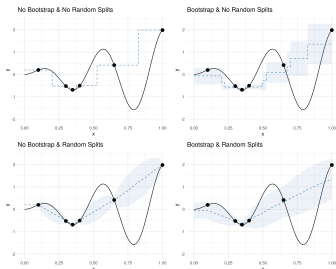


Optimization in Machine Learning

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 size)
- 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

The posterior mean and the posterior variance for a GP can be derived analytically:

- Let $\mathcal{D}^{[t]} = \{(\mathbf{x}^{[i]}, y^{[i]})\}_{i=1, \dots, t}$ be the data we are fitting the GP on
- Let $\mathbf{y} := (y^{[1]}, \dots, y^{[t]})$ be the vector of observed outputs
- For a covariance kernel $k(\mathbf{x}, \mathbf{x}')$, let $\mathbf{K} := (k(\mathbf{x}^{[i]}, \mathbf{x}^{[j]}))_{i,j}$ denote the **kernel (Gram) matrix** and $k(\mathbf{x}) := (k(\mathbf{x}, \mathbf{x}^{[1]}), \dots, k(\mathbf{x}, \mathbf{x}^{[t]}))^\top$
- Further, we assume a zero-mean GP prior

GAUSSIAN PROCESS

Example kernel functions:

- Radial basis function kernel (also known as Gauss kernel):

$$k(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{d(\mathbf{x}, \mathbf{x}')^2}{2l^2} \right)$$

- l length scale; $d(\cdot, \cdot)$ Euclidean distance
- infinitely differentiable - very “smooth”

- Matérn kernels:

$$k(\mathbf{x}, \mathbf{x}') = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\frac{\sqrt{2\nu}}{l} d(\mathbf{x}, \mathbf{x}') \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}}{l} d(\mathbf{x}, \mathbf{x}') \right)$$

- l length scale; $d(\cdot, \cdot)$ Euclidean distance; $K_\nu(\cdot)$ modified Bessel function; $\Gamma(\cdot)$ Gamma function
- for $\nu = 3/2$ once differentiable, for $\nu = 5/2$ twice differentiable
- Popular choice as a kernel function when using a GP as SM

GAUSSIAN PROCESS

The posterior predictive distribution for a new test point $\mathbf{x} \in \mathcal{S}$ under a GP is

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

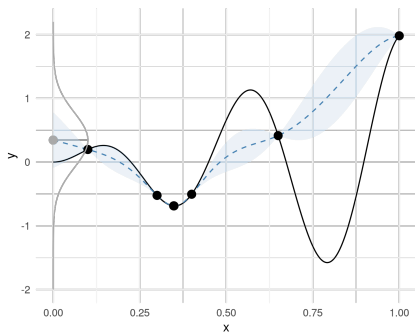
with

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

GAUSSIAN PROCESS

Pros:

- GPs yield well calibrated uncertainty estimates
- The posterior predictive distribution under a GP is normal



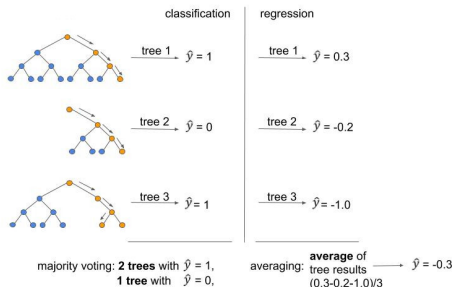
GAUSSIAN PROCESS

Cons:

- Vanilla GPs scale cubic in the number of data points
- GPs can natively only handle numeric features
Categorical features and dependencies require special handling via custom kernel
- GPs aren't that robust
In practice, “white-noise” models can occur where the posterior mean and variance is constant (except for the interpolation of training points)
- Performance can be sensitive to the choice of kernel and hyperparameters

RANDOM FOREST

- Bagging ensemble
- Fit B decision trees on bootstrap samples using different features sampled at random



“extratrees” / Random Splits:

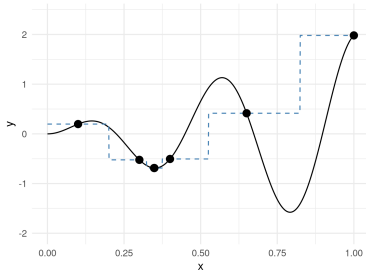
- Choose split location uniformly at random
- Results in a “smoother” mean prediction

RANDOM FOREST - MEAN AND VARIANCE

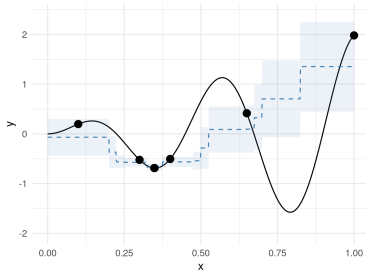
- Let $\hat{f}_b : \mathcal{S} \rightarrow \mathbb{R}$ be the mean prediction of a decision tree b (mean of all data points in the same node as observation $\mathbf{x} \in \mathcal{S}$)
- Let $\hat{s}_b^2 : \mathcal{S} \rightarrow \mathbb{R}$ be the variance prediction (variance of all data points in 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

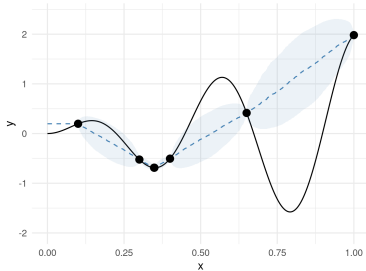
No Bootstrap & No Random Splits



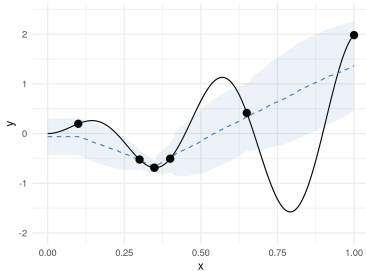
Bootstrap & No Random Splits



No Bootstrap & Random Splits



Bootstrap & Random Splits



RANDOM FOREST

Pros:

- Cheap 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

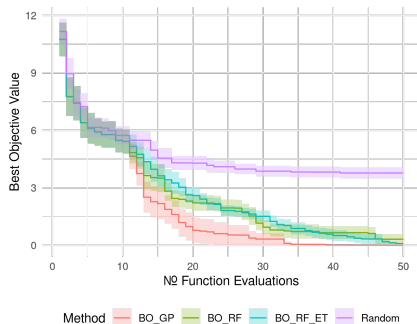
RANDOM FOREST

Cons:

- Poor 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.