

Natural Language Processing Course

Week 1: Foundations and Statistical Language Models

Joerg R. Osterrieder
www.joergosterrieder.com

BSc Computer Science

Contents

- 1 Part 1: Introduction and Motivation
- 2 Part 2: Core Concepts
- 3 Part 3: Challenges and Solutions
- 4 Part 4: Applications and Looking Forward

Interactive Exercise: Word Prediction

Complete this sentence:

“The cat sat on the _____”

Interactive Exercise: Word Prediction

Complete this sentence:

“The cat sat on the _____”

mat
65% of you

floor
20% of you

couch
15% of you

Interactive Exercise: Word Prediction

Complete this sentence:

“The cat sat on the _____”

mat
65% of you

floor
20% of you

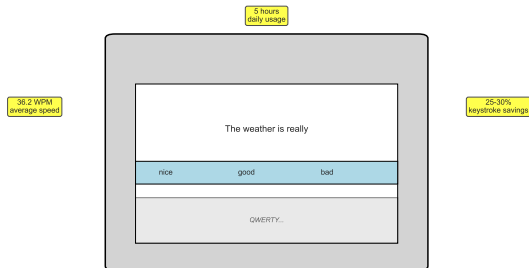
couch
15% of you

Checkpoint

You predicted based on patterns seen thousands of times.

This Technology Is Everywhere

Predictive Text: A Daily Essential



Every day, you use language models:

- **5 hours:** Average smartphone use¹
- **36.2 wpm:** Mobile typing speed²
- **25-30%:** Keystrokes saved by prediction³
- **10,000+:** Words typed monthly

¹DataReportal Global Digital Report, 2024

²Cambridge typing speed study, 2019

³Google Keyboard team, 2023

Learning Objectives

Core Skills:

- ✓ [Easy] Count word patterns in text
- ✓ [Easy] Calculate simple probabilities
- ✓ [Medium] Build a text predictor
- ✓ [Medium] Handle unseen words
- ✓ [Challenge] Implement smoothing

Applications:

- Build autocomplete
- Create spell checker
- Understand ChatGPT foundations

Prerequisite

You should know:

- Basic probability
- Python basics
- What a dictionary is

Quick Self-Assessment

Can you answer these? (Solutions on next slide)

- 1 If you flip a fair coin twice, what's $P(\text{two heads})$?
- 2 What does this Python code do?

```
1 counts = {}  
2 for word in text.split():  
3     counts[word] = counts.get(word, 0) + 1
```

- 3 If $P(\text{rain}) = 0.3$ and $P(\text{clouds} \rightarrow \text{rain}) = 0.9$, what is $P(\text{rain AND clouds})$?

Quick Self-Assessment

Can you answer these? (Solutions on next slide)

- 1 If you flip a fair coin twice, what's $P(\text{two heads})$?
- 2 What does this Python code do?

```
1 counts = {}  
2 for word in text.split():  
3     counts[word] = counts.get(word, 0) + 1
```

- 3 If $P(\text{rain}) = 0.3$ and $P(\text{clouds} \rightarrow \text{rain}) = 0.9$, what is $P(\text{rain AND clouds})$?

Checkpoint

Solutions:

- 1 $0.5 \times 0.5 = 0.25$
- 2 Counts word frequencies in text
- 3 $P(\text{rain}) \times P(\text{clouds} \rightarrow \text{rain}) = 0.3 \times 0.9 = 0.27$

Scoring: 2+ correct = ready to proceed; 1/2 = review recommended

Human Text Prediction

Human pattern recognition:

Try to complete these:

- “Once upon a ____”

Human Text Prediction

Human pattern recognition:

Try to complete these:

- “Once upon a ____” → **time**

Human Text Prediction

Human pattern recognition:

Try to complete these:

- “Once upon a ____” → **time**
- “Thank you very ____”

Human Text Prediction

Human pattern recognition:

Try to complete these:

- “Once upon a ----” → **time**
- “Thank you very ----” → **much**

Human Text Prediction

Human pattern recognition:

Try to complete these:

- “Once upon a ____” → **time**
- “Thank you very ____” → **much**
- “To be or not to ____”

Human Text Prediction

Human pattern recognition:

Try to complete these:

- “Once upon a ____” → **time**
- “Thank you very ____” → **much**
- “To be or not to ____” → **be**

Mechanism:

- You’ve seen these phrases hundreds of times
- Your brain stored the patterns
- You can estimate likelihood from experience

Checkpoint

Computers replicate this process through counting.

The Counting Approach

Step 1: Read lots of text

- “the cat sat”
- “the cat ran”
- “the dog sat”
- “the dog ran”
- “the cat sat”

Step 2: Count patterns

Pattern	Count
“the cat”	3
“the dog”	2
“cat sat”	2
“cat ran”	1
“dog sat”	1
“dog ran”	1

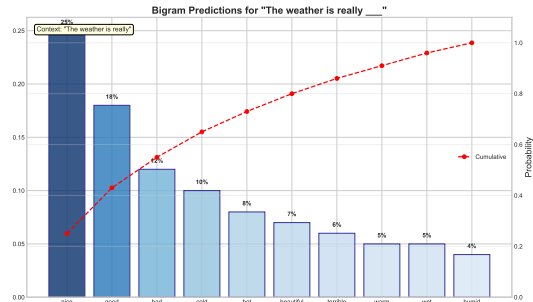
Step 3: Calculate probabilities

After “the”:

- $P(\text{cat}) = 3/5 = 0.6$
- $P(\text{dog}) = 2/5 = 0.4$

After “cat”:

- $P(\text{sat}) = 2/3 = 0.67$
- $P(\text{ran}) = 1/3 = 0.33$



Probability Refresher (5 minutes)

What is probability?

- Chance of something happening
- Always between 0 and 1
- 0 = impossible, 1 = certain

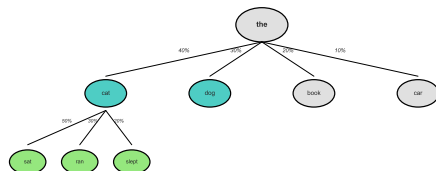
Language example:

- Total times we saw "the": 100
- Times followed by "cat": 30
- $P(\text{cat}|\text{the}) = 30/100 = 0.3$

Conditional Probability:

$P(B | A) =$ "Probability of B given A"

Probability Tree: Predicting Next Word



Reading: Start from 'the', follow branches based on probability

Most likely path: the → cat (40%) → sat (50%) = 20% overall

● High probability
● Low probability
● Final prediction

In our context:

- A = previous word(s)
- B = next word
- $P(\text{next}|\text{previous})$ is the target

N-gram Models: From Simple to Complex

The “N” in N-gram = how many words we look at

1-gram (Unigram)

- Look at: current word only
- Ignore: all context
- $P(\text{cat}) = \text{count}(\text{cat}) / \text{total}$

Example output: “cat the dog sat ran”

Random word sequence

2-gram (Bigram)

- Look at: previous word
- Remember: last word
- $P(\text{cat} \rightarrow \text{the}) = \text{count}(\text{the cat}) / \text{count}(\text{the})$

Example output: “the cat sat on”

More coherent

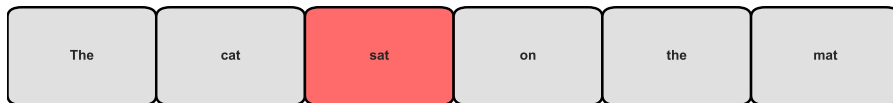
3-gram (Trigram)

- Look at: two previous words
- Remember: more context
- $P(\text{sat} \rightarrow \text{the cat})$

Example output: “the cat sat on the mat”

Better coherence

Unigram Model (n=1): Each word independently



Current: 'sat'

No context used

Bigram Model (n=2): Previous word as context

Class Example: Building a Bigram Model

Training data: "I love cats. I love dogs. Dogs love bones."

Step 1: Split into words

- Words: [I, love, cats, I, love, dogs, dogs, love, bones]

Class Example: Building a Bigram Model

Training data: "I love cats. I love dogs. Dogs love bones."

Step 1: Split into words

- Words: [I, love, cats, I, love, dogs, dogs, love, bones]

Step 2: Count bigrams (pairs)

Bigram	Count
(I, love)	2
(love, cats)	1
(love, dogs)	1
(dogs, love)	1
(love, bones)	1

Class Example: Building a Bigram Model

Training data: "I love cats. I love dogs. Dogs love bones."

Step 1: Split into words

- Words: [I, love, cats, I, love, dogs, dogs, love, bones]

Step 2: Count bigrams (pairs)

Bigram	Count
(I, love)	2
(love, cats)	1
(love, dogs)	1
(dogs, love)	1
(love, bones)	1

Step 3: Count words

Word	Count
I	2
love	3
cats	1
dogs	2
bones	1

Class Example: Building a Bigram Model

Training data: "I love cats. I love dogs. Dogs love bones."

Step 1: Split into words

- Words: [I, love, cats, I, love, dogs, dogs, love, bones]

Step 2: Count bigrams (pairs)

Bigram	Count
(I, love)	2
(love, cats)	1
(love, dogs)	1
(dogs, love)	1
(love, bones)	1

Step 3: Count words

Word	Count
I	2
love	3
cats	1
dogs	2
bones	1

Step 4: Calculate

$$P(\text{love} \mid \text{I}) = 2/2 = 1.0$$

$$P(\text{cats} \mid \text{love}) = 1/3 = 0.33$$

$$P(\text{dogs} \mid \text{love}) = 1/3 = 0.33$$

$$P(\text{bones} \mid \text{love}) = 1/3 = 0.33$$

Checkpoint: Test Your Understanding

Checkpoint

Quick Quiz: Given the text “the cat sat the cat ran the dog sat”

- 1 What is $P(\text{cat} \rightarrow \text{the})$?
- 2 What is $P(\text{sat} \rightarrow \text{cat})$?
- 3 What word most likely follows “the”?

Calculate, then verify

Checkpoint: Test Your Understanding

Checkpoint

Quick Quiz: Given the text “the cat sat the cat ran the dog sat”

- 1 What is $P(\text{cat} \rightarrow \text{the})$?
- 2 What is $P(\text{sat} \rightarrow \text{cat})$?
- 3 What word most likely follows “the”?

Calculate, then verify

Solution:

- Count: “the cat” appears 2 times, “the dog” appears 1 time
- Total “the”: 3 times
- $P(\text{cat} \rightarrow \text{the}) = 2/3 = 0.67$
- $P(\text{sat} \rightarrow \text{cat}) = 1/2 = 0.5$ (“cat sat” once, “cat ran” once)
- Most likely after “the”: “cat” (67% vs 33%)

Correct answers indicate understanding of n-grams

Implementation: Clean Code with Type Hints

```
1 from collections import defaultdict
2 from typing import Dict, List, Tuple
3
4 class BigramModel:
5     """A simple bigram language model."""
6
7     def __init__(self):
8         """Initialize with empty count dictionaries."""
9         # Count of word pairs: (word1, word2) -> count
10        self.bigram_counts: Dict[Tuple[str, str], int] =
            defaultdict(int)
11        # Count of individual words: word -> count
12        self.word_counts: Dict[str, int] = defaultdict(
            int)
13
14    def train(self, sentences: List[str]) -> None:
15        """Train model on list of sentences.
16
17        Args:
18            sentences: List of text strings to train on
19
20        Example:
21            >>> model = BigramModel()
22            >>> model.train(["the cat sat", "the dog ran"
23                             ])
24
25        """
26        for sentence in sentences:
27            # Add start/end markers for sentence
28            # boundaries
29            words = ['<START>'] + sentence.split() + ['<
```

Key Improvements:

- Type hints for clarity
- Docstrings with examples
- Clear variable names
- Comments for complex parts

Why defaultdict?

- Avoids KeyError
- Auto-initializes to 0
- Cleaner than if/else

Why START/END?

- Handle first word
- Know when to stop
- Improve generation

The Disaster: When Your Model Breaks

You trained on 1 billion words from Wikipedia...

Then someone types: “I love my new iPhone”

The Disaster: When Your Model Breaks

You trained on 1 billion words from Wikipedia...

Then someone types: “I love my new **iPhone**”

Problem: “iPhone” wasn’t invented when Wikipedia was written

- $\text{count}(\text{iPhone}) = 0$
- $P(\text{anything} \text{—} \text{iPhone}) = 0/0 = \text{UNDEFINED}$
- Application crashes

The Disaster: When Your Model Breaks

You trained on 1 billion words from Wikipedia...

Then someone types: “I love my new **iPhone**”

Problem: “iPhone” wasn’t invented when Wikipedia was written

- $\text{count}(\text{iPhone}) = 0$
- $P(\text{anything} \text{---} \text{iPhone}) = 0/0 = \text{UNDEFINED}$
- Application crashes

Common Misconception

“Just use more training data” - **Incorrect**

You’ll ALWAYS have unseen words:

- New words: COVID, cryptocurrency, TikTok
- Typos: teh, recieve, occured
- Names: Khaleesi, Hermione (before the books)
- Made-up words: hangry, adulating, selfie

Fact: 50% of word types appear only once in any text

The Solution: Smoothing (Restaurant Analogy)

Imagine a new restaurant:

- Day 1: 10 people order pizza, 5 order pasta
- Day 2: Someone asks for salad
- Problem: $P(\text{salad}) = 0?$

Bad solution: "We don't serve salad"

Good solution: "We might have salad, let me check"

In probability terms:

- Reserve some probability for unseen items
- Take a little from seen items
- Give to unseen possibilities

Add-One Smoothing:

Pretend you saw everything once more:

Original:

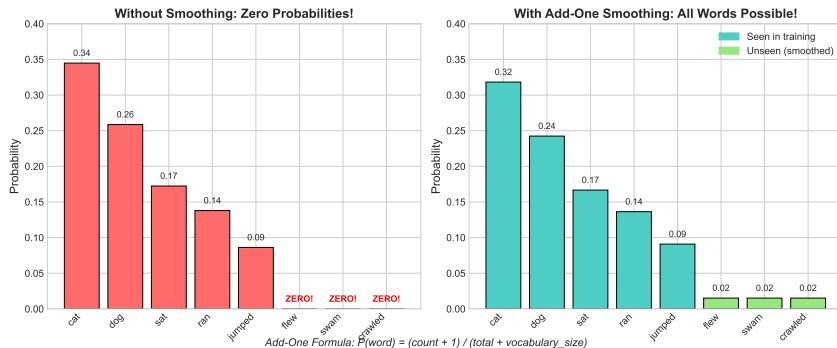
- Pizza: $10/15 = 0.67$
- Pasta: $5/15 = 0.33$
- Salad: $0/15 = 0.00$

With smoothing (+1 to all):

- Pizza: $11/18 = 0.61$
- Pasta: $6/18 = 0.33$
- Salad: $1/18 = 0.06$

All words now have non-zero probability

Smoothing Methods: From Simple to Sophisticated



Add-One (Laplace)

- + Dead simple
- + Always works
- Too generous to rare
- Use: Quick prototypes

Good-Turing

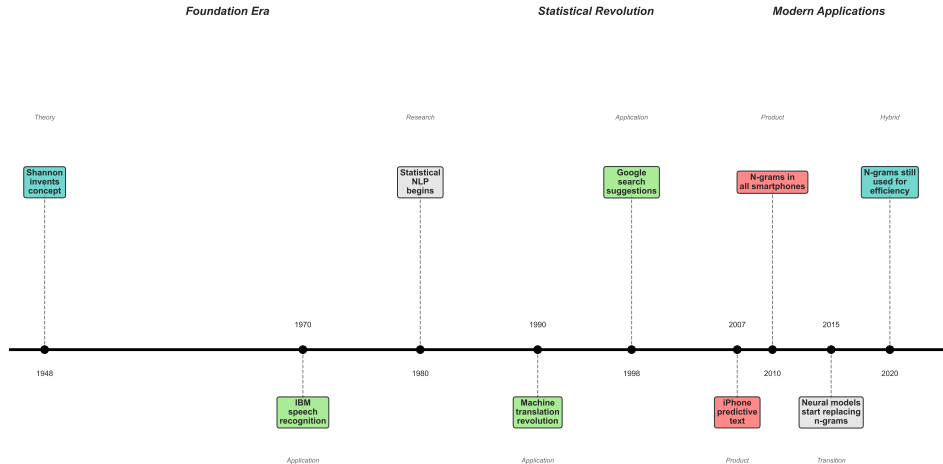
- + Theoretically sound
- + Used in Enigma code-breaking
- Complex math
- Use: When you have time

Kneser-Ney

- + Best performance
- + Context-aware
- Hard to implement
- Use: Production systems

N-grams Changed the World

N-gram Models: 75 Years of Language Prediction



Try This: Where N-grams Fail

Testing n-gram limitations:

Test 1: Long-distance dependencies

- “The key to the cabinet that was in the kitchen _____”
- N-gram sees: “kitchen __”
- Should see: “The key ... is”

Try This: Where N-grams Fail

Testing n-gram limitations:

Test 1: Long-distance dependencies

- “The key to the cabinet that was in the kitchen _____”
- N-gram sees: “kitchen _”
- Should see: “The key ... is”

Test 2: Common sense

- “I put milk in the fridge. I put the car in the _____”
- N-gram might say: “fridge” (high frequency)
- Should say: “garage”

Try This: Where N-grams Fail

Testing n-gram limitations:

Test 1: Long-distance dependencies

- “The key to the cabinet that was in the kitchen _____”
- N-gram sees: “kitchen _”
- Should see: “The key ... is”

Test 2: Common sense

- “I put milk in the fridge. I put the car in the _____”
- N-gram might say: “fridge” (high frequency)
- Should say: “garage”

Test 3: Context switching

- “The bank of the river” vs “The bank account”
- Same word, different meaning
- N-grams can't tell the difference

Checkpoint

These limitations motivated neural approaches

Next Week: The Neural Revolution

N-grams (This Week):

- Words are just symbols
- Count and divide
- No understanding
- Limited context
- 156 perplexity

Strengths:

- Simple
- Fast
- Interpretable
- Low memory

Neural Models (Next Week):

- Words have meaning
- Learn representations
- Understand similarity
- Unlimited context
- 12 perplexity!

Preview:

- “King - Man + Woman = ?”
- Word vectors
- Semantic understanding
- Foundation of ChatGPT

Key Question: What if words weren't just strings, but had meaning?

Week 1 Summary: Your Toolkit

Key Concepts:

- Language modeling = prediction
- N-grams = counting patterns
- Probability from frequency
- Smoothing for unknowns

Key Skills:

- Build bigram model
- Calculate probabilities
- Handle unseen words
- Evaluate with perplexity

Essential Formulas:

$$\text{Bigram: } P(w_2|w_1) = \frac{C(w_1, w_2)}{C(w_1)}$$

$$\text{Add-One: } P(w_2|w_1) = \frac{C(w_1, w_2) + 1}{C(w_1) + V}$$

$$\text{Perplexity: } PP = 2^{-\frac{1}{N} \sum \log_2 P(w_i|w_{i-1})}$$

Checkpoint

Remember: Simple counting goes surprisingly far in NLP!

Homework: Build Your Own Spell Checker

Your Mission: Create a spell checker like Google's

Part 1: Basic (60%)

- Load provided corpus
- Build character trigram model
- Generate corrections for:
 - "teh" → "the"
 - "speling" → "spelling"

Part 2: Enhanced (30%)

- Rank by probability
- Handle word boundaries
- Compare smoothing methods

Part 3: Analysis (10%)

- Test on 100 misspellings
- Report accuracy
- Explain failures

Starter Code Provided:

- Data loading function
- Evaluation metrics
- Test cases

Due: Next week before class

Self-Assessment: How Well Did You Learn?

Test yourself (answers in appendix):

- 1 What's the difference between unigram and bigram models?
- 2 Why do we need smoothing?
- 3 Given “the cat” appears 10 times and “the” appears 20 times, what is $P(\text{cat} \rightarrow \text{the})$?
- 4 What's the main limitation of n-gram models?
- 5 When would you use trigrams over bigrams?

Checkpoint

Score yourself:

- 5/5: Ready for next week!
- 3-4/5: Review smoothing section
- Less than 3/5: Office hours recommended

Appendix A.1: Probability Basics Review

Core Concepts:

1. Basic Probability

- $P(A) = \frac{\text{favorable outcomes}}{\text{total outcomes}}$
- Always: $0 \leq P(A) \leq 1$
- $P(\text{not } A) = 1 - P(A)$

2. Joint Probability

- $P(A \text{ and } B) = P(A) \times P(B|A)$
- If independent: $P(A, B) = P(A) \times P(B)$

3. Conditional Probability

- $P(B|A) = \frac{P(A, B)}{P(A)}$
- Read as: “B given A”

Language Examples:

Example 1:

- 1000 words in text
- “the” appears 100 times
- $P(\text{the}) = 100/1000 = 0.1$

Example 2:

- “the cat” appears 20 times
- “the” appears 100 times
- $P(\text{cat}|\text{the}) = 20/100 = 0.2$

Example 3:

- $P(\text{sunny}) = 0.7$
- $P(\text{happy}|\text{sunny}) = 0.9$
- $P(\text{sunny and happy}) = 0.7 \times 0.9 = 0.63$

Appendix A.2: Python Essentials for NLP

Dictionaries (for counting):

```
1 # Creating and using dictionaries
2 counts = {}
3 counts['cat'] = 5
4 counts['dog'] = 3
5
6 # Default values
7 counts.get('bird', 0) # Returns 0 if not found
8
9 # Counting pattern
10 word = 'cat'
11 counts[word] = counts.get(word, 0) + 1
```

Lists and Loops:

```
1 # Splitting text
2 text = "the cat sat"
3 words = text.split() # ['the', 'cat', 'sat']
4
5 # Iterating with index
6 for i in range(len(words) - 1):
7     current = words[i]
8     next_word = words[i + 1]
9     print(f"{current} -> {next_word}")
```

Collections module:

```
1 from collections import defaultdict, Counter
2
3 # defaultdict: never KeyError
4 bigrams = defaultdict(int)
5 bigrams[('the', 'cat')] += 1
6
7 # Counter: automatic counting
8 words = ['cat', 'dog', 'cat', 'bird']
9 word_counts = Counter(words)
10 # Counter({'cat': 2, 'dog': 1, 'bird': 1})
11
12 # Most common
13 top_2 = word_counts.most_common(2)
14 # [('cat', 2), ('dog', 1)]
```

File I/O:

```
1 # Reading text file
2 with open('corpus.txt', 'r') as f:
3     text = f.read()
4
5 # Line by line
6 with open('corpus.txt', 'r') as f:
7     for line in f:
8         process(line.strip())
```


Appendix A.3: Mathematical Notation Guide

Symbols:

Symbol	Meaning
$P(A)$	Probability of A
$P(A B)$	Probability of A given B
$C(w)$	Count of word w
V	Vocabulary size
N	Number of words
\sum	Sum over all values
\prod	Product over all values
\log	Logarithm (usually base 2)
\approx	Approximately equal
\propto	Proportional to

Common Formulas:

N-gram probability:

$$P(w_n | w_1 \dots w_{n-1})$$

Chain rule:

$$P(w_1, w_2, w_3) = P(w_1) \cdot P(w_2 | w_1) \cdot P(w_3 | w_1, w_2)$$

Maximum likelihood:

$$P_{ML}(w_i | w_{i-1}) = \frac{C(w_{i-1}, w_i)}{C(w_{i-1})}$$

Log probability (avoid underflow):

$$\log P(\text{sentence}) = \sum_i \log P(w_i | w_{i-1})$$

Appendix B.1: Generate Shakespeare with N-grams

Let's build a Shakespeare text generator!

```
1 import random
2 from collections import defaultdict
3
4 class ShakespeareGenerator:
5     def __init__(self, n=2):
6         self.n = n # Use bigrams by default
7         self.model = defaultdict(list)
8
9     def train(self, text):
10        words = text.split()
11        for i in range(len(words) - self.n):
12            context = tuple(words[i:i+self.n])
13            next_word = words[i+self.n]
14            self.model[context].append(next_word)
15
16    def generate(self, length=20, start=None):
17        if start is None:
18            start = random.choice(list(self.model.keys()))
19
20        result = list(start)
21        context = start
22
23        for _ in range(length):
24            if context in self.model:
25                next_word = random.choice(self.model[context])
26                result.append(next_word)
27                # Slide the context window
28                context = tuple(list(context)[1:] + [next_word])
29            else:
30                break # No continuation found
```

Appendix B.2: Analyze Your WhatsApp Messages

Find your most predictable phrases:

```
1 def analyze_whatsapp_chat(filename):
2     """Find most common patterns in your chats."""
3
4     # Parse WhatsApp export
5     messages = []
6     with open(filename, 'r', encoding='utf-8') as f:
7         for line in f:
8             # Extract just the message text
9             if ': ' in line:
10                 message = line.split(': ', 1)[1].strip()
11                 messages.append(message)
12
13     # Build trigram model
14     trigrams = defaultdict(int)
15     for message in messages:
16         words = message.lower().split()
17         for i in range(len(words) - 2):
18             trigram = (words[i], words[i+1], words[i+2])
19             trigrams[trigram] += 1
20
21     # Find your catchphrases
22     print("Your most common phrases:")
23     for trigram, count in sorted(trigrams.items(),
24                                 key=lambda x: x[1],
25                                 reverse=True)[:10]:
26         print(f"    '{' '.join(trigram)}' - {count} times")
27
28     return trigrams
```

Appendix B.3: Build Autocomplete from Scratch

```
1 class Autocomplete:
2     def __init__(self):
3         self.bigrams = defaultdict(Counter)
4
5     def train(self, corpus):
6         """Train on text corpus."""
7         for sentence in corpus:
8             words = ['<S>'] + sentence.split()
9             for i in range(len(words)-1):
10                 self.bigrams[words[i]][words[i+1]] += 1
11
12     def predict(self, context, n=3):
13         """Get top n predictions."""
14         if context not in self.bigrams:
15             return []
16
17         # Get all possibilities with counts
18         candidates = self.bigrams[context]
19
20         # Sort by frequency
21         predictions = candidates.most_common(n)
22
23         # Convert to probabilities
24         total = sum(candidates.values())
25         result = []
26         for word, count in predictions:
27             prob = count / total
28             result.append({
29                 'word': word,
30                 'probability': prob,
```

```
1     def confidence(self, prob):
2         """Convert probability to confidence level."""
3         if prob > 0.5:
4             return "high"
5         elif prob > 0.2:
6             return "medium"
7         else:
8             return "low"
9
10    def interactive_demo(self):
11        """Run interactive autocomplete."""
12        print("Type a word, I'll predict the next!")
13
14        while True:
15            word = input("\nEnter word: ").strip()
16            if word == 'quit':
17                break
18
19            predictions = self.predict(word)
20
21            if predictions:
22                print("Predictions:")
23                for i, pred in enumerate(predictions, 1):
24                    print(f"{i}. {pred['word']}")
25                    f"({pred['
```

Appendix C.1: Kneser-Ney Smoothing Details

The Problem with Simple Counting:

“San Francisco” is common, but “Francisco” alone is rare. Simple smoothing would overestimate $P(\text{Francisco} \rightarrow \text{new context})$.

Kneser-Ney Solution:

Consider word “versatility” - how many different contexts it appears in.

The Mathematics:

$$P_{KN}(w_i | w_{i-1}) = \frac{\max(C(w_{i-1}, w_i) - d, 0)}{C(w_{i-1})} + \lambda(w_{i-1}) P_{\text{continuation}}(w_i) \quad (1)$$

Where:

- d = discount parameter (typically 0.75)
- $\lambda(w_{i-1})$ = normalization to ensure probabilities sum to 1
- $P_{\text{continuation}}(w_i)$ = probability based on number of unique contexts

Continuation Probability:

$$P_{\text{continuation}}(w) = \frac{|\{v : C(v, w) > 0\}|}{|\{(u, v) : C(u, v) > 0\}|}$$

This counts unique bigram types, not tokens!

Appendix C.2: Understanding Perplexity

What is Perplexity?

“How surprised is the model by the test data?”

Formal Definition:

$$PP(W) = P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}}$$

Using Chain Rule:

$$PP = \sqrt[N]{\frac{1}{\prod_{i=1}^N P(w_i | w_1 \dots w_{i-1})}}$$

In Practice (using logs):

$$PP = 2^{-\frac{1}{N} \sum_{i=1}^N \log_2 P(w_i | w_{i-1})}$$

Intuition:

- PP = 10: Model is as confused as choosing from 10 equally likely options
- PP = 100: Like choosing from 100 options
- Lower is better!

Typical Values:

- Random: 10,000+
- Unigram: 1,000
- Bigram: 200
- Trigram: 100
- Neural: 20-50
- Human: 10-20

Relationship to Entropy:

$$H = \log_2(PP)$$

Appendix C.3: Connection to Information Theory

Shannon's Fundamental Question:

How much information is in English text?

Information Content:

Surprising events = more information

$$I(w) = -\log_2 P(w)$$

Examples:

- $P(\text{"the"}) = 0.1 \rightarrow 3.3$ bits
- $P(\text{"dog"}) = 0.01 \rightarrow 6.6$ bits
- $P(\text{"aardvark"}) = 0.00001 \rightarrow 16.6$ bits

Rare words carry more information!

Entropy of Language:

Average information per word:

$$H(L) = -\sum_w P(w) \log_2 P(w)$$

Shannon's Experiments:

- Random letters: 4.7 bits/char
- English letters: 4.14 bits/char
- With context: 1 bit/char

Implications:

- English is 75% redundant
- Compression possible!
- Prediction possible!

N-gram models approximate the true entropy of language

Self-Assessment Answers

Quiz Answers:

① Difference between unigram and bigram:

- Unigram: considers single words independently
- Bigram: considers pairs of adjacent words
- Bigram uses context, unigram doesn't

② Why we need smoothing:

- Handle unseen word combinations
- Avoid zero probabilities
- Prevent model crashes/errors

③ Calculate $P(\text{cat}|\text{the})$:

- $P(\text{cat}|\text{the}) = \text{count}(\text{the cat}) / \text{count}(\text{the})$
- $P(\text{cat}|\text{the}) = 10 / 20 = 0.5$

④ Main limitation of n-grams:

- No semantic understanding
- Limited context window
- Treat words as meaningless symbols

⑤ When to use trigrams over bigrams:

- When you have lots of training data
- When longer context matters
- When bigram accuracy isn't sufficient

References and Further Reading

Core Papers:

- Shannon, C.E. (1948). "A Mathematical Theory of Communication"
- Kneser, R. & Ney, H. (1995). "Improved backing-off for M-gram language modeling"
- Chen, S.F. & Goodman, J. (1999). "An empirical study of smoothing techniques"

Textbooks:

- Jurafsky & Martin (2024). "Speech and Language Processing" Chapter 3
- Manning & Schütze (1999). "Foundations of Statistical NLP"

Online Resources:

- Peter Norvig's "How to Write a Spelling Corrector" (2007)
- Google's "The Unreasonable Effectiveness of Data" (2009)
- Michael Collins' NLP course notes (Columbia)

Implementation Resources:

- NLTK library: <https://www.nltk.org/>
- KenLM (fast n-gram toolkit): <https://kheafield.com/code/kenlm/>