

LECTURE 8:
PREDICTIONS

Jan Zouhar

Introductory Econometrics

Two kinds of predictions

2

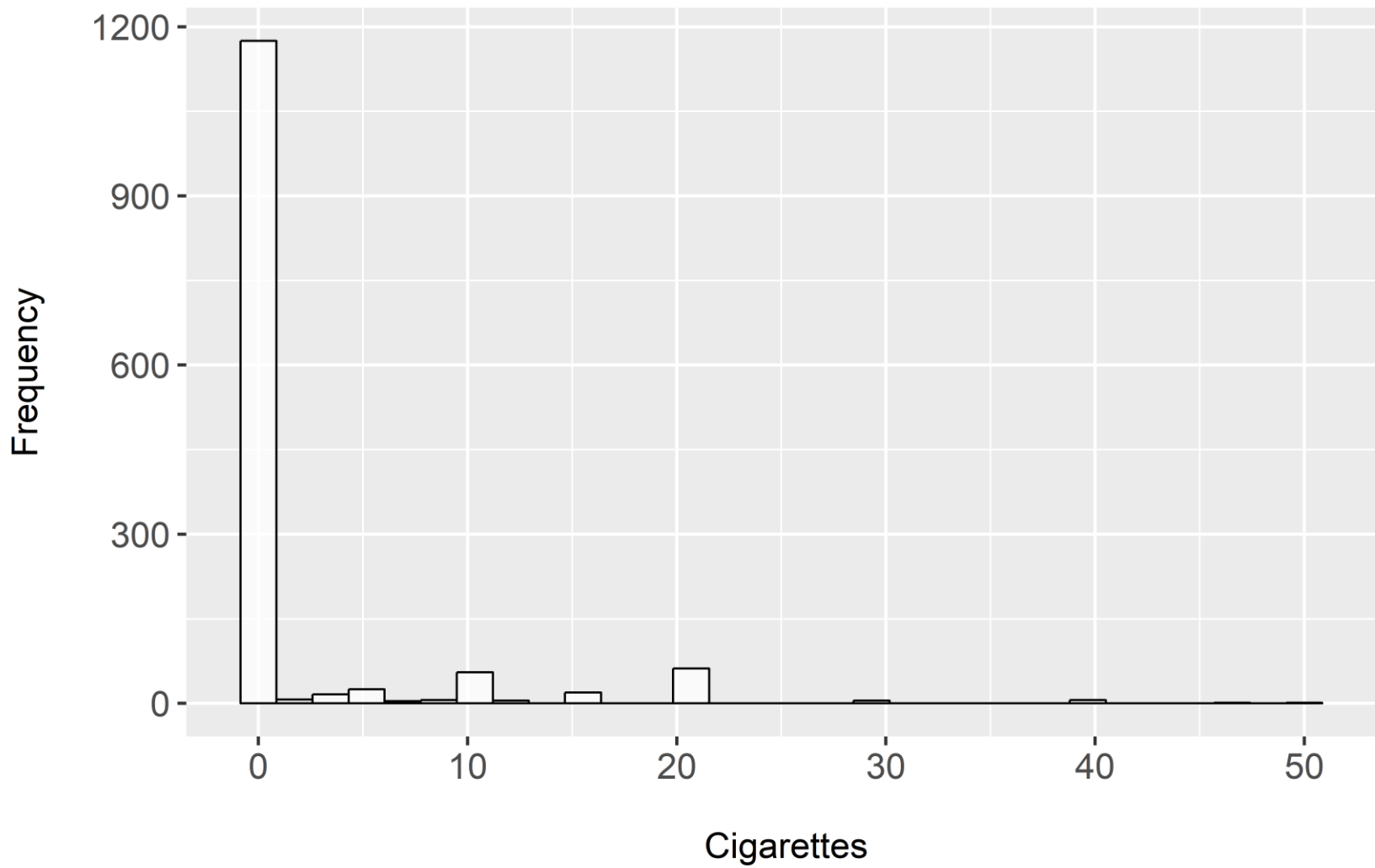
- consider the model

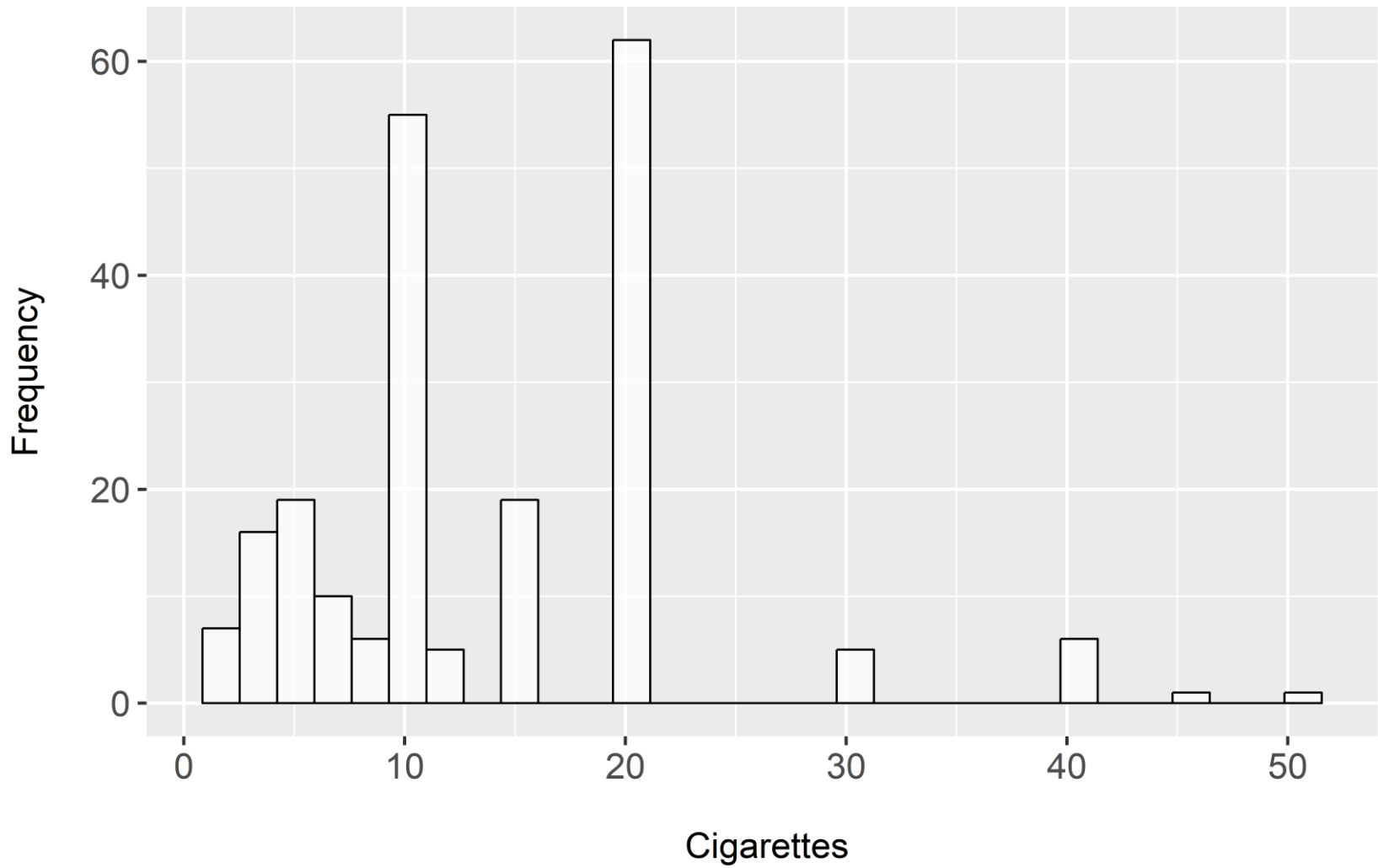
$$bweight = \beta_0 + \beta_1 cigs + u$$

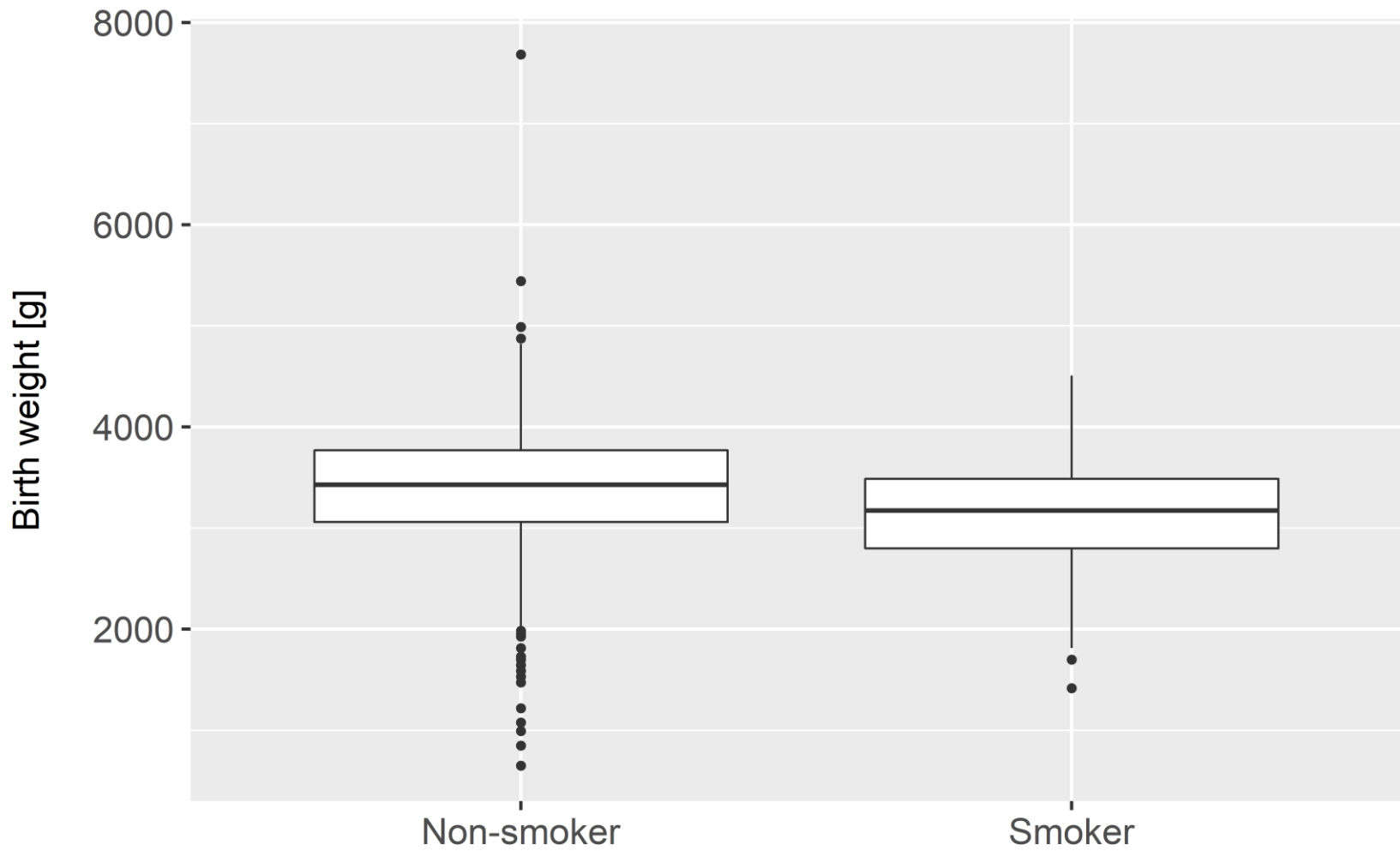
estimated as

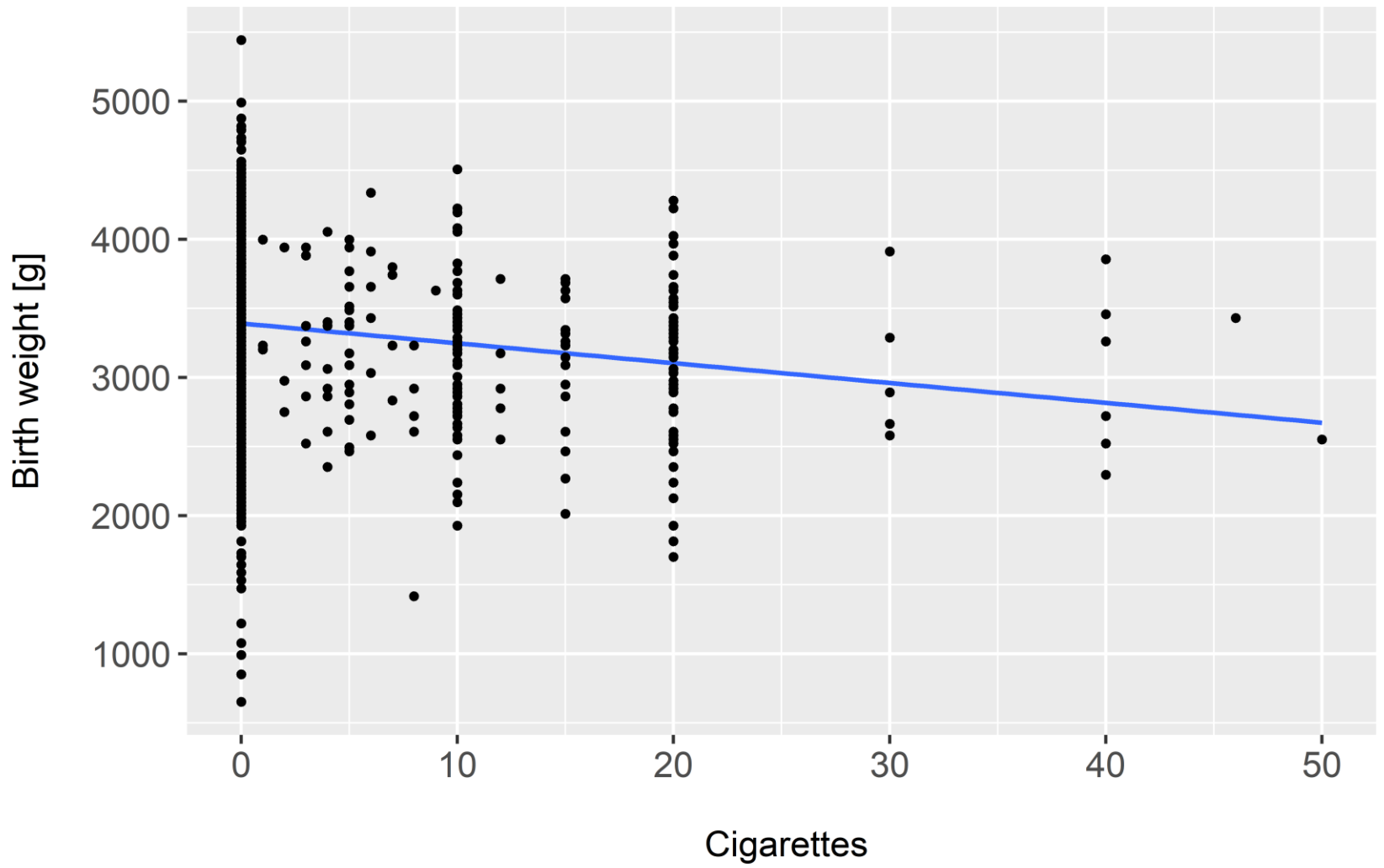
$$\hat{bweight} = 3,395 - 14.57 cigs$$

- assume we know a pregnant woman who smokes 10 cigarettes a day
- there are two kinds of predictions we might be interested in:
 1. predict the actual weight of the woman's baby, denoted $bweight_p$
 2. predict $E[bweight \mid cigs = 10]$, the average birth weight for a mother smoking 10 cigarettes a day, denoted θ
- point prediction is the same in both cases: $3,395 - 14.57 \times 10$, denoted $\hat{\theta}$
- what differs is the 95% CI (or, the standard errors)
 1. 95% CI for the birth weight of a baby of a particular mother (smoking 10 cigarettes a day), typically called the **prediction interval**
 2. 95% CI for a mean in the category of mothers (smoking 10 cigarettes a day), i.e. 95% CI for θ









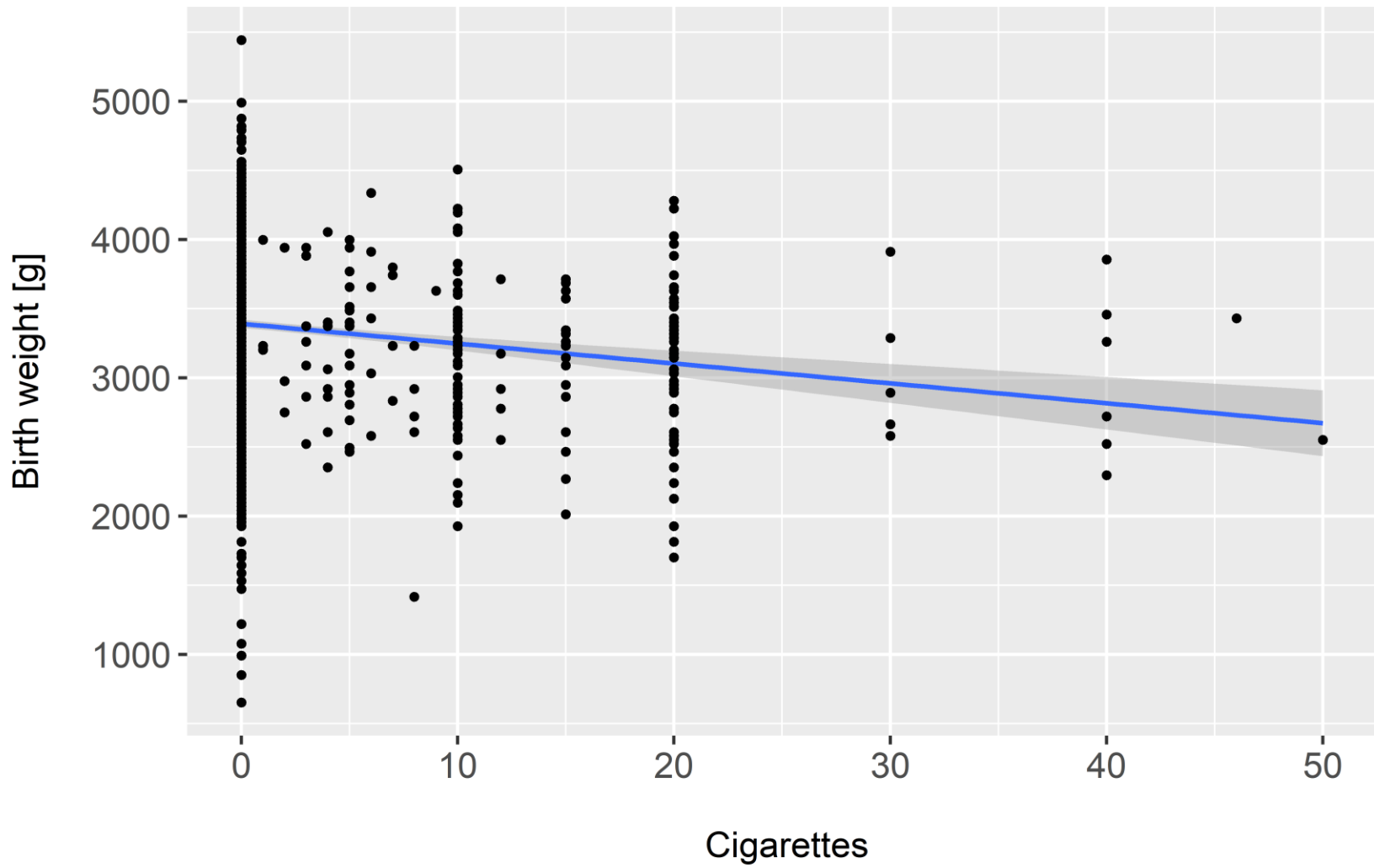
95% CI for ϑ

7

- we are interested in the 95% CI for $\theta = \beta_0 + 10\beta_1$
- it is useful to rewrite the estimated equation as

$$\begin{aligned} bweight &= \beta_0 + \beta_1 cigs + u \\ &= \underbrace{\beta_0 + 10\beta_1}_{\theta} + \beta_1 \underbrace{(cigs - 10)}_A + u \\ &= \theta + \beta_1 A + u \end{aligned}$$

- this gives us a simple procedure to find the 95% for $E[bweight | cigs = 10]$
 1. Create a new variable $A = cigs - 10$.
Gretl: (Add \rightarrow Define new variable $\rightarrow A = cigs - 10$)
 2. Regress $bweight$ on A . The intercept (constant) in this equation is θ , and the 95% CI for the intercept is constructed as usual, i.e. $\hat{\theta} \pm c \cdot se(\hat{\theta})$, where c is either the number 2, or, if more precision is required, the 97.5th percentile of t with $n - k - 1$ degrees of freedom.
Gretl: (Analysis \rightarrow Confidence intervals for coefficients)



Prediction intervals

9

- predicted value: $bweight^P = \beta_0 + \beta_1(10) + u = \theta + u$
- point prediction: $\hat{\theta}$
 - this makes sense as $E u = 0$
- **prediction error:** $\hat{e}^P = bweight^P - bweight_P (= \theta + u - \hat{\theta})$
- two sources of the predictions error:
 1. the population regression function is not estimated precisely; simply put, $\hat{\theta} \neq \theta$
 2. random variation around the mean: u
- we have: $\text{var}(\hat{e}^P) = \text{var}(\theta + u - \hat{\theta}) = \text{var} \hat{\theta} + \text{var} u = \text{var} \hat{\theta} + \sigma^2$
- therefore, a natural estimator of the standard deviation of prediction error is
$$\text{se}(\hat{e}^P) = \sqrt{[\text{se}(\hat{\theta})]^2 + \hat{\sigma}^2}$$
- and the 95% CI is $\hat{\theta} \pm c \cdot \text{se}(\hat{e}^P)$

Prediction intervals

10

- obtaining the prediction interval for a pregnant woman who smokes 10 cigarettes a day:

1. Create a new variable $A = \text{cigs} - 10$.

Gretl: (Add → Define new variable → $A = \text{cigs} - 10$)

2. Regress $bweight$ on A . The intercept (constant) in this equation is θ , its std. error is $se(\hat{\theta})$. The regression output will probably contain either $\hat{\sigma}^2$ or $\hat{\sigma}$. (In Gretl, $\hat{\sigma}$ is called S.E. of regression).

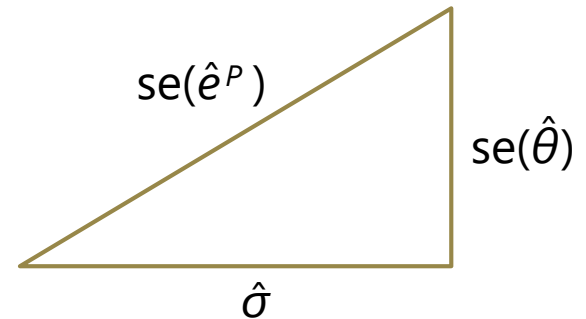
3. Calculate $se(\hat{e}^P) = \sqrt{[se(\hat{\theta})]^2 + \hat{\sigma}^2}$.

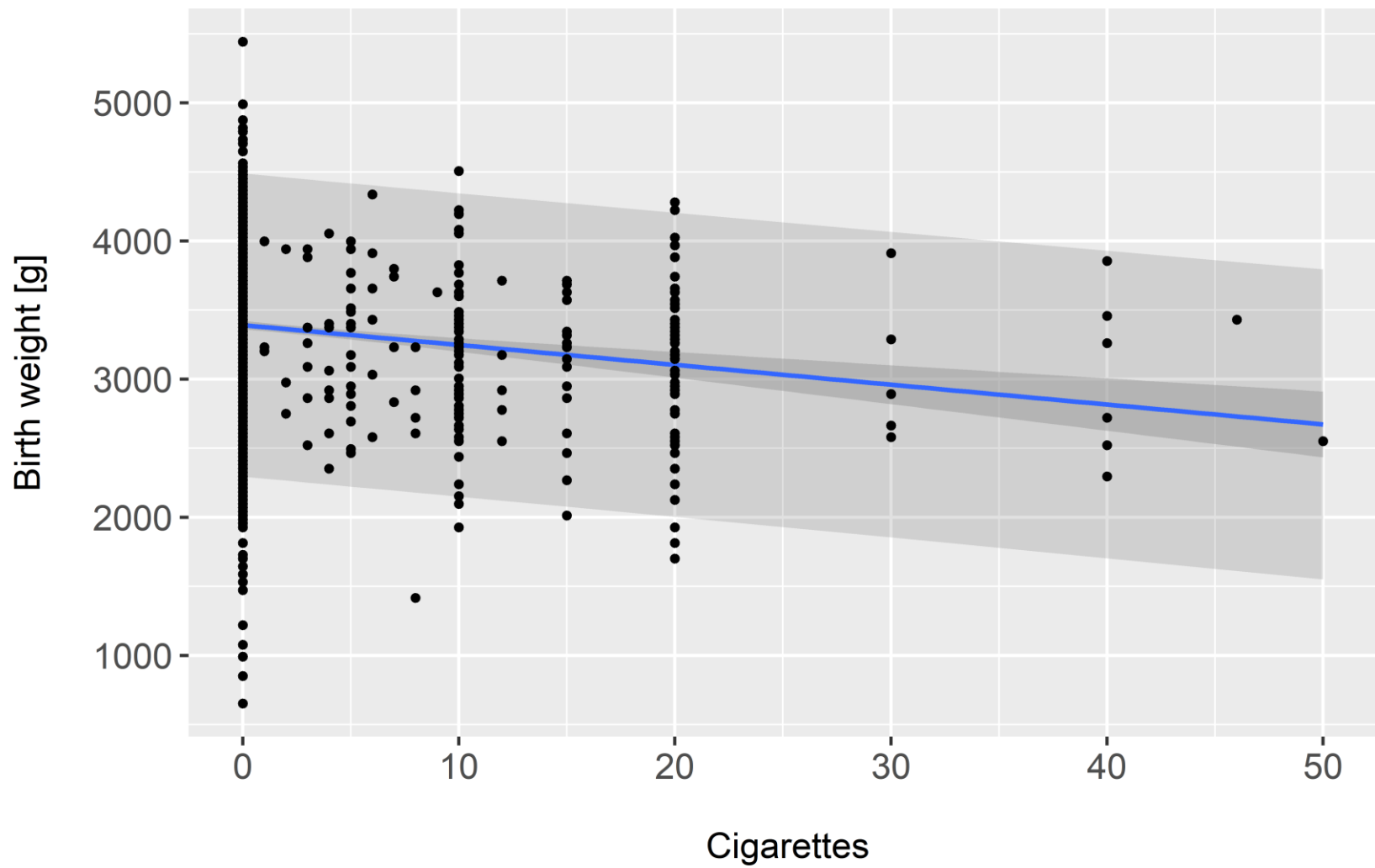
4. Calculate the 95% prediction interval as $\hat{\theta} \pm c \cdot se(\hat{e}^P)$.

- with multiple regression models the procedure is analogous:

- all explanatory variables need to be specified, say $x_1 = c_1, \dots, x_k = c_k$

- we regress y on $(x_1 - c_1), \dots, (x_k - c_k)$ in step 2, the rest is the same





Predicting y when $\log(y)$ is the dependent variable

12

- consider the model

$$\log y = \beta_0 + \beta_1 x + u \quad (1)$$

- how do we predict, say, $E[y | x = 10]$?
- the model implies that $y = e^{\beta_0 + \beta_1 x + u}$
- therefore:

$$\begin{aligned} E[y | x] &= E[e^{\beta_0 + \beta_1 x + u} | x] \\ &= \underbrace{e^{\beta_0 + \beta_1 x}}_A \cdot \underbrace{E[e^u | x]}_B \end{aligned}$$

- A is consistently estimated as $e^{\hat{\beta}_0 + \hat{\beta}_1 x}$, where $\hat{\beta}_0$ and $\hat{\beta}_1$ are parameter estimates from (1) and x is replaced with the specified value
- B is more tricky; here we'll discuss two options:
 1. If we assume that u is normally distributed, then $B = \exp(\sigma^2 / 2)$
 2. Duan's (1983) estimator: estimate B as the sample mean of exponentiated residuals from (1), i.e. $B = \frac{1}{n} \sum_{i=1}^n \exp(\hat{u}_i)$

$$\hat{B} = \frac{1}{n} \sum_{i=1}^n \exp(\hat{u}_i) = 1.0807$$

The screenshot shows the gretl software interface for 'model 2'. The 'Save' menu is open, with 'Fitted values' selected. The regression summary table is visible, showing various statistics. The 'S.E. of regression' value, 0.396542, is circled in red. Below the screenshot, an arrow points from this circled value to the equation $\hat{B} = \exp(0.396542^2 / 2) = 1.0818$.

	Ratio	p-value
const	0.94	0.0002 ***
educ	0.07	2.60e-025 ***
exper	0.59	2.79e-012 ***
expersq	0.56	2.48e-010 ***
tenure	0.68	7.08e-08 ***
nonwhite	0.483	0.6542
female	0.56	1.33e-016 ***
smsa	0.57	8.64e-05 ***

Mean dependent var	0.531538
Sum squared resid	81.45340
R-squared	0.450863
F(7, 518)	60.75682
Log-likelihood	-255.7956
Schwarz criterion	561.7137
S.E. of regression	0.396542
Adjusted R-squared	0.433442
P-value (F)	1.72e-63
Akaike criterion	527.5912
Hannan-Quinn	540.9517

Log-likelihood for wage = -1109.63

Excluding the constant, p-value was highest for variable 5 (nonwhite)

$$\hat{B} = \exp(0.396542^2 / 2) = 1.0818$$

LECTURE 8:
PREDICTIONS

Jan Zouhar

Introductory Econometrics