

Problem Statement

- **Setup:** A *slate* consists of N slots. At each round $t \in [T]$, the learner receives N sets of items \mathcal{X}_t^i such that $|\mathcal{X}_t^i| = K$.
- **Learner's play:** The learner chooses an item for each slot i , denoted by \mathbf{x}_t^i , and plays the slate $\mathbf{x}_t = (\mathbf{x}_t^1, \dots, \mathbf{x}_t^N)$.
- **Learner's feedback:** The learner receives reward $y_t(\mathbf{x}_t)$ such that $\Pr[y_t(\mathbf{x}_t) = 1 \mid \mathbf{x}_t] = \frac{\exp(\mathbf{x}_t^\top \boldsymbol{\theta}^*)}{1 + \exp(\mathbf{x}_t^\top \boldsymbol{\theta}^*)}$ where $\boldsymbol{\theta}^* \in \mathbb{R}^{Nd}$ is unknown to the learner.
- **Learner's goal:** Minimize *expected cumulative regret* $R_T = \sum_{t=1}^T \max_{\mathbf{x} \in \mathcal{X}_t} \mathbb{E}[y_t(\mathbf{x})] - \mathbb{E}[y_t(\mathbf{x}_t)]$.
- We wish to develop algorithms with optimal regret guarantees that can achieve (1) **Computational Efficiency**, and can handle (2) **Bandit feedback**.
 - Computational Efficiency: Should not enumerate over candidate set with size $2^{\Omega(N)}$.
 - Bandit feedback: A single reward/feedback for the slate played.

Contributions

Algorithms

Experiments

1. **Slate-GLM-OFU:** based on OFU paradigm. Incurs $\tilde{O}(Nd\sqrt{T})$ regret with high probability.
 2. **Slate-GLM-TS:** based on Thompson Sampling principle.
 3. **Slate-GLM-TS-Fixed:** Fixed-Arm version of **Slate-GLM-TS**. Incurs $\tilde{O}((Nd)^{3/2}\sqrt{T})$ regret with high probability.
1. **Slate-GLM-OFU** incurs lowest regret and **Slate-GLM-TS** is competitive with baselines.
 2. Exponential decrease in run time compared to baselines (attributed to *pulling* time).
 3. **Prompt Tuning:** Choose in-context examples for language models for SST2 and Yelp Review. Achieve $\sim 80\%$ accuracy.

Relevant Techniques from Prior Works

1. **Online-Proxy Hessian Matrix** [Fauray et al. 2022]: Replace the optimal Hessian matrix $\mathbf{H}_t = \lambda \mathbf{I} + \sum_{s \in [t]} \dot{\mu}(\mathbf{x}_s^\top \boldsymbol{\theta}^*) \mathbf{x}_s \mathbf{x}_s^\top$ with an online proxy matrix, $\mathbf{W}_t = \lambda \mathbf{I} + \sum_{s \in [t]} \dot{\mu}(\mathbf{x}_s^\top \boldsymbol{\theta}_{s+1}) \mathbf{x}_s \mathbf{x}_s^\top$.
2. **Data-driven condition** [Fauray et al. 2022]: An adaptive condition $\dot{\mu}(\mathbf{x}_s^\top \boldsymbol{\theta}_s) \leq 2\dot{\mu}(\mathbf{x}_s^\top \boldsymbol{\theta}_t^*)$ that helps maintain control over the diameter of the permissible set of parameters.
3. **Distribution for Thompson Sampling** [Abeille et al. 2017]: A distribution \mathcal{D} which enables the noise $\boldsymbol{\eta}$ to explore enough (*anti-concentration*) but not too much (*concentration*).

Algorithm and Our Approach

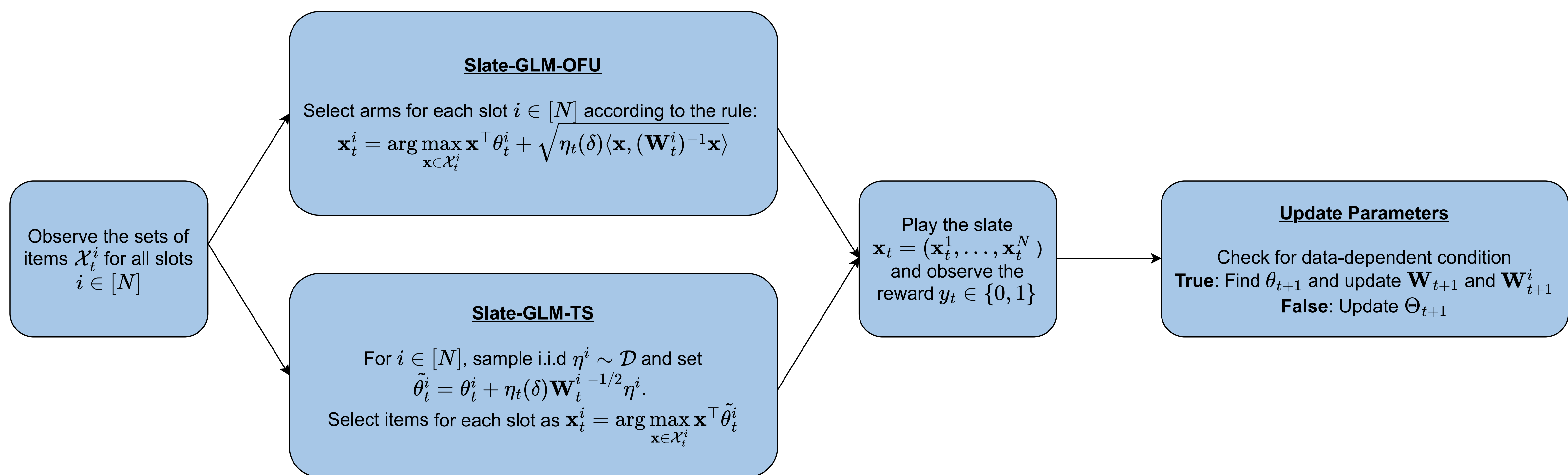


Figure 1. Skeleton of **Slate-GLM-OFU** and **Slate-GLM-TS**

1. To enable independent selection of items for each slot, we require guarantees with respect to the *slot-level online proxy Hessians* $\mathbf{W}^i \forall i \in [N]$. The inner products are manageable.
2. Note that the proxy matrix \mathbf{W} is of the form $\mathbf{x}\mathbf{x}^\top$:

$$\mathbf{x}\mathbf{x}^\top = \begin{bmatrix} \mathbf{x}^1 \mathbf{x}^{1\top} & \mathbf{x}^1 \mathbf{x}^{2\top} & \dots & \mathbf{x}^1 \mathbf{x}^{N\top} \\ \mathbf{x}^2 \mathbf{x}^{1\top} & \mathbf{x}^2 \mathbf{x}^{2\top} & \dots & \mathbf{x}^2 \mathbf{x}^{N\top} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}^N \mathbf{x}^{1\top} & \mathbf{x}^N \mathbf{x}^{2\top} & \dots & \mathbf{x}^N \mathbf{x}^{N\top} \end{bmatrix} = \text{diag}(\mathbf{x}^1 \mathbf{x}^{1\top}, \dots, \mathbf{x}^N \mathbf{x}^{N\top}) + \mathbf{A}$$

Achieving an equivalence between \mathbf{W} and \mathbf{W}^i requires handling the matrix consisting of the off-diagonal terms, denoted by \mathbf{A} .

3. Using the diversity conditions ($\mathbb{E}[\mathbf{x}_t^i \mid \mathcal{F}_t] = 0$ and $\mathbb{E}[\mathbf{x}_t^i \mathbf{x}_t^{i\top} \mid \mathcal{F}_t] \succeq \rho \kappa \mathbf{I}$) enables us to show that $\mathbf{A} \preceq C \cdot \text{diag}(\mathbf{x}^1 \mathbf{x}^{1\top}, \dots, \mathbf{x}^N \mathbf{x}^{N\top})$.
4. Hence, we end up with a *multiplicative equivalence* between \mathbf{W} and $\mathbf{W}^i \forall i \in [N]$, i.e., $(1 - C) \cdot \text{diag}(\mathbf{W}^1, \dots, \mathbf{W}^N) \preceq \mathbf{W} \preceq (1 + C) \cdot \text{diag}(\mathbf{W}^1, \dots, \mathbf{W}^N)$.
5. We thus have

$$\|\mathbf{x}\|_{\mathbf{W}^{-1}} = \left\| \sum_{i=1}^N \mathbf{x}^i \otimes \mathbf{e}_i \right\|_{\mathbf{W}^{-1}} \leq \frac{1}{1 - C} \sum_{i=1}^N \|\mathbf{x}^i \otimes \mathbf{e}_i\|_{(\text{diag}(\mathbf{W}^1, \dots, \mathbf{W}^N))^{-1}} = \frac{1}{1 - C} \sum_{i=1}^N \|\mathbf{x}^i\|_{(\mathbf{W}^i)^{-1}}$$

allowing us to select items for each slot independent of the others.

Experiments

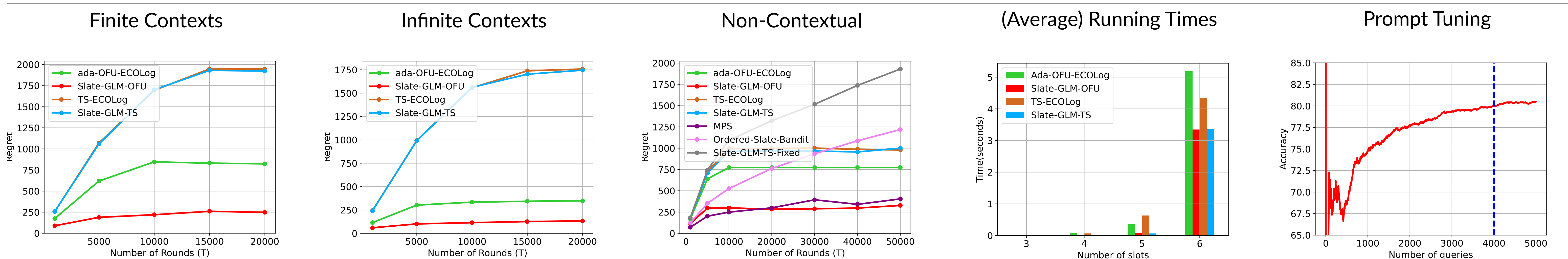


Figure 2. $d = 5, K = 5, N = 3$

Figure 3. $d = 5, K = 5, N = 3$

Figure 4. $d = 5, K = 5, N = 3$

Figure 5. $d = 5, K = 7$

Figure 6. Accuracy v/s. T