

**Longitudinal Analysis of Mathematical Achievement of National Longitudinal Survey of
Youth 1997 (NLSY 1997) Participants**

Jonathan King

UCLA Fielding School of Public Health

BIOSTAT M236: Longitudinal Data

Dr. Robert Weiss

June 6, 2022

The National Longitudinal Survey of Youth 1997 (NLSY 1997) is an ongoing longitudinal survey sponsored by the United States Bureau of Labor Statistics (BLS) of the U.S. Department of Labor that follows 8,984 American men and women born between 1980 and 1984. Participants were interviewed annually from 1997 to 2011 and biennially since, for a total of 19 rounds to date (BLS, 2021). This analysis determines if mathematical achievement, measured by scores on the Peabody Individual Achievement Test (PIAT) math subtest, a standardized test administered during the first six years of the study from 1997-2002, with 1997 defined as baseline (BLS, n.d.), follows a definitive time trend and if this time trend is affected by certain demographic characteristics of the survey participants through exploratory data analysis, fixed effects analysis, and covariance model selection. The potential covariates that are investigated include years since baseline (the time variable), age at baseline, sex, whether English is the primary reading language, annual gross household income at baseline category, baseline region, race/ethnicity, and whether parents were born in the United States. Since missing data is an issue due to implementation of age restrictions after the first round of testing (BLS, n.d.), only individuals who completed at least the first two assessments and had no missing data in the covariates will be included in the analysis, limiting the sample size to 930 subjects with a total of 4,628 observations.

Exploratory data analysis for this study includes profile plots and empirical summary plots. Since constructing profile plots for the whole dataset is not useful due to the large sample size, two randomly selected subsets of 50 subjects are plotted (Figure 1). These plots show that few participants took the final PIAT test 5 years after baseline and that there are quite a few bivariate and univariate outliers. These subsets of the data also clearly have a random slope and random intercept, so it would be reasonable to assume that this is the case for the full sample. Empirical

within-subjects residual profile plots of the same subsets show that PIAT scores do slightly decrease from baseline to year 2 but increase from years 2-4 (Figure 2). The trend at year 5 cannot be determined due to high rate of dropout. Overall, changes in PIAT scores appear to be non-linear, so higher-order polynomial time trends should be considered. To help determine which covariates could have a significant effect on the PIAT score time trend while accounting for unequal group sizes, empirical summary plots were created for each covariate (Figure 3). Mean PIAT scores do not appear to differ between sexes or depend on whether parents were born in the US or age at baseline (Figure 3a-c). However, PIAT score means are clearly higher for individuals who reported higher annual gross household income at baseline (Figure 3d). Comparing by race/ethnicity also shows significant differences in mean PIAT scores, with mixed race and non-Black/non-Hispanic individuals clearly scoring higher than both Blacks and Hispanics and Blacks scoring lower than Hispanics (Figure 3e). Baseline region is another covariate that shows differences between groups, with individuals from the southern United States scoring lower on average than those from the northeastern or north central United States (Figure 3f). Lastly, average PIAT scores are much higher with lower standard errors for participants who primarily read in English compared to those who do not (Figure 3g).

Fixed effects analysis consists of two parts: population time trend determination and fixed effects models selection. Time trend selection using the random intercept (RI) covariance model reveals that the best-fitting population time trend is a cubic time trend because the highest degree term in each model is significant ($p \leq .017$) until the quartic model ($p = 0.74$; Table 1). To determine the correct fixed effects model, forward selection was conducted using the RI covariance model with the initial base model including the cubic time trend, sex, and intercept because sex is an essential demographic variable for this analysis. The selected model includes

the linear fixed effects of race/ethnicity, income category, region, primary reading language, as well as interactions between the linear time term and both income category and race/ethnicity (Table 2; all $p \leq .029$).

Covariance model selection using restricted maximum likelihood (REML) leads to two possible models that best fit the selected fixed effects: Unstructured, which has the lowest AIC, and random intercept and slope (RIAS), which has the lowest BIC (Table 3). Ultimately, the unstructured covariance model is selected because the likelihood-ratio test against the RIAS model is highly significant ($p < .0001$). The fitted time trend of PIAT scores for the reference group, defined as non-Black and non-Hispanic females who primarily read in English, are from households that had an annual gross income of less than \$10,000 at baseline, and lived in the western United States at baseline, is $-0.22t^3 + 1.89t^2 - 2.84t + 94.80$, where t is time in years since baseline (Table 4). Estimated mean PIAT scores does not significantly depend on sex, controlling for all other variables (Table 4). Subjects who live in the north central ($p = .009$), northeastern ($p = .0003$), or southern United States ($p = .044$), or report a gross household income of \$50,000 or more during the baseline year have higher estimated mean PIAT scores ($p < .0001$), controlling for all other variables (Table 4). However, subjects who identify as Black ($p < .0001$) or Hispanic ($p = .04$), or do not primarily read in English ($p = .0001$), have lower estimated mean PIAT scores, holding all other variables constant. The estimated linear time coefficients for subjects who report a gross household income of \$50,000 or more at baseline are higher ($p \leq .0015$), while the estimated linear time coefficients for Black ($p < .0001$) and Hispanic subjects ($p = .04$) are lower, adjusting for all other variables. The full estimated covariance matrix shows that model fit for each subject generally worsen over time, except at the year 5 timepoint and that residual errors within the same subject all have the same sign (Table 5).

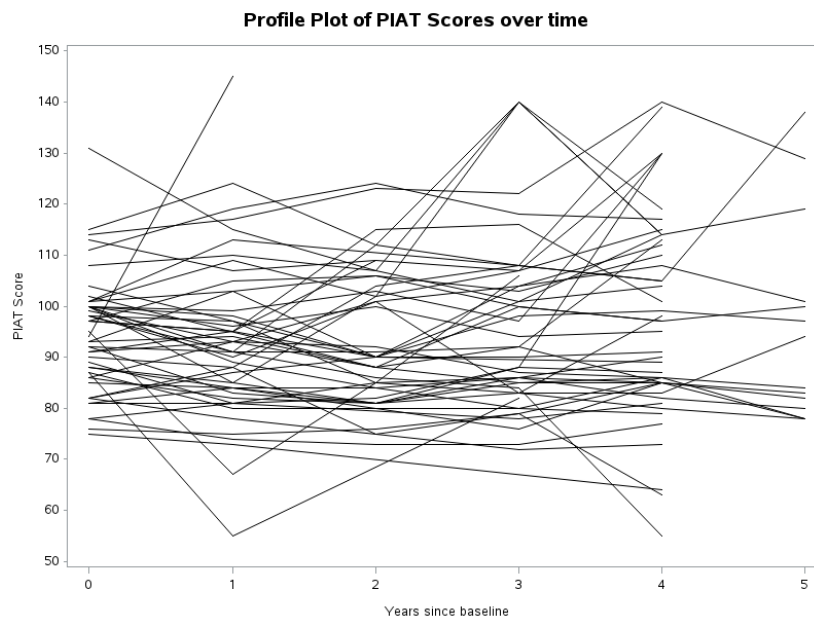
References

U.S. Bureau of Labor Statistics. (2021, December 1). *NLSY97 Data Overview*. U.S. Bureau of Labor Statistics. Retrieved May 29, 2022, from <https://www.bls.gov/nls/nlsy97.htm>.

U.S. Bureau of Labor Statistics. (n.d.). *PIAT math test*. National Longitudinal Surveys.

Appendix

a)



b)

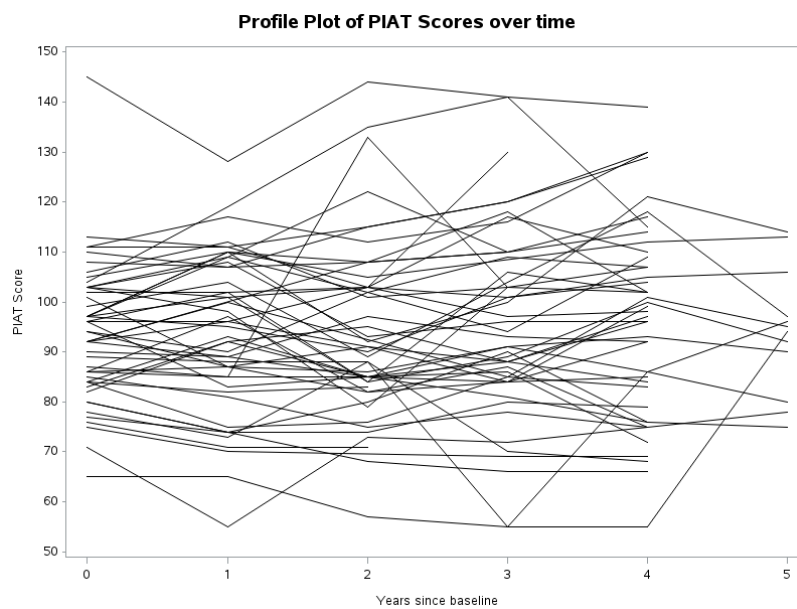
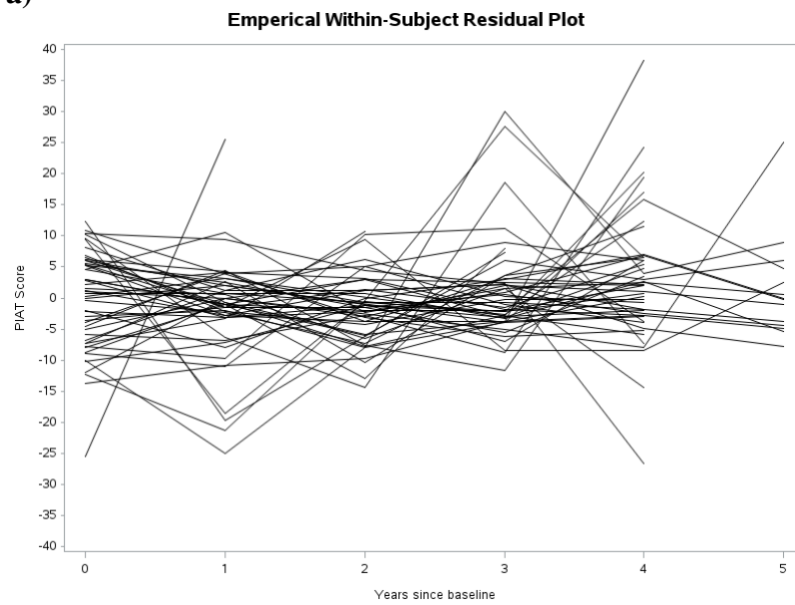


Figure 1. Profile plot of PIAT scores against time for two randomly selected subsets of 50 subjects (a-b). Consecutive observations within the same subject are connected by line segments.

a)



b)

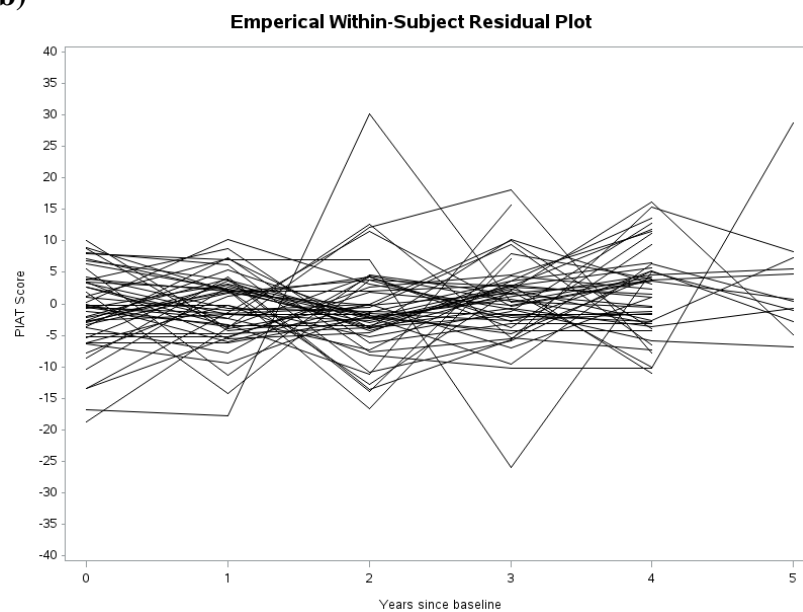
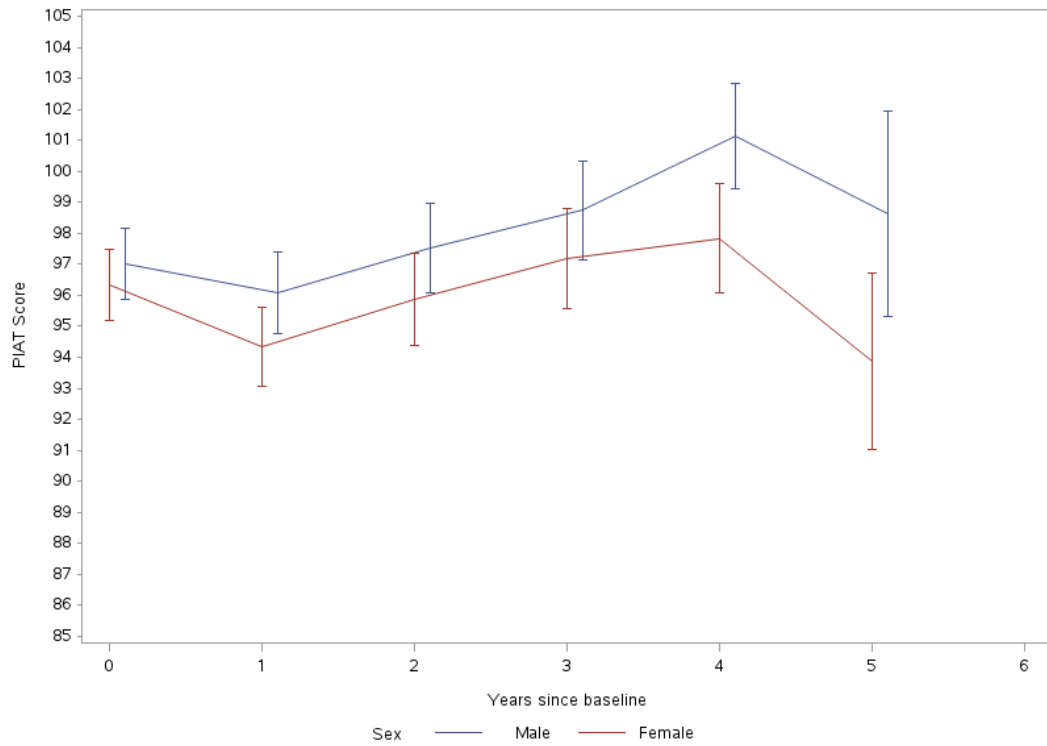


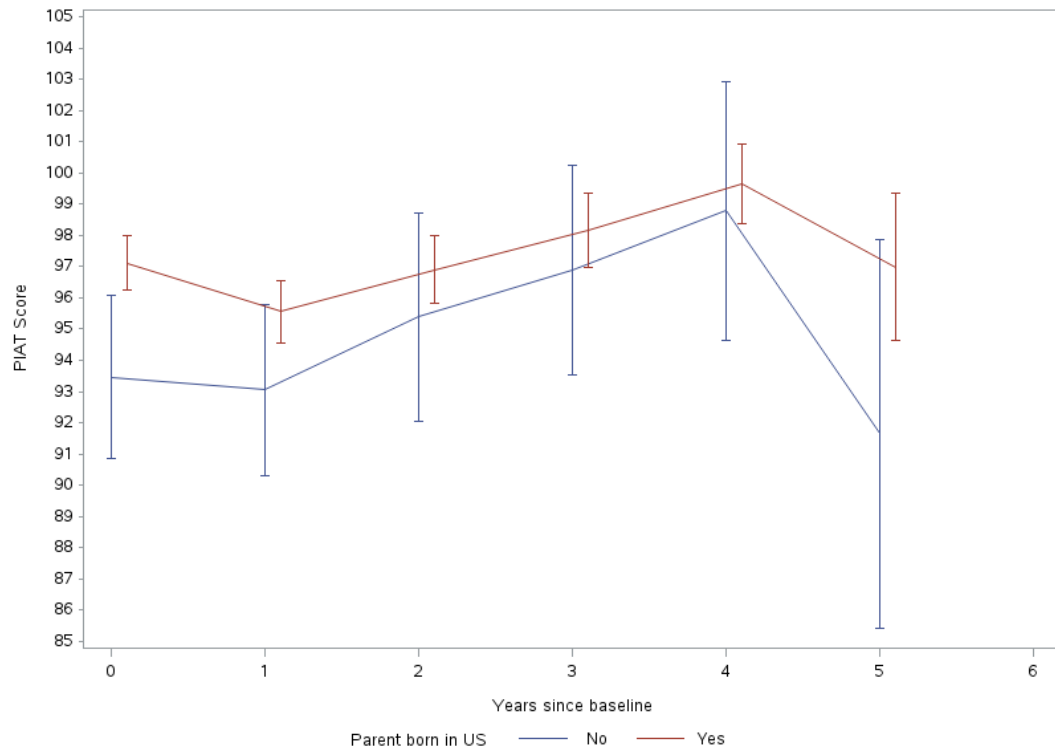
Figure 2. Profile plot of empirical within-subject residuals of PIAT scores against time for two randomly selected subsets of 50 subjects (a-b). Residuals are calculated as the difference between each observation and the subject-specific mean. Consecutive observations within the same subject are connected by line segments.

Figure 3. Empirical summary plots of PIAT score against time by **a)** sex, **b)** parent birthplace, **c)** age at baseline, **d)** annual household income at baseline, **e)** baseline region, **f)** race/ethnicity, and **g)** primary reading language. Mean PIAT score at each timepoint are plotted and connected by line segments. Error bars at each point represent two standard errors in both directions from the mean. Mean profiles representing different groups are offset to avoid overplotting. **a)** The mean profile for males is blue; female mean profile is red. **b)** Subjects with parents born outside of the United States have a blue mean profile; subjects with parents born in the United States have a red mean profile **c)** Subjects who were 12 years-old at baseline have a blue mean profile; subjects who were 13 years-old at baseline have a red mean profile. **d)** Subjects who reported a baseline annual household income less than \$10,000, from \$10,000 to \$49,999, from \$50,000 to \$99,999 and more than \$100,000 have mean profiles that are blue, red, green, and brown, respectively. **e)** Subjects who reside in the northeastern, north central, south, and western United States have mean profiles that are blue, red, green, and brown, respectively. **f)** Subjects who identify as Black, Hispanic, mixed race, and non-Black/non-Hispanic have mean profiles that are blue, red, green, and brown, respectively. **g)** The mean profile for subjects whose primary reading language is English is blue; mean profiles for subjects whose primary reading language is not English is red.

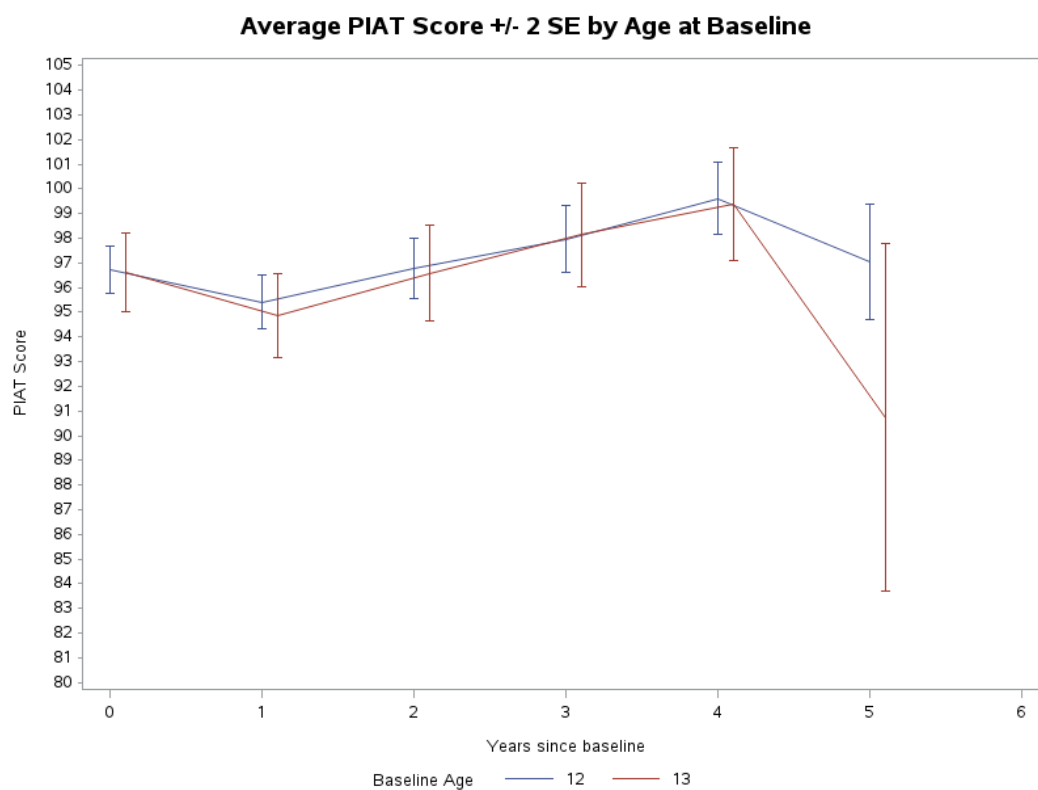
a)

Average PIAT Score \pm 2 SE by Sex

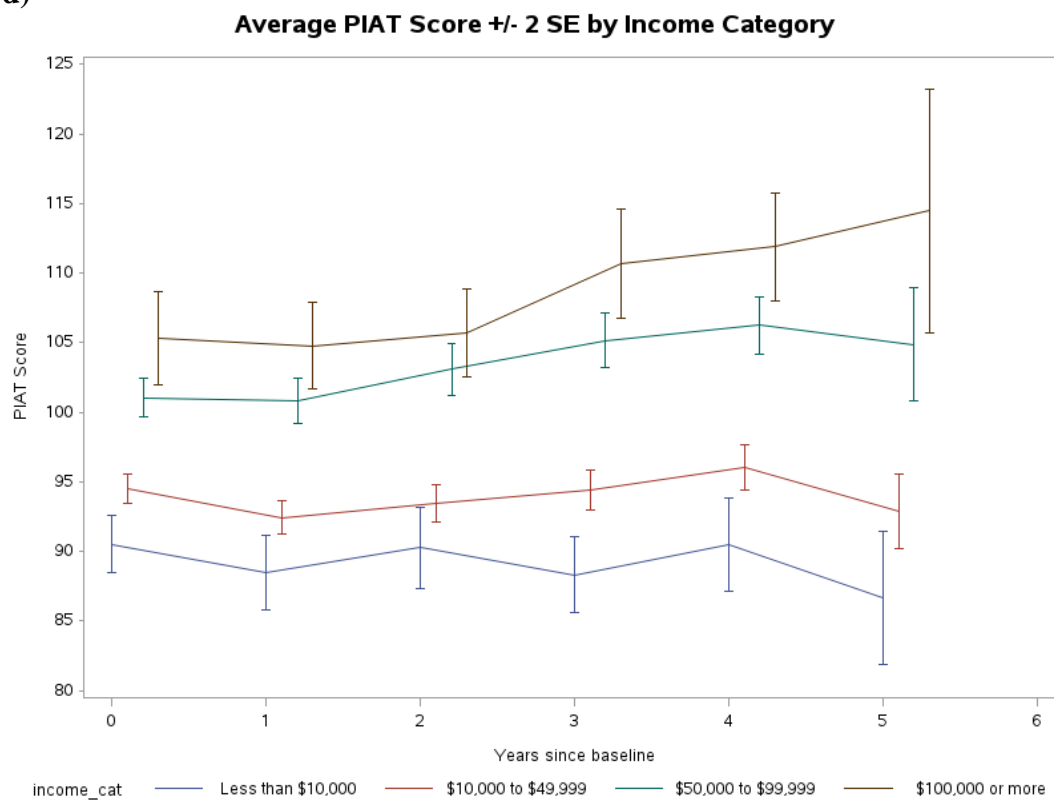
b)

Average PIAT Score \pm 2 SE by Parent Birthplace

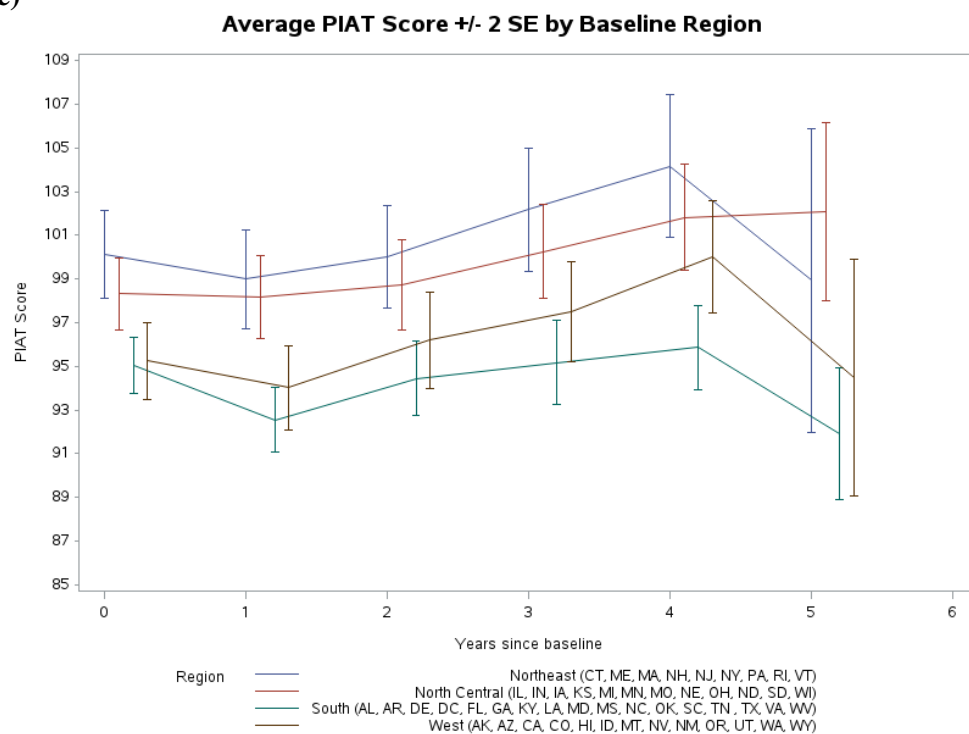
c)



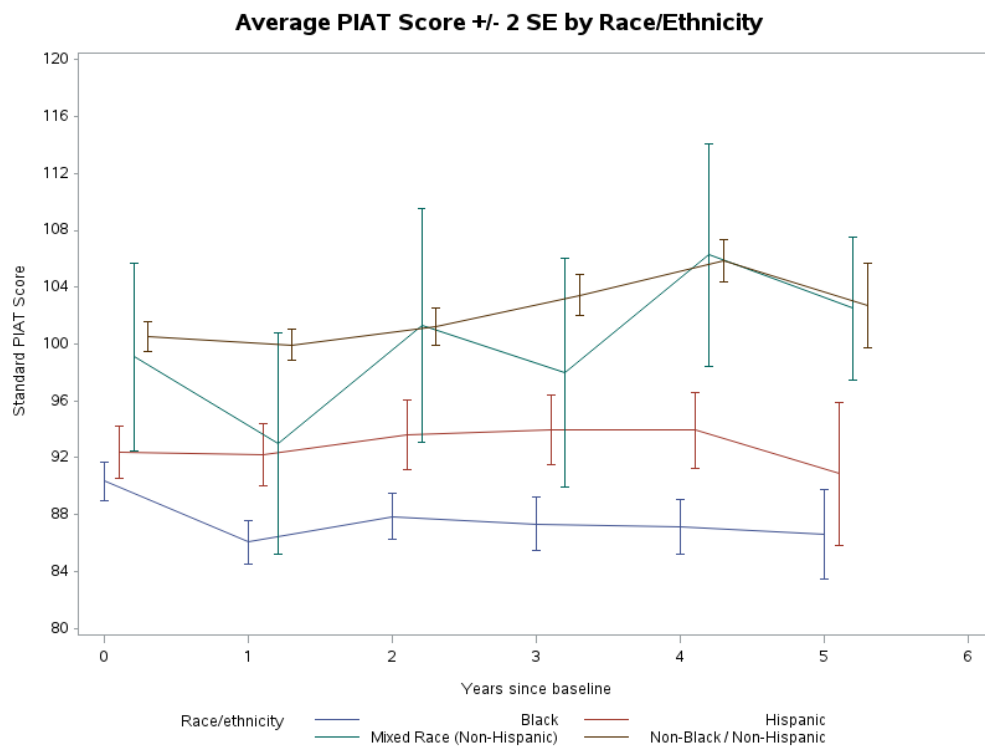
d)



e)



f)



g)

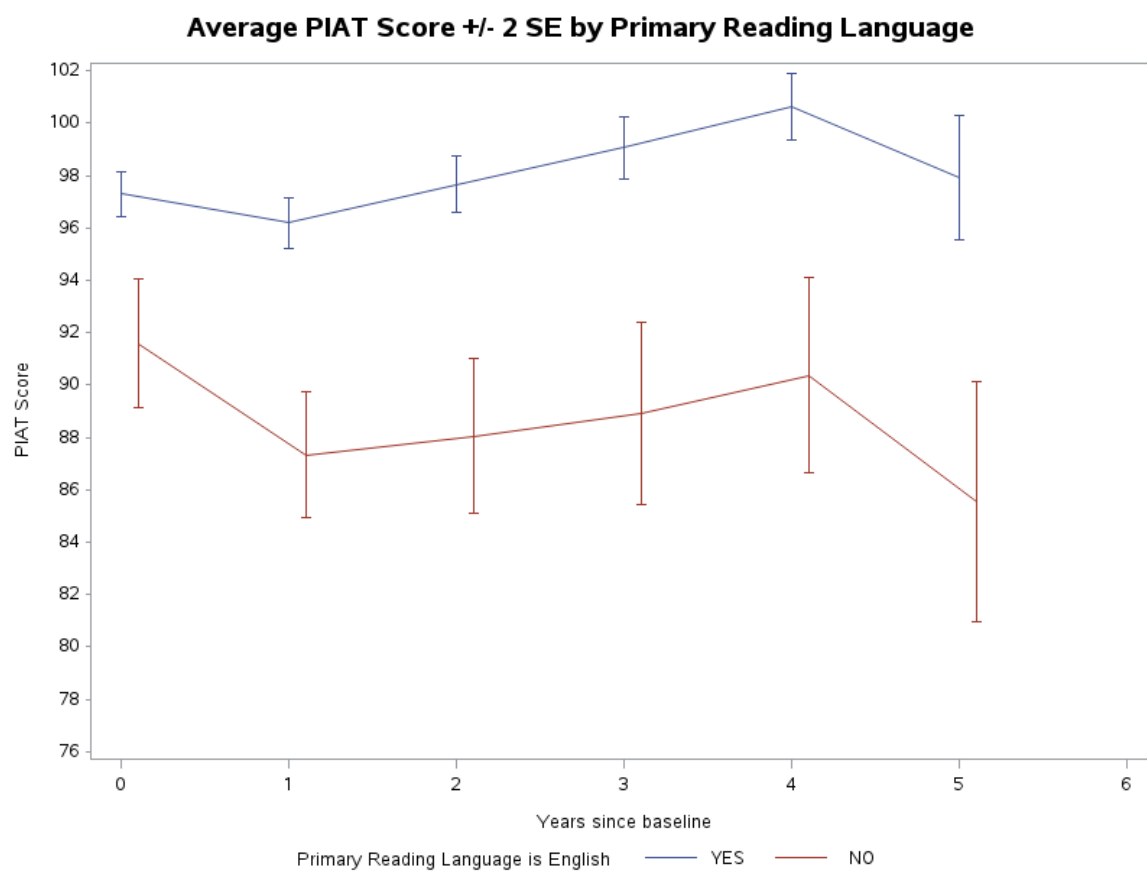


Table 1. Model output for the linear, quadratic, cubic, and quartic models fit using the random intercept (RI) covariance model and the restricted maximum likelihood (REML) estimation.

Each row contains a parameter estimate, standard error, and corresponding *t*-statistic and *p*-value from testing the null hypothesis that the parameter estimate equals 0.

Model	Parameter	Estimate	Standard Error	<i>t</i> -statistic	<i>p</i> -value
Linear	Intercept	95.77	0.47	203.86	<.0001
	Time	0.62	0.08	7.55	<.0001
Quadratic	Intercept	96.11	0.49	195.82	<.0001
	Time	0.009	0.27	0.03	.97
	Time*time	0.14	0.06	2.38	.017
Cubic	Intercept	96.66	0.50	193.85	<.0001
	Time	-2.83	0.55	-5.15	<.0001
	Time*time	1.82	0.29	6.37	<.0001
	Time*time*time	-0.24	0.04	-5.91	<.0001
Quartic	Intercept	83.88	0.50	192.98	<.0001
	Time	-3.12	1.02	-3.05	.002
	Time*time	2.14	1.00	2.14	.032
	Time*time*time	-0.35	0.32	-1.08	.28
	Time*time*time*time	0.01	0.03	0.34	.74

Table 2. *F*-tests for the selected fixed effects model fit using the random intercept (RI) covariance model and the maximum likelihood (ML) estimation. Each row contains the fixed effect name, *F*-statistic, and corresponding *p*-value.

Effect	<i>F</i> -statistic	<i>p</i> -value
Time	17.72	<.0001
Time*time	36.59	<.0001
Time*time*time	31.50	<.0001
Sex	4.79	.029
Region	3.45	.016
Race/Ethnicity	19.29	<.0001
Household income category	3.45	<.0001
Primary reading language	18.04	<.0001
time* household income category	8.51	<.0001
time*race/ethnicity	14.57	<.0001

Table 3. Results of fitting different covariance models using restricted maximum likelihood (REML) estimation. Each row includes the covariance model name, number of covariance parameters, -2 log restricted maximum likelihood, Akaike information criteria (AIC) , Bayesian information criteria (BIC), and test statistic and p -value from the likelihood ratio test versus the unstructured model. AIC and BIC are in “smaller is better” form. Abbreviations: AR(1) = autoregressive; ARH(1) = heterogenous autoregressive; ARMA (1,1) = autoregressive moving average; ANTE(1) = antedependence; CSH = heterogenous compound symmetry; FA(q) = factor analytic; IND = independence; LRT = likelihood ratio test; RI = random intercept; RIAS = random intercept and slope; RIASAQ = random intercept, slope, and quadratic; UN = unstructured.

Covariance Model	# of Covariance Parameters	REML: -2 Restricted Log Likelihood	AIC (smaller is better)	BIC (smaller is better)	LRT vs. UN	
					Test Statistic	p -value
RI	2	34618	34622	34632	227	<.0001
RIAS	4	34442	34450	34469	51	<.0001
RIASAQ	7	34426	34440	34473	35	.0017
CSH	7	34476	34490	34524	85	<.0001
IND	1	37075	37077	37082	2684	<.0001
ARH(1)	7	34851	34865	34899	460	<.0001
AR(1)	2	34912	34916	34926	521	<.0001
ARMA(1,1)	3	34552	34558	34573	161	<.0001
ANTE(1)	11	34822	34844	34897	431	<.0001
UN	21	34391	34433	34534	.	.

Table 3 (continued).

FA(q)	12	34433	34457	34515	42	<.0001
-------	----	-------	-------	-------	----	--------

Table 4. Fixed effect parameter estimates for the selected fixed effects model (cubic time trend, intercept, sex, race/ethnicity, income category, region, primary reading language, and interactions between the linear time term and both income category and race/ethnicity) using restricted maximum likelihood (REML) estimation and the unstructured (UN) covariance model. Each row contains the parameter estimate, standard error of the estimate, and the corresponding *t*-statistic and *p*-value from a *t*-test testing that the parameter estimates equals to 0. Reference region of residence at baseline is western United States; reference sex is female; reference race/ethnicity is non-Black/non-Hispanic; reference annual gross income category at baseline is less than \$10,000; reference category for English is primary reading language is yes.

Effect	Estimate	Standard Error	t-statistic	p-value
Intercept	98.86	1.90	62.05	<.0001
Time	-1.41	0.62	-2.26	<.0001
time*time	1.69	0.28	6.00	<.0001
time*time*time	-0.22	0.04	-5.48	<.0001
Sex (male)	1.01	0.68	1.48	.14
Region (north central)	2.78	1.06	2.63	.009
Region (northeast)	4.18	1.15	3.64	.0003
Region (south)	1.97	0.99	2.02	.044
Race/ethnicity (Black)	-8.18	0.94	-8.72	<.0001
Race/ethnicity (Hispanic)	-4.20	1.03	-4.08	<.0001
Race/ethnicity (mixed)	-1.41	3.20	-0.44	.66
Income category (\$10,000 to \$49,999)	1.71	1.13	1.51	.13
Income category (\$100,000 or more)	9.17	1.71	5.37	<.0001
Income category (\$50,000 to \$99,999)	6.38	1.27	5.01	<.0001

Table 4 (continued).

English primary reading language (no)	-4.30	1.12	-3.83	.0001
time*income category (\$10,000 to \$49,999)	0.25	0.29	0.87	.39
Time*income category (\$100,000 or more)	1.44	0.45	3.18	.0015
time*income category (\$50,000 to \$99,999)	1.05	0.33	3.21	.0014
time*race/ethnicity (Black)	-1.31	0.24	-5.55	<.0001
time*race/ethnicity (Hispanic)	-0.54	0.26	-2.06	.04
time*race/ethnicity (mixed)	0.87	0.85	1.02	.31

Table 5. Estimated full covariance matrix from the unstructured (UN) covariance model using the restricted maximum likelihood (REML) estimation. Each time in row j has an estimated covariance with the time in column k , with both j and k representing indices of the time vector [0, 1, 2, 3, 4, 5], starting from index 1. This is an estimated variance when j and k are equal. Each entry is the estimated residual covariance between observations or residual variance of observations within the same subject at the corresponding pairs of times. All estimated variances and covariances are highly significant ($p < .0001$).

	Baseline	Year1	Year2	Year3	Year4	Year5
Baseline	126.37	81.73	88.32	97.72	100.30	103.57
Year1	81.73	150.85	110.26	108.32	114.46	109.91
Year2	88.32	110.26	182.93	130.42	135.04	122.15
Year3	97.72	108.32	130.42	202.77	158.70	145.20
Year4	100.30	114.46	135.04	158.70	233.49	159.65
Year5	103.57	109.91	122.15	145.20	159.65	216.79