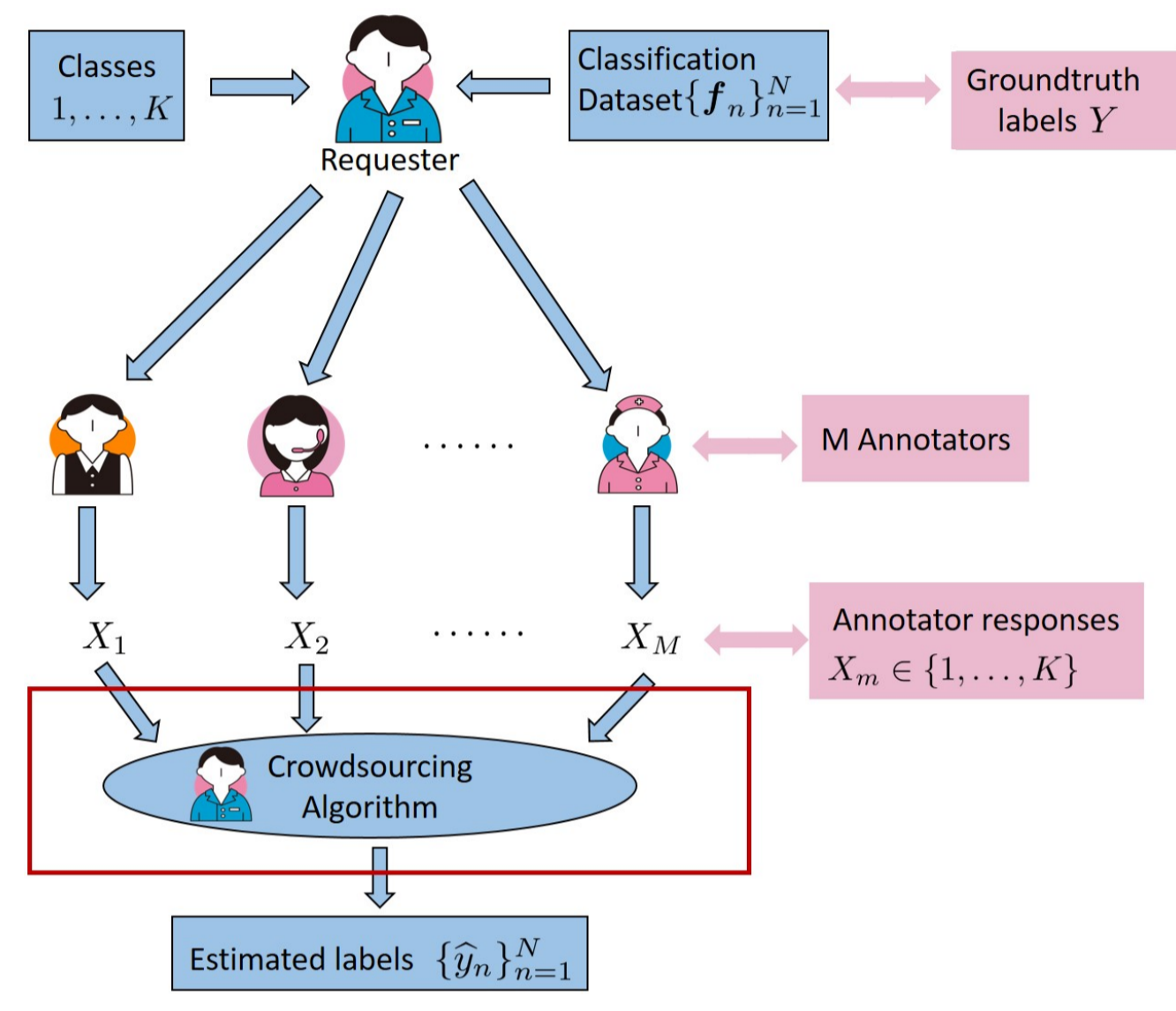


## Crowdsourcing

### ► Crowdsourcing techniques

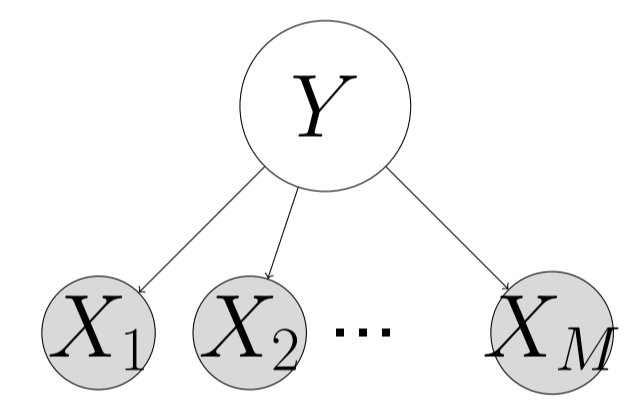
- employ a group of annotators to label the data items
- integrate the acquired labels



### Dawid-Skene Model

### ► Dawid and Skene [1979] formulated label integration as model identification

### ► Under Dawid & Skene (D&S) model,



$$\Pr(X_1 = k_1, \dots, X_M = k_M) = \sum_{k=1}^K \Pr(Y = k) \prod_{m=1}^M \Pr(X_m = k_m | Y = k)$$

### ► Define the confusion matrix $A_m \in \mathbb{R}^{K \times K}$ for each annotator and the prior probability vector $\lambda \in \mathbb{R}^K$ such that

$$A_m(k_m, k) := \Pr(X_m = k_m | Y = k) \quad \lambda(k) := \Pr(Y = k)$$

### ► One can build a MAP estimator for $y_n$ after identifying $A_m$ 's and $\lambda$

### Prior Approaches with Dawid-Skene Model

### ► Dawid-Skene (D&S) Model & EM Algorithm [Dawid and Skene, 1979]:

- No model identifiability & algorithm tractability
- Bayesian Methods [Whitehill et al., 2009; Zhou et al., 2012]:
  - Extended D&S model considering "item difficulty" and "annotator ability"
  - No model identifiability
- Tensor Methods [Zhang et al., 2016; Traganitis et al., 2018]:

### ► Using third-order co-occurrences of annotator responses, for e.g.,

$$\Pr(X_m = k_m, X_\ell = k_\ell, X_j = k_j)$$

### ► Established model identifiability

- High sample complexity due to third-order statistics
- High computational cost from the tensor decomposition

### ► Coupled NMF (CNMF)-based Approach [Ibrahim et al., 2019]:

$$\Pr(X_m = k_m, X_j = k_j) = \sum_{k=1}^K \underbrace{\Pr(Y = k)}_{\lambda(k)} \underbrace{\Pr(X_m = k_m | Y = k)}_{A_m(k_m, k)} \underbrace{\Pr(X_j = k_j | Y = k)}_{A_j(k_j, k)}$$

- less sample complexity compared to third-order statistics

### ► If annotators $m$ and $j$ co-label some items, $R_{m,j}$ can be estimated via sample averaging

## CNMF Approach - A Deeper Look

### ► The CNMF criterion in [Ibrahim et al., 2019]:

$$\text{find } \{A_m\}_{m=1}^M, \lambda$$

$$\text{s.t. } R_{m,j} = A_m D A_j^T, (m, j) \in \Omega, \leftarrow \text{observed set}$$

$$A_m \geq 0, \mathbf{1}^T A_m = \mathbf{1}^T, \mathbf{1}^T \lambda = 1, \lambda \geq 0.$$

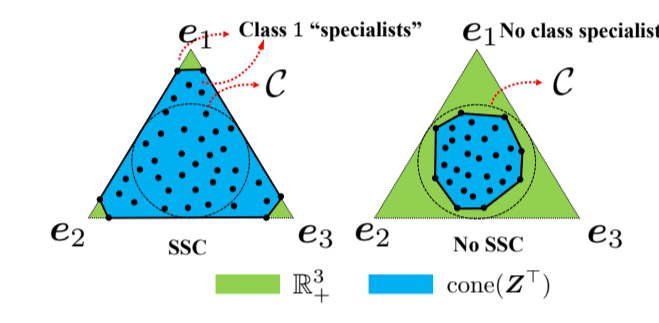
### ► Identifiability under the assumption that there exist two subsets of annotators $\mathcal{P}_1$ and $\mathcal{P}_2$ , where $\mathcal{P}_1 \cap \mathcal{P}_2 = \emptyset$ and $\mathcal{P}_1 \cup \mathcal{P}_2 \subseteq [M]$ ,

$$H^{(1)} := [A_{m_1}^T, \dots, A_{m_{|\mathcal{P}_1|}}^T]^T, \quad H^{(2)} := [A_{j_1}^T, \dots, A_{j_{|\mathcal{P}_2|}}^T]^T,$$

### s.t. $H^{(1)}$ and $H^{(2)}$ satisfy the sufficiently scattered condition (SSC)

### Def. 1: (SSC) [Fu et al., 2015]

Any nonnegative matrix  $Z \in \mathbb{R}_+^{I \times K}$  satisfies the SSC if the conic hull of  $Z^T$  (i.e.,  $\text{cone}(Z^T)$ ) satisfies  $\mathcal{C} \subseteq \text{cone}\{Z^T\}$  where  $\mathcal{C} = \{x \in \mathbb{R}^K \mid x^T \mathbf{1} \geq \sqrt{K-1} \|x\|_2\}$ .



### ► A row of $H^{(i)}$ (i.e., a row of $A_m$ ) close to $k$ th unit vector implies $A_m(k, k) \approx 1$ and $A_m(k, k_m) \approx 0, k_m \neq k$ (class specialists),

i.e., annotator  $m$  rarely confuses data from other classes with those from class  $k$

### ► Identifiability Challenge:

Both  $H^{(1)}$  and  $H^{(2)}$  satisfy the SSC  $\implies$  the disjoint  $\mathcal{P}_1$  and  $\mathcal{P}_2$  both contain "class specialists" for all  $K$  classes (somewhat restrictive condition)

### ► Computational Challenges:

- The CNMF criterion in [Ibrahim et al., 2019] is handled using KL-divergence based model fitting problem with constraints (hardly scalable)
- Unclear convergence guarantee even if there is no noise
- Unclear identifiability guarantee when there is noise

## Proposed Approach - SymNMF Framework

### ► Assume that all $R_{m,j} = A_m D A_j^T$ are available for all $m, j \in [M]$

### Symmetric Non-negative Matrix Factorization (SymNMF)

$$X = \begin{bmatrix} R_{1,1} & \dots & R_{1,M} \\ \vdots & \ddots & \vdots \\ R_{M,1} & \dots & R_{M,M} \end{bmatrix} = \underbrace{[A_1^T, \dots, A_M^T]^T}_{H} D^{1/2} \underbrace{[A_1^T, \dots, A_M^T]}_{H^T}$$

### ► If $H$ satisfies SSC, the SymNMF model is unique [Huang et al., 2014], i.e., $A_m$ 's and $\lambda$ can be identified upto common column permutations

### ► SSC of $H \implies$ only one set of "class specialists" is needed

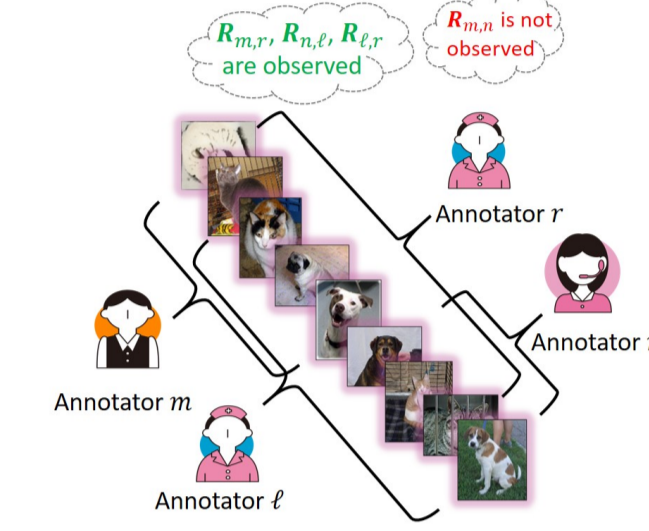
- the CNMF framework needs two disjoint sets of annotators  $\mathcal{P}_1$  and  $\mathcal{P}_2$  both contain "class specialists" for all  $K$  classes
- much easier to satisfy compared to the CNMF framework case

### ► The challenge in SymNMF framework is that many $R_{m,j}$ 's may be missing:

- $R_{m,m} = A_m D A_m^T, \forall m$  do not have physical meaning and thus cannot be observed
- if annotators  $m, j$  never co-labeled any items,  $R_{m,j}$  is missing

## Designated Annotators-based Imputation

### ► In crowdsourcing, some annotators may be designated to co-label items with other annotators.



$$1. C \leftarrow [R_{m,r}^T, R_{\ell,r}^T]^T$$

$$2. C \xrightarrow{\text{thin SYD}} [U_m^T, U_\ell^T]^T \Sigma_{m,\ell,r} V_r^T$$

$$3. R_{m,n} \leftarrow U_m U_\ell^{-1} R_{n,\ell}^T$$

The diagonal blocks  $R_{m,m}$ 's can be estimated by observing  $R_{m,\ell}, R_{m,r}$ , and  $R_{\ell,r}$

### Theorem 1: (Informal)

Assume that  $R_{m,r}, R_{n,\ell}$  and  $R_{\ell,r}$  are estimated using at least  $S$  items and that  $\kappa(A_m) \leq \gamma$  and  $\text{rank}(A_m) = \text{rank}(D) = K$  for all  $m$ . Suppose that  $S$  is above certain threshold. Then, any unobserved  $R_{m,n}$  can be estimated via (1)-(3), with probability of at least  $1 - \delta$  such that  $\|R_{m,n} - \hat{R}_{m,n}\|_F = O(K^2 \gamma^3 \sqrt{\log(1/\delta)/S})$ .

### ► What if we do not have designated annotators?

## Robust Co-occurrence Imputation Criterion

$$\text{minimize } \sum_{(m,j) \in \Omega} \|R_{m,j} - U_m U_j^T\|_F$$

$$\text{subject to } \|U_m\|_F \leq D_{\text{(certain upper bound)}}, \forall m$$

- $\hat{R}_{m,j}$ 's are estimated using unequal no of samples
- the formulation is robust under such unbalanced estimates

### ► block $\ell_2/\ell_1$ -mixed norm based criterion

### Theorem 2: Stability under Finite Samples

Assume that  $\hat{R}_{m,j}$ 's are estimated with  $S_{m,j}$  samples,  $\forall (m, j) \in \Omega$  and each  $\hat{R}_{m,j}$  is observed with the same probability. Let  $\{U_m^*, U_j^*\}$  be any optimal solution of the above. Then we have

$$\frac{1}{L} \sum_{m < j} \|U_m^* (U_j^*)^T - R_{m,j}\|_F \leq C \sqrt{\frac{MK^2 \log(M)}{|\Omega|} + \left(\frac{1}{|\Omega|} + \frac{1}{L}\right) \sum_{(m,j) \in \Omega} \frac{1 + \sqrt{M}}{S_{m,j}}}$$

with probability of at least  $1 - 3 \exp(-M)$ , where  $L = M(M-1)/2$  and  $C > 0$ .

## Shifted ReLU Empowered SymNMF

### Assuming that $X$ is observed after co-occurrence imputation:

$$X = H H^T \xrightarrow{\text{square root decomp.}} X \rightarrow U U^T \implies U = H Q^T, Q \text{ is orthogonal}$$

### Estimation Criterion:

$$\text{minimize } \|H - U Q\|_F^2$$

$$\text{subject to } H \geq 0, Q^T Q = I$$

### Proposed Algorithm:

$$H_{(t+1)} \leftarrow \text{ReLU}_{\alpha(t)}(U Q_{(t)}) \text{ (Orth. proj. of each element of } U Q_{(t)} \text{ to } [\alpha(t), +\infty))$$

$$W_{(t+1)} \Sigma_{(t+1)} V_{(t+1)}^T \leftarrow \text{svd}(H_{(t+1)}^T U)$$

$$Q_{(t+1)} \leftarrow V_{(t+1)} W_{(t+1)}^T \text{ (Procrustes projection)}$$

### ► reminiscent of the SymNMF algorithm proposed in [Huang et al., 2014]:

- always uses  $\alpha(t) = 0$ ; convergence w/o noise is unclear
- Convergence analysis for SymNMF algorithms is challenging
- most existing SymNMF works showed only stationary point convergence [Huang et al., 2014; He et al., 2011]

## Convergence of the Proposed SymNMF Algorithm

### Theorem 3: (Informal)

Consider  $\hat{U} = H Q^T + N$ . Denote  $\nu = \|N\|_F, \sigma = \|H\|_F, h_{(t)} = \|H_{(t)} - H H^T\|_F^2$  and  $q_{(t)} = \|Q_{(t)} - Q H^T\|_F^2$ , where  $H^T$  is any permutation matrix. Under the assumptions that,  $H$  is full rank and sparse; the energy of range space of  $H$  is well spread over its rows; the noise term  $\nu$  and the initial error  $q_{(0)}$  are small enough; there exists  $\alpha_{(t)} = \alpha > 0, \eta > 0$  and  $0 < \rho < 1$  such that with high probability,

$$q_{(t)} \leq \rho q_{(t-1)} + O(K \sigma^2 \nu^2), \quad h_{(t)} \leq 2\eta \sigma^2 q_{(t-1)} + 2\nu^2 \leftarrow \text{linear convergence}$$

- Shifted ReLU operator is crucial for guaranteeing the convergence
- The rate parameter  $\rho$  is smaller (faster convergence) if  $H$  is sparser

## Experiment Results

### ► Experiments - UCI Data:

- Each annotator (MATLAB classifiers) is allowed to label an item with prob.  $p_m \in (0, 1]$ ; randomly choosing two annotators and letting them label with higher prob. (i.e.,  $p_d$ )

Table: UCI Connect4 dataset ( $N = 20,561, M = 10, K = 3$ )

Algorithms	$p_m = 0.3$	$p_m \in (0.3, 0.5), p_d = 0.8$	$p_m \in (0.5, 0.7), p_d = 0.8$	Time(s)
RobSymNMF	33.26	33.06	32.16	0.142
RobSymNMF-EM	34.27	33.20	32.11	0.191
DesSymNMF	33.45	32.18	31.42	0.061
DesSymNMF-EM	33.94	32.50	31.40	0.128
CNMF	36.26	39.55	34.70	4.741
TensorADMM	36.20	34.34	35.18	5.183
Spectral-D&S	64.28	66.95	71.97	20.388
MV-EM	34.14	34.17	34.19	0.107
MinimaxEntropy	36.20	36.17	35.46	27.454
Majority Voting	37.76	36.88	36.75	-

### ► Experiments - Amazon Mechanical Turk (AMT) Data:

Table: AMT datasets "RTE" and "TREC"

Algorithms	RTE		TREC	
	( $N = 800, M = 164, K = 2$ )	( $N = 19,033, M = 762, K = 2$ )	( $N = 19,033, M = 762, K = 2$ )	( $N = 19,033, M = 762, K = 2$ )
	Error (%)	Time (s)	Error (%)	Time (s)
RobSymNMF	7.25	2.31	30.68	64.99
RobSymNMF-EM	7.12	2.4	29.62	67.39
DesSymNMF	13.87	3.32	36.75	71.31
DesSymNMF-EM	7.25	3.43	29.36	72.13
CNMF	7.12	18.12	29.84	536.86
TensorADMM	N/A	N/A	N/A	N/A
Spectral-D&S	7.12	6.34	29.58	919.98
MV-EM	7.25	0.09	30.02	3.12
MinimaxEntropy	7.5	6.4	30.89	356.32
Majority Voting	10.31	N/A	34.85	N/A

## References

- A. P. Dawid and A. M. Skene. Maximum likelihood estimation of observer error-rates using the EM algorithm. Applied statistics, pp 20-28, 1979.
- S. Ibrahim, X.Fu, N. Kargas, and K. Huang. Crowdsourcing via pairwise co-occurrences: Identifiability and algorithms. In Advances in NeurIPS, vol 32, pp 7847-7857, 2019.