

Vandex

Yandex

Practice of Efficient Data Collection via Crowdsourcing: Aggregation, Incremental Relabelling, and Pricing

Alexey Drutsa, Valentina Fedorova, Dmitry Ustalov, Olga Megorskaya, Evfrosiniya Zerminova, Daria Baidakova

WSDM 2020 hands-on tutorial



Theory on aggregation

Valentina Fedorova, Research analyst

Yandex.Toloka is a service of Swiss company Yandex Services AG

Tutorial outline



Part VI: 25 min Theory on Aggregation

Coffee break: 30 min

Part VI: 60 min Set & Run Projects cont.

Part VII: 20 min Incremental relabeling and pricing

Part VIII: 10 min Results & Conclusions

Labelling data with crowdsourcing



> How to choose a reliable label?

- > How many workers per object?
- > How much to pay to workers?



Evaluation of labelling approaches

Accuracy



Labels with a maximal level of accuracy for a given budget or Labels of a <u>chosen accuracy level</u> for a <u>minimal budget</u>



Key components of labelling with crowds



based pricing

Aggregation

Labelling data with crowds



Upload multiple copies of each object to label

Workers assign noisy labels to objects

Aggregate multiple labels for each object into a more reliable one

Process results

| Yandex Toloka | Projects | Users | Skills | Profile | Messages |
|---|----------|-------------|-----------------|---------|-----------|
| Projects -> Does the image contains traffic lights? -> pool | | | | | |
| pool – clo | sed | | | | |
| POOL TASKS (File example for task uploading (tsv, UTF-8)) 🕜 | | | | | |
| 🛨 Upload 🖺 files | | Edit | | | • Preview |
| 30 task suites | | 0 ti | raining ask | | |
| 90 tasks | | 10 | control task | | |
| | | | | | |



Multiclass labels

Project 1: Filter images



Are there shoes in the picture?



Notation

- Categories $k \in \{1, ..., K\}$. E.g.: O Cat \triangle Dog \Box Other >
- > Objects $j \in \{1, ..., J\}$. E.g.:



Workers: $w \in \{1, ..., W\}$. E.g.:



• $W_j \subseteq \{1, \dots, W\}$ - workers labelled object j



The simplest aggregation: Majority Vote (MV)

- The problem of aggregation: **Observe noisy labels** $\mathbf{y} = \{y_j^w | j = 1, ..., J \text{ and } w = 1, ..., W\}$
- Recover true labels $\mathbf{z} = \{z_j | j = 1, ..., J\}$ A straightforward solution:



 $\hat{z}_j^{MV} = \arg \max_{y=1,...,K} \sum_{w \in W_j} \delta(y = y_j^w)$, where $\delta(A) = 1$ if A is true and 0 otherwise

Performance of MV vs other methods



Zhou D. et al. Regularized minimax conditional entropy for crowdsourcing. 2015



Properties of MV

> All workers are treated similarly



> All objects are treated similarly



Advanced aggregation: workers and objects

Parameterize expertise of workers
 Parameterize difficulty of objects
 by e^w
 by d_j





Advanced aggregation: latent label models



Latent label models: noisy label model



> A noisy label model $M_j^w = M(e^w, d_j)$

is a matrix of size *K*×*K* with elements

$$M_j^w[c,k] = Pr(\underline{Y_j^w} = k \mid \underline{Z_j} = c)$$



Latent label models: generative process



- > Noisy labels generation:
- 1. Sample z_j from a distribution $P_Z(p)$
- 2. Sample y_j^w from a distribution $P_Y(M_j^w[z_j,\cdot])$

In multiclassification, a standard choice for $P_Z(\cdot)$ and $P_Y(\cdot)$ is a Multinomial distribution Mult(\cdot)

Latent label models: parameters optimization

- The likelihood of y and z under the latent label model:



Estimate parameters and true labels by maximizing L(...)

Assumption: y_i^w is cond. independent of everything else given z_i , d_i , e^w

observed noisy label

$$\prod_{i \in J} \sum_{\substack{z_j \in \{1, \dots, K\}}} \Pr(z_j | p) \prod_{w \in W_j} \Pr(\frac{y_j^{w} | z_j, d_j, e^w}{y_j^{w} | z_j, d_j, e^w})$$

likelihood of noisy and true labels for object *j*

Latent label models: EM algorithm

> Maximization of the expectation of log-likelihood (LL)*

$$\mathbb{E}_{\mathbf{z}}\log\Pr(\mathbf{y}, \mathbf{z}) = \sum_{j \in J} \sum_{z_j \in \{1, \dots, K\}} \Pr(z_j)$$

> **E-step**: Use Bayes' theorem for posterior distribution of \hat{z} given p, d, e:

$$\hat{z}_j[c] = \Pr(Z_j = c | \mathbf{y}, p, \mathbf{d}, \mathbf{e}) \propto \Pr(Z_j = c | p) \prod_{w \in W_j} \Pr(\mathbf{y}_j^w | Z_j = c, \mathbf{d}_j, e^w)$$

> M-step: Maximize the expectation of LL with respect to the posterior distribution of \hat{z} :

 $(p, \mathbf{d}, \mathbf{e}) = \operatorname{argmax} \mathbb{E}_{\hat{\mathbf{z}}} \log$

- Analytical solutions
- Gradient descent

 $Z_j|p)\log\prod_{w\in W_j} \Pr(Z_j|p)\Pr(Y_j^w|Z_j,d_j,e^w)$

$$\operatorname{SPr}(z_j|p) \prod_{w \in W_j} \operatorname{Pr}(y_j^w|z_j, d_j, e^w)$$



Dawid and Skene model (DS):

- categories are different
- objects are similar >
- workers are different >

Generative model of labels, abilities, and difficulties (GLAD):

- categories are similar
- objects are different
- workers are different

Minimax conditional entropy model (MMCE):

- categories are different
- objects are different
- workers are different >

Dawid and Skene model (DS)¹



1. Dawid and Skene, "Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm", 1979

- LLM with parameters:
 - > p vector of length K: p[i] = Pr(Z = c)
 - > e^{w} matrix of size $K \times K$: $e^{w}[c,k] = \Pr(Y^{w} = k | Z = c)$



DS: parameters optimization

> E-step:

$$\widehat{z_j}[c] = \frac{p[c] \prod_{w \in W_j} e^w[c, y_j^w]}{\sum_k p[k] \prod_{w \in W_j} e^w[k, y_j^w]}, \qquad c = 1, \dots, K$$

> M-step: Analytical solution

$$e^{w}[c,k] = \frac{\sum_{j \in J} \widehat{z_{j}}[c] \delta(\underline{y_{j}^{w}} = k)}{\sum_{q=1}^{K} \sum_{j \in J} \widehat{z_{j}}[c] \delta(\underline{y_{j}^{w}} = q)}, \qquad k, c = 1, \dots, K$$
$$p[c] = \frac{\sum_{j \in J} \widehat{z_{j}}[c]}{L}, \qquad c = 1, \dots, K$$

Generative model of Labels, Abilities, and Difficulties (GLAD)²



difficulty

2. Whitehill et al., Whose vote should count more: Optimal integration of labels from labelers of unknown expertise, 2009

- LLM with parameters:
 - > scalar $d_i \in (0, \infty)$
 - > scalar $e^w \in (-\infty, \infty)$
 - Model:

$$\left(\frac{Y_{j}^{w}}{K}=k|Z_{j}=c\right) = \begin{cases} a(w,j), & c=k\\ \frac{1-a(w,j)}{K-1}, c\neq k\\ 1\\ \end{array}$$
where $a(w,j) = \frac{1}{1+\exp(-e^{w}d_{j})}$

GLAD: parameters optimization

> Let
$$a(w,j) = \frac{1}{1 + \exp(-e^w d_j)}$$
 and $P(z)$

E-step:

$$\widehat{z_j}[c] \propto \mathbb{P}(Z_j = c) \prod_{w \in W_j} a(w, j)^{\delta\left(\frac{y_j^w}{j} = c\right)} \left(\frac{1 - a(w, j)}{K - 1}\right)^{\delta\left(\frac{y_j^w}{j} \neq c\right)}, \ c = 1, \dots, K$$

M-step: estimate (d, e) for given \hat{z} using gradient descent $(\mathbf{d}^{\mathsf{t}}, \mathbf{e}^{\mathsf{t}}) = \operatorname{argmax} \sum_{i \in I} \left[\mathbb{E}_{\widehat{z_i}} \right]_{i \in I}$

 z_i) be a predefined prior (e.g., $P(z_i) = \frac{1}{K}$)

$$\log P(z_j) + \sum_{w \in W_j} \mathbb{E}_{\widehat{z_j}} \log \Pr(y_j^w | z_j) \right]$$

MiniMax Conditional Entropy model (MMCE)³

Find parameters that minimize the maximum conditional entropy of observed labels:



$$\operatorname{ax}_{P} - \sum_{\substack{j \in J \\ c \in \{1, \dots, K\}}} Q(Z_{j} = c) \sum_{\substack{w \in W \\ k \in \{1, \dots, K\}}} P(Y_{j}^{w} = k | Z_{j} = c) \log P(Y_{j}^{w} = k | Z_{j})$$

LLM with parameters:

> d_i – matrix of size $K \times K$

> e^{w} – matrix of size $K \times K$

Noisy label model: $\Pr(\underline{Y_j^w} = k | \underline{Z_j} = c) = \exp(\underline{d_j}[c, k] + e^w[c, k])$





Summary of aggregation methods



Pairwise comparisons

Project 4: Compare items



Which shoes look more similar to the one in the picture?











Notation

- Answers: Left or Right >
- Items $d_i \in \{1, ..., N\}$. E.g.: >



Tasks:



> Workers: $w \in \{1, ..., W\}$. E.g.:





Formalization

Ranking from pairwise comparisons: > Given pairwise comparisons for items in *D*: $P = \{(w_k, d_i, d_j): i \succ_k j\}$

> Obtain **a ranking** π over items D – based on answers in P



$$\rightarrow \{1, \dots, N\}$$



Difference from multiclassification

The latent label assumption is not satisfied when comparing complex items >



Different tasks may contain common items







Bradley and Terry model (BT)

Assume that each item $d_i \in D$ has a latent "quality" score $s_i \in \mathbb{R}$



The probability that $d_i \in D$ will be preferred in a comparison over $d_i \in D$ $\Pr(i \succ j) = f(s_i - s_j),$

where $f(x) = \frac{1}{1 + \rho^{-x}}$.

The model assumes that all workers are equally good and truthful



- Probability that w_k reads a task is
 - If w_k reads the task, she answers according to scores: $(f(s_i - s_j), f(s_j - s_i))$
 - Probability to choose Left if compares items

Probability to choose **Left** if <u>answers randomly</u>

Logistic function $Pr(w_k \text{ reads a task}) = f(\gamma_k)$



If w_k does not read the task, she answers according to her bias: $(f(\boldsymbol{q}_k), f(-\boldsymbol{q}_k))$

NoisyBT: likelihood of workers' answers

The likelihood of $i \succ_k j$ is

$$\Pr(i \succ_k j) = \underbrace{f(\gamma_k)f(s_i - s_j)}_{(\gamma_k)} + \underbrace{(1 - f(\gamma_k))f((-1)^{(1 - \mathbb{I}(d_i \text{ was left}))}q_k)}_{(\gamma_k)},$$

Truthful answer

where $I(d_i \text{ was left})$ is the indicator for the order of d_i and d_j

$$d_i$$
 d_j

 $\mathbb{I}(d_i \text{ was left}) = 1$

Random answer



 $\mathbb{I}(d_i \text{ was left}) = 0$

NoisyBT: parameters optimization Likelihood of observed comparisons:

$$T(s, \mathbf{q}, \mathbf{\gamma}) = \int_{(w_k, d_k)} \left[\log[f(\mathbf{\gamma}_k)f(\mathbf{q}_k)] \right] ds$$

 $\sum_{\substack{d_i,d_j \in P}} \log \Pr(i \succ_k j) =$ $(s_i - s_i) + (1 - f(\gamma_k))g(q_k)$

- > $\{s_i\}_{i=1,...,N}$ and $\{\gamma_k, q_k\}_{k=1,...,W}$ are inferred by maximizing the log-likelihood:
 - $T(s,q,\gamma) \rightarrow \max_{\{s_i,\gamma_k,q_k\}}$
 - To obtain a ranking π over items, sort items according to their scores

Summary about pairwise comparisons

Latent scores models for ranking from pairwise comparisons:



To reduce bias from unreliable answers parameterize workers



 \rightarrow "reliability" γ_k and "bias" q_k



Thank you! Questions?

Valentina Fedorova

Research analysts



valya17@yandex-team.ru



https://research.yandex.com/tutorials/crowd/wsdm-2020

