

[Supplementary Document] Configuration Space Reduction in Automated Machine Learning Using Relative Landmarking

Tien-Dung Nguyen¹, Bogdan Gabrys², and Katarzyna Musial²

Advanced Analytics Institute, University of Technology Sydney, Sydney, Australia

TienDung.Nguyen-2@student.uts.edu.au¹

{Bogdan.Gabrys, Katarzyna.Musial-Gabrys}@uts.edu.au²

This document provides more details about the experiments for the paper entitled *Configuration Space Reduction in Automated Machine Learning Using Relative Landmarking*.

Table 1, 2 and 3 presents the best pipelines found using different configuration spaces.

To facilitate the study, we use prior evaluations of AutoWeka4MCPS with 2 hours optimisation time, 1GB memory, using SMAC as the ML pipeline composition and optimisation method over 20 datasets. We extract the mean error rate of predictors based on the error rate of their ML pipelines within 20 datasets in the prior evaluations. Based on the mean error rate of the predictors, we generate the ranking of predictors across 20 datasets, as shown in Figure 1. We can see that there is no predictor has the best ranking across all datasets. For example, the component *Logistic* has the top ranking in the cases of the datasets *abalone*, *semeion* and *waveform*. However, this component ranks 19th and 26th in the case of the datasets *convex* and *secom*.

Table 1. The best ML pipelines found using different methods to design configuration spaces.

Dataset	r30	no avatar	avatar
abalone	CustomReplaceMissingValues → RandomSubset → Resample → Logistic → Bagging	SimpleLogistic	SMO
adult	J48	PART	PART
amazon	CustomReplaceMissingValues → Normalize → RandomSubset → NaiveBayesMultinomial → RandomSubSpace	NaiveBayes	NaiveBayesMultinomial
car	SMO → MultiClassClassifier	SMO	SMO
cifar10small	RandomForest → MultiClassClassifier	DecisionStump	RandomForest
convex	RandomForest → AdaBoostM1	VotedPerceptron	RandomForest
dexter	DecisionStump → AdaBoostM1	NaiveBayesMultinomial	SGD
dorothea	OneR → RandomSubSpace	PART	DecisionStump
gcredit	LMT → Bagging	SMO	SMO
gisette	VotedPerceptron → RandomSubSpace	VotedPerceptron	VotedPerceptron
kddcup	-	DecisionStump	ClassBalancer → RemoveOutliers → InterquartileRange → AttributeSelection → Resample → PART
krvsnp	Jrip → AdaBoostM1	RandomForest	J48
madelon	PrincipalComponents → IBk → LogitBoost	Jrip	Jrip
mnist	CustomReplaceMissingValues → Center → J48 → AdaBoostM1	-	PART
secom	J48 → AdaBoostM1	ClassBalancer → EMImputation → Normalize → PrincipalComponents → Kstar → MultiClassClassifier	SimpleLogistic
semeion	CustomReplaceMissingValues → PrincipalComponents → SMO	SMO	SMO
shuttle	RandomForest → AdaBoostM1	RandomForest	RandomForest
waveform	RemoveOutliers → InterquartileRange → Normalize → SMO → AttributeSelectedClassifier	LMT	SimpleLogistic
winequality	RandomForest → AdaBoostM1	RandomForest	Kstar
yeast	RandomForest → Bagging	RandomForest	RandomForest

Table 2. The best ML pipelines found using different methods to design configuration spaces.

Dataset	L-k1	L-k4	L-k8	L-k10	L-k19
abalone	CustomReplaceMissingValues → Normalize → RandomSubset → SimpleLogistic	ClassBalancer → RemoveOutliers → InterquartileRange → Normalize → RandomSubset → SimpleLogistic	REPTree	RandomForest	PART
adult	RemoveOutliers → InterquartileRange → Center → Logistic	Logistic	NaiveBayes	PART	PART
amazon	ClassBalancer → RemoveOutliers → InterquartileRange → Center → Logistic	SMO	NaiveBayesMultinomial	NaiveBayesMultinomial	NaiveBayesMultinomial
car	SMO	SMO	SMO	SMO	SMO
cifar10small	RandomSubset → PeriodicSampling → Logistic	NaiveBayesMultinomial	RandomForest	NaiveBayes	-
convex	-	Jrip	SMO	RandomForest	RandomForest
dexter	SMO	SMO	SMO	SMO	SimpleLogistic
dorothea	-	-	-	DecisionStump	DecisionStump
gcredit	LMT	NaiveBayes	NaiveBayes	SMO	SMO
gisette	Logistic	-	Jrip	RandomForest	Logistic
kddcup	-	-	-	-	-
krvskp	SMO	SMO	J48	J48	RandomForest
madelon	SpreadSubsample → CustomReplaceMissingValues → Normalize → RandomSubset → Resample → VotedPerceptron	Jrip	Jrip	Jrip	Jrip
mnist	-	-	-	SMO	SimpleLogistic
secom	-	ZeroR	ZeroR	DecisionTable	LMT
semeion	SMO	SMO	SMO	SMO	SMO
shuttle	SMO	Jrip	PART	RandomForest	RandomForest
waveform	SMO	SMO	SimpleLogistic	SimpleLogistic	LMT
winequality	SMO	Kstar	SMO	RandomForest	Kstar
yeast	Logistic	SMO	RandomForest	RandomForest	RandomForest

Table 3. The best ML pipelines found using different methods to design configuration spaces.

Dataset	O-k1	O-k4	O-k8	O-k10	O-k19
abalone	Logistic	DecisionTable	SimpleLogistic	PART	MultilayerPerceptron
adult	-	PART	J48	J48	PART
amazon	-	NaiveBayesMultinomial	NaiveBayesMultinomial	NaiveBayesMultinomial	NaiveBayesMultinomial
car	LMT	SMO	SMO	SMO	SMO
cifar10small	-	RandomForest	RandomForest	NaiveBayes	RandomForest
convex	-	RandomForest	SMO	RandomForest	RandomForest
dexter	SGD	SGD	VotedPerceptron	SGD	SimpleLogistic
dorothea	SpreadSubsample → CustomReplaceMissingValues → Center → DecisionStump	NaiveBayes	NaiveBayes	NaiveBayes	OneR
gcredit	NaiveBayes	MultilayerPerceptron	SMO	MultilayerPerceptron	MultilayerPerceptron
gisette	VotedPerceptron	VotedPerceptron	VotedPerceptron	VotedPerceptron	VotedPerceptron
kddcup	MultilayerPerceptron	Ibk	DecisionStump	-	DecisionStump
krvsnp	LMT	J48	J48	J48	J48
madelon	Jrip	Jrip	Jrip	Jrip	Jrip
mnist	SMO	-	SMO	RandomForest	SMO
secom	ClassBalancer → CustomReplaceMissingValues → RemoveOutliers → InterquartileRange → Normalize → PeriodicSampling → DecisionStump	SpreadSubsample → Standardize → RandomSubset → PeriodicSampling → ZeroR	Kstar	VotedPerceptron	Kstar
semeion	Logistic	SMO	SMO	SMO	SMO
shuttle	RandomForest	RandomForest	RandomForest	RandomForest	RandomForest
waveform	Logistic	SimpleLogistic	SimpleLogistic	SimpleLogistic	SimpleLogistic
winequality	RandomForest	RandomForest	RandomForest	Kstar	Kstar
yeast	RandomForest	RandomForest	RandomForest	RandomForest	RandomForest

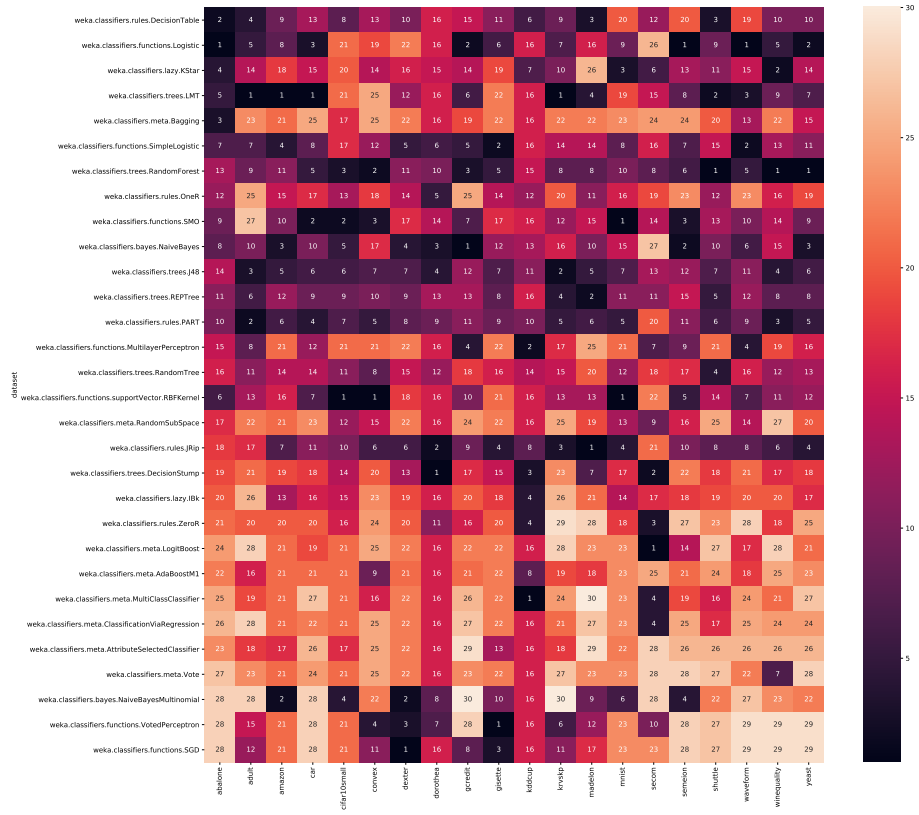


Fig. 1. Ranking of predictors based on mean error rate of their pipelines that is extracted from historical runs of 20 datasets within 2 hours optimisation time.