

YIEPI meets DNA: on vehicle datasets analysis

on-going progress, findings and record-matching strategies

YIEPI meets DNA: on vehicle datasets analysis

let's go!

vehicle datasets analysis

**we all want to
identify sold
vehicles...**

identify

we all want to

ok, a glimpse on the data!

vehicles...

± 260K
●
b2b data

± 6.25Mio.

b2c data

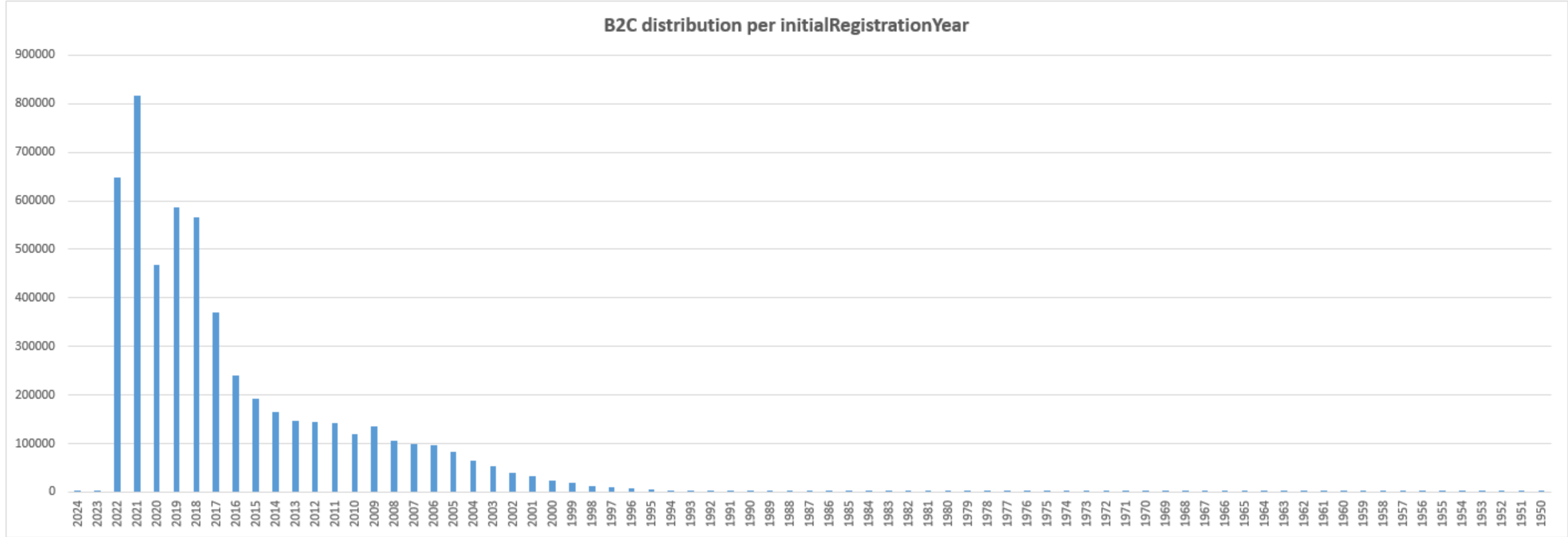
*b2b not even
5% of the volume!*

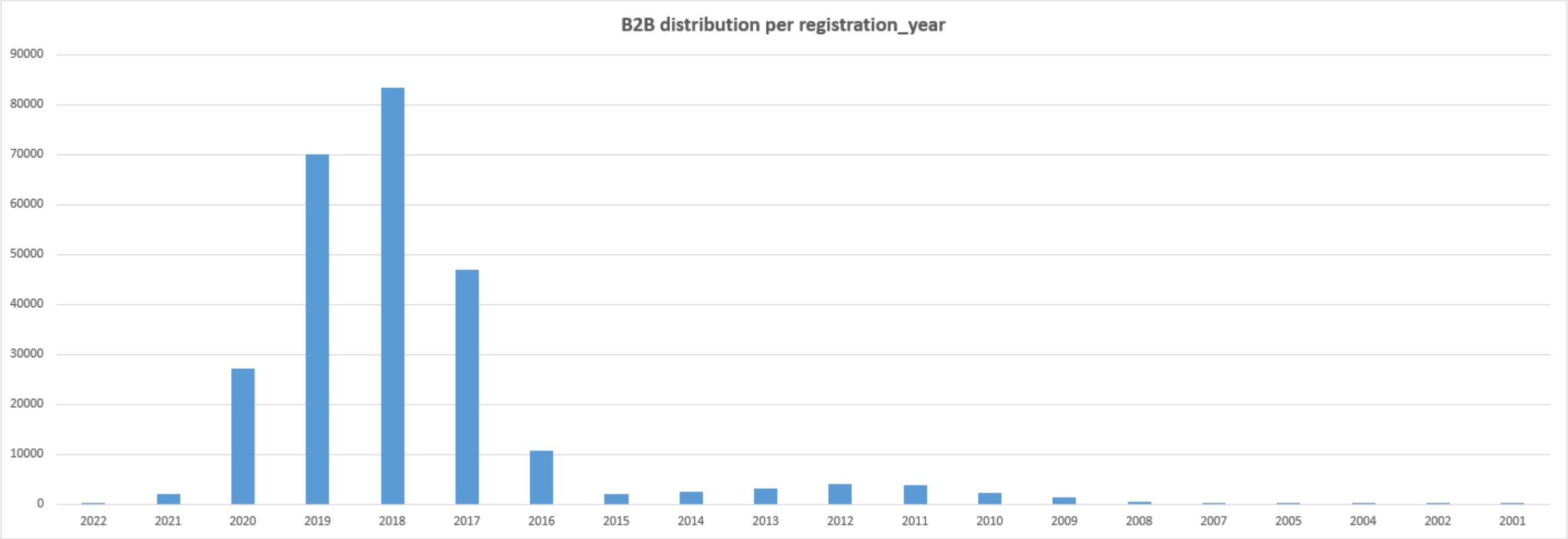
± 260K
●
b2b data

± 6.25Mio.

b2c data

**some other quick
insights...**





**hey, what about the
vin number?**

hey, what about the vin number?

null in approx. 85%
of b2c records! →

**we all want to
identify sold
vehicles...**

we all want to

ok, so let's match the vins!

vehicles...

± 60K
•
**matching
vins**

± 260K
•
b2b data

± 6.25Mio.

b2c data

± 60K
•
**matching
vins**

± 260K
•
b2b data

± 6.25Mio.

b2c data

→ 23%, and around 1.000
duplicated

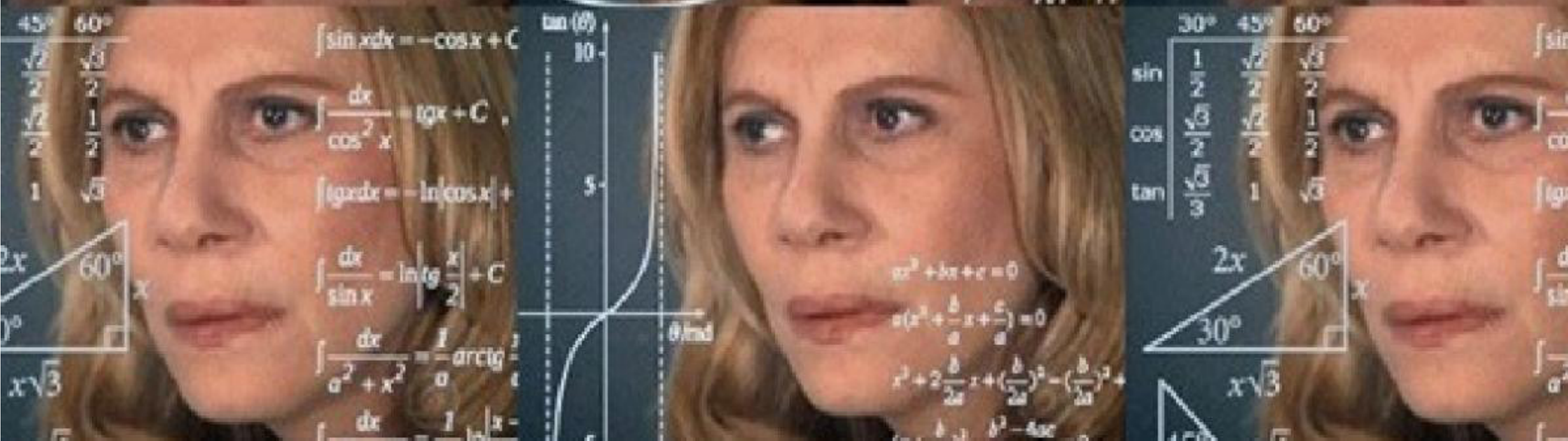
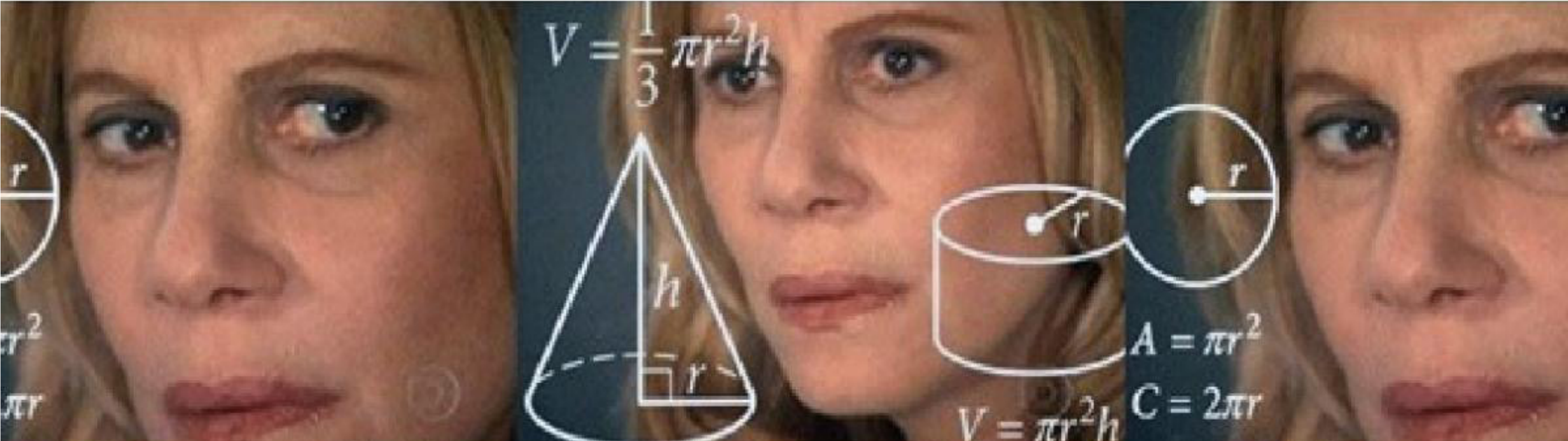


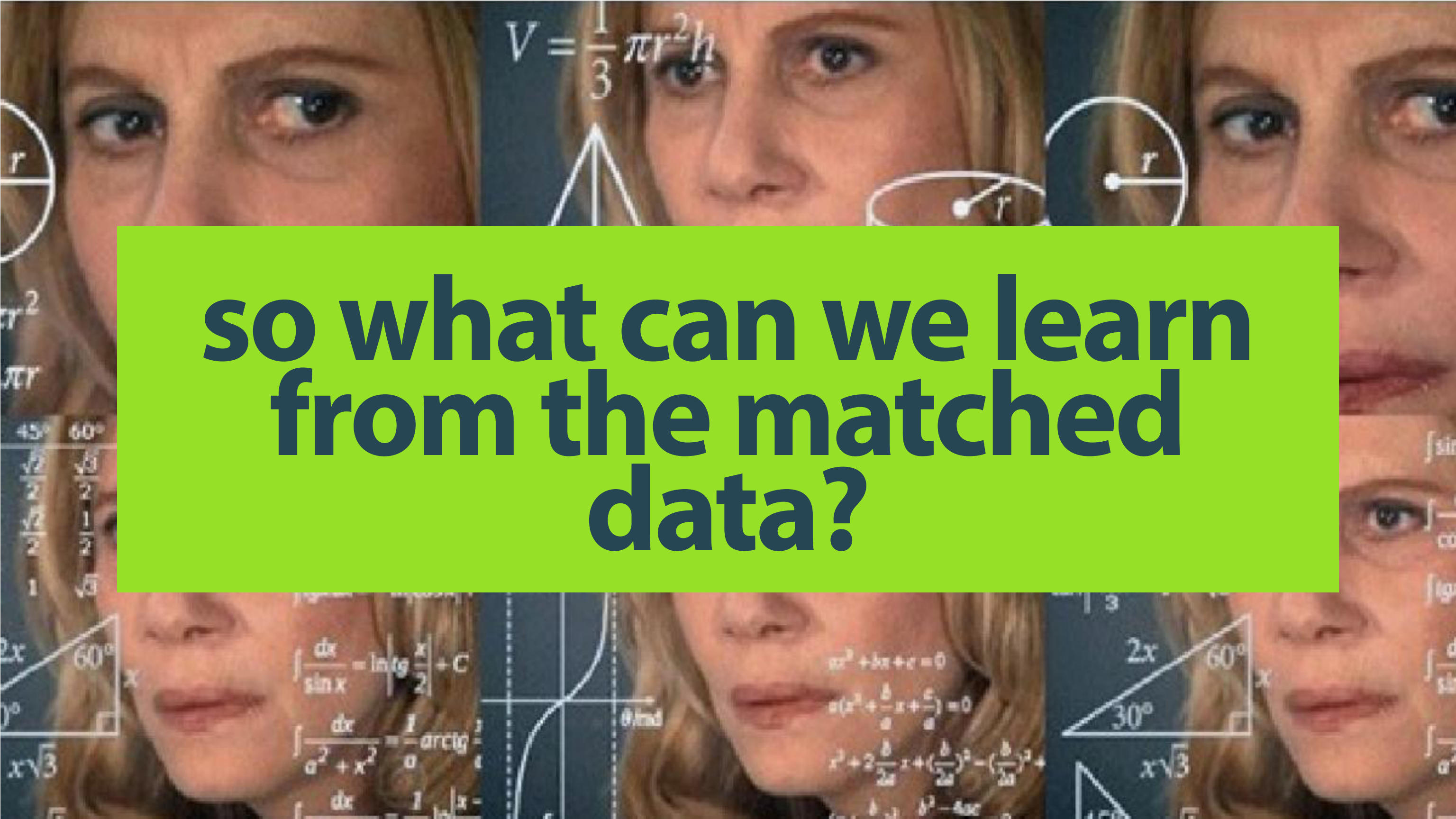
**ok, Super VIN is our
first hero**



ok, Super VIN is our
first hero

but we still have
75% of the b2b
left to tackle





**so what can we learn
from the matched
data?**



\pm 60K

matched vins



± 60K

matched vins

let's call it our
"golden dataset"

quick hypotheses:

quick hypotheses:

- maybe model matches;

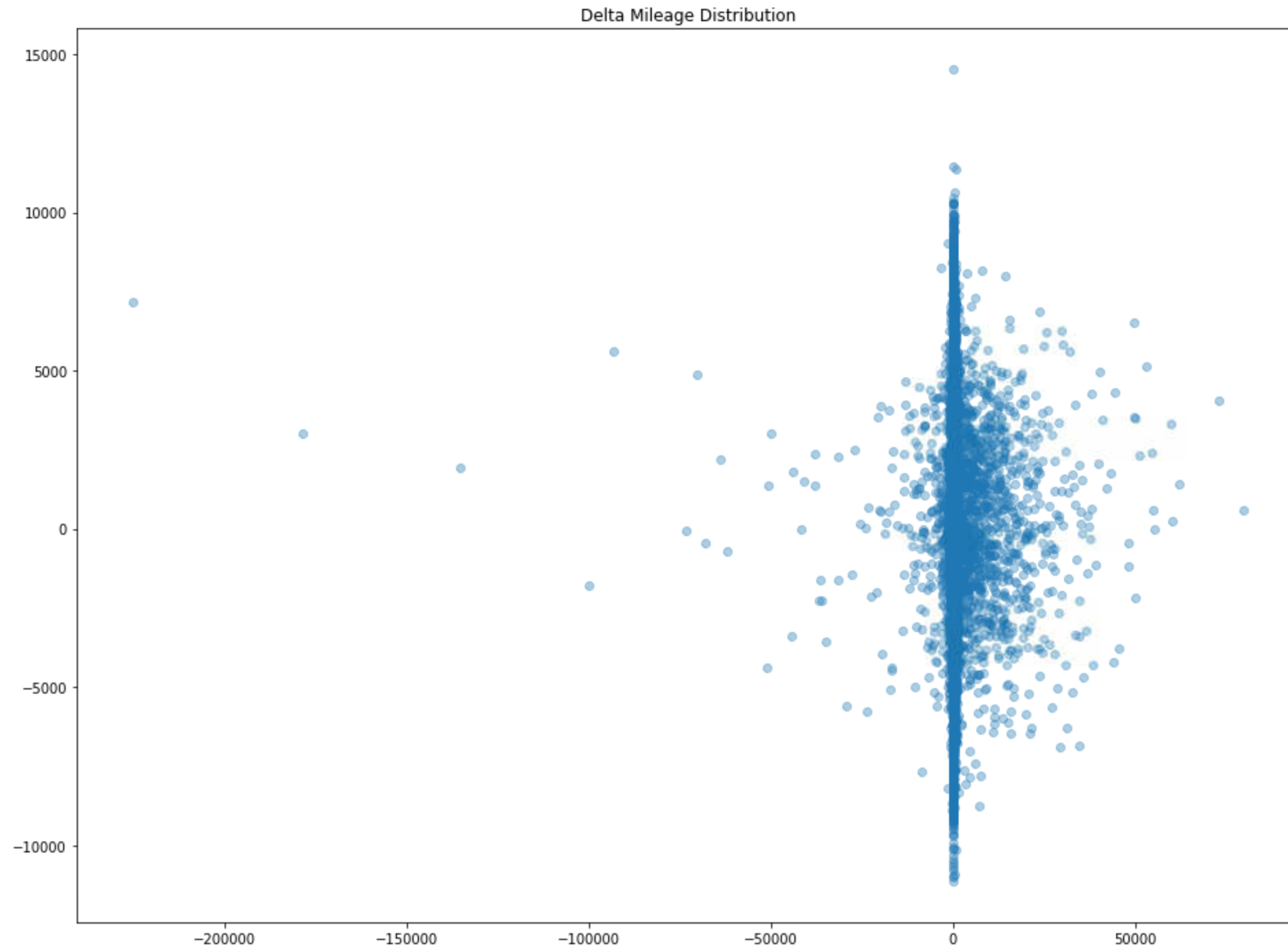
quick hypotheses:

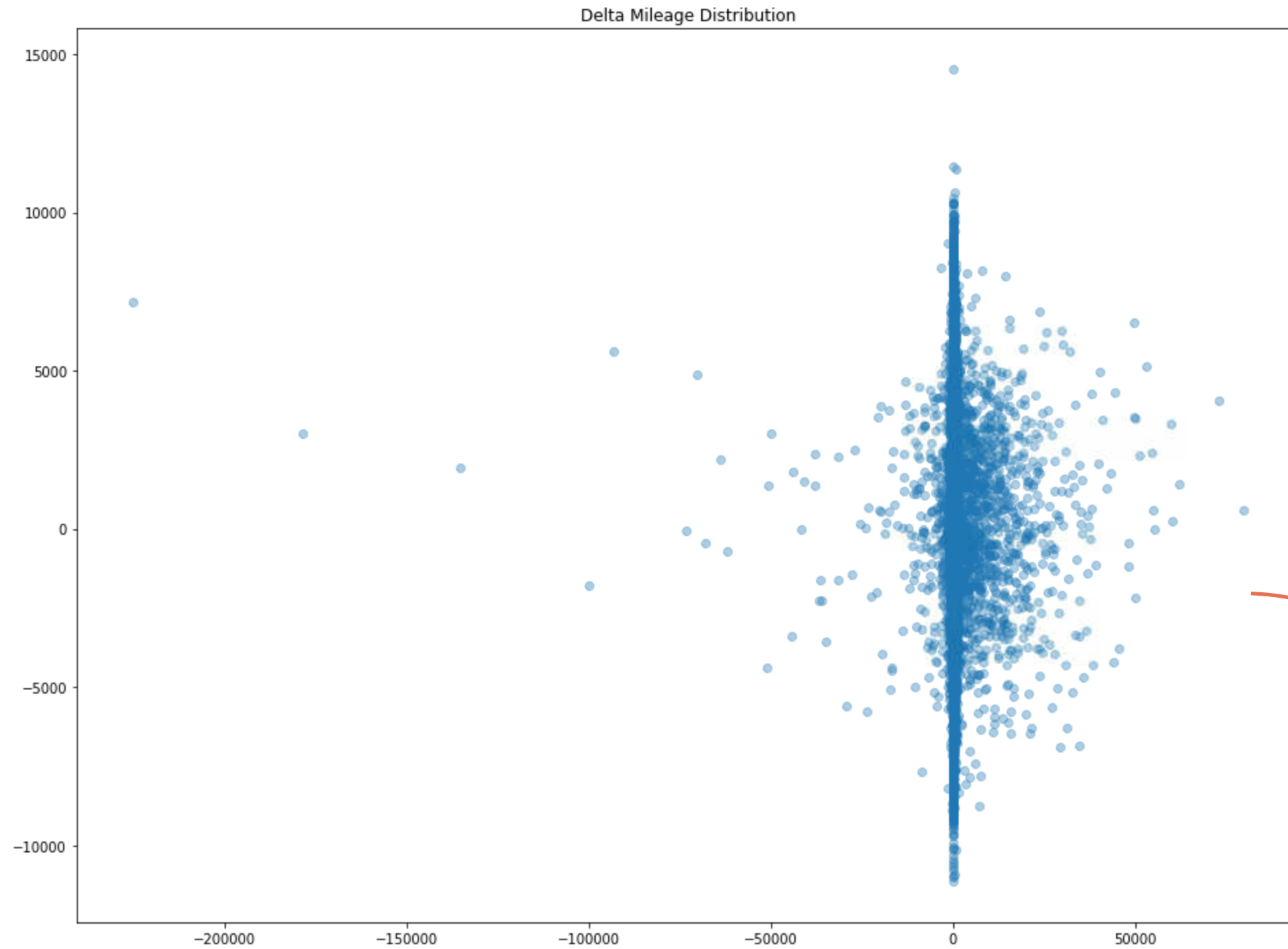
- maybe model matches;

98.70% of the cases!

quick hypotheses:

- maybe model matches;
- **maybe mileage matches;**





50% have a delta
of zero

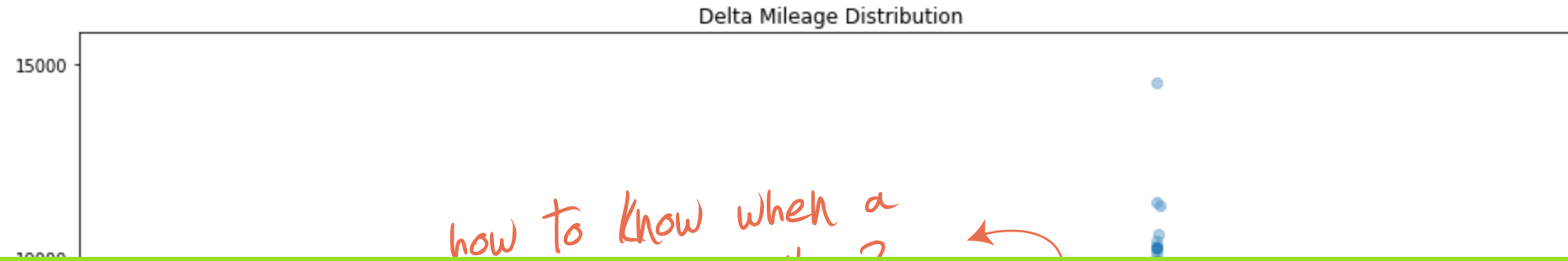




**ok, let's take a deeper
look here**

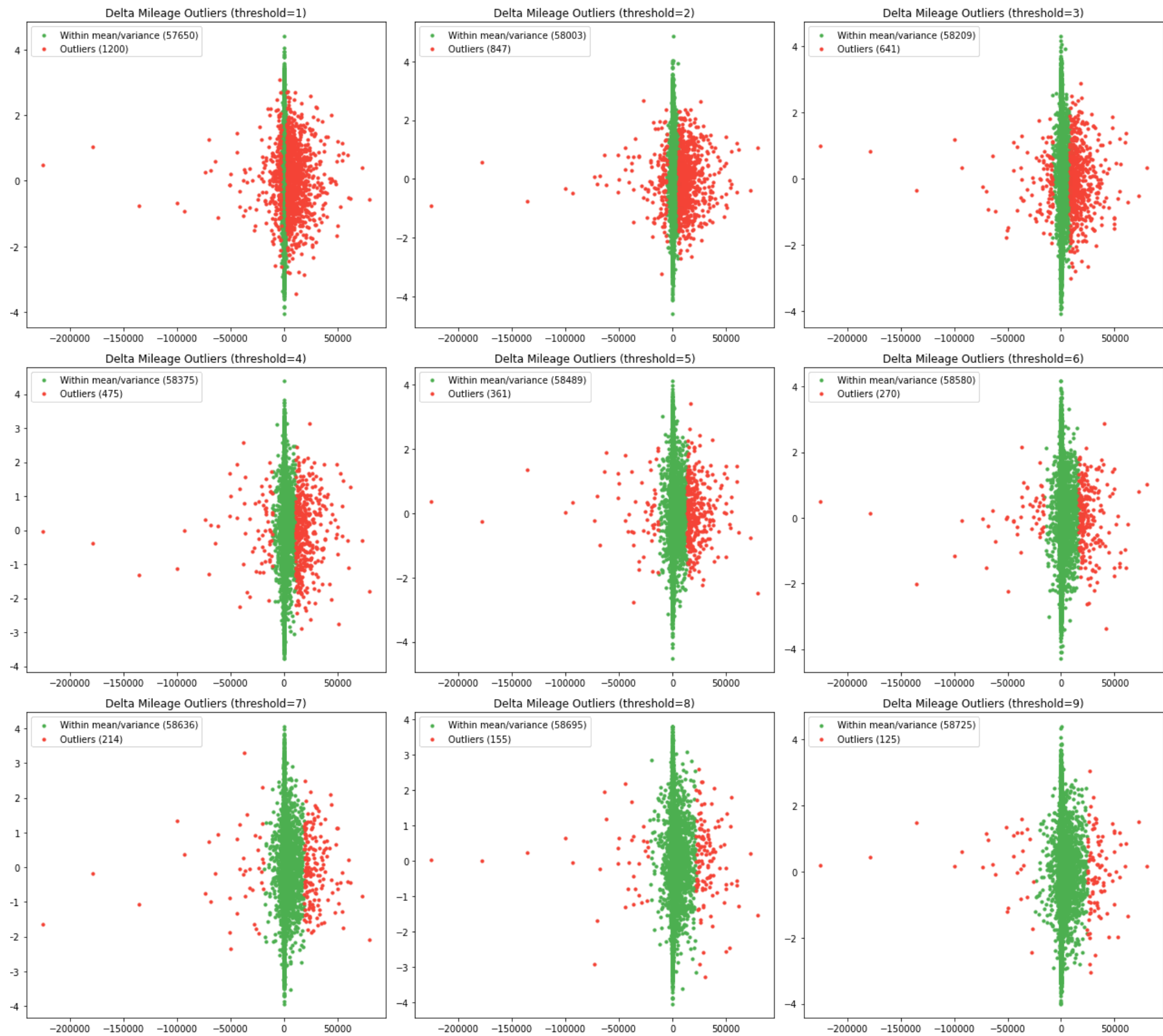


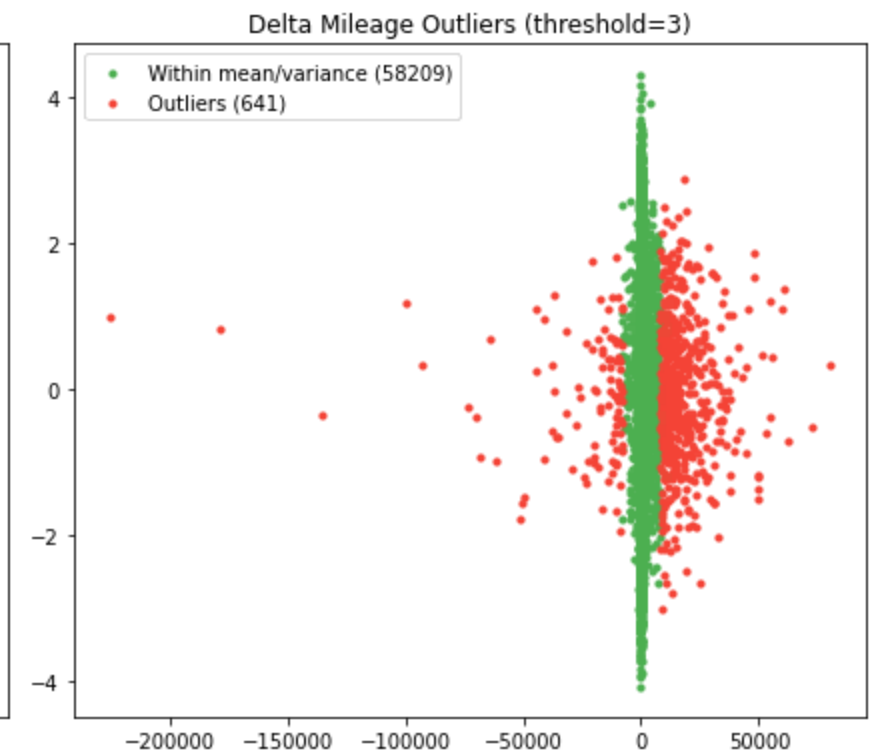
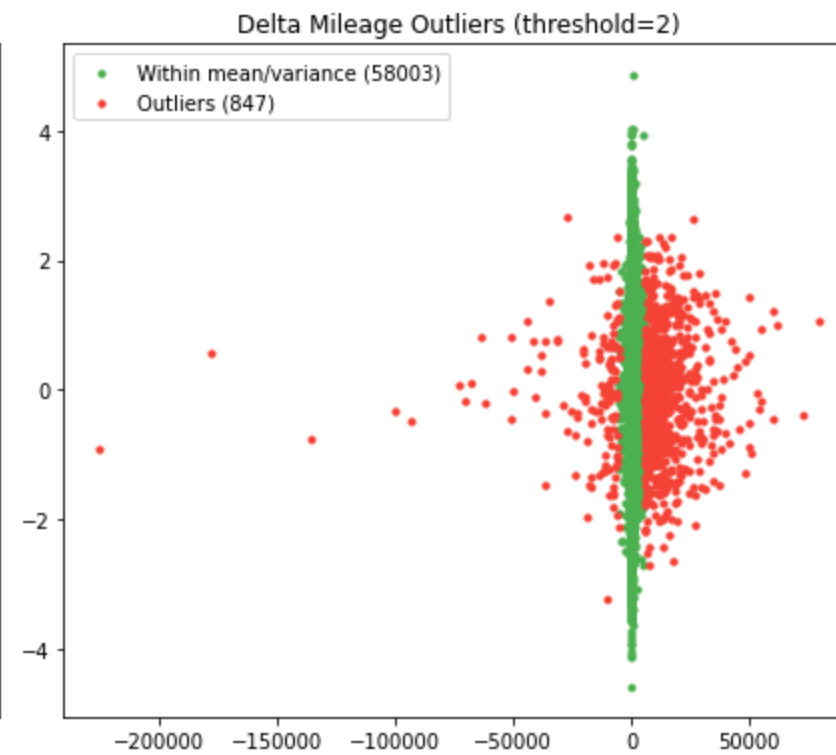
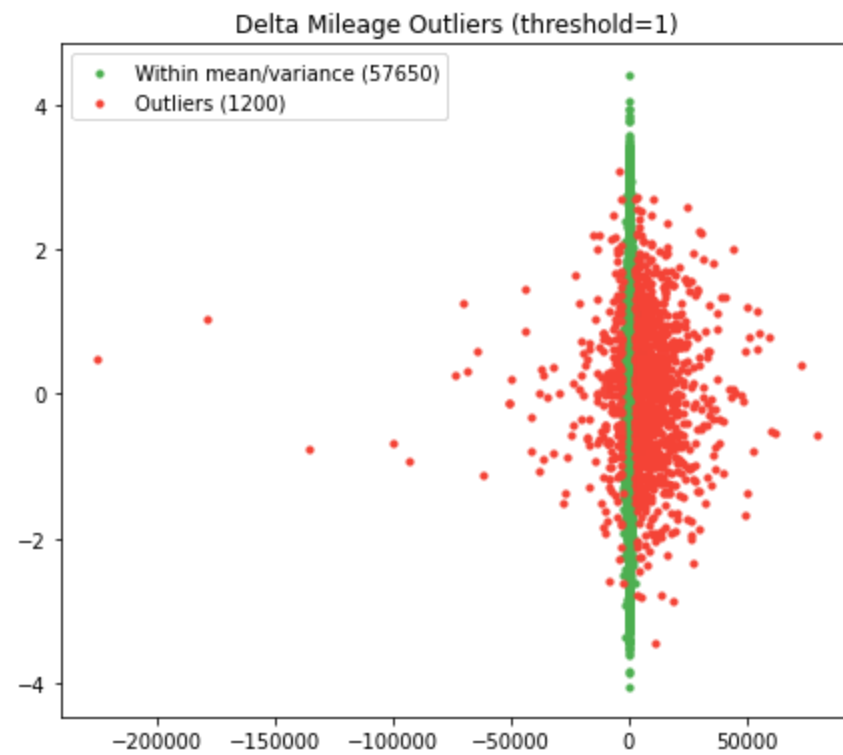




z-score to the rescue

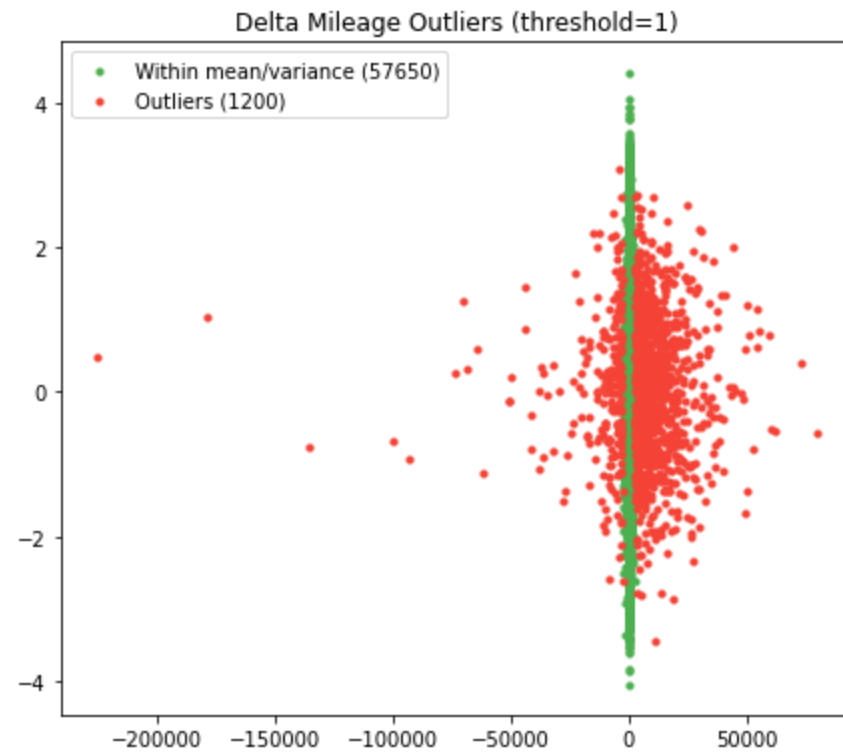
**nothing fancy, just how many std deviations away from mean a point can be,
before it is seen as “outlier”**



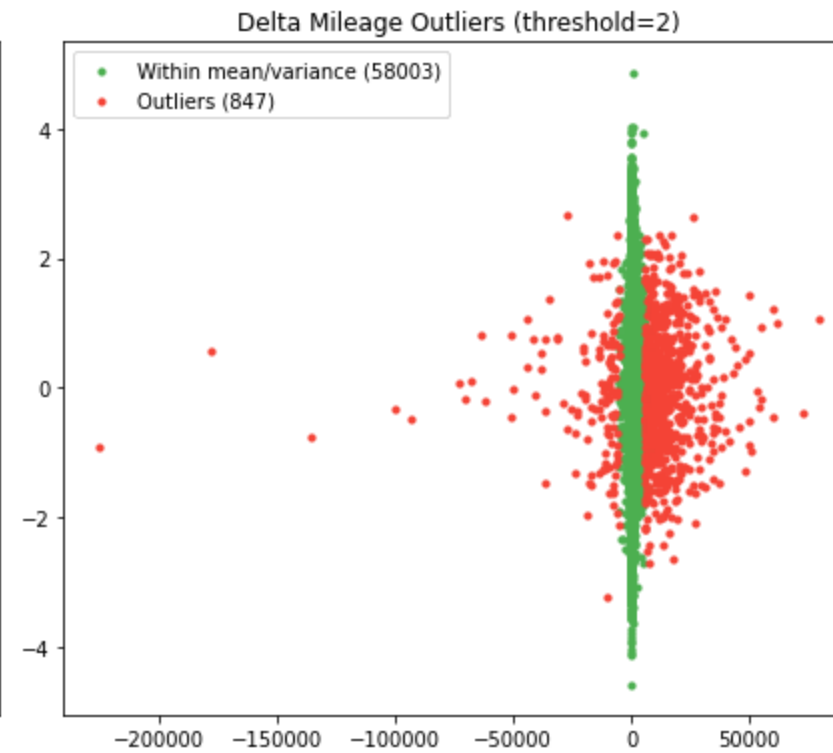


rule of thumb, z-score in [1,3]

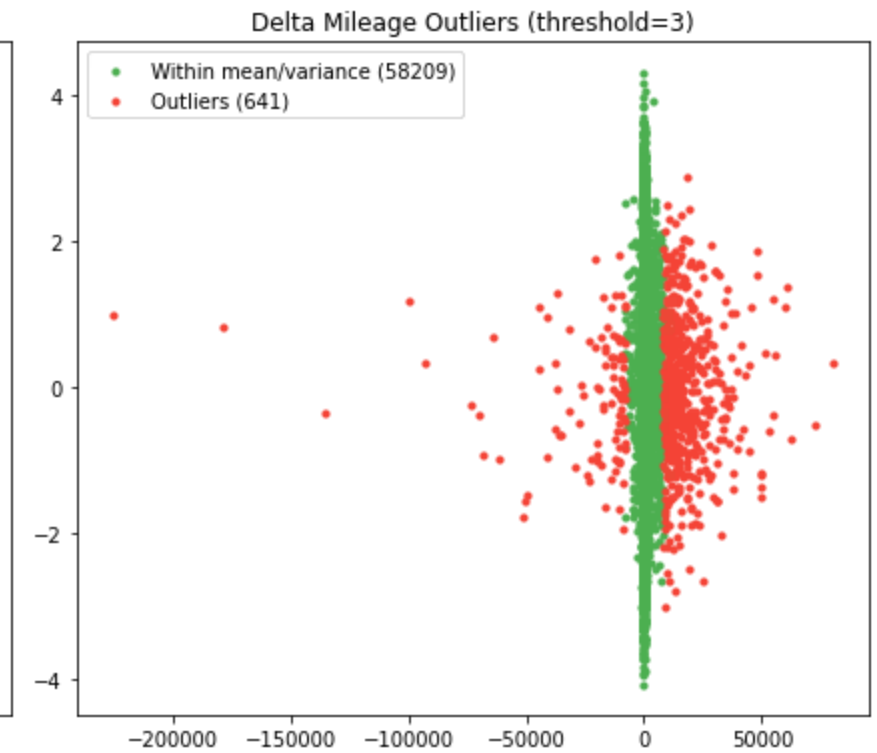
97.96%



98.56%



98.91%



rule of thumb, z-score in [1,3]

quick hypotheses:

- maybe model matches;
- maybe mileage matches;
- **maybe registration year matches;**

quick hypotheses:

- maybe model matches;
- maybe mileage matches;
- **maybe registration year matches;**

99.80% of the cases!

quick hypotheses:

- maybe model matches;
- maybe mileage matches;
- maybe registration year matches;
- **maybe more features match;**

quick hypotheses:

- maybe model matches;
- maybe mileage matches;
- maybe registration year matches;

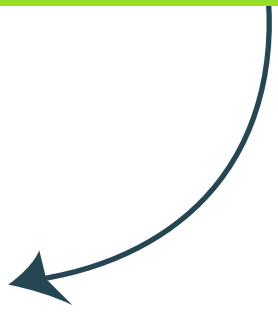
correlation matrix as a rescue tool

quick hypotheses:

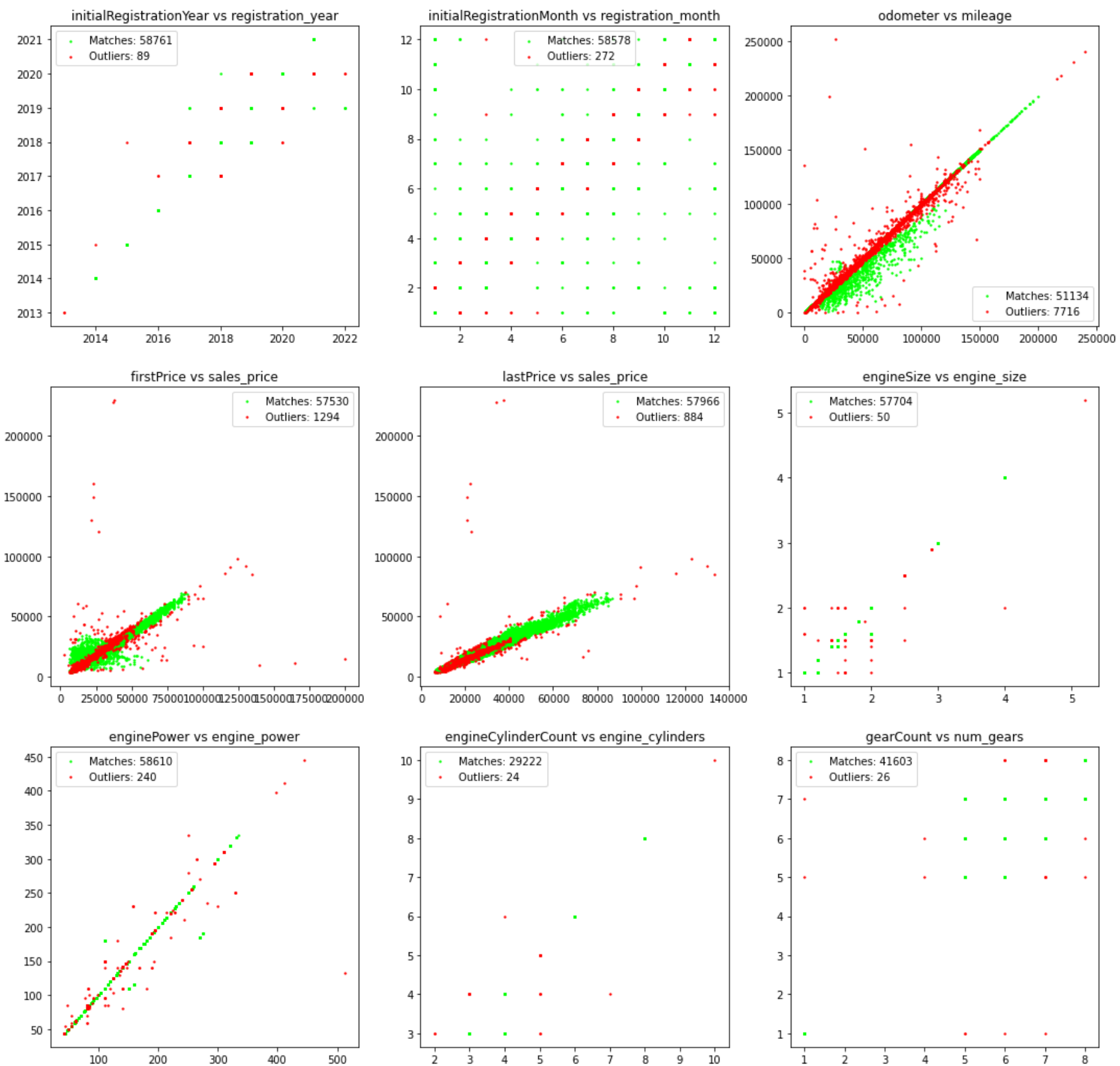
- maybe model matches;
- maybe mileage matches;
- maybe registration year matches;

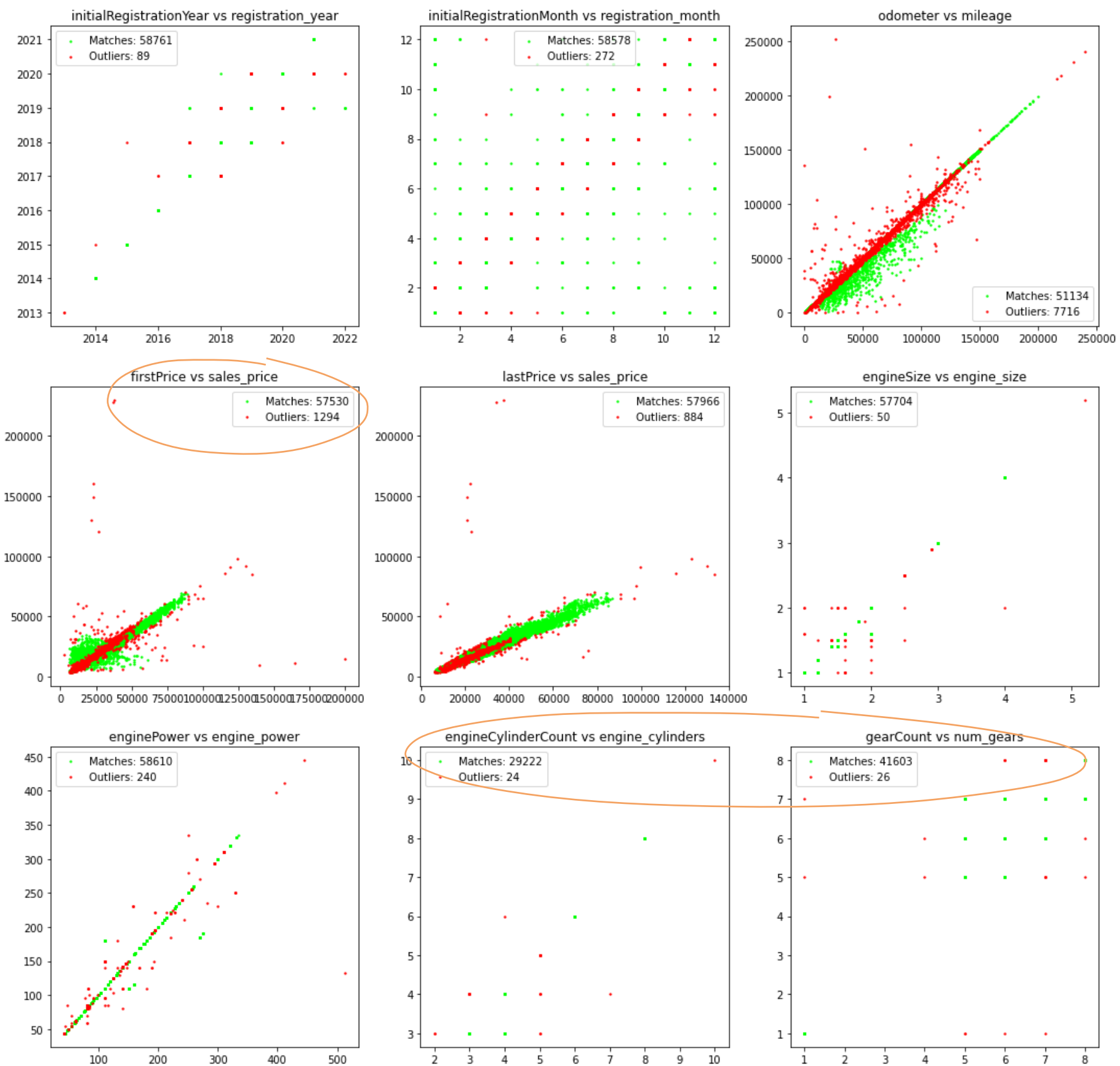
correlation matrix as a rescue tool

*better analyzed with
Spearman method*



B2C		B2B
initialRegistrationYear	99.85%	registration_year
initialRegistrationMonth	99.25%	registration_month
odometer	99.60%	mileage
firstPrice	96.98%	sales_price
lastPrice	97.86%	sales_price
engineSize	99.87%	engine_size
enginePower	99.72%	engine_power
engineCylinderCount	99.51%	engine_cylinders
gearCount	97.63%	num_gears



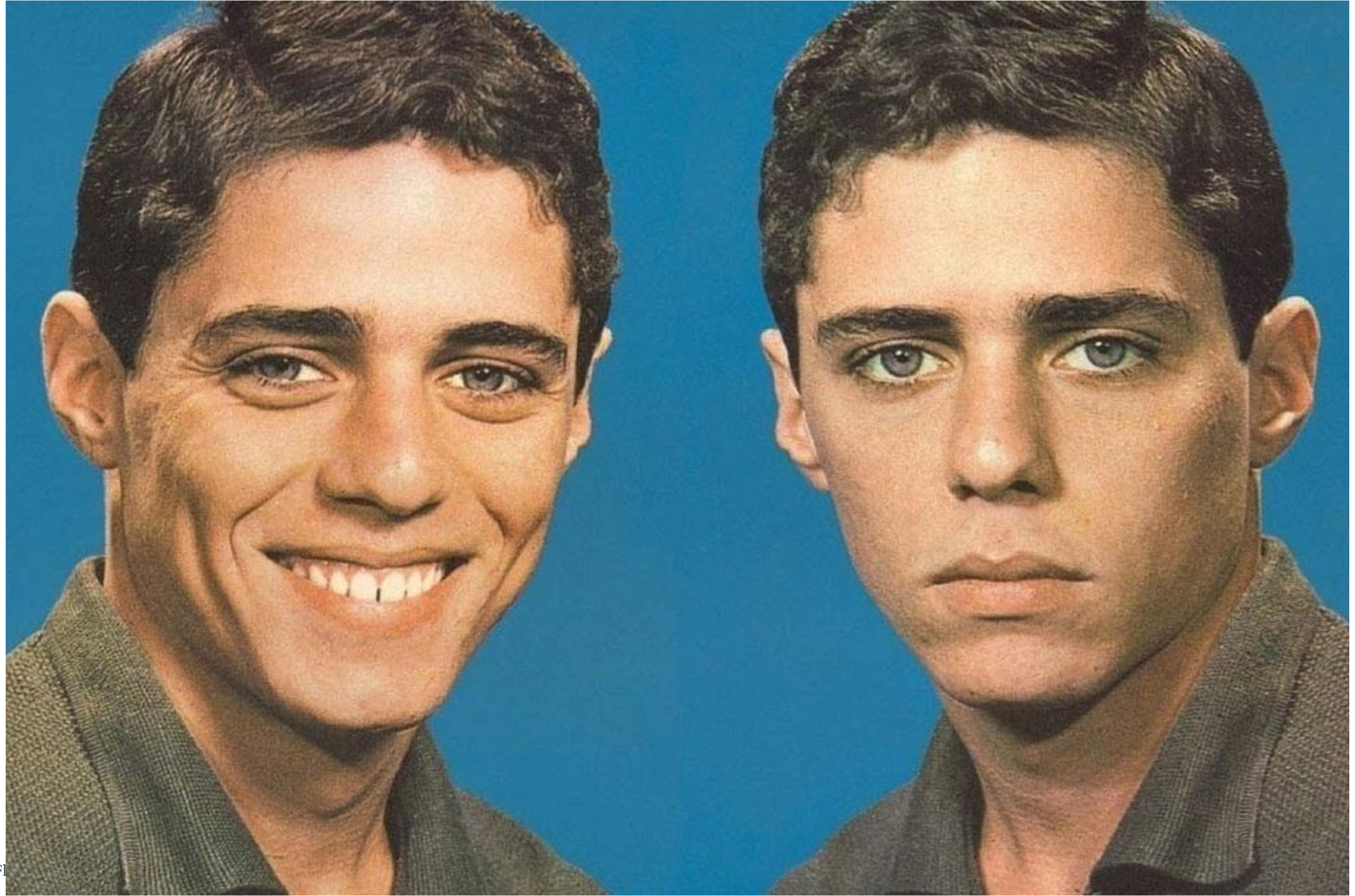


**time for a research
question:**

RQ1: Due to the proximity of these features, can the matched 60K points be reidentified against B2C but without VIN?

RQ1: Due to the proximity of these features, can the matched 60K points be reidentified against B2C but without VIN?

→ after all, we still have records with no label





good idea, bad results.

euclidean distance

euclidean distance

mileage, engine_size, engine_power, sales_price, registration_year, registration_month, model*

stacked B2C and B2B, all normalized (RobustScaler)

Hard Filtered euclidean distance

mileage, engine_size, engine_power, sales_price, registration_year, registration_month, model*

stacked B2C and B2B, all normalized (RobustScaler)

**on a first attempt to
rematch the
“golden dataset”:**

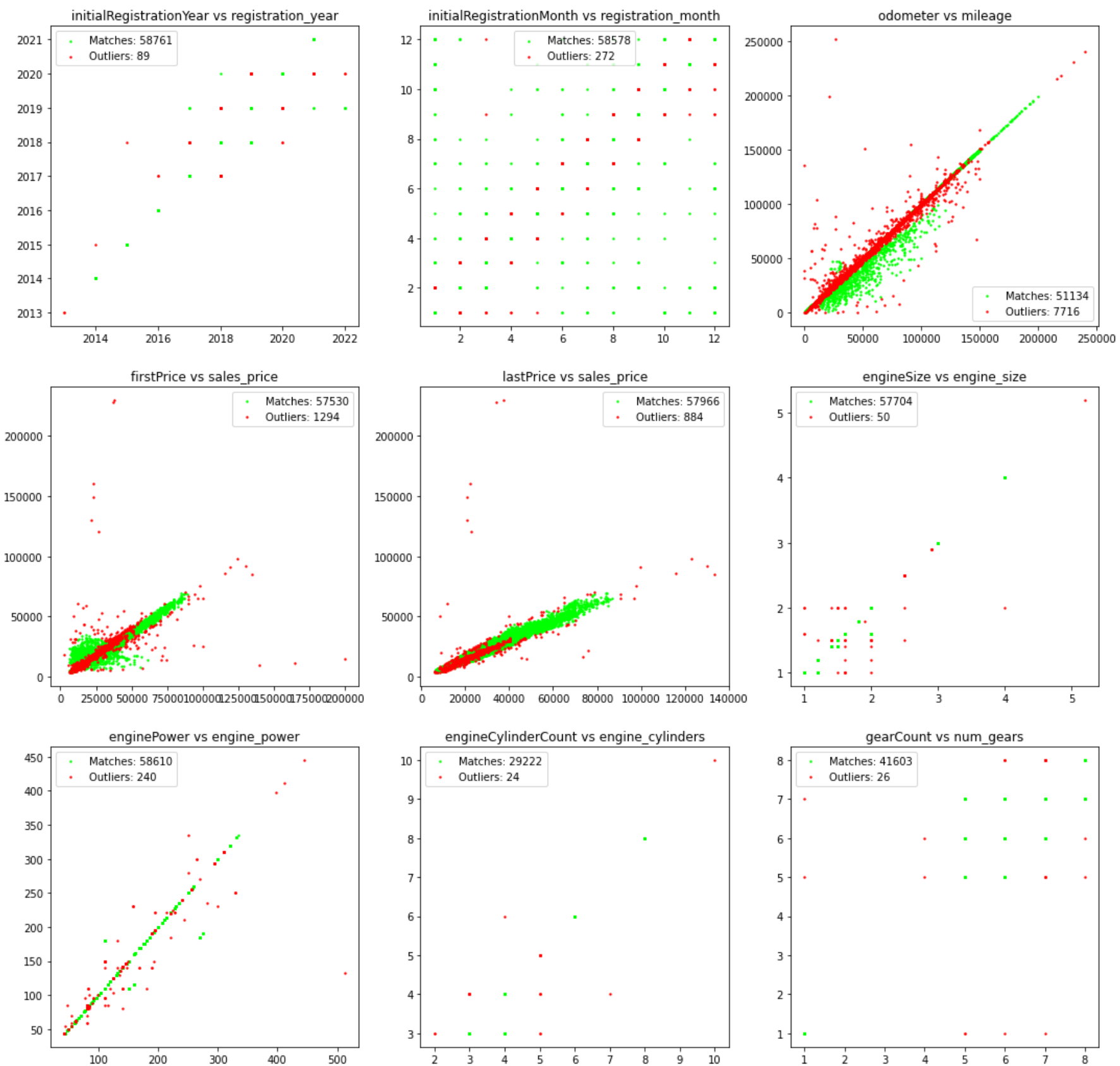
1%

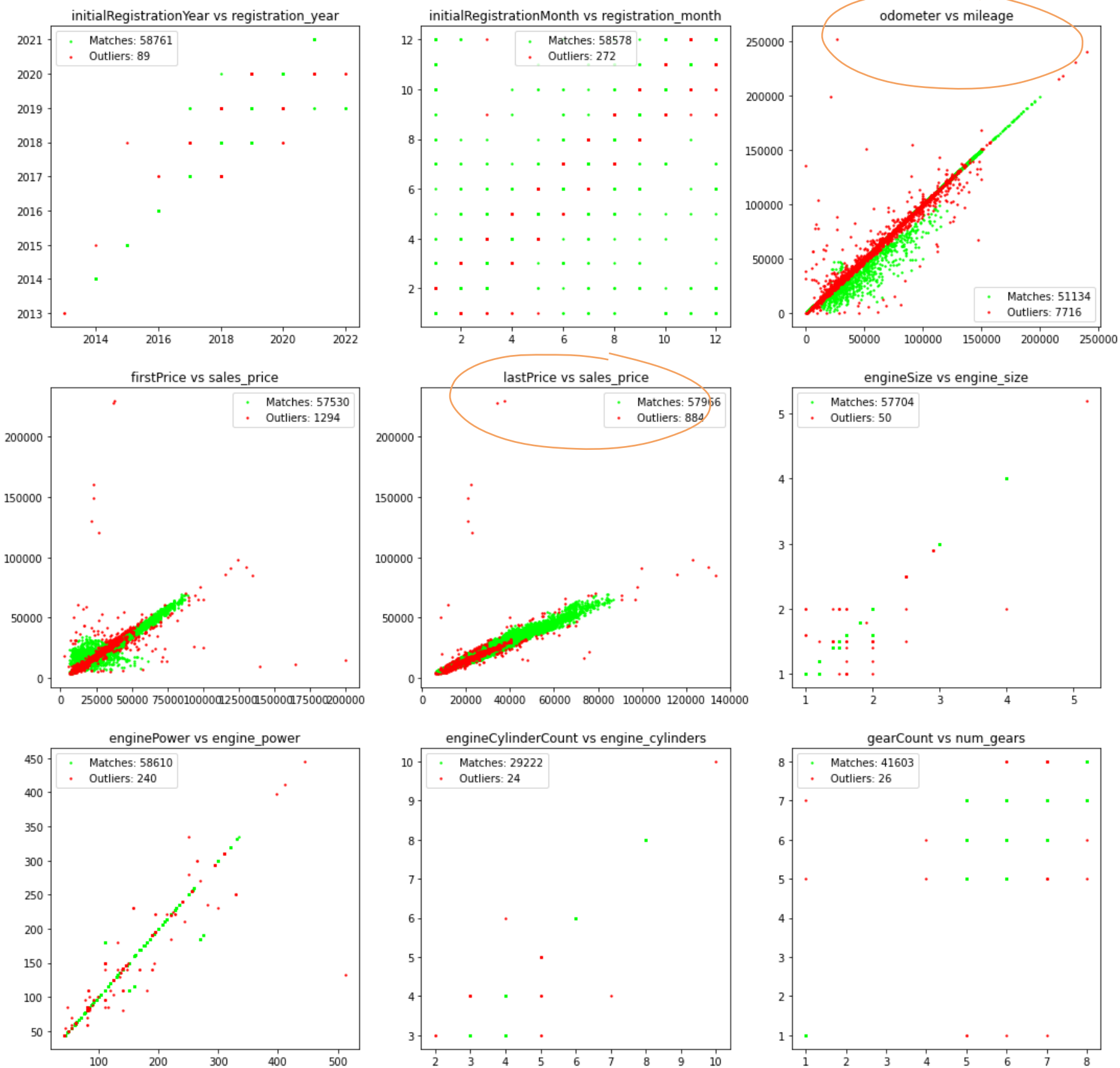
**when searching for closest
point against full B2C dataset**

10%

**when searching for closest
point against 60K B2C dataset**

but wait a minute!





these are continuous features!!!

RQ1: Due to the proximity of these features, can the matched 60K points be reidentified against B2C but without VIN?

RQ1: Due to the proximity of these features, can the matched 60K points be reidentified against B2C but without VIN?

Answer: Apparently not. **(RQ1.1)** But what if we boost the continuous variables to distantiate from the non-matching points?



the incredible Mileage Booster

Hard Filtered euclidean distance

mileage, engine_size, engine_power, sales_price, registration_year, registration_month, model*

stacked B2C and B2B, all normalized (RobustScaler)

Hard Filtered Boosted euclidean distance

MILEAGE, engine_size, engine_power, sales_price, registration_year, registration_month, model*

stacked B2C and B2B, all normalized (RobustScaler)

methodology:

methodology:

- grid search;

methodology:

- grid search;

arbitrary values inside interval [2, 40.000]

methodology:

- grid search;
- **matches;**

methodology:

- grid search;
- **matches;**

against hardly-filtered B2C dataset! (year / month / model)

methodology:

- grid search;
- matches;
- **likelihood;**

methodology:

- grid search;
- matches;
- **likelihood;**

different TOP N rankings, useful to assist on labeling unknown data
(± 200K B2B records)

methodology:

- grid search;
- matches;
- likelihood;
- **acceptance threshold.**

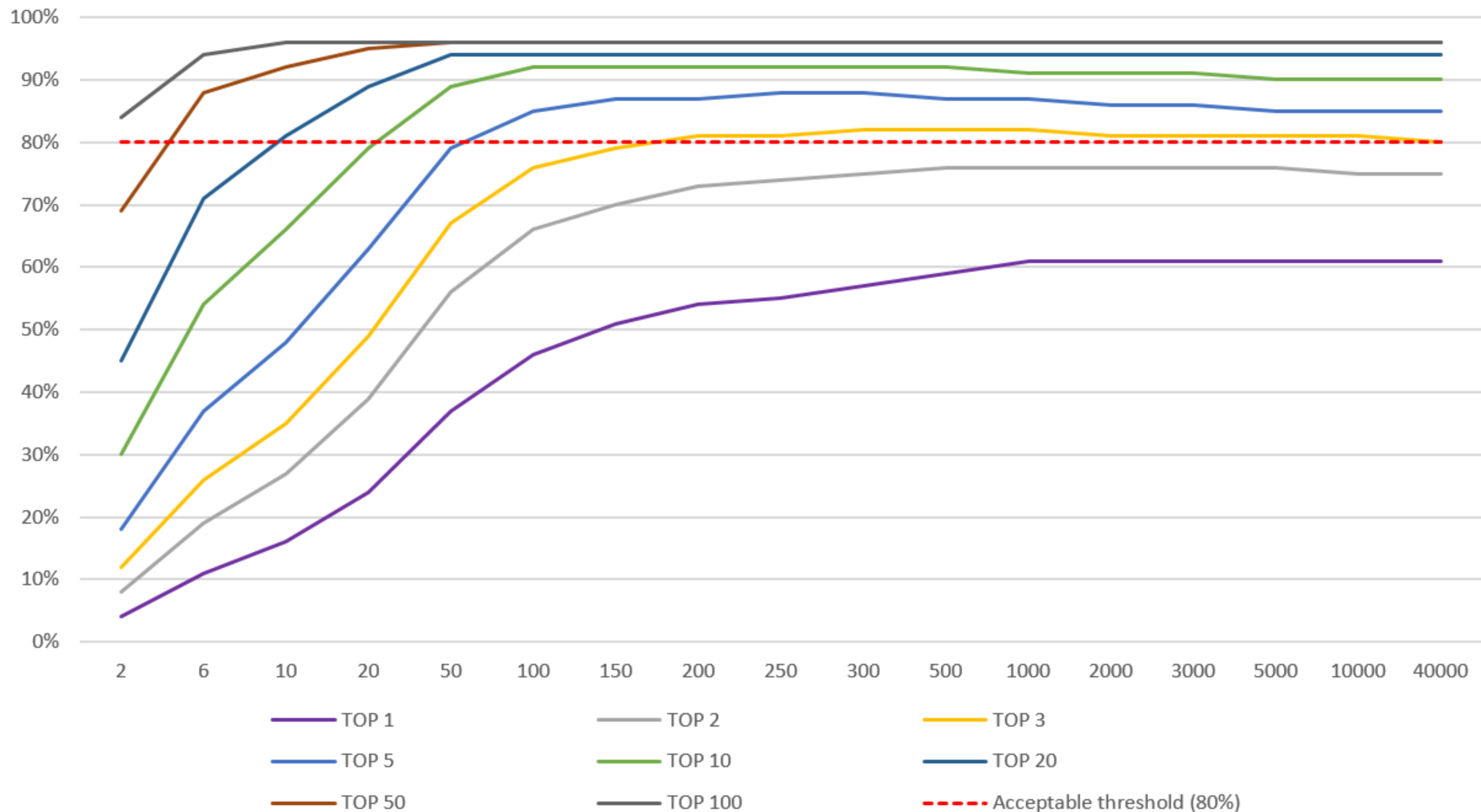
methodology:

- grid search;
- matches;
- likelihood;
- **acceptance threshold.**

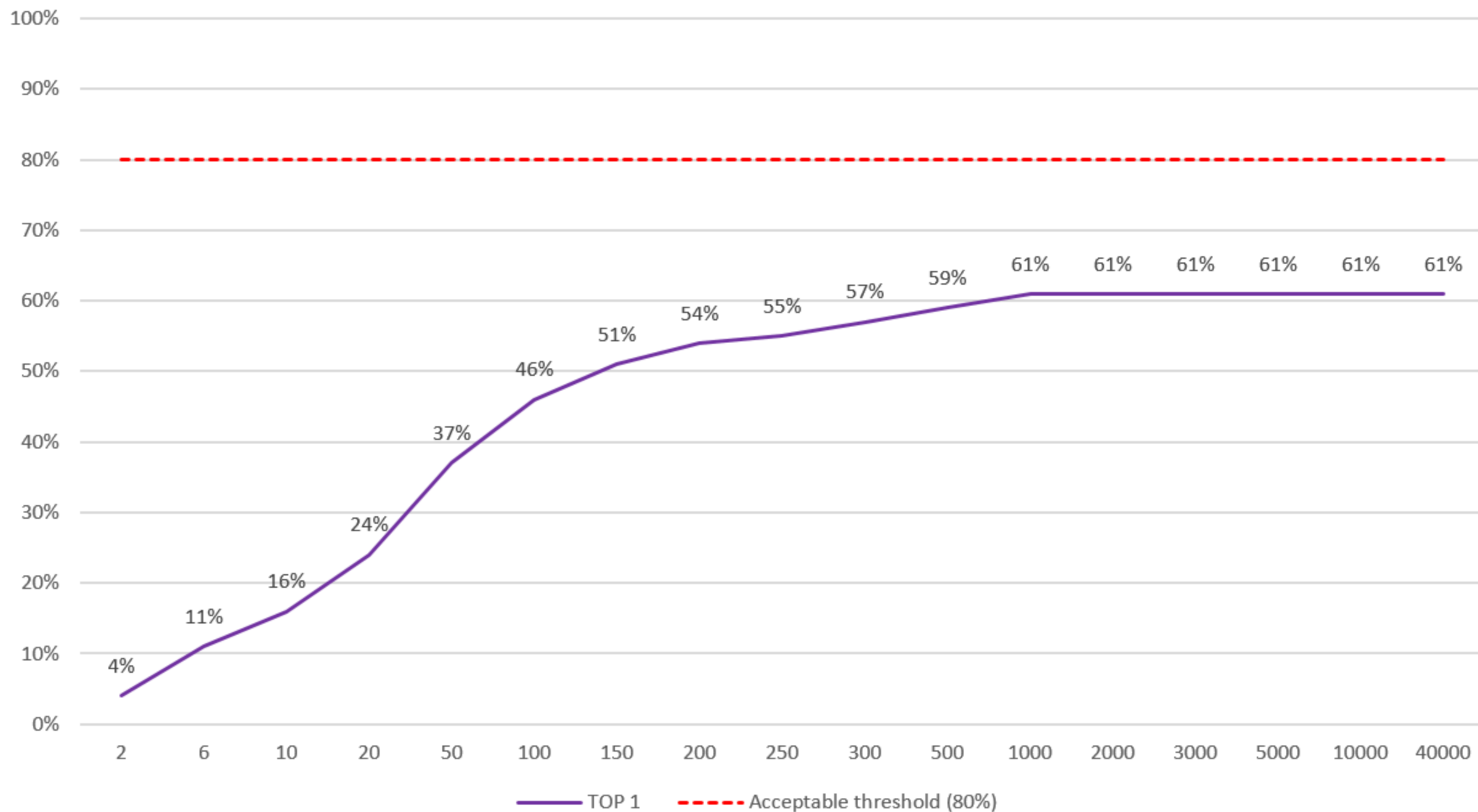
80% as challenge for useful results
(for exact matches, or for likelihood — target in TOP N)

this is quite a lot to process!

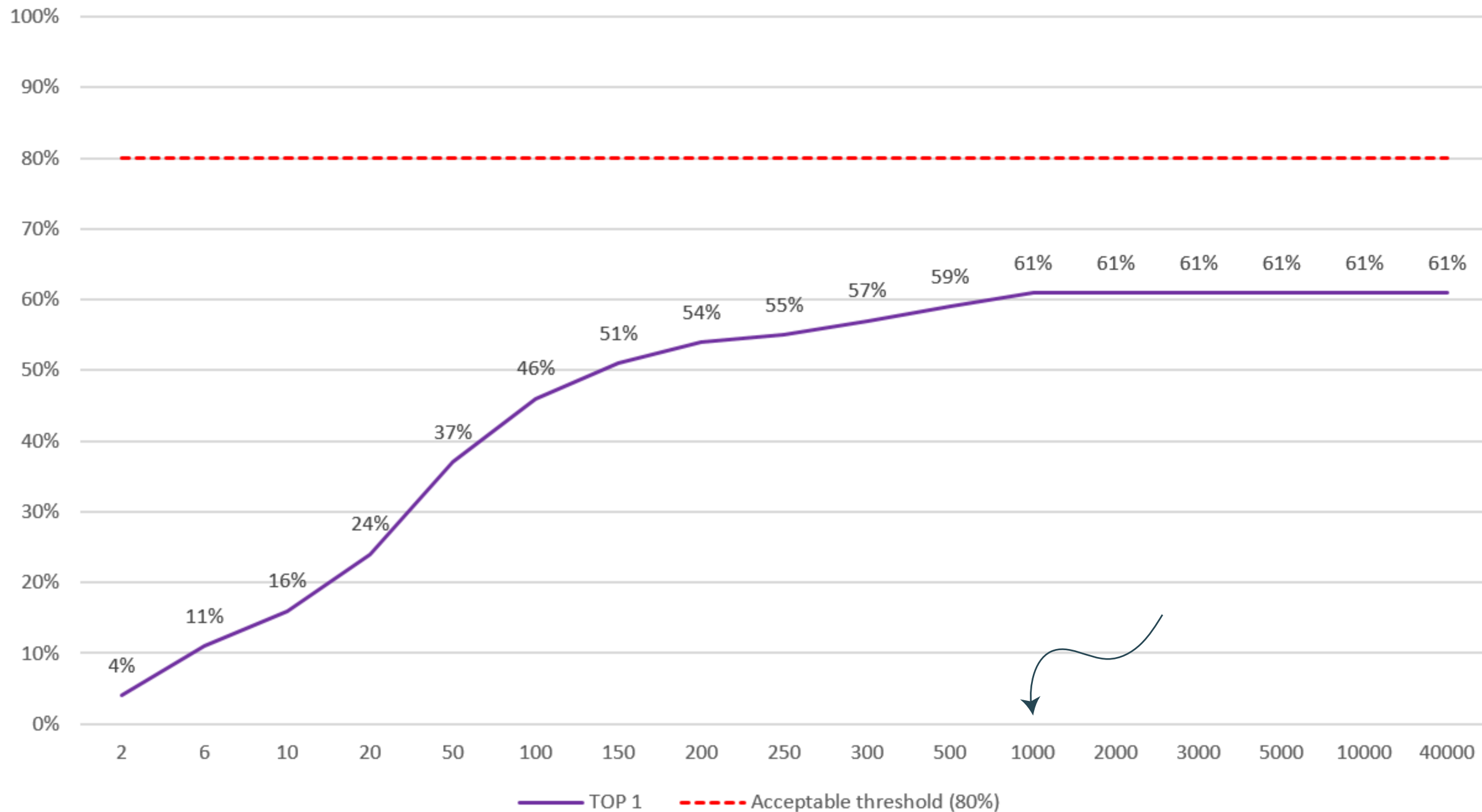
Presence of B2C Target — TOP N over MILEAGE boost factor



Presence of B2C Target — TOP N over MILEAGE boost factor

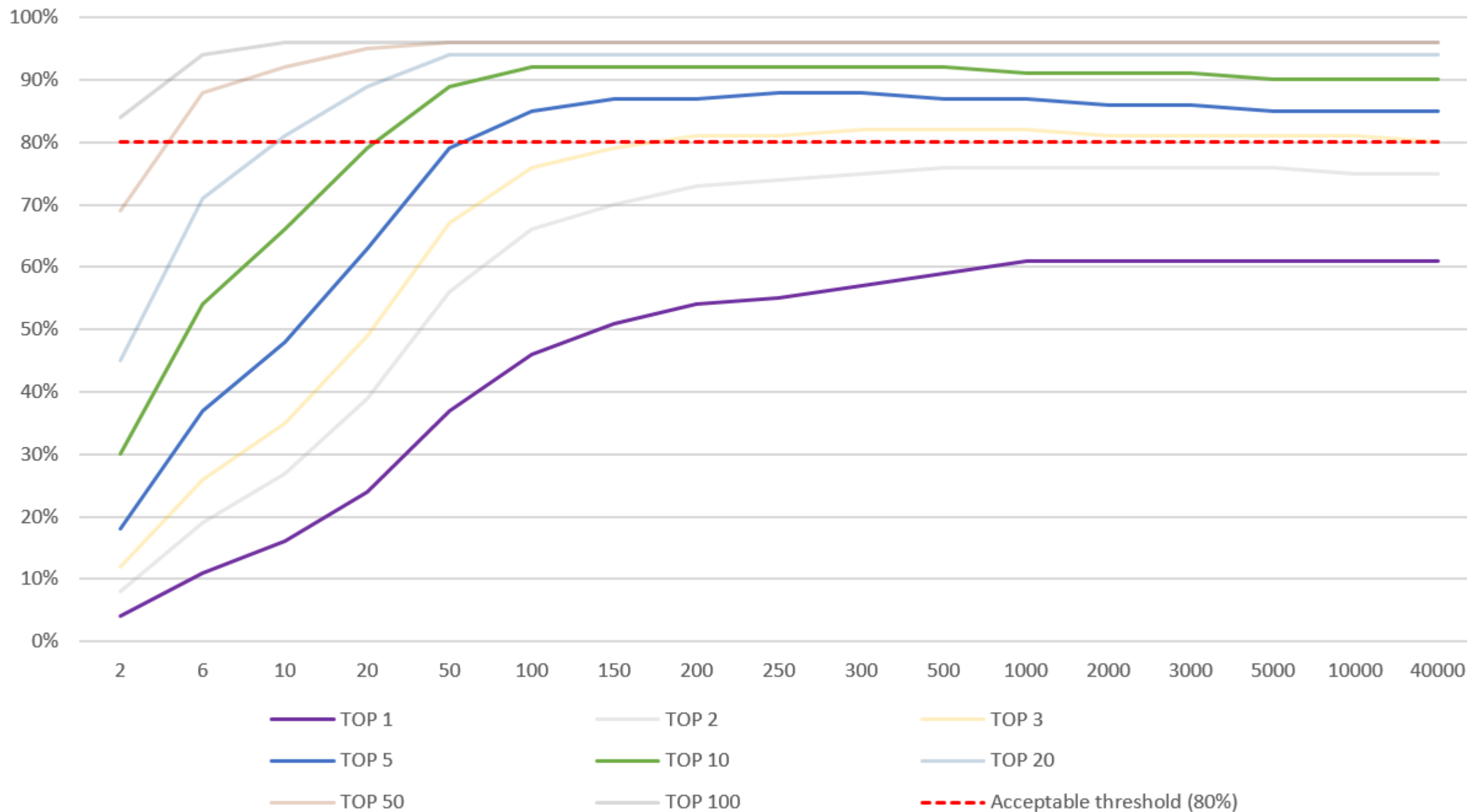


Presence of B2C Target — TOP N over MILEAGE boost factor

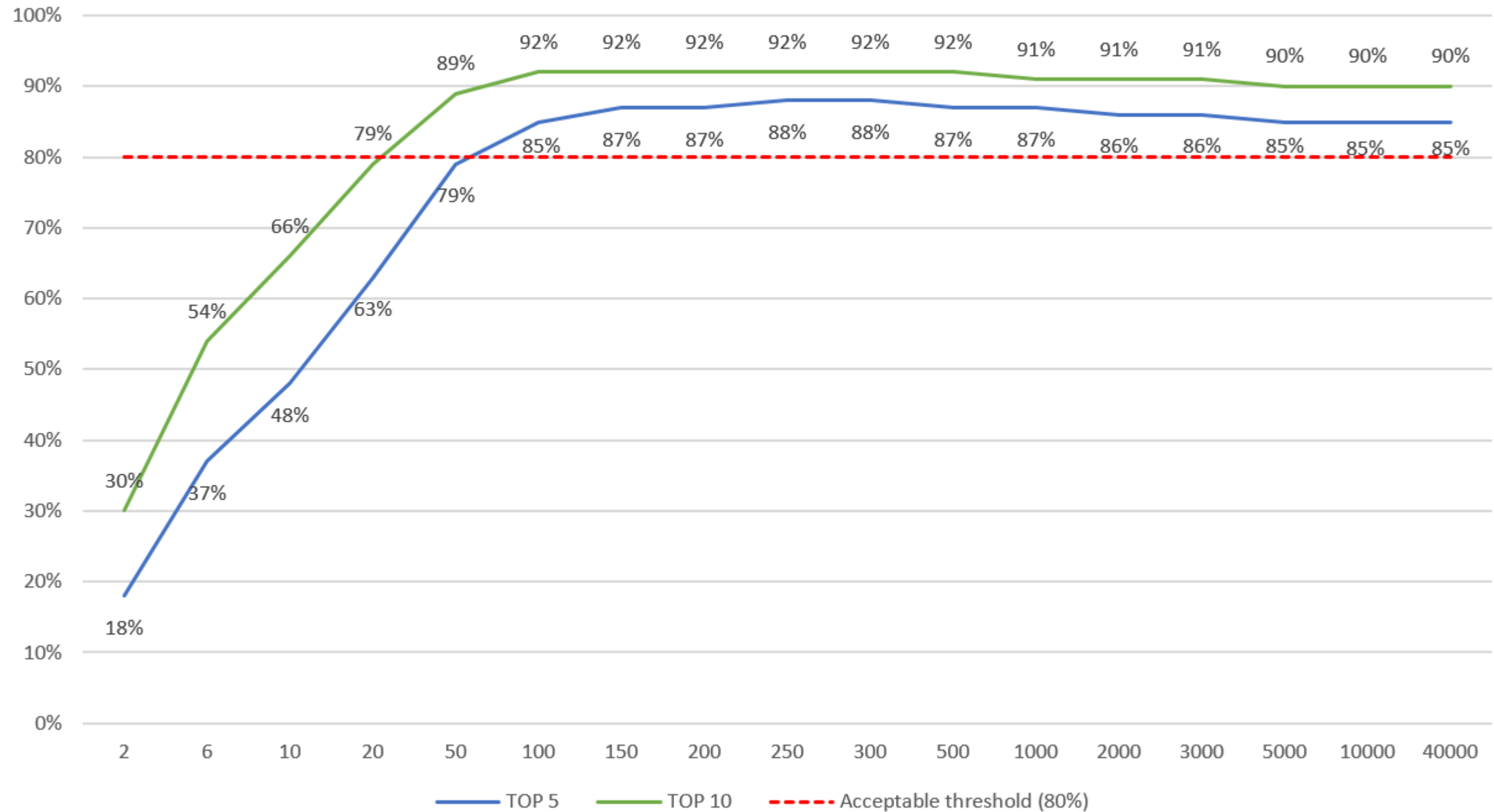


**let's not forget the
amount of records
with zero mileage
delta!**

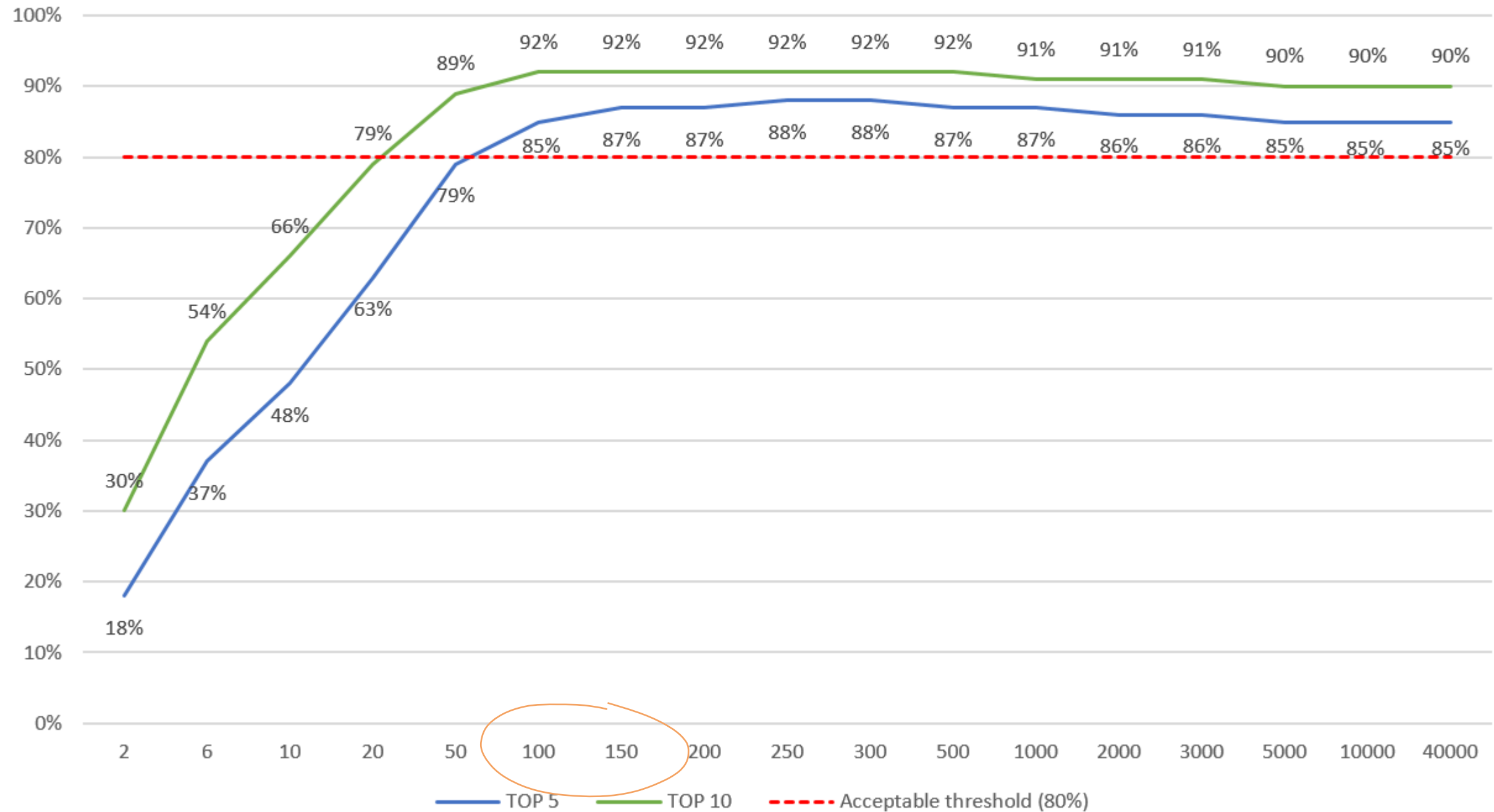
Presence of B2C Target — TOP N over MILEAGE boost factor



Presence of B2C Target — TOP N over MILEAGE boost factor



Presence of B2C Target — TOP N over MILEAGE boost factor



Mileage Boost	TOP 1	TOP 2	TOP 3	TOP 5	TOP 10	TOP 20	TOP 50	TOP 100
2	4%	8%	12%	18%	30%	45%	69%	84%
6	11%	19%	26%	37%	54%	71%	88%	94%
10	16%	27%	35%	48%	66%	81%	92%	96%
20	24%	39%	49%	63%	79%	89%	95%	96%
50	37%	56%	67%	79%	89%	94%	96%	96%
100	46%	66%	76%	85%	92%	94%	96%	96%
150	51%	70%	79%	87%	92%	94%	96%	96%
200	54%	73%	81%	87%	92%	94%	96%	96%
250	55%	74%	81%	88%	92%	94%	96%	96%
300	57%	75%	82%	88%	92%	94%	96%	96%
500	59%	76%	82%	87%	92%	94%	96%	96%
1000	61%	76%	82%	87%	91%	94%	96%	96%
2000	61%	76%	81%	86%	91%	94%	96%	96%
3000	61%	76%	81%	86%	91%	94%	96%	96%
5000	61%	76%	81%	85%	90%	94%	96%	96%
10000	61%	75%	81%	85%	90%	94%	96%	96%
40000	61%	75%	80%	85%	90%	94%	96%	96%

Mileage Boost	TOP 1	TOP 2	TOP 3	TOP 5	TOP 10	TOP 20	TOP 50	TOP 100
2	4%	8%	12%	18%	30%	45%	69%	84%
6	11%	19%	26%	37%	54%	71%	88%	94%
10	16%	27%	35%	48%	66%	81%	92%	96%

short info: boosting sales_price did not prove successful :-)

3000	61%	76%	81%	86%	91%	94%	96%	96%
5000	61%	76%	81%	85%	90%	94%	96%	96%
10000	61%	75%	81%	85%	90%	94%	96%	96%
40000	61%	75%	80%	85%	90%	94%	96%	96%

Mileage Boost	TOP 1	TOP 2	TOP 3	TOP 5	TOP 10	TOP 20	TOP 50	TOP 100
2	4%	8%	12%	18%	30%	45%	69%	84%
6	11%	19%	26%	37%	54%	71%	88%	94%
10	16%	27%	35%	48%	66%	81%	92%	96%

**short info: boosting
sales_price did not
prove successful :-)**

what if we “degraded” this feature, then? (idea)

40000

61%

75%

80%

85%

90%

94%

96%

96%

**pause for a
quick breath...**

**can we sketch a
pipeline?**

Pipeline



Pipeline

- TOP 1 == VIN;



Pipeline

- TOP 1 == VIN;
- **Reduce dataset for non-matching VINS;**



**of course, the set of
features can be
revisited!**

Pipeline

- TOP 1 == VIN;
- Reduce dataset for non-matching VINS;
- **Use boosted distance to output likelihood — TOP N;**



Pipeline

- TOP 1 == VIN;
- Reduce dataset for non-matching VINS;
- **Use boosted distance to output likelihood — TOP N;**

focusing on statistical approach, no complex n.n.



Pipeline

- TOP 1 == VIN;
- Reduce dataset for non-matching VINS;
- Use boosted distance to output likelihood — TOP N;
- **Input likelihood on a second (or more) model(s): best probable B2C match.**



Pipeline

- TOP 1 == VIN;
- Reduce dataset for non-matching VINS;
- Use boosted distance to output likelihood — TOP N;
- **Input likelihood on a second (or more) model(s): best probable B2C match.**

*maybe Filtering mileage? or extra
parameters for other features*



Pipeline

- TOP 1 == VIN;
- Reduce dataset for non-matching VINS;
- Use boosted distance to output likelihood — TOP N;
- **Input likelihood on a second (or more) model(s): best probable B2C match.**

also, dealer info can
be handy



to be continued...



Further superheroes to be unlocked...

~~to be continued...~~

