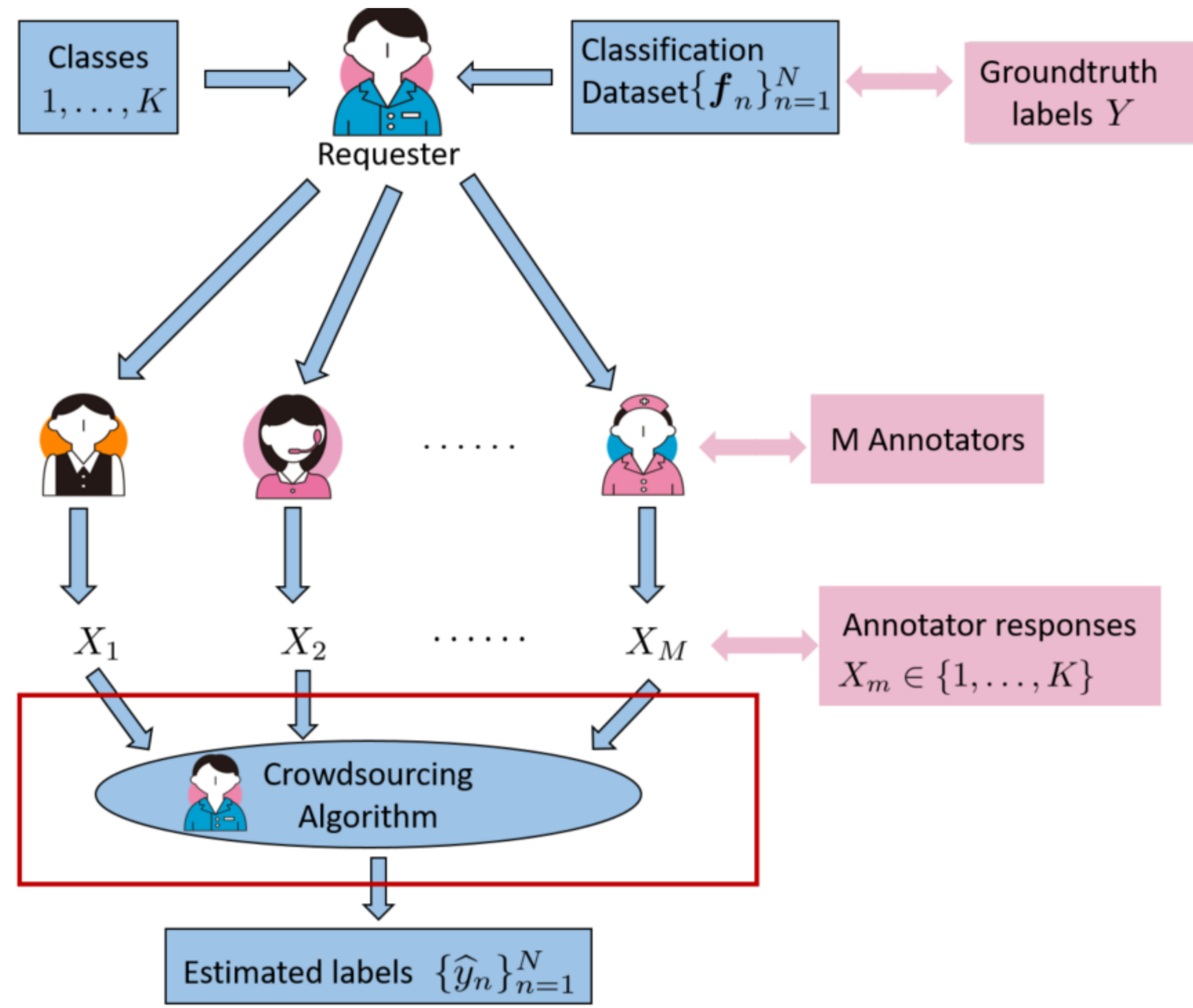


Data Labeling and Crowdsourcing

- ▶ Massive labeled data is a key performance booster of deep networks.
- ▶ **Crowdsourcing** is widely used for data labeling.



Dawid-Skene Model

- ▶ The **confusion matrix** $A_m \in \mathbb{R}^{K \times K}$ for each annotator m and the **prior probability vector** $d \in \mathbb{R}^K$ are the Dawid-Skene model parameters.

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

$$d(k) := \Pr(Y = k)$$

- ▶ The goal is to estimate A_m for $m = 1, \dots, M$ and d .

Prior Art

- ▶ **Dawid-Skene Model** [Dawid & Skene, 1979]:
 - ▶ Proposed expectation maximization (EM) algorithm for ML estimation.
 - ▶ Widely used, but **model identifiability is unclear**.
- ▶ **Spectral Method** [Zhang et al., 2014]:
 - ▶ Established identifiability using orthogonal and symmetric tensor decomposition.
 - ▶ Employed third-order co-occurrences of responses; may have **high sample complexity**.

Pairwise Co-occurrences of Annotator Responses

- ▶ The joint PMF of any two annotator responses,

$$R_{m,\ell}(k_m, k_\ell) = \sum_{k=1}^K \Pr(Y = k) \Pr(X_m = k_m | Y = k) \Pr(X_\ell = k_\ell | Y = k),$$

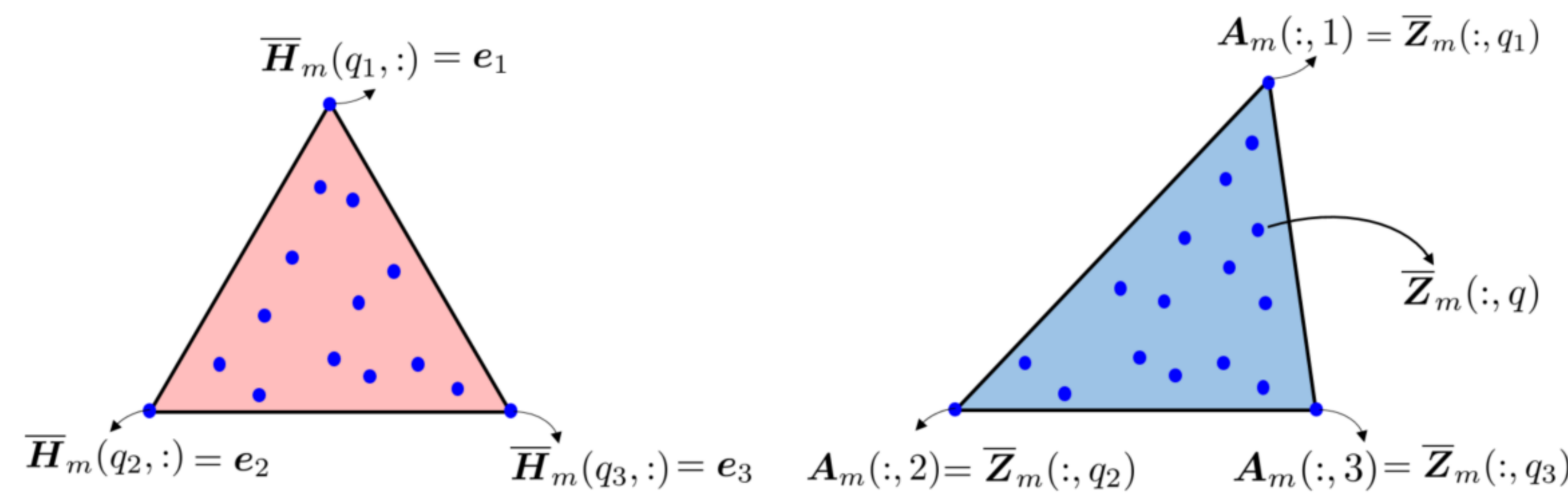
$$\Rightarrow R_{m,\ell} = A_m D A_\ell^\top \in \mathbb{R}^{K \times K}, \quad D = \text{Diag}(d).$$

- ▶ $R_{m,\ell}$'s can be estimated via sample averaging.
- ▶ $R_{m,\ell}$'s are second-order statistics; easier to estimate than third-order ones.

Proposed Approach

- ▶ Consider an annotator m who co-labels with annotators $m_1, \dots, m_{T(m)}$,
- $$Z_m = [R_{m,m_1}, R_{m,m_2}, \dots, R_{m,m_{T(m)}}] = A_m [D A_{m_1}^\top, \dots, D A_{m_{T(m)}}^\top].$$

- ▶ ℓ_1 -normalize the columns of Z_m to get $\bar{Z}_m = A_m \bar{H}_m^\top$ where \bar{H}_m^\top is row normalized.
- ▶ Assume that there exists an index set $A_q = \{q_1, \dots, q_K\}$ such that $\bar{H}_m(A_q, :) = I_K$ (known as **seperability**) [Donoho & Stodden, 2003].



- ▶ Estimating A_m boils down to identifying index set A_q which can be achieved by **successive projection algorithm (SPA)** [Araújo et al. 2001].
- ▶ Index identification via SPA is repeated for every A_m (named as MultiSPA).

Model Identifiability

- ▶ If each class k has an annotator who can perfectly identify class k , then $\bar{H}_m(A_q, :) = I_K$ can be satisfied.

$$Z_m = A_m \underbrace{D [A_{m_1}^\top, \dots, A_{m_{e_1}}^\top, \dots, A_{m_{e_2}}^\top, \dots, A_{m_{e_3}}^\top, \dots, A_{m_{T(m)}}^\top]}_{H_m^\top}$$

Theorem 1: Assume that annotators m and t co-label at least S samples $\forall t \in \{m_1, \dots, m_{T(m)}\}$. Also assume that the constructed \bar{Z}_m satisfies $\|\bar{Z}_m(:, l)\|_1 \geq \eta, \forall l \in \{1, \dots, KT(m)\}$, where $\eta \in (0, 1]$. Suppose that for every class index $k \in \{1, \dots, K\}$, there exists an annotator $m_{t(k)} \in \{m_1, \dots, m_{T(m)}\}$ such that

$$\Pr(X_{m_{t(k)}} = k | Y = k) \geq (1 - \epsilon) \sum_{j=1}^K \Pr(X_{m_{t(k)}} = k | Y = j), \quad \epsilon \in [0, 1]$$

Then, if $\epsilon \leq \mathcal{O}(\max(K^{-1}\kappa^{-3}(A_m), \sqrt{\ln(1/\delta)}(\sigma_{\max}(A_m)\sqrt{S\eta})^{-1}))$, with probability greater than $1 - \delta$, the SPA algorithm can estimate an \hat{A}_m from $Z_m = A_m D H_m^\top$ with the estimation error bounded by $\mathcal{O}(\sqrt{K}\kappa^2(A_m) \max(\sigma_{\max}(A_m)\epsilon, \sqrt{\ln(1/\delta)}(\sqrt{S\eta})^{-1}))$ where $\sigma_{\max}(A_m)$ is the largest singular value of A_m , and $\kappa(A_m)$ is the condition number of A_m .

- ▶ **Implication:** Even if there are no perfect annotators for each class, MultiSPA estimates A_m .

Do we favour more annotators?

Theorem 2: Let $\rho > 0, \epsilon > 0$, and assume that the rows of \bar{H}_m are generated within the $(K-1)$ -probability simplex uniformly at random. If $M \geq \Omega(\frac{\epsilon^{-2(K-1)}}{K} \log(\frac{K}{\rho}))$, then with probability greater than or equal to $1 - \rho$, there exists rows of \bar{H}_m indexed by q_1, \dots, q_K such that

$$\|\bar{H}_m(q_k, :) - e_k^\top\|_2 \leq \epsilon, \quad k = 1, \dots, K.$$

- ▶ **Implication:** If more number of annotators are available, there exists high chance for seperability condition.

Enhanced Identifiability

- ▶ The model can be identified under a relaxed assumption by solving

$$\text{find } \{A_m\}_{m=1}^M, D \quad (1a)$$

$$\text{subject to } R_{m,\ell} = A_m D A_\ell^\top, \quad \forall m, \ell \in \{1, \dots, M\} \quad (1b)$$

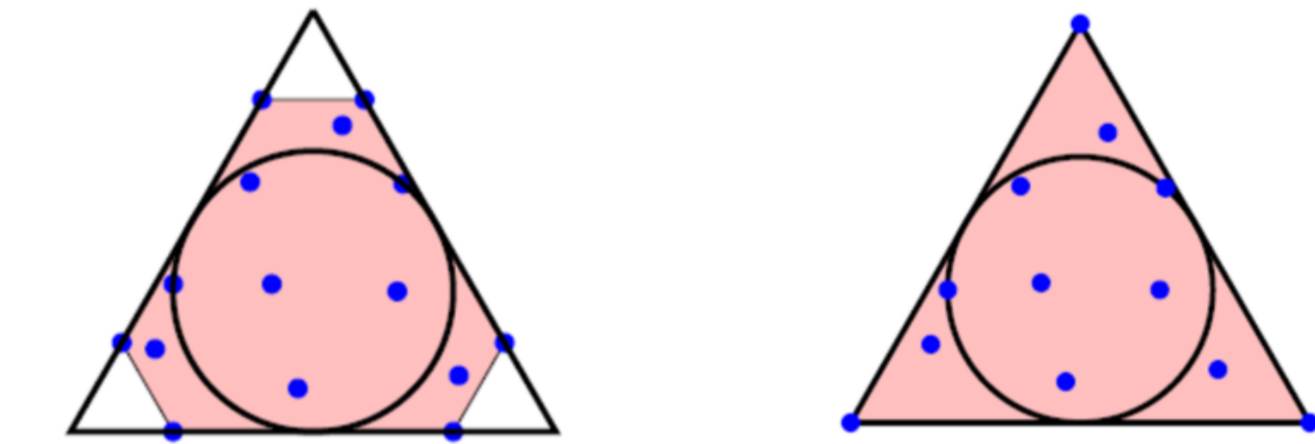
$$\mathbf{1}^\top A_m = \mathbf{1}^\top, \quad A_m \geq 0, \quad \forall m, \quad \mathbf{1}^\top d = 1, \quad d \geq 0. \quad (1c)$$

Theorem 3: Assume that $\text{rank}(D) = \text{rank}(A_m) = K$ for all $m = 1, \dots, M$, and that there exist two subsets of the annotator, indexed by \mathcal{P}_1 and \mathcal{P}_2 , where $\mathcal{P}_1 \cap \mathcal{P}_2 = \emptyset$ and $\mathcal{P}_1 \cup \mathcal{P}_2 \subseteq \{1, \dots, M\}$. Suppose that from \mathcal{P}_1 and \mathcal{P}_2 the following two matrices can be constructed:

$$\hat{R} = \begin{bmatrix} R_{m_1, \ell_1} & R_{m_1, \ell_2} & \dots & R_{m_1, \ell_{|\mathcal{P}_2|}} \\ \vdots & \vdots & \dots & \vdots \\ R_{m_{|\mathcal{P}_1|}, \ell_1} & R_{m_{|\mathcal{P}_1|}, \ell_2} & \dots & R_{m_{|\mathcal{P}_1|}, \ell_{|\mathcal{P}_2|}} \end{bmatrix} = \begin{bmatrix} A_{m_1} \\ \vdots \\ A_{m_{|\mathcal{P}_1|}} \end{bmatrix} D \underbrace{[A_{\ell_1}^\top, \dots, A_{\ell_{|\mathcal{P}_2|}}^\top]}_{H^{(2)\top}}$$

Denote $H^{(1)} = [A_{m_1}^\top, \dots, A_{m_{|\mathcal{P}_1|}}^\top]^\top$, $H^{(2)} = [A_{\ell_1}^\top, \dots, A_{\ell_{|\mathcal{P}_2|}}^\top]^\top$, where $m_t \in \mathcal{P}_1$ and $\ell_j \in \mathcal{P}_2$. Furthermore, assume that both $H^{(1)}$ and $H^{(2)}$ are **sufficiently scattered**. Then, solving **Problem (1)** recovers A_m for $m = 1, \dots, M$ and $D = \text{diag}(d)$ up to identical column permutation.

- ▶ Extremely well trained annotators for each class are not required to satisfy sufficiently scattered condition.



Left: Sufficiently scattered H ; Right: Separable H

- ▶ Problem (1) is solved by a BCD algorithm with KL divergence as the fitting criterion (used MultiSPA as initialization, thus named as MultiSPA-KL).

Amazon Mechanical Turk (AMT) Experiment Results

- ▶ The datasets annotated by AMT workers are used.

Algorithms	TREC		Bluebird		RTE		Web		Dog	
	(%)Error	(sec)Time	(%)Error	(sec)Time	(%)Error	(sec)Time	(%)Error	(sec)Time	(%)Error	(sec)Time
MultiSPA	31.47	50.68	13.88	0.07	8.75	0.28	15.22	0.54	17.09	0.07
MultiSPA-KL	29.23	536.89	11.11	1.94	7.12	17.06	14.58	12.34	15.48	15.88
MultiSPA-D&S	29.84	53.14	12.03	0.09	7.12	0.32	15.11	0.84	16.11	0.12
Spectral-D&S	29.58	919.98	12.03	1.97	7.12	6.40	16.88	179.92	17.84	51.16
TensorADMM	N/A	N/A	12.03	2.74	N/A	N/A	N/A	N/A	17.96	603.93
MV-D&S	30.02	3.20	12.03	0.02	7.25	0.07	16.02	0.28	15.86	0.04
Minmax-entropy	91.61	352.36	8.33	3.43	7.50	9.10	11.51	26.61	16.23	7.22
EigenRatio	43.95	1.48	27.77	0.02	9.01	0.03	N/A	N/A	N/A	N/A
KOS	51.95	9.98	11.11	0.01	39.75	0.03	42.93	0.31	31.84	0.13
GhoshSVD	43.03	11.62	27.77	0.01	49.12	0.03	N/A	N/A	N/A	N/A
Majority Voting	34.85	N/A	21.29	N/A	10.31	N/A	26.93	N/A	17.91	N/A

References

- ▶ Dawid, A. P. and Skene, A. M. *Maximum likelihood estimation of observer error-rates using the em algorithm*. Applied statistics, pp. 20–28, 1979.
- ▶ Zhang, Y., Chen, X., Zhou, D., and Jordan, M. I. *Spectral methods meet em: A provably optimal algorithm for crowdsourcing*. In Advances in Neural Information Processing Systems, pp. 1260–1268, 2014.
- ▶ Donoho, D. and Stodden, V. *When does non-negative matrix factorization give the correct decomposition into parts?* In Advances in Neural Information Processing Systems, pp. 1141–1148 2004.