

Example 1 for Attack textfooler

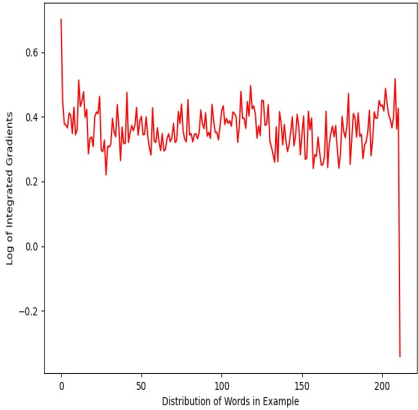
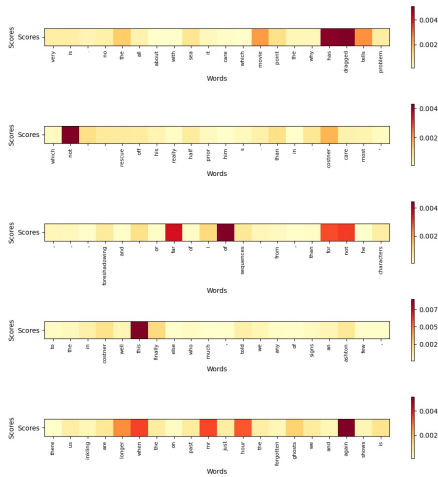
Example Text

Once again Mr. Costner has dragged out a movie for far longer than necessary. Aside from the terrific sea rescue sequences, of which there are very few I just did not care about any of the characters. Most of us have ghosts in the closet, and Costner's character are realized early on, and then forgotten until much later, by which time I did not care. The character we should really care about is a very cocky, overconfident Ashton Kutcher. The problem is he comes off as kid who thinks he's better than anyone else around him and shows no signs of a cluttered closet. His only obstacle appears to be winning over Costner. Finally when we are well past the half way point of this stinker, Costner tells us all about Kutcher's ghosts. We are told why Kutcher is driven to be the best with no prior inkling or foreshadowing. No magic here, it was all I could do to keep from turning it off an hour in.

Analysis Begins

Known Example Type:	Clean
Example original label:	0
Surrogate Model predicted label:	0
Confidence on Label:	0.99975914

Explainability Analysis



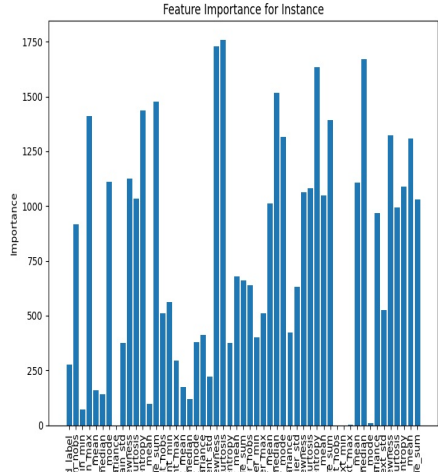
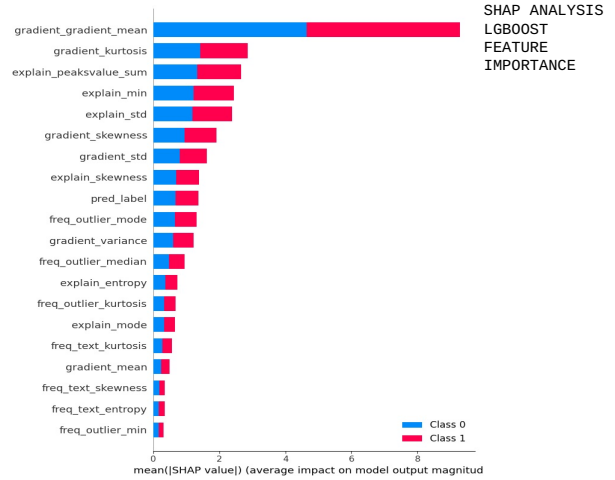
ATTENTION MAPS INTEGRATED GRADIENTS

Detector Inference

Example is Detected as :Clean

Analysis finish

Feature Importance Analysis of Adv Detector



Analysis finished

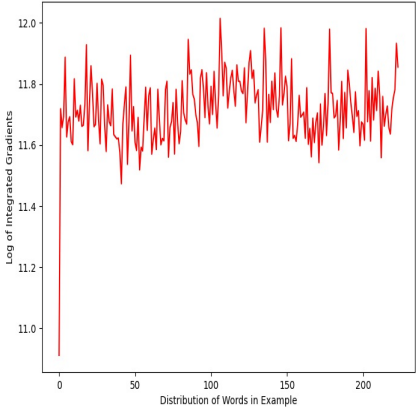
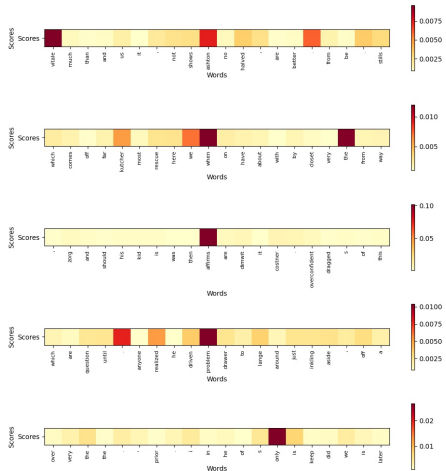
Example Text

Whenever again Mr. Costner has dragged out a stills for far lange than vitale. Aside from the wondrous sea rescue sequential, of which there are very few I just did not zorg about any of the peculiarity. Most of us have poltergeists in the closet, and Costner's character are realized early on, and then forgotten until much later, by which time I did not care. The character we should really care about is a very cocky, overconfident Ashton Kutcher. The problem is he comes off as kid who thinks he's better than anyone else around him and shows no signs of a cluttered drawer. His only obstacle appears to be winning over Costner. Finally when we are well past the halved way question of this dimwit, Costner affirms us all about Kutcher's ghosts. We are told why Kutcher is driven to be the best with no prior inkling or foreshadowing. No magic here, it was all I could do to keep from turning it off an hour in.

Analysis Begins

Known Example Type:	Adversarial
Example original label:	0
Surrogate Model predicted label:	1
Confidence on Label:	0.8150268

Explainability Analysis



ATTENTION MAPS INTEGRATED GRADIENTS

Detector Inference

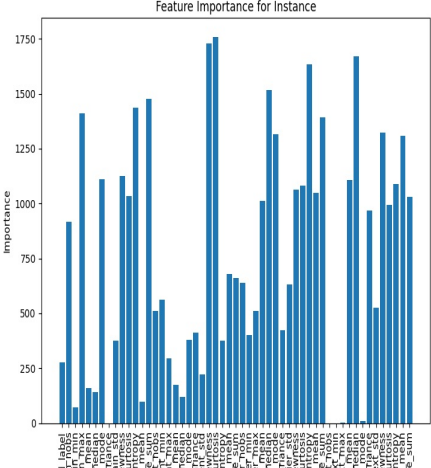
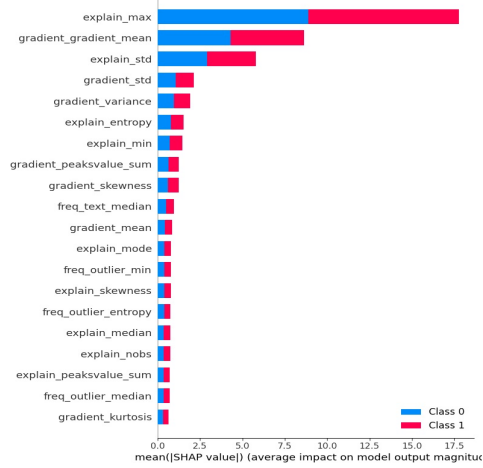
Example is Detected as :Adversarial

Identification and Transformation

Identified Words:	['costner', 'vitale', 'wondrous', 'sea', 'sequential', 'zorg', 'peculiarity', 'poltergeists', 'costner', 'character', 'character', 'cocky', 'overconfident', 'ashton', 'kutcher', 'problem', 'drawer', 'appears', 'winning', 'costner', 'finally', 'halved', 'dimwit', 'costner', 'affirms', 'kutcher', 'told', 'kutcher', 'best', 'inkling', 'foreshadowing', 'aside', 'wondrous', 'sequential', 'poltergeists', 'closet', 'time', 'signs', 'obstacle', 'finally', 'past', 'halved', 'way', 'question', 'affirms', 'magic']
Possible Replacements:	{'wondrous': 'marvellous', 'sea': 'ocean', 'sequential': 'computation', 'peculiarity': 'shortcoming', 'poltergeists': 'commonplaces', 'character': 'protagonist', 'cocky': 'smug', 'overconfident': 'presumptuous', 'problem': 'serious', 'drawer': 'briefcase', 'appears': 'seems', 'winning': 'winner', 'finally': 'eventually', 'halved': 'trimmed', 'dimwit': 'plainspoken', 'affirms': 'tells'}
The predicted label after Transform:	0
Confidence on Label:	0.99636793
Human intervention is required:	NO
Message:	Converted to non-ADVERSARIAL EXAMPLE with newscore
Transformed Texted:	whenever again mr . costner has dragged out a stills for far lange than vitale . aside from the marvellous ocean rescue computation , of which there are very few i just did not org about any of the shortcoming . most of us have commonplaces in the closet , and costner 's protagonist are realized early on , and then forgotten until much later , by which time i did not care . the character we should really care about is a very smug , presumptuous ashton kutcher . the serious is he comes off as kid who thinks he 's better than anyone else around him and shows no signs of a cluttered briefcase . his only obstacle seems to be winner over costner . eventually when we are well past the trimmed way question of this plainspoken , costner tells us all about kutcher 's ghosts . we are told why kutcher is driven to be the best with no prior inkling or foreshadowing . no magic here , it was all i could do to keep from turning it off an hour in .

Feature Importance Analysis of Adv Detector

SHAP ANALYSIS LGB00ST FEATURE IMPORTANCE



Analysis finished

Example 2 for Attack textfooler

Example Text

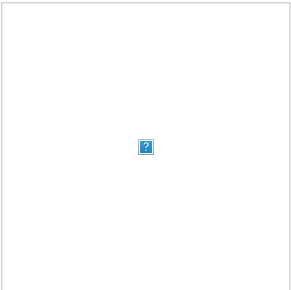
This is an example of why the majority of action films are the same. Generic and boring, there's really nothing worth watching here. A complete waste of the then barely-tapped talents of Ice-T and Ice Cube, who've each proven many times over that they are capable of acting, and acting well. Don't bother with this one, go see New Jack City, Ricochet or watch New York Undercover for Ice-T, or Boyz n the Hood, Higher Learning or Friday for Ice Cube and see the real deal. Ice-T's horribly cliched dialogue alone makes this film grate at the teeth, and I'm still wondering what the heck Bill Paxton was doing in this film? And why the heck does he always play the exact same character? From Aliens onward, every film I've seen with Bill Paxton has him playing the exact same irritating character, and at least in Aliens his character died, which made it somewhat gratifying...

Overall, this is second-rate action trash. There are countless better films to see, and if you really want to see this one, watch Judgement Night, which is practically a carbon copy but has better acting and a better script. The only thing that made this at all worth watching was a decent hand on the camera - the cinematography was almost refreshing, which comes close to making up for the horrible film itself - but not quite. 4/10.

Analysis Begins

Known Example Type:	Clean
Example original label:	0
Surrogate Model predicted label:	0
Confidence on Label:	0.99981695

Explainability Analysis



Detector Inference

Idenfication and Transformation

Example 2

Example # 1 Analysis

K+Pbbb/Hyyy9LW609zz/0oaGhckdQPLVajfT0dIiiCBsb61y9ehUNGzaU0xZpcerUKaSkpEi3v/jiC7Rr106WLG5ubpg6dSr/UK4mQkNDYWnjg/j4eERERCA60lqaFiWHbdu2YdKkSRAEAveuXMGff/6JDz/8EL/++qtsmZTC3N width="500" height="500" />

Detector Inference

Example is Detected as :Clean

Analysis finish

Feature Importance Analysis of Adv Detector

Analysis finished

Example Text

These is an example of why the majority of action films are the same. Generic and boring, there's really nothing worth watching here. A complete detritus of the then barely Overall, this is second-rate endeavor trash. There are countless better films to see, and if you really want to see this one, watch Judgement Night, which is practically a

Analysis Begins

Explainability Analysis

ATTENTION MAPS INTEGRATED GRADIENTS

Detector Inference

Example is Detected as :Clean

Analysis finish

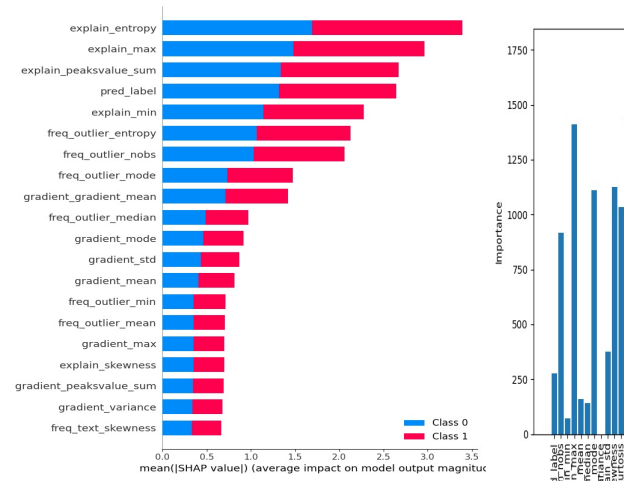
Feature Importance Analysis of Adv Detector

SHAP ANALYSIS LGBBOOST FEATURE IMPORTANCE

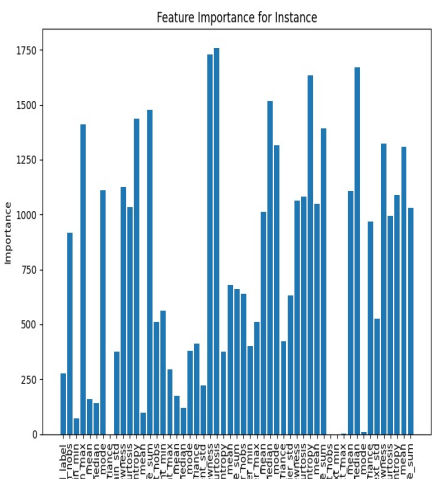
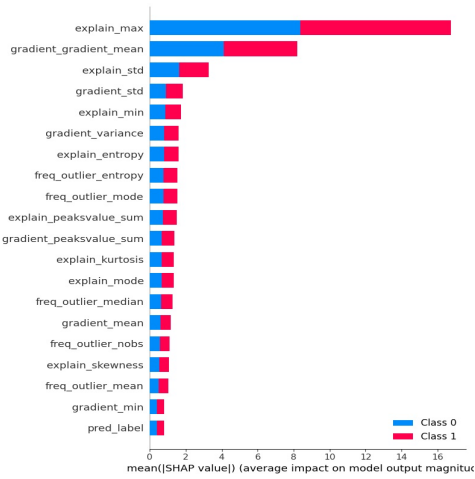
Analysis finished

Example 3 for Attack textfooler

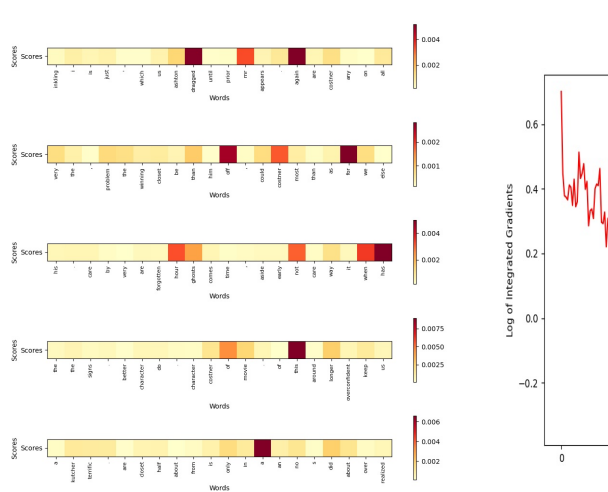
Example Text



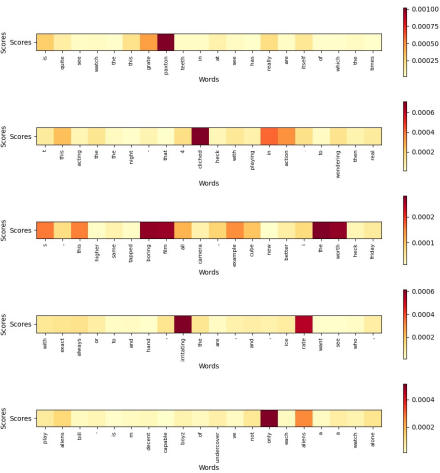
Analysis Begins



Explainability Analysis



Explainability Analysis



ATTENTION MAPS INTEGRATED GRADIENTS

Detector Inference

Example is Detected as :Clean

Analysis finish

Feature Importance Analysis of Adv Detector

SHAP ANALYSIS LGBBOOST FEATURE IMPORTANCE

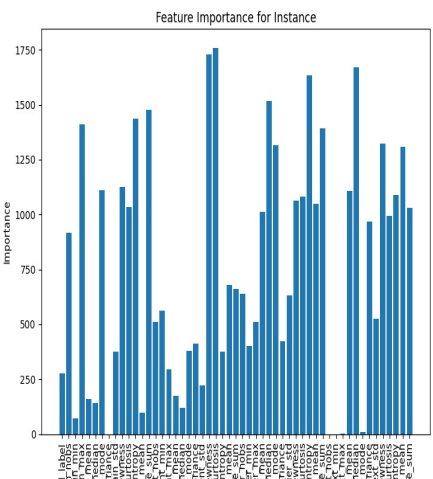
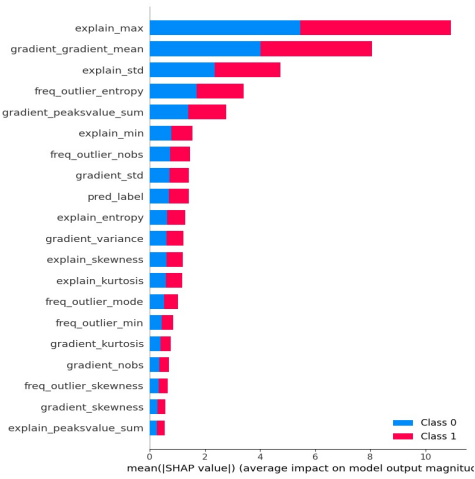
Analysis finished

Example Text

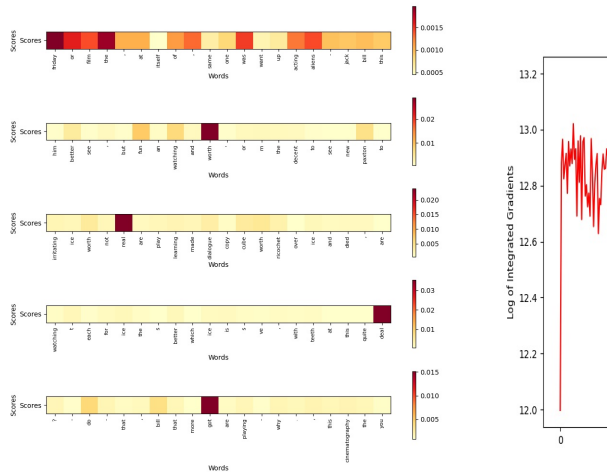
This is an example of why the lots of action films are the same. Ge

Overall, this is second-rate action trash. There are numerous better

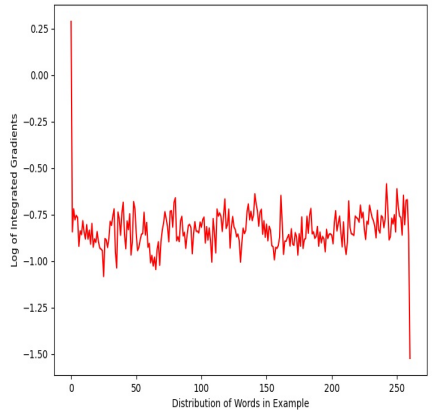
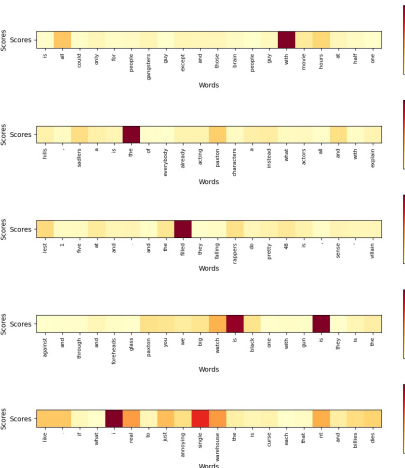
Analysis Begins



Explainability Analysis



Explainability Analysis



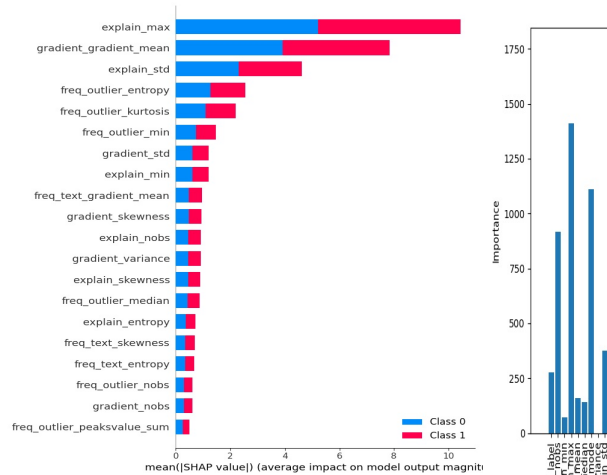
ATTENTION MAPS INTEGRATED GRADIENTS

Detector Inference

Example is Detected as :Clean

Analysis finish

Feature Importance Analysis of Adv Detector



Explainability Analysis

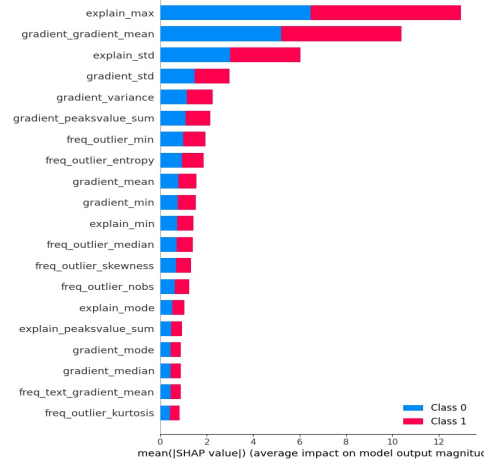
ATTENTION MAPS INTEGRATED GRADIENTS

Detector Inference

Example is Detected as :Adversarial

Identification and Transformation

Feature Importance Analysis of Adv Detector



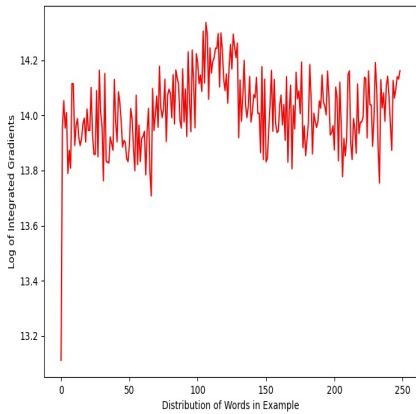
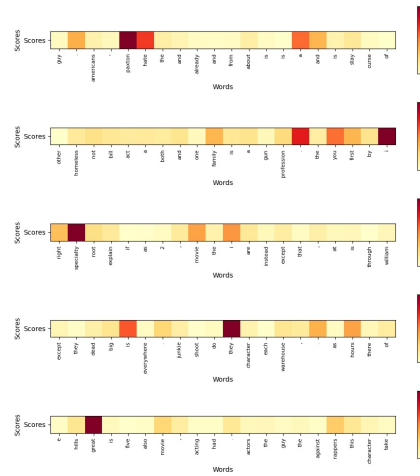
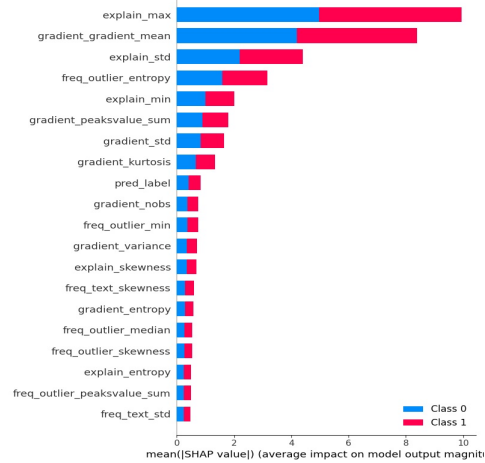
ATTENTION MAPS INTEGRATED GRADIENTS

Detector Inference

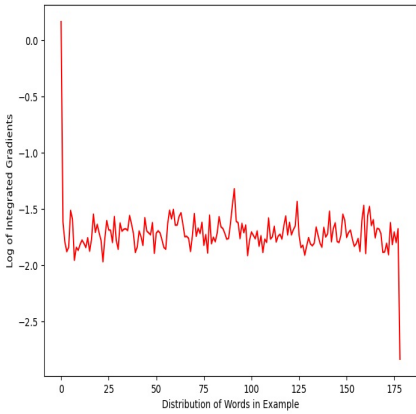
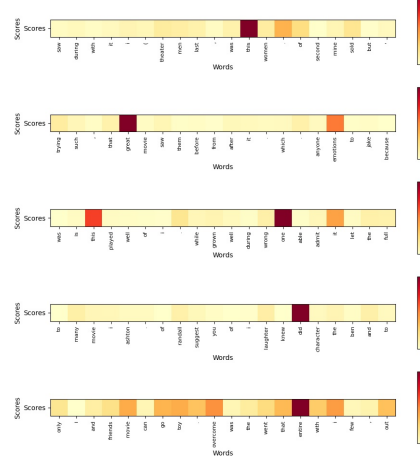
Example is Detected as :Clean

Analysis finish

Feature Importance Analysis of Adv Detector

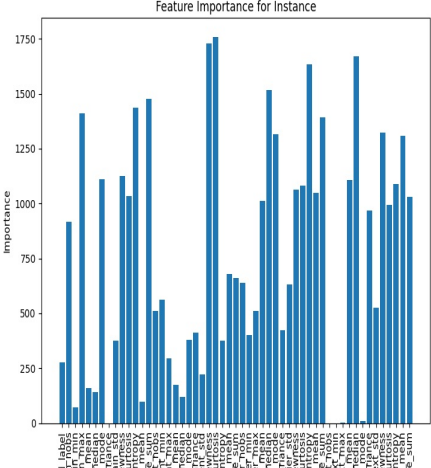
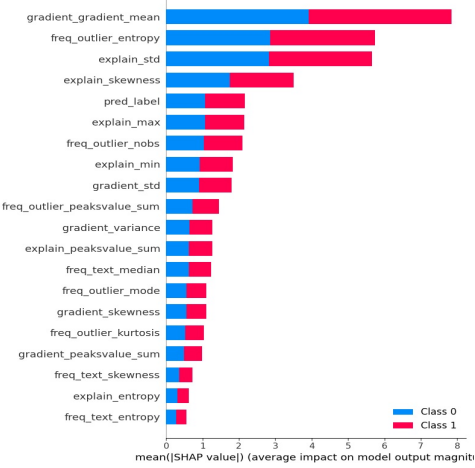


Explainability Analysis



Analysis finished

Example 2 for Attack deepwordbug

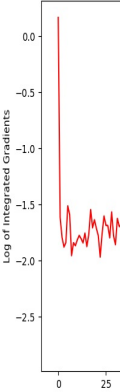
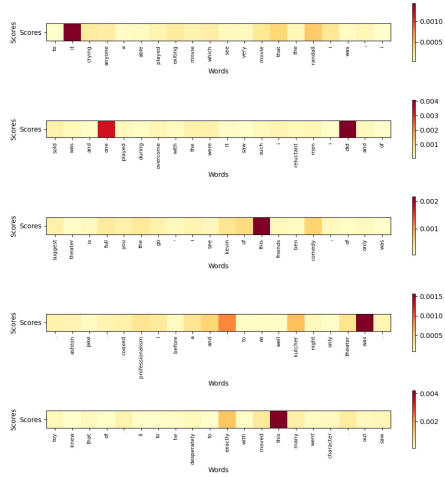


Example Text

I went and saw this movie last night after being coaxed to by a few

Analysis Begins

Explainability Analysis

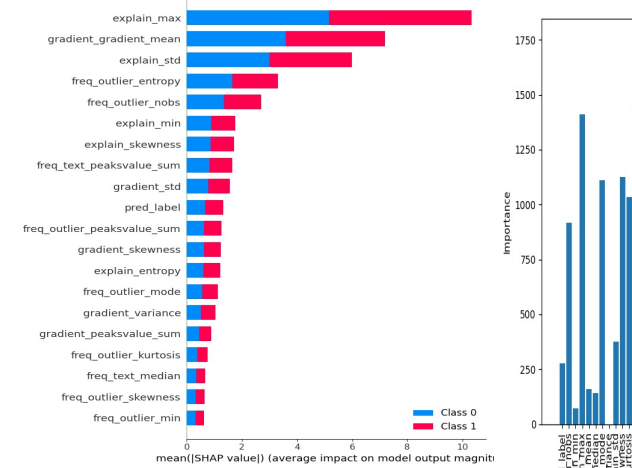
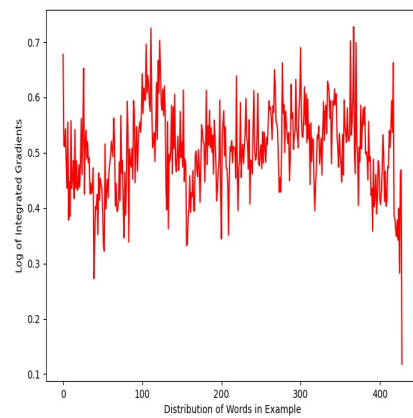


Explainability Analysis

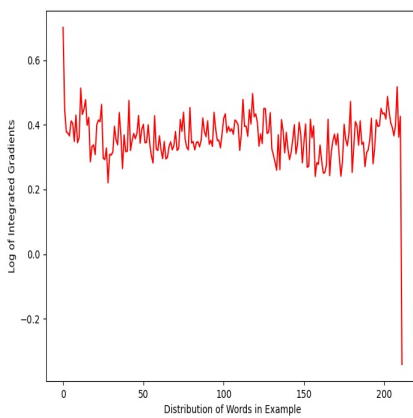
ATTENTION MAPS INTEGRATED GRADIENTS

Detector Inference

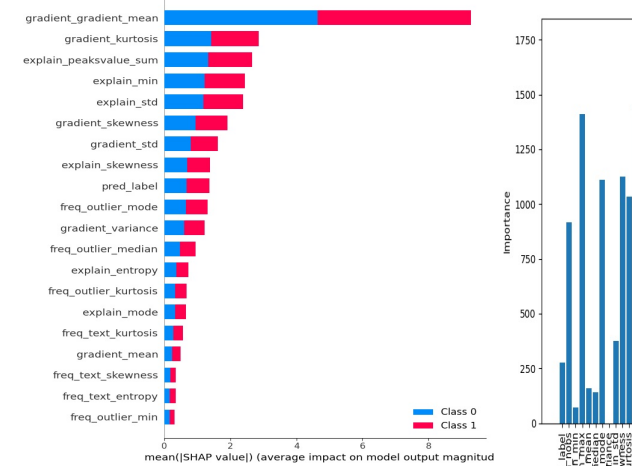
Feature Importance Analysis of Adv Detector



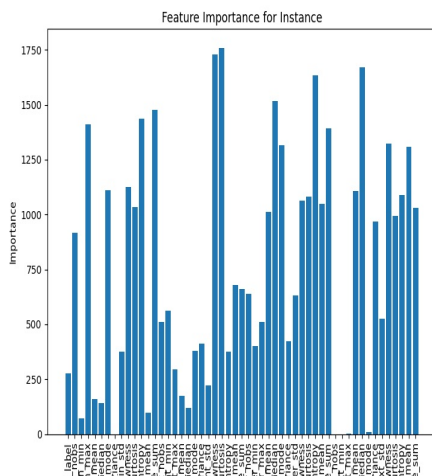
Analysis Begins

[illegible]

Feature Importance Analysis of Adv Detector



Analysis finished



Analysis Begins