

Precept 2: Word Embeddings

COS 484

Sijia Liu (slides borrowed from Tianyu Gao, Spring 2022 and Simon Park, Spring 2025)

2/6/2026

Today's Plan

1. Overview
2. Predict-based methods (count-based methods skipped)
3. Evaluation
4. Matrix Calculus
5. Exercises
6. Additional topic: Cursor

Overview

Overview - Word Embeddings

- Represent words as vectors
 - e.g., apple -> [0.1, 0.2, 0.5]
 - Encode semantic information
 - Useful for downstream NLP tasks

QUESTION: how can we get good word vectors

Overview - Distributional Hypothesis

- Words that occur in similar contexts have similar meaning
- EXAMPLE
 - A is the capital of ...
 - B is the capital of ...
- A, B should have similar meaning
- Word vectors for A, B should be “similar”

Overview - Different Approaches

- **Count-based** methods: PMI, ...
 - Use statistics
 - Covered in great details in lectures, we will skip today

- **Predict-based** methods: word2vec, GloVe, ...
 - Use ML

Predict-based Methods

Overview

- Learn an ML model
- Input: **corpus**, dictionary **V**, and desired dimension **d**
- Output: learned model parameters
 - **embedding vector** of dimension **d** for each word

Word2vec / skip-gram

- Learn an ML model
- Input: corpus, dictionary V , and desired dimension d
- Output: learned model parameters
 - **two embedding vectors** of dimension d for each word
 - \mathbf{u} when the word is a target word
 - \mathbf{v} when the word is a context word

Word2vec / skip-gram

- Learn an ML model
- Input: corpus, dictionary V , and desired dimension d
- Output: learned model parameters
 - **two embedding vectors** of dimension d for each word
 - \mathbf{u} when the word is a target word
 - \mathbf{v} when the word is a context word
 - **Quick poll: How many model parameters total?**

Word2vec / skip-gram

- Learn an ML model
- Input: corpus, dictionary V , and desired dimension d
- Output: learned model parameters
 - **two embedding vectors** of dimension d for each word
 - \mathbf{u} when the word is a target word
 - \mathbf{v} when the word is a context word
 - **Quick poll: How many model parameters total?** $2d |V|$

Word2vec / skip-gram: Learning Objective Part 1

- Input: target word \mathbf{t} , context word \mathbf{c}
- Current model parameters: \mathbf{u} , \mathbf{v} for each word in dictionary
- Model's behavior:
 - Predict $\mathbb{P}[\mathbf{c} | \mathbf{t}]$
 - Probability that \mathbf{c} (out of all possible words) is in context of \mathbf{t}
 - How? Compute logits $\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}'}$ for all possible context words
 - Normalize with softmax function

- $$\mathbb{P}[\mathbf{c} | \mathbf{t}] = \frac{\exp(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\sum_{\mathbf{c}'} \exp(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}'})}$$

Word2vec / skip-gram: Learning Objective Part 1

- Input: target word \mathbf{t} , context word \mathbf{c}
- Current model parameters: \mathbf{u} , \mathbf{v} for each word in dictionary
- Model's behavior:

- $$\mathbb{P}[\mathbf{c} | \mathbf{t}] = \frac{\exp(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\sum_{\mathbf{c}'} \exp(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}'})}$$

- Loss = cross-entropy / negative log likelihood
- $L_{\mathbf{t},\mathbf{c}} = -\log \mathbb{P}[\mathbf{c} | \mathbf{t}]$

Word2vec / skip-gram: Learning Objective Part 2

- Input: a sequence of words w_1, w_2, \dots, w_T
- Model's behavior:
 - Choose a context window size m
 - Choose a word w_t and consider it a target word
 - For each word $w_{t-m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m}$ consider it a context word
 - Compute loss $L_{w_t, w_{t+j}}$ for all $-m \leq j \leq m, j \neq 0$
 - Sum it all up $L_t = \sum_{-m \leq j \leq m, j \neq 0} L_{w_t, w_{t+j}}$
- Model's loss for predicting context words of a particular word w_t

Word2vec / skip-gram: Learning Objective Part 3

- Input: a sequence of words w_1, w_2, \dots, w_T

- Model's behavior:

- $$L_t = \sum_{-m \leq j \leq m, j \neq 0} L_{w_t, w_{t+j}}$$

- Model's loss for predicting context words of a particular word w_t

- Take the average to get the final loss
$$L = \frac{1}{T} \sum_{t=1}^T L_t$$

Word2vec / skip-gram: Learning Objective Part 3

- Input: a sequence of words w_1, w_2, \dots, w_T

- Model's behavior:

- $$L = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{\mathbf{c} \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{\mathbf{c}})}$$

- Note $\mathbf{c} \in V$ takes the sum over all possible words in the dictionary
- Not just the words that appear in the corpus / sequence of words

Word2vec / skip-gram: How to Train

- Input: a sequence of words w_1, w_2, \dots, w_T

- Compute loss:

$$L = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{c \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_c)}$$

- Compute gradient of the loss with respect to all $\mathbf{u}, \mathbf{v} : \frac{\partial L}{\partial \mathbf{u}}, \frac{\partial L}{\partial \mathbf{v}}$
- Update model parameters via Gradient Descent
- Problem?

Word2vec / skip-gram: How to Train

- Input: a sequence of words w_1, w_2, \dots, w_T

- Compute loss:

$$L = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{\mathbf{c} \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{\mathbf{c}})}$$

- Compute gradient of the loss with respect to all $\mathbf{u}, \mathbf{v} : \frac{\partial L}{\partial \mathbf{u}}, \frac{\partial L}{\partial \mathbf{v}}$
- Update model parameters via Gradient Descent
- Problem? For every pair (\mathbf{t}, \mathbf{c}) , you need to update $2d \mid V \mid$ parameters
- ONE SOLUTION: instead of $\mathbf{c} \in V$, only sample K (5-20) alternatives

Word2vec / skip-gram: Learning Objective Part 1 (recap)

- Input: target word \mathbf{t} , context word \mathbf{c}
- Current model parameters: \mathbf{u} , \mathbf{v} for each word in dictionary
- Model's behavior:
 - Predict $\mathbb{P}[\mathbf{c} | \mathbf{t}]$
 - Probability that \mathbf{c} (out of all possible words) is in context of \mathbf{t}
 - How? Compute logits $\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}'}$ for all possible context words
 - Normalize with softmax function

- $$\mathbb{P}[\mathbf{c} | \mathbf{t}] = \frac{\exp(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\sum_{\mathbf{c}'} \exp(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}'})}$$

Word2vec / skip-gram: Negative Sampling Part 1

- Input: target word \mathbf{t} , context word \mathbf{c} , **alternative context words** $\mathbf{c}_1, \dots, \mathbf{c}_K$
- Current model parameters: \mathbf{u}, \mathbf{v} for each word in dictionary
- Model's behavior:

- Predict $\mathbb{P}[+ | \mathbf{t}, \mathbf{c}] \cdot \prod_{i=1}^K \mathbb{P}[- | \mathbf{t}, \mathbf{c}_i]$

- Probability that \mathbf{c} is in context of \mathbf{t} **AND** $\mathbf{c}_1, \dots, \mathbf{c}_K$ are not in context of \mathbf{t}
- How? Consider these **independent** binary logistic regression
- $\mathbb{P}[+ | \mathbf{t}, \mathbf{c}] = \sigma(\mathbf{u}_t \cdot \mathbf{v}_c)$
- $\mathbb{P}[- | \mathbf{t}, \mathbf{c}_i] = \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$

Word2vec / skip-gram: Negative Sampling Part 1

- Input: target word \mathbf{t} , context word \mathbf{c} , **alternative context words** $\mathbf{c}_1, \dots, \mathbf{c}_K$
- Current model parameters: \mathbf{u}, \mathbf{v} for each word in dictionary
- Model's behavior:

- $\mathbb{P}[+ | \mathbf{t}, \mathbf{c}] \cdot \prod_{i=1}^K \mathbb{P}[- | \mathbf{t}, \mathbf{c}_i] = \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) \cdot \prod_{i=1}^K \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$

- Loss = cross-entropy / negative log likelihood

- $L_{\mathbf{t}, \mathbf{c}} = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^K \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i}))$

- Recall that $\mathbf{c}_1, \dots, \mathbf{c}_K$ randomly sampled

Word2vec / skip-gram: Negative Sampling Part 1

- Input: target word \mathbf{t} , context word \mathbf{c} , **alternative context words** $\mathbf{c}_1, \dots, \mathbf{c}_K$
- Current model parameters: \mathbf{u}, \mathbf{v} for each word in dictionary
- Model's behavior:

- $\mathbb{P}[+ | \mathbf{t}, \mathbf{c}] \cdot \prod_{i=1}^K \mathbb{P}[- | \mathbf{t}, \mathbf{c}_i] = \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) \cdot \prod_{i=1}^K \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$

- Loss = cross-entropy / negative log likelihood

- $L_{\mathbf{t}, \mathbf{c}} = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^K \mathbb{E}_{\mathbf{c}_i \sim V} \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i}))$

- **Q: why do we need to take the expectation?**

Word2vec / skip-gram: Negative Sampling Part 1

- Input: target word \mathbf{t} , context word \mathbf{c} , **alternative context words** $\mathbf{c}_1, \dots, \mathbf{c}_K$
- Current model parameters: \mathbf{u}, \mathbf{v} for each word in dictionary
- Model's behavior:

- $\mathbb{P}[+ | \mathbf{t}, \mathbf{c}] \cdot \prod_{i=1}^K \mathbb{P}[- | \mathbf{t}, \mathbf{c}_i] = \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) \cdot \prod_{i=1}^K \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$

- Loss = cross-entropy / negative log likelihood

- $L_{\mathbf{t}, \mathbf{c}} = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^K \mathbb{E}_{\mathbf{c}_i \sim V} \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i}))$

- **Q: why do we need to take the expectation?**
- $\mathbf{c}_1, \dots, \mathbf{c}_K$ are random variables so the second term is a Monte-Carlo estimate using K alternative samples rather than the ground truth negative set.
 - In practice, you can just compute with K random samples without bootstrapping

Word2vec / skip-gram: Negative Sampling Part 1

- Input: target word \mathbf{t} , context word \mathbf{c} , **alternative context words** $\mathbf{c}_1, \dots, \mathbf{c}_K$
- Current model parameters: \mathbf{u}, \mathbf{v} for each word in dictionary
- Model's behavior:

- $\mathbb{P}[+ | \mathbf{t}, \mathbf{c}] \cdot \prod_{i=1}^K \mathbb{P}[- | \mathbf{t}, \mathbf{c}_i] = \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) \cdot \prod_{i=1}^K \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$

- Loss = cross-entropy / negative log likelihood

- $L_{\mathbf{t}, \mathbf{c}} = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^K \mathbb{E}_{\mathbf{c}_i \sim V} \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i}))$

- **Quick Poll: How many parameters contribute to this loss value?**

Word2vec / skip-gram: Negative Sampling Part 1

- Input: target word \mathbf{t} , context word \mathbf{c} , **alternative context words** $\mathbf{c}_1, \dots, \mathbf{c}_K$
- Current model parameters: \mathbf{u}, \mathbf{v} for each word in dictionary
- Model's behavior:

- $\mathbb{P}[+ | \mathbf{t}, \mathbf{c}] \cdot \prod_{i=1}^K \mathbb{P}[- | \mathbf{t}, \mathbf{c}_i] = \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) \cdot \prod_{i=1}^K \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$

- Loss = cross-entropy / negative log likelihood

- $L_{t,c} = -\log \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) - \sum_{i=1}^K \mathbb{E}_{\mathbf{c}_i \sim V} \log \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$

- **Quick Poll: How many parameters contribute to this loss value?**

- $(K + 2)d$ (d for each of $\mathbf{u}_t, \mathbf{v}_c, \mathbf{v}_{c_1}, \mathbf{v}_{c_2}, \dots, \mathbf{v}_{c_K}$) (much more efficient compared to $2d |V|$ w/o negative sampling)

Word2vec / skip-gram: Train with Negative Sampling

- Input: a sequence of words w_1, w_2, \dots, w_T

- Compute loss:

$$L = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} L_{w_t, w_{t+j}}$$

- Compute gradient of the loss with respect to all $\mathbf{u}, \mathbf{v} : \frac{\partial L}{\partial \mathbf{u}}, \frac{\partial L}{\partial \mathbf{v}}$
- Update model parameters via Gradient Descent

Evaluation

Overview

- Recall earlier QUESTION: **how can we get good word vectors**
- How do we know if we got **good** word vectors?

- Intrinsic Evaluation
 - evaluate word vectors directly
 - “similarity” based tasks

- Extrinsic Evaluation:
 - evaluate ML model built on top of the word vectors
 - other tasks

Matrix Calculus

Matrix Calculus / Vectorized Gradients

- Go through this note:
- <http://web.stanford.edu/class/cs224n/readings/gradient-notes.pdf>
- Make sure that you can understand all the cases in Section 2, 3

Basic Definition

$$f: \mathbb{R}^n \longrightarrow \mathbb{R}^m$$

$$f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))$$

$$\frac{\partial f}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Gradient w.r.t vector - matrix * vector

Can directly take gradients of a vector with respect to a vector

$$\mathbf{x} \in \mathbb{R}^n \quad \mathbf{z} \in \mathbb{R}^m \quad \mathbf{W} \in \mathbb{R}^{m \times n}$$

$$\mathbf{z} = \mathbf{W}\mathbf{x}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \mathbf{W}$$

Gradient w.r.t vector - vector * matrix

Can directly take gradients of a vector with respect to a vector

$$\mathbf{x} \in \mathbb{R}^n \quad \mathbf{z} \in \mathbb{R}^m \quad \mathbf{W} \in \mathbb{R}^{n \times m}$$

$$\mathbf{z} = \mathbf{x}\mathbf{W}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \mathbf{W}^T$$

Gradient w.r.t vector - elementwise operation

Can directly take gradients of a vector with respect to a vector

$$\mathbf{x} \in \mathbb{R}^n \quad \mathbf{z} \in \mathbb{R}^n \quad \mathbf{f} \in \mathbb{R}^n \longrightarrow \mathbb{R}^n$$

$$\mathbf{z} = f(\mathbf{x}) = (f(x_1), \dots, f(x_n))$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \text{diag}(f'(\mathbf{x})) = \begin{bmatrix} f'(x_1) & 0 & \dots & 0 \\ 0 & f'(x_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f'(x_n) \end{bmatrix}$$

Gradient w.r.t matrix - matrix * vector

Cannot directly take gradients of a vector with respect to a matrix

Need to take gradient of the final loss value (scalar) with respect to matrix

$$\mathbf{x} \in \mathbb{R}^n \quad \mathbf{z} \in \mathbb{R}^m \quad \mathbf{W} \in \mathbb{R}^{m \times n} \quad L \in \mathbb{R}$$

$$\mathbf{z} = \mathbf{W}\mathbf{x} \quad L = f(\mathbf{z})$$

$$\frac{\partial L}{\partial \mathbf{W}} = \left(\frac{\partial L}{\partial \mathbf{z}} \right)^T \mathbf{x}^T$$

Gradient w.r.t matrix - vector * matrix

Cannot directly take gradients of a vector with respect to a matrix

Need to take gradient of the final loss value (scalar) with respect to matrix

$$\mathbf{x} \in \mathbb{R}^n \quad \mathbf{z} \in \mathbb{R}^m \quad \mathbf{W} \in \mathbb{R}^{n \times m} \quad L \in \mathbb{R}$$

$$\mathbf{z} = \mathbf{x}\mathbf{W} \quad L = f(\mathbf{z})$$

$$\frac{\partial L}{\partial \mathbf{W}} = \mathbf{x}^T \left(\frac{\partial L}{\partial \mathbf{z}} \right)$$

Exercises

Q: Gradients for Skip-gram with Negative Sampling

Recall: loss for target \mathbf{t} , context \mathbf{c} , alternative context $\mathbf{c}_1, \dots, \mathbf{c}_K$

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

(a) $\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{u}_{\mathbf{t}}} = ?$

(b) $\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}}} = ?$

(c) $\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}_j}} = ?$

A particular index j (not the same index i in the sum)

Q: (a)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{u}_{\mathbf{t}}} =$$

Q: (a)

$$L_{t,c} = -\log \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$$

$$\frac{\partial L_{t,c}}{\partial \mathbf{u}_t} = -\frac{\frac{\partial \sigma(\mathbf{u}_t \cdot \mathbf{v}_c)}{\partial \mathbf{u}_t}}{\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)} - \sum_{i=1}^K \frac{\frac{\partial \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})}{\partial \mathbf{u}_t}}{\sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})}$$

1. $d \log(x) / dx = 1/x$
2. chain rule

Q: (a)

$$L_{t,c} = -\log \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$$

$$\frac{\partial L_{t,c}}{\partial \mathbf{u}_t} = -\frac{\frac{\partial \sigma(\mathbf{u}_t \cdot \mathbf{v}_c)}{\partial \mathbf{u}_t}}{\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)} - \sum_{i=1}^K \frac{\frac{\partial \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})}{\partial \mathbf{u}_t}}{\sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})}$$

1. $d \text{ sigmoid}(x) / dx = \text{sig}(x) * (1 - \text{sig}(x))$
2. chain rule

$$= -\frac{\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)(1 - \sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) \frac{\partial(\mathbf{u}_t \cdot \mathbf{v}_c)}{\partial \mathbf{u}_t}}{\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)} - \sum_{i=1}^K \frac{\sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})(1 - \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})) \frac{\partial(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})}{\partial \mathbf{u}_t}}{\sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})}$$

Q: (a)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{u}_{\mathbf{t}}} = -\frac{\frac{\partial \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\frac{\partial \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

d (x * y) / dx = y

$$= -\frac{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})(1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})) \frac{\partial(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})(1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})) \frac{\partial(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

$$= -\frac{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})(1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}))\mathbf{v}_{\mathbf{c}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})(1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i}))(-\mathbf{v}_{\mathbf{c}_i})}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

Q: (a)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{u}_{\mathbf{t}}} = -\frac{\frac{\partial \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\frac{\partial \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

$$= -\frac{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})(1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})) \frac{\partial(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})(1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})) \frac{\partial(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

$$= -\frac{\cancel{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}(1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}))\mathbf{v}_{\mathbf{c}}}{\cancel{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}} - \sum_{i=1}^K \frac{\cancel{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}(1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i}))(-\mathbf{v}_{\mathbf{c}_i})}{\cancel{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}}$$

Q: (a)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{u}_{\mathbf{t}}} = -\frac{\frac{\partial \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\frac{\partial \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

$$= -\frac{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})(1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})) \frac{\partial(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})(1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})) \frac{\partial(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

$$= -(1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}))\mathbf{v}_{\mathbf{c}} - \sum_{i=1}^K (1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i}))(-\mathbf{v}_{\mathbf{c}_i})$$

Q: (a)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{u}_{\mathbf{t}}} = -\frac{\frac{\partial \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\frac{\partial \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

$$= -\frac{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})(1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}))\frac{\partial(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})} - \sum_{i=1}^K \frac{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})(1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i}))\frac{\partial(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}{\partial \mathbf{u}_{\mathbf{t}}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})}$$

$$= -(1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}))\mathbf{v}_{\mathbf{c}} - \sum_{i=1}^K (1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i}))(-\mathbf{v}_{\mathbf{c}_i})$$

$$= (\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - 1)\mathbf{v}_{\mathbf{c}} + \sum_{i=1}^K \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})\mathbf{v}_{\mathbf{c}_i}$$

sigmoid(-x) = 1 - sigmoid(x)

Q: (a) FINAL VERSION

$$\frac{\partial L_{t,c}}{\partial \mathbf{u}_t} = (\sigma(\mathbf{u}_t \cdot \mathbf{v}_c) - 1)\mathbf{v}_c + \sum_{i=1}^K \sigma(\mathbf{u}_t \cdot \mathbf{v}_{c_i})\mathbf{v}_{c_i}$$

Q: (b)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}}} =$$

Q: (b)

$$L_{t,c} = -\log \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c_i})$$

$$\frac{\partial L_{t,c}}{\partial \mathbf{v}_c} = -\frac{\frac{\partial \sigma(\mathbf{u}_t \cdot \mathbf{v}_c)}{\partial \mathbf{v}_c}}{\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)}$$

Summands do not depend on \mathbf{v}_c

Q: (b)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}}} = - \frac{\frac{\partial \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}{\partial \mathbf{v}_{\mathbf{c}}}}{\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})}$$

$$= - (1 - \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}})) \mathbf{u}_{\mathbf{t}}$$

$$= (\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - 1) \mathbf{u}_{\mathbf{t}}$$

Same as before

Q: (b) FINAL VERSION

$$\frac{\partial L_{t,c}}{\partial \mathbf{v}_c} = (\sigma(\mathbf{u}_t \cdot \mathbf{v}_c) - 1)\mathbf{u}_t$$

Q: (c)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}_j}} =$$

Q: (c)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}_j}} = - \frac{\frac{\partial \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_j})}{\partial \mathbf{v}_{\mathbf{c}_j}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_j})}$$

All other summands do not depend on j

Q: (c)

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$\frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}_j}} = - \frac{\frac{\partial \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_j})}{\partial \mathbf{v}_{\mathbf{c}_j}}}{\sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_j})}$$

$$= - (1 - \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_j}))(-\mathbf{u}_{\mathbf{t}})$$

$$= \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_j})\mathbf{u}_{\mathbf{t}}$$

Same as before

Q: (c) FINAL VERSION

$$\frac{\partial L_{t,c}}{\partial \mathbf{v}_{c_j}} = \sigma(\mathbf{u}_t \cdot \mathbf{v}_{c_j}) \mathbf{u}_t$$

Q: FINAL VERSION

Recall: loss for target \mathbf{t} , context \mathbf{c} , **alternative context** $\mathbf{c}_1, \dots, \mathbf{c}_K$

$$L_{\mathbf{t},\mathbf{c}} = -\log \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - \sum_{i=1}^K \log \sigma(-\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})$$

$$(a) \frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{u}_{\mathbf{t}}} = (\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - 1)\mathbf{v}_{\mathbf{c}} + \sum_{i=1}^K \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_i})\mathbf{v}_{\mathbf{c}_i}$$

$$(b) \frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}}} = (\sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}}) - 1)\mathbf{u}_{\mathbf{t}}$$

$$(c) \frac{\partial L_{\mathbf{t},\mathbf{c}}}{\partial \mathbf{v}_{\mathbf{c}_j}} = \sigma(\mathbf{u}_{\mathbf{t}} \cdot \mathbf{v}_{\mathbf{c}_j})\mathbf{u}_{\mathbf{t}}$$

Additional topic: Cursor Demo

- What's cursor?
- A very powerful coding agent!
- Now it's free for students for 1 year: <https://cursor.com/students>

Additional topic: Cursor Demo

- **Key concepts:**

- Tab
- Agent (cmd + I)
- Inline Edit (cmd + K)
- Chat (cmd + N)
- Rules
- Semantic search
- MCP
- Context
- Models

Additional topic: Cursor Demo

- **Best practices with cursor agents**
 - **Use plan mode before coding**
 - **Let the agent find the context**
 - **Start with tests and use as feedback loop**
 - **Customize with rules and skills**
 - Tip: Start simple! Add rules only when you notice the agent making the same mistake repeatedly.
 - **Use debug mode for tricky bugs**
- See more in <https://cursor.com/blog/agent-best-practices>