# Added-variable plots for panel-data estimation

John Luke Gallup
Portland State University
Portland, OR
jlgallup@pdx.edu

**Abstract.**   In this article, I extend the theory of added-variable plots to three panel-data estimation methods: fixed effects, between effects, and random effects. An added-variable plot is an effective way to show the correlation between an independent variable and a dependent variable conditional on other independent variables. In a multivariate context, a simple scatterplot showing $x$ versus $y$ is not adequate to show the relationship of $x$ with $y$, because it ignores the impact of the other covariates. Added-variable plots are also useful for spotting influential outliers in the data that affect the estimated regression parameters. Stata can display added-variable plots with the command `avplot`, but it can be used only after `regress`. My new command, `xtavplot`, is a postestimation command that creates added-variable plots after `xtreg` estimates. Unlike `avplot`, `xtavplot` can display a confidence interval around the fitted regression line.

**Keywords:** gr0082, xtavplot, xtavplots, added-variable plot, panel data, postestimation diagnostics, xtreg

## 1 Introduction

An added-variable plot displays a scatterplot of a transformation of an independent variable (say, $x_1$) and the dependent variable ($y$) that nets out the influence of all the other independent variables. The fitted regression line between these transformed variables has the same slope as the coefficient on $x_1$ in the full regression model, which includes all the independent variables.

An added-variable plot is a visually compelling method for showing a partial correlation between $x_1$ and $y$. A confidence interval shows how precisely the sample data fit that correlation. An added-variable plot is the multivariate analogue of using a simple scatterplot with a regression fit in a univariate context.

The main purpose of the panel-data estimation methods in `xtreg` is to control for individual effects. If it is important to control for them in regressions, it is also important to control for them in graphs of the relationship of a covariate with the dependent variable. `xtavplot` controls for the influence of individual effects as well as other covariates on the partial correlation of $x_1$ and $y$.

Outliers in a simple scatterplot of $x_1$ versus $y$ may no longer be outliers when other covariates are included in the model. An added-variable plot is a handy visual diagnostic for spotting influential outliers after conditioning on the other covariates in the model.

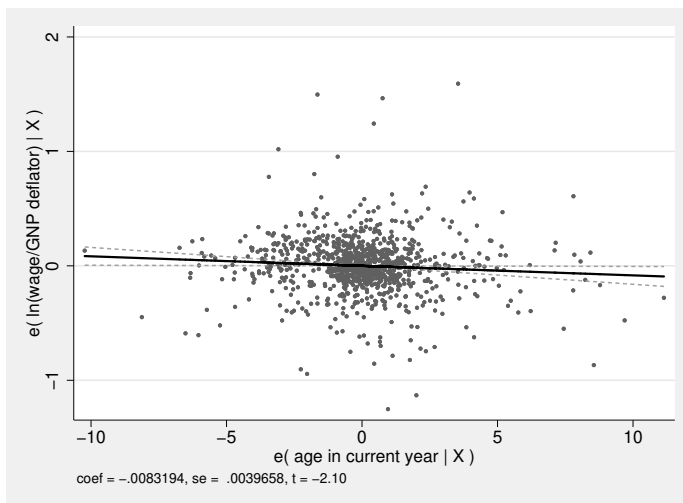# 2 Why do we need added-variable plots, and where do they come from?

The purpose of multivariate regression is to assess the influence of each independent variable on the dependent variable while accounting for the influence of all the other independent variables. The regression coefficient quantifies the partial correlation of an independent variable ($x_1$) on the dependent variable ($y$), controlling for the other independent variables ($\mathbf{x}$). A simple scatterplot is an effective visual presentation of the unconditional correlation of $x_1$ with $y$, but an added-variable plot is needed to display the partial correlation of $x_1$ with $y$ conditional on other $\mathbf{x}$ variables. The partial correlation typically has a different magnitude and may even have a different sign than the unconditional correlation.

For example, there is a positive correlation between the log of wages and worker age in the National Longitudinal Study of Young Women Stata dataset. This is clear to the eye from a scatterplot of the data with a regression line:



However, in a fixed-effects regression that includes age as well as a quadratic in job tenure and total years of labor market experience, age has a negative partial correlation with log wages in this sample. We can graphically display this relationship—the partial correlation of age with log wages controlling for the other independent variables—with `xtavplot`:

```
. xtreg ln_w age tenure c.tenure#c.tenure ttl_exp, fe
  (output omitted)
. xtavplot age
```



coef = −.0083194, se = .0039658, t = −2.10

The added-variable plot provides a graphical representation of the relationship be-tween age and wages when other regressors are also included in the model, which is dramatically different from the unconditional relationship of age and wages. The pos-itive unconditional correlation of age with wages becomes a negative correlation when it is conditional on the other included regressors. The slope of the fitted regression line in the added-variable plot is equal to the estimated coefficient on $x_1$ in the fixed-effects regression.[1]

The next subsections explain the statistical basis for added-variable plots. If that is not your interest, please skip to section 3.1—the syntax of `xtavplot`—and to detailed examples of its use in section 5.

## 2.1   Partial regression

The statistical basis for an added-variable plot is partial regression. Partial regression shows that the partial correlation of $x_1$, one of multiple independent variables, with the dependent variable $y$ can be found by "partialing out" the influence of the other independent variables on both $x_1$ and $y$ first and then regressing the partialed $x_1$ on the partialed $y$.

---

1. Note that the added-variable plot is not a good method for evaluating the functional form of the relationship between $x_1$ and $y$, because its validity depends on the assumed linear relationship between $y$ and all the $x$'s, as shown in the next section.

Take the standard linear regression equation relating the dependent variable, $y$, to $K-1$ independent variables $x_1, \ldots, x_{K-1}$, and an intercept term and an error term $\varepsilon$:

$$y_i = \beta_1 x_{1i} + \cdots + \beta_{K-1} x_{K-1,i} + \beta_K + \varepsilon_i$$

The intercept term is placed after the $x$ variables for notational convenience.

If we draw a sample of $N$ observations of data that conform to this relationship, we have $n \times 1$ data vectors of the dependent variable $\mathbf{y}$ and the $K$ independent variables (including $\mathbf{x}_K \equiv \mathbf{1}$, a vector of 1s, for the intercept $\beta_K$), $\mathbf{x}_1, \ldots, \mathbf{x}_K$. Combining all the independent variables into an $n \times K$ matrix $\mathbf{X}$, the data fit the equation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

where $\boldsymbol{\beta}$ is a $K \times 1$ vector of unknown parameters and $\boldsymbol{\varepsilon}$ is an $n \times 1$ vector of the unobserved errors.

The ordinary least-squares (OLS) estimator $\mathbf{b}$ is derived by minimizing the sum of squared residuals ($\widehat{\boldsymbol{\varepsilon}}'\widehat{\boldsymbol{\varepsilon}}$, where $\widehat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{Xb}$) and solving the first-order normal equation

$$\mathbf{X}'\mathbf{Xb} = \mathbf{X}'\mathbf{y} \tag{1}$$

We can partition the $\mathbf{X}$ matrix into $\mathbf{X} = [\mathbf{x}_1 \mathbf{X}_2]$, where $\mathbf{X}_2 = [\mathbf{x}_2 \ldots \mathbf{x}_K]$; partition $\mathbf{b}$ into $\mathbf{b} = \begin{bmatrix} b_1 \\ \mathbf{b}_2 \end{bmatrix}$, where $\mathbf{b}_2 = \begin{bmatrix} b_2 \\ \vdots \\ b_K \end{bmatrix}$; and rewrite (1) as

$$\begin{bmatrix} \mathbf{x}_1'\mathbf{x}_1 & \mathbf{x}_1'\mathbf{X}_2 \\ \mathbf{X}_2'\mathbf{x}_1 & \mathbf{X}_2'\mathbf{X}_2 \end{bmatrix} \begin{bmatrix} b_1 \\ \mathbf{b}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1'\mathbf{y} \\ \mathbf{X}_2'\mathbf{y} \end{bmatrix}$$

With some manipulation, we can solve for $b_1 = (\mathbf{x}_1'\mathbf{M}_2\mathbf{x}_1)^{-1}\mathbf{x}_1'\mathbf{M}_2\mathbf{y}$, where $\mathbf{M}_2 = (\mathbf{I} - \mathbf{X}_2(\mathbf{X}_2'\mathbf{X}_2)^{-1}\mathbf{X}_2')$. Because $\mathbf{M}_2$ is symmetric and idempotent, we can rewrite $b_1$ as

$$b_1 = (\mathbf{x}_1'\mathbf{M}_2'\mathbf{M}_2\mathbf{x}_1)^{-1}\mathbf{x}_1'\mathbf{M}_2'\mathbf{M}_2\mathbf{y} = (\mathbf{e}_{\mathbf{x}_1}'\mathbf{e}_{\mathbf{x}_1})^{-1}\mathbf{e}_{\mathbf{x}_1}'\mathbf{e}_{\mathbf{y}} \tag{2}$$

where $\mathbf{e}_{\mathbf{x}_1} = \mathbf{M}_2\mathbf{x}_1$ and $\mathbf{e}_{\mathbf{y}} = \mathbf{M}_2\mathbf{y}$.

By inspecting the equation for $\mathbf{M}_2$, we can see that $\mathbf{e}_{\mathbf{y}} = \mathbf{M}_2\mathbf{y}$ is the vector of residuals from the regression of $\mathbf{y}$ on $\mathbf{X}_2$, and likewise, $\mathbf{e}_{\mathbf{x}_1} = \mathbf{M}_2\mathbf{x}_1$ is the vector of residuals from the regression of $\mathbf{x}_1$ on $\mathbf{X}_2$.

$\mathbf{e}_{\mathbf{y}}$ and $\mathbf{e}_{\mathbf{x}_1}$ can be interpreted as $\mathbf{y}$ and $\mathbf{x}_1$ purged of the influence of the $\mathbf{X}_2$ variables. $\mathbf{e}_{\mathbf{y}} = \mathbf{y} - \widehat{\mathbf{y}}_{\mathbf{x}_2}$, where $\widehat{\mathbf{y}}_{\mathbf{x}_2}$ is the predicted value of $\mathbf{y}$ from the regression of $\mathbf{y}$ on $\mathbf{X}_2$. That is, $\mathbf{e}_{\mathbf{y}}$ is what is left over when all the variation in $\mathbf{y}$ that can be predicted by $\mathbf{X}_2$ has been subtracted out. The process is similar for $\mathbf{e}_{\mathbf{x}_1}$. So the correlation of $\mathbf{e}_{\mathbf{y}}$ and $\mathbf{e}_{\mathbf{x}_1}$ is the partial correlation $\mathbf{y}$ and $\mathbf{x}$ conditional on $\mathbf{X}_2$.

This decomposition gives rise to the added-variable plot. A scatterplot of the values in $\mathbf{e}_{\mathbf{x}_1}$ versus $\mathbf{e}_{\mathbf{y}}$ will show the correlation of the $x_1$ variable with the $y$ variable, controlling for the influence of the other independent variables in the multiple regression.

From (2), we can see that the OLS estimator $b_1$ of $\beta_1$ is the result of regressing $\mathbf{e_y}$ on $\mathbf{e_{x_1}}$ (with no intercept term). Thus, the OLS linear fit of the data in the scatterplot of $\mathbf{e_{x_1}}$ versus $\mathbf{e_y}$ is equal to $b_1$, the estimated partial effect of $x_1$ on $y$.

This is what we were seeking: a way of displaying the relationship between $x_1$ and $y$, controlling for the effect of the other independent variables in the regression. An added-variable plot creates a scatterplot of $\mathbf{e_{x_1}}$ versus $\mathbf{e_y}$ and displays the linear fit line with confidence interval boundaries above and below the regression line. The regression line has a slope of $b_1$.

## 2.2   Partial regression of transformed variables

The derivation of partial regression above applies only to OLS estimation because it results from the OLS normal equation (1). However, we can derive a partial-regression formula for non-OLS estimation methods if their estimating equations can be transformed so that they meet OLS assumptions.[2] The fixed-effects, between-effects, and random-effects panel-estimation methods can each be represented as transformations of the original model, which can then be fit by OLS yielding the $\boldsymbol{\beta}$ coefficient estimates we are seeking.

If the transformed variables $\mathbf{y}^*$, $\mathbf{x}_1^*$, and $\mathbf{X}_2^*$ conform to OLS assumptions, the equation

$$\mathbf{y}^* = \mathbf{x}_1^*\beta_1 + \mathbf{X}_2^*\boldsymbol{\beta}_2 + \varepsilon^*$$

results in the OLS normal equation

$$\begin{bmatrix} \mathbf{x}_1^{*\prime}\mathbf{x}_1^* & \mathbf{x}_1^{*\prime}\mathbf{X}_2^* \\ \mathbf{X}_2^{*\prime}\mathbf{x}_1^* & \mathbf{X}_2^{*\prime}\mathbf{X}_2^* \end{bmatrix} \begin{bmatrix} b_1 \\ \mathbf{b}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^{*\prime}\mathbf{y}^* \\ \mathbf{X}_2^{*\prime}\mathbf{y}^* \end{bmatrix}$$

As above,

$$b_1 = (\mathbf{x}_1^{*\prime}\mathbf{M}_2^*\mathbf{x}_1^*)^{-1}\mathbf{x}^{*\prime}\mathbf{M}_2^*\mathbf{y}^* = (\mathbf{e}_{\mathbf{x}_1^*}^{\prime}\mathbf{e}_{\mathbf{x}_1^*})^{-1}\mathbf{e}_{\mathbf{x}_1^*}^{\prime}\mathbf{e}_{\mathbf{y}^*}$$
$$\text{for } \mathbf{M}_2^* = \mathbf{I} - \mathbf{X}_2^*(\mathbf{X}_2^{*\prime}\mathbf{X}_2^*)^{-1}\mathbf{X}_2^{*\prime} \qquad \mathbf{e}_{\mathbf{x}_1^*} = \mathbf{M}_2^*\mathbf{x}_1^* \quad \text{and} \quad \mathbf{e}_{\mathbf{y}^*} = \mathbf{M}_2^*\mathbf{y}^*$$

$$(3)$$

The next three subsections apply the partial-regression formula for a transformed estimating equation to three panel-data estimation methods: fixed effects, between effects, and random effects.

## 2.3   Fixed-effects estimation

Fixed-effects estimation is just a computationally efficient way of estimating OLS coefficients incorporating a separate intercept for each cross-sectional unit in the panel-data sample. Direct computation using OLS with dummy variables for each unit is straightforward but cumbersome. In the typical situation, where the number of cross-sectional

---

2. This is the idea behind the typical proof of the properties of generalized least-squares (GLS) estimation.

units $n$ is large and the number of time-series observations per unit $T_i$ is small, unit-specific intercepts result in many dummy variables, and their coefficients are usually not of interest in themselves (or consistently estimated). Fixed-effects estimation transforms the estimating equation to eliminate the numerous intercept terms. Estimating the transformed equation via OLS still delivers the same coefficients and standard errors (after a degrees-of-freedom adjustment) as direct computation, making the estimation faster and more convenient.

Given panel data on individuals or units indexed by $i \in \{1, \ldots, n\}$ for multiple time periods $t \in \{1, \ldots, T_i\}$, consider the linear model

$$y_{it} = \mathbf{x}_{it}\boldsymbol{\beta} + \upsilon_i + \varepsilon_{it} \tag{4}$$

where $\mathbf{x}_{it}$ is a $1 \times K$ row vector of independent variables and $\upsilon_i$ is an individual or unit-specific intercept term that is assumed to be uncorrelated with the error term $\varepsilon_{it}$. The advantage of including the individual intercepts is that they control for all characteristics of the individual that do not change over time. Without panel data, one could not control for fixed individual characteristics without gathering data on each of the characteristics. This model can be fit using OLS by including dummy variables for each individual in the sample. Because the individual intercepts are not typically of interest, however, one can save time and effort by subtracting out their effects.

Taking the average of the observations over each individual, (4) becomes

$$\overline{y}_i = \overline{\mathbf{x}}_i\boldsymbol{\beta} + \upsilon_i + \overline{\varepsilon}_i \tag{5}$$

where $\overline{y}_i = 1/T_i \sum_{T_i} y_{it}$, $\overline{\mathbf{x}}_i = 1/T_i \sum_{T_i} \mathbf{x}_{it}$, and $\overline{\varepsilon}_i = 1/T_i \sum_{T_i} \varepsilon_{it}$. Subtracting (5) from (4),

$$y_{it} - \overline{y}_i = (\mathbf{x}_{it} - \overline{\mathbf{x}}_i)\boldsymbol{\beta} + \varepsilon_{it} - \overline{\varepsilon}_i$$

which cancels out all the $\upsilon_i$ terms, dramatically reducing the dimensionality of the estimation when $n$ is large. This can be rewritten as

$$y_{it}^* = \mathbf{x}_{it}^*\boldsymbol{\beta} + \varepsilon_{it}^* \tag{6}$$

where $y_{it}^* = y_{it} - \overline{y}_{it}$, $\mathbf{x}_{it}^* = \mathbf{x}_{it} - \overline{\mathbf{x}}_{it}$, and $\varepsilon_{it}^* = \varepsilon_{it} - \overline{\varepsilon}_{it}$.

Fixed-effects estimation applies OLS to (6) to estimate the $\boldsymbol{\beta}$ coefficients efficiently.[3]

One could apply the partitioned regression formula in (3) to (6) to derive residuals $\mathbf{e}_{\mathbf{y}^*}$ and $\mathbf{e}_{\mathbf{x}_1^*}$. These could be plotted, and the slope of their linear fit would be $b_1$. However, the meaning of the residuals is not intuitive. $\mathbf{e}_{\mathbf{y}^*}$ is a vector of $y_{it}^*$ controlling for $\mathbf{x}_{2it}^*$ (where $\mathbf{x}_{it}^* = [x_{1it}^* \; \mathbf{x}_{2it}^*]$), not $y_{it}$ controlling for $\mathbf{x}_{2it}$. Similarly, $\mathbf{e}_{\mathbf{x}_1^*}$ is a vector of $x_{1it}^*$ controlling for $\mathbf{x}_{2it}^*$, not $x_{1it}$ controlling for $\mathbf{x}_{2it}$.

---

3. Although not often mentioned, the fixed-effects transformation of the error terms $\varepsilon_{it}^* = \varepsilon_{it} - \overline{\varepsilon}_i$ violates OLS assumptions because it introduces both serial correlation and heteroskedasticity (if the $T_i$ are not identical) into the transformed error. Nonetheless, OLS estimation of the transformed equation provides efficient estimates of $\boldsymbol{\beta}$ because the transformed $x_{it}^*$ cancel out the problem. See Wooldridge (2010, 305).

It is straightforward, however, to calculate the OLS $\mathbf{e_y}$ and $\mathbf{e_{x_1}}$ from the fixed-effects $\mathbf{e_{y^*}}$ and $\mathbf{e_{x_1^*}}$. $\mathbf{e_{y^*}}$ is the fixed-effects residual from the regression of $\mathbf{x_2^*}$ on $\mathbf{y^*}$, producing the coefficient $\mathbf{b}_{y^*|\mathbf{x_2^*}}$. An element of $\mathbf{e_{y^*}}$ is $e_{y_{it}^*} = y_{it} - \overline{y}_i - (\mathbf{x}_{2it} - \overline{\mathbf{x}}_{2i})\mathbf{b}_{y^*|\mathbf{x_2^*}}$. The fixed-effects coefficient $\mathbf{b}_{y^*|\mathbf{x_2^*}}$ is exactly equal to the OLS coefficient $\mathbf{b}_{y|\mathbf{x_2}}$ from regressing $\mathbf{x}_{2it}$ and $v_i$ on $y_{it}$.[4] So, $e_{y_{it}^*} = (y_{it} - \mathbf{x}_{2it}\mathbf{b}_{y|\mathbf{x_2}}) - (\overline{y}_i - \overline{\mathbf{x}}_{2i}\mathbf{b}_{y|\mathbf{x_2}})$. The second term, $(\overline{y}_i - \overline{\mathbf{x}}_{2i}\mathbf{b}_{y|\mathbf{x_2}}) = u_{(y|\mathbf{x_2})i}$, is the OLS estimate of the individual effect. Hence, $e_{y_{it}^*} = e_{y_{it}} - u_{(y|\mathbf{x_2})i}$ and $\mathbf{e_y} = \mathbf{e_{y^*}} + \mathbf{u_y}$, where $\mathbf{u_y}$ is an $(N = \sum_i T_i) \times 1$ vector of $u_{(y|\mathbf{x_2})i}$. Similarly, $\mathbf{e_{x_1}} = \mathbf{e_{x_1^*}} + \mathbf{u_{x_1}}$. That means that one can readily calculate the more intuitive OLS residuals $\mathbf{e_{x_1}}$ and $\mathbf{e_y}$ from the fixed-effects estimates.

So, in the case of fixed effects, the estimation of the transformed (6) produces $\mathbf{b}$ coefficients identical to those from a direct OLS estimation of (4). The fixed-effects estimates are used to transform the fixed-effects residuals $\mathbf{e_{x_1^*}}$ and $\mathbf{e_{y^*}}$ into the OLS residuals of $\mathbf{e_{x_1}}$ and $\mathbf{e_y}$ to create an added-variable plot whose fitted regression line has slope $b_1$.

## 2.4    Between-effects estimation

Between-effects estimation applies OLS to the $n$ unique individual mean values of (5), taking $v_i$ as part of the error term because it is not separately identifiable.

The per-individual averages are transformations of the original $y$ and $\mathbf{x}$ variables, so we can apply the partial regression of transformed variables in (3), where

$$\mathbf{y}^* = \begin{bmatrix} \overline{y}_1 \\ \vdots \\ \overline{y}_n \end{bmatrix} \quad \text{and} \quad \mathbf{X}^* = \begin{bmatrix} \overline{\mathbf{x}}_1 \\ \vdots \\ \overline{\mathbf{x}}_n \end{bmatrix}$$

Then, $\mathbf{e_{y^*}}$ and $\mathbf{e_{x_1^*}}$ provide the data points for the added-variable plot. In this case, $\mathbf{e_{y^*}}$ and $\mathbf{e_{x_1^*}}$ are rather intuitive. The plot shows the relationship of the individual means of $y$ versus the means of $x_1$ controlling for the influence of the means of $\mathbf{x_2}$.

## 2.5    Random-effects estimation

Random-effects estimation considers the same model as fixed-effects estimation in (4) but interprets the individual effects $v_i$ as belonging to the error term. This means the error terms $v_i + \varepsilon_{it}$ are not independent and identically distributed as required for efficient estimation by OLS. The model, however, reveals the structure of the errors, so it can be estimated by generalized least squares (GLS). GLS is estimated by applying OLS estimation to transformations of the observed variables, which renders the transformed errors independent and identically distributed.

---

4. One can show that $\mathbf{b}_{y^*|\mathbf{x_2^*}} = \mathbf{b}_{y|\mathbf{x_2}}$ by applying the partial-regression formula in (2) because $y_{it}^* = y_{it} - \overline{y}_i$ and $\mathbf{x}_{2it}^* = \mathbf{x}_{2it} - \overline{\mathbf{x}}_{2i}$ are the residuals from regressing the $v_i$ individual dummy variables on $y_{it}$ and $\mathbf{x}_{2it}$. That is, fixed-effects regression itself is an application of partial regression.

The appropriate transformation of the panel-data model in (4) for feasible GLS estimation is

$$y_{it} - \widehat{\theta}_i \overline{y}_i = \left( \mathbf{x}_{it} - \widehat{\theta}_i \overline{\mathbf{x}}_i \right) \boldsymbol{\beta} + \left( 1 - \widehat{\theta}_i \right) \upsilon_i + \varepsilon_{it} - \widehat{\theta}_i \overline{\varepsilon}_i$$

where $\widehat{\theta}_i = 1 - \{\widehat{\sigma}_\varepsilon^2/(T_i \widehat{\sigma}_\upsilon^2 + \widehat{\sigma}_\varepsilon^2)\}$. $\widehat{\sigma}_\upsilon^2$ and $\widehat{\sigma}_\varepsilon^2$ are estimates of the variances of $\upsilon_i$ and $\varepsilon_{it}$, respectively.

We can apply the partial regression of transformed variables in (3), where

$$\mathbf{y}^* = \begin{bmatrix} y_{11} - \widehat{\theta}_1 \overline{y}_1 \\ \vdots \\ y_{1T_1} - \widehat{\theta}_1 \overline{y}_1 \\ \vdots \\ y_{n1} - \widehat{\theta}_n \overline{y}_n \\ \vdots \\ y_{nT_n} - \widehat{\theta}_n \overline{y}_n \end{bmatrix} \quad \text{and} \quad \mathbf{X}^* = \begin{bmatrix} \mathbf{x}_{11} - \widehat{\theta}_1 \overline{\mathbf{x}}_1 \\ \vdots \\ \mathbf{x}_{1T_1} - \widehat{\theta}_1 \overline{\mathbf{x}}_1 \\ \vdots \\ \mathbf{x}_{n1} - \widehat{\theta}_n \overline{\mathbf{x}}_n \\ \vdots \\ \mathbf{x}_{nT_n} - \widehat{\theta}_n \overline{\mathbf{x}}_n \end{bmatrix} \tag{7}$$

enabling us to construct $\mathbf{e}_{\mathbf{y}^*}$ and $\mathbf{e}_{\mathbf{x}_1^*}$. Regressing $\mathbf{e}_{\mathbf{x}_1^*}$ on $\mathbf{e}_{\mathbf{y}^*}$ produces the coefficient $b_1$, but unlike fixed-effects estimates, the residuals cannot be converted into OLS residuals $\mathbf{e}_{\mathbf{y}}$ and $\mathbf{e}_{\mathbf{x}_1}$ and still have a fitted regression slope of $b_1$. Therefore, we make the added-variable plot out of $\mathbf{e}_{\mathbf{y}^*}$ and $\mathbf{e}_{\mathbf{x}_1^*}$, which have a somewhat intuitive interpretation as heteroskedasticity-corrected residuals.[5]

The added-variable plot of $\mathbf{e}_{\mathbf{y}^*}$ and $\mathbf{e}_{\mathbf{x}_1^*}$ presents the contribution of each data point $(x_{1it}, y_{it})$ to the estimated coefficient $b_1$, so the plot is a good visual diagnostic for outlier observations having a large influence on the estimated relationship, just as in the OLS, fixed-effects, or between-effects cases.

## 2.6 Maximum-likelihood random-effects and population-averaged model

The maximum likelihood estimation of neither the random-effects (`xtreg, mle`) nor the population-averaged model (`xtreg, pa`) can be represented as a transformed partial-regression in the form of (3) in the way OLS and GLS estimators can. `xtavplot` cannot be used after these estimation methods. This may not be much of a loss in the case of `xtreg, mle`. The *Methods and formulas* section of [XT] **xtreg** notes that it yields "essentially the same results" as `xtreg, re` except when the sample is small ($\leq 200$ observations) and unbalanced.

---

5. The errors are not independently and identically distributed because of autocorrelation between the errors for each individual caused by the individual effects (a clustering effect), as well as heteroskedasticity across individuals if the time spans $T_i$ vary across $i$.

# 3  The xtavplot and xtavplots commands

## 3.1  Syntax

xtavplot *indepvar* [ , *options* ]

xtavplots [ , *options* ]

| *options* | Description |
|-----------|-------------|
| *marker_options* | change look of markers (color, size, etc.) |
| *marker_label_options* | add marker labels; change look or position |
| | |
| <u>rl</u>opts(*cline_options*) | affect rendition of the regression line |
| <u>noco</u>ef | turn off display of coefficient below graph |
| | |
| ciopts(*cline_options*) | affect rendition of the confidence interval line |
| noci | turn off confidence interval |
| <u>ciu</u>nder | graph confidence interval underneath scatterplot |
| <u>level</u>(#) | specify the confidence level |
| ciplot(*plottype*) | how to plot confidence intervals; |
| | default is ciplot(rline); |
| | a common alternative is ciplot(rarea) |
| | |
| *twoway_options* | any options documented in [G-3] ***twoway_options***, |
| | except for by() |
| | |
| <u>addm</u>eans | rescale the residuals, regression line, and |
| | confidence intervals to be centered on |
| | the means of $x$ and $y$ instead of zero |
| **xtavplot-only options** | |
| xlim(#[ # ]), ylim(#[ # ]) | limit the ranges of the $x$ and $y$ residuals displayed |
| <u>g</u>enerate(*exvar eyvar*) | save the values of $x$ and $y$ residuals in new variables |
| <u>nod</u>isplay | suppress display of the plot |
| addplot(*plot*) | add other plots to the generated graph |
| **xtavplots-only option** | |
| *combine_options* | any of the options documented in |
| | [G-2] ***graph combine*** |

## 3.2 Description of xtavplot and xtavplots

xtavplot creates an added-variable plot (also known as a partial-regression leverage plot, a partial-regression plot, or an adjusted partial-residual plot) after xtreg, fe (fixed-effects estimation), xtreg, re (random-effects estimation), or xtreg, be (between-effects estimation). xtavplot cannot be used after xtreg, mle or xtreg, pa.

xtavplots creates a matrix of added-variable plots of all the *indepvars*.

*indepvar* is an independent ($x$) variable (also known as a predictor, carrier, or covariate) that may or may not have been included in the preceding estimation. The user would choose an *indepvar* not already in the estimation to evaluate whether to include it.

xtavplot shows the partial correlation between one *indepvar* and the *depvar* from a multivariate panel regression.

Besides showing the relationship between the *indepvar* and the *depvar* controlling for the other regressors, xtavplot is useful for visually identifying which outlier observations have a big effect on the estimated coefficient.

After fixed-effects estimation, the plotted e(x|X) values are the residuals from the regression of $x_1$ on the other $\mathbf{x}_2$ variables in the original regression, and the plotted e(y|X) values are the residuals from the regression of $y$ on the other $\mathbf{x}_2$ variables.

After between-effects estimation, e(av.x|av.X) and e(av.y|av.X) are the residuals from the regression of per-unit means $\overline{x}_{1i}$ and $\overline{y}_i$ on the per-unit means $\overline{\mathbf{x}}_{2i}$ of the other independent variables.

After random-effects estimation, e(x*|X*) and e(y*|X*) are the residuals from the regression of heteroskedasticity-corrected $x_1^*$ and heteroskedasticity-corrected $y^*$ on the other heteroskedasticity-corrected independent $\mathbf{x}_2^*$ variables.

The fitted line shown in the graph is the least-squares fit between the residuals. For each of the three panel-data estimation methods, the fitted line has the same slope as the estimated coefficient on the *indepvar* in the preceding regression.

Because of their construction, the residuals each have a mean of zero, and the regression line fit between them passes exactly through e(x|X)=0 and e(y|X)=0. At that point, the confidence interval has zero width, giving it an unfamiliar shape.[6]

## 3.3 Options for xtavplot and xtavplots

*marker_options* affect the rendition of markers drawn at the plotted points, including their shape, size, color, and outline; see [G-3] **marker_options**.

---

6. The confidence interval for a conventional regression with no constant term also has this shape at the point where all the independent variables have a value of zero.

*marker_label_options* specify if and how markers are to be labeled; see
   [G-3] **marker_label_options**.

rlopts(*cline_options*) affects the rendition of the regression (fitted) line; see
   [G-3] **cline_options**.

nocoef turns off the display below the graph of the values of the regression coefficient,
   standard error, and *t* statistic.

ciopts(*cline_options*) affects how the upper and lower confidence interval lines are
   rendered; see [G-3] **cline_options**. If you specify ciplot(), then rather than using
   *cline_options*, you should specify what options are appropriate for the *plottype*.

noci turns off the display of the confidence interval on the graph.

ciunder causes the confidence interval to be graphed underneath the scatterplot (that
   is, the scatter points are visible on top of the confidence interval). This is mainly
   useful when graphing a solid confidence interval with the option ciplot(rarea).

level(*#*) specifies the confidence level, as a percentage, for confidence intervals
   around the regression line. The default is level(95) or as set by set level; see
   [U] **20.8 Specifying the width of confidence intervals**.

ciplot(*plottype*) specifies how the confidence interval is to be plotted. The default
   is ciplot(rline), meaning that the prediction will be plotted by graph twoway
   rline.

   A common alternative is ciplot(rarea), which will substitute shading around the
   prediction line. See [G-2] **graph twoway** for a list of *plottype* choices. You may
   choose any *plottype*s that expect two *y* variables and one *x* variable.

*twoway_options* are any of the options documented in [G-3] **twoway_options**, excluding
   by(). These include options for titling the graph (see [G-3] **title_options**) and saving
   the graph to disk (see [G-3] **saving_option**).

addmeans rescales the scatterplot values, the regression line, and the confidence inter-
   vals to be centered on the mean values of the *x* and *y* variables instead of being
   centered on zero by default. This may make the plot more visually intuitive, but it
   is important to make clear to viewers that the graph is showing conditioned values
   rather than the original data.

## 3.4   Options only for xtavplot

xlim(*#*[ *#* ]) and ylim(*#*[ *#* ]) constrain the range of the *indepvar* and *depvar* resid-
   uals displayed. If only one number is specified, residuals with a value below that
   number will not be displayed in the scatterplot. If two numbers are specified, resid-
   uals below the first number and above the second number will not be displayed.

   Excluding observations of the residuals does not affect the slope of the regression line
   in the graph. The purpose of these options is to avoid a few outlying observations

dramatically extending the range of the $x$ or $y$ axis, thus obscuring the display of the relationship between the variables. Because panel datasets are typically large, it is common to have a few distant outliers that do not significantly affect the estimates. Make sure that the undisplayed observations are not important to the estimated relationship and that their exclusion is noted in the text.

generate(*exvar eyvar*) saves the values of the $x$ and $y$ residuals in variables named by the user. The user must specify two variable names for *exvar* and *eyvar*. These residuals can be used for subsequent calculations or graphing commands. See sections 3.6 and 4 below for how to access the estimate $b_1$ and its standard error and how to calculate the regression fit and confidence intervals.

nodisplay suppresses display of the plot. This is mainly useful for users creating their own plots from variables created with generate().

addplot(*plot*) provides a way to add other plots to the generated graph; see [G-3] ***addplot_option***.

## 3.5    Options only for xtavplots

*combine_options* are any of the options documented in [G-2] **graph combine** for arranging a matrix of plots in a single image.

## 3.6    Stored results

xtavplot stores the following in r():

Scalars
    r(coef)               estimated coefficient of the added variable
    r(se)                 standard error of the estimated coefficient

    After the addmeans option:

Scalars
    r(ybar)               (possibly weighted) mean of the *depvar* **y**
    r(xbar)               (possibly weighted) mean of the added variable $\mathbf{x}_1$

# 4    Methods and formulas

Because xtavplot is an xtreg postestimation command, the preceding xtreg command will have the form

$$\text{xtreg } y \ \ x_1 \ \mathbf{x}_2, model \tag{8}$$

where $y$ is the *depvar*, $x_1$ is one of the *indepvars*, $\mathbf{x}_2$ is a vector of the other *indepvars*, and *model* is a choice of fe, be, or re. This will be followed by the command

$$\text{xtavplot } x_1, options$$

xtavplot allows for $x_1$ not to be included in the preceding xtreg *indepvars*. In that case, there is some adjustment to these formulas, principally to fit the full xtreg model including $x_1$.

## 4.1   After xtreg, fe

xtavplot calculates residuals $\mathbf{e_y}$ and $\mathbf{e_{x_1}}$ in (2) from

$$\text{xtreg } y \text{ } \mathbf{x}_2, \text{fe}$$
$$\text{predict } \mathbf{e_y}, \text{xbu}$$
$$\text{xtreg } x_1 \text{ } \mathbf{x}_2, \text{fe}$$
$$\text{predict } \mathbf{e_{x_1}}, \text{xbu}$$

using the same weights and sample restrictions as (8).

## 4.2   After xtreg, be

xtavplot forms the $n$ individual means $\overline{y}$, $\overline{x}_1$, and $\overline{\mathbf{x}}_2$ as defined in (5). Residuals $\mathbf{e_{y^*}}$ and $\mathbf{e_{x_1^*}}$ in (3) are calculated from

$$\text{regress } \overline{y} \text{ } \overline{\mathbf{x}}_2$$
$$\text{predict } \mathbf{e_{y^*}}, \text{residuals}$$
$$\text{regress } \overline{x}_1 \text{ } \overline{\mathbf{x}}_2$$
$$\text{predict } \mathbf{e_{x_1^*}}, \text{residuals}$$

using the weights and sample of (8).

## 4.3   After xtreg, re

xtavplot forms the weighted deviations from the mean variables $\mathbf{y}^*$, $\mathbf{x}_1^*$, and $\mathbf{X}_2^*$ as defined in (7), where $\mathbf{X}^* = \begin{bmatrix} \mathbf{x}_1^* & \mathbf{X}_2^* \end{bmatrix}$. The weights $\widehat{\theta}_i = 1 - (\widehat{\sigma}_\varepsilon^2)/(T_i\widehat{\sigma}_v^2 + \widehat{\sigma}_\varepsilon^2)$ are calculated from $\widehat{\sigma}_\varepsilon^2 = $ e(sigma_e)^2 and $\widehat{\sigma}_v^2 = $ e(sigma_u)^2 from the preceding xtreg, re command. Define the $(N = \sum_i T_i) \times 1$ vector

$$(\mathbf{1} - \theta) = \begin{bmatrix} 1 - \widehat{\theta}_1 \\ \vdots \\ 1 - \widehat{\theta}_1 \\ \vdots \\ 1 - \widehat{\theta}_n \\ \vdots \\ 1 - \widehat{\theta}_n \end{bmatrix}$$

where each $1 - \widehat{\theta}_i$ is repeated $T_i$ times.

$\mathbf{e_{y^*}}$ and $\mathbf{e_{x_1^*}}$ are calculated from

$$\texttt{regress } \mathbf{y}^* \; (\mathbf{1}-\theta) \; \mathbf{X}_2^*, \texttt{ noconstant}$$
$$\texttt{predict } \mathbf{e_{y^*}}, \texttt{residuals}$$
$$\texttt{regress } \mathbf{x}_1^* \; (\mathbf{1}-\theta) \; \mathbf{X}_2^*, \texttt{ noconstant}$$
$$\texttt{predict } \mathbf{e_{x_1^*}}, \texttt{residuals}$$

using the sample of (8) (weights are not allowed in $\texttt{xtreg, re}$ estimation).

Note that it does not work to use $\texttt{xtreg } y \; \mathbf{x}_2, \texttt{ re}$ and $\texttt{xtreg } x_1 \; \mathbf{x}_2, \texttt{ re}$ to generate residuals, because they will estimate different values for $\widehat{\sigma}_\varepsilon^2$ and $\widehat{\sigma}_v^2$, which vary depending on the included *indepvars*.

## 4.4    Confidence interval

The preceding subsections explain how to calculate the residuals $\mathbf{e_y}$ and $\mathbf{e_{x_1}}$ (or $\mathbf{e_{y^*}}$ and $\mathbf{e_{x_1^*}}$, as appropriate throughout this section). It is not necessary to regress one residual on the other to calculate the coefficient $b_1$ and its standard error $\widehat{\sigma}_{b_1}$, because they are already available from the preceding $\texttt{xtreg}$ command.[7]

By default, $\texttt{xtavplot}$ displays a confidence interval around the predicted fit from the regression of $\mathbf{e_{x_1}}$ on $\mathbf{e_y}$. The fitted values of $\mathbf{e_y}$ are $\widehat{\mathbf{e}}_\mathbf{y} = \mathbf{e_{x_1}} b_1$. The confidence interval boundaries are $\widehat{\mathbf{e}}_\mathbf{y} \pm t_{\alpha/2} \mathbf{e_{x_1}} \widehat{\sigma}_{b_1}$ for fixed-effects and between-effects estimates and $\widehat{\mathbf{e}}_\mathbf{y} \pm z_{\alpha/2} \mathbf{e_{x_1}} \widehat{\sigma}_{b_1}$ for random-effects estimates, where $t_{\alpha/2}$ is the $\alpha/2$ percentile of the cumulative $t$ distribution, $z_{\alpha/2}$ is the $\alpha/2$ percentile of the cumulative standard normal distribution, and $\alpha = 1 - \texttt{level}/100$.

## 4.5    The addmeans option

The $\texttt{addmeans}$ option recenters the graph on the mean values of $\mathbf{y}$ and $\mathbf{x}_1$, instead of the default of zero. The mean $\overline{y}$ of $\mathbf{y}$ and $\overline{x}_1$ of $\mathbf{x}_1$ are calculated using the weights and sample restrictions in the preceding $\texttt{xtreg}$ command. $\overline{x}_1$ is added to the residuals $\mathbf{e_{x_1}}$, and $\overline{y}$ is added to $\mathbf{e_y}$, the predicted values, and the confidence interval boundaries

---

7. This also eliminates the need to worry about heteroskedasticity corrections that may have been implemented in the preceding regression because they affect only the standard errors of the estimates, not the values of the residuals $\mathbf{e_{x_1}}$ and $\mathbf{e_y}$. If the user is interested in verifying that the residuals are calculated correctly (consistent with the coefficient $b_1$ in the preceding regression), there is an otherwise undocumented $\texttt{xtavplot}$ option, $\texttt{debug}$, that calculates $b_1$ as the coefficient on $\mathbf{e_{x_1}}$ from

$$\texttt{regress } \mathbf{e_y} \; \mathbf{e_{x_1}}, \texttt{ noconstant}$$

and stores the result in $\texttt{r(b\_check)}$. This regression does not calculate the correct standard error for $b_1$, which requires an adjustment for the additional degrees of freedom taken up by controlling for the influence of the other covariates. The correct standard errors can be calculated using the undocumented $\texttt{regress}$ option $\texttt{dof()}$ to change the degrees of freedom:

$$\texttt{regress } \mathbf{e_y} \; \mathbf{e_{x_1}}, \texttt{ noconstant dof}(\mathit{df})$$

where $df = \texttt{e(N)} - \texttt{e(df\_m)} - 1$ after $\texttt{xtreg, fe}$ and $\texttt{xtreg, re}$ and $df = \texttt{e(df\_r)}$ after $\texttt{xtreg, be}$.

before the graph is displayed. The means are not added to the values of $\mathbf{e_{x_1}}$ and $\mathbf{e_y}$ saved by the generate() option, but $\overline{y}$ and $\overline{x}_1$ are saved as r(ybar) and r(xbar) in the return values.

# 5 Examples of xtavplot and xtavplots in use

Because xtavplot and xtavplots are xtreg postestimation commands, we first load an example Stata panel dataset, nlswork.dta. We keep only the first 1,000 observations of the large dataset so that the graphs display more quickly. Use xtreg to fit a fixed-effects model of the correlates of wages. The specification of the model is discussed in help xtreg.

```
. webuse nlswork
(National Longitudinal Survey.  Young Women 14-26 years of age in 1968)

. keep in 1/1000
(27,534 observations deleted)

. xtreg ln_w tenure c.tenure#c.tenure ttl_exp not_smsa south, fe
```

| Fixed-effects (within) regression | | | | Number of obs | = | 989 |
| Group variable: idcode | | | | Number of groups | = | 163 |

| R-sq: | | | Obs per group: | | |
| within  = 0.1840 | | | min = | 1 |
| between = 0.2753 | | | avg = | 6.1 |
| overall = 0.2004 | | | max = | 15 |
| | | | F(5,821) | = | 37.03 |
| corr(u_i, Xb)  = 0.1490 | | | Prob > F | = | 0.0000 |

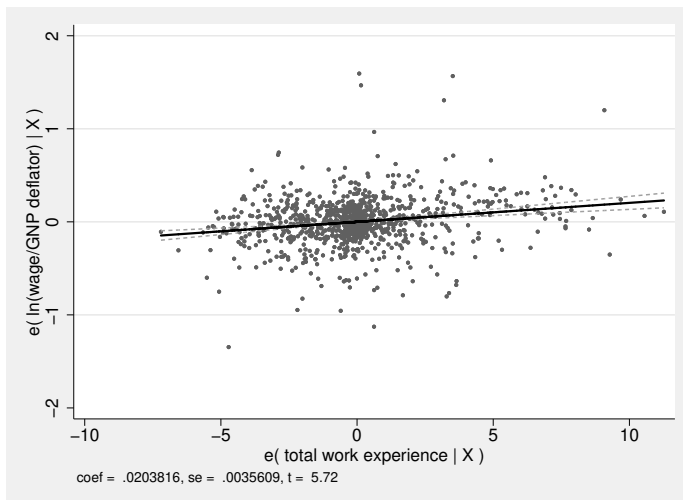| ln_wage | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| tenure | .0379093 | .0076476 | 4.96 | 0.000 | .0228981 | .0529206 |
| c.tenure#c.tenure | -.0014394 | .0004394 | -3.28 | 0.001 | -.0023018 | -.000577 |
| ttl_exp | .0203816 | .0035609 | 5.72 | 0.000 | .013392 | .0273712 |
| not_smsa | -.0450833 | .0707906 | -0.64 | 0.524 | -.1840351 | .0938685 |
| south | -.0727986 | .0986778 | -0.74 | 0.461 | -.2664892 | .1208919 |
| _cons | 1.626667 | .0202064 | 80.50 | 0.000 | 1.587005 | 1.666329 |
| sigma_u | .34239623 | | | | | |
| sigma_e | .27343844 | | | | | |
| rho | .61058793 | (fraction of variance due to u_i) | | | | |

```
F test that all u_i=0: F(162, 821) = 7.68                Prob > F = 0.0000
```

Invoking the command xtavplot ttl_exp will display a graph of the partial correlation between ttl_exp and ln_wage, giving a sense of how closely the individual observations fit this relationship. The slope of the regression of residuals e(ttl_exp|X) on e(ln_wage|X) is shown as a solid line, and the limits of its confidence interval are shown as dashed lines.
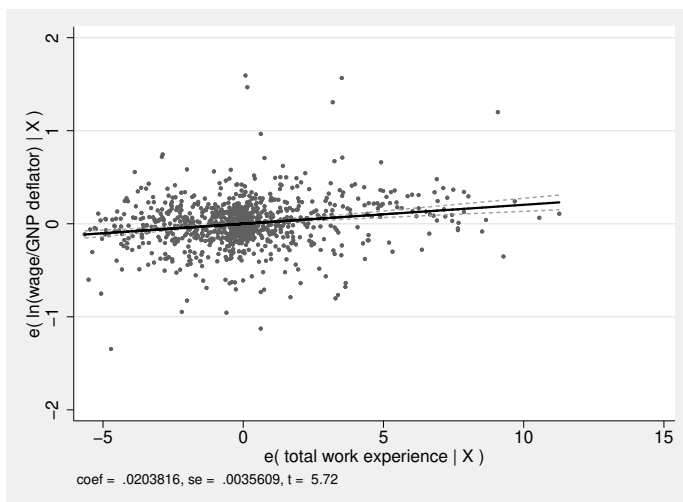
```
. xtavplot ttl_exp
```



The graph has excessive white space to the left of the data because of one observation with a value of e(ttl_exp|X) equal to −6.2. When we add the option xlim(-6), the graph is better situated:
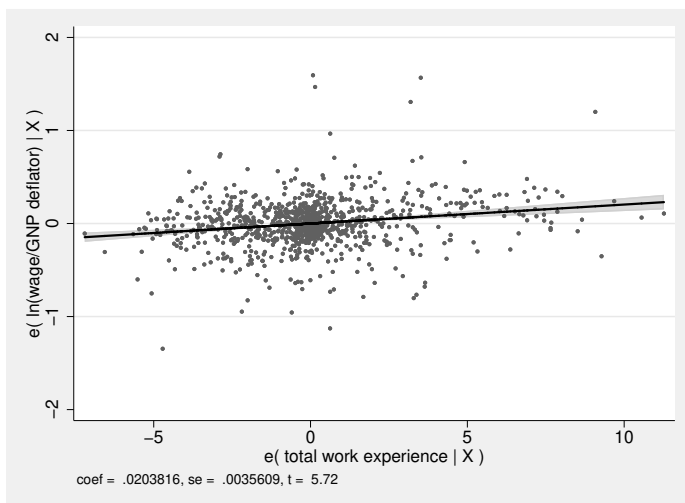
```
. xtavplot ttl_exp, xlim(-6)
```



In this particular case, the source of the problem is the label algorithm, which could be better solved with the option xlabel(-5(5)10), causing no observations to

be omitted, as in the graph below. However, if the value of this outlier had been $-10$, the `xlim()` option would be helpful because the problem could not be solved with an `xlabel()` option. Omitting the value of $-10$ would probably warrant a footnote.
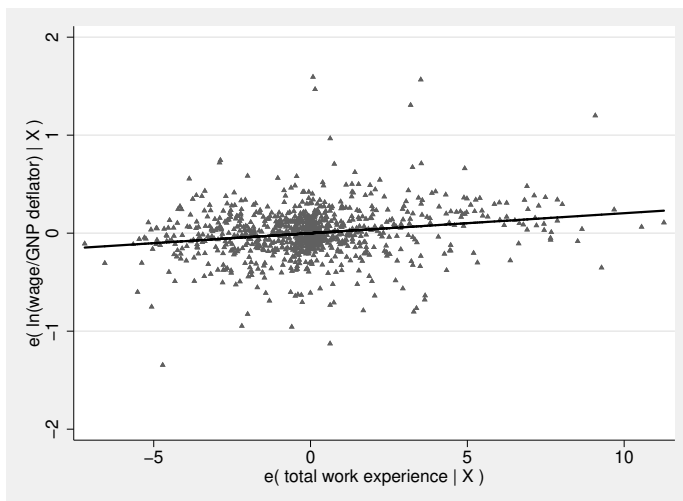
The confidence interval can be displayed as an area plot with the `ciplot(rarea)` option, as displayed in the command `lfitci`. The `ciunder` option causes the confidence interval to appear underneath the scatterplot. By default, the confidence interval would be above the scatter, obscuring some of the data points.

```
. // With solid confidence interval area
. xtavplot ttl_exp, ciplot(rarea) ciunder xlabel(-5(5)10)
```



The graph below changes the scatterplot marker symbol to triangles, does not display a confidence interval around the fitted line, and removes the value of the `ttl_exp` coefficient, standard error, and $t$ statistic from the bottom of the graph.

```
. xtavplot ttl_exp, msymbol(t) noci nocoef xlabel(-5(5)10)
```
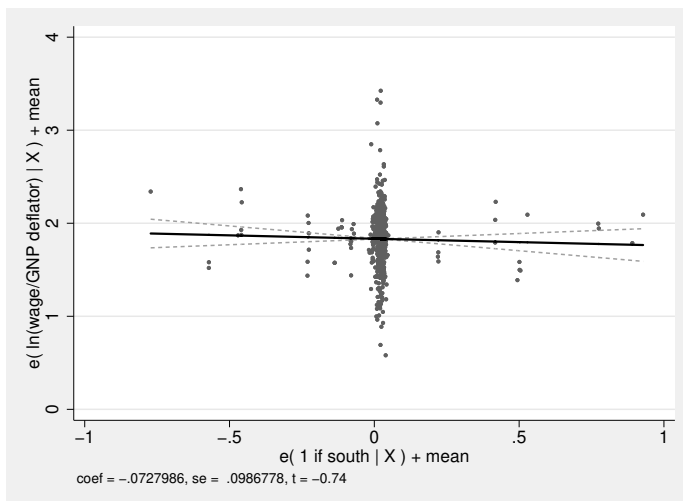


The `addmeans` option rescales the graph to be centered on the actual means of **y** and **x**$_1$ instead of the zero means of the residuals $\mathbf{e_y}$ and $\mathbf{e_{x_1}}$. This may be more intuitive for the reader by conveying the central values of **y** and **x**$_1$. Note that the graph shows the conditional values $\mathbf{e_y}$ and $\mathbf{e_{x_1}}$, not the actual values **y** and **x**$_1$.

The graph below shows the added-variable plot of `south` centered on its mean value of 0.02 and the mean `ln_wage` of 1.83. The mean value of `south`, close to 0, shows that there are few southerners in the sample.

Note that added-variable plots can be an intuitive way of graphing the relationship of dummy variables like `south` to the dependent variable because the values of the residuals $\mathbf{e_{x_1}}$ are continuous even though the unconditional values of `south` are 0 or 1.
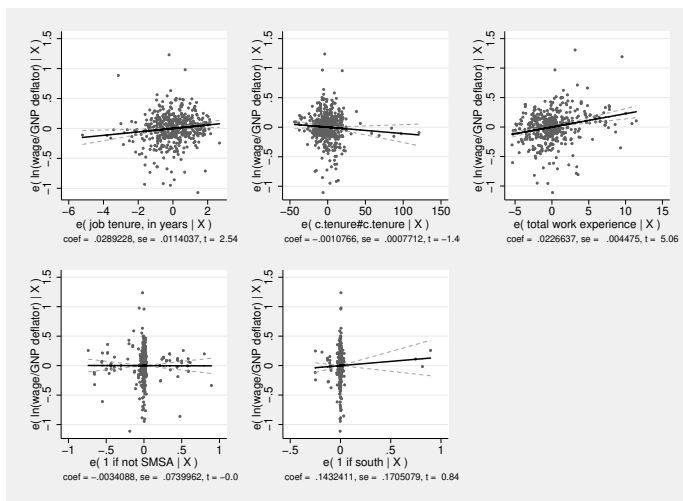
```
. xtavplot south, addmeans
```



coef = −.0727986, se = .0986778, t = −0.74

## 5.1   xtavplots

The command xtavplots with an s on the end creates all possible added-variable plots
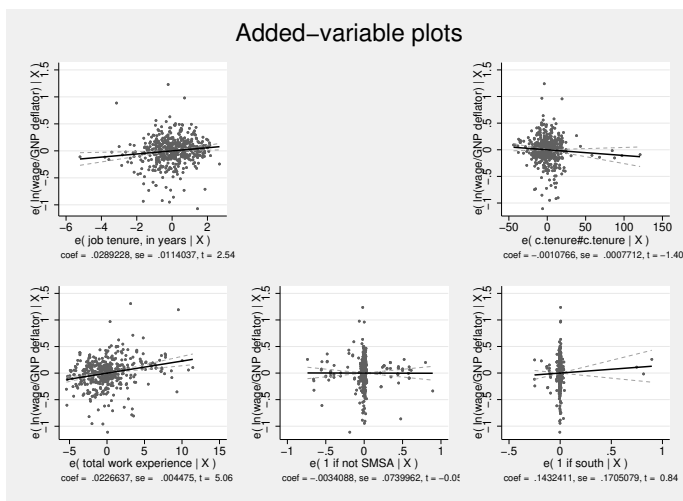of the *indepvars* in a matrix as a single image.

```
. keep in 1/500
(500 observations deleted)
. xtreg ln_w tenure c.tenure#c.tenure ttl_exp not_smsa south, fe
  (output omitted)
```

```
. xtavplots
```



Adding a title and shifting the position of the plots with the `holes()` option make the image look better.

```
. xtavplots, title(Added-variable plots) holes(2)
```



The examples above have focused on graphing options to change the appearance of the graphs created by `xtavplot` after fixed-effects estimation. `xtavplot` can also be employed after between-effects and random-effects estimation. The conceptual issues

involved in creating added-variable plots after these other estimation methods are discussed in previous sections, but the visual considerations when creating these graphs are the same as after fixed-effects estimation.

# 6  Programs and supplemental materials

To install a snapshot of the corresponding software files as they existed at the time of publication of this article, type

```
. net sj 20-1
. net install gr0082      (to install program files, if available)
. net get gr0082          (to install ancillary files, if available)
```

# 7  Reference

Wooldridge, J. M. 2010. *Econometric Analysis of Cross Section and Panel Data.* 2nd ed. Cambridge, MA: MIT Press.

**About the author**

John Luke Gallup is an associate professor in the Department of Economics at Portland State University. He wrote the command `outreg`, among others.