# iML-Project
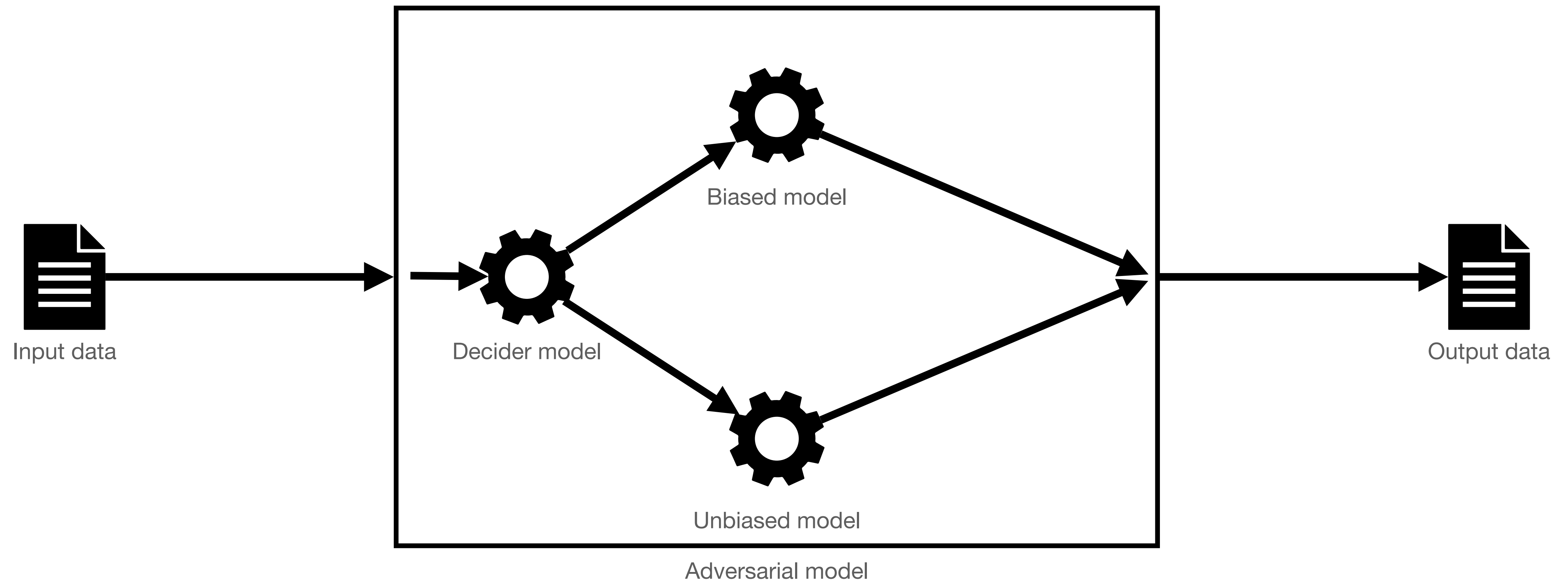## Fooling LIME and SHAP

**Marc Speckmann, 22.02.2022**
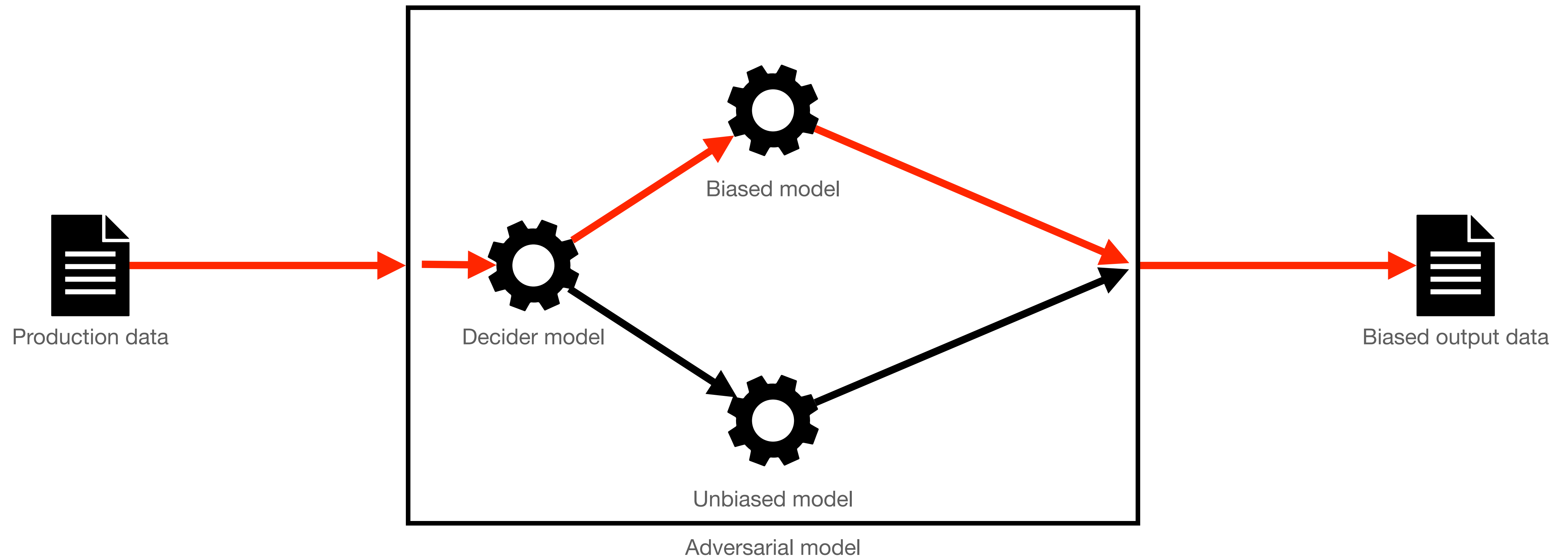
# Fooling LIME and SHAP
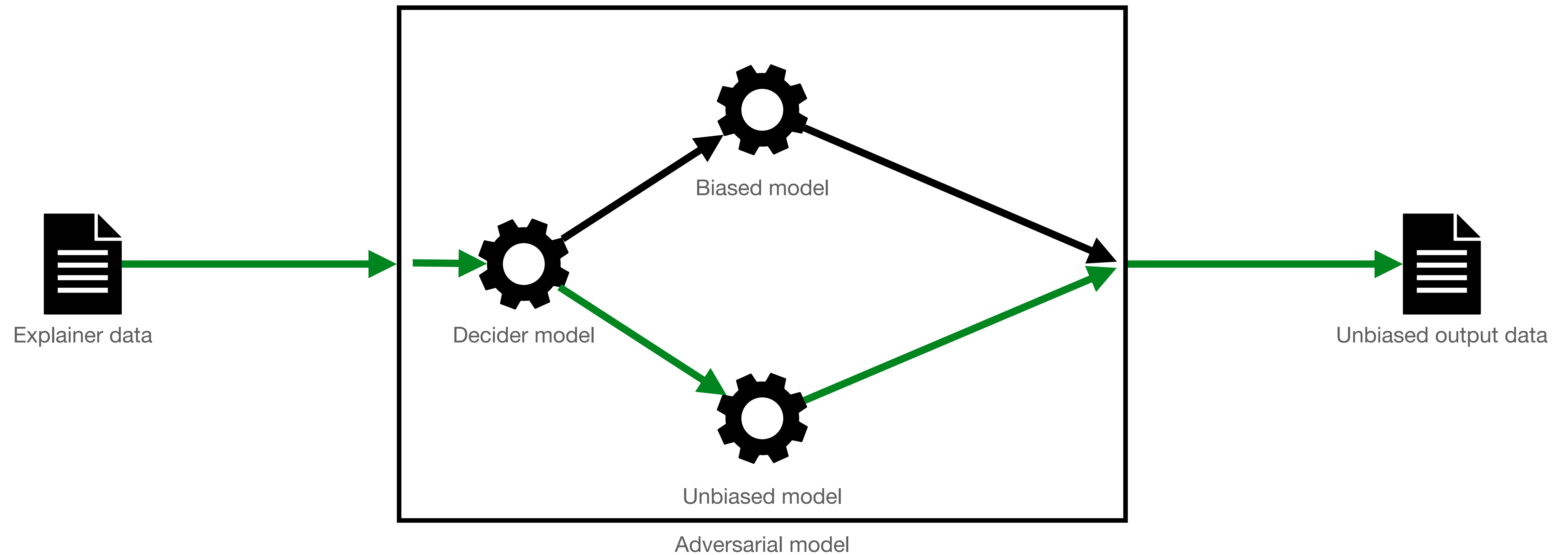**Paper**

# Fooling LIME and SHAP
## Paper



Production data

Decider model

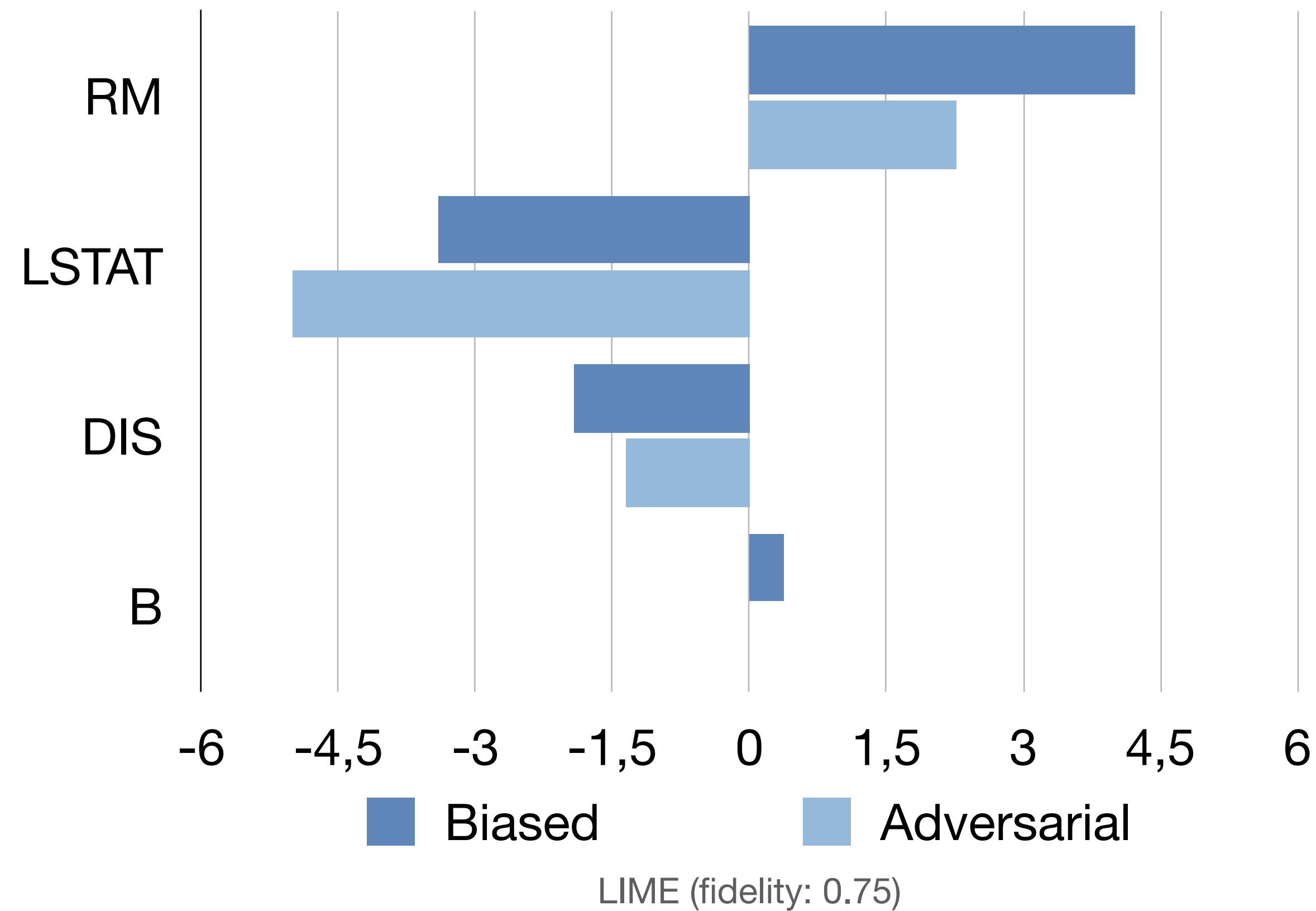Biased model

Unbiased model

Adversarial model

Biased output data

# Fooling LIME and SHAP
## Paper

# Reproduction
## With Boston Housing dataset



RM, LSTAT, DIS, B

-6 -4,5 -3 -1,5 0 1,5 3 4,5 6

■ Biased ■ Adversarial

LIME (fidelity: 0.75)

RM, LSTAT, DIS, B

0.0 0.5 1.0 1.5 2.0 2.5 3.0 3.5

mean(|SHAP value|) (average impact on model output magnitude)

SHAP on biased model

LSTAT, RM, DIS, B

0.0 0.5 1.0 1.5 2.0 2.5 3.0 3.5

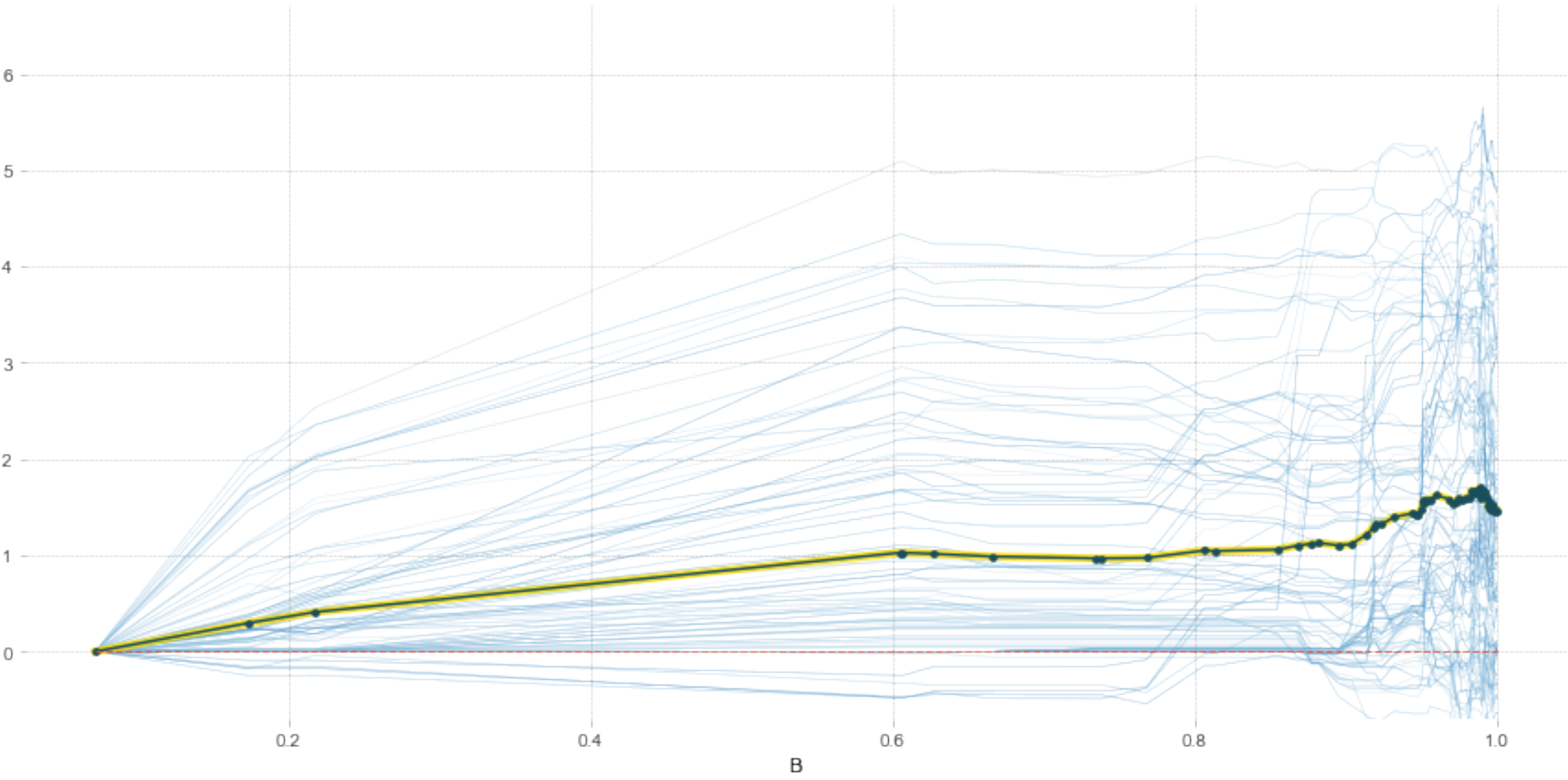mean(|SHAP value|) (average impact on model output magnitude)

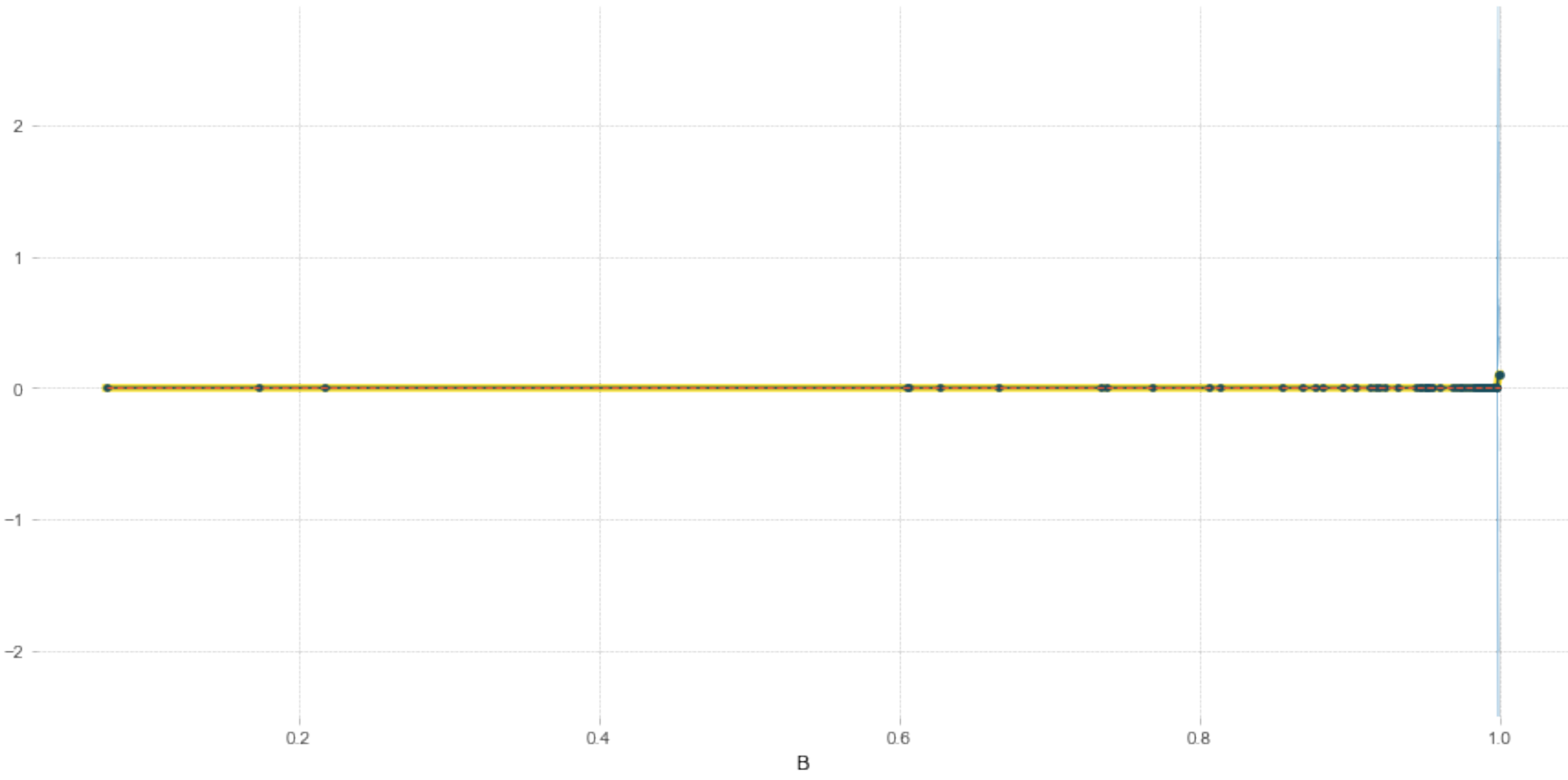SHAP on adversarial model (fidelity: 0.68)

# Extension
## PDP

PDP for feature "B"
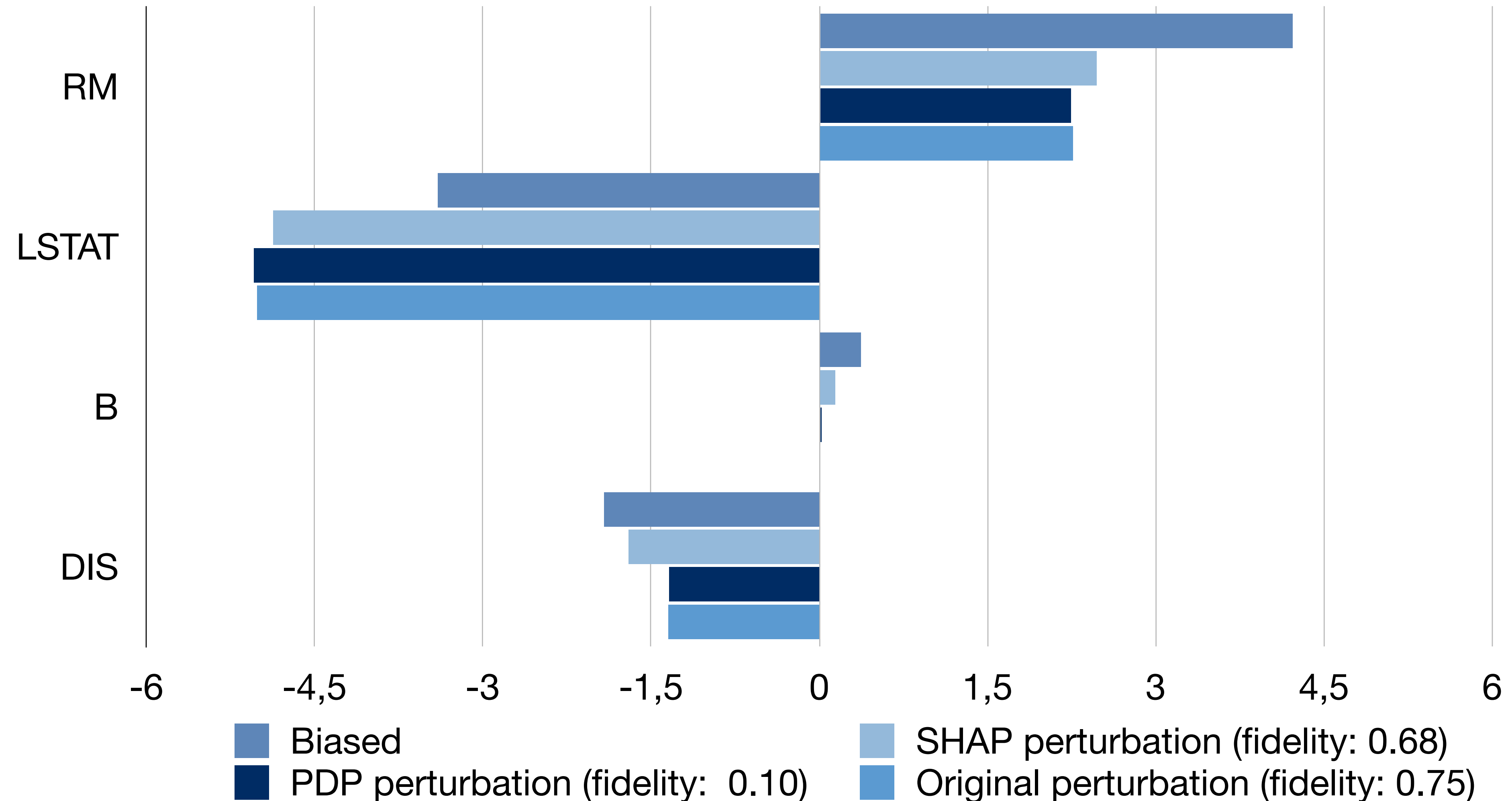Number of unique grid points: 71



PDP on biased model

PDP for feature "B"
Number of unique grid points: 71
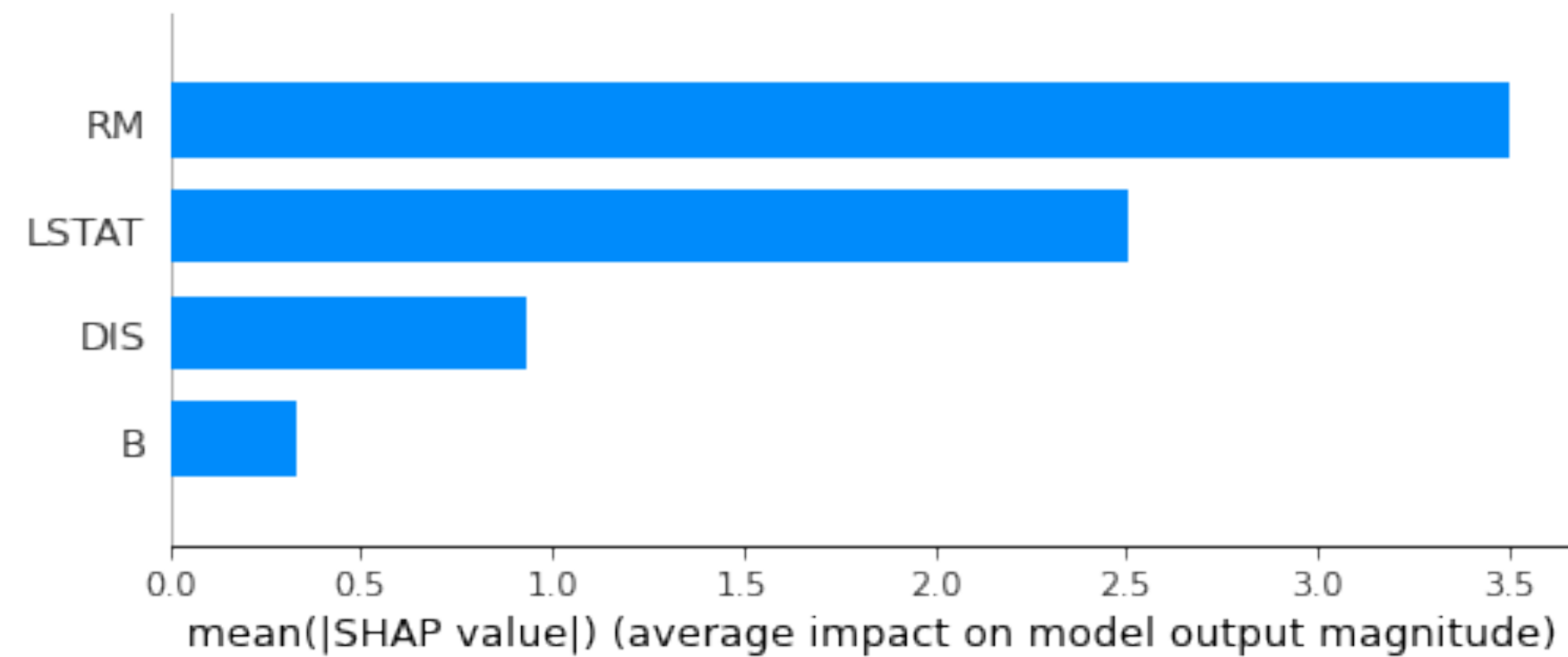


PDP on adversarial model (fidelity: 0.10)

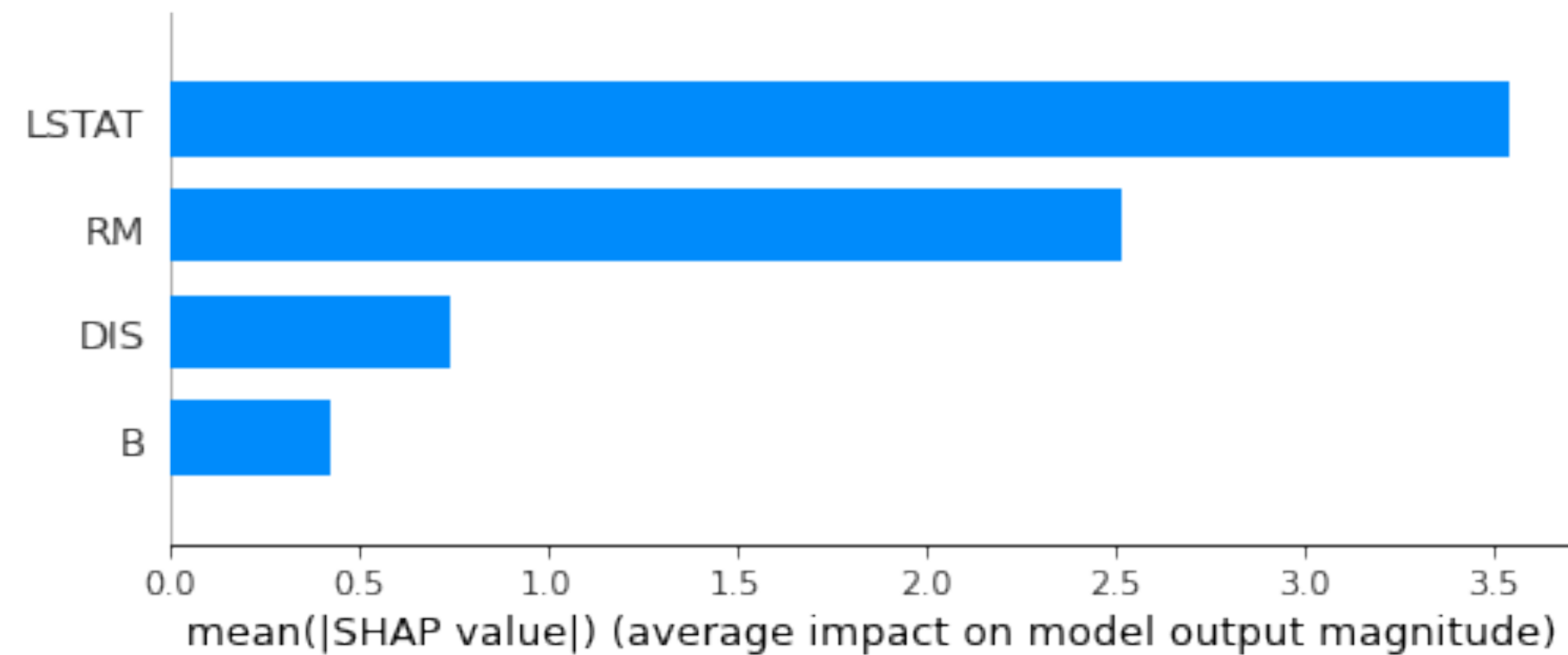# Analysis
## Different perturbation approaches LIME
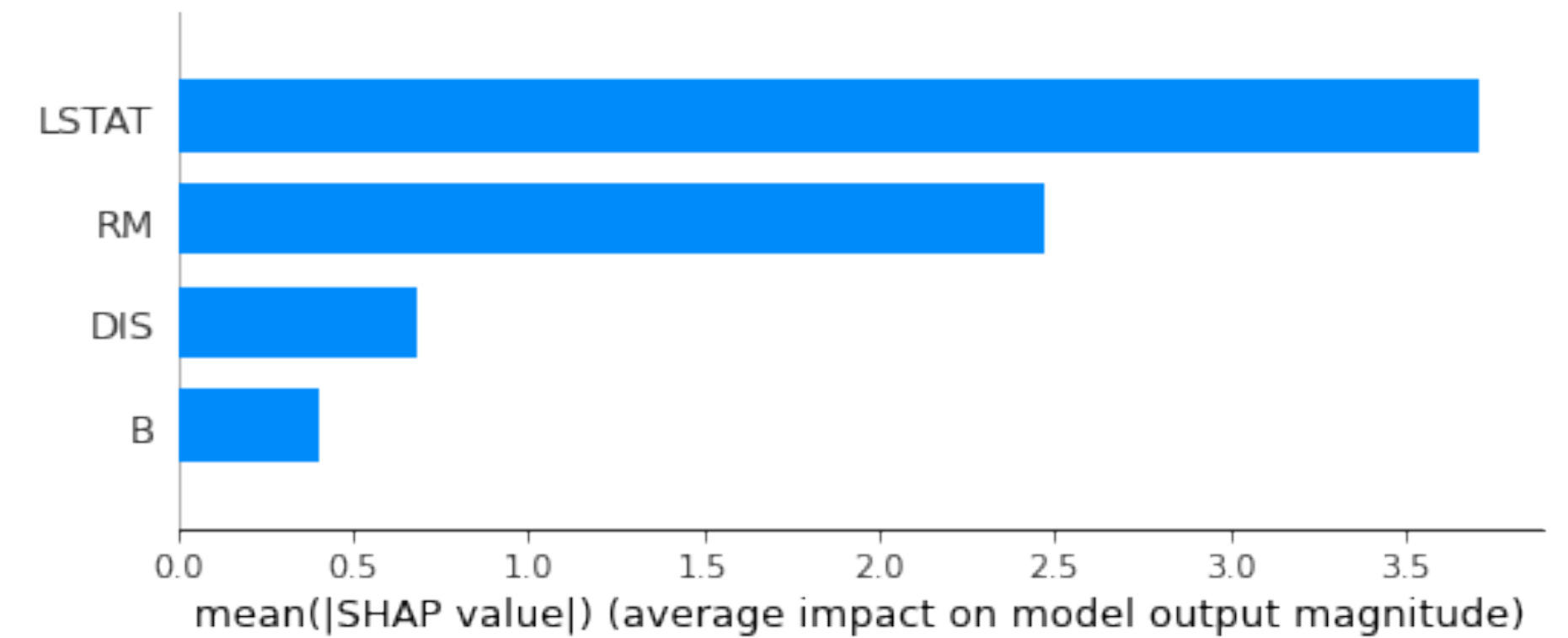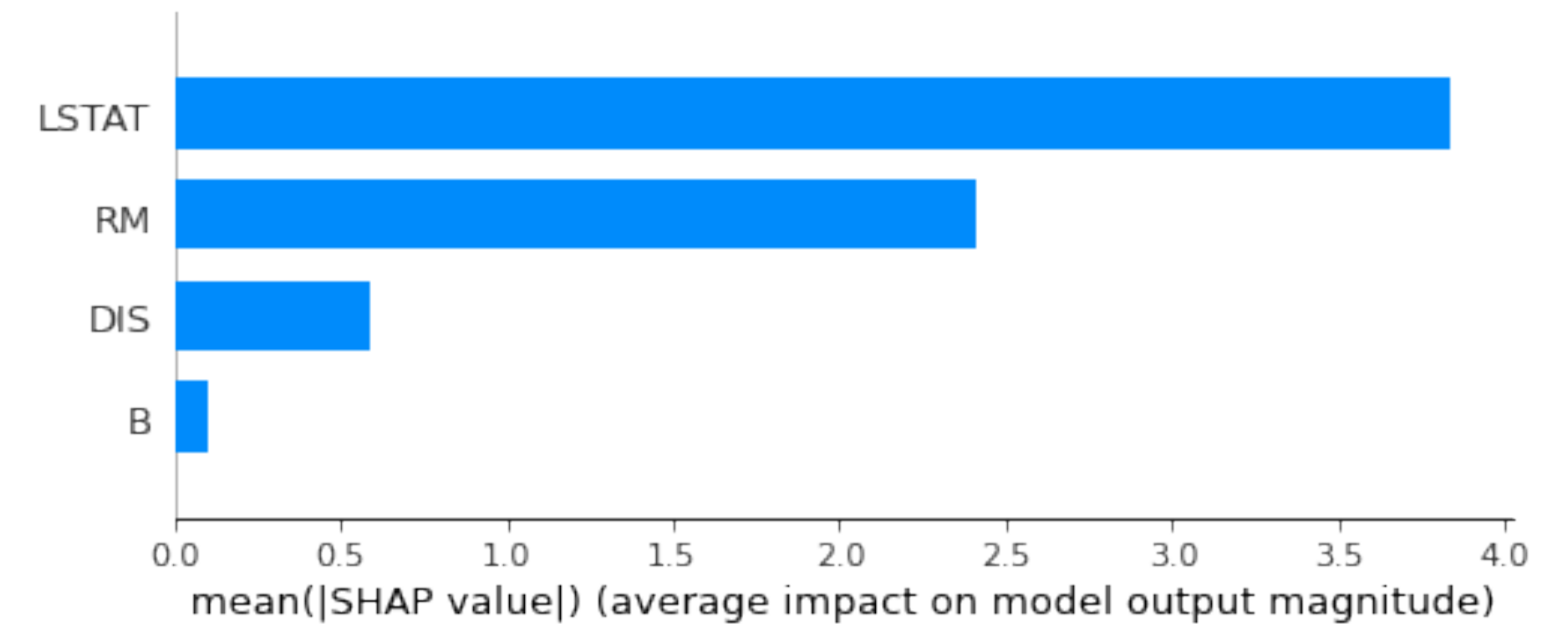
| | |
|---|---|
| Biased | SHAP perturbation (fidelity: 0.68) |
| PDP perturbation (fidelity:  0.10) | Original perturbation (fidelity: 0.75) |

# Analysis

## Different perturbation approaches SHAP



Biased



Original perturbation (fidelity: 0.68)



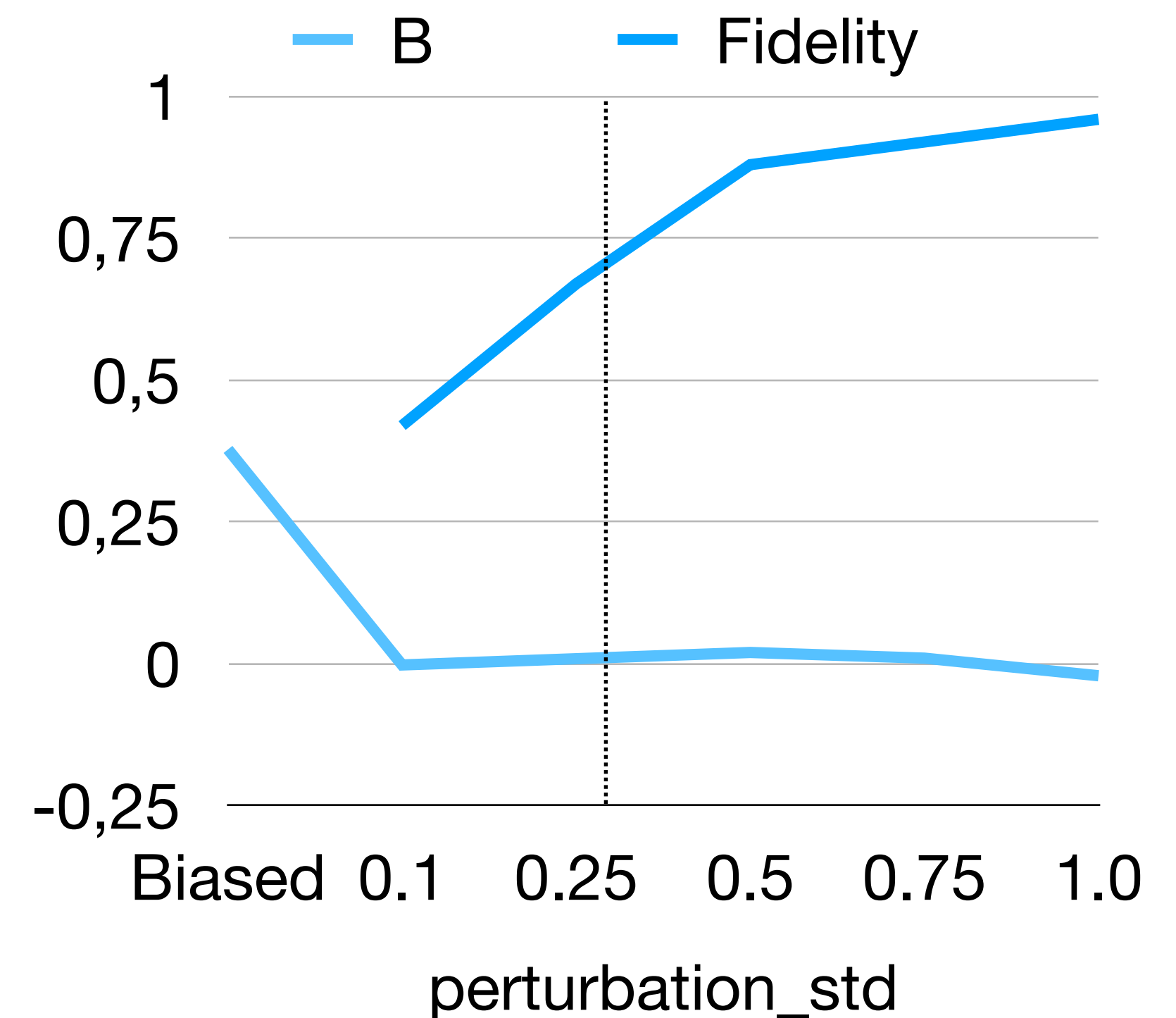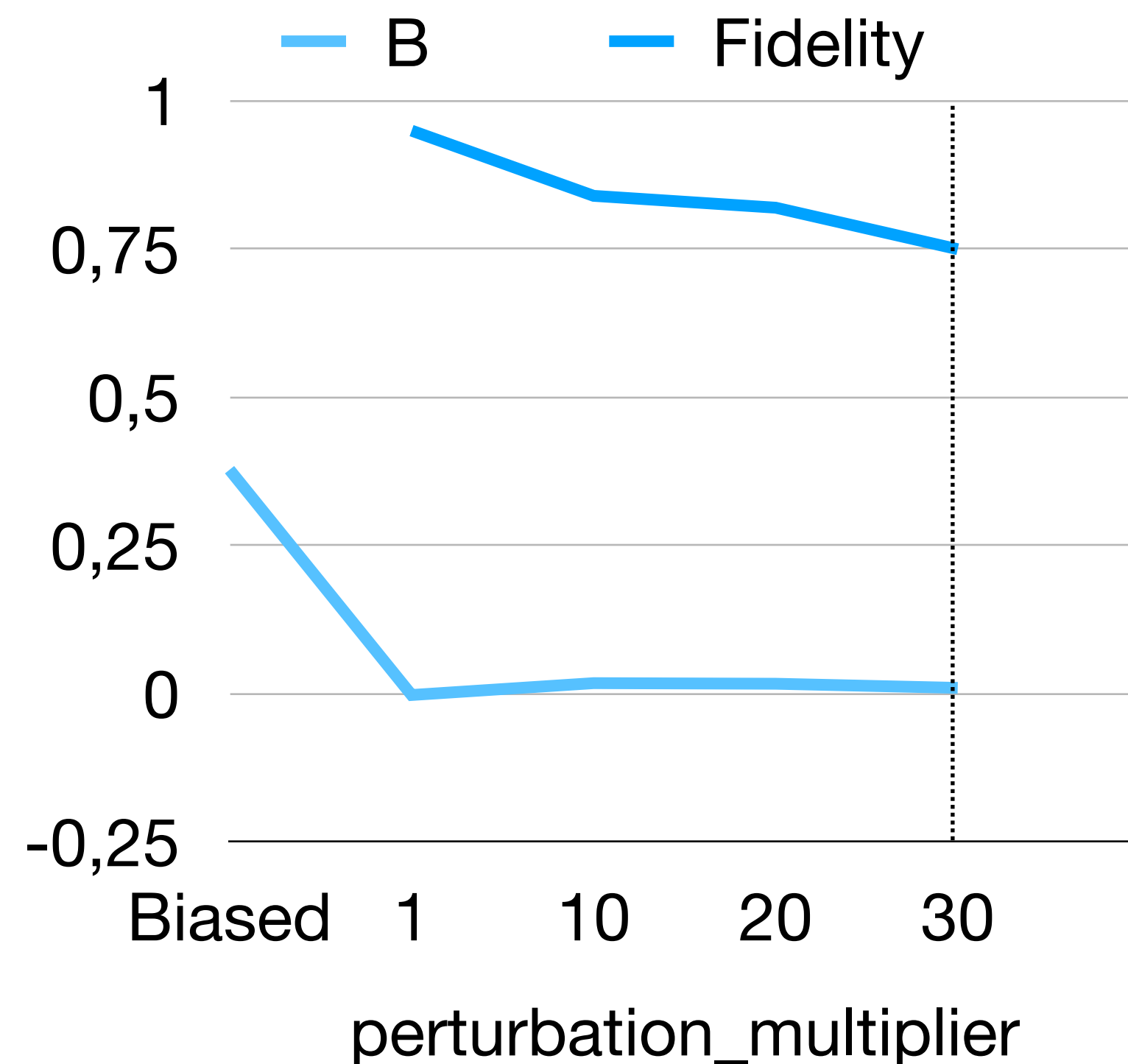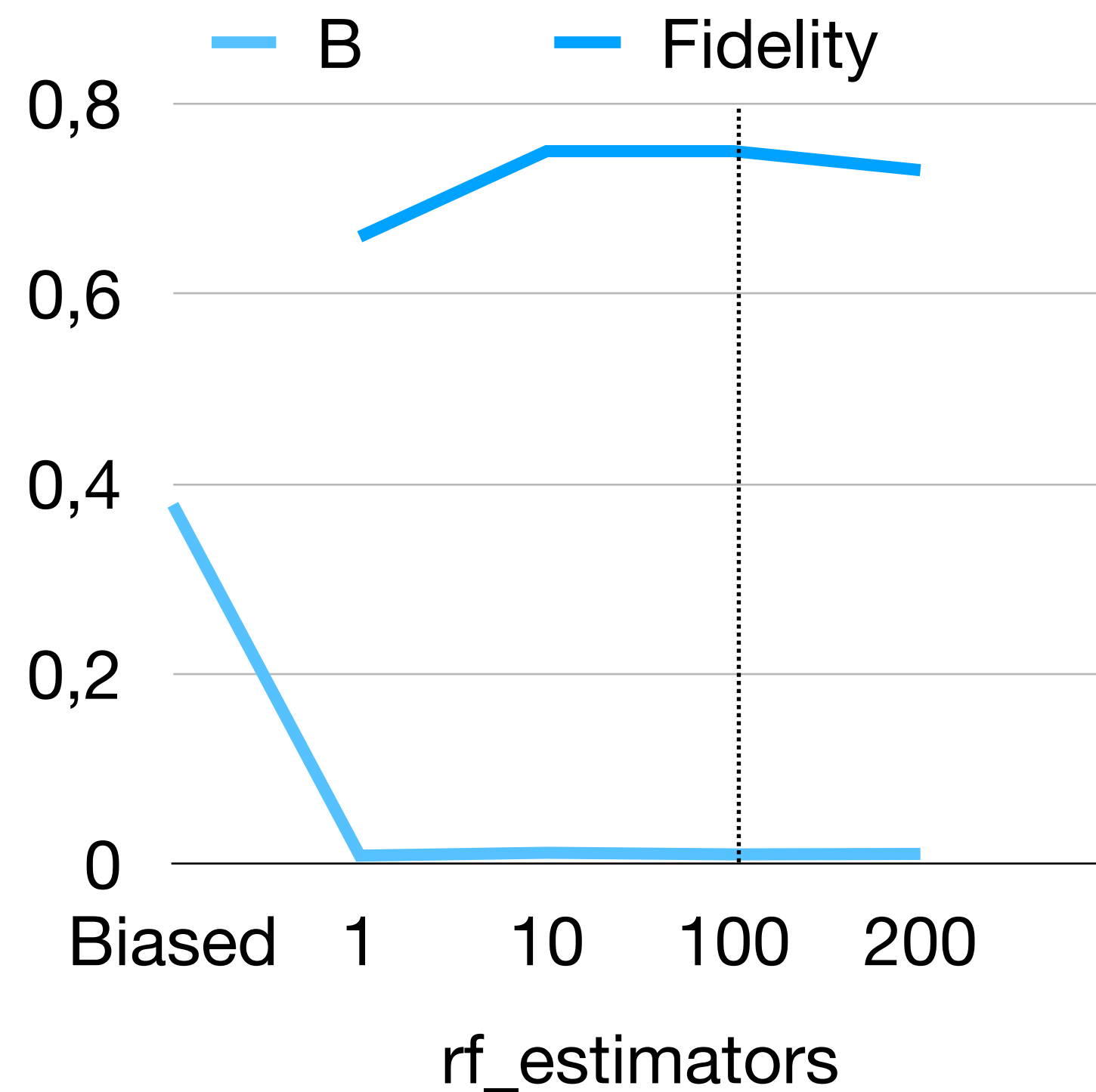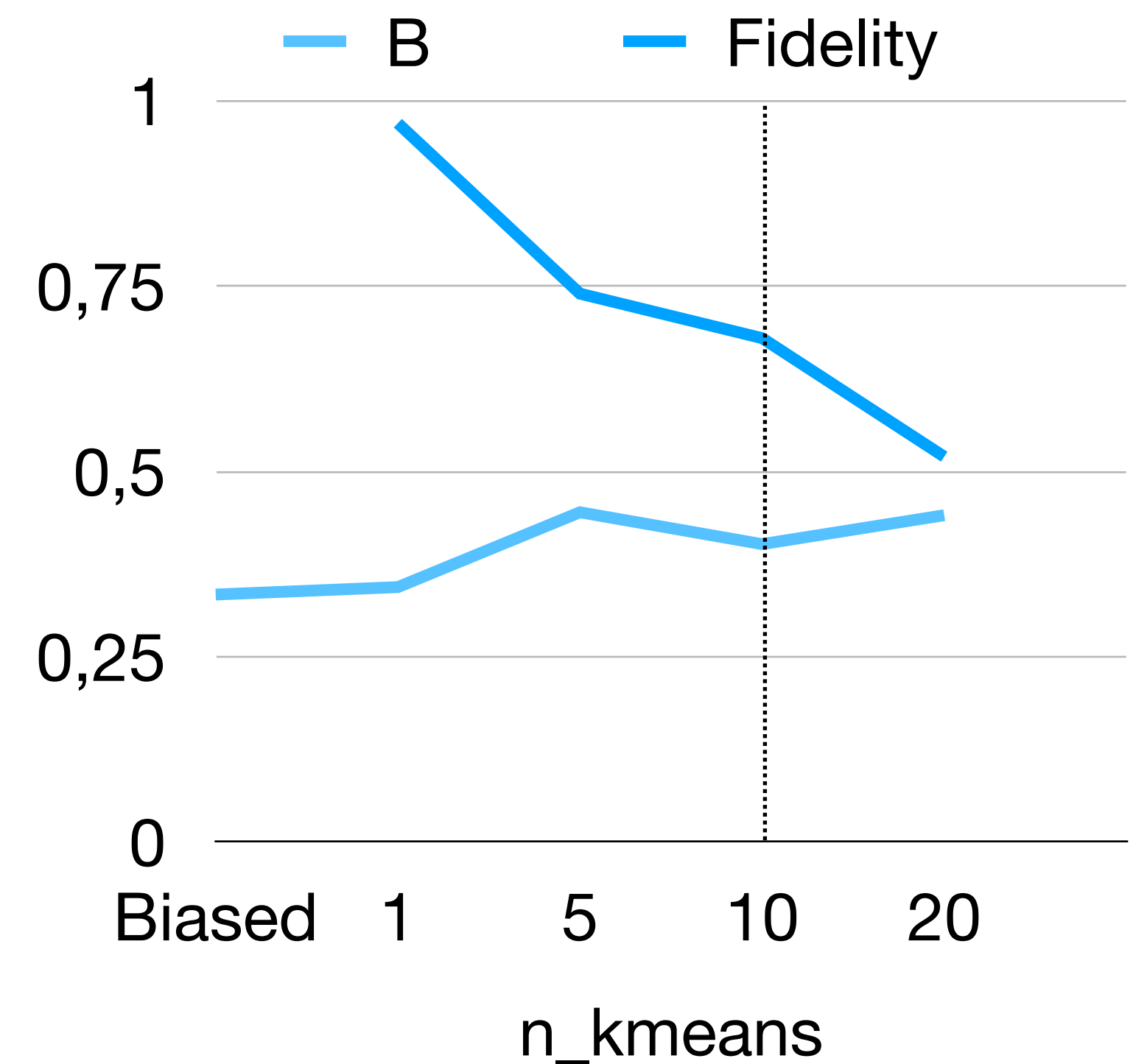LIME perturbation (fidelity: 0.75)
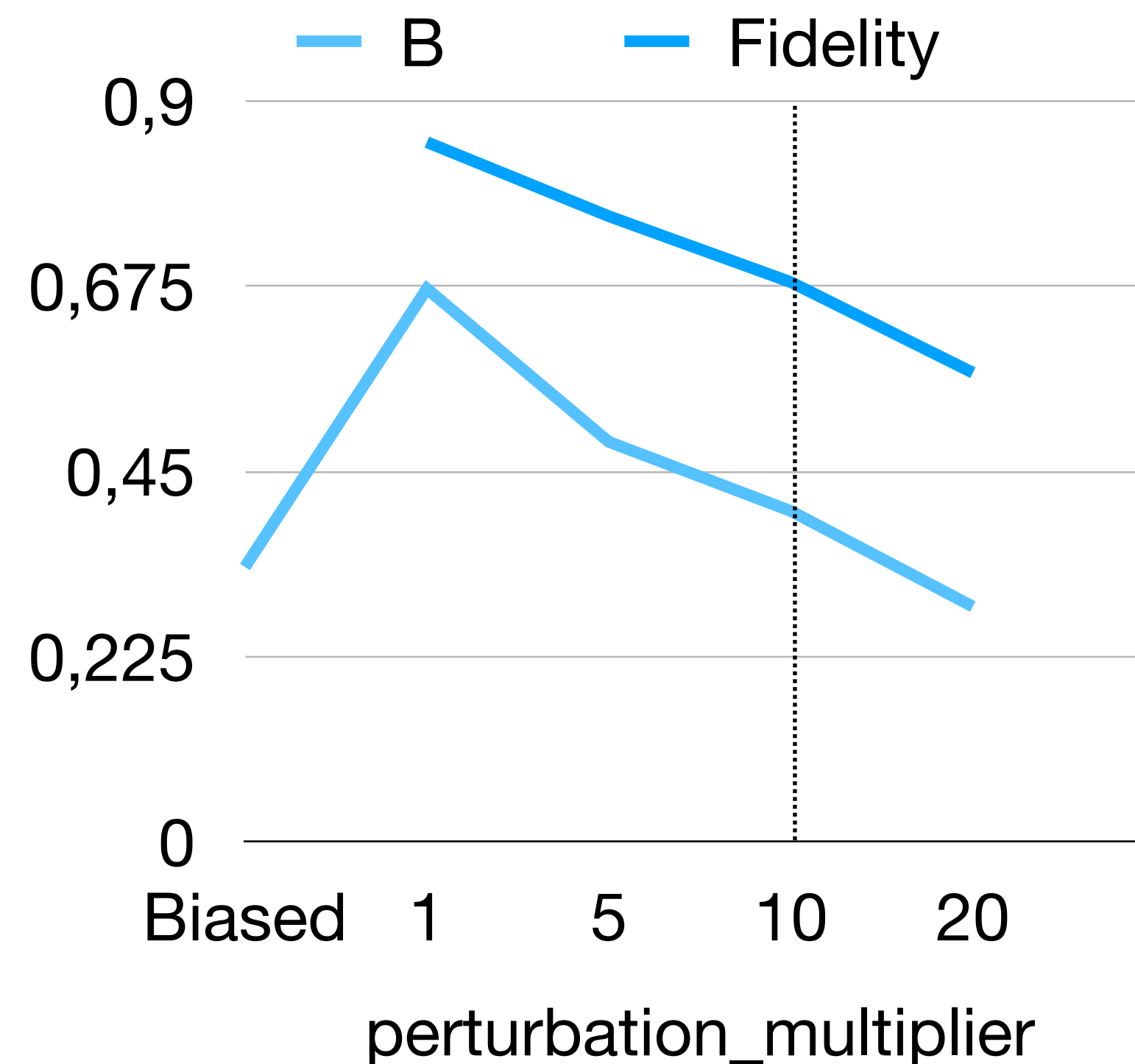


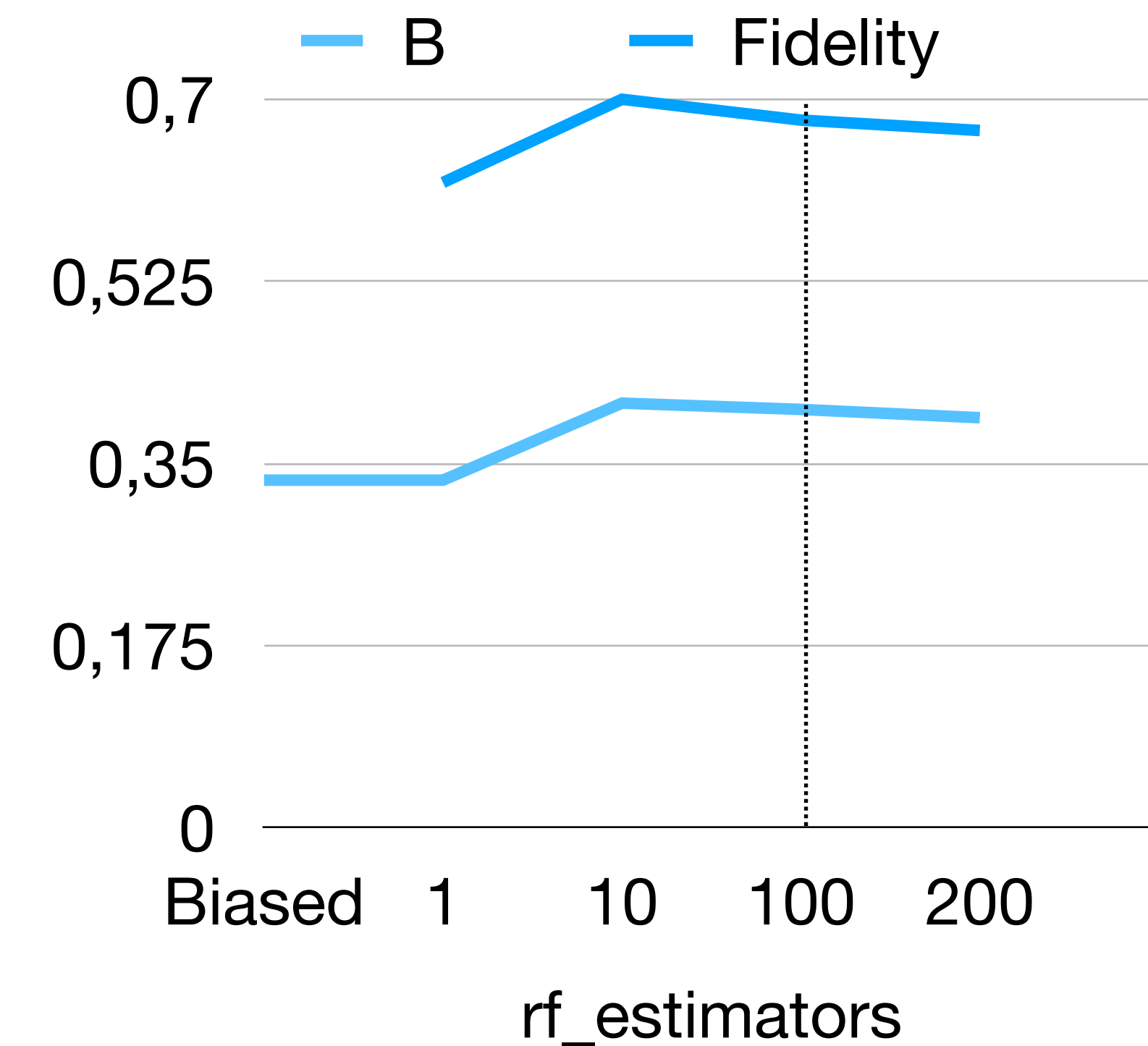PDP perturbation (fidelity: 0.10)

# Hyperparameter Sensitivity
## LIME
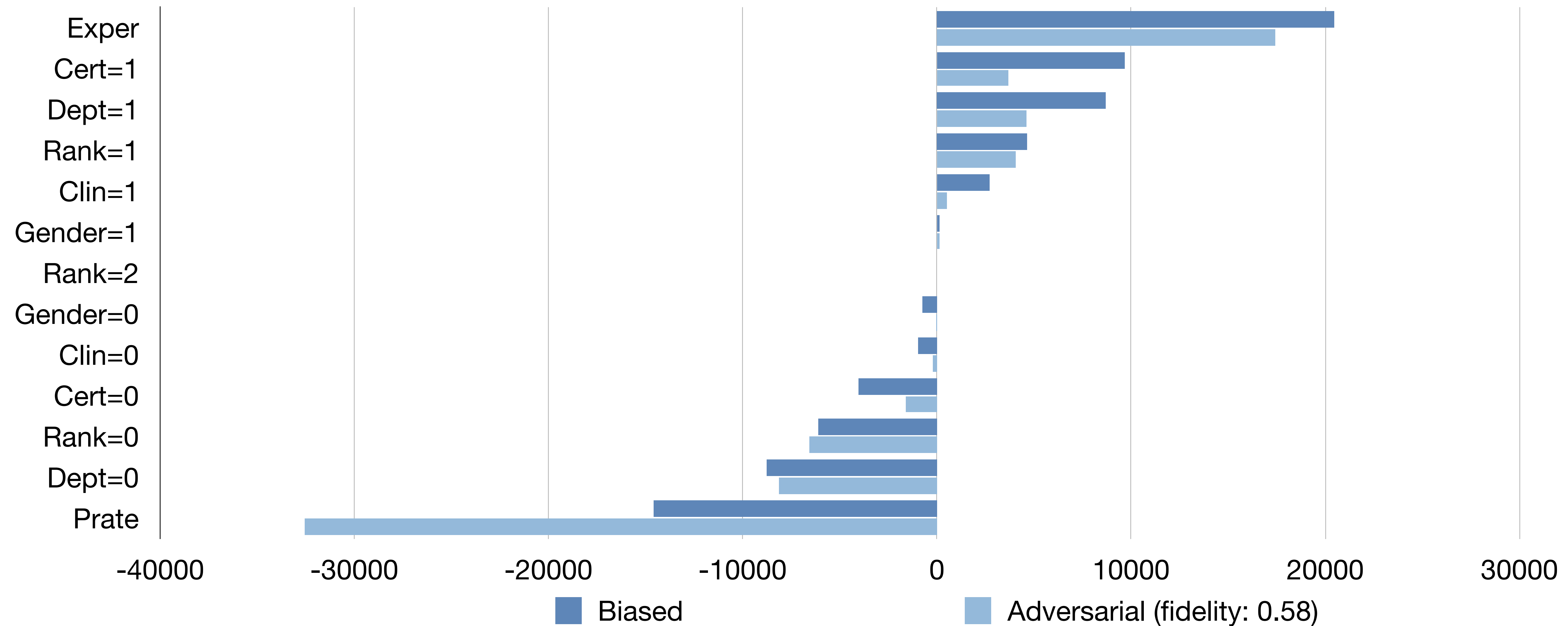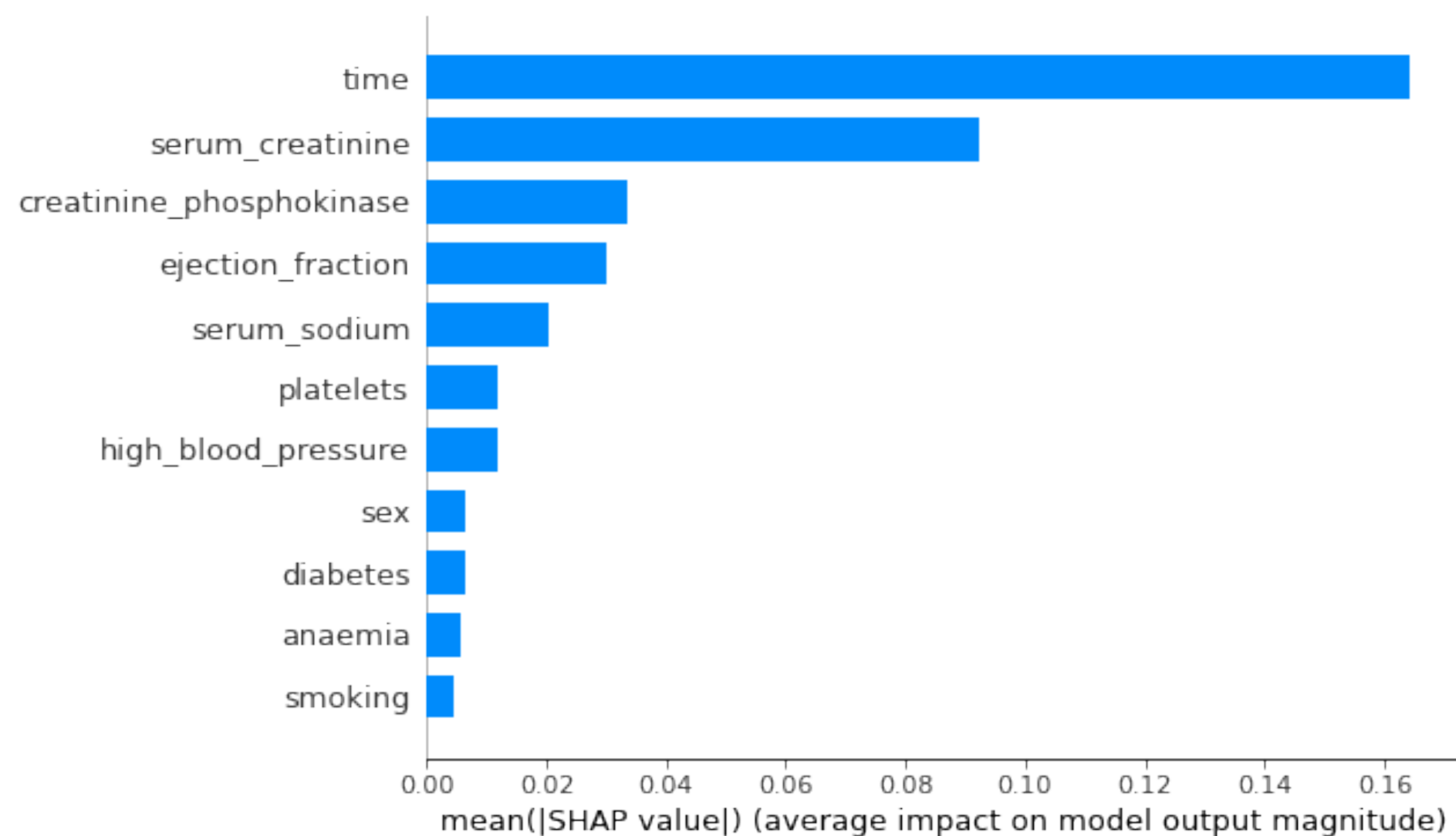
# Hyperparameter Sensitivity
## SHAP

# New Datasets

- Gender discrimination dataset

  - All female doctors at Houston College of Medicine who claimed that the College has engaged in a pattern and practice of discrimination against women in giving promotions and setting salaries.

- Heart failure prediction dataset

  - Cardiovascular diseases (CVDs) are the number 1 cause of death globally, taking an estimated 17.9 million lives each year, which accounts for 31% of all deaths worlwide.
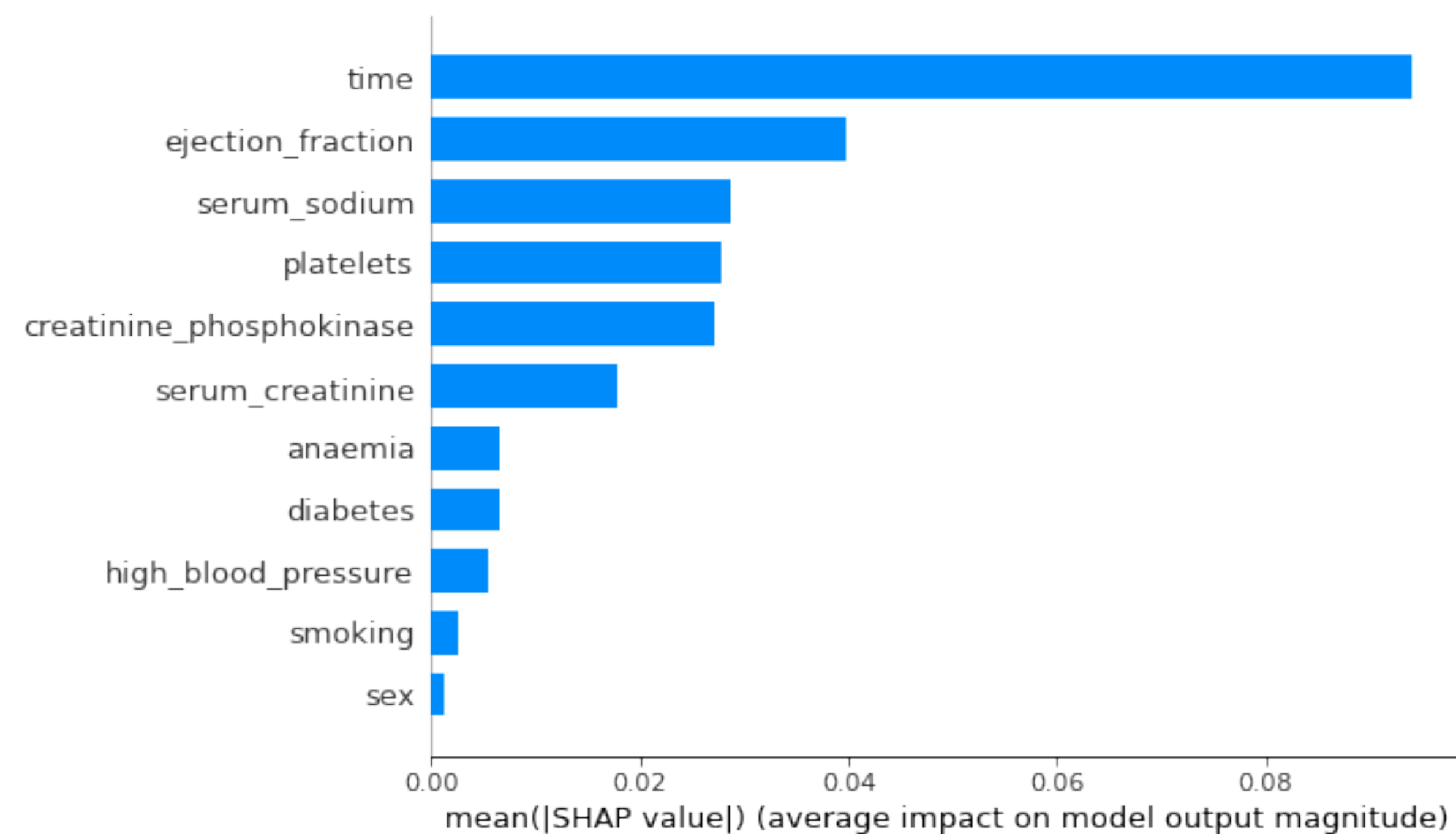
# Gender discrimination dataset

| | Biased | Adversarial (fidelity: 0.58) |

# Heart failure prediction dataset



Biased model

Adversarial model (fidelity: 0.92)

# Thanks for the attention