HARNESSING THE POWER OF AI IN RESEARCH AND DATA SCIENCE: TOOLS, TECHNIQUES, AND LIVE DEMOS

Minha Hwang



AGENDA

AI IN RESEARCH AI TOOLS AI IN DATA SCIENCE

APPENDIX AI-ASSISTED CODING AITOOLS





AI IN RESEARCH

AI IN RESEARCH (1/2)

Research Process	What Can Be Used
Problem Formulation	Brainstorming with a strong reasoning model (e.g., GPT-03)
Literature Review and Reference	Deep Research: Your RA
 Research Design and Method Paper to Code: Reproducing Past Research Simulation Study Code Generation 	AI Chatbot with strong coding capabilities: GPT-03, Claude Sonnet 3.7

Data Collection and Acquisition

- Synthetic Data Generation from Prompt Engineering
- Open-Source Data Collection

AI Chatbot with strong coding capabilities: GPT-03, Claude Sonnet 3.7

Deep Research: Your RA

AI IN RESEARCH (2/2)

Research Process	What Can Be Used
Data Analysis and Modeling	Al-assisted coding and data science agent: e.g., coding, EDA, interactive app
Interpretation and Insight Generation	AI Chatbot with strong coding capabilities: GPT-03, Claude Sonnet 3.7
Writing and Visualization	Custom GPT <u>ChatGPT - Creative Writing Coach</u>
Paper Review and Critique	AI Chatbot with strong coding capabilities: GPT-03, Claude Sonnet 3.7
Research Paper Evaluation Publication Likelihood 	Custom GPT <u>ChatGPT - Research Paper Evaluation Framework</u>

- Publication Likelihood ٠
- Journal Outlet ٠

DEEP RESEARCH

	Key Characteristics
ChatGPT (Plus)	 Best Deep Research on the Market Once you make a request, ChatGPT asks clarifying questions. Usually takes 1 – 20 minutes to generate synthesized report with references Useful for initial exploratory research, literature review, references, paper critique / review
Gemini (Paid)	 2nd Best Deep Research on the Market Gemini shows "step-by-step research plan." You can modify this as you want. Usually takes I – 20 minutes Useful for similar tasks above. Integration with Google services (YouTube, Google Search) is a plus.
Grok 3 (Paid)	 Okay Deep Research on the Market You can use this to complement / validate deep research from above two services. Usually takes I – 20 minutes Useful for similar tasks above.
Perplexity	 Free Deep Research on the Market You can use this to complement / validate deep research from above two services. Usually takes 1 – 20 minutes Useful for similar tasks above. Less comprehensive compared to first two.
M365 Copilot	 Only Enterprise Deep Research on Market: Secure Enterprise Data from Knowledge Graph Add "Researcher" on M365 Copilot Need M365 Enterprise subscription Usually takes I – 20 minutes

LIVE DEMO



AI IN DATA SCIENCE

GEN AI DATA SCIENTIST VS. TRADITIONAL DATA SCIENTIST

	Gen Al Data Scientist	(Traditional) Data Scientist
Coding	 Leverages coding and app development agents such as GitHub Copilot, Windsurf, Rue, Cursor.ai, Replit, Lovable, v0, Designs code at a high-level; Coding in English with step-by-step plans Focuses on efficiency by utilizing Al chatbot assistant and agentic development 	 Writes code end-to-end in a familiar programming language. Invests significant time learning specific languages, libraries, and frameworks. Seeks solutions through peer consultation or resources like Stack Overflow.
Data analysis	 Utilizes conversational data science agents for exploratory data analys (EDA). High-level design of data pipeline and analysis plan Uses Al-assisted visualization tools for efficient chart and graph creation Handles unstructured data (text, image, voice) with Al-powered tools 	 Either skips EDA or relies on expert knowledge. Writes the entire data analysis pipeline manually. Spends significant time fine-tuning visualizations (fonts, colors, layout Primarily focuses on tabular (numerical) data.
Focus	 Emphasizes model evaluation and data curation, including synthetic data generation. Works with pre-trained models as a service, enabling limited adaptation. Utilizes pre-trained models from platforms such as Hugging Face, Azure AI Foundry, Google AI Studio, and AWS Bedrock. 	 Concentrates on model training, including relevance engineering and recommendation algorithms. Develops models from scratch, focusing on model architecture, hyper-parameter optimization and training recipes.
Key Skills	 Strong problem definition and problem-solving abilities. Structured planning and efficient execution. Excellent communication skills. Deep understanding of Al mechanisms and functionalities. 	 Strong technical expertise in model training, statistics, and machine learning/deep learning (ML/DL). Limited emphasis on soft skills and communication.

DATA SCIENCE AGENT



In addition, use agentic code IDEs

- GitHub Copilot (Agent mode)
- Cursor
- Windsurf



GITHUB COPILOT: HOW TO INSTALL GITHUB COPILOT AND GITHUB COPILOT CHAT EXTENSION ON VS CODE

Are you ready to supercharge your coding experience with AI-powered assistance? Follow these simple steps to install GitHub Copilot and the GitHub Copilot Chat extension on Visual Studio Code (VS Code):

Step 1: Install Visual Studio Code

- 1. Download and install Visual Studio Code from the [official website](https://lnkd.in/g7y9giDJ).
- 2. Launch Visual Studio Code after installation.

Step 2: Install GitHub Copilot Extension

- 1. Open Visual Studio Code.
- 2. Go to the Extensions view by clicking on the Extensions icon in the Activity Bar on the side of the window or by pressing `Ctrl+Shift+X`.
- 3. In the search bar, type "GitHub Copilot" and select the GitHub Copilot extension from the list.
- 4. Click the "Install" button to install the extension.

Step 3: Sign In to GitHub

1. After installing the GitHub Copilot extension, you will be prompted to sign in to your GitHub account.

2. Follow the on-screen instructions to sign in and authorize GitHub Copilot.

Step 4: Install GitHub Copilot Chat Extension

- 1. Go to the Extensions view again by clicking on the Extensions icon or pressing `Ctrl+Shift+X`.
- 2. In the search bar, type "GitHub Copilot Chat" and select the GitHub Copilot Chat extension from the list.
- 3. Click the "Install" button to install the extension.

Step 5: Start Using GitHub Copilot and Copilot Chat

- 1. Open a new or existing project in Visual Studio Code.
- 2. Start typing code, and GitHub Copilot will provide code suggestions as you type.
- 3. To use the Copilot Chat, open the Chat view by pressing `Ctrl+Alt+I` (Windows/Linux) or `^ 光I` (Mac).
- 4. Enter your coding-related questions or tasks in the chat input field and let Copilot Chat assist you.

LIVE DEMO



SAMPLE PROMPT: KAGGLE HOUSE PRICE PREDICTION

Adding Context: File or Folder \House_Price\data_description.txt

Prompt

Solve the problem below with the desired goal and metric. You can follow the specified key steps below. This will be provided one by one

Goal:

It is your job to predict the sales price for each house. For each Id in the test set, you must predict the value of the SalePrice variable.

Metric:

Submissions are evaluated on <u>Root-Mean-Squared-Error (RMSE)</u> between the logarithm of the predicted value and the logarithm of the observed sales price. (Taking logs means that errors in predicting expensive houses and cheap houses will affect the result equally.)

Key Steps:

- (1) Import train.csv
- (2) Import required libraries for the analyses
- (3) Use the data description to understand the train.csv data file
- (4) Generate summary statistics and describe data
- (5) Conduct exploratory analyses
- (6) Investigate key drivers of housing prices
- (7) Create a regression model to explain housing prices
- (8) Synthesize key insights from the analyses

WHAT CAN YOU DO? (1/5)

Task	Prompt	Output
Data Cleaning	"Clean the dataset by removing rows with missing values and encoding categorical variables."	The agent returns a cleaned dataset with missing rows removed and categorical variables encoded using one-hot encoding.
Exploratory Data Analysis (EDA)	"Provide summary statistics, correlations across variables, histograms, and box plots for the 'sales' dataset."	The agent generates descriptive statistics (mean, median, standard deviation), correlations, and visualizations like histograms and box plots for the 'sales' dataset.
Feature Engineering	"Create a new feature representing the interaction between 'age' and 'income'."	The agent adds a new feature to the dataset by multiplying 'age' and 'income' columns, capturing their interaction effect.

- You just need to know "what to do"
- Do not need to look up PyPl or Stack Overflow for technical details and syntax

WHAT CAN YOU DO? (2/5)

Task	Prompt	Output
Model Building & Evaluation	"Train a logistic regression model to predict customer churn and evaluate its accuracy.	The agent trains a logistic regression model, outputs the accuracy score, and provides a confusion matrix to assess performance.
Natural Language Processing (NLP)	"Analyze the sentiment of customer reviews in the 'reviews' column."	The agent processes the text data, assigns sentiment scores (positive, negative, neutral) to each review, and summarizes the overall sentiment distribution.
Time Series Analysis	"Forecast the next 12 months of sales using the historical 'monthly_sales' data."	The agent fits a time series model (e.g., ARIMA), forecasts future sales, and plots the predicted values alongside historical data.

- You just need to know "what to do"
- Do not need to look up PyPI or Stack Overflow for technical details and syntax

WHAT	CANYOU	DO?	(3/5)
------	--------	-----	-------

Task	Prompt	Output
Interactive Dashboard Creation	"Create an interactive dashboard showing key performance indicators (KPIs) such as conversion rate, click- through rate, and return on investment for the marketing campaign."	The agent develops a dashboard with visualizations like bar charts and line graphs, displaying KPIs such as conversion rate, click-through rate, and return on investment.
Data Summarization & Reporting	"Summarize the key findings from the quarterly sales data and generate a report highlighting trends and anomalies."	The agent analyzes the sales data, identifies significant trends (e.g., a 15% increase in Q2 sales), detects anomalies (e.g., a sudden drop in a specific region), and compiles a comprehensive report with visualizations and executive summaries.
Data Integration & Merging	"Merge the customer demographic dataset with the transaction history dataset to create a unified view for analysis."	The agent performs data cleaning, resolves discrepancies between datasets, and merges them based on common identifiers, resulting in a consolidated dataset ready for further analysis.

- You just need to know "what to do"
- Do not need to look up PyPI or Stack Overflow for technical details and syntax

WHAT	CANYOU	DO?	(4/5)
------	--------	-----	-------

Task	Prompt	Output
Model Optimization	"Optimize the hyperparameters of the random forest model to improve prediction accuracy."	The agent conducts hyperparameter tuning using techniques like grid search or random search, evaluates model performance, and outputs the optimal parameters that enhance accuracy.
Text Classification	"Classify customer feedback into categories: Positive, Negative, or Neutral."	The agent processes the textual feedback, applies natural language processing techniques, and assigns each feedback entry to one of the specified categories, providing a summary of the distribution.
Anomaly Detection	"Identify any unusual patterns or outliers in the website traffic data over the past month."	The agent analyzes the traffic data, detects anomalies such as sudden spikes or drops in visits, and highlights potential causes or correlations with external events.

Data analysis in English : Prompting is All You Need You just need to know "what to do"

- Do not need to look up PyPI or Stack Overflow for technical details and syntax •

WHAT CAN YOU DO? (5/5)

Task	Prompt	Output
Trend Analysis	"Analyze the trend of product returns over the last year and identify any seasonal patterns."	The agent examines the return data, identifies trends (e.g., increased returns during holiday seasons), and provides visualizations to illustrate these patterns.
Statistical Testing	"Conduct a t-test to determine if there's a significant difference in average purchase amounts between two customer groups."	The agent performs the t-test, calculates the p- value, and interprets the results to indicate whether the difference is statistically significant.
Data Pipeline Automation	"Automate the data ingestion and preprocessing steps for the daily sales data pipeline."	The agent sets up automated scripts or workflows that fetch daily sales data, clean and preprocess it, and store it in a designated database or data warehouse.

- You just need to know "what to do"
- Do not need to look up PyPI or Stack Overflow for technical details and syntax

APPENDIX

AITOOLS

AITOOLS (1/4)

- Chatbots · CHATGPT: <u>HTTPS://CHATGPT.COM/</u> DEEP RESEARCH (RESEARCH ASSISTANCE), GPT STORE, OPERATOR (AGENT), AGENTS SDK, REASONING MODEL (O1, O3, O1-PRO)
 - MICROSOFT COPILOT: <u>HTTPS://COPILOT.MICROSOFT.COM/</u> OFFICE INTEGRATION (E.G. POWER-POINT SLIDE GENERATION FROM WORD FILE, GITHUB COPILOT), FREE REASONING MODEL ACCESS (O3-MINI-HIGH)
 - CLAUDE.AI: <u>HTTPS://CLAUDE.AI/</u> FRONT-END WEB/APP DEVELOPMENT, STRONG CODING, CREATIVE WRITING, HYBRID-REASONING MODEL (SONNET 3.7), COMPUTER USE
 - GOOGLE GEMINI: <u>HTTPS://GEMINI.GOOGLE.COM/</u> (TRUE) MULTI-MODAL MODEL, GOOGLE SERVICE INTEGRATION E.G.YOUTUBE), REASONING MODEL WITH TRANSPARENT INTERMEDIATE STEPS (THINK EXPERIMENTAL), DEEP RESEARCH
 - DEEPSEEK-RI: <u>HTTPS://DEEPSEEK.COM/</u> OPEN-WEIGHT REASONING MODEL WITH INTERMEDIATE STEPS (CENSORING), LOW-COST TOKENS, DISTILLED MODELS
 - **GROK 3**: <u>HTTPS://GROK.COM/</u>– **DEEP RESEARCH**, **REASONING MODEL**, X PLATFORM INTEGRATION
 - ERPLEXITY: <u>HTTPS://WWW.PERPLEXITY.AI/</u> **ANSWER ENGINE**, MODEL CHOICES, **STRONG RAG**, **DEEP**RESEARCH, CUSTOM DISTILLED MODEL (SONAR)
 - LIQUD.AI: <u>HTTPS://PLAYGROUND.LIQUID.AI/</u>- NON-TRANSFORMER-BASED ARCHITECTURE LIQUID NN: FAST INFERENCE AND ON-DEVICE FOCUS, MULTIMODAL

AITOOLS (2/4)

Coding	 CURSOR.AI: <u>CURSOR - THE AI CODE EDITOR</u> GITHUB COPILOT: <u>GITHUB COPILOT · YOUR AI PAIR PROGRAMMER</u> WINDSURF (AGENTIC IDE): <u>WINDSURF EDITOR BY CODEIUM</u> ROOCODE: <u>ROO CODE – YOUR AI-POWERED DEV TEAM IN VS CODE</u> CLINE (AGENTIC IDE): <u>CLINE - AI AUTONOMOUS CODING AGENT FOR VS CODE</u> AUGMENT CODE: <u>HTTPS://WWW.AUGMENTCODE.COM/</u> CODY: <u>HTTPS://CODDY.TECH/</u> (DUOLINGO FOR CODING)
Web / App Dev	 REPLIT: <u>HTTPS://REPLIT.COM/</u> LOVABLE: <u>HTTPS://LOVABLE.DEV/</u> CLAUDE CODE: <u>HTTPS://DOCS.ANTHROPIC.COM/EN/DOCS/AGENTS-AND-TOOLS/CLAUDE-CODE/OVERVIEW/</u> V0 BY VERCEL: <u>V0 BY VERCEL</u>
Research Assistance	 CHATGPT DEEP RESEARCH PAID SUBSRIPTION): <u>HTTPS://CHATGPT.COM/</u> GEMINI DEEP RESEARCH (PAID SUBSCRIPTION): <u>HTTPS://GEMINI.GOOGLE.COM/APP</u> GROK 3 BETA DEEP RESEARCH: <u>HTTPS://GROK.COM/</u> PERPLEXITY DEEP RESEARCH: <u>PERPLEXITY</u> HUGGING FACE OPEN-SOURCE DEEP RESEARCH: <u>HTTPS://HUGGINGFACE.CO/BLOG/OPEN-DEEP-RESEARCH</u>

AITOOLS (3/4)



AITOOLS (4/4)

	•	OPEN AI OPERATOR: INTRODUCING OPERATOR OPENAI
Agent	•	OPENAI AGENTS SDK: OPENAI AGENTS SDK
	•	LANG CHAIN AND LANG GRAPH: <u>LANGGRAPH</u>
	•	MICROSOFT AUTOGEN: <u>HTTPS://GITHUB.COM/MICROSOFT/AUTOGEN</u>
	•	CREWAI: <u>HTTPS://GITHUB.COM/CREWAIINC</u>
	•	OPEN-SOURCE PRE-TRAINED MODELS: <u>MODELS - HUGGING FACE</u>
Online	•	LLM ONLINE COURSES: <u>COURSES - DEEPLEARNING.AI</u>
Learning	•	HUGGING FACE LEARN: <u>HUGGING FACE – LEARN</u>
	•	OPENAI ACADEMY: <u>OPENAI ACADEMY</u>
	•	5-DAY GEN AI INTENSIVE COURSE – KAGGLE/GOOGLE: <u>5-DAY GEN AI INTENSIVE</u>
		<u>COURSEWITH GOOGLE LEARN GUIDE KAGGLE</u>
	•	NOTEBOOK LM: Google NotebookLM Note Taking & Research Assistant Powered by AI
Application	s .	CAREER DREAMER: Explore Your Possibilities with Career Dreamer - Grow with Google
	•	HELIX: A Vision-Language-Action Model for Generalist Humanoid Control
	•	TYPEFACE: Typeface - Enterprise Generative AI Platform for Marketing & Content Creation
	•	SYNTHESIA – Text-to-video with Avatars: <u>https://www.synthesia.io/</u>

AI-ASSITED CODING DEMO

- COLAB NOTEBOOK:

LIVE DEMO

<u>AI_CODING/AI_ASSISTED_CODING.IPYNBAT</u> <u>MAIN · DRSQUARE/AI_CODING</u>