

Summary of Agency Compliance Reporting of Algorithmic Tools CY 2023

Overview

This summary communicates the results of New York City’s fourth annual process for reporting on algorithmic tools. Under Local Law 35 (“LL 35”)¹, the city continues its commitment to provide the public with a transparent view of these applications of agency data and technology. Pursuant to mayoral Executive Order 3 of 2022 (“EO 3”)², the city’s Office of Technology & Innovation (“OTI”) manages this process, providing guidance to agencies and ensuring that agency materials are prepared for the public.

In October 2023, the city published the New York City Artificial Intelligence (“AI”) Action Plan³, which contains a commitment to expand public reporting of algorithmic tools, many of which are considered types of AI, to “make resulting reports readily accessible to the public through the Open Data platform.” This year, reporting is also available as an open data set at <https://data.cityofnewyork.us/d/jaw4-yuem>.

Update 3/27/2024

The Department of Social Services updated their reporting to include changes made to the Homebase Risk Assessment Questionnaire after initial publication of the report.

Key Changes for 2023

Reporting for 2023 added four new elements. A full list of reportable elements is available in the Appendix.

Summary of Agency Reports

The following table (continued on next page) summarizes the reporting from City agencies for 2023.

| Agency | Number of Tools Reported |
|---|--------------------------|
| Administration for Children's Services (ACS) | 5 |
| Business Integrity Commission (BIC) | 0 |
| Commission on Human Rights (CCHR) | 0 |
| Civilian Complaint Review Board (CCRB) | 0 |
| Conflicts of Interest Board (COIB) | 0 |
| Department of Citywide Administrative Services (DCAS) | 0 |
| Department of Cultural Affairs (DCLA) | 0 |
| Department of City Planning (DCP) | 0 |
| Department of Consumer and Worker Protection (DCWP) | 1 |
| Department of Design and Construction (DDC) | 0 |
| Department of Environmental Protection (DEP) | 1 |
| Department for the Aging (DFTA) | 0 |

¹ For the full text of LL 35, see:

<https://legistar.council.nyc.gov/LegislationDetail.aspx?ID=4265421&GUID=FBA29B34-9266-4B52-B438-A772D81B1CB5>

² For the full text of EO 3, see: <https://www.nyc.gov/office-of-the-mayor/news/003-002/executive-order3>

³ For the full text of the AI Action Plan, see: <https://www.nyc.gov/assets/oti/downloads/pdf/reports/artificial-intelligence-action-plan.pdf>

| | |
|--|-----------|
| Department of Buildings (DOB) | 0 |
| Department of Correction (DOC) | 0 |
| Department of Education (DOE) | 5 |
| Department of Finance (DOF) | 0 |
| Department of Health and Mental Hygiene (DOHMH) | 16 |
| Department of Investigation (DOI) | 1 |
| Department of Probation (DOP) | 0 |
| Department of Records and Information Services (DORIS) | 0 |
| Department of Transportation (DOT) | 0 |
| Department of Parks & Recreation (DPR) | 0 |
| Department of Sanitation (DSNY) | 0 |
| Department of Social Services (DSS)* | 3 |
| Department of Veterans' Services (DVS) | 0 |
| Department of Youth and Community Development (DYCD) | 0 |
| Economic Development Corporation (EDC) | 0 |
| Fire Department (FDNY) | 4 |
| Health + Hospitals (H+H) | 1 |
| Department of Housing Preservation and Development (HPD) | 0 |
| Law Department (LAW) | 0 |
| Landmarks Preservation Commission (LPC) | 0 |
| Mayor's Office (MO) | 3 |
| NYC Emergency Management (NYCEM) | 0 |
| New York City Housing Authority (NYCHA) | 0 |
| New York Police Department (NYPD) | 3 |
| Office of Administrative Trials and Hearings (OATH) | 0 |
| Office of Chief Medical Examiner (OCME) | 1 |
| Office of Technology and Innovation (OTI) | 1 |
| Public Design Commission (PDC) | 0 |
| Department of Small Business Services (SBS) | 0 |
| School Construction Authority (SCA) | 1 |
| Taxi & Limousine Commission (TLC) | 0 |
| Total | 46 |

*Note that the set of offices included in the Department of Social Services' reporting has been updated for 2023 to align with other citywide reporting efforts. Notably, the Public Engagement Unit, which previously reported under the Mayor's Office, now reports as part of the Department of Social Services.

Table of Contents

- 1.0 Administration for Children’s Services (ACS).....1
- 2.0 Department of Consumer and Worker Protection (DCWP) 7
- 3.0 Department of Environmental Protection (DEP).....8
- 4.0 Department of Education (DOE).....9
- 5.0 Department of Health and Mental Hygiene (DOHMH)14
- 6.0 Department of Investigations (DOI)26
- 7.0 Department of Social Services (DSS)27
- 8.0 Fire Department (FDNY)31
- 9.0 Health + Hospitals (H+H).....34
- 10.0 Mayor’s Office (MO)35
- 11.0 Police Department (NYPD).....38
- 12.0 Office of the Chief Medical Examiner (OCME)41
- 13.0 Office of Technology & Innovation (OTI)42
- 14.0 School Construction Authority (SCA).....43
- Appendix45

1.0 Administration for Children’s Services (ACS)
1.1 Accelerated Safety Analysis Protocol (ASAP) Tool

First Used: May 2018

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: No |

Tool Description

Predictions of Severe Harm (identifying likelihood of substantiated allegations of physical or sex abuse within the next 24 months) are based on a machine learning methodology and are calculated for all children involved in active investigations early in the investigation (day 10). An investigation is assigned a numeric likelihood of this outcome based on the child in the case with the highest likelihood. The ACS Quality Assurance unit in the Division of Child Protection reviews about 3,000 active investigations annually, selecting those with the highest likelihood of severe harm.

Tool Purpose & Description of Use

The Quality Assurance Unit in the Division of Child Protection at ACS has the capacity to review about 3,000 investigation cases out of about 50,000 investigations annually. ACS developed this predictive model to support the selection of cases for Quality Assurance review. Open investigations involving children with the greatest likelihood to experience future severe harm – substantiated allegations of physical or sex abuse in the following 24 months – are selected for review. The tool does not support decisions about services or interventions for individuals or families involved with ACS, beyond the selection of the case for this additional Quality Assurance review.

If the Quality Assurance review team identifies gaps in routine, required documentation or practice, the team speaks with the field office conducting the investigation and follows up to make certain these gaps have been addressed. Scores are not shared with staff in the Quality Assurance unit or the investigative unit. The model only supports the decision about which investigations are prioritized for review by the Quality Assurance (QA) unit.

Data Analyzed

- Training Data** ACS trained the model on ACS historic administrative data about closed investigations from April 2014 to April 2016. The training set included about 142,026 observations. The model was tested on closed investigations from April 2016 to April 2017 with 53,477 observations.
- Input Data** Predictions are based on administrative data about prior and current child welfare involvement including investigations triggered by a New York State Central Register (SCR) call and time spent in foster care. Only ACS administrative data are used in the model.
- Output Data** Rank ordered list of open investigation cases involving children with the highest likelihood to experience future severe harm, defined as substantiated allegations of physical or sex abuse in the following 24 months to be reviewed by a special QA Review Team.

Vendor Involvement

Vendor Name: N/A

1.2 Service Termination Conference (STC)

First Used: July 2017

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: No |

Tool Description

Predictions of Repeat Maltreatment (identifying the likelihood of being involved in a future indicated investigation within the following 24 months towards the end of service) are based on a machine learning methodology and are calculated for all children receiving prevention services from ACS prevention service providers. The prediction is made with the assumption that the case is closed the day

the model is run. A prevention case is assigned a numeric likelihood of an indicated investigation based on a New York State Central Register (SCR) within 24 months from the end of a prevention service.

Tool Purpose & Description of Use

The Repeat Maltreatment model was known as the Service Termination Conference (STC) model and was used by Preventive Services managers to identify whether or not a case should have ACS or provider-agency facilitation at the service termination conference. The model was initially used to prioritize ACS facilitation of prevention termination conferences so that ACS could be certain that services had been provided in these cases. When a family is ready to exit ACS prevention services, an end-of-services conference is required (known as a “Service Termination Conference”). However, Service Termination Conferences are generally no longer facilitated by ACS and this prioritization is no longer necessary.

In November 2022, a pilot was initiated with the STC list to recommend low-risk families to prevention service providers in case they were interested in discontinuing service. This pilot was stopped in January 2023.

Data Analyzed

- Training Data** ACS trained the model on ACS historic administrative data about closed investigations from July 2009 to December 2015. Training set included about 130,982 observations. The model was tested on closed investigations from January 2016 to December 2018 with 48,771 observations.
- Input Data** Predictions are based on administrative data about prior and current child welfare involvement at the end of a prevention service case. It includes SCR investigations and time spent in foster care. Only ACS administrative data are used in the model.
- Output Data** As of September 2021, ACS is no longer required to facilitate prevention termination conferences (all conferences are facilitated by the contract prevention programs) and the STC model is no longer being used for this purpose.

Vendor Involvement

Vendor Name: N/A

1.3 Prevention Score Card

First Used: September 2021

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: Yes |

Tool Description

Predictions of Repeat Maltreatment (identifying the likelihood of being involved in a future indicated investigation within the next 24 months at the start of service) are based on a machine learning methodology and are calculated for all children receiving prevention services from ACS prevention service providers.

Tool Purpose & Description of Use

The Repeat Maltreatment model is used to make predictions on day 10 from the start of the prevention case to assess the risk of the family at the beginning of the service. A prevention case is assigned a numeric likelihood of an indicated investigation based on a New York State Central Register (SCR) within 24 months from the start of a prevention service.

The prevention providers are assessed for their performance based on the service needs/risk levels of the families they've served during the previous fiscal year.

The programs were sorted and ranked based on their average risk, and then divided into four quartiles by rank order: the top 25 percent of programs are classified as the Very High-Risk Cohort, the next 25 percent of programs as the High-Risk Cohort, the next 25 percent as the Medium-Risk Cohort, and the lowest 25 percent as the Low-Risk Cohort. Assignment to a cohort is not a way of performance assessment of the program but to group prevention service providers for fair comparisons based on the risk level of families they served.

Data Analyzed

Training Data ACS trained the model on ACS historic administrative data about closed investigations from July 2009 to June 2016. Training set included about 158,787 observations. The model was tested on closed investigations from July 2016 to June 2018 with 46,969 observations.

Input Data Predictions are based on administrative data about prior and current child welfare involvement at the start of a case. This includes SCR investigations and time spent in foster care. Only ACS administrative data are used in the model.

Output Data The model is used for generating a scorecard of prevention service providers by categorizing prevention programs based on the average risk profile of the cases they served during the assessment year. These groupings of program cohorts provide context for understanding the Scorecard, as it allows for performance comparison of programs that accepted and served families with similar risk profiles.

Vendor Involvement

Vendor Name: N/A

Update Description

In August 2023, ACS retrained the model on recent ACS historic administrative data about closed investigations with an 80/20 split in data from July 2009 to June 2018. The training set included 80

percent (about 338, 467) observations. The model was tested on 20 percent of closed investigations from July 2009 to June 2018 with 84,494 observations.

1.4 Un-Involvement Model

First Used: May 2023

Analysis Type: Predictive modeling

Population Type: Individuals

Identifying Information: Yes

Updated in 2023: No

Tool Description

Predictions of “Un-Involvement” (identifying the likelihood of no future involvement with ACS within 24 months beyond the current investigation) are based on a machine learning methodology and are calculated for all children in an ongoing investigation. This future engagement for families may be in the form of ACS prevention services, court-ordered supervision, or foster care services. There are two models run on day 10 and day 40 to ensure that a case that was recommended for early closure on day 10 is still eligible and recommended on day 40, right before actually closing the case.

Tool Purpose & Description of Use

The model helps Child Protection managers identify low-risk cases that are likely not to require further ACS involvement beyond the “current” investigation and therefore could be considered for early closure, sooner than the typical 60 days.

The model generates initial risk predictions on the 10th day of a new investigation to identify low-risk cases. Subsequently, on the 40th day of the investigation, the model re-evaluates these cases with a new set of risk predictions to determine if they continue to be classified as low-risk. If new information collected during the intervening 30 days suggests that a case is no longer eligible or is no longer considered low-risk, the recommendation to close the case will be revised.

Data Analyzed

| | |
|----------------------|--|
| Training Data | ACS trained the model on ACS historic administrative data about closed investigations from 2012 to 2017. An 80/20 split of data to train on 80 percent and test on 20 percent ensuring that no family appears in both sets. The training set contains 381,649 children from 183,516 investigations ending between Jan 2012 and December 2017. The test set consisted of 101,369 children from 48,794 investigations ending between Jan 2012 and December 2017. |
| Input Data | Predictions are based on administrative data about prior and current child welfare involvement at the start of a case. This includes SCR investigations and time spent in foster care. Only ACS administrative data are used in the model. |
| Output Data | A recommendation and description are displayed on both day 10 and day 40 of all open cases via a reporting platform, viewable to only deputy directors. Upon discussion with the |

deputy directors on day 40, the case worker makes a determination to close the case or not. The case workers are not aware that the recommendation is generated by a machine learning model, to not bias the decision-making process. Alternatively, the deputy directors use the list to inform their workload planning conversations with managers and staff, when caseloads are high.

Vendor Involvement

Vendor Name: N/A

1.5 Housing Prioritization

First Used: April 2023

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: No |

Tool Description

The city has allocated 100 housing vouchers to families receiving ACS Prevention Services. The shelter application model identifies the likelihood of a family in prevention services applying for homeless shelter within 12 months beyond the current prevention case. The model uses a machine learning methodology and is calculated for all children in a prevention case. ACS Prevention Services reaches out to the Service Providers assisting the families with the highest risk for applying for shelter.

Tool Purpose & Description of Use

The model estimates a risk score for a child receiving prevention services whose family will apply and be eligible for a homeless shelter within 12 months from the start of service (day 14 of the prevention case).

The model helps predict the risk of application for homeless shelters among families receiving prevention services. With a limited number of vouchers available, the risk model helps ACS prioritize housing assistance for those families at greatest risk of becoming homeless.

The Service Provider meets with the family to conduct a qualitative assessment of the family’s housing needs and vouchers are offered based on their findings.

Data Analyzed

Training Data ACS trained the model on ACS historic administrative data regarding preventive services started between 2014 and 2020. An 80/20 split of data to train on 80 percent and test on 20 percent ensuring that no family appears in both sets. The training set contains 140,242 observations between Jan 2014 and December 2020. The test set consisted of 34,508 observations between Jan 2014 and December 2020.

Input Data Predictions are based on administrative data about prior and current child welfare involvement at the start of a case. This includes SCR investigations and time spent in foster care. Only ACS administrative data are used in the model.

Output Data Rank ordered list of open prevention cases involving children whose families have the highest likelihood of applying for a homeless shelter within 12 months of starting a prevention service.

Vendor Involvement

Vendor Name: N/A

2.0 Department of Consumer and Worker Protection (DCWP)

2.1 Route Automation

First Used: July 2020

| | |
|------------------------------------|--|
| Analysis Type: Optimization | Population Type: Group, organization, or business |
| Identifying Information: No | Updated in 2023: No |

Tool Description

Inspection Supervisor selects an inspector, enters a date and the number of businesses to be inspected, and the geographic area to be considered. The system identifies businesses in the selected area and assigns them to the route based on inspection priority until the number of businesses entered has been reached. Then the tool runs a Simulated Annealing Algorithm to optimize the order businesses appear on the route based on proximity and method of travel.

Tool Purpose & Description of Use

DCWP inspectors conduct inspections based on a route, or list of businesses to be inspected on a specific day, which must be pre-approved by their supervisor. The Route Automation tool generates a route for an inspector on a specific date based on configuration variables and geographic area. All routes generated by the tool still require supervisor review and approval.

Data Analyzed

Training Data The tool is not trained in an AI or Machine Learning sense. The tool makes decisions based on configuration tables, businesses and their licenses (if any), and inspection and violation history, and uses a Simulated Annealing Algorithm to optimize the order in which businesses appear on a route.

Input Data Inspection date, business category and address, licenses held (if any), last inspection date and type, violation history (if any), date of inspection request or license application/renewal (if applicable), Inspection Unit (of the inspector), and geographic area.

Output An ordered list of businesses to be inspected on a given day by a given inspector.

Data

Vendor Involvement

Vendor Name: PruTech

The tool was designed and built by PruTech, an outside vendor contracted to design and build DCWP’s Automated Inspection Management System (AIMS) and its accompanying Mobile Enforcement platform. The tool is part of the AIMS system.

3.0 Department of Environmental Protection (DEP)

3.1 Idling Complaints Program

First Used: August 2022

| | |
|---------------------------------------|-------------------------------------|
| Analysis Type: Computer vision | Population Type: Individuals |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

A contractor helped create an AI tool that analyzes the audio and visual aspects of pictures and videos submitted by citizens of alleged car idling complaint occurrences that are in violation of New York City air pollution laws.

Tool Purpose & Description of Use

The analysis from the tool makes a recommendation to staff reviewers whether the submitted evidence support an occurrence of car idling in violation of New York City laws. The tool also provides a level of confidence in its recommendation. The tool does not make the review decision in the Idling Complaints system. It is still entirely up to the staff to decide whether to take the tool’s recommendation or not.

Data Analyzed

- Training Data** Videos and pictures of cars idling submitted by citizens, along with staff decisions on whether the picture/video constituted as an idling violation.
- Input Data** Videos and pictures submitted by citizens through our web portal.
- Output Data** Recommendation, confidence level, description of its decision from the tool.

Vendor Involvement

Vendor Name: Acuvate

Acuvate developed the AI tool that performs the automated analysis of the submitted evidence.

Update Description

Retraining model with newer data, modified the description text generated to be more useful for reviewer staff.

4.0 Department of Education (DOE)

4.1 MySchools

First Used: August 2018

| | |
|------------------------------------|-------------------------------------|
| Analysis Type: Matching | Population Type: Individuals |
| Identifying Information: No | Updated in 2023: No |

Tool Description

The tool utilizes the Gale-Shapley deferred acceptance algorithm to match applicants to schools. This algorithm has been in existence for many years, used internationally for various purposes. Perhaps most common is its use in the National Resident Matching Program for medical school students.

Deferred acceptance works as an iterative series of steps: students and programs are tentatively matched in each step, but nothing is finalized until the algorithm terminates (hence the deferred).

1. Each student “proposes” to their first choice.

- Programs assign seats to students one at a time.
- When all seats are filled, programs may reject previously accepted students in favor of new applications from students they prefer (e.g., students with a better lottery number).
- Remaining students are rejected.

2. Students rejected in the last step “propose” to the next choice on their list.

3. The algorithm terminates when all students are matched or have proposed to all the programs they listed.

Tool Purpose & Description of Use

MySchools is an application used to house online school directories, collect application choices, and run the admissions matching algorithm that is used for all centralized admissions processes (3K, pre-K, Gifted & Talented, middle school, and high school). The tool encompasses a family-facing portal, a school-facing portal, and an administrative portal.

Data Analyzed

Training Data The algorithm was already widely recognized for its advantages prior to adoption in New York City. The DOE consulted with a team of researchers at MIT who had been closely

involved in its initial creation when we adopted it.

- Input Data** Student biographical information (e.g., home address, poverty status, home language), student academic information (e.g., course grades, state test scores), and student school records (e.g., sending school).
- Output Data** The algorithm outputs a school match for each student.

Vendor Involvement

Vendor Name: Blenderbox

We have a five year contract with the agency Blenderbox who designed the application and implemented the algorithmic matching functionality. The work is meant to transition to be run in-house, by the Division of Instructional and Information Technology (DIIT) within the Department of Education, by the end of the contract. The team at DIIT has already begun to takeover maintenance and development of the tool.

4.2 NYCDOE APPR Measures of Student Learning (MOSL) Growth Model

First Used: September 2013

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: No | Updated in 2023: No |

Tool Description

The growth model uses a variety of student-level (assessment scores; English language learner, disability, and economic disadvantage indicators), classroom-level (e.g. percent students with disabilities), and school-level data (e.g. percent English language learners, percent students with disability, average prior achievement, school type) to estimate/predict a student’s score on one of many possible course-culminating assessments. These predicted scores are used to either 1) identify “peer groups” of students, from which student growth percentiles (SGPs) are determined, or 2) compared to actual scores to determine student credit values. These units (SGPs or credit values) are then weight-averaged to generate an educator-level result - the Measures of Student Learning (MOSL) rating. The MOSL Rating is combined with the Measures of Teaching/Leadership Practice rating (below) to produce an Overall Rating. Per state law 3012-d, annual ratings “shall be a significant factor in HR decisions.” This is often implemented by making ratings a qualifying/disqualifying element in decision-making concerning employment, tenure, salary, and other professional opportunities.

Tool Purpose & Description of Use

In accordance with New York State law and New York State Education Department (NYSED) regulations, the Department developed and maintains a “growth model” to produce Measures of Student Learning (MOSL) ratings for use in annual professional performance reviews (APPR) for teachers and principals.

The MOSL ratings are combined with Measures of Teaching/Leadership Practice (MOTP/MOLP) ratings to produce an annual Overall Rating for each eligible educator.

Data Analyzed

- Training Data** The growth model process is employed in both retrospective and prospective ways. In the retrospective version, the results are determined entirely within-sample. In the prospective version, the coefficients of the model are estimated on multiple prior years of data.
- Input Data** The growth model makes use of three types of data: (1) students’ end-of-year assessment scores, (2) enrollment and attendance records that link students to teachers and schools, and (3) historical academic and demographic information used to identify groups of similar students.
- Output Data** The model outputs an estimate of a student's score on a course-culminating assessment.

Vendor Involvement

Vendor Name: Education Analytics

Education Analytics provides technical assistance and quality assurance for the growth model.

4.3 NYCDOE APPR Measures of Teaching/Leadership Practice (MOTP/MOLP) Calculation

First Used: September 2013

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: No | Updated in 2023: No |

Tool Description

Throughout a school year, evaluators observe teachers/principals multiple times and use a rubric to provide a numerical rating on one or more rubric components. These rubric component scores are then weight-averaged according to collectively bargained rules to produce a Measures of Teaching/Leadership Practice (MOTP/MOLP) rating. The MOTP/MOLP Rating is combined with the Measures of Student Learning rating (above) to produce an Overall Rating for each eligible educator. Per state law 3012-d, annual ratings “shall be a significant factor in HR decisions.” This is often implemented by making ratings a qualifying/disqualifying element in decision-making concerning employment, tenure, salary, and other professional opportunities.

Tool Purpose & Description of Use

In accordance with New York State law and New York State Education Department (NYSED) regulations, the Department developed and maintains databases and calculation rules to produce Measures of Teaching/Leadership Practice (MOTP/MOLP) ratings for use in annual professional performance reviews

(APPR) for teachers and principals. The MOTP/MOLP ratings are combined with Measures of Student Learning (MOSL) ratings to produce an annual Overall Rating for each eligible educator.

Data Analyzed

Training Data Pilot data prior to program launch was used to inform the weights assigned to various rubric components. However, the weights are ultimately determined via collective bargaining.

Input Data Rubric component numerical ratings.

Output Data The model outputs a score for teachers and principals.

Data

Vendor Involvement

Vendor Name: N/A

4.4 Eureka! Chatbot

First Used: August 2023

| | |
|--|-------------------------------------|
| Analysis Type: Speech and language processing | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: No |

Tool Description

The Azure Cognitive Services technology and chatbot (internally branded as "Eureka!") has been configured and deployed in August 2023 to be the first response to calls to the NYC Public Schools (NYCPS) IT Service Desk. It accesses scripts to handle four common reasons for a user to call or contact the service desk: Password Reset, Create a Ticket, Ticket Status, Request for Information. The chatbot accesses pre-defined scripts to respond to user voice or text input. The user's request is either serviced, completed and closed by the chatbot, or the user is given the option (at any time) to connect to a live agent.

Tool Purpose & Description of Use

The tool is used to respond to common IT service desk requests: Password Reset, Create a Ticket, Ticket Status, Request for Information. Users can access the tool by phone, by computer through the DOE Support Hub application, and from links from MS Teams and other DOE systems, such as TeachHub.

Data Analyzed

Training Data Pre-defined scripts designed to respond to four common requests to the IT Service Desk.

Input Data A voice call or text-based chat session initiated by a user and responded to by the Eureka! chatbot before being handled by a human Service Desk agent.

Output Data The chatbot generates responses to user-entered prompts based on the training data, or forwards the inquiry to a human Service Desk agent. Since its launch in August 2022, the chatbot handles an average of 1,500 calls and 300 web-based inquiries each day. Approximately 30 percent of the voice calls and 10 percent of the web-based queries have been handled completely by Eureka! without being forwarded to a human Service Desk agent.

Vendor Involvement

Vendor Name: Nagarro and Microsoft

Developed by an IT services vendor (Nagarro) using Microsoft Cognitive services.

4.5 Open Gen AI and Teaching Assistant Tool

First Used: May 2023

| | |
|--|-------------------------------------|
| Analysis Type: Speech and language processing | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: No |

Tool Description

The generative AI system using large language models was a system custom-built by DIIT using advanced Microsoft technologies to create a set of generative AI tools. To date, two tools have been built. One is named "Open Gen AI" - it accesses a large language model (currently GPT 3.5) to provide responses to a broad range of prompts. The other is named "Teaching Assistant for Algebra" - it accesses specific Algebra-focused content to provide responses to prompts related to Algebra.

Tool Purpose & Description of Use

The tool is used to generate responses to prompts entered by a student or teacher, requesting the generative AI tool to compose a text response to a text input.

Data Analyzed

Training Data For the Open Gen AI tool, the ChatGPT large language model is used as the "trained data." For the Teaching Assistant for Algebra tool, the LLM has been trained exclusively on curriculum from Illustrative Math.

Input Data Prompts provided by the users of the system.

Output Data The output data for the Open Gen AI tool is the response generated by the ChatGPT LLM. The output data for the Teaching Assistant is the response generated by specifically developed LLM using the Illustrative Math curriculum.

Vendor Involvement

Vendor Name: Microsoft

Microsoft provided technical guidance for their emerging generative AI technology and built some small module of code for the specific NYCPS Gen AI and Teaching Assistant use cases.

5.0 Department of Health and Mental Hygiene (DOHMH)

5.1 Improving Foodborne Disease Outbreak Detection by Incorporating Complaints Identified in Social Media Data

First Used: November 2016

| | |
|--|---|
| Analysis Type: Speech and language processing | Population Type: Individuals, Group, organization, or business |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

Restaurant associated foodborne disease outbreaks are often identified through complaints received via New York City’s 311 non-emergency information system, however not all individuals report to 311. The New York City Department of Health and Mental Hygiene (NYC DOHMH) in collaboration with Columbia University developed a text classifier program which monitors Yelp and Twitter data to identify complaints of foodborne illness which was supported by grants from the Alfred P Sloan Foundation and the National Science Foundation. As of April 2023, the tool no longer uses data from Twitter/X due to API changes.

Tool Purpose & Description of Use

The model uses data from Yelp restaurant reviews and previously used Twitter data that was available on Twitter’s publicly available API. Twitter (X) removed free access to their publicly available API in April 2023, so these data are no longer included in our analyses. The classifiers assign a “sick score” to each Yelp review or tweet indicating the likelihood that the review or tweet pertains to foodborne illness. The sick score is based on whether the review/tweet contains key words indicative of foodborne illness (“e.g. vomit”); the Yelp classifier also incorporates if the review indicates that multiple people became sick and if the review indicates a time between eating at a restaurant and illness onset (incubation period) that is consistent with foodborne illness. Each review and tweet with a sick score greater than or equal to a threshold value are reviewed and annotated by DOHMH foodborne disease epidemiology and environmental health staff to determine if the review/tweet was actually reporting foodborne illness possibly associated with a New York City restaurant; if yes, Yelp messages are sent to Yelp reviewers, requesting that they contact DOHMH, and a Twitter message with a survey link was tweeted back to Twitter users to confirm foodborne illness. Data from annotations are used to improve classifier performance. Foodborne disease complaints identified through Yelp and previously Twitter are

combined with foodborne disease complaints reported to 311 to improve efficiency of outbreak detection.

Data Analyzed

- Training Data** Training data was used in the development of both the Yelp and Twitter classifiers. The training data consisted of restaurant reviews and tweets obtained, respectively, from Yelp and Twitter by Columbia University; a subset of these data were joined with annotations provided by DOHMH staff. The annotations of restaurant reviews focused on the following: 1) if the review indicated foodborne illness, 2) if the incident occurred in the past 30 days, 3) if multiple people were sick and 4) if the incubation period was consistent with foodborne illness. For tweets, the annotations focused on if the tweet was indicating foodborne illness and if the incident occurred in New York City. The training data is periodically updated (with annotations from DOHMH) to improve the classifiers.
- Input Data** Yelp reviews of New York City restaurants are pulled from a privately available application programming interface (API) provided by Yelp. Publicly available data from Twitter’s API was used through April 2023.
- Output Data** The output data includes a “sick score” that the classifiers assign to each Yelp review or tweet (when that data were being used) indicating the likelihood that the review or tweet pertains to foodborne illness.

Vendor Involvement

Vendor Name: Columbia University

DOHMH staff, including Bureau of Communicable Disease, Office of Environmental Investigations, and Division of Informatics and Information Technology & Telecommunications and Columbia University are involved in making decisions about the tool. Columbia University Department of Computer Science professors and doctoral students maintain the classifier. The project was previously funded by the Alfred P Sloan Grant, for which The Fund for Public Health in New York provided administrative support and grant management to DOHMH. This support and management ended at the completion of the grant in 2021.

Update Description

As of April 2023, the tool no longer uses data from Twitter/X due to API changes.

5.2 Guppy

First Used: June 2020

| | |
|---|---|
| Analysis Type: Predictive modeling | Population Type: Individuals, Other: Sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: No |

Tool Description

Converts electric signals to predict a nucleotide and enables filtering of low-quality calls.

Tool Purpose & Description of Use

This is a tool designed specifically for Oxford Nanopore Technology (ONT) data. This is a neural network based basecaller, a tool that determines nucleotide bases of a genetic material, that converts electric signals into strings to represent genomic data. In addition to basecalling, the tool also performs filtering of low-quality reads, a stretch of sequenced genetic material. This is the initial step that converts electric signals to fragments of sequence data, which can then be used for COVID-19 sequencing analysis.

Data Analyzed

- Training Data** The default models within Guppy are trained on a mixture of native and amplified DNA/RNA, from multiple organisms including plant, animal, bacterial and viral genomes.
- Input Data** DNA/RNA strand passing through the nanopore. Raw data is stored as .fast5 files
- Output Data** .fast5 files, fastq, or BAM files.

Vendor Involvement

Vendor Name: Oxford Nanopore Technologies

Developed and maintains the tool.

5.3 BWA

First Used: July 2017

| | |
|---|---|
| Analysis Type: Predictive modeling, Matching | Population Type: Individuals, Other: Sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: No |

Tool Description

Aligns sequencing data to a reference sequence.

Tool Purpose & Description of Use

Burrows-Wheeler Aligner (BWA) is aligning sequence data to reference using Burrows-Wheeler transformations. This tool is optimal for low-divergent genomic data and short read data, such as Illumina sequence data. This tool is used to predict the order in which the fragments generated by sequencers are pieced together to form a complete genomic sequence data. This tool is used for Legionella and PulseNet sequencing analyses.

Data Analyzed

Training N/A

Data

Input Data Sequence reads (fastq) for single or paired-end runs (sequence reads can be considered strings).

Output Data Aligned reads in SAM format.

Vendor Involvement

Vendor Name: N/A

5.4 Minimapp2

First Used: May 2020

| | |
|---|---|
| Analysis Type: Predictive modeling, Matching | Population Type: Individuals, Other: Sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

Aligns sequencing data to a reference sequence.

Tool Purpose & Description of Use

Minimapp2 uses optimal chaining scores to align sequencing data to reference genomes. This tool is faster and more optimal for long read sequences, such as Oxford Nanopore Technologies (ONT) data. This tool is used to predict the order in which the fragments generated by sequencers are pieced together to form a complete genomic sequence data. This tool is used for COVID-19 and monkeypox virus (MPXV) sequencing analyses.

Data Analyzed

Training N/A

Data

Input Data Sequence reads (fastq) for single or paired-end runs (sequence reads can be considered strings)

Output Data Aligned reads in SAM format

Vendor Involvement

Vendor Name: N/A

Update Description

Fixed the broken Python package. Updated variable weights/bug fixes.

5.5 Pangolin

First Used: July 2021

| | |
|---|---|
| Analysis Type: Predictive modeling | Population Type: Individuals, Other: sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

Assigns lineage names to SARS-CoV-2.

Tool Purpose & Description of Use

Pangolin uses a combination of several methods, including random forest tree, classification methods, and maximum parsimony to assign lineage names to SARS-CoV-2 genomic sequences to bin sequences that are more likely to be similar. This is a tool that designates a name based on a nomenclature for COVID-19 sequence data.

Data Analyzed

- Training Data** Trained on a data set of genomes that have been designated to Pango lineages using whole genome information.
- Input Data** Fasta files.
- Output Data** .csv file with taxon name and lineage assigned.

Vendor Involvement

Vendor Name: N/A

Update Description

Updated variable weights/bug fixes/dependencies.

5.6 MAFFT

First Used: January 2021

| | |
|---|---|
| Analysis Type: Predictive modeling, Matching | Population Type: Individuals, Other: sequence data can belong to any species |
|---|---|

| | |
|------------------------------------|----------------------------|
| Identifying Information: No | Updated in 2023: No |
|------------------------------------|----------------------------|

Tool Description

Aligns multiple sequencing data.

Tool Purpose & Description of Use

MAFFT (for Multiple Alignment using Fast Fourier Transform) includes several algorithmic methods, including guided tree, scoring matrices, and sequence alignment algorithms to realign multiple genomic sequencing data. The realignment tool is used to locally re-arrange sequence data to make all sequences comparable but the same genomic coordinates. This is used in all sequencing analysis prior to building a phylogenetic tree or distance tree.

Data Analyzed

Training N/A

Data

Input Data Sequences can be in GCG, FASTA, EMBL (Nucleotide only), GenBank, PIR, NBRF, PHYLIP or UniProtKB/Swiss-Prot (Protein only) format

Output Data Fasta or Clustalw

Vendor Involvement

Vendor Name: N/A

5.7 Bowtie2

First Used: June 2022

| | |
|---|---|
| Analysis Type: Predictive modeling, Matching | Population Type: Individuals, Other: sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

Aligns sequencing data to a reference sequence. Bowtie2 aligns sequencing data to a reference using Burrows-Wheeler transformations. It is geared to use with Illumina sequencing data.

Tool Purpose & Description of Use

Bowtie2 is an intermediate step in the workflow to analyze COVID variants in wastewater.

Data Analyzed

Training N/A

Data

Input Data Sequence reads (fastq) for single or paired-end runs (sequence reads can be considered strings).

Output Data Aligned reads in SAM format.

Vendor Involvement

Vendor Name: N/A

Update Description

Updated variable weights/bug fixes/dependencies.

5.8 Vsearch

First Used: June 2022

| | |
|---|---|
| Analysis Type: Predictive modeling | Population Type: Individuals, Other: sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

Vsearch uses the Needleman-Wunsch algorithm to merge read pairs and align and dereplicate sequences to detect chimeric genomic sequences.

Tool Purpose & Description of Use

Vsearch is an intermediate step in the workflow to analyze COVID variants in wastewater.

Data Analyzed

Training N/A

Data

Input Data Sequence reads (fastq, Fasta) for single or paired-end runs (sequence reads can be considered strings).

Output Data FASTA, FASTQ, tables, alignments, SAM.

Vendor Involvement

Vendor Name: N/A

Update Description

Updated variable weights/bug fixes/dependencies.

5.9 IQTREE

First Used: May 2020

| | |
|---|---|
| Analysis Type: Predictive modeling, Matching | Population Type: Individuals, Other: sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

IQTREE uses maximum-likelihood regression to create phylogenetic trees from genomes.

Tool Purpose & Description of Use

Produced phylogenetic trees are used to help rule in or out outbreaks of COVID or other organisms.

Data Analyzed

Training Data N/A

Input Data FASTA, NEXUS, CLUSTALW, PHYLIP.

Output Data Readable report, ML tree in NEWICH format, log file for entire run.

Vendor Involvement

Vendor Name: N/A

Update Description

Updated variable weights/bug fixes/dependencies.

5.10 kSNP3

First Used: March 2022

| | |
|---|---|
| Analysis Type: Predictive modeling | Population Type: Individuals, Other: sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: No |

Tool Description

kSNP3 can use multiple algorithms (maximum-likelihood, parsimony, neighbor-joining) to infer phylogenetic trees from genomes.

Tool Purpose & Description of Use

Produced phylogenetic trees are used to help rule in or out outbreaks of bacteria.

Data Analyzed

Training Data N/A
Input Data Fasta
Output Data ML tree in NEWICH format, log & configuration files, Fasta file

Vendor Involvement

Vendor Name: N/A

5.11 Spades

First Used: October 2017

| | |
|---|---|
| Analysis Type: Predictive modeling, Matching | Population Type: Individuals, Other: sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: No |

Tool Description

Spades uses several algorithms to simplify genomic read data into de Bruijn graphs and finds overlaps to assemble genomes.

Tool Purpose & Description of Use

Spades is an intermediate step in the workflows of bacterial analyses.

Data Analyzed

Training Data N/A
Data
Input Data Fastq.
Output Data Fastas and other files for corrected reads, scaffolds, contigs, paths in GFA format, fastg assembly graph.

Vendor Involvement

Vendor Name: N/A

5.12 GATK

First Used: October 2017

| | |
|---|--|
| Analysis Type: Predictive modeling | Population Type: Individuals and Other; sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

A suite of tools for variant calling and filtering after sequence alignment. It uses naive Bayesian to qualify aligned bases as sequence or erroneous data, which would be excluded from the final genomic sequence.

Tool Purpose & Description of Use

Used to identify mutations and call upon differences from the reference, which is used to generate the predicted complete sequence.

Data Analyzed

- Training Data** Sets of known variant sites.
- Input Data** Fasta, uBam, SAM/BAM/CRAM, VCF.
- Output Data** Bam, txt, vcf.

Vendor Involvement

Vendor Name: N/A

Update Description

Updated variable weights/bug fixes/dependencies.

5.13 PHYLOViZ

First Used: October 2017

| | |
|------------------------------------|--|
| Analysis Type: Optimization | Population Type: Individuals and Other: sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: No |

Tool Description

For representing the possible evolutionary relationships between strains, PHYLOViZ uses the goeBURST algorithm, a refinement of eBURST algorithm by Feil et al., and its expansion to generate a complete minimum spanning tree (MST).

Tool Purpose & Description of Use

Used to generate the minimum spanning tree relationships.

Data Analyzed

Training Data N/A

Input Data txt, NEWICK, FASTA.

Output Data Minimum spanning tree.

Vendor Involvement

Vendor Name: N/A

5.14 Bionumerics

First Used: September 2017

| | |
|---|---|
| Analysis Type: Predictive modeling, Matching | Population Type: Individuals, Other: sequence data can belong to any species |
| Identifying Information: No | Updated in 2023: No |

Tool Description

A suite of tools used to align and analyze bacterial genomes.

Tool Purpose & Description of Use

BioNumerics is used to 1) re-assemble the bacterial genome (since the sequencing process involves fragmenting the bacterial DNA and then amplifying it into millions of pieces) 2) identify the genus, species, and serotype of the bacterial isolate 3) perform quality control checks to ensure the sequence meets certain quality standards 4) perform whole genome multi-locus sequence typing (wgMLST, a technique used to type bacteria based on their genetic code) 5) perform cluster analysis for cases related to one another based upon case definitions recommended by the Centers for Disease Control and Prevention (CDC). This information is then communicated to partners including foodborne epidemiologists at the Bureau of Communicable Disease, who investigate all reported cases of foodborne disease, with those investigations potentially resulting in restaurant inspections, closures, and food recalls.

Data Analyzed

Training Data N/A
Input Data Fastq.
Output Data txt, Excel.

Vendor Involvement

Vendor Name: Bionumerics

Developed and maintains the tool.

5.15 ChoiceMaker (CM)

First Used: June 2003

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling, Optimization, Matching | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: No |

Tool Description

ChoiceMaker (CM) is a record-matching tool that identifies duplicate records belonging to the same individual.

Tool Purpose & Description of Use

CM is used by BOI and Healthy Homes to identify duplicate immunization and lead records. The outputs produced by CM are used in ongoing manual and automated deduplication processes (record merging).

Data Analyzed

Training Data N/A
Input Data CM uses demographic data (e.g., names, date of birth, address, identifiers) and health event data (e.g., date and type of event) from BOI’s Citywide Immunization Registry (CIR) and Healthy Homes’ LeadQuest registry in its evaluation.
Output Data The program outputs a series of record pairs and a match probability for each pair.

Vendor Involvement

Vendor Name: HLN Consulting

A vendor was involved in the development of the program initially. CM is now available as an open-source program. The DOHMH implementation is maintained by HLN Consulting.

5.16 ICE - Immunization Calculation Engine

First Used: 1997

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: Yes |

Tool Description

The Immunization Calculation Engine (ICE) is an immunization evaluation and forecasting system, whose default immunization schedule supports all routine childhood, adolescent, and adult immunizations based on the recommendations of the Advisory Committee on Immunization Practices (ACIP). ICE is free and open-source available through <https://cdsframework.atlassian.net/wiki/spaces/ICE/overview>.

Tool Purpose & Description of Use

ICE is used by the Bureau of Immunization to evaluate a patient’s immunization history and generate appropriate immunization recommendations.

Data Analyzed

Training Data N/A

Data

Input Data ICE uses demographic data (e.g. date of birth) and vaccination data (e.g., immunization date, vaccine group and type) in the evaluation process. Data used are stored in the Citywide Immunization Registry (CIR).

Output Data The program returns recommendations on whether a patient has completed a vaccine series or is due for vaccines.

Vendor Involvement

Vendor Name: HLN Consulting

A vendor was involved in the development of the program and continues to be involved in program enhancements. ICE is also available as an open-source program. The DOHMH implementation is maintained by HLN Consulting.

Update Description

New vaccine groups and recommendations were added.

6.0 Department of Investigations (DOI)

6.1 Facial Recognition Technology

First Used: March 2019

| | |
|---|-------------------------------------|
| Analysis Type: Computer vision, Optimization, Matching | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: Yes |

Tool Description

The tool analyzes an uploaded image or video and searches and compares it with lawfully possessed images to generate a pool of possible matches. If possible matches are identified, a trained DOI examiner visually analyzes and evaluates potential matches to assess reliability of a match consistent with agency policy and applicable laws. A match serves as an investigative lead for additional investigative steps and does not constitute a positive identification.

Tool Purpose & Description of Use

Facial recognition is a digital technology that DOI uses to analyze uploaded images or videos of people and objects obtained during an investigation by comparison with lawfully possessed images. Facial recognition generates possible matches of an object or individual from this analysis and comparison. The purpose of the tool is to assist DOI investigations of matters within its jurisdiction including fraud and other criminal activity.

Data Analyzed

Training Data Self-trained in system usage.

Input Data Images.

Output Data Images.

Vendor Involvement

Vendor Name: Not disclosable

Out-of-the-box products. The vendors provide ongoing technical assistance. Confidentiality agreements are in place with the vendors.

Update Description

DOI acquired additional facial recognition technology in the past year for the above-described purpose. DOI’s facial recognition tools receive product updates and data sets on an ongoing basis.

7.0 Department of Social Services (DSS)

7.1 Homebase Risk Assessment Questionnaire (RAQ)

First Used: June 2012

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: Yes |

Tool Description

Homebase applicants answer screening questions about their current housing situation, history of disruptive experiences, shelter history, and other domains. Each of the answers is assigned a number of points, and applicants that reach a certain point threshold are eligible for deeper Homebase services, such as financial assistance and case management. Workers are able to override a limited number of model decisions with permission of a supervisor.

Tool Purpose & Description of Use

The Homebase program was created to prevent households from entering the DHS shelter system. Since NYC has a range of antipoverty programs and the number of households entering shelter is small compared to the pool of New Yorkers who have an eviction filing each year, the Agency had to ensure that the households who most needed additional homelessness prevention services were being enrolled in Homebase programs. Research showed that staff were not accurately able to predict who would or would not enter the DHS shelter system and that using a risk assessment would provide a better way to match resources to the families who would benefit the most.

Data Analyzed

- Training Data** The RAQ was developed based on analysis of data on Homebase enrollees from 2004 to 2008, conducted in conjunction with a team of academic researchers, to determine predictive factors for those entering shelter. It was updated in 2023 based on analysis led by NYC DSS researchers of 2013-2016 Homebase data.
- Input Data** Factors include, among others: personal characteristics such as age and pregnancy; educational attainment and employment status; housing issues such as eviction, discord, a move in the past year; past and recent experience of homelessness.
- Output Data** The tool produces a score that is used to assess eligibility for full versus brief Homebase services.

Vendor Involvement

Vendor Name: Multiple researchers

DHS contracted with researchers to evaluate years of Homebase administrative data to develop a risk assessment. The DSS research team then led an updated analysis that led to tool revisions. The published research papers are listed below:

<https://ajph.aphapublications.org/doi/10.2105/AJPH.2013.301468>

<https://www.journals.uchicago.edu/doi/abs/10.1086/686466?mobileUi=0&journalCode=ssr>

<https://www.tandfonline.com/doi/abs/10.1080/10511482.2022.2077801>

Update Description

Fresh training data. Started using on 12/21/2023.

7.2 SmartVAN / TargetSmart

First Used: November 2019

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling, Matching | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: Yes |

Tool Description

The Mayor’s Public Engagement Unit (PEU) uses SmartVAN to manage outreach across a range of projects. SmartVAN provides functionality to create lists of potential clients to contact, collect personal information and survey responses from clients, and conduct outreach via phone banks and canvassing. SmartVAN also contains a frequently updated commercial dataset, provided by TargetSmart, of New York City residents and their demographic, contact, and other information. PEU uses this preloaded data to create outreach lists when data on existing clients or from partner agencies is unavailable or insufficient to meet the scope of the outreach project.

Tool Purpose & Description of Use

In 2023, PEU has the TargetSmart data within SmartVAN on a number of projects. PEU frequently uses the data to create lists of residents who live within certain zip codes that PEU wants to target for outreach. For example, PEU created lists using TargetSmart data to conduct door-knocking and phone banking outreach to New Yorkers identified as potentially eligible for the DOF Rent Freeze program based on TargetSmart data. In cases like these, TargetSmart’s determination of who lives in which zip codes as well as estimated income affects whether New Yorkers receive PEU outreach. Additionally, the algorithm that TargetSmart uses to match phone numbers to individuals impacts the type of outreach that New Yorkers receive.

Data Analyzed

Training Data Training data is part of vendor’s proprietary processes.

Data

Input Data Input data is part of vendor’s proprietary processes.

Data

Output Data The algorithmically-derived data that PEU accesses is the output of proprietary algorithmic processes developed and operated by TargetSmart. These algorithmic processes include matching multiple input datasets to determine residency, contact information, and demographics on New York City residents. SmartVAN also includes a number of algorithmically-determined likelihood scores, including scores for the likelihood that a household contains children under 18, etc.

Vendor Involvement

Vendor Name: EveryAction and TargetSmart

EveryAction and TargetSmart jointly provide the SmartVAN product. EveryAction is the software provider. TargetSmart is the data provider. TargetSmart is the entity who applies algorithmic techniques. EveryAction provides access to this data through their platform.

Update Description

TargetSmart data available in SmartVAN is updated on a regular basis by the vendors.

7.3 Splink

First Used: March 2023

| | |
|-------------------------------------|-------------------------------------|
| Analysis Type: Matching | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: N/A |

Tool Description

Splink is Python package used for entity resolution (i.e., deduplication) of records in which there is no unique identifier. It helps users implement probabilistic matching.

Tool Purpose & Description of Use

In 2023, PEU moved the Tenant Helpline onto a live caller system using the Virtual Call Center and Salesforce. Prior to migrating data from the original system (which involved clients leaving voicemails), PEU needed to identify call records from the same client. PEU used the Splink package to identify records from the same caller even in cases with relatively sparse information and/or inconsistent data entry (e.g., different spellings of the same name).

Data Analyzed

- Training Data** PEU fine-tuned our use of the algorithm using a sample of the original Tenant Helpline data.
- Data** There are a number of user-defined parameters that affect the results of the matching. PEU tested which parameters were most appropriate for our use case of very sparse data and conducted human quality control over its performance.
- Input Data** The input data was information on individual Tenant Helpline callers. This included their names, phone numbers, address information, etc.
- Output Data** The output of PEU’s usage of Splink was a set of callers identified as unique linked to one or more calls to the Tenant Helpline. This information was transformed and migrated to the new Salesforce system.

Vendor Involvement

Vendor Name: N/A

8.0 Fire Department (FDNY)

8.1 RBIS (Risk Based Inspection Program); ALARM (A Learning Approach to Risk Modeling)

First Used: November 2019

| | |
|---|--|
| Analysis Type: Predictive modeling | Population Type: Property, Other: civilian fatalities |
| Identifying Information: No | Updated in 2023: No |

Tool Description

A Learning Approach to Risk Modeling (ALARM) creates risk scores for each building in the city. These scores are used to schedule our Fire Operations building inspections within the inspectable population of buildings in the city (~330,000 Building Identification Numbers (BINs)), as a part of the Risk-Based Inspection Program (RBIS).

Tool Purpose & Description of Use

ALARM is a combined approach using machine learning and risk ratios to assess the risk of a building for structural fire ignition (probability) and civilian fire injury/death (impact). The machine learning algorithm takes incident data, housing characteristics, and 311 data and creates a probability of structural fire ignition. This is combined with a civilian injury or death risk ratio for the building which is based on building characteristics, incident data and nearby felony crimes to create a risk score (range is one to nine), with one being highest risk and nine being lowest risk. Buildings are prioritized within each of the nine risk scores according to the residential population in each building.

Data Analyzed

- Training Data** In order to create the models, the team utilized a five-year incident dataset and reserved 99 percent of the data to train the probability model and 80 percent of the data to train the impact model.
- Input Data** The ALARM risk score utilizes data from our fire and Emergency Medical Services (EMS) dispatch system, building characteristic data, 311 calls, felony crimes, census data and civilian injury data.
- Output Data** The tool outputs a risk score from one (highest risk) to nine (lowest risk).

Vendor Involvement

Vendor Name: N/A

ALARM was built in-house by a team of analysts from the Management, Analysis and Planning Bureau.

8.2 EMS Hospital Suggestion Algorithm

First Used: March 2007

| | |
|---|--|
| Analysis Type: Predictive modeling | Population Type: Geographic space |
| Identifying Information: No | Updated in 2023: No |

Tool Description

The EMS Hospital Suggestion Algorithm is used to determine the closest, appropriate hospital to the incident location based on the needs of a patient requiring transport.

Tool Purpose & Description of Use

The algorithm computes a list of hospitals in order of closest to furthest in time for each medical condition category as currently established. (For example, there is a list of hospitals computed in order of closest in time for all hospitals that accept General Emergency Department patients and for all hospitals that accept special conditions, such as burns). Depending on the medical needs category of the patient, the algorithm produces a pre-determined list of hospitals which is based on the location of the patient and then made available to the crew as a list of “closest, most appropriate hospitals.”

Data Analyzed

Training Data The EMS Hospital Suggestion algorithm relies on telematics data from the Department of Citywide Administrative Services city-owned vehicles collected between 2015 and mid-2016 to calibrate a network analysis model that derives incident to hospital transport times. The order of suggested hospitals are then compared with five years of historical EMS hospital transport data from before the COVID-19 pandemic (2015-2019) to validate and correct the network model.

Input Data The inputs for the algorithm include the location and medical condition of the patient.

Output Data The algorithm outputs the closest, most appropriate hospitals.

Vendor Involvement

Vendor Name: N/A

This algorithm and the resulting output file that is used in our EMS CAD system to suggest hospitals was provided by Deccan International, until September 2020. The Department currently creates this file using a new algorithm, developed in-house by the Geographic Information Science (GIS) unit in conjunction with engineers from Columbia University.

8.3 EMS Unit Suggestion Algorithm

First Used: March 2007

| | |
|---|--|
| Analysis Type: Predictive modeling | Population Type: Geographic space |
| Identifying Information: No | Updated in 2023: No |

Tool Description

The Emergency Medical Services (EMS) Unit Suggestion Algorithm is used to determine which order of geographic regions (known as atoms) to search in order for the EMSCAD system to select an appropriate EMS unit for dispatch to an incident.

Tool Purpose & Description of Use

The algorithm computes a list of geographic regions (known as atoms) in order of closest to furthest in travel time for each atom in the city. This list of ordered atoms is the output of an algorithm that relies on a calibrated network model to derive travel time estimates. The output is an excel file which is converted into an EMSCAD-compatible file and loaded into the system for real-time unit selection capabilities. The file is generated and implemented as a 24/7 source file, meaning, the recommended search order is not currently varying by time of day. The Department is intending to implement time-of day search orders in the near future.

Data Analyzed

- Training Data** The EMS Unit Suggestion algorithm relies on historical FDNY CAD trip time data which is used to calibrate a network analysis model which derives atom-to-atom transport times.
- Input Data** The input for the algorithm is a geographic location.
- Output Data** The algorithm outputs a recommended EMS unit for dispatch.

Vendor Involvement

Vendor Name: Deccan International

This algorithm and the resulting output file that is used in our EMS CAD system to suggest atom order for unit search is currently provided by a vendor, Deccan International.

8.4 EMD Schedule Optimization Tool

First Used: June 2021

| | |
|---|---|
| Analysis Type: Predictive modeling | Population Type: Other: FDNY radio and assignment dispatcher employees |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

The purpose of the tool is to provide Emergency Medical Dispatchers (EMD) staff a tool to optimally allocate call takers during a 24-hour period. The tool uses an expected number of incoming calls and the number of personnel scheduled to work in order to allocate the call takers to different shifts such that the supply of call takers exceeds the demand for call takers.

Tool Purpose & Description of Use

The algorithm requires two datasets. First, the tool requires the average number of medical calls per hour for a 24-hour period. Second, the tool requires a user to specify the number of call takers assigned to each tour. Based on these two inputs, the tool provides a projection of supply (call takers) versus demand (medical calls). Additionally, the tool can take the total number of available staff and optimally allocate them across tours to maximize the minimum difference between supply and demand. Based on these outputs, EMD officers can identify times during the day when call taker utilization is high and reallocate staff to accommodate.

Data Analyzed

- Training Data** This is an optimization model and was not “trained” using training data. The algorithm relies on actual historical data to determine average hourly medical calls.
- Input Data** The tool requires an hourly count of medical calls arriving during a 24-hour period. Additional “data” requirements are input from the user depending on user-driven scenarios. For example, a user could specify five eight-hour tours per day (at different start times) rather than existing four tours (two eight-hour tours and two 12-hour tours).
- Output Data** The algorithm outputs a projection of supply (call takers) versus demand (medical calls). Additionally, the tool can take the total number of available staff and optimally allocate them across tours to maximize the minimum difference between supply and demand.

Vendor Involvement

Vendor Name: N/A

The tool was developed internally at FDNY in partnership with Columbia University’s Industrial Engineering and Operations Research Department.

Update Description

Front-end user interface added.

9.0 Health + Hospitals (H+H)

9.1 Adult Risk of IP/ED Utilization Score

First Used: January 2019

| | |
|---|-------------------------------------|
| Analysis Type: Predictive Modeling | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: Yes |

Tool Description

The Adult Risk of Inpatient (IP)/Emergency Department (ED) Utilization Score predicts the number of days in the ED or IP setting that a patient may have in the coming year.

Tool Purpose & Description of Use

The Adult Risk of IP/ED Utilization Score predicts the number of days in the ED or IP setting that a patient may have in the coming year. It uses internal electronic medical record data covering past utilization, diagnoses, and documented behavioral health risk factors. Patients in the top 5 percent may be eligible to participate in the health home program, even if they are not Medicaid eligible.

Data Analyzed

Training Data This model was trained on NYC Health + Hospital population to create a utilization prediction method that was responsive to the need of safety-net patients. Most national algorithms are tuned on large claims dataset and may not be generalizable to uninsured patients. The model was developed with 70+ predictor variable and tuned using a LASSO continuous model. We assessed 70+ predictor variables on electronic health record data from CY Q3 2016-Q2 2017 to train the model (n=833,969). Data from CY Q3 2017-Q2 2018 was used to validate. This study was approved by NYC Health + Hospital Institutional Review Board (IRB) partner, BRANY IRB.

Input Data The final model had 17 predictors. The top binary predictors for the final model were psychosis diagnosis ($\beta=1.17$), history of incarceration ($\beta=0.47$), antipsychotic medication prescription ($\beta=0.40$), and substance use disorder diagnosis ($\beta=0.38$). Top continuous predictors were inpatient visits ($\beta=0.36$), ED visits ($\beta=0.34$), and number of chronic conditions ($\beta=0.21$).

Output Data Outputs of the model (top 5 percent flag) allow clinicians to connect patients with social work support or referrals to community organizations.

Vendor Involvement

Vendor Name: N/A

Update Description

Updated our homelessness variable to incorporate additional sources of documentation within our system, and validated the algorithm against fiscal year 2021 data to evaluate continued fit.

10.0 Mayor’s Office (MO)

10.1 Scorecard Blockface Sampling Algorithm - MO - Operations

First Used: March 2022

| | |
|------------------------------------|--|
| Analysis Type: Optimization | Population Type: Geographic space |
| Identifying Information: No | Updated in 2023: No |

Tool Description

The Scorecard program sends inspectors across New York City to rate street and sidewalk cleanliness. The sampling algorithm creates a monthly list of blocks for inspectors to visit and rate.

Tool Purpose & Description of Use

The primary goal of the algorithm is to produce a sample of blockfaces that is statistically sound and geographically representative. This list is used to rate street and sidewalk cleanliness citywide, as well as by borough and DSNY district.

Data Analyzed

Training Data N/A

Input Data

The blockface sample is selected from the Pavement Edge File, which is part of the New York City Planimetric Database managed by the Office of Technology and Innovation. Sampling is weighted towards blockfaces in high-density areas and includes extra sampling of blockfaces in Business Improvement Districts (BIDs). It also takes into account the linear miles of street within a DSNY District.

Output Data A count of blockfaces that are statistically representative of our target areas throughout the city.

Vendor Involvement

Vendor Name: Legacy Mayor’s Office of the Chief Technology Officer

The sampling algorithm was developed by the former Mayor’s Office of the Chief Technology Officer in partnership with the Mayor’s Office of Operations.

10.2 Methodology for Poll Site Language Assistance - MO - CEC

First Used: November 2020

| | |
|---|-------------------------------------|
| Analysis Type: Predictive modeling | Population Type: Individuals |
| Identifying Information: No | Updated in 2023: No |

Tool Description

Since no dataset is currently available that reliably captures the number of limited English proficient (LEP) registered voters for all program languages, the Civic Engagement Commission (CEC) uses the percentage of LEP citizens of voting age (CVALEP) as a substitute or proxy measure of need. CEC ranks the program-eligible languages in order of magnitude of CVALEP and distributes poll sites to each language based on its ranking (excluding CVALEP persons that speak languages served by the NYC Board of Elections (NYCBOE) in certain New York City counties). The number of poll sites that will receive services in any given language will depend on each language’s share of the total CVALEP in the population eligible to be served. For example, according to U.S. Census data, approximately 207,926 New Yorkers are CVALEP and speak a language that is served by this program. This proportionality approach allows CEC to balance goals of including diverse language communities as well as fair access to the total number of eligible voters within each language community. The program provides interpreters in program-eligible languages at poll sites based on U.S. Census data showing concentrations of CVALEP individuals who speak these languages and reside around each poll site. For each language, poll sites are chosen in descending order of concentration of CVALEP, until the language’s share is met. This process is repeated for each language, thereby including the poll sites with the highest concentration of CVALEP for each program-eligible language until that language’s share is met, and the total number of poll sites for which resources are allocated is reached. It may be possible, based on analysis of data, to reassign poll sites to languages with greater need; however, each language will receive a minimum of at least one poll site. Models used included the Thiessen polygon method to create a Voronoi diagram to determine CVALEP estimates.

Tool Purpose & Description of Use

This is a methodology for determining how the New York City Civic Engagement Commission (NYCCEC) will provide interpretation services at poll sites for limited English proficient voters. The methodology explains how the NYCCEC will identify the languages and locations in which interpretation services will be offered during the November 2020 election and beyond. These services supplement the interpretation assistance provided by NYC Board of Elections in several languages. Under the Charter, the NYCCEC can only provide interpretation services in a language if: (1) it is a designated citywide language; or (2) it is spoken by a greater number of LEP New Yorkers than the lowest ranked designated citywide language and at least one poll site has a significant concentration of speakers of such language with LEP. This methodology ensures service for all languages that are eligible under the Charter.

Data Analyzed

Training Data N/A

Input Data For citywide estimates, this methodology uses current data from the American Community Survey (ACS) 2016-2020 five-year estimates. This methodology also uses the American Community Survey Census Tract 2016-2020 five-year Public Use Microdata Samples for poll site level analysis; this is the most current and accurate data available on resident New Yorkers at the neighborhood level. In addition, the methodology uses data from the Board of Elections on the location of election districts and poll sites.

Output Data The algorithm estimates the number of citizens of voting age with Limited English Proficiency for each program-eligible language who could report to each polling site.

Vendor Involvement

Vendor Name: N/A

The tool was designed with the support of the Mayor’s Office of Data Analytics which is currently part of the Office of Technology and Innovation.

10.3 ElevenLabs Speech Synthesis - MO - COS

First Used: June 2023

| | |
|--|-------------------------------------|
| Analysis Type: Speech and language processing | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: N/A |

Tool Description

ElevenLabs creates the most realistic, versatile and contextually-aware AI audio, providing the ability to generate speech in hundreds of new and existing voices in over 20 languages.

Tool Purpose & Description of Use

The tool was used to generate audio recordings of mayor’s voice delivering hiring hall/RiseUp concert messages in various languages to be used for hiring hall/RiseUp concert robo-calls.

Data Analyzed

Training Data Submitted audio recordings of mayor’s speech, as well as a live language sample, into the program.
Input Data Script translated in desired languages (Spanish, Yiddish, Haitian Creole).
Output Data Audio recording (MP3) of mayor’s voice speaking the script in the desired language (Spanish, Yiddish, Haitian Creole).

Vendor Involvement

Vendor Name: ElevenLabs

Utilized two existing features offered by ElevenLabs (VoiceLab and Speech Synthesis).

11.0 Police Department (NYPD)

11.1 Facial Recognition Technology

First Used: October 2011

| | |
|---|-------------------------------------|
| Analysis Type: Computer vision, Matching | Population Type: Individuals |
| Identifying Information: Yes | Updated in 2023: No |

Tool Description

Tool may help investigators identify unknown subjects in law enforcement investigations.

Tool Purpose & Description of Use

Facial recognition is a digital technology that NYPD uses to compare images obtained during investigations with lawfully possessed arrest and parole photos. The tool analyzes an uploaded image, known as a probe image, and searches and compares against the image repository. The purpose of the tool is to enhance law enforcement’s ability to investigate criminal activity as well as identify deceased persons and missing persons. When used in combination with human analysis and additional investigation, facial recognition technology is a valuable tool in solving crimes and increasing public safety.

Data Analyzed

Training Data Training data is proprietary to the vendor.

Data

Input Data If NYPD investigators obtain a still image depicting a face of an unknown individual during an investigation, the image can be submitted for facial recognition analysis in accordance with NYPD facial recognition policy. Known as a probe image, NYPD facial recognition software compares the image to a controlled and limited group of lawfully obtained photos called the photo repository.

Output Data The facial recognition software will generate a pool of possible match candidates for review by trained Facial Identification Section investigators.

Vendor Involvement

Vendor Name: Dataworks

Software developed and maintained by vendor.

11.2 ShotSpotter

First Used: March 2015

| | |
|------------------------------------|--|
| Analysis Type: Matching | Population Type: Geographic space |
| Identifying Information: No | Updated in 2023: Yes |

Tool Description

Provides acoustic gunshot detection to assist with emergency call response.

Tool Purpose & Description of Use

Provides acoustic gunshot detection to assist with emergency call response. The tool supports patrol operations in alerting units to potential gunfire and enhances investigations involving firearms.

Data Analyzed

Training Training data is proprietary to the vendor.

Data

Input Specialized software analyzes audio signals for potential gunshots.

Data

Output The tool determines the location of the sound source, and once classified as potential gunfire sends the incident to acoustic experts for additional analysis. Notifications are sent for confirmed gunfire. ShotSpotter activations may result in evidence collection that can enhance case investigations. Problematic locations identified through alerts may require additional resource deployment and/or investigations.

Vendor Involvement

Vendor Name: ShotSpotter

Software developed and maintained by vendor.

Update Description

Routine maintenance.

11.3 Patternizr

First Used: December 2016

| | |
|------------------------------------|--|
| Analysis Type: Matching | Population Type: Individuals, Property, Geographic space, Other: Crime Classification |
| Identifying Information: No | Updated in 2023: No |

Tool Description

Aids crime analysis in detection of potential crime patterns.

Tool Purpose & Description of Use

Patternizr compares features of crimes and finds ones that are similar and may be part of a crime pattern. Analysts will look at the candidate crimes and suggest the formation of crime patterns to a

pattern identification module. If a pattern is formed, detectives often consolidate the investigative efforts (e.g., one detective investigates all the crimes in the pattern.) The report filters non-normal trends into a spreadsheet and displays year-over-year counts of crimes that have non-normal trends. The tool requires a human user to evaluate the output data to see if complaints identified as similar are, in fact, connected to a pattern.

Data Analyzed

Training Data Separate models were trained for each of three different crime types (burglaries, robberies, and grand larcenies). These crime types have a sufficient corpus of prior manually identified patterns for use as training examples. This corpus consists of approximately 10,000 patterns between 2006 and 2015 from each crime type. A portion of this corpus includes complaint records where the same individual was arrested for multiple crimes of the same type within a span of two days.

Input Data The input data is a candidate crime and its features. A complaint describes details of the crime, including the date and time (which can be a range if the precise time of occurrence is unknown), location, crime subcategory, modus operandi, and suspect information. This information is used to calculate the five types of crime-to-crime similarities used as features by Patternizr: location, date-time, categorical, suspect and unstructured text.

Output Data Probability that a complaint is connected to a pattern.

Vendor Involvement

Vendor Name: N/A

The tool was developed by data scientists and analysts at NYPD. Contractors and NYPD personnel integrated it into the Domain Awareness System. Personnel in Crime Control Strategies work with the Information Technology Bureau to maintain the tool.

12.0 Office of the Chief Medical Examiner (OCME)

12.1 STRMix

First Used: January 2017

| | |
|--|--|
| <p>Analysis Type: Other: STRmix is a forensic DNA analysis software program that uses a probabilistic genotyping algorithm to interpret complex DNA profiles, such as those from mixed samples that contain DNA from multiple contributors. Specifically, STRmix uses a continuous probabilistic modeling approach called Markov chain Monte Carlo (MCMC) analysis.</p> | <p>Population Type: Individuals</p> |
|--|--|

| | |
|-------------------------------------|----------------------------|
| Identifying Information: Yes | Updated in 2023: No |
|-------------------------------------|----------------------------|

Tool Description

STRmix™ combines sophisticated biological modeling and standard mathematical processes to interpret a wide range of complex DNA profiles. Using well-established statistical methods, the software builds millions of conceptual DNA profiles. It grades them against the evidential sample, finding the combinations that best explain the profile. A range of likelihood ratio options are provided for subsequent comparisons to reference profiles. Using a Markov Chain Monte Carlo engine, STRmix™ models any types of allelic and stutter peak heights as well as drop-in and drop-out behavior. It does this rapidly, accessing evidential information previously out of reach with traditional methods. STRmix™ is supported by comprehensive empirical studies with its mathematics readily accessible to DNA analysts, so results are easily explained in court.

Tool Purpose & Description of Use

STRMix is a probabilistic genotyping tool that is used to analyze mixtures of DNA profiles to help associate the crime scene evidence to potential victims or suspects of crimes.

Data Analyzed

- Training Data** Training data was not used in the sense of AI software. The OCME performed thousands of tests using the software to validate it for optimum use with our current laboratory standard operating procedures and genetic analyzers.
- Input Data** Forensic DNA profiles from crime scenes as well as the DNA profiles from victims and suspects of crimes.
- Output Data** The output is a deconvolution of genotype probability distribution that lists all of the accepted genotype sets and their associated weights. These weights can take any value from 0 to 1.

Vendor Involvement

Vendor Name: NicheVision Forensics, LLC

The software has been developed by New Zealand Crown Institute of Environmental Science and Research (ESR) with Forensic Science South Australia. The developer assisted OCME in analyzing and interpreting our data during the validation of the software.

13.0 Office of Technology & Innovation (OTI)

13.1 MyCity Chatbot

First Used: September 2023

| | |
|--|---|
| Analysis Type: Speech and language processing | Population Type: Individuals, Group, organization, or business |
|--|---|

| | |
|------------------------------------|-----------------------------|
| Identifying Information: No | Updated in 2023: N/A |
|------------------------------------|-----------------------------|

Tool Description

The NYC MyCity chatbot is a beta AI-powered chatbot that provides information and access to services for residents and businesses in New York City.

Tool Purpose & Description of Use

The NYC MyCity chatbot is a beta AI-powered chatbot that provides information and access to services for residents and businesses in New York City. It’s currently focused on two main areas: Business Services and MyCity Basics. The chatbot provides information on starting or operating a business in New York City, answers questions about permits, licenses, regulations, and other business requirements, and connects users with relevant resources and support services. It also offers information on various city services and benefits, and helps users find resources related to childcare, career, and other areas. The chatbot is using Microsoft’s Azure AI technology and OpenAI’s ChatGPT 3.5 Turbo LLM.

Data Analyzed

Training Training data is proprietary to the vendor.

Data

Input Data Text queries are input by the user on the MyCity portal.

Output Data The tool produces text responses with references based on information from Business Services and MyCity Basics.

Vendor Involvement

Vendor Name: Microsoft, Nuvalence

Microsoft provides Cloud-based ChatGPT services and Nuvalence was the professional services vendor for implementation.

14.0 School Construction Authority (SCA)

14.1 GitHub Copilot

First Used: May 2023

| | |
|--|-------------------------------------|
| Analysis Type: Speech and language processing | Population Type: Individuals |
| Identifying Information: No | Updated in 2023: No |

Tool Description

GitHub Copilot is an AI-powered code assistant that provides suggestions for whole lines or blocks of code in a wide range of programming languages. It leverages a vast codebase and machine learning to

improve coding efficiency, helping programmers by autocompleting code snippets and offering context-appropriate code suggestions.

Tool Purpose & Description of Use

GitHub Copilot is primarily used by some of our software developers as an advanced coding assistant within our agency. Its role is to augment and streamline the coding process for our software development projects. By providing real-time code suggestions and completions, it reduces the time developers spend on routine coding tasks, allowing them to focus on more complex aspects of software development.

The tool functions by analyzing the context of the code being written and suggesting relevant, syntactically correct code snippets. This includes generating code for standard programming patterns, filling in boilerplate code, and offering solutions to simple programming queries. It's important to note that while GitHub Copilot assists in the coding process, final decisions on the code's implementation and its use in any software or application rest solely with our human developers. The tool's suggestions are always reviewed and potentially modified by our team to ensure they meet our specific requirements and standards. Therefore, GitHub Copilot acts as a support tool in the decision-making process of software development rather than a decisive entity.

Data Analyzed

| | |
|----------------------|---|
| Training Data | GitHub Copilot was developed by OpenAI and trained using a large corpus of public source code available on GitHub. This training data includes a wide variety of code in multiple programming languages, along with associated comments and documentation. The data encompasses a broad range of coding styles, patterns, and solutions across different software development projects. |
| Input Data | When in use, GitHub Copilot analyzes the code that a developer is currently writing. This input data consists of the programming language syntax, structure, and any comments or context within the code file. The tool also takes into account the specific coding task, patterns, and functions that the developer is working on. This real-time data is essential for the tool to provide relevant and context-appropriate coding suggestions. |
| Output Data | The output data from GitHub Copilot includes suggested lines or blocks of code that align with the input data provided by the developer. These suggestions are generated based on the patterns, structures, and coding practices learned from the training data. The output is designed to seamlessly integrate with the existing code, offering syntactically correct and contextually relevant code completions. |

Vendor Involvement

Vendor Name: GitHub, a subsidiary of Microsoft

Appendix

Reportable Elements

1. **Tool name.** The name or commercial name and a brief description of the algorithmic tool.
2. **Tool description.** A description of how the information received from the algorithmic tool is used, including the purpose for which the agency is using the algorithmic tool.
3. **Training data.** The type of data collected or analyzed by the algorithmic tool and the source of the data.
4. **Vendor name and involvement.** Whether a vendor or contractor was involved in the development or ongoing use of the algorithmic tool, a description of such involvement, and the name of such vendor or contractor when feasible.
5. **Date of first use.** The month and year in which the algorithmic tool began to be used, if known.
6. **Identifying Information.** Whether the tool collects or analyzes "identifying information" as defined under New York City's Identifying Information Law (Local Laws 245 and 247 of 2017).
7. **Analysis type [all that apply].** Type(s) of analysis used by the tool:
 - a. Predictive modelling
 - b. Speech and language processing
 - c. Computer vision
 - d. Optimization
 - e. Matching
 - f. Other (please describe)
8. **Population type [all that apply].** Population(s) are reviewed, assessed, or directly affected by the tool:
 - a. Individuals
 - b. Property
 - c. Group, organization, or business
 - d. Geographic space
 - e. Other (please describe)
9.
 - a. **Updated tool.** If reported previously, whether the tool has been updated in the past year.
 - b. **Update description.** If yes, how the tool was updated (e.g., fresh training data, variable weights).