

CHAPTER 10

localhost:8888/notebooks/Untitled47.ipynb?kernel_name=py38

jupyter Untitled47 Last Checkpoint: 3 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3.8 (py38) Logout

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# Make plots a bit bigger
plt.rcParams["figure.figsize"] = (6, 4)

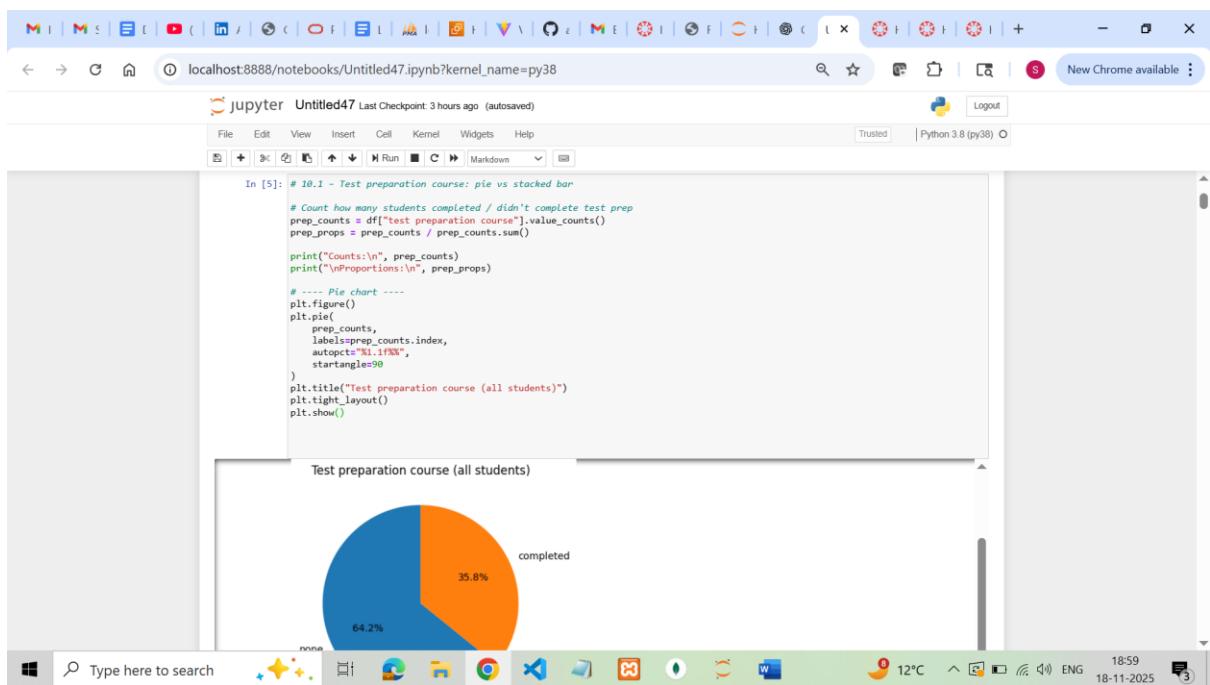
# Load StudentsPerformance.csv directly from GitHub
url = "https://raw.githubusercontent.com/abhijitpaul0212/DS-ML-DataSet/main/StudentsPerformance.csv"
df = pd.read_csv(url)

# Quick check
print(df.head())
print(df.columns)
```

	gender	race/ethnicity	parental level of education	lunch	test preparation course	math score	reading score	writing score
0	female	group B	bachelor's degree	standard	none	72	72	74
1	female	group C	some college	standard	completed	69	90	88
2	female	group B	master's degree	standard	none	98	95	93
3	male	group A	associate's degree	free/reduced	none	47	57	44
4	male	group C	some college	standard	none	76	78	75

Index(['gender', 'race/ethnicity', 'parental level of education', 'lunch', 'test preparation course', 'math score', 'reading score', 'writing score'], dtype='object')

Type here to search 12°C ENG 18-11-2025



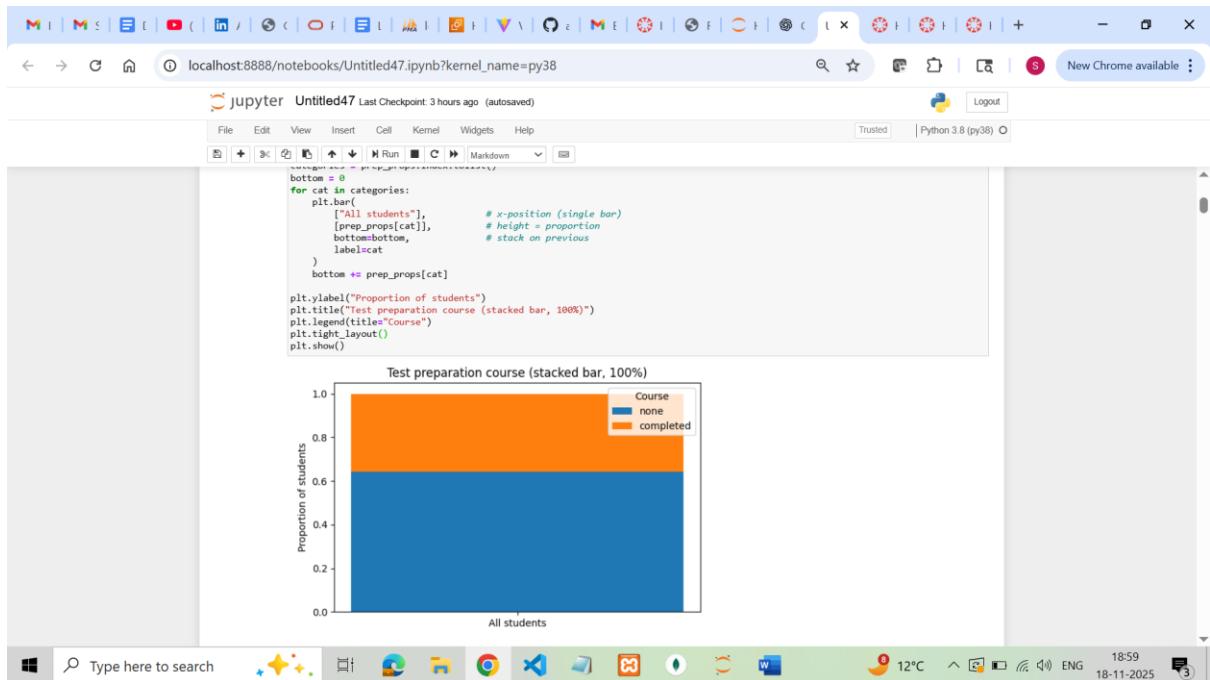
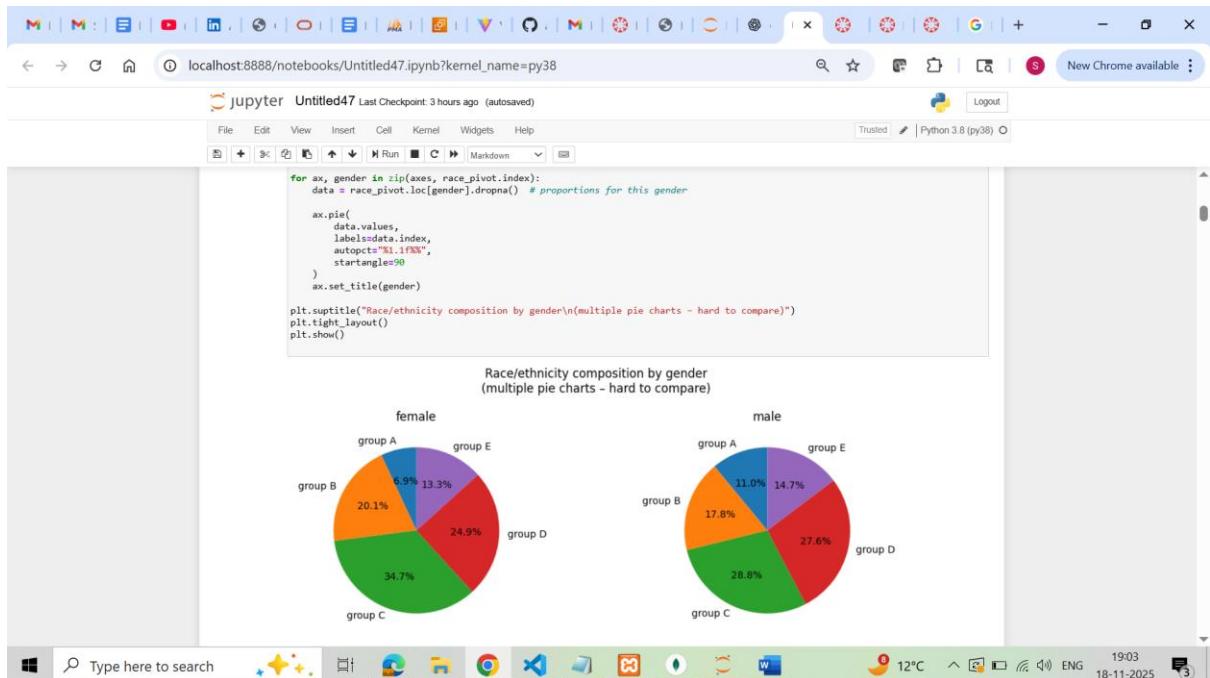


Figure 10.1a – Pie chart of test preparation course

This pie chart shows the overall proportion of students who completed the test-preparation course versus those who did not. Most students did not complete test prep, while a smaller slice completed it. With only two categories and a single set of proportions, the pie chart works reasonably well and emphasizes simple fractions of the whole.

Figure 10.1b – 100% stacked bar of test preparation course

This stacked bar encodes the same information as the pie chart: the bar height is 100% of students and the segments show the share with and without test prep. It reinforces that stacked bars and pies both show “part of a whole,” but for one single group the stacked bar feels less natural than the pie.



```

In [3]: # 10.2 - Race/ethnicity proportions by gender: side-by-side bars

race_by_gender = (
    df.groupby(["gender", "race/ethnicity"])
    .size()
    .reset_index(name="count")
)

race_by_gender["prop"] = (
    race_by_gender
    .groupby("gender")["count"]
    .transform(lambda x: x / x.sum())
)

print(race_by_gender.head())

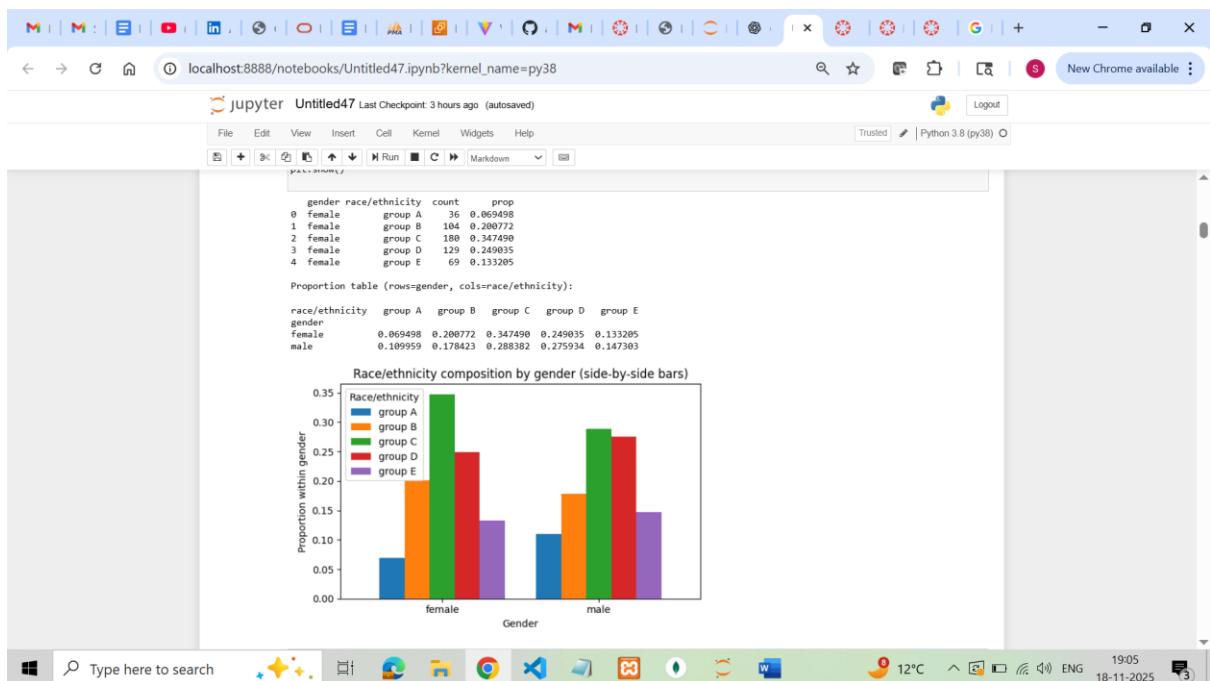
# Pivot to gender (rows) x race/ethnicity (columns)
race_pivot = race_by_gender.pivot(
    index="gender",
    columns="race/ethnicity",
    values="prop"
)

print("\nProportion table (rows=gender, cols=race/ethnicity):\n")
print(race_pivot)

# Plot: side-by-side (grouped) bars
ax = race_pivot.plot(
    kind="bar",
    stacked=False,
    width=0.8
)

ax.set_xlabel("Proportion within gender")
ax.set_ylabel("Gender")
ax.set_title("Race/ethnicity composition by gender (side-by-side bars)")

```



10.2 – Multiple pies (bad) vs side-by-side bars (good)

Figure 10.2a – Multiple pie charts by gender (bad design)

Here each pie shows the race/ethnicity composition within one gender. While each individual pie is readable, it is very hard to compare slices across pies (e.g., “Is group B more common for males or females?”). This matches the lecture example where multiple pies over years are labeled as a poor design for comparisons.

Figure 10.2b – Side-by-side bars of race/ethnicity by gender (good design)

This grouped bar chart shows the same information, but now the race/ethnicity categories are bars on a common scale within each gender. Differences between genders are easy to compare by looking

at bar heights. This illustrates why side-by-side bars are preferred when the goal is to compare proportions across groups.

localhost:8888/notebooks/Untitled47.ipynb?kernel_name=py38

jupyter Untitled47 Last Checkpoint: 3 hours ago (autosaved)

In [8]: # 10.3(g) - Stacked bars: test prep vs math score bands

```
# Define math score bands
bins = [0, 50, 70, 90, 100]
labels = ["<50", "50-69", "70-89", "90-100"]

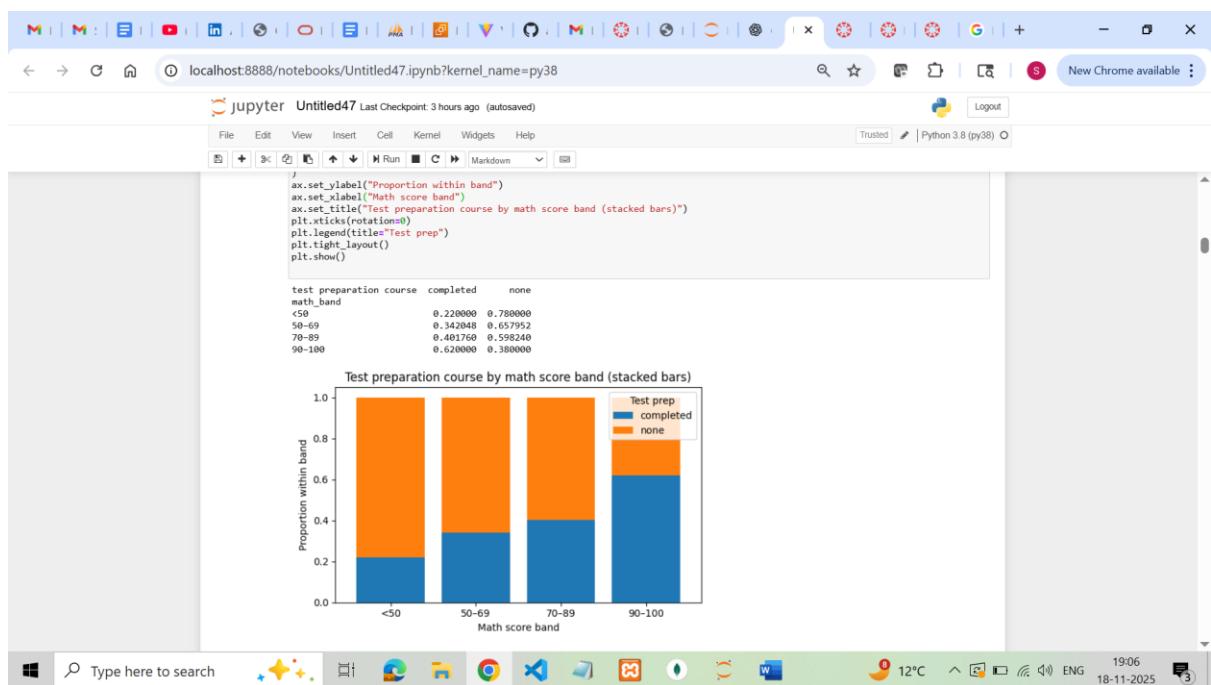
df["math_band"] = pd.cut(
    df["math score"],
    bins=bins,
    labels=labels,
    include_lowest=True
)

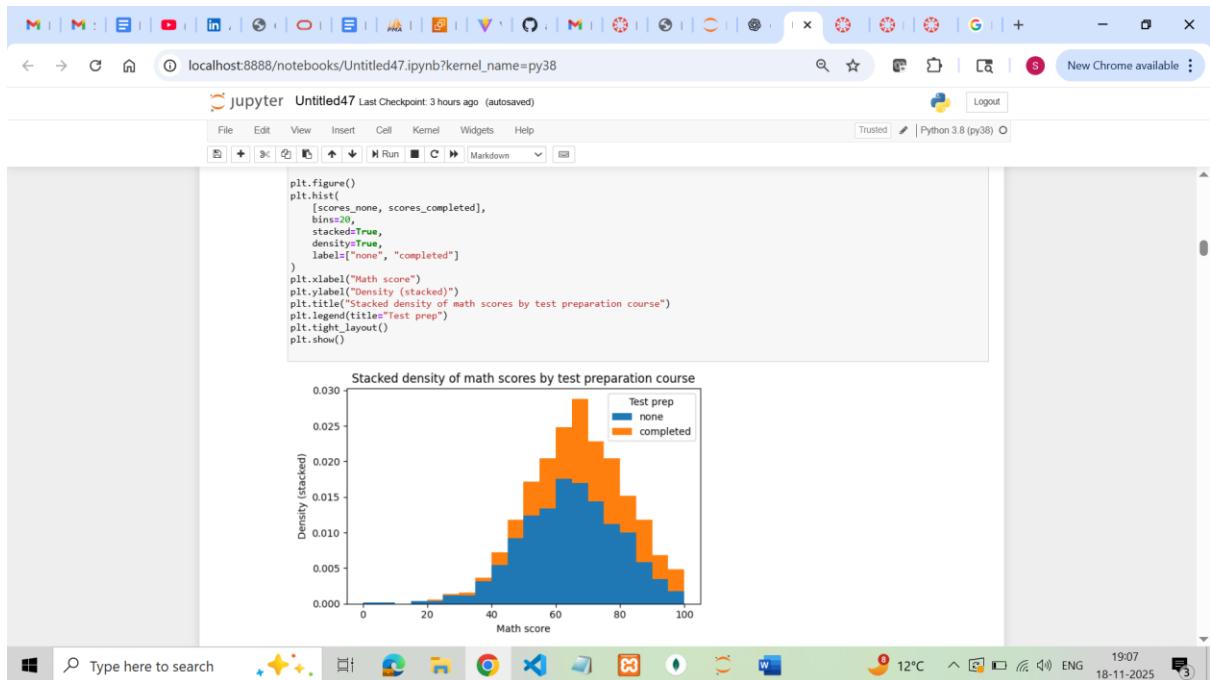
# Crosstab: proportion of test-prep categories within each band
band_prepable = pd.crosstab(
    df["math_band"],
    df["test preparation course"],
    normalize="index"
)

print(band_prepable)

# Plot stacked bars (100% within each band)
ax = band_prepable.plot(
    kind="bar",
    stacked=True,
    width=0.8
)
ax.set_xlabel("Proportion within band")
ax.set_xlabel("Math score band")
ax.set_title("Test preparation course by math score band (stacked bars)")
plt.xticks(rotation=0)
plt.legend(title="Test prep")
plt.tight_layout()
plt.show()
```

test preparation course completed none





10.3 – Stacked bars and stacked densities

Figure 10.3a – 100% stacked bars of test prep by math score band

Math scores are split into bands (<50, 50–69, 70–89, 90–100), and within each band the bar is stacked into “none” vs “completed” test prep. Each bar sums to 100% of students in that band. Higher score bands have a larger share of students who completed test prep, while lower bands have more students with no prep. This is similar to the textbook’s example of stacked bars showing composition over conditions.

Figure 10.3b – Stacked histogram (stacked density) of math scores

The stacked histogram shows the overall density of math scores, with contributions from the “none” and “completed” test-prep groups stacked on top of each other. At low scores the density is dominated by students with no test prep, whereas at high scores it is dominated by students who completed test prep. This is the stacked-density idea from the chapter, showing how category composition changes along a continuous axis.

```

In [10]: # 10.4 - Partial distributions:
# show each test-prep group as part of the total math-score distribution

N = len(df)
values = df["math score"].values

# Common binning for all histograms
bin_edges = np.linspace(0, 100, 21) # 20 bins from 0 to 100
bin_centers = 0.5 * (bin_edges[1:] + bin_edges[:-1])
bin_width = bin_edges[1] - bin_edges[0]

# Overall distribution: proportion of ALL students in each bin
overall_counts, _ = np.histogram(values, bins=bin_edges)
overall_prop = overall_counts / N # sums to 1

# Helper: partial proportions for a subset (still divided by TOTAL N)
def partial_prop(subset):
    sub_vals = subset["math score"].values
    sub_counts, _ = np.histogram(sub_vals, bins=bin_edges)
    return sub_counts / N

groups = ["none", "completed"]
fig, axes = plt.subplots(1, len(groups), figsize=(10, 4), sharey=True)

for ax, g in zip(axes, groups):
    subset = df[df["test preparation course"] == g]
    part_prop = partial_prop(subset)

    # Overall distribution in background
    ax.bar(
        bin_centers,
        overall_prop,
        width=bin_width,
        edgecolor="none",
        label="All students"
    )

    # Group-specific partial distribution on top (narrower)
    ax.bar(
        bin_centers,
        part_prop,
        width=bin_width * 0.6,
        edgecolor="none",
        label=g
    )

    ax.set_title(f"Test prep = {g}")
    ax.set_xlabel("Math score")

axes[0].set_ylabel("Proportion of ALL students in bin")
axes[0].legend()
plt.suptitle("Math scores by test preparation course\n(n=partial distributions of total)")
plt.tight_layout()
plt.show()

```

The figure consists of two side-by-side histograms. The left histogram is titled "Test prep = none" and the right is titled "Test prep = completed". Both histograms plot the "Proportion of ALL students in bin" on the y-axis (ranging from 0.00 to 0.14) against "Math score" on the x-axis (ranging from 0 to 100). Each histogram contains two types of bars: gray bars representing the "All students" distribution and colored bars representing the "none" and "completed" groups. In the "none" group histogram, the colored bars are shifted towards lower and mid scores. In the "completed" group histogram, the colored bars are concentrated at higher scores.

10.4 – Partial distributions as parts of the total

Figure 10.4 – Partial distributions of math scores by test preparation course

These two panels show partial distributions for the two test-prep groups. In each panel, the gray bars represent the overall math-score distribution (all students), and the colored bars show the proportion of all students who are in that score bin and in that specific test-prep group. For the “completed” group, the colored bars are concentrated at higher scores; for the “none” group they are shifted toward lower and mid scores. This follows the chapter’s idea of viewing each group as a part of the total distribution, rather than stacking everything in a single figure.

CHAPTER 11

```
print("Parental education percentages:\n", parent_props)
print("\nTest prep percentages:\n", prep_props)

# -----
# 11.1(a) BAD EXAMPLE - Combined pie that double-counts
# -----
# Build one big list of slices:
# - all parental education categories
# - all test prep categories
labels = list(parent_props.index) + [f"prep: {x}" for x in prep_props.index]
sizes = list(parent_props.values) + list(prep_props.values)

print("\nSum of all slice percentages (SHOULD NOT be >100 for a valid pie):",
      sum(sizes))

plt.figure(figsize=(6, 6))
plt.pie(
    sizes,
    labels=labels,
    autopct="%1.1f%%",
    startangle=90
)
plt.title("BAD: Combined pie for parental education + test prep\n(doubles counts, >100%)")
plt.tight_layout()
plt.show()

# -----
# 11.1(b) BAD / MISLEADING - Mixed bar chart
# -----
```

```
# -----
# 11.1(b) BAD / MISLEADING - Mixed bar chart
# -----
# Put both sets of percentages into one table
bad_bar = (
    pd.concat([
        parent_props.rename("percentage"),
        prep_props.rename("percentage")
    ])
    .reset_index()
)

bad_bar.rename(columns={"index": "category"}, inplace=True)

# Tag which group each category belongs to
bad_bar["group"] = np.where(
    bad_bar["category"].isin(parent_props.index),
    "Parental education",
    "Test prep"
)

print("\nData for mixed bar chart:\n", bad_bar)

plt.figure(figsize=(8, 4))
plt.bar(bad_bar["category"], bad_bar["percentage"])
plt.xticks(rotation=45, ha="right")
plt.xlabel("Percentage of students")
plt.title("BAD: Mixed categories in one bar plot\n(hides overlap between groups)")
plt.tight_layout()
plt.show()
```

Category	Percentage
associate's degree	22.2
bachelor's degree	11.8
high school	19.6
master's degree	5.9
some college	22.6
some high school	17.9

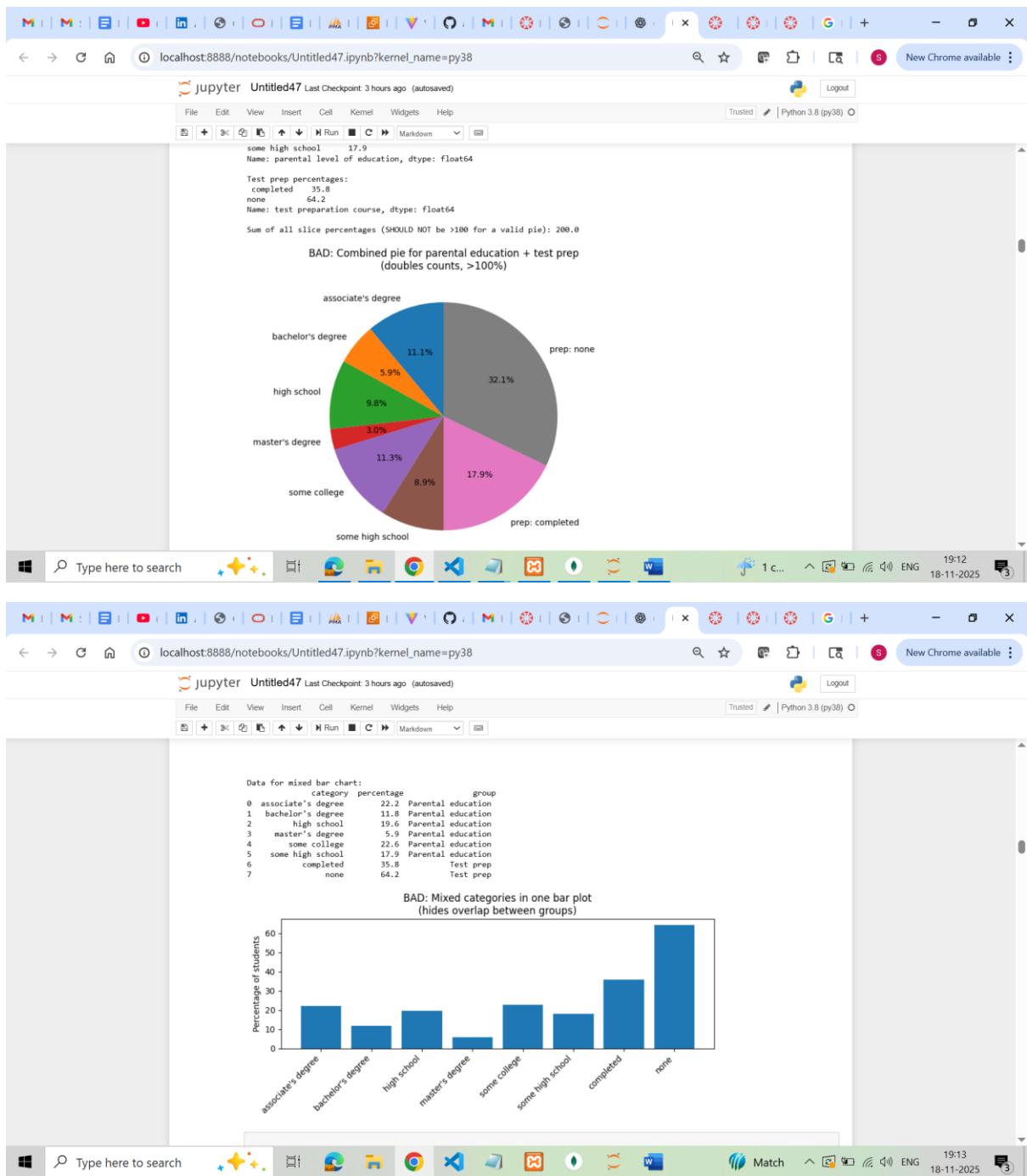


Figure 11.1(a) – BAD combined pie

I combined the percentages of parental education levels and test-preparation completion into a single pie chart. The slices sum to more than 100%, because each student is counted once in a parental-education category and again in a test-prep category. A pie chart implies mutually exclusive categories that add up to one whole, which is violated here, so this is a “nested proportions gone wrong” example.

Figure 11.1(b) – BAD / misleading bar chart

The bar chart shows the same percentages as separate bars. This plot is not mathematically wrong, since bar heights don't need to add to 100%. However, it is conceptually misleading: it suggests many independent categories, even though all test-prep bars are just another way of slicing the same students that already appear in the parental-education bars. The overlap between groups is completely hidden.

```

In [13]: # =====
# CHAPTER 11.2 - Mosaic plots and treemaps
# 11.2(a) Mosaic-style plot
# =====

top_var = "parental level of education"
sub_var = "test preparation course"

# Cross-tabulation: counts for each combination
ct = pd.crosstab(df[top_var], df[sub_var])
print("Crosstab of parental education x test prep:\n")
print(ct)

total = ct.to_numpy().sum()

# Order of categories
top_levels = ct.index.tolist()
sub_levels = ct.columns.tolist()

# Colors for the nested variable
color_cycle = plt.rcParams["axes.prop_cycle"].by_key()["color"]
colors = [lev: color_cycle[i % len(color_cycle)] for i, lev in enumerate(sub_levels)]

fig, ax = plt.subplots(figsize=(8, 4))

x_offset = 0.0
x_center = []
x_labels = []

for top in top_levels:
    # width of strip = proportion of ALL students with this parental education
    parent_total = ct.loc[top].sum()
    if parent_total == 0:
        continue
    width = parent_total / total
    y_offset = 0.0

    for sub in sub_levels:
        count = ct.loc[top, sub]
        x_center.append(x_offset + width / 2)
        x_labels.append(sub)
        ax.bar(x_offset, count, width=width, color=colors[sub])
        y_offset += count
        x_offset += width

    x_offset += width
    x_center.append(x_offset)
    x_labels.append("")

ax.legend(handles, sub_levels, title=sub_var, bbox_to_anchor=(1.05, 1), loc="upper left")

plt.tight_layout()
plt.show()

```

parental level of education	test preparation course	completed	none
associate's degree	completed	82	149
associate's degree	none	46	73
bachelor's degree	completed	56	140
bachelor's degree	none	20	39
high school	completed	77	149
high school	none	77	182
master's degree	completed		
some college	completed		
some high school	completed		

Mosaic plot: parental education x test preparation course

Jupyter Untitled47 Last Checkpoint: 3 hours ago (unsaved changes)

```
In [15]: # 11.2(b) Treemap-style plot: parental.education x test_prep
# (pure Matplotlib, no squarify needed)

top_var = "parental level of education"
sub_var = "test preparation course"

# Crosstab (in case it's not defined in this cell)
ct = pd.crosstab(ufftop_var, offsub_var)

# Flatten to long format: one row per (parent ed, test prep) combo
treemap_df = (
    ct.stack()
    .reset_index(names="count")
)

print("Data used for treemap-style plot:\n", treemap_df)

# ----- simple recursive treemap layout (slice-and-dice) -----
def treemap_layout(items, x=0.0, y=0.0, w=1.0, h=1.0, horizontal=True):
    """
    items: list of (id, size)
    returns: list of dicts {id, x, y, w, h}
    """

    if not items:
        return []
    if len(items) == 1:
        id_, size = items[0]
        return [{"id": id_, "x": x, "y": y, "w": w, "h": h}]

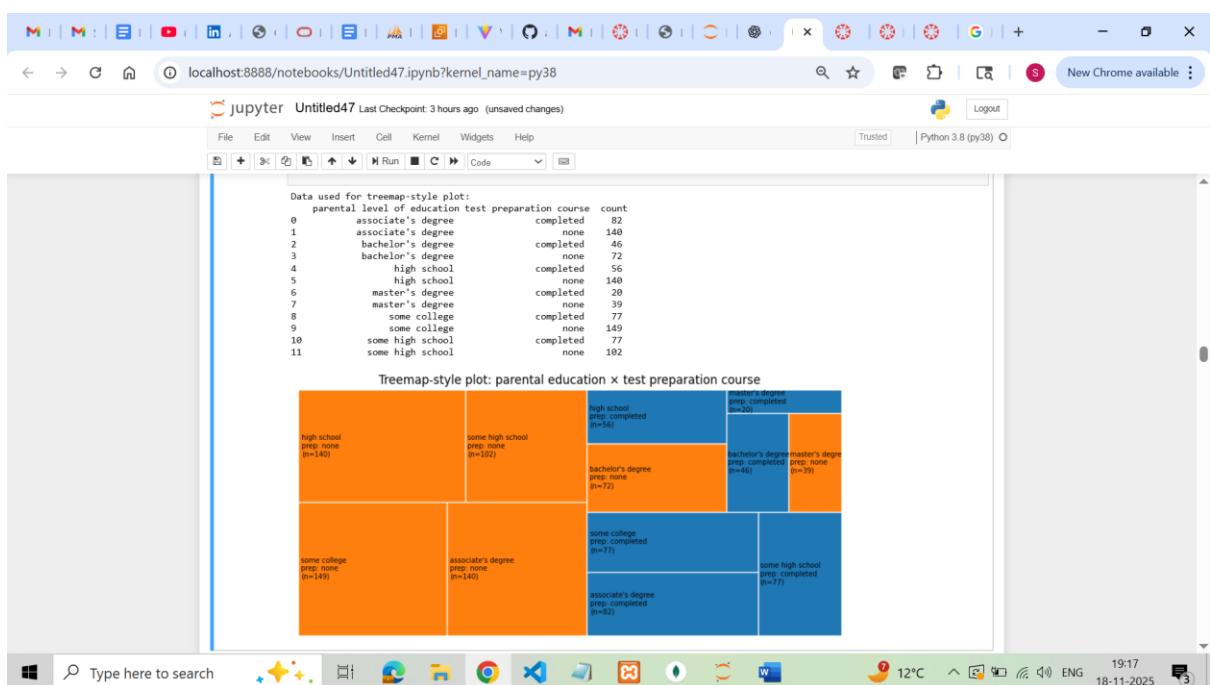
    # sort by size (largest first) for nice layout
    items = sorted(items, key=lambda t: t[1], reverse=True)
    total = sum(size for _, size in items)

    # split into two groups with roughly half the total size each
    acc = 0.0
    split_index = 0
    for i, (id_, s) in enumerate(items):
        if acc >= total / 2 and i > 0:
            break
        acc += s

    # recursive call
    left_items = items[:split_index]
    right_items = items[split_index:]

    left_rects = treemap_layout(left_items, x=x, y=y, w=w/2, h=h, horizontal=horizontal)
    right_rects = treemap_layout(right_items, x=x+w/2, y=y, w=w/2, h=h, horizontal=horizontal)

    return left_rects + right_rects
```



The treemap-style plot divides the full plotting area into rectangles whose area is proportional to the number of students in each combination of parental education level and test-preparation status. All rectangles together represent the entire sample. Colors encode test-preparation status, while text labels show the parent-education category and the group size. Like the mosaic plot, this correctly visualizes nested proportions, but here the rectangles are arranged in a free-form tiling rather than a strict grid.

```

In [16]: # CHAPTER 11.3 - Nested pies
# gender (inner ring) x test preparation course (outer ring)
# =====

top_var = "gender"
sub_var = "test preparation course"

# Drop rows with missing values in these columns
data = df.dropna(subsets=[top_var, sub_var])

# Top-level counts (inner ring)
top_counts = (
    data[top_var]
    .value_counts()
    .sort_index()
)
top_levels = top_counts.index.tolist()
print("Top-level counts (gender):\n", top_counts, "\n")

# Nested counts (outer ring): gender x test prep
nested_counts = (
    data.groupby([top_var, sub_var])
    .size()
    .sort_index()
)
print("Nested counts (gender x test prep):\n", nested_counts, "\n")

# Make sure we have a consistent order for subcategories
sub_levels = (
    data[sub_var]
    .value_counts()
    .sort_index()
    .index
    .tolist()
)

# Build outer ring sizes in order: for each gender, each test-prep Level
outer_sizes = []

```

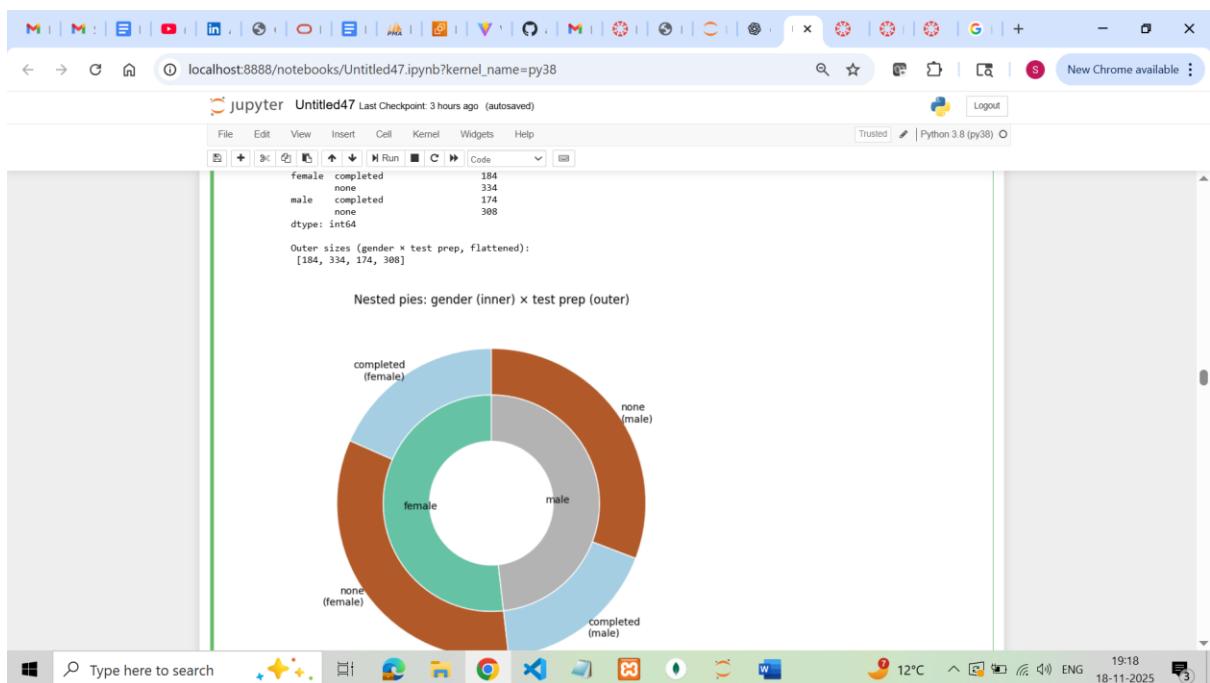


Figure 11.3 – Nested pies for gender and test preparation course.

The inner ring shows the proportion of female versus male students in the sample. The outer ring breaks each gender slice into students who completed the test-preparation course and those who did not. The area and angle of each outer wedge represent the fraction of all students in that gender-test-prep combination. This is a nested-pies visualization of the same nested proportions we showed earlier with mosaic and treemap plots.

localhost:8888/notebooks/Untitled47.ipynb?kernel_name=py38

jupyter Untitled47 Last Checkpoint: 3 hours ago (autosaved)

```
In [17]: # CHAPTER 11.4 - Parallel sets
# parental Level of education + test preparation course
# =====

import numpy as np
import matplotlib.pyplot as plt
from matplotlib.patches import Rectangle, PathPatch
from matplotlib.path import Path

var_left = "parental level of education"
var_right = "test preparation course"

data = df.dropna(subset=[var_left, var_right])

total = len(data)
print("Total students used in parallel sets:", total)

# -----
# 1) Category order and segment sizes on each axis
# -----

# left axis categories (parental education)
left_counts = (
    data[var_left]
    .value_counts()
    .sort_index()
)
left_levels = left_counts.index.tolist()
left_heights = (left_counts / total).tolist()
print("\nLeft axis (parental education) counts:\n", left_counts)

# Right axis categories (test prep)
right_counts = (
    data[var_right]
    .value_counts()
    .sort_index()
)
right_levels = right_counts.index.tolist()
right_heights = (right_counts / total).tolist()
```

localhost:8888/notebooks/Untitled47.ipynb?kernel_name=py38

jupyter Untitled47 Last Checkpoint: 3 hours ago (autosaved)

```
# -----
# 3) Plot parallel sets
# -----

fig, ax = plt.subplots(figsize=(8, 4))

# Positions of the two axes
x_left = 0.0
x_right = 1.0
bar_width = 0.05

# ---- Draw category bars on left axis ----
for lev in left_levels:
    y0, y1 = left_segments[lev]
    rect = Rectangle(
        (x_left - bar_width / 2, y0),
        bar_width,
        y1 - y0,
        facecolor="lightgray",
        edgecolor="black"
    )
    ax.add_patch(rect)
    ax.text(
        x_left - 0.07,
        (y0 + y1) / 2,
        lev,
        va="center",
        ha="right",
        fontsize=8
    )

# ---- Draw category bars on right axis ----
for lev in right_levels:
    y0, y1 = right_segments[lev]
    rect = Rectangle(
        (x_right + bar_width / 2, y0),
        bar_width,
        y1 - y0,
        facecolor="lightgray",
        edgecolor="black"
    )
    ax.add_patch(rect)
```

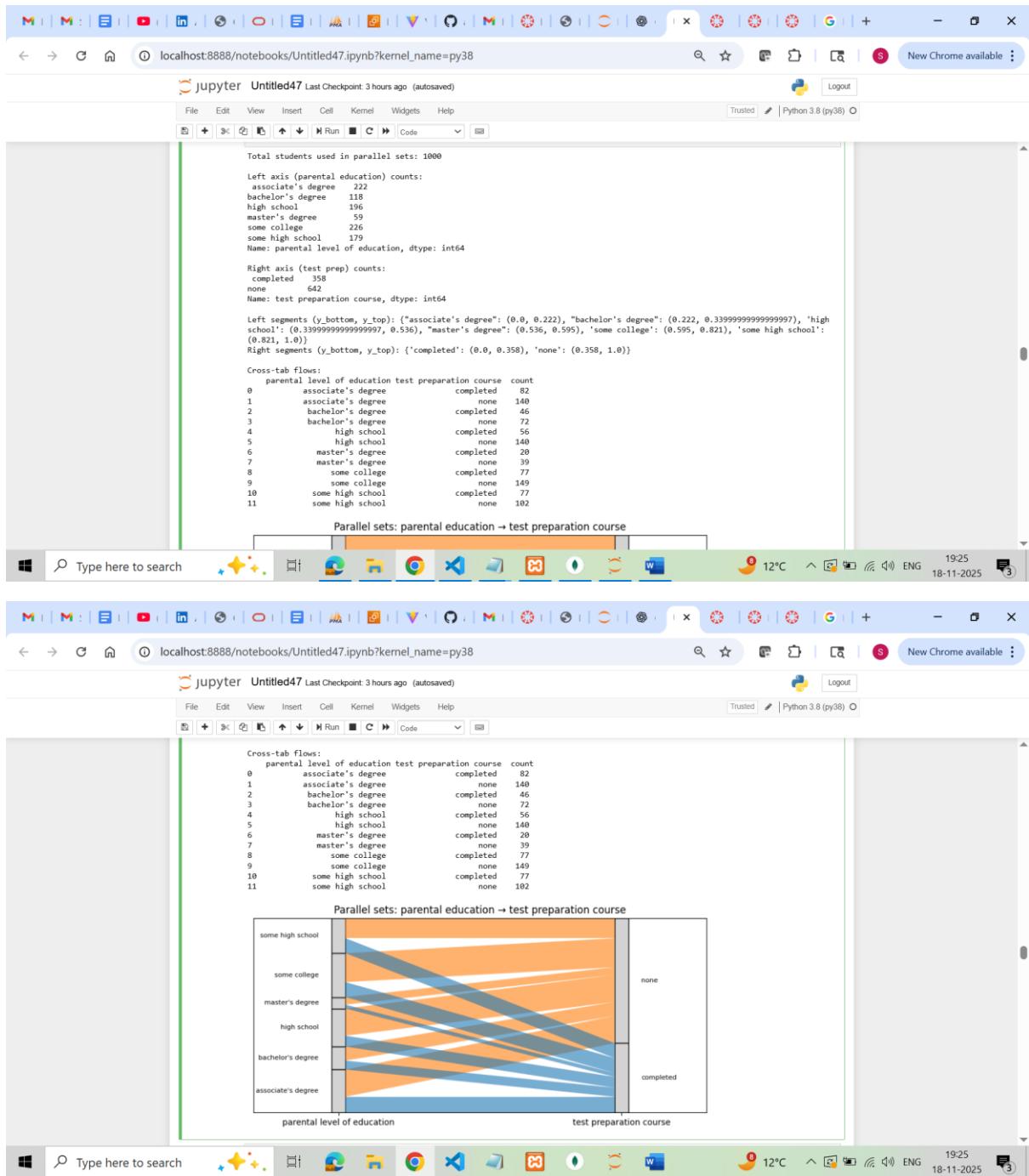


Figure 11.4 – Parallel sets for parental education and test preparation course.

The left axis shows how all students are distributed across parental education levels, and the right axis shows how they are distributed across test-preparation status. The shaded bands connect categories on the left to categories on the right. The thickness of each band is proportional to the number of students in that combination (e.g., “some college” + “completed test prep”). Coloring the bands by test-preparation status highlights where most of the “completed” students come from. This is a simple two-axis parallel sets plot; the same idea can be extended to three or more categorical variables, as in the textbook example with bridge material, length, era, and river.

localhost:8888/notebooks/Untitled47.ipynb?kernel_name=py38

jupyter Untitled47 Last Checkpoint: 3 hours ago (autosaved)

In [18]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

plt.rcParams["figure.figsize"] = (6, 4)

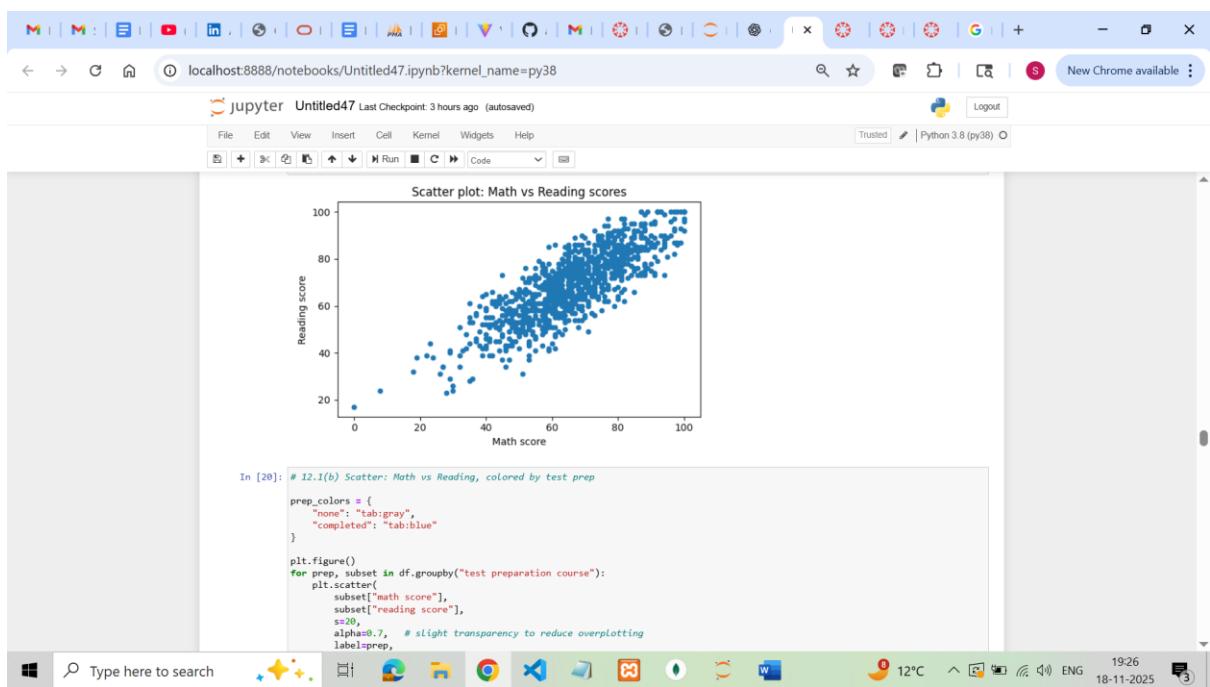
# df already loaded from URL earlier
# df = pd.read_csv(url)
```

In [19]:

```
# *****
# CHAPTER 12 - Scatter plots
# 12.1(a) Basic scatter: math vs reading
# *****

plt.figure()
plt.scatter(
    df["math score"],
    df["reading score"],
    s=20 # marker size
)
plt.xlabel("Math score")
plt.ylabel("Reading score")
plt.title("Scatter plot: Math vs Reading scores")
plt.tight_layout()
plt.show()
```

Scatter plot: Math vs Reading scores



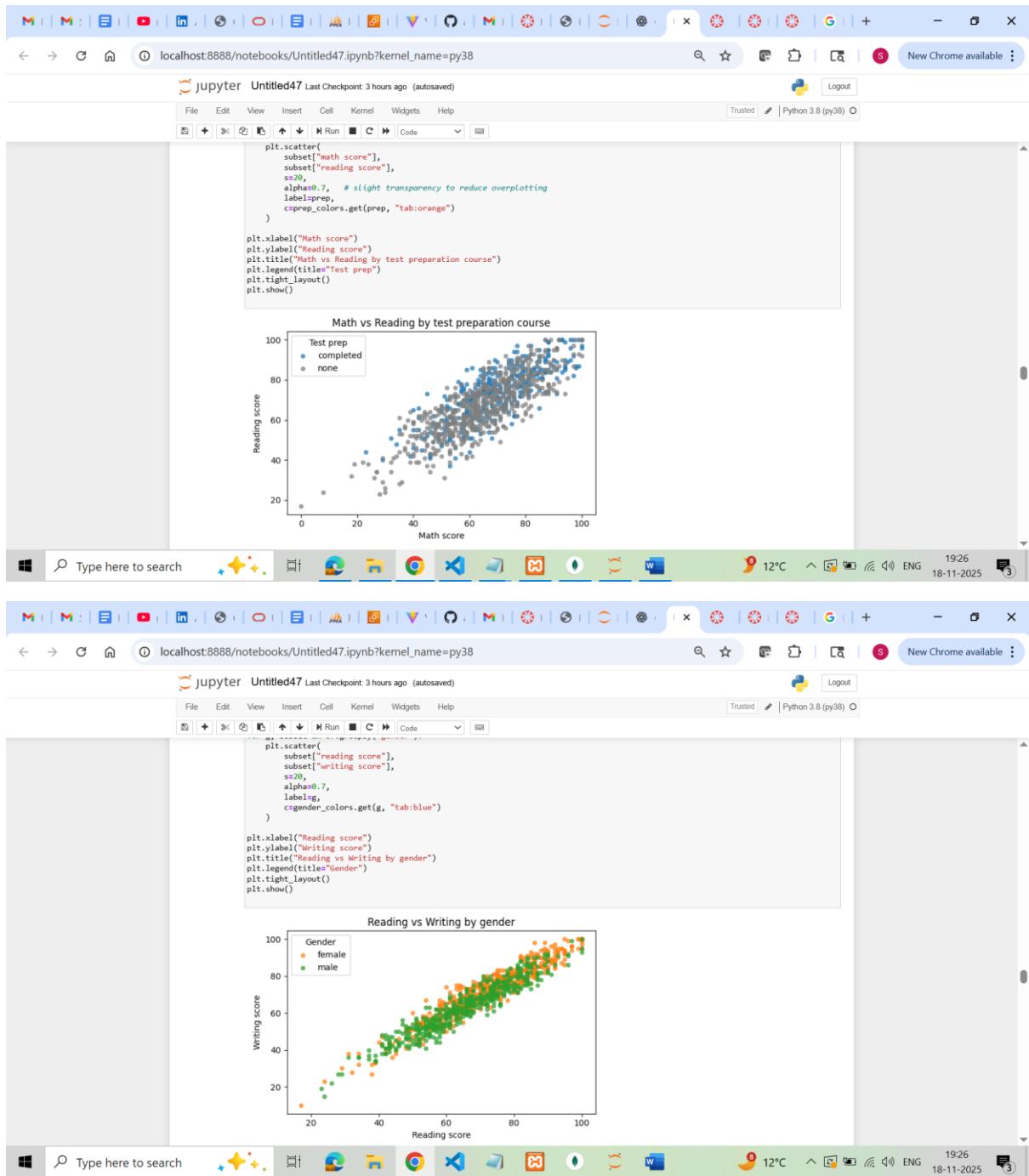


Figure 12.1(a) shows a basic scatter plot of math score versus reading score.

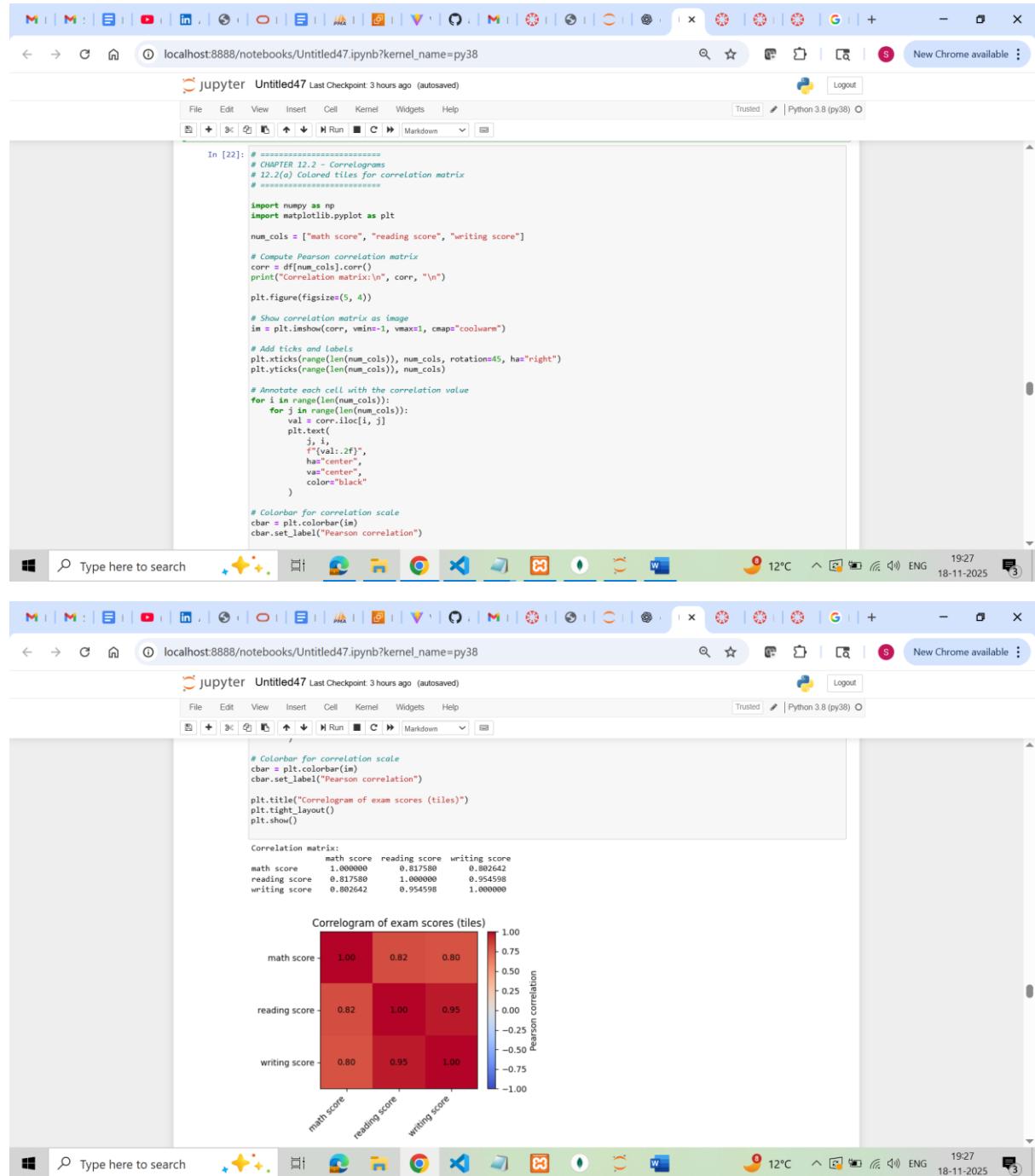
The points form a clear upward cloud, indicating a strong positive association: students with higher math scores also tend to have higher reading scores.

Figure 12.1(b) colors the same scatter plot by test-preparation course.

Students who completed the test-preparation course (blue) tend to cluster in the upper-right region (high math and reading scores), while students with no prep (gray) are more spread out and more common at lower scores. This shows how scatter plots can incorporate a third variable via color to reveal group differences.

Figure 12.1(c) shows reading score versus writing score, colored by gender.

The relationship between reading and writing scores is almost perfectly linear and positive, with both genders following the same general pattern. Any gender differences are subtle compared to the strong overall association between the two quantitative variables.



The screenshot shows two consecutive screenshots of a Jupyter Notebook interface running in a Chrome browser. The notebook title is "Untitled47".

Screenshot 1: The code cell contains Python code to generate a correlation matrix and a heatmap. The code imports numpy and matplotlib.pyplot, defines three columns ("math score", "reading score", "writing score"), computes the Pearson correlation matrix, and then creates a heatmap where each cell's value is annotated with its numerical value. A colorbar on the right indicates the scale from -1.00 (blue) to 1.00 (red).

```
In [22]: # CHAPTER 12.2 - Correlograms
# 12.2(a) Colored tiles for correlation matrix
# *****

import numpy as np
import matplotlib.pyplot as plt

num_cols = ["math score", "reading score", "writing score"]

# Compute Pearson correlation matrix
corr = df[num_cols].corr()
print("Correlation matrix:\n", corr, "\n")

plt.figure(figsize=(5, 4))

# Show correlation matrix as image
im = plt.imshow(corr, vmin=-1, vmax=1, cmap="coolwarm")

# Add ticks and labels
plt.xticks(range(len(num_cols)), num_cols, rotation=45, ha="right")
plt.yticks(range(len(num_cols)), num_cols)

# Annotate each cell with the correlation value
for i in range(len(num_cols)):
    for j in range(len(num_cols)):
        val = corr.iat[i, j]
        plt.text(j, i,
                str(val).ljust(2),
                ha="center",
                va="center",
                color="black"
                )

# Colorbar for correlation scale
cbar = plt.colorbar(im)
cbar.set_label("Pearson correlation")
```

Screenshot 2: The code cell has been executed, and the output shows the correlation matrix and the heatmap. The heatmap is titled "Correlogram of exam scores (tiles)". The correlation matrix table is as follows:

	math score	reading score	writing score
math score	1.00	0.82	0.80
reading score	0.82	1.00	0.95
writing score	0.80	0.95	1.00

The heatmap shows a strong positive correlation between all three scores, with values ranging from approximately 0.80 to 1.00.

```
# Colorbar for correlation scale
cbar = plt.colorbar(im)
cbar.set_label("Pearson correlation")

plt.title("Correlogram of exam scores (tiles)")
plt.tight_layout()
plt.show()
```

Correlation matrix:

	math score	reading score	writing score
math score	1.000000	0.817580	0.802642
reading score	0.817580	1.000000	0.954598
writing score	0.802642	0.954598	1.000000

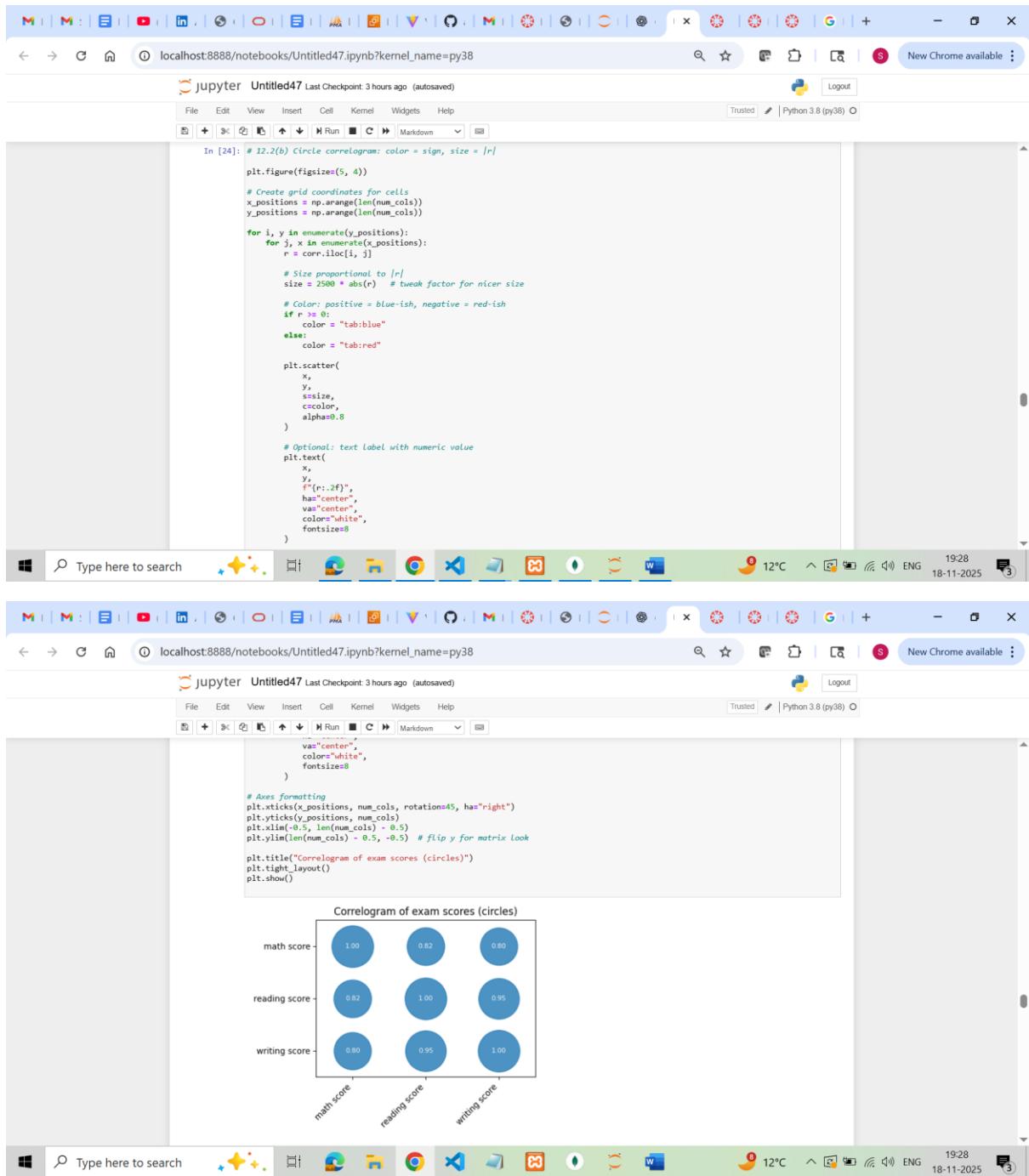


Figure 12.2(a) – Correlogram with colored tiles.

I computed the Pearson correlation matrix for math, reading, and writing scores and displayed it as a tile plot. All off-diagonal correlations are strongly positive (around 0.8–0.95), meaning that students who do well on one exam tend to do well on the others. The color scale from blue to red encodes correlation strength and sign, similar to the correlogram example in the textbook.

Figure 12.2(b) – Circle correlogram.

This plot shows the same correlation matrix, but each cell is represented by a circle whose color indicates the sign of the correlation and whose size scales with the absolute value of the correlation coefficient. Large blue circles highlight strong positive relationships (e.g., reading vs. writing), while

any correlations near zero would appear as very small circles, visually deemphasizing weak associations as recommended in the chapter

In [25]:

```
# CHAPTER 12.3 - Dimension reduction (PCA)
# Using math, reading, writing scores
#
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
import numpy as np
import matplotlib.pyplot as plt

num_cols = ["math score", "reading score", "writing score"]

# Drop rows with missing values in these columns
X = df[num_cols].dropna().values

# 1) Standardize the variables (mean=0, std=1)
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# 2) Fit PCA with up to 3 components (same as number of variables)
pca = PCA(n_components=3)
X_pca = pca.fit_transform(X_scaled)

print("Explained variance ratio by component:", pca.explained_variance_ratio_)
print("Cumulative explained variance:", np.cumsum(pca.explained_variance_ratio_))

# 12.3(a) Screen-like bar plot: how much variance each PC explains
#
plt.figure(figsize=(5, 4))
pcs = np.arange(1, 4)

plt.bar(pcs, pca.explained_variance_ratio_)
plt.plot(
    pcs,
    np.cumsum(pca.explained_variance_ratio_),
    markers="o"
)
```

```
plt.show()

# 12.3(b) Scatter plot of first two PCs, colored by test prep
#
pc1 = X_pca[:, 0]
pc2 = X_pca[:, 1]

prep_colors = {
    "none": "tab:gray",
    "completed": "tab:blue"
}

plt.figure(figsize=(6, 4))

# Need matching subset of df without missing exam scores
df_pca = df[num_cols + ["test preparation course"]].dropna()

for prep, subset in df_pca.groupby("test preparation course"):
    # indices of this group in the PCA output (same order as dropna)
    idx = subset.index
    plt.scatter(
        pc1[idx], df_pca.index[0], # align indices if they don't start at 0
        pc2[idx] - df_pca.index[0],
        s=20,
        alpha=0.7,
        label=prep,
        c=prep_colors.get(prep, "tab:orange")
    )

plt.xlabel("PC1")
plt.ylabel("PC2")
plt.title("PCA of exam scores (PC1 vs PC2)\ncolored by test preparation course")
plt.legend(title="Test prep")
plt.tight_layout()
plt.show()
```

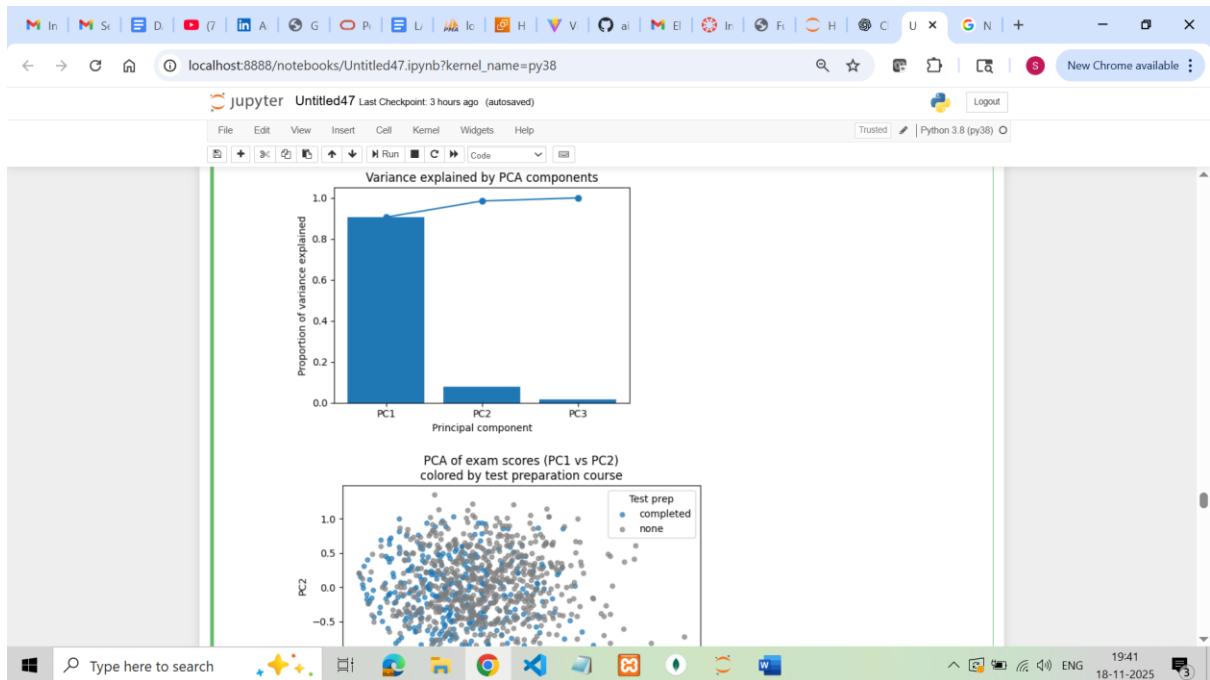
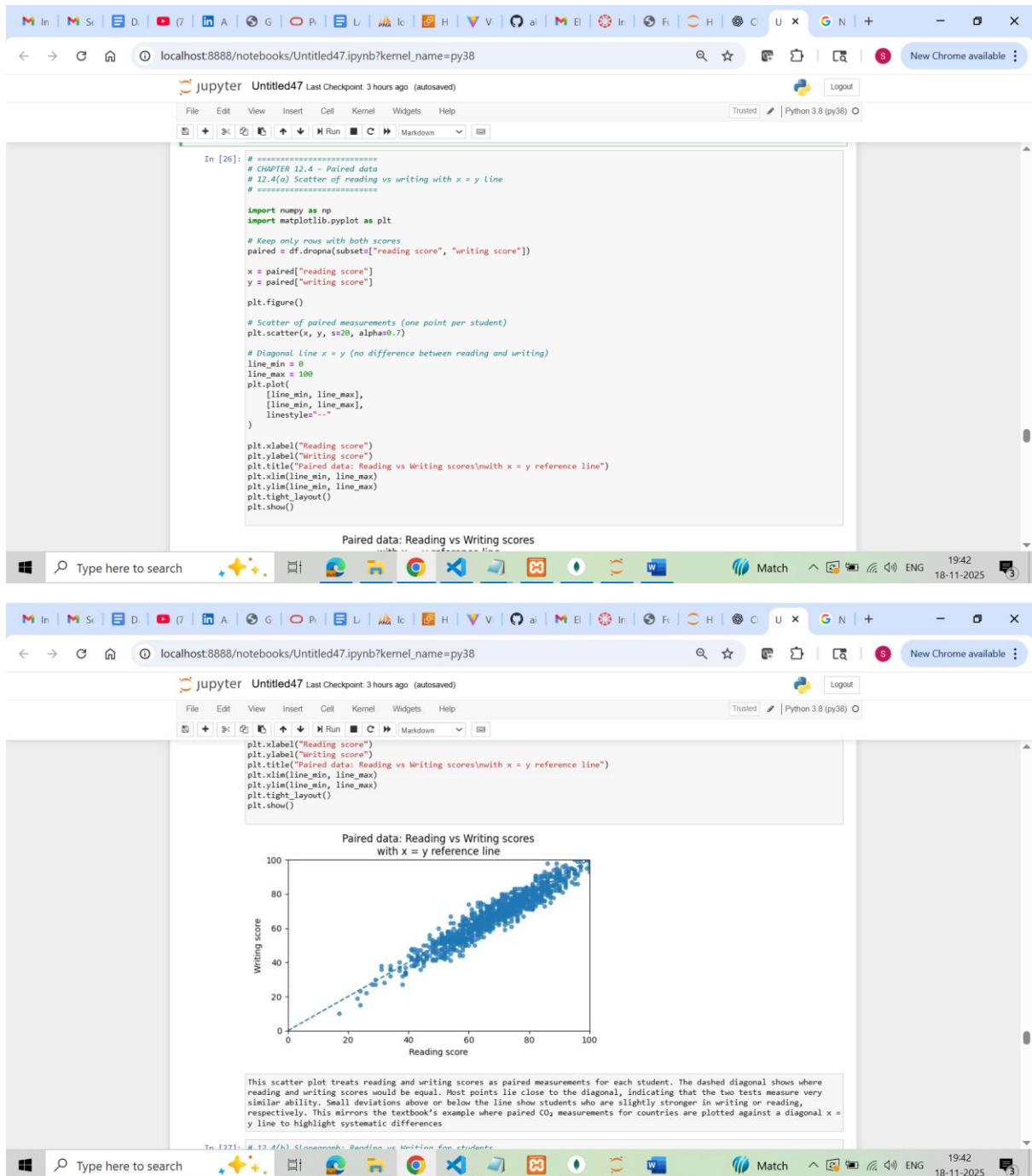


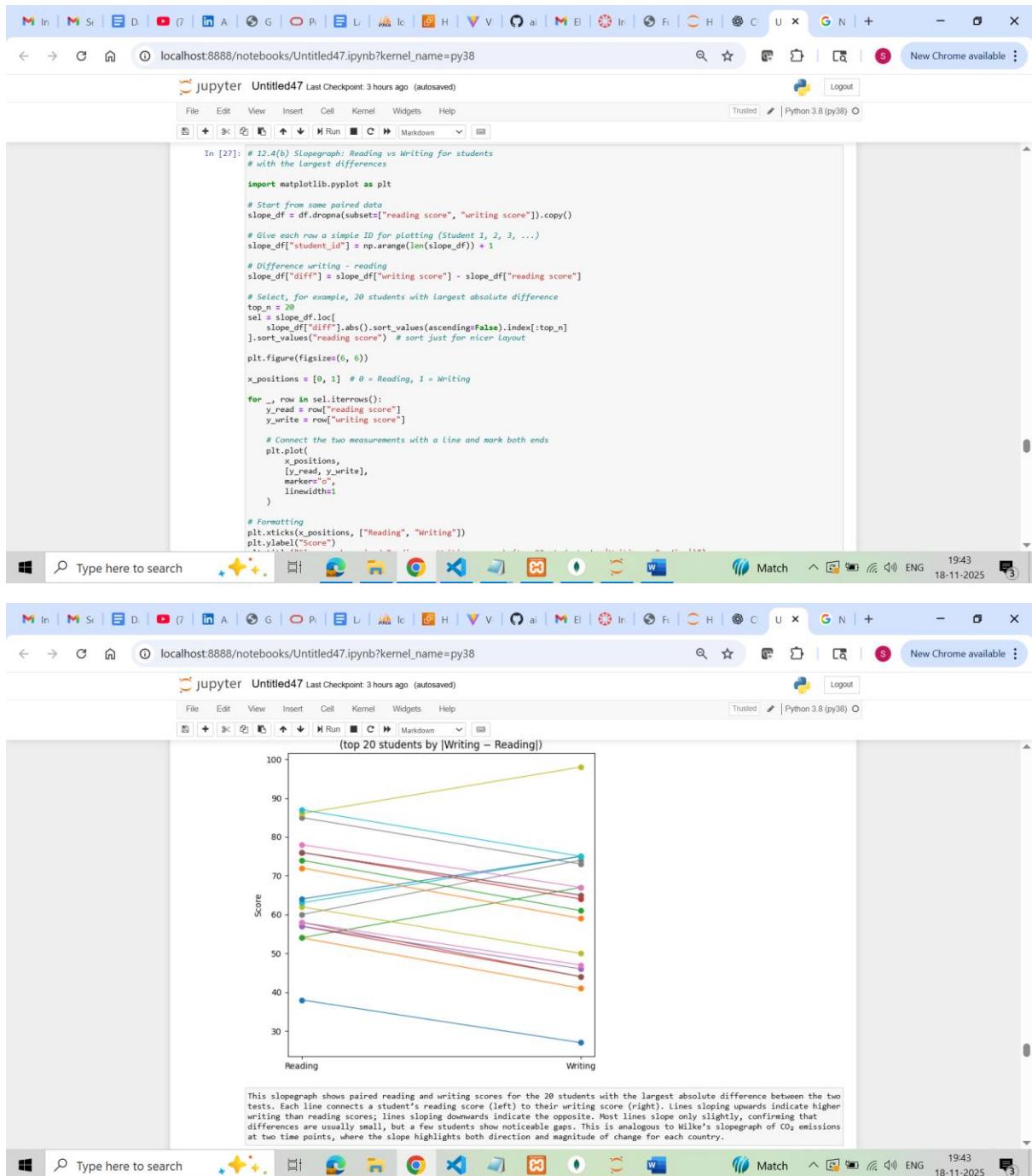
Figure 12.3(a) – Variance explained by PCA.

I applied principal component analysis (PCA) to the three standardized exam scores (math, reading, writing). The bar/line plot shows how much variance each principal component explains. The first component (PC1) explains most of the total variance, while PC2 and PC3 contribute much less. This indicates that a single underlying dimension—roughly “overall exam performance”—captures most of the variation in these three scores.

Figure 12.3(b) – PC1 vs PC2 scatter plot colored by test-preparation course.

The scatter plot shows the first two principal components. Each point is a student; color indicates whether they completed the test-preparation course. PC1 separates students largely by overall performance: higher PC1 values correspond to higher scores on all three exams. Students who completed test prep tend to lie toward the higher-PC1 region, reinforcing what we saw in the original scatter plots. PC2 mostly captures smaller differences between the individual exams. PCA thus reduces the three correlated exam scores to two new variables that summarize performance while retaining most of the structure in the data.





CHAPTER 13

M In M Sk D YouTube LinkedIn A G P L H V ai M B Ir F C U X localhost:8888/notebooks/Untitled47.ipynb?kernel_name=py38

jupyter Untitled47 Last Checkpoint: 3 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3.8 (py38)

CHAPTER 13

```
In [28]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

plt.rcParams["figure.figsize"] = (8, 4)

# Daily minimum temperatures in Melbourne (1981-1990)
url_temp = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/daily-min-temperatures.csv"
temp_df = pd.read_csv(url_temp)

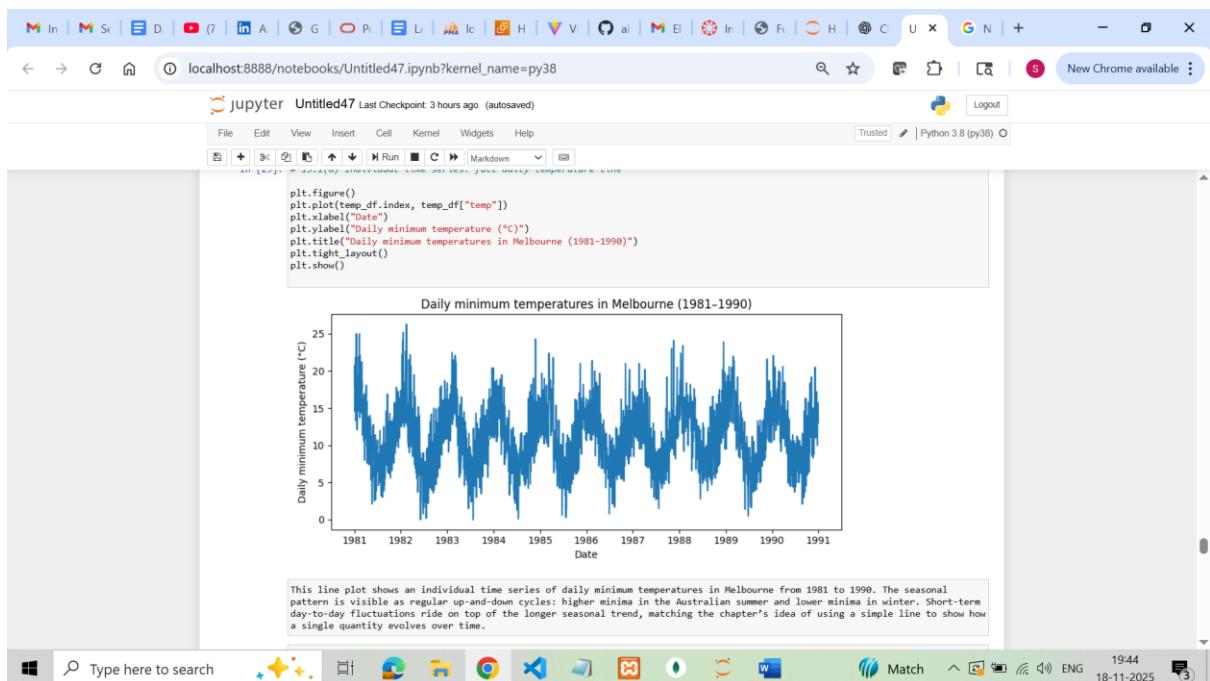
# Parse dates
temp_df["Date"] = pd.to_datetime(temp_df["Date"])
temp_df = temp_df.rename(columns={"Temp": "temp"})
temp_df = temp_df.set_index("Date").sort_index()

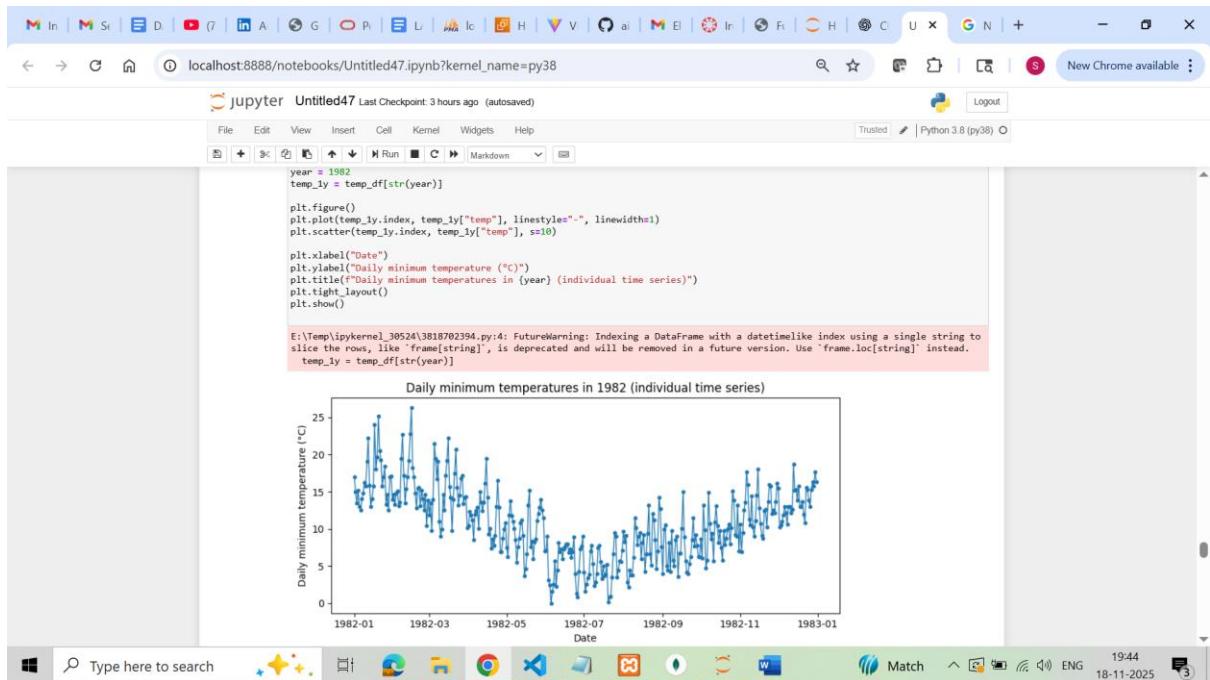
print(temp_df.head())
print(temp_df.index.min(), "to", temp_df.index.max())

      temp
Date
1981-01-01  20.7
1981-01-02  17.9
1981-01-03  18.8
1981-01-04  14.6
1981-01-05  15.8
1981-01-01 00:00:00 + 1990-12-31 00:00:00
```

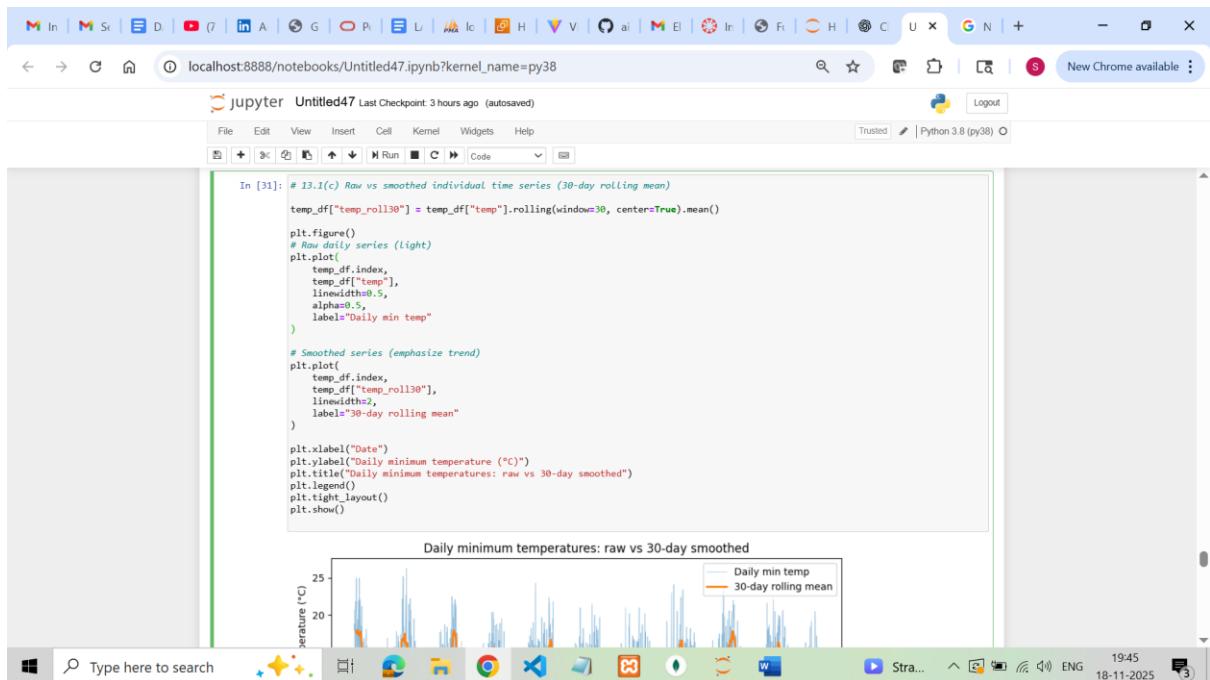
```
In [29]: # 13.1(a) Individual time series: full daily temperature line

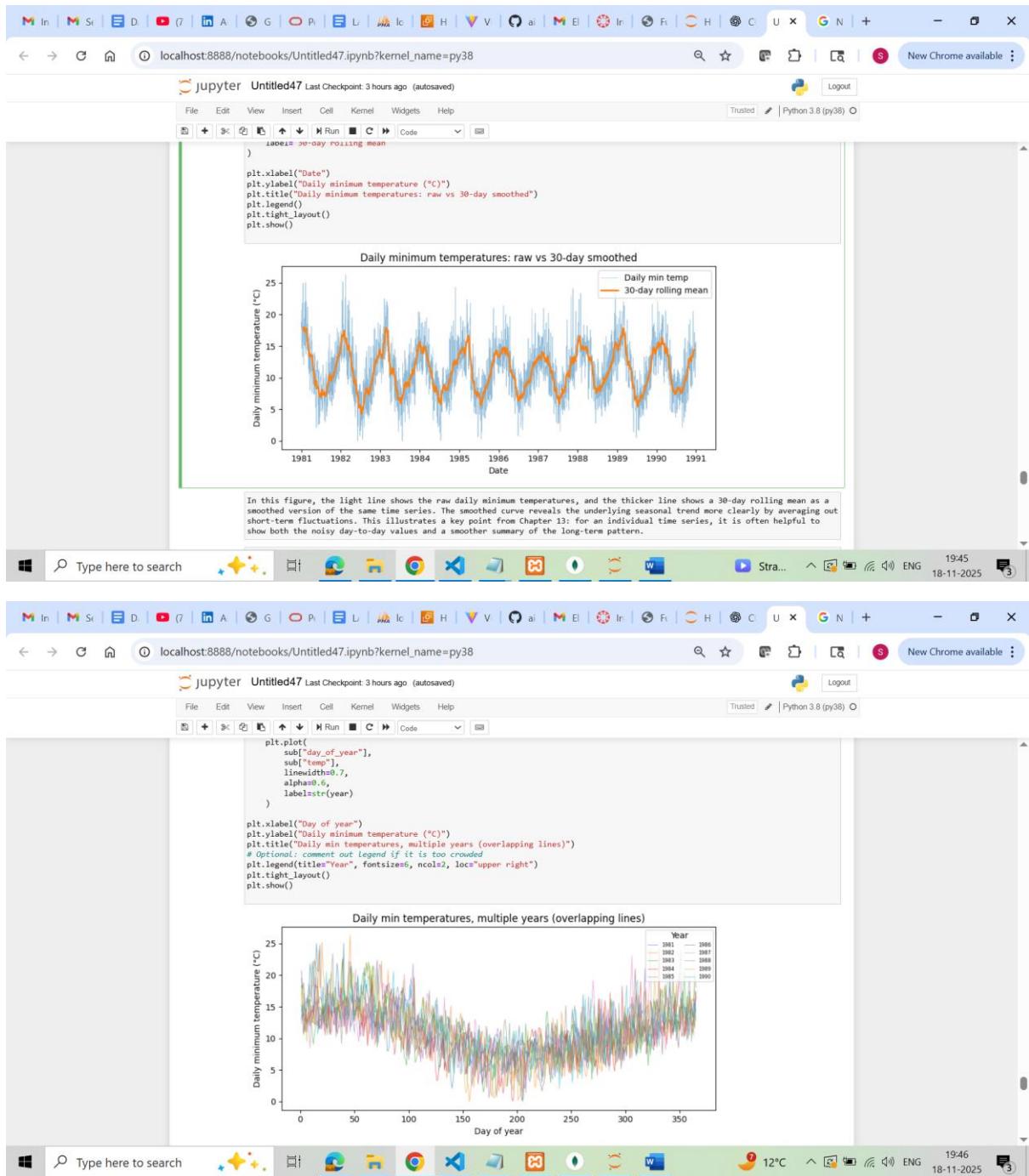
plt.figure()
plt.plot(temp_df.index, temp_df["temp"])
plt.xlabel("Date")
plt.ylabel("Daily minimum temperature (°C)")
plt.title("Daily minimum temperatures in Melbourne (1981-1990)")
plt.tight_layout()
```





Focusing on 1982 alone, I plotted the daily minimum temperatures as both a line and individual points. The line reveals the smooth seasonal rise and fall across the year, while the points emphasize that each dot is a separate daily observation. This matches the “individual time series” examples in the textbook, where both the trajectory and the discreteness of the data are important.





This figure plots a separate time series for each year (1981–1990) on the same axes. While it technically shows multiple time series, the result is a “spaghetti plot” where lines overlap heavily. It is hard to follow any single year, and comparing years is difficult. This is a typical “bad” multiple time series design mentioned in the book, where adding many lines to one panel makes the visualization cluttered instead of informative.

In [35]: # 13.2(b) BETTER: small multiples (one time series per panel, same scales)

```

unique_years = sorted(temp_df["year"].unique())
n_years = len(unique_years)

# We have 10 years; 2 rows x 5 columns works nicely
n_rows, n_cols = 2, 5

fig, axes = plt.subplots(n_rows, n_cols, figsize=(14, 5), sharex=True, sharey=True)

for ax, year in zip(axes.flatten(), unique_years):
    sub = temp_df[temp_df["year"] == year]
    ax.plot(
        sub["day_of_year"],
        sub["temp"],
        linewidth=0.8
    )
    ax.set_title(str(year), fontsize=9)
    ax.set_xticks([1, 91, 182, 274, 365])
    ax.set_xticklabels(["Jan", "Apr", "Jul", "Oct", "Dec"])

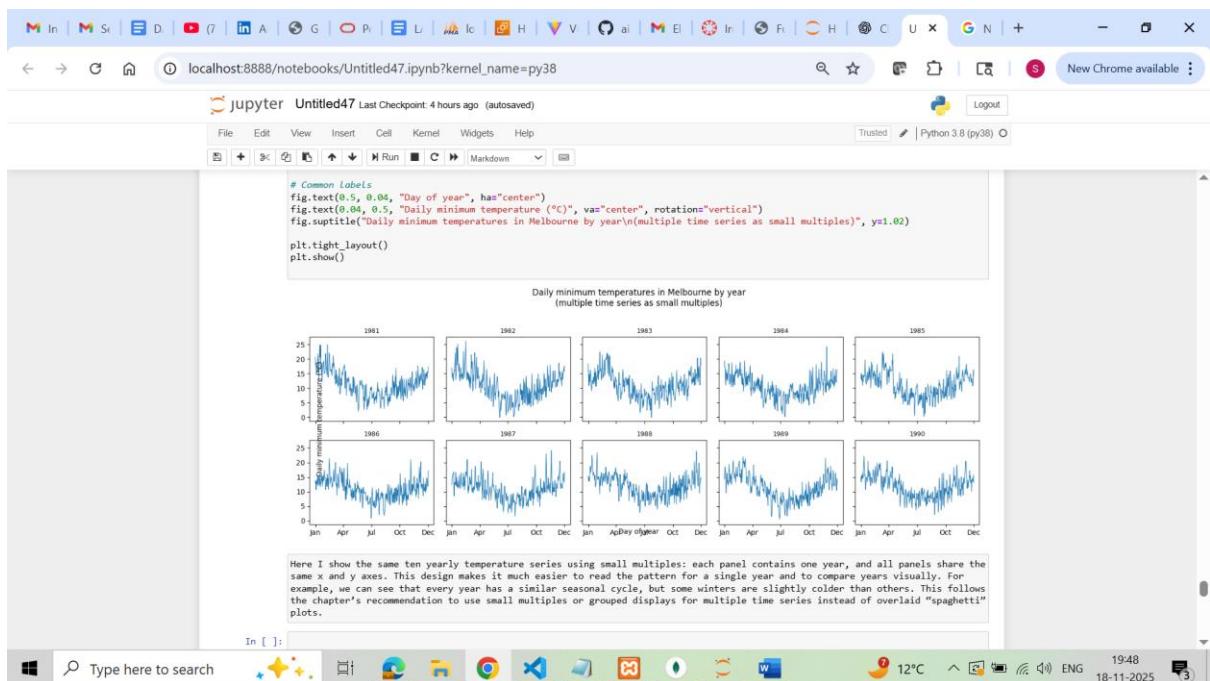
# Hide any unused axes (if any)
for ax in axes.flatten()[n_years:]:
    ax.set_visible(False)

# Common Labels
fig.text(0.5, 0.04, "Day of year", ha="center")
fig.text(0.04, 0.5, "Daily minimum temperature (""C)"", va="center", rotation="vertical")
fig.suptitle("Daily minimum temperatures in Melbourne by year\n(multiple time series as small multiples)", y=1.02)

plt.tight_layout()
plt.show()

```

Daily minimum temperatures in Melbourne by year
(multiple time series as small multiples)



localhost:8888/notebooks/Untitled47.ipynb?kernel_name=py38

jupyter Untitled47 Last Checkpoint: 4 hours ago (autosaved)

```
In [42]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from matplotlib.collections import LineCollection

plt.rcParams["figure.figsize"] = (8, 4)

# Daily minimum temperatures (same dataset as 13.1)
url_temp = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/daily-min-temperatures.csv"
temp_df = pd.read_csv(url_temp)

temp_df["Date"] = pd.to_datetime(temp_df["Date"])
temp_df = (
    temp_df.rename(columns={"Temp": "temp"})
    .set_index("Date")
    .sort_index()
)

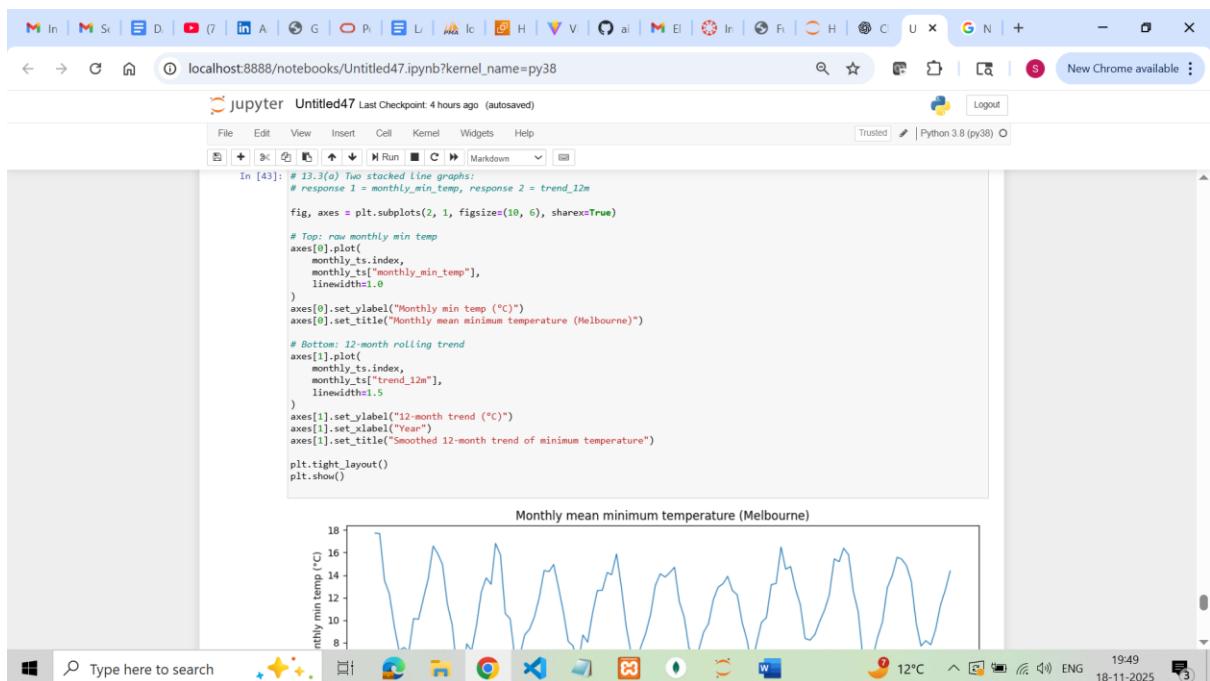
print("temp_df columns:", temp_df.columns)

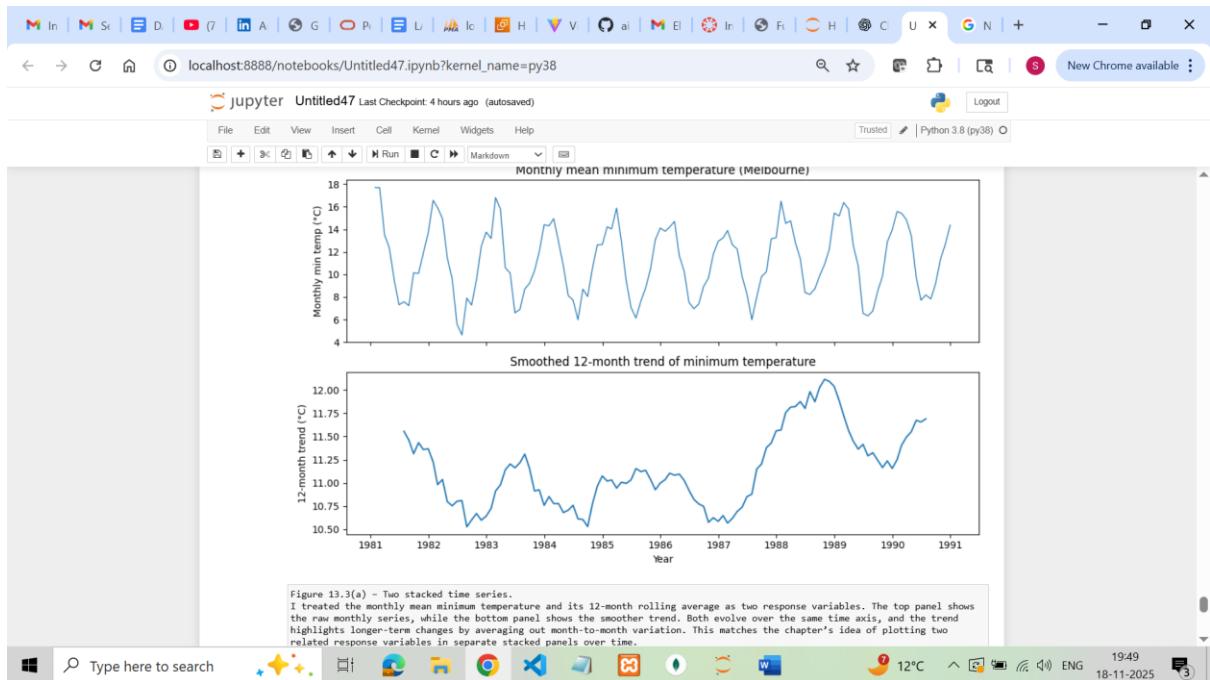
# Monthly mean of daily minimum temps
monthly_ts = temp_df["temp"].resample("M").mean().to_frame("monthly_min_temp")

# Second response variable: 12-month rolling mean (trend)
monthly_ts["trend_12m"] = monthly_ts["monthly_min_temp"].rolling(12, center=True).mean()

print("monthly_ts columns:", monthly_ts.columns)
print(monthly_ts.head())
```

Date	monthly_min_temp	trend_12m
1981-01-31	17.712093	NaN
1981-02-28	17.678571	NaN
1981-03-31	13.500000	NaN
1981-04-30	12.356667	NaN
1981-05-31	9.490323	NaN





```
In [44]: # 13.3(b) Connected scatter: monthly_min_temp vs trend_12m

# Drop NAs from the rolling trend at start/end
plot_df = monthly_ts.dropna(subset=["trend_12m"]).copy()

x = plot_df["monthly_min_temp"].values
y = plot_df["trend_12m"].values
n = len(plot_df)

# Build line segments between consecutive months
points = np.column_stack([x, y])
segments = np.stack([points[:-1], points[1:]], axis=1)

# Color by time index (0 = earliest, n-1 = latest)
t = np.arange(n)
norm = plt.Normalize(t.min(), t.max())
fig, ax = plt.subplots(figsize=(6, 6))

# Line segments with color gradient
lc = LineCollection(segments, cmap="viridis", norm=norm)
lc.set_array(t[1:-1])
lc.set_linewidth(2)
ax.add_collection(lc)

# Scatter points
sc = ax.scatter(
    x,
    y,
    c=t,
    cmap="viridis",
    s=20,
    edgecolor="none"
)

# Colorbar for time
cbar = plt.colorbar(sc, ax=ax)
cbar.set_label("Time (earlier = later)")


12°C 19:49 18-11-2025
```

