

Deep Learning - Embedding

深度学习 - 嵌入

Jian Li

IIS, Tsinghua

Word2Vec

“You shall know a word by the company it keeps”

(J. R. Firth 1957: 11)

Build a cooccurrence matrix (using a moving window)

- I like deep learning.
- I like NLP.
- I enjoy flying.

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

- high dimensional, and sparse

- we want relatively low-dim representation

Word Embedding

- Map each word to a vector (in relatively low-dim space)

- $X_{apple} - X_{apples} \approx X_{car} - X_{cars} \approx X_{family} - X_{families}$

- $X_{shirt} - X_{clothing} \approx X_{chair} - X_{furniture}$

- $X_{king} - X_{man} \approx X_{queen} - X_{woman}$

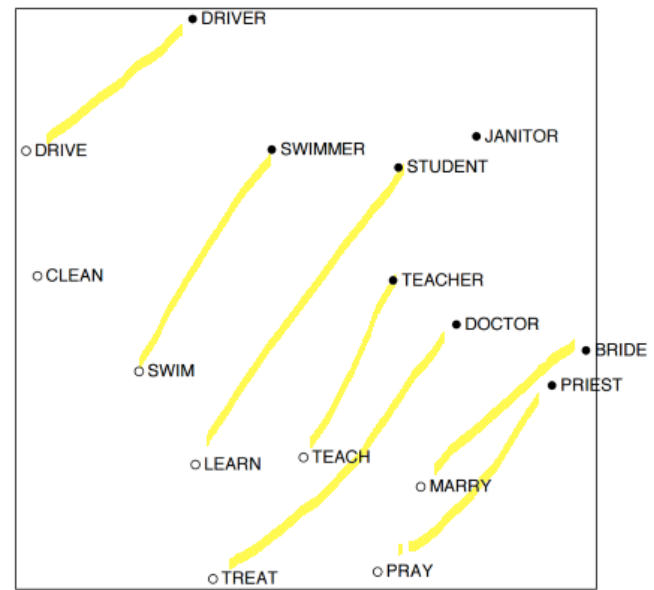
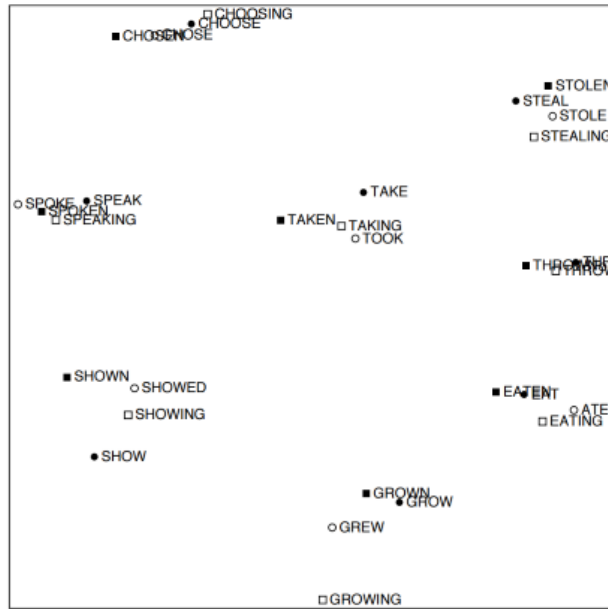
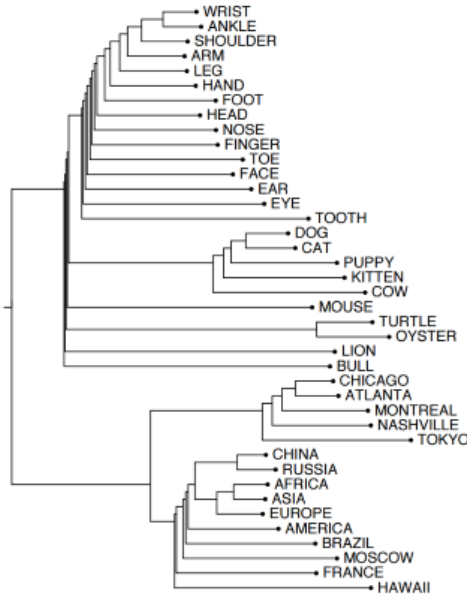
Approach 1 SVD/PCA

- Some notes: remove "of" "the" (Syntactic words)

- Use Pearson correlations (rather than count)

- set neg values to zero

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - E_X)(Y - E_Y)]}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$



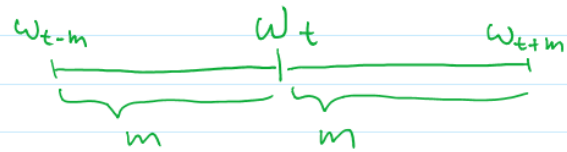
problem: SVD expensive for large vocabulary

special algorithm (doesn't fit into DL pipeline)

Word2Vec

Word 2 Vec (Mikolov et al. 2013)

Predict surrounding words in a window (of len m) of every word



First define a probabilistic model

each word w has two vectors $\begin{cases} \mu_w & \text{(the embedding of } w \text{ when } w \text{ is outside)} \\ \nu_w & \text{(center)} \end{cases}$

We want to learn μ_w, ν_w for $w \in W$

$$P(o|c) = \frac{\exp(\langle \mu_o, \nu_c \rangle)}{\sum_{w=1}^W \exp(\langle \mu_w, \nu_c \rangle)} \leftarrow \text{softmax / logistic regression}$$

↑ outside word o ↑ center word c



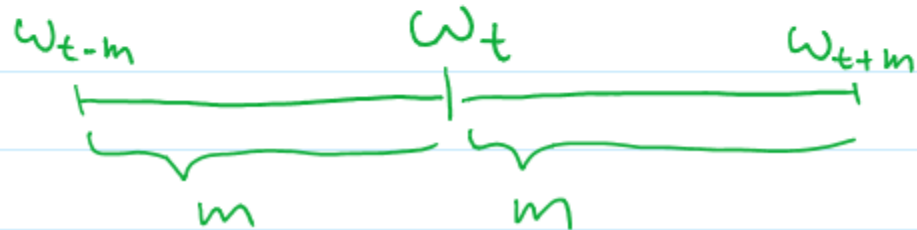
if μ_o, ν_c in the same direction. $P(o|c)$ is large.

Objective of Word2Vec

Objective:
$$J = \frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(\omega_{t+j} | \omega_t)$$

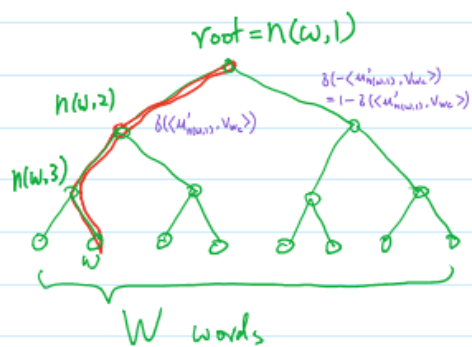
Sum
over
all word

log Prob of surrounding word of ω_t



training: $\nabla \log P(w_o | w_c)$ is expensive ($O(W)$ time)

① Hierarchical Softmax (it is an approximation)



$n(w,j)$: j th node on the root \rightarrow w path

$L(w)$: length of root \rightarrow w path

for every inner node, we also have a vector $u'_{n(w,j)}$

for every leaf word w , only one vector v_w

W words

± 1 indicator

left child of $n(w,j)$

$$P(w | w_c) = \prod_{j=1}^{L(w)-1} \delta \left(\mathbb{1} \left\{ n(w, j+1) = \text{left child of } n(w, j) \right\} \langle u'_{n(w, j)}, v_{w_c} \rangle \right)$$

$$\delta(x) = \frac{1}{1 + e^{-x}}$$

$$\sum_{w=1}^W P(w | w_c) = 1 \quad (\text{easy to check})$$

- not necessarily a full binary tree. e.g. the paper uses Huffman tree

(frequent word \rightarrow shorter path)

- time. len of the path.

Training method 2

Negative Sampling.

Noise Contrastive Estimation (NCE) (Gutmann & Hyvarinen)

a good model should be able to differentiate data from noise (by logistic regression)

Negative sampling objective

($\sum \log \delta(\cdot)$ vs $\log \sum \delta(\cdot)$)

replace by $P(w_o|w_c)$ by $\log \delta(\langle u_{w_o}, v_{w_c} \rangle) + \sum_{i=1}^k \log \delta(-\langle u_{w_i}, v_{w_c} \rangle)$

$w_i \sim p(w) \leftarrow$ noise distributions $\left\{ \begin{array}{l} \text{uniform distr} \\ \text{Unigram distr } U(w) \text{ (i.e., word frequency)} \\ U(w)^{\frac{3}{4}} / Z \end{array} \right.$

($K \sim 5-20$ typically. if data is large. $K \sim 2-5$)

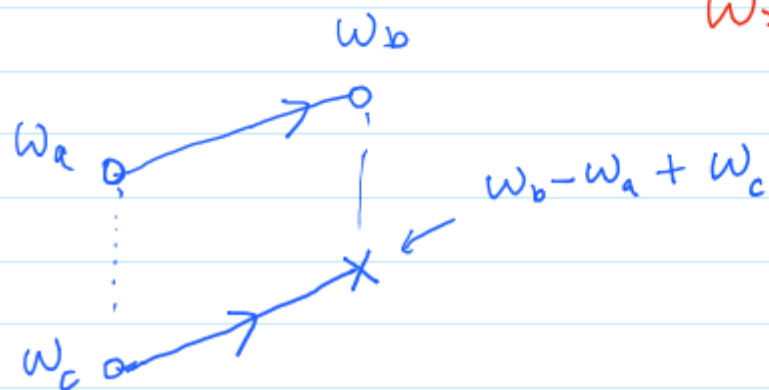
- $X_{apple} - X_{apples} \approx X_{car} - X_{cars} \approx X_{family} - X_{families}$
- $X_{shirt} - X_{clothing} \approx X_{chair} - X_{furniture}$
- $X_{king} - X_{man} \approx X_{queen} - X_{woman}$

a:b :: c:?



$$d = \arg \max_x \frac{(w_b - w_a + w_c)^T w_x}{\|w_b - w_a + w_c\|}$$

w_x produces the largest $\langle \nearrow, w_x \rangle$
 direction of $w_b - w_a + w_c$



Notes

- Word2Vec: unsupervised learning
 - Huge amount of training data (no label is needed)
- Can be incorporated to deep learning pipeline
 - The corresponding layers is usually called the **embedding layer**
 - The resulting vectors obtained from word2vec can be used to initialize the parameters of NN
 - We can also get the embedding from training a specified DNN (for a specific task)
 - Computational complexity too high (much higher than word2vec)
 - The embedding may not be useful in other tasks
 - On the other hand, word2vec captures a lot of semantic information, which is useful in a variety of tasks

According to Mikolov:

CBOW (Continuous Bag of Words): Use context to predict the current word.

--several times faster to train than the skip-gram, slightly better accuracy for the frequent words

Skip-gram: Use the current word to predict the context.

--works well with small amount of the training data, represents well even rare words or phrases.

GloVe

Two methodologies

① global matrix factorization

② local context window (skip-gram)

predicting the context given a word

Ratio Matters

Co-occurrence matrix $X = \left[\begin{array}{c} X_{ij} \end{array} \right]$ ← # times word j occurs in the context of word i

$$P_{ij} = P(j|i) = \frac{X_{ij}}{X_i} = \frac{X_{ij}}{\sum_j X_{ij}}$$

Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$	$k = \text{fashion}$
$P(k \text{ice})$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k \text{steam})$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k \text{ice})/P(k \text{steam})$	8.9	8.5×10^{-2}	1.36	0.96

↑
Solid is more relevant to ice than to steam

↑
gas is less relevant to ice than to steam

Derivation

$$\text{Define } F(\underbrace{w_i}_{\in \mathbb{R}^d}, \underbrace{w_j}_{\in \mathbb{R}^d}, \underbrace{\tilde{w}_k}_{\in \mathbb{R}^d}) = \frac{P_{ik}}{P_{jk}}$$

Using the property of F , we try to derive the form of $\langle w_i, \tilde{w}_k \rangle$

$$\text{Want } \textcircled{1} F(w_i, w_j, \tilde{w}_k) = F(\langle w_i - w_j, \tilde{w}_k \rangle)$$

enforce a linear structure

② Imagine we switch the role of a word & a context

$$\text{i.e. } X \rightarrow X^T$$

then we'd better have $w_i \rightarrow \tilde{w}_i$ (symmetric)

So we choose $F = \exp$

$$\left(\exp(\langle w_i - w_j, \tilde{w}_k \rangle) = \frac{\exp(\langle w_i, \tilde{w}_k \rangle)}{\exp(\langle w_j, \tilde{w}_k \rangle)} = \frac{P_{ik}}{P_{jk}} = \frac{X_{ik}/X_i}{X_{jk}/X_j} \right)$$

$$\Rightarrow \boxed{\langle w_i, \tilde{w}_k \rangle + b_i + \tilde{b}_k = \log(X_{ik})}$$

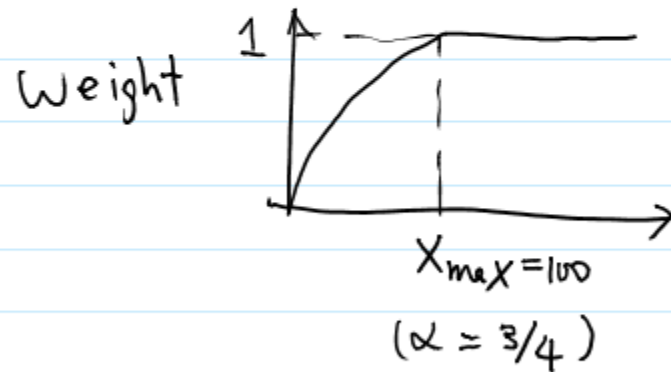
$$\text{(Verify } \langle w_i - w_j, \tilde{w}_k \rangle = \log\left(\frac{X_{ik}}{X_{jk}}\right) - (\overset{\log X_i}{b_i} - \overset{\log X_j}{b_j}))$$

We try to factorize $\log X$ (i.e. $w \begin{pmatrix} \tilde{w} \\ \tilde{b} \end{pmatrix} = \log \left(\begin{matrix} X \\ X \end{matrix} \right)$)

Objective

Objective (weighted least square)

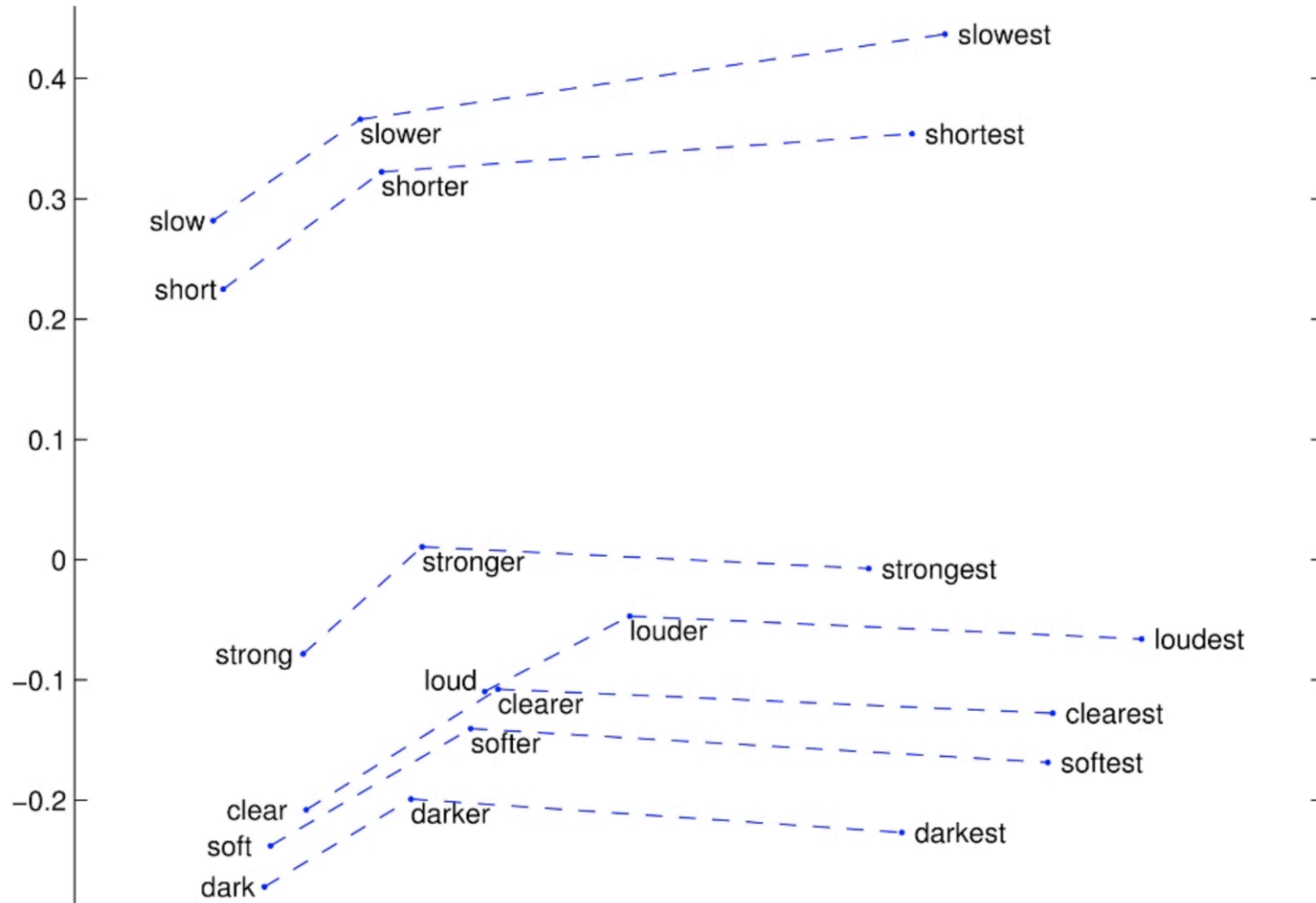
$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^W \underbrace{f(X_{ij})}_{\text{weight}} (\langle w_i, \tilde{w}_j \rangle + b_i + \tilde{b}_j - \log X_{ij})^2$$



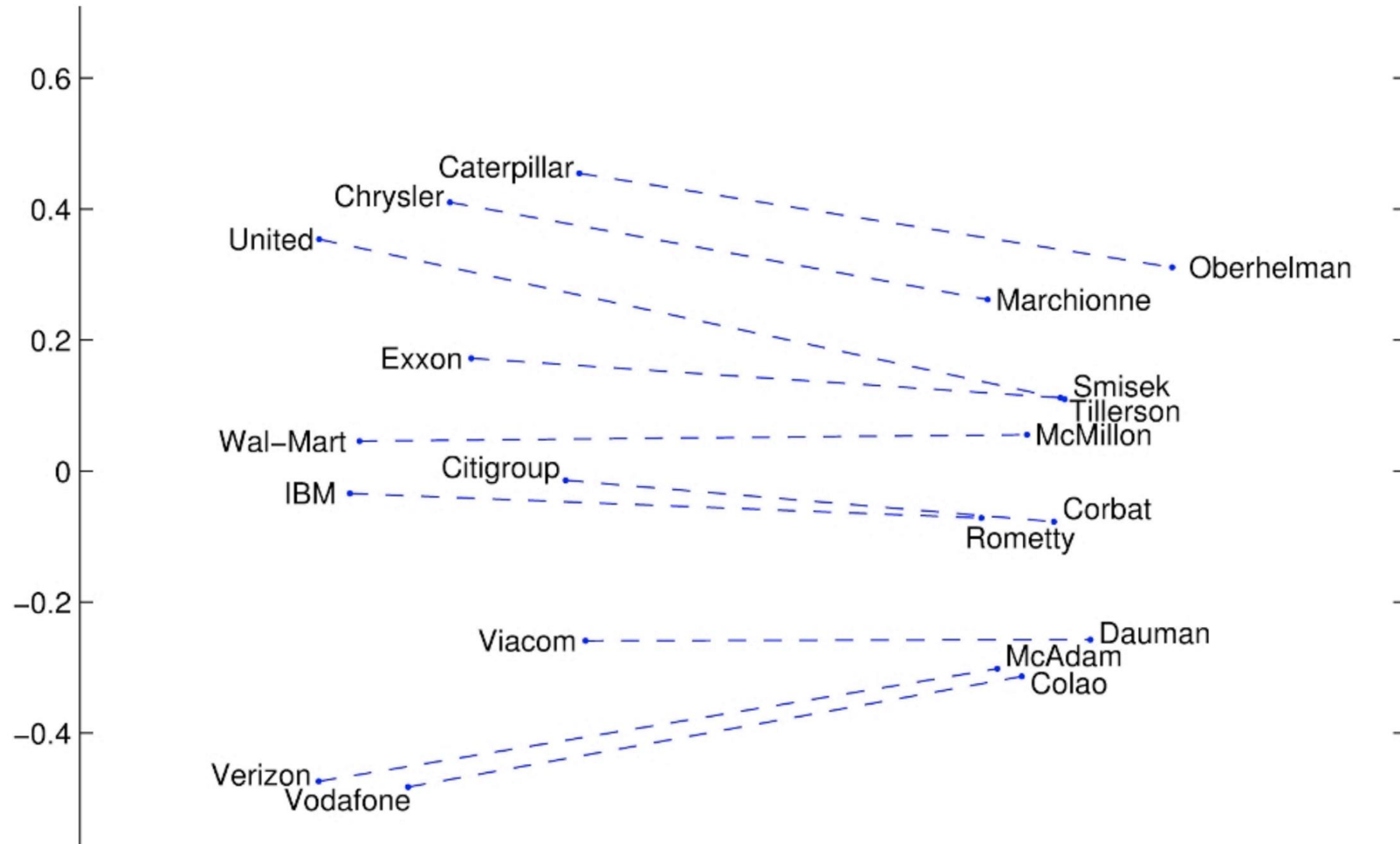
$$f(x) = \begin{cases} (x/x_{max})^\alpha & x < x_{max} \\ 1 & \text{o.w.} \end{cases}$$

Training: SGD, Adagrad

Glove Visualizations: Superlatives



Glove Visualizations: Company - CEO



Glove results

Nearest words to
frog:

1. frogs
2. toad
3. litoria
4. leptodactylidae
5. rana
6. lizard
7. eleutherodactylus



litoria



leptodactylidae



rana



eleutherodactylus

Reference

- Baroni, Marco, Georgiana Dinu, and Germán Kruszewski. Don't count, predict! A systematic comparison of context-counting vs. context-predicting semantic vectors. ACL 14
 - An early paper claims that the prediction formulation (like word2vec) is better than factorizing a co-occurrence matrix
- Levy, O., Goldberg, Y., & Dagan, I. (2015). Improving Distributional Similarity with Lessons Learned from Word Embeddings.

(comparing several SNSG, GloVe, and SVD)

A very reasonable blog discussing the relations between different models

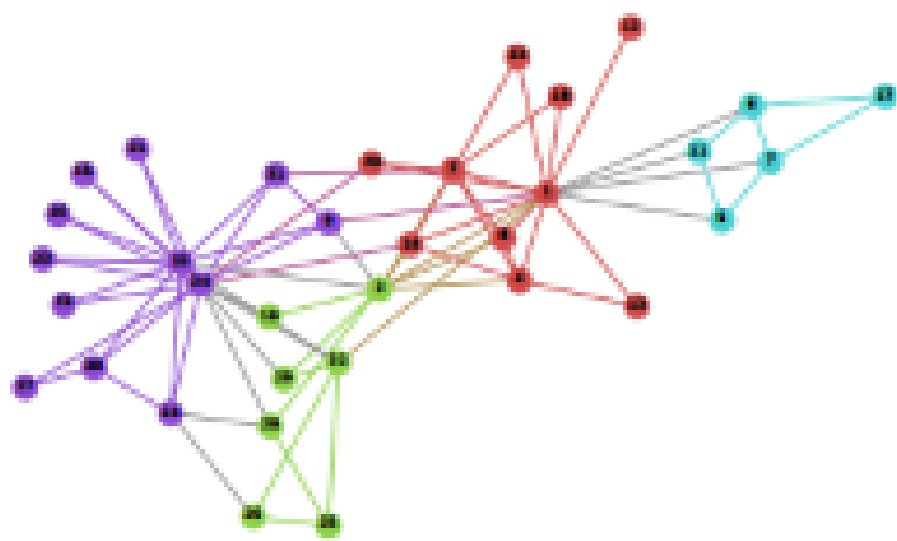
<http://sebastianruder.com/secret-word2vec/index.html>

- O Levy, Y Goldberg. Neural Word Embedding as Implicit Matrix Factorization. NIPS 14. (in blackboard)
- Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and node2vec. Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, Jie Tang. The 11th ACM International Conference on Web Search and Data Mining (WSDM 2018).

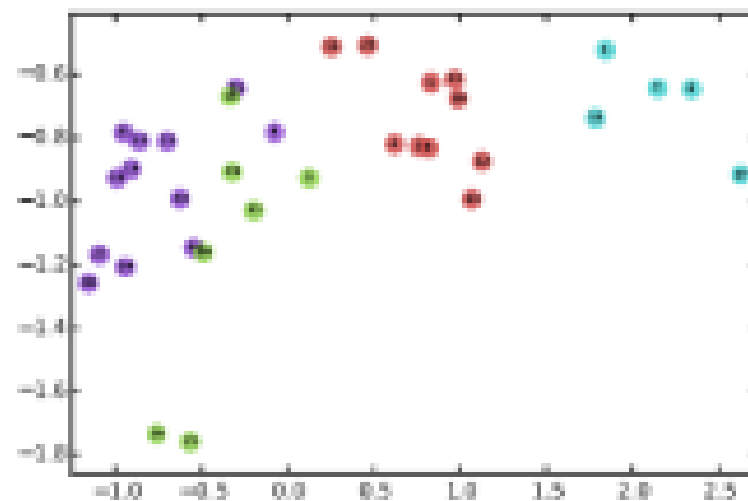
Deep Walk

-embedding node in a social network

Embedding node (using pairwise relations)



(a) Input: Karate Graph



(b) Output: Representation

Key Idea

treat vertex as words , random walks as sentences

Algorithm 1 DEEPWALK(G, w, d, γ, t)

Input: graph $G(V, E)$
window size w
embedding size d
walks per vertex γ
walk length t

Output: matrix of vertex representations $\Phi \in \mathbb{R}^{|V| \times d}$

1: Initialization: Sample Φ from $\mathcal{U}^{|V| \times d}$

2: Build a binary Tree T from V

3: **for** $i = 0$ to γ **do**

4: $\mathcal{O} = \text{Shuffle}(V)$

5: **for each** $v_i \in \mathcal{O}$ **do**

6: $\mathcal{W}_{v_i} = \text{RandomWalk}(G, v_i, t)$

7: $\text{SkipGram}(\Phi, \mathcal{W}_{v_i}, w)$

8: **end for**

9: **end for**

← all parameters

$$\Phi = \begin{matrix} | \\ \vdots \\ | \end{matrix} \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}$$

← random shuffle of all vertices. Speed up SGD

← Generate a random walk from v_i of length t

← SGD update step

Algorithm 2 SkipGram($\Phi, \mathcal{W}_{v_i}, w$)

1: **for each** $v_j \in \mathcal{W}_{v_i}$ **do**

2: **for each** $u_k \in \mathcal{W}_{v_i}[j - w : j + w]$ **do**

3: $J(\Phi) = -\log \Pr(u_k | \Phi(v_j))$

4: $\Phi = \Phi - \alpha * \frac{\partial J}{\partial \Phi}$

5: **end for**

6: **end for**

} → SGD

Prob of u_k (context word)
conditioning on the vector of
 v_j (center word)

Use Hierarchical softmax to approximate.

Multimodal representation learning

---Image Caption 2

Kires et al. Unifying Visual-Semantic Embeddings with
Multimodal Neural Language Models



there is a cat sitting on a shelf .



a plate with a fork and a piece of cake .



a black and white photo of a window .



a young boy standing on a parking lot next to cars .



a wooden table and chairs arranged in a room .



a kitchen with stainless steel appliances .



this is a herd of cattle out in the field .



a car is parked in the middle of nowhere .



a ferry boat on a marina with a group of people .



a little boy with a bunch of friends on the street .



a giraffe is standing next to a fence in a field .
(hallucination)



the two birds are trying to be seen in the water .
(counting)



a parked car while driving down the road .
(contradiction)



the handlebars are trying to ride a bike rack .
(nonsensical)



a woman and a bottle of wine in a garden .
(gender)

Figure 1: Sample generated captions. The bottom row shows different error cases. Additional results can be found at http://www.cs.toronto.edu/~rkiros/lstm_scn1m.html

Overview

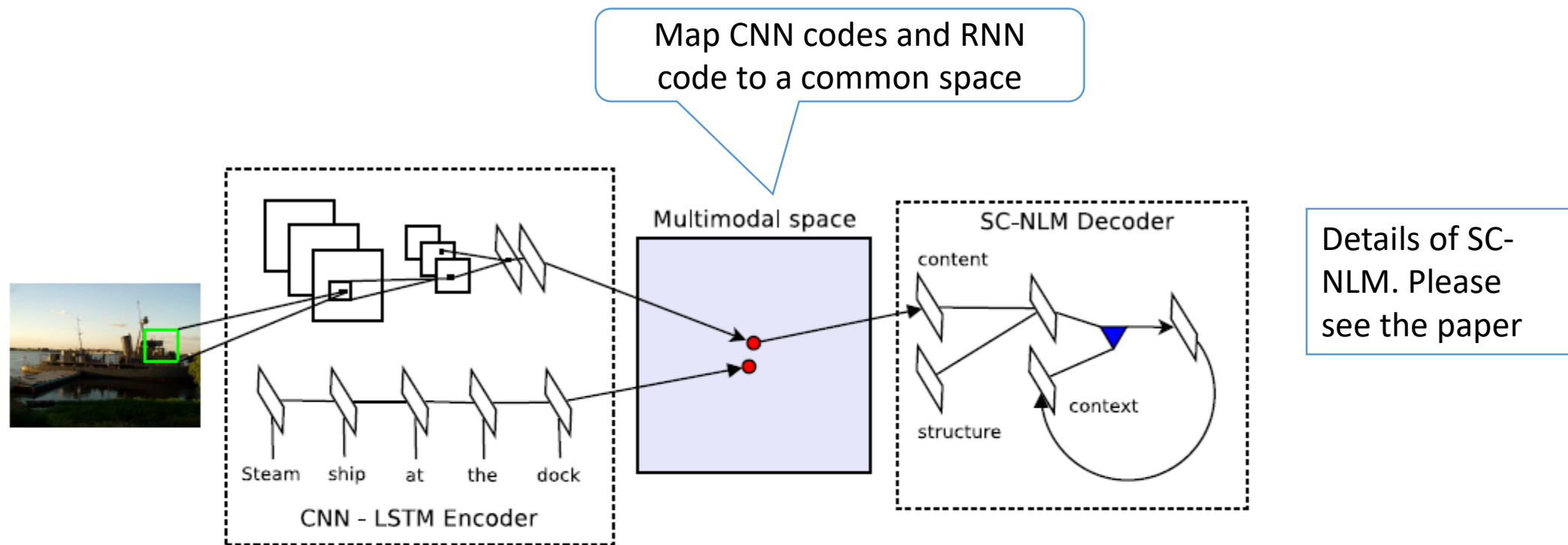


Figure 2: **Encoder:** A deep convolutional network (CNN) and long short-term memory recurrent network (LSTM) for learning a joint image-sentence embedding. **Decoder:** A new neural language model that combines structure and content vectors for generating words one at a time in sequence.

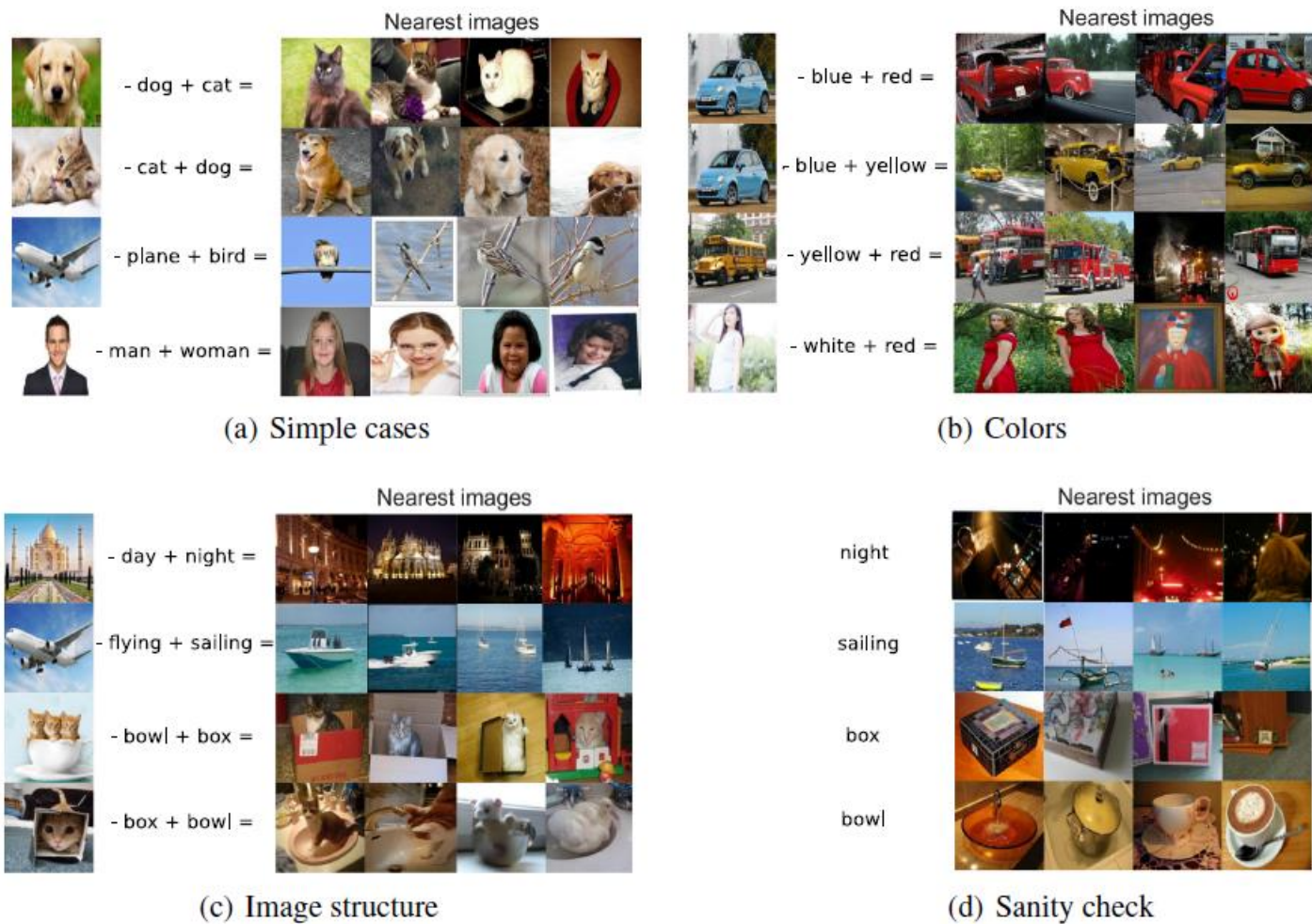


Figure 4: Multimodal vector space arithmetic. Query images were downloaded online and retrieved images are from the SBU dataset.

$$\begin{aligned}
 \mathbf{V}_{car} &\approx \mathbf{I}_{bcar} - \mathbf{V}_{blue} \\
 \mathbf{V}_{red} + \mathbf{V}_{car} &\approx \mathbf{I}_{bcar} - \mathbf{V}_{blue} + \mathbf{V}_{red} \\
 \mathbf{I}_{rcar} &\approx \mathbf{I}_{bcar} - \mathbf{V}_{blue} + \mathbf{V}_{red}
 \end{aligned}$$

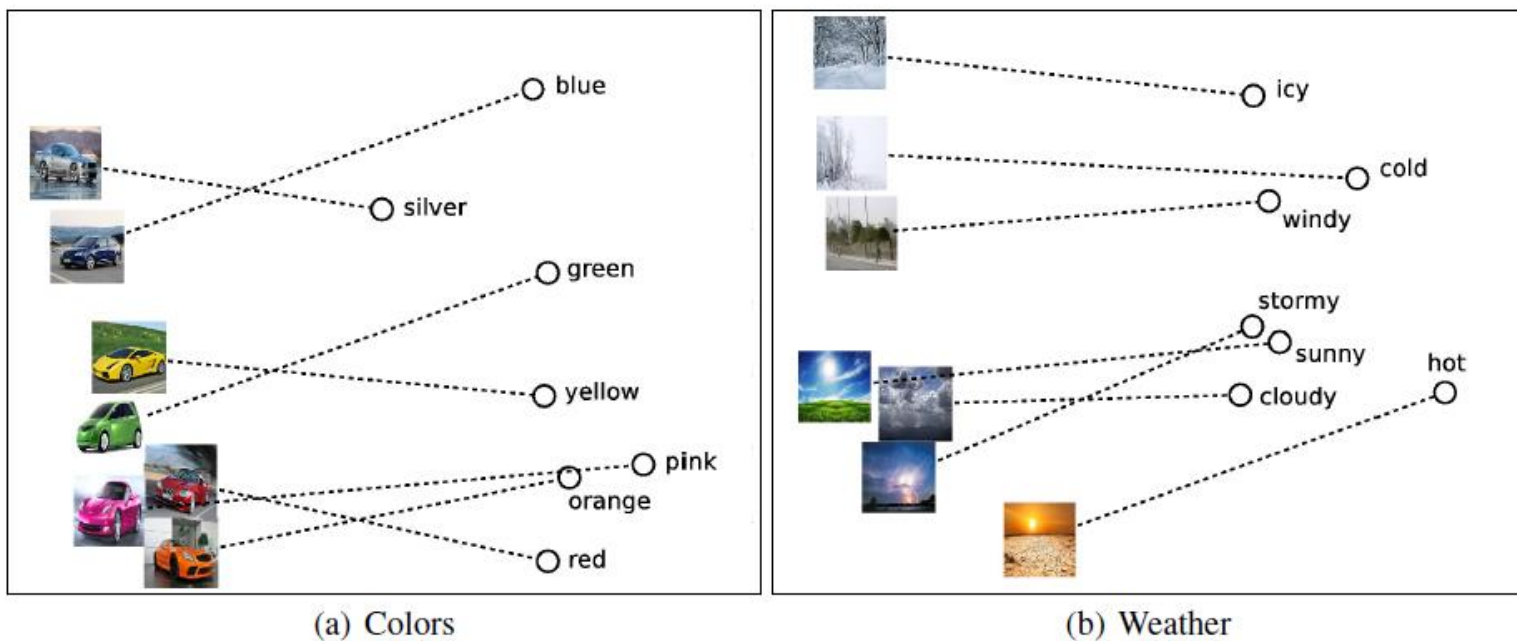


Figure 5: PCA projection of the 300-dimensional word and image representations for (a) cars and colors and (b) weather and temperature.

Details

- LSTM notations used in this work

Let \mathbf{X}_t denote a matrix of training instances at time t . In our case, \mathbf{X}_t is used to denote a matrix of word representations for the t -th word of each sentence in the training batch. Let $(\mathbf{I}_t, \mathbf{F}_t, \mathbf{C}_t, \mathbf{O}_t, \mathbf{M}_t)$ denote the input, forget, cell, output and hidden states of the LSTM at time step t . The LSTM architecture in this work is implemented using the following equations:

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \cdot \mathbf{W}_{xi} + \mathbf{M}_{t-1} \cdot \mathbf{W}_{hi} + \mathbf{C}_{t-1} \cdot \mathbf{W}_{ci} + \mathbf{b}_i) \quad (1)$$

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \cdot \mathbf{W}_{xf} + \mathbf{M}_{t-1} \cdot \mathbf{W}_{hf} + \mathbf{C}_{t-1} \cdot \mathbf{W}_{cf} + \mathbf{b}_f) \quad (2)$$

$$\mathbf{C}_t = \mathbf{F}_t \bullet \mathbf{C}_{t-1} + \mathbf{I}_t \bullet \tanh(\mathbf{X}_t \cdot \mathbf{W}_{xc} + \mathbf{M}_{t-1} \cdot \mathbf{W}_{hc} + \mathbf{b}_c) \quad (3)$$

$$\mathbf{O}_t = \sigma(\mathbf{X}_t \cdot \mathbf{W}_{xo} + \mathbf{M}_{t-1} \cdot \mathbf{W}_{ho} + \mathbf{C}_t \cdot \mathbf{W}_{co} + \mathbf{b}_o) \quad (4)$$

$$\mathbf{M}_t = \mathbf{O}_t \bullet \tanh(\mathbf{C}_t) \quad (5)$$

where (σ) denotes the sigmoid activation function, (\cdot) indicates matrix multiplication and (\bullet) indicates component-wise multiplication. 1

Details

Let $\mathbf{q} \in \mathbb{R}^D$ denote an image feature vector

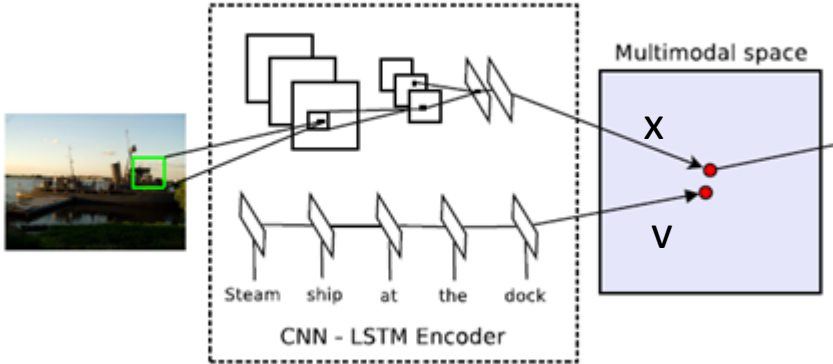
- D: length of the CNN code (CNN can be AlexNet, VggNet, or ResNet)

$\mathbf{x} = \mathbf{W}_I \cdot \mathbf{q} \in \mathbb{R}^K$ be the image embedding.

image description $S = \{w_1, \dots, w_N\}$ with words w_1, \dots, w_N

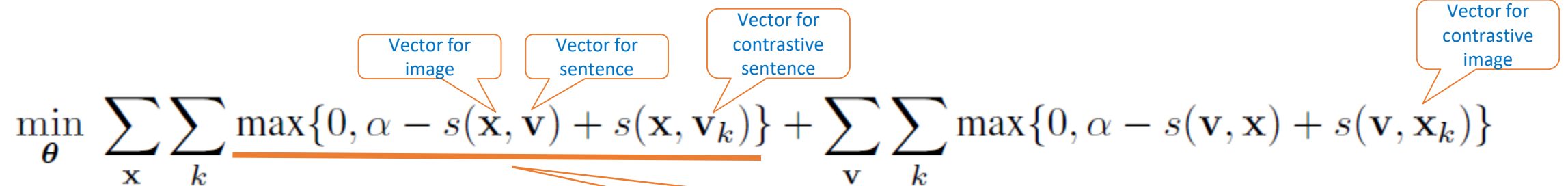
$\{\mathbf{w}_1, \dots, \mathbf{w}_N\}, \mathbf{w}_i \in \mathbb{R}^K, i = 1, \dots, n$ denote the corresponding word representations to words w_1, \dots, w_N (entries in the matrix \mathbf{W}_T). The representation of a sentence \mathbf{v} is the hidden state of the LSTM at time step N (i.e. the vector \mathbf{m}_t).

\mathbf{W}_T : precomputed using e.g. word2vec



Details

- Optimize pairwise rank loss (θ : parameters needed to be learnt: W_I and LSTM parameters) (similar to negative sampling in spirit)

$$\min_{\theta} \sum_{\mathbf{x}} \sum_k \max\{0, \alpha - s(\mathbf{x}, \mathbf{v}) + s(\mathbf{x}, \mathbf{v}_k)\} + \sum_{\mathbf{v}} \sum_k \max\{0, \alpha - s(\mathbf{v}, \mathbf{x}) + s(\mathbf{v}, \mathbf{x}_k)\}$$


Max-margin formulation. α margin
We have nonzero loss if
 $s(\mathbf{x}, \mathbf{v})$ is less than $s(\mathbf{x}, \mathbf{v}_k) + \alpha$

scoring function $s(\mathbf{x}, \mathbf{v}) = \mathbf{x} \cdot \mathbf{v}$,

\mathbf{v}_k is a contrastive (non-descriptive) sentence for image embedding \mathbf{x} , and vice-versa with \mathbf{x}_k .