

Two-stage transfer surrogate model for automatic hyperparameter optimization.

- 作者: Martin Wistuba, Nicolas Schilling, and Lars Schmidt-Thieme
- 机构: University of Hildesheim
- 会议: ECML2016
- 地址: <http://arxiv.org/abs/1802.02219>
- 代码: <https://github.com/wistuba/TST>

论文主要内容

摘要

通过在其他data sets上获得的knowledge来提升在新data sets上的贝叶斯优化。文章通过两个stage来提升transfer knowledge。第一个stage拟合hyperparameter response function（理解成一个回归模型，超参=>score）。第二个stage利用第一个stage的所有model组装成一个surrogate model.

研究内容

定义:

$$\mathcal{A}_{\lambda} \left(D^{\text{train}} \right) = \arg \min_{M_{\lambda} \in \mathcal{M}} \mathcal{L} \left(M_{\lambda}, D^{\text{train}} \right) + \mathcal{R} \left(M_{\lambda} \right)$$

hyperparameter response function: $f_D(\lambda) = \mathcal{L} \left(\mathcal{A}_{\lambda} \left(D^{\text{train}} \right), D^{\text{valid}} \right)$

Algorithm 1 Sequential Model-based Optimization

Input: Hyperparameter space Λ , observation history \mathcal{H} , number of trials T , acquisition function a , surrogate model Ψ .

Output: Best hyperparameter configuration found.

- 1: **for** $t = 1$ to T **do**
 - 2: Fit Ψ_{t+1} to \mathcal{H}_t
 - 3: $\lambda \leftarrow \arg \max_{\lambda \in \Lambda} a \left(\mu \left(\Psi_{t+1} \left(\lambda \right) \right), \sigma \left(\Psi_{t+1} \left(\lambda \right) \right), f^{\min} \right)$
 - 4: Evaluate $f \left(\lambda \right)$
 - 5: $\mathcal{H}_{t+1} \leftarrow \mathcal{H}_t \cup \{ (\lambda, f \left(\lambda \right)) \}$
 - 6: **if** $f \left(\lambda \right) < f^{\min}$ **then**
 - 7: $\lambda^{\min}, f^{\min} \leftarrow \lambda, f \left(\lambda \right)$
 - 8: **return** λ^{\min}
-

方法

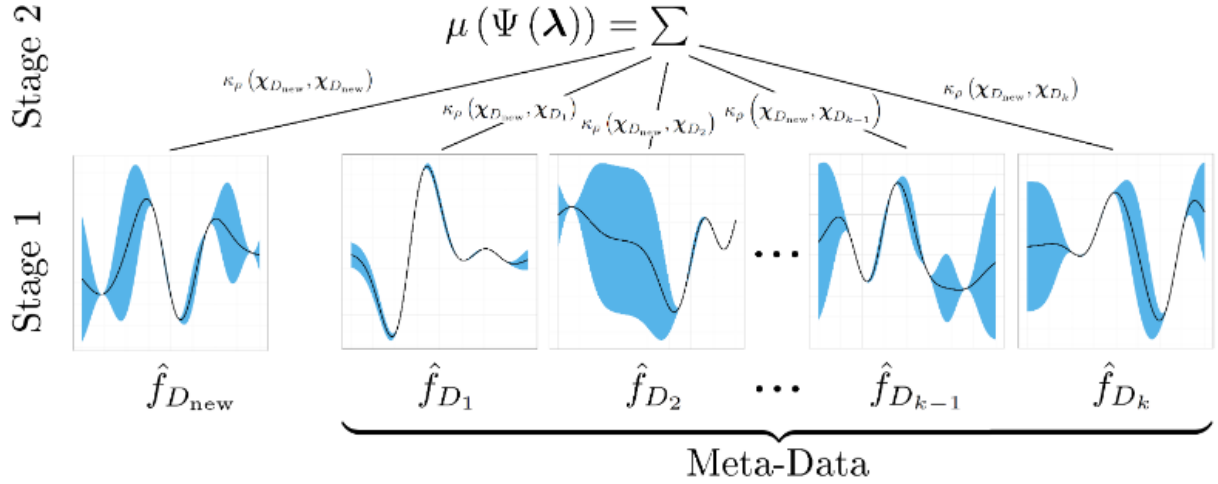


Fig. 1. At the first stage the hyperparameter response functions of the new data set D_{new} as well as data sets $\mathcal{D} = \{D_1, \dots, D_k\}$ used for previous experiments are approximated using known evaluations. At the second stage the predictions of each individual model \hat{f}_D are taken into account weighted by the similarity between D and D_{new} to determine the final predicted score.

第一个stage

meta-dataset

每个data set有288个Meta-instances（超参，score）、归一化label，平衡不同数据任务的scale
拟合meta-instance，使用高斯过程模型（可以使用任何其他ML model）

model: \hat{f}_D

第二个stage

组合第一个stage的所有model: Nadaraya-Watsonkernel-weighted average

代理模型（surrogate model）的均值函数、核函数：取candidate最大的

$$\mu(\Psi(\lambda)) = \frac{\sum_{D \in \mathcal{D} \cup \{D_{\text{new}}\}} \kappa_\rho(\mathbf{x}_{D_{\text{new}}}, \mathbf{x}_D) \hat{f}_D(\lambda)}{\sum_{D \in \mathcal{D} \cup \{D_{\text{new}}\}} \kappa_\rho(\mathbf{x}_{D_{\text{new}}}, \mathbf{x}_D)}$$

$$\sigma(\Psi(\lambda)) = \sigma(\hat{f}_{D_{\text{new}}}(\lambda))$$

其中，

$$\kappa_\rho(\mathbf{x}_D, \mathbf{x}_{D'}) = \delta\left(\frac{\|\mathbf{x}_D - \mathbf{x}_{D'}\|_2}{\rho}\right)$$

$$\delta(t) = \begin{cases} \frac{3}{4}(1 - t^2) & \text{if } t \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

χ_D 是用来描述data的向量

后面就是使用EI（Expected Improvement） acquisition function过程

Data Set Description

1. 使用Meta-Features：(TST-M)

Table 1. The list of all meta-features used by us.

Meta-Features	
Number of Classes	Class Probability Max
Number of Instances	Class Probability Mean
Log Number of Instances	Class Probability Standard Deviation
Number of Features	Kurtosis Min
Log Number of Features	Kurtosis Max
Data Set Dimensionality	Kurtosis Mean
Log Data Set Dimensionality	Kurtosis Standard Deviation
Inverse Data Set Dimensionality	Skewness Min
Log Inverse Data Set Dimensionality	Skewness Max
Class Cross Entropy	Skewness Mean
Class Probability Min	Skewness Standard Deviation

2. Pairwise Hyperparameter Performance Rankings: (TST-R)

$$(\chi_D)_{j+(i-1)t} = \begin{cases} 1 & \text{if } \hat{f}_D(\lambda_i) > \hat{f}_D(\lambda_j) \\ 0 & \text{otherwise} \end{cases}$$

t个obsversion内的超参的rank来描述一个data

实验结果

metric:

1. average rank

2. average distance to the global minimum: $\text{ADTM}(\Lambda_t, \mathcal{D}) = \sum_{D \in \mathcal{D}} \min_{\lambda \in \Lambda_t} \frac{f_D(\lambda) - f_D^{\min}}{f_D^{\max} - f_D^{\min}}$

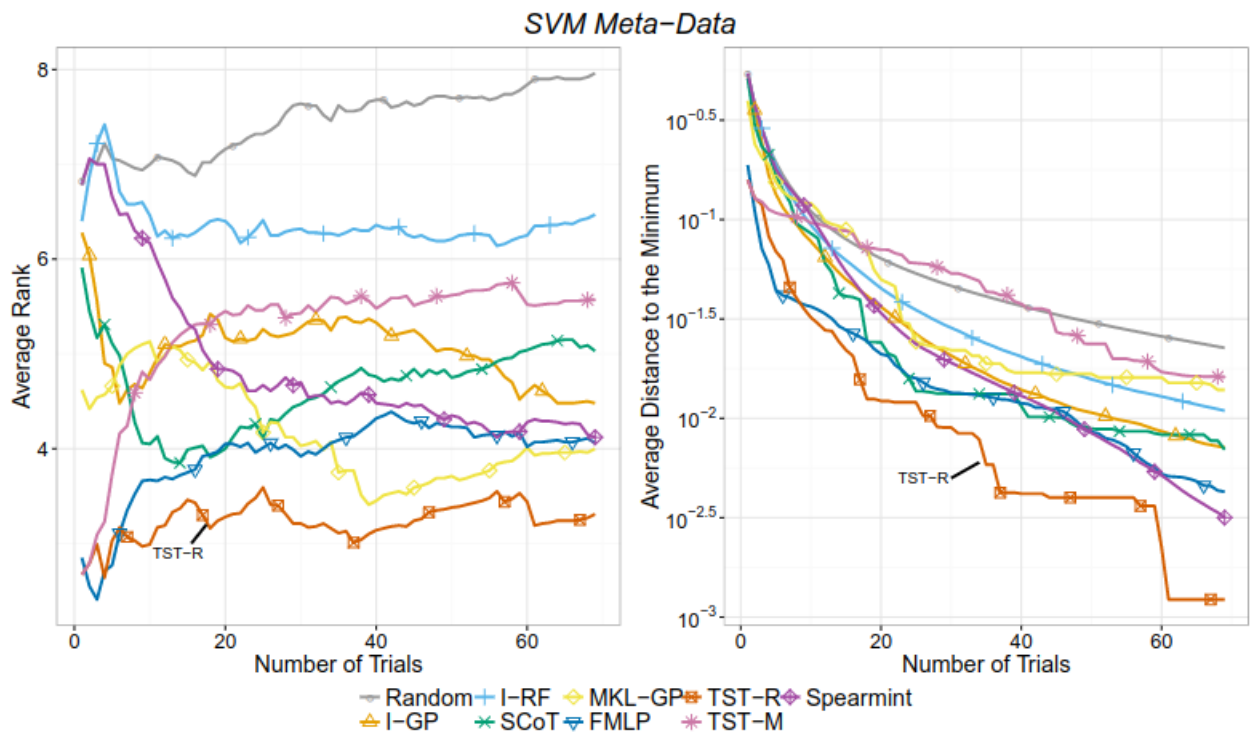


Fig. 2. Our proposed transfer surrogate model TST-R provides the best performance with respect to both evaluation measures for the task of hyperparameter tuning. For both metrics, the smaller the better.

less descriptive? NO

TST-R的distance会在每次增加observe之后更新，只关注在knowledge在新数据的某时间点上 meta-feature后期没有把新数据的信息全考虑进去。

Weka Meta-Data

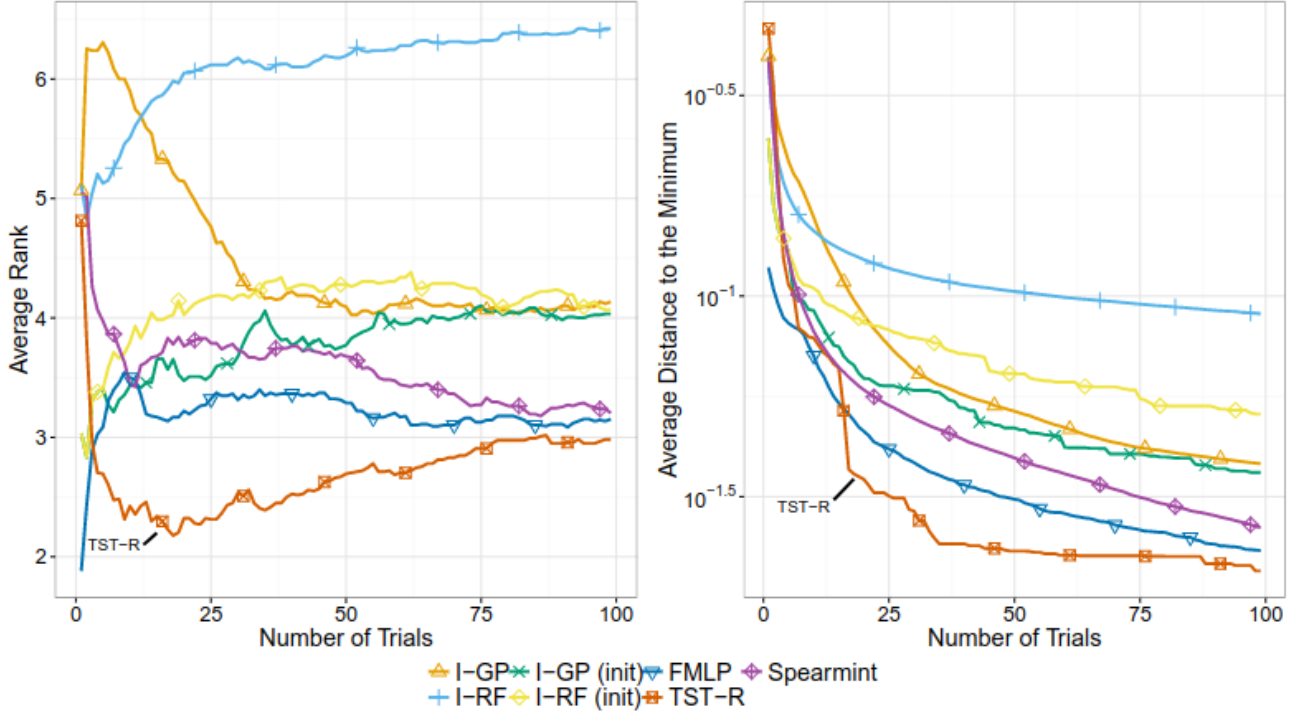


Fig. 3. Our approach TST-R also outperforms the competitor methods for the task of combined algorithm selection and hyperparameter tuning. Surrogate models that use Gaussian processes that train over the whole meta-data are not feasible for this data set [31]. Therefore, we consider I-GP and I-RF with meta-learning initialization.