Some Progress from Derivative-Free Optimization to Experienced Derivative-Free Optimization

Yi-Qi Hu

LAMDA Group, Nanjing University, China 4Paradigm Inc., China

2018.08.28







Automatic Machine Learning



Automatic Machine Learning

Purpose of AutoML:

AutoML aims to make machine learning configuration involve less human expert knowledge.



earning And Mining from Da

Handigm Copyright ©2018 4Paradigm All Rights Reserved





Classification-Based Optimization (I)





Classification-Based Optimization (I)





Classification-Based Optimization (I)





Classification-Based Optimization (II)

Querying complexity:

$$O\left(\frac{1}{|D_{\epsilon}|}\left((1-\lambda)+\frac{\lambda}{\gamma T}\sum_{t=1}^{T}\frac{1-Q\cdot R_{\mathcal{D}_{t}}-\theta}{|D_{\alpha_{t}}|}\right)^{-1}\ln\frac{1}{\delta}\right),$$

where $Q=1/(1-\lambda).$

On Local Lipschitz continuous functions:

I

binary space

$$L_2 \|x - x^*\|_H^{\beta_2} \le f(x) - f(x^*) \le L_1 \|x - x^*\|_H^{\beta_1}$$

continuous space

$$L_2 \|x - x^*\|_2^{\beta_2} \le f(x) - f(x^*) \le L_1 \|x - x^*\|_2^{\beta_1}$$

 $\rightarrow poly(\frac{1}{\epsilon}, n, \frac{1}{\beta_1}, \beta_2, \ln L_1, \ln \frac{1}{L_2}) \cdot \ln \frac{1}{\delta}$



classification-based optimization is efficient



Classification-Based Optimization (II)

Querying complexity:

$$O\left(\frac{1}{|D_{\epsilon}|}\left((1-\lambda)+\frac{\lambda}{\sqrt{2}}\sum_{t=1}^{T}\frac{1-Q\cdot R_{\mathcal{D}_{t}}-\theta}{|D_{\alpha_{t}}|}\right)^{-1}\ln\frac{1}{\delta}\right),$$

where $Q = 1/(1-\lambda).$

Smaller Θ the better: the classifier should be highly randomized Smaller Υ the better: the learnt positive area should be small



Classification-Based Optimization (III)

Classification model design:

Considerations:

1. a classifier with a sampled positive area

Implementation:

learn an axis-parallel region







Classification-Based Optimization (III)

Classification model design:

Considerations:

- 1. a classifier with a sampled positive area
- 2. smaller Θ -> less dependent

Implementation:

learn an axis-parallel region

with randomness



Haradigm Copyright ©2018 4 Paradigm All Rights Reserved



Classification-Based Optimization (III)

Classification model design:

Considerations:

- 1. a classifier with a sampled positive area
- 2. smaller Θ -> less dependent
- 3. smaller Υ -> small positive area

Implementation:

learn an axis-parallel region

with randomness as small as possible



RAndomized COordinate Shrinking (RACOS)



Sequential-Model Optimization (I)

Motivation



RAndomized COordinate Shrinking (RACOS)



Sequential-Model Optimization (II)



Sequential-Model Optimization (III)

Sequential RACOS (SRACOS)

with fixed size of training data set



Strategies of replacing solution:

- replace the solution with the worst evaluation value (WR)
- replace a randomly selected solution (RR)
- replace the solution having the largest distance to the bestso-fat solution (LM)



Sequential-Model Optimization (III)

Querying complexity:

$$O\left(\frac{1}{|D_{\epsilon}|}\left((1-\lambda) + \frac{\lambda}{\gamma T}\sum_{t=1}^{T}\frac{1-Q\cdot R_{\mathcal{D}_{t}}-\theta}{|D_{\alpha_{t}}|}\right)^{-1}\ln\frac{1}{\delta}\right),$$

where $Q = 1/(1-\lambda).$

$$O\left(\max\left\{\frac{1}{|D_{\epsilon}|}\left((1-\lambda)+\frac{\lambda}{\gamma(N-r)}\sum_{t=r+1}^{N}\Phi_{t}\right)^{-1}\ln\frac{1}{\delta},N\right\}\right),\$$
where $\Phi_{t} = \left(1-R_{\mathcal{D}_{t}}-\#X\sqrt{\frac{1}{2}D_{KL}(\mathcal{D}_{t}||\mathcal{U}_{X})}-\theta\right)\cdot$
 $|D_{\alpha_{t}}|^{-1}$ and $\#X$ is the volume of X .

RACOS

SRACOS

THEOREM 2

Ignoring the constant factor and fixing θ and γ , the sequential classification-based optimization algorithm can have a better (or worse) query complexity upper bound than the batch-mode if for any iteration t

$$R_{\mathcal{D}_t^S} < (or >) \frac{1}{1-\lambda} R_{\mathcal{D}_t^B} - \# X \sqrt{\frac{1}{2} D_{KL}(\mathcal{D}_t^S \| \mathcal{U}_X)}.$$

sequential can be better

Aparadigm Copyright ©2018 4Paradigm All Rights Reserved

Experiments (I)

Compared methods:

CMA-ES [Hansen, 2013]: evolutionary strategies in state-of-the-art; DE [Storn, 1997]: differential evolution; CE [Rubinstein, 1997]: cross entropy; IMGPO [Kawaguchi, 2015]: a Bayesian optimization method; RACOS [Yu, 2016]: batch-based RACOS.

Selected functions:











Experiments (II)

Convergence:



Experiments (III)

Direct policy search on reinforcement learning:





Experiments (IV)

Gym tasks:





Experiments (IV)

Half-Cheetah:



SRACOS



CE



RACOS







CMA-ES



IMGPO





Is SRACOS efficient enough?



Why DFO Needs Many Evaluations

Exploration in derivative-free optimization

- the points in i-th iteration
- the best point in i-th iteration
- the points in (i+1)-th iteration
- the points which is better point in i-th iteration



Without gradient, DFO should spend many evaluations on exploration.

If we know the search direction, the exploration on DFO can be avoided.

How can we know the direction?

Haradigm Copyright ©2018 4Paradigm All Rights Reserved



Consider AutoML Problems' Property



We consider a problem distribution, but not a single problem.

We can get directions from experience! [Hu, et al. 2018]

Paradigm Copyright ©2018 4Paradigm All Rights Reserved



Get Experience (I)

In i-th iteration of DFO

We can log the stored samples and label the direction to get experience.

- the points in i-th iteration
 the best point in i-th iteration
- the points in (i+1)-th iteration
- the points which is better point in i-th iteration
- → the bad sample direction
- → the good sample direction



- 1. sample a bad point, and label the direction is negative
- 2. sample a good point, and label the direction is positive

Can this experience be used on new problems?

Not directly!



Get Experience (II)

For a simple example





Get Experience (III)

In 2-D problems

 $m{\kappa}_t = \left[egin{array}{c} m{x}_{t,1}^- - m{x}_t^+ \ m{x}_{t,2}^- - m{x}_t^+ \ m{x}_{t,2}^- - m{x}_t^+ \ m{\vdots} \ m{x}_{t,m}^- - m{x}_t^+ \end{array}
ight]$ centralization feature: x'_{t} $oldsymbol{x}_{t,1}^- - oldsymbol{x}_t^+$ x_t^+ $\ell_t\left(\left[\boldsymbol{\kappa}_t; \boldsymbol{x}_t'\right]\right) = \begin{cases} 1, & f(\boldsymbol{x}_t') < f(\tilde{\boldsymbol{x}}_t); \\ 0, & f(\boldsymbol{x}_t') \ge f(\tilde{\boldsymbol{x}}_t). \end{cases}$ label: $\mathcal{D}_{F_e} = \{ ([\kappa_1; x_1'], \ell_1), ([\kappa_2; x_2'], \ell_2), \dots \}$ \mathcal{D}_{F_e} -Learning the directional model on the experience dataset. directional

Paradigm Copyright ©2018 4Paradigm All Rights Reserved



MLP

model

Utilize Experience

Utilizing the directional model to predict the direction of the next sample.



Discussion: The directional model can predict which sample is most valuable to be evaluated. In this way, ExpSRacos improves efficiency by avoiding wasting many evaluations for exploration.

Experiments (I)

Compared method

We compared ExpSRACOS with CMAES [Hansen et al., 2003], SMAC [Hutter et al., 2011], TPE [Bergstra et al., 2011] and SRACOS [Hu et al., 2017].

Exp 1: On synthetic function

Ackley function:
$$f(\boldsymbol{x}) = -20e^{(-\frac{1}{5}\sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_i - x_i^*)^2})} - e^{\frac{1}{n}\sum_{i=1}^{n}\cos 2\pi(x_i - x_i^*)} + e + 20.$$





Experiments (II)



Average performance:

We test all compared methods with just only 50 evaluation budget.

Settings		EXPSRACOS	SRACOS	CMAES	SMAC	TPE	
X^{shift}	n						
$[-0.1, 0.1]^n$	10	1.67 ±0.42	2.37±0.36	2.85 ± 0.36	3.01±0.21	2.69 ± 0.26	
$[-0.1, 0.1]^n$	20	2.76 ±0.26	$2.87 {\pm} 0.20$	3.41 ± 0.25	3.31 ± 0.13	3.09 ± 0.17	
$[-0.5, 0.5]^n$	10	2.26 ±0.41	2.42 ± 0.36	$2.93 {\pm} 0.39$	$2.95 {\pm} 0.25$	2.71 ± 0.29	
$[-0.5, 0.5]^n$	20	2.89 ±0.27'	2.98 ± 0.22	3.56 ± 0.28	$3.35 {\pm} 0.18$	3.18 ± 0.18	



Experiments (III)

Exp 2: On AutoML tasks

Problem distribution: We optimize hyper-parameters for a specific learning model on different datasets.

Criterion: k-fold cross validation error rate on training data.

Support: This experiment is based on Scikit-learn [Pedregosa et al., 2011].



Learning And Mining from Da

Paradigm Copyright ©2018 4Paradigm All Rights Reserved

Experiments (IV)

Result	s:	Dataset	Optimiza	Generalization performance on testing samples					ples	Default				
			EXPSRACOS	SRACOS	SMAC	TPE	Best C	EXPSRACOS	SRA	COS	SMAC	TPE	Best C	$\mathbf{C}_{\mathbf{V}}$
		Annealing	.0401•	.04600	.0877	.0522	.0590	•0000	.01	000	.0200	.01000	.01000	.0400
		Arcene	.13860	.1656	.2395	.1478	.1257•	.18660	.22	66	.2600	.3000	.1600•	.3700
		Balance S.	.0822•	.08480	.1142	.1042	.1063	.0834•	.23	54	.1666	.1904	.09680	.2300
		Banknote	•0000	•0000	•0000	.0018	•0000	•0000	.00	00•	•0000	•0000	● 0000.	.0000
		Breast C. W.	.0357•	.04050	.0429	.0538	.0466	.06380	.06	85	.0780	.0567•	.0936	.0638
		Car	.02600	.0231•	.0744	.0758	.1439	.3410•	.37	86	.3526	.3757	.34210	.3872
		Chess	.0089•	.00960	.0462	.0497	.0755	.0671•	.11	300	.1312	.1453	.1375	.1109
		CMC	.4315•	.43550	.4451	.4462	.4520	.4155•	.43	460	.4459	.4391	.4391	.4290
		CNAE9	.0439•	.0444	.0477	.04410	.0454	.0416•	.04	93	.0416•	.0509	.0527	.0555
	sets	Credit	.0889•	.09380	.1243	.1244	.0981	.2571•	.27	57	.2661	.25890	.3741	.2517
	tas	Cylinder	.1370•	.15060	.3827	.3189	.4100	.40360	.39	57•	.4220	.4220	.4128	.4128
Ι	Dataset	Opt	imization per	formance	e on trair	ning sam	ples	Generalization performance on testing samples						
		EXPSR	ACOS SRA	cos Sl	MAC	TPE	Best C	EXPSRAC	COS	SR	ACOS	SMAC	TPE	Best C
Source	1 st/2nd/3rc	d 20/4 /	/0 4/1	3/5 1	/1/6	0/2/9	3/2/4	15/6/3		5	5/5/5	3/4/7	2/5/7	4/9/4
dataset	Avg. rank	1.166	57 2.20	083 3.	8750	3.6250	3.6667	1.5000		2.	9583	3.2500	3.1667	2.7917
Taroet		J 9/1/	0 2/6	5/2 0)/0/4	0/0/3	0/2/1	6/3/1		1	14/2	0/1/3	1/1/2	2/1/3
datasat	Avg rank		2	000 3	6000	3 0000	1 2000	1 5000		2	0000	3 0000	3 5000	3 1000
uataset	Avg. runk	1.100	2.0	JUU J.	0000	5.9000	4.2000	1.5000		۷.	9000	5.9000	5.5000	5.1000
			5000	.2370	.2933	.2701	.25100	0(7(.2003	.23910	.23910	.2900
		Seismic	.06240	.061/•	.0763	.06//	.0657	.06/60	.21	21	.0793	.2030	.065/0	.0831
		WDBC	.05510	.05390	.0590	.0440	.0440	.03750	.00	95 20	1051	.0809	.0095	.0782
		WFDC	.1/13•	.1909	.10550	.1057	.1971	.15/0	.20	52	.1931	.1931	.14030	.2439
		Ava rank	20/4/0	2 2083	3 8750	3 6250	3 6667	1 5000	20	583	3 2500	21511	4/9/4 2 7017	-
		Mushroom	0000	0000	0001	0000	0407	15110	2.9.	10	1642	1623	1023	1642
	s		0030	00320	0147	0085	0301	03090	04	68	0816	0392	0768	0509
	set	Snambase	0543	05570	0750	0782	0801	0588	.04	84o	0716	0694	0738	0966
	ata	Statlog S.	.0205•	.02160	.0233	.0238	.0254	.0238	.03	39	.02810	.0432	.0311	.0389
	st d	Wilt	.0165•	.01720	.0204	.0204	.0202	.03140	.01	46•	.1120	.0420	.1060	.1000
	rge	Wine O. R.	.32410	.3177•	.4217	.4209	.4258	.4409•	.48	13	.4720	.4751	.45650	.4472
	Ta	Yeast	.4025•	.41030	.4556	.4561	.4640	.4186•	.43	130	.4417	.4481	.4320	.4983
		Gisette	.0190•	.02030	.0216	.0210	.0270	.0180•	.02	59	.0214	.02100	.0214	.0270
		Jsbach	.2334•	.2535	.2537	.2578	.24740	.3135•	.31	540	.3317	.3408	.3723	.4047
		Nursery	.0064•	.0310	.0340	.0713	.02950	.3141	.30	970	.3350	.3151	.3028•	.3149
		1st/2nd/3rd	9/1/0	2/6/2	0/0/4	0/0/3	0/2/1	6/3/1	1/4	1/2	0/1/3	1/1/2	2/1/3	-
		Avg. rank	1.1000	2.0000	3.6000	3.9000	4.2000	1.5000	2.9	000	3.9000	3.5000	3.1000	-



Experiments (V)

Exp 2.1: On neural network architecture optimization



optimizing CNN architecture for image classification tasks.

Hyper-parameter number: 19

Datasets: MNIST and SVHN. MINIST is the source dataset and SVHN is the target dataset.

Experience dataset: We apply SRacos to optimize hyper-parameters on the source dataset with 200 evaluation budget and repeat for 5 times.

Results:

[Springenberg et al., 2015]

Dataset	0	ptimizatior	n error	ge	HC-Net					
	ExpSRacos	SRacos	SMAC	TPE	ExpSRacos	SRacos	SMAC	TPE		
Source dataset (MNIST)	.0069	.0079	.0103	.0093	.0078	.0081	.0095	.0091	.0083	
Target dataset (SVHN)	.0557	.0617	.0782	.0772	.0567	.0664	.0759	.0796	.0634	





Conclusion



Paradigm Copyright ©2018 4Paradigm All Rights Reserved



The End

Support by:

Joint work with:





Paradigm 第回范式

Al for everyone.

Prof. Yang Yu LAMDA Group, Nanjing University Ph.D. Hong Qian

LAMDA Group, Nanjing University 4Paradigm Inc., China

Thank you for your attention!

Haradigm Copyright ©2018 4Paradigm All Rights Reserved

