

Problem2

1. text , 'It is a', self.mask, 'sentence.'

The verbalizers are:

0: good, kind

1: toxic, terrible

Example: [his mother is a slutty whore.] It is a [terrible] sentence.

2. text, 'This comment is', self.mask, '.'

The verbalizers are:

0: positive, kind

1: toxic, aggressive

Example: [dad gum. you fuckers dont like me editing do you. dont you be dissing charlie bray now. hes just a fatass kid.] This comment is [toxic].

3. text, "The commenter is very", self.mask, 'to others.'

The verbalizers are:

0: neutral, kind

1: hatred, rude

Example: [Hey, you fascist swine, you're gonna get raped in the ass. Nice job being Jimbo's bitch, you pathetic piece of shit.] The commenter is very [hatred] to others.

4. text, 'That is too ', self.mask, '.'

The verbalizers are:

0: neutral, positive

1: terrible, toxic

Example: [Prick. Hey Ryan, Fuck you faggot, how do you like those apples bitch?]. That is too [terrible].

5. text, 'It is obviously', self.mask, 'to others.'

The verbalizers are:

0: great, kind

1: rude, offensive

Example: [BECAUSE HE IS A FAGOT.] It is obviously [offensive] to others.

6. text, 'From my point of view, this is a', self.mask, 'comment.'

The verbalizers are:

0: neutral, great

1: horrible, toxic

Example: [There was something boring here before.] From my point of view, this is a [neutral] comment.

7. text_a, 'What a', self.mask, 'comment!'

The verbalizers are:

0: good, great

1: terrible, offensive

Example: [Neither was your mum, when I had sex with her.]What a [terrible] comment!

8. 'In my opinion,', text_a,'is a',self.mask,'comment.'

The verbalizers are:

0: positive, kind

1: aggressive, rude

Example: In my opinion, [up, HES A FRAUD, A FUCKING FRAUD!!!!!!] is somewhat [aggressive].

9. text_a, 'in my understanding, is', self.mask, '.'

The verbalizers are:

0: kind,great

1: terrible, offensive

Example: [It's alright, I'm happy with the name I've got. But hey, thanks for your concern. 86.133.101.139], in my understanding, is [kind].

10. 'The sentence', text_a, "is quite ",self.mask, '.'

The verbalizers are:

0: neutral,positive

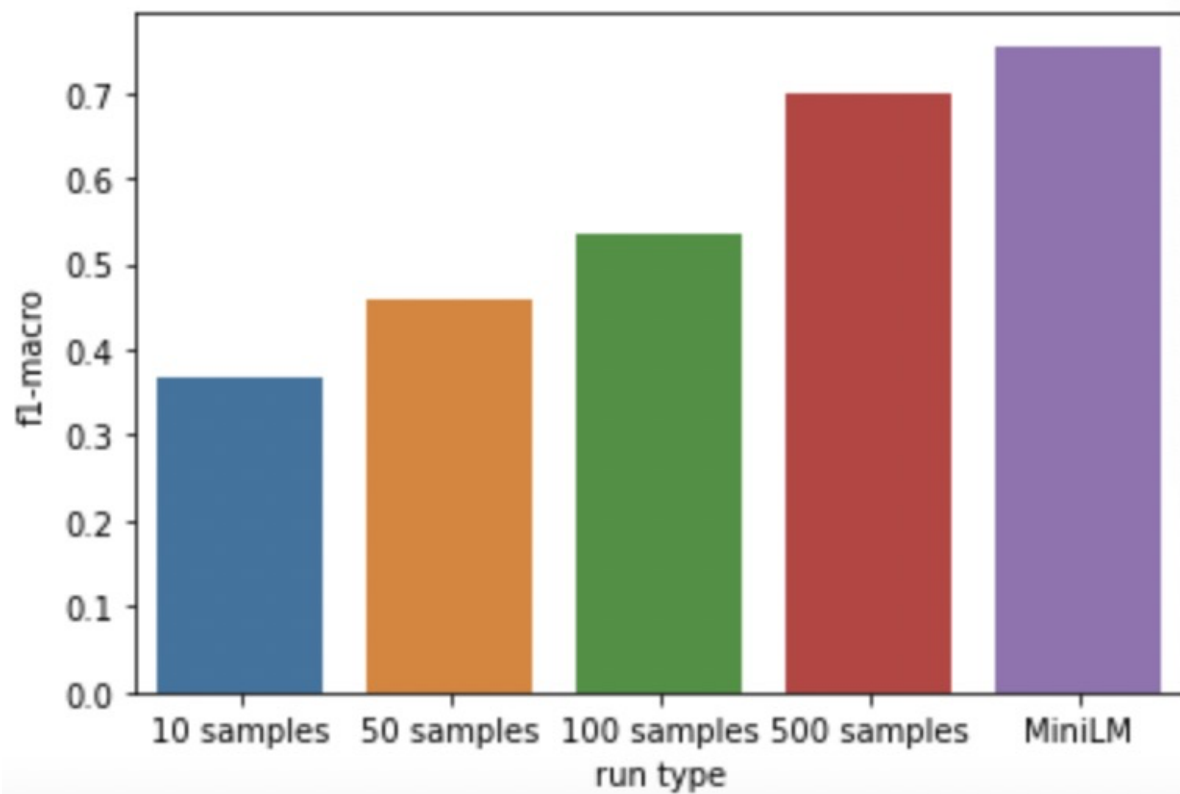
1: toxic, offensive

Example: The sentence [This does not seem to make any sense.] is quite [neutral].

Problem 4

The scores are:

Method	Score
Pet 10 samples	0.3662667123383554
Pet 50 samples	0.45964581053008463
Pet 100 samples	0.5344730447728572
Pet 500 samples	0.6975612118162938
MiniLM	0.754844



Training on around 700-1000 samples will result in a similar score to MiniLM model.