

Objectives {.unnumbered}

1. Load and analyze the **Lightcast dataset** in **Spark DataFrame**.
2. Create **five easy and three medium-complexity visualizations** using **Plotly**.
3. Explore **salary distributions, employment trends, and job postings**.
4. Analyze **skills in relation to NAICS/SOC/ONET codes and salaries**.
5. Customize **colors, fonts, and styles** in all visualizations (**default themes result in a 2.5-point deduction**).
6. Follow **best practices for reporting on data communication**.

Step 1: Load the Dataset {unnumbered}

```
In [7]: import pandas as pd
import plotly.express as px
import plotly.io as pio

# Set Plotly renderer for Google Colab
pio.renderers.default = "colab"

from pyspark.sql import SparkSession
from pyspark.sql.functions import col, count, expr, percentile_approx

# Initialize Spark Session
spark = SparkSession.builder.appName("LightcastData").getOrCreate()

# Load Data
df = spark.read.option("header", "true").option("inferSchema", "true").option("multiLine", "true").option("escape")

# Show Schema and Sample Data
df.printSchema()
df.show(5)
```

```

root
|-- ID: string (nullable = true)
|-- LAST_UPDATED_DATE: string (nullable = true)
|-- LAST_UPDATED_TIMESTAMP: timestamp (nullable = true)
|-- DUPLICATES: integer (nullable = true)
|-- POSTED: string (nullable = true)
|-- EXPIRED: string (nullable = true)
|-- DURATION: integer (nullable = true)
|-- SOURCE_TYPES: string (nullable = true)
|-- SOURCES: string (nullable = true)
|-- URL: string (nullable = true)
|-- ACTIVE_URLS: string (nullable = true)
|-- ACTIVE_SOURCES_INFO: string (nullable = true)
|-- TITLE_RAW: string (nullable = true)
|-- BODY: string (nullable = true)
|-- MODELED_EXPIRED: string (nullable = true)
|-- MODELED_DURATION: integer (nullable = true)
|-- COMPANY: integer (nullable = true)
|-- COMPANY_NAME: string (nullable = true)
|-- COMPANY_RAW: string (nullable = true)
|-- COMPANY_IS_STAFFING: boolean (nullable = true)
|-- EDUCATION_LEVELS: string (nullable = true)
|-- EDUCATION_LEVELS_NAME: string (nullable = true)
|-- MIN_EDULEVELS: integer (nullable = true)
|-- MIN_EDULEVELS_NAME: string (nullable = true)
|-- MAX_EDULEVELS: integer (nullable = true)
|-- MAX_EDULEVELS_NAME: string (nullable = true)
|-- EMPLOYMENT_TYPE: integer (nullable = true)
|-- EMPLOYMENT_TYPE_NAME: string (nullable = true)
|-- MIN_YEARS_EXPERIENCE: integer (nullable = true)
|-- MAX_YEARS_EXPERIENCE: integer (nullable = true)
|-- IS_INTERNSHIP: boolean (nullable = true)
|-- SALARY: integer (nullable = true)
|-- REMOTE_TYPE: integer (nullable = true)
|-- REMOTE_TYPE_NAME: string (nullable = true)
|-- ORIGINAL_PAY_PERIOD: string (nullable = true)
|-- SALARY_TO: integer (nullable = true)
|-- SALARY_FROM: integer (nullable = true)
|-- LOCATION: string (nullable = true)
|-- CITY: string (nullable = true)
|-- CITY_NAME: string (nullable = true)
|-- COUNTY: integer (nullable = true)
|-- COUNTY_NAME: string (nullable = true)
|-- MSA: integer (nullable = true)
|-- MSA_NAME: string (nullable = true)
|-- STATE: integer (nullable = true)
|-- STATE_NAME: string (nullable = true)
|-- COUNTY_OUTGOING: integer (nullable = true)
|-- COUNTY_NAME_OUTGOING: string (nullable = true)
|-- COUNTY_INCOMING: integer (nullable = true)
|-- COUNTY_NAME_INCOMING: string (nullable = true)
|-- MSA_OUTGOING: integer (nullable = true)
|-- MSA_NAME_OUTGOING: string (nullable = true)
|-- MSA_INCOMING: integer (nullable = true)
|-- MSA_NAME_INCOMING: string (nullable = true)
|-- NAICS2: integer (nullable = true)
|-- NAICS2_NAME: string (nullable = true)
|-- NAICS3: integer (nullable = true)
|-- NAICS3_NAME: string (nullable = true)
|-- NAICS4: integer (nullable = true)
|-- NAICS4_NAME: string (nullable = true)
|-- NAICS5: integer (nullable = true)
|-- NAICS5_NAME: string (nullable = true)
|-- NAICS6: integer (nullable = true)
|-- NAICS6_NAME: string (nullable = true)
|-- TITLE: string (nullable = true)
|-- TITLE_NAME: string (nullable = true)
|-- TITLE_CLEAN: string (nullable = true)
|-- SKILLS: string (nullable = true)
|-- SKILLS_NAME: string (nullable = true)
|-- SPECIALIZED_SKILLS: string (nullable = true)
|-- SPECIALIZED_SKILLS_NAME: string (nullable = true)
|-- CERTIFICATIONS: string (nullable = true)
|-- CERTIFICATIONS_NAME: string (nullable = true)
|-- COMMON_SKILLS: string (nullable = true)
|-- COMMON_SKILLS_NAME: string (nullable = true)
|-- SOFTWARE_SKILLS: string (nullable = true)
|-- SOFTWARE_SKILLS_NAME: string (nullable = true)
|-- ONET: string (nullable = true)

```

```

--- ONET_NAME: string (nullable = true)
--- ONET_2019: string (nullable = true)
--- ONET_2019_NAME: string (nullable = true)
--- CIP6: string (nullable = true)
--- CIP6_NAME: string (nullable = true)
--- CIP4: string (nullable = true)
--- CIP4_NAME: string (nullable = true)
--- CIP2: string (nullable = true)
--- CIP2_NAME: string (nullable = true)
--- SOC_2021_2: string (nullable = true)
--- SOC_2021_2_NAME: string (nullable = true)
--- SOC_2021_3: string (nullable = true)
--- SOC_2021_3_NAME: string (nullable = true)
--- SOC_2021_4: string (nullable = true)
--- SOC_2021_4_NAME: string (nullable = true)
--- SOC_2021_5: string (nullable = true)
--- SOC_2021_5_NAME: string (nullable = true)
--- LOT_CAREER_AREA: integer (nullable = true)
--- LOT_CAREER_AREA_NAME: string (nullable = true)
--- LOT_OCCUPATION: integer (nullable = true)
--- LOT_OCCUPATION_NAME: string (nullable = true)
--- LOT_SPECIALIZED_OCCUPATION: integer (nullable = true)
--- LOT_SPECIALIZED_OCCUPATION_NAME: string (nullable = true)
--- LOT_OCCUPATION_GROUP: integer (nullable = true)
--- LOT_OCCUPATION_GROUP_NAME: string (nullable = true)
--- LOT_V6_SPECIALIZED_OCCUPATION: integer (nullable = true)
--- LOT_V6_SPECIALIZED_OCCUPATION_NAME: string (nullable = true)
--- LOT_V6_OCCUPATION: integer (nullable = true)
--- LOT_V6_OCCUPATION_NAME: string (nullable = true)
--- LOT_V6_OCCUPATION_GROUP: integer (nullable = true)
--- LOT_V6_OCCUPATION_GROUP_NAME: string (nullable = true)
--- LOT_V6_CAREER_AREA: integer (nullable = true)
--- LOT_V6_CAREER_AREA_NAME: string (nullable = true)
--- SOC_2: string (nullable = true)
--- SOC_2_NAME: string (nullable = true)
--- SOC_3: string (nullable = true)
--- SOC_3_NAME: string (nullable = true)
--- SOC_4: string (nullable = true)
--- SOC_4_NAME: string (nullable = true)
--- SOC_5: string (nullable = true)
--- SOC_5_NAME: string (nullable = true)
--- LIGHTCAST_SECTORS: string (nullable = true)
--- LIGHTCAST_SECTORS_NAME: string (nullable = true)
--- NAICS_2022_2: integer (nullable = true)
--- NAICS_2022_2_NAME: string (nullable = true)
--- NAICS_2022_3: integer (nullable = true)
--- NAICS_2022_3_NAME: string (nullable = true)
--- NAICS_2022_4: integer (nullable = true)
--- NAICS_2022_4_NAME: string (nullable = true)
--- NAICS_2022_5: integer (nullable = true)
--- NAICS_2022_5_NAME: string (nullable = true)
--- NAICS_2022_6: integer (nullable = true)
--- NAICS_2022_6_NAME: string (nullable = true)

```

	ID LAST_UPDATED_DATE LAST_UPDATED_TIMESTAMP DUPLICATES	POSTED	EXPIRED DURATION	SOUR
CE_TYPES	SOURCES	URL ACTIVE_URLS ACTIVE_SOURCES_INFO	TITLE_RAW	
BODY MODELED_EXPIRED MODELED_DURATION	COMPANY	COMPANY_NAME COMPANY_RAW COMPANY_IS_STAFFING	EDUCATION_LEVE	LS
EDUCATION_LEVELS NAME MIN EDULEVELS	MIN EDULEVELS NAME MAX EDULEVELS MAX EDULEVELS NAME EMPLOYMENT TYPE EMPLOY			

4/18

5/18

-----+
only showing top 5 rows

Salary Distribution by Employment Type

- Identify salary trends across different employment types.
- Filter the dataset**
 - Remove records where **salary is missing or zero**.
- Aggregate Data**
 - Group by **employment type** and compute salary distribution.
- Visualize results**
 - Create a **box plot** where:
 - X-axis** = EMPLOYMENT_TYPE_NAME
 - Y-axis** = SALARY_FROM
 - Customize **colors, fonts, and styles** to avoid a **2.5-point deduction**.
- Explanation:** Write two sentences about what the graph reveals.

```
In [8]: from pyspark.sql import SparkSession
from pyspark.sql.functions import col

# Initialize Spark Session
spark = SparkSession.builder.appName("LightcastData").getOrCreate()

# Load Data
df = spark.read.option("header", "true") \
    .option("inferSchema", "true") \
    .option("multiline", "true") \
    .option("escape", "\\") \
    .csv("lightcast_job_postings.csv")

# Check unique values
employment_types = df.select("EMPLOYMENT_TYPE_NAME").distinct().dropna()
employment_types.show(20, truncate=False)
```

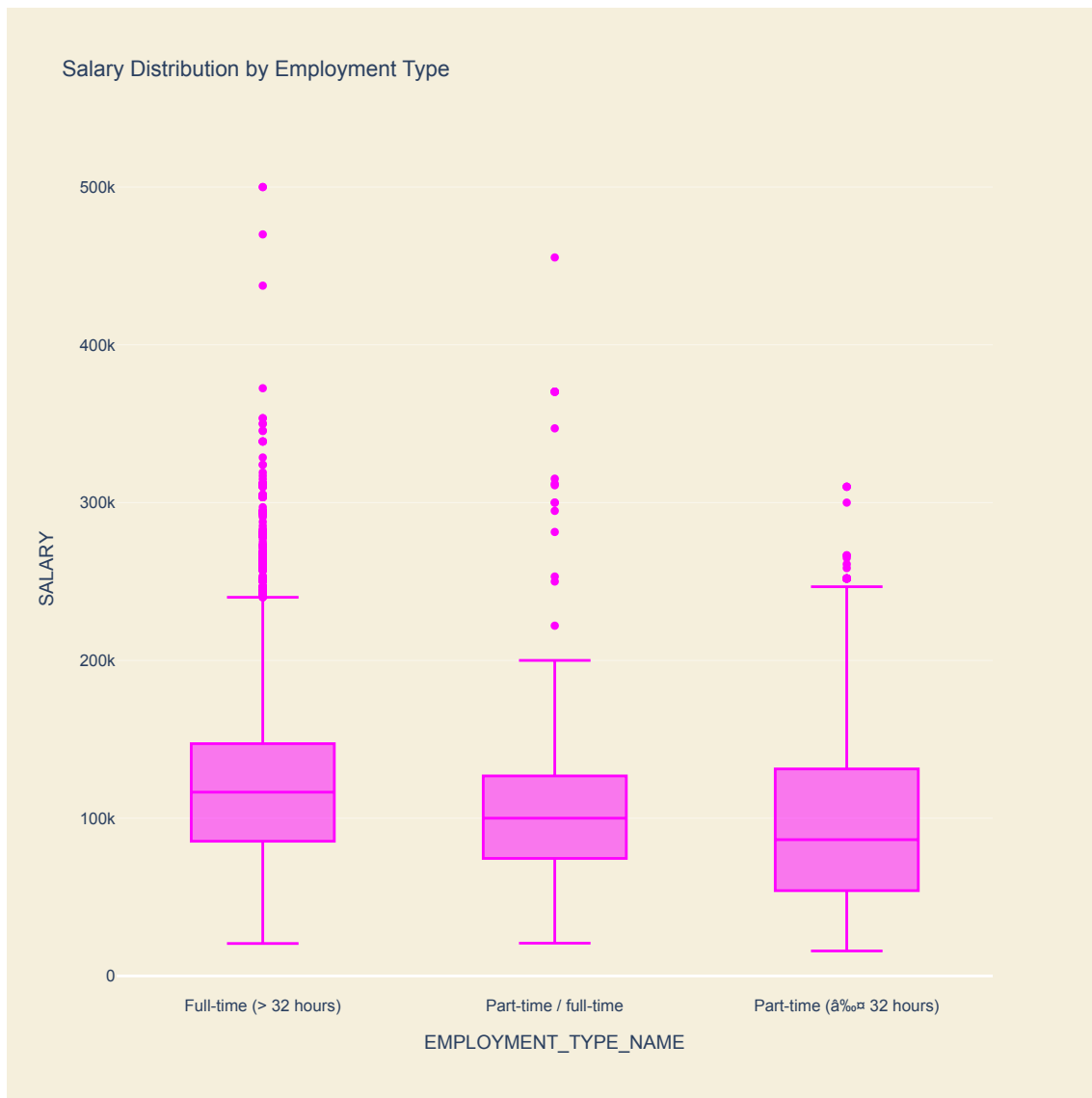
```
[Stage 14:> (0 + 1) / 1]
+-----+
|EMPLOYMENT_TYPE_NAME|
+-----+
|Part-time / full-time|
|Part-time (≈ 32 hours)|
|Full-time (> 32 hours)|
+-----+
```

```
In [9]: # Select needed columns and convert to Pandas
pdf = df.select("EMPLOYMENT_TYPE_NAME", "SALARY").toPandas()

# box plot
fig = px.box(pdf,
             x="EMPLOYMENT_TYPE_NAME",
             y="SALARY",
             title="Salary Distribution by Employment Type",
             color_discrete_sequence=["Magenta"],
             width=800,
             height=800)

# Customizing
fig.update_layout(
    font_family="Arial",
    title_font_size=16,
    plot_bgcolor="#f5efdc", # inner plot area
    paper_bgcolor="#f5efdc" # outer background area
)

fig.show()
```



Interpretation

The box plot shows that full-time jobs (over 32 hours) offer the highest median salary, around

116.5k, and also have a wider salary range with many high-paying outliers. Jobs that are both part-time and full-time come next with a median of about 100k, while part-time roles (under 32 hours) have the lowest median salary, around \$86.4k. Overall, full-time positions not only pay more on average, but also show more variation in salary compared to other job types.

Salary Distribution by Industry

- Compare salary variations across industries.
- **Filter the dataset**
 - Keep records where **salary is greater than zero**.
- **Aggregate Data**
 - Group by **NAICS industry codes**.
- **Visualize results**
 - Create a **box plot** where:
 - **X-axis** = `NAICS2_NAME`
 - **Y-axis** = `SALARY_FROM`
 - Customize colors, fonts, and styles.
- **Explanation:** Write two sentences about what the graph reveals.

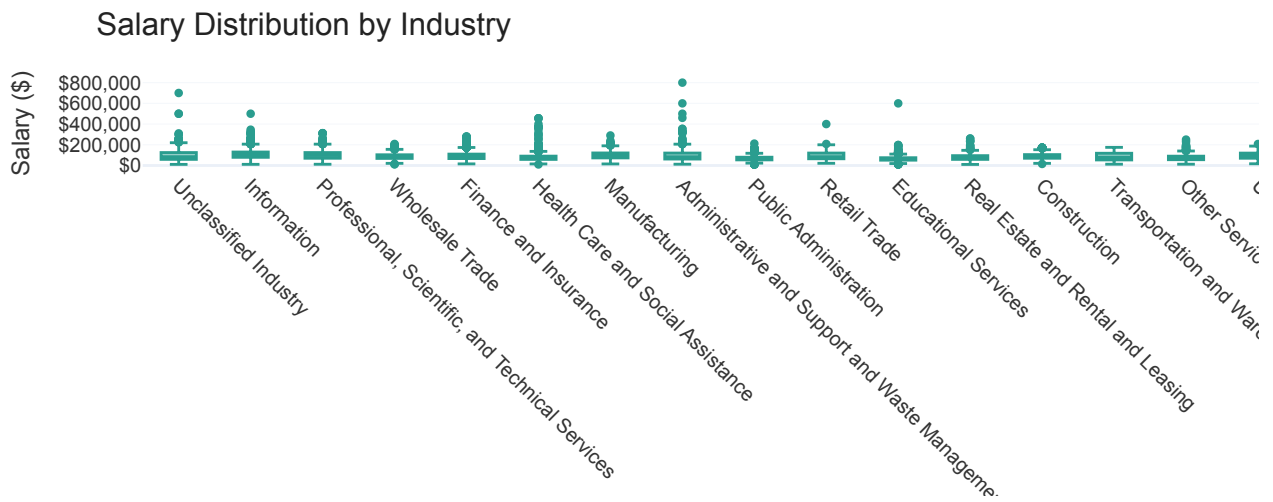
```
In [10]: # Step 1: Filter and convert to Pandas
filtered_df = df.filter((df.SALARY_FROM > 0) & df.NAICS2_NAME.isNotNull())
salary_industry_df = filtered_df.select("NAICS2_NAME", "SALARY_FROM").toPandas()

# Step 2: Create box plot
import plotly.express as px

fig = px.box(
    salary_industry_df,
    x="NAICS2_NAME",
    y="SALARY_FROM",
    title="Salary Distribution by Industry",
    color_discrete_sequence=["#2a9d8f"],
    template="plotly_white"
)

# Step 3: Layout styling
fig.update_layout(
    title_font=dict(size=22, family="Arial", color="#222"),
    font=dict(family="Arial", size=14, color="#333"),
    xaxis_title="Industry",
    yaxis_title="Salary ($)",
    yaxis_tickprefix="$",
    yaxis_tickformat=",",
    xaxis_tickangle=45,
    margin=dict(t=60, l=60, r=40, b=100),
    plot_bgcolor="#ffffff",
    paper_bgcolor="#ffffff"
)

fig.show()
```



Interpretation

The box plot shows that industries like Administrative, Manufacturing, and Information tend to offer higher salary ranges and include more high-paying outliers, suggesting there's good potential for top-tier salaries in these fields. On the other hand, industries like Arts, Entertainment, and Accommodation usually have lower and more tightly packed salaries, meaning the pay is generally more limited and consistent in those sectors.

Job Posting Trends Over Time

- Analyze how job postings fluctuate over time.
- **Aggregate Data**
 - Count job postings per **posted date (POSTED)**.
- **Visualize results**
 - Create a **line chart** where:
 - **X-axis** = POSTED
 - **Y-axis** = Number of Job Postings
 - Apply custom colors and font styles.
- **Explanation:** Write two sentences about what the graph reveals.

```
In [11]: # 1. Import Spark functions and convert POSTED to date
from pyspark.sql.functions import to_date, count

# 2. Group by date and count job postings
df_posting_trend = df.withColumn("POSTED_DATE", to_date("POSTED")) \
    .groupBy("POSTED_DATE") \
    .agg(count("*").alias("count")) \
    .orderBy("POSTED_DATE")

# 3. Convert to Pandas for Plotly
pdf_posting_trend = df_posting_trend.toPandas()
```

```
In [12]: # Group by POSTED date and count number of job postings
df_posting_trend = df.groupBy("POSTED").count().orderBy("POSTED")

# Convert to Pandas for plotting
pdf_posting_trend = df_posting_trend.toPandas()

# Create a line chart
import plotly.express as px

fig = px.line(pdf_posting_trend,
              x="POSTED",
              y="count",
              title="Job Posting Trends Over Time",
              markers=True,
              color_discrete_sequence=["steelblue"],
              width=1000,
              height=500)

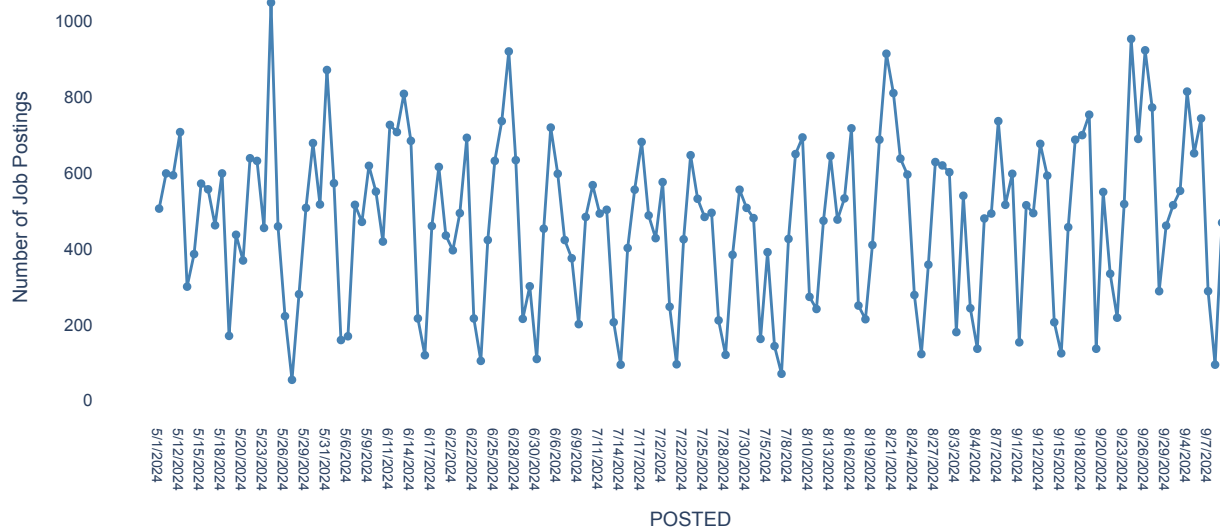
# Customization
fig.update_layout(
    font_family="Arial",
    title_font_size=16,
    plot_bgcolor="white",
    paper_bgcolor="white",
)

fig.update_xaxes(
    tickformat="%b %d", # e.g., May 05
    tickfont=dict(size=10),
    showgrid=False
)

fig.update_yaxes(
    title_text="Number of Job Postings",
    showgrid=True
)

fig.show()
```

Job Posting Trends Over Time



Interpretation

The line chart shows that job postings go up and down quite a bit day to day, which suggests that hiring activity can be pretty unpredictable. There's a clear spike in late August to mid-September, likely because companies are hiring after the summer slowdown. Even with all the ups and downs, there doesn't seem to be a consistent long-term trend — things stay fairly balanced over time.

Top 10 Job Titles by Count

- Identify the most frequently posted job titles.
- **Aggregate Data**
 - Count the occurrences of each **job title (TITLE_NAME)**.
 - Select the **top 10 most frequent titles**.
- **Visualize results**
 - Create a **bar chart** where:
 - **X-axis** = **TITLE_NAME**
 - **Y-axis** = **Job Count**
 - Apply custom colors and font styles.
- **Explanation:** Write two sentences about what the graph reveals.

```
In [13]: # Group by TITLE_NAME, count, and get top 10
df_titles = df.groupby("TITLE_NAME").count().orderBy("count", ascending=False).limit(10)

# Convert to Pandas
pdf_titles = df_titles.toPandas()
```

```
In [14]: import plotly.express as px

fig = px.bar(
    pdf_titles,
    x="TITLE_NAME",
    y="count",
    title="Top 10 Job Titles by Count",
    color_discrete_sequence=["#8B4513"], # Hex for "dark brown"
    width=900,
    height=500
)

fig.update_layout(
```

```

title_font=dict(size=22, family="Arial", color="#222"),
font=dict(family="Arial", size=14, color="#333"),
xaxis_title="Job Title",
yaxis_title="Job Count",
plot_bgcolor="#ffffff",
paper_bgcolor="#ffffff",
margin=dict(t=60, l=60, r=40, b=100)
)

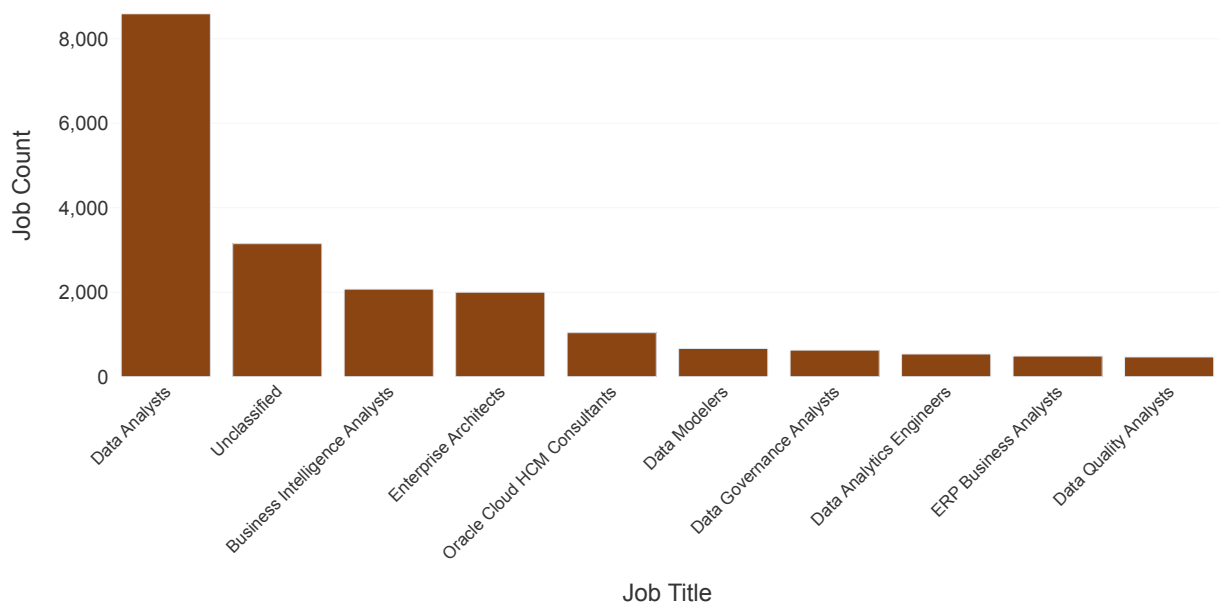
fig.update_xaxes(
    tickangle=-45,
    tickfont=dict(size=12)
)

fig.update_yaxes(
    tickformat="," ,
    gridcolor="#eeeeee"
)

fig.show()

```

Top 10 Job Titles by Count



Interpretation

The job title "Data Analysis" stands out as the most posted role by far, way ahead of the others. Titles like Business Intelligence Analyst and Enterprise Architect also appear frequently, showing that tech and data roles are in high demand. After the top few titles, the number of postings drops off quickly, which means employers are mainly focused on hiring for a small set of key positions.

Remote vs On-Site Job Postings

- Compare the proportion of remote and on-site job postings.
- **Aggregate Data**
 - Count job postings by **remote type** (**REMOTE_TYPE_NAME**).
- **Visualize results**
 - Create a **pie chart** where:
 - **Labels** = **REMOTE_TYPE_NAME**
 - **Values** = **Job Count**
 - Apply custom colors and font styles.
- **Explanation:** Write two sentences about what the graph reveals.

```
In [15]: from pyspark.sql.functions import col
```

```
# Filter for valid remote types and count them
df_remote_filtered = df.filter(
    col("REMOTE_TYPE_NAME").isin(["Remote", "Hybrid Remote", "Not Remote"])
).groupBy("REMOTE_TYPE_NAME").count()

# Convert to Pandas
pdf_remote_filtered = df_remote_filtered.toPandas()
```

```
In [16]: import plotly.express as px

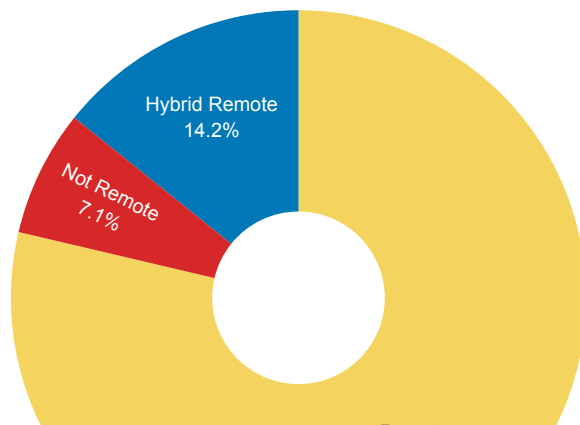
fig = px.pie(
    pdf_remote_filtered,
    names="REMOTE_TYPE_NAME",
    values="count",
    title="Remote vs On-Site Job Postings",
    color_discrete_sequence=["#f4d35e", "#0077b6", "#d62828"]
,
    hole=0.3 # makes it a donut chart - optional but stylish
)

# Custom styling
fig.update_traces(
    textinfo='percent+label',
    textfont_size=14
)

fig.update_layout(
    title_font=dict(size=22, family="Arial", color="#222"),
    font=dict(family="Arial", size=14, color="#333"),
    plot_bgcolor="#ffffff",
    paper_bgcolor="#ffffff",
    margin=dict(t=60, l=60, r=40, b=60)
)

fig.show()
```

Remote vs On-Site Job Postings



Interpretation

Most of the job postings are for remote positions, which shows that working from home has become the new normal for many companies. There are still some hybrid and fully on-site jobs out there, but they make up a much smaller part of the market compared to remote roles.

Skill Demand Analysis by Industry (Stacked Bar Chart)

- Identify which skills are most in demand in various industries.
- **Aggregate Data**
 - Extract **skills** from job postings.
 - Count occurrences of skills grouped by **NAICS industry codes**.
- **Visualize results**
 - Create a **stacked bar chart** where:
 - **X-axis** = `Industry`
 - **Y-axis** = `Skill Count`
 - **Color** = `Skill`
 - Apply custom colors and font styles.
- **Explanation:** Write two sentences about what the graph reveals.

```
In [17]: from pyspark.sql.functions import explode, split, trim

# Splitting the SKILLS_NAME column into an array of skills
df_split = df.withColumn("Skill", explode(split("SKILLS_NAME", ",")))

# Cleaning whitespace from skill names
df_split = df_split.withColumn("Skill", trim(df_split["Skill"]))

# Filtering nulls if needed
df_split = df_split.filter(df_split["Skill"].isNotNull())
```

```
In [18]: # Register the cleaned & exploded DataFrame as a new SQL view
df_split.createOrReplaceTempView("job_data")

# Updated Spark SQL query
query = """
WITH SkillCounts AS (
    SELECT
        NAICS_2022_2_NAME AS Industry,
        Skill,
        COUNT(*) AS Skill_Count
    FROM job_data
    WHERE Skill IS NOT NULL AND NAICS_2022_2_NAME IS NOT NULL
    GROUP BY Industry, Skill
),
RankedSkills AS (
    SELECT *,
        ROW_NUMBER() OVER (PARTITION BY Industry ORDER BY Skill_Count DESC) AS rank
    FROM SkillCounts
)
SELECT * FROM RankedSkills WHERE rank <= 10
"""

# Execute and store results
df_top_skills = spark.sql(query)
```

```
In [19]: pdf_top_skills = df_top_skills.toPandas()
```

```
In [20]: import plotly.express as px

fig = px.bar(
    pdf_top_skills,
    x="Industry",
    y="Skill_Count",
    color="Skill",
    title="Top 10 In-Demand Skills by Industry",
    text_auto=True,
    barmode='stack',
    color_discrete_sequence=px.colors.qualitative.Set2, # brighter categorical palette
    width=1200,
    height=700
)

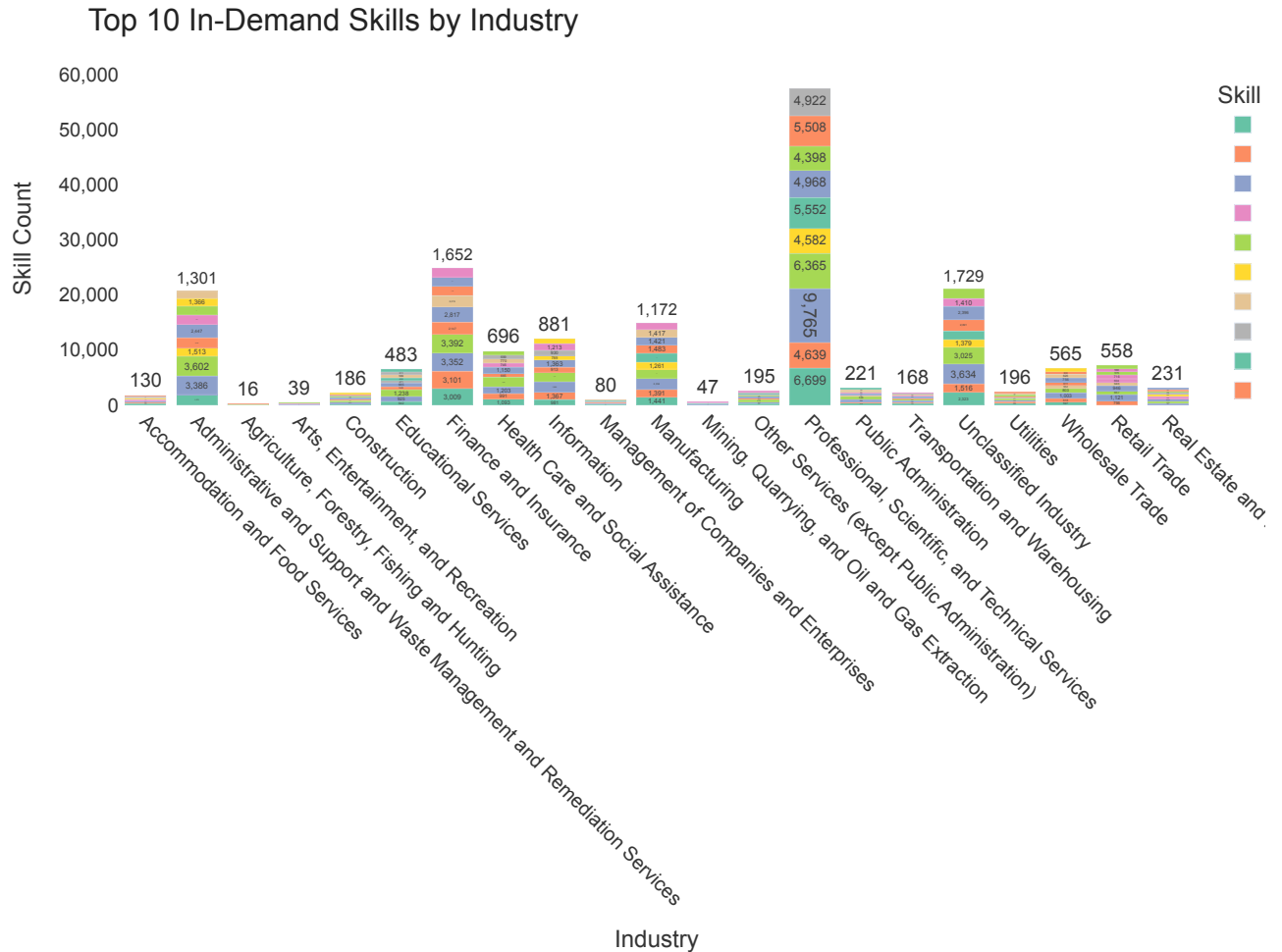
fig.update_layout(
    title_font=dict(size=22, family="Arial", color="#222"),
    font=dict(family="Arial", size=14, color="#333"),
```

```

axis_title="Industry",
yaxis_title="Skill Count",
plot_bgcolor="#ffffff",
paper_bgcolor="#ffffff",
margin=dict(t=60, l=60, r=40, b=120),
xaxis_tickangle=45
)

fig.update_yaxes(tickformat=",")
fig.show()

```



Interpretation

Skills like "Management" and "Leadership" are in high demand across almost every industry, showing how important they are no matter the field. The Professional, Scientific, and Technical Services sector stands out with a wide range of skill needs, especially for things like Data Analysis, Agile Methodology, and Communication. More strategic skills like "Business Requirements", "Project Management", and "Planning" are mostly found in Unclassified and Professional industries, where coordination and planning are key. While industries like Retail Trade and Real Estate have fewer overall postings, they still consistently ask for customer service and problem-solving skills.

Salary Analysis by ONET Occupation Type (Bubble Chart)

- Analyze how salaries differ across ONET occupation types.
- Aggregate Data**
 - Compute **median salary** for each occupation in the **ONET taxonomy**.
- Visualize results**
 - Create a **bubble chart** where:
 - X-axis** = `ONET_NAME`
 - Y-axis** = `Median Salary`

- **Size** = Number of job postings
- Apply custom colors and font styles.
- **Explanation:** Write two sentences about what the graph reveals.

```
In [21]: df.createOrReplaceTempView("job_data")
```

```
In [22]: query_salary = """
WITH SalaryStats AS (
    SELECT
        TITLE_NAME,
        percentile_approx(SALARY, 0.5) AS Median_Salary,
        COUNT(*) AS Job_Count
    FROM job_data
    WHERE SALARY IS NOT NULL AND TITLE_NAME IS NOT NULL
    GROUP BY TITLE_NAME
)
SELECT * FROM SalaryStats
"""
df_salary_stats = spark.sql(query_salary)
```

```
In [23]: pdf_salary_stats = df_salary_stats.toPandas()
```

```
In [24]: import plotly.express as px

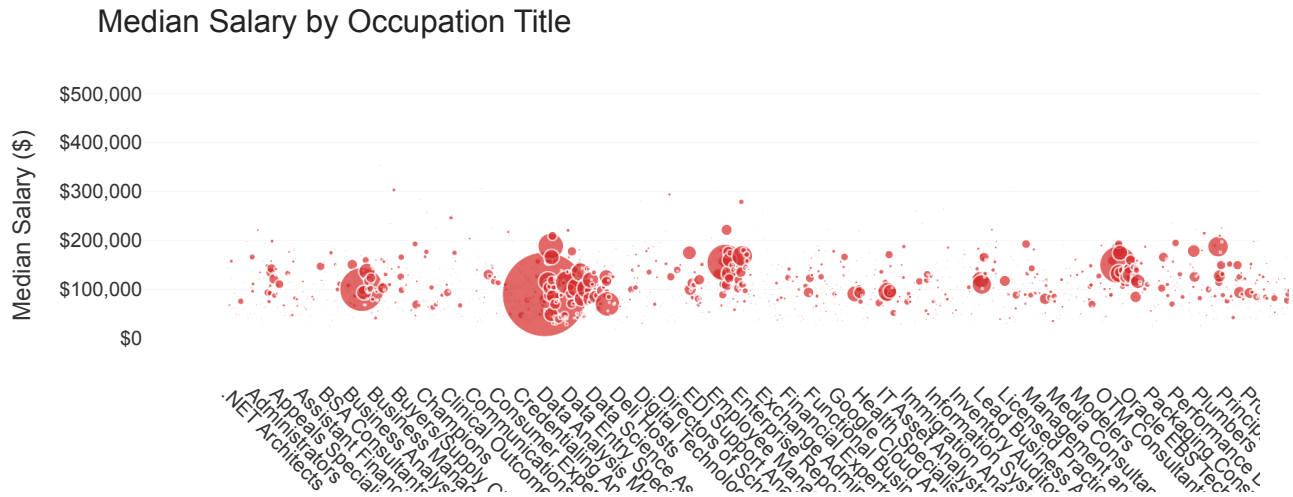
fig = px.scatter(
    pdf_salary_stats,
    x="TITLE_NAME",          # X-axis: Job Title
    y="Median_Salary",       # Y-axis: Median Salary
    size="Job_Count",        # Bubble size = Number of postings
    color_discrete_sequence=["#d62828"], # Bold red
    title="Median Salary by Occupation Title",
    labels={
        "TITLE_NAME": "Occupation Title",
        "Median_Salary": "Median Salary ($)",
        "Job_Count": "Number of Job Postings"
    },
    size_max=60 # Controls the max bubble size
)

fig.update_layout(
    title_font=dict(size=22, family="Arial", color="#222"),
    font=dict(family="Arial", size=14, color="#333"),
    xaxis_title="Occupation Title",
    yaxis_title="Median Salary ($)",
    xaxis_tickangle=45,
    plot_bgcolor="#ffffff",
    paper_bgcolor="#ffffff",
    margin=dict(t=60, l=60, r=40, b=120)
)

fig.update_yaxes(
    tickformat=",$.0f",      # e.g. $50,000
    gridcolor="#e0e0e0"
)

fig.update_traces(marker=dict(opacity=0.7, line=dict(width=1, color='white')))

fig.show()
```



```
In [25]: import plotly.express as px

fig = px.scatter(
    pdf_salary_stats,
    x="TITLE_NAME",          # X-axis: Job Title
    y="Median_Salary",       # Y-axis: Median Salary
    size="Job_Count",        # Bubble size = Number of postings
    color="Median_Salary",    # Bubble color = Median Salary value
    color_continuous_scale="Viridis", # Gradient: light to dark (low to high salary)
    title="Median Salary by Occupation Title",
    labels={
        "TITLE_NAME": "Occupation Title",
        "Median_Salary": "Median Salary ($)",
        "Job_Count": "Number of Job Postings"
    },
    size_max=60
)

# Layout styling
fig.update_layout(
    title_font=dict(size=22, family="Arial", color="#222"),
    font=dict(family="Arial", size=14, color="#333"),
    xaxis_title="Occupation Title",
    yaxis_title="Median Salary ($)",
    xaxis_tickangle=45,
    plot_bgcolor="#ffffff",
    paper_bgcolor="#ffffff",
    margin=dict(t=60, l=60, r=40, b=120),
    coloraxis_colorbar=dict(
        title="Median Salary",
        tickprefix="$",
        ticks="outside"
    )
)

# Y-axis formatting and bubble borders
fig.update_yaxes(
    tickformat=",$.0f",
    gridcolor="#e0e0e0"
)

fig.update_traces(marker=dict(opacity=0.75, line=dict(width=1, color='white')))

fig.show()
```



```
# Create mapping to integer index
label_to_index = {label: idx for idx, label in enumerate(all_labels)}

# Map source and target to indices
pdf_transitions['source_id'] = pdf_transitions['source'].map(label_to_index)
pdf_transitions['target_id'] = pdf_transitions['target'].map(label_to_index)
```

```
In [29]: import plotly.graph_objects as go

fig = go.Figure(data=[go.Sankey(
    node=dict(
        pad=15,
        thickness=20,
        line=dict(color="black", width=0.5),
        label=all_labels,
        color="#ffb3b3" # Softer red-pink for node fill
    ),
    link=dict(
        source=pdf_transitions['source_id'],
        target=pdf_transitions['target_id'],
        value=pdf_transitions['value'],
        color="rgba(214, 39, 40, 0.5)" # Semi-transparent red
    )
)])

fig.update_layout(
    title_text="Career Pathway Trends (SOC Code Transitions)",
    title_font=dict(size=22, family="Arial", color="#222"),
    font=dict(family="Arial", size=13, color="#333"),
    paper_bgcolor="#ffffff", # Clean white background
    margin=dict(t=60, l=60, r=40, b=60)
)

fig.show()
```

Career Pathway Trends (SOC Code Transitions)



Interpretation: Career Pathway Trends (SOC Code Transitions)

The diagram shows that many people are moving from Computer and Mathematical jobs into Mathematical Science roles, which means there's a lot of movement between these two closely related fields. This suggests that the skills in these areas often overlap, making it easier for professionals to switch or grow their careers within the tech and math space. The thick red band highlights how strong that connection is — it's one of the most common transitions in the data. Employers and educators can use this insight to create training programs that help tech workers transition into roles in mathematical sciences.