

CROWDSOURCING

Stephen Mayhew, Daniel Khashabi

May 5, 2015

University of Illinois at Urbana Champaign

1. Majority Voting
2. Hubs and Authorities
3. Singular vector approach
4. EM
5. EM with priors
6. Iterative Weighted Majority Voting
7. Simplified BP
8. Discretized BP
9. Tree-reweighted message passing¹
10. Spectral Meets EM¹

¹never finished...

1. Majority Voting
2. Hubs and Authorities
3. Singular vector approach
4. EM
5. **EM with priors**
6. **Iterative Weighted Majority Voting**
7. Simplified BP
8. **Discretized BP**
9. Tree-reweighted message passing¹
10. Spectral Meets EM¹

¹never finished...

Define

$$g_{ij}(t_i, p_j) = p_j \mathbf{I}\{A_{ij} = t_i\} + (1 - p_j) \mathbf{I}\{A_{ij} \neq t_i\}$$

Define

$$g_{ij}(t_i, p_j) = p_j \mathbf{I}\{A_{ij} = t_i\} + (1 - p_j) \mathbf{I}\{A_{ij} \neq t_i\}$$

Full likelihood (with Beta priors):

$$\mathcal{L}(p; t) = \prod_{i,j} g_{ij}(t_i, p_j) \prod_j \left(c + \frac{1-c}{B(\alpha, \beta)} p_j^{\alpha-1} (1-p_j)^{\beta-1} \right)$$

Define

$$g_{ij}(t_i, p_j) = p_j \mathbf{I}\{A_{ij} = t_i\} + (1 - p_j) \mathbf{I}\{A_{ij} \neq t_i\}$$

Full likelihood (with Beta priors):

$$\mathcal{L}(p; t) = \prod_{i,j} g_{ij}(t_i, p_j) \prod_j \left(c + \frac{1-c}{B(\alpha, \beta)} p_j^{\alpha-1} (1-p_j)^{\beta-1} \right)$$

M-Step becomes:

$$p_j = \frac{\alpha - 1 + \sum_j q(A_{ij})}{\alpha + \beta - 2 + \sum_j q(A_{ij}) + \sum_j q(-A_{ij})}$$

(Note: this same result is a lower bound on the situation where priors come from $p_j = 0.1 + 0.9Z$)

From *Error Rate Bounds and Iterative Weighted Majority Voting for Crowdsourcing* by Hongwei Li and Bin Yu²

Algorithm 1 The iterative weighted majority voting algorithm (IWMV)

Input: Number of workers= M ; Number of items= N ; data matrix: $Z \in [L]^{M \times N}$;

Output: the predicted labels $\{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N\}$

Initialization: $\nu_i = 1, \forall i \in [M]$; $T_{ij} = \mathbf{I}(Z_{ij} \neq 0), \forall i \in [M], \forall j \in [N]$.

repeat

$$\hat{y}_j \leftarrow \operatorname{argmax}_{k \in [L]} \sum_{i=1}^M \nu_i \mathbf{I}(Z_{ij} = k), \quad \forall j \in [N].$$

$$\hat{w}_i \leftarrow \frac{\sum_{j=1}^N \mathbf{I}(Z_{ij} = \hat{y}_j)}{\sum_{j=1}^N T_{ij}}, \quad \forall i \in [M].$$

$$\nu_i \leftarrow L \hat{w}_i - 1, \quad \forall i \in [M].$$

until converges or reaches S iterations.

Output the predictions $\{\hat{y}_j\}_{j \in [N]}$ by $\hat{y}_j = \operatorname{argmax}_{k \in [L]} \sum_{i=1}^M \nu_i \mathbf{I}(Z_{ij} = k)$.

²<http://arxiv.org/pdf/1411.4086v1.pdf>

Original BP updates are:

$$y_{a \rightarrow i}^{(k)}(p_a) \propto \mathcal{F}(p_a) \prod_{j \in \partial a \setminus i} \left\{ (p_a + \bar{p}_a + (p_a - \bar{p}_a)A_{ja})x_{j \rightarrow a}^{(k)}(+1) + (p_a + \bar{p}_a - (p_a - \bar{p}_a)A_{ja})x_{j \rightarrow a}^{(k)}(-1) \right\}$$

$$x_{i \rightarrow a}^{(k+1)}(\hat{t}_i) \propto \prod_{b \in \partial i \setminus a} \int \left(y_{b \rightarrow i}^{(k)}(p_b)(p_b \mathbb{I}_{(A_{ib}=\hat{t}_i)} + \bar{p}_b \mathbb{I}_{(A_{ib} \neq \hat{t}_i)}) \right) dp_b$$

Original BP updates are:

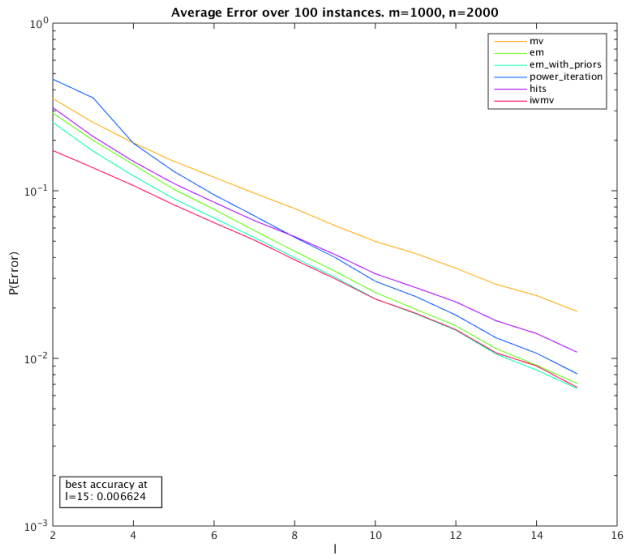
$$y_{a \rightarrow i}^{(k)}(p_a) \propto \mathcal{F}(p_a) \prod_{j \in \partial a \setminus i} \left\{ (p_a + \bar{p}_a + (p_a - \bar{p}_a)A_{ja})x_{j \rightarrow a}^{(k)}(+1) + (p_a + \bar{p}_a - (p_a - \bar{p}_a)A_{ja})x_{j \rightarrow a}^{(k)}(-1) \right\}$$

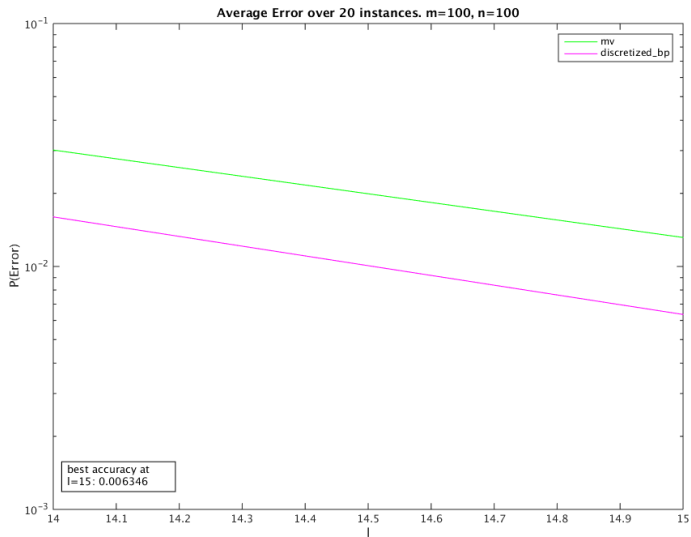
$$x_{i \rightarrow a}^{(k+1)}(\hat{t}_i) \propto \prod_{b \in \partial i \setminus a} \int \left(y_{b \rightarrow i}^{(k)}(p_b)(p_b \mathbb{I}_{(A_{ib}=\hat{t}_i)} + \bar{p}_b \mathbb{I}_{(A_{ib} \neq \hat{t}_i)}) \right) dp_b$$

Discretize p_j :

$$x_{i \rightarrow a}^{(k+1)}(\hat{t}_i) \propto \prod_{b \in \partial i \setminus a} \sum_{p_b \in P} \left(y_{b \rightarrow i}^{(k)}(p_b)(p_b \mathbb{I}_{(A_{ib}=\hat{t}_i)} + \bar{p}_b \mathbb{I}_{(A_{ib} \neq \hat{t}_i)}) \right)$$

For $P = \{p_1, p_2, p_3, \dots, p_k\}$





Algorithm	m, n	$\ell = 15$
Discretized BP	100,100	0.006346
EM with Priors	100,100	0.008478
EM with Priors	1000,2000	0.006624
EM with Priors	250,2000	0.005512
Simplified BP	250,2000	0.005466
IWMV	250,2000	0.0057