

Parameter-Free Online Linear Optimization with Side Information via Universal Coin Betting

Jongha (Jon) Ryu¹, Alankrita Bhatt¹, Young-Han Kim^{1,2}
¹University of California, San Diego ²Gauss Labs Inc.

PROBLEM: ONLINE LINEAR OPTIMIZATION (OLO)

- Assume Hilbert space V with norm $\|\cdot\|$
- In each round $t = 1, 2, \dots$
 - Learner picks action $\mathbf{w}_t \in V$
 - Receives a vector $\mathbf{g}_t \in V$ such that $\|\mathbf{g}_t\| \leq 1$
 - Gains reward $\langle \mathbf{g}_t, \mathbf{w}_t \rangle$
- Goal:** maximize the cumulative reward $\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{w}_t \rangle$
- The standard metric:** regret with respect to the best static competitor in hindsight

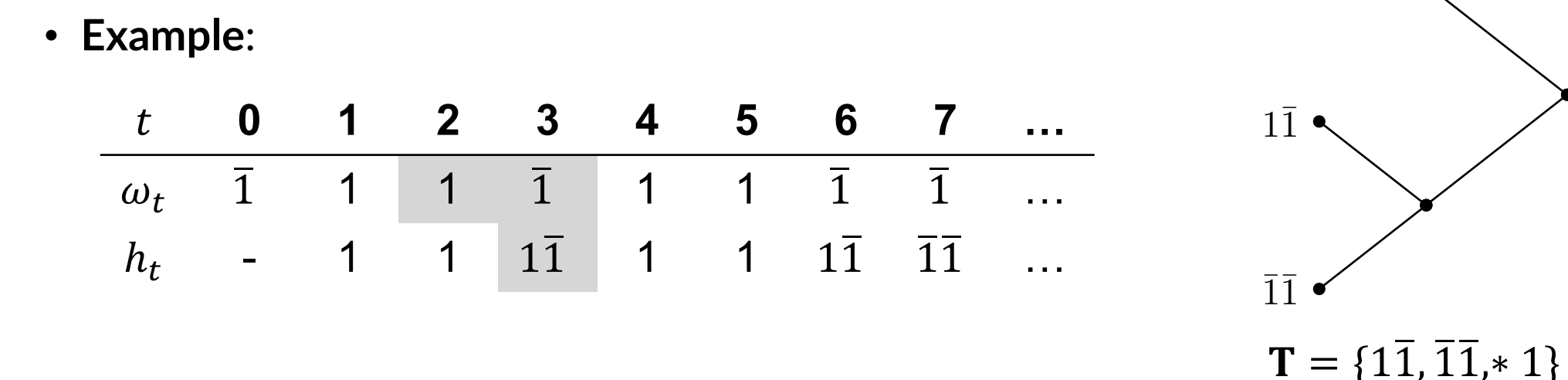
$$\text{Reg}(\mathbf{u}; \mathbf{g}^T) := \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{u} \rangle - \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{w}_t \rangle \text{ for } \mathbf{u} \in V$$
- Two issues**
 - Learning rate tuning requires a priori knowledge on $\|\mathbf{u}\|$
 - Static competitors are weak

Parameter-Free OLO via Universal Coin Betting

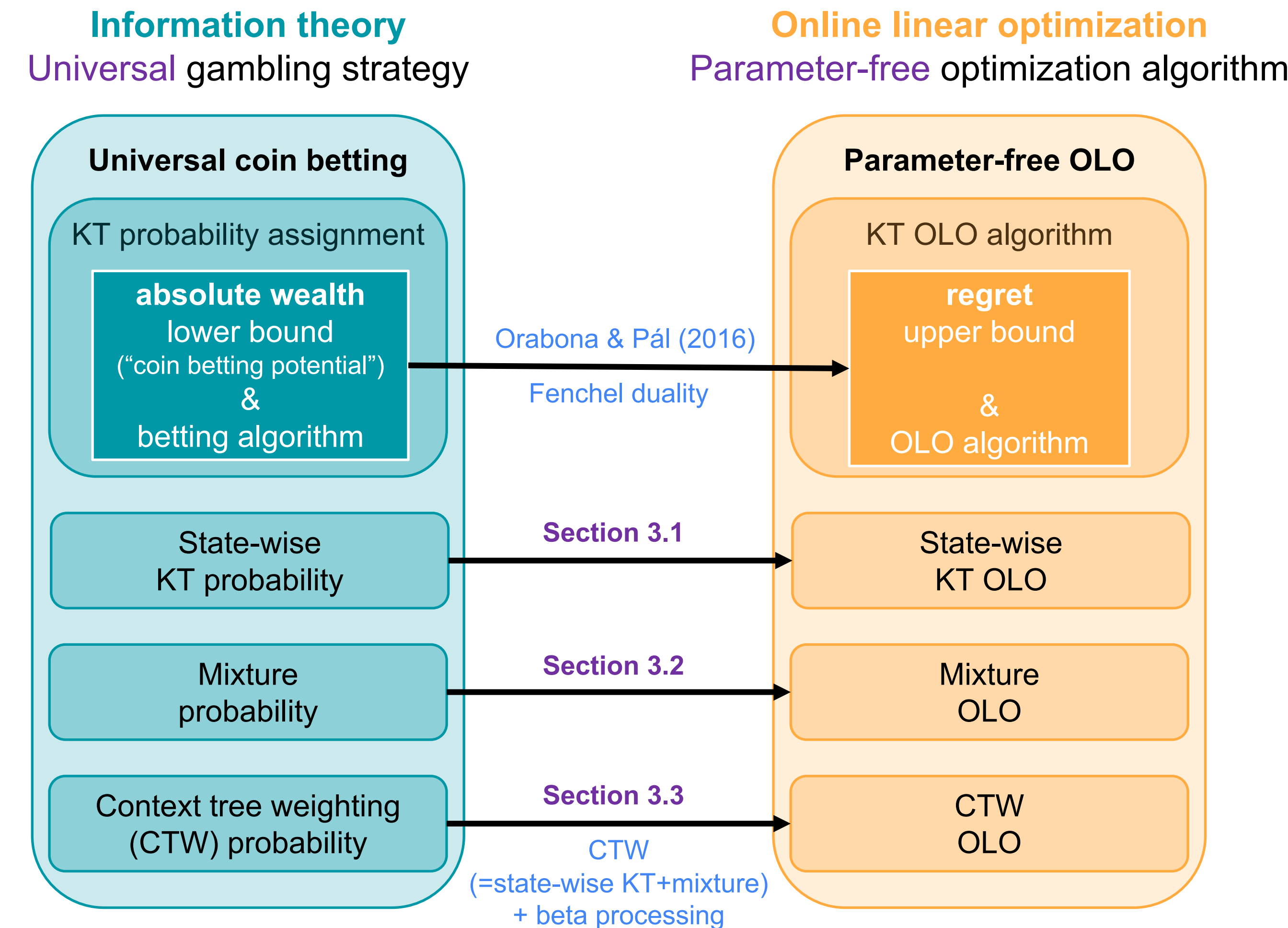
- To attain optimal rates, naive approaches require $\|\mathbf{u}\|$
- Q.** Can we attain optimal regret w/o the need of tuning parameters?
- A.** Orabona and Pál (2016) showed that a universal coin betting algorithm can be converted to a near-optimal-regret parameter-free OLO algorithm!
- Key tool:** Fenchel duality
- Note:** Other parameter-free algorithms exist

OLO with Side Information

- Static competitors $\{\mathbf{u}; \mathbf{u} \in V\}$ are weak**
 - Example:** for $\mathbf{g}, -\mathbf{g}, \mathbf{g}, -\mathbf{g}, \dots$, the best reward with $\mathbf{u} \in V$ is zero
 - In general, $\langle \sum_{t=1}^T \mathbf{g}_t, \mathbf{u}_t \rangle$ can be large iff $\|\sum_{t=1}^T \mathbf{g}_t\|$ is large
- Q.** Can we leverage a possible structure in $(\mathbf{g}_t)_{t \geq 1}$?
- Our approach:** Provided that we have access to a side information sequence $(h_t \in \{1, \bar{1}\})_{t \geq 1}$ which may potentially capture a structure, develop a method that adapts to side information!
- Example:** $h_t = \text{sgn}(\langle \mathbf{g}_{t-1}, \mathbf{f} \rangle)$ (quantization with $\mathbf{f} \in V$)
- To capture a more complex structure, we consider:
- Def (tree side information):** Given a suffix tree \mathbf{T} and an auxiliary sequence $\Omega = (\omega_t \in \{1, \bar{1}\})_{t \geq 1}$, the tree side information $(h_t)_{t \geq 1}$ is $h_t = (\text{the matching suffix of } \omega_t \text{ w.r.t. } \mathbf{T})$



PARAMETER-FREE OLO WITH SIDE INFORMATION VIA UNIVERSAL COIN BETTING



- Idea:** coin betting wealth lower bound can be translated into OLO algorithm AND regret bound
- In general, parameter-freeness incurs additional multiplicative logarithmic factors
- Building blocks**
 - To adapt to a single side information sequence: state-wise KT OLO
 - To adapt to any one of multiple side information sequences: mixture OLO
- Application:** tree side information
 - Goal:** given Ω , adapt to any tree side information sequence of depth $\leq D$
 - Approach:** Take a mixture of state-wise KT OLOs for all tree side information sequences, following CTW (Willems et al. 1995) for universal tree source compression
 - Challenge:** The mixture over all subtrees of depth $\leq D$ involves $O(2^{2^D})$ summands
 - CTW OLO algorithm**
 - can adapt the beta processing algorithm (Willems et al. 2006) for CTW;
 - runs in $O(D)$ time complexity per step, with $O(D)$ storage complexity

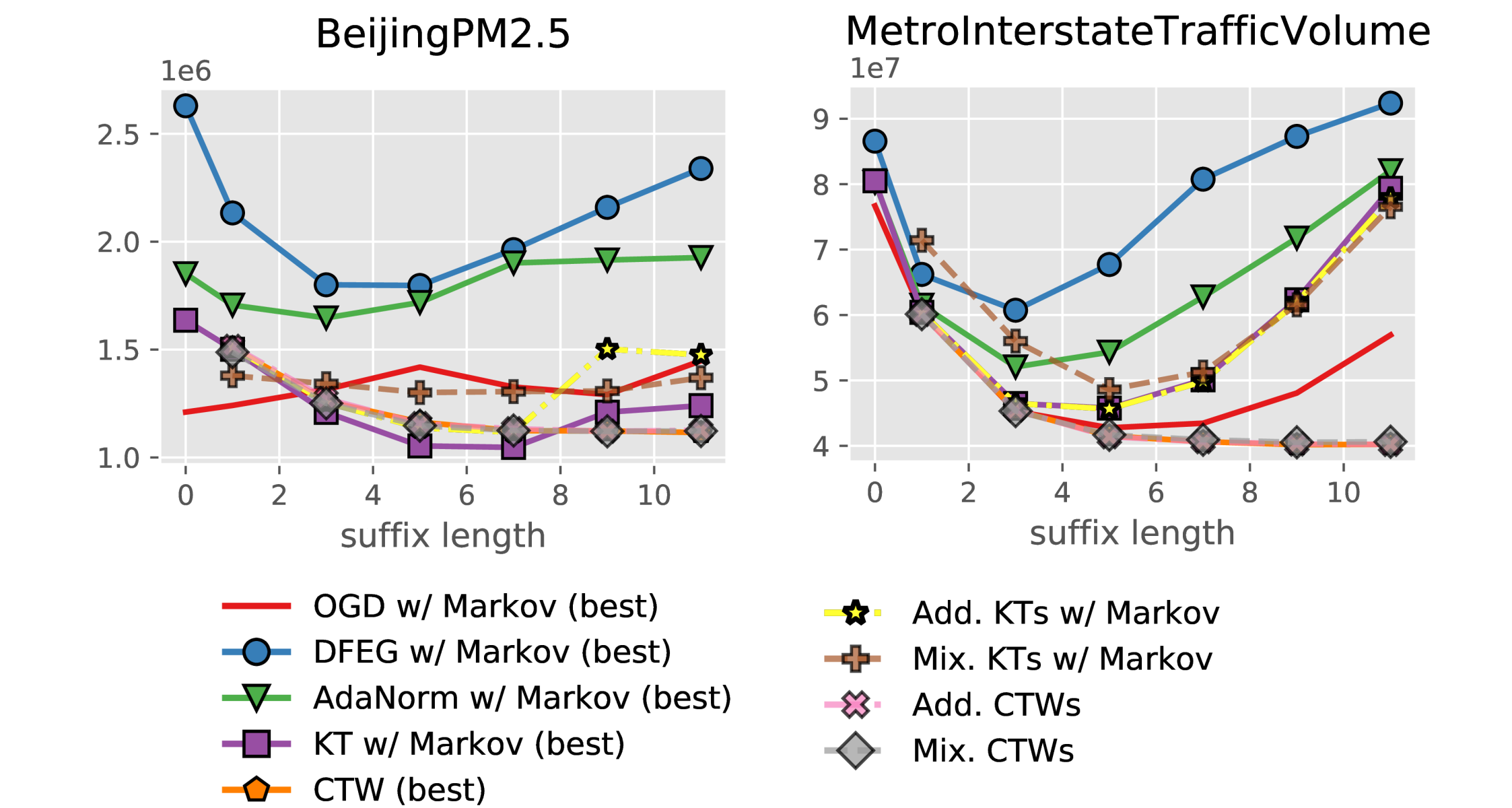


EXPERIMENT

- Online linear regression with absolute loss: $\mathbf{x} \in \mathbb{R}^d, y \in \mathbb{R}$

$$\ell_t(\mathbf{w}_t) \triangleq \ell(\hat{y}_t, y_t) \quad \hat{y}_t \triangleq \langle \mathbf{w}_t, \mathbf{x}_t \rangle$$

$$\partial \ell_t(\mathbf{w}_t) = \text{sgn}(\langle \mathbf{w}_t, \mathbf{x}_t \rangle - y_t) \mathbf{x}_t$$
- Two real-world temporal datasets
 - Beijing PM2.5 (air pollution dataset)
 - Metro Inter State Traffic Volume (traffic volume dataset)
- Auxiliary sequence construction: for each dimension $i \in [d]$, apply canonical binary quantizer Q_{e_i} for each symbol (\mathbf{e}_i = the i -th standard vector)
- Run algorithms with tree side information of depth $D \in \{0, 1, 3, 5, 7, 9, 11\}$
- Since we do not know which depth is best a priori, apply mixture (or addition)



Observations

- The performance of OGD, DFEG, AdaNorm with Markov side information get worse as the side information depth increases
- The best of CTWs over dimensions achieves incurs almost the lowest losses
- The addition or mixture of CTWs over the dimensions attain the performance of the best of optimally tuned OGDs, KT, and CTWs

References

Willems, F. M., Shtarkov, Y. M., and Tjalkens, T. J. (1995). "The context-tree weighting method: Basic properties." In: IEEE Trans. Inf. Theory, 41(3):653-664.

Willems, F. M., Tjalkens, T. J., and Ignatenko, T. (2006). Context-tree weighting and maximizing: Processing betas. In: Proc. UCSD Inf. Theory Appl. Workshop.

Orabona, F. and Pál, D. (2016). Coin betting and parameter-free online learning. In: Adv. Neural Inf. Proc. Syst., volume 29. Curran Associates, Inc.