

MACHINE LEARNING

PART 2

Lecture

SOURCES OF ERROR

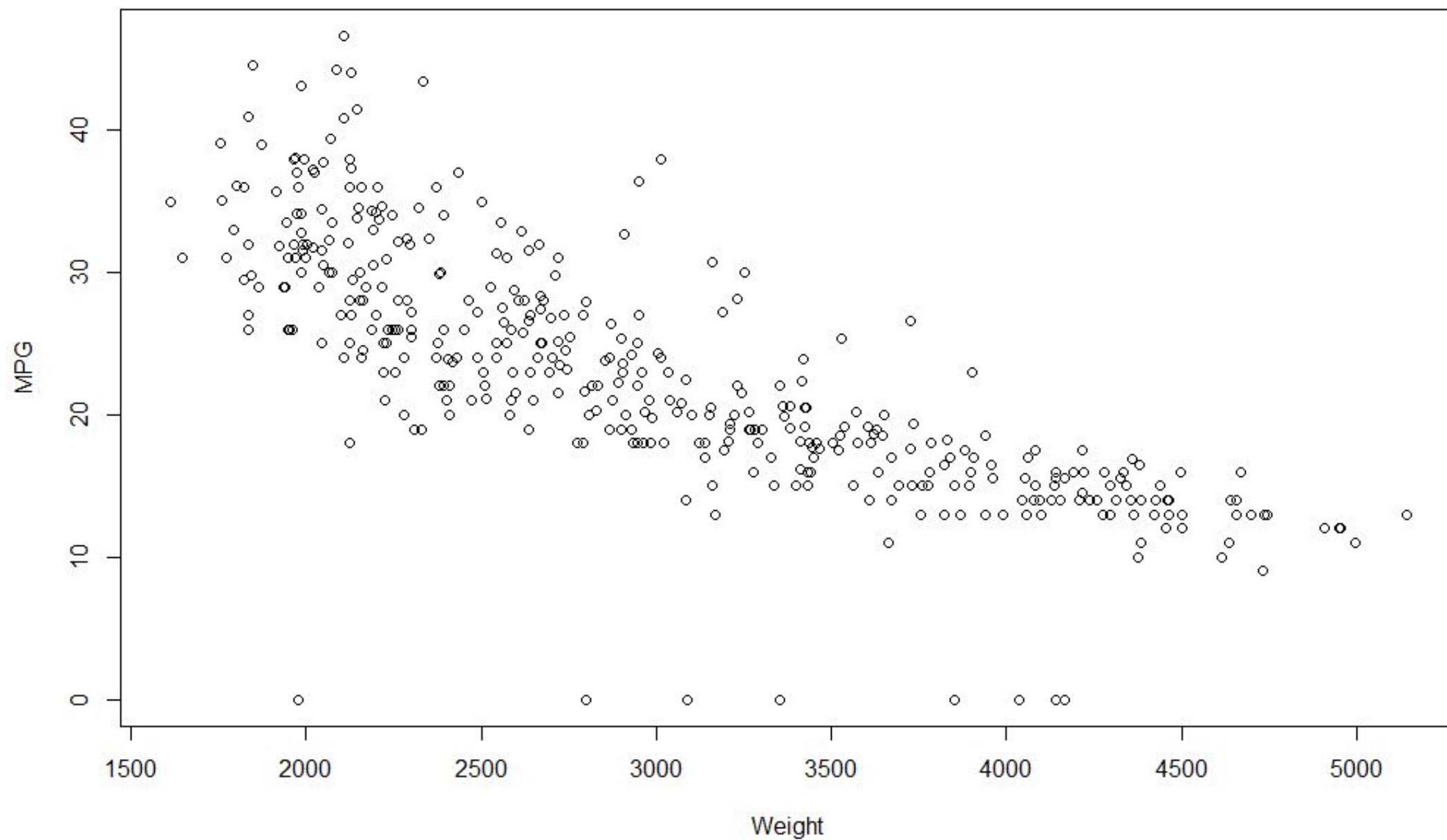
By

Gourab Nath

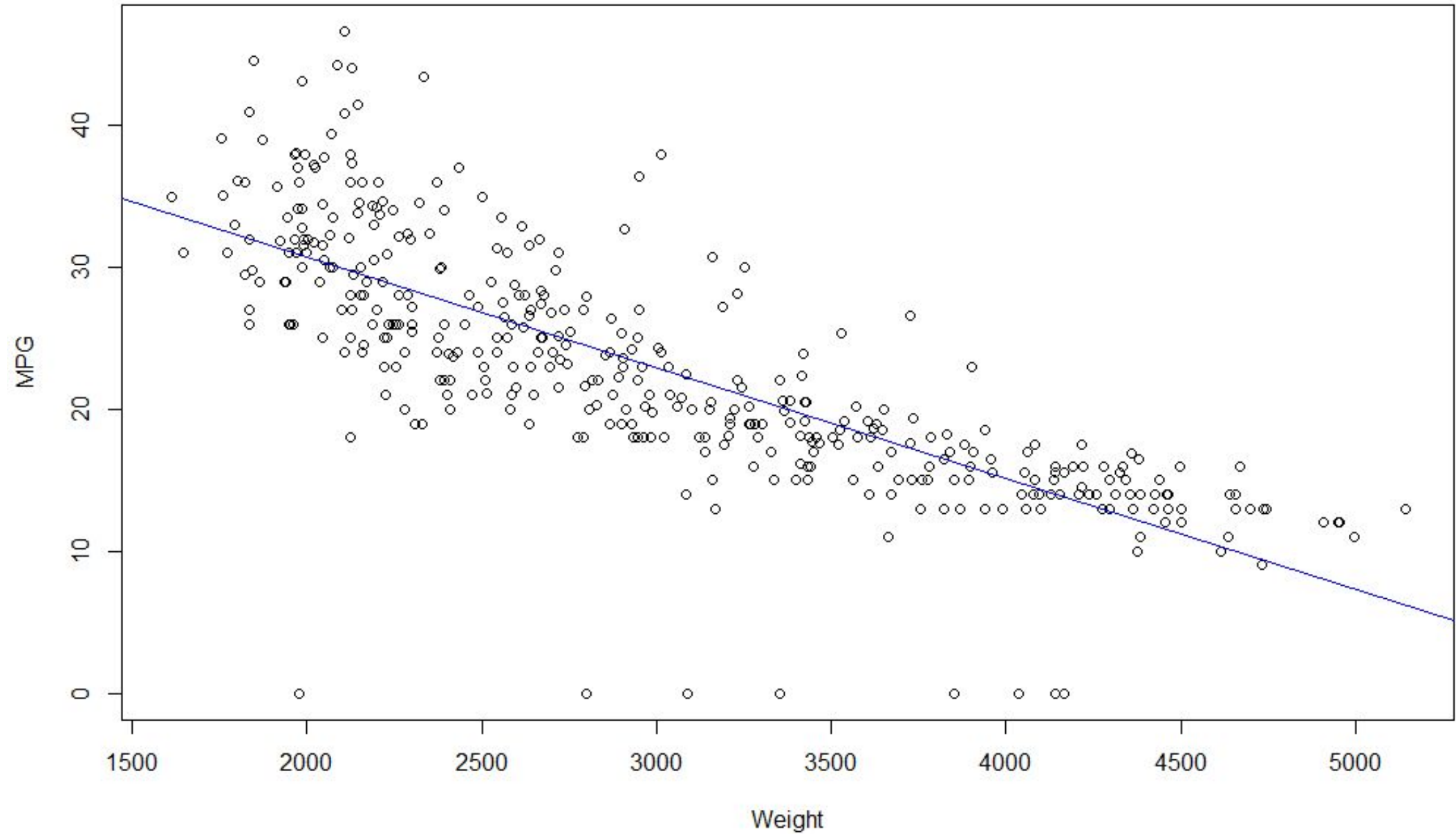
Assistant Professor of Data Science

Praxis Business School

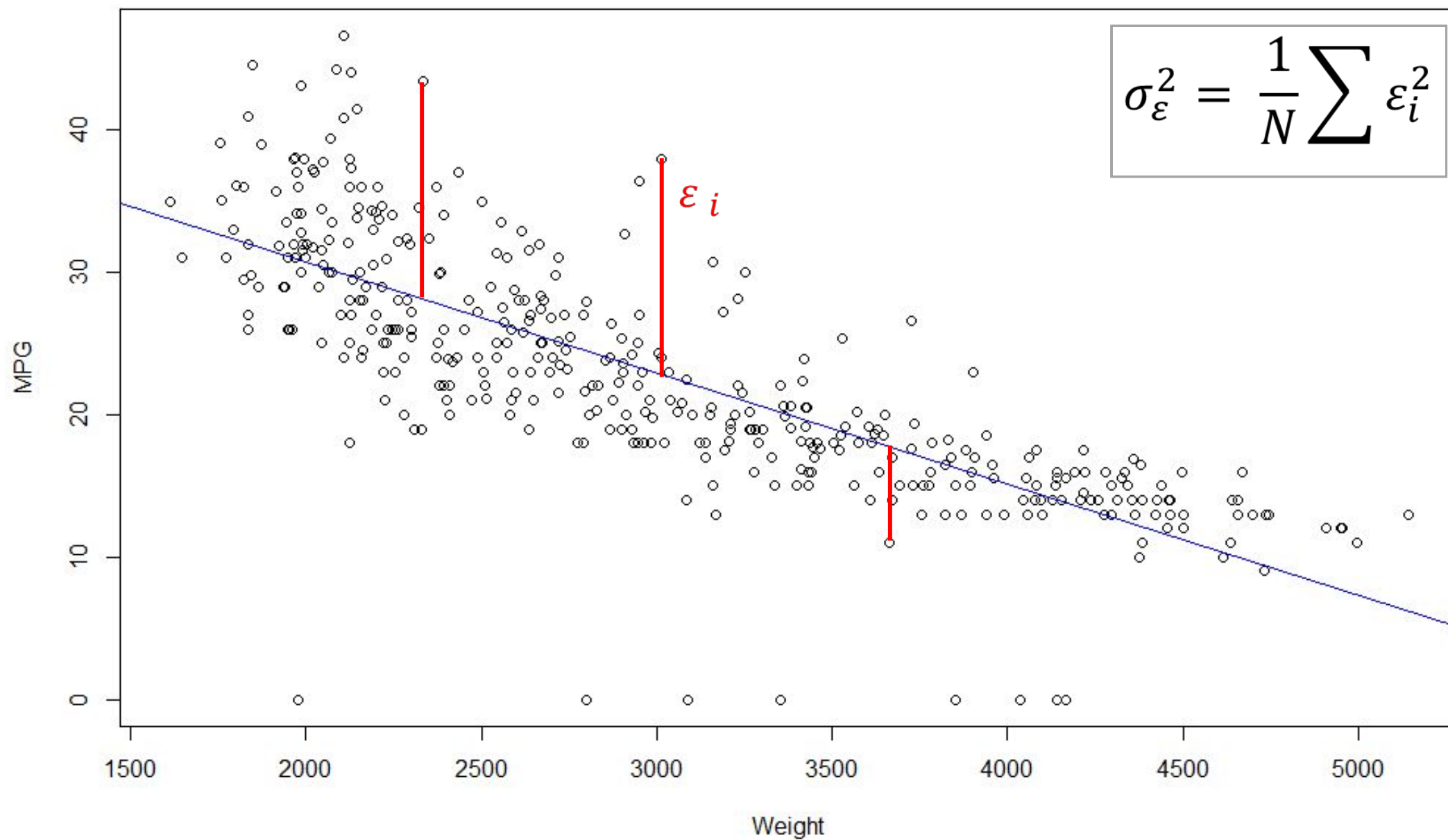
POPULATION DATA



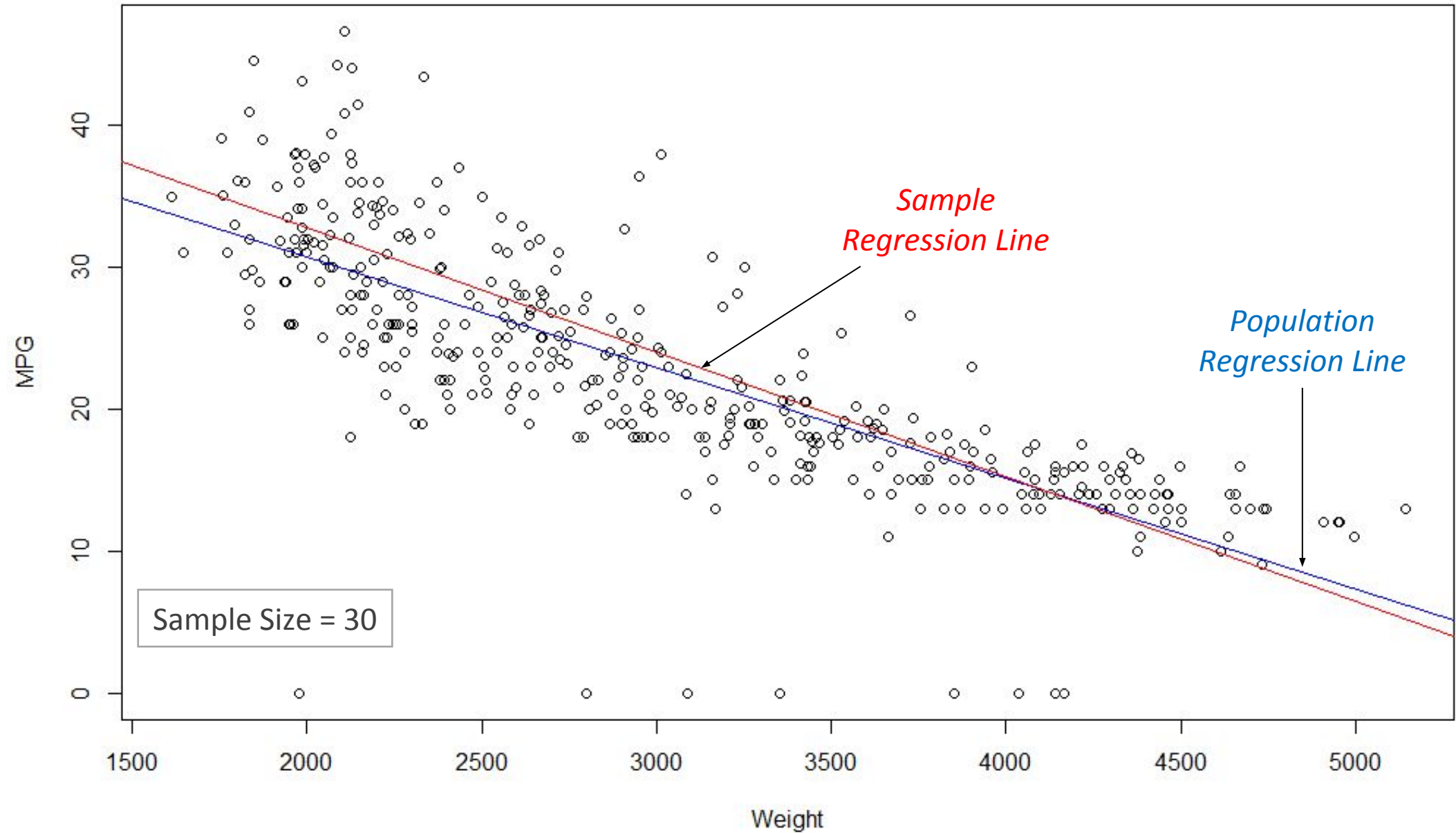
POPULATION REGRESSION LINE



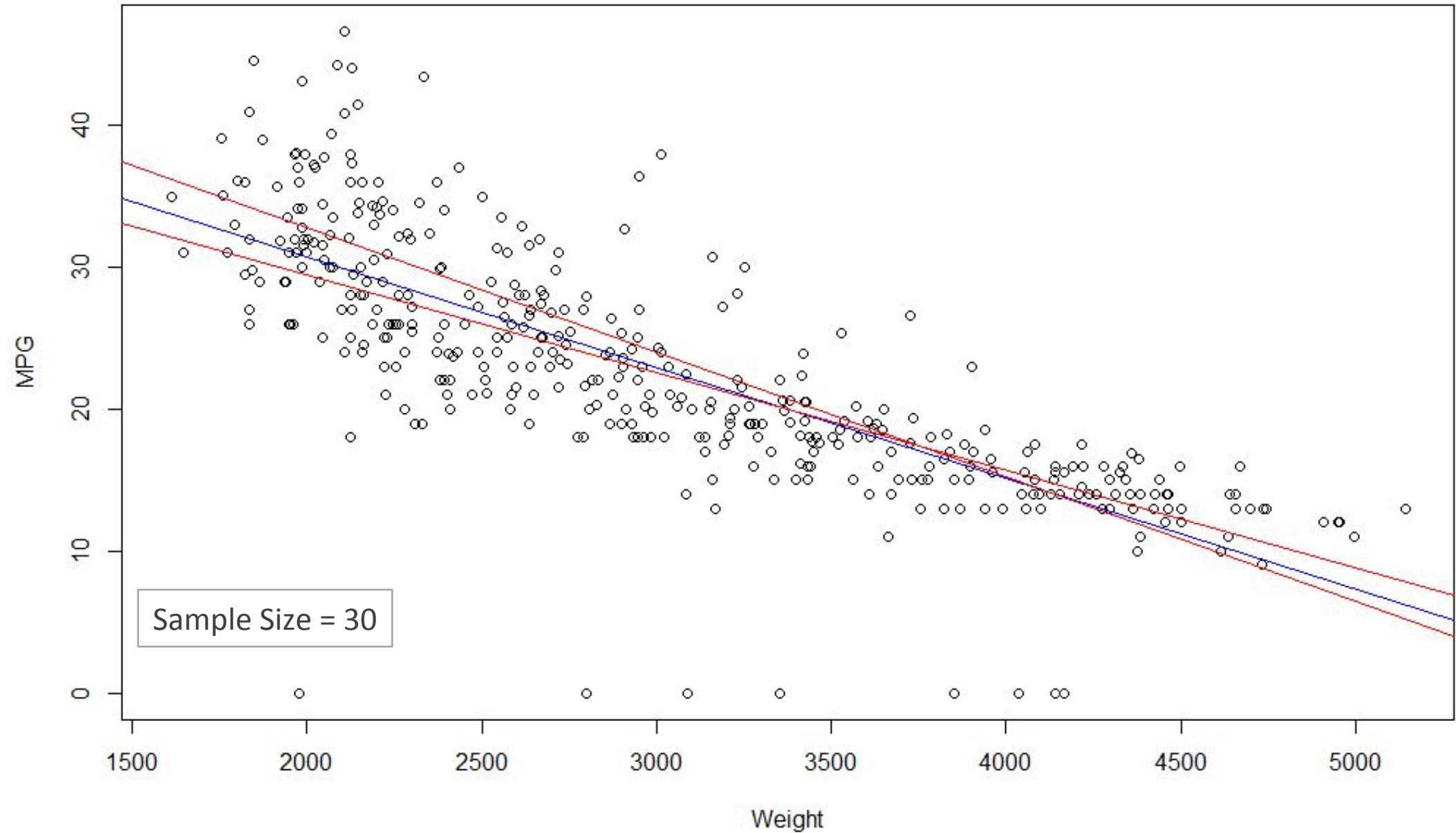
IRREDUCIBLE ERROR



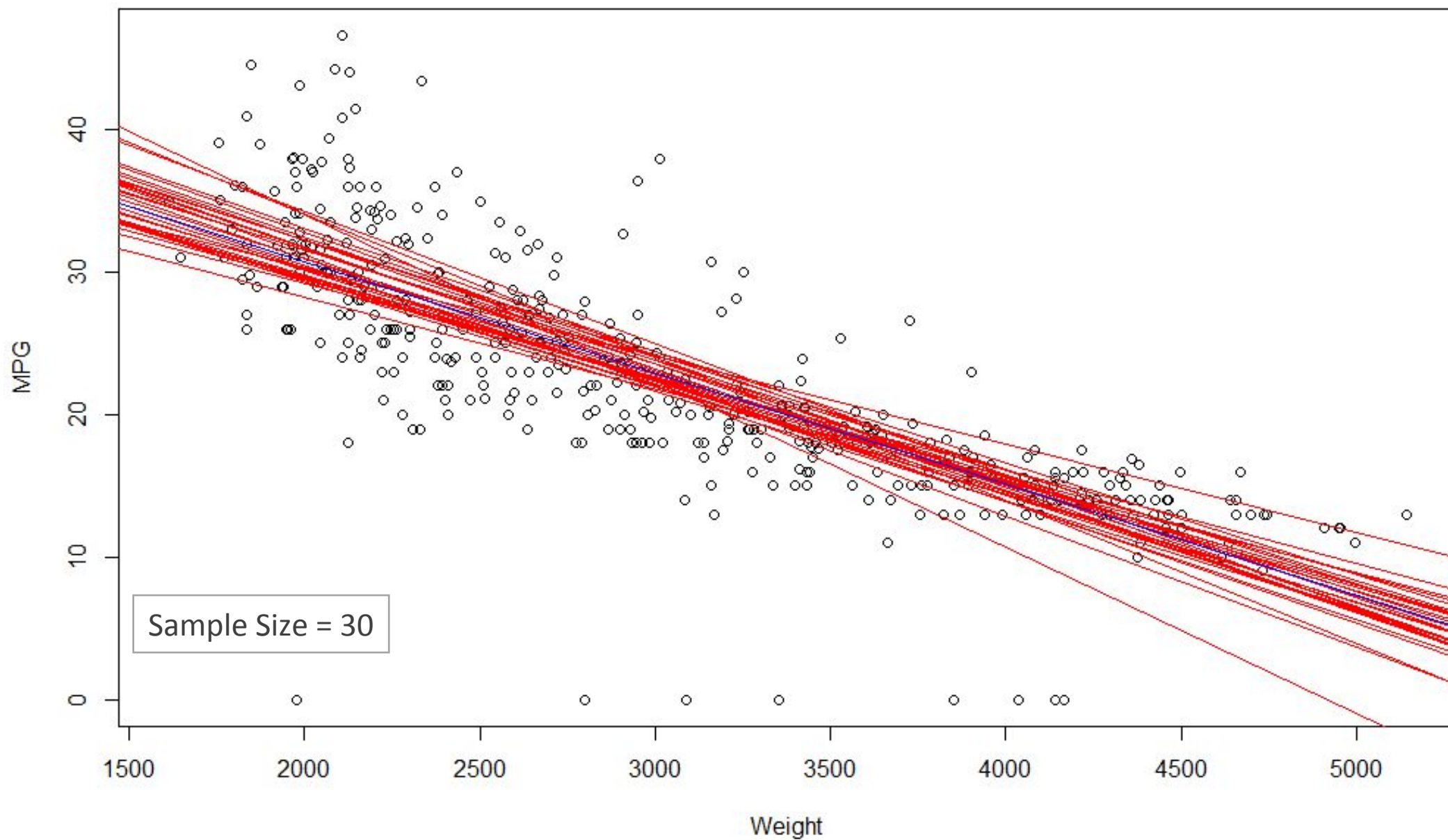
REGRESSION LINE ON 1 SAMPLE



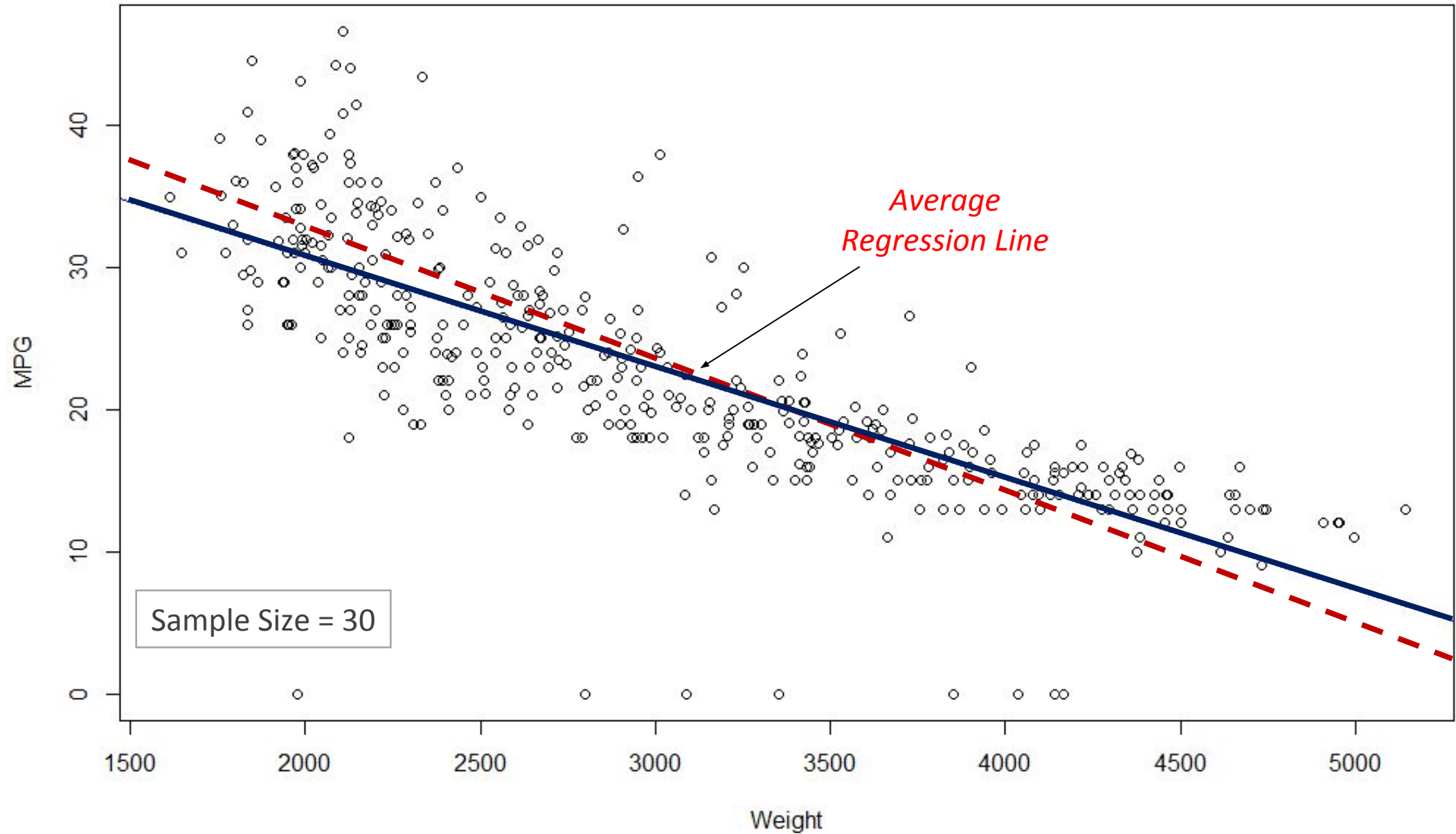
REGRESSION LINES FITTED ON 2 DIFFERENT SAMPLE



REGRESSION LINES FITTED ON 30 DIFFERENT SAMPLE



THE AVERAGE REGRESSION LINE



NOTATIONS

Let \mathbf{x} be an observation.

$f(\mathbf{x}) =$ True regression line. That is, the regression line obtained by fitting on the population data.

$\widehat{f_{sample}}(\mathbf{x}) =$ The fitted sample regression line. That is, the regression line obtained by fitting on a sample.

$\bar{f}(\mathbf{x}) =$ The average of all the regression line fitted over all possible sample.

OUTPUTS OF THE FUNCTIONS

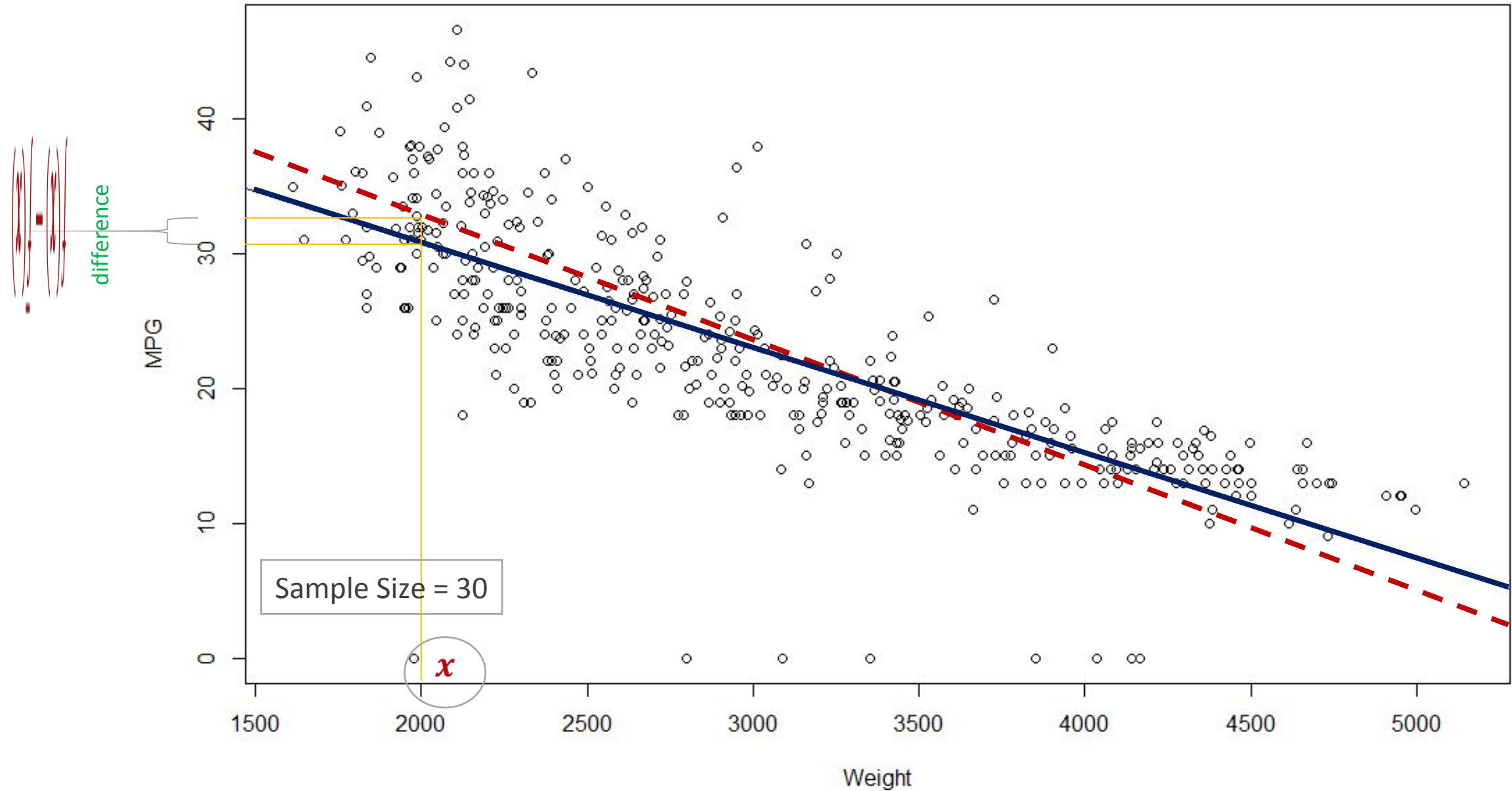
Let \mathbf{x} curl be an observation.

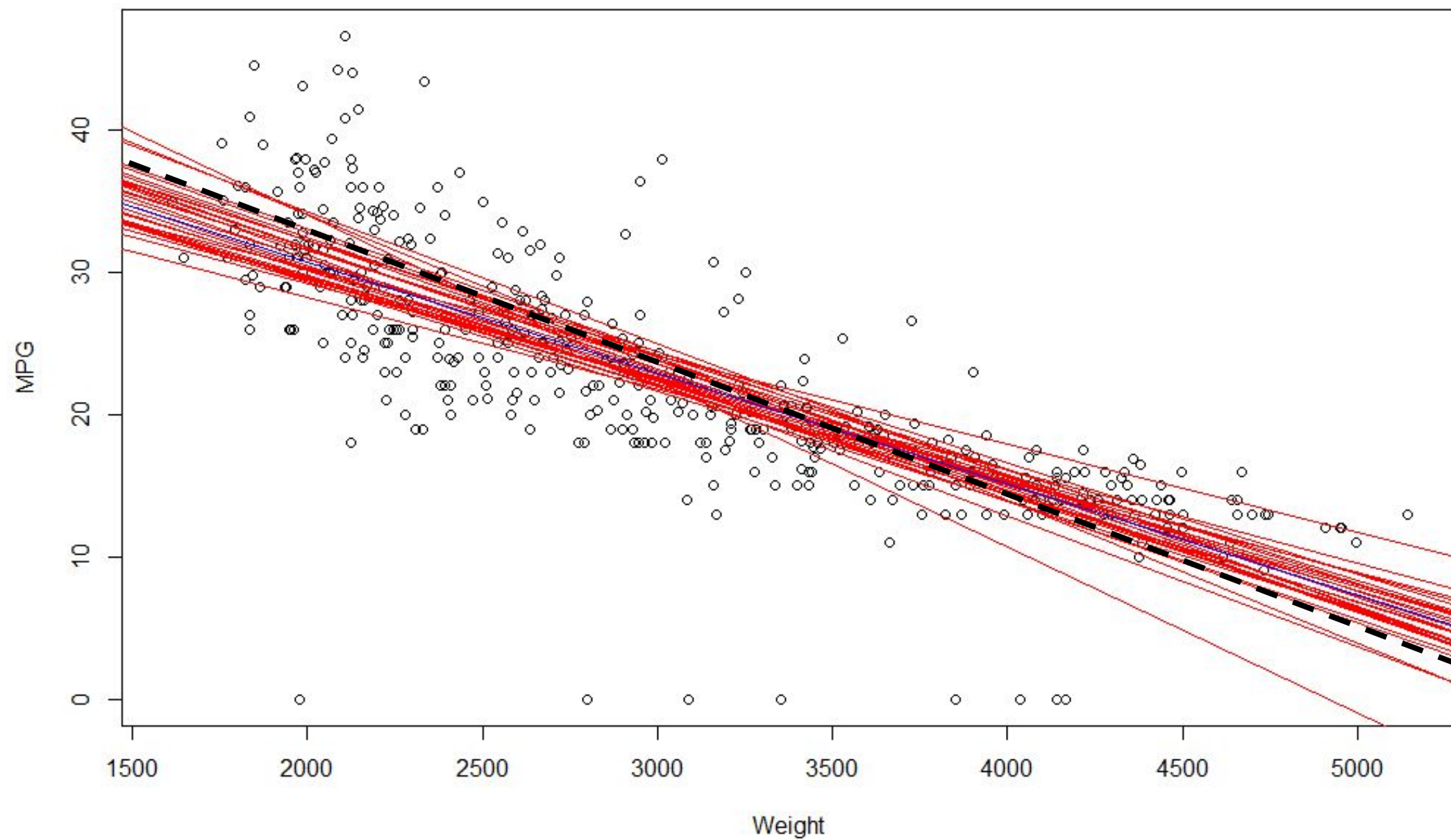
$f(\mathbf{x}) =$ Returns the predicted value of \mathbf{x} using the fitted population regression model

$\widehat{f_{sample}}(\mathbf{x}) =$ Returns the predicted value of \mathbf{x} using the estimated regression model based on the sample.

$\bar{f}(\mathbf{x}) =$ Returns the average of all the predicted value of \mathbf{x} predicted using all possible regression models fitted over all possible samples.

THE DIFFERENCE BETWEEN ACTUAL AND PREDICTED





BIAS AND VARIANCE

$$Bias(\mathbf{x}) = f(\mathbf{x}) - \bar{f}(\mathbf{x})$$

How far the average of all the models is from the true model?

$$Variance(\mathbf{x}) = E_{sample} [\widehat{f_{sample}}(\mathbf{x}) - \bar{f}(\mathbf{x})]^2$$

How much does the outputs of each models (fitted over different samples of same size) varies from the average of all the models?

MSE at \mathbf{x}

$$MSE = Irreducible\ Error + bias^2(\mathbf{x}) + variance(\mathbf{x})$$

$$\bar{x} = \frac{1}{n} \sum x_i$$

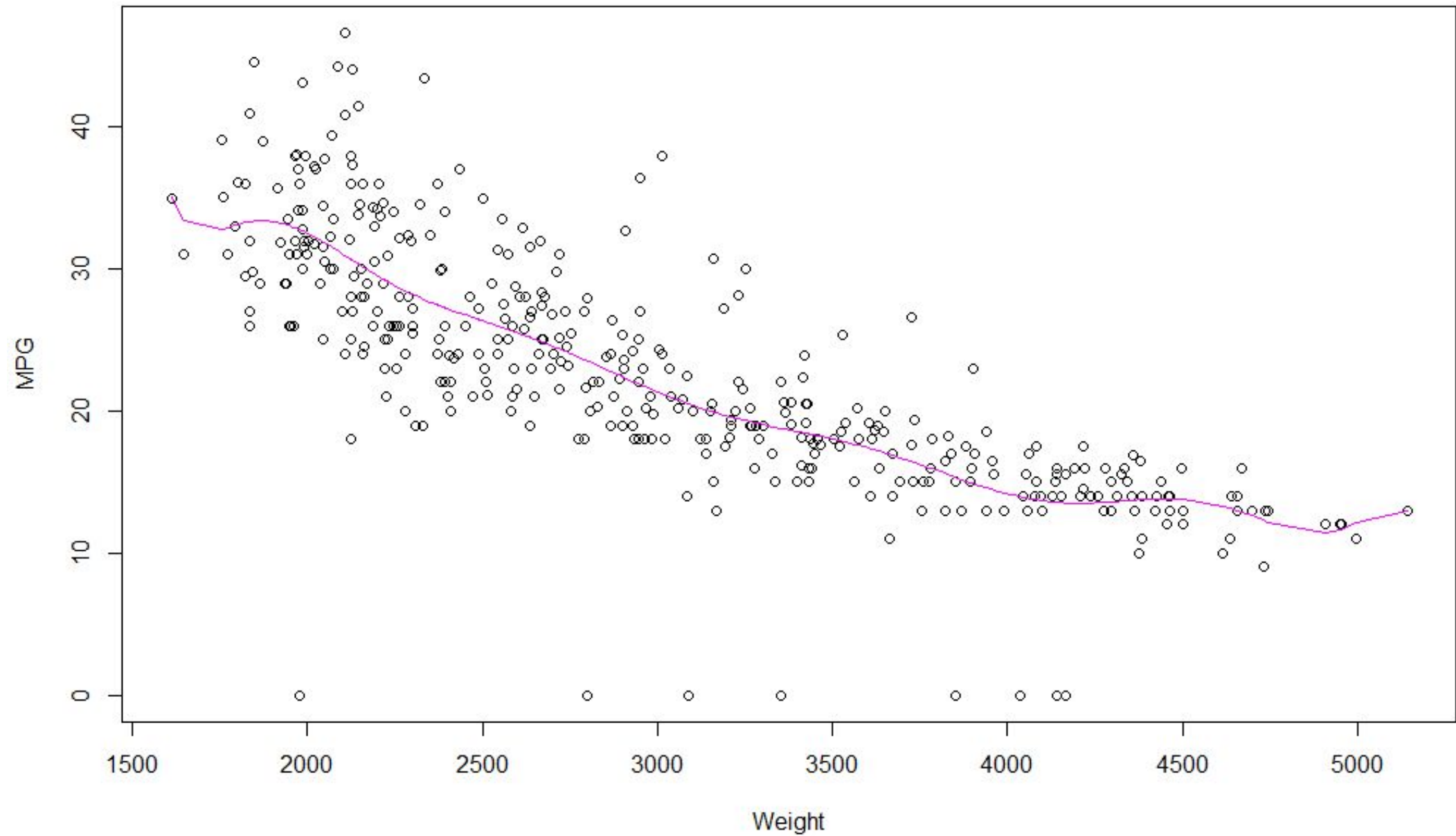
$$variance = \frac{1}{n} \sum (x_i - \bar{x})^2$$

$$E(X) = \sum x_i p_i$$

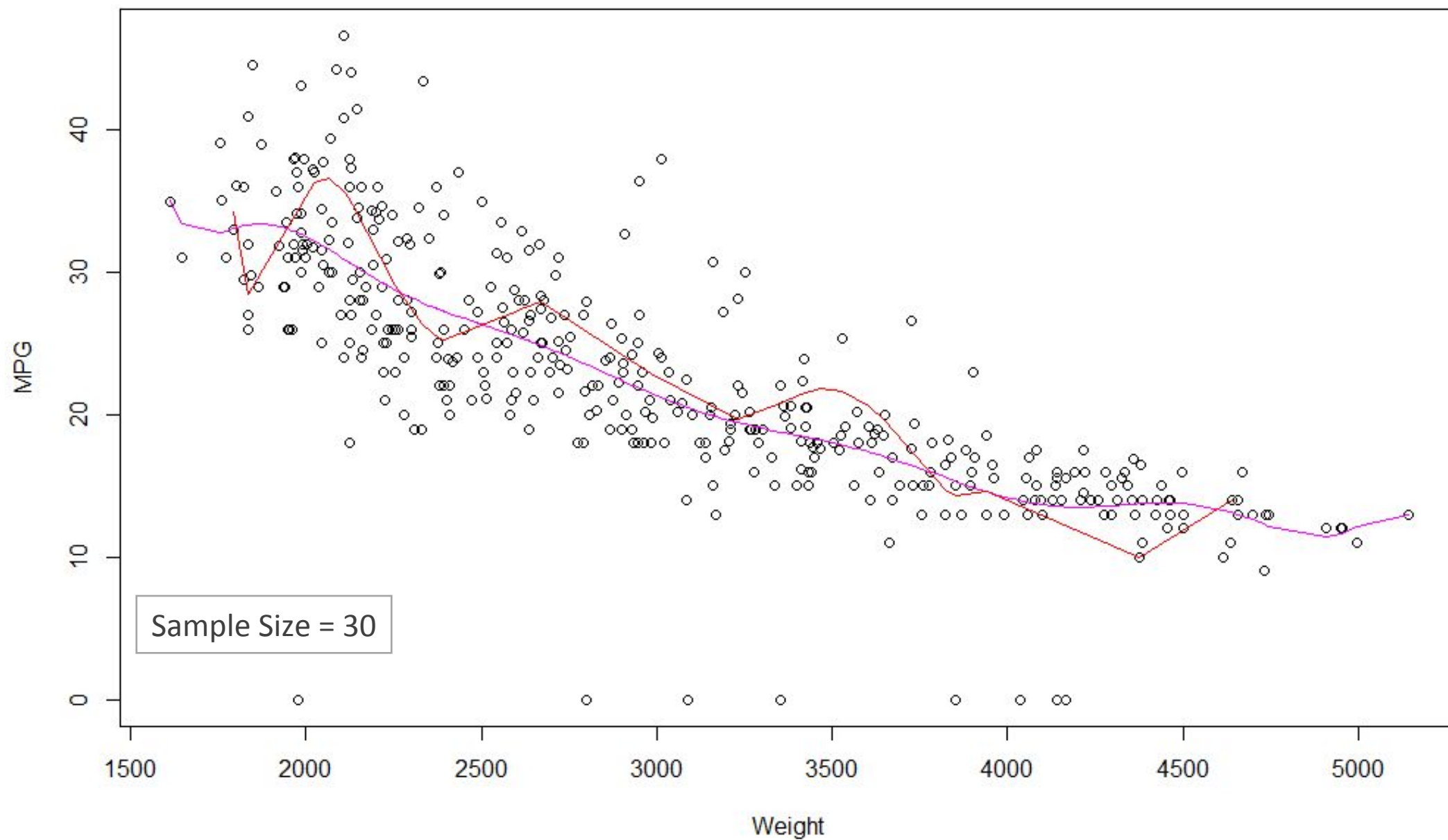
$$Var(X) = E[(X - E(X))^2]$$

MODEL COMPLEXITY AND VARIANCE

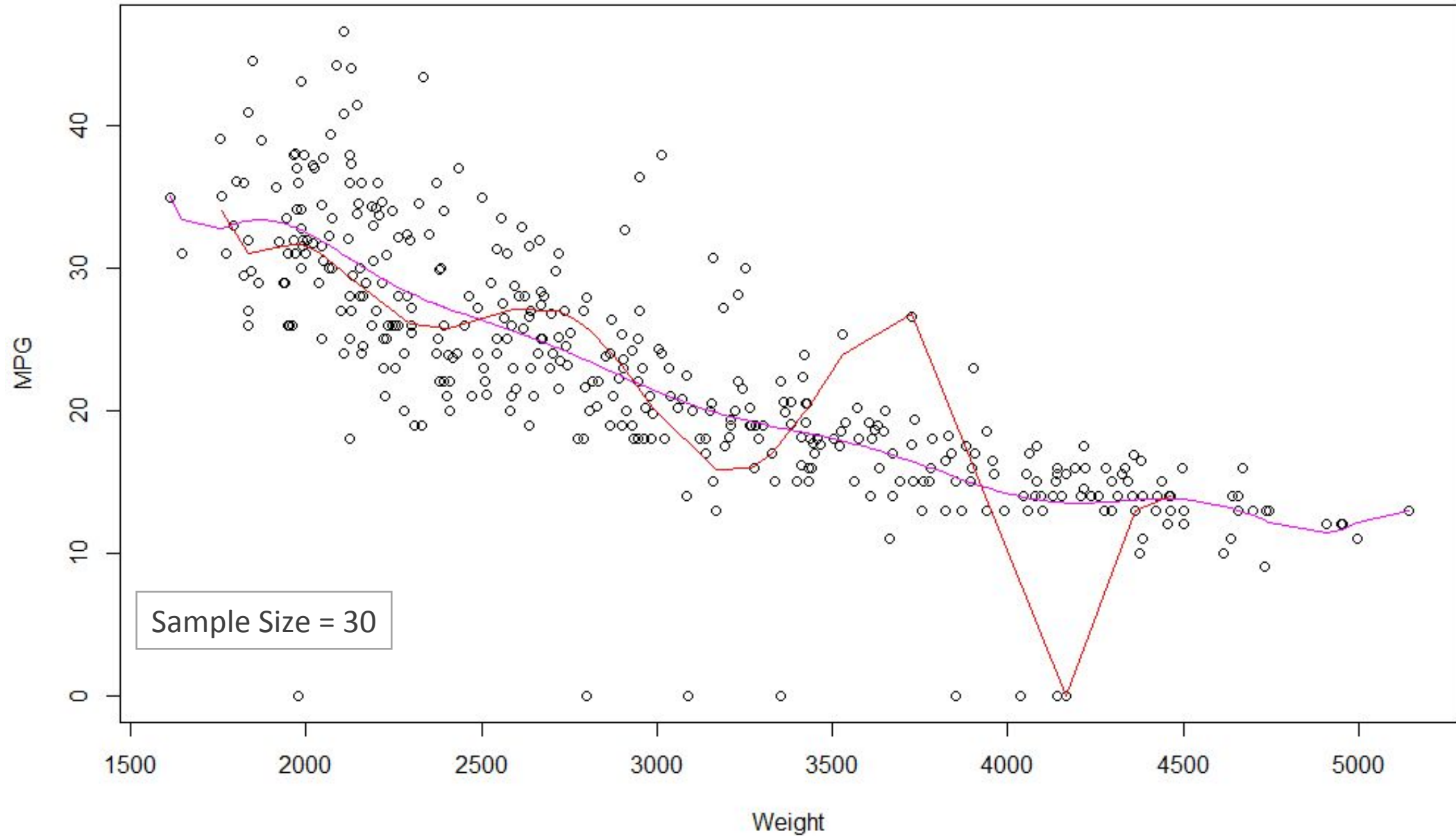
POLYNOMIAL OF ORDER 12



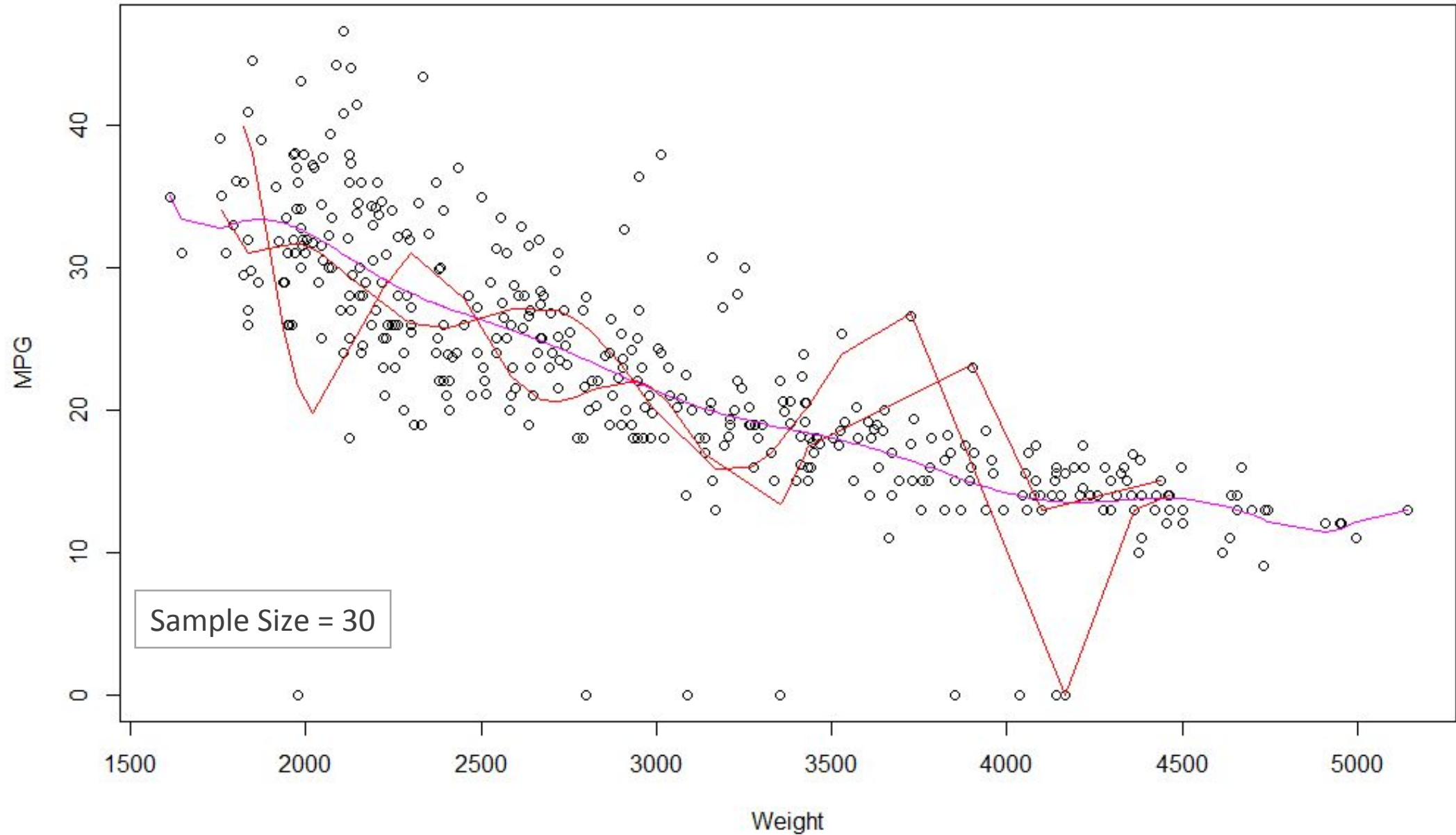
POLYNOMIAL OF ORDER 12 FITTED ON 1 SAMPLE



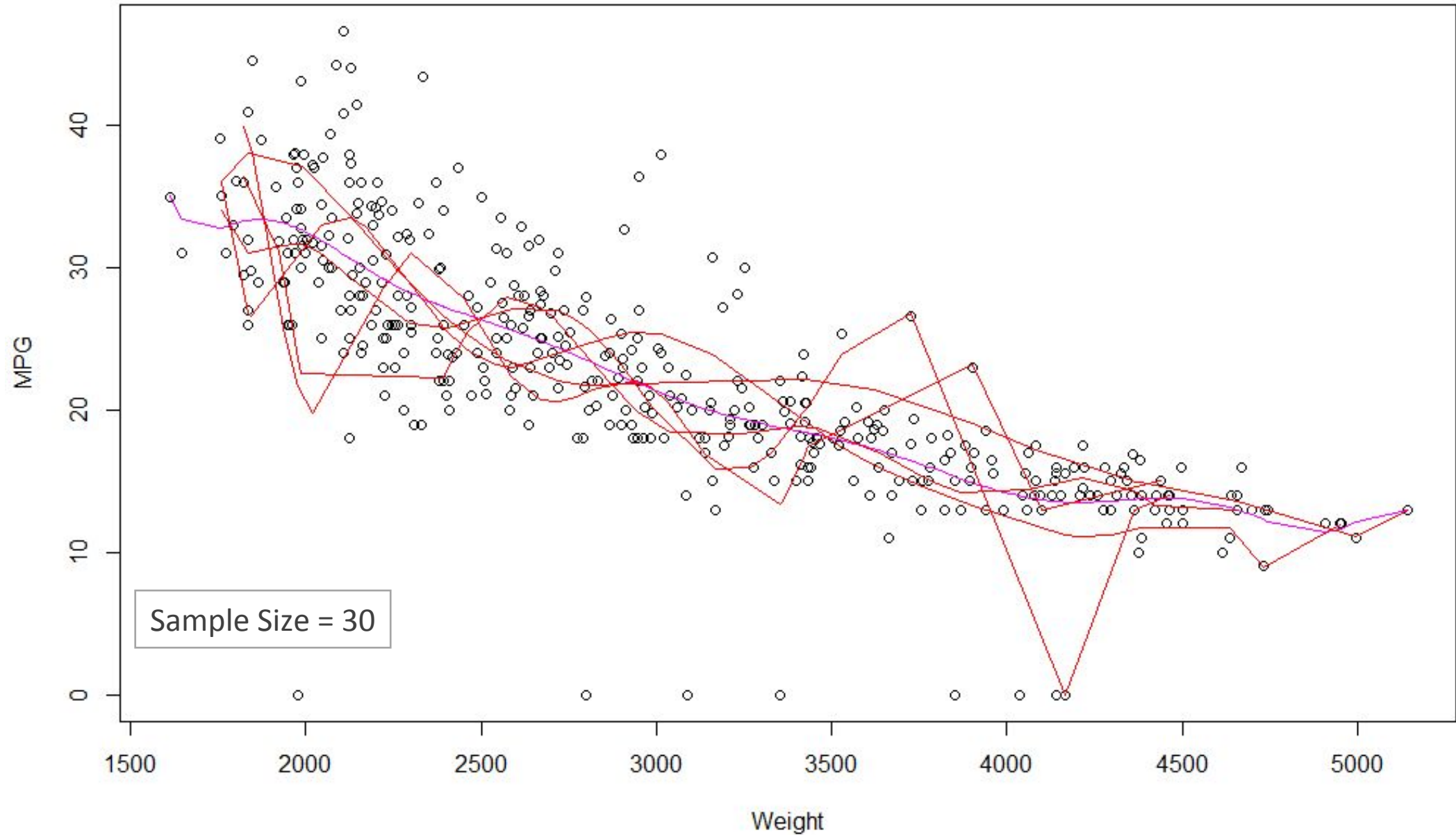
POLYNOMIAL OF ORDER 12 FITTED ON A DIFFERENT SAMPLE



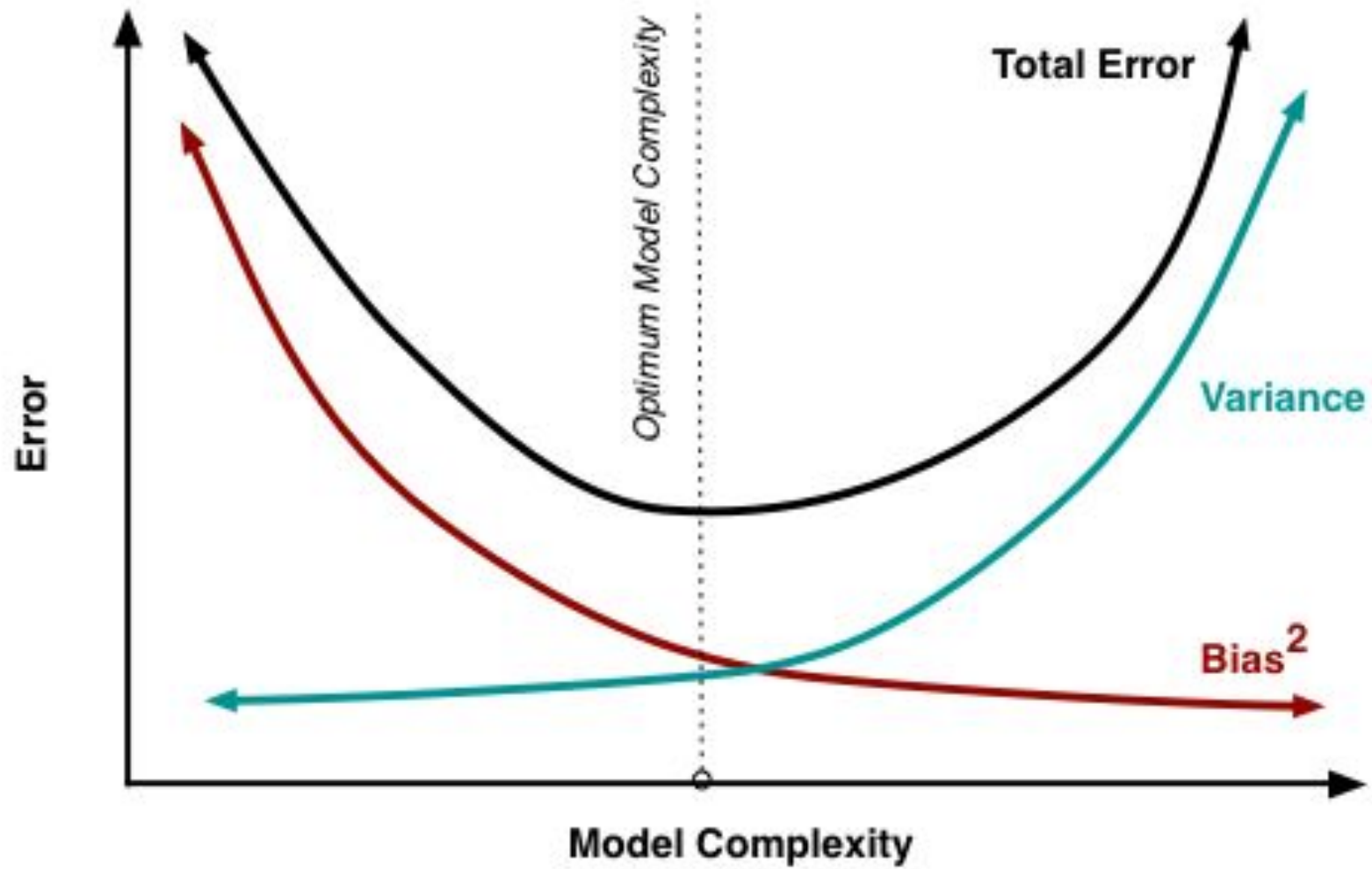
POLYNOMIAL OF ORDER 12 FITTED ON 2 DIFFERENT SAMPLES



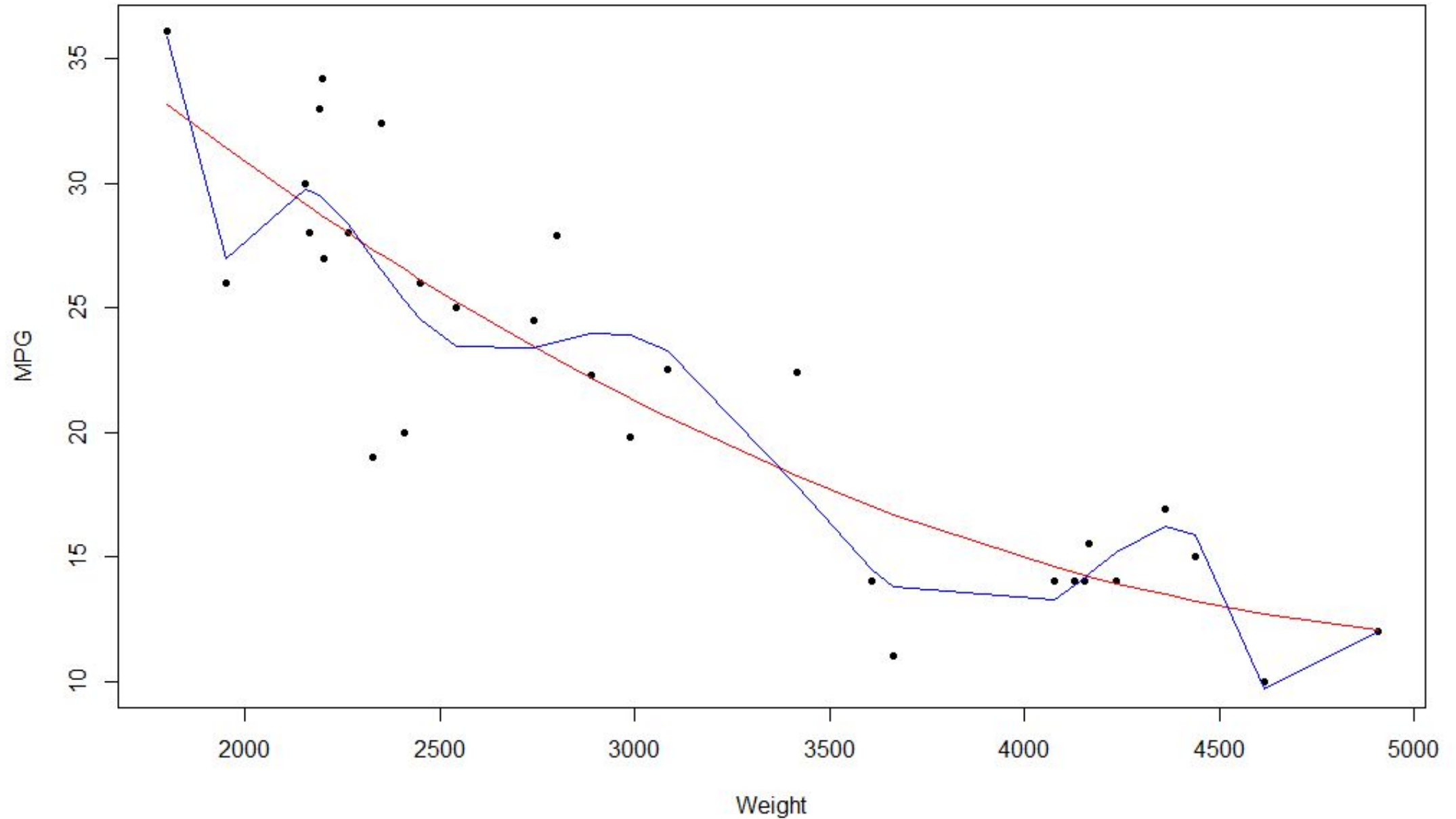
POLYNOMIAL OF ORDER 12 FITTED ON 5 DIFFERENT SAMPLES



MODEL COMPLEXITY – BIAS - VARIANCE

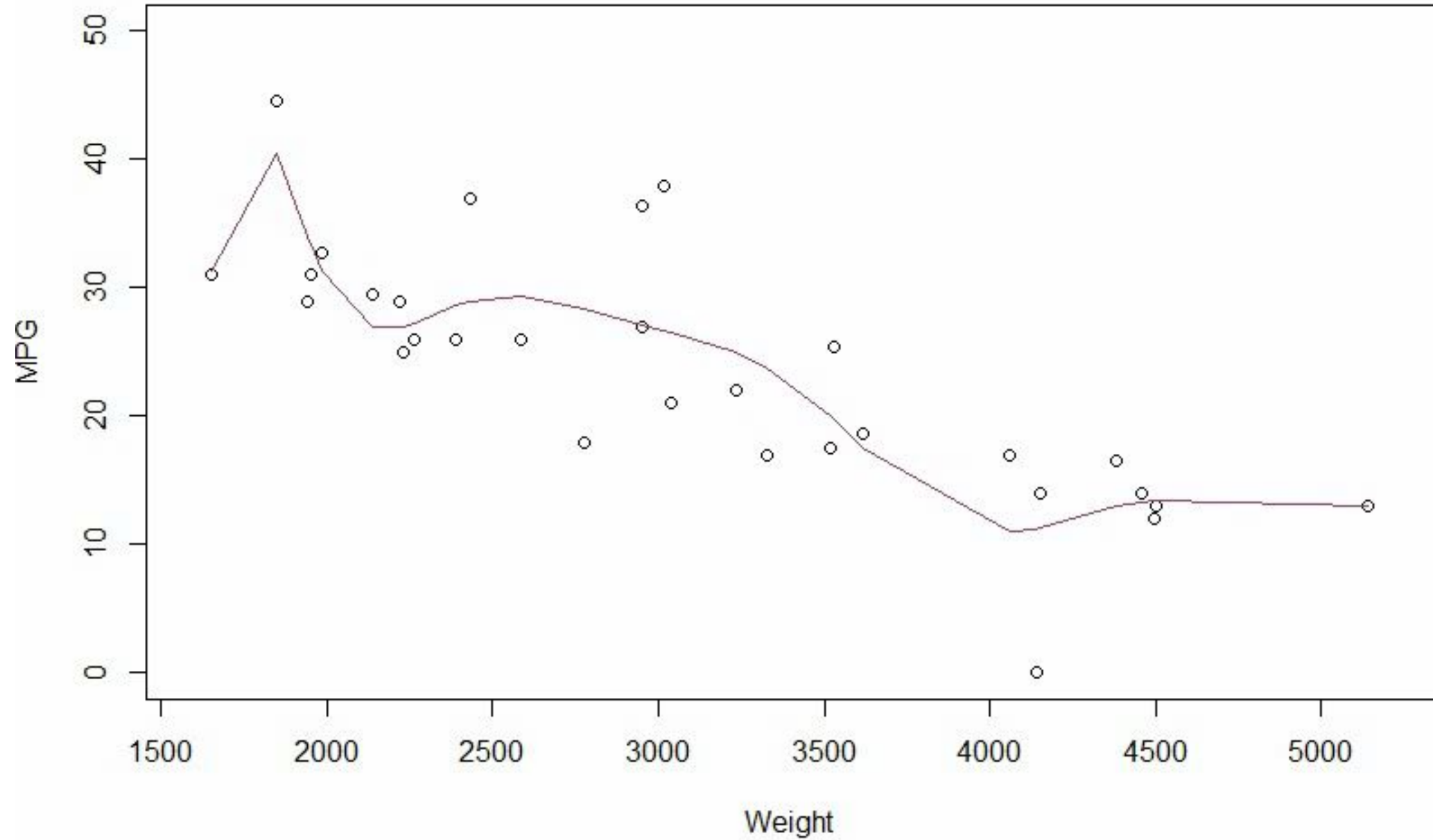


COMPLEXITY VS OVERFITTING

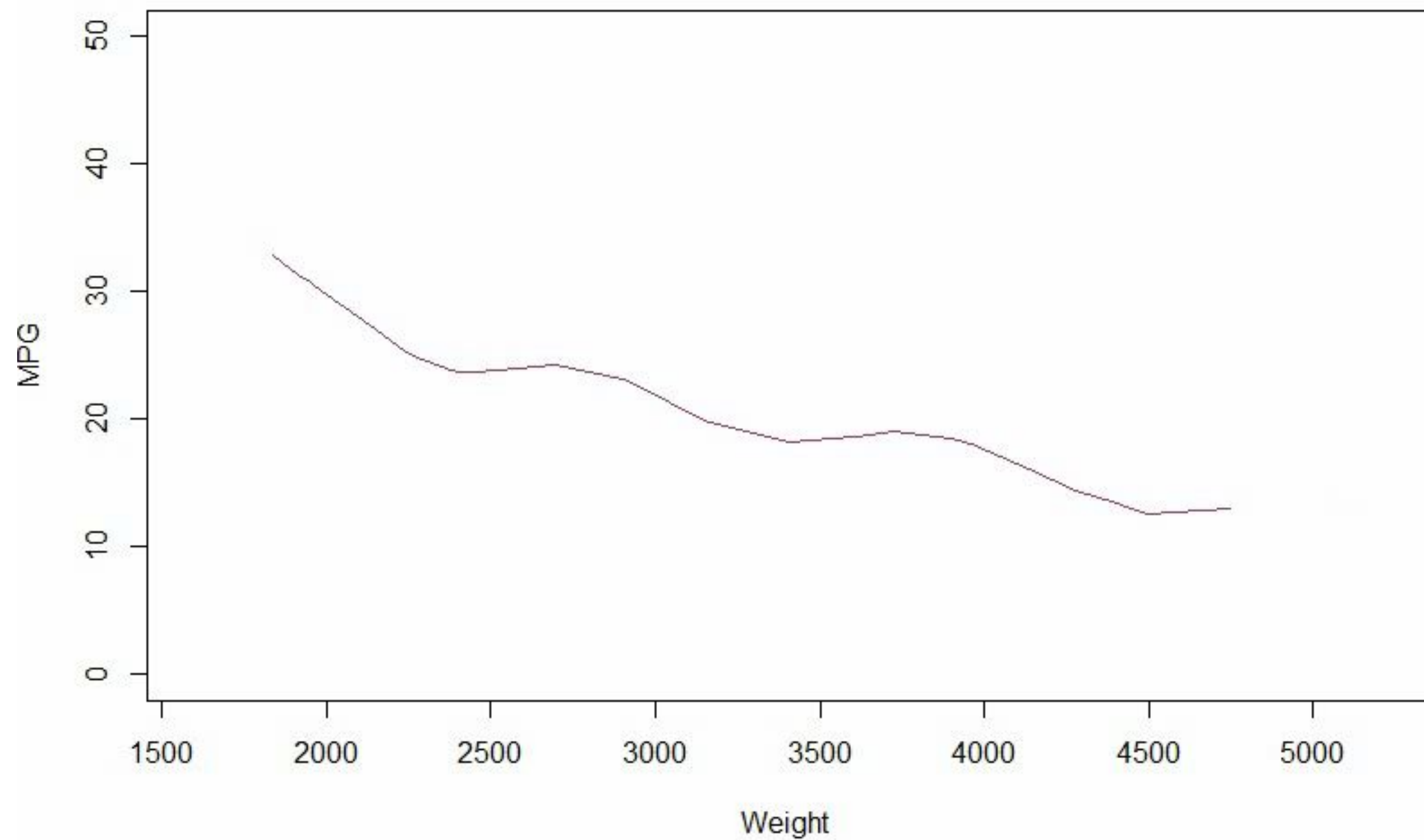


MODEL VARIANCE AGAINST SAMPLE SIZE

Sample Size = 30



Sample Size = 30



EXPECTED ERROR VS SAMPLE SIZE

