

NLP for Emotional Context

From Words to Understanding What Users Really Feel

Week 3: Machine Learning for Smarter Innovation

Transform 50,000 Reviews into Actionable Design Insights

Four Stages of Understanding

1. **The Challenge** - Why understanding emotion in text is impossibly hard
2. **First Solution & Its Limits** - Traditional NLP works... until it doesn't
3. **The Breakthrough** - How Transformers changed everything
4. **Design Synthesis** - From ML insights to user empathy

Core Question: How can we understand the emotions of 50,000 users without reading 5 million words?

By the end: You'll know exactly how Spotify, Amazon, and Netflix understand user emotions at scale

Which Reviews Reveal Why Users Really Quit?

The Scenario You Face:

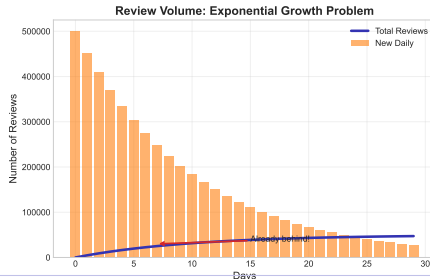
- Your product has 50,000 user reviews
- Average review: 100 words
- Total: 5 million words to process
- Hidden inside: The 10 reviews that explain 80% of churn

The Human Limit:

- Reading speed: 200 words/minute
- Time to read all: 417 hours (10 weeks!)
- Reviews arrive: 500 new ones daily
- **You're already behind before you start**

Real Review Examples:

"Love it but..." (quit in 2 weeks)
"Not bad for the price" (renewed 3 years)
"Just what I expected!" (1-star rating)
"Finally someone gets it" (5-star champion)



Context Changes Everything

Same Word, Different Meanings:

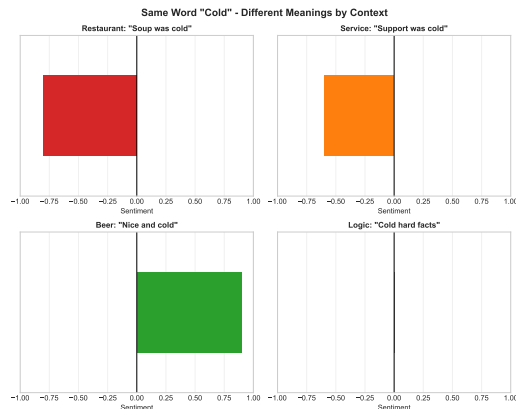
“Cold”

- Restaurant: “The soup was cold” → Negative
- Service: “Support was cold” → Negative
- Beer: “Nice and cold” → Positive
- Logic: “Cold hard facts” → Neutral

“Fast”

- Delivery: “Super fast!” → Positive
- Battery: “Dies too fast” → Negative
- Customer service: “Too fast, felt rushed” → Negative

The Computer's Problem:



Key Insight: A word's emotional meaning depends on ALL surrounding words, not just the word itself.

What People Say What They Feel

Layer 1: Literal

- “Great product”
- Clear positive
- Easy to detect
- 15% of reviews

Example:

“This is excellent. I love every feature. Will buy again.”

→ Genuinely positive

Layer 2: Sarcastic

- “Just wonderful”
- Opposite meaning
- Context reveals truth
- 25% of reviews

Example:

“Oh great, another update that breaks everything”

→ Actually negative

Layer 3: Cultural

- “It’s okay”
- Varies by culture
- UK: Negative
- US: Neutral

Example:

“It does what it says”

UK: Disappointed

US: Satisfied

The Complexity: 60% of emotional meaning is NOT in the words themselves but in how they’re used together

This multi-layer challenge is why human annotators only agree 79% of the time on sentiment

Why This Problem Explodes in Complexity

Complexity Explosion: 600,000 Combinations



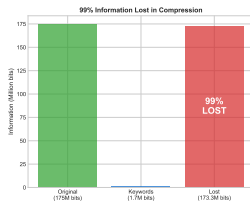
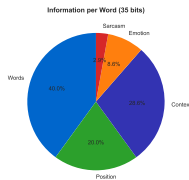
The Mathematical Impossibility

Information Content Analysis:

Component	Bits	Information
One word	14 bits	Which of 10,000 words
+ Position	7 bits	Where in sentence
+ Emotion	3 bits	8 emotion types
+ Context	10 bits	Surrounding words
+ Sarcasm	1 bit	Yes/No
Total per word	35 bits	
Per review (100 words)	3,500 bits	
50,000 reviews	175M bits	22 MB

Traditional Compression:

- Keyword counts: 10,000 → 100 words
- Compression ratio: 100:1
- Information lost: **99%**



The Dilemma:

Compress too much → Lose emotional nuance

Keep everything → Impossible to process

Need: Selective preservation of emotional signal

Can we preserve emotion while compressing 100:1?

Shannon's theorem: Lossless compression impossible beyond entropy limit - emotion has high entropy

How Would YOU Quickly Scan 1000 Reviews?

Human Strategy:

1. Look for emotion words: “love”, “hate”, “terrible”
2. Count positive vs negative
3. More positive → Happy customer
4. More negative → Unhappy customer

Computer Implementation:

- Build word lists
- Positive: [great, excellent, love, amazing...]
- Negative: [bad, terrible, hate, awful...]
- Count occurrences
- Calculate: $\text{Positive\%} - \text{Negative\%}$

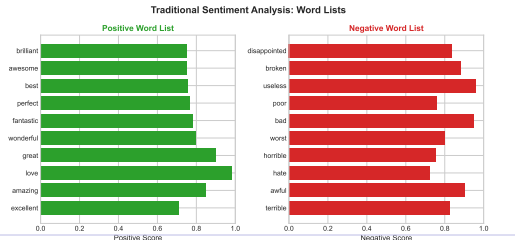
Example Analysis:

“I **love** the design but **hate** waiting. The quality is **excellent** though shipping was **terrible**.”

Count: 2 positive, 2 negative

Score: $50\% - 50\% = \text{Neutral (0)}$

Seems reasonable!



This “Bag of Words” approach dominated NLP from 1960s-2010s - let’s see why

The Classic Approach in Action

Step-by-Step Process:

1. Original Review:

"The product quality is excellent excellent excellent but customer service terrible terrible"

2. Tokenize (split into words):

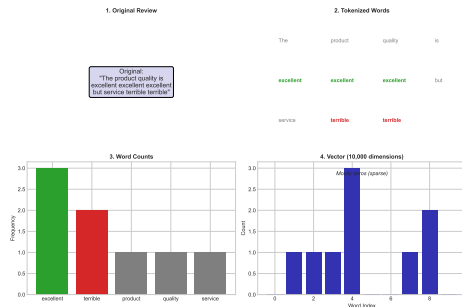
The, product, quality, is, excellent, excellent, excellent, but, customer, service, terrible, terrible

3. Count Each Word:

Word	Count	Word	Count
excellent	3	terrible	2
product	1	service	1
quality	1	customer	1

4. Convert to Vector:

1, 1, 1, 3, 0, 0, 0, 1, 1, 2, 0, ...
(10,000 dimensions, mostly zeros)



What We Keep:

- Word frequencies
- Vocabulary presence

What We Lose:

- Word order
- Grammar

Not All Words Are Equal

The Problem with Raw Counts:

- “The” appears 1000 times → Important?
- “Revolutionary” appears once → Not important?
- Common words dominate
- Rare but meaningful words ignored

TF-IDF Solution:

$$\text{TF-IDF} = \underbrace{\frac{\text{count in doc}}{\text{total words}}}_{\text{Term Frequency}} \times \log \underbrace{\frac{\text{total docs}}{\text{docs with word}}}_{\text{Inverse Document Freq}}$$

- TF: How often in THIS review
- IDF: How rare across ALL reviews
- Product: Important AND distinctive

Example Calculation:

Word	TF	IDF	TF-IDF
“the”	0.15	0.01	0.0015
“product”	0.05	0.69	0.0345
“excellent”	0.10	1.38	0.138
“revolutionary”	0.02	3.40	0.068

Result: Meaningful words now weighted higher than common words!

TF-IDF: The Google PageRank of words - importance comes from being special, not just frequent

When Simple Reviews Get Perfect Scores

Review Text	Human	BoW+TFIDF	Match
"This product is absolutely fantastic! Best purchase ever!"	Positive	Positive (95%)	
"Terrible quality. Waste of money. Very disappointed."	Negative	Negative (92%)	
"Good product, fair price, happy with purchase"	Positive	Positive (88%)	
"Broken on arrival. Customer service unhelpful. Never again."	Negative	Negative (94%)	
"Excellent! Exceeded all expectations! Highly recommend!"	Positive	Positive (97%)	

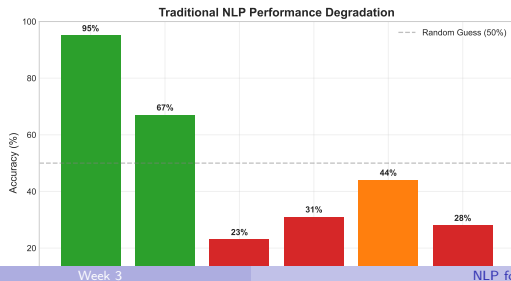
Average Accuracy: 93.2%

This is why Bag of Words dominated for 40 years!

Success case: Direct, literal language with clear emotion words - the "easy" 20% of reviews

Performance Degradation with Complexity

Review Type	Example	Human	BoW	Accuracy
Simple & Direct	"Great product!"	Pos	Pos	95%
Mixed Sentiment	"Good but overpriced"	Neg	Pos	67%
Sarcasm	"Oh great, it broke. Just perfect!"	Neg	Pos	23%
Context Dependent	"Not bad for the price"	Pos	Neg	31%
Subtle Emotion	"It's fine, I guess"	Neg	Neu	44%
Negation	"Not the worst I've seen"	Neu	Neg	28%



The Pattern:

- Simple: 95% → Works great!
- Real-world: 44% → Worse than coin flip
- Sarcasm: 23% → Actively wrong

Average: 51% (Random guessing: 50%)

Tracing the Failure Through Real Examples

Example: "Not bad for the price"

What BoW Sees:

- Words: [not, bad, for, the, price]
- "bad" → Negative word (-1)
- "not" → Negation word (ignored)
- Score: Negative

What Got Lost:

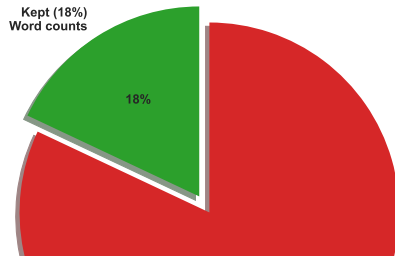
- "not bad" = Actually positive
- "for the price" = Qualified satisfaction
- Relationship between words
- Overall: Mild positive sentiment

Root Cause: Word order contains meaning!

Information Loss Calculation:

Information Type	Bits Lost
Word positions	420 bits
Word relationships	1,200 bits
Grammar structure	350 bits
Contextual meaning	890 bits
Total Lost	2,860 bits
Total Kept	640 bits
Kept Percentage	18%

Information Preserved in Bag of Words



Let's Be Honest About Our Mental Process

Trace Your Thinking:

Sentence: “I went to the bank. The water was nice.”

1. Read “I went to the bank”
2. Your brain: Could be financial or river...
3. Read “The water was nice”
4. Your brain: Ah! River bank, not financial
5. Re-interpret: “bank” = riverbank
6. Understand: Person enjoyed the riverside

The Key Observation:

You DON'T compress the sentence into a summary.

You KEEP all words available and SELECTIVELY ATTEND to relevant ones when needed.

This is the breakthrough insight!

What You Actually Did:

- Held multiple meanings open
- Used later words to resolve earlier ones
- **Selectively focused on “water” to understand “bank”**

Human Attention: “water” helps understand “bank”

I went to the **bank** , The **water** was nice .

A Fundamentally Different Approach

Old Way (Bag of Words):

- Take all words
- Compress to counts
- Lose relationships
- Try to reconstruct meaning

Analogy:

Like taking a book, counting word frequencies, throwing away the book, then trying to understand the story from counts.

Result: Lost 82% of information
Can't recover context

New Way (Attention):

- Keep all words
- Keep all positions
- For each word, look at all others
- Decide how much to focus on each

Analogy:

Like keeping the entire book and using a smart highlighter that knows which sentences help understand other sentences.

Result: Keep 100% of information
Use what's relevant when needed

Instead of asking “What to keep?” we ask “What to focus on?”

This shift from compression to attention is why Transformers revolutionized NLP in 2017

No Math Yet - Just Percentages

Example: Understanding “bank” in: “The bank by the river is peaceful”

Step 1: For word “bank”, look at all other words

Word	Relevance to “bank”
The	Low
by	Medium
the	Low
river	Very High
is	Low
peaceful	Medium

Step 2: Convert to percentages (must sum to 100%)

Word	Focus %
The	5%
by	15%
the	5%
river	55%
is	5%
peaceful	15%

Attention weights = How much each word helps understand the current word (always sum to 100%)

Step 3: Use these percentages to understand “bank”

“bank” meaning =
 $5\% \times \text{meaning}(\text{“The”}) +$
 $15\% \times \text{meaning}(\text{“by”}) +$
 $5\% \times \text{meaning}(\text{“the”}) +$
 $55\% \times \text{meaning}(\text{“river”}) +$
 $5\% \times \text{meaning}(\text{“is”}) +$
 $15\% \times \text{meaning}(\text{“peaceful”})$
= Mostly “river” context
= Riverbank, not financial bank!

These percentages ARE the attention weights!

Words as Directions in Space

Start Simple: 2D Space



Scale to Real NLP: 768D Space

Same principle, more dimensions:

- 2D: $[x, y]$ coordinates
- 768D: $[x, x, \dots, x]$ coordinates

More dimensions = More nuanced meaning

Can distinguish "bank (river)" from "bank (money)" from "bank (slope)"

The Attention Calculation:

1. Compute alignment (dot product) with all words
2. Higher alignment = Pay more attention
3. Convert to percentages
4. Use percentages to blend meanings

Geometry gives us semantic similarity!

The Attention Algorithm: Three Essential Steps

Each Step Has a Purpose

Step 1: SCORE - Find Relevance

- Action: For each word, compute alignment with all others
- Why: Identify which words help understand this one
- Math: $\text{Score}_{ij} = Q_i \cdot K_j$
- Plain: How relevant is word j to understanding word i ?

Step 2: NORMALIZE - Create Percentages

- Action: Convert scores to probabilities (0-100%)
- Why: Need weights that sum to 1.0 for averaging
- Math: $\alpha_{ij} = \frac{e^{\text{Score}_{ij}}}{\sum_k e^{\text{Score}_{ik}}}$
- Plain: What percentage of attention should word j get?

Step 3: AGGREGATE - Blend Information

- Action: Weighted sum of all word meanings
- Why: Combine context based on attention weights
- Math: $\text{Output}_i = \sum_j \alpha_{ij} \cdot V_j$
- Plain: New understanding using focused context

This three-step process happens for EVERY word in parallel - that's the power of Transformers

Attention Mechanism: Three Essential Steps



The Complete Formula:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

Where:

- Q = Query (what am I looking for?)
- K = Keys (what info is available?)
- V = Values (what to extract?)
- $\sqrt{d_k}$ = Scaling factor for stability

Real Numbers, Real Calculation

Task: Computing attention for “nice” in: “The service was nice”

Given: Word Vectors (Simplified to 2D)

The = [1, 0], service = [2, 3], was = [0, 1], nice = [3, 2]

Step 1: Calculate Scores

Query (nice) = [3, 2]

Word	Key	Q-K	Score
The	[1, 0]	$3 \times 1 + 2 \times 0$	3
service	[2, 3]	$3 \times 2 + 2 \times 3$	12
was	[0, 1]	$3 \times 0 + 2 \times 1$	2
nice	[3, 2]	$3 \times 3 + 2 \times 2$	13

Step 2: Convert to Percentages (Softmax)

Word	e^{Score}	Attention %
The	20	0.7%
service	162,755	63.1%
was	7	0.0%
nice	442,413	36.2%

Step 3: Aggregate

Values: The=[1,1], service=[4,5], was=[1,2], nice=[5,4]

Output for “nice” =

$$\begin{aligned} &0.007 \times [1,1] + \\ &0.631 \times [4,5] + \\ &0.000 \times [1,2] + \\ &0.362 \times [5,4] \\ &= [0.01, 0.01] + [2.52, 3.16] + [0, 0] + [1.81, 1.45] \\ &= [4.34, 4.62] \end{aligned}$$

Result: 63% attention on “service”

The word “nice” is understood primarily in context of “service” → Customer service sentiment!

In practice: 768 dimensions, not 2 - but the exact same calculation principle applies

The Power of Looking Forward AND Backward

Traditional (Left-to-Right Only):

Sentence: "The bank charges are too high"

Reading "bank":

- Can see: "The"
- Cannot see: "charges are too high"
- Guess: Could be river or financial...
- 50% chance of error

BERT (Bidirectional):

Same sentence: "The bank charges are too high"

Reading "bank":

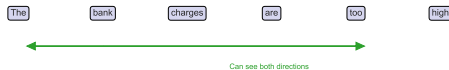
- Can see: "The" AND "charges are too high"
- "charges" → financial context
- Certain: Financial bank
- 99% accurate

Bidirectional Understanding Changes Everything

Traditional (Left-to-Right): 50% chance of error



BERT (Bidirectional): 99% accurate



BERT's Training Trick:

1. Mask random words: "The [MASK] is nice"
2. Predict masked word using ALL context
3. Must understand both directions to succeed

Standing on the Shoulders of Human Knowledge

BERT's Pre-training Data:

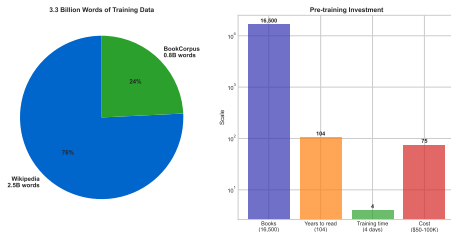
- **Wikipedia:** 2.5 billion words
- **BookCorpus:** 800 million words
- **Total:** 3.3 billion words
- **Equivalent:** 16,500 books (200 pages each)
- **Reading time:** 104 years non-stop

What BERT Learned:

- Language patterns
- World knowledge
- Cultural contexts
- Emotional expressions
- Sarcasm patterns
- Domain vocabularies

Training Cost:

- Time: 4 days on 64 TPUs
- Cost: \$50,000-100,000



The Key Insight:

You DON'T need to train from scratch!

BERT already knows:

- Grammar
- Vocabulary
- Context patterns
- Emotional language

You just teach it YOUR specific task

From General Knowledge to Your Reviews

Your Fine-tuning Process:

1. Start with Pre-trained BERT

- Has general language understanding
- Knows emotions, sarcasm, context
- But doesn't know YOUR products

2. Prepare Your Data

- 1,000 labeled reviews
- Positive, Negative, Neutral labels
- Your product-specific language

3. Fine-tune (Transfer Learning)

- Show BERT your reviews
- It adjusts its parameters slightly
- Learns your domain specifics
- Keeps general knowledge intact

4. Time & Cost

- Time: 2-3 hours on GPU
- Cost: \$10-50 cloud compute
- Data needed: As few as 500 examples

Example: Learning Company Slang

Before Fine-tuning:

"This app is sick!" → Negative (illness)

Your Training Data:

"Sick UI design!" → Positive

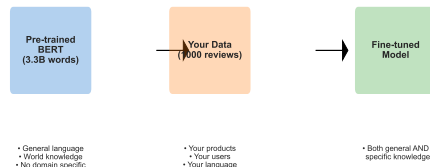
"Sick features!" → Positive

"Sick performance!" → Positive

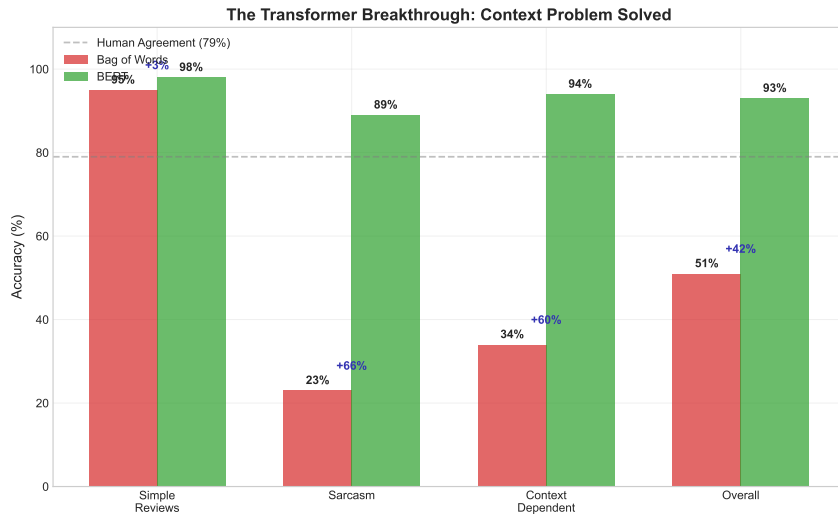
After Fine-tuning:

"This app is sick!" → Positive (cool)

Fine-tuning: Teaching BERT Your Specific Task

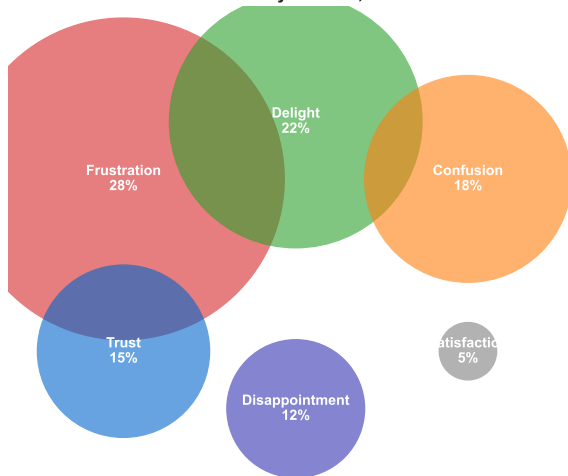


Transformers Solve the Context Problem



What BERT Found in 50,000 Reviews

Emotion Taxonomy from 50,000 Reviews



6 Core Emotion Clusters:

1. Frustration (28%)

- Long wait times
- Complex interfaces
- Missing features

2. Delight (22%)

- Unexpected features
- Beautiful design
- Fast performance

3. Confusion (18%)

- Unclear instructions
- Hidden functions
- Inconsistent behavior

4. Trust (15%)

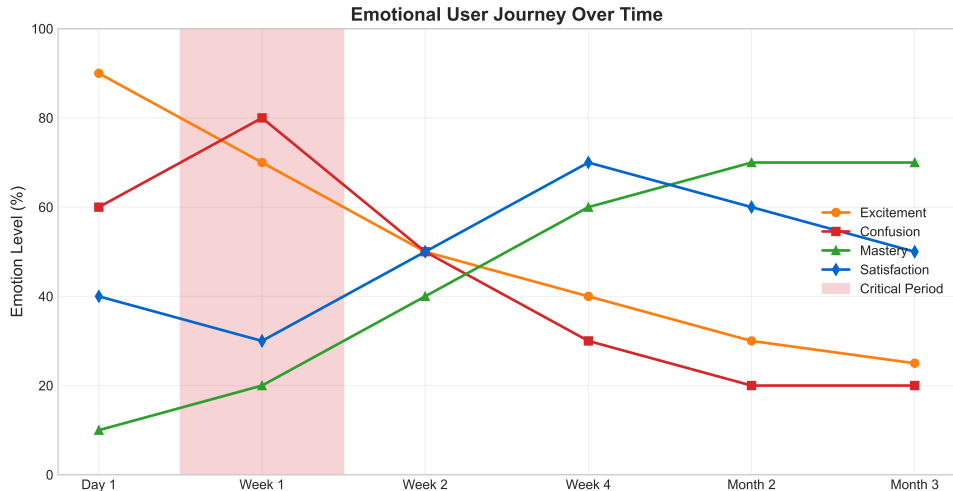
- Data security
- Reliable service
- Transparent pricing

5. Disappointment (12%)

- Unmet expectations
- Quality issues
- Broken promises

6. Satisfaction (5%)

When and Why Emotions Shift

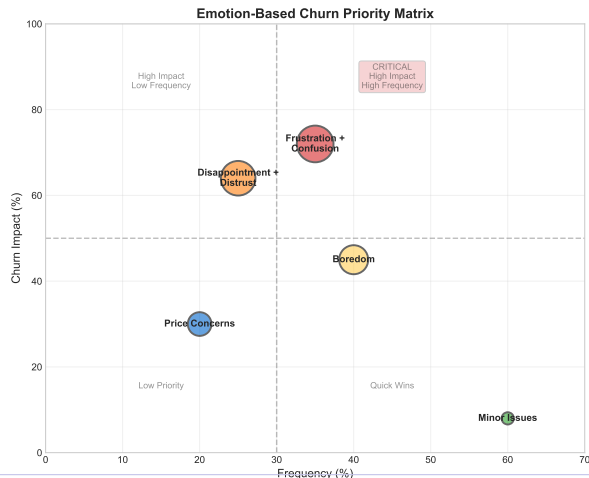


Onboarding (Day 1-7):

Learning (Week 2-4):

Routine (Month 2+):

Which Emotions Predict User Loss?



Emotion combinations matter more than individual emotions - BERT captures these interactions

High Impact on Churn:

1. Frustration + Confusion (72% quit)
"Can't figure it out and support doesn't help"
Design Action: Better onboarding + in-app help

2. Disappointment + Distrust (64% quit)
"Not what was promised, feels sketchy"
Design Action: Align marketing with reality

Low Impact on Churn:

3. Minor Frustrations (8% quit)
"Annoying but I deal with it"
Design Action: Fix in regular updates

Real Users, Not Imagined Ones

The Enthusiast

15% of users

Language:

- "Love it!"
- "Game changer"
- "Can't wait for..."

Emotions:

- Delight: 78%
- Anticipation: 22%

Design Need:

Advanced features

The Struggler

35% of users

Language:

- "Trying to..."
- "Can't find..."
- "How do I..."

Emotions:

- Confusion: 61%
- Frustration: 39%

Design Need:

Better guidance

The Pragmatist

30% of users

Language:

- "It works"
- "Does the job"
- "Fair price"

Emotions:

- Satisfaction: 82%
- Neutral: 18%

Design Need:

Reliability

The Critic

20% of users

Language:

- "Should have..."
- "Compared to X..."
- "Missing..."

Emotions:

- Disappointment: 54%
- Frustration: 46%

Design Need:

Feature parity

These personas emerged from BERT clustering - not designer assumptions

BERT groups users by emotional language patterns, revealing authentic user segments

Personalized Understanding, Automated Delivery

The Scale Challenge:

- 1 million users
- 100 reviews each
- 100 million opinions
- Impossible to read manually

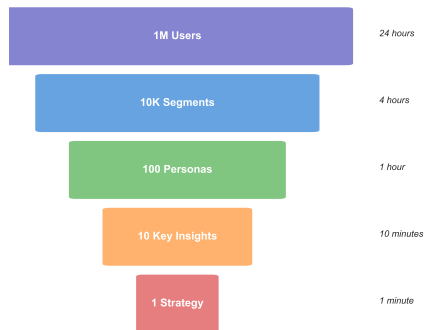
BERT's Solution:

- Process all in 24 hours
- Understand each individually
- Group by emotional need
- Generate targeted responses

Personalization Examples:

- Frustrated user → Proactive support offer
- Confused user → Tutorial recommendation
- Delighted user → Feature beta invitation
- Disappointed user → Feedback survey

Empathy at Scale: From Millions to One



Real Impact:

Emotion-Driven Engagement at Scale

The Challenge:

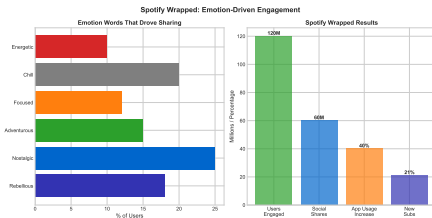
- 400 million users
- Make each feel special
- Drive social sharing
- Increase engagement

NLP Analysis Revealed:

- Users want validation of taste
- Nostalgia drives sharing
- Uniqueness matters most
- Discovery excites users

Design Response:

- Personal emotion words: “Your year was Rebellious”
- Unique statistics: “Top 0.5% of fans”
- Nostalgic moments: “You played X 47 times in March”
- Social proof: “Share your unique taste”



Results:

- 120M users engaged
- 60M social shares
- 40% increase in app usage
- 21% increase in subscriptions

ROI: 400% on NLP investment

Emotion understanding → Personalization → Engagement → Business value

Apply These Techniques to Real Data

Workshop Exercise (45 minutes):

Dataset:

- 5,000 app store reviews
- Your choice of app category
- Mix of ratings (1-5 stars)
- Real user language

Your Tasks:

1. Load pre-trained BERT model
2. Fine-tune on 500 labeled reviews
3. Analyze remaining 4,500 reviews
4. Discover emotion clusters
5. Identify top 3 pain points
6. Generate design recommendations

Tools Provided:

- Jupyter notebook template
- Pre-processed data

Expected Outputs:

- 1. Emotion Distribution Chart**
Show percentages of each emotion
- 2. Pain Point Priority List**
Ranked by impact on ratings
- 3. User Segment Personas**
Data-driven, not assumed
- 4. Design Recommendations**
Specific, actionable, prioritized

Learning Objectives:

- Use BERT for emotion analysis
- Interpret attention weights
- Convert ML insights to design actions
- Experience the power of scale

From 5,000 reviews to 5 key insights in 45 minutes!

Three Levels of Hands-on Learning

Exercise 1: Basic Sentiment Analysis

Time: 20 minutes

Difficulty: Beginner

Task:

- Use pre-trained model
- Analyze 100 reviews
- Classify positive/negative
- Calculate accuracy

Learning Goal:

Understand basic NLP pipeline

Tools:

- TextBlob or
- Hugging Face pipeline

Exercise 2: Intermediate Emotion Detection

Time: 45 minutes

Difficulty: Medium

Task:

- Fine-tune BERT
- 6-class emotions
- Analyze attention weights
- Visualize results

Learning Goal:

Work with transformers

Tools:

- Transformers library
- Pre-labeled dataset

Exercise 3: Advanced Full Pipeline

Time: 90 minutes

Difficulty: Challenging

Task:

- Scrape real reviews
- Preprocess text
- Fine-tune model
- Generate personas
- Create dashboard

Learning Goal:

End-to-end implementation

Deliverable:

Emotion insights report

Resources: Notebooks at github.com/ml-design-course/week3-nlp

Each exercise builds on the previous - complete all three for mastery

From Words to Design Insights

Technical Understanding:

- Why context matters in language
- How Bag of Words loses information
- What attention mechanisms do
- How BERT reads bidirectionally
- Why pre-training + fine-tuning works

Practical Skills:

- Identify emotion in text at scale
- Use pre-trained models effectively
- Fine-tune for specific domains
- Interpret attention visualizations
- Convert ML outputs to insights

Design Applications:

- Data-driven persona creation
- Emotion-based user journeys
- Priority matrices from ML analysis
- Scalable empathy systems
- Personalization strategies

Remember:

The Goal: Not to read faster, but to understand deeper
The Method: Selective attention to what matters
The Result: True user understanding at scale

Next Week: Classification & Problem Definition

You now have the tools to understand what millions of users really feel

Key Takeaway

Understanding emotion at scale is not about reading faster—it's about attending smarter

Next Week: Classification & Problem Definition
From emotions to actionable categories