# word2vec

Yan Song

# Outline

- Neural Representation Model for Words

- word2vec
  - Concept
  - Basics
  - Details with Implementation

- Evaluation
  - Word similarity measurement

# Neural Representation Model for Words

- Words are the most fundamental units in language
  - Basic semantics
- Can be used in a flexible way in many tasks
- Easy to learn

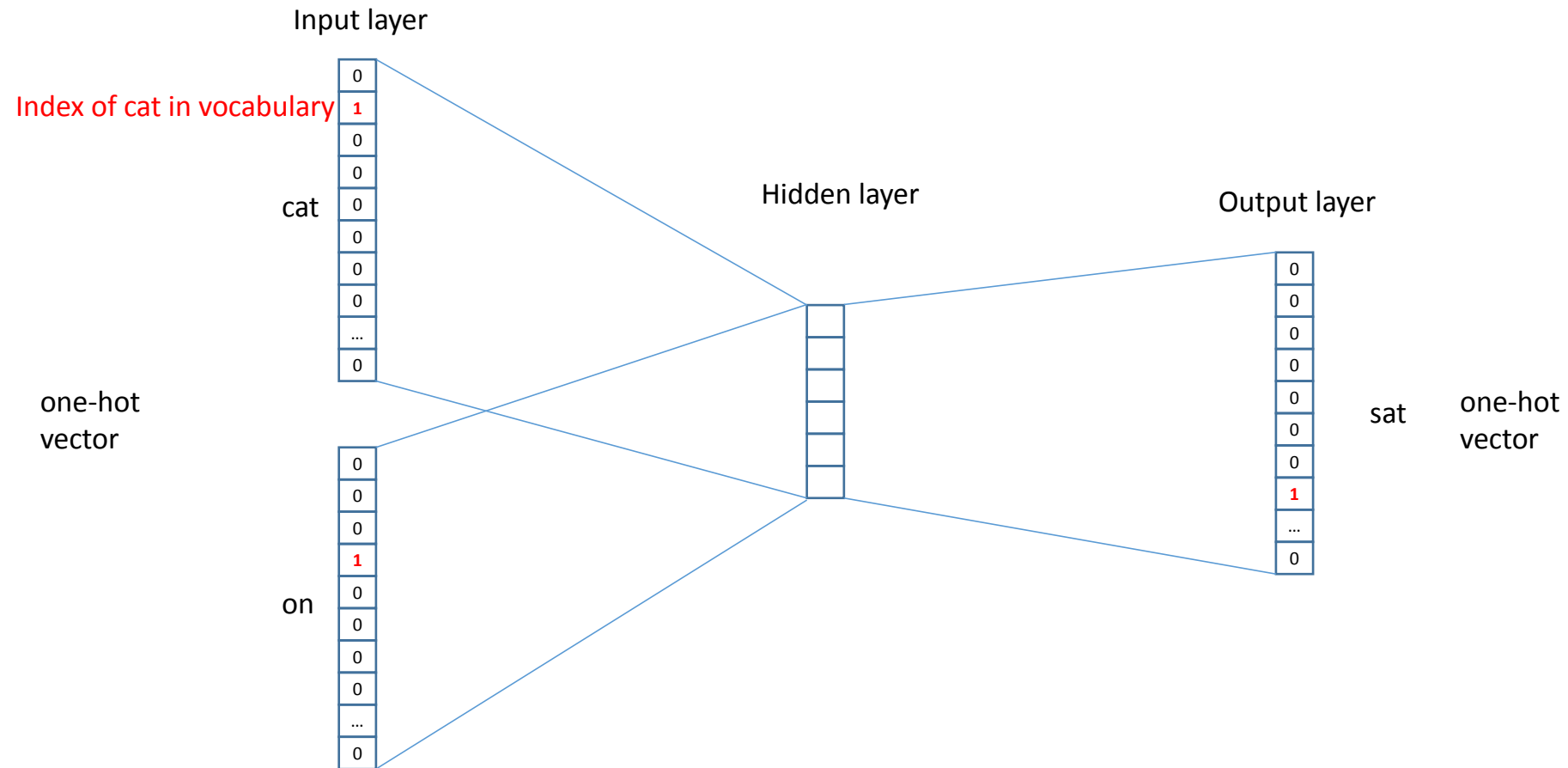- Of course there are more studies on sub-words ;-)

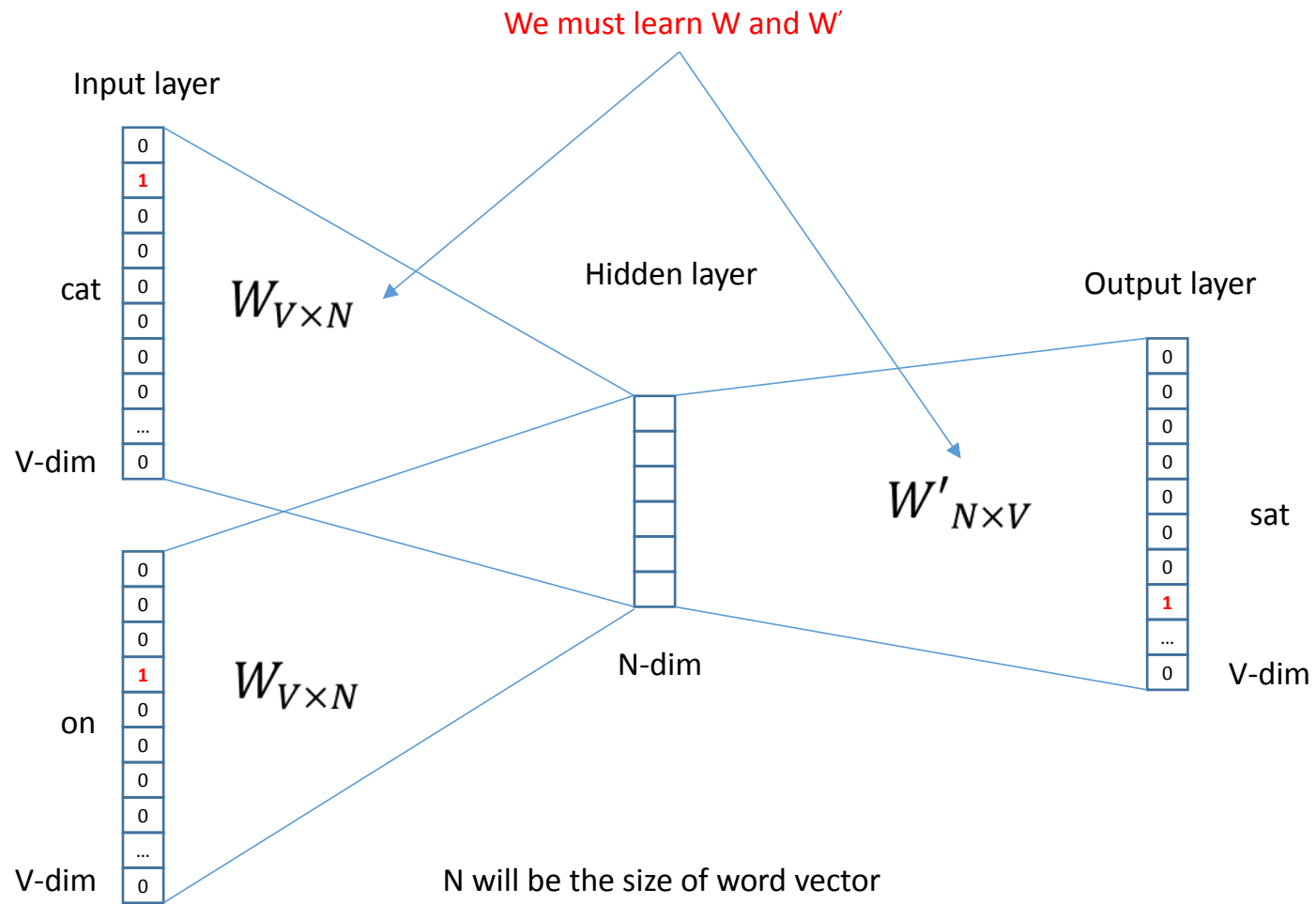# Neural Representation Model for Words

- What information can be leveraged?
  - Context!
- N-gram (word association) language modeling is the most effecient way for building word-word relations
- Two ways:
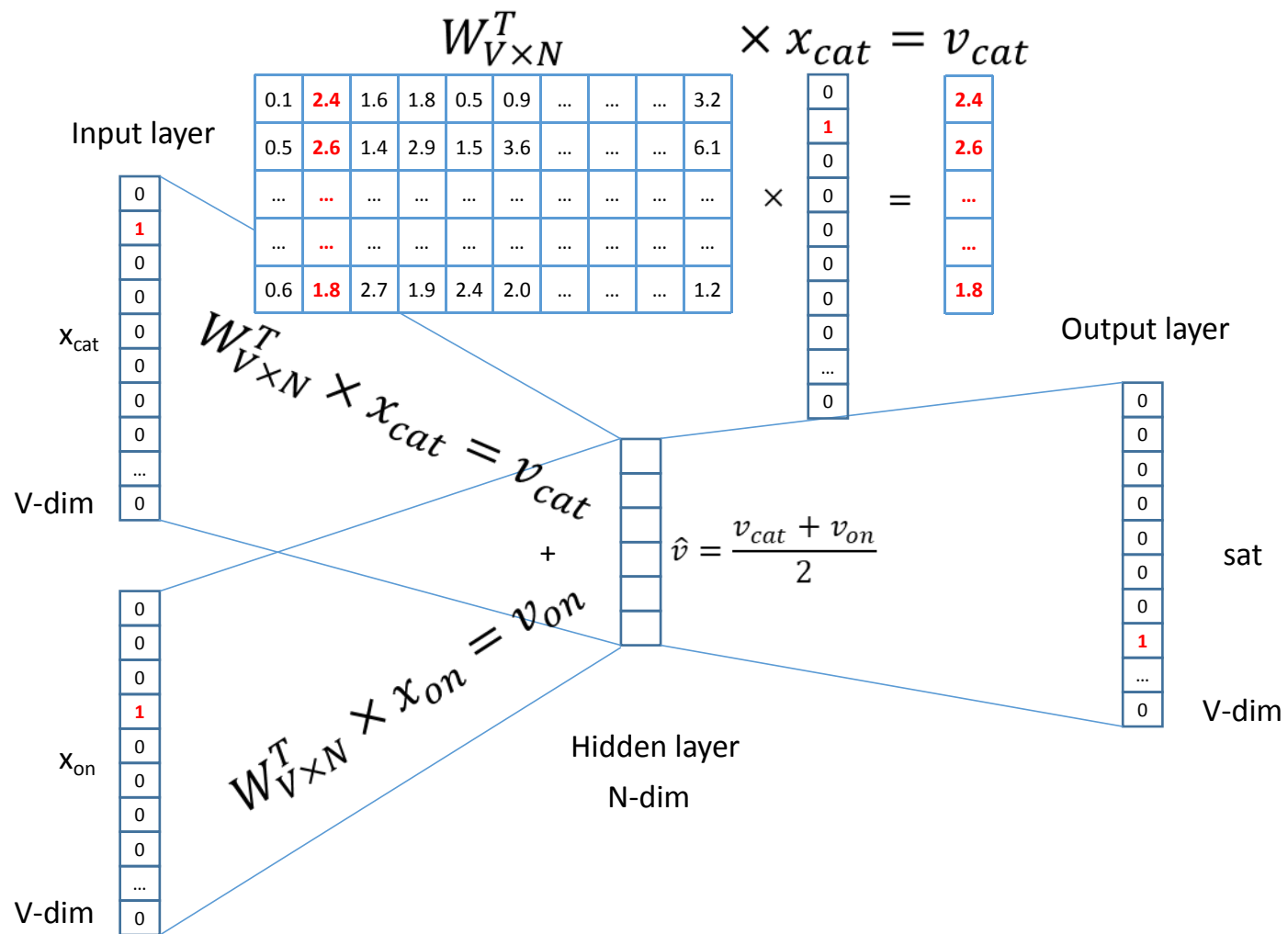  - Top-down (GloVe)
  - Bottom-up (word2vec)

# word2vec - Concept

- Key idea: Predict surrounding words of every word
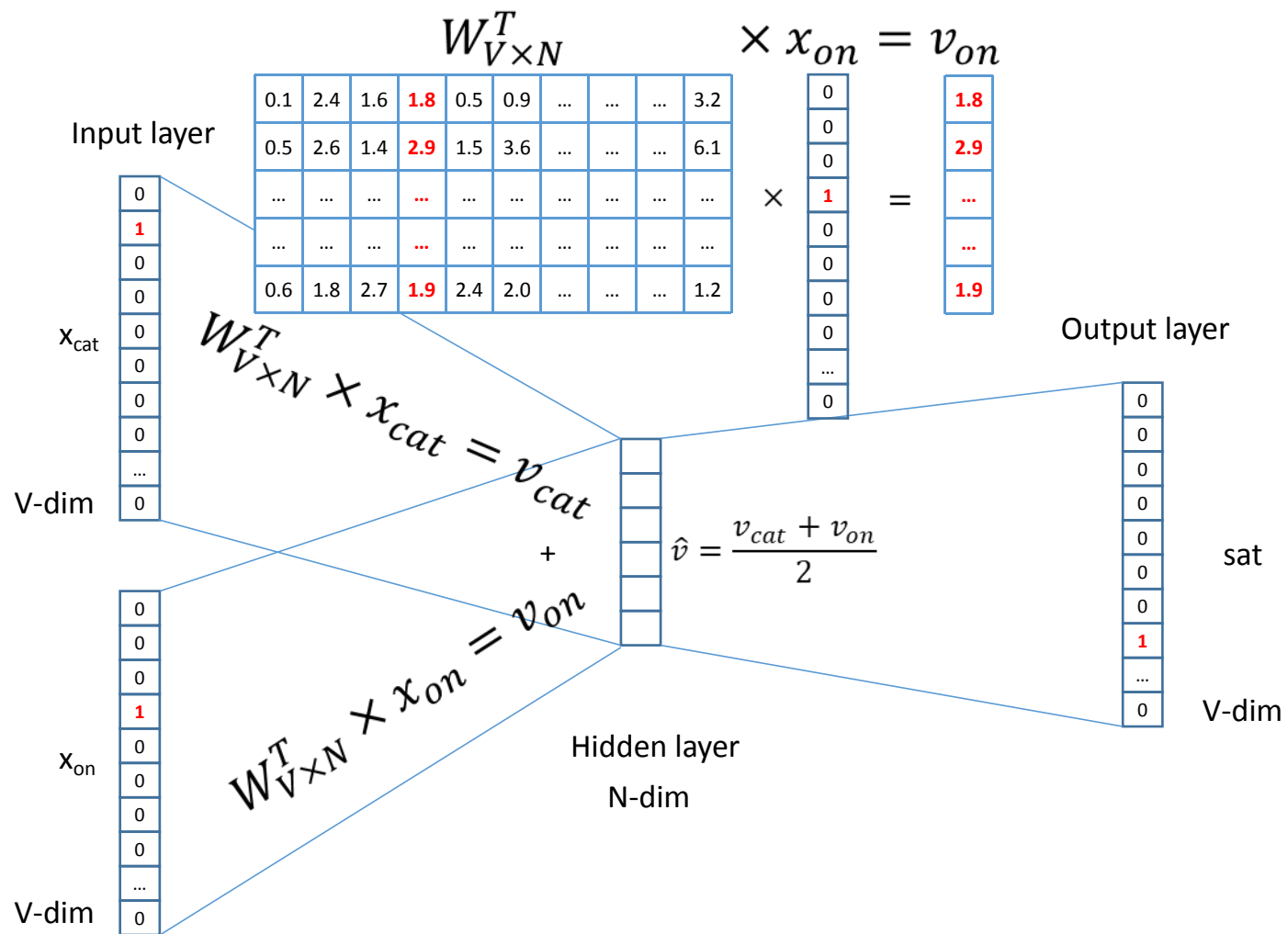
- Faster and can easily incorporate a new sentence/document or add a word to the vocabulary

- An indirect matrix decomposition way of learning dense representations (compared with GloVe)

- Let's start from a typical process of neural language modeling

# Illustration

Input layer

Index of cat in vocabulary

Hidden layer

Output layer

cat

sat

one-hot vector

one-hot vector

on

We must learn W and W′

Input layer

| 0 |
|---|
| **1** |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| … |
| 0 |

cat

V-dim

$W_{V \times N}$

Hidden layer

| 0 |
|---|
| 0 |
| 0 |
| **1** |
| 0 |
| 0 |
| 0 |
| 0 |
| … |
| 0 |

on

V-dim

$W_{V \times N}$

N-dim

$W'_{N \times V}$

Output layer

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| **1** |
| … |
| 0 |

sat

V-dim

N will be the size of word vector

$$W_{V \times N}^{T} \times x_{cat} = v_{cat}$$

Input layer

| 0.1 | 2.4 | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.5 | 2.6 | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | 1.8 | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

$\times$

| 0 |
|---|
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

$=$

| 2.4 |
|-----|
| 2.6 |
| ... |
| ... |
| 1.8 |

$x_{cat}$

| 0 |
|---|
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

V-dim

$W_{V \times N}^{T} \times x_{cat} = v_{cat}$

Output layer

$\hat{v} = \dfrac{v_{cat} + v_{on}}{2}$

+

$x_{on}$

| 0 |
|---|
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

V-dim

$W_{V \times N}^{T} \times x_{on} = v_{on}$

Hidden layer

N-dim

sat

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| ... |
| 0 |

V-dim

Input layer

0
1
0
0
cat  0
0
0
0
...
V-dim  0

$W_{V \times N}$

Hidden layer

Output layer

$W'_{V \times N} \times \hat{v} = z$

0
0
0
0
0
0
0
1
...
0

$\hat{y} = softmax(z)$

$\hat{v}$

N-dim

$\hat{y}_{\text{sat}}$

V-dim

0
0
0
1
on  0
0
0
0
...
V-dim  0

$W_{V \times N}$

N will be the size of word vector

Input layer

| 0 |
| 1 |
| 0 |
| 0 |
cat | 0 |
| 0 |
| 0 |
| 0 |
| ... |
V-dim | 0 |

| 0 |
| 0 |
| 0 |
| 1 |
on | 0 |
| 0 |
| 0 |
| 0 |
| ... |
V-dim | 0 |

$W_{V \times N}$

$W_{V \times N}$

Hidden layer

$\hat{v}$

N-dim

N will be the size of word vector

Output layer

$W'_{V \times N} \times \hat{v} = z$
$\hat{y} = softmax(z)$

| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| ... |
| 0 |

$\hat{y}_{\text{sat}}$

V-dim

We would prefer $\hat{y}$ close to $\hat{y}_{sat}$

| 0.01 |
| 0.02 |
| 0.00 |
| 0.02 |
| 0.01 |
| 0.02 |
| 0.01 |
| 0.7 |
| ... |
| 0.00 |

$\hat{y}$

11

$$W_{V \times N}^T$$

| 0.1 | **2.4** | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|---------|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.5 | **2.6** | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | **1.8** | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

Contain word's vectors

Input layer

| 0 |
|---|
| **1** |
| 0 |
| 0 |

$x_{cat}$

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| ... |

V-dim | 0 |

$W_{V \times N}$

| 0 |
|---|
| 0 |
| 0 |
| **1** |

$W_{V \times N}$

$x_{on}$

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| ... |

V-dim | 0 |

Hidden layer

N-dim

Output layer

$W'_{V \times N}$

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |

sat

| **1** |
|---|
| ... |
| 0 |

V-dim

We can consider either W or W' as the word's representation. Or even take the average.

# word2vec - Basics

- Two basic structures:
  - Continuous Bag of Word (CBOW): use a window of word to predict the middle word
  - Skip-gram (SG): use a word to predict the surrounding ones in window.



CBOW

Skip-gram

# word2vec - Basics

- Challenges of Implementation
  - Sparse output layer
  - Matrix (vector) computation over a huge corpus
  - Deal with words with different frequencies
  - Effecient gradient updating
  - …

# word2vec - Basics

- Two prediction strategies
  - Hierarchical Softmax
    - Several pipelined softmaxes go through a tree structure
    - Effective to capture words' information when they have similar patterns in frequencies
  - Negative sampling
    - A binary classification to distinguish if a word is in context or not
    - Effective when training data is huge (less softmax operations)

# Hierarchical Softmax

- An illustration

- Softmax is done at each node along the path from root to leaves

- Sibling nodes may share more frequency patterns

- Effectively convert a high-dimensional classification problem into a series of small ones

# word2vec - Implementation

- Workflow
  - Load data
  - Prepare data structure
  - Training (CBOW or SG)
    - if with HS
      - predict go through hierarchical path
      - update
    - else: with NS
      - randomly choose negative samples
      - predict and update
  - Save model

# word2vec - Implementation

- Load data
  - The goal is to create a vocabulary
  - Related Functions
    - AddWordToVocab
    - GetWordHash
    - LearnVocabFromTrainFile
    - ReadVocab
    - ReadWord
    - ReadWordIndex
    - ReduceVocab
    - SearchVocab
    - SortVocab

# word2vec - Implementation

- Some tricks
  - All embedddings are stored in a very long 1-D array.

```
341    a = posix_memalign((void **)&syn0, 128, (long long)vocab_size * layer1_size * sizeof(real));
342    if (syn0 == NULL) {printf("Memory allocation failed\n"); exit(1);}
```

  - Sigmoid function is precomputed

```
695    expTable = (real *)malloc((EXP_TABLE_SIZE + 1) * sizeof(real));
696    for (i = 0; i < EXP_TABLE_SIZE; i++) {
697      expTable[i] = exp((i / (real)EXP_TABLE_SIZE * 2 - 1) * MAX_EXP); // Precompute the exp() table
698      expTable[i] = expTable[i] / (expTable[i] + 1);                   // Precompute f(x) = x / (x + 1)
699    }
```

- important variables: syn0, syn1, neu1, neu1e

# word2vec - Implementation

- Prepare data structure

```c
195  // Create binary Huffman tree using the word counts
196  // Frequent words will have short uniqe binary codes
197  void CreateBinaryTree() {
198      long long a, b, i, min1i, min2i, pos1, pos2, point[MAX_CODE_LENGTH];
199      char code[MAX_CODE_LENGTH];
200      long long *count = (long long *)calloc(vocab_size * 2 + 1, sizeof(long long));
201      long long *binary = (long long *)calloc(vocab_size * 2 + 1, sizeof(long long));
202      long long *parent_node = (long long *)calloc(vocab_size * 2 + 1, sizeof(long long));
203      for (a = 0; a < vocab_size; a++) count[a] = vocab[a].cn;
204      for (a = vocab_size; a < vocab_size * 2; a++) count[a] = 1e15;
205      pos1 = vocab_size - 1;
206      pos2 = vocab_size;
207      // Following algorithm constructs the Huffman tree by adding one node at a time
208      for (a = 0; a < vocab_size - 1; a++) {
209          // First, find two smallest nodes 'min1, min2'
210          if (pos1 >= 0) {
211              if (count[pos1] < count[pos2]) {
212                  min1i = pos1;
213                  pos1--;
214              } else {
215                  min1i = pos2;
216                  pos2++;
217              }
218          } else {
219              min1i = pos2;
220              pos2++;
221          }
222          if (pos1 >= 0) {
223              if (count[pos1] < count[pos2]) {
224                  min2i = pos1;
225                  pos1--;
226              } else {
227                  min2i = pos2;
228                  pos2++;
229              }
230          } else {
231              min2i = pos2;
232              pos2++;
233          }
234          count[vocab_size + a] = count[min1i] + count[min2i];
235          parent_node[min1i] = vocab_size + a;
236          parent_node[min2i] = vocab_size + a;
237          binary[min2i] = 1;
238      }
239      // Now assign binary code to each vocabulary word
240      for (a = 0; a < vocab_size; a++) {
241          b = a;
242          i = 0;
243          while (1) {
244              code[i] = binary[b];
245              point[i] = b;
246              i++;
247              b = parent_node[b];
248              if (b == vocab_size * 2 - 2) break;
249          }
250          vocab[a].codelen = i;
251          vocab[a].point[0] = vocab_size - 2;
252          for (b = 0; b < i; b++) {
253              vocab[a].code[i - b - 1] = code[b];
254              vocab[a].point[i - b] = point[b] - vocab_size;
255          }
256      }
257      free(count);
258      free(binary);
259      free(parent_node);
260  }
261
```

# word2vec - Implementation

- CBOW



```
422   if (cbow) {  //train the cbow architecture
423     // in -> hidden
424     cw = 0;
425     for (a = b; a < window * 2 + 1 - b; a++) if (a != window) {
426       c = sentence_position - window + a;
427       if (c < 0) continue;
428       if (c >= sentence_length) continue;
429       last_word = sen[c];
430       if (last_word == -1) continue;
431       for (c = 0; c < layer1_size; c++) neu1[c] += syn0[c + last_word * layer1_size];
432       cw++;
433     }
434     if (cw) {
435       for (c = 0; c < layer1_size; c++) neu1[c] /= cw;
436       if (hs) for (d = 0; d < vocab[word].codelen; d++) {
437         f = 0;
438         l2 = vocab[word].point[d] * layer1_size;
439         // Propagate hidden -> output
440         for (c = 0; c < layer1_size; c++) f += neu1[c] * syn1[c + l2];
441         if (f <= -MAX_EXP) continue;
442         else if (f >= MAX_EXP) continue;
443         else f = expTable[(int)((f + MAX_EXP) * (EXP_TABLE_SIZE / MAX_EXP / 2))];
444         // 'g' is the gradient multiplied by the learning rate
445         g = (1 - vocab[word].code[d] - f) * alpha;
446         // Propagate errors output -> hidden
447         for (c = 0; c < layer1_size; c++) neu1e[c] += g * syn1[c + l2];
448         // Learn weights hidden -> output
449         for (c = 0; c < layer1_size; c++) syn1[c + l2] += g * neu1[c];
450       }
451       // NEGATIVE SAMPLING
452       if (negative > 0) for (d = 0; d < negative + 1; d++) {
453         if (d == 0) {
454           target = word;
```

Projection

```
455           label = 1;
456         } else {
457           next_random = next_random * (unsigned long long)25214903917 + 11;
458           target = table[(next_random >> 16) % table_size];
459           if (target == 0) target = next_random % (vocab_size - 1) + 1;
460           if (target == word) continue;
461           label = 0;
462         }
463         l2 = target * layer1_size;
464         f = 0;
465         for (c = 0; c < layer1_size; c++) f += neu1[c] * syn1neg[c + l2];
466         if (f > MAX_EXP) g = (label - 1) * alpha;
467         else if (f < -MAX_EXP) g = (label - 0) * alpha;
468         else g = (label - expTable[(int)((f + MAX_EXP) * (EXP_TABLE_SIZE / MAX_EXP / 2))]) * alpha;
469         for (c = 0; c < layer1_size; c++) neu1e[c] += g * syn1neg[c + l2];
470         for (c = 0; c < layer1_size; c++) syn1neg[c + l2] += g * neu1[c];
471       }
472       // hidden -> in
473       for (a = b; a < window * 2 + 1 - b; a++) if (a != window) {
474         c = sentence_position - window + a;
475         if (c < 0) continue;
476         if (c >= sentence_length) continue;
477         last_word = sen[c];
478         if (last_word == -1) continue;
479         for (c = 0; c < layer1_size; c++) syn0[c + last_word * layer1_size] += neu1e[c];
480       }
481     }
```

# word2vec - Implementation

- SG



```
482  } else {  //train skip-gram
483    for (a = b; a < window * 2 + 1 - b; a++) if (a != window) {
484      c = sentence_position - window + a;
485      if (c < 0) continue;
486      if (c >= sentence_length) continue;
487      last_word = sen[c];
488      if (last_word == -1) continue;
489      l1 = last_word * layer1_size;
490      for (c = 0; c < layer1_size; c++) neu1e[c] = 0;
491      // HIERARCHICAL SOFTMAX
492      if (hs) for (d = 0; d < vocab[word].codelen; d++) {
493        f = 0;
494        l2 = vocab[word].point[d] * layer1_size;
495        // Propagate hidden -> output
496        for (c = 0; c < layer1_size; c++) f += syn0[c + l1] * syn1[c + l2];
497        if (f <= -MAX_EXP) continue;
498        else if (f >= MAX_EXP) continue;
499        else f = expTable[(int)((f + MAX_EXP) * (EXP_TABLE_SIZE / MAX_EXP / 2))];
500        // 'g' is the gradient multiplied by the learning rate
501        g = (1 - vocab[word].code[d] - f) * alpha;
502        // Propagate errors output -> hidden
503        for (c = 0; c < layer1_size; c++) neu1e[c] += g * syn1[c + l2];
504        // Learn weights hidden -> output
505        for (c = 0; c < layer1_size; c++) syn1[c + l2] += g * syn0[c + l1];
506      }
507      // NEGATIVE SAMPLING
508      if (negative > 0) for (d = 0; d < negative + 1; d++) {
509        if (d == 0) {
510          target = word;
511          label = 1;
512        } else {
513          next_random = next_random * (unsigned long long)25214903917 + 11;
514          target = table[(next_random >> 16) % table_size];
```

Projection

```
515          if (target == 0) target = next_random % (vocab_size - 1) + 1;
516          if (target == word) continue;
517          label = 0;
518        }
519        l2 = target * layer1_size;
520        f = 0;
521        for (c = 0; c < layer1_size; c++) f += syn0[c + l1] * syn1neg[c + l2];
522        if (f > MAX_EXP) g = (label - 1) * alpha;
523        else if (f < -MAX_EXP) g = (label - 0) * alpha;
524        else g = (label - expTable[(int)((f + MAX_EXP) * (EXP_TABLE_SIZE / MAX_EXP / 2))]) * alpha;
525        for (c = 0; c < layer1_size; c++) neu1e[c] += g * syn1neg[c + l2];
526        for (c = 0; c < layer1_size; c++) syn1neg[c + l2] += g * syn0[c + l1];
527      }
528      // Learn weights input -> hidden
529      for (c = 0; c < layer1_size; c++) syn0[c + l1] += neu1e[c];
530    }
531  }
532  sentence_position++;
533  if (sentence_position >= sentence_length) {
534    sentence_length = 0;
535    continue;
536  }
537 }
```

# word2vec - Implementation

- A closer look of HS

```
491    // HIERARCHICAL SOFTMAX
492    if (hs) for (d = 0; d < vocab[word].codelen; d++) {
493      f = 0;
494      l2 = vocab[word].point[d] * layer1_size;
495      // Propagate hidden -> output
496      for (c = 0; c < layer1_size; c++) f += syn0[c + l1] * syn1[c + l2];
497      if (f <= -MAX_EXP) continue;
498      else if (f >= MAX_EXP) continue;
499      else f = expTable[(int)((f + MAX_EXP) * (EXP_TABLE_SIZE / MAX_EXP / 2))];
500      // 'g' is the gradient multiplied by the learning rate
501      g = (1 - vocab[word].code[d] - f) * alpha;
502      // Propagate errors output -> hidden
503      for (c = 0; c < layer1_size; c++) neu1e[c] += g * syn1[c + l2];
504      // Learn weights hidden -> output
505      for (c = 0; c < layer1_size; c++) syn1[c + l2] += g * syn0[c + l1];
506    }
```

Go through the path

# word2vec - Implementation

- A closer look of NS

```
507    // NEGATIVE SAMPLING
508    if (negative > 0) for (d = 0; d < negative + 1; d++) {
509        if (d == 0) {
510            target = word;
511            label = 1;
512        } else {
513            next_random = next_random * (unsigned long long)25214903917 + 11;
514            target = table[(next_random >> 16) % table_size];
515            if (target == 0) target = next_random % (vocab_size - 1) + 1;
516            if (target == word) continue;
517            label = 0;
518        }
519        l2 = target * layer1_size;
520        f = 0;
521        for (c = 0; c < layer1_size; c++) f += syn0[c + l1] * syn1neg[c + l2];
522        if (f > MAX_EXP) g = (label - 1) * alpha;
523        else if (f < -MAX_EXP) g = (label - 0) * alpha;
524        else g = (label - expTable[(int)((f + MAX_EXP) * (EXP_TABLE_SIZE / MAX_EXP / 2))]) * alpha;
525        for (c = 0; c < layer1_size; c++) neu1e[c] += g * syn1neg[c + l2];
526        for (c = 0; c < layer1_size; c++) syn1neg[c + l2] += g * syn0[c + l1];
527    }
```

Positive

Negative

# word2vec - Implementation

- Update the input embeddings
  - CBOW

```
472    // hidden -> in
473    for (a = b; a < window * 2 + 1 - b; a++) if (a != window) {
474      c = sentence_position - window + a;
475      if (c < 0) continue;
476      if (c >= sentence_length) continue;
477      last_word = sen[c];
478      if (last_word == -1) continue;
479      for (c = 0; c < layer1_size; c++) syn0[c + last_word * layer1_size] += neu1e[c];
480    }
```

  - SG

```
528    // Learn weights input -> hidden
529    for (c = 0; c < layer1_size; c++) syn0[c + l1] += neu1e[c];
```

# Evaluation

- We already know cosine similarity is the main metric for intrinsic evaluation

- But, how to sysmatically evaluate?
  - Prepare a human annotated similarity/relatedness dataset
    - A list of word pairs
  - Compute the similarities of all pairs using the resulted embeddings
  - Compute the correlation between human annotations and the similarities

# Evaluation

- Spearman's correlation

$$\rho = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}}.$$

| | | x | y |
|---|---|---|---|
| tiger | cat | 7.35 | 0.517977544 |
| tiger | tiger | 10.00 | 1.0 |
| plane | car | 5.77 | 0.217732013636 |
| train | car | 6.31 | 0.19087634791 |
| television | radio | 6.77 | 0.257970738766 |
| media | radio | 7.42 | 0.212014745796 |
| bread | butter | 6.19 | 0.7017263618 |
| cucumber | potato | 5.92 | 0.396596853365 |
| doctor | nurse | 7.00 | 0.382799131727 |
| professor | doctor | 6.62 | 0.288006966845 |
| student | professor | 6.81 | 0.19973114775 |
| smart | stupid | 5.81 | 0.228546918905 |
| wood | forest | 7.73 | 0.229380996826 |

# Hw4

- Train your word2vec embedding using the Text8 corpus
  - with HS (CBOW, SG)
  - with NS (CBOW, SG)
- Evaluate using spearman's correlation
- Submit a result.txt to me with
  - The results from five embeddings (4 word2vec and 1 GloVe)
  - Your observations on HS v.s. NS on CBOW and SG