#### Modeling Concepts: Inference vs Prediction Justin Post

#### What do we want to be able to do?

Data Science!

- Read in raw data and manipulate it
- Combine data sources
- Summarize data to glean insights
- Apply common analysis methods
- Communicate Effectively

## Modeling Ideas

What is a (statistical) model?

- A mathematical representation of some phenomenon on which you've observed data
- Form of the model can vary greatly!

##	#	A tibble: 1,061 x 9						
##		<pre>log_km_driven log_s</pre>	elling_price	name	<pre>selling_price</pre>	year	<pre>seller_type</pre>	owner
##		<dbl></dbl>	<dbl></dbl>	<chr></chr>	<dbl></dbl>	<dbl></dbl>	<chr></chr>	<chr></chr>
##	1	5.86	12.1	Royal E~	175000	2019	Individual	1st ~
##	2	8.64	10.7	Honda D~	45000	2017	Individual	1st ~
##	3	9.39	11.9	Royal E~	150000	2018	Individual	1st ~
##	4	10.0	11.1	Yamaha ~	65000	2015	Individual	1st ~
##	5	9.95	9.90	Yamaha ~	20000	2011	Individual	2nd ~
##	#	i 1,056 more rows						
##	#	i 2 more variables:	km_driven <	dbl>, ex_s	showroom_price	<dbl></dbl>		

```
ggplot(bike_data, aes(x = log_km_driven, y = log_selling_price)) +
geom_point() +
geom_smooth(method = "lm")
```

```
## `geom_smooth()` using formula = 'y ~ x'
```



```
ggplot(bike_data, aes(x = log_km_driven, y = log_selling_price)) +
geom_point() +
stat_smooth(method = "lm", formula = y ~ x + I(x^2))
```



```
ggplot(bike_data, aes(x = log_km_driven, y = log_selling_price)) +
geom_point() +
stat_smooth(method = "lm", formula = y ~ x + I(x^2) + I(x^3))
```



• First a visual on motorcycle sales data

```
ggplot(bike_data, aes(x = log_km_driven, y = log_selling_price)) +
  geom_point() +
  geom_smooth()
```

## `geom\_smooth()` using method = 'gam' and formula = 'y ~ s(x, bs = "cs")'



## Warning: package 'tree' was built under R version 4.1.3

```
ggplot(bike_data, aes(x = log_km_driven, y = log_selling_price)) +
geom_point() +
geom_line(data = preds, aes(x = log_km_driven, y = Tree_Predictions), color = "Blue", linewidth = 2)
```



## Modeling Ideas

What is a (statistical) model?

- A mathematical representation of some phenomenon on which you've observed data
- Form of the model can vary greatly!

**Statistical learning** - Inference, prediction/classification, and pattern finding

- Supervised learning a variable (or variables) represents an **output** or **response** of interest
  - $\circ~$  May model response and
    - Make **inference** on the model parameters
    - predict a value or classify an observation

Our Goal: Understand what it means to be a good predictive model (not make inference)

## Training a Model

• Once a class of models is chosen, we must define some criteria to **fit** (or train) the model

**Simple Linear Regression (SLR) Model** 

 $E(Y_i|x_i)=eta_0+eta_1x_i$ 

## Training a Model

• Once a class of models is chosen, we must define some criteria to **fit** (or train) the model

**Simple Linear Regression (SLR) Model** 

$$E(Y_i|x_i)=eta_0+eta_1x_i$$

• Loss function - Criteria used to fit or train a model

 $\circ\,$  For a given **numeric** response value,  $y_i$  and prediction,  ${\hat y}_i$ 

$${y_i} - {{\hat y}_i},({y_i} - {{\hat y}_i})^2,|{y_i} - {{\hat y}_i}|$$

## Training a Model

• Once a class of models is chosen, we must define some criteria to **fit** (or train) the model

**Simple Linear Regression (SLR) Model** 

$$E(Y_i|x_i) = eta_0 + eta_1 x_i$$

• Loss function - Criteria used to fit or train a model

 $\circ\,$  For a given **numeric** response value,  $y_i$  and prediction,  ${\hat y}_i$ 

$${y_i} - {{\hat y}_i},({y_i} - {{\hat y}_i})^2,|{y_i} - {{\hat y}_i}|$$

• We try to optimize the loss over all the observations used for training

$$\sum_{i=1}^n (y_i - {\hat y}_i)^2 \qquad \qquad \sum_{i=1}^n |y_i - {\hat y}_i|$$

## Training (Fitting) the SLR Model

- Often use squared error loss (least squares regression)
- Nice solutions for our estimates exist!

$$\hat{eta}_0 = ar{y} - ar{x} \hat{eta}_1 \ \hat{eta}_1 = rac{\sum_{i=1}^n (x_i - ar{x})(y_i - ar{y})}{\sum_{i=1}^n (x_i - ar{x})^2}$$

## Training (Fitting) the SLR Model

- Often use squared error loss (least squares regression)
- Nice solutions for our estimates exist!

$${\hat eta}_0 = ar y - ar x {\hat eta}_1$$

$${\hat eta}_1 = rac{\sum_{i=1}^n (x_i - ar x)(y_i - ar y)}{\sum_{i=1}^n (x_i - ar x)^2}$$

y <- bike\_data\$log\_selling\_price x <- bike\_data\$log\_km\_driven b1\_hat <- sum((x-mean(x))\*(y-mean(y)))/sum((x-mean(x))^2) b0\_hat <- mean(y)-mean(x)\*b1\_hat c(round(b0\_hat, 4), round(b1\_hat, 4))

## [1] 14.6356 -0.3911

- Now we can find a prediction! Denoted as  $\hat{y}$ 

## Training (Fitting) the SLR Model in R

- Use lm() function to fit in R
- Utilizes formula notation: y ~ x -> response ~ model terms

```
slr_fit <- lm(log_selling_price ~ log_km_driven, data = bike_data)
slr_fit
##
## Call:
## lm(formula = log_selling_price ~ log_km_driven, data = bike_data)
##
## Coefficients:
## (Intercept) log_km_driven
## 14.6356 -0.3911</pre>
```

- If we assume iid errors that are Normally distributed with the same variance, we can conduct inference!
  - Confidence intervals and hypothesis tests around the slope parameter

• Use summary() (generic function!) on our fitted model

```
summary(slr_fit)
##
## Call:
## lm(formula = log_selling_price ~ log_km_driven, data = bike_data)
##
## Residuals:
                                                              Min
                                                                                                                                            10 Median
##
                                                                                                                                                                                                                                                                                        30
                                                                                                                                                                                                                                                                                                                                                    Max
## -1.9271 -0.3822 -0.0337 0.3794 2.5656
##
## Coefficients:
                                                                                                                                                    Estimate Std. Error t value Pr(>|t|)
##
## (Intercept) 14.63557 0.18455 79.31 <2e-16 ***
## log_km_driven -0.39109 0.01837 -21.29 <2e-16 ***</pre>
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
 ## Residual standard error: 0.5953 on 1059 degrees of freedom
    ++ \mathbf{M} \cdot \mathbf{I} + \mathbf{M} \cdot \mathbf{M} \cdot \mathbf{I} + \mathbf{M} \cdot \mathbf{M} + \mathbf{M} \cdot \mathbf{M} + \mathbf{M} \cdot \mathbf{M} + \mathbf{M}
```

- If we assume iid errors that are Normally distributed with the same variance, we can conduct inference!
  - Can use anova() to get Analysis of Variance information

anova(slr\_fit)
## Analysis of Variance Table
##
## Response: log\_selling\_price
## Df Sum Sq Mean Sq F value Pr(>F)
## log\_km\_driven 1 160.58 160.583 453.19 < 2.2e-16 \*\*\*
## Residuals 1059 375.24 0.354
## --## Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1</pre>

- If we assume iid errors that are Normally distributed with the same variance, we can conduct inference!
  - Residual diagnostics can be found via plot() on the fitted model



- If we assume iid errors that are Normally distributed with the same variance, we can conduct inference!
  - Residual diagnostics can be found via plot() on the fitted model



#### Prediction Using the SLR Model in R

• Can use the line for prediction with predict()!

• Another generic function in R

predict

```
## function (object, ...)
## UseMethod("predict")
## <bytecode: 0x000000015bef250>
## <environment: namespace:stats>
 predict.lm
## function (object, newdata, se.fit = FALSE, scale = NULL, df = Inf,
       interval = c("none", "confidence", "prediction"), level = 0.95,
##
       type = c("response", "terms"), terms = NULL, na.action = na.pass,
##
##
       pred.var = res.var/weights, weights = 1, ...)
## {
       tt <- terms(object)</pre>
##
       if (!inherits(object, "lm"))
##
##
           warning("calling predict.lm(<fake-lm-object>) ...")
       if (missing(newdata) || is.null(newdata)) {
##
##
           mm <- X <- model.matrix(object)</pre>
##
           mmDone <- TRUE</pre>
##
           offset <- object$offset</pre>
##
       else {
##
```

#### Prediction Using the SLR Model in R

- Can use the line for prediction with predict()!
  - Should supply fitted object and newdata
    - An optional data frame in which to look for variables with which to predict. If omitted, the fitted values are used.

```
predict(slr_fit, newdata = data.frame(log_km_driven = c(log(1000), log(10000), log(10000))))
## 1 2 3
## 11.93404 11.03353 10.13302
exp(predict(slr_fit, newdata = data.frame(log_km_driven = c(log(1000), log(10000), log(10000)))))
## 1 2 3
## 152365.60 61915.64 25160.19
```

#### Quantifying How Well the Model Predicts

We use a **loss** function to fit the model. We use a **metric** to evaluate the model!

- Often use the same loss function for fitting and as the metric
- For a given **numeric** response value,  $y_i$  and prediction,  $\hat{y}_i$

$$(y_i-{\hat y}_i)^2, |y_i-{\hat y}_i|$$

• Incorporate all points via

$$rac{1}{n}\sum_{i=1}^n (y_i - {\hat y}_i)^2, rac{1}{n}\sum_{i=1}^n |y_i - {\hat y}_i|$$

#### **Metric Function**

• For a numeric response, we commonly use squared error loss as our metric to evaluate a prediction

$$L(y_i, {\hat y}_i) = (y_i - {\hat y}_i)^2$$

• Use Root Mean Square Error as a **metric** across all observations

$$RMSE = \sqrt{rac{1}{n}\sum_{i=1}^{n}L(y_i, \hat{y}_i)} = \sqrt{rac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}$$

#### **Commonly Used Metrics**

For prediction (numeric response)

- Mean Squared Error (MSE) or Root Mean Squared Error (RMSE)
- Mean Absolute Error (MAE or MAD deviation)

$$L(y_i, {\hat y}_i) = |y_i - {\hat y}_i|$$

• Huber Loss

• Doesn't penalize large mistakes as much as MSE

#### **Commonly Used Metrics**

For prediction (numeric response)

- Mean Squared Error (MSE) or Root Mean Squared Error (RMSE)
- Mean Absolute Error (MAE or MAD deviation)

$$L(y_i, {\hat y}_i) = |y_i - {\hat y}_i|$$

• Huber Loss

• Doesn't penalize large mistakes as much as MSE

For classification (categorical response)

- Accuracy
- log-loss
- AUC
- F1 Score

#### Evaluating our SLR Model

• We could find our metric for our SLR model using the training data

```
• Called training error
```

```
head(predict(slr_fit))
## 1 2 3 4 5 6
## 12.34461 11.25681 10.96222 10.70779 10.74337 10.33280
mean((bike_data$log_selling_price-predict(slr_fit))^2)
## [1] 0.3536708
sqrt(mean((bike_data$log_selling_price-predict(slr_fit))^2))
```

## [1] 0.5947023

• Doesn't tell us how well we do on data we haven't seen!

## Training vs Test Sets

Ideally we want our model to predict well for observations **it has yet to see**!

- For *multiple* linear regression models, our training MSE will always decrease as we add more variables to the model...
- We'll need an independent **test** set to predict on (more on this shortly!)

## **Big Picture Modeling**

Supervised Learning methods try to relate predictors to a response variable through a model

- Lots of common models
  - Regression models
  - $\circ~$  Tree based methods
  - $\circ~$  Naive Bayes
  - k Nearest Neighbors
  - o ...
- For a set of predictor values, each will produce some prediction we can call  $\hat{y}$
- Evaluate model via a metric
- Will use an independent test set or cross-validation to more accurately judge our model