

Reverse Regression and Orthogonal Regression in Employment Discrimination

Analysis

Abstract: In a recent review article, White and Piette provide an overview of the use of reverse regressions in discrimination-related litigation. They explain the technique, provide a model application, summarize its advantages and disadvantages, and identify litigation in which it has been used. We point out weaknesses in common uses of reverse regression, some of which might cause serious misinterpretations of the data. We suggest that typical interpretations of reverse-regression results are incorrect. We also question the practice of conducting both direct and reverse regressions when studying employment discrimination, since the two approaches make mutually-inconsistent assumptions about the nature of the stochastic error; these assumptions generally can not be corroborated from the data and bias the regressions' results. We suggest orthogonal regression as a potential alternative to the problems associated with direct and reverse regressions in some cases. We then provide a model application of orthogonal regression to discrimination-related investigations.

Reverse Regression and Orthogonal Regression in Employment Discrimination

Analysis

White and Piette (1998) provide a helpful review of the reverse-regression technique's use in employment discrimination analysis. After an explanation of the technique, White and Piette offer an application as a template for the use of reverse regressions, analyzing data on 572 employees from a department of a manufacturing company. They find that direct regression leads to a conclusion that discrimination exists within the department, while the reverse-regression analysis leads to "results (that) are completely opposite from those of the direct regression approach." (White and Piette, 1998, 133) They argue that this paradox is not an artifact of their particular data set, finding that "it is common to find the implications of the direct regression results opposite of those from the reverse regression." (White and Piette, 1998, 133) They then review arguments for and against the use of reverse regressions, and summarize several court cases in which the technique was influential.

This would all seem to make reverse regression a litigant's dream: an established statistical technique that will produce defensible evidence for either side of a discrimination lawsuit. A litigant's dream, and a philosopher's nightmare. One hopes that there really is an objective truth "out there" to be explored via statistical work. Yet a review of the literature might leave the impression that, in discrimination lawsuits at least, mathematics is not much more than a rhetorical device in the service of advocacy.

We would like to point out several shortcomings in common uses of reverse regressions. These shortcomings raise doubts about common interpretations of reverse regressions. We then suggest orthogonal regression as a potential alternative method for investigating some forms of discriminatory practices.

Interpretation of Coefficients

Like nearly all work of this sort, the basic result in White and Piette's application section is produced in the following way:

1. First, a standard regression estimates the significance of the various factors that might influence salaries. A log-linear regression of salary on productivity measures (performance ratings, years of experience, full-time/part-time status) and race (a binary variable, equal to 1 for whites in their sample) is estimated:

$$\ln \text{Salary} = a + b (\text{Performance Ratings}) + c (\text{Experience}) + d (\text{FT/PT status}) + e (\text{Race})^1$$

White and Piette's regression yields a statistically-significant, negative estimate (coefficient=-0.0963, P value = .0037) for the coefficient e on race. Following common practice, this is taken by the authors to indicate discrimination.²

2. Now a reverse regression is constructed, in which the dependent and independent variables from the first regression are reversed. Since regressions by nature require that there be a single dependent variable, this method requires that an index of productivity be calculated to serve as the left-hand side of the new equation. As usual, White and Piette index productivity by the forecast value of $\ln[\text{Salary}]$ from the first regression, using forecasts that ignore the contribution of race to salary. A reverse regression of the productivity index on log salary is then completed:

$$\text{Forecast-}\ln(\text{Salary})|_{\text{race assumed non-white}} = \alpha + \beta (\ln (\text{Salary})) + \gamma (\text{Race})$$

This regression yields a coefficient γ on Race whose *sign* would be consistent with the result from the first regression (i.e., γ is positive 0.0138, e is negative). However, γ is judged to not be statistically significant (P value = .6307), and this is taken by the authors (and the larger literature that they cite) as evidence that the two regressions have reached opposite conclusions about discrimination.³

¹ For simplicity, we will speak throughout the paper of discrimination based upon race. But of course the same techniques and results are easily extended to other forms of discrimination, like gender-based salary or qualification differences.

² The paper interprets this as evidence of discrimination against non-white employees. However, the binary race variable in the first regression takes the value 1 for white employees, 0 for others. Since this regression finds that e is *negative*, the authors have misinterpreted their findings, which in this case indicate discrimination against white employees. Taking the anti-log of the reported coefficient, it would appear that whites earn about \$1101 less per year than non-whites with the same performance ratings, experience and FT/PT status

³ This estimate assumes that salary was measured in thousands, which seems the only reasonable possibility. Since γ in the second regression is reported to be positive (0.0138), the reverse regression *also* indicates that if discrimination exists it is harming the white employees: At a given log-salary, whites generally have a productivity index that is 0.0138 units larger than non-whites' productivity. This would translate into an estimate that salary for whites that is approximately \$1339 lower than salaries for equally-qualified non-whites. (Taking the inverse of the slope of the regression ($1/.0473 = 21.14$) times the race coefficient (the amount by which the intercept for white is smaller than the intercept for others, .0138) gives us .2917 for the implied difference in the log of salary of the two

In what sense are these two regressions reaching “opposite” conclusions? Only in the sense that the second regression’s results do not pass a standard t-test for statistical significance. This phenomenon shouldn’t surprise us. Maddala’s (1992, 461) textbook discussion of reverse regression indicates that “since the direct regression gives biased estimates of (the qualification coefficients), what we have here is a biased index of qualifications (serving as the reverse regression’s dependent variable),” so that “one should not make (statistical-test) inferences... but obtain bounds for (the discrimination coefficient) from the direct regression and reverse regression estimates.” In other words, since the dependent variable in the reverse regression is subject to estimation errors (which we will explore soon), statistical tests in the reverse regression are not valid. In the present example, it does not matter that γ is not statistically significant. Instead we should think of the two discrimination-related coefficients, e and γ , as upper- and lower-bound measures of the true extent of discrimination.

Maddala (1992, 459-461, 71-74) concedes that there are cases in which both direct and reverse regressions must be employed, since one approach may tend to overestimate the crucial coefficient while the other may tend to underestimate it. But some have interpreted this as a *carte blanche* endorsement of the simultaneous use of both regressions. This is unwise. Maddala (1992, 75) provides several guidelines:

1. If we know the direction of causation--which variable, salary or qualifications, was used to sort applicants-- then the variables should go on their proper axis in a single regression. In this case, “the opposite regression does not make sense.” (75) If positions were advertised requiring specific qualifications, and salaries were then influenced by race, the direct regression is appropriate; if instead race affected the qualification thresholds for a job that paid everyone similar salaries, the reverse regression should be paramount.
2. If the direction of causation is not known (both salary and qualifications may have been affected simultaneously by race, with both being joint-normally distributed), or if both variables are measured with error, then both a direct and a reverse regression may be necessary in order to get “bounds” on the race coefficient. The estimates then indicate the extremes of a range of possibilities that exists in the population. If one regression indicates discrimination exists against whites and the other indicates discrimination against blacks, it would be incorrect to say that the two regressions “disagree” or reach opposite

groups; taking the anti-log yields 1.339 for the implied difference in salary.) Given the problems inherent in reverse

conclusions. Together, the two regressions have reached a single conclusion: We do not yet know if there is discrimination or not. So far, the legitimate boundaries for the discrimination coefficient include the number zero. If we want to settle the question of whether discrimination exists, we will have to look to other kinds of information.⁴

Alternatives to Reverse Regression

As White and Piette point out, many other problems have been associated with the use of reverse regression. Econometricians might normally think of instrumental-variables estimation or orthogonal regression rather than reverse regression when the direction of causation is unknown or when both variables are measured with error. (Maddala, 74). Instrumental-variable regression (for cases with measurement error) is well discussed in standard texts, but orthogonal regression (for cases in which the direction of causation is unclear or there is measurement error) is less commonly understood. In fact it is difficult to find a textbook introduction to the method, even among graduate econometrics texts. We present a brief introduction to orthogonal regression, then an application of both reverse and orthogonal regression to an employment-discrimination dataset.⁵

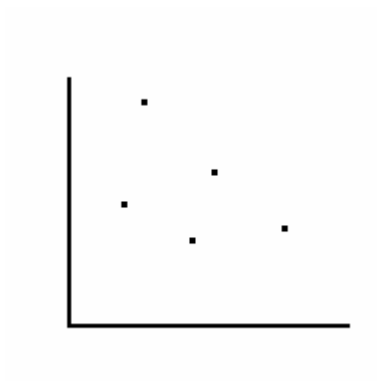
Consider the simplest case: a regression involving a single independent variable. Say that the data look like those in Figure One.

regressions, this estimate seems to be quite close to the direct regression's finding.

⁴ Of course, this was not the case in the White and Piette example, which consistently finds evidence of race-related discrimination.

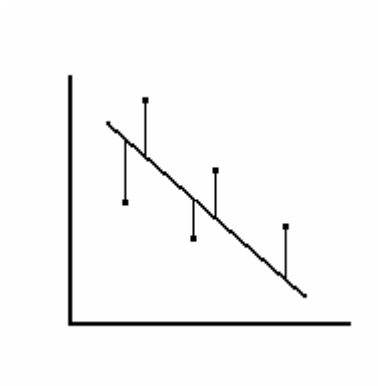
⁵ Unfortunately, White and Piette indicate that their data are confidential; they were not available from the authors for replication. We instead use a standard employee data set of approximately the same size, the EmployeeData.Sav file provided with SPSS version 10.0.

Figure One: Data



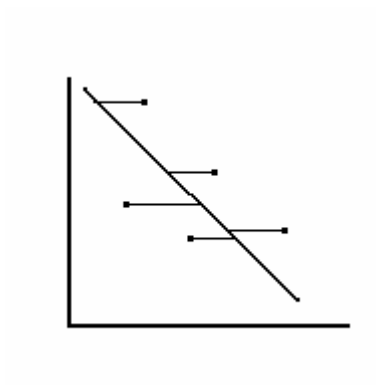
The standard regression model thinks of the independent variable as non-stochastic and measured without error, whereas the dependent variable is stochastic and therefore measured with some error from the true population regression line. We normally proceed by fitting a line that minimizes the squared vertical distances—the estimated squared (vertical) errors--between the data and the regression line, as in Figure Two.

Figure Two: Direct Regression



A reverse regression using the same data would think of the formerly independent variable as the dependent variable—this is what has been “reversed.” Now the variable that has been on the vertical axis is presumed to be non-stochastic and measured without error, whereas the variable on the horizontal axis is presumed to be stochastic and measured with error. The line should now be fit by minimizing the squared *horizontal* distances between the data and the line, as in Figure Three.

Figure Three: Reverse Regression



This brings us to the usual problem in discrimination-related data: The underlying factors that should legitimately affect salaries (like productivity differences) are only imperfectly measured by proxies (like years of education or previous performance evaluation scores.) Thus in a standard regression of salary on productivity proxies, *both* variables are measured with some stochastic error: those on the horizontal axis (because they are represented by imperfect proxies), and salary on the vertical axis (for the usual reasons; stochastic error affects the y-axis variable because not all the variables affecting it have been included in the study). The literature has typically handled this by estimating *both of* the last two models and comparing the results.

In a perfect world we would hope that the slope estimate of the second model would equal the inverse of the slope estimate of the first model, and that statistical tests within each model would yield identical results. But this is generally not so in discrimination cases. We have already indicated that statistical tests are problematic in reverse regressions. Now we can see that the two models' parameter estimates may also not be reciprocals of each other, because the independent variable in the second model is actually not identical to the dependent variable in the first. Since several variables are typically used to proxy productivity in the first model, there is no single candidate to serve as the dependent variable in the second model. So one constructs a proxy-for-the-proxies, using the forecast value of salary from the first regression to serve as the dependent variable in the second regression. Thus we have the rather bizarre reverse-regression circumstance in which forecast salary is being regressed on actual salary. This process introduces a bias into both models' estimation of the slope coefficient, so that neither of the two estimates is normally unbiased; the two estimates form upper- and lower-bounds for the true population parameter. (Maddala, 460-61)

Since both dependent and independent variables in either regression are typically measured with some stochastic error that we can not observe, we could think of the two regressions' slope estimates as approximations that each assume something different about the relative size of the measurement errors. Model One, the direct regression, assumes that the errors in measuring the "vertical" variable are very large relative to those pertaining to the "horizontal" variable. Model Two, the reverse regression, assumes the opposite.⁶ If one knew the relative size of

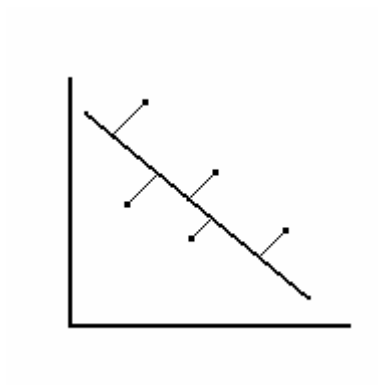
⁶ Kennedy, 141.

the measurement errors, one could obtain a single estimate of the slope by weighting the two models' estimates accordingly.⁷ But normally we have no idea of the relative measurement errors' sizes.

One way to proceed, which would eliminate the acknowledged problems with the reverse regression approach at the expense of making a simplifying assumption, would be to estimate the slope coefficient under the assumption that the variance of the vertical error is equal to the variance of the horizontal error. This is the assumption behind so-called "orthogonal" regression, illustrated in Figure Four. Orthogonal regressions minimize the squared distance of data from the regression line, defining "distance" in the common, right-angle, Pythagorean way.

⁷ Kmenta, 317-21.

Figure Four: Orthogonal Regression



Without knowing the relative measurement-error variances one can't be certain that the assumption behind this approach is called for. This agnosticism probably explains why orthogonal regressions are not more common,⁸ though this is weak logic indeed when it is known that the variables on both axes are being measured with error--one can't know that the assumptions of either of the previous models are called for either. In fact each of them—direct and reverse regression--makes a more extreme assumption about measurement errors than the orthogonal approach does, since they each assume that one set of variables is measured without error. And the statistical tests of significance that are crucial in establishing a discrimination finding are straightforward in orthogonal regression, but are not legitimate in most discrimination-related reverse regression estimations. Furthermore, Monte Carlo simulation has indicated that orthogonal least squares “performs quite well relative to ordinary least squares” estimation.⁹

One might propose that investigators routinely do both a direct and a reverse regression and compare the results; if the two estimates are close to each other, one might conclude that measurement errors must not be a significant problem.¹⁰ One would then invoke orthogonal regression only as a last resort if a comparison of direct and reverse regression results indicates that there is a problem to be solved. But this approach is not as easy as it sounds. Since one regression is log-linear and the other log-log, and since the independent variable in the second model is not identical to the dependent variable in the first, a direct comparison of the two slope estimates is not straightforward. And if the comparison were more straightforward, one would still face the difficulty of deciding how close to each other the estimates must be before measurement errors are judged to be significant.

⁸ Direct and reverse regressions are also more computationally easy than orthogonal regressions—an advantage that no longer matters in a day of inexpensive computing technology, but likely put orthogonal regressions at a distinct disadvantage in the era in which forensic econometric practice was developing its formative habits.

⁹ Boggs et. al.; abstracted in Kennedy, 148. This is a particularly important finding. One is unsure which of these three approaches—direct, reverse, or orthogonal regression—is appropriate because the relative variance of the measurement errors is unknown. In Monte Carlo simulations, artificial databases are constructed for which the investigator *does know* the relative variance of measurement errors; then all three models are applied to these data to see which approach correctly identifies the properties that are known to be present. Thus one can have some confidence that even when “flying blind” with unknown measurement error variances, the orthogonal approach embodies assumptions that seem to be generally defensible.

¹⁰ Kennedy, 148.

Since orthogonal regression has been relatively infrequently used¹¹ and is infrequently discussed in texts,¹² we model its use via an example.

Direct, Reverse and Orthogonal Regression Applied to a Discrimination Case

Because the White and Piette data are confidential and not available from the authors for replication, we instead use a standard employee data set of approximately the same size, the EmployeeData.Sav file provided with SPSS version 10.0. It has approximately the same number of observations (N=474) as the White and Piette data, and includes similar variables. We list the variable names and descriptions below, followed by the variable mean, standard deviation, maximum and minimum:

GENDER	Male/female employee gender, coded male=1 (Mean: 0.54; Standard deviation: 0.50; Maximum: 1.00; Minimum: 0.00)
EDUC	Educational level, years (13.49; 2.88; 21.00; 8.00)
JOB CAT	Employment category (1.41; 0.77; 3.00; 1.00) 1 = Clerical 2 = Custodial 3 = Managerial
SALARY	Current salary, dollars (34,419.57; 17,075.66; 135,000.00; 15,750.00)
SALBEGIN	Beginning salary, dollars (17,016.09; 7,870.64; 79,980.00; 9,000.00)
JOBTIME	Months since hire (81.11; 10.06; 98.00; 63.00)
PREVEXP	Previous work experience, months (95.86; 104.59; 476.00; 0.00)
MINORITY	Minority status; coded 1=minority (0.22; 0.41; 1.00; 0.00)

Table One presents a standard direct regression of the natural log of salary on the usual explanatory variables, and Table Two presents a standard reverse regression analysis.¹³ The direct regression indicates that males and non-

¹¹ Kennedy, 26.

¹² See e.g. Malinvaud, 7-11.

minorities have a salary advantage after all other explanatory variables have been accounted for. These discrimination-related coefficients, like all of the others in the regression, are significant at a 99% level of confidence. The reverse regression indicates that males have lower qualifications at any given salary level than females, but this coefficient fails to achieve a 99% level of confidence. Thus the gender coefficients here show the same pattern as the discrimination coefficients in the White and Piette regressions—signs consistent with discrimination across the regressions, but declining statistical significance in the reverse regression. The reverse regression yields a positive coefficient on minority status, indicating that non-majority workers have higher qualifications at any given salary level, though this coefficient bears a terrible t-test for significance (P-value = .978). As we have indicated earlier, these results would typically be interpreted as conflicting information, though we argue that they constitute evidence that discrimination exists against minority-status employees. In sum, these data give us examples of the sort of reverse-regression paradox that is commonly suggested in the literature—reverse-regression results with signs that are consistent with the direct-regression result, but which lack the direct regression’s statistical significance.

¹³ We have not included beginning salary as an explanatory variable, since if discrimination exists it probably affects starting salary, which would affect current salary indirectly. So inclusion of the starting salary variable would bias our measure of the discrimination coefficients downward. To keep things simple, we have not interacted the job category variables with the gender and minority variables, though that would be an interesting way to explore occupational segregation.

One might split the JOBCAT variable into separate categorical variables for each job category, but we found that such a regression yields the expected increasing positive coefficients for each category in the 1-2-3 sequence with direct- and reverse-regression results that were virtually identical to those we produced using categorical variables for job categories. The reverse regression still produces a statistically-insignificant coefficient on minority status. We judged that consistent use of a single job-category variable throughout the paper was more parsimonious.

Table One: Direct Regression Results

Source	SS	df	MS	Number of obs = 474		
-----+-----				F(6, 467) = 262.87		
Model	57.6151857	6	9.60253096	Prob > F = 0.0000		
Residual	17.0594343	467	.036529838	R-squared = 0.7715		
-----+-----				Adj R-squared = 0.7686		
Total	74.67462	473	.157874461	Root MSE = .19113		

lsalary	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
-----+-----						
educ	.0420381	.0039229	10.72	0.000	.0343293	.0497469
prevexp	-.0003981	.0000925	-4.31	0.000	-.0005798	-.0002164
jobcat	.2704028	.0140235	19.28	0.000	.2428458	.2979597
jobtime	.0025214	.0008772	2.87	0.004	.0007976	.0042451
gender	.1822294	.0202249	9.01	0.000	.1424863	.2219725
minority	-.0741008	.0219154	-3.38	0.001	-.1171657	-.0310359
_cons	9.158715	.0844174	108.49	0.000	8.99283	9.3246

Table Two: Reverse Regression Results

Source	SS	df	MS	Number of obs = 474		
Model	30.6005274	3	10.2001758	F(3, 470)	=	423.02
Residual	11.3331213	470	.024113024	Prob > F	=	0.0000
				R-squared	=	0.7297
				Adj R-squared	=	0.7280
Total	41.9336487	473	.088654648	Root MSE	=	.15528

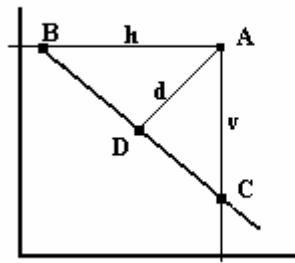
qualif	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lsalary	.6643316	.0217821	30.50	0.000	.6215293	.7071338
gender	-.0392546	.0170984	-2.30	0.022	-.0728534	-.0056559
minority	.0004977	.0179318	0.03	0.978	-.0347387	.0357341
_cons	3.414776	.221865	15.39	0.000	2.978806	3.850746

We have made the case that an orthogonal regression might be a legitimate way to resolve the difficulties raised by these direct and reverse regressions. Such an orthogonal regression could be conducted in several ways. Some regression packages will allow a form of weighted least squares estimation in which the programmer may maintain that the measurement-error variances of the dependent and independent variables are equal. Another approach would involve the use maximum likelihood estimation; the investigator would derive a formula for the orthogonal projections from each data point to the estimated regression hyperplane, then maximize a likelihood function that assumes these projections are normally distributed. A third option, especially useful in cases of non-linear regression or when the relative variances of the dependent and independent variables are known to not be equal, would be use of the Fortran-based public-access software ODRPACK (Boggs *et al.*, 1987).¹⁴

Since ODRPACK is fully documented elsewhere, and since a description of weighted-least-squares estimation would be bound to the particulars of a specific regression software program, we will present results for maximum-likelihood estimations that measure the orthogonal projections from data points to the estimated regression hyperplane. To motivate the discussion, we first consider the simplest case: univariate regression, with a single independent variable. Recall our original orthogonal-regression illustration (Figure Four). Now consider a close-up figure of a single data point (see Figure Five), from which we will derive the appropriate likelihood function for the estimations.

¹⁴ The software is available from <http://www.netlib.org/odrpack/>.

Figure Five: Orthogonal Regression Detail



Let the regression line in Figure Five follow the form

$$\hat{y} = \hat{\alpha} + \hat{\beta}x. \quad (1)$$

We then want to write the distance AD (which we will denote as d) in terms of the coefficients of the regression line and the (x, y) observation. Under the assumption that these errors are normally distributed, we can then minimize the sum of the squared distances between the data and the line via maximum likelihood estimation.

Let Point A represent an observation, (x, y) . Then Point C represents the point

$$(x, \hat{\alpha} + \hat{\beta}x), \quad (2)$$

since the y-coordinate of Point C is the regression's forecast of the dependent variable at this level of the independent variable. Distance AC (which we denote as v for "vertical") is therefore equal to

$$v = y - \hat{y} = y - \hat{\alpha} - \hat{\beta}x, \quad (3)$$

the usual error term in an OLS regression.

Point B's y-coordinate is the same as Point A's. Point B's x-coordinate is, by happy circumstance, the level of x consistent with a regression forecast \hat{y} that is equal to Point A's level of y . Solving Equation 1 for x yields this x-coordinate, so that Point B's coordinates are

$$\left(\frac{y - \hat{\alpha}}{\hat{\beta}}, y\right). \quad (4)$$

Therefore the distance BA (which we denote as h for "horizontal") is equal to the difference between the x coordinates of Points B and A,

$$h = x - \frac{y - \hat{\alpha}}{\hat{\beta}}. \quad (5)$$

Now we can use Equations (4) and (2) to find the distance d . d is an altitude of right triangle ABC. The altitude of any right triangle is equal to the geometric mean of the two non-hypotenuse sides. (Burrill, 360-61). The geometric mean of any two numbers j and k is defined to be the number m for which

$$\frac{j}{m} = \frac{m}{k}, \quad (6)$$

or, solving for the mean,

$$m = \sqrt{j \cdot k} . \quad (7)$$

Therefore altitude d , the orthogonal error term, is equal to

$$d = \sqrt{v \cdot h} . \quad (8)$$

The squared (orthogonal) distance between our data point and the regression line is therefore the product of its vertical and horizontal distances from the line. Stating this in symbols, the squared orthogonal error is equal to

$$d^2 = [(y - \hat{\alpha} - \hat{\beta}x) \cdot (x - \frac{y - \hat{\alpha}}{\hat{\beta}})], \quad (9)$$

the simple product of h and v . Under the assumption that the regression errors follow a normal distribution, it is not difficult to perform a maximum-likelihood estimation of the regression parameters for the univariate regression.

The multivariate case requires matrix algebra, but is not different in kind from the univariate case we have presented.¹⁵ Let the vector

$$A_i = (x_{1,i}, x_{2,i}, \dots, x_{k,i}, [y_i - \hat{\alpha}]) \quad (10)$$

represent any observation i , and let the vector

$$D_i = (x_{1,i}, x_{2,i}, \dots, x_{k,i}, [\hat{y}_i - \hat{\alpha}]) \quad (11)$$

represent the forecast value of \hat{y}_i for the same observation—that is, the regression-line vector. Then the squared orthogonal distance from the point to the line is¹⁶

$$\frac{(A^T A)(D^T D) - (A^T D)^2}{(D^T D)} . \quad (12)$$

We present our Stata maximum-likelihood regression program in Appendix A, and its regression results in Table Three.¹⁷ The coefficients and their levels of significance are extremely close to those of the original direct regression. Thus there is strong evidence for the presence of gender- and race-related salary discrimination.

¹⁵ For a full discussion, see e.g. Strang, Chapter Three.

¹⁶ Strang, 107.

¹⁷ To limit the possibility of rounding errors, our program is written to minimize the compounding of errors carried from one iteration to the next, and requires double precision in all calculations.

Table Three: Maximum-Likelihood Orthogonal Regression Results

Iteration 8: log likelihood = 115.36781

Number of obs = 474

Wald chi2(6) = 1561.55

Log likelihood = 115.36781

Prob > chi2 = 0.0000

```

-----
      lsalary |      Coef.   Std. Err.      z    P>|z|      [95% Conf. Interval]
-----+-----
LS
      educ |      .042046   .0040729    10.32   0.000   .0340632   .0500287
     prevexp |     -.0003981   .0000927    -4.29   0.000   -.0005797  -.0002164
      jobcat |      .2704053   .0140464    19.25   0.000   .242875    .2979357
     jobtime |      .0025246   .0009037     2.79   0.005   .0007535   .0042958
      gender |      .1822158   .0201784     9.03   0.000   .1426668   .2217648
  minority |     -.0740971   .0217702    -3.40   0.001   -.116766   -.0314283
      _cons |      9.158341   .0880861   103.97   0.000   8.985696   9.330987
-----+-----
sigma
      _cons |      .1896964   .0061611    30.79   0.000   .1776209   .2017718
-----

```

Conclusion

The academic literature has raised concerns about the use of reverse regression analysis in employment-discrimination studies, but the greatest legitimate concern—the absence of usable tests of significance in reverse regression—has often gone unnoticed. In the typical case, both direct and reverse regression may be making extreme, mutually-inconsistent and untestable assumptions concerning regression error terms.

We argue that the choice concerning the regression's form—the decision between direct, reverse, and orthogonal regression-- should be a principled choice. If we know the direction of causation--which variable, salary or qualifications, was used to sort applicants--then the variables should go on their proper axes in a single regression. If the direction of causation is not known (both salary and qualifications may have been affected simultaneously by race, with both being joint-normally distributed), or if both variables are measured with error, then both a direct and a reverse regression might be suggested in order to get "bounds" on the race coefficient. If the discrimination-related coefficients in these two regressions do not bear consistent signs, one must seek other sources of information. We have argued, based on Monte Carlo studies and the reasonableness of implicit assumptions concerning error terms, that orthogonal regression analysis is a legitimate method for resolving discrepancies between direct- and reverse-regression results, and may even prove to be the method of preference when there is no dispute between the direct and reverse regressions.

Appendix A: Stata Programs

*** Stata/Windows program to conduct direct and reverse regressions**

version 6.0

* get Employee Data Recode ASCII data set:

insheet using "F:\My Documents\ReverseRegression\edra.dat"

* variable names in row one:

*ID(1-474) GENDER (1=male) BDATE (X/XX/XXXX) EDUC (8-19)

* JOBCAT (1-3 for clerical/custodial/managerial) SALARY (XXXXX)

* SALBEGIN (XXXXX) JOBTIME (63-98=months since hire)

* PREVEXP(0-476 months) MINORITY (0=white)

* BDATE is a string; others are numeric

*conduct direct regression:

generate lsalary=ln(salary)

regress lsalary educ prevexp jobcat jobtime gender minority

* generate the qualification variable for reverse regression:

predict qual

matrix myb=e(b)

scalar bgen=myb[1,5]

scalar bmin=myb[1,6]

generate qualif = qual- (gender*bgen) - (minority*bmin)

*do reverse regression:

regress qualif lsalary gender minority

*** Stata/Windows program to conduct m-1 orthogonal regressions**

version 7.0

set memory 50000

set more off

set matsize 500

```

* get Employee Data Recode ASCII data set:
insheet using "F:\My Documents\ReverseRegression\edra.dat"

* variable names:
  *ID(1-474) GENDER (1=male) BDATE (X/XX/XXXX) EDUC (8-19)
  * JOBCAT (1-3 for clerical/custodial/managerial) SALARY (XXXXX)
  * SALBEGIN (XXXXX) JOBTIME (63-98=months since hire)
  * PREVEEXP(0-476 months) MINORITY (0=white)
* BDATE is a string; others are numeric

generate lsalary=ln(salary)

*write subroutine to evaluate likelihoods, ending with "end"
program define myodrreg
version 7.0
args todo b lmf

*create scalars for coefficient estimates
tempname b1 b2 b3 b4 b5 b6 b7 a myb sig1 lsum bee aa bb cc beebec bee2
tempname abeeabee abee abeed abeed2 abee2 n1 n2 n aa a2 ata sig
tempvar yhat likefn z zsq yhatsh ysh y num denom

*generate coefficient estimates
matrix `myb'=`b'
scalar `b7'=`myb'[1,7]
scalar `sig1'=`myb'[1,8]
scalar `sig'=abs(`sig1')

*calculate predicted values of y
mlevel `yhat'=`b', eq(1)

* generate variable for ln(salary) observations
generate double `y'=$ML_y1

*prepare to calculate squared orthogonal distances z2

```

```

mkmat educ prevexp jobcat jobtime gender /*
    */ minority , matrix(`aa')
gen double `yhatsh'=`yhat'-`b7'
mkmat `yhatsh', matrix(`bb')
matrix `a'=(`aa',`bb')
generate double `ysh'=`y'-`b7'
mkmat `ysh', matrix(`cc')
matrix `bee'=(`aa',`cc')
matrix `beebee'=`bee'*`bee'
matrix `bee2'=vecdiag(`beebee')
matrix `ata'=`a'*`a''
matrix `a2'=vecdiag(`ata')
matrix `abeeabee'=`a'*`bee'
matrix `abee'=vecdiag(`abeeabee')
matrix `abeed'=diag(`abee')
matrix `abeed2'=`abeed'*`abeed'
matrix `abee2'=vecdiag(`abeed2')
matrix `n1'=`bee2'*`a2'
matrix `n2'=vecdiag(`n1')
matrix `n'=`n2'-`abee2'

*Make n and a2 column matrixes before making them variables
tempname ntr a2tr
matrix `ntr'=`n''
matrix `a2tr'=`a2''

*save these as variables
svmat double `ntr', name(`num')
svmat double `a2tr', name(`denom')

*calculate squared orthogonal distances; take square roots
generate double `zsq'=(`num'/`denom')
quietly replace `zsq'=abs(`zsq')
generate double `z'=`zsq'^(.5)

*evaluate likelihood function; see p. 40, ML Est. w Stata
tempvar zz

```

```

generate double `zz'=`z'/`sig'
tempvar prob pone ptwo
gen double `prob'=normden(`zz')
gen double `pone'=log(`prob')
gen double `ptwo'=log(`sig')
generate double `likefn'=`pone' - `ptwo'

*send likelihood back to parent program
mlsum `lnf'=`likefn'

end

*parent m-l program

ml model d0 myodrreg (LS:lsalary = educ prevexp jobcat /*
    */ jobtime gender minority)(sigma:)
ml check
ml search
ml init LS:_cons=9.1944 LS:educ=.0371 LS:prevexp=-.0005 /*
    */LS:jobcat=.25 LS:jobtime=.0025 /*
    */ LS:gender=.2138 LS:minority=-.0926 sigma:_cons=.1801
ml report
ml trace
ml maximize
ml graph
ml display

```

References

- Boggs, P.T. et al. "A Computational Examination of Orthogonal Distance Regression." *Journal of Econometrics* 38 (1988), 169-201.
- Boggs, P.T., R.H. Byrd, J.R. Donaldson and R.B. Schnabel. 1987. "ODRPACK--Software for Weighted Orthogonal Distance Regression." *Technical Report No. CU-CS-360-87*. University of Colorado, Department of Computer Science, Boulder, Colorado.
- Burrill, Gail F., Timothy D. Kanold, Jerry J. Cummings, Lee E. Yunker. *Geometry: Applications and Connections*. Columbus, Ohio: Glencoe/McGraw Hill, 1995.
- Greene, William H. *Econometric Analysis*, Fourth Edition. Upper Saddle River, NJ: Prentice Hall, 2000.
- Kennedy, Peter. *A Guide to Econometrics*, Fourth Edition. Cambridge, Massachusetts: The MIT Press, 1998.
- Kmenta, Jan. *Elements of Econometrics*. New York: Macmillan, 1971
- Maddala, G. S. *Introduction to Econometrics*. New York: Macmillan, 1992.
- Malinvaud, E. *Statistical Methods of Econometrics*. Amsterdam: North Holland, 1966.
- Strang, Gilbert. *Linear Algebra and its Applications*, Second Edition. New York: Academic Press, 1980.
- White, Paul F. and Michael J. Piette, "The Use of 'Reverse Regression' in Employment Discrimination Analysis," *Journal of Forensic Economics*, 1998, 11(2), 127-138.