# Feature engineering using Local Interpretable Model-agnostic Explanations (LIME)

## Recap

A binary classification model may be complex and non-linear, as shown by the blue/pink boundary in the Fig 1. below. A specific instance that needs to be explained is shown as a bright red X. LIME samples instances in the vicinity of this specific instance and their actual class is used to fit a simple linear classifier, without dealing with the complications of global fit (local fidelity). The key feature values of the specific instance, and their contribution to its classification are shown in Fig 2. Details about LIME are published here.
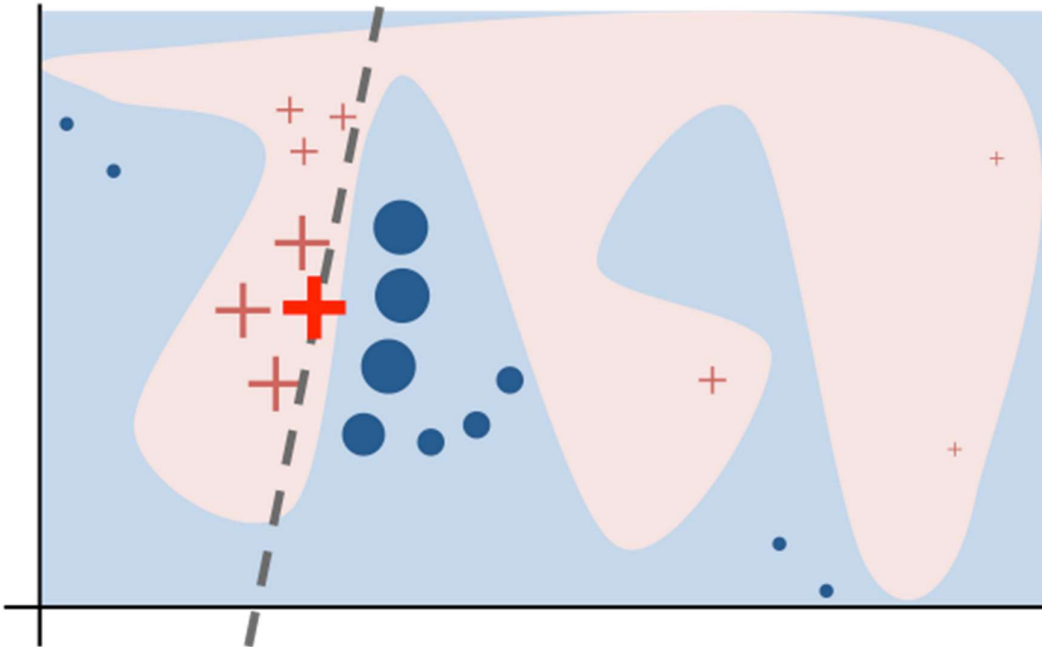


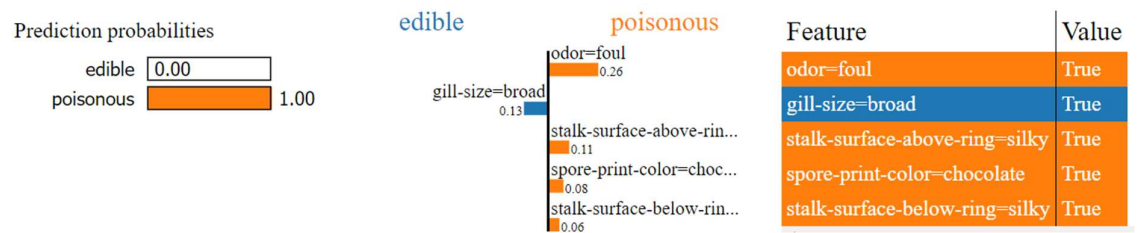Fig 1: Simple local linear classifier by LIME



Fig 2: Explaining a prediction as to why a mushroom is poisonous and not edible

While this local classifier can be used explain an individual prediction, it can also be used to improve features. Playing around with the LIME author's work, here are some observations shared as learnings and not as any original work.
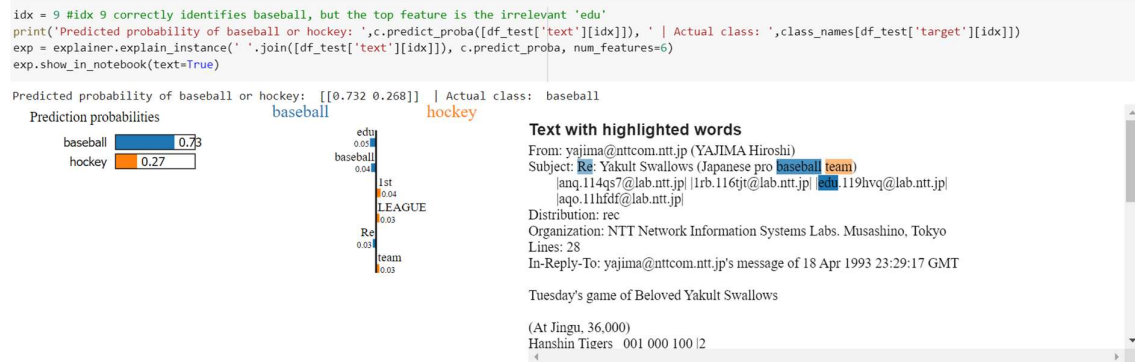
## Feature Engineering

Taking the example of a Random Forest classifier that is trying to label TF-IDF vectors of text related to baseball or hockey and going by its metrics, the model appears to perform quite well.

Code:
https://colab.research.google.com/drive/1LtOEWFJ0PUbvGJgPXgWTZLgYpiGnFp5B#scrollTo=lkPnhvrdE2ZG

```
sklearn.metrics.precision_recall_fscore_support(df_test['target'], pred, average='binary')
```

```
(0.9536784741144414, 0.8771929824561403, 0.9138381201044387, None)
```

But when we examine the features used, irrelevant details like the sender's name and mail id are used as features (see comments in green). Clearly, the model cannot be generalized to unseen text and classify correctly. This is an important insight that LIME provides.



A simple regex function written to strip out irrelevant details, improving features fed to the model.

```python
def remUninformative(text):
    import re
    next_text = re.sub(r'[A-Za-z0-9.]*@[A-Za-z0-9.]*', "", text)#email adress
    next_text = re.sub(r'^From:.*', "", next_text, flags=re.M)#From line completely
    next_text = re.sub(r'^Organization:.*', "", next_text, flags=re.M)#Organization line completely
    next_text = re.sub(r'Re:', "", next_text)#Just the tag
    next_text = re.sub(r'Subject:', "", next_text)#Just the tag
    return next_text
```

The retrained classifier provides greater confidence that it is using relevant features and that it can generalize better.



But the metrics haven't changed/ improved much.

```
sklearn.metrics.precision_recall_fscore_support(df_test['target'], pred, average='binary')
```

```
(0.9885057471264368, 0.8621553884711779, 0.9210174029451138, None)
```

## Conclusion

The similarity in model metrics before and after feature engineering, demonstrates the need to go beyond metrics to choose features. Using LIME, the impact of features used can be understood better, enabling improved feature engineering and model generalize-ability.

Code:
https://colab.research.google.com/drive/1LtOEWFJ0PUbvGJgPXgWTZLgYpiGnFp5B#scrollTo=lkPnhvrdE2ZG