

Modeling Annotator Expertise

-Learning when everybody knows a bit of something

Y. Yan¹ R. Rosales² G. Fung² M. Schmidt³
G. Hermosillo² L. Bogoni² L. Moy⁴ J. Dy¹

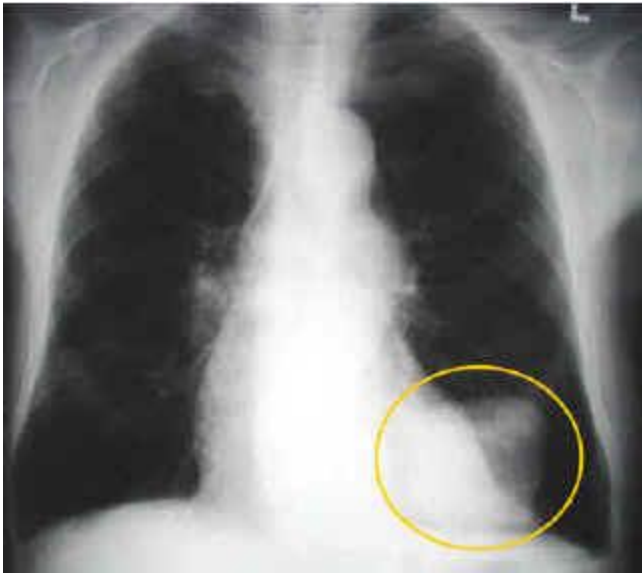
¹Northeastern Univ., Boston, MA USA ²Siemens Healthcare, Malvern, PA USA

³Univ. of British Columbia, Vancouver, BC Canada

⁴Univ. of Penn., Philadelphia, PA USA

Motivation

- *Multiple Expert Diagnoses*
- *Amazon Mechanical Turk*



amazonmechanical turk
Artificial Intelligence

YOUR HITS
Design Publish Manage Resource Center

Resource Center > Use Cases

Mechanical Turk is the new way to outsource information work.

Mechanical Turk Use Cases

- < Get Things Done Faster
- < Staff Projects Instantly
- < Data Management
- < Enforce Site Guidelines
- < Create Content
- < You're In Control
- < Faster, Cheaper...Better

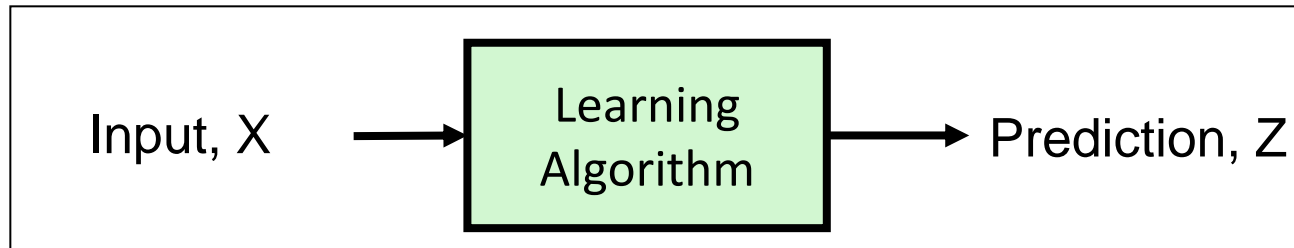
1. *How should the patients be diagnosed when doctors disagree?*
2. *How do we evaluate the doctors' diagnoses?*

Model Assumptions

- 1. Multiple yet unreliable annotators.
- 2. Varying performance on types of data.
 - *Due to different expertise.*
 - *Due to quality of data.*

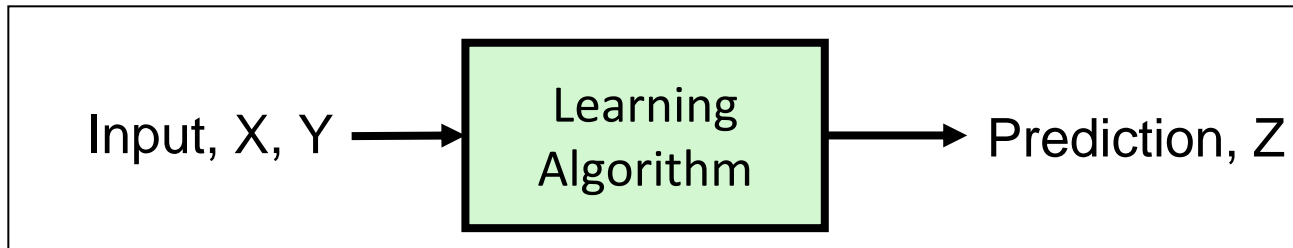
Typical Classification Problem

	Age	Temp.	Symptoms...	Z
Patient 1	1	96	...	not sick
Patient 2	50	102	...	sick
...			...	
Patient N	65	95	...	not sick

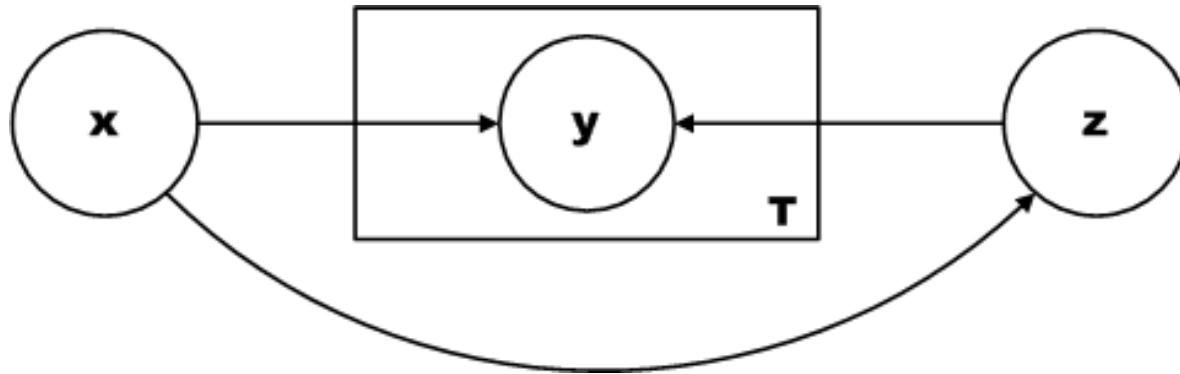


The Multiple Expert Problem

	Age	Temp.	Symptoms...	Ann. Y_1	Ann. Y_2	Ann. ...	Ann. Y_T
Patient 1	1	96	...	not sick	sick	...	sick
Patient 2	50	102	...	sick	sick	...	sick
...			...				
Patient N	65	95	...	not sick	not sick	...	sick



Probabilistic Model



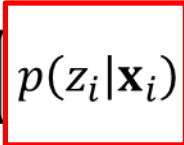
$$p(Y, Z|X) = \prod_i p(z_i|\mathbf{x}_i) \prod_t p(y_i^{(t)}|\mathbf{x}_i, z_i)$$

- $\mathbf{x}_i \in \mathbf{R}^D, i = 1, 2, \dots, N$ observations;

- $y_i^{(t)} \in \mathbf{R}, i = 1, 2, \dots, N; t = 1, 2, \dots, T$ annotation by t for sample i ;

- $z_i \in \mathbf{R}, i = 1, 2, \dots, N$ true (hidden) label for sample i .


Probabilistic Model

$$p(Y, Z|X) = \prod_i p(z_i|\mathbf{x}_i) \prod_t p(y_i^{(t)}|\mathbf{x}_i, z_i)$$


$$p(z = 1|\mathbf{x}) = (1 + \exp(-\boldsymbol{\alpha}^T \mathbf{x} - \beta))^{-1}$$

Classifier: Logistic regression model

Probabilistic Model

$$p(Y, Z|X) = \prod_i p(z_i|\mathbf{x}_i) \prod_t p(y_i^{(t)}|z_i)$$


Bernoulli Model:

$$p(y_i^{(t)}|z_i) = (1 - \eta_t)^{|y_i^{(t)} - z_i|} \eta_t^{1 - |y_i^{(t)} - z_i|}$$

η_t : Probability of labeler t to be correct

Gaussian Model:

$$p(y_i^{(t)}|z_i) = N(y_i^{(t)}; z_i, \sigma_t)$$

σ_t : How labeler t deviates from the true label z

Probabilistic Model

when annotator's performance vary with data

$$p(Y, Z|X) = \prod_i p(z_i|\mathbf{x}_i) \prod_t p(y_i^{(t)}|\mathbf{x}_i, z_i)$$

Bernoulli Model:

$$p\left(y_i^{(t)}|\mathbf{x}_i, z_i\right) = (1 - \eta_t(\mathbf{x}_i))^{|y_i^{(t)} - z_i|} \eta_t(\mathbf{x}_i)^{1 - |y_i^{(t)} - z_i|}$$

$$\eta_t(\mathbf{x}) = (1 + \exp(-\mathbf{w}_t^T \mathbf{x} - \gamma_t))^{-1}$$

Gaussian Model:

$$p\left(y_i^{(t)}|\mathbf{x}_i, z_i\right) = N\left(y_i^{(t)}; z_i, \sigma_t(\mathbf{x}_i)\right)$$

$$\sigma_t(\mathbf{x}) = (1 + \exp(-\mathbf{w}_t^T \mathbf{x} - \gamma_t))^{-1}$$

Implementation

Maximum Likelihood Estimation:

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \prod_t \prod_i p(y_i^{(t)} | \mathbf{x}_i; \theta) \\ &= \arg \max_{\theta} \sum_t \sum_i \log \sum_{z_i} p(y_i^{(t)}, z_i | \mathbf{x}_i; \theta)\end{aligned}$$

Since z_i are *hidden*, **EM** algorithm is utilized:

E-step Compute:

$$\begin{aligned}\tilde{p}(z_i) &\triangleq p(z_i | \mathbf{x}_i, y_i) \\ &\propto p(z_i, y_i | \mathbf{x}_i) \\ &= \prod_t p(y_i^{(t)} | \mathbf{x}_i, z_i) p(z_i | \mathbf{x}_i)\end{aligned}$$

M-step Maximize:

$$\sum_t \sum_i E_{\tilde{p}(z_i)} [\log p(y_i^{(t)}, z_i | \mathbf{x}_i)]$$

to update $\tilde{\theta} = \{\alpha, \beta, \{\mathbf{w}_t\}, \{\gamma_t\}\}$

Insights on Classification Model

$$LLR(\{y^{(t)}\}, \mathbf{x}) = \log \frac{p(z = 1 | \{y^{(t)}\}, \mathbf{x})}{p(z = 0 | \{y^{(t)}\}, \mathbf{x})}$$

- **Bernoulli Case**

$$LLR = \alpha^T \mathbf{x} + \beta + \sum_t (-1)^{(1-y^{(t)})} (\mathbf{w}_t^T \mathbf{x} + \gamma_t)$$

by general learnt classifier

by each annotator

- **Gaussian Case**

$$LLR = \alpha^T \mathbf{x} + \beta + T^+ - T^- + \sum_t (-1)^{(1-y^{(t)})} \exp(-\mathbf{w}_t^T \mathbf{x} - \gamma_t)$$

by general learnt classifier

by each annotator

Missing Annotators

- *When not all annotators provided a label for a particular sample, the true label is predicted based on:*

- 1.
$$p(z|\{y^{t \setminus k}\}, \mathbf{x}) = \frac{\prod_{t \setminus k} p(y^{(t)}|z, \mathbf{x})p(z|\mathbf{x})}{\sum_z \prod_{t \setminus k} p(y^{(t)}|z, \mathbf{x})p(z|\mathbf{x})}$$

- 2.
$$p(z = 1|\mathbf{x}) = (1 + \exp(-\boldsymbol{\alpha}^T \mathbf{x} - \beta))^{-1}$$

Predicting Ground Truth without Observation

- *Estimate hidden label purely on annotations when observation is not available.*
- $$p(z|\{y^{(t)}\}) = \int \prod_t p(y^{(t)}|z, \mathbf{x}) p(z|\mathbf{x}) d p(\mathbf{x})$$
$$\approx \frac{1}{S} \sum_{s=1}^S p(z|\mathbf{x}_s) \prod_t p(y^{(t)}|z, \mathbf{x}_s)$$

Approximation is reached by sampling.

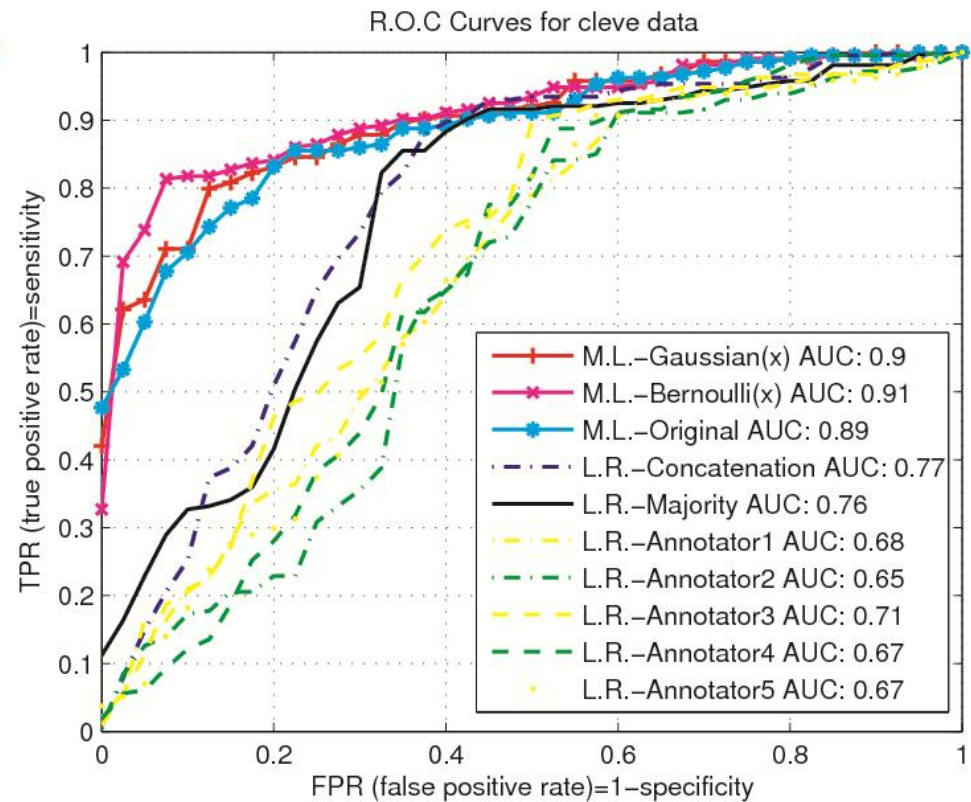
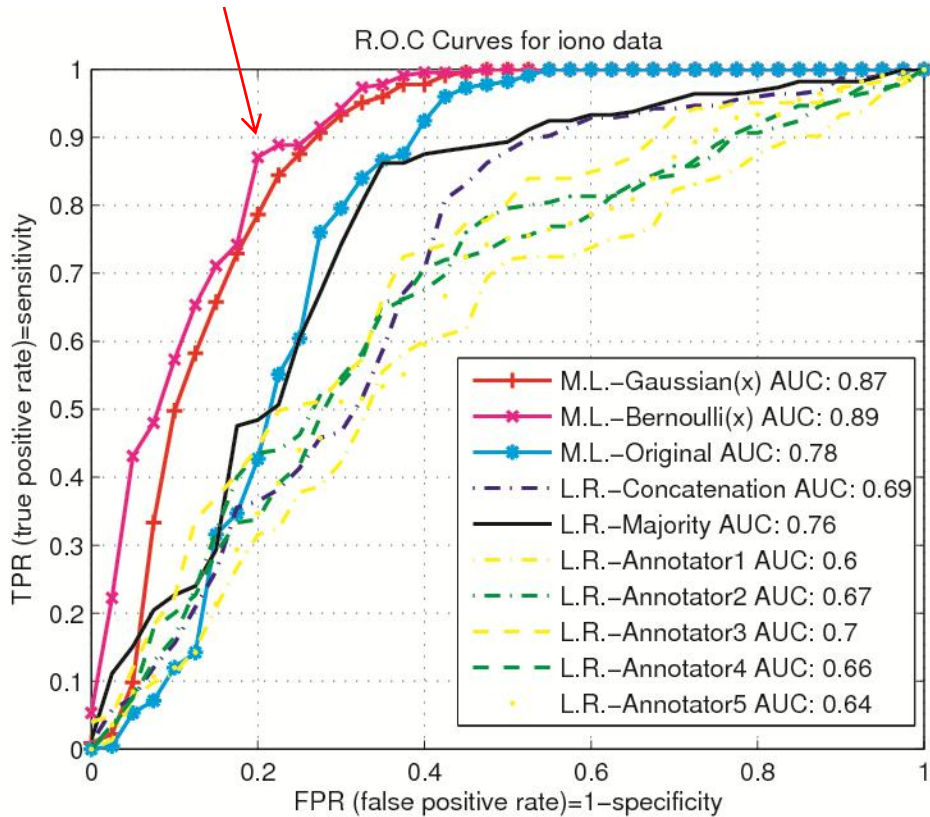
Evaluating Annotators

- *Is it possible to evaluate annotators without ground truth?*

- $$p(y^{(k)} | \{y^{(t \setminus k)}\}, \mathbf{x}) = \frac{p(\{y^{(t)}\} | \mathbf{x})}{p(\{y^{(t \setminus k)}\} | \mathbf{x})}$$
$$= \frac{\sum_z p(\{y^{(t)}\} | z, \mathbf{x}) p(z | \mathbf{x})}{\sum_z p(\{y^{(t \setminus k)}\} | z, \mathbf{x}) p(z | \mathbf{x})}$$

UCI Data Classification

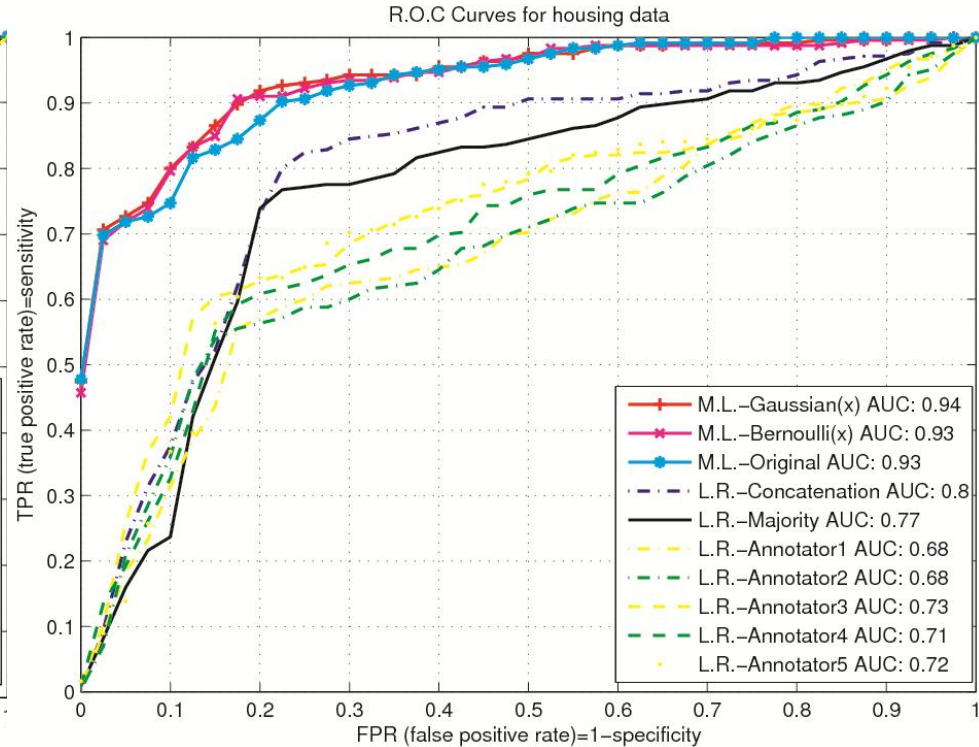
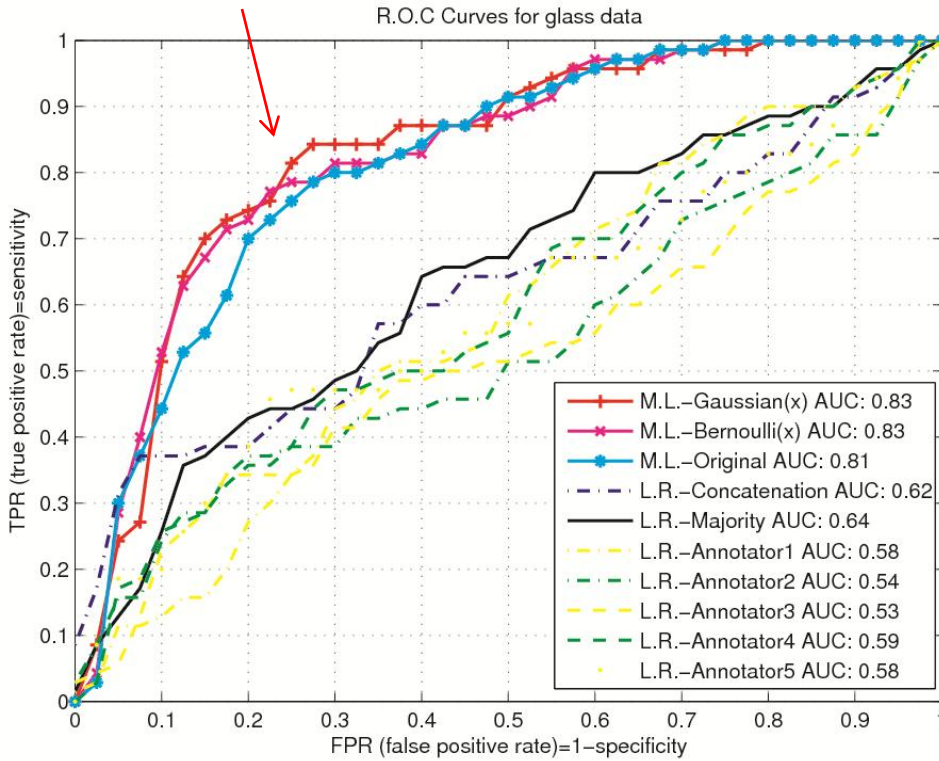
Our method (x)



Data tested: Ionosphere, Cleveland Heart.

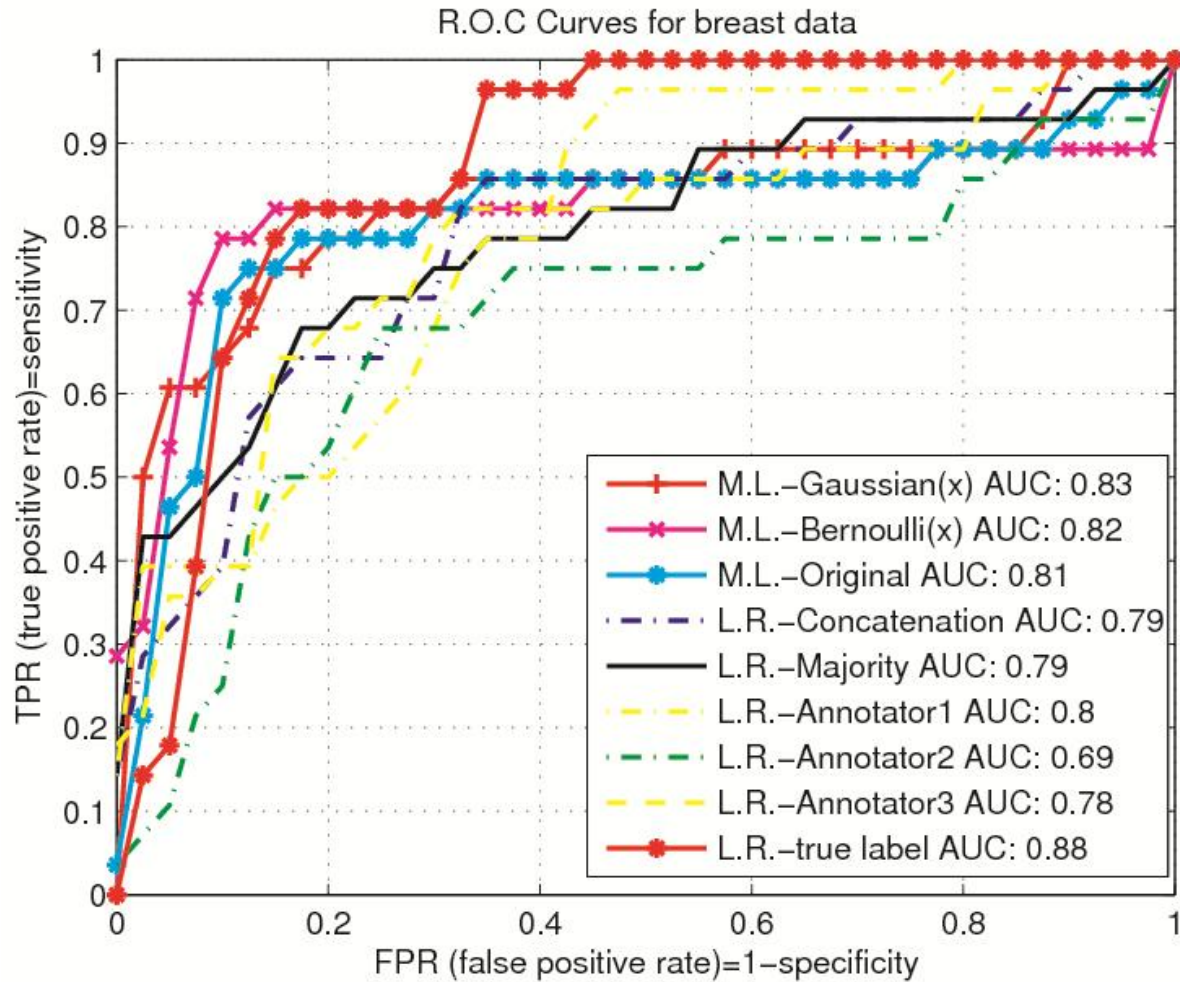
UCI Data Classification

Our method (x)



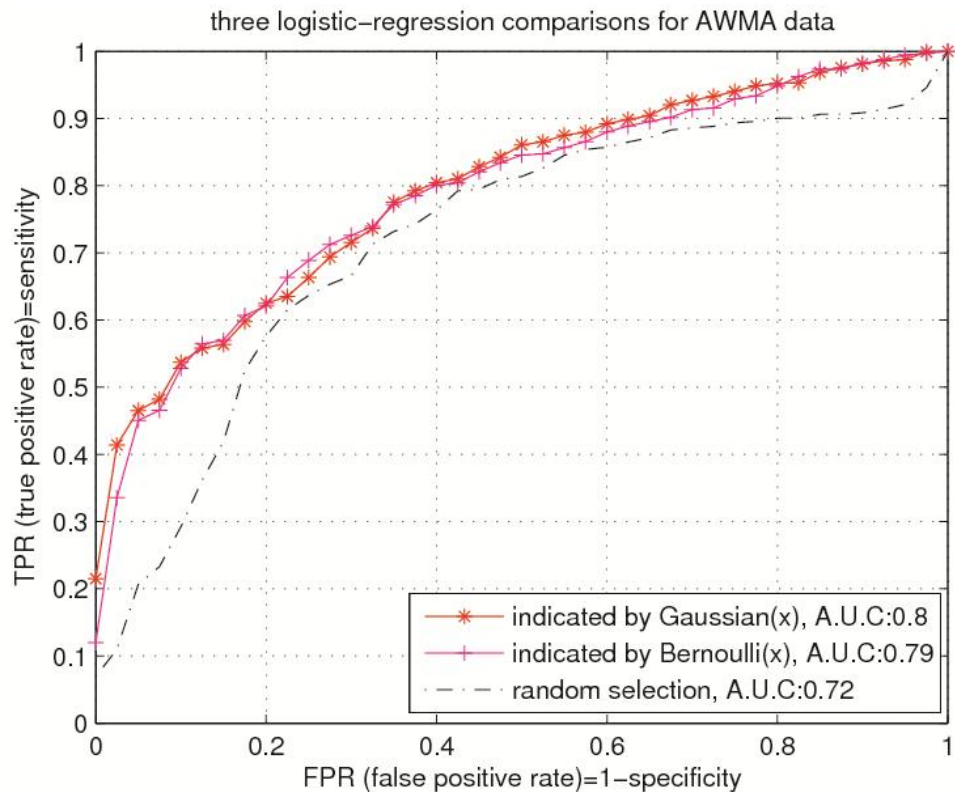
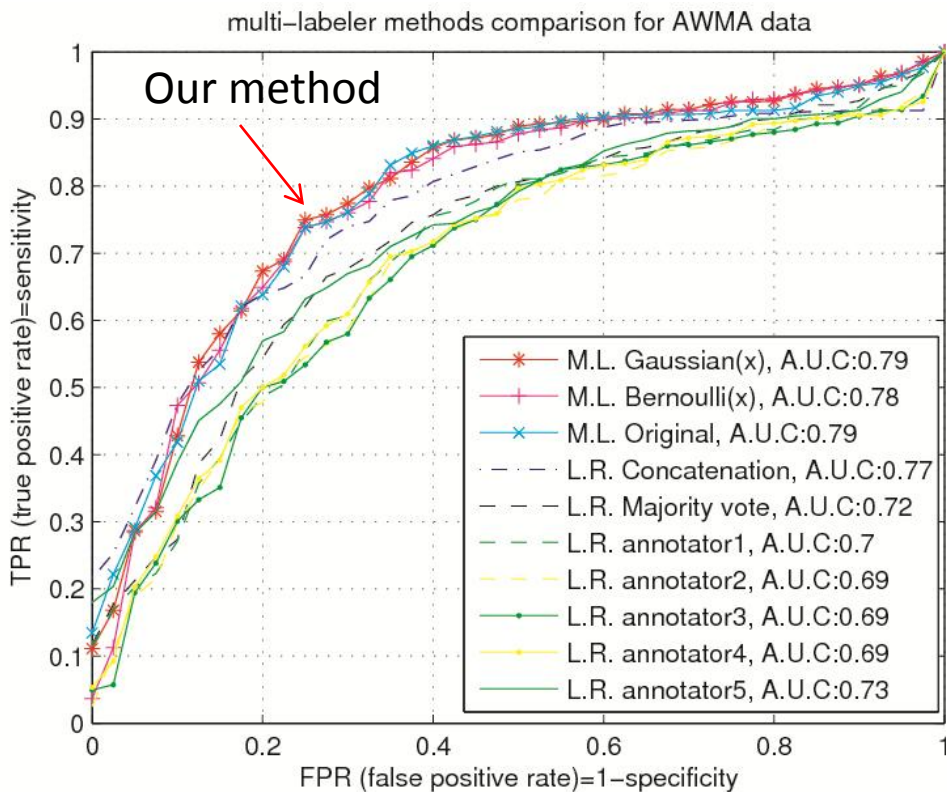
Data tested: Glass, and Housing.

Breast Cancer Detection



75 cases, 8 morphological features, 3 annotators (radiologists)

Cardiac Wall Motion Anomaly Detection



220 cases, 16 LV heart wall segments, 5 annotators (doctors), binary labels (-/+1)

Conclusions

- We provided a probabilistic model that allows learning from multiple annotators whose annotations may be noisy;
- Our model takes into account that the quality of annotation may vary with data;
- We show that this model can deal with missing annotators/data;
- Our model can also be utilized to evaluate annotators even when ground truth is not available; and
- We can also utilize our model to select the most trustworthy/accurate annotator for each new instance labeling.

Thanks for Listening

Questions?

References

- K. Crammer, et. al. (2008). Learning from multiple sources. *J. of Machine Learning Research*, 9: 1757-1774
- O. Dekel, and O. Shamir (2009). Good learners for evil teachers. In *Int. Conf. on Machine Learning*.
- J. Howe (2008). *Crowdsourcing: why the power of the crowd is driving the future of business*. Crown Business
- R. Jin, and Z. Ghahramani (2003). Learning with multiple labels. In *Adv. Neural Inf. Processing Systems*.
- V. Raykar, and et. al. (2009). Supervised learning from multiple experts: whom to trust when everyone lies a bit. In *Int. Conf. on Machine Learning*.
- V. S. Sheng, and et. al. (2008). Get another label? Improve data quality and data mining using multiple, noisy labelers. In *Knowledge Discovery and Data Mining (KDD)*.
- R. Snow, and et. al. (2008). Cheap and fast – but is it good? Evaluating non-expert annotations for natural language tasks. In *Conf. Empirical Methods on Natural Language Processing (EMNLP)*.