

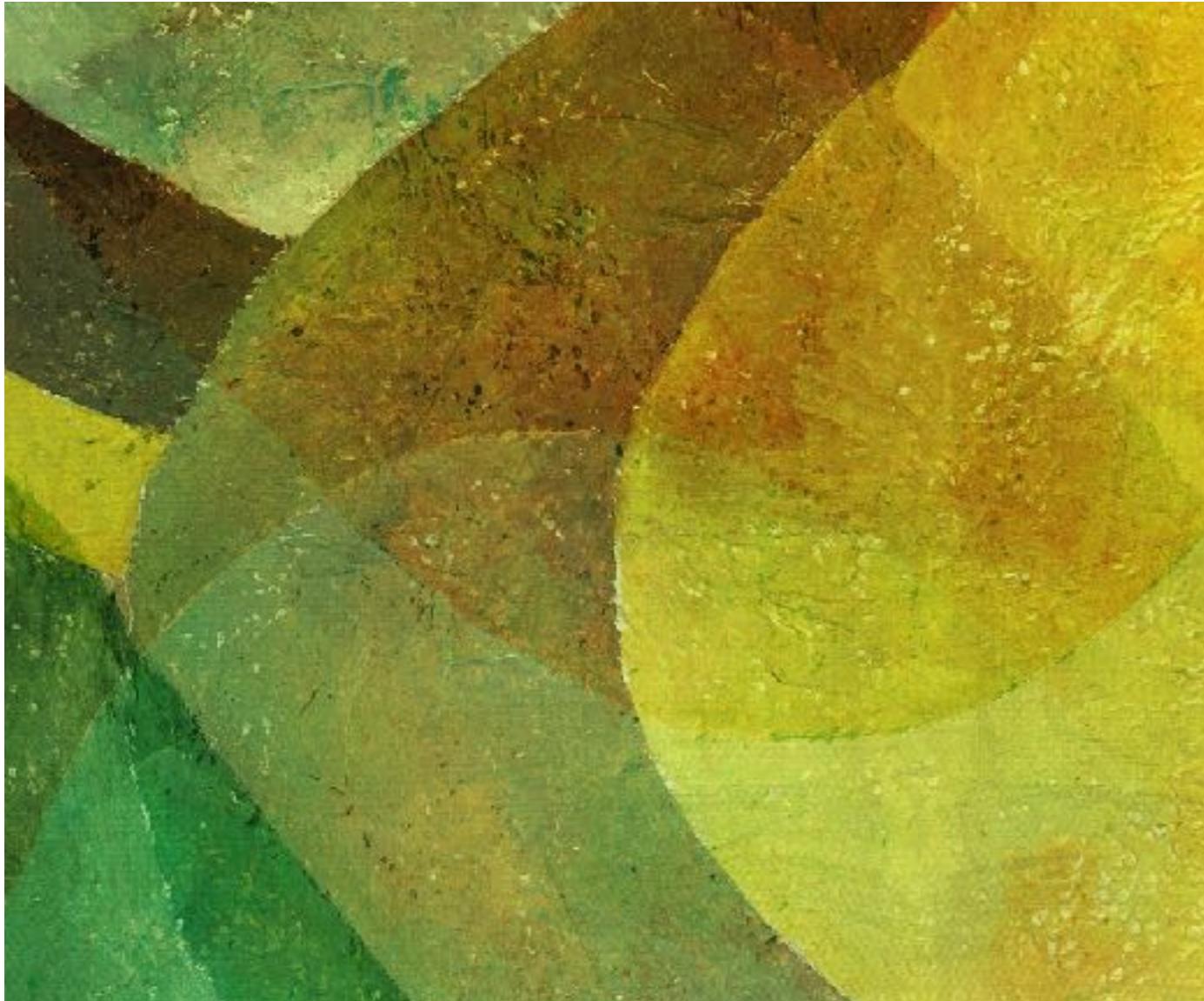
# HERMENEUTIC INVESTIGATIONS

---

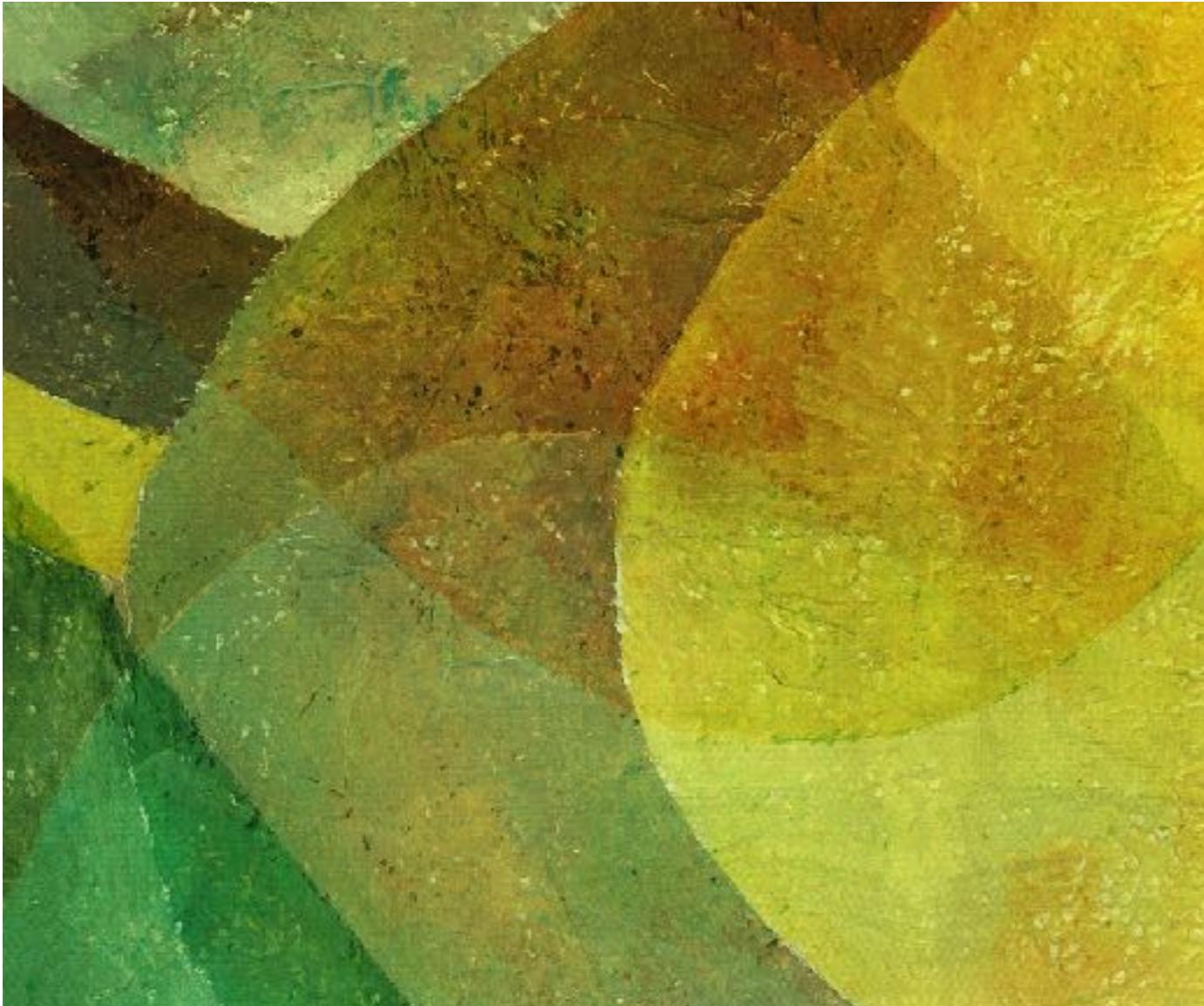
Dean Allsopp

PyData London

November 6, 2018

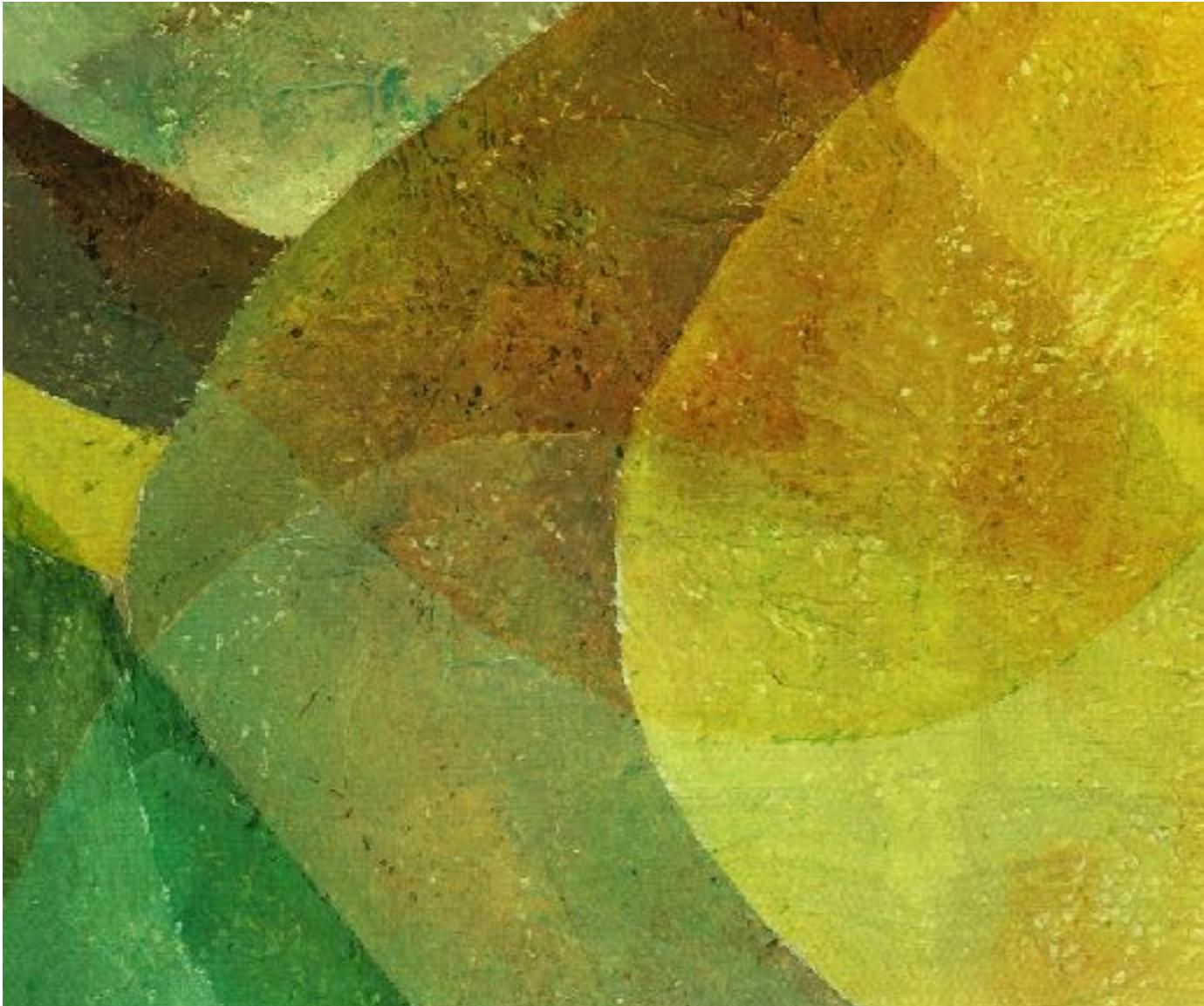


# **INTERPRETABLE MACHINE LEARNING - WHY?**



- *COMMERCIAL*
- *REGULATORY*
- *SOCIAL*

# **INTERPRETABLE MACHINE LEARNING - WHY?**



- EPISTEMIC
- MOTIVATIONAL

- COMMERCIAL
- REGULATORY
- SOCIAL

# INTERPRETABLE MACHINE LEARNING - WHY?



# HERMENEUTICS

‘The branch of knowledge that deals with interpretation’



# UNDERSTANDING VIA

---

➤ ITERATIVE  
INTERACTION



# UNDERSTANDING VIA

---

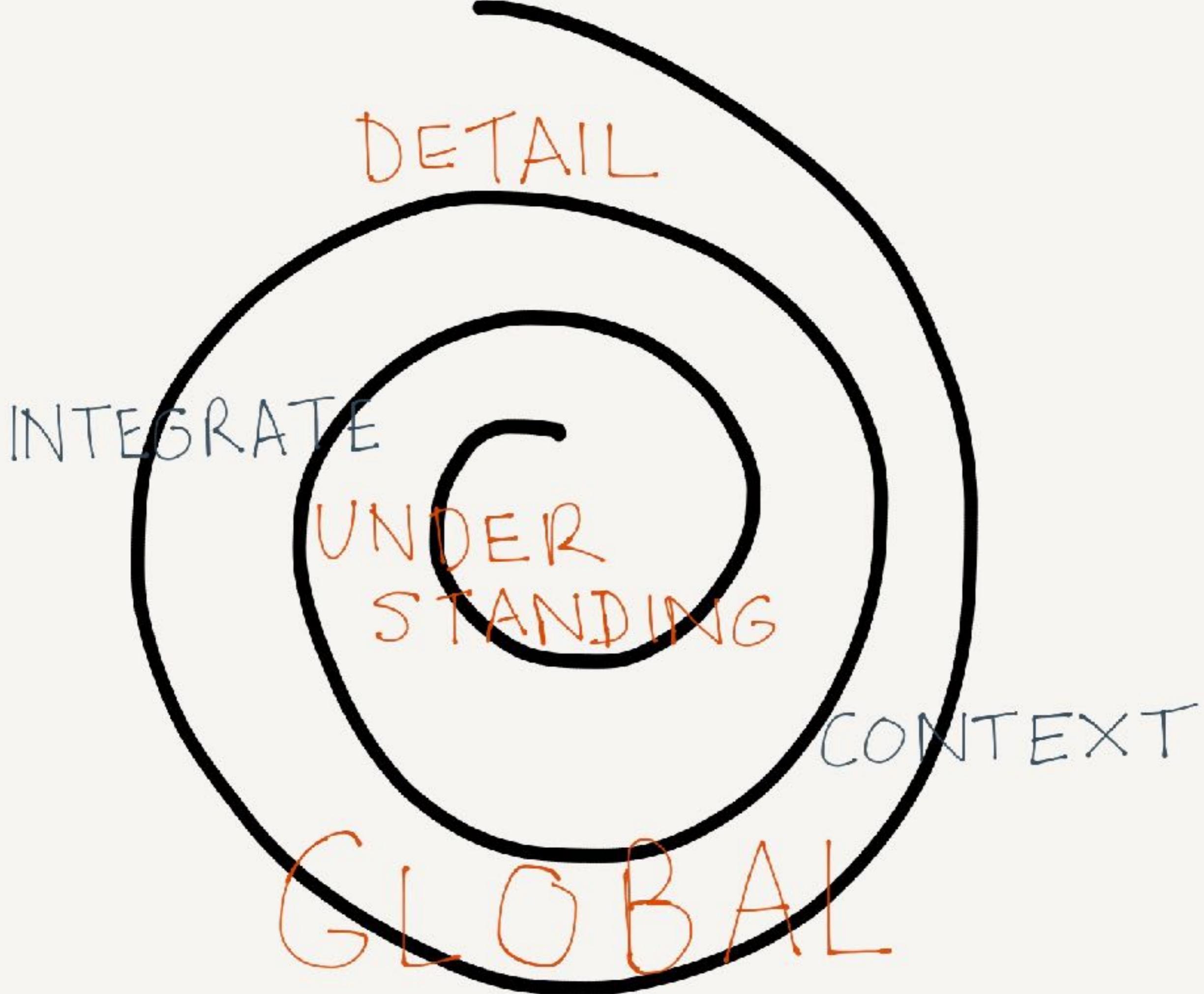
- ITERATIVE  
INTERACTION
- COMPARISON OF  
WHOLE & PART



# UNDERSTANDING VIA

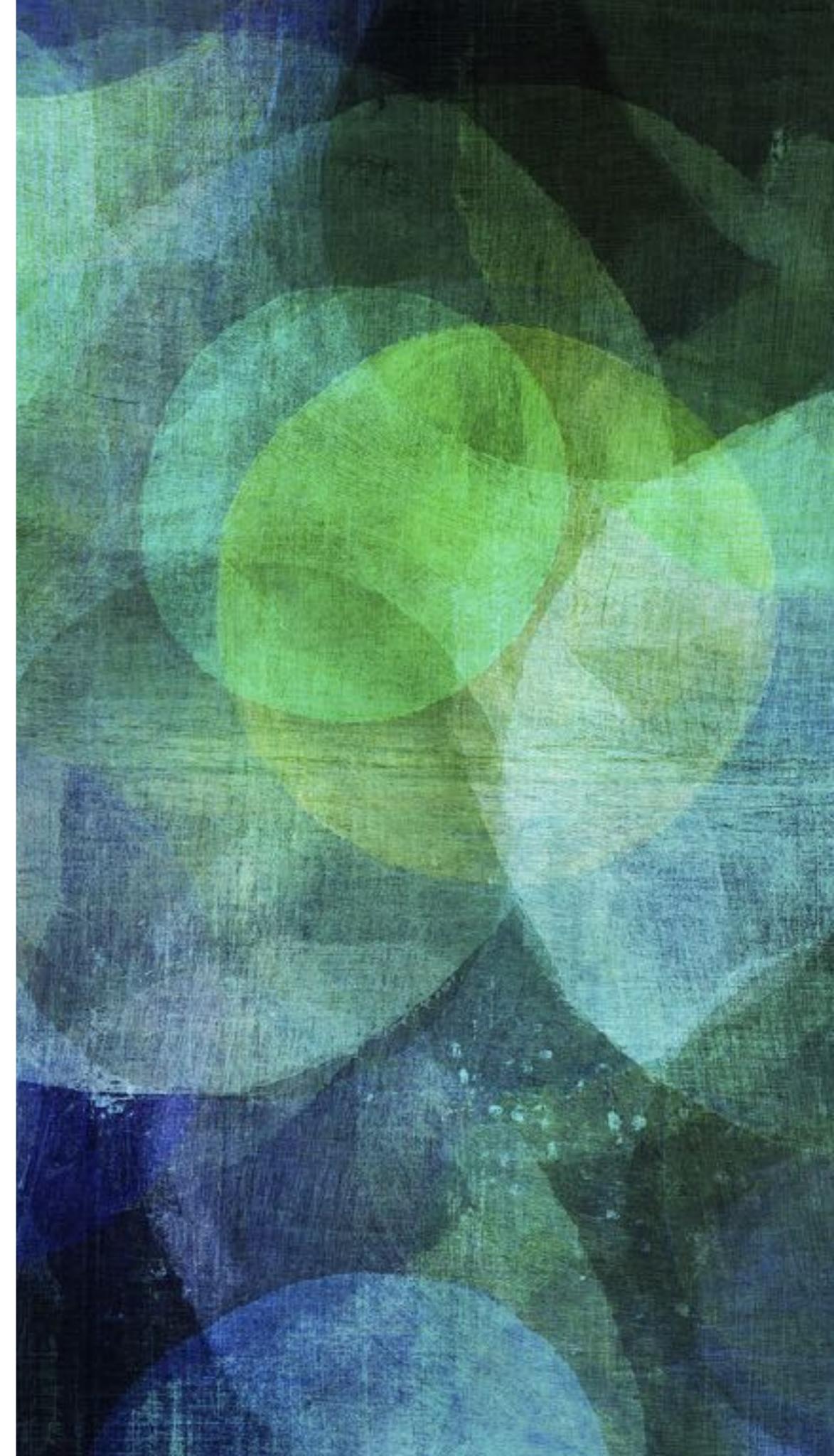
---

- ITERATIVE  
INTERACTION
  
- COMPARISON OF  
WHOLE & PART
  
- FUSION OF  
HORIZONS



# INTERPRETABLE MACHINE LEARNING

---



66

**Global ->**

**Model characteristics**

**Local ->**

**Individual predictions**

*Terms*

66

**Interpretability ->**

**Improving understanding**

**Explainability ->**

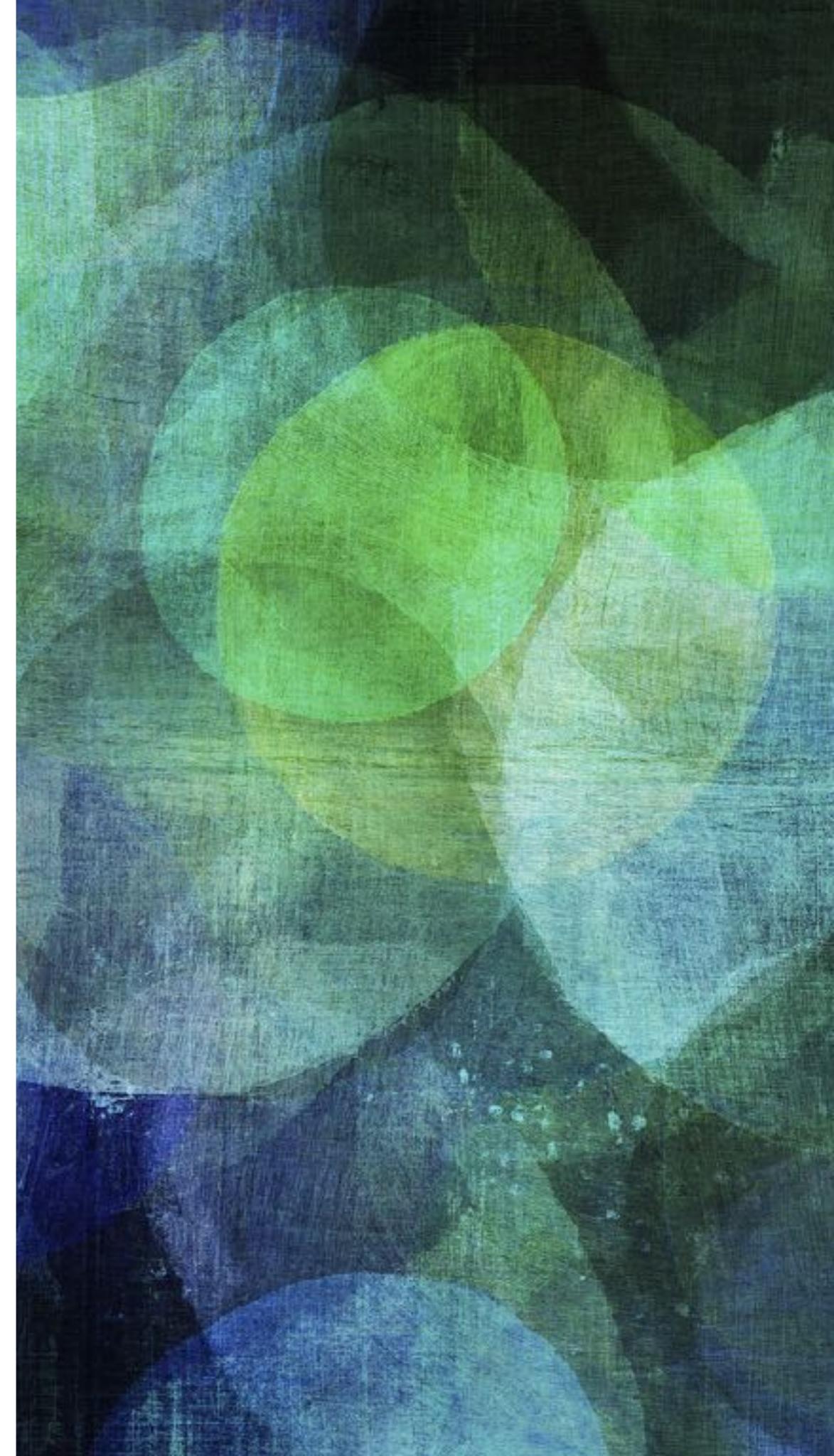
**Getting towards why**

*Terms*

# LIME

---

*Linearity in Complexity*



# *LIME slices (local; agnostic)*

*Generate data points based on training data*

*Compute complex model predictions from the generated data to find the ‘most useful features’*

*Fit a local linear model for the ‘most useful features’ and use the feature coefficients as reason codes*

# A few lines of code...

```
In [2]: # import lime tools
import lime
import lime.lime_tabular

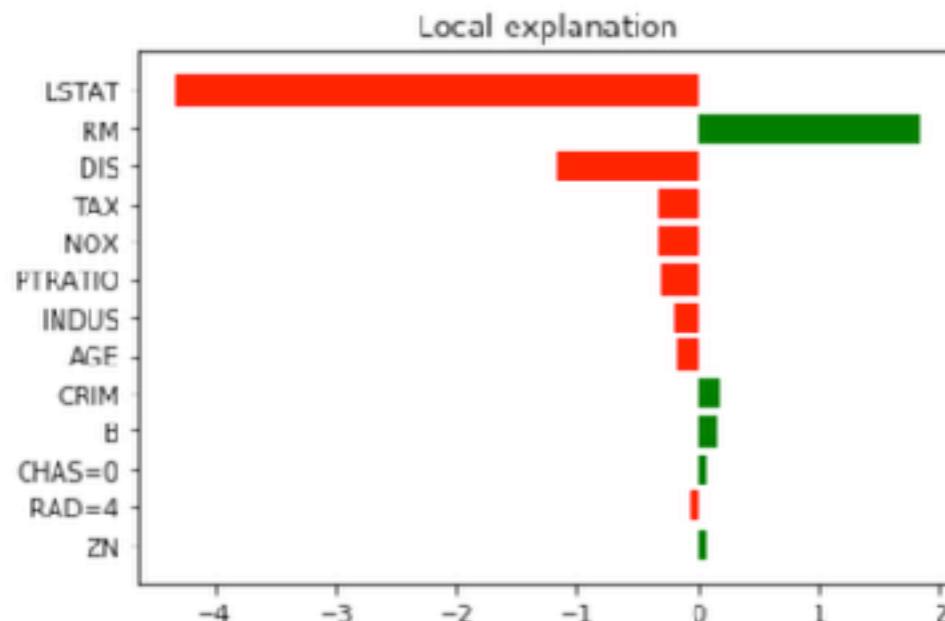
# generate an "explainer" object
categorical_features = np.argwhere(np.array([len(set(boston.data[:,x])) for x in range(boston.data.shape[1]))]) <= 10)
explainer = lime.lime_tabular.LimeTabularExplainer(train, feature_names=boston.feature_names, class_names=['price'], categ
```

```
In [3]: #generate an explanation
i = 13
exp = explainer.explain_instance(test[i], rf.predict, num_features=14)
```

```
In [4]: %matplotlib inline
fig = exp.as_pyplot_figure();
```



# *Limits of LIME*

*Local non-linearity and kernel width...*

*Linearity versus feature interactions?*

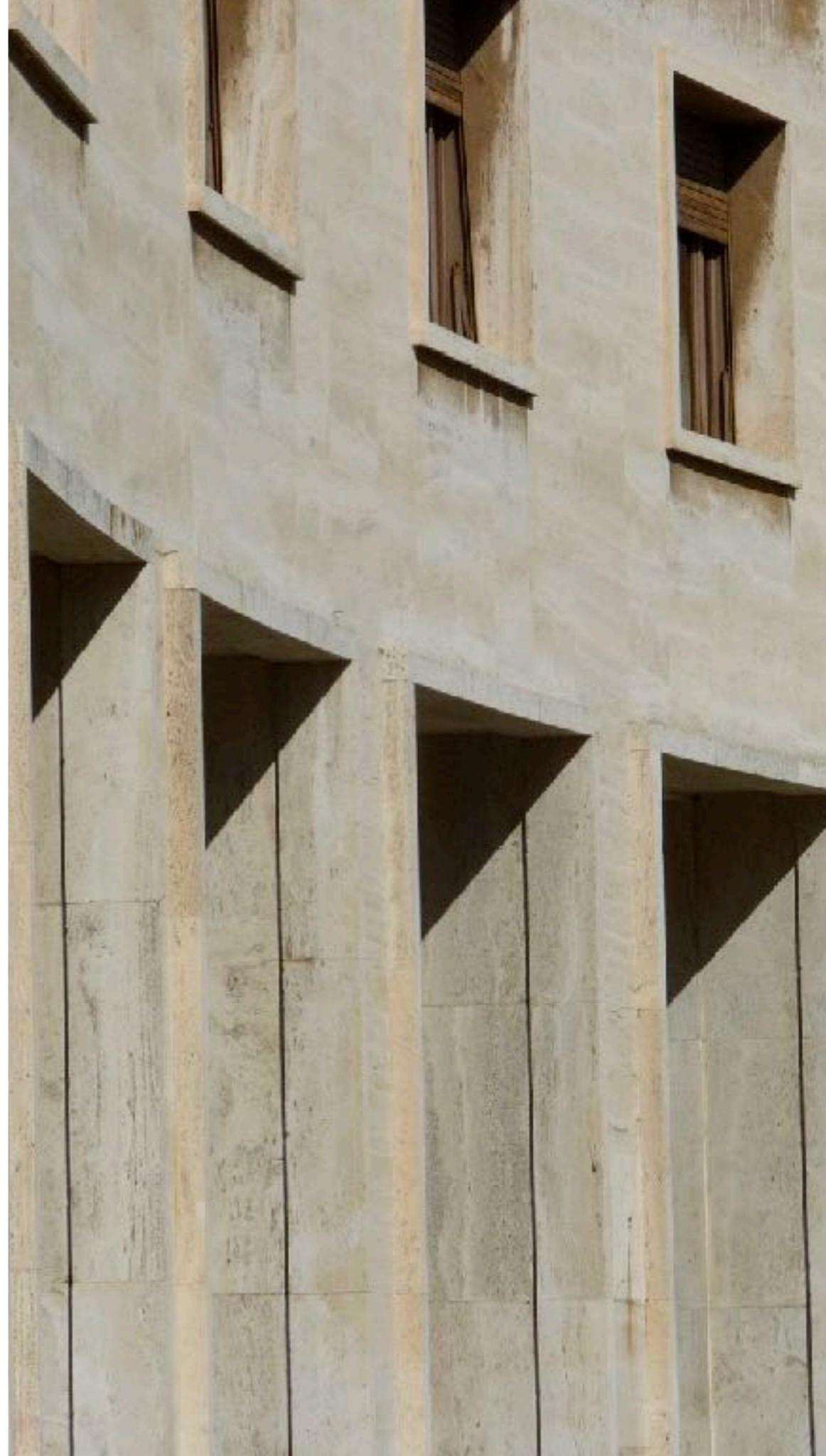
*Latency of results due to computation*

*Check the LIME prediction probability\**

# SKATER

---

*GENERATE TREES  
RULES  
LIME*

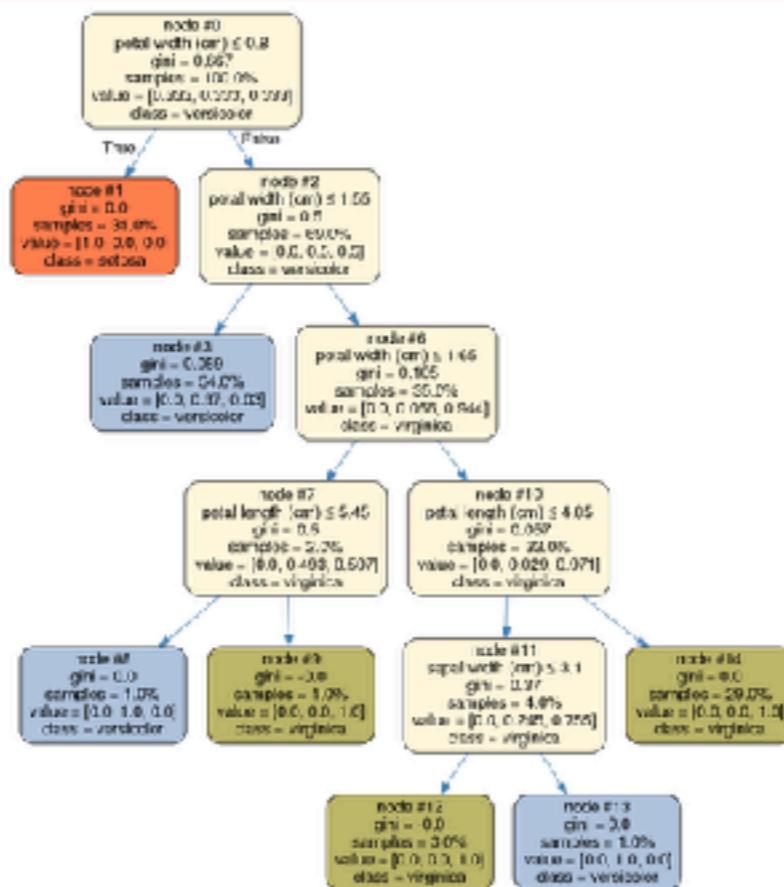


# Skater tree surrogate (regression)

# Skater tree surrogate (classification)

```
In [28]: surrogate_explainer2.plot_global_decisions(colors=['coral', 'lightsteelblue', 'darkkhaki'],
                                                file_name='simple_tree_post.png')
show_in_notebook('simple_tree_post.png', width=400, height=300)
```

2018-09-23 18:31:23,476 - skater.util.dataops - INFO - File Name: simple\_tree\_post.png



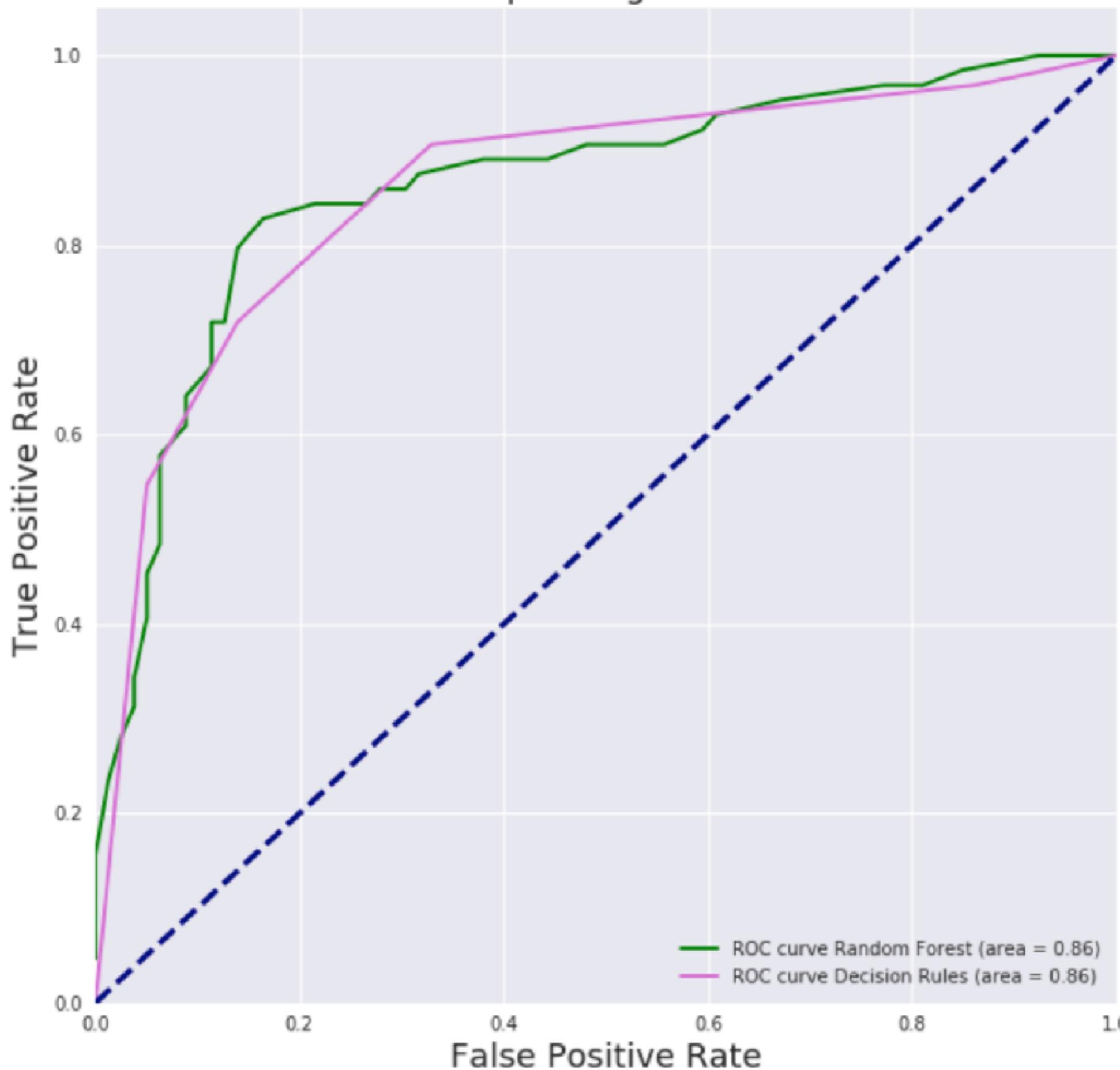
# *Skater Bayesian rules list (classification)*

```
In [31]: sbrl_big.print_model()
```

```
The rules list is :
```

```
If      {Sex_Encoded=1} (rule[276]) then positive probability = 0.18709677
else if {Pclass=3} (rule[237]) then positive probability = 0.32352941
else (default rule)  then positive probability = 0.94029851
```

## Receiver operating characteristic



# SHAP

---

*INTERPRETING TREES  
KERNELS  
GRADIENTS*



# *Cooperative Game Theory!*

*Shapley Values (approximations)*

*Calculate the value of a member to a group by comparing combinations of coalitions*

*Beware the effect of multicollinearity...*

# Tree SHAP(*local, groups*)

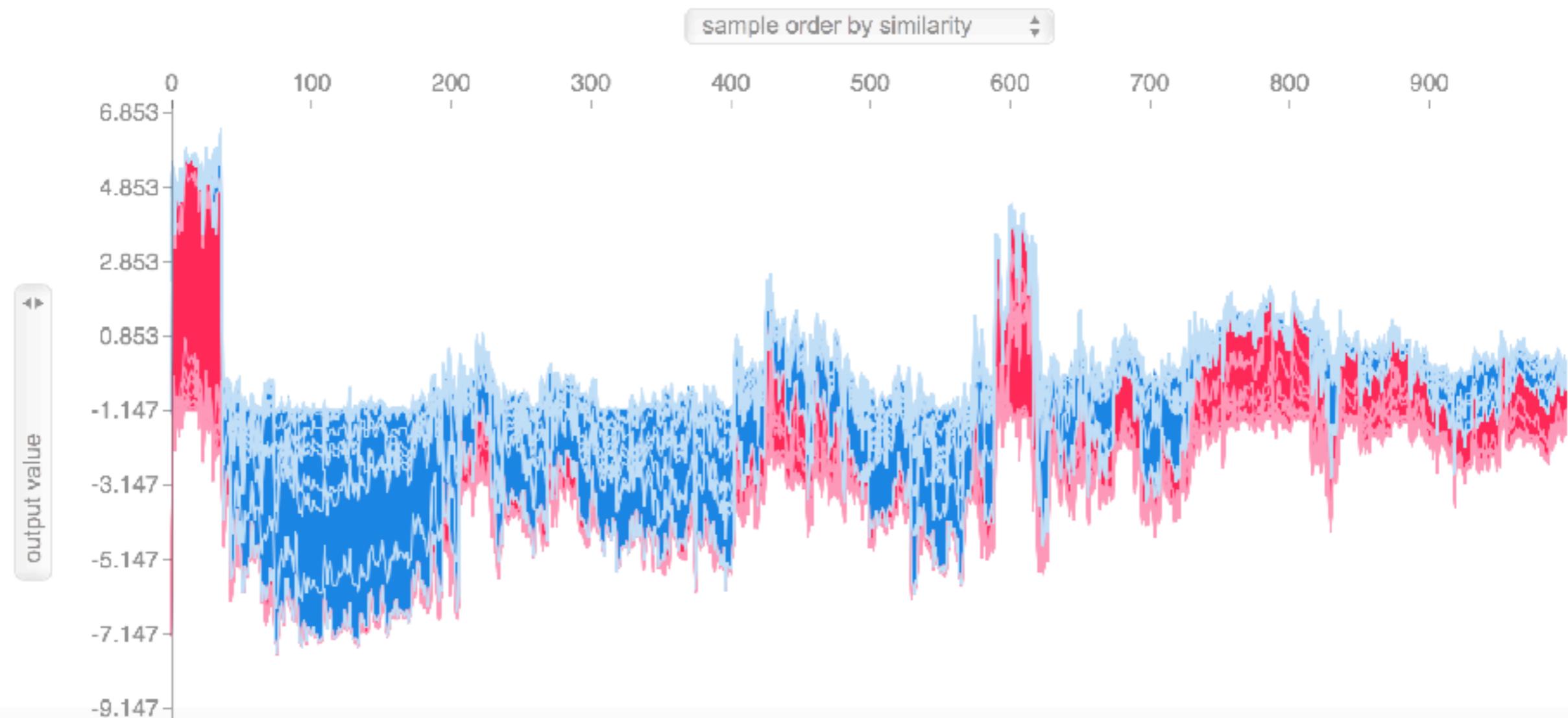
```
In [9]: shap.force_plot(explainer.expected_value, shap_values[0,:], X_display.iloc[0,:])
```

Out[9]:



```
In [10]: shap.force_plot(explainer.expected_value, shap_values[:1000,:], X_display.iloc[:1000,:])
```

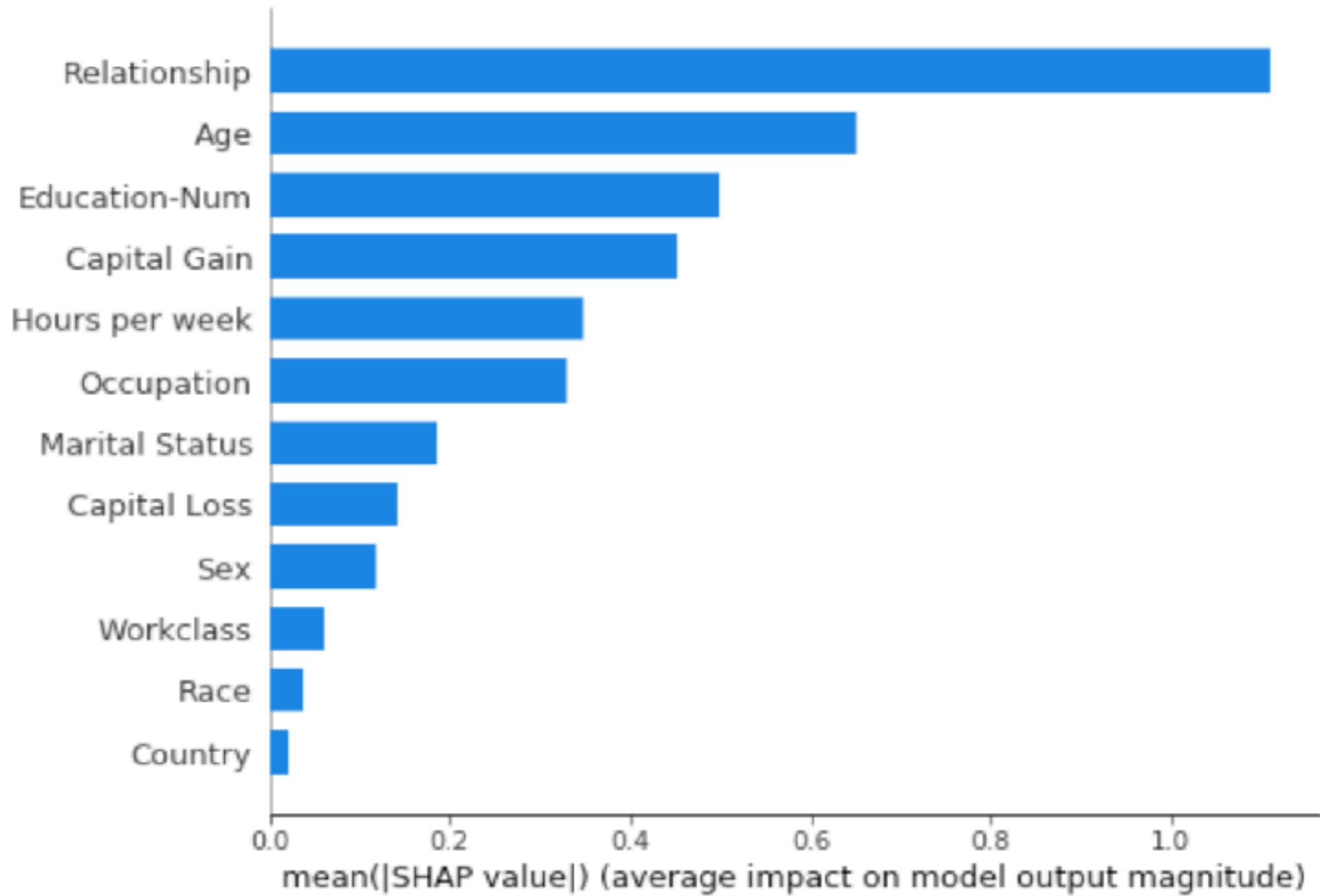
Out[10]:



# *Tree SHAP(global)*

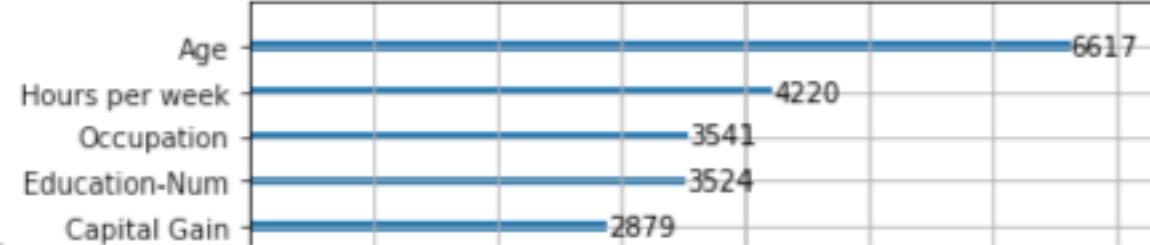
## *Variable Importance*

```
In [11]: shap.summary_plot(shap_values, x_display, plot_type="bar")
```

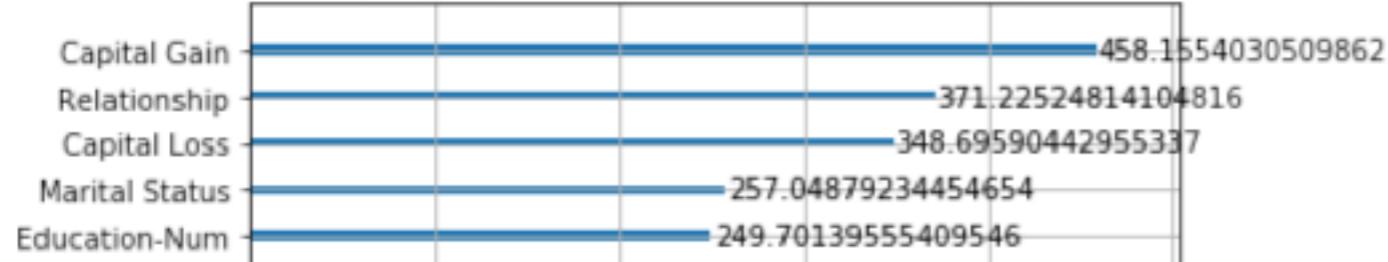


# Variable Importances (global...)

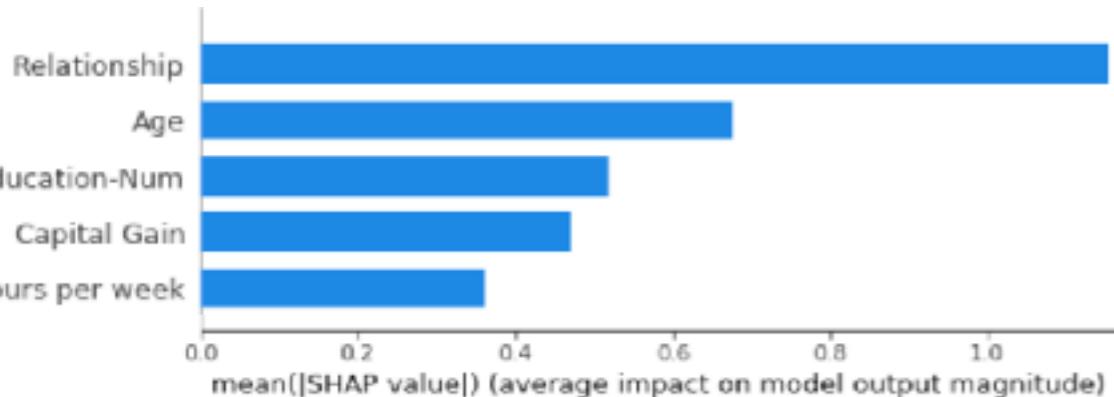
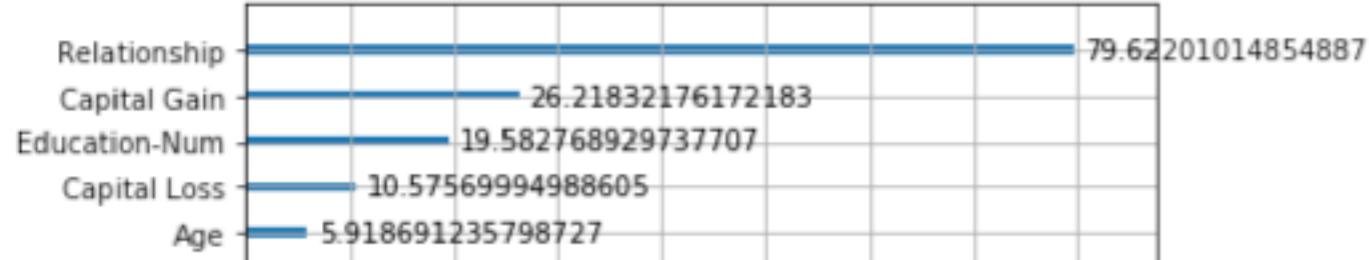
xgboost.plot\_importance(model)



xgboost.plot\_importance(model, importance\_type="cover")



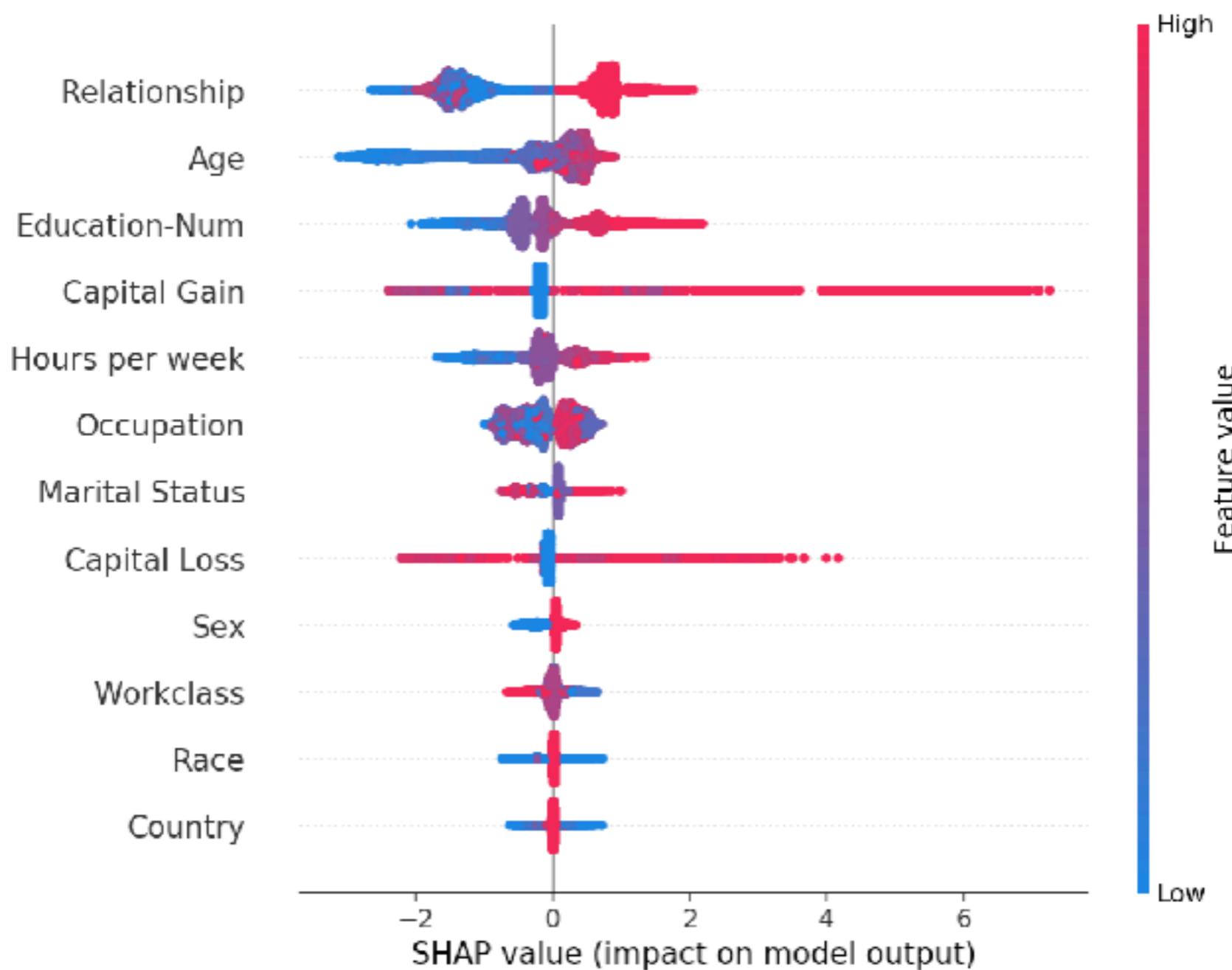
xgboost.plot\_importance(model, importance\_type="gain")



# *Tree SHAP (global and local)*

## *Overall and Individual*

```
In [9]: shap.summary_plot(shap_values, X)
```



**NOTE: XGBOOST/CATBOOST  
TREE SHAP INSIDE**

# CONTRASTIVE COUNTERFACTUAL

---



“

He who differentiates well teaches well

*-John Amos Comenius*

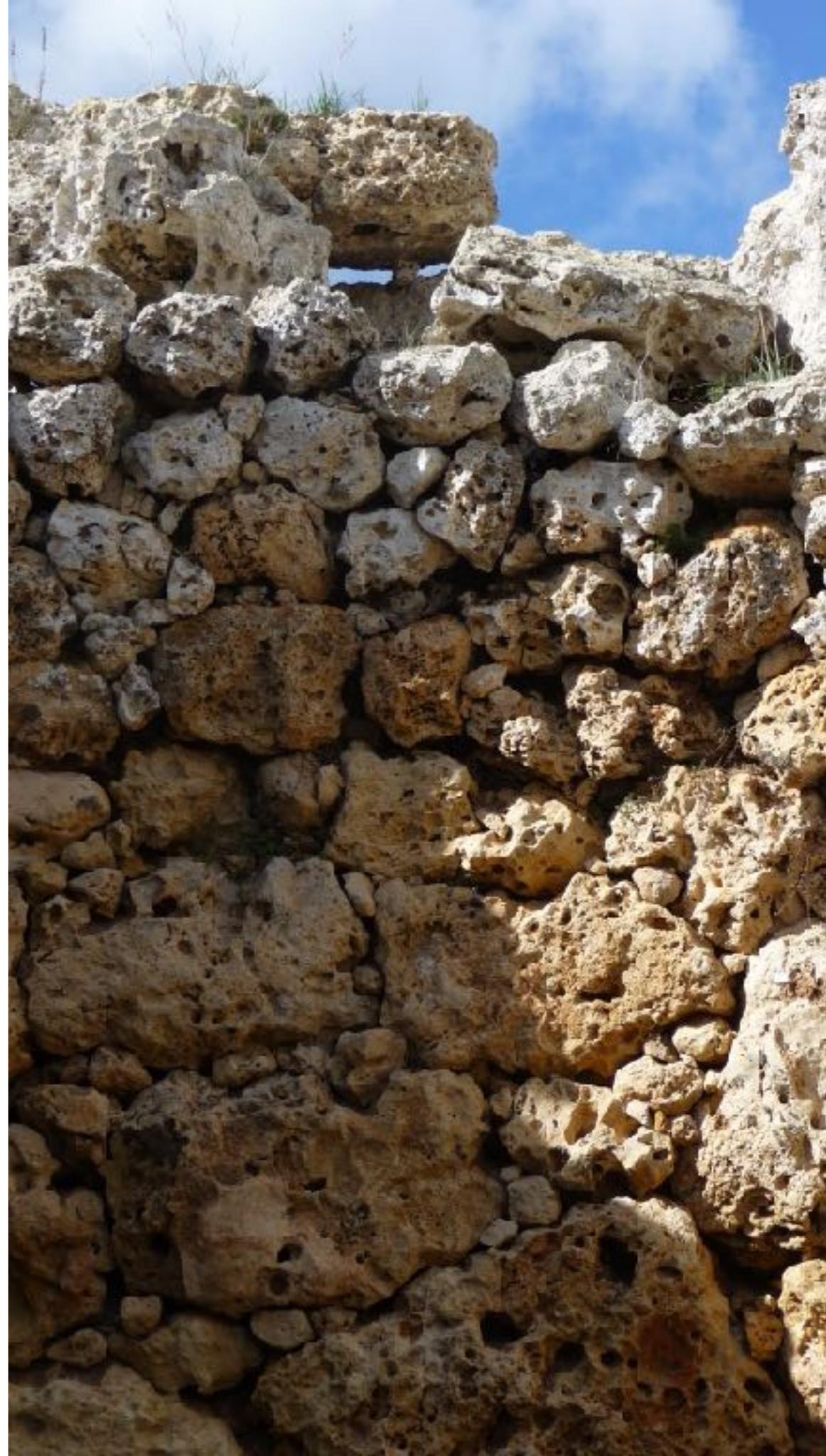
# *ContrastiveExplanation (local)*

The model predicted '92.414' instead of 'more than 92.414' because 's3 <= 0.1  
The model predicted '92.414' because 'bmi <= 0.006 and s5 <= -0.003 and s3 <= 0.1 s1 <= -0.058 and s1 > 0.023 and age > -0.03 and s4 > 0.016 and s2 <= -0.002

of 'more than 92.414' because 's3 <= 0.106 and bmi <= -0.021'",  
'bmi <= 0.006 and s5 <= -0.003 and s3 <= 0.02 and bp > 0.054 and s4 <= 0.024 > -0.03 and s4 > 0.016 and s2 <= -0.002 and s1 > 0.113 and bp <= -0.018'")

# PYBREAKDOWN

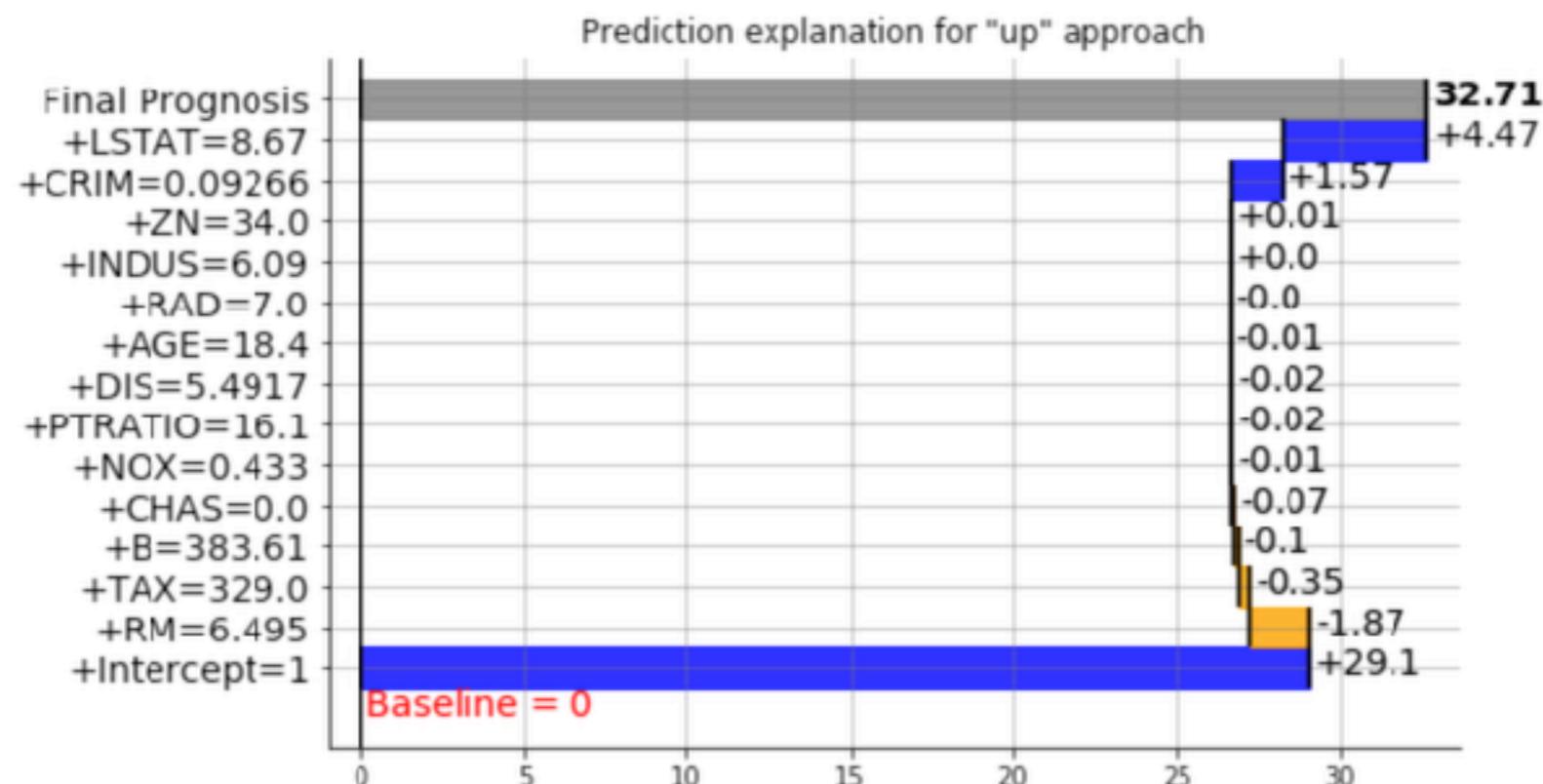
---



# *Breakdown (Global)*

```
In [7]: # visualisation
```

```
explanation.visualize(figsize=(8,5),dpi=100)
```



# RULEFIT

---

*linearity with interactions*



# *RuleFit*

*Fit a sparse linear model with the training set  
and additional features which are decision rules*

*The decision rules are generated by decision trees and reveal interactions  
between the original features*

# *RuleFit*

```
rf = RuleFit(tree_size=4, sample_fract='default', max_rules=2000,  
             memory_par=0.01,  
             tree_generator=None,  
             rfmode='classify', lin_trim_quantile=0.025,  
             lin_standardise=True, exp_rand_tree_size=True, random_state=1  
rf.fit(X, y_class, feature_names=features)
```

```
665  lstat > 19.775001525878906 & lstat <= 19.82999...      rule  8.465416  
1582 lstat <= 7.71999979019165 & lstat > 7.68499994...      rule -6.059011
```

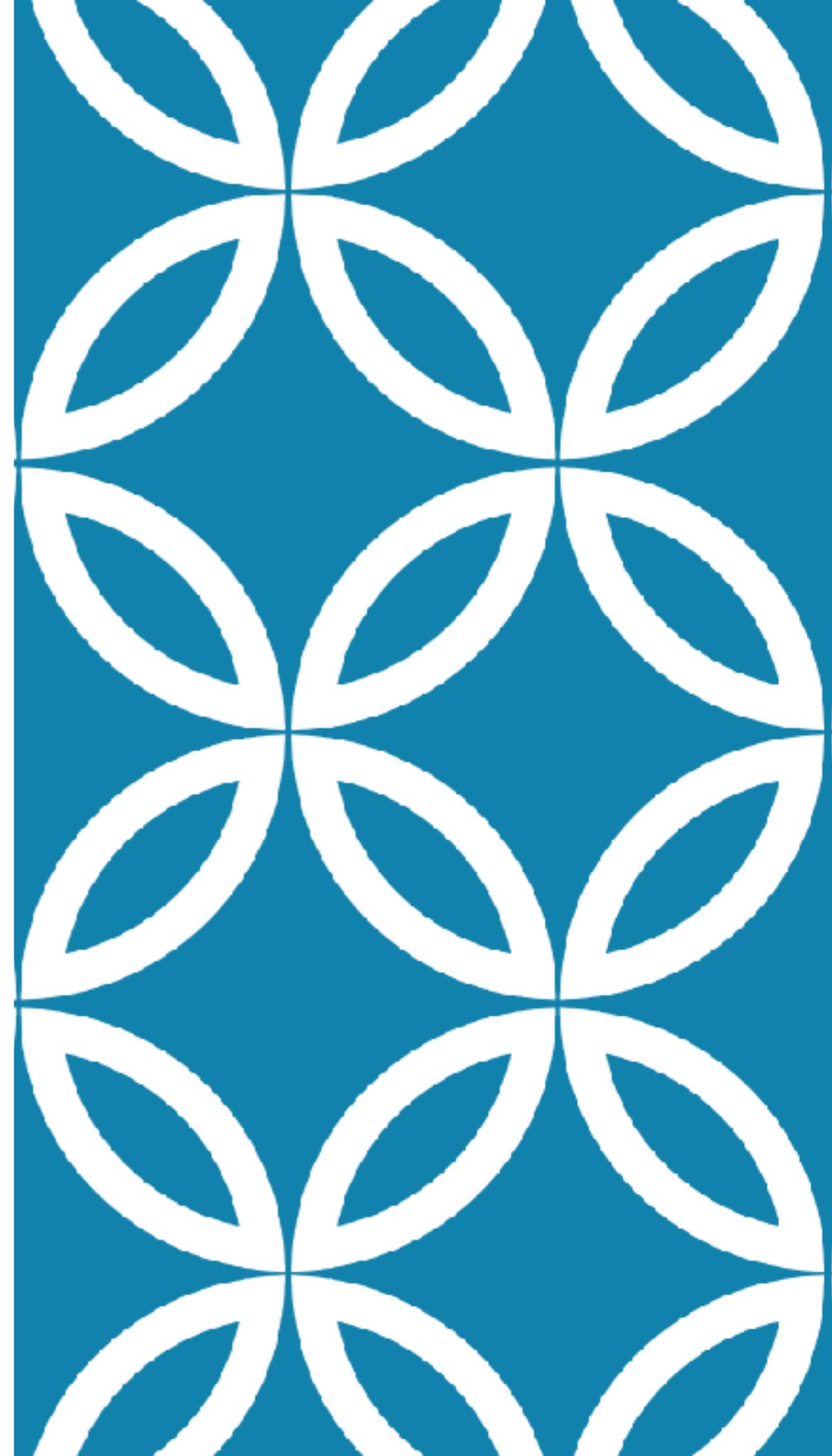
# NONCONFORMIST

---

*INFORMATION THEORETIC  
LOCAL EXPLANATIONS  
CONFORMAL*

L2X

LORE



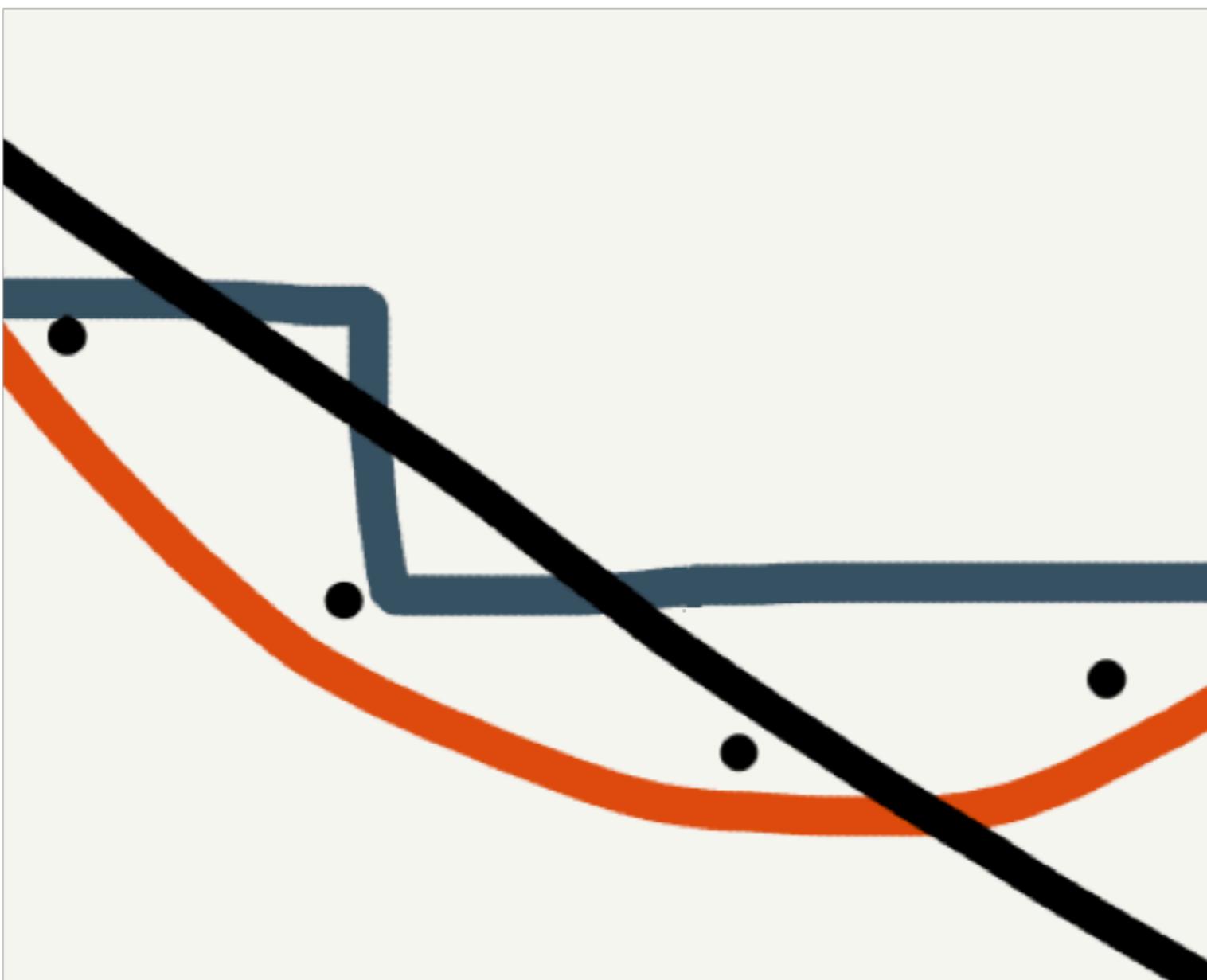
# XGBOOST

---

*Model Monotonicity*



# *Monotonic models*



# INTO THE DEEP

---

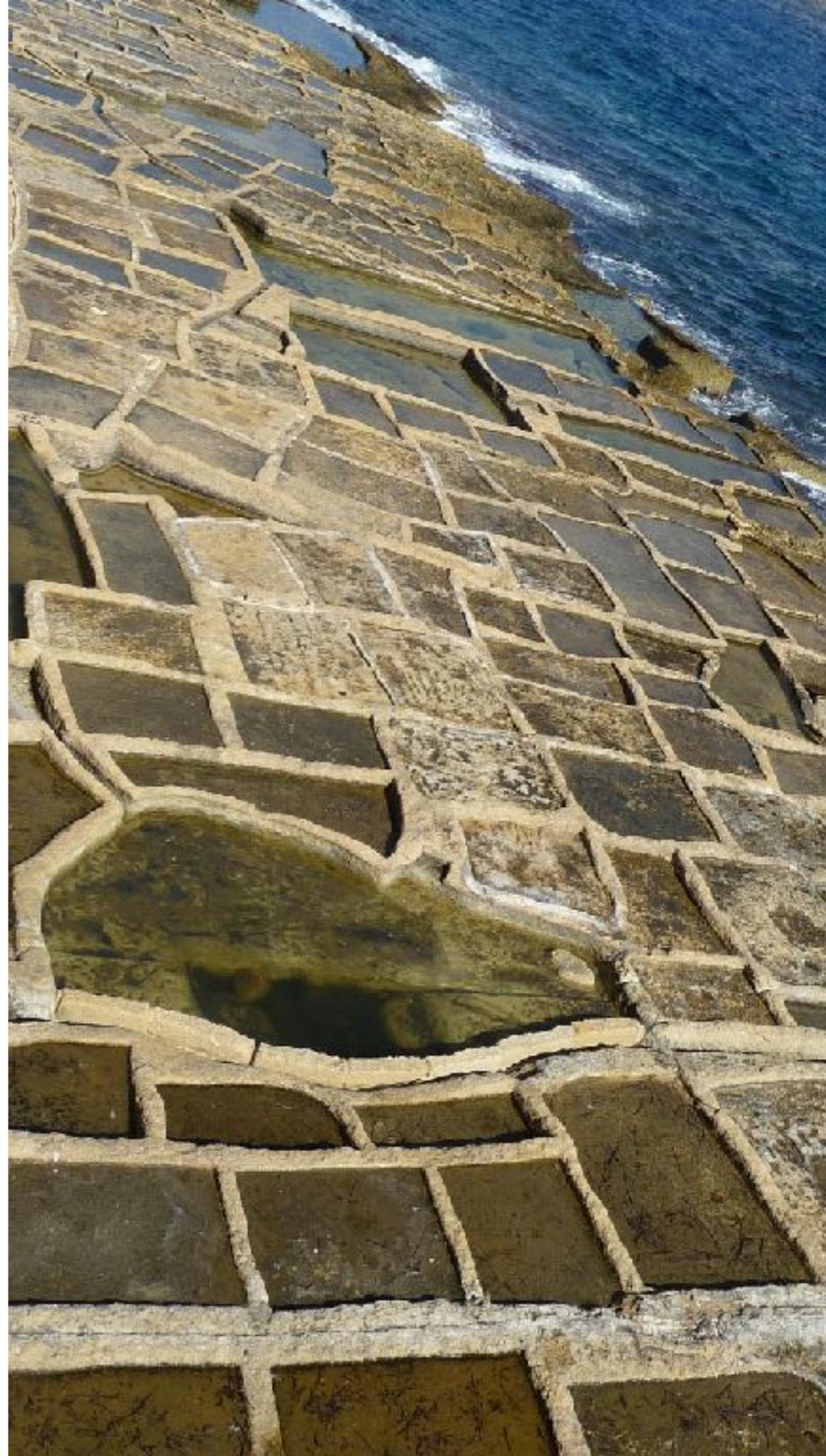
*'what about neural networks?'*



# TENSORFLOW LATTICE

---

*Model Monotonicity*



# DKNN

# DARKSIGHT

# INVESTIGATE!

*DEEP  $k$ -NEAREST NEIGHBOURS*

*MODEL DISTILLATION*

*DEEP TAYLOR DECOMPOSITION*



# *Deep Taylor Decomposition*

*Saliency*

*Sensitivity*

*Relevance\**

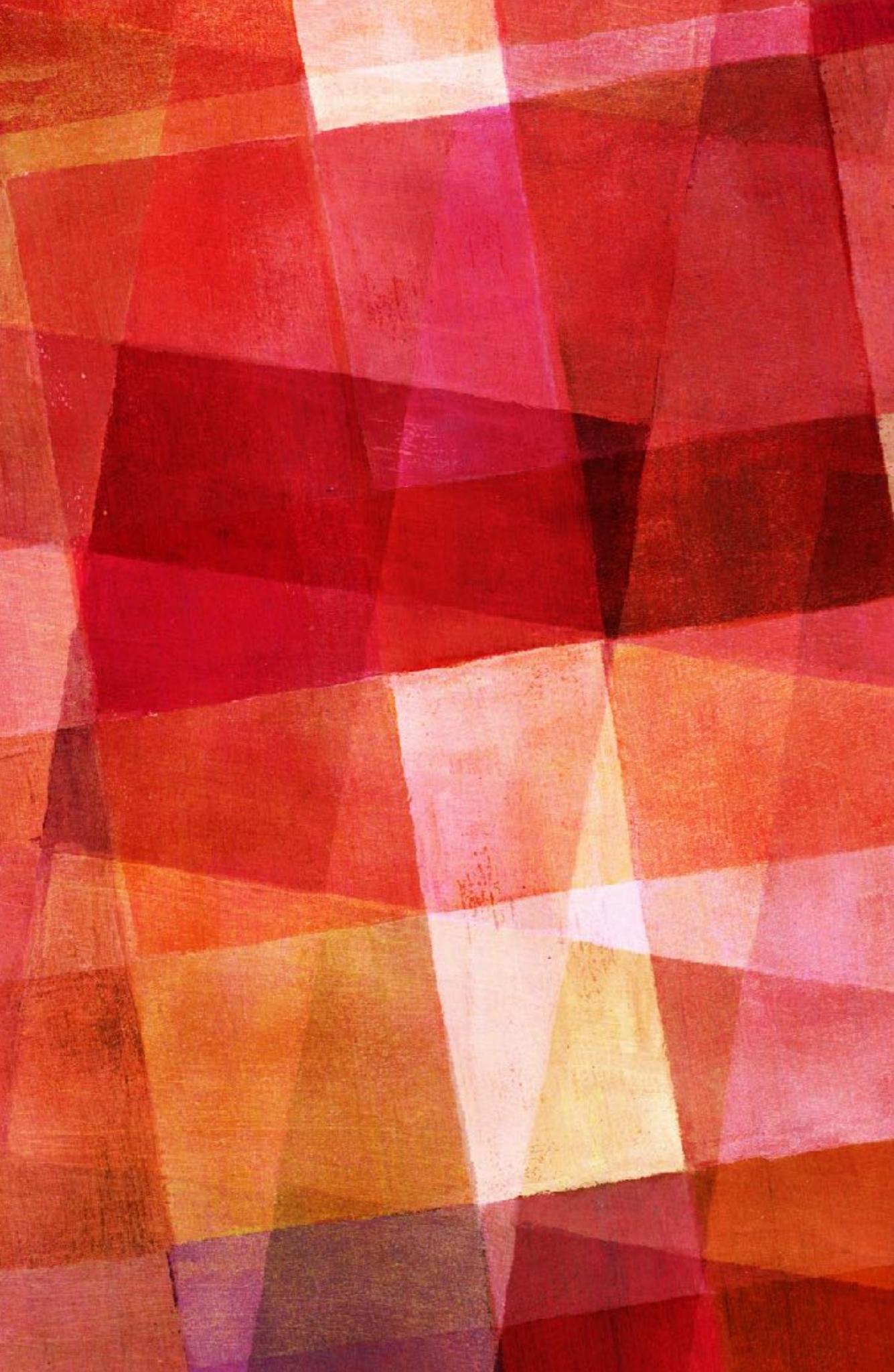
*To the reading of DL research papers  
there is no end...*



# VISUALISATION

---

*'biological neuron activation'*



# PLOTTING CONTRIBUTIONS

---

- **PDPBOX**
- **PARTIAL  
DEPENDENCE  
(AVERAGE)**
- **INDIVIDUAL  
CONDITIONAL  
EXPECTATION**



# RELATIONSHIPS

---

- YELLOWBRICK  
(DIAGNOSTICS)
  
- KEPLER MAPPER  
(TOPOLOGICAL  
DATA ANALYSIS)



# ANOTHER ANGLE

---

*'R there other options?'*





# PACKAGES

---

➤ **FFTREES**  
FAST AND  
FRUGAL TREES

➤ **GAMSEL**  
GENERAL  
ADDITIVE  
MODEL  
SELECTION



# PLOTS

---

- **DALEX**  
**(DESCRIPTIVE  
MACHINE  
LEARNING  
EXPLANATIONS)**
- ACCUMULATED  
LOCAL EFFECTS
- CETERIS PARIBUS

# ONE MORE THING

.....

*'All that coding...'*

# WHAT IF?

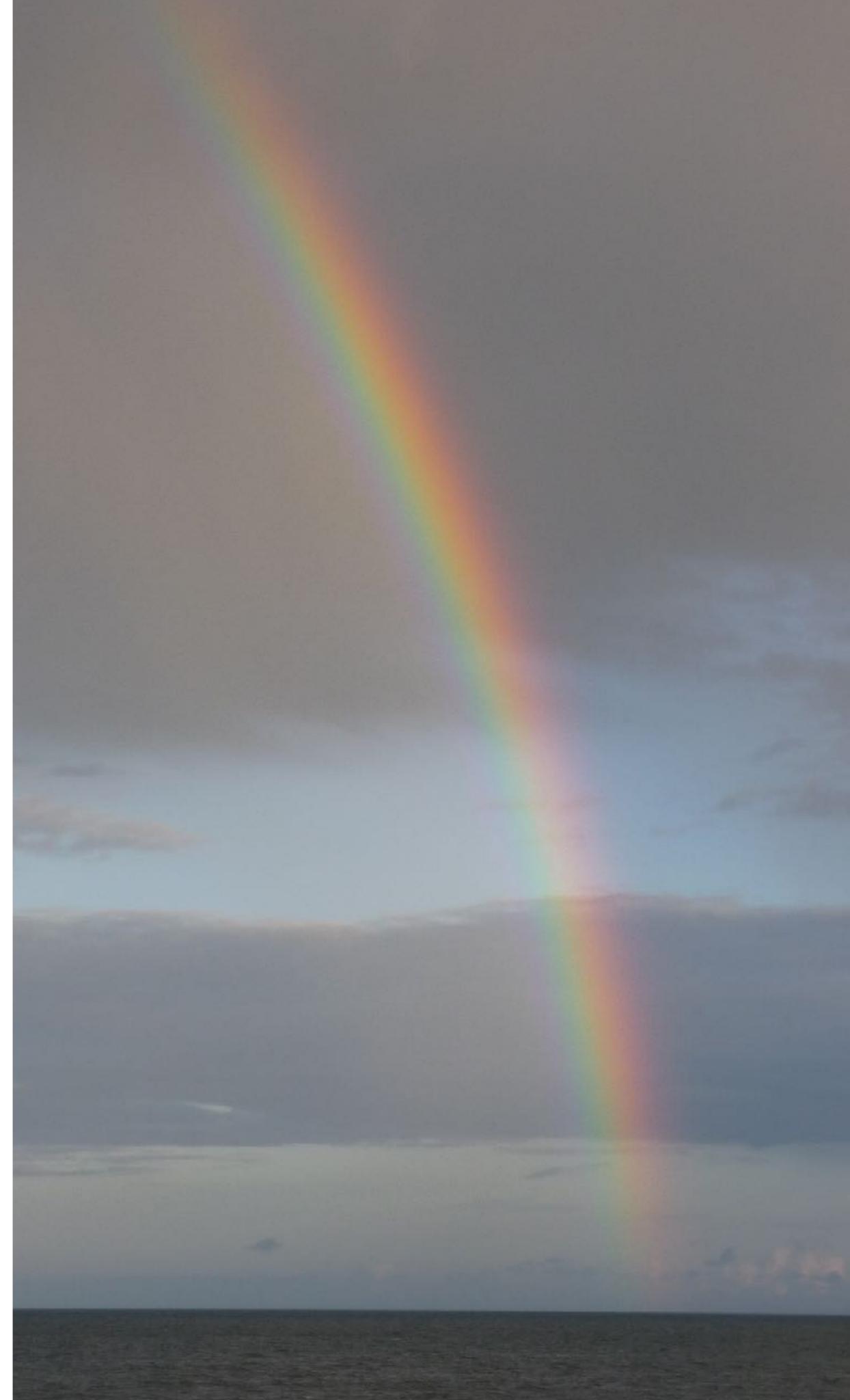
---

VIEW SLICES OF  
PREDICTIONS

VISUALISE INTERACTIVE  
EFFECTS

COMPARE COUNTERFACTUAL  
EXAMPLES

APPLY ALGORITHMIC  
FAIRNESS CONSTRAINTS



Partial dependence plots

Compute distance

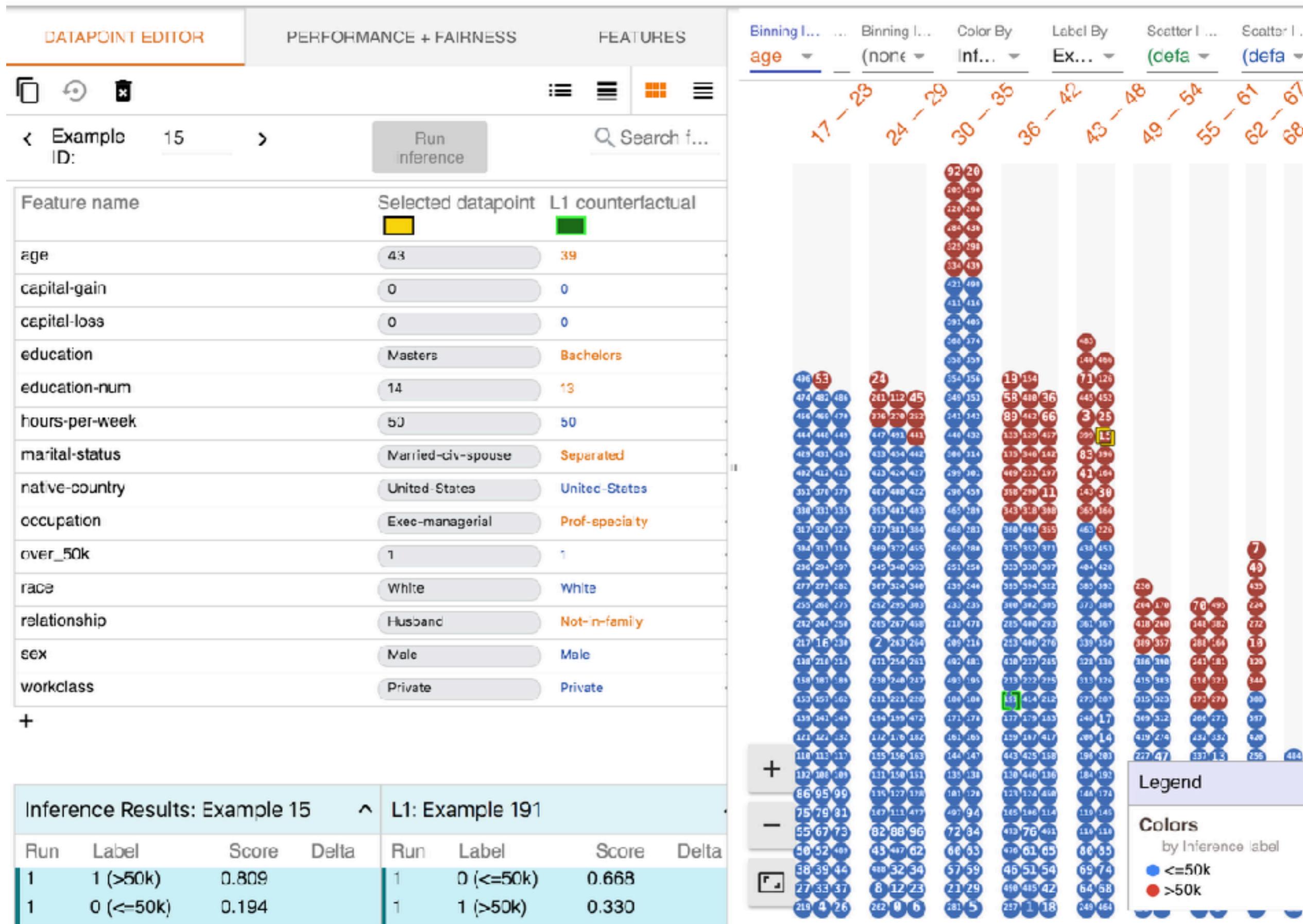
Show nearest counterfactual:

L1

L2

?

500 examples loaded





**STOP**

A red octagonal sign with the word "STOP" in large, white, sans-serif capital letters. The sign has a thick black border and is mounted on a light-colored wall.

**REVIEW**



# UNDERSTANDING VIA

---

➤ COMBINING  
MULTIPLE  
TECHNIQUES



# UNDERSTANDING VIA

---

- COMPARISON OF GLOBAL & LOCAL CONTRIBUTIONS
- COMBINING MULTIPLE TECHNIQUES



# UNDERSTANDING VIA

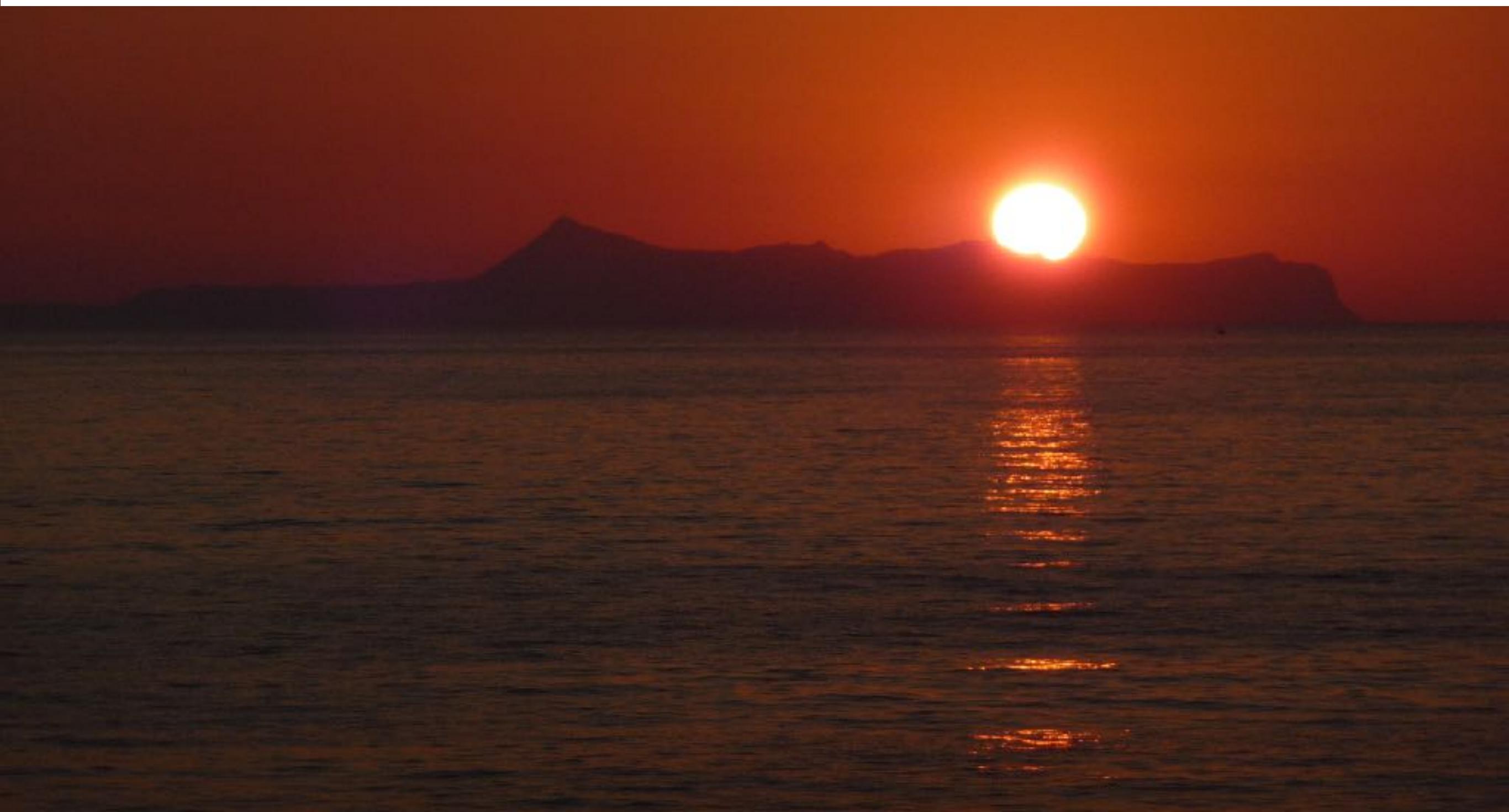
---

- ITERATIVE  
INTERACTION
- COMPARISON OF  
GLOBAL & LOCAL  
CONTRIBUTIONS
- COMBINING  
MULTIPLE  
TECHNIQUES

# HOW FAR DOES INTERPRETABILITY GO?

---

*Data, Preprocessing, Feature Engineering,  
Model Building, Model Selection, Drift Monitoring*



# INTERPRETABILITY - THE BEST OF TIMES?

---

*Model: Sensitivity? Fidelity? Specific? Implementation compatible?*

*Code: Robustness? Licensing? Scalability?*

*dean\_allsopp@hotmail.com*

→ iml ls

ContrastiveExplanation PDPbox  
LZX SHAP  
LORE SKATER

anchor  
defragTrees  
kepler-mapper

lime  
manual\_pdp\_ice  
monotonic

nonconformist  
pyBreakDown  
rulefit

Questions?

THANKS