

Click to verify





























[illegible]



We covered free datasets sources and discussed common ways to download datasets from them. Through this article, we learned how to download and use those datasets in Python and Pandas. We covered different Python libraries which offer public datasets for learning. Finally, we covered how to create test datasets with fake data. Those datasets and ideas should be sufficient for practicing the machine learning and data science. By using DataScienceTool, you agree to our Cookie Policy. Wherewith you're a student, exploring new concepts or a seasoned professional, you'll find an untapped wealth of information with which we all benefit from the importance of quality data. We've all heard the saying, "garbage in, garbage out," and it's a reminder that our projects are only as good as the data we feed them. Seeking robust and relevant datasets is key to the success of any endeavor. And let's face it, there's no better way to excel and make an impact than by rolling up your sleeves and diving into projects, headfirst. This is especially true when it comes to data-centric work, where each new dataset presents a unique challenge and an opportunity to sharpen your skills. I'm here to let you in on a secret: there's a treasure trove of fascinating and accessible datasets just waiting for your discovery. No more struggling to find that perfect data source or feeling limited by data constraints. If you're eager to take your projects to the next level and dive into some exciting resources, you've come to the right place. In this article, we will guide you through a carefully curated list of websites and sources that offer an abundance of datasets, perfect for any idea you want to bring to life. Get ready to bookmark these go-to websites and embark on a journey of endless project possibilities, where you can focus on honing your craft, gaining valuable experience, and maybe even creating something that will revolutionize the industry. So, without further ado, let's unlock the door to a world of data-driven exploration! Table of contents Kaggle is a prime platform for accessing datasets due to its vast repository covering diverse topics like astronomy, diabetes, and more. With user-friendly features for filtering by license type and topic, Kaggle ensures easy access to high-quality datasets, making it ideal for both beginners and experienced data scientists seeking valuable resources for their projects. Find Kaggle Datasets here. To master Kaggle, read this article. AWS Data Exchange is a robust platform for data exchange, offering a wide range of datasets from various providers, including government agencies and companies. It is particularly useful for projects as it provides a one-stop shop for diverse data needs. With data from multiple sources, it saves time and offers consistency. The platform's reliability and the ability to discover unique datasets make it a valuable resource for data scientists. Click here to explore AWS Data Exchange. Data.world is a fantastic resource for datasets due to its community-driven nature. It offers a vast collection of user-uploaded datasets, ensuring variety and specialty. The platform fosters collaboration and data exchange, making it ideal for finding unique and specific datasets. The website's search functionality and categorization make it easy to navigate and discover relevant data quickly. Click here to explore Data.world! GitHubHub, while primarily known as a code repository, has evolved into a valuable resource for datasets. Its vast community of users often shares datasets alongside their code, providing a unique perspective. The website is ideal for projects as it offers a one-stop shop for code and data, with a simple interface. GitHub's search functionality and filtering options make finding relevant datasets efficient and straightforward. Click here to checkout this website to find datasets. OpenDataSoft is a reliable data-sharing platform, offering a comprehensive directory of open datasets. Its strength lies in its focus on data sharing and collaboration. The website is ideal for projects as it provides a one-stop source for diverse data needs, with a user-friendly interface. The platform's commitment to open data and its global reach make it a valuable tool for anyone looking for transparent and accessible datasets. Click here to checkout Open Data Soft. DataHub is an excellent platform for accessing free and open datasets on various topics. Its strength lies in its comprehensive collection of data from different sources, making it a one-stop shop for projects. The website is user-friendly, well-organized, and efficient for finding relevant data quickly. DataHub's commitment to open data and its partnership with the data analytics company Qlik ensure a reliable and valuable resource for any data-driven project. Click here to checkout this website to find datasets. The Google Public Data Explorer is a unique tool that provides access to a vast array of public data from international organizations and academic institutions. While it is not a direct dataset source, the platform offers a user-friendly interface to explore and visualize data. This makes it ideal for gaining insights and understanding trends. The tool's strength lies in its ability to make complex data accessible and its suitability for projects requiring dynamic data representation. Click here to explore Google Public Data Explorer! Data.gov is an extensive open data platform provided by the US government. It centralizes data from a wide range of federal, state, and local agencies, covering diverse topics. The website is well-structured, allowing users to easily search and filter data by format, topic, and agency. For projects, Data.gov offers a reliable and authoritative source of data, ensuring consistency and authenticity. Its regular updates and commitment to open data make it an invaluable resource for anyone seeking US-specific datasets. Click here to explore this website to find datasets. Data.gov.uk is the UK government's comprehensive open data portal. It provides access to a wide range of data from UK public bodies and agencies, ensuring transparency and accessibility. The website is well-organized and user-friendly, making it easy to navigate and find relevant data. For projects, it offers a reliable source of UK-specific data, covering various topics such as economics, health, and education. The platform's commitment to open data and its regular updates make it a valuable resource for data-driven projects. Click here to explore data.gov. The Census Data of India website offers a rich collection of demographic, economic, and social data about the country. It provides valuable insights into India's population, including various statistics and indicators. The website is essential for projects focusing on India, offering detailed information at the national and regional levels. Click here to explore datasets at Census Data of India. Open Government Data Platform India: This is a centralized platform for open data from the Indian government, covering various sectors. India Data Portal: It provides open data on various themes, including agriculture, education, and energy. National Data Portal: This portal offers a comprehensive collection of datasets, publications, and maps from the Indian government. Ministry of Statistics and Program Implementation: The official website of the ministry provides data on economic indicators, demographics, and more. Open Data Initiative India: This initiative aims to provide open data on governance, infrastructure, and other key areas. The World Bank Open Data platform is a comprehensive source of global development data. It offers extensive datasets covering a wide range of topics, including economic indicators, social development, and more. The platform's extensive coverage of economic and social development data make it a go-to resource for researchers, policymakers, and data analysts alike. Click here to explore this website to find more datasets. UN Data, maintained by the United Nations, is a rich repository of global data covering a wide range of topics. The platform offers data on areas such as population, environment, trade, and human development, providing valuable insights into global trends and issues. UN Data is well-structured and user-friendly, allowing users to search and filter data by country, region, and theme. For projects with a global perspective, this website is essential, offering reliable and authoritative data. The UN's commitment to data transparency and its comprehensive coverage of socio-economic and environmental data make it a trusted source for researchers, policymakers, and anyone working on international development initiatives. Click here to explore the datasets at UN Data. Eurostat is the European Union's statistical office, providing comprehensive data on EU member states. It offers a wide range of data covering economic, social, and agricultural topics, among others. The website is user-friendly, allowing easy navigation and data exploration by country, indicator, and theme. Eurostat is particularly valuable for projects focused on Europe, offering reliable and up-to-date information. The office's commitment to data transparency and its extensive coverage of EU-specific data make it an essential resource for understanding the economic and social dynamics of the region. Eurostat also provides data visualization tools and analytical reports, further enhancing its usefulness for researchers and policymakers alike. Click here to explore the websites to find datasets. FRED Economic Data, hosted by the Federal Reserve Bank of St. Louis, is an extensive database of US and international economic data. It offers thousands of economic time series, covering various indicators such as inflation, employment, interest rates, and more. The platform is user-friendly, providing powerful search and filtering tools to navigate the vast dataset. FRED Economic Data is ideal for projects requiring economic analysis, offering a one-stop shop for historical and current economic indicators. The Federal Reserve Bank's commitment to data transparency and its regular updates make this platform a trusted source for researchers, economists, and anyone interested in economic trends and forecasting. Click here to checkout the FRED Economic Data here. The UCI Machine Learning Repository is a well-known and trusted source for machine learning and artificial intelligence research and education. Maintained by the University of California, Irvine, it offers a diverse collection of datasets specifically curated for ML and AI applications. The repository is user-friendly, providing a comprehensive dataset search and filtering system. It is valuable for projects as it covers various data types, from text and images to time series and sensor data. The platform's commitment to supporting ML research and education, along with its regular updates, makes it an indispensable resource for students, researchers, and practitioners in the field of machine learning and artificial intelligence. Click here to checkout datasets in UCI ML Repository. OpenML is a collaborative platform designed specifically for the machine learning community. It offers a unique approach to data sharing by providing not just datasets but also machine-learning tasks and flows. This platform allows researchers and practitioners to share and reproduce experiments easily. OpenML is well-organized, with a user-friendly interface, making it simple to search and explore datasets, tasks, and flows. The platform fosters reproducibility and transparency in ML research, making it ideal for projects requiring a more comprehensive approach to data and experimentation. OpenML's commitment to openness and its active community make it a valuable resource for advancing machine-learning practices. Click here to explore this website to find datasets. CMU StatLib is a renowned statistical database provided by Carnegie Mellon University. It offers a rich collection of datasets specifically curated for statistical and machine learning research and education. The database is well-organized, providing a comprehensive search and browsing experience. CMU StatLib is valuable for projects requiring statistical analysis, offering a diverse range of data types and indicators. The platform's association with a leading university ensures the reliability and quality of the datasets. CMU StatLib's regular updates and commitment to supporting statistical research and education make it an indispensable resource for students, researchers, and practitioners in statistics and machine learning. Click here to explore this website to find datasets. Google Dataset Search is a specialized search engine by Google, launched in 2018, designed to help researchers find freely available online data. It allows filtering by data type and is based on schema.org metadata standards. The service complements Google Scholar and offers a user-friendly interface accessible on mobile devices. Click here to checkout Google Datasets Search! Open Data Monitor is a unique website that aggregates open data portals from around the world. It serves as a discovery platform, making it easier for users to find datasets from different countries and regions. The website is well-designed, providing a centralized search engine for exploring global open data initiatives. Open Data Monitor is valuable for projects requiring international data, offering a diverse and comprehensive collection of sources. The platform promotes transparency and accessibility, ensuring that users can quickly locate relevant data portals and gain insights into global open data practices. Its continuous updates and expansion make it a dynamic resource for anyone working with international data. Click here to explore this website to find datasets. DataPortals.org is a comprehensive global registry specifically designed to help users discover open data portals from cities, regions, and countries around the world. It serves as a centralized platform, providing easy access to a wide range of datasets offered by governments and organizations. DataPortals.org is valuable for projects and research that require diverse and localized data. The website promotes transparency and open data practices, ensuring users can quickly locate and utilize datasets that align with their specific needs. With regular updates and a growing community, DataPortals.org has established itself as a dynamic and trusted resource for anyone seeking open data sources, offering a unique perspective on global data initiatives. Click here to explore datasets at DataPortals.org. Data Is Plural is a unique initiative that curates interesting and diverse datasets from various sources on the web. It takes the form of both a newsletter and an archive, providing a regular stream of data-related content. Data Is Plural offers a broad range of topics, covering areas that are often overlooked by traditional data platforms. This makes it ideal for projects requiring unique and specialized data. The newsletter format provides a convenient way to discover new datasets, while the archive ensures a growing collection of valuable resources. Data Is Plural's commitment to exploring the "plural" nature of data and its focus on lesser-known datasets make it a dynamic and intriguing resource for data enthusiasts, researchers, and anyone seeking fresh perspectives in their projects. Click here to explore this website to find datasets. Nasdaq is a renowned global electronic marketplace for buying and selling securities, particularly known for its focus on technology stocks. While primarily a stock exchange, Nasdaq also offers a wealth of data and analytics tools on its website. This includes real-time market data, company profiles, financial news, and investment analysis. For projects involving stock market analysis or financial research, Nasdaq is a valuable resource, providing authoritative data and insights. The platform offers various data products, APIs, and solutions tailored to different user needs. Nasdaq's reputation, combined with its commitment to innovation and data transparency, makes it a trusted source of financial information for investors, traders, and researchers worldwide. Click here to access Nasdaq datasets. Yelp is a well-known crowd-sourced review platform that periodically releases large datasets of its business and review data. The Yelp Dataset is valuable for academic research and data science competitions, offering insights into consumer behavior and preferences. It provides rich information, including business details, user reviews, and ratings, allowing for a wide range of analytical projects. The dataset is unique due to its scale and real-world applicability. Yelp's commitment to data transparency and its impact on local businesses make this dataset a valuable resource for researchers and data scientists, offering a window into consumer trends and behavior patterns. Checkout datasets at yelp here. The Pew Research Center is a non-profit think tank that conducts surveys and research on a wide range of topics, including social issues, media usage, and political attitudes. The user is known for its commitment to providing unbiased and reliable data to the public. Its website offers easy access to a wealth of datasets, making it a valuable resource for projects requiring public opinion and demographic information. The Pew Research Center's data covers a diverse range of subjects, such as technology adoption, social trends, and global attitudes. The platform provides user-friendly data exploration tools and detailed methodology explanations, ensuring transparency and understanding. Researchers, journalists, and anyone interested in societal insights will find the center's data and analyses invaluable, offering a window into the beliefs and behaviors of diverse populations. Click here to checkout this website to find datasets. NASA Open Data is a fascinating portal that provides access to a wide range of scientific data from NASA's various missions and research. It offers a unique opportunity to explore space science, Earth science, and aerospace research data. The platform is user-friendly, allowing easy discovery and download of datasets, images, and even software. NASA Open Data is ideal for projects requiring scientific and space-related information, providing authoritative and detailed insights. The platform's commitment to data transparency and its continuous updates ensure that researchers, students, and enthusiasts can access the latest findings and contribute to further exploration. With data from NASA's renowned missions, this portal offers a window into the universe, inspiring innovation and discovery. Click here to explore NASA Open Data. Figshare is a trusted repository designed for hosting and sharing scientific research outputs, including datasets, code, and other research artifacts. It provides a platform for researchers to make their work openly accessible and citable. Figshare is valuable for projects requiring scientific data, offering a wide range of disciplines, such as life sciences, social sciences, and computer sciences. The platform ensures proper credit and attribution to researchers, promoting open science practices. Figshare's user-friendly interface allows easy search and download of datasets, fostering collaboration and reproducibility. With a commitment to long-term data preservation and an expanding community, it has become an indispensable resource for researchers, institutions, and anyone seeking open scientific data and resources. Click here to explore this website to find datasets. BuzzFeed News Data is a unique initiative by the BuzzFeed data journalism team, where they release the datasets used in their investigative journalism and articles. The platform offers a range of datasets covering topics like politics, social issues, and media. BuzzFeed News Data provides valuable insights into the data-driven stories that shape our world. The datasets are often accompanied by explanatory articles, providing context and understanding. This initiative promotes data transparency and accountability, allowing researchers and the public to explore and analyze the information themselves. BuzzFeed News Data is ideal for projects requiring real-world, contemporary datasets with a focus on current affairs. It bridges the gap between data and storytelling, offering a dynamic resource for data journalists and researchers alike. Checkout the following links to find the datasets: Reddit's "/r/datasets" community, or "subreddit," is a vibrant and unique collection of datasets shared and discussed by its members. It offers a diverse range of datasets covering various topics, from science and technology to social sciences and hobby projects. The community-driven nature of "/r/datasets" provides a dynamic and engaging space for data enthusiasts to collaborate and explore. The subreddit is valuable for finding specialized and niche datasets that may not be easily accessible elsewhere. It fosters a culture of data sharing and discussion, with members offering insights, feedback, and suggestions. For projects requiring unique or specific data, "/r/datasets" is a valuable resource, providing a combination of crowd-sourced data and expert advice. It bridges the gap between data enthusiasts and experts, creating a collaborative environment for data exploration and discovery. Click here to explore datasets. I hope that this list of resources would prove extremely useful for people looking out for doing pet projects or side projects. For the starters, this is definitely a gold mine. Make sure you pick a few side projects and continue to work on them. If you can think of any application of these datasets or know of any popular resources which I have missed, please feel free to share them with me via the comments below. Looking forward to your feedback! To understand other domains, it is important to wear a thinking cap and... Gautam Vermani Data Consultant at Confidential Not sure what you are looking for? View All Projects This section compiles a selection of data sets that will spark your imagination for data visualization and help you create engaging stories with data. These datasets will inspire you to produce data visualizations that effectively communicate the necessary insights- whether your goal is to visualize global trends, analyze social media trends, or dig into the complexities of financial markets. Source- Wikipedia FiveThirtyEight (now known as 538) is a data journalism website that publishes stories on various topics, including sports, politics, and science. They also maintain a repository of open-source data sets on GitHub that they use for their analyses and are available for others to use. These data sets provide valuable insights across several domains, making them ideal resources for data visualization projects. Data science experts often use these interesting data sets to visualize data, such as their election data sets covering historical voting trends, candidate performance, etc. Their sports datasets include player statistics, game outcomes, and team performances and can help create engaging visualizations showcasing trends within the sports world. Source- FiveThirtyEight Datasets Below are some valuable datasets offered by FiveThirtyEight- Source- HDX The Humanitarian Data Exchange (HDX) platform provides access to open-source datasets related to humanitarian crises and development. It offers a vast collection of datasets covering various topics, including natural disasters, food security, health, and displacement. The platform hosts 20,340 datasets from 254 locations combined from 1947 sources. These data sets can be used for several data visualization projects, such as creating maps to show the spatial distribution of humanitarian crises, and creating interactive visualizations to inform humanitarian decision-making. Source- Humanitarian Data Exchange Datasets Sample Datasets by The HDX Below are some valuable datasets offered by the Humanitarian Data Exchange- Source- WHO The World Health Organization (WHO) is a specialized United Nations agency providing global public health leadership. They maintain a vast repository of healthcare data and data on several health-related topics, including data on diseases, health systems, and health outcomes. These data sets can be used for various data visualization projects, such as creating maps to show the geographical distribution of health indicators, visualizing trends in health indicators over time, and creating interactive visualizations to inform public health policies and interventions. Source- World Health Organization Datasets Sample Datasets by WHO Below are some valuable datasets offered by WHO- Global Health Estimates (Offer the latest available data on causes of mortality and disability worldwide by WHO region and country, age, sex, and income group). Immunization Data and Statistics (Provides data on immunization coverage and trends for various diseases worldwide) Source- NASA The National Aeronautics and Space Administration, or NASA, offers data on space exploration and Earth science. NASA's Open Data Portal provides access to vast data sets covering various topics, including astronomy, planetary science, the Earth observations, and aerospce. You can use these NASA data sets for any typical data visualization project, such as creating maps to show the geographical distribution of data, visualizing trends in climate data over time, and comparing different regions or entities based on their data. Source- NASA Databases Sample Datasets by NASA Below are some valuable datasets offered by NASA- Global Surface Water Datasets (ProjectPro's Global coverage of surface water extent from 1993 to present, derived from satellite observations) Air Quality Data (Provides air quality data for various locations worldwide, including measurements of PM2.5, ozone, and nitrogen dioxide). Get confident to build end-to-end projects Access to a curated library of 250+ end-to-end industry projects with solution code, videos and tech support. Request a demo Did you know that the top five happiest countries in 2023 were Finland, Denmark, Iceland, Israel, and the Netherlands? Source- World Happiness Report 2023 The World Happiness Report is a landmark survey of the global state of happiness. It is published annually by the Sustainable Development Solutions Network (SDSN) and is based on a study of over 150 countries. The report primarily uses the Gallup World Poll data and ranks countries based on their citizens' reported happiness levels, using factors such as life expectancy, social support, and income. You can use the World Happiness Report data for various data visualization projects, such as creating maps to show the geographical distribution of happiness scores, visualizing trends in happiness scores over time, and comparing different countries or regions based on their happiness scores. Source- World Happiness Report Data Best Datasets For Data Processing Projects Data processing is the backbone of data science, transforming raw data into a clean and organized state ready for analysis. This section showcases a collection of carefully selected data sets that will test your data processing skills. From handling missing values and dealing with outliers to normalizing data and feature engineering, these datasets provide ample opportunities to hone your data-wrangling skills. Source- Amazon Amazon Web Services (AWS) provides a comprehensive collection of publicly accessible datasets through its Open Data Registry. The datasets are available in CSV, JSON, and parquet formats, making them suitable for various data processing tasks. AWS Open Data Registry project datasets are ideal for various data processing projects, including data cover classification, and also support scientific research in fields like climate change and renewable energy development. Source- AWS Open Data Registry Datasets Sample Datasets by AWS Below are some valuable datasets offered by AWS- Open Surface Water Datasets (ProjectPro's Global coverage of surface water extent from 1993 to present, derived from satellite observations) Sample Data by Stack Exchange Datasets Here are a few sample datasets you can use from Stack Exchange datasets- Stack Overflow Posts (Provides over 25 million questions, 35 million answers, 90 million comments, and 25K tags from Stack Overflow) Ask Ubuntu Posts (Provides over 412K questions, 520K answers, 1.5 million comments, and 3.2K tags from Ask Ubuntu) Data Processing Projects For Practice Here are some projects you can add to your data science portfolio that use data processing techniques- This section will prepare you to tidy up your data like a pro with our roundup of popular data sets for data-cleaning projects. Let us explore a few widely used datasets specifically chosen to help you practice data-cleaning techniques, such as exploratory data analysis (EDA), ensuring your datasets are error-free and ready for analysis. Source- Academic Torrents Academic Torrents is a platform that provides a curated collection of over 127.15TB of research data from scientific papers across various disciplines. These exciting data sets are available through the BitTorrent protocol, enabling efficient and distributed access to large datasets. The platform is a valuable resource for researchers and data scientists working on various data-cleaning projects. Leveraging these datasets to build a data cleaning project can involve cleaning datasets by performing exploratory data analysis to identify incomplete or inaccurate entries, standardizing data formats to ensure consistency for further analysis, and enhancing the dataset by adding relevant information from external sources. Source- Academic Torrents Sample Data by Academic Torrents Datasets Here are a few sample datasets you can use from Academic Torrents datasets- YASP 3.5 Million Data Dump (Provides an extensive collection of over 3.5 million interesting data-driven articles) OpenAcl Snapshot (Offers metadata for 209M works (journal articles, books, etc.), 2013M disambiguated authors, 124K venues (journals and online repositories), 109k institutions, and 65k Wikidata concepts) Source- Reddit With approximately 1.1 billion users worldwide, Reddit is a popular community-driven online forum and discussion platform. It offers a separate section called the datasets subreddit, or /r/datasets, a valuable resource for data cleaning projects due to its vast collection of real-world data, diverse topics, and active community involvement. The subreddit provides various datasets from various domains and sources, suitable for multiple data-cleaning needs. The datasets are typically open-access and freely available, making them accessible to many users. This further encourages discussion and collaboration among users, facilitating the exchange of data-cleaning techniques and best practices. Source- Reddit Datasets Sample Data by Reddit Datasets Here are a few sample datasets you can use from Reddit datasets- Source- Data.world Data.world (or the 'social network for data people') is a platform that provides access to a vast repository of open data (nearly 129K datasets) from various sources, including multiple US government agencies, academic institutions, and non-profit organizations. The platform offers a user-friendly interface for searching, exploring, and downloading datasets, and it also provides tools for data analysis and cleaning. The datasets are curated by experts and undergo rigorous quality checks, ensuring the reliability of the data. The datasets are freely accessible and maintained by a community of researchers, and include detailed metadata and documentation, providing context and facilitating data understanding. Source- Data.world Datasets Sample Data by Data.world Datasets You can acquire data sets from Data.world platform such as Climate Change Data (Offers data from World Development Indicators and Climate Change Knowledge Portal on climate systems, greenhouse gas emissions, and energy use) UN Human Development Index (Provides data on the United Nations Human Development Index (HDI) for various countries over time for 2022) Another useful dataset is the Billion Word Language Modeling Benchmark Dataset, which can be found on Google Research and is a valuable resource for data cleaning projects in natural language processing (NLP). It provides a massive collection of text data from various sources, including books and news articles. Various use cases of the Billion Word Language Modeling Benchmark Dataset for data cleaning projects include developing data cleaning algorithms on a large corpus of real-world text data, and assessing the impact of data cleaning on downstream NLP tasks, such as sentiment analysis and text summarization. Source- Billion Word Language Modeling Benchmark Dataset Data Cleaning Projects For Practice Here are some projects you can add to your data science portfolio that use data-cleaning methods, such as exploratory data analysis and others- Interesting Datasets For Machine Learning Projects In machine learning, datasets are the secret tools that power algorithms, enabling them to discover hidden patterns and make predictions. This section curates a collection of intriguing datasets that will challenge your machine-learning skills and inspire you to implement ML in real-world scenarios. Whether you want to predict customer churn, analyze medical trends, or optimize marketing campaigns, these datasets will provide the perfect platform to hone your machine-learning expertise. Source- Nasdaq Data Link Trusted by over 800K professionals, Nasdaq Data Link is a platform that provides access to over 250 datasets covering economic and financial data, such as historical and real-time stock prices, economic data, and company financials. The platform offers data sets in various formats, including CSV, XML, and JSON, and provides a user-friendly API for accessing data programmatically. You can use the Nasdaq Data Link free datasets for various machine-learning tasks. For instance, you use these datasets to analyze stock price data and indices, predict economic indicators to identify trends and patterns, and develop trading algorithms based on technical analysis and machine learning. Source- Nasdaq Data Link Datasets Sample Data by Nasdaq Data Link Datasets You can find free datasets on the Nasdaq Data Link platform, such as US Federal Reserve Data Releases (Provides official US figures on money supply, government finances, bank assets and debt, exchange rates, etc.) World Bank Data (Provides data from hundreds of countries and regions worldwide, from multiple categories such as finance, economy, climate change, government expenditures, etc.) With a global community of 262K quants, researchers, data scientists, and engineers, QuantConnect provides a platform for developing and deploying trading algorithms and access to a comprehensive collection of financial datasets. These datasets cover various asset classes, including stocks, futures, options, and foreign exchange, ideal for various ML applications in the financial domain. Data scientists can employ QuantConnect datasets for various ML tasks in finance to train ML models to identify profitable trading signals and predict market movements, evaluate trading algorithms using ML techniques, and optimize portfolio allocation using ML models Source- QuantConnect Datasets Sample Data by QuantConnect Datasets You can find free datasets on the Nasdaq Data Link platform, such as The Statlog Shuttle Landing Discovery Dataset is a challenging and valuable dataset for ML projects offered by the UCI Machine Learning Repository, which offers a collection of shuttle flight-related measurements, including altitude, velocity, acceleration, and sensor readings. Data enthusiasts can leverage this dataset to identify anomalous patterns in sensor data to predict potential shuttle landing failures, build predictive models to predict equipment failures and evaluate various ML algorithms for classifying shuttle flights based on their sensor readings. Source- Statlog Shuttle Landing Discovery Dataset Sample Data by Statlog Shuttle Landing Discovery Dataset You can pick sample data from the Statlog Shuttle Landing Discovery dataset, such as Shuttle Test Set (A subset of the data that can be used to evaluate the performance of trained ML models) Shuttle Landing Failure Reason Codes (A table of codes that indicate the specific reasons for shuttle landing failures) The House Prices: Advanced Regression Techniques dataset, hosted on Kaggle, offers detailed information about houses sold in Ames, Iowa, along with their sale prices. This dataset aims to train ML models to accurately predict a house's sale price based on various features, such as its size, location, amenities, and other characteristics. You can use this dataset to develop ML models to predict the houses' sale price of houses based on their features, analyze factors that impact real estate prices in different regions, and build recommendation systems that recommend suitable houses to potential buyers as per their preferences and budget. Source- House Prices Advanced Regression Techniques Dataset Machine Learning Projects For Practice Here are some fascinating projects for you that use ML methods- Bonus Public Datasets For Data Science Projects As data science projects grow in complexity, the need for real-time data sources has become increasingly important. While traditional data science projects rely on pre-existing datasets, an increasing amount of valuable data online services generate in real-time, such as social media feeds, financial market data, and sensor readings. Here are a few bonus streaming data sources that can be used for your next streaming data project- Source- GitHubHub GitHub is a valuable resource for real-time data science projects, offering a curated collection of open-source datasets available in streaming or real-time formats. These datasets provide a continuous data stream, enabling researchers and developers to analyze and respond to real-time events. You can use streaming datasets from GitHub to analyze social media streams to gain insights into public opinion on trending topics or any organization, develop real-time models to predict stock prices based on market data or news sentiment, and several other innovative data science projects. Source- GitHubHub Picture this: a constant flow of live tweets at your fingertips. That's what Twitter's real-time streaming API offers! Twitter offers access to a vast real-time data stream through its Twitter Streaming API. It provides a valuable resource for data science projects that require the latest insights into public sentiment, trending topics, and breaking news. These real-time datasets offer a window into millions of users' collective thoughts and opinions worldwide. You can use real-time Twitter data for various data science tasks, such as analyzing real-time tweets to understand public opinion on trending topics or political discussions or to develop real-time models to predict market trends, identify potential crises, or detect unusual patterns in public sentiment. You can even use this data to analyze real-time user behavior and sentiment to optimize social media campaigns, target advertising, and personalize user experiences. Source- Twitter API Ready to go on a data expedition with Pew Research Center's real-time streaming datasets? From Political Typology survey data to the News Interest Index data, the Pew Research Center provides access to a collection of real-time datasets through its American Trends Panel (ATP), offering insights into public opinion, political attitudes, and social trends in real-time. These fascinating datasets allow you to analyze sentiments during significant events, track sudden shifts in public opinion, analyze emerging trends in social attitudes, or even forecast cultural movements. With Pew's real-time datasets, you are not just observing trends; you are navigating the currents of societal shifts in real-time! Source- Pew Research Center Datasets Bonus Projects For Your Data Science Portfolio Series Here are some exciting projects you can add to your data science portfolio- Explore Interesting Datasets For Data Science Projects With ProjectPro In data science, datasets are the foundation for significant discoveries and innovative solutions. But how can you leverage the power of these datasets to solve real-world business problems? The key lies in understanding how to use these datasets effectively, transforming raw data into actionable insights. This is where ProjectPro comes in, offering a vast collection of data science projects designed to provide hands-on experience in implementing datasets and building effective data solutions. Through these projects, aspiring data scientists can explore data in real-world scenarios, learn how to tackle challenges, and gain valuable insights from various datasets. Whether analyzing customer behavior, predicting market trends, or optimizing operational efficiency, the ProjectPro repository provides the perfect platform to hone data science skills and gain invaluable experience. So, jumpstart your data science journey with ProjectPro, where datasets are not just numbers but tools to transform businesses and shape the future. FAQs on Datasets For Data Science Projects What are some factors to consider when choosing a dataset for data science projects? When choosing a data set for a data science project, it is essential to consider the following factors- Relevance- The data set should be relevant to the problem you are trying to solve. Quality- The data set should be of high quality, with minimal errors and missing values. Size- The data set should be large enough to train your ML models or perform your analysis. Accessibility- The data set should be easy to access and download. What are the different types of datasets for data science projects? Data science projects employ various types of datasets, including- Structured data- Organized data stored in tables, spreadsheets, or databases. Unstructured data- Text, images, audio, and video that lack a predefined format. Time-series data- Data collected over time, such as stock prices or sensor readings. Real-time data- Data that is continuously generated and streamed, such as social media feeds or financial market data.