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
```
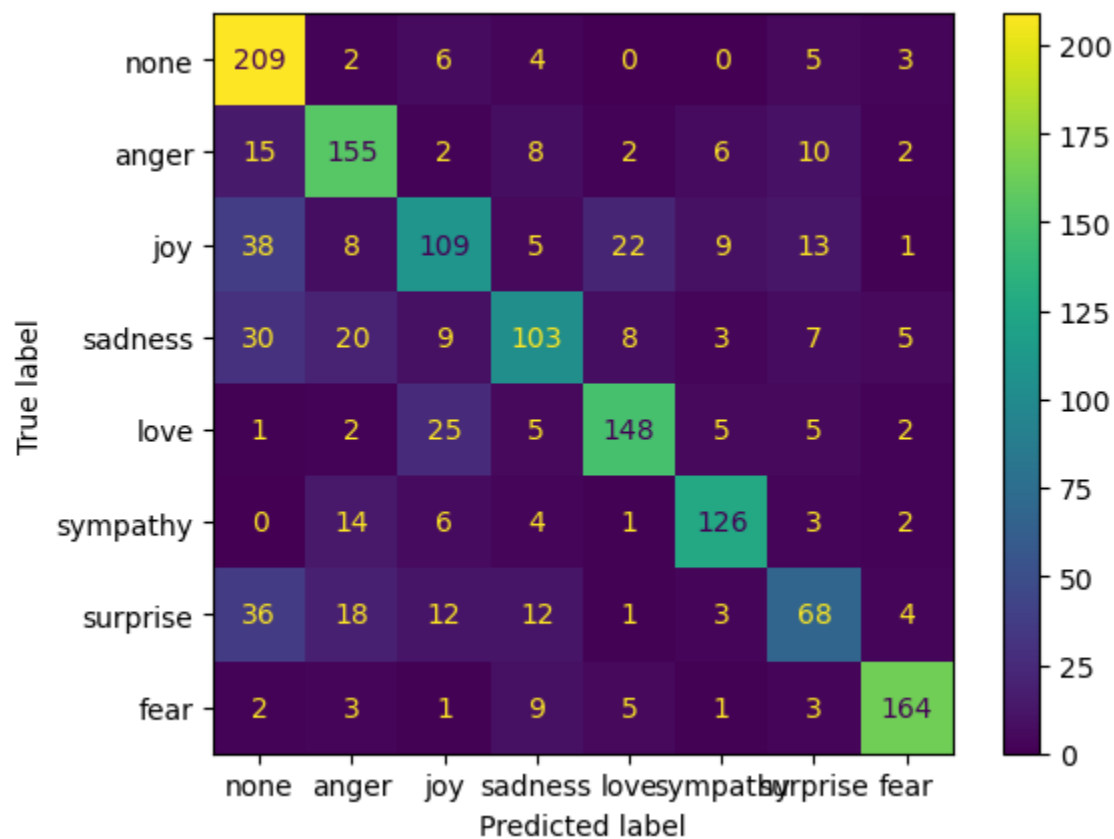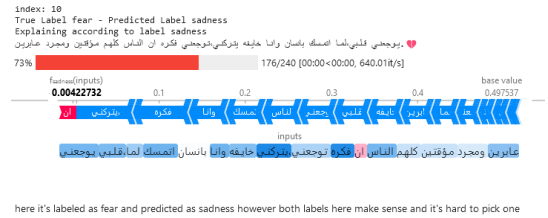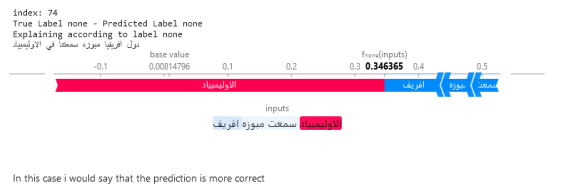
# Light Stemming vs Raw Data:



Here the model made the correct prediction with light stemming.



The model stopped using stopwords as explanations.

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:



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

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



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   0.1   0.3   0.5   0.7   f_fear(inputs)
0.877994

inputs

حزن عابر ومجرد مؤق الناس فكره يتركني،توجع خايفه وانا باسان اتمسك قليل يوجع

**Prediction probabilities**

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

NOT fear          fear

خايفه 0.31
ومجرد 0.08
يوجع 0.06
فكره 0.06
اتمسك 0.04
عابر 0.04
باس 0.03
وانا 0.03
الناس 0.01
يتركني،توجع 0.01

Text with highlighted words

يوجع قليل اتمسك باس وانا خايفه يتركني،توجع فكره الناس مؤق ومجرد عابر حزن

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
0.2   0.3   0.315782   0.4   f_joy(inputs)   0.5   0.6
0.487558

inputs

مشرف تمثيل الاوليمبياد داخل المصر

**Prediction probabilities**

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

NOT joy          joy

مشرف 0.29
داخل 0.22
الاوليمبياد 0.14
تمثيل 0.10
المصر 0.05

Text with highlighted words

المصر داخل الاوليمبياد تمثيل مشرف

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

---

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

Mortada د؛ الفاجر التحكيمي الظلم الم،!!!!!حرا

base value
0   0.0580826   0.1   0.2   0.3   f_anger(inputs)   0.4   0.5
0.378002

inputs

الفاجر التحكيمي الظلم حرام

**Prediction probabilities**

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

NOT anger          anger

الفاجر 0.25
التحكيمي 0.10
حرام 0.05
الظلم 0.00

Text with highlighted words

حرام الظلم التحكيمي الفاجر

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 اخطبه من ابوه بس خايفه يرفض 😊 😊



base value: 0.0148965

f_love(inputs): 0.450807

inputs: يرفض خايفه ابوه اخطبه ودي اروح اخطبه 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
اخطبه 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 اخطبه من ابوه خايقه 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 بـ السمنه علي كثر ما يسوون تكميم



base value: 0.0374464

f_anger(inputs): 0.475184

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
السعوديه 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

```
explain_example(test1)
```

2)



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

base value: 0.0374464

f_sad(inputs): 0.367688

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 0.02

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 دہ الفاجر التحكيمي الظلم حراااااام
```
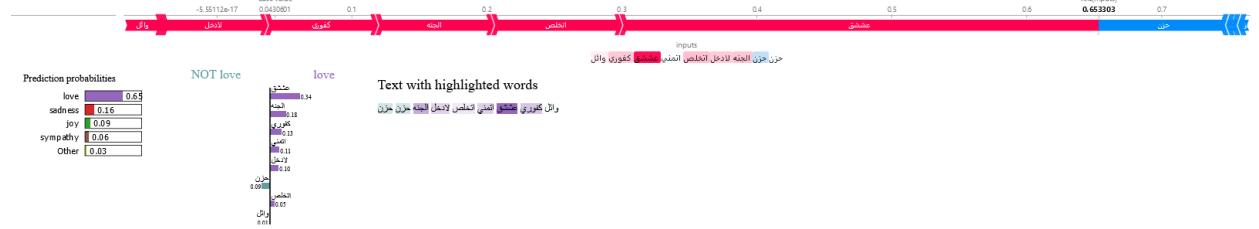
base value: 0.0374464

f_sad(inputs): 0.367688

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 عشششششق قلبي انتمي التخلص منه لادخل الجنه حزن ♥ 😍 🏠 Admin 👤
```

base value: 0.0430801

f_love(inputs): 0.653303

وائل    لادخل    كفوري    الجنه    التخلص    عششق

inputs
حزن حزن الجنه لادخل التخلص انتمي عششق كفوري وائل

Prediction probabilities
| love | 0.65 |
| sadness | 0.16 |
| joy | 0.09 |
| sympathy | 0.06 |
| Other | 0.03 |

NOT love    love

عششق 0.33
الجنه 0.18
كفوري 0.13
انتمي 0.11
لادخل 0.09
حزن 0.09
التخلص 0.05
وائل 0.02

Text with highlighted words
وائل كفوري عششق انتمي التخلص لادخل الجنه حزن حزن

```
explain_example(testi)
```

```
index: 15
True Label love - Predicted Label love
Explaining according to label love
وائل كفوري عششششق قلبي انتمي التخلص منه لادخل الجنه حزن ♥ 😍 🏠 Admin 👤
```

base value: 0.0430801

f_love(inputs): 0.653303

وائل    لادخل    كفوري    الجنه    التخلص    عششق

inputs
حزن حزن الجنه لادخل التخلص انتمي عششق كفوري وائل

Prediction probabilities
| love | 0.65 |
| sadness | 0.16 |
| joy | 0.09 |
| sympathy | 0.06 |
| Other | 0.03 |

NOT love    love

عششق 0.34
الجنه 0.18
كفوري 0.13
انتمي 0.11
لادخل 0.10
حزن 0.09
التخلص 0.05
وائل 0.01

Text with highlighted words
وائل كفوري عششق انتمي التخلص لادخل الجنه حزن حزن

4)

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

base value        f_fear(inputs)
0   0.0393674        0.2        0.4        0.6        0.8 **0.830301**

منك أكون / عنا   ايديك   ايدي   ياحبيبي      خايف       لامسه

inputs
بيك بحلم ياحبيبي اكون خايف ايديك لامسه ايدي معقول مش

**Prediction probabilities**    NOT fear    fear

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

خايف   0.71
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 0.01

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

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

base value        f_fear(inputs)
0   0.0393674        0.2        0.4        0.6        0.8 **0.830301**

منك أكون / عنا   ايديك   ايدي   ياحبيبي      خايف       لامسه

inputs
بيك بحلم ياحبيبي اكون خايف ايديك لامسه ايدي معقول مش

**Prediction probabilities**    NOT fear    fear

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

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

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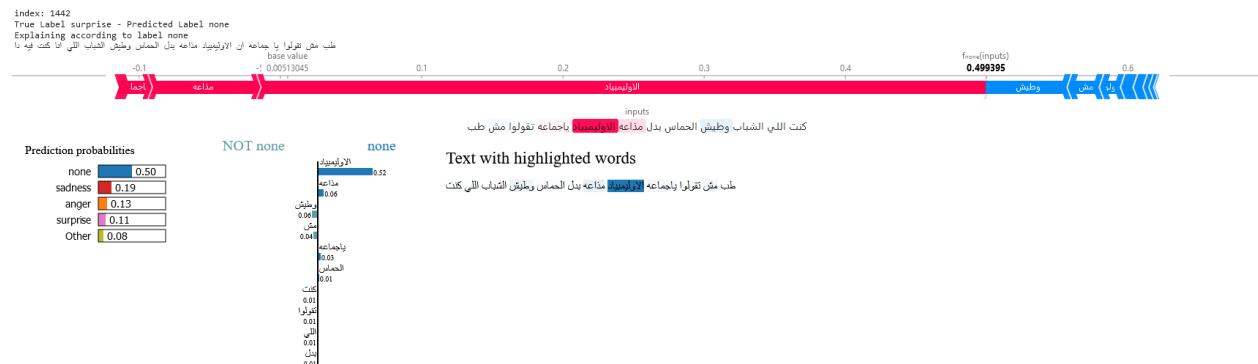مش معقول ايدي لامسه ايديك خايف اكون ياحبيبي بحلم بيك

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

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

جما   مذاعه      الاولمبياد       وطيش   مش   بل

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

**Prediction probabilities**    NOT none    none

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

الاولمبياد   0.52
مذاعه 0.06
وطيش 0.06
مش 0.04
ياجماعه 0.03
الحماس 0.01
كنت 0.01
تقولوا 0.01
اللي 0.01
بدل 0.01

**Text with highlighted words**
طب مش تقولوا ياجماعه الاولمبياد مذاعه بدل الحماس وطيش الشباب اللي كنت

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 با

base value
-0.1    -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 الاولمبياد

Here the order of significance is different.

index: 233
True Label surprise - Predicted Label none
Explaining according to label none
!!!!هو حماده طلعت لو كان لقي علم داعش (بدلا من علم السعوديه) و مكتوب عليه الشهادتين واقع على الارض كان شاله و لوح بيه في الاولمبياد برضه؟؟

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 | غا | ة | بيه | برضه | الاولمبياد | داعش | ع

inputs
برضه الاولمبياد بيه ولوح شاله الارض 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

index: 233
True Label surprise - Predicted Label none
Explaining according to label none
!!!!هو حماده طلعت لو كان لقي علم داعش (بدلا من علم السعوديه) و مكتوب عليه الشهادتين واقع على الارض كان شاله و لوح بيه في 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 | غا | ة | بيه | برضه | الاولمبياد | داعش | ع

inputs
برضه الاولمبياد بيه ولوح شاله الارض 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.

index: 82
True Label none - Predicted Label none
Explaining according to label none
نول عيال طربه بس كان جيل فشيخ كان ممكن يعمل حاجه في الاولمبياد

base value
-0.1    0.00513045    0.1    0.2    0.3    0.4    0.5    f none(inputs) 0.591645    0.7

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

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 طرزيه بس كان جيل فشيخ كان ممكن يعمل حاجه في 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 :
```

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 :
```

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.

Observe this 21-words example.

index: 697
True Label anger - Predicted Label sympathy
Explaining according to label sympathy
البعض لم يجد له موضوع يغرد بشانه سوي الدعاء علي العالم الكيمائي احمد زويل رحمه الله بزعم انه صهيوني! واتسائل هل من صنع ( الفورد ) مسلم ؟؟
Error displaying widget: model not found

Partition explainer: 2it [00:11, 11.21s/it]

**Prediction probabilities**

| | |
|---|---|
| sympathy | 0.38 |
| sadness | 0.26 |
| anger | 0.14 |
| surprise | 0.10 |
| Other | 0.12 |

NOT sympathy          sympathy

Text with highlighted words

2)



index: 117
True Label love - Predicted Label love
Explaining according to label love
صوتك يضحك بزوايا البيت و يضوّي ايامي ♥ 🔪 بوح كلمات جرح البحرين الكويت السعوديه الخليج الصوت الجريح اخبار رساله مع زجاجه عطر
Error displaying widget: model not found

Partition explainer: 2it [00:10, 10.35s/it]

**Prediction probabilities**

| | |
|---|---|
| love | 0.80 |
| sadness | 0.10 |
| surprise | 0.04 |
| joy | 0.03 |
| Other | 0.04 |

NOT love          love

Text with highlighted words

index: 117
True Label love - Predicted Label love
Explaining according to label love
صوتك يضحك بزوايا البيت    و يضوّي ايامي ♥ 🎀 بوح كلمات جرح البحرين الكويت السعوديه الخليج الصوت الجريح رساله مع زجاجه عطر
Error displaying widget: model not found

Partition explainer: 2it [00:10, 10.39s/it]

base value                                                    f love(inputs)
0.0197016          0.2          0.4          0.6          0.801182          1

عطر زجاجه رساله اخبار الجريح الصوت الخليج السعوديه الكويت البحر جرح كلم بوح حب ايامي ويضوي البيت بزوا يضحك صوتك

Prediction probabilities          NOT love          love          Text with highlighted words

love        0.80
sadness     0.10
surprise    0.04
joy         0.03
Other       0.04

صوتك يضحك بزوا البيت ويضوي ايامي حب بوح كلم جرح البحر الكويت السعوديه الخليج الصوت الجريح اخبار رساله زجاجه عطر

حب 0.15
عطر 0.15
بوح 0.12
صوتك 0.10
ايامي 0.07
ويضوي 0.05
الجريح 0.05
اخبار 0.04
بزوا 0.04
زجاجه 0.03

3)

index: 838
True Label love - Predicted Label love
Explaining according to label love
كنت احب وكنت اخاف من الغياب صرت اخاف اني احب وصرت اغيب عياب بوح عتاب احب وله جرح مشتاق احبك اشتياق
Error displaying widget: model not found

Partition explainer: 2it [00:11, 11.62s/it]

base value                                                    f love(inputs)
-0.3     -0.1     0.0197016     0.1     0.3     0.5     0.674428     0.9

اشتياق احبك مشتاق جرح فراق حزن خواطر بوح عتاب احب اني اخاف وصرت اغيب صرت الغياب اخاف وكنت احب كنت

inputs

Prediction probabilities          NOT love          love          Text with highlighted words

love        0.67
fear        0.25
sadness     0.07
sympathy    0.00
Other       0.00

كنت احب وكنت اخاف الغياب صرت اخاف اني احب عتاب بوح خواطر حزن فراق جرح مشتاق احبك اشتياق

احب 0.24
اخاف 0.20
احبك 0.20
بوح 0.07
خواطر 0.06
مشتاق 0.06
اشتياق 0.04
حزن 0.03
الغياب 0.03
صرت 0.03

index: 838
True Label love - Predicted Label love
Explaining according to label love
كنت احب وكنت اخاف من الغياب صرت اخاف اني احب وصرت اغيب عياب بوح عتاب احب وله جرح مشتاق احبك اشتياق
Error displaying widget: model not found

Partition explainer: 2it [00:11, 11.36s/it]

base value                                                    f love(inputs)
-0.3     -0.1     0.0197016     0.1     0.3     0.5     0.674428     0.9

اشتياق احبك مشتاق جرح فراق حزن خواطر بوح عتاب احب اني اخاف وصرت اغيب صرت الغياب اخاف وكنت احب كنت

inputs

Prediction probabilities          NOT love          love          Text with highlighted words

love        0.67
fear        0.25
sadness     0.07
sympathy    0.00
Other       0.00

كنت احب وكنت اخاف الغياب صرت اخاف اني احب عتاب بوح خواطر حزن فراق جرح مشتاق احبك اشتياق

احب 0.24
اخاف 0.22
احبك 0.20
بوح 0.07
خواطر 0.07
مشتاق 0.05
اشتياق 0.04
حزن 0.03
صرت 0.03
كنت 0.03

4)

اه لو نبطل نفسه من الي كان فقير وربنا كرمه او من الي معاه فلوس او اي حاجه نفسنا فيها ومش قادرين نجيبها    حب لاخيك ما تحب لنفسك

Error displaying widget: model not found
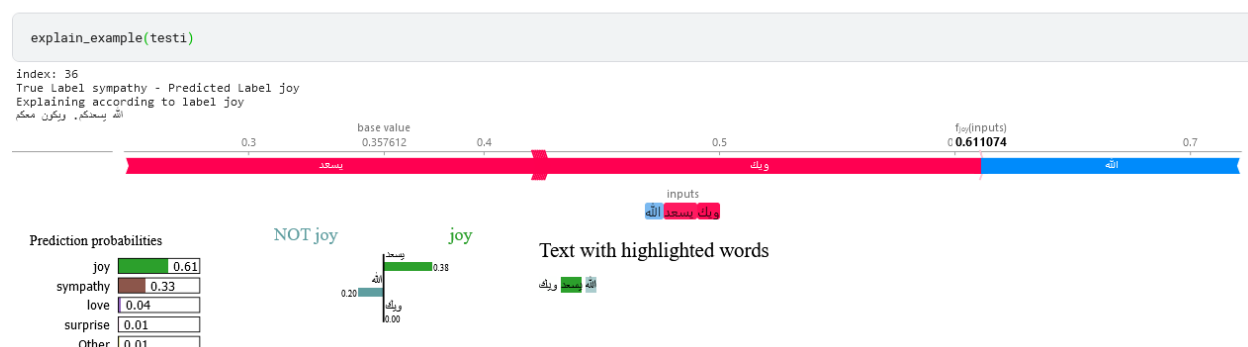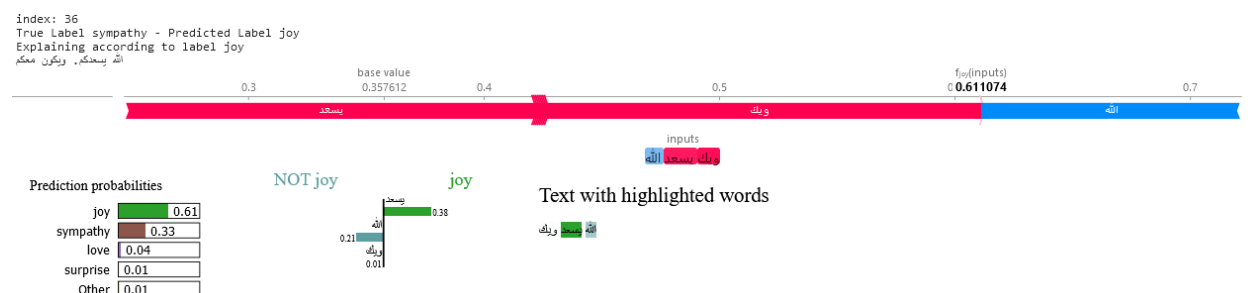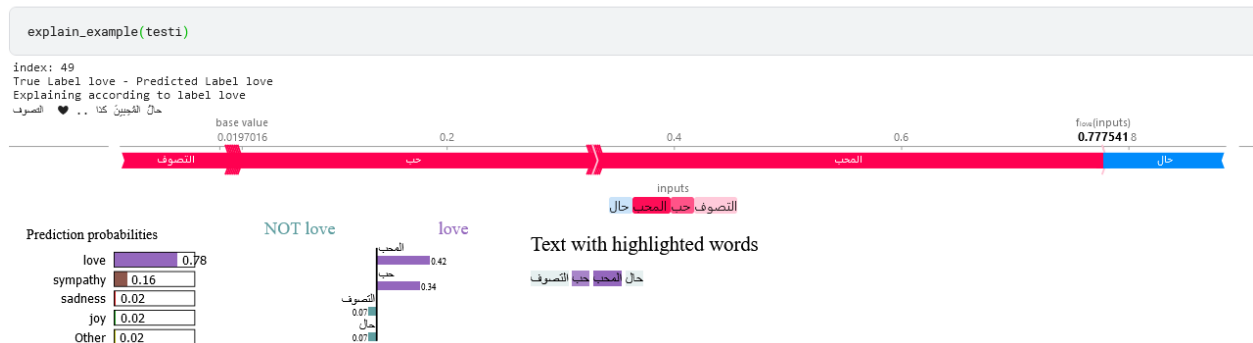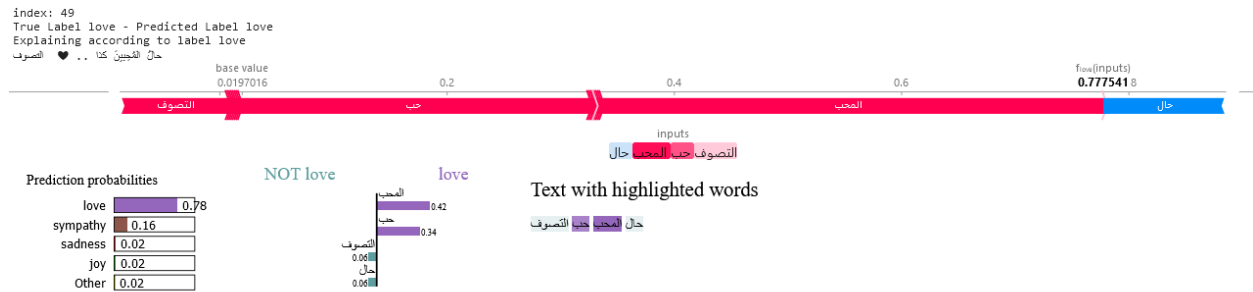
Partition explainer: 2it [00:10, 10.38s/it]



Prediction probabilities

sympathy   0.44
anger      0.28
sadness    0.12
love       0.05
Other      0.11

NOT sympathy    sympathy

Text with highlighted words

اه نبطل نفسنه الي فقير ورب كرمه او الي معاه فلوس او اي حاجه ومش قادر نجيب حب لاخيك تحب

فقير 0.20
نفسنه 0.15
قادر 0.11
الي 0.09
كرمه 0.09
ومش 0.08
لاخيك 0.08
حب 0.07
ورب 0.07
نبطل 0.07

اه لو نبطل نفسه من الي كان فقير وربنا كرمه او من الي معاه فلوس او اي حاجه نفسنا فيها ومش قادرين نجيبها    حب لاخيك ما تحب لنفسك

Error displaying widget: model not found

Partition explainer: 2it [00:10, 10.72s/it]



Prediction probabilities

sympathy   0.44
anger      0.28
sadness    0.12
love       0.05
Other      0.11

NOT sympathy    sympathy

Text with highlighted words

اه نبطل نفسنه الي فقير ورب كرمه او الي معاه فلوس او اي حاجه ومش قادر نجيب حب لاخيك تحب

فقير 0.18
نفسنه 0.15
قادر 0.10
الي 0.09
كرمه 0.08
ومش 0.08
لاخيك 0.08
ورب 0.06
حب 0.06
معاه 0.05

# Let's observe small examples now:

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



```
base value                                                    f_anger(inputs)
0    0.0576145    0.1         0.2         0.3    0.370942    0.4         0.5
```

Prediction probabilities

| | |
|---|---|
| anger | 0.37 |
| surprise | 0.23 |
| sadness | 0.23 |
| fear | 0.12 |
| Other | 0.05 |

NOT anger          anger

inputs
الفاجر التحكيمي الظلم حرام

Text with highlighted words

حرام الظلم التحكيمي الفاجر

الفاجر 0.24
التحكيمي 0.10
حرام 0.04
الظلم 0.00

```
explain_example(testi)
```

## 2)

```
index: 36
True Label sympathy - Predicted Label joy
Explaining according to label joy
الله يسعدكم. ويكون معكم
```



```
base value                                                    f_joy(inputs)
0.3         0.357612    0.4              0.5    0.611074    0.7
```

Prediction probabilities

| | |
|---|---|
| joy | 0.61 |
| sympathy | 0.33 |
| love | 0.04 |
| surprise | 0.01 |
| Other | 0.01 |

NOT joy          joy

inputs
ويك يسعد الله

Text with highlighted words

الله يسعد ويك

يسعد 0.38
الله 0.21
ويك 0.01

```
explain_example(testi)
```

## 3)

index: 49
True Label love - Predicted Label love
Explaining according to label love

حالُ المُحِبِّ كذا .. ♥ التصوف

base value
0.0197016

0.2

0.4

0.6

f_love(inputs)
0.7775418

التصوف حب المحب حال

inputs
التصوف حب المحب حال

Prediction probabilities

love 0.78
sympathy 0.16
sadness 0.02
joy 0.02
Other 0.02

NOT love        love

المحب 0.42
حب 0.34
التصوف 0.06
حال 0.06

Text with highlighted words

حال المحب حب التصوف

```
explain_example(testi)
```

index: 49
True Label love - Predicted Label love
Explaining according to label love

حالُ المُحِبِّ كذا .. ♥ التصوف

base value
0.0197016

0.2

0.4

0.6

f_love(inputs)
0.7775418

التصوف حب المحب حال

inputs
التصوف حب المحب حال

Prediction probabilities

love 0.78
sympathy 0.16
sadness 0.02
joy 0.02
Other 0.02

NOT love        love

المحب 0.42
حب 0.34
التصوف 0.07
حال 0.07

Text with highlighted words

حال المحب حب التصوف

4)

index: 104
True Label fear - Predicted Label fear
Explaining according to label fear

جمهور مرعب

base value
-5.55112e-17  0.0338351

0.1

0.2

0.3

0.4

0.5

f_fear(inputs)
0.541999

0.6

مرعب

جمهور

inputs
مرعب جمهور

Prediction probabilities

fear 0.54
surprise 0.23
joy 0.16
sympathy 0.03
Other 0.05

NOT fear        fear

مرعب 0.34
جمهور 0.09

Text with highlighted words

جمهور مرعب

+ Code    + Markdown

```
explain_example(testi)
```

index: 104
True Label fear - Predicted Label fear
Explaining according to label fear

جمهور مرعب

base value
-5.55112e-17  0.0338351

0.1

0.2

0.3

0.4

0.5

f_fear(inputs)
0.541999

0.6

مرعب

جمهور

inputs
مرعب جمهور

Prediction probabilities

fear 0.54
surprise 0.23
joy 0.16
sympathy 0.03
Other 0.05

NOT fear        fear

مرعب 0.34
جمهور 0.09

Text with highlighted words

جمهور مرعب

It seems that lime is unable to explain it properly and It may be because Lime generates 5000 random samples and while those examples are enough to cover most sentences, it appears that it is not enough to cover all permutations possible in a long sentence so, the changes between different runs are a bit higher.

The highest change in short examples tend to be 0.01 but in long sentences it can reach 0.02 of change in multiple words at once.

**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:

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
(:) الاولمبياد كانت عباره عن باص من تريكه يوصل صلاح بالجول وجول البرازيل وشاره الكابتن .. وحشتنا ياعم

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

base value 0.248209
f_joy(inputs) 0.463062

inputs

فرح ياعم وحشتنا الكابتن وشاره البرازيل وجول بالجول صلاح يوصل تريكه باص عباره كانت الاولمبياد

**Prediction probabilities**

| | |
|---|---|
| joy | 0.46 |
| none | 0.33 |
| surprise | 0.09 |
| sadness | 0.06 |
| Other | 0.06 |

NOT joy | joy

وحشتنا 0.20
الاولمبياد 0.18
تريكه 0.09
يوصل 0.07
البرازيل 0.07
فرح 0.04
صلاح 0.03
الكابتن 0.02
ياعم 0.01
وشاره 0.01

**Text with highlighted words**

الاولمبياد كانت عباره عن باص من تريكه يوصل صلاح بالجول وجول البرازيل وشاره الكابتن وحشتنا ياعم فرح

---

index: 100
True Label none - Predicted Label joy
Explaining according to label joy
شنو يعجبني بالكويت ماعمري 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 ماعمري 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
ماعمري 0.04
السيار 0.02
بالكويت 0.02
احمد 0.02
زرتها 0.01

**Text with highlighted words**

شنو يعجبك بالكويت ماعمري زرتها يعجبني احمد السيار

---

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

base value 0.248209
f_joy(inputs) 0.479533

ازر واخلاقهم اسم الاولمبياد باداهم يرفعوا مصر مصر شجعوا

inputs

الاولمبياد مصر شجعوا واخلاقهم باداهم 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
شجعوا 0.10
مصر 0.10
يرفعوا 0.06
باداهم 0.05
الاولمبياد 0.02
يارب 0.00

**Text with highlighted words**

يارب يرفعوا اسم مصر باداهم واخلاقهم شجعوا مصر الاولمبياد

---

index: 801
True Label none - Predicted Label joy
Explaining according to label joy
احتفالات الاولمبياد 1 احتفالات قناة 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 اعظم رياضي في التاريخ خارج التنس 28

base value 0.248209
f_joy(inputs) 0.45988

الاعظم واعظم دهب ميداليه خارج الاولمبياد التاريخ

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
واش لو استضفنا الاولمبياد وشغلنا مهرجانات هتبقي فكره كريتف

base value | f_joy(inputs)
0.248209 | 0.404496

-5.55112e-17   0.1   0.2   0.3   0.4   0.5   0.6   0.7

والله ← كريتف ← مهرجانات ← فكره → الاولمبياد → هتبقي → استضفنا → وشغلنا

inputs
كريتف فكره هتبقي مهرجانات وشغلنا الاولمبياد استضفنا والله

Prediction probabilities
NOT joy    joy
joy     0.40
none    0.34
sadness 0.12
anger   0.07
Other   0.06

الاولمبياد 0.22
مهرجانات 0.16
فكره 0.08
وشغلنا 0.08
استضفنا 0.07
والله 0.06
هتبقي
كريتف 0.04

Text with highlighted words
واش استضفنا الاولمبياد وشغلنا هتبقي فكره كريتف

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

base value | f_joy(inputs)
0.248209 | 0.673445

0.2   0.3   0.4   0.5   0.6   0.7

واحد ← الاولمبياد ← فرح ← فرح → الخامس → مصر → مركز

inputs
فرح فرح الاولمبياد الخامس المركز واحد رقم مصر

Prediction probabilities
NOT joy    joy
joy      0.67
none     0.19
surprise 0.11
sadness  0.02
Other    0.01

فرح 0.39
الاولمبياد 0.12
المركز
مصر
واحد 0.06
واحد 0.05
رقم 0.03
الخامس 0.00

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

index: 917
True Label none - Predicted Label joy
Explaining according to label joy
Rio2016 التجنيس ده حلو يا فخري حلو اوي حلو يجيب ميداليات في الاولمبياد اجدع

base value | f_joy(inputs)
0.248209 | 0.620837

0.1   0.2   0.3   0.4   0.5   0.6   0.7   0.8

التجنيس ← اوي ← يافخري ← حلو ← حلو → الاولمبياد → ياجيب → يجيب → ميداليات

inputs
اجدع الاولمبياد ميداليات يجيب ياجدع اوي حلو يافخري حلو التجنيس

Prediction probabilities
NOT joy    joy
joy      0.62
surprise 0.19
none     0.14
anger    0.03
Other    0.02

حلو 0.40
يافخري 0.12
يجيب 0.07
ياجدع 0.06
التجنيس 0.05
اوي 0.04
الاولمبياد 0.04
ميداليات 0.04
اجدع 0.00

Text with highlighted words
التجنيس حلو يافخري حلو اوي ياجدع يجيب ميداليات الاولمبياد اجدع

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

base value | f_joy(inputs)
0.248209 | 0.450508

0.1   0.2   0.3   0.4   0.5   0.6

عزاء ← سودانا ← للمتشائمين ← العاشره → الاولمبياد → اما → اولادنا → حقتنا

inputs
للمتشائمين عزاء اولادنا والاولاد حقتنا الاولمبياد اما سودانا العاشره

Prediction probabilities
NOT joy    joy
joy      0.45
none     0.31
sympathy 0.12
sadness  0.07
Other    0.05

عزاء 0.10
الاولمبياد 0.08
للمتشائمين 0.08
سودانا 0.07
اما 0.07
العاشره 0.06
حقتنا 0.06
اولادنا 0.03
والاولاد 0.01

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

index: 1315
True Label none - Predicted Label joy
Explaining according to label joy
🔢 .. احنا عندنا ف الاولمبياد ناس زي الذهب

base value | f_joy(inputs)
0.248209 | 0.502388

0.2   0.3   0.4   0.5   0.6

زي ← الذهب → الاولمبياد

inputs
الذهب زي ناس الاولمبياد

Prediction probabilities
NOT joy    joy
joy      0.50
none     0.27
surprise 0.14
sadness  0.05
Other    0.03

الذهب 0.26
زي 0.21
الاولمبياد 0.13
ناس 0.06

Text with highlighted words
🔢 الاولمبياد 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.

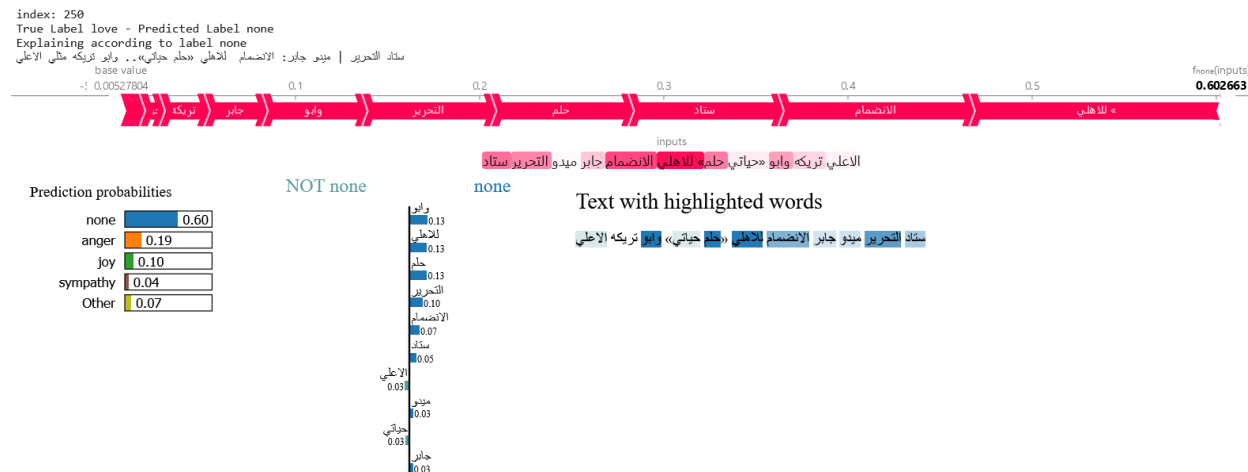So, let's examine something a bit less confusing.

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

index: 1197
True Label none - Predicted Label anger
Explaining according to label anger
اصلا ضغط الاولمبياد كفايه علي الولد حرام تحمله زياده بدل ما ينزل مش عارف يلعب من اساسه



index: 1317
True Label none - Predicted Label anger
Explaining according to label anger
وزاره التجاره تدعو الشركات المتاثره بارتفاع البنزين لتقديم ادلتها مشدده علي ضروره تطبيق القانون المخالفين الذين يرفعون الاسعاربشكل مصطنع

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

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 طروش ... اما انت حمار صحيح 🙂

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

-0.4    -0.2    base value 0 0.0391545    0.2    0.4   f<sub>anger</sub>(inputs) 0.450145    0.6    0.8

inputs

**Prediction probabilities**

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

NOT anger     anger

**Text with highlighted words**

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

ستاد التحرير | ميدو جابر: الانضمام للاهلي «حلم حياتي».. وابو تريكه مثلي الاعلي

base value -0.00527804    0.1    0.2    0.3    0.4    0.5   f<sub>none</sub>(inputs) 0.602663

inputs

**Prediction probabilities**

| | |
|---|---|
| none | 0.60 |
| anger | 0.19 |
| joy | 0.10 |
| sympathy | 0.04 |
| Other | 0.07 |

NOT none     none

**Text with highlighted words**

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: كفو عيال ؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟؟ديرتييييي وبنت ديرتيي قبيلتي بنت مبروووووك الف الف الكويت فخر العنزي احلام
الديره كل يوم 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 يلاد :
المسلمين ♡

Cleaned: الدوله العثمانيه الاسلاميه ستبقي رغم انوف الاعداء ربي احفظ تركيا المسلمه مكر الماكرين واحفظ يلاد 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: الهلال الساعه 7:30 بالتوفيق شباب الاتحاد &amp الاربعاء ديربي.. اتحاد الشرس

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

All of them are due to "الاولىمبياد"

ALL OF THEM!!!