Bayesian Optimization with Categorical and Continuous Variables

Dr Vu Nguyen Machine Learning Scientist **amazon**





Outline

• Bayesian optimisation and Illustration

• Mixed categorical-continuous optimisation $y = f(x_1)$



- Mixed optimization in high dimensional space
- Mixed optimization in population-based AutoRL



Hyperparameters Optimization

- ML algorithm's performances depend on hyperparameters.
- Finding the best hyperparameters for the highest performance



Traditional Hyperparameters Tuning

• Grid Search:

- Create a list of values for each parameter.
- Consider all possible combinations of these values.
- Exhaustively evaluate the model and choose the best parameter.

• Random Search:

- Randomly select a parameter to evaluate.
- Select the best parameter.





Grid vs Random vs Bayesian Optimization





Random Search



Bayesian Optimization

5

Grid vs Random vs Bayesian Optimization

Missed



Missed

Grid Search



Random Search



Found

Bayesian Optimization







Blackbox optimisation competition at NeurIPS20

• The top 20 teams uses Bayes Opt.

blac

< 00	x	Login / Sign Up	Leaderboard	Alternate Le	eaderboard	Virtua	al Room	Submit yo	ur entry
	LEADERBOARD								
	Team				Score	Prize	Paper	GitHub	
	Huawei Noah's Ark Lab				93.519	\$6,000	link	repo	
	NVIDIA RAPIDS.AI				92.928	\$4,000	link	repo	
	AutoML.org & IOHprofiler, fe	aturing the switching sq	uirrel (*)		92.551	(*)	link	repo	
	JetBrains Research				92.509	\$3,000	link	repo	
	Duxiaoman DI				92.212	\$1,000	link	repo	
	Optuna Developers (Preferre	d Networks & CyberAge	nt)		91.806	\$1,000	link	repo	
	Ambitious Audemer				91.107				
	jumpshot				91.089				
	KAIST OSI				90.872		link		
	Able Anteater				90.302				
	Oxford BXL				90.143				
	Innovatrics				90.081		link	repo	
	IBM AI RBFOpt				90.050				
	Jim Liu				89.996				
	Jzkay	ttps://bt	<u>oochal</u>	lenge	89.660	<u>m/le</u>	eade	erboa	ard
	Better call Bayes				89.846		link	repo	
	dannynguyen				89.706			repo	
	AlexLekov				89.403				
	ABO				89.354				
	a2i2team				89.237				
	Tiny, Shiny & Don				89.229				

Black-box Optimization

• The relationship from x to y is through the black-box.



Properties of Black-box Function

$$f \colon X \in \mathcal{R}^d \to Y \in \mathcal{R}$$





No derivative form



Expensive to evaluate (in time and cost)

Nothing is known about the function, except a few evaluations y = f(x)

Bayesian Optimization Overview



- Make a series of evaluations $x_1, x_2, \dots x_T$
- Find the optimum using few evaluations



Illustration of Bayes Opt (3 points)

• Given 3 initial observations



Illustration of Bayes Opt (3 points)

• Given 3 initial observations



Illustration of Bayes Opt (4 points)



13

Illustration of Bayes Opt (5 points)

Accuracy f(x)



14

Illustration of Bayes Opt (6 points)



Accuracy f(x)

Illustration of Bayes Opt (7 points)



16

Illustration of Bayes Opt (8 points)



Outline

• Bayesian optimisation and Illustration

• Mixed categorical-continuous optimisation

- Mixed optimization in high dimensional space
- Mixed optimization in population-based AutoRL



Bayes Opt Mixed Categorical – Continuous Input

• Tuning hyperparameters for deep neural network







Bayes Opt Mixed Categorical – Continuous Input

- One-hot encoding:
 - Red: [1,0,0] Green: [0,1,0] Blue: [0,0,1]
- Drawbacks:
 - Make the search space *large*.
 - if C = 4 categories, each has V = 5 choices=> 20 extra dimensions.
 - Non-continuous and non-differentiable space
- Challenging in optimizing mixed-type: categorical continuous



Bayes Opt Mixed Categorical – Continuous Input



1. Continuous variable is specific to categorical variable

2. Continuous is sharing across categorical variables

1. Continuous variable is specific to categorical variable



Kernel parameter is specific to the *kernel type*.

dimensions for the continuous: 1 or 2?

1. Continuous variable is specific to categorical variable



2. Continuous is sharing across categorical variables

Learning rate in *DL* is sharing across *activation types*



categorical f(x, h)continuous

1. Continuous variable is specific to categorical variable



$f_1(x) \quad f_2(x) \quad f_3(x)$

• Independent/local functions

2. Continuous is sharing across categorical variables



continuous

• Global function f(x, h)

Algorithm overview



Optimize h by multi-armed bandits







Setting (1): continuous is categorical-specific





Setting (2): continuous shared across categories





- Using Bayes Opt to find *x* in a global space across all categories *h*.
- Joint kernel for both x and h
 Additive k(x, x') + k(h, h')
 - Multiplicative $k(\mathbf{x}, \mathbf{x}') \times k(\mathbf{h}, \mathbf{h}')$
 - $k(z, z') = (1 \lambda)k(\mathbf{x}, \mathbf{x}')k(\mathbf{h}, \mathbf{h}') + \lambda[k(\mathbf{x}, \mathbf{x}') + k(\mathbf{h}, \mathbf{h}')]$
 - λ can be estimated from the data.

Visualization of the Algorithm





Outline

- Bayesian optimisation and Illustration
- Mixed categorical-continuous optimization

• Mixed optimization in high dimensional space

• Mixed optimization in population-based AutoRL



Mixed optimization with 200 dimensions ?

• High dimension causes problem for optimization.

- <u>Statistical challenge</u>: the search space grows exponentially
- <u>Computational challenge</u>: global optimizers fail to return an optimum within limited time and resource



Local Trust Optimization

• Build multiple trust regions

- Perform local optimization in each trust region
- The local region is narrower for the local optimizer to be successful



Local Trust Optimization

Iteratively

- expand the trust region for exploration if successes
- shrink the trust region for exploitation if failures
- terminate and randomly jump to another region if local search is exhausted

exploration



Given limited observations, the GP estimation is more accurate within a local region

CASMOPOLITAN

- Extending local trust idea for mixed categorical-continuous variables
 - Separate trust regions: <u>categorical</u> and <u>continuous</u>, defined from the best seen location



• After a local search is exhausted, a new trust region is selected using the UCB

exploit promising region

$$\begin{split} \mathbf{h}_{i}^{(0)} &= \arg \max_{\mathbf{h} \in \mathcal{H}} \mu_{gl}(\mathbf{h}; D_{i-1}^{*}) + \sqrt{\beta_{i}} \sigma_{gl}(\mathbf{h}; D_{i-1}^{*}) \\ & \text{explore region with high uncertainty} \end{split}$$

Illustration of Updating Trust Region Center



At unseen nodes: *(i) Predict the expected function* value and uncertainty (ii) Calculate acq func value

Exploit/explore to select the highest value node to guery the black-box

Illustration of the algorithm in the mixed space

Iteratively expand, shrink and restart the local search
Natural extension to handle parallelism



a) *Conditional on <u>categorical</u> inputs,* update the <u>continuous</u> variables for 1 step



Step 2: Acquisition optimisation

b) Conditional on <u>continuous</u> inputs, update the **categorical** variables for 1 step

-0.5

0.0

0.5



c) Repeat a) and b) until convergence for the next objective function query





Outline

- Bayesian optimisation and Illustration
- Mixed categorical-continuous optimisation
- Mixed optimization in high dimensional space

Mixed optimization in population-based AutoRL



Population Based Training (PBT)

Science

f

in

Contents -

News -

Population Based Training of Neural Networks Max Jaderberg Valentin Dalibard Simon Osindero Wojciech M. Czarnecki Jeff Donahue Ali Razavi Oriol Vinyals Tim Green Iain Dunning Karen Simonyan Chrisantha Fernando Koray Kayukcuoglu DeepMind, London, UK

Abstract

Neural networks dominate the modern machine learning landscape, but their training and success still suffer from sensitivity to empirical choices of hyperparameters such as model architecture, loss function, and optimisation algorithm. In this work we present Population Based Training (PBT), a simple asynchronous optimisation algorithm which effectively utilises a fixed computational budget to jointly optimise a population of models and their hyperparameters to maximise performance. Importantly, PBT discovers a schedule of hyperparameter settings rather than following the generally sub-optimal strategy of trying to find a single fixed set to use for the whole course of training. With just a small modification to a typical distributed hyperparameter training framework, our method allows robust and reliable training of models. We demonstrate the effectiveness of PBT on deep reinforcement learning problems, showing faster wall-clock convergence and higher final performance of agents by optimising over a suite of hyperparameters. In addition, we show the same method can be applied to supervised learning for machine translation, where PBT is used to maximise the BLEU score directly, and also to training of Generative Adversarial Networks to maximise the Inception score of generated images. In all cases PBT results in the automatic discovery of hyperparameter schedules and model selection which results in stable training and better final performance.

Read our COVID-19 research and news SHARE **RESEARCH ARTICLE** Human-level performance in 3D multiplayer games with population-based reinforcement learning 🐵 Max Jaderberg*,†, 🐵 Wojciech M. Czarnecki*,†, 🐵 lain Dunning†, Luke Marris, 🕲 Guy Lever, 💿 Antonio Garcia Castañed...

Careers -

Journals

+ See all authors and affiliations Science 31 May 2019: Vol. 364, Issue 6443, pp. 859-865

Article Figures & Data Info & Metrics eLetters PDF

OOLS		
	*	Download Powe
	0	Request Permiss

ARTICLE T

Semail

Print

Alerts

Science

31 May 2019

Vol 364, Issue 6443

Table of Contents

Advertising (PDF)

Classified (PDF)

Masthead (PDF)

Citation tools

Print Table of Conten

Performance Agent 1 (CPU 1) Hyperparameters (Neural network 1) Model Agent 2 (CPU 2) (Neural network 2) Agent 3 (CPU 3) (Neural network 3) Parallel agents

DOI: 10.1126/science.aau6249

39

Population Based Training (PBT)



Population Based Training

Terminates **poor-performing** agents and restarts by:

- Copies a network **weight** from the **best-performing** agent.
- Generate a new hyperparameter by randomly perturb from the best agent.

At the end of this single training process, we obtain the "well-trained" model.



Two Key Advantages of PBT

- Learning a schedule of hyperparameters, such as selecting large learning rate at the beginning and smaller at the later stage.
 - Existing HPO: a single set of hyperparameters are fixed during training.
 - PBT: perfectly adapts the hyperparameters.
- Time efficiency:
 - Existing HPO: repeatedly evaluate the black-box with different hypers.
 - PBT: uses a single training run.

Population Based Bandit (PB2)

- PB2 is an extension of Population Based Training (PBT).
- The hyperparameters update is controlled by time-varying Gaussian process bandit optimization.



J. Parker-Holder, V. Nguyen, and S. Roberts. "Provably efficient online hyperparameter optimization with population-based bandits." NeurIPS'20

Population Based Bandit (PB2)

- PB2 is deployed in Ray Tune library (18k users).
 - Import PB2 from Ray and try in a few line of codes.



PB2-Mix

- PB2 is only applicable for continuous variables.
- PB2-Mix extends PB2 to handle mixed categ-cont variables



J. Parker-Holder, V. Nguyen, S. Desai, and S. Roberts. "Tuning Mixed Input Hyperparameters on the Fly for Efficient Population Based AutoRL." NeurIPS'21

Time-varying Parallel Extension of MAB and BO



Propose *Time-varying* and *Parallel* MAB to select categorical variables *Time-varying parallel BO* has been available in *PB2*

TV.EXP3.M

• *Time-varying* and *Parallel* TV.EXP3.M to select categorical variable

- EXP3 is a popular algorithm in multi-armed bandits
- <u>EXP3.M</u> is the *Parallel* extension for <u>EXP3</u>
- *Time-varying* extension for <u>EXP3.M</u>

Algorithm 5 TV.EXP3.M algorithm Input: $\gamma = \sqrt{\frac{C \ln(C/B)}{(e-1)BT}}$, $\alpha = \frac{1}{T}$, C #categorical choice, T #max iteration, B #multiple play 1: Init $\omega_c = 1, \forall c = 1...C$ and denote $\eta = (\frac{1}{B} - \frac{\gamma}{C})\frac{1}{1-\gamma}$ 2: for t = 1 to T do Iteratively if $\arg \max_{c \in [C]} \omega_c \ge \eta \sum_{c=1}^C \omega(c)$ then ν s.t. $\frac{\nu}{\eta} = \sum_{\omega_t(c) \ge \nu} \nu + \sum_{\omega_t(c) < \omega_t(c)} \omega_t(c)$ Set $S_0 = (c : \omega_t(c) \ge \nu)$ and $\omega_t(S_0) = \nu$ 3: 5: else 6: Set $S_0 = \emptyset$ 7: 8: end if Compute $p_t^c = B\left((1-\gamma)\frac{\omega_c}{\sum_{i=1}^C \omega_c} + \frac{\gamma}{C}\right), \forall c$ 9: $S_t = \text{DepRound} \left(B, \left[p_t^1 p_t^2 \dots p_t^{c=1} \right] \right)$ Observe the reward $g_t(c) = f(h_t = c)$ for $c \in S_t$ 10: Select a batch of arms 11: (categorical variables) $\hat{g}_t(c) = \frac{g_t(c)}{p_t^c}, \forall c \in S_t \text{ and } \hat{g}_t(c) = 0 \text{ otherwise}$ 12: $\forall c \notin S_0$: update $\omega_c = \omega_c \times \exp\left(B\gamma \hat{g}_t(c)/C\right) + \frac{e\alpha}{C} \sum_{i=1}^C \omega_c$ 13: $\forall c = S_0$: update $\omega_c = \omega_c + \frac{e\alpha}{C} \sum_{i=1}^C \omega_c$ 14: 15: end for reward is time-varying

Theorem 1. Set $\alpha = \frac{1}{T}$ and $\gamma = \min\left\{1, \sqrt{\frac{C \ln(C/B)}{(e-1)BT}}\right\}$, we assume the reward distributions changes at arbitrary instances, but the total number of change points is no more than $V \ll \sqrt{T}$ times. The expected regret of TV.EXP3.M satisfies the following sublinear bound

$$\mathbb{E}[R_{TB}] \le [1+e+V]\sqrt{(e-1)\frac{CT}{B}\ln\frac{CT}{B}}.$$

Sublinear rate on the cumulative regret

Uchiya, T., Nakamura, A. and Kudo, M., Algorithms for adversarial bandit problems with multiple plays. In *International Conference on Algorithmic Learning Theory*. 2010 J. Parker-Holder, V. Nguyen, S. Desai, and S. Roberts. "Tuning Mixed Input Hyperparameters on the Fly for Efficient Population Based AutoRL." *NeurIPS'21*

Data Augmentation and Hyper-optimization for AutoRL

• Augmentation types are categorical variables

• Hyperparameters are continuous

• Regularization, PPO clip, learning rate and entropy exploration.



J. Parker-Holder, V. Nguyen, S. Desai, and S. Roberts. "Tuning Mixed Input Hyperparameters on the Fly for Efficient Population Based AutoRL." NeurIPS'21

Takeaway: mixed categorical-continuous Bayes opt



• MAB/BO for mixed optimization

• Local trust region for handling mixed optimization in high dimension

 Parallel time-varying bandit/BO for population-based AutoRL





References

- S. Gopakumar, S. Gupta, S. Rana, V. Nguyen, & S. Venkatesh. Algorithmic assurance: An active approach to algorithmic testing using BO. *NeurIPS'18*
- B. Ru, A. Alvi, A., V. Nguyen, M. Osborne, & S. Roberts. Bayesian optimisation over multiple continuous and categorical inputs. *ICML'20*
- Eriksson et al. "Scalable global optimization via local Bayesian optimization." NeurIPS'19
- X. Wan, V. Nguyen, H. Ha, B. Ru, C. Lu, and M. Osborne. "Think Global and Act Local: Bayes Opt over High-Dimensional Categorical and Mixed Search Spaces." *ICML'21*
- Bogunovic, I., Scarlett, J., & Cevher, V. Time-varying Gaussian process bandit optimization. *AISTATS 2016*
- Jaderberg, M., Dalibard, V., Osindero, S., Czarnecki, W.M., Donahue, J., Razavi, A., Vinyals, O., Green, T., Dunning, I., Simonyan, K. and Fernando, C. *Population based training of neural networks. arXiv preprint arXiv:1711.09846.*
- J. Parker-Holder, V. Nguyen, and S. Roberts. "Provably efficient online hyperparameter optimization with population-based bandits." *NeurIPS'20*
- Uchiya, T., Nakamura, A. and Kudo, M., Algorithms for adversarial bandit problems with multiple plays. In *International Conference on Algorithmic Learning Theory* 2010
- J. Parker-Holder, V. Nguyen, S. Desai, and S. Roberts. "Tuning Mixed Input Hyperparameters on the Fly for Efficient Population Based AutoRL." *NeurIPS'21*



Cocabo: https://github.com/rubinxin/CoCaBO_code

Casmopolitan: https://github.com/xingchenwan/Casmopolitan

- PB2: <u>https://github.com/jparkerholder/PB2</u>
- PB2-Mix: <u>https://github.com/jparkerholder/procgen_autorl</u>

Question and Answer







<u>@nguyentienvu</u>



vu-nguyen.org



github.com/ntienvu