# Language Understanding

14 - Information Extraction and Knowledge Graphs

Hossein Zeinali

# What is Knowledge Graph?

- *"A huge knowledge graph of interconnected entities and their attributes".* Amit Singhal, Senior Vice President at Google

- *"A knowledge based used by Google to enhance its search engine's results with semantic-search information gathered from a wide variety of sources"* http://en.wikipedia.org/wiki/Knowledge_Graph

- Based on information derived from many sources including *Freebase, CIA World Factbook, Wikipedia*

- By May 2020, contains 5 billion entities and more than 500 billion facts about and relationships between these different objects

# Who Constructs Knowledge Graphs?

- Google

- Amazon

- Microsoft

- Bloomberg (business intelligence)

- National Cancer Institute (part of NIH)

- National Center for Biomedical Ontology (funded by NIH)

- Specialized knowledge graphs:
  - Health / life sciences (health-lifesci.schema.org)
  - Earth science
  - Agriculture
  - NCI Enterprise Vocabulary Services (evs.nci.nih.gov)

# Google Knowledge Graph (GKG)

GKG enhances Google Search in three main ways:

- **Find the right thing:** Deals with the ambiguity of the language

# Google Knowledge Graph (GKG)

GKG enhances Google Search in three main ways:

- **Summaries:** summarize relevant content around that topic, including key facts about the entity

# Google Knowledge Graph (GKG)

GKG enhances Google Search in three main ways:

- **Deeper and broader information:** reveal new facts, anticipate what the next questions and provide the information beforehand

# How GKG is used?

- Explore your search

AUT, Language Understanding Course, Fall 2022, Hossein Zeinali

# Sampling of GKG Entities

- Book
- BookSeries
- EducationalOrganization
- Event
- GovernmentOrganization
- LocalBusiness
- Movie
- MovieSeries
- MusicAlbum
- MusicGroup
- MusicRecording

- Organization
- Periodical
- Person
- Place
- SportsTeam
- TVEpisode
- TVSeries
- VideoGame
- VideoGameSeries
- WebSite

# "Person" Schema

| Property | Expected Type | Description |
|---|---|---|
| **Properties from Person** | | |
| `additionalName` | Text | An additional name for a Person, can be used for a middle name. |
| `address` | PostalAddress or Text | Physical address of the item. |
| `affiliation` | Organization | An organization that this person is affiliated with. For example, a school/university, a club, or a team. |
| `alumniOf` | EducationalOrganization or Organization | An organization that the person is an alumni of. Inverse property: alumni. |
| `award` | Text | An award won by or for this item. Supersedes awards. |
| `birthDate` | Date | Date of birth. |
| `birthPlace` | Place | The place where the person was born. |
| `brand` | Brand or Organization | The brand(s) associated with a product or service, or the brand(s) maintained by an organization or business person. |
| `children` | Person | A child of the person. |
| `colleague` | Person or URL | A colleague of the person. Supersedes colleagues. |
| `contactPoint` | ContactPoint | A contact point for a person or organization. Supersedes contactPoints. |
| `deathDate` | Date | Date of death. |
| `deathPlace` | Place | The place where the person died. |
| `duns` | Text | The Dun & Bradstreet DUNS number for identifying an organization or business person. |
| `email` | Text | Email address. |
| `familyName` | Text | Family name. In the U.S., the last name of an Person. This can be used along with givenName instead of the name property. |
| `faxNumber` | Text | The fax number. |
| `follows` | Person | The most generic uni-directional social relation. |

# Data Sources Used for GKG

# CIA World Factbook

- CIA World Factbook is a reference resource produced by the Central Intelligence Agency of the United States with almanac-style information about the countries of the world.

- GKG integrates information about geography, government, economy, etc. from CIA World Factbook.



Austria

Country

Austria, officially the Republic of Austria, is a federal republic and a landlocked country of roughly 8.47 million people in Central Europe. Wikipedia

**Capital:** Vienna
**Currency:** Euro
**Chancellor:** Werner Faymann
**National anthem:** Land der Berge, Land am Strome
**Official language:** German Language
**Government:** Federal republic, Parliamentary republic, Federation

Points of interest

Zugspitze   Schönbrunn Palace   Vienna State Opera   Prater   Kitzsteinh...

# Freebase

- Freebase is large collaborative knowledge base, developed by Metaweb and acquired by Google in 2010.

- GKG uses UIDs directly from the Freebase; detective work of Andreas Thalhammer showing how to get from GKG UIDs to Freebased UIDs using base64 and gzip

- Check the "Knowledge Graph links to Freebase" thread on w3c mailing list
http://lists.w3.org/Archives/Public/semantic-web/2012Jun/0028.html

# Wikipedia

- For most search results first sentences come from **Wikipedia**



**WIKIPEDIA**
*The Free Encyclopedia*



Dieter Fensel
Innsbruck
○ Friends

## Dieter Fensel

Dieter Fensel is a researcher in the field of formal languages and the semantic web. He is University Professor at the University of Innsbruck, where he directs the Semantic Technologies Institute ...
Wikipedia

**Born:** October 10, 1960 (age 52), Nuremberg

**Books:** The knowledge acquisition and representation language, KARL

People also search for

Rudi Studer | Frank van Harmelen | James Hendler | Ian Horrocks | Deborah McGuinne...

# Other Sources

- GKG also considers the information Google retrieves from the volume of queries done by the users and the links those users have clicked on the results presented for those queries.

- GKG is also integrated with other Google products e.g. Google+

# Web of Data



Web of Data

Semantic Web

Semantic Web

Semantic Annotations

Web

Picture from [4]

?

Hypermedia

Hypertext

"As We May Think", 1945

Picture from http://www.theatlantic.com/doc/194507/bush

# Web of Data

- Characteristics:
  - Links between arbitrary things (e.g., persons, locations, events, buildings)
  - Structure of data on Web pages is made explicit
  - Things described on Web pages are named and get URIs
  - Links between things are made explicit and are typed

- HTML defines how text should look when presented to humans. Semantic web markup defines how information should be organized to be interpretable by machines.

- Web of Data

Typed Links

"Things"

# Web of Data

- A **closed** implementation of *Web of Data* principles
  - is not about documents, but objects such as people, places and things
  - objects are interlinked in the GKG
  - objects have structured information which is obtained from the web

- The Google Knowledge Graph is the basis for transforming Google' core search product from an *information engine* to a *knowledge engine (entity search engine)*

# Building a KG Through ML

# NELL: Never-Ending Language Learning

- Tom Mitchell et al. (CMU), 2010 to present.

- Learning to "read the web" 24 hours/day.

- Training data includes a collection of 1.2 billion web pages.

- Access to additional data through search engine APIs (100K calls/day).

- KB has 2.8 million instances over 1186 different categories.

- KB is freely available for download.

- You can help train NELL via Twitter.

Source: Never-Ending Learning (2018)

# NELL Knowledge Fragment

# Motivation for NELL

Thesis: "we will never truly understand human or machine learning until we can build computer programs that, like people,

- Learn many different types of knowledge or functions

- From years of diverse, mostly self-supervised experience

- In a staged curricular fashion, where previously learned knowledge enables learning further types of knowledge

- Where self-reflection and the ability to formulate new representations and new learning tasks enable the learning to avoid stagnation and performance plateaus."

# NELL Basic Idea

- NELL learns several things:
  - Categories
  - Triples: noun phrase 1 - relation - noun phrase 2
  - New relations

- Multiple inference algorithms propose triples and gather evidence for them.
  - Linguistic information
  - Word co-occurrence
  - Image labeling
  - Etc.

- Categories and triples supported by multiple sources of evidence grow in confidence.

# NELL Architecture

# Never-Ending Learning

- Set of learning tasks $L = \{L_i\}$
- Task $L_i = <T_i, P_i, E_i>$
    - $T_i$ is a task $<X_i, Y_i>$ specifying the domain of a function $f_i^* : X_i \rightarrow Y_i$
    - $Pi$ is a performance metric $P_i : f \rightarrow \mathbb{R}$
    - $E_i$ is an experience
- Coupling constraints $C = \{<\phi_k, V_k>\}$
    - $\phi_k$ specifies degree of satisfaction of the coupling constraint among tasks
    - $V_k$ is a vector of indices over learning tasks specifying the arguments to $\phi_k$
- $f_i^* = \min\limits_{f \in F_i} P_i(f)$

✓Goal: improve the quality of the task functions $f_i$ as measured by the $P_i$.

✓NELL faces over 4100 distinct learning tasks.

# Category Classification Tasks

1. Character string features of the noun phrase: Coupled Morphological Classifier system (CMC)

2. Distribution of text contexts found around this noun phrase in the 1.2 billion page database: Coupled Pattern Learner system (CPL)

3. Distribution of text contexts found through active web search (OpenEval).

4. HTML structure of web pages that mention the noun phrase: Set Expander for Any Language system (SEAL)

5. Visual images associated with the noun phrase: Never Ending Image Learner (NEIL)

6. Learned vector embeddings (feature vectors) of the noun phrase: LE (Learned Embeddings)

# Relation Classification

Does "Pittsburgh" + "US" satisfy the relation CityLocatedInCountry(x,y) ?

There are 461 relations in the ontology.

Four methods are used for relation classification:

1. Distribution of text contexts from CPL
2. Distribution of text context from OpenEval
3. HTML structure from SEAL
4. Learned vector embeddings from LE

# Entity Resolution

- Functions to classify whether pairs of noun phrases are synonyms.

- Noun phrases are kept distinct from the entities to which they refer.

- Necessary to deal with polysemy.
  - "Coach" can be either a person or a vehicle.

- Two methods are used:
  - String similarity
  - Similarities in beliefs about the entities

- NELL learns for each category what are the good types of knowledge to take as evidence for synonymy.

# Inference Rules Among Belief Triples

- Functions that propose new beliefs to be added to the KB.

- For each relation, the corresponding function is represented by a collection of restricted Horn Clause rules learned by the Path Ranking Algorithm (PRA)
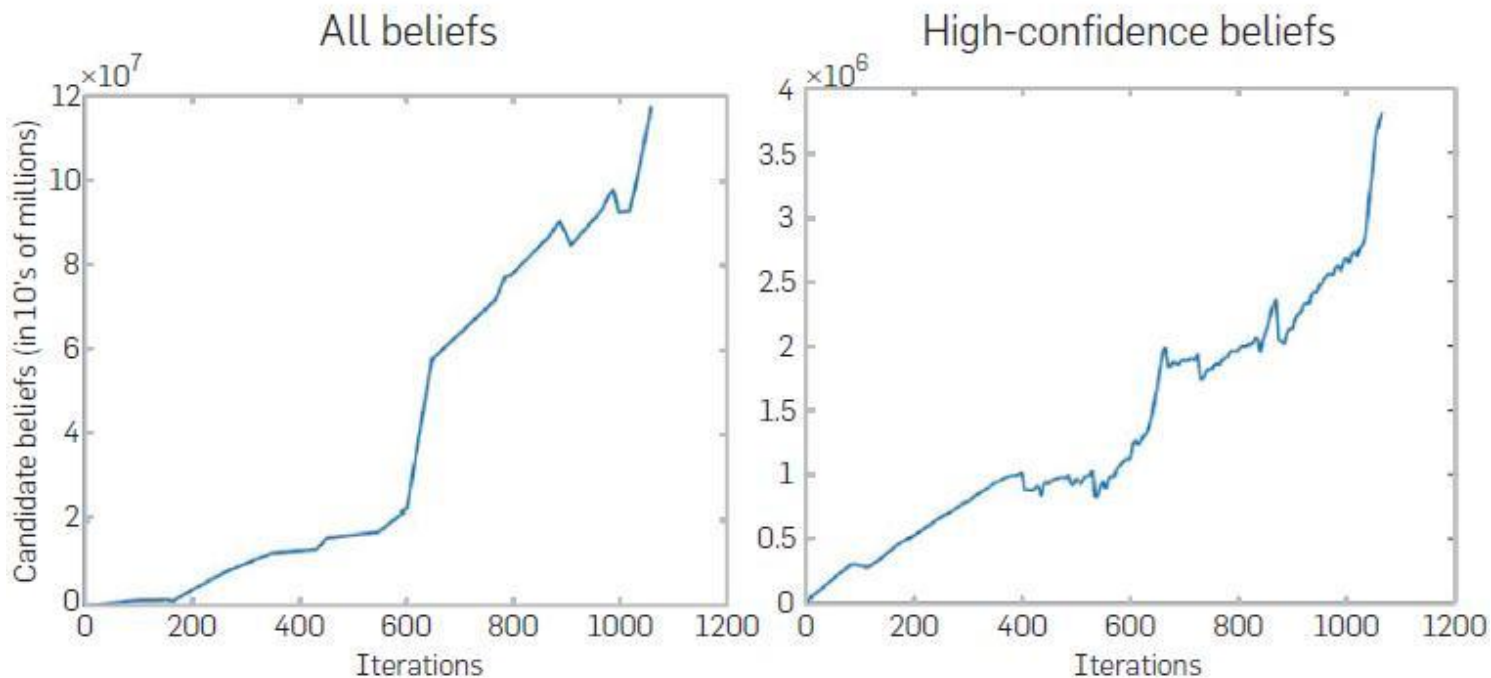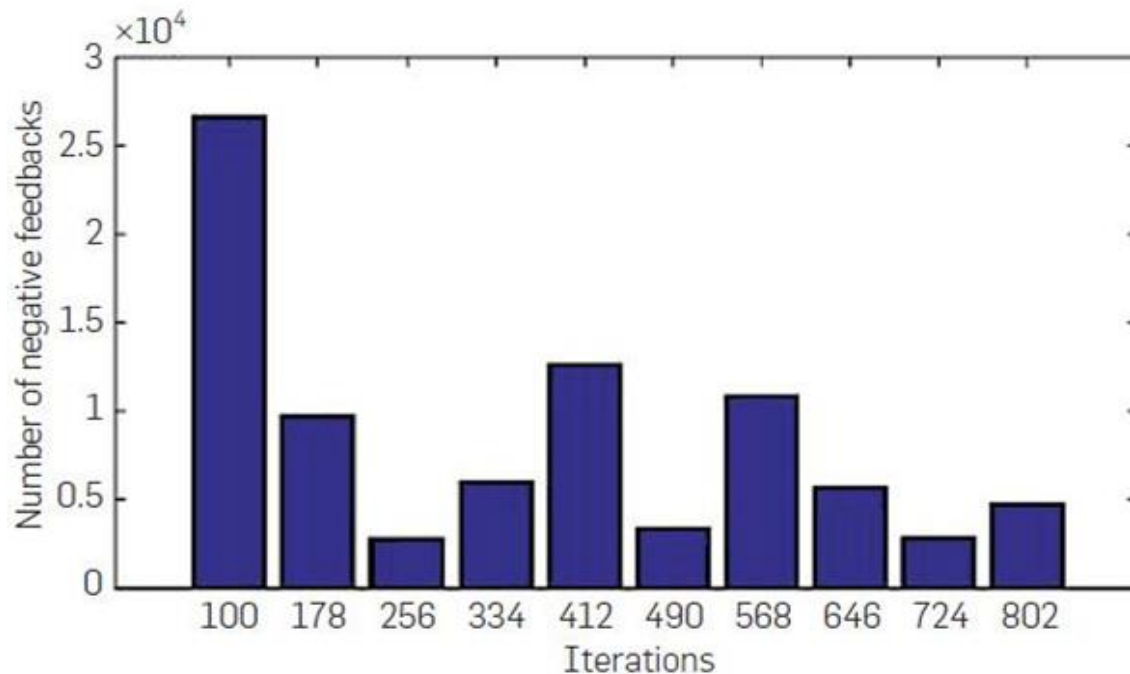
# Samples of Self-discovered Relations

- athleteWonAward
- animalEatsFood
- languageTaughtInCity
- clothingMadeFromPlant
- beverageServedWithFood
- fishServedWithFood
- athleteBeatAthlete
- athleteInjuredBodyPart
- arthropodFeedsOnInsect
- animalEatsVegetable
- plantRepresentsEmotion
- foodDecreasesRiskOfDisease

- clothingGoesWithClothing
- bacteriaCausesPhysCondition
- buildingFeatureMadeFromMaterial
- emotionAssociatedWithDisease
- foodCanCauseDisease
- agriculturalProductAttractsInsect
- arteryArisesFromArtery
- countryHasSportsFans
- bakedGoodServedWithBeverage
- beverageContainsProtein
- animalCanDevelopDisease
- beverageMadeFromBeverage

# Growth of the KB Over Time

# Human Correction of NELL



Average 2.4 negative feedback items per month per predicate.

# Thanks for your attention

# References and IP Notice

- Some slides from Dave Touretzky's slides on knowledge graph.

- The begging slides are mostly based on Ioan Toma's slides.

- Some graphics from Slidesgo online template.