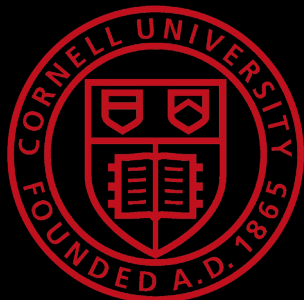# Batch Learning from Bandit Feedback

CS7792 Counterfactual Machine Learning – Fall 2018

Thorsten Joachims

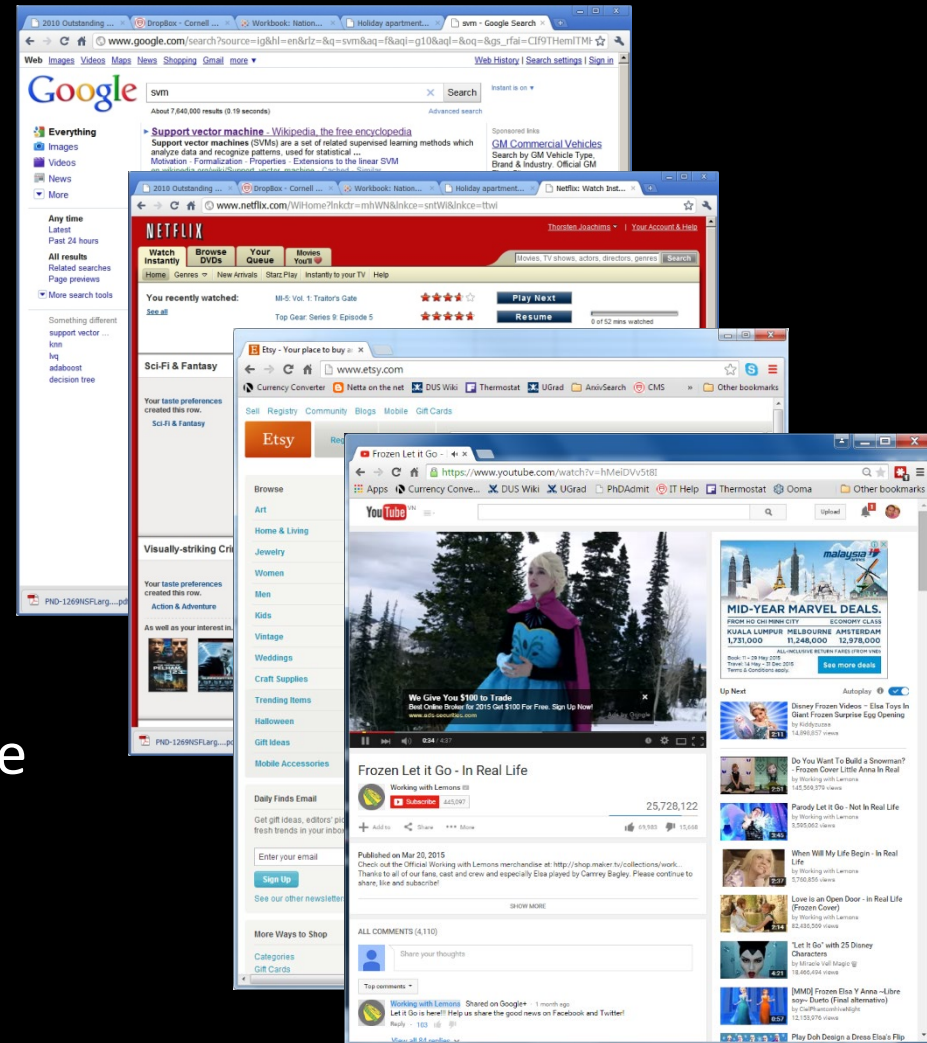Departments of Computer Science and Information Science

Cornell University

- A. Swaminathan, T. Joachims, Batch Learning from Logged Bandit Feedback through Counterfactual Risk Minimization, JMLR Special Issue in Memory of Alexey Chervonenkis, 16(1):1731-1755, 2015.
- T. Joachims, A. Swaminathan, M. de Rijke. Deep Learning with Logged Bandit Feedback. In ICLR, 2018.
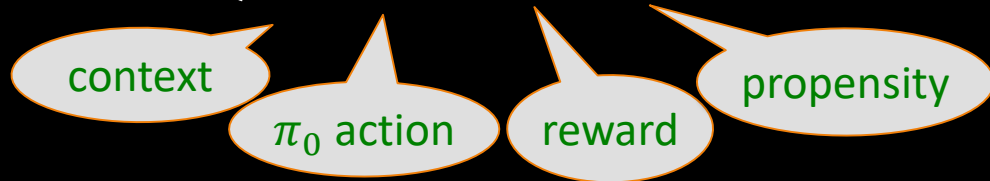
# Interactive Systems

- Examples
  - Ad Placement
  - Search engines
  - Entertainment media
  - E-commerce
  - Smart homes
- Log Files
  - Measure and optimize performance
  - Gathering and maintenance of knowledge
  - Personalization

# Batch Learning from Bandit Feedback

- Data

$$S = \left((x_1, y_1, \delta_1, p_1), \ldots, (x_n, y_n, \delta_n, p_n)\right)$$

context
$\pi_0$ action
reward
propensity

→ "Bandit" Feedback

- Properties
  - Contexts $x_i$ drawn i.i.d. from unknown $P(X)$
  - Actions $y_i$ selected by existing system $\pi_0: X \rightarrow Y$
  - Loss $\delta_i$ drawn i.i.d. from unknown $P(\delta_i|x_i, y_i)$

- Goal of Learning
  - Find new system $\pi$ that selects $y$ with better $\delta$

[Zadrozny et al., 2003] [Langford & Li], [Bottou, et al., 2014]

# Learning Settings

|  | **Full-Information (Labeled) Feedback** | **Partial-Information (e.g. Bandit) Feedback** |
|---|---|---|
| Online Learning | • Perceptron<br>• Winnow<br>• Etc. | • EXP3<br>• UCB1<br>• Etc. |
| Batch Learning | • SVM<br>• Random Forests<br>• Etc. | ? |

# Comparison with Supervised Learning

|  | **Batch Learning from Bandit Feedback** | **Conventional Supervised Learning** |
|---|---|---|
| Train example | $(x, y, \delta)$ | $(x, y^*)$ |
| Context $x$ | drawn i.i.d. from unknown $P(X)$ | drawn i.i.d. from unknown $P(X)$ |
| Action $y$ | selected by existing system $h_0: X \rightarrow Y$ | N/A |
| Feedback $\delta$ | Observe $\delta(x, y)$ only for $y$ chosen by $h_0$ | Assume known loss function $\Delta(y, y^*)$ $\rightarrow$ know feedback $\delta(x, y)$ for every possible y |

# Outline of Lecture

- Batch Learning from Bandit Feedback (BLBF)

$$S = \left((x_1, y_1, \delta_1), \dots, (x_n, y_n, \delta_n)\right)$$

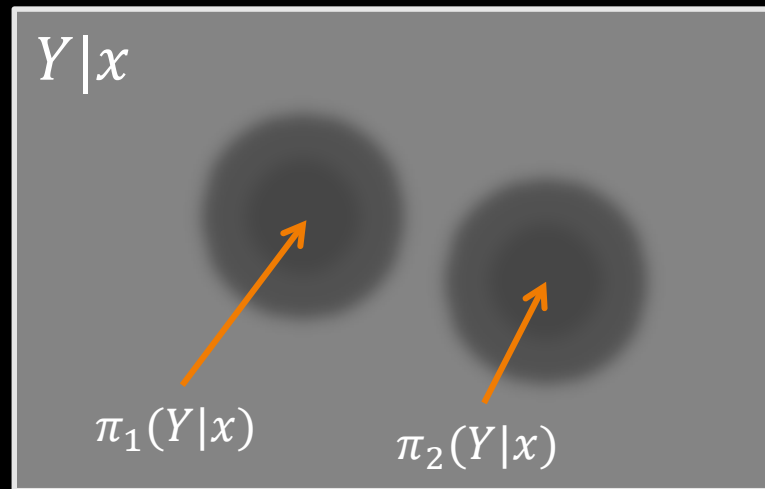→ Find new policy $\pi$ that selects $y$ with better $\delta$

- Learning Principle for BLBF
  - Hypothesis Space, Risk, Empirical Risk, and Overfitting
  - Learning Principle: Counterfactual Risk Minimization
- Learning Algorithms for BLBF
  - POEM: Bandit training of CRF policies for structured outputs
  - BanditNet: Bandit training of deep network policies

# Hypothesis Space

Definition [Stochastic Hypothesis / Policy]:

Given context $x$, hypothesis/policy $\pi$ selects action $y$ with probability $\pi(y|x)$
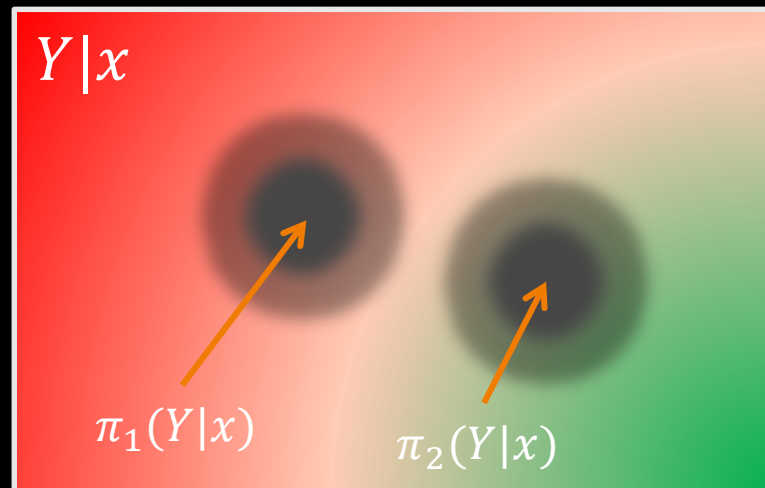
$Y|x$

$\pi_1(Y|x)$     $\pi_2(Y|x)$

Note: stochastic prediction rules $\supset$ deterministic prediction rules

Definition [Expected Loss (i.e. Risk)]:
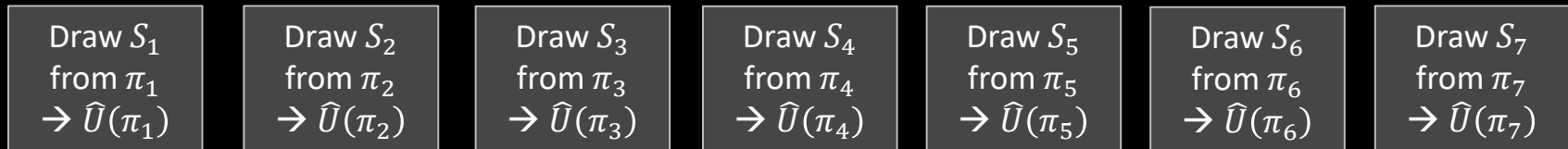
The expected loss / risk $R(\pi)$ of policy $\pi$ is

$$R(\pi) = \int \int \delta(x, y)\pi(y|x)P(x) \, dx \, dy$$

# Evaluating Online Metrics Offline

- Online: On-policy A/B Test

| Draw $S_1$ from $\pi_1$ $\rightarrow \hat{U}(\pi_1)$ | Draw $S_2$ from $\pi_2$ $\rightarrow \hat{U}(\pi_2)$ | Draw $S_3$ from $\pi_3$ $\rightarrow \hat{U}(\pi_3)$ | Draw $S_4$ from $\pi_4$ $\rightarrow \hat{U}(\pi_4)$ | Draw $S_5$ from $\pi_5$ $\rightarrow \hat{U}(\pi_5)$ | Draw $S_6$ from $\pi_6$ $\rightarrow \hat{U}(\pi_6)$ | Draw $S_7$ from $\pi_7$ $\rightarrow \hat{U}(\pi_7)$ |

- Offline: Off-policy Counterfactual Estimates

Draw $S$ from $\pi_0$ $\rightarrow$

$\hat{U}(\pi_6)$   $\hat{U}(\pi_{12})$   $\hat{U}(\pi_{18})$   $\hat{U}(\pi_{24})$   $\hat{U}(\pi_{30})$

# Approach 1: Direct Method

- Data:
  $$S = \left((x_1, y_1, \delta_1), \dots, (x_n, y_n, \delta_n)\right)$$

1. Learn reward predictor
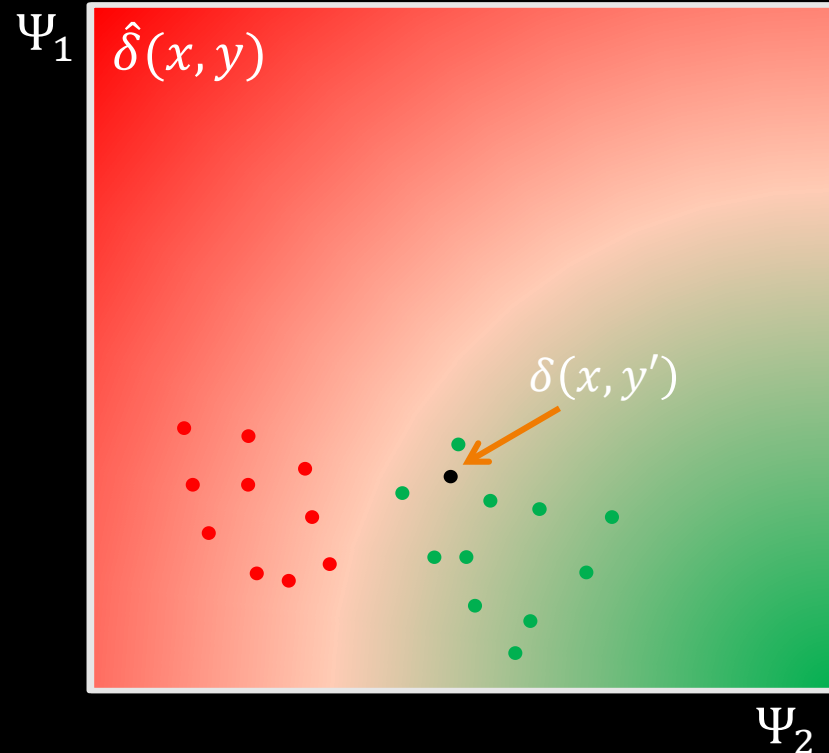   $$\hat{\delta}: x \times y \to \Re$$

   Represent via features $\Psi(x, y)$

   Learn regression based on $\Psi(x, y)$ from $S$ collected under $\pi_0$

2. Derive policy $\pi(x)$
   $$\pi(x) \stackrel{\text{def}}{=} \underset{y}{\text{argmax}}\left[\hat{\delta}(x, y)\right]$$

# Approach 2:
# Off-Policy Risk Evaluation

Given $S = \big((x_1, y_1, \delta_1), \ldots, (x_n, y_n, \delta_n)\big)$ collected under $\pi_0$,

$$\hat{R}(\pi) = \frac{1}{n} \sum_{i=1}^{n} \delta_i \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)}$$

Propensity $p_i$



$\pi_0(Y|x)$

$\pi(Y|x)$

→ Unbiased estimate of risk, if propensity nonzero everywhere (where it matters).

[Horvitz & Thompson, 1952] [Rubin, 1983] [Zadrozny et al., 2003] [Langford, Li, 2009.]

# Partial Information Empirical Risk Minimization



$\pi_0(Y|x)$

$\pi_1(Y|x)$

$\pi_0(Y|x)$

$\pi_{237}(Y|x)$

- Training

$$\hat{\pi} := \operatorname{argmin}_{\pi \in H} \sum_i^n \frac{\pi(y_i|x_i)}{p_i} \, \delta_i$$

[Zadrozny et al., 2003] [Langford & Li], [Bottou, et al., 2014]

# Generalization Error Bound for BLBF

- Theorem [Generalization Error Bound]
  - For any hypothesis space $H$ with capacity $C$, and for all $\pi \in H$ with probability $1 - \eta$

$$\mathrm{R}(\pi) \leq \hat{R}(\pi) + O\left(\sqrt{\widehat{Var}(\pi)/n}\right) + O(C)$$

Unbiased Estimator    Variance Control    Capacity Control

$$\hat{R}(\pi) = \widehat{Mean}\left(\frac{\pi(y_i|x_i)}{p_i}\delta_i\right)$$

$$\widehat{Var}(\pi) = \widehat{Var}\left(\frac{\pi(y_i|x_i)}{p_i}\delta_i\right)$$

→ Bound accounts for the fact that variance of risk estimator can vary greatly between different $\pi \in \mathrm{H}$

[Swaminathan & Joachims, 2015]

# Counterfactual Risk Minimization

- Theorem [Generalization Error Bound]

$$\mathrm{R}(\pi) \leq \hat{R}(\pi) + O\left(\sqrt{\widehat{Var}(\pi)/n}\right) + O(C)$$

→ Constructive principle for designing learning algorithms

$$\pi^{crm} = \underset{\pi \in H_i}{\mathrm{argmin}}\; \hat{R}(\pi) + \lambda_1\left(\sqrt{\widehat{Var}(\pi)/n}\right) + \lambda_2 C(H_i)$$

$$\hat{R}(\pi) = \frac{1}{n}\sum_i^n \frac{\pi(y_i|x_i)}{p_i}\delta_i \qquad \widehat{Var}(\pi) = \frac{1}{n}\sum_i^n \left(\frac{\pi(y_i|x_i)}{p_i}\delta_i\right)^2 - \hat{R}(\pi)^2$$

# Outline of Lecture

- Batch Learning from Bandit Feedback (BLBF)
  $$S = \big( (x_1, y_1, \delta_1), \dots, (x_n, y_n, \delta_n) \big)$$
  $\rightarrow$ Find new policy $\pi$ that selects $y$ with better $\delta$

- Learning Principle for BLBF
  - Hypothesis Space, Risk, Empirical Risk, and Overfitting
  - Learning Principle: Counterfactual Risk Minimization

- Learning Algorithms for BLBF
  - POEM: Bandit training of CRF policies for structured outputs
  - BanditNet: Bandit training of deep network policies

# POEM Hypothesis Space

Hypothesis Space: Stochastic policies

$$\pi_w(y|x) = \frac{1}{Z(x)}\exp\big(w \cdot \Phi(x,y)\big)$$

with

- $w$: parameter vector to be learned
- $\Phi(x,y)$: joint feature map between input and output
- Z(x): partition function

Note: same form as CRF or Structural SVM

# POEM Learning Method

- Policy Optimizer for Exponential Models (POEM)
  - Data: $S = \big((x_1, y_1, \delta_1, p_1), \dots, (x_n, y_n, \delta_n, p_n)\big)$
  - Hypothesis space: $\pi_w(y|x) = \exp\big(w \cdot \phi(x, y)\big)/Z(x)$
  - Training objective: Let $z_i(w) = \pi_w(y_i|x_i)\delta_i/p_i$

$$w = \underset{w \in \Re^N}{\operatorname{argmin}} \left[ \frac{1}{n}\sum_{i=1}^{n} z_i(w) + \lambda_1 \sqrt{\left(\frac{1}{n}\sum_{i=1}^{n} z_i(w)^2\right) - \left(\frac{1}{n}\sum_{i=1}^{n} z_i(w)\right)^2} + \lambda_2 ||w||^2 \right]$$
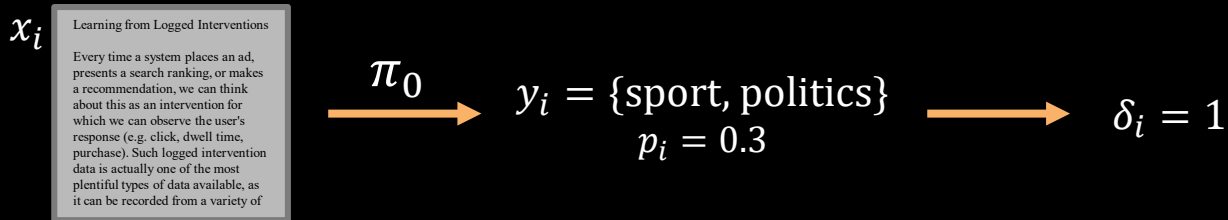
Unbiased Risk Estimator

Variance Control

Capacity Control
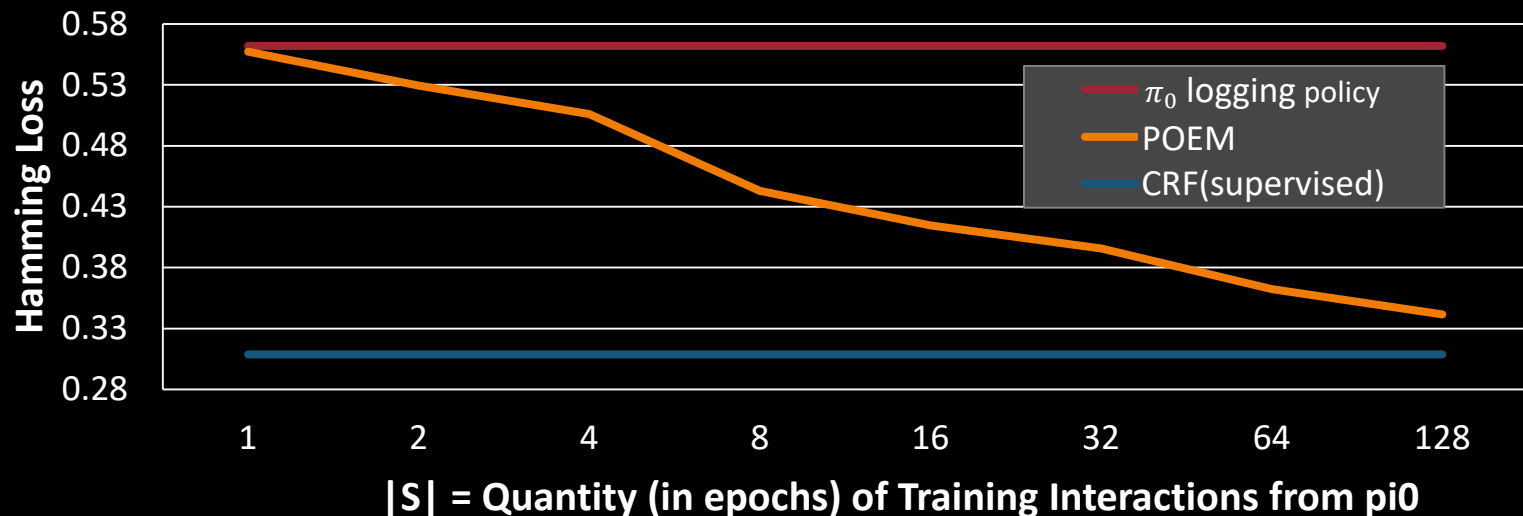
[Swaminathan & Joachims, 2015]

# POEM Experiment
# Multi-Label Text Classification

- Data: $S = \big((x_1, y_1, \delta_1, p_1), \ldots, (x_n, y_n, \delta_n, p_n)\big)$

Reuters LYRL RCV1 (top 4 categories)

$x_i$

> Learning from Logged Interventions
>
> Every time a system places an ad, presents a search ranking, or makes a recommendation, we can think about this as an intervention for which we can observe the user's response (e.g. click, dwell time, purchase). Such logged intervention data is actually one of the most plentiful types of data available, as it can be recorded from a variety of

$\pi_0$ → $y_i = \{\text{sport, politics}\}$ $\quad$ $p_i = 0.3$ → $\delta_i = 1$

- Results: POEM with H isomorphic to CRF with one weight vector per label



(Chart) Hamming Loss vs |S| = Quantity (in epochs) of Training Interactions from pi0

Legend:
- $\pi_0$ logging policy
- POEM
- CRF(supervised)

y-axis: Hamming Loss (0.28, 0.33, 0.38, 0.43, 0.48, 0.53, 0.58)
x-axis: 1, 2, 4, 8, 16, 32, 64, 128

[Swaminathan & Joachims, 2015]

# Does Variance Regularization Improve Generalization?

- IPS: $w = \underset{w \in \Re^N}{\operatorname{argmin}} \left[ \widehat{R}(w) + \lambda_2 ||w||^2 \right]$

- POEM: $w = \underset{w \in \Re^N}{\operatorname{argmin}} \left[ \widehat{R}(w) + \lambda_1 \left( \sqrt{\widehat{Var}(w)/n} \right) + \lambda_2 ||w||^2 \right]$

| Hamming Loss | Scene | Yeast | TMC | LYRL |
|---|---|---|---|---|
| $\pi_0$ | 1.543 | 5.547 | 3.445 | 1.463 |
| IPS | 1.519 | 4.614 | 3.023 | 1.118 |
| POEM | **1.143** | **4.517** | **2.522** | **0.996** |
| # examples | 4*1211 | 4*1500 | 4*21519 | 4*23149 |
| # features | 294 | 103 | 30438 | 47236 |
| # labels | 6 | 14 | 22 | 4 |

# POEM Efficient Training Algorithm

- Training Objective:

$$OPT = \min_{w \in \Re^N} \left[ \frac{1}{n} \sum_{i=1}^{n} z_i(w) + \lambda_1 \sqrt{\left( \frac{1}{n} \sum_{i=1}^{n} z_i(w)^2 \right) - \left( \frac{1}{n} \sum_{i=1}^{n} z_i(w) \right)^2} \right]$$
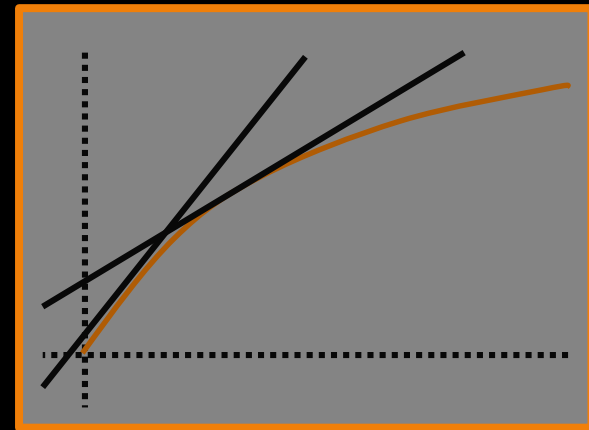
- Idea: First-order Taylor Majorization
  - Majorize $\sqrt{\phantom{x}}$ at current value
  - Majorize $-(\phantom{x})^2$ at current value

$$OPT \leq \min_{w \in \Re^N} \left[ \frac{1}{n} \sum_{i=1}^{n} A_i \, z_i(w) + B_i \, z_i(w)^2 \right]$$

- Algorithm:
  - Majorize objective at current $w_t$
  - Solve majorizing objective via Adagrad to get $w_{t+1}$

[De Leeuw, 1977+] [Groenen et al., 2008] [Swaminathan & Joachims, 2015]

# Counterfactual Risk Minimization

- Theorem [Generalization Error Bound]

$$\text{R}(\pi) \leq \hat{R}(\pi) + O\left(\sqrt{\widehat{Var}(\pi)/n}\right) + O(C)$$

$\rightarrow$ Constructive principle for designing learning algorithms

$$\pi^{crm} = \underset{\pi \in H_i}{\text{argmin}} \; \hat{R}(\pi) + \lambda_1\left(\sqrt{\widehat{Var}(\pi)/n}\right) + \lambda_2 C(H_i)$$

$$\hat{R}(\pi) = \frac{1}{n}\sum_i^n \frac{\pi(y_i|x_i)}{p_i}\delta_i \qquad \widehat{Var}(\pi) = \frac{1}{n}\sum_i^n \left(\frac{\pi(y_i|x_i)}{p_i}\delta_i\right)^2 - \hat{R}(\pi)^2$$

[Swaminathan & Joachims, 2015]

# Propensity Overfitting Problem

- Example
  - Instance Space $X = \{1, \ldots, k\}$
  - Label Space $Y = \{1, \ldots, k\}$
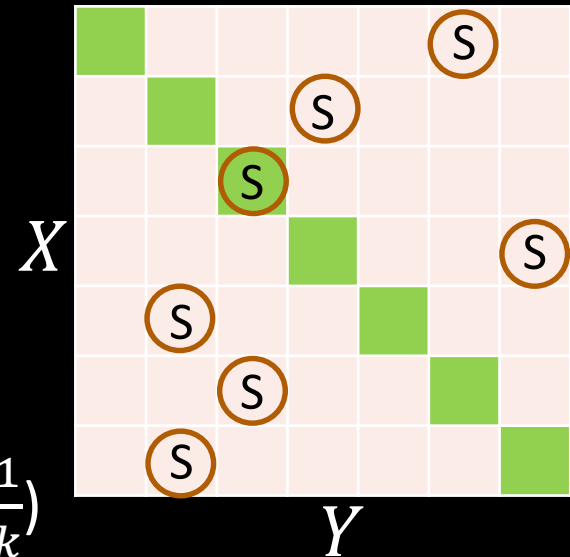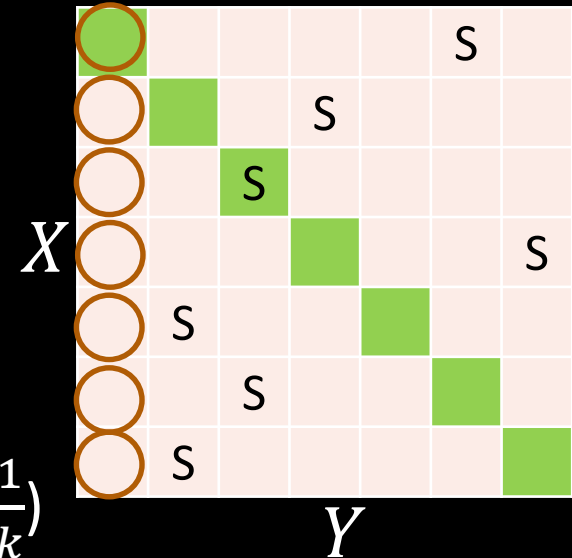  - Loss $\delta(x, y) = \begin{cases} -2 & if\ y == x \\ -1 & otherwise \end{cases}$
  - Training data: uniform x,y sample $(p_i = \frac{1}{k})$
  - Hypothesis space: all deterministic functions
    $\rightarrow \pi_{opt}(x) = x$ with risk $R(\pi_{opt}) = -2$

$$R(\hat{\pi}) = \min_{\pi \in H} \frac{1}{n} \sum_{i}^{n} \frac{\pi(y_i|x_i)}{p_i} \delta_i =$$

$\rightarrow$ Problem 1: Unbounded risk estimate!

# Propensity Overfitting Problem

- Example
  - Instance Space $X = \{1, \dots, k\}$
  - Label Space $Y = \{1, \dots, k\}$
  - Loss $\delta(x, y) = \begin{cases} \cancel{-2}\ 0\ if\ y == x \\ \cancel{-1}\ 1\ otherwise \end{cases}$
  - Training data: uniform x,y sample $\left(p_i = \frac{1}{k}\right)$
  - Hypothesis space: all deterministic functions
    $\rightarrow \pi_{opt}(x) = x$ with risk $R(\pi_{opt}) = \cancel{-2}\ 0$

$$R(\hat{\pi}) = \min_{\pi \in H} \frac{1}{n} \sum_i^n \frac{\pi(y_i | x_i)}{p_i} \delta_i =$$

$\rightarrow$ Problem 2: Lack of equivariance!

# Control Variate

- Idea: Inform estimate when expectation of correlated random variable is known.

  - Estimator:
    $$\hat{R}(\pi) = \frac{1}{n} \sum_{i}^{n} \frac{\pi(y_i|x_i)}{p_i} \delta_i$$

  - Correlated RV with known expectation:
    $$\hat{S}(\pi) = \frac{1}{n} \sum_{i}^{n} \frac{\pi(y_i|x_i)}{p_i}$$

$$E\big[\hat{S}(\pi)\big] = \frac{1}{n} \sum_{i}^{n} \int \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)} \pi_0(y_i|x_i) P(x) dy_i dx_i = 1$$

→ Alternative Risk Estimator: Self-normalized estimator
$$\hat{R}^{SN}(\pi) = \frac{\hat{R}(\pi)}{\hat{S}(\pi)}$$

[Hesterberg, 1995] [Swaminathan & Joachims, 2015]

# SNIPS Learning Objective

- Method:
  - Data: $S = \big((x_1, y_1, \delta_1, p_1), \ldots, (x_n, y_n, \delta_n, p_n)\big)$
  - Hypothesis space: $\pi_w(y|x) = \exp\big(w \cdot \phi(x, y)\big)/Z(x)$
  - Training objective:

$$w = \underset{w \in \Re^N}{\operatorname{argmin}} \left[ \hat{R}^{SNIPS}(w) + \lambda_1 \sqrt{\widehat{Var}\big(\hat{R}^{SNIPS}(w)\big)} + \lambda_2 ||w||^2 \right]$$

Self-Normalized Risk Estimator

Variance Control

Capacity Control

[Swaminathan & Joachims, 2015]

# How well does NormPOEM generalize?

| Hamming Loss | Scene | Yeast | TMC | LYRL |
|---|---:|---:|---:|---:|
| $\pi_0$ | 1.511 | 5.577 | 3.442 | 1.459 |
| POEM (IPS) | 1.200 | 4.520 | 2.152 | 0.914 |
| POEM (SNIPS) | **1.045** | **3.876** | **2.072** | **0.799** |
| # examples | 4*1211 | 4*1500 | 4*21519 | 4*23149 |
| # features | 294 | 103 | 30438 | 47236 |
| # labels | 6 | 14 | 22 | 4 |

# Outline of Lecture

- Batch Learning from Bandit Feedback (BLBF)
$$S = \big( (x_1, y_1, \delta_1), \ldots, (x_n, y_n, \delta_n) \big)$$
  → Find new policy $\pi$ that selects $y$ with better $\delta$

- Learning Principle for BLBF
  - Hypothesis Space, Risk, Empirical Risk, and Overfitting
  - Learning Principle: Counterfactual Risk Minimization

- Learning Algorithms for BLBF
  - POEM: Bandit training of CRF policies for structured outputs
  - BanditNet: Bandit training of deep network policies

# BanditNet: Hypothesis Space

Hypothesis Space: Stochastic policies

$$\pi_w(y|x) = \frac{1}{Z(x)} \exp(DeepNet(x, y|w))$$

with

- $w$: parameter tensors to be learned
- Z(x): partition function

Note: same form as Deep Net with softmax output

[Joachims et al., 2017]

# BanditNet: Learning Method

- Method:
  - Data: $S = \big((x_1, y_1, \delta_1, p_1), \ldots, (x_n, y_n, \delta_n, p_n)\big)$
  - Hypotheses: $\pi_w(y|x) = \exp\big(DeepNet(x|w)\big)/Z(x)$
  - Training objective:

$$w = \underset{w \in \mathfrak{R}^N}{\operatorname{argmin}} \left[ \hat{R}^{SNIPS}(w) + \lambda_1 \sqrt{\widehat{Var}\big(\hat{R}^{SNIPS}(w)\big)} + \lambda_2 ||w||^2 \right]$$

Self-Normalized Risk Estimator

Variance Control

Capacity Control

[Joachims et al., 2017]

# BanditNet: Learning Method

- Method:
  - Data: $S = \left( (x_1, y_1, \delta_1, p_1), \dots, (x_n, y_n, \delta_n, p_n) \right)$
  - Representation: Deep Network Policies

$$\pi_w(y|x) = \frac{1}{Z(x,w)} \exp\left( DeepNet(y|x,w) \right)$$

  - SNIPS Training Objective:

$$w = \underset{w \in \Re^N}{\operatorname{argmin}} \left[ \hat{R}_{SNIPS}(w) + \lambda ||w||^2 \right]$$

$$= \underset{w \in \Re^N}{\operatorname{argmin}} \left[ \frac{1}{\sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i}} \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i} \delta_i + \lambda ||w||^2 \right]$$

# Optimization via SGD

- Problem: SNIPS objective not suitable for SGD
- Step 1: Discretize over values in denominator

$$\hat{w} = \underset{S_j}{\mathrm{argmin}} \left[ \underset{w}{\mathrm{argmin}} \left[ \frac{1}{S_j} \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i} \delta_i \right] \text{ subject to } \frac{1}{n} \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i} = S_j \right]$$

- Step 2: View as series of constrained OP

$$\hat{w}_j = \underset{w}{\mathrm{argmin}} \left[ \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i} \delta_i \right] \text{ subject to } \frac{1}{n} \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i} = S_j$$

- Step 3: Eliminate constraint via Lagrangian

$$\hat{w}_j = \underset{w}{\mathrm{argmin}} \underset{\lambda}{\mathrm{max}} \left[ \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i} (\delta_i - \lambda) + \lambda S_j \right]$$

# Optimization via SGD

- Step 4: Search grid over $\lambda$ instead of $S_j$
  - Hard: Given $S_j$, find $\lambda_j$.
  - Easy: Given $\lambda_j$, find $S_j$.

→Solve

$$\widehat{w}_j = \underset{w}{\operatorname{argmin}} \left[ \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i} (\delta_i - \lambda_j) + \lambda_j S_j \right]$$

→Compute

$$S_j = \frac{1}{n} \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i}$$

# BanditNet: Training Algorithm

- Given:
  - Data: $S = \big((x_1, y_1, \delta_1, p_1), \ldots, (x_n, y_n, \delta_n, p_n)\big)$
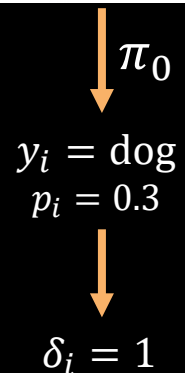  - Lagrange Multipliers: $\lambda_j \in \{\lambda_1, \ldots, \lambda_k\}$
- Compute:
  - For each $\lambda_j$ solve:
  $$\widehat{w}_j = \operatorname*{argmin}_{\mathrm{w}} \left[ \sum_{i=1}^{n} \frac{\pi_w(y_i|x_i)}{p_i} (\delta_i - \lambda_j) \right]$$

  - For each $\widehat{\mathrm{w}}_j$ compute:
  $$S_j = \frac{1}{n} \sum_{i=1}^{n} \frac{\pi_{\widehat{w}_j}(y_i|x_i)}{p_i}$$
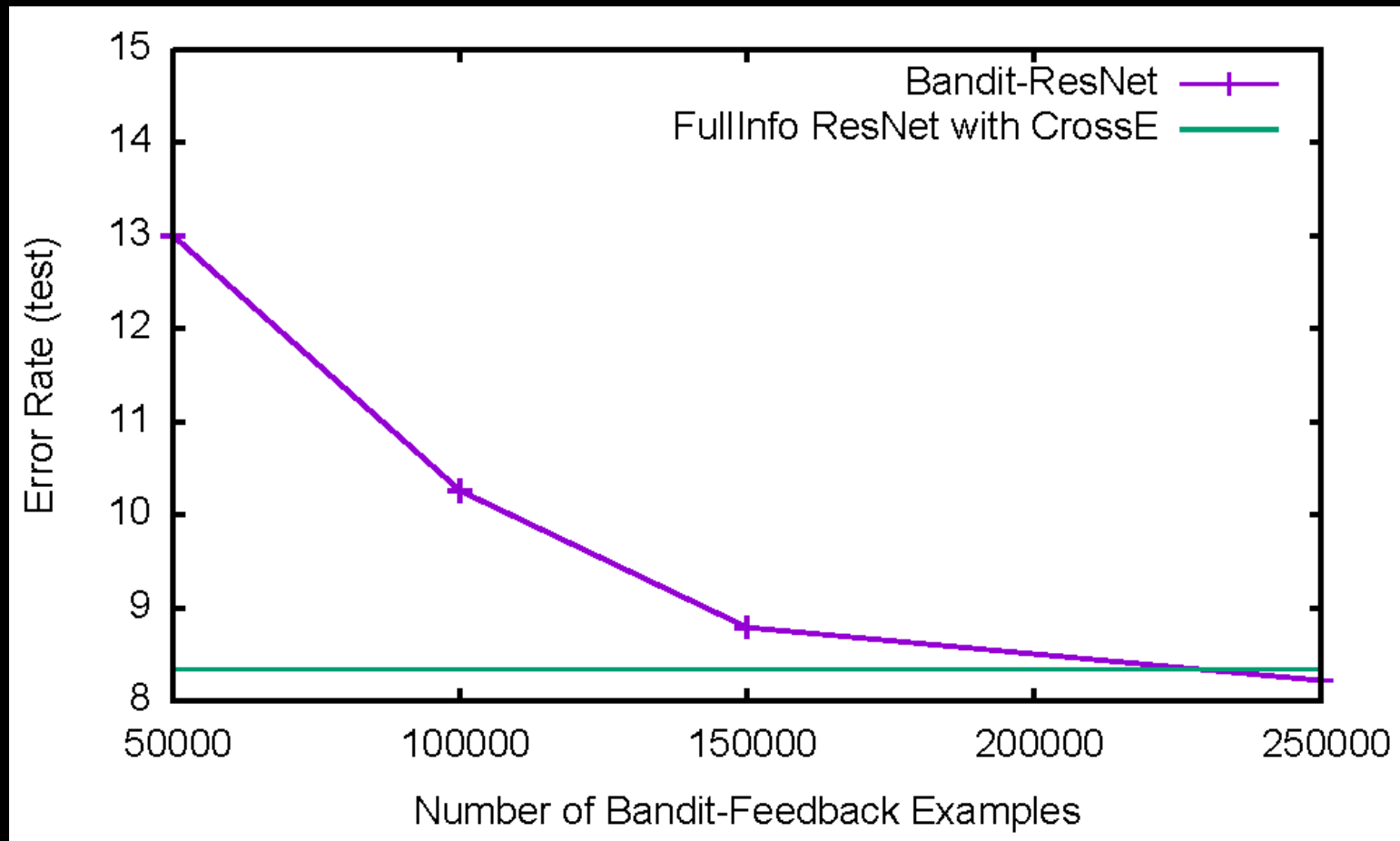
  - Find overall $\widehat{w}$:
  $$\widehat{w} = \operatorname*{argmin}_{\widehat{w}_j, S_j} \left[ \frac{1}{S_j} \sum_{i=1}^{n} \frac{\pi_{\widehat{w}_j}(y_i|x_i)}{p_i} \delta_i \right]$$

# Object Recognition: Data and Setup

- Data: CIFAR-10 (fully labeled)
  $\rightarrow S^* = \big((x_1, y_1^*), \dots, (x_m, y_m^*)\big)$

- Bandit feedback generation:
  - Draw image $x_i$
  - Use logging policy $\pi_0(Y|x_i)$ to predict $y_i$
    - Record propensity $\pi_0(Y = y_i|x_i)$
  - Observe loss $\delta_i = [y_i \neq y_i^*]$
  $\rightarrow S = \big((x_1, y_1, \delta_1, p_1), \dots, (x_n, y_n, \delta_n, p_n)\big)$

- Network architecture: ResNet20 [He et al., 2016]

$\pi_0$

$y_i = \text{dog}$
$p_i = 0.3$

$\delta_i = 1$

[Beygelzimer & Langford, 2009] [Joachims et al., 2017]
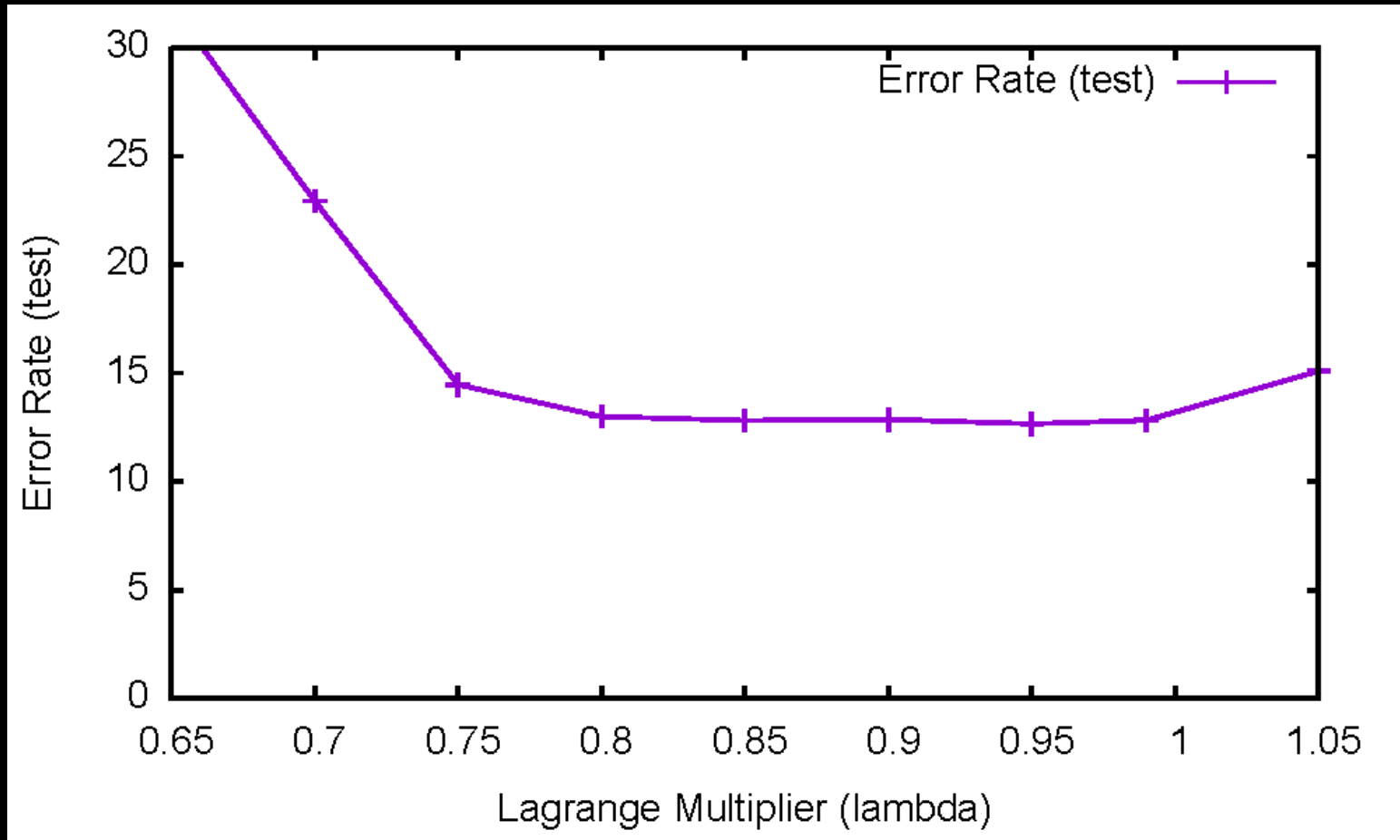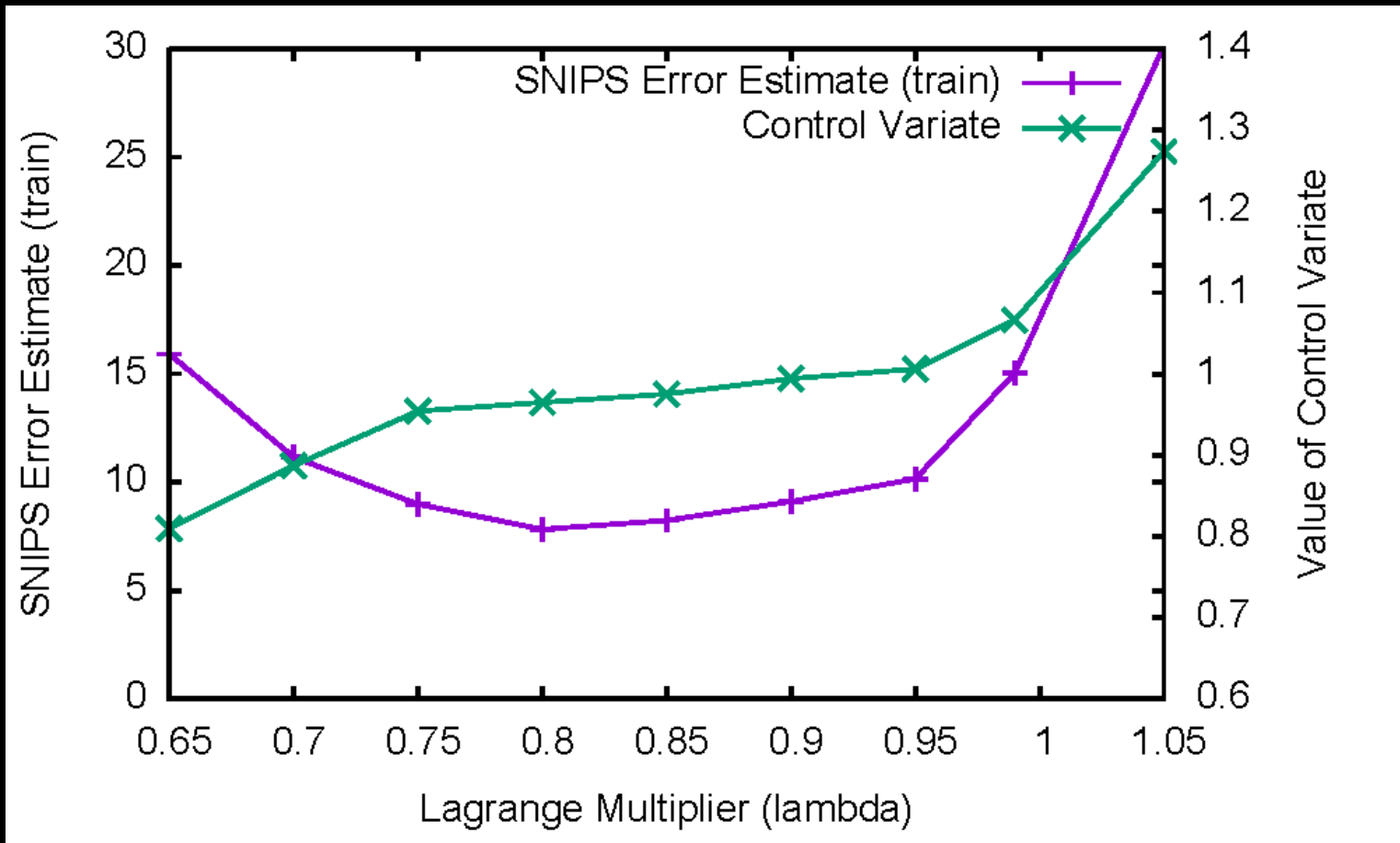
# Bandit Feedback vs. Test Error



Logging Policy $\pi_0$:              49% error rate
Bandit-ResNet with naïve IPS:   >49% error rate

[Joachims et al., 2017]

# Lagrange Multiplier vs. Test Error



Large basin of optimality far away from naïve IPS.

# Analysis of SNIPS Estimate



Control variate responds to the Lagrange multiplier monotonically.
SNIPS training error resembles test error.

# Conclusions and Future

- Batch Learning from Bandit Feedback
  - Feedback for only presented action
    $$S = \big((x_1, y_1, \delta_1, p_1), \dots, (x_n, y_n, \delta_n, p_n)\big)$$
  - Goal: Find new system $\pi$ that selects $y$ with better $\delta$
  - Learning Principle for BLBF: Counterfactual Risk Minimization
- Learning from Logged Interventions: BLBF and Beyond
  - POEM: [Swaminathan & Joachims, 2015c]
  - NormPOEM: [Swaminathan & Joachims, 2015c]
  - BanditNet: [Joachims et al., 2018]
  - SVM PropRank [Joachims et al., 2017a]
  - DeepPropDCG: [Agarwal et al., 2018]
  - Unbiased Matrix Factorization: [Schnabel et al. 2016]
- Future Research
  - Other learning algorithms? Other partial-information settings?
  - How to handle new bias-variance trade-off in risk estimators?
  - Applications
- Software, Papers, SIGIR Tutorial, Data: www.joachims.org