

# Relational Probabilistic Models

David Poole

Department of Computer Science,  
University of British Columbia

Work with: <http://minervaintelligence.com>, <https://treatment.com/>

March 20, 2019

# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 Unique properties of relational models
  - Learning general knowledge vs learning about a data set
  - Varying Populations
  - What can be observed?
- 4 Conclusions and Challenges

# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 Unique properties of relational models
  - Learning general knowledge vs learning about a data set
  - Varying Populations
  - What can be observed?
- 4 Conclusions and Challenges

# Choosing Individuals and Relations in Logic

First-order logical languages allow many different ways of representing facts.

E.g., How to represent: “Pen #7 is red.”

# Choosing Individuals and Relations in Logic

First-order logical languages allow many different ways of representing facts.

E.g., How to represent: “Pen #7 is red.”

- $red(pen_7)$ . It's easy to ask “What's red?”  
Can't ask “what is the color of  $pen_7$ ?”

# Choosing Individuals and Relations in Logic

First-order logical languages allow many different ways of representing facts.

E.g., How to represent: “Pen #7 is red.”

- $red(pen_7)$ . It’s easy to ask “What’s red?”  
Can’t ask “what is the color of  $pen_7$ ?”
- $color(pen_7, red)$ . It’s easy to ask “What’s red?”  
It’s easy to ask “What is the color of  $pen_7$ ?”  
Can’t ask “What property of  $pen_7$  has value  $red$ ?”

# Choosing Individuals and Relations in Logic

First-order logical languages allow many different ways of representing facts.

E.g., How to represent: “Pen #7 is red.”

- $red(pen_7)$ . It’s easy to ask “What’s red?”  
Can’t ask “what is the color of  $pen_7$ ?”
- $color(pen_7, red)$ . It’s easy to ask “What’s red?”  
It’s easy to ask “What is the color of  $pen_7$ ?”  
Can’t ask “What property of  $pen_7$  has value  $red$ ?”
- $prop(pen_7, color, red)$ . It’s easy to ask all these questions.

# Choosing Individuals and Relations in Logic

First-order logical languages allow many different ways of representing facts.

E.g., How to represent: “Pen #7 is red.”

- $red(pen_7)$ . It’s easy to ask “What’s red?”  
Can’t ask “what is the color of  $pen_7$ ?”
- $color(pen_7, red)$ . It’s easy to ask “What’s red?”  
It’s easy to ask “What is the color of  $pen_7$ ?”  
Can’t ask “What property of  $pen_7$  has value  $red$ ?”
- $prop(pen_7, color, red)$ . It’s easy to ask all these questions.

With a single relation it can be implicit  $\rightarrow$  triples:

$\langle pen_7, color, red \rangle$ .



# Universality of *prop*

To represent “a is a parcel”

# Universality of *prop*

To represent “a is a parcel”

- $prop(a, type, parcel)$ , where *type* is a special property and *parcel* is a class.
- $prop(a, parcel, true)$ , where *parcel* is a Boolean property (characteristic function of the class *parcel*).

# Triples

- To represent *tutorial("StarAI", nips2017, 1045, hallC)*. "the Star AI tutorial at NIPS 2017 is scheduled at 10:45 in Hall C."

# Triples

- To represent *tutorial("StarAI", nips2017, 1045, hallC)*. “the Star AI tutorial at NIPS 2017 is scheduled at 10:45 in Hall C.”
- Let *t123* name the offering of the tutorial:

*prop(t123, type, tutorial)*.

*prop(t123, title, "StarAI")*.

*prop(t123, event, nips2017)*.

*prop(t123, time, 1045)*.

*prop(t123, room, hallC)*.

# Triples

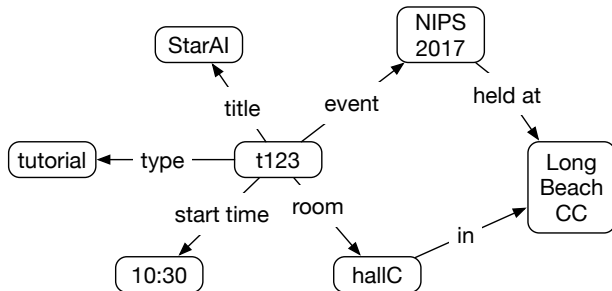
- To represent *tutorial("StarAI", nips2017, 1045, hallC)*. “the Star AI tutorial at NIPS 2017 is scheduled at 10:45 in Hall C.”
- Let *t123* name the offering of the tutorial:
  - prop(t123, type, tutorial)*.
  - prop(t123, title, "StarAI")*.
  - prop(t123, event, nips2017)*.
  - prop(t123, time, 1045)*.
  - prop(t123, room, hallC)*.
- We have **reified** the booking.
- Reify means: to make into an individual.

# Triples

- To represent *tutorial("StarAI", nips2017, 1045, hallC)*. "the Star AI tutorial at NIPS 2017 is scheduled at 10:45 in Hall C."
- Let *t123* name the offering of the tutorial:
  - prop(t123, type, tutorial)*.
  - prop(t123, title, "StarAI")*.
  - prop(t123, event, nips2017)*.
  - prop(t123, time, 1045)*.
  - prop(t123, room, hallC)*.
- We have **reified** the booking.
- Reify means: to make into an individual.
- How can we add extra arguments (e.g., presenters, chair)?

# Triples and Knowledge Graphs

- Subject–verb–object  
Individual–property–value  
triples can be depicted as edges in graphs



- A modeller or learner needs to invent (reified) objects.

# Triples are universal representations of relations

All relations can be represented in terms of **triples**:

	...	$P_j$	...
$r_i$	...	...	...
	...	$v_{ij}$	...
	...	...	...
	...	...	...

can be represented as



# Triples are universal representations of relations

All relations can be represented in terms of **triples**:

	...	$P_j$	...
$r_i$	...	...	...
	...	$v_{ij}$	...
	...	...	...

can be represented as the triple  $\langle r_i, P_j, v_{ij} \rangle$ .

- $r_i$  becomes a **reified** individual.

# Triples are universal representations of relations

All relations can be represented in terms of **triples**:

	...	$P_j$	...
$r_i$	...	...	...
	...	$v_{ij}$	...
	...	...	...

can be represented as the triple  $\langle r_i, P_j, v_{ij} \rangle$ .

- $r_i$  becomes a **reified** individual.
- **Examples of reified individuals**: a booking, a marriage, a talk, a lab report, an event, a party, a meeting, a drink

# Triples are universal representations of relations

All relations can be represented in terms of **triples**:

	...	$P_j$	...
$r_i$	...	...	...
	...	$v_{ij}$	...
	...	...	...

can be represented as the triple  $\langle r_i, P_j, v_{ij} \rangle$ .

- $r_i$  becomes a **reified** individual.
- **Examples of reified individuals**: a booking, a marriage, a talk, a lab report, an event, a party, a meeting, a drink
- Challenge for learning: each reified individual has limited data to learn from; (at most) one value for each property.

# Triples are universal representations of relations

All relations can be represented in terms of **triples**:

	...	$P_j$	...
$r_i$	...	...	...
	...	$v_{ij}$	...
	...	...	...

can be represented as the triple  $\langle r_i, P_j, v_{ij} \rangle$ .

- $r_i$  becomes a **reified** individual.
- **Examples of reified individuals**: a booking, a marriage, a talk, a lab report, an event, a party, a meeting, a drink
- Challenge for learning: each reified individual has limited data to learn from; (at most) one value for each property.

$\text{prop}(\text{Individual}, \text{Property}, \text{Value})$  is the only relation needed:  
 **$\langle \text{Individual}, \text{Property}, \text{Value} \rangle$  triples, Semantic network, entity relationship model, knowledge graphs, ...**

# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 Unique properties of relational models
  - Learning general knowledge vs learning about a data set
  - Varying Populations
  - What can be observed?
- 4 Conclusions and Challenges

# Vector & Tensor Representations of Entities & Relations

- Knowledge graphs (semantic networks) are databases of triples
- We want to estimate  $P(\langle h, r, t \rangle)$

# Vector & Tensor Representations of Entities & Relations

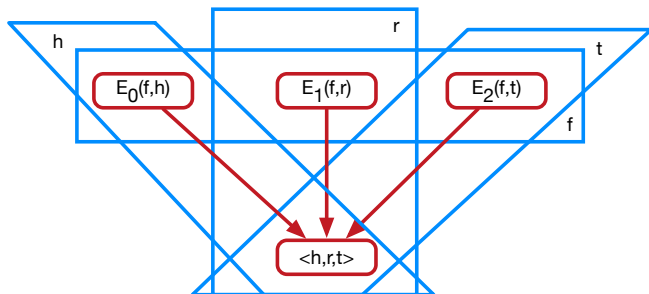
- Knowledge graphs (semantic networks) are databases of triples
- We want to estimate  $P(\langle h, r, t \rangle)$
- Polyadic decomposition model (1927): two vector embeddings for each entity, and one for each relation

$$P(\langle h, r, t \rangle) = \text{sigmoid}\left(\sum_f E_0(f, h) * E_1(f, r) * E_2(f, t)\right)$$

# Vector & Tensor Representations of Entities & Relations

- Knowledge graphs (semantic networks) are databases of triples
- We want to estimate  $P(\langle h, r, t \rangle)$
- Polyadic decomposition model (1927): two vector embeddings for each entity, and one for each relation

$$P(\langle h, r, t \rangle) = \text{sigmoid}\left(\sum_f E_0(f, h) * E_1(f, r) * E_2(f, t)\right)$$





# Polyadic decomposition variations

- Polyadic decomposition doesn't work very well...
  - Consider  $\langle p123, likes, m53 \rangle$  and  $\langle m53, directed\_by, p534 \rangle$ .

# Polyadic decomposition variations

- Polyadic decomposition doesn't work very well...
  - Consider  $\langle p123, likes, m53 \rangle$  and  $\langle m53, directed\_by, p534 \rangle$ .
  - Requires two embeddings per entity, but head embeddings and tail embeddings do not interact.

# Polyadic decomposition variations

- Polyadic decomposition doesn't work very well...
  - Consider  $\langle p123, likes, m53 \rangle$  and  $\langle m53, directed\_by, p534 \rangle$ .
  - Requires two embeddings per entity, but head embeddings and tail embeddings do not interact.
- DistMult: share same embedding for head and tail.

# Polyadic decomposition variations

- Polyadic decomposition doesn't work very well...
  - Consider  $\langle p123, likes, m53 \rangle$  and  $\langle m53, directed\_by, p534 \rangle$ .
  - Requires two embeddings per entity, but head embeddings and tail embeddings do not interact.
- DistMult: share same embedding for head and tail.  
Problem:

# Polyadic decomposition variations

- Polyadic decomposition doesn't work very well...
  - Consider  $\langle p123, likes, m53 \rangle$  and  $\langle m53, directed\_by, p534 \rangle$ .
  - Requires two embeddings per entity, but head embeddings and tail embeddings do not interact.
- DistMult: share same embedding for head and tail.  
Problem: can only represent symmetric relations.

# Polyadic decomposition variations

- Polyadic decomposition doesn't work very well...
  - Consider  $\langle p123, likes, m53 \rangle$  and  $\langle m53, directed\_by, p534 \rangle$ .
  - Requires two embeddings per entity, but head embeddings and tail embeddings do not interact.
- DistMult: share same embedding for head and tail.  
Problem: can only represent symmetric relations.
- ComplexX: the embeddings are complex numbers, tail is the conjugate of the head embedding

# Polyadic decomposition variations

- Polyadic decomposition doesn't work very well...
  - Consider  $\langle p123, likes, m53 \rangle$  and  $\langle m53, directed\_by, p534 \rangle$ .
  - Requires two embeddings per entity, but head embeddings and tail embeddings do not interact.
- DistMult: share same embedding for head and tail.  
Problem: can only represent symmetric relations.
- ComplexX: the embeddings are complex numbers, tail is the conjugate of the head embedding
- SimpleE: have an embedding for  $r^{-1}$  and learn to predict both  $\langle h, r, t \rangle$  and  $\langle t, r^{-1}, h \rangle$

# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom.



# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom. Options:
  - learn general knowledge, e.g., transitivity, how male and female friendships work, how location affect friendship...

# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom. Options:
  - learn general knowledge, e.g., transitivity, how male and female friendships work, how location affect friendship...
  - learn specific knowledge about who is friends with who; e.g., which particular group of people are generally friends with each other.

# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom. Options:
  - learn general knowledge, e.g., transitivity, how male and female friendships work, how location affect friendship...
  - learn specific knowledge about who is friends with who; e.g., which particular group of people are generally friends with each other.
- The specific knowledge will tend to be more accurate on that population, but doesn't generalize to different populations.
- The general knowledge will tend to transfer better.
- Which is better depends on the goals and how success is measured.

# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom. Options:
  - learn general knowledge, e.g., transitivity, how male and female friendships work, how location affect friendship...
  - learn specific knowledge about who is friends with who; e.g., which particular group of people are generally friends with each other.
- The specific knowledge will tend to be more accurate on that population, but doesn't generalize to different populations.
- The general knowledge will tend to transfer better.
- Which is better depends on the goals and how success is measured.
- Many of the methods try to do both; learn about specific individuals and general knowledge.

# Challenges of learning knowledge graphs

- Evaluating predictions when only positive examples are provided

Consider the following relations:

- Married to
- Friend of
- Knows about
  
- Would get along with

# Challenges of learning knowledge graphs

- Evaluating predictions when only positive examples are provided

Consider the following relations:

- Married to — each person related to 0 or 1 other persons (with a few exceptions)
- Friend of — each person related to tens or hundreds of others
- Knows about — each person might know about hundreds or thousands of others. Some people may be known by millions or billions of others.
- Would get along with — almost everyone gets along with almost everyone else, but with some exceptions.

# Challenges of learning knowledge graphs

- Evaluating predictions when only positive examples are provided

Consider the following relations:

- Married to — each person related to 0 or 1 other persons (with a few exceptions)
  - Friend of — each person related to tens or hundreds of others
  - Knows about — each person might know about hundreds or thousands of others. Some people may be known by millions or billions of others.
  - Would get along with — almost everyone gets along with almost everyone else, but with some exceptions.
- Often we compare rankings (ordering), but what if the answer is “no”?

# Challenges of learning knowledge graphs

- Evaluating predictions when only positive examples are provided

Consider the following relations:

- Married to — each person related to 0 or 1 other persons (with a few exceptions)
- Friend of — each person related to tens or hundreds of others
- Knows about — each person might know about hundreds or thousands of others. Some people may be known by millions or billions of others.
- Would get along with — almost everyone gets along with almost everyone else, but with some exceptions.
- Often we compare rankings (ordering), but what if the answer is “no”?
- Difficult to learn about reified entities.



# Predicting Properties

- Tensor factorization models work well for predicting relations, but not for predicting properties.  
Tensor factorization relies on lower-dimensional representations, and there isn't one for properties.

# Predicting Properties

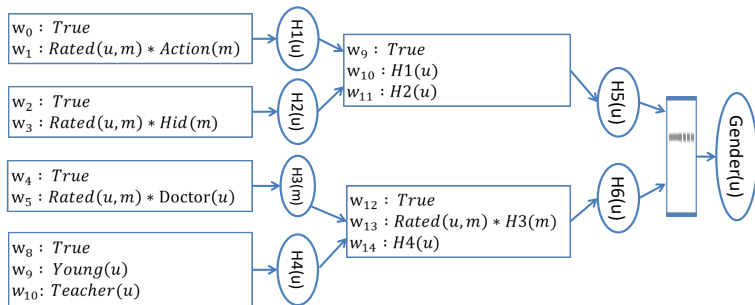
- Tensor factorization models work well for predicting relations, but not for predicting properties.  
Tensor factorization relies on lower-dimensional representations, and there isn't one for properties.
- Imagine trying to predict  $gender(P)$ , the gender of person  $P$ , and  $rating(P, M)$  the rating of person  $P$  on movie  $M$ .  
One of the embeddings of each person can just represent the gender — no generalization!  
— there are too many parameters

# Predicting Properties

- Tensor factorization models work well for predicting relations, but not for predicting properties.  
Tensor factorization relies on lower-dimensional representations, and there isn't one for properties.
- Imagine trying to predict  $gender(P)$ , the gender of person  $P$ , and  $rating(P, M)$  the rating of person  $P$  on movie  $M$ .  
One of the embeddings of each person can just represent the gender — no generalization!  
— there are too many parameters
- We can build relational neural networks to solve this

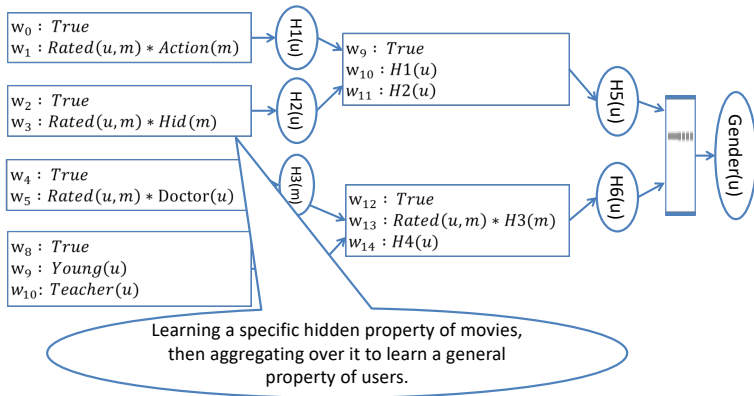
[Kazemi & Poole, AAAI 2017]

# Relational Neural Nets (ReINNs)



[Kazemi & Poole AAAI 2017]

# Relational Neural Nets (ReINNs)



[Kazemi & Poole AAAI 2017]

# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 Unique properties of relational models
  - Learning general knowledge vs learning about a data set
  - Varying Populations
  - What can be observed?
- 4 Conclusions and Challenges

# Desiderata for a Representation

- **Expressiveness:**  
Is it expressive enough to solve problem at hand?

# Desiderata for a Representation

- **Expressiveness:**  
Is it expressive enough to solve problem at hand?
- **Efficient Inference:**  
Is it efficient in the worst case or average case?  
Can it exploit structure (e.g., independencies and symmetries)



# Desiderata for a Representation

- **Expressiveness:**  
Is it expressive enough to solve problem at hand?
- **Efficient Inference:**  
Is it efficient in the worst case or average case?  
Can it exploit structure (e.g., independencies and symmetries)
- **Understandability or explainability:**  
Can people understand the model?  
Can a particular prediction be explained?

# Desiderata for a Representation

- **Expressiveness:**  
Is it expressive enough to solve problem at hand?
- **Efficient Inference:**  
Is it efficient in the worst case or average case?  
Can it exploit structure (e.g., independencies and symmetries)
- **Understandability or explainability:**  
Can people understand the model?  
Can a particular prediction be explained?
- **Learnability:** Can it be learned from:
  - heterogeneous data
  - prior knowledge

# Desiderata for a Representation

- **Expressiveness:**  
Is it expressive enough to solve problem at hand?
- **Efficient Inference:**  
Is it efficient in the worst case or average case?  
Can it exploit structure (e.g., independencies and symmetries)
- **Understandability or explainability:**  
Can people understand the model?  
Can a particular prediction be explained?
- **Learnability:** Can it be learned from:
  - heterogeneous data
  - prior knowledge
- **Modularity:**  
Can independently developed parts be combined to form larger model?  
Can a larger model be decomposed into smaller parts?

# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 Unique properties of relational models
  - Learning general knowledge vs learning about a data set
  - Varying Populations
  - What can be observed?
- 4 Conclusions and Challenges

# Directed vs Undirected Probabilistic Graphical Models

- **Undirected models** (Markov networks, factor graphs) represent probability distributions in terms of factors.
  - a factor is a non-negative function of a set of variables
  - variables in a factor are neighbours of each other
  - each variable is independent of its non-neighbours given its neighbours.

# Directed vs Undirected Probabilistic Graphical Models

- **Undirected models** (Markov networks, factor graphs) represent probability distributions in terms of factors.
  - a factor is a non-negative function of a set of variables
  - variables in a factor are neighbours of each other
  - each variable is independent of its non-neighbours given its neighbours.
- In **directed models**, factors represent conditional probabilities:
  - each variable is independent of its non-descendants given its parents.

# Directed vs Undirected Probabilistic Graphical Models

- **Undirected models** (Markov networks, factor graphs) represent probability distributions in terms of factors.
  - a factor is a non-negative function of a set of variables
  - variables in a factor are neighbours of each other
  - each variable is independent of its non-neighbours given its neighbours.
- In **directed models**, factors represent conditional probabilities:
  - each variable is independent of its non-descendants given its parents.
- $\{\textit{directed\_models}\} \subset \{\textit{undirected\_models}\}$

# Directed vs Undirected Probabilistic Graphical Models

- **Undirected models** (Markov networks, factor graphs) represent probability distributions in terms of factors.
  - a factor is a non-negative function of a set of variables
  - variables in a factor are neighbours of each other
  - each variable is independent of its non-neighbours given its neighbours.
- In **directed models**, factors represent conditional probabilities:
  - each variable is independent of its non-descendants given its parents.
- $\{\textit{directed\_models}\} \subset \{\textit{undirected\_models}\}$
- Directed (and undirected) models are universal: can represent any probability distribution over a finite set of variables



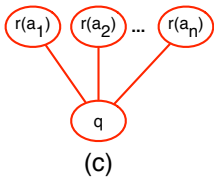
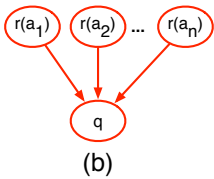
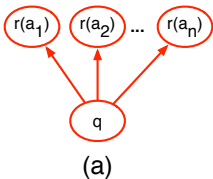
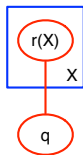
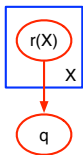
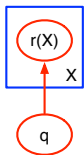
# Directed vs Undirected Probabilistic Graphical Models

- **Undirected models** (Markov networks, factor graphs) represent probability distributions in terms of factors.
  - a factor is a non-negative function of a set of variables
  - variables in a factor are neighbours of each other
  - each variable is independent of its non-neighbours given its neighbours.
- In **directed models**, factors represent conditional probabilities:
  - each variable is independent of its non-descendants given its parents.
- $\{\textit{directed\_models}\} \subset \{\textit{undirected\_models}\}$
- Directed (and undirected) models are universal: can represent any probability distribution over a finite set of variables
- Algorithms developed for undirected models work for both.

# Directed vs Undirected Probabilistic Graphical Models

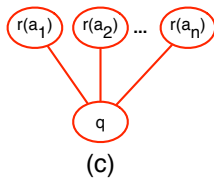
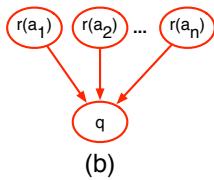
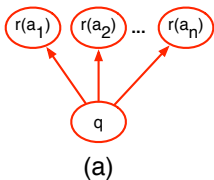
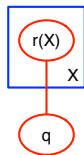
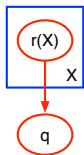
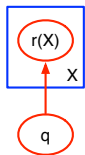
- **Undirected models** (Markov networks, factor graphs) represent probability distributions in terms of factors.
  - a factor is a non-negative function of a set of variables
  - variables in a factor are neighbours of each other
  - each variable is independent of its non-neighbours given its neighbours.
- In **directed models**, factors represent conditional probabilities:
  - each variable is independent of its non-descendants given its parents.
- $\{\textit{directed\_models}\} \subset \{\textit{undirected\_models}\}$
- Directed (and undirected) models are universal: can represent any probability distribution over a finite set of variables
- Algorithms developed for undirected models work for both.
- That does **not** mean that representations for undirected models can represent directed models.

# Three Elementary Models



(a)

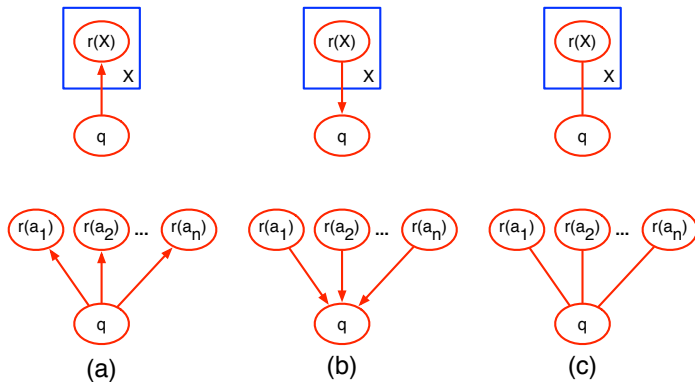
# Three Elementary Models



(a) Naïve Bayes

(b)

# Three Elementary Models

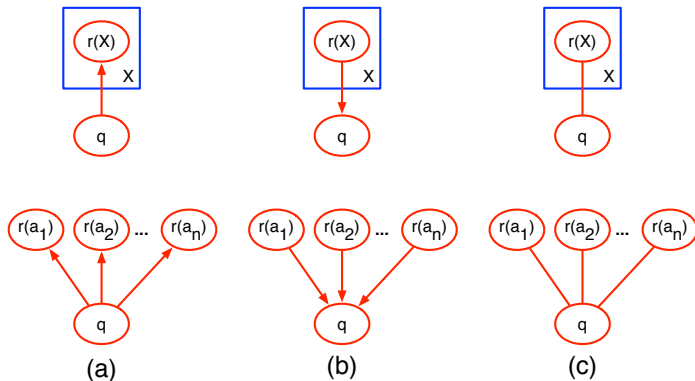


(a) Naïve Bayes

(b) (Relational) Logistic Regression

(c)

# Three Elementary Models



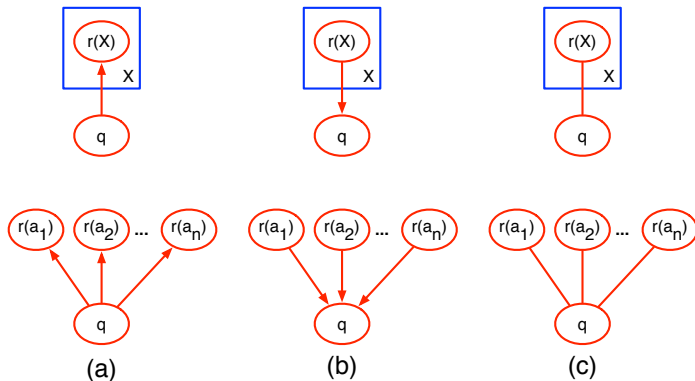
(a) Naïve Bayes

(b) (Relational) Logistic Regression

(c) Markov (Logic) network

} described by unary and pairwise factors

# Three Elementary Models

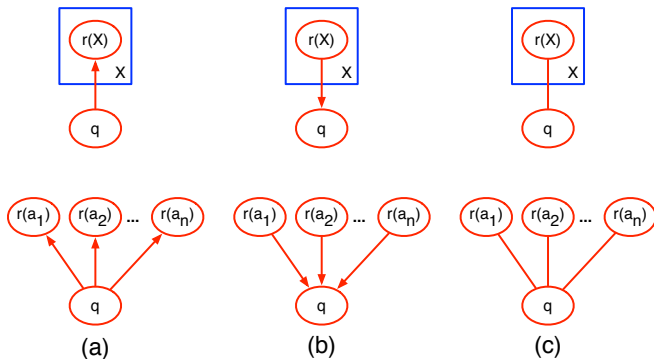


(a) Naïve Bayes  
 (b) (Relational) Logistic Regression  
 (c) Markov (Logic) network

} described by unary and pairwise factors

— They are identical models when all  $r$ 's are observed.

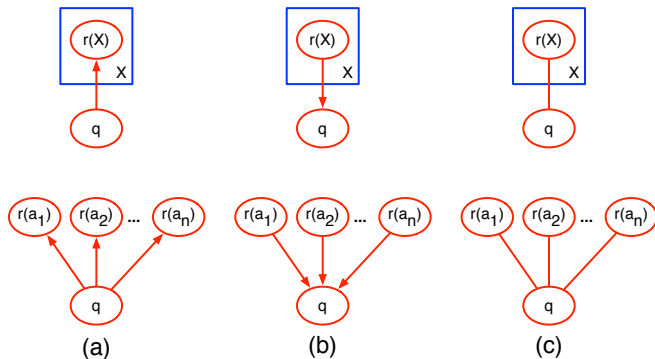
# Independence Assumptions



- Naïve Bayes (a) and Markov network (c):  $r(a_i)$  and  $r(a_j)$ 
  - are independent

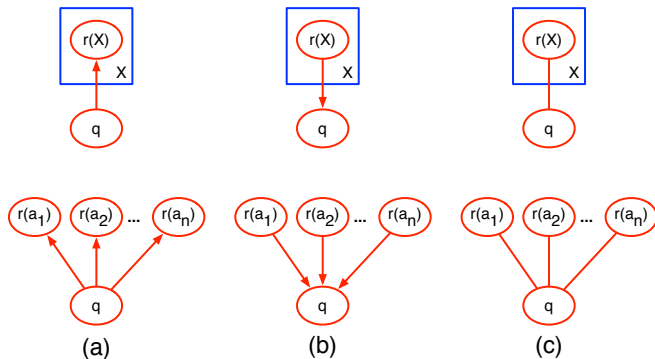


# Independence Assumptions



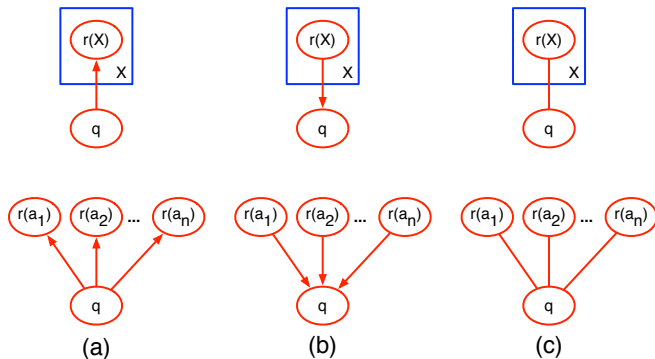
- Naïve Bayes (a) and Markov network (c):  $r(a_i)$  and  $r(a_j)$ 
  - are independent given  $Q$
  - are dependent

# Independence Assumptions



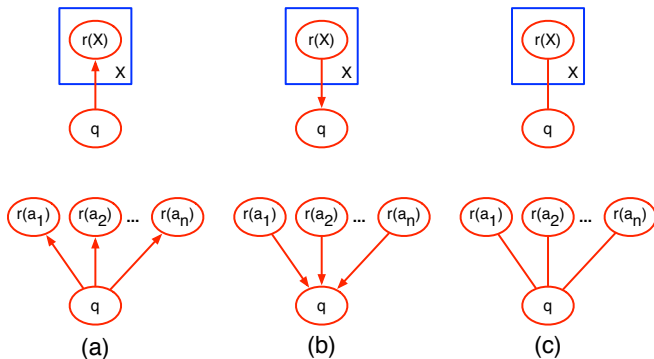
- Naïve Bayes (a) and Markov network (c):  $r(a_i)$  and  $r(a_j)$ 
  - are independent given  $Q$
  - are dependent not given  $Q$ .
- Directed model with aggregation (b):  $r(a_i)$  and  $r(a_j)$ 
  - are dependent

# Independence Assumptions



- Naïve Bayes (a) and Markov network (c):  $r(a_i)$  and  $r(a_j)$ 
  - are independent given  $Q$
  - are dependent not given  $Q$ .
- Directed model with aggregation (b):  $r(a_i)$  and  $r(a_j)$ 
  - are dependent given  $Q$ ,
  - are independent

# Independence Assumptions



- Naïve Bayes (a) and Markov network (c):  $r(a_i)$  and  $r(a_j)$ 
  - are independent given  $Q$
  - are dependent not given  $Q$ .
- Directed model with aggregation (b):  $r(a_i)$  and  $r(a_j)$ 
  - are dependent given  $Q$ ,
  - are independent not given  $Q$ .

# Modularity

- Directed models are inherently modular.  
 $P(q \mid r(X))$  is defined so that distribution over  $r(a_1) \dots r(a_n)$  is not affected.

# Modularity

- Directed models are inherently modular.  
 $P(q \mid r(X))$  is defined so that distribution over  $r(a_1) \dots r(a_n)$  is not affected.
- MLNs are provably not modular: If there is a distribution over  $r(a_1) \dots r(a_n)$  (e.g., they are independent),  $P(q \mid r(X))$  **cannot** be defined in an MLN so that
  - $q$  depends on the  $r$ 's ( $P(q \mid r(X)) \neq P(q)$ ) and
  - if  $q$  is summed out, the distribution over  $r(a_1) \dots r(a_n)$  is not changed.

# Modularity

- Directed models are inherently modular.  
 $P(q \mid r(X))$  is defined so that distribution over  $r(a_1) \dots r(a_n)$  is not affected.
- MLNs are provably not modular: If there is a distribution over  $r(a_1) \dots r(a_n)$  (e.g., they are independent),  $P(q \mid r(X))$  **cannot** be defined in an MLN so that
  - $q$  depends on the  $r$ 's ( $P(q \mid r(X)) \neq P(q)$ ) and
  - if  $q$  is summed out, the distribution over  $r(a_1) \dots r(a_n)$  is not changed.
  - **Why?** requires factors on arbitrary subsets of  $r(a_1) \dots r(a_k)$ 
    - pairwise (or 3-wise or ...) interactions are not adequate
    - can't marry the parents

# Cyclic Models

*Whether people smoke depends on whether their friends smoke.*

- MLN:

$$w : \text{smokes}(X) \leftarrow \text{friends}(X, Y) \wedge \text{smokes}(Y)$$



# Cyclic Models

*Whether people smoke depends on whether their friends smoke.*

- MLN:

$$w : \text{smokes}(X) \leftarrow \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

(where  $\leftarrow$  is material implication) is equivalent to

$$w : \text{true}(X) \wedge \text{true}(Y)$$

$$-w : \neg \text{smokes}(X) \wedge \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

# Cyclic Models

*Whether people smoke depends on whether their friends smoke.*

- MLN:

$$w : \text{smokes}(X) \leftarrow \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

(where  $\leftarrow$  is material implication) is equivalent to

$$w : \text{true}(X) \wedge \text{true}(Y)$$

$$\neg w : \neg \text{smokes}(X) \wedge \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

- Problog

$$w : \text{smokes}(X) \leftarrow \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

# Cyclic Models

*Whether people smoke depends on whether their friends smoke.*

- MLN:

$$w : \text{smokes}(X) \leftarrow \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

(where  $\leftarrow$  is material implication) is equivalent to

$$w : \text{true}(X) \wedge \text{true}(Y)$$

$$\neg w : \neg \text{smokes}(X) \wedge \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

- Problog

$$w : \text{smokes}(X) \leftarrow \exists Y \text{ friends}(X, Y) \wedge \text{smokes}(Y)$$

# Cyclic Models

*Whether people smoke depends on whether their friends smoke.*

- MLN:

$$w : \text{smokes}(X) \leftarrow \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

(where  $\leftarrow$  is material implication) is equivalent to

$$w : \text{true}(X) \wedge \text{true}(Y)$$

$$\neg w : \neg \text{smokes}(X) \wedge \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

- Problog

$$w : \text{smokes}(X) \leftarrow \exists Y \text{ friends}(X, Y) \wedge \text{smokes}(Y)$$

- Probability of smokes goes up as the number of friends increases!

# Cyclic Models

*Whether people smoke depends on whether their friends smoke.*

- MLN:

$$w : \text{smokes}(X) \leftarrow \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

(where  $\leftarrow$  is material implication) is equivalent to

$$w : \text{true}(X) \wedge \text{true}(Y)$$

$$-w : \neg \text{smokes}(X) \wedge \text{friends}(X, Y) \wedge \text{smokes}(Y)$$

- Problog

$$w : \text{smokes}(X) \leftarrow \exists Y \text{ friends}(X, Y) \wedge \text{smokes}(Y)$$

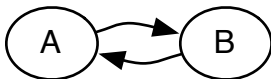
- Probability of smokes goes up as the number of friends increases!
- Problog cannot represent negative effects: someone is less likely to smoke if their friends smoke (without there being a non-zero probability of logical inconsistency)

# Cyclic Directed Models

- Make model acyclic, by totally ordering variables.  
Destroys exchangeability. Symmetries are not preserved.

# Cyclic Directed Models

- Make model acyclic, by totally ordering variables.  
Destroys exchangeability. Symmetries are not preserved.
- (Relational) dependency networks: directed model,



- $P(A, B)$  has 3 degrees of freedom,
- $P(A | B), P(B | A)$ , uses 4 numbers; typically inconsistent.
- resulting distribution means stationary (equilibrium) distribution of Markov chain.

# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 Unique properties of relational models
  - Learning general knowledge vs learning about a data set
  - Varying Populations
  - What can be observed?
- 4 Conclusions and Challenges



# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom.

# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom. Options:
  - learn general knowledge, e.g., transitivity, how male and female friendships work, how location affect friendship...

# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom. Options:
  - learn general knowledge, e.g., transitivity, how male and female friendships work, how location affect friendship...
  - learn specific knowledge about who is friends with who; e.g., which particular group of people are generally friends with each other.

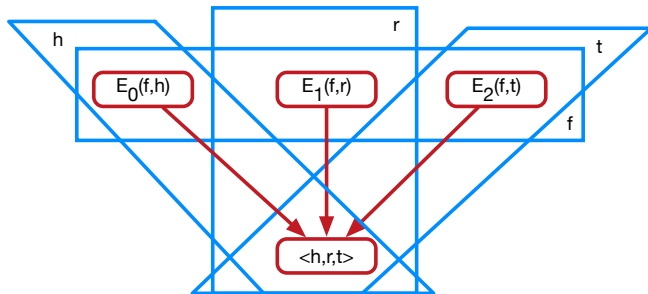
# Learning general knowledge vs learning about a data set

- Suppose you want to create a model of who is friends with whom. Options:
  - learn general knowledge, e.g., transitivity, how male and female friendships work, how location affect friendship...
  - learn specific knowledge about who is friends with who; e.g., which particular group of people are generally friends with each other.
- The specific knowledge will tend to be more accurate on that population, but doesn't generalize to different populations.
- The general knowledge will tend to transfer better.
- Which is better depends on the goals and how success is measured.

# Learning general knowledge vs learning about a data set

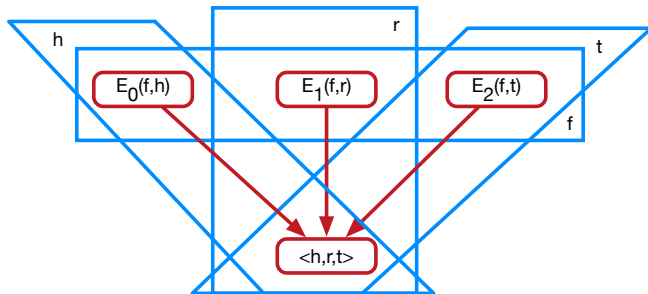
- Suppose you want to create a model of who is friends with whom. Options:
  - learn general knowledge, e.g., transitivity, how male and female friendships work, how location affect friendship...
  - learn specific knowledge about who is friends with who; e.g., which particular group of people are generally friends with each other.
- The specific knowledge will tend to be more accurate on that population, but doesn't generalize to different populations.
- The general knowledge will tend to transfer better.
- Which is better depends on the goals and how success is measured.
- Many of the methods try to do both; learn about specific individuals and general knowledge.

# Canonical polyadic tensor factorization model



$$P(\langle I, P, V \rangle) = \text{sigmoid}\left(\sum_d h_s(I, d)h_v(P, d)h_o(V, d)\right)$$

# Canonical polyadic tensor factorization model



$$P(\langle I, P, V \rangle) = \text{sigmoid}\left(\sum_d h_s(I, d)h_v(P, d)h_o(V, d)\right)$$

- Doesn't learn generalized knowledge but about the particular individuals
- There are many more and less sophisticated models

# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 Unique properties of relational models
  - Learning general knowledge vs learning about a data set
  - **Varying Populations**
  - What can be observed?
- 4 Conclusions and Challenges



# What happens as Population size $n$ Changes: Simplest case

$$\alpha_0 : q$$

$$\alpha_1 : q \wedge \neg r(x)$$

$$\alpha_2 : q \wedge r(x)$$

$$\alpha_3 : r(x)$$

Weighted formulae define distribution:

$$P_{MLN}(q \mid n) = \text{sigmoid}( \alpha_0 + n \log(e^{\alpha_2} + e^{\alpha_1 - \alpha_3}) )$$

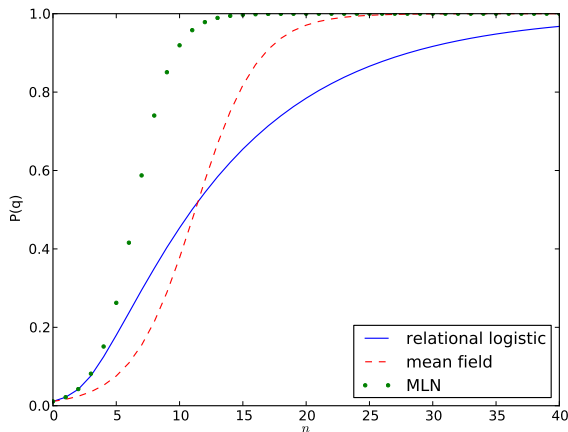
Weighted formulae define conditionals:

$$P_{RLR}(q \mid n) = \sum_{i=0}^n \binom{n}{i} \text{sigmoid}(\alpha_0 + i\alpha_1 + (n-i)\alpha_2) (1-p_r)^i p_r^{n-i}$$

Mean-field approximation:

$$P_{MF}(q \mid n) = \text{sigmoid}(\alpha_0 + np_r\alpha_1 + n(1-p_r)\alpha_2)$$

# Population Growth: $P(q | n)$



All:  $P(q | n) \rightarrow 0$  or  $1$  as  $n \rightarrow \infty$

# Challenges of varying populations

- **Example:** The Movielens 100k dataset contains data about  $rated(P, M, R, T)$  meaning person  $P$  gave movie  $M$  a rating of  $R$  at time  $T$ .  
Plus user demographic and movie information.
- Number of ratings per user is between 20 (arbitrary threshold) and 737; average of 106
- Number of ratings per movie is between 1 and 583; average of 59

# Challenges of varying populations

- **Example:** The Movielens 100k dataset contains data about  $rated(P, M, R, T)$  meaning person  $P$  gave movie  $M$  a rating of  $R$  at time  $T$ .  
Plus user demographic and movie information.
- Number of ratings per user is between 20 (arbitrary threshold) and 737; average of 106
- Number of ratings per movie is between 1 and 583; average of 59
- **Predicting age from ratings** is difficult:
  - With some high counts we can't assume movies produce independent evidence

# Challenges of varying populations

- **Example:** The Movielens 100k dataset contains data about  $rated(P, M, R, T)$  meaning person  $P$  gave movie  $M$  a rating of  $R$  at time  $T$ .  
Plus user demographic and movie information.
- Number of ratings per user is between 20 (arbitrary threshold) and 737; average of 106
- Number of ratings per movie is between 1 and 583; average of 59
- **Predicting age from ratings** is difficult:
  - With some high counts we can't assume movies produce independent evidence
  - With many low counts we need to regularize (and can't measure dependence)

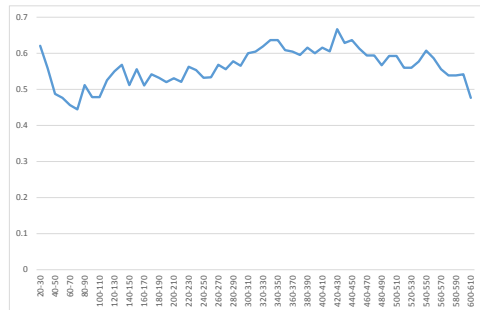
# Challenges of varying populations

- **Example:** The Movielens 100k dataset contains data about  $rated(P, M, R, T)$  meaning person  $P$  gave movie  $M$  a rating of  $R$  at time  $T$ .  
Plus user demographic and movie information.
- Number of ratings per user is between 20 (arbitrary threshold) and 737; average of 106
- Number of ratings per movie is between 1 and 583; average of 59
- **Predicting age from ratings** is difficult:
  - With some high counts we can't assume movies produce independent evidence
  - With many low counts we need to regularize (and can't measure dependence)
  - There is 100 times as much ratings information as age information

# Challenges of varying populations

- **Example:** The Movielens 100k dataset contains data about  $rated(P, M, R, T)$  meaning person  $P$  gave movie  $M$  a rating of  $R$  at time  $T$ .  
Plus user demographic and movie information.
- Number of ratings per user is between 20 (arbitrary threshold) and 737; average of 106
- Number of ratings per movie is between 1 and 583; average of 59
- **Predicting age from ratings** is difficult:
  - With some high counts we can't assume movies produce independent evidence
  - With many low counts we need to regularize (and can't measure dependence)
  - There is 100 times as much ratings information as age information
- Bigger datasets have even more variability.

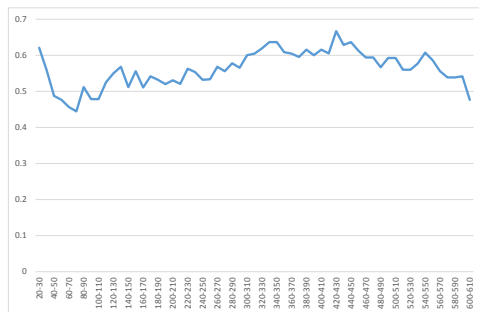
# Real Data



Observed  $P(25 < \text{Age}(u) < 45 \mid n)$ , where  $n$  is number of movies watched from the Movielens dataset.



# Real Data



Observed  $P(25 < \text{Age}(u) < 45 \mid n)$ , where  $n$  is number of movies watched from the Movielens dataset.

Dont use:

$$w : \text{middle\_age}(U) \leftarrow \text{rated}(U, M) \wedge \text{foo}(M)$$

then  $P(\text{middle\_age}(U)) \rightarrow 0$  or  $1$  as number of movies increases.

# Example of polynomial dependence of population

$$\alpha_0 : q$$

$$\alpha_1 : q \wedge \text{true}(X)$$

$$\alpha_2 : q \wedge r(X)$$

$$\alpha_3 : \text{true}(X)$$

$$\alpha_4 : r(X)$$

$$\alpha_5 : q \wedge \text{true}(X) \wedge \text{true}(Y)$$

$$\alpha_6 : q \wedge r(X) \wedge \text{true}(Y)$$

$$\alpha_7 : q \wedge r(X) \wedge r(Y)$$

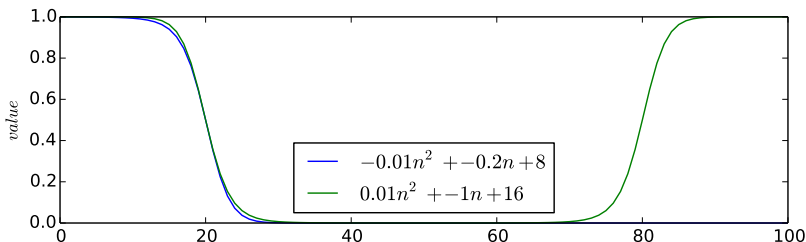
In RLR and in MLN, if all  $r(a_i)$  are observed:

$$P(q \mid \text{obs}) = \text{sigmoid}(\alpha_0 + n\alpha_1 + n_1\alpha_2 + n^2\alpha_5 + n_1n\alpha_6 + n_1^2\alpha_7)$$

$r(X)$  is true for  $n_1$  individuals out of a population of  $n$ .

# Danger of fitting to data without understanding the model

- RLR can fit sigmoid of any polynomial.
- Consider a polynomial of degree 2:



# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 **Unique properties of relational models**
  - Learning general knowledge vs learning about a data set
  - Varying Populations
  - **What can be observed?**
- 4 Conclusions and Challenges

# Observation Protocols

- Example: we want to predict the probability that Sam will like a apartment.
- Observation: there is a pink bedroom.

# Observation Protocols

- Example: we want to predict the probability that Sam will like a apartment.
- Observation: there is a pink bedroom.

The protocol of how the observation was obtained matters:

- Was it a pink sensor that reported what was pink?

# Observation Protocols

- Example: we want to predict the probability that Sam will like a apartment.
- Observation: there is a pink bedroom.

The protocol of how the observation was obtained matters:

- Was it a pink sensor that reported what was pink?
- Was it reporting the first thing it observed?

# Observation Protocols

- Example: we want to predict the probability that Sam will like a apartment.
- Observation: there is a pink bedroom.

The protocol of how the observation was obtained matters:

- Was it a pink sensor that reported what was pink?
- Was it reporting the first thing it observed?
- Was it reporting the most unusual feature of the apartment?



# Observation Protocols

- Example: we want to predict the probability that Sam will like a apartment.
- Observation: there is a pink bedroom.

The protocol of how the observation was obtained matters:

- Was it a pink sensor that reported what was pink?
- Was it reporting the first thing it observed?
- Was it reporting the most unusual feature of the apartment?
- Was it telling us the most positive aspect of the apartment?

# Observation Protocols

- Example: we want to predict the probability that Sam will like a apartment.
- Observation: there is a pink bedroom.

The protocol of how the observation was obtained matters:

- Was it a pink sensor that reported what was pink?
- Was it reporting the first thing it observed?
- Was it reporting the most unusual feature of the apartment?
- Was it telling us the most positive aspect of the apartment?
- Did it know that Sam was interested in whether there was a pink bedroom?

# Observation Protocols

- Example: we want to predict the probability that Sam will like a apartment.
- Observation: there is a pink bedroom.

## The protocol of how the observation was obtained matters:

- Was it a pink sensor that reported what was pink?
- Was it reporting the first thing it observed?
- Was it reporting the most unusual feature of the apartment?
- Was it telling us the most positive aspect of the apartment?
- Did it know that Sam was interested in whether there was a pink bedroom?
- Was it reporting the colour of each bedroom?

# Observation Protocols

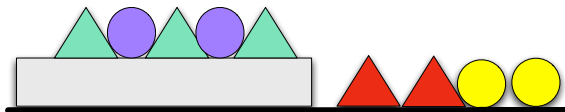
- Example: we want to predict the probability that Sam will like a apartment.
- Observation: there is a pink bedroom.

The protocol of how the observation was obtained matters:

- Was it a pink sensor that reported what was pink?
- Was it reporting the first thing it observed?
- Was it reporting the most unusual feature of the apartment?
- Was it telling us the most positive aspect of the apartment?
- Did it know that Sam was interested in whether there was a pink bedroom?
- Was it reporting the colour of each bedroom?

There are unboundedly many possible relations in a real-world object such as a house.

# Observation Protocols

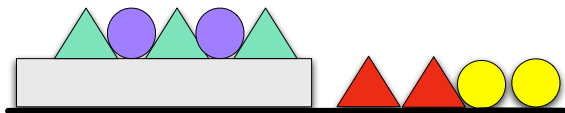


Observe a triangle and a circle touching. What is the probability the triangle is green?

$$P(\text{green}(x) \mid \text{triangle}(x) \wedge \exists y \text{ circle}(y) \wedge \text{touching}(x, y))$$

The answer depends on how the  $x$  and  $y$  were chosen!

# Protocol for Observing



$$P(\text{green}(x) \mid \text{triangle}(x) \wedge \exists y \text{ circle}(y) \wedge \text{touching}(x, y))$$

$$\begin{array}{c} | \\ \text{select}(x) \end{array}$$

$$\begin{array}{c} | \\ \text{select}(y) \end{array}$$

$$\begin{array}{c} | \\ 3/4 \end{array}$$

$$\begin{array}{c} | \\ \text{select}(y) \end{array}$$

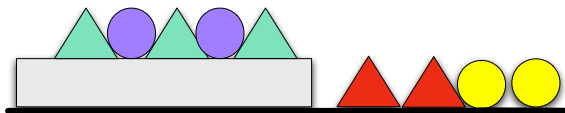
$$\begin{array}{c} | \\ \text{select}(x) \end{array}$$

$$\begin{array}{c} | \\ 2/3 \end{array}$$

$$\begin{array}{c} | \\ \text{select}(x, y) \end{array}$$

$$\begin{array}{c} | \\ 4/5 \end{array}$$

# Protocol for Observing



$$P(\text{green}(x) \mid \text{triangle}(x) \wedge \exists y \text{ circle}(y) \wedge \text{touching}(x, y))$$

$$\begin{array}{c} | \\ \text{select}(x) \end{array}$$

$$\begin{array}{c} | \\ \text{select}(y) \end{array}$$

$$\begin{array}{c} | \\ 3/4 \end{array}$$

$$\begin{array}{c} | \\ \text{select}(y) \end{array}$$

$$\begin{array}{c} | \\ \text{select}(x) \end{array}$$

$$\begin{array}{c} | \\ 2/3 \end{array}$$

$$\begin{array}{c} | \\ \text{select}(x, y) \end{array}$$

$$\begin{array}{c} | \\ 4/5 \end{array}$$

A logical formula does not provide enough information to determine the probabilities.

# Data

Real data is messy!

- Multiple levels of abstraction
- Multiple levels of detail
- Sometimes observations are abstract and lifted  
e.g., “3 people out of 300 in the audience asked a question”.
- Uses the vocabulary from many ontologies
- Rich meta-data:
  - Who collected each datum? (identity and credentials)
  - Who transcribed the information?
  - What was the protocol used to collect the data? (Chosen at random or chosen because interesting?)
  - What were the controls — what was manipulated, when?
  - What sensors were used? What is their reliability and operating range?
  - What is the provenance of the data; what was done to it when?
- Errors, forgeries, . . .



# Outline

- 1 Knowledge Graphs
  - Tensor Factorization and Neural Network Models
- 2 Representation Issues
  - Desiderata
  - How do relational models relate to probabilistic graphical models
- 3 Unique properties of relational models
  - Learning general knowledge vs learning about a data set
  - Varying Populations
  - What can be observed?
- 4 Conclusions and Challenges

## Take Home

- Exchangeability and dependence on population size distinguish relational models from non-relational models

# Take Home

- Exchangeability and dependence on population size distinguish relational models from non-relational models
- Relational models are different from normal graphical models

# Take Home

- Exchangeability and dependence on population size distinguish relational models from non-relational models
- Relational models are different from normal graphical models

Challenges that relational models tackle:

- Heterogeneity: information about individuals varies greatly in kind and amount

# Take Home

- Exchangeability and dependence on population size distinguish relational models from non-relational models
- Relational models are different from normal graphical models

Challenges that relational models tackle:

- Heterogeneity: information about individuals varies greatly in kind and amount
- Representations should
  - let people state their prior knowledge,
  - let them understand what they stated, and
  - let them understand the posterior models (given evidence).

# Take Home

- Exchangeability and dependence on population size distinguish relational models from non-relational models
- Relational models are different from normal graphical models

Challenges that relational models tackle:

- Heterogeneity: information about individuals varies greatly in kind and amount
- Representations should
  - let people state their prior knowledge,
  - let them understand what they stated, and
  - let them understand the posterior models (given evidence).
- Learn general knowledge as well as about particular individuals

# Take Home

- Exchangeability and dependence on population size distinguish relational models from non-relational models
- Relational models are different from normal graphical models

Challenges that relational models tackle:

- Heterogeneity: information about individuals varies greatly in kind and amount
- Representations should
  - let people state their prior knowledge,
  - let them understand what they stated, and
  - let them understand the posterior models (given evidence).
- Learn general knowledge as well as about particular individuals
- Use the meta-data of how data was collected
- Model protocol used to generate the observations
- Also model what is not observed (e.g., because it was redundant information, unimportant, false or unknown)