

Enhancing NJCCIC Threat Detection: A Machine Learning Approach to Anomalous Login Behavior

Problem Statement and Goals

The New Jersey Cybersecurity and Communications Integration Cell leads and coordinates the cybersecurity efforts of New Jersey to protect all its citizens. Up to 80 percent of all breaches, including unauthorized access and human operated ransomware attacks, are the results of stolen or compromised login credentials¹. The NJCCIC seeks to improve their current threat detection methodologies. I assisted their efforts by experimenting with machine learning capabilities for anomalous login behavior.

Goals:

- Leverage data engineering techniques for feature preparation.
- Apply Machine Learning for classification of anomalous logins.



Rosemary James is a hard-working New Jersey state employee and mother of three children. Her family has access to many smart devices: televisions, phones, tablets, Google Home, and watches. Rosemary is always on the move trying to keep up with work and her kids. She often accesses her state account from different devices and locations daily. Sometimes she forgets to log out of her account when she is not working. Her kids love to watch educational videos on her iPhone. If Rosemary is not careful, her kids may click a phishing link and have her login credentials breached.

	Awareness	Consideration	Decision	Retention
Touchpoints	User logs onto network as usual.	Account password changed. Someone hacked the account and obtained her login credentials.	Following tips on account security given by the customer service representative of the website.	Going to NJCCIC for tips on best practices for personal security online.
Actions	Leaves her tablet/phone logged-in in reach of her children. Ignoring google warning of password security.	Consider asking customer service for help in reinstating the account. Implications of critical changes in the account/lost critical information. Consult NJCCIC for reporting stolen login.	Changing and adding passwords. Instructing children about internet safety.	Add MFA (Multi-factor authentication). Not using the same passwords for everything.
Experience	Happy	Disappointed	Interested	Happy
Solution	Continues as usual.	Contact customer service representatives.	Being vigilant.	Advocate for best cybersecurity practices for everyone.

Approach

General Approach:

The NJCCIC team provided a subset of raw log files. Feature engineering techniques, such as aggregation and binning, were applied to create new features aimed at catching behavior patterns of users. There were three sets of features including our initial set, our larger 4/4a set, and our minimized 4b set. A correlation matrix was used to remove initial features that correlated too strongly with each other to avoid multicollinearity. Feature selection and hyperparameter tuning was performed when training our machine learning models to identify the best performing model.

The Autoencoder Model:

The project leveraged BigQueryML's native Autoencoder ML model specialized for anomaly detection. An autoencoder is a type of artificial neural network that learns to represent data in a compressed form and then reconstructs it as closely as possible to the original input.

Autoencoders consists of two components:

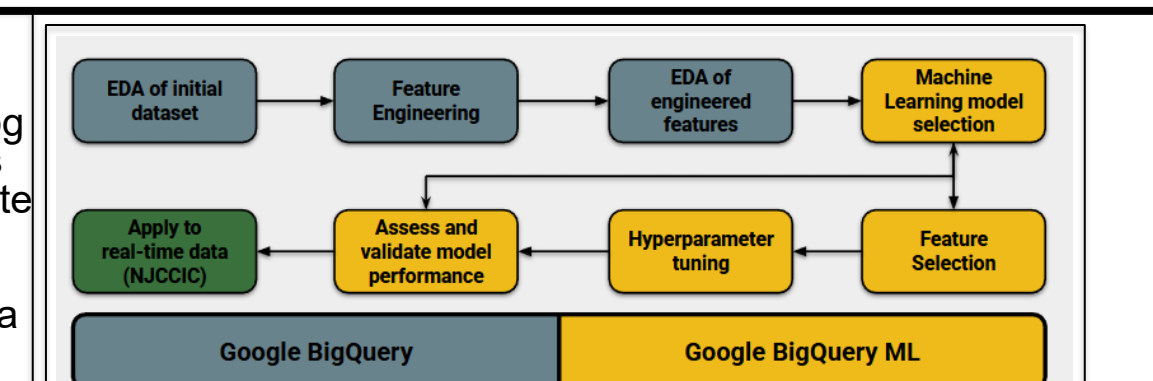
Encoder: This compresses the input into a compact representation and capture the most relevant information.

Decoder: It reconstructs the input data from this compressed form to make it as similar as possible to the original input.

The MSE values are a calculation of the difference between the input and the output based on similarity.

Autoencoder4/4a Features

Field Name	Type	Media	Description	gsearch_order	NTSL	NULLABLE
baseurlid_0m	INT64	NULLABLE	Google Analytics baseurlid_0m		NTSL	NULLABLE
MFA_enabled_reentry_0m	INT64	NULLABLE	MFA_enabled_reentry_0m		NTSL	NULLABLE
MFA_enabled_2m_registration_0m	INT64	NULLABLE	MFA_enabled_2m_registration_0m		NTSL	NULLABLE
Registration_0m	INT64	NULLABLE	Registration_0m		NTSL	NULLABLE
account_age_0	INT64	NULLABLE	account_age_0		NTSL	NULLABLE
city_of_work_0	INT64	NULLABLE	city_of_work_0		NTSL	NULLABLE
day_of_week_0	INT64	NULLABLE	day_of_week_0		NTSL	NULLABLE
device_id_device_per_user_0_0	INT64	NULLABLE	device_id_device_per_user_0_0		NTSL	NULLABLE
device_id_device_per_user_0_1	INT64	NULLABLE	device_id_device_per_user_0_1		NTSL	NULLABLE
device_id_device_per_user_0_2	INT64	NULLABLE	device_id_device_per_user_0_2		NTSL	NULLABLE
device_id_device_per_user_0_3	INT64	NULLABLE	device_id_device_per_user_0_3		NTSL	NULLABLE
device_id_device_per_user_0_4	INT64	NULLABLE	device_id_device_per_user_0_4		NTSL	NULLABLE
device_id_device_per_user_0_5	INT64	NULLABLE	device_id_device_per_user_0_5		NTSL	NULLABLE
device_id_device_per_user_0_6	INT64	NULLABLE	device_id_device_per_user_0_6		NTSL	NULLABLE
device_id_device_per_user_0_7	INT64	NULLABLE	device_id_device_per_user_0_7		NTSL	NULLABLE
device_id_device_per_user_0_8	INT64	NULLABLE	device_id_device_per_user_0_8		NTSL	NULLABLE
device_id_device_per_user_0_9	INT64	NULLABLE	device_id_device_per_user_0_9		NTSL	NULLABLE
device_id_device_per_user_0_10	INT64	NULLABLE	device_id_device_per_user_0_10		NTSL	NULLABLE
device_id_device_per_user_0_11	INT64	NULLABLE	device_id_device_per_user_0_11		NTSL	NULLABLE
device_id_device_per_user_0_12	INT64	NULLABLE	device_id_device_per_user_0_12		NTSL	NULLABLE
device_id_device_per_user_0_13	INT64	NULLABLE	device_id_device_per_user_0_13		NTSL	NULLABLE
device_id_device_per_user_0_14	INT64	NULLABLE	device_id_device_per_user_0_14		NTSL	NULLABLE
device_id_device_per_user_0_15	INT64	NULLABLE	device_id_device_per_user_0_15		NTSL	NULLABLE
device_id_device_per_user_0_16	INT64	NULLABLE	device_id_device_per_user_0_16		NTSL	NULLABLE
device_id_device_per_user_0_17	INT64	NULLABLE	device_id_device_per_user_0_17		NTSL	NULLABLE
device_id_device_per_user_0_18	INT64	NULLABLE	device_id_device_per_user_0_18		NTSL	NULLABLE
device_id_device_per_user_0_19	INT64	NULLABLE	device_id_device_per_user_0_19		NTSL	NULLABLE
device_id_device_per_user_0_20	INT64	NULLABLE	device_id_device_per_user_0_20		NTSL	NULLABLE
device_id_device_per_user_0_21	INT64	NULLABLE	device_id_device_per_user_0_21		NTSL	NULLABLE
device_id_device_per_user_0_22	INT64	NULLABLE	device_id_device_per_user_0_22		NTSL	NULLABLE
device_id_device_per_user_0_23	INT64	NULLABLE	device_id_device_per_user_0_23		NTSL	NULLABLE
device_id_device_per_user_0_24	INT64	NULLABLE	device_id_device_per_user_0_24		NTSL	NULLABLE
device_id_device_per_user_0_25	INT64	NULLABLE	device_id_device_per_user_0_25		NTSL	NULLABLE
device_id_device_per_user_0_26	INT64	NULLABLE	device_id_device_per_user_0_26		NTSL	NULLABLE
device_id_device_per_user_0_27	INT64	NULLABLE	device_id_device_per_user_0_27		NTSL	NULLABLE
device_id_device_per_user_0_28	INT64	NULLABLE	device_id_device_per_user_0_28		NTSL	NULLABLE
device_id_device_per_user_0_29	INT64	NULLABLE	device_id_device_per_user_0_29		NTSL	NULLABLE
device_id_device_per_user_0_30	INT64	NULLABLE	device_id_device_per_user_0_30		NTSL	NULLABLE
device_id_device_per_user_0_31	INT64	NULLABLE	device_id_device_per_user_0_31		NTSL	NULLABLE
device_id_device_per_user_0_32	INT64	NULLABLE	device_id_device_per_user_0_32		NTSL	NULLABLE
device_id_device_per_user_0_33	INT64	NULLABLE	device_id_device_per_user_0_33		NTSL	NULLABLE
device_id_device_per_user_0_34	INT64	NULLABLE	device_id_device_per_user_0_34		NTSL	NULLABLE
device_id_device_per_user_0_35	INT64	NULLABLE	device_id_device_per_user_0_35		NTSL	NULLABLE
device_id_device_per_user_0_36	INT64	NULLABLE	device_id_device_per_user_0_36		NTSL	NULLABLE
device_id_device_per_user_0_37	INT64	NULLABLE	device_id_device_per_user_0_37		NTSL	NULLABLE
device_id_device_per_user_0_38	INT64	NULLABLE	device_id_device_per_user_0_38		NTSL	NULLABLE
device_id_device_per_user_0_39	INT64	NULLABLE	device_id_device_per_user_0_39		NTSL	NULLABLE
device_id_device_per_user_0_40	INT64	NULLABLE	device_id_device_per_user_0_40		NTSL	NULLABLE
device_id_device_per_user_0_41	INT64	NULLABLE	device_id_device_per_user_0_41		NTSL	NULLABLE
device_id_device_per_user_0_42	INT64	NULLABLE	device_id_device_per_user_0_42		NTSL	NULLABLE
device_id_device_per_user_0_43	INT64	NULLABLE	device_id_device_per_user_0_43		NTSL	NULLABLE
device_id_device_per_user_0_44	INT64	NULLABLE	device_id_device_per_user_0_44		NTSL	NULLABLE
device_id_device_per_user_0_45	INT64	NULLABLE	device_id_device_per_user_0_45		NTSL	NULLABLE
device_id_device_per_user_0_46	INT64	NULLABLE	device_id_device_per_user_0_46		NTSL	NULLABLE
device_id_device_per_user_0_47	INT64	NULLABLE	device_id_device_per_user_0_47		NTSL	NULLABLE
device_id_device_per_user_0_48	INT64	NULLABLE	device_id_device_per_user_0_48		NTSL	NULLABLE
device_id_device_per_user_0_49	INT64	NULLABLE	device_id_device_per_user_0_49		NTSL	NULLABLE
device_id_device_per_user_0_50	INT64	NULLABLE	device_id_device_per_user_0_50		NTSL	NULLABLE
device_id_device_per_user_0_51	INT64	NULLABLE	device_id_device_per_user_0_51		NTSL	NULLABLE
device_id_device_per_user_0_52	INT64	NULLABLE	device_id_device_per_user_0_52		NTSL	NULLABLE
device_id_device_per_user_0_53	INT64	NULLABLE	device_id_device_per_user_0_53		NTSL	NULLABLE
device_id_device_per_user_0_54	INT64	NULLABLE	device_id_device_per_user_0_54		NTSL	NULLABLE
device_id_device_per_user_0_55	INT64	NULLABLE	device_id_device_per_user_0_55		NTSL	NULLABLE
device_id_device_per_user_0_56	INT64	NULLABLE	device_id_device_per_user_0_56		NTSL	NULLABLE
device_id_device_per_user_0_57	INT64	NULLABLE	device_id_device_per_user_0_57		NTSL	NULLABLE
device_id_device_per_user_0_58	INT64	NULLABLE	device_id_device_per_user_0_58		NTSL	NULLABLE
device_id_device_per_user_0_59	INT64	NULLABLE	device_id_device_per_user_0_59		NTSL	NULLABLE
device_id_device_per_user_0_60	INT64	NULLABLE	device_id_device_per_user_0_60		NTSL	NULLABLE
device_id_device_per_user_0_61	INT64	NULLABLE	device_id_device_per_user_0_61		NTSL	NULLABLE
device_id_device_per_user_0_62	INT64	NULLABLE	device_id_device_per_user_0_62		NTSL	NULLABLE
device_id_device_per_user_0_63	INT64	NULLABLE	device_id_device_per_user_0_63		NTSL	NULLABLE
device_id_device_per_user_0_64	INT64	NULLABLE	device_id_device_per_user_0_64		NTSL	NULLABLE
device_id_device_per_user_0_65	INT64	NULLABLE	device_id_device_per_user_0_65		NTSL	NULLABLE
device_id_device_per_user_0_66	INT64	NULLABLE	device_id_device_per_user_0_66		NTSL	NULLABLE
device_id_device_per_user_0_67	INT64	NULLABLE	device_id_device_per_user_0_67		NTSL	NULLABLE
device_id_device_per_user_0_68	INT64	NULLABLE	device_id_device_per_user_0_68		NTSL	NULLABLE
device_id_device_per_user_0_69	INT64	NULLABLE	device_id_device_per_user_0_69		NTSL	NULLABLE
device_id_device_per_user_0_70	INT64	NULLABLE	device_id_device_per_user_0_70		NTSL	NULLABLE
device_id_device_per_user_0_71	INT64	NULLABLE	device_id_device_per_user_0_71		NTSL	NULLABLE
device_id_device_per_user_0_72	INT64	NULLABLE	device_id_device_per_user_0_72		NTSL	NULLABLE
device_id_device_per_user_0_73	INT64	NULLABLE	device_id_device_per_user_0_73		NTSL	NULLABLE
device_id_device_per_user_0_74	INT64	NULLABLE	device_id_device_per_user_0_74		NTSL	NULLABLE
device_id_device_per_user_0_75	INT64	NULLABLE	device_id_device_per_user_0_75		NTSL	NULLABLE
device_id_device_per_user_0_76	INT64	NULLABLE	device_id_device_per_user_0_76		NTSL	NULLABLE
device_id_device_per_user_0_77	INT64	NULLABLE	device_id_device_per_user_0_77		NTSL	NULLABLE
device_id_device_per_user_0_78	INT64	NULLABLE	device_id_device_per_user_0_78		NTSL	NULLABLE
device_id_device_per_user_0_79	INT64	NULLABLE	device_id_device_per_user_0_79		NTSL	NULLABLE
device_id_device_per_user_0_80	INT64	NULLABLE	device_id_device_per_user_0_80		NTSL	NULLABLE
device_id_device_per_user_0_81	INT64	NULLABLE	device_id_device_per_user_0_81		NTSL	NULLABLE
device_id_device_per_user_0_82	INT64	NULLABLE	device_id_device_per_user_0_82		NTSL	NULLABLE
device_id_device_per_user_0_83	INT64	NULLABLE	device_id_device_per_user_0_83		NTSL	NULLABLE
device_id_device_per_user_0_84	INT64	NULLABLE	device_id_device_per_user_0_84		NTSL	NULLABLE
device_id_device_per_user_0_85	INT64	NULLABLE	device_id_device_per_user_0_85		NTSL	NULLABLE
device_id_device_per_user_0_86	INT64	NULLABLE	device_id_device_per_user_0_86		NTSL	NULLABLE
device_id_device_per_user_0_87	INT64	NULLABLE	device_id_device_per_user_0_87		NTSL	NULLABLE
device_id_device_per_user_0_88	INT64	NULLABLE	device_id_device_per_user_0_88		NTSL	NULLABLE
device_id_device_per_user_0_89	INT64	NULLABLE	device_id_device_per_user_0_89		NTSL	NULLABLE
device_id_device_per_user_0_90	INT64	NULLABLE	device_id_device_per_user_0_90		NTSL	NULLABLE
device_id_device_per_user_0_91	INT64	NULLABLE	device_id_device_per_user_0_91		NTSL	NULLABLE
device_id_device_per_user_0_92	INT64	NULLABLE	device_id_device_per_user_0_92		NTSL	NULLABLE
device_id_device_per_user_0_93	INT64	NULLABLE	device_id_device_per_user_0_93		NTSL	NULLABLE
device_id_device_per_user_0_94	INT64	NULLABLE	device_id_device_per_user_0_94		NTSL	NULLABLE
device_id_device_per_user_0_95	INT64	NULLABLE	device_id_device_per_user_0_95		NTSL	NULLABLE
device_id_device_per_user_0_96	INT64	NULLABLE	device_id_device_per_user_0_96		NTSL	NULLABLE
device_id_device_per_user_0_97	INT64	NULLABLE	device_id_device_per_user_0_97		NTSL	NULLABLE
device_id_device_per_user_0_98	INT64	NULLABLE	device_id_device_per_user_0_98		NTSL	NULLABLE
device_id_device_per_user_0_99	INT64	NULLABLE	device_id_device_per_user_0_99		NTSL	NULLABLE



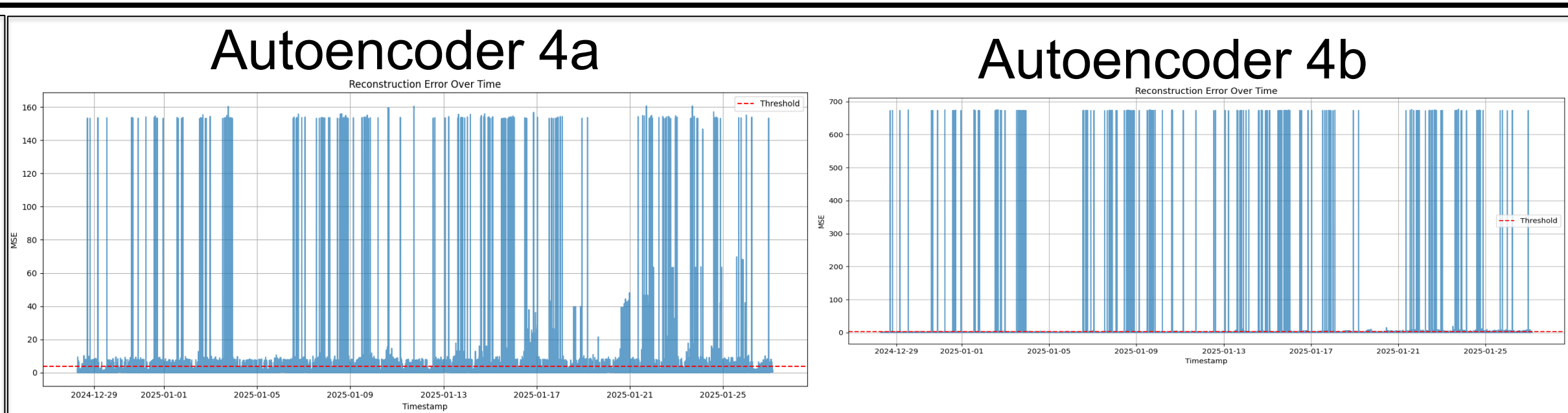
Initial Features: The initial feature set was used to derive new feature sets. Autoencoder 4/4a used the full set of features while Autoencoder 4b used a greatly reduced set of features.

Initial Training Results

Model	Mean absolute error	Mean squared error	Mean squared log error
Autoencoder4 model	0.3581	0.5473	0.0310
Autoencoder4a model	0.1355	0.2150	0.0207
Autoencoder4b model	0.1239	0.2055	0.0132

Results

The Autoencoder model operates through the DETECT_ANOMALIES method within Google's BigQueryML platform which calculates the MSE values for each login event in the dataset. A contamination value is needed to determine what is anomalous. The contamination value was determined by examining the percentile distribution of the MSE values so that 99% of the data was below the threshold. The model works best when it only marks the exceptionally high MSE values. The MSE values and contamination threshold were visualized over time to determine which model performed better. It was determined that model 4b, run on the 4b feature set, performed the best due to the lack of noise at the threshold. The charts below were derived from the results of model 4b.

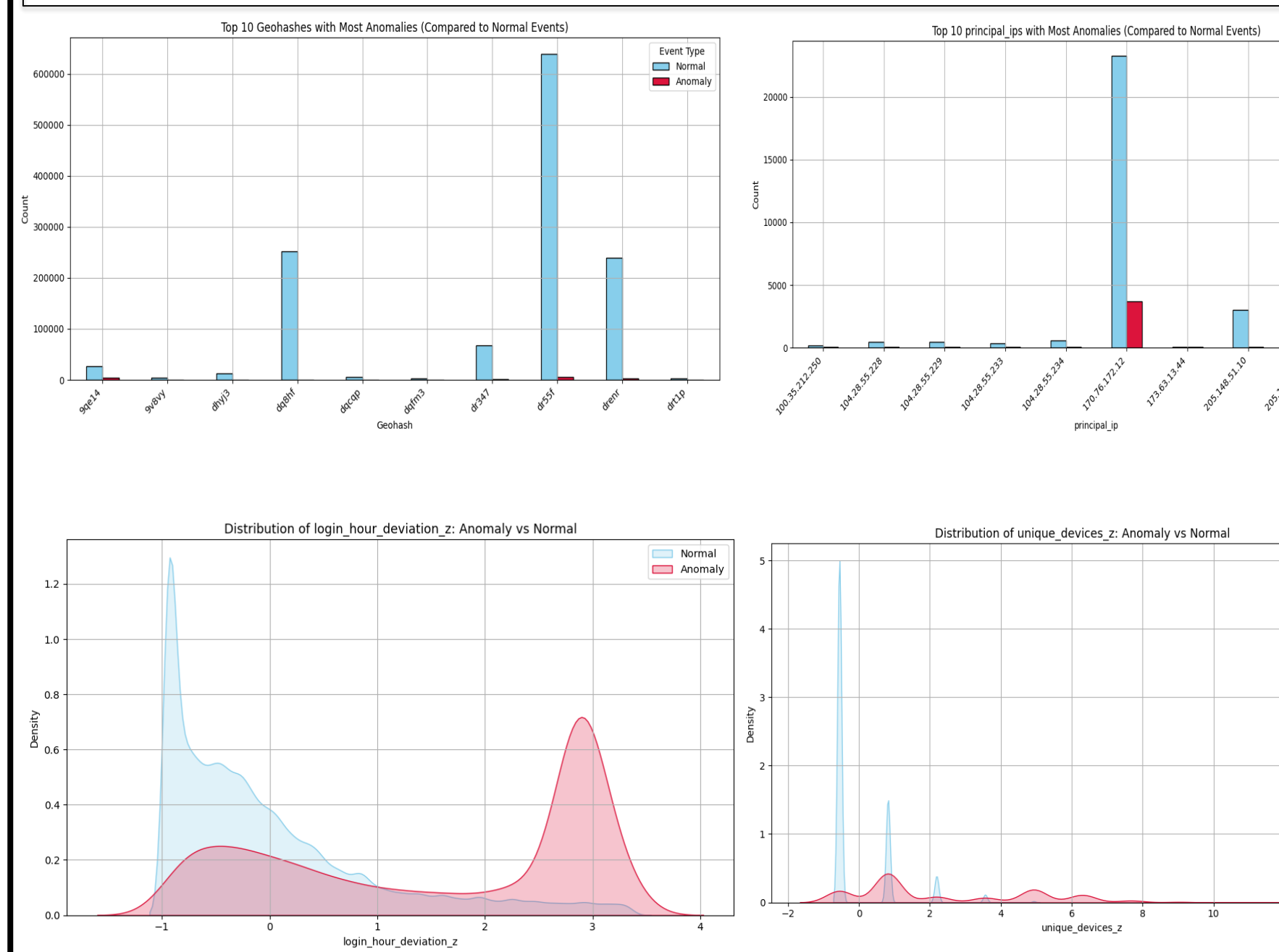


Autoencoder4a Percentiles:

90th percentile MSE: 0.77643
95th percentile MSE: 1.65557
97.5th percentile MSE: 2.55098
99th percentile MSE: 3.81254

Autoencoder4b Percentiles:

90th percentile MSE: 0.56305
95th percentile MSE: 1.35461
97.5th percentile MSE: 1.86643
99th percentile MSE: 2.21938



Future Recommendations:
 Use a select number of features - Identify a subset of features that are easily understood by people. Including too many features harms the interpretability and effectiveness of the model.
 Verify the Model - Apply our model to an unseen dataset with instances that are specifically marked as an anomalous login. This will provide the best determination of model efficacy and accuracy.
 Verify Anomalous Results - From our work in this project, things like when and how someone logs in, plus whether they use multi-factor authentication, help us spot unusual login activity.

References:
 1. New Jersey Cybersecurity & Communications Integration Cell. (n.d.). New Jersey Cybersecurity & Communications Integration Cell. Retrieved from <https://www.cyber.nj.gov/>
 2. Gupta, A. (2025, December 23). Autoencoders in Machine Learning. GeeksforGeeks. Retrieved from <https://www.geeksforgeeks.org/machine-learning/auto-encoders/>

Acknowledgements:
 This project was conducted through the Rutgers MBS Externship Exchange in collaboration with the New Jersey Cybersecurity & Communications Integration Cell. Special thanks to Dr. Karen Bemis, Abbe Rosenthal, Rob Bruder, Aditi Shah, Aditya Patil, Indira Rivera, Harin Ponna, and Vatsal Malkari.

Logan Buddenbaum
Rutgers University



Enhancing NJCCIC Threat Detection: A Machine Learning Approach to Anomalous Login Behavior

Abstract

State-level cybersecurity operations centers (SOCs) face increasing volumes of authentication data, making manual review of potential compromises impossible. This project, conducted in partnership with the New Jersey Cybersecurity & Communications Integration Cell (NJCCIC), developed a machine learning-based approach to detect anomalous login behaviors. By leveraging Google Cloud Platform for data engineering and applying anomaly detection algorithms to custom feature sets, including MFA status and geolocation, the team successfully identified potential malicious activity patterns that static rules and the human eye often miss. The results demonstrate that behavioral analytics provides a higher fidelity signal for threat detection, allowing SOC analysts to prioritize high-risk events more effectively.

Logan Buddenbaum
Rutgers University

