

# **Evaluating NLP Models**

Mayank Singh

ACM-IKDD Summer School on Data Science, IIT Gandhinagar, 8<sup>th</sup> July 2022

# A Brief Recap

- **NLP Tasks**
- **The Traditional NLP pipeline**
- **The Modern NLP pipeline**
- **Cleaning and Processing Text**
- **A walkthrough of Modern NLP pipeline (lab session)**

**Let's look at some popular evaluation metrics**

# Text Classification Tasks

S.No.	Text	Actual Class
1	Very user friendly...long lasting battery....very clear display	1
2	Worse graphics, won't keep a wireless connection,overall not satisfied	0
3	We have had no issues with this tablet. Love it!TY	1
4	so far. A few not so great things, modt of them managed to resolve them though.	0
5	I brought this tablet during the Black Friday sale fast forward to Christmas when my 6 year opens it he was happy I set it up for him two days later it crashed never came back on. Blank screen. He was disappointed and so was I .	0
6	Great tablet/e-reader. You can do a lot on this bad boy!!	1
7	This is an older tablet. I probably expected too much. If I'm going to deal with lagginess I prefer my ipad2.	0
8	Really expected a chord with it for the price we paid. Very dissapointed.	0
9	I bought 3, better than we ever dreamed for money spent	1
10	We bought this for my mom. She is very satisfied with this product.	1

**Ten reviews from Amazon with sentiment class labels**

**1: Good Product/Service**

**0: Bad Product/Service**

# Text Classification Tasks

S.No.	Text	Actual Class	Predicted Class
1	Very user friendly...long lasting battery....very clear display	1	1
2	Worse graphics, won't keep a wireless connection,overall not satisfied	0	1
3	We have had no issues with this tablet. Love it!TY	1	0
4	so far. A few not so great things, modt of them managed to resolve them though.	0	0
5	I brought this tablet during the Black Friday sale fast forward to Christmas when my 6 year opens it he was happy I set it up for him two days later it crashed never came back on. Blank screen. He was disappointed and so was I .	0	1
6	Great tablet/e-reader. You can do a lot on this bad boy!!	1	0
7	This is an older tablet. I probably expected too much. If I'm going to deal with lagginess I prefer my ipad2.	0	0
8	Really expected a chord with it for the price we paid. Very dissapointed.	0	1
9	I bought 3, better than we ever dreamed for money spent	1	0
10	We bought this for my mom. She is very satisfied with this product.	1	1

**We passed these reviews through a trained sentiment classifier and got the predictions**

# Text Classification Tasks

S.No.	Text	Actual Class	Predicted Class	Category
1	Very user friendly...long lasting battery....very clear display	1	1	TP
2	Worse graphics, won't keep a wireless connection,overall not satisfied	0	1	FP
3	We have had no issues with this tablet. Love it!TY	1	0	FN
4	so far. A few not so great things, modt of them managed to resolve them though.	0	0	TN
5	I brought this tablet during the Black Friday sale fast forward to Christmas when my 6 year opens it he was happy I set it up for him two days later it crashed never came back on. Blank screen. He was disappointed and so was I .	0	1	FP
6	Great tablet/e-reader. You can do a lot on this bad boy!!	1	0	FN
7	This is an older tablet. I probably expected too much. If I'm going to deal with lagginess I prefer my ipad2.	0	0	TN
8	Really expected a chord with it for the price we paid. Very dissapointed.	0	1	FP
9	I bought 3, better than we ever dreamed for money spent	1	0	FN
10	We bought this for my mom. She is very satisfied with this product.	1	1	TP

**True Positive:** Correct prediction for class 1

**True Negative:** Correct prediction for class 0

**False Positive:** Incorrect prediction for class 0

**False Negative:** Incorrect prediction for class 1

# Text Classification Tasks

S.No.	Text	Actual Class	Predicted Class	Category	Metrics
1	Very user friendly...long lasting battery....very clear display	1	1	TP	
2	Worse graphics, won't keep a wireless connection,overall not satisfied	0	1	FP	
3	We have had no issues with this tablet. Love it!TY	1	0	FN	<b>Precision</b> = $TP/(TP+FP) = 2/(2+3) = 2/5$
4	so far. A few not so great things, modt of them managed to resolve them though.	0	0	TN	<b>Recall</b> = $TP/(TP+FN) = 2/(2+3) = 2/5$
5	I brought this tablet during the Black Friday sale fast forward to Christmas when my 6 year opens it he was happy I set it up for him two days later it crashed never came back on. Blank screen. He was disappointed and so was I .	0	1	FP	<b>F-Score</b> = $2*P*R/(P+R) = 2/5$
6	Great tablet/e-reader. You can do a lot on this bad boy!!	1	0	FN	<b>Accuracy</b> = $(TP + TN)/(\text{no. of examples}) = (2+2)/10 = 2/5$
7	This is an older tablet. I probably expected too much. If I'm going to deal with lagginess I prefer my ipad2.	0	0	TN	
8	Really expected a chord with it for the price we paid. Very dissapointed.	0	1	FP	
9	I bought 3, better than we ever dreamed for money spent	1	0	FN	
10	We bought this for my mom. She is very satisfied with this product.	1	1	TP	

**True Positive:** Correct prediction for class 1

**True Negative:** Correct prediction for class 0

**False Positive:** Incorrect prediction for class 0

**False Negative:** Incorrect prediction for class 1

# Summarization Tasks

- **Extrinsic:** the summary quality is judged on the basis of how helpful summaries are for a given task.
- **Intrinsic:** based on analysis of the summary itself
  - It involves a comparison with the source document
  - How many main ideas of the source document are covered by the summary in comparison with an abstract written by a human.



# Summarization Tasks

## ROUGE (Recall Oriented Understudy for Gisting Evaluation)

Given a document  $D$ , and an automatic summary  $X$ :

- Have  $N$  humans produce a set of reference summaries of  $D$
- What percentage of the bigrams from the reference summaries appear in  $X$ ?

$$ROUGE - 2 = \frac{\sum_{S \in \{RefSummaries\}} \sum_{bigrams \in S} count_{match}(bigrams)}{\sum_{S \in \{RefSummaries\}} \sum_{bigrams \in S} count(bigrams)}$$

# Summarization Tasks

**ROUGE (Recall Oriented Understudy for Gisting Evaluation)**

**Automatic Summary:** the cat was found under the bed

**Reference Summary:** the cat was under the bed

# Summarization Tasks

## ROUGE (Recall Oriented Understudy for Gisting Evaluation)

**Automatic Summary:** the cat was found under the bed

**Reference Summary:** the cat was under the bed

**Generated Summary Bigrams:** the cat, cat was, was found, found under, under the, the bed

**Reference Summary Bigrams:** the cat, cat was, was under, under the, the bed

$$\text{ROUGE-2} = \frac{4}{5}$$

# Summarization Tasks

## ROUGE (Recall Oriented Understudy for Gisting Evaluation)

**Automatic Summary:** the cat was found under the bed

**Reference Summary:** the cat was under the bed

**Generated Summary Bigrams:** the cat, cat was, was found, found under, under the, the bed

**Reference Summary Bigrams:** the cat, cat was, was under, under the, the bed

$$\text{ROUGE-2} = \frac{4}{5}$$

**ROUGE-1, ROUGE-L, ROUGE-S, ROUGE-SU??**

# Machine Translation Tasks

## BLEU (Recall Oriented Understudy for Gisting Evaluation)

Given a machine generated text D (**candidates**), and an reference translations X:

- Have N humans produce a set of reference translations X
- What percentage of the n-grams in D are present in reference translations X?

$$p_n = \frac{\sum_{C \in \{Candidates\}} \sum_{n\text{-gram} \in C} Count_{clip}(n\text{-gram})}{\sum_{C' \in \{Candidates\}} \sum_{n\text{-gram}' \in C'} Count(n\text{-gram}')}.$$

# Machine Translation Tasks

**BLEU:** Precision based metric

**ROUGE:** Recall-based metric

# Text Generation

## Perplexity Metric

**Perplexity is the inverse probability of the test data, normalized by the number of words**

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

# Text Generation

## Perplexity Metric

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

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

How do you compute the probability of a sentence?



# The Probability Of A Sentence

- A sentence is a sequence of tokens

$$W = w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n$$

- $P(W) = P(w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n)$

- $P(w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_2 w_1) \dots P(w_n|w_{n-1} \dots w_1)$

# The Probability Of A Sentence

- A sentence is a sequence of tokens

$$W = w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n$$

- $P(W) = P(w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n)$
- $P(w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_2 w_1) \dots P(w_n|w_{n-1} \dots w_1)$
- **A bigram assumption:** A word only depends on its previous word
- $P(w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_2) \dots P(w_n|w_{n-1})$

# The Probability Of A Sentence

- A sentence is a sequence of tokens

$$W = w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n$$

- $P(W) = P(w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n)$
- $P(w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_2 w_1) \dots P(w_n|w_{n-1} \dots w_1)$
- A bigram assumption: A word only depends on its previous word
- $P(w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots w_n) = P(w_1)P(w_2|w_1)P(w_3|w_2) \dots P(w_n|w_{n-1})$

**How do you compute the probability of a token?**

# The Maximum Likelihood Estimate

## Bigram Probability

$$P(w_i|w_{i-1}) = \frac{\textit{count}(w_{i-1}, w_i)}{\textit{count}(w_{i-1})}$$

# Revisiting Perplexity

## Perplexity Metric

**Perplexity is the inverse probability of the test data, normalized by the number of words**

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

## Perplexity Metric with Bigram Assumption




**Perplexity is the inverse probability of the test data, normalized by the number of words**

$$PP(W) = \left( \prod \frac{1}{P(w_i | w_{i-1})} \right)^{\frac{1}{N}}$$

# Behavioral Testing of NLP Models




# Behavioral Testing of NLP Models

## Sentiment Analysis

Test <i>TYPE</i> and Description		Failure Rate (%)					Example test cases & expected behavior
			<b>G</b>			RoB	
Vocab.+POS	<b>MFT</b> : Short sentences with neutral adjectives and nouns	0.0	7.6	4.8	94.6	81.8	The company is Australian. <b>neutral</b> That is a private aircraft. <b>neutral</b>
	<b>MFT</b> : Short sentences with sentiment-laden adjectives	4.0	15.0	2.8	0.0	0.2	That cabin crew is extraordinary. <b>pos</b> I despised that aircraft. <b>neg</b>
	<b>INV</b> : Replace neutral words with other neutral words	9.4	16.2	12.4	10.2	10.2	@Virgin should I be concerned <b>that</b> → <b>when</b> I'm about to fly ... <b>INV</b> @united <b>the</b> → <b>our</b> nightmare continues... <b>INV</b>
	<b>DIR</b> : Add positive phrases, fails if sent. goes down by > 0.1	12.6	12.4	1.4	0.2	10.2	@SouthwestAir Great trip on 2672 yesterday... <b>You are extraordinary.</b> ↑ @AmericanAir AA45 ... JFK to LAS. <b>You are brilliant.</b> ↑
	<b>DIR</b> : Add negative phrases, fails if sent. goes up by > 0.1	0.8	34.6	5.0	0.0	13.2	@USAirways your service sucks. <b>You are lame.</b> ↓ @JetBlue all day. <b>I abhor you.</b> ↓
Robust.	<b>INV</b> : Add randomly generated URLs and handles to tweets	9.6	13.4	24.8	11.4	7.4	@JetBlue that selfie was extreme. <b>@pi9QDK</b> <b>INV</b> @united stuck because staff took a break? Not happy 1K.... <b>https://t.co/PWK1jb</b> <b>INV</b>
	<b>INV</b> : Swap one character with its neighbor (typo)	5.6	10.2	10.4	5.2	3.8	<b>@JetBlue</b> → <b>@JeBtlue</b> I cri <b>INV</b> @SouthwestAir no <b>thanks</b> → <b>thakns</b> <b>INV</b>
NER	<b>INV</b> : Switching locations should not change predictions	7.0	20.8	14.8	7.6	6.4	@JetBlue I want you guys to be the first to fly to # <b>Cuba</b> → <b>Canada</b> ... <b>INV</b> @VirginAmerica I miss the #nerdbird in <b>San Jose</b> → <b>Denver</b> <b>INV</b>
	<b>INV</b> : Switching person names should not change predictions	2.4	15.1	9.1	6.6	2.4	...Airport agents were horrendous. <b>Sharon</b> → <b>Erin</b> was your saviour <b>INV</b> @united 8602947, <b>Jon</b> → <b>Sean</b> at http://t.co/58tuTgli0D, thanks. <b>INV</b>

# Behavioral Testing of NLP Models


## Sentiment Analysis

Test <i>TYPE</i> and Description		Failure Rate (%)					Example test cases & expected behavior
			<b>G</b>			RoB	
Temporal	<b>MFT:</b> Sentiment change over time, present should prevail	41.0	36.6	42.2	18.8	11.0	I used to hate this airline, although now I like it. <b>pos</b> In the past I thought this airline was perfect, now I think it is creepy. <b>neg</b>
	<b>MFT:</b> Negated negative should be positive or neutral	18.8	54.2	29.4	13.2	2.6	The food is not poor. <b>pos or neutral</b> It isn't a lousy customer service. <b>pos or neutral</b>
Negation	<b>MFT:</b> Negated neutral should still be neutral	40.4	39.6	74.2	98.4	95.4	This aircraft is not private. <b>neutral</b> This is not an international flight. <b>neutral</b>
	<b>MFT:</b> Negation of negative at the end, should be pos. or neut.	100.0	90.4	100.0	84.8	7.2	I thought the plane would be awful, but it wasn't. <b>pos or neutral</b> I thought I would dislike that plane, but I didn't. <b>pos or neutral</b>
	<b>MFT:</b> Negated positive with neutral content in the middle	98.4	100.0	100.0	74.0	30.2	I wouldn't say, given it's a Tuesday, that this pilot was great. <b>neg</b> I don't think, given my history with airplanes, that this is an amazing staff. <b>neg</b>
	<b>MFT:</b> Author sentiment is more important than of others	45.4	62.4	68.0	38.8	30.0	Some people think you are excellent, but I think you are nasty. <b>neg</b> Some people hate you, but I think you are exceptional. <b>pos</b>
SRL	<b>MFT:</b> Parsing sentiment in (question, "yes") form	9.0	57.6	20.8	3.6	3.0	Do I think that airline was exceptional? Yes. <b>neg</b> Do I think that is an awkward customer service? Yes. <b>neg</b>
	<b>MFT:</b> Parsing sentiment in (question, "no") form	96.8	90.8	81.6	55.4	54.8	Do I think the pilot was fantastic? No. <b>neg</b> Do I think this company is bad? No. <b>pos or neutral</b>




# Behavioral Testing of NLP Models

## Quora Question Pair

Test <i>TYPE</i> and Description		Failure Rate		Example Test cases & expected behavior
		 RoB		
Vocab.	<b>MFT</b> : Modifiers changes question intent	78.4	78.0	{ Is Mark Wright a photographer?   Is Mark Wright an accredited photographer? } $\neq$
Taxonomy	<b>MFT</b> : Synonyms in simple templates	22.8	39.2	{ How can I become more vocal?   How can I become more outspoken? } =
	<b>INV</b> : Replace words with synonyms in real pairs	13.1	12.7	Is it necessary to follow a religion? Is it necessary to follow an <b>organized</b> $\rightarrow$ <b>organised</b> religion? } <b>INV</b>
	<b>MFT</b> : More X = Less antonym(X)	69.4	100.0	{ How can I become more optimistic?   How can I become less pessimistic? } =
Robust.	<b>INV</b> : Swap one character with its neighbor (typo)	18.2	12.0	{ Why am I <b>getting</b> $\rightarrow$ <b>gettnig</b> lazy?   Why are we so lazy? } <b>INV</b>
	<b>DIR</b> : Paraphrase of question should be duplicate	69.0	25.0	Can I gain weight from not eating enough? <b>Can I</b> $\rightarrow$ <b>Do you think I can</b> gain weight from not eating enough? } =
NER	<b>INV</b> : Change the same name in both questions	11.8	9.4	Why isn't <b>Hillary Clinton</b> $\rightarrow$ <b>Nicole Perez</b> in jail? Is <b>Hillary Clinton</b> $\rightarrow$ <b>Nicole Perez</b> going to go to jail? } <b>INV</b>
	<b>DIR</b> : Change names in one question, expect $\neq$	35.1	30.1	What does India think of Donald Trump? What India thinks about <b>Donald Trump</b> $\rightarrow$ <b>John Green</b> ? } $\neq$
	<b>DIR</b> : Keep first word and entities of a question, fill in the gaps with RoBERTa; expect $\neq$	30.0	32.8	Will it be difficult to get a US Visa if Donald Trump gets elected? Will the US accept Donald Trump? } $\neq$

# Behavioral Testing of NLP Models

## Quora Question Pair

	Test <i>TYPE</i> and Description	Failure Rate		Example Test cases & expected behavior
		 RoB		
Temporal	<b>MFT:</b> Is $\neq$ used to be, non-duplicate	61.8	96.8	{ Is Jordan Perry an advisor?   Did Jordan Perry use to be an advisor? } $\neq$
	<b>MFT:</b> before $\neq$ after, non-duplicate	98.0	34.4	{ Is it unhealthy to eat after 10pm?   Is it unhealthy to eat before 10pm? } $\neq$
	<b>MFT:</b> before becoming $\neq$ after becoming	100.0	0.0	What was Danielle Bennett's life before becoming an agent? } $\neq$ What was Danielle Bennett's life after becoming an agent? }
Negation	<b>MFT:</b> simple negation, non-duplicate	18.6	0.0	{ How can I become a person who is not biased?   How can I become a biased person? } $\neq$
	<b>MFT:</b> negation of antonym, should be duplicate	81.6	88.6	{ How can I become a positive person?   How can I become a person who is not negative } $\neq$
Coref	<b>MFT:</b> Simple coreference: he $\neq$ she	79.0	96.6	If Joshua and Chloe were alone, do you think he would reject her? } $\neq$ If Joshua and Chloe were alone, do you think she would reject him? }
	<b>MFT:</b> Simple resolved coreference, his and her	99.6	100.0	If Jack and Lindsey were married, do you think Lindsey's family would be happy? } $\neq$ If Jack and Lindsey were married, do you think his family would be happy? }
SRL	<b>MFT:</b> Order is irrelevant for comparisons	99.6	100.0	{ Are tigers heavier than insects?   What is heavier, insects or tigers? } =
	<b>MFT:</b> Orders is irrelevant in symmetric relations	81.8	100.0	{ Is Nicole related to Heather?   Is Heather related to Nicole? } =
	<b>MFT:</b> Order is relevant for asymmetric relations	71.4	100.0	{ Is Sean hurting Ethan?   Is Ethan hurting Sean? } $\neq$
	<b>MFT:</b> Active / passive swap, same semantics	65.8	98.6	{ Does Anna love Benjamin?   Is Benjamin loved by Anna? } =
	<b>MFT:</b> Active / passive swap, different semantics	97.4	100.0	{ Does Danielle support Alyssa?   Is Danielle supported by Alyssa? } $\neq$

# Behavioral Testing of NLP Models

<https://github.com/marcotcr/checklist>

```
In [27]: ▶ editor.visual_suggest('This is {a:mask} movie.')
```



This is **a:mask** movie .

FILL IN WITH...

- ☐ Check All
- ☐ a good
- ☐ an amazing
- ☐ an excellent
- ☐ an awful



Preview



No Data

```
In [26]: ▶ editor.selected_suggestions
```

Wordnet



**ANY  
QUESTIONS?**

**Email:** [singh.mayank@iitgn.ac.in](mailto:singh.mayank@iitgn.ac.in)

**Webpage:** <https://mayank4490.github.io/>