# Exploiting Structure within Data for Accurate Labeling using Conditional Random Fields

Aman Goel, Craig A. Knoblock, Kristina Lerman

Department of Computer Science & Information Sciences Institute

University of Southern California

# Outline

- Problem
- Existing approaches
- Our approach
- Experiments
- Real word application
- Conclusion
- Q&A (5 mins)
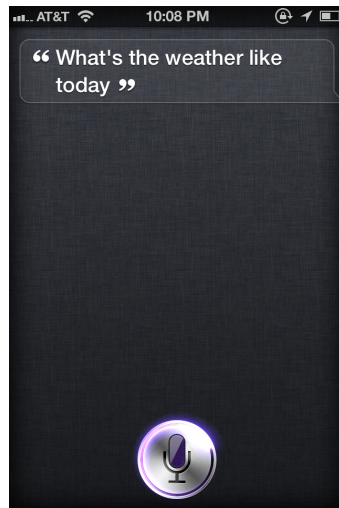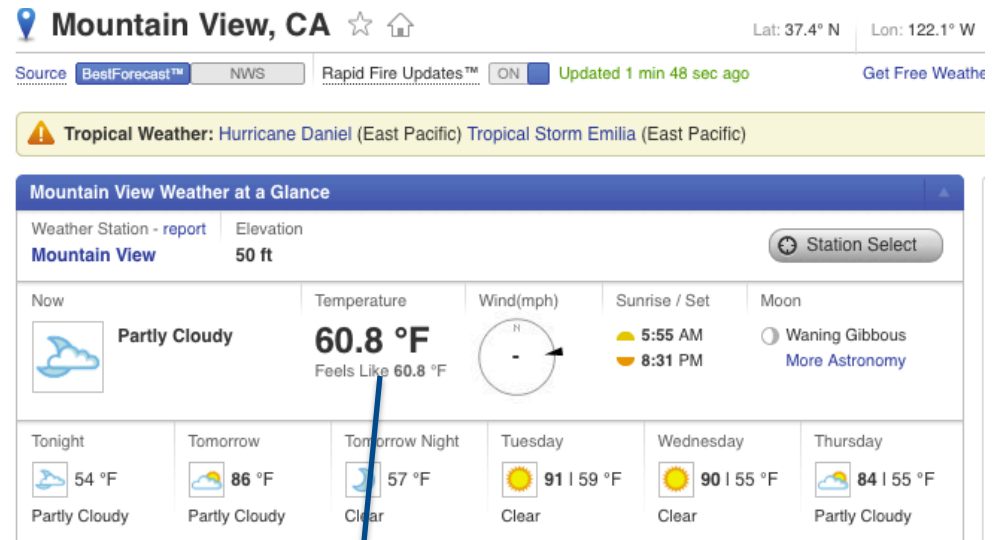
15 mins

# Problem definition

Aman Goel    =>    Name

2667 Ellendale Place, Los Angeles, CA 90007    =>    Address

(323) 246-7180    =>    Telephone Number

# Motivation



Visit a website

Can't use !!!

Extract strings

Mountain View,          Partly Cloudy,          60.8 $^o$F,          5:56 AM

- Vertical approach (schema matching)
- Horizontal approach (hidden markov model)

# Vertical approach (schema matching)

| Name | Address | Telephone |
|------|---------|-----------|
| Aman Goel | 4676 Admiralty Way, Marina del Rey, CA 90292 | (967) 123-9835 |
| ... | ... | ... |

? → 
? → 

Semantic types
1) Contact
2) First name
3) Residence
4) ...

Approach
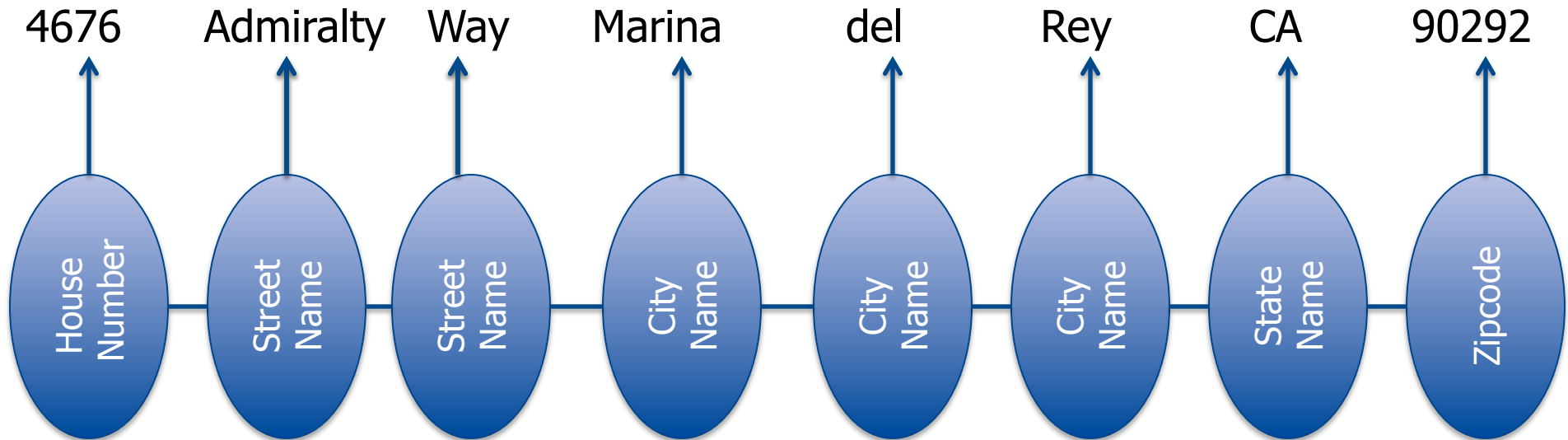• Using dictionaries and string edit distances

Problems
• Assume relational data
• Missing column names

6

# Horizontal approach (HMMs)

Automatic segmentation of text into structured records:
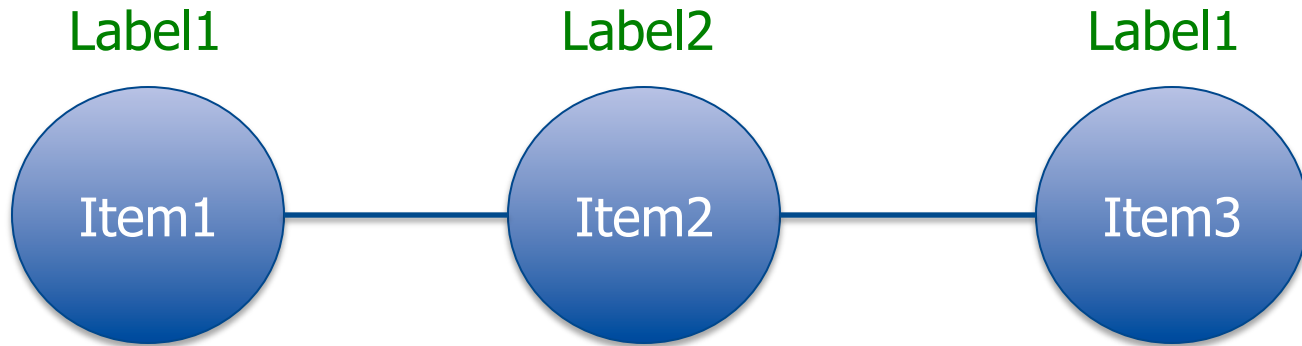
*V. Borkar, K. Deshmukh and S. Sarawagi, SIGMOD, 2001*

| 4676 | Admiralty | Way | Marina | del | Rey | CA | 90292 |
|------|-----------|-----|--------|-----|-----|-----|-------|
| House Number | Street Name | Street Name | City Name | City Name | City Name | State Name | Zipcode |

Problems

- Complex semantic type is already assumed

7

# Our approach

Observations

- Ordering relationship between complex types
  - Name before Address
- Ordering relationship between the tokens
  - House number before Street name
- Relationship between tokens in different fields
  - (mph, inches) vs (kmph, mm)

# CRF

Label1           Label2           Label1

**Item1**           **Item2**           **Item3**

Feature1=True      Feature1=False      Feature1=True

...           ...           ...

FeatureN=False      FeatureN=True      FeatureN=True

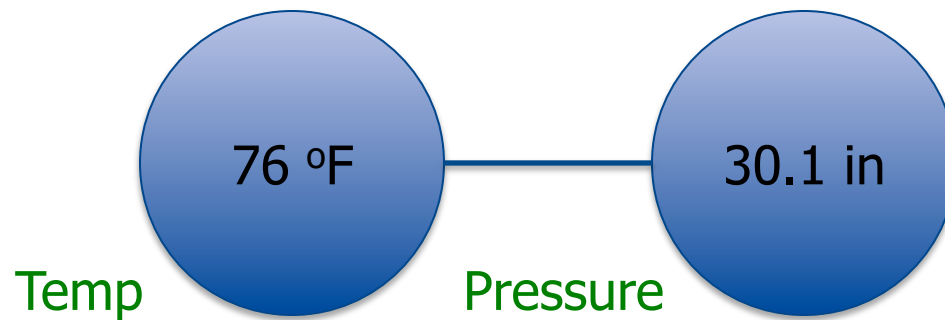<u>Feature functions</u>                           <u>Weights</u>

- Label1 is followed by Label2    - - - - - - - - - -      0.7
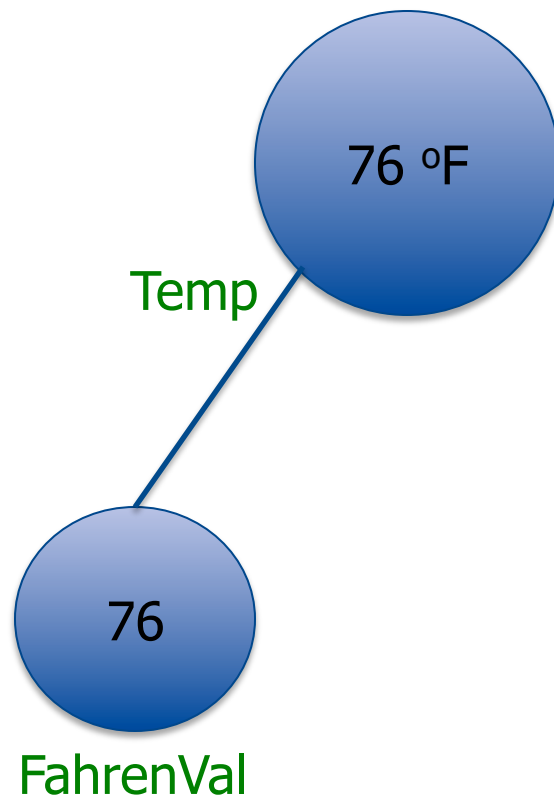- Label1 has Feature1=True      - - - - - - - - - -      1.5

# CRF model from data

# Features

Token features:

- Is_capitalized
- 7_characters_in_token
- Starts_with_B
- Number_is_in_100s
- Has_2_precision_digits
- Is_dollar_sign
- Is_negative
- Number_starts_with_9
- Is_all_caps
- Is_percent_sign
- ...

- Relationship between field nodes



Temp_field_is_followed_by_Pressure_field

# Feature functions
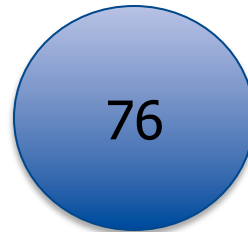
- Relationship between field and token node

76 °F

Temp

Temp_field_has_token_FahrenVal

76

FahrenVal

- Relationship between neighboring token nodes



FahrenVal_token_is_followed_by_DegreeSym_token

# Feature functions

- Relationship between token label and its feature



76

FahrenVal

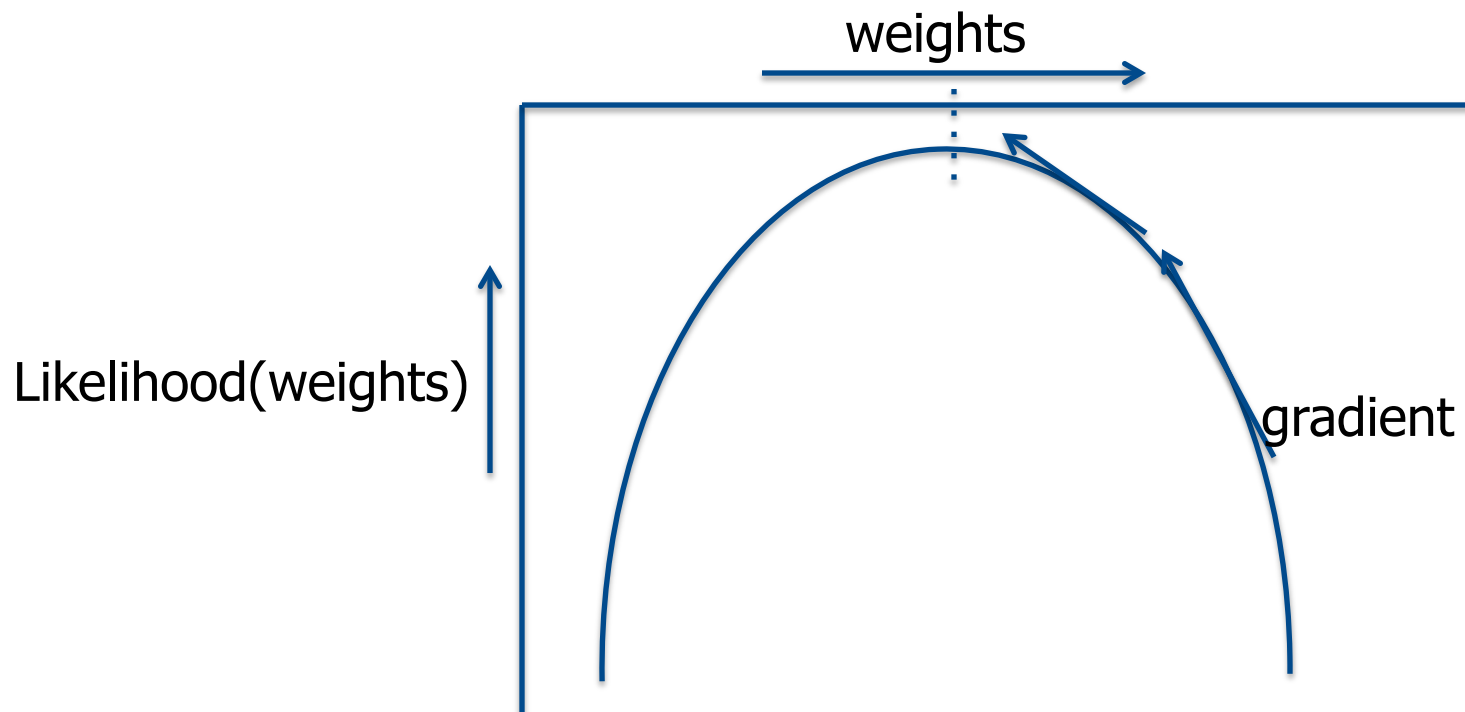FahrenVal_token_is_in_10s

FahrenVal_token_has_0_precision_digits

labels

features

$$1/Z(x) * exp(\sum_c (\sum_k w_k f_k(y_c, x)))$$

weights

feature functions
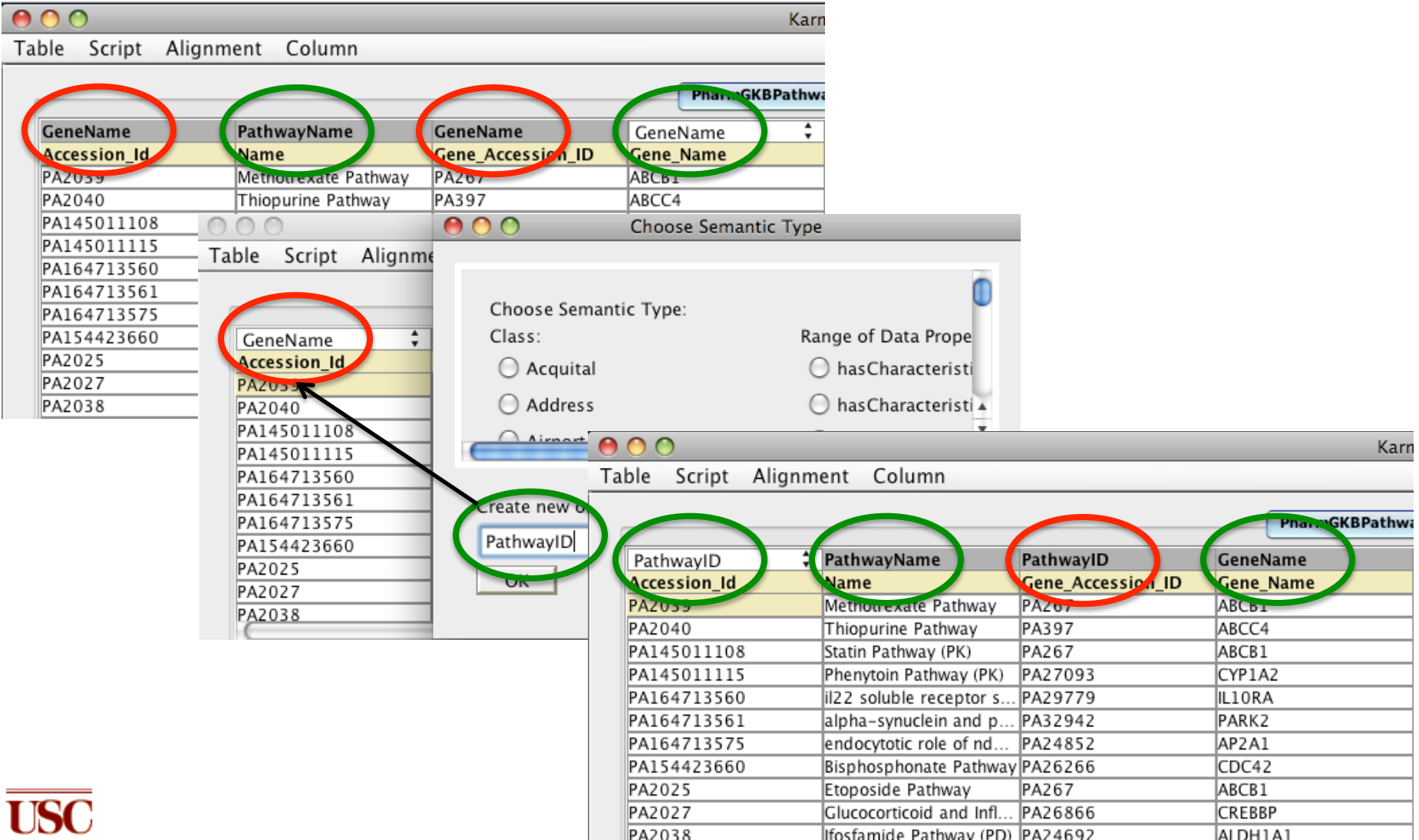
- Convex optimization problem

weights

Likelihood(weights)

gradient

# Experiment setup

- Collected 400 tuples from 4 websites
- Trained on 300 tuples from 3 sites
- Tested on 100 labeled examples from held-out site

|              | Weather | Flight | Geocoding |
|--------------|---------|--------|-----------|
| # field types | 15      | 8      | 5         |
| # token types | 37      | 17     | 12        |

# Results

| | Weather | Flights | Geocoding |
|---|---|---|---|
| Field labeling accuracy | 0.89 | 0.97 | 0.98 |
| Token labeling accuracy | 0.86 | 0.87 | 0.90 |
| Labeling accuracy using regular expressions* | 0.65 | 0.42 | 0.36 |

* Semantic labeling of online information sources:

   K. Lerman, A. Plangrasopchok, and C. A. Knoblock, IJSWIS, 2006

# Karma uses this approach

# Conclusion

Contributions

- Accurately identifying complex semantic types
- Also identify token types
- Fast