The Vizier Gaussian Process Bandit Algorithm

Xingyou (Richard) Song xingyousong@google.com

On behalf of the Vizier Team

Vizier + Extended Team



Google Vizier (2017)



- Tunes many of Google's research + products
- O(1K) monthly users, O(70M+) objectives
- Invented before Tensorflow!





Notable Users / Downstream Wins

Production

• Search, Ads, Youtube

Tuning Research Results:

- Hardware Design, Robotics
- Protein Design

Backend for Evolution:

- <u>Neural Architecture Search</u>
- <u>Symbolic Algorithm Search</u>









Vertex/Cloud Vizier



- Prod service for external users / businesses
- Shared client API: Easily switch b/w OSS or Cloud



Fundamentally Black-box Optimization

• Optimize unknown function f(x)

- Evaluating f(x) is expensive.
- Small evaluation budget (~100-100K)



Table of Contents

1. Why a Service?

Service

Algorithm

Experiments

- 2. User/Client API: Distributed Tuning
- 3. Bayesian Optimization Intro
- 4. Gaussian Process, Acquisitions
- 5. Batched + Multi-Objective
- 6. Baselines + Benchmarks
- 7. End-to-End Results
- 8. Ablations

Questions?

Why a Service?



The Wide Variety of Scenarios

- Tuning large ML model hyperparameters
- <u>Chemical/Biological processes</u>
- Optimizing cookie recipes

Very different workflows!



Google's new AI learns by baking tasty machine learning cookies

The system "designs excellent cookies", according to its creators



Workflow Possibilities

- Eval Latency: Seconds to Weeks
- Eval Budget: 10¹ to 10⁷ Trials
- Asynchronous or Synchronous (Batched)
- Failed evaluations: Retried or abandoned
- Early Stopping



Benefits of a Service: No Evaluation Assumptions!

- Users have freedom of when to:
 - Request trials
 - Evaluate Trials
 - Report results

Metadata

- Service can preserve data on prior usage
 - Led to OptFormer, OmniPred, and more offline data research!







Setting up Client

Algorithm, search space, and metrics. study_config = vz.StudyConfig(algorithm='GAUSSIAN_PROCESS_BANDIT') study_config.search_space.root.add_float_param('w', 0.0, 5.0) study_config.search_space.root.add_int_param('x', -2, 2) study_config.search_space.root.add_discrete_param('y', [0.3, 7.2]) study_config.search_space.root.add_categorical_param('z', ['a', 'g', 'k'])

study = clients.Study.from_study_config(study_config, owner='my_name', study_id='example')

Tuning Loop

Loop involves:

- Client obtains suggestions from server
- Evaluating suggestions
- Completing suggestions + updating server

```
for i in range(10):
    suggestions = study.suggest(count=1)
    for suggestion in suggestions:
        params = suggestion.parameters
        objective = evaluate(params['w'], params['x'], params['y'], params['z'])
        suggestion.complete(vz.Measurement({'metric_name': objective}))
```

Suggestion Animation (Full)





OSS Vizier: 2022



- Standalone + customizable Python codebase
- User can host service

Open Source Vizier: Reliable and Flexible Black-Box Optimization.
pypi package 0.1.18 C ci passing C docs_test passing
Google AI Blog Getting Started Documentation Installation Citing and Highlights

Questions?



Google DeepMind

Algorithm

Bayesian Optimization Basics

Regressor / Model: Predict outcome y = f(x) from history

- Common Examples
 - Gaussian Process
 - Random Forest
- Guides search: Explore-Exploit
 - Acquisition Function



GP-UCB

Gaussian Process + Upper Confidence Bound (Srinivas, 2010)

• Provably convergent, sublinear regret

Algorithm 1 Gaussian Process Upper Confidence Bound (GP-UCB) with parameter $\beta > 0$

for $t \in \{1, 2, ...\}$ do Compute posteriors to obtain $\mu_t(x)$ and $\sigma_t(x)$ $x_t = \arg \max_{x \in X} \mu_t(x) + \sqrt{\beta} \cdot \sigma_t(x)$ Evaluate $y_t = f(x_t)$ end for

Key Components



Search Space

Core:

- Double: Continuous range [a,b]
- Integer: Integer range [a,b]
- Discrete: Finite set of floats.
- Categorical: Finite set of strings.

Each ParameterSpec also contains:

• Scaling Type (uniform, log)



X-Normalization

Enforce feature coordinates in [0,1]

- DOUBLE, INTEGER, DISCRETE by min/max bounds
- CATEGORICAL with one-hot

 $\widehat{x} \in [0, 1]^D$

Search Space



Y-Warping

Raw Y-values can be huge range (e.g. 10^{-7} to 10^{7}), need normalization

- Half-Rank Warper: Fit poor y-values to normal curve
- **Log-Warper:** Stretch good-values, compress bad-values
- Infeasible Warper: Replace with unpromising values



Gaussian Process Model

Define kernel (distance metric between two points)

• Matern-5/2 is Euclidean-based

$$\delta(\widehat{x}_{i},\widehat{x}_{j})^{2} = 5 \cdot \sum_{d=1}^{D} \frac{\left(\widehat{x}_{i}^{(d)} - \widehat{x}_{j}^{(d)}\right)^{2}}{\lambda^{(d)}}$$
$$K(\widehat{x}_{i},\widehat{x}_{j}) = \alpha^{2} \cdot \left(1 + \delta + \frac{\delta^{2}}{3}\right) \cdot \exp\left(-\delta\right)$$
$$\overline{f} \sim \mathcal{GP}(0, K)$$
$$\widehat{y} \sim \mathcal{N}\left(\overline{f}(\widehat{x}), \exp(\varepsilon_{\log})\right)$$



MAP Estimation

Fit GP hyperparameters to maximize likelihood: $Pr(\{\widehat{x}_s, \widehat{y}_s\}_{s=1}^t, \alpha_{\log}, \overrightarrow{\lambda}_{\log}, \varepsilon_{\log})$

• L-BFGS-B optimizer

$$\log \Pr(\alpha_{\log}) + \log \Pr(\overrightarrow{\lambda}_{\log}) + \log \Pr(\varepsilon_{\log}) + \log \Pr(\{\widehat{x}_s, \widehat{y}_s\}_{s=1}^t; \alpha_{\log}, \overrightarrow{\lambda}_{\log}, \varepsilon_{\log})$$

$\begin{aligned} & \frac{\text{GP Hyperparameter Distribution}}{\alpha_{\log} \sim \mathcal{N}_{[-3,1]}(\log 0.039, 50)} \\ & \lambda_{\log}^{(d)} \sim \mathcal{N}_{[-2,1]}(\log 0.5, 50) \\ & \varepsilon_{\log} \sim \mathcal{N}_{[-10,0]}(\log 0.0039, 50) \end{aligned}$



UCB Acquisition

- Mean + standard deviation
- Explicit exploration-exploitation tradeoff

$$\mathrm{UCB}(x) = \mu_t(x) + \sqrt{eta} \cdot \sigma_t(x)$$



Acquisition: Trust Region

UCB Coefficient is high: $\sqrt{\beta} = 1.8$

UCB(x) is highest around corners.

$$\mathrm{UCB}(x) = \mu_t(x) + \sqrt{eta} \cdot \sigma_t(x)$$

Use trust region around previous "trusted" points!

 $x \mapsto \begin{cases} \text{UCB}(x) & \text{if dist}(\widehat{x}, \text{trusted}) \leq \text{radius} \\ -10^{12} - \text{dist}(\widehat{x}, \text{trusted}) & \text{if dist}(\widehat{x}, \text{trusted}) > \text{radius} \end{cases}$



Acquisition Function Optimization

Before AutoGrad: **Everything in C++.**

• Multithreaded, not GPU-accelerated

Zeroth-Order Optimization (No gradients!)

• "Hopping up a hill"



AF Optimization: Zeroth-Order Benefits

"ls a(x) > a(x')?"

• Invariant to function scale.



AF Optimization: Firefly Algorithm

Pick two random "fireflies" and move "dimmer" towards "brighter" based on distance

$$\widehat{x}_{low} \leftarrow \widehat{x}_{low} + \eta \exp(-\gamma \cdot r^2) \cdot (\widehat{x}_{high} - \widehat{x}_{low}) + \mathcal{N}(0, \omega^2 I_D)$$



AF Optimization: Vectorized Firefly

Write Algorithm in JAX for crucial speedups!

- Vectorization (simultaneous forces from other fireflies)
- Support CATEGORICALs
- Repel from dimmer fireflies



Questions?

Batched Optimization: Motivation

Users request:

- Batch of X's
- More X's w/o evaluating previous

Naively sample from sequential algorithm:

• Argmax(UCB(history)) leads to duplicate trials!



Batched Optimization: Pure Exploration

1. Give dummy values (0) to unevaluated objectives.

$$\mathcal{U}_t = \{x_v, 0\}_{v=1}$$

2. Maximize standard deviation over real + "dummy" history:

$$\sigma_t(x|\mathcal{D}_t \cup \mathcal{U}_t)$$

3. Make sure UCB is still better than historical maximum

$$\sigma_t(x|\mathcal{D}_t \cup \mathcal{U}_t) + \left| \rho \cdot \min\left(\text{UCB}(x|\mathcal{D}_t, \emptyset, \beta_e) - \tau_t, 0 \right) \right.$$
$$\tau_t := \mu_t(x_t^*|\mathcal{D}_t)$$

Multi-Objective Optimization

M metrics, find points which maximize hypervolume

- volume of Pareto frontier
- w.r.t. reference point

Exact computation given points is #P-Hard!



Scalarization: Intuition

Suppose we want to maximize:

$$f(x) = (f^{(1)}(x), f^{(2)}(x))$$

Try weighted sum, ex:

$$s_{(0.4,0.6)}(f(x)) = 0.4 \cdot f^{(1)}(x) + 0.6 \cdot f^{(2)}(x)$$

Hypervolume is integral over polar coordinates!



Hypervolume is integral over polar coordinates!



Hypervolume is integral over polar coordinates!





$$HV(\{y_1,\ldots,y_k\})pprox \mathbb{E}_w[s_w(y-y_{ref})]$$



Metric 1 (higher is better)

Metric 1 (higher is better)

UCB + Hypervolume

Compute UCB's over all metrics:

$$\overrightarrow{\mathrm{UCB}}(x) = (\mathrm{UCB}^{(1)}(x), \ldots, \mathrm{UCB}^{(M)}(x))$$

Scalarize Acquisition:

$$\mathbb{E}_w\left[s_w(\overrightarrow{\mathrm{UCB}}(x))
ight]$$

Hypervolume Improvement:

$$\mathbb{E}_w\left[\max_{y\in\mathcal{D}_t}\{s_w(y)\},s_w(\overrightarrow{\mathrm{UCB}}(x))\}
ight]$$

Questions?

Experiments



Algorithmic Baselines

- Ax / BoTorch (Meta)
- BayesianOptimization
- HEBO
- HyperOpt
- Optuna
- Scikit-Optimize





OPTUN





A



~50-90% of all use-cases!

Method Similarities

TPE-based: HyperOpt, Optuna

GP-based: Ax, BayesianOptimization, HEBO, SkOpt

- Matern-5/2 Kernel (Continuous)
- L-BFGS-B MAP Estimation

$$\delta(\widehat{x}_i, \widehat{x}_j)^2 = 5 \cdot \sum_{d=1}^{D} \frac{\left(\widehat{x}_i^{(d)} - \widehat{x}_j^{(d)}\right)^2}{\lambda^{(d)}}$$

$$K(\widehat{x}_i, \widehat{x}_j) = \alpha^2 \cdot \left(1 + \delta + \frac{\delta^2}{3}\right) \cdot \exp\left(-\delta\right)$$



GP Method Differences

Algorithm	Acquisition Function	Acquisition Optimizer
Ax	Expected Improvement (EI)	L-BFGS-B (Continuous) Hill-Climb (Mixed/Categorical)
Bayesian- Optimization	UCB (Coeff = 2.5)	L-BFGS-B
HEBO	Max(EI, PI, UCB)	NSGA2
SkOpt	Random(EI, LCB, PI)	L-BFGS-B (Continuous + Mixed) Random Search (Categorical)
Vizier (ours)	UCB (Coeff = 1.8)	Firefly

Benchmark Functions

Emphasis on different shapes

- Blackbox Optimization Benchmark (BBOB)
- COMBO
- Multi-Objective (DTLZ, WFG, ZDT)





Most Common: Full Continuous (20D)

 $LogEfficiency(y \mid f, \mathcal{A}) = log\left(\frac{RequiredBudget(y \mid f, Vizier)}{RequiredBudget(y \mid f, \mathcal{A})}\right)$



Full Continuous (20D) Individual Plots



Continuous (Varied Dimension)

Giant "kitchen-sink" spaces





Continuous (Varied Dimension)

Giant "kitchen-sink" spaces



High-Dim Degradation



Pure Categorical

Vizier > Everyone Else significantly.



Pure Categorical

Vizier > Everyone Else significantly.

HEBO crashes on non-boolean spaces.



Hybrid (Continuous + Categorical)



Hybrid (Continuous + Categorical)

Massive drop in baselines!



Batched Case



Batched Case

High batch settings still robust e.g. against HEBO



Multi-Objective (Individual Plots)



Multi-Objective (Individual Plots)

Instability with e.g. HEBO



Multi-Objective (Aggregate)



Questions?

Ablation: Acquisition Optimization

Firefly > L-BFGS-B



Ablation: AF-Optimization vs E2E

- Firefly > L-BFGS-B
- Trust Region (TR) > Without Trust Region



Ablation: What if Ax used UCB = 1.8?

Equalize acquisition functions.

Ax's median still the same

Suggests AF optimizer is very important!



Ablation: Noisy Objectives

Ranking same as before.



Ablation: Latency + GPU Acceleration

All components in JAX:

- Gaussian Process
- Zero-Order Acq. Optimizer
 - Others don't GPU-accelerate



Links

Code: https://github.com/google/vizier

Documentation: https://oss-vizier.readthedocs.io/en/latest/index.html

Al Blog:

https://ai.googleblog.com/2023/02/open-source-vizier-towards-reliable-and.html

Algorithm Paper: https://arxiv.org/abs/2408.11527

Systems Paper: https://arxiv.org/abs/2207.13676

Thanks!