NETS 213: CROWDSOURCING AND HUMAN COMPUTATION

Quality Control part 3





Different Mechanisms for Quality Control

Aggregation and redundancy Embedded gold standard data Economic incentives Reputation systems Statistical models



Expectation Maximization algorithm

EM is an algorithm for finding the probabilities of unobserved variables

We will use it to estimate how accurate workers' labels are, and infer how good each worker is

This is more sophisticated than voting



Dawid and Skene (1977)

Maximum Likelihood Estimation of Observer Error-rates using the EM Algorithm

Dawid and Skene examined application of EM to medical diagnosis.

Patients are sometimes treated by multiple physicians, who can give different diagnoses

Why? Doctors may ask different questions. Patient may describe history differently. Doctors may classify symptoms differently.

Observer Error

Given that different doctors have different opinions, they can't all be right.

How often do individual physicians suffer from "observer error"? Are their errors systematic?

Answers depend on the "true" diagnosis.

Observer Error

Observer error would be easy to calculate if we had ground truth.

Simply count the misdiagnoses and divide by the total number of diagnoses.

However, sometimes it is impossible to know what diagnosis is correct. Same set of symptoms can arise from multiple root causes.

"I know it when I see it"

I shall not today attempt further to define "hardcore pornography"; and perhaps I could never succeed in intelligibly doing so. But I know it when I see it. -Justice Potter Stewart



url	worker 1	worker 2	worker 3	worker 4	worker 5
sunnyfun.com	porn	not	not	not	porn
sex- mission.com	porn	porn	porn	porn	porn
google.com	not	porn	not	not	porn
youporn.com	porn	porn	porn	porn	not
yahoo.com	porn	not	not	not	porn

Solution?

Can't have Justice Stewart rule on everything

Instead, we will apply Dawid and Skene's EM algorithm, which iteratively

- 1. Estimates the correct answers, using labels from multiple workers, and accounts for the quality of each worker
- 2. Estimates the quality of the workers by comparing the submitted answers to the inferred correct answers

Inputs

a set of **N** objects o**1** ... o_N sunnyfun.com, sex-mission.com, google.com, youporn.com, yahoo.com

a set of *L* possible labels: {porn, not porn}

Labels for each object by **K** workers worker1, worker2, worker3, worker4, worker5



Recover the true class label $T(o_n)$ for each object o_n when "gold" truth is unknown Since the true labels are not known / never directly observed, they are called **latent** variables

Goal 2

For each worker who contributed labels, calculate their accuracy or reliability

To calculate accuracy show how often they mistakenly choose one label when a different one is the actual truth

Chicken and egg problem

If we knew what the **true class labels** were for each object for each object, then we could compute each Turker's accuracy

If we had **accuracies for every Turker**, then we could infer what the true label for each object should be

Input: Labels l[k][n] from worker (k) to object o_n ,

- **Output**: Confusion matrix $\pi_{ij}^{(k)}$ for each worker (k), Correct labels $T(o_n)$ for each object o_n , Class priors $Pr\{C\}$ for each class C
- 1 Initialize error rates $\pi_{ij}^{(k)}$ for each worker (k) (e.g., assume each worker is perfect);
- 2 Initialize correct label for each object $T(o_n)$ (e.g., using majority vote);
- 3 while not converged do
- 4 Estimate the correct label $T(o_n)$ for each object, using the labels $l[\cdot][n]$ assigned to o_n by workers, weighting the votes using the error rates $\pi_{ij}^{(k)}$;
- 5 Estimate the error rates $\pi_{ij}^{(k)}$, for each worker (k), using the correct labels $T(o_n)$ and the assigned labels l[k][n];
- **6** Estimate the class priors $Pr\{C\}$, for each class C;
- **7** end
- 8 return Estimated error rates π^(k)_{ij}, Estimated correct labels T(o_n), Estimated class priors Pr{C}
 Algorithm 1: The EM algorithm for worker quality estimation.

😽 Penn Engineering

	worker1	worker2	worker3	worker4	worker5
sunnyfun.com	porn	not	not	not	porn
sex- mission.com	porn	porn	porn	porn	porn
google.com	not	porn	not	not	porn
youporn.com	porn	porn	porn	porn	not
yahoo.com	porn	not	not	not	porn



Output: "True" Labels

url	True Labels
sunnyfun.com	not
sex-mission.com	porn
google.com	not
youporn.com	porn
yahoo.com	not

Repeat until convergence

You can continue to iterate until your values converge For this example, we converge after the first iteration

	worker1	worker2	worker3	worker4	worker5
sunnyfun.com	porn	not	not	not	porn
sex- mission.com	porn	porn	porn	porn	porn
google.com	not	porn	not	not	porn
youporn.com	porn	porn	porn	porn	not
yahoo.com	porn	not	not	not	porn



	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

	porn	not
sunnyfun	?	?
sex-mission	?	?
google	?	?
youporn	?	?
yahoo	?	?

worker1	porn	not
porn	?	?
not	?	?

worker2	porn	not
porn	?	?
not	?	?

worker3	porn	not
porn	?	?
not	?	?

worker4	porn	not
porn	?	?
not	?	?

worker5	porn	not
porn	?	?
not	?	?



	worker1	worker2	worker3	wor	ker4	worker5
sunnyfun	porn	not	not	not		porn
sex- mission	porn	porn	porn	porn		porn
google	not	porn	not	not		porn
youporn	porn	porn	porn	pori	n	not
yahoo	porn	not	not	not		porn
			porn		not	
	sunnyfur	1	?		?	
	sex-missi	on	?		?	
	google youporn		?		?	
			?		?	
	yahoo		?		?	

worker1	porn	not
porn	1	0
not	0	1
worker2	porn	not
porn	1	0
not	0	1
worker3	porn	not
porn	1	0
not	0	1
worker4	porn	not
porn	1	0
not	0	1
worker5	porn	not
porn	1	0
not	0	1

	els using
Compute law	ie
majority	

porn	not	not	not		porn
		porn		not	
sunnyfun		2		?	
sex-missi	on	?		?	
google		?		?	
youporn		?		?	
yahoo		?		?	

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

worker1	porn	not
porn	1	0
not	0	1
worker2	porn	not

porn	1	0
not	0	1
worker3	porn	not
porn	1	0
not	0	1
worker4	porn	not

vorker4	porn	not
oorn	1	0
not	0	1

worker5	porn	not
porn	1	0
not	0	1

	hels using
Compute la	ote
majority	Sandraw Street Str

porn	not	not	not		porn
		porn		not	
sunnyfun		2		3	
sex-missi	on	?		?	
google		?		?	
youporn		?		?	
yahoo		?		?	

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

worker1	porn	not
porn	1	0
not	0	1

worker2	porn	not
porn	1	0
not	0	1
		-
worker3	porn	not
worker3 porn	porn 1	not 0

vorker4	porn	not
orn	1	0
not	0	1

worker5	porn	not	
porn	1	0	
not	0	1	

		-
worker4	porn	not
porn	1	0
not	0	1

Treat the correct laber as the one with the most votes	

orn	not	not	not		porn
		porn		not	
sunnyfun		0		1	
sex-missi	on	1		0	
google		0		1	
/ouporn		1		0	
yahoo		0		1	

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

worker1	porn	not
porn	0	0
not	0	0
worker2	porn	not
porn	0	0
not	0	0
worker3	porn	not
porn	0	0
not	0	0
worker4	porn	not
porn	0	0
not	0	0
worker5	porn	not
porn	0	0

not

0

0

rn	0	0	
t	0	0	
orker3	porn	not	
rn	0	0	

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

	porn	not
sunnyfun	0	1
sex-mission	1	0
google	0	1
youporn	1	0
yahoo	0	1



worker4	porn	not
porn	0	0
not	0	0

worker5	porn	not
porn	0	0
not	0	0

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

	porn	not
sunnyfun	0	1
sex-mission	1	0
google	0	1
youporn	1	0
yahoo	0	1



worker4	porn	not
porn	0	0
not	0	0

worker5	porn	not
porn	0	0
not	0	0

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

	porn	not
sunnyfun	0	1
sex-mission	1	0
google	0	1
youporn	1	0
yahoo	0	1



worker4	porn	not
porn	0	0
not	0	0

worker5	porn	not
porn	0	0
not	0	0

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

	porn	not
sunnyfun	0	1
sex-mission	1	0
google	0	1
youporn	1	0
yahoo	0	1



worker4	porn	not
porn	0	0
not	0	0

worker5	porn	not
porn	0	0
not	0	0

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

	porn	not
sunnyfun	0	1
sex-mission	1	0
google	0	1
youporn	1	0
yahoo	0	1

worker1	porn	not
porn	1	0
not	0.67	0.33
worker2	porn	not
porn	0	0 R
not	0	0
worker3	porn	not
oorn	0	0
not	0	0
worker4	porn	not
oorn	0	0
not	0	0
worker5	porn	not
oorn	0	0
not	0	0

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn
			porn	not	

	porn	not
sunnyfun	0	1
sex-mission	1	0
google	0	1
youporn	1	0
yahoo	0	1

worker1	porn	not
porn	1	0
not	0.67	0.33
worker2	porn	not
porn	1	0
not	0.33	0.67
worker3	porn	not
porn	1	0
not	0	1
worker4	porn	not
porn	1	0
not	0	1
		Ī
worker5	porn	not
porn	0.5	0.5
not	1	0

فتجحد



Donn	Lingin	oring
	LIQUI	
\bigcirc	0	0

	sunnyfun	porn	not	not	not		porn
	sex- mission	porn	porn	porn	porr	ו	porn
	google	not	porn	not	not		porn
	youporn	porn	porn	porn	porr	٦	not
	yahoo	porn	not	not	not		porn
	-			porn		not	
label	5	sunnyfun		1.5			
Recompute need		sexmissi	on				
using majority vo	and the second se	google					
		youporn					
Penn Engineering		yahoo					

worker1 worker2 worker3 worker4 worker5

worker1	porn	not
porn	1	0
not	0.67	0.33

worker2	porn	not
porn	1	0
not	0.33	0.67
worker3	porn	not

worker3	porn	not
oorn	1	0
not	0	1

vorker4	porn	not
oorn	1	0
not	0	1

worker5	porn	not
porn	0.5	0.5
not	1	0

Donn	Lingin	oning
А ГЕШІ		
V	8	8

	sunnyfun	porn	not	not	not		porn
	sex- mission	porn	porn	porn	porr	า	porn
	google	not	porn	not	not		porn
	youporn	porn	porn	porn	porr	٦	not
	yahoo	porn	not	not	not		porn
	-			porn		not	
alabel	5	sunnyfun		1.5		4.34	
Recompute heighted	and the second second	sexmissi	on				
majority vo		google					
Landon and a second		youporn					
enn Engineering		yahoo					

worker1 worker2 worker3 worker4 worker5

worker1	porn	not
porn	1	0
not	0.67	0.33
worker2	porn	not
porn	1	0
not	0.33	0.67
worker3	porn	not
porn	1	0
not	0	1
worker4	porn	not
porn	1	0
not	0	1
worker5	porn	not
porn	0.5	0.5
not	1	0

	🐯 Penn Engineering	
--	--------------------	--

	youporn	porn	porn	porn	porr	ſ	not
	yahoo	porn	not	not	not		porn
	-			porn		not	
labe	15	sunnyfun		0.26		0.74	
Recompute need		sexmissi	on	0.69		0.31	
majority voe	No. of Concession, Name	google		0.29		0.71	
		youporn		0.82		0.18	
Penn Engineering		yahoo		0.26		0.74	
0 0							

worker2 worker3 worker4 worker5

not

porn

not

porn

porn

porn

not

porn

not

worker1

not

porn

porn

porn

porn

not

sunnyfun

sex-

mission

google

worker1	porn	not
porn	1	0
not	0.67	0.33

worker2	porn	not
porn	1	0
not	0.33	0.67
workor?	0050	nat

worker3	porn	not
porn	1	0
not	0	1

worker4	porn	not
oorn	1	0
not	0	1

worker5	porn	not
porn	0.5	0.5
not	1	0

1		-	
not	0	1	
worker4	porn	not	
porn	1	0	
	0	4	

Treat the correct label as the one with the most votes	
Lanna and a second seco	

orn	not	not	not		porn
		porn		not	
sunnyfun		0		1	
sex-missi	on	1		0	
google		0		1	
/ouporn		1		0	
/ahoo		0		1	

	worker1	worker2	worker3	worker4	worker5
sunnyfun	porn	not	not	not	porn
sex- mission	porn	porn	porn	porn	porn
google	not	porn	not	not	porn
youporn	porn	porn	porn	porn	not
yahoo	porn	not	not	not	porn

worker1	porn	not
porn	1	0
not	0.67	0.33
worker2	porp	not

workerz	porn	not			
porn	1	0			
not	0.33	0.67			
worker3	porn	not			
porn	1	0			
not	0	1			

vorker4	porn	not
oorn	1	0
not	0	1

worker5	porn	not
porn	0.5	0.5
not	1	0

	worker1	worker2 worker3		worker4	worker5	
sunnyfun	porn	not	not	not	porn	
sex- mission	porn	porn	porn	porn	porn	
google	not	porn	not	not	porn	
youporn	porn	porn	porn	porn	not	
yahoo	porn	not	not	not	porn	

	porn	not
sunnyfun	0	1
sex-mission	1	0
google	0	1
youporn	1	0
yahoo	0	1

	-	
worker1	porn	not
porn	1	0
not	0.67	0.33
worker2	porn	not
porn	1	0
not	0.33	0.67
worker3	porn	Red
porn	1	0 t
not	0	1
worker4	porn	not
porn	1	0
not	0	1
worker5	porn	not
porn	0.5	0.5
not	1	0



		worker1	worker2	worker3	wor	ker4	worker5
	sunnyfun	porn	not	not	not		porn
	sex- mission	porn	porn	porn	por	n	porn
	google	not	porn	not	not		porn
	youporn	porn	porn	porn	por	n	not
	vatioo	porn	not	not	not		porn
				porn		not	
Iterate until		sunnyfur	1	0.26		0.74	
convero	and the second	sex-missi	on	0.69		0.31	
		google		0.29		0.71	
		youporn		0.82		0.18	
🛪 Penn Engineering		yahoo		0.26		0.74	

worker1	porn	not
porn	1	0
not	0.67	0.33

worker2	porn	not
porn	1	0
not	0.33	0.67

worker3	porn	not
porn	1	0
not	0	1

-

worker4	porn	not
porn	1	0
not	0	1

worker5	porn	not
porn	0.5	0.5
not	1	0



How would you use gold standard data in the EM process?



EM Algorithm

Re-Calculate Worker Scores over two steps:

- 1. Estimate the probability that each answer is correct, using labels from multiple workers weighted by the probability that they are correct
- 2. Estimate the quality of the workers by comparing their submitted answers to the inferred correct answers

Confusion Matrix Gives us Worker Error

From the confusion matrix we can measure the overall error rate for each worker Sum of the non-diagonal elements of the confusion matrix (weighted by the priors) This results in a single, scalar value as the quality score for each worker

Worker error

worker1	porn	not	
porn	1	0	.67
not	0.67	0.33	

worker3	porn	not	
porn	1	0	0
not	0	1	

worker5	porn	not
porn	0.5	0.5
not	1	0

worker2	porn	not	
porn	1	0	22
not	0.33	0.67	

worker4	porn	not	
porn	1	0	0
not	0	1	

1.5

🐯 Penn Engineering

Advanced Topics

Bias versus error

How noisy can the workers be and still allow us to still converge to a correct solution?

Bias versus error

Error rate alone is not sufficient to measure the inherent value of a worker.

For example, workers may be careful but biased

In a non-binary case, this is more apparent

What if instead of asking our workers to label sites porn or not porn, we asked them to label the G, PG, R, X?

Bias versus error

Parents with young children tend to be more conservative They tend to classify PG-rated sites as R-rated sites, and R-rated sites as X-rated. Such workers give consistently and predictably incorrect answers It is possible to automatically correct for bias

Implications

Unlike with spammers, with biased workers it is possible to "reverse" the errors

We can recover a label assignment of much higher quality

In the presence of systematic bias, the naive measurement of error rate results in underestimates of the true quality of the worker

This potentially leads to incorrect rejections and blocks of legitimate workers



For more details

Check out two papers by Panos Ipeirotis and his collaborators

Managing Crowdsourcing Workers discusses separating error and bias

Get Another Label? Improving Data Quality and Data Mining Using Multiple, Noisy Labelers

discusses how noisy judgements can be, with us still getting good quality results

NETS 213: CROWDSOURCING AND HUMAN COMPUTATION

Term Project





Self-Designed Term Project

The term project is a self-designed team project.

Team size must be 4-6.

Project is open ended but must do the following:

- 1. use a crowd,
- 2. have a quality control and an aggregation component,
- 3. include either a machine learning component or human-computer interaction component,
- 4. include an analysis.

There will be deliverables throughout the remainder of the semester to help you build up your project incrementally.



Examples of past projects

- <u>PictureThis</u> uses crowdsourcing to have the crowd write new version of picture books.
- <u>Shoptimum</u> is a crowdsourced fashion portal to get links to cheaper alternatives for celebrity clothes and accessories.
- <u>Crowd Art</u> capturing the beauty and majesty of horses through art
- <u>Picture Perfect</u> let the crowd help you pick the best photo from your vacation pics
- <u>Textel</u> a collaboratively written text adventure game