

NIDDK-CR Resources for Research

Data Science and Meet the Expert Webinar Series



March 27, 2025



NIDDK Central Repository Overview

Central Repository

Mission

Established in 2003 to **facilitate sharing of data**, **biospecimens**, **and other resources** generated from studies supported by NIDDK and within NIDDK's mission by making these **resources available for request to the broader scientific and research community**.

- Supports receipt and distribution of data and biospecimens in a manner that is ethical, equitable, and efficient
- Enables investigators not involved with the original work to test new hypotheses without the need to collect new data or biospecimens
- Promotes FAIR (Findable, Accessible, Interoperable, and Reusable) and TRUST (Transparency, Responsibility, User focus, Sustainability, and Technology) principles



Recorded past tutorials, webinars, and other educational resources can be found on the NIDDK-CR website







NIDDK Data Sharing Ecosystem

The NIDDK-CR is a part of the broader NIH-funded biomedical data ecosystem and plays a key role in NIH's FAIRness and TRUSTworthiness goals. The NIDDK-CR houses a broad range of data types for secondary research, provides access to biospecimens, and direct links to other repositories with additional resources such as genomics data.





Future Functionality: Analytics Workbench

Central Repository

Streamlining end-to-end data science lifecycle and discovery of data-driven biomedical insights.

Innovation and ease of use

A cloud-based analytics environment where researchers and data scientists can access a suite of integrated analytics tools and cloud computing resources to participate in data challenges and Al innovation.

Expected Benefits of Analytics Workbench:







NIDDK-CR Data Science Centric Challenge Series

Goals of NIDDK-CR Data-science centric challenge series

- Develop tools, approaches, models and/or methods to increase data interoperability and usability for artificial intelligence (AI) and machine learning (ML) applications
- Augment and enhance existing data for future secondary research, including data-driven discovery by AI/ML researchers
- Discover innovative approaches to enhance the utility of datasets for AI/ML applications



Visit our website for more information on our data-centric movement and to learn more about our past data-challenges



Secondary Data Science and Meet the Expert Webinar Series

Central Repository

About the Series

- Aims to accelerate data science and Al-driven biomedical research by fostering collaboration between biomedical researchers and experts in the field
- Monthly webinar held on the last Thursday of each month

Upcoming Webinars

- Today Artificial Intelligence fundamentals
- April 24 Artificial Intelligence fundamentals applications
- May 29 FAIR data sharing
- June 26 Different privacy preserving techniques and implications for researchers
- July 31 Challenges, opportunities, and considerations for secondary researchers using electronic health records and real-world data sources
- August 28 Impact and innovations realized



Learn more about the webinar series, register for future webinars, and access past webinars materials and recordings



Central Repository

Meet the Experts



Gordon Aiello is a Lead Scientist at Booz Allen Hamilton with a PhD in Applied Mathematical and Computational Sciences. He works full-time developing and delivering specialized data science, artificial intelligence, machine learning, and Python trainings for clients in the Navy and Intelligence Community. Prior to joining Booz Allen Hamilton, Dr. Aiello worked in the Office of Macroeconomic Affairs at the U.S. Department of State, using machine learning techniques to analyze developing and emerging market economies. Additionally, he has taught courses on data science and the R programming language for the Foundation for Advanced Education in the Sciences (FAES) at the NIH. He is passionate about working with others to expand their understanding of data science techniques and their applications.



Al Fundamentals

NIDDK-CR Data Science Meet the Experts Webinar Series March 27th, 2025



Presented by: Booz Allen Hamilton



Central Repository

Training Guidance

- Avoid CUI/PII/PHI conversations
- Questions in Teams Chat are encouraged
- Due to size of class, stay on mute until end of class



Data Science Learning Program

If you're new to data science, start your learning journey with the **Foundations** courses. A more in-depth learning track starts with the **Data Science Fundamentals** course and continues to the **Data Science Labs**. Those interested in more specialized topics can explore courses in the **Select Topics** track.



Introduction to Data Visualization

- Telling a story with your data
- How to create more impactful briefings

ûdemy

• Not product specific

NAVWAR



() 3 hours

SELECT TOPICS

Python Fundamentals for Data Science

- Foundational Python syntax
- Develop essential analytic skills
- Machine Learning and Artificial Intelligence

Live **JUPITER** () 7 hours (2 sessions)

Artificial Intelligence Fundamentals

- Al initiatives and foundational Al
- Al ecosystems and Al operations
- Responsible and Ethical AI
- Neural Networks

Live Training () 7 hours (2 sessions)

Data Science for Managers

Developed in partnership with NGA

- Management responsibilities in Data Science Projects
- Ethical considerations in Data Science
- Data Science and AI Opportunities

In Person Training



Diabetes and Digestive and Kidney Diseases





- 1. Al Demystified
- 2. Responsible Al
- 3. Al Operations
- 4. Al Implementations





AI Demystified



Central Repository

Data Science Defined

The goal of data science is to extract meaningful insights from data.

Data – any kind of qualitative or quantitative set of values

- Common examples in data science today:
 - Natural text: "I'm cold," "I'm not very cold"
 - Categories: "yellow," "green," "red"
 - o Numbers: 1, 2.53, -4

Images:



• Sometimes you have the data, sometimes you need to procure the data

Science – a systematic approach to building knowledge by testing hypotheses

- Think Scientific Method:
 - Define a hypothesis \rightarrow Collect the data \rightarrow Analyze results \rightarrow Draw conclusions
- Hypotheses must be testable, and experiments must be reproducible



Data Science Is an Interdisciplinary Field

Central Repository

Data science lies at the intersection of mathematics and statistics, computer science and programming, and domain knowledge and expertise.



Source: Booz Allen Hamilton



Data Science Is an Interdisciplinary Field

Data science lies at the intersection of mathematics and statistics, computer science and programming, and domain knowledge and expertise.



Source: Booz Allen Hamilton



Data-Information-Knowledge-Wisdom Model





How DS Intersects with AI & ML

The various subfields of data science intersect with one another:

Central Repository

Artificial Intelligence:

• The theory and development of computer systems able to perform tasks normally requiring human intelligence

Machine Learning:

• A field of computer science that uses statistical techniques to give computer systems the ability to "learn" (e.g., progressively improve performance on a specific task) with data, without being explicitly programmed

Deep Learning:

 Is part of a broad family of machine learning methods based on learning data representations, as opposed to task-specific algorithms (predominately based on artificial neural networks)







Why Now?



Central Repository

Why Now? Big Data



- Businesses, governments, and societies are only starting to tap data's vast potential
- Calls for individuals to "strengthen data and statistics for accountability and decisionmaking purposes"



Why Now? Data Is Universal

Central Repository

- Data has replaced oil as the world's most valuable resource, according to the Economist
- In 2000, 5% of the world was connected to the internet
- In 2024, 67.5% of the world is connected to the internet, causing massive increases in the volume of data being produced daily
- In 2024, approximately 400 quintillion* bytes of data are created every day

*400 quintillion = 400,000,000,000,000,000





Why Now? Computer Capabilities

Central Repository

Computer capabilities have greatly improved over the years, allowing us to perform computations that would take longer than the age of the universe on previous models.









Brief History of AI & Its Rapid Ascension





Brief History of AI & Its Rapid Ascension







Limitations of Data Science & Al



Central Repository

Data Science Myths

Myth 1: It's as simple as pushing a button



Myth 3: It's too soon to talk



Myth 2: Data science will make analysts obsolete



Myth 4: Everybody needs to learn coding





Al Is Statistics Over Larger Datasets





What Is AI?

Al is a field concerned with producing machines able to autonomously perform tasks that would normally require human intelligence by giving them the ability to perceive, learn from, abstract (classify, conceptualize, and generate rules), and act using data.



What is AI?

- Interchangeable with the term machine intelligence or "MI"
- Able to perform certain narrowly defined tasks as well or better than humans
- A substitute for human intelligence in certain tasks within jobs

What isn't AI?

- Interchangeable with the term "data science"
- Able to perform wide-ranging tasks as well as or better than humans
- A substitute for any human's entire job



Tasks Well Suited to AI

- Simple, requires little to no context
- Involves finding patterns in data
- Characterized by large amounts of data preferably labeled
- Occurs in a statis environment or one with little uncertainty

Tasks Poorly Suited to Al

- Complex, requires contextual understanding
- Requires explaining patterns in data
- Little or no data to characterize the problem
- Occurs in a dynamic environment with a lot of uncertainty



How Is AI Used in Medical Care and Biomedical Research?





• Al Uses in Healthcare

- **Radiology**: The ability of AI to interpret imaging results may aid in detecting a minute change in an image that a clinician might accidentally miss.
- **Early Detection of Kidney Disease**: Al can help in early detection by analyzing patterns in patient data (e.g., blood tests, medical history, and imaging scans). For instance, Al systems are used to predict the onset of chronic kidney disease (CKD) by analyzing blood markers like creatinine levels or urine albumin-to-creatinine ratio.
- Diabetes Prediction and Risk Assessment: AI can analyze patient data, including family history, lifestyle factors, and clinical records, to predict the risk of developing Type 2 diabetes. By identifying high-risk individuals, early interventions like lifestyle changes or medications can be recommended.
- **Continuous Glucose Monitoring (CGM) and Al-Driven Insulin Dosing**: For individuals with Type 1 diabetes, continuous glucose monitors (CGM) provide real-time glucose level data. Al models can predict glucose fluctuations and automatically adjust insulin doses to maintain optimal glucose levels.





How Machines Learn



How Machines Learn

- We can get an intuition for how machines learn by experimenting with an interactive applet demonstrated on the upcoming slides.
- A helpful analogy for thinking of this process is to relate machine learning with baking, as illustrated in the correspondence table below.

Machine Learning	Baking	
Data – Human Provided	Ingredients to be mixed together (e.g., flour, sugar, butter, cream, etc.)	
Desired Model Output – Human Provided	Tasty treat (e.g., cake, cookie, biscuit, etc.)	
Model Parameters – Computer Determines These	Quantities of ingredients used in recipe (e.g., 3 cups sugar, 4 tbsp butter, etc.)	



Fitting a Linear Model





Fitting a Quadratic Model





Fitting a Cubic Model





Deep Learning

- Neural Networks: A machine learning approach modeled on the human brain in which algorithms process signals via interconnected layers of artificial neurons.
 - Mimicking biological nervous systems, artificial neural networks have been used to recognize and predict patterns of neural signals involved in brain function, especially in differentiating typographic characters and in recognizing facial features.
- **Deep Learning**: A form of machine learning that uses many layers of computation to form a more complex neural network — often called a deep neural network that is capable of learning from large amounts of complex, unstructured data.
 - Deep neural networks enable such devices as voicecontrolled virtual assistants; self-driving vehicles, which learn to recognize traffic signs; and computer-vision systems used by doctors to more quickly and accurately assess X-rays and other medical images.

Deep Neural Network







Responsible Al



Al at the NIH

• The NIH makes a wealth of biomedical data and Al opportunities available to research communities and aims to make these data findable, accessible, interoperable, and reusable—or FAIR.





NIH AI Policy Guidance

Central Repository

Protection of Human Subjects (45 CFR 46): Outlines basic provisions for Institutional Review Boards, informed consent, and assurance of compliance for NIH-supported research involving human participants and their data, including considerations of risks & benefits.
 Establishes the requirement to submit a DMS Plan and comply with NIH-approved plans. In addition, NIH Institutes, Centers, and Offices can request additional or specific information be included within the plan to support programmatic priorities.

	Research Participant Protections	Data Management and Sharing	NOT-OD-23-149: Informs the
Health Insurance Portability and Accountability Act (HIPAA) helps protect the privacy and security of health data used in research, including research involving AI, thereby fostering trust in healthcare research activities.	Health Information Privacy Biosec Bio	Peer Review urity and safety	extramural community that the NIH prohibits NIH scientific peer reviewers from using natural language processors, large language models, or other generative AI technologies for analyzing and formulating peer review critiques for grant applications and R&D contract
Research funded by NIH, includi	ng research using the tools and tech	nologies enabled or informed by AI, fall	under this oversight framework.



NIH AI Policy Resources

- <u>Use of Generative AI in Peer Review FAQs (NIH Office of Extramural Research)</u>
- <u>NIH Office of Data Science Strategy</u>
- <u>US Department of Health and Human Services Artificial Intelligence Use Cases Inventory</u>
- <u>Executive Order on the Safe, Secure and Trustworthy Development and Use of Artificial Intelligence</u>
- <u>PCAST Report to the President Supercharging Research: Harnessing Artificial Intelligence to Meet Global</u> <u>Challenges</u>
- <u>NIH STRIDES Initiative | NIH STRIDES</u>



RAI Lifecycle: Planning

Central Repository



The planning phase is the process to map out, identify the objective, and define the functionality of the AI system to solve a given problem.



RAI Lifecycle: Planning

Phase I: Planning

Questions to answer throughout the Planning phase are:

- Have you evaluated other methods of implementing a similar objective?
- What do the success metrics look like?
- Do you have the right and appropriate data for the problem?
- What are the resources required to create the solution?
- Is ownership of the AI solution clearly defined?
- Who are the stakeholders and are they informed of the process and solutions?
- Are the risks identified and the process laid out for a system malfunction?



RAI Lifecycle: Development

Central Repository



In the development phase your team will iterate on the process to build the planned-out AI system through data modeling, output verification, and data governance.



RAI Lifecycle: Development

Central Repository

Phase II: Development

Questions to answer throughout the Development phase are:

- What are the impacts of data or model biases and how can you mitigate them?
- What are the metrics needed for post deployment monitoring?
- Who is the owner for maintaining the AI system?
- Are the AI system and outputs explainable?
- What is the process for routine maintenance?



RAI Lifecycle: Deployment

Central Repository



In the Deployment phase, the AI system is put to use in an operational setting to solve the problem initially posed.



RAI Lifecycle: Deployment

Central Repository

Phase III: Deployment

Questions to answer throughout the Deployment phase are:

- Is there a process in place to continuously monitor and validate the outputs from the AI system?
- Does the capability still meet the needs of the desired solution and tasking?
- Are the potential biases and negative impacts of the outputs considered and mitigated?





Al Operations



AI Ops





AlOps Processes – Data Ops

Central Repository



Data Ops – a collaborative data management practice focused on improving the communication, integration, and automation of data flows between data managers and data consumers across an organization

Data readiness can be measured through several key drivers:

•Accuracy: The degree to which data correctly describes the "real world" object or event being described.

•**Usability**: A measure of data's reliability. (Data that is messy, is missing substantial records, is difficult to ingest, or is otherwise prohibitive would be less usable.)

•**Provenance:** A record that describes entities and processes involved in producing and delivering or otherwise influencing that resource. Provenance provides a critical foundation for assessing authenticity, enabling trust, and allowing reproducibility.

•Security: How safe the data is from loss or unwanted manipulation.



Data Governance Priorities

- Creating policies and procedures tailored to support data needs for AI at enterprise scale.
- Addressing data acquisition policies and procedures to include legal processes, agreements, and data brokers/data services.
- Empowering key stakeholders to facilitate and lead your organization's data governance.
- Developing data security plans and data privacy policies in accordance with appropriate regulatory frameworks (e.g., Health Insurance Portability and Accountability Act [HIPAA], Federal Information Security Modernization Act [FISMA], European Union's General Data Protection Regulation [GDPR], etc.).
- Prioritizing metadata management to ensure compliance and troubleshooting support.





Machine Learning Operations

Central Repository



Machine Learning Ops – focused primarily on the governance and life cycle management of a wide range of operationalized artificial intelligence (AI) and decision models

MLOps governs the development and sustainment of an AI solution's analytic backbone throughout its operational lifecycle. Not all ML solutions have the same requirements in terms of:

- Required data
- Development techniques
- Deployment considerations.



Machine Learning Lifecycle





DevSecOps



DevOps – a set of processes and techniques devoted to increasing the efficiency of communication and coordination between developers and operations staff,
 DevSecOps expands this relationship further to include security and the need to make systems robust against intentional meddling by a motivated adversary





AlOps Integration





Al Adoption Blueprint





Responsible AI: addresses risk evaluation and management concerns to promote responsible workforce adoption of AI solutions
Analytic frameworks: leveraging advanced modeling and MLOps frameworks to implement analytics.
Digital architecture: outlines reference architectures and infrastructure to enable component integration and DevSecOps.
Data engineering: realizes enterprise data strategy through implementation of DataOps, utilizing data as a strategic asset.





AI Solutions



Guidelines for Approaching Al Projects





Design

- The design phase represents the beginning of the AI lifecycle. The inputs to this phase are typically either a business or mission **need** or an innovative **idea**.
- Critical Activities:
 - Requirements Analysis decomposition of end-user needs
 - Reference Architecture infrastructure elements that the AI solution requires
- Useful Questions to ask:
 - What is the problem the solution must address?
 - What relevant data is available to the team?
 - Is there necessary data that the team lacks access to?
 - Do any relevant pre-trained models exist?
 - How should available data be explored, understood, and prepared?







- The Develop phase picks up where design ends and seeks to turn a detailed idea into a functioning prototype.
- In practice, these steps are not sequential but rather cyclical and iterative.
- Use Agile development where cross-functional teams focus on continuous improvement, in which planning, design, and development all run in parallel
- Focus on a Minimum Viable Product
 - An MVP is a product with **just enough features** to attract early-adopter users and **validate a product idea** early in the development cycle.
 - Focus on what matters most.
 - Get small wins early in the development process.
 - Generate user feedback early and often.





Deploy

- The Deploy phase is all about putting your AI models to use. Often, this means deploying models to users who
 will employ them in their day-to-day work streams.
- AIOps for Deployment
 - Activities needed to manage the AI Lifecycle
 - Seeks to take the repetitive, low-value work of deploying and maintaining models off the plate of AI experts, allowing them to focus on innovation tasks that cannot be automated
- If you don't have a solid operations process in place on your AI project, your team will have the following challenges:
 - Reproducing experimental results will require manual record-keeping and data archiving.
 - Updating an AI model to a new version will be time consuming, and users might experience downtime.
 - All of your critical processes will be prone to human error.
 - Retraining models will be cumbersome and time-consuming.
 - Your team will spend a lot of time on repetitive, low-value tasks, detracting from innovation opportunities.





People-Process-Technology (PPT) Framework

Golden Triangle





People

- **People:** Workforce knowledge, skills, and abilities to perform required tasks across the AI solution lifecycle.
 - Technical
 - Non-technical
- Leadership and project/resource management: Who is driving the overarching strategic AI objective, product, or service forward?
- **Domain and business expertise**: Who is analyzing investments and risks? Who is familiar with the data?
- **Data engineering and stewardship**: Who is defining and maintaining processes for reliable and clean data?
- **Data science and applied mathematics**: Who is analyzing the data in detail to identify the optimal models/algorithms?
- **Machine learning engineering**: Who is researching, selecting, optimizing models and the software/hardware configurations to run them?
- Software development: Who is building and managing the software infrastructure?
- **Security and deployment engineering**: Who is ensuring secure implementation of the product or service?
- Visualization and reporting expertise: Who is facilitating efficient content digestion for project stakeholders?





Processes

- Al processes seek to orchestrate the interactions between the workforce (people) and the platform or tools (technology) required to engage Al solutions.
- From an operational perspective, you can consider AI Operations in terms of its technical component processes:
 - DataOps
 - MLOps
 - DevSecOps
- You might also leverage AI lifecycle frameworks to inform processes around high-level phases: design, develop, and deploy.





Technology

- Technical components, in contrast to people and processes, are the fairly tangible components of AI-enabling infrastructure. However, these technologies are extremely diverse in their intended use, capabilities, and application to your unique AI challenge.
- **Platforms**: A digital environment (hardware or software) that houses centralized technology including development tools, applications, services, and user interfaces.
 - EX: Amazon Web Services (AWS), Microsoft Azure, Google Cloud Platform (GCP), Advana/Jupiter, etc.
- **Development Tools**: A tool that supports the developer to create, test, and share repositories to enable collaborative AI solutions.
 - EX: Integrated Development Environments (IDEs), GitHub, Databricks, Jupyter Notebooks, Google Colab, etc.
- Infrastructure as Code: Software that manages and provisions data centers through machine-readable definition files, rather than physical hardware configuration.
 - EX: Terraform, Ansible, Puppet, Chef, Otter, etc.
- **Physical Infrastructure**: Hardware used to house and provide secure access to platforms.
 - EX: servers, storage media, personal devices, edge devices, sensors, communications technologies, etc





Central Repository

Technology Demo







- Al Demystified
- Responsible AI at the NIH
- Al Operations
- Al Implementations



Questions?





Arica Christensen – <u>Christensen_arica@bah.com</u> Dr. Gordon Aiello – <u>Aiello_Gordon@bah.com</u>

April 24 – Al Fundamentals Applications



Thank You!