

GRETEL: THE GOOD, THE BAD AND THE UGLY

Bardh Prenkaj

Post-doc in AI

prenkaj@di.uniroma1.it



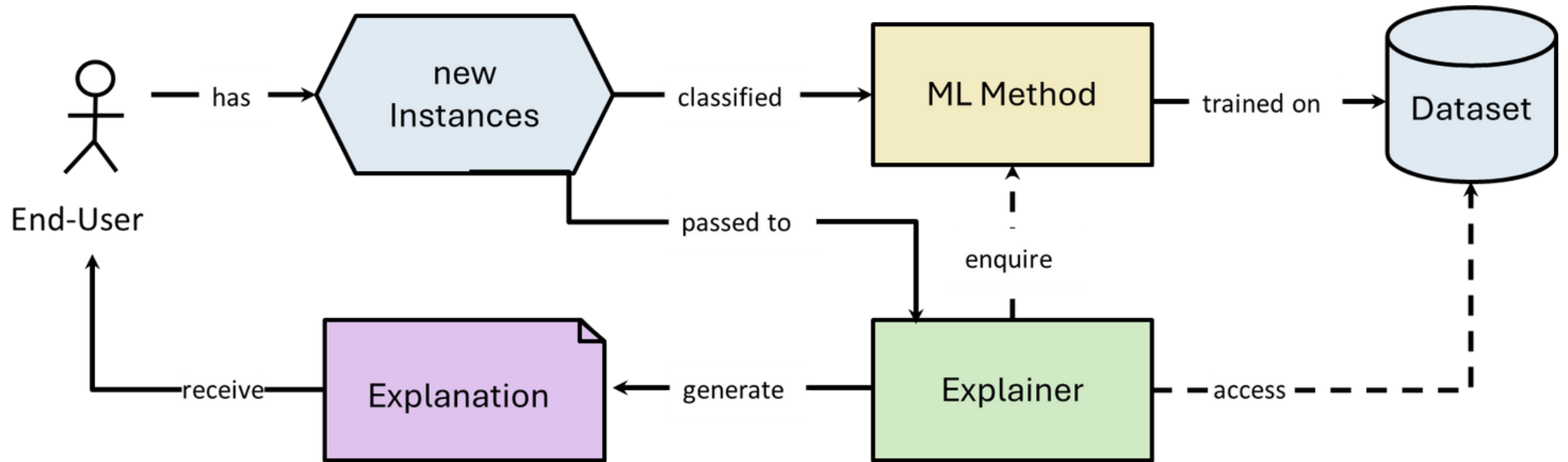
SAPIENZA
UNIVERSITÀ DI ROMA



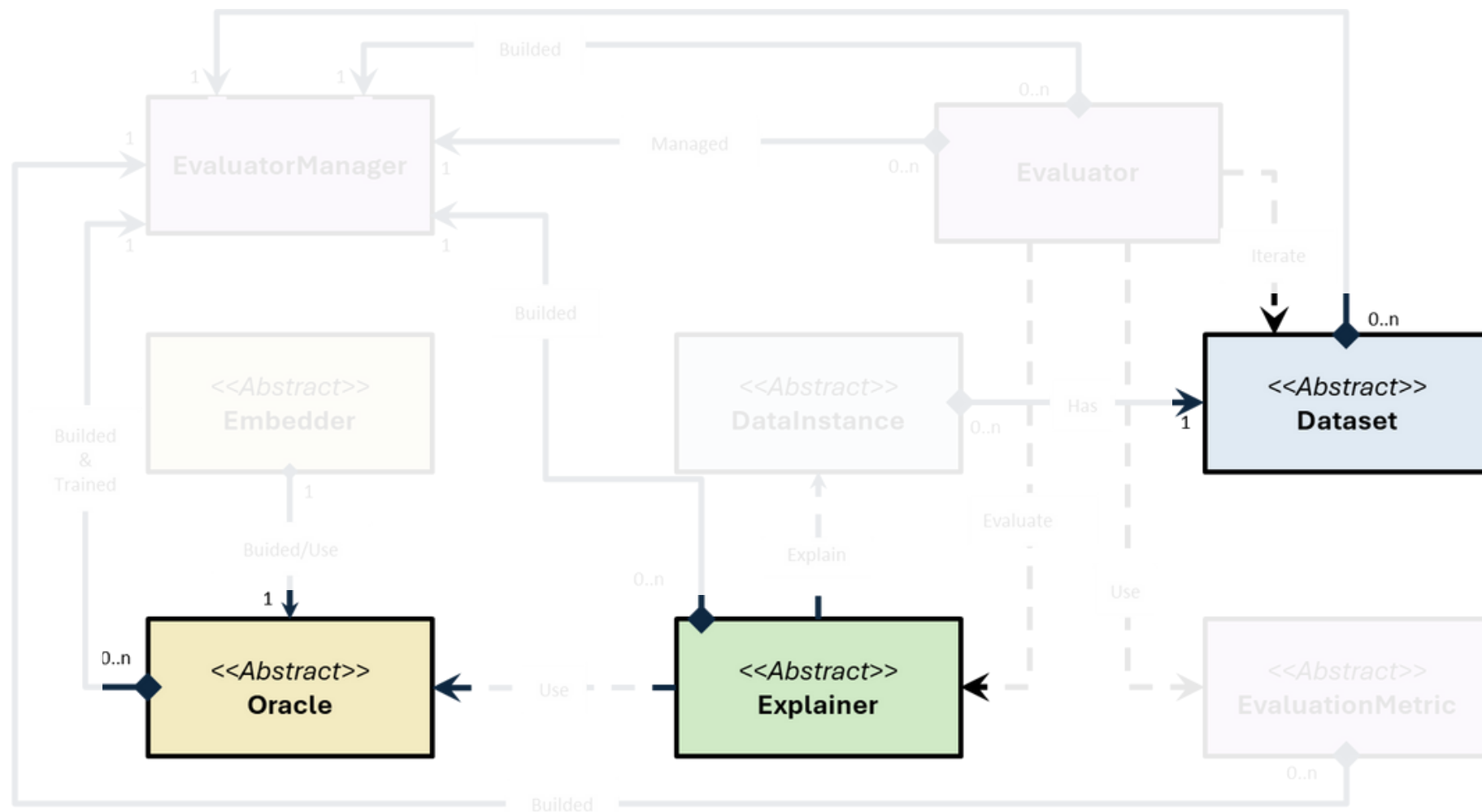
aiim



Typical Explainability Workflow



GRETEL



THE GOOD



What's **good** in GRETEL?

- Highly modular framework
- **Core philosophy:** users can implement their explainers by only extending the Explainer class without having to worry about the entire XAI pipeline
- 80% of graph classification SoA methods are covered inside
- Configuration opportunities for already supported methods, datasets, evaluations



Example of **good**: Ad-hoc explainer

```
giovanni, 4 days ago | 1 author (giovanni)
9  class DCESearchExplainer(Explainer):
10     """The Distribution Compliant Explanation Search Explainer performs a search of
11     the minimum counterfactual instance in the original dataset instead of generating
12     a new instance"""
13
14     def init(self):
15         super().init()
16         self._gd = GraphEditDistanceMetric()
17         self.fold_id=-1
18         self.dist_mat = np.full((len(self.dataset.instances), len(self.dataset.instances)), -1)
19         self.cls_mat = np.full((len(self.dataset.instances), len(self.dataset.instances)), -1)
20
21
22     def explain(self, instance):
23         l_input_inst = self.oracle.predict(instance)
24
25         # if the method does not find a counterfactual example returns the original graph
26         min_counterfactual = instance
27
28         min_counterfactual_dist = sys.float_info.max
29
```



Example of **good**: Ad-hoc explainer

```
giovanni, 4 days ago | 1 author (giovanni)
9  class DCESearchExplainer(Explainer):
10     """The Distribution Compliant Explanation Search Explainer performs a search of
11     the minimum counterfactual instance in the original dataset instead of generating
12     a new instance"""
13
14     def init(self):
15         super().init()
16         self._gd = GraphEditDistanceMetric()
17         self.fold_id=-1
18         self.dist_mat = np.full((len(self.dataset.instances), len(self.dataset.instances)), -1)
19         self.cls_mat = np.full((len(self.dataset.instances), len(self.dataset.instances)), -1)
20
21     def explain(self, instance):
22         l_input_inst = self.oracle.predict(instance)
23
24         # if the method does not find a counterfactual example returns the original graph
25         min_counterfactual = instance
26
27         min_counterfactual_dist = sys.float_info.max
28
29
```




Just extend the
abstract Explainer class



Example of **good**: Ad-hoc explainer

```
giovanni, 4 days ago | 1 author (giovanni)
9 class DCESearchExplainer(Explainer):
10     """The Distribution Compliant Explanation Search Explainer performs a search of
11     the minimum counterfactual instance in the original dataset instead of generating
12     a new instance"""
13     giovanni, 4 days ago • FIRST IMPORT GOLD
14
15     def init(self):
16         super().init()
17         self._gd = GraphEditDistanceMetric()
18         self.fold_id=-1
19         self.dist_mat = np.full((len(self.dataset.instances), len(self.dataset.instances)), -1)
20         self.cls_mat = np.full((len(self.dataset.instances), len(self.dataset.instances)), -1)
21
22     def explain(self, instance):
23         l_input_inst = self.oracle.predict(instance)
24
25         # if the method does not find a counterfactual example returns the original graph
26         min_counterfactual = instance
27
28         min_counterfactual_dist = sys.float_info.max
29
```



Define the abstract
method explain



Example of **good**: Ad-hoc dataset

```
giovanni, 4 days ago | 1 author (giovanni)
11 class ADHD(Generator):      giovanni, 4 days ago • FIRST IMPORT GOLD
12
13     def init(self):
14         base_path = self.local_config['parameters']['data_dir']
15         # Path to the instances of the "Attention Deficit Hyperactivity Disorder class"
16         self.adhd_class_path = join(base_path, 'adhd_dataset')
17         self._td_file_path = join(base_path, 'td')
18         self.generate_dataset()
19
20     def get_num_instances(self):
21         return len(self.dataset.instances)
22
23     def generate_dataset(self):
24         if not len(self.dataset.instances):
25             self.read_adjacency_matrices()
26
```



Example of **good**: Ad-hoc dataset

```
giovanni, 4 days ago | 1 author (giovanni)
11 class ADHD(Generator):      giovanni, 4 days ago • FIRST IMPORT GOLD
12
13     def init(self):
14         base_path = self.local_config['parameters']['data_dir']
15         # Path to the instances of the "Attention Deficit Hyperactivity Disorder class"
16         self.adhd_class_path = join(base_path, 'adhd_dataset')
17         self._td_file_path = join(base_path, 'td')
18         self.generate_dataset()
19
20     def get_num_instances(self):
21         return len(self.dataset.instances)
22
23     def generate_dataset(self):
24         if not len(self.dataset.instances):
25             self.read_adjacency_matrices()
26
```

Just extend the abstract
Generator class



Example of **good**: Ad-hoc dataset

```
giovanni, 4 days ago | 1 author (giovanni)
11 class ADHD(Generator):      giovanni, 4 days ago • FIRST IMPORT GOLD
12
13     def init(self):
14         base_path = self.local_config['parameters']['data_dir']
15         # Path to the instances of the "Attention Deficit Hyperactivity Disorder class"
16         self.adhd_class_path = join(base_path, 'adhd_dataset')
17         self._td_file_path = join(base_path, 'td')
18         self.generate_dataset()
19
20     def get_num_instances(self):
21         return len(self.dataset.instances)
22
23     def generate_dataset(self):
24         if not len(self.dataset.instances):
25             self.read_adjacency_matrices()
26
```

Implement the
`generate_dataset` abstract
method to contain the
dataset creation logic



Example of **good**: Extend Oracle

```
giovanni, 4 days ago | 1 author (giovanni)
9 class Oracle(Trainable,metaclass=ABCMeta):      giovanni, 4 days ago • FIRST IMPORT GOLD
10     def __init__(self, context:Context, local_config) -> None:
11         super().__init__(context, local_config)
12         self._call_counter = 0
13
14     @final
15     def predict(self, data_instance):             59
16         """predicts the label of a given data instance
17         -----
18         INPUT:
19         |   data_instance : The instance whose class is going to be predicted
20         |   -----
21         OUTPUT:
22         |   The predicted label for the data instance
23         |   """
24         self._call_counter += 1
25
26         return self._real_predict(data_instance)
27
28     @final
29     def predict_proba(self, data_instance):
30         """predicts the probability estimates for a given data instance
31         -----
```

+

```
60     @abstractmethod
61     def _real_predict(self, data_instance):
62         pass
63
64     @abstractmethod
65     def _real_predict_proba(self, data_instance):
66         pass
67
```



Example of **good**: Extend Oracle

```
giovanni, 4 days ago | 1 author (giovanni)
9 class Oracle(Trainable,metaclass=ABCMeta):      giovanni, 4 days ago * FIRST IMPORT GOLD
10     def __init__(self, context:Context, local_config) -> None:
11         super().__init__(context, local_config)
12         self._call_counter = 0
13
14     @final
15     def predict(self, data_instance):
16         """predicts the label of a given data instance
17         -----
18         INPUT:
19         |   data_instance : The instance whose class is going to be predicted
20         -----
21         OUTPUT:
22         |   The predicted label for the data instance
23         """
24         self._call_counter += 1
25
26         return self._real_predict(data_instance)
27
28     @final
29     def predict_proba(self, data_instance):
30         """predicts the probability estimates for a given data instance
31         -----
```

+

```
59
60     @abstractmethod
61     def _real_predict(self, data_instance):
62         pass
63
64     @abstractmethod
65     def _real_predict_proba(self, data_instance):
66         pass
67
```

When creating a new
Oracle (predictor model)
you need to **define** these
abstract methods



Example of **good**: Extend Oracle

```
giovanni, 4 days ago | 1 author (giovanni)
6 class TreeCyclesOracle(Oracle):      giovanni, 4 days ago • FIRST IMPORT C
7
8     def init(self):
9         super().init()
10        self.model = ""
11
12    def real_fit(self):
13        pass
14
15    def _real_predict(self, data_instance):
16        try:
17            nx.find_cycle(data_instance.get_nx(), orientation='ignore')
18            return 1
19        except nx.exception.NetworkXNoCycle:
20            return 0
21
22    def _real_predict_proba(self, data_instance):
23        # softmax-style probability predictions
24        try:
25            nx.find_cycle(data_instance.get_nx(), orientation='ignore')
26            return np.array([0,1])
27        except nx.exception.NetworkXNoCycle:
28            return np.array([1,0])
```



THE BAD



What's **bad** in GRETEL?

- High learning curve
- Only graph classification is supported for now
- Continuous development; you might find changes happening in real-time
- Some bad/ugly choices of using the same naming of built-in methods
- WILL NOT SUPPORT TENSORFLOW.... EVER



Example of **bad**: init()

```
giovanni, 4 days ago | 1 author (giovanni)
6 class TreeCyclesOracle(Oracle):      giovanni, 4 days ago • FIRST IMPORT C
7
8     def init(self):
9         super().init()
10        self.model = ""
11
12    def real_fit(self):
13        pass
14
15    def _real_predict(self, data_instance):
16        try:
17            nx.find_cycle(data_instance.get_nx(), orientation='ignore')
18            return 1
19        except nx.exception.NetworkXNoCycle:
20            return 0
21
22    def _real_predict_proba(self, data_instance):
23        # softmax-style probability predictions
24        try:
25            nx.find_cycle(data_instance.get_nx(), orientation='ignore')
26            return np.array([0,1])
27        except nx.exception.NetworkXNoCycle:
28            return np.array([1,0])
```



AND THE UGLY



What's **ugly** in GRETEL?

- A complex project structure
- Configuration files can be daunting
- There's a Context object passing around all the create objects dynamically; there's a `check_configuration` method that can set default hyperparameters
- Some `__call__` methods are overridden to obfuscate the torch-like functionality to the end-user



Example of **ugly**: check_configuration()

```
62 def check_configuration(self):
63     super().check_configuration()
64     self.local_config['parameters']['n_labels'] = self.local_config['parameters'].get('n_labels', 2)
65     self.local_config['parameters']['batch_size_ratio'] = self.local_config['parameters'].get('batch_size_ratio', .1)
66     self.local_config['parameters']['h_dim'] = self.local_config['parameters'].get('h_dim', 10)
67     self.local_config['parameters']['z_dim'] = self.local_config['parameters'].get('z_dim', 10)
68     self.local_config['parameters']['dropout'] = self.local_config['parameters'].get('dropout', .1)
69     self.local_config['parameters']['encoder_type'] = self.local_config['parameters'].get('encoder_type', 'gcn')
70     self.local_config['parameters']['graph_pool_type'] = self.local_config['parameters'].get('graph_pool_type', 'mean')
71     self.local_config['parameters']['disable_u'] = self.local_config['parameters'].get('disable_u', False)
72     self.local_config['parameters']['epochs'] = self.local_config['parameters'].get('epochs', 200)
73     self.local_config['parameters']['alpha'] = self.local_config['parameters'].get('alpha', 5)
74     self.local_config['parameters']['lr'] = self.local_config['parameters'].get('lr', 1e-3)
75     self.local_config['parameters']['weight_decay'] = self.local_config['parameters'].get('weight_decay', 1e-5)
76     self.local_config['parameters']['lambda_sim'] = self.local_config['parameters'].get('lambda_sim', 1)
77     self.local_config['parameters']['lambda_kl'] = self.local_config['parameters'].get('lambda_kl', 1)
78     self.local_config['parameters']['lambda_cfe'] = self.local_config['parameters'].get('lambda_cfe', 1)
79     self.local_config['parameters']['beta_x'] = self.local_config['parameters'].get('beta_x', 10)
80     self.local_config['parameters']['beta_adj'] = self.local_config['parameters'].get('beta_adj', 10)
81
82     n_nodes = self.local_config['parameters'].get('n_nodes', None)
83     if not n_nodes:
84         n_nodes = max([x.num_nodes for x in self.dataset.instances])
85     self.local_config['parameters']['n_nodes'] = n_nodes
86
87     self.local_config['parameters']['feature_dim'] = len(self.dataset.node_features_map)
88
```



Example of **ugly**: check_configuration()

```
62 def check_configuration(self):
63     super().check_configuration()
64     self.local_config['parameters']['n_labels'] = self.local_config['parameters'].get('n_labels', 2)
65     self.local_config['parameters']['batch_size_ratio'] = self.local_config['parameters'].get('batch_size_ratio', .1)
66     self.local_config['parameters']['h_dim'] = self.local_config['parameters'].get('h_dim', 10)
67     self.local_config['parameters']['z_dim'] = self.local_config['parameters'].get('z_dim', 10)
68     self.local_config['parameters']['dropout'] = self.local_config['parameters'].get('dropout', .1)
69     self.local_config['parameters']['encoder_type'] = self.local_config['parameters'].get('encoder_type', 'gcn')
70     self.local_config['parameters']['graph_pool_type'] = self.local_config['parameters'].get('graph_pool_type', 'mean')
71     self.local_config['parameters']['disable_u'] = self.local_config['parameters'].get('disable_u', False)
72     self.local_config['parameters']['epochs'] = self.local_config['parameters'].get('epochs', 200)
73     self.local_config['parameters']['alpha'] = self.local_config['parameters'].get('alpha', 5)
74     self.local_config['parameters']['lr'] = self.local_config['parameters'].get('lr', 1e-3)
75     self.local_config['parameters']['weight_decay'] = self.local_config['parameters'].get('weight_decay', 1e-5)
76     self.local_config['parameters']['lambda_sim'] = self.local_config['parameters'].get('lambda_sim', 1)
77     self.local_config['parameters']['lambda_kl'] = self.local_config['parameters'].get('lambda_kl', 1)
78     self.local_config['parameters']['lambda_cfe'] = self.local_config['parameters'].get('lambda_cfe', 1)
79     self.local_config['parameters']['beta_x'] = self.local_config['parameters'].get('beta_x', 10)
80     self.local_config['parameters']['beta_adj'] = self.local_config['parameters'].get('beta_adj', 10)
81
82     n_nodes = self.local_config['parameters'].get('n_nodes', None)
83     if not n_nodes:
84         n_nodes = max([x.num_nodes for x in self.dataset.instances])
85     self.local_config['parameters']['n_nodes'] = n_nodes
86
87     self.local_config['parameters']['feature_dim'] = len(self.dataset.node_features_map)
88
```

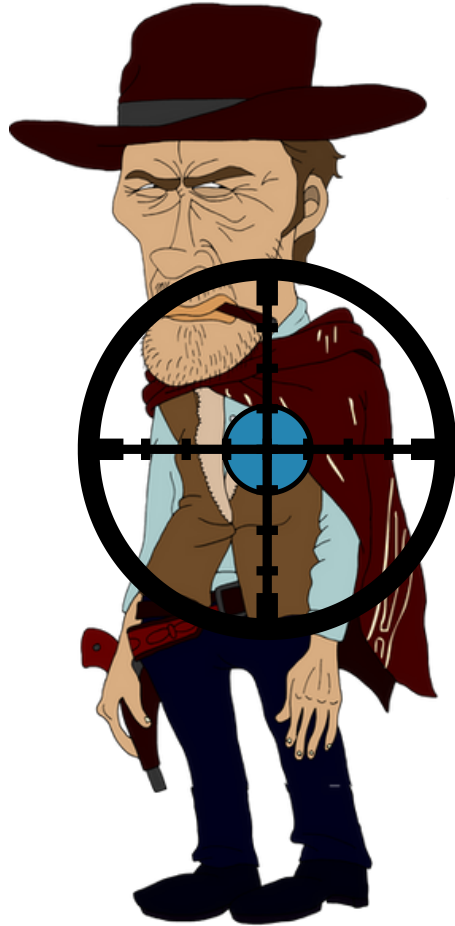


Example of ugly: `__call__`

```
107
108     def __call__(self, *args: Tuple[GraphInstance], **kwargs: Any) -> Any:
109         torch_data = torch.from_numpy(args[0].data[None, None, :, :]).float()
110         return self.generator(torch_data)    You, 22 hours ago • Uncommitted changes
```



**THE
GOOD**



**THE
BAD**



**AND THE
UGLY**



Let's code



Open research questions

- Are generative counterfactual explainers worth it?
- What's the difference between counterfactual explanations and adversarial attacks?
- How sure are we about the counterfactuals generated?
Can we incorporate uncertainty in them?
- How can we backtrack near the decision boundary once we overshoot on the other side to produce minimally perturbed counterfactuals?