

# Week 4

## Conditional Expectation and the Best Linear Predictor

---

Paul Laskowski and Alex Hughes

January 12, 2023

UC Berkeley, School of Information

# Motivating Examples

---

## MOTIVATING EXAMPLES

Given a random variable  $\mathbf{X}$  and another random variable  $\mathbf{Y}$ , you often want to make predictions about  $\mathbf{Y}$  that are *as good as possible*.

- Given some budget, what will an organization's head count be next year?
- Given a citizen's age and income, how likely are they to vote?
- Given a set of sounds picked up by a microphone, what are the chances that someone in the room has fallen down (Sound Flux, 2019)?
- Given a set of voxels, does an image contain a malignant tumor?

# PLAN FOR THE WEEK

Two sections:

1. Conditional expectations
2. Best predictors and best linear predictors

## PLAN FOR THE WEEK (CONT.)

At the end of this week, you will be able to:

- Derive statements describing relationships among random variables
- Understand how to divide variance into a part that is explained and an error component that you have not explained
- Understand how a best linear predictor uses information from a set of random variables to help us predict the value of an outcome

# Conditional Expectation

---

# Introduction to Conditional Expectation

---

# INTRODUCTION TO CONDITIONAL EXPECTATION, PART I

To this point, we have characterized random variables with three concepts:

1. **The Expected Value**,  $E[X]$
2. **The Variance**,  $V[X]$
3. **The Covariance**,  $\text{cov}[X, Y]$

But, when there is dependence between two random variables  $X$  and  $Y$ , then if we know something about  $X$ , we might be able to make a new contextualized statement about the *conditional distribution*.



# INTRODUCTION TO CONDITIONAL EXPECTATION, PART II

We began the week with the questions:

- Given some budget, what will an organization's head count be next year?
- Given some attributes about a citizen, how likely are they to vote?
- Given a set of sounds picked up by a microphone, what are the chances that someone in the room has fallen down (Sound Flux, 2019)?
- Given a set of voxels, does this image contain a malignant tumor?

## INTRODUCTION TO CONDITIONAL EXPECTATION, PART III

But, those were all conditional statements!

Rather than saying, “On average, 78% of the electorate votes in any given election,” we can instead say:

- “Among people 29 or younger, 41% voted in the last presidential election”
- “Among people who *have* voted before, 90% will vote in any given election”

This is our goal in data science! Use information to make more informed, *better* statements.

# Conditional Expectation Operator

---

# INTRODUCTION TO CONDITIONAL EXPECTATION, PART I

- Recall expected values:

$$E[Y] = \int_{-\infty}^{\infty} y \cdot f_Y dy$$

- The expected value of  $Y$ ,  $E[Y]$ , is the integral of the product of {the value  $y$  and the probability density function for that value}
- Expectation is an operator that maps from  $\mathbb{R}^n \rightarrow \mathbb{R}$

# INTRODUCTION TO CONDITIONAL EXPECTATION, PART II

- The conditional expectation operator of  $Y$  given  $X$ ,  $E[Y|X]$  holds a similar form:

$$E[Y|X] = \int_{-\infty}^{\infty} y \cdot f_{Y|X}(y|x) dy, \forall X \in \text{Supp}[X]$$

- Recall also that expectation is an operator.
- So, too, is conditional expectation.

# INTRODUCTION TO CONDITIONAL EXPECTATION, PART III

# Conditional Expectation Demo

---

# SOFTWARE DEMO: CONDITIONAL EXPECTATION

**Note: This is a lecture + software demo. We're just placing this here for organization.**



# **Learnosity: Conditional Expectation**

---

## CONDITIONAL EXPECTATION EXAMPLE

The following table shows the joint probability distribution for discrete random variables  $X$  and  $Y$ .

		X		
		1	2	3
Y	1	0.1	0.0	0.0
	2	0.2	0.2	0.1
	3	0.1	0.2	0.1

1. Compute:  $E(Y|X = 1)$
2. Compute:  $E(Y|X = 2)$

# Reading: Conditional Expectation Operator

---

## READING: CONDITIONAL EXPECTATION OPERATOR

Read page 67 through the first half of 69 (including Theorem 2.2.14)

# Properties of Conditional Operators

---

# UNDERSTANDING CONDITIONAL OPERATORS

Conditional Expectation  $\longrightarrow E[\text{Conditional Distribution}]$

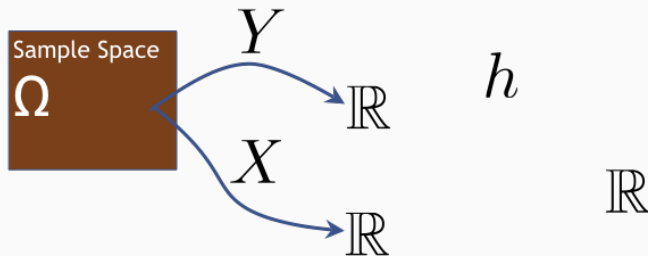
Conditional Variance  $\longrightarrow V[\text{Conditional Distribution}]$

## FORMULAS FOR CONDITIONAL VARIANCE

$$V[Y|X = x] = E[(Y - E[Y|X = x])^2|X = x]$$

$$V[Y|X = x] = E[Y^2|X = x] - E[Y|X = x]^2$$

# UNDERSTANDING CONDITIONAL LOTUS



- Condition on  $X = x$ .



## Theorem: Conditional LOTUS

Given discrete random variables  $X$  and  $Y$  with joint pmf  $f$ , and  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,

$$\mathbb{E}[h(X, Y)|X = x] = \sum_y h(x, y)f_{Y|X}(y|x),$$

For all  $x \in \text{Supp}[X]$ .

Given continuous random variables  $X$  and  $Y$  with joint pdf  $f$ ,

$$\mathbb{E}[h(X, Y)|X = x] = \int_{-\infty}^{\infty} h(x, y)f_{Y|X}(y|x)dy,$$

for all  $x \in \text{Supp}[X]$ .

## Theorem: Linearity of Conditional Expectation

Given random variables  $X$  and  $Y$  and  $a, b \in \mathbb{R}$ , then for all  $x \in \text{Supp}[X]$ ,

$$E[aY + b|X = x] = aE[Y|X = x] + b$$

Moreover, given  $g, h : \mathbb{R} \rightarrow \mathbb{R}$ , then for all  $x \in \text{Supp}[X]$ ,

$$E[g(X)Y + h(X)|X = x] = g(x)E[Y|X = x] + h(x)$$

# CONDITIONAL EXPECTATION EXAMPLE

Simplify:  $E[XY|X = x]$

# Demonstration of the Conditional Expectation Function

---

# DEMONSTRATION OF THE CEF

**Note: This is a Lecture and Software Demo, we're just placing it here for organization.**

# Reading: Conditional Expectation Function

---

## READING: CONDITIONAL EXPECTATION FUNCTION

- Read the second half of page 69, through page 71, stopping before the Law of Iterated Expectations.
- These pages are tough, so we're going to whiteboard theorems 2.2.11 and 2.2.12 on the other side of your reading.
- You can probably skip theorem 2.2.14 without great problem.
- Focus specific attention on definition 2.2.15.

# **Learnosity: Write a Conditional Expectation Function**

---



# **Introduction to the Law of Iterated Expectations**

---

# LAW OF ITERATED EXPECTATIONS

Earlier, we worked with **theorem 1.1.13**, *The law of total probability*, which stated:

## **Theorem 1.1.13: the law of total probability**

If  $\{A_1, A_2, A_3, \dots\}$  partition  $\Sigma$ ,  $B \in S$ , and  $P(A_i) > 0, \forall i$  then:

$$P(B) = \sum_i P(B|A_i)P(A_i)$$

- Useful because we can break the overall probability into smaller “chunks” whose probabilities might be easier to know.
- The Law of Iterated Expectations is a generalization and reapplication of this principle.

# **Reading: Law of Iterated Expectations and Total Variance**

---

## READING: LAW OF ITERATED EXPECTATIONS

Read pages 72, 73, and the top of 74, stopping before theorem 2.2.19.

# **Lightboard: Law of Iterated Expectations**

---

# LAW OF ITERATED EXPECTATIONS

# Law of Total Variance

---

# SPLITTING VARIANCE

## Theorem 2.2.18: the law of total variance

For random variables  $X$  and  $Y$ :

$$V[Y] = V[E[Y|X]] + E[V[Y|X]]$$

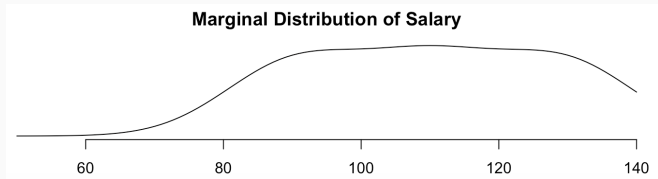
**Key question:** How much of the variance in  $Y$  can be explained by  $X$ ?



## SPLITTING VARIANCE: EXAMPLE

$S$  represents salary.

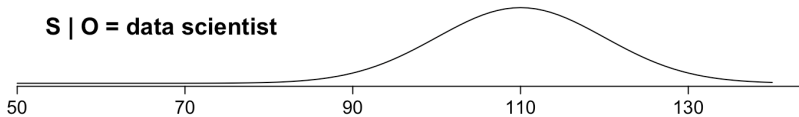
$O$  represents occupation (1 = data scientist, 2 = data engineer, 3 = data analyst).



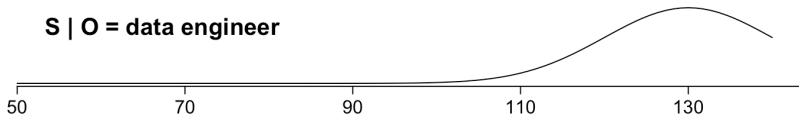
$V[S] = 366$ . How much of this is explained by occupation?

# SALARY CONDITIONED ON OCCUPATION

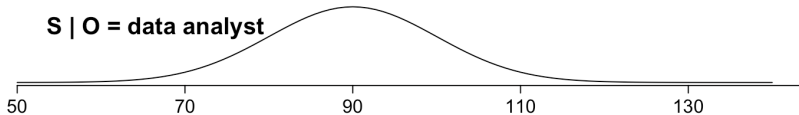
**$S \mid O = \text{data scientist}$**



**$S \mid O = \text{data engineer}$**



**$S \mid O = \text{data analyst}$**



# COMPONENTS OF VARIANCE

$V[S]$  is the sum of two components.

1. **Explained variance:**  $V[E[S|O]] = 266$

- Measures how much information  $O$  gives about  $S$
- Can also be called systematic variance

# COMPONENTS OF VARIANCE

$V[S]$  is the sum of two components.

1. **Explained variance:**  $V[E[S|O]] = 266$

- Measures how much information  $O$  gives about  $S$
- Can also be called systematic variance

2. **Unexplained variance:**  $E[V[S|O]] = 100$

- Measures how much extra variation is there in  $S$  that we can't explain with  $O$
- Can also be called error variance

## LAW OF TOTAL VARIANCE IMPLICATIONS

Suppose that you can divide data into groups along (possibly) two dimensions.

1. A dimension that does not explain any variation in outcomes
2. A dimension that explains variation in outcomes

As you're going to see in the next coding exercise, through this identity, if you can produce groups that have different means, you will necessarily produce a smaller  $E[V[Y|X]]$ .

# **Learnosity: Variance Breakdown**

---

# LEARNOSITY: VARIANCE BREAKDOWN

**Note: This is a learnosity activity. We're just including it here for organization.**

# Best Predictors

---



# Deviations From the CEF

---

## MORE COMMENTS FOR ALEX

# JOINT DENSITY FUNCTION

For this section, we will work first with the same joint distribution function throughout.

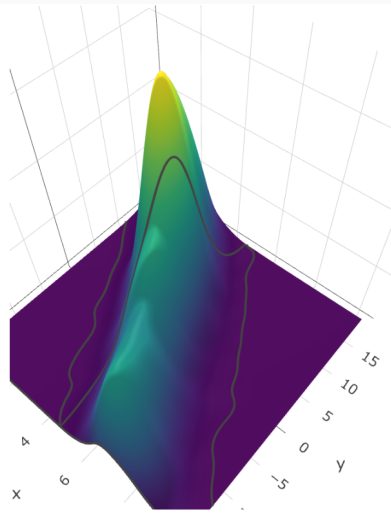
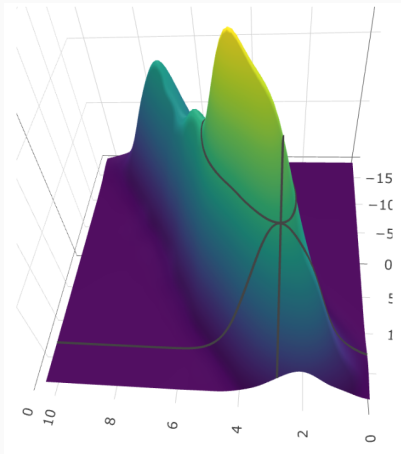
$$X \sim U(\text{min} = 0, \text{max} = 10)$$

$$W \sim N(\text{mean} = 0, \text{sd} = 1)$$

$$Y = 10 - 2 \cdot X + W$$

Where  $X$  and  $W$  are independent.

# JOINT DENSITY FUNCTION (CONT.)



# DEVIATIONS FROM THE CEF

## Deviations from the CEF

Suppose that  $\epsilon = Y - E[Y|X]$  is the distance between my estimate within a grid and the actual value.

- Inside each of the grids, how far, on average, will I be from the true elevation?  $E[\epsilon|X] = ?$
- Across the whole ridgeline, how far, on average, will I be from the true elevation?  $E[\epsilon] = E[E[\epsilon|X]] = ?$
- Within each of the grids, what features will shape how close you are, on average?  $V[\epsilon|X] = ?$
- Across the whole ridgeline, what features will shape how close you are, on average?  $V[\epsilon] = ?$

## **Reading: Deviations from the CEF**

---

## DEVIATIONS FROM THE CEF

- Begin by reading from where we left off previously on page 74 to the middle of page 76, stopping when the discussion turns to linear restrictions.
- Try to work through, on your own, the proof of why the CEF is the best (minimum MSE) estimator of  $Y$ .
- In the proof of theorem 2.2.20, the authors add  $E[\epsilon|X]$  seemingly from nowhere. This is a legal move because in 2.2.19, you have derived that  $E[\epsilon|X] = 0$ .

# Best Predictors

---



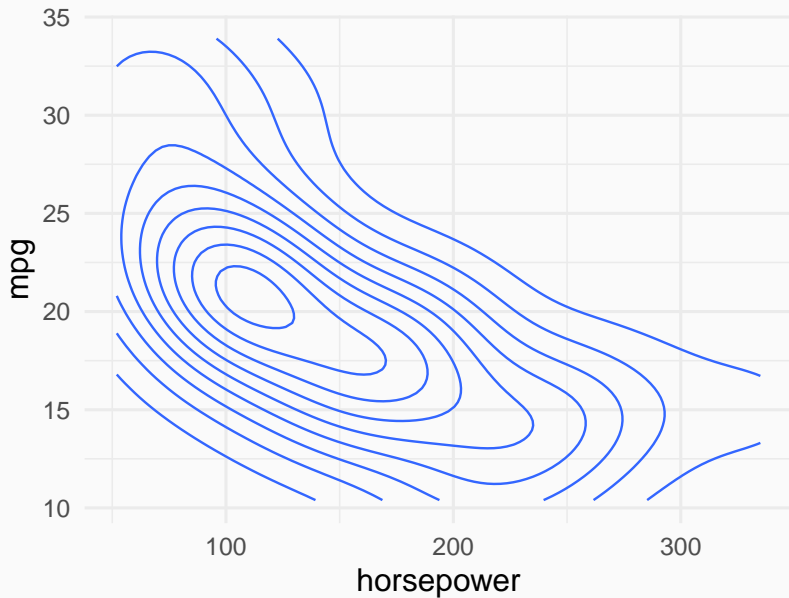
# PREDICTORS ARE FUNCTIONS

## Definition: Predictor

Given random variables  $X$  and  $Y$ , a predictor for  $Y$  is a function,  $g : \mathbb{R} \rightarrow \mathbb{R}$ , such that  $g(X)$  is regarded as a guess for  $Y$ .

- Given a value  $X = x$ , a predictor suggests a single value for  $Y$ .
- Define error as  $\epsilon = Y - g(X)$ .

# PREDICTORS AND ERRORS



# CHOOSING A GOOD PREDICTOR

**Idea:** Minimize mean squared error,  $E[\epsilon^2]$ .

# IMPORTANCE OF THE CEF

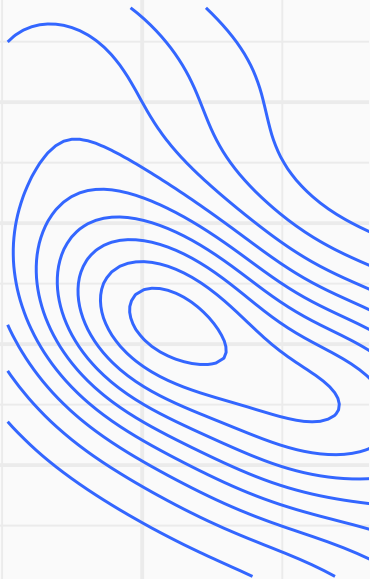
## **Theorem: The CEF Minimizes MSE**

Given random variables  $X$  and  $Y$ , the CEF  $E[Y|X]$  has the smallest MSE out of all predictors of  $Y$ .

# Lightboard: The CEF Minimizes MSE

---

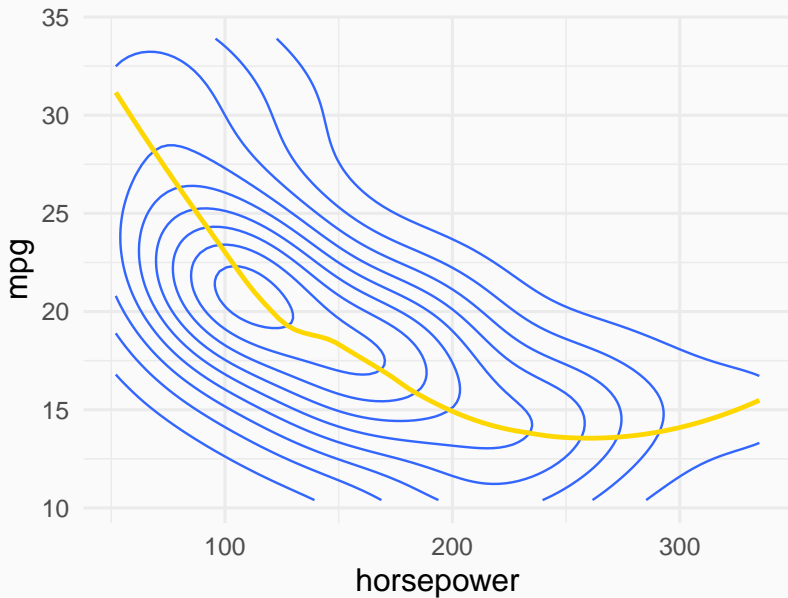
# THE CEF MINIMIZES MSE



# Best Linear Predictors

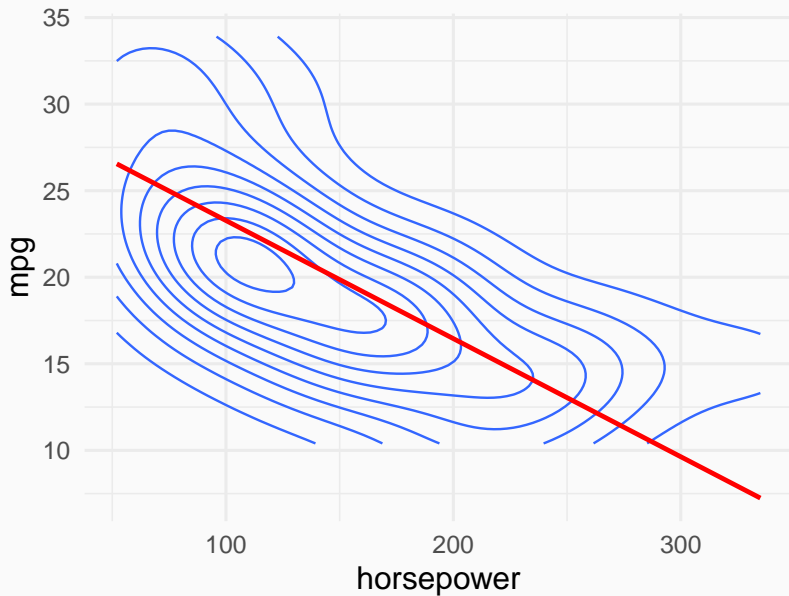
---

# COMPLEXITY OF THE CEF





# LINEARIZING THE PREDICTOR



# CHOOSING A LINEAR PREDICTOR

## Definition: The Best Linear Predictor (BLP)

Given random variables  $X$  and  $Y$ , the *best linear predictor* (BLP) for  $Y$  is the function  $g : \mathbb{R} \rightarrow \mathbb{R}$ ,  $g(X) = \alpha + \beta X$ , with coefficients given by,

$$\operatorname{argmin}_{\alpha, \beta} \mathbb{E} \left[ \left( Y - (\alpha + \beta X) \right)^2 \right]$$

# SOLVING FOR THE BLP

## Theorem: The BLP Solution

Given random variables  $X$  and  $Y$ , the best linear predictor for  $Y$  is given by  $g(X) = \alpha + \beta X$ , where,

$$\beta = \frac{\text{cov}[X, Y]}{V[X]}, \quad \alpha = E[Y] - \frac{\text{cov}[X, Y]}{V[X]}E[X]$$

# **Reading: The Best Linear Predictor**

---

## READING ASSIGNMENT: LINEAR APPROXIMATION

Read to the middle of page 79, stopping before you get to example 2.2.23.

# Moment Conditions

---

## MOMENT CONDITIONS

The BLP is given by  $\operatorname{argmin}_{\alpha, \beta} E\left[\left(Y - (\alpha + \beta X)\right)^2\right]$

# MOMENT CONDITIONS - SOLUTION

The BLP is given by  $\operatorname{argmin}_{\alpha, \beta} E \left[ \left( Y - (\alpha + \beta X) \right)^2 \right]$

$$0 = \frac{\partial E[\epsilon^2]}{\partial \alpha} = E \left[ \frac{\partial \epsilon^2}{\partial \alpha} \right] = E \left[ 2\epsilon \frac{\partial \epsilon}{\partial \alpha} \right] = -2E[\epsilon]$$

$$0 = \frac{\partial E[\epsilon^2]}{\partial \beta} = E \left[ \frac{\partial \epsilon^2}{\partial \beta} \right] = E \left[ 2\epsilon \frac{\partial \epsilon}{\partial \beta} \right] = -2E[\epsilon X]$$

Version 1:

1.  $E[\epsilon] = 0$
2.  $E[\epsilon X] = 0$

Version 2:

1.  $E[\epsilon] = 0$
2.  $\operatorname{cov}[\epsilon, X] = E[\epsilon X] - E[\epsilon]E[X] = 0$

What about  $E[\epsilon|X]$ ?



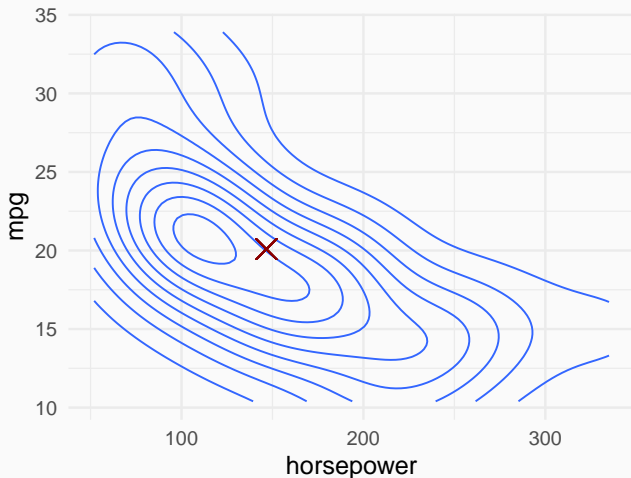
# **Lightboard: Understanding Moment Conditions**

---

# UNDERSTANDING MOMENT CONDITIONS

1.  $0 = E[\epsilon]$

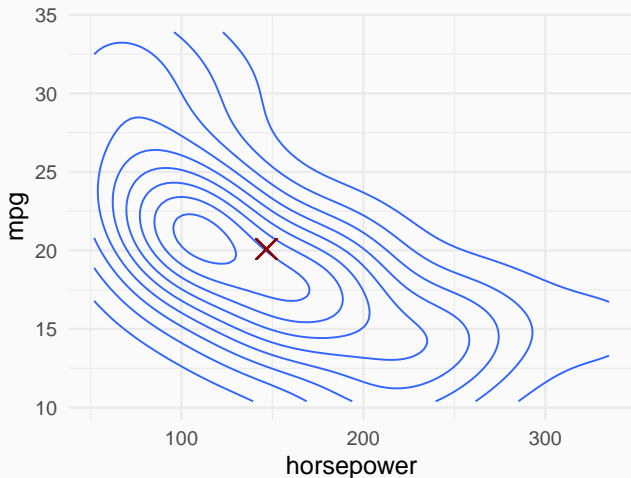
2.  $0 = \text{cov}[\epsilon, X]$



# UNDERSTANDING MOMENT CONDITIONS - SOLUTION

1.  $0 = E[\epsilon] = E[Y - (\alpha + \beta X)] \implies E[Y] = \alpha + \beta E[X]$

2.  $0 = \text{cov}[\epsilon, X]$



# Lightboard: The BLP Solution

---

# THE BLP SOLUTION

The BLP is defined by the moment conditions:

1.  $0 = E[\epsilon]$

2.  $0 = E[\epsilon X]$

# THE BLP SOLUTION - SOLUTION

The BLP is defined by the moment conditions:

$$1. \ 0 = E[\epsilon] = E[Y - (\alpha + \beta X)] = E[Y] - \alpha - \beta E[X].$$

$$\alpha = E[Y] - \beta E[X].$$

$$\begin{aligned} 2. \ 0 &= E[\epsilon X] = E[(Y - (\alpha + \beta X))X] = \\ &E[XY] - \alpha E[X] - \beta E[X^2] \\ &= E[XY] - (E[Y] - \beta E[X])E[X] - \beta E[X^2] \\ &= E[XY] - E[X]E[Y] + \beta E[X]^2 - \beta E[X^2] \\ &= \text{cov}[X, Y] - \beta V[X] \end{aligned}$$

$$\beta = \frac{\text{cov}[X, Y]}{V[X]}$$

$$\alpha = E[Y] - \frac{\text{cov}[X, Y]}{V[X]} E[X]$$

# **Reading: Multivariate Generalizations**

---

## READING: MULTIVARIATE GENERALIZATIONS

Read section 2.3, Multivariate Generalizations.



# Best Linear Predictors in Higher Dimensions

---

# LINEAR PREDICTORS IN HIGHER DIMENSIONS

## Definition: Linear Predictor

Given random variable  $Y$ , and random vector

$$\mathbf{X} = (1, X_1, X_2, X_3, \dots, X_k)$$

A *linear predictor* for  $Y$  is a function,  $g : \mathbb{R}^k \rightarrow \mathbb{R}$  of the form,

$$g(x_0, x_1, x_2, \dots, x_k) = \beta_0 x_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

such that  $\hat{Y} = g(1, X_1, X_2, \dots, X_k)$  is regarded as a guess for  $Y$ .

# INTERPRETING MODEL COEFFICIENTS

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$
$$\frac{\partial \hat{Y}}{\partial X_i} = \beta_i, \quad \Delta \hat{Y} = \beta_i \cdot \Delta X_i$$

## Ceteris Paribus: All else equal

- If  $X_i$  changes by  $\Delta X_i$ , the prediction  $\hat{Y}$  changes by  $\beta_i \cdot \Delta X_i$ , *if the other  $X$ 's are held constant.*

# INTERPRETING MODEL COEFFICIENTS - EXAMPLE

Does this model say peacocks with longer tails fly slower?

$$\widehat{\text{air\_speed}} = 4.3 - 1.2 \cdot \text{tail\_length} + 0.8 \cdot \text{muscle\_mass}$$



Photo by Thimindu Goonatillake CC BY-SA 2.0

# BEST LINEAR PREDICTORS IN HIGHER DIMENSIONS

## Definition: Best Linear Predictor

Given random variable  $Y$ , and random vector  $\mathbf{X} = (1, X_1, X_2, X_3, \dots, X_k)$

The *best linear predictor* for  $Y$  is the function,

$$g : \mathbb{R}^{k+1} \rightarrow \mathbb{R},$$

$g(x_0, x_1, x_2, \dots, x_k) = \beta_0 x_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$ , with coefficients given by,

$$\operatorname{argmin}_{\beta_0, \dots, \beta_k} \mathbb{E} \left[ \left( Y - (\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k) \right)^2 \right]$$

# MOMENT CONDITIONS IN HIGHER DIMENSIONS

Apply first-order conditions:

$$0 = \frac{\partial E[\epsilon^2]}{\partial \beta_0} = E \left[ \frac{\partial \epsilon^2}{\partial \beta_0} \right] = E \left[ 2\epsilon \frac{\partial \epsilon}{\partial \beta_0} \right] = -2E[\epsilon]$$

$$\text{For } j > 0, 0 = \frac{\partial E[\epsilon^2]}{\partial \beta_j} = E \left[ \frac{\partial \epsilon^2}{\partial \beta_j} \right] = E \left[ 2\epsilon \frac{\partial \epsilon}{\partial \beta_j} \right] = -2E[\epsilon X_j]$$

Version 1:

1.  $E[\epsilon] = 0$
2.  $E[\epsilon X_j] = 0$

Version 2:

1.  $E[\epsilon] = 0$
2.  $\text{cov}[\epsilon, X_j] = 0$

# THE BLP SOLUTION IN HIGHER DIMENSIONS

$$\beta = E[(\mathbf{X}^T \mathbf{X})^{-1}] E[\mathbf{X}^T \mathbf{Y}]$$

## LOOKING AHEAD

- The Best Linear Predictor is also called the Population Regression Function.
- Linear Regression: An algorithm for selecting a linear predictor given a sample of data.
- Ordinary Least Squares Regression: An algorithm for estimating the best linear predictor.