

Simple Linear Regression with Normality

Assumed model: $\underline{Y_i} = \beta_0 + \beta_1 x_i + \epsilon_i$, $\epsilon_i \sim \text{i.i.d. } N(0, \sigma^2)$,
 $i = 1, \dots, n$.

Note: The values x_1, \dots, x_n of predictor X are known and fixed (i.e., non-random), and are assumed to be measured without error.

Properties:

- Properties:**

 - $E(Y_i|x_i) = E[\underbrace{\beta_0 + \beta_1 x_i}_\text{fixed} + \epsilon_i] = \beta_0 + \beta_1 x_i + \underbrace{E[\epsilon_i]}_\text{||} = \beta_0 + \beta_1 x_i$
 - $\text{var}(Y_i|x_i) = \text{var}[\underbrace{\beta_0 + \beta_1 x_i}_\text{fixed} + \epsilon_i] = \text{var}[\epsilon_i] = \sigma^2$.
 - $Y_i|x_i \sim \text{independent } N(\beta_0 + \beta_1 x_i, \sigma^2)$ i.e. and in X -

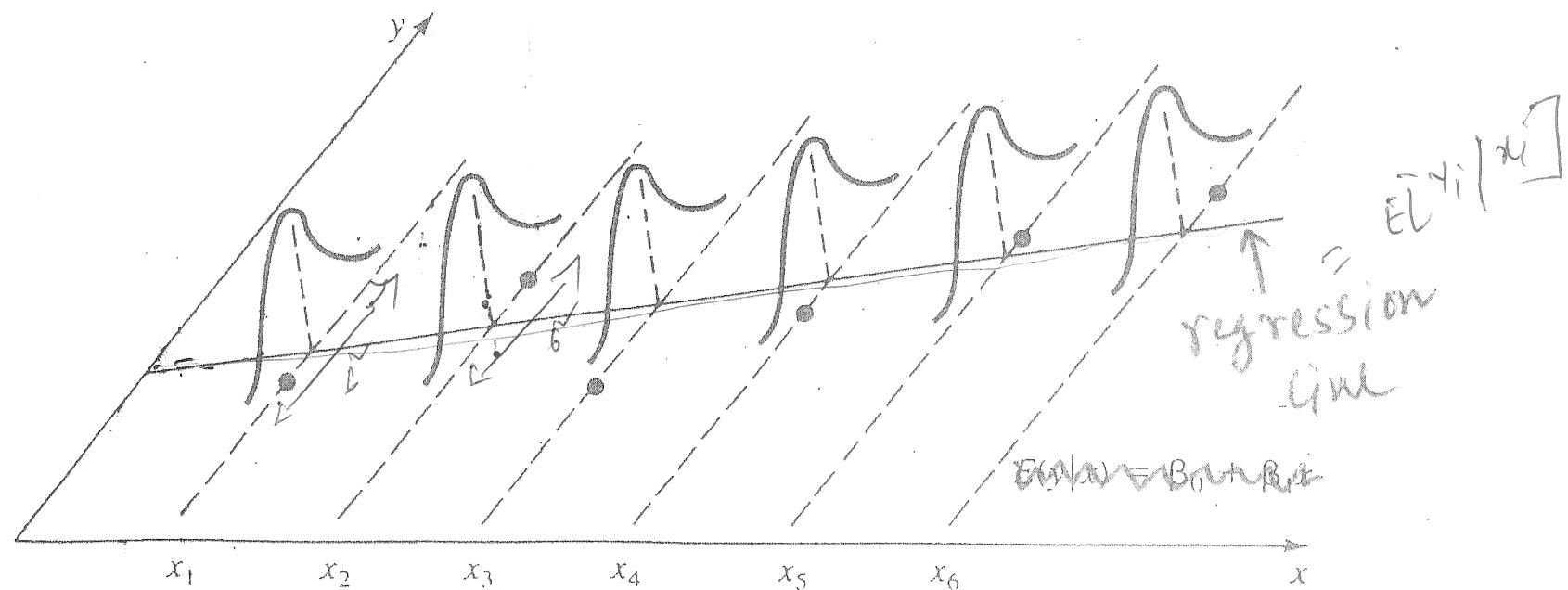
→ Since $E[y_i/x_i]$ changes with i , the y_i 's are not identically distributed.

$$\underbrace{(\beta_0 + \beta_1 x_i, \sigma^2)}_{E[y_i/x_i]} \quad \begin{array}{l} \text{constant, i.e. depend on } x \\ \text{don't depend on } x \\ \text{(homoscedasticity assumption).} \end{array}$$

Note that when $\beta_1 = 0$, we get the

One-Sample Setup

Simple linear regression model



Each y_i is a draw from $N[\beta_0 + \beta_1 x_i, \sigma^2]$.

What are the model parameters: $\beta_0, \beta_1, \sigma^2$ — 3 parameters in total.

- The least squares estimators $(\hat{\beta}_0, \hat{\beta}_1)$ of (β_0, β_1) are also maximum likelihood estimators.
- $\hat{\beta}_1 \sim N(\beta_1, \sigma^2 / \{(n-1)S_x^2\})$
- $E[\hat{\beta}_1] = \beta_1 \Rightarrow$ unbiased —
- Define: $\hat{\sigma}^2 = SS_{ERR}/(n-2)$. Then, $E(\hat{\sigma}^2) = \sigma^2 \Rightarrow \hat{\sigma}^2$ is unbiased.
- An unbiased estimator of σ^2 is $\hat{\sigma}^2 \rightarrow$ # mean related parameters - reg. coeff. - that are being estimated.
- Note: The sample variance S_y^2 is no longer unbiased for σ^2 . This is because the y_i 's are not identically distributed — they have different means, being estimated.
- $SS_{ERR}/\sigma^2 = (n-2)\hat{\sigma}^2/\sigma^2$ follows a χ^2 distribution with $(n-2)$ degrees of freedom, and is indep. of the distributions of the estimated regression coefficients.

$$e_i = y_i - \hat{y}_i$$

~~SS~~ SS_{ERR} (\equiv SS due to residuals) = $\sum_{i=1}^n (e_i - 0)^2$

Since sample variance of the residuals, with division by $n-2$.

ANOVA table: A standard summary of regression fit. Here we have “simple linear regression” *i.e.* two regression coefficients, β_0 and β_1 .

Source	SS	d.f.	MS	F
Model	SS_{REG}	1	$MS_{REG} = \frac{SS_{REG}}{1}$	$\frac{MS_{REG}}{MS_{ERR}}$
Error	SS_{ERR}	$n - 2$	$MS_{ERR} = \frac{SS_{ERR}}{n-2}$	
Total	SS_{TOT}	$n - 1$		↑ see later.

Recall that:

- $SS_{TOT} = \sum_{i=1}^n (Y_i - \bar{Y})^2$ — measures total variation in the response
- $SS_{REG} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$ — " " the predicted value \Rightarrow variation explained by the regression
- $SS_{ERR} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$

$$SS_{TOT} - SS_{REG}$$

$$\sum_{i=1}^n (e_i - \bar{e})^2$$

remaining variation, variation in the residuals.

$$E[Y|X] = \beta_0 + \beta_1 X.$$

Inference about slope β_1

Issue: Is the predictor X “significant”, i.e., does it really help in predicting the response Y ? $\nearrow x \text{ is not useful for predicting } Y$

Approach 1: Test $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$. This is equivalent to testing $H_0 : \rho = 0$ vs. $H_1 : \rho \neq 0$. (Why?) $\rightarrow x \text{ is useful for predicting } Y$.

Test statistic:

$$t = \frac{\hat{\beta}_1 - 0}{\hat{SE}(\hat{\beta}_1)}, \quad \hat{SE}(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{(n-1) s_x^2}}$$

↑ Correlation $y_w (y, x)$

Recall: Test statistic has the general form:

$$\frac{\hat{\theta} - \theta_0}{\hat{SE}(\hat{\theta})}$$

Null distribution:

$$t_{(n-2)} \rightarrow \text{err d.f.}$$

- A two-sided t -test.

— know how to do a t -int.

100(1 - α)% Confidence Interval for β_1 : $\hat{\beta}_1 \pm t_{\alpha/2, n-2} \hat{SE}(\hat{\beta}_1)$

Approach 2: Test for model significance. In simple linear regression, this is equivalent to testing $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$.

Test statistic:

Reject H_0 if
~~F \leq statistic is large.~~ $F = \frac{MS_{\text{REG}}}{MS_{\text{ERR}}}$

→ will see that
~~F~~ F is large if R^2 is large.

Null distribution: This F statistic follows an F distribution with numerator d.f. 1 and denominator d.f. $n - 2$.

- An F -test.
- Equivalent to the t -test seen before because $t^2 = F$ (verify).

d.f. for the model
from Anova table

d.f. associated
with error

level α F-test: Reject H_0 if $F > F_{\alpha, 1, n-2}$ ← upper α th percentile of F -dist. of 1 and $n-2$ d.f.
 31/74

Model evaluation

Issue: Is the fitted model a good representation of the data?

Approach: Examine the residuals, $e_i = Y_i - \hat{Y}_i$, $i = 1, \dots, n$, and verify the key assumptions, namely,

- Errors have mean zero and constant variance
- Errors are normally distributed
- Errors are independent — often an issue when the data are collected over time.

Key Graphical Tools:



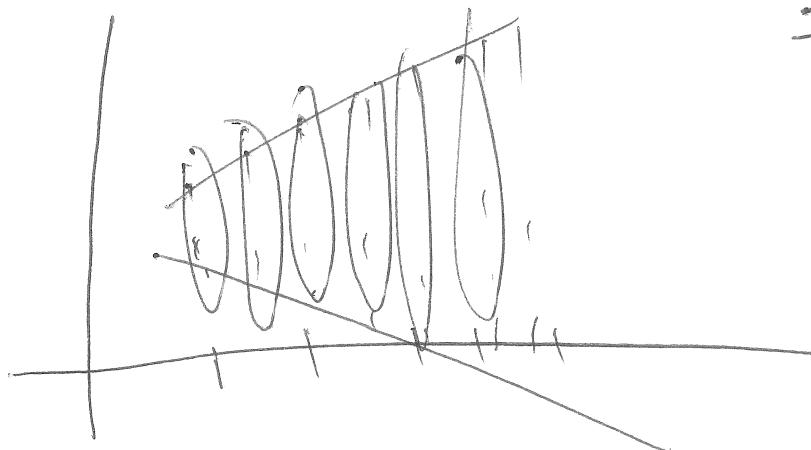
Two common deviations:
• trend in the plot
⇒ regression model
is not correct.

- **Residual plot:** Plot of residuals e_i against fitted values \hat{Y}_i . In the ideal plot, the points are scattered around zero and there is no pattern. This verifies the first assumption.
- **Normal QQ plot:** This verifies the normality assumption.
- **Time series plot:** Plot e_i against i . In the ideal plot, there should be no dependence, which verifies the independence assumption. More sophisticated tools exist.

NOT
so much
of
invent.
D&C

Another common deviation:

- Non-constant vertical scatter in the plot



⇒ Err don't have
a common variance.

Ex: House price data, continued.

```
x <- house$size  
y <- house$price
```

$$\textcircled{n = 58}$$

```
house.reg <- lm (y ~ x)
```

```
# ANOVA table
```

```
> (anova(house.reg))
```

Analysis of Variance Table

Response: y

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	→ F-stat p-value
→ x	1	71534	71534	184.62	< 2.2e-16	***
→ Residuals	56	21698	387			p-value for testing the null hyp. of model significance
	---	58-2				Here $H_0: \beta_1 = 0$.

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

→ Reject $H_0 \Rightarrow$ useful prediction

>

Testing for zero slope

> summary(house.reg)

Call:

lm(formula = y ~ x)

Residuals:

Min	1Q	Median	3Q	Max
-38.489	-14.512	-1.422	14.919	54.389

5- # summary
of the residuals.
Easier to interpret
after dividing by
 \bar{x}

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept) $\hat{\beta}_0$	5.432	$s_e(\hat{\beta}_0)$	8.191	0.663
x	56.083	4.128	13.587	$<2e-16$ ***

except these
to y_w (23,3)

$0.51 \rightarrow H_0: \beta_0 = 0$

/
Generally
not
of
interest.

$H_0: \beta_1 = 0$

~~not 10
SE~~

34 / 74

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 19.68 on 56 degrees of freedom

Multiple R-squared: 0.7673, Adjusted R-squared: 0.7631

F-statistic: 184.6 on 1 and 56 DF, p-value: < 2.2e-16

> ANOVA
Task

R²

Latenz

Confidence interval for slope

> confint(house.reg)

	2.5 %	97.5 %
(Intercept)	-10.97619	21.83933
x	47.81473	64.35183

} t-intervale

>

Prediction at a new x

x.new <- data.frame(x=3) more than one val: $\hat{y} = c(3, 2, \dots)$.

> predict(house.reg, newdata=x.new)

1

173.6814

>

Use fitted(house.reg) to get the fitted values

Use resid(house.reg) to get the residuals

↓ ~~resid~~ resid(house.reg, type = "pear")

divide the residuals
by $\hat{\sigma}$.

Residual plot

plot(fitted(house.reg), resid(house.reg))
abline(h=0)

QQ plot

```
qqnorm(resid(house.reg))
qqline(resid(house.reg))

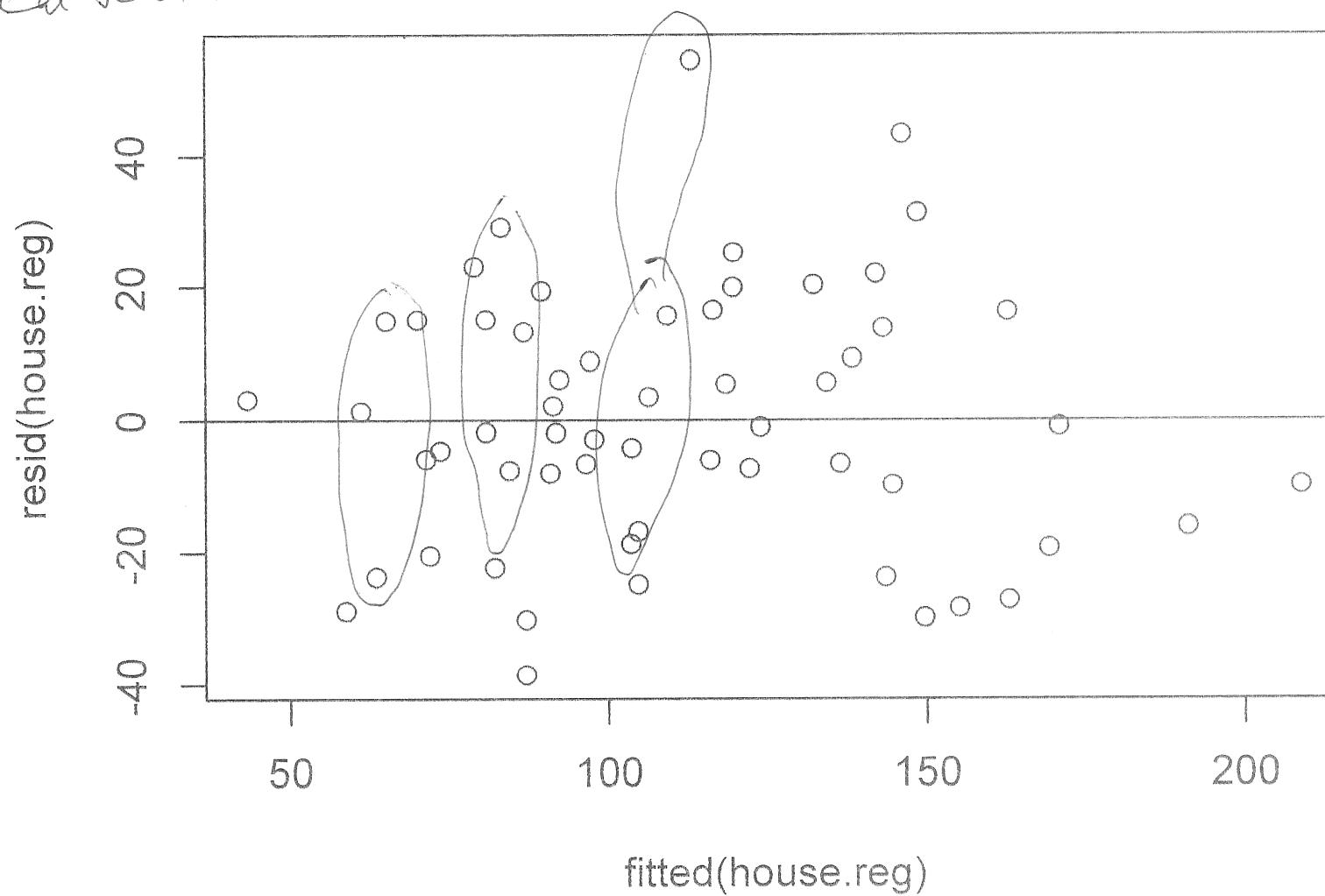
# Time series plot of residuals

plot(resid(house.reg), type="l")
abline(h=0)
```

Here:

- No trend
- No change in vertical scatter.

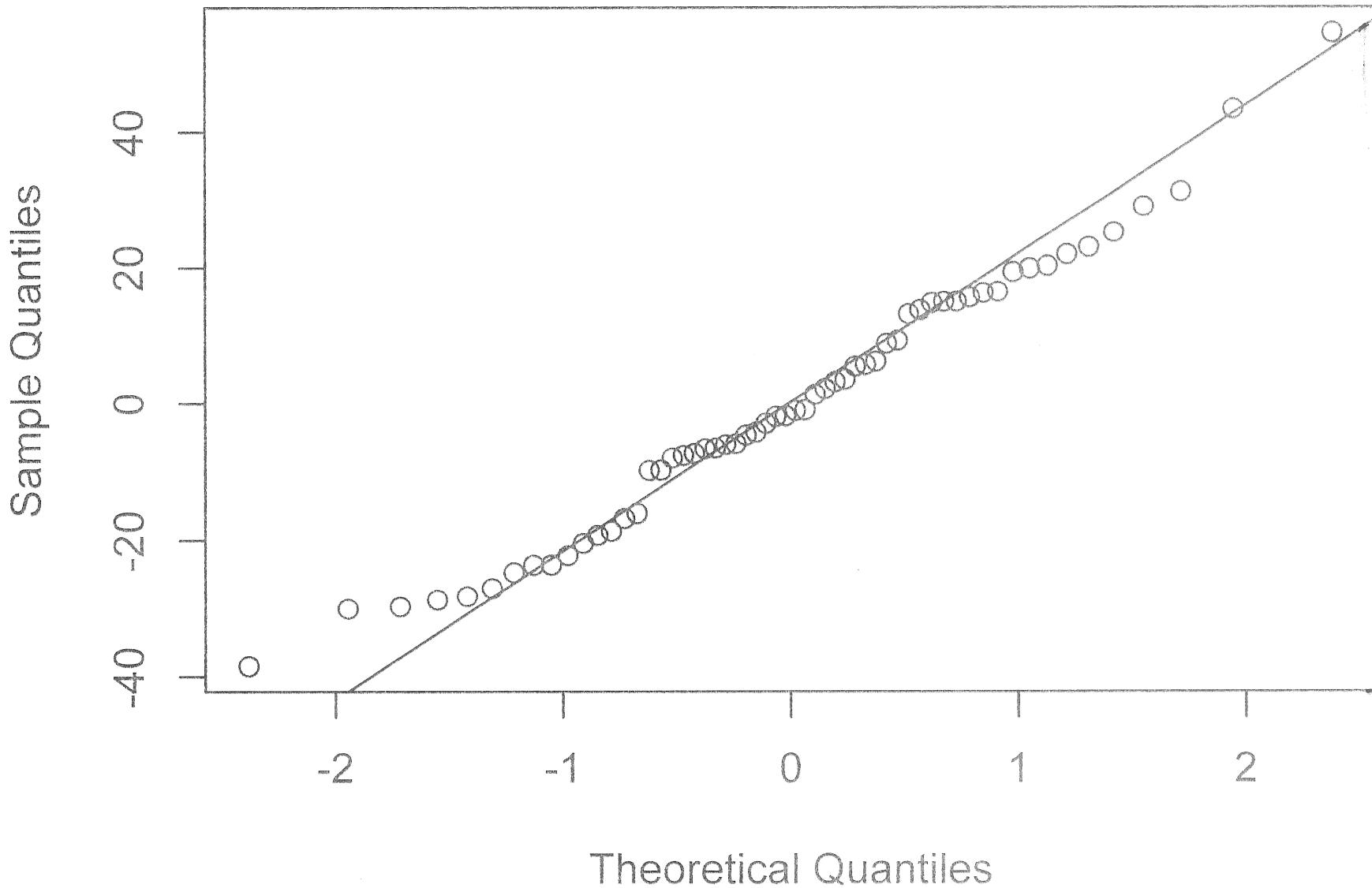
Residual plot



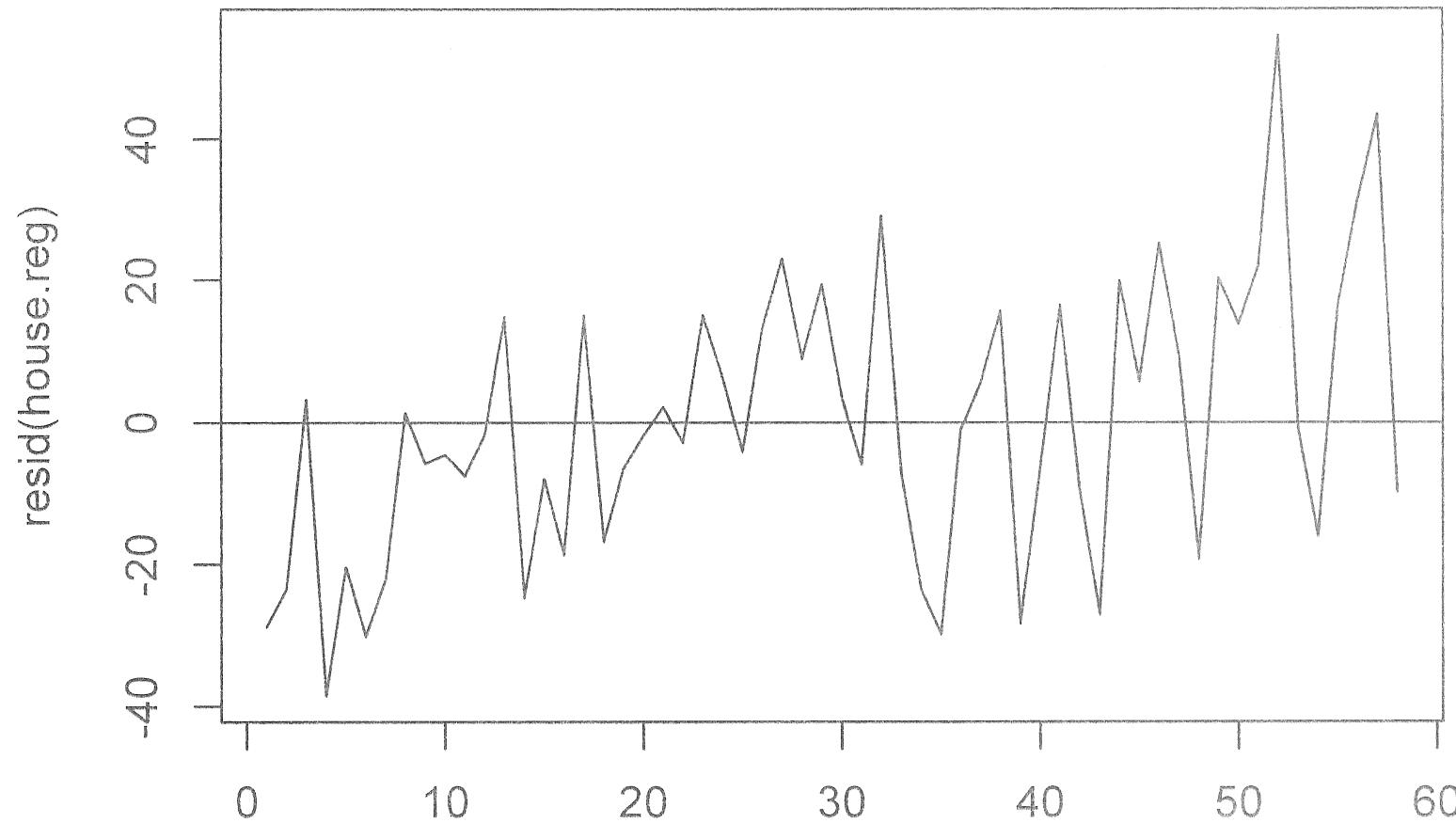
- What if a trend is seen?
 - Assumption about $E[Y|X]$ not correct.
 - Add higher order terms
 - Bring in more predictors
 - may be a transformation of X and/or Y is needed. — e.g., square-root, cube-root, and log —
- What if non-constant vertical scatter is seen?
 - Try a transformation

The normality assumption
seems reasonable

Normal Q-Q Plot



Time series plot



The pattern here is not
totally ~~random~~ random. There seems
to be some dependence over time,
but we will ignore it in this class.