

Lecture 5

Today's agenda – more on OLS regressions

1. goodness of fit measures
2. standard errors of the OLS regression coefficients
3. some examples

$$y_i = X_i \beta + u_i$$

classical case \rightarrow ① $u_i \sim N(0, \sigma^2)$
② more general case

$$\frac{1}{N} \sum [x_i x_i'] \beta$$

$\underbrace{\quad}_{K \times K}$

pop: $E[u_i(y_i - x_i' \beta^*)] = 0$

sample: $\frac{1}{N} \sum [x_i(y_i - x_i' \hat{\beta})] = 0$

We have a regression model for an outcome y with explanatory variables x . Write the population model as:

$$y_i = x_i' \beta^* + u_i.$$

In a sample (size N) the FOC for the OLS estimator $\hat{\beta}$ is:

$$\sum_{i=1}^N x_i(y_i - x_i' \hat{\beta}) = 0 \Rightarrow \frac{1}{N} \sum_{i=1}^N x_i(y_i - x_i' \hat{\beta})$$

assume

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N x_i y_i = \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right) \hat{\beta}$$

$$\Rightarrow \hat{\beta} = \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1} \frac{1}{N} \sum_{i=1}^N x_i y_i$$

$$X'X = \sum_{i=1}^N x_i x_i'$$

$\begin{bmatrix} x_1' & x_2' & \dots & x_N' \end{bmatrix}$

invertible

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1N} \\ \vdots & \vdots & & \vdots \\ x_{21} & x_{22} & \dots & x_{2N} \\ \vdots & \vdots & & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NN} \end{bmatrix}$$

$$y_i \geq c$$

$$E(y_i) \quad \bar{y} = \frac{1}{N} \sum y_i$$

How well does the model fit? A "null" model is one with just a constant. If we fit that model, we know the estimated coefficient will be the mean of y_i . The "sum of squares" of y_i is:

if $\hat{u}_i \neq 0$ more likely

$$SS = \sum_{i=1}^N (y_i - \bar{y})^2$$

$$X_i \hat{\beta} \text{ vs } \bar{y}$$

which is the sum of the squared prediction errors from a model with only a constant. From our OLS model, the corresponding "sum of squared residuals" is

$$SSR = \sum_{i=1}^N (y_i - x_i \hat{\beta})^2 = \sum_{i=1}^N \hat{u}_i^2.$$

We define:

$$R^2 = 1 - (SSR/SS)$$

only constants

$$X_j = \begin{bmatrix} 1 \\ x_{1j} \\ x_{2j} \end{bmatrix}$$

from
FOC

$$\hat{\beta} = \begin{bmatrix} \bar{y} \\ 0 \\ 0 \end{bmatrix}$$

$$\sum_i \frac{x_i}{1} (y_i - \frac{x_i \hat{\beta}}{1}) = 0$$

$$\downarrow$$

$$\frac{1}{N} \sum_i (y_i - \hat{\beta}_1) = 0$$

Then,

$$SSR = SS$$

$$R^2 = 1.$$

$$R^2 = 1 - (SSR/SS)$$

$$= 1 - \frac{\frac{1}{N} \sum_{i=1}^N \hat{u}_i^2}{\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2}$$

The lower bound is 0: in this case $SSR = SS$ and the model is no better than a constant. The upper bound is 1: in this case $SSR = 0$ and the model explains y perfectly.

One problem with R^2 is that it can only get better as you add more x' s – there is no penalty for adding extra variables even if they don't add much. Adjusted R^2 , denoted by \bar{R}^2 , corrects for the use of extra explanatory variables. The idea is to replace the two terms in R^2 with “degrees-of-freedom-corrected” terms:

$$R^2 = 1 - \frac{\frac{1}{N} \sum_{i=1}^N \hat{u}_i^2}{\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2}$$

$$\bar{R}^2 = 1 - \frac{\frac{1}{N-K} \sum_{i=1}^N \hat{u}_i^2}{\frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2}$$

$$\bar{R}^2 = 1 - \frac{\frac{1}{N-K} \sum_{i=1}^N \hat{u}_i^2}{\frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2}$$

The term $\frac{1}{N-K} \sum_{i=1}^N \hat{u}_i^2$ is usually called the “mean squared error” (MSE) of the model. The intuition for dividing by $N - K$ is this: if you have K explanatory variables, and you had a sample of size $N = K$, then you could perfectly predict all the observations, and $SSR = 0$. So the “degrees of freedom” in the SSR is $N - K$. It is also conventional to estimate the variance of the regression residual by the MSE:

$$\widehat{Var}[u_i] = \frac{1}{N - K} \sum_{i=1}^N \hat{u}_i^2$$

Finally, people call \sqrt{MSE} the “standard error of the regression”. It is an estimate of the *standard deviation* of u_i .

Inference

After we estimate the coefficients in a model, we want to know how likely it is that our estimates are “close” to the truth (where “truth” means the population regression coefficients). This will depend on the sample size, how much variability there is in y , and on how easy it is to separate out the effects of the various elements of x . We begin with the “classic” case where the X 's are “fixed” and we have normally distributed residuals.

indep
variance

We assume the population model is:

$$y_i = x_i' \beta^* + u_i.$$

We also assume:

1. we have an independent sample of size N

2. the elements of x_i are not linearly dependent $\begin{bmatrix} X_i & X_i' \end{bmatrix}$

3. given the x_i 's each of the u_i are iid normals: $u_i \sim N(0, \sigma_u^2)$
 \Rightarrow invertible

Under the previous assumptions, the vector of OLS coefficients is distributed normally conditional on the sample x' s:

$$\hat{\beta} - \beta^* \sim \underline{N(0, V)} \quad \text{Variance of matrix}$$

$$V = \frac{\sigma_u^2}{N} S_{xx}^{-1}, \quad \text{and} \quad S_{xx} = \frac{1}{N} \sum_{i=1}^N x_i x_i'.$$

For $\hat{\beta}$ to be precise, N big,
 σ_u^2 small, S_{xx} small [large range
of x' s] 8

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T y & y &= X\beta^* + u \\ &= (X^T X)^{-1} X^T [X\beta^* + u] \\ &= \beta^* + (X^T X)^{-1} X^T u\end{aligned}$$

To prove this we write:

$$\begin{aligned}\hat{\beta} &= \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1} \frac{1}{N} \sum_{i=1}^N x_i y_i & \hat{\beta} - \beta^* &= \left(\frac{1}{N} X^T X \right)^{-1} \frac{1}{N} X^T u \\ &= S_{xx}^{-1} \frac{1}{N} \sum_{i=1}^N x_i (x_i' \beta^* + u_i) & &= \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1} \frac{1}{N} \sum_{i=1}^N x_i u_i \quad \leftarrow \text{rows} \\ &= S_{xx}^{-1} \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right) \beta^* + S_{xx}^{-1} \frac{1}{N} \sum_{i=1}^N x_i u_i \\ &= \beta^* + \sum_{i=1}^N a_i(X) u_i\end{aligned}$$

where $a_i(X) = \frac{1}{N} S_{xx}^{-1} x_i$ is a $K \times 1$ vector that is a function of all the x_i 's (X).

$$\hat{\beta} - \beta^* = \sum_{i=1}^N a_i(X) u_i$$

So conditional on X , each row of $\hat{\beta} - \beta^*$ is just a weighted sum of the u_i 's. Thus, conditional on X , $\hat{\beta} - \beta^*$ is normally distributed! We have:

$$E[\hat{\beta} - \beta^* | X] = E\left[\sum_{i=1}^N a_i(X) u_i | X\right] = \sum_{i=1}^N a_i(X) E[u_i | X] = 0$$

What about the variance of $\hat{\beta} - \beta^*$?

To calculate the variance of $\hat{\beta}$, note:

$$\hat{\beta} - \beta^* = \sum_{i=1}^N a_i u_i$$

There are K rows of this expression:

$$\begin{array}{rcl} \hat{\beta}_1 - \beta_1^* & \xrightarrow{\downarrow} & \sum_{i=1}^N a_{1i} u_i \\ \hat{\beta}_2 - \beta_2^* & = & \sum_{i=1}^N a_{2i} u_i \\ \dots & & \dots \\ \hat{\beta}_K - \beta_K^* & \xrightarrow{\downarrow} & \sum_{i=1}^N a_{Ki} u_i \end{array}$$

$$(\hat{\beta}_j - \beta_j^*) = \sum_i \overbrace{a_{ji}}^{\text{known #'s}} u_i \quad \text{Cov}(\hat{\beta}_j - \beta_j^*, \hat{\beta}_k - \beta_k^*)$$

$$\text{Since } u_i \sim N(0, \sigma_u^2) | X, \quad \sum_i a_{ji} u_i \quad \sum_i a_{ki} u_i$$

So:

* relevant on

heteroskedasticity

$$\text{Var}[\hat{\beta}_j - \beta_j^* | X] = \sigma_u^2 \sum_{i=1}^N a_{ji}^2 \quad \begin{aligned} &\geq E \left[\left(\sum_i a_{ji} u_i \right) \cdot \left(\sum_i a_{ki} u_i \right) \right] \\ &\geq \sigma_u^2 a_{j1} a_{k1} + \sigma_u^2 a_{j2} a_{k2} \end{aligned}$$

outer product

$$\text{Cov}[\hat{\beta}_j - \beta_j^*, \hat{\beta}_k - \beta_k^* | X] = \sigma_u^2 \sum_{i=1}^N a_{ji} a_{ki}$$

since the u_i 's are independent of each other. In matrix form:

$$\text{Var}[\hat{\beta} - \beta^* | X] = \sigma_u^2 \sum_{i=1}^N a_i a_i'$$

Assume 2 coefficients

$$\text{Var} \begin{pmatrix} \hat{\beta}_1 - \beta_1^* \\ \hat{\beta}_2 - \beta_2^* \end{pmatrix} \mid X$$

Thus

it's a matrix

$$\text{Var}[\hat{\beta} - \beta^* \mid X]$$

variance of the individual coefficients with the others

$$\begin{aligned} &= \sigma_u^2 \sum_{i=1}^N a_i a_i' \\ &= \sigma_u^2 \sum_{i=1}^N \left(\frac{1}{N^2} S_{xx}^{-1} x_i x_i' S_{xx}^{-1} \right) \\ &= \frac{1}{N} \sigma_u^2 S_{xx}^{-1} \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right) S_{xx}^{-1} \\ &= \frac{1}{N} \underbrace{\sigma_u^2}_{\text{estimated}} S_{xx}^{-1} S_{xx} S_{xx}^{-1} \end{aligned}$$

$$\begin{pmatrix} \text{Var}(\hat{\beta}_1 - \beta_1^*) & \text{Cov}(\hat{\beta}_1 - \beta_1^*, \hat{\beta}_2 - \beta_2^*) \\ \text{Cov}(\hat{\beta}_1 - \beta_1^*, \hat{\beta}_2 - \beta_2^*) & \text{Var}(\hat{\beta}_2 - \beta_2^*) \end{pmatrix}$$

A very important example: $x'_i = (1, x_{2i})$ – we have a constant and 1 other regressor.

$$\begin{aligned}
 S_{xx} &= \frac{1}{N} \sum_{i=1}^N \begin{pmatrix} 1 \\ x_{2i} \end{pmatrix} (1 \quad x_{2i}) \quad \begin{bmatrix} 1 \\ x_{2i} \end{bmatrix} \\
 &= \frac{1}{N} \sum_{i=1}^N \begin{pmatrix} 1 & x_{2i} \\ x_{2i} & x_{2i}^2 \end{pmatrix} \quad \geq \begin{bmatrix} 1 & x_{2i} \\ x_{2i} & x_{2i}^2 \end{bmatrix} \\
 &= \begin{pmatrix} 1 & \bar{x}_2 \\ \bar{x}_2 & \frac{1}{N} \sum_{i=1}^N x_{2i}^2 \end{pmatrix} \quad \text{outer product}
 \end{aligned}$$

$$S_{xx}^{-1} = \frac{1}{\frac{1}{N} \sum_{i=1}^N x_{2i}^2 - \bar{x}_2^2} \begin{pmatrix} \frac{1}{N} \sum_{i=1}^N x_{2i}^2 & -\bar{x}_2 \\ -\bar{x}_2 & 1 \end{pmatrix}$$

$$S_{xx}^{-1} = \frac{1}{\frac{1}{N} \sum_{i=1}^N x_{2i}^2 - \bar{x}_2^2} \begin{pmatrix} \frac{1}{N} \sum_{i=1}^N x_{2i}^2 & -\bar{x}_2 \\ -\bar{x}_2 & 1 \end{pmatrix}$$

The (2,2) element is:

$$\begin{aligned} (S_{xx}^{-1})_{2,2} &= \frac{1}{\frac{1}{N} \sum_{i=1}^N x_{2i}^2 - \bar{x}_2^2} \\ &= \frac{1}{\frac{1}{N} \sum_{i=1}^N (x_{2i} - \bar{x}_2)^2} \\ &= \frac{1}{\hat{\sigma}^2(x_2)} \end{aligned}$$

So, conditional on X ,

$$\hat{\beta}_2 - \beta_2^* \sim N\left(0, \frac{1}{N} \frac{\sigma_u^2}{\hat{\sigma}^2(x_2)}\right)$$

which is a formula you should recall from 140/141. Notice that we don't know σ_u^2 : we estimate it using the MSE:

$$\hat{\sigma}_u^2 = \frac{1}{N-2} \sum_{i=1}^N \hat{u}_i^2$$

which makes sense because we are assuming each u_i has the same variance.

In a standard regression package, the default standard errors (or default “estimated sampling errors”) of the OLS coefficients are calculated “as if” the classic normal model were true. In particular, the reported variance/covariance matrix of the estimated OLS coefficients is:

$$\frac{1}{N} \hat{\sigma}_u^2 S_{xx}^{-1}$$

$$\hat{\beta} - \beta^* = (X^T X)^{-1} X^T v$$

$$E[X^T v] = 0$$

where $S_{xx} = \frac{1}{N} \sum_i x_i x_i'$. Notice the three terms:

$$\text{Var}(X^T \beta | X) = \left[X^T (I \sigma^2) X \right]$$

$$= (X^T X) \sigma^2$$

$$\text{Var}(A\beta) = A \text{Var}(\beta) A^T$$

- the $1/N$ term reflects the sample size

$$\text{Var}(X^T X)^{-1} (X^T u)$$

- the $\hat{\sigma}_u^2$ term reflects our estimate of the variability in y_i conditional on x_i

$$= (X^T X)^{-1} (X^T X) \sigma^2 (X^T X)^{-1}$$

- the S_{xx}^{-1} term reflects the difficulty of pulling apart the contributions of the different rows of x_i

So we have $\hat{\beta} - \beta^* \sim N(0, V)$, and an estimate of V , $\hat{V} = \frac{1}{N} \hat{\sigma}_u^2 S_{xx}^{-1}$.
How do we conduct inference?

1. To test the null $\beta_k^* = 0$: we form the ratio: $\hat{\beta}_k / \sqrt{\hat{V}_{kk}}$ and compare this to a $N(0, 1)$

If $|\hat{\beta}_k|/\sqrt{V_{kk}} > 1.96$ then we “reject the null” under a 2-sided test at 95% confidence

2. How do we test $b\beta_k^* - c\beta_j^* = 0$? $\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$

Note $Var[b\beta_k^* - c\beta_j^*] = b^2Var[\beta_k^*] + c^2Var[\beta_j^*] - 2bcCov[\beta_k^*, \beta_j^*]$

So we pull out the corresponding elements of \hat{V} , and form the ratio:

$$\frac{b\hat{\beta}_k - c\hat{\beta}_j}{b^2\hat{V}_{kk} + c^2\hat{V}_{jj} - 2bc\hat{V}_{kj}}$$

$\beta_1 \neq 0$

$\beta_2 \neq 0$

$\beta_{n-1} \neq 0$

$\beta_n = 0$

mean

0

There is a potentially serious limitation of the “classic normal” model. In lots of cases, we think that $E[y_i|x_i]$ is not really linear. The classic case says that

$$E[y_i|x_i] = x_i'\beta^*$$

When that is not true, we know that the “true error” for the i th observation is

$$\begin{aligned} u_i &= y_i - E[y_i|x_i] + E[y_i|x_i] - x_i'\beta^* \\ &= \varepsilon_i + v_i \end{aligned}$$

In this case the error at each value of x_i can have a non-zero mean, and possibly a variance that depends on x_i . So we'll have to “tech up” our sampling errors!

Example:

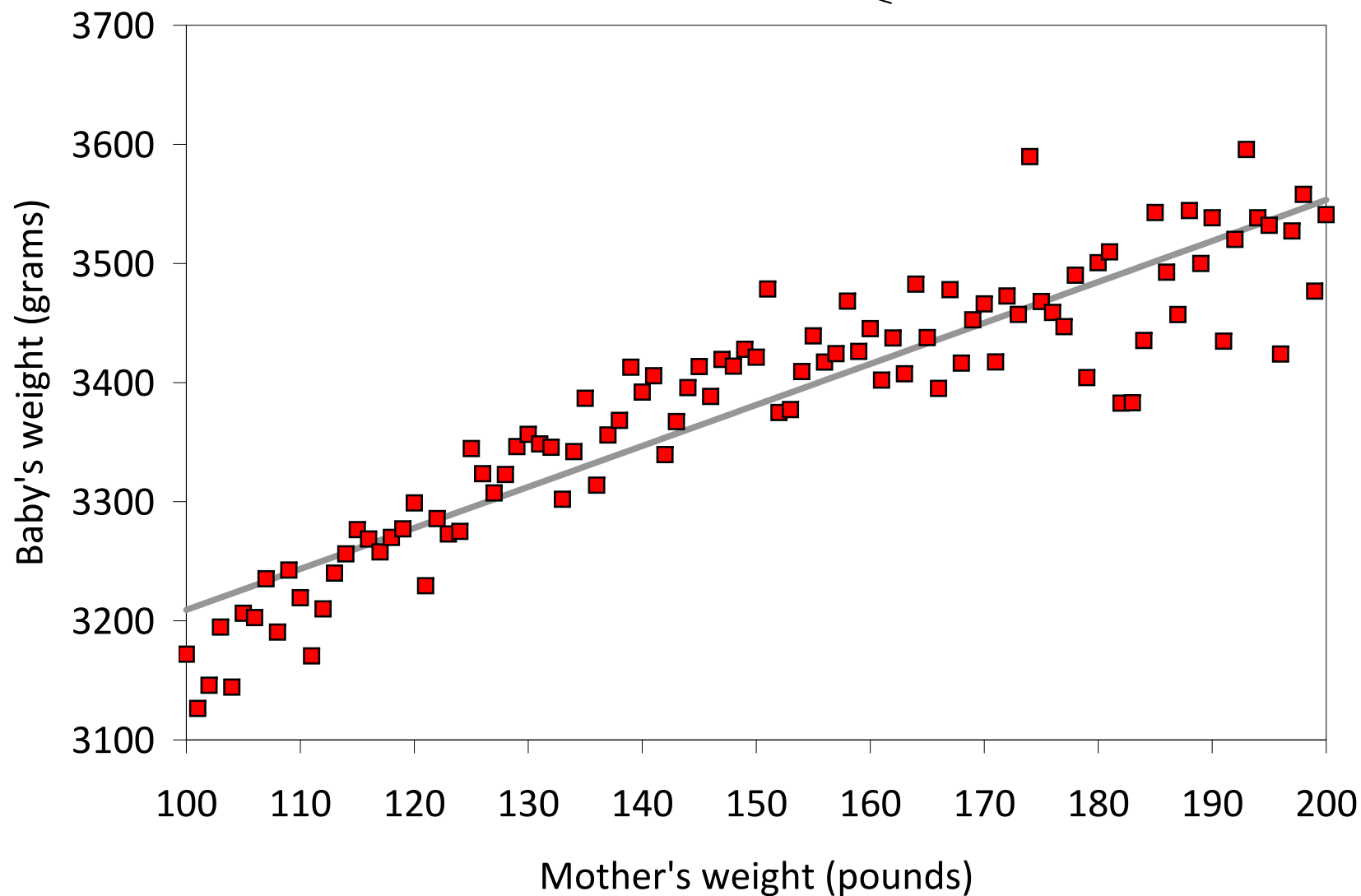
- mom-baby data from problem set #1
- can fit a linear model BUT the R-sq is pretty low - lots of dispersion
- looking more closely we see that there is a "specification error"

$$u_i = y_i - E[y_i|x_i] + E[y_i|x_i] - x_i'\beta^* = \varepsilon_i + v_i$$

- some curvature in simple linear fit
- and heteroskedasticity

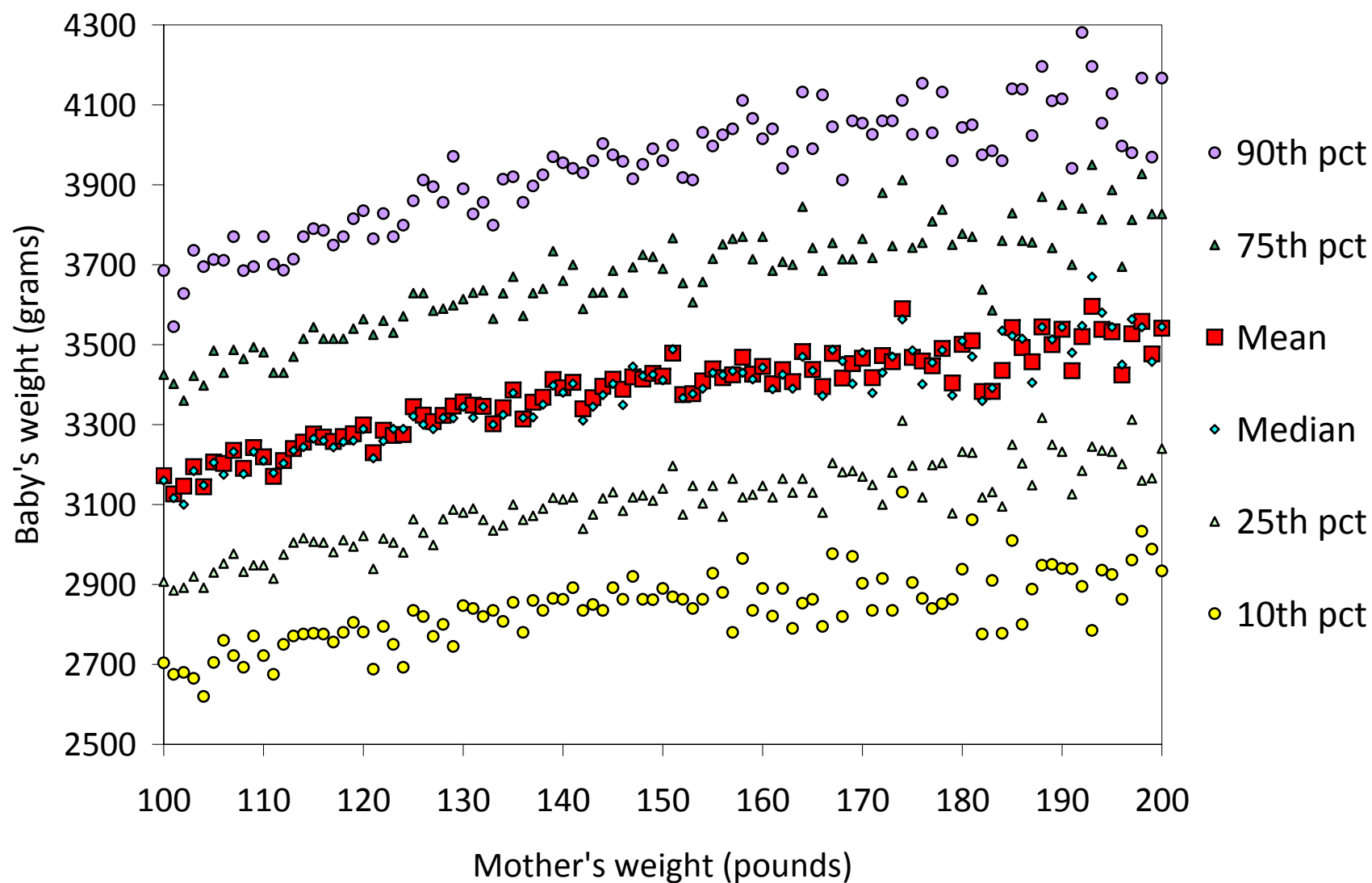
Baby's birth weight versus mother's weight

(very slightly curved)



Baby's birth weight versus mother's weight

$$\begin{aligned} E(u_i|x_i) &\neq 0 \\ \text{Var}(u_i|x_i) &\neq \sigma_u^2 \end{aligned}$$



Variation in baby's birth weight versus mother's weight

