



Hyperparameter Optimization for Software Effort Estimation

Tianpei Xia

txia4@ncsu.edu

Nov 27th, 2018

Overview

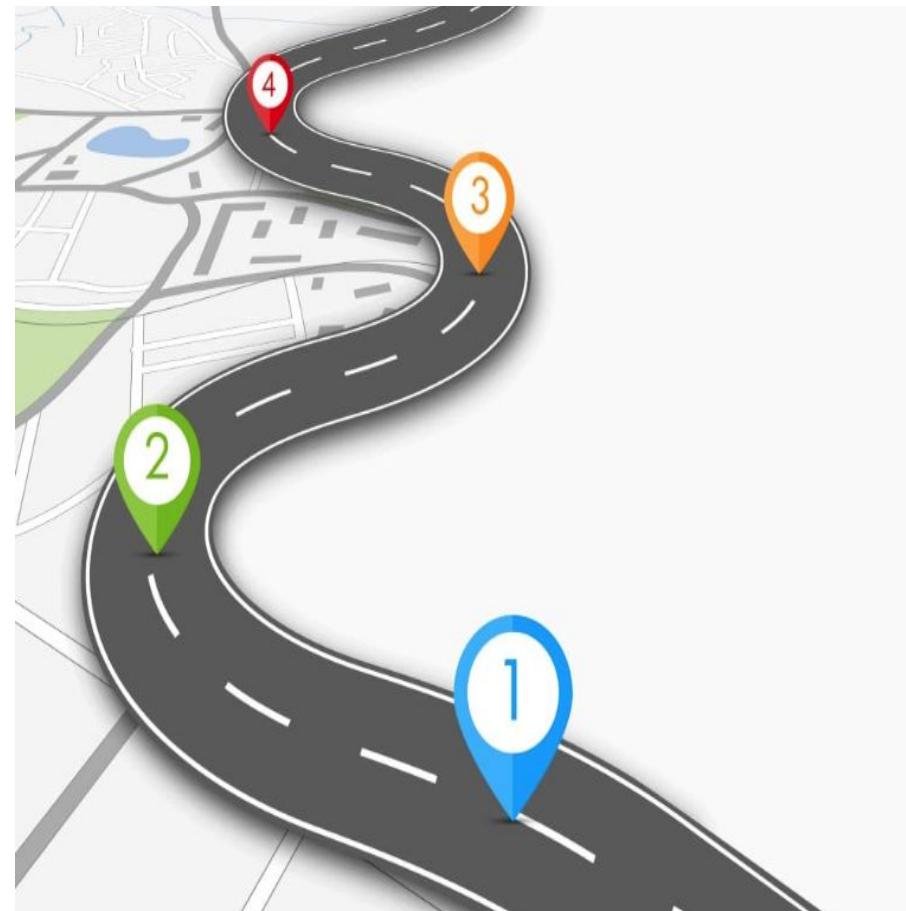
Proposal (supported by current evidence):

RATE (Rapid Automatic Tuning for
Effort-estimation)

- A novel configuration tool for effort estimation
- Based on the combination of regression tree learners and optimizers (e.g. Different Evolution [**Storn' 1997**], FLASH [**Nair' 2017**])

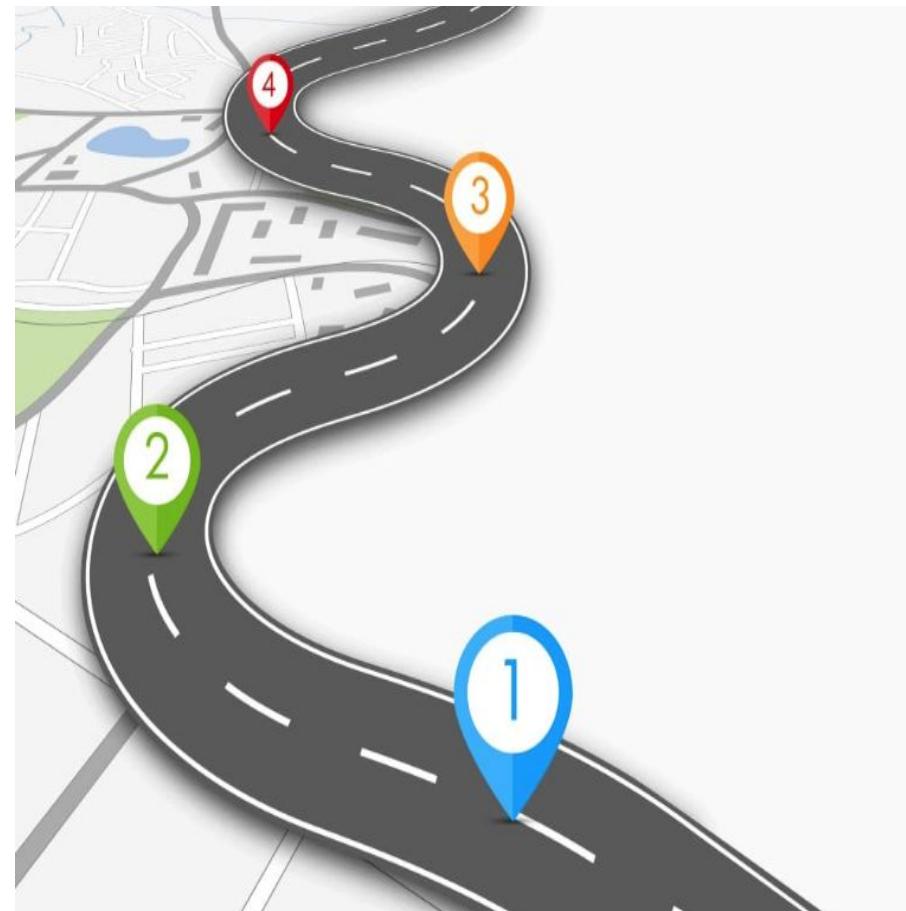
Roadmap

- **Introduction**
 - Effort Estimation
- **Background**
 - Hyperparameter Tuning
- **Empirical Study**
 - RATE
 - Dataset
 - Experiment Design
- **Results**
- **Future Work**



Roadmap

- **Introduction**
 - Effort Estimation
- **Background**
 - Hyperparameter Tuning
- **Empirical Study**
 - RATE
 - Dataset
 - Experiment Design
- **Results**
- **Future Work**



Introduction

What is effort estimation?

- Predicting human effort to plan and develop software projects.
 - Based on the information collected in previous related software projects.
 - Usually expressed in terms of hours, days or months of human work.



Introduction

Why effort estimation?

Can be used for

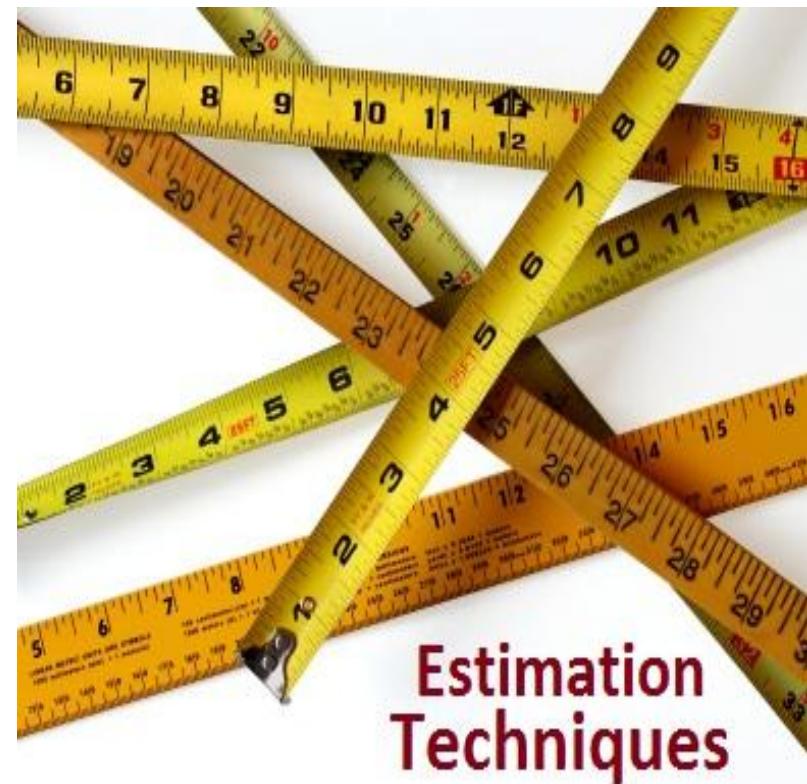
- Request for proposal
 - Contract negotiations
 - Scheduling
 - Monitoring and control



Introduction

Ways to estimate effort

- Human-based Methods
[Jørgensen' 2004]
- Algorithm-based Methods
 - COCOMO [Boehm' 1981]
 - ATLM [Whigham' 2015]
 - LP4EE [Sarro' 2018]
 - Regression Trees (e.g. CART)
 - Analogy-Based Estimation (ABE)



Introduction

Accuracy is critical

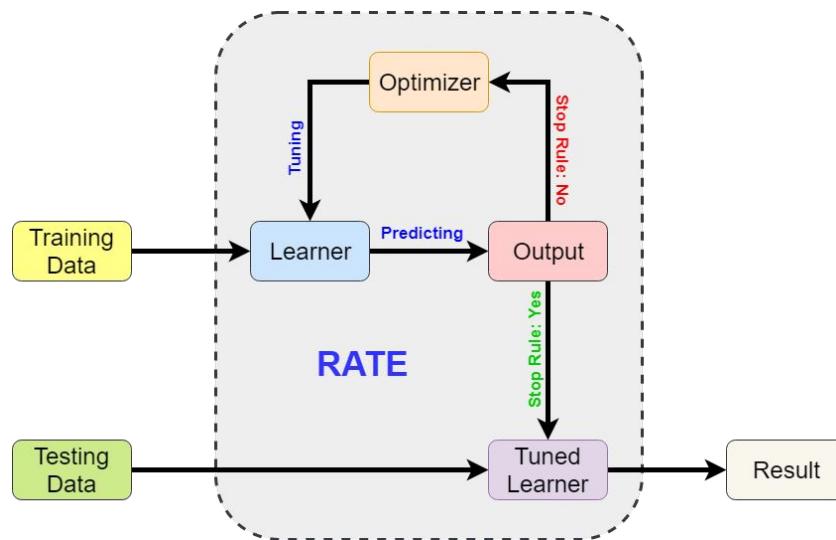
Inaccurate resource allocation can cause a considerable waste of resource and time.

- NASA canceled its Launch Control System project after the initial \$200M estimate was exceeded by another \$200M. [**Cowing' 2002**]



Introduction

Improve accuracy by hyperparameter tuning



RATE

(Rapid Automatic
Tuning for
Effort-estimation)

A hyperparameter
tuning based method.

RATE will be discussed later in the talk.

Introduction

Performance improvement

Some SE communities (e.g. defect prediction) endorse hyperparameter tuning to improve the performance:

- For open source JAVA system datasets, tunings improved the defect detection precision from 0% to 60%. [Fu' 2016]
- For defect prediction models, tuning improves 40% of AUC performance.
[Tantithamthavorn' 2016]



So why not use it on effort estimation?

Introduction

So why not effort estimation?

CPU intensive and running time?



- Arcuri et al reported their tuning may require **weeks, or more**, of CPU time. [**Arcuri' 2011**]
- Wang et al. needed **15 years of CPU** to explore 9.3 million candidate configurations for software clone detectors. [**Wang' 2013**]

Introduction

So why not effort estimation?

Are these default hyperparameter values already good enough?



Maybe
not!

Introduction

Applying hyperparameter tuning to effort estimation

Open issues:

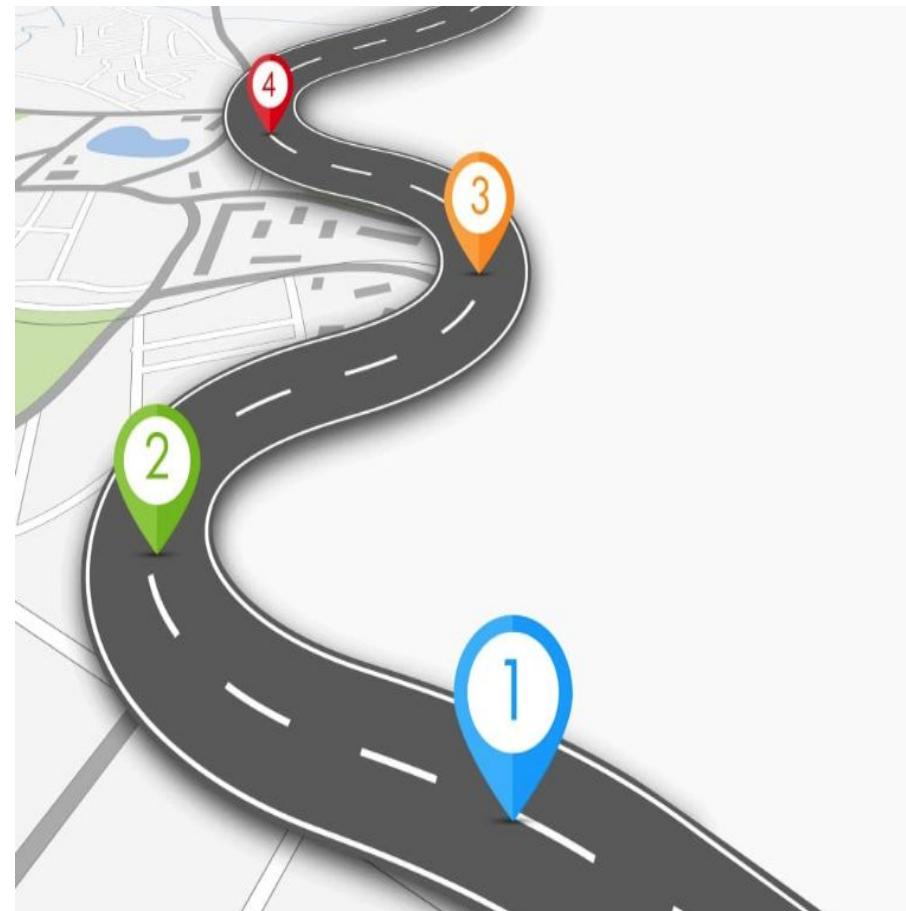
- Is it best to just use “off-the-shelf” defaults?
- Can we replace the old defaults with new defaults?
- What hyperparameter optimizers to use for effort estimation?
- Can we avoid slow hyperparameter optimization?



We'll get back to these questions later in the talk.

Roadmap

- **Introduction**
 - Effort Estimation
- **Background**
 - Hyperparameter Tuning
- **Empirical Study**
 - RATE
 - Dataset
 - Experiment Design
- **Results**
- **Future Work**

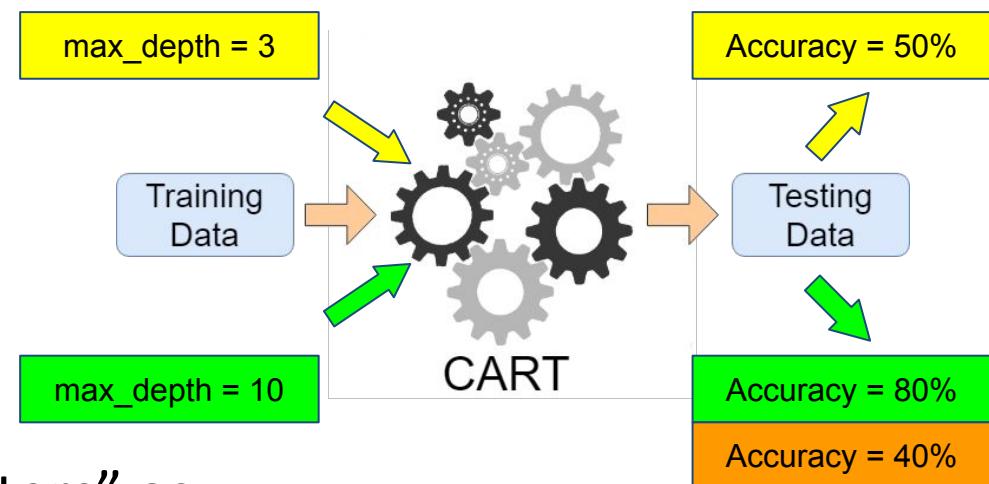


Background

What is hyperparameter?

Data mining algorithms usually have some “magic parameters” to control their behavior.

- Decision Tree (CART):
 - max_depth
 - min_samples_leaf
 -



We call these “magic parameters” as hyperparameters.

Background

Tuning for Optimization

Finding suitable hyperparameters of data miners for better results.

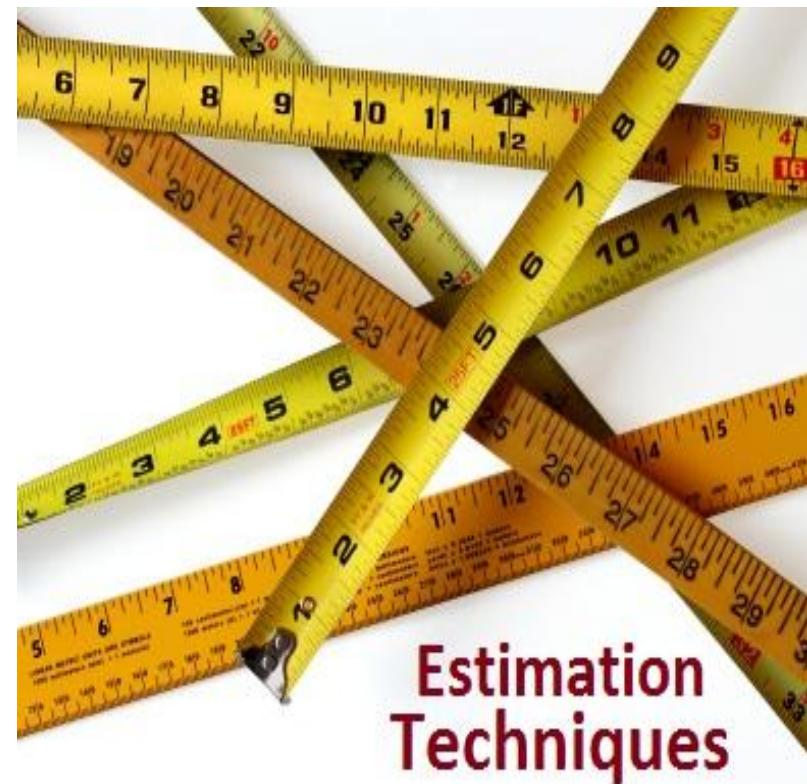
When tuning, data miners will use different heuristics and generate different models.



Background

Ways to tune

- Mathematical Optimization
- Grid Search
- Evolutionary Algorithms
 - Differential Evolution [**Storn' 1997**]
 - NSGA-II [**Deb' 2002**]
 - MOEA/D [**Zhang' 2007**]



Background

Mathematical Optimization

Objective: $\min f_0(x)$

Constraints: s.t. $f_i(x) \leq b_i, i = 1, \dots, m.$

Based on the property of objective function and constraint function:

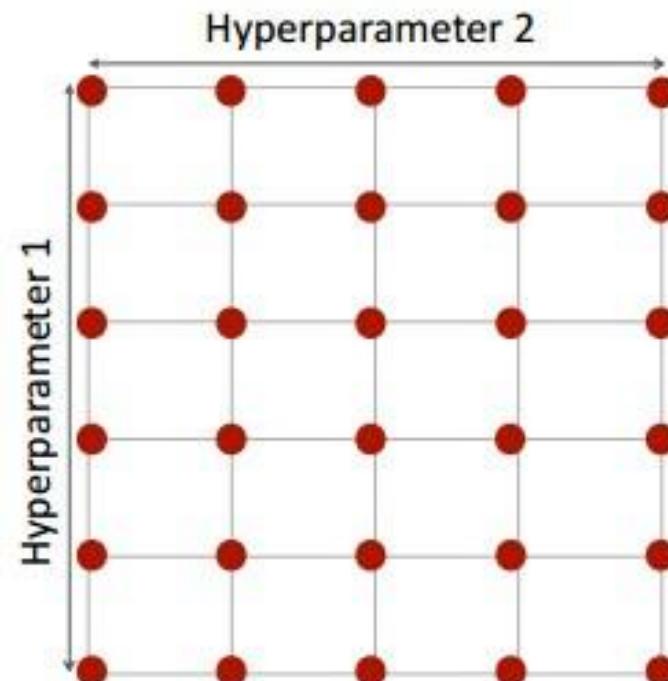
- linear programming
- non-linear programming
-

Not used much in SE since differential functions usually are not simple, and computational complexity issue. [Harman' 2001]

Background

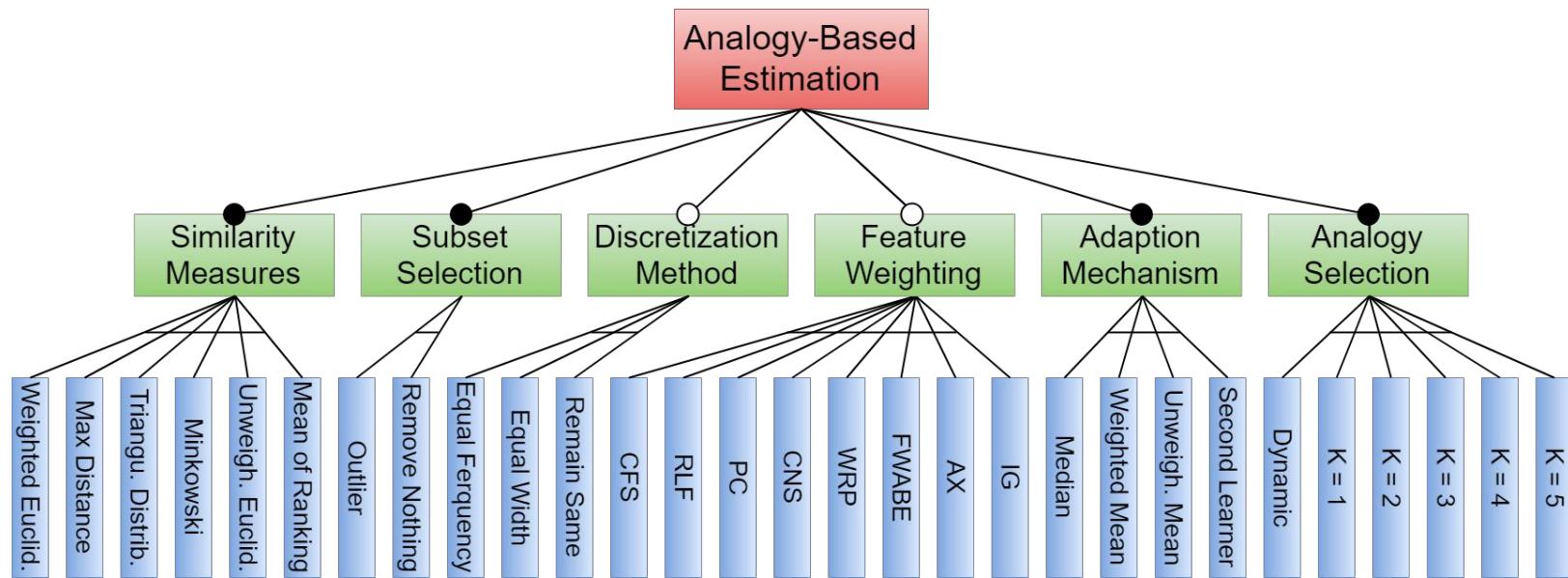
Grid Search

- Divide all C configuration options into N values.
- Evaluate N^C different combinations.
- Can be very slow!



Background

Grid Search

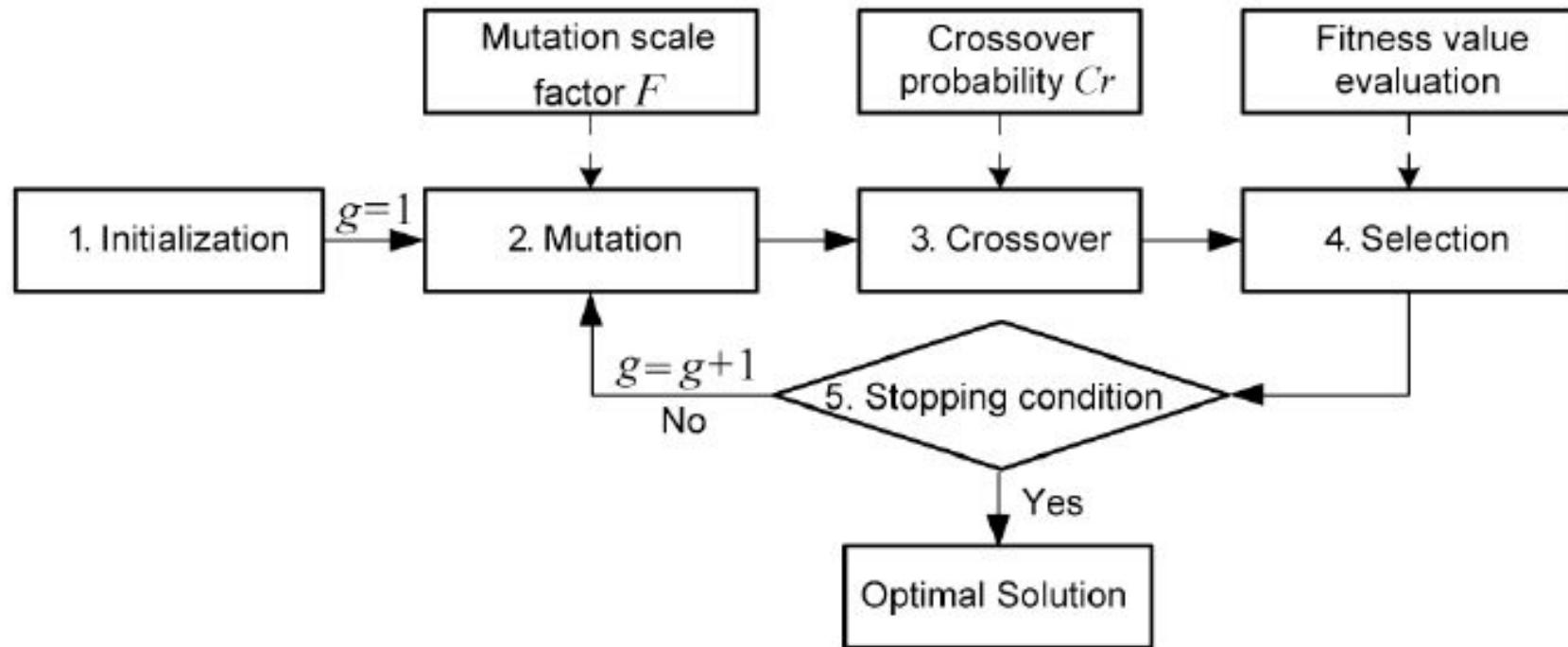


For ABEN methods, there are

$$6 * 2 * 3 * 8 * 4 * 6 = 6,912 \text{ different variants!}$$

Background

Differential Evolution (DE)



Background

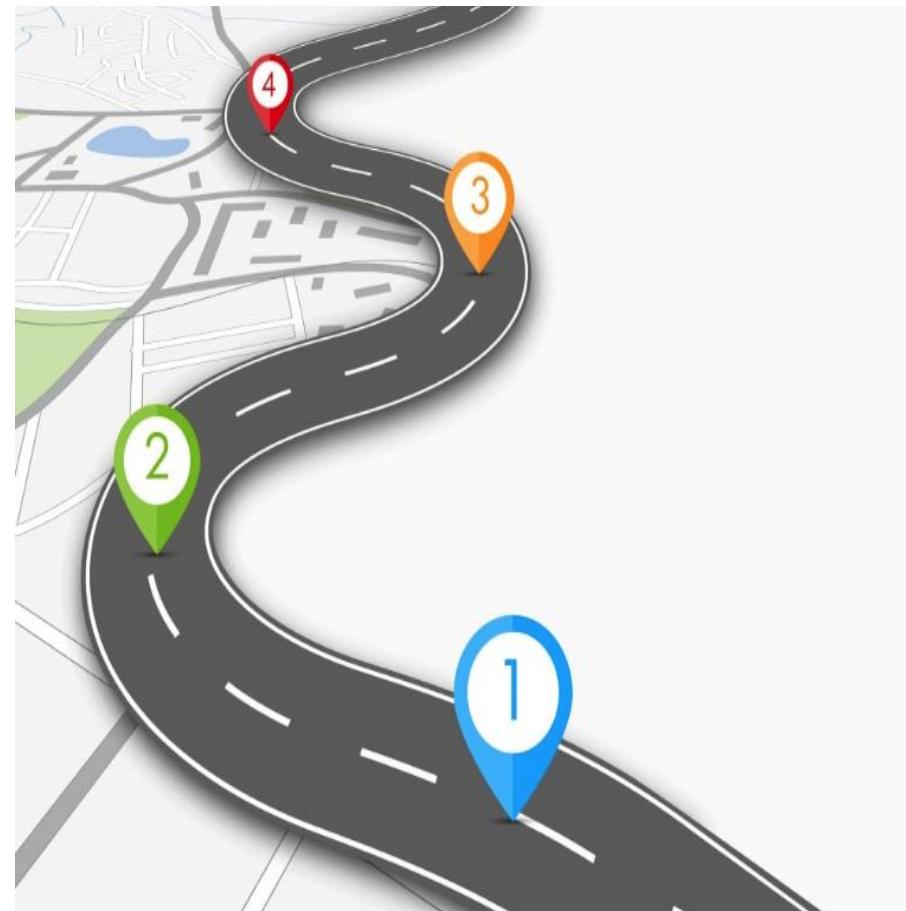
FLASH [Nair' 2017]

Pick a number of data into build_pool, evaluate the build_pool, and put the rest into rest_pool

```
while life > 0:  
    build CART model by using build_pool  
    next_point = max( model.predict( rest_pool ) )  
    build_pool += next_point  
    rest_pool -= next_point  
    if model.evaluate( next_point ) < max( build_pool ):  
        life -= 1  
return max( build_pool )
```

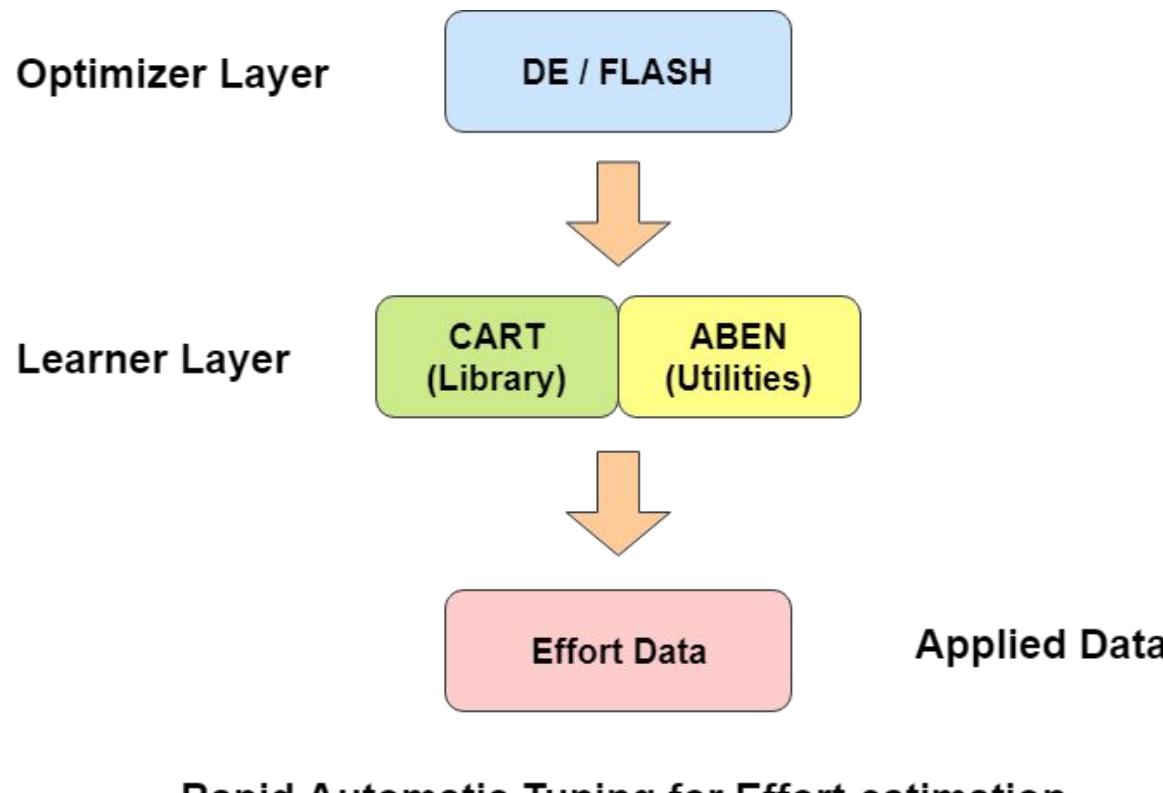
Roadmap

- **Introduction**
 - Effort Estimation
- **Background**
 - Hyperparameter Tuning
- **Empirical Study**
 - RATE
 - Dataset
 - Experiment Design
- **Results**
- **Future Work**



Empirical Study

RATE Architecture

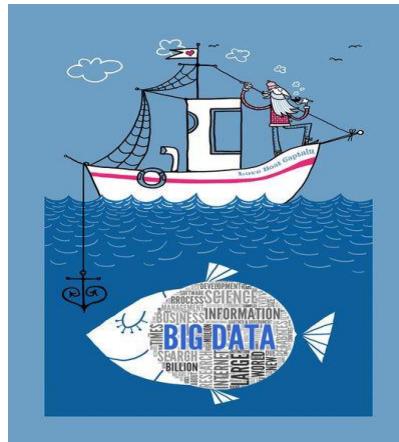


Empirical Study

Software effort dataset

SEACRAFT

A research repository for SE datasets

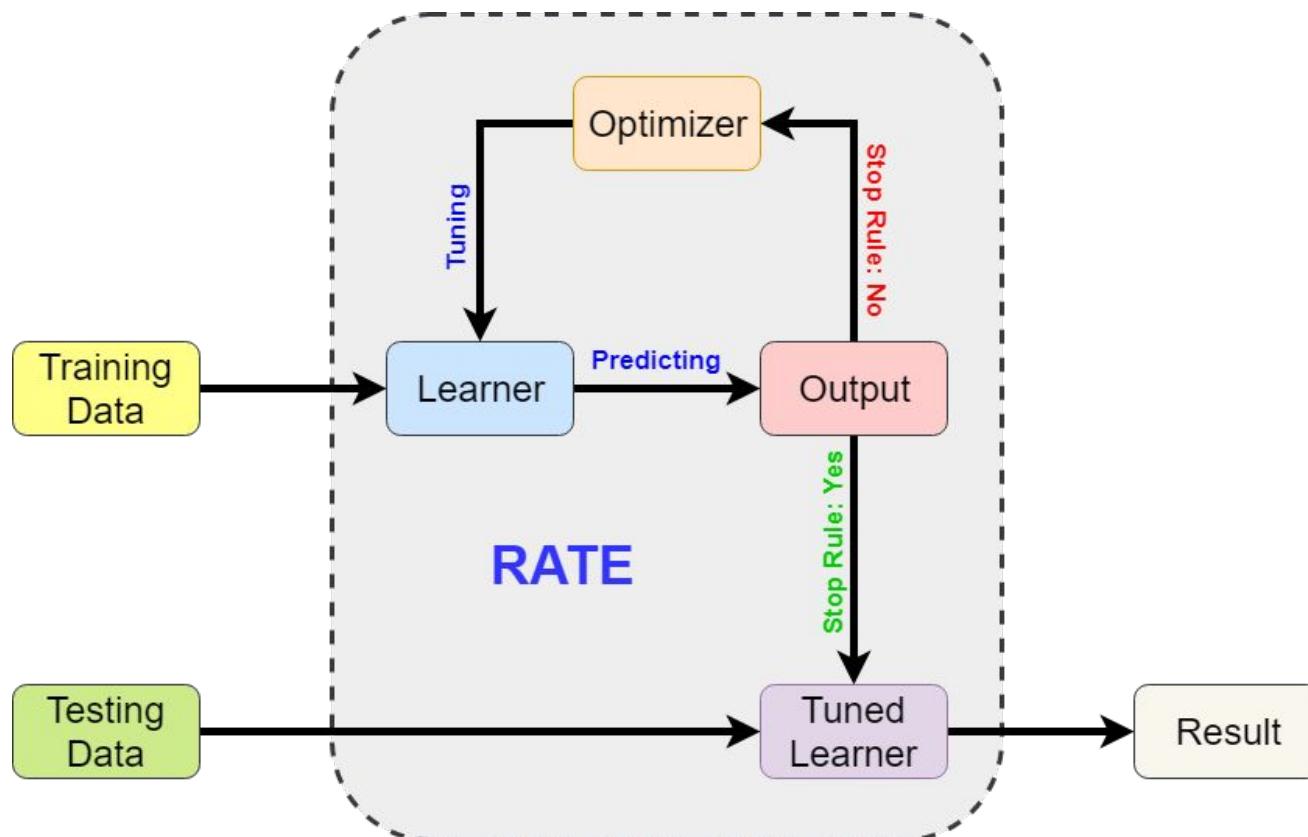


	feature	min	max	mean	std
desharnais	TeamExp	0	4	2.3	1.3
	MngExp	0	7	2.6	1.5
	Length	1	36	11.3	6.8
	Trans.s	9	886	177.5	146.1
	Entities	7	387	120.5	86.1
	AdjPts	73	1127	298.0	182.3
	Effort	546	23940	4834	4188

Datasets	#Projects	#Features
kemerer	15	6
albrecht	24	7
isbsg10	37	11
finnish	38	7
miyazaki	48	7
maxwell	62	25
desharnais	77	6
kitchenham	145	6
china	499	16
Total#	945	

Empirical Study

Experiment design



Empirical Study

Performance metrics

Magnitude of Relative Error (MRE) [Conte' 1986]

$$\text{MRE} = \left(\frac{|Actual - Predicted|}{Actual} \right) \times 100$$

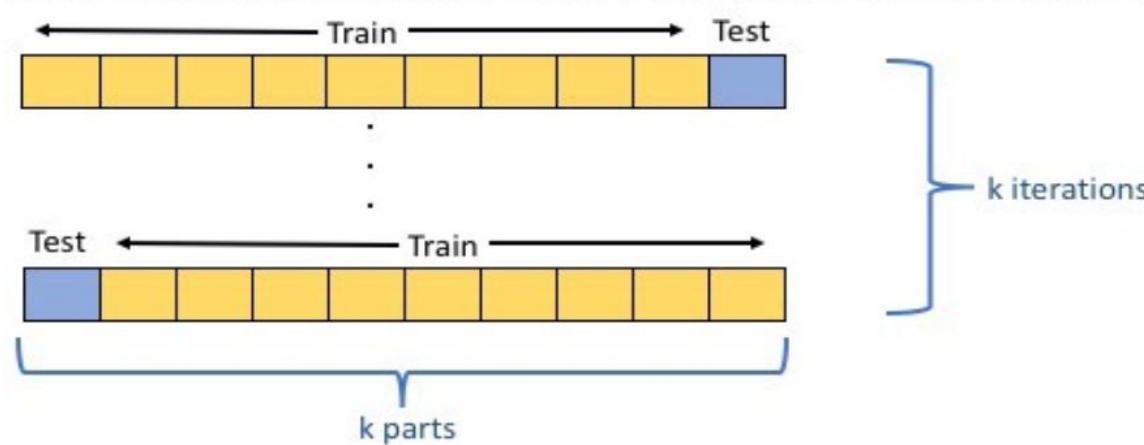
Standardized Accuracy (SA) [Shepperd' 2012]

$$\text{SA} = \left(1 - \frac{|Actual - Predicted|}{\text{mean}(Actual)} \right) \times 100$$

Note: MRE: Lower the better; SA: Higher the better.

Empirical Study

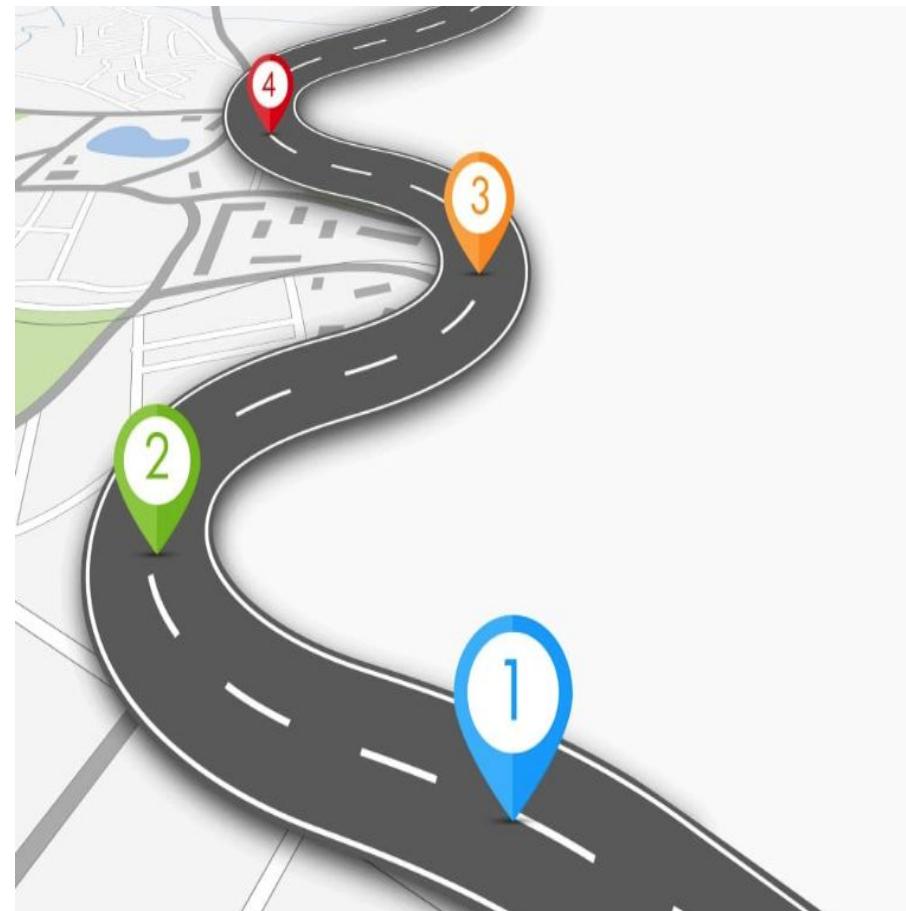
K-Folds Cross Validation



1. Divide the sample data into k parts.
2. Use $k-1$ of the parts for training, and 1 for testing.
3. Repeat the procedure k times, rotating the test set.
4. Determine an expected performance metric based on the results across the iterations.

Roadmap

- **Introduction**
 - Effort Estimation
- **Background**
 - Hyperparameter Tuning
- **Empirical Study**
 - RATE
 - Dataset
 - Experiment Design
- **Results**
- **Future Work**



Results

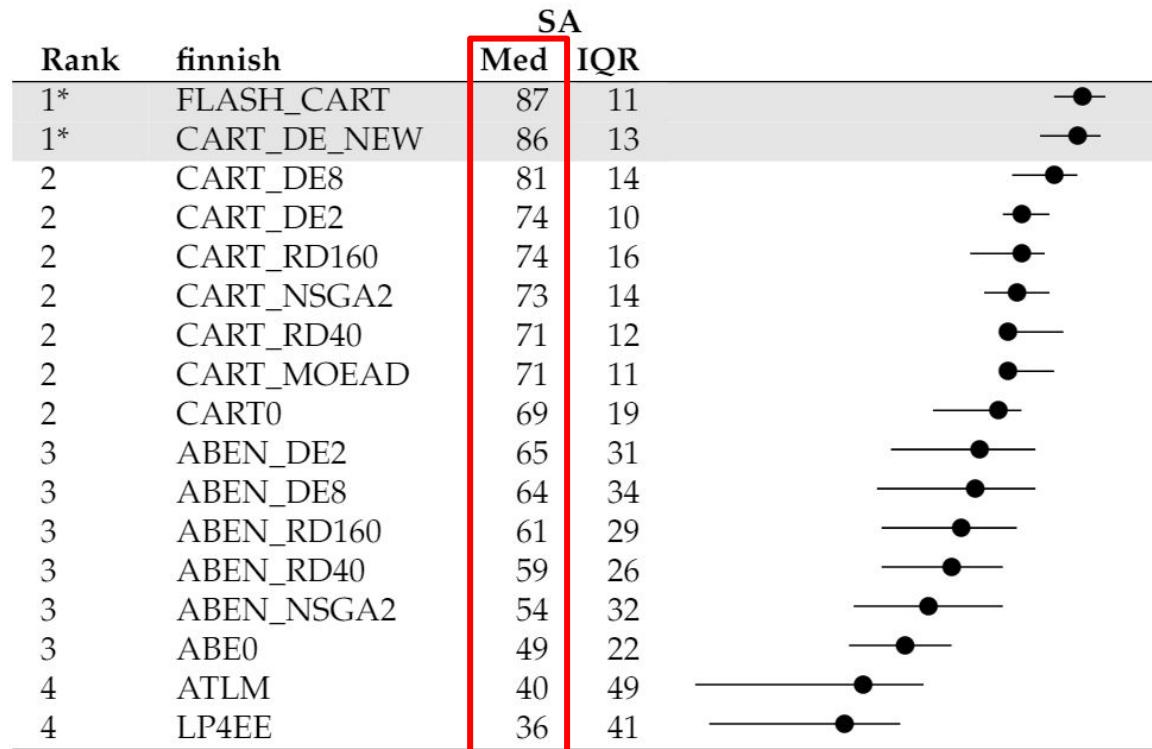
Performance ranking (SA scores)

Rank	finnish	SA		
		Med	IQR	
1*	FLASH_CART	87	11	
1*	CART_DE_NEW	86	13	
2	CART_DE8	81	14	
2	CART_DE2	74	10	
2	CART_RD160	74	16	
2	CART_NSGA2	73	14	
2	CART_RD40	71	12	
2	CART_MOEAD	71	11	
2	CART0	69	19	
3	ABEN_DE2	65	31	
3	ABEN_DE8	64	34	
3	ABEN_RD160	61	29	
3	ABEN_RD40	59	26	
3	ABEN_NSGA2	54	32	
3	ABE0	49	22	
4	ATLM	40	49	
4	LP4EE	36	41	

Example of SA scores for finnish dataset, using 10-way cross validation, 20 times experiment repeats

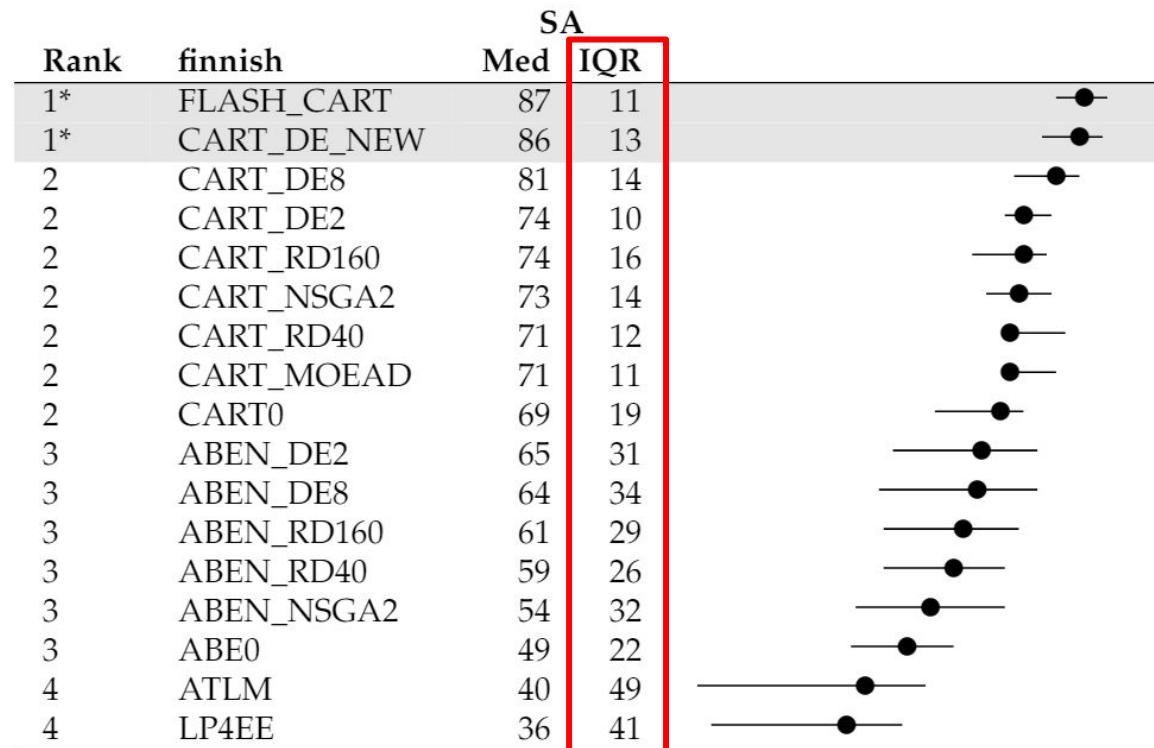
Results

Performance ranking



Results

Performance ranking



Results

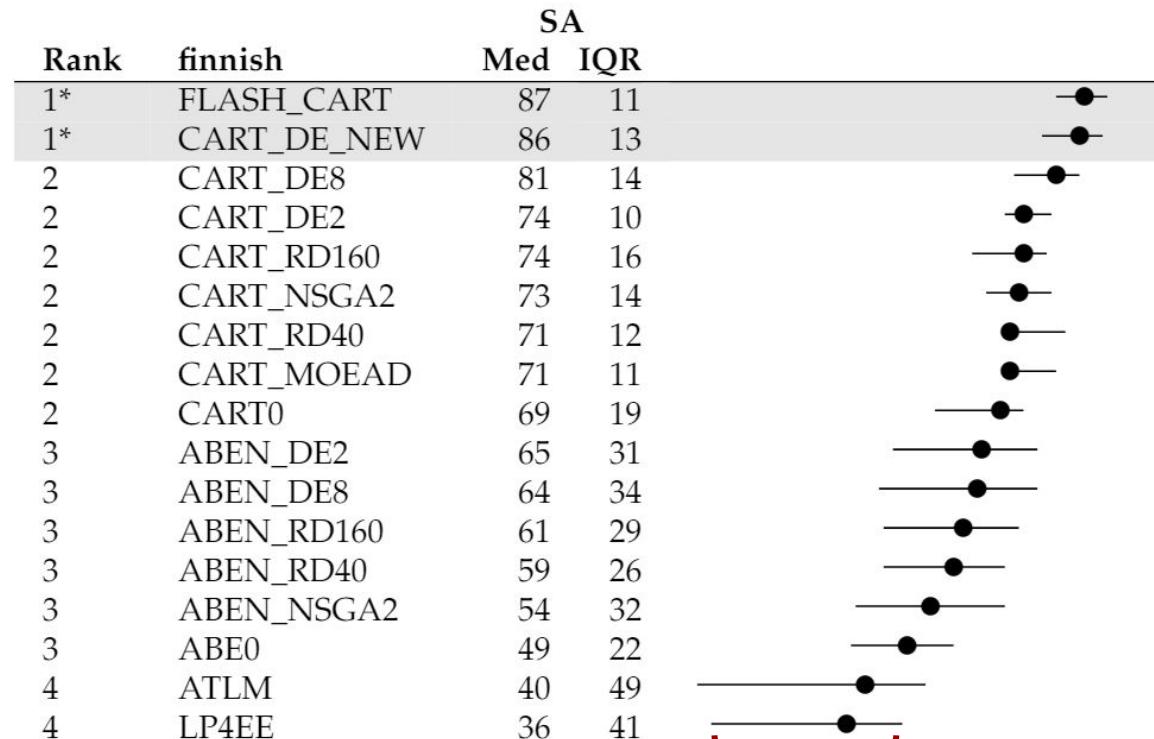
Performance ranking

Rank	finnish	SA		
		Med	IQR	
1*	FLASH_CART	87	11	
1*	CART_DE_NEW	86	13	
2	CART_DE8	81	14	
2	CART_DE2	74	10	
2	CART_RD160	74	16	
2	CART_NSGA2	73	14	
2	CART_RD40	71	12	
2	CART_MOEAD	71	11	
2	CART0	69	19	
3	ABEN_DE2	65	31	
3	ABEN_DE8	64	34	
3	ABEN_RD160	61	29	
3	ABEN_RD40	59	26	
3	ABEN_NSGA2	54	32	
3	ABE0	49	22	
4	ATLM	40	49	
4	LP4EE	36	41	



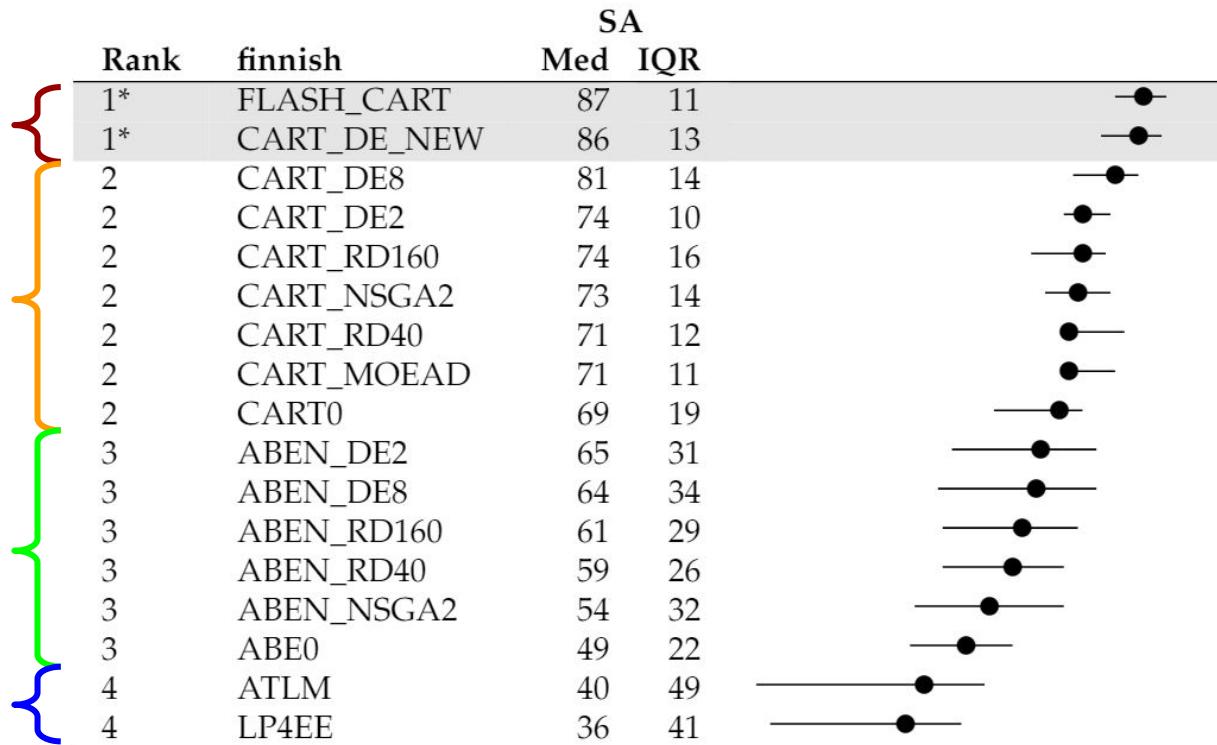
Results

Performance ranking



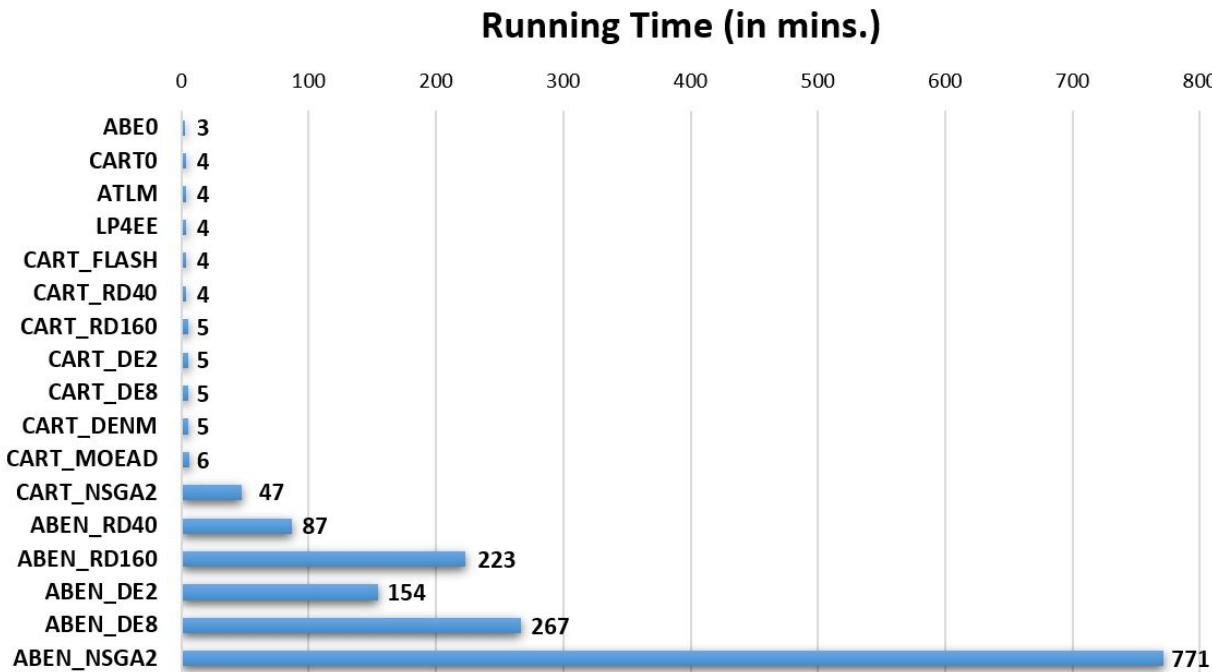
Results

Performance ranking



Results

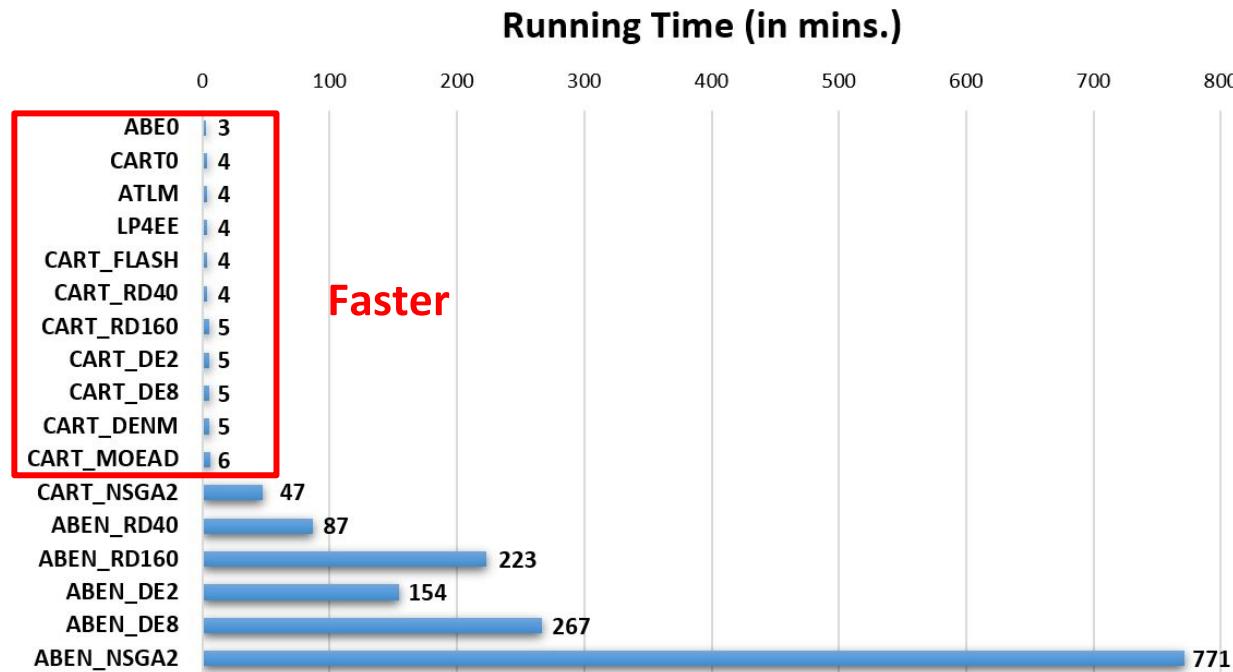
Running time



Average runtime (in minutes), for one-way out of an $N \times M$ cross-validation experiment, in a local machine.

Results

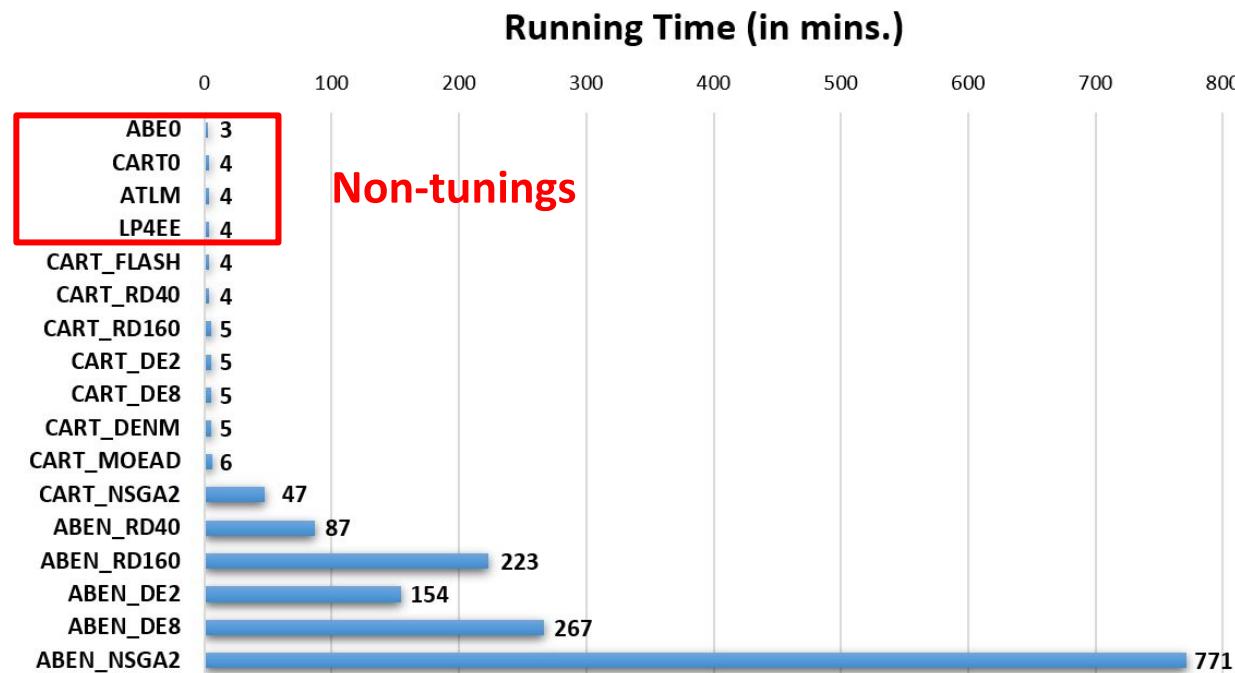
Running time



Average runtime (in minutes), for one-way out of an $N \times M$ cross-validation experiment, in a local machine.

Results

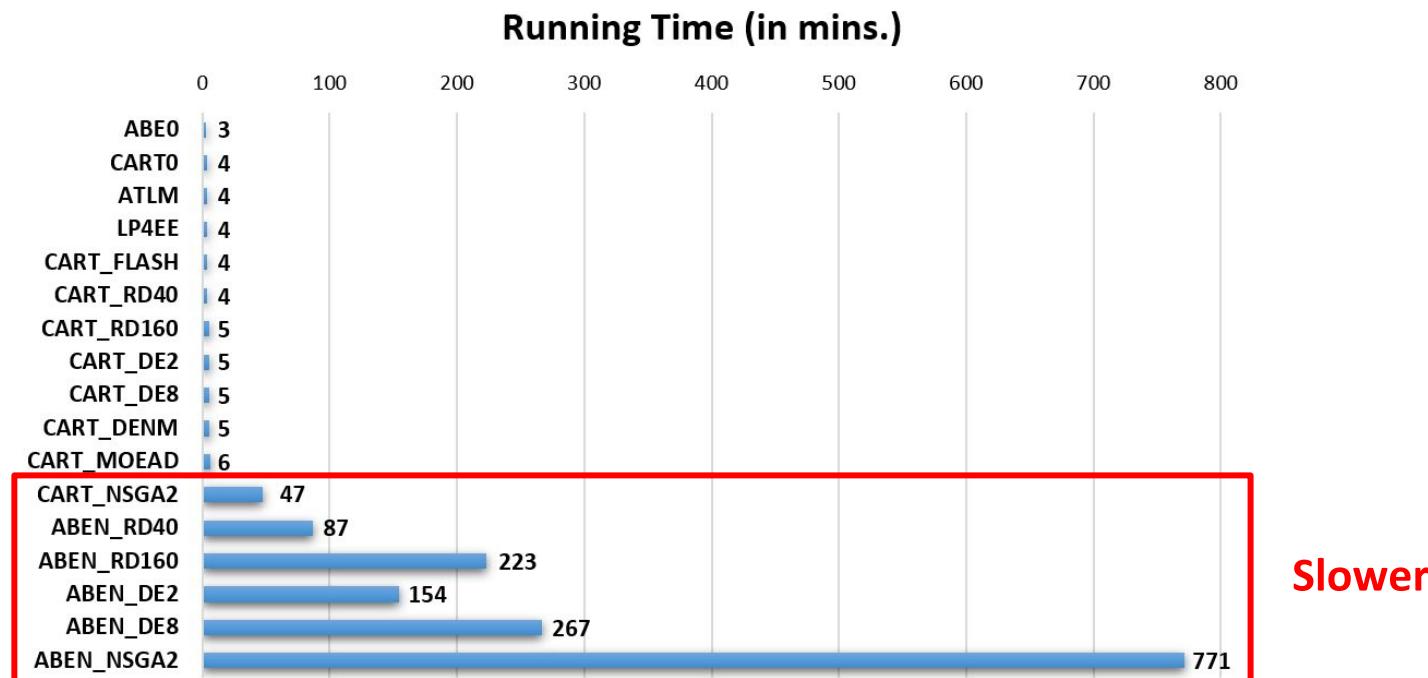
Running time



Average runtime (in minutes), for one-way out of an $N \times M$ cross-validation experiment, in a local machine.

Results

Running time



Average runtime (in minutes), for one-way out of an $N \times M$ cross-validation experiment, in a local machine.

Results

Performance ranking

Rank	albrecht	MRE		
		Med	IQR	
1*	FLASH_CART	46	21	
	ABEN_RD160	48	29	
	ABEN_DE2	48	20	
1*	ABE0	51	28	
	ABEN_DE8	51	30	
1*	CART_DENM	51	33	
1*	CART_DE8	52	28	
2	CART_DE2	56	23	
	CART_MOEAD	56	23	
	CART0	59	21	
4	LP4EE	143	87	<i>out-of-range</i>
4	ATLM	178	122	<i>out-of-range</i>

Results

Performance ranking

Rank	albrecht	MRE			
		Med	IQR		
1*	FLASH_CART	46	21		
→ 1	ABEN_RD160	48	29		
→ 1	ABEN_DE2	48	20		
1*	ABE0	51	28		
→ 1	ABEN_DE8	51	30		
1*	CART_DENM	51	33		
1*	CART_DE8	52	28		
2	CART_DE2	56	23		
2	CART_MOEAD	56	23		
2	CART0	59	21		
4	LP4EE	143	87	<i>out-of-range</i>	
4	ATLM	178	122	<i>out-of-range</i>	

Research Questions

- Is it best to just use “off-the-shelf” defaults?
- Can we replace the old defaults with new defaults?
- Can we avoid slow hyperparameter optimization?
- What hyperparameter optimizers to use for effort estimation?



Research Questions

- Is it best to just use “off-the-shelf” defaults?
- Can we replace the old defaults with new defaults?
- Can we avoid slow hyperparameter optimization?
- What hyperparameter optimizers to use for effort estimation?



Is it best to just use “off-the-shelf” defaults?

Rank	Using	Med.	IQR	
1*	CART_MOEAD	65	21	
1	ABEN_DE2	63	11	
1	ABEN_DE8	62	18	
1	ABEN_RD160	60	18	
1*	ABE0	60	16	
1*	FLASH_CART	59	39	
2	CART_DENM	54	32	
2	CART_DE8	53	19	
2	CART_DE2	50	52	
2	CART0	45	40	
3	ATLM	18	30	
3	LP4EE	17	34	
china				
1*	CART_DE8	93	2	
1*	CART_DE2	93	2	
2	FLASH_CART	89	5	
2	CART_DENM	87	9	
3	CART0	85	7	
4	ABE0	61	11	
5	CART_MOEAD	54	19	
5	LP4EE	51	13	
5	ATLM	48	7	
desharnais				
1*	FLASH_CART	47	16	
1*	CART_DENM	44	24	
1	ABEN_DE8	44	24	
1	ABEN_NSGA2	43	27	
1*	ATLM	42	13	
1	ABEN_RD40	42	27	
1*	CART_DE8	42	27	
1	ABEN_DE2	41	28	
1*	CART_DE2	40	26	
1*	LP4EE	40	24	
1	ABEN_RD160	40	31	
2	CART_MOEAD	35	31	
2	ABE0	30	42	
2	CART0	23	18	
finnish				
1*	FLASH_CART	87	11	
1*	CART_DENM	86	13	
2	CART_DE8	81	14	
2	CART_DE2	74	10	
2	CART_MOEAD	71	11	
2	CART0	69	19	
3	ABE0	49	22	
4	ATLM	40	49	
4	LP4EE	36	41	

Overall, “off-the-shelf” defaults rarely have good performance.

Just use “Off-the-shelf” defaults should be deprecated.

Research Questions

- Is it best to just use “off-the-shelf” defaults?

- Can we replace the old defaults with new defaults?
- Can we avoid slow hyperparameter optimization?
- What hyperparameter optimizers to use for effort estimation?



Research Questions

- Is it best to just use “off-the-shelf” defaults?

- Can we replace the old defaults with new defaults?
- Can we avoid slow hyperparameter optimization?
- What hyperparameter optimizers to use for effort estimation?



Can we replace the old defaults with new defaults?

	max_features (selected at random; 100% means "use all")				max_depth (of trees)				min_sample_split (continuation criteria)				min_samples_leaf (termination criteria)			
	25%	50%	75%	100%	≤03	≤06	≤09	≤12	≤5	≤10	≤15	≤20	≤03	≤06	≤09	≤12
kemerer	18	32	23	27	57	37	05	00	95	02	03	00	92	02	05	02
albrecht	13	23	20	43	63	28	08	00	68	32	00	00	83	15	02	00
isbsg10	12	35	28	25	57	33	08	00	47	23	15	15	60	27	10	03
finnish	07	03	27	63	32	56	12	00	73	18	05	03	78	17	05	00
miyazaki	10	22	27	40	31	46	20	03	42	24	18	16	78	13	07	02
maxwell	04	16	40	40	18	60	20	02	44	27	17	12	50	33	14	04
desharnais	25	23	27	25	40	46	11	02	36	26	13	25	32	26	24	19
kitchenham	01	12	32	56	03	42	45	10	43	30	17	10	48	35	12	04
china	00	04	25	71	00	00	25	75	56	30	10	02	68	28	04	00

KEY: 10 20 30 40 50 60 70 80 90 100 %

Tunings discovered by hyperparameter selections for CART_DE8

Overall, there are no “best” default settings.

Research Questions

- Is it best to just use “off-the-shelf” defaults?
NO
- Can we replace the old defaults with new defaults?
NO
- Can we avoid slow hyperparameter optimization?
- What hyperparameter optimizers to use for effort estimation?

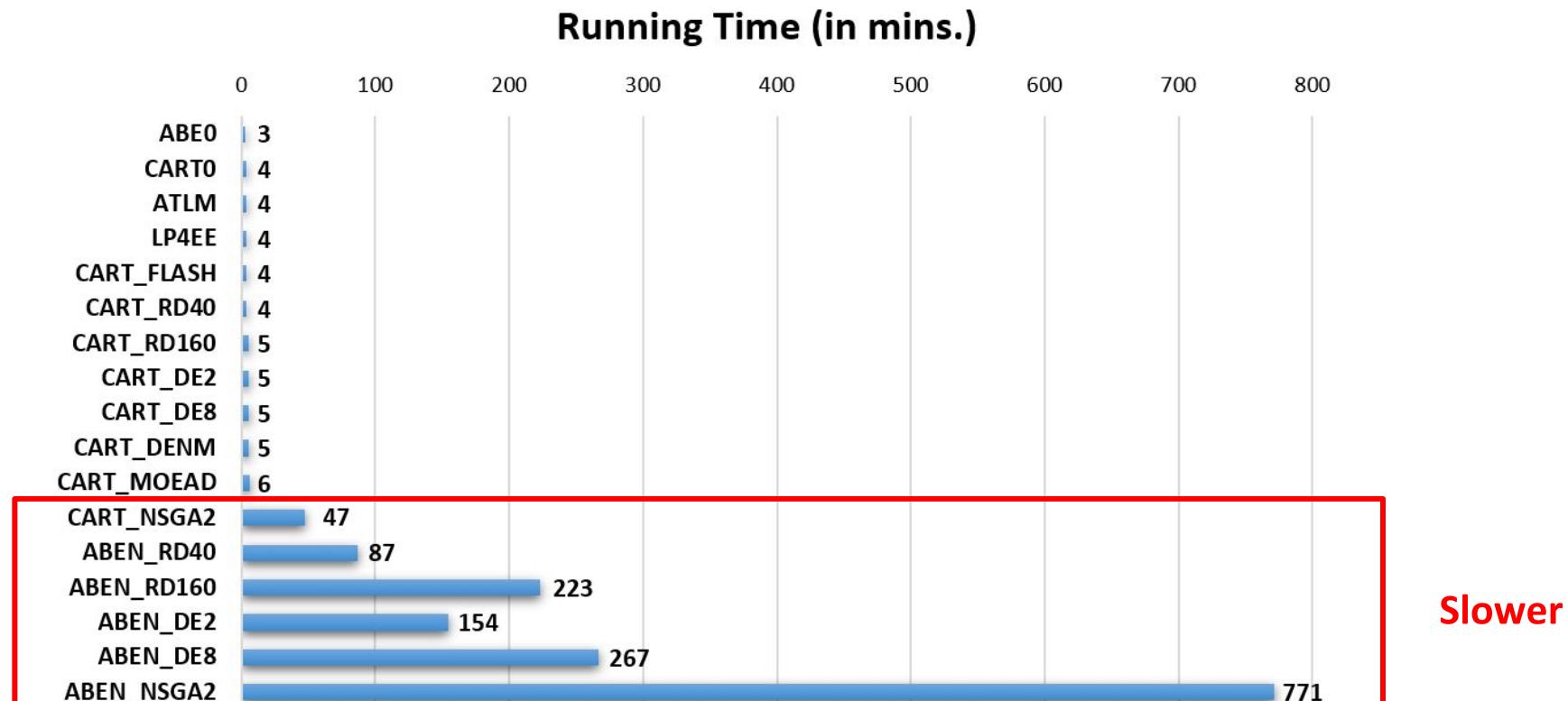


Research Questions

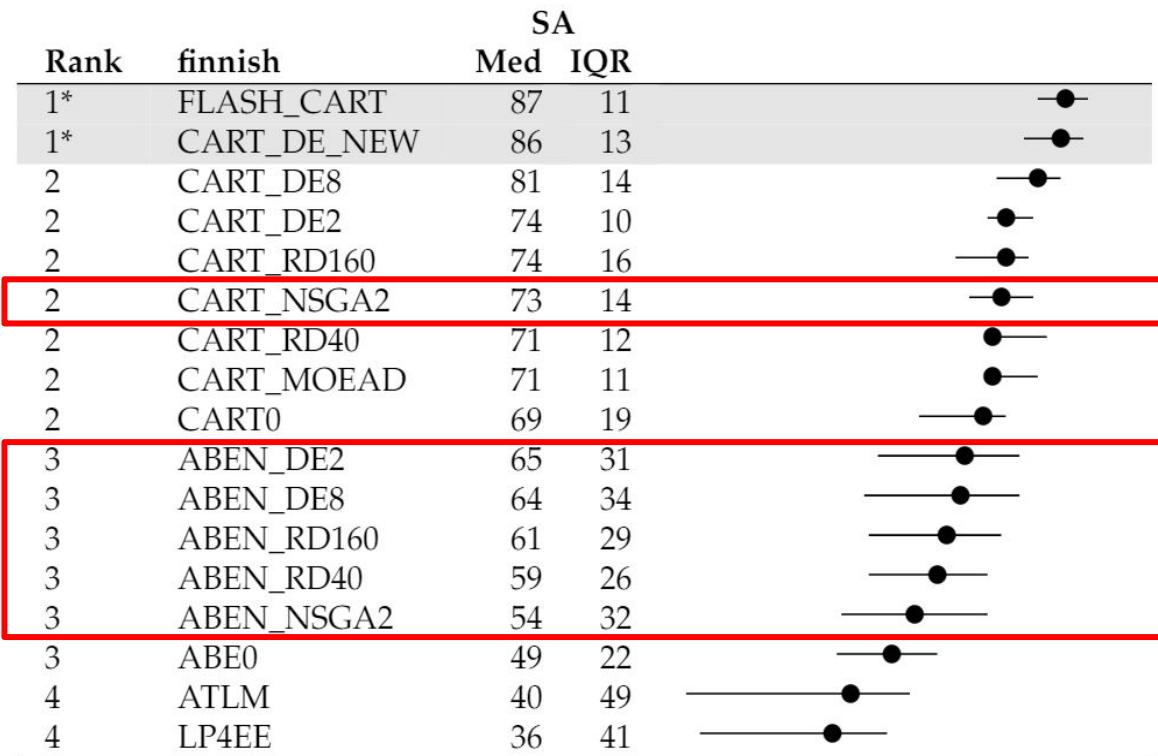
- Is it best to just use “off-the-shelf” defaults?
NO
- Can we replace the old defaults with new defaults?
NO
- Can we avoid slow hyperparameter optimization?
- What hyperparameter optimizers to use for effort estimation?



Can we avoid slow hyperparameter optimization?



Can we avoid slow hyperparameter optimization?



Overall, our slowest optimizers perform no better than certain faster ones.

In all cases (9 datasets * 2 metrics), slower ones only win (Rank as “1*”) in 1/18.

Research Questions

- Is it best to just use “off-the-shelf” defaults?
NO
- Can we replace the old defaults with new defaults?
NO
- Can we avoid slow hyperparameter optimization?
YES
- What hyperparameter optimizers to use for effort estimation?



Research Questions

- Is it best to just use “off-the-shelf” defaults?
NO
- Can we replace the old defaults with new defaults?
NO
- Can we avoid slow hyperparameter optimization?
YES
- What hyperparameter optimizers to use for effort estimation?



What hyperparameter optimizers to use for effort estimation?

In all 18 experiment cases (9 datasets, 2 metrics for each),

If we use only one method, after counting the frequencies in the top-ranks (Rank as “1*”), then we get:

9/18: CART_DE2, CART_DENM

10/18: CART_DE8

12/18: CART_FLASH

What hyperparameter optimizers to use for effort estimation?

In all 18 experiment cases (9 datasets, 2 metrics for each),

If we use two methods combinations, then we get:

16/18: CART_DE2 + CART_FLASH

16/18: CART_DE8 + CART_FLASH

What hyperparameter optimizers to use for effort estimation?

In all 18 experiment cases (9 datasets, 2 metrics for each),

And if we use three methods combinations, then:

17/18: ABEO + CART_DE2 + CART_FLASH

17/18: ABEO + CART_DE8 + CART_FLASH

17/18: ATLM + CART_DE2 + CART_FLASH

17/18: ATLM + CART_DE8 + CART_FLASH

17/18: CART_MOEAD + CART_FLASH + CART_DE2

17/18: CART_NSGA2 + CART_FLASH + CART_DE2

17/18: CART_MOEAD + CART_FLASH + CART_DE8

17/18: CART_NSGA2 + CART_FLASH + CART_DE8

What hyperparameter optimizers to use for effort estimation?

- To simplify the implementation of our methods, we will avoid CART_DE8, MOEA/D and NSGA-II.
- Since 2-methods combo have similar performance as 3-methods combo (16/18 vs 17/18) with faster running time, the recommended choice will be:

CART_DE2 + CART_FLASH

Hence, for effort datasets, try a combination of CART with the optimizers Differential Evolution and FLASH.

Research Questions

- Is it best to just use “off-the-shelf” defaults?
NO
- Can we replace the old defaults with new defaults?
NO
- Can we avoid slow hyperparameter optimization?
YES
- What hyperparameter optimizers to use for effort estimation?
YES

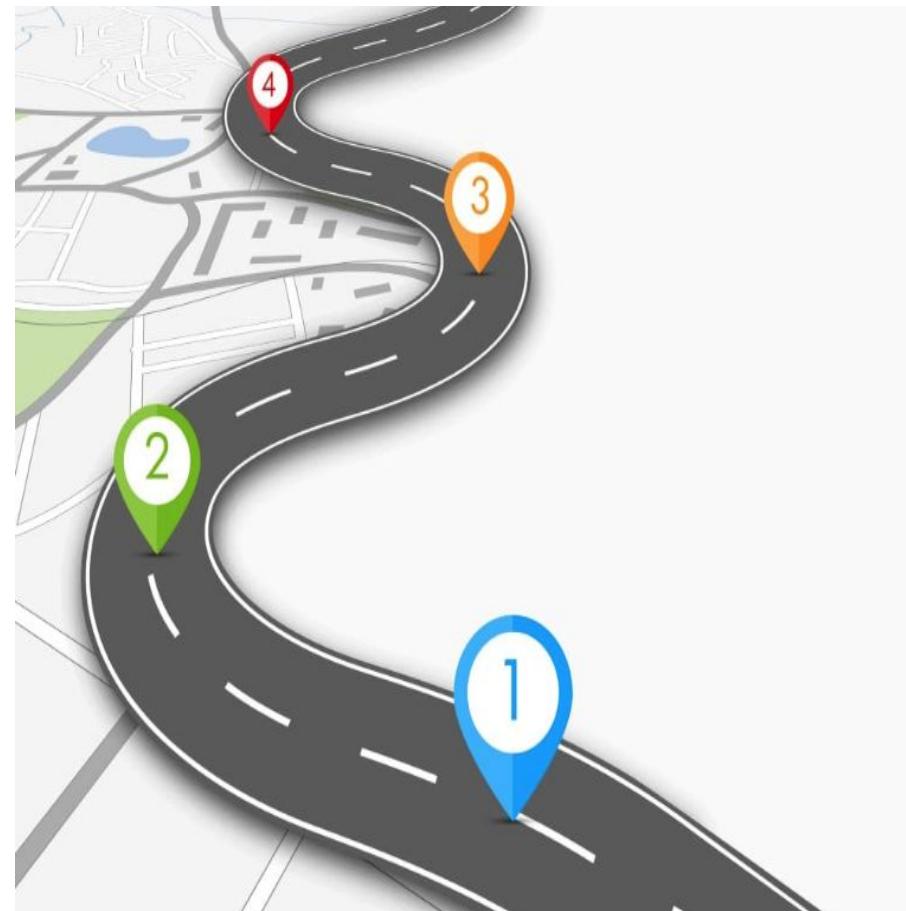


Contributions

- A demonstration that default settings are not the best way to perform effort estimation.
- A recognition of the inherent difficulty associated with effort estimation.
- A new criteria for assessing effort estimators.
- The identification of a combination of learner and optimizer that works good.
- An extensible open-source architecture.

Roadmap

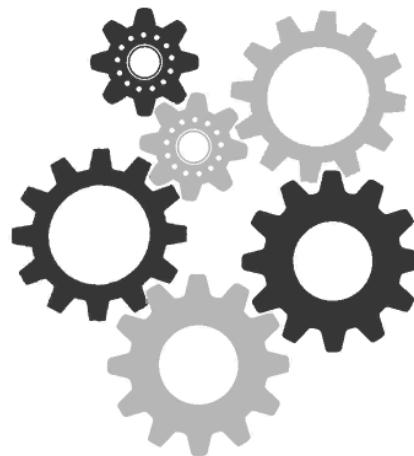
- **Introduction**
 - Effort Estimation
- **Background**
 - Hyperparameter Tuning
- **Empirical Study**
 - RATE
 - Dataset
 - Experiment Design
- **Results**
- **Future Work**



Future Work

More Exploration on RATE

- Better Optimizers
- Better Learners



Future Work

New Data Collections

- Github
- Stack Overflow
-



Future Work

Other Domains

Deep Learning

- State-of-the-art but computing intensive
- What to tune?
 - number of hidden layer
 - learning rate
 - amount of regularization
 -



Future Work

Other Domains

In a 2018 TSE paper [**Huang' 2018**]:



Automating Intention Mining

ISSN: 0098-5589

TABLE 11
The accuracy achieved by different approaches using the email and issue report dataset.

Approaches	Email	Email to Issue	Issue to Email
Ours	0.791	0.579	0.658
CNN	0.710	0.522	0.569
NLP	0.575	0.460	0.483
SMO	0.707	0.370	0.450
LibSVM	0.273	0.231	0.173
NBM	0.607	0.349	0.347
RF	0.693	0.329	0.383
kNN	0.604	0.239	0.370
DE-SVM	0.714	0.374	0.402
Auto-Weka	0.707	0.349	0.329

TABLE 12
Training time cost of each approach under different experiment setting

Approach	Issue	Email	Issue to Email	Email to Issue
Our (with BN)	30min	6min	38min	7min
Our (without BN)	5.6h	1h	7.2h	1.2h
Kim's CNN	20min	4min	24min	5min
NLP	4min	1min	6min	1min
SMO	11s	2s	12s	3s
LibSVM	9s	1s	10s	2s
NBM	1s	1s	1s	1s
RF	3min	1min	4min	1min
kNN	1s	1s	1s	1s
DE-SVM	20min	10min	16min	11min
Auto-Weka	16min	15min	15min	15min

Future Work

Beyond Hyperparameter

- The optimizers like DE or FLASH in RATE also have their own magic parameters
- Hyper-hyperparameter optimization could be a challenge.



Thank You!



Reference

[Storn' 1997] Storn, R. and Price, K., 1997. Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces. *Journal of global optimization*, 11(4), pp.341-359.

[Nair' 2017] Nair, V., Yu, Z. and Menzies, T., 2017. Flash: A faster optimizer for sbse tasks. *arXiv preprint arXiv:1705.05018*.

[Jørgensen' 2004] Jørgensen, M., 2004. A review of studies on expert estimation of software development effort. *Journal of Systems and Software*, 70(1-2), pp.37-60.

[Boehm' 1981] Boehm, B.W., 1981. *Software engineering economics* (Vol. 197). Englewood Cliffs (NJ): Prentice-hall.

[Whigham' 2015] Whigham, P.A., Owen, C.A. and Macdonell, S.G., 2015. A baseline model for software effort estimation. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 24(3), p.20.

[Sarro' 2018] Sarro, F. and Petrozziello, A., 2018. Linear Programming as a Baseline for Software Effort Estimation. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 27(3), p.12.

[Cowing' 2002] Cowing, K., 2002. NASA to Shut Down Checkout & Launch Control System. SpaceRef.

[Fu' 2016] Fu, W., Menzies, T. and Shen, X., 2016. Tuning for software analytics: Is it really necessary?. *Information and Software Technology*, 76, pp.135-146.

[Tantithamthavorn' 2016] Tantithamthavorn, C., McIntosh, S., Hassan, A.E. and Matsumoto, K., 2016, May. Automated parameter optimization of classification techniques for defect prediction models. In *Software Engineering (ICSE), 2016 IEEE/ACM 38th International Conference on* (pp. 321-332). IEEE.

Reference

[Arcuri' 2011] Arcuri, A. and Fraser, G., 2011, September. On parameter tuning in search based software engineering. In *International Symposium on Search Based Software Engineering* (pp. 33-47). Springer, Berlin, Heidelberg.

[Wang' 2013] Wang, T., Harman, M., Jia, Y. and Krinke, J., 2013, August. Searching for better configurations: a rigorous approach to clone evaluation. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering* (pp. 455-465). ACM.

[Dev' 2002] Deb, K., Pratap, A., Agarwal, S. and Meyarivan, T.A.M.T., 2002. A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE transactions on evolutionary computation*, 6(2), pp.182-197.

[Zhang' 2007] Zhang, Q. and Li, H., 2007. MOEA/D: A multiobjective evolutionary algorithm based on decomposition. *IEEE Transactions on evolutionary computation*, 11(6), pp.712-731.

[Harman' 2001] Harman, M. and Jones, B.F., 2001. Search-based software engineering. *Information and software Technology*, 43(14), pp.833-839.

[Conte' 1986] Conte, S.D., Dunsmore, H.E. and Shen, Y.E., 1986. *Software engineering metrics and models*. Benjamin-Cummings Publishing Co., Inc..

[Shepperd' 2012] Shepperd, M. and MacDonell, S., 2012. Evaluating prediction systems in software project estimation. *Information and Software Technology*, 54(8), pp.820-827.

[Huang' 2018] Huang, Q., Xia, X., Lo, D. and Murphy, G.C., 2018. Automating Intention Mining. *IEEE Transactions on Software Engineering*.