# Economic Policy Institute

# Race and ethnicity in empirical analysis

## How should we interpret the race variable?

By Trevon D. Logan • June 15, 2022

**Summary:** In trying to understand racial and ethnic groups well enough to write policy that improves their economic outcomes, we have to have a clear understanding of what "race" means in statistical analysis and how the effect of race is measured. Race factors into economic outcomes in complicated ways that even more sophisticated statistical models can't capture. We need to carefully interpret the effect or predictive power of race in measured disparities—in both descriptive and more sophisticated statistical models—because our assumptions affect how we design policy to address racial disparities.

Researchers who seek to more equitably grapple with questions that arise when analyzing the effect of race on economic outcomes and policymakers and advocates who use research to inform their work should heed five key points about race and ethnicity in empirical analysis:

- **Racial and ethnic categorizations are problematic but necessary components of consistent comparative analysis over time.** We continue to use the federal government's broad classification of racial and ethnic groups in many cases because to not use some standard set of definitions would make generalizability across time much more difficult.

- **Average or median measures by race are commonly used to compare groups but are not always all that edifying.** Descriptive comparisons across groups inform much of policy research, with distributional analyses (like measures at the median or 90th percentile) providing more information than simple averages. But descriptive analysis cannot shed light on the factors that may be associated with measures of economic well-being for a group.

- **Regression analysis is more sophisticated than descriptive comparison and can tell us more about the effects of race on economic outcomes by controlling for other factors related to economic outcomes. However, a regression equation leaves the meaning of the race variable coefficient up to interpretation.** Some characteristics are hard to measure, or otherwise not included in the model, which means we can't assume the coefficient measures the direct effect of discrimination on individuals.

- **How we interpret the race variable affects how we design policy to address racial economic disparities.** Whether the race variable in a regression is thought of as a quality of the individual, or as a marker of a group-specific process happening to people of a certain race, is key to how the results of an empirical analysis are interpreted.

- **One promising approach to interpreting gaps in economic outcomes by race is to take a stratification economics lens, where we rely on history and context to interpret the race variable coefficient in regression analysis.** The approach considers how the race variable is shaped by structural factors of the economy that affect large groups, and moves away from difficult-to-measure individual deficits that are somehow assumed to be embedded in disadvantaged group members on average.

# Introduction: Why researchers, policymakers, and advocates all have an interest in the race variable

This essay is targeted to researchers seeking to more equitably grapple with questions that arise when interpreting race and/or ethnicity variables in statistical analysis. But it also lays out important considerations for policymakers and advocates who use research to inform their work. The essay begins with the problems posed by definitions of race in available government data sets, followed by a discussion of the ways race is typically encountered in quantitative research and analysis. In particular, it discusses the limits of interpreting race coefficients (i.e., interpreting the effect or predictive power of race in measured disparities) and the additional considerations needed to contextualize findings.

# Racial and ethnic categorizations are problematic but necessary components of consistent comparative analysis over time

Most data for empirical policy and social science research adheres to federal practice in classifying race and ethnicity. The Office of Management and Budget (OMB) has issued regulations regarding the classifications of race by federal agencies, including the U.S. Census Bureau, which conducts the major household and business surveys used by researchers. The current classifications are white; Black or African American; Asian; Native American, American Indian, or Alaska native; and native Hawaiian or other Pacific Islander, in addition to the "some other race" category in certain circumstances. There are also two ethnicity classifications, Hispanic and non-Hispanic. As such, everyone is a member of both a race and an ethnicity. The OMB *requires* that race data be collected for a minimum of five groups: White; Black or African American; Native American, American Indian, or Alaska native; Asian; and native Hawaiian or other Pacific Islander. They *permit* the U.S. Census Bureau to use a sixth category, which is "some other race." And respondents may now (as of the 2000 Census) report membership in more than one race.

These categories and how the selection of categories for survey respondents is made have changed over time. Looking first at categories, until 1850 the only options in the U.S. census were white and nonwhite and so, historically, we tend to use and assume that all of the nonwhites recorded were Black. Then in 1860, Native American was added as a category. In 1870, Chinese was added, and in 1890 Japanese was added. Filipino, Korean, and Hindu were added in 1920. Mexican was added as a category in 1930. For the 2000 U.S. Census, the Asian or Pacific Islander category was separated into two categories, "Asian" and "Native Hawaiian or Other Pacific Islander." Although the use in federal data is largely standardized, this same taxonomy is used in a large number of social surveys conducted by state and local entities and, by definition, is contained in administrative data (for example, in a state's unemployment insurance system data).

Looking at selection, until the 1960 census, race was determined in census records and in many surveys, by the surveyor. Now, in the census and most other surveys, race and ethnicity are self-reported by the interviewee.

Categorizations are certainly problematic, for example they lump several groups together in one category, but they are also systematic. Having no racial categorization might not necessarily be an improvement over what we currently have, even if it is flawed. To perform comparative analysis over time in between groups it is important to have some consistent measure of race that is in some way time invariant, so that people of a particular classification can always be noted as such. Researchers and nonfederal institutions that collect data by race take these racial classifications from the census definitions of race. If

we recorded race separately and differently, or allowed for complete self-categorization of race, it would really be difficult to have generalizable claims across different data sources or even over time, about racial classification.

# Average or median measures by race are common but not always all that edifying

Race and ethnicity are most often used in empirical analysis as an explanatory variable. In economics, we are typically looking to do one of two things with a racial variable: Describe the group in comparison with some other group (descriptive analysis) or look at relationships between race and other factors to assess how race factors into economic outcomes.

With regard to descriptive analysis, we can look to see the average of some measure by race or ethnicity itself. If we are looking at the fraction of, say, African Americans who are unemployed, we are looking at, for a given a racial classification, the fraction of people who are unemployed. Or looking at wages of Asian American, we collect data on the wages reported by those who self-report themselves to be Asian. We take the average of those, and that is the average wage for those who are Asian. This is descriptive analysis (i.e., it only describes a population versus, say, looking into the factors that may be associated with measures of economic well-being for a group), but it informs much of the work on racial disparities and inequities.

Such descriptions may or may not be what researchers are interested in. This is an aggregation of individuals and we typically interpret this as a group effect, but is it actually a group effect? In the statistical sense this is a group effect by definition—we have taken the average by group membership, and tautologically this is the group's average. But when we start to interpret this in terms of policy, it may or may not be appropriate to think of it this way. A group average does not tell us much about *heterogeneity* (diversity) in these racial experiences or classifications. For example, all white people are placed together in the same category, which includes people who are from Europe, from North Africa, and from the Middle East. These groups would have different historical and social experiences in the United States, and some groups would contain, on average, Americans of different immigration-descent status (second or third or fourth generation, for example). So given the classification itself, we may want to know more about the heterogeneity of those within that racial group, and that would require additional information.

As another, more specific example, suppose one was going to analyze Asians in the Minneapolis-St. Paul, Minnesota metropolitan statistical area (MSA). A significant portion of that group will be Hmong Americans. The Minneapolis-St. Paul community has one of the largest Hmong American populations in the United States. This Hmong population is quite distinct from others who would also be categorized as Asian. Around a quarter of the Hmong Americans in the United States are found to be in poverty, and that is twice the

national average, which itself is higher than the Asian national average for those in poverty (Budiman 2021).

Research has found very high poverty rates among the Hmong population specifically in the Minneapolis-St. Paul MSA (Minnesota Historical Society 2014). These Hmong Americans (primarily from Laos), have a very different historical and political circumstance than those Asians who are of South Asian and East Asian descent. All of these groups fall under the Asian census definition, but have very heterogeneous experiences which make averages limited in the inferences we can draw from them.

Averages derived by sorting populations into broad racial categories have other limitations. For example, consider how classic racial definitions can obscure differences within and between races in wealth. The concentration of Asians varies widely by geographic location, even within cities. If Asians are geographically concentrated in areas where housing prices are relatively high, they will appear to be wealthier simply for geographic reasons. We would want to take these geographic distributions into account, particularly when looking at national data, when describing wealth differences. That is, we would want to hold these "exogenous" factors such as geographic location constant, leaving out the within-group differences.

The limitations of averages and broad descriptive analysis lead researchers to more nuanced analyses that consider other factors that affect the measure in question, whether it be wealth or wages or some other outcome. For example, we could look at wages of white and American Indian or Alaskan Native high school graduates, and compare the difference in those average wages with wage differences among those who have not completed high school. One thing we could infer from such analysis is whether the racial gaps are (or are not) constant over the level of education. Or, going back to the wealth by racial group example, we could expand our analysis to control for geographic location. Analyzing racial gaps in an outcome (say, wages) as mediated by some other condition (say level of education) is common in analysis and forms a central basis for challenging false narratives that we have around racial disparities. For example, policymakers in the past have focused narrowly on promoting education as the means of closing racial wage and wealth gaps. But researchers have found that gaps in wealth by race are not mediated by education. That is, highly educated Blacks and whites have large gaps in wealth, as do those with less education (Hamilton et al. 2018).

Yet even with a more nuanced approach to research, averages have several problems in addition to hiding within-group differences. They can be skewed by outliers, for example, and they might not necessarily get us the answer to the questions that we want answered. They do not tell us anything about the distribution itself by race and ethnicity when we do analysis of averages since we are looking at differences in a measure of centrality. For example, the median—the measure at the 50th percentile of the distribution of the outcome in question—could be more informative than the mean, although we tend to pay more attention to the mean. By telling us the value at the center of the distribution by race, the median would be much less sensitive to outliers than the mean. Knowing that half of the Asian American and Pacific Islander population is above/below some level of wages, income, wealth, etc., is potentially more useful information. In some descriptions, the

distribution could be far more important depending on the question that we are attempting to answer. As such, analysis with distributional measures would be preferred to averages alone.

# Regression analysis can tell us more about the effects of race on economic outcomes but still must be designed and interpreted with care

In more sophisticated analyses involving race, we use regression analysis to look at relationships (i.e., correlation) between race as an explanatory (or independent) variable and one or more response (dependent) variables. In lay terms, we look at how something like what an average hourly wage would be if the worker were a given race, holding other factors known to affect wages (for example, education) constant. In technical terms, we seek to model an outcome $y$ as a linear function (or transformed linear function) of a set of $x$ measures (variables), of which race and/or ethnicity would be one. Such a model can be written as $y = a + b_1x_1 + b_2x_2 + .... + b_nx_n + e$, where $e$ is an error term and $b_1$, $b_2$, etc. are the coefficients for the variables (i.e., a measure of the degree to which that variable affects the outcome measure $y$). Even in this more sophisticated analysis, race is still fundamentally telling us about the average for people in that racial classification, holding the other measures constant (that is, ignoring the way that they might also vary by race). If one is using a linear regression model and there is an outcome variable $y$ and $a$ is the intercept and $x_1$ is the race variable, and you discretely categorize it for all of these races and control for a variety of factors in the other $x$'s, what we estimate with race is the average level effect, intercept shifter for the race variable, $b_1$, conditional on all the other controls.

As just one example, this kind of a model could tell us that being Black reduces wages on average by a certain percent and having a college education increases wages by a given percent but what this model doesn't do is tell us how simultaneously being Black and having a college education affects wages.

Given those shortcomings, we could choose to be more sophisticated at analyzing the distribution. In our wage example, let's say we suspect that getting an extra year of education would not provide as big of a boost to Black workers as it does to white workers. We could do an analysis that encompasses not just how much education boosts wages for the average worker, but how that boost varies by race. In technical terms, instead of having an intercept there would be a model with a slope difference. Such a model would estimate an equation $y = a + b_1x_1 + b_2x_2 + .... + b_n{}^*x_n + c_1(x_1{}^*m) + e$. This is more sophisticated race analysis, as here the model allows for how $y$ changes as $m$ increases or decreases for each race. What would this look like in an analysis? How would one interpret these sorts of models? One is to look at the $b$'s as intercepts and the $c$'s as slopes. The estimate $c_1$ is the slope effect for race. For example, if we know that education is positively

related to wages, and race is negatively related to wages, we might be interested in the wage gradient as a function of race—that is, what is the slope of the wage schedule for Black Americans as a function of their education? In lay terms, how much would each additional year of education increase wages for a Black worker? This sort of specification can lead to an answer to this question.

There is a way in which economists have interpreted the race coefficients, the betas, very differently over time. In most traditional economic analysis, theory would predict that the coefficient on $b_1$ (were $b_1$ Black race) would be close to zero. In other words, once there is a control for all of the choice variables and inputs, say, into a human capital production function, that would explain all the differences in wages that we would observe and there would not be any residual racial effect itself (in other words, that discrimination plays no role in any of the observed differences in wages and other measures). There is very little empirical evidence that the race effect is actually zero. We have very little theory to tell us why race would have an effect on the outcome, outside of theories of discrimination in the labor market or something else that would vary by race (as opposed to by individual) which would lead to such group-based differences. A traditional way of detecting discrimination, in fact, would be to take the beta one coefficient as evidence of discrimination. The controls themselves are the choices that people would actually make and the characteristics that should not be related to labor market outcomes, such as race, should not matter to labor market outcomes if there is no discrimination in the market so that $b_1$ should be zero.

Interpretations of the coefficient on the race variable in wage regressions as being face-value evidence of discrimination have fallen out of favor almost completely in empirical analysis, at least in economics. This is because we now believe that there are many factors related to wages that would be in a theoretical model (for example, someone's marginal productivity) but that are not included in wage regressions because they are difficult to measure and not observed in the data. The thought is that these omitted factors could be correlated with race or ethnicity or drive that level effect that we are seeing by race or ethnicity. These could be omitted noncognitive skills, for example, someone's teamwork, their persistence, their drive, their independence, their personality, their communications skills, etc. Or these could be omitted cognitive skills, school or education inequality, etc. There is a problem with making these arguments as being about race because we are running *individual* level regressions and what you then have to argue is that these noncognitive or cognitive skills or other omitted factors *vary by group* because we have estimated these with an individual regression and we are looking at a *group* effect. The only way that these omitted factors can work to explain these racial differences is if they are correlated with the group itself. That is extremely important and that has to be true statistically.

We have a difficult time in economics thinking about racial inequality, which should explain some of these unobserved factors themselves. For example, the quality of education may be lower in under-resourced school districts serving students of color than in affluent districts where plentiful property taxes fund the latest in facilities and resources. What we could think about doing is parsing out, say that *b1* coefficient into factors that are related to racial inequality and those that are not. But that typically is not the way that economists

proceed: They would like for the market to have none of these sort of structural factors, such as educational inequality; the market should bid them all away.

# How we interpret the race variable shapes how we develop policy

The examples above show that there are different ways to interpret the race coefficient in regression analyses. Is it evidence of racial inequalities, of discrimination, or is it actually caused by some omitted factor or factors? If it did reflect those omitted factors, those factors would have to be correlated with race or ethnicity to drive the result. But when we talk about these measures that vary by race (wages, for example), typically we do not acknowledge that whatever is omitted (thus weakening the predictive power of the model) itself has to be strongly correlated with race or ethnicity to drive the analysis. For example, if we criticize a model for failing to account for the effect of education quality on wages, and take the position that were the model to include that effect, the observed gap between white and Black workers would shrink, we are ignoring the problematic fact that the quality of education for Black children is lower.

There are deeper conceptual problems with the traditional ways of interpreting race variables. Prominent is the framing of measures as the averages of all the individuals who belong to a specific racial category. As Kyle Moore notes in his essay in this collection, that can lead to or enable the misguided notion that disparities are caused by the failures of individual Black people to get enough education or develop wealth-building cultural habits or other factors. One new approach, detailed in Moore's essay, has been to think of outcomes by race not as the average of all the individuals who belong to a specific racial category, but as the outcome of a process that is group specific by nature. In economics, thinking and putting people theoretically into groups in a way that marries empirical analysis with theoretical development models is referred to as stratification economics. It was pioneered by William Darity, Darrick Hamilton, and James Stewart (see Darity, Hamilton, and Stewart 2015) and others, who sought to look at race and ethnic outcomes as arising from group conflict over resources. That somewhat addresses the problem that we currently have in empirical economics because it explicitly provides a role for racial inequality to be related to these omitted factors that we think drive the race coefficient, but that are systemic in nature (e.g., unequal education) rather than individual deficits (e.g., the failure to get enough education). Stratification economists start with the assumption that gaps in outcomes by race are at least partially the result of policy that explicitly took race into account. In this approach, the reason why the race variable matters is because race played and continues to play a role in how we divide resources and how we develop policy.

# Conclusion

The inference and interpretation of the descriptive or even the more sophisticated

statistical models we have used to explain outcomes by race is not a given. We first have to consider how race is defined and then how it is constructed. Race is an aggregation that we may want to divide to understand what is behind the "racial group effect," as a heterogenous range of experiences may be hidden there. Similarly, when moving to more sophisticated analysis at the individual level, the average race effect again requires serious interrogation. If we are analyzing individuals what does the group effect imply about group-based experiences? Are these factors truly omitted and varying at the group level, or do they reflect group processes that are the combination of policy and dynamic processes over time, or perhaps a combination of the two? Thinking more critically about the sources of such differences would enhance our analysis and our design of policy.

# Additional reading and resources

Readers interested in delving deeper into the issues touched on in this chapter are encouraged to explore the following resources suggested by the author.

### Articles

Cook, L.D., T.D. Logan, and J.M. Parman. 2014. "Distinctively Black Names in the American Past." *Explorations in Economic History* 53, issue C: 64–82.

Gullickson, A. 2019. "The Racial Identification of Young Adults in a Racially Complex Society." *Emerging Adulthood* 7, no. 2: 150–161.

Gullickson, A., and A. Morning. 2011. "Choosing Race: Multiracial Ancestry and Identification." *Social Science Research* 40, no. 2: 498–512.

Ward, Z. 2021. "Intergenerational Mobility in American History: Accounting for Race and Measurement Error." National Bureau of Economic Research Working Paper no. w29256.

### Books

Dietz, T., and L. Kalof. 2010. *Introduction to Social Statistics: The Logic of Statistical Reasoning*. Wiley-Blackwell.

Frankfort-Nachmias, C., A. Leon-Guerrero, and G. Davis. 2020. *Social Statistics for a Diverse Society*, 9th ed. SAGE.

Roberts, D. 2011. *Fatal Invention: How Science, Politics, and Big Business Re-*

*Create Race in the Twenty-First Century*. New Press/ORIM.

## Subject matter experts

**William A. Darity Jr.** • Duke University

**Aaron Gullickson** • University of Oregon

**Hedwig Lee** • Washington University in St. Louis

**Dorothy Roberts** • University of Pennsylvania

# References

Budiman, Abby. 2021. "Hmong in the U.S. Fact Sheet." Pew Research Center. April 2019.

Darity, Jr., William A., Darrick Hamilton, and James B. Stewart. 2015. "A Tour de Force in Understanding Intergroup Inequality: An Introduction to Stratification Economics." *Review of Black Political Economy* 42, nos. 1–2. https://doi.org/10.1007/s12114-014-9201-2.

Hamilton, Darrick, William Darity, Jr., Anne E. Price, Vishnu Sridharan, and Rebecca Tippett. 2018. *Umbrellas Don't Make It Rain: Why Studying and Working Hard Isn't Enough for Black Americans*. The New School, Duke Center for Social Equity, and Insight Center for Community Economic Development. August, 2018.

Minnesota Historical Society. 2014. "We Are Hmong Minnesota: FACT SHEET." December 8, 2014.