## Data Profiling

### Dallas

```
[2]: import pandas as pd
     import numpy as np

[11]: df = pd.read_csv('Food_Inspections_20240225.tsv.tsv', delimiter ='\t')

      ---------------------------------------------------------------------------
      FileNotFoundError                         Traceback (most recent call last)
      Cell In[11], line 1
      ----> 1 df = pd.read_csv('Food_Inspections_20240225.tsv.tsv', delimiter = '\t')

      File ~\anaconda3\Lib\site-packages\pandas\io\parsers\readers.py:948, in read_csv(filepath_or_buffer, sep, delimiter, header, names, index_col, usecol
      s, dtype, engine, converters, true_values, false_values, skipinitialspace, skiprows, skipfooter, nrows, na_values, keep_default_na, na_filter, verbos
      e, skip_blank_lines, parse_dates, infer_datetime_format, keep_date_col, date_parser, date_format, dayfirst, cache_dates, iterator, chunksize, compress
      ion, thousands, decimal, lineterminator, quotechar, quoting, doublequote, escapechar, comment, encoding, encoding_errors, dialect, on_bad_lines, delim
      _whitespace, low_memory, memory_map, float_precision, storage_options, dtype_backend)
          935 kwds_defaults = _refine_defaults_read(
          936     dialect,
          937     delimiter,
          (...)
          944     dtype_backend=dtype_backend,
          945 )
          946 kwds.update(kwds_defaults)
      --> 948 return _read(filepath_or_buffer, kwds)

[5]: from ydata_profiling import ProfileReport

[6]: Dallas = ProfileReport(df, title="Report for Food Inspection")

[7]: Dallas.to_file("report_for_food_inspection.html")

      C:\Users\durge\anaconda3\Lib\site-packages\ydata_profiling\model\pandas\duplicates_pandas.py:40: PerformanceWarning: DataFrame is highly fragmented.
      This is usually the result of calling `frame.insert` many times, which has poor performance.  Consider joining all columns at once using pd.concat(axi
      s=1) instead. To get a de-fragmented frame, use `newframe = frame.copy()`
        .reset_index(name=duplicates_key)
      C:\Users\durge\anaconda3\Lib\site-packages\ydata_profiling\model\pandas\duplicates_pandas.py:40: PerformanceWarning: DataFrame is highly fragmented.
      This is usually the result of calling `frame.insert` many times, which has poor performance.  Consider joining all columns at once using pd.concat(axi
      s=1) instead. To get a de-fragmented frame, use `newframe = frame.copy()`
        .reset_index(name=duplicates_key)
      C:\Users\durge\anaconda3\Lib\site-packages\ydata_profiling\model\pandas\duplicates_pandas.py:40: PerformanceWarning: DataFrame is highly fragmented.
      This is usually the result of calling `frame.insert` many times, which has poor performance.  Consider joining all columns at once using pd.concat(axi
      s=1) instead. To get a de-fragmented frame, use `newframe = frame.copy()`
        .reset_index(name=duplicates_key)

      Generate report structure: 100% ████████████  1/1 [02:31<00:00, 151.67s/it]

      Render HTML: 100% ██████  1/1 [00:08<00:00, 8.02s/it]

      Export report to file: 100% ██████  1/1 [00:00<00:00, 3.25it/s]
```
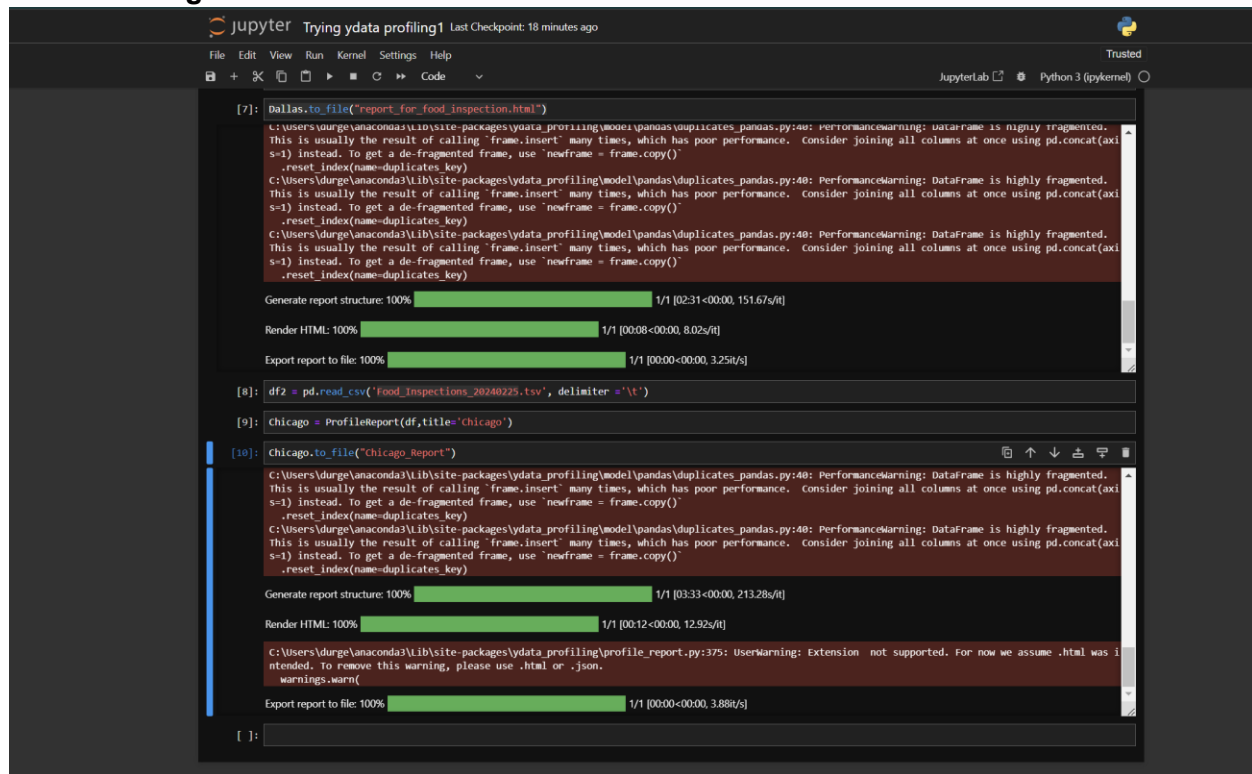
Here we used ydata_profiling for profiling of **"Restaurant and Food Establishment Inspections (October 2016 to Present) | Dallas OpenDataLinks to an external site"**.

We stored the file to a dataframe and then using ydata profiling generated a "Profile Report" in html.

**Chicago**



Just like dallas, we have used YdataProfiling for Chicago as well to generate a data report for Chicago.

The following is the Dimensional Model that was created in ER Studio:

Workflow for Chicago Dataset.



After data profiling, we used the file input delimited component to get the file with tmap component and database out as the component to stage the data in MySql.



Later we created the above ETL pipeline. Here we took the staged table as an input and used tnormalize to normalize the data then mapped with tmap and used it to map out the violations.

Later we used tAggregateRow and then used tmap to get the output for "Inspection Fact", "Facility Dimension", "Address Dimension" and "Date Dimension".

The following is the tmap configuration used to map the normalized data which was used to sperate the violations



The following is the tmap configuration used to Aggregate the data and then populated in different dimensional and one fact table.

## Staging of Dallas dataset



Here we took an input mapped using tmap and then the mapped output is written to database.

s



the above is the ETL pipeline that was used where the input is the staged data of dallas and then this was mapped and the output was passed to different dimensions and one fact table.

The above is the tmap used for putting the data in different dimensions and the fact table.



This is the tmap configuration used to stage the dallas dataset.

The following are the count row queries for all the dimensions and one fact table

**Screenshot 1**

MyDB ×

File  Edit  View  Query  Database  Server  Tools  Scripting  Help

Tabs: ate | damg7370.dim_insepction | damg7370.dim_violations | damg7370.facility_dim | damg7370.fact_insepction | address_dim | dim_date | damg7370.dim_insepction | dim_insepction | dim_violations

Navigator

SCHEMAS

Filter objects

- address_dim
- dim_date
- dim_insepction
  - Columns
  - Indexes
  - Foreign Keys
  - Triggers
- dim_violations
- facility_dim
- fact_insepction
  - Columns
  - Indexes
  - Foreign Keys
  - Triggers
- stg_chicago
- stg_chicagofoodinspectiontable
- stg_chicagotable
- stg_chicagotable2
- stg_dallas_address_dimension_table
- stg_dallas_date_dimension_table
- stg_dallas_facility_dimension_table
- stg_dallas_inspection_fact_table
- stg_dallasconnection
- stg_dallasconnectiontable

Administration  Schemas

Information

Table: dim_violations

Columns:
- Violations_Sk    int PK
- Inspection_ID    int
- Violation_codes  varchar(4000)
- Violations_description  varchar(8000)

Object Info  Session

Query:
```
1  SELECT count(*) FROM damg7370.dim_violations;
```

Result Grid | Filter Rows:        | Export: | Wrap Cell Content:

count(*)
1343443

Result 2 ×                                               Read Only

Output

Action Output

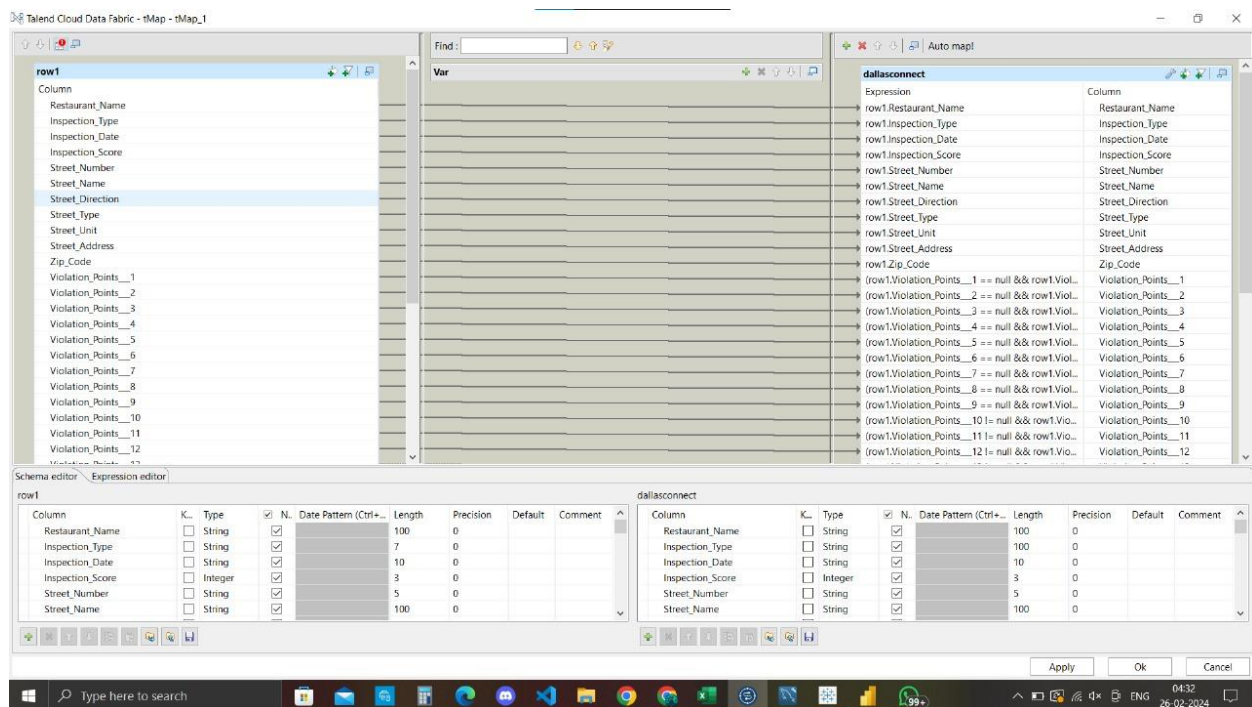| # | Time | Action | Message | Duration / Fetch |
|---|------|--------|---------|------------------|
| 30 | 04:37:10 | SELECT * FROM damg7370.dim_insepction | 267908 row(s) returned | 0.000 sec / 3.031 sec |
| 31 | 04:37:22 | SELECT count(*) FROM damg7370.dim_insepction | 1 row(s) returned | 0.078 sec / 0.000 sec |
| 32 | 04:37:36 | SELECT * FROM damg7370.dim_violations | 1343443 row(s) returned | 0.000 sec / 4.578 sec |
| 33 | 04:37:51 | SELECT count(*) FROM damg7370.dim_violations | 1 row(s) returned | 3.125 sec / 0.000 sec |

Type here to search     ENG  04:37  26-02-2024

Result Grid / Form Editor / Field Types

---

**Screenshot 2**

MySQL Workbench

MyDB ×

File  Edit  View  Query  Database  Server  Tools  Scripting  Help

Tabs: dim_insepction | damg7370.dim_violations | damg7370.facility_dim | damg7370.fact_insepction | address_dim | dim_date | damg7370.dim_insepction | dim_insepction | dim_violations | facility_dim

Navigator

SCHEMAS

Filter objects

- address_dim
- dim_date
- dim_insepction
  - Columns
  - Indexes
  - Foreign Keys
  - Triggers
- dim_violations
- facility_dim
- fact_insepction
  - Columns
  - Indexes
  - Foreign Keys
  - Triggers
- stg_chicago
- stg_chicagofoodinspectiontable
- stg_chicagotable
- stg_chicagotable2
- stg_dallas_address_dimension_table
- stg_dallas_date_dimension_table
- stg_dallas_facility_dimension_table
- stg_dallas_inspection_fact_table
- stg_dallasconnection
- stg_dallasconnectiontable

Administration  Schemas

Information

Table: facility_dim

Columns:
- Facility_Sk    int PK
- DBA_Name    varchar(100)
- AKA_Name    varchar(100)
- License    int
- Facility_Type    varchar(100)
- DI_WorkFlowFileName    varchar(20)
- Date_Sk    int

Object Info  Session

Query:
```
1  SELECT count(*) FROM damg7370.facility_dim;
```

Result Grid | Filter Rows:        | Export: | Wrap Cell Content:

count(*)
346308

Result 2 ×                                               Read Only

Output

Action Output

| # | Time | Action | Message | Duration / Fetch |
|---|------|--------|---------|------------------|
| 32 | 04:37:36 | SELECT * FROM damg7370.dim_violations | 1343443 row(s) returned | 0.000 sec / 4.578 sec |
| 33 | 04:37:51 | SELECT count(*) FROM damg7370.dim_violations | 1 row(s) returned | 3.125 sec / 0.000 sec |
| 34 | 04:38:04 | SELECT * FROM damg7370.facility_dim | 346308 row(s) returned | 0.016 sec / 0.890 sec |
| 35 | 04:38:10 | SELECT count(*) FROM damg7370.facility_dim | 1 row(s) returned | 0.078 sec / 0.000 sec |

Type here to search     ENG  04:38  26-02-2024

Result Grid / Form Editor / Field Types

The following are the DDL scripts for the above tables:
**Address_Dim**

CREATE TABLE `address_dim` (

  `Address_Sk` int NOT NULL,

  `Address` varchar(100) DEFAULT NULL,

  `City` varchar(100) DEFAULT NULL,

  `State` varchar(100) DEFAULT NULL,

  `Zip` int DEFAULT NULL,

  `Latitude` varchar(100) DEFAULT NULL,

  `Longitude` varchar(100) DEFAULT NULL,

  `Location` varchar(100) DEFAULT NULL,

  `DI_WorkFlowFileName` varchar(20) DEFAULT NULL,

  `DI_CreateDate` datetime DEFAULT NULL,

  PRIMARY KEY (`Address_Sk`)

) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci


**dim_Date**

```
CREATE TABLE `dim_date` (
  `Date_Sk` int NOT NULL,
  `year` varchar(500) DEFAULT NULL,
  `day` varchar(500) DEFAULT NULL,
  `month` varchar(500) DEFAULT NULL,
  PRIMARY KEY (`Date_Sk`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci
```

**dim_Inspection**
```
CREATE TABLE `dim_insepction` (
  `inspection_sk` int NOT NULL,
  `inspection_ID` int DEFAULT NULL,
  `facility_sk` int NOT NULL,
  `address_sk` int NOT NULL,
  `Risk` varchar(50) DEFAULT NULL,
  `Inspection_Type` varchar(50) DEFAULT NULL,
  `Results` varchar(20) DEFAULT NULL,
  `TotalViolationScore` int DEFAULT NULL,
  `DI_WorkFlowFileName` varchar(20) DEFAULT NULL,
  `Inspection_Result` varchar(100) DEFAULT NULL,
  `Date_Sk` int NOT NULL,
  PRIMARY KEY (`inspection_sk`)
) ENGINE=InnoDB DEFAULT CHARFSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci
```

**dim_violations**
```
CREATE TABLE `dim_violations` (
```

```
  `Violations_Sk` int NOT NULL,

  `Inspection_ID` int DEFAULT NULL,

  `Violation_codes` varchar(4000) DEFAULT NULL,

  `Violations_description` varchar(8000) DEFAULT NULL,

  `DI_WorkFlowFileName` varchar(100) DEFAULT NULL,

  PRIMARY KEY (`Violations_Sk`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci
```

**Facility_Dim**

```
CREATE TABLE `facility_dim` (

  `Facility_Sk` int NOT NULL,

  `DBA_Name` varchar(100) DEFAULT NULL,

  `AKA_Name` varchar(100) DEFAULT NULL,

  `License` int DEFAULT NULL,

  `Facility_Type` varchar(100) DEFAULT NULL,

  `DI_WorkFlowFileName` varchar(20) DEFAULT NULL,

  `Date_Sk` int NOT NULL,

  `DI_CreateDate` datetime DEFAULT NULL,

  PRIMARY KEY (`Facility_Sk`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci
```

**Fact_Inspection**

```
CREATE TABLE `fact_insepction` (

  `Inspection_Sk` int NOT NULL,
```

```
  `Inspection_ID` int DEFAULT NULL,

  `Facility_Sk` int NOT NULL,

  `Address_Sk` int NOT NULL,

  `Risk` varchar(50) DEFAULT NULL,

  `Inspection_Type` varchar(50) DEFAULT NULL,

  `Results` varchar(20) DEFAULT NULL,

  `Total_Violation_Score` int DEFAULT NULL,

  `DI_WorkFlowFileName` varchar(20) DEFAULT NULL,

  `Inspection_result` varchar(100) DEFAULT NULL,

  PRIMARY KEY (`Inspection_Sk`)

) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci
```

**Dashboard  Screenshots**

**Pages**

**Columns** | Longitude (generated)

**Rows** | Latitude (generated)

**Filters**
- Measure Names
- Zip Code
- Inspection Result

**Marks**

Detail | Tooltip

- Inspection Result
- Measure Values
- Zip Code
- Measure Values

**Measure Values**
- CNT(dimaddress)
- AVG(Latitude)
- AVG(Longitude)

### Inspection Result Vs Location



© 2024 Mapbox © OpenStreetMap

24 nulls

**Inspection Result**
- Pass
- Pass with Warning

**Measure Values**
- · -98
- 1,000
- 2,000
- 3,000
- 4,000
- 4,674

### Sheet 2



Canvass | License | Canvass Re-Inspection

Complaint | Complaint | Short Form

License Re-Inspection

CNT(Inspection ID)

0      1,39,241

### Sheet 3



**Risk**
- Null
- All
- Risk 1 (High)
- Risk 2 (Medium)
- Risk 3 (Low)

## Columns
CNT(Inspection ID)

## Rows
Inspection result

## Sheet 1

**Inspection result**

| | |
|---|---|
| Fail | |
| Pass | |
| Pass with warning | |

0K 10K 20K 30K 40K 50K 60K 70K 80K 90K 100K 110K 120K 130K

Count of Inspection ID

## Columns
CNT(Risk)

## Rows
City

## Sheet 4

**City**

312CHICAGO
ALGONQUIN
ALSIP
BANNOCKBURNDEE...
BERWYN
BLOOMINGDALE
BLUE ISLAND
BOLINGBROOK
BRIDGEVIEW
BROADVIEW
BURBANK
BURNHAM
CALUMET CITY
CCHICAGO
CH
CHARLES A HAYES
CHCHICAGO
CHCICAGO
CHICAGO
CHICAGO HEIGHTS
CHICAGO.
chicagoBEDFORD PA...
CHICAGOC
CHICAGOCHICAGO
CHICAGOI
CHICAGOO
CICERO
COUNTRY CLUB HILLS
Dallas
DES PLAINES
EAST HAZEL CREST

0K 20K 40K 60K 80K 100K 120K 140K 160K 180K 200K 220K 240K 260K 280K

Count of Risk

**Sheet 2**

| Canvass | License |
| | Complaint |

**Sheet 1**

Inspection result

Fail
Pass
Pass with warning

0K    50K    100K

Count of Inspection ID

Count of Inspection ID

0 ▭ 1,39,241

Risk
■ Null
■ All
■ Risk 1 (High)
■ Risk 2 (Medium)
■ Risk 3 (Low)

**Sheet 5**

Facility Type

15 monts to 5 years ..
1005 NURSING HOME
1023
1023 CHILDERN'S SE..
1023 CHILDERN'S SE..
1023 CHILDREN'S SE..
1023-CHILDREN'S S..
1475 LIQUOR
1584-DAY CARE ABO..
(convenience store)
(gas station)
A-Not-For-Profit Che..
ADULT DAYCARE
Adult Family Care Ce..
AFTER SCHOOL CARE
AFTER SCHOOL PRO..
Airport Lounge
ALTERNATIVE SCHO..
Animal Shelter Cafe ..
ARCHDIOCESE

Facility Type:
Count of Inspection ID: 5,119

50K 100K

Count of Insp..

**Sheet 4**

City

312CHICAGO
ALGONQUIN
ALSIP
BANNOCKBURNDEE..
BERWYN
BLOOMINGDALE
BLUE ISLAND
BOLINGBROOK
BRIDGEVIEW
BROADVIEW
BURBANK
BURNHAM

0K    100K    200K

Count of Risk

**Sheet 3**

**Columns** CNT(Inspection ID)

**Rows** Facility Type

## Sheet 5

Facility Type

| Facility Type | |
|---|---|
| 15 monts to 5 years .. | |
| 1005 NURSING HOME | |
| 1023 | |
| 1023 CHILDERN'S SE.. | |
| 1023 CHILDERN'S SE.. | |
| 1023 CHILDREN'S SE.. | |
| 1023-CHILDREN'S S.. | |
| 1475 LIQUOR | |
| 1584-DAY CARE ABO.. | |
| (convenience store) | |
| (gas station) | |
| A-Not-For-Profit Che.. | |
| ADULT DAYCARE | |
| Adult Family Care Ce.. | |
| AFTER SCHOOL CARE | |
| AFTER SCHOOL PRO.. | |
| Airport Lounge | |
| ALTERNATIVE SCHO.. | |
| Animal Shelter Cafe .. | |
| ARCHDIOCESE | |
| ART CENTER | |
| ART GALLERY | |
| ART GALLERY W/WI.. | |
| ASSISSTED LIVING | |
| ASSISTED LIVING | |
| Assisted Living Seni.. | |
| Bakery | |
| BAKERY/ RESTAURA.. | |
| BAKERY/DELI | |
| BAKERY/GROCERY | |
| bakery/restaurant | |

0K   10K   20K   30K   40K   50K   60K   70K   80K   90K   100K   110K   120K   130K   140K

Count of Inspection ID

**PowerBI Screenshots**

## Location and Risk

**Risk** ● ● All ● Risk 1 (High)

## Inspection ID vs Inspection Result

46.09K (17.2%)

130.06K (48.55%)

91.76K (34.25%)

**Inspection_Result**
- Pass with warning
- Fail
- Pass

## Count of inspection_ID by Risk

| | | |
|---|---|---|
| Risk 1 (High) | Risk 2 (Medium) | Risk 3 (Low) |

Total Inspections

Risk

## Facility Type by Inspection ID

- Restaurant
- Grocery Store
- School
- Children's Services Facility
- Bakery
- Daycare Above and Under 2 Years
- Daycare (2 - 6 Years)
- Long Term Care
- Catering

Facility Type

Count of inspection_ID

## Facility Type by Inspection ID



| Facility Type | Count of inspection_ID |
|---|---|
| Restaurant | |
| Grocery Store | |
| School | |
| Children's Services Facility | |
| Bakery | |
| Daycare Above and Under 2 Years | |
| Daycare (2 - 6 Years) | |
| Long Term Care | |
| Catering | |
| Liquor | |
| Mobile Food Preparer | |
| Mobile Food Dispenser | |
| Hospital | |
| Golden Diner | |
| Daycare Combo 1586 | |
| Wholesale | |
| tavern | |
| Shared Kitchen User (Long Term) | |
| Daycare (Under 2 Years) | |
| Special Event | |
| Banquet Hall | |
| Shared Kitchen | |
| GAS STATION | |
| CHARTER SCHOOL | |
| Shelter | |
| Mobile Prepared Food Vendor | |
| Live Poultry | |
| BANQUET | |
| KIOSK | |

Count of inspection_ID (axis: 0K, 50K, 100K, 150K, 200K)

| DI_WorkFlowFileName | Restaurant_Name | Results | Count of inspection_id | TotalViolationScore |
|---|---|---|---|---|
| Dallas Profile | AIN | Business Not Located | 1 | 32 |
| Dallas Profile | ALOFT/ELEMENT CATERING KITCHEN | Business Not Located | 1 | 22 |
| Dallas Profile | CENTERPLATE GO NATURAL #1 KIOSK | Business Not Located | 1 | 2 |
| Dallas Profile | CHICK-FIL-A | Business Not Located | 1 | 22 |
| Dallas Profile | CONCRETE COWBOY | Business Not Located | 1 | 12 |
| Dallas Profile | EATALY DALLAS - GROCERY (GROCERY/BUTCHER/SEAFOOD) | Business Not Located | 1 | 12 |
| Dallas Profile | EL RANCHO #8 / BAKERY | Business Not Located | 1 | 32 |
| Dallas Profile | MARROSSO CAFE | Business Not Located | 1 | 2 |
| Dallas Profile | MOCKINGBIRD ELEMENTARY SCHOOL | Business Not Located | 1 | 32 |
| Dallas Profile | PAESANOS CAFE | Business Not Located | 1 | 12 |
| Dallas Profile | PASTERERIA DEL NORTE | Business Not Located | 1 | 32 |
| **Total** | | | **78400** | |

| DI_WorkFlowFileName | DBA_Name | Results | Count of Results | TotalViolationScore |
|---|---|---|---|---|
| Chicago Profile | #1 CHINA EXPRESS, LTD. | Out of Business | 1 | 32 |
| Chicago Profile | #1 CHINA EXPRESS, LTD. | Pass | 1 | 42 |
| Chicago Profile | #1 CHOP SUEY | Fail | 1 | 102 |
| Chicago Profile | #1 CHOP SUEY | Fail | 2 | 108 |
| Chicago Profile | #1 CHOP SUEY | Fail | 3 | 110 |
| Chicago Profile | #1 CHOP SUEY | Fail | 1 | 112 |
| Chicago Profile | #1 CHOP SUEY | Fail | 1 | 116 |
| Chicago Profile | #1 CHOP SUEY | Fail | 1 | 124 |
| Chicago Profile | #1 CHOP SUEY | No Entry | 2 | 32 |
| Chicago Profile | #1 CHOP SUEY | Out of Business | 2 | 32 |
| Chicago Profile | #1 CHOP SUEY | Pass | 8 | 32 |
| **Total** | | | **267908** | |