

# Course Summary

## Fundamental Techniques in Data Science



**Utrecht  
University**

Kyle M. Lang

Department of Methodology & Statistics  
Utrecht University

1. *Journal of the American Medical Association*, 1990; 263: 1025-1028.

— **1998** —

### Rehabilitation Odds

— *Journal of the American Medical Association*

# Linear Regression

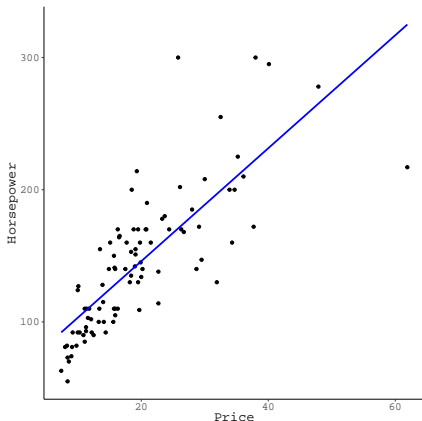


# Simple Linear Regression

In linear regression, we want to find the best fit line:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

- For any  $X_n$ , the corresponding  $\hat{Y}_n$  represents the model-implied, conditional mean of  $Y$ .



# Simple Linear Regression

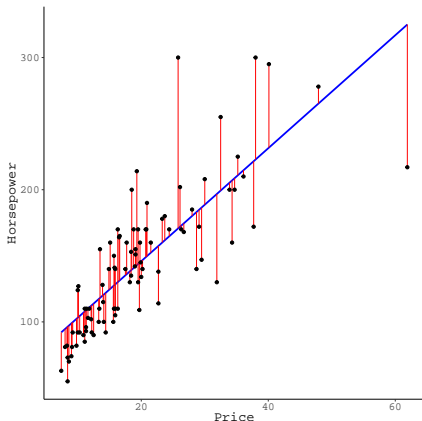
In linear regression, we want to find the best fit line:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

- For any  $X_n$ , the corresponding  $\hat{Y}_n$  represents the model-implied, conditional mean of  $Y$ .

After accounting for the estimation error, we get the full regression equation:

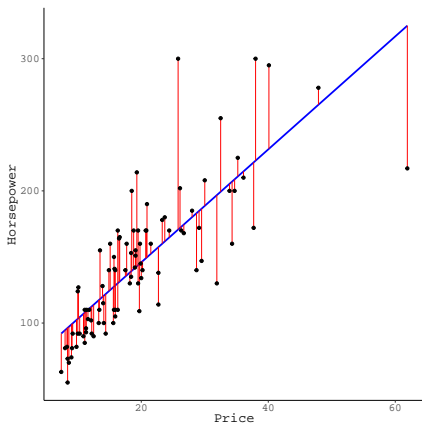
$$Y = \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\varepsilon}$$



# Residuals as the Basis of Estimation

We use the residuals,  $\hat{\varepsilon}_n$ , to estimate the model.

$$\begin{aligned}RSS &= \sum_{n=1}^N \hat{\varepsilon}_n^2 = \sum_{n=1}^N (Y_n - \hat{Y}_n)^2 \\ &= \sum_{n=1}^N (Y_n - \hat{\beta}_0 - \hat{\beta}_1 X_n)^2\end{aligned}$$



# Example

```
## Read in the 'diabetes' dataset:
```

```
diabetes <- readRDS("../data/diabetes.rds")
```

```
## Estimate and summarize a regression model:
```

```
lm(bp ~ age + ldl + hdl + sex, data = diabetes) %>% partSummary(-1)
```

Residuals:

Min	1Q	Median	3Q	Max
-34.195	-8.734	-1.011	7.945	42.186

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	78.18713	4.29453	18.206	< 2e-16
age	0.30043	0.04789	6.273	8.52e-10
ldl	0.03887	0.02079	1.870	0.06220
hdl	-0.09063	0.05124	-1.769	0.07763
sexmale	4.07606	1.32803	3.069	0.00228

Residual standard error: 12.72 on 437 degrees of freedom

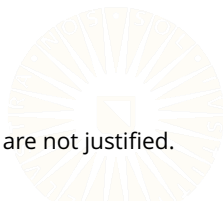
Multiple R-squared: 0.162, Adjusted R-squared: 0.1543

F-statistic: 21.12 on 4 and 437 DF, p-value: 6.163e-16

# Assumptions

---

1. The model is linear in the parameters.
  - *Otherwise:* We are not working with linear regression.
2. The predictor matrix is *full rank*.
  - *Otherwise:* The model is not estimable.
3. The predictors are strictly exogenous.
  - *Otherwise:* The estimated regression coefficients will be biased.
4. The errors have constant, finite variance.
  - *Otherwise:* Standard errors will be biased.
5. The errors are uncorrelated.
  - *Otherwise:* Standard errors will be biased.
6. The errors are normally distributed.
  - *Otherwise:* Small-sample inferences and some estimates are not justified.



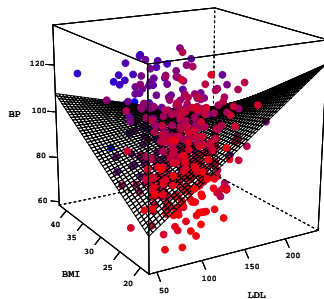
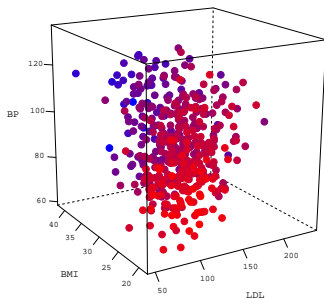


# Moderation



# Moderated Regression

The effect of  $X$  on  $Y$  varies **as a function** of  $Z$ .



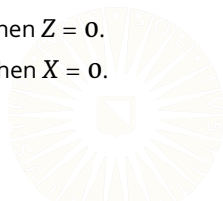
# Interpretation

---

Given the following equation:

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 Z + \hat{\beta}_3 XZ + \hat{\varepsilon}$$

- $\hat{\beta}_3$  quantifies the effect of  $Z$  on the focal effect (the  $X \rightarrow Y$  effect).
  - For a unit change in  $Z$ ,  $\hat{\beta}_3$  is the expected change in the effect of  $X$  on  $Y$ .
- $\hat{\beta}_1$  and  $\hat{\beta}_2$  are *conditional effects*.
  - Interpreted where the other predictor is zero.
  - For a unit change in  $X$ ,  $\hat{\beta}_1$  is the expected change in  $Y$ , when  $Z = 0$ .
  - For a unit change in  $Z$ ,  $\hat{\beta}_2$  is the expected change in  $Y$ , when  $X = 0$ .



# Continuous Moderators

---

```
## Moderated Model:
```

```
out2 <- lm(bp ~ bmi * ldl, data = diabetes)
partSummary(out2, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	14.480616	14.291677	1.013	0.311514
bmi	2.867825	0.541312	5.298	1.86e-07
ldl	0.448771	0.127160	3.529	0.000461
bmi:ldl	-0.015352	0.004716	-3.255	0.001221

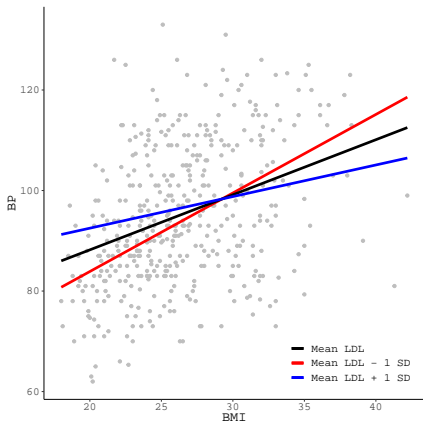
Residual standard error: 12.54 on 438 degrees of freedom

Multiple R-squared: 0.1834, Adjusted R-squared: 0.1778

F-statistic: 32.78 on 3 and 438 DF, p-value: < 2.2e-16

# Visualizing the Interaction

We can get a better idea of the patterns of moderation by plotting the focal effect at conditional values of the moderator.



# Categorical Moderators

```
## Load data:  
socSup <- readRDS("../data/social_support.rds")
```

```
## Estimate the moderated regression model:  
out4 <- lm(bdi ~ tanSat * sex, data = socSup)  
partSummary(out4, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	20.8478	6.2114	3.356	0.00115
tanSat	-0.5772	0.3614	-1.597	0.11372
sexmale	14.3667	12.2054	1.177	0.24223
tanSat:sexmale	-0.9482	0.7177	-1.321	0.18978

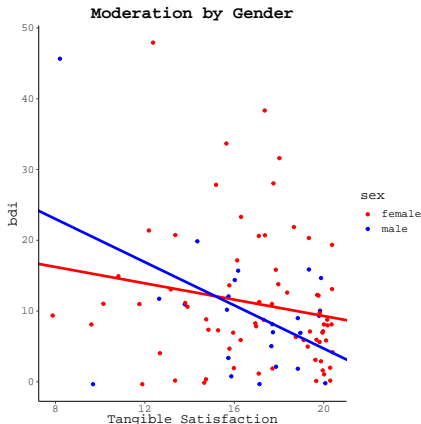
Residual standard error: 9.267 on 91 degrees of freedom

Multiple R-squared: 0.08955, Adjusted R-squared: 0.05954

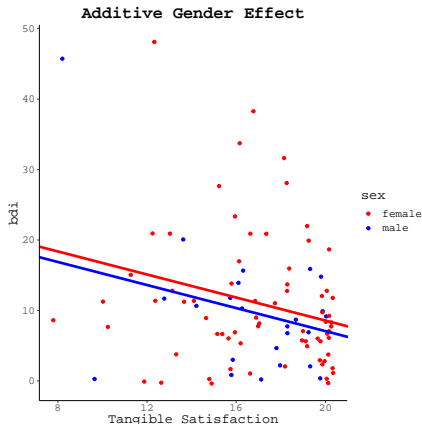
F-statistic: 2.984 on 3 and 91 DF, p-value: 0.03537

# Visualizing Categorical Moderation

$$\hat{Y}_{BDI} = 20.85 - 0.58X_{tsat} + 14.37Z_{male} - 0.95X_{tsat}Z_{male}$$



$$\hat{Y}_{BDI} = 24.91 - 0.82X_{tsat} - 1.50Z_{male}$$



# Prediction





# Prediction Example

---

Let's fit the following model using the *diabetes* data:

$$Y_{LDL} = \beta_0 + \beta_1 X_{BP} + \beta_2 X_{gluc} + \beta_3 X_{BMI} + \varepsilon$$

Training this model on the first  $N = 400$  patients' data produces the following fitted model:

$$\hat{Y}_{LDL} = 22.135 + 0.089X_{BP} + 0.498X_{gluc} + 1.48X_{BMI}$$



# Prediction Example

---

Let's fit the following model using the *diabetes* data:

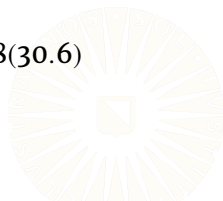
$$Y_{LDL} = \beta_0 + \beta_1 X_{BP} + \beta_2 X_{gluc} + \beta_3 X_{BMI} + \varepsilon$$

Training this model on the first  $N = 400$  patients' data produces the following fitted model:

$$\hat{Y}_{LDL} = 22.135 + 0.089X_{BP} + 0.498X_{gluc} + 1.48X_{BMI}$$

Suppose a new patient presents with  $BP = 121$ ,  $gluc = 89$ , and  $BMI = 30.6$ . We can predict their *LDL* score by:

$$\begin{aligned}\hat{Y}_{LDL} &= 22.135 + 0.089(121) + 0.498(89) + 1.48(30.6) \\ &= 122.463\end{aligned}$$



# Interval Estimates Example

---

Two flavors of interval to quantify prediction uncertainty:

1. Confidence intervals
2. Prediction intervals

In our example, we get the following 95% interval estimates:

$$95\% CI_{\hat{Y}} = [115.6; 129.33]$$

$$95\% PI = [66.56; 178.37]$$

- We can be 95% confident that the average LDL of patients with *Glucose* = 89, *BP* = 121, and *BMI* = 30.6 will be somewhere between 115.6 and 129.33.
- We can be 95% confident that the LDL of a specific patient with *Glucose* = 89, *BP* = 121, and *BMI* = 30.6 will be somewhere between 66.56 and 178.37.

# Model Fit



# Model Fit

---

We quantify the proportion of the outcome's variance that is explained by our model using the  $R^2$  statistic:

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$

where

$$TSS = \sum_{n=1}^N (Y_n - \bar{Y})^2 = \text{Var}(Y) \times (N - 1)$$

For the model we estimated in the above prediction example, we get:

$$R^2 = 1 - \frac{315383}{361704} \approx 0.13$$



# Model Fit for Prediction

We use the *mean squared error* (MSE) to assess predictive performance.

$$\begin{aligned} \text{MSE} &= \frac{1}{N} \sum_{n=1}^N (Y_n - \hat{Y}_n)^2 \\ &= \frac{1}{N} \sum_{n=1}^N \left( Y_n - \hat{\beta}_0 - \sum_{p=1}^P \hat{\beta}_p X_{np} \right)^2 \\ &= \frac{\text{RSS}}{N} \end{aligned}$$

For our example problem, we get:

$$\text{MSE}_{\text{train}} = \frac{315383}{400} \approx 788.46$$

$$\text{MSE}_{\text{train}} = \frac{48092.04}{42} \approx 1145.05$$



# Information Criteria

---

We can use *information criteria* to quickly compare *non-nested* models while accounting for model complexity.

- Akaike's Information Criterion (AIC)

$$AIC = 2K - 2\hat{\ell}(\theta|X)$$

- Bayesian Information Criterion (BIC)

$$BIC = K \ln(N) - 2\hat{\ell}(\theta|X)$$

For our example, we get the following estimates of AIC and BIC:

$$\begin{aligned} AIC &= 2(3) - 2(-1901.59) \\ &= 3813.18 \end{aligned}$$

$$\begin{aligned} BIC &= 3 \ln(400) - 2(-1901.59) \\ &= 3833.14 \end{aligned}$$



# Cross Validation

---

To train a model that best predicts new data, we can use *cross-validation* to evaluate the expected predictive performance on new data.

1. Split the sample into two, disjoint sub-samples
  - *Training* data
  - *Testing* data
2. Estimate a candidate model,  $f(\mathbf{X})$ , on the training data.
3. Check the predictive performance of  $\hat{f}(\mathbf{X})$  on the testing data.





# Cross Validation

---

To train a model that best predicts new data, we can use *cross-validation* to evaluate the expected predictive performance on new data.

1. Split the sample into two, disjoint sub-samples
  - *Training* data
  - *Testing* data
2. Estimate a candidate model,  $f(\mathbf{X})$ , on the training data.
3. Check the predictive performance of  $\hat{f}(\mathbf{X})$  on the testing data.

We can use this idea to select the best model from a pool of candidate models,  $\mathcal{F} = \{f_1(\mathbf{X}), f_2(\mathbf{X}), \dots, f_J(\mathbf{X})\}$

1. Repeat Steps 2 and 3 for all candidate models in  $\mathcal{F}$ .
2. Pick the  $\hat{f}_j(\mathbf{X})$  that best predicts the testing data.



# K-Fold Cross-Validation

---

1. Partition the data into  $K$  disjoint subsets  $C_k = C_1, C_2, \dots, C_K$ .



# K-Fold Cross-Validation

---

1. Partition the data into  $K$  disjoint subsets  $C_k = C_1, C_2, \dots, C_K$ .
2. Conduct  $K$  training replications.
  - For each training replication, collapse  $K - 1$  partitions into a set of training data, and use this training data to estimate the model.
  - Compute the test MSE for the  $k$ th partition,  $MSE_k$ , by using subset  $C_k$  as the test data for the  $k$ th fitted model.



# K-Fold Cross-Validation

---

1. Partition the data into  $K$  disjoint subsets  $C_k = C_1, C_2, \dots, C_K$ .
2. Conduct  $K$  training replications.
  - For each training replication, collapse  $K - 1$  partitions into a set of training data, and use this training data to estimate the model.
  - Compute the test MSE for the  $k$ th partition,  $MSE_k$ , by using subset  $C_k$  as the test data for the  $k$ th fitted model.
3. Compute the overall  $K$ -fold cross-validation error as:

$$CVE = \sum_{k=1}^K \frac{N_k}{N} MSE_k,$$



# Logistic Regression



# Probabilities & Odds

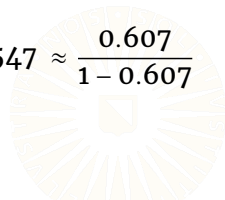
Sex	Complete	
	No	Yes
Female	95	147
Male	753	1540

$$P(C|M) = \frac{1540}{1540 + 753} = 0.672$$

$$O(C|M) = \frac{1540}{753} = 2.045 \approx \frac{0.672}{1 - 0.672}$$

$$P(C|F) = \frac{147}{147 + 95} = 0.607$$

$$O(C|F) = \frac{147}{95} = 1.547 \approx \frac{0.607}{1 - 0.607}$$



# The Generalized Linear Model

---

Every GLM is built from three components:

1. The systematic component,  $\eta$ .
  - A linear function of the predictors,  $\{X_p\}$ .
  - Describes the association between  $\mathbf{X}$  and  $Y$ .
2. The link function,  $g(\mu_Y)$ .
  - Transforms  $\mu_Y$  so that it can take any value on the real line.
3. The random component,  $P(Y|g^{-1}(\eta))$ 
  - The distribution of the observed  $Y$ .
  - Quantifies the error variance around  $\eta$ .



# The Logistic Regression Model

---

The logistic regression model can be represented as:

$$Y \sim \text{Bin}(\pi, 1)$$

$$\text{logit}(\pi) = \beta_0 + \sum_{p=1}^P \beta_p X_p$$

The fitted model can be represented as:

$$\text{logit}(\hat{\pi}) = \hat{\beta}_0 + \sum_{p=1}^P \hat{\beta}_p X_p$$

To convert fitted values,  $\hat{\eta} = \hat{\beta}_0 + \sum_{p=1}^P \hat{\beta}_p X_p$ , from a logit scale to a probability scale, we apply the *logistic* function:

$$\text{logistic}(\hat{\eta}) = \frac{e^{\hat{\eta}}}{1 + e^{\hat{\eta}}}$$





# Logistic Regression Example

```
## Coarsen the blood glucose variable:
diabetes %<>% mutate(highGlu = as.numeric(glu > 90))

## Estimate the model:
out1 <- glm(highGlu ~ age + bmi + bp, data = diabetes, family = binomial())
partSummary(out1, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-6.479104	0.912899	-7.097	1.27e-12
age	0.034597	0.008635	4.007	6.16e-05
bmi	0.106852	0.026660	4.008	6.12e-05
bp	0.022691	0.008560	2.651	0.00803

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 610.42 on 441 degrees of freedom  
Residual deviance: 538.18 on 438 degrees of freedom  
AIC: 546.18

Number of Fisher Scoring iterations: 4

# Assumptions

---

We can state the assumptions of logistic regression as follows:

1. The outcome follows a binomial distribution.
2. The predictor matrix is full-rank.
3. The predictors are linearly related to *logit*( $\pi$ ).
4. The observations are independent after accounting for the predictors.

Unlike linear regression, we don't need to assume

- Constant, finite error variance
- Normally distributed errors

For computational reasons, we also need the following:

- Large sample
- Relatively well-balance outcome
- No highly influential cases



# Classification



# Classification Example

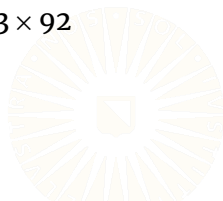
---

Say we want to classify a new patient into either the “high glucose” group or the “not high glucose” group using the model fit above.

- Assume this patient has the following characteristics:
  - They are 57 years old
  - Their BMI is 28
  - Their average blood pressure is 92

First we plug their predictor data into the fitted model to get their model-implied  $\eta$ :

$$\begin{aligned}\hat{\eta} &= -6.479 + 0.035 \times 57 + 0.107 \times 28 + 0.023 \times 92 \\ &= 0.572\end{aligned}$$



# Classification Example

---

Next we convert the predicted  $\eta$  value into a model-implied success probability by applying the logistic function:

$$\hat{\pi} = \text{logistic}(0.572) = \frac{e^{0.572}}{1 + e^{0.572}} = 0.639$$

Finally, to make the classification, assume a threshold of  $\hat{\pi} = 0.5$  as the decision boundary.

- Because  $0.639 > 0.5$  we would classify this patient into the “high glucose” group.



# Confusion Matrix

True	Predicted	
	Low	High
Low	123	82
High	62	175

Confusion Matrix of Blood Glucose Level

$$\text{Sensitivity} = \frac{175}{175 + 62} = 0.738$$

$$\text{Specificity} = \frac{123}{123 + 82} = 0.6$$

$$\text{Accuracy} = \frac{175 + 123}{175 + 123 + 62 + 82} = 0.674$$

