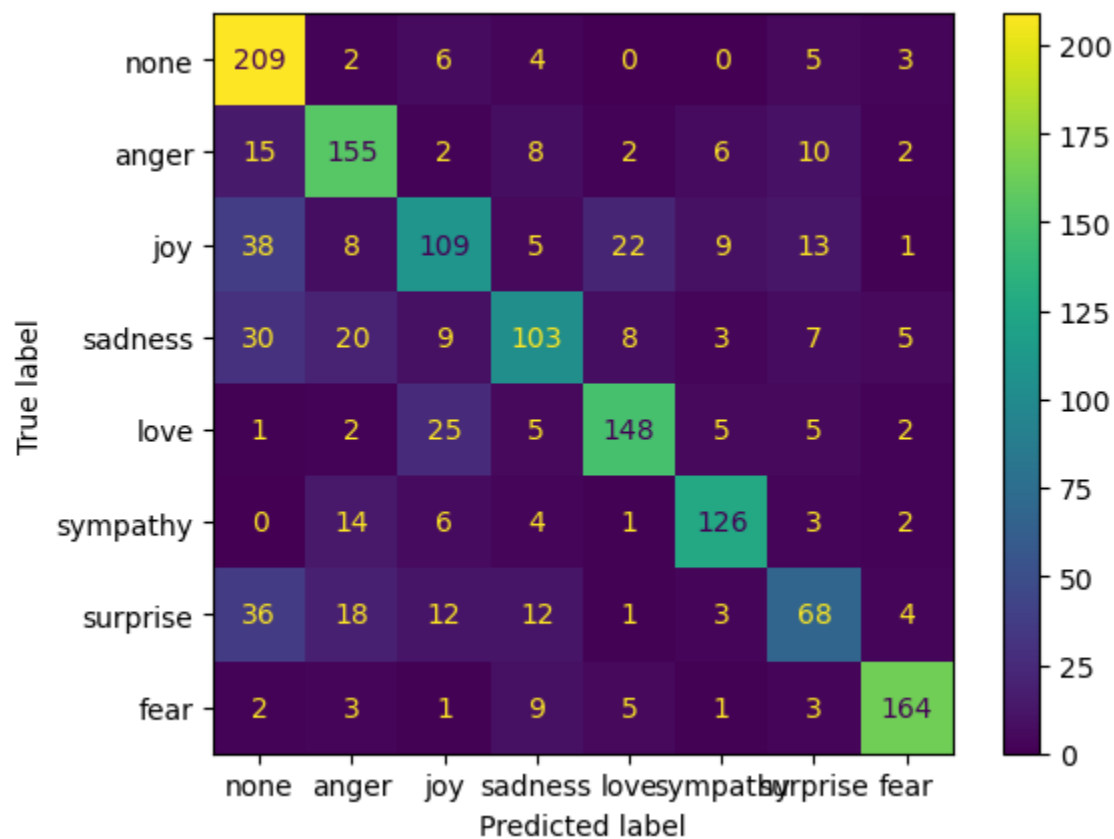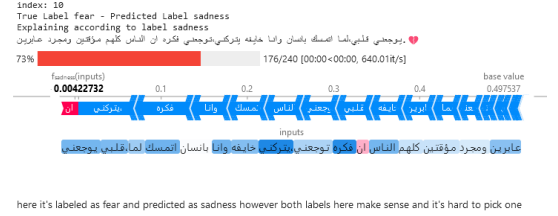We will be using the light stemming model from now.

|           | precision | recall | f1-score | support |
|-----------|-----------|--------|----------|---------|
| none      | 0.63      | 0.91   | 0.75     | 229     |
| anger     | 0.70      | 0.78   | 0.73     | 200     |
| joy       | 0.64      | 0.53   | 0.58     | 205     |
| sadness   | 0.69      | 0.56   | 0.61     | 185     |
| love      | 0.79      | 0.77   | 0.78     | 193     |
| sympathy  | 0.82      | 0.81   | 0.82     | 156     |
| surprise  | 0.60      | 0.44   | 0.51     | 154     |
| fear      | 0.90      | 0.87   | 0.88     | 188     |
|           |           |        |          |         |
| accuracy  |           |        | 0.72     | 1510    |
| macro avg | 0.72      | 0.71   | 0.71     | 1510    |
| weighted avg | 0.72   | 0.72   | 0.71     | 1510    |

# Light Stemming vs Raw Data:

index: 10
True Label fear - Predicted Label fear
Explaining according to label fear

index: 10
True Label fear - Predicted Label sadness
Explaining according to label sadness

Here the model made the correct prediction with light stemming.

index: 14
True Label sadness - Predicted Label anger
Explaining according to label anger

index: 14
True Label sadness - Predicted Label anger
Explaining according to label anger

explain_example(incorrectly_classified[6])

This example, i would label it as anger myself but it is labeled sadness

This example, i would label it as anger myself but it is labeled sadness

The model stopped using stopwords as explanations.

index: 74
True Label none - Predicted Label none
Explaining according to label none

explain_example(incorrectly_classified[21])

index: 74
True Label none - Predicted Label sadness
Explaining according to label sadness

In this case i would say that the prediction is more correct

This was a case which I agreed with the prediction more than the label, but now the model got the label correctly. The explanation is bad too. ***(try the correct word)



No significant change after correction?

Only that the word turned from being not-none to none but its score is not significant enough to change the outcome in either cases

Let's us try without the problematic word:



```
explanation = explainer.explain_instance(df_test[stemtype].iloc[testi], model_predict, num_features=10, labels = range(8))

# showing the explanation
explanation.show_in_notebook()
```

Now let us fix the word:

[0]

outputs

none anger joy **sadness** love sympathy **surprise** fear

| 0.3 | f_sadness(inputs) **0.378346** | 0.4 | | 0.5 | base value 0.563556 | 0.6 |

مموزه | سمعتنا | افريقيا

inputs

سمعتنا مموزه أفريقيا

```
+ Code    + Markdown
```

```
explanation = explainer.explain_instance(df_test[stemtype].iloc[testi], model_predict, num_features=10, labels = range(8))

    # showing the explanation
explanation.show_in_notebook()
```

Prediction probabilities

| sadness | 0.38 |
| surprise | 0.27 |
| anger | 0.18 |
| joy | 0.12 |
| Other | 0.05 |

NOT none — none

افريقيا 0.01

NOT anger — **anger**

افريقيا 0.07
مموزه 0.06
سمعتنا 0.01

NOT joy — joy

سمعتنا 0.05
افريقيا 0.04
مموزه 0.04

NOT sadness — **sadness**

مموزه 0.33
افريقيا 0.19
سمعتنا 0.14

NOT love — love

مموزه 0.03
افريقيا 0.01
سمعتنا 0.01

NOT sympathy — sympathy

سمعتنا 0.06
مموزه 0.05
افريقيا 0.01

NOT surprise — surprise

مموزه 0.20
سمعتنا 0.10
افريقيا 0.09

NOT fear — fear

افريقيا 0.02
سمعتنا 0.01
مموزه 0.00

Text with highlighted words

أفريقيا مموزه سمعتنا

We returned to sadness once more.

This one just got more confusing.





The model got confused in this case, as light stemming loses some of the context, so the mistake is understandable.

**Conclusion:**

The explanations got better but some of the context was lost.

# Shap vs Lime

**Interruption:**

Some errors were found in light stemming during the comparison of lime and shap.

Raw data:

💔 .يوجعني قلبي،لما اتمسك بانسان وانا خايفه يتركني،توجعني فكره ان الناس كلهم مؤقتين ومجرد عابرين

Light Stemming:

يوجع قلبيل اتمسك بانس وانا خايفه يتركنيتوجع فكره الناس مؤق ومجرد عابر حزن

```
index: 10
True Label fear - Predicted Label fear
Explaining according to label fear
يوجعني قلبي،لما اتمسك باسان وانا خايفه يتركني،توجعني فكره ان الناس كلهم مؤقتين ومجرد عابرين. ♥
```

base value
0.0382008 ... 0.877994

**Prediction probabilities**

| fear | 0.88 |
| sadness | 0.08 |
| love | 0.02 |
| sympathy | 0.02 |
| Other | 0.01 |

NOT fear          fear

Text with highlighted words

While the main contributor to the label is the same, the other words are assigned different importance in different explainers.



```
index: 12
True Label sadness - Predicted Label joy
Explaining according to label joy
المصريين داخلين الاوليمبياد تمثيل مشرف
```

base value
0.315782 ... 0.487558

**Prediction probabilities**

| joy | 0.49 |
| none | 0.29 |
| surprise | 0.16 |
| sadness | 0.04 |
| Other | 0.02 |

NOT joy          joy

Text with highlighted words

The same words contribute to different labels according to different explainers.



```
index: 14
True Label sadness - Predicted Label anger
Explaining according to label anger
حرا!!!!!!ام الظلم التحكيمي الفاجر ده Mortada
```

base value
0.0580826 ... 0.378002

**Prediction probabilities**

| anger | 0.38 |
| surprise | 0.25 |
| sadness | 0.19 |
| fear | 0.13 |
| Other | 0.05 |

NOT anger          anger

Text with highlighted words

Here, both have mostly the same explanation, but the words are assigned different weights.

index: 69
True Label fear - Predicted Label love
Explaining according to label love

😊 😊 😊 من كثر ما احبه ودي اروح احطيه من ايده بس خايفه يرفض

base value         f_love(inputs)

-0.1    0.0148965    0.1    0.2    0.3    0.4    **0.450807**    0.5    0.6

ودي    ابوه    احطيه    احبه    خايفه    يرفض

inputs
يرفض خايفه ابوه احطيه ودي اروح احبه كثر

**Prediction probabilities**

| | |
|---|---|
| love | 0.45 |
| fear | 0.42 |
| anger | 0.04 |
| joy | 0.04 |
| Other | 0.05 |

NOT love      love

احبه 0.25
احطيه 0.14
خايفه 0.14
ودي 0.11
كثر 0.11
يرفض 0.07
ابوه 0.03
اروح 0.01

**Text with highlighted words**

كثر احبه ودي اروح احطيه ابوه خايفه يرفض

Here, LIME gave better insight into the workings of the model by providing the probabilities. We can see that even the model understands that there are mixed feelings.

Observe the weight of each words in the following examples:

1)

index: 13
True Label anger - Predicted Label anger
Explaining according to label anger

؟؟؟؟؟؟؟؟ السعوديه الثالثه في السمنه مفروض الاولي ب السمنه علي كثر ما يسون تكميم

base value      f_anger(inputs)

-5.55112e-17    0.0374464    0.1    0.2    0.3    0.4    **0.475184**    0.5

السمنه    مفروض    السمنه    تكميم    يسون    السعوديه    الاولي    الثالثه    كثر

inputs
تكميم يسون كثر السمنه الاولي مفروض السمنه الثالثه السعوديه

**Prediction probabilities**

| | |
|---|---|
| anger | 0.48 |
| surprise | 0.18 |
| fear | 0.17 |
| sympathy | 0.06 |
| Other | 0.12 |

NOT anger      anger

السمنه 0.19
يسون 0.13
مفروض 0.11
تكميم
الثالثه 0.09
الثالثه 0.06
السعوديه 0.04
الاولي 0.03
كثر 0.03

**Text with highlighted words**

السعوديه الثالثه السمنه مفروض الاولي السمنه كثر يسون تكميم

+ Code    + Markdown

```
explain_example(test1)
```

2)

```
index: 14
True Label sadness - Predicted Label anger
Explaining according to label anger
Mortada نه الفاجر التحكيمي الظلم حرااااام
```

base value | f_anger(inputs)
0.0374464 | 0.367688

-0.1    0    0.1    0.2    0.3    0.4    0.5

التحكيمي    الفاجر    الظلم

inputs
حرام الظلم التحكيمي الفاجر

**Prediction probabilities**

| | |
|---|---|
| anger | 0.37 |
| surprise | 0.25 |
| fear | 0.20 |
| sadness | 0.15 |
| Other | 0.04 |

NOT anger      anger

الفاجر 0.27
التحكيمي 0.11
حرام 0.03
الظلم 0.07

**Text with highlighted words**

حرام الظلم التحكيمي الفاجر

---

```
explain_example(testi)
```

```
index: 14
True Label sadness - Predicted Label anger
Explaining according to label anger
Mortada نه الفاجر التحكيمي الظلم حرااااام
```

base value | f_anger(inputs)
0.0374464 | 0.367688

-0.1    0    0.1    0.2    0.3    0.4    0.5

التحكيمي    الفاجر    الظلم

inputs
حرام الظلم التحكيمي الفاجر

**Prediction probabilities**

| | |
|---|---|
| anger | 0.37 |
| surprise | 0.25 |
| fear | 0.20 |
| sadness | 0.15 |
| Other | 0.04 |

NOT anger      anger

الفاجر 0.26
التحكيمي 0.11
حرام 0.03
الظلم 0.03

**Text with highlighted words**

حرام الظلم التحكيمي الفاجر

3)

```
index: 15
True Label love - Predicted Label love
Explaining according to label love
وائل كفوري خششششق الفاجر ❤ 😍 🏠 تمني الخلص منه لادخل الجنه Admin 👤
```

base value | f_love(inputs)
-5.55112e-17 | 0.0430801 | 0.653303

0.1   0.2   0.3   0.4   0.5   0.6   0.7

وائل   لادخل   كفوري   الجنه   الخلص   عشقو   حزن

inputs
حزن حزن الجنه لادخل اتخلص اتمني كفوري عشقو وائل

**Prediction probabilities**

| | |
|---|---|
| love | 0.65 |
| sadness | 0.16 |
| joy | 0.09 |
| sympathy | 0.06 |
| Other | 0.03 |

NOT love      love

عشقو 0.33
الجنه 0.18
كفوري 0.13
اتمني 0.11
لادخل 0.09
حزن 0.09
الخلص 0.05
وائل 0.02

**Text with highlighted words**

وائل كفوري عشقو اتمني اتخلص لادخل الجنه حزن حزن

---

```
explain_example(testi)
```

```
index: 15
True Label love - Predicted Label love
Explaining according to label love
وائل كفوري خششششق الفاجر ❤ 😍 🏠 تمني الخلص منه لادخل الجنه Admin 👤
```

base value | f_love(inputs)
-5.55112e-17 | 0.0430801 | 0.653303

0.1   0.2   0.3   0.4   0.5   0.6   0.7

وائل   لادخل   كفوري   الجنه   الخلص   عشقو   حزن

inputs
حزن حزن الجنه لادخل اتخلص اتمني كفوري عشقو وائل

**Prediction probabilities**

| | |
|---|---|
| love | 0.65 |
| sadness | 0.16 |
| joy | 0.09 |
| sympathy | 0.06 |
| Other | 0.03 |

NOT love      love

عشقو 0.34
الجنه 0.18
كفوري 0.13
اتمني 0.11
لادخل 0.10
حزن 0.09
الخلص 0.05
وائل 0.01

**Text with highlighted words**

وائل كفوري عشقو اتمني اتخلص لادخل الجنه حزن حزن

4)

index: 16
True Label fear - Predicted Label fear
Explaining according to label fear
مُش معقول انا ايدي لامسه ايديك .. خايف اكون يا حبيبي بحلم بيك :

Prediction probabilities

| | |
|---|---|
| fear | 0.83 |
| surprise | 0.08 |
| joy | 0.04 |
| love | 0.02 |
| Other | 0.02 |

NOT fear    fear

Text with highlighted words

مش معقول ايدي لامسه ايديك خايف اكون ياحبيبي بحلم بيك

index: 16
True Label fear - Predicted Label fear
Explaining according to label fear
مُش معقول انا ايدي لامسه ايديك .. خايف اكون يا حبيبي بحلم بيك :

Prediction probabilities

| | |
|---|---|
| fear | 0.83 |
| surprise | 0.08 |
| joy | 0.04 |
| love | 0.02 |
| Other | 0.02 |

NOT fear    fear

Text with highlighted words
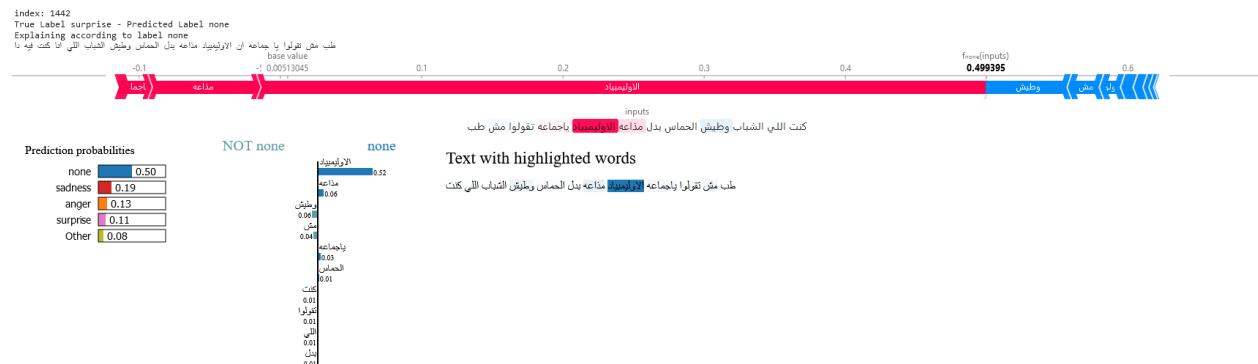
مش معقول ايدي لامسه ايديك خايف اكون ياحبيبي بحلم بيك

Due to the nature of its calculations, SHAP is very consistent in its output. Lime on the other hand is inconsistent due to its random nature.

Explaining the same example using Lime multiple times gives different weights each time but it's usually generally correct.

The most significant word is almost always the same and the probabilities don't appear to change but it would make a difference if the values were close to each other such that a 0.01 is enough to make a difference.

index: 1442
True Label surprise - Predicted Label none
Explaining according to label none
طب مش تقولوا يا جماعه ان الاولمبياد مذاعه بدل الحماس وطيش الشباب اللي انا كنت فيه دا

Prediction probabilities

| | |
|---|---|
| none | 0.50 |
| sadness | 0.19 |
| anger | 0.13 |
| surprise | 0.11 |
| Other | 0.08 |

NOT none    none

Text with highlighted words

طب مش تقولوا ياجماعه الاولمبياد مذاعه بدل الحماس وطيش الشباب اللي كنت

base value
-0.00513045
0.1
0.2
0.3
0.4
f_none(inputs)
**0.499395**
0.6

اجما | مذاعه | الاولمبياد | وطيش | بل | مش

inputs
كنت اللي الشباب وطيش الحماس بدل مذاعه الاولمبياد ياجماعه تقولوا مش طب

**Prediction probabilities**

| none | 0.50 |
| sadness | 0.19 |
| anger | 0.13 |
| surprise | 0.11 |
| Other | 0.08 |

NOT none          none

الاولمبياد 0.52
وطيش 0.08
مذاعه 0.05
مش 0.04
ياجماعه 0.03
الحماس 0.01
تقولوا 0.01
كنت 0.01
الشباب 0.01
اللي 0.01

**Text with highlighted words**

طب مش تقولوا ياجماعه مذاعه بدل الحماس وطيش الشباب اللي كنت

Here the order of significance is different.

76%    208/272 [00:02<00:00, 772.65it/s]

base value
-0.00513045
0.1
0.2
0.3
0.4
0.5
f_none(inputs)
**0.569874**
0.6

ولوح | عا | ة | بيه | برضه | الاولمبياد | ع | داعش

inputs
برضه الاولمبياد بيه ولوح شاله واقع الشهادتين ومكتوب السعوديه علم بدلا داعش علم لقي طلعت حماده

**Prediction probabilities**

| none | 0.57 |
| surprise | 0.15 |
| anger | 0.15 |
| sadness | 0.09 |
| Other | 0.04 |

NOT none          none

الاولمبياد 0.53
بيه 0.05
داعش 0.03
الارض 0.03
طلعت 0.03
برضه 0.02
ولوح 0.02
ومكتوب 0.02
الشهادتين 0.02
شاله 0.02

**Text with highlighted words**

حماده طلعت لقي علم داعش بدلا علم السعوديه ومكتوب الشهادتين واقع الارض شاله ولوح بيه الاولمبياد برضه

76%    208/272 [00:01<00:00, 781.19it/s]

base value
-0.00513045
0.1
0.2
0.3
0.4
0.5
f_none(inputs)
**0.569874**
0.6

ولوح | عا | ة | بيه | برضه | الاولمبياد | ع | داعش

inputs
برضه الاولمبياد بيه ولوح شاله واقع الشهادتين ومكتوب السعوديه علم بدلا داعش علم لقي طلعت حماده

**Prediction probabilities**

| none | 0.57 |
| surprise | 0.15 |
| anger | 0.15 |
| sadness | 0.09 |
| Other | 0.04 |

NOT none          none

الاولمبياد 0.53
بيه 0.04
داعش 0.03
طلعت 0.03
الارض 0.03
حماده 0.02
ومكتوب 0.02
برضه 0.02
ولوح 0.02
شاله 0.01

**Text with highlighted words**

حماده طلعت لقي علم داعش بدلا علم السعوديه ومكتوب الشهادتين واقع الارض شاله ولوح بيه الاولمبياد برضه

Here we can observe that some words appeared in the first sentence but not the second sentence and vice versa.

base value
-0.00513045
0.1
0.2
0.3
0.4
0.5
f_none(inputs)
**0.591645**
0.7

عيال | يعمل | ممكن | الاولمبياد | حاجه | طربه | جيل

inputs
الاولمبياد حاجه يعمل ممكن فشيخ جيل طربه عيال

**Prediction probabilities**

| none | 0.59 |
| joy | 0.20 |
| sadness | 0.11 |
| surprise | 0.04 |
| Other | 0.06 |

NOT none          none

الاولمبياد 0.54
ممكن 0.17
يعمل 0.10
حاجه 0.06
فشيخ 0.06
طربه 0.05
عيال 0.01
جيل 0.00

**Text with highlighted words**

عيال طربه جيل فشيخ ممكن يعمل حاجه الاولمبياد

index: 82
True Label none - Predicted Label none
Explaining according to label none

نزل عيال طزيه بس كان جيل فشيخ كان ممكن يعمل حاجه في الاوليمبياد

base value

f_none(inputs)
**0.591645**

-0.1          0.00513045        0.1        0.2        0.3        0.4        0.5        0.6        0.7

عيال | يعمل | مش | الاوليمبياد | جيل | طزيه | حاجه | مش

inputs

الاوليمبياد حاجه يعمل ممكن فشيخ جيل طزيه عيال

Prediction probabilities

none     0.59
joy      0.20
sadness  0.11
surprise 0.04
Other    0.06

NOT none            none

الاوليمبياد 0.54
ممكن 0.17
يعمل 0.10
فشيخ 0.07
حاجه 0.06
طزيه 0.05
عيال 0.03
جيل 0.00

Text with highlighted words

عيال طزيه جيل ممكن فشيخ حاجه يعمل الاوليمبياد

As we can see here, a small change was enough to make a difference.

The last word switched from None to Not-none.

Now, we attempted to figure out why some words don't have weight in shap but have high weight in lime:

index: 16
True Label fear - Predicted Label fear
Explaining according to label fear

مش معقول انا ايدي لامسه ايديك .. خايف اكون يا حبيبي بحلم بيك :

base value

f_fear(inputs)
**0.830301**

0          0.0393674          0.2          0.4          0.6          0.8

مش | اكون | ايديك | ايدي | ياحبيبي | خايف | لامسه

inputs

-0.007

بيك بحلم ياحبيبي اكون خايف ايديك لامسه ايديمعقول مش

Prediction probabilities

fear     0.83
surprise 0.08
joy      0.04
love     0.02
Other    0.02

NOT fear            fear

خايف 0.72
معقول 0.09
ياحبيبي 0.07
بحلم 0.06
اكون 0.06
ايدي 0.05
بيك 0.03
ايديك 0.02
مش 0.01
لامسه 0.00

Text with highlighted words

مش معقول ايدي لامسه ايديك خايف اكون ياحبيبي بحلم بيك

Let's return to this example once more.

The word 'معقول' has a weight of -0.09 in lime (it was verified by multiple runs) but appears to have no significant weight with shape (-0.007)

Across multiple runs I tried to remove words that has little to no weight in both lime and shap as following:

index: 16
True Label fear - Predicted Label fear
Explaining according to label fear

مش معقول انا ايدي لامسه ايديك .. خايف اكون يا حبيبي بحلم بيك :

base value

f_fear(inputs)
**0.844052**

0          0.0393674          0.2          0.4          0.6          0.8          1

ايدي | اكون | ياحبيبي | ايديك | خايف | لامسه | بحلم

inputs

بيك بحلم ياحبيبي اكون خايف ايديك لامسه ايدي معقول

Prediction probabilities

fear     0.84
surprise 0.08
joy      0.03
love     0.02
Other    0.02

NOT fear            fear

خايف 0.71
معقول 0.09
ايدي 0.07
ياحبيبي 0.06
اكون 0.05
بحلم 0.03
بيك 0.04
لامسه 0.01
ايديك 0.01

Text with highlighted words

معقول ايدي لامسه ايديك خايف اكون ياحبيبي بحلم بيك

```
index: 16
True Label fear - Predicted Label fear
Explaining according to label fear
مش معقول انا ايدي لاسه اديك .. خايف- اكون يا حبيبي بطم بيك ;
```

As we can observe in lime, other than the removed word, the words have the same order and a little difference in weight.

But in Shap, each removed word has a significant effect on the rest of the weights even if the removed word's weight is not significant itself.

**Conclusion:**

Due to the way of calculation, Lime has little in the way of stability when running the same example multiple times but is surprisingly robust when some of the input is changed or dropped.

Shap on the other hand is very consistent in its calculations when running the same example, but it is greatly affected by small changes in the input and seemingly unrelated words may have significant changes in the output.

**Interruption (again):**

Some errors were found in light stemming during the comparison of lime and shap.

Raw data:

دول افريفيا مبوزه سمعتنا في الاوليمبياد

Light Stemming:

افريف مبوزه سمعت الاوليمبياد

There is a spelling mistake here, and some removed stopwords that need to be looked into.
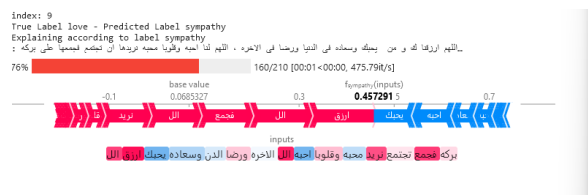
And another one:

Raw data:

محدش عارف هما تعبوا اد ايه عشان يوصلوا الاوليمبياد و هم مش ناقصين تعليقات الجهله في اللعبه لادائهم و لا تعليقات المصريين عموما علي لبسهم
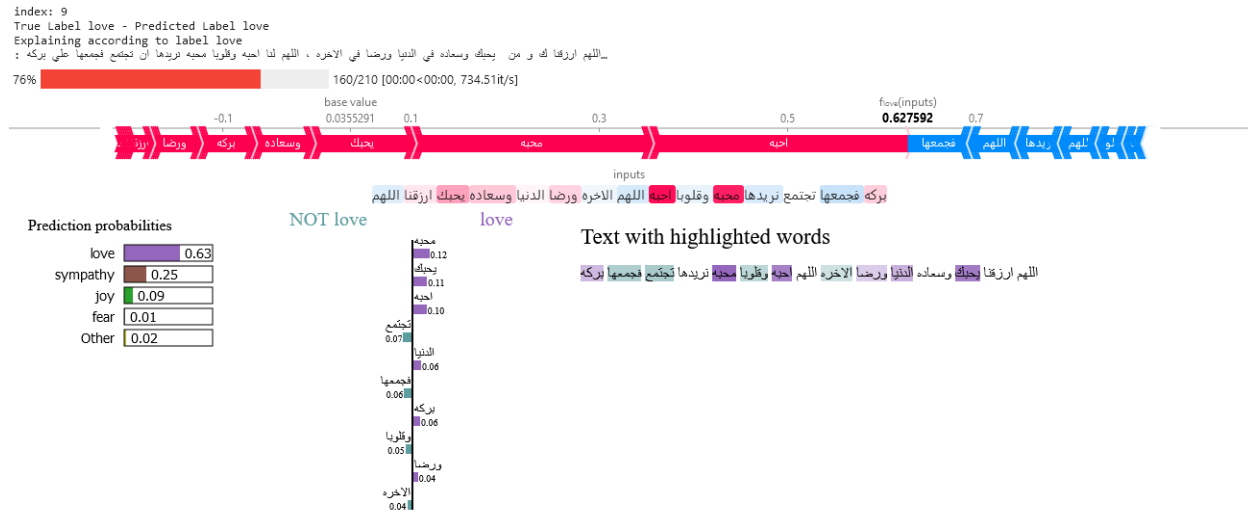
Light Stemming:

محدش عارف تعب اد عشان يوصل الاوليمبياد مش ناقص تعليق الجهله اللعبه لادائ تعليق المصر عمو لبس

There are some removed stopwords that need to be looked into.

Using Cleaned data without stemming:



There was this error that I attributed to a loss of context but without stemming, the model got it correctly once more:

index: 9
True Label love - Predicted Label love
Explaining according to label love

اللهم ارزقنا اك و من يحبك وسعاده في الدنيا ورضا في الاخره ، اللهم لنا احبه وقلوبا محبه تريدها ان تجتمع فجمعها طي بركه :

76% ▮▮▮▮▮▮▮▮▮  160/210 [00:00<00:00, 734.51it/s]



| | base value | | | | f love(inputs) | |
|---|---|---|---|---|---|---|
| -0.1 | 0.0355291 | 0.1 | 0.3 | 0.5 | **0.627592** | 0.7 |

inputs

بركه فجمعها تجتمع تريدها محبه وقلوبا احبه اللهم الاخره ورضا الدنيا وسعاده يحبك ارزقنا اللهم

NOT love          love

Text with highlighted words

اللهم ارزقنا يحبك وسعاده الدنيا ورضا الاخره اللهم احبه وقلوبا محبه تريدها تجتمع فجمعها بركه

Prediction probabilities

| | |
|---|---|
| love | 0.63 |
| sympathy | 0.25 |
| joy | 0.09 |
| fear | 0.01 |
| Other | 0.02 |

محبه 0.12
يحبك 0.11
احبه 0.10
تجتمع 0.07
الدنيا 0.06
فجمعها 0.06
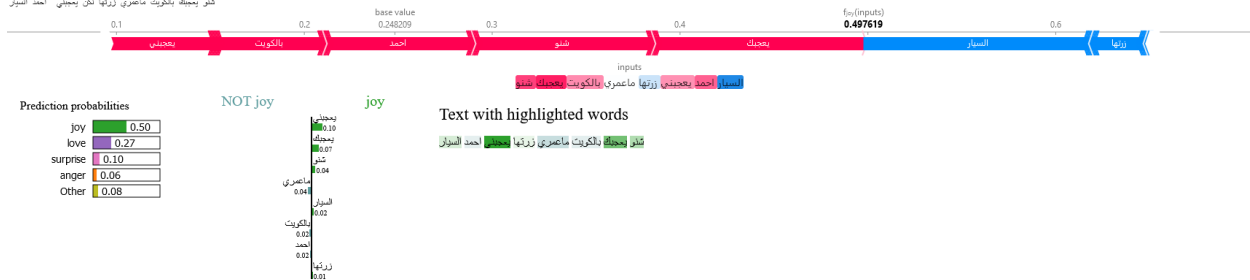بركه 0.06
وقلوبا 0.05
ورضا 0.04
الاخره 0.04

So, context can be lost with stemming but it appears that without stemming, it may get better.

Since we will be using cleaned data now, here are the metrics:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| none | 0.64 | 0.89 | 0.74 | 229 |
| anger | 0.71 | 0.79 | 0.75 | 200 |
| joy | 0.62 | 0.56 | 0.58 | 205 |
| sadness | 0.68 | 0.55 | 0.61 | 185 |
| love | 0.79 | 0.76 | 0.77 | 193 |
| sympathy | 0.85 | 0.83 | 0.84 | 156 |
| surprise | 0.59 | 0.44 | 0.50 | 154 |
| fear | 0.91 | 0.87 | 0.89 | 188 |
| | | | | |
| accuracy | | | 0.72 | 1510 |
| macro avg | 0.72 | 0.71 | 0.71 | 1510 |
| weighted avg | 0.72 | 0.72 | 0.71 | 1510 |

By comparing the confusion matrix with the previous one, we can observe that the box that indicates (true = none, pred = joy) increased from 6 to 10.

That was despite the fact that context appears to be more complete without stemming so lets examine those 10 samples.

index: 26
True Label none - Predicted Label joy
Explaining according to label joy

(() : الاوليمبياد كانت عباره عن باص من تريكه يوصل صلاح بالجول وجول البرازيل وشاره الكابتن .. وحشتها يامع

76%    160/210 [00:00<00:00, 738.97it/s]

base value      f_joy(inputs)
0.248209         0.463062

NOT joy          joy

Text with highlighted words

Prediction probabilities
joy 0.46
none 0.33
surprise 0.09
sadness 0.06
Other 0.06

وحشتنا 0.20
الاوليمبياد 0.16
تريكه 0.09
يوصل 0.07
البرازيل 0.07
فرح 0.06
صلاح 0.03
الكابتن 0.02
يامع 0.01
وشاره 0.01

**index: 100**
True Label none - Predicted Label joy
Explaining according to label joy
شنو يعجبك بالكويت ماعصري زرتها لكن يعجبني احمد السيار

base value: 0.248209
f_{joy}(inputs): 0.497619

يعجبني | بالكويت | احمد | شنو | يعجبك | السيار | زرتها

inputs
السيار احمد يعجبني زرتها ماعصري بالكويت يعجبك شنو

Prediction probabilities
joy 0.50
love 0.27
surprise 0.10
anger 0.06
Other 0.08

NOT joy | joy

يعجبني 0.10
يعجبك 0.07
شنو 0.04
ماعصري 0.04
السيار 0.02
بالكويت 0.02
احمد 0.02
زرتها 0.01

Text with highlighted words
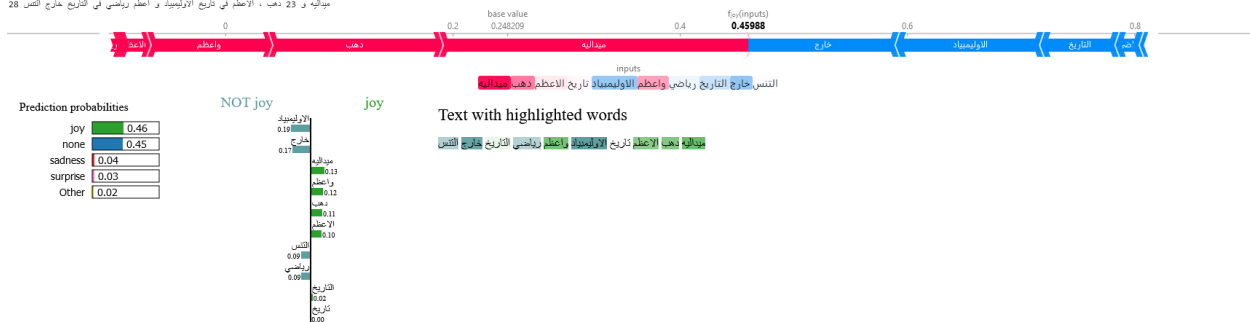شنو يعجبك بالكويت ماعصري زرتها لكن يعجبني احمد السيار

---

**index: 325**
True Label none - Predicted Label joy
Explaining according to label joy
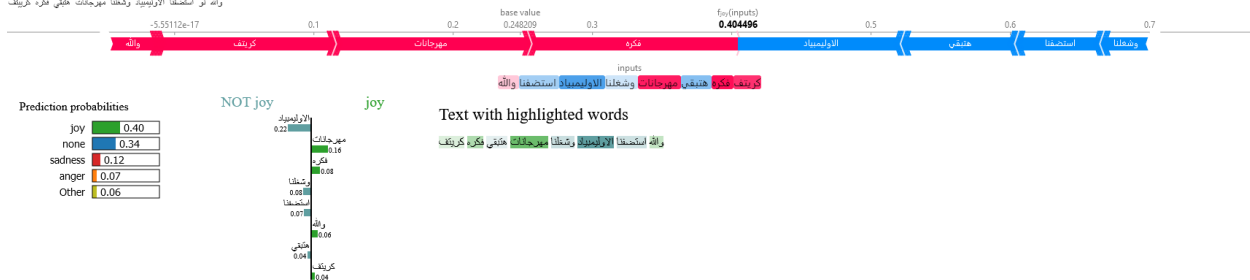يارب يرفعوا اسم مصر بادائهم واخلاقهم شجموا مصر في الاولمبياد

base value: 0.248209
f_{joy}(inputs): 0.479533

زرت | واخلاقهم | اسم | الاولمبياد | بادائهم | يرفعوا | مصر | مصر | شنجوا

inputs
الاولمبياد مصر شجموا واخلاقهم بادائهم مصر اسم يرفعوا يارب

Prediction probabilities
joy 0.48
none 0.37
sadness 0.08
surprise 0.04
Other 0.03

NOT joy | joy

واخلاقهم 0.16
اسم 0.11
شجموا 0.10
مصر 0.10
يرفعوا 0.06
بادائهم 0.05
الاولمبياد 0.02
يارب 0.00

Text with highlighted words
يارب يرفعوا اسم مصر بادائهم واخلاقهم شجموا مصر الاولمبياد

---

**index: 801**
True Label none - Predicted Label joy
Explaining according to label joy
احتفالات الاولمبياد 1 احتفالات قناه السويس 0

base value: 0.248209
f_{joy}(inputs): 0.463415

احتفالات | السويس | احتفالات | الاولمبياد

inputs
السويس قناه احتفالات الاولمبياد احتفالات

Prediction probabilities
joy 0.46
none 0.36
sadness 0.11
surprise 0.04
Other 0.02

NOT joy | joy

احتفالات 0.34
السويس 0.17
الاولمبياد 0.15
قناة 0.06

Text with highlighted words
احتفالات الاولمبياد احتفالات قناة السويس

---

**index: 815**
True Label none - Predicted Label joy
Explaining according to label joy
ميداليه و 23 دهب ، الاعظم في تاريخ الاولمبياد و اعطم رياضي في التاريخ خارج التنس 28

base value: 0.248209
f_{joy}(inputs): 0.45988

الاعظ | واعظم | دهب | ميداليه | خارج | الاولمبياد | التاريخ | التن

inputs
التنس خارج التاريخ رياضي واعظم الاولمبياد تاريخ الاعظم دهب ميداليه

Prediction probabilities
joy 0.46
none 0.45
sadness 0.04
surprise 0.03
Other 0.02

NOT joy | joy

الاولمبياد 0.19
خارج 0.17
ميداليه 0.13
واعظم 0.12
دهب 0.11
الاعظم 0.10
التنس 0.09
رياضي 0.09
التاريخ 0.02
تاريخ 0.00

Text with highlighted words
ميداليه دهب الاعظم تاريخ الاولمبياد واعظم رياضي التاريخ خارج التنس

---

**index: 874**
True Label none - Predicted Label joy
Explaining according to label joy
والله لو استضفنا الاولمبياد وشعلنا مهرجانات هنيقى فكره كريتف

base value: 0.248209
f_{joy}(inputs): 0.404496

والله | كريتف | مهرجانات | فكرة | الاولمبياد | هنيقى | استضفنا | وشعلنا

inputs
كريتف فكره هنيقى مهرجانات وشعلنا الاولمبياد استضفنا والله

Prediction probabilities
joy 0.40
none 0.34
sadness 0.12
anger 0.07
Other 0.06

NOT joy | joy

الاولمبياد 0.22
مهرجانات 0.16
فكرة 0.08
وشعلنا 0.08
استضفنا 0.07
والله 0.06
هنيقى 0.04
كريتف 0.04

Text with highlighted words
والله استضفنا الاولمبياد وشعلنا مهرجانات هنيقى فكره كريتف

index: 904
True Label none - Predicted Label joy
Explaining according to label joy
😊😀 مصر رقم واحد في المركز الخامس في الاوليمبياد

base value f_joy(inputs)
0.248209 0.673445

inputs
فرح فرح الاوليمبياد الخامس المركز واحد رقم مصر

Prediction probabilities
joy 0.67
none 0.19
surprise 0.11
sadness 0.02
Other 0.01

NOT joy          joy

Text with highlighted words
مصر رقم واحد المركز الخامس الاوليمبياد فرح فرح

index: 917
True Label none - Predicted Label joy
Explaining according to label joy
Rio2016 التجنيس ده حلو يا فخري حلو يا فخري يا جدع بيجيب ميداليات في الاوليمبياد اجدع

base value f_joy(inputs)
0.248209 0.620837

inputs
اجدع الاوليمبياد ميداليات بيجيب ياجدع اوى حلو يافخري حلو التجنيس

Prediction probabilities
joy 0.62
surprise 0.19
none 0.14
anger 0.03
Other 0.02

NOT joy          joy

Text with highlighted words
التجنيس حلو يافخري حلو اوي ياجدع بيجيب ميداليات الاوليمبياد اجدع

index: 1269
True Label none - Predicted Label joy
Explaining according to label joy
العاشره سودانا اما بعد: الاوليمبياد حقتنا والاولاد اولادنا ولا عزاء للمتشائمين

base value f_joy(inputs)
0.248209 0.450508

inputs
للمتشائمين عزاء اولادنا والاولاد حقتنا الاوليمبياد اما سودانا العاشره

Prediction probabilities
joy 0.45
none 0.31
sympathy 0.12
sadness 0.07
Other 0.05

NOT joy          joy

Text with highlighted words
العاشره سودانا اما الاوليمبياد حقتنا والاولاد اولادنا عزاء للمتشائمين

index: 1315
True Label none - Predicted Label joy
Explaining according to label joy
🏅🏅 .. احنا عدنا ف الاوليمبياد ناس زي الذهب

base value f_joy(inputs)
0.248209 0.502388

inputs
الذهب زي ناس الاوليمبياد

Prediction probabilities
joy 0.50
none 0.27
surprise 0.14
sadness 0.05
Other 0.03

NOT joy          joy

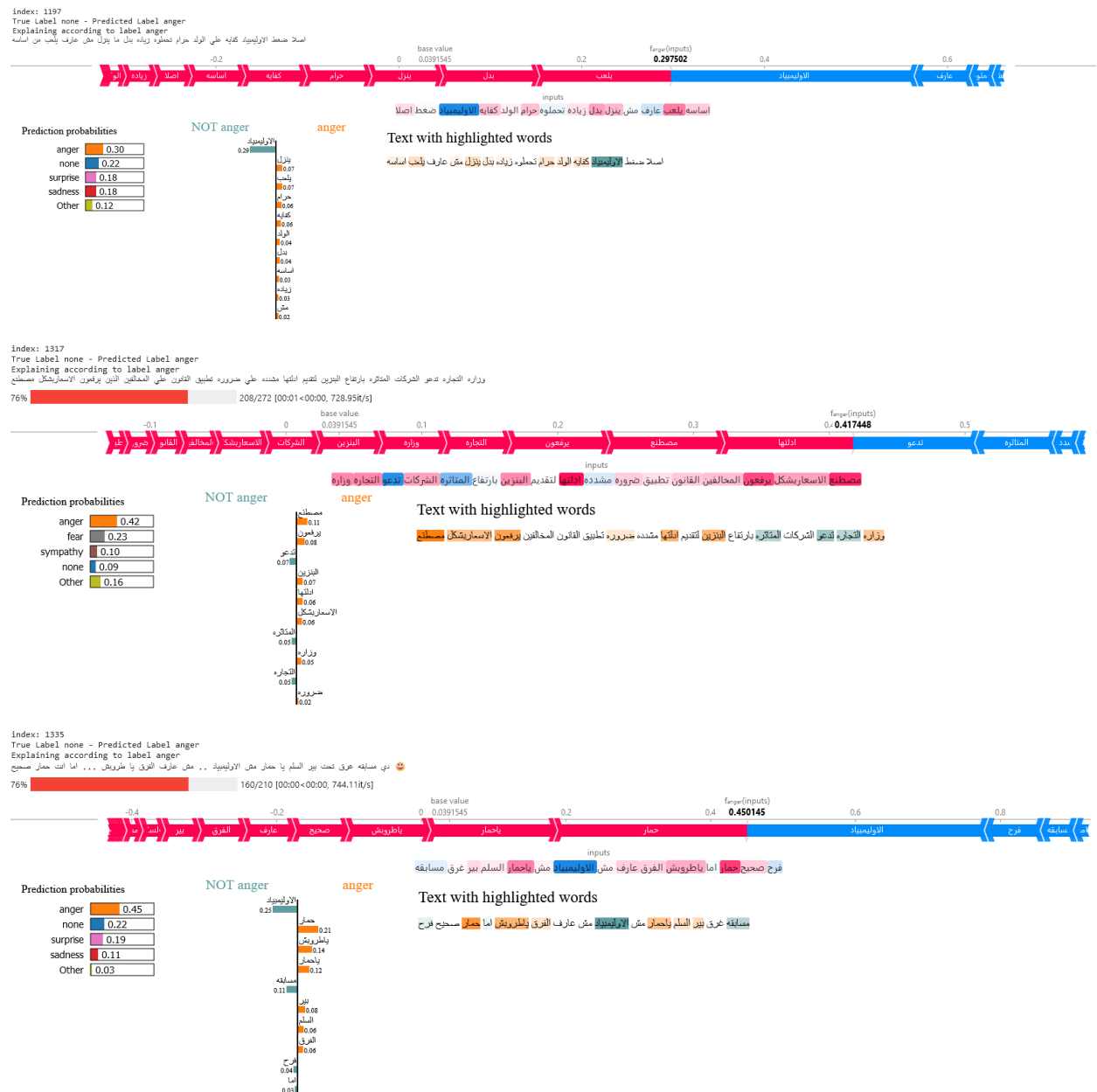Text with highlighted words
🏅🏅 الاوليمبياد ناس زي الذهب

Some examples are understandable why the model predicted them as joy, while others are just confusing to me and the whole context did not help either.

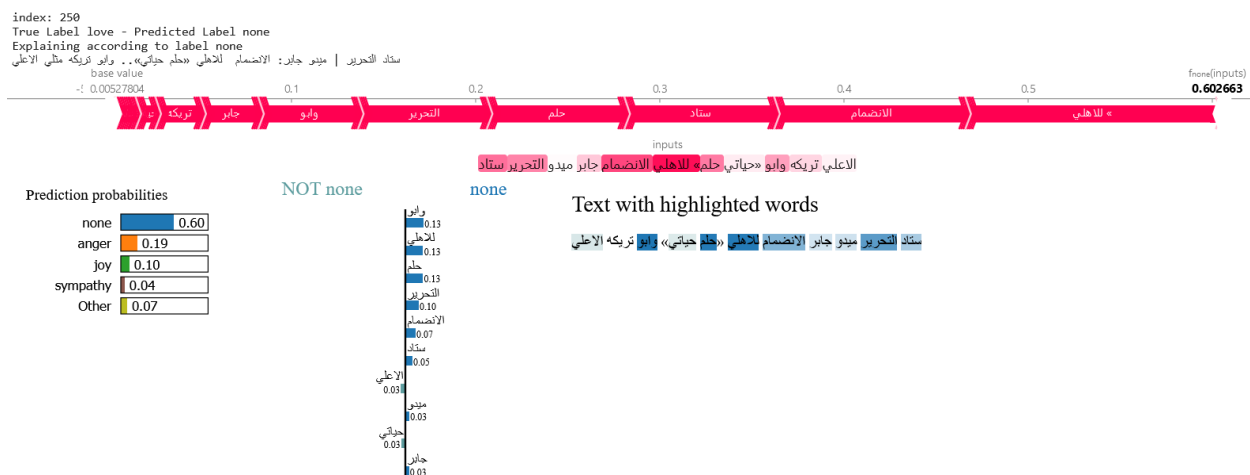The just context made the model more confused.

So, let's examine something a bit less confusing.

(true = none, pred = anger, total error = 3)

index: 1197
True Label none - Predicted Label anger
Explaining according to label anger

اصلا ضغط الاوليمبياد كفايه على الولد حرام تحملوه زياده بدل ما ينزل مش عارف يلعب من اساسه

base value

f(inputs)
**0.297502**

inputs

اساسه يلعب عارف مش ينزل بدل زياده تحمله حرام الولد كفايه الاوليمبياد ضغط اصلا

Prediction probabilities

| anger | 0.30 |
| none | 0.22 |
| surprise | 0.18 |
| sadness | 0.18 |
| Other | 0.12 |

NOT anger        anger

Text with highlighted words

اصلا ضغط الاوليمبياد كفايه الولد حرام تحمله زياده بدل ينزل مش عارف يلعب اساسه

الاوليمبياد 0.29
ينزل 0.07
يلعب 0.07
حرام 0.06
كفايه 0.06
الولد 0.04
بدل 0.04
اساسه 0.03
زياده 0.03
مش 0.02

index: 1317
True Label none - Predicted Label anger
Explaining according to label anger

وزاره التجاره تدعو الشركات المتأثره بارتفاع البنزين لتقديم ادلتها مشدده على ضروره تطبيق القانون الذين يرفعون الاسعاربشكل مصطنع

76%        208/272 [00:01<00:00, 728.95it/s]

base value
0.0391545

f(inputs)
**0.417448**

inputs

مصطنع الاسعاربشكل يرفعون المخالفين القانون تطبيق ضروره مشدده ادلتها لتقديم البنزين بارتفاع المتأثره الشركات تدعو التجاره وزاره

Prediction probabilities

| anger | 0.42 |
| fear | 0.23 |
| sympathy | 0.10 |
| none | 0.09 |
| Other | 0.16 |

NOT anger        anger

Text with highlighted words

وزاره التجاره تدعو الشركات المتأثره بارتفاع البنزين لتقديم ادلتها مشدده ضروره تطبيق القانون المخالفين يرفعون الاسعاربشكل مصطنع

مصطنع 0.11
يرفعون 0.08
تكبير 0.07
البنزين 0.07
ادلتها 0.06
الاسعاربشكل 0.06
المتأثره 0.05
وزاره 0.05
التجاره 0.05
ضروره 0.02

index: 1335
True Label none - Predicted Label anger
Explaining according to label anger

😊 دي مسابقه عرق تحت بير السلم يا حمار مش الاوليمبياد .. مش عارف الفرق يا طروريش ... اما انت حمار صحيح

76%        160/210 [00:00<00:00, 744.11it/s]

base value
0.0391545

f(inputs)
**0.450145**

inputs

فرح صحيح حمار اما ياطروريش الفرق عارف مش الاوليمبياد مش ياحمار السلم بير عرق مسابقه

Prediction probabilities

| anger | 0.45 |
| none | 0.22 |
| surprise | 0.19 |
| sadness | 0.11 |
| Other | 0.03 |

NOT anger        anger

Text with highlighted words

مسابقه عرق بير السلم ياحمار مش الاوليمبياد مش عارف الفرق ياطروريش اما حمار صحيح فرح

الاوليمبياد 0.25
حمار 0.21
ياطروريش 0.14
ياحمار 0.12
مسابقه 0.11
بير 0.08
السلم 0.06
الفرق 0.06
فرح 0.04
اما 0.03

The first and last example are understandable why the model would label them as anger, and I am inclined to agree with its explanation. Especially the last one, as 3 curse words in one sentence is a convincing argument for why it is predicted anger.

(true = love, pred = none, total error = 1)



```
index: 250
True Label love - Predicted Label none
Explaining according to label none
```

Most of the love samples in the dataset are the romantic kind of love, so the problem here is a lack of data in this very specific area.

(true = love, pred = joy, total error = 30)

index: 203

True Label love - Predicted Label joy

Original: احلام العنزي فخر الكويت الف الف مبرووووك بنت قبيلتي وبنت ديرتيييي؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟ كفو عيال الديره كل يوم اسم رافع اسمنا بشاره خير ان شاء الله

Cleaned: احلام العنزي فخر الكويت الف الف مبروك بنت قبيلتي وبنت ديرتي كفو عيال الديره يوم اسم رافع اسمنا بشاره خير شاء الله

I would say this one is mixed feelings.

index: 594

True Label love - Predicted Label joy

Original: الدوله العثمانيه الاسلاميه ستبقي رغم انوف الاعداء ربي احفظ تركيا المسلمه من مكر الماكرين واحفظ كل بلاد المسلمين ♡

Cleaned: الدوله العثمانيه الاسلاميه ستبقي رغم انوف الاعداء ربي احفظ تركيا المسلمه مكر الماكرين واحفظ بلاد المسلمين ♡

This one supports my earlier argument about romantic love and love for one's country.

index: 823

True Label love - Predicted Label joy

Original:  لون حياتك بحب وطاعه الله فهي طويله قصيره ...  ايمان  تقوه

Cleaned:  لون حياتك بحب وطاعه الله طويله قصيره ايمان تقوه

Same as the above.

**Conclusion on True = love:**

It appears that most of them are predicted wrongly due to a few reasons:

1.  The difference between romantic and non-romantic love
2.  Mixed feelings
3.  Misunderstood context
4.  Generally confusing such as the following example:

index: 1430

True Label love - Predicted Label joy

Original:  اتحاد الشرس  ..الاربعاء ديربي &amp الاربعاء شباب الاتحاد بالتوفيق 7:30 الهلال الساعه

Cleaned:  الاربعاء ديربي اتحاد الشرس الهلال الساعه بالتوفيق شباب الاتحاد

(true = surprise, pred = none, total error = 35)

All of them are due to "الاوليمبياد"

ALL OF THEM!!!