

\propto SVM for Learning with Label Proportions

Felix X. Yu[†] Dong Liu[†] Sanjiv Kumar[§] Tony Jebara[†] Shih-Fu Chang[†]

[†]Columbia University, New York, NY 10027

[§]Google Research, New York, NY 10011

\propto is the symbol for “proportional-to”.



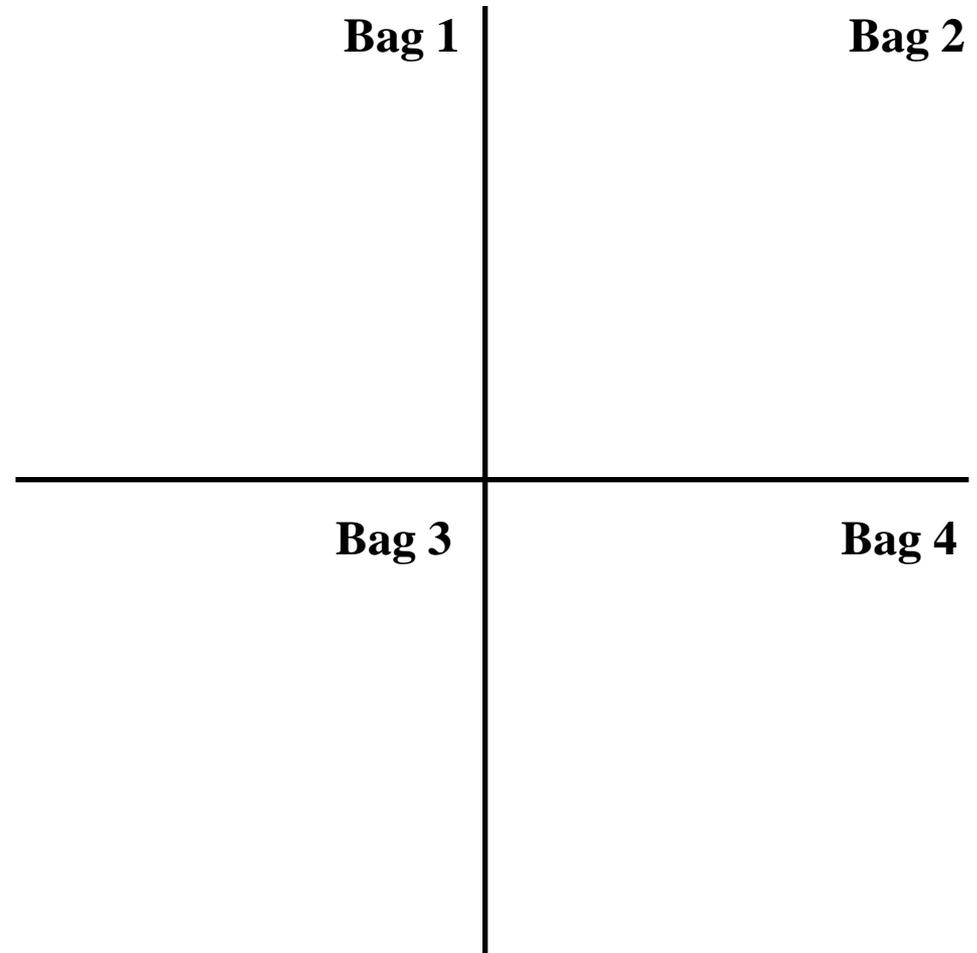
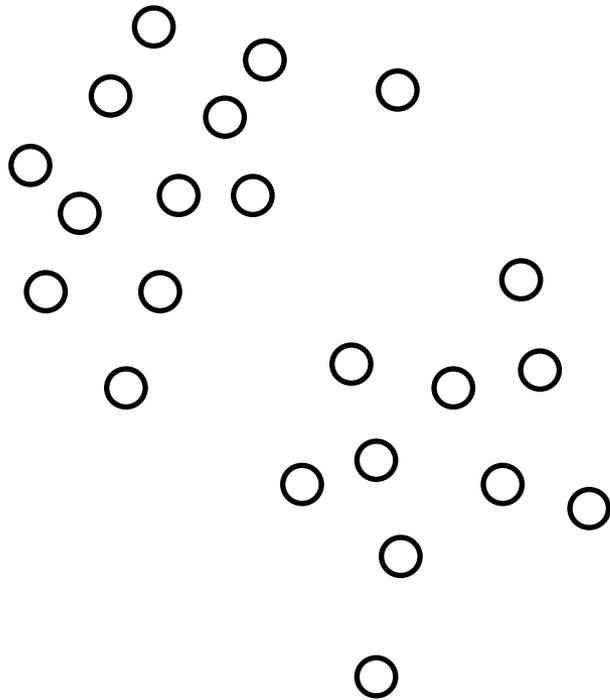
COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

Google™

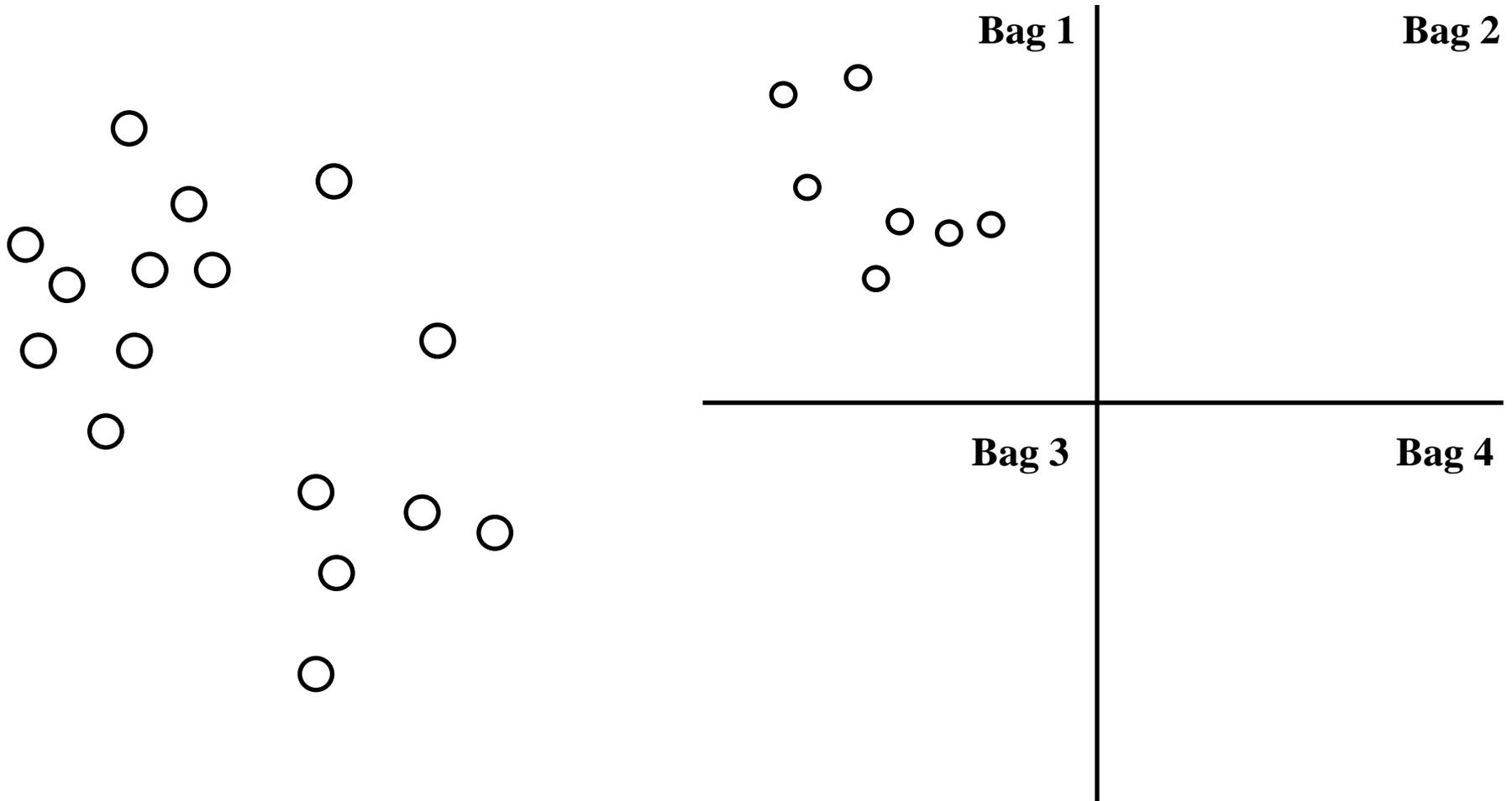
Outline

- Learning Setting and Applications
- Related Works
- Formulation
- Algorithms
- Experiments

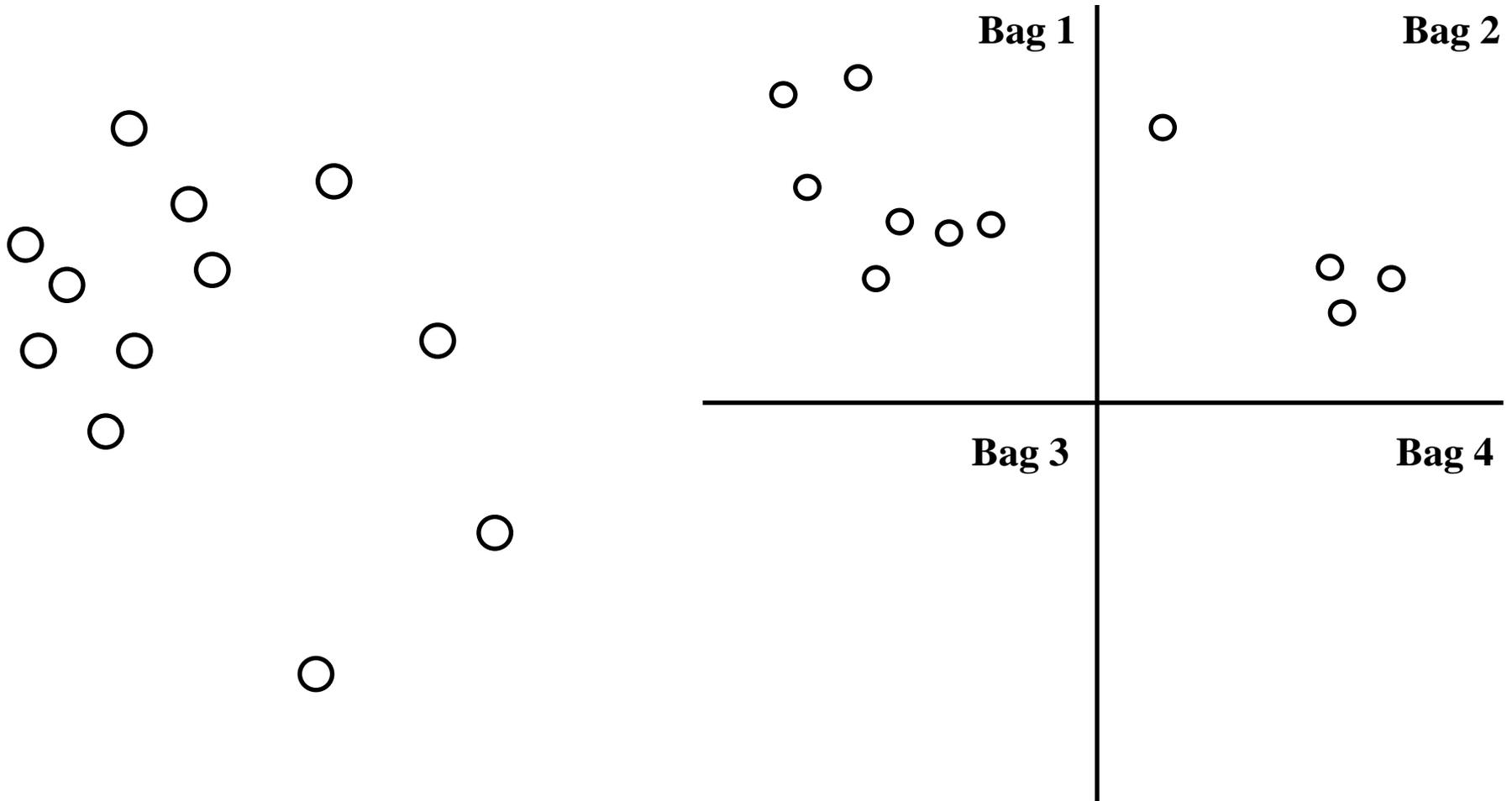
Learning with label proportions



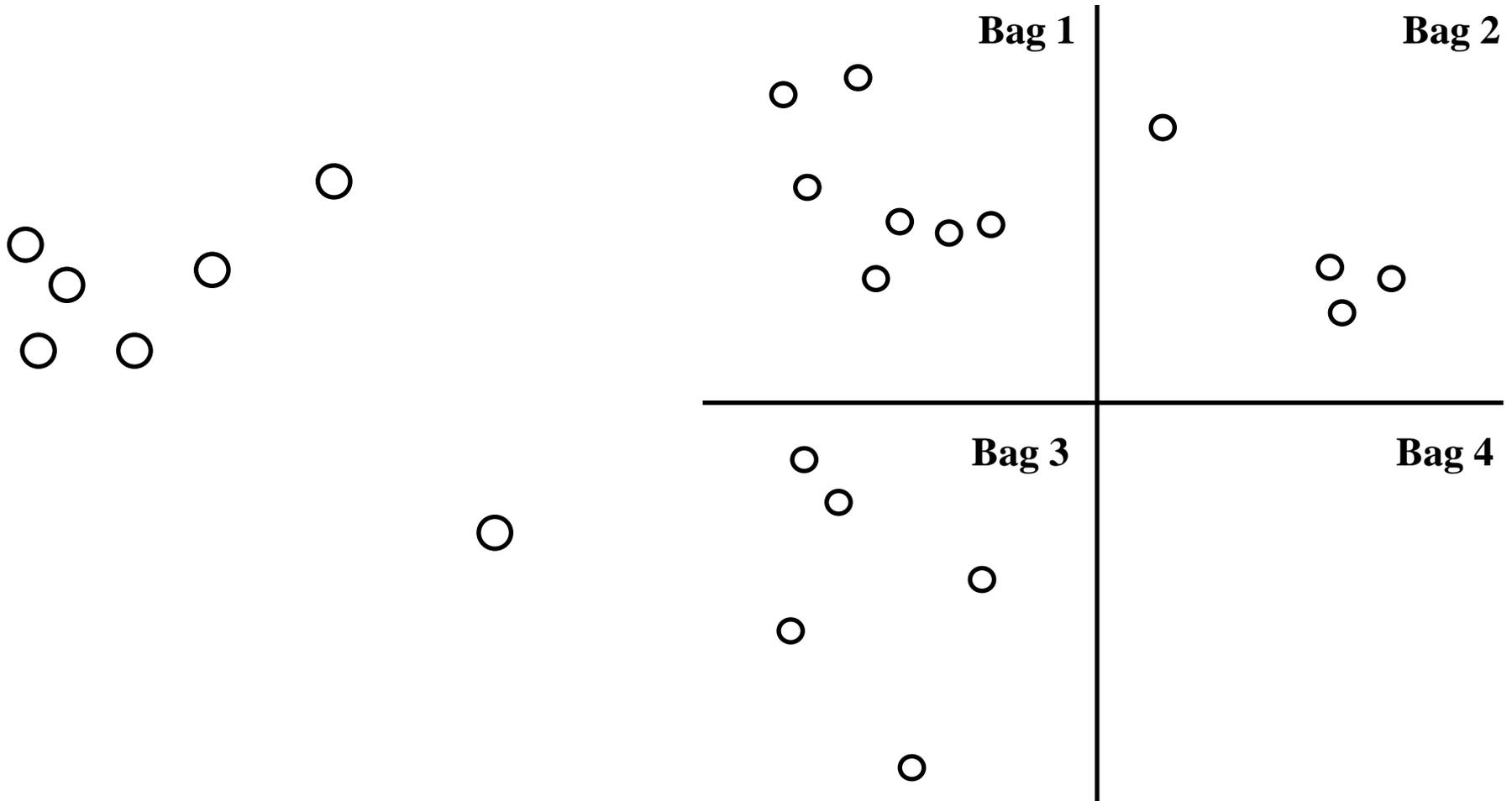
Learning with label proportions



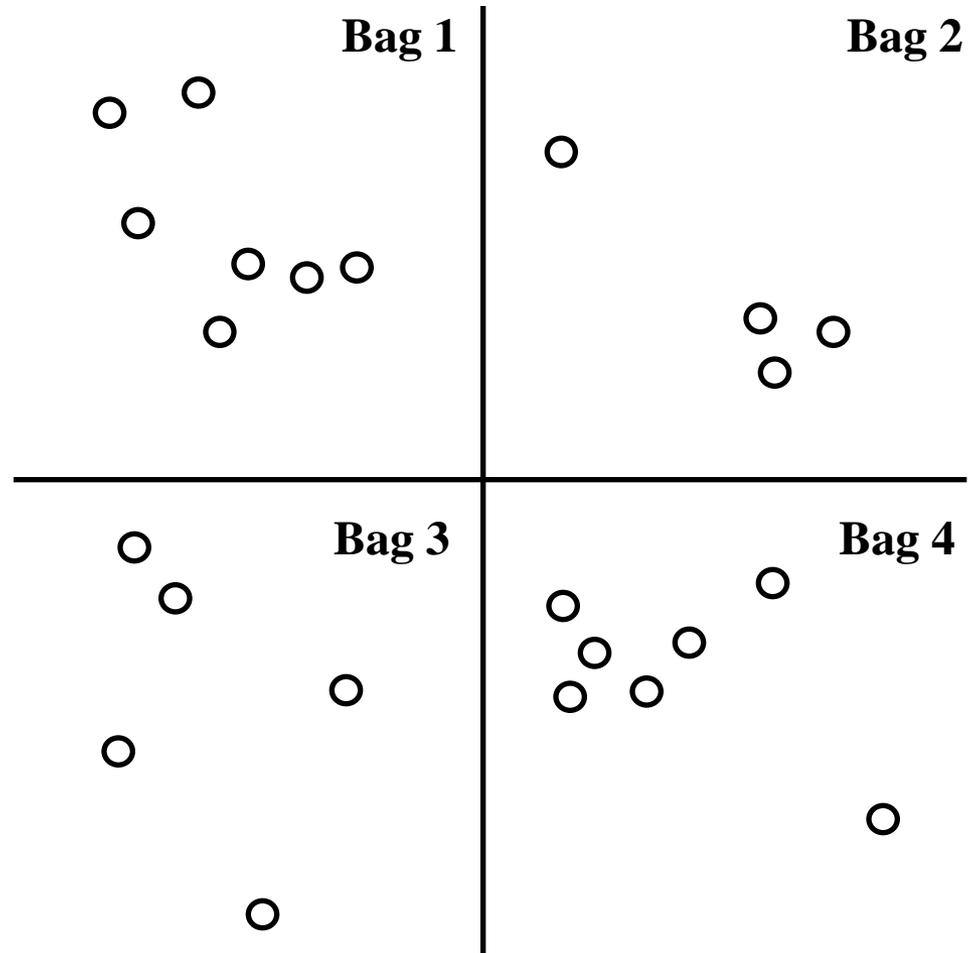
Learning with label proportions



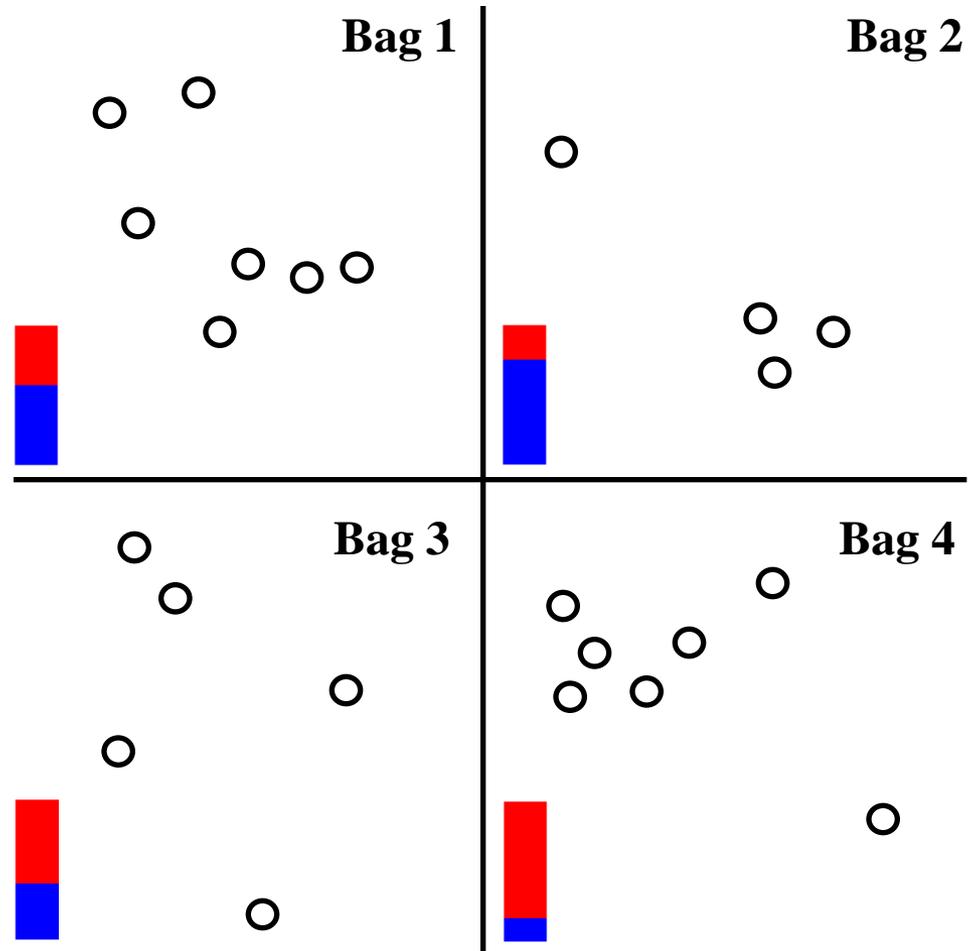
Learning with label proportions



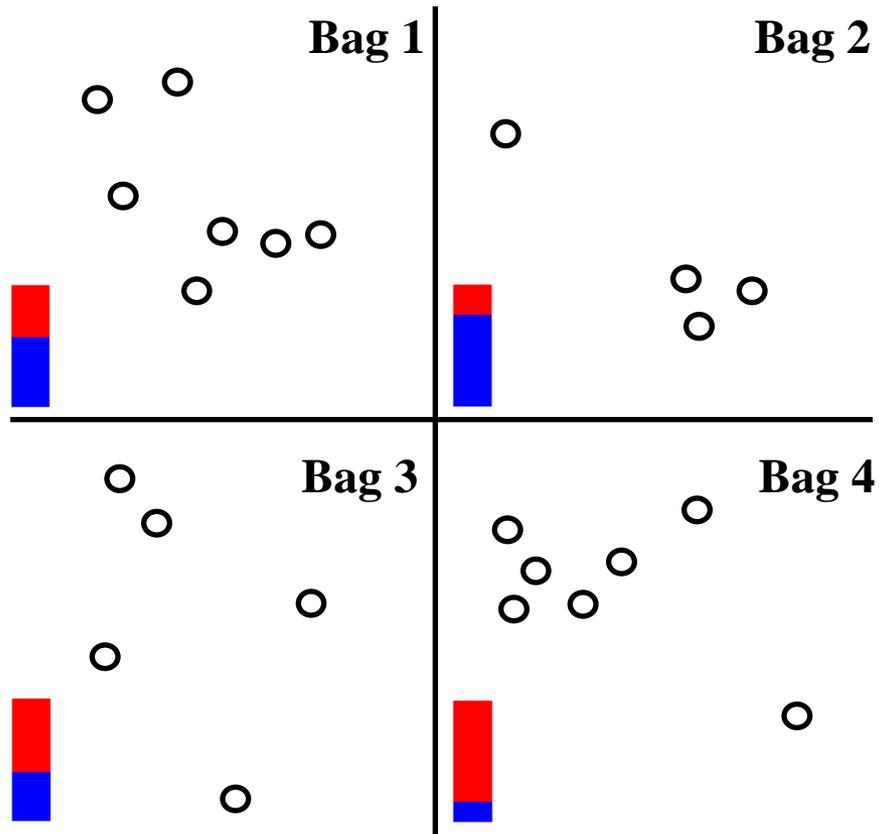
Learning with label proportions



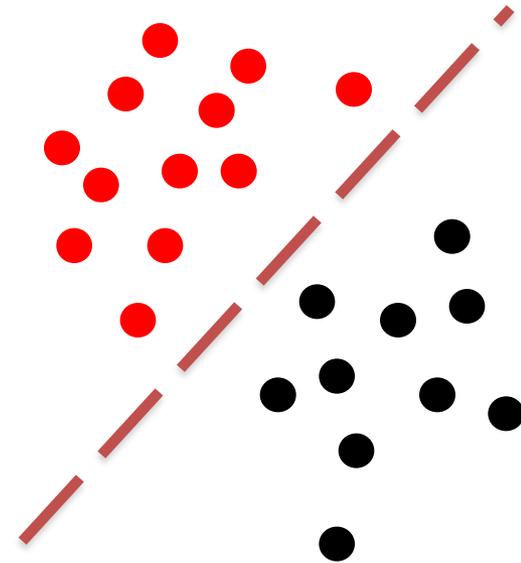
Learning with label proportions



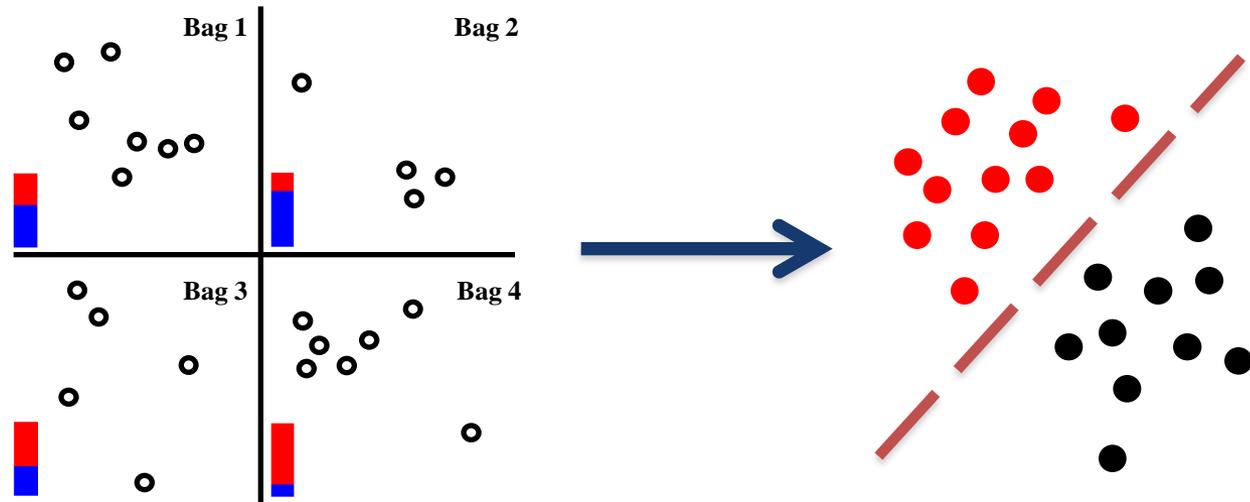
Learning with label proportions



Learning
with Label
Proportions

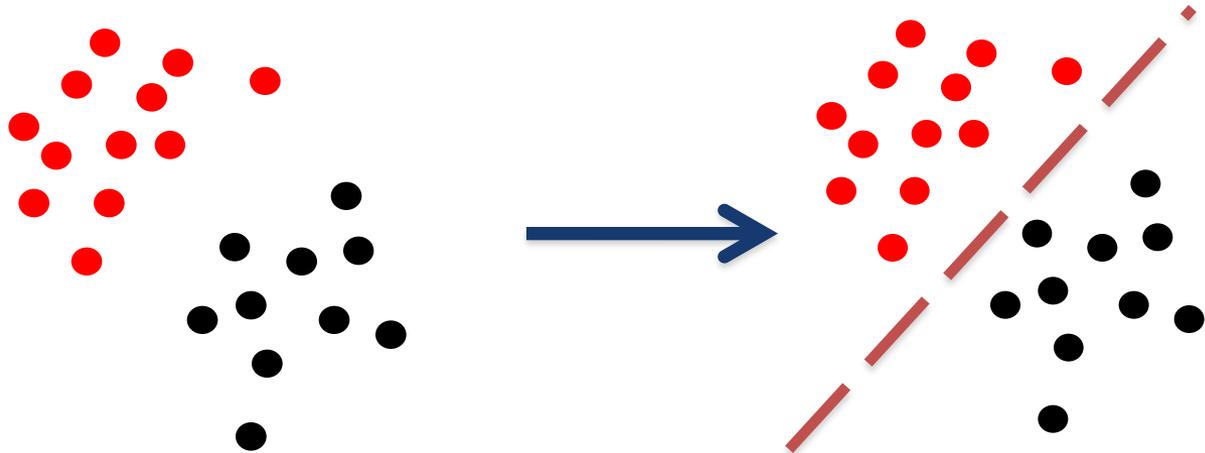


Learning with label proportions



Learning with Label Proportions

Supervised Learning



Applications

- Healthcare
 - Medical record
- Social science
 - Voting behavior
 - Census data
 - Energy consumption
 - Marketing
- Computer vision
 - Facial attributes (“90 percent of Asians have black hair”)



Applications

- Healthcare
 - Medical record
- Social science
 - Voting behavior
 - Census data
 - Energy consumption
 - Marketing
- Computer vision
 - Facial attributes (“90 percent of Asians have black hair”)

Privacy issues

Applications

- Healthcare
 - Medical record
- Social science
 - Voting behavior
 - Census data
 - Energy consumption
 - Marketing
- Computer vision
 - Facial attributes (“90 percent of Asians have black hair”)

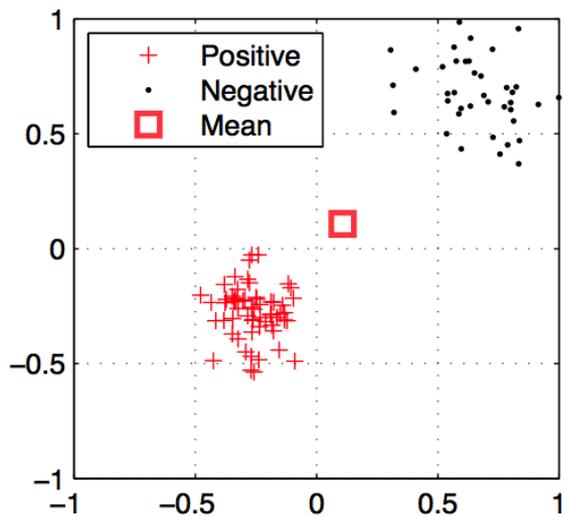
Easier to get
label proportions

Outline

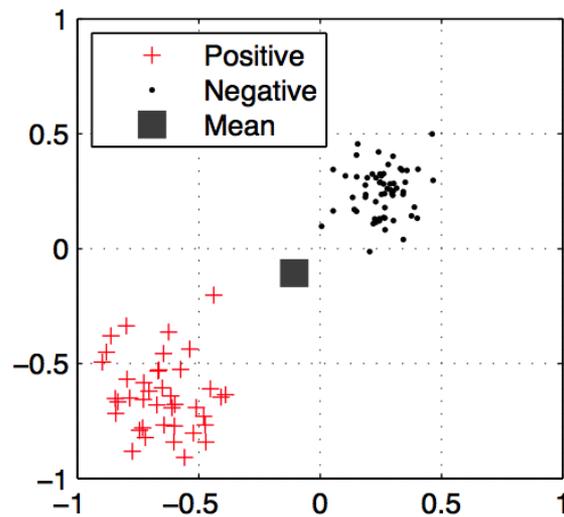
- Learning Setting and Applications
- **Related Works**
- Formulation
- Algorithms
- Experiments

Related Works

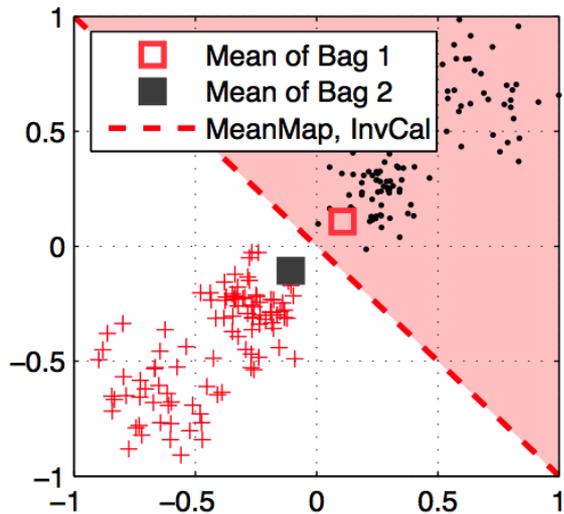
- Related learning settings: semi-supervised learning, clustering, multi-instance learning etc.
- Learning with label proportions: former works rely heavily on the “mean of each bag”
 - MeanMap (Quadranto et al., 2009)
 - exponential model
 - class-conditional distribution of data is independent of the bag
 - Inverse Calibration (Rueping, 2011)
 - large-margin regression
 - mean of each bag has a soft label corresponding to its label proportion



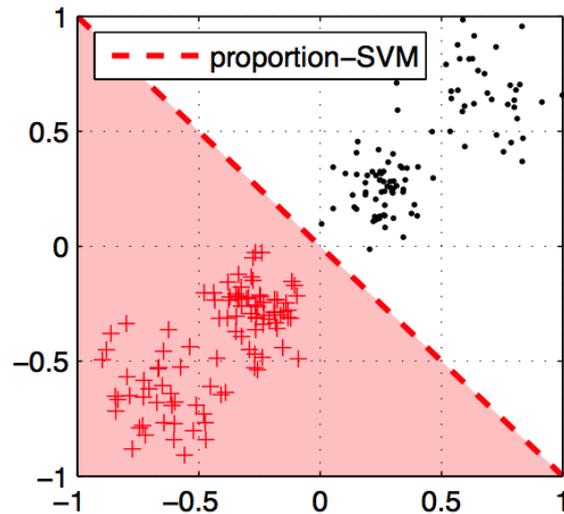
(a) Bag 1, with $p_1 = 0.6$



(b) Bag 2, with $p_2 = 0.4$



(c) MeanMap and InvCal with 0% accuracy.



(d) \propto SVM with 100% accuracy.

Contributions

- Introduce α SVM which **explicitly** models the **unknown instance labels**.
- Alleviates the need for making restrictive assumptions on the data.
- Two optimization algorithms based on alternating minimization and convex relaxation.
- Outperforms existing methods under various settings/datasets.

Outline

- Learning Setting and Applications
- Related Works
- **Formulation**
- Algorithms
- Experiments

Formulation (Learning Setting)

- The training set $\{\mathbf{x}_i\}_{i=1}^N$ is given in the form of K non-overlapping bags:

$$\{\mathbf{x}_i | i \in \mathcal{B}_k\}_{k=1}^K, \quad \cup_{k=1}^K \mathcal{B}_k = \{1 \cdots N\}.$$

- The k -th bag is with label proportion p_k :

$$\forall_{k=1}^K, \quad p_k := \frac{|\{i | i \in \mathcal{B}_k, y_i^* = 1\}|}{|\mathcal{B}_k|}.$$

$y_i^* \in \{1, -1\}$: the *unknown* ground-truth label of \mathbf{x}_i , $\forall_{i=1}^N$

Formulation

- Prediction model:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \varphi(\mathbf{x}) + b).$$

- Explicitly model the unknown instance labels as

$$\mathbf{y} = (y_1, \dots, y_N)^T. \quad y_i \in \{1, -1\}, \quad \forall_{i=1}^N.$$

- The label proportion of the k -th bag can be modeled as

$$\tilde{p}_k(\mathbf{y}) = \frac{|\{i | i \in \mathcal{B}_k, y_i = 1\}|}{|\mathcal{B}_k|} = \frac{\sum_{i \in \mathcal{B}_k} y_i}{2|\mathcal{B}_k|} + \frac{1}{2}.$$

α SVM Formulation

- Large-margin framework:

$$\begin{aligned} \min_{\mathbf{y}, \mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N L(y_i, \mathbf{w}^T \varphi(\mathbf{x}_i) + b) + C_p \sum_{k=1}^K L_p(\tilde{p}_k(\mathbf{y}), p_k) \\ \text{s.t.} \quad & \forall_{i=1}^N, \quad y_i \in \{-1, 1\}. \end{aligned}$$

α SVM Formulation

- Large-margin framework:

$$\begin{aligned} \min_{\mathbf{y}, \mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N L(y_i, \mathbf{w}^T \varphi(\mathbf{x}_i) + b) + C_p \sum_{k=1}^K L_p(\tilde{p}_k(\mathbf{y}), p_k) \\ \text{s.t.} \quad & \forall_{i=1}^N, \quad y_i \in \{-1, 1\}. \end{aligned}$$

- Generalizes the classic SVM.
- Naturally spans supervised/semi-supervised learning and clustering.

α SVM Formulation

- Large-margin framework:

$$\begin{aligned} \min_{\mathbf{y}, \mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N L(y_i, \mathbf{w}^T \varphi(\mathbf{x}_i) + b) + C_p \sum_{k=1}^K L_p(\tilde{p}_k(\mathbf{y}), p_k) \\ \text{s.t.} \quad & \forall_{i=1}^N, \quad y_i \in \{-1, 1\}. \end{aligned}$$

- Generalizes the classic SVM.
- Naturally spans supervised/semi-supervised learning and clustering.

However:

- A non-convex integer programming problem.

Outline

- Learning Setting and Applications
- Related Works
- Formulation
- **Algorithms**
- Experiments

The alter- α SVM Algorithm

- For a fixed \mathbf{y} , the optimization *w.r.t* \mathbf{w} and b is an SVM problem.
- When \mathbf{w} and b are fixed:

$$\begin{aligned} \min_{\mathbf{y}} \quad & \sum_{i=1}^N L(y_i, \mathbf{w}^T \varphi(\mathbf{x}_i) + b) + \frac{C_p}{C} \sum_{k=1}^K L_p(\tilde{p}_k(\mathbf{y}), p_k) \\ \text{s.t.} \quad & \forall_{i=1}^N, \quad y_i \in \{1, -1\}. \end{aligned}$$

The alter- α SVM Algorithm

$$\begin{aligned} \min_{\mathbf{y}} \quad & \sum_{i=1}^N L(y_i, \mathbf{w}^T \varphi(\mathbf{x}_i) + b) + \frac{C_p}{C} \sum_{k=1}^K L_p(\tilde{p}_k(\mathbf{y}), p_k) \\ \text{s.t.} \quad & \forall_{i=1}^N, \quad y_i \in \{1, -1\}. \end{aligned}$$

- Consider each bag separately.
- For the k -th bag: sorting, $\mathcal{O}(|\mathcal{B}_k| \log |\mathcal{B}_k|)$ time).

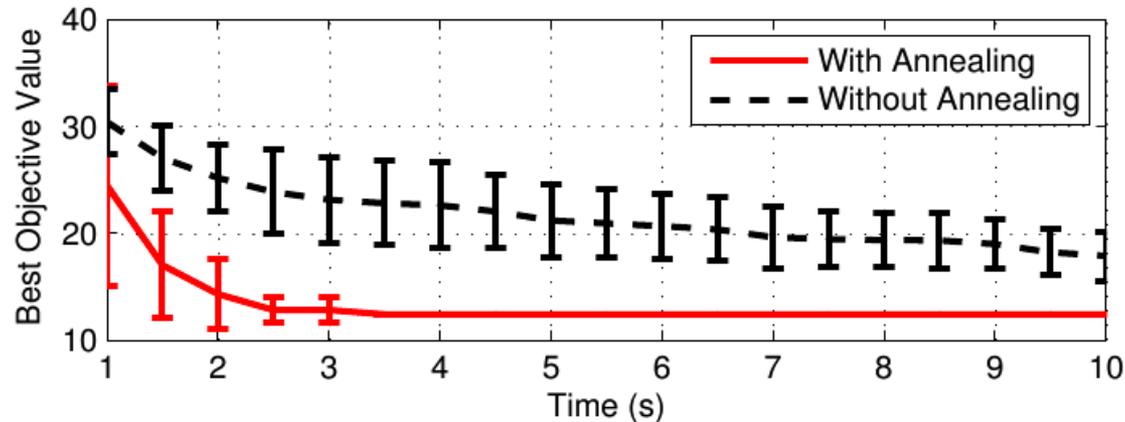
Proposition:

The above can be solved in $\mathcal{O}(N \log(J))$ time, $J = \max_{k=1 \dots K} |\mathcal{B}_k|$.

The alter- α SVM Algorithm

To alleviate the problem of local solutions:

- Multiple initializations.
- An additional annealing loop to gradually increase C .



The conv- α SVM Algorithm

- Does not require multiple initializations.
- Motivated by large-margin clustering (Xu et al., 2004) (Li et al., 2009).

The conv- α SVM Algorithm

- Reformulation:

$$\min_{\mathbf{y} \in \mathcal{Y}, \mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N L(y_i, \mathbf{w}^T \varphi(\mathbf{x}_i))$$
$$\mathcal{Y} = \left\{ \mathbf{y} \mid |\tilde{p}_k(\mathbf{y}) - p_k| \leq \epsilon, y_i \in \{-1, 1\}, \forall_{k=1}^K \right\}$$

- Write the inner problem as its dual (with hinge loss):

$$\min_{\mathbf{y} \in \mathcal{Y}} \max_{\boldsymbol{\alpha} \in \mathcal{A}} -\frac{1}{2} \boldsymbol{\alpha}^T (\mathcal{K} \odot \mathbf{y} \mathbf{y}^T) \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{1}$$

$$\boldsymbol{\alpha} \in \mathbb{R}^N$$

$$\mathcal{A} = \{ \boldsymbol{\alpha} \mid 0 \leq \boldsymbol{\alpha} \leq C \}$$

The conv- α SVM Algorithm

$$\min_{\mathbf{y} \in \mathcal{Y}} \max_{\boldsymbol{\alpha} \in \mathcal{A}} -\frac{1}{2} \boldsymbol{\alpha}^T (\mathcal{K} \odot \mathbf{y}\mathbf{y}^T) \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{1}$$

$$\boldsymbol{\alpha} \in \mathbb{R}^N$$
$$\mathcal{A} = \{\boldsymbol{\alpha} \mid 0 \leq \boldsymbol{\alpha} \leq C\}$$

- Convex in $\mathbf{M} := \mathbf{y}\mathbf{y}^T$.
- Relax the feasible space of \mathbf{M} to get a convex problem.

$$\mathcal{M}_0 = \{\mathbf{y}\mathbf{y}^T \mid \mathbf{y} \in \mathcal{Y}\}$$

↓ Relaxation

$$\mathcal{M} = \left\{ \sum_{\mathbf{y} \in \mathcal{Y}} \mu_{(\mathbf{y})} \mathbf{y}\mathbf{y}^T \mid \boldsymbol{\mu} \in \mathcal{U} \right\},$$
$$\mathcal{U} = \left\{ \boldsymbol{\mu} \mid \sum_{\mathbf{y} \in \mathcal{Y}} \mu_{(\mathbf{y})} = 1, \mu_{(\mathbf{y})} \geq 0 \right\}$$

The conv- α SVM Algorithm

- Solving the relaxed M is identical to finding μ :

$$\min_{\mu \in \mathcal{U}} \max_{\alpha \in \mathcal{A}} -\frac{1}{2} \alpha^T \left(\sum_{\mathbf{y} \in \mathcal{Y}} \mu_{(\mathbf{y})} \mathcal{K} \odot \mathbf{y} \mathbf{y}^T \right) \alpha + \alpha^T \mathbf{1}.$$

- Multiple Kernel Learning (MKL).
- $|\mathcal{Y}|$ is very large. Not tractable to solve directly.
- Primal variables \rightarrow dual constraints.
- Cutting plane method (Li et al., 2009) (Joachims et al., 2009).

Outline

- Learning Setting and Applications
- Related Works
- Formulation
- Algorithms
- Experiments

Experiments

- Performance of different techniques on 12 datasets from the UCI/LibSVM repository.

Dataset	Size	Attributes	Classes
heart	270	13	2
heart-c	303	13	2
colic	366	22	2
vote	435	16	2
breast-cancer	683	10	2
australian	690	14	2
credit-a	690	15	2
breast-w	699	9	2
ala	1,605	119	2
dna	2,000	180	3
satimage	4,435	36	6
cod-rna.t	271,617	8	2

- Follow the experimental setting of (Rueping, 2011):
 - Random bag generation (with different bag sizes). Performance of 5-fold cross validation.
 - Linear and RBF kernels.



Experiments

Dataset	Method	2	4	8	16	32	64
heart	MeanMap	82.69±0.71	80.80±0.97	79.65±0.82	79.44±1.21	80.03±2.05	77.26±0.85
	InvCal	83.15±0.56	81.06±0.70	80.26±1.32	79.61±3.84	76.36±3.72	73.90±3.00
	alter- α SVM	83.15±0.85	82.89±1.30	81.51±0.54	80.07±1.21	78.10±0.96	78.63±1.85
colic	MeanMap	82.45±0.88	81.38±1.26	81.71±1.16	79.94±1.33	76.36±2.43	77.84±1.69
	InvCal	82.20±0.61	81.20±0.87	81.17±1.74	78.59±2.19	74.09±5.26	72.81±4.80
	alter- α SVM	83.28±0.50	82.97±0.39	82.03±0.44	81.62±0.46	81.53±0.21	81.39±0.34
vote	MeanMap	91.15±0.33	90.32±0.62	91.54±0.20	90.28±1.63	89.58±1.09	89.38±1.33
	InvCal	95.68±0.19	94.77±0.44	93.95±0.43	93.03±0.37	87.79±1.64	86.63±4.74
	alter- α SVM	95.80±0.20	95.54±0.25	94.88±0.94	92.44±0.60	90.72±1.11	90.93±1.30
australian	MeanMap	88.97±0.72	85.88±0.34	85.34±1.01	85.36±2.04	83.12±1.52	80.58±5.41
	InvCal	86.06±0.30	86.11±0.26	86.32±0.45	84.13±1.62	82.73±1.70	81.87±3.29
	alter- α SVM	85.74±0.22	85.71±0.21	86.26±0.61	85.65±0.43	83.63±1.83	83.62±2.21
dna-1	MeanMap	91.53±0.25	90.58±0.34	86.00±1.04	80.77±3.69	77.35±3.59	68.47±4.30
	InvCal	89.32±3.39	92.73±0.53	87.99±1.65	81.05±3.14	74.77±2.95	67.75±3.86
	alter- α SVM	95.67±0.40	94.65±0.52	93.71±0.47	92.52±0.63	91.85±1.42	90.64±1.32
satimage-2	MeanMap	97.08±0.48	96.82±0.38	96.50±0.43	96.45±1.16	95.51±0.73	94.26±0.22
	InvCal	97.53±1.33	98.33±0.13	98.38±0.23	97.99±0.14	96.97±0.18	94.47±0.97
	alter- α SVM	98.83±0.36	98.69±0.37	98.62±0.27	98.51±0.13	98.47±0.13	98.47±0.13

Horizontal: bag size

Dataset	Method	2	4	8	16	32	64
australian	MeanMap	85.97±0.72	85.88±0.34	85.34±1.01	83.36±2.04	83.12±1.52	80.58±5.41
	InvCal	86.06±0.30	86.11±0.26	86.32±0.45	84.13±1.62	82.73±1.70	81.87±3.29
	alter- α SVM	85.74±0.22	85.71±0.21	86.26±0.61	85.65±0.43	83.63±1.83	83.62±2.21
	conv- α SVM	85.97±0.53	86.46±0.23	85.30±0.70	84.18±0.53	83.69±0.78	82.98±1.32
dna-1	MeanMap	91.53±0.25	90.58±0.34	86.00±1.04	80.77±3.69	77.35±3.59	68.47±4.30
	InvCal	89.32±3.39	92.73±0.53	87.99±1.65	81.05±3.14	74.77±2.95	67.75±3.86
	alter- α SVM	95.67±0.40	94.65±0.52	93.71±0.47	92.52±0.63	91.85±1.42	90.64±1.32
	conv- α SVM	93.36±0.53	86.75±2.56	81.03±3.58	75.90±4.56	76.92±5.91	77.94±2.48

- Our methods outperform MeanMap and InvCal.
- The gains from α SVM are typically even more significant when the bag size is large.

-- on the dna-1 dataset, with RBF kernel and bag size 64, alter- α SVM outperforms the former works by 22%.

- Also shown in the paper: less sensitive to bag proportion variations.

Conclusion

- ∞ SVM
- Two optimization algorithms
- State-of-the-art result

Future works/ Open issues

- Robustness to proportion noise.
- Bag generation, bags with overlappings.
- ∞ SVM for semi-supervised learning, and learning with label errors.