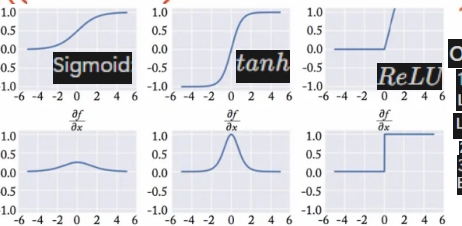
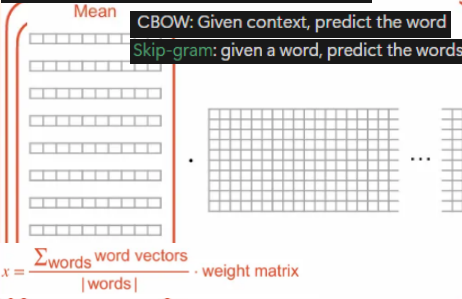


Word2vec algorithm — CBOW



| 方法    | 含义   |
|-------|--|
| IO    | I: Inside (Inside Entity), O: Outside (Outside Entity)             |
| BIO   | B: Begin, I: Inside, O: Outside                                    |
| BIOES | B: Begin, I: Inside, O: Outside, E: End, S: Single (singel Entity) |

micro-P  
Many/TN  
Possible outputs = |Labels| \* Length \* Length

5 keys components of NLP

- Data
  - not every data can be easily split
- Model
  - According to input get corresponding output
- Inference Method
  - find a high scoring option
- Metric
  - Loss function
- Learning Method
  - How to update model

context window  
"Small context (+/- 2 words) → more syntactically"  
"Larger context (+/- 5 words) → more semantically"

- Overfitting Solution
1. Modify loss with cost,
    - L1 driving unimportant feature weights to zero
    - L2 penalizes large weights; keep weights small
  2. during training, randomly set some weights to zero Dropout
  3. don't backpropagate error further once it gets small Early stopping

| Words      | IO Label | BIO Label | BIOES Label |
|------------|----------|-----------|-------------|
| Jane       | I-PER    | B-PER     | B-PER       |
| Villanueva | I-PER    | I-PER     | I-PER       |
| of         | O        | O         | O           |
| United     | I-ORG    | B-ORG     | B-ORG       |
| Airlines   | I-ORG    | I-ORG     | I-ORG       |
| Holding    | I-ORG    | I-ORG     | I-ORG       |
| discussed  | O        | O         | O           |
| the        | O        | O         | O           |
| Chicago    | I-LOC    | B-LOC     | B-LOC       |
| route      | O        | O         | O           |

If your goal is to maximize the total number of correct predictions, regardless of class distribution, then **micro** averaging is more appropriate.

If you are concerned about performance on rare or minority classes, then **macro** averaging is a better fit because it treats all classes equally!

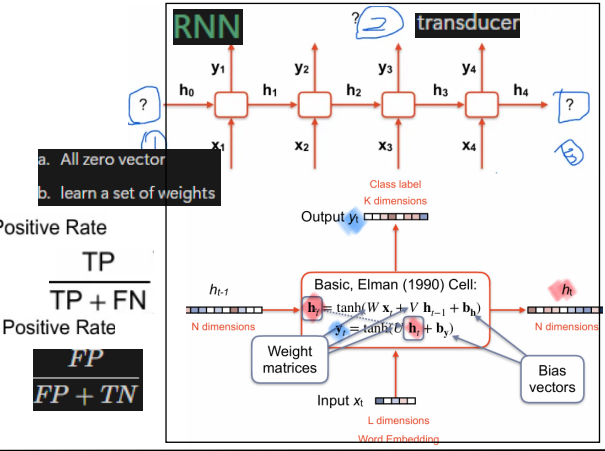
$$Y_{micro} = \frac{\sum TP}{\sum (TP + Y)} X_{macro} = \frac{\sum X}{\text{Number of classes}}$$

when there is no 'None' class  $P_{micro} = R_{micro}$

The closer the **ROC** curve is to the top-left corner (i.e., high TPR and low FPR), the better the model's performance.

**AUC** (Area Under the Curve) measures the model's ability to distinguish between positive and negative samples. A higher AUC indicates better separability.

The closer the Precision-Recall Curve (**PRC**) is to the top-right corner, the better the model's performance.



$n$  represents the number of input tokens  $k$  represents the number of labels Exhaustive —  $k^n$

**Greedy** —  $k * n$  Make choice one at a time. Sometimes the answer can match the Exhaustive method. For example, if every part of the output is independent

- Top-1** At each step, choose the highest scoring option
- Random** — full distribution  
At each step, choose using random sampling from the probability distribution
- Top-K**  
At each step, filter to the top  $K$  options, then choose using random sampling from the probability distribution.  
The probabilities need to be redistributed according to the proportion
- Top-P**  
At each step, filter to the options that cover  $P\%$  of the probability distribution, then choose using random sampling.
- Contrastive**  
At each step, adjust scores based on similarity with recent outputs, then choose the highest scoring option.

**Beam Search** —  $k * m * n$   
Beam Search maintains the top  $m$  highest-scoring sequences (called the beam width).  
For each of these  $m$  sequences, it explores all possible next tokens.  
It then selects the new top  $m$  sequences with the highest overall scores (including past steps).  
This process continues until the end of the sequence.

| Method             | Advantages   | Disadvantages  |
|--------------------|--|--|
| Top-1 (Greedy)     | - Deterministic: same input always gives the same output<br>- Fast and simple  | - Only considers locally optimal choices<br>- May miss better global sequences<br>- Tends to produce repetitive or generic output              |
| Random             | - Can escape local optima<br>- Encourages high diversity   | - May select low-probability, low-quality outputs<br>- Results often incoherent or meaningless   |
| Top-k              | - Avoids very unlikely tokens<br>- Good for problems where local choices don't lead to global optimum                                  | - Still randomly selects among top $k$ , may pick suboptimal tokens<br>- Total probability mass of top $k$ may be low, reducing output quality |
| Top-p (Nucleus)    | - Dynamically selects candidates based on cumulative probability $\geq p$<br>- Avoids small-probability regions more safely than Top-k | - More computation due to sorting and accumulation at each step<br>- Might exclude rare but meaningful tokens in skewed distributions          |
| Contrastive Search | - Balances diversity and coherence<br>- Reduces repetition<br>- Produces fluent and high-quality text                                  | - More complex to implement<br>- Depends heavily on model quality and contrastive scoring<br>- Computationally expensive                       |

How to select the  $m$  answer? Usually  $m$  is small enough that any of these are fine

Simple Approach —  $O(nm)$

- For each option, go through the list to find where it goes
- Heap Approach —  $O(n \log m)$
- Update it with each new item
- Quick Select —  $O(n)$

When candidates have very similar scores, Beam Search may struggle to differentiate

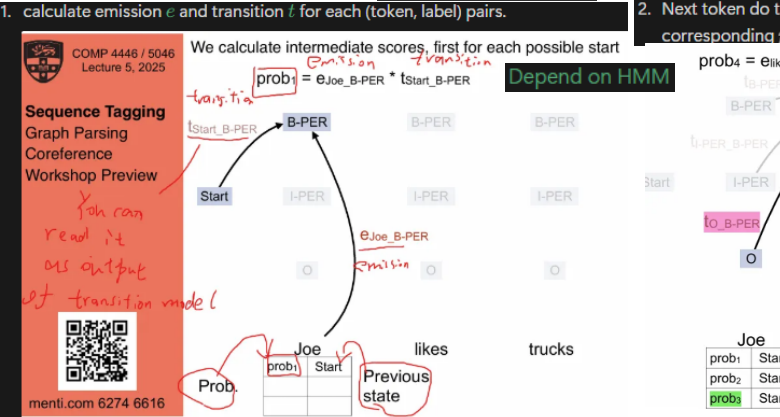
Increase beamwidth  $m$

Advantages of Viterbi Algorithm

- Finds the globally optimal sequence (maximum joint probability)
- Uses dynamic programming to efficiently explore all possible label sequences without enumerating

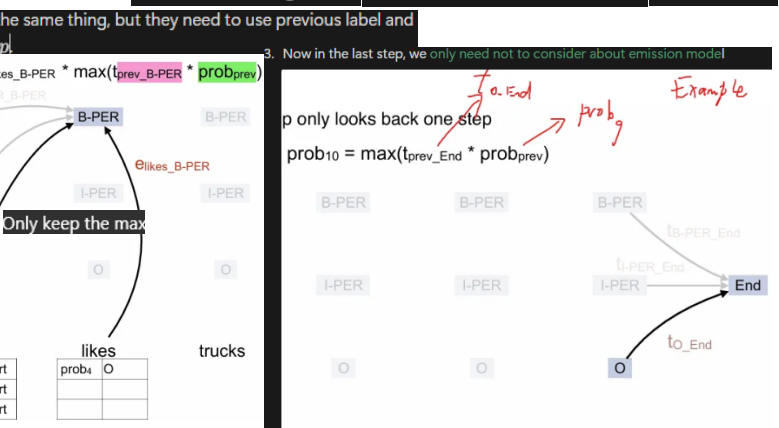
Viterbi algorithm Dynamic Programming 1 previous label  $O(|words| * |labels|^2)$

2 previous label  $O(|words| * |labels|^3)$  Emission Does not change



we can use in any model as long as it has previous input and current input

Linear model (e.g. CRF) Feedforward Networks RNNs Transformers



CKY algorithm

