

OBJECTIVES

Regularize the Scene Graph learning deep network via multi-relational tensor factorization for robust visual relationship learning.



PROBLEM: VISUAL RELATIONSHIP DETECTION

► Given:

- 1. Image with **Object** *i* and **Object** *j* of interest
- 2. Bounding boxes and features (i.e., Faster R-CNN)
- 3. **Object labels** (ground truths or detection results)
- ► Goal:

Predict the visual relationship of the objects: [Object *i*, ?, Object *j*] corresponding k^{th} predicate.









Figure: (a): A relationship instance in a training set. (b): An unknown relationship to predict. (c): The interactions of the objects (i.e., motorcycle and horse are both 'ridable') can be used to infer the correct relationship.

CHALLENGES IN SEMANTIC INFERENCE

Large observation space:

> N object categories & M possible predicates $\Rightarrow N^2 M$ possible combinations

Sparse observations:

Visual Genome has 1M relationship instances \Rightarrow But observed only \sim 2% of possible combinations

Zero-shot learning:

Inferring cases unobserved in train set

$\Rightarrow \sim$ 98% of possible cases *NOT* observed in train set

Zero-shot Learning: Unobserved Visual Relationship Detection







ee, behind, bear mouse, on, cabinet horse, wear, hat Figure: The *unobserved* observed relationships are potentially much harder to detect. Conference on Computer Vision and Pattern Recognition (CVPR) 2018





Figure: Each 'slice' X_k encodes possible relationships involving its





MULTI-RELATIONAL TENSOR FACTORIZATION



Tensorize, Factorize and Regularize: Robust Visual Relationship Learning

Seong Jae Hwang Sathya N. Ravi Zirui Tao Hyunwoo J. Kim Maxwell D. Collins Vikas Singh

STEP 1: TENSORIZE THE VISUAL RELATIONSHIPS

► Multi-relational tensor $X \in \mathbb{R}^{n \times n \times m}$ given *n* object categories and *m* possible predicates

 \blacktriangleright X(i, j, k): number of [Object i, Predicate k, Object j] in train set $ho \sim 2\%$ is non-zero \Rightarrow extremely sparse tensor

STEP 2: FACTORIZE THE RELATIONAL TENSOR

Multi-relational Tensor Factorization

 \blacktriangleright Based on the multi-relational tensor X from train set, derive . **Common** latent representation of objects $A \in \mathbb{R}^{n \times r}$ 2. **Relationship-specific** factor matrix $R_k \in \mathbb{R}^{r \times r}$ for each rela-

tionship $k \in \{1, \ldots, m\}$

such that $X_k \approx AR_k A^T$

Figure: Multi-relational tensor factorization $X_k \approx AR_k A^T$ for $k \in \{1, \ldots, m\}$.

$$\min_{A,R_k} \sum_{k=1}^{m} ||X_k - AR_k A^T||_F^2$$

. 4th-order term A: Use auxiliary variables to decouple A and A^{T}

$$\min_{A,R_k} \sum_{k=1}^{m} ||X_k - B_k A^T||_F^2 \quad \text{s.t.} \quad B_k = AR_k.$$

2. Low-rank Initialization via SVD: For m = 1, $U\Sigma V^T = X$. Initialize with the "basin of attraction" (Luo et al., Tu et al.):

$$\mathbf{A} = V \Sigma^{1/2}, \quad B = U \Sigma^{1/2}.$$

3. **Restrict degenerate cases**: $A' = AP^{-T}$ and $B' = BP^{-T}$ for any invertible P. Normalize A and B_k :

$$\lambda_{p} \sum_{k=1}^{m} ||A^{T}A - B_{k}^{T}B_{k}||_{F}^{2}.$$
(3)

$$|A_k - B_k A^T||_F^2 + \gamma \sum_{k=1}^m ||B_k - AR_k||_F^2 + \lambda_p \sum_{k=1}^m ||A^T A - B_k^T B_k||_F^2$$
 (4)

- 3: $\overline{X} \leftarrow \sum_{k=1}^m X_k$

- $R_k \leftarrow (A^T A)^{-1} (A^T B_k)$
- end while



bounding boxes (red/blue).



(a) Input image b) Scene Graph Figure: Scene graph detection tasks. Check marks indicate required prediction components. The tasks become incrementally more demanding from top (PredCls) to bottom (SgGen).

ACKNOWLEDGMENT

SJH was supported by a University of Wisconsin CIBM fellowship (5T15LM007359-14). We acknowledge support from NIH R01 AG040396 (VS), NSF CCF 1320755 (VS), NSF CAREER award 1252725 (VS), UW ADRC AG033514, UW ICTR 1UL1RR025011, Waisman Core grant P30 HD003352-45 and UW CPCP AI117924 (SNR).

http://pages.cs.wisc.edu/~sjh



$$k^* = k^*_{SG} \odot D(k^*_{RL}, \theta) \tag{5}$$