# Prediction Games:
## From Maximum Likelihood Estimation to Active Learning, Fair Machine Learning, and Structured Prediction

# Brian Ziebart

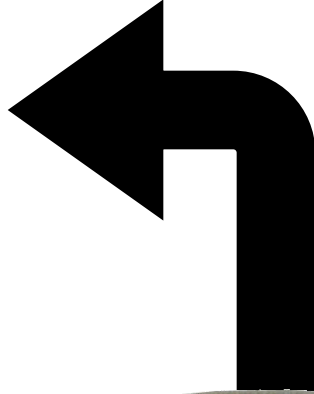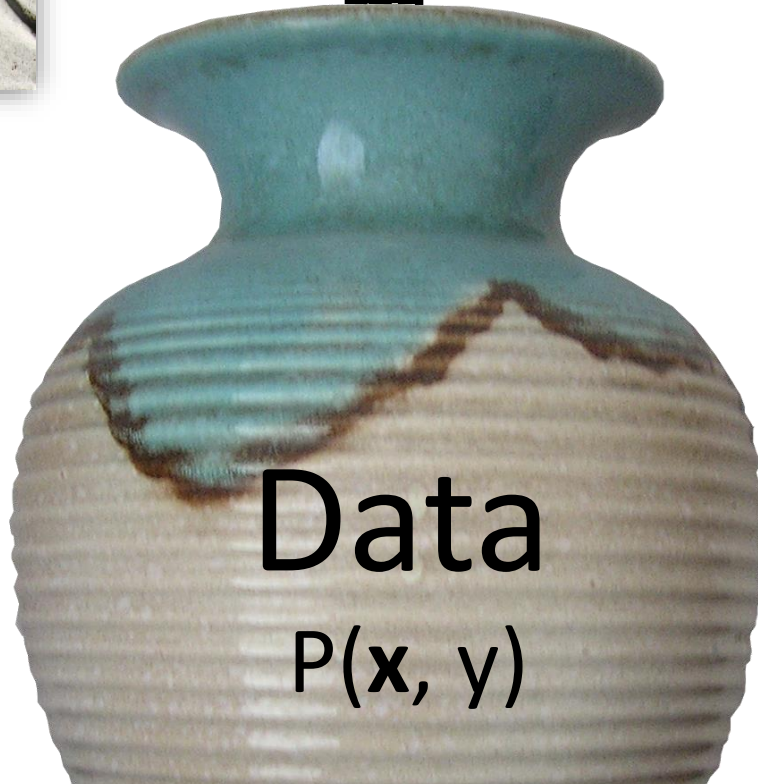**Training**

Viagra Cialis cheap — spam

dog

**x** y

**Sample dist.**
$\tilde{P}(\mathbf{x}, y)$
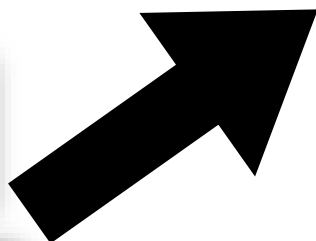
**M samples**

Data
$P(\mathbf{x}, y)$

# Training

**Predictor** f: $X \rightarrow Y$

Viagra
Cialis
cheap
spam

dog

**x**      y

## Sample dist.
$\tilde{P}(\mathbf{x}, y)$

Data
$P(\mathbf{x}, y)$

**Predictor** f: $X \rightarrow Y$  **Testing**

|   |   | y |   |
|---|---|---|---|
|   | **Dog** | **Cat** | **Car** |
| **Dog** | 0 | 1 | 1 |
| **Cat** | 1 | 0 | 1 |
| **Car** | 1 | 1 | 0 |

$\hat{y}$ (left label), y (below table)

**Prediction:** $\hat{y} = f(\mathbf{x})$

dog

**Loss:** $loss(\hat{y}, y)$

**Expected Loss:**
$E_P[loss(f(\mathbf{X}), Y)]$

**dog**

Data
$P(\mathbf{x}, y)$

# Predictor f: $X \rightarrow Y$

## Standard Idea:
**Empirical Risk Minimization**

1. Restrict predictor to some set. f ⊆ Y (or regularize)
2. Minimize (surrogate) loss over predictor set

**Approximate loss**
**Exact training data**

## Dist. Robust Idea:
**Adversarial Risk Minimization**

1. Approximate true label dist. using sample dist. statistics
2. Minimize loss over worst case approximation

**Exact loss**
**Approx. train labels**

Hinge loss   Squared loss

Log loss

0-1 loss

0   w·x

$$\min_{\{\hat{p}_x\}} \max_{\{\check{p}_x\} \in \Xi} \sum_i \hat{p}_{x_i}^\top \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix} \check{p}_{x_i}$$

Data
P(**x**, y)

# Adversarial Supervised Learning Formulation

Construct predictor robust to worst label distribution:

$$\min_{\hat{P}(\hat{y}|\mathbf{x}) \in \Delta} \max_{\check{P}(\check{y}|\mathbf{x}) \in \Delta \cap \Xi} - \sum_{\mathbf{x},y} \tilde{P}(\mathbf{x}) \check{P}(y|\mathbf{x}) \log \hat{P}(y|\mathbf{x})$$

**Logarithmic loss** measures surprise from labels
**Predictor** $\hat{P}$ minimizes loss (in probability simplex $\Delta$)
**Adversary** $\check{P}$ maximizes loss, but must be similar to
  available data (set $\Xi$, e.g., $\mathbb{E}_{\substack{\mathbf{X} \sim \tilde{P} \\ Y|\mathbf{X} \sim \hat{P}}}[\phi(\mathbf{X}, Y)] = \mathbb{E}_{\mathbf{X}, Y \sim \tilde{P}}[\phi(\mathbf{X}, Y)]$)

Reduces to **maximizing entropy** ($\hat{P} = \check{P}$) (Topsøe 1979)
**Robust Bayesian Games** (Grünwald & Dawid 2004)
**DRO** with **expectations constraints** (Wieseman et al. 2014)

$\rightarrow$ **Standard logistic regression and MLE**: $\hat{P}(y|\mathbf{x}) \propto e^{\theta \cdot \phi(\mathbf{x},y)}$

# Part 1: Covariate Shift & Active Learning

Joint work with Anqi Liu (NeurIPS 2014),
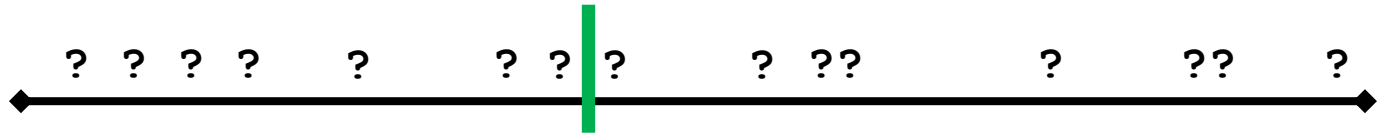Anqi Liu & Lev Reyzin (AAAI 2015)

# Active Learning

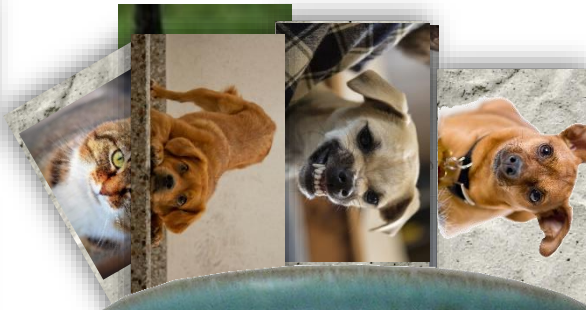Learn linear thresholds: x > c → **+**   (**o** otherwise)

**Passive**

o o o o   o   o o  **+**   **+ ++**   **+**   **++**   **+**

**Active**

? ? ? ?   ?   ? ? ?   ? ??   ?   ??   ?

**4** vs. **15** <u>labels</u>
**O(log n)** vs. **O(n)**

# Training

dog

dog

cat

**Chosen sample dist.**
$$\tilde{P}(\mathbf{x}, y)$$

Data

**Training**
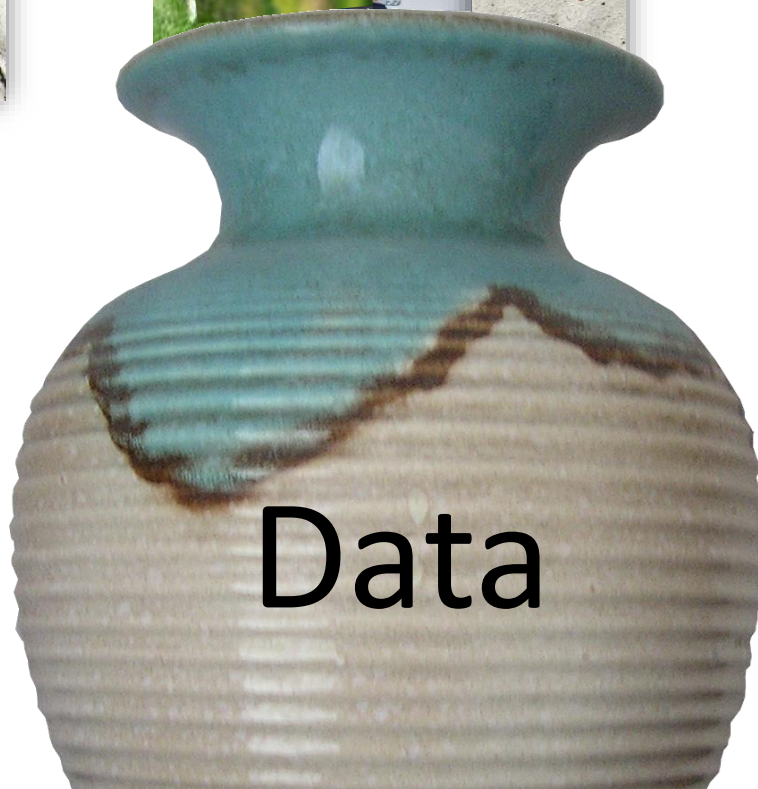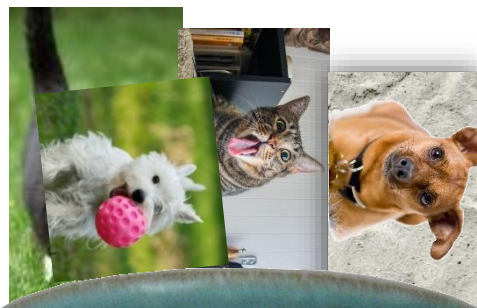
dog

dog

cat

**Predictor** $\hat{P}(y|\mathbf{x})$

**Expected Loss:**

$$\mathbb{E}_{\mathbf{X},Y\sim\tilde{P}}\left[-\log\hat{P}(Y|\mathbf{X})\right]$$

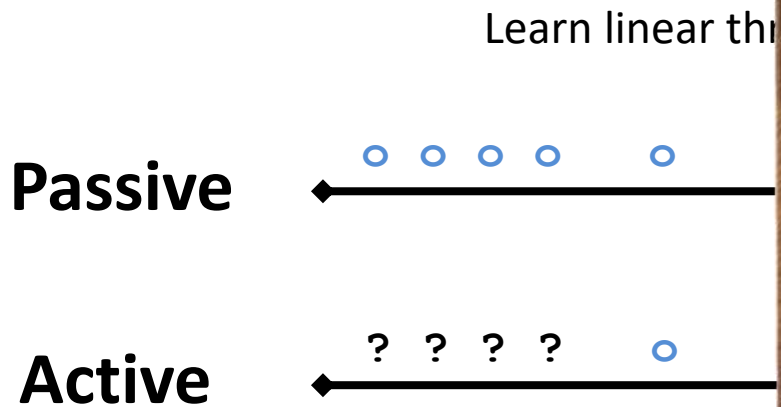**Chosen sample dist.** $\tilde{P}(\mathbf{x},y)$

Data

# Active Learning

Learn linear thr...

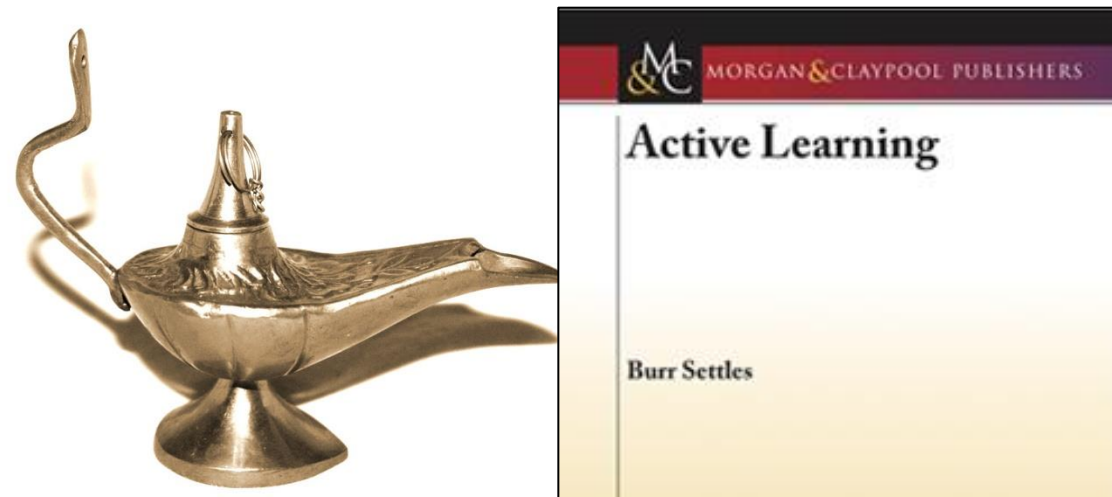**Passive** ○ ○ ○ ○    ○

**Active** ? ? ? ?    ○

## Active Learning

1. Train predictor from labeled data
2. Label most "useful" datapoint
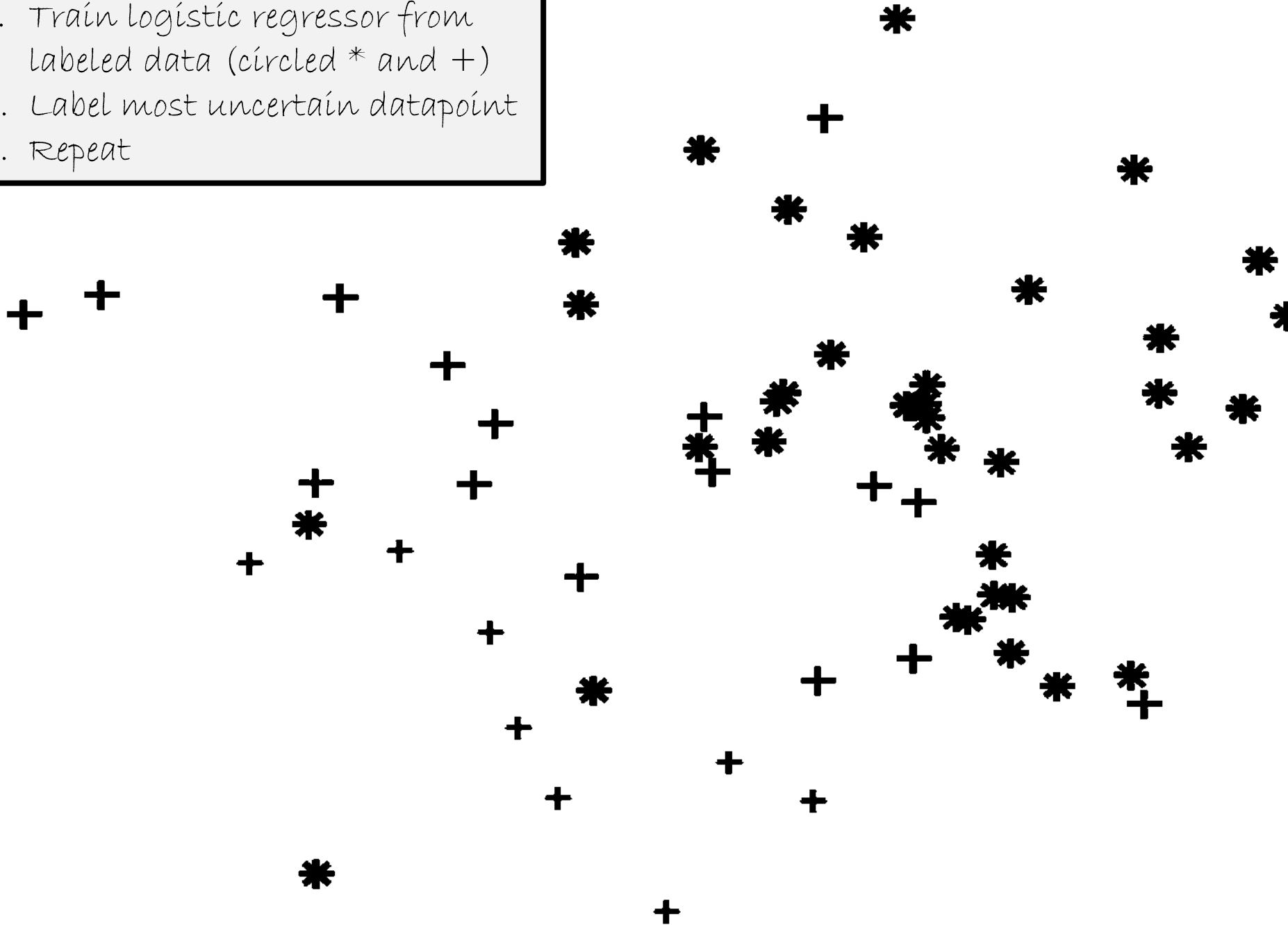3. Repeat

# Results frequently **<u>worse</u>** than passive learning!

"random sampling ... may be more advisable than taking one's chances on active learning with an inappropriate learning model"

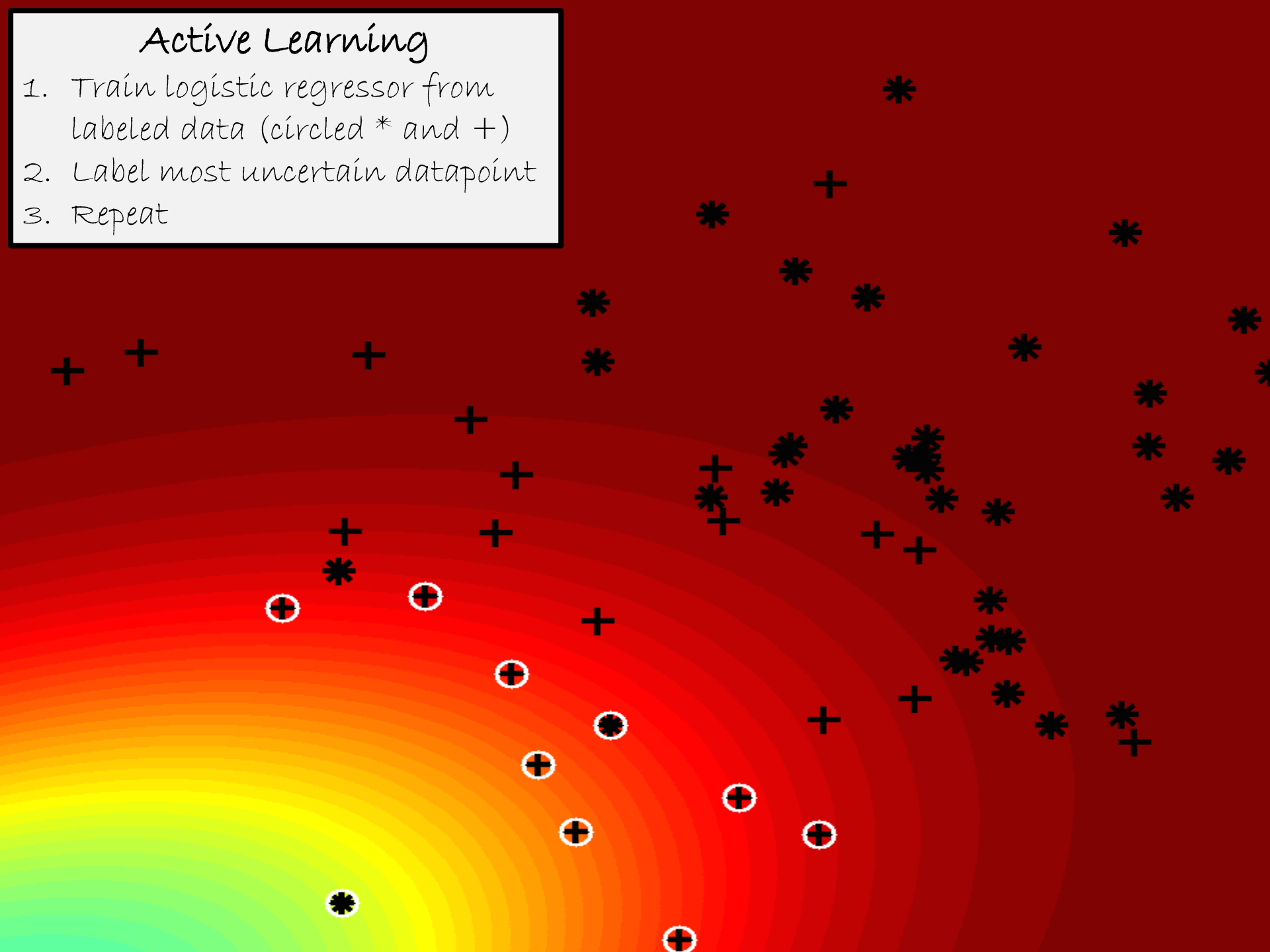Active Learning
1. Train logistic regressor from labeled data (circled * and +)
2. Label most uncertain datapoint
3. Repeat

Active Learning
1. Train logistic regressor from labeled data (circled * and +)
2. Label most uncertain datapoint
3. Repeat

# What's wrong with this recipe?

> **Active Learning**
> 1. Train logistic regressor from labeled data
> 2. Label most uncertain datapoint
> 3. Repeat

# What's wrong with this recipe?

**Active Learning**

1. Train logistic regressor from labeled data
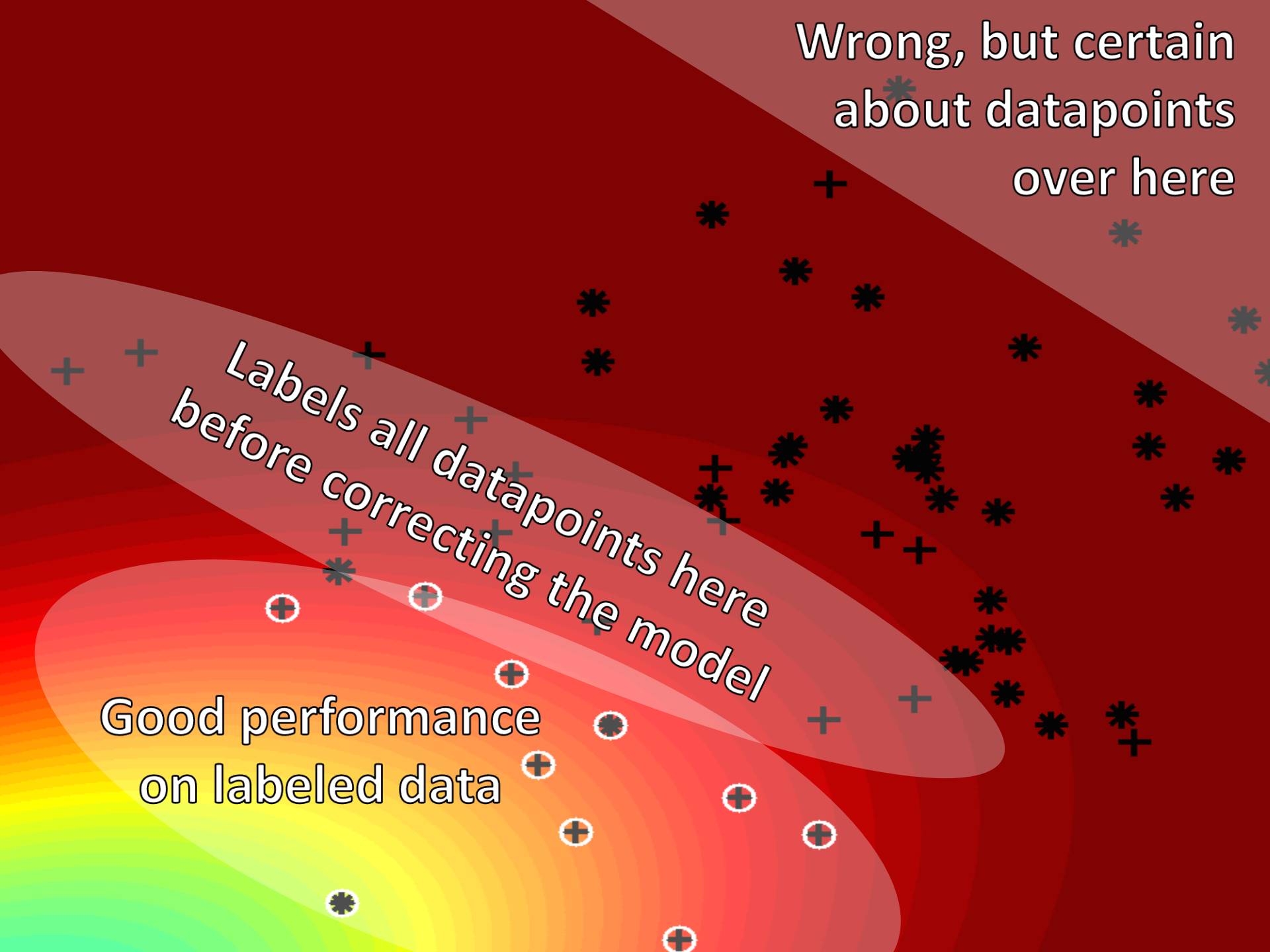2. Label most uncertain datapoint
3. Repeat

Assumes IID data

Produces non-IID data

# Re-Weighted Empirical Risk Minimization

(Shimodaira 2000, Kanamori and Shimodaira 2003)

$$\max_\theta \mathbb{E}_{\tilde{P}(\mathbf{x},y)} \left[ \log \hat{P}_\theta(Y|\mathbf{X}) \right] + \lambda ||\theta||$$

$$\frac{P_{\text{test}}(\mathbf{x})}{P_{\text{train}}(\mathbf{x})}$$

If $\dfrac{P_{\text{test}}(\mathbf{x})}{P_{\text{train}}(\mathbf{x})} > 0$, asymptotically unbiased estimator!

(a) Iris

(b) Seed

(c) Banknote

(d) E. coli

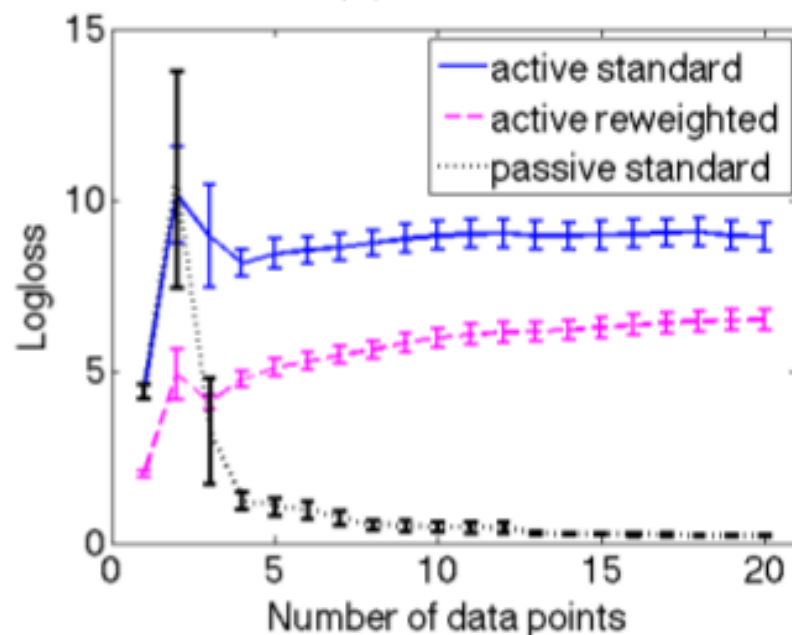# Re-Weighted Empirical Risk Minimization

(Shimodaira 2000, Kanamori and Shimodaira 2003)



**Issues:**

- <u>High variance</u> estimates
- Slow (or <u>no</u>) convergence (Cortes et al. 2010)
- Especially bad for small sample sizes

# Adversarial Prediction for Sample Selection Bias

$$\min_{\hat{P}} \max_{\check{P}} \mathbb{E}_{\tilde{P}(\mathbf{x})\check{P}(\check{y}|\mathbf{x})} \left[ -\log \hat{P}(\check{Y}|\mathbf{X}) \right]$$

$$\text{such that: } \mathbb{E}_{\tilde{P}(\mathbf{x})\check{P}(\check{y}|\mathbf{x})} \left[ \phi(\mathbf{X}, \check{Y}) \right] = \tilde{\phi}$$

IID: $\quad \hat{P}_\theta(y|\mathbf{x}) \propto e^{\theta \cdot \phi(\mathbf{x},y)}$

# Adversarial Prediction for Sample Selection Bias

$$\min_{\hat{P}} \max_{\check{P}} \mathbb{E}_{P_{\text{test}}(\mathbf{x})\check{P}(\check{y}|\mathbf{x})} \left[ -\log \hat{P}(\check{Y}|\mathbf{X}) \right]$$

$$\text{such that: } \mathbb{E}_{P_{\text{train}}(\mathbf{x})\check{P}(\check{y}|\mathbf{x})} \left[ \phi(\mathbf{X}, \check{Y}) \right] = \tilde{\check{\phi}}$$

IID: $\hat{P}_{\theta}(y|\mathbf{x}) \propto e^{\theta \cdot \phi(\mathbf{x},y)}$

Covariate Shift: $\hat{P}_{\theta}(y|\mathbf{x}) \propto e^{\frac{P_{\text{train}}(\mathbf{x})}{P_{\text{test}}(\mathbf{x})} \theta \cdot \phi(\mathbf{x},y)}$

# Adversarial Prediction for Sample Selection Bias

(a) Iris

(b) Seed

(c) Banknote

(d) E. coli

(a) Iris

(b) Seed

(c) Banknote

(d) E. coli

Legend:
- active robust
- active density robust
- passive robust
- active standard
- passive standard
- active reweighted
- passive reweighted

# Part 2: Group Fairness

Joint work with: Ashkan Rezaei, Rizal Fathony, Omid Memarrast (AAAI 2020)

# Fairness for data-driven decision making

<span style="color:red">Group 1</span>  <span style="color:blue">Group 2</span>  (Hardt et al. 2016)

|  | Qualified | Unqualified |
|---|---|---|
| **Accept** | True Positive ($TP_1$) | False Positive ($FP_1$) |
| **Reject** | False Negative ($FN_1$) | True Negative ($TN_1$) |

|  | Qualified | Unqualified |
|---|---|---|
| **Accept** | True Positive ($TP_2$) | False Positive ($FP_2$) |
| **Reject** | False Negative ($FN_2$) | True Negative ($TN_2$) |



**Demographic Parity:** Decision ⫫ Group

$$(TP_1+FP_1)/N_1 = (TP_2+FP_2)/N_2$$

**Equalized Opportunity:** Decision ⫫ Group|Qualified=True

$$TP_1/(TP_1+FN_1) = TP_2/(TP_2+FN_2)$$

**Equalized Odds:** Decision ⫫ Group|Qualified

$$TP_1/(TP_1+FN_1) = TP_2/(TP_2+FN_2); \quad FP_1/(FP_1+TN_1) = FP_1/(FP_1+TN_1)$$

# Fair and Robust Log Loss Predictor

$$\min_{\mathbb{P}\in\Delta\cap\Gamma} \max_{\mathbb{Q}\in\Delta\cap\Xi} \mathbb{E}_{\substack{\widetilde{P}(\mathbf{x},a,y)\\ \mathbb{Q}(\widehat{y}|\mathbf{x},a,y)}} \left[ -\log\mathbb{P}(\widehat{Y}|\mathbf{X},A,Y) \right].$$

$$\Xi : \left\{ \mathbb{Q} \mid \mathbb{E}_{\widetilde{P}(\mathbf{x});\mathbb{Q}(\widehat{y}|\mathbf{x})}[\phi(\mathbf{X},\widehat{Y})] = \mathbb{E}_{\widetilde{P}(\mathbf{x},y)}[\phi(\mathbf{X},Y)] \right\},$$

$\Gamma : \mathbb{P}$ is fair

$$\Gamma : \left\{ \mathbb{P} \mid \frac{1}{p_{\gamma_1}}\mathbb{E}_{\substack{\widetilde{P}(,a,y)\\ \mathbb{P}(\widehat{y}|\mathbf{x},a,y)}}[\mathbb{I}(\widehat{Y}=1 \wedge \gamma_1(A,Y))] = \frac{1}{p_{\gamma_0}}\mathbb{E}_{\substack{\widetilde{P}(\mathbf{x},a,y)\\ \mathbb{P}(\widehat{y}|\mathbf{x},a,y)}}[\mathbb{I}(\widehat{Y}=1 \wedge \gamma_0(A,Y))] \right\}$$

$$\Gamma_{\mathsf{dp}} \iff \gamma_j(A,Y) = \mathbb{I}(A=j);$$
$$\Gamma_{\mathsf{e.opp}} \iff \gamma_j(A,Y) = \mathbb{I}(A=j \wedge Y=1);$$
$$\Gamma_{\mathsf{e.odd}} \iff \gamma_j(A,Y) = \begin{bmatrix} \mathbb{I}(A=j \wedge Y=1) \\ \mathbb{I}(A=j \wedge Y=0) \end{bmatrix}.$$

(Agarwal et al. 2018)

46

$$\widehat{P}_{\theta,\lambda}(\widehat{y}=1|\mathbf{x},a,y) =$$

$$\begin{cases} \min\left\{\frac{\exp(\theta^\top\phi(\mathbf{x},1))}{Z_\theta(\mathbf{x})}, \frac{p_{\gamma_1}}{\lambda}\right\} & \text{if } \gamma_1(a,y) \\ \max\left\{\frac{\exp(\theta^\top\phi(\mathbf{x},1))}{Z_\theta(\mathbf{x})}, 1-\frac{p_{\gamma_0}}{\lambda}\right\} & \text{if } \gamma_0(a,y) \\ \frac{\exp(\theta^\top\phi(\mathbf{x},1))}{Z_\theta(\mathbf{x})} & \text{otherwise}; \end{cases}$$

$$\check{P}_{\theta,\lambda}(\widehat{y}=1|\mathbf{x},a,y) = \widehat{P}_{\theta,\lambda}(\widehat{y}=1|\mathbf{x},a,y) \times$$

$$\begin{cases} \left(1+\frac{\lambda}{p_{\gamma_1}}\widehat{P}_{\theta,\lambda}(\widehat{y}=0|\mathbf{x},a,y)\right) & \text{if } \gamma_1(a,y) \\ \left(1-\frac{\lambda}{p_{\gamma_0}}\widehat{P}_{\theta,\lambda}(\widehat{y}=0|\mathbf{x},a,y)\right) & \text{if } \gamma_0(a,y) \\ 1 & \text{otherwise}. \end{cases}$$

# Benefits of jointly optimizing θ and λ:

# Benefits of jointly optimizing θ and λ:



**Fair Logistic Regression 50%**

# Experiments

|            | Label              | Protected Attribute | Features | Examples |
|------------|--------------------|---------------------|----------|----------|
| **UCI Adult** | Income > $50k   | Gender              | 12       | 45,222   |
| **COMPAS**    | recidivism      | Race                | 10       | 6,167    |

- Evaluation on 20 random splits (70%/30% train/test)
- Baselines
  - Unconstrained (unfair) logistic regression
  - Reweighting approach (Kamiran & Calders 2012)
  - Cost-sensitive reduction approach (Agarwal et al. 2018)
  - Post-processing (Hardt et al. 2016)

Adult

COMPAS

# Part 3: Structured Prediction

Joint work with: Rizal Fathony, Sima Behpour, Xinhua Zhang (ICML 2018)

# Bipartite Matching Task



$\psi_1(\pi_1 = 1)$

$\psi_2(\pi_2 = 1)$

$\psi_1(\pi_1 = 2)$

$\psi_4(\pi_4 = 4)$

A      B

$\pi = [4, 3, 1, 2]$

Maximum weighted bipartite matching:

$$\max_{\pi \in \Pi} \psi(\pi) = \max_{\pi \in \Pi} \sum_i \psi_i(\pi_i)$$

Machine learning task:

Learn appropriate weights $\psi_i(\cdot)$

Objective:

Minimize a loss metric, e.g., the Hamming loss

$$\text{loss}_{\text{Ham}}(\pi, \pi') = \sum_{i=1}^{n} 1(\pi_i' \neq \pi_i)$$

# Bipartite Matching Applications

## Word alignment

(Taskar et. al., 2005; Pado & Lapta, 2006; Mac-Cartney et. al., 2008)

natürlich   ist   das   haus   klein

of   course   the   house   is   small

## Correspondence between images

(Belongie et. al., 2002; Dellaert et. al., 2003)



## Learning to rank documents

(Dwork et. al., 2001; Le & Smola, 2007)

# Learning Bipartite Matchings

## 1 Conditional Random Field

(Petterson et. al., 2009; Volkovs & Zemel, 2012)

$$P_\psi(\pi) = \frac{1}{Z_\psi} \exp\left(\sum_{i=1}^{n} \psi_i(\pi_i)\right)$$

$$Z_\psi = \sum_\pi \prod_{i=1}^{n} \exp\left(\psi_i(\pi_i)\right) = \mathrm{perm}(\mathbf{M})$$

where $M_{i,j} = \exp\left(\psi_i(j)\right)$

✓ **Fisher consistent**
Produces Bayes optimal prediction in ideal case

✗ **Computationally intractable**
Normalization term requires matrix permanent computation (a #P-hard problem). Approximation is needed.

## 2 Structured SVM

(Tsochantaridis et. al., 2005)

Based on CS hinge loss
solved using constraint generation

$$\min_\psi \mathbb{E}_{\pi \sim \tilde{P}}\left[\max_{\pi'} \{\mathrm{loss}(\pi, \pi') + \psi(\pi')\} - \psi(\pi)\right]$$

$\tilde{P}$ is the empirical distribution

✓ **Computationally efficient**
Hungarian algorithm for computing the maximum violated constraints

✗ **No Fisher consistency guarantee**
Not consistent for distribution with no majority label

# Adversarial Bipartite Matchings

[Fathony et al., ICML 2018]

Primal:

$$\min_{\hat{P}(\hat{\pi}|x)} \max_{\check{P}(\check{\pi}|x)} \mathbb{E}_{x \sim \tilde{P}; \hat{\pi}|x \sim \hat{P}; \check{\pi}|x \sim \check{P}} \left[ \text{loss}(\hat{\pi}, \check{\pi}) \right]$$

$$\text{s.t. } \mathbb{E}_{x \sim \tilde{P}; \check{\pi}|x \sim \check{P}} \left[ \sum_{i=1}^{n} \phi_i(x, \check{\pi}_i) \right] = \mathbb{E}_{(x,\pi) \sim \tilde{P}} \left[ \sum_{i=1}^{n} \phi_i(x, \pi_i) \right]$$

**Fisher consistency guaranteed**

Dual:

$$\min_{\theta} \mathbb{E}_{x, \pi \sim \tilde{P}} \min_{\hat{P}(\hat{\pi}|x)} \max_{\check{P}(\check{\pi}|x)} \mathbb{E}_{\substack{\hat{\pi}|x \sim \hat{P} \\ \check{\pi}|x \sim \check{P}}} \left[ \text{loss}(\hat{\pi}, \check{\pi}) + \theta \cdot \sum_{i=1}^{n} \left( \phi_i(x, \check{\pi}_i) - \phi_i(x, \pi_i) \right) \right]$$

Augmented Hamming loss matrix for $n = 3$ permutations

|  | $\check{\pi} = 123$ | $\check{\pi} = 132$ | $\check{\pi} = 213$ | $\check{\pi} = 231$ | $\check{\pi} = 312$ | $\check{\pi} = 321$ |
|---|---|---|---|---|---|---|
| $\hat{\pi} = 123$ | $0 + \delta_{123}$ | $2 + \delta_{132}$ | $2 + \delta_{213}$ | $3 + \delta_{231}$ | $3 + \delta_{312}$ | $2 + \delta_{321}$ |
| $\hat{\pi} = 132$ | $2 + \delta_{123}$ | $0 + \delta_{132}$ | $3 + \delta_{213}$ | $2 + \delta_{231}$ | $2 + \delta_{312}$ | $3 + \delta_{321}$ |
| $\hat{\pi} = 213$ | $2 + \delta_{123}$ | $3 + \delta_{132}$ | $0 + \delta_{213}$ | $2 + \delta_{231}$ | $2 + \delta_{312}$ | $3 + \delta_{321}$ |
| $\hat{\pi} = 231$ | $3 + \delta_{123}$ | $2 + \delta_{132}$ | $2 + \delta_{213}$ | $0 + \delta_{231}$ | $3 + \delta_{312}$ | $2 + \delta_{321}$ |
| $\hat{\pi} = 312$ | $3 + \delta_{123}$ | $2 + \delta_{132}$ | $2 + \delta_{213}$ | $3 + \delta_{231}$ | $0 + \delta_{312}$ | $2 + \delta_{321}$ |
| $\hat{\pi} = 321$ | $2 + \delta_{123}$ | $3 + \delta_{132}$ | $3 + \delta_{213}$ | $2 + \delta_{231}$ | $2 + \delta_{312}$ | $0 + \delta_{321}$ |

size: $n! \times n!$

Intractable for modestly-sized $n$

# Adversarial Bipartite Matchings
## [Fathony et al., ICML 2018]

Dual:
$$\min_{\theta} \mathbb{E}_{(x,\pi)\sim\tilde{P}} \min_{\hat{P}(\hat{\pi}|x)} \max_{\check{P}(\check{\pi}|x)} \mathbb{E}_{\hat{\pi}|x\sim\hat{P};\check{\pi}|x\sim\check{P}} \left[ \sum_{i=1}^{n} I(\pi_i' \neq \pi_i) + \theta \cdot \sum_{i=1}^{n} (\phi_i(x,\check{\pi}_i) - \phi_i(x,\pi_i)) \right]$$

**Marginal Distribution Matrices:**

Predictor
$\mathbf{P} =$

|  | 1 | 2 | 3 |
|---|---|---|---|
| $\hat{\pi}_1$ | $p_{1,1}$ | $p_{1,2}$ | $p_{1,3}$ |
| $\hat{\pi}_2$ | $p_{2,1}$ | $p_{2,2}$ | $p_{2,3}$ |
| $\hat{\pi}_3$ | $p_{3,1}$ | $p_{3,2}$ | $p_{3,3}$ |

$$p_{i,j} = \hat{P}(\hat{\pi}_i = j)$$

Adversary
$\mathbf{Q} =$

|  | 1 | 2 | 3 |
|---|---|---|---|
| $\check{\pi}_1$ | $q_{1,1}$ | $q_{1,2}$ | $q_{1,3}$ |
| $\check{\pi}_2$ | $q_{2,1}$ | $q_{2,2}$ | $q_{2,3}$ |
| $\check{\pi}_3$ | $q_{3,1}$ | $q_{3,2}$ | $q_{3,3}$ |

$$q_{i,j} = \check{P}(\check{\pi}_i = j)$$

**Birkhoff − Von Neumann theorem:**



convex polytope whose points are doubly stochastic matrices

$$\mathbf{P}\mathbf{1} = \mathbf{P}^{\top}\mathbf{1} = \mathbf{Q}\mathbf{1} = \mathbf{Q}^{\top}\mathbf{1} = \mathbf{1}$$

reduce variables from $O(n!)$ to $O(n^2)$

**Marginal Formulation:**

Rearrange the optimization order and add regularization and smoothing penalties

$$\max_{\mathbf{Q}\geq\mathbf{0}} \min_{\theta} \frac{1}{m} \sum_{i=1}^{m} \min_{\mathbf{P}_i\geq\mathbf{0}} \left[ \langle \mathbf{Q}_i - \mathbf{Y}_i, \sum_k \theta_k \mathbf{X}_{i,k} \rangle - \langle \mathbf{P}_i, \mathbf{Q}_i \rangle + \frac{\mu}{2}\|\mathbf{P}_i\|_F^2 - \frac{\mu}{2}\|\mathbf{Q}_i\|_F^2 \right] + \frac{\lambda}{2}\|\theta\|_2^2$$

$$\text{s.t.} : \mathbf{P}_i\mathbf{1} = \mathbf{P}_i^{\top}\mathbf{1} = \mathbf{Q}_i\mathbf{1} = \mathbf{Q}_i^{\top}\mathbf{1} = \mathbf{1}, \quad \forall i$$

projected Quasi-Newton (Schmidt, et.al., 2009) for Q; closed-form for $\theta$;
projection to doubly-stochastic matrix for P using ADMM

# Adversarial Bipartite Matchings
## [Fathony et al., ICML 2018]

### Application: Video Tracking



### Datasets

*Table 3.* Dataset properties

| DATASET | # ELEMENTS | # EXAMPLES |
|---|---|---|
| TUD-CAMPUS | 12 | 70 |
| TUD-STADTMITTE | 16 | 178 |
| ETH-SUNNYDAY | 18 | 353 |
| ETH-BAHNHOF | 34 | 999 |
| ETH-PEDCROSS2 | 30 | 836 |

### Empirical runtime (until convergence)

*Table 5.* Running time (in seconds) of the model for various number of elements $n$ with fixed number of samples ($m = 50$)

| DATASET | # ELEMENTS | ADV MARG. | SSVM |
|---|---|---|---|
| CAMPUS | 12 | 1.96 | 0.22 |
| STADTMITTE | 16 | 2.46 | 0.25 |
| SUNNYDAY | 18 | 2.75 | 0.15 |
| PEDCROSS2 | 30 | 8.18 | 0.26 |
| BAHNHOF | 34 | 9.79 | 0.31 |

Adversarial. Marginal Formulation: grows (roughly) quadratically in $n$

CRF: impractical even for $n = 20$
(Petterson et. al., 2009)

# Adversarial Bipartite Matchings
## [Fathony et al., ICML 2018]

Table 1: The mean and standard deviation (in parenthesis) of the average accuracy (1 - the average Hamming loss) for the adversarial bipartite matching model compared with Structured-SVM.

| TRAINING/ TESTING | ADV. BIPARTITE MATCHING | STRUCTURED SVM |
|---|---|---|
| CAMPUS/ STADTMITTE | 0.662 (0.08) | 0.662 (0.08) |
| STADTMITTE/ CAMPUS | 0.667 (0.11) | 0.660 (0.12) |
| BAHNHOF/ SUNNYDAY | **0.754** (0.10) | 0.729 (0.15) |
| PEDCROSS2/ SUNNYDAY | **0.750** (0.10) | 0.736 (0.13) |
| SUNNYDAY/ BAHNHOF | **0.751** (0.18) | 0.739 (0.20) |
| PEDCROSS2/ BAHNHOF | **0.763** (0.16) | 0.731 (0.21) |
| BAHNHOF/ PEDCROSS2 | **0.714** (0.16) | 0.701 (0.18) |
| SUNNYDAY/ PEDCROSS2 | **0.712** (0.17) | 0.700 (0.18) |

6 pairs of dataset

significantly outperforms SSVM

2 pairs of dataset competitive with SSVM

# Adversarial Bipartite Matchings
## [Fathony et al., ICML 2018]

| | Efficient? | Consistent? | Performs well? |
|---|:---:|:---:|:---:|
| **Conditional Random Field (CRF)**<br>(Petterson et. al., 2009; Volkovs & Zemel, 2012) | ✖ | ✔ | ? |
| **Structured SVM**<br>(Tsochantaridis et. al., 2005) | ✔ | ✖ | ▬ |
| **Adversarial Bipartite Matching**<br>(our approach) | ✔ | ✔ | ✔ |

# Summary & Conclusions

$$\min_{\hat{P}(\hat{y}|\mathbf{x}) \in \Delta \cap \Gamma} \max_{\check{P}(\check{y}|\mathbf{x}) \in \Delta \cap \Xi} \mathbb{E}_{\substack{\mathbf{x} \sim \tilde{P} \\ \hat{y}|\mathbf{x} \sim \hat{P} \\ \check{y}|\mathbf{x} \sim \check{P}}} \left[ \text{loss}(\hat{\mathbf{Y}}, \check{\mathbf{Y}}) \right]$$

**Covariate Shift/Active Learning:** $P_{\text{train}}(\mathbf{x}) \neq P_{\text{test}}(\mathbf{x})$

→ **Avoids harmful extrapolations**

**Fairness:** Minimizer also satisfies fairness requirements ($\Gamma$)

→ **Robust/smooth group fairness**

**Structured Prediction:** Structured objects **y**, bilinear loss

→ **Consistency and Computational Tractability**

Foundational framework for parametric predictors

Versatile for a wide range of settings

# Questions?

- Liu, Ziebart. "*Robust Classification Under Sample Selection Bias*." NeurIPS 2014.
- Liu, Reyzin, Ziebart. "*Shift-Pessimistic Active Learning Using Robust Bias-Aware Prediction*." AAAI 2015.
- Rezaei, Fathony, Memmarest, Ziebart. "*Fairness for robust log loss classification*." AAAI 2020.
- Fathony, Behpour, Zhang, Ziebart. "*Efficient and Consistent Adversarial Bipartite Matching*." ICML 2018.
- Fathony, Liu, Asif, Ziebart. "*Adversarial Multiclass Classification: A Risk Minimization Perspective*." NeurIPS 2016.
- Fathony, Bashiri, Ziebart. "*Adversarial Surrogate Losses for Ordinal Regression*." NeurIPS 2017.
- Wang, Xing, Asif, Ziebart. "*Adversarial Prediction Games for Multivariate Losses*." NeurIPS 2015.
- Tirinzoni, Chen, Petrik, Ziebart. "*Policy-conditioned uncertainty sets for robust Markov Decision Processes*." NeurIPS 2018.
- Bashiri, Ziebart, Zhang. "*Distributionally Robust Imitation Learning*." NeurIPS 2021.