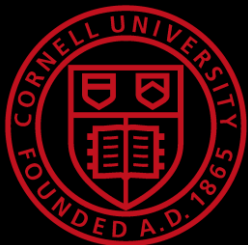# Counterfactual Evaluation and Learning

Adith Swaminathan, Thorsten Joachims

Department of Computer Science & Department of Information Science

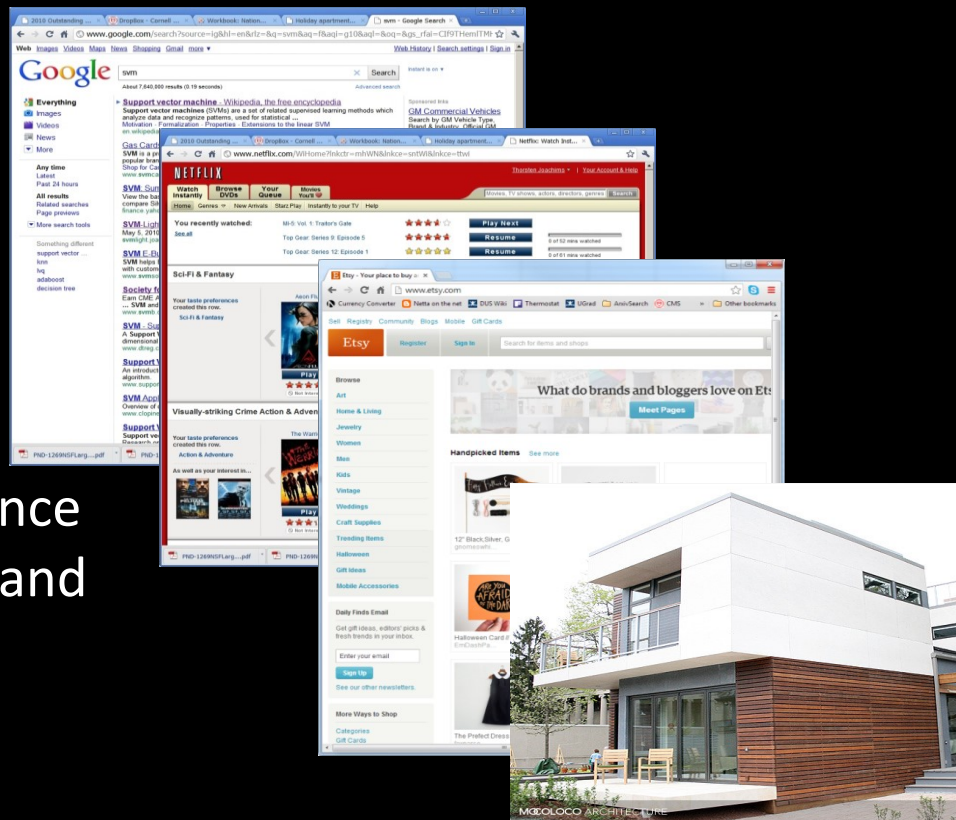Cornell University
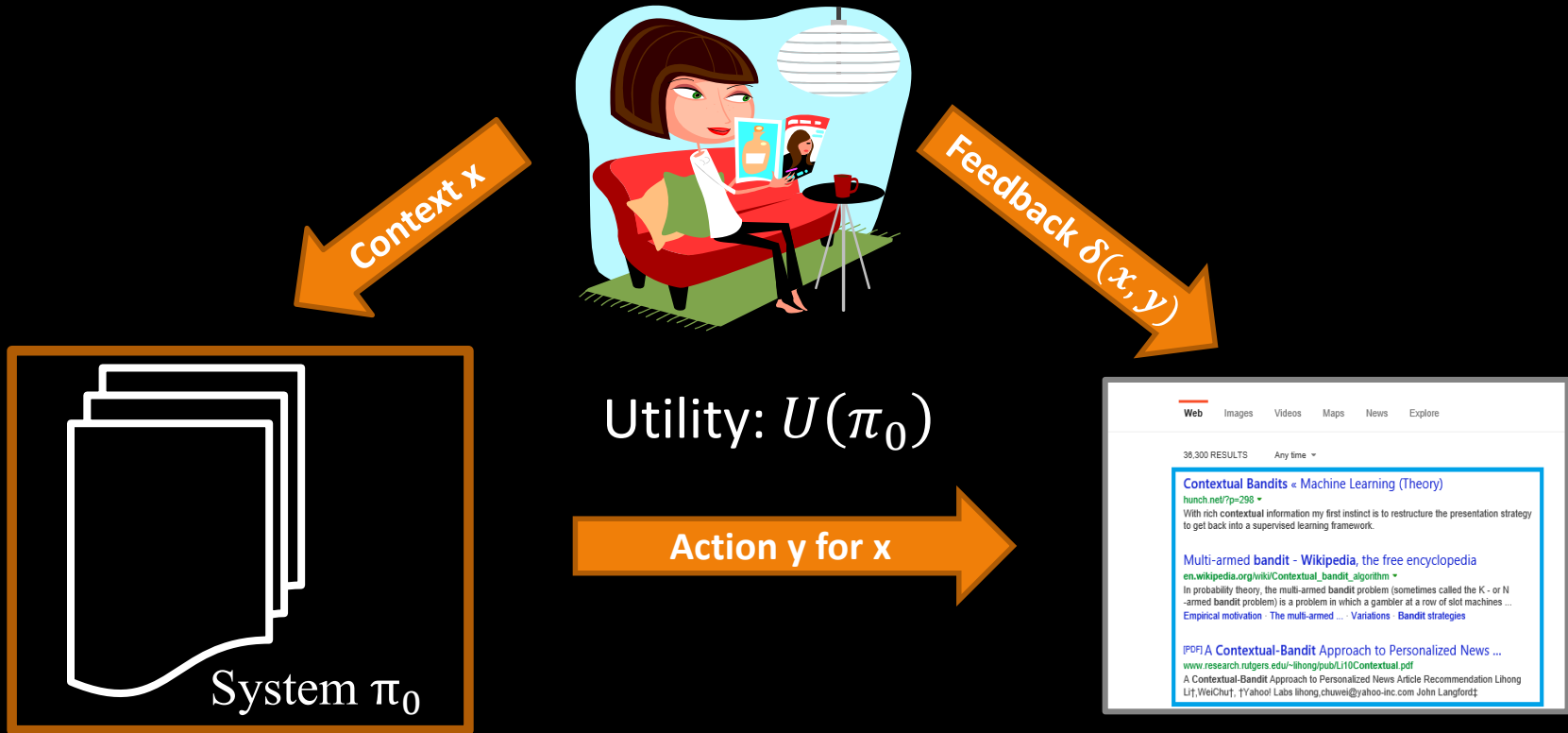
Website: http://www.cs.cornell.edu/~adith/CfactSIGIR2016/

# User Interactive Systems

Examples
- Search engines
- Entertainment media
- E-commerce
- Smart homes, robots, etc.
→ Logs of User Behavior for
- Evaluating system performance
- Learning improved systems and gathering knowledge
- Personalization

# Interactive System Schematic



Context x

Feedback $\delta(x, y)$

Utility: $U(\pi_0)$

Action y for x

System $\pi_0$

# Ad Placement

- ## Context $x$:
  - User and page
- ## Action $y$:
  - Ad that is placed
- ## Feedback $\delta(x, y)$:
  - Click / no-click

# News Recommender

- ## Context $x$:
  - User

- ## Action $y$:
  - Portfolio of newsarticles

- ## Feedback $\delta(x, y)$:
  - Reading time in minutes

# Search Engine

- ## Context $x$:
  - Query
- ## Action $y$:
  - Ranking
- ## Feedback $\delta(x, y)$:
  - win/loss against baseline in interleaving

# Log Data from Interactive Systems

- Data

  context     $\pi_0$ action     reward / loss

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

  → Partial Information (aka "Contextual Bandit")
    Feedback

- Properties
  - Contexts $x_i$ drawn i.i.d. from unknown $P(X)$
  - Actions $y_i$ selected by existing system $\pi_0 : X \to Y$
  - Feedback $\delta_i$ from unknown function $\delta : X \times Y \to \Re$

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

# Goals for this Tutorial

- Use interaction log data
$$S = \big((x_1, y_1, \delta_1), \dots, (x_n, y_n, \delta_n)\big)$$
  for
  - Evaluation:
    - Estimate online measures of some system $\pi$ offline.
    - System $\pi$ is typically different from $\pi_0$ that generated log.
  - Learning:
    - Find new system $\pi$ that improves performance over $\pi_0$.
    - Do not rely on interactive experiments like in online learning.

SIGIR 2016 Tutorial
Counterfactual Evaluation and Learning

# PART 1: EVALUATION

# Evaluation: Outline

- Evaluating Online Metrics Offline
  - A/B Testing (on-policy) → Counterfactual estimation from logs (off-policy)
- Approach 1: "Model the world"
  - Estimation via reward prediction
- Approach 2: "Model the bias"
  - Counterfactual Model
  - Inverse propensity scoring (IPS) estimator
- Advanced Estimators
  - Self-normalized IPS estimator
  - Doubly robust estimator
  - Slates estimator
- Case Studies
- Summary & Demonstration with code samples

# Online Performance Metrics

Example metrics
- CTR
- Revenue
- Time-to-success
- Interleaving
- Etc.

→ Correct choice depends on application and is not the focus of this tutorial.

This tutorial:

Metric encoded as $\delta(x, y)$     [click/payoff/time for (x,y) pair]

# System

- Definition [Deterministic Policy]: Function

$$y = \pi(x)$$

  that picks action $y$ for context $x$.

- Definition [Stochastic Policy]: Distribution

$$\pi(y|x)$$

  that samples action $y$ given context $x$

# System Performance

Definition [Utility of Policy]:

The expected reward / utility $U(\pi)$ of policy $\pi$ is

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



$Y|x_i$

$\pi(Y|x_i)$

e.g. reading time of user x for portfolio y

$\cdots$

$Y|x_j$

$\pi(Y|x_j)$

# Online Evaluation: A/B Testing

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

$$\hat{U}(\pi_0) = \frac{1}{n} \sum_{i=1}^{n} \delta_i$$

→  A/B Testing

Deploy $\pi_1$: Draw $x \sim P(X)$, predict $y \sim \pi_1(Y|x)$, get $\delta(x, y)$

Deploy $\pi_2$: Draw $x \sim P(X)$, predict $y \sim \pi_2(Y|x)$, get $\delta(x, y)$

$\vdots$

Deploy $\pi_{|H|}$: Draw $x \sim P(X)$, predict $y \sim \pi_{|H|}(Y|x)$, get $\delta(x, y)$

# Pros and Cons of A/B Testing

- Pro
  - User centric measure
  - No need for manual ratings
  - No user/expert mismatch
- Cons
  - Requires interactive experimental control
  - Risk of fielding a bad or buggy $\pi_i$
  - Number of A/B Tests limited
  - Long turnaround time

# 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 $\pi_0$ $\rightarrow$

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

# Evaluation: Outline

- Evaluating Online Metrics Offline
  - A/B Testing (on-policy) → Counterfactual estimation from logs (off-policy)
- Approach 1: "Model the world"
  - Estimation via reward prediction
- Approach 2: "Model the bias"
  - Counterfactual Model
  - Inverse propensity scoring (IPS) estimator
- Advanced Estimators
  - Self-normalized IPS estimator
  - Doubly robust estimator
  - Slates estimator
- Case Studies
- Summary & Demonstration with code samples

# Approach 1: Reward Predictor

- Idea:
  - Use $S = \big((x_1, y_1, \delta_1), \dots, (x_n, y_n, \delta_n)\big)$ from $\pi_0$ to estimate reward predictor $\hat{\delta}(x, y)$



- Deterministic $\pi$: Simulated A/B Testing with predicted $\hat{\delta}(x, y)$
  - For actions $y_i' = \pi(x_i)$ from new policy $\pi$, generate predicted log
  $$S' = \left(\left(x_1, y_1', \hat{\delta}(x_1, y_1')\right), \dots, \left(x_n, y_n', \hat{\delta}(x_n, y_n')\right)\right)$$
  - Estimate performace of $\pi$ via $\widehat{U}_{rp}(\pi) = \frac{1}{n} \sum_{i=1}^{n} \hat{\delta}(x_i, y_i')$

- Stochastic $\pi$: $\widehat{U}_{rp}(\pi) = \frac{1}{n} \sum_{i=1}^{n} \sum_y \hat{\delta}(x_i, y) \, \pi(y|x_i)$

# Regression for Reward Prediction

## Learn $\hat{\delta}: x \times y \rightarrow \Re$

1. Represent via features $\Psi(x, y)$
2. Learn regression based on $\Psi(x, y)$ from $S$ collected under $\pi_0$
3. Predict $\hat{\delta}(x, y')$ for $y' = \pi(x)$ of new policy $\pi$

# News Recommender: Exp Setup

- Context x: User profile

- Action y: Ranking
  – Pick from 7 candidates to place into 3 slots

- Reward $\delta$: "Revenue"
  – Complicated hidden function



- Logging policy $\pi_0$: Non-uniform randomized logging system
  – Placket-Luce "explore around current production ranker" (see case study)

# News Recommender: Results



REVENUE 3 slots, 7 candidates

Avg. Error over 10 trials 3 slots, 7 candidates

RP is inaccurate even with more training and logged data

# Problems of Reward Predictor

- Modeling bias
  - choice of features and model
- Selection bias
  - $\pi_0$'s actions are over-represented

$$\rightarrow \hat{U}_{rp}(\pi) = \frac{1}{n}\sum_i \hat{\delta}(x_i, \pi(x_i))$$

Can be unreliable and biased

# Evaluation: Outline

- Evaluating Online Metrics Offline
  - A/B Testing (on-policy) → Counterfactual estimation from logs (off-policy)
- Approach 1: "Model the world"
  - Estimation via reward prediction
- Approach 2: "Model the bias"
  - Counterfactual Model
  - Inverse propensity scoring (IPS) estimator
- Advanced Estimators
  - Self-normalized IPS estimator
  - Doubly robust estimator
  - Slates estimator
- Case Studies
- Summary & Demonstration with code samples

# Approach "Model the Bias"

- Idea:

  Fix the mismatch between the distribution $\pi_0(Y|x)$ that generated the data and the distribution $\pi(Y|x)$ we aim to evaluate.

$$\mathrm{U}(\pi_0^{\pi}) = \int \int \delta(x, y) \pi_0(y|x)^{\pi(y|x)} P(x) \, dx \, dy$$

# Counterfactual Model

- Example: Treating Heart Attacks
  - Treatments: $Y$
    - Bypass / Stent / Drugs
  - Chosen treatment for patient $x_i$: $y_i$
  - Outcomes: $\delta_i$
    - 5-year survival: 0 / 1
  - Which treatment is best?

$$
\text{Patients } x_i \in \{1, \ldots, n\}
\begin{array}{ccc}
\text{Bypass} & \text{Stent} & \text{Drugs} \\
\end{array}
\begin{bmatrix}
0 & & \\
& 1 & \\
& & 1 \\
& 0 & \\
& & 1 \\
1 & & \\
& & 1 \\
& & 0 \\
0 & & \\
& 1 & \\
1 & &
\end{bmatrix}
$$

# Counterfactual Model

- Example: ~~Treating Heart Attacks~~
  - Treatments: $Y$
    - ~~Bypass / Stent / Drugs~~ Pos 1 / Pos 2/ Pos 3
  - Chosen treatment for patient $x_i$: $y_i$
  - Outcomes: $\delta_i$
    - ~~5-year survival: 0 / 1~~ Click / no Click on SERP
  - Which treatment is best?
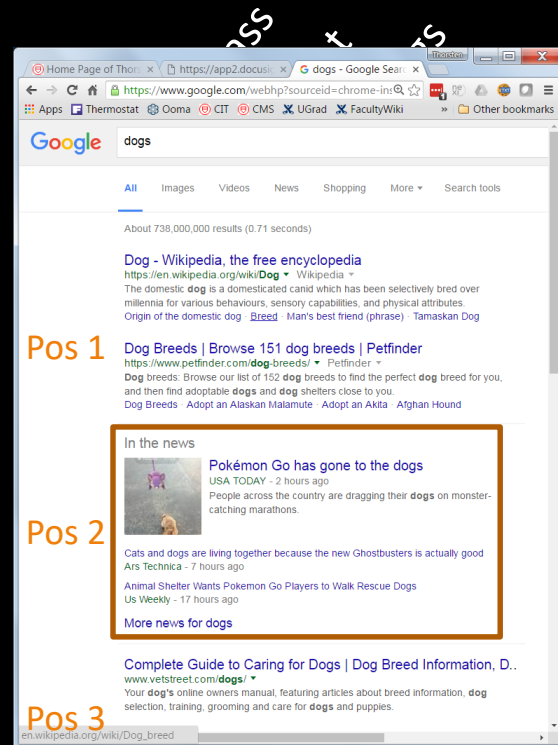
# Counterfactual Model

- Example: Treating Heart Attacks
  - Treatments: $Y$
    - Bypass / Stent / Drugs
  - Chosen treatment for patient $\mathrm{x}_i$: $\mathrm{y}_i$
  - Outcomes: $\delta_i$
    - 5-year survival: 0 / 1
  - Which treatment is best?
    - Everybody Drugs
    - Everybody Stent
    - Everybody Bypass
    - → Drugs 3/4, Stent 2/3, Bypass 2/4 – really?

Patients $\mathrm{x}_i, i \in \{1, \ldots, n\}$

$$
\begin{array}{ccc}
\text{Bypass} & \text{Stent} & \text{Drugs} \\
0 & & \\
 & 1 & \\
 & & 1 \\
 & 0 & \\
 & & 1 \\
1 & & \\
 & & 1 \\
 & & 0 \\
0 & & \\
 & 1 & \\
1 & &
\end{array}
$$

# Treatment Effects

- Average Treatment Effect of Treatment $y$
  - $\mathsf{U}(y) = \frac{1}{n}\sum_i \delta(x_i, y)$
- Example
  - $\mathsf{U}(bypass) = \frac{5}{11}$
  - $\mathsf{U}(stent) = \frac{7}{11}$
  - $\mathsf{U}(drugs) = \frac{4}{11}$

Bypass    Stent    Drugs

Factual Outcome

Counterfactual Outcomes

Patients

$$
\begin{bmatrix}
0 & 1 & 0 \\
1 & 1 & 0 \\
0 & 0 & 1 \\
0 & 0 & 0 \\
0 & 1 & 1 \\
1 & 0 & 0 \\
1 & 0 & 1 \\
0 & 1 & 0 \\
0 & 1 & 0 \\
1 & 1 & 0 \\
1 & 1 & 1
\end{bmatrix}
$$

# Assignment Mechanism

- Probabilistic Treatment Assignment
  - For patient i: $\pi_0(Y_i = y|x_i)$
  - Selection Bias
- Inverse Propensity Score Estimator
  - $\widehat{U}_{ips}(y) = \dfrac{1}{n}\sum_i \dfrac{\mathbb{I}\{y_i = y\}}{p_i}\delta(x_i, y_i)$
  - Propensity: $p_i = \pi_0(Y_i = y_i|x_i)$
  - Unbiased: $E\left[\widehat{U}(y)\right] = U(y)$,
    if $\pi_0(Y_i = y|x_i) > 0$ for all $i$
- Example
  - $\widehat{U}(drugs) = \dfrac{1}{11}\left(\dfrac{1}{0.8} + \dfrac{1}{0.7} + \dfrac{1}{0.8} + \dfrac{0}{0.1}\right)$
    $= 0.36 < 0.75$

$\pi_0(Y_i = y|x_i)$

| | | | Bypass | Stent | Drugs |
|---|---|---|---|---|---|
| 0.3 | 0.6 | 0.1 | 0 | 1 | 0 |
| 0.5 | 0.4 | 0.1 | 1 | 1 | 0 |
| 0.1 | 0.1 | 0.8 | 0 | 0 | 1 |
| 0.6 | 0.3 | 0.1 | 0 | 0 | 0 |
| 0.2 | 0.5 | 0.7 | 0 | 1 | 1 |
| 0.7 | 0.2 | 0.1 | 1 | 0 | 0 |
| 0.1 | 0.1 | 0.8 | 1 | 0 | 1 |
| 0.1 | 0.8 | 0.1 | 0 | 1 | 0 |
| 0.3 | 0.3 | 0.4 | 0 | 1 | 0 |
| 0.3 | 0.6 | 0.1 | 1 | 1 | 0 |
| 0.4 | 0.4 | 0.2 | 1 | 1 | 1 |

Patients

# Experimental vs Observational

- Controlled Experiment
  - Assignment Mechanism under our control
  - Propensities $p_i = \pi_0(Y_i = y_i | x_i)$ are known by design
  - Requirement: $\forall y \colon \pi_0(Y_i = y | x_i) > 0$ (probabilistic)
- Observational Study
  - Assignment Mechanism not under our control
  - Propensities $p_i$ need to be estimated
  - Estimate $\hat{\pi}_0(Y_i | z_i) = \pi_0(Y_i | x_i)$ based on features $z_i$
  - Requirement: $\hat{\pi}_0(Y_i | z_i) = \hat{\pi}_0(Y_i | \delta_i, z_i)$ (unconfounded)

# Conditional Treatment Policies

- Policy (deterministic)
  - Context $x_i$ describing patient
  - Pick treatment $y_i$ based on $x_i$: $y_i = \pi(x_i)$
  - Example policy:
    - $\pi(A) = drugs, \pi(B) = stent, \pi(C) = bypass$
- Average Treatment Effect
  - $U(\pi) = \frac{1}{n}\sum_i \delta(x_i, \pi(x_i))$
- IPS Estimator
  - $\widehat{U}_{ips}(\pi) = \frac{1}{n}\sum_i \frac{\mathbb{I}\{y_i = \pi(x_i)\}}{p_i}\delta(x_i, y_i)$

# Stochastic Treatment Policies

- Policy (stochastic)
  - Context $x_i$ describing patient
  - Pick treatment $y$ based on $x_i$: $\pi(Y|x_i)$
- Note
  - Assignment Mechanism is a stochastic policy as well!
- Average Treatment Effect
  - $U(\pi) = \frac{1}{n} \sum_i \sum_y \delta(x_i, y) \pi(y|x_i)$
- IPS Estimator
  - $\hat{U}(\pi) = \frac{1}{n} \sum_i \frac{\pi(y_i|x_i)}{p_i} \delta(x_i, y_i)$

# Counterfactual Model = Logs



| | |
|---|---|
| Context $x_i$ | |
| Treatment $y_i$ | |
| Outcome $\delta_i$ | |
| Propensities $p_i$ | |
| New Policy $\pi$ | |
| T-effect $U(\pi)$ | Average quality of new policy. |

Recorded in Log

# Evaluation: Outline

- Evaluating Online Metrics Offline
  - A/B Testing (on-policy) → Counterfactual estimation from logs (off-policy)
- Approach 1: "Model the world"
  - Estimation via reward prediction
- Approach 2: "Model the bias"
  - Counterfactual Model
  - Inverse propensity scoring (IPS) estimator
- Advanced Estimators
  - Self-normalized IPS estimator
  - Doubly robust estimator
  - Slates estimator
- Case Studies
- Summary & Demonstration with code samples

# System Evaluation via Inverse Propensity Scoring

Definition [IPS Utility Estimator]:
Given $S = \big((x_1, y_1, \delta_1), \dots, (x_n, y_n, \delta_n)\big)$ collected under $\pi_0$,

$$\widehat{U}_{ips}(\pi) = \frac{1}{n}\sum_{i=1}^{n} \delta_i \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)}$$

Propensity
$p_i$

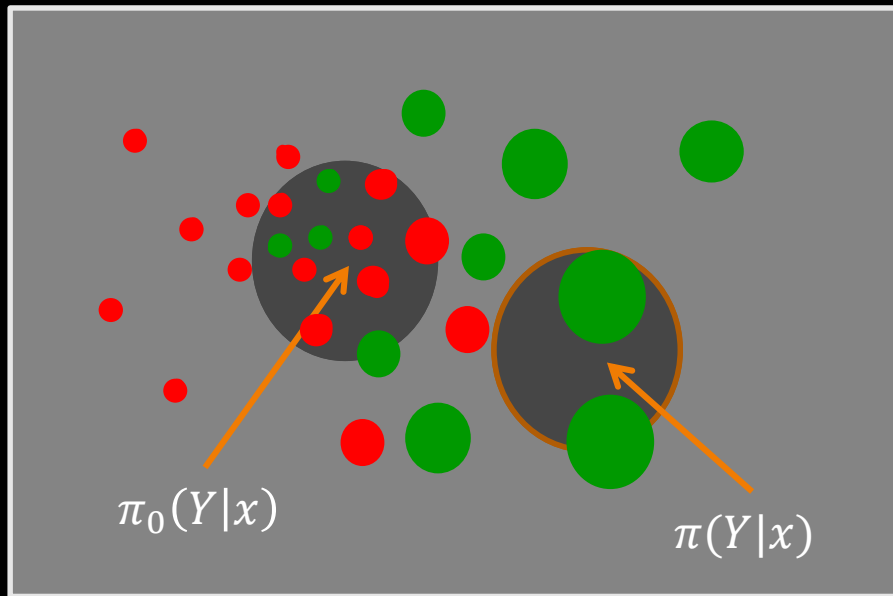→ Unbiased estimate of utility for any $\pi$, if propensity nonzero whenever $\pi(y_i|x_i) > 0$.

Note:

If $\pi = \pi_0$, then online A/B Test with $\widehat{U}_{ips}(\pi_0) = \frac{1}{n}\sum_i \delta_i$

→ Off-policy vs. On-policy estimation.

[Horvitz & Thompson, 1952] [Rubin, 1983] [Zadrozny et al., 2003] [Li et al., 2011]

# Illustration of IPS

IPS Estimator:

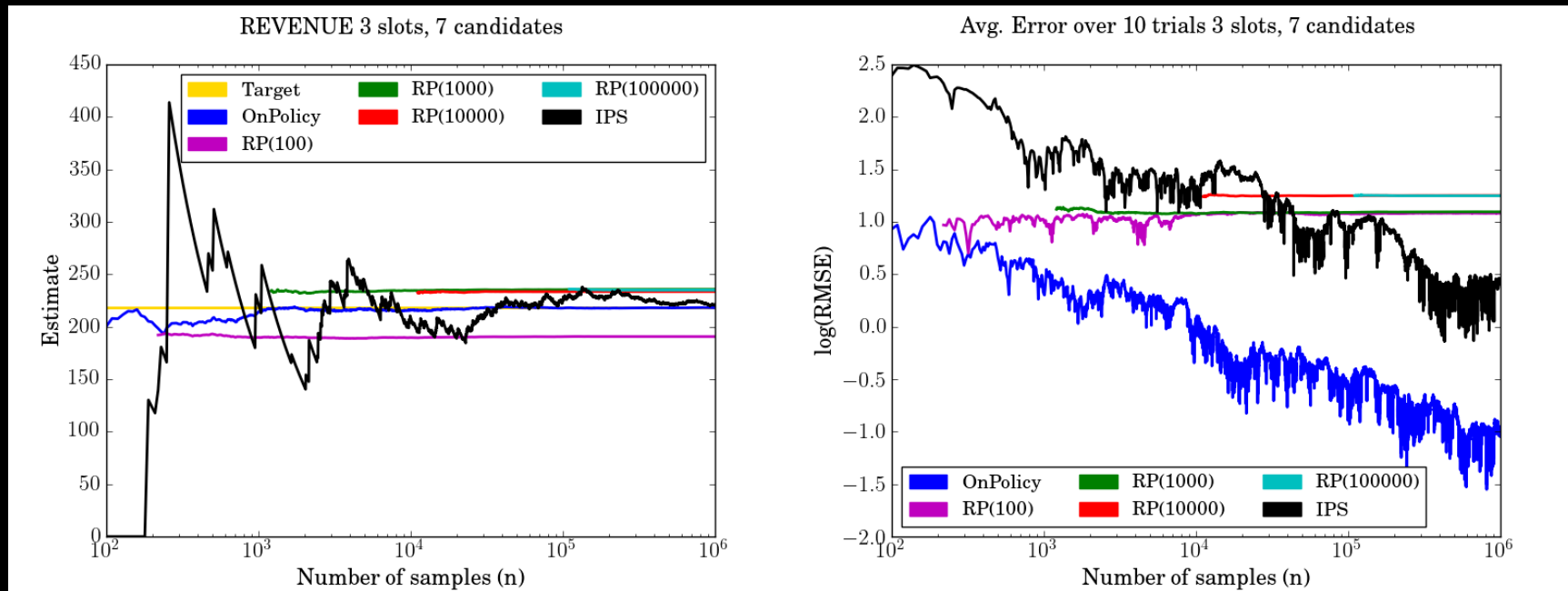$$\widehat{U}_{IPS}(\pi) = \frac{1}{n} \sum_i \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)} \delta_i$$

$\pi_0(Y|x)$

$\pi(Y|x)$

# IPS Estimator is Unbiased

$$E\big[\hat{U}(\pi)\big] = \frac{1}{n} \sum_{x_1,y_1} \ldots \sum_{x_n,y_n} \left[ \sum_i \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)} \delta(x_i,y_i) \right] \pi_0(y_1|x_1) \ldots \pi_0(y_n|x_n) P(x_1) \ldots P(x_n)$$

$$= \frac{1}{n} \sum_{x_1,y_1} \pi_0(y_1|x_1)P(x_1) \ldots \sum_{x_n,y_n} \pi_0(y_n|x_n)P(x_n) \left[ \sum_i \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)} \delta(x_i,y_i) \right]$$

$$= \frac{1}{n} \sum_i \sum_{x_1,y_1} \pi_0(y_1|x_1)P(x_1) \ldots \sum_{x_n,y_n} \pi_0(y_n|x_n)P(x_n) \left[ \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)} \delta(x_i,y_i) \right]$$

$$= \frac{1}{n} \sum_i \sum_{x_i,y_i} \pi_0(y_i|x_i)P(x_i) \left[ \frac{\pi(y_i|x_i)}{\pi_0(y_i|x_i)} \delta(x_i,y_i) \right]$$

Probabilistic Assignment

$$= \frac{1}{n} \sum_i \sum_{x_i,y_i} \pi(y_i|x_i)P(x_i)\delta(x_i,y_i) \quad = \frac{1}{n} \sum_i U(\pi) = U(\pi)$$

# News Recommender: Results



IPS eventually beats RP; variance decays as $O\left(\frac{1}{\sqrt{n}}\right)$

Adith takes over