

Hyperparameter Optimization

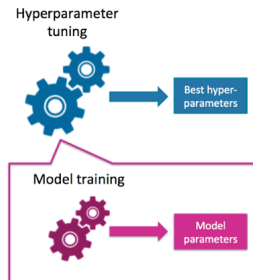
Dr. Josif Grabocka, Nicolas Schilling, Dr. Martin Wistuba

ISMLL, University of Hildesheim

Business Analytics

Hyperparameter Optimization

- ▶ Model parameters are learned to minimize the loss on a training dataset
- ▶ Yet several hyper-parameters exist that control model complexity, characteristics and optimization choices



Courtesy: [dato.com](https://www.dato.com)

Learning Problems

The ultimate task of machine learning is to learn a model

$$\mathcal{A}(D^{\text{train}}) = \arg \min_{\hat{y} \in \mathcal{M}} \mathcal{L}(\hat{y}, D^{\text{train}}) + \mathcal{R}(\hat{y})$$

- ▶ \mathcal{A} is a learning algorithm that returns a model
- ▶ \mathcal{M} is the given model space
- ▶ \mathcal{L} is a loss function, D^{train} is training data
- ▶ \mathcal{R} is a regularizer
- ▶ \mathcal{A} depends on the correct setting of hyperparameters λ , which *usually* cannot be determined analytically

Hyperparameter Optimization

We are interested in finding hyperparameters $\lambda^* \in \Lambda$ for which holds:

$$\lambda^* := \arg \min_{\lambda \in \Lambda} \mathcal{L}(\mathcal{A}_\lambda(D^{\text{train}}), D^{\text{val}}) := \arg \min_{\lambda \in \Lambda} f(\lambda, D) .$$

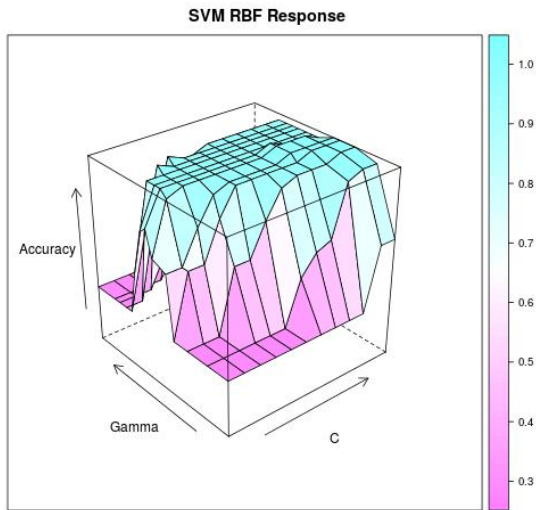
- ▶ search for the configuration yielding a learned model with best validation performance
- ▶ replace the loss of a learned model given λ by the black-box function $f(\lambda, D)$
- ▶ f is called **the response function**

Example

- ▶ Model: SVM
- ▶ Loss/Response: Misclassification Rate / Accuracy
- ▶ Hyperparameters:
 - ▶ tradeoff C , choice of Kernel, kernel width γ , kernel degree d

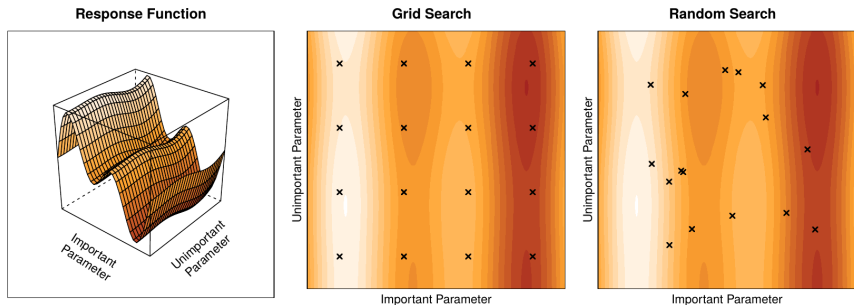
- ▶ Model: Neural Network
- ▶ Loss/Response: Misclassification Rate / Accuracy, MSE
- ▶ Hyperparameters:
 - ▶ Number of hidden layers, number of neurons per layer, choice of activation function, dropout rate, regularization penalty on weights, learning rate, choice of optimization algorithms

Example (II)



Grid and Random Search

- ▶ Need to try as few $\lambda \in \Lambda$, e.g.: in a grid, or randomly:



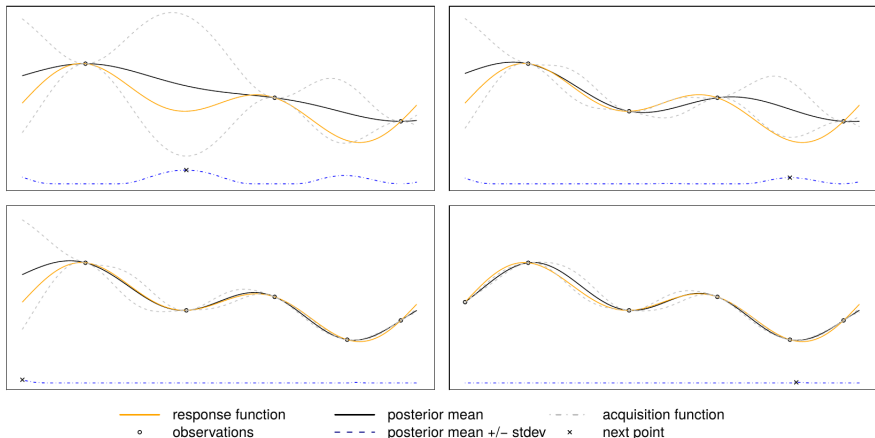
- ▶ and keep the hyper-parameters λ^* yielding smallest $f(\lambda^*, D)$.
- ▶ Why are those approaches not sufficient?

Sequential Model Based Optimization

Optimization of f is costly as it involves many runs of the learning algorithm, therefore SMBO is used

- ▶ learn a **surrogate model** on known observations $\Psi \approx f$
- ▶ query Ψ for regions of Λ where **expected improvement** is high
- ▶ evaluate f in these regions (run the learning algorithm)
- ▶ repeat the process until f is optimized

A sneak peek of SMBO ...



The Surrogate Model Ψ :

- ▶ learns from the history of N previous hyper-parameter evaluations

$$\mathcal{H} = \{(\lambda_1, f(\lambda_1, D)), (\lambda_2, f(\lambda_2, D)), \dots, (\lambda_N, f(\lambda_N, D))\}$$

- ▶ is a parametric model $\Psi(\lambda; \theta)$ that approximates the true response $f(\lambda, D)$ over observed history \mathcal{H} :

$$\arg \min_{\theta} \sum_{n=1}^N \mathcal{L}_{\Psi}(\Psi(\lambda_n; \theta), f(\lambda_n, D)) + \mathcal{R}(\theta)$$

- ▶ Ψ is often a Gaussian Process regressor

Gaussian Processes (GP)

- ▶ Given our history of configurations λ and responses $f(\lambda, D)$:

$$\Psi \sim \mathcal{N}(\mu(\lambda), k(\lambda, \lambda))$$

- ▶ And the covariance expressed as the squared exponential kernel:

$$k(\lambda, \lambda') = \exp\left(\frac{-\|\lambda - \lambda'\|^2}{2\sigma^2}\right)$$

- ▶ GP produces the mean and variance of a configuration λ_* :

$$\begin{aligned}\mu(\Psi(\lambda_*)) &= k_*^T K^{-1} f(\lambda) \\ \sigma^2(\Psi(\lambda_*)) &= k_{**} - k_*^T K^{-1} k_*\end{aligned}$$

- ▶ $k = k(\lambda, \lambda)$, $k_* = k(\lambda, \lambda_*)$ and $k_{**} = k(\lambda_*, \lambda_*)$

Sequential Model Based Optimization

Require: initial observation history \mathcal{H} , surrogate model Ψ , acquisition function a .

- 1: **return** Best hyperparameter configuration λ^*
- 2: **for** $t = 1$ to T **do**
- 3: Fit Ψ to \mathcal{H}
- 4: $\lambda^{\text{new}} = \arg \max_{\lambda \in \Lambda} a(\lambda, \Psi(\lambda))$
- 5: Evaluate $f(\lambda^{\text{new}}, D)$
- 6: $\mathcal{H} = \mathcal{H} \cup (\lambda^{\text{new}}, f(\lambda^{\text{new}}, D))$
- 7: **end for**
- 8: **return** $\lambda^* = \arg \min_{\lambda} f(\lambda, D)$

-
- Note: f can also be classification accuracy and the problem becomes a maximization, i.e. $\lambda^* = \arg \max_{\lambda} f(\lambda, D)$

How to choose next hyper-parameter candidate?

- ▶ Knowing the mean $\mu(\Psi(\lambda))$ and $\sigma(\Psi(\lambda))$ of the surrogate model
- ▶ We use an acquisition function a over all possible hyper-parameter configurations $\lambda^{\text{new}} = \arg \max_{\lambda \in \Lambda} a(\lambda, \Psi(\lambda))$
- ▶ It is important to balance **exploration vs. exploitation** of the search space
- ▶ Alternative 1: The Lower Confidence Bound acquisition

$$a(\lambda) = -\mu(\Psi(\lambda)) + \beta\sigma(\Psi(\lambda))$$

Acquisition Choice: Expected Improvement

- ▶ Define the standardized improvement of a configuration λ as :

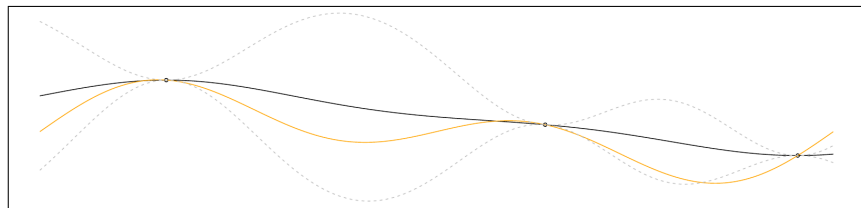
$$Z(\lambda) = \frac{f^{\min} - \mu(\Psi(\lambda))}{\sigma(\Psi(\lambda))}$$

- ▶ The expected improvement is:

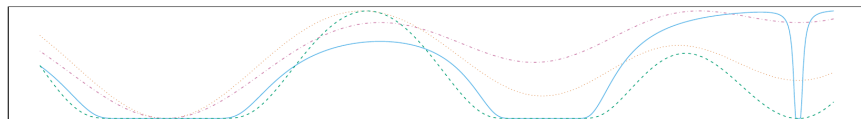
$$EI(\lambda) = \begin{cases} \sigma(\Psi(\lambda)) (Z(\lambda) \cdot \Phi(Z(\lambda)) + \phi(Z(\lambda))) & \sigma^2(\Psi(\lambda)) > 0 \\ 0 & \text{otherwise} \end{cases}$$

- ▶ where $\Phi(Z(\lambda))$: Gaussian density, $\phi(Z(\lambda))$: Gaussian cumulative distribution function

Choice of acquisitions

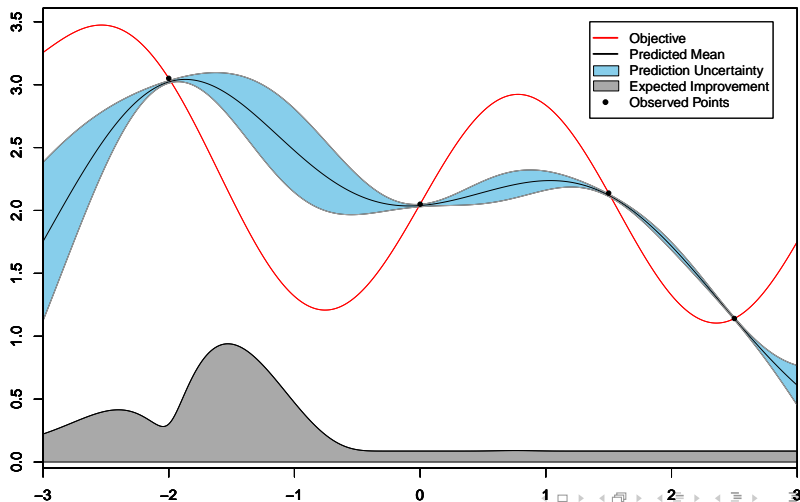


— response function ◦ observations — posterior mean - - - posterior mean +/- stdev

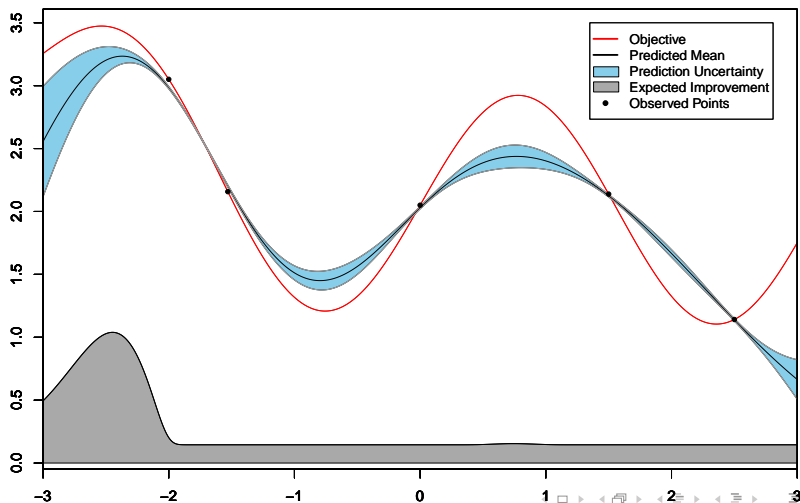


— PI - - - EI GP-LCB ($\beta=2$) - - - GP-LCB ($\beta=0.5$)

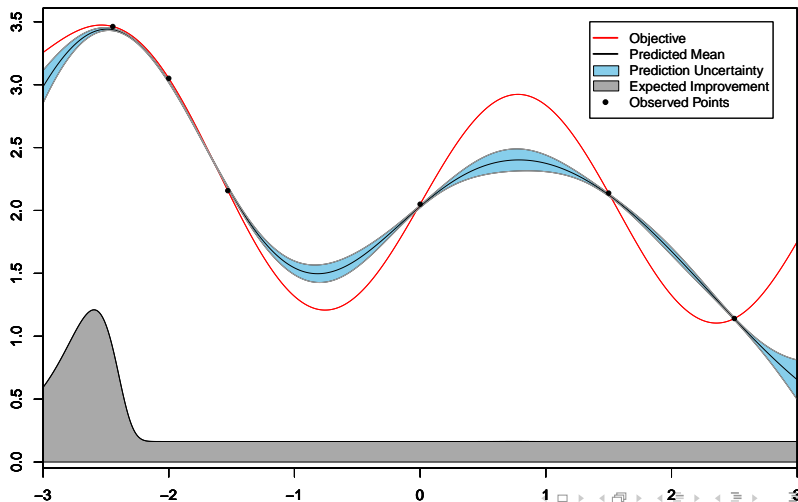
Gaussian Process Learning with SMBO (I)



Gaussian Process Learning with SMBO (II)



Gaussian Process Learning with SMBO (III)



Gaussian Process Learning with SMBO (IV)

