

Counterfactual Model for Learning 2

CS6780 – Advanced Machine Learning
Spring 2019

Thorsten Joachims
Cornell University

Reading:

G. Imbens, D. Rubin, Causal Inference for Statistics ..., 2015. Chapters 1,3,12.

From Evaluation to Learning

Setting: Batch Learning from Bandit Feedback (BLBF)

- Naïve “Model the World” Learning:

- Learn: $\hat{\delta}: x \times y \rightarrow \mathfrak{R}$

- Derive Policy:

$$\pi(y|x) = \operatorname{argmin}_{y'} [\hat{\delta}(x, y')]$$

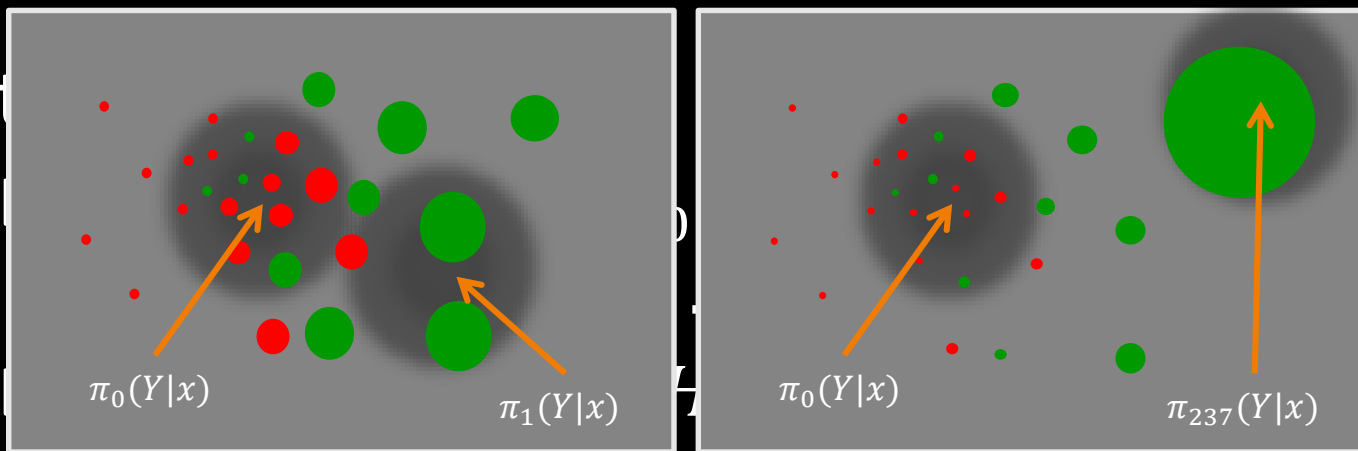
- • Naïve “Model the Bias” Learning:

- Find policy that optimizes IPS training error

$$\pi = \operatorname{argmin}_{\pi'} \left[\sum_i \frac{\pi'(y_i|x_i)}{\pi_0(y_i|x_i)} \delta_i \right]$$

Partial-Information ERM

- Set



- Training

$$\hat{\pi} := \operatorname{argmax}_{\pi \in H} \left[\sum_{i=1}^n \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)} \delta_i \right]$$

Learning Theory for BLBF

Theorem [Generalization Error Bound]

For any policy space H with capacity C , and for all $\pi \in H$ with probability $1 - \eta$

$$U(\pi) \geq \hat{U}(\pi) - O\left(\sqrt{\frac{\widehat{\text{var}}(\hat{U}(\pi))}{n}}\right) - O(C)$$

Unbiased
Estimator

Variance
Overfitting

Capacity
Overfitting

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

Counterfactual Risk Minimization

Constructive principle for learning algorithms

→ Maximize learning theoretical bound

$$\pi^{crm} = \operatorname{argmax}_{\pi \in H} \left[\hat{U}(\pi) - \lambda_1 \left(\sqrt{\widehat{\operatorname{Var}}(\hat{U}(\pi)) / n} \right) - \lambda_2 C(H) \right]$$

Training
Error/Utility (IPS)

Variance
Regularization

Capacity
Regularization

POEM: Policy Space

Policy space

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

with

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

Note: same form as CRF or Structural SVM

POEM: Learning Method

Policy Optimizer for Exponential Models (POEM)

- Data: $S = ((x_1, y_1, \delta_1, p_1), \dots, (x_n, y_n, \delta_n, p_n))$
- Policy space: $\pi_w(y|x) = \exp(w \cdot \phi(x, y)) / Z(x)$

$$w = \operatorname{argmax}_{w \in \mathbb{R}^N} \left[\hat{U}(\pi_w) - \lambda_1 \left(\sqrt{\widehat{\operatorname{Var}}(\hat{U}(\pi_w))} \right) - \lambda_2 \|w\|^2 \right]$$

IPS Estimator

Variance
Regularization

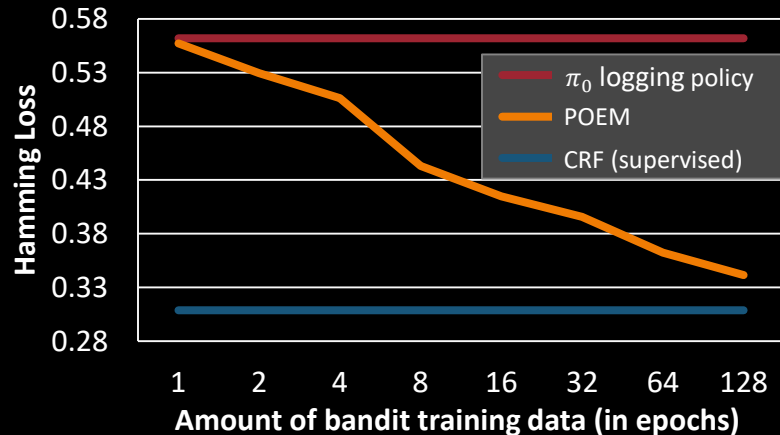
Capacity
Regularization

POEM: Text Classification

Data: Reuters Text Classification

- $S^* = ((x_1, y_1^*), \dots, (x_m, y_m^*))$
- Label vectors $y^* = (y^1, y^2, y^3, y^4)$

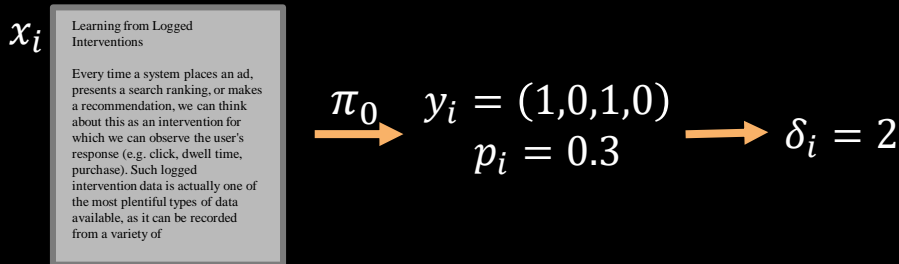
Results:



Bandit feedback generation:

- Draw document x_i
- Pick y_i via logging policy $\pi_0(Y|x_i)$
- Observe loss $\delta_i = \text{Hamming}(y_i, y_i^*)$

$$\rightarrow S = ((x_1, y_1, \delta_1, p_1), \dots, (x_n, y_n, \delta_n, p_n))$$



BanditNet: Policy Space

Policy space

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

with

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

Note: same form as Deep Net with softmax output

BanditNet: Learning Method

- Deep networks with bandit feedback (BanditNet):
 - Data: $S = ((x_1, y_1, \delta_1, p_1), \dots, (x_n, y_n, \delta_n, p_n))$
 - Hypotheses: $\pi_w(y|x) = \exp(\text{DeepNet}(x|w))/Z(x)$

$$w = \operatorname{argmax}_{w \in \mathcal{R}^N} \left[\hat{U}(\pi_w) - \lambda_1 \left(\sqrt{\widehat{\text{Var}}(\hat{U}(\pi_w))} \right) - \lambda_2 \|w\|^2 \right]$$

Self-Normalized
IPS Estimator

Variance
Regularization

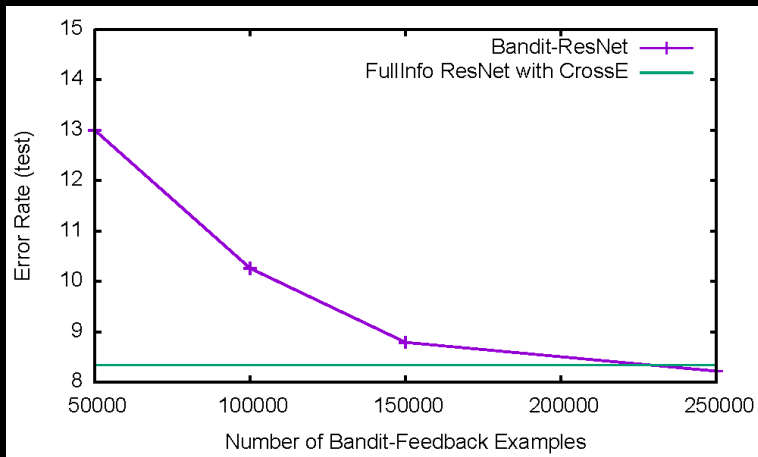
Capacity
Regularization

BanditNet: Object Recognition

- Data: CIFAR-10

- $S^* = ((x_1, y_1^*), \dots, (x_m, y_m^*))$
- ResNet20 [He et al., 2016]

- Results



- Bandit feedback generation:

- Draw image x_i
- Pick y_i via logging policy $\pi_0(Y|x_i)$
- Observe loss $\delta_i = [y_i \neq y_i^*]$

$$\rightarrow S = ((x_1, y_1, \delta_1, p_1), \dots, (x_n, y_n, \delta_n, p_n))$$



$$\begin{array}{c} \pi_0 \\ \rightarrow \end{array} \begin{array}{c} y_i = \text{dog} \\ p_i = 0.3 \end{array} \rightarrow \delta_i = 1$$