

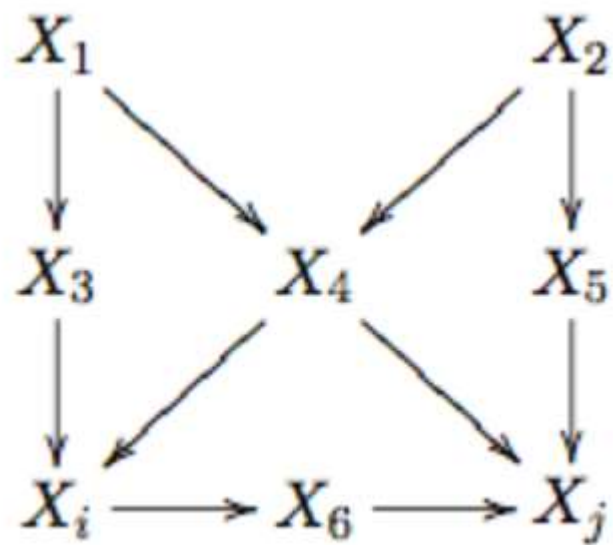
Counterfactual fairness

Shen chuyun

Case: The effect of the treatment

- Z_i : i 个体是否接受治疗
- $\{Y_i(1), Y_i(0)\}$ 表示一组潜在结果, 1是接受治疗, 0是不接受治疗
- 实际观察结果 $Y_i = Z_i * Y_i(1) + (1 - Z_i) * Y_i(0)$
- 对 Z 做随机化
- 平均因果作用 $ACE(Z \rightarrow Y) = E\{Y_i(1) - Y_i(0)\}$

Causal Diagram



箭头表示存在因果关系
有向无环图DAG

$$P(x_1, \dots, x_n) = \prod_{i=1}^p P(x_i \mid pa_i),$$

Do 算子

- do 的意思可以理解成“干预” (intervention)
- 在DAG中 $do(X_i)=x'$, 表示如下的操作: 将DAG中指向 X_i 的所有有向边全部切断, 且将 X_i 的取值固定为常数 x_i'

$$P(x_1, \dots, x_n \mid do(X_i) = x'_i) = \frac{P(x_1, \dots, x_n)}{P(x_i \mid pa_i)} I(x_i = x'_i).$$

- $ACE(Z \rightarrow Y) = E\{Y \mid do(Z)=1\} - E\{Y \mid do(Z)=0\}$

Fairness Background

- V : 可观察的特征
- A : 敏感特征, 如性别, 种族
- X : V 中除了 A 之外的特征
- U : 不可观测的特征

The definition of Fairness

- Fairness Through Unawareness (FTU): 在学习过程中不用A
- Individual Fairness (IF): 相似的人, 相似的预测
- Demographic Parity (DP): $P(\hat{Y} | A = 0) = P(\hat{Y} | A = 1)$
- (Equality of Opportunity (EO): $P(\hat{Y} = 1 | A = 0, Y = 1) = P(\hat{Y} = 1 | A = 1, Y = 1)$

Causal Models and Counterfactuals

- (U, V, F) 定义一个因果模型
- 基于假设, 不存在 W 属于 U , 且 W 是 V 中任意 M 的果
- $V_i = f_i(\text{pai}, U_{\text{pai}})$ pai 是 i 的父节点集合
- Do操作就是 $\text{do}(Z)=z'$, 表示把 F_z 换成 z'
- 给定 $U=u$, Do操作记为 $Y_{Z \leftarrow z}(u)$

Counterfactual

- 反事实: $P(Y_{Z \leftarrow z}(U) \mid W = w)$
- 1. Abduction: for a given prior on U , compute the posterior distribution of U given the evidence $W = w$;
- 2. Action: substitute the equations for Z with the interventional values z , resulting in the modified set of equations F_z ;
- 3. Prediction: compute the implied distribution on the remaining elements of V using F_z and the posterior $P(U \mid W = w)$.

Counterfactual Fairness

$$P(\hat{Y}_{A \leftarrow a}(U) = y \mid X = x, A = a) = P(\hat{Y}_{A \leftarrow a'}(U) = y \mid X = x, A = a),$$

为什么不用这个 因为下面这个算的是一个平均

$$P(\hat{Y} = 1 \mid do(A = a)) = P(\hat{Y} = 1 \mid do(A = a'))$$

Lemma

Lemma 1. *Let \mathcal{G} be the causal graph of the given model (U, V, F) . Then \hat{Y} will be counterfactually fair if it is a function of the non-descendants of A .*

Proof. Let W be any non-descendant of A in \mathcal{G} . Then $W_{A \leftarrow a}(U)$ and $W_{A \leftarrow a'}(U)$ have the same distribution by the three inferential steps in Section 2.2. Hence, the distribution of any function \hat{Y} of the non-descendants of A is invariant with respect to the counterfactual values of A . \square

4.1 Algorithm

Let $\hat{Y} \equiv g_\theta(U, X_{\neq A})$ be a predictor parameterized by θ , such as a logistic regression or a neural network, and where $X_{\neq A} \subseteq X$ are non-descendants of A . Given a loss function $l(\cdot, \cdot)$ such as squared loss or log-likelihood, and training data $\mathcal{D} \equiv \{(A^{(i)}, X^{(i)}, Y^{(i)})\}$ for $i = 1, 2, \dots, n$, we define $L(\theta) \equiv \sum_{i=1}^n \mathbb{E}[l(y^{(i)}, g_\theta(U^{(i)}, x_{\neq A}^{(i)})) \mid x^{(i)}, a^{(i)}] / n$ as the empirical loss to be minimized with respect to θ . Each expectation is with respect to random variable $U^{(i)} \sim P_{\mathcal{M}}(U \mid x^{(i)}, a^{(i)})$ where $P_{\mathcal{M}}(U \mid x, a)$ is the conditional distribution of the background variables as given by a causal model \mathcal{M} that is available by assumption. If this expectation cannot be calculated analytically, Markov chain Monte Carlo (MCMC) can be used to approximate it as in the following algorithm.

- 1: **procedure** FAIRLEARNING(\mathcal{D}, \mathcal{M}) \triangleright Learned parameters $\hat{\theta}$
- 2: For each data point $i \in \mathcal{D}$, sample m MCMC samples $U_1^{(i)}, \dots, U_m^{(i)} \sim P_{\mathcal{M}}(U \mid x^{(i)}, a^{(i)})$.
- 3: Let \mathcal{D}' be the augmented dataset where each point $(a^{(i)}, x^{(i)}, y^{(i)})$ in \mathcal{D} is replaced with the corresponding m points $\{(a^{(i)}, x^{(i)}, y^{(i)}, u_j^{(i)})\}$.
- 4: $\hat{\theta} \leftarrow \operatorname{argmin}_{\theta} \sum_{i' \in \mathcal{D}'} l(y^{(i')}, g_\theta(U^{(i')}, x_{\neq A}^{(i')}))$.
- 5: **end procedure**

At prediction time, we report $\tilde{Y} \equiv \mathbb{E}[\hat{Y}(U^*, x_{\neq A}^*) \mid x^*, a^*]$ for a new data point (a^*, x^*) .

Three level

Level 1. Build \hat{Y} using only the observable non-descendants of A . This only requires partial causal ordering and no further causal assumptions, but in many problems there will be few, if any, observables which are not descendants of protected demographic factors.

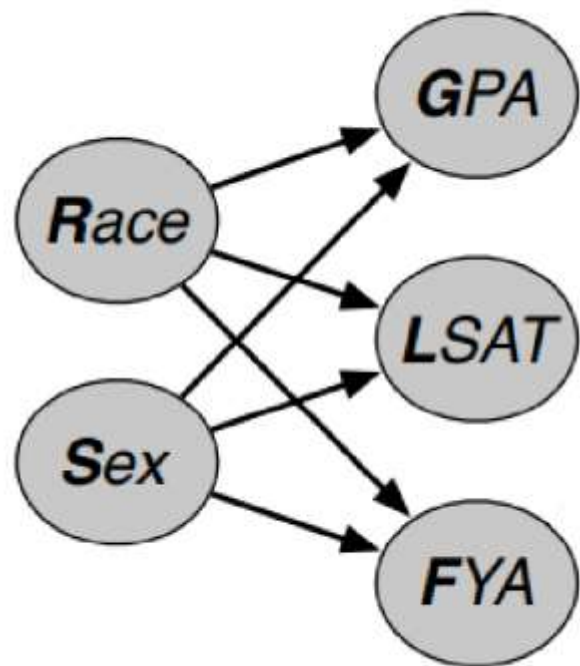
Level 2. Postulate background latent variables that act as non-deterministic causes of observable variables, based on explicit domain knowledge and learning algorithms⁵. Information about X is passed to \hat{Y} via $P(U \mid x, a)$.

Level 3. Postulate a fully deterministic model with latent variables. For instance, the distribution $P(V_i \mid pa_i)$ can be treated as an additive error model, $V_i = f_i(pa_i) + e_i$ [31]. The error term e_i then becomes an input to \hat{Y} as calculated from the observed variables. This maximizes the information extracted by the fair predictor \hat{Y} .

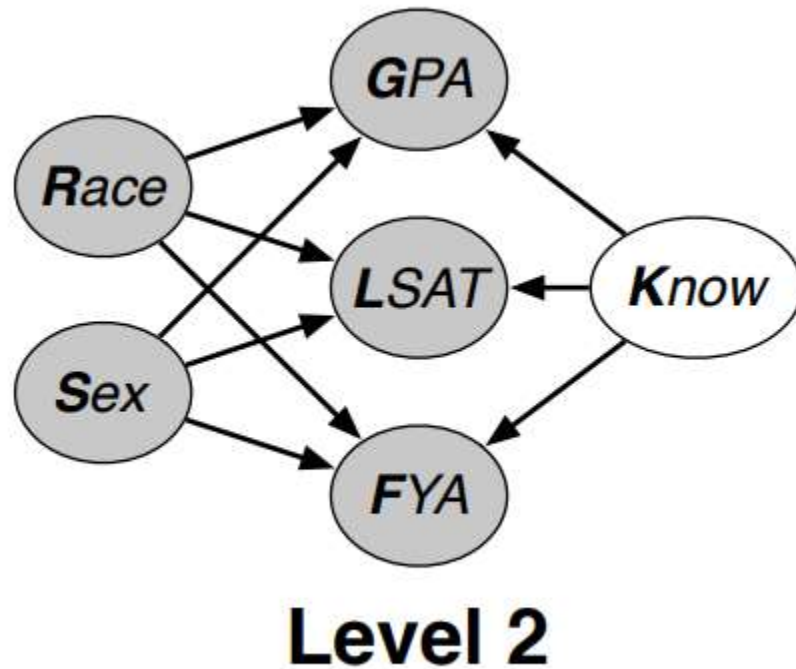
Law School

- entrance exam scores (LSAT)
- grade-point average (GPA) collected prior to law school
- their first year average grade (FYA)
- Given this data, a school may wish to predict if an applicant will have a high FYA.

Law School



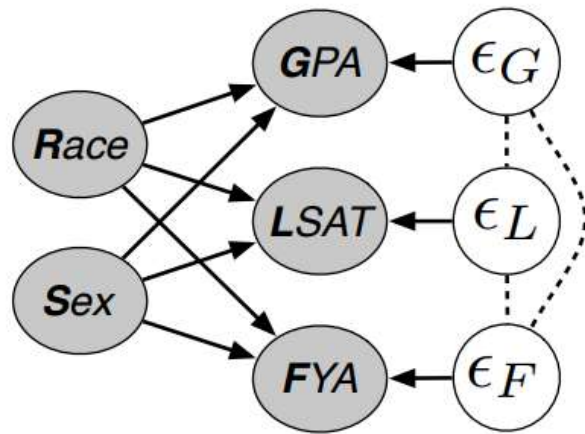
Law School



$$\begin{aligned} \text{GPA} &\sim \mathcal{N}(b_G + w_G^K K + w_G^R R + w_G^S S, \sigma_G), \\ \text{LSAT} &\sim \text{Poisson}(\exp(b_L + w_L^K K + w_L^R R + w_L^S S)), \end{aligned}$$

$$\begin{aligned} \text{FYA} &\sim \mathcal{N}(w_F^K K + w_F^R R + w_F^S S, 1), \\ K &\sim \mathcal{N}(0, 1) \end{aligned}$$

Law School

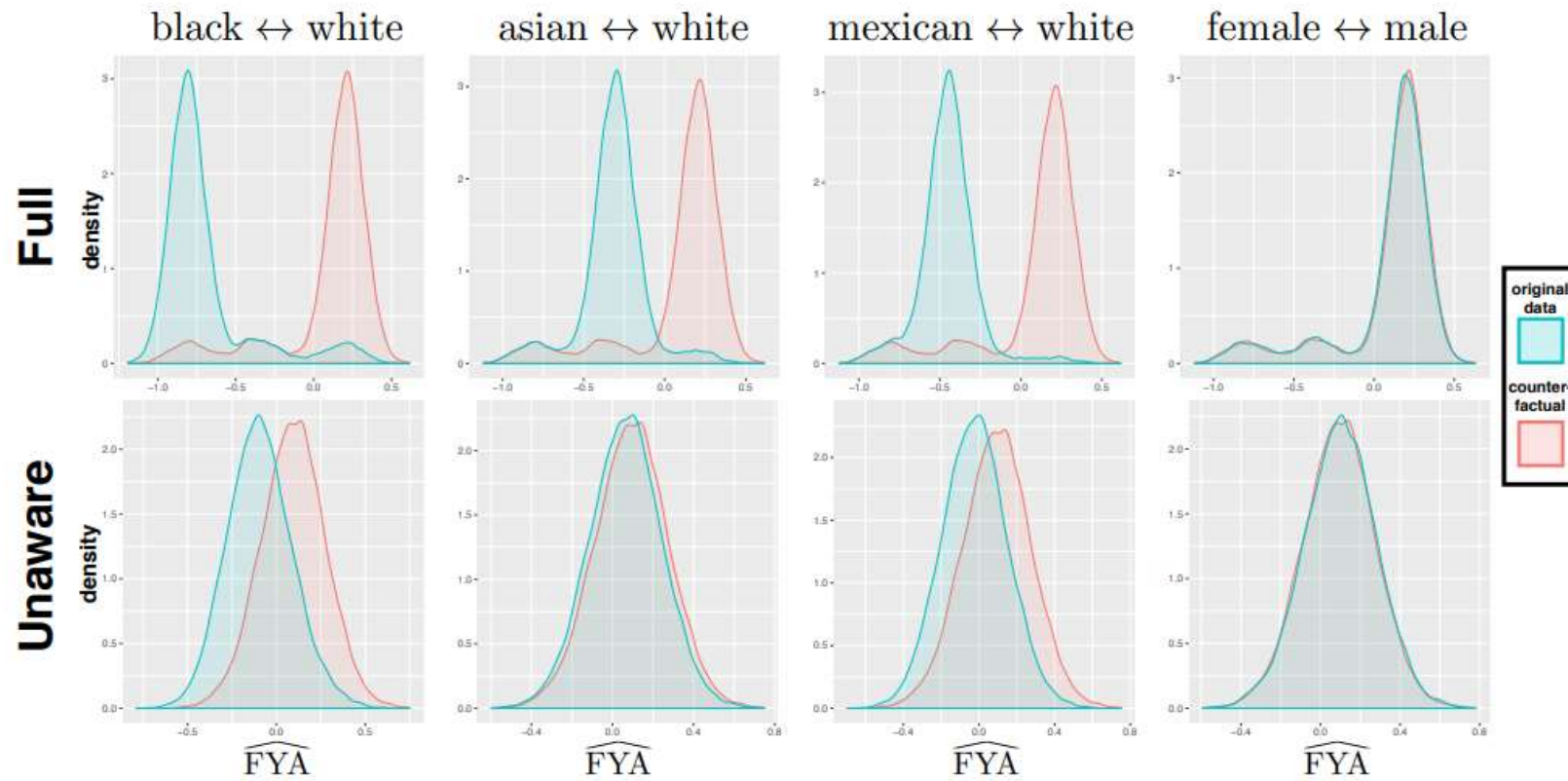


Level 3

$$\begin{aligned}\text{GPA} &= b_G + w_G^R R + w_G^S S + \epsilon_G, \quad \epsilon_G \sim p(\epsilon_G) \\ \text{LSAT} &= b_L + w_L^R R + w_L^S S + \epsilon_L, \quad \epsilon_L \sim p(\epsilon_L) \\ \text{FYA} &= b_F + w_F^R R + w_F^S S + \epsilon_F, \quad \epsilon_F \sim p(\epsilon_F)\end{aligned}$$

We estimate the error terms ϵ_G, ϵ_L by first fitting two models that each use race and sex to individually predict GPA and LSAT. We then compute the residuals of each model (e.g., $\epsilon_G = \text{GPA} - \hat{Y}_{\text{GPA}}(R, S)$). We use these residual estimates of ϵ_G, ϵ_L to predict FYA. We call this *Fair Add*.

Law School



Law School

	Full	Unaware	Fair K	Fair Add
RMSE	0.873	0.894	0.929	0.918