

Statistical Relation Learning for 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

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

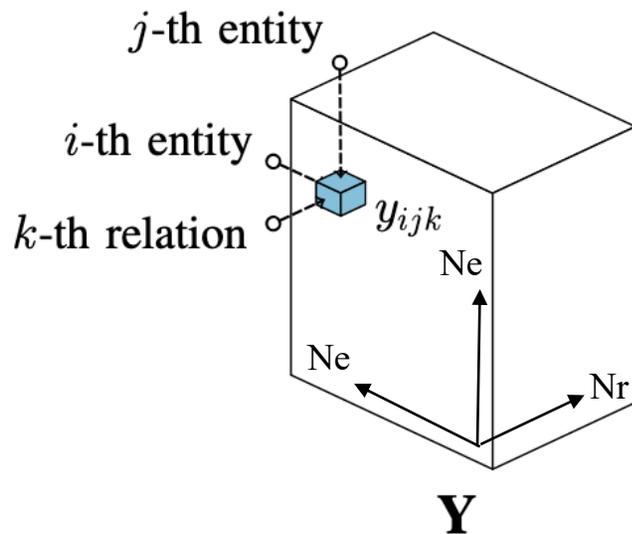
- Let $E = \{e_1, e_2, \dots, e_{N_e}\}$ be the set of all entities
- $R = \{r_1, r_2, \dots, r_{N_r}\}$ the set of all relations in a KG
 - Each possible triple is defined as $X_{ijk} = (e_i, r_k, e_j)$
 - over set of entities E and relations R
 - as a binary r.v. $y_{ijk} \in \{0,1\}$ that indicates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

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

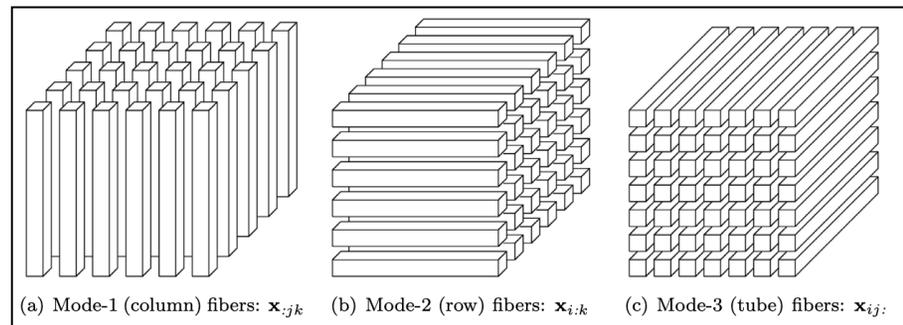
 - Each possible realization of Y can be interpreted as a possible world.

Tensor Representation

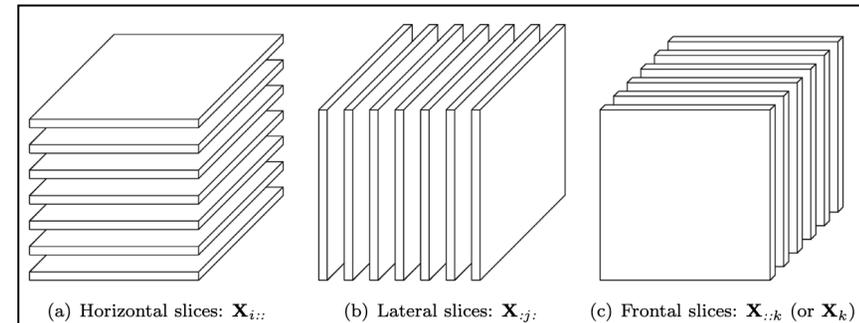
- Binary relational data is represented as a tensor



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 $ExRxE$ 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

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

2. Other tensor factorization models

3. Matrix factorization methods

4. Multi-layer perceptrons

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

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

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

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

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

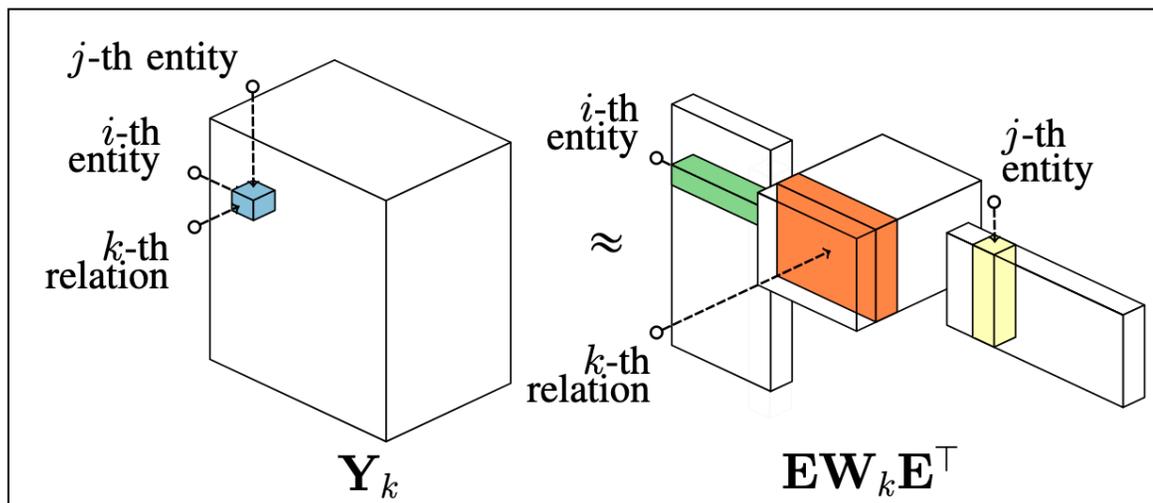
$$f_{ijk}^{\text{RESCAL}} := \mathbf{e}_i^\top \mathbf{W}_k \mathbf{e}_j = \sum_{a=1}^{H_e} \sum_{b=1}^{H_e} w_{abk} e_{ia} e_{jb}$$

- where $\mathbf{W}_k \in \mathbb{R}^{H_e \times H_e}$ 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

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

can be written compactly as

$$F_k = \mathbf{E} \mathbf{W}_k \mathbf{E}^T$$

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 $\mathbf{E} \in R^{Ne \times He}$ 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

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

– 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

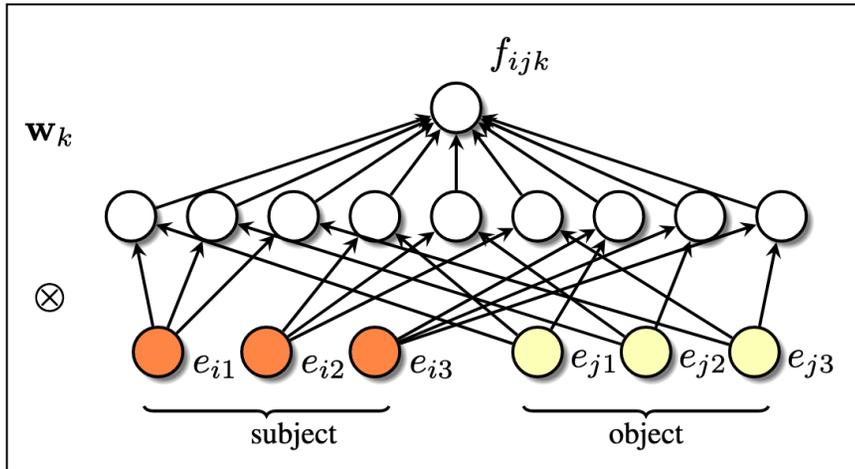
MLP based on *entities and relationships* (ER)

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

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

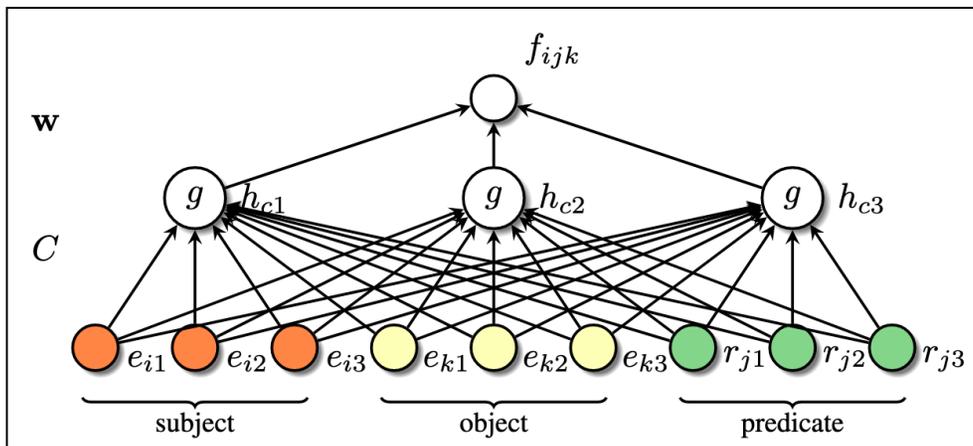
RESCAL and ER-MLP as Neural Nets



$$f_{ijk}^{\text{E-MLP}} := \mathbf{w}_k^\top \mathbf{g}(\mathbf{h}_{ijk}^a)$$

$$\mathbf{h}_{ijk}^a := \mathbf{A}_k^\top \phi_{ij}^{\text{E-MLP}}$$

$$\phi_{ij}^{\text{E-MLP}} := [\mathbf{e}_i; \mathbf{e}_j]$$



$$f_{ijk}^{\text{ER-MLP}} := \mathbf{w}^\top \mathbf{g}(\mathbf{h}_{ijk}^c)$$

$$\mathbf{h}_{ijk}^c := \mathbf{C}^\top \phi_{ijk}^{\text{ER-MLP}}$$

$$\phi_{ijk}^{\text{ER-MLP}} := [\mathbf{e}_i; \mathbf{e}_j; \mathbf{r}_k].$$

$H_e=H_r=3$ 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

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

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 → Anne ← 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

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