# Information Extraction.

① Obtain structured text/ from unstructured data.

**Objectives:** IE system.

① Find & understand [limited] relevant parts of text.

ⅱ Gather info from many different texts.

ⅲ Produce a structured representation of the relevant info gathered.

Relations (database)
& Knowledge base / Knowledge graph.

**2nd goal:**

- Make the data "useful" for users
- Organize the info in a semantically precise manner. $\xrightarrow[\text{by}]{\text{used}}$ computer algorithms.

Input (unstructured/semi-structured data) $\longrightarrow$ structured data.

Information: ① exact
             ② factual

Roughly :      IE system

[Who] did [What] to [whom] and [when]?

The headquarters of BHP Billion Limited, and the
global headquarters of the combined BHP
Billion Group are located in Melbourne
Australia.

[headquarters] ( "BHP Billion Limited", "Melbourne,
                                                Australia)

relation        entity1                        entity2.

In 1998 Larry Page and Sergey Brin founded Google Inc.

FounderOf (Larry Page, Google Inc.)
FounderOf (Sergey Brin, Google Inc.)
Founded In (Google Inc., 1998).

Applications?

Biomedical domain.

① Large no. of scientific publications from the biomed domain

② Discover easily innovations (i) particular genes (ii) proteins & other entities.

③ Biomed entities ← ① many ambiguous names.

② lot of ~~synonymous~~ synonymous words.

Automatically identify the bio medical
entities and __link.__ them in
knowledge bases.
$\hookrightarrow$ WikiData.

| subject | relation | object |
|---|---|---|
| p53 | is-a | protein |
| p53 | has function | apptosis |
| apoptosis | involved-in | cell death |
| ⋮ | ⋮ | ⋮ |

# News articles.

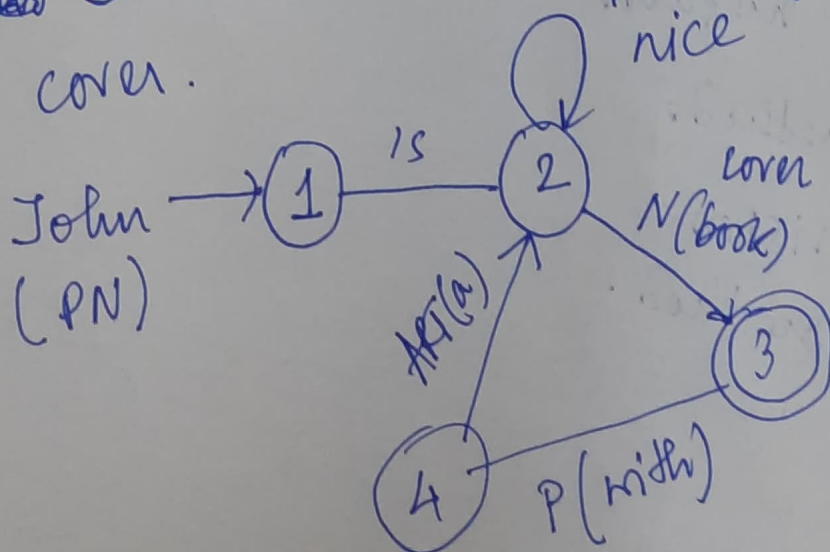| Relations. | | Examples | Types. |
|---|---|---|---|
| Affiliations | → Personal | → married to, mother of | PER → PER |
| | → Organisational | → president of, spokesman for | PER → ORG |
| | → Artifactual. | → owns, invented. | (PER\|ORG) → ART |
| Geospatial | → Proximity | → near, on outskirts | LOC → LOC |
| | → Directional | → south of | LOC → LOC |
| Part - of | → Organisational | → a unit of, parent of | ORG → ORG |
| | → Political. | → ~~annexed~~ annexed acquired | GPE → GPE |

geopolitical entity.

- Hand written rules ( REGEX)
- Boot strapping methods.
- Supervised methods
- Distant supervision
- Unsupervised method.

# RELEX based IE.

Use finite automaton for noun groups:



(proper noun)
PN

ART (article)

N (noun)

IS

P

PN

IS

ADJ (adjective).

ART

→ John's interesting book with a nice cover.

interesting } ADJ
nice



John → (1) —IS→ (2) —N(book)→ (3)
ART(a)
(4) P(with)

John (PN)

cover

many such RELEX to extract info.

→ Determine which person holds what position in what organization.

[person], [position] of [org].←

Vuk Draskovic, (leader of) the }
Serbian Renewal Movement }

[org] (named, appointed, etc.) [person] PREP [office].

(NATO) appointed Wesley Clark (as) (Commander in Chief)

org        person     Preposition.   office

→ Determine where an organization is located

[org] in [loc].
NATO headquarters in Brussels.
[org], [loc] (division, branch, headquarters, etc.).
KFOR, Kosovo headquarters -

(Agar) is a substance prepared from a mixture of red algae [Such as] (Gelidium) for lab or industry use.