

ECON 0150 | Economic Data Analysis

The economist's data analysis pipeline.

Part 2.1 | Data Cleaning

Data Cleaning

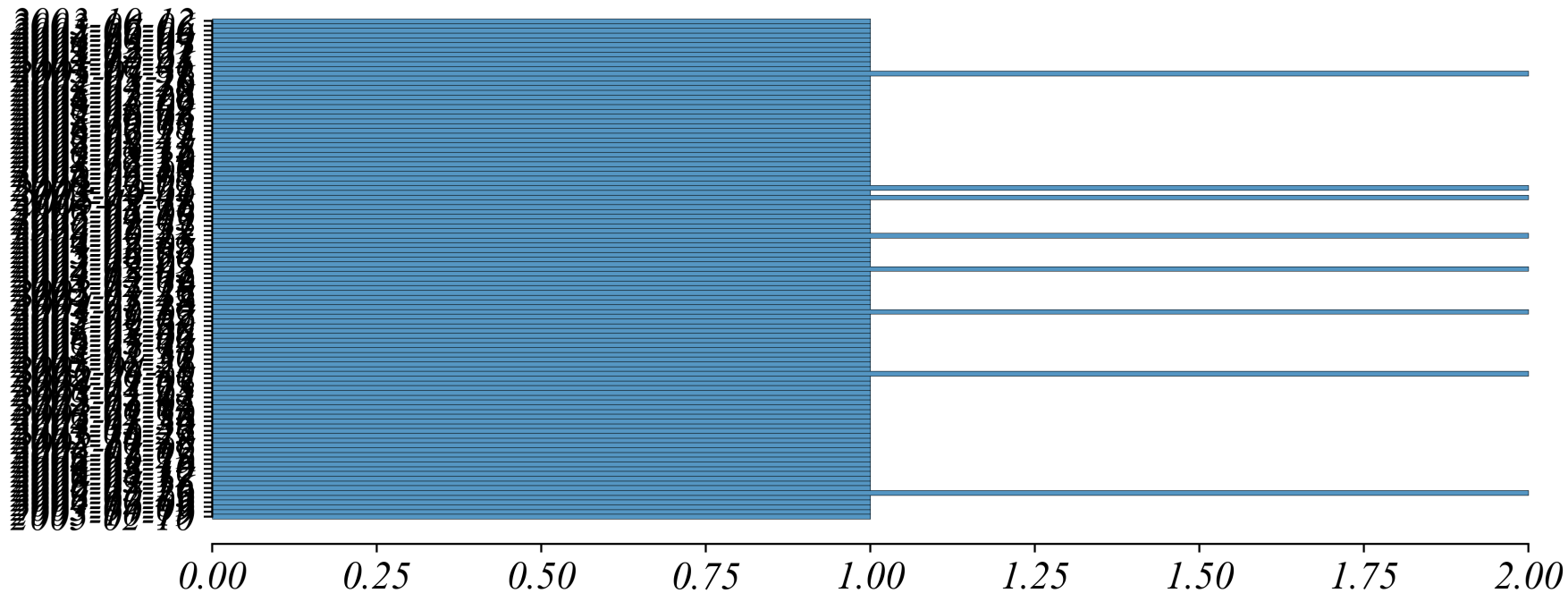
Q. Are students who live further away older?

Data Cleaning

Q. Are students who live further away older?

Let's examine age and distance from Pittsburgh.

When is your birthday?

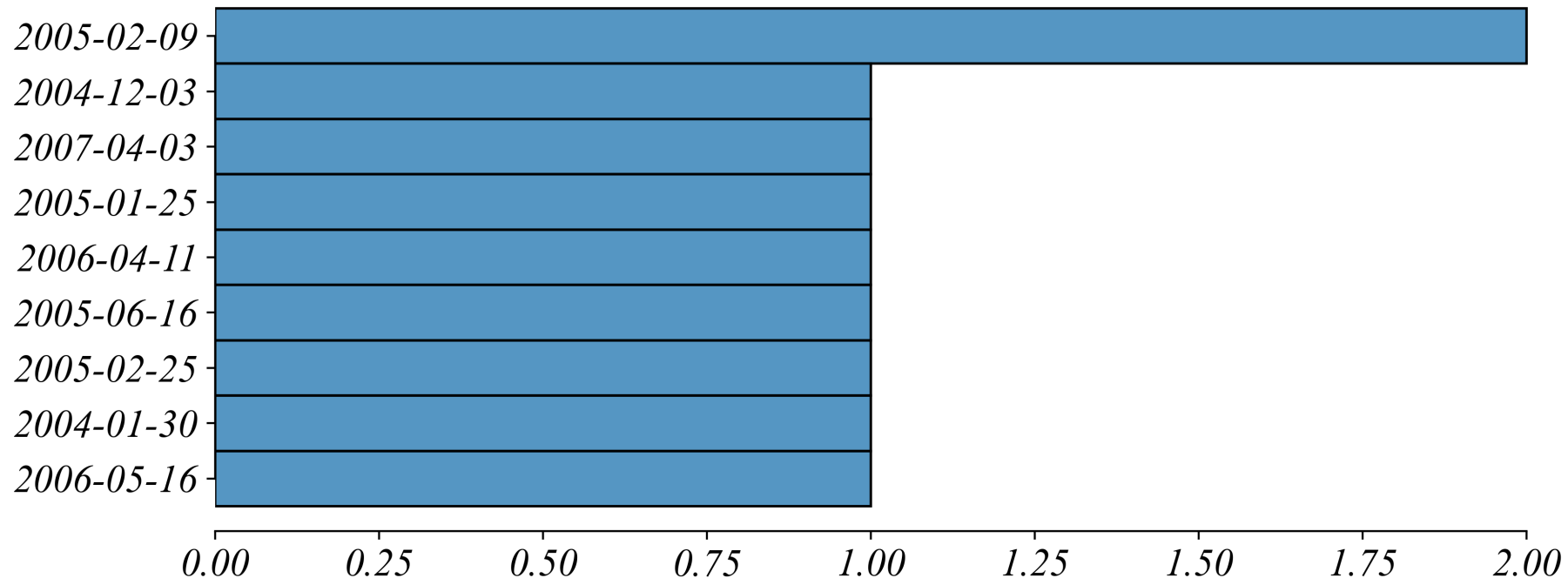


Data Cleaning

Q. Are students who live further away older?

Let's examine age and distance from Pittsburgh.

When is your birthday?



> *the birthday data is stored as text: “08/15/2005”*

> *we need to extract the year to calculate age*

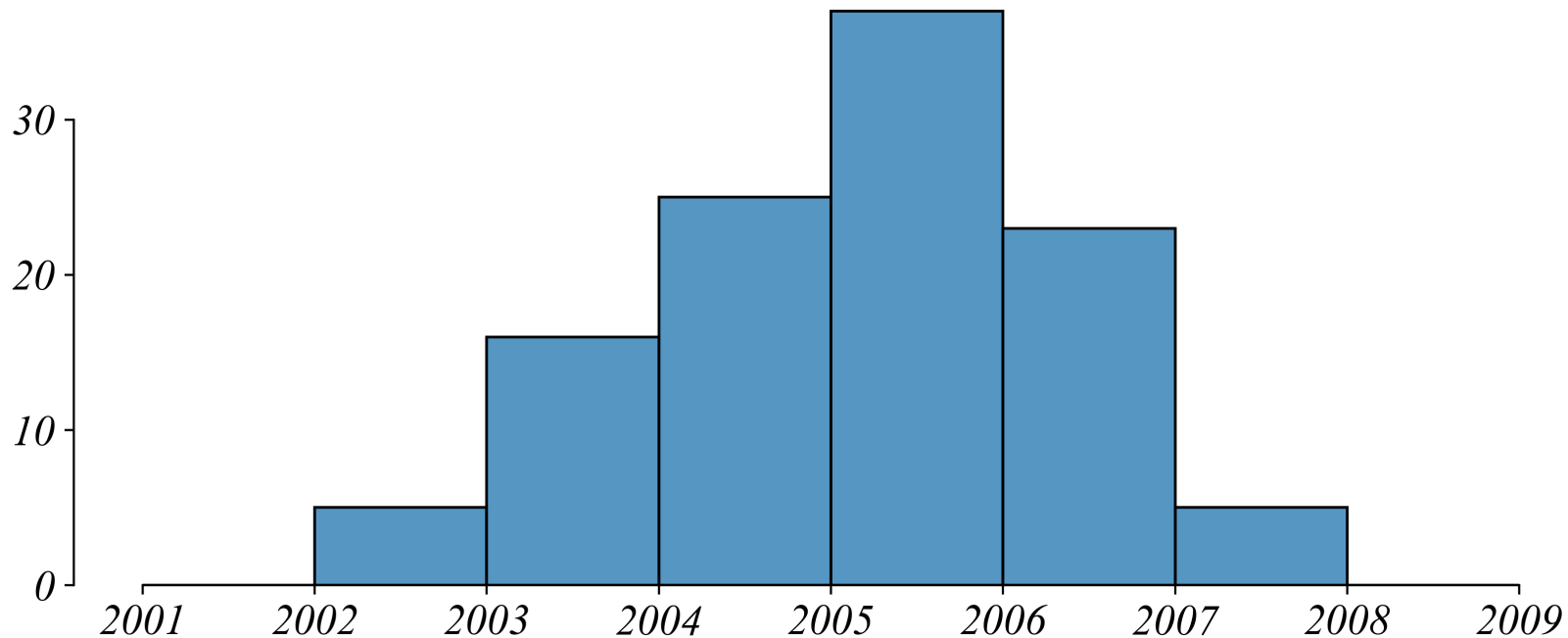
String Parsing

Extracting useful information from text

What we have: “08/15/2005”

What we need: 2005

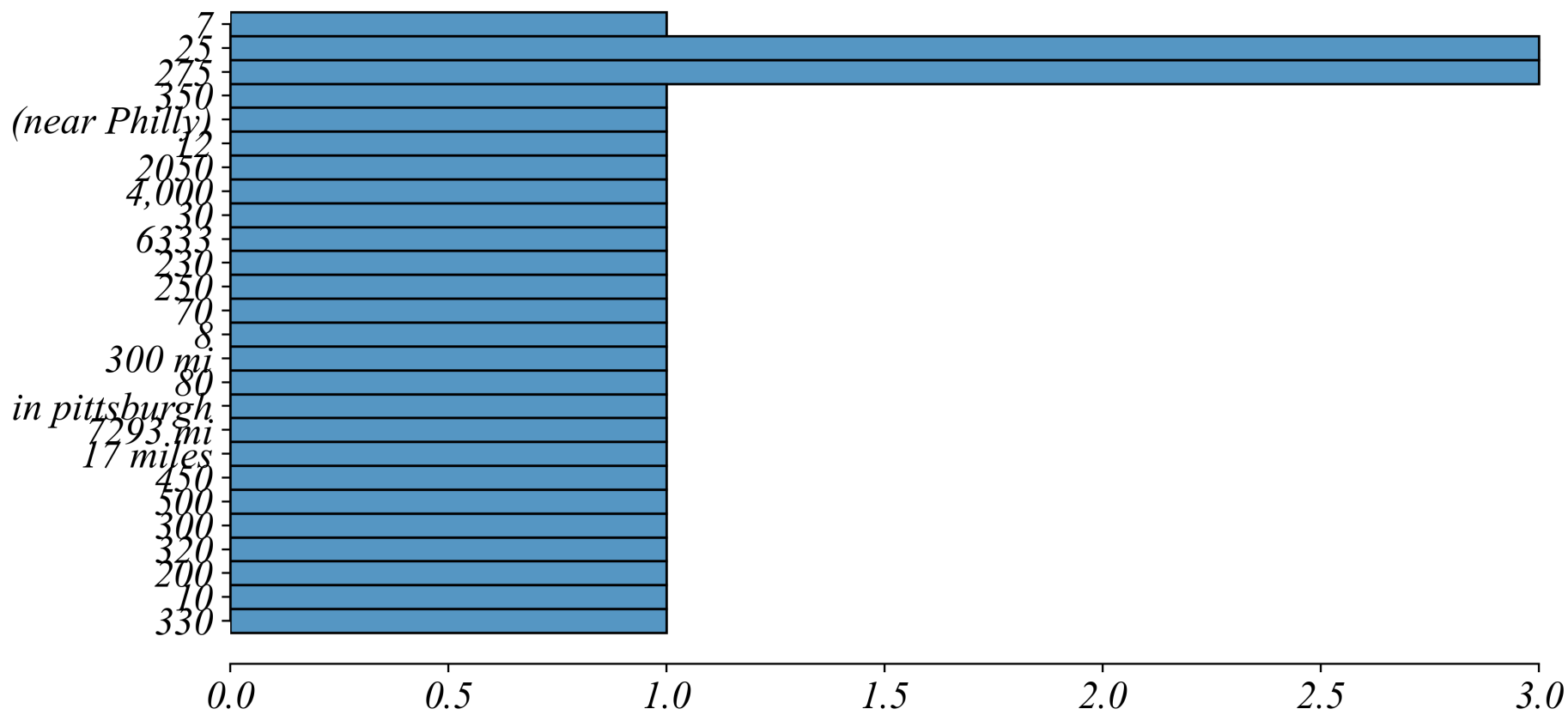
When is your birthday?



Distance from Pittsburgh

Q. Are students who live further away older?

How many miles away from Pittsburgh is your hometown?



> *lots of different formats!*

Distance from Pittsburgh

Answers can be in many creative forms...

- *“0 miles”*
- *“~500”*
- *“about 1000”*
- *“2.5 hours”*
- *“very far”*

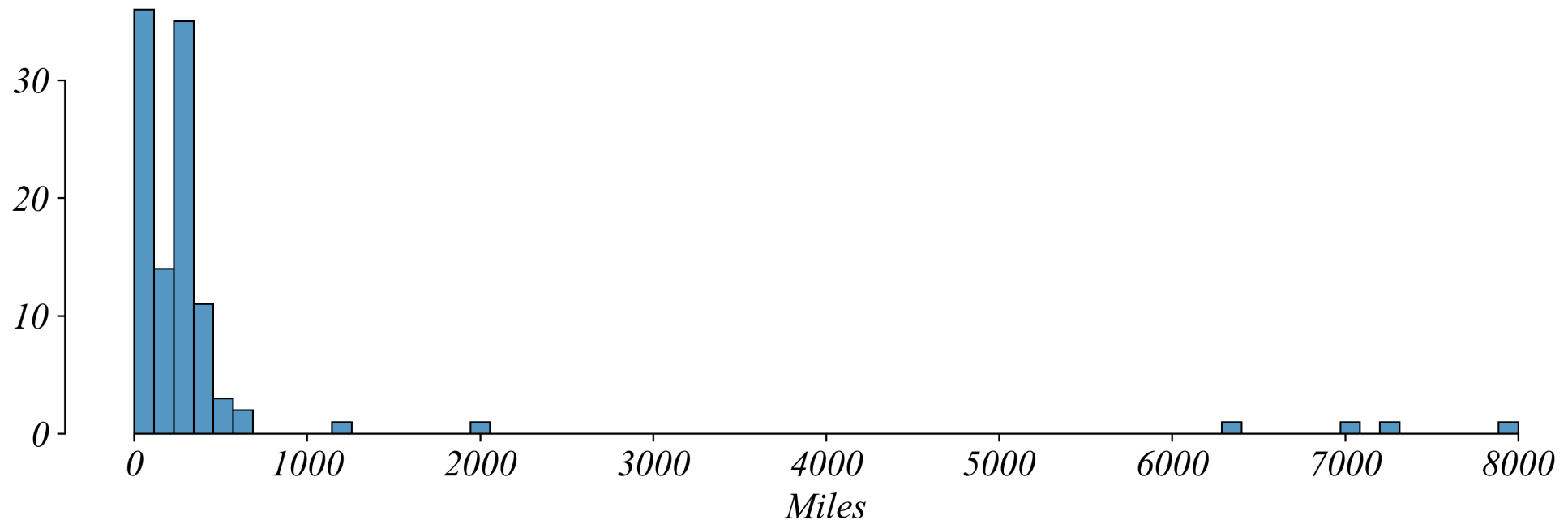
> computers can't do math with text

Type Conversion

Converting text to numbers

We can convert text to numbers, forcing errors to become NA.

How many miles away from Pittsburgh is your hometown?



> *entries like “very far” become NA*

> *entries like “500” become 500.0*

Missing Values

What happened to the non-numeric entries?

new Approximately how many miles away from Pittsburgh is your hometown?		
0	400.0	400
1	16.0	16
2	300.0	300
3	300.0	300
4	400.0	400

Missing Values

What happened to the non-numeric entries?

new **Approximately how many miles away from Pittsburgh is your hometown?**

6	NaN	176 miles away
---	-----	----------------

17	NaN	0 (it's Pittsburgh)
----	-----	---------------------

18	NaN	400-450ish miles
----	-----	------------------

22	NaN	350 miles
----	-----	-----------

23	NaN	240 miles
----	-----	-----------

> *they all became NaN (Not Available)*

> *we need to decide what to do with them*

Handling Missing Values

Two main approaches

After replacing problematic values, there are generally two options.

Option 1: Drop the missing values

- *Removes entire rows with NA*
- *Reduces sample size*
- *Simple and clean*

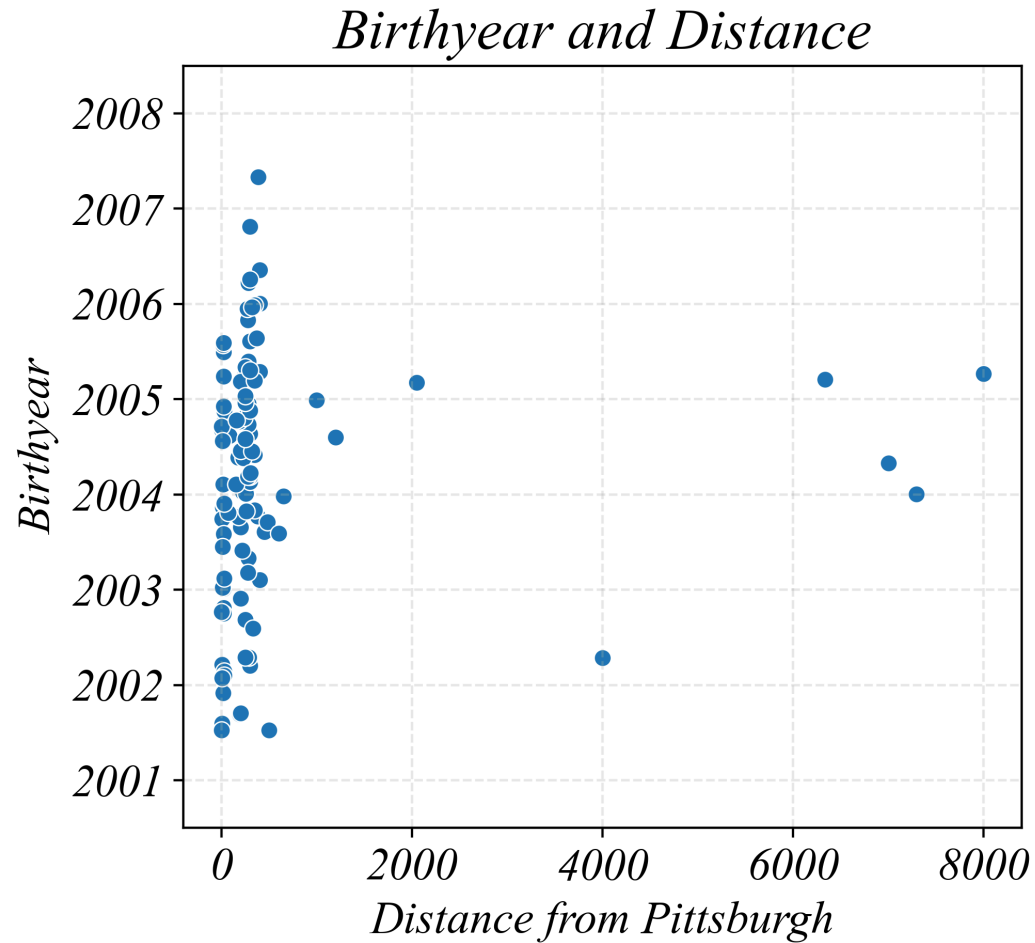
Option 2: Replace with a value

- *Fill with 0, mean, or median*
- *Keeps sample size*
- *May introduce bias*

> for distance, dropping makes sense - we can't guess locations

After Cleaning

Q. Are students who live further away older?



> as expected, there does not seem to be much of a relationship

Summary

Some common data cleaning operations

- *String Parsing: Extract information from text*
- *Type Conversion: Change text to numbers*
- *Missing Values: Drop or replace NAs*

Exercise 2.1 | Data Cleaning

Let's find the median birthyear and the mean hometown distance from Pittsburgh.

- *Data:* `Fall_2025_Survey_raw.csv`

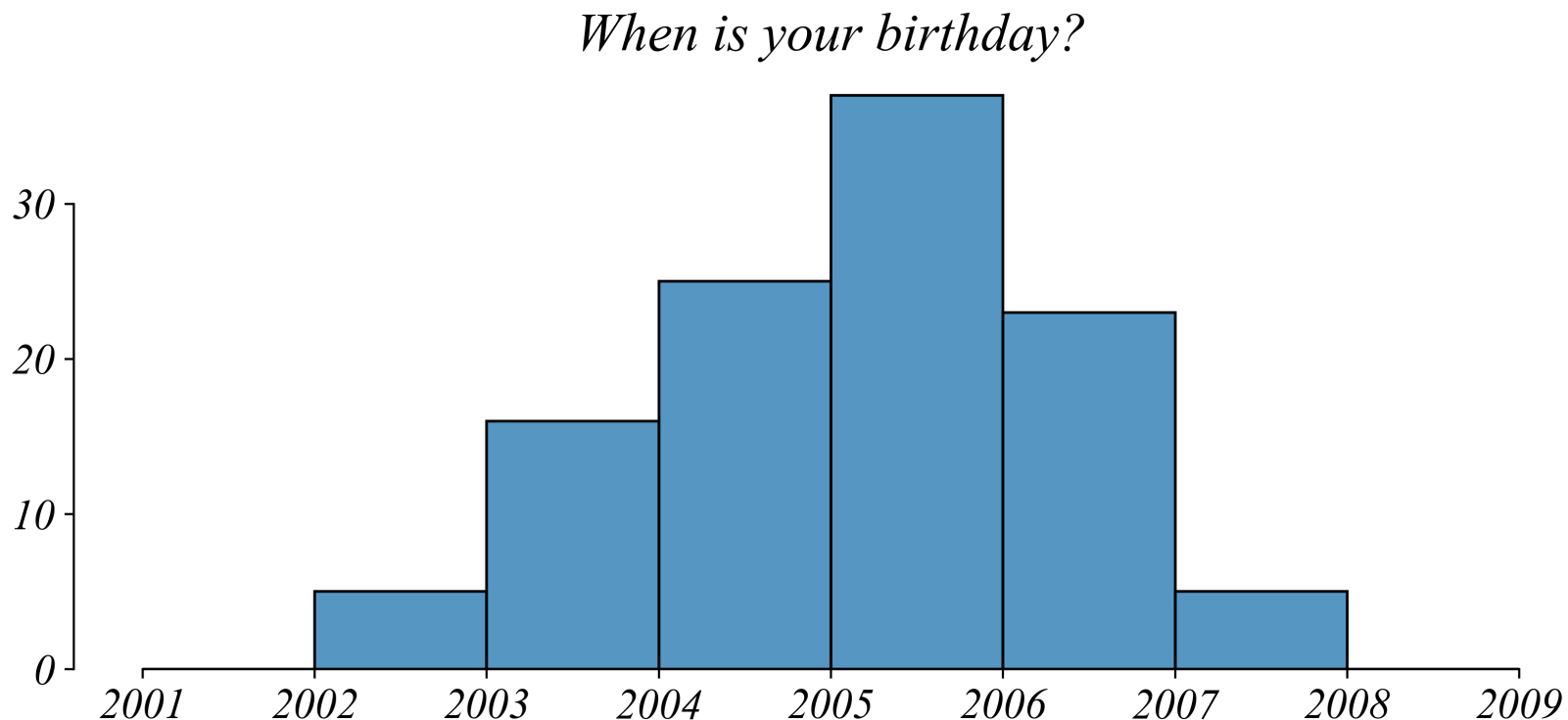
Exercise 2.1 | Birthday to Birthyear

Extract year from birthday text

```
1 # Convert birthday to datetime
2 survey['birthday'] = pd.to_datetime(survey['When is your birthday?'])
```

```
1 # Extract year from date
2 survey['birthyear'] = survey['birthday'].dt.year
```

```
1 sns.histplot(survey, x='birthyear')
```



Exercise 2.1 | Distance Conversion (Simple)

Convert distance text to numbers

```
1 # Convert to numeric, errors become NA
2 survey['distance'] = pd.to_numeric(survey['Approximately how many miles away from Pittsburgh i
```

```
1 # Check how many became NA
2 survey['distance'].isna().sum()
```


Exercise 2.1 | Handle Missing Values

Two approaches to NAs

Drop missing values:

```
1 # Remove rows where distance is NA
2 survey_dropped = survey.dropna(subset=['distance_clean'])
```

Replace with a value:

```
1 # Replace NA with median distance
2 median_dist = survey['distance_clean'].median()
3 survey['distance_filled'] = survey['distance_clean'].fillna(median_dist)
```

```
1 # Or replace with 0
2 survey['distance_zero'] = survey['distance_clean'].fillna(0)
```

Exercise 2.1 | Distance Conversion (Replace)

Convert distance text to numbers

```
1 # Replace non-numeric
2 replacements = {
3     '400-450ish miles ': 400,
4     'live in pittsburgh': 0,
5     '176 miles away': 176,
6     '0 (it's Pittsburgh)': 0,
7     '350 miles': 350,
8     '240 miles': 240,
9     '388 miles': 388,
10    '17 miles': 17,
11    '300 miles': 300,
12    '7293 mi': 7293,
13    '4 miles ': 4,
14    '27 miles': 27,
15    '255 (near Philly)': 255,
16    '4,000': 4000,
17    '650 miles': 650,
18    '250 miles': 250,
19    '318 mi': 318,
20    '300 mi': 300,
21    '1000+': 1000,
22    '305 miles': 305
```

```
1 # Check how many became NA
2 survey['distance'].isna().sum()
```

Exercise 2.1 | Scatterplot

Check that it worked

```
1 # Create scatterplot
2 sns.scatterplot(data, x='distance', y='age')
```

