Combining Counterfactuals With Shapley Values To Explain Image Models

Aditya Lahiri¹, Kamran Alipour¹, Ehsan Adeli² and Babak Salimi¹

¹ University of California San Diego, ² Stanford University,

adlahiri@ucsd.edu,kalipour@eng.ucsd.edu, bsalimi@ucsd.edu,² eadeli@stanford.edu

Algorithmic Decision Making

Individual Harms	Collective Social Harms
Hiring, Housing	Loss Of Opportunity
Credit, Goods' Pricing	Economic Loss
Liberty Loss, Surveillance	Social Stigmatization
Source : Megan Smith, Former CTO of United States.	

Shapley For Images: Challenges

Shift Predictor Algorithm









Shapley Values

3 desirable axioms :

• Efficiency Feature attributions add up the to difference between the individual's outcome and the average outcome.

Experiments And Results

Young

Makeup

Blond

Bald

Male



















• Null Features that do not affect model outcome in any way obtain a shapley value of 0.

• **Symmetry** Pair of features that interact with the model in a similar way, get same attributions.