



# Evaluation and Perplexity

---

Natural Language Processing

Dr. Imran Ihsan  
[www.imranihsan.com](http://www.imranihsan.com)

# Evaluation: How good is our model?

Does our language model prefer good sentences to bad ones?

Assign higher probability to “real” or “frequently observed” sentences

Than “ungrammatical” or “rarely observed” sentences?

We train parameters of our model on a **training set**.

We test the model’s performance on data we haven’t seen.

A **test set** is an unseen dataset that is different from our training set, totally unused.

An **evaluation metric** tells us how well our model does on the test set.

# Extrinsic evaluation of N-gram models

## Best evaluation for comparing models A and B

Put each model in a task

spelling corrector, speech recognizer, MT system

Run the task, get an accuracy for A and for B

How many misspelled words corrected properly

How many words translated correctly

Compare accuracy for A and B

# Difficulty of extrinsic (in-vivo) evaluation of N-gram models

## Extrinsic evaluation

Time-consuming; can take days or weeks

So

Sometimes use **intrinsic** evaluation: **perplexity**

Bad approximation

unless the test data looks **just** like the training data

So **generally, only useful in pilot experiments**

But is helpful to think about.

# Intuition of Perplexity

## The Shannon Game:

How well can we predict the next word?

I always order pizza with cheese and \_\_\_\_\_

The 33<sup>rd</sup> President of the US was \_\_\_\_\_

I saw a \_\_\_\_\_

Unigrams are terrible at this game. (Why?)

mushrooms 0.1

pepperoni 0.1

anchovies 0.01

....

fried rice 0.0001

....

and 1e-100

A better model of a text

is one which assigns a higher probability to the word that actually occurs

# Perplexity

The best language model is one that best predicts an unseen test set

- Gives the highest  $P(\text{sentence})$

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

Perplexity is the inverse probability of the test set, normalized by the number of words:

$$= \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}}$$

Chain rule:

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

For bigrams:

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

Minimizing perplexity is the same as maximizing probability

# The Shannon Game intuition for perplexity

From Josh Goodman

Perplexity is weighted equivalent branching factor

How hard is the task of recognizing digits '0,1,2,3,4,5,6,7,8,9'

Perplexity 10

How hard is recognizing (30,000) names at Microsoft.

Perplexity = 30,000

Let's imagine a call-routing phone system gets 120K calls and has to recognize

"Operator" (let's say this occurs 1 in 4 calls)

"Sales" (1 in 4)

"Technical Support" (1 in 4)

30,000 different names (each name occurring 1 time in the 120K calls)

What is the perplexity? Next slide

# The Shannon Game intuition for perplexity

Josh Goodman: imagine a call-routing phone system gets 120K calls and must recognize

"Operator" (let's say this occurs 1 in 4 calls)

"Sales" (1 in 4)

"Technical Support" (1 in 4)

30,000 different names (each name occurring 1 time in the 120K calls)

We get the perplexity of this sequence of length 120K by first multiplying 120K probabilities (90K of which are 1/4 and 30K of which are 1/120K), and then taking the inverse 120,000th root:

$$\text{Perp} = (1/4 * 1/4 * 1/4 * 1/4 * 1/4 * \dots * 1/120K * 1/120K * \dots)^{(-1/120K)}$$

But this can be arithmetically simplified to just  $N = 4$ : the operator (1/4), the sales (1/4), the tech support (1/4), and the 30,000 names (1/120,000):

$$\text{Perplexity} = ((1/4 * 1/4 * 1/4 * 1/120K)^{(-1/4)}) = 52.6$$



# Perplexity as branching factor

Let's suppose a sentence consisting of random digits

What is the perplexity of this sentence according to a model that assign  $P=1/10$  to each digit?

$$\begin{aligned} PP(W) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \left(\frac{1}{10}\right)^{-\frac{1}{N}} \\ &= \frac{1}{10}^{-1} \\ &= 10 \end{aligned}$$

# Lower perplexity = better model

Training 38 million words, test 1.5 million words, WSJ

N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109

# The Shannon Visualization Method

Choose a random bigram  
( $\langle s \rangle$ ,  $w$ ) according to its probability

$\langle s \rangle$  I

I want

want to

Now choose a random bigram  
( $w$ ,  $x$ ) according to its probability

to eat

eat Chinese

Chinese food

And so on until we choose  $\langle /s \rangle$

food  $\langle /s \rangle$

Then string the words together

I want to eat Chinese food

# Approximating Shakespeare

## Unigram

To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have  
Every enter now severally so, let  
Hill he late speaks; or! a more to leg less first you enter  
Are where exeunt and sighs have rise excellency took of.. Sleep knave we. near; vile like

## Bigram

What means, sir. I confess she? then all sorts, he is trim, captain.  
Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.  
What we, hath got so she that I rest and sent to scold and nature bankrupt, nor the first gentleman?

## Trigram

Sweet prince, Falstaff shall die. Harry of Monmouth's grave.  
This shall forbid it should be branded, if renown made it empty.  
Indeed the duke; and had a very good friend.  
Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.

## Quadrigram

King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;  
Will you not tell me who I am?  
It cannot be but so.  
Indeed the short and the long. Marry, 'tis a noble Lepidus.

# Shakespeare as corpus

$N=884,647$  tokens,  $V=29,066$

Shakespeare produced 300,000 bigram types out of  $V^2= 844$  million possible bigrams.

So 99.96% of the possible bigrams were never seen (have zero entries in the table)

Quadrigrams worse: What's coming out looks like Shakespeare because it **is** Shakespeare

# The Wall Street Journal is not Shakespeare (no offense)

## **Unigram**

Months the my and issue of year foreign new exchange's september were recession ex-  
change new endorsed a acquire to six executives

## **Bigram**

Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor  
would seem to complete the major central planners one point five percent of U. S. E. has  
already old M. X. corporation of living on information such as more frequently fishing to  
keep her

## **Trigram**

They also point to ninety nine point six billion dollars from two hundred four oh six three  
percent of the rates of interest stores as Mexico and Brazil on market conditions

Can you guess the training set author of the LM that generated these random 3-gram sentences?

They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and gram Brazil on market conditions

This shall forbid it should be branded, if renown made it empty.

“You are uniformly charming!” cried he, with a smile of associating and now and then I bowed and they perceived a chaise and four to wish for.

# The perils of overfitting

N-grams only work well for word prediction if the test corpus looks like the training corpus

In real life, it often doesn't

We need to train robust models that generalize!

One kind of generalization: Zeros!

Things that don't ever occur in the training set

But occur in the test set



# Zeros

Training set:

- ... denied the allegations
- ... denied the reports
- ... denied the claims
- ... denied the request

$P(\text{"offer"} \mid \text{denied the}) = 0$

Test set

- ... denied the offer
- ... denied the loan

# Zero probability bigrams

Bigrams with zero probability

mean that we will assign 0 probability to the test set!

And hence we cannot compute perplexity (can't divide by 0)!