

Learning from Aggregates

Giorgio Patrini

Australian National University, NICTA

Summary

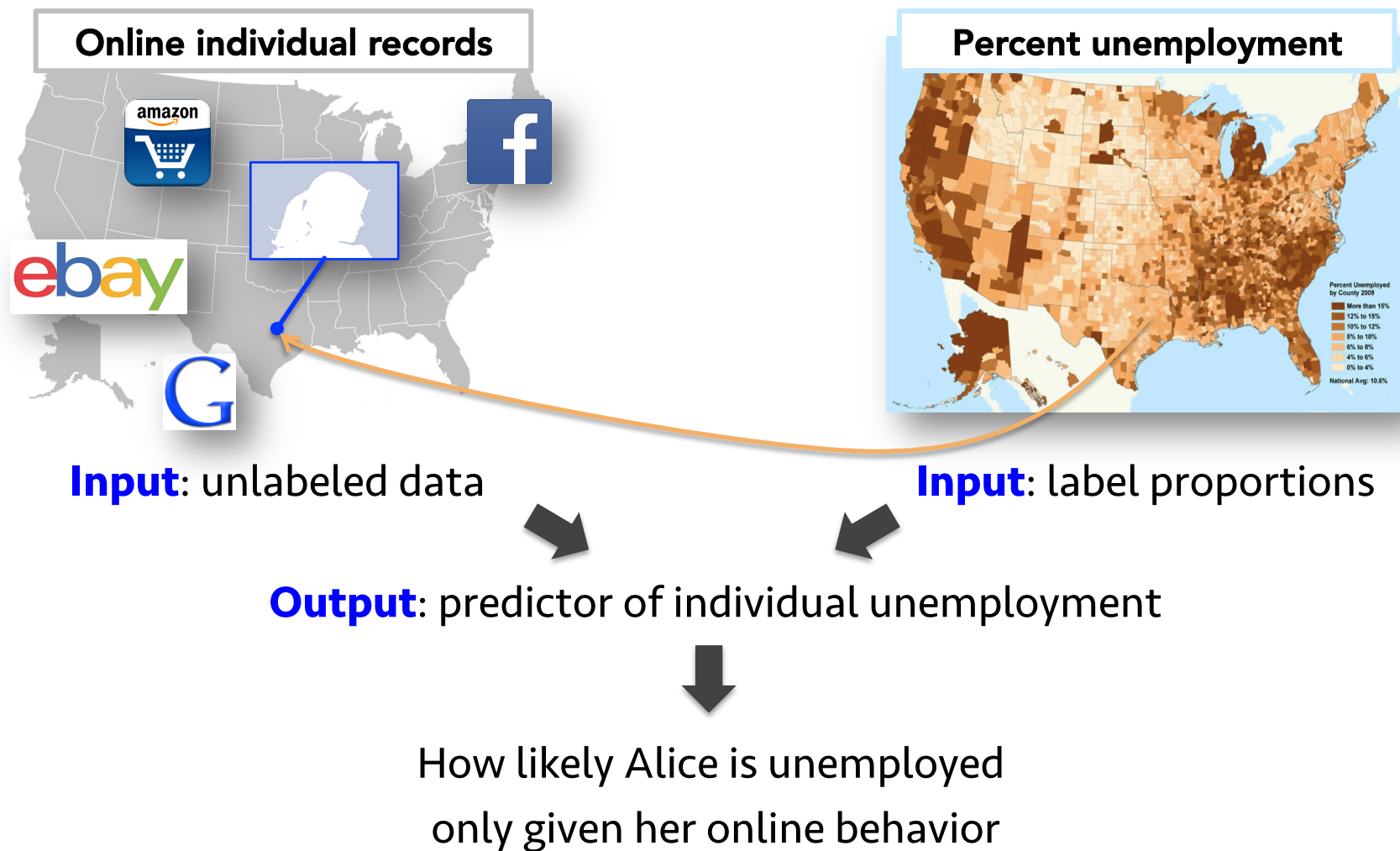
- Learning from label proportions
- Laplacian Mean Map algorithm

G.Patrini, R.Nock, P.Rivera, T.Caetano, (Almost) no label no cry, NIPS'14

- Do we need individual feature vectors?

R.Nock, G.Patrini, A.Friedman, Rademacher observations, private data, and boosting, ICML'15

Learning from Label Proportions (LLP)



Learning from Label Proportions (LLP)

Other applications:

- Bags of images/pixels in Computer Vision
- Classify sentences as positive/negative based on overall review score
- Data comes from physical measurements which are technically **feasible only in aggregated form**
- Potentially, applications already explored by Multiple Instance Learning (MIL)

Learning setting

- Sample $\mathcal{S} = \{(\mathbf{x}_i, y_i), i \in [m]\}$, on $\mathbb{R}^d \supseteq \mathcal{X} \times \{-1, +1\}$
- **No label is observed**
- Known: partition of bags $\cup_j \mathcal{S}_j = \mathcal{S}, j \in [n]$, and relative **label proportions** π_j
- (No assumption on how the bags were made)

Goal: learn a binary (linear) classifier θ for individual feature vectors x to predict the label as $\text{sgn } \langle \theta, x \rangle$

Our solution, step 1: factorisation theorem

Def (Altun&Smola COLT'06): the **mean operator**

$$\mu = 1/m \sum_{i=1}^m y_i x_i$$

Thm (**proper losses factorisation**): μ is **sufficient** for the label variable for most proper losses:

$$\text{PROPER-LOSS} = \text{LOSS w/o LABELS}(\theta) - \frac{1}{2} \langle \theta, \mu \rangle$$

Our solution, step 1: factorisation theorem

Def (Altun&Smola COLT'06): the **mean operator**

$$\mu = 1/m \sum_{i=1}^m y_i x_i$$

Thm (**proper losses factorisation**): μ is **sufficient** for the label variable for most proper losses:

$$\text{PROPER-LOSS} = \text{LOSS w/o LABELS}(\theta) - \frac{1}{2} \langle \theta, \mu \rangle$$


**e.g., classic
logistic loss**

$$\operatorname{argmin}_{\theta} \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y\theta^\top x_i}) =$$

$$\operatorname{argmin}_{\theta} \frac{1}{m} \sum_{i=1}^m \log \sum_{y \in \{-1,1\}} e^{-y\theta^\top x_i} - \langle \theta, \frac{1}{2m} \sum_{i=1}^m y_i x_i \rangle$$

Our solution, step 2: estimate the mean operator

$$\begin{aligned}\mu &= \sum_{j=1}^n p(j) \mu_j = \sum_{j=1}^n p(j) \sum_{y \in \{-1, 1\}} y p(y|j) \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j, y] \\ &= \sum_{j=1}^n p(j) (\pi_j \mathbf{b}_j^+ - (1 - \pi_j) \mathbf{b}_j^-) \end{aligned}$$

 $\mathbf{b}_j^y = \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j, y]$

Then, come up with a system of equations with \mathbf{b}_j^y as only unknowns:

$$\mathbf{b}_j = \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j] = \sum_{y \in \{-1, 1\}} p(y|j) \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j, y] = \sum_{y \in \{-1, 1\}} \pi_j \mathbf{b}_j^y$$

$$\mathbf{b}_j = \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j] = \sum_{y \in \{-1,1\}} p(y|j) \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j, y] = \sum_{y \in \{-1,1\}} \pi_j \mathbf{b}_j^y$$

2 variables for each equation!

Quadrianto et al. JMLR'09

$$\mathbf{b}_j = \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j] = \sum_{y \in \{-1,1\}} p(y|j) \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j, y] = \sum_{y \in \{-1,1\}} \pi_j \mathbf{b}_j^y$$

2 variables for each equation!

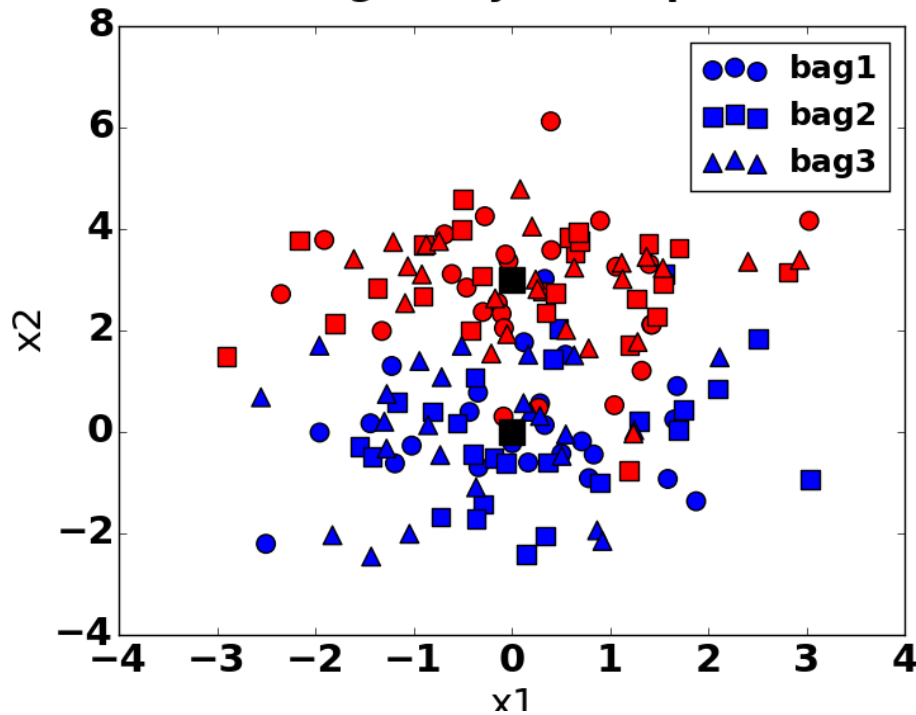
Solution of Quadrianto et al. JMRL'09 with Mean Map, **homogeneity** assumption:

$$\forall_j \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j, y] = \mathbb{E}_{\mathcal{S}}[\mathbf{x}|y]$$

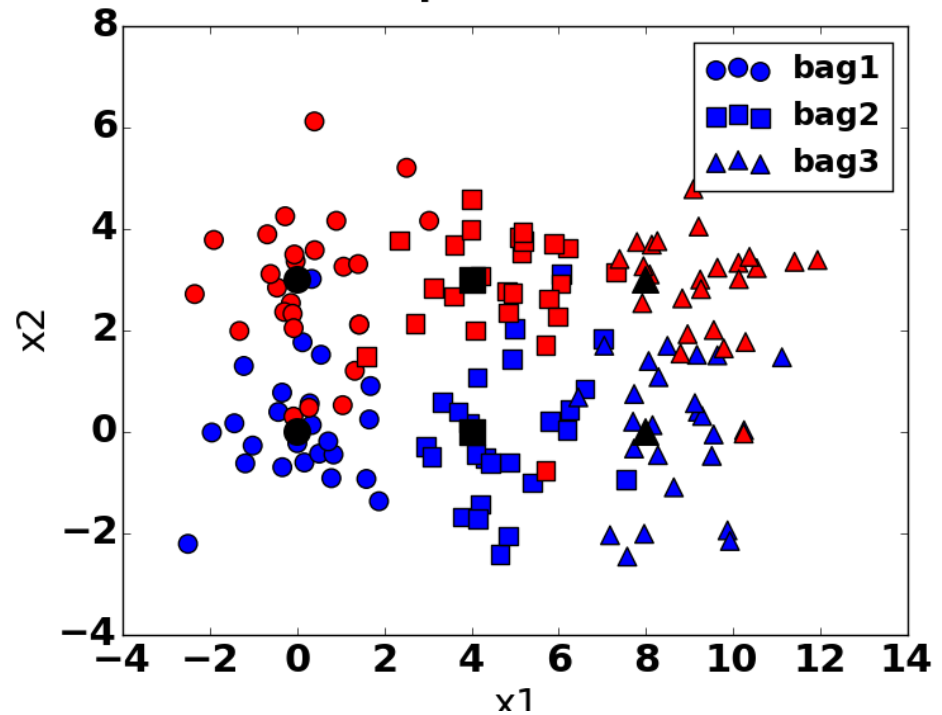
“Unemployed people in all the counties behave online in the same way, in average”

Homogeneity assumption: $\forall_j \mathbb{E}_{\mathcal{S}}[x|j, y] = \mathbb{E}_{\mathcal{S}}[x|y]$

Homogeneity assumption



Assumption violation



We relax it

$$\mathbf{b}_j = \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j] = \sum_{y \in \{-1,1\}} p(y|j) \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j, y] = \sum_{y \in \{-1,1\}} \pi_j \mathbf{b}_j^y$$

We only asks **smoothness** on “similar” bags:

$$\forall_{j,j'} \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j] \approx \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j'] \implies \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j, y] \approx \mathbb{E}_{\mathcal{S}}[\mathbf{x}|j', y]$$

“The more similar the counties, the more similar the online behaviour of the people unemployed there”

Our solution, step 3: Laplacian regularization

Let $v_{j,j'}$ be the similarity between bags. Then we can encode our assumption in a **regularized least square**:

$$\operatorname{argmin}_{\mathbf{b}_j^y} \sum_j (\mathbf{b}_j - \sum_{y \in \{-1,1\}} \pi_j \mathbf{b}_j^y)^2 + \gamma \sum_{j,j'} v_{j,j'} [(\mathbf{b}_j^+ - \mathbf{b}_{j'}^+)^2 + (\mathbf{b}_j^- - \mathbf{b}_{j'}^-)^2]$$

Then, in matrix form:

$$\begin{aligned} \mathbf{B} &= [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n]^\top, \quad \mathbf{B}^\pm = [\mathbf{b}_1^+, \mathbf{b}_2^+, \dots, \mathbf{b}_n^+, \mathbf{b}_1^-, \mathbf{b}_2^-, \dots, \mathbf{b}_n^-]^\top, \\ \Pi &= [\text{DIAG}(\boldsymbol{\pi}) | \text{DIAG}(\mathbf{1} - \boldsymbol{\pi})] \end{aligned}$$

$$\operatorname{argmin}_{\mathbf{B}^\pm} \operatorname{tr} ((\mathbf{B} - \Pi \mathbf{B}^\pm)^\top (\mathbf{B} - \Pi \mathbf{B}^\pm)) + \gamma \operatorname{tr} ((\mathbf{B}^\pm)^\top \mathbf{L} \mathbf{B}^\pm)$$

**Laplacian
matrix on** $v_{j,j'}$



Our solution: Laplacian Mean Map algorithm

(steps in reverse order)

Laplacian Mean Map (LMM)

Input $\mathcal{S}_j, \pi_j, j \in [n]; \lambda, \gamma > 0; \mathbf{V};$

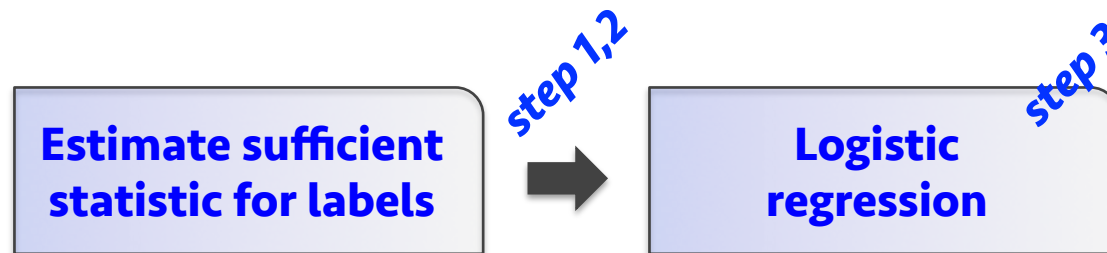
Step 1 : let $\mathbf{B}^\pm \leftarrow (\mathbf{\Pi}\mathbf{\Pi}^T + \gamma\mathbf{L})^{-1}\mathbf{\Pi}\mathbf{B}$

Step 2 : let $\boldsymbol{\mu} \leftarrow \sum_j p_j (\pi_j \mathbf{b}_j^+ - (1 - \pi_j) \mathbf{b}_j^-)$

Step 3 : let $\boldsymbol{\theta}_* \leftarrow \arg \min_{\boldsymbol{\theta}} \text{LOSS W/O LABEL}(\boldsymbol{\theta}) + \frac{1}{2} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle + \lambda \|\boldsymbol{\theta}\|_2^2;$

Return $\boldsymbol{\theta}^*$

Scalability: Step 1 is only $O(n^3) \ll O(m^3)$



Approximation of the mean operator

Theorem 1 Suppose that γ satisfies $\gamma\sqrt{2} \leq \max_{j \neq j'} v_{jj'}$. Let $M \doteq [\boldsymbol{\mu}_1 | \boldsymbol{\mu}_2 | \dots | \boldsymbol{\mu}_n]^\top \in \mathbb{R}^{n \times d}$, $\tilde{M} \doteq [\tilde{\boldsymbol{\mu}}_1 | \tilde{\boldsymbol{\mu}}_2 | \dots | \tilde{\boldsymbol{\mu}}_n]^\top \in \mathbb{R}^{n \times d}$ and $\psi(V, B^\pm) \doteq (\max_{j \neq j'} v_{jj'})^2 \|B^\pm\|_F$. The following holds:

$$\|M - \tilde{M}\|_F \leq \sqrt{n/2} \times \psi(V, B^\pm) .$$

(Assuming homogeneity with Mean Map, the norm is unbounded.)

Choose the similarity $v_{jj'}^G \doteq \exp(-\|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2)$

Under mild conditions, it holds, w.r.t. the max norm of $\mathbf{b}_j^y = \mathbb{E}_S[x|j, y]$:

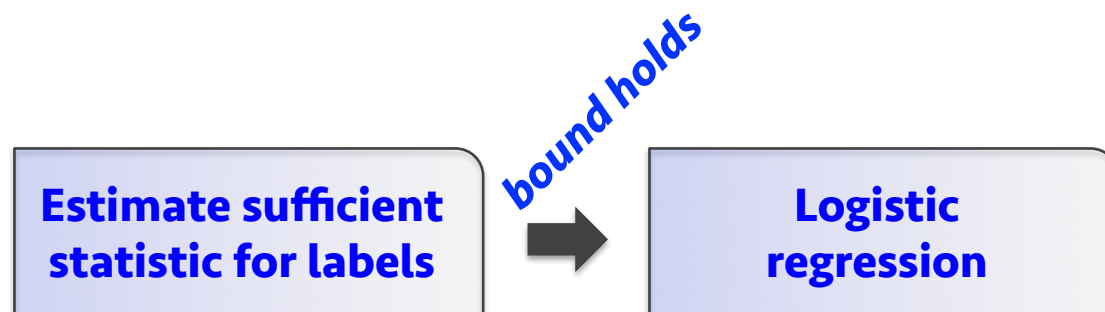
$$\psi(V^G, B^\pm) = o(1)$$

Approximation of the model

Theorem 1 *Let θ_* be the model computed with the true mean operator μ . Let $\tilde{\mu}$, $\tilde{\theta}_*$ be the respective estimates. For any proper loss L_2 -regularized with parameter $\lambda > 0$, there exists $q > 0$ such that:*

$$\|\tilde{\theta}_* - \theta_*\|_2^2 \leq 1/(2\lambda + q) \|\tilde{\mu} - \mu\|_2^2$$

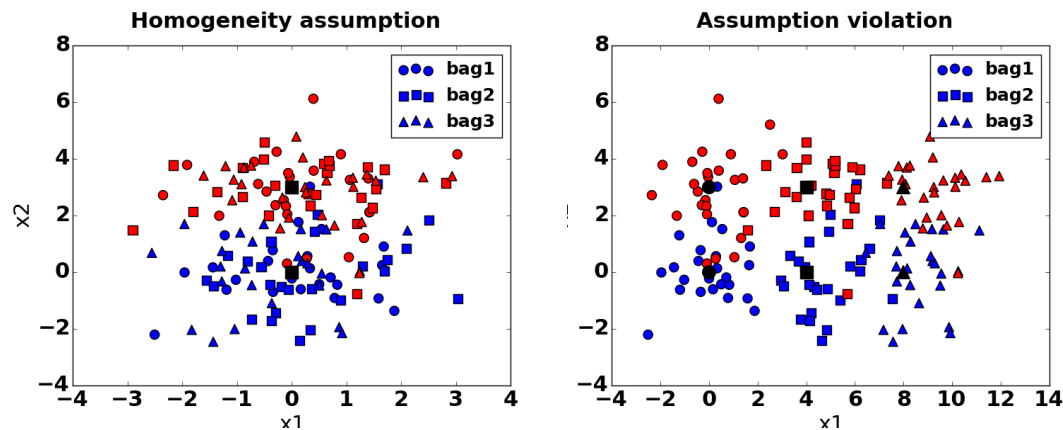
This holds for **any estimator of μ** , even outside the LLP setting



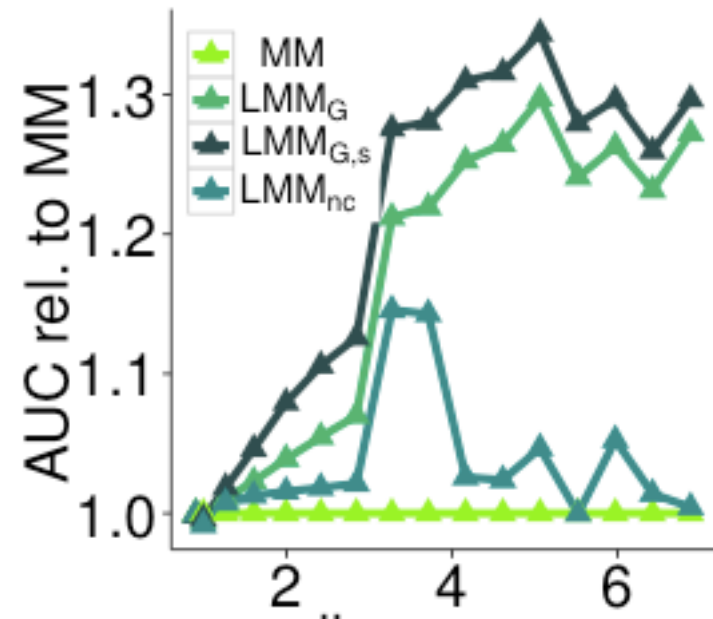
And more in the paper

- **Alternating Mean Map**: use LMM as initialization and optimizes further, inferring labels as latent variables (similar to Expectation Maximization)
- We also provide **generalization bounds** based on **Rademacher Complexity**.

Experiments: homogeneity assumption



**gradual violation
of homogeneity**



Experiments: comparative tests

14 UCI datasets **converted to LLP** (up to ~300K examples)

- Select a categorical feature, use its value to assign bags and proportions; then remove the feature.
- Compare with SVMs (*Yu et al. ICML'13*) and InvCal (*Rueping ICML'10*)

Table 1: 10 small domains results. #win/#lose for row vs column on 50 tests; ties not reported. Bold faces when $p\text{-val} < .001$ for Wilcoxon signed-rank tests.

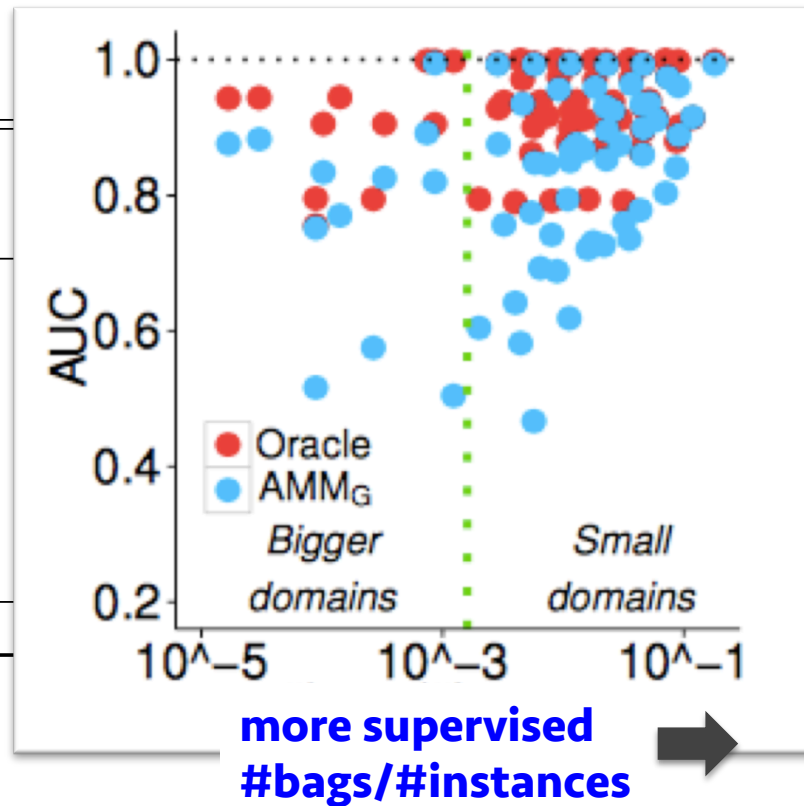
algorithm		MM	G	LMM G,s	nc	InvCal	MM	AMM ^{min}		10ran	conv- \propto SVM
LMM	G	36/4									
	G,s	38/3	30/6								
	nc	28/12	3/37	2/37							
	InvCal	4/46	3/47	4/46	4/46						
AMM ^{min}	MM	33/16	26/24	25/25	32/18	46/4					
	G	38/11	35/14	30/20	37/13	47/3	31/7				
	G,s	35/14	33/17	30/20	35/15	47/3	24/11	7/15			
	10ran	27/22	24/26	22/28	26/24	44/6	20/ 30	16/ 34	19/ 31		
SVM	conv- \propto	21/ 29	2/ 48	2/ 48	2/ 48	2/ 48	4/ 46	3/ 47	3/ 47	4/ 46	
	alter- \propto	0/ 50	0/ 50	0/ 50	0/ 50	20/ 30	0/ 50	0/ 50	0/ 50	3/ 47	27/23

Experiments: no label no cry

algorithm	<i>adult</i> : 48842 × 89			<i>marketing</i> : 45211 × 41			<i>census</i> : 299285 × 381			
	IV(5)	V(16)	VI(42)	V(4)	VII(4)	VIII(12)	IV(5)	VIII(9)	VI(42)	
MM	80.93	76.65	74.01	54.64	50.71	49.70	75.21	90.37	75.52	
LMM _G	81.79	78.40	78.78	54.66	51.00	51.93	75.80	71.75	76.31	
LMM _{G,s}	84.89	78.94	80.12	49.27	51.00	65.81	84.88	60.71	69.74	
AMM ^{min}	AMM _{MM}	83.73	77.39	80.67	52.85	75.27	58.19	89.68	84.91	68.36
	AMM _G	83.41	82.55	81.96	51.61	75.16	57.52	87.61	88.28	76.99
	AMM _{G,s}	81.18	78.53	81.96	52.03	75.16	53.98	89.93	83.54	52.13
	AMM ₁	81.32	75.80	80.05	65.13	64.96	66.62	89.09	88.94	56.72
AMM ^{max}	AMM _{MM}	82.57	71.63	81.39	48.46	51.34	56.90	50.75	66.76	58.67
	AMM _G	82.75	72.16	81.39	50.58	47.27	34.29	48.32	67.54	77.46
	AMM _{G,s}	82.69	70.95	81.39	66.88	47.27	34.29	80.33	74.45	52.70
	AMM ₁	75.22	67.52	77.67	66.70	61.16	71.94	57.97	81.07	53.42
Oracle	90.55	90.55	90.50	79.52	75.55	79.43	94.31	94.37	94.45	

Experiments: no label no cry

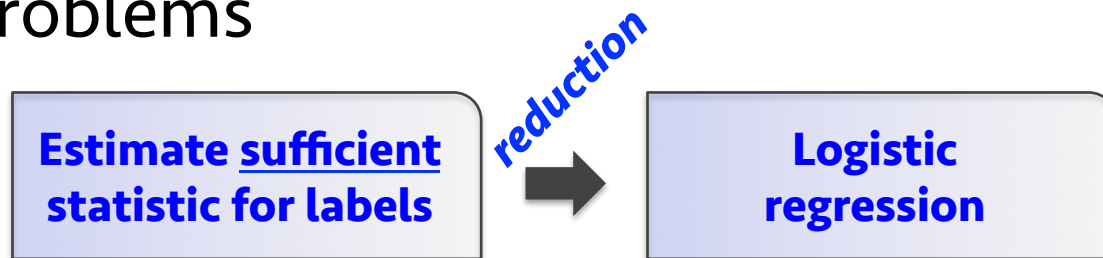
algorithm		<i>adult</i> : IV(5)
MM		80.93
LMM _G		81.79
LMM _{G,s}		84.89
AMM ^{min}	AMM _{MM}	83.73
	AMM _G	83.41
	AMM _{G,s}	81.18
	AMM ₁	81.32
AMM ^{max}	AMM _{MM}	82.57
	AMM _G	82.75
	AMM _{G,s}	82.69
	AMM ₁	75.22
Oracle		90.55



<i>census</i> : 299285 × 381			
IV(5)	VIII(9)	VI(42)	
75.21	90.37	75.52	
75.80	71.75	76.31	
84.88	60.71	69.74	
89.68	84.91	68.36	
87.61	88.28	76.99	
89.93	83.54	52.13	
89.09	88.94	56.72	
50.75	66.76	58.67	
48.32	67.54	77.46	
80.33	74.45	52.70	
57.97	81.07	53.42	
94.31	94.37	94.45	

Take-home messages (until here)

- **(Almost) no label no cry**: few proportions can suffice to learn. Privacy threat?
- **Sufficiency of mean operator**: any “*weakly-supervised*” learner can exploit the same trick, *e.g.* semi-supervised, MIL, noisy labels. Bound for the classifier holds.
- **Do not reinvent the wheel**: *reduction* between ML problems



But what about individual *feature vectors*?

- Is there an analogue of the mean operator that allows us to learn with *aggregate feature vectors*?
- YES. Define a **Rademacher observation** as a (non-normalized) mean operator restricted to a subsample $s \in \mathcal{S}$:

$$\mu_s = \sum_{i: (\mathbf{x}_i, y_i) \in s} y_i \mathbf{x}_i$$

Rademacher observations and logistic loss

$$\operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y\boldsymbol{\theta}^\top x_i}) =$$

$$\operatorname{argmin}_{\boldsymbol{\theta}} \log(2) + \frac{1}{m} \log \left(\frac{1}{2^m} \sum_{s \subseteq \mathcal{S}} e^{-\boldsymbol{\theta}^\top \boldsymbol{\mu}_s} \right) \leftarrow \text{they are all aggregated here}$$

The number of $\boldsymbol{\mu}_s$ is exponential in m , but we can still learn on a small subset of Rademacher observations. See our *ICML'15* for details.

Yes, but why?

- When we have all the data but do not want to share it entirely with the learner, but still want to learn good models. **Privacy** constraints.
 - Can prove differential privacy
 - Properties of non-reconstruct-ability of the original data (NP-hardness and algebraic impossibility)

Conclusion

Learning from aggregate data is possible, with unexpected applications on

- weakly-supervised learning
- privacy
- distributed learning - one μ_s per cluster?
- and social sciences, *e.g.* the ecological inference

➤ **NIPS'15 workshop** on "*Learning and privacy with incomplete data and weak supervision*"