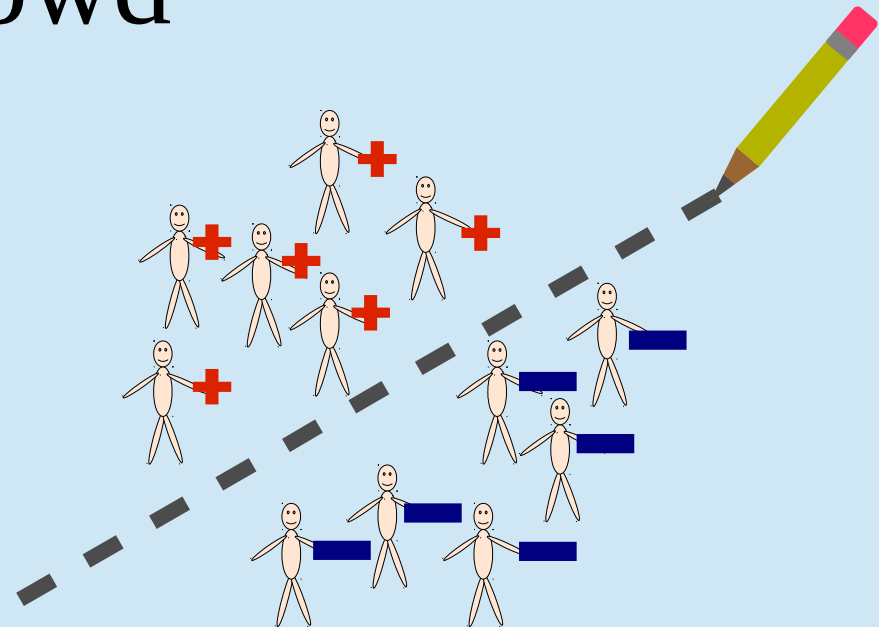


Toward Buying Labels From the Crowd

Jacob Abernethy Michigan
Yiling Chen Harvard
Chien-Ju Ho UCLA/Harvard
Bo Waggoner Harvard

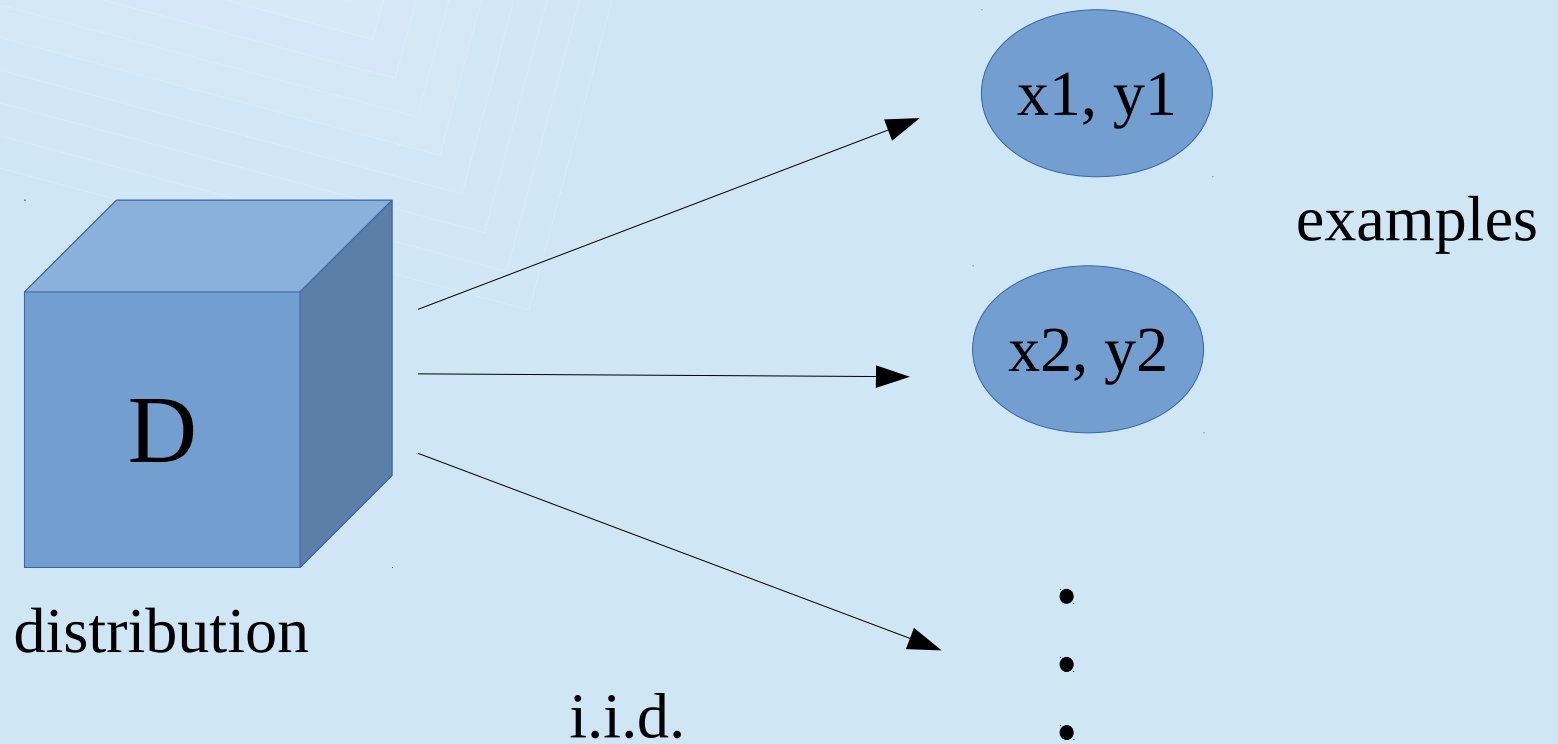


Indo-US Lectures Week in Machine Learning, Game Theory and Optimization
9th January 2014

Outline

- General setting
- Related work
- Our approach

Learning Setting



Learning Setting

x_1 $\xrightarrow{\text{hypothesis}}$ $h(x_1)$

$h(x_1), y_1$ $\xrightarrow{\text{loss function}}$ $\text{Loss}(h(x_1), y_1)$

x_1, y_1

x_2, y_2

•
•
•

examples

Learning Setting

x_1 \longrightarrow $h(x_1)$
hypothesis

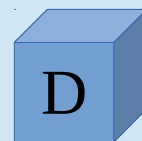
$h(x_1), y_1$ \longrightarrow $\text{Loss}(h(x_1), y_1)$
loss function

x_1, y_1

x_2, y_2

examples

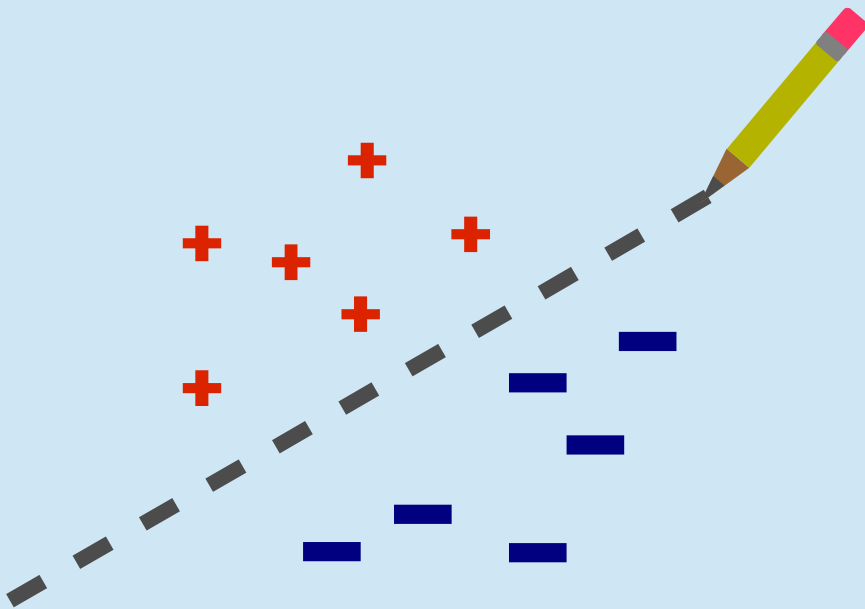
Goal: from **few** examples,
pick a hypothesis with **small loss**
(in expectation, with high probability) w.r.t.



-
-
-

Example 1: Classification

x = point in the plane
 y = “+” or “-”
hypothesis = line
loss = 0 if correct, 1 if incorrect
or in $[-1, 1]$ weighted by distance



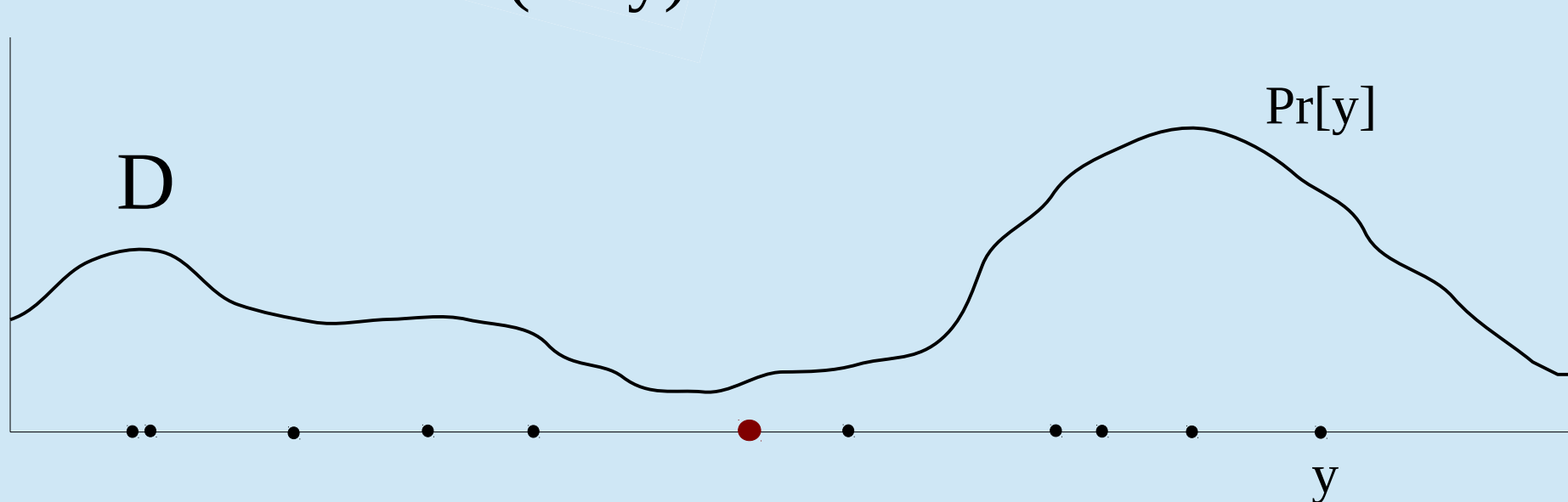
Example 2: Estimate the mean

x = doesn't matter (e.g. always zero)

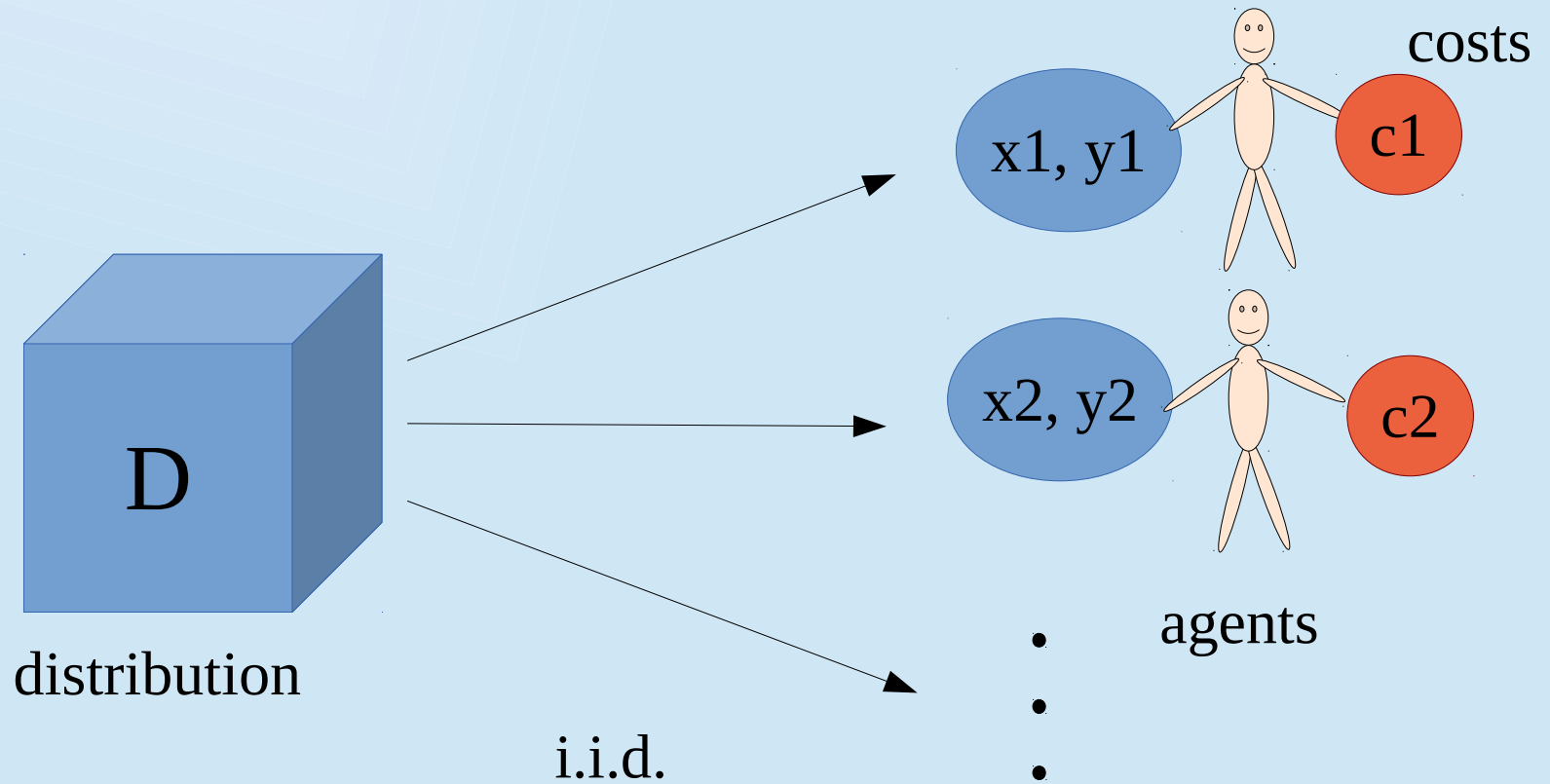
y = real number in $[0,1]$

hypothesis = real number in $[0,1]$

loss = $(h - y)^2$

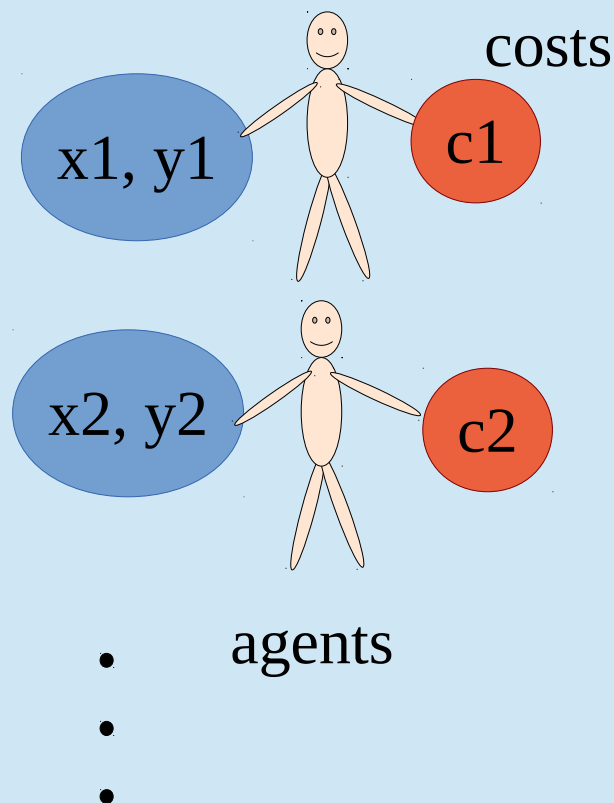


Adding incentives



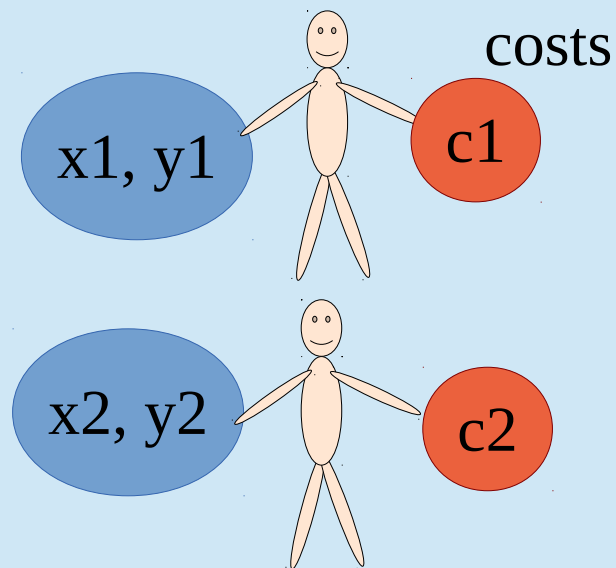
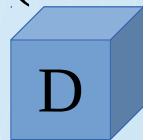
Incentives Setting

- (x_1, y_1, c_1) drawn from D
- Must design **mechanism** and **learning algorithm** together
- Many possible assumptions:
 - costs in $[0,1]$
 - agents cannot misreport (x,y)
 - ...



Goal

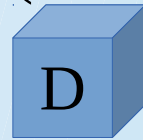
Goal: with **small budget**, purchase data and pick a hypothesis with **small loss** (in expectation, with high probability) w.r.t.



- agents
-
-

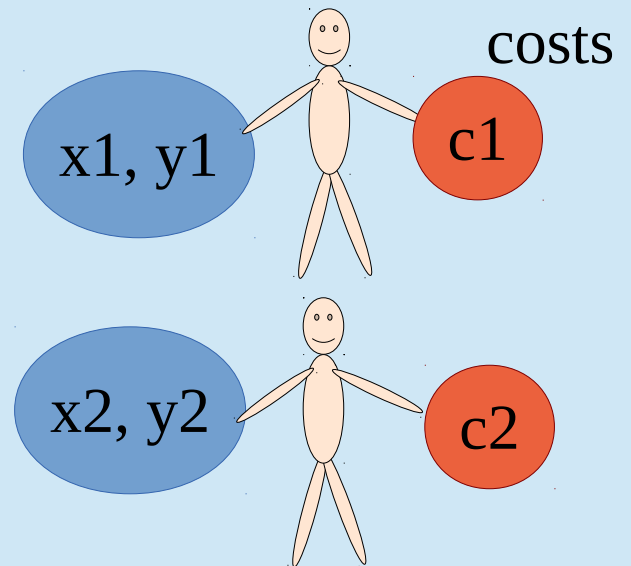
Goal

Goal: with **small budget**, purchase data and pick a hypothesis with **small loss** (in expectation, with high probability) w.r.t.



Naive approach:
Offer B of the agents
a price of 1 (maximum).

→ Seems non-obvious how to
improve on this!



• agents

-
-

Three possible avenues

1. **Centralized/simultaneous:**

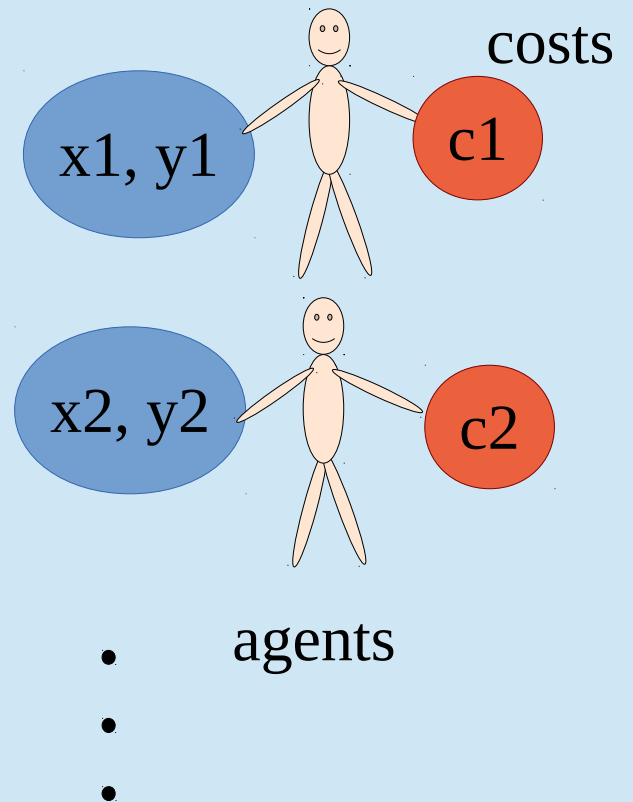
auction of some sort.

2. **Decentralized/simultaneous:**

survey offered to all agents.

→ Both miss interactions in the data!

3. **Iterative** (but perhaps myopic).



Digression: Importance Weighting

Goal: compute sum of y_1, y_2, \dots, y_n .

Twist: each y_i is observed independently with probability p_i .

So: estimate sum = $\frac{y_1}{p_1} + \frac{y_3}{p_3} + \frac{y_4}{p_4} + \dots$

Digression: Importance Weighting

Goal: compute sum of y_1, y_2, \dots, y_n .

Twist: each y_i is observed independently with probability p_i .

So: estimate sum = $\frac{y_1}{p_1} + \frac{y_3}{p_3} + \frac{y_4}{p_4} + \dots$

Can apply *Hoeffding*: Given independent Y_1, \dots, Y_n , with Y_i in $[0, b_i]$:

Let $d = \Pr[|\sum_i Y_i - \text{expectation}| > \text{eps}]$,

Then $d < 2\exp[-2 \text{eps}^2 / \sum_i b_i^2]$.

Or, if I want probability $1-d$, then I get error $\text{eps} < \sqrt{\frac{\ln(2/d) \sum_i b_i^2}{2}}$

Outline

- General setting
- Related work
- Our approach

Conducting Truthful Surveys, Cheaply

Roth and Schoenebeck, EC 2011.

Problem: Estimate the mean.

Assumptions:

- marginal on costs, F , is known.
- **decentralized/simultaneous (survey) approach.**

Goal: unbiased estimator with minimum (or close to minimum) *worst-case expected variance*.

(*worst-case*: over all distributions D whose cost marginal is F .)

(*expected*: over the data points drawn from D .)

(*variance*: over the randomization of the mechanism.)

Conducting Truthful Surveys, Cheaply

Roth and Schoenebeck, EC 2011.

Results:

- WLOG to consider “Take-It-Or-Leave-It” posted price mechanisms.
→ **Reduces the problem to picking a single posted-price distribution.**
- Must assume agents then report true costs!
- Describes posted-price distribution giving unbiased estimator with close to minimum *worst-case expected variance*.

Conducting Truthful Surveys, Cheaply

Roth and Schoenebeck, EC 2011.

What we want to do differently:

- More complex learning problems.
- ***Iterative*** rather than their ***simultaneous/decentralized*** approach.
- Generalization-error type bounds.

Importance-Weighted Active Learning

Beygelzimer, Dasgupta, and Langford, ICML 2009.

Problem: Learn while buying a small *number of labels*.

Assumptions:

- All costs are 1.
- Algorithm **can observe x** before deciding.
- **Iterative approach!**

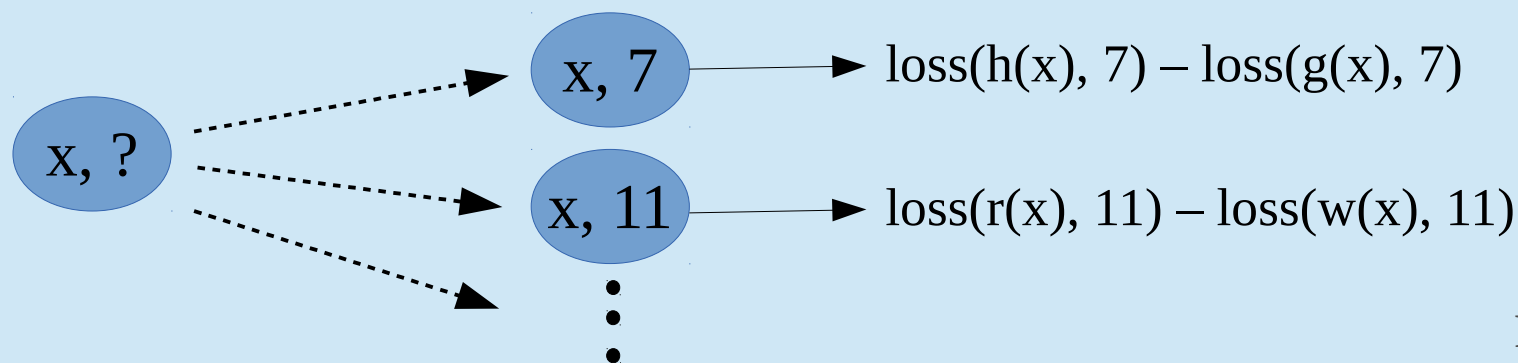
Goal: Buy few labels, compare to if we'd bought *all* labels.

Importance-Weighted Active Learning

Beygelzimer, Dasgupta, and Langford, ICML 2009.

Results:

- **IWAL framework:** for each arriving point, set probability of sampling, then importance-weight losses to get unbiased estimators of expected loss.
- Instantiation: continuously narrow hypothesis set; sampling probability = possibility to distinguish within hypothesis set



Importance-Weighted Active Learning

Beygelzimer, Dasgupta, and Langford, ICML 2009.

What we'd like to do differently:

- Modify **existing** learning algorithms and (hopefully) leverage their guarantees.
 - We'll use no-regret algorithms.
- Agents have costs in $[0,1]$.
- Not just worst-case guarantees, but understanding when we do well.

Outline

- General setting
- Related work
- Our approach

Our approach

Ideal world:

Here's my **learning** problem, and here's a good online learning algorithm for it!

Abra Kadabra Alakazam!

...

OK, here is a **mechanism** for you to use!

Our approach

Ideal world:

Here's my **learning** problem, and here's a good online learning algorithm for it!

By the way, here's a **regret bound** for that learning algorithm!

Abra Kadabra Alakazam!

...

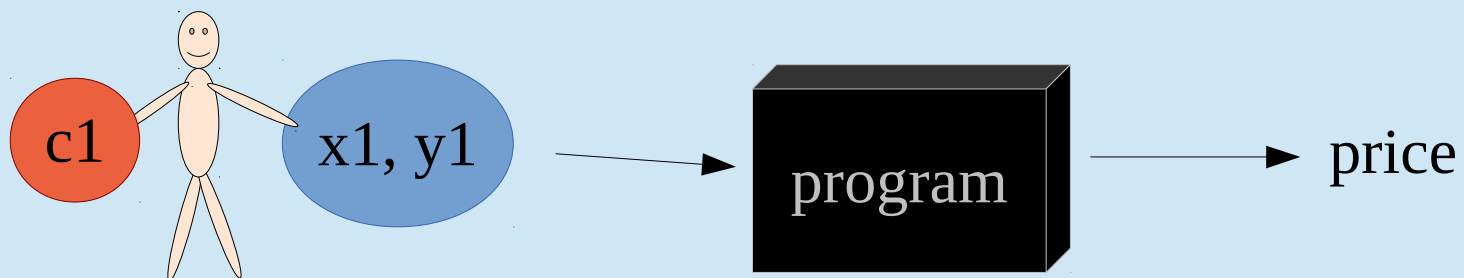
OK, here is a **mechanism** for you to use!

...

OK, here is a **generalization error and budget bound** for that mechanism!

Our approach

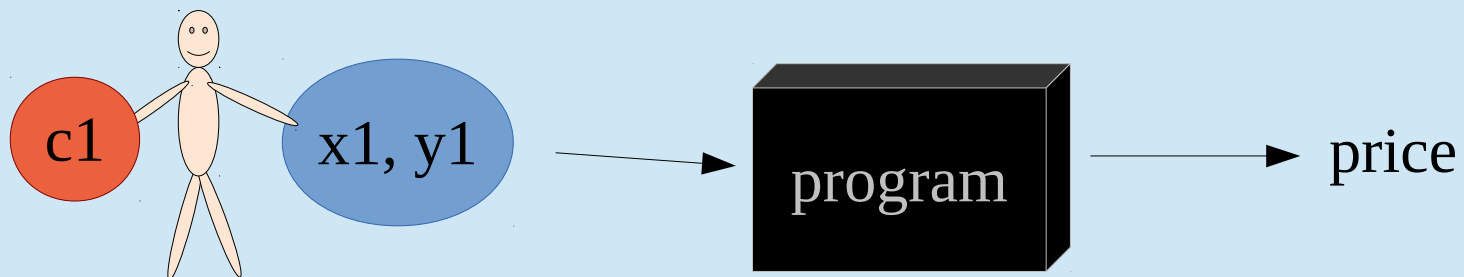
- Key assumption: mechanism can set price based on **both** x and y ! (and agents cannot misreport x,y)
- Example: medical data (difficult to misreport).
- Implementation: give agents a price-calculating program.



General Framework

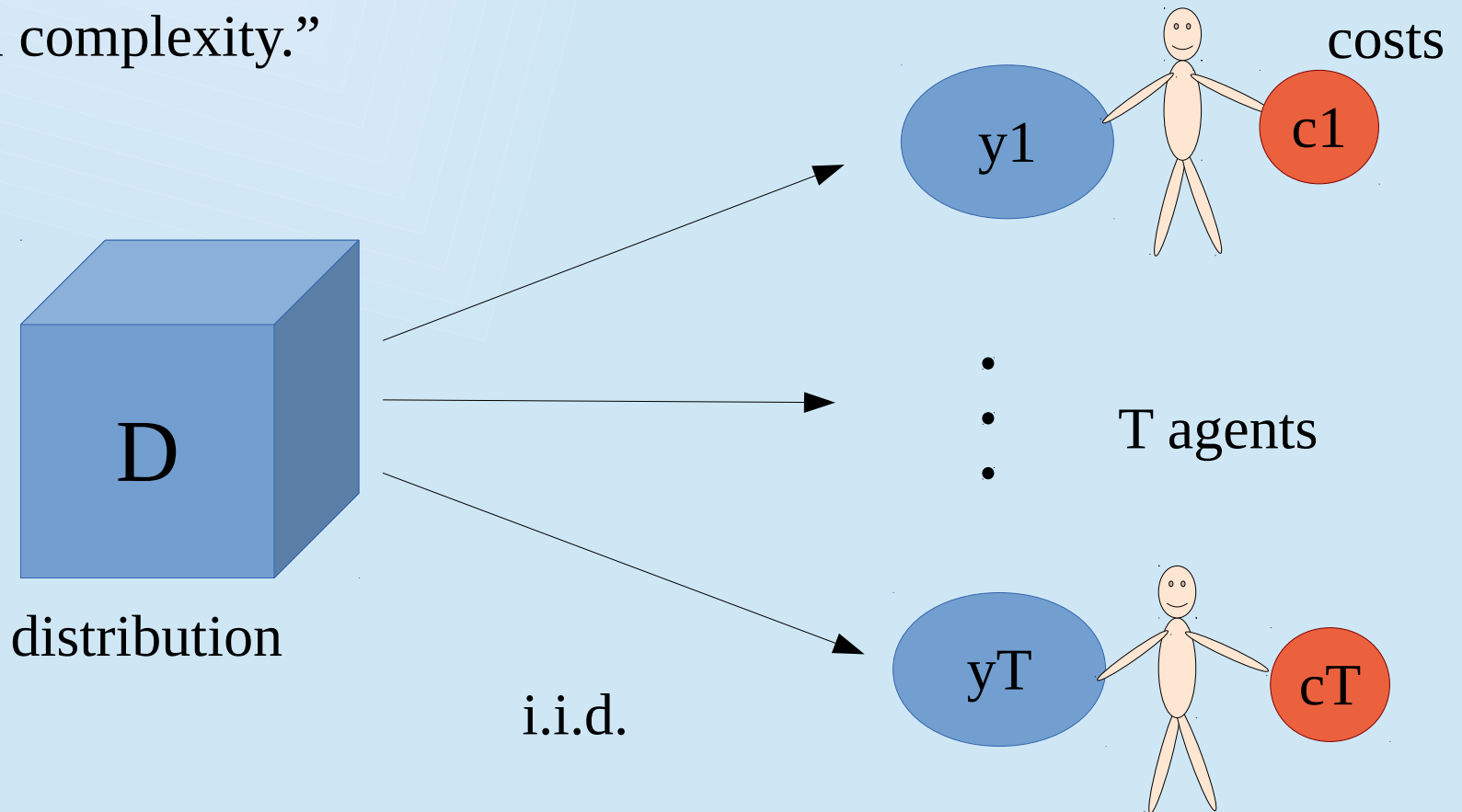
Given a no-regret algorithm for the problem:

1. Decide the **“value”** of the next agent's data point.
2. (Randomly) set a **posted price** based on this value and the marginal cost distribution.
3. If taken, **importance-weight the loss** based on the probability the random price would've been accepted. Update the no-regret algorithm.
4. Repeat.



Simple example: estimate the mean

Assume all costs are 1.
→ “Label complexity.”



Simple example: estimate the mean

Assume all costs are 1.
→ “Label complexity.”

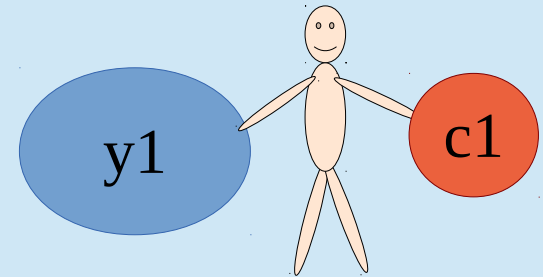
No-regret algorithm: h = sample mean.

Benchmark: buy all T labels.

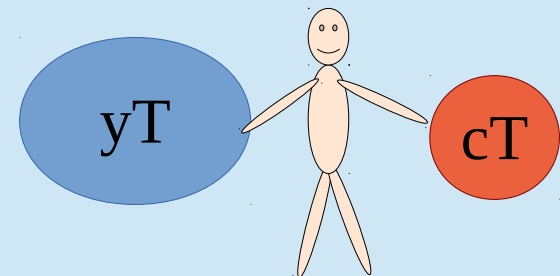
Let u = true mean.

→ with prob. $1-d$, $|h - u| = O\left(\sqrt{\frac{\ln(2/d)}{T}}\right)$

Can we improve somehow??



•
•
•



Applying our framework

1. Decide the “**value**” of the next data point.
2. (Randomly) set a **posted price**.
3. If taken, **importance-weight** and update.
4. Repeat.

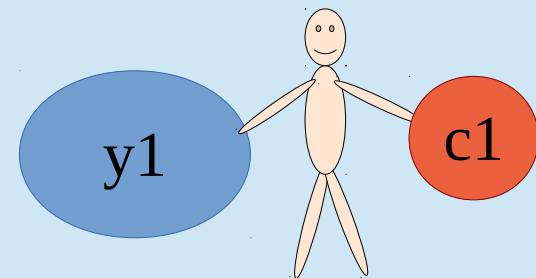
Scheme A:

Set value $p_t = y_t$.

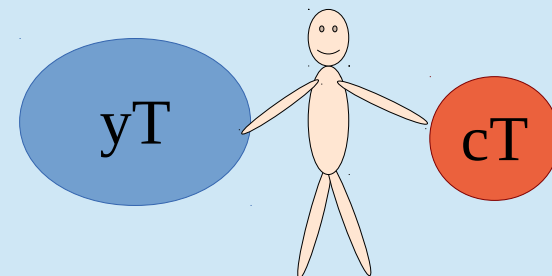
Buy with probability p_t .

→ Error within constant factor of benchmark!

→ Purchase $\sim u_T$ labels! ($u = \text{mean}$)



•
•
•



Applying our framework

1. Decide the “**value**” of the next data point.
2. (Randomly) set a **posted price**.
3. If taken, **importance-weight** and update.
4. Repeat.

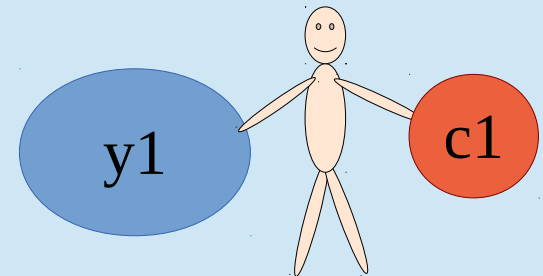
Scheme B:

$$\text{Set value } p_t = |h_t - y_t| + \sqrt{\frac{\ln(T)}{t}}$$

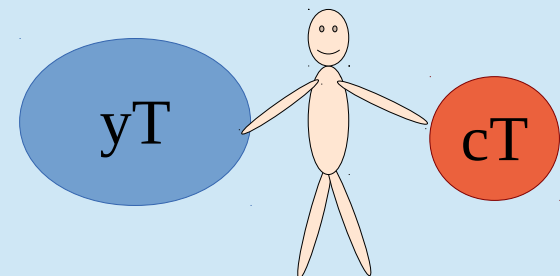
Buy with probability p_t .

(I think) this should give:

- Error “close” to benchmark
- Purchase $\sim o(T)$ labels (o = std deviation)



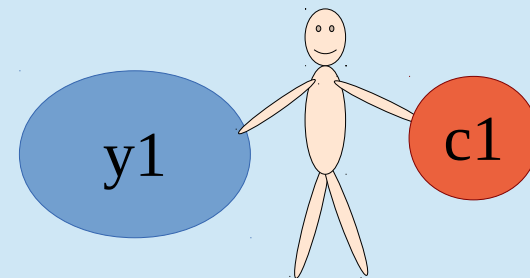
•
•
•



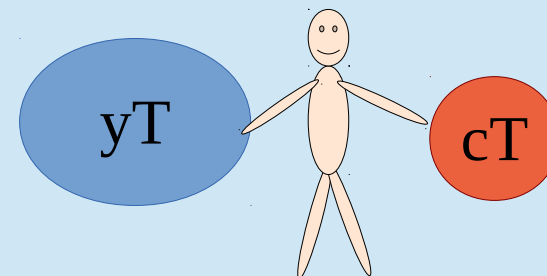
What about costs in $[0,1]$?

- Could compose our mechanism with Roth-Schoenebeck.
- Guarantees? (e.g. spend $\sim uTc$, where c = average cost?)

Seems hard to tell from their analysis, may want another approach.

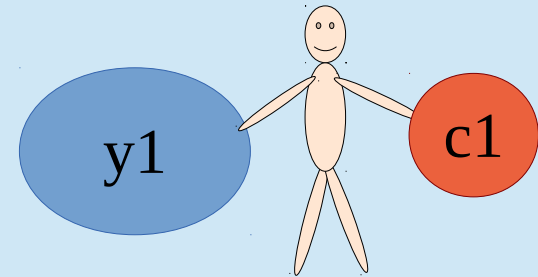


•
•
•

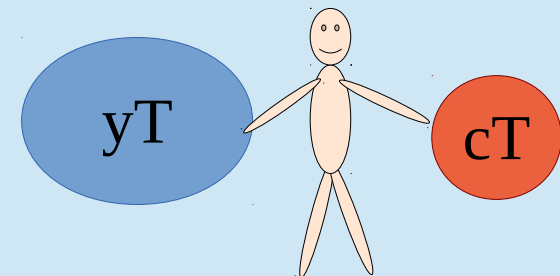


Why this might hopefully work in general

- No-regret algorithms guarantee average **regret** of $1/\sqrt{T}$ or better.
- When drawing examples i.i.d., only want **generalization error** $1/\sqrt{T}$.
- If problem has regret guarantee better than $1/\sqrt{T}$, try to **convert to budget guarantee** while keeping acceptable g.e.



-
-
-



Wrapup of talk

- Problem: seemingly natural but tricky!
- Need to think carefully about **assumptions**.
- Our approach: tweak existing no-regret algorithms, use them to set prices and probabilities.
- When regret is smaller than needed for good generalization error, trade off **regret** and **budget** using **importance-weighting**.
- Todo: understand/prove this “generally”!

