

Applied Regression Modeling:
A Business Approach
Chapter 2: Simple Linear Regression
Sections 2.4–2.7

by Iain Pardoe

2.4 Model assumptions	2
Regression model assumptions	2
Viewing the assumptions on a scatterplot	3
Checking the model assumptions	4
Residual plots which pass	5
Residual plots which fail	6
Histograms of residuals.	7
QQ-plots of residuals	8
Assessing assumptions in practice.	9
2.5 Model interpretation	10
Homes example model results	10
Interpreting model results	11
Regression summary plot	12
2.6 Estimation and prediction	13
Estimation and prediction	13
Confidence interval for population mean, $E(Y)$	14
Prediction interval for an individual Y -value.	15
Confidence and prediction intervals.	16
2.7 Chapter summary	17
Steps in a simple linear regression analysis.	17

Regression model assumptions

Four assumptions about random errors, $e = Y - E(Y) = Y - b_0 - b_1X$:

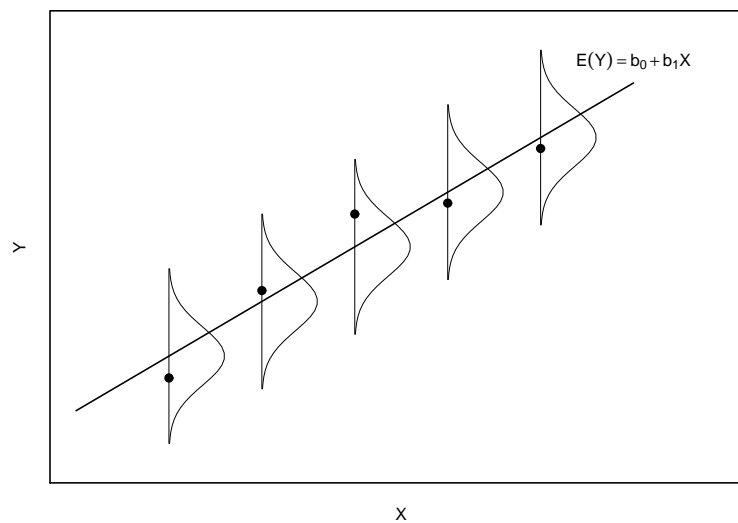
- Probability distribution of e at each value of X has a **mean of zero**;
- Probability distribution of e at each value of X has **constant variance**;
- Probability distribution of e at each value of X is **normal**;
- Value of e for one observation is **independent** of the value of e for any other observation.

© Iain Pardoe, 2006

2 / 17

Viewing the assumptions on a scatterplot

Random error probability distributions.



© Iain Pardoe, 2006

3 / 17

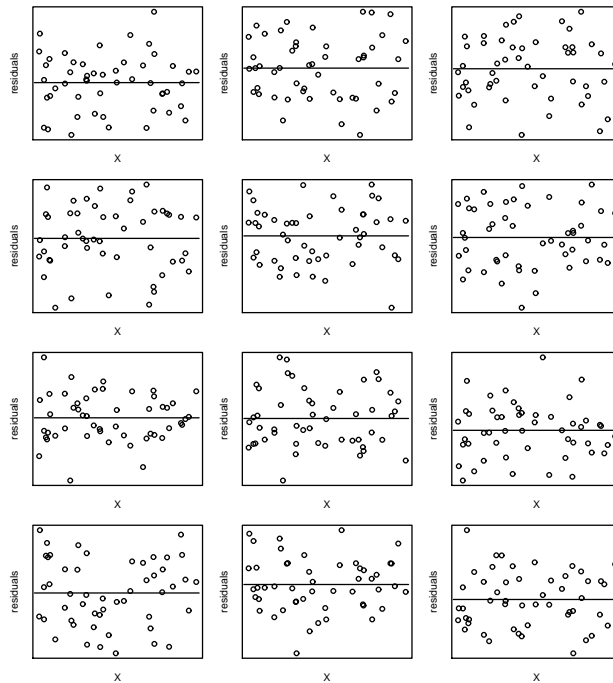
Checking the model assumptions

- Calculate residuals, $\hat{e} = Y - \hat{Y} = Y - \hat{b}_0 - \hat{b}_1 X$.
- Draw a residual plot with \hat{e} along the vertical axis and X along the horizontal axis.
 - Assess **zero mean** assumption—do the residuals average out to zero as we move across the plot from left to right?
 - Assess **constant variance** assumption—is the (vertical) variation of the residuals similar as we move across the plot from left to right?
 - Assess **independence** assumption—do residuals look “random” with no systematic patterns?
- Draw a histogram and QQ-plot of the residuals.
 - Assess **normality** assumption—does histogram look approximately bell-shaped and symmetric and do QQ-plot points lie close to line?

© Iain Pardoe, 2006

4 / 17

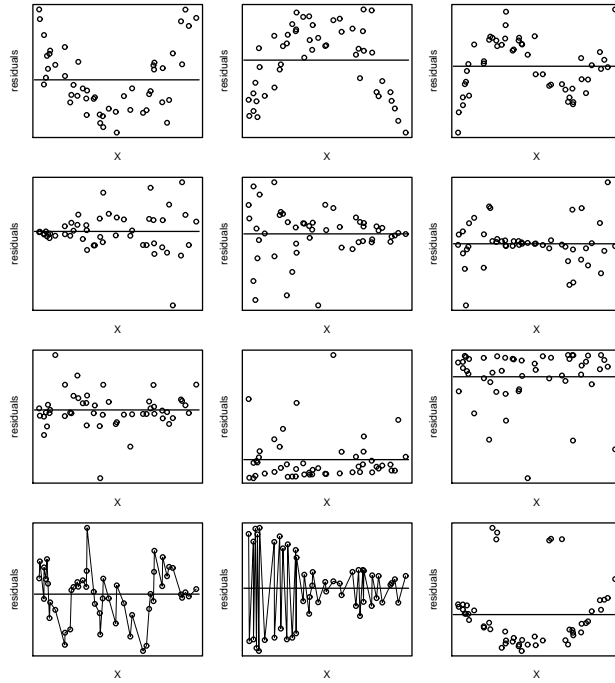
Residual plots which pass



© Iain Pardoe, 2006

5 / 17

Residual plots which fail

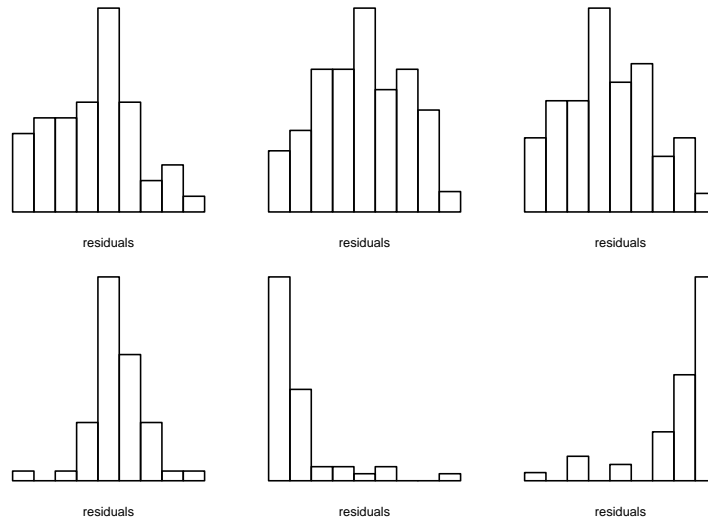


© Iain Pardoe, 2006

6 / 17

Histograms of residuals

Upper three pass, lower three fail

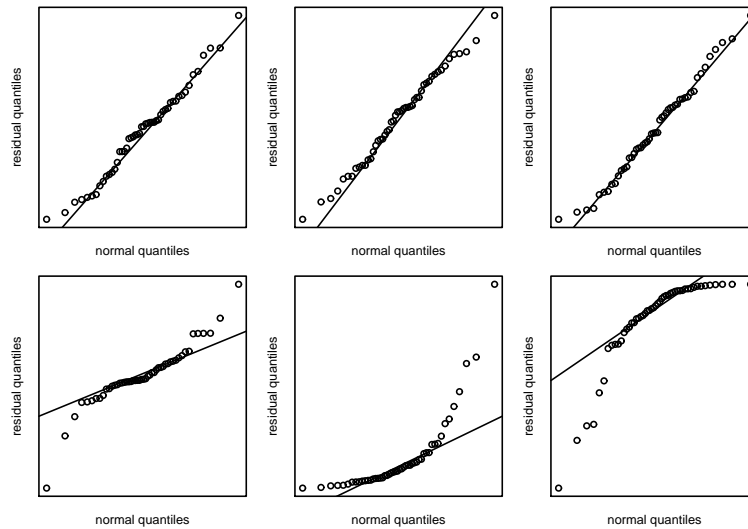


© Iain Pardoe, 2006

7 / 17

QQ-plots of residuals

Upper three pass, lower three fail



© Iain Pardoe, 2006

8 / 17

Assessing assumptions in practice

- Assessing assumptions in practice can be difficult and time-consuming.
- Taking the time to check the assumptions is worthwhile and can provide additional support for any modeling conclusions.
- *Clear* violation of one or more assumptions could mean results are questionable and should probably not be used (possible remedies to come in Chapters 3 and 4).
- Regression results tend to be quite robust to *mild* violations of assumptions.
- Checking assumptions when n is very small (or very large) can be particularly challenging.
- Example: **CARS2** data file—is weight or horsepower better for predicting cost?

© Iain Pardoe, 2006

9 / 17

Homes example model results

Model Summary

Model	Multiple R	R Squared	Adjusted R Squared	Regression Std. Error
1	0.972 ^a	0.945	0.927	2.7865

^a Predictors: (Intercept), X.

Parameters^a

Model		Estimate	Std. Error	t-stat	Pr(> t)
1	(Intercept)	190.318	11.023	17.266	0.000
	X	40.800	5.684	7.179	0.006

95% Confidence Interval

Model		Lower Bound	Upper Bound
1	(Intercept)	155.238	225.398
	X	22.712	58.888

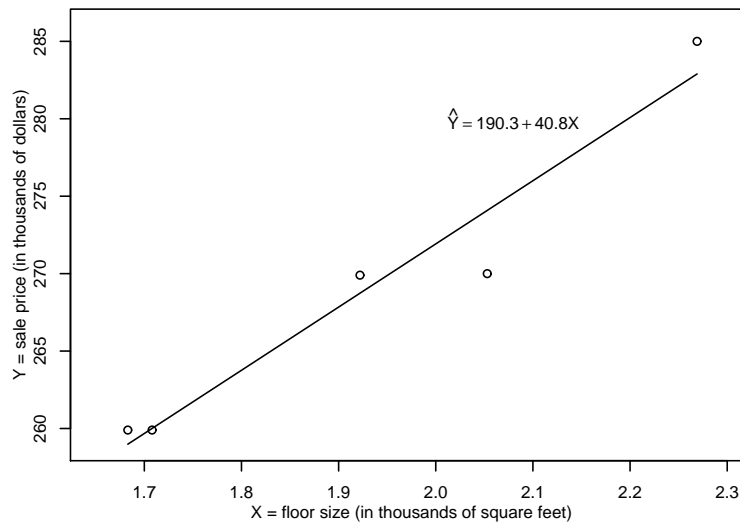
^a Response variable: Y.

Interpreting model results

- We found a statistically significant straight-line relationship (at a 5% significance level) between $Y = \text{sale price (\$k)}$ and $X = \text{floor size (k sq. feet)}$.
- Estimated equation: $\hat{Y} = \hat{b}_0 + \hat{b}_1 X = 190.3 + 40.8X$.
- $X = 0$ does not make sense for this application, nor do we have data close to $X = 0$, so we cannot meaningfully interpret $\hat{b}_0 = 190.3$.
- Expect sale price to increase \$4080 when floor size increases 100 sq. feet, for 1683–2269 sq. feet homes (95% confident sale price increases between \$2270 and \$5890 when floor size increases 100 sq. feet).
- Can expect a prediction of an unobserved sale price from a particular floor size to be accurate to within approximately $\pm \$5570$ (with 95% confidence).
- 94.5% of the variation in sale price (about its mean) can be explained by a straight-line relationship between sale price and floor size.

Regression summary plot

Simple linear regression model.



© Iain Pardoe, 2006

12 / 17

2.6 Estimation and prediction

13 / 17

Estimation and prediction

- Recall the confidence interval for a univariate population mean, $E(Y)$:
 $m_Y \pm t\text{-percentile}(s_Y/\sqrt{n})$.
- Also, a prediction interval for an individual univariate Y -value:
 $m_Y \pm t\text{-percentile}(s_Y\sqrt{1+1/n})$.
- Similar distinction between confidence and prediction intervals for simple linear regression.
- Confidence interval for the population mean, $E(Y)$, at a particular X -value is
 $\hat{Y} \pm t\text{-percentile}(s_{\hat{Y}})$.
- Prediction interval for an individual Y -value at a particular X -value is $\hat{Y}^* \pm t\text{-percentile}(s_{\hat{Y}^*})$.
- Which should be wider? Is it harder to estimate a mean or predict an individual value?

© Iain Pardoe, 2006

13 / 17

Confidence interval for population mean, E(Y)

- Formula: $\hat{Y} \pm t\text{-percentile}(s_{\hat{Y}})$
where $s_{\hat{Y}} = s \sqrt{\frac{1}{n} + \frac{(X_p - m_X)^2}{\sum_{i=1}^n (X_i - m_X)^2}}$.
- Interval is narrower:
 - when n is large;
 - when X_p is close to its sample mean, m_X ;
 - when the regression standard error, s , is small;
 - for lower levels of confidence.
- Example: for home prices–floor size dataset, the 95% confidence interval for E(Y) when $X = 2$ is (267.7, 276.1).
- Interpretation: we're 95% confident that average sale price is between \$267,700 and \$276,100 for 2000 square foot homes.

© Iain Pardoe, 2006

14 / 17

Prediction interval for an individual Y-value

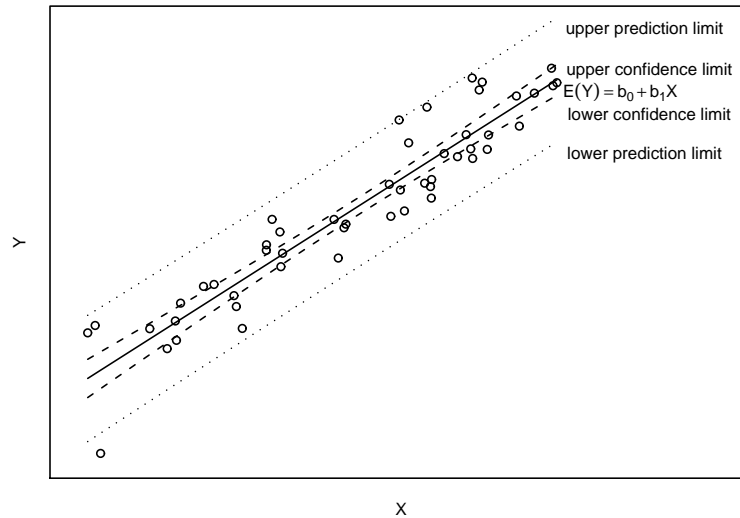
- Formula: $\hat{Y}^* \pm t\text{-percentile}(s_{\hat{Y}^*})$
where $s_{\hat{Y}^*} = s \sqrt{1 + \frac{1}{n} + \frac{(X_p - m_X)^2}{\sum_{i=1}^n (X_i - m_X)^2}}$.
- Interval is narrower:
 - when n is large;
 - when X_p is close to its sample mean, m_X ;
 - when the regression standard error, s , is small;
 - for lower levels of confidence.
- Since $s_{\hat{Y}^*} > s_{\hat{Y}}$, prediction interval is wider than confidence interval.
- Example: home prices–floor size dataset, the 95% prediction interval for Y^* at $X = 2$ is (262.1, 281.7).
- Interpretation: we're 95% confident that the sale price for an individual 2000 square foot home is between \$262,100 and \$281,700.
- What is a 95% prediction interval for large n ?

© Iain Pardoe, 2006

15 / 17

Confidence and prediction intervals

Compare widths of confidence and prediction intervals.



© Iain Pardoe, 2006

16 / 17

2.7 Chapter summary

17 / 17

Steps in a simple linear regression analysis

- Formulate model.
- Construct a scatterplot of Y versus X .
- Estimate model using least squares.
- Evaluate model:
 - Regression standard error, s ;
 - Coefficient of determination, R^2 ;
 - Population slope, b_1 .
- Check model assumptions.
- Interpret model.
- Estimate $E(Y)$ and predict Y .

© Iain Pardoe, 2006

17 / 17