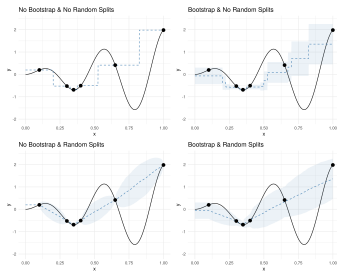
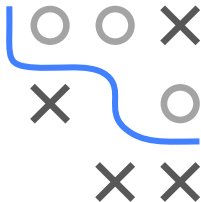


Optimization in Machine Learning

Bayesian Optimization Important Surrogate Models



Learning goals

- Search space / input data peculiarities in black box problems
- Gaussian process
- Random forest

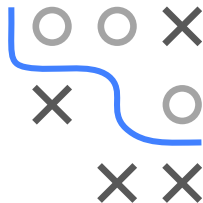
SURROGATE MODELS

Desiderata:

- Regression model (there are also classification approaches)
- Non-linear local model
- Accurate predictions (especially for small sample sizes)
- Often: uncertainty estimates
- Robust, works often well without human modeler intervention

Depending on the application:

- Can handle different types of inputs (numerical and categorical)
- Can handle dependencies (i.e., hierarchical input)



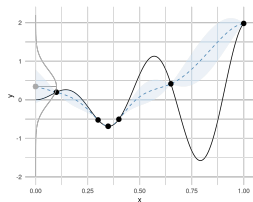
GAUSSIAN PROCESS

Posterior predictive distribution for test point $\mathbf{x} \in \mathcal{S}$ under zero mean:

$$Y(\mathbf{x}) \mid \mathbf{x}, \mathcal{D}^{[t]} \sim \mathcal{N} \left(\hat{f}(\mathbf{x}), \hat{s}^2(\mathbf{x}) \right)$$

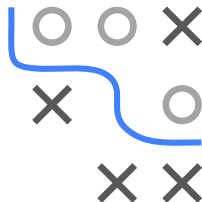
with

$$\begin{aligned}\hat{f}(\mathbf{x}) &= \mathbf{k}(\mathbf{x})^\top \mathbf{K}^{-1} \mathbf{y} \\ \hat{s}^2(\mathbf{x}) &= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^\top \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})\end{aligned}$$



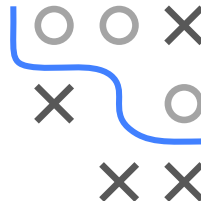
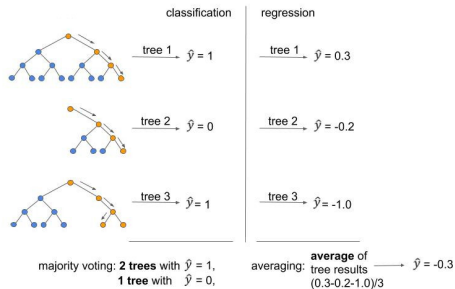
Note: \mathbf{x} here denotes the test input. $\mathbf{k}(\mathbf{x}) := (k(\mathbf{x}, \mathbf{x}^{[1]}), \dots, k(\mathbf{x}, \mathbf{x}^{[t]}))^\top$. $\mathbf{y} := (y^{[1]}, \dots, y^{[t]})^\top$.

Kernel / Gram matrix $\mathbf{K} := \left(k(\mathbf{x}^{[i]}, \mathbf{x}^{[j]}) \right)_{i,j}$ where $i, j \in \{1, \dots, t\}$.



RANDOM FOREST

- Bagging ensemble
- Fit B decision trees on bootstrap samples
- Feature subsampling

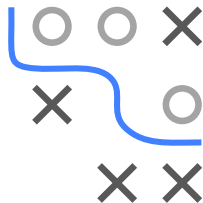


“extratrees” / random splits:

- Choose split location uniformly at random
- Results in a “smoother” mean prediction

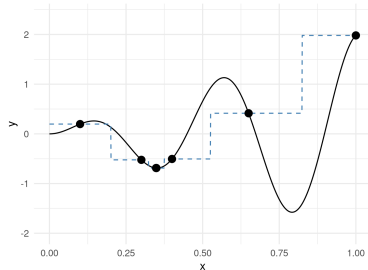
RANDOM FOREST - MEAN AND VARIANCE

- Let $\hat{f}_b : \mathcal{S} \rightarrow \mathbb{R}$ be the mean prediction of a decision tree b (mean of all data points in the same node as observation $\mathbf{x} \in \mathcal{S}$)
- Let $\hat{s}_b^2 : \mathcal{S} \rightarrow \mathbb{R}$ be the variance prediction (variance of all data points in the same node as observation $\mathbf{x} \in \mathcal{S}$)
- Mean prediction of forest: $\hat{f} : \mathcal{S} \rightarrow \mathbb{R}, \mathbf{x} \mapsto \frac{1}{B} \sum_{b=1}^B \hat{f}_b(\mathbf{x})$
- Variance prediction of forest: $\hat{s}^2 : \mathcal{S} \rightarrow \mathbb{R},$
 $\mathbf{x} \mapsto \left(\frac{1}{B} \sum_{b=1}^B \hat{s}_b^2(\mathbf{x}) + \hat{f}_b(\mathbf{x})^2 \right) - \hat{f}(\mathbf{x})^2$
(law of total variance assuming a mixture of B models)
- Alternative variance estimator:
 - (infinitesimal) Jackknife
- Variance prediction derived from randomness of individual trees
 - Bagging / bootstrap samples
 - Features sampled at random
 - (randomized split locations in the case of “extratrees”)

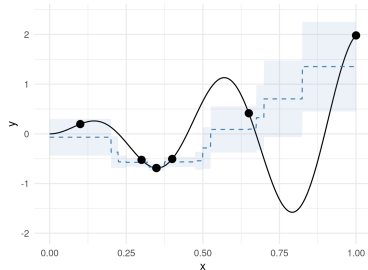


RANDOM FOREST - DIFFERENT CHOICES

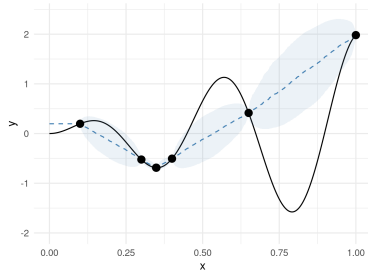
No Bootstrap & No Random Splits



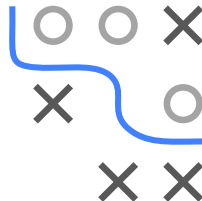
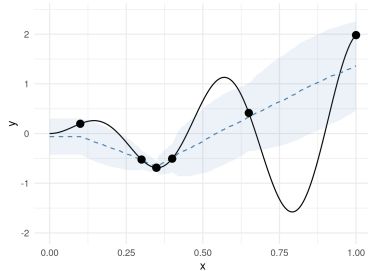
Bootstrap & No Random Splits



No Bootstrap & Random Splits



Bootstrap & Random Splits



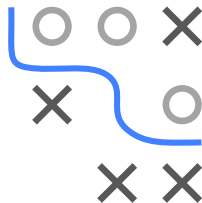
RANDOM FOREST

Pros:

- Cheap(er) to train
- Scales well with the number of data points
- Scales well with the number of dimensions
- Can easily handle hierarchical mixed spaces. Either via imputation or directly respecting dependencies in the tree structure
- Robust

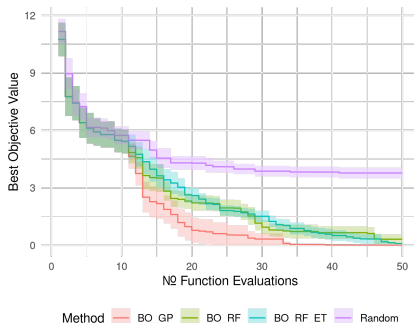
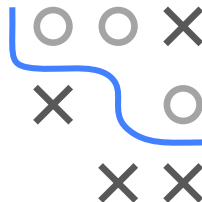
Cons:

- Suboptimal uncertainty estimates
- Not really Bayesian (no real posterior predictive distribution)
- Poor extrapolation



EXAMPLE

Minimize the 2D Ackley Function using BO_GP (GP with Matérn 3/2, EI), BO_RF (standard Random Forest, EI), BO_RF_ET (Random Forest with extratrees, EI) or a random search:



Strong BO_GP performance. BO_RF and BO_RF_ET not too bad either. BO_RF_ET maybe slightly better final performance than BO_RF.