

# **Towards a more transparent AI - Decrypting ML models using LIME**

Laisha Wadhwa

# Session Agenda



**Introduction**



**Importance of  
interpretable AI**



**LIME: What  
and How?**



**Hands on model  
interpretation**



**Final thoughts**



## About Me

 @laishawadhwa

 /laishawadhwa

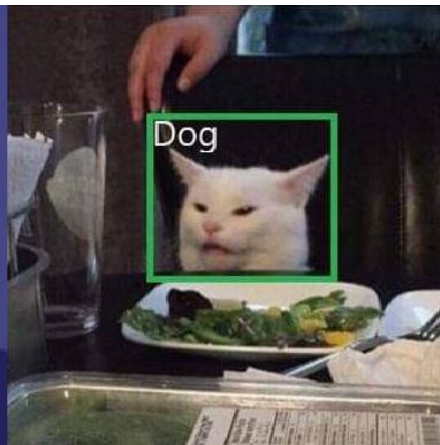
- Data Engineer, Couture.ai
- Microsoft AI challenge 2018 winner
- Amex & Techgig Geek Goddess 2019 AIML hackathon winner
- Sabre Hack 2019 Winner
- Mercedes Benz Digital Challenge winner
- Icertis AI and blockchain RU
- Women Tech Network, Global Ambassador
- Podcast Host: Co-Learning Lounge
- Technical Writer, Omdena
- Tech Speaker

# We all have been here!

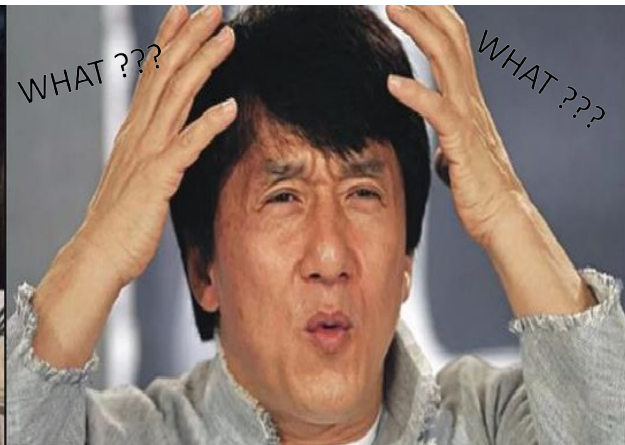
When your classifier  
achieves  
98% accuracy



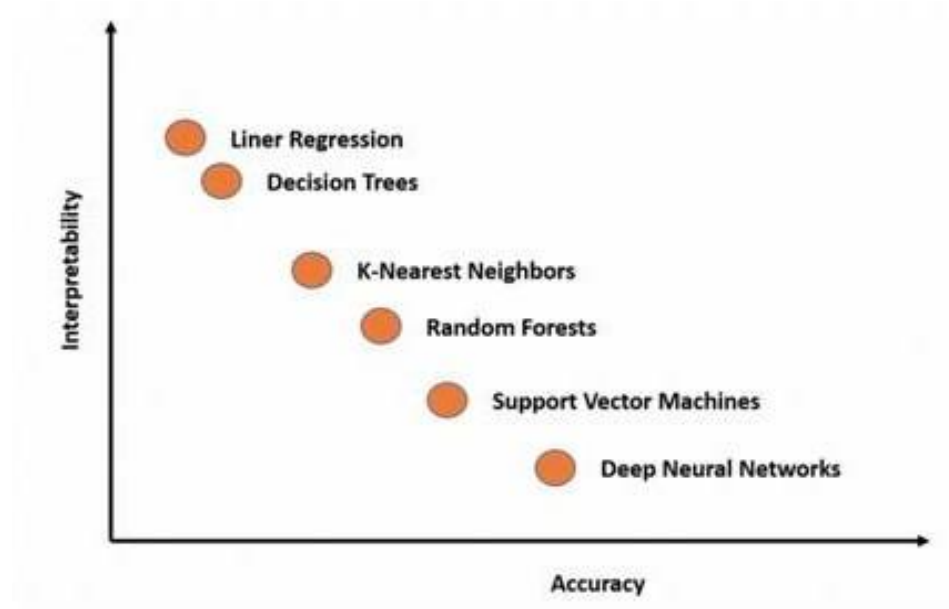
Your classifier's  
Output



YOU



# Accuracy vs Interpretability



# ML Algorithmic Trade-Off



- The higher the interpretability of a model, the easier it is for someone to comprehend why certain decisions (predictions) were made.
- A model has better interpretability than another model, if its decisions are easier to comprehend for a human than decisions from the second model.

## White Box

- Regression, Decision Trees, Association Rule Mining, Linear SVMs

## Grey Box

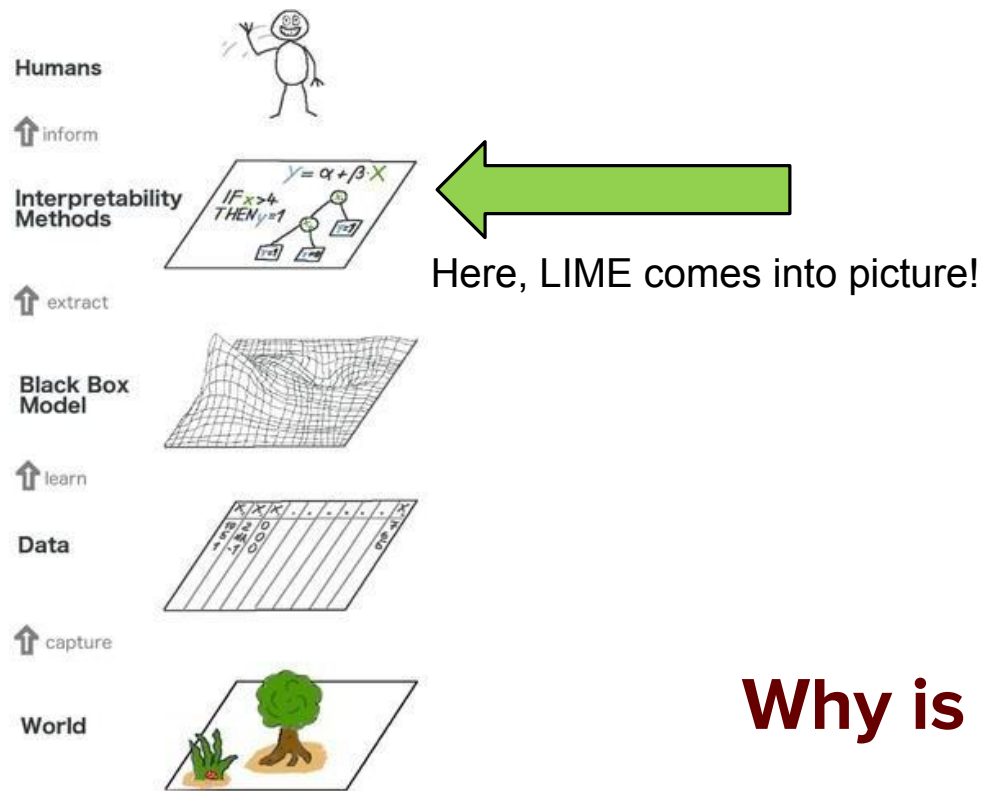
- Clustering, Bayesian Nets, Genetic algorithms, Logic programming

## Black Box

- DNNs, Non-linear matrix factorizations, Non-linear dimensionality reduction

Slide credits: **Bring in the Lawyers: Explainable AI Driven Decision-making by Vikas Agrawal at [ODSC India](#)**

# ML and DL models have become a black-box



*The idea is to interpret the models and understand why the classifier chose a particular class or why the model gives more False Positives ??*

## Why is it important?

# Creators of LIME



**Riberio Marco**



**Sameer Singh**



**Guestrin Carlos**



# What is LIME??

**L**ocal

**I**nterpretable

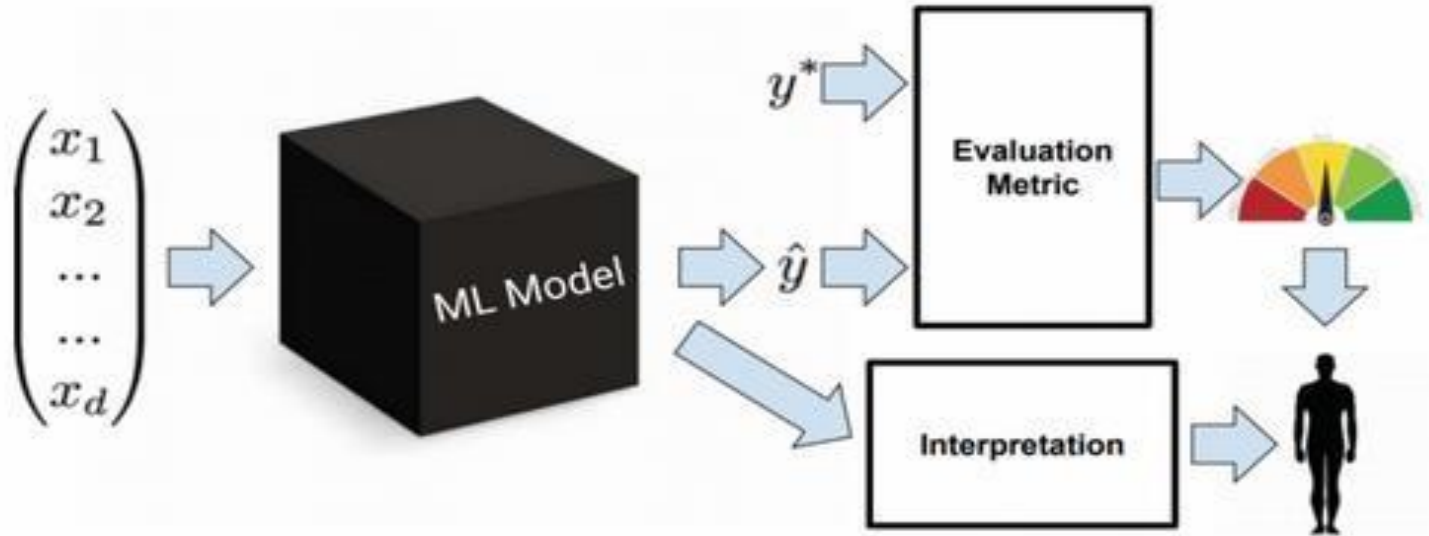
**M**odel agnostic

**E**xplanations



# LIME helps to interpret models, HOW?

---



# Interpretable Machine Learning with LIME

Machine Learning Model

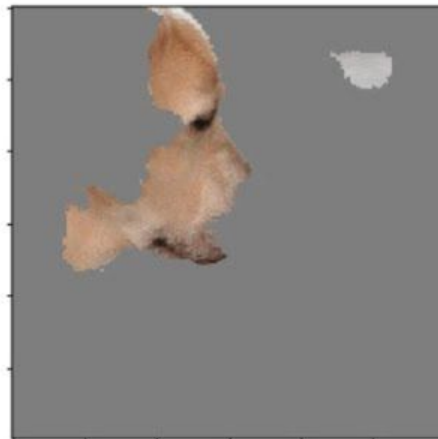
**This is a  
“labrador”**



**Why?**  
→

**LIME**

**Because:**



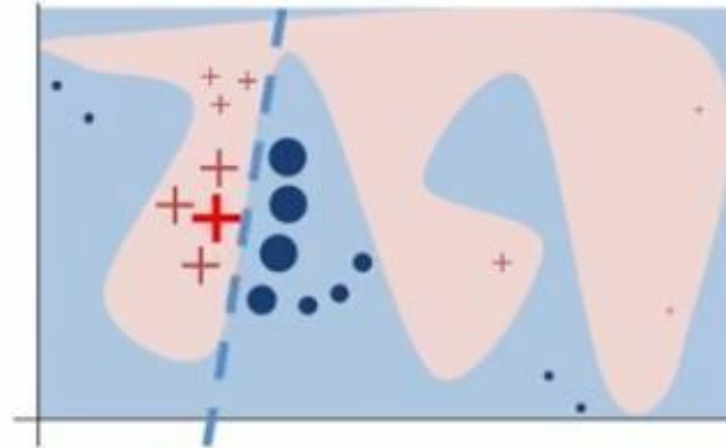
# The Secret Sauce



- LIME minimises the **locality-aware loss** which is measured using a mathematical formula (It's basic math!) which takes into consideration measure of **unfaithfulness of the model** in approximating the target variable and measure of model complexity ( $\rho$ ) of explanation.
- E.g. if the model to be explained is decision tree,  $\rho$  can be depth of the tree or in case of linear explanation models it can be number of non zero weights.

## Using LIME to explain a complex model's prediction for input $x_i$

1. Sample points around  $x_i$
2. Use complex model to predict labels for each sample
3. Weigh samples according to distance to  $x_i$
4. Learn new simple model on weighted samples
5. Use simple model to explain

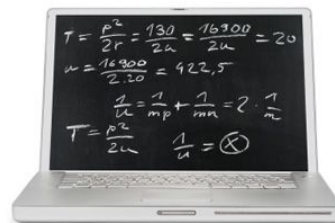


Ack: Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin

53

Slide credits: **Bring in the Lawyers: Explainable AI Driven Decision-making** by Vikas Agrawal at [ODSC India](#)

# The Math - locality-aware loss



$$\xi(x) = \operatorname{argmin}_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

**First Term** : measure of unfaithfulness of  $g$  in approximating  $f$  in the locality defined by  $\pi_x$ . This is termed as **locality-aware loss** in the original paper

**Last term** : measure of model complexity of explanation  $g$ . For example if your explanation model is decision tree it can be depth of the tree or in case of linear explanation models it can be number of non zero weights

Learn another model on top of the ML Model



Generate new data by permutation of the input data

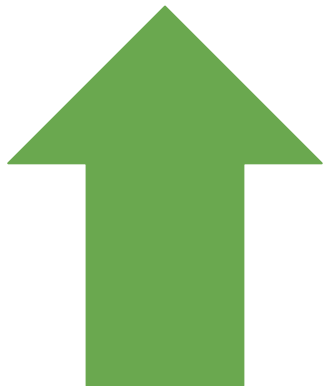


Prediction on this new data on the the simple (local fit: an approximation)



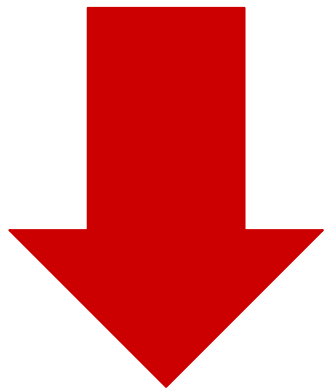
**More weight is being put  
on data that is similar to  
the original data)**

# Why should you use it?



## LIME

- > Interpretability.
- > Explains decision boundaries of the model in human understandable form.
- > Local Fidelity & Model Agnostic
- > Locally approximate global model

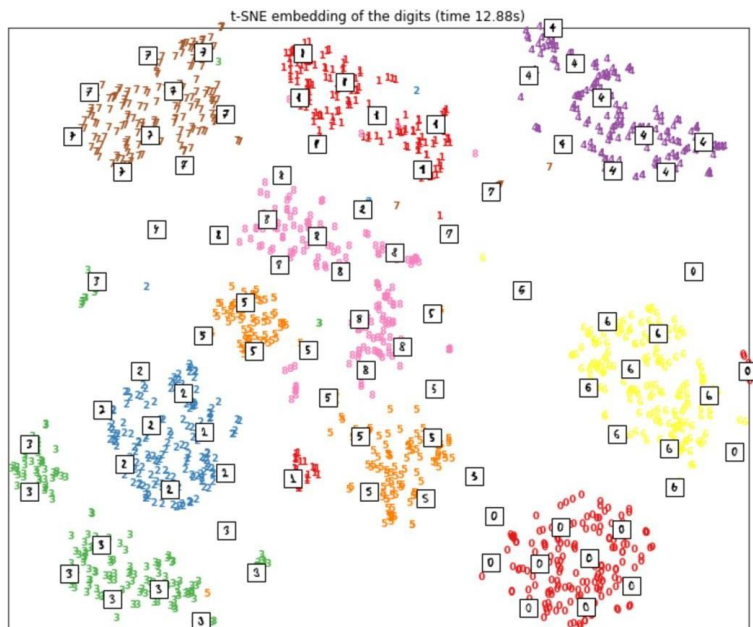


## Traditional Approach

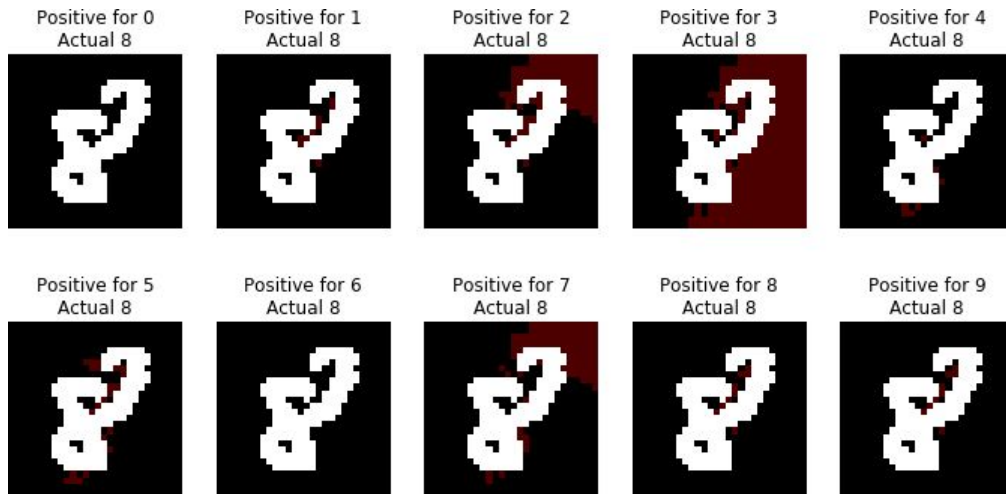
- > EDA: Visualisation [t-SNE
- > Model Performance Evaluation Metrics.
- > Can't comprehend the decision boundary with changing nature of input points.
- > No insight on feature imp for new data pts.



# Traditional



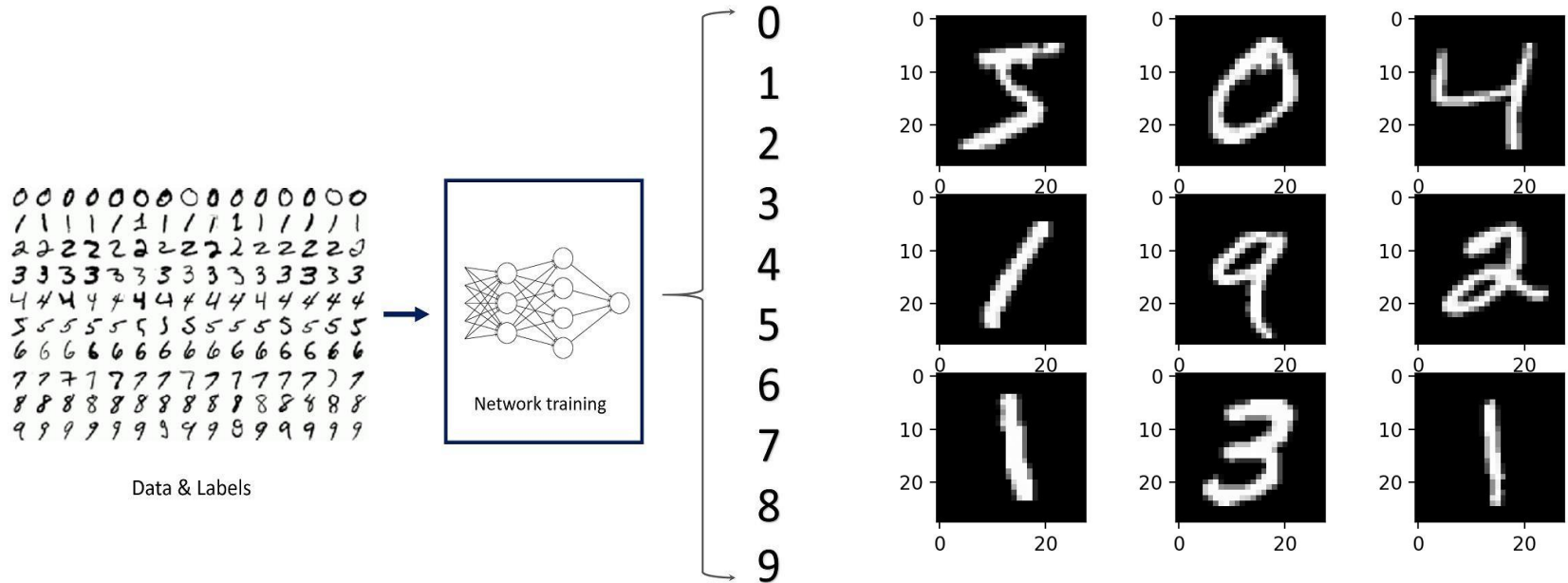
# LIME




# What all can LIME do?

- ❖ It can explain models with any kind of data:
  - a) Text : It represents presence/absence of words.  
**LimeTextExplainer**
  - b) Image : It represents presence/absence of super pixels ( contiguous patch of similar pixels ). **LimeImageExplainer**
  - c) Tabular data : It is a weighted combination of columns. **LimeTabularExplainer**

# DIGIT classification



# How do we do it??



```
simple_dt_pipeline = Pipeline([
    ('Make Gray', makegray_step),
    ('Flatten Image', flatten_step),
    ('DT', DecisionTreeClassifier(criterion = "gini",
                                  random_state = 100, max_depth=7, min_samples_leaf=5) )
])
```

# LIME: explain\_instance and ImageExplainer

[illegible]

# Let's see it in Action!

- LIME tests out what happens to your black box model's predictions when you feed variations or perturbations of your dataset.
- How? By generating a new set comprising of perturbed samples and the corresponding black box model's predictions.
- LIME then trains an interpretable model weighted by the proximity of the sampled instances to the instance of interest.

Let's see how well can a decision tree differentiate the numbers!!

# Let's delve deeper

Let's say we are trying to classify wolves and dogs.



Wolf



Dog

What if we have wolves and dogs in entirely different backgrounds? This is entirely possible as wolves are mostly found in the wild (in snow, jungles, etc.) while dogs are in completely different backgrounds (households) generally.



We observe that, our model is actually learning the background pretty well!!

Let's see this in action! Does it work well for DL?]



# Image classification- How to?



```
#Loading model
checkpoints_dir = '/content/tf-models/slim/pretrained'
init_fn = slim.assign_from_checkpoint_fn(
    os.path.join(checkpoints_dir, 'inception_v3.ckpt'),
    slim.get_model_variables('InceptionV3'))
init_fn(session)
```

# Running inference on a set of images




```
#prediction Function for running inference
```

```
def predict_fn(images):  
    return session.run(probabilities, feed_dict={processed_images: images})
```

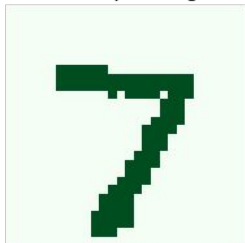
```
preds = predict_fn(images)  
for x in preds.argsort()[0][-5:]:  
    print (x, names[x], preds[0,x])
```

# Which LIME class to use?

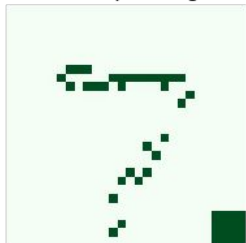


```
explainer = lime_image.LimeImageExplainer()  
explanation = explainer.explain_instance(image, predict_fn,  
                                       top_labels=10, hide_color=0, num_samples=1000)
```

LIME explaining 7



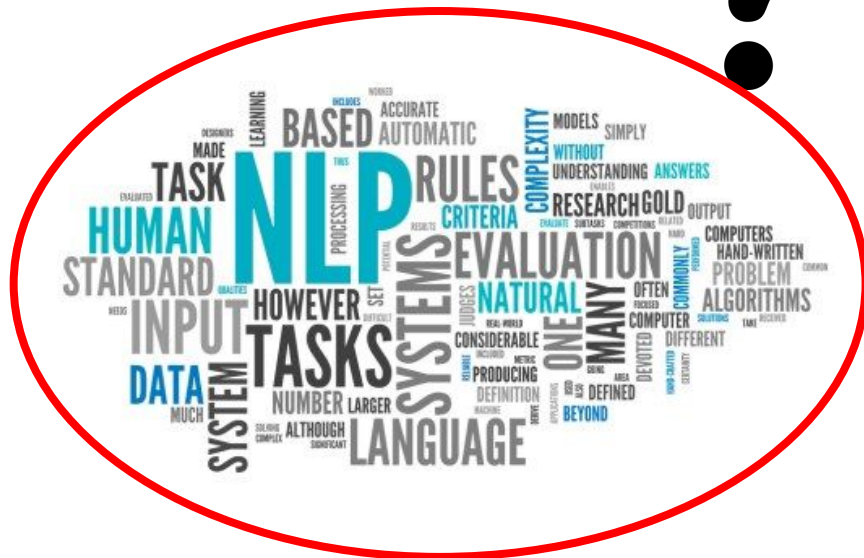
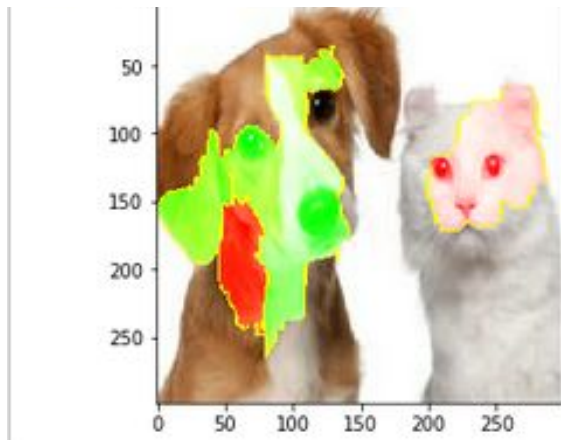
LIME explaining 3



Input gradients for 7



Input gradients for 3



# We have seen digits, cats and dogs..

Let's see Text classification on stackoverflow post and tags dataset

	post	tags
0	what is causing this behavior in our c# datet...	c#
1	have dynamic html load as if it was in an ifra...	asp.net
2	how to convert a float value in to min:sec i ...	objective-c
3	.net framework 4 redistributable just wonderi...	.net
4	trying to calculate and print the mean and its...	python

**How will we do it?**

**- We'll first convert words to vector using  
CountVectorizer()**

# Next , we build a Logistic regression model for text data



```
Vectorizer = CountVectorizer(analyser='word', token_pattern=r'\w(1,)', ngram_range=(1, 3), stop_words =  
'english', binary=True) train_vectors = vectorizer.fit_transform(X_train)  
#iniatializing the model  
logreg = LogisticRegression(n_jobs=1, C=1e5)  
  
#fitting the model with train_vectors  
logreg.fit(train_vectors, y_train)
```

# LimeTextExplainer

```
'''
    LimeTextExplainer
    LimeTextExplainer(kernel_width=25, verbose=False, class_names=None, feature_selection=u'auto',
split_expression=u'\W+', bow=True)
    Explains text classifiers. It uses an exponential kernel on cosine distance,
'''
explainer = LimeTextExplainer(class_names=class_names)

'''
    explain_instance(text_instance, classifier_fn, labels=(1, ), top_labels=None, num_features=10,
num_samples=5000, distance_metric=u'cosine', model_regressor=None)
'''

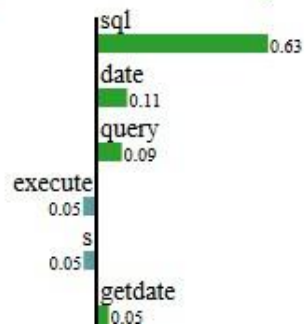
exp = explainer.explain_instance(X_test[idx], c.predict_proba, num_features=6, labels=[5, 3])
```

### Prediction probabilities

sql	1.00
python	0.00
mysql	0.00
angularjs	0.00
Other	0.00

NOT sql

sql



NOT python

python





# Highlight the words for more context

## Prediction probabilities

html	<div><div></div></div>	0.95
.net	<div><div>0.01</div></div>	
c#	<div><div>0.00</div></div>	
jquery	<div><div>0.00</div></div>	
Other	<div><div>0.00</div></div>	

## NOT html

html	<div><div></div></div>	0.96
document	<div><div>0.03</div></div>	
element	<div><div>0.02</div></div>	
specific	<div><div>0.02</div></div>	
project	<div><div>0.02</div></div>	
available	<div><div>0.01</div></div>	

## html

## Text with highlighted words

user friendly tools for describing an element on a **html** document i ve taken on a new project that s already been in development for some time the development team and project managers have a good issue tracking system in place and are generally good people to work with but i m finding it hard to follow legacy issues and understand exactly what people are talking about this is partly because i am new to the project and it s of a fairly large magnitude but also partly due to vague descriptions of technical problems i want to know if there are any simple tools available to exactly describe an element that exists in a **html** document or languages that help describe specific problems for developers xpath seems like a good starting point this must be easy for project managers to work with not just developers usage scenarios are another tool i can think of that might be relevant

All colab notebooks and slide deck are here

<https://github.com/laishawadhwa/PyConIndia2020>

# Things to Explore

- You can check out similar library for interpretability called: SHAP (<https://shap.readthedocs.io/en/latest/>)
- **A Guide for Making Black Box Models Explainable** by **Christoph Molnar** is a great book to refer to if you want to delve deep and make machine learning more interpretable.

# Conclusion



- Trust
- Transparency
- Public perception
- Improvement through feedback
- Making informed decisions
- Explanations are short, selective and possibly contrastive.



- Sampling from a Gaussian distribution.
- Complexity of the explanation model has to be defined in advance.

# What have we learnt?

- Don't trust models without explanations !
- Interpreting the machine's intelligence is tough ( For Humans as well).
- For Business centric use cases we must ensure that model is learning the scene description well irrespective of the accuracy. (Dog wolf example!).

# The Future of Interpretability

- ✓ The focus will be on model-agnostic interpretability tools.
- ✓ Machine learning will be automated and, with it, interpretability.
- ✓ We do not analyze data, we analyze models.

Test whether A or B is better: For this we can also use partial dependence functions. What we do not have yet are statistical tests for arbitrary black box models.

- ✓ The data scientists will automate themselves.
- ✓ Interpretability could boost machine intelligence research.

# References

- 1 "Why Should I Trust You?" Explaining the Predictions of Any Classifier — Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin.
- 2 Blog by the authors,  
<https://homes.cs.washington.edu/~marcotcr/blog/lime/>
- 3 Github link : <https://github.com/marcotcr/lime>:
- 4 Link to original paper : <https://arxiv.org/abs/1602.04938>



# Questions?



@laishawadhwa



/laishawadhwa





---

# Thank You!

---