



data
iku

Unsupervised Anomaly Detection with Advanced Analytics: your next steps

London Meetup 24/1/2018



Speaker Bio: Harizo Rajaona

Data Scientist @Dataiku

Harizo has been working at Dataiku for one year, working with customers from financial services and CPG. Before joining Dataiku, he worked at the French Alternative Energies and Atomic Energy Commission, where he focused on designing statistical methods to reconstruct air pollution sources with operational applications in accidental cases and in emission monitoring contexts. He holds a PhD in Mathematics from University Lille 1.



data
iku

Unsupervised anomaly detection

Jan.24th 2018 • Harizo Rajaona (Data Scientist)

Analytics & Data science meetup - London

Outline



Summary &
conclusion



Which algorithms? How to
benchmark them?



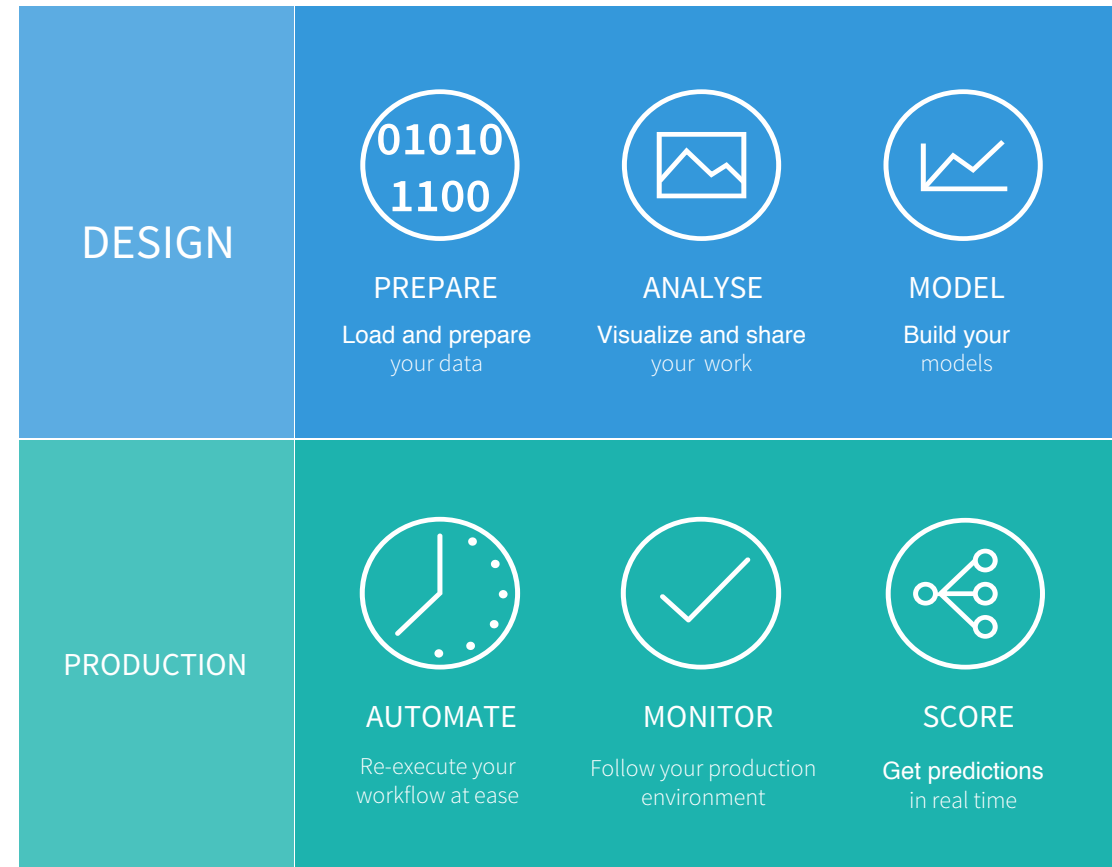
How to properly
report outliers ?



Introduction: context,
supervised vs. unsupervised

About Dataiku

- Founded in 2013 (5y. birthday party 1 week ago!)
- 110+ employees, 120+ clients
- Paris, NYC, London, Munich



Introduction: context, supervised vs. unsupervised

Anomaly detection (AD) in the wild

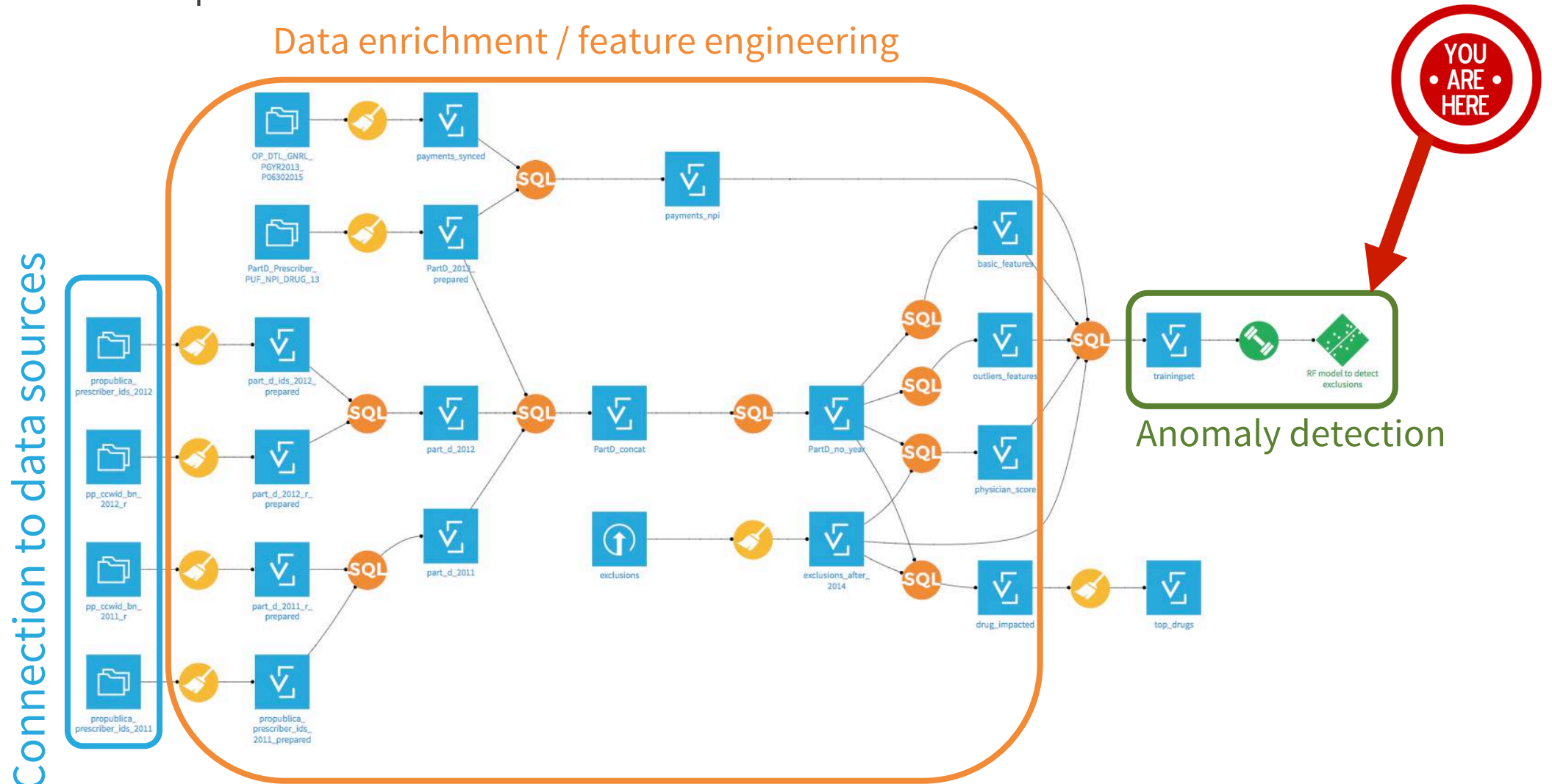
Some examples

- **Fraud detection:**
 - Detecting fraudulent credit card transactions
 - Flagging suspicious warranty/service claims
- **Predictive maintenance:**
 - Tracking a machine's abnormal behaviour in a manufacturing plant
- **Healthcare:**
 - Looking for unusual/unobserved symptoms in a group of patients
- **Network monitoring:**
 - Identify non-idle traffic surges/drops
 - Trigger alerts in case of breaches/intrusions



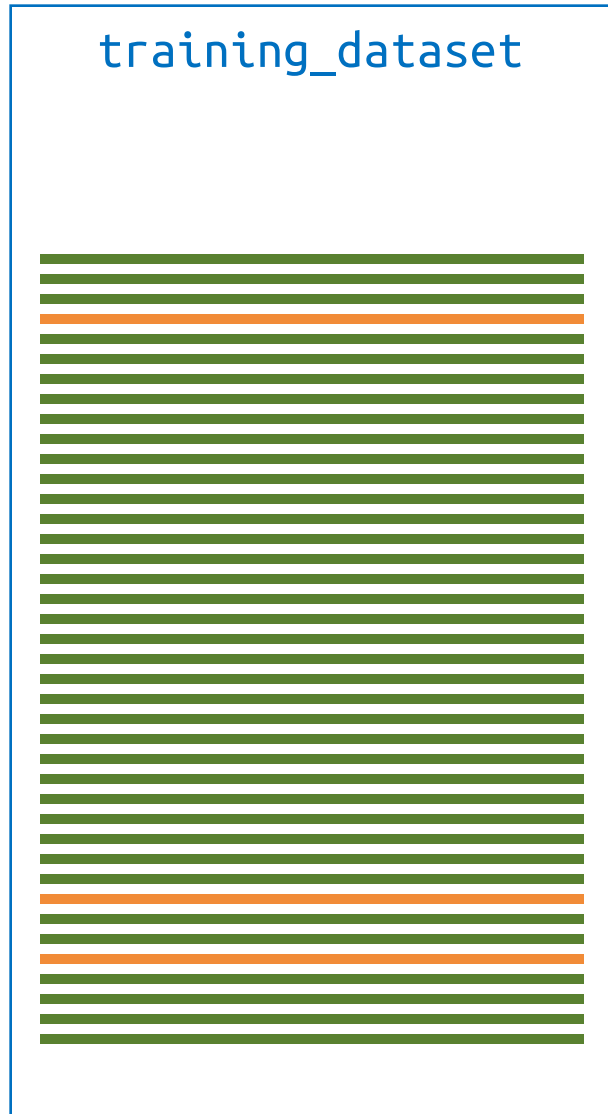
AD: one step in the data science workflow

Data enrichment / feature engineering



(Example of a flow in Dataiku DSS for Medicare fraud detection)

Definitions & the supervised case



Inliers:

- Vast majority of the data
- Represent an “expected behaviour”

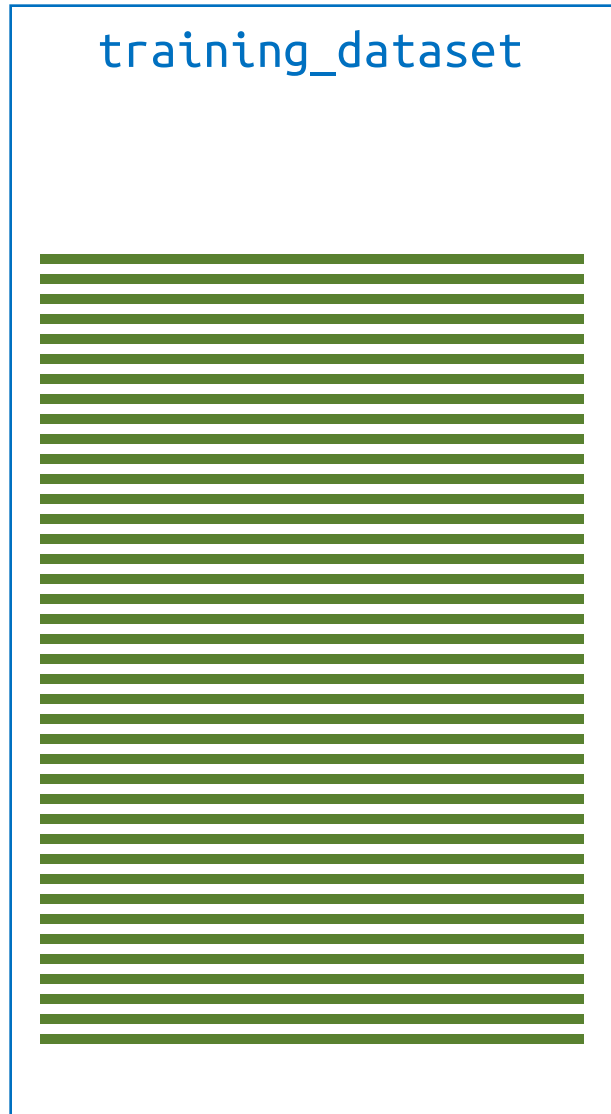
Outliers:

- (Very!) small portion of the data
- Abnormal cases

Supervised AD: data is labelled (inlier/outlier)

- Binary classification problem
- Account for high class imbalance
- Cross-validate carefully (stratified splits)

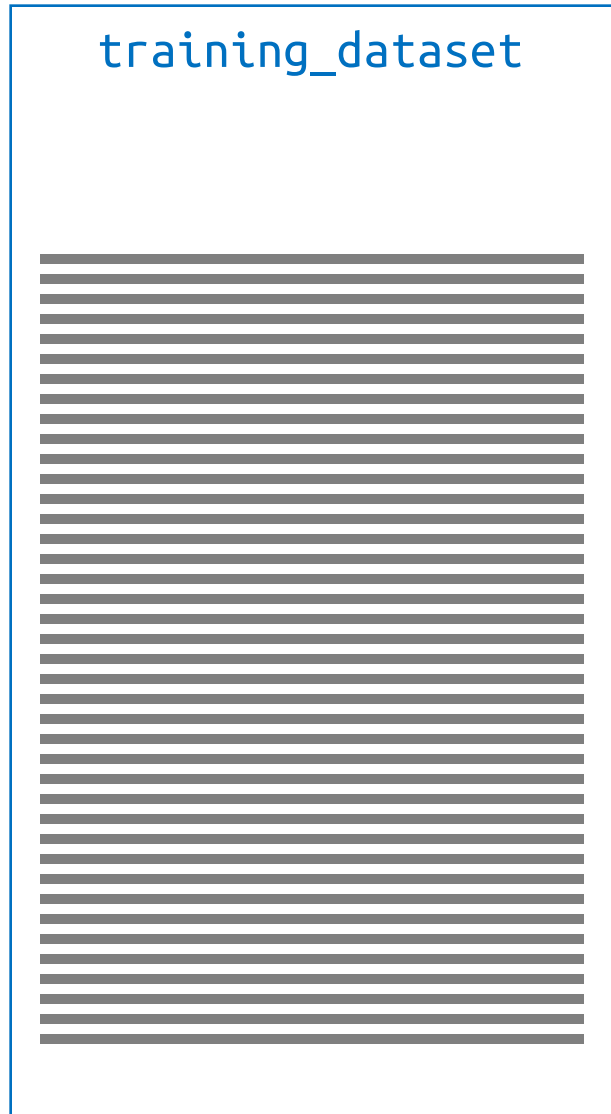
The semi-supervised case



Semi-supervised AD: data contains only inliers

- “Learning normality”(one-class)
- Outliers that appear *at test time* are flagged (novelty detection).

The unsupervised case

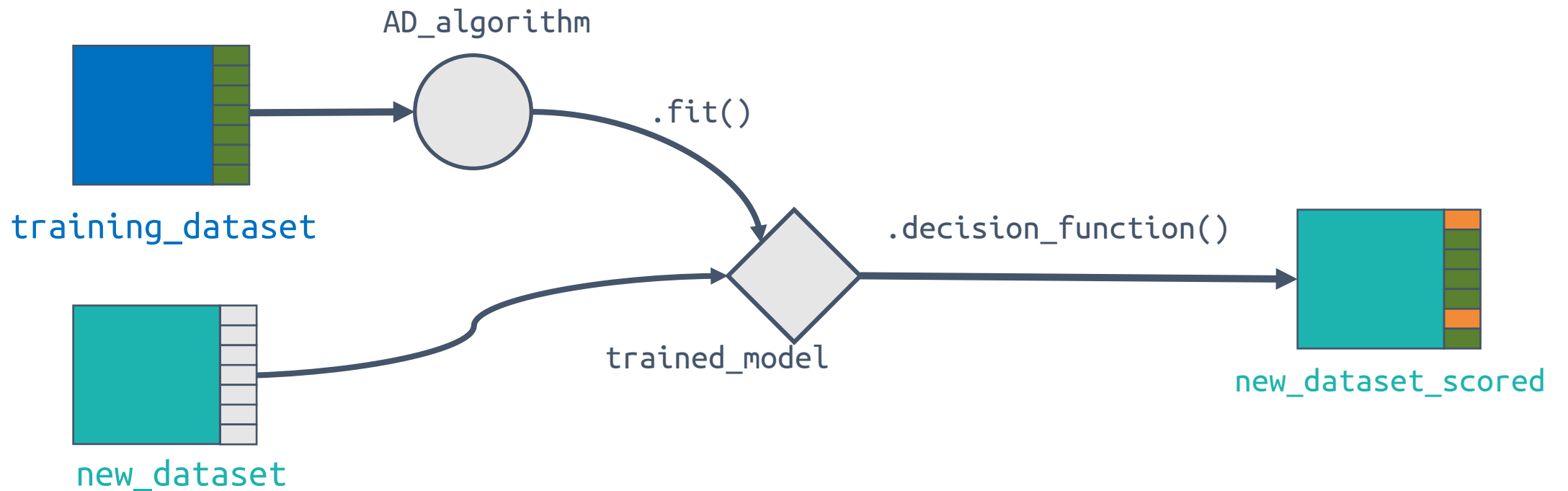


Unsupervised AD (UAD) / outlier detection:

- No labels at all
- Data contains a **sufficiently small amount of outliers**
- Happens quite often !
 - Costly label processes
 - Unreliable labels
 - Potential overfitting on existing sample

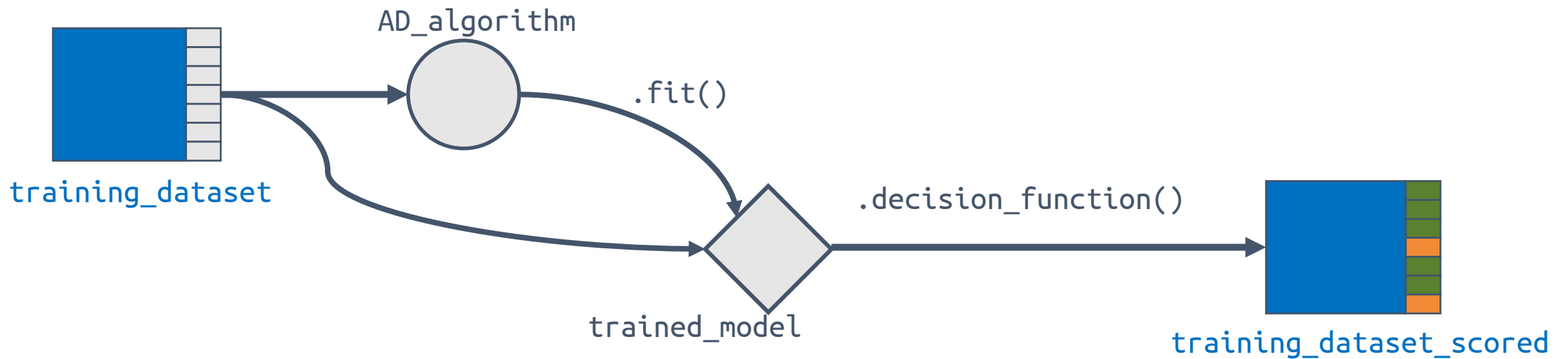
Anomaly detection workflows

Semi-supervised



Anomaly detection workflows

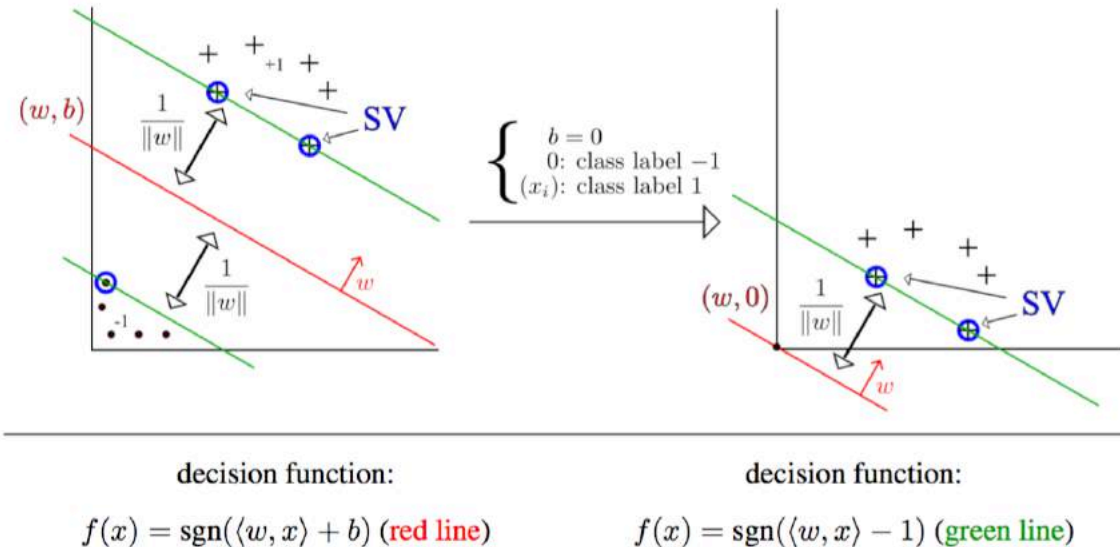
Fully unsupervised



Which algorithms? How to benchmark them?

One-class support vector machines (OCSVM)

A kernel-based method

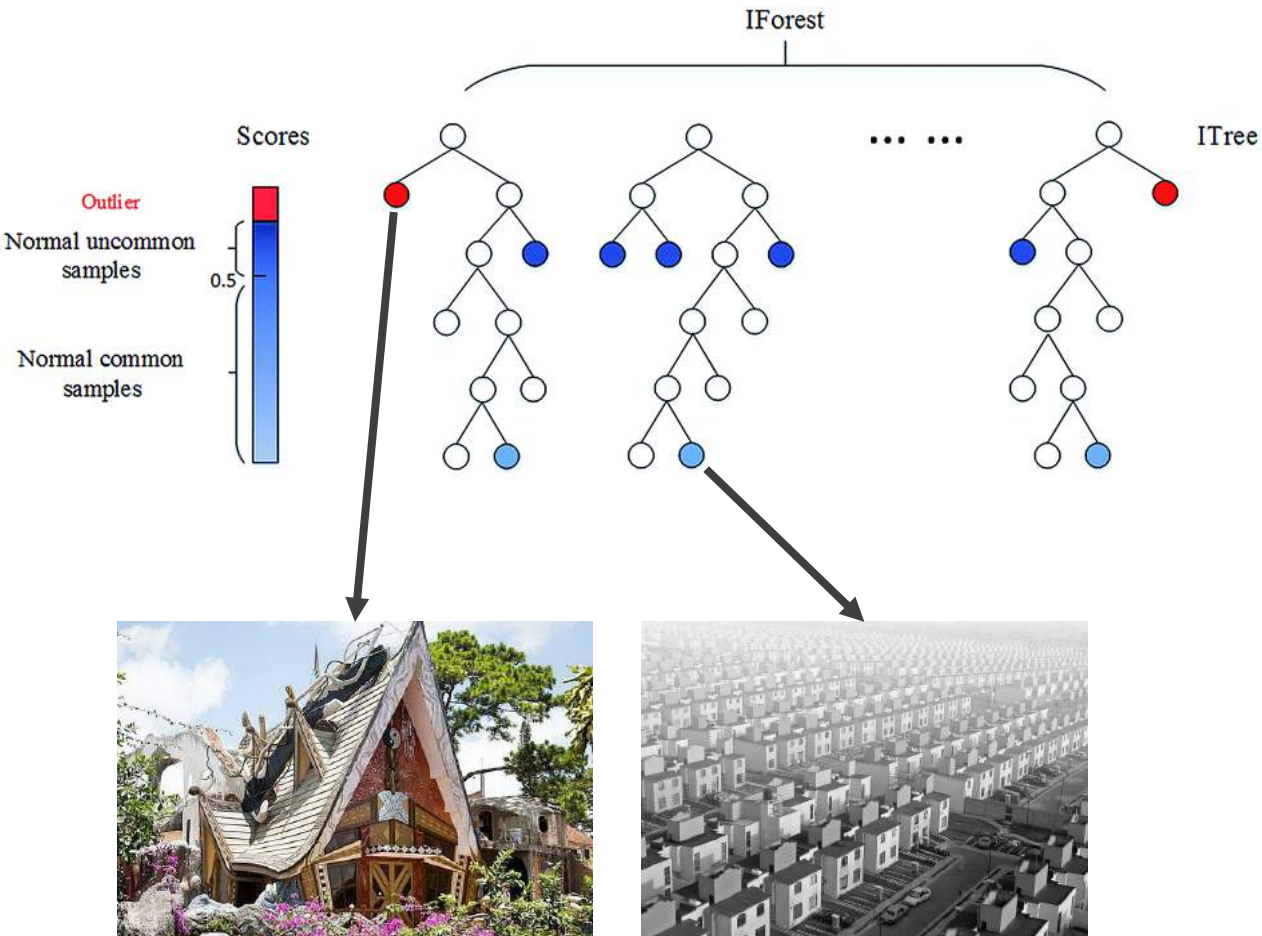


(Machine learning and extremes for anomaly detection, N. Goix, 2016)

- SVM classification problem: find the optimal separation between classes with maximum margin rule
- One-class SVM: find the optimal separation **between the inlier class and the origin** with maximum margin rule
- Can handle highly non-linear cases
- Non-trivial settings (kernel, hyperparams), struggles with high dimensionality

Isolation forest

A tree-based method

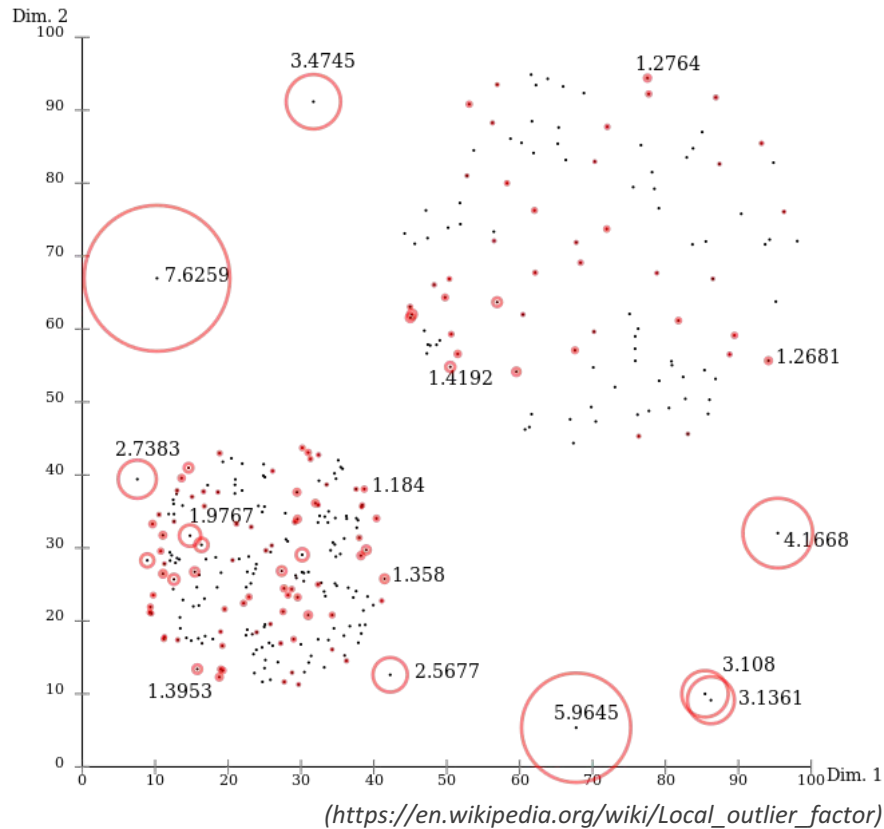


- Ensemble of randomized decision trees (random features and splits)
- Outliers have the **shortest average path length** over the whole set of trees
- Intuitive interpretation
- Scalable implementation (parallelizable)

(Representative subset selection and outlier detection via isolation forest, Chen et al, 2016)

Local outlier factor (LOF)

A distance-based method

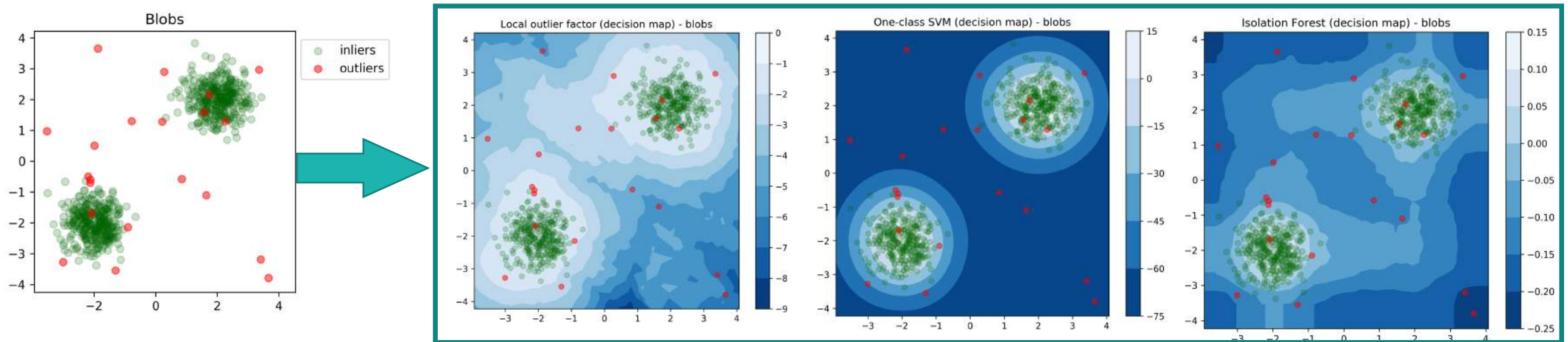


- Compares the local density of a point with the densities of its neighbors.
- « If I'm lonely and my neighbors are too, then I'm normal. »
- Allows to define relative outliers which **may not reflect a totally abnormal state**
- Distance-based -> computation load increases with dimensionality
- Full pass over dataset at evaluation time (costly)

Performance evaluation (simple dataset)

How well are we doing ?

- To evaluate the performance, we actually need labels (that are **NOT** used at train time)
- Decision maps (only for simple 2D cases): apply `.decision_function()` to a grid over the feature space



- Not appropriate for more complex datasets !

Performance evaluation (real datasets)

How well are we doing ?

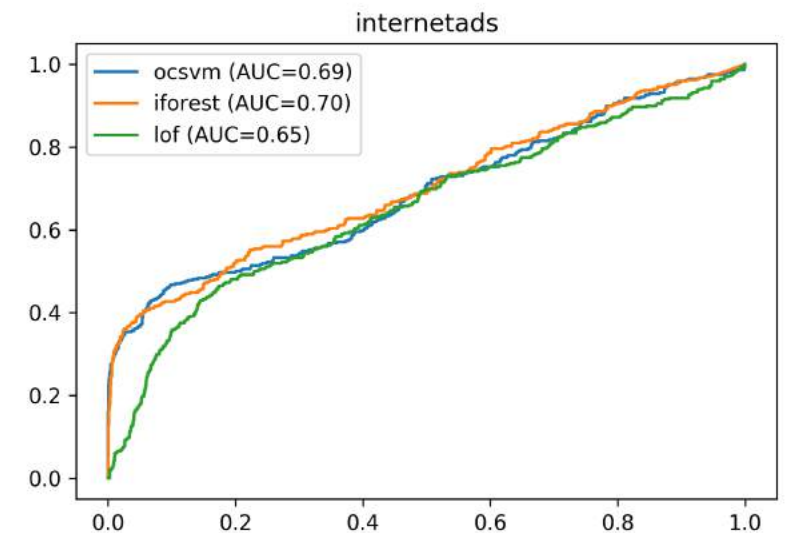
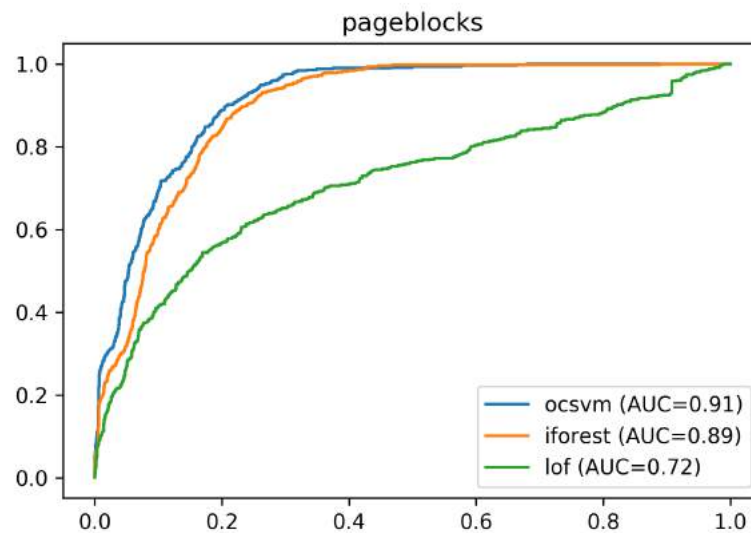
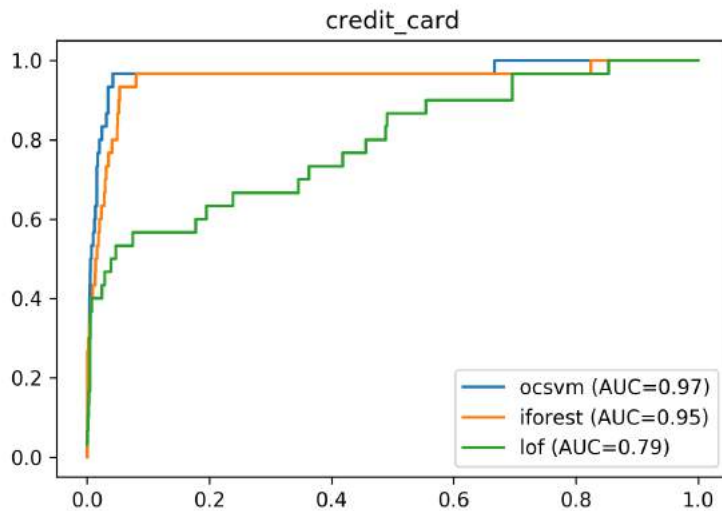
- Datasets from Kaggle (credit card fraud detection) and academic benchmarks (internetads, pageblock)
- Cover some configurations often met in « real-life » datasets (categorical features, high dimensionality)

Dataset	#rows	#columns	Outlier ratio	Feature types
Credit card (Kaggle), 6% sample	17088	29	0.002	Numerical
pageblocks	4883	510	0.095	Numerical
internetads	1966	1554	0.187	Binary

Performance evaluation (real datasets)

How well are we doing ?

- Receiver Operator Characteristic (ROC) curve:
 - Same logic as a supervised binary classification problem.
 - Labels are available but **NOT** used at training time!
 - Area Under the Curve (AUC) provides a numerical metric of performance



Performance evaluation (real datasets)

How well are we doing?

- Computation times matter as well !
- Fit / score times in seconds:

Dataset (shape)	IsolationForest(n_jobs=-1)	OneClassSVM	LocalOutlierFactor
Credit_card (17088, 29)	1.33 0.65	19.31 6.92	32.56 32.99
Pageblocks (4883, 510)	0.51 0.16	0.90 0.41	0.27 0.24
Internetads (1966, 1554)	3.26 1.80	5.64 3.64	10.38 10.30

How to properly report outliers?

Anomaly scoring

One does not simply **predict_proba()**

- Anomaly scores are purely mathematical and not easily interpretable!

$$s(x, \psi) = 2^{\frac{-\mathbb{E}[h(x)]}{c(\psi)}}$$

(isolation forest)

$$\text{lof}_k(x) = \frac{\sum_{o \in N_k(x)} \left(\frac{\text{lrd}_k(o)}{\text{lrd}_k(x)} \right)}{|N_k(x)|}$$

(local outlier factor)

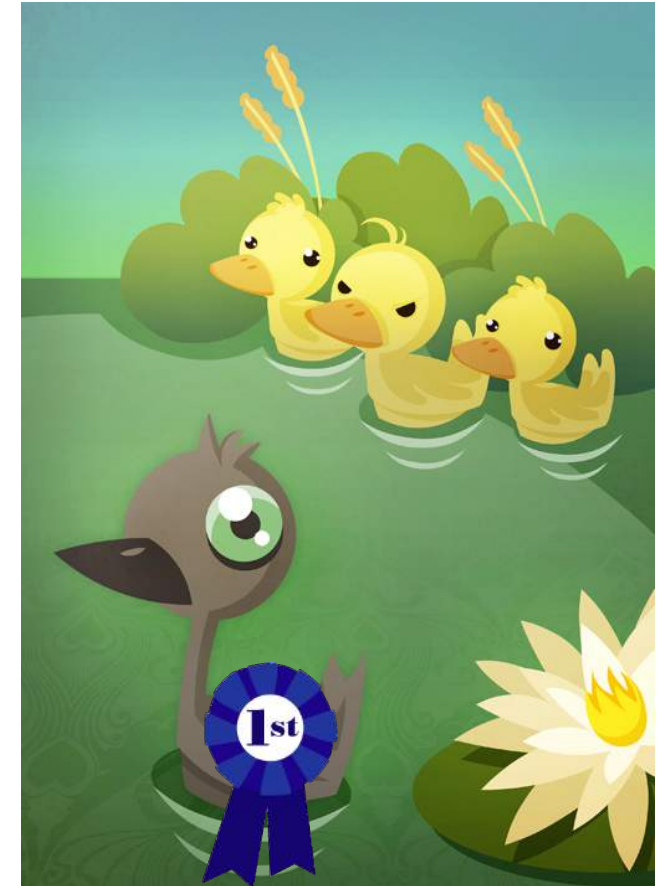
- Scores are algorithm-specific -> not the same scale !
- ... So how do we compare predictions across multiple algorithms ?

It's actually a ranking problem !

An interpretable prediction system

One possible solution: **run multiple algorithms and rank data points by anomaly score :**

1. Choose a reference algorithm
2. Compute dense ranks for each algorithm score column
3. Sort by increasing rank **for the reference algorithm**
4. Compare top N points and **see if other algorithms « agree » on ranking values**



It's actually a ranking problem !

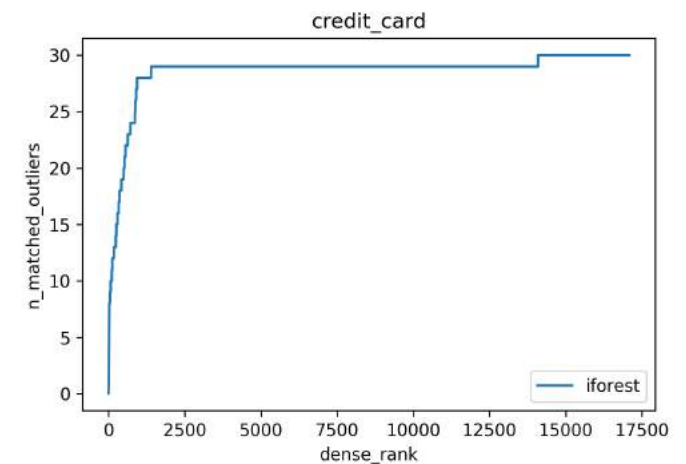
An interpretable prediction system

- Example on isolation forest's top 10 on the credit_card dataset:

	y_true	iforest_score	ocsvm_score	lof_score
3948	0	0.192481	1763.174401	3.699082
7723	1	0.189296	1763.174400	3.035509
1903	1	0.187674	1763.174401	4.620851
14861	1	0.172651	1763.174399	2.526187
1269	1	0.167669	1763.162092	2.201619
16129	0	0.161433	1763.174362	3.475972
9752	0	0.161160	1763.174401	3.970376
4548	1	0.159525	1763.174401	2.765399
2947	0	0.154045	1762.652671	1.782176
6966	0	0.152313	1763.174401	2.683288

dense_rank + sort

	y_true	iforest_rank	ocsvm_rank	lof_rank
3948	0	1.0	5.0	5.0
7723	1	2.0	7.0	13.0
1903	1	3.0	1.0	1.0
14861	1	4.0	9.0	39.0
1269	1	5.0	16.0	85.0
16129	0	6.0	10.0	9.0
9752	0	7.0	1.0	2.0
4548	1	8.0	1.0	27.0
2947	0	9.0	61.0	285.0
6966	0	10.0	4.0	32.0



~ ROC curve !

Prediction errors

Dealing with false positives/negatives

Types of errors and their business impact:

- **False positives** (high rank but NOT actual outliers)
 - Bad customer experience
 - Costly investigation (time + money)
- **False negatives** (low rank but actual outliers)
 - Potential (high) financial damage

	y_true	iforest_rank	ocsvm_rank	lof_rank	
3948	0	1.0	5.0	5.0	☹️ FP
7723	1	2.0	7.0	13.0	
1903	1	3.0	1.0	1.0	😊
14861	1	4.0	9.0	39.0	
1269	1	5.0	16.0	85.0	😐
16129	0	6.0	10.0	9.0	
9752	0	7.0	1.0	2.0	☹️ FP
4548	1	8.0	1.0	27.0	
2947	0	9.0	61.0	285.0	
6966	0	10.0	4.0	32.0	

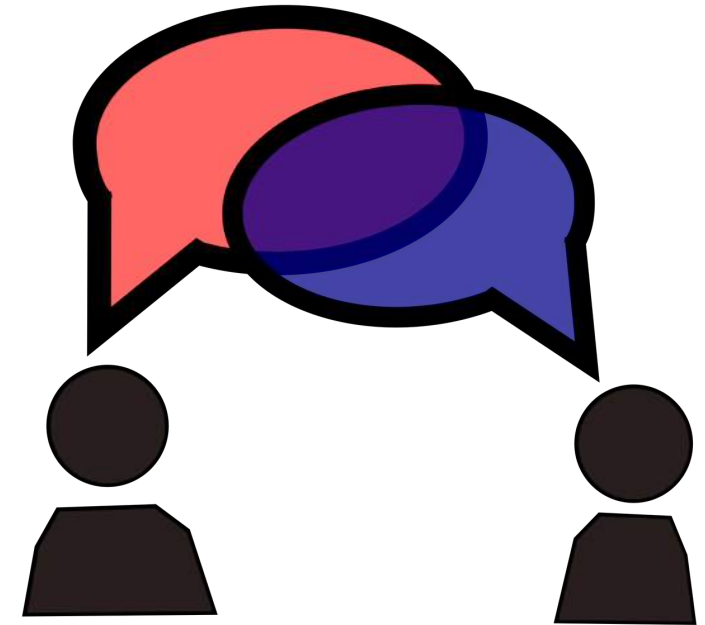
	y_true	iforest_rank	ocsvm_rank	lof_rank	
839	1	12263.0	11356.0	14544.0	☹️ FN
12708	1	1338.0	744.0	6213.0	
6103	1	573.0	559.0	102.0	
3187	1	571.0	434.0	424.0	
14770	1	523.0	615.0	101.0	

Prediction errors

Dealing with false positives/negatives

Solutions:

- Average rankings -> more robust estimation
- Build **custom filters/rules** (e.g. amount threshold, projected damage...)
- Rely heavily on domain knowledge experts !



Communication between data scientists and domain experts is paramount !

Summary & conclusion

Summary & conclusion

- Outlier detection is a **first step** to take when looking for suspicious points in an **unlabelled dataset**.
 - Provides a « **short list** » that can be further investigated
 - Several algorithms and implementations available (Python, R)
- Evaluating outlier detection algorithms is **difficult**:
 - Often requires **actual label availability**
 - Choice of hyperparameters
 - Area of **active research** (check out AD workshop @ICML)
- Reporting anomalies should be done in a **transparent way**:
 - Actions taken on outliers are **business-driven**
 - Algorithms are meant to **help domain experts**, not replace them

Summary & conclusion

Unsupervised anomaly detection is evolving fast:

- Other types of algorithms available:
 - **Distance-based**: variations of the LOF (e.g. COF, ODIN...)
 - **Statistical methods** (e.g. elliptic envelope fitting)
 - **Clustering-based**: flag points that don't fit with inferred clusters (e.g. DBSCAN)
- Novel approaches using **deep learning** for anomaly detection:
 - Images: *Learning discriminative reconstructions for unsupervised outlier removal*, Xia et al., 2015
 - Time series: *Engineering extreme event forecasting at Uber with recurrent neural networks*, Uber, 2017

Thank you for your attention !





Speaker Bio: Silviu Tofan

Data Scientist @Dataiku

Silviu is working as a data scientist at Dataiku. Coming from a business-oriented background, with an MSc in Business Analytics from the University of Manchester, he transitioned to a more technical focus while working on an optimization problem together with ARM. Before coming to Dataiku, Silviu was involved with social housing organizations to develop their data science capabilities.



data
iku

Anomaly Detection in Aviation

Silviu Tofan - Data Scientist

London Meetup - 24 January 2018

Outline

Problem Introduction

What exactly is our
airplane problem?

Deployment

How do we deploy this
into a production
environment?

Next Steps

Where do we go from
here?

Methodology

What are we trying to
solve? How do we
actually do it?

Monitoring

How do we showcase
the work and track
performance?

Problem Introduction

or what is this guy even talking about



Flight Tracking Data

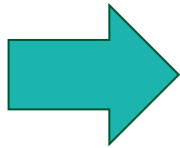
ETL explained for Aviation Enthusiasts



HPE Edgeline Server



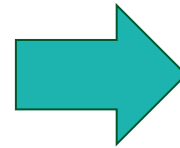
Real-time Data Capture



DUMP1090

A screenshot of a DUMP1090 data dump, showing a table of flight tracking data. The table has columns for various fields including hex, mode, type, flight, alt, speed, heading, and more. The data is presented in a grid-like format with multiple rows of information.

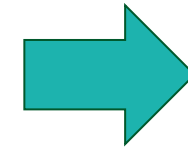
Decoding Signals &
Single Stream



Apache Kafka



Publish Messages
Over Kafka Topics



Microfocus Vertica



Micro-batch Into Tables

Collected Data



Viewing dataset sample

[Configure sample](#)

10000 rows, 16 cols

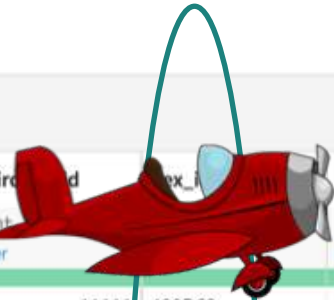


DISPLAY



10000 matching

dump1090_sitename	record_type	record_type_number	sg_session_id	sg_aircraft_id	sg_flight_id	msg_log_ts	altitude	latitude	longitude
string Natural lang.	string Text	bigint Integer	bigint Integer	bigint Integer	bigint Integer	string Date (unparsed)	bigint Integer		
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-20 20:20:08.088	2017-07-20 20:20:04.053	25100	47.24579	7.44398
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-20 20:20:09.138	2017-07-20 20:20:09.108	25575	47.23879	7.43355
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-20 20:20:10.038	2017-07-20 20:20:10.026	25550	47.23727	7.43809
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-20 21:37:15.083	2017-07-20 21:37:15.072	25400	46.70151	7.15498
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-20 21:37:17.682	2017-07-20 21:37:17.635	25475	46.70474	7.16109
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 00:09:49.756	2017-07-21 00:09:49.749	38025	47.43547	7.29689
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-21 08:16:47.317	2017-07-21 08:16:47.299	25900	47.24162	7.46459
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-21 08:16:48.328	2017-07-21 08:16:48.287	25850	47.24047	7.46324
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-21 08:16:51.088	2017-07-21 08:16:51.039	25750	47.2367	7.45859
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-21 08:16:52.167	2017-07-21 08:16:52.149	25700	47.23526	7.45683
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-21 08:16:52.648	2017-07-21 08:16:52.61	25675	47.23441	7.45577
pennardel01	MSG	3	111	11111 400D69	111111 2017-07-21 08:16:55.589	2017-07-21 08:16:55.558	25575	47.23064	7.45112
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:50:54.112	2017-07-21 08:50:54.095	38025	46.33358	5.46721
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:50:55.186	2017-07-21 08:50:55.148	38025	46.33177	5.46975
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:50:55.583	2017-07-21 08:50:55.541	38025	46.33116	5.47064
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:50:56.083	2017-07-21 08:50:56.064	38025	46.33056	5.47147
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:53:14.909	2017-07-21 08:53:14.879	38025	46.12172	5.76645
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:53:15.329	2017-07-21 08:53:15.328	38025	46.12112	5.76727
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:53:15.878	2017-07-21 08:53:15.855	38025	46.12019	5.76856
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:53:16.289	2017-07-21 08:53:16.252	38025	46.1196	5.76943
pennardel01	MSG	3	111	11111 4054A7	111111 2017-07-21 08:53:17.789	2017-07-21 08:53:17.758	38050	46.11717	5.77285



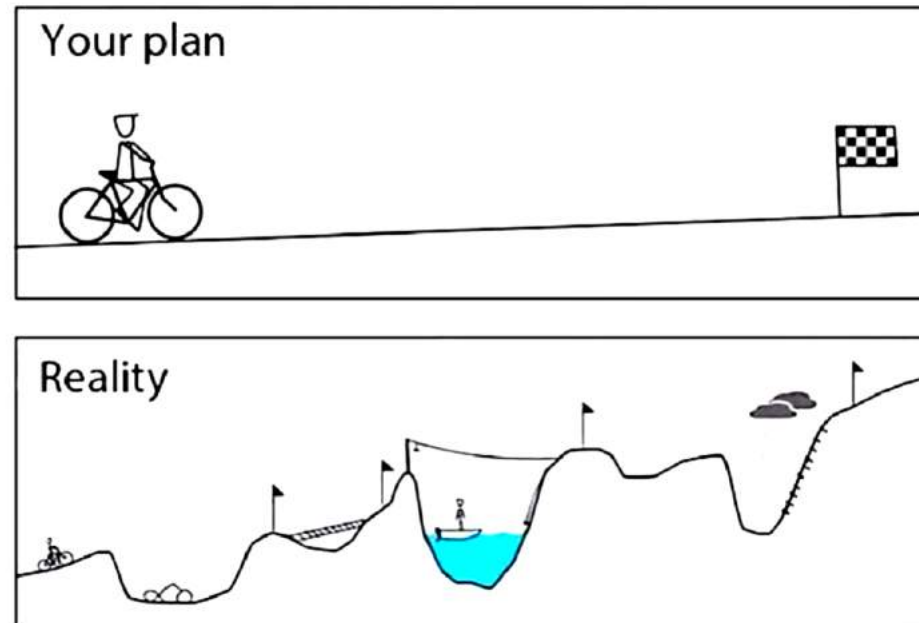
Methodology

or how we came up with all this stuff



Defining a Goal

“Can we say that the current **path** followed by a particular **flight** is **abnormal**?”



Separation of Flights

Some aviation knowledge

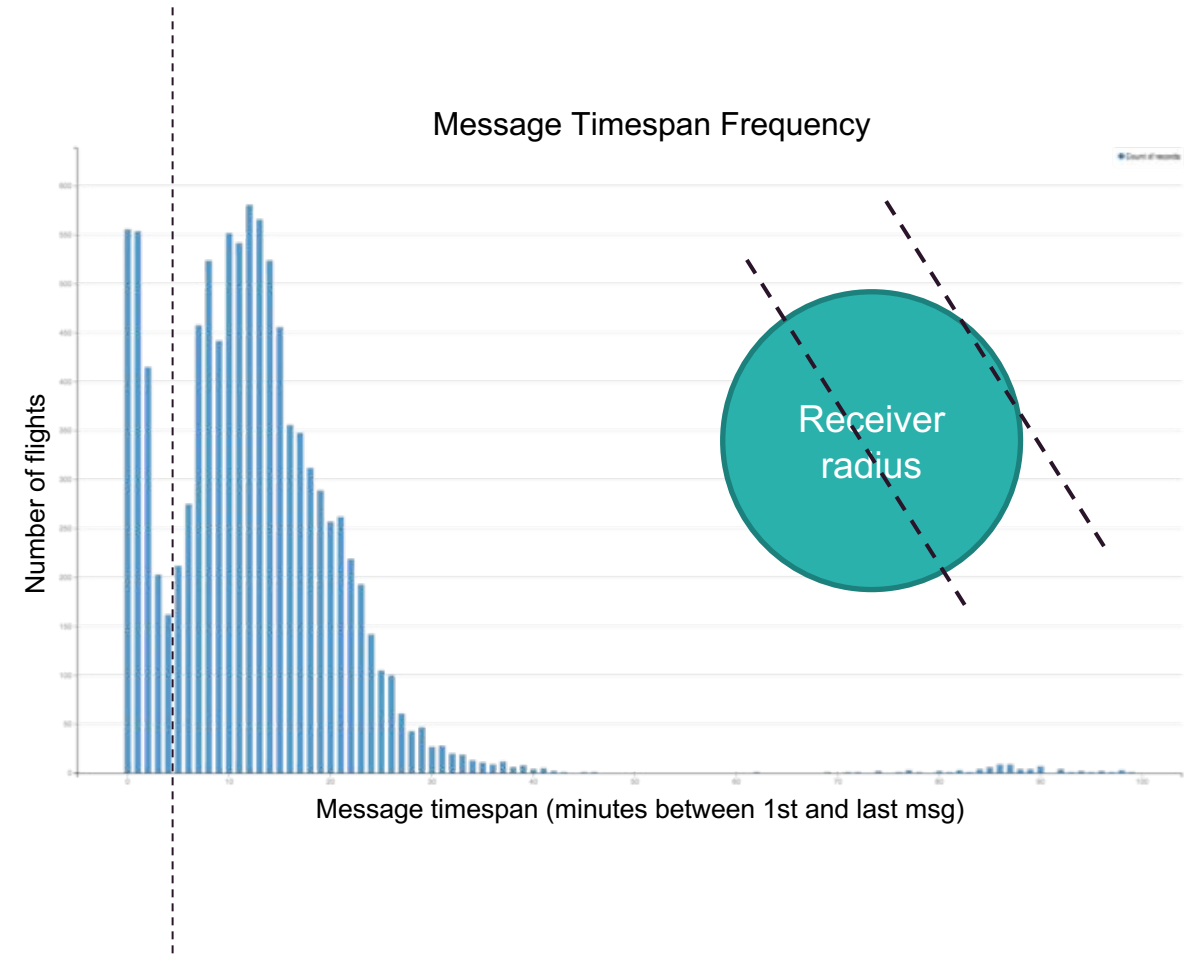
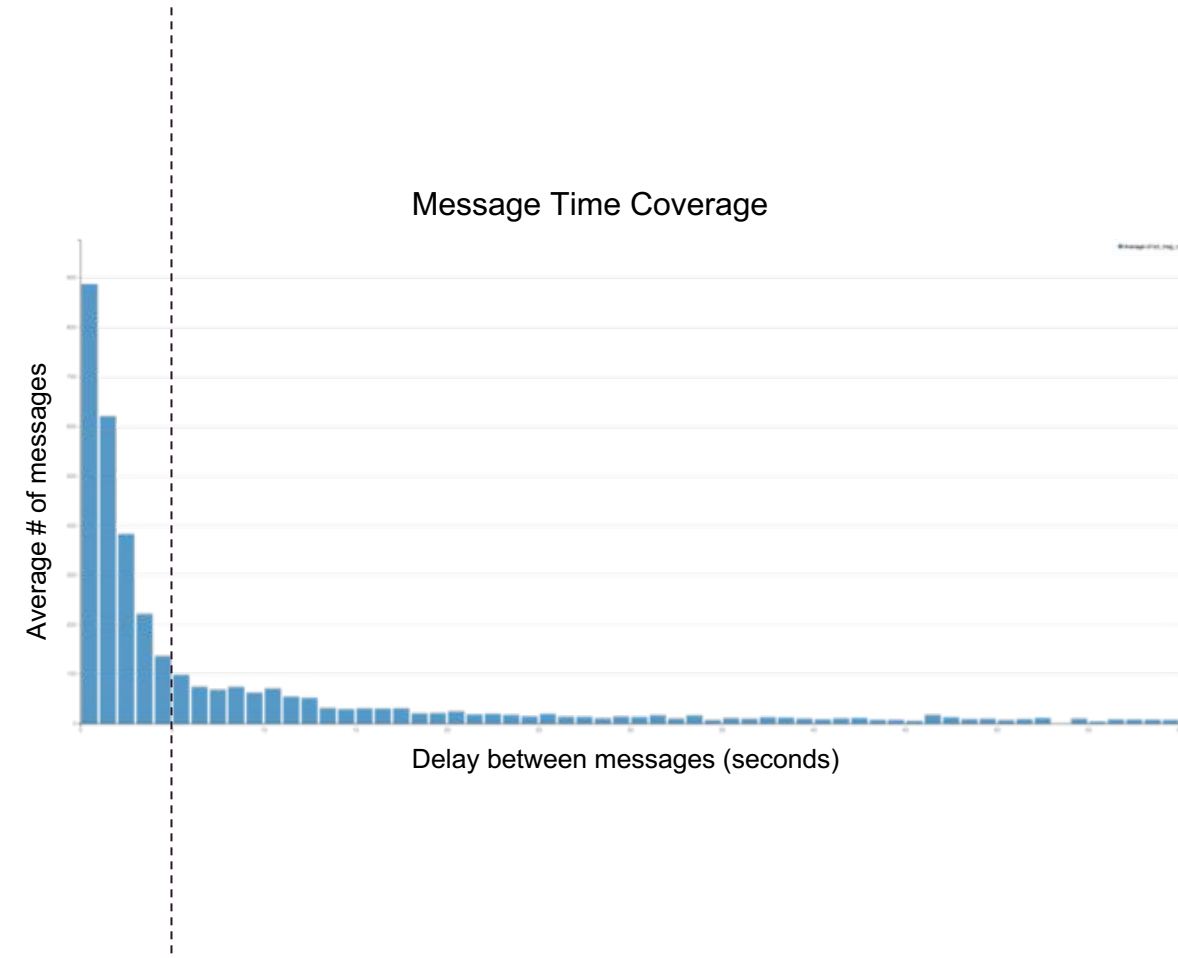
Some data checks

Some subjectivity

Some Vertica SQL

```
SELECT *,  
    hex_ident||'_'||conditional_true_event(msg_gen_ts - LAG(msg_gen_ts) > '1 HOUR')  
    OVER (PARTITION BY hex_ident ORDER BY msg_gen_ts) AS flight_id  
FROM "dump1090_kafka"."dump1090_msg_3"
```


Feature Creation



Prepared Dataset



Viewing dataset sample

Configure sample

DISPLAY

10000 rows, 66 cols

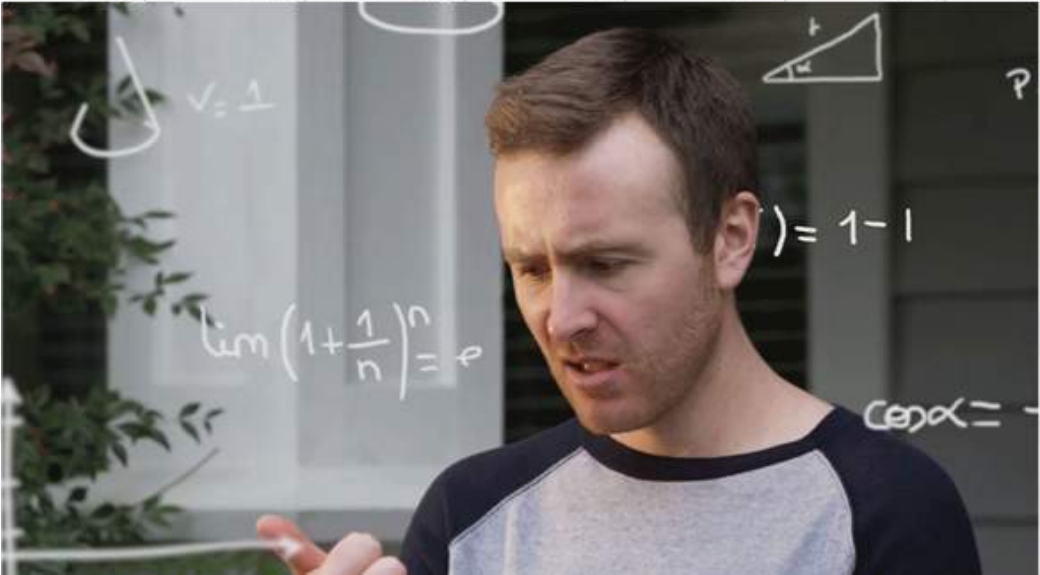
10000 matching rows

flight_id	ts_step	ts_step_partition	ts_rank	alt	lag1	lag2	lag3	lag4	lag5	lag6	lag7	lag8	lag9
string	date	date	bigint	string	string	string	string	string	string	string	string	string	string
Text	Date	Date	Integer	Integer	Integer	Integer	Integer	Integer	Integer	Integer	Integer	Integer	Integer
0A004B_0	2017-12-17T15:47:05.000Z	2017-12-17T15:47:05.000Z		36000	35975	35975	36000	36000	36000	36000	36000	36000	36000
0A004B_0	2017-12-17T15:47:10.000Z	2017-12-17T15:47:10.000Z		36000	36000	36000	35975	36000	36000	36000	36000	36000	36000
0A004B_0	2017-12-17T15:47:15.000Z	2017-12-17T15:47:15.000Z		36000	36000	36000	36000	35975	36000	36000	36000	36000	36000
0A004B_0	2017-12-17T15:47:20.000Z	2017-12-17T15:47:20.000Z		36000	36000	36000	36000	35975	35975	36000	36000	36000	36000
0A004B_0	2017-12-17T15:47:25.000Z	2017-12-17T15:47:25.000Z		36000	36000	36000	36000	36000	35975	36000	36000	36000	36000
0A004B_0	2017-12-17T15:47:30.000Z	2017-12-17T15:47:30.000Z		36000	36000	36000	36000	36000	36000	36000	36000	36000	36000
0A004B_0	2017-12-17T15:47:35.000Z	2017-12-17T15:47:35.000Z	145	36000	36000	36000	36000	36000	36000	36000	35975	35975	36000
0A004B_0	2017-12-17T15:47:40.000Z	2017-12-17T15:47:40.000Z	146	36000									
06A0B2_1	2017-12-17T18:20:00.000Z	2017-12-17T18:20:00.000Z	147	36000									
06A0B2_1	2017-12-17T18:20:05.000Z	2017-12-17T18:20:05.000Z	61	12850									
06A0B2_1	2017-12-17T18:20:10.000Z	2017-12-17T18:20:10.000Z	62	12925									
06A0B2_1	2017-12-17T18:20:15.000Z	2017-12-17T18:20:15.000Z	63	13150									
06A0B2_1	2017-12-17T18:20:20.000Z	2017-12-17T18:20:20.000Z	64	13325									
06A0B2_1	2017-12-17T18:20:25.000Z	2017-12-17T18:20:25.000Z	65	13500									
06A0B2_1	2017-12-17T18:20:30.000Z	2017-12-17T18:20:30.000Z	66	13725									
06A0B2_1	2017-12-17T18:20:35.000Z	2017-12-17T18:20:35.000Z	67	13925									
06A0B2_1	2017-12-17T18:20:40.000Z	2017-12-17T18:20:40.000Z	69	14300									
06A0B2_1	2017-12-17T18:20:45.000Z	2017-12-17T18:20:45.000Z	70	14375									
06A0B2_1	2017-12-17T18:20:50.000Z	2017-12-17T18:20:50.000Z	71	14600									
06A0B2_1	2017-12-17T18:20:55.000Z	2017-12-17T18:20:55.000Z	72	14800									

Current altitude

Altitude 5 seconds ago

Altitude 30 seconds ago



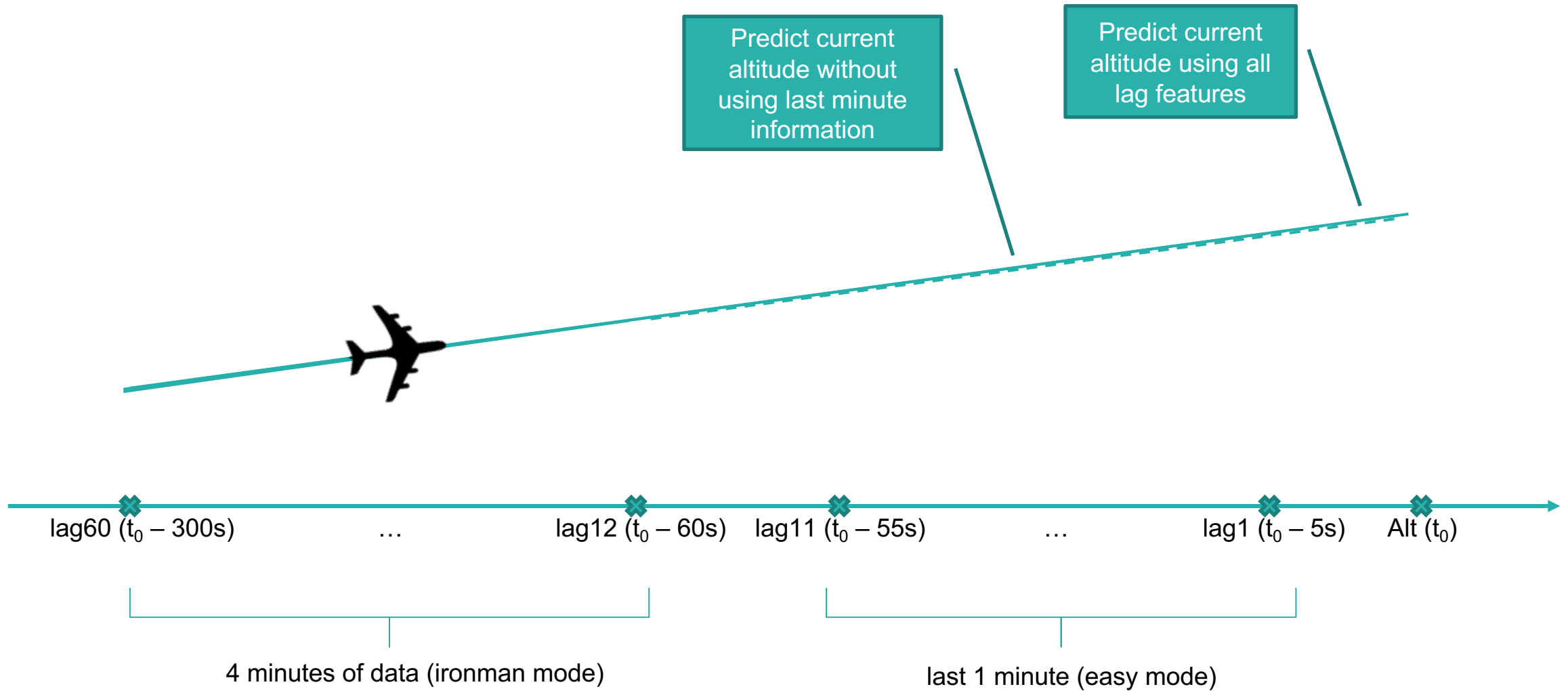
We Have to Adapt



“The difficulty lies, not in the new ideas, but in escaping from the old ones, which ramify, for those brought up as most of us have been, into every corner of our minds.”

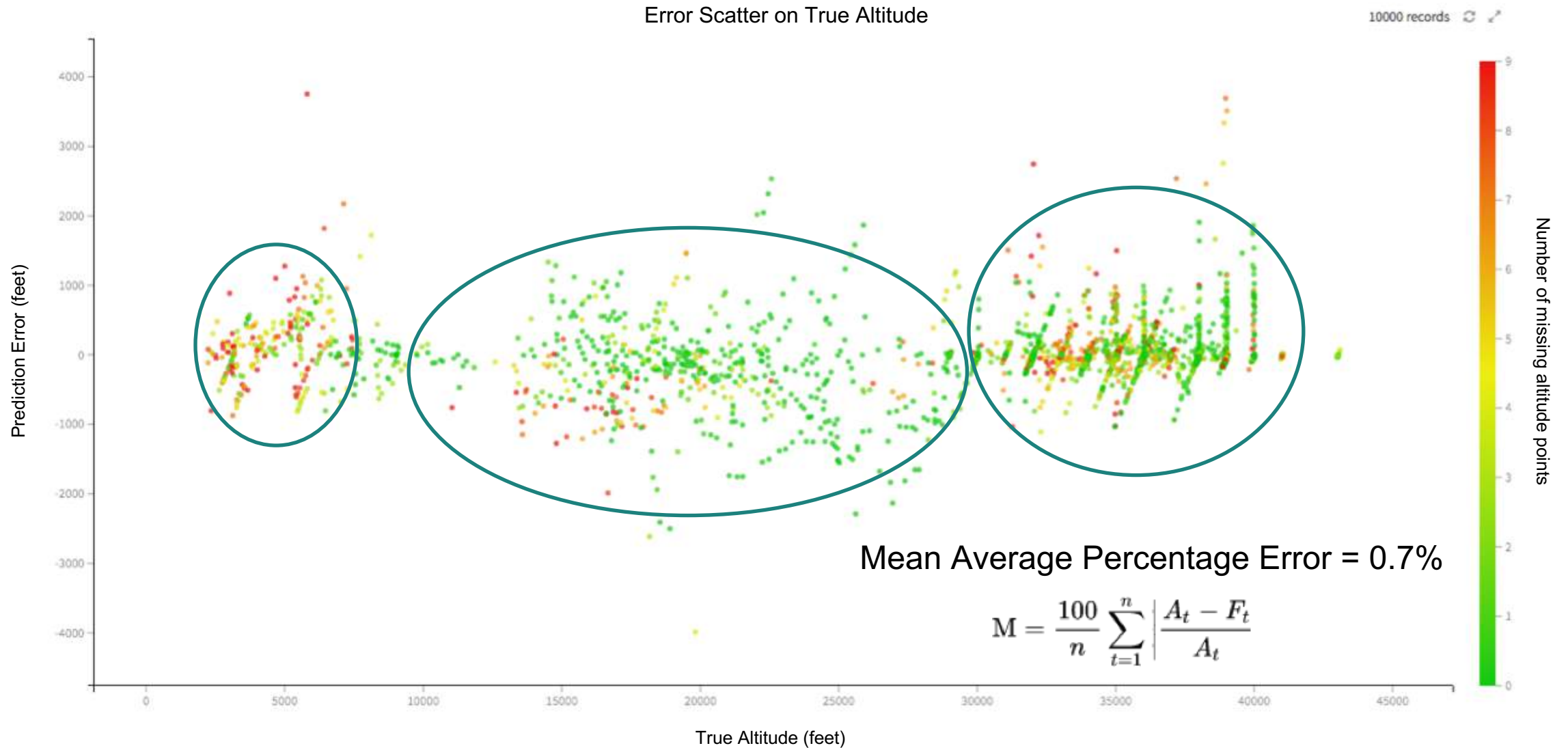
John Maynard Keynes, Economist

Altitude Timeline

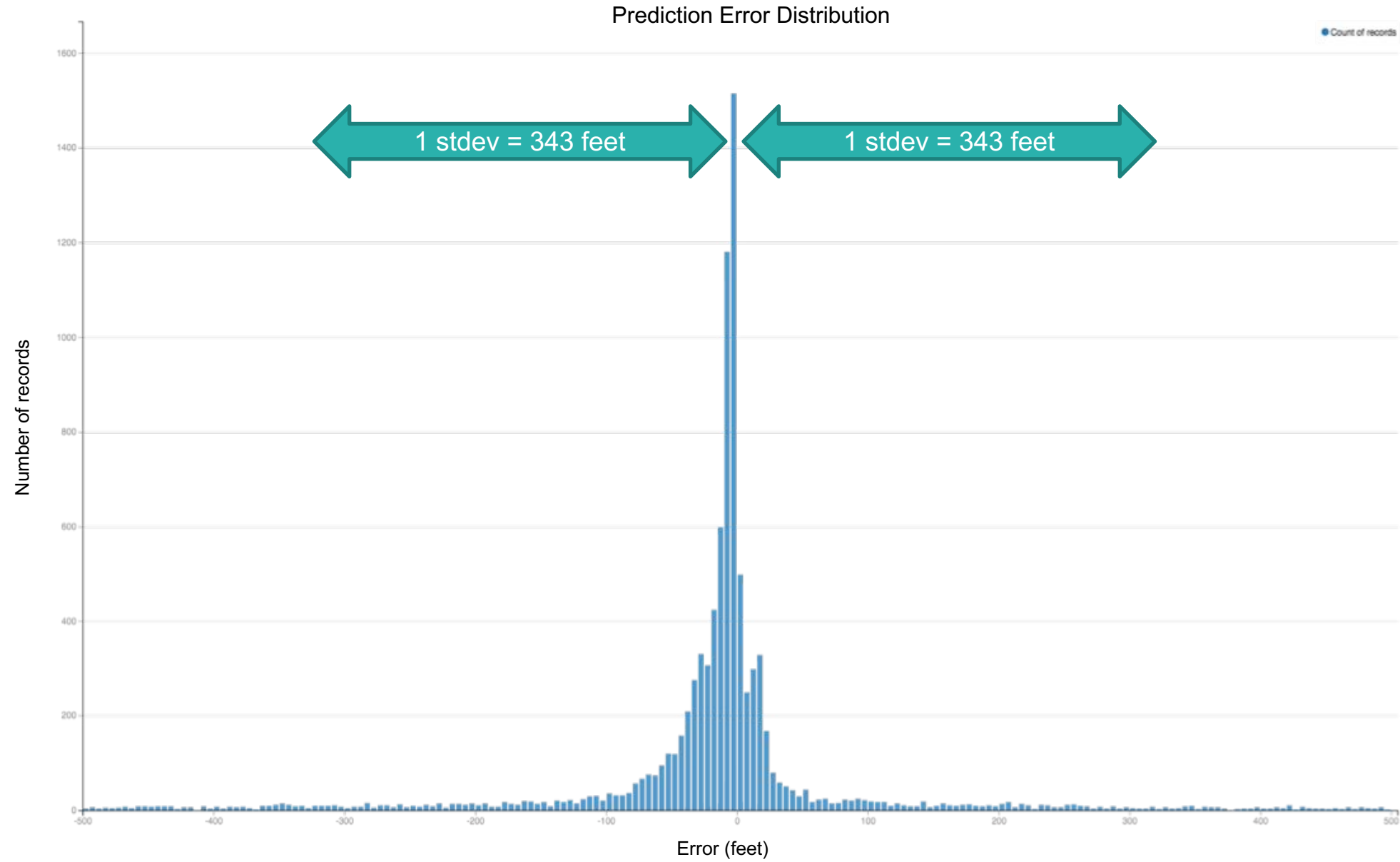


Altitude Predictions

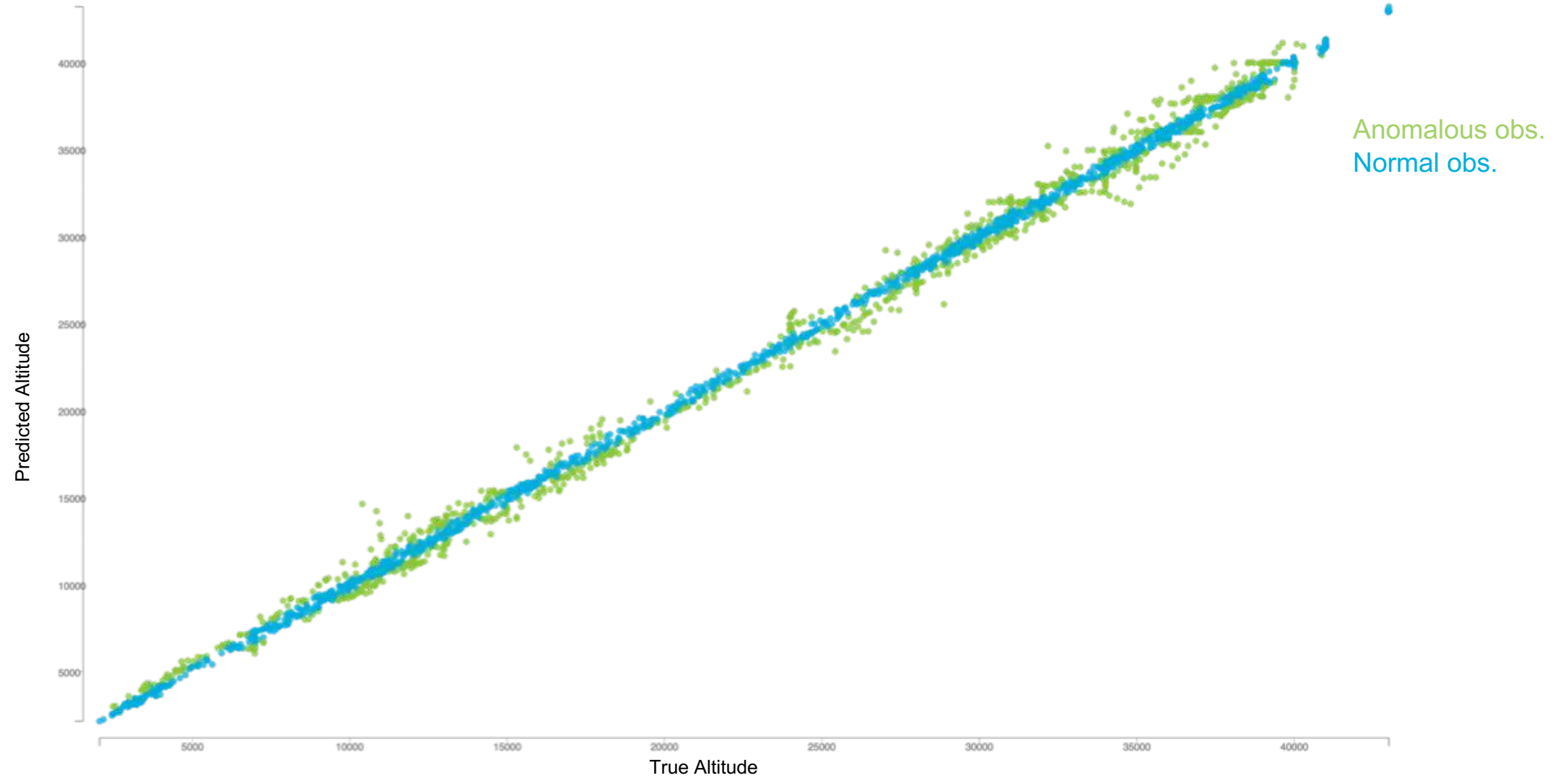
Error Scatter on True Altitude



Prediction Errors



Anomalies



Implemented Improvements

Missing Values

Removed rows with 10 or more missing historical data in the last 5 minutes

Train/Control Split

- Totaling 5 weeks of data
- First 4 weeks train
- Final week control

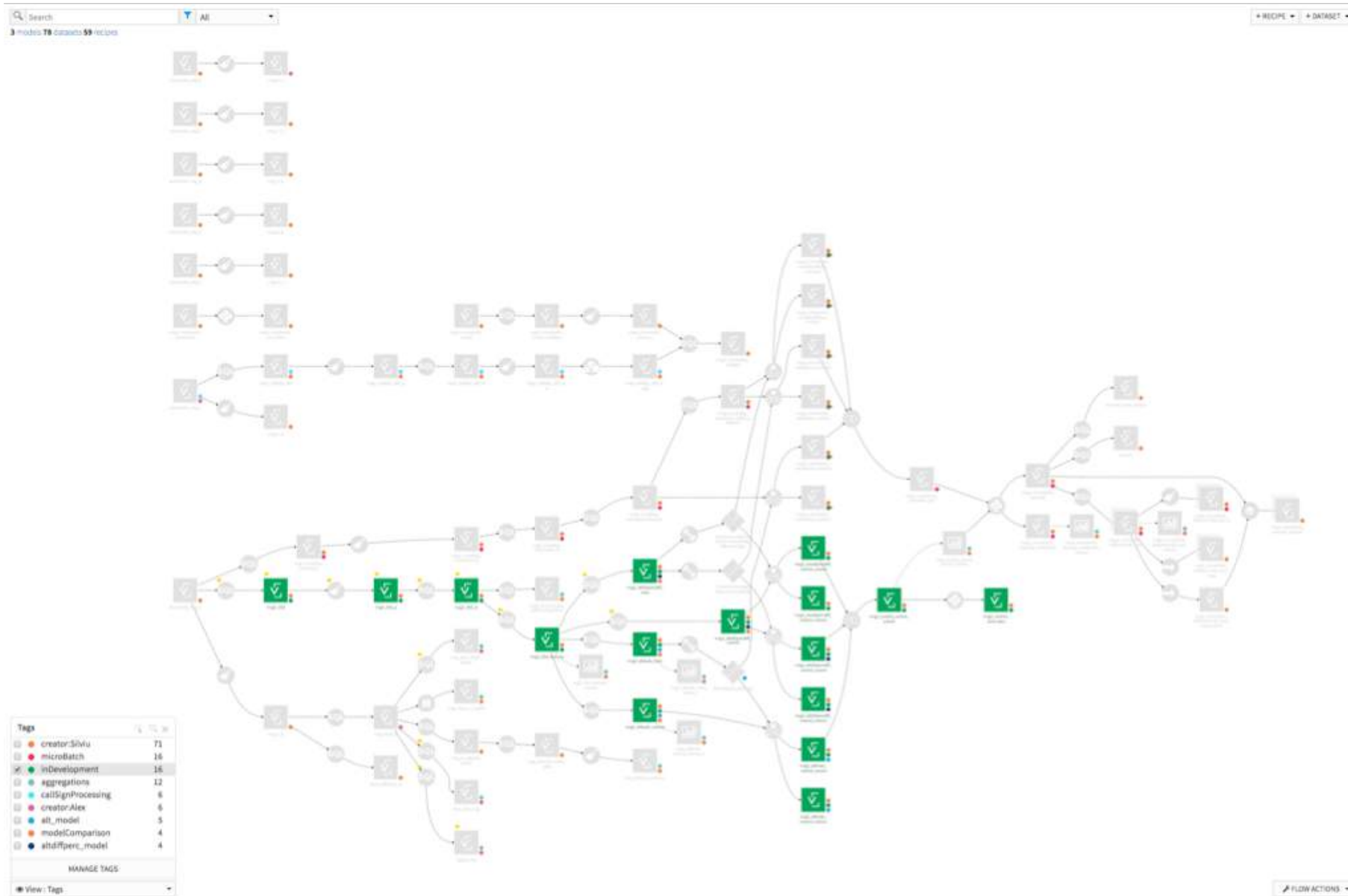
Multiple Feature/Prediction Models

- Raw altitude
- Percentage difference in altitude
- Combination of the two

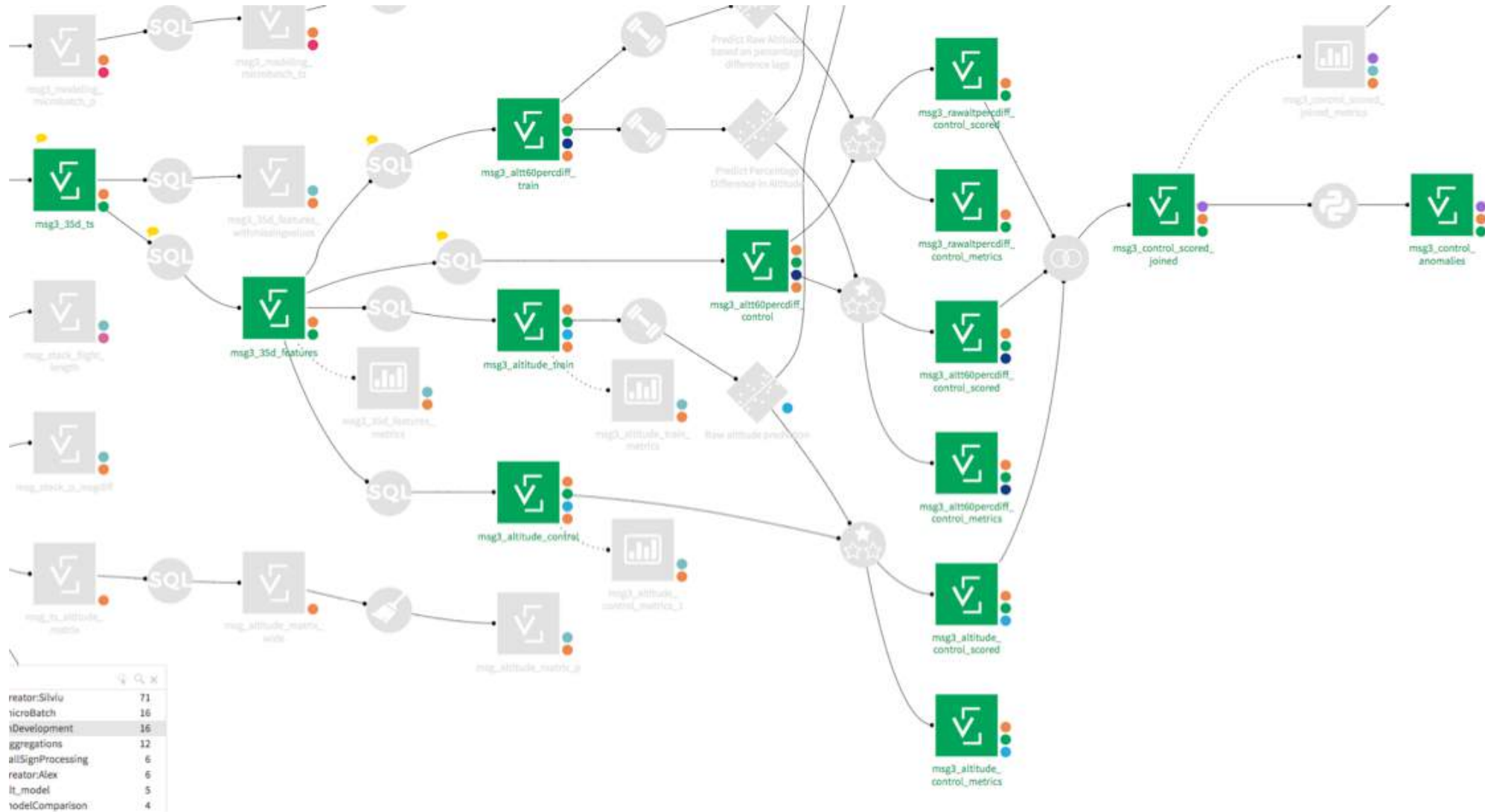
Overall Flow Cleanliness

Continuously reformatted the flow for maintainability

Final Modeling Flow



Final Modeling Flow

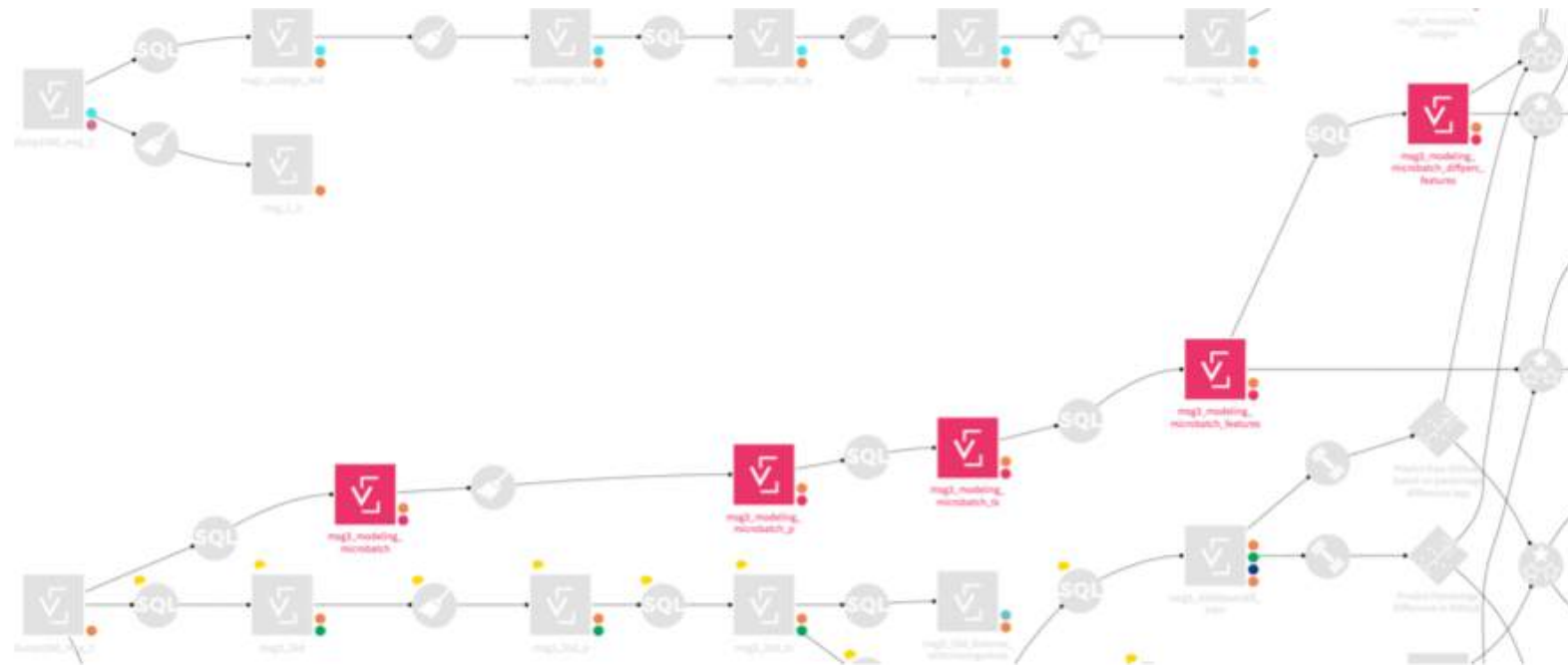


Deployment

or how this is using our partner's IT resources



Microbatches



Triggers

ADD TRIGGER

Time-based Active ☒

Frequency

Run every minutes Interval in minutes between each run of the scenario

Last triggered on 2018/01/23/09:13:25

AirPlane

compute_msg3_modeling_microbatch

Datasets Variables

Inputs

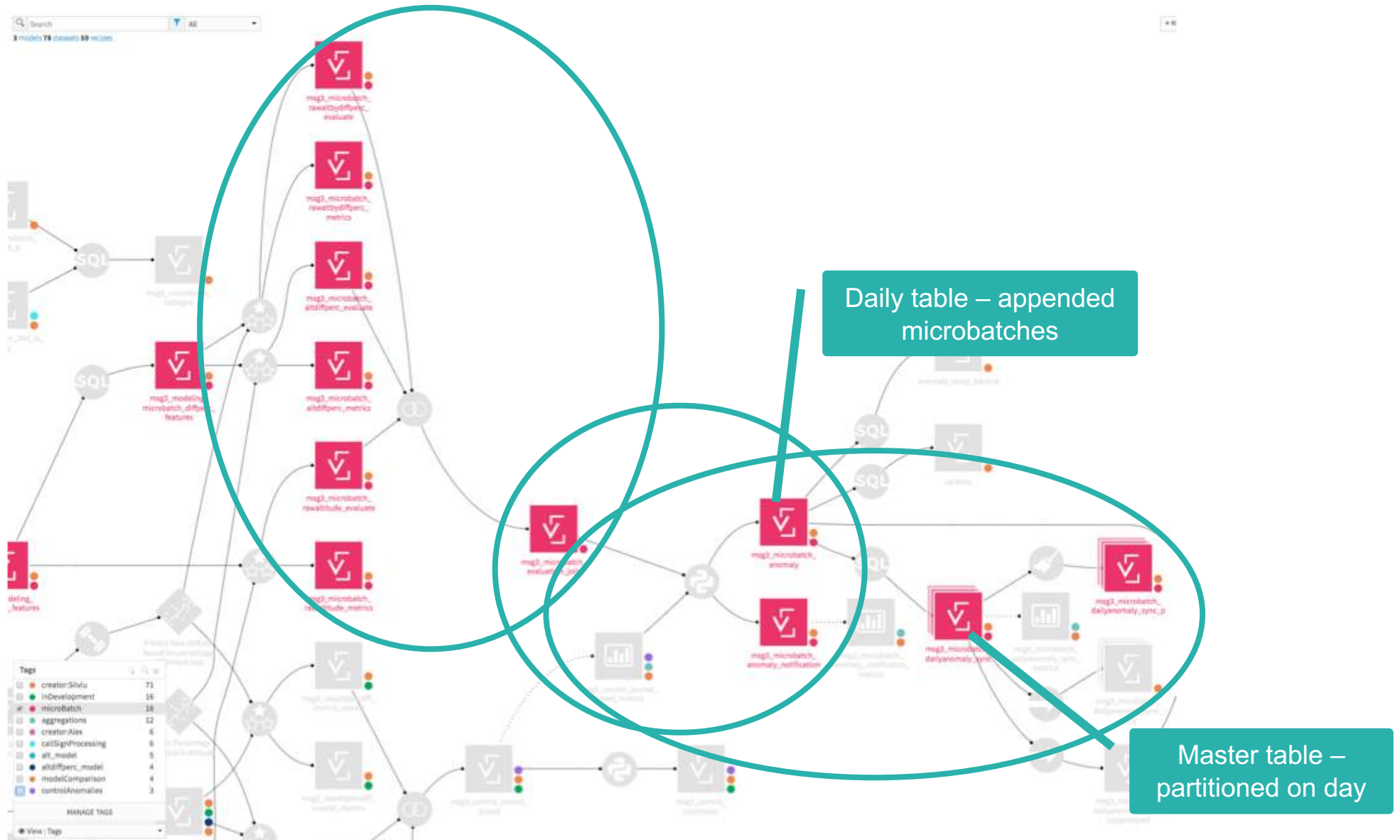
- dump1090_msg_3
Table: dump1090_msg_3

Outputs

- msg3_modeling_microbatch
Table: msg3_modeling_microbat...

```
1 select *,hex_id||'|'||conditional_true_event( msg_gen_ts - LAG(msg_gen_ts) > '1 HOUR') over (partition by hex_id order by msg_gen_ts) as flight_id
2 from "dump1090_kafka"."dump1090_msg_3"
3 where CAST(msg_gen_ts AS DATETIME) >= TIMESTAMPADD (minute, -6, GETDATE())
```


Microbatches



Scenario Management



S1-S4) Get Data, Check Control, Retrain, Update

(id: S1S4GETDATAHECKCONTROLRETRAINUPDATE)

S5) Build and Eval Microbatch

(id: S5BUILDANDEVALMICROBATCH)

S6) Push and Clear Microbatch

(id: S6PUSHANDCLEARMICROBATCH)

Auto-triggers ☒ ON
1 am run

Auto-triggers ☒ ON
Running since 2018/01/22 10:13:10

Auto-triggers ☒ ON
Midnight batch push

S5) Build and Eval Microbatch

Build microbatch features

Build evaluate microbatches

Build joined microbatch evals

Build append microbatch anomalies

Build

microbatch features

INFO

Build mode

Force-rebuild dataset and dependencies

Sync. Hive metastore



Ignore failure



Continue the scenario even if the build fails (warning: the datasets' states may be incoherent)

Item

mag3_modeling_microbatch_diffperc_features

+ DATASET

+ FOLDER

+ MODEL

S1-S4) Get Data, Check Control, Retrain, Upd...

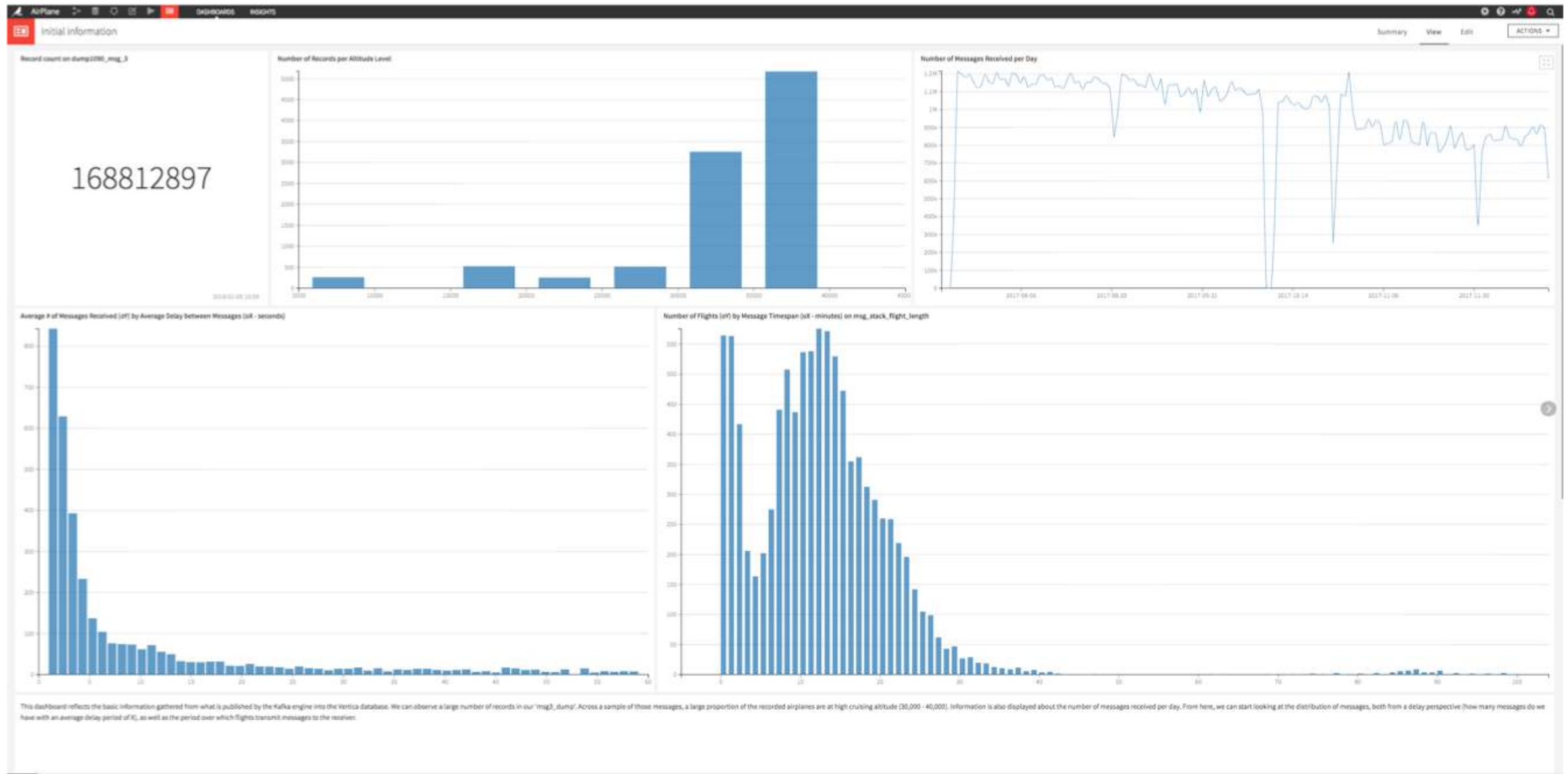
```
1 from dataiku.scenario import Scenario
2 import dataiku
3 import pandas as pd
4 scenario = Scenario()
5
6 #Function that returns the latest 'metric' from a 'table'
7 def get_metric(table, metric):
8     df = dataiku.Dataset(table).get_dataframe()
9     df.sort_values(by='date', ascending=False, inplace=True)
10    val = df[metric].iloc[0]
11    return val
12
13
14 #Force build "mag3_35d_features" and all dependencies
15 scenario.build_dataset("mag3_35d_features", build_mode="RECURSIVE_FORCED_BUILD")
16
17 #Build equivalent but with missing values (for reporting purposes)
18 scenario.build_dataset("mag3_35d_features_withmissingvalues", build_mode="NON_RECURSIVE_FORCED_BUILD")
19
20 #Get current MAPE / RMSE for the previous control samples
21 rawalt_mape = get_metric("mag3_altitude_control_metrics", 'mape')
22 rawaltbydiffperc_mape = get_metric("mag3_rawaltperdiff_control_metrics", 'mape')
23 altdiffperc_rmse = get_metric("mag3_altd60perdiff_control_metrics", 'rmse')
24
25 #Build control samples - only "mag3_altitude_control" and "mag3_altd60perdiff_control"
26 scenario.build_dataset("mag3_altitude_control", build_mode="NON_RECURSIVE_FORCED_BUILD")
27 scenario.build_dataset("mag3_altd60perdiff_control", build_mode="NON_RECURSIVE_FORCED_BUILD")
28
29 #Evaluate new control samples
30 scenario.build_dataset("mag3_altitude_control_scored", build_mode="NON_RECURSIVE_FORCED_BUILD")
31 scenario.build_dataset("mag3_rawaltperdiff_control_scored", build_mode="NON_RECURSIVE_FORCED_BUILD")
32 scenario.build_dataset("mag3_altd60perdiff_control_scored", build_mode="NON_RECURSIVE_FORCED_BUILD")
33
34 #Get MAPE/RMSE for the newly created controls
35 rawalt_mape_new = get_metric("mag3_altitude_control_metrics", 'mape')
36 rawaltbydiffperc_mape_new = get_metric("mag3_rawaltperdiff_control_metrics", 'mape')
37 altdiffperc_rmse_new = get_metric("mag3_altd60perdiff_control_metrics", 'rmse')
38
39 #Check control performance against current MAPE/RMSE (within 10%) as well as fixed values
40 #Do we retrain when the mape decreases too?
41
42 #Normal zeros now
43
44 if rawalt_mape_new >= 0.015 or rawalt_mape_new >= 1.1 * rawalt_mape:
45
46     #Get metric in current active version
47     for model in dataiku.Model('GaoJN868').list_versions():
48         if model['active'] == True:
49             active_mape = model['snippet']['mape']
50
51     train_ret = scenario.train_model("GaoJN868")
52     trained_model = train_ret.get_trained_model()
53     performance = trained_model.get_new_version_metrics().get_performance_values()
54     #performance = get_version_metrics().get_performance_values()
55     if performance["MAPE"] < active_mape:
56         trained_model.activate_new_version()
57
58 #Rebuild scored/metrics
59 scenario.build_dataset("mag3_altitude_control_scored", build_mode="NON_RECURSIVE_FORCED_BUILD")
60
```

Monitoring

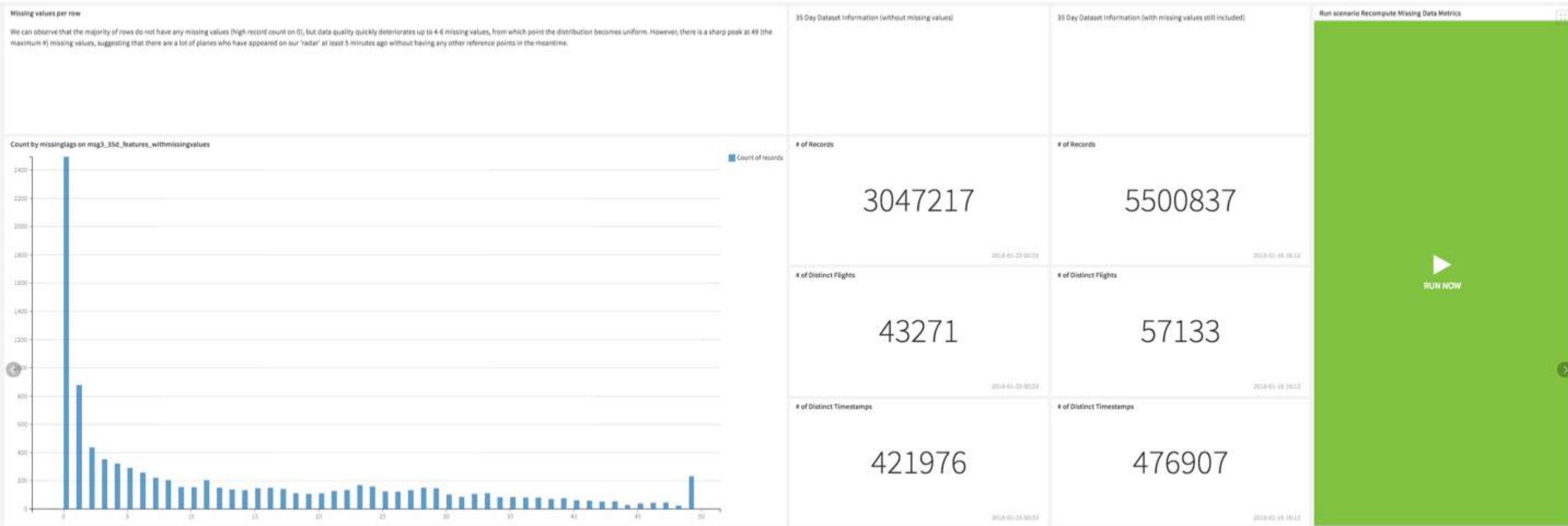
or how we remember we did all of this



Exploratory Info



Missing Data Analysis

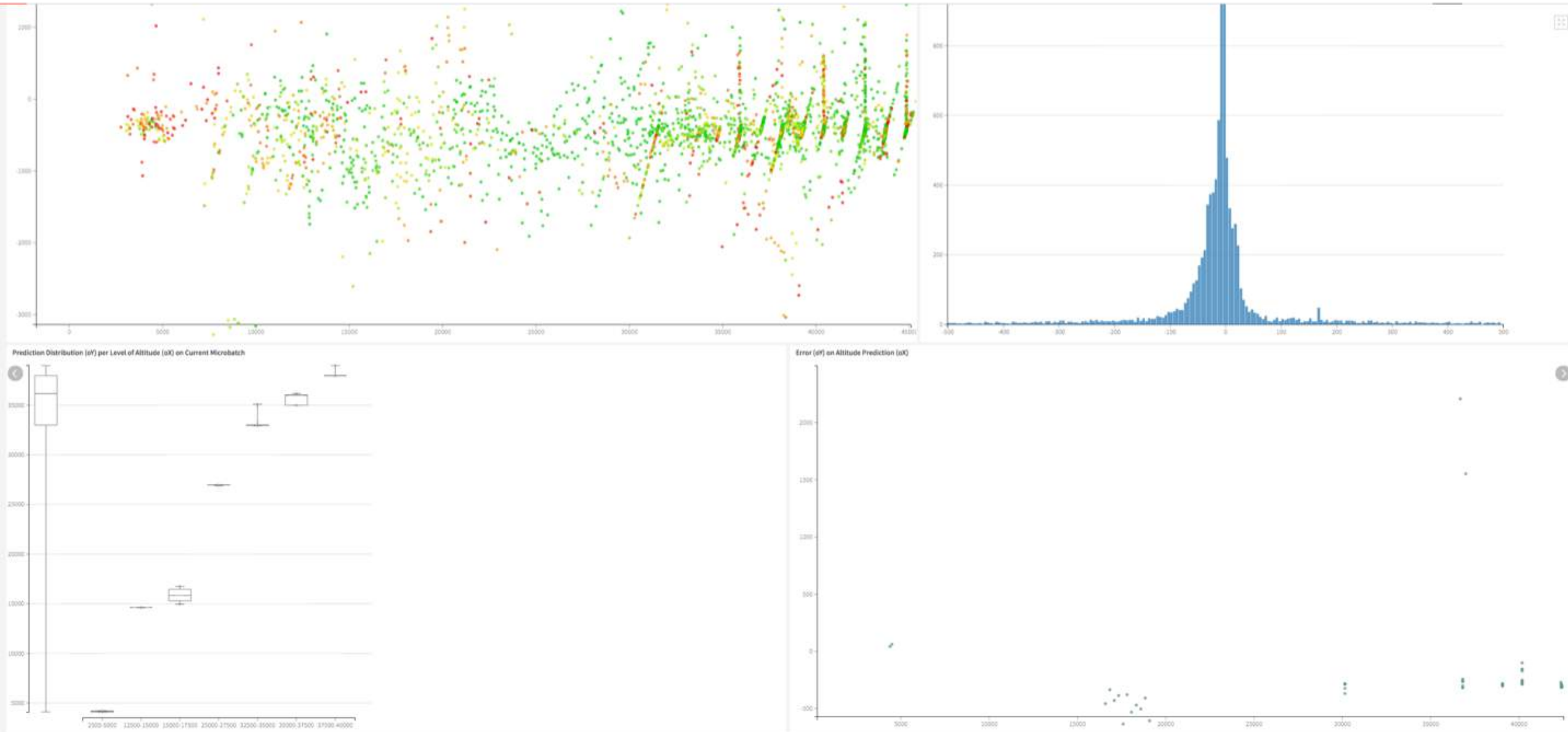


Splitting Methodology Review

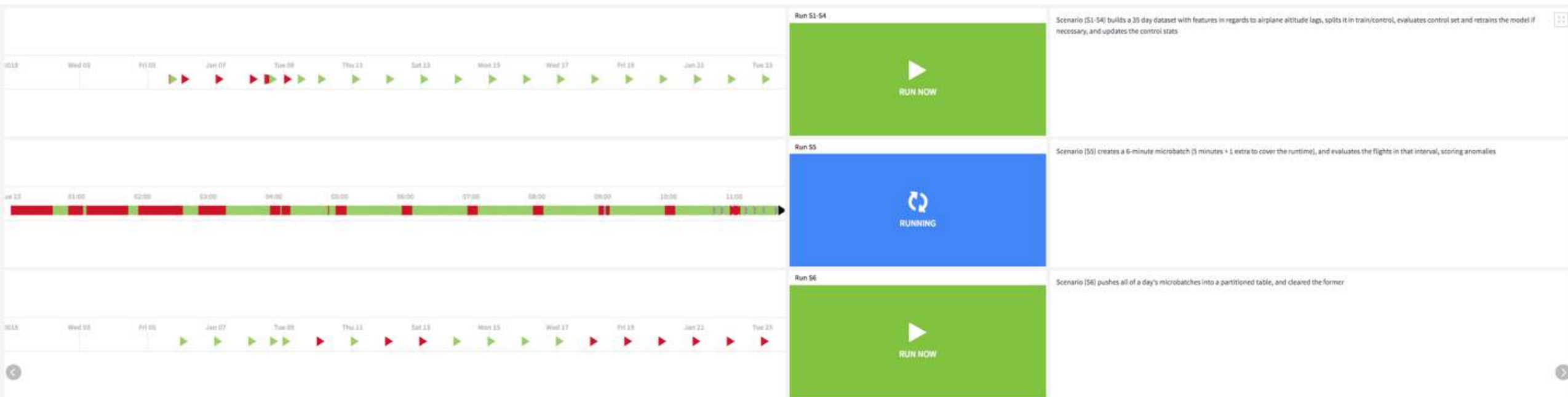


Train + Control	Train	Control	Current Microbatch	Aggregated Microbatch	Run scenario to recompute metrics
Information about stacked train and control sets (35 days in the past from yesterday - at point of running)	Information about the training set (a period of 4 weeks from 1 week ago - at point of running)	Information about the control set (a period of 1 week from now - at point of running)	Information about our current microbatch set	Information about the aggregate number of microbatches we have scored so far	
# of Records	# of Records	# of Records	# of Records	# of Partitions	
3047217	2512476	648500	93	5	
# of Distinct Flights	# of Distinct Flights	# of Distinct Flights	# of Distinct Flights	# of Distinct Flights	
43271	35602	8927	12	2257	
Min Timestamp	Min Timestamp	Min Timestamp	Min Timestamp	Min Timestamp	
2017-12-19 01:05:35	2017-12-04 01:14:10	2018-01-16 00:51:25	2018-01-23 11:24:25	2018-01-10 10:11:00	
Max Timestamp	Max Timestamp	Max Timestamp	Max Timestamp	Max Timestamp	
2018-01-23 00:51:25	2017-12-31 12:05:40	2018-01-23 00:51:25	2018-01-23 11:25:15	2018-01-16 23:54:55	

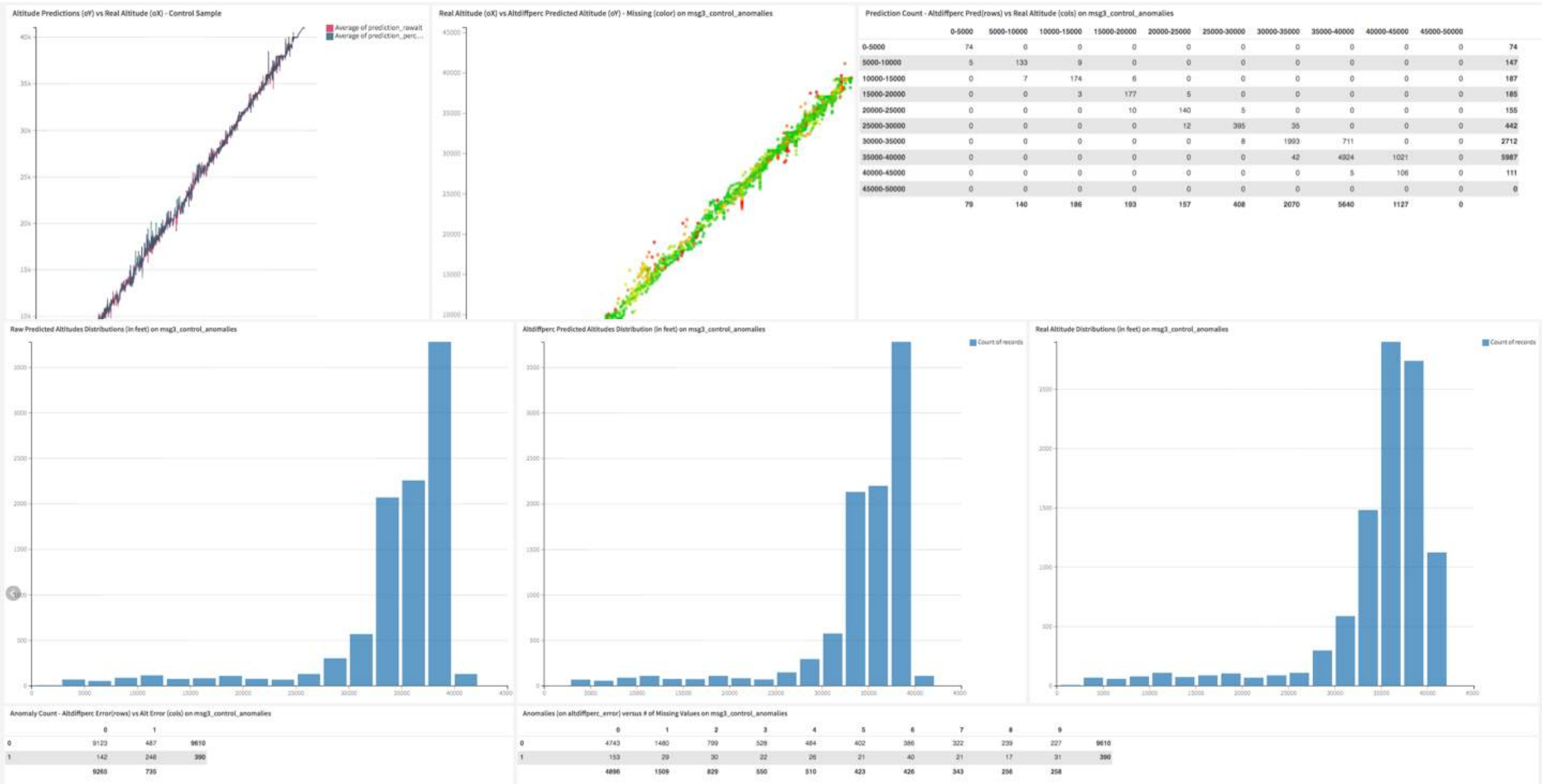
Model Performance



Scenario Tracking



Model Comparison



Next Steps

or if we have anything else to do

(I couldn't find a funny gif)

We can do LOADS more



Great story so far, but how about:

- latitude
- longitude
- speed
- verticality

Similar train of thought, different approach

- regression target becomes an array
- deep learning involved (5 features x 60 timestamps x N rows)
- add anomaly definition on top
- make it run in near real-time

Key Takeaways



- Don't always adhere to the standard
- Start simple
- Keep stakeholders involved
- Think (a lot) and don't rush

