

# Causality-based Explanation of Classification Outcomes

Leopoldo Bertossi, Jordan Li, Maximilian Schleich, **Dan Suciu**, Zografoula Vagena

2020 Workshop on Data Management for End-to-End Machine Learning (DEEM)

# Background

- Tackle interpretability of black-box models: We have seen LIME [1] and we also learned SHAP in the lectures
- This paper is about interpreting BBox models from the perspective of Causality

[1] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). " Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 1135-1144).

# Background : Causality

Let's review some of the concepts we discussed in class and also used in this paper.

- Counterfactual cause : remove/change a **feature** would flip the prediction

A score of measuring the degree to which a feature is affecting “flipping the prediction”

$$\text{COUNTER}(\mathbf{e}^\star, F_i) \stackrel{\text{def}}{=} L(\mathbf{e}^\star) - \mathbf{E} \left[ L(\mathbf{e}) | \mathbf{e}_{\mathcal{F} - \{F_i\}} = \mathbf{e}_{\mathcal{F} - \{F_i\}}^\star \right]$$

$\mathbf{e}^\star$ : an entity

$L(\mathbf{e}^\star)$ : outcome prediction by a model  $L$  (assume to be 1 in this paper)

$\mathbf{e}$ : any entity that is not  $\mathbf{e}^\star$  that shares the same feature values except  $F_i$

# Background : Causality

Counterfactual example: Loan Approval

| Name  | Education  | Credit Score | Place of Origin | Approved |    |
|-------|------------|--------------|-----------------|----------|----|
| Alice | Master's   | High         | China           | 1        | T1 |
| Bob   | Master's   | High         | USA             | 1        | T2 |
| Carol | Master's   | Low          | USA             | 0        | T3 |
| David | Master's   | Low          | USA             | 0        | T4 |
| Eve   | Master's   | High         | China           | 1        | T5 |
| Grace | Master's   | Low          | China           | 0        | T6 |
| Frank | Bachelor's | Middle       | USA             | 0        | T7 |

Why is Eve's loan approved by Model?  
(1: approved, 0: denied)

# Background : Causality

| Name  | Education I | Credit Score | Place of Origin | Approved |    |
|-------|-------------|--------------|-----------------|----------|----|
| Alice | Master's    | High         | China           | 1        | T1 |
| Bob   | Master's    | High         | USA             | 1        | T2 |
| Carol | Master's    | Low          | USA             | 0        | T3 |
| David | Master's    | Low          | USA             | 0        | T4 |
| Eve   | Master's    | High         | China           | 1        | T5 |
| Grace | Master's    | Low          | China           | 0        | T6 |
| Frank | Bachelor's  | Middle       | USA             | 0        | T7 |

$$\text{COUNTER}(\mathbf{e}^*, F_i) \stackrel{\text{def}}{=} L(\mathbf{e}^*) - \mathbf{E} \left[ L(\mathbf{e}) | \mathbf{e}_{\mathcal{F} - \{F_i\}} = \mathbf{e}_{\mathcal{F} - \{F_i\}}^* \right]$$

Candidate exp 1: POO=china

$$\text{COUNTER}(t5, \text{POO}) = 1 - (1+1)/2 = 0$$

# Background : Causality

Counterfactual example: Loan Approval

| Name  | Education  | Credit Score | Place of Origin | Approved |    |
|-------|------------|--------------|-----------------|----------|----|
| Alice | Master's   | High         | China           | 1        | T1 |
| Bob   | Master's   | High         | USA             | 1        | T2 |
| Carol | Master's   | Low          | USA             | 0        | T3 |
| David | Master's   | Low          | USA             | 0        | T4 |
| Eve   | Master's   | High         | China           | 1        | T5 |
| Grace | Master's   | Low          | China           | 0        | T6 |
| Frank | Bachelor's | Middle       | USA             | 0        | T7 |

$$\text{COUNTER}(\mathbf{e}^*, F_i) \stackrel{\text{def}}{=} L(\mathbf{e}^*) - \mathbf{E} \left[ L(\mathbf{e}) | \mathbf{e}_{\mathcal{F} - \{F_i\}} = \mathbf{e}_{\mathcal{F} - \{F_i\}}^* \right]$$

Candidate exp 2: CS=High

$$\text{COUNTER}(t5, \text{CS}) = 1 - (1+0)/2 = 0.5$$

# Background : Causality

- Actual Cause : A pair  $(Fi, v)$  is called an actual cause with contingency  $(\Gamma, w)$ , where  $\Gamma$  is a set of features and  $w$  is a set of values, if  $(Fi, v)$  is a counterfactual cause for  $e \star [\Gamma := w]$

Alice: 30 years old, bachelors -> Denied

Alice: 35 years old, bachelors -> Denied

Alice: 35 years old, masters -> Approved

We say “masters” is actual cause with contingency set (age=35 years old), the responsibility is 1/2 since  $|\Gamma| = 1$

# Background : Causality

$$\text{COUNTER}(\mathbf{e}^\star, F_i) \stackrel{\text{def}}{=} L(\mathbf{e}^\star) - \mathbf{E} \left[ L(\mathbf{e}) | \mathbf{e}_{\mathcal{F}-\{F_i\}} = \mathbf{e}_{\mathcal{F}-\{F_i\}}^\star \right]$$

$$\mathbf{e}' \stackrel{\text{def}}{=} \mathbf{e}^\star[\Gamma := \mathbf{w}]$$

$$\mathbf{e}'' \stackrel{\text{def}}{=} \mathbf{e}'[F_i := v]$$

Extend with contingency set and  
actual cause

$$\text{RESP}(\mathbf{e}^\star, F_i, \Gamma, \mathbf{w}) \stackrel{\text{def}}{=} \frac{L(\mathbf{e}') - \mathbf{E} \left[ L(\mathbf{e}'') | \mathbf{e}_{\mathcal{F}-\{F_i\}}'' = \mathbf{e}_{\mathcal{F}-\{F_i\}}' \right]}{1 + |\Gamma|}$$

$$\varphi_i(v) = \frac{1}{\text{number of players}} \sum_{\text{coalitions excluding } i} \frac{\text{marginal contribution of } i \text{ to coalition}}{\text{number of coalitions excluding } i \text{ of this size}}$$

- **SHAP-score is based on the shapley value we discussed in the lecture**
- **A refresher:**

Fix an entity  $\mathbf{e}^\star$  and a feature  $F_i$ . Let  $\pi$  be a permutation on the set of features  $\mathcal{F}$ ; in other words,  $\pi$  fixes a total order on the set of features. Denote by  $\pi^{<F_i}$  the set of features  $F_j$  that come before  $F_i$  in the order  $\pi$ ; similarly,  $\pi^{\leq F_i}$  denotes  $\pi^{<F_i} \cup \{F_i\}$ . The *contribution* of the feature  $F_i$  is defined as:

“Player” here is defined as feature

$$c(\mathbf{e}^\star, F_i, \pi) \stackrel{\text{def}}{=} \mathbf{E} \left[ L(\mathbf{e}) | \mathbf{e}_{\pi^{\leq F_i}} = \mathbf{e}_{\pi^{\leq F_i}}^\star \right] - \mathbf{E} \left[ L(\mathbf{e}) | \mathbf{e}_{\pi^{< F_i}} = \mathbf{e}_{\pi^{< F_i}}^\star \right]$$

$$\text{SHAP}(\mathbf{e}^\star, F_i) \stackrel{\text{def}}{=} \frac{1}{n!} \sum_{\pi} c(\mathbf{e}^\star, F_i, \pi)$$

# Background: Connection between SHAP and COUNTER

$S \subseteq F - \{F_i\}$ , and define the contribution of  $F_i$  w.r.t.  $S$  as

$$c'(e^*, F_i, S) \stackrel{\text{def}}{=} E \left[ L(e) | e_{S \cup \{F_i\}} = e_{S \cup \{F_i\}}^* \right] - E \left[ L(e) | e_S = e_S^* \right]$$

For a number  $0 \leq \ell < n$ , We define the SHAP-score at level  $\ell$  as

$$\text{SHAP}(e^*, F_i, \ell) \stackrel{\text{def}}{=} \frac{\ell!(n - \ell - 1)!}{n!} \sum_{S \in \binom{F - \{F_i\}}{\ell}} c'(e^*, F_i, S) \quad (2)$$

$$\text{SHAP}(e^*, F_i) = \sum_{\ell=0, n-1} \text{SHAP}(e^*, F_i, \ell)$$

$$\text{COUNTER}(e^*, F_i) \stackrel{\text{def}}{=} L(e^*) - E \left[ L(e) | e_{F - \{F_i\}} = e_{F - \{F_i\}}^* \right]$$

SHAP value at level  $(n-1)$  is actually directly related to COUNTER value

$$\cdot \text{SHAP}(e^*, F_i, n - 1) = \frac{1}{n} \text{COUNTER}(e^*, F_i).$$

# Enough formulas, let's start calculating!

- However, we need to address some other issues
  - We have a lot of conditional expectations in both RESP and SHAP, dataset at hand is just a sample.
  - **2 Probability spaces: approximate the actual distribution of data**
    - **Product space** : assuming independence (which is often not true in reality), prob for a tuple to appear in the distribution is

$$p(\langle x_1, \dots, x_n \rangle) \stackrel{\text{def}}{=} \prod_i p(F_i = x_i)$$

As if we have more data by replacing data itself with probabilities

## Actual Implementation

- Start by calculating COUNTER score
- If all the counterfactual scores are 0, then from bottom up enumerating contingency sets until find a solution

# Enough formulas, let's start calculating!

- Turns out SHAP calculation in product space is #P-Hard:
  - A reduction from SHAP-scores to probability calculation problem
- **2nd probability space:**
  - **The Empirical Distribution:**
    - lets just focus on the data we have **and only** the data we have (we don't claim any probability of existence of any unseen tuples)
    - This is more restrictive, but with one advantage: we can do early stopping in SHAP-score

$$\text{SHAP}(\mathbf{e}^\star, F_i, \ell) \stackrel{\text{def}}{=} \frac{\ell!(n - \ell - 1)!}{n!} \sum_{S \in \binom{\mathcal{F} - \{F_i\}}{\ell}} c'(\mathbf{e}^\star, F_i, S) \quad (2)$$

$$c(\mathbf{e}^\star, F_i, \pi) \stackrel{\text{def}}{=} \mathbf{E} \left[ L(\mathbf{e}) | \mathbf{e}_{\pi \leq F_i} = \mathbf{e}_{\pi \leq F_i}^\star \right] - \mathbf{E} \left[ L(\mathbf{e}) | \mathbf{e}_{\pi < F_i} = \mathbf{e}_{\pi < F_i}^\star \right]$$

As S size increases, it is more and more unlikely we will have any tuples left that fulfills

However, distribution model doesn't work well when computing RESP because of lack of data. Sigh

## 2 Probability spaces: a summary

### Product Space

- Generalize data set to “generate” more data
- Good to use in computing COUNTER and RESP
- Doesn't work well when dealing with SHAP

### Empirical Distribution

- Focus on the data and only the data
- Makes the SHAP feasible because of early stopping opportunities
- Doesn't work well when dealing with COUNTER and RESP because of limited data

# Experiment evaluation

- Based on the results + tradeoffs from 2 probability spaces:
  - RESP score : calculated on product space
  - SHAP score : calculated on empirical space
- Datasets:
  - FICO data : loan grant risk
  - Kaggle Credit Card Fraud: fraud or not fraud

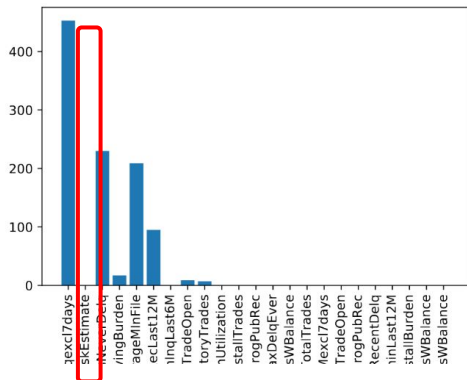
# FICO

- Compare RESP, SHAP and FICO scores from a logistic regression based model from [1] ([show demo briefly](#))
- FICO-score:
  - First find top M “subscales”
  - And find the top N features within each of the M subscales
  - Final sortorder is determined by sorting [subscale-score, feature-score-within-subscale]

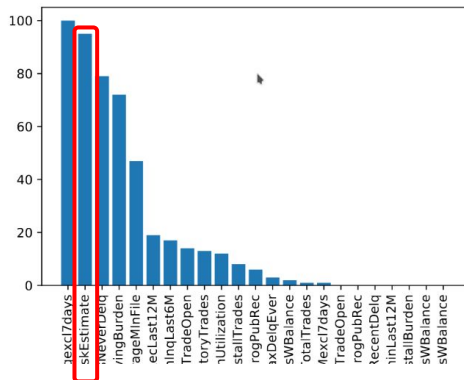
[1] Chaofan Chen, et al. An interpretable model with globally consistent explanations for credit risk. CoRR, abs/1811.12615, 2018

# FICO

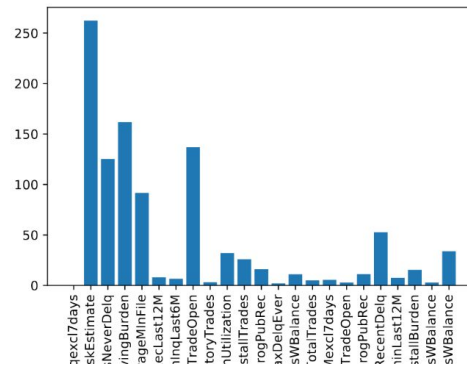
FICO-explanation



RESP-explanation



SHAP-explanation

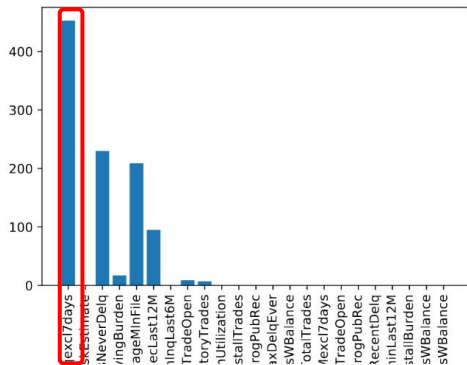


ExternalRiskEstimate

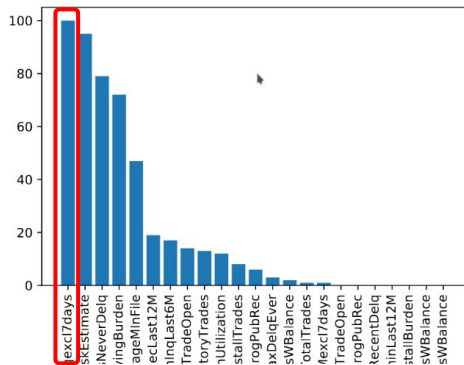
1. Due to the way FICO scores are sorted hierarchically
2. Model makes predictions on single cases and ignored the rest of the data
3. This also result in the fact that FICO-scores are less diverse (again, because of how the features are ranked)

# FICO

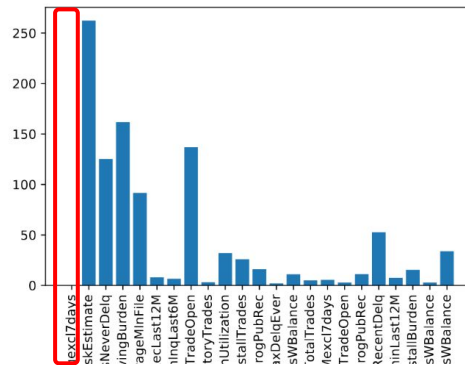
FICO-explanation



RESP-explanation



SHAP-explanation



MSinceMostRecentInqexcl7(MMR7)

Weights are head and tail heavy, MMR7 is actually evenly distributed among 2 classes in test data

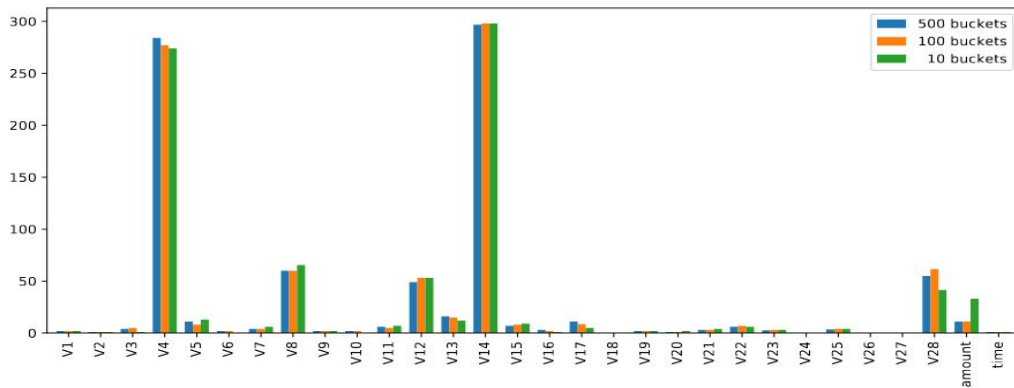
$$\text{SHAP}(\mathbf{e}^\star, F_i, \ell) \stackrel{\text{def}}{=} \frac{\ell!(n - \ell - 1)!}{n!} \sum_{S \in \binom{\mathcal{F} - \{F_i\}}{\ell}} c'(\mathbf{e}^\star, F_i, S) \quad (2)$$

$$\text{SHAP}(\mathbf{e}^\star, F_i) = \sum_{\ell=0, n-1} \text{SHAP}(\mathbf{e}^\star, F_i, \ell)$$

# Kaggle Credit Card Fraud

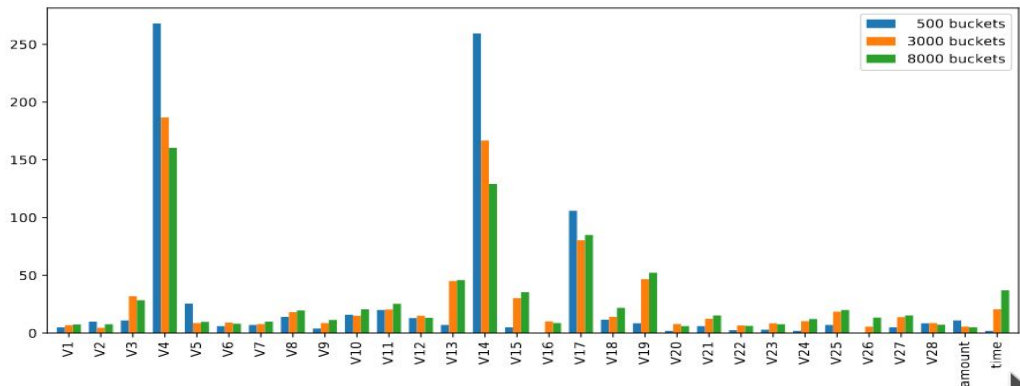
- Logistic regression as black-box-model
- Compare RESP and SHAP
- Feature values are also bucketized

# Kaggle Credit Card Fraud



RESP:

More buckets- longer runtime (why?)  
Insensitive to bucket size



SHAP:

More buckets- shorter runtime (why?)  
sensitive to bucket size

# Conclusion

- Investigated the feasibility and performance of using causality to approach explanation on black box models
  - Because of the complexity of naively computing scores, some compromises have to be made : 2 probability spaces
  - Experimentally show the results of RESP and SHAP and reasonably explained the reasons behind
- I am wondering practically speaking, which kind of explanation is more convincing for a real customer, thoughts?