

Generative Models for NLP

Language Models

Nadi Tomeh

24/1/25

Navigation icons: back, forward, search, etc.

Nadi Tomeh

Generative Models for NLP

24/1/25

1 / 84

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Navigation icons: back, forward, search, etc.

Nadi Tomeh

Generative Models for NLP

24/1/25

2 / 84

Notes

Outline

- Introduction to Language Models
 - Vocabulary and Tokenization
 - Applications
 - Parametrization and Estimation
 - n-Gram Language Models
 - Addressing Data Sparsity in n-Gram Models
 - Evaluation Metrics for Language Models
 - Toy Example
 - Generation Strategies for Language Models
 - Feed-Forward Neural Language Models
 - Training
 - Recurrent Neural Networks (RNNs)
 - LSTMs and GRUs
 - What's Next?

Notes

What is a Language Model?

Definition (Language Model)

A **language model** is a function that defines a joint probability distribution $p(w_1, w_2, \dots, w_n)$, over an ordered sequence of tokens $\mathbf{w} = (w_1, w_2, \dots, w_n)$. Each $w_k \in \mathcal{V}$, a finite set of tokens called the **vocabulary**. A valid language model must satisfy the constraint:

$$\sum_{\mathbf{w} \in \mathcal{W}} p(\mathbf{w}) = 1,$$

where $\mathcal{W} \subseteq \mathcal{V}^*$ is the (possibly infinite) set of all token sequences (recall the definition of a **formal language**).

Chain Rule of Probability

Using the chain rule, we can factorize this joint probability as:

$$p(w_1, w_2, \dots, w_n) = \prod_{k=1}^n p(w_k \mid w_1, \dots, w_{k-1}).$$

Each term $p(w_k \mid w_1, \dots, w_{k-1})$ is a conditional probability of the current token given all previous tokens.

- **Interpretation:**
 - The model measures how “natural” or likely a sequence is.
 - Each factor $p(w_k \mid w_1, \dots, w_{k-1})$ represents how likely the next token w_k is given the context (w_1, \dots, w_{k-1}) .

Notes

Generative vs. Discriminative Models: Basic Concepts

Generative Models

- A **generative model** aims to learn the joint probability $p(\mathbf{x}, \mathbf{y})$, where \mathbf{x} represents the observed data (e.g., a sequence of tokens) and \mathbf{y} represents labels, latent variables, or outputs (can be *structured*).
- A **language model** is generative because it learns $p(w_1, \dots, w_n)$, i.e., the probability of entire sequences.
 - Once a generative model is learned, you can derive $p(\mathbf{x} \mid \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})}$, and $p(\mathbf{y} \mid \mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}$.
 - The marginal probability of \mathbf{x} , necessary for computing $p(\mathbf{y} \mid \mathbf{x})$, is obtained by summing (or integrating) over all possible values of \mathbf{y} :

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) \quad (\text{discrete case})$$

Discriminative Models

- A **discriminative model** directly learns the conditional probability $p(\mathbf{y} \mid \mathbf{x})$, without modeling the joint distribution $p(\mathbf{x}, \mathbf{y})$ or the data likelihood $p(\mathbf{x})$.
- A discriminative **text classifier** takes an input sequence $\mathbf{x} = (w_1, w_2, \dots, w_n)$ and predict a class label y (e.g., *positive* or *negative* sentiment) and directly models $p(y \mid \mathbf{x})$.

Notes

Outline

- Introduction to Language Models
- **Vocabulary and Tokenization**
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Tokenization: Definitions and Approaches

What is Tokenization?

- **Tokenization** is the process of splitting text into smaller units, called *tokens*, which serve as the atomic input to language models.
- A *token* can be:
 - A full word
 - A subword or morpheme
 - A single character (especially in low-resource or highly morphologically rich languages)

Key Considerations

- **Vocabulary Size:**
 - Large vocabulary \implies fewer unknowns or Out-Of-Vocabulary (OOV) tokens, but increases parameter count.
 - Small vocabulary \implies risk of high OOV rates, or reliance on subword tokens.
- **Handling Unknown Words:**
 - Use a special `<unk>` token, or fallback to character-level tokens.
- **Granularity:**
 - **Word-Level:** Simplest, but OOV issues can be severe.
 - **Subword-Level** (BPE, WordPiece, SentencePiece): Balances coverage and vocabulary size.
 - **Character-Level:** No OOVs, but leads to longer sequences and sometimes slower training.

Notes

Byte-Pair Encoding (BPE): Algorithm and Vocabulary Evolution I

Core Idea of BPE

- **Byte-Pair Encoding (BPE)** is a data compression technique adapted for tokenization.
- Iteratively merges the most frequent pair of symbols (characters or subwords) into a single token.
- Produces a subword-based vocabulary that reduces out-of-vocabulary issues while controlling vocabulary size.

Algorithm (High-Level Steps)

1. **Initialize Vocabulary \mathcal{V}_0 :**
 - Each unique character is its own token (e.g., `l`, `o`, `v`, `e`, `c`, `a`, `t`, `s`, plus any spaces or special markers).
2. **Count Pair Frequencies:**
 - Scan the training text for adjacent token pairs (e.g., `l+o`, `o+v`, etc.).
3. **Merge Most Frequent Pair:**
 - Combine that pair into a single token (e.g., `o_v`).
 - Update your text (i.e., each occurrence of `o v` becomes the new merged token).
 - Add this newly merged token to your vocabulary \mathcal{V}_1 .
4. **Repeat** for n merges or until desired vocabulary size is reached.

Notes

Byte-Pair Encoding (BPE): Algorithm and Vocabulary Evolution II

Vocabulary & Text Evolution (Simplified Example)

Training Text (repeated twice): i love love cats

Initial vocabulary \mathcal{V}_0 (characters only):

$$\mathcal{V}_0 = \{i, l, o, v, e, c, a, t, s\}$$

Step 1: Most frequent adjacent pair is l + o.

$$\text{Merge } (l, o) \rightarrow l_o. \quad \mathcal{V}_1 = \mathcal{V}_0 \cup \{l_o\}.$$

Text now becomes: i l_o v e l_o v e c a t s

Step 2: Next frequent pair might be l_o + v.

$$\text{Merge } (l_o, v) \rightarrow l_o_v. \quad \mathcal{V}_2 = \mathcal{V}_1 \cup \{l_o_v\}.$$

Text now becomes: i l_o_v e l_o_v e c a t s

Step 3: Merge l_o_v + e to form l_o_v_e, etc.

Over several merges, common subwords like cat or love end up as single tokens.

Final Vocabulary \mathcal{V}_n (after n merges):

$$\mathcal{V}_n = \{i, l_o_v_e, c_a_t_s, \dots\}$$

Why BPE is Useful

- **Reduced OOVs:** Rare words can be decomposed into known subwords (e.g., homework \rightarrow home + work).
- **Tunable Vocab Size:** Stop merges early for a smaller vocab, or merge more pairs for fewer but larger tokens.
- **Empirical Success:** Widely adopted in modern NLP (e.g., GPT, RoBERTa) for balancing coverage and memory footprint.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Applications

Evaluating Text Likelihood

- Given a sequence $\mathbf{w} = (w_1, w_2, \dots, w_n)$, compute its probability: $p(\mathbf{w}) = \prod_{k=1}^n p(w_k \mid w_1, \dots, w_{k-1})$.
- **Use Cases:**
 - **Speech Recognition & Machine Translation:** Re-rank candidate outputs based on their probabilities.
 - **Error Correction:** Identify unlikely sequences as potential errors.
 - **Quality Assessment:** Evaluate fluency and coherence of text in various applications.

Text Generation

- **Next-Token Prediction:** Iteratively extend the sequence $(w_1, w_2, \dots, w_{k-1}) \rightarrow (w_1, w_2, \dots, w_{k-1}, w_k)$ until a stopping criterion is met.
- Used for dialogue systems, creative content creation, and auto-completion.

Generalization

Used for modeling any kind of sequences: code, time series, etc.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- **Parametrization and Estimation**
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Parameterization of Language Models

Parameterized Probability

Rather than directly specifying $p(w_k \mid w_1, \dots, w_{k-1})$, language models introduce a set of parameters θ to define:

$$p_{\theta}(w_k \mid w_1, \dots, w_{k-1}).$$

- The joint probability over a sequence is then parameterized as:

$$p_{\theta}(w_1, \dots, w_n) = \prod_{k=1}^n p_{\theta}(w_k \mid w_1, \dots, w_{k-1}).$$

- Parameterization allows using various model architectures:
 - **n-gram** models: Use fixed-context frequency counts with parameters derived from observed counts.
 - **Neural networks**: Use parameters θ to encode complex dependencies (e.g., in RNNs, Transformers).
- The goal: Find θ that best captures the underlying language patterns.

Notes

Estimation of Parameters from Data

Statistical Estimation

In **probability theory**, the distribution $p(x)$ is assumed known, and we derive properties (e.g., expectations, variances) from that distribution. In **statistics**, the distribution is *unknown*, and we **estimate** its parameters or form based on observed data \mathcal{D} .

Maximum Likelihood Estimation (MLE)

Given a training corpus $\mathcal{D} = \{\mathbf{w}^{(i)}\}_{i=1}^N$, estimate parameters by maximizing the likelihood:

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \prod_{i=1}^N p_{\theta}(\mathbf{w}^{(i)}).$$

Equivalently, maximize the log-likelihood:

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \sum_{i=1}^N \log p_{\theta}(\mathbf{w}^{(i)}) = \arg \max_{\theta} \sum_{i=1}^N \sum_{k=1}^{n^{(i)}} \log p_{\theta}(w_k^{(i)} \mid w_1^{(i)}, \dots, w_{k-1}^{(i)}).$$

where $n^{(i)}$ is the length of sequence i . Optimization is typically performed using gradient-based methods (e.g., stochastic gradient descent) and backpropagation for neural models.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- **n-Gram Language Models**
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Markov Assumption in Language Modeling

Full Conditional Probability

Recall the chain rule for a sequence $\mathbf{w} = (w_1, w_2, \dots, w_n)$:

$$p(\mathbf{w}) = \prod_{k=1}^n p(w_k \mid w_1, \dots, w_{k-1}).$$

Markov Assumption

The **Markov assumption** simplifies this by assuming that the probability of the next token depends only on a finite history of previous tokens:

$$p(w_k \mid w_1, \dots, w_{k-1}) \approx p(w_k \mid w_{k-n+1}, \dots, w_{k-1}),$$

where n is the **order** of the Markov model.

- This *finite memory* assumption reduces computational complexity and makes estimation from data feasible.
- It introduces conditional independence: w_k is independent of tokens beyond the last $n - 1$ given the recent history.
- Leads directly to **n-gram models**, where probabilities are estimated based on limited context of length $n - 1$.

Notes

Parametrization of n -Gram Models Using Categorical Distributions

Parametrization

- For each possible $(n - 1)$ -gram context $\mathbf{c} = (w_{k-n+1}, \dots, w_{k-1})$, define a **categorical distribution**:

$$p_{\theta}(w_k \mid \mathbf{c}) = \theta_{\mathbf{c}, w_k}, \quad \text{where} \quad \sum_{w_k \in \mathcal{V}} \theta_{\mathbf{c}, w_k} = 1.$$

- $\theta_{\mathbf{c}, w_k}$ represents the probability of observing w_k given the history \mathbf{c} .
- For each context \mathbf{c} , the model stores a parameter vector:

$$\theta_{\mathbf{c}} = (\theta_{\mathbf{c}, w_1}, \theta_{\mathbf{c}, w_2}, \dots, \theta_{\mathbf{c}, w_{|\mathcal{V}|}}),$$

which lies in the $|\mathcal{V}|$ -dimensional **probability simplex**.

Notes

Parameters of n -Gram Models

Number of Parameters

- Total parameters:

$$|\mathcal{V}|^{n-1} \cdot (|\mathcal{V}| - 1),$$

where:

- $|\mathcal{V}|^{n-1}$: Number of possible $(n - 1)$ -token contexts.
- $|\mathcal{V}| - 1$: Free parameters per context (due to the simplex constraint).

Parameter Estimation

Parameters are estimated using **maximum likelihood** in *closed form*:

$$\hat{\theta}_{\mathbf{c}, w_k} = \frac{\text{count}(\mathbf{c}, w_k)}{\text{count}(\mathbf{c})},$$

where:

- $\text{count}(\mathbf{c}, w_k)$: Number of times (\mathbf{c}, w_k) appears in the training corpus \mathcal{D} .
- $\text{count}(\mathbf{c}) = \sum_{w_k \in \mathcal{V}} \text{count}(\mathbf{c}, w_k)$: Total occurrences of \mathbf{c} .

Notes

Deriving the MLE for n -Gram Language Models I

Log-Likelihood for n -Gram Models

Given a training corpus $\mathcal{D} = \{\mathbf{w}^{(i)}\}_{i=1}^M$, where each sequence $\mathbf{w}^{(i)} = (w_1^{(i)}, \dots, w_{n^{(i)}}^{(i)})$, the log-likelihood of the parameters θ is:

$$\mathcal{L}(\theta; \mathcal{D}) = \sum_{i=1}^M \log p_{\theta}(\mathbf{w}^{(i)}) = \sum_{i=1}^M \sum_{k=1}^{n^{(i)}} \log p_{\theta}(w_k^{(i)} \mid \mathbf{c}_k^{(i)}),$$

where $\mathbf{c}_k^{(i)} = (w_{k-n+1}^{(i)}, \dots, w_{k-1}^{(i)})$ is the $(n-1)$ -token context.

Maximizing the Log-Likelihood

Substitute $p_{\theta}(w_k \mid \mathbf{c}) = \theta_{\mathbf{c}, w_k}$:

$$\mathcal{L}(\theta; \mathcal{D}) = \sum_{\mathbf{c} \in \mathcal{V}^{n-1}} \sum_{w \in \mathcal{V}} \text{count}(\mathbf{c}, w) \log \theta_{\mathbf{c}, w}.$$

Subject to the constraint that for each context \mathbf{c} ,

$$\sum_{w \in \mathcal{V}} \theta_{\mathbf{c}, w} = 1.$$

Navigation icons: back, forward, search, etc.

Notes

Deriving the MLE for n -Gram Language Models II

Solving with Lagrange Multipliers

Define the Lagrangian:

$$\mathcal{L}'(\theta, \lambda) = \sum_{\mathbf{c} \in \mathcal{V}^{n-1}} \sum_{w \in \mathcal{V}} \text{count}(\mathbf{c}, w) \log \theta_{\mathbf{c}, w} + \sum_{\mathbf{c} \in \mathcal{V}^{n-1}} \lambda_{\mathbf{c}} \left(1 - \sum_{w \in \mathcal{V}} \theta_{\mathbf{c}, w} \right).$$

Taking the derivative w.r.t. $\theta_{\mathbf{c}, w}$ and setting to zero:

$$\frac{\partial \mathcal{L}'}{\partial \theta_{\mathbf{c}, w}} = \frac{\text{count}(\mathbf{c}, w)}{\theta_{\mathbf{c}, w}} - \lambda_{\mathbf{c}} = 0 \quad \implies \quad \theta_{\mathbf{c}, w} = \frac{\text{count}(\mathbf{c}, w)}{\lambda_{\mathbf{c}}}.$$

Enforce the normalization constraint:

$$\sum_{w \in \mathcal{V}} \theta_{\mathbf{c}, w} = 1 \quad \implies \quad \lambda_{\mathbf{c}} = \text{count}(\mathbf{c}).$$

Substitute $\lambda_{\mathbf{c}}$ back to get the MLE estimate for $\theta_{\mathbf{c}, w}$:

$$\hat{\theta}_{\mathbf{c}, w} = \frac{\text{count}(\mathbf{c}, w)}{\text{count}(\mathbf{c})}.$$

Navigation icons: back, forward, search, etc.

Notes

Order of n -Gram Models

Definition of **Order**

- An **n -gram model** uses the last $n - 1$ tokens to predict the next token:

$$p(w_k \mid w_1, \dots, w_{k-1}) \approx p(w_k \mid w_{k-n+1}, \dots, w_{k-1}).$$

- The integer n is called the **order** of the model. For example:
 - $n = 1$: **Unigram** model (context-free).
 - $n = 2$: **Bigram** model (1-token context).
 - $n = 3$: **Trigram** model (2-token context).

Impact of Model Order

- **Higher order** (n large):
 - Captures longer-range dependencies in text.
 - Increases the number of parameters dramatically, leading to potential *data sparsity*.
- **Lower order** (n small):
 - Fewer parameters, simpler to estimate from limited data.
 - May miss important context (lacks expressive power).

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- **Addressing Data Sparsity in n-Gram Models**
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Data Sparsity in n-Gram Models

Where Does Sparsity Come From?

- The vocabulary \mathcal{V} can be large (tens or hundreds of thousands of tokens).
- As N grows, so does the number of possible $(n - 1)$ -token contexts: $|\mathcal{V}|^{n-1}$.
- Many valid $(n - 1)$ -gram contexts may appear *zero or very few times* in the training data \mathcal{D} .

Consequences of Data Sparsity

- **Zero Counts:** Some $(n - 1)$ -gram contexts are never observed, leading to

$$\hat{\theta}_{\mathbf{c},w} = \frac{\text{count}(\mathbf{c}, w)}{\text{count}(\mathbf{c})} = 0 \quad (\text{no observed tokens}).$$

- **Poor Generalization:** A context not seen in training has probability 0, causing the entire probability of any sentence containing such a context to also become 0.
- **Need for Smoothing:** Techniques like Laplace, Kneser–Ney, or Good–Turing adjust counts to avoid assigning zero probability.
- **Memory and Computation:** Large $|\mathcal{V}|^{n-1}$ means storing and computing vast tables for $\theta_{\mathbf{c},w}$.

Notes

Laplace Smoothing (Add-One Smoothing)

Motivation

- Pure MLE often assigns *zero* probability to unseen $(n - 1)$ -gram contexts.
- **Smoothing** redistributes probability mass to ensure every event has a nonzero probability.

Formula for Add-One Smoothing

- Original MLE estimate:

$$\hat{\theta}_{\mathbf{c},w} = \frac{\text{count}(\mathbf{c}, w)}{\text{count}(\mathbf{c})}.$$

- Add-One smoothing (Laplace):

$$\hat{\theta}_{\mathbf{c},w}^{\text{Laplace}} = \frac{\text{count}(\mathbf{c}, w) + 1}{\text{count}(\mathbf{c}) + |\mathcal{V}|}.$$

- Each (\mathbf{c}, w) is treated as if it appeared at least once.
- Denominator adds $|\mathcal{V}|$ to account for adding 1 for each possible token w .
- Eliminates zero probabilities.

Notes

Problems with Laplace Smoothing

Uniform Distribution for Unobserved Contexts and Over-Smoothing Rare Contexts

- For unobserved $(n - 1)$ -gram contexts \mathbf{c} ($\text{count}(\mathbf{c}) = 0$), Laplace smoothing assigns:

$$\hat{\theta}_{\mathbf{c},w}^{\text{Laplace}} = \frac{1}{|\mathcal{V}|},$$

resulting in a uniform distribution across the vocabulary.

- This fails to capture any linguistic structure or dependencies in the data.
- For rare contexts (e.g., $\text{count}(\mathbf{c}) = 2$), smoothing redistributes too much probability to unseen tokens.

High Sensitivity to Vocabulary Size and Model's Order

- The denominator ($\text{count}(\mathbf{c}) + |\mathcal{V}|$) grows with $|\mathcal{V}|$, making the smoothed probabilities heavily dependent on the vocabulary size.
- As n increases, the number of possible $(n - 1)$ -gram contexts grows exponentially: $|\mathcal{V}|^{n-1}$.
- Even large corpora cannot cover this space, leading to unrealistic distributions for unseen or rare contexts.

Notes

Handling Unknown Words

Out-of-Vocabulary (OOV) Words

- Even large training corpora cannot cover every word form or proper noun.
- Any word ω *not observed* in training is **out-of-vocabulary** (OOV).
- Issue:** If OOV word appears in testing (or real-world usage), the n -gram model has zero probability for any sequence containing ω .

<UNK> Token

- A common approach is to **preemptively** replace low-frequency words in the training data with a special symbol <UNK>.
- This maps all rare or unobserved words to a single <UNK> token, effectively reducing vocabulary size.
- <UNK> is then treated like any other token in the n -gram model, allowing the model to handle previously unseen words during inference.
- Threshold Method:**
 - If $\text{count}(\omega) < \tau$, replace ω with <UNK> in training.
 - Choose τ (e.g., 1, 2, 5) based on data scale and performance.
- Vocabulary Pruning:**
 - Keep only the top $\alpha\%$ most frequent words and map the rest to <UNK>.

Notes

Linear Interpolation of n -gram Models

Motivation

- Pure **MLE** or simple smoothing (e.g., Laplace) can still suffer from zero probabilities for higher-order n -grams with low counts.
- **Interpolation** combines multiple context lengths (orders) rather than “backing off” only when higher-order counts are insufficient.
- Offers a continuous blend of *all* available contexts, reducing the abruptness of pure backoff.

General Interpolation Formula (Trigram Example)

Suppose you want to interpolate among unigram ($N = 1$), bigram ($N = 2$), and trigram ($N = 3$) models:

$$p_{\text{interp}}(w_k \mid w_{k-2}, w_{k-1}) = \lambda_3 p_{\text{MLE}}(w_k \mid w_{k-2}, w_{k-1}) + \lambda_2 p_{\text{MLE}}(w_k \mid w_{k-1}) + \lambda_1 p_{\text{MLE}}(w_k),$$

where:

- $\sum_{i=1}^3 \lambda_i = 1$.
- $p_{\text{MLE}}(\cdot)$ are the standard MLE estimates for each context size, can use smoothing.
- $\{\lambda_i\}$ can be tuned on a held-out *validation set* (e.g., maximize likelihood or minimize perplexity).
- Often, λ_i depend on context counts so that higher-order models get more weight when data is sufficient.

Notes

Special Tokens: $\langle s \rangle$ and $\langle /s \rangle$

Purpose of Special Tokens

- $\langle s \rangle$: Marks the **start of a sentence**.
- $\langle /s \rangle$: Marks the **end of a sentence**.

Motivation

- **Defining Sentence Boundaries:** $\langle s \rangle$ and $\langle /s \rangle$ provide explicit delimiters for sequences.
- **Context Padding for n -Gram Models:**
 - For n -gram models, prepend $(n - 1)$ ' $\langle s \rangle$ ' tokens to the beginning of a sentence.
 - Example (Bigram Model): $p(w_1, w_2, w_3) \approx p(w_1 \mid \langle s \rangle)p(w_2 \mid w_1)p(w_3 \mid w_2)p(\langle /s \rangle \mid w_3)$.
- **Termination in Generation:** Models recognize $\langle /s \rangle$ as the endpoint for generated sequences, preventing infinite loops.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- **Evaluation Metrics for Language Models**
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Entropy of a Discrete Distribution I

Intuition

- **Entropy** measures the average uncertainty or surprise of a random variable.
- In language modeling, it reflects how predictable or unpredictable the tokens are under a distribution.

Formal Definition

Let X be a discrete random variable with a probability mass function $p(x)$ over some set \mathcal{V} . The **entropy** $H(X)$ is defined as:

$$H(X) = - \sum_{x \in \mathcal{V}} p(x) \log p(x).$$

- The base of the logarithm determines the units:
 - Base 2: Entropy is measured in **bits**.
 - Base e : Entropy is measured in **nats**.
- **Bits**: The number of binary (yes/no) questions needed, on average, to identify an outcome of X .
- High entropy \Rightarrow high unpredictability; low entropy \Rightarrow more predictability.

Notes

Entropy of a Discrete Distribution II

Example

- Suppose $\mathcal{V} = \{\text{cat}, \text{dog}, \text{mouse}\}$ with $p(\text{cat}) = 0.5$, $p(\text{dog}) = 0.3$, $p(\text{mouse}) = 0.2$. Then

$$H(X) = -[0.5 \log 0.5 + 0.3 \log 0.3 + 0.2 \log 0.2].$$

- If base 2, $H(X) \approx 1.485$ bits.

Interpreting Binary Questions

- Suppose X represents a random word from $\{\text{cat}, \text{dog}, \text{mouse}\}$:
- To identify the outcome of X using yes/no questions:
 - Q1: Is it cat? ($p(\text{cat}) = 0.5$) - If yes, stop (probability 0.5). - If no, proceed (probability 0.5).
 - Q2: Is it dog? ($p(\text{dog}) = 0.3$) - If yes, stop (probability 0.3). - If no, stop at mouse (probability 0.2).
- Expected Number of Questions:**

$$H(X) \approx 1.485 \text{ bits} \quad (\text{on average, slightly fewer than 2 binary questions}).$$

Notes

Cross-Entropy and the Derivation of KL-Divergence I

Cross-Entropy Definition

Let $p(x)$ be the *true* distribution and $q(x)$ be a *model* distribution over the same set \mathcal{V} . The **cross-entropy** $H(p, q)$ is:

$$H(p, q) = - \sum_{x \in \mathcal{V}} p(x) \log q(x).$$

- Measures how well the model q “fits” the true data p .
- If $q = p$, then $H(p, q) = H(p)$, the entropy of p .

Notes

Cross-Entropy and the Derivation of KL-Divergence II

Derivation of KL-Divergence

Starting from the cross-entropy:

$$H(p, q) = - \sum_{x \in \mathcal{V}} p(x) \log q(x),$$

we can rewrite:

$$- \sum_{x \in \mathcal{V}} p(x) \log q(x) = - \sum_{x \in \mathcal{V}} p(x) \log p(x) - \sum_{x \in \mathcal{V}} p(x) \log \left(\frac{q(x)}{p(x)} \right).$$

Observe that

$$\log q(x) = \log p(x) + \log \left(\frac{q(x)}{p(x)} \right).$$

Therefore,

$$H(p, q) = \underbrace{- \sum_{x \in \mathcal{V}} p(x) \log p(x)}_{= H(p)} + \underbrace{\sum_{x \in \mathcal{V}} p(x) \log \left(\frac{p(x)}{q(x)} \right)}_{= D_{\text{KL}}(p \parallel q)}.$$

Notes

Cross-Entropy and the Derivation of KL-Divergence III

KL-Divergence

We define the **Kullback–Leibler (KL) divergence** as

$$D_{\text{KL}}(p \parallel q) = \sum_{x \in \mathcal{V}} p(x) \log \frac{p(x)}{q(x)} \geq 0.$$

Hence, we obtain the well-known relationship:

$$H(p, q) = H(p) + D_{\text{KL}}(p \parallel q).$$

- $D_{\text{KL}}(p \parallel q) = 0$ if and only if $p = q$.
- Minimizing cross-entropy \Leftrightarrow Minimizing KL-divergence.

Notes

Using Cross-Entropy for LM Evaluation I

Empirical vs. Model Distribution

- We have a **test set** of sequences:

$$\mathcal{D}_{\text{test}} = \left\{ (w_1^{(i)}, w_2^{(i)}, \dots, w_{n(i)}^{(i)}) \right\}_{i=1}^M.$$

- Let $N = \sum_{i=1}^M n^{(i)}$ be the total number of tokens across all sequences.
- The *empirical distribution* \hat{p} places probability $\frac{1}{N}$ on each token $w_k^{(i)}$ in $\mathcal{D}_{\text{test}}$.
- Our **language model** is a distribution $p_{\theta}(w_k \mid w_{1:k-1})$ over the next token given its context.

Notes

Using Cross-Entropy for LM Evaluation II

Per-Token Cross-Entropy and Negative Log-Likelihood

- **Cross-Entropy:**

$$H(\hat{p}, p_{\theta}) = - \sum_{i=1}^N \frac{1}{N} \log p_{\theta}(w^{(i)}),$$

where each $w^{(i)}$ is treated as an i.i.d. sample from \hat{p} .

- Equivalently,

$$H(\hat{p}, p_{\theta}) = - \frac{1}{N} \sum_{i=1}^M \sum_{k=1}^{n(i)} \log(p_{\theta}(w_k^{(i)} \mid w_{1:k-1}^{(i)})).$$

- This **per-token cross-entropy** is exactly the **average negative log-likelihood** of the test set under p_{θ} .
- ↓ **Lower cross-entropy** \Rightarrow the model assigns *higher probability* to the observed tokens.

Notes

Perplexity as a Measure of LM Quality

Definition of Perplexity

- Perplexity is an exponentiation of the cross-entropy, providing a more intuitive scale.
- If using natural logs,

$$PP(p_{\theta}) = \exp\left(H(\hat{p}, p_{\theta})\right).$$

- If using base-2 logs,

$$PP(p_{\theta}) = 2^{H(\hat{p}, p_{\theta})}.$$

Why Perplexity is Intuitive

- **Average Branching Factor:**
 - Imagine each token prediction as choosing among equally likely options.
 - Perplexity says “on average, how many distinct choices does the model effectively consider?”
 - A perplexity of 1 means the model is *never* uncertain; larger values indicate greater uncertainty.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- **Toy Example**
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Example: Toy Corpus and Tri-Gram Model Setup I

Corpus & Vocabulary

Toy Corpus \mathcal{D} consists of three sentences, each prepended with two $\langle s \rangle$:

```
<s> <s> i love cats </s>
<s> <s> i love dogs </s>
<s> <s> cats chase mice </s>
```

Vocabulary \mathcal{V} : $\{\langle s \rangle, i, love, cats, dogs, chase, mice, \langle /s \rangle\}$.

Trigram Model Assumption

- For each position k , we model $p_{\theta}(w_k \mid w_{k-2}, w_{k-1})$.
- Example: In $\langle s \rangle \langle s \rangle i \text{ love cats } \langle /s \rangle$, the third token i is predicted by $p_{\theta}(i \mid \langle s \rangle, \langle s \rangle)$.
- We will collect all (2-token context, next token) counts from \mathcal{D} and apply MLE:

$$\hat{\theta}_{c,w} = \frac{\text{count}(c, w)}{\text{count}(c)}.$$

Notes

Example: Toy Corpus and Tri-Gram Model Setup II

Context-Next Token Counts

Below is a **partial** table of contexts ($c = (w_{k-2}, w_{k-1})$) and how often each next token appears:

Context (w_{k-2}, w_{k-1})	Next Token	Count	Sum Over Next Toks	MLE Probability
$\langle s \rangle, \langle s \rangle$	i	2	3	$\frac{2}{3} \approx 0.67$
$\langle s \rangle, \langle s \rangle$	cats	1		$\frac{1}{3} \approx 0.33$
$\langle s \rangle, i$	love	2	2	$\frac{2}{2} = 1.0$
$i, love$	cats	1	2	$\frac{1}{2} = 0.5$
$i, love$	dogs	1		$\frac{1}{2} = 0.5$
$love, cats$	$\langle /s \rangle$	1	1	1.0
...				

Note: Fill out this table for *all* observed 2-token contexts in the corpus (omitting zero-count contexts not observed, or using smoothing).

Notes

Example: Toy Corpus and Tri-Gram Model Setup III

Probability of a New Sentence

Test Sentence: <s> <s> i love mice </s>

- Using chain rule for trigrams:

$$p_{\theta}(\text{<s> <s> i love mice </s>}) = p_{\theta}(i \mid \text{<s> <s>}) \times p_{\theta}(\text{love} \mid \text{<s> i}) \\ \times p_{\theta}(\text{mice} \mid \text{i love}) \times p_{\theta}(\text{</s>} \mid \text{love mice}).$$

- Since $p_{\theta}(\text{mice} \mid \text{i, love}) = 0$, then the entire product is zero *unless* we apply smoothing.

Cross-Entropy & Perplexity Computation

- Let N be total tokens in <s> <s> i love mice </s> (which is 6).
- Per-token cross-entropy** = $-\frac{1}{6} \sum_{k=1}^5 \log p_{\theta}(w_k \mid w_{k-2}, w_{k-1})$.
- Perplexity** = $\exp(\text{cross-entropy})$.
- Example:** If $p_{\theta}(\text{mice} \mid \text{i, love}) = 0$, perplexity is infinite.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models**
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Generation Strategies for Language Models I

Decoding Algorithm: Greedy vs. Sampling (with Temperature)

Input:

- Trained language model $p_{\theta}(w_k \mid w_1, \dots, w_{k-1})$. Initial context c_{init} (e.g., $\langle s \rangle, \langle s \rangle$ for a trigram model).
- Decoding strategy: choose either greedy or sampling. (Optional) Temperature T for sampling.

Algorithm:

1. Initialize context $c \leftarrow c_{\text{init}}$ and set $sequence \leftarrow []$.
2. Repeat
 - 2.1 Compute $p_{\theta}(w \mid c)$ for all $w \in \mathcal{V}$.
 - 2.2 if strategy is greedy:
$$w^* \leftarrow \arg \max_{w \in \mathcal{V}} p_{\theta}(w \mid c).$$
 - 2.3 else if strategy is sampling:
$$w^* \sim p_{\theta}^{(T)}(w \mid c) = \frac{p_{\theta}(w \mid c)^{1/T}}{\sum_{w' \in \mathcal{V}} p_{\theta}(w' \mid c)^{1/T}}.$$
 - 2.4 Append w^* to $sequence$ and update context c .
3. Until $w^* = \langle /s \rangle$.
4. Return $sequence$ (optionally excluding special tokens like $\langle s \rangle$ and $\langle /s \rangle$).

Navigation icons: back, forward, search, etc.

Notes

Generation Strategies for Language Models II

Greedy vs. Sampling

- **Sampling:**
 - At each step, sample the next token w_k from $p_{\theta}(w_k \mid w_1, \dots, w_{k-1})$.
 - **Pros:** Can produce diverse, creative outputs.
 - **Cons:** May generate nonsensical or low-probability tokens if distribution is broad.
- **Greedy Decoding:**
 - Always pick the token w_k with the highest probability $\arg \max p_{\theta}(w_k \mid w_1, \dots, w_{k-1})$.
 - **Pros:** Fastest method, easy to implement.
 - **Cons:** Often gets stuck in repetitive or sub-optimal sequences (lack of diversity).
- **Temperature Scaling**
 - Effects of T : $T > 1$: Flattens the distribution, increasing randomness. $T < 1$: Sharpens the distribution.
 - **Pros:** Fine-grained control over output randomness.
 - **Cons:** Requires careful tuning of T for desired behavior.

Other Sampling Strategies

- **Top-k:** Restrict sampling to the k most probable tokens at each step.
- **Nucleus (Top- p):** Sample from the smallest set of tokens whose cumulative probability exceeds p .

Navigation icons: back, forward, search, etc.

Notes

Generation Strategies for Language Models III

Beam Search Algorithm

Input:
Trained language model $p_{\theta}(w_k \mid w_1, \dots, w_{k-1})$.
Initial context \mathbf{c}_{init} (e.g., $\langle \text{s} \rangle, \langle \text{s} \rangle$).
Beam size B (number of parallel hypotheses to maintain) and maximum length L .

Algorithm:
Initialize $\text{candidates} \leftarrow \{(\mathbf{c}_{\text{init}}, 0)\}$, where each candidate is a tuple of context and log-probability.
Initialize $\text{final_sequences} \leftarrow []$.

Repeat:
For each candidate $(\mathbf{c}, \text{score})$ in candidates :
 Compute $p_{\theta}(w \mid \mathbf{c})$ for all $w \in \mathcal{V}$.
 Extend \mathbf{c} with each w , forming new candidates:
 $(\mathbf{c} + w, \text{score} + \log p_{\theta}(w \mid \mathbf{c}))$.
 If $w = \langle \text{s} \rangle$:
 Move $(\mathbf{c} + w, \text{score})$ to final_sequences .
Retain the top B candidates by score for the next step.

Until: All B candidates end with $\langle \text{s} \rangle$ or maximum length is reached.

Return: The highest-scoring sequence from final_sequences .

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- **Feed-Forward Neural Language Models**
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Problems with Categorical n -Gram Parametrization

Exponential Growth of Parameters

- A categorical n -gram model requires a unique parameter $\theta_{\mathbf{c},w}$ for each context \mathbf{c} (of length $n - 1$) and next word w . Total number of parameters is exponential in the context length $n - 1$:

$$|\mathcal{V}|^{n-1} \cdot (|\mathcal{V}| - 1).$$

Sparsity and Zero Probabilities

- For most possible n -grams, the count $\text{count}(\mathbf{c}, w) \approx 0$, causing $\theta_{\mathbf{c},w} \approx 0$ for many (\mathbf{c}, w) if no smoothing is used to modify the counts.

Lookup Table Representation $p(w \mid \mathbf{c}) = \theta_{\mathbf{c},w}$

- Input:** Each word in the key n -gram can be seen as a **one-hot vector** $\mathbf{1}_w \in \{0, 1\}^{|\mathcal{V}|}$. The n -gram $(\mathbf{c}_1, \dots, \mathbf{c}_{n-1}, w)$ can be seen as a concatenation of n one-hot word vectors:

$$\mathbf{x} = [\mathbf{1}_{\mathbf{c}_1}; \dots; \mathbf{1}_{\mathbf{c}_{n-1}}; \mathbf{1}_w] \in \mathbb{R}^{|\mathcal{V}| \cdot n}.$$

- There is no intrinsic notion of similarity between different contexts in this representation.

Navigation icons: back, forward, search, etc.

Notes

Feed-Forward Neural Network Parametrization

Input Mapping

- Each word w in the context \mathbf{c} is mapped to an embedding $\mathbf{e}_w \in \mathbb{R}^d$, with $d \ll |\mathcal{V}|$.
- Computed by projecting one-hot vectors through an embedding matrix $\mathbf{E} \in \mathbb{R}^{|\mathcal{V}| \times d}$: $\mathbf{e}_w = \mathbf{E}^\top \mathbf{1}_w$.
- Embeddings of the $n - 1$ context words are concatenated

$$\mathbf{x} = [\mathbf{e}_{\mathbf{c}_1}; \dots; \mathbf{e}_{\mathbf{c}_{n-1}}] \in \mathbb{R}^{d \cdot (n-1)}.$$

Hidden Layer

$$\mathbf{h} = \tanh(\mathbf{W}^{(h)} \mathbf{x} + \mathbf{b}^{(h)}), \quad \mathbf{W}^{(h)} \in \mathbb{R}^{m \times (d \cdot (n-1))}, \quad \mathbf{b}^{(h)} \in \mathbb{R}^m.$$

- Computes a continuous context embedding $\mathbf{h} \in \mathbb{R}^m$. Learns to mix features from the word embeddings.

Output Layer

$$p(w \mid \mathbf{c}) = \mathbf{p}_w \quad \text{where} \quad \mathbf{p} \in \mathbb{R}^{|\mathcal{V}|}, \quad \mathbf{p} = \text{softmax}(\mathbf{W}^{(o)} \mathbf{h} + \mathbf{b}^{(o)}), \quad \mathbf{W}^{(o)} \in \mathbb{R}^{|\mathcal{V}| \times m}, \quad \mathbf{b}^{(o)} \in \mathbb{R}^{|\mathcal{V}|}.$$

Navigation icons: back, forward, search, etc.

Notes

Comparison: Categorical vs. Neural Parametrization I

Parameter Sets

Categorical n -gram:

$$\theta_{\text{cat}} = \left\{ \theta_{c,w} \mid c \in \mathcal{V}^{n-1}, w \in \mathcal{V} \right\}$$

Neural LM:

$$\theta_{\text{nn}} = \left\{ \mathbf{E} \in \mathbb{R}^{|\mathcal{V}| \times d}, \mathbf{W}^{(h)} \in \mathbb{R}^{m \times (d \cdot (n-1))}, \mathbf{b}^{(h)} \in \mathbb{R}^m, \mathbf{W}^{(o)} \in \mathbb{R}^{|\mathcal{V}| \times m}, \mathbf{b}^{(o)} \in \mathbb{R}^{|\mathcal{V}|} \right\}.$$

- Here, $\mathbf{h} \in \mathbb{R}^m$ is the hidden context embedding. Neural LM uses far fewer parameters due to sharing across contexts.

Notes

Comparison: Categorical vs. Neural Parametrization II

Practical Example of Parameter Sizes

Assume a vocabulary size $|\mathcal{V}| = 10,000$, embedding dimension $d = 300$, hidden layer size $m = 500$, and $n = 3$ (trigram).

▪ **Categorical Trigram:**

$$\text{Parameters} \approx |\mathcal{V}|^2 \cdot (|\mathcal{V}| - 1) \approx 10,000^2 \cdot 9,999 \approx 10^{12}.$$

▪ **Neural Trigram:**

$$\begin{aligned} & |\mathcal{V}| \times d + m \times ((n-1) \cdot d) + m + |\mathcal{V}| \times m + |\mathcal{V}| \\ & \approx 10,000 \times 300 + 500 \times (2 \times 300) + 500 + 10,000 \times 500 + 10,000 \\ & \approx 3 \times 10^6 + 300,000 + 500 + 5 \times 10^6 + 10,000 \\ & \approx 8.3 \times 10^6 \text{ parameters.} \end{aligned}$$

Notes

Feature Extraxtion

Feature Mixing in Hidden Layers

$$\underbrace{\left(\dots \underbrace{W_{j,i_1}^{(h)}}_{\text{large weight}} \dots \underbrace{W_{j,i_2}^{(h)}}_{\text{large weight}} \dots \right)}_{\text{row } j \text{ of } \mathbf{W}^{(h)}} \times \underbrace{\begin{pmatrix} \vdots \\ \underbrace{x_{i_1}}_{\text{dim } i_1} \\ \vdots \\ \underbrace{x_{i_2}}_{\text{dim } i_2} \\ \vdots \\ \underbrace{x}_{\mathbf{x}} \end{pmatrix}}_{\mathbf{x}} \rightarrow h_j = \tanh \left(\sum_{i=1}^{d \cdot (n-1)} W_{j,i}^{(h)} x_i + b_j^{(h)} \right).$$

- Each row $\mathbf{W}_{j,\cdot}^{(h)}$ selectively combines specific dimensions of the input \mathbf{x} .
- Larger weights $|W_{j,i}^{(h)}|$ amplify embedding dimensions (e.g., those tied to nouns or adjectives).
- Thus, h_j can learn a particular pattern by focusing on relevant parts of \mathbf{x} .

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- **Training**
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Training a Feed-Forward Neural LM I

Training Corpus and Empirical Distribution

- Let $\mathcal{D} = \{\mathbf{w}^{(i)}\}_{i=1}^N$ be a set of N sentences (or sequences), each $\mathbf{w}^{(i)} = (w_1^{(i)}, \dots, w_{T_i}^{(i)})$.
- The **empirical distribution** $\hat{p}(\mathbf{w})$ places probability $\frac{1}{N}$ on each training sentence $\mathbf{w}^{(i)}$.

Cross-Entropy \Leftrightarrow Maximum Likelihood

- Our model $p_{\theta}(\mathbf{w})$ assigns a probability to any sentence \mathbf{w} . **Cross-entropy** between \hat{p} and p_{θ} :

$$H(\hat{p}, p_{\theta}) = - \sum_{i=1}^N \frac{1}{N} \log p_{\theta}(\mathbf{w}^{(i)}).$$

- Minimizing $H(\hat{p}, p_{\theta}) \Leftrightarrow \arg \max_{\theta} \prod_{i=1}^N p_{\theta}(\mathbf{w}^{(i)})$, i.e. **maximum likelihood estimation (MLE)**.
- This objective is also known as the **negative log-likelihood (NLL)**:

$$\ell(\theta) = - \sum_{i=1}^N \log p_{\theta}(\mathbf{w}^{(i)}).$$

Navigation icons: back, forward, search, etc.

Notes

Training a Feed-Forward Neural LM II

Chain Rule and n -grams

- In a feed-forward LM with context size $n - 1$:

$$p_{\theta}(\mathbf{w}^{(i)}) = \prod_{k=1}^{T_i} p_{\theta}(w_k^{(i)} \mid w_{k-n+1}^{(i)}, \dots, w_{k-1}^{(i)}).$$

- Each term $p_{\theta}(w_k \mid \mathbf{c}_k)$ is computed via:

$$\mathbf{c}_k \mapsto \mathbf{x} \mapsto \mathbf{h} \mapsto \mathbf{p} = \text{softmax}(\mathbf{W}^{(o)} \mathbf{h} + \mathbf{b}^{(o)}), \quad p_{\theta}(w_k \mid \mathbf{c}_k) = p_{w_k}.$$

Loss Over the Entire Corpus

$$\ell(\theta) = - \sum_{i=1}^N \log p_{\theta}(\mathbf{w}^{(i)}) = - \sum_{i=1}^N \sum_{k=1}^{T_i} \log p_{\theta}(w_k^{(i)} \mid \mathbf{c}_k^{(i)}).$$

- Minimizing $\ell(\theta)$ sums the negative log-probabilities over all context-target pairs (\mathbf{c}_k, w_k) .
- **Single Pair Loss:** $\ell(\theta; \mathbf{c}, w) = - \log p_{\theta}(w \mid \mathbf{c})$.

Navigation icons: back, forward, search, etc.

Notes

Training a Feed-Forward Neural LM III

Parameter Update Rule

- Use gradient-based methods (SGD, Adam, etc.) to update parameters:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \ell(\theta).$$

- η is the learning rate. In practice, Adam or RMSProp handle adaptive step sizes and momentum.

Forward, Loss, and Backprop

Forward Pass:

$$\mathbf{x} = [\mathbf{e}_{w_{k-n+1}}; \dots; \mathbf{e}_{w_{k-1}}], \mathbf{h} = \tanh(\mathbf{W}^{(h)} \mathbf{x} + \mathbf{b}^{(h)}), \mathbf{z} = \mathbf{W}^{(o)} \mathbf{h} + \mathbf{b}^{(o)}, \mathbf{p} = \text{softmax}(\mathbf{z}).$$

$$\ell(\theta; \mathbf{c}, \mathbf{w}) = -\log \mathbf{p}_w.$$

Backward Pass:

- Derive $\nabla_{\mathbf{z}} \ell$ from the softmax derivative, propagate to \mathbf{h} and \mathbf{x} via chain rule (through tanh, matrix multiplies).
- Accumulate gradients for $\mathbf{W}^{(h)}, \mathbf{b}^{(h)}, \mathbf{W}^{(o)}, \mathbf{b}^{(o)}$, and \mathbf{E} (embedding matrix).

Navigation icons: back, forward, search, etc.

Notes

Training a Feed-Forward Neural LM IV

Mini-Batch Training and Vectorization

- Instead of processing one (\mathbf{c}, w) at a time, we group examples into mini-batches (e.g., size 32).
- **Vectorization:**
 - Stack the \mathbf{x} vectors of multiple examples into a matrix \mathbf{X} .
 - Compute $\mathbf{W}^{(h)} \mathbf{X}$ (and subsequent layers) in parallel for the whole batch.
- Average gradients over the mini-batch, then update parameters, resulting in more stable training and GPU efficiency.

Navigation icons: back, forward, search, etc.

Notes

Detailed Gradient Derivations I

Setup:

$$\ell(\boldsymbol{\theta}; \mathbf{c}, w) = -\log p_{\boldsymbol{\theta}}(w \mid \mathbf{c}), \quad p_{\boldsymbol{\theta}}(w \mid \mathbf{c}) = \text{softmax}(\mathbf{z})_w, \quad \mathbf{z} = \mathbf{W}^{(o)} \mathbf{h} + \mathbf{b}^{(o)},$$
$$\mathbf{h} = \tanh(\mathbf{W}^{(h)} \mathbf{x} + \mathbf{b}^{(h)}), \quad \mathbf{x} = [\mathbf{e}_{w_{k-n+1}}, \dots, \mathbf{e}_{w_{k-1}}].$$

where $\mathbf{p} = \text{softmax}(\mathbf{z})$ and $\mathbf{y} \in \{0, 1\}^{|\mathcal{V}|}$ is the one-hot vector for the correct word w . Then:

1. Gradient w.r.t. output logits \mathbf{z} :

$$\frac{\partial \ell}{\partial z_j} = \frac{\partial}{\partial z_j} [-\log(\mathbf{p}_w)] = p_j - y_j \quad (\text{for } j = 1, \dots, |\mathcal{V}|).$$

2. Output layer parameters:

$$\nabla_{\mathbf{W}^{(o)}} \ell = (\mathbf{p} - \mathbf{y}) \mathbf{h}^\top, \quad \nabla_{\mathbf{b}^{(o)}} \ell = \mathbf{p} - \mathbf{y}.$$

3. Hidden layer gradient:

$$\nabla_{\mathbf{h}} \ell = (\mathbf{W}^{(o)})^\top (\mathbf{p} - \mathbf{y}).$$

Then apply chain rule for tanh:

$$\nabla_{\mathbf{z}^{(h)}} \ell = (1 - \tanh^2(\mathbf{z}^{(h)})) \odot \nabla_{\mathbf{h}} \ell,$$

where $\mathbf{z}^{(h)} = \mathbf{W}^{(h)} \mathbf{x} + \mathbf{b}^{(h)}$.

Notes

Detailed Gradient Derivations II

4. Hidden layer parameters:

$$\nabla_{\mathbf{W}^{(h)}} \ell = \nabla_{\mathbf{z}^{(h)}} \ell \mathbf{x}^\top, \quad \nabla_{\mathbf{b}^{(h)}} \ell = \nabla_{\mathbf{z}^{(h)}} \ell.$$

5. Embedding matrix \mathbf{E} : Backprop through \mathbf{x} (the concatenation of each context word's embedding). Each relevant row in \mathbf{E} is updated according to $\frac{\partial \ell}{\partial \mathbf{e}_{w_i}}$.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Motivation: Moving Beyond Fixed Context Size

Limitations of Feed-Forward LM

- **Fixed Window:** A feed-forward LM uses a context of size $n - 1$. Any dependency beyond $n - 1$ tokens is *not* captured.
- **Long-Distance Dependencies in Language:**
 - Example:
*The **car** that I drove yesterday **broke down** this morning.*
The mention of “car” is quite far from the point where we describe what happened to it.

Recurrent Neural Networks (RNNs)

- Designed to capture *variable-length* contexts and long-distance dependencies by maintaining a **hidden state** that updates at each time step.
- The RNN hidden state plays the role of **memory**, combining information from all previous tokens.

Notes

Elman RNN: Detailed Equations I

Notation and Setup

- Let $\mathbf{w} = (w_1, w_2, \dots, w_T)$ be a tokenized sequence.
- At each time step t , the RNN processes the *embedding* $\mathbf{x}_t \in \mathbb{R}^d$ of the current token w_t .
- Maintains a hidden state $\mathbf{h}_t \in \mathbb{R}^m$ capturing *all previously seen* tokens, thus overcoming the fixed-window limitation.

Forward Pass of an Elman RNN

$$\mathbf{h}_t = \tanh\left(\mathbf{W}_{xh} \mathbf{x}_t + \mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{b}_h\right), \quad \mathbf{h}_0 = \mathbf{0} \text{ (or learned)}.$$

- $\mathbf{W}_{xh} \in \mathbb{R}^{m \times d}$: transforms current input \mathbf{x}_t (as in feed-forward LMs).
- $\mathbf{W}_{hh} \in \mathbb{R}^{m \times m}$: **new** recurrent connection, combining the previous state \mathbf{h}_{t-1} .
- $\mathbf{b}_h \in \mathbb{R}^m$: bias term.
- tanh: typical nonlinear activation; other choices (ReLU, etc.) are possible.

Notes

Elman RNN: Detailed Equations II

Key Difference vs. Feed-Forward LM

- Unlike a feed-forward LM (which sees only a fixed window of size $n - 1$), the RNN **recurrently** incorporates \mathbf{h}_{t-1} through \mathbf{W}_{hh} .
- This enables the network to (in principle) use an unbounded context.

Output and Next-Word Distribution

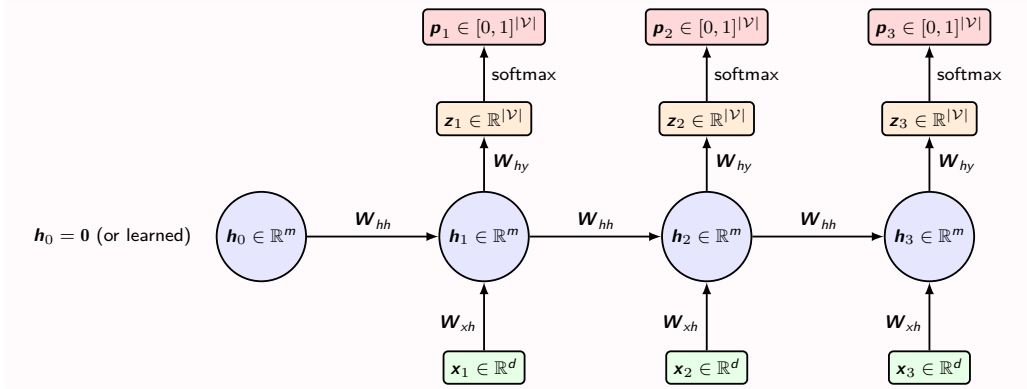
$$\mathbf{z}_t = \mathbf{W}_{hy} \mathbf{h}_t + \mathbf{b}_y, \quad \mathbf{p}_t = \text{softmax}(\mathbf{z}_t), \quad p_{\theta}(w_{t+1} \mid w_{1 \dots t}) = \mathbf{p}_{t, w_{t+1}}.$$

- $\mathbf{z}_t \in \mathbb{R}^{|\mathcal{V}|}$: output logits for next token at time t .
- $\mathbf{p}_t \in \mathbb{R}^{|\mathcal{V}|}$: next-token probability distribution via softmax.
- $\mathbf{W}_{hy} \in \mathbb{R}^{|\mathcal{V}| \times m}$, $\mathbf{b}_y \in \mathbb{R}^{|\mathcal{V}|}$.

Notes

Unrolling RNN in time I

A Diagram



Notes

Unrolling RNN in time II

In Equations

- For inputs x_1, x_2, x_3, x_4 , each $x_t \in \mathbb{R}^d$.
- The hidden state at the final time step ($h_4 \in \mathbb{R}^m$) unfolds as:

$$h_4 = \tanh(W_{hh} \tanh(W_{hh} \tanh(W_{hh} \tanh(W_{hh} h_0 + W_{xh} x_1 + b_h) + W_{xh} x_2 + b_h) + W_{xh} x_3 + b_h) + W_{xh} x_4 + b_h)$$

- Logits at time step 4 ($z_4 \in \mathbb{R}^{|V|}$): $z_4 = W_{hy} h_4 + b_y$.
- Notice that h_4 depends on h_0 and all prior inputs x_1, \dots, x_4 , each influencing the hidden state through multiple nested \tanh transformations.

Notes

Comparison with Feed-Forward LM Architecture

Drawbacks of Feed-Forward LM

- **Fixed window:** $(w_{k-n+1}, \dots, w_{k-1}) \rightarrow \text{concat embeddings} \rightarrow \text{hidden layer} \rightarrow \text{softmax}(\dots)$.
- **Limitations:**
 - Cannot look beyond $(n - 1)$ tokens of context.
 - Parameter explosion if n is large.
 - No built-in mechanism to capture long-distance or variable-length dependencies.

RNN LM Advantages

- **Implicitly unbounded context:** \mathbf{h}_t in principle encodes all previous tokens (w_1, \dots, w_{t-1}) .
- **Shared parameters over time steps:** leads to statistical strength and fewer parameters for large contexts than a large-window feed-forward LM.
- **Recurrent updating:** \mathbf{h}_t evolves recursively, capturing sequential correlations in language.

Notes

Training an RNN LM and Its Challenges I

Objective and Unrolling in Time

- Similar to feed-forward LMs, we define a training set $\mathcal{D} = \{\mathbf{w}^{(i)}\}_{i=1}^N$ of sequences $\mathbf{w}^{(i)} = (w_1^{(i)}, \dots, w_{T_i}^{(i)})$.
- Our RNN LM factorizes $p_{\theta}(\mathbf{w})$ via:

$$p_{\theta}(w_1, \dots, w_T) = \prod_{t=1}^T p_{\theta}(w_t \mid w_1, \dots, w_{t-1}).$$

- **Unrolled Computation:**
 - A hidden state $\mathbf{h}_t \in \mathbb{R}^m$ is computed at each time t : $\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)$ (e.g. Elman update with \tanh).
 - Output logits $\mathbf{z}_t = \mathbf{W}_{hy} \mathbf{h}_t + \mathbf{b}_y$, probabilities $\mathbf{p}_t = \text{softmax}(\mathbf{z}_t)$.
- **Loss over entire sequence:**

$$\ell(\theta; \mathbf{w}) = - \sum_{t=1}^T \log p_{\theta}(w_t \mid w_{1:t-1}).$$

Notes

Backprop Through Time (BPTT)

- We sum (or average) over all time steps and all sequences:

$$\ell(\theta) = \sum_{i=1}^N \sum_{t=1}^{T_i} -\log p_{\theta}(w_t^{(i)} \mid w_{1:t-1}^{(i)}).$$

- **Gradient Computation:**
 - We *unroll* the RNN across time steps $1 \dots T$.
 - Apply backprop to each unrolled connection, known as **BPTT**.
 - Accumulate gradients $\nabla w_{xh}, \nabla w_{hh}, \nabla w_{hy}, \dots$

- **Parameter Updates:**

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \ell(\theta),$$

typically in mini-batches for efficiency.

Notes

Challenges: Vanishing/Exploding Gradients

- **Vanishing Gradients:**
 - When $\|W_{hh}\| < 1$, backprop terms can decay exponentially over many steps.
 - The model struggles to learn long-term dependencies.
- **Exploding Gradients:**
 - When $\|W_{hh}\| > 1$, gradients can grow exponentially, causing instability.
 - Common solutions: gradient clipping, careful initialization.
- Both issues arise because gradients repeatedly multiply through W_{hh} across time.
- **Recurrent Architectures (LSTM/GRU)** partially address these challenges with gating.

Notes

Detailed Gradients for one sequence: $\nabla_{\mathbf{h}_t} \ell(\boldsymbol{\theta}; \mathbf{w})$

Goal: Gradient w.r.t. Hidden State

- Let us call $L = \ell(\boldsymbol{\theta}; \mathbf{w}) = -\sum_{t=1}^T \log p_{\boldsymbol{\theta}}(w_t \mid w_{1:t-1})$.
- We want $\frac{\partial L}{\partial \mathbf{h}_t}$, the gradient of the total sequence loss L wrt. the hidden state \mathbf{h}_t :

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L_t}{\partial \mathbf{h}_t} + \frac{\partial L_{t+1}}{\partial \mathbf{h}_t} + \dots + \frac{\partial L_T}{\partial \mathbf{h}_t}.$$

- Summing direct and indirect contributions:

$$\frac{\partial L}{\partial \mathbf{h}_t} = \underbrace{\frac{\partial L_t}{\partial \mathbf{h}_t}}_{\text{direct from step } t} + \sum_{k=t+1}^T \underbrace{\frac{\partial L_k}{\partial \mathbf{h}_t}}_{\text{indirect from future steps } k>t} = \frac{\partial L_t}{\partial \mathbf{h}_t} + \sum_{k=t+1}^T \frac{\partial L_k}{\partial \mathbf{h}_{t+1}} \frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t}.$$

- Often simplified as a **recursive formula**:

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L_t}{\partial \mathbf{h}_t} + \frac{\partial L}{\partial \mathbf{h}_{t+1}} \frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t}.$$

Notes

Computing the Recursive Gradient I

Hidden State Update

- Recall the simple Elman RNN:

$$\mathbf{h}_{t+1} = \tanh(\mathbf{a}_{t+1}), \quad \mathbf{a}_{t+1} = \mathbf{W}_{hh} \mathbf{h}_t + \mathbf{W}_{xh} \mathbf{x}_{t+1} + \mathbf{b}_h.$$

- We compute:

$$\frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t} = \underbrace{\frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{a}_{t+1}}}_{\text{diag}(1 - \tanh^2(\mathbf{a}_{t+1}))} \cdot \underbrace{\frac{\partial \mathbf{a}_{t+1}}{\partial \mathbf{h}_t}}_{\mathbf{W}_{hh}}.$$

Notes

Computing the Recursive Gradient II

Chain Rule in Detail

- Since $\mathbf{h}_{t+1} = \tanh(\mathbf{a}_{t+1})$ and $\mathbf{W}_{xh} \mathbf{x}_{t+1}$ and \mathbf{b}_h are constants wrt. \mathbf{h}_t :

$$\frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{a}_{t+1}} = \text{diag}(1 - \tanh^2(\mathbf{a}_{t+1})), \quad \frac{\partial \mathbf{a}_{t+1}}{\partial \mathbf{h}_t} = \mathbf{W}_{hh}$$

- Combine:

$$\frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t} = \text{diag}(1 - \tanh^2(\mathbf{a}_{t+1})) \mathbf{W}_{hh}.$$

- Then the gradient update:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{h}_t} &= \frac{\partial L_t}{\partial \mathbf{h}_t} + \frac{\partial L}{\partial \mathbf{h}_{t+1}} \frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t}. \\ &= \frac{\partial L_t}{\partial \mathbf{h}_t} + \frac{\partial L}{\partial \mathbf{h}_{t+1}} \text{diag}(1 - \tanh^2(\mathbf{a}_{t+1})) \mathbf{W}_{hh}. \end{aligned}$$

- Adjust for shape (often a transpose factor). Final form:

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L_t}{\partial \mathbf{h}_t} + \mathbf{W}_{hh}^\top \left[\text{diag}(1 - \tanh^2(\mathbf{a}_{t+1})) \frac{\partial L}{\partial \mathbf{h}_{t+1}} \right].$$

Navigation icons: back, forward, search, etc.

Notes

Unrolling the Recursion

Repeated Application

- Applying the recurrence from t to $t+1$, $t+2$, etc. yields:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{h}_t} &= \frac{\partial L_t}{\partial \mathbf{h}_t} + \mathbf{W}_{hh}^\top \Phi'_{t+1} \frac{\partial L}{\partial \mathbf{h}_{t+1}} \\ &= \frac{\partial L_t}{\partial \mathbf{h}_t} + \mathbf{W}_{hh}^\top \Phi'_{t+1} \left(\frac{\partial L_{t+1}}{\partial \mathbf{h}_{t+1}} + \mathbf{W}_{hh}^\top \Phi'_{t+2} \frac{\partial L}{\partial \mathbf{h}_{t+2}} \right) \\ &\vdots \\ &= \sum_{k=t}^T \left(\left(\prod_{j=t+1}^k \mathbf{W}_{hh}^\top \Phi'_j \right) \frac{\partial L_k}{\partial \mathbf{h}_k} \right) \end{aligned}$$

- Φ'_j denotes $\text{diag}(1 - \tanh^2(\mathbf{a}_j))$.
- This product across many steps can **vanish** if $\|\mathbf{W}_{hh}\| < 1$ or **explode** if $\|\mathbf{W}_{hh}\| > 1$.

Navigation icons: back, forward, search, etc.

Notes

Vanishing & Exploding Gradients (Recap)

Vanishing Gradients

- If $\|W_{hh}\|_2 < 1$, repeated multiplication shrinks gradients **exponentially** with distance:

$$\left\| \frac{\partial L}{\partial h_t} \right\| \leq \left(\|W_{hh}\|_2 \gamma \right)^{(k-t)} \left\| \frac{\partial L}{\partial h_k} \right\|.$$

- Hard to learn long-term dependencies.

Exploding Gradients

- If $\|W_{hh}\|_2 > 1$, norms can blow up:

$$\left\| \frac{\partial L}{\partial h_t} \right\| \geq \left(\|W_{hh}\|_2 \gamma \right)^{(k-t)} \left\| \frac{\partial L}{\partial h_k} \right\|.$$

- Causes numerical instability; we often do **gradient clipping**.

Notes

Mitigating Gradient Problems

Common Strategies

- **Gradient Clipping:**
 - Restricts the norm $\|\nabla_{\theta} \ell\|$ to a predefined threshold.
 - Prevents numeric overflow when gradients become large (exploding gradients).
- **Initialization Techniques:**
 - Properly initializing W_{hh} , W_{xh} etc. to maintain stable gradient propagation.
 - Use orthogonal or unitary matrices for W_{hh} , e.g. $W_{hh} W_{hh}^T = I$.
 - Preserves the norm: $\|W_{hh} x\| = \|x\|$.
 - Helps combat vanishing/exploding gradients.
- **Activation Functions:**
 - ReLU or similar (e.g. Leaky ReLU) can reduce gradient decay compared to tanh.
 - For instance, $\text{ReLU}(x) = \max(0, x)$, derivative is 1 for $x > 0$, allowing large gradient flow.
- **Advanced RNN Architectures:**
 - **LSTM** Introduces a *cell state* and gating mechanisms to preserve long-term information.
 - **GRU** A simpler variant of LSTM, also addresses gradient issues through gating.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- **LSTMs and GRUs**
- What's Next?

Notes

LSTM Architecture: Scalar and Vector Forms I

Scalar Equations (Conceptual)

- For each time t , an LSTM maintains c_t (the *cell state*) and h_t (the *hidden state*).
- Example (scalar version):

$c_t = f_t \cdot c_{t-1} + i_t \cdot z_t,$	cell state
$h_t = o_t \psi(c_t),$	hidden output
$z_t = \varphi(\tilde{z}_t),$	$\tilde{z}_t = w_z^\top x_t + r_z h_{t-1} + b_z,$
$i_t = \sigma(\tilde{i}_t),$	$\tilde{i}_t = w_i^\top x_t + r_i h_{t-1} + b_i,$
$f_t = \sigma(\tilde{f}_t),$	$\tilde{f}_t = w_f^\top x_t + r_f h_{t-1} + b_f,$
$o_t = \sigma(\tilde{o}_t),$	$\tilde{o}_t = w_o^\top x_t + r_o h_{t-1} + b_o.$

- σ is the logistic sigmoid, φ could be tanh. This form highlights the gating logic: i_t (input gate), f_t (forget gate), and o_t (output gate).

Notes

LSTM Architecture: Scalar and Vector Forms II

Vector Form (Practical Implementation)

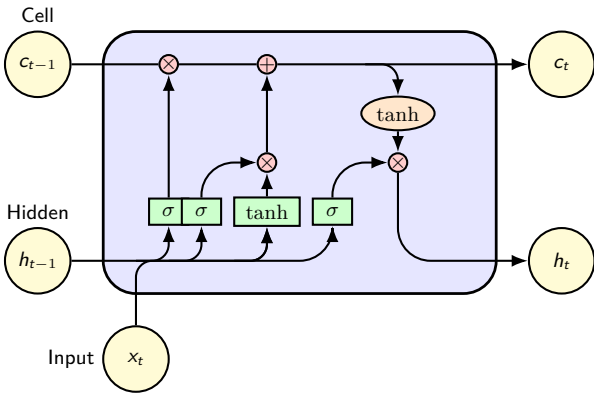
- In practice, we combine scalar gates into vector/matrix operations. For each time t :

$$\begin{aligned} \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{z}_t, & \mathbf{h}_t &= \mathbf{o}_t \odot \psi(\mathbf{c}_t), \\ \mathbf{z}_t &= \varphi(\mathbf{W}_z \mathbf{x}_t + \mathbf{R}_z \mathbf{h}_{t-1} + \mathbf{b}_z), \\ \mathbf{i}_t &= \sigma(\mathbf{W}_i \mathbf{x}_t + \mathbf{R}_i \mathbf{h}_{t-1} + \mathbf{b}_i), \\ \mathbf{f}_t &= \sigma(\mathbf{W}_f \mathbf{x}_t + \mathbf{R}_f \mathbf{h}_{t-1} + \mathbf{b}_f), \\ \mathbf{o}_t &= \sigma(\mathbf{W}_o \mathbf{x}_t + \mathbf{R}_o \mathbf{h}_{t-1} + \mathbf{b}_o). \end{aligned}$$

- $\mathbf{x}_t \in \mathbb{R}^d, \mathbf{h}_t, \mathbf{c}_t \in \mathbb{R}^m$.
- $\mathbf{W}_* \in \mathbb{R}^{m \times d}, \mathbf{R}_* \in \mathbb{R}^{m \times m}, \mathbf{b}_* \in \mathbb{R}^m$.
- Each gate $\mathbf{i}_t, \mathbf{f}_t, \mathbf{o}_t \in \mathbb{R}^m$ controls how info flows in/out of the cell state \mathbf{c}_t .

Notes

Illustration of an LSTM Cell Structure



Notes

How the LSTM's Constant Error Carousel (CEC) Addresses Vanishing Gradients I

Constant Error Carousel in LSTM

- The key update rule in LSTMs for the cell state $\mathbf{c}_t \in \mathbb{R}^m$:

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{z}_t,$$

where \odot is element-wise multiplication.

- **Additive** updates (rather than purely multiplicative) avoid exponential shrinking of gradients.
- Each component $\mathbf{f}_t, \mathbf{i}_t, \mathbf{z}_t$ is computed via gates (e.g. σ or \tanh).

Notes

How the LSTM's Constant Error Carousel (CEC) Addresses Vanishing Gradients II

Gradient Flow Through CEC

- **Recursive Gradient Equation:**

$$\frac{\partial L}{\partial \mathbf{c}_t} = \frac{\partial L_t}{\partial \mathbf{c}_t} + \left(\frac{\partial L}{\partial \mathbf{c}_{t+1}} \odot \mathbf{f}_{t+1} \right).$$

- **Unrolling the Recursion:**

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{c}_t} &= \frac{\partial L_t}{\partial \mathbf{c}_t} + \left[\frac{\partial L_{t+1}}{\partial \mathbf{c}_{t+1}} + \left(\frac{\partial L}{\partial \mathbf{c}_{t+2}} \odot \mathbf{f}_{t+2} \right) \right] \odot \mathbf{f}_{t+1} \\ &= \frac{\partial L_t}{\partial \mathbf{c}_t} + \left(\frac{\partial L_{t+1}}{\partial \mathbf{c}_{t+1}} \odot \mathbf{f}_{t+1} \right) + \left(\frac{\partial L}{\partial \mathbf{c}_{t+2}} \odot \mathbf{f}_{t+2} \odot \mathbf{f}_{t+1} \right) + \dots \\ &= \sum_{k=t}^T \left(\frac{\partial L_k}{\partial \mathbf{c}_k} \odot \prod_{j=t+1}^k \mathbf{f}_j \right). \end{aligned}$$

- Each term is modulated by the product of forget gates $\mathbf{f}_j \in [0, 1]^m$, which can preserve gradient flow if $\mathbf{f}_j \approx 1$. This prevents the exponential decay of gradients, thus solving the vanishing gradient problem.

Notes

Gated Recurrent Units (GRUs)

Motivation

- **Simplify the LSTM architecture:** Reduce the number of gates and parameters while still addressing vanishing gradients.
- **Combine Forget and Input gates** into a single *update* gate to decide how much past information to keep or overwrite.
- Often yields comparable performance to LSTM with a simpler structure and sometimes trains faster.

Key Differences from LSTM

- **No separate cell state c_t .** GRU keeps a single hidden state vector h_t .
- **Two main gates:**
 - z_t (*update gate*): controls how much of the previous hidden state to retain.
 - r_t (*reset gate*): decides how strongly to forget the old hidden state.
- **Fewer parameters** than LSTM, potentially faster convergence.

Notes

GRU Architecture I

Gate Definitions

$$\begin{aligned} z_t &= \sigma(W_z x_t + R_z h_{t-1} + b_z) \quad (\text{update gate}), \\ r_t &= \sigma(W_r x_t + R_r h_{t-1} + b_r) \quad (\text{reset gate}), \\ \tilde{h}_t &= \tanh(W_h x_t + R_h (r_t \odot h_{t-1}) + b_h), \\ h_t &= (1 - z_t) \odot \tilde{h}_t + z_t \odot h_{t-1}. \end{aligned}$$

- z_t blends old vs. new information: when $z_t \approx 1$, we preserve more of h_{t-1} .
- r_t gates how much of h_{t-1} is used in creating \tilde{h}_t .

Parameter Shapes

- $W_* \in \mathbb{R}^{m \times d}$, $R_* \in \mathbb{R}^{m \times m}$, $b_* \in \mathbb{R}^m$.
- Each gate has its own W_* , R_* , b_* , e.g. $W_z, W_r, W_h, R_z, R_r, R_h$.

Notes

Outline

- Introduction to Language Models
- Vocabulary and Tokenization
- Applications
- Parametrization and Estimation
- n-Gram Language Models
- Addressing Data Sparsity in n-Gram Models
- Evaluation Metrics for Language Models
- Toy Example
- Generation Strategies for Language Models
- Feed-Forward Neural Language Models
- Training
- Recurrent Neural Networks (RNNs)
- LSTMs and GRUs
- What's Next?

Notes

Next Lecture: Attention and Transformer models

Attention Mechanisms

- **Purpose:** Enable models to dynamically focus on relevant parts of the input.
- **Types of Attention:**
 - Additive (Bahdanau) Attention
 - Multiplicative (Dot-Product) Attention
 - Scaled Dot-Product Attention
- **Key Equation:**
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$
- **Applications:** Machine translation, text summarization, question answering.

Transformer Architectures

- **Core Components:**
 - Encoder-Decoder Structure
 - Multi-Head Self-Attention
 - Position-wise Feed-Forward Networks
 - Positional Encoding
- **Multi-Head Attention:**
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$
where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$
- **Advantages:**
 - Parallelization over sequence length
 - Captures long-range dependencies effectively
 - Scalable to large datasets and models
- **Impact:** Foundation for state-of-the-art models like BERT, GPT, and more.

Notes
