# Auto-correct and Dynamic Programming

## Auto-correct

is an application that changes misspelled words into the correct ones.

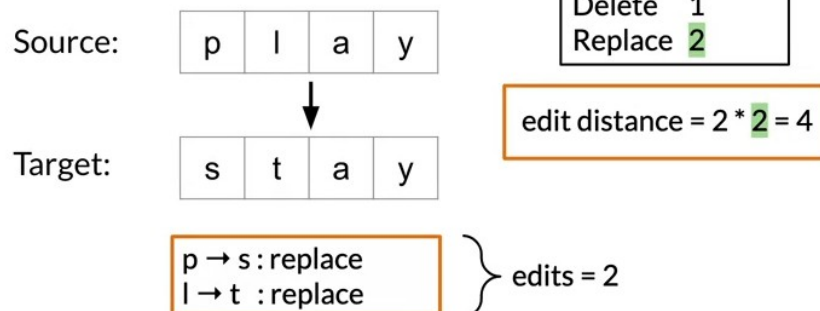- Example: Happy birthday deah friend! ==> dear

- How it works:

      1.Identify a misspelled word

      2. Find strings n edit distance away

      3. Filter candidates

      4. Calculate word probabilities

## Minimum edit distance

    - Evaluate the similarity between 2 strings.

    - Minimum number of edits needed to transform 1 string into another.

- the algorithm try to minimize the edit cost.



- Example:

Source: | p | l | a | y |

Target: | s | t | a | y |

Edit cost:
Insert   1
Delete   1
Replace  2

edit distance = 2 * 2 = 4

p → s : replace
l → t : replace

edits = 2

## Applications:

    - Spelling correction

    - document similarity

    - machine translation

    - DNA sequencing

# Part of Speech Tagging and Hidden Markov Models

## Part of Speech Tagging

The category of words or the lexical terms in the language.

**Tags:** Noun, Verb, adjective, preposition, adverb,...

**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 |
| ... | ... | |

Why not learn something ?

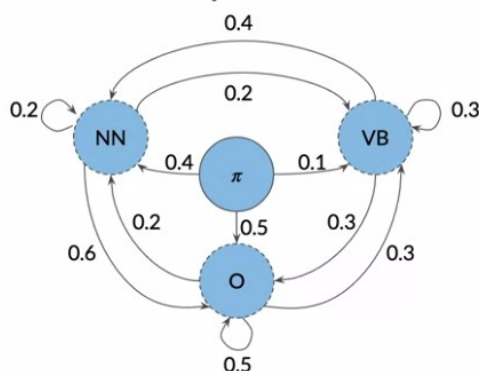WRB    RB      VB              NN          .

**Applications:**

- Named entities

- Co-reference resolution

- Speech recognition

---

## Markov Chains

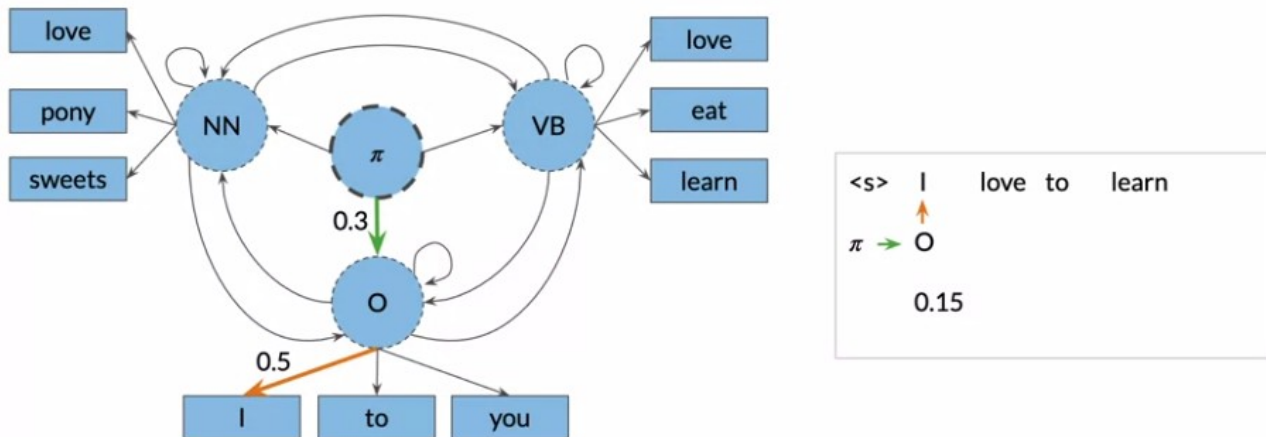A stochastic model describing a sequence of possible events.



|  | NN | VB | O |
|---|---|---|---|
| $\pi$ (initial) | 0.4 | 0.1 | 0.5 |
| NN (noun) | 0.2 | 0.2 | 0.6 |
| VB (verb) | 0.4 | 0.3 | 0.3 |
| O (other) | 0.2 | 0.3 | 0.5 |

$A =$

# The Viterbi Algorithm

A graph algorithm that finds the sequence of hidden states or parts of speech tags that have the highest probability for a sequence.



Viterbi algorithm – a graph algorithm

Given your transition and emission probabilities, we first populates and then use the auxiliary matrices C and D

matrix **C** holds the intermediate optimal probabilities.

matrix **D** holds the indices of the visited states (tags).

$$c_{i,j} = \max_{k} c_{k,j-1} * a_{k,i} * b_{i,cindex(w_j)}$$

$$d_{i,j} = \underset{k}{\mathrm{argmax}}\, c_{k,j-1} * a_{k,i} * b_{i,cindex(w_j)}$$

$C =$

|       | $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ |
|-------|-------|-------|-------|-------|-------|
| $t_1$ | 0.25  | 0.125 | 0.025 | 0.0125| 0.01  |
| $t_2$ | 0.1   | 0.025 | 0.05  | 0.01  | 0.003 |
| $t_3$ | 0.3   | 0.05  | 0.025 | 0.02  | 0.0000|
| $t_4$ | 0.2   | 0.1   | 0.000 | 0.0025| 0.0003|

C[1,5] = 0.01

D[1,5] = 1

$$s = \underset{i}{\mathrm{argmax}}\, c_{i,K} = 1$$

# Word Embeddings

## Basic Word Representations

1. **Integers**

   + Simple

   - Ordering: no semantic sense

2. **One-hot vectors**

   + Simple

   + No implied ordering

   - Takes a lot of time and space

   - No embedded meaning

3. **Word embedding vectors**

   + Low dimension

   - e.g. semantic distance: forest $\approx$ tree

   + Embed meaning

   - e.g. analogies: Paris:France :: Rome:?

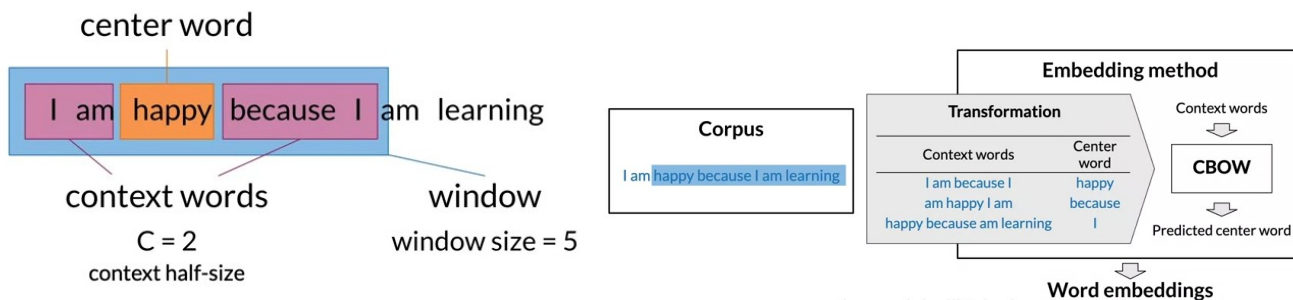## Basic word embedding methods

- **word2vec** (Google, 2013)
  - <u>Continuous bag-of-words (CBOW):</u> Which predict the missing word just giving the surround word.
  - <u>Continuous skip-gram (SGNS):</u> which does the reverse of the <u>CBOW</u> method, <u>SGNS</u> learns to predict the word surrounding a given input word.
- **Global Vectors (GloVe)** (Stanford, 2014)
- **FastText** (Facebook, 2016): based on the skip-gram model.

## Advanced word embedding methods

- **BERT** (Google, 2018)
- **ELMo** (Allen Institute for AI, 2018)
- **GPT-2** (OpenAI, 2018)
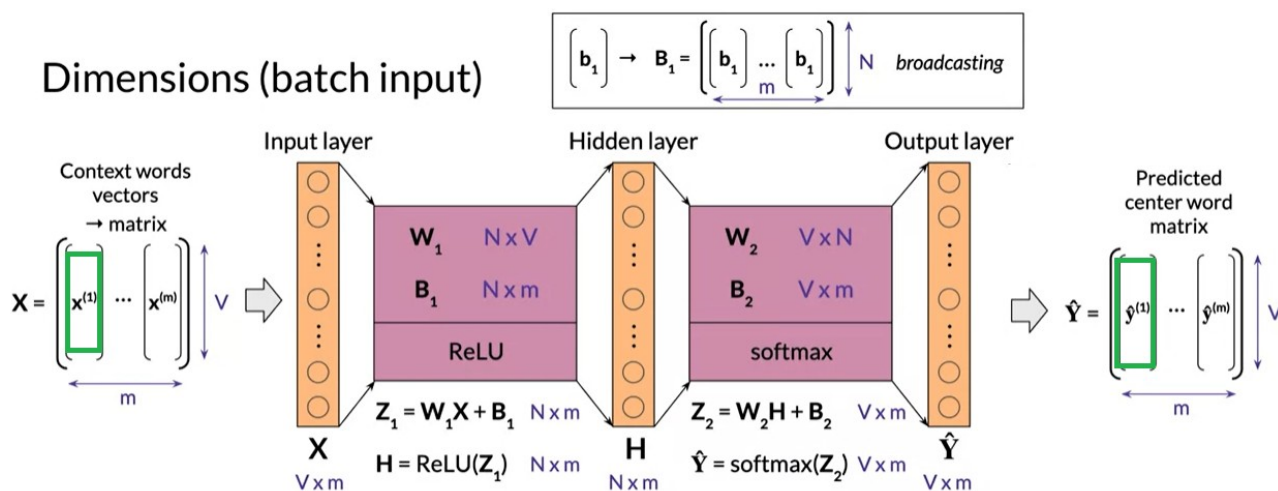
  Note: Tunable pre-trained models.

# Continuous Bag-of-Words Model



center word

I am **happy** because I am learning

context words   window
C = 2   window size = 5
context half-size

**Embedding method**

| Transformation | | |
|---|---|---|
| Context words | Center word | |
| I am because I | happy | |
| am happy I am | because | |
| happy because am learning | I | |

**Corpus**

I am happy because I am learning

Context words → CBOW → Predicted center word

Word embeddings

## Architecture

- CBOW model is based on the shallow dense neural network with an input layer, a single hidden layer, and output layer.
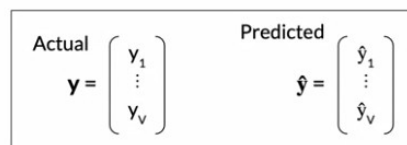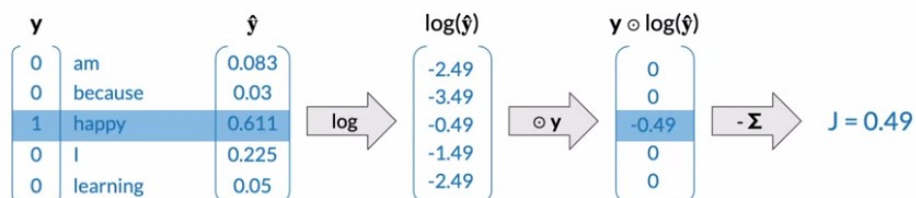


### Dimensions (batch input)

$$\begin{bmatrix} b_1 \end{bmatrix} \rightarrow B_1 = \begin{bmatrix} \begin{bmatrix} b_1 \end{bmatrix} \cdots \begin{bmatrix} b_1 \end{bmatrix} \end{bmatrix} \quad N \quad broadcasting$$

Context words vectors → matrix

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

Input layer

$W_1$  N x V
$B_1$  N x m
ReLU
$Z_1 = W_1 X + B_1$  N x m
$H = ReLU(Z_1)$  N x m
**X**  V x m

Hidden layer

$W_2$  V x N
$B_2$  V x m
softmax
$Z_2 = W_2 H + B_2$  V x m
$\hat{Y} = softmax(Z_2)$  V x m
**H**  N x m

Output layer

Predicted center word matrix

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

**Ŷ**  V x m

## Cost Function: Cross-entropy loss (log loss)



### Cross-entropy loss

$$J = -\sum_{k=1}^{V} y_k \log \hat{y}_k$$

Actual $y = \begin{bmatrix} y_1 \\ \vdots \\ y_V \end{bmatrix}$   Predicted $\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_V \end{bmatrix}$

I am happy because I am learning

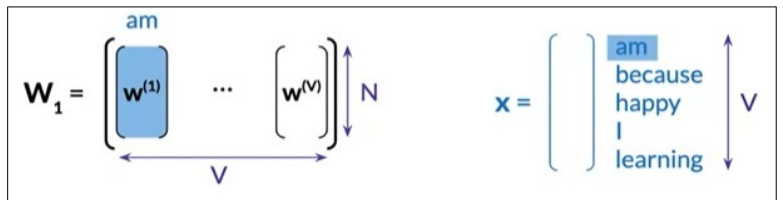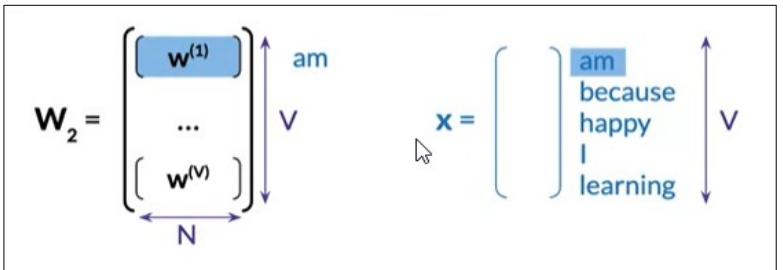| y | | ŷ | log(ŷ) | y ⊙ log(ŷ) |
|---|---|---|---|---|
| 0 | am | 0.083 | -2.49 | 0 |
| 0 | because | 0.03 | -3.49 | 0 |
| 1 | happy | 0.611 | -0.49 | -0.49 |
| 0 | I | 0.225 | -1.49 | 0 |
| 0 | learning | 0.05 | -2.49 | 0 |

log → ⊙y → -Σ → J = 0.49

# Extracting Word Embedding Vectors

After we have trained the neural network, we can extract three alternative word embedding representations.

**i.** Consider each column of W_1 as the column vector embedding vector of a word of the vocabulary



**ii.** Use each row of W_2 as the word embedding row vector for the corresponding word.



**iii.** Average W_1 and the transpose of W_2 to obtain W_3, a new n by v matrix.