

slides:

dsc40a.com

Lecture 1

Introduction to Modeling

DSC 40A, Fall 2024

Agenda

- Introductions.
- What is DSC 40A about?
- Logistics.
- Modeling.
- The constant model.

Introductions

Instructor: Gal Mishne

- Assistant Professor at Halicioğlu Data Science Institute since summer 2019.
- Undergrad: EE and Physics at Technion.
- Grad school: EE PhD at Technion.
- Postdoc: Applied math at Yale
- Outside interests: traveling/hiking, cooking, reading, painting.

Course staff

We have 1 TA and 7 tutors, all of whom are excited to help you in discussion and office hours!

Sawyer Robertson

Rebecca (Jiaying) Chen

Zoe Ludena

Brighten Hayama

Utkarsh Lohia

Varun Pabreja

Javier Ponce

Owen Miller

Read more about us at dsc40a.com/staff.

Throughout lecture, ask questions!

- You're always free to ask questions during lecture, and I'll try and stop for them frequently. But still, you may not feel like asking your question out loud.
- You can **type your questions anonymously** at the following link and I'll try and answer them.

q.dsc40a.com

- You'll also use this form to answer questions that I ask you during lecture.
- If the direct link doesn't work, use the 🤔 **Lecture Questions** link in the top right corner of dsc40a.com.



Ed



Gradescope



Practice



Lecture Questions

What is DSC 40A about?

Theoretical Foundations of Data Science I

What have you *heard* about DSC 40A?

Here are some responses from the Welcome Survey.

That it was very math and proof heavy, not much coding involved at all.

It's a hard class that is very conceptual and math heavy... linear algebra based... very hard statistics class...

I heard the homework has been very difficult and getting help is pretty hard. I also heard exams are scary difficult with little practice given.

→ over 15 hours of OH

→ check practice.dsc40a.com

I have heard that it is very math intensive and very difficult, but well worth it because you learn so much content and information for the future. !

Why do we need to study theoretical foundations?



Machine learning is about **automatically learning patterns from data.**

Humans are good at understanding handwriting – but how do we get computers to understand handwriting?

Course overview

Part 1: Learning from Data (Weeks 1-5)

- Summary statistics and loss functions; empirical risk minimization.
- Linear regression (including multiple variables); linear algebra.

Part 2: Probability (Weeks 6-10)

- Set theory and combinatorics; probability fundamentals.
- Conditional probability and independence.
- The Naive Bayes classifier.

Learning objectives

After this class, you'll...

- understand the basic principles underlying almost every machine learning and data science method.
- be better prepared for the math in upper division: vector calculus, linear algebra, and probability.

What do DSC 80 students have to say about DSC 40A?

Here are some responses from the End-of-Quarter Survey last quarter in DSC 80.

study hardy, pay attention in DSC 40A and start work early :)

40A and Math 18 is super important for this class. Don't wait till the last minute too!

I think DSC40[A] was the most important prerequisite for this class.

Logistics

Getting started

- The course website, dsc40a.com, contains all content. **Read the syllabus carefully!**
 - Click around; you'll find other helpful resources.
- Other sites you'll need to use:
 - **Gradescope** is where you'll submit all assignments. You'll be automatically added within 24 hours of enrolling.
 - **Ed** is where all announcements will be made. If you're not enrolled, there's a join link in the syllabus.
 - We aren't using Canvas.
- Make sure to fill out the **Welcome Survey** ASAP.

Lectures

Pepper Canyon 106

- Lecture is here, ~~Cantor Hall 105~~, MWF 4:00-4:50p.
- Lecture slides will be posted on the course website before class, and annotated slides will be posted after class.
- Lecture will be podcasted.
- **The value of lecture is interaction and discussion, so even though attendance isn't required, it's highly, highly recommended.**

Discussions

- Discussion weekly on Mondays, Center Hall 212, 6:00-6:50p.
- Discussion will primarily be used for **groupwork** – that is, working on problems in small groups of size 2-4.
 - You may work in a self-organized group outside of a discussion section for 80% credit, but no matter what, **you cannot work alone.**
- Groupwork worksheets are due to Gradescope on **Mondays at 11:59PM.**
 - Only one group member needs to submit, and should add the rest of the group to the submission.
- **The value of attending is getting support from the TA and tutors.**

Grading

- **Homeworks (40%):** Due to Gradescope **Fridays at 11:59p.**
 - Graded for correctness. Lowest score is dropped.
- **Groupworks (10%):** Due to Gradescope on **Mondays at 11:59p.**
 - Graded for effort. Lowest score is dropped.
- **Midterm Exam (20%):** Monday, November 4th, in class.
- **Final Exam (30%):** Tuesday, December 10th, 3pm. See the [syllabus](#) for the redemption policy.

Support

>15 hours of OH

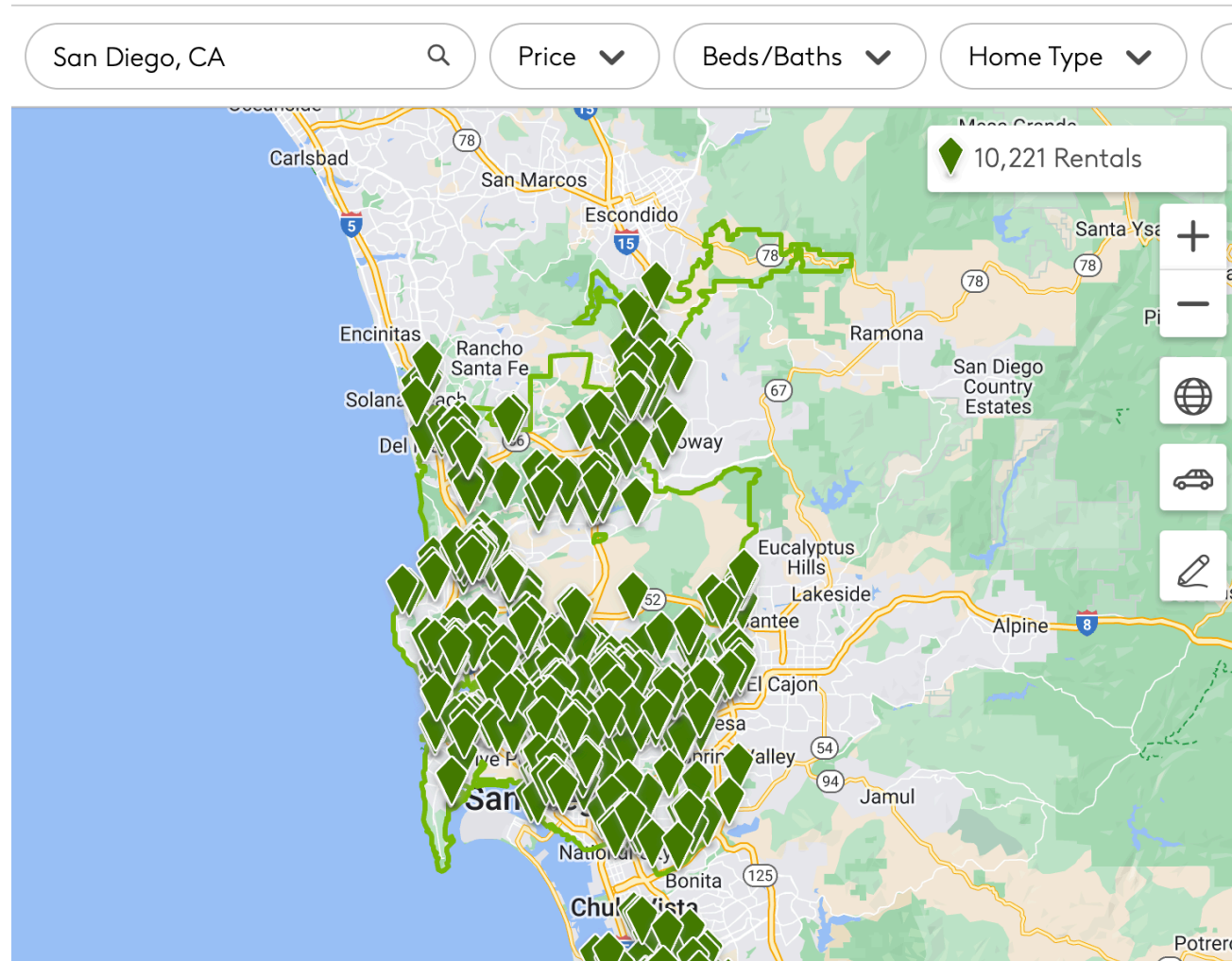
We know this is a challenging class, and we're here to help:

- **Office hours:** In-person in HDSI 155. Plan to attend at least once a week for homework help.
- **Ed:** Use it! We're here to help you. Post conceptual questions publicly – just don't post answers to homework questions.

A bunch of new-ish things to improve the student experience:

- practice.dsc40a.com to give you access to practice exam problems, categorized by topic.
- Walkthrough videos to show you our thought process when answering questions.
- More time reviewing linear algebra.

Modeling



You might be starting to look for off-campus apartments, none of which are affordable.

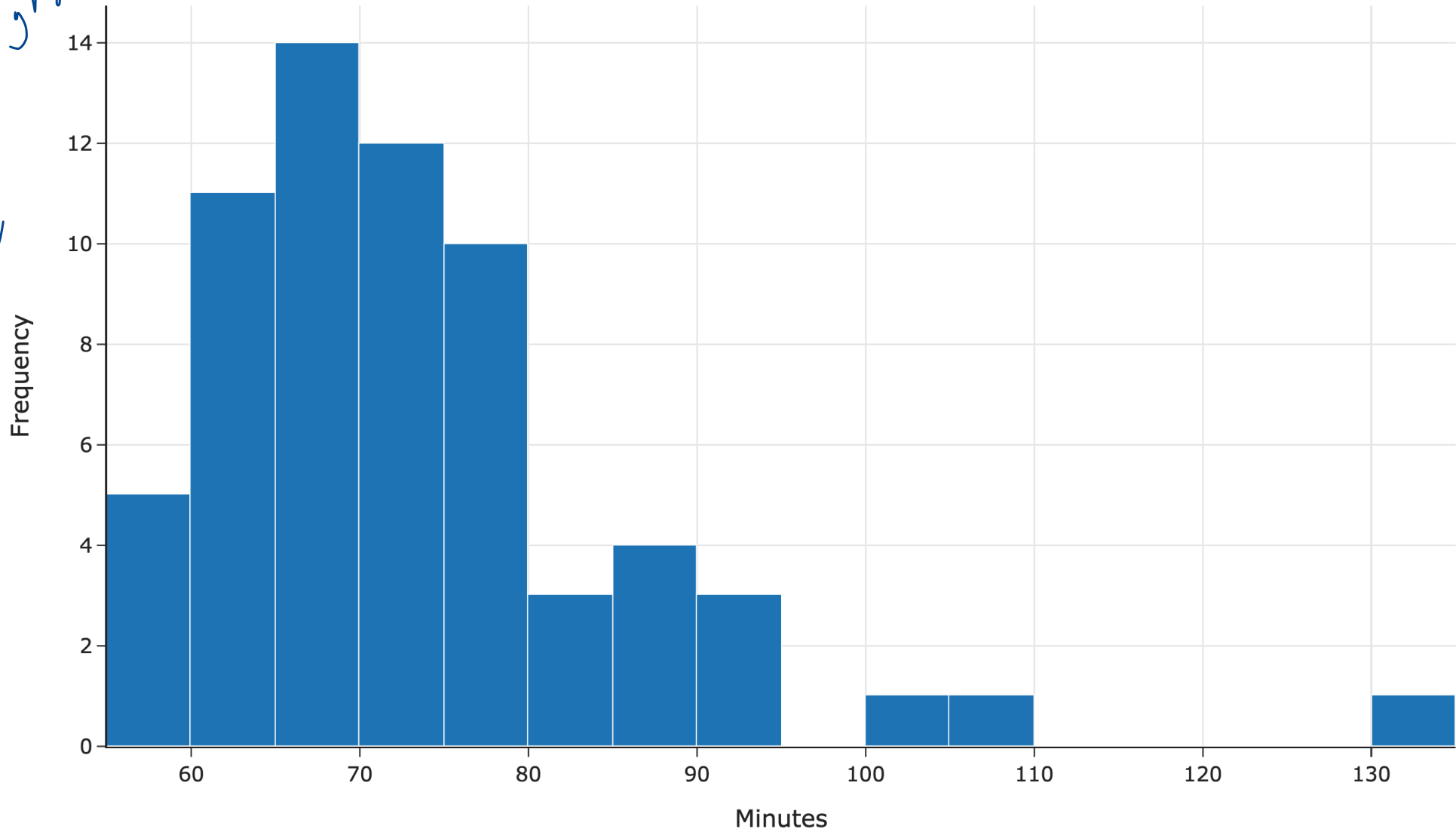
	date	day	departure_hour	minutes
0	5/22/2023	Mon	8.450000	63.0
1	9/18/2023	Mon	7.950000	75.0
2	10/17/2023	Tue	10.466667	59.0
3	11/28/2023	Tue	8.900000	89.0
4	2/15/2024	Thu	8.083333	69.0

...

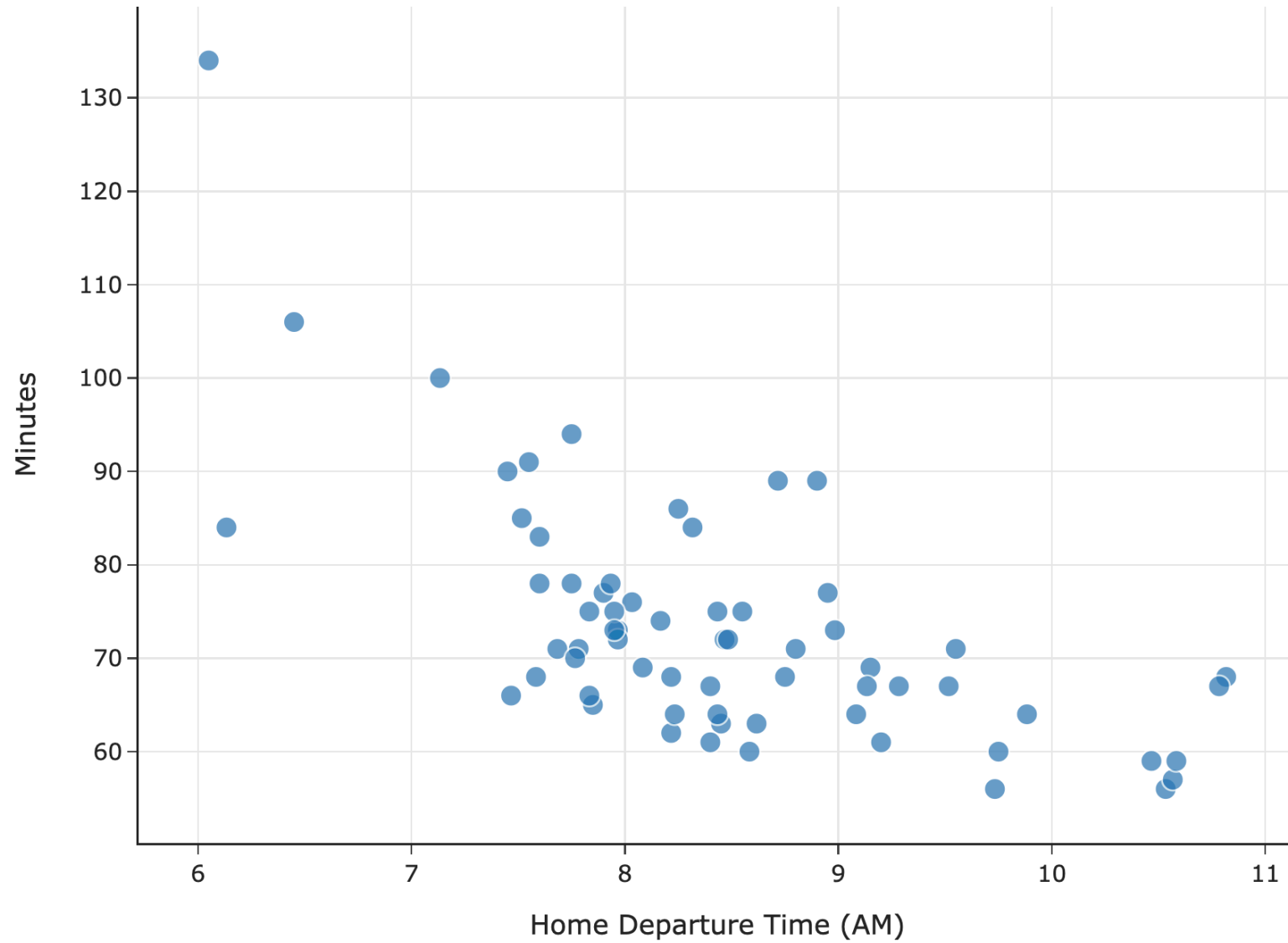
You decide to live with your parents in Orange County and commute.
You keep track of how long it takes you to get to school each day.

not a density histogram (doesn't sum to 1)

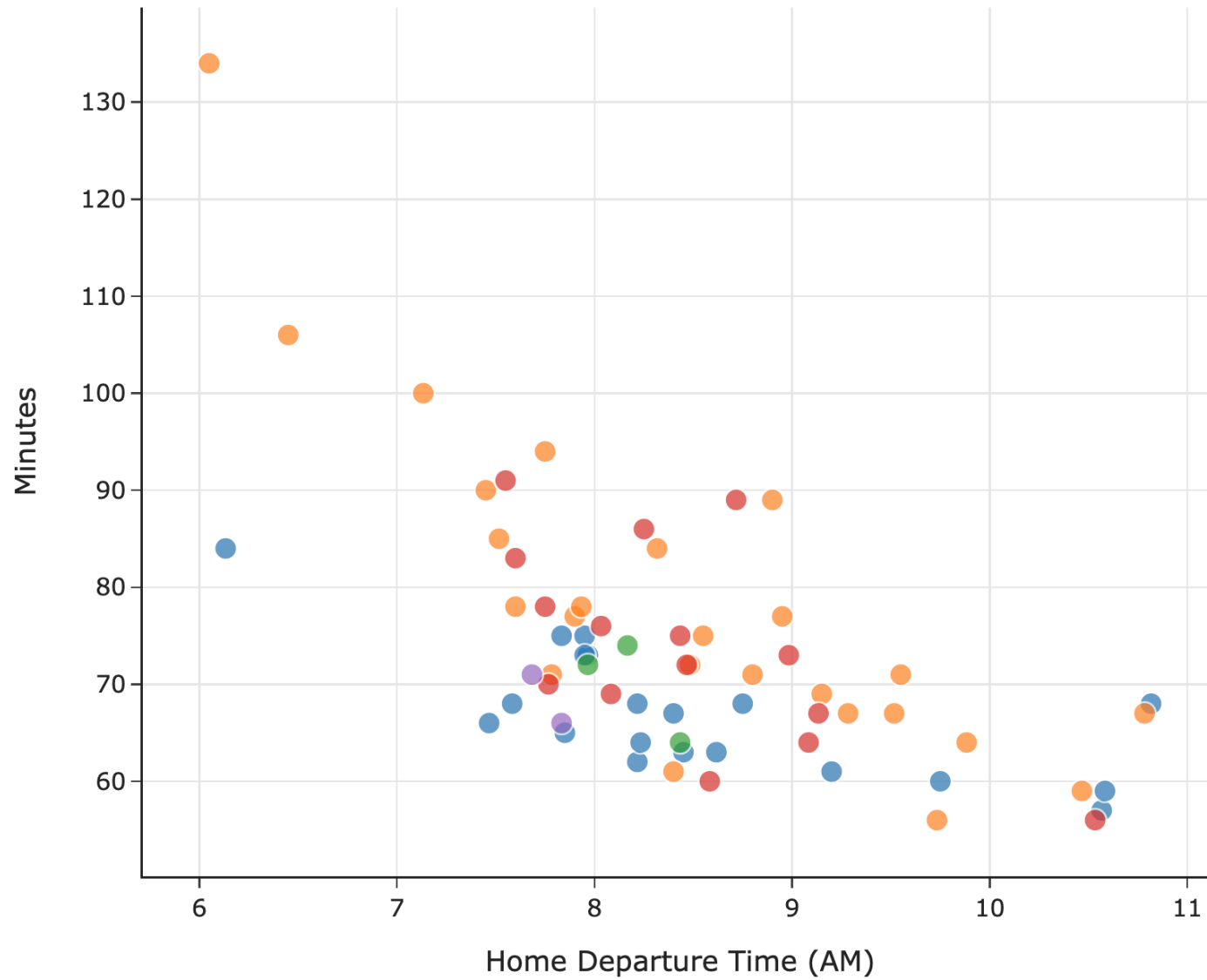
Distribution of Commuting Time



Commuting Time vs. Home Departure Time



Commuting Time vs. Home Departure Time



day

- Mon
- Tue
- Wed
- Thu
- Fri

↑
Commute patterns
may vary by
day of the week

Goal: Predict your commute time.

That is, predict how long it'll take to get to school.

data in the
future looks like
data in the
past

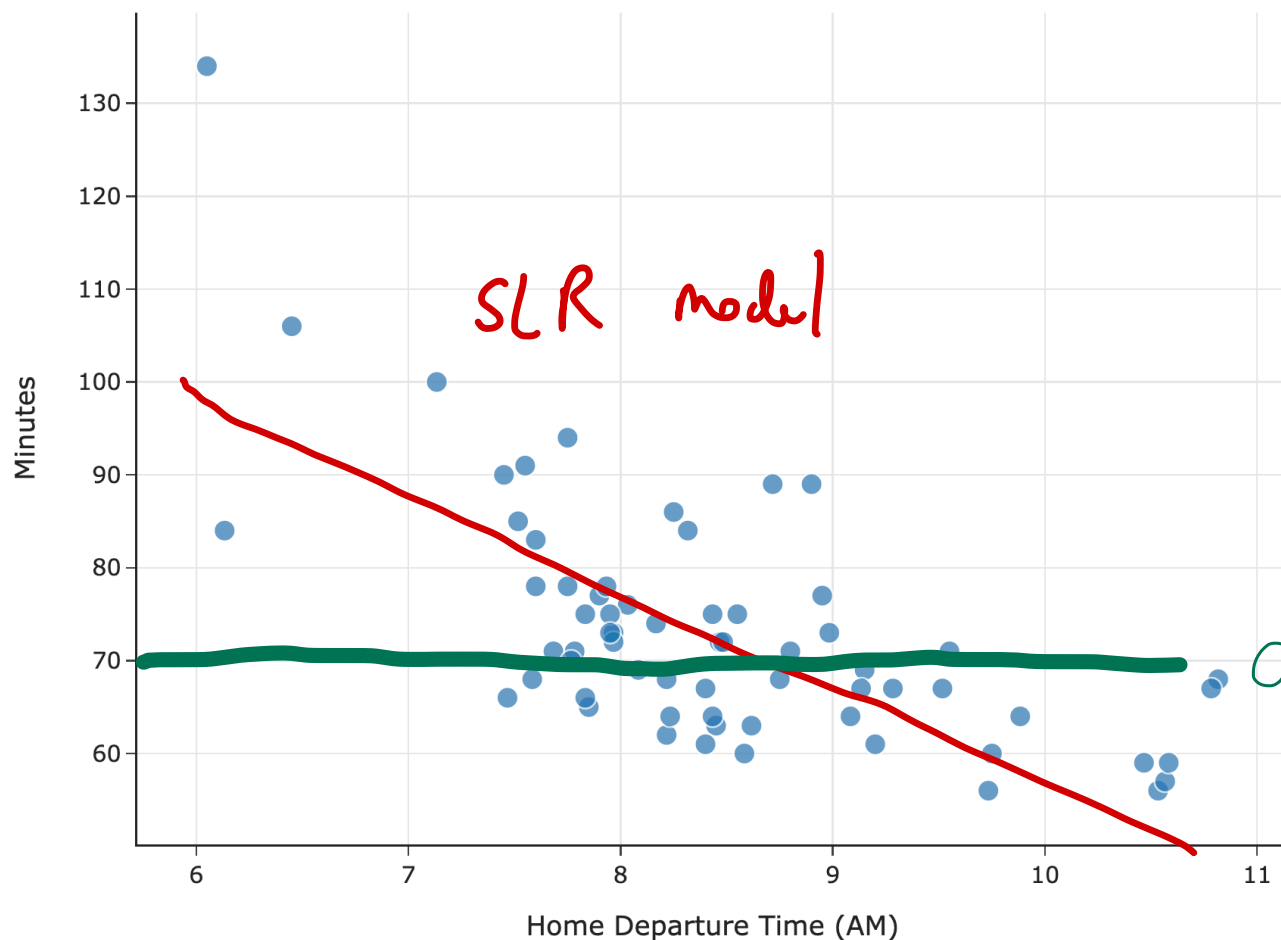
How can we do this?
What will we need to assume?

learn patterns from
data

A **model** is a set of assumptions about how data were generated.

Possible models

Commuting Time vs. Home Departure Time



SLR - simple linear regression model

constant model

constant