# ISE 599
# Special Topics Applied Predictive Analytics

Mayank Kejriwal

Research Assistant Professor/Research Lead

Department of Industrial and Systems Engineering
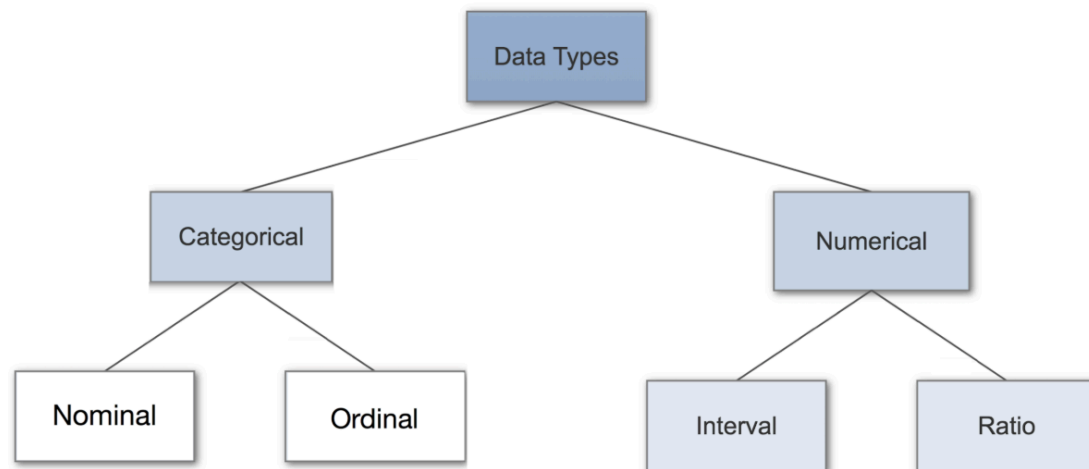
Information Sciences Institute

USC Viterbi School of Engineering

kejriwal@isi.edu

# Data: how we 'make sense' of it (them?)

- First, important to understand the different types of data
  - An 'ontology' of data types
- We've already (kind of) seen one example!

# In practice, such formal definitions are rarely useful (except in labs)

- Article in Forbes shows 13 different types of data
- It's really broad, and some of the categories overlap, but useful as a framework

1 - Big data

2 - Structured, unstructured, semi-structured data

3 - Time-stamped data

4 - Machine data

5 - Spatiotemporal data

6 - Open data

7 - Dark data

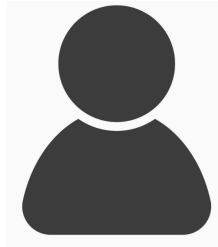- Real time data

9 - Genomics data

10 - Operational data

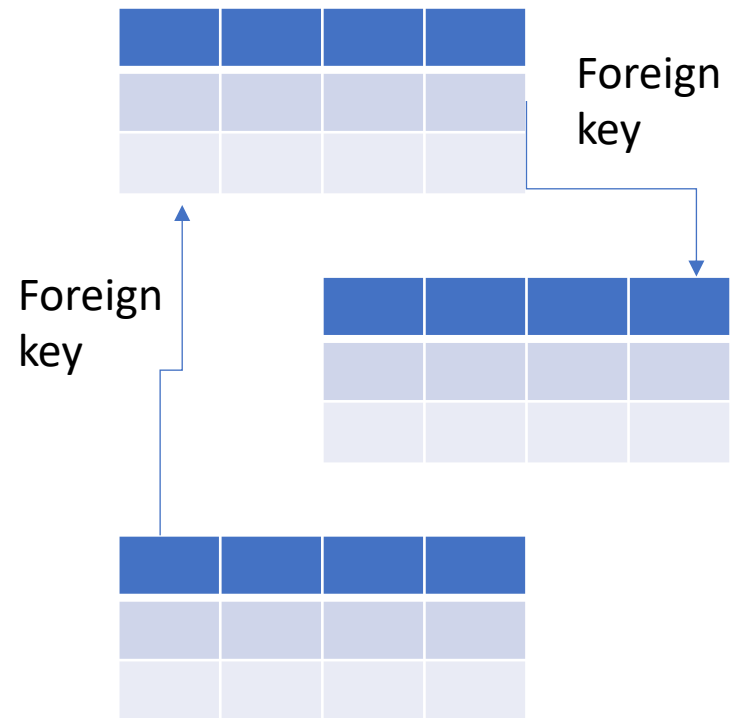11 - High-dimensional data

12 - Unverified outdated data

13 - Translytic Data

# Structured vs. 'unstructured' data

# Two 'extremes'?



Natural language, social media data

Foreign key

Foreign key

# Is data ever *really* unstructured?

- Many computer scientists would call English 'unstructured'
  - You can substitute English for any 'natural' language spoken by humans in society
- The great philosopher Gottlieb Frege, like so many others, felt English was woefully imprecise
- Thought of logic as one way to address these difficulties
- Never panned out in the AI community, too many irregularities in English and other languages!

# Structure (beauty) is in the eye of the application (beholder)

- Unfortunately, 'natural language' data is often called unstructured data by many practitioners
    - I encourage the phrase 'natural language' vs. unstructured, since it has an impact on how you think about the data

Demo: spacy
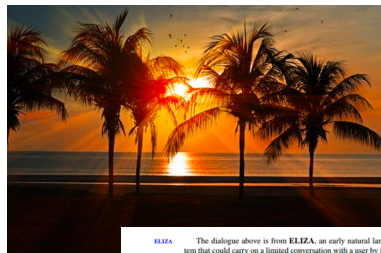
# https://demos.explosion.ai/displacy-ent

Kejriwal, Szel

# In practice, IE is usually more complex

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Blomkvist visits Henrik Vanger at *same* te on the *same* and of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of V *owns* niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist beco *uncleOf* inted with the men *hires* the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

Af *same* overing that Salander has hacked into his computer, he persuade *same* assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou *headOf* gations for Dragan Armansky, who, in tu *same* ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

# In practice, IE is usually more complex

It's about the disappearance forty years ago of **Harriet Vanger**, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

**Blomkvist** visits **Henrik Vanger** at ... te on the ... and of Hedeby. The old man ... Blomkvist in by promising solid evidence against **Wennerström**. Blomkvist ag... spend a year writing the **Vanger family** history as a cover for the real assignment: the disappearance of V **owns** niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in **Vanger Enterprises**. Blomkvist beco **uncleOf** ... ted with the men **hires** the extended Vanger family, most of whom resent his pres...

**Why should this make it 'easier' for machines?**

A ... **same** ... ing that **Salander** has hacked into his computer, he persuade **same** assist him with research. They eventually become lovers, but **Blomkvist** has trouble getting close to **Lisbeth** who treats virtually everyone she meets with hostility. Ultimately the two discover that **Harriet's brother Martin,** CEO of **Vanger Industries** is secretly a serial killer.

A **24-year-old computer hacker** sporting an ... ment of tattoos and body piercings su... herself by doing deep backgrou **headOf** gations for Dragan Armansky, who, in tu **same** ... ies that **Lisbeth Salander** is "the perfect victim for anyone who wished her ill."

# Machines like certain kinds of structure and technologies like IE can 'parse' that structure from human-centric structure

**Humans like…**



**Machines like…**

XML

<empinfo>
  <employees>
    <employee>
      <name>James Kirk</name>
      <age>40></age>
    </employee>
    <employee>
      <name>Jean-Luc Picard</name>
      <age>45</age>
    </employee>
    <employee>
      <name>Wesley Crusher</name>
      <age>27</age>
    </employee>
  </employees>
</empinfo>

JSON

{ "empinfo" :
  {
    "employees" : [
      {
        "name" : "James Kirk",
        "age" : 40,
      },
      {
        "name" : "Jean-Luc Picard",
        "age" : 45,
      },
      {
        "name" : "Wesley Crusher",
        "age" : 27,
      }
    ]
  }
}

# Many other kinds of 'structure' out there…

**Text paragraphs without formatting**

Astro Teller is the CEO and co-founder of BodyMedia. Astro holds a Ph.D. in Artificial Intelligence from Carnegie Mellon University, where he was inducted as a national Hertz fellow. His M.S. in symbolic and heuristic computation and B.S. in computer science are from Stanford University. His work in science, literature and business has appeared in international media from the New York Times to CNN to NPR.

**Grammatical sentences plus some formatting & links**



**Non-grammatical snippets, rich formatting & links**



**Tables**



**Charts**

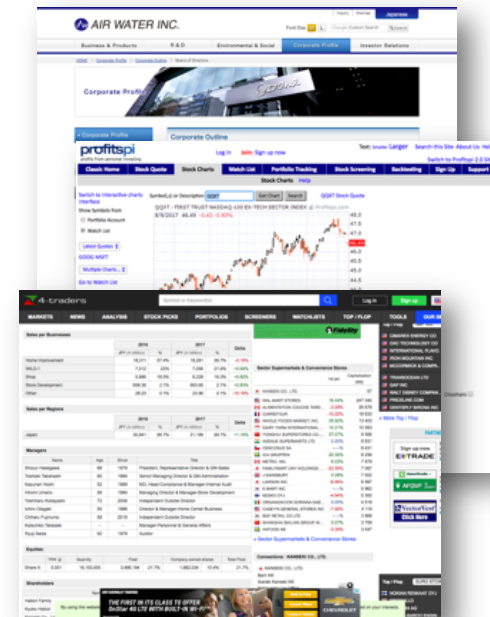# Structure (and IEs) can depend on **scope**

**Web site specific**

**Genre specific
(e.g., forums)**

**Wide, non-specific**

# Structure (and IEs) can depend on **pattern complexity**

### Closed set

U.S. states

> He was born in <u>Alabama</u>...

> The big <u>Wyoming</u> sky...

### Regular set

U.S. phone numbers

> Phone: <u>(413) 545-1323</u>

> The CALD main office can be reached at <u>412-268-1299</u>

### Complex pattern

U.S. postal addresses

> University of Arkansas
> <u>P.O. Box 140</u>
> <u>Hope, AR  71802</u>

> Headquarters:
> <u>1128 Main Street, 4th Floor</u>
> <u>Cincinnati, Ohio 45210</u>

### Ambiguous patterns, needing context and many sources of evidence

Person names

> ...was among the six houses sold by <u>Hope Feldman</u> that year.

> <u>Pawel Opalinski</u>, Software Engineer at WhizBang Labs.