

# Data Cleaning with Pandas

Cell 1: Import Required Libraries

```
import pandas as pd  
import numpy as np
```

Cell 2: Load the CSV File

```
df = pd.read_csv("/content/Messy_Employment_India_Dataset.csv")
```

Cell 3: View the Dataset

```
df.head(110)
```

```
{"summary": {"name": "df", "rows": 2000, "fields": [{}], "column": "Status", "properties": {"dtype": "category", "num_unique_values": 6, "samples": [{"EMPLOYED": "UNEMPLOYED"}, {"Employed": "Unemployed"}], "semantic_type": "\\", "description": "\n"}, "column": "Age Group", "properties": {"dtype": "category", "num_unique_values": 7, "samples": [{"25_34": "18-24"}, {"25-34": "25-34"}], "semantic_type": "\\", "description": "\n"}, "column": "Education", "properties": {"dtype": "category", "num_unique_values": 9, "samples": [{"Bachelors": "Diploma"}, {"Bachelors": "Master"}, {"Diploma": "Master"}], "semantic_type": "\\", "description": "\n"}, "column": "Industry", "properties": {"dtype": "category", "num_unique_values": 8, "samples": [{"Fintech": "Healthcare"}, {"Technology": "Healthcare"}, {"Fintech": "Technology"}, {"Technology": "Fintech"}], "semantic_type": "\\", "description": "\n"}, "column": "Location", "properties": {"dtype": "category", "num_unique_values": 8, "samples": [{"rural": "Mumbai"}, {"rural": "Urban"}, {"Mumbai": "Urban"}], "semantic_type": "\\", "description": "\n"}, "column": "AI Risk", "properties": {"dtype": "category", "num_unique_values": 6, "samples": [{"moderate": "Low"}, {"moderate": "High"}], "semantic_type": "\\", "description": "\n"}, "column": "Years of Experience", "properties": {"dtype": "category", "num_unique_values": 11, "samples": [{"moderate": "Low"}, {"moderate": "High"}, {"moderate": "Very Low"}, {"moderate": "Very High"}, {"moderate": "Extremely Low"}, {"moderate": "Extremely High"}, {"moderate": "Low"}, {"moderate": "High"}, {"moderate": "Very Low"}, {"moderate": "Very High"}, {"moderate": "Extremely Low"}, {"moderate": "Extremely High"}], "semantic_type": "\\", "description": "\n"}]}
```

```

  "properties": {
    "Status": {
      "dtype": "number",
      "min": 0.0,
      "max": 30.0,
      "num_unique_values": 31,
      "samples": [
        19.0,
        8.0,
        27.0
      ],
      "semantic_type": "\",
      "description": """
        },
      "column": "Monthly Salary (INR)",
      "properties": {
        "dtype": "number",
        "std": 41628.00805409873,
        "min": 5100.0,
        "max": 149900.0,
        "num_unique_values": 978,
        "samples": [
          126900.0,
          117100.0,
          31500.0
        ],
        "semantic_type": "\",
        "description": """
          },
        "column": "Date Recorded",
        "properties": {
          "dtype": "object",
          "num_unique_values": 2000,
          "samples": [
            "2/4/2028",
            "12/20/2023",
            "8/26/2026"
          ],
          "semantic_type": "\",
          "description": """
            }
          }
        }
      }
    }
  }
}

```

#### Cell 4: Understand Dataset Structure

```

df.shape      # number of rows and columns
(2000, 9)

df.info()     # data types and missing values

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2000 entries, 0 to 1999
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Status            1732 non-null    object 
 1   Age Group         1768 non-null    object 
 2   Education         1804 non-null    object 
 3   Industry          1799 non-null    object 
 4   Location          1787 non-null    object 
 5   AI Risk           1716 non-null    object 
 6   Years of Experience 980 non-null    float64
 7   Monthly Salary (INR) 1613 non-null    float64
 8   Date Recorded    2000 non-null    object 
dtypes: float64(2), object(7)
memory usage: 140.8+ KB

df.describe() # statistical summary

{
  "summary": {
    "name": "df",
    "rows": 8,
    "fields": [
      {
        "column": "Years of Experience",
        "properties": {
          "dtype": "number",
          "std": 341.5416910202622,
          "min": 0.0,
          "max": 980.0,
          "num_unique_values": 8,
          "samples": [

```

```

15.244897959183673,\n          16.0,\n          980.0\n      ],\n      {\\"semantic_type\\": \"\",\\n      \\"description\\": \"\\n      }\n    },\\n    {\\"column\\": \"Monthly Salary (INR)\\\",\\n      \\"properties\\\": {\\"n      \\"dtype\\": \"number\\\",\\n      \\"std\\\":\n51354.52706818114,\\n      \\"min\\\": 1613.0,\\n      \\"max\\\":\n149900.0,\\n      \\"num_unique_values\\\": 8,\\n      \\"samples\\\": [\n76886.3608183509,\\n          77200.0,\\n          1613.0\\n      ],\\n      {\\"semantic_type\\\": \"\",\\n      \\"description\\\": \"\\n      }\n    }\\n  }\\n}","type":"dataframe"}]
```

## Cell 5: Check Missing Values

```

df.isnull()

{"summary": {"name": "df", "rows": 2000, "fields": [
  {"column": "Status", "properties": {"dtype": "boolean", "num_unique_values": 2, "samples": [true, false], "semantic_type": "", "description": ""}, "column": "Age Group", "properties": {"dtype": "boolean", "num_unique_values": 2, "samples": [true, false], "semantic_type": "", "description": ""}, "column": "Education", "properties": {"dtype": "boolean", "num_unique_values": 2, "samples": [true, false], "semantic_type": "", "description": ""}, "column": "Industry", "properties": {"dtype": "boolean", "num_unique_values": 2, "samples": [true, false], "semantic_type": "", "description": ""}, "column": "Location", "properties": {"dtype": "boolean", "num_unique_values": 2, "samples": [true, false], "semantic_type": "", "description": ""}, "column": "AI Risk", "properties": {"dtype": "boolean", "num_unique_values": 2, "samples": [false, true], "semantic_type": "", "description": ""}, "column": "Years of Experience", "properties": {"dtype": "boolean", "num_unique_values": 2, "samples": [false, true], "semantic_type": "", "description": ""}, "column": "Monthly Salary (INR)", "properties": {"dtype": "boolean", "num_unique_values": 2, "samples": [true, false], "semantic_type": "", "description": ""}}]}
```

```

  "description": """",
  "Date Recorded": {
    "properties": {
      "num_unique_values": 1,
      "samples": [
        false
      ],
      "semantic_type": ""
    }
  }
}, "type": "dataframe"
}

df.isnull().sum()

Status          268
Age Group       232
Education        196
Industry         201
Location          213
AI Risk           284
Years of Experience 1020
Monthly Salary (INR) 387
Date Recorded        0
dtype: int64

```

## Cell 6: Handle Missing Values

Fill missing numerical values with mean

```

df['Years of Experience'] = df['Years of Experience'].fillna(df['Years of Experience'].mean())

df

{
  "summary": {
    "name": "df",
    "rows": 2000,
    "fields": [
      {
        "column": "Status",
        "properties": {
          "dtype": "category",
          "num_unique_values": 6,
          "samples": [
            "EMPLOYED",
            "UNEMPLOYED",
            "Employed"
          ],
          "semantic_type": ""
        },
        "description": ""
      },
      {
        "column": "Age Group",
        "properties": {
          "dtype": "category",
          "num_unique_values": 7,
          "samples": [
            "25_34",
            "18-24",
            "25-34"
          ],
          "semantic_type": ""
        },
        "description": ""
      },
      {
        "column": "Education",
        "properties": {
          "dtype": "category",
          "num_unique_values": 9,
          "samples": [
            "Bachelor's",
            "Diploma",
            "Master"
          ],
          "semantic_type": ""
        },
        "description": ""
      },
      {
        "column": "Industry",
        "properties": {
          "dtype": "category",
          "num_unique_values": 8,
          "samples": [
            "Fintech",
            "Healthcare",
            "Technology"
          ],
          "semantic_type": ""
        },
        "description": ""
      },
      {
        "column": "Location",
        "properties": {
          "dtype": ""
        },
        "description": ""
      }
    ]
  }
}

```

```

\"category\", \n      \"num_unique_values\": 8, \n      \"samples\": \n      [\n        {\"rural\": \"Mumbai\", \"Urban\": \"\n      ], \n      \"semantic_type\": \"\", \n      \"description\": \"\"\n    }, \n    { \n      \"column\": \"AI Risk\", \n      \"properties\": { \n        \"dtype\": \"category\", \n        \"num_unique_values\": 6, \n        \"samples\": [ \n          \"moderate\", \n          \"Low\", \n          \"High\", \n          \"High\", \n          \"High\", \n          \"High\" \n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      }, \n      \"column\": \"Years of Experience\", \n      \"properties\": { \n        \"dtype\": \"number\", \n        \"std\": 6.371099410521485, \n        \"min\": 0.0, \n        \"max\": 30.0, \n        \"num_unique_values\": 32, \n        \"samples\": [ \n          29.0, \n          27.0, \n          27.0 \n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      }, \n      \"column\": \"Monthly Salary (INR)\", \n      \"properties\": { \n        \"dtype\": \"number\", \n        \"std\": 41628.00805409873, \n        \"min\": 5100.0, \n        \"max\": 149900.0, \n        \"num_unique_values\": 978, \n        \"samples\": [ \n          126900.0, \n          117100.0, \n          31500.0 \n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      }, \n      \"column\": \"Date Recorded\", \n      \"properties\": { \n        \"dtype\": \"object\", \n        \"num_unique_values\": 2000, \n        \"samples\": [ \n          \"2/4/2028\", \n          \"12/20/2023\", \n          \"8/26/2026\" \n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\"\n      } \n    } \n  }, \n  {\"type\": \"dataframe\", \"variable_name\": \"df\"}

```

Fill missing categorical values

```

df['Location'] = df['Location'].fillna("Unknown")

df

{ "summary": { \n    \"name\": \"df\", \n    \"rows\": 2000, \n    \"fields\": [\n      { \n        \"column\": \"Status\", \n        \"properties\": { \n          \"dtype\": \"category\", \n          \"num_unique_values\": 6, \n          \"samples\": [ \n            \"EMPLOYED\", \n            \"UNEMPLOYED\", \n            \"Employed\", \n            \"Unemployed\", \n            \"Unknown\" \n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\"\n        }, \n        \"column\": \"Age Group\", \n        \"properties\": { \n          \"dtype\": \"category\", \n          \"num_unique_values\": 7, \n          \"samples\": [ \n            \"25_34\", \n            \"18-24\", \n            \"25-34\", \n            \"35-44\", \n            \"45-54\", \n            \"55-64\", \n            \"65+\" \n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\"\n        }, \n        \"column\": \"Education\", \n        \"properties\": { \n          \"dtype\": \"category\", \n          \"num_unique_values\": 9, \n          \"samples\": [ \n            \"Bachelors\", \n            \"Diploma\", \n            \"Master\", \n            \"Postgraduate\", \n            \"Tertiary\", \n            \"Secondary\", \n            \"Primary\", \n            \"Unknown\", \n            \"Other\" \n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\"\n        } \n      } \n    ] \n  }

```

```

"Industry",\n      "properties": {\n          "dtype":\n          "category",\n              "num_unique_values": 8,\n          "samples": [\n                  "Fintech",\n                  "Healthcare",\n                  "Technology"],\n                  "semantic_type": "\",\n                  "description": \"\\n          },\n          {\n              "column":\n                  "Location",\n                      "properties": {\n                          "category",\n                              "num_unique_values": 9,\n                          "samples": [\n                              "urban",\n                              "rural"],\n                              "semantic_type": "\",\n                              "description": \"\\n          },\n          {\n              "column":\n                  "AI Risk",\n                      "properties": {\n                          "dtype": "category",\n                          "num_unique_values": 6,\n                          "samples": [\n                              "moderate",\n                              "Low",\n                              "High"],\n                          "semantic_type": "\",\n                          "description": \"\\n          },\n          {\n              "column":\n                  "Years of Experience",\n                      "properties": {\n                          "dtype": "number",\n                          "std": 6.371099410521485,\n                          "min": 0.0,\n                          "max": 30.0,\n                          "num_unique_values": 32,\n                          "samples": [\n                              29.0,\n                              27.0\n                          ],\n                          "semantic_type": "\",\n                          "description": \"\\n          },\n          {\n              "column":\n                  "Monthly Salary (INR)",\n                      "properties": {\n                          "dtype": "number",\n                          "std": 41628.00805409873,\n                          "min": 5100.0,\n                          "max": 149900.0,\n                          "num_unique_values": 978,\n                          "samples": [\n                              126900.0,\n                              117100.0,\n                              31500.0\n                          ],\n                          "semantic_type": "\",\n                          "description": \"\\n          },\n          {\n              "column":\n                  "Date Recorded",\n                      "properties": {\n                          "dtype": "object",\n                          "num_unique_values": 2000,\n                          "samples": [\n                              "2/4/2028",\n                              "12/20/2023",\n                              "8/26/2026"],\n                          "semantic_type": "\",\n                          "description": \"\\n          }\n          ]\n      },\n      "type": "dataframe",\n      "variable_name": "df"

```

Cell 7: Remove Duplicate Rows

```
df = df.drop_duplicates()
```

Cell 8: Check and Fix Data Types

```
df.dtypes
```

Status	object
Age Group	object
Education	object
Industry	object
Location	object
AI Risk	object
Years of Experience	float64
Monthly Salary (INR)	float64

```
Date Recorded          object
dtype: object

df['Age Group'] = df['Age Group'].astype(int)

-----
-----
ValueError           Traceback (most recent call
last)
/tmp/ipython-input-2511504638.py in <cell line: 0>()
----> 1 df['Age Group'] = df['Age Group'].astype(int)

/usr/local/lib/python3.12/dist-packages/pandas/core/generic.py in
astype(self, dtype, copy, errors)
    6641         else:
    6642             # else, only a single dtype is given
-> 6643             new_data = self._mgr.astype(dtype=dtype,
copy=copy, errors=errors)
    6644             res = self._constructor_from_mgr(new_data,
axes=new_data.axes)
    6645             return res.__finalize__(self, method="astype")

/usr/local/lib/python3.12/dist-packages/pandas/core/internals/managers
.py in astype(self, dtype, copy, errors)
    428         copy = False
    429
--> 430         return self.apply(
    431             "astype",
    432             dtype=dtype,

/usr/local/lib/python3.12/dist-packages/pandas/core/internals/managers
.py in apply(self, f, align_keys, **kwargs)
    361             applied = b.apply(f, **kwargs)
    362         else:
--> 363             applied = getattr(b, f)(**kwargs)
    364             result_blocks = extend_blocks(applied,
result_blocks)
    365

/usr/local/lib/python3.12/dist-packages/pandas/core/internals/blocks.p
y in astype(self, dtype, copy, errors, using_cow, squeeze)
    756             values = values[0, :] # type: ignore[call-
overload]
    757
--> 758             new_values = astype_array_safe(values, dtype,
copy=copy, errors=errors)
    759
    760             new_values = maybe_coerce_values(new_values)

/usr/local/lib/python3.12/dist-packages/pandas/core/dtypes/astype.py
```

```

in astype_array_safe(values, dtype, copy, errors)
  235
  236     try:
--> 237         new_values = astype_array(values, dtype, copy=copy)
  238     except (ValueError, TypeError):
  239         # e.g. _astype_nansafe can fail on object-dtype of
strings

/usr/local/lib/python3.12/dist-packages/pandas/core/dtypes/astype.py
in astype_array(values, dtype, copy)
  180
  181     else:
--> 182         values = _astype_nansafe(values, dtype, copy=copy)
  183
  184     # in pandas we don't store numpy str dtypes, so convert to
object

/usr/local/lib/python3.12/dist-packages/pandas/core/dtypes/astype.py
in _astype_nansafe(arr, dtype, copy, skipna)
  131     if copy or arr.dtype == object or dtype == object:
  132         # Explicit copy, or required since NumPy can't view
from / to object.
--> 133         return arr.astype(dtype, copy=True)
  134
  135     return arr.astype(dtype, copy=copy)

ValueError: invalid literal for int() with base 10: '18-24'

```

## Cell 9: Rename Columns (Optional)

```

df = df.rename(columns={
    'Years of Experience': 'Y0F',
    'Location': 'Loc'
})

df

{
  "summary": {
    "name": "df",
    "rows": 2000,
    "fields": [
      {
        "column": "Status",
        "properties": {
          "dtype": "category",
          "num_unique_values": 6,
          "samples": ["EMPLOYED", "UNEMPLOYED"],
          "semantic_type": "\",
          "description": "\n        }}, {
            "column": "Age Group",
            "properties": {
              "dtype": "category",
              "num_unique_values": 7,
              "samples": ["25_34", "18-24", "25-34"],
              "semantic_type": "\",
              "description": "\n        }}, {
            "column": "Education",
            "properties": {
              "dtype": "category",
              "num_unique_values": 9,
              "samples": [""

```

```

    \"Bachelors\", \n          \"Diploma\", \n          \"Master\" \\"
], \n      \"semantic_type\": \"\", \n      \"column\": \n      \"description\": \"\\n      }\\n      }, \n      {\n      \"column\": \n      \"Industry\", \n      \"properties\": {\n      \"category\", \n      \"num_unique_values\": 8, \n      \"samples\": [\n      \"Fintech\", \n      \"Healthcare\", \n      \"Technology\", \n      ], \n      \"semantic_type\": \"\", \n      \"column\": \n      \"Loc\", \n      \"properties\": {\n      \"category\", \n      \"num_unique_values\": 9, \n      \"samples\": [\n      \"urban\", \n      \"rural\", \n      \"Mumbai\" \n      ], \n      \"semantic_type\": \"\", \n      \"description\": \"\\n      }\\n      }, \n      {\n      \"column\": \n      \"AI Risk\", \n      \"properties\": {\n      \"category\", \n      \"num_unique_values\": 6, \n      \"samples\": [\n      \"moderate\", \n      \"Low\", \n      \"High\" \n      ], \n      \"semantic_type\": \"\", \n      \"description\": \"\\n      }\\n      }, \n      {\n      \"column\": \n      \"YOF\", \n      \"properties\": {\n      \"number\", \n      \"std\": 6.371099410521485, \n      \"min\": 0.0, \n      \"max\": 30.0, \n      \"num_unique_values\": 32, \n      \"samples\": [\n      29.0, \n      2.0, \n      27.0 \n      ], \n      \"semantic_type\": \"\", \n      \"description\": \"\\n      }\\n      }, \n      {\n      \"column\": \n      \"Monthly Salary (INR)\", \n      \"properties\": {\n      \"number\", \n      \"std\": 41628.00805409873, \n      \"min\": 5100.0, \n      \"max\": 149900.0, \n      \"num_unique_values\": 978, \n      \"samples\": [\n      126900.0, \n      117100.0, \n      31500.0 \n      ], \n      \"semantic_type\": \"\", \n      \"description\": \"\\n      }\\n      }, \n      {\n      \"column\": \n      \"Date Recorded\", \n      \"properties\": {\n      \"object\", \n      \"num_unique_values\": 2000, \n      \"samples\": [\n      \"2/4/2028\", \n      \"12/20/2023\", \n      \"8/26/2026\" \n      ], \n      \"semantic_type\": \"\", \n      \"description\": \"\\n      }\\n      } \n    }, \n    \"type\": \"dataframe\", \n    \"variable_name\": \"df\" 

```

Cell 10: Handle Outliers (Simple Filtering)

```

df_sort = df[df['Monthly Salary (INR)'] < 10000]
df_sort.head()

{"repr_error": "0", "type": "dataframe", "variable_name": "df_sort"} 

```

Cell 11: Clean Text Data

```

#Removes extra spaces from the beginning and end of each city name.
df['Education'] = df['Education'].str.strip()
#Converts all city names to lowercase letters.
df['Loc'] = df['Loc'].str.lower() 

```

## Cell 12: Data Validation

```
df = df[df['YOF'] > 0]
#df = df[(df['Age'] > 0) & (df['Age'] <= 120)]
df

{"summary": {"name": "df", "rows": 1965, "fields": [{"column": "Status", "properties": {"dtype": "category", "num_unique_values": 6, "samples": [{"EMPLOYED": "UNEMPLOYED"}, {"Employed": "Unemployed"}], "semantic_type": "\\", "description": "\n"}, {"column": "Age Group", "properties": {"dtype": "category", "num_unique_values": 7, "samples": [{"25_34": "18-24"}, {"25-34": "25-34"}, {"semantic_type": "\\", "description": "\n"}, {"Education": "Diploma"}, {"Bachelor": "Master"}, {"rural": "mumbai"}, {"urban": "urban"}], "semantic_type": "\\", "description": "\n"}, {"column": "Industry", "properties": {"dtype": "category", "num_unique_values": 8, "samples": [{"Fintech": "Healthcare"}, {"Technology": "Healthcare"}, {"Loc": "Properties"}, {"rural": "urban"}, {"semantic_type": "\\", "description": "\n"}, {"AI Risk": "moderate"}, {"Low": "Medium"}, {"number": "std": 6.090880707692265, "min": 1.0, "max": 30.0, "num_unique_values": 31, "samples": [{"YOF": 19.0}, {"YOF": 8.0}, {"YOF": 27.0}], "semantic_type": "\\", "description": "\n"}, {"column": "Monthly Salary (INR)", "properties": {"dtype": "number", "std": 41699.44279962968, "min": 5100.0, "max": 149900.0, "num_unique_values": 965, "samples": [{"YOF": 34600.0}, {"YOF": 116700.0}, {"YOF": 31500.0}], "semantic_type": "\\", "description": "\n"}, {"column": "Date Recorded", "properties": {"dtype": "object", "num_unique_values": 1965, "samples": []}]}]}}, "semantic_type": "\\", "description": "\n"}]
```

```

\"5/16/2028\", \"semantic_type\": \"\", \"description\": \"\"\n],\n}]\n}, \"type\": \"dataframe\", \"variable_name\": \"df\"}

```

## Cell 13: Final Check

```

df.head()

{
  "summary": {
    "name": "df",
    "rows": 1965,
    "fields": [
      {
        "column": "Status",
        "properties": {
          "dtype": "category",
          "num_unique_values": 6,
          "samples": [
            {"EMPLOYED": "UNEMPLOYED", "semantic_type": "", "description": ""},
            {"Employed": "Age Group", "semantic_type": "", "description": ""},
            {"25_34": "18-24", "semantic_type": "", "description": ""},
            {"Education": "Bachelor", "semantic_type": "", "description": ""},
            {"Education": "Diploma", "semantic_type": "", "description": ""},
            {"Education": "Master", "semantic_type": "", "description": ""},
            {"Industry": "Fintech", "semantic_type": "", "description": ""},
            {"Industry": "Healthcare", "semantic_type": "", "description": ""},
            {"Industry": "Technology", "semantic_type": "", "description": ""},
            {"Loc": "mumbai", "semantic_type": "", "description": ""},
            {"Loc": "urban", "semantic_type": "", "description": ""},
            {"Loc": "rural", "semantic_type": "", "description": ""},
            {"Risk": "AI Risk", "semantic_type": "", "description": ""},
            {"Risk": "moderate", "semantic_type": "", "description": ""},
            {"Risk": "Low", "semantic_type": "", "description": ""},
            {"Risk": "Medium", "semantic_type": "", "description": ""},
            {"Year": "Y0F", "properties": {
              "dtype": "number",
              "std": 6.090880707692265,
              "min": 1.0,
              "max": 30.0,
              "num_unique_values": 31,
              "samples": [
                {"Year": 19.0, "Value": 8.0, "Count": 27.0}
              ],
              "semantic_type": "",
              "description": ""
            }},
            {"Salary": "Monthly Salary (INR)", "properties": {
              "dtype": "number",
              "std": 41699.44279962968,
              "min": 5100.0,
              "max": 149900.0,
              "num_unique_values": 965,
              "samples": [
                {"Salary": 34600.0, "Value": 116700.0, "Count": 31500.0}
              ],
              "semantic_type": "",
              "description": ""
            }}
          ]
        }
      }
    ]
  }
}

```

```
    },\n    },\n    {\n        "column": "Date Recorded",\n        "properties": {\n            "dtype": "object",\n            "num_unique_values": 1965,\n            "samples": [\n                "5/16/2028",\n                "10/22/2025",\n                "2/28/2023"\n            ],\n            "semantic_type": "",\n            "description": ""\n        }\n    }]\n},\n{"type": "dataframe", "variable_name": "df"}\n\ndf.info()\n\n<class 'pandas.core.frame.DataFrame'>\nIndex: 1965 entries, 0 to 1999\nData columns (total 9 columns):\n #   Column           Non-Null Count  Dtype \n--- \n 0   Status            1702 non-null    object \n 1   Age Group         1738 non-null    object \n 2   Education         1775 non-null    object \n 3   Industry          1767 non-null    object \n 4   Loc               1965 non-null    object \n 5   AI Risk            1684 non-null    object \n 6   YOF               1965 non-null    float64 \n 7   Monthly Salary (INR) 1582 non-null    float64 \n 8   Date Recorded     1965 non-null    object \n dtypes: float64(2), object(7)\n memory usage: 218.1+ KB
```

Cell 14: Save Cleaned Data

```
df.to_csv("cleaned_data.csv", index=False)
```