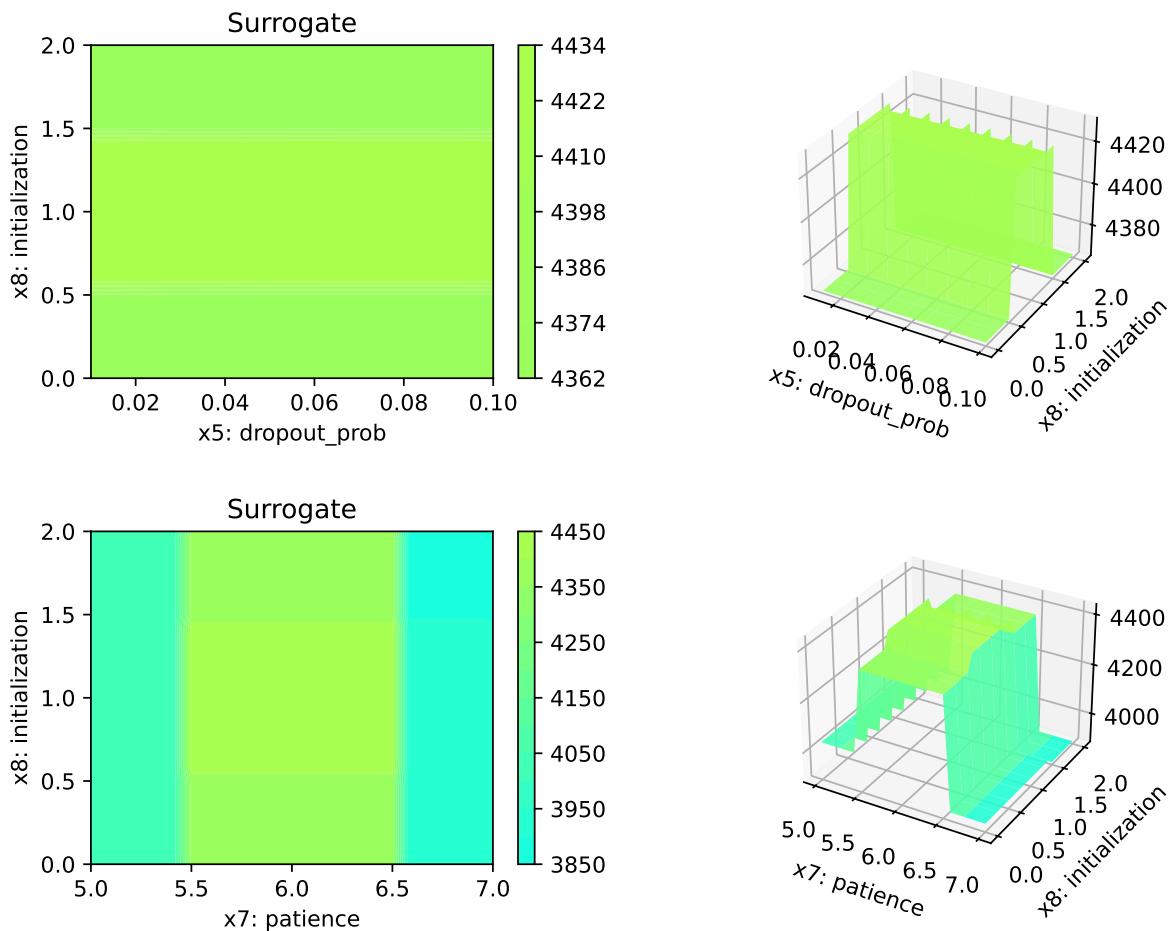
Figure 15.3: Contour plots.











### 15.10.2 Parallel Coordinates Plot

```
spot_tuner.parallel_plot()
```

Unable to display output for mime type(s): text/html

Parallel coordinates plots

Unable to display output for mime type(s): text/html

### 15.10.3 Cross Validation With Lightning

- The KFold class from `sklearn.model_selection` is used to generate the folds for cross-validation.
- These mechanism is used to generate the folds for the final evaluation of the model.
- The `CrossValidationDataModule` class [SOURCE] is used to generate the folds for the hyperparameter tuning process.
- It is called from the `cv_model` function [SOURCE].

```
from spotPython.light.cvmodel import cv_model
set_control_key_value(control_dict=fun_control,
                      key="k_folds",
                      value=2,
                      replace=True)
set_control_key_value(control_dict=fun_control,
                      key="test_size",
                      value=0.6,
                      replace=True)
cv_model(config, fun_control)
```

```
k: 0
Train Dataset Size: 221
Val Dataset Size: 221
train_model result: {'val_loss': 4612.650390625, 'hp_metric': 4612.650390625}
k: 1
Train Dataset Size: 221
Val Dataset Size: 221
train_model result: {'val_loss': 3470.23681640625, 'hp_metric': 3470.23681640625}
```

Validate metric	DataLoader 0
hp_metric	4612.650390625
val_loss	4612.650390625

Validate metric	DataLoader 0
hp_metric	3470.23681640625
val_loss	3470.23681640625

4041.443603515625

#### 15.10.4 Plot all Combinations of Hyperparameters

- Warning: this may take a while.

```
PLOT_ALL = False
if PLOT_ALL:
    n = spot_tuner.k
    for i in range(n-1):
        for j in range(i+1, n):
            spot_tuner.plot_contour(i=i, j=j, min_z=min_z, max_z = max_z)
```

#### 15.10.5 Visualizing the Activation Distribution (Under Development)

 Reference:

- The following code is based on [\[PyTorch Lightning TUTORIAL 2: ACTIVATION FUNCTIONS\]](#), Author: Phillip Lippe, License: [\[CC BY-SA\]](#), Generated: 2023-03-15T09:52:39.179933.

After we have trained the models, we can look at the actual activation values that find inside the model. For instance, how many neurons are set to zero in ReLU? Where do we find most values in Tanh? To answer these questions, we can write a simple function which takes a trained model, applies it to a batch of images, and plots the histogram of the activations inside the network:

```
from spotPython.torch.activation import Sigmoid, Tanh, ReLU, LeakyReLU, ELU, Swish
act_fn_by_name = {"sigmoid": Sigmoid, "tanh": Tanh, "relu": ReLU, "leakyrelu": LeakyReLU, "elu": ELU}
```

```
from spotPython.hyperparameters.values import get_one_config_from_X
X = spot_tuner.to_all_dim(spot_tuner.min_X.reshape(1,-1))
config = get_one_config_from_X(X, fun_control)
model = fun_control["core_model"](**config, _L_in=64, _L_out=11)
model
```

```
NetLightRegression(
(layers): Sequential(
(0): Linear(in_features=64, out_features=32, bias=True)
(1): Swish()
```

```
(2): Dropout(p=0.08951095646892425, inplace=False)
(3): Linear(in_features=32, out_features=16, bias=True)
(4): Swish()
(5): Dropout(p=0.08951095646892425, inplace=False)
(6): Linear(in_features=16, out_features=16, bias=True)
(7): Swish()
(8): Dropout(p=0.08951095646892425, inplace=False)
(9): Linear(in_features=16, out_features=8, bias=True)
(10): Swish()
(11): Dropout(p=0.08951095646892425, inplace=False)
(12): Linear(in_features=8, out_features=11, bias=True)
)
)

# from spotPython.utils.eda import visualize_activations
# visualize_activations(model, color=f"C{0}")
```

# 16 HPT PyTorch Lightning: Diabetes Using a Recurrent Neural Network

In this tutorial, we will show how `spotPython` can be integrated into the PyTorch Lightning training workflow for a regression task.

This chapter describes the hyperparameter tuning of a PyTorch Lightning network on the Diabetes data set. This is a PyTorch Dataset for regression. A toy data set from scikit-learn. Ten baseline variables, age, sex, body mass index, average blood pressure, and six blood serum measurements were obtained for each of  $n = 442$  diabetes patients, as well as the response of interest, a quantitative measure of disease progression one year after baseline.

## 16.1 Step 1: Setup

- Before we consider the detailed experimental setup, we select the parameters that affect run time, initial design size, etc.
- The parameter `MAX_TIME` specifies the maximum run time in seconds.
- The parameter `INIT_SIZE` specifies the initial design size.
- The parameter `WORKERS` specifies the number of workers.
- The prefix `PREFIX` is used for the experiment name and the name of the log file.
- The parameter `DEVICE` specifies the device to use for training.

```
from spotPython.utils.device import getDevice
from math import inf
MAX_TIME = 1
FUN_EVALS = inf
INIT_SIZE = 5
WORKERS = 0
PREFIX="032"
DEVICE = getDevice()
```



Caution: Run time and initial design size should be increased for real experiments

- MAX\_TIME is set to one minute for demonstration purposes. For real experiments, this should be increased to at least 1 hour.
- FUN\_EVALS is set to infinity.
- INIT\_SIZE is set to 5 for demonstration purposes. For real experiments, this should be increased to at least 10.
- WORKERS is set to 0 for demonstration purposes. For real experiments, this should be increased. See the warnings that are printed when the number of workers is set to 0.
- PREFIX is set to “032”. This is used for the experiment name and the name of the log file.
- DEVICE is set to the device that is returned by `getDevice()`, e.g., `gpu`.



Note: Device selection

- Although there are no `.cuda()` or `.to(device)` calls required, because Lightning does these for you, see [LIGHTNINGMODULE](#), we would like to know which device is used. Therefore, we imitate the LightningModule behaviour which selects the highest device.
- The method `spotPython.utils.device.getDevice()` returns the device that is used by Lightning.

## 16.2 Step 2: Initialization of the `fun_control` Dictionary

`spotPython` uses a Python dictionary for storing the information required for the hyperparameter tuning process.

```
from spotPython.utils.init import fun_control_init
import numpy as np

fun_control = fun_control_init(
    _L_in=10,
    _L_out=1,
    PREFIX=PREFIX,
    TENSORBOARD_CLEAN=True,
    device=DEVICE,
    enable_progress_bar=False,
    fun_evals=FUN_EVALS,
    log_level=10,
```

```
max_time=MAX_TIME,
num_workers=WORKERS,
show_progress=True,
test_size=0.1,
tolerance_x=np.sqrt(np.spacing(1)),
verbosity=1
)
```

```
Moving TENSORBOARD_PATH: runs/ to TENSORBOARD_PATH_OLD: runs_OLD/runs_2024_01_09_20_43_24
Created spot_tensorboard_path: runs/spot_logs/032_p040025_2024-01-09_20-43-24 for SummaryWriter
```

## 16.3 Step 3: Loading the Diabetes Data Set

```
from spotPython.hyperparameters.values import set_control_key_value
from spotPython.data.diabetes import Diabetes
dataset = Diabetes()
set_control_key_value(control_dict=fun_control,
                      key="data_set",
                      value=dataset,
                      replace=True)
print(len(dataset))
```

442

**i** Note: Data Set and Data Loader

- As shown below, a DataLoader from `torch.utils.data` can be used to check the data.

```

# Set batch size for DataLoader
batch_size = 5
# Create DataLoader
from torch.utils.data import DataLoader
dataloader = DataLoader(dataset, batch_size=batch_size, shuffle=False)

# Iterate over the data in the DataLoader
for batch in dataloader:
    inputs, targets = batch
    print(f"Batch Size: {inputs.size(0)}")
    print(f"Inputs Shape: {inputs.shape}")
    print(f"Targets Shape: {targets.shape}")
    print("-----")
    print(f"Inputs: {inputs}")
    print(f"Targets: {targets}")
    break

Batch Size: 5
Inputs Shape: torch.Size([5, 10])
Targets Shape: torch.Size([5])
-----
Inputs: tensor([[ 0.0381,  0.0507,  0.0617,  0.0219, -0.0442, -0.0348, -0.0434, -0.0026,
                 0.0199, -0.0176],
               [-0.0019, -0.0446, -0.0515, -0.0263, -0.0084, -0.0192,  0.0744, -0.0395,
                -0.0683, -0.0922],
               [ 0.0853,  0.0507,  0.0445, -0.0057, -0.0456, -0.0342, -0.0324, -0.0026,
                 0.0029, -0.0259],
               [-0.0891, -0.0446, -0.0116, -0.0367,  0.0122,  0.0250, -0.0360,  0.0343,
                 0.0227, -0.0094],
               [ 0.0054, -0.0446, -0.0364,  0.0219,  0.0039,  0.0156,  0.0081, -0.0026,
                 -0.0320, -0.0466]]])
Targets: tensor([151.,  75., 141., 206., 135.])

```

## 16.4 Step 4: Preprocessing

Preprocessing is handled by Lightning and PyTorch. It is described in the [LIGHTNING-DATAMODULE](#) documentation. Here you can find information about the `transforms` methods.

## 16.5 Step 5: Select the Core Model (algorithm) and core\_model\_hyper\_dict

spotPython includes the `NetLightRegression` class [SOURCE] for configurable neural networks. The class is imported here. It inherits from the class `Lightning.LightningModule`, which is the base class for all models in `Lightning`. `Lightning.LightningModule` is a subclass of `torch.nn.Module` and provides additional functionality for the training and testing of neural networks. The class `Lightning.LightningModule` is described in the [Lightning documentation](#).

- Here we simply add the NN Model to the `fun_control` dictionary by calling the function `add_core_model_to_fun_control`:

```
from spotPython.light.regression.rnnlightregression import RNNLightRegression
from spotPython.hyperdict.light_hyper_dict import LightHyperDict
from spotPython.hyperparameters.values import add_core_model_to_fun_control
add_core_model_to_fun_control(fun_control=fun_control,
                             core_model=RNNLightRegression,
                             hyper_dict=LightHyperDict)
```

The hyperparameters of the model are specified in the `core_model_hyper_dict` dictionary [SOURCE].

 Note: User specified models and hyperparameter dictionaries

- The user can specify a model and a hyperparameter dictionary in a subfolder, e.g., `userRNN` in the current working directory.
- The model and the hyperparameter dictionary are imported with the following code:

```
from spotPython.hyperparameters.values import add_core_model_to_fun_control
import sys
sys.path.insert(0, './userRNN')
import userrnn
import user_hyper_dict
add_core_model_to_fun_control(fun_control=fun_control,
                             core_model=userrnn.RNNLightRegression,
                             hyper_dict=user_hyper_dict.UserHyperDict)
```

- Example files can be found in the `userRNN` folder.
- These files can be modified by the user.
- They can be used without re-compilation of the `spotPython` source code, if they

are located in a subfolder of the current working directory.

## 16.6 Step 6: Modify `hyper_dict` Hyperparameters for the Selected Algorithm aka `core_model`

`spotPython` provides functions for modifying the hyperparameters, their bounds and factors as well as for activating and de-activating hyperparameters without re-compilation of the Python source code.

 Caution: Small number of epochs for demonstration purposes

- `epochs` and `patience` are set to small values for demonstration purposes. These values are too small for a real application.
- More resonable values are, e.g.:
  - `set_control_hyperparameter_value(fun_control, "epochs", [7, 9])`  
and
  - `set_control_hyperparameter_value(fun_control, "patience", [2, 7])`

```
from spotPython.hyperparameters.values import set_control_hyperparameter_value

set_control_hyperparameter_value(fun_control, "l1", [3, 8])
set_control_hyperparameter_value(fun_control, "epochs", [7, 9])
set_control_hyperparameter_value(fun_control, "batch_size", [2, 6])
set_control_hyperparameter_value(fun_control, "optimizer", [
    "Adadelta",
    "Adagrad",
    "Adam",
    "Adamax"])
set_control_hyperparameter_value(fun_control, "dropout_prob", [0.01, 0.25])
set_control_hyperparameter_value(fun_control, "lr_mult", [0.5, 5.0])
set_control_hyperparameter_value(fun_control, "patience", [3, 9])
set_control_hyperparameter_value(fun_control, "act_fn", ["ReLU"])
set_control_hyperparameter_value(fun_control, "initialization", ["Default"])
```

Now, the dictionary `fun_control` contains all information needed for the hyperparameter tuning. Before the hyperparameter tuning is started, it is recommended to take a look at the experimental design. The method `gen_design_table` [\[SOURCE\]](#) generates a design table as follows:

```
from spotPython.utils.eda import gen_design_table
print(gen_design_table(fun_control))
```

name	type	default	lower	upper	transform
l1	int	3	3	8	transform_power_2_int
epochs	int	4	7	9	transform_power_2_int
batch_size	int	4	2	6	transform_power_2_int
act_fn	factor	ReLU	0	0	None
optimizer	factor	SGD	0	3	None
dropout_prob	float	0.01	0.01	0.25	None
lr_mult	float	1.0	0.5	5	None
patience	int	2	3	9	transform_power_2_int
initialization	factor	Default	0	0	None

This allows to check if all information is available and if the information is correct.

**i** Note: Hyperparameters of the Tuned Model and the `fun_control` Dictionary

The updated `fun_control` dictionary can be shown with the command `fun_control["core_model_hyper_dict"]`.

## 16.7 Step 7: Data Splitting, the Objective (Loss) Function and the Metric

### 16.7.1 Evaluation

The evaluation procedure requires the specification of two elements:

1. the way how the data is split into a train and a test set
2. the loss function (and a metric).

**🔥 Caution:** Data Splitting in Lightning

The data splitting is handled by **Lightning**.

### 16.7.2 Loss Function

The loss function is specified in the configurable network class [\[SOURCE\]](#) We will use MSE.

### 16.7.3 Metric

- Similar to the loss function, the metric is specified in the configurable network class [\[SOURCE\]](#).

 Caution: Loss Function and Metric in Lightning

- The loss function and the metric are not hyperparameters that can be tuned with `spotPython`.
- They are handled by `Lightning`.

## 16.8 Step 8: Calling the SPOT Function

### 16.8.1 Preparing the SPOT Call

```
from spotPython.utils.init import design_control_init, surrogate_control_init
design_control = design_control_init()
set_control_key_value(control_dict=design_control,
                      key="init_size",
                      value=INIT_SIZE,
                      replace=True)

surrogate_control = surrogate_control_init()
set_control_key_value(control_dict=surrogate_control,
                      key="noise",
                      value=True,
                      replace=True)
set_control_key_value(control_dict=surrogate_control,
                      key="n_theta",
                      value=2,
                      replace=True)
```

### 16.8.2 The Objective Function `fun`

The objective function `fun` from the class `HyperLight` [\[SOURCE\]](#) is selected next. It implements an interface from PyTorch's training, validation, and testing methods to `spotPython`.

```
from spotPython.fun.hyperlight import HyperLight
fun = HyperLight(log_level=10).fun
```

### 16.8.3 Showing the fun\_control Dictionary

```
import pprint
pprint.pprint(fun_control)
```

```
{'CHECKPOINT_PATH': 'runs/saved_models/',
'DATASET_PATH': 'data/',
'RESULTS_PATH': 'results/',
'TENSORBOARD_PATH': 'runs/',
'_L_in': 10,
'_L_out': 1,
'accelerator': 'auto',
'core_model': <class 'spotPython.light.regression.rnnlightregression.RNNLightRegression'>,
'core_model_hyper_dict': {'act_fn': {'class_name': 'spotPython.torch.activation',
                                      'core_model_parameter_type': 'instance()',
                                      'default': 'ReLU',
                                      'levels': ['ReLU'],
                                      'lower': 0,
                                      'transform': 'None',
                                      'type': 'factor',
                                      'upper': 0},
                           'batch_size': {'default': 4,
                                         'lower': 2,
                                         'transform': 'transform_power_2_int',
                                         'type': 'int',
                                         'upper': 6},
                           'dropout_prob': {'default': 0.01,
                                            'lower': 0.01,
                                            'transform': 'None',
                                            'type': 'float',
                                            'upper': 0.25},
                           'epochs': {'default': 4,
                                       'lower': 7,
                                       'transform': 'transform_power_2_int',
                                       'type': 'int',
                                       'upper': 9},
                           'initialization': {'core_model_parameter_type': 'str',
```

```

        'default': 'Default',
        'levels': ['Default'],
        'lower': 0,
        'transform': 'None',
        'type': 'factor',
        'upper': 0},
    'l1': {'default': 3,
            'lower': 3,
            'transform': 'transform_power_2_int',
            'type': 'int',
            'upper': 8},
    'lr_mult': {'default': 1.0,
                'lower': 0.5,
                'transform': 'None',
                'type': 'float',
                'upper': 5.0},
    'optimizer': {'class_name': 'torch.optim',
                  'core_model_parameter_type': 'str',
                  'default': 'SGD',
                  'levels': ['Adadelta',
                             'Adagrad',
                             'Adam',
                             'Adamax'],
                  'lower': 0,
                  'transform': 'None',
                  'type': 'factor',
                  'upper': 3},
    'patience': {'default': 2,
                  'lower': 3,
                  'transform': 'transform_power_2_int',
                  'type': 'int',
                  'upper': 9}},
    'counter': 0,
    'data': None,
    'data_dir': './data',
    'data_module': None,
    'data_set': <spotPython.data.diabetes.Diabetes object at 0x2b1d6e5d0>,
    'design': None,
    'device': 'mps',
    'devices': 1,
    'enable_progress_bar': False,
    'eval': None,
    'fun_evals': inf,

```

```

'fun_repeats': 1,
'infill_criterion': 'y',
'k_folds': 3,
'log_level': 10,
'loss_function': None,
'lower': array([3. , 4. , 1. , 0. , 0. , 0. , 0.1, 2. , 0. ]),
'max_time': 1,
'metric_params': {},
'metric_river': None,
'metric_sklearn': None,
'metric_torch': None,
'model_dict': {},
'n_points': 1,
'n_samples': None,
'noise': False,
'num_workers': 0,
'ocba_delta': 0,
'optimizer': None,
'path': None,
'prep_model': None,
'save_model': False,
'seed': 123,
'show_batch_interval': 1000000,
'show_models': False,
'show_progress': True,
'shuffle': None,
'sigma': 0.0,
'spot_tensorboard_path': 'runs/spot_logs/032_p040025_2024-01-09_20-43-24',
'spot_writer': <torch.utils.tensorboard.SummaryWriter object at 0x2ae288490>,
'target_column': None,
'task': None,
'test': None,
'test_seed': 1234,
'test_size': 0.1,
'tolerance_x': 1.4901161193847656e-08,
'train': None,
'upper': array([ 8. ,  9. ,  4. ,  1. , 11. ,  0.25, 10. ,  6. ,  2. ]),
'var_name': ['l1',
             'epochs',
             'batch_size',
             'act_fn',
             'optimizer',
             'dropout_prob',

```

```

    'lr_mult',
    'patience',
    'initialization'],
'verbose': ['int',
            'int',
            'int',
            'factor',
            'factor',
            'float',
            'float',
            'int',
            'factor'],
'verbosity': 1,
'weights': 1.0}

pprint.pprint(design_control)

{'init_size': 5, 'repeats': 1}

pprint.pprint(surrogate_control)

{'max_theta': 3,
'min_theta': -3,
'model_fun_evals': 10000,
'model_optimizer': <function differential_evolution at 0x175bb4040>,
'n_p': 1,
'n_theta': 2,
'noise': True,
'optim_p': False,
'seed': 124,
'theta_init_zero': True,
'var_type': None}

```

#### 16.8.4 Starting the Hyperparameter Tuning

The `spotPython` hyperparameter tuning is started by calling the `Spot` function [\[SOURCE\]](#).

```

from spotPython.spot import spot
spot_tuner = spot.Spot(fun=fun,
                      fun_control=fun_control,

```

```
        design_control=design_control,
        surrogate_control=surrogate_control)
spot_tuner.run()
```

```
In fun(): config:
{'act_fn': ReLU(),
'batch_size': 64,
'dropout_prob': 0.19355651674791854,
'epochs': 256,
'initialization': 'Default',
'l1': 16,
'lr_mult': 1.5691149440098038,
'optimizer': 'Adam',
'patience': 32}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 64
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 64
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 5832.97998046875, 'hp_metric': 5832.97998046875}

In fun(): config:
{'act_fn': ReLU(),
'batch_size': 16,
'dropout_prob': 0.09424169914869776,
'epochs': 256,
'initialization': 'Default',
'l1': 128,
'lr_mult': 3.35818256351233,
'optimizer': 'Adadelta',
'patience': 512}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 16
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 16
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 3652.889404296875, 'hp_metric': 3652.889404296875}

In fun(): config:
```

```
{'act_fn': ReLU(),
 'batch_size': 4,
 'dropout_prob': 0.21164199382623602,
 'epochs': 512,
 'initialization': 'Default',
 'l1': 128,
 'lr_mult': 0.9336514668325573,
 'optimizer': 'Adamax',
 'patience': 16}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 4
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 4
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 2414.205078125, 'hp_metric': 2414.205078125}
```

```
In fun(): config:
{'act_fn': ReLU(),
 'batch_size': 8,
 'dropout_prob': 0.05728504399550885,
 'epochs': 128,
 'initialization': 'Default',
 'l1': 64,
 'lr_mult': 4.575980093998586,
 'optimizer': 'Adam',
 'patience': 32}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 8
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 8
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 5862.02099609375, 'hp_metric': 5862.02099609375}
```

```
In fun(): config:
{'act_fn': ReLU(),
 'batch_size': 16,
 'dropout_prob': 0.14352914208400058,
 'epochs': 256,
 'initialization': 'Default',
 'l1': 8,
 'lr_mult': 2.4204853123355816,
```

```

'optimizer': 'Adagrad',
'patience': 128}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 16
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 16
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 3681.71923828125, 'hp_metric': 3681.71923828125}

In fun(): config:
{'act_fn': ReLU(),
 'batch_size': 4,
 'dropout_prob': 0.21164827556769356,
 'epochs': 512,
 'initialization': 'Default',
 'l1': 128,
 'lr_mult': 0.9336469468487545,
 'optimizer': 'Adamax',
 'patience': 16}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 4
LightDataModule: train_dataloader(). Training set size: 359
LightDataModule: train_dataloader(). batch_size: 4
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 3021.425048828125, 'hp_metric': 3021.425048828125}
spotPython tuning: 2414.205078125 [#####--] 80.86%

In fun(): config:
{'act_fn': ReLU(),
 'batch_size': 4,
 'dropout_prob': 0.21160083967554538,
 'epochs': 512,
 'initialization': 'Default',
 'l1': 128,
 'lr_mult': 0.9336847184495153,
 'optimizer': 'Adamax',
 'patience': 16}
train_model(): Test set size: 45
train_model(): Train set size: 359
train_model(): Batch size: 4
LightDataModule: train_dataloader(). Training set size: 359

```

```
LightDataModule: train_dataloader(). batch_size: 4
LightDataModule: train_dataloader(). num_workers: 0
train_model result: {'val_loss': 3326.468505859375, 'hp_metric': 3326.468505859375}
spotPython tuning: 2414.205078125 [#####] 100.00% Done...
```

Validate metric	DataLoader 0
hp_metric	5832.97998046875
val_loss	5832.97998046875

Validate metric	DataLoader 0
hp_metric	3652.889404296875
val_loss	3652.889404296875

Validate metric	DataLoader 0
hp_metric	2414.205078125
val_loss	2414.205078125

Validate metric	DataLoader 0
hp_metric	5862.02099609375
val_loss	5862.02099609375

Validate metric	DataLoader 0
hp_metric	3681.71923828125
val_loss	3681.71923828125

```
Validate metric           DataLoader 0
hp_metric                3021.425048828125
val_loss                  3021.425048828125
```

```
Validate metric           DataLoader 0
hp_metric                3326.468505859375
val_loss                  3326.468505859375
```

```
<spotPython.spot.spot.Spot at 0x2b400b6d0>
```

## 16.9 Step 9: Tensorboard

The textual output shown in the console (or code cell) can be visualized with Tensorboard.

```
tensorboard --logdir="runs/"
```

Further information can be found in the [PyTorch Lightning documentation](#) for Tensorboard.

## 16.10 Step 10: Results

After the hyperparameter tuning run is finished, the results can be analyzed.

```
spot_tuner.plot_progress(log_y=False,
                         filename=".//figures/" + PREFIX + "_progress.png")
```

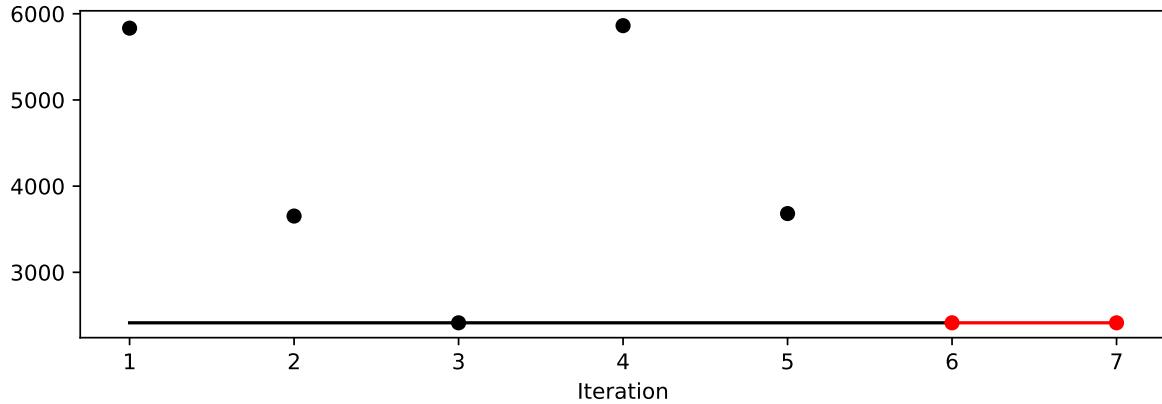


Figure 16.1: Progress plot. *Black* dots denote results from the initial design. *Red* dots illustrate the improvement found by the surrogate model based optimization.

```
from spotPython.utils.eda import gen_design_table
print(gen_design_table(fun_control=fun_control, spot=spot_tuner))
```

name	type	default	lower	upper	tuned	transform
l1	int	3	3.0	8.0	7.0	transform_l1
epochs	int	4	7.0	9.0	9.0	transform_l1
batch_size	int	4	2.0	6.0	2.0	transform_l1
act_fn	factor	ReLU	0.0	0.0	ReLU	None
optimizer	factor	SGD	0.0	3.0	Adamax	None
dropout_prob	float	0.01	0.01	0.25	0.21164199382623602	None
lr_mult	float	1.0	0.5	5.0	0.9336514668325573	None
patience	int	2	3.0	9.0	4.0	transform_l1
initialization	factor	Default	0.0	0.0	Default	None

```
spot_tuner.plot_importance(threshold=0.025,
                           filename=".//figures/" + PREFIX + "_importance.png")
```

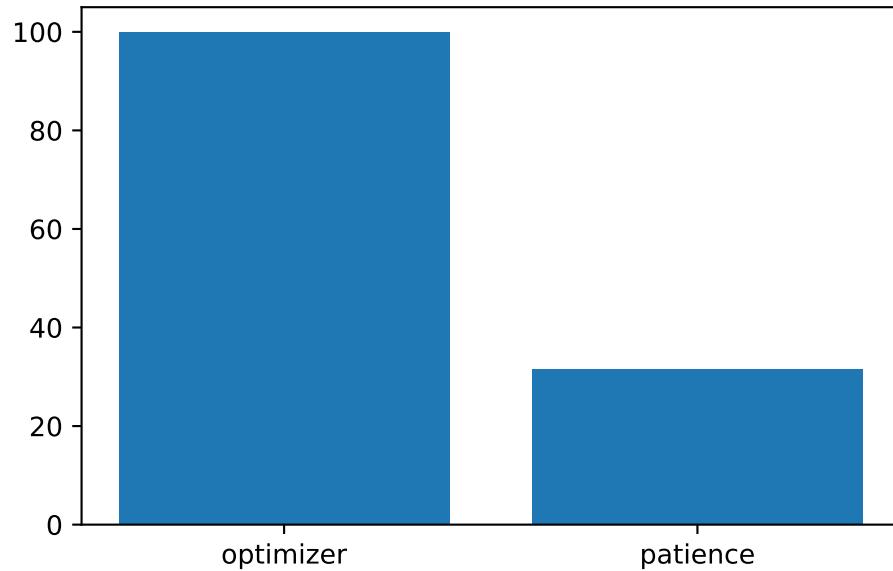


Figure 16.2: Variable importance plot, threshold 0.025.

### 16.10.1 Get the Tuned Architecture

```
from spotPython.hyperparameters.values import get_tuned_architecture
config = get_tuned_architecture(spot_tuner, fun_control)
print(config)
```

```
{'l1': 128, 'epochs': 512, 'batch_size': 4, 'act_fn': ReLU(), 'optimizer': 'Adamax', 'dropout': 0.2}
```

- Test on the full data set

```
from spotPython.light.testmodel import test_model
test_model(config, fun_control)
```

```
LightDataModule: train_dataloader(). Training set size: 71
LightDataModule: train_dataloader(). batch_size: 4
LightDataModule: train_dataloader(). num_workers: 0
LightDataModule: test_dataloader(). Training set size: 266
LightDataModule: test_dataloader(). batch_size: 4
LightDataModule: test_dataloader(). num_workers: 0
test_model result: {'val_loss': 3671.13623046875, 'hp_metric': 3671.13623046875}
```

Test metric	DataLoader 0
hp_metric	3671.13623046875
val_loss	3671.13623046875

(3671.13623046875, 3671.13623046875)

```
from spotPython.light.loadmodel import load_light_from_checkpoint
model_loaded = load_light_from_checkpoint(config, fun_control)
```

```
config: {'l1': 128, 'epochs': 512, 'batch_size': 4, 'act_fn': ReLU(), 'optimizer': 'Adamax',
Loading model with 128_512_4_ReLU_Adamax_0.2116_0.9337_16_Default_TEST from runs/saved_models
Model: RNNLightRegression(
    (rnn_layer): RNN(10, 128, batch_first=True)
    (fc): Linear(in_features=128, out_features=128, bias=True)
    (output_layer): Linear(in_features=128, out_features=1, bias=True)
    (dropout1): Dropout(p=0.21164199382623602, inplace=False)
    (dropout2): Dropout(p=0.0, inplace=False)
    (dropout3): Dropout(p=0.0, inplace=False)
    (activation_fct): ReLU()
)
```

```
filename = "./figures/" + PREFIX
spot_tuner.plot_important_hyperparameter_contour(filename=filename)
```

```
optimizer: 100.0
patience: 31.653018801863347
```

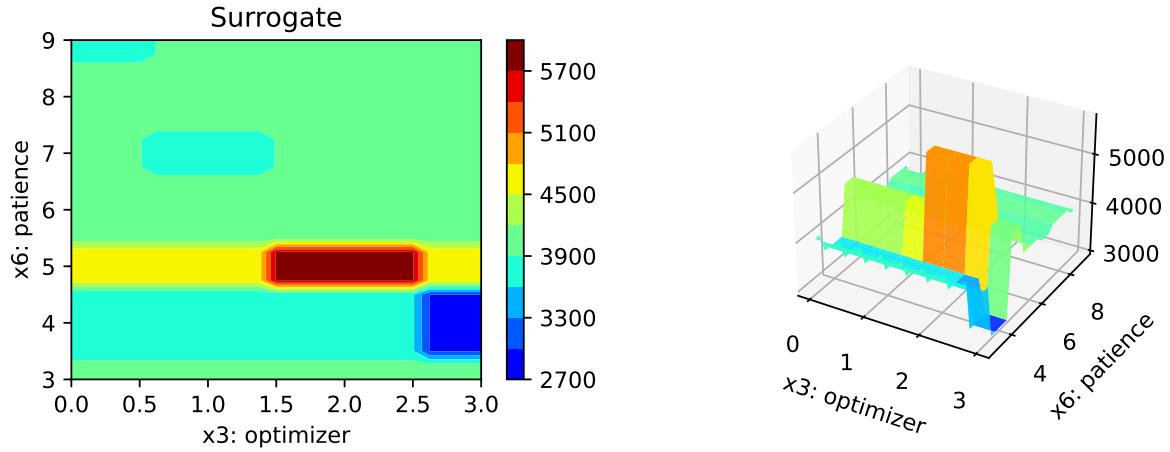


Figure 16.3: Contour plots.

### 16.10.2 Parallel Coordinates Plot

```
spot_tuner.parallel_plot()
```

Unable to display output for mime type(s): text/html

Parallel coordinates plots

Unable to display output for mime type(s): text/html

### 16.10.3 Cross Validation With Lightning

- The KFold class from `sklearn.model_selection` is used to generate the folds for cross-validation.
- These mechanism is used to generate the folds for the final evaluation of the model.
- The `CrossValidationDataModule` class [SOURCE] is used to generate the folds for the hyperparameter tuning process.
- It is called from the `cv_model` function [SOURCE].

```

from spotPython.light.cvmodel import cv_model
set_control_key_value(control_dict=fun_control,
                      key="k_folds",
                      value=2,
                      replace=True)
set_control_key_value(control_dict=fun_control,
                      key="test_size",
                      value=0.1,
                      replace=True)
cv_model(config, fun_control)

```

```

k: 0
Train Dataset Size: 221
Val Dataset Size: 221
train_model result: {'val_loss': 3144.087158203125, 'hp_metric': 3144.087158203125}
k: 1
Train Dataset Size: 221
Val Dataset Size: 221
train_model result: {'val_loss': 3290.426513671875, 'hp_metric': 3290.426513671875}

```

Validate metric	DataLoader 0
hp_metric	3144.087158203125
val_loss	3144.087158203125

Validate metric	DataLoader 0
hp_metric	3290.426513671875
val_loss	3290.426513671875

3217.2568359375

#### 16.10.4 Plot all Combinations of Hyperparameters

- Warning: this may take a while.

```

PLOT_ALL = False
if PLOT_ALL:
    n = spot_tuner.k
    for i in range(n-1):
        for j in range(i+1, n):
            spot_tuner.plot_contour(i=i, j=j, min_z=min_z, max_z = max_z)

```

### 16.10.5 Visualizing the Activation Distribution (Under Development)

 Reference:

- The following code is based on [\[PyTorch Lightning TUTORIAL 2: ACTIVATION FUNCTIONS\]](#), Author: Phillip Lippe, License: [\[CC BY-SA\]](#), Generated: 2023-03-15T09:52:39.179933.

After we have trained the models, we can look at the actual activation values that find inside the model. For instance, how many neurons are set to zero in ReLU? Where do we find most values in Tanh? To answer these questions, we can write a simple function which takes a trained model, applies it to a batch of images, and plots the histogram of the activations inside the network:

```

from spotPython.torch.activation import Sigmoid, Tanh, ReLU, LeakyReLU, ELU, Swish
act_fn_by_name = {"sigmoid": Sigmoid, "tanh": Tanh, "relu": ReLU, "leakyrelu": LeakyReLU, "elu": ELU, "swish": Swish}

```

```

from spotPython.hyperparameters.values import get_one_config_from_X
X = spot_tuner.to_all_dim(spot_tuner.min_X.reshape(1,-1))
config = get_one_config_from_X(X, fun_control)
model = fun_control["core_model"](**config, _L_in=64, _L_out=11)
model

```

```

RNNLightRegression(
    (rnn_layer): RNN(64, 128, batch_first=True)
    (fc): Linear(in_features=128, out_features=128, bias=True)
    (output_layer): Linear(in_features=128, out_features=11, bias=True)
    (dropout1): Dropout(p=0.21164199382623602, inplace=False)
    (dropout2): Dropout(p=0.0, inplace=False)
    (dropout3): Dropout(p=0.0, inplace=False)
    (activation_fct): ReLU()
)

```

```
# from spotPython.utils.eda import visualize_activations
# visualize_activations(model, color=f"C{0}")
```

# A Introduction to Jupyter Notebook

Jupyter Notebook is a widely used tool in the Data Science community. It is easy to use and the produced code can be run per cell. This has a huge advantage, because with other tools e.g. (pycharm, vscode, etc.) the whole script is executed. This can be a time consuming process, especially when working with huge data sets.

## A.1 Different Notebook cells

There are different cells that the notebook is currently supporting:

- code cells
- markdown cells
- raw cells

As a default, every cells in jupyter is set to “code”

### A.1.1 Code cells

The code cells are used to execute the code. They are following the logic of the chosen kernel. Therefore, it is important to keep in mind which programming language is currently used. Otherwise one might yield an error because of the wrong syntax.

The code cells are executed my be **Run** button (can be found in the header of the notebook).

### A.1.2 Markdown cells

The markdown cells are a usefull tool to comment the written code. Especially with the help of headers can the code be brought in a more readable format. If you are not familiar with the markdown syntax, you can find a usefull cheat sheet here: [Markdown Cheat Sheet](#)

### A.1.3 Raw cells

The “Raw NBConvert” cell type can be used to render different code formats into HTML or LaTeX by Sphinx. This information is stored in the notebook metadata and converted appropriately.

#### A.1.3.1 Usage

To select a desired format from within Jupyter, select the cell containing your special code and choose options from the following dropdown menus:

1. Select “Raw NBConvert”
2. Switch the Cell Toolbar to “Raw Cell Format” (The cell toolbar can be found under View)
3. Choose the appropriate “Raw NBConvert Format” within the cell

Data Science is fun

## A.2 Install Packages

Because python is a heavily used programming language, there are many different packages that can make your life easier. Sadly, there are only a few standard packages that are already included in your python environment. If you have the need to install a new package in your environment, you can simply do that by executing the following code snippet in a **code cell**

```
!pip install numpy
```

- The `!` is used to run the cell as a shell command
- `pip` is package manager for python packages.
- `numpy` is the package you want to install

**Hint:** It is often useful to restart the kernel after installing a package, otherwise loading the package could lead to an error.

## A.3 Load Packages

After successfully installing the package it is necessary to import them before you can work with them. The import of the packages is done in the following way:

```
import numpy as np
```

The imported packages are often abbreviated. This is because you need to specify where the function is coming from.

The most common abbreviations for data science packages are:

Table A.1: Abbreviations for data science packages

Abbreviation	Package	Import
np	numpy	import numpy as np
pd	pandas	import pandas as pd
plt	matplotlib	import matplotlib.pyplot as plt
px	plotly	import plotly.express as px
tf	tensorflow	import tensorflow as tf
sns	seaborn	import seaborn as sns
dt	datetime	import datetime as dt
pkl	pickle	import pickle as pkl

## A.4 Functions in Python

Because python is not using Semicolon's it is import to keep track of indentation in your code. The indentation works as a placeholder for the semicolons. This is especially important if your are defining loops, functions, etc. ...

**Example:** We are defining a function that calculates the squared sum of its input parameters

```
def squared_sum(x,y):  
    z = x**2 + y**2  
    return z
```

If you are working with something that needs indentation, it will be already done by the notebook.

**Hint:** Keep in mind that is good practice to use the *return* parameter. If you are not using *return* and a function has multiple paramaters that you would like to return, it will only return the last one defined.

## A.5 List of Useful Jupyter Notebook Shortcuts

Table A.2: List of useful Jupyter Notebook Shortcuts

Function	Keyboard Shortcut	Menu Tools
Save notebook	Esc + s	File → Save and Checkpoint
Create new Cell	Esc + a (above), Esc + b (below)	Insert → Cell above; Insert → Cell below
Run Cell	Ctrl + enter	Cell → Run Cell
Copy Cell	c	Copy Key
Paste Cell	v	Paste Key
Interrupt Kernel	Esc + i i	Kernel → Interrupt
Restart Kernel	Esc + 0 0	Kernel → Restart

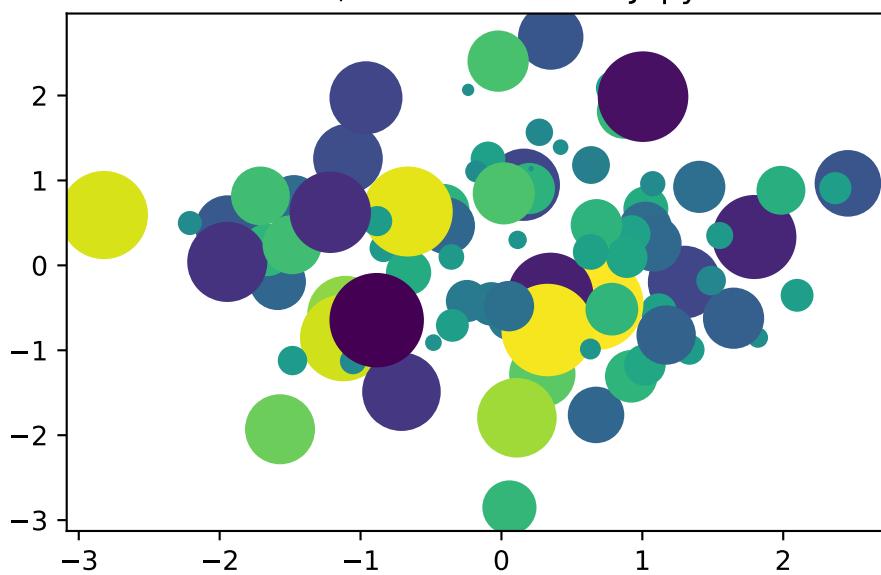
If you combine everything you can create beautiful graphics

```
import matplotlib.pyplot as plt
import numpy as np

# Generate 100 random data points along 3 dimensions
x, y, scale = np.random.randn(3, 100)
fig, ax = plt.subplots()

# Map each onto a scatterplot we'll create with Matplotlib
ax.scatter(x=x, y=y, c=scale, s=np.abs(scale)*500)
ax.set(title="Some random data, created with the Jupyter Notebook!")
plt.show()
```

Some random data, created with the Jupyter Notebook!



# B Git Introduction

## B.1 Learning Objectives

In this learning unit, you will learn how to set up Git as a version control system for a project. The most important Git commands will be explained. You will learn how to track and manage changes to your projects with Git. Specifically:

- Initializing a repository: `git init`
- Ignoring files: `.gitignore`
- Adding files to the staging area: `git add`
- Checking status changes: `git status`
- Reviewing history: `git log`
- Creating a new branch: `git branch`
- Switching to the current branch: `git switch` and `git checkout`
- Merging two branches: `git merge`
- Resolving conflicts
- Reverting changes: `git revert`
- Uploading changes to GitLab: `git push`
- Downloading changes from GitLab: `git pull`
- Advanced: `git rebase`

## B.2 Basics of Git

### B.2.1 Initializing a Repository: `git init`

To set up Git as a version control system for your project, you need to initialize a new Git repository at the top-level folder, which is the working directory of your project. This is done using the `git init` command.

All files in this folder and its subfolders will automatically become part of the repository. Creating a Git repository is similar to adding an all-powerful passive observer of all things to your project. Git sits there, observes, and takes note of even the smallest changes, such as a single character in a file within a repository with hundreds of files. And it will tell you where these changes occurred if you forget. Once Git is initialized, it monitors all changes made

within the working directory, and it tracks the history of events from that point forward. For this purpose, a historical timeline is created for your project, referred to as a “branch,” and the initial branch is named `main`. So, when someone says they are on the `main branch` or working on the `main branch`, it means they are in the historical main timeline of the project. The Git repository, often abbreviated as `repo`, is a virtual representation of your project, including its history and branches, a book, if you will, where you can look up and retrieve the entire history of the project: you work in your working directory, and the Git repository tracks and stores your work.

### B.2.2 Ignoring Files: `.gitignore`

It’s useful that Git watches and keeps an eye on everything in your project. However, in most projects, there are files and folders that you don’t need or want to keep an eye on. These may include system files, local project settings, libraries with dependencies, and so on.

You can exclude any file or folder from your Git repository by including them in the `.gitignore` file. In the `.gitignore` file, you create a list of file names, folder names, and other items that Git should not track, and Git will ignore these items. Hence the name “gitignore.” Do you want to track a file that you previously ignored? Simply remove the mention of the file in the `gitignore` file, and Git will start tracking it again.

### B.2.3 Adding Changes to the Staging Area: `git add`

The interesting thing about Git as an all-powerful, passive observer of all things is that it’s very passive. As long as you don’t tell Git what to remember, it will passively observe the changes in the project folder but do nothing.

When you make a change to your project that you want Git to include in the project’s history to take a snapshot of so you can refer back to it later, your personal checkpoint, if you will, you need to first stage the changes in the staging area. What is the staging area? The staging area is where you collect changes to files that you want to include in the project’s history.

This is done using the `git add` command. You can specify which files you want to add by naming them, or you can add all of them using `-A`. By doing this, you’re telling Git that you’ve made changes and want it to remember these particular changes so you can recall them later if needed. This is important because you can choose which changes you want to stage, and those are the changes that will eventually be transferred to the history.

**Note:** When you run `git add`, the changes are not transferred to the project’s history. They are only transferred to the staging area.

 Example of git add from the beginning

```
# Create a new directory for your
# repository and navigate to that directory:

mkdir my-repo
cd my-repo

# Initialize the repository with git init:

git init

# Create a .gitignore file for Python code.
# You can use a template from GitHub:

curl https://raw.githubusercontent.com/github/gitignore/master/Python.gitignore -o .gitigno

# Add your files to the repository using git add:

git add .
```

This adds all files in the current directory to the repository, except for the files listed in the `.gitignore` file.

#### B.2.4 Transferring Changes to Memory: `git commit`

The power of Git becomes evident when you start transferring changes to the project history. This is done using the `git commit` command. When you run `git commit`, you inform Git that the changes in the staging area should be added to the history of the project so that they can be referenced or retrieved later.

Additionally, you can add a commit message with the `-m` option to explain what changes were made. So when you look back at the project history, you can see that you added a new feature.

`git commit` creates a snapshot, an image of the current state of your project at that specific time, and adds it to the branch you are currently working on.

As you work on your project and transfer more snapshots, the branch grows and forms a timeline of events. This means you can now look back at every transfer in the branch and see what your code looked like at that time.

You can compare any phase of your code with any other phase of your code to find errors, restore deleted code, or do things that would otherwise not be possible, such as resetting the project to a previous state or creating a new timeline from any point.

So how often should you add these commits? My rule of thumb is not to commit too often. It's better to have a Git repository with too many commits than one with too few commits.

**i** Continuing the example from above:

After adding your files with `git add`, you can create a commit to save your changes. Use the `git commit` command with the `-m` option to specify your commit message:

```
git commit -m "My first commit message"
```

This creates a new commit with the added files and the specified commit message.

### B.2.5 Check the Status of Your Repository: `git status`

If you're wondering what you've changed in your project since the last commit snapshot, you can always check the Git status. Git will list every modified file and the current status of each file.

This status can be either:

- Unchanged (**unmodified**), meaning nothing has changed since you last transferred it, or
- It's been changed (**changed**) but not staged (**staged**) to be transferred into the history, or
- Something has been added to staging (**staged**) and is ready to be transferred into the history.

When you run `git status`, you get an overview of the current state of your project.

**i** Continuing the example from above:

The `git status` command displays the status of your working directory and the staging area. It shows you which files have been modified, which files are staged for commit, and which files are not yet being tracked:

```
git status
```

`git status` is a useful tool to keep track of your changes and ensure that you have added all the desired files for commit.

## B.2.6 Review Your Repository's History: `git log`

 Continuing the example from above:

You can view the history of your commits with the `git log` command. This command displays a list of all the commits in the current branch, along with information such as the author, date, and commit message:

```
git log
```

There are many options to customize the output of `git log`. For example, you can use the `--pretty` option to change the format of the output:

```
git log --pretty=oneline
```

This displays each commit in a single line.

## B.3 Branches (Timelines)

### B.3.1 Creating an Alternative Timeline: `git branch`

In the course of developing a project, you often reach a point where you want to add a new feature, but doing so might require changing the existing code in a way that could be challenging to undo later.

Or maybe you just want to experiment and be able to discard your work if the experiment fails. In such cases, Git allows you to create an alternative timeline called a `branch` to work in.

This new `branch` has its own name and exists in parallel with the `main branch` and all other branches in your project.

During development, you can switch between branches and work on different versions of your code concurrently. This way, you can have a stable codebase in the `main branch` while developing an experimental feature in a separate `branch`. When you switch from one `branch` to another, the code you're working on is automatically reset to the latest commit of the branch you're currently in.

If you're working in a team, different team members can work on their own branches, creating an entire universe of alternative timelines for your project. When features are completed, they can be seamlessly merged back into the `main branch`.

**i** Continuing the example from above:

To create a new `branch`, you can use the `git branch` command with the name of the new `branch` as an argument:

```
git branch my-tests
```

### B.3.2 The Pointer to the Current Branch: `HEAD`

How does Git know where you are on the timeline, and how can you keep track of your position?

You're always working at the tip (`HEAD`) of the currently active branch. The `HEAD` pointer points there quite literally. In a new project archive with just a single `main` branch and only new commits being added, `HEAD` always points to the latest commit in the `main` branch. That's where you are.

However, if you're in a repository with multiple branches, meaning multiple alternative timelines, `HEAD` will point to the latest commit in the branch you're currently working on.

### B.3.3 Switching to an Alternative Timeline: `git switch`

As your project grows, and you have multiple branches, you need to be able to switch between these branches. This is where the `switch` command comes into play.

At any time, you can use the `git switch` command with the name of the branch you want to switch to, and `HEAD` moves from your current branch to the one you specified.

If you've made changes to your code before switching, Git will attempt to carry those changes over to the branch you're switching to. However, if these changes conflict with the target branch, the switch will be canceled.

To resolve this issue without losing your changes, return to the original branch, add and commit your recent changes, and then perform the `switch`.

### B.3.4 Switching to an Alternative Timeline and Making Changes: `git checkout`

To switch between branches, you can also use the `git checkout` command. It works similarly to `git switch` for this purpose: you pass the name of the branch you want to switch to, and `HEAD` moves to the beginning of that branch.

But `checkout` can do more than just switch to another timeline. With `git checkout`, you can also move to any commit point in any timeline. In other words, you can travel back in time and work on code from the past.

To do this, use `git checkout` and provide the commit ID. This is an automatically generated, random combination of letters and numbers that identifies each commit. You can retrieve the commit ID using `git log`. When you run `git log`, you get a list of all the commits in your repository, starting with the most recent ones.

When you use `git checkout` with an older commit ID, you check out a commit in the middle of a branch. This disrupts the timeline, as you're actively attempting to change history. Git doesn't want you to do that because, much like in a science fiction movie, altering the past might also alter the future. In our case, it would break the version control branch's coherence.

To prevent you from accidentally disrupting time and altering history, checking out an earlier commit in any branch results in the warning "Detached Head," which sounds rather ominous. The "Detached Head" warning is appropriate because it accurately describes what's happening. Git literally detaches the head from the branch and sets it aside.

Now, you're working outside of time in a space unbound to any timeline, which again sounds rather threatening but is perfectly fine in reality.

To continue working on this past code, all you need to do is reattach it to the timeline. You can use `git branch` to create a new branch, and the detached head will automatically attach to this new branch.

Instead of breaking the history, you've now created a new alternative timeline that starts in the past, allowing you to work safely. You can continue working on the branch as usual.

**i** Continuing the example from above:

To switch to a new branch, you can use the `git checkout` command:

```
git checkout meine-tests
```

Now you're using the new branch and can make changes independently from the original branch.

### B.3.5 The Difference Between `checkout` and `switch`

What is the difference between `git switch` and `git checkout`? `git switch` and `git checkout` are two different commands that both serve the purpose of switching between branches. You can use both to switch between branches, but they have an important distinction. `git switch` is a new command introduced with Git 2.23. `git checkout` is an older command that has existed since Git 1.6.0. So, `git switch` and `git checkout` have

different origins. `git switch` was introduced to separate the purposes of `git checkout`. `git checkout` has two different purposes: 1. It can be used to switch between branches, and 2. It can be used to reset files to the state of the last commit.

Here's an example: In my project, I made a change since the last commit, but I haven't staged it yet. Then, I realized that I actually don't want this change. I want to reset the file to the state before the last commit. As long as I haven't committed my changes, I can do this with `git checkout` by targeting the specific file. So, if that file is named `main.js`, I can say: `git checkout main.js`. And the file will be reset to the state of the last commit, which makes sense. I'm checking out the file from the last commit.

But that's quite different from switching between the beginning of one branch to another. `git switch` and `git restore` were introduced to separate these two operations. `git switch` is for switching between branches, and `git restore` is for resetting the specified file to the state of the last commit. If you try to restore a file with `git switch`, it simply won't work. It's not intended for that. As I mentioned earlier, it's about separating concerns.

:::{.callout-note} ##### Difference Between `checkout` and `switch` `git checkout` and `git switch` are both commands for switching between branches in a Git repository. The main difference between the two commands is that `git switch` is a newer command specifically designed for branch switching, while `git checkout` is an older command that can be used for various tasks, including branch switching.

Here's an example demonstrating how to initialize a repository and switch between branches:

```
# Create a new directory for your repository
# and navigate to that directory:
mkdir my-repo
cd my-repo

# Initialize the repository with git init:
git init

# Create a new branch with git branch:
git branch my-new-branch

# Switch to the new branch using git switch:
git switch my-new-branch

# Alternatively, you can also use git checkout
# to switch to the new branch:

git checkout my-new-branch
```

Both commands lead to the same result: You are now on the new branch.

## B.4 Merging Branches and Resolving Conflicts

### B.4.1 git merge: Merging Two Timelines

Git allows you to split your development work into as many branches or alternative timelines as you like, enabling you to work on many different versions of your code simultaneously without losing or overwriting any of your work.

This is all well and good, but at some point, you need to bring those various versions of your code back together into one branch. That's where `git merge` comes in.

Consider an example where you have two branches, a `main` branch and an experimental branch called `experimental-branch`. In the experimental branch, there is a new feature. To merge these two branches, you set `HEAD` to the branch where you want to incorporate the code and execute `git merge` followed by the name of the branch you want to merge. `HEAD` is a special pointer that points to the current branch. When you run `git merge`, it combines the code from the branch associated with `HEAD` with the code from the branch specified by the branch name you provide.

```
# Initialize the repository
git init

# Create a new branch called "experimental-branch"
git branch experimental-branch

# Switch to the "experimental-branch"
git checkout experimental-branch

# Add the new feature here and
# make a commit
# ...

# Switch back to the "main" branch
git checkout main

# Perform the merge
git merge experimental-branch
```

During the merge, matching pieces of code in the branches overlap, and any new code from the branch being merged is added to the project. So now, the main branch also contains the code from the experimental branch, and the events of the two separate timelines have been merged into a single one. What's interesting is that even though the experimental branch was merged

with the main branch, the last commit of the experimental branch remains intact, allowing you to continue working on the experimental branch separately if you wish.

### B.4.2 Resolving Conflicts When Merging

Merging branches where there are no code changes at the same place in both branches is a straightforward process. It's also a rare process. In most cases, there will be some form of conflict between the branches – the same code or the same code area has been modified differently in the different branches. Merging two branches with such conflicts will not work, at least not automatically.

In this case, Git doesn't know how to merge this code. So, when such a situation occurs, it's marked as a conflict, and the merging process is halted. This might sound more dramatic than it is. When you get a conflict warning, Git is saying there are two different versions here, and Git needs to know which one you want to keep. To help you figure out the conflict, Git combines all the code into a single file and automatically marks the conflicting code as the current change, which is the original code from the branch you're working on, or as the incoming change, which is the code from the file you're trying to merge.

To resolve this conflict, you'll edit the file to literally resolve the code conflict. This might mean accepting either the current or incoming change and discarding the other. It could mean combining both changes or something else entirely. It's up to you. So, you edit the code to resolve the conflict. Once you've resolved the conflict by editing the code, you add the new conflict-free version to the staging area with `git add` and then commit the merged code with `git commit`. That's how the conflict is resolved.

A merge conflict occurs when Git struggles to automatically merge changes from two different branches. This usually happens when changes were made to the same line in the same file in both branches. To resolve a merge conflict, you must manually edit the affected files and choose the desired changes. Git marks the conflict areas in the file with special markings like `<<<<<`, `=====`, and `>>>>>`. You can search for these markings and manually select the desired changes. After resolving the conflicts, you can add the changes with `git add` and create a new commit with `git commit` to complete the merge.

Here's an example:

```
# Perform the merge (this will cause a conflict)
git merge experimenteller-branch

# Open the affected file in an editor and manually resolve the conflicts
# ...

# Add the modified file
git add <filename>
```

```
# Create a new commit  
git commit -m "Resolved conflicts"
```

### B.4.3 git revert: Undoing Something

One of the most powerful features of any software tool is the “Undo” button. Make a mistake, press “Undo,” and it’s as if it never happened. However, that’s not quite as simple when an all-powerful, passive observer is watching and recording your project’s history. How do you undo something that you’ve added to the history without rewriting the history?

The answer is that you can overwrite the history with the `git reset` command, but that’s quite risky and not a good practice.

A better solution is to work with the historical timeline and simply place an older version of your code at the top of the branch. This is done with `git revert`. To make this work, you need to know the commit ID of the commit you want to go back to.

The commit ID is a machine-generated set of random numbers and letters, also known as a hash. To get a list of all the commits in the repository, including the commit ID and commit message, you can run `git log`.

```
# Show the list of all operations in the repository  
git log
```

By the way, it’s a good idea to leave clear and informative commit messages for this reason. This way, you know what happened in your previous commits. Once you’ve found the commit you want to revert to, call that commit ID with `git revert`, and then the ID. This will create a new commit at the top of the branch with the code from the reference commit. To transfer the code to the branch, add a commit message and save it. Now, the last commit in your branch matches the commit you’re reverting to, and your project’s history remains intact.

**i** An example with `git revert`

```
# Initialize a new repository
git init

# Create a new file
echo "Hello, World" > file.txt

# Add the file to the repository
git add file.txt

# Create a new commit
git commit -m "First commit"

# Modify the file
echo "Goodbye, World" > file.txt

# Add the modified file
git add file.txt

# Create a new commit
git commit -m "Second commit"

# Use git log to find the commit ID of the second commit
git log

# Use git revert to undo the changes from the second commit
git revert <commit-id>
```

To download the `students` branch from the repository `git@git-ce.rwth-aachen.de:spotseven-lab/numerisoc` to your local machine, add a file, and upload the changes, you can follow these steps:

**i** An example with `git clone`, `git checkout`, `git add`, `git commit`, `git push`

```
# Clone the repository to your local machine:  
git clone git@git-ce.rwth-aachen.de:spotseven-lab/numerische-mathematik-sommersemester2023  
  
# Change to the cloned repository:  
cd numerische-mathematik-sommersemester2023  
  
# Switch to the students branch:  
git checkout students  
  
# Create the Test folder if it doesn't exist:  
mkdir Test  
  
# Create the Testdatei.txt file in the Test folder:  
touch Test/Testdatei.txt  
  
# Add the file with git add:  
git add Test/Testdatei.txt  
  
# Commit the changes with git commit:  
git commit -m "Added Testdatei.txt"  
  
# Push the changes with git push:  
git push origin students
```

This will upload the changes to the server and update the students branch in the repository.

## B.5 Downloading from GitLab

To download changes from a GitLab repository to your local machine, you can use the `git pull` command. This command downloads the latest changes from the specified remote repository and merges them with your local repository.

Here is an example:

### An example with `git pull`

```
# Navigate to the local repository  
# linked to the GitHub repository:  
cd my-local-repository  
  
# Make sure you are in the correct branch:  
git checkout main  
  
# Download the latest changes from GitHub:  
git pull origin main
```

This downloads the latest changes from the main branch of the remote repository named “origin” and merges them with your local repository.

If there are conflicts between the downloaded changes and your local changes, you will need to resolve them manually before proceeding.

## B.6 Advanced

### B.6.1 `git rebase`: Moving the Base of a Branch

In some cases, you may need to “rewrite history.” A common scenario is that you’ve been working on a new feature in a feature branch, and you realize that the work should have actually happened in the `main branch`.

To resolve this issue and make it appear as if the work occurred in the `main branch`, you can reset the experimental branch. “Rebase” literally means detaching the base of the experimental branch and moving it to the beginning of another branch, giving the branch a new base, thus “rebasing.”

This operation is performed from the branch you want to “rebase.” You use `git rebase` and specify the branch you want to use as the new base. If there are no conflicts between the experimental branch and the branch you want to rebase onto, this process happens automatically.

If there are conflicts, Git will guide you through the conflict resolution process for each commit from the rebase branch.

This may sound like a lot, but there’s a good reason for it. You are literally rewriting history by transferring commits from one branch to another. To maintain the coherence of the new version history, there should be no conflicts within the commits. So, you need to resolve

them one by one until the history is clean. It goes without saying that this can be a fairly labor-intensive process. Therefore, you should not use `git rebase` frequently.

### An example with `git rebase`

`git rebase` is a command used to change the base of a branch. This means that commits from the branch are applied to a new base, which is usually another branch. It can be used to clean up the repository history and avoid merge conflicts.

Here is an example showing how to use `git rebase`:

- In this example, we initialize a new Git repository and create a new file. We add the file to the repository and make an initial commit. Then, we create a new branch called “feature” and switch to that branch. We make changes to the file in the feature branch and create a new commit.
- Then, we switch back to the main branch and make changes to the file again. We add the modified file and make another commit.
- To rebase the feature branch onto the main branch, we first switch to the feature branch and then use the `git rebase` command with the name of the main branch as an argument. This applies the commits from the feature branch to the main branch and changes the base of the feature branch.

```

# Initialize a new repository
git init
# Create a new file
echo "Hello World" > file.txt
# Add the file to the repository
git add file.txt
# Create an initial commit
git commit -m "Initial commit"
# Create a new branch called "feature"
git branch feature
# Switch to the "feature" branch
git checkout feature
# Make changes to the file in the "feature" branch
echo "Hello Feature World" > file.txt
# Add the modified file
git add file.txt
# Create a new commit in the "feature" branch
git commit -m "Feature commit"
# Switch back to the "main" branch
git checkout main
# Make changes to the file in the "main" branch
echo "Hello Main World" > file.txt
# Add the modified file
git add file.txt
# Create a new commit in the "main" branch
git commit -m "Main commit"
# Use git rebase to rebase the "feature" branch
# onto the "main" branch
git checkout feature
git rebase main

```

## B.7 Exercises

In order to be able to carry out this exercise, we provide you with a functional working environment. This can be accessed [here](#). You can log in using your GMID. If you do not have one, you can generate one [here](#). Once you have successfully logged in to the server, you must open a terminal instance. You are now in a position to carry out the exercise.

Alternatively, you can also carry out the exercise locally on your computer, but then you will need to install git.

### **B.7.1 Create project folder**

First create the `test-repo` folder via the command line and then navigate to this folder using the corresponding command.

## **B.8 Initialize repo**

Now initialize the repository so that the future project, which will be saved in the `test-repo` folder, and all associated files are versioned.

### **B.8.1 Do not upload / ignore certain file types**

In order to carry out this exercise, you must first download a file which you then have git ignore. To do this, download the current examination regulations for the Bachelor's degree program in Electrical Engineering using the following command `curl -o pruefungsordnung.pdf https://www.th-koeln.de/mam/downloads/deutsch/studium/studiengaenge/f07/ordnungen_plaene/f07...`

The PDF file has been stored in the root directory of your repo and you must now exclude it from being uploaded so that no changes to this file are tracked. Please note that not only this one PDF file should be ignored, but all PDF files in the repo.

### **B.8.2 Create file and stage it**

In order to be able to commit a change later and thus make it traceable, it must first be staged. However, as we only have a PDF file so far, which is to be ignored by git, we cannot stage anything. Therefore, in this task, a file `test.txt` with some string as content is to be created and then staged.

### **B.8.3 Create another file and check status**

To understand the status function, you should create the file `test2.txt` and then call the status function of git.

### **B.8.4 Commit changes**

After the changes to the `test.txt` file have been staged and these are now to be transferred to the project process, they must be committed. Therefore, in this step you should perform a corresponding commit in the current branch with the message `test-commit`. Finally, you should also display the history of the commits.

### **B.8.5 Create a new branch and switch to it**

In this task, you are to create a new branch with the name `change-text` in which you will later make changes. You should then switch to this branch.

### **B.8.6 Commit changes in the new branch**

To be able to merge the new branch into the main branch later, you must first make changes to the `test.txt` file. To do this, open the file and simply change the character string in this file before saving the changes and closing the file. Before you now commit the file, you should reset the file to the status of the last commit for practice purposes and thus undo the change. After you have done this, open the file `test.txt` again and change the character string again before saving and closing the file. This time you should commit the file `test.txt` and then commit it with the message `test-commit2`.

### **B.8.7 Merge branch into main**

After you have committed the change to the `test.txt` file, you should merge the `change-text` branch including the change into the main branch so that it is also available there.

### **B.8.8 Resolve merge conflict**

To simulate a merge conflict, you must first change the content of the `test.txt` file before you commit the change. Then switch to the branch `change-text` and change the file `test.txt` there as well before you commit the change. Now you should try to merge the branch `change-text` into the main branch and solve the problems that occur in order to be able to perform the merge successfully.

# C Python Introduction

## C.1 Recommendations

[Beginner's Guide to Python](#)

# D Documentation of the Sequential Parameter Optimization

This document describes the `Spot` features. The official `spotPython` documentation can be found here: <https://sequential-parameter-optimization.github.io/spotPython/>.

## D.1 Example: `spot`

```
import numpy as np
from math import inf
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
from scipy.optimize import shgo
from scipy.optimize import direct
from scipy.optimize import differential_evolution
import matplotlib.pyplot as plt
```

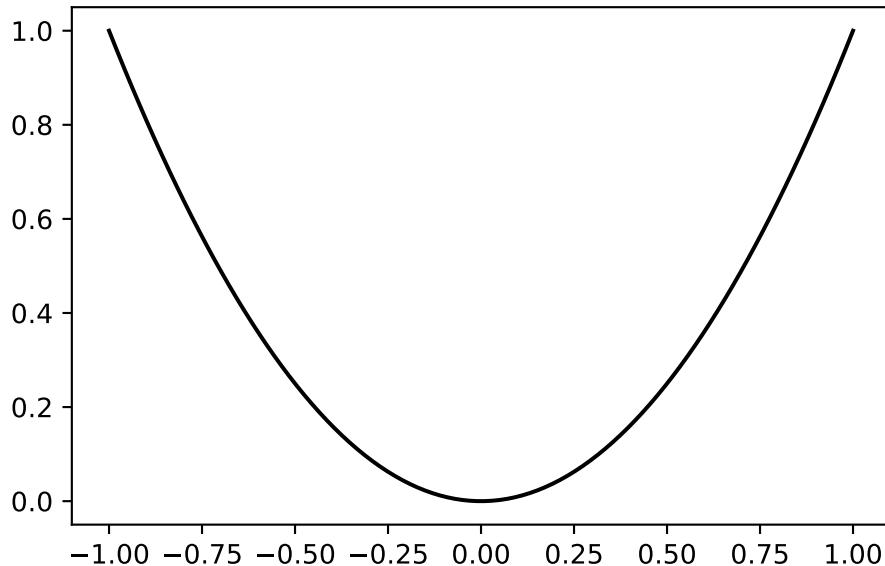
### D.1.1 The Objective Function

The `spotPython` package provides several classes of objective functions. We will use an analytical objective function, i.e., a function that can be described by a (closed) formula:

$$f(x) = x^2$$

```
fun = analytical().fun_sphere
```

```
x = np.linspace(-1,1,100).reshape(-1,1)
y = fun(x)
plt.figure()
plt.plot(x,y, "k")
plt.show()
```



```
from spotPython.utils.init import fun_control_init, design_control_init, surrogate_control_init
spot_1 = spot.Spot(fun=fun,
                    fun_control=fun_control_init(
                        lower = np.array([-10]),
                        upper = np.array([100]),
                        fun_evals = 7,
                        fun_repeats = 1,
                        max_time = inf,
                        noise = False,
                        tolerance_x = np.sqrt(np.spacing(1)),
                        var_type=["num"],
                        infill_criterion = "y",
                        n_points = 1,
                        seed=123,
                        log_level = 50),
                    design_control=design_control_init(
                        init_size=5,
                        repeats=1),
                    surrogate_control=surrogate_control_init(
                        noise=False,
                        min_theta=-4,
                        max_theta=3,
                        n_theta=1,
                        model_optimizer=differential_evolution,
                        model_fun_evals=10000))
```

`spot`'s `__init__` method sets the control parameters. There are two parameter groups:

1. external parameters can be specified by the user
2. internal parameters, which are handled by `spot`.

### D.1.2 External Parameters

external parameter	type	description	default	mandatory
<code>fun</code>	object	objective function		yes
<code>lower</code>	array	lower bound		yes
<code>upper</code>	array	upper bound		yes
<code>fun_evals</code>	int	number of function evaluations	15	no
<code>fun_evals</code>	int	number of function evaluations	15	no
<code>fun_control</code>	dict	noise etc.	{}	n
<code>max_time</code>	int	max run time budget	<code>inf</code>	no
<code>noise</code>	bool	if repeated evaluations of <code>fun</code> results in different values, then <code>noise</code> should be set to True.	<code>False</code>	no

external parameter	type	description	default	mandatory
<b>tolerance_x</b>	float	tolerance for new x solutions. Minimum distance of new solutions, generated by <code>suggest_new_X</code> , to already existing solutions. If zero (which is the default), every new solution is accepted.	0	no
<b>var_type</b>	list	list of type information, can be either "num" or "factor"	["num"]	no
<b>infill_criterion</b>	string	Can be "y", "s", "y" "ei" (negative expected improvement), or "all"		no
<b>n_points</b>	int	number of infill points	1	no
<b>seed</b>	int	initial seed. If <code>Spot.run()</code> is called twice, different results will be generated. To reproduce results, the <code>seed</code> can be used.	123	no

external parameter	type	description	default	mandatory
log_level	int	log level with the following settings: NOTSET (0), DEBUG (10: Detailed information, typically of interest only when diagnosing problems.), INFO (20: Confirmation that things are working as expected.), WARNING (30: An indication that something unexpected happened, or indicative of some problem in the near future (e.g. ‘disk space low’). The software is still working as expected.), ERROR (40: Due to a more serious problem, the software has not been able to perform some function.), and CRITICAL (50: A serious error, indicating that the program itself may be unable to continue running.)	50	no

external parameter	type	description	default	mandatory
<code>show_models</code>	bool	Plot model. Currently only 1-dim functions are supported	<code>False</code>	no
<code>design</code>	object	experimental design	<code>None</code>	no
<code>design_control</code>	dict	control parameters	see below	no
<code>surrogate</code>		surrogate model	<code>kriging</code>	no
<code>surrogate_control</code>		control parameters	see below	no
<code>optimizer</code>	object	optimizer	see below	no
<code>optimizer_control</code>		control parameters	see below	no

- Besides these single parameters, the following parameter dictionaries can be specified by the user:

- `fun_control`
- `design_control`
- `surrogate_control`
- `optimizer_control`

## D.2 The `fun_control` Dictionary

external parameter	type	description	default	mandatory
<code>sigma</code>	float	noise: standard deviation	0	yes
<code>seed</code>	int	seed for rng	124	yes

## D.3 The `design_control` Dictionary

external parameter	type	description	default	mandatory
<code>init_size</code>	int	initial sample size	10	yes

external parameter	type	description	default	mandatory
<b>repeats</b>	int	number of repeats of the initial samples	1	yes

## D.4 The surrogate\_control Dictionary

external parameter	type	description	default	mandatory
<b>noise</b>				
<b>model_optimizer</b>	object	optimizer	<b>differential_evolution</b>	
<b>model_fun_evals</b>				
<b>min_theta</b>			-3.	
<b>max_theta</b>			3.	
<b>n_theta</b>			1	
<b>n_p</b>			1	
<b>optim_p</b>			<b>False</b>	
<b>cod_type</b>			"norm"	
<b>var_type</b>				
<b>use_cod_y</b>	bool		<b>False</b>	

## D.5 The optimizer\_control Dictionary

external parameter	type	description	default	mandatory
<b>max_iter</b>	int	max number of iterations. Note: these are the cheap evaluations on the surrogate.	1000	no

## D.6 Run

```
spot_1.run()
```

```
spotPython tuning: 2.0106521524877827 [#####--] 85.71%
spotPython tuning: 0.01033163973935242 [#####---] 100.00% Done...
```

```
<spotPython.spot.spot.Spot at 0x1066ef410>
```

## D.7 Print the Results

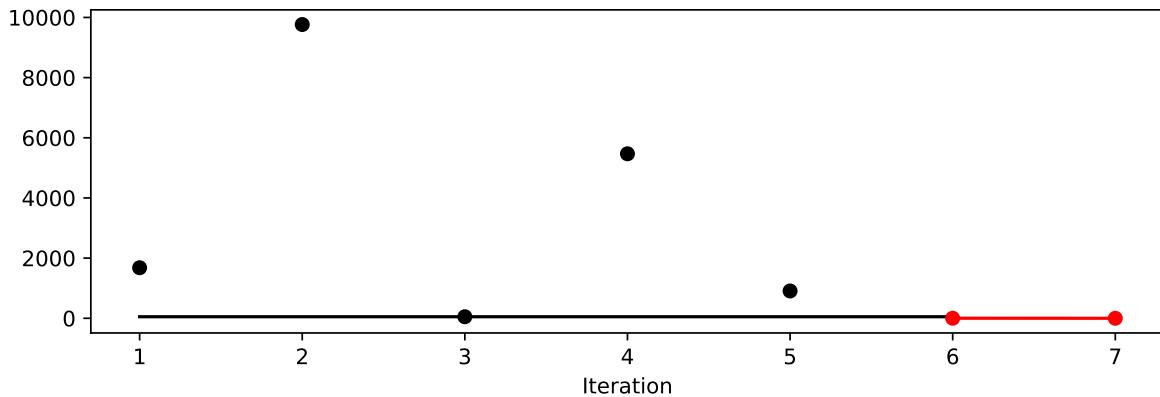
```
spot_1.print_results()
```

```
min y: 0.01033163973935242
x0: -0.10164467393499976
```

```
[['x0', -0.10164467393499976]]
```

## D.8 Show the Progress

```
spot_1.plot_progress()
```

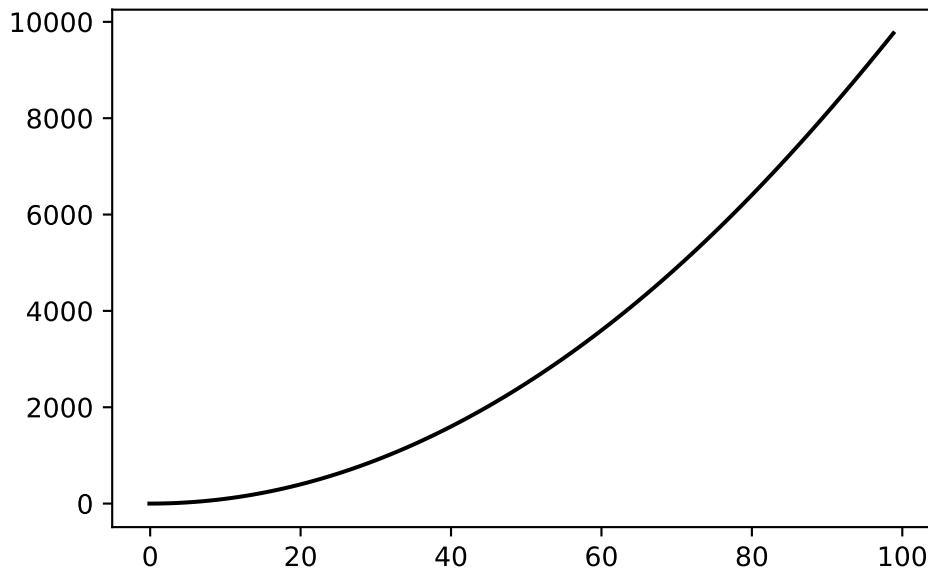


## D.9 Visualize the Surrogate

- The plot method of the `kriging` surrogate is used.
- Note: the plot uses the interval defined by the ranges of the natural variables.

```
spot_1.surrogate.plot()
```

<Figure size 2700x1800 with 0 Axes>



## D.10 Run With a Specific Start Design

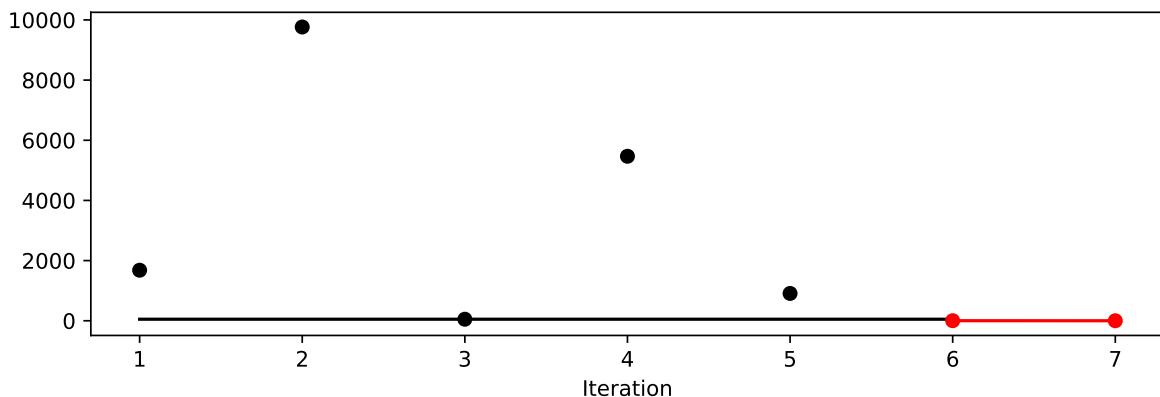
```
spot_x0 = spot.Spot(fun=fun,
                     fun_control=fun_control_init(
                         lower = np.array([-10]),
                         upper = np.array([100]),
                         fun_evals = 7,
                         fun_repeats = 1,
                         max_time = inf,
                         noise = False,
                         tolerance_x = np.sqrt(np.spacing(1)),
                         var_type=["num"],
```

```

        infill_criterion = "y",
        n_points = 1,
        seed=123,
        log_level = 50),
    design_control=design_control_init(
        init_size=5,
        repeats=1),
    surrogate_control=surrogate_control_init(
        noise=False,
        min_theta=-4,
        max_theta=3,
        n_theta=1,
        model_optimizer=differential_evolution,
        model_fun_evals=10000))
spot_x0.run(X_start=np.array([0.5, -0.5]))
spot_x0.plot_progress()

```

spotPython tuning: 2.0106521524877827 [#####-] 85.71%  
spotPython tuning: 0.01033163973935242 [#####] 100.00% Done...



## D.11 Init: Build Initial Design

```

from spotPython.design.spacefilling import spacefilling
from spotPython.build.kriging import Kriging
from spotPython.fun.objectivefunctions import analytical
gen = spacefilling(2)

```

```

rng = np.random.RandomState(1)
lower = np.array([-5,-0])
upper = np.array([10,15])
fun = analytical().fun_branin
fun_control = {"sigma": 0,
               "seed": 123}

X = gen.scipy_lhd(10, lower=lower, upper = upper)
print(X)
y = fun(X, fun_control=fun_control)
print(y)

```

```

[[ 8.97647221 13.41926847]
 [ 0.66946019  1.22344228]
 [ 5.23614115 13.78185824]
 [ 5.6149825  11.5851384 ]
 [-1.72963184  1.66516096]
 [-4.26945568  7.1325531 ]
 [ 1.26363761 10.17935555]
 [ 2.88779942  8.05508969]
 [-3.39111089  4.15213772]
 [ 7.30131231  5.22275244]]
[128.95676449 31.73474356 172.89678121 126.71295908 64.34349975
 70.16178611 48.71407916 31.77322887 76.91788181 30.69410529]

```

## D.12 Replicability

Seed

```

gen = spacefilling(2, seed=123)
X0 = gen.scipy_lhd(3)
gen = spacefilling(2, seed=345)
X1 = gen.scipy_lhd(3)
X2 = gen.scipy_lhd(3)
gen = spacefilling(2, seed=123)
X3 = gen.scipy_lhd(3)
X0, X1, X2, X3

```

```

(array([[0.77254938, 0.31539299],
       [0.59321338, 0.93854273],

```

```

[0.27469803, 0.3959685 ]]),
array([[0.78373509, 0.86811887],
       [0.06692621, 0.6058029 ],
       [0.41374778, 0.00525456]]),
array([[0.121357 , 0.69043832],
       [0.41906219, 0.32838498],
       [0.86742658, 0.52910374]]),
array([[0.77254938, 0.31539299],
       [0.59321338, 0.93854273],
       [0.27469803, 0.3959685 ]]))

```

## D.13 Surrogates

### D.13.1 A Simple Predictor

The code below shows how to use a simple model for prediction. Assume that only two (very costly) measurements are available:

1.  $f(0) = 0.5$
2.  $f(2) = 2.5$

We are interested in the value at  $x_0 = 1$ , i.e.,  $f(x_0 = 1)$ , but cannot run an additional, third experiment.

```

from sklearn import linear_model
X = np.array([[0], [2]])
y = np.array([0.5, 2.5])
S_lm = linear_model.LinearRegression()
S_lm = S_lm.fit(X, y)
X0 = np.array([[1]])
y0 = S_lm.predict(X0)
print(y0)

```

[1.5]

Central Idea: Evaluation of the surrogate model  $S_{lm}$  is much cheaper (or / and much faster) than running the real-world experiment  $f$ .

## D.14 Demo/Test: Objective Function Fails

SPOT expects `np.nan` values from failed objective function values. These are handled. Note: SPOT's counter considers only successful executions of the objective function.

```
import numpy as np
from spotPython.fun.objectivefunctions import analytical
from spotPython.spot import spot
import numpy as np
from math import inf
# number of initial points:
ni = 20
# number of points
n = 30

fun = analytical().fun_random_error
fun_control=fun_control_init(
    lower = np.array([-1]),
    upper= np.array([1]),
    fun_evals = n,
    show_progress=False)
design_control=design_control_init(init_size=ni)

spot_1 = spot.Spot(fun=fun,
                    fun_control=fun_control,
                    design_control=design_control)
spot_1.run()
# To check whether the run was successfully completed,
# we compare the number of evaluated points to the specified
# number of points.
assert spot_1.y.shape[0] == n
```

```
[      nan      nan -0.02203599 -0.21843718  0.78240941      nan
-0.3923345  0.67234256  0.31802454 -0.68898927 -0.75129705  0.97550354
 0.41757584      nan  0.82585329      nan -0.49274073      nan
-0.17991251  0.1481835 ]
[-1.]
[nan]
[-0.14624037]
[0.166475]
[nan]
[-0.3352401]
```

$[-0.47259301]$   
 $[0.95541987]$   
 $[0.17335968]$   
 $[-0.58552368]$   
 $[-0.20126111]$   
 $[-0.60100809]$   
 $[-0.97897336]$   
 $[-0.2748985]$   
 $[0.8359486]$   
 $[0.99035591]$   
 $[0.01641232]$   
 $[0.5629346]$

## D.15 PyTorch: Detailed Description of the Data Splitting

### D.15.1 Description of the "train\_hold\_out" Setting

The "train\_hold\_out" setting is used by default. It uses the loss function specified in `fun_control` and the metric specified in `fun_control`.

1. First, the method `HyperTorch().fun_torch` is called.
2. `fun_torch()`, which is implemented in the file `hypertorch.py`, calls `evaluate_hold_out()` as follows:

```
df_eval, _ = evaluate_hold_out(
    model,
    train_dataset=fun_control["train"],
    shuffle=self.fun_control["shuffle"],
    loss_function=self.fun_control["loss_function"],
    metric=self.fun_control["metric_torch"],
    device=self.fun_control["device"],
    show_batch_interval=self.fun_control["show_batch_interval"],
    path=self.fun_control["path"],
    task=self.fun_control["task"],
    writer=self.fun_control["writer"],
    writerId=config_id,
)
```

Note: Only the data set `fun_control["train"]` is used for training and validation. It is used in `evaluate_hold_out` as follows:

```
trainloader, valloader = create_train_val_data_loaders(
    dataset=train_dataset, batch_size=batch_size_instance, shuffle=shuffle
)
```

`create_train_val_data_loaders()` splits the `train_dataset` into `trainloader` and `valloader` using `torch.utils.data.random_split()` as follows:

```
def create_train_val_data_loaders(dataset, batch_size, shuffle, num_workers=0):
    test_abs = int(len(dataset) * 0.6)
    train_subset, val_subset = random_split(dataset, [test_abs, len(dataset) - test_abs])
    trainloader = torch.utils.data.DataLoader(
        train_subset, batch_size=int(batch_size), shuffle=shuffle, num_workers=num_workers
    )
    valloader = torch.utils.data.DataLoader(
```

```

        val_subset, batch_size=int(batch_size), shuffle=shuffle, num_workers=num_workers
    )
    return trainloader, valloader

```

The optimizer is set up as follows:

```

optimizer_instance = net.optimizer
lr_mult_instance = net.lr_mult
sgd_momentum_instance = net.sgd_momentum
optimizer = optimizer_handler(
    optimizer_name=optimizer_instance,
    params=net.parameters(),
    lr_mult=lr_mult_instance,
    sgd_momentum=sgd_momentum_instance,
)

```

3. `evaluate_hold_out()` sets the `net` attributes such as `epochs`, `batch_size`, `optimizer`, and `patience`. For each epoch, the methods `train_one_epoch()` and `validate_one_epoch()` are called, the former for training and the latter for validation and early stopping. The validation loss from the last epoch (not the best validation loss) is returned from `evaluate_hold_out`.
4. The method `train_one_epoch()` is implemented as follows:

```

def train_one_epoch(
    net,
    trainloader,
    batch_size,
    loss_function,
    optimizer,
    device,
    show_batch_interval=10_000,
    task=None,
):
    running_loss = 0.0
    epoch_steps = 0
    for batch_nr, data in enumerate(trainloader, 0):
        input, target = data
        input, target = input.to(device), target.to(device)
        optimizer.zero_grad()
        output = net(input)
        if task == "regression":
            target = target.unsqueeze(1)

```

```

        if target.shape == output.shape:
            loss = loss_function(output, target)
        else:
            raise ValueError(f"Shapes of target and output do not match:
                            {target.shape} vs {output.shape}")
    elif task == "classification":
        loss = loss_function(output, target)
    else:
        raise ValueError(f"Unknown task: {task}")
    loss.backward()
    torch.nn.utils.clip_grad_norm_(net.parameters(), max_norm=1.0)
    optimizer.step()
    running_loss += loss.item()
    epoch_steps += 1
    if batch_nr % show_batch_interval == (show_batch_interval - 1):
        print(
            "Batch: %5d. Batch Size: %d. Training Loss (running): %.3f"
            % (batch_nr + 1, int(batch_size), running_loss / epoch_steps)
        )
        running_loss = 0.0
    return loss.item()

```

5. The method `validate_one_epoch()` is implemented as follows:

```

def validate_one_epoch(net, valloader, loss_function, metric, device, task):
    val_loss = 0.0
    val_steps = 0
    total = 0
    correct = 0
    metric.reset()
    for i, data in enumerate(valloader, 0):
        # get batches
        with torch.no_grad():
            input, target = data
            input, target = input.to(device), target.to(device)
            output = net(input)
            # print(f"target: {target}")
            # print(f"output: {output}")
            if task == "regression":
                target = target.unsqueeze(1)
                if target.shape == output.shape:
                    loss = loss_function(output, target)

```

```

        else:
            raise ValueError(f"Shapes of target and output
                            do not match: {target.shape} vs {output.shape}")
        metric_value = metric.update(output, target)
    elif task == "classification":
        loss = loss_function(output, target)
        metric_value = metric.update(output, target)
        _, predicted = torch.max(output.data, 1)
        total += target.size(0)
        correct += (predicted == target).sum().item()
    else:
        raise ValueError(f"Unknown task: {task}")
    val_loss += loss.cpu().numpy()
    val_steps += 1
loss = val_loss / val_steps
print(f"Loss on hold-out set: {loss}")
if task == "classification":
    accuracy = correct / total
    print(f"Accuracy on hold-out set: {accuracy}")
# metric on all batches using custom accumulation
metric_value = metric.compute()
metric_name = type(metric).__name__
print(f"{metric_name} value on hold-out data: {metric_value}")
return metric_value, loss

```

#### D.15.1.1 Description of the "test\_hold\_out" Setting

It uses the loss function specified in `fun_control` and the metric specified in `fun_control`.

1. First, the method `HyperTorch().fun_torch` is called.
2. `fun_torch()` calls `spotPython.torch.traintest.evaluate_hold_out()` similar to the `"train_hold_out"` setting with one exception: It passes an additional `test` data set to `evaluate_hold_out()` as follows:

```
test_dataset=fun_control["test"]
```

`evaluate_hold_out()` calls `create_train_test_data_loaders` instead of `create_train_val_data_loaders`. The two data sets are used in `create_train_test_data_loaders` as follows:

```

def create_train_test_data_loaders(dataset, batch_size, shuffle, test_dataset,
    num_workers=0):
    trainloader = torch.utils.data.DataLoader(
        dataset, batch_size=int(batch_size), shuffle=shuffle,
        num_workers=num_workers
    )
    testloader = torch.utils.data.DataLoader(
        test_dataset, batch_size=int(batch_size), shuffle=shuffle,
        num_workers=num_workers
    )
    return trainloader, testloader

```

3. The following steps are identical to the "train\_hold\_out" setting. Only a different data loader is used for testing.

#### D.15.1.2 Detailed Description of the "train\_cv" Setting

It uses the loss function specified in `fun_control` and the metric specified in `fun_control`.

1. First, the method `HyperTorch().fun_torch` is called.
2. `fun_torch()` calls `spotPython.torch.traintest.evaluate_cv()` as follows (Note: Only the data set `fun_control["train"]` is used for CV.):

```

df_eval, _ = evaluate_cv(
    model,
    dataset=fun_control["train"],
    shuffle=self.fun_control["shuffle"],
    device=self.fun_control["device"],
    show_batch_interval=self.fun_control["show_batch_interval"],
    task=self.fun_control["task"],
    writer=self.fun_control["writer"],
    writerId=config_id,
)

```

3. In 'evaluate\_cv()', the following steps are performed: The optimizer is set up as follows:

```

optimizer_instance = net.optimizer
lr_instance = net.lr
sgd_momentum_instance = net.sgd_momentum
optimizer = optimizer_handler(optimizer_name=optimizer_instance,
    params=net.parameters(), lr_mult=lr_mult_instance)

```

`evaluate_cv()` sets the `net` attributes such as `epochs`, `batch_size`, `optimizer`, and `patience`. CV is implemented as follows:

```
def evaluate_cv(
    net,
    dataset,
    shuffle=False,
    loss_function=None,
    num_workers=0,
    device=None,
    show_batch_interval=10_000,
    metric=None,
    path=None,
    task=None,
    writer=None,
    writerId=None,
):
    lr_mult_instance = net.lr_mult
    epochs_instance = net.epochs
    batch_size_instance = net.batch_size
    k_folds_instance = net.k_folds
    optimizer_instance = net.optimizer
    patience_instance = net.patience
    sgd_momentum_instance = net.sgd_momentum
    removed_attributes, net = get_removed_attributes_and_base_net(net)
    metric_values = {}
    loss_values = {}

    try:
        device = getDevice(device=device)
        if torch.cuda.is_available():
            device = "cuda:0"
            if torch.cuda.device_count() > 1:
                print("We will use", torch.cuda.device_count(), "GPUs!")
                net = nn.DataParallel(net)
        net.to(device)
        optimizer = optimizer_handler(
            optimizer_name=optimizer_instance,
            params=net.parameters(),
            lr_mult=lr_mult_instance,
            sgd_momentum=sgd_momentum_instance,
        )
        kfold = KFold(n_splits=k_folds_instance, shuffle=shuffle)
        for fold, (train_ids, val_ids) in enumerate(kfold.split(dataset)):
```

```

print(f"Fold: {fold + 1}")
train_subampler = torch.utils.data.SubsetRandomSampler(train_ids)
val_subampler = torch.utils.data.SubsetRandomSampler(val_ids)
trainloader = torch.utils.data.DataLoader(
    dataset, batch_size=batch_size_instance,
    sampler=train_subampler, num_workers=num_workers
)
valloader = torch.utils.data.DataLoader(
    dataset, batch_size=batch_size_instance,
    sampler=val_subampler, num_workers=num_workers
)
# each fold starts with new weights:
reset_weights(net)
# Early stopping parameters
best_val_loss = float("inf")
counter = 0
for epoch in range(epochs_instance):
    print(f"Epoch: {epoch + 1}")
    # training loss from one epoch:
    training_loss = train_one_epoch(
        net=net,
        trainloader=trainloader,
        batch_size=batch_size_instance,
        loss_function=loss_function,
        optimizer=optimizer,
        device=device,
        show_batch_interval=show_batch_interval,
        task=task,
    )
    # Early stopping check. Calculate validation loss from one epoch:
    metric_values[fold], loss_values[fold] = validate_one_epoch(
        net, valloader=valloader, loss_function=loss_function,
        metric=metric, device=device, task=task
    )
    # Log the running loss averaged per batch
    metric_name = "Metric"
    if metric is None:
        metric_name = type(metric).__name__
        print(f"{metric_name} value on hold-out data:
            {metric_values[fold]}")
    if writer is not None:
        writer.add_scalars(

```

```

        "evaluate_cv fold:" + str(fold + 1) +
        ". Train & Val Loss and Val Metric" + writerId,
        {"Train loss": training_loss, "Val loss":
            loss_values[fold], metric_name: metric_values[fold]},
        epoch + 1,
    )
    writer.flush()
if loss_values[fold] < best_val_loss:
    best_val_loss = loss_values[fold]
    counter = 0
    # save model:
    if path is not None:
        torch.save(net.state_dict(), path)
else:
    counter += 1
    if counter >= patience_instance:
        print(f"Early stopping at epoch {epoch}")
        break
df_eval = sum(loss_values.values()) / len(loss_values.values())
df_metrics = sum(metric_values.values()) / len(metric_values.values())
df_preds = np.nan
except Exception as err:
    print(f"Error in Net_Core. Call to evaluate_cv() failed. {err=},
          {type(err)=}")
    df_eval = np.nan
    df_preds = np.nan
add_attributes(net, removed_attributes)
if writer is not None:
    metric_name = "Metric"
    if metric is None:
        metric_name = type(metric).__name__
    writer.add_scalars(
        "CV: Val Loss and Val Metric" + writerId,
        {"CV-loss": df_eval, metric_name: df_metrics},
        epoch + 1,
    )
    writer.flush()
return df_eval, df_preds, df_metrics

```

4. The method `train_fold()` is implemented as shown above.
5. The method `validate_one_epoch()` is implemented as shown above. In contrast to the hold-out setting, it is called for each of the  $k$  folds. The results are stored in a

dictionaries `metric_values` and `loss_values`. The results are averaged over the  $k$  folds and returned as `df_eval`.

#### D.15.1.3 Detailed Description of the "test\_cv" Setting

It uses the loss function specified in `fun_control` and the metric specified in `fun_control`.

1. First, the method `HyperTorch().fun_torch` is called.
2. `fun_torch()` calls `spotPython.torch.traintest.evaluate_cv()` as follows:

```
df_eval, _ = evaluate_cv(  
    model,  
    dataset=fun_control["test"],  
    shuffle=self.fun_control["shuffle"],  
    device=self.fun_control["device"],  
    show_batch_interval=self.fun_control["show_batch_interval"],  
    task=self.fun_control["task"],  
    writer=self.fun_control["writer"],  
    writerId=config_id,  
)
```

Note: The data set `fun_control["test"]` is used for CV. The rest is the same as for the "train\_cv" setting.

#### D.15.1.4 Detailed Description of the Final Model Training and Evaluation

There are two methods that can be used for the final evaluation of a Pytorch model:

1. "train\_tuned" and
2. "test\_tuned".

`train_tuned()` is just a wrapper to `evaluate_hold_out` using the `train` data set. It is implemented as follows:

```
def train_tuned(  
    net,  
    train_dataset,  
    shuffle,  
    loss_function,  
    metric,  
    device=None,  
    show_batch_interval=10_000,
```

```

        path=None,
        task=None,
        writer=None,
    ):
        evaluate_hold_out(
            net=net,
            train_dataset=train_dataset,
            shuffle=shuffle,
            test_dataset=None,
            loss_function=loss_function,
            metric=metric,
            device=device,
            show_batch_interval=show_batch_interval,
            path=path,
            task=task,
            writer=writer,
        )

```

The `test_tuned()` procedure is implemented as follows:

```

def test_tuned(net, shuffle, test_dataset=None, loss_function=None,
               metric=None, device=None, path=None, task=None):
    batch_size_instance = net.batch_size
    removed_attributes, net = get_removed_attributes_and_base_net(net)
    if path is not None:
        net.load_state_dict(torch.load(path))
        net.eval()
    try:
        device = getDevice(device=device)
        if torch.cuda.is_available():
            device = "cuda:0"
            if torch.cuda.device_count() > 1:
                print("We will use", torch.cuda.device_count(), "GPUs!")
                net = nn.DataParallel(net)
        net.to(device)
        valloader = torch.utils.data.DataLoader(
            test_dataset, batch_size=int(batch_size_instance),
            shuffle=shuffle,
            num_workers=0
        )
        metric_value, loss = validate_one_epoch(
            net, valloader=valloader, loss_function=loss_function,

```

```
        metric=metric, device=device, task=task
    )
df_eval = loss
df_metric = metric_value
df_preds = np.nan
except Exception as err:
    print(f"Error in Net_Core. Call to test_tuned() failed. {err=},
          {type(err)=}")
    df_eval = np.nan
    df_metric = np.nan
    df_preds = np.nan
add_attributes(net, removed_attributes)
print(f"Final evaluation: Validation loss: {df_eval}")
print(f"Final evaluation: Validation metric: {df_metric}")
print("-----")
return df_eval, df_preds, df_metric
```

# References

- Bartz, Eva, Thomas Bartz-Beielstein, Martin Zaefllerer, and Olaf Mersmann, eds. 2022. *Hyperparameter Tuning for Machine and Deep Learning with R - A Practical Guide*. Springer.
- Bartz-Beielstein, Thomas. 2023. “PyTorch Hyperparameter Tuning with SPOT: Comparison with Ray Tuner and Default Hyperparameters on CIFAR10.” [https://github.com/sequential-parameter-optimization/spotPython/blob/main/notebooks/14\\_spot\\_ray\\_hpt\\_torch\\_cifar10.ipynb](https://github.com/sequential-parameter-optimization/spotPython/blob/main/notebooks/14_spot_ray_hpt_torch_cifar10.ipynb).
- Forrester, Alexander, András Sóbester, and Andy Keane. 2008. *Engineering Design via Surrogate Modelling*. Wiley.
- Santner, T J, B J Williams, and W I Notz. 2003. *The Design and Analysis of Computer Experiments*. Berlin, Heidelberg, New York: Springer.