

Evaluating and Aggregating Feature-based Model Explanations

Umang Bhatt

Carnegie Mellon University and University of Cambridge

Joint Work with Adrian Weller (Cambridge) and José Moura (CMU)

usb20@cam.ac.uk

Overview

- 1 Why are feature level explanations important (in medicine)?
- 2 What are existing feature-based explanation techniques?
- 3 How do we evaluate feature-based explanations?
- 4 How do we aggregate feature-based explanations?
- 5 Future Work

Medical Diagnostic Tasks

If we understand how a doctor reasons to a diagnosis, then we can build models that mimic that decision making.

Medical Diagnostic Tasks

If we understand how a doctor reasons to a diagnosis, then we can build models that mimic that decision making.

- Doctors leverage **vital signs** and **indicators** to diagnose

Medical Diagnostic Tasks

If we understand how a doctor reasons to a diagnosis, then we can build models that mimic that decision making.

- Doctors leverage **vital signs** and **indicators** to diagnose
 - This translates to a semantically meaningful feature vector x

Medical Diagnostic Tasks

If we understand how a doctor reasons to a diagnosis, then we can build models that mimic that decision making.

- Doctors leverage **vital signs** and **indicators** to diagnose
 - This translates to a semantically meaningful feature vector x
- Doctors leverage insights from **past patients** to better diagnose current patients

Medical Diagnostic Tasks

If we understand how a doctor reasons to a diagnosis, then we can build models that mimic that decision making.

- Doctors leverage **vital signs** and **indicators** to diagnose
 - This translates to a semantically meaningful feature vector x
- Doctors leverage insights from **past patients** to better diagnose current patients
 - This translates to influential points in the training distribution

Medical Diagnostic Tasks

If we understand how a doctor reasons to a diagnosis, then we can build models that mimic that decision making.

- Doctors leverage **vital signs** and **indicators** to diagnose
 - This translates to a semantically meaningful feature vector x
- Doctors leverage insights from **past patients** to better diagnose current patients
 - This translates to influential points in the training distribution

Medical Diagnostic Tasks

If we understand how a doctor reasons to a diagnosis, then we can build models that mimic that decision making.

- Doctors leverage **vital signs** and **indicators** to diagnose
 - This translates to a semantically meaningful feature vector x
- Doctors leverage insights from **past patients** to better diagnose current patients
 - This translates to influential points in the training distribution

We ask: can we validate the above intuition on a trained predictor, f , using feature importance and sample importance?

Common feature-based explanations

SHAP (Lundberg and Lee. NeurIPS 2017)

$$g_{ci} = g(\mathbf{f}, \mathbf{x})_i = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (\mathbf{f}(\mathbf{x}_{S \cup \{i\}}) - \mathbf{f}(\mathbf{x}_S))$$

Provides a “fair” distribution of contribution over all features, since Shapley values satisfy **efficiency**, symmetry, additivity, and dummy (zero).

Common feature-based explanations

SHAP (Lundberg and Lee. NeurIPS 2017)

$$\mathbf{g}_{ci} = \mathbf{g}(\mathbf{f}, \mathbf{x})_i = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (\mathbf{f}(\mathbf{x}_{S \cup \{i\}}) - \mathbf{f}(\mathbf{x}_S))$$

Provides a “fair” distribution of contribution over all features, since Shapley values satisfy **efficiency**, symmetry, additivity, and dummy (zero).

Integrated Gradients (Sundarajan et al. ICML 2017)

Accumulates the gradients along a straight line path between \mathbf{x} and $\bar{\mathbf{x}}$, where $\mathbf{f}(\bar{\mathbf{x}}) \approx 0$, and satisfies **completeness**, $\sum_{i=1}^d \mathbf{g}(\mathbf{f}, \mathbf{x})_i = \mathbf{f}(\mathbf{x}) - \mathbf{f}(\bar{\mathbf{x}})$.

Common feature-based explanations

SHAP (Lundberg and Lee. NeurIPS 2017)

$$\mathbf{g}_{ci} = \mathbf{g}(\mathbf{f}, \mathbf{x})_i = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (\mathbf{f}(\mathbf{x}_{S \cup \{i\}}) - \mathbf{f}(\mathbf{x}_S))$$

Provides a “fair” distribution of contribution over all features, since Shapley values satisfy **efficiency**, symmetry, additivity, and dummy (zero).

Integrated Gradients (Sundarajan et al. ICML 2017)

Accumulates the gradients along a straight line path between \mathbf{x} and $\bar{\mathbf{x}}$, where $\mathbf{f}(\bar{\mathbf{x}}) \approx 0$, and satisfies **completeness**, $\sum_{i=1}^d \mathbf{g}(\mathbf{f}, \mathbf{x})_i = \mathbf{f}(\mathbf{x}) - \mathbf{f}(\bar{\mathbf{x}})$.

LIME (Ribeiro et al. KDD 2016)

$$\mathbf{g}(\mathbf{f}, \mathbf{x})_i = \arg \min_{\mathbf{g} \in \mathcal{G}} \mathcal{L}(\mathbf{f}, \mathbf{g}, \pi_{\mathbf{x}}) + \Omega(\mathbf{g})$$

Local surrogate model, \mathbf{g} , to approximate original model, f , in some kernelized region $\pi_{\mathbf{x}}$, and encourages **sparsity** by keeping model complexity, $\Omega(\mathbf{g})$, low

Evaluating explanations

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

$$\mu_{\text{AVG}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \int_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z}$$

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

$$\mu_{\text{AVG}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \int_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z}$$

$$\mu_{\text{MAX}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \max_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z}))$$

Let D be the distance between explanations and ρ be the distance between inputs

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

$$\mu_{\text{AVG}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \int_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z}$$

$$\mu_{\text{MAX}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \max_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z}))$$

Let D be the distance between explanations and ρ be the distance between inputs

Faithfulness

Does the explanation capture features important to the prediction?

Evaluating explanations

Sensitivity

Do similar inputs have similar explanations?

$$\mu_{\text{AVG}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \int_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z})) \mathbb{P}_{\mathbf{x}}(\mathbf{z}) d\mathbf{z}$$

$$\mu_{\text{MAX}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, r) = \max_{\rho(\mathbf{x}, \mathbf{z}) \leq r} D(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{z}))$$

Let D be the distance between explanations and ρ be the distance between inputs

Faithfulness

Does the explanation capture features important to the prediction?

$$\mu_{\text{F}}(\mathbf{f}, \mathbf{g}, \mathbf{x}, S) = \text{corr}\left(\frac{1}{|S|} \sum_{i \in S} \mathbf{g}(\mathbf{f}, \mathbf{x})_i, \mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_{[x_s = \bar{x}_s]})\right)$$

Fix a subset size and randomly sample subsets of that size from \mathbf{x} to estimate the Pearson Correlation Coefficient

Evaluating explanations (cont.)

Evaluating explanations (cont.)

Complexity

Is the explanation digestible?

Evaluating explanations (cont.)

Complexity

Is the explanation digestible? We define an attribution contribution distribution:

$$\mathbb{P}_A = \left\{ \frac{|\mathbf{g}(\mathbf{x})_1|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|}, \frac{|\mathbf{g}(\mathbf{x})_2|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|}, \dots, \frac{|\mathbf{g}(\mathbf{x})_d|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|} \right\}$$

$$\mu_C(\mathbf{f}, \mathbf{g}, \mathbf{x}) = H(\mathbf{x}) = \mathbb{E}_i \left[-\ln(\mathbb{P}_A(i)) \right] = -\sum_{i=1}^d \mathbb{P}_A(i) \ln(\mathbb{P}_A(i))$$

Evaluating explanations (cont.)

Complexity

Is the explanation digestible? We define an attribution contribution distribution:

$$\mathbb{P}_A = \left\{ \frac{|\mathbf{g}(\mathbf{x})_1|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|}, \frac{|\mathbf{g}(\mathbf{x})_2|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|}, \dots, \frac{|\mathbf{g}(\mathbf{x})_d|}{\sum_{j \in [d]} |\mathbf{g}(\mathbf{x})_j|} \right\}$$

$$\mu_C(\mathbf{f}, \mathbf{g}, \mathbf{x}) = H(\mathbf{x}) = \mathbb{E}_i \left[-\ln(\mathbb{P}_A(i)) \right] = -\sum_{i=1}^d \mathbb{P}_A(i) \ln(\mathbb{P}_A(i))$$

The least complex explanation is one where $\mathbf{g}(\mathbf{x})_i = 1$ and the most complex explanation is one where $\mathbf{g}(\mathbf{x})_i = \frac{1}{d}$.

Aggregating Existing Techniques

Can we learn an aggregate explanation of existing techniques that does better with respect to a criterion of interest? An approach to study \mathbf{g}_{agg} can be to set the problem up as follows:

$$\mathbf{g}_{\text{agg}} = \arg \max_{\mathbf{g} \in \mathcal{G}} \mu(\mathbf{f}, \mathbf{g}), \text{ s.t. } \mathbf{g} = h(\mathcal{G}_m)$$

Aggregating Existing Techniques

Can we learn an aggregate explanation of existing techniques that does better with respect to a criterion of interest? An approach to study \mathbf{g}_{agg} can be to set the problem up as follows:

$$\mathbf{g}_{agg} = \arg \max_{\mathbf{g} \in \mathcal{G}} \mu(\mathbf{f}, \mathbf{g}), \text{ s.t. } \mathbf{g} = h(\mathcal{G}_m)$$

Three candidate methods for $h(\cdot)$.

- Convex Combination: $\mathbf{g}_{agg} = w\mathbf{g}_1 + (1 - w)\mathbf{g}_2$

Aggregating Existing Techniques

Can we learn an aggregate explanation of existing techniques that does better with respect to a criterion of interest? An approach to study \mathbf{g}_{agg} can be to set the problem up as follows:

$$\mathbf{g}_{agg} = \arg \max_{\mathbf{g} \in \mathcal{G}} \mu(\mathbf{f}, \mathbf{g}), \text{ s.t. } \mathbf{g} = h(\mathcal{G}_m)$$

Three candidate methods for $h(\cdot)$.

- Convex Combination: $\mathbf{g}_{agg} = w\mathbf{g}_1 + (1 - w)\mathbf{g}_2$
- Centroid Aggregation: $\mathbf{g}_{agg} \in \arg \min_{\mathbf{g} \in \mathcal{G}} \sum_{i=1}^m d(\mathbf{g}, \mathbf{g}_i)$

Aggregating Existing Techniques

Can we learn an aggregate explanation of existing techniques that does better with respect to a criterion of interest? An approach to study \mathbf{g}_{agg} can be to set the problem up as follows:

$$\mathbf{g}_{agg} = \arg \max_{\mathbf{g} \in \mathcal{G}} \mu(\mathbf{f}, \mathbf{g}), \text{ s.t. } \mathbf{g} = h(\mathcal{G}_m)$$

Three candidate methods for $h(\cdot)$.

- Convex Combination: $\mathbf{g}_{agg} = w\mathbf{g}_1 + (1-w)\mathbf{g}_2$
- Centroid Aggregation: $\mathbf{g}_{agg} \in \arg \min_{\mathbf{g} \in \mathcal{G}} \sum_{i=1}^m d(\mathbf{g}, \mathbf{g}_i)$
- Bayesian Optimization: $\max_{\mathbf{g}_{agg} \in \mathcal{G}} \mu(\mathbf{g}_{agg})$ where

$$k(\mathbf{g}_i, \mathbf{g}_j) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_x} \left[k(\mathbf{g}_i(\mathbf{x}), \mathbf{g}_j(\mathbf{x})) \right] = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_x} \left[e^{-\frac{1}{2} \|\mathbf{g}_i(\mathbf{x}) - \mathbf{g}_j(\mathbf{x})\|^2} \right]$$

Generalized Aggregation

Convex Combination

$$\mathbf{g}_{agg} = \mathbf{w}^T \mathbf{G}$$

$$\mathbf{G}^T = \left(\begin{array}{c|c|c|c} \text{SHAP}_1 & \text{LIME}_2 & \dots & \text{IG}_m \\ \hline & & & \end{array} \right)$$

Generalized Aggregation

Convex Combination

$$\mathbf{g}_{agg} = w^T G$$

$$G^T = \left(\begin{array}{c|c|c|c} \text{SHAP}_1 & \text{LIME}_2 & \dots & \text{IG}_m \\ \hline & & & \end{array} \right)$$

$$w_{agg} \in \arg \max_w \sum_{i=1}^m \mu(w^T G_i)$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

$$\mathbf{g}_c^S = [1 \quad -2 \quad 7] \rightarrow \mathbf{g}_m^S = [.1 \quad .2 \quad .7] \rightarrow \text{rank}^S = [C \quad B \quad A]$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

$$\mathbf{g}_c^S = [1 \quad -2 \quad 7] \rightarrow \mathbf{g}_m^S = [.1 \quad .2 \quad .7] \rightarrow \text{rank}^S = [C \quad B \quad A]$$

$$\text{rank}^{S_1} = [C \quad B \quad A] \quad \text{rank}^{S_2} = [C \quad A \quad B] \quad \text{rank}^{S_3} = [A \quad B \quad C]$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

$$\mathbf{g}_c^S = [1 \quad -2 \quad 7] \rightarrow \mathbf{g}_m^S = [.1 \quad .2 \quad .7] \rightarrow \text{rank}^S = [C \quad B \quad A]$$

$$\text{rank}^{S_1} = [C \quad B \quad A] \quad \text{rank}^{S_2} = [C \quad A \quad B] \quad \text{rank}^{S_3} = [A \quad B \quad C]$$

$$\text{Borda Count: } \mathbf{g}_{agg} = \text{rank}^{agg} = [C \quad A \quad B]$$

Generalized Aggregation

Classical Rank Aggregation

$$\mathcal{G}_m = \{SHAP_1, LIME_2, \dots, IG_m\}$$

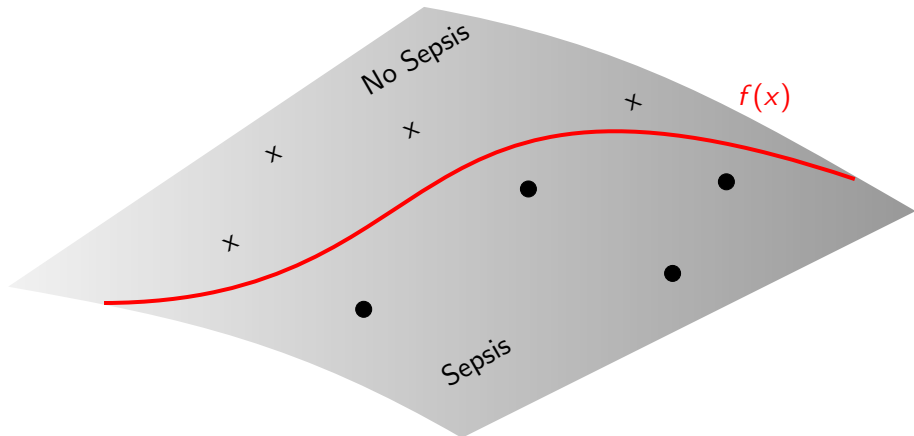
$$\mathbf{g}_c^S = [1 \quad -2 \quad 7] \rightarrow \mathbf{g}_m^S = [.1 \quad .2 \quad .7] \rightarrow \text{rank}^S = [C \quad B \quad A]$$

$$\text{rank}^{S_1} = [C \quad B \quad A] \quad \text{rank}^{S_2} = [C \quad A \quad B] \quad \text{rank}^{S_3} = [A \quad B \quad C]$$

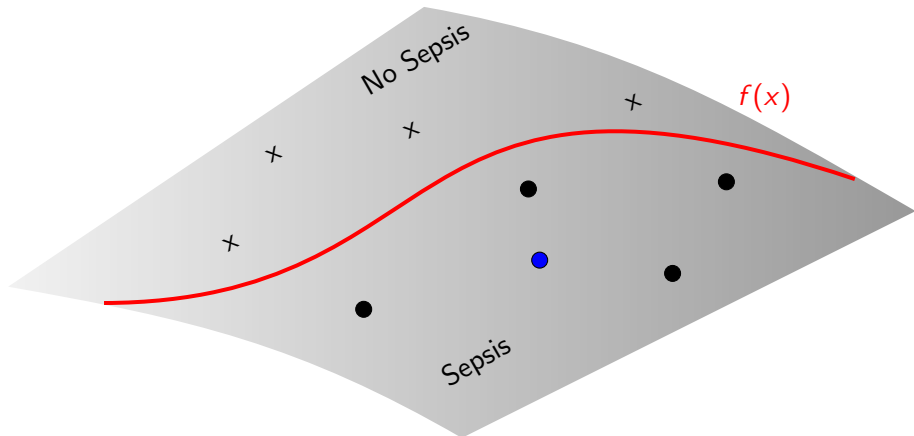
$$\text{Borda Count: } \mathbf{g}_{agg} = \text{rank}^{agg} = [C \quad A \quad B]$$

$$\mathbf{g}_{agg} \in \arg \min_{\mathbf{g}} \sum_{\mathbf{g}_i \in \mathcal{G}_m} d(\mathbf{g}, \mathbf{g}_i)$$

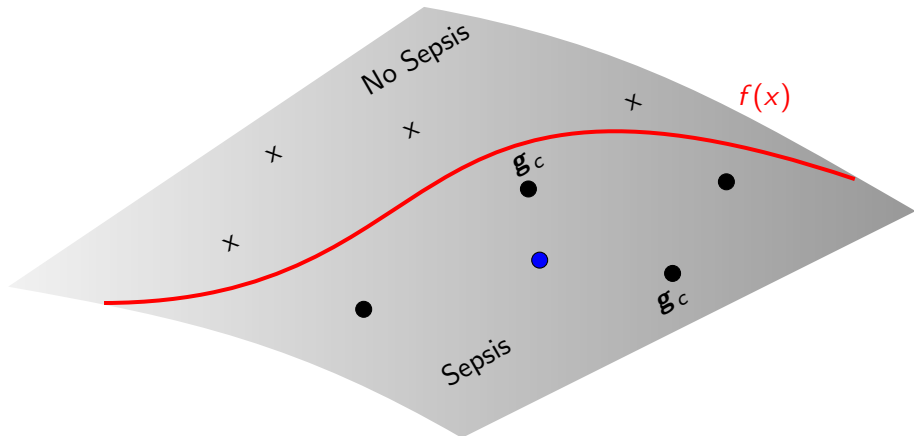
Aggregating Local Explanations



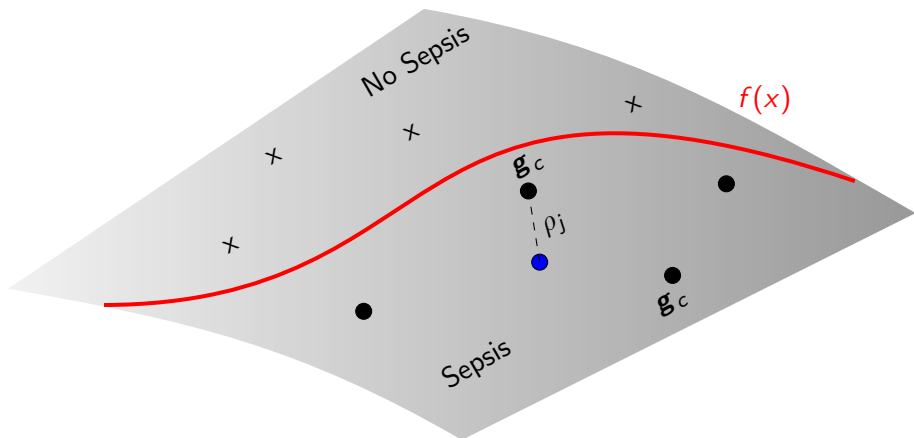
Aggregating Local Explanations



Aggregating Local Explanations



Aggregating Local Explanations



AVA: Aggregate Valuation of Antecedents

Can we use weighted Shapley values (Kalai et al. Journal of Game Theory 1987) to aggregate feature-based explanations with lower sensitivity?

AVA: Aggregate Valuation of Antecedents

Can we use weighted Shapley values (Kalai et al. Journal of Game Theory 1987) to aggregate feature-based explanations with lower sensitivity?

- 1 Find k nearest neighbors, \mathcal{N}_k , of x_{test} and their weights, ρ_j

$$\rho_j = \frac{d}{d\epsilon} \mathcal{L}(f_{\epsilon, x^{(j)}}, x_{\text{test}}) \Big|_{\epsilon=0}$$

$$\mathcal{N}_k(x_{\text{test}}, \mathcal{D}) = \arg \max_{\mathcal{N} \subset \mathcal{D}, |\mathcal{N}|=k} \sum_{x^{(j)} \in \mathcal{N}} \rho_j$$

AVA: Aggregate Valuation of Antecedents

Can we use weighted Shapley values (Kalai et al. Journal of Game Theory 1987) to aggregate feature-based explanations with lower sensitivity?

- 1 Find k nearest neighbors, \mathcal{N}_k , of x_{test} and their weights, ρ_j

$$\rho_j = \frac{d}{d\epsilon} \mathcal{L}(f_{\epsilon, x^{(j)}}, x_{\text{test}}) \Big|_{\epsilon=0}$$

$$\mathcal{N}_k(x_{\text{test}}, \mathcal{D}) = \arg \max_{\mathcal{N} \subset \mathcal{D}, |\mathcal{N}|=k} \sum_{x^{(j)} \in \mathcal{N}} \rho_j$$

- 2 Calculate the attributions, \mathbf{g}_{ci} , for all points in \mathcal{N}_k

$$\mathbf{g}_{ci} = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (f(x_{S \cup \{i\}}) - f(x_S))$$

AVA: Aggregate Valuation of Antecedents

Can we use weighted Shapley values (Kalai et al. Journal of Game Theory 1987) to aggregate feature-based explanations with lower sensitivity?

- 1 Find k nearest neighbors, \mathcal{N}_k , of x_{test} and their weights, ρ_j

$$\rho_j = \frac{d}{d\epsilon} \mathcal{L}(f_{\epsilon, x^{(j)}}, x_{\text{test}}) \Big|_{\epsilon=0}$$

$$\mathcal{N}_k(x_{\text{test}}, \mathcal{D}) = \arg \max_{\mathcal{N} \subset \mathcal{D}, |\mathcal{N}|=k} \sum_{x^{(j)} \in \mathcal{N}} \rho_j$$

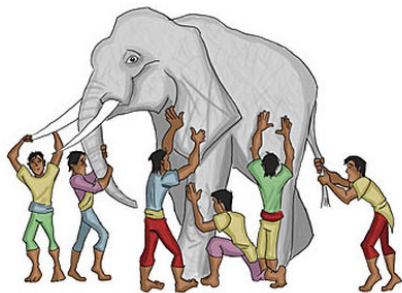
- 2 Calculate the attributions, \mathbf{g}_c , for all points in \mathcal{N}_k

$$\mathbf{g}_{ci} = \phi_i = \frac{1}{|F|} \sum_{S \subseteq F \setminus \{i\}} \binom{F-1}{S}^{-1} (f(x_{S \cup \{i\}}) - f(x_S))$$

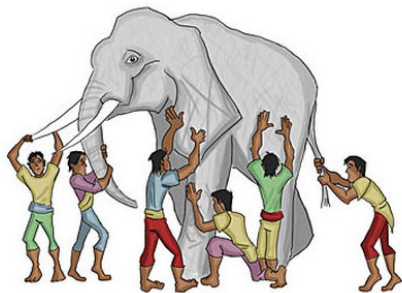
- 3 Aggregate the k explanations into a consensus, \mathbf{g}_{agg}

$$\mathbf{g}_{\text{agg}} = \sum_{x^{(j)} \in \mathcal{N}_k} \frac{\rho_j}{\rho} \mathbf{g}_c^j$$

Why aggregate?

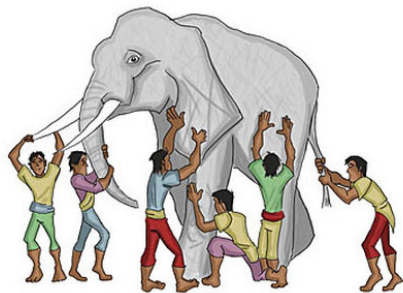


Why aggregate?



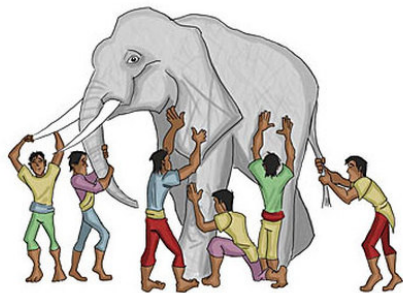
- Each training point has its own “learned” attribution

Why aggregate?



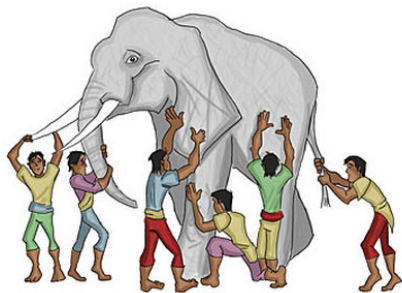
- Each training point has its own “learned” attribution
- Aggregate explanation now has lower sensitivity

Why aggregate?



- Each training point has its own “learned” attribution
- Aggregate explanation now has lower sensitivity
- Resulting attribution uses motivating reasoning of a doctor

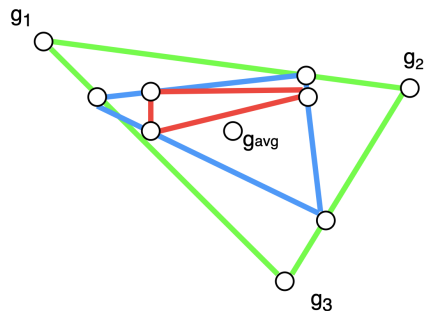
Why aggregate?



- Each training point has its own “learned” attribution
- Aggregate explanation now has lower sensitivity
- Resulting attribution uses motivating reasoning of a doctor
- **SUPER SUPER** cheap to compute

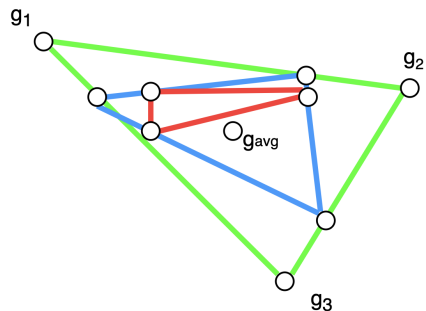
Minimizing Complexity

Region Shrinking Method

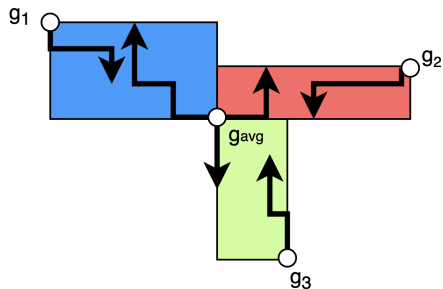


Minimizing Complexity

Region Shrinking Method



Gradient-Descent Style Method



Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?
 - Consider counterfactuals, natural language explanations, etc.

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?
 - Consider counterfactuals, natural language explanations, etc.
- 3 Working with medical experts to find a \mathbf{g}^*

Conclusion

Summary

- 1 Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- 2 Aggregating Shapley value explanations results in a Shapley value
- 3 We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- 1 Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- 2 Are feature-based explanations even useful?
 - Consider counterfactuals, natural language explanations, etc.
- 3 Working with medical experts to find a \mathbf{g}^*
- 4 Multi-Objective optimization

Conclusion

Summary

- ① Aggregate local explanations with classical rank aggregation or via convex combination can be useful
- ② Aggregating Shapley value explanations results in a Shapley value
- ③ We can learn aggregate explanations to lower sensitivity and complexity

Future Work

- ① Axiomatic Aggregation
 - If $\mathbf{g}_1, \dots, \mathbf{g}_n$ satisfy Axiom R, then $\mathbf{g}_{\mathcal{A}}$ satisfies R.
- ② Are feature-based explanations even useful?
 - Consider counterfactuals, natural language explanations, etc.
- ③ Working with medical experts to find a \mathbf{g}^*
- ④ Multi-Objective optimization
 - Resulting Setup

$$\max \text{faithfulness}(\mathbf{g}_{agg}) + \text{sensitivity}(\mathbf{g}_{agg})$$