DOCUMENT RESUME

ED 058 279                                    TM 000 992

AUTHOR        Angoff, William H.; Ford, Susan F.
TITLE         Item-Race Interaction on a Test of Scholastic
              Aptitude.
INSTITUTION   Educational Testing Service, Princeton, N.J.
REPORT NO     RB-71-59
PUB DATE      Oct 71
NOTE          26p.

EDRS PRICE    MF-$0.65 HC-$3.29
DESCRIPTORS   *Aptitude Tests; Caucasian Students; *Item Analysis;
              Mathematics; Negro Students; *Racial Differences;
              *Rural Urban Differences; Standardized Tests; *Test
              Bias; Verbal Ability
IDENTIFIERS   *Preliminary Scholastic Aptitude Test

ABSTRACT
        Several samples of Black and White students were
drawn from the 1970 PSAT administration in Georgia and studied for
item x race interaction on both the verbal and mathematical sections
of the test. When subsamples of candidates were drawn from their
respective racial groups, matched on mathematical for the study of
verbal items and matched on verbal for the study of mathematical
items, there was an observable decrease in the size of the item x
race interaction, suggesting that one factor contributing to that
interaction was simply the difference in performance levels on the
test shown by the two races. Further analyses demonstrated a moderate
item x group interaction for Blacks native to different cities and a
moderate item x group interaction for Blacks native to areas of
different population density. (Author)

# Item-Race Interaction
## on a Test of Scholastic Aptitude

•

William H. Angoff and Susan F. Ford

# ITEM-RACE INTERACTION ON A TEST OF SCHOLASTIC APTITUDE

## Abstract

Several samples of Black and White students were drawn from the 1970 PSAT administration in Georgia and studied for item x race interaction on both the verbal and mathematical sections of the test. When subsamples of candidates were drawn from their respective racial groups, matched on mathematical for the study of verbal items and matched on verbal for the study of mathematical items, there was an observable decrease in the size of the item x race interaction, suggesting that one factor contributing to that interaction was simply the difference in performance levels on the test shown by the two races. Further analyses demonstrated a moderate item x group interaction for Blacks native to different cities and a moderate item x group interaction for Blacks native to areas of different population density.

# ITEM-RACE INTERACTION ON A TEST OF SCHOLASTIC APTITUDE

In their attempt to study the interaction of items and race on the Preliminary Scholastic Aptitude Test (PSAT), Cleary and Hilton (1968) carried out a two-factor analysis of variance with multiple measurements (Items) on one factor (Race), following a pattern previously described by Cardall and Coffman (1964). Within Race three levels of Socioeconomic Status were nested and within each level Individuals were nested. The results of the Cleary-Hilton study indicated that even within SES the item x race interaction was significant. However, plots of item difficulties, Black versus White, failed to reveal particular items that could be diagnosed from knowledge of their content and position on the plot as "biased" toward one race or the other.

One of the difficulties in pursuing studies of this sort--studies to compare the performance of different racial groups on educational and psychological tests--is the stricture against asking questions of individuals regarding their racial (or ethnic or religious) background. Originally these strictures were erected in an attempt to thwart efforts of employers, for example, to exercise negative biases in differential selection and treatment of minority groups. However, as a matter of law the same strictures were of necessity also applied to questionnaires that were not intended for use in selection or differential treatment biased against minority groups but for purposes of research, for example, and even for purposes of differential selection favoring minority groups.

In October of 1970 a pilot form of a Student Descriptive Questionnaire was administered in a tryout administration in Georgia high schools along with the regular administration of the PSAT. Students were asked to identify their racial-ethnic background, but were given the option of omitting their

responses to the question if they chose to. The data emerging from this administration, even though restricted to a small subgroup of the national population of PSAT-takers, afforded the first opportunity to conduct an item x race interaction study on a broader basis than was originally carried out by Cleary and Hilton. Even so, it was anticipated that the results of such a study would be more provocative than informative, and would yield information as much related to methodological problems and considerations as to matters of substantive psychological and social content.

It was planned in the analysis of the Georgia data that a series of analyses would be undertaken to examine not only the general item x race interaction but other interactions as well. The questions to be investigated were these:

1. What is the nature of the item x race (i.e., Black-White) interaction for unselected groups? Is this interaction larger than one might expect from two random unselected groups of Blacks? Of Whites? If there is a substantially larger between-race interaction than within-race interaction, how much larger is it?

2. Are there items that stand out from the others, and do they have noticeably more appropriate content for one race than for the other?

3. How much of the interaction is attributable to racial differences, and how much to literal ability differences, as measured by the PSAT? Assuming that unselected racial groups showed marked differences in mean performance on the PSAT, would the interaction be reduced if the racial groups studied were matched on ability first?

4. Is there item x city interaction? That is, do Blacks from one city perform in the same way on the individual test items as do Blacks from another city?

5. Finally, do urban Blacks perform in the same way on the test items as do nonurban Blacks?

## Samples

In an attempt to answer the foregoing questions, 10 samples were selected from those taking the PSAT in October 1970:

Two mutually exclusive random samples of Atlanta Blacks. (Samples 1 and 2)

Two specially selected (matched with Whites; see below) samples of Atlanta Blacks, both drawn from the examinee population of Atlanta Blacks, exclusive of Sample 1 (essentially Sample 2), but not mutually exclusive. (Samples 3 and 4)

Two mutually exclusive random samples of Atlanta Whites. (Samples 5 and 6)

Two specially selected (matched with Blacks; see below) samples of Atlanta Whites, both drawn from the examinee population of Atlanta Whites, exclusive of Sample 5, but not mutually exclusive. (Samples 7 and 8)

One random sample of Savannah Blacks. (Sample 9)

One random sample of nonurban Blacks, nonurban defined as an area with population less than 25,000. (Sample 10)

Sample 3 and Sample 7 were drawn respectively from the Atlanta Black and Atlanta White populations, exclusive of Samples 1 and 5, but matched frequency for frequency on PSAT-verbal scores and used to study performance on PSAT-mathematical items; Sample 4 and Sample 8 were also drawn respectively from the Atlanta Black and Atlanta White populations, exclusive of Samples 1 and 5, but matched frequency for frequency on PSAT-mathematical scores and used to study performance on the PSAT-verbal items. The method of matching essentially involved placing the two random samples side by side, lined up by score level, and removing cases at random at each score level from the distribution that

had the larger frequency. This procedure was continued until the frequencies in both distributions at each level were equal. Thus, to the extent that PSAT-verbal and mathematical scores are correlated (correlations for the national administrations of the PSAT generally run .65-.70 and were about .65 in the four unselected samples of Blacks and Whites), this type of matching is effective. However, although the matching on verbal does reduce the differences between the samples in means and variances on mathematical--and vice versa--the matching is necessarily incomplete; if, for example, unselected Whites score higher than unselected Blacks on PSAT-mathematical, the matching on PSAT-verbal will reduce that difference but will not remove it.

## Method

Once the samples were defined, each was subjected to an item analysis, in which p-values (proportions of the sample answering correctly) for each item were calculated. These p-values were then transformed, by reference to a table of the normal curve, to normal deviates $(z)$ and from the normal deviates to deltas by the linear transformation, $\Delta = 4z + 13$ . These item deltas were calculated for both PSAT-V and PSAT-M items for all samples except those specially matched. For Sample 3 (Black) and Sample 7 (White) which were matched on verbal scores, deltas were calculated only for the PSAT-mathematical items; for Samples 4 (Black) and 8 (White), which were matched on mathematical scores, deltas were calculated only for the PSAT-verbal items. (Mention should also be made of the fact that the item analysis program used to obtain the deltas contains the restrictions that deltas are calculated only for those items for which $.05 \leq p \leq .95$, and for which the percentage attempting the item equals or exceeds 50. Because of these restrictions not all of the 70 PSAT-verbal items and not all of the 50 PSAT-mathematical items were analyzed for all samples.)

6

In the comparison of any two samples of examinees, a plot was made of
the points, $\Delta_x$ vs. $\Delta_y$ , one point for each of the items under considera-
tion for which deltas were available. The plot of these points normally
appears in the form on an ellipse extending from lower left to upper right;
and if the samples are drawn from the same type of population, the scatter-
plot will be a long, narrow one, often representing a correlation as high as
.98 or .99. When the samples are different in level of performance, the
points will still fall in a narrow ellipse, displaced vertically or horizon-
tally, depending on which group is the abler one. Even when the groups differ
in degree of dispersion, the points will still fall in the same type of
ellipse, but the ellipse will be tilted at an angle more or less steeply than
$45^\circ$, depending on which sample is the more dispersed. However, when the
groups are different in type, the item difficulties will not fall in precisely
the same rank order for the two samples and the correlation represented by the
delta-points will be lower. The items falling at some distance from the plot
may then be regarded as contributing to the item x group interaction. They
are the items that are especially more difficult for one group than for the
other, relative to the other items. Illustrations of these delta plots are
seen in Figures 1-6.

The method that was developed for summarizing the significant features
of each plot involved the determination of the equation for the major axis of
the ellipse represented by the plot and calculating the perpendicular distance
$(D_i)$ from each point to that line. The standard deviation of the distribution
of these distances is a function of the item x group interaction. The correla-
tion coefficient represented by the ellipse represents the degree to which the

items have the same rank order of difficulty in the two groups--also a representation (inversely) of the item x group interaction.

The equation used for the major axis of the ellipse is $Y = AX + B$, where

$$A = \frac{(\sigma_y^2 - \sigma_x^2) \pm \sqrt{(\sigma_y^2 - \sigma_x^2)^2 + 4r_{xy}^2 \sigma_x^2 \sigma_y^2}}{2r_{xy}\sigma_x\sigma_y}$$

and

$$B = M_y - AM_x \quad .$$

(It is recalled that the variables, x and y, are, respectively, the delta values for the two groups under consideration.) The formula for the perpendicular distance, $D_i$, of each point, i, in the plot to the line is given as:

$$D_i = \frac{AX_i - Y_i + B}{\sqrt{A^2 + 1}} \quad .$$

## Results

Means and standard deviations of scores and of item deltas for the 10 samples on both the verbal and mathematical sections of the PSAT are given in Table 1. As may be observed there, the means for the unselected White samples

---------------------------
Insert Table 1 about here
---------------------------

(Samples 5 and 6) on both PSAT-verbal and PSAT-mathematical are about one and one-third standard deviations higher than the means for the unselected Black samples (Samples 1 and 2). The White samples are also more heterogeneous;

their standard deviations are about one and one-fifth times greater than those
of the Black samples. However, after matching on the alternate score, the
means and standard deviations on both verbal and mathematical come closer into
line, the means for the White groups dropping and the means for the Black
groups rising to a point where they are found to be only half a standard devi-
ation apart on verbal (Samples 4 and 8) and about two-thirds of a standard
deviation apart on mathematical (Samples 3 and 7). The standard deviations
themselves also come closer, as expected.

The highest-scoring of the 10 groups are the unselected Whites (Samples 5
and 6); the lowest by far (also the most homogeneous) are the nonurban Blacks,
who scored 20 raw score points lower on verbal and 16 raw score points lower
on mathematical than the unselected Whites. These differences represent about
one and two-thirds times the White standard deviation and two and one-half
times the nonurban Black standard deviation.

Table 2 gives a statistical summary of the results of plotting the item
deltas for selected pairs of groups. The correlations between deltas for

--------------------------

Insert Table 2 about here

--------------------------

paired samples are all higher than .9, indicating a high correspondence in
general between the rank orders of item difficulties for the paired groups.
Correspondingly, the standard deviations of the D-values--the distances from
the points on the ellipse to the major axis of the ellipse--are small, and in
general, smaller when the correlations are higher. Among the highest correla-
tions within the verbal and mathematical columns are those between the two
unselected random samples of Blacks (.978 for Plot I-Verbal and .972 for

Plot I-Mathematical) and between the two unselected random samples of Whites

(.987 for Plot II-Verbal and .987 for Plot II-Mathematical). The plots relat-

ing the delta values for the unselected Black sample and for the unselected

White sample (III-Verbal and III-Mathematical) yielded the lowest correlations

in their respective columns: .929 for verbal and .901 for mathematical. As

the values of $\sigma_D$ also attest, these plots are more dispersed than one would

expect from randomly drawn samples within race. The reader's attention has

already been called to the graphical representations of these six plots

(Plots I, II, and III; Verbal and Mathematical) in Figures 1-6.

---------------------------------

Insert Figures 1-6 about here

---------------------------------

These two sets of plots--one, comparing unselected samples within races

(Plot I and Plot II) and two, comparing unselected samples between races

(Plot III)--describe the extremes of agreement. With these extremes established,

the purpose of the study was to determine whether an attempt to match samples

across races on a related ability would reduce the observed disparity, not only

in the mean delta level but in the rank order of the deltas. Accordingly, a

plot of deltas was made on PSAT-verbal items for Blacks (Sample 4) vs. Whites

(Sample 8) after these groups were selected from their respective populations

and matched frequency for frequency on PSAT-mathematical scores. This plot

(Plot IV-Verbal) yielded a correlation of .959, intermediate between the

extremes of agreement referred to above. Similarly, a plot of deltas was

made on PSAT-mathematical items for Blacks (Sample 3) vs. Whites (Sample 7)

after these groups were selected from their respective populations and matched

frequency for frequency on PSAT-verbal scores. This plot (Plot IV-Mathematical)

yielded a correlation of .923, also intermediate between the extremes of

agreement. Since the item x group interaction is reduced by matching the groups on a related ability, it may be concluded that at least part of the disparity observed in the item x race interaction is due in some way to the fact that the races are so different in their performance on these items. Very likely, a better match, say one based on verbal scores for the preparation of verbal plots and one based on mathematical scores for the preparation of mathematical plots, would be more effective than the method employed here, of matching on verbal for mathematical plots and matching on mathematical for verbal plots. With correlations between verbal and mathematical as low as .65 or so, as were observed in Samples 1, 2, 5 and 6, it would be unreasonable to expect a very close match.

Some slight additional evidence that the matching operation tends to produce a greater agreement in the rank orders of deltas may be seen in Table 2 in the comparison of plots of selected vs. unselected samples within race with plots of unselected vs. unselected samples within race. Plot V-Verbal-- the selected Black sample (4) vs. the unselected Black sample (1)--showed a correlation (.972) slightly lower than that (.978) between the two unselected Black samples (Plot I-Verbal). Although this difference is indeed a small one, the direction of the difference is repeated in the PSAT-M plots for Blacks (.967 in Plot V-Mathematical as compared with .972 in Plot I-Mathematical), and again in the plots for Whites (.968 in Plot VI-Verbal as compared with .987 in Plot II-Verbal and .976 in Plot VI-Mathematical as compared with .987 in Plot II-Mathematical). These comparisons do indeed suggest that item x group interactions are greater (the plot correlations are lower) even within race, when there is a disparity in the general level of performance of the two groups.

11

Detailed examinations of the statistics of the plots and of the
particular items that departed most extremely from the concentrations of
points in the plots were made in a search for clues that would suggest reasons
for their aberrant behavior and more specific hypotheses for the nature of the
item x race interaction. Calculation of the mean D-values separately by item
type (the PSAT-verbal consists of four item types: antonyms, analogies,
sentence completion, and reading comprehension) revealed that the D-values
for the reading comprehension items were more often positive and higher than
the D-values for the other item types, indicating that they were especially
difficult for the Blacks, probably suggestive of the special reading disabil-
ities often characteristic of Black secondary school students. On the other
hand, the antonyms items, which represent a less complex cognitive task,
calling only for the identification of the word most clearly opposite in mean-
ing to the one given in the stem of the item, had relatively high negative D-
values, indicating that they were relatively easier than other item types for
the Black students. It is also interesting to note that both of these observa-
tions, for the reading comprehension items and for the antonyms items, were
more clearly visible when an attempt was made to match the groups on ability.
Further editorial examination of the items that were especially harder for the
Blacks suggested, as one would expect, particular difficulties with vocabulary
and concepts pertaining to unfamiliar places and experiences, and possibly
also to confusions with special meanings and significances characteristic of
the ghetto.

Brief attempts were also made in this study to investigate the hypothesis
that there is a similarity between groups of Blacks that cuts across urban
lines. Accordingly, plots of deltas were made for Blacks in Atlanta versus

12

Blacks in the nearby city of Savannah, also in Georgia. The results of these

plots are somewhat inconsistent. The verbal plot (Plot VII-Verbal) showed a

relatively low correlation (.937), even lower than the correlation between

selected Blacks and selected Whites (.959), and almost as low as the correla-

tion between unselected Blacks and Whites (.929), thus possibly giving some

question to the assumption that there is a homogeneous urban Black culture

cutting across these two cities, at least insofar as their responses to these

verbal items are concerned. It is possible that the low correlation of .937

is due in part to the relative unreliability of the Savannah deltas, which

are based on a sample of only 125 cases. The mathematical plot (Plot VII-

Mathematical) does not yield a correlation quite as low. There the correla-

tion between the Atlanta and Savannah Blacks (.967) runs considerably higher

than either the correlation between selected Blacks and Whites (.923) or

between unselected Blacks and Whites (.901), a finding which may possibly

point to a consistent difference across cities in the nature, as well as the

quality, of mathematics education for Blacks and Whites in the early grades

of school in Georgia.

The final comparison is that made in Plot VIII-Verbal and Plot VIII-

Mathematical, the plots of deltas for urban (i.e., Atlanta) Blacks versus

nonurban Blacks. Here the correlations between deltas are relatively low,

much lower than between random (unselected) samples of Blacks (Plot I-Verbal

and Plot I-Mathematical), but not as low as between unselected samples between

races (Plot III-Verbal and Plot III-Mathematical). In Plot VIII-Verbal the

correlation of .946 is even lower than the correlation of .959 in Plot IV-

Verbal, the plots between matched (selected) samples across races. The

corresponding mathematical plot also has a low correlation (.941) but not as

low as the correlation in the matched cross-racial plot (Plot IV-Mathematical: .923). Again there is evidence that could support the hypothesis that there are differences cutting across urban-nonurban lines, not only in the quality but in the relative emphases given to different aspects of early mathematical education in Georgia for Blacks and Whites.

## Summary

Several samples of Black and White secondary school students in Atlanta who took the PSAT in October 1970 were selected on the basis of their responses to a Student Descriptive Questionnaire administered at that time for tryout purposes. Four of these samples--two Black and two White--were drawn at random from their respective racial populations. Four other samples--two Black and two White--were constructed by matching the frequencies of their verbal scores in order to study the behavior of the mathematical items and by matching the frequencies of their mathematical scores in order to study the behavior of the verbal items. Two additional samples of Blacks were drawn, one from Savannah and the other from nonurban areas of Georgia (population less than 25,000).

Item analyses were conducted, sample by sample, and item deltas were plotted for selected pairs of samples, separately for verbal and mathematical, to determine the degree of item x group interaction, as reflected by the correlation between deltas ($r_{\Delta_x \Delta_y}$) for pairs of groups. (The larger the interaction, the lower the correlation.) The principal findings were that the correlations between deltas were relatively low for unselected samples of Blacks vs. unselected samples of Whites ($r_{\Delta_x \Delta_y}$ was .929 for verbal and .901 for mathematical), noticeably lower than corresponding within-race correlations ($r_{\Delta_x \Delta_y}$ for random unselected samples of Blacks were .978 and .972 for verbal and mathematical, respectively, and .987 for both verbal and mathematical for

random samples of Whites). The between-race correlations were raised sub-
stantially when deltas were correlated for groups that had been matched for
test performance on the alternate section of the PSAT. These correlations
were found to be .959 (verbal) and .923 (mathematical), suggesting that one
possible psychological/educational factor contributing to the item x race
interaction is the difference in levels of performance shown by the two races
on the test. (It is noted that this effect is not a statistical one in the
usual sense, since the measure of interaction used here is the correlation
between item deltas, which is a pure number, independent of level or disper-
sion of performance.)

Further analyses demonstrated a moderate item x group interaction for
Blacks who are native to different cities and for Blacks who are native to
areas of different population density.

Although the present study was carried out only as a pilot, based on
data collected in a tryout administration of the Student Descriptive Question-
naire, the findings have been sufficiently provocative to deserve more detailed
study than was possible here. Some possibilities that come to mind for inves-
tigating the variable that will have an effect on the correlation between plots
(e.g., reducing the item x race interaction) are: (1) that the matching be
done on more directly related variables than the ones used here, perhaps short
forms of the tests under study; (2) that additional variables, like socio-
economic level, be used to match samples for the between-race plots; and
(3) that more detailed studies of interurban and interarea comparisons be made
for Blacks than was possible in the present study, comparisons among cities
as widely separated as Boston, Washington, Detroit, Chicago, St. Louis,
Birmingham, and Los Angeles, for example. With larger samples, drawn from

15

geographically dispersed urban concentrations of ghetto Blacks, there would be more interesting possibilities than were available in these data for studying the homogeneity of the Black culture.

## References

Cardall, C., & Coffman, W. E.  A method for comparing the performance of
different groups on the items in a test.  Research Bulletin 64-61.
Princeton, N. J.:  Educational Testing Service, 1964.

Cleary, A. T., & Hilton, T. L.  An investigation of item bias.  Educational
and Psychological Measurement, 1968, 28, 61-75.

Table 1

Means and Standard Deviations of Scores and Item Deltas on PSAT for Ten Samples

| Sample | Sample No. | No. of Cases | PSAT-Verbal | | | | | PSAT-Mathematical | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Scores | | Deltas | | | Scores | | Deltas | | |
| | | | Mean | S.D. | No. of Items | Mean | S.D. | Mean | S.D. | No. of Items | Mean | S.D. |
| Unselected Blacks | 1 | 300 | 13.51 | 10.64 | 68 | 14.89 | 2.24 | 8.88 | 8.00 | 47 | 15.00 | 2.12 |
| Unselected Blacks | 2 | 340 | 13.18 | 10.10 | 67 | 14.83 | 2.21 | 8.07 | 7.58 | 46 | 15.11 | 2.11 |
| Selected Blacks | 3 | 280 | -----Matched on Verbal----- | | | | | 9.09 | 7.67 | 45 | 14.86 | 2.26 |
| Selected Blacks | 4 | 275 | 14.45 | 10.15 | 68 | 14.71 | 2.33 | -----Matched on Mathematical----- | | | | |
| Unselected Whites | 5 | 300 | 28.32 | 12.17 | 70 | 12.96 | 2.79 | 21.11 | 9.54 | 50 | 12.76 | 2.65 |
| Unselected Whites | 6 | 300 | 28.45 | 12.58 | 70 | 12.91 | 2.72 | 21.57 | 9.65 | 50 | 12.67 | 2.66 |
| Selected Whites | 7 | 280 | -----Matched on Verbal----- | | | | | 14.80 | 8.98 | 49 | 13.98 | 2.43 |
| Selected Whites | 8 | 275 | 18.75 | 9.94 | 70 | 14.26 | 2.52 | -----Matched on Mathematical----- | | | | |
| Savannah Blacks | 9 | 125 | 13.00 | 9.48 | 69 | 15.01 | 2.36 | 7.74 | 8.00 | 49 | 15.12 | 1.96 |
| Nonurban Blacks | 10 | 300 | 8.36 | 8.07 | 70 | 15.52 | 1.90 | 5.18 | 6.35 | 50 | 15.65 | 1.77 |

Table 2

Statistical Summary of Delta Plots

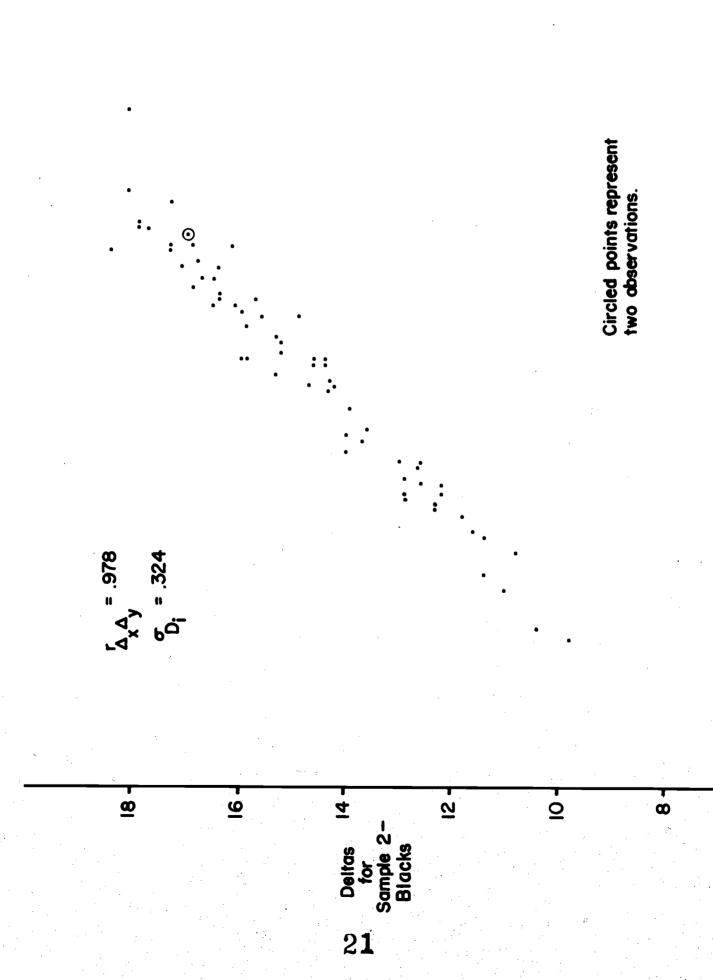| Plot | Sample | Sample No. | | Sample | Sample No. | PSAT-Verbal $r_{\Delta_x\Delta_y}$ | PSAT-Verbal $\sigma_{D_i}$ | PSAT-Mathematical $r_{\Delta_x\Delta_y}$ | PSAT-Mathematical $\sigma_{D_i}$ |
|------|--------|-----------|---|--------|-----------|-----------|-----------|-----------|-----------|
| I | Unselected Blacks | 1 | versus | Unselected Blacks | 2 | .978 | .324 | .972 | .320 |
| II | Unselected Whites | 5 | versus | Unselected Whites | 6 | .987 | .317 | .987 | .308 |
| III | Unselected Blacks | 1 | versus | Unselected Whites | 5 | .929 | .641 | .901 | .706 |
| IV | Selected Blacks | 4 | versus | Selected Whites | 8 | .959 | .490 | .923 | .585 |
| IV | Selected Blacks | 3 | versus | Selected Whites | 7 | | --- | | |
| V | Selected Blacks | 4 | versus | Unselected Blacks | 1 | .972 | .376 | .967 | .358 |
| V | Selected Blacks | 3 | versus | Unselected Blacks | 1 | | --- | | |
| VI | Selected Whites | 8 | versus | Unselected Whites | 5 | .968 | .474 | .976 | .375 |
| VI | Selected Whites | 7 | versus | Unselected Whites | 5 | | --- | | |
| VII | Atlanta Blacks | 1 | versus | Savannah Blacks | 9 | .937 | .558 | .967 | .338 |
| VIII | Atlanta Blacks | 1 | versus | Nonurban Blacks | 10 | .946 | .464 | .941 | .443 |

## Figure Captions

Figure 1.  Plot I-Verbal

Figure 2.  Plot I-Mathematical

Figure 3.  Plot II-Verbal

Figure 4.  Plot II-Mathematical

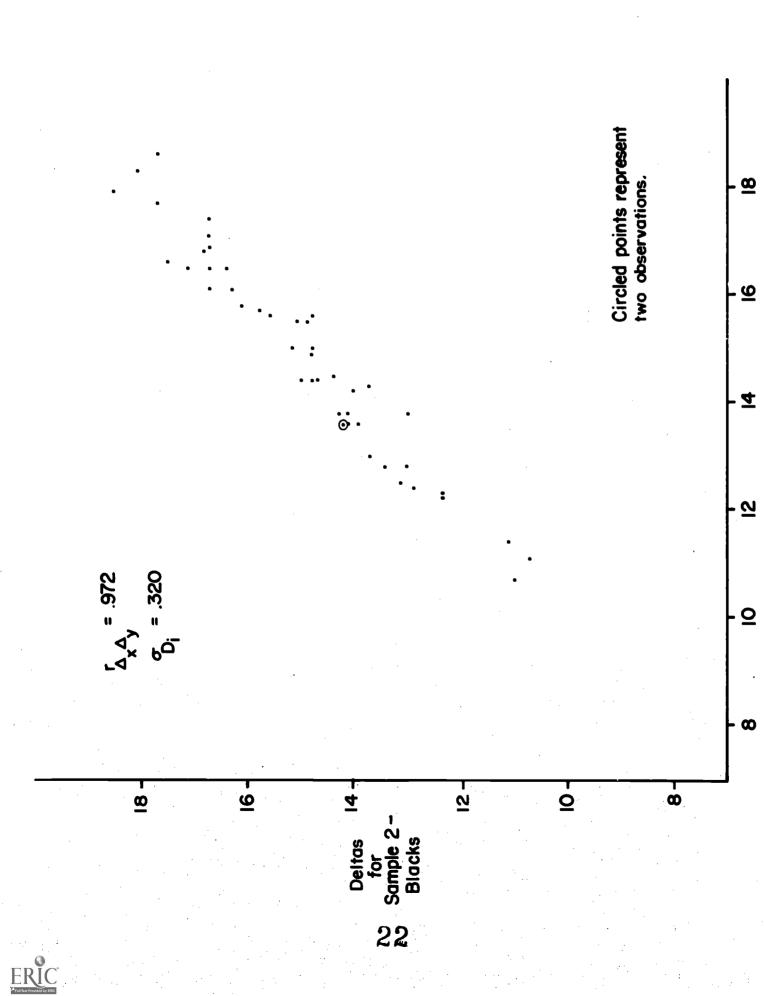Figure 5.  Plot III-Verbal

Figure 6.  Plot III-Mathematical
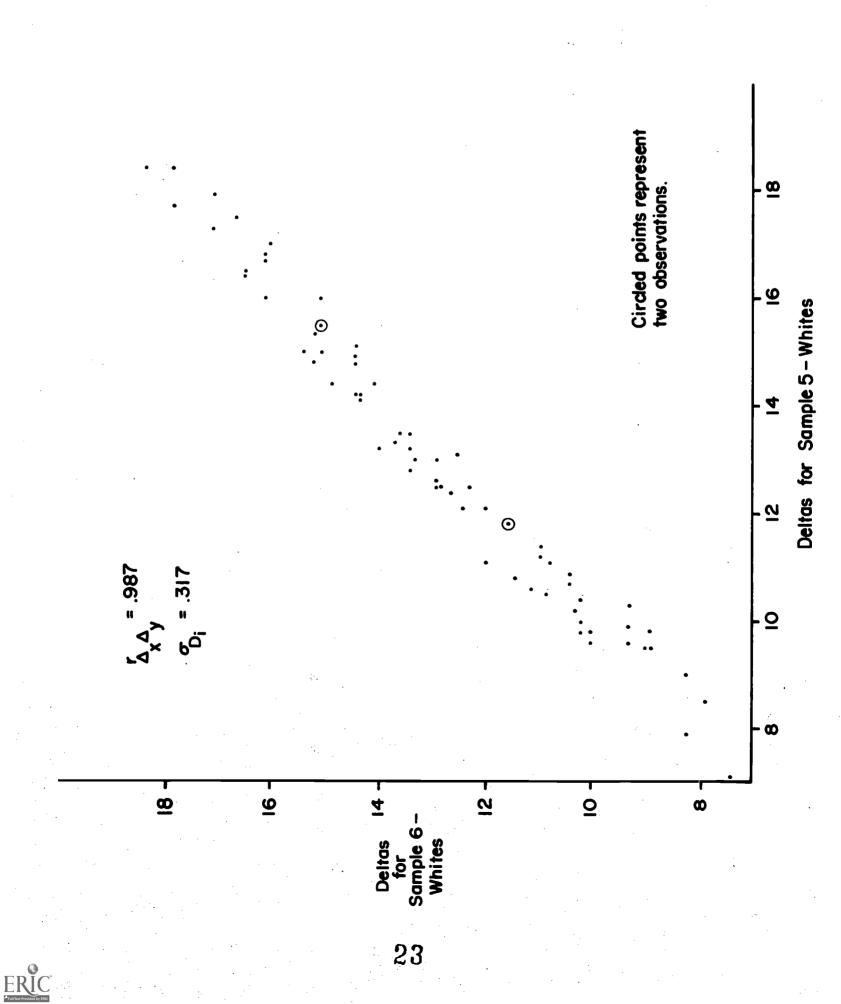
$r_{\Delta_x \Delta_y} = .978$

$\sigma_{D_i} = .324$

Deltas
for
Sample 2 –
Blacks

Deltas for Sample I – Blacks

Circled points represent
two observations.

$r_{\Delta_x \Delta_y} = .972$

$\sigma_{D_i} = .320$

Deltas for Sample 2 – Blacks

18

16

14

12

10

8

Circled points represent two observations.

Deltas for Sample I – Blacks

8   10   12   14   16   18

22

$r_{\Delta_x \Delta_y} = .987$

$\sigma_{D_i} = .317$

Deltas for Sample 6 – Whites

Deltas for Sample 5 – Whites

Circled points represent two observations.

Circled points represent
two observations.

Deltas for Sample 5 – Whites

$r_{\Delta_x \Delta_y} = .987$

$\sigma_{D_i} = .308$

Deltas
for
Sample 6 –
Whites

$r_{\Delta_x \Delta_y} = .929$

$\sigma_{D_i} = .641$

Deltas for Sample 5 – Whites

Deltas for Sample I – Blacks

Circled points represent two observations.

$r_{\Delta_x \Delta_y} = .901$

$\sigma_{D_i} = .706$

Deltas for Sample 1 – Blacks

Deltas for Sample 5 – Whites