# Using CrowdSourcing for Data Analytics

Hector Garcia-Molina

(work with Steven Whang, Peter Lofgren, Aditya Parameswaran and others)

*Stanford University*

1

---

- Big Data Analytics
- CrowdSourcing

# CrowdSourcing

# Real World Examples

Categorizing Images

Search Relevance

Data Gathering

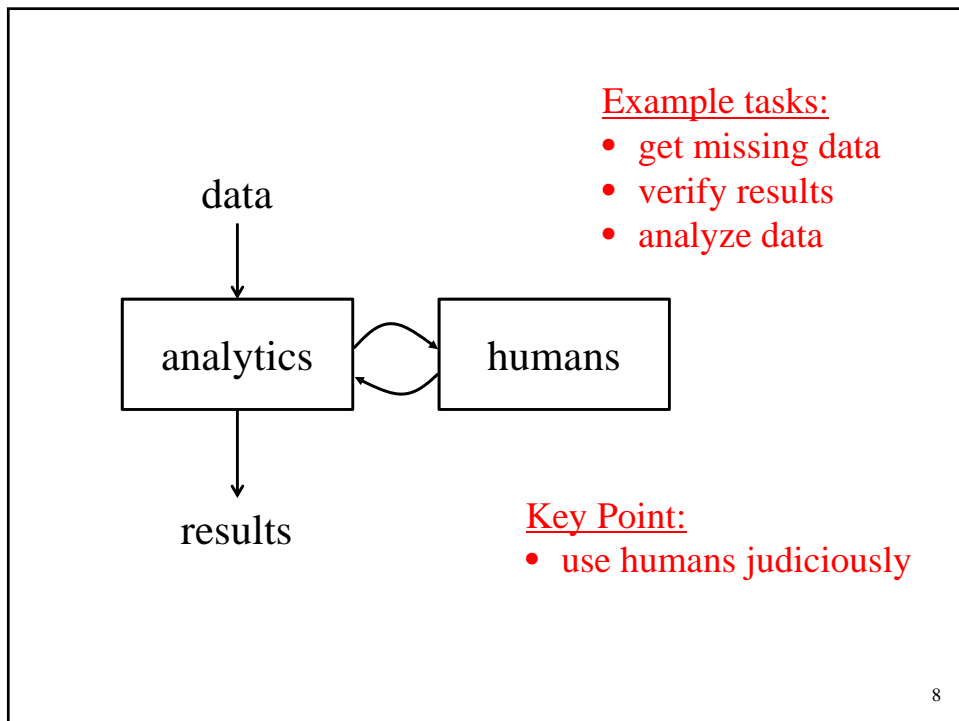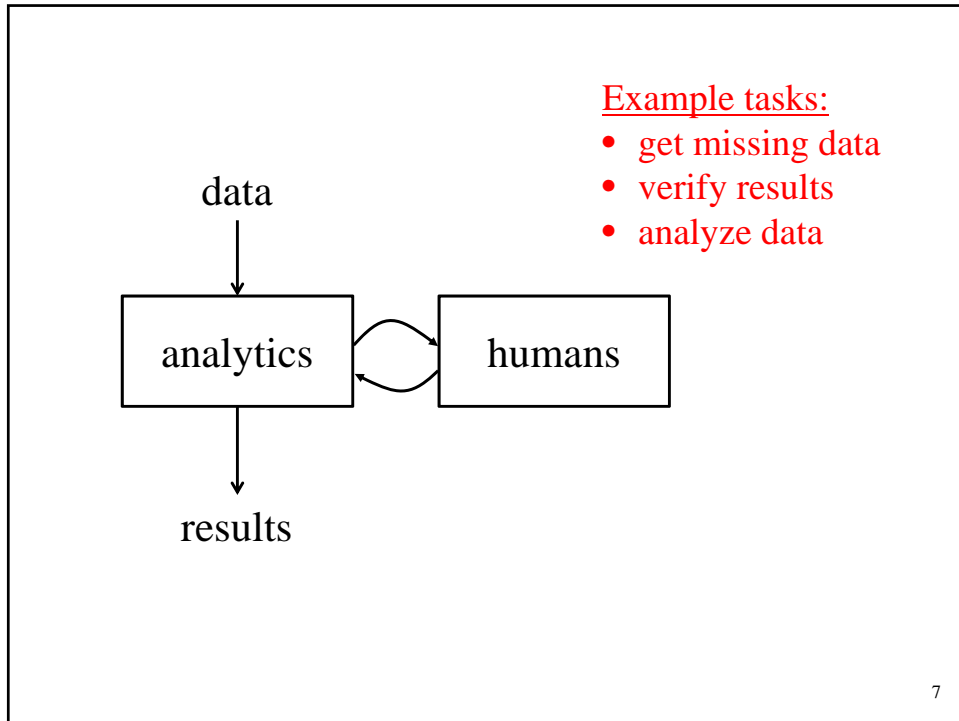CrowdFl wer

Image Matching
Translation

Mission 4636
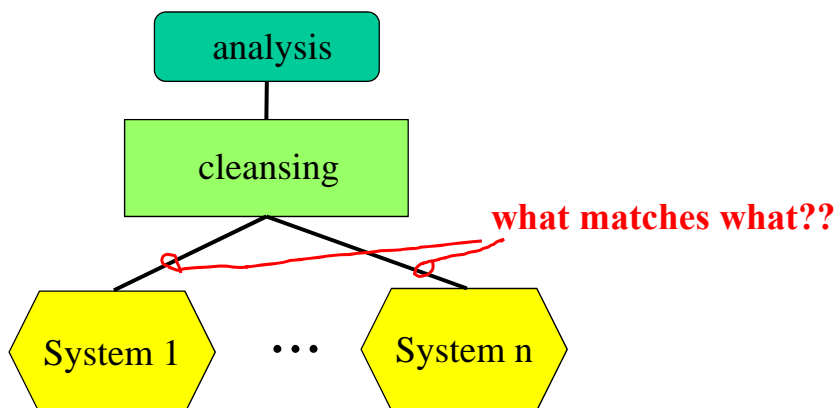
# Many Crowdsourcing Marketplaces!



# Many Research Projects!

Example tasks:
- get missing data
- verify results
- analyze data

data

analytics ⟷ humans

results

7



Example tasks:
- get missing data
- verify results
- analyze data

data

analytics ⟷ humans

results

Key Point:
- use humans judiciously

8

# Today will illustrate with

- Entity Resolution
- (may cover another topic briefly)

# Traditional Entity Resolution

```
          ┌──────────┐
          │ analysis │
          └──────────┘
               │
        ┌─────────────┐
        │  cleansing  │        what matches what??
        └─────────────┘
          ╱        ╲
    ┌──────────┐    ┌──────────┐
    │ System 1 │ ···│ System n │
    └──────────┘    └──────────┘
```

# Why is ER Challenging?

- Huge data sets
- No unique identifiers
- Missing data
- Lots of uncertainty
- Many ways to skin the cat

# Simple ER Example
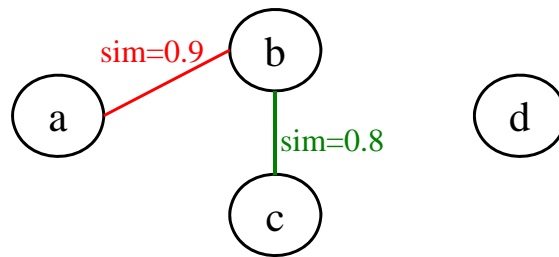
| Record | Name | Address | Phone | Passport | Credit Card |
|--------|------|---------|-------|----------|-------------|
| a | Robert | 123 Main | 1234 | | |
| b | Bob | 123 Main | 1234 | abcd | 777 |
| c | Rob | | | abcd | 777 |
| d | Sally | | 5678 | efgh | 999 |

## Simple ER Example

| Record | Name | Address | Phone | Passport | Credit Card |
|--------|------|---------|-------|----------|-------------|
| a | Robert | 123 Main | 1234 | | |
| b | Bob | 123 Main | 1234 | abcd | 777 |
| c | Rob | | | abcd | 777 |
| d | Sally | | 5678 | efgh | 999 |



13

## Simple ER Example

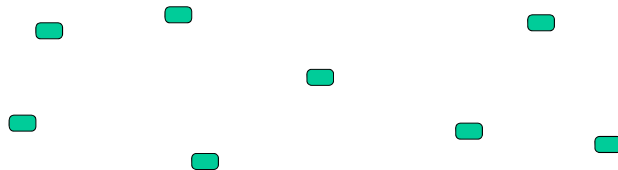| Record | Name | Address | Phone | Passport | Credit Card |
|--------|------|---------|-------|----------|-------------|
| a | Robert | 123 Main | 1234 | | |
| b | Bob | 123 Main | 1234 | abcd | 777 |
| c | Rob | | | abcd | 777 |
| d | Sally | | 5678 | efgh | 999 |



14

# ER: Exact vs Approximate



15

# Simple ER Algorithm

- Compute pairwise similarities
- Apply threshold
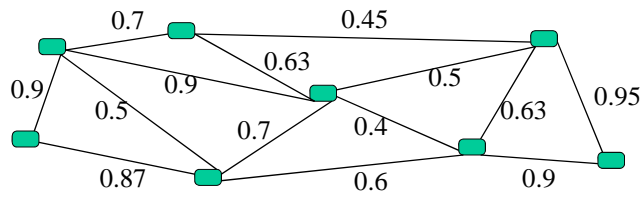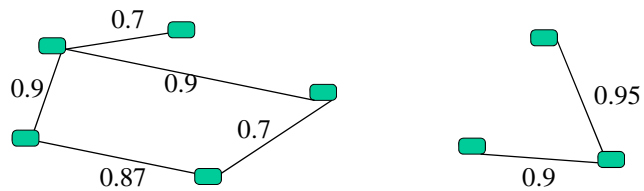- Perform transitive closure



16

8

# Simple ER Algorithm

- Compute pairwise similarities
- Apply threshold
- Perform transitive closure



17

# Simple ER Algorithm

- Compute pairwise similarities
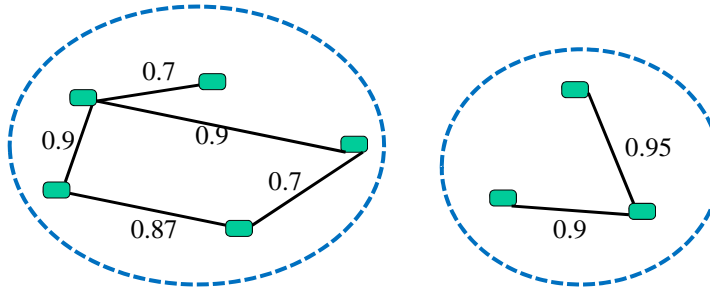- Apply threshold
- Perform transitive closure



threshold = 0.7

18

# Simple ER Algorithm

- Compute pairwise similarities
- Apply threshold
- Perform transitive closure



19

# Crowd ER



20

10

## Same as this?

## Crowd ER

- First Cut: For every pair of records,
  ask workers if they match (i.e., get similarity)

# Crowd ER

- First Cut: For every pair of records,
  ask workers if they match (i.e., get similarity)
- Too expensive!

# Crowd ER

- Second Cut: Compute similarities;
  workers verify "critical" pairs



critical??

# Crowd ER

- Second Cut: Compute similarities; workers verify "critical" pairs



25

# Crowd ER

- Second Cut: Compute similarities; workers verify "critical" pairs



26

13

records

new evidence

| pairwise analysis | generate questions | crowd |

global analysis

Key Point:
• use humans judiciously

clusters

27

# Key Issue: Semantics of Crowd Answer



28

14

Key Issue: Semantics of Crowd Answer

C  D  E

B  A

?

29



Also issue: Similarities as Probabilities

sim(a,b) →  prob(a,b)

30

# Strategy

a

0.9 ╱╲ 0.2

b ─── c
0.5

current state

↓

ER result

use any given ER algorithm

31

# Strategy

a

0.9 ╱╲ 0.2

b ─── c
0.5

current state

Q(a,b)

Q(b,c)

consider ALL possible questions (three in this example)

Q(a,c)

32

16

## Strategy

0.9  a  0.2

b ──── c
   0.5

current state

Q(a,b) — Y → new state → ER result
       — N → new state → ER result

Q(b,c) — Y → new state → ER result
       — N → new state → ER result

Q(a,c) — Y → new state → ER result
       — N → new state → ER result

consider possible outcomes

33

## Strategy

0.9  a  0.2

b ──── c
   0.5

current state → Q(b,c) — Y → new state → ER result

0.9  a  0.2

b ──── c
   1.0

a

b ──── c

example

34

## Strategy



**Two Remaining Issues**

- How do we score an ER result?



F score

- Efficiency?

Gold Standard?

a
0.9 ⌃ 0.2
b —— c
0.5

a
1.0 ⌃ 0.2
b —— c
0.6

→ sim to prob

37



Gold Standard?

possible worlds

a
0.9 ⌃ 0.2
b —— c
0.5

a
1.0 ⌃ 0.2
b —— c
0.6

→ sim to prob

a
b —— c      0.12

a
b —— c      0.48

a
b      c      0.08

a
b      c      0.32

38

19

Gold Standard?

possible worlds

possible clustering (via ER algorithm)

sim to prob

39



Strategy

current state

Q(a,b)

Q(b,c)

Q(a,c)

new state → ER result → score vs GS?

40

20

# Evaluating Efficiently

- <u>See:</u> Steven E. Whang, Peter Lofgren, and H. Garcia-Molina. Question Selection for Crowd Entity Resolution. To appear in Proc. 39th Int'l Conf. on Very Large Data Bases (PVLDB), Trento, Italy, 2013.

41

# Sample Result



: Cora Results using High Threshold

## Summary

data

analytics ⟷ humans

results

Key Point:
- use humans judiciously

43

# Now for something completely different!

analytics

DBMS

big
data

44

22

# Now for something completely different!

analytics

DBMS

big data

humans

45

# DeCo: Declarative CrowdSourcing

what is best price for Nikon DSLR cameras? End user

DBMS

data

humans

46

# DeCo: Declarative CrowdSourcing

what is best price for Nikon DSLR cameras? End user

DBMS

data

humans

| model | type | brand |
|-------|------|-------|
| D7100 | DSLR | Nikon |
| 7D | DSLR | Canon |
| P5000 | comp | Nikon |
| ••• | ••• | ••• |

47

# DeCo: Declarative CrowdSourcing

what is best price for Nikon DSLR cameras? End user

DBMS

data

humans

| model | type | brand |
|-------|------|-------|
| D7100 | DSLR | Nikon |
| 7D | DSLR | Canon |
| P5000 | comp | Nikon |
| ••• | ••• | ••• |

what is best price for Nikon D7100 camera? Crowd

48

# Example with a bit more detail:

User view

| restaurant | rating | cuisine |
|---|---|---|
| Chez Panisse | 4.9 | French |
| Chez Panisse | 4.9 | California |
| Bytes | 3.8 | California |
| ••• | ••• | ••• |

# Example with a bit more detail:

User view

| restaurant | rating | cuisine |
|---|---|---|
| Chez Panisse | 4.9 | French |
| Chez Panisse | 4.9 | California |
| Bytes | 3.8 | California |
| ••• | ••• | ••• |

⋈

| restaurant |
|---|
| Chez Panisse |
| Bytes |
| ••• |

*Anchor*

| restaurant | rating |
|---|---|
| Chez Panisse | 4.8 |
| Chez Panisse | 5.0 |
| Chez Panisse | 4.9 |
| Bytes | 3.6 |
| Bytes | 4.0 |
| ••• | ••• |

*Dependent*

| restaurant | cuisine |
|---|---|
| Chez Panisse | French |
| Chez Panisse | California |
| Bytes | California |
| Bytes | California |
| ••• | ••• |
| ••• | ••• |

*Dependent*

50

25

## Example with a bit more detail:

**User view**

| restaurant | rating | cuisine |
|---|---|---|
| Chez Panisse | 4.9 | French |
| Chez Panisse | 4.9 | California |
| Bytes | 3.8 | California |
| ••• | ••• | ••• |

⋈

| restaurant |
|---|
| Chez Panisse |
| Bytes |
| ••• |

*Anchor*
*fetch rule*

| restaurant | rating |
|---|---|
| Chez Panisse | 4.8 |
| Chez Panisse | 5.0 |
| Chez Panisse | 4.9 |
| Bytes | 3.6 |
| Bytes | 4.0 |
| Bytes | ••• |

*fetch rule*
*Dependent*

| restaurant | cuisine |
|---|---|
| Chez Panisse | French |
| Chez Panisse | California |
| Bytes | California |
| Bytes | California |
| Chez Panisse | ••• |
| ••• | ••• |

*Dependent*
*fetch rule*

51

---

## Example with a bit more detail:

**User view**

| restaurant | rating | cuisine |
|---|---|---|
| Chez Panisse | 4.9 | French |
| Chez Panisse | 4.9 | California |
| Bytes | 3.8 | California |
| ••• | ••• | ••• |

⋈

| restaurant |
|---|
| Chez Panisse |
| Bytes |
| ••• |

*Anchor*
*fetch rule*

| restaurant | rating |
|---|---|
| Chez Panisse | 4.8 |
| Chez Panisse | 5.0 |
| Chez Panisse | 4.9 |
| Bytes | 3.6 |
| Bytes | 4.0 |
| Bytes | ••• |

*fetch rule*
*Dependent*

| restaurant | cuisine |
|---|---|
| Chez Panisse | French |
| Chez Panisse | California |
| Bytes | California |
| Bytes | California |
| Chez Panisse | ••• |
| ••• | French |

*fetch rule*
*Dependent*
*fetch rule*

52

26

# Example with a bit more detail:



User view

| restaurant | rating | cuisine |
|---|---|---|
| Chez Panisse | 4.9 | French |
| Chez Panisse | 4.9 | California |
| Bytes | 3.8 | California |
| ••• | ••• | ••• |

*resolution rule*     *resolution rule*

| restaurant |
|---|
| Chez Panisse |
| Bytes |
| ••• |

*Anchor*

| restaurant | rating |
|---|---|
| Chez Panisse | 4.8 |
| Chez Panisse | 5.0 |
| Chez Panisse | 4.9 |
| Bytes | 3.6 |
| Bytes | 4.0 |
| Bytes | ••• |

*Dependent*

| restaurant | cuisine |
|---|---|
| Chez Panisse | French |
| Chez Panisse | California |
| Bytes | California |
| Bytes | California |
| Chez Panisse | ••• |
| ••• | ••• |

*Dependent*

53

# Example with a bit more detail:



User view

| restaurant | rating | cuisine |
|---|---|---|
| Chez Panisse | 4.9 | French |
| Chez Panisse | 4.9 | California |
| Bytes | 3.8 | California |
| ••• | ••• | ••• |

**1. Fetch**
**2. Resolve**
**3. Join**

| restaurant |
|---|
| Chez Panisse |
| Bytes |
| ••• |

*Anchor*

| restaurant | rating |
|---|---|
| Chez Panisse | 4.8 |
| Chez Panisse | 5.0 |
| Chez Panisse | 4.9 |
| Bytes | 3.6 |
| Bytes | 4.0 |
| ••• | ••• |

*Dependent*

| restaurant | cuisine |
|---|---|
| Chez Panisse | French |
| Chez Panisse | California |
| Bytes | California |
| Bytes | California |
| ••• | ••• |
| ••• | ••• |

*Dependent*

54

27

## Many Query Processing Challenges

SELECT n,l,c
FROM country
WHERE l = 'Spanish'
ATLEAST 8

AtLeast [8]

Join

Filter [l='Spanish']

Resolve[m3]

Join

Scan
D2(n,c)

Fetch
[n⇒l,c]

Resolve[d.e]

Resolve[m3]

Scan
A(n)

Fetch
[l⇒n,c]

Scan
D1(n,l)

Fetch
[n⇒l,c]

55

# Deco Prototype V1.0



56

28

# Conclusion

- Crowdsourcing is important for managing data!
- Still many challenges ahead!

57



58