lmod <- lm(Life.Exp ~ ., statedata)</pre>

Lab4

Lab 4

```
data(state)
statedata <- data.frame(state.x77, row.names = state.abb)</pre>
head(statedata)
```

	Population <dbl></dbl>	Income <dbl></dbl>	Illiteracy <dbl></dbl>	Life.Exp <dbl></dbl>	Murder <dbl></dbl>	HS.Grad <dbl></dbl>	Frost <dbl></dbl>	Area <dbl></dbl>
AL	3615	3624	2.1	69.05	15.1	41.3	20	50708
AK	365	6315	1.5	69.31	11.3	66.7	152	566432
AZ	2212	4530	1.8	70.55	7.8	58.1	15	113417
AR	2110	3378	1.9	70.66	10.1	39.9	65	51945
CA	21198	5114	1.1	71.71	10.3	62.6	20	156361
CO	2541	4884	0.7	72.06	6.8	63.9	166	103766
6 rows								

summary(lmod) ## Call: ## lm(formula = Life.Exp ~ ., data = statedata)

Can use the . to indicate to include all the other variables in your model

```
## Residuals:
        Min
                 1Q Median
                                  3Q
                                         Max
 ## -1.48895 -0.51232 -0.02747 0.57002 1.49447
 ## Coefficients:
               Estimate Std. Error t value Pr(>|t|)
 ## (Intercept) 7.094e+01 1.748e+00 40.586 < 2e-16 ***
 ## Population 5.180e-05 2.919e-05 1.775 0.0832 .
               -2.180e-05 2.444e-04 -0.089
                                          0.9293
 ## Income
 ## Illiteracy 3.382e-02 3.663e-01 0.092 0.9269
              -3.011e-01 4.662e-02 -6.459 8.68e-08 ***
 ## Murder
 ## HS.Grad
            4.893e-02 2.332e-02 2.098
                                          0.0420 *
 ## Frost
              -5.735e-03 3.143e-03 -1.825 0.0752
              -7.383e-08 1.668e-06 -0.044 0.9649
 ## Area
 ## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
 ## Residual standard error: 0.7448 on 42 degrees of freedom
 ## Multiple R-squared: 0.7362, Adjusted R-squared: 0.6922
 ## F-statistic: 16.74 on 7 and 42 DF, p-value: 2.534e-10
Accuracy of the model
```

summary output R^2 summary(lmod)\$r.squared

[1] 0.7361563

 $R^{2} = 1 - \frac{\sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}} = 1 - \frac{SSR}{SST}$

```
# calculate R^2 by hand
y <- statedata$Life.Exp # Response values
y_hat <- fitted(lmod) # Fitted Values</pre>
y_bar <- mean(y)</pre>
SSR <-sum((y - y_hat)^2)
SST \leftarrow sum((y - y_bar)^2)
r 2 \leftarrow 1 - SSR/SST
r_2
## [1] 0.7361563
# R^2=cor(y_hat,y)^2
cor(y_hat,y)^2
## [1] 0.7361563
```

```
Hypothesis Testing
T-test
```

• the t-test here can test the significance for a single covariate β_i . • The hypothesis here is: $H_0: \beta_i = 0 \text{ vs } H_1: \beta_i \neq 0$. If corresponding p-value is less than significane level (say, 0.05), we reject H_0 .

Income

[1] 2.097882

[1] 0.04197175

n <- dim(statedata)[1] # number of observations, or equivalently use nrow(statedata)</pre> p <- 7 # number of predictors

Is a specific predictor useful in predicting Y? t-test

round(coefficients(summary(lmod)), 5)

Illiteracy 0.03382 0.36628 0.09233 0.92687

-0.00002 0.00024 -0.08921 0.92934

Estimate Std. Error t value Pr(>|t|)## (Intercept) 70.94322 1.74798 40.58594 0.00000 ## Population 0.00005 0.00003 1.77477 0.08318

```
-0.30112 0.04662 -6.45900 0.00000
 ## Murder
 ## HS.Grad 0.04893 0.02332 2.09788 0.04197
 ## Frost -0.00574
                          0.00314 -1.82456 0.07519
             0.00000
 ## Area
                          0.00000 - 0.04426 \quad 0.96491
Let's double check HS.Grad t-value and p-value
 # summary output t - value
 coefficients(summary(lmod))[6,3]
```

t value <- coefficients(summary(lmod))[6,1]/coefficients(summary(lmod))[6,2] t value

calculate t - value by hand

```
## [1] 2.097882
# summary output p - value
coefficients(summary(lmod))[6,4]
```

```
# calculate p - value by hand
p_value = pt(q = -t_value, df = n - p - 1) * 2
p_value
## [1] 0.04197175
```

```
• The F-test can test the difference between two nested models. nested models mean that: Model_1 = f(covariates_1) and
     Model_2 = f(covariates_2), we have covariates_1 \subset covariates_2.
   • The hypothesis here is: H_0: Model_1 is same as Model_2 vs H_1: Model_1 is not same as Model_2
Global F test
```

mod_M <- lm(Life.Exp ~ ., statedata) # Larger model with all the predictors</pre> mod_m <- lm(Life.Exp ~ 1, statedata) # Smaller model with only intercept</pre>

2 rows

Coefficients:

Population

Illiteracy

Income

Murder

HS.Grad

F-test

anova(mod_m, mod_M) # Global F - Test Res.Df RSS Df Sum of Sq <dbl> <dbl> <dpl> <dpl> <dbl>

Pr(>F)

<dbl>

Pr(>F)

<dbl>

0.9992586

.resid

<dbl>

3.853661

-35.962052

-31.561556

-26.511670

NA

<dbl>

0.00655296

NA

.fitted

<dbl>

1025.1463

969.9621

975.5616

1031.5117

takers

takers

2.534328e-10

NA

NA NA 49 88.29900 NA 7 42 23.29714 65.00186 16.74073

summary(mod_M)
##
Call:
Call:
lm(formula = Life.Exp ~ ., data = statedata)
##
Residuals:

```
-5.735e-03 3.143e-03 -1.825
 ## Frost
                                                  0.0752 .
                 -7.383e-08 1.668e-06 -0.044
                                                  0.9649
 ## Area
 ## Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
 ## Residual standard error: 0.7448 on 42 degrees of freedom
 ## Multiple R-squared: 0.7362, Adjusted R-squared: 0.6922
 ## F-statistic: 16.74 on 7 and 42 DF, p-value: 2.534e-10
Partial F Test
   • Want to test the null hypothesis that \beta_{HS.Grad} = \beta_{Frost} = 0
 mod_M <- lm(Life.Exp ~ ., statedata) # Larger model with all the predictors</pre>
 mod_m <- lm(Life.Exp ~ Population +</pre>
 Income +
 Illiteracy +
 Murder +
 Area, statedata) # smaller model without HS.Grad and Frost
 anova(mod_m, mod_M)
                                   RSS
                                                                                          F
             Res.Df
                                             Df
                                                             Sum of Sq
                                                                                                                  Pr(>F)
              <dbl>
                                          <dbl>
                                                                  <dbl>
                                  <dbl>
                                                                                      <dbl>
                                                                                                                   <dbl>
                 44
                                29.30301
                                            NA
                                                                   NA
                                                                                        NA
                                                                                                                     NA
                 42
                                23.29714
                                              2
                                                               6.005868
                                                                                   5.413679
                                                                                                             0.008094657
 2 rows
```

```
2 rows
The lab uses the now-familiar SAT data throughout.
data from faraway on 1998 per capita income for each U.S. state and the proportion of residents of each state born in the U.S. as of the 1990
census. By now, this should be a familiar plot:
```

Sum of Sq

0.01090466

<dbl>

NA

Arizona 4.778 19.3 32.175 4.459 17.1 28.934 Arkansas 4 rows | 1-10 of 14 columns

Recall from lecture that the three 'classic' diagnostic plots are:

1. Residuals versus fitted values (check overall linearity)

3. Quantile-quantile plot (Normality assumption)

ggplot(aes(y = .resid, x = value)) +

geom_hline(aes(yintercept = 0))

augment(fit naive, sat) %>%

geom point() +

So let's add that term:

add quadratic term in expenditure

Rationales behind QQ plot: page 93 - 95

ggplot(aes(y = .resid, x = value)) +

geom_hline(aes(yintercept = 0)) +

facet_wrap(~ name, scales = 'free_x') +

pivot_longer(cols = c(.fitted, takers, expend)) %>%

.fitted

that sometimes non-significant predictors should still be included in a model.

fit <- lm(total ~ poly(expend, 2, raw = T) + poly(takers, 2, raw = T), data = sat)</pre>

 $geom_smooth(method = 'loess', formula = 'y ~ x', se = F, span = 1)$

facet_wrap(~ name, scales = 'free_x') +

2. Residuals versus predictors (check linearity w.r.t. that predictors)

pivot_longer(cols = c(.fitted, takers, expend)) %>%

expend

<dpl>

4.405

8.963

augment(fit_naive, sat) %>% head(4)

.rownames

<chr>

Alabama

More on residual plot

panel of residual plots augment(fit_naive, sat) %>%

geom_point() +

Alaska

- It is convenient to present plots (1) and (2) in a panel, since all of these are scatterplots with the residuals on the y axis. To do so, pivot the predictors and the fitted values, and then facet. Adding a horizontal line at zero helps, as ideally we'd like to see the residuals spread evenly around that line.
 - .fitted expend 50 **-**
- 850 900 950 1000 1050 10 40 60 80 value Notice that the non-linearity in takers appears as a pattern in both the residual-fit plot and the residual-predictor plot. Ostensibly, the leftmost panel indicates there is some nonlinearity, and then the rightmost panel points to which variable is the culprit. It won't always work out so nicely, but sometimes these patterns are really clear and unambigous – there's definitely a parabolic shape to the residuals in takers, so that predictor should probably enter quadratically into the model. Sometimes it can be useful to add a smoothed trend line to help visualize the pattern: # sometimes a smoother helps (but beware the span!)

50 resid -50 **-**1000 1050 850 950 10

value

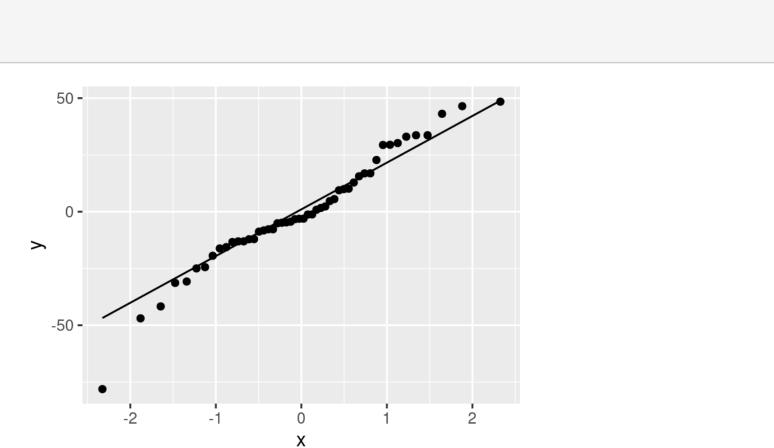
It's perhaps questionable, but we can observe a slight parabolic pattern to expend. You may recall adding this term and finding that it was not a

statistically significant predictor; however, despite that, it does appear to be needed for correct model specification. It's interesting to consider

expend

Let's assume that there are no additional problems in these plots ((1) and (2)). Given, then, that the model appears adequately specified and there are no obvious problems with the constant variance assumption, we can check the normality assumption. The quantile-quantile plot is simple to construct:

```
geom_qq() +
geom qq line()
                                     50 -
```



 $(F_value \leftarrow ((88.299 - 23.297)/7)/(23.297/42))$ ## [1] 16.74087

-1.48895 -0.51232 -0.02747 0.57002 1.49447

5.180e-05 2.919e-05

3.382e-02 3.663e-01

4.893e-02 2.332e-02

(Intercept) 7.094e+01 1.748e+00 40.586 < 2e-16 ***

-2.180e-05 2.444e-04 -0.089

• Now lets test the null hypothesis that $eta_{income} = eta_{Area} = eta_{Illiteracy} = 0$

mod_m <- lm(Life.Exp ~ Population +</pre>

Murder + HS.Grad +

anova(mod m, mod M)

Res.Df

<dpl>

45

42

mod_M <- lm(Life.Exp ~ ., statedata) # Larger model with all the predictors</pre>

Frost, statedata) # smaller model without Income, Area, and Illiteracy

RSS

<dbl>

23.30804

23.29714

augment function adds residuals, fitted, and case influence stats

salary

<dbl>

31.144

47.951

ratio

<dbl>

17.2

17.6

Df

NA

3

<dpl>

Estimate Std. Error t value Pr(>|t|)

-3.011e-01 4.662e-02 -6.459 8.68e-08 ***

1.775

0.092

2.098

• Is at least one of the predictors useful in predicting Y? F-test

```
## Residuals:
       Min
                1Q Median
                                  3Q
```

0.0832 . 0.9293

0.9269

0.0420 *

Checking model assumptions (a.k.a. Diagnosis of linea nodel)
 Assumptions for linear model (ordered by importance): The model is correct, i.e. linearity is satisfied (The residuals "bounce randomly" around the 0 line) response are uncorrelated equal variance (The residuals roughly form a "horizontal band" around the 0 line.) normality (qq-plot)
naive fit maybe linear is good enough fit_naive <- lm(total ~ takers + expend, data = sat)

takers

<int>

8

47

27

6

verbal

<int>

491

445

448

482

math

<int>

538

489

496

523

total

<int>

1029

934

944

1005

