Data Pipeline Selection and Optimization DOLAP 2019



Alexandre Quemy IBM, Data and Al Politechnika Poznańska





Poland Software Lab



The usual workflow



The hard life of data scientists

- \rightarrow Dealing with missing value:
 - \rightarrow Discarding? Row? Column?
 - → Imputation? What imputation? Mean? Median? Model-based? What model?
- \rightarrow Imbalanced datasets:
 - \rightarrow Downsampling? Oversampling?
 - \rightarrow Nothing? What bias it implies?
- \rightarrow Data too large:
 - \rightarrow Dimensional reductions: what algorithm? PCA? normalization or not?
 - \rightarrow Subsampling: what technique? what bias?
- \rightarrow Outliers detection and curation:
 - \rightarrow What threshold? What deviation measure?
 - \rightarrow Trimming? Truncating? Censoring? Winsorizing?
- \rightarrow Encoding for method domain requirements:
 - \rightarrow Discretization? Grid? What step? Cluster? What method? What hyperparameter?
 - \rightarrow Categorial encoder? Binary? Hot-One? Helmert? Backward Difference?
- \rightarrow NLP:
 - \rightarrow How many tokens?
 - \rightarrow Size of m-grams?

The usual workflow



The workflow proposed in the paper



The workflow proposed in the paper



Pipeline prototype



Rebalance: 4 operators Normalize: 5 operators Features: 4 operators

Configuration space: 4750 configurations

Baseline: (Id, Id, Id)

Protocol

- •Datasets: Breast, Iris, Wine.
- •Methods: SVM, Random Forest, Neural Network, Decision Tree.
- •Dataset split: 60% for training set, 40% for test set.
- •Pipeline configuration space size: 4750 configurations.
- •Performance metric: Cross-validation accuracy
- •Metaoptimizer: Tree Parzen Estimator (hyperopt)
- •Budget: 100 configurations (~2% of the space)

No algorithm hyperparameter tuning!

 \Rightarrow We want to quantify the influence of data pipeline

Exhaustive search to compare between baseline and max score.

Results

Density for Random Forest on Breast









Accuracy for SVM on Wine

	Baseline	Exhaustive	SMBO	SMBO (norm.)	Imp. Inter.		
Iris							
SVM	0.9667	0.9889	0.9778	0.9831	[11, 11]		
Random Forest	0.9222	0.9778	0.9667	0.9828	[8, 27]		
Neural Net	0.9667	0.9889	0.9778	0.9831	[17, 17]		
Decision Tree	0.9222	0.9889	0.9889	1.0000	[1, 83]		
		Brea	st				
SVM	0.9501	0.9765	0.9765	1.0000	[12, 20]		
Random Forest	0.9384	0.9619	0.9560	0.9780	[4, 19]		
Neural Net	0.9326	0.9765	0.9707	0.9903	[1, 7]		
Decision Tree	0.9296	0.9619	0.9589	0.9900	[0, 67]		
		Win	e				
SVM	0.9151	1.0000	0.9906	0.9811	[3, 13]		
Random Forest	0.9623	0.9906	0.9811	0.9818	[5, 20]		
Neural Net	0.9057	0.9906	0.9906	1.0000	[1, 25]		
Decision Tree	0.9057	0.9811	0.9811	1.0000	[5, 35]		

In average, with 20 iterations (0.42% of the search space): decrease of error by 58% compared to the baseline 1. 2. **98.92%** in the normalized score space)

How close are we from the optimal pipeline?



A solution for Euclidian space

Definition 3.2 (Mean Absolute Deviation).

MAD
$$(p^*, r) = \frac{1}{N} \sum_{i=1}^{N} |p_i^* - r|$$

r: a reference point

For each optimal configuration r:

Definition 3.3 (Normalized Mean Absolute Deviation).

$$NMAD(\mathbf{p}^*, r) = \frac{1}{K} ||MAD(\mathbf{p}^*, r)||_{2}$$

2.Express the sample in normalized conf. space 3.Calculate the NMAD on the sample

N: number of algorithms **K**: dimension of the configuration space **p***: sample of optimal configurations

1. Build the sample w.r.t. to the algorithms:

 \rightarrow For each algorithm, select the optimal point that is the closest from the reference point.

Results on two datasets for text classification

Method	(n,k)	accuracy				
ECHR			ECH	ŦR	Newsg	roup
Decision Tree	(5, 50000)	0.900				loup
Neural Network	(5, 50000)	0.960	Point	NMAD	Point	NMAD
Random Forest	(3, 10000), (4, 10000), (5, 50000)	0.910	(5, 50000)	0	(4, 5000)	0.306
Linear SVM	(3, 50000), (4, 50000), (5, 50000)	0.921	(3, 10000)	0.275	(4, 3000) (4, 100000)	0.300
	Newsgroup		(4, 10000)	0.213	(5, 50000)	0.356
Decision Tree	(4, 5000), (4, 100000)	0.889	(3, 50000)	0.175	(3, 10000)	0.294
Neural Network	(5, 50000)	0.953	(4, 50000)	0.094	(2, 100000)	0.362
Random Forest	(3, 10000)	0.931				
Linear SVM	(2, 100000)	0.946				

Future work

Work in progress:

- \rightarrow Tests on larger configuration spaces.
- \rightarrow Online architecture.
- \rightarrow Metric between pipelines



Thank you



Don't forget the poster session!

