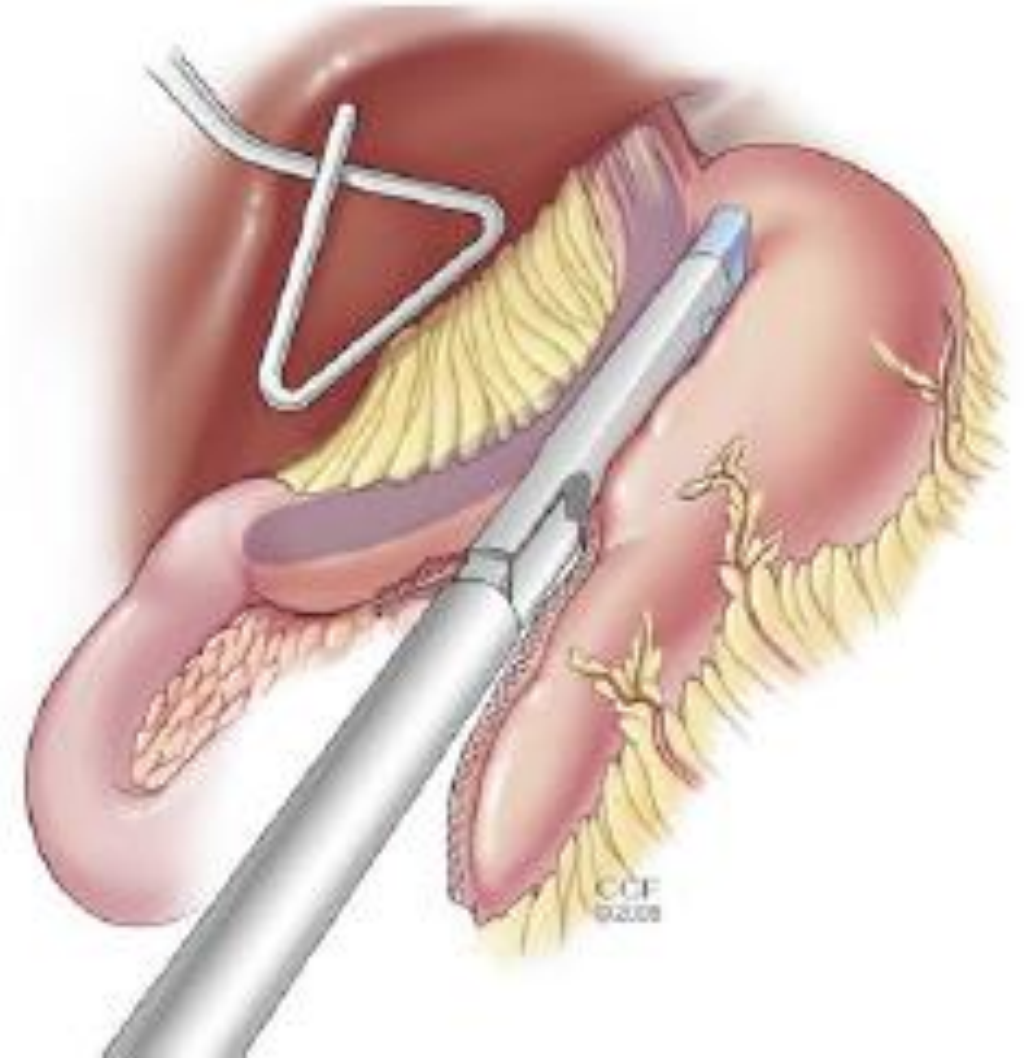


# TPOT's performance for Biomedical Data

Tim Beishuizen

Supervisor: Joaquin Vanschoren



# Introduction

- Computational Biology
  - Data analysis framework

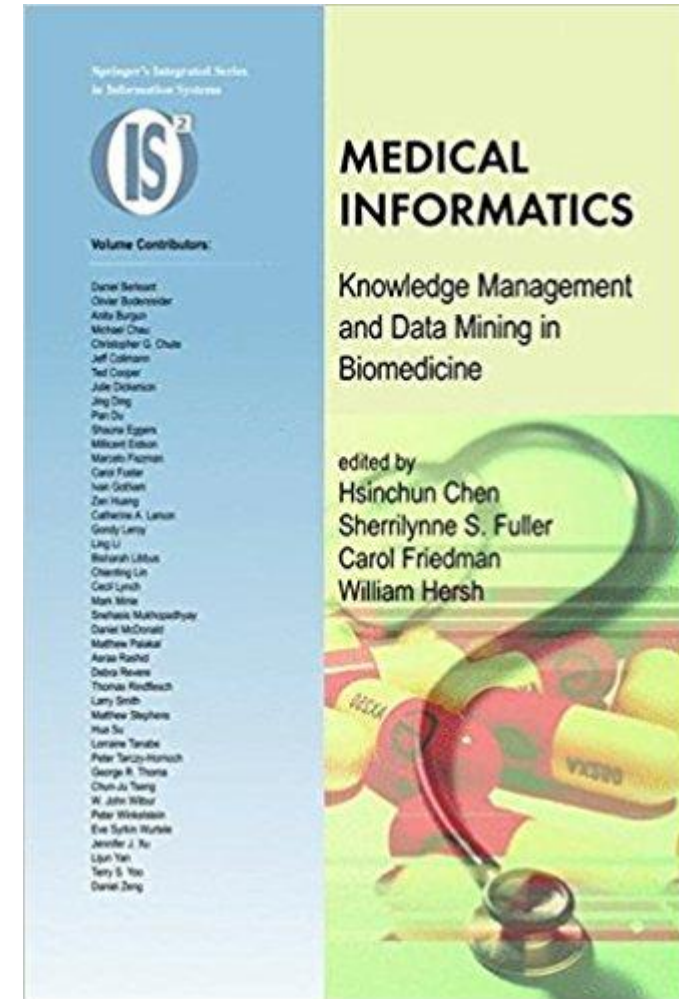
## Content

- Biomedical Data
  - Bariatric dataset
- Automated Machine learning
  - TPOT



# Biomedical Data

- Challenges
  - Volume
  - Dimensionality
  - Complexity
  - Heterogeneity
  - Quality



Chen, Hsinchun, et al., eds. *Medical informatics: knowledge management and data mining in biomedicine*. Vol. 8. Springer Science & Business Media, 2006.

# Bariatric Data set

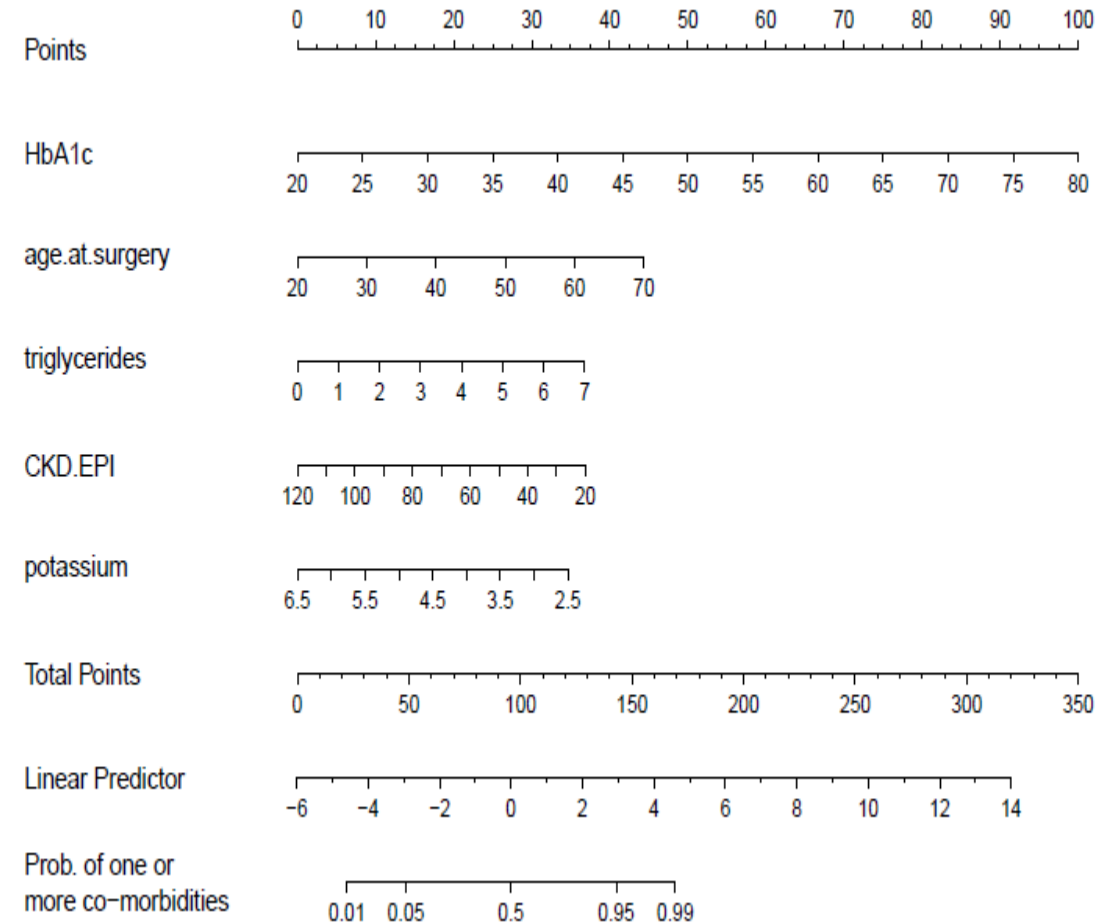
- Bariatric surgery
  - Catharina Hospital
- Two data sets
  - DATO set
    - Co-morbidity presence
    - Basic health variables
  - Lab set
    - Biomarkers
- Challenges
  - Combining two data sets
  - Missing values
  - High error rate (5%)



catharina  
ziekenhuis

# Bariatric Data set

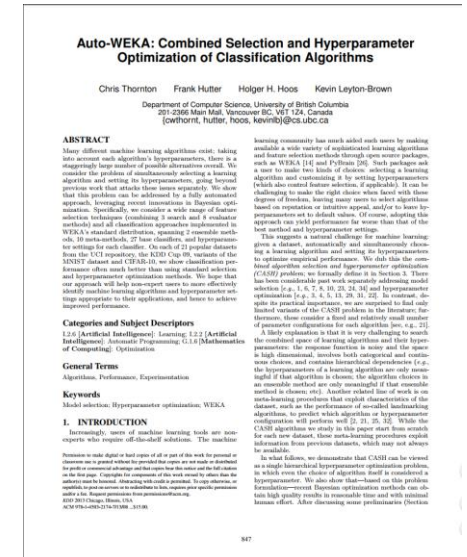
- Previous project
- Challenges
  - Dropped non matching data
  - Dropped missing values
  - Ignored missing data
- Logistic regression
- Model nomogram



# Automated Machine Learning

- Machine learning pipelines
  - Algorithm selection
  - Hyperparameter Optimization
  - Preprocessing
- Combined Algorithm Selection and Hyperparameter Optimization (CASH)
  - Find the best automatically
  - Meta-learning approach
  - Hyperparameter optimization
  - Preprocessing

- Auto-WEKA



Thornton, Chris, et al. "Auto-WEKA: Combined selection and hyperparameter optimization of classification algorithms." *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2013.

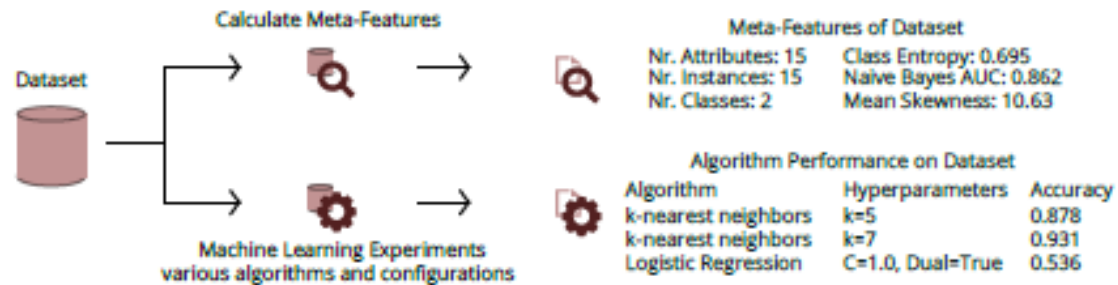
# Automated Machine Learning

- Meta-learning
  - Find performance changes
  - Time and space constraints

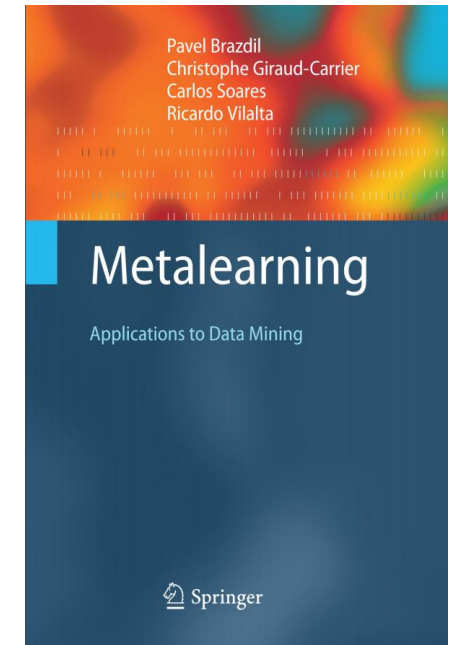
## 1. Collect Datasets



## 2. Compute metadata for each dataset



## 3. Create meta-dataset and learn a meta-model

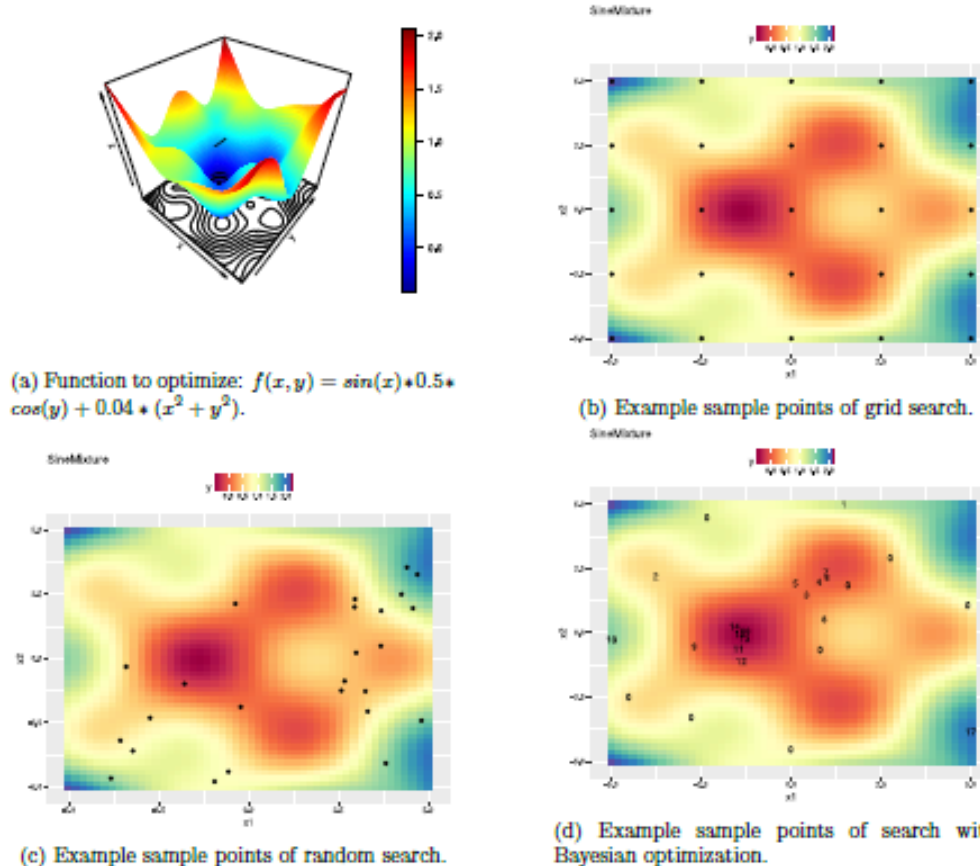


Brazdil, Pavel, et al. "Development of metalearning systems for algorithm recommendation." *Metalearning: Applications to Data Mining* (2009): 31-59.



# Automated Machine Learning

- Hyperparameter optimization
  - Grid search
  - Random search
  - Bayesian optimization



## Practical Bayesian Optimization of Machine Learning Algorithms

Jasper Snoek  
Department of Computer Science  
University of Toronto  
jasper@cs.toronto.edu

Hugo Larochelle  
Department of Computer Science  
University of Sherbrooke  
hugo.larochelle@usherbrooke.ca

Ryan P. Adams  
School of Engineering and Applied Sciences  
Harvard University  
rpads@seas.harvard.edu

### Abstract

The use of machine learning algorithms frequently involves careful tuning of learning parameters and model hyperparameters. Unfortunately, this tuning is often a "black art" requiring expert experience, rules of thumb, or sometimes brute-force search. There is therefore great appeal for automatic approaches that can optimize the performance of any given learning algorithm to the problem at hand. In this work, we consider this problem through the framework of Bayesian optimization, in which a learning algorithm's generalization performance is modeled as a sample from a Gaussian process (GP). We show that certain choices for the nature of the GP, such as the type of kernel and the treatment of its hyperparameters, can play a crucial role in obtaining a good optimizer that can achieve expert-level performance. We describe new algorithms that take into account the variable cost (duration) of learning algorithm experiments and that can leverage the presence of multiple cores for parallel experimentation. We show that these proposed algorithms improve on previous automatic procedures and can reach or surpass human expert-level optimization for many algorithms including latent Dirichlet allocation, structured SVMs and convolutional neural networks.

### 1 Introduction

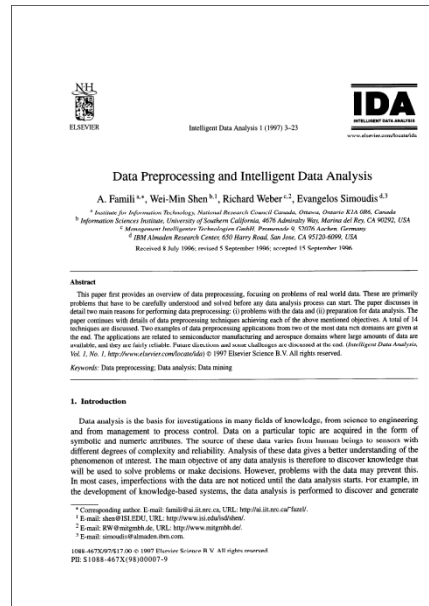
Machine learning algorithms are rarely parameter-free: parameters controlling the rate of learning or the capacity of the underlying model must often be specified. These parameters are often considered nuisances, making it appealing to develop machine learning algorithms with fewer of them. Another, more flexible take on this issue is to view the optimization of such parameters as a procedure to be automated. Specifically, we could view such tuning as the optimization of an unknown black-box function and invoke algorithms developed for such problems. A good choice is Bayesian optimization [1], which has been shown to outperform other state of the art global optimization algorithms on a number of challenging optimization benchmark functions [2]. For continuous functions, Bayesian optimization typically works by assuming the unknown function was sampled from a Gaussian process and maintains a posterior distribution for this function as observations are made or, in our case, as the results of running learning algorithm experiments with different hyperparameters are observed. To pick the hyperparameters of the next experiment, one can optimize the expected improvement (EI) [1] over the current best result or the Gaussian process upper confidence bound (UCB) [3]. EI and UCB have been shown to be efficient in the number of function evaluations required to find the global optimum of many multimodal black-box functions [4, 3].

Snoek, Jasper, Hugo Larochelle, and Ryan P. Adams. "Practical bayesian optimization of machine learning algorithms." *Advances in neural information processing systems*. 2012.

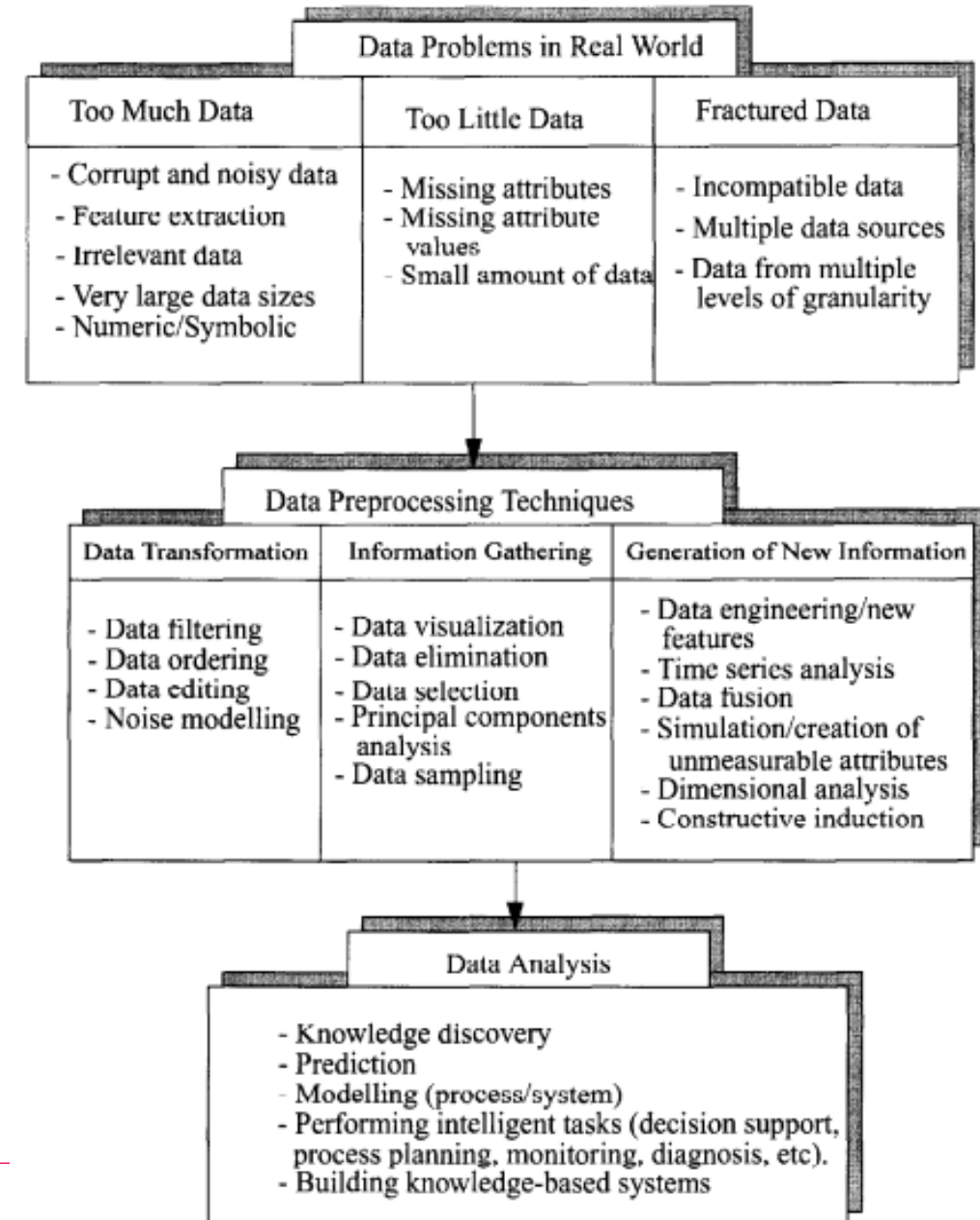


# Automated Machine Learning

- Preprocessing
  - Problems
  - Techniques

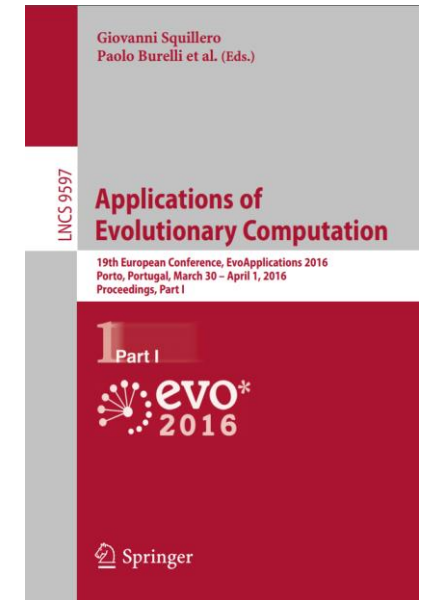
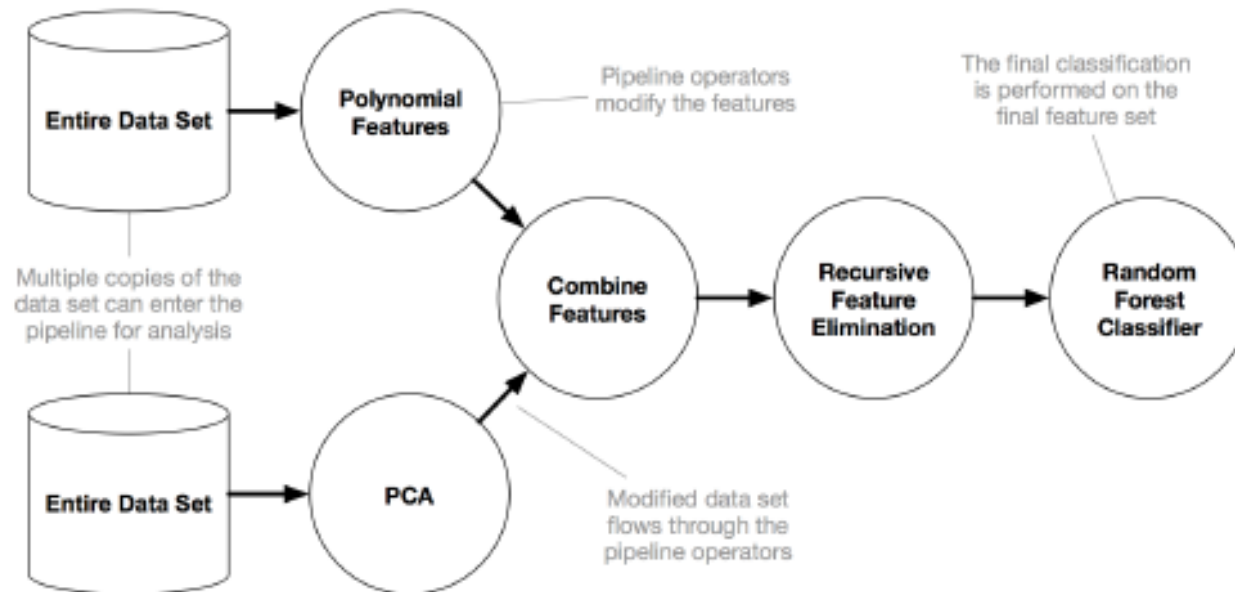


Famili, A., et al. "Data preprocessing and intelligent data analysis." *Intelligent data analysis* 1.1-4 (1997): 3-23.



# TPOT

- Tree-based Pipeline Optimization Tool
  - Automated Machine Learning tool
  - Scikit-learn
  - Evaluation, selection and mutations



Olson, Randal S., et al.  
"Automating biomedical data science through tree-based pipeline optimization." *European Conference on the Applications of Evolutionary Computation*. Springer, Cham, 2016.

# Research Question

- How does TPOT perform on specific biomedical data set problems and how can it be improved?
- Hypothesis
  - Combining data sets
  - Erroneous values
  - Missing values
- Improvement on missing values
  - Extrapolation
  - Nearest Neighbours