# Interpretation of Models
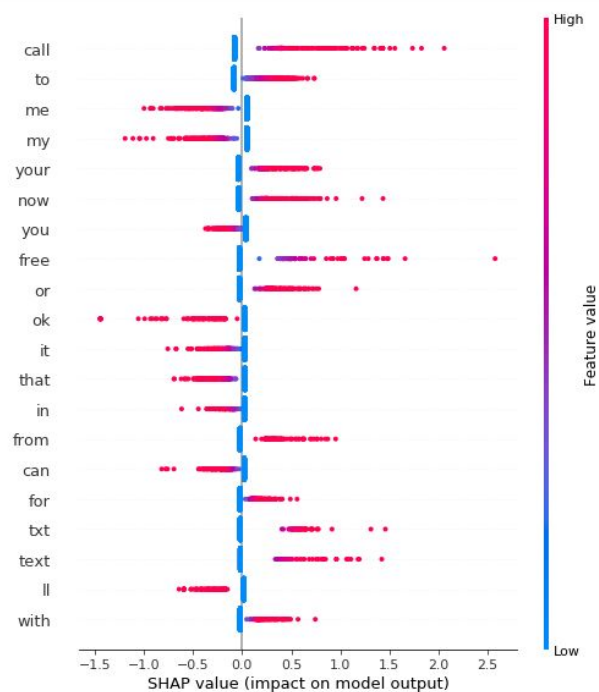
## 1 . Shap (SHapley Additive exPlanations) :

It illustrates the model dependencies with the help of values called **shapley_values.** These values describe the feature's impact on the model. (either collectively or Independently).

**Pros:**

- **Summary_plot :** Shap has very informative plot called summary plot which gives the overall picture of features and their impact on model

```
In [85]: shap.summary_plot(shap_values, X_test_tfid.toarray(), feature_names=tf_id_trans.get_feature_names())
```



The words `call` `to` `me` `my` `your` ... has the highest SHAP values. That means they have the highest impact on model to choose between spam ar ham

- The shap values generated help to interpret the values easily and identify the cost of impact
- Force plot provides the flexibility for the user to plot the feature impact on models for required predictions

```
In [86]: print(X_test[id])
         shap.force_plot(
             explainer.expected_value, shap_values[idx,:], X_test_array[idx,:],
             feature_names=tf_id_trans.get_feature_names()
         )
```

Congratulations ur awarded either å£500 of CD gift vouchers & Free entry 2 our å£100 weekly draw txt MUSIC to 87066 TnCs www.Ldew.com1win150ppmx3age16

Out[86]:

- LinearExplainer is fast when compared to other Shap models
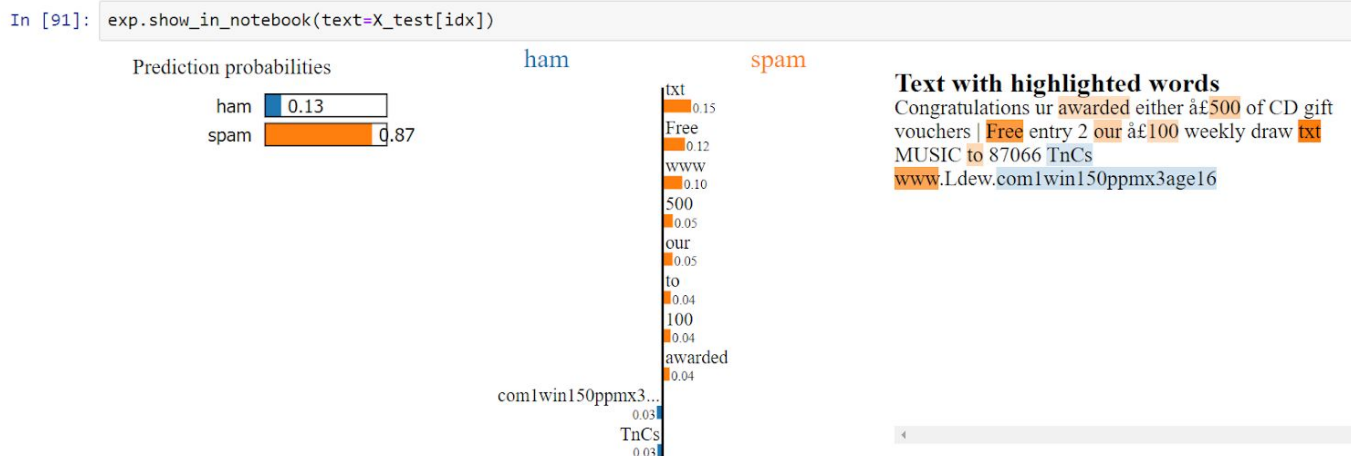
**Cons :**

- KernalShap is computationally very slow

## 2 . Lime (Local Interpretable Model-Agnostic Explanations) :

Lime takes a classifier and raw text or numpy array into a function and explains the reason behind predicting the output. At present lime works on the classifiers with two or more classes.

**Pros:**

- Lime works very well on text classifiers
- It is very simple to implement and very quick in response, even for huge datasets
- For textual data Lime highlights each word which is responsible for the prediction

**Cons :**

- Lime is not consistent when compared to Shap

## 3 . InterpretML :

InterpretML helps find the model's global behaviour and also individual predictions. It brings all the interpretable techniques such as Shap, Lime into one roof.

**Pros:**

- It provides a dashboard to play with around the features and its impact on model



- It brings all the techniques in to one place

**Cons :**

- It consumes more resources as it extracts all the possibilities of feature relations and its coefficients.


## 4 . ELI5:

It inspects model parameters based on weights and tries to figure out how the model works globally. It inspects an individual prediction of a model, and identifies the reason behind the prediction.

**Pros :**

- We can see the most weighted features which has the impact on the model globally



- We can view the individual prediction along with the highlighted words that are responsible for the prediction

**Cons :**

- Multi class availability is limited

```python
import eli5
```

```python
eli5.show_weights(text_clf.named_steps["clf"], vec=tf_id_trans, top=20)
```

Error: only binary libsvm-based classifiers are supported