

#### Cuprins

1 SETUL DE DATE

O2 PREPROCESAREA DATELOR

03 EXTRAGEREA CARACTERISTICILOR

**04** MODEL





#### SETUL DE DATE

În urma efectuării unei analize a mai multor seturi de date, am optat pentru setul de date Fig-qa descărcat de pe Hugging Face. Pentru fiecare metaforă din dataset, se mai găsesc 2 posibilități de interpretare a acesteia, alături de label-ul care corespunde traducerii fiecărei metafore ( 0 pentru ending1, 1 pentru ending2).

	Α	В	С	D	Е	F	G	н		J	K	L	М	N	0
1	startphras	e,ending1,	ending2,la	abels, valid											
2	The future	e is as brigh	nt as the su	ın,The futu	re is bright	t,The futu	e is not br	ight,0,1							
3	The future	is as brigh	nt as ink,Th	ne future is	bright,The	future is	not bright,	1,1							
4	She think	of herself	as the mo	ther of hur	nanity.,Sh	e is self-ri	ghteous ar	d entitled	She is hur	mble and r	espectful	towards th	e different	t people.,0,	1
5	She think	of herself	as a partic	le of sand i	in the dese	ert.,She is	self-righte	ous and er	titled.,She	e is humble	e and resp	ectful towa	ards the di	fferent peo	ple.,1,1
6	His sacrifi	ce showed	he had a h	eart like a	lion,His sa	crifice was	very bold	His sacrifi,	ce was tim	id and wea	ak,0,1				
	His sacrifi	ce showed	he had a h	eart like a	mouse,His	sacrifice	was very b	old,His sac	rifice was t	imid and v	veak,1,1				
8	The gover	nment ran	with the e	efficiency o	f a well-oi	led machii	ne,The gov	ernment v	as very ef	ficient.,Th	e governm	ent was no	ot efficient	t at all.,0,1	
9	The government ran with the efficiency of a sputtering engine, The government was very efficient., The government was not efficient at all., 1,1														
10	She was a	s flaky as w	ood shavii	ngs,She wa	s extreme	ly flaky.,S	he wasn't f	laky at all.	0,1						
11	She was a	s flaky as s	teel,She w	as extreme	ely flaky.,S	he wasn't	flaky at all	.,1,1							
12	The book is an complex as the night sky,The book is very complicated.,The book is easy to understand.,0,1														
13	The book	is an comp	lex as a chi	ildren's puz	zle,The bo	ook is very	complicat	ed.,The bo	ok is easy t	to underst	and.,1,1				
14	The fog h	s the thick	ness of a h	neavy wom	an's thighs	The fog is	thick,The	fog is thin	0,1						
15	The fog h	s the thick	ness of Ch	ristian Bale	in The Ma	achinist,Th	e fog is th	ick,The fog	is thin,1,1						
16	The preschool teacher was a grim reaper.,The preschool teacher was dour.,The preschool teacher was funny.,0,1														
17	The presc	hool teach	er was a clo	own.,The p	reschool t	eacher wa	s dour.,Th	e preschoo	l teacher w	vas funny.,	1,1				
18	Hollywoo	d is as genu	uine as Silk	,Hollywoo	d is the rea	al deal.,Ho	llywood is	fake.,0,1							
19	Hollywoo	d is as genu	uine as Ray	on,Hollyw	ood is the	real deal.,	Hollywood	l is fake.,1,	1						
20	He ate lik	e a horse,H	e ate a lot,	He barely	ate at all,0	,1									



# PREPROCESAREA DATELOR

Înainte de preprocesarea propriu-zisă, am aplicat oprații precum





REDUNUMIREA COLOANELOR

TRATAREA
ETICHETELOR
LIPSA

pentru a asigura coerența în structura dataset-ului, cât și între dataset-uri.

## PREPROCESAREA DATELOR

Principalul pas de prelucrare a datelor pe care l-am folisit implică tokenizarea folosind tokenizerul BERT. Funcția tokenizează datele de intrare în *input\_ids* și *attention\_mask*, care sunt intrările necesare pentru modelele BERT.

```
def preprocess_function(examples):
    first_sentences = [[context] * 2 for context in examples["startphrase"]]
    question_headers = [examples['ending1'], examples['ending2']]
    choices = list(map(list, zip(*question_headers)))

first_sentences = sum(first_sentences, [])
    choices = sum(choices, [])

tokenized_examples = tokenizer(
        first_sentences,
        choices,
        truncation=True,
        padding="max_length",
        max_length=128
)

return {
        'input_ids': [tokenized_examples['input_ids'][i:i + 2] for i in range(0, len(tokenized_examples['attention_mask'][i:i + 2] for i in r
```





**DATALOADER** 

```
def load_figga_dataset(train_path, dev_path, test_path):
    train_df = pd.read_csv(train_path)
   dev_df = pd.read_csv(dev_path)
    test_df = pd.read_csv(test_path)
   train_df.rename(columns={'labels': 'label'}, inplace=True)
   dev_df.rename(columns={'labels': 'label'}, inplace=True)
    test_df.rename(columns={'labels': 'label'}, inplace=True)
   if 'label' not in test_df.columns:
        test_df['label'] = -1
   train_dataset = Dataset.from_pandas(train_df)
    dev_dataset = Dataset.from_pandas(dev_df)
    test_dataset = Dataset.from_pandas(test_df)
    return DatasetDict({
        'train': train_dataset,
        'dev': dev_dataset,
        'test': test dataset
    })
```



Am folosit acestă metodă pentru a forma obiecte propriu-zise de tip Dataset din input-ul din fișierele CSV. După ce ne-am asigurat de corectitudinea denumirii coloanelor, returnăm un dicționar care conține cele 3 dataset-uri – train, dev și test.



DEFINIREA ARGUMENTELOR

```
training_args = TrainingArguments(
   output_dir='./results',
   evaluation_strategy='epoch',
   learning_rate=2e-5,
   per_device_train_batch_size=8,
   per_device_eval_batch_size=8,
   num_train_epochs=3,
   weight_decay=0.01,
   logging_dir='./logs',
    logging_steps=10.
    load_best_model_at_end=True,
   metric_for_best_model='accuracy',
   save_strategy='epoch',
   save_total_limit=2,
   report_to=[] # Disable W&B
```



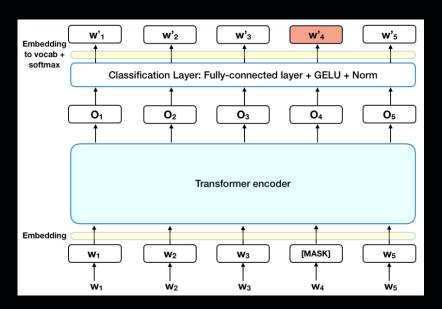
DEFINIREA ARGUMENTELOR

Aici am definit argumentele de antrenare. Pentru acest model, evaluarea se va face la finalul fiecăuria dintre cele 3 epochs. Fiecare dataframe este stocat in batch-uri de mărime 8. De asemenea, după fiecare epoch se va căuta salvarea celui mai bun model, cu o limită de 2 modele salvate.



#### **MODEL**

Am împărțit setul de date în trei dataframe-uri: train, test și validare. Pentru acest proiect, noi am folosit modelul BERT.





#### **REZULTATE**

Am antrenat modelul pentru 3 de epoci, iar pentru fiecare epocă am afișat loss-ul pentru a le putea compara.

În urma antrenării modelului, am obținut o acuratețe de 86% pe datele de test

Epoch	Training Loss	Validation Loss	Accuracy		
1	0.701000	0.652408	0.723035		
2	0.524200	0.397242	0.845521		
3	0.494300	0.316244	0.869287		
				[137/137 00:08]	
				0061569214, 'eval_accuracy': 0.8692870201096892, 'eval_runtime': 8.4474, s_per_second': 16.218, 'epoch': 3.0}	'eva

#### MULŢUM!!

#### Proiect realizat de:

Furdui Vlad-Rareş Ivan Bogdan Mihai Andrei-Alexandru