



Efficient Gaussian Process Bandits via Believing only Informative Actions

Alec Koppel

U.S. Army Research Laboratory

Joint with: Amrit Singh Bedi (ARL), Dheeraj Peddireddy (Purdue), Vaneet Aggarwal (Purdue)

Uncertainty Representations in Bandits and Reinforcement Learning
INFORMS Annual Meeting

Nov. 10, 2020



Multi-Armed Bandits (MAB)



Setting: $\mathcal{A} = \{1, \dots, A\}$ possible *arms* (actions) ¹

→ suppose we select arm a_t at time t

⇒ then reward $r_t(a_t)$ is sequentially revealed by environment.

⇒ $r_t(a_t) \sim$ unknown distribution

→ e.g., Bernoulli with unknown mean

¹ Auer, P., Cesa-Bianchi, N., Freund, Y., & Schapire, R. E. (2002). The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1), 48-77.
Bubeck, S., & Cesa-Bianchi, N. (2012). Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems. *Machine Learning*, 5(1), 1-122.



Multi-Armed Bandits (MAB)



Setting: $\mathcal{A} = \{1, \dots, A\}$ possible *arms* (actions) ¹

→ suppose we select arm a_t at time t

→ then reward $r_t(a_t)$ is sequentially revealed by environment.

→ $r_t(a_t) \sim$ unknown distribution

→ e.g., Bernoulli with unknown mean

→ **Goal:** select $\{a_t\} \Rightarrow$ max. cumulative return $\sum_{t=1}^{\infty} r_t(a_t)$

→ want *regret* growth sublinearly in T , i.e., **Reg_T/T** → 0 w/ T

$$\mathbf{Reg}_T = \max_{\pi} \sum_{t=1}^T r_t(\pi) - \sum_{t=1}^T r_t(a_t) \Rightarrow \text{in some probabilistic sense}$$

¹Auer, P., Cesa-Bianchi, N., Freund, Y., & Schapire, R. E. (2002). The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1), 48-77.

Bubeck, S., & Cesa-Bianchi, N. (2012). Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems. *Machine Learning*, 5(1), 1-122.



Multi-Armed Bandits (MAB)



Setting: $\mathcal{A} = \{1, \dots, A\}$ possible *arms* (actions) ¹

→ suppose we select arm a_t at time t

→ then reward $r_t(a_t)$ is sequentially revealed by environment.

→ $r_t(a_t) \sim$ unknown distribution

→ e.g., Bernoulli with unknown mean

→ **Goal:** select $\{a_t\} \Rightarrow$ max. cumulative return $\sum_{t=1}^{\infty} r_t(a_t)$

→ want *regret* growth sublinearly in T , i.e., **Reg_T/T → 0 w/ T**

$$\mathbf{Reg}_T = \max_{\pi} \sum_{t=1}^T r_t(\pi) - \sum_{t=1}^T r_t(a_t) \Rightarrow \text{in some probabilistic sense}$$

$\pi = \pi_{\theta} \mapsto \mathcal{A}$ is time-invariant, e.g, parameterized by $\theta \in \mathbb{R}^d$

¹Auer, P., Cesa-Bianchi, N., Freund, Y., & Schapire, R. E. (2002). The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1), 48-77.

Bubeck, S., & Cesa-Bianchi, N. (2012). Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems. *Machine Learning*, 5(1), 1-122.



Technological Context



Healthcare ²:

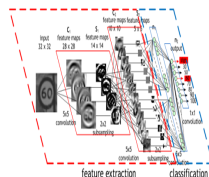
- clinical trials: arm \Rightarrow patient choice, reward \Rightarrow treatment efx.
- personalized dosing: arm \Rightarrow dosage, reward \Rightarrow no side efx.

Recommender systems ³:

- arm \Rightarrow ad/article/movie on page, reward \Rightarrow time on page
- dynamic pricing: arm \Rightarrow price, reward \Rightarrow revenue

Hyperparameter Search in Machine Learning Models ⁴:

- arm \Rightarrow regularizer/step-size, reward \Rightarrow validation accuracy



²Audrey Durand, Charis Achilleos, Demetris Iacovides, Katerina Strati, Georgios D Mitsis, and Joelle Pineau. Contextual bandits for adapting treatment in a mouse model of de novo carcinogenesis. In Machine Learning for Healthcare Conference, pages 67–82, 2018

³Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. A contextual-bandit approach to personalized news article recommendation. CoRR, 2010.

⁴Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., & Talwalkar, A. (2017). Hyperband: A novel bandit-based approach to hyperparameter optimization. The Journal of Machine Learning Research, 18(1), 6765-6816.



Technological Context



Healthcare ²:

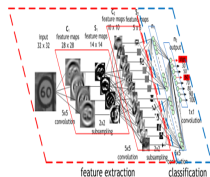
- clinical trials: arm \Rightarrow patient choice, reward \Rightarrow treatment efx.
- personalized dosing: arm \Rightarrow dosage, reward \Rightarrow no side efx.

Recommender systems ³:

- arm \Rightarrow ad/article/movie on page, reward \Rightarrow time on page
- dynamic pricing: arm \Rightarrow price, reward \Rightarrow revenue

Hyperparameter Search in Machine Learning Models ⁴:

- arm \Rightarrow regularizer/step-size, reward \Rightarrow validation accuracy



²Audrey Durand, Charis Achilleos, Demetris Iacovides, Katerina Strati, Georgios D Mitsis, and Joelle Pineau. Contextual bandits for adapting treatment in a mouse model of de novo carcinogenesis. In Machine Learning for Healthcare Conference, pages 67–82, 2018

³Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. A contextual-bandit approach to personalized news article recommendation. CoRR, 2010.

⁴Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., & Talwalkar, A. (2017). Hyperband: A novel bandit-based approach to hyperparameter optimization. The Journal of Machine Learning Research, 18(1), 6765-6816.



Conceptual Context



- Classical solutions to MAB \Rightarrow track statistics of $r_t(a)$ for each a
- \rightarrow UCB ⁵: track mean & std. dev., actions
 - \Rightarrow select arms according to upper conf. bound

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}} \hat{\mu}_t + \beta \hat{\sigma}_t$$

⁵ Lai, T. L. (1987). "Adaptive treatment allocation and the multi-armed bandit problem". The Annals of Statistics, 1091-1114.
Lai, T. L., & Robbins, H. (1985). "Asymptotically efficient adaptive allocation rules." Advances in applied mathematics, 6(1), 4-22.

⁶ Wang, Z., & de Freitas, N. Theoretical analysis of Bayesian optimisation with unknown Gaussian process hyper-parameters. NIPS Workshop on Bayesian Optimization, 2014

⁷ Thompson, W. 1933. "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples". Biometrika. 25(3/4): 285-294.

⁸ Gittins, J. & D. Jones. 1979. "A dynamic allocation index for the discounted multiarmed bandit problem". Biometrika. 66(3): 561-565.



Conceptual Context



- Classical solutions to MAB \Rightarrow track statistics of $r_t(a)$ for each a
- \rightarrow UCB ⁵: track mean & std. dev., actions
 - \Rightarrow select arms according to upper conf. bound

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}} \hat{\mu}_t + \beta \hat{\sigma}_t$$

- \rightarrow Expected Improvement ⁶ operates similarly
- \rightarrow Thompson Sampling ⁷ & Gittins index ⁸:
 - \Rightarrow construct distribution over rewards
 - \Rightarrow select max estimated cond. mean of cumulative return

⁵ Lai, T. L. (1987). "Adaptive treatment allocation and the multi-armed bandit problem". The Annals of Statistics, 1091-1114.
Lai, T. L., & Robbins, H. (1985). "Asymptotically efficient adaptive allocation rules." Advances in applied mathematics, 6(1), 4-22.

⁶ Wang, Z., & de Freitas, N. Theoretical analysis of Bayesian optimisation with unknown Gaussian process hyper-parameters. NIPS Workshop on Bayesian Optimization, 2014

⁷ Thompson, W. 1933. "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples". Biometrika. 25(3/4): 285-294.

⁸ Gittins, J. & D. Jones. 1979. "A dynamic allocation index for the discounted multiarmed bandit problem". Biometrika. 66(3): 561-565.



Conceptual Context



- Classical solutions to MAB \Rightarrow track statistics of $r_t(a)$ for each a
- \rightarrow UCB ⁵: track mean & std. dev., actions
 - \Rightarrow select arms according to upper conf. bound

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}} \hat{\mu}_t + \beta \hat{\sigma}_t$$

- \rightarrow Expected Improvement ⁶ operates similarly
- \rightarrow Thompson Sampling ⁷ & Gittins index ⁸:
 - \Rightarrow construct distribution over rewards
 - \Rightarrow select max estimated cond. mean of cumulative return
- \rightarrow These approaches yield sublinear regret **Reg_T/T \rightarrow 0**
 - \Rightarrow exhibit **computational challenges when # arms A large**
 - \Rightarrow that is, need statistics/posterior with **complexity $\propto A$ or T**

⁵ Lai, T. L. (1987). "Adaptive treatment allocation and the multi-armed bandit problem". The Annals of Statistics, 1091-1114.
Lai, T. L., & Robbins, H. (1985). "Asymptotically efficient adaptive allocation rules." Advances in applied mathematics, 6(1), 4-22.

⁶ Wang, Z., & de Freitas, N. Theoretical analysis of Bayesian optimisation with unknown Gaussian process hyper-parameters. NIPS Workshop on Bayesian Optimization, 2014

⁷ Thompson, W. 1933. "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples". Biometrika. 25(3/4): 285-294.

⁸ Gittins, J. & D. Jones. 1979. "A dynamic allocation index for the discounted multiarmed bandit problem". Biometrika. 66(3): 561-565.



When \mathcal{A} is either continuous, or discrete but A is large scale

→ tracking conditional mean/variance or density is costly

→ **Lipschitz bandits**: discretize space & form bins

⇒ balance **regret**, **number of parameters = # bins** $\propto T^{9/10}$

⁹Magureanu, S., Combes, R., & Proutiere, A. (2014, May). Lipschitz Bandits: Regret Lower Bound and Optimal Algorithms. In Conference on Learning Theory (pp. 975-999).

¹⁰Bubeck, S., Stoltz, G., & Yu, J. Y. (2011, October). Lipschitz bandits without the Lipschitz constant. In International Conference on Algorithmic Learning Theory (pp. 144-158). Springer, Berlin, Heidelberg.

¹¹Srinivas, N., Krause, A., Kakade, S. M., & Seeger, M. W. (2012). Information-theoretic regret bounds for gaussian process optimization in the bandit setting. IEEE Transactions on Information Theory, 58(5), 3250-3265.

¹²Wang, Z., & de Freitas, N. Theoretical analysis of Bayesian optimisation with unknown Gaussian process hyper-parameters. NIPS Workshop on Bayesian Optimization, 2014

¹³Nguyen, V., Gupta, S., Rana, S., Li, C., & Venkatesh, S. (2017, November). Regret for expected improvement over the best-observed value and stopping condition. In Asian Conference on Machine Learning (pp. 279-294).



When \mathcal{A} is either continuous, or discrete but A is large scale

- tracking conditional mean/variance or density is costly
- **Lipschitz bandits**: discretize space & form bins
 - ⇒ balance **regret**, **number of parameters** = # bins $\propto T^{9, 10}$
- **Gaussian Process** ⇒ define posterior directly over \mathcal{A}
 - ⇒ posterior dist. used to define UCB¹¹, EI¹², MPI¹³
 - ⇒ complexity of computing posterior parameters $\mathcal{O}(T^3)$
- **Central question**: can we **define no-regret bandit algorithm**
 - ⇒ whose **complexity remains moderate** as $A, T \rightarrow \infty$

⁹Magureanu, S., Combes, R., & Proutiere, A. (2014, May). Lipschitz Bandits: Regret Lower Bound and Optimal Algorithms. In Conference on Learning Theory (pp. 975-999).

¹⁰Bubeck, S., Stoltz, G., & Yu, J. Y. (2011, October). Lipschitz bandits without the Lipschitz constant. In International Conference on Algorithmic Learning Theory (pp. 144-158). Springer, Berlin, Heidelberg.

¹¹Srinivas, N., Krause, A., Kakade, S. M., & Seeger, M. W. (2012). Information-theoretic regret bounds for gaussian process optimization in the bandit setting. IEEE Transactions on Information Theory, 58(5), 3250-3265.

¹²Wang, Z., & de Freitas, N. Theoretical analysis of Bayesian optimisation with unknown Gaussian process hyper-parameters. NIPS Workshop on Bayesian Optimization, 2014

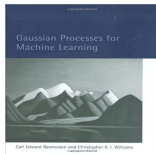
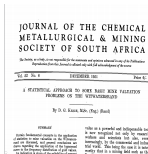
¹³Nguyen, V., Gupta, S., Rana, S., Li, C., & Venkatesh, S. (2017, November). Regret for expected improvement over the best-observed value and stopping condition. In Asian Conference on Machine Learning (pp. 279-294).



Gaussian Processes



- GPs \Rightarrow nonparametric Bayesian method ($\mathcal{A} \subset \mathbb{R}^p, \mathcal{Y} \subset \mathbb{R}$)
 - $\Rightarrow \hat{y} = f(\mathbf{a}) \Rightarrow$ capture relationship of $(\mathbf{a}, y) \in \mathcal{A} \times \mathcal{Y}$
 - \Rightarrow estimate f via $T - 1$ training examples $\mathcal{S} = \{\mathbf{a}_t, y_t\}_{t=1}^{T-1}$.
 - \rightarrow Suppose $f(\mathbf{a})$ in a parameterized family \Rightarrow estimate params
 - \rightarrow Prior $\mathbf{f}_S = [f(\mathbf{a}_1), \dots, f(\mathbf{a}_{T-1})] \Rightarrow$ Gaussian: $\mathbf{f}_S \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{T-1})$
 - \Rightarrow Covariance $\mathbf{K}_{T-1} = [\kappa(\mathbf{a}_{t'}, \mathbf{a}_t)]_{m,t=1}^{T-1, T-1}$ via kernel $\kappa : \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}$
 - \Rightarrow Kernel \Rightarrow distance, e.g., $\kappa(\mathbf{a}_{t'}, \mathbf{a}_t) = \exp\{-\|\mathbf{a}_{t'} - \mathbf{a}_t\|^2 / c^2\}$
 - \rightarrow Standard GP $\Rightarrow \mathbf{f}_S$ observed in noise $\mathbb{P}(\mathbf{y} | \mathbf{f}_S) = \mathcal{N}(\mathbf{f}_S, \sigma^2 \mathbf{I})$
 - \Rightarrow where σ^2 is some variance parameter.





Gaussian Processes



- GPs \Rightarrow nonparametric Bayesian method ($\mathcal{A} \subset \mathbb{R}^p$, $\mathcal{Y} \subset \mathbb{R}$)
- $\Rightarrow \hat{y} = f(\mathbf{a}) \Rightarrow$ capture relationship of $(\mathbf{a}, y) \in \mathcal{A} \times \mathcal{Y}$
 - \Rightarrow estimate f via $T - 1$ training examples $\mathcal{S} = \{\mathbf{a}_t, y_t\}_{t=1}^{T-1}$.
 - \rightarrow Suppose $f(\mathbf{a})$ in a parameterized family \Rightarrow estimate params
 - \rightarrow Upon receiving new sample \mathbf{a}_t , form posterior for \hat{y}_t as

$$\mathbb{P}(y_t | \mathcal{S} \cup \mathbf{a}_t) = \mathcal{N}(\boldsymbol{\mu}_t |_{\mathcal{S}}, \boldsymbol{\Sigma}_t |_{\mathcal{S}})$$

- \Rightarrow where the mean and covariance are given by

$$\begin{aligned}\boldsymbol{\mu}_t |_{\mathcal{S}} &= \mathbf{k}_{\mathcal{S}}(\mathbf{a}_t) [\mathbf{K}_{T-1} + \sigma^2 \mathbf{I}]^{-1} \mathbf{y}_{T-1} \\ \boldsymbol{\Sigma}_t |_{\mathcal{S}} &= \kappa(\mathbf{a}_t, \mathbf{a}_t) \\ &\quad - \mathbf{k}_{\mathcal{S}}^T(\mathbf{a}_t) [\mathbf{K}_{T-1} + \sigma^2 \mathbf{I}]^{-1} \mathbf{k}_{\mathcal{S}}(\mathbf{a}_t)\end{aligned}$$

- $\Rightarrow \mathbf{k}_{\mathcal{S}}(\mathbf{a}) = [\kappa(\mathbf{a}_1, \mathbf{a}); \dots \kappa(\mathbf{a}_{T-1}, \mathbf{a})] \Rightarrow$ empirical kernel map



Gaussian Process Bandits



Gaussian Process posterior \Rightarrow use w/ various action selections

\rightarrow incorporates all past samples into posterior at present

\Rightarrow Upper-Confidence Bound (UCB)¹¹: \Rightarrow via posterior stats.

$$\mathbf{a}_{t+1} = \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} \underbrace{\mu_T |_{S_{T-1}} + \sqrt{\beta_t \Sigma_T |_{S_{T-1}}}}_{:= \alpha_{\text{UCB}}(\mathbf{a})}$$

\rightarrow Expected Improvement (EI)¹³: \Rightarrow action selected as

$$\mathbf{a}_{t+1} = \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} \underbrace{\sigma_{t-1}(\mathbf{a}) \phi(z) + [\mu_{t-1}(\mathbf{a}) - y_{t-1}^{\max}] \Phi(z)}_{:= \alpha_{\text{EI}}(\mathbf{a})},$$

\Rightarrow where $y_{t-1}^{\max} = \max\{y_u\}_{u \leq t}$

$\Rightarrow z = z_{t-1}(\mathbf{a}) = (\mu_{t-1}(\mathbf{a}) - y_{t-1}^{\max}) / \sigma_{t-1}(\mathbf{a})$

$\Rightarrow \phi(z) / \Phi(z)$ denote density/distribution standard Gaussian



Curse of Dimensionality



- Posterior \Rightarrow complexity scales at least $\mathcal{O}(T)$
 - \Rightarrow As time $T \rightarrow \infty \Rightarrow$ GP bandits require **infinite complexity**
- \rightarrow Approaches to compress GPs¹⁴
 - \Rightarrow forward selection¹⁵
 - \Rightarrow variational approx. GP likelihood¹⁶
- \rightarrow Fix **memory M** , “project” onto fixed “subspace,” may **diverge**

¹⁴Williams, C. K., & Rasmussen, C. E. (2006). Gaussian processes for machine learning (Vol. 2, No. 3, p. 4). MIT press.

¹⁵Csat, L., & Opper, M. (2002). Sparse on-line Gaussian processes. *Neural computation*, 14(3), 641-668.

Seeger, Matthias, Christopher Seeger, Matthias, Christopher KI Williams, and Neil D. Lawrence. "Fast forward selection to speed up sparse gaussian process regression." *Workshop on AI and Stats*. 9. 2003.

¹⁶Titsias, M. (2009, April). Variational learning of inducing variables in sparse Gaussian processes. In *AISTATS* (pp. 567-574). Snelson, E., & Ghahramani, Z. (2006). Sparse Gaussian processes using pseudo-inputs. In *NeurIPS* (pp. 1257-1264).

¹⁷Calandriello, D., Carratino, L., Lazaric, A., Valko, M., & Rosasco, L. (2019, June). Gaussian Process Optimization with Adaptive Sketching: Scalable and No Regret. In *Conference on Learning Theory* (pp. 533-557).

¹⁸Chowdhury, S. R., & Gopalan, A. (2019, April). Online learning in kernelized markov decision processes. *AISTATS* (pp. 3197-3205).

¹⁹Calandriello, D., Carratino, L., Valko, M., Lazaric, A., & Rosasco, L. (2020). Near-linear Time Gaussian Process Optimization with Adaptive Batching and Resparsification. *arXiv preprint arXiv:2002.09954*.



Curse of Dimensionality



- Posterior \Rightarrow complexity scales at least $\mathcal{O}(T)$
 \Rightarrow As time $T \rightarrow \infty \Rightarrow$ GP bandits require **infinite complexity**
- \rightarrow Approaches to compress GPs ¹⁴
 - \Rightarrow forward selection ¹⁵
 - \Rightarrow variational approx. GP likelihood ¹⁶
 - \rightarrow Fix **memory M** , “project” onto fixed “subspace,” may **diverge**
 - \rightarrow **Key Idea:** **posterior grows/shrinks w.r.t. data importance**
 - \Rightarrow **quantified by cond. entropy** \Rightarrow **information-theoretic regret**
 - \Rightarrow directly tunable tradeoff between **model complexity/regret**
 - \rightarrow Alternatives use eff. prob. dim., an **offline statistic** ^{17 18 19}

¹⁴Williams, C. K., & Rasmussen, C. E. (2006). Gaussian processes for machine learning (Vol. 2, No. 3, p. 4). MIT press.

¹⁵Csat, L., & Opper, M. (2002). Sparse on-line Gaussian processes. Neural computation, 14(3), 641-668.

Seeger, Matthias, Christopher Seeger, Matthias, Christopher KI Williams, and Neil D. Lawrence. “Fast forward selection to speed up sparse gaussian process regression.” Workshop on AI and Stats. 9. 2003.

¹⁶Titsias, M. (2009, April). Variational learning of inducing variables in sparse Gaussian processes. In AISTATS (pp. 567-574). Snelson, E., & Ghahramani, Z. (2006). Sparse Gaussian processes using pseudo-inputs. In NeurIPS (pp. 1257-1264).

¹⁷Calandriello, D., Carratino, L., Lazaric, A., Valko, M., & Rosasco, L. (2019, June). Gaussian Process Optimization with Adaptive Sketching: Scalable and No Regret. In Conference on Learning Theory (pp. 533-557).

¹⁸Chowdhury, S. R., & Gopalan, A. (2019, April). Online learning in kernelized markov decision processes. AISTATS (pp. 3197-3205).

¹⁹Calandriello, D., Carratino, L., Valko, M., Lazaric, A., & Rosasco, L. (2020). Near-linear Time Gaussian Process Optimization with Adaptive Batching and Resparsification. arXiv preprint arXiv:2002.09954.



Proposed Algorithm



for $t = 1, 2, \dots$ **do**

Select action \mathbf{a}_t via **UCB** or **EI** :

$$\mathbf{a}_t = \arg \max_{\mathbf{a} \in \mathcal{X}} \alpha(\mathbf{a})$$

Sample: $y_t = f(\mathbf{a}_t) + \epsilon_t$, i.e., arm $r_t(\mathbf{a}_t)$

If cond. entropy exceeds ϵ threshold $\mathbf{H}(y_t | \mathbf{y}_{t-1}) > \epsilon$

Augment dict. $\mathbf{D}_t = [\mathbf{D}_{t-1}; \mathbf{a}_t]$, append target $\mathbf{y}_{\mathbf{D}_t} = [\mathbf{y}_{\mathbf{D}_{t-1}}; y_t]$

Update posterior mean $\mu_{\mathbf{D}_t}(\mathbf{a})$ & variance $\sigma_{\mathbf{D}_t}(\mathbf{a})$

$$\mu_{\mathbf{D}_t}(\mathbf{a}) = \mathbf{k}_{\mathbf{D}_t}(\mathbf{a})^T (\mathbf{K}_{\mathbf{D}_t} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}_{\mathbf{D}_t}$$

$$\sigma_{\mathbf{D}_t}^2(\mathbf{a}) = \kappa(\mathbf{a}, \mathbf{a}') - \mathbf{k}_{\mathbf{D}_t}(\mathbf{a})^T (\mathbf{K}_{\mathbf{D}_t, \mathbf{D}_t} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{\mathbf{D}_t}(\mathbf{a}')$$

else

Fix dict. $\mathbf{D}_t = \mathbf{D}_{t-1}$, target $\mathbf{y}_{\mathbf{D}_t} = \mathbf{y}_{\mathbf{D}_{t-1}}$, & GP.

$$(\mu_{\mathbf{D}_t}(\mathbf{a}), \sigma_{\mathbf{D}_t}(\mathbf{a}), \mathbf{D}_t) = (\mu_{\mathbf{D}_{t-1}}(\mathbf{a}), \sigma_{\mathbf{D}_{t-1}}(\mathbf{a}), \mathbf{D}_{t-1})$$

end for

Conditional entropy of GP can be evaluated in closed form as $\mathbf{H}(y_t | \mathbf{y}_{t-1}) = \frac{1}{2} \log (2\pi e(\sigma^2 + \sigma_{\mathbf{D}_{t-1}}^2(\mathbf{a}_t)))$.



Compression Intuition



Matrix of past actions $\mathbf{A}_t = [\mathbf{a}_1; \dots; \mathbf{a}_{t-1}]$

⇒ dictionary \mathbf{D}_t ⇒ subset of columns

⇒ add \mathbf{a}_t only if “significant”

If $\mathbf{H}(y_t | \hat{\mathbf{y}}_{t-1}) > \epsilon$

update $\mathbf{D}_t = [\mathbf{D}_{t-1}; \mathbf{a}_t]$

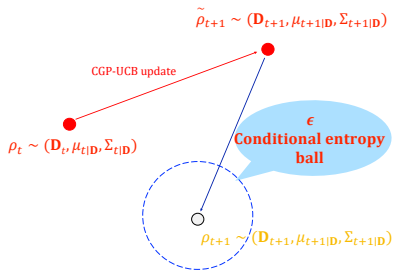
else

update $\mathbf{D}_t = \mathbf{D}_{t-1}$,

→ where

$$\mathbf{H}(y_t | \hat{\mathbf{y}}_{t-1}) = \frac{1}{2} \log(2\pi e(\sigma^2 + \sigma_{\mathbf{D}_{t-1}}^2(\mathbf{a}_t)))$$

⇒ cond. entropy of GP





Regret Bound for GP-UCB



Theorem

For $\delta \in (0, 1)$, the proposed CGP-UCB achieves:

(i) for finite decision set

$$\mathbb{P} \left\{ \mathbf{Reg}_T \leq \sqrt{C_1 T \beta_T \hat{\gamma}_T} + \sqrt{\epsilon T} \right\} \geq 1 - \delta,$$

(ii) for general decision set

$$\mathbb{P} \left\{ \mathbf{Reg}_T \leq \sqrt{C_1 T \beta_T \hat{\gamma}_T} + \sqrt{\epsilon T} + \frac{\pi^2}{6} \right\} \geq 1 - \delta.$$

Theorem

Suppose that the entropy $H(\{y_t\})$ is bounded for all t . Then, the number of elements in the dictionary \mathbf{D}_T denoted by $M_T(\epsilon)$ is finite as $T \rightarrow \infty$ for fixed compression threshold ϵ .



Regret Bound for GP-EI



Theorem

For the finite decision set, with $\delta \in (0, 1)$, the proposed Compressed EI achieves

$$\mathbb{P}\left\{\mathbf{Reg}_T \leq \sqrt{\frac{2T(\gamma_T + \epsilon T)}{\log(1 + \sigma^{-2})}} \left[\sqrt{3(\beta_T + 1 + R^2)} + \sqrt{\beta_T} \right]\right\} \geq 1 - \delta,$$

where

$$R := \sup_{t \geq 0} \sup_{\mathbf{a} \in \mathcal{X}} \frac{|\mu_{t-1}(\mathbf{a}) - y^{\max}|}{\sigma_{t-1}(\mathbf{a})}$$



Experimentation

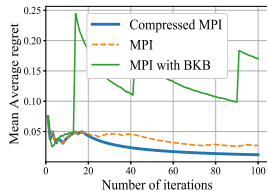
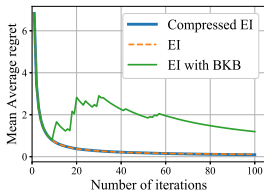
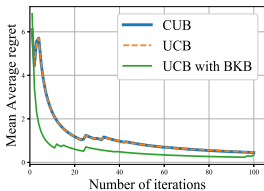


We employ proposed scheme for ML hyperparameter tuning

- ⇒ multi-class classification using CNN on MNIST data set
- ⇒ learning rate, batch size, dropout of inputs, ℓ_2 regularizer
- Instantaneous reward is statistical accuracy

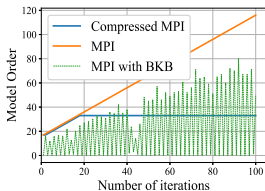
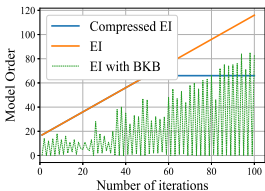
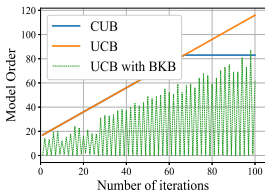
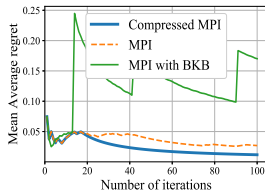
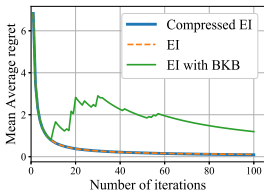
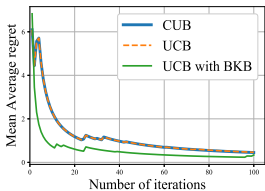


Hyper-parameter Tuning



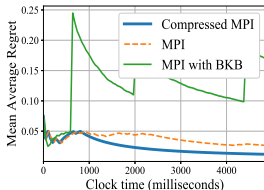
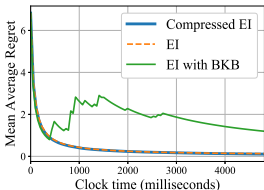
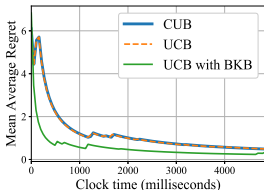
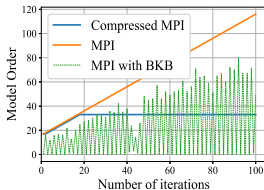
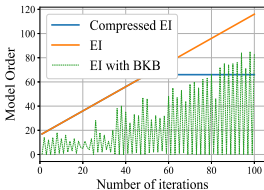
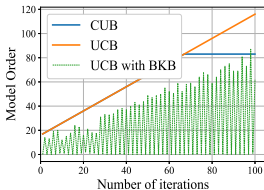
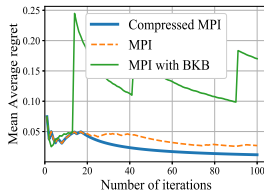
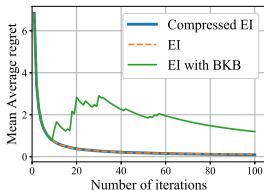
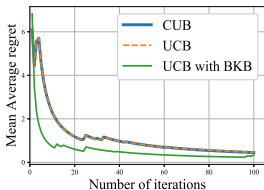


Hyper-parameter Tuning





Hyper-parameter Tuning





Conclusion



We consider the problems with large actions spaces

- Parameterize the action distribution as the GP
- Unfortunately, GP exhibit **complexity challenges**
 - ⇒ **grows cubically with the time**
 - ⇒ Memory requirement grows indefinitely
- Designed entropy-base **compression**
 - ⇒ balance **regret** and **posterior complexity** on the fly
- Future directions :
 - ⇒ extensions to non-stationary (restless) bandits,
 - ⇒ information theoretic compression of CNNs



References



⇒ A. S. Bedi, D. Peddireddy, V. Aggarwal, and A. Koppel, "Efficient Large-Scale Gaussian Process Bandits by Believing only Informative Actions," in Learning for Dynamics and Control (L4DC), University of California, Berkeley, CA, June 2020.

→ A. S. Bedi, D. Peddireddy, V. Aggarwal, and A. Koppel, "Efficient Gaussian Process Bandits by Believing only Informative Actions," arXiv preprint (2020).