

Basic aggregate functions

TIME SERIES ANALYSIS IN SQL SERVER



Maham Faisal Khan

Senior Data Science Content Developer

Key aggregation functions

Counts

`COUNT()`

`COUNT_BIG()`

`COUNT(DISTINCT)`

Other Aggregates

`SUM()`

`MIN()`

`MAX()`

What counts with COUNT()

Number of Rows

```
COUNT(*)
```

```
COUNT(1)
```

```
COUNT(1/0)
```

Non-NULL Values

```
COUNT(d.YR)
```

```
COUNT(NULLIF(d.YR, 1990))
```

Distinct counts

```
SELECT
    COUNT(DISTINCT c.CalendarYear) AS Years,
    COUNT(DISTINCT NULLIF(c.CalendarYear, 2010)) AS Y2
FROM dbo.Calendar c;
```

Years	Y2
50	49

Filtering aggregates with CASE

```
SELECT
    MAX(CASE WHEN ir.IncidentTypeID = 1
        THEN ir.IncidentDate
        ELSE NULL
    END) AS I1,
    MAX(CASE WHEN ir.IncidentTypeID = 2
        THEN ir.IncidentDate
        ELSE NULL
    END) AS I2,
FROM dbo.IncidentRollup ir;
```

I1	I2
2020-06-30	2020-06-29

Let's practice!

TIME SERIES ANALYSIS IN SQL SERVER

Statistical aggregate functions

TIME SERIES ANALYSIS IN SQL SERVER



Maham Faisal Khan

Senior Data Science Content Developer

Statistical aggregate functions

AVG()

Mean

STDEV()

Standard Deviation

STDEVP()

Population Standard Deviation

VAR()

Variance

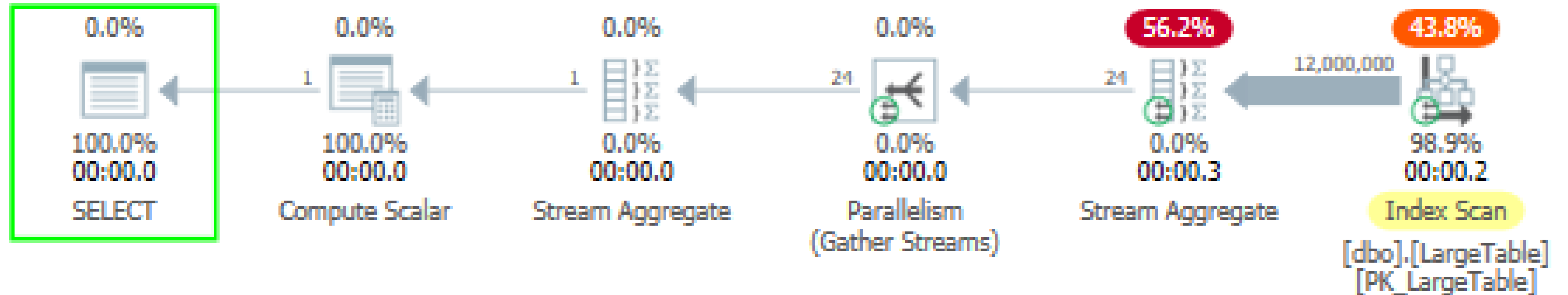
VARP()

Population Variance

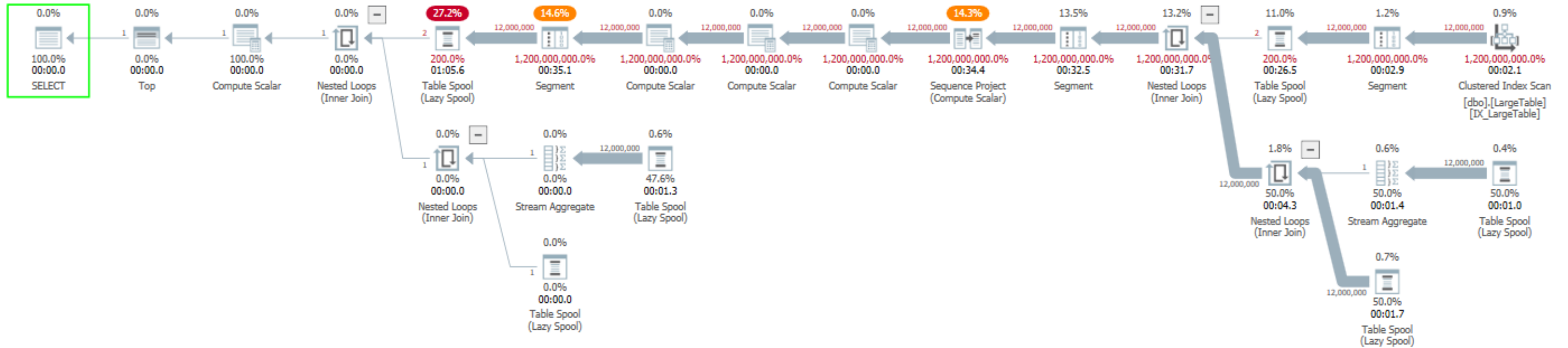
What about median?

```
SELECT TOP(1)
    PERCENTILE_CONT(0.5)
        WITHIN GROUP (ORDER BY l.SomeVal DESC)
        OVER () AS MedianIncidents
FROM dbo.LargeTable l;
```

But how bad is it?



This bad



The cost of median

	Median	Mean
Est. Cost	95.7%	4.3%
Duration	68.5s	0.37s
CPU	68.5s	8.1s
Reads	72,560,946	39,468
Writes	87,982	0

Let's practice!

TIME SERIES ANALYSIS IN SQL SERVER

Downsampling and upsampling data

TIME SERIES ANALYSIS IN SQL SERVER

SQL

Maham Faisal Khan

Senior Data Science Content Developer

Data in nature

SELECT

SomeDate

FROM dbo.SomeTable

SomeDate

2019-08-11 06:14:29.990

2019-08-11 11:07:37.633

2019-08-11 14:08:00.337

Downsampling data

```
SELECT  
    CAST(SomeDate AS DATE) AS SomeDate  
FROM dbo.SomeTable
```

SomeDate
2019-08-11
2019-08-11
2019-08-11

Further downsampling

SELECT

DATEADD(**HOUR**, DATEDIFF(**HOUR**, 0, SomeDate), 0) **AS** SomeDate

FROM dbo.SomeTable

DATEDIFF(HOUR, 0, '2019-08-11 06:21:16') = 1,048,470

DATEADD(HOUR, 1048470, 0) = 2019-08-11 06:00:00

SomeDate

2019-08-11 06:00:00.000

2019-08-11 11:00:00.000

2019-08-11 14:00:00.000

What about upsampling?

Downsampling

- Aggregate data
- Can usually sum or count results
- Provides a higher-level picture of the data
- Acceptable for most purposes

Upsampling

- Disaggregate data
- Need an allocation rule
- Provides artificial granularity
- Acceptable for data generation, calculated averages

Let's practice!

TIME SERIES ANALYSIS IN SQL SERVER

Grouping by ROLLUP, CUBE, and GROUPING SETS

TIME SERIES ANALYSIS IN SQL SERVER

SQL

Maham Faisal Khan

Senior Data Science Content Developer

Hierarchical rollups with ROLLUP

```
SELECT
    t.Month,
    t.Day,
    SUM(t.Events) AS Events
FROM Table
GROUP BY
    t.Month,
    t.Day
WITH ROLLUP
ORDER BY
    t.Month,
    t.Day;
```

Month	Day	Events
NULL	NULL	100
1	NULL	60
1	1	3
1	2	4
...
2	NULL	40
2	1	8

Cartesian aggregation with CUBE

```
SELECT
    t.IncidentType,
    t.Office,
    SUM(t.Events) AS Events
FROM Table
GROUP BY
    t.IncidentType,
    t.Office
WITH CUBE
ORDER BY
    t.IncidentType,
    t.Office;
```

IncidentType	Office	Events
NULL	NULL	250
NULL	NY	70
NULL	CT	180
T1	NULL	55
T1	NY	30
T1	CT	25

Define grouping sets with GROUPING SETS

```
SELECT
    t.IncidentType,
    t.Office,
    SUM(t.Events) AS Events
FROM Table
GROUP BY GROUPING SETS
(
    (t.IncidentType, t.Office),
    ()
)
ORDER BY
    t.IncidentType,
    t.Office;
```

IncidentType	Office	Events
NULL	NULL	250
T1	NY	30
T1	CT	25
T2	NY	10
T2	CT	110
T3	NY	30
T3	CT	45

Let's practice!

TIME SERIES ANALYSIS IN SQL SERVER