

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import re
```

▼ Data Understanding

```
1 df_time = pd.read_excel("Sample Month End Report for Assesment.xlsx", sheet_name="Time spent Analysis")
2 df_key = pd.read_excel("Sample Month End Report for Assesment.xlsx", sheet_name="Key Accounts")
3 df_meet = pd.read_excel("Sample Month End Report for Assesment.xlsx", sheet_name="Meeting Log")
4 df_inquiry = pd.read_excel("Sample Month End Report for Assesment.xlsx", sheet_name="Inquiry Tracker")
5
```

```
1 df_time.describe()
2 df_time.info()
3 df_time.head()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 142 entries, 0 to 141
Data columns (total 6 columns):
 #   Column                                Non-Null Count  Dtype
---  -
0   Company                              142 non-null    object
1   Category                             142 non-null    object
2   Total Hours                          142 non-null    float64
3   Estimated Cost                       142 non-null    int64
4   3500                                 0 non-null     float64
5   inr per hour is cost assumed        0 non-null     float64
dtypes: float64(3), int64(1), object(2)
memory usage: 6.8+ KB
```

	Company	Category	Total Hours	Estimated Cost	3500	inr per hour is cost assumed
0	Chennai	Unknown	32.00	112000	NaN	NaN
1	Pidilite	Distributor	23.75	83125	NaN	NaN
2	Kalpataru Projects	Contractor	23.50	82250	NaN	NaN
3	Asg win	Applicator	15.50	54250	NaN	NaN
4	Internal/Review (MRM, Hubspot, etc.)	Admin	12.50	43750	NaN	NaN

Next steps: [Generate code with df_time](#) [New interactive sheet](#)

```
1 df_key.head(10)
2 df_key.info()
3 df_key.describe()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 42 entries, 0 to 41
Data columns (total 6 columns):
 #   Column                                Non-Null Count  Dtype
---  -
0   Location                              42 non-null    object
1   Key Account Names (company name)      42 non-null    object
2   Types of Projects                     42 non-null    object
3   Number of Decision Makers/ Sub Teams  41 non-null    object
4   Estimated Vol of Boq Per Year with EJC 42 non-null    object
5   Target number of Boq from the account  42 non-null    int64
dtypes: int64(1), object(5)
memory usage: 2.1+ KB
```

Target number of Boq from the account	
count	42.000000
mean	2.500000
std	2.520646
min	1.000000
25%	1.000000
50%	1.000000
75%	3.000000
max	13.000000

```
1 df_meet.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 468 entries, 0 to 467
Data columns (total 13 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                -
0   Date of Activity                      467 non-null   object
1   Activity Type                        159 non-null   object
2   Company                             466 non-null   object
3   Project                             450 non-null   object
4   Hours spent                         447 non-null   object
5   Accomplishment / Summary            404 non-null   object
6   Remarks                             301 non-null   object
7   Unnamed: 7                          0 non-null     float64
8   Unnamed: 8                          0 non-null     float64
9   Unnamed: 9                          2 non-null     object
10  Unnamed: 10                         12 non-null    object
11  Unnamed: 11                         2 non-null     object
12  Unnamed: 12                         1 non-null     object
dtypes: float64(2), object(11)
memory usage: 47.7+ KB
```

```
1 df_inquiry.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30 entries, 0 to 29
Data columns (total 22 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                -
0   Inquiry Date                        26 non-null    object
1   Projects                           25 non-null    object
2   Number of Hours                    15 non-null    float64
3   Value/Hours                        8 non-null     float64
4   Feasability Ratio                  0 non-null     float64
5   Expected Close date                25 non-null    object
6   Unnamed: 6                        0 non-null     float64
7   Construction type                  25 non-null    object
8   Stage                             25 non-null    object
9   Source                             25 non-null    object
10  Quantity in LM                     9 non-null     object
11  Value ($)                          26 non-null    float64
12  Value (INR)                       26 non-null    float64
13  Quoted Via Applicator              25 non-null    object
14  Applicator                         23 non-null    object
15  Lost Reason                        0 non-null     float64
16  Unnamed: 16                       0 non-null     float64
17  Unnamed: 17                       2 non-null     object
18  Unnamed: 18                       7 non-null     object
19  Unnamed: 19                       7 non-null     float64
20  Unnamed: 20                       2 non-null     object
21  Unnamed: 21                       2 non-null     object
dtypes: float64(9), object(13)
memory usage: 5.3+ KB
```

▼ Data Cleaning

```
1 def normalize_name(s):
2     if pd.isna(s):
3         return s
4     s = str(s).lower().replace("&", " and ")
5     s = re.sub(r'[\\(\).,;:/\-\ ]', ' ', s)
6     # collapse multiple spaces into one
7     s = re.sub(r'\s+', ' ', s).strip()
8     return s
```

```
1 #Here we have to clean the data of the first sheet that is Time spent analysis.
2 #Here we have to clean the data of the first sheet that is Time spent analysis.
3 df_time = pd.read_excel("Sample Month End Report for Assesment.xlsx", sheet_name="Time spent Analysis")
4 df_time = df_time.dropna(axis=1, how='all')
5 df_time.columns = ['company', 'category', 'total_hours', 'estimated_cost']
6
7 df_time['category']=df_time['category'].str.strip().str.lower().str.replace(r'\s+', ' ', regex=True)
8 df_time['company']=df_time['company'].str.strip().str.lower().str.replace(r'\s+', ' ', regex=True)
9
10 category_mapping = {
11     'consultant/architect': 'Consultant',
12     'developer': 'Developer',
13     'developers/projects': 'Developer',
14     'developers/office/projects': 'Developer',
15     'developers/groups': 'Developer',
16     'developers/group': 'Developer',
17     'estates/developers': 'Developer',
```

```

18     'distributor': 'Distributor',
19     'contractor': 'Contractor',
20     'applicator': 'Applicator',
21     'admin': 'Admin',
22     'unknown': 'Unknown'
23 }
24
25 df_time['category'] = df_time['category'].map(category_mapping).fillna(df_time['category'])
26
27 # Calculate cost per hour
28 df_time['cost_per_hour'] = df_time['estimated_cost'] / df_time['total_hours']
29
30 # Apply the normalize function to the column
31 df_time['company'] = df_time['company'].apply(normalize_name)
32 df_time['category'] = df_time['category'].apply(normalize_name)

```

```

1 # Cleaning of the second sheet
2 df_key = pd.read_excel("Sample Month End Report for Assesment.xlsx", sheet_name="Key Accounts")
3
4 df_key_Raw = df_key.copy()
5 rename_columns = {
6     'Key Account Names (company name)': 'company_name',
7     'Types of Projects': 'types_of_projects',
8     'Number of Decision Makers/ Sub Teams': 'num_decision_makers',
9     'Estimated Vol of Boq Per Year with EJC': 'estimated_boq_per_year',
10    'Target number of Boq from the account': 'target_boq'
11 }
12 df_key = df_key.rename(columns=rename_columns)
13
14 def normalize_to_lower(s):
15     if pd.isna(s):
16         return s
17     # Moved the processing lines outside the 'if' block so they are always executed for non-NaN values
18     s = str(s).lower().replace('or', ' ').replace('to', ' ').strip()
19     s = re.sub(r'\s+', ' ', s).strip()
20
21     match = re.search(r'(\d+)', s)
22     if match:
23         return int(match.group(1))
24     return s
25
26
27 df_key['estimated_boq_per_year'] = df_key['estimated_boq_per_year'].apply(normalize_to_lower)
28 df_key['target_boq'] = df_key['target_boq'].apply(normalize_to_lower)
29
30 df_key['estimated_boq_int'] = pd.to_numeric(df_key['estimated_boq_per_year'], errors='coerce')
31 df_key['target_boq'] = pd.to_numeric(df_key['target_boq'], errors='coerce').astype('Int64')
32
33 df_key['account_norm'] = df_key['company_name'].apply(normalize_name)

```

```

1 df_key[['company_name', 'account_norm', 'estimated_boq_per_year', 'estimated_boq_int']].head(10)
2

```

	company_name	account_norm	estimated_boq_per_year	estimated_boq_int	
0	AECOM	aecom	5	5	
1	AEDBM	aedbm	1	1	
2	Aedium Designs	aedium designs	2	2	
3	Archetype	archetype	3	3	
4	Ashok Narayan & Eshwar	ashok narayan and eshwar	5	5	
5	Base	base	1	1	
6	Besten Consultants	besten consultants	1	1	
7	Bureau Engineers (structural)	bureau engineers structural	5	5	
8	Bureau Hepolt	bureau hepolt	1	1	
9	Chandavarkar and Thacker Architects	chandavarkar and thacker architects	2	2	

```
1 df_meet.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 468 entries, 0 to 467
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Date of Activity       467 non-null   object
1   Activity Type          159 non-null   object

```

```

2 Company          466 non-null    object
3 Project          450 non-null    object
4 Hours spent      447 non-null    object
5 Accomplishment / Summary  404 non-null    object
6 Remarks          301 non-null    object
7 Unnamed: 7       0 non-null      float64
8 Unnamed: 8       0 non-null      float64
9 Unnamed: 9       2 non-null      object
10 Unnamed: 10     12 non-null     object
11 Unnamed: 11     2 non-null     object
12 Unnamed: 12     1 non-null     object
dtypes: float64(2), object(11)
memory usage: 47.7+ KB

```

```

1 df_meet = pd.read_excel("Sample Month End Report for Assesment.xlsx", sheet_name="Meeting Log")
2
3 df_meet_raw = df_meet.copy()
4 df_meet = df_meet.dropna(axis=1, how='all')
5 df_meet = df_meet.loc[:, ~df_meet.columns.str.contains('^Unnamed')]
6
7 keep_cols = [
8     'Date of Activity',
9     'Activity Type',
10    'Company',
11    'Project',
12    'Hours spent',
13    'Accomplishment / Summary',
14    'Remarks'
15 ]
16
17 df_meet = df_meet[keep_cols]
18
19 df_meet['Date of Activity'] = pd.to_datetime(df_meet['Date of Activity'], errors='coerce')
20
21
22 df_meet['Hours spent'] = pd.to_numeric(df_meet['Hours spent'], errors='coerce')
23
24 df_meet = df_meet[~df_meet['Project'].isna()]
25
26 df_meet['Activity Type'] = df_meet['Activity Type'].fillna('Unknown')
27
28 df_meet['project_norm'] = df_meet['Project'].apply(normalize_name)
29 df_meet['company_norm'] = df_meet['Company'].apply(normalize_name)
30

```

```

1 #SHEET-4 CLEANING
2
3
4 df_inquiry = pd.read_excel("Sample Month End Report for Assesment.xlsx",
5                             sheet_name="Inquiry Tracker")
6
7 df_inquiry_raw = df_inquiry.copy() # backup
8
9
10 # ----- 3. Remove totally empty columns + header junk rows -----
11 # Drop completely empty columns
12 df_inquiry = df_inquiry.dropna(axis=1, how='all')
13
14 # Drop first 2 rows if they are header/notes (as you observed earlier)
15 df_inquiry = df_inquiry.iloc[2:].reset_index(drop=True)
16
17 # Drop Unnamed columns (Excel artifacts)
18 df_inquiry = df_inquiry.loc[:, ~df_inquiry.columns.str.contains('^Unnamed')]
19
20
21 # ----- 4. Quick check of columns (just to see real names) -----
22 print("Columns in df_inquiry:")
23 print(df_inquiry.columns)
24
25
26 # ----- 5. Convert dates -----
27 # Inquiry Date
28 if 'Inquiry Date' in df_inquiry.columns:
29     df_inquiry['Inquiry Date'] = pd.to_datetime(df_inquiry['Inquiry Date'], errors='coerce')
30
31
32 # ----- 6. Clean Expected Close date -----
33 def clean_expected_close(v):
34     """
35     Handle values like:
36     - '2026' -> 2026-12-31
37     - 'oct 2026' -> 2026-10-01

```

```

38 - 'project on hold' -> NaT, on_hold=True
39 """
40 if pd.isna(v):
41     return pd.NaT, False
42
43 s = str(v).strip().lower()
44
45 # On hold / no date
46 if 'hold' in s:
47     return pd.NaT, True
48
49 # Pure year '2026'
50 if re.fullmatch(r'\d{4}', s):
51     year = int(s)
52     return pd.Timestamp(year=year, month=12, day=31), False
53
54 # Month + year like 'oct 2026'
55 try:
56     parts = s.split()
57     if len(parts) == 2 and re.search(r'[a-zA-Z]', parts[0]) and re.fullmatch(r'\d{4}', parts[1]):
58         s2 = '1 ' + s # assume first day of month
59         dt = pd.to_datetime(s2, errors='coerce', dayfirst=True)
60     else:
61         dt = pd.to_datetime(s, errors='coerce', dayfirst=True)
62 except Exception:
63     dt = pd.NaT
64
65 return dt, False
66
67 if 'Expected Close date' in df_inquiry.columns:
68     tmp = df_inquiry['Expected Close date'].apply(clean_expected_close)
69     df_inquiry['expected_close_clean'] = tmp.apply(lambda x: x[0])
70     df_inquiry['on_hold_flag'] = tmp.apply(lambda x: x[1])
71 else:
72     df_inquiry['expected_close_clean'] = pd.NaT
73     df_inquiry['on_hold_flag'] = False
74
75
76 # ----- 7. Convert numeric columns -----
77 numeric_cols = [
78     'Number of Hours',
79     'Value/Hours',
80     'Feasability Ratio',
81     'Quantity in LM',
82     'Value ($)',
83     'Value (INR)'
84 ]
85
86 for col in numeric_cols:
87     if col in df_inquiry.columns:
88         df_inquiry[col] = pd.to_numeric(df_inquiry[col], errors='coerce')
89
90
91 # ----- 8. Basic missing value handling -----
92 # Drop rows with missing Projects (no project = not useful for project-level analysis)
93 if 'Projects' in df_inquiry.columns:
94     df_inquiry = df_inquiry[~df_inquiry['Projects'].isna()]
95
96 # Fill Stage nulls with 'Unknown'
97 if 'Stage' in df_inquiry.columns:
98     df_inquiry['Stage'] = df_inquiry['Stage'].fillna('Unknown')
99
100 # If you have a Construction type column, fill nulls
101 for col in df_inquiry.columns:
102     if col.strip().lower() == 'construction type':
103         df_inquiry[col] = df_inquiry[col].fillna('Unknown')
104
105
106 # ----- 9. Normalize project name for joining -----
107 if 'Projects' in df_inquiry.columns:
108     df_inquiry['project_norm'] = df_inquiry['Projects'].apply(normalize_name)
109 else:
110     df_inquiry['project_norm'] = pd.NA
111

```

```

Columns in df_inquiry:
Index(['Inquiry Date', 'Projects', 'Number of Hours', 'Value/Hours',
      'Expected Close date', 'Construction type', 'Stage', 'Source',
      'Quantity in LM', 'Value ($)', 'Value (INR)', 'Quoted Via Applicator',
      'Applicator'],
      dtype='object')

```

▼ Data Analysis

```

1 # =====
2 # SIMPLE KPIs - SHEET 1 (TIME SPENT)
3 # =====
4
5 # KPI 1: Companies consuming the most hours
6 kpi_hours_by_company = (
7     df_time.groupby('company')['total_hours']
8     .sum()
9     .sort_values(ascending=False)
10 )
11 print("KPI 1: Hours by company")
12 print(kpi_hours_by_company.head(10), "\n")
13
14
15 # KPI 2: Per-hour return per company
16 df_time['return_per_hour'] = df_time['estimated_cost'] / df_time['total_hours']
17 kpi_return_per_hour = (
18     df_time.groupby('company')['return_per_hour']
19     .mean()
20     .sort_values(ascending=False)
21 )
22 print("KPI 2: Return per hour by company")
23 print(kpi_return_per_hour.head(10), "\n")
24
25
26 # KPI 3: Hours by category
27 kpi_hours_by_category = (
28     df_time.groupby('category')['total_hours']
29     .sum()
30     .sort_values(ascending=False)
31 )
32 print("KPI 3: Hours by category")
33 print(kpi_hours_by_category, "\n")
34
35
36 # KPI 4: Revenue by category
37 kpi_revenue_by_category = (
38     df_time.groupby('category')['estimated_cost']
39     .sum()
40     .sort_values(ascending=False)
41 )
42 print("KPI 4: Revenue by category")
43 print(kpi_revenue_by_category, "\n")
44
45
46 # KPI 5: Number of entries per company
47 kpi_entries_by_company = df_time['company'].value_counts()
48 print("KPI 5: Entry count per company")
49 print(kpi_entries_by_company.head(10), "\n")
50
51
52 # KPI 6: Revenue per company
53 kpi_revenue_by_company = (
54     df_time.groupby('company')['estimated_cost']
55     .sum()
56     .sort_values(ascending=False)
57 )
58 print("KPI 6: Revenue by company")
59 print(kpi_revenue_by_company.head(10), "\n")
60
61 # KPI 7
62 avg_cost_per_hour = df_time['cost_per_hour'].mean()
63 df_time['CEI'] = df_time['cost_per_hour'] / avg_cost_per_hour
64
65 kpi_cei = df_time[['company', 'cost_per_hour', 'CEI']].sort_values('CEI', ascending=False)
66 print("KPI 7: Cost Efficiency Index")
67 print(kpi_cei.head(10))
68
69 #KPI 8
70 total_hours = df_time['total_hours'].sum()
71 pareto_hours_pct = (df_time.groupby('category')['total_hours'].sum() / total_hours) * 100
72
73 print("\nKPI 8: Pareto Distribution")
74 print(pareto_hours_pct.sort_values(ascending=False))
75
76 #KPI 9
77
78 total_revenue = df_time['estimated_cost'].sum()
79 category_revenue_pct = (df_time.groupby('category')['estimated_cost'].sum() / total_revenue) * 100

```

```

70 category_revenue_pct = df_time.groupby('category')['estimated_cost'].sum() / total_revenue
80
81 print("\nKPI 9: Revenue Distribution")
82 print(category_revenue_pct.sort_values(ascending=False))
83
84 #kpi10
85 df_time['efficiency'] = df_time['estimated_cost'] / df_time['total_hours']
86
87 high_hours = df_time['total_hours'].quantile(0.75)
88 low_eff = df_time['efficiency'].quantile(0.25)
89
90 time_draining_accounts = df_time[
91     (df_time['total_hours'] >= high_hours) &
92     (df_time['efficiency'] <= low_eff)
93 ][['company', 'total_hours', 'estimated_cost', 'efficiency']]
94
95 print("\nKPI 10: Time Draining Accounts")
96 print(time_draining_accounts)
97
98

```

16	assetz home	7.25	25375
17	brigade and prestige	7.00	24500
18	sriram properties	7.00	24500
19	innotech	6.50	22750
20	symmetric architects	6.50	22750
21	bps	6.25	21875
22	rmz developers	6.00	21000
23	snn estates developers	6.00	21000
24	sattva developers groups	5.75	20125
25	abb india pvt ltd	5.25	18375
26	prestige groups	5.25	18375
27	optimus	5.00	17500
28	dhruva engineering	4.50	15750
29	shapoorji and pallonji constructions	4.50	15750
30	sundaram architects	4.50	15750
31	dsp designs	4.00	14000
32	hubspot discussion with jayanth	4.00	14000
33	inform architects	4.00	14000
34	magcor pvt ltd	4.00	14000
35	somerset pharma	4.00	14000

efficiency	
0	3500.0
1	3500.0
2	3500.0
3	3500.0
4	3500.0
5	3500.0
6	3500.0
7	3500.0
8	3500.0
9	3500.0
10	3500.0
11	3500.0
12	3500.0
13	3500.0
14	3500.0
15	3500.0
16	3500.0
17	3500.0
18	3500.0
19	3500.0
20	3500.0
21	3500.0
22	3500.0
23	3500.0
24	3500.0
25	3500.0
26	3500.0
27	3500.0
28	3500.0
29	3500.0
30	3500.0
31	3500.0
32	3500.0
33	3500.0
34	3500.0
35	3500.0

```

1 # =====
2 # SHEET 2 – KEY ACCOUNTS KPI ANALYSIS
3 # =====
4
5 # Ensure 'num_decision_makers' is numeric before calculations
6 df_key['num_decision_makers'] = df_key['num_decision_makers'].apply(normalize_to_lower)
7 df_key['num_decision_makers'] = pd.to_numeric(df_key['num_decision_makers'], errors='coerce').fillna(0).astype(int)
8

```

```

9 # KPI 1 – Estimated BOQ per account (highest potential accounts)
10 kpi_estimated_boq = (
11     df_key.groupby('company_name')['estimated_boq_int']
12     .sum()
13     .sort_values(ascending=False)
14 )
15 print("\nKPI 1 – Estimated BOQ per account:")
16 print(kpi_estimated_boq.head(10))
17
18
19 # KPI 2 – Target BOQ per account (highest assigned goals)
20 kpi_target_boq = (
21     df_key.groupby('company_name')['target_boq']
22     .sum()
23     .sort_values(ascending=False)
24 )
25 print("\nKPI 2 – Target BOQ per account:")
26 print(kpi_target_boq.head(10))
27
28
29 # KPI 3 – BOQ Gap (Estimated - Target)
30 df_key['gap_boq'] = df_key['estimated_boq_int'] - df_key['target_boq']
31
32 kpi_gap_boq = (
33     df_key.groupby('company_name')['gap_boq']
34     .mean()
35     .sort_values(ascending=False)
36 )
37 print("\nKPI 3 – BOQ Gap (Estimated - Target):")
38 print(kpi_gap_boq.head(10))
39
40
41 # KPI 4 – Complexity Score (Estimated BOQ × # decision makers)
42 df_key['opportunity_index'] = (
43     df_key['estimated_boq_int'] / df_key['num_decision_makers']
44 )
45 df_key['complexity'] = df_key['estimated_boq_int'] * df_key['num_decision_makers']
46 df_key['opportunity_index'].replace([np.inf, -np.inf], np.nan, inplace=True)
47
48 kpi_complexity = (
49     df_key.groupby('company_name')['complexity']
50     .sum()
51     .sort_values(ascending=False)
52 )
53 print("\nKPI 4 – Complexity Score:")
54 print(kpi_complexity.head(10))
55
56
57 # KPI 5 – Opportunity Index (Estimated BOQ / # decision makers)
58
59
60 kpi_opportunity = (
61     df_key.groupby('company_name')['opportunity_index']
62     .mean()
63     .sort_values(ascending=False)
64 )
65 print("\nKPI 5 – Opportunity Index:")
66 print(kpi_opportunity.head(10))
67
68
69 # KPI 6 – BOQ by project type
70 kpi_project_type = (
71     df_key.groupby('types_of_projects')['estimated_boq_int']
72     .sum()
73     .sort_values(ascending=False)
74 )
75 print("\nKPI 6 – BOQ by Project Type:")
76 print(kpi_project_type)
77
78
79 # KPI 7 – Tier Classification (High / Medium / Low potential accounts)
80 boq = df_key.groupby('company_name')['estimated_boq_int'].sum()
81
82 high_potential = boq.quantile(0.75)
83 low_potential = boq.quantile(0.25)
84
85 df_segments = boq.apply(
86     lambda x: 'High' if x >= high_potential else
87             ('Low' if x <= low_potential else 'Medium')
88 )
89

```



```
90 print("\nKPI 7 – Account Potential Tiers:")
91 print(df_segments.head(20))
```

KPI 1 – Estimated BOQ per account:

```
company_name
Design Tree          15
CRN Architects        10
Innotech             10
Ashok Narayan & Eshwar 5
Bureau Engineers (structural) 5
AECOM                5
RSB Consultants       5
Natraj And Venkat.    5
Archetype            3
Jayam Consultants     3
Name: estimated_boq_int, dtype: int64
```

KPI 2 – Target BOQ per account:

```
company_name
Design Tree          13
Innotech             9
CRN Architects        8
Ashok Narayan & Eshwar 5
Bureau Engineers (structural) 5
AECOM                5
RSB Consultants       5
Jayam Consultants     3
Archetype            3
Natraj And Venkat.    3
Name: target_boq, dtype: Int64
```

KPI 3 – BOQ Gap (Estimated - Target):

```
company_name
CRN Architects        2.0
Design Tree          2.0
Natraj And Venkat.    2.0
Innotech             1.0
Ducon                1.0
Geostructural         1.0
Chandavarkar and Thacker Architects 1.0
Gensler.             1.0
Murthy and Manyam.   1.0
Nadig consultants     1.0
Name: gap_boq, dtype: Float64
```

KPI 4 – Complexity Score:

```
company_name
Design Tree          1050
Bureau Engineers (structural) 80
AECOM                75
Ashok Narayan & Eshwar 60
RSB Consultants       40
Dius Consultants      36
Innotech             30
CRN Architects        30
cheralathan Architect 30
Archetype            21
Name: complexity, dtype: int64
```

""" KPI 5 – Total hours per company """

```
1 # =====
2 # SHEET 3 – MEETING LOG KPIs
3
4 # safety checks
5 required = ['Date of Activity', 'Activity Type', 'Company', 'Project', 'Hours spent']
6 missing = [c for c in required if c not in df_meet.columns]
7 if missing:
8     raise KeyError(f"Missing columns in df_meet: {missing}")
9
10 # DATE CHECKING
11 df_meet['Date of Activity'] = pd.to_datetime(df_meet['Date of Activity'], errors='coerce')
12
13 # KPI M1: Total hours per company
14 hours_by_company = df_meet.groupby('Company')['Hours spent'].sum().sort_values(ascending=False)
15 print("\nKPI M1 – Total hours by company (top 10):")
16 print(hours_by_company.head(10))
17
18 # KPI M2: Total hours per project
19 hours_by_project = df_meet.groupby('Project')['Hours spent'].sum().sort_values(ascending=False)
20 print("\nKPI M2 – Total hours by project (top 10):")
21 print(hours_by_project.head(10))
22
23 # KPI M3: Total hours per activity type
24 hours_by_activity = df_meet.groupby('Activity Type')['Hours spent'].sum().sort_values(ascending=False)
25 print("\nKPI M3 – Total hours by activity type:")
```

```

26 print(hours_by_activity)
27
28 # KPI M4: Number of meetings per company
29 meetings_by_company = df_meet.groupby('Company').size().sort_values(ascending=False)
30 print("\nKPI M4 - Number of meetings per company (top 10):")
31 print(meetings_by_company.head(10))
32
33 # KPI M5: Number of meetings per project
34 meetings_by_project = df_meet.groupby('Project').size().sort_values(ascending=False)
35 print("\nKPI M5 - Number of meetings per project (top 10):")
36 print(meetings_by_project.head(10))
37
38 # KPI A (advanced): Average hours per meeting type
39 avg_hours_by_activity = df_meet.groupby('Activity Type')['Hours spent'].mean().sort_values(ascending=False)
40 print("\nKPI A - Average hours per activity type:")
41 print(avg_hours_by_activity)
42
43 # KPI B (advanced): Meeting intensity per company (hours / meetings)
44 hours = df_meet.groupby('Company')['Hours spent'].sum()
45 counts = df_meet.groupby('Company').size()
46 meeting_intensity = (hours / counts).sort_values(ascending=False)
47 print("\nKPI B - Meeting intensity (avg hours per meeting) per company (top 10):")
48 print(meeting_intensity.head(10))
49
50 # KPI C: Meetings & hours by month (trend)
51 df_meet['month'] = df_meet['Date of Activity'].dt.to_period('M')
52 meetings_per_month = df_meet.groupby('month').size().sort_index()
53 hours_per_month = df_meet.groupby('month')['Hours spent'].sum().sort_index()
54 print("\nKPI C - Meetings per month (chronological):")
55 print(meetings_per_month)
56 print("\nKPI C - Hours per month (chronological):")
57 print(hours_per_month)
58
59 # End of Sheet 3 KPIs
60

```

KPI M1 - Total hours by company (top 10):

Company	
Chennai	32.00
Pidilite	22.75
Bagmane Developers	18.25
Kalpataru projects International Limited	17.75
Hyderabad Travel	16.00
Design Tree	13.25
CRN Architects	12.75
Birla Groups	12.00
Innotech	10.75
Colliers	10.40

Name: Hours spent, dtype: float64

KPI M2 - Total hours by project (top 10):

Project	
Multiple projects	45.25
Multiple Project	32.00
Multiple Projects	18.75
Event	16.00
Chennai	16.00
Pidilite	15.00
Birla Tisya	14.00
Sap Labs	12.00
MRM	11.00
Somerset pharma	8.50

Name: Hours spent, dtype: float64

KPI M3 - Total hours by activity type:

Activity Type	
Unknown	448.65
Introduction pitch	57.00
Techno commercial	42.00
Site Visit	32.25
Techno - Commercial	22.75
Internal team meeting	21.75
Techno - commercial	21.75
Site Visit	11.00
Introduction pitch	10.00
Introduction Pitch	8.50
Technical presentation	8.00
Admin	2.25
Applicator Meeting	2.00
Group Presntation	2.00
Applicator Techno Commercial	1.50
Internal meeting	1.25
Site visit	1.00
Activity Type	0.00
Leave	0.00
Festival - Diwali	0.00

Name: Hours spent, dtype: float64

KPI M4 – Number of meetings per company (top 10):

Company	
Kalpataru projects International Limited	12
Bagmane Developers	11
Birla Groups	10

```

1 # =====
2 # SHEET 4 – INQUIRY TRACKER KPIS
3
4
5 # safety checks
6 required = ['Inquiry Date', 'Projects', 'Stage', 'Value (INR)']
7 missing = [c for c in required if c not in df_inquiry.columns]
8 if missing:
9     raise KeyError(f"Missing required columns in df_inquiry: {missing}")
10
11 # Ensure dates and numeric types
12 df_inquiry['Inquiry Date'] = pd.to_datetime(df_inquiry['Inquiry Date'], errors='coerce')
13 if 'expected_close_clean' not in df_inquiry.columns:
14     if 'Expected Close date' in df_inquiry.columns:
15         df_inquiry['expected_close_clean'] = pd.to_datetime(df_inquiry['Expected Close date'], errors='coerce')
16     else:
17         df_inquiry['expected_close_clean'] = pd.NaT
18
19 df_inquiry['Value (INR)'] = pd.to_numeric(df_inquiry['Value (INR)'], errors='coerce')
20
21 # KPI I1: Inquiry count per Stage
22 inquiries_by_stage = df_inquiry['Stage'].value_counts()
23 inquiries_by_stage_pct = df_inquiry['Stage'].value_counts(normalize=True) * 100
24 print("\nKPI I1 – Inquiry count by Stage:")
25 print(inquiries_by_stage)
26 print("\nKPI I1 – Inquiry distribution (%) by Stage:")
27 print(inquiries_by_stage_pct)
28
29 # KPI I2: Inquiry count per Construction type
30 if 'Construction type' in df_inquiry.columns:
31     inquiries_by_construction = df_inquiry['Construction type'].value_counts()
32     print("\nKPI I2 – Inquiry count by Construction type:")
33     print(inquiries_by_construction)
34 else:
35     print("\nKPI I2 – 'Construction type' column not found.")
36
37 # KPI I3: Total Value (INR) per Stage
38 value_by_stage = df_inquiry.groupby('Stage')['Value (INR)'].sum().sort_values(ascending=False)
39 print("\nKPI I3 – Total Value (INR) by Stage:")
40 print(value_by_stage)
41
42 # KPI I4: Inquiry volume per month & KPI I5: Inquiry value per month
43 df_inquiry['month'] = df_inquiry['Inquiry Date'].dt.to_period('M')
44 inquiries_per_month = df_inquiry.groupby('month').size().sort_index()
45 value_per_month = df_inquiry.groupby('month')['Value (INR)'].sum().sort_index()
46 print("\nKPI I4 – Inquiries per month (chronological):")
47 print(inquiries_per_month)
48 print("\nKPI I5 – Value (INR) per month (chronological):")
49 print(value_per_month)
50
51 # KPI I6: Stage conversion - counts and percent
52 stage_counts = df_inquiry['Stage'].value_counts()
53 stage_percent = df_inquiry['Stage'].value_counts(normalize=True) * 100
54 print("\nKPI I6 – Stage funnel counts:")
55 print(stage_counts)
56 print("\nKPI I6 – Stage funnel %:")
57 print(stage_percent)
58
59 # KPI I7: Mean lead duration (expected_close_clean - Inquiry Date)
60 df_inquiry['lead_duration_days'] = (df_inquiry['expected_close_clean'] - df_inquiry['Inquiry Date']).dt.days
61 mean_lead_duration = df_inquiry['lead_duration_days'].mean()
62 print("\nKPI I7 – Mean lead duration (days):")
63 print(mean_lead_duration)
64
65 # KPI I8: High-value inquiries stuck in early stages (top 25% by Value)
66 early_stages = ['New', 'Feasibility', 'Commercial Review'] # adjust if needed
67 value_series = df_inquiry['Value (INR)'].dropna()
68 if len(value_series)>0:
69     high_value_threshold = value_series.quantile(0.75)
70 else:
71     high_value_threshold = np.nan
72
73 mask_stage = df_inquiry['Stage'].isin(early_stages)
74 mask_value = df_inquiry['Value (INR)'] >= high_value_threshold
75 df_high_stuck = df_inquiry[mask_stage & mask_value].copy()

```

```

76
77 print("\nKPI I8 – High-value inquiries stuck in early stages (summary):")
78 print("High-value threshold (75th percentile):", high_value_threshold)
79 print("Count high-value & early-stage:", len(df_high_stuck))
80 print("Total value stuck (INR):", df_high_stuck['Value (INR)'].sum())
81
82 print("\nTop high-value stuck rows (if any):")
83 show_cols = [c for c in ['Inquiry Date', 'Projects', 'project_norm', 'Stage', 'Value (INR)', 'Number of Hours', 'expected_close_date']]
84 print(df_high_stuck.sort_values('Value (INR)', ascending=False)[show_cols].head(20))
85
86 # KPI I10: Overdue high-value inquiries (>60 days)
87 if 'lead_duration_days' in df_inquiry.columns:
88     df_overdue = df_inquiry[(df_inquiry['lead_duration_days'] > 60) & (df_inquiry['Value (INR)'] >= high_value_threshold)]
89     print("\nKPI I10 – Overdue (>60 days) high-value inquiries:")
90     print("Count:", len(df_overdue))
91     print("Total overdue value (INR):", df_overdue['Value (INR)'].sum())
92     show_cols2 = [c for c in ['Inquiry Date', 'Projects', 'Stage', 'Value (INR)', 'lead_duration_days', 'expected_close_date']]
93     print(df_overdue.sort_values('lead_duration_days', ascending=False)[show_cols2].head(20))
94 else:
95     print("\nKPI I10 – lead_duration_days not available; compute KPI I7 first.")
96
97 # End of Sheet 4 KPIs
98

```

KPI I1 – Inquiry count by Stage:

Stage	
Commercial Review	20
Consultant Review	1
Deal won	1
Technical Review	1
Tender Awarded	1

Name: count, dtype: int64

KPI I1 – Inquiry distribution (%) by Stage:

Stage	
Commercial Review	83.333333
Consultant Review	4.166667
Deal won	4.166667
Technical Review	4.166667
Tender Awarded	4.166667

Name: proportion, dtype: float64

KPI I2 – Inquiry count by Construction type:

Construction type	
New	24

Name: count, dtype: int64

KPI I3 – Total Value (INR) by Stage:

Stage	
Commercial Review	132345000.0
Consultant Review	10200000.0
Deal won	8500000.0
Technical Review	3995000.0
Tender Awarded	3400000.0

Name: Value (INR), dtype: float64

KPI I4 – Inquiries per month (chronological):

month	
2024-11	2
2024-12	2
2025-01	2
2025-03	1
2025-04	2
2025-05	3
2025-06	2
2025-07	1
2025-08	2
2025-09	2
2025-10	3

Freq: M, dtype: int64

KPI I5 – Value (INR) per month (chronological):

month	
2024-11	12325000.0
2024-12	15895000.0
2025-01	21930000.0
2025-03	10200000.0
2025-04	8585000.0
2025-05	23630000.0
2025-06	10795000.0
2025-07	8500000.0

▼ Data Problem

```

1 # assume existing cleaned DataFrames from your notebook:
2 # df_time (has 'company'), df_key (has 'account_norm'), df_meet (has 'company_norm','project_norm'), df_inquiry (has 'project_norm')
3
4 # Basic sizes
5 print("df_time rows:", len(df_time))
6 print("df_key rows:", len(df_key))
7 print("df_meet rows:", len(df_meet))
8 print("df_inquiry rows:", len(df_inquiry))
9
10 # Exact overlap – company-level between df_time and df_key
11 time_vs_key = pd.merge(
12     df_time[['company']].drop_duplicates().rename(columns={'company':'key'}),
13     df_key[['account_norm']].drop_duplicates().rename(columns={'account_norm':'key'}),
14     on='key',
15     how='outer',
16     indicator=True
17 )
18 print("Company join counts (df_time <-> df_key):")
19 print(time_vs_key['_merge'].value_counts())
20
21 # Exact overlap – company-level between df_time and df_meet
22 time_vs_meet = pd.merge(
23     df_time[['company']].drop_duplicates().rename(columns={'company':'key'}),
24     df_meet[['company_norm']].drop_duplicates().rename(columns={'company_norm':'key'}),
25     on='key',
26     how='outer',
27     indicator=True
28 )
29 print("Company join counts (df_time <-> df_meet):")
30 print(time_vs_meet['_merge'].value_counts())
31
32 # Exact overlap – project-level between df_meet and df_inquiry
33 proj_vs_inq = pd.merge(
34     df_meet[['project_norm']].drop_duplicates().rename(columns={'project_norm':'key'}),
35     df_inquiry[['project_norm']].drop_duplicates().rename(columns={'project_norm':'key'}),
36     on='key',
37     how='outer',
38     indicator=True
39 )
40 print("Project join counts (df_meet <-> df_inquiry):")
41 print(proj_vs_inq['_merge'].value_counts())
42

```

```

df_time rows: 142
df_key rows: 42
df_meet rows: 450
df_inquiry rows: 24
Company join counts (df_time <-> df_key):
_merge
left_only    131
right_only   34
both         8
Name: count, dtype: int64
Company join counts (df_time <-> df_meet):
_merge
both         106
right_only   101
left_only    33
Name: count, dtype: int64
Project join counts (df_meet <-> df_inquiry):
_merge
left_only    254
right_only   19
both         5
Name: count, dtype: int64

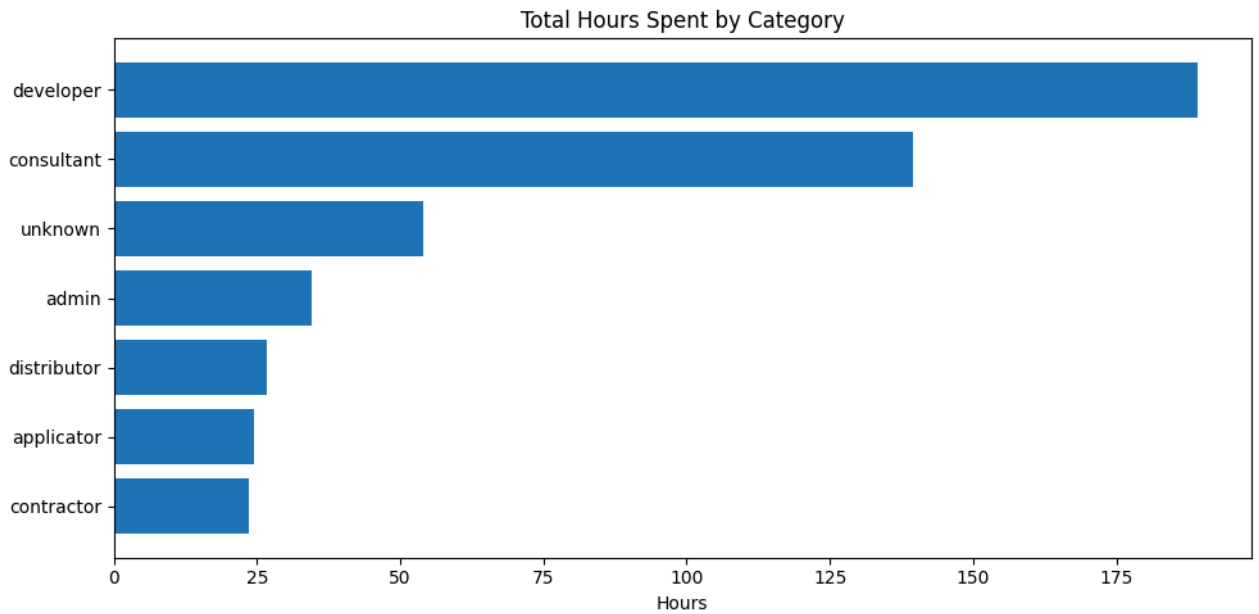
```

▼ Data Visualization

```

1 import matplotlib.pyplot as plt
2
3 category_hours = df_time.groupby('category')['total_hours'].sum().sort_values()
4
5 plt.figure(figsize=(10,5))
6 plt.barh(category_hours.index, category_hours.values)
7 plt.title("Total Hours Spent by Category")
8 plt.xlabel("Hours")
9 plt.tight_layout()
10 plt.show()
11

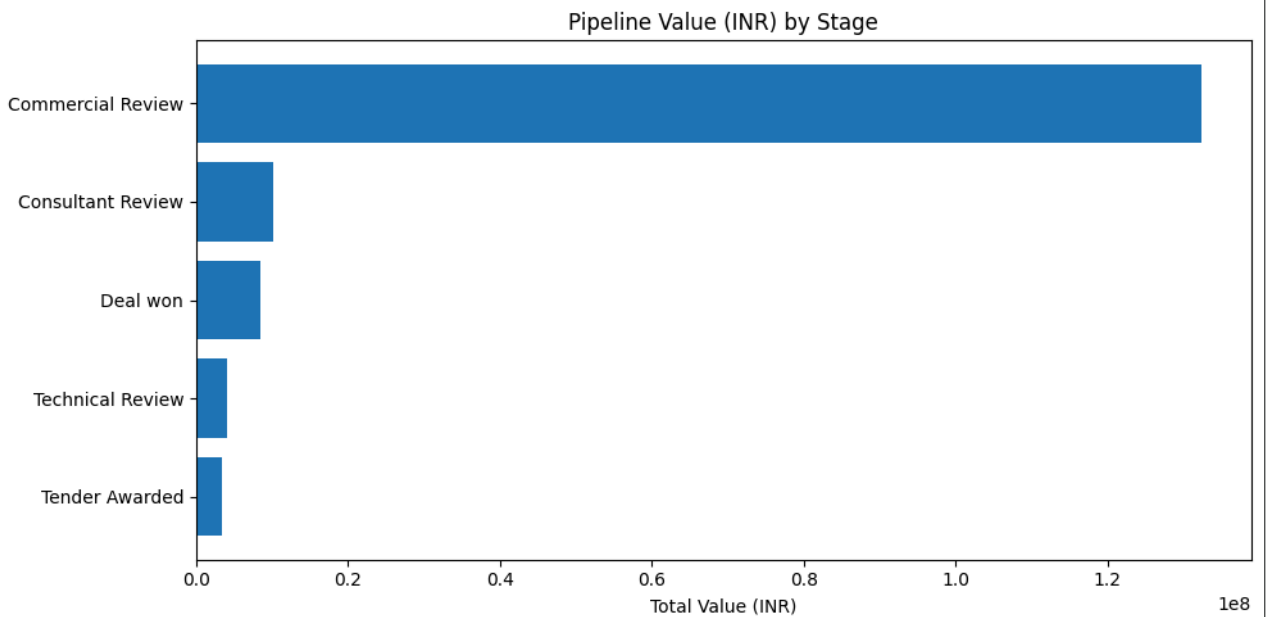
```



```

1 value_by_stage = df_inquiry.groupby('Stage')['Value (INR)'].sum().sort_values()
2
3 plt.figure(figsize=(10,5))
4 plt.barh(value_by_stage.index, value_by_stage.values)
5 plt.title("Pipeline Value (INR) by Stage")
6 plt.xlabel("Total Value (INR)")
7 plt.tight_layout()
8 plt.show()
9

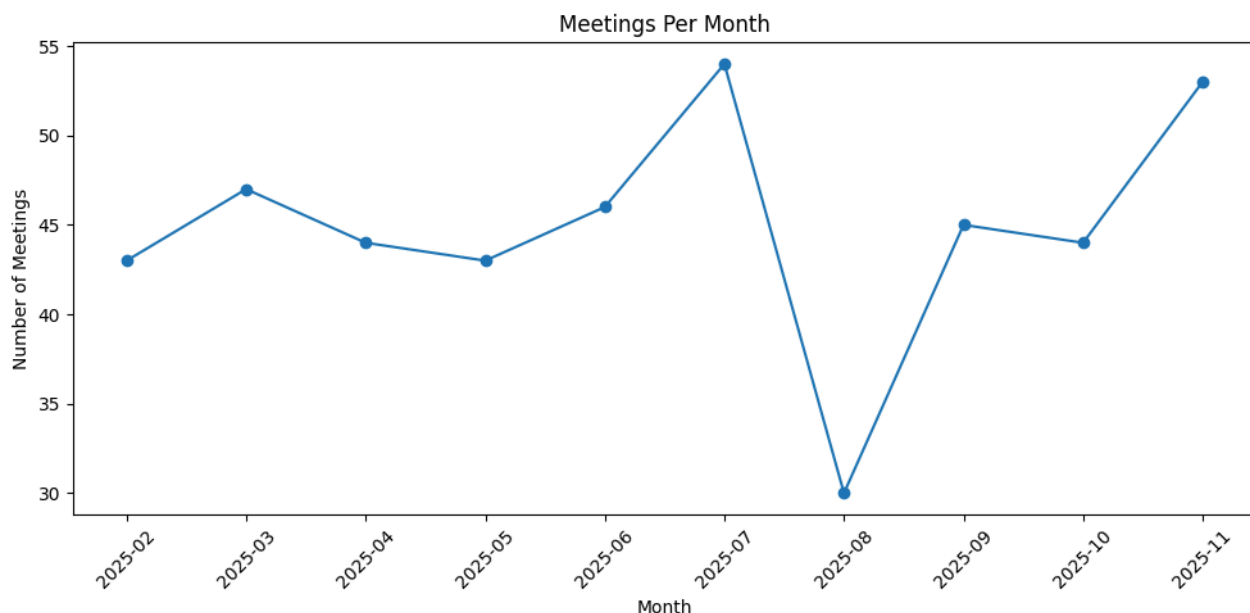
```



```

1 df_meet['month'] = df_meet['Date of Activity'].dt.to_period('M')
2 meetings_per_month = df_meet.groupby('month').size()
3
4 plt.figure(figsize=(10,5))
5 plt.plot(meetings_per_month.index.astype(str), meetings_per_month.values, marker='o')
6 plt.title("Meetings Per Month")
7 plt.xlabel("Month")
8 plt.ylabel("Number of Meetings")
9 plt.xticks(rotation=45)
10 plt.tight_layout()
11 plt.show()
12

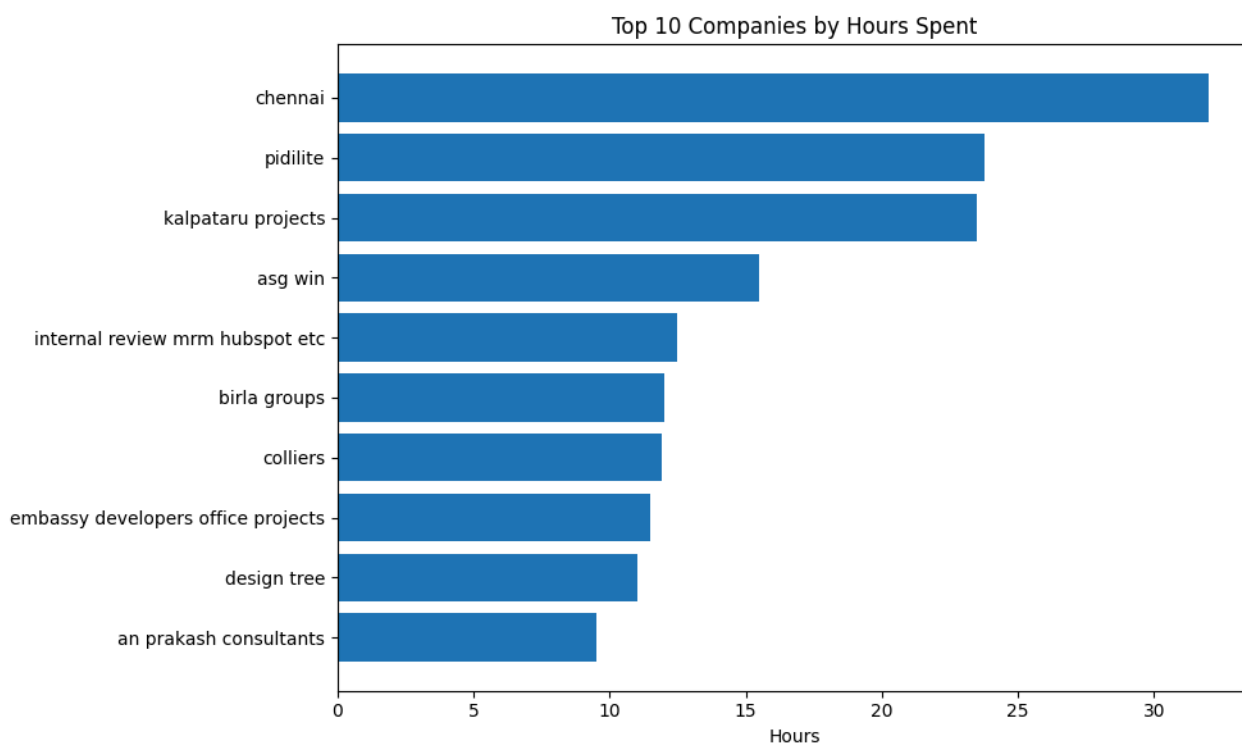
```



```

1 top_companies = df_time.groupby('company')['total_hours'].sum().nlargest(10)
2
3 plt.figure(figsize=(10,6))
4 plt.barh(top_companies.index[::-1], top_companies.values[::-1])
5 plt.title("Top 10 Companies by Hours Spent")
6 plt.xlabel("Hours")
7 plt.tight_layout()
8 plt.show()
9

```



```

1 plt.figure(figsize=(10,5))
2 plt.hist(df_inquiry['lead_duration_days'].dropna(), bins=15)
3 plt.title("Distribution of Deal Aging (Days)")
4 plt.xlabel("Days in Pipeline")
5 plt.ylabel("Number of Deals")
6 plt.tight_layout()
7 plt.show()
8

```

