INTRODUCTION TO DATA ANALYSIS



OUTLINE

fill me





Case study murder rates

MURDER RATES

murder_data

GGally::ggpairs(murder_data, title = "Murder rate data")

			x 4	A tibble: 20	# #	##
	population	unemployment	low_income	murder_rate	ŧ	##
0.03-	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	¢	##
0.00	587000	6.2	16.5	11.2	# 1	##
0.02 -	643000	6.4	20.5	13.4	¢ 2	##
0.01 -	635000	9.3	26.3	40.7	¢ 3	##
0.00 -	692000	5.3	16.5	5.3	# 4	##
	1248000	7.3	19.2	24.8	¢ 5	##
24 -	643000	5.9	16.5	12.7	≠ 6	##
21-	1964000	6.4	20.2	20.9	¢ 7	##
18 -	1531000	7.6	21.3	35.7	¢ 8	##
15 -	713000	4.9	17.2	8.7	¢ 9	##
9 -	749000	6.4	14.3	9.6	¢ 10	##
8-	7895000	6	18.1	14.5	<i>†</i> 11	##
7 -	762000	7.4	23.1	26.9	<i>‡</i> 12	##
6-	2793000	5.8	19.1	15.7	# 13	##
5-	741000	8.6	24.7	36.2	# 14	##
8e+06-	625000	6.5	18.6	18.1	# 15	##
6e+06-	854000	8.3	24.9	28.9	# 16	##
10+06-	716000	6.7	17.9	14.9	# 17	##
46100	921000	8.6	22.4	25.8	¢ 18	##
2e+06 -	595000	8.4	20.2	21.7	<i>†</i> 19	##
	3353000	6.7	16.9	25.7	<i>‡</i> 20	##

Murder rate data

9 -

7 -

10 20



annual murders per million inhabitants

percentage inhabitants with low income

percentage inhabitants who are unemployed

total population





linear regression what? why? how?

MURDER RATE IN EACH CITY



We know the murder rates of all cities. A city is randomly drawn from the data set (or the population from which the data is sampled). We are to predict the murder rate of that randomly drawn city. What's a best guess if we know nothing more about that city?

The mean of all murder rates!

20



MURDER RATE IN EACH CITY WITH GRAND MEAN AS PREDICTOR



MURDER RATE IN EACH CITY WITH GRAND MEAN AS PREDICTOR



MURDER RATE IN EACH CITY WITH GRAND MEAN AS PREDICTOR





```
y <- murder_data %>% pull(murder_rate)
n <- length(y)</pre>
tss_simple <- sum((y - mean(y))^2)</pre>
tss_simple
```

[1] 1855.202





MURDER RATE IN EACH CITY AS A FUNCTION OF UNEMPLOYMENT RATE



9

As before, we are to predict the murder rate of a randomly drawn city. But now **we also get to know that city's unemployment rate.** What's a best guess if we know nothing more about that city?

Let's just assume the following linear relationship:

$$\hat{y}_i = 2x_i + 4$$



MURDER RATE IN EACH CITY AS A FUNCTION OF UNEMPLOYMENT RATE



Overal prediction error: $RSS = \sum_{i=1}^{N} (y_i - \hat{y})^2$ i=1[residual sum of squares]

```
y <- murder_data %>% pull(murder_rate)
x <- murder_data %>% pull(unemployment)
predicted_y <- 2 * x + 4
n <- length(y)</pre>
rss_guesswork <- sum((y - predicted_y)^2)</pre>
rss_guesswork
```

[1] 1327.74



SIMPLE LINEAR REGRESSION [ORDINARY LEAST-SQUARES REGRESSION]

linear prediction function

prediction error

best-fitting parameters





 $\hat{y}_i = \beta_0 + \beta_1 x_i$

 $\mathbf{RSS}_{\langle \beta_0, \beta_1 \rangle} = \sum_{i=1}^{k} (y_i - \hat{y}_i)^2$ i=1

 $\langle \hat{\beta}_0, \hat{\beta}_1 \rangle = \arg \min_{\langle \beta_0, \beta_1 \rangle} \mathsf{RSS}_{\langle \beta_0, \beta_1 \rangle}$ $\langle \beta_0, \beta_1 \rangle$

BEST LINEAR PREDICTOR





finding best-fits for OLS regression

USING `OPTIM`

```
# data to be explained / predicted
y <- murder_data %>% pull(murder_rate)
                                                             ##
# data to use for prediction / explanation
                                                             ##
x <- murder_data %>% pull(unemployment)
# function to calculate residual sum of squares
get_rss = function(y, x, beta_0, beta_1) {
  yPred = beta_0 + x * beta_1
  sum((y-yPred)^2)
# finding best-fitting values for TSS
fit_rss = optim(par = c(0, 1),
  fn = function(par) {
                                                             rate
    get_rss(y, x, par[1], par[2])
                                                            20-
# output the results
message(
  "Best fitting parameter values:",
  "\n\tIntercept: ", fit_rss$par[1] %>% signif(5),
  "\n\tSlope: ", fit_rss$par[2] %>% signif(5),
  "\nRSS for best fit: ", fit_rss$value %>% signif(5)
```

Best fitting parameter values:

- Intercept: -28.528
- Slope: 7.0795

RSS for best fit: 467.6



USING `LM`

```
# fit an OLS regression
fit_lm <- lm(
 # the formula argument specifies dependent and independent variables
 formula = murder_rate ~ unemployment,
 # we also need to say where the data (columns) should come from
 data = murder_data
# output the fitted object
fit_lm
```

##

Call:

```
## lm(formula = murder_rate ~ unemployment, data = murder_data)
```

##

Coefficients:

(Intercept) unemployment ##

-28.53 ## 7.08







USING `LM`

summary(fit_lm)

Call: ## lm(formula = murder_rate ~ unemployment, data = murder_data) ## ## Residuals: 10 Median ## Min Max ЗQ ## -9.2415 -3.7728 0.5795 3.2207 10.4221 ## ## Coefficients: Estimate Std. Error t value Pr(>|t|) ## ## (Intercept) -28.5267 6.8137 -4.187 0.000554 *** ## unemployment 7.0796 7.309 8.66e-07 *** 0.9687 ## ---## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 ## ## Residual standard error: 5.097 on 18 degrees of freedom ## Multiple R-squared: 0.748, Adjusted R-squared: 0.7339 ## F-statistic: 53.41 on 1 and 18 DF, p-value: 8.663e-07



USING MATH

Theorem 3.1 (OLS solution) For a simple linear

regression model with just one predictor, the solution for:

$$rgmin_{\langleeta_0,eta_1
angle}\sum_{i=1}^k(y_i-(eta_0+eta_1x_i))^2$$

is given by:

$$\hat{eta_1} = rac{Cov(x,y)}{Var(x)} \qquad \hat{eta_0} = ar{y} - \hat{eta_1}ar{x}$$

```
tibble(
    beta_1 = cov(x,y) / var(x),
    beta_0 = mean(y) - beta_1 * mean(x)
)
```

```
## # A tibble: 1 x 2
## beta_1 beta_0
## <dbl> <dbl>
## 1 7.08 -28.5
```





finding best-fits for a likelihood-based approach

FREQUENTIST MODEL [LIKELIHOOD-BASED APPROACH]

 $egin{aligned} ext{ikelihood-based regression} \ & [ext{explicit version}] \ & y_{ ext{pred}} = eta_0 + eta_1 x \ & y_i = y_{ ext{pred}} + \epsilon_i \ & \epsilon_i \sim ext{Normal}(0, \sigma) \end{aligned}$

 $egin{aligned} \mathbf{ikelihood-based\ regression}\ & [\mathbf{compact\ version}]\ & y_{\mathrm{pred}} &= eta_0 + eta_1 x\ & y \sim \mathrm{Normal}(\mu = y_{\mathrm{pred}}, \sigma) \end{aligned}$



FITTING WITH `OPTIM`

```
# data to be explained / predicted
y <- murder_data %>% pull(murder_rate)
# data to use for prediction / explanation
x <- murder_data %>% pull(unemployment)
# function to calculate negative log-likelihood
get_nll = function(y, x, beta_0, beta_1, sd) {
  if (sd <= 0) {return( Inf )}</pre>
  yPred = beta_0 + x * beta_1
  nll = -dnorm(y, mean=yPred, sd=sd, log = T)
  sum(nll)
# finding MLE
fit_lh = optim(par = c(0, 1, 1)),
  fn = function(par) {
   get_nll(y, x, par[1], par[2], par[3])
# output the results
message(
  "Best fitting parameter values:",
  "\n\tIntercept: ", fit_lh$par[1] %>% signif(5),
  "\n\tSlope: ", fit_lh$par[2] %>% signif(5),
  "\nNegative log-likelihood for best fit: ", fit_lh$value %>% signif(5)
```

```
Best fitting parameter values:
Intercept: -28.517
Slope: 7.0783
Negative log-likelihood for best fit: 59.898
```



FITTING WITH `GLM`

```
fit_glm <- glm(murder_rate ~ unemployment, data = murder_data)
fit_glm</pre>
```

##	
##	Call: glm(formula = murder_rate ~ unemployment, data = murd
##	
##	Coefficients:
##	(Intercept) unemployment
##	-28.53 7.08
##	
##	Degrees of Freedom: 19 Total (i.e. Null); 18 Residual
##	Null Deviance: 1855
##	Residual Deviance: 467.6 AIC: 125.8





FITTING WITH `GLM`

```
summary(fit_glm)
```

##							
##	Call:						
##	glm(formula = murder_rate ~ unemployment, data = murder_da						
##							
##	Deviance Re	siduals	5:				
##	Min	1Q	Median	3Q	Max		
##	-9.2415 -3	7728	0.5795	3.2207	10.4221		
##							
##	Coefficients:						
##		Estin	nate Std.	Error t	value Pr(> t)	
##	(Intercept)	-28.5	267	6.8137 -	4.187 0.0	00554 ***	
##	unemploymen	t 7.0	796	0.9687	7.309 8.6	6e-07 ***	
##							
##	Signif. cod	es: 0	'***' 0.	001 '**'	0.01 '*'	0.05 '.'	0.1 '



CONNECTION OLS & MLE

Theorem 3.2 (MLE solution) For a simple linear regression model with just

one predictor, the solution for:

$$rg\max_{\langleeta_0,eta_1,\sigma
angle} \prod_{i=1}^k \mathrm{Normal}(\mu=eta_0+eta_1x_i,\sigma)$$

is the same as for the OLS approach:

$$\hat{eta_1} = rac{Cov(x,y)}{Var(x)} \qquad \hat{eta_0} = ar{y} - \hat{eta}_1 ar{x}$$



Bayesian approach

BAYESIAN SIMPLE LINEAR REGRESSION MODEL



- $\sigma \sim \text{Trunc-Norm}(...)$
- $\beta_0 \sim \text{Student-t}(\ldots)$
- $\beta_1 \sim \text{Student-t}(\ldots)$
- $y_i \sim \text{Normal}(\beta_0 + \beta_1 x_i, \sigma)$

data to be explained / predicted y <- murder_data %>% pull(murder_rate) # data to use for prediction / explanation x <- murder_data %>% pull(unemployment) <- as_data(y) y_greta x_greta <- as_data(x) # latent variables and priors intercept <- student(df= 1, mu = 0, sigma = 10)</pre> slope <- student(df= 1, mu = 0, sigma = 10) <- normal(0, 5, truncation = c(0, Inf)) sigma # derived latent variable (linear model) y_pred <- intercept + slope * x_greta</pre> # likelihood distribution(y) <- normal(y_pred, sigma)</pre> # finalize model, register which parameters to monitor murder_model <- model(intercept, slope, sigma)</pre>



BAYESIAN SIMPLE LINEAR REGRESSION MODEL





Hypothesis tests for regression coefficients

BAYESIAN APPROACH [ESTIMATION-BASED]

```
# get means and 95% HDI
Bayes_estimates <- tidy_draws_murder_data %>%
group_by(Parameter) %>%
summarise(
    mean = mean(value),
    '|95%' = HDInterval::hdi(value)[1],
    '95|%' = HDInterval::hdi(value)[2]
)
Bayes_estimates
```

##	#	A tibble:	3 x 4		
##		Parameter	mean	` 95%`	`95 %`
##		<fct></fct>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
##	1	intercept	-24.6	-39.1	-9.51
##	2	sigma	5.36	3.78	7.20
##	3	slope	6.53	4.34	8.53



FREQUENTIST T-TESTS FOR REGRESSION COEFFICIENTS



$$\frac{\operatorname{Cov}(x, y)}{\operatorname{Var}(y)} \qquad \qquad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Null-hypothesi
$$\beta_0 = \beta_1 = 0$$

$$\frac{\sum_{i} (y_{i} - (\hat{\beta}_{0} + \hat{\beta}_{1}x_{i}))^{2}}{(N-2) \cdot \sum_{i} (x_{i} - \bar{x})^{2}}$$

$$= \sqrt{\sum_{i}^{i} (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2 \cdot \left(\frac{1}{(N-2)} + \frac{\bar{x}^2}{\sum_{i}^{i} (x_i - \bar{x})^2}\right)}$$

$$\hat{\beta}_1$$

 E_{β_1}

$$t_0 = \frac{\hat{\beta}_0}{\mathrm{SE}_{\beta_0}}$$

Sampling distribution: $t_0, t_1 \sim \text{Student-t}(\nu = N - 2)$





FREQUENTIST T-TESTS FOR REGRESSION COEFFICIENTS

```
# observed data
y_obs <- murder_data %>% pull(murder_rate)
x_obs <- murder_data %>% pull(unemployment)
n_obs <- length(y_obs)</pre>
# best-fitting coefficients
beta_1_hat <- cov(x_obs, y_obs) / var(x_obs)</pre>
beta 0 hat <- mean(y obs) - beta 1 hat * mean(x obs)</pre>
# calculating t-scores
MSE <- sum((y_obs - beta_0_hat - beta_1_hat * x_obs)^2) / (n_obs-2)</pre>
S_x < - sum((x_obs - mean(x_obs))^2)
SE_beta_1_hat <- sqrt(MSE / S_xx)</pre>
SE_beta_0_hat <- sqrt(MSE * (1/n_obs + mean(x_obs)^2 / S_xx))</pre>
t_slope = (beta_1_hat) / SE_beta_1_hat
t_intercept = beta_0_hat / SE_beta_0_hat
tibble(t_slope, t_intercept)
```

A tibble: 1 x 2
t_slope t_intercept
<dbl> <dbl>
1 7.31 -4.19

p_value_intercept = pt(t_intercept, df = n_obs -2) + 1-pt(-t_intercept, df = n_obs -2)
p_value_slope = pt(-t_slope, df = n_obs -2) + 1-pt(t_slope, df = n_obs -2)
tibble(p_value_intercept, p_value_slope)

summary(glm(murder_rate ~ unemployment, data = murder_data))

Coefficients:

Estimate Std. Error t value Pr(>|t|)
(Intercept) -28.5267 6.8137 -4.187 0.000554 ***
unemployment 7.0796 0.9687 7.309 8.66e-07 ***
---## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

