

Testing hypotheses: Additionally assume that ϵ_i are normal, i.e., $\epsilon_i \sim \text{independent } N(0, \sigma^2)$.

- Each $\hat{\beta}_j$ follows a normal distribution
- $(n - k - 1)\hat{\sigma}^2/\sigma^2$ follows a χ^2_{n-k-1} distribution

$$E[Y|X] = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

Testing significance of j th predictor:

- $H_0: \beta_j \leq 0$ vs. $H_1: \beta_j \neq 0$
- True: Keep all the other predictors in model.
⇒ check whether X_j is useful or not after adjusting for the effects of the other predictors in model.
- H_0 : Predictor X_j is not useful for predicting response after adjusting for the other predictors
 - Test statistic: $(\hat{\beta}_j - 0) / \hat{SE}(\hat{\beta}_j)$
 - Null distribution: $\sim t$ w.r.t. d.f. when H_0 is true.
 - Rejection region:
 - p -value:
 - $100(1 - \alpha)\%$ CI for β_j :

→ just a t-test and know how to do this.

Testing model significance:

$$H_0 : \beta_1 = \dots = \beta_k = 0 \text{ vs. } H_1 : \text{at least one } \beta_j \neq 0$$

- H_0 : None of the predictors is useful for predicting response
- Test statistic: $\frac{MS_{\text{reg}}}{MS_{\text{error}}} \sim F_{K, n-K-1}$ if H_0 is true.
- Null distribution: \downarrow from d.f.
- Rejection region: $\# \text{ signs}$
- p -value:

→ an F-test.

know
how to
do this.

Coefficient of determination: As before,

$$R^2 = \frac{SS_{REG}}{SS_{TOT}} = 1 - \frac{SS_{ERR}}{SS_{TOT}}$$

- R^2 = proportion of total variation explained by regression
- R^2 = square of correlation between (Y_i, \hat{Y}_i) , $i = 1, \dots, n$ — can be verified
- When we add new predictors to the model, $R^2 \uparrow$.
- Even adding useless predictors will also ~~also~~ $\uparrow R^2$.
- Not a fair criterion for comparing models with different numbers of predictors

Beware of overfitting a model:

- Overfitting = Having too many predictors in the model
- An overfitted model will provide a good fit to the data at hand, but it will be terrible at predicting future observations.
- Estimated regression coefficients have large standard errors.
- See Figure 11.4 on page 369

if
estimates are
reliable

if
predictions are
not reliable in
the sense that
they will have
large SE's -

Adjusted R^2 :

$$R^2 = 1 - \frac{SS_{\text{ERR}}}{SS_{\text{TOT}}}.$$

$$R_{\text{adj}}^2 = 1 - \frac{\overbrace{SS_{\text{ERR}}/(n-k-1)}^{\text{MS}_{\text{ERR}}}}{\overbrace{SS_{\text{TOT}}/(n-1)}^{\text{MS}_{\text{TOT}}}}$$

- Unlike R^2 , R_{adj}^2 rewards adding a predictor only if it reduces the error SS considerably
- Imagine adding a useless predictor. In this case, SS_{REG} and hence SS_{ERR} does not change ^{much}. However, $SS_{\text{ERR}}/(n-k-1)$ increases, which in turn decreases R_{adj}^2
- A more fair measure of goodness-of-fit than R^2
- Can be used to compare two models with different numbers of predictors — choose the model with the highest R_{adj}^2

Comparing two Nested Models

Nested models: Model 2 is nested within Model 1 if the predictors of Model 2 are a subset of predictors of Model 1.

Issue: How to compare two nested models?

Full model: Predictors X_1, \dots, X_m *common to both* (full model)

Reduced model: Predictors X_1, \dots, X_k , i.e., it does not have predictors $(X_{k+1}, \dots, X_m) \rightarrow (m-k)$ in #. *full model = reduced model*

Hypotheses: $H_0: \beta_{k+1} = \dots = \beta_m = 0$, vs., $H_1: \text{at least one slope } \neq 0$ *some of the last (m-k) predictors are useful.*

Extra sum of squares: Difference in variation explained by the two models

$$\begin{aligned} SS_{EX} &= \boxed{SS_{REG}(\text{full}) - SS_{REG}(\text{reduced})} \\ &= \boxed{SS_{ERR}(\text{reduced}) - SS_{ERR}(\text{full})} \end{aligned}$$

variability explained by the last (m-k) predictors.

regression slopes set to zero under H_0 .

- SS_{EX} has $m - k$ degrees of freedom. It equals the number of regression coefficients set to zero under H_0 .

Test statistic:

$$F = \frac{MS_{EX}}{MS_{ERR}(\text{full})} = \frac{SS_{EX}/(m - k)}{SS_{ERR}(\text{full})/(n - m - 1)}$$

↑ note: This is not $MS_{ERR}(\text{reduced})$.

$\sim F_{m-k, n-m-1}$.
if H_0 is true.

Null distribution:

Rejection region:

p-value:



- aka “partial F-test”
- Used for designing stepwise model selection procedures (see pages 392-394)
- To test $H_0: \beta_j = 0 \Leftrightarrow H_1: \beta_j \neq 0$ — can do partial F test (with num d.f. 1 and denom d.f. = num d.f.)
can see that: $F = t_j^2$ — two tests are equivalent
(their p-values will be the same)

Example: Home price data. These data come from a sample of homes sold in Maplewood, NJ in 2001.

```
# Read the home price data
```

```
home <- read.table("homeprice_multiple_predictors.txt",  
sep=",", header=T)
```

```
> str(home)
```

```
'data.frame': 29 obs. of 7 variables:  
 $ list : num 80 151 310 295 339 ...  
 $ sale : num 118 151 300 275 340 ...  
 $ full : int 1 1 2 2 2 1 3 1 1 1 ...  
 $ half : int 0 0 1 1 0 1 0 1 2 0 ...  
 $ bedrooms : int 3 4 4 4 3 4 3 3 3 1 ...  
 $ rooms : int 6 7 9 8 7 8 7 7 7 3 ...  
 $ neighborhood: int 1 1 3 3 4 3 2 2 3 2 ...
```

```
>
```

```
# Attach the dataset in R's memory so that we can  
# directly use the names of the variables
```

```
attach(home)
```

*Ans + step: Exploratory analysis
(boxplots, scatterplots, etc.)*

```
# Look at distributions of some predictors
```

```
> table(bedrooms)
```

```
bedrooms
```

```
1 2 3 4 5
```

```
1 3 16 8 1
```

```
>
```

```
> table(full)
```

```
full
```

```
1 2 3
```

```
13 11 5
```

```
>
```

```
> table(half)
half
  0   1   2
13 13  3
>
> table(neighborhood)
neighborhood
  1   2   3   4   5
  2   8  12  5   2
>

# Regress sale price on # bedrooms and neighborhood

fit1 <- lm(sale ~ bedrooms + neighborhood)

> summary(fit1)
```

Call:

lm(formula = sale ~ bedrooms + neighborhood)

Residuals:

Min	1Q	Median	3Q	Max
-90.871	-39.861	0.636	28.815	107.660

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
β_0 (Intercept)	-132.057	40.341	-3.273	0.003001 **	$H_0: \beta_0 = 0$ $H_1: \beta_0 \neq 0$
β_1 bedrooms	42.483	11.446	3.712	0.000987 ***	$H_0: \beta_1 = 0$ $H_1: \beta_1 \neq 0$
β_2 neighborhood	93.493	9.101	10.273	1.21e-10 ***	$H_0: \beta_2 = 0$ $H_1: \beta_2 \neq 0$

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 $H_1: \beta_2 \neq 0$

$\hat{6}$ err df.

Residual standard error: 47.3 on 26 degrees of freedom

Multiple R-squared: 0.8491, Adjusted R-squared: 0.8375

F-statistic: 73.16 on 2 and 26 DF, p-value: 2.1e-11

> $H_0: \beta_1 = \beta_2 = 0$ vs. $H_1: \beta_1$ or β_2 is not zero.

keeping all the remaining prediction in the model.

```
# Add # full and half baths
```

Either 'lm' ~
update,

```
fit2 <- update(fit1, . ~ . + full + half)
```

```
> summary(fit2)
```

Call:

```
lm(formula = sale ~ bedrooms + neighborhood + full + half)
```

Residuals:

Min	1Q	Median	3Q	Max
-56.554	-38.067	6.027	26.998	53.311

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-125.121	33.136	-3.776	0.000926 ***	
bedrooms	29.513	10.091	2.925	0.007419 **	

neighborhood	78.724	9.669	8.142	2.31e-08	***
full	27.345	13.604	2.010	0.055785	.
half	45.553	12.129	3.756	0.000974	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '1'
 ↓ Initial F-test
 for $\beta_{full} = 0$
 will have
 p-value = 0.055785

Residual standard error: 38.79 on 24 degrees of freedom

Multiple R-squared: 0.9063, Adjusted R-squared: 0.8907

F-statistic: 58.05 on 4 and 24 DF, p-value: 5.425e-12

>

Drop # full baths

fit3 <- update(fit2, . ~ . - full)

> summary(fit3)

Call:

```
lm(formula = sale ~ bedrooms + neighborhood + half)
```

Residuals:

Min	1Q	Median	3Q	Max
-67.55	-42.27	7.17	26.93	68.83

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-127.348	35.073	-3.631	0.00127 **
bedrooms	35.649	10.187	3.500	0.00177 **
neighborhood	90.982	7.947	11.449	1.95e-11 ***
half	37.004	12.030	3.076	0.00503 **

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

Residual standard error: 41.08 on 25 degrees of freedom

Multiple R-squared: 0.8905, Adjusted R-squared: 0.8774

F-statistic: 67.8 on 3 and 25 DF, p-value: 3.808e-12

>

```
# Compare the nested models
```

```
Check {\tt ?anova.lm}
```

Important note: When comparing two models using anova the results are as expected from the partial F -test. However, when more than two models are compared using anova, the F -statistic and p -value may not be what we would like. The reason for this is that the F -statistic compares the mean SS for a row to the MS_{ERR} for the largest model considered.

> `anova(fit1, fit3, fit2)`

be Careful in interpretation.

Analysis of Variance Table

Model 1: sale ~ bedrooms + neighborhood

Model 2: sale ~ bedrooms + neighborhood + half

Model 3: sale ~ bedrooms + neighborhood + full + half

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	26	58164				
2	25	42194	1	15970.1	10.6132	0.003338 **
3	24	36114	1	6080.1	4.0406	0.055785 .

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘

> `anova(fit1, fit2)`

Analysis of Variance Table

Model 1: sale ~ bedrooms + neighborhood

reduced model

$H_0: \beta_{full} = \beta_{half} = 0$ vs $H_1: \text{At least one is } \neq 0$.

Model 2: sale ~ bedrooms + neighborhood + full + half

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
26	58164				P-value of
2	36114	(2)	22050	7.3269	0.003283 **
---	d.f.	num d.f.	TSS _{EX}	(S _{EX} /2)	MS _{ERR} (full)
Signif. codes:	0	'***'	0.001	'**'	0.01 '*' 0.05 '.' 0.1 '
>					

> anova(fit3, fit2)

Analysis of Variance Table

reduced

Model 1: sale ~ bedrooms + neighborhood + half

Model 2: sale ~ bedrooms + neighborhood + full + half

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
25	42194				same n)
24	36114	1	6080.1	4.0406	0.05579 → the P-value for t-test.

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

>

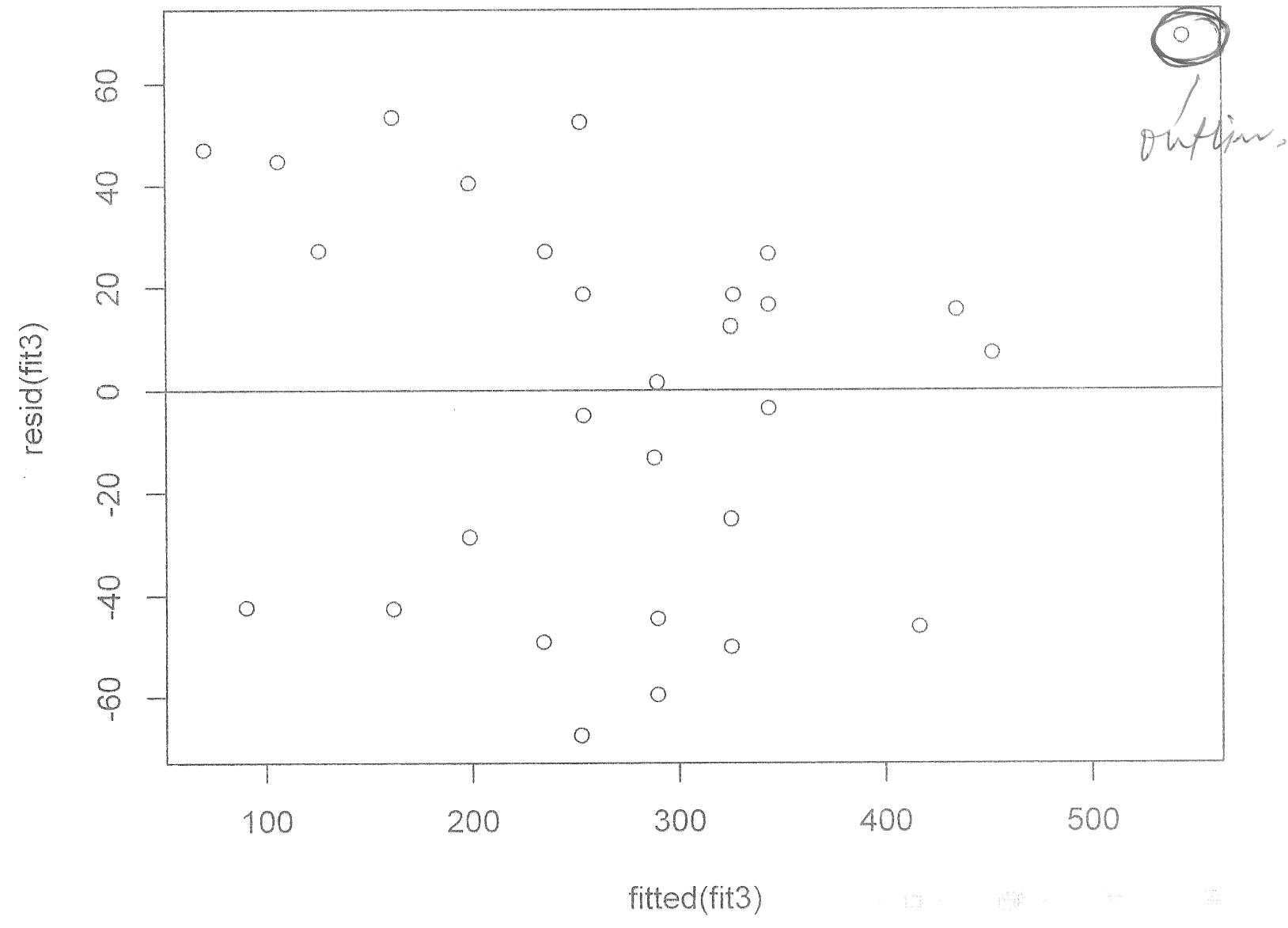
```
# Residual plot
```

```
plot(fitted(fit3), resid(fit3))
abline(h=0)
```

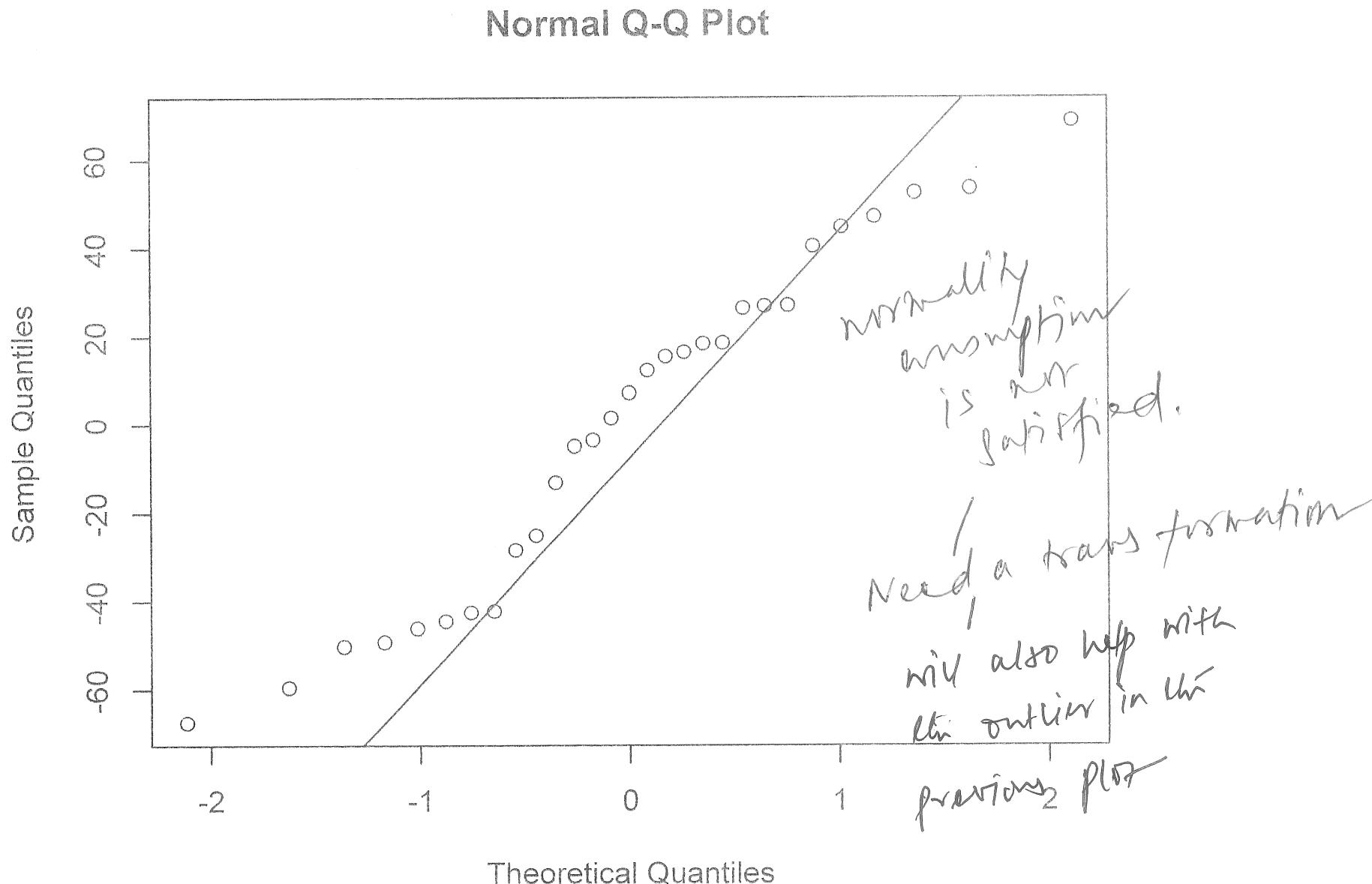
```
# QQ plot
```

```
qqnorm(resid(fit3))
qqline(resid(fit3))
```

Residual plot



Normal QQ plot



```
# Take sqrt(sale) rather than sale as response  
  
fit4 <- update(fit3, sqrt(sale) ~ .)  
  
# New QQ plot  
  
qqnorm(resid(fit4))  
qqline(resid(fit4))
```

sqrt transformation of sale

Normal QQ plot for transformed data

