

Probabilistic Knowledge Graphs

Sargur N. Srihari
srihari@buffalo.edu

Topics

- Knowledge Graphs (KGs)
- Statistical Relation Learning (SRL) for KGs
- Latent Feature Models
 - RESCAL, ER-MLP, Latent distance
 - Training SRL
- Markov Random Fields from KGs
- References

Google Knowledge Graph (KG)

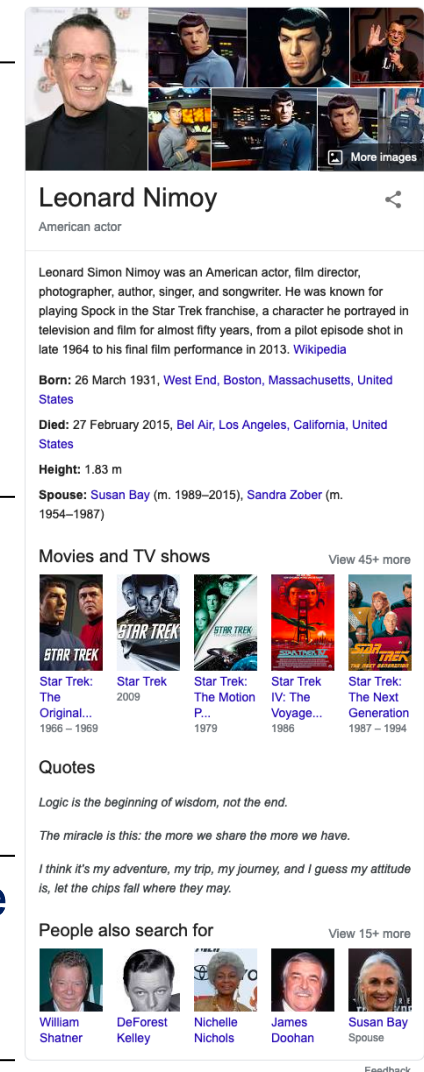
Knowledge panels

information boxes that appear when you search for entities (people, places, organizations, things) that are in the [Knowledge Graph](#)

They are meant to help get a quick snapshot of information on a topic based on available content on the web.

Google's KG, stores 18 billion facts about 570 million entities, For its search engine

KGs provide semantically structured information interpretable by computers
an important ingredient to build more intelligent machines



The screenshot shows a Google Knowledge Panel for Leonard Nimoy. At the top is a collage of images including a portrait of Nimoy, his character Spock, and various Star Trek scenes. Below the images is the name "Leonard Nimoy" and the subtitle "American actor". A brief biography follows, stating he was an actor, film director, photographer, author, singer, and songwriter, known for playing Spock. It mentions his tenure from 1964 to 2013. Key facts listed include his birth (26 March 1931, West End, Boston, Massachusetts, United States) and death (27 February 2015, Bel Air, Los Angeles, California, United States). His height is noted as 1.83 m, and his spouse is listed as Susan Bay (m. 1989–2015), with Sandra Zober (m. 1954–1987) also mentioned. A section titled "Movies and TV shows" displays a row of Star Trek titles: "Star Trek: The Original Series" (1966–1969), "Star Trek: The Motion Picture" (1979), "Star Trek: IV: The Voyage..." (1986), and "Star Trek: The Next Generation" (1987–1994). Below this is a "Quotes" section with three quotes from Nimoy. At the bottom, a "People also search for" section lists related names: William Shatner, DeForest Kelley, Nichelle Nichols, James Doohan, and Susan Bay (Spouse). A "Feedback" link is visible at the very bottom right.

Knowledge Graphs use Relations

- Binary Relation
 - Set of ordered pairs of objects
 - Pairs in set have the relation while those not in set do not
- Relation in the context of AI
 - The relation plays the role of a verb, while two arguments to the relation play the role of its subject and object
 - These sentences take the form of a triplet of tokens (subject, verb, object) or (subject, predicate, object)
 - with values ($\text{entity}_i, \text{relation}_j, \text{entity}_k$)

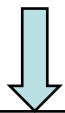
Knowledge Graph Example

- Fact:** Leonard Nimoy was an actor who played the character Spock in the science-fiction movie Star Trek

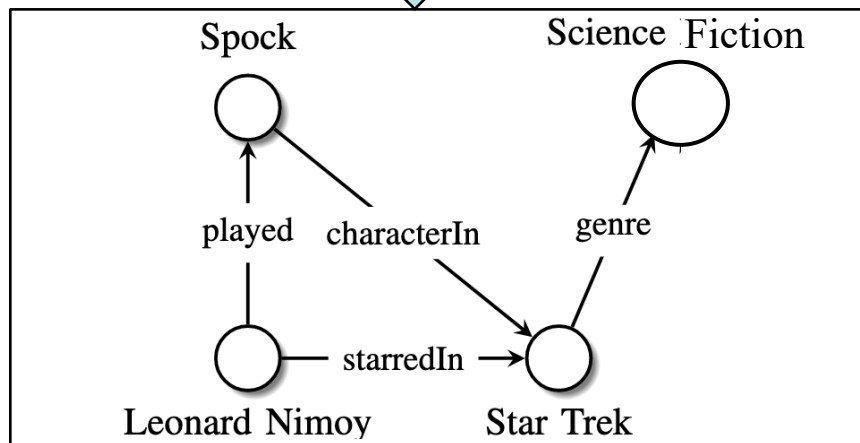


SPO triples

Subject	Predicate	Object
(LeonardNimoy,	profession,	Actor)
(LeonardNimoy,	starredIn,	StarTrek)
(LeonardNimoy,	played,	Spock)
(Spock,	characterIn,	StarTrek)
(StarTrek,	genre,	ScienceFiction)



Knowledge Graph



Subject and Object are *entities* and Predicate is the relation between them

[RDF standard: represent facts in the form of binary relationships]

Combine SPO triples to form a multi-graph:

Nodes represent entities

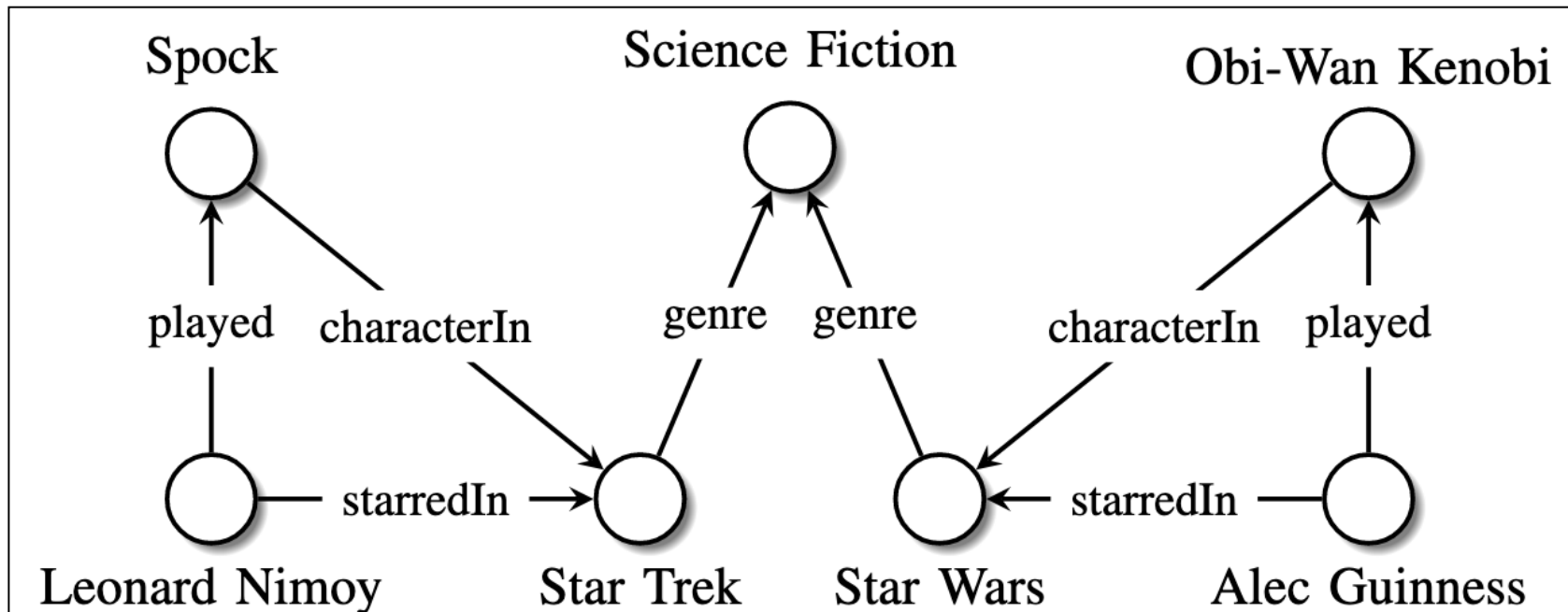
(all subjects and objects)

Directed edges represent relationships.

Edge direction indicates whether entities occur as subjects or objects, (edge points from subject to object)

Different relations are represented via different types of edges (also called edge labels)

Expanded Knowledge Graph



Nodes represent entities

Edge labels represent types of relations

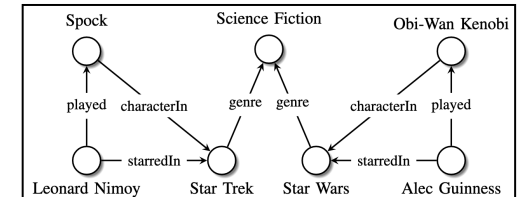
Edges represent existing relationships

In addition to being a collection of facts, knowledge graphs provide type hierarchies (Leonard Nimoy is an actor, which is a person, which is a living thing)

and type constraints (e.g., a person can only marry another person, not a thing)

Open vs Closed-world assumption

- In *closed world* assumption (CWA)
 - a non-existing triple is a false relationship
 - No starredIn edge from Leonard Nimoy to Star Wars means that Nimoy did not star in this movie
- In *open world* assumption (OWA)
 - a non-existing triple is interpreted as unknown
 - Corresponding relationship is either true or false
 - Missing edge does not mean that Nimoy did not star in Star Wars
- RDF and the Semantic Web make the OWA
- CWA often used for training relational models



Definition of Attribute

- We can also define an attribute, a concept analogous to a relation, but taking only one argument:

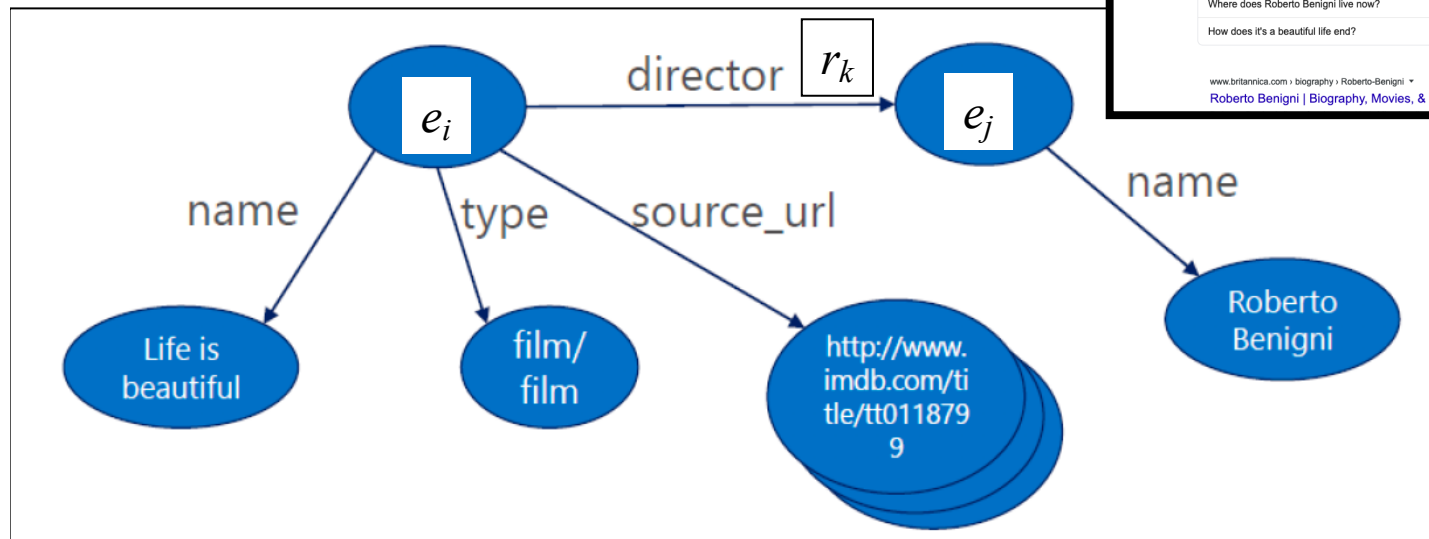
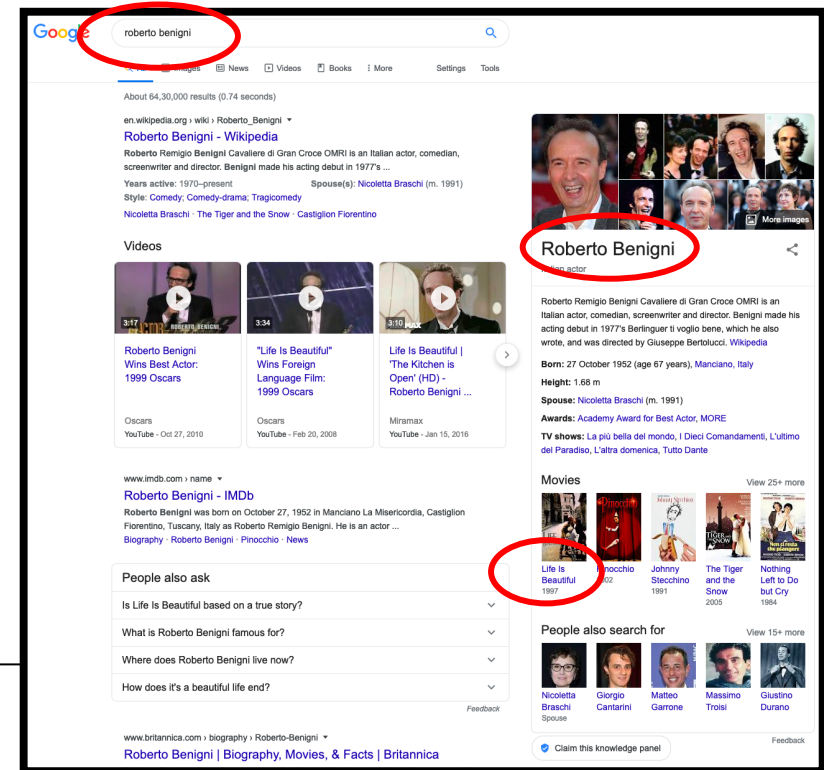
$(\text{entity}_i, \text{attribute}_j)$

- For example, we could define the `has_fur` attribute, and apply it to entities like `dog`

Knowledge Graph with attributes

Part of a semantic KG
representing relationships
between the entities

Roberto Benigni (e_1) and
the film Life is Beautiful (e_2)



Training Datasets

- Inferring relations between entities from data
- From Knowledge Bases
 1. Databases conveying commonsense knowledge about everyday life
 - Freebase, OpenCyc, WordNet, or Wikibase, etc
 2. Databases with expert knowledge about an application AI system
 - GeneOntology
- From Unstructured natural language

KG Construction Approaches

- **Curated:**
 - Triples created by experts
 - Cyc, WordNet, UMLS
- **Collaborative:**
 - Triples created by volunteers
 - Wikidata, Freebase
- **Automated semi-structured**
 - Triples from text via rules:
 - YAGO, DBPedia, Freebase
- **Automated unstructured**
 - Triples text via ML, NLP
 - Knowledge Vault, Nell

Knowledge Graph Sizes

Knowledge Graph	Number of		
	Entities	Relation Types	Facts
Freebase ³	40 M	35,000	637 M
Wikidata ⁴	18 M	1,632	66 M
DBpedia (en) ⁵	4.6 M	1,367	538 M
YAGO2 ⁶	9.8 M	114	447 M
Google Knowledge Graph ⁷	570 M	35,000	18,000 M

Neural Models of Relations/Entities

- In addition to training data, we also need to define a model family to train
- Extend neural language models to model entities and relations
- Neural language models learn:
 1. A vector that provides a distributed representation of each word
 2. About interactions between words
 - such as which word is likely to come after a sequence of words, by learning functions of these vectors
- Extend approach to entities and relations by learning an embedding vector for each relation

Applications of neural knowledge models

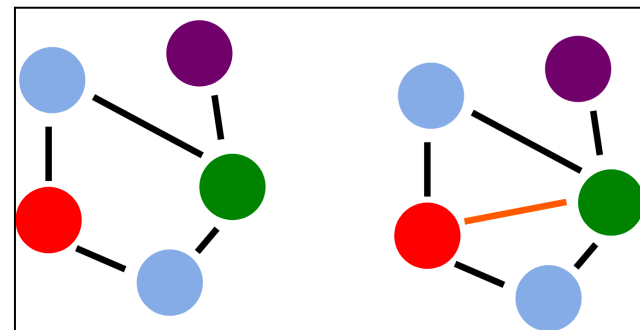
- Link prediction:
 - Predicting missing arcs in the knowledge graph
- Entity Resolution
 - Is A. Guinness the same as Alec Guinness?
- Link-based Clustering
 - Extend feature-based clustering to relation based on similarity of entities and relations
- Word sense disambiguation
 - Deciding which of the senses of a word is the appropriate one, in some context

Link Prediction

- Predicting existence (probability of correctness) of (typed) edges in graph (i.e., triples)
 - KGs often miss facts, and some edges are incorrect

 $G=(V,E)$ $e=(x,y)$ Similarity score s_{xy}

Local Similarity measure:
Common Neighbors (CN)

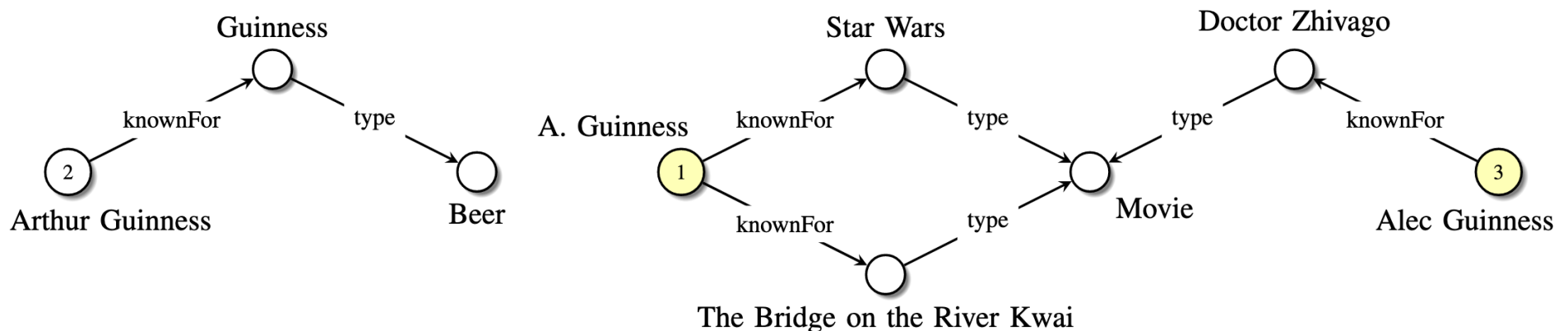


$$s_{xy}^{\text{CN}} = |\Gamma(x) \cap \Gamma(y)|$$

Set of NN. Cardinality is the degree of the node

Entity Resolution

- Problem of identifying which objects in relational data refer to the same underlying entities



Nodes 1 and 3 refer to the identical entity, the actor Alec Guinness.

Node 2 on the other hand refers to Arthur Guinness, the founder of the Guinness brewery.

The surface name of node 2 ("A. Guinness") alone would not be sufficient to perform a correct matching as it could refer to both Alec Guinness and Arthur Guinness.

However, since links in the graph reveal the occupations of the persons, a relational approach can perform the correct matching

Link-based Clustering

- Extend feature-based clustering to relational learning
 - Group entities in relational database using similarity
- Entities are not only grouped by similarity of their features but also by similarity of their links
- In social networks, it is community detection

Statistical Relation Learning (SRL)

- Creation of statistical models for relational data
- Triples are assumed to be incomplete and noisy
- Entities and relation types may contain duplicates

Probabilistic Knowledge Graphs

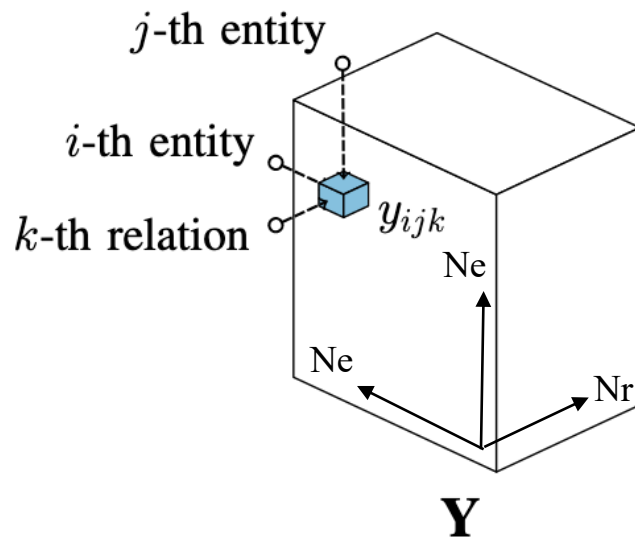
- $E = \{e_1, e_2, \dots, e_{Ne}\}$ is the set of entities
- $R = \{r_1, r_2, \dots, r_{Nr}\}$ the set of all relations in a KG
 - Each possible triple is defined as $X_{ijk} = (e_i, r_k, e_j)$
 - over set of entities E and relations R
 - as a binary r.v. $y_{ijk} \in \{0,1\}$ that indicates its existence
- All possible triples in $E \times R \times E$ can be grouped in an *Adjacency tensor* (a three-way array)
 - $Y \in \{0,1\}^{Ne \times Ne \times Nr}$, whose entries are set such that

$$y_{ijk} = \begin{cases} 1, & \text{if the triple } (e_i, r_k, e_j) \text{ exists} \\ 0, & \text{otherwise.} \end{cases}$$

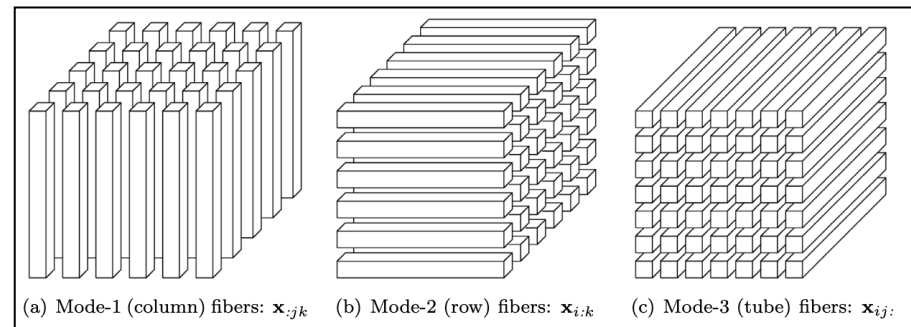
 - Each possible realization of Y can be interpreted as a possible world.

Tensor Representation

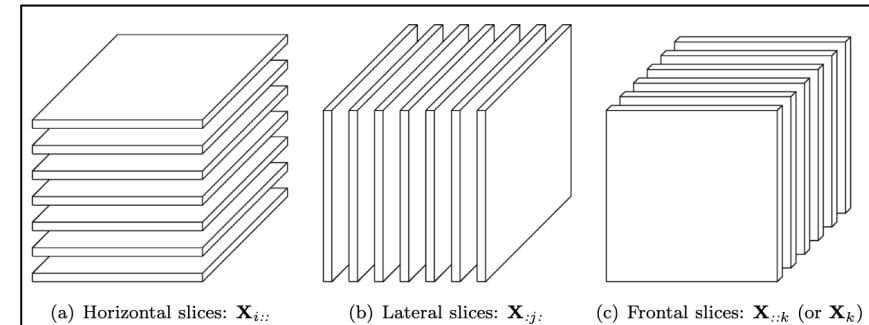
- Binary relational data is represented as a tensor



Fibers of a 3rd order tensor



Slices of a 3rd order tensor



Model for Knowledge Graph

- We are interested in estimating the joint distribution $P(Y)$, from a subset $D \subseteq E \times R \times E \times \{0,1\}$ of observed triples
- In doing so, we are estimating a probability distribution over possible worlds, which allows us to predict the probability of triples based on the state of the entire knowledge graph

Size of the Adjacency Tensor

- Y can be enormous for large knowledge graphs
 - E.g., in Freebase, with over 40 million entities and 35,000 relations, the no. of possible triples $E \times R \times E$ exceeds 10^{19} elements
- Of course, type constraints reduce this number considerably
- Even among syntactically valid triples, only a tiny fraction are likely to be true.
 - For example, there are over 450,000 actors and over 250,000 movies stored in Freebase
 - But each actor stars only in a small no. of movies

Statistical Properties of KGs

- KGs typically adhere to some deterministic rules, such as type constraints and transitivity
 - Ex: if Leonard Nimoy was born in Boston, and Boston is located in the USA, then we can infer that Leonard Nimoy was born in the USA
- Also various “softer” statistical patterns or regularities, which are not universally true but nevertheless have useful predictive power
 - Homophily: tendency of entities to be related to other entities with similar characteristics
 - US-born actors are more likely to star in US-made movies

Types of SRL Models

- Triples correlated with certain other triples
 - i.e., random variables y_{ijk} correlated with each other
- Three main ways to model these correlations:
 - 1.M1: Latent feature models
 - y_{ijk} are conditionally independent given latent features associated with SPO type and additional parameters
 - 2.M2: Graph feature models
 - y_{ijk} are conditionally independent given observed graph features and additional parameters
 - 3.M3: Markov random fields
 - Assume all y_{ijk} have local interactions

Probability Models M1 and M2

- Model classes M1 and M2 predict the existence of a triple x_{ijk} via a score function $f(x_{ijk}; \theta)$
 - which represents the model's confidence that a triple exists given the parameters θ
- The conditional independence assumptions of M1 and M2 allow the probability model to be written as follows:

$$P(\underline{\mathbf{Y}}|\mathcal{D}, \Theta) = \prod_{i=1}^{N_e} \prod_{j=1}^{N_e} \prod_{k=1}^{N_r} \text{Ber}(y_{ijk} | \sigma(f(x_{ijk}; \Theta))) \quad (1)$$

where $\sigma(u) = 1/(1 + e^{-u})$ is the sigmoid (logistic) function, and

$$\text{Ber}(y|p) = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases} \quad (2)$$

is the Bernoulli distribution.

Latent Feature Models

- Variables y_{ijk} are conditionally independent given global latent features and parameters:

$$P(\underline{\mathbf{Y}}|\mathcal{D}, \Theta) = \prod_{i=1}^{N_e} \prod_{j=1}^{N_e} \prod_{k=1}^{N_r} \text{Ber}(y_{ijk} | \sigma(f(x_{ijk}; \Theta))) \quad (1)$$

where $\sigma(u) = 1/(1 + e^{-u})$ is the sigmoid (logistic) function, and

$$\text{Ber}(y|p) = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases} \quad (2)$$

is the Bernoulli distribution.

- We discuss forms for score function $f(x_{ijk}; \theta)$ below
 - All models explain triples via latent features

Example of Latent Features

- **Fact:** Alec Guinness **received** Academy Award
- **Possible Explanation:** he is a good actor
 - Uses latent features of entities (good actor)
 - “Latent” because not directly observed in the data
- **Task:** infer these features from data
 - Denote latent features of entity e_i by vector $\mathbf{e}_i \in R^{H_e}$ where H_e is no. of latent features in the model
 - A model for:
 - Alec Guinness **is a good actor and**
 - Academy Award **is prestigious**

$$\mathbf{e}_{\text{Guinness}} = \begin{bmatrix} 0.9 \\ 0.2 \end{bmatrix}, \quad \mathbf{e}_{\text{AcademyAward}} = \begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$$

Types of Latent Feature Models

- Intuition behind relational latent feature models:
 - Relationships between entities can be derived from interactions of their latent features
- Methods to model interactions, and to derive the existence of a relationship from them

1. Bilinear model

$$f_{ijk}^{\text{RESICAL}} := \mathbf{e}_i^\top \mathbf{W}_k \mathbf{e}_j = \sum_{a=1}^{H_e} \sum_{b=1}^{H_e} w_{abk} e_{ia} e_{jb}$$

2. Other tensor factorization models

3. Matrix factorization methods

4. Multi-layer perceptrons

- Based on entities (E) and entity relationships (ER)

$$\begin{aligned} f_{ijk}^{\text{E-MLP}} &:= \mathbf{w}_k^\top \mathbf{g}(\mathbf{h}_{ijk}^a) \\ \mathbf{h}_{ijk}^a &:= \mathbf{A}_k^\top \boldsymbol{\phi}_{ij}^{\text{E-MLP}} \\ \boldsymbol{\phi}_{ij}^{\text{E-MLP}} &:= [\mathbf{e}_i; \mathbf{e}_j] \end{aligned}$$

\mathbf{h}^a is an additive hidden layer

$$\begin{aligned} f_{ijk}^{\text{ER-MLP}} &:= \mathbf{w}^\top \mathbf{g}(\mathbf{h}_{ijk}^c) \\ \mathbf{h}_{ijk}^c &:= \mathbf{C}^\top \boldsymbol{\phi}_{ijk}^{\text{ER-MLP}} \\ \boldsymbol{\phi}_{ijk}^{\text{ER-MLP}} &:= [\mathbf{e}_i; \mathbf{e}_j; \mathbf{r}_k]. \end{aligned}$$

RESCAL

- RESCAL is an embedding method for learning from knowledge graphs
 - For tasks like link prediction and entity resolution
 - Scalable to KGs with millions of entities and billions of facts
 - Provide access to relational information for deep learning methods
- RESCAL = Relational + Scalable

RESCAL (Bilinear Model)

- Explains triples via pairwise latent features
 - The score of a triple x_{ijk} is modeled as

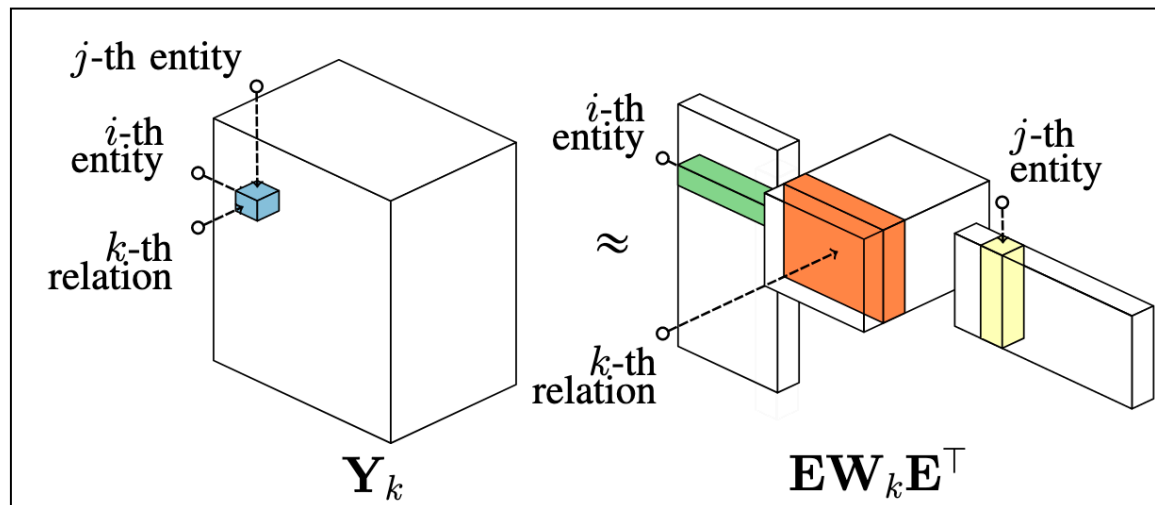
$$f_{ijk}^{\text{RESCAL}} := \mathbf{e}_i^\top \mathbf{W}_k \mathbf{e}_j = \sum_{a=1}^{H_e} \sum_{b=1}^{H_e} w_{abk} e_{ia} e_{jb}$$

- where $\mathbf{W}_k \in R^{H_e \times H_e}$ with entries w_{abk} specify how latent features a and b interact in the k^{th} relation
 - It is bilinear, since interactions between entity vectors are multiplicative
- Ex: Model the pattern that good actors receive prestigious awards via

$$\mathbf{W}_{\text{receivedAward}} = \begin{bmatrix} 0.1 & 0.9 \\ 0.1 & 0.1 \end{bmatrix}$$

RESCAL as Tensor Factorization

RESCAL is similar to methods used in recommendation systems, and to traditional tensor factorization methods



Equation

$$f_{ijk}^{\text{RESCAL}} := \mathbf{e}_i^\top \mathbf{W}_k \mathbf{e}_j = \sum_{a=1}^{H_e} \sum_{b=1}^{H_e} w_{abk} e_{ia} e_{jb}$$

can be written compactly as

$$F_k = E W_k E^T$$

where $F_k \in R^{Ne \times Ne}$ is the matrix holding all scores for the k^{th} relation and the i^{th} row of $E \in R^{Ne \times He}$ holds the latent representation of e_i

Multi-Layer Perceptrons

- We can interpret RESCAL as creating composite representations of triples and predicting their existence from this representation
- In particular, we can rewrite RESCAL as

$$\begin{aligned} f_{ijk}^{\text{RESCAL}} &:= \mathbf{w}_k^\top \phi_{ij}^{\text{RESCAL}} \\ \phi_{ij}^{\text{RESCAL}} &:= \mathbf{e}_j \otimes \mathbf{e}_i, \end{aligned}$$

– Where $\mathbf{w}_k = \text{vec}(\mathbf{W}_k)$

RESCAL as an MLP

MLP based on *entities* (E)

$$\begin{aligned} f_{ijk}^{\text{E-MLP}} &:= \mathbf{w}_k^\top \mathbf{g}(\mathbf{h}_{ijk}^a) \\ \mathbf{h}_{ijk}^a &:= \mathbf{A}_k^\top \phi_{ij}^{\text{E-MLP}} \\ \phi_{ij}^{\text{E-MLP}} &:= [\mathbf{e}_i; \mathbf{e}_j] \end{aligned}$$

One disadvantage of the E-MLP is that it has to define a vector \mathbf{w}_k and a matrix \mathbf{A}_k for every possible relation, which requires $H_a + (H_a \times 2H_e)$ parameters per relation.

An alternative is to embed the relation itself, using a H_r -dimensional vector r_k

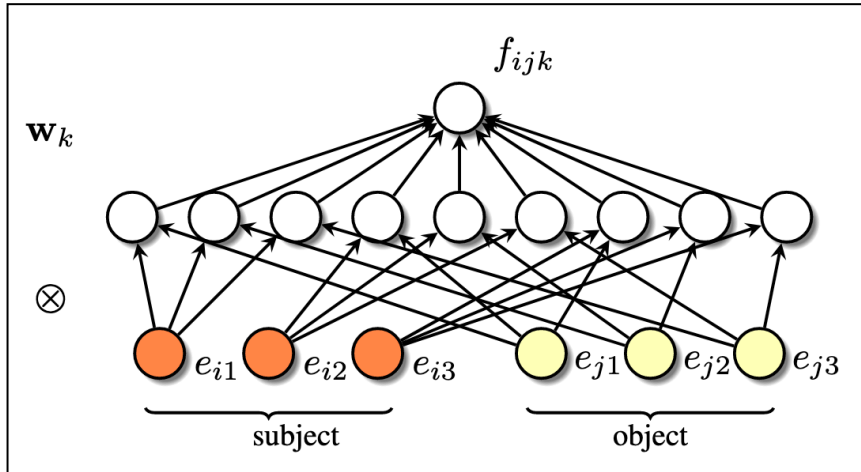
MLP based on *entities and relationships* (ER)

$$\begin{aligned} f_{ijk}^{\text{ER-MLP}} &:= \mathbf{w}^\top \mathbf{g}(\mathbf{h}_{ijk}^c) \\ \mathbf{h}_{ijk}^c &:= \mathbf{C}^\top \phi_{ijk}^{\text{ER-MLP}} \\ \phi_{ijk}^{\text{ER-MLP}} &:= [\mathbf{e}_i; \mathbf{e}_j; \mathbf{r}_k]. \end{aligned}$$

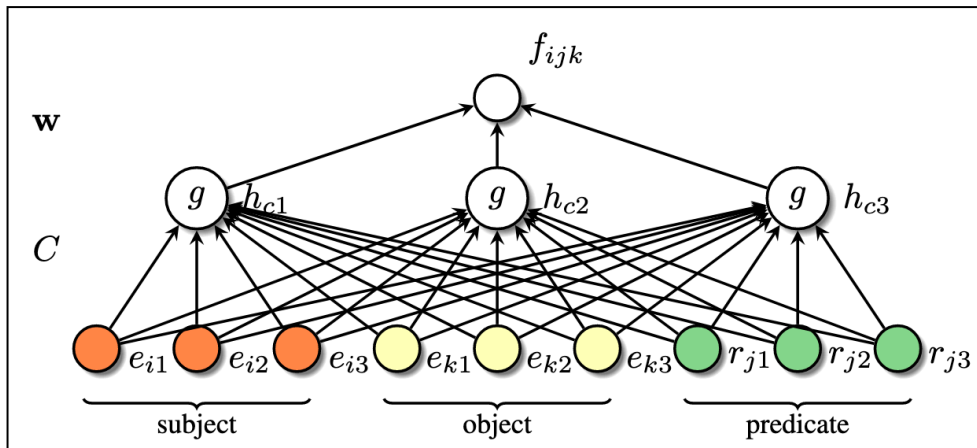
ER-MLP uses a global weight vector for all relations.

This model was used in the KV project since it has many fewer parameters than the E-MLP; the reason is that \mathbf{C} is independent of the relation k

RESCAL and ER-MLP as Neural Nets



$$\begin{aligned}
 f_{ijk}^{\text{E-MLP}} &:= \mathbf{w}_k^\top \mathbf{g}(\mathbf{h}_{ijk}^a) \\
 \mathbf{h}_{ijk}^a &:= \mathbf{A}_k^\top \phi_{ij}^{\text{E-MLP}} \\
 \phi_{ij}^{\text{E-MLP}} &:= [\mathbf{e}_i; \mathbf{e}_j]
 \end{aligned}$$



$$\begin{aligned}
 f_{ijk}^{\text{ER-MLP}} &:= \mathbf{w}^\top \mathbf{g}(\mathbf{h}_{ijk}^c) \\
 \mathbf{h}_{ijk}^c &:= \mathbf{C}^\top \phi_{ijk}^{\text{ER-MLP}} \\
 \phi_{ijk}^{\text{ER-MLP}} &:= [\mathbf{e}_i; \mathbf{e}_j; \mathbf{r}_k].
 \end{aligned}$$

$H_e=H_r=3$ and $H_a=3$. Note, that the inputs are latent features.
The symbol g denotes the application of the function $g(\cdot)$

Embedding of relations in ER-MLP

- Nearest neighbors of latent representations of selected relations computed with a 60 dimensional model on Freebase
 - Nos. represent squared Euclidean distances.
- Semantically related relations near each other
 - Closest relations to children relation are:
 - parents, spouse, and birth-place

Relation	Nearest Neighbors					
children	parents	(0.4)	spouse	(0.5)	birth-place	(0.8)
birth-date	children	(1.24)	gender	(1.25)	parents	(1.29)
edu-end ¹⁰	job-start	(1.41)	edu-start	(1.61)	job-end	(1.74)

Neural Tensor Networks (NTNs)

- NTN is a combination of traditional neural networks with bilinear models
 - The NTN model:

$$\begin{aligned}f_{ijk}^{\text{NTN}} &:= \mathbf{w}_k^\top \mathbf{g}([\mathbf{h}_{ijk}^a; \mathbf{h}_{ijk}^b]) \\ \mathbf{h}_{ijk}^a &:= \mathbf{A}_k^\top [\mathbf{e}_i; \mathbf{e}_j] \\ \mathbf{h}_{ijk}^b &:= [\mathbf{e}_i^\top \mathbf{B}_k^1 \mathbf{e}_j, \dots, \mathbf{e}_i^\top \mathbf{B}_k^{H_b} \mathbf{e}_j]\end{aligned}$$

- Here \mathbf{B}_k is a tensor, where the l^{th} slice \mathbf{B}_k^l has size $H_e \times H_e$, and there are H_b slices
 - \mathbf{h}_{ijk}^b is a bilinear hidden layer, since it is derived from a weighted combination of multiplicative terms
- With more parameters than E-MLP or RESCAL models it tends to overfit

Latent Distance Models

- Derive probability relationships from distance between latent representations of entities
 - Entities are likely related if representations are near
 - model probability x_{ij} via score $f(e_i, e_j) = -d(e_i, e_j)$

1. Structure Embeddings Model

- A_k^s, A_k^o transform global latent features of entities to model relationships specifically for the k^{th} relation

$$f_{ijk}^{\text{SE}} := -\|A_k^s \mathbf{e}_i - A_k^o \mathbf{e}_j\|_1 = -\|\mathbf{h}_{ijk}^a\|_1$$

2. TransE Model

- Translates latent features via a relation-specific offset instead of matrix multiplications.

- Score of a triple x_{ijk} is defined as: $f^{\text{TransE}}_{ijk} := -d(e_i + r_k, e_j)$

Summary of Latent Feature Models

Method	f_{ijk}	\mathbf{A}_k	\mathbf{C}	\mathbf{B}_k	Num. Parameters
RESCAL [64]	$\mathbf{w}_k^\top \mathbf{h}_{ijk}^b$	-	-	$[\delta_{1,1}, \dots, \delta_{H_e, H_e}]$	$N_r H_e^2 + N_e H_e$
E-MLP [92]	$\mathbf{w}_k^\top \mathbf{g}(\mathbf{h}_{ijk}^a)$	$[\mathbf{A}_k^s; \mathbf{A}_k^o]$	-	-	$N_r(H_a + H_a \times 2H_e) + N_e H_e$
ER-MLP [28]	$\mathbf{w}_k^\top \mathbf{g}(\mathbf{h}_{ijk}^c)$	-	\mathbf{C}	-	$H_c + H_c \times (2H_e + H_r) + N_r H_r + N_e H_e$
NTN [92]	$\mathbf{w}_k^\top \mathbf{g}([\mathbf{h}_{ijk}^a; \mathbf{h}_{ijk}^b])$	$[\mathbf{A}_k^s; \mathbf{A}_k^o]$	-	$[\mathbf{B}_k^1, \dots, \mathbf{B}_k^{H_b}]$	$N_e^2 H_b + N_r(H_b + H_a) + 2N_r H_e H_a + N_e H_e$
Structured Embeddings [93]	$-\ \mathbf{h}_{ijk}^a\ _1$	$[\mathbf{A}_k^s; -\mathbf{A}_k^o]$	-	-	$2N_r H_e H_a + N_e H_e$
TransE [94]	$-(2h_{ijk}^a - 2h_{ijk}^b + \ \mathbf{r}_k\ _2^2)$	$[\mathbf{r}_k; -\mathbf{r}_k]$	-	\mathbf{I}	$N_r H_e + N_e H_e$

Best model will be dataset dependent

ER-MLP model outperformed the NTN model on a particular dataset

RESCAL worked best on two link prediction tasks.

Graph Feature Models

- Existence of an edge predicted by extracting features from observed edges in the graph
 - Consider existence of the path
John -- parentOf \rightarrow Anne \leftarrow parentOf -- Mary
representing a common child
 - So we could predict the triple (John, marriedTo, Mary)
- Explains triples directly from observed triples in the KG

Training Statistical Relation Learning

- We have a set of N_d observed triples and let the n^{th} triple be denoted by x_n
- Each observed triple is either true (denoted $y_n=1$) or false (denoted $y_n=0$)
- The labeled dataset is $D = \{(x_n, y_n) \mid n=1, \dots, N_d\}$
- Given this, a natural way to estimate the parameters θ is to compute the maximum a posteriori (MAP) estimate:

$$\max_{\Theta} \sum_{n=1}^{N_d} \log \text{Ber}(y^n \mid \sigma(f(x^n; \Theta))) + \log p(\Theta \mid \lambda)$$

- where λ controls the strength of the prior

Loss Function

- We can equivalently state this as a regularized loss minimization problem

$$\min_{\Theta} \sum_{n=1}^N \mathcal{L}(\sigma(f(x^n; \Theta)), y^n) + \lambda \text{reg}(\Theta)$$

- Where $L(p, y) = -\log \text{Ber}(y|p)$ is the log loss
- Another loss is squared loss, $L(p, y) = (p - y)^2$
 - Using the squared loss can be especially efficient in combination with a closed-world assumption (CWA)
 - Minimization for RESCAL becomes

$$\min_{\mathbf{E}, \{\mathbf{W}_k\}} \sum_k \|\mathbf{Y}_k - \mathbf{E}\mathbf{W}_k\mathbf{E}^\top\|_F^2 + \lambda_1 \|\mathbf{E}\|_F^2 + \lambda_2 \sum_k \|\mathbf{W}_k\|_F^2$$

Markov Random Fields

- We drop the assumption that the random variables y_{ijk} in Y are conditionally independent
- However, in the case of relational data and without the conditional independence assumption, each y_{ijk} can depend on any of the other $N_e \times N_e \times N_r - 1$ random variables in Y
- Due to this enormous number of possible dependencies, it is intractable to estimate the joint distribution $P(Y)$ without further constraints

MRFs are Dependency graphs

- Each random variable (in our case, a possible fact y_{ijk}) is represented as a node in the graph, while each dependency between random variables is represented as an edge
- To distinguish MRFs from knowledge graphs, we refer to them as *dependency graphs*
 - While knowledge graphs encode the existence of facts, dependency graphs encode statistical dependencies between random variables

Definition of MRF

- MRF has the form

$$P(\underline{\mathbf{Y}}|\boldsymbol{\theta}) = \frac{1}{Z} \prod_c \psi(\mathbf{y}_c|\boldsymbol{\theta})$$

- where $\psi(\mathbf{y}_c|\boldsymbol{\theta}) \geq 0$ is a *potential function* on the c^{th} clique in the dependency graph, and

$$Z = \sum_{\mathbf{y}} \prod_c \psi(\mathbf{y}_c|\boldsymbol{\theta})$$

- is the partition function to ensure distribution sums to one
- potential functions capture local correlations between variables in each clique c
- Defines a distribution over “possible worlds”
 - i.e., over random variables \mathbf{Y}
- Structure of MRF derived using Markov logic
 - a template language based on logical formulae

Markov Logic Networks

- Combining probability and first-order logic in a single representation
 - PGMs enable efficient handling of uncertainty
 - First-order logic enables compact representation of a wide variety of knowledge

First Order Knowledge Base

- A first-order knowledge base(KB) is a set of sentences or formulas in first order logic
- Formulas are constructed using four types of symbols: constants, variables, functions, and predicates.
- Constant symbols represent objects in the domain of interest (e.g., people: Anna, Bob, Chris, etc.).
- Variable symbols range over objects in the domain
- Function symbols (e.g.,MotherOf) represent mappings from tuples of objects to objects.
- Predicate symbols represent relations among objects in the domain (e.g.,Friends) or attributes of objects(e.g.,Smokes)

First-Order Logic

- An atomic formula or atom is a predicate symbol applied to a tuple of terms e.g., $\text{Friends}(x; \text{MotherOf}(\text{Anna}))$
- Formulas are recursively constructed from atomic formulas using logical connectives and quantifiers.
- If F_1 and F_2 are formulas, the following are also formulas $\neg F_1$ (negation), which is true iff F_1 is false;
- $F_1 \wedge F_2$ (conjunction), which is true iff both F_1 and F_2 are true
- $F_1 \vee F_2$ (disjunction), which is true iff F_1 or F_2 is true; F_1)
- $F_1 \rightarrow F_2$ (implication), which is true iff F_1 is false or F_2 is true
- $F_1 \leftrightarrow F_2$ (equivalence), true iff F_1 and F_2 have same truth value
- $\forall x F_1$ (universal quantification), which is true iff F_1 is true for every object in the domain; and
- $\exists x F_1$ (existential quantification), which is true iff F_1 is true for at least one object x in the domain.

Definition of MLN

- An MLN L is a set of pairs $(F_i; w_i)$, where F_i is a formula in first-order logic and w_i is a real no.
- Together with constants $C = \{c_1; c_2; \dots; c_{|C|}\}$, it defines MRF $M_{L;C} \left(\boxed{P(X=x) = \frac{1}{Z} \exp \left(\sum_j w_j f_j(x) \right)} \right)$ as follows:
 - 1. $M_{L;C}$ contains one binary node for each possible grounding of each predicate appearing in L
 - Value of node is 1 if ground atom is true, and 0 otherwise
 - 2. $M_{L;C}$ contains one feature for each possible grounding of each formula F_i in L
 - The value of this feature is 1 if the ground formula is true, and 0 otherwise. The weight of the feature is the w_i associated with F_i in L

Dependency Graph from Markov Logic

- Given formulae $\mathcal{F} = \{F_i\}_{i=1}^L$, create an edge between nodes if corresponding facts occur in at least one grounded formula
 - A grounding of formula F_i is given by the assignment of entities to the variables in F_i
- Furthermore, we define $\psi(y_c|\theta)$ such that

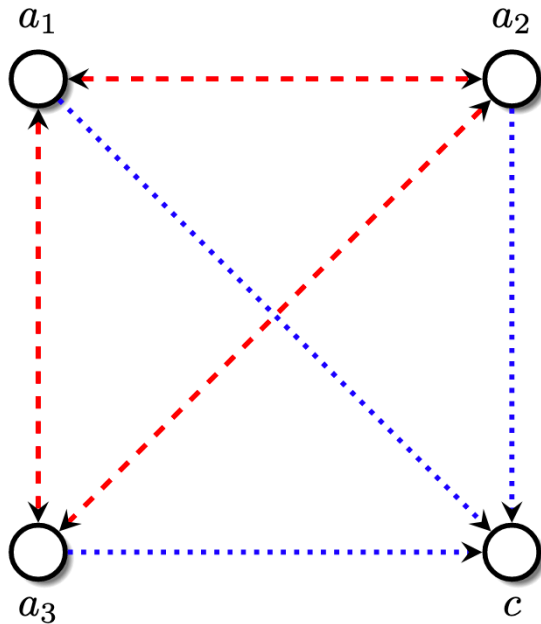
$$P(\underline{\mathbf{Y}}|\underline{\boldsymbol{\theta}}) = \frac{1}{Z} \prod_c \exp(\theta_c x_c)$$

- where x_c denotes no. of true groundings of F_c in \mathbf{Y} , and θ_c denotes the weight for formula F_c
 - If $\theta_c > 0$, we prefer worlds where formula F_c is satisfied
 - If $\theta_c < 0$, we prefer worlds where formula F_c is violated
 - If $\theta_c = 0$, then formula F_c is ignored

An example

- Two types of entities: adults and children
- Two types of relations: parentOf and marriedTo
 - These relations (edges) are correlated
 - since people who share a common child are often married, while people rarely marry their own children
- In Markov logic, we represent these dependencies using formulae such as:
 - $F_1: (x, \text{parentOf}, z) \wedge (y, \text{parentOf}, z) \Rightarrow (x, \text{marriedTo}, y)$
 - $F_2: (x, \text{marriedTo}, y) \Rightarrow \neg (y, \text{parentOf}, x)$

Simple Knowledge Graph



There are 4 entities (circles): 3 adults (a_1, a_2, a_3) and 1 child c

There are 2 types of edges:

adults may or may not be married to each other, as indicated by the red dashed edges,

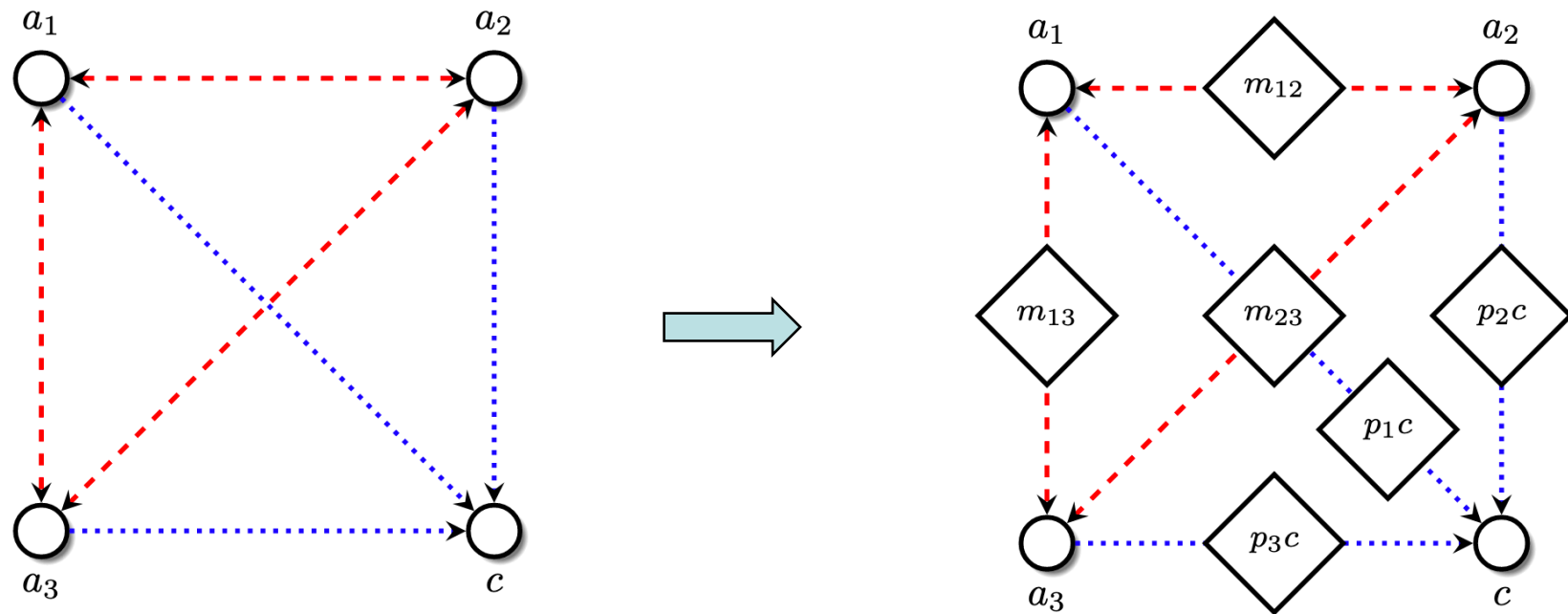
and the adults may or may not be parents of the child, as indicated by the blue dotted edges.

$$F_1: (x, \text{parentOf}, z) \wedge (y, \text{parentOf}, z) \Rightarrow (x, \text{marriedTo}, y)$$

$$F_2: (x, \text{marriedTo}, y) \Rightarrow \neg (y, \text{parentOf}, x)$$

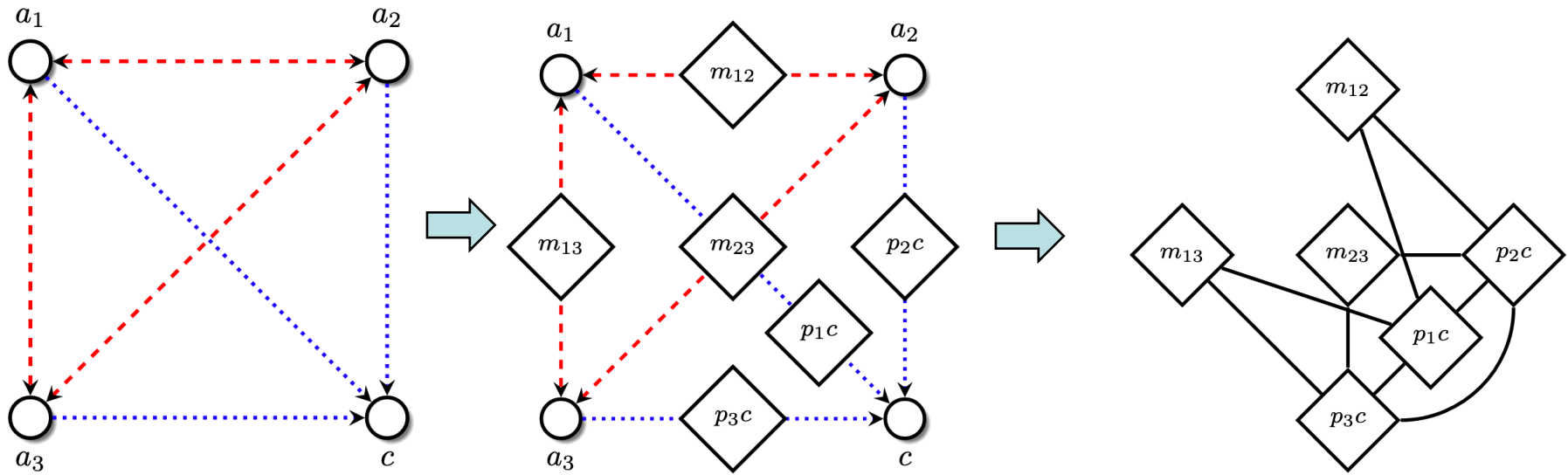
Six possible facts corresponding to the edges

Knowledge Graph relations



- We add binary random variables (represented by diamonds) to each KG edge
 - m : married, p : parent

Knowledge Graph to Dependency Graph



Create edges between these nodes if the corresponding facts occur in grounded formulae

$$F_1: (x, \text{parentOf}, z) \wedge (y, \text{parentOf}, z) \Rightarrow (x, \text{marriedTo}, y)$$

$$F_2: (x, \text{marriedTo}, y) \Rightarrow \neg (y, \text{parentOf}, x)$$

- For instance, grounding F_1 with $x=a_1, y=a_3$, and $z=c$, creates the edges $m_{13} \rightarrow p_{1c}$, $m_{13} \rightarrow p_{3c}$, and $p_{1c} \rightarrow p_{3c}$

MRF graph differs from KG

- The process of generating the MRF graph by applying templated rules to a set of entities is known as grounding or instantiation
- The topology of the resulting graph is quite different from the original KG
 - In particular, we have one node per possible KG edge, and these nodes are densely connected
 - This can cause computational difficulties

MRF Inference Tasks

- The inference problem consists of estimating:
 1. The most probable configuration $\mathbf{y}^* = \arg \max_{\mathbf{y}} p(\mathbf{y}|\theta)$
 2. The posterior marginals $p(\mathbf{y}_i|\theta)$
- Both problems are computationally intractable so heuristic approximations must be used
 - Solutions
 - For MAP Task: *Max Product Linear Programming*
 - For posterior marginals: *Gibbs sampling*
 - If potential functions are restricted to disjunctions (OR and NOT, but no AND), then we get a hinge loss MRF (HL-MRF)
 - For which efficient convex algorithms can be applied, based on a continuous relaxation of the binary random variables

MRF Learning Task

- Deals with specifying:
 - Form of potential functions (called *structure learning*)
 - Values for the numerical parameters θ
- In the case of MRFs for KGs, potential functions are often specified in the form of logical rules
 - In this case, structure learning is equivalent to rule learning
 - The parameter estimation problem (which is usually cast as MAP estimation), although convex, is in general quite expensive, since it needs to call inference as a subroutine

Conclusion on MRFs

- Approaches based on MRFs are very flexible
- But harder to make scalable inference and devise learning algorithms
 - Compared to methods based on observable or latent feature models

References

1. M. Nickel, K. Murphy, V. Tresp, E. Gabrilovich, “A Review of Relational Machine Learning for Knowledge Graphs”
Proceedings of the IEEE 104(1): 11-33 (2016)
2. M. Richardson, P. Domingos, “Markov logic networks”, *Machine Learning* 62, 107–136 (2006)