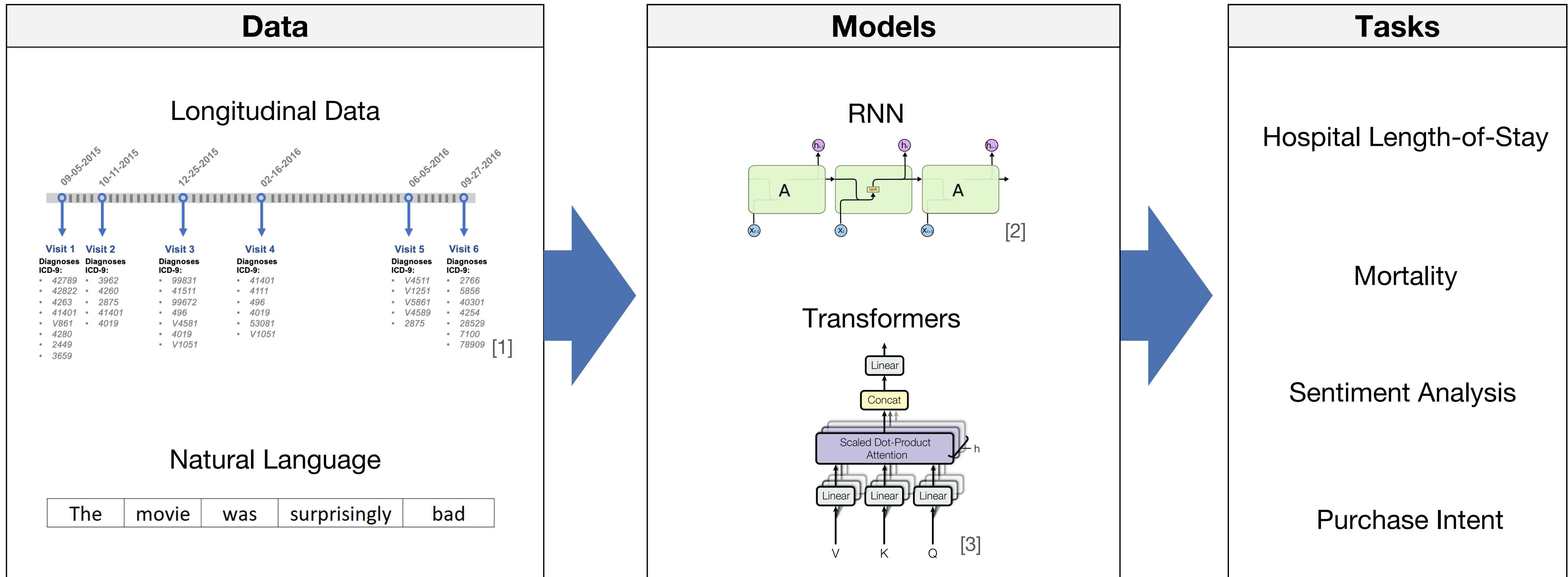


OrdShap: Feature Position Importance for Sequential Black-Box Models

Davin Hill, Brian L. Hill, Aria Masoomi, Vijay S. Nori, Robert E. Tillman & Jennifer Dy

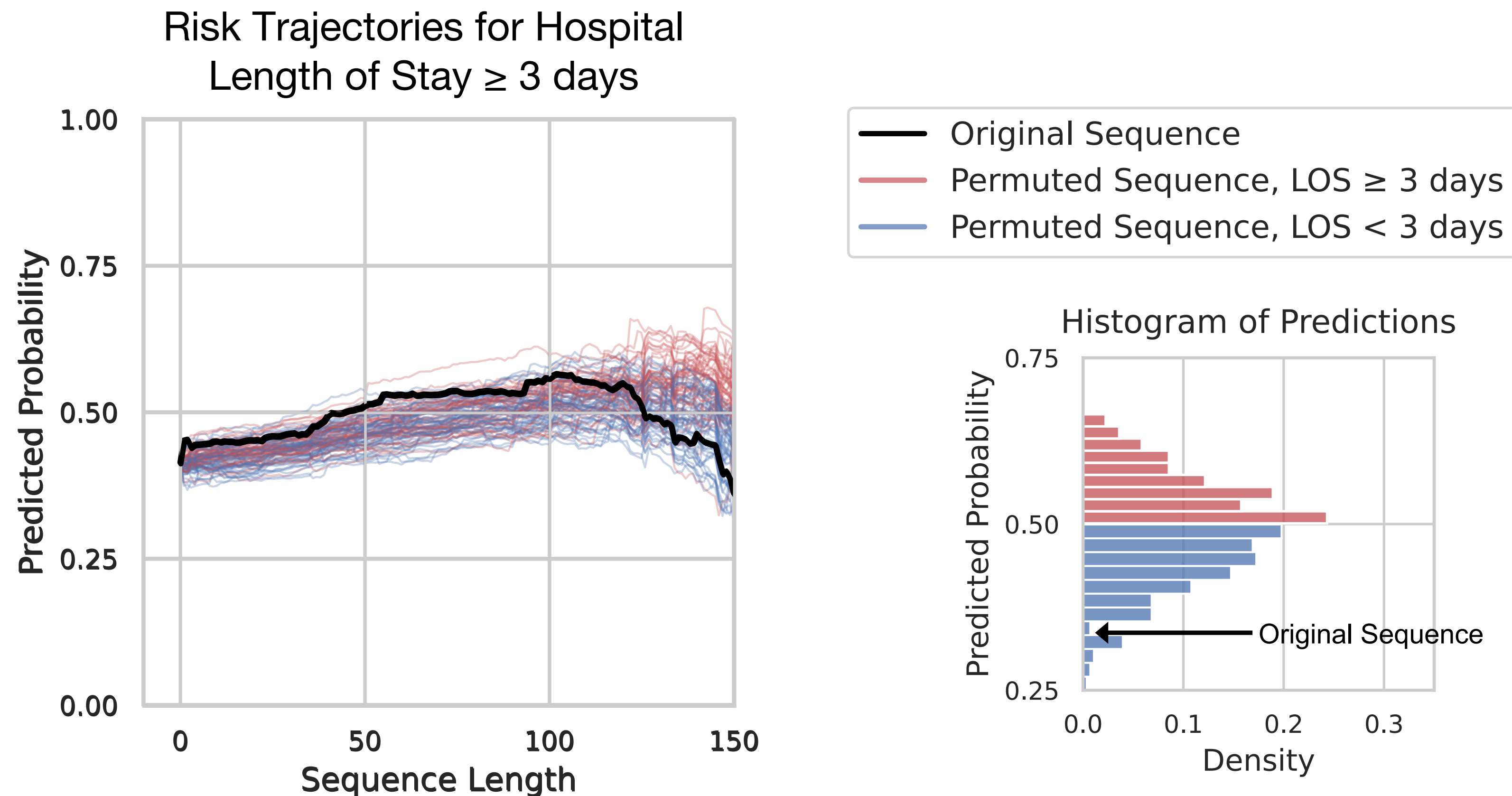


Making Predictions with Sequential Models



[1] Baytas et al., Patient Subtyping via Time-Aware LSTM Networks. [2] <https://colah.github.io/posts/2015-08-Understanding-LSTMs/> [3] Vaswani et al., Attention is all you need.

Permuting Ordering of Data Samples



Existing Attributions do not capture feature ordering, however permuting order significantly affects model output

Overview of OrdShap

Challenges

- How do we distentangle importance due to **1) Feature Value**, and **2) Feature Position**?
- How do we efficiently approximate this attribution?

Contributions

- We propose **OrdShap**, an attribution for Sequential models
- We establish a theoretical connection between OrdShap and Sanchez-Bergantinos Values
- We propose 2 algorithms to approximate OrdShap

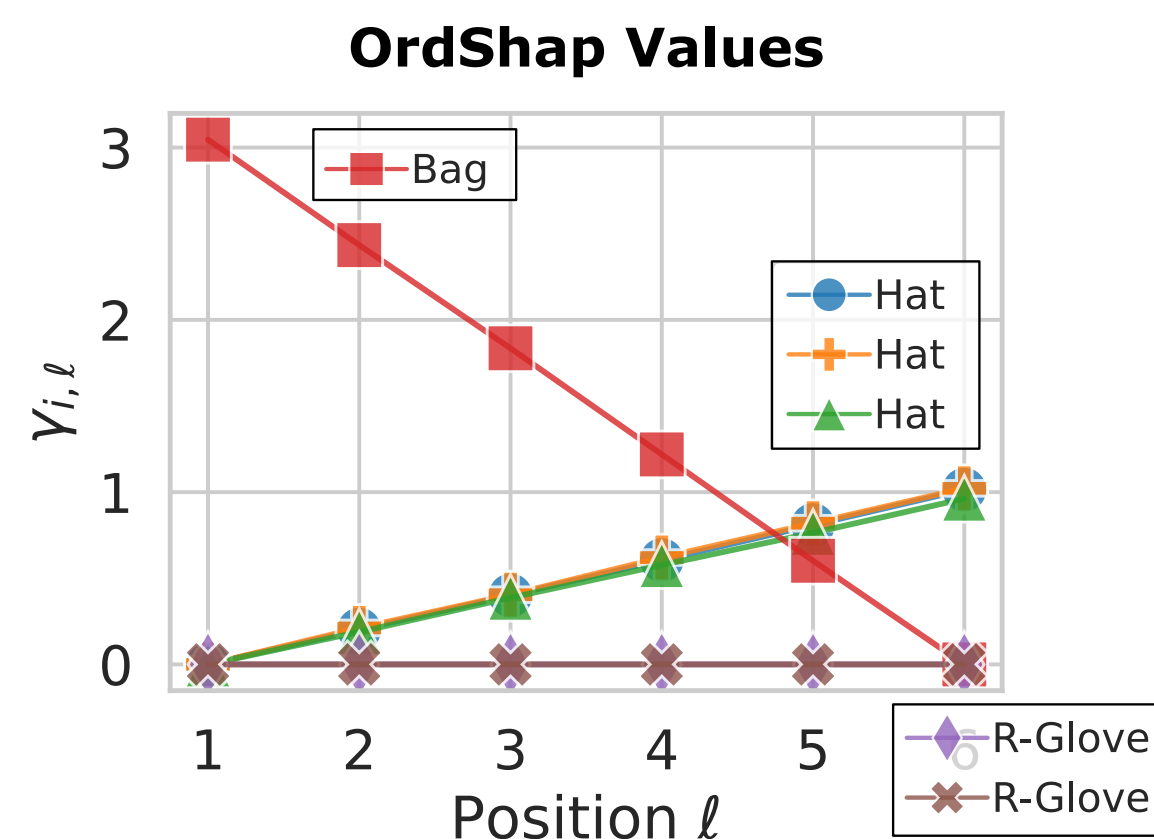
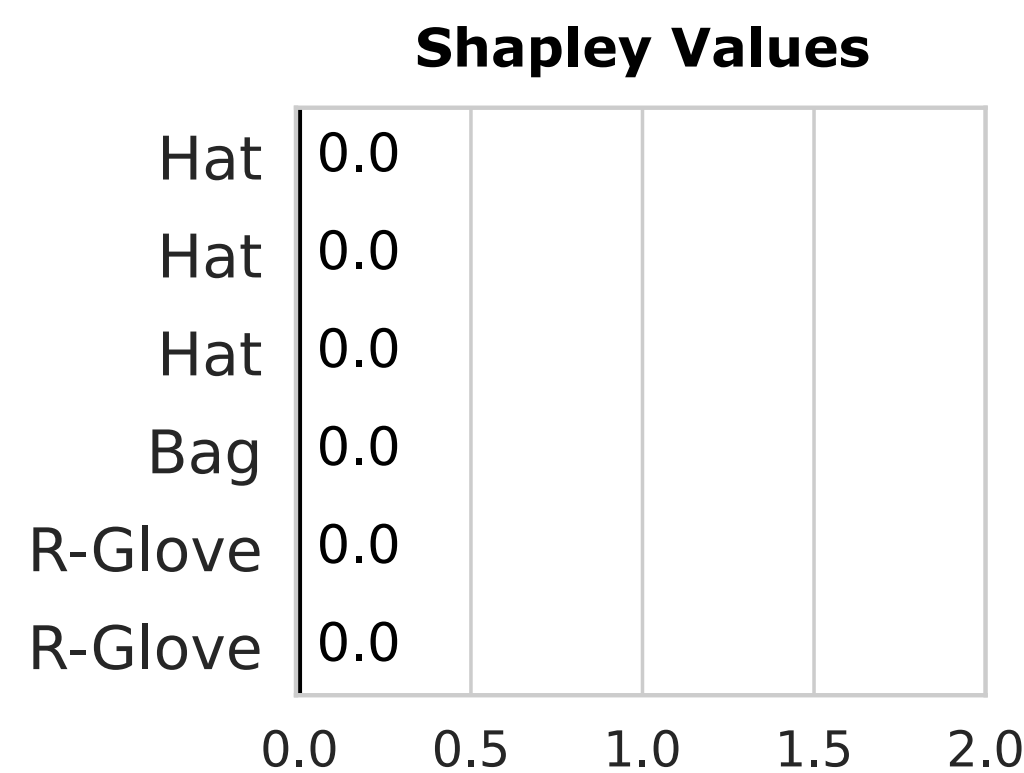
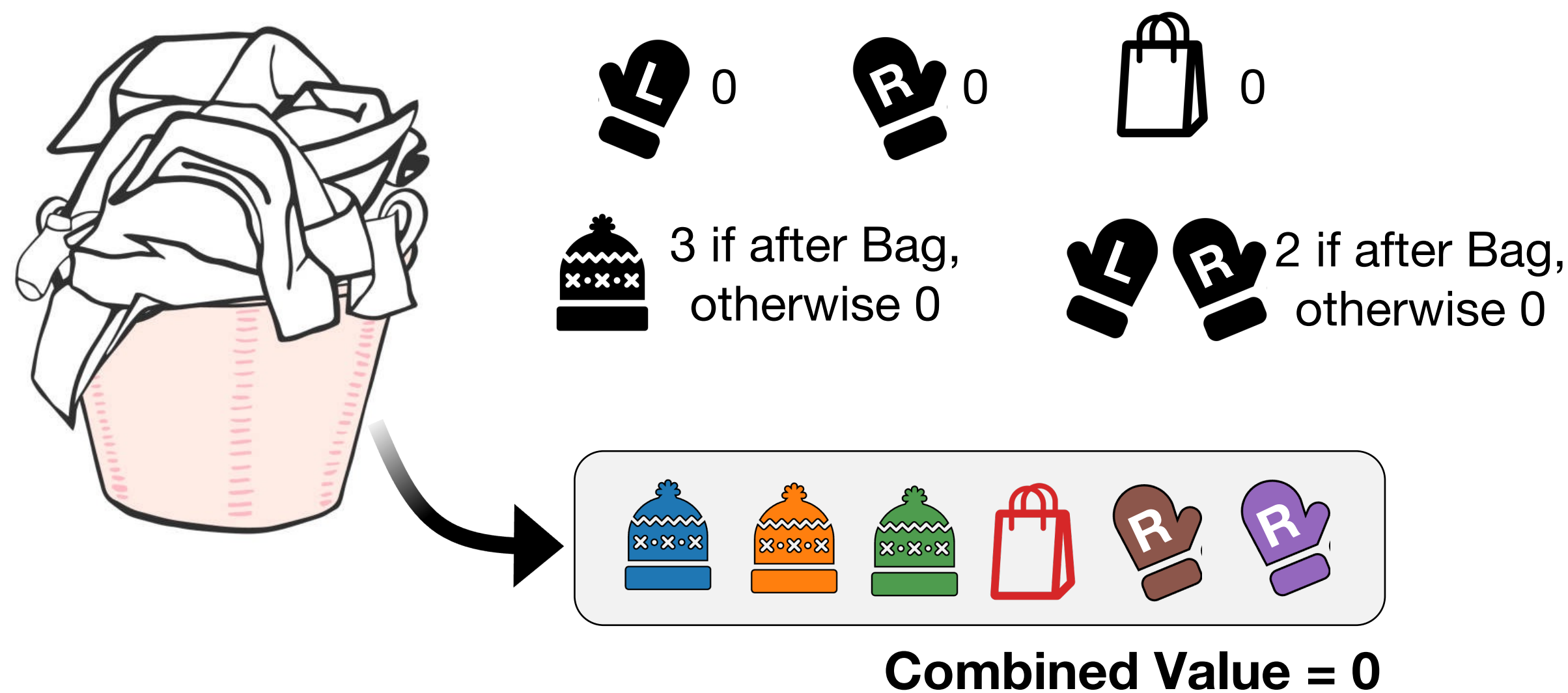
Quantifying Positional Dependency

$$\begin{aligned}
 \gamma_{i,\ell}(N, \tilde{\omega}) = & \underbrace{\sum_{\substack{S \subseteq N \\ i \in S}}}_{\text{Average over Subsets}} \underbrace{\sum_{\substack{\sigma \in \mathfrak{S}_N \\ \sigma^{-1}(i) = \ell}}}_{\substack{\text{Average over} \\ \text{Permutations} \\ \text{Conditioned} \\ \text{on Permuting} \\ i \rightarrow \ell}} \underbrace{\frac{(|S| - 1)! (|N| - |S|)!}{(|N| - 1)! |N|!}}_{\text{Weighting}} \underbrace{[\tilde{\omega}(S, \sigma) - \tilde{\omega}(S \setminus \{i\}, \sigma)]}_{\text{Marginal decrease from removing } i}
 \end{aligned}$$

\mathfrak{S}_S Symmetric Group on $S = \{1, \dots, |S|\}$

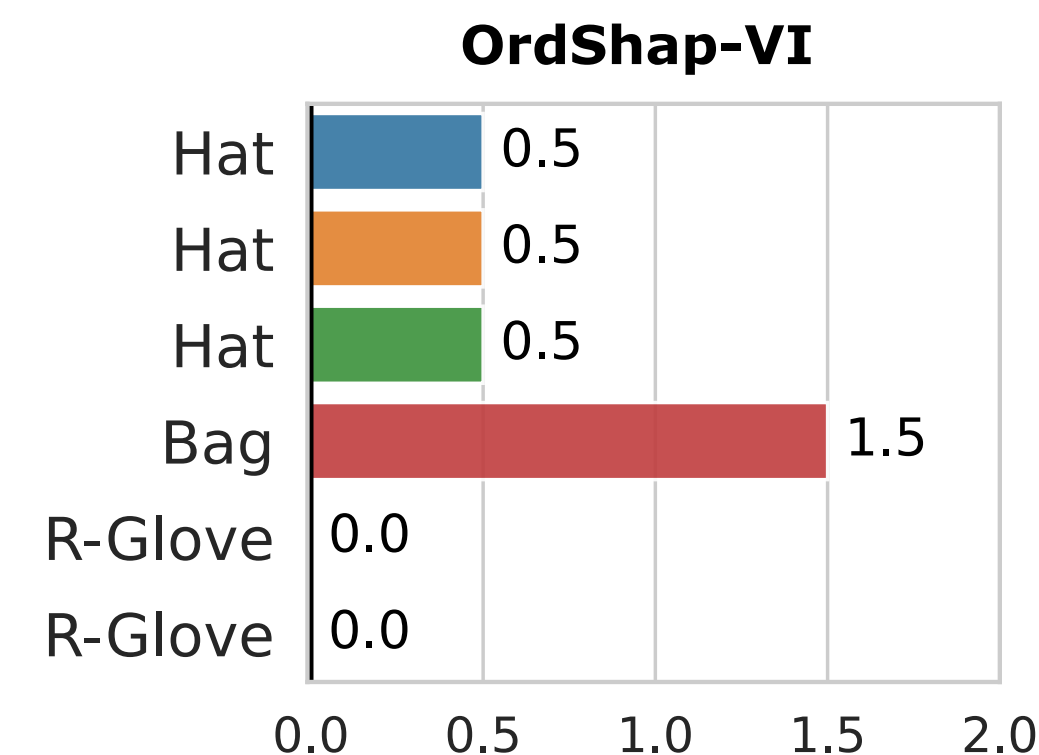
ω Characteristic function on Permutations

Toy Example: Position and Value Importance



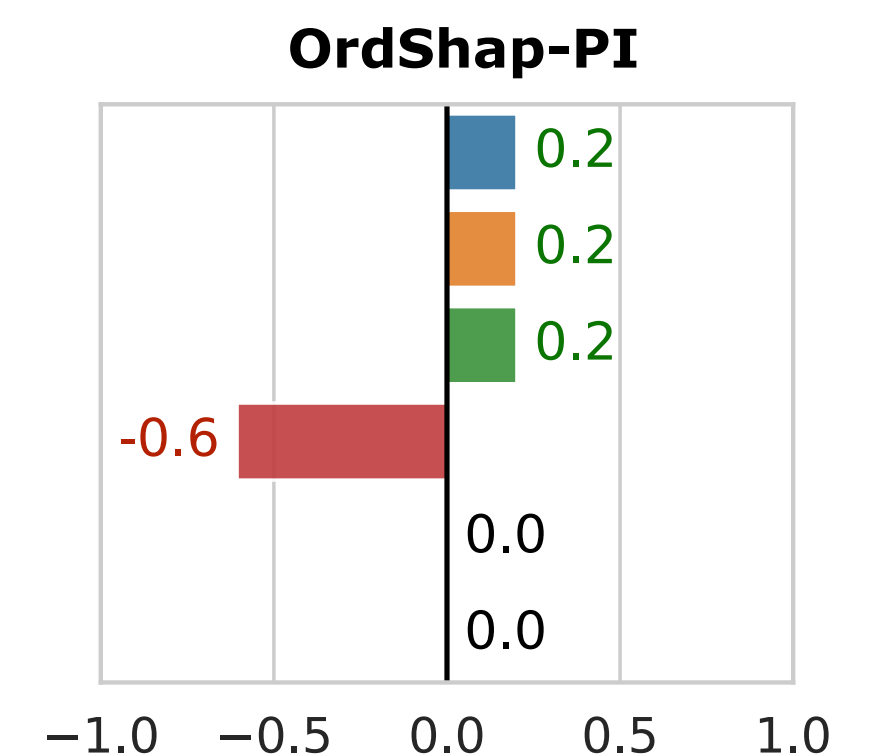
1. Value Importance

$$\bar{\gamma}_i(N, \tilde{\omega}) = \frac{1}{|N|} \sum_{\ell \in N} \gamma_{i,\ell}(N, \tilde{\omega})$$



2. Position Importance

$$\arg \min_{\beta_i} [(\ell - \bar{\ell})\beta_i - (\gamma_{i,\ell}(N, v) - \bar{\gamma}_i)]^2$$



Positive OrdShap-PI:
More important when later
in the sequence

Theorem 1: OrdShap-VI fulfills Shapley Axioms of Efficiency*, Symmetry*, Null Player*, Additivity*

OrdShap Approximation

Naïve Calculation: $\mathcal{O}(d! 2^d \delta_f)$

Monte Carlo Sampling Algorithm: $\mathcal{O}(dKL\delta_f + d^2KL)$

Least-Squares Algorithm with Monte Carlo Sampling: $\mathcal{O}(KL\delta_f + d^2KL + d^3)$

d : # features

δ_f : Model Evaluation

K, L : # MC samples

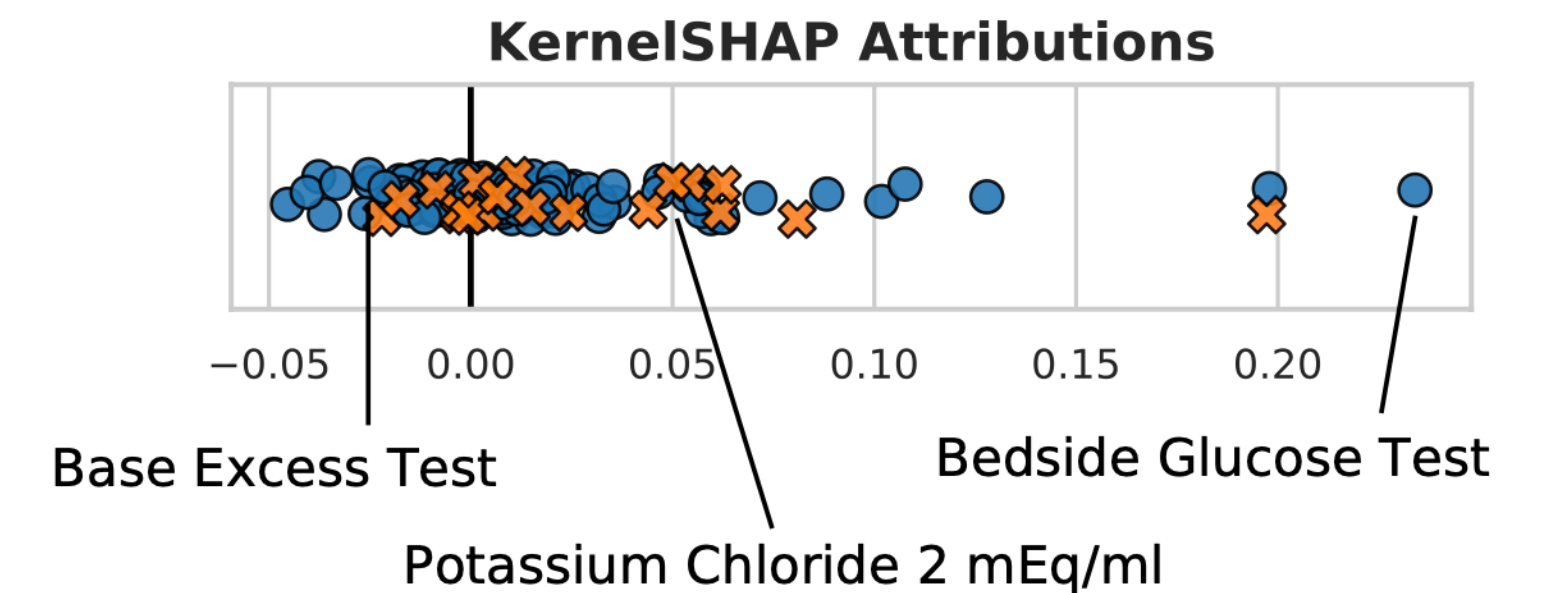
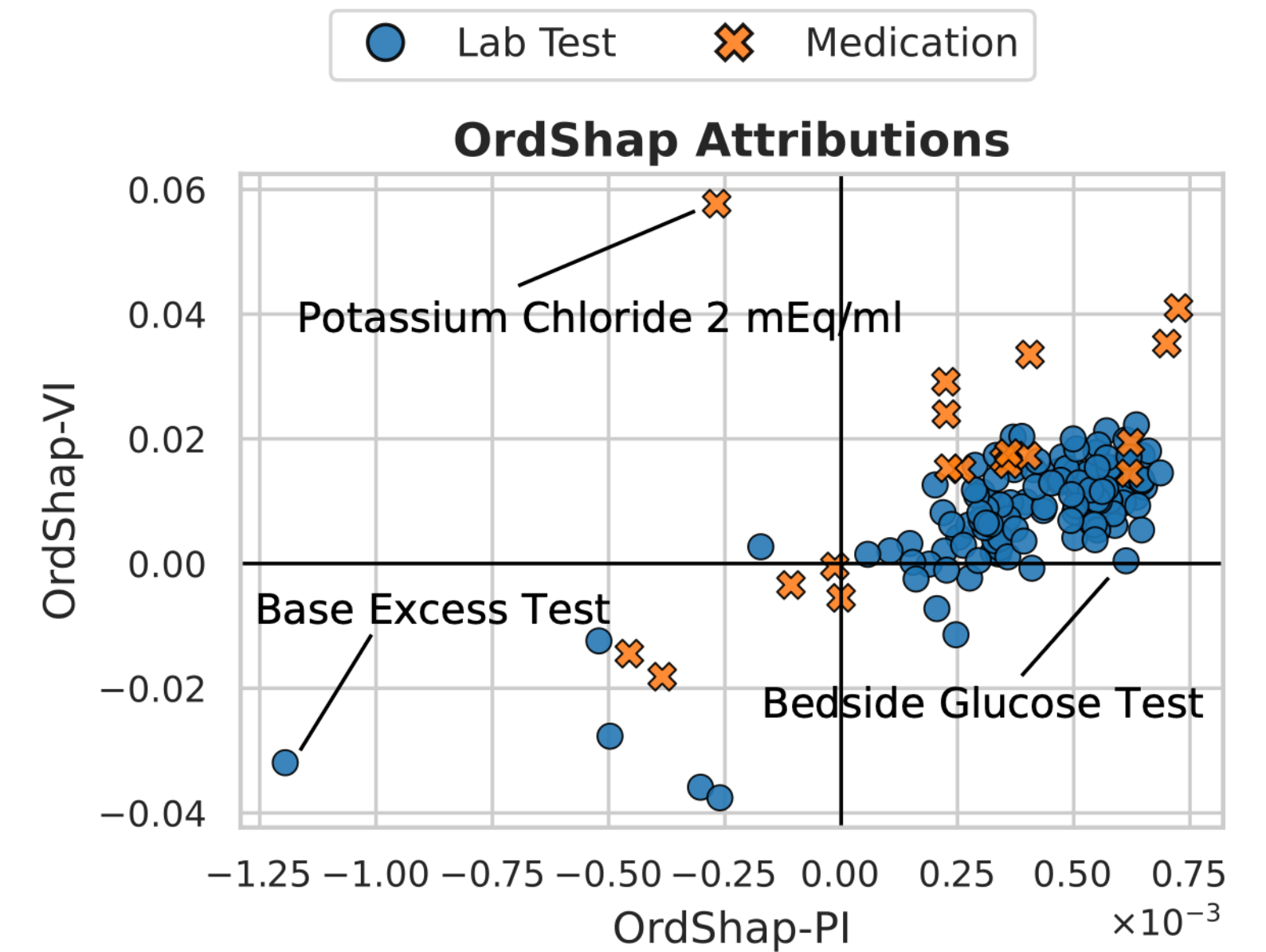
$$\min_{\alpha_1, \dots, \alpha_n, \beta_1, \dots, \beta_n} \sum_{\substack{S \subseteq N \\ S \neq \emptyset, N}} \sum_{\sigma \in \mathfrak{S}_N} \mu(|S|) \left[\underbrace{\sum_{i \in S} \alpha_i}_{\text{Value Importance}} + \underbrace{\sum_{i \in S} [\sigma^{-1}(i) - \bar{\ell}] \beta_i}_{\text{Position Importance}} - \underbrace{[\tilde{\omega}(S, \sigma) - \tilde{\omega}(\emptyset, \sigma_N^{id})]}_{\text{Model Output}} \right]^2$$

$$s.t. \sum_{i \in N} \alpha_i = \frac{1}{|N|!} \sum_{\sigma \in \mathfrak{S}_N} \tilde{\omega}(N, \sigma) + \tilde{\omega}(\emptyset, \sigma_N^{id})$$

Theorem 2: Optimal α, β corresponds to OrdShap

Experiments

- Quantitative Comparisons using Inclusion / Exclusion AUC
- Evaluation on Synthetic Data
- Qualitative Examples on Medical Tokens
- Execution Time Results
- Sensitivity Analysis



Thanks for Listening!

OrdShap: Feature Position Importance for Sequential Black-Box Models

Davin Hill, Brian L. Hill, Aria Masoomi, Vijay S. Nori, Robert E. Tillman & Jennifer Dy

