

# Statistical Fit and Algorithmic Fairness in Risk Adjustment for Health Policy

Sherri Rose and Thomas G. McGuire

Department of Health Care Policy, Harvard Medical School, Boston, MA

---

Sherri Rose is Associate Professor, Department of Health Care Policy, Harvard Medical School, Boston, MA, 01201 (E-mail: [rose@hcp.med.harvard.edu](mailto:rose@hcp.med.harvard.edu)). Thomas McGuire is Professor, Department of Health Care Policy, Harvard Medical School, Boston, MA 01201 and Research Associate at NBER, Cambridge, MA, 02138. This work was supported by NIH grant DP2-MD012722 and the Laura and John Arnold Foundation. The authors thank Monica Farid for data preparation.

## Abstract

While risk adjustment is pervasive in the health care system, relatively little attention has been devoted to studying the fairness of these formulas for individuals who may be harmed by them. In practice, risk adjustment algorithms are often built with respect to statistical fit, as measured by p-values or  $R^2$  statistics. The main goal of a health plan payment risk adjustment system is to convey incentives to health plans such that they provide health care services efficiently, a component of which is not to discriminate in access or care for persons or groups likely to be expensive. In an attempt to accomplish this, risk adjustment formulas include indicators for the presence of health conditions associated with higher costs. The salient issue is that incentives mainly operate at a group level, not an individual level; plans can discriminate at the group level in ways they cannot at the person level. Because health plans providing sparse care for certain illnesses is a key policy concern, group-level fit is arguably one of the most important metrics for formula evaluation. Giving primacy on the basis of individual fit when group fit may be the larger concern can lead to harmful decision making. We therefore discuss the role of p-values and statistical fit for this policy problem while considering the fairness of the risk adjustment algorithm for vulnerable groups. Enrollees with mental health and substance use disorders have been found to be subject to the adverse incentives noted above. We apply our ideas to this vulnerable group with a group-level net compensation metric of the incentives to health plans to underprovide services.

Key words: risk adjustment, p-value, health policy, statistical fit, fairness

## 1. INTRODUCTION

Risk adjustment for health policy spans many areas, including public reporting, quality measurement, and payment models. Fair formulas are especially critical in healthcare systems where risk adjustment may have a direct or indirect impact on human health. However, building an “appropriate” risk adjustment formula in a policy environment is not as simple as defining a single overall statistical metric for evaluation. Overreliance on statistical measures of global fit in observational individual-level data does not consider the inequalities created or exacerbated in potentially vulnerable groups. In this context, p-values and  $R^2$  may not be “morally neutral,” to quote Professor Unsworth’s use of this phrase for the subjectivity and ethics of statistical algorithms (Reyes 2016).

For example, valid evaluation of the relative performance of hospitals with respect to quality and outcomes requires risk adjustment for the health of each hospital’s patient population, often referred to as “case-mix” (Shahian and Normand 2008). Hospitals with sicker patients along dimensions not included in the risk adjustment formula will perform poorly in these assessments, leading to payment sanctions in public health insurance programs and possible closure by state regulators. Some approaches to case-mix adjustment for hospital quality explicitly use p-values to select variables in more parsimonious models (e.g., O’Malley et al. 2005).

In this paper, we demonstrate with a straightforward example how potentially misleading and even harmful a focus on improvement in fit at the individual level can be in decisions about how to design risk adjustment in another key setting: health plan payment schemes. While we stay within the realm of statistical fit, we take an expanded view to consider aspects of fairness. We explain that more important than fit at the individual level is fit at the group level, where

groups are considered according to the means a health plan has at their disposal to discourage enrollment and underserve. We therefore implement a metric to assess fair predictions for individuals in vulnerable (i.e., undercompensated) groups, defined by their billed medical conditions. Such a group might be, for instance, individuals with mental health and substance use disorders (MHSUD).

## 2. PLAN PAYMENT RISK ADJUSTMENT

One important way prediction models are used in health policy is for forecasting a health plan's annual spending on an individual enrollee. These predictions feed into determination of the payment the health plan receives in exchange for accepting responsibility for paying for the individual's health care costs, referred to as a capitation payment. Reasonably accurate prediction is necessary to ensure that private health insurance plans, many of which, in the U.S., are part of for-profit firms, will be willing to accept and provide good care to enrollees with high costs. In the U.S., Medicare Advantage (for older and some disabled adults), Medicaid Managed Care (for people with lower incomes) and the Marketplaces (for otherwise uninsured, created as part of the Affordable Care Act) all operate as individual health insurance markets with competing health plans. This form of health insurance sector organization, otherwise known as regulated competition, is universal in Belgium, Germany, the Netherlands, Switzerland, and Israel, among other countries. In all of these settings, some form of risk adjustment is used to determine plan payments.

The Marketplaces, organized at the state level, are the smallest of the three individual health insurance markets in the U.S. structured according to principles of regulated competition. Many participating plans are quite small – on the order of 10,000 enrollees. Inappropriate

incentives are particularly dangerous in this market because plans' conformation to regulations about benefits (for example, breadth of provider networks) is difficult to monitor. Furthermore, certain Marketplaces are served by few plans, and these plans move into and out of the market as they are attracted by profit and deterred by risk, putting added pressure on the design of the payment model.

Plan payment in the Marketplaces is complex, involving adjustments for geographic factors, premiums the plans collect from enrollees, and market shares of very high-cost cases, among other factors (Layton et al. 2018). At the core of the payment scheme, however, is the Department of Health and Human Services Hierarchical Condition Category (HSS-HCC) prediction model, which determines the base payment for each individual before adjustments (Centers for Medicare and Medicaid Services 2016). The HHS-HCC model is a linear least squares regression built using health insurance claims data where the outcome is equal to plan spending on a person per year. The right-hand side variables (risk adjustors) are age x gender cells, disease indicators (i.e., the HCCs), and selected interactions.

Each HCC is the result of a mapping from some of the thousands of the five-digit International Classification of Disease and Related Health Problems (ICD) diagnoses reported on claims to a much smaller number of categories. For example, the HCC for "major depressive, bipolar, and paranoid disorders" is generated from over 50 ICD-9 flags. It is important to note that not all ICD-9 (or ICD-10, adopted in 2015) codes map to an HCC used for payment. This has been shown to be problematic with respect to accurate payments, especially for MHSUD in the Marketplaces (Montz et al. 2016). Montz et al. (2016) found systematic underpayment of enrollees with MHSUD; 80% of individuals with MHSUD are not recognized by the

Marketplace system, contributing to undercompensation for individuals with MHSUD on average.

The HHS-HCC model undergoes regular evaluation, as does its progenitor, the Centers for Medicare and Medicaid (CMS) version used for Medicare Advantage, including consideration of the HCC diagnostic adjusters. Notably, the HCCs used in the HHS-HCC model are a subset of the full 264 HCCs defined in the full system mapping. Adding or subtracting HCCs from the right-hand side of the risk adjustment formula is a component in the evaluation of the prediction function. Government reports lay out the criteria used in defining HCCs, with the first two being the HCC should be “clinically meaningful” and “predictive” (Pope et al. 2004; Ellis et al. 2018). Other risk adjustment formulas consider additional variable types for special populations, such as measures of functional status in the CMS frailty model (Kautter and Pope 2004).

### 3. STATISTICAL FIT AND POLICY CRITERIA

Although researchers and regulators recognize a number of criteria for inclusion of variables in risk adjustment formulas, in practice, the most influential statistical criterion appears to be whether a given variable improves performance, as measured by p-value or  $R^2$  statistic. These metrics are automatically produced by statistical programs, but no metric is automatically generated to capture performance in other, potentially more important dimensions. Simply put, using a p-value or  $R^2$  as the statistical criterion to decide variable inclusion in the prediction function can lead to mistakes, in the sense that improving fit at the person level may not improve fit in the ways that may have more impact, at the group level.

We do not reiterate many of the arguments on the disadvantages of p-values, for example, as enumerated in Wasserstein and Lazar (2016). However, before defining our group-level fit metric to evaluate fairness, we do highlight two features of p-values that are particularly relevant for plan payment risk adjustment. Firstly, as is well-known, large sample sizes will yield significant p-values for variable coefficients that are of trivial magnitude (e.g., Chatfield 1995; van der Laan and Rose 2010). Many plan payment risk adjustment formulas, including those built for the Marketplaces, are created using millions of observations; assessing a variable's importance for predicting health spending using p-values is effectively useless. The second is regarding the relationship between  $R^2$  and p-values. We keep in mind that substantial improvements in  $R^2$  are accompanied by significant p-values for the added variable, but also that miniscule improvements in  $R^2$  may result from the addition of a variable with a nonsignificant *or* significant p-value.

The most important criteria for evaluating a risk adjustment scheme follow from the efficiency or fairness problem risk adjustment is trying to fix (Layton et al. 2017). In all of the health care systems listed earlier, health plans are prohibited by law from discriminating in enrollment or services against *individuals*. For example, in the sectors mentioned, plans must accept any individual who applies for membership. However, plans can and do discriminate against *groups* of individuals, such as those with MHSUD, by: (a) limiting provider networks treating this disorder, (b) setting low provider payments to mental health providers to discourage supply, (c) providing less favorable coverage of drugs in the plan drug formula, and other means. Consequently, individual-level fit is secondary to group-level fit as a metric for alternative plan payment schemes.

In risk adjustment research in the U.S., group-level fit is often measured by predictive ratios, equal to the ratio of predicted values over actual values for a group. A predictive ratio less than one indicates that the prediction function underpredicts and will therefore underpay for the group. Layton et al. (2017) report predictive ratios for four chronic illness groups: persons with cancer, diabetes, heart disease, and mental illness using an updated version of the data used to estimate the HHS-HCC model for the Marketplaces (and an earlier version of the data used in this paper). Underpayment as measured by the predictive ratio was most severe for the mental illness group. This is concerning because, as noted, although a Marketplace plan must accept all applicants, the plan can provide poor care for mental illness, discouraging individuals with mental illness from seeking enrollment in the first place.

In Europe, it is more common to measure group fit by net compensation, equal to the average difference between predicted values for a group and actual values:

$$\text{Net Compensation} = \frac{\sum_{i \in g} \hat{Y}_i}{n_g} - \frac{\sum_{i \in g} Y_i}{n_g},$$

where  $\hat{Y}_i$  is predicted spending for individual  $i$ ,  $Y_i$  is observed spending for individual  $i$ , and  $n_g$  is the sample size for the group of interest (Layton et al. 2017). The sums are taken over all individuals in the group. Net compensation measures incentives to a plan to provide good service to a group. We use this metric in our demonstration as it is on the same scale as our outcome of interest, and therefore has an easy interpretation. Other metrics of health plan performance have been implemented or proposed in the case of ensuring equal access to mental health care, but these are recognized as incomplete and of doubtful effectiveness as a basis for monitoring plan services (McGuire 2016). While many fairness metrics focus on classification problems, individual vs. group-based notions of fairness have been studied (e.g., Zemel et al. 2013; Hu and Chen 2017), as well as general frameworks that include non-classification problems (e.g., Kusner



et al. 2017). We refer to Mitchell (2017) for a didactic summary of fairness metrics in the machine learning and computer science literature.

#### 4. APPLICATION OF THE CRITERIA

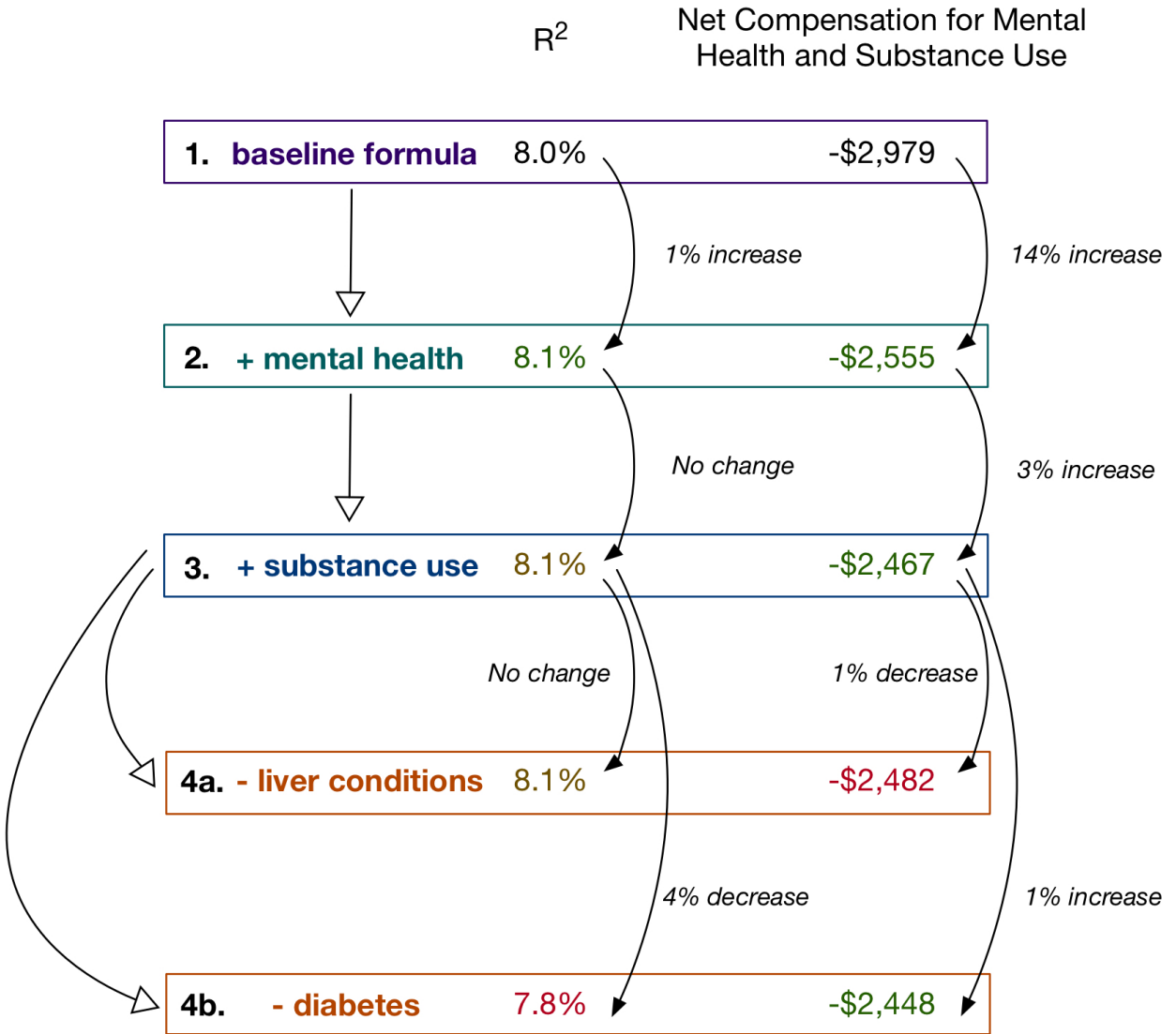
We include an instructive example using the Truven MarketScan data for risk adjustment, which mirrors the approach implemented in the Marketplaces with some modification. We studied a prospective risk adjustment model, using this year's disease indicators to predict next year's spending, whereas the Marketplace risk adjustment model is concurrent, i.e., using disease indicators to predict the same year spending. While our empirical model diverged from actual practice in the Marketplaces, it aligned more closely with the way risk adjustment is commonly done. In the U.S. Medicare Advantage program, as well as all the countries mentioned above with universal systems, risk adjustment is prospective. Our sample contained all those continuously enrolled from January 2015 to December 2016 who had prescription drug coverage and mental health coverage. We excluded enrollees with missing geographic region or capitated claims information, as well as those with negative capitated claims. With these restrictions, we had a total eligible sample size of 5,495,135 that we used in our analysis. The outcome, total annual expenditures in 2016, was calculated by summing all inpatient, outpatient, and drug payments. Median total spending for an individual adult not eligible for Medicare was \$2,363.

The baseline formula included 75 HCCs, as well as age and sex, as predictor variables in a main terms parametric regression. We additionally considered two mental health HCCs and two substance use disorder HCCs. These 79 HCCs followed the CMS-HCC risk adjustment formula in place for Medicare Advantage (Pope et al. 2011). We defined the MHSUD group for calculation of net compensation using Clinical Classification Software (CCS) categories, a more

comprehensive set of variables compared to the HCCs. Each ICD flag maps to a CCS category, unlike the mapping from ICD to HCC described in Section 2. Therefore, our calculation of net compensation will capture the impact of the risk adjustment formula for those with MHSUD recognized *and* unrecognized by the formula. (We emphasize here that individuals with MHSUD but no MHSUD ICD flags will not be captured by HCCs or CCS categories.) The MHSUD group contains 17.3% of the sample, compared to the 2.7% of the sample identified using the four HCCs. Median total spending for an individual in the MHSUD group was \$3,947, which was 67% higher than median spending for an individual in the total sample.

We describe possible iterative decision-making processes in Figure 1 guided by different metrics (that could be defined a priori and within cross-validation, e.g., to avoid cherry picking and overfitting). As noted earlier, with a sample size of over five million enrollees, p-values are a largely ineffective metric. All but four variables in the baseline formula were significant, even those with small event rates, such as “pressure ulcer of the skin, with necrosis through to muscle, tendon, bone” with 171 enrollees. The most prevalent HCC was “diabetes without complications,” occurring among 5.2% of individuals in the sample. This baseline formula had an adjusted  $R^2$  of 8.0%, which is lower than the 12%  $R^2$  typically achieved with the same specification in Medicare. The CMS-HCC formula was originally created to optimize fit on the older, sicker population in Medicare, not the younger, generally healthier population in our MarketScan data. Net compensation for MHSUD was a nontrivial underpayment, reflected in the negative value of -\$2,979, which was 75% of the median total spending in the MHSUD group. Therefore, enrollees with MHSUD are vastly underpaid relative to those without MHSUD, giving insurers a strong incentive to distort their plan offerings to avoid these enrollees.

**Figure 1. Decision-Making Flow Chart for Plan Payment Formula with Differing Metrics**



A formula building exercise that aims to maximize R<sup>2</sup> or optimize MHSUD net compensation toward zero might proceed from the baseline formula (formula 1 in Figure 1) to formula 2, which added two mental health HCCs (“schizophrenia” and “major depressive, bipolar, and paranoid disorders”). This provided a small 1% increase in R<sup>2</sup> and a 14% increase in MHSUD net compensation; p-values for both variables were statistically significant. Because MHSUD net compensation calculated based on formula 1 was a large negative number,

*increases* in MHSUD net compensation that did not lead to a value greater than zero were improvements. Moving from formula 2 to formula 3, we added two substance use HCCs (“drug/alcohol psychosis” and “drug/alcohol dependence”). The  $R^2$  did not change; an algorithm maximizing based on  $R^2$  (and parsimony) would revert back to formula 2. However, there was a continued increase in MHSUD net compensation, thus an algorithm optimizing this metric would keep formula 3 and continue.

Moving from formula 3, which contained the complete set of 79 HCCs used in the Medicare Advantage risk adjustment formula, the algorithm might consider so-called deletion steps (Miller 2002; Sinisi and van der Laan 2004) of statistically significant variables to improve fit and fairness. Formula 4a removed the HCCs associated with liver conditions with no impact on  $R^2$ , but a 1% decrease in MHSUD net compensation (a movement in the wrong direction). An algorithm driven by  $R^2$  and parsimony would prefer such a formula, despite being worse than formula 3 for MHSUD net compensation (and likely worse for those with liver conditions). An alternative deletion step from formula 3 might be to remove diabetes HCCs, represented in formula 4b. Here,  $R^2$  dropped by 4%, yet MHSUD net compensation improved by 1%, and was indeed the best MHSUD net compensation of all the formulas in Figure 1. An algorithm that is searching based solely on improving MHSUD net compensation selects formula 4b, despite its overall poorer global fit (as assessed by  $R^2$ ) and impact on those with diabetes.

## 5. DISCUSSION

There are many technical and ethical issues to address when considering the use of p-values and other measures of statistical fit for algorithms that have a real impact on human lives. One should also be transparent about the goals of their work. To paraphrase Jackie (2018) and apply

it to risk adjustment for plan payment: are we seeking to (1) justify the use of these tools, (2) solve a technical issue with the algorithms, or (3) propose they be dismantled? In this paper, we highlighted the *technical issue* of metric selection and how decision rules driven by global statistical fit can disadvantage vulnerable groups. It is indeed a valid question to take a step back and ask if we could *dismantle* plan payment risk adjustment. What other system might be fairer? Dismantling risk adjustment but leaving the rest of the regulation of individual health insurance markets untouched would create strong incentives to health plans to discourage the enrollment of sicker persons. Health plans have ways to do just this, such as underfunding services used by these enrollees. However, when it comes to *dismantling*, health policy discussions go beyond *dismantling risk adjustment* to questioning the overall policy of individual health insurance markets and regulated competition. In the U.S., alternative models based on a single payer (“Medicare for all”) or organizing competition at the group rather than an individual level (mimicking employer-based health insurance) would address incentives for underservice for certain groups by other means and would not require risk adjustment. These policy alternatives have other disadvantages, however, and we refer readers to Cutler (2002) and Newhouse (2002) for a comprehensive discussion of these tradeoffs.

Within the context of individual health insurance and the *existing* risk adjustment scheme, we assert that there is a pressing need to consider a formal ensemble of metrics for evaluation of plan payment risk adjustment that balances both global fit and multiple fairness metrics. It is well-accepted that evaluation of risk adjusted capitation models for plan payment involves numerous criteria. While regulators can and do make good faith efforts to examine vulnerable groups using predictive ratios and other measures (Layton et al. 2016), most of these procedures are ad hoc. Our plan payment risk adjustment demonstration represented a

simplification of a deeply difficult policy problem and considered only one vulnerable group – those with MHSUD. Even in this simplified example with only one fairness metric, the “best” MHSUD net compensation was still a large underpayment relative to median total spending, further highlighting the challenge of deploying a more comprehensive system. While unlikely to solve the totality of the fairness issues plaguing plan payment risk adjustment, we also note that there are numerous planned changes for the 2019 Medicare Advantage plan payment risk adjustment formula, including adding additional MHSUD HCCs (Centers for Medicare & Medicaid Services 2017). Fairness issues will also not be addressed by simply using machine learning to estimate the risk adjustment formula or perform variable selection (e.g., Rose 2016; Shrestha et al. 2017) if standard statistical fit metrics are still the basis of evaluation for those tools.

The broad issues facing plan payment risk adjustment are not entirely unique given the pervasive use of risk adjustment in health policy. Consideration of what makes a formula fair for the lives these health care algorithms touch has thus far been comparatively underdeveloped. Important work has been done in the areas of predictive policing (e.g., Lum and Isaac 2016), recidivism (e.g., Chouldechova 2017), and hospital ratings (e.g., Phillips 2018), for example. Leveraging and adapting these vital advances while expanding fairness approaches for the distinctive needs of health plan payment is a crucial issue moving forward. The Association for Computing Machinery recently issued a statement on automated decision-making describing seven principles for algorithmic transparency and accountability: awareness, access and redress, responsibility, explanation, data provenance, auditability, and validation and testing (Association for Computing Machinery 2017). It will be especially fruitful to bring these principles to bear in

a context where the algorithms have a real impact on the welfare of vulnerable groups in the health care system.

## References

- Association for Computing Machinery (2017), "Statement on Algorithmic Transparency and Accountability," [online]. Available at [https://www.acm.org/binaries/content/assets/public-policy/2017\\_usacm\\_statement\\_algorithms.pdf](https://www.acm.org/binaries/content/assets/public-policy/2017_usacm_statement_algorithms.pdf).
- Centers for Medicare & Medicaid Services (2016), "HHS-Operated Risk Adjustment Methodology Meeting Discussion Paper," [online]. Available at <https://www.cms.gov/CCIIO/Resources/Forms-Reports-and-Other-Resources/Downloads/RA-March-31-White-Paper-032416.pdf>.
- Centers for Medicare & Medicaid Services (2017), "Advance Notice of Methodological Changes for Calendar Year (CY) 2019 for the Medicare Advantage (MA) CMS-HCC Risk Adjustment Model," [online]. Available at <https://www.cms.gov/Medicare/Health-Plans/MedicareAdvtgSpecRateStats/Downloads/Advance2019Part1.pdf>.
- Chatfield, C. (1995), "Model uncertainty, data mining and statistical inference," *Journal of the Royal Statistical Society. Series A (Statistics in Society)*, 158(3), 419-466.
- Chouldechova, A. (2017), "Fair prediction with disparate impact: A study of bias in recidivism prediction instruments," *Big Data*, 5(2), 153-163.
- Cutler, D. (2002), "Equity, Efficiency, and Market Fundamentals: The Dynamics of International Medical Care Reform," *Journal of Economic Literature*, 40(3), 881-906.
- Ellis, R., Martins, B., and Rose, S. (2018), "Risk adjustment for health plan payment." In McGuire, T.G., and Van Kleef, R.C. (eds.), *Risk Adjustment, Risk Sharing and Premium Regulation in Health Insurance Markets: Theory and Practice*, New York, NY: Elsevier.
- Hu, L., and Chen, Y. (2017), "Fairness at equilibrium in the labor market," in *Proceedings of Fairness, Accountability, and Transparency in Machine Learning*, pp. 1-5.
- Jackie (hatfinisher). (2018), "like what is your take on risk scoring, predictive policing, etc? do you see those projects as in need of justification, in need of solving some technical hurdles, or in need of being dismantled / not implemented?" [online Tweet sent 28 February 2:41PM]. Available at <https://twitter.com/hatfinisher/status/968934234181197825>.
- Kautter, J., and Pope, G. C. (2004), "CMS frailty adjustment model," *Health Care Financing Review*, 26(2), 1.
- Kusner, M.J., Loftus, J., Russell, C., and Silva, R. (2017), "Counterfactual fairness," in *Advances in Neural Information Processing Systems*, pp. 4069-4079.
- Layton, T.J., McGuire, T.G., Van Kleef, R.C., (2016), "Deriving Risk Adjustment Payment Weights to Maximize Efficiency of Health Insurance Markets." NBER Working Paper 22642.



Layton, T.J., Ellis, R.P., McGuire, T.G., and Van Kleeef, R.C. (2017), "Measuring Efficiency of Health Plan Payment Systems in Managed Competition Health Insurance Markets," *Journal of Health Economics*, 56, 237-255.

Layton, T.J., Montz, E., and Shepard, M. (2018), "Health Plan Payment in U.S. Marketplaces: Regulated Competition with a Weak Mandate," in McGuire, T.G., and Van Kleeef, R.C. (eds.), *Risk Adjustment, Risk Sharing and Premium Regulation in Health Insurance Markets: Theory and Practice*, New York, NY: Elsevier.

Lum, K., and Isaac, W. (2016), "To predict and serve?" *Significance*, 13(5), 14-19.

McGuire, T.G. (2016), "Achieving Mental Health Care Parity Might Require Changes in Payment and Competition," *Health Affairs*, 35(6), 1029-1035.

Miller, A.J. (2002), "Subset Selection in Regression," Norwell, MA: CRC Press.

Mitchell, S. (2017). "Fairness: Notation, definitions, data, legality," [online]. Available at <https://speak-statistics-to-power.github.io/fairness/old.html>.

Montz, E., Layton, T.J., Busch, A., Ellis, R., Rose, S., and McGuire T.G. (2016), "Risk adjustment simulation: Plans may have incentives to distort mental health and substance use coverage," *Health Affairs*, 35(6), 1022–28.

Newhouse, J.P. (2002), *Pricing the Priceless: A Health Care Conundrum*, Cambridge, MA: MIT Press.

O'Malley, A. J., Zaslavsky, A. M., Elliott, M. N., Zaborski, L., and Cleary, P. D. (2005), "Case-Mix Adjustment of the CAHPS® Hospital Survey," *Health Services Research*, 40, 2162-2181.

Phillips, D. (2018), "At Veterans Hospital in Oregon, a Push for Better Ratings Puts Patients at Risk, Doctors Say," *New York Times* [online]. Available at <https://www.nytimes.com/2018/01/01/us/at-veterans-hospital-in-oregon-a-push-for-better-ratings-puts-patients-at-risk-doctors-say.html>.

Pope, G.C., Kautter, J., Ellis, R.P., Ash, A.S., Ayanian, J.Z., Ingber, M.J., Levy, J.M., and Robst, J. (2004), "Risk Adjustment Of Medicare Capitation Payments Using The CMS-HCC Model," *Health Care Financing Review*, 25(4), 119-141.

Pope, G. C., Kautter, J., Ingber, J.J., Freeman, S., Sekar, R., and Newhart, C. (2011), "Evaluation of the CMS-HCC Risk Adjustment Model," [online]. Available at [http://www.nber.org/risk-adjustment/2011/Evaluation2011/Evaluation\\_Risk\\_Adj\\_Model\\_2011.pdf](http://www.nber.org/risk-adjustment/2011/Evaluation2011/Evaluation_Risk_Adj_Model_2011.pdf)

Reyes J. (2016), "Technologists must do better: Drexel prof on the ethics of algorithms." *Technical.ly Philly* [online]. Available at <https://technical.ly/philly/2016/09/30/kris-unsworth-ethics-algorithms>.

Rose, S. (2016), "A Machine Learning Framework for Plan Payment Risk Adjustment." *Health Services Research*, 51(6), 2358-2374.

Shahian, D.M., Normand, S.L. (2008), "Comparison of 'risk-adjusted' hospital outcomes," *Circulation*, 117, 1955-1963.

Shrestha, A., Bergquist, S., Montz, E., Rose, S. (2017), "Mental health risk adjustment with clinical categories and machine learning," *Health Services Research*, advance online publication. doi: 10.1111/1475-6773.12818.

Sinisi, S.E., and van der Laan, M.J. (2004), "Deletion/substitution/addition algorithm in learning with applications in genomics," *Statistical Applications in Genetics and Molecular Biology*, 3, 1-38.

van der Laan, M.J., and Rose, S. (2010), "Statistics ready for a revolution: Next generation of statisticians must build tools for massive data sets," *Amstat News*, 399, 38-39.

Wasserstein, R.L., and Lazar, N. (2016), "The ASA's statement on p-values: context, process, and purpose," *The American Statistician*, 70(2), 129-133.

Zemel, R., Wu, Y., Swersky, K., Pitassi, T., and Dwork, C. (2013), "Learning fair representations," in *International Conference on Machine Learning*, pp. 325-333.