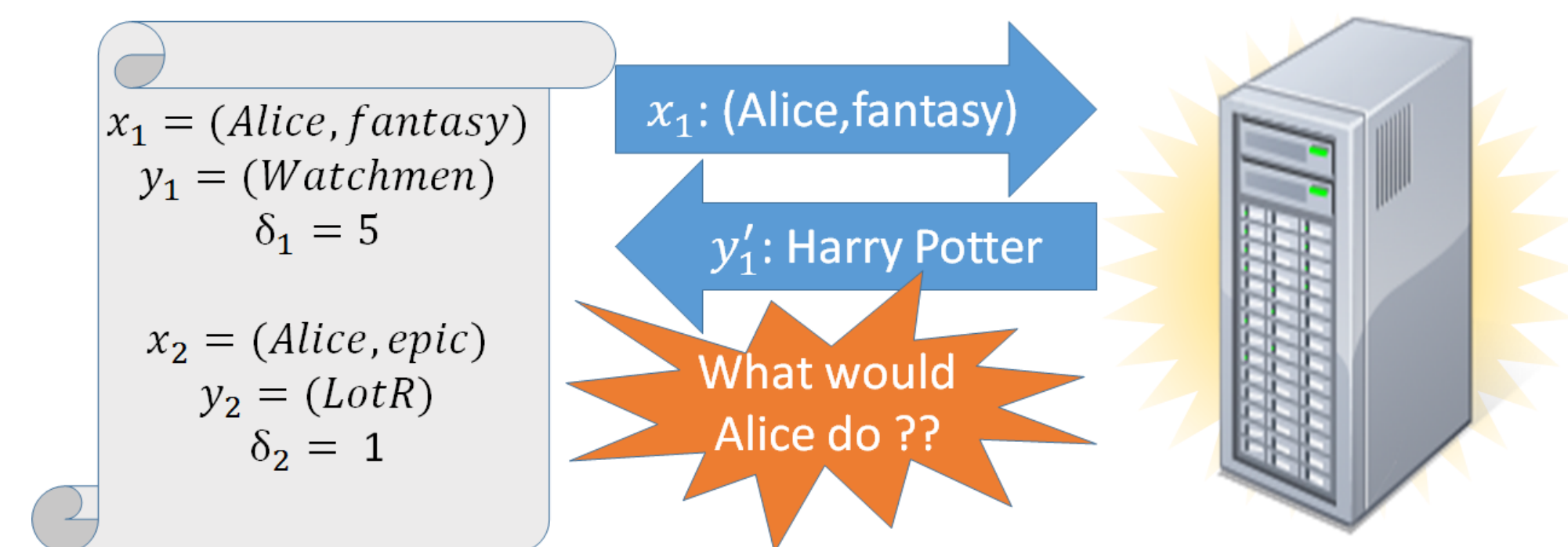


Counterfactual risk minimization: Learning from logged bandit feedback

Adith Swaminathan and Thorsten Joachims
Department of Computer Science, Cornell University

Aim: Offline learning for interactive systems

Can we re-use the interaction logs of deployed on-line systems (e.g. search engines, recommendation systems) to train better models offline?



- Logs are **biased** (actions favored by deployed system will be over-represented),
- and **incomplete** (no feedback for other plausible actions).

Our contribution

A learning principle — **Counterfactual Risk Minimization** — and an efficient algorithm — **Policy Optimizer for Exponential Models** — for this learning setting [1]. Our solution is to

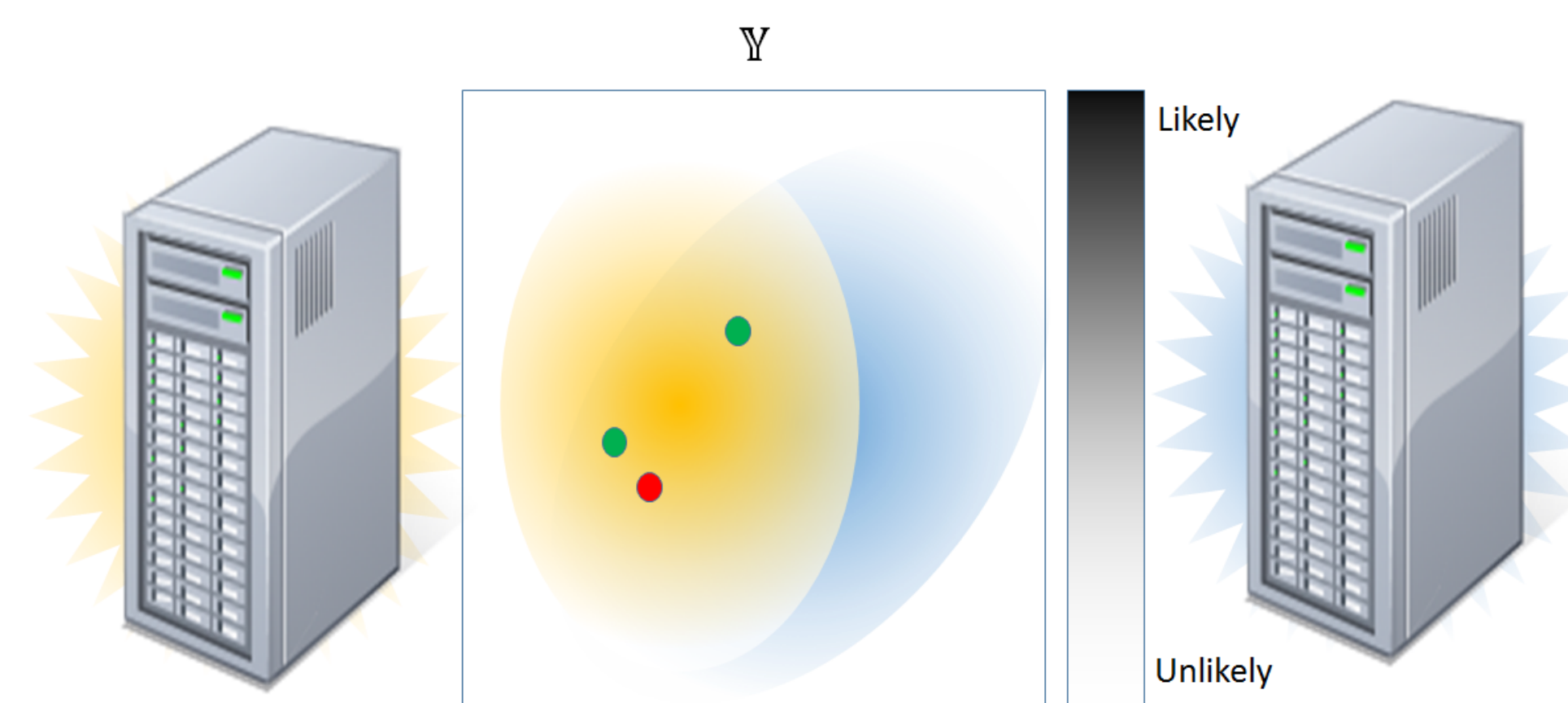
- predict by **sampling** and log **propensities**,
- use **counterfactual** risk estimators to fix bias,
- regularize the variance**,
- and optimize a conservative bound using **majorization minimization**.

POEM

POEM is a simple, fast, stochastic optimizer for structured output prediction available at <http://www.cs.cornell.edu/~adith/poem>

It is as fast and expressive as Conditional Random Fields (CRFs), and trains using logged bandit feedback, without any supervised labels.

Counterfactual estimators



Learning from logged data without exploration is not possible. Suppose the deployed system *sampled* $y \sim h_0(\mathcal{Y} | x)$.

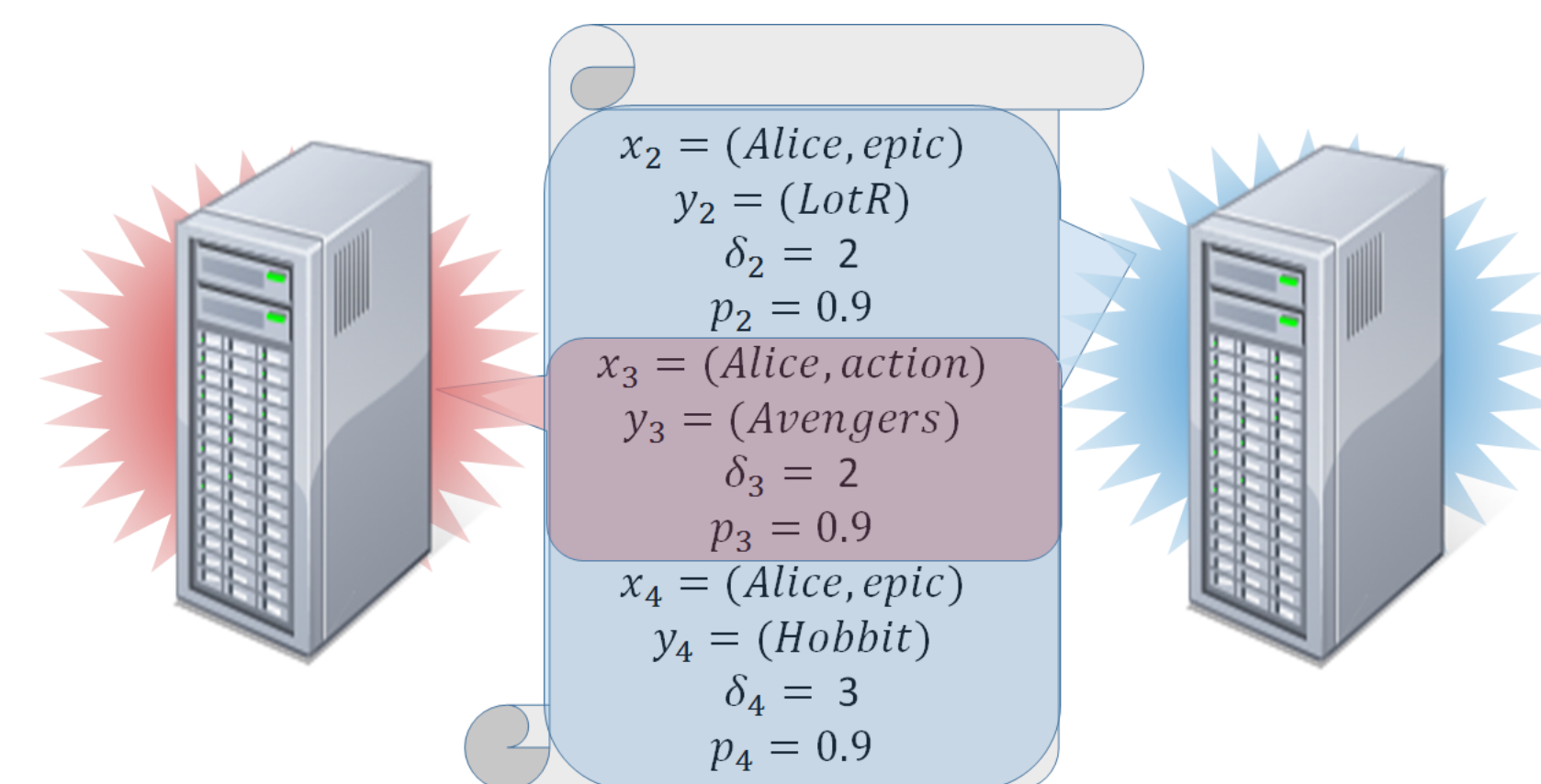
$$\underbrace{\mathbb{E}_x \mathbb{E}_{y \sim h(x)} [\delta(x, y)]}_{R(h), \text{ Risk of } h} = \underbrace{\mathbb{E}_x \mathbb{E}_{y \sim h_0(x)} [\delta(x, y) \frac{h(y | x)}{h_0(y | x)}]}_{\text{Samples from deployed } h_0} \underbrace{\frac{h(y | x)}{h_0(y | x)}}_{\text{Importance weight}}.$$

With $\mathcal{D} = \{(x_i, y_i, \delta_i, p_i)\}_{i=1}^n$, $p_i \equiv h_0(y_i | x_i)$,

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \delta_i \frac{h(y_i | x_i)}{p_i}.$$

This unbiased estimator has issues:

- Unbounded variance (think $p_i \simeq 0$).
- Degenerate minimizer (think $\delta_i \geq 0$).
- Importance sampling introduces variance.



Different effective sample sizes for different h !

Inverse propensity scoring,

$$h^{IPS} = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (\delta_i - \delta^{max}) \underbrace{\min\{M, \frac{h(y_i | x_i)}{p_i}\}}_{\hat{R}^M(h)}.$$

fixes the first two issues. For the variance issue, we employ an empirical Bernstein argument [3].

Variance regularization

With high probability in $\mathcal{D} \sim h_0$, $\forall h \in \mathcal{H}$,

$$\underbrace{R(h)}_{\text{True risk}} \leq \underbrace{\hat{R}^M(h)}_{\text{Empirical risk}} + \underbrace{\mathcal{O}\left(\sqrt{\frac{\hat{Var}(\mathbf{u}_h)}{n}}\right)}_{\text{Variance control}} + \underbrace{M \cdot \mathcal{O}\left(\frac{\mathcal{N}_\infty(\mathcal{H})}{n}\right)}_{\text{Capacity control}}.$$

Counterfactual risk minimization

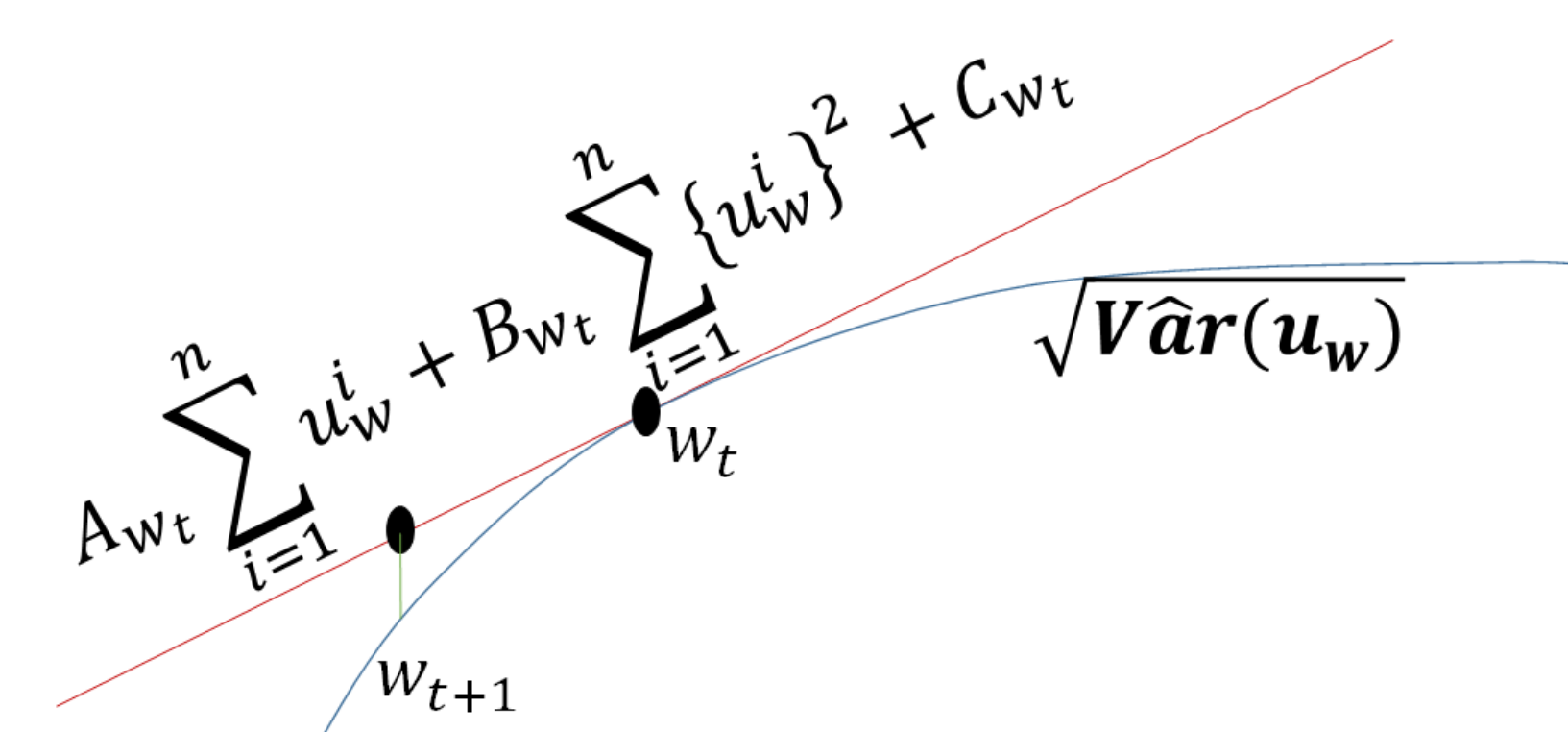
$$h^{CRM} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}^M(h) + \lambda \sqrt{\frac{\hat{Var}(h)}{n}}.$$

Deriving POEM from CRM

For CRFs, $h_w(y | x) \propto e^{w \cdot \phi(x, y)}$.

$$w^{CRM} = \operatorname{argmin}_w \frac{1}{n} \sum_{i=1}^n u_w^i + \lambda \sqrt{\frac{\hat{Var}(\mathbf{u}_w)}{n}} + \mu \|w\|^2.$$

To optimize at scale, we Taylor-approximate.



- Adagrad with $\nabla u_w^i \{1 + \lambda \sqrt{n} (A_{w_t} + 2B_{w_t} u_w^i)\}$.
- After epoch, $w_{t+1} \leftarrow w$, compute $A_{w_{t+1}}, B_{w_{t+1}}$.

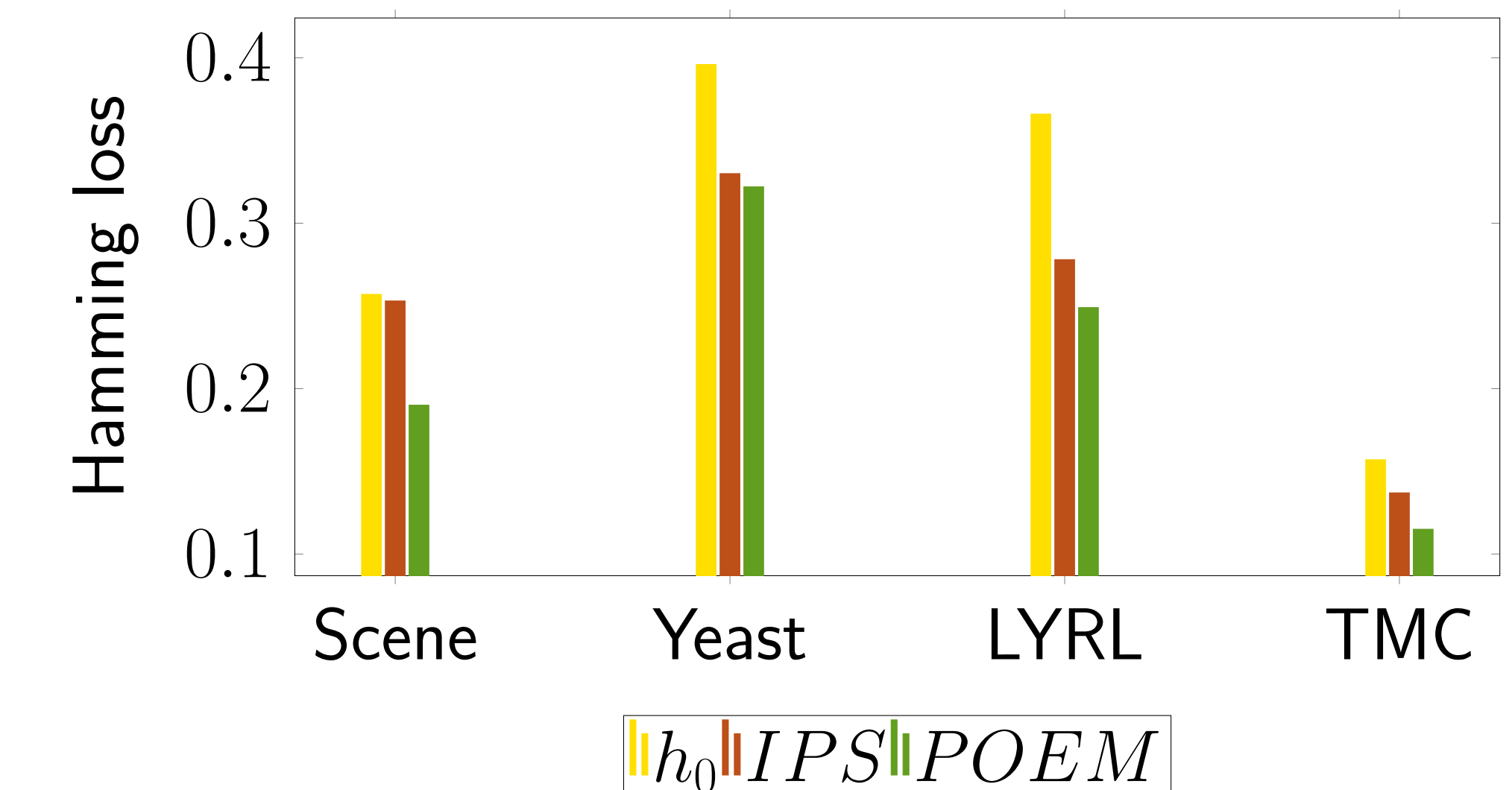
Experiments

Supervised \leftrightarrow Bandit Multi-Label classification with $\delta \equiv$ Hamming loss on four datasets.

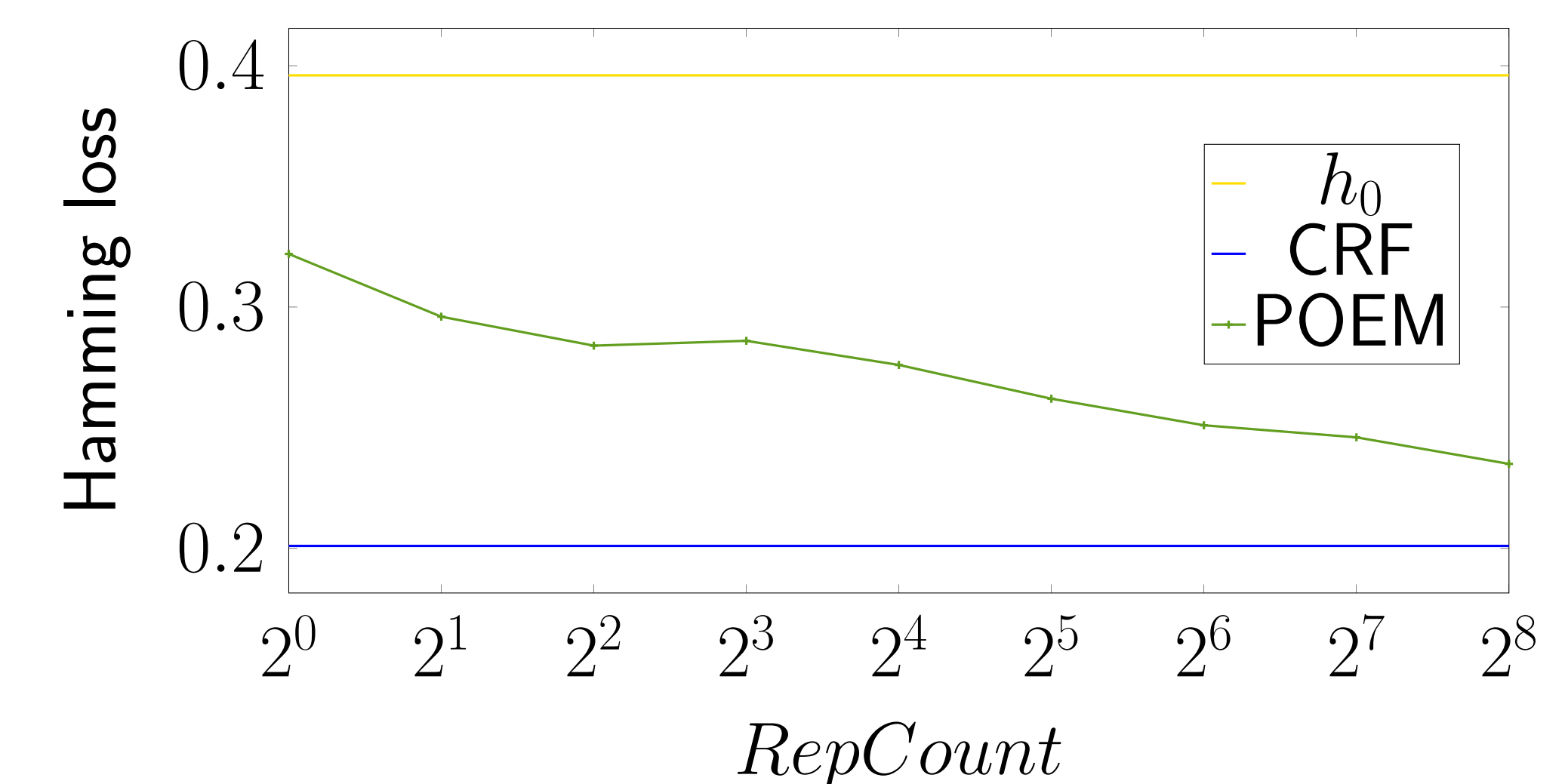
POEM is computationally efficient versus batch L-BFGS and compares favorably with CRF of *scikit-learn*.

Avg. Time (s)	Scene	Yeast	TMC	LYRL
POEM(L-BFGS)	75.20	94.16	949.95	561.12
POEM	4.71	5.02	276.13	120.09
CRF	4.86	3.28	99.18	62.93

POEM is statistically significantly better ($p = 0.05$) than IPS and h_0 (CRF trained on 5% of train set) on all datasets.



POEM recovers supervised performance as $n \rightarrow \infty$ (simulated by $\mathcal{D} \sim h_0$ *RepCount* many times on the *Yeast* dataset).



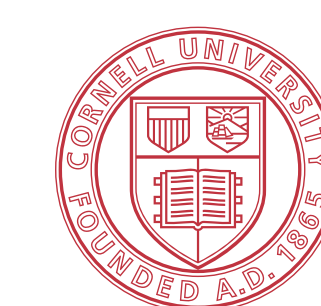
- Non-convex objective, but good local optima.
- Even with poor h_0 , POEM achieves good loss.
- Sweet spot for stochasticity of h_0 .
- MAP predictions from POEM work well.

References

- Adith Swaminathan and Thorsten Joachims. Counterfactual risk minimization: Learning from logged bandit feedback. *ICML*, 2015.
- Léon Bottou, Jonas Peters, Joaquin Q. Candela, Denis X. Charles, Max Chickering, Elon Portugaly, Dipankar Ray, Patrice Y. Simard, and Ed Snelson. Counterfactual reasoning and learning systems: the example of computational advertising. *Journal of Machine Learning Research*, 14(1):3207–3260, 2013.
- Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample-variance penalization. *COLT*, 2009.

Acknowledgment

This research was funded through NSF Award IIS- 1247637, IIS-1217686, JTCII Cornell-Technion Research Fund, and a gift from Bloomberg.



Cornell University