

- Crowd is large, anonymous, transient
  - Impossible to build up a trust relationship
  - Difficult to condition payment on correct answers (no groundtruth, any delay is annoying workers, ...)
  - Need to exploit redundancy:
    - One tasks assigned to multiple workers
    - One worker realising more than one task
- This presentation: some (classical) ML approach to Learning with Crowds
  - Classical, not alicia(l)!

# Models

- First Level: estimate the ground truth  $y_{i \uparrow}$  from a set of  $\{y_{i \uparrow k}\}$ 
  - Ex: Majority Voting
- Second Level: models the worker's expertise
  - Simple “accuracy” parameter
  - “Confusion-matrix” set of parameters
  - Intra/Inter-consistency parameters
- Third Level: also models task difficulty
  - One parameter for worker's expertise and one parameter for task difficulty
- Of course, all parameters are unknown and should be identified together with:
  - Either the correct labels ( $y_{i \uparrow}$ )
  - Classification models to predict the correct labels for new crowdsourced instances

# Level 1 : Estimating the correct label $y_i$

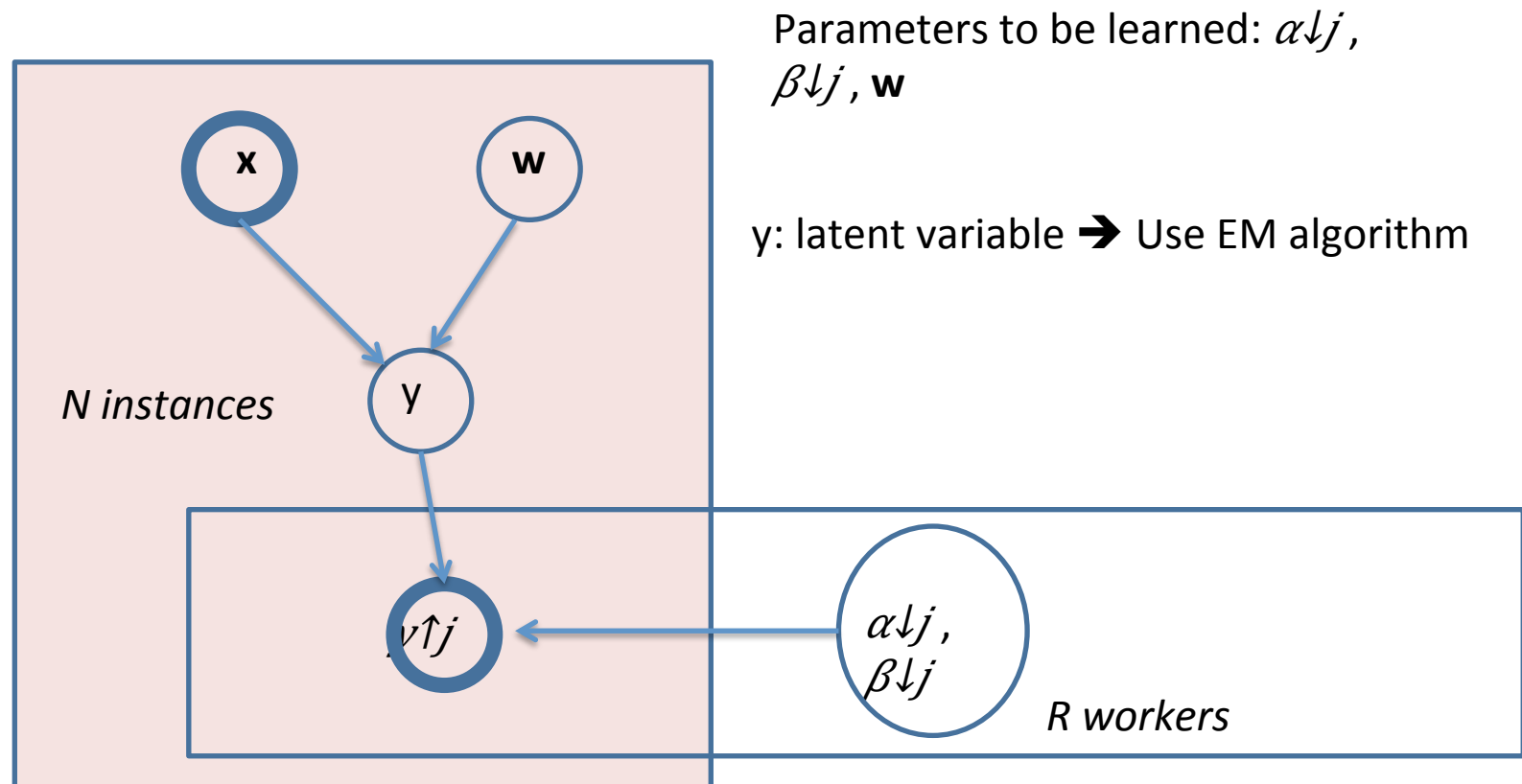
- Assume  $m$  tasks,  $n$  workers, 1 task to  $w$  workers, 1 workers to  $t$  tasks ( $\rightarrow m.w=n.t$ )
  - One-shot setting! (anonymous, transient workers)
- Basic inputs:  $A_{i,j}$  (binary labels  $y_i : -1 \text{ or } +1$ )
  - $A_{i,j} = 0$  if worker  $j$  didn't label task  $i$
  - $A_{i,j} = y_i$  with probability  $p_j$
  - $A_{i,j} = -y_i$  with probability  $(1-p_j)$
  - $p_j$  not modelled but assumed i.i.d. from some unknown distributions with some properties (for convergence)
- Bipartite graph – HITS algorithm
  - Starts with  $v_j = 1/w$
  - Iterate until convergence
    - $u_i = \sum_j A_{i,j} \cdot v_j$
    - $v_j = \sum_i A_{i,j} \cdot u_i$
  - Output:  $y_i = \text{sign}(u_i)$

# Level 1 : HITS Algo

- Let  $z = E[(2 p_{\downarrow j} - 1)]$  and  $q = E[(2 p_{\downarrow j} - 1)^2]$ 
  - $q=0 \rightarrow$  all random workers
  - $q=1 \rightarrow$  all workers are diligent
- Then, if  $z > 0$  and  $q^2(w-1) \cdot (t-1) > 1$   
 $P(y_{\downarrow i} \neq y_{\uparrow i}) \leq e^{-wq/k}$
- NB. Converges towards the leading left and right singular vectors of  $A_{\downarrow i, j}$
- The bounds can be translated in the  $w$  requirements to achieve a given accuracy, given the average quality  $q$ .

## Level II : Individual Worker Model

- $p(y^j) = (\beta^j)^y (1 - \beta^j)^{1-y}$



# EM Algorithm (I)

- Conditional Log-Likelihood (iid assumptions):

$$\text{Log } p(\mathbf{D}, \boldsymbol{\theta}, \mathbf{x}) = \sum_{i=1}^n \log(a_i p_i + b_i (1 - p_i))$$

With  $p_i = \sigma(\mathbf{w}^T \cdot \mathbf{x}_i)$

$$a_i = \prod_{j=1}^R \alpha_j^{y_{ij}} (1 - \alpha_j)^{(1 - y_{ij})}$$

$$b_i = \prod_{j=1}^R \beta_j^{(1 - y_{ij})} (1 - \beta_j)^{y_{ij}}$$

Find ML estimator of  $\boldsymbol{\theta} = [\boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{w}]$  as  $\text{argmax}\{\text{Log } p(\mathbf{D}, \boldsymbol{\theta}, \mathbf{x})\}$

→ Efficiently solved by EM algo

- E-step: fixed  $\boldsymbol{\theta} \rightarrow$  expectation of  $y_{ij} \triangleq \mu_{ij}$
- M-step: fix  $\mu_{ij}$ , and estimate  $\boldsymbol{\theta}$

## EM Algorithm (II)

- E-step:  $\mu_i = p(y_i = 1 | y_i \in \{1, \dots, R\}, x_i, \theta)$   
 $= a_i p_i / (a_i p_i + b_i (1 - p_i))$

- **M-step:** find that maximizes

$$\sum_{i=1}^N \mu_i \log(a_i p_i) + (1 - \mu_i) \log(b_i (1 - p_i))$$

$$\rightarrow \alpha_j = \sum_{i=1}^N \mu_i y_i^j / \sum_{i=1}^N \mu_i ; \beta_j = \sum_{i=1}^N (1 - \mu_i) (1 - y_i^j) / \sum_{i=1}^N (1 - \mu_i)$$

and w solved by logistic regression with soft labels (or equivalently weighted logistic regression), with soft label  $= \mu_i$

# Notes

- Could be bayesian, by giving a gaussian prior on  $w$ , and beta priors on  $\alpha, \beta$
- $\rightarrow \alpha_{\downarrow j} = a_1 + \sum_{i \in N} \mu_{\downarrow i} y_{\downarrow i \uparrow j} / a_1 + a_2 + \sum_{i \in N} \mu_{\downarrow i}$  ,  $\beta_{\downarrow j} = b_1 + \sum_{i \in N} (1 - \mu_{\downarrow i})(1 - y_{\downarrow i \uparrow j}) / b_1 + b_2 + \sum_{i \in N} (1 - \mu_{\downarrow i})$   
and L2-regularized soft log regr.
- In practice, start with M-step and  $\mu_{\downarrow i}$  given by average vote.
- Similar developments for regression
- Could be extended to the “no-feature” case ( $p_{\downarrow i}$  replaced by  $p_{\downarrow 0}$  )
  - E-step:  $\mu_{\downarrow i} = a_{\downarrow i} p_{\downarrow 0} / a_{\downarrow i} p_{\downarrow 0} + b_{\downarrow i} (1 - p_{\downarrow 0})$
  - M-step: estimate  $\alpha, \beta$  as before and  $p_{\downarrow 0} = \sum_{i \in N} \mu_{\downarrow i} / N$