Mini-Survey of Attribution



Why?

- Interpretability of Deep Learning
- Explain what input explains decisions



Background

Some Types of Attribution

- Data Analysis
- Relationship between inputs and outputs
 - Pertinent parts of the dataset
 - Dataset relationship to the output
 - Relationship between model layers and the output
- Feature Analysis



Number of data points per driving strategy





Data Analysis





Bias? Distribution

Attribution Methods

- Model relationship tests
- Perturbation Methods
- Gradient methods
- Other methods

Model relationship tests

- Model Parameter Randomization Test
 - Cascading Randomization
 - Independant Randomization
- Data Randomization Test

Cascading Randomization

- The weights of a model are randomized over time starting from the top layers and moving down to the bottom ones.
- This test shows the sensitivity of an attribution method to the model's parameters.



Independant Randomization

- This is done by performing randomization layer-bylayer rather than by weight
- This gives a more granular indication of dependency for an attribution method by the order of the layer.



Data randomization results

Perturbation

- Type of function which compares two networks
 - The original network
 - A network trained on the dataset where relevant features have been altered
 - Masked
 - Obscured
 - Removed
 - Biased

Gradient Methods

• Make use of the gradient based methods in the backpropagation step to attempt to extract spatial information captured in the network

Saliency Maps

• Compute the absolute partial derivative of the output neuron with respect to the input features

 $rac{\partial \mathrm{output}}{\partial input}$



Issues with Gradient Based Methods

- Many methods partially reconstruct the input data
- Brittle to noise and interference (misleading results)
- Many of the advanced guided methods dont have an adequate relationship between the input data and output nodes of a network
- Some methods (like some saliency maps) may not work with features that have a negative effect on the output

Other Methods

- Tree approximations
- Uncertainty of the models (Sensitivity Analysis)
- Prototype Selection

Issues with Attribution

- The outputs of networks may have dimensions that are much smaller than the inputs
- Noise becomes a factor which may cause artifacting
- The boundaries of visualisation techniques may become distorted due to network shape and feature extractor size
- The final step of attribution may introduce its own bias i.e. upscaling, resizing etc
- Computational cost

Evaluating Attribution

• Take the result of an attribution method and find reasoning from it

Evaluating Attribution

- Generate visualisations which are then interpreted by a human
- Perturbation techniques on top of other techniques to find localised importance from the dataset
- Model performance metric and perturbation
- Calculating similarity between a truth dataset feature map and an extracted feature map (Something like KL-divergence)

KL-Divergence

Kullback and Leibler (1951) (and Belov and Armstrong (2011)) define the KL Divergence as:

$$D(g\|b) = \int_{-\infty}^{+\infty} g(x) \ln \frac{g(x)}{b(x)} dx$$

where g(x) and h(x) are probability density functions.

