# How it works

1. Identify a misspelled word
2. Find strings n edit distance away
3. Filter candidates
4. Calculate word probabilities

deah → dear ✅
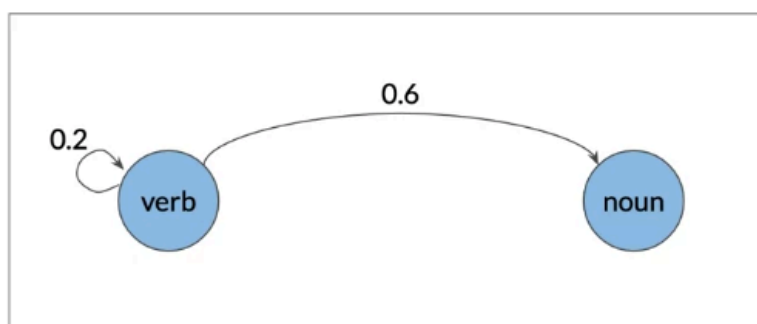yeah
dear
dean
... *etc*

# Part of speech (POS) tagging

## Part of speech tags:

| lexical term | tag | example |
|---|---|---|
| noun | NN | something, nothing |
| verb | VB | learn, study |
| determiner | DT | the, a |
| w-adverb | WRB | why, where |
| ... | ... | |

Markov Chains

# The transition matrix



$$A = \begin{array}{c|ccc} & \text{NN} & \text{VB} & \text{O} \\ \hline \text{NN (noun)} & 0.2 & 0.2 & 0.6 \\ \text{VB (verb)} & 0.4 & 0.3 & 0.3 \\ \text{O (other)} & 0.2 & 0.3 & 0.5 \end{array}$$

$$\sum_{j=1}^{N} a_{ij} = 1$$

# Transition probabilities



1. Count occurrences of tag pairs

$$C(t_{i-1}, t_i)$$

2. Calculate probabilities using the counts

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{\sum_{j=1}^{N} C(t_{i-1}, t_j)}$$

# Populating the transition matrix

$$A = \begin{array}{c|cccc} & \text{NN} & \text{VB} & \text{O} & \\ \hline \pi & 1 & 0 & 2 & 3 \\ \text{NN} & 0 & 0 & 6 & 6 \\ \text{VB} & 0 & 0 & 0 & 0 \\ \text{O} & 6 & 0 & 8 & 14 \end{array}$$

$$P(\text{NN|O}) = \frac{C(\text{O, NN})}{\sum_{j=1}^{N} C(\text{O}, t_j)} = \frac{6}{14}$$

# Smoothing

|   | NN | VB | O |
|---|---|---|---|
| $\pi$ | 0.3333 | 0.0003 | 0.6663 |
| NN | 0.0001 | 0.0001 | 0.9996 |
| VB | 0.3333 | 0.3333 | 0.3333 |
| O | 0.4285 | 0.0000 | 0.5713 |

$A =$ (to the left of the table)

$$P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i) + \epsilon}{\sum_{j=1}^{N} C(t_{i-1}, t_j) + N * \epsilon}$$

# Viterbi algorithm – a graph algorithm

- Create **language model (LM)** from text corpus to
  - Estimate probability of word sequences
  - Estimate probability of a word following a sequence of words
- Apply this concept to **autocomplete a sentence** with most likely suggestions



## Learning objectives

- Process text corpus to N-gram language model
- Out of vocabulary words
- Smoothing for previously unseen N-grams
- Language model evaluation

Sentence auto-complete

## N-gram

An N-gram is a sequence of N words

Corpus: I am happy because I am learning

Unigrams: { I , am , happy , because , learning }

Bigrams: { I  am , am happy , happy because … }        ❌ I happy

# Trigram Probability

Corpus: I am happy because I am learning

$$P(happy|I\ am) = \frac{C(I\ am\ happy)}{C(I\ am)} = \frac{1}{2}$$

Probability of a trigram:
$$P(w_3|w_1^2) = \frac{C(w_1^2\ w_3)}{C(w_1^2)}$$

# N-gram probability

Probability of N-gram:
$$P(w_N|w_1^{N-1}) = \frac{\boxed{C(w_1^{N-1}\ w_N)}}{C(w_1^{N-1})}$$

$$C(w_1^{N-1}\ w_N) = C(w_1^N)$$

# Approximation of sequence probability

- Markov assumption: only last N words matter

- Bigram $\quad P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-1})$
- N-gram $\quad P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-N+1}^{n-1})$

- Entire sentence modeled with bigram $\quad P(w_1^n) \approx \prod_{i=1}^{n} P(w_i|w_{i-1})$

# Perplexity

$$PP(W) = P(s_1, s_2, ..., s_m)^{-\frac{1}{m}}$$

**W** → test set containing **m** sentences **s**
$s_i$ → i-th sentence in the test set, each ending with </s>
**m** → number of all words in entire test set **W** including
</s> but not including <s>

# Perplexity for bigram models

$$PP(W) = \sqrt[m]{\prod_{i=1}^{m} \prod_{j=1}^{|s_i|} \frac{1}{P(w_j^{(i)}|w_{j-1}^{(i)})}}$$

$w_j^{(i)}$ → j-th word in i-th sentence

- concatenate all sentences in W

$$PP(W) = \sqrt[m]{\prod_{i=1}^{m} \frac{1}{P(w_i|w_{i-1})}}$$

$w_i$ → i-th word in test set

# Smoothing

- Add-one smoothing (Laplacian smoothing)

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}, w_n) + 1}{\sum_{w \in V}(C(w_{n-1}, w) + 1)} = \frac{C(w_{n-1}, w_n) + 1}{C(w_{n-1}) + V}$$

- Add-k smoothing

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}, w_n) + k}{\sum_{w \in V}(C(w_{n-1}, w) + k)} = \frac{C(w_{n-1}, w_n) + k}{C(w_{n-1}) + k * V}$$
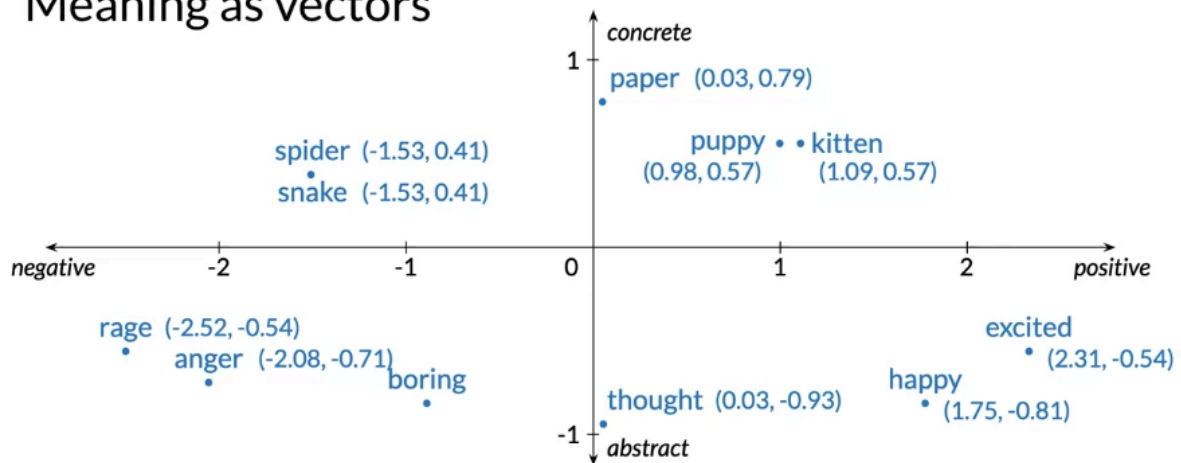
# Summary

- N-Grams and probabilities
- Approximate sentence probability from N-Grams
- Build language model from corpus
- Fix missing information
  - Out of vocabulary words with <UNK>
  - Missing N-Gram in corpus with smoothing, backoff and interpolation
- Evaluate language model with perplexity

## Meaning as vectors

rage anger spider boring     paper          kitten     happy excited

negative ←——————————————————————————————————————→ positive
          -2          -1         0          1          2

## Meaning as vectors

*concrete*
1 —
paper (0.03, 0.79)

spider (-1.53, 0.41)
snake (-1.53, 0.41)

puppy • • kitten
(0.98, 0.57)    (1.09, 0.57)

negative ←————————————————————————————————————→ positive
          -2          -1         0          1          2

rage (-2.52, -0.54)                                    excited
    anger (-2.08, -0.71)                            • (2.31, -0.54)
            boring                         happy
                         thought (0.03, -0.93)   • (1.75, -0.81)
-1 —
*abstract*

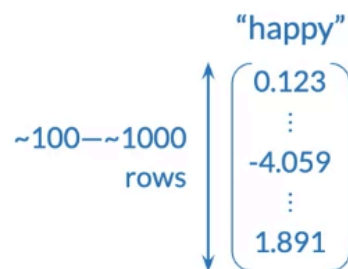## Word embedding vectors

+ Low dimension

+ Embed meaning
  ○ e.g. semantic distance

    forest ≈ tree     forest ≢ ticket

  ○ e.g. analogies

    Paris:France :: Rome:?

"happy"

~100—~1000
rows

$\begin{pmatrix} 0.123 \\ \vdots \\ -4.059 \\ \vdots \\ 1.891 \end{pmatrix}$

# Word embedding process



# Basic word embedding methods

- word2vec (Google, 2013)
  - Continuous bag-of-words (CBOW)
  - Continuous skip-gram / Skip-gram with negative sampling (SGNS)

- Global Vectors (GloVe) (Stanford, 2014)

- fastText (Facebook, 2016)
  - Supports out-of-vocabulary (OOV) words

# Advanced word embedding methods

Deep learning, contextual embeddings

- BERT (Google, 2018)

- ELMo (Allen Institute for AI, 2018)

- GPT-2 (OpenAI, 2018)

Tunable pre-trained models available

# Creating a training example

center word

I am **happy** because I am learning

context words
C = 2
context half-size

window
window size = 5

# Final prepared training set

| Context words | Context words vector | Center word | Center word vector |
|---|---|---|---|
| I am because I | [0.25; 0.25; 0; 0.5; 0] | happy | [0; 0; 1; 0; 0] |

# Architecture of the CBOW model

**Hyperparameters**
N: Word embedding size    ...

Input layer          Hidden layer          Output layer

Context words
vector

$x =$
V

"I am happy
because I am
learning"
$\rightarrow V = 5$

$W_1$
weights
$b_1$
biases
ReLU

$W_2$
weights
$b_2$
biases
softmax

Center word
vector

$\hat{y} =$
V

V    x         N    h         V    $\hat{y}$

# Dimensions (single input)

Input layer          Hidden layer          Output layer

$\mathbf{W}_1$   $N \times V$

$\mathbf{b}_1$   $N \times 1$

ReLU

$\mathbf{W}_2$   $V \times N$

$\mathbf{b}_2$   $V \times 1$

softmax

$\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$   $N \times 1$

$\mathbf{h} = \text{ReLU}(\mathbf{z}_1)$   $N \times 1$

$\mathbf{z}_2 = \mathbf{W}_2 \mathbf{h} + \mathbf{b}_2$   $V \times 1$

$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z}_2)$   $V \times 1$

$\mathbf{x}$   $V \times 1$

$\mathbf{h}$   $N \times 1$

$\hat{\mathbf{y}}$   $V \times 1$

# Dimensions (batch input)

Input layer    Hidden layer    Output layer

Context words matrix

Predicted center word matrix

$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)} & \cdots & \mathbf{x}^{(m)} \end{bmatrix}$   $V$   $m$

$\mathbf{W}_1$   $\mathbf{B}_1$   ReLU

$\mathbf{W}_2$   $\mathbf{B}_2$   softmax

$\hat{\mathbf{Y}} = \begin{bmatrix} \hat{\mathbf{y}}^{(1)} & \cdots & \hat{\mathbf{y}}^{(m)} \end{bmatrix}$   $V$   $m$

$\mathbf{X}$   $V \times m$

$\mathbf{H}$   $N \times m$

$\hat{\mathbf{Y}}$   $V \times m$

1.00

# Cross-entropy loss

$$J = -\sum_{k=1}^{V} y_k \log \hat{y}_k$$
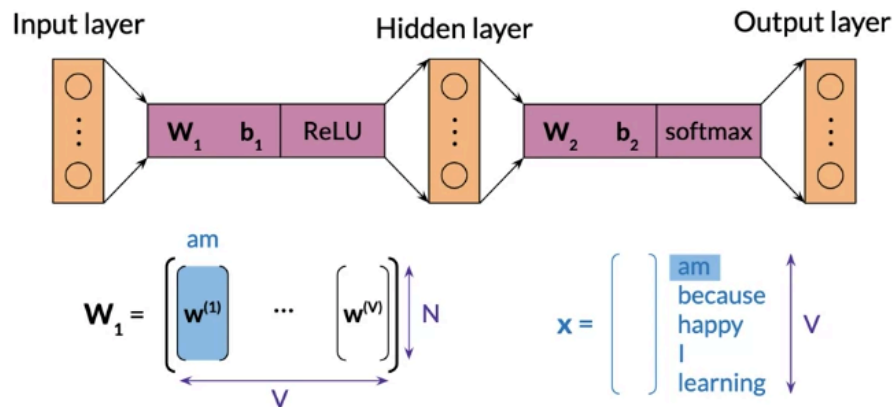
Actual $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_V \end{pmatrix}$     Predicted $\hat{\mathbf{y}} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_V \end{pmatrix}$

I am happy because I am learning

| $\mathbf{y}$ | | $\hat{\mathbf{y}}$ | | $\log(\hat{\mathbf{y}})$ | | $\mathbf{y} \odot \log(\hat{\mathbf{y}})$ | | |
|---|---|---|---|---|---|---|---|---|
| 0 | am | 0.083 | | -2.49 | | 0 | | |
| 0 | because | 0.03 | | -3.49 | | 0 | | |
| 1 | happy | 0.611 | log | -0.49 | $\odot \mathbf{y}$ | -0.49 | $-\Sigma$ | J = 0.49 |
| 0 | I | 0.225 | | -1.49 | | 0 | | |
| 0 | learning | 0.05 | | -2.49 | | 0 | | |

# Extracting word embedding vectors: option 1

Input layer           Hidden layer          Output layer

| $W_1$ | $b_1$ | ReLU |

| $W_2$ | $b_2$ | softmax |

am

$$W_1 = \begin{bmatrix} w^{(1)} & \cdots & w^{(V)} \end{bmatrix} \bigg] N$$

V

$$x = \begin{bmatrix} \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{bmatrix} \bigg] V$$

# Intrinsic evaluation

Test relationships between words

- Analogies

- Clustering

- Visualization

             village

    city  town

  gas    country

oil           happy

petroleum    sad  joyful

# Extrinsic evaluation

Test word embeddings on external task
e.g. named entity recognition, parts-of-speech tagging

+ Evaluates actual usefulness of embeddings

- Time-consuming

- More difficult to troubleshoot