

Can Agentic AI Match the Performance of Human Data Scientists?

An Luo*, Jin Du*, Fangqiao Tian*, Xun Xian†, Robert Specht*, Ganghua Wang‡, Xuan Bi§, Charles Fleming¶, Jayanth Srinivasa¶, Ashish Kundu¶, Mingyi Hong†, Jie Ding*

*School of Statistics, University of Minnesota, Minneapolis, MN, USA

†Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN, USA

‡Data Science Institute, University of Chicago, Chicago, IL, USA

§Carlson School of Management, University of Minnesota, Minneapolis, MN, USA

¶Cisco Research, San Jose, CA, USA

Abstract—Data science plays a critical role in transforming complex data into actionable insights across numerous domains. Recent developments in large language models (LLMs) have significantly automated data science workflows, but a fundamental question persists: Can these agentic AI systems truly match the performance of human data scientists who routinely leverage domain-specific knowledge? We explore this question by designing a prediction task where a crucial latent variable is hidden in relevant image data instead of tabular features. As a result, agentic AI that generates generic codes for modeling tabular data cannot perform well, while human experts could identify the important hidden variable using domain knowledge. We demonstrate this idea with a synthetic dataset for property insurance. Our experiments show that agentic AI that relies on generic analytics workflow falls short of methods that use domain-specific insights. This highlights a key limitation of the current agentic AI for data science and underscores the need for future research to develop agentic AI systems that can better recognize and incorporate domain knowledge.

Index Terms—Agents, automated data science, human-AI teaming, large language models, synthetic data.

I. INTRODUCTION

Data science is a central interdisciplinary field that blends statistics, computer science, and domain expertise to extract actionable insights from complex, heterogeneous data [1], [2]. By transforming raw information into knowledge and value, data science drives innovation and shapes decision-making in science, industry, healthcare, finance, and beyond [3], [4].

Recent advancements in large language models (LLMs) have significantly accelerated the automation of data science workflows. LLMs such as GPT-4 [5] and Claude [6] have demonstrated impressive capabilities in automating code generation and executing regular machine learning tasks [7]–[10]. These developments offer promising potential for streamlining common analytical processes and reducing the manual workload of human data scientists.

Despite these advancements, there is still a lack of understanding of whether agentic AI really perform as good as human data scientists. In practice, human data scientists consistently rely on specialized knowledge about the data

or task and incorporate crucial nuances that enhance model performance [11]–[15]. Such domain-driven decisions are often subtle yet essential, as they address complexities not captured by typical analytics workflows. However, current research on LLM-driven data science has largely focused on generating generic code and pipeline executions [7], [10]. These approaches often neglect the domain-specific knowledge needed for complex, real-world problems. Meanwhile, existing evaluation benchmarks such as MLE-bench [16] and DSbench [17] aim to assess predictive performance, but do not test whether agentic AI can effectively leverage domain insights outside tabular data. The above observations motivate a fundamental question: *Can agentic AI, which typically relies on generic code generation, truly match the performance of human data scientists who could apply domain knowledge?*

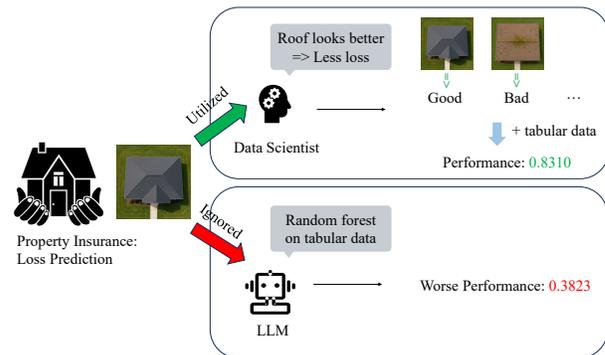


Fig. 1. **Comparison of data scientists and agentic AI approaches to loss prediction in property insurance.** The upper panel shows how a human data scientist leverages domain knowledge: by interpreting roof images to infer the critical latent variable (Roof Health) and incorporating it with tabular data, they can achieve substantially high predictive performance (normalized Gini = 0.8310). The lower panel depicts an agentic AI’s approach, which applies standard tabular modeling while ignoring the image modality and domain-specific cues, resulting in much worse performance (normalized Gini = 0.3823). This demonstrates the empirical gap between human data scientist and agentic AI performance when domain knowledge is necessary.

To address this question, our paper presents an experimental study on a carefully curated dataset that mimics the complexity of real-world data science problems. The use of synthetic

data allows us to control key latent variables that introduce complexity behind observed feature variables. This controlled setup reveals important differences: human data scientists are often able to explicitly identify and leverage domain-specific cues, whereas agentic AI may rely on generic algorithms that do not fully capture the influence of latent factors. Figure 1 illustrates this idea in our property insurance setup.

Our main contributions are summarized below.

- We design a synthetic dataset that clearly illustrates a fundamental gap between agentic AI (which generates generic, tabular-focused code) and human data scientists (who leverage domain knowledge embedded in images).
- We empirically quantify this performance gap and demonstrate the importance of domain knowledge in achieving excellent prediction performance.

Through this work, we aim to highlight the need for future research to improve agentic AI’s ability to identify and use domain-specific knowledge from multimodal data sources.

II. DESIGN OF SYNTHETIC DATA

A. General Design Principles

To evaluate the limits of both human data scientists and agentic AI, we generate a controlled synthetic dataset with a hidden latent factor that affects the prediction target but is not present in the tabular features. Instead, this latent variable is embedded in a secondary modality (here, overhead images) to ensure it can be accessed only through domain knowledge and not by generic code. While the approach could generalize to other modalities such as text or audio, we focus on images to illustrate the mechanism. Importantly, these images are crafted so that a knowledgeable data scientist can interpret the latent variable in the context of property insurance, making the challenge meaningful and realistic.

To accomplish this, we use a text-to-image model with engineered prompts to ensure that the generated images faithfully reflect the intended values of the latent variable. This design allows us to examine the gap between generic AI pipelines and human data scientists that can look for the incorporation of domain knowledge.

B. Data Curation

The data science task specifies a policy table as tabular dataset, and each policy is associated with an image. Their goal is to predict each home’s total insured loss in the next policy year, Y_p . The key latent variable is *RoofHealth*, a three-level variable—Good, Fair, or Bad. This variable is never shown in the policy table but can be inferred from the image. Below we present how we create this synthetic dataset.

Step 1: Generate Structured Policy Features

For each policy $p = 1, \dots, n$ we draw:

- 1) PolicyID: “POL-000001”, ...
- 2) HouseValue: $X_{\text{val},p} \sim \text{LogNormal}(12.9, 0.45)$ (median $\approx \$403\text{k}$).
- 3) HouseAge: $X_{\text{age},p} \sim 120 \text{Beta}(4, 3)$ (≈ 40 yr median).

- 4) WallType: $X_{\text{wall},p} \sim \text{Bernoulli}(\text{Wood, Brick})$ with probabilities $\{0.9, 0.1\}$.
- 5) AreaRisk: $X_{\text{risk},p} \sim \text{Beta}(2, 5)$ (0–1 storm exposure).
- 6) CreditScore: $X_{\text{cred},p}$ drawn from the US FICO distribution (300–850).
- 7) **RoofHealth** (latent): compute

$$S_p = 0.02 X_{\text{age},p} + 3.0 X_{\text{risk},p} - 2.0(X_{\text{cred},p}/850) + \varepsilon_p, \quad \varepsilon_p \sim N(0, 1),$$

then assign Good, Fair, or Bad by partitioning S_p at the 55th and 80th percentiles of all scores.

Only columns 1–6 are released as features in the policy table.

Step 2: Create Roof Images

Each policy gets one 1024×1024 PNG image synthesised with `gpt-image-1` using the prompt template below:

Prompt Template For Roof Image Generation

```
Realistic straight-down aerial photo of a detached house, full roof and surrounding lawn in view, {roof_style} roof with {shingle_color} shingles, {surface_descriptor}, {edge_descriptor}, {extra_descriptor}.
```

In the template, the roof style is sampled from $\{gable, hip, flat, mansard, shed\}$; shingle color is from $\{dark-gray, light-gray, brown, black, red-tile\}$. Descriptors are from examples as shown in Table I so that Good/Fair/Bad roofs differ in surface integrity, edge condition, and so on. In this way, each image faithfully represents the intended roof condition. See examples of generated images in Figure 2.

TABLE I
DESCRIPTOR EXAMPLES USED IN PROMPT TEMPLATE FOR IMAGES GENERATION UNDER EACH ROOF HEALTH CATEGORY

Roof Health	Surface Example	Edge Example
Good	even rows of intact shingles	well-sealed ridge lines
Fair	slightly faded shingles	ridge line with mild wear
Bad	multiple missing shingles	damaged or sagging ridge

Step 3: Simulate Next-Year Loss Y_p

a) *Claim count*: With $\alpha_{\text{rh}} = \{0, 1.2, 2.4\}$ for Good/Fair/Bad,

$$\lambda_p = \exp\left(-3.0 + 0.03 \ln \frac{X_{\text{val},p}}{250,000} + 0.01 X_{\text{age},p} + 0.05 X_{\text{risk},p} + \alpha_{\text{rh}}(\text{RoofHealth}_p)\right), \quad (1)$$

$$N_p \sim \text{NegBinom}(r = 10, \text{mean} = \lambda_p).$$

b) *Claim loss*: With $\beta_{\text{rh}} = \{0, 1.0, 2.0\}$,

$$\mu_p = 7.0 + 0.02 \mathbb{1}(X_{\text{wall},p} = \text{Wood}) + 0.02 X_{\text{risk},p} + \beta_{\text{rh}}(\text{RoofHealth}_p), \quad (2)$$

$$Z_{p,j} \sim \Gamma(k = 2, \theta = \exp(\mu_p)/2).$$

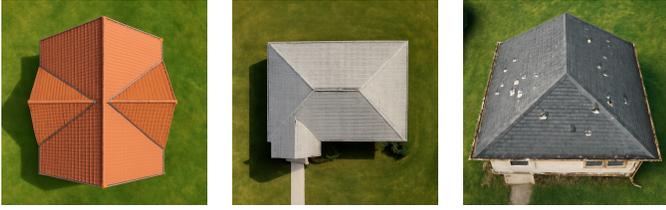


Fig. 2. Example overhead roof images generated for our synthetic property insurance dataset. Each image visually encodes a key latent variable, `RoofHealth`, with three possible states: (a) Good, (b) Fair, and (c) Bad. This variable is never released directly in the tabular data but can be inferred from domain-specific visual cues, such as surface, edge, and extra details such as flashing condition and debris. To create this setting, we use a text-to-image model with carefully designed prompts to ensure each image faithfully represents the intended roof condition. This design allows us to rigorously compare standard agentic AI pipelines (which use only tabular data) against approaches or human experts capable of incorporating additional domain knowledge from the image modality.

c) *Total loss*:

$$Y_p = \sum_{j=1}^{N_p} Z_{p,j}, \quad \text{if } N_p = 0 \text{ then } Y_p = 0.$$

The training file exposes `NextYearLoss = Yp`; the test file omits it for evaluation. See Figure 3 for an illustration of how the outcome variable, `Next-Year Loss`, is generated.

The construction of our synthetic property insurance dataset is grounded in established actuarial practice and empirical research. Roof condition is an important factor in property risk and claims, but it is often not directly available in tabular data [18], [19]. Our use of roof images is intended to reflect this real-world limitation. The target outcome, next-year loss, is generated using a compound frequency-severity model. This approach is a standard actuarial method for property insurance loss modeling [20], [21]. Together, these design choices ensure our dataset realistically captures the complexities of property insurance prediction.

III. PERFORMANCE STUDY: GENERIC PIPELINE FROM AI VS. DOMAIN KNOWLEDGE USAGE FROM HUMAN

A. Experimental Setting

Our goal is to evaluate how the use of domain knowledge impacts predictive performance when crucial information is embedded in the roof images. To do this, we compare three groups of modeling strategies. The first group simulates a generic agentic AI approach that uses only tabular data. The second group represents methods a human data scientist might use, combining tabular data with different ways of extracting information from roof images. The final group is an oracle model that has access to the true latent variables and the underlying data generation process, serving as the best achievable benchmark. By comparing these approaches, we can quantify the value of domain knowledge and highlight the limitations of generic AI workflows.

Below, we outline the specific modeling strategies in each group and how they reflect use of domain knowledge:

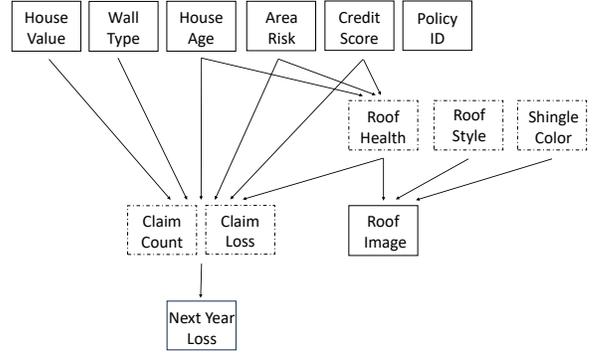


Fig. 3. Illustration of data generating process for property insurance. The diagram shows how each policy’s outcome, `Next-Year Loss`, is generated. Dotted lines surround latent variables that are hidden. The process unfolds as follows: (1) Structured policy features (e.g., `House Value`, `House Age`, `Wall Type`, `Area Risk`, `Credit Score`) are generated for each policy. (2) A latent variable, `Roof Health` (Good, Fair, Bad), is determined by a function of selected features, but is not included in the released tabular data. Instead, it is visually encoded in an accompanying roof image, which is generated for each policy using a random combination of roof style and shingle color. (3) Claim count and claim loss are simulated using both policy features and the latent `RoofHealth`. The total insured loss for the next year (Y_p) is calculated as the sum of all claim loss. To achieve optimal prediction, the hidden `Roof Health` must be inferred from the roof image.

- 1) **Agentic AI (Generic pipeline)** Only the tabular features are used and no image data is included. This matches how agentic AI would typically apply generic code to a standard tabular prediction problem.
- 2) **Data Scientists (Image use).** Both tabular features and image data are utilized. We consider several practical ways a human data scientist might incorporate image information. One approach is to extract features from the images using a pretrained CLIP model [22] and either use these features directly or cluster them into categories for use in the predictive model. Another approach is to apply a vision-language model `gpt-4o-mini` to extract the `RoofHealth` label from the images. Finally, we include an ideal scenario where the data scientist perfectly labels the true `RoofHealth` for each image. This represents the best possible use of domain expertise.
- 3) **Oracle (Best achievable).** This method uses the exact data-generation formulas and the true `RoofHealth`. It calculates the predicted loss as the exact product of expected claim counts and severities, i.e.

$$\hat{Y}_p = \lambda_p \times \exp(\mu_p),$$

where λ_p and μ_p are computed from Equation 1 and 2 respectively. This gives the Bayes-optimal expected loss and any remaining error reflects only inherent randomness in claims.

All experiments use the synthetic data with 2000 policies generated as described in Section II. Among them, 1000 policies are for training and the other 1000 policies are held

out for evaluation. Each policy carries six tabular features and a 1024×1024 overhead roof image.

B. Predictive Performance Evaluation: Normalized Gini

We measure predictive performance using *normalized Gini coefficient*, a standard practice and widely used metric in the insurance domain to evaluate predictive models [23]–[26]. It is a rank-based metric that captures how well predicted scores prioritize higher value observations and is appropriate for loss outcomes, which are often heavy-tailed.

Let $\{(y_i, \hat{y}_i)\}_{i=1}^n$ be the true responses and model predictions. Sort the pairs by descending predicted value, yielding $(y_{(1)}, \hat{y}_{(1)}), \dots, (y_{(n)}, \hat{y}_{(n)})$. Define the cumulative true sum

$$C_k = \sum_{i=1}^k y_{(i)}, \quad Y = \sum_{i=1}^n y_i.$$

The *raw* Gini coefficient is then

$$G_{\text{raw}}(y, \hat{y}) = \frac{1}{n} \sum_{k=1}^n \frac{C_k}{Y} - \frac{n+1}{2n}.$$

To make this metric in $[-1, 1]$ and comparable across datasets, it is normalized by the “perfect” Gini achieved when $\hat{y}_i = y_i$:

$$G_{\text{norm}}(y, \hat{y}) = \frac{G_{\text{raw}}(y, \hat{y})}{G_{\text{raw}}(y, y)}.$$

$G_{\text{norm}} = 1$ indicates a perfect ranking, and $G_{\text{norm}} = 0$ correspond to a random ordering. $G_{\text{norm}} < 0$ would mean predictions are worse than random. The higher normalized Gini signals better model performance.

C. Performance Gap: From Generic Pipeline to Oracle

In Table II, we compare predictive performance across the modeling approaches described above. This highlights the gap between generic agentic AI pipelines that use only tabular data and methods used by human data scientists that incorporate domain knowledge from images. For methods using image features, the “Corr.” column shows how well the extracted variable aligns with the true underlying roof health.

The first row, **Agentic AI (Generic pipeline)**, reflects typical agentic AI workflows that use only tabular data and ignore image and domain-specific information. This approach represents the performance of a standard pipeline LLMs would generate without domain insight, achieving a normalized Gini of 0.3823. In this case, when important information is hidden in images, standard pipelines struggle to achieve good results.

The next group, **Data Scientists (Image use)**, includes several practical strategies for incorporating image information, just as a human data scientist might try. Using naive clustering of CLIP embeddings as categorical features provides some improvement (Gini 0.5042), but does not fully capture the signal (correlation with true roof health is 0.40). Feeding the full CLIP features into the model yields much better results (Gini 0.7719). Extracting the RoofHealth label from images with a vision-language model (gpt-4o-mini) also boosts performance (Gini 0.7271), with a much higher correlation to the true latent variable (0.81). When the model is given the true

RoofHealth label as if a human labels the images perfectly, the performance almost matches the best possible (Gini 0.8310). The clear trend is that methods using image-based domain knowledge achieve much higher predictive performance.

The last row, **Oracle (Best achievable)** represents the optimal achievable performance (Gini 0.8379), where predictions utilize the exact underlying generative mechanism and the true RoofHealth labels. This tier’s result reflects only inherent randomness in claims data and sets a practical upper bound for predictive performance.

The improvements observed across these levels clearly demonstrate the importance of domain-specific knowledge in data science and highlight the limitations of generic, tabular-only approaches typically employed by current agentic AI.

TABLE II

NORMALIZED GINI FOR DIFFERENT MODELING APPROACHES. ROWS ARE GROUPED BY METHOD TYPE: (1) **AGENTIC AI (GENERIC PIPELINE)**: RANDOM FOREST USING ONLY TABULAR DATA; (2) **DATA SCIENTISTS (IMAGE USE)**: RANDOM FOREST LEVERAGING IMAGE FEATURES OR ROOF-HEALTH LABELS DERIVED FROM IMAGES, SIMULATING VARYING LEVELS OF DOMAIN KNOWLEDGE; (3) **ORACLE (BEST ACHIEVABLE)**: MODEL WITH ACCESS TO THE TRUE LATENT VARIABLE AND THE GENERATIVE PROCESS. CORRELATION (**CORR.**) INDICATES HOW WELL THE EXTRACTED VARIABLE ALIGNS WITH THE TRUE UNDERLYING ROOF HEALTH. RF = RANDOM FOREST; CLIP = IMAGE FEATURE EXTRACTOR.

Method	Corr.	Normalized Gini
Agentic AI (Generic pipeline)		
RF (tabular only)	—	0.3823
Data Scientists (Image use)		
RF + CLIP clustered as 3 labels	0.4009	0.5042
RF + CLIP features of images	—	0.7719
RF + RoofHealth extracted by gpt-4o-mini	0.8062	0.7271
RF + true RoofHealth	1.0000	0.8310
Oracle (Best achievable)		
Oracle (Bayes-optimal expected loss)	1.0000	0.8379

IV. CONCLUSION

In this work, we illustrate that agentic AI cannot match the performance of human data scientists in our controlled setting, using a carefully designed synthetic dataset. The dataset is constructed so that an important latent variable is hidden within the image data. As a result, generic algorithms that rely solely on tabular data become insufficient, which is precisely the approach typically employed by agentic AI. In contrast, a human data scientist equipped with domain knowledge can correctly identify and utilize this latent information from the images, resulting in substantially improved performance. This underscores a limitation of current agentic AI for data science: they typically generate generic algorithms without adequately incorporating domain-specific insights. We hope this work will inspire further research into building agentic AI that can critically incorporate and utilize domain-specific knowledge, thereby bridging the gap between automated workflows and expert human performance.

REFERENCES

- [1] L. Cao, “Data science: A comprehensive overview,” *ACM Computing Surveys*, 2017.

- [2] M. L. Brodie, "Defining data science: a new field of inquiry," *arXiv preprint arXiv:2306.16177*, 2023.
- [3] V. Grossi, F. Giannotti, D. Pedreschi, P. Manghi, P. Pagano, and M. Assante, "Data science: a game changer for science and innovation," *International Journal of Data Science and Analytics*, 2021.
- [4] G. S. Blair, P. A. Henrys, A. A. Leeson, J. Watkins, E. F. Eastoe, S. G. Jarvis, and P. J. Young, "Data science of the natural environment: A research roadmap," *Frontiers in Environmental Science*, 2019.
- [5] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.
- [6] Anthropic, "Claude 3.7 sonnet and claude code," 2025, accessed: 2025-05-18. [Online]. Available: <https://www.anthropic.com/news/claude-3-7-sonnet>
- [7] Z. Jiang, D. Schmidt, D. Srikanth, D. Xu, I. Kaplan, D. Jacenko, and Y. Wu, "AIDE: AI-driven exploration in the space of code," *arXiv preprint arXiv:2502.13138*, 2025.
- [8] S. Hong, Y. Lin, B. Liu, B. Wu, D. Li, J. Chen, J. Zhang, J. Wang, L. Zhang, M. Zhuge, T. Guo, T. Zhou, W. Tao, W. Wang, X. Tang, X. Lu, X. Liang, Y. Fei, Y. Cheng, Z. Gou, Z. Xu, C. Wu, L. Zhang, M. Yang, and X. Zheng, "Data Interpreter: An LLM agent for data science," in *Findings of the Association for Computational Linguistics: ACL 2025*, 2025.
- [9] Z. Liang, F. Wei, W. Xu, L. Chen, Y. Qian, and X. Wu, "I-MCTS: Enhancing agentic AutoML via introspective monte carlo tree search," *arXiv preprint arXiv:2502.14693*, 2025.
- [10] Z. Li, Q. Zang, D. Ma, J. Guo, T. Zheng, M. Liu, X. Niu, Y. Wang, J. Yang, J. Liu, W. Zhong, W. Zhou, W. Huang, and G. Zhang, "AutoKaggle: A multi-agent framework for autonomous data science competitions," *arXiv preprint arXiv:2410.20424*, 2024.
- [11] Y. Mao, D. Wang, M. J. Muller, I. Baldini, and C. Dugan, "How data scientists work together with domain experts in scientific collaborations," *Proceedings of the ACM on Human-Computer Interaction*, vol. 3, pp. 1 – 23, 2019.
- [12] A. X. Zhang, M. J. Muller, and D. Wang, "How do data science workers collaborate? roles, workflows, and tools," *Proceedings of the ACM on Human-Computer Interaction*, vol. 4, pp. 1 – 23, 2020.
- [13] Z. Lin, A. Marin-Llobet, J. Baek, Y. He, J. Lee, W. Wang, X. Zhang, A. J. Lee, N. Liang, J. Du, J. Ding, N. Li, and J. Liu, "Spike sorting AI agent," *bioRxiv*, 2025.
- [14] Z. Lin, W. Wang, A. Marin-Llobet, Q. Li, S. D. Pollock, X. Sui, A. Aljovic, J. Lee, J. Baek, N. Liang, X. Zhang, C. K. Wang, J. Huang, M. Liu, Z. Gao, H. Sheng, J. Du, S. J. Lee, B. Wang, Y. He, J. Ding, X. Wang, J. R. Alvarez-Dominguez, and J. Liu, "Spatial transcriptomics AI agent charts hPSC-pancreas maturation in vivo," *bioRxiv*, 2025.
- [15] A. Luo, X. Xian, J. Du, F. Tian, G. Wang, M. Zhong, S. Zhao, X. Bi, Z. Liu, J. Zhou, J. Srinivasa, A. Kundu, C. Fleming, M. Hong, and J. Ding, "AssistedDS: Benchmarking how external domain knowledge assists LLMs in automated data science," in *The 2025 Conference on Empirical Methods in Natural Language Processing*, 2025.
- [16] J. S. Chan, N. Chowdhury, O. Jaffe, J. Aung, D. Sherburn, E. Mays, G. Starace, K. Liu, L. Maksin, T. Patwardhan, A. Madry, and L. Weng, "MLE-bench: Evaluating machine learning agents on machine learning engineering," in *Thirteenth International Conference on Learning Representations*, 2025.
- [17] L. Jing, Z. Huang, X. Wang, W. Yao, W. Yu, K. Ma, H. Zhang, X. Du, and D. Yu, "DSBench: How far are data science agents from becoming data science experts?" in *Thirteenth International Conference on Learning Representations*, 2025.
- [18] A. Alzarrad, I. Awolusi, M. T. Hatamleh, and S. Terreno, "Automatic assessment of roofs conditions using artificial intelligence (AI) and unmanned aerial vehicles (UAVs)," in *Frontiers in Built Environment*, 2022.
- [19] T. M. Brown, W. H. Pogorzelski, and I. M. Giammanco, "Evaluating hail damage using property insurance claims data," *Weather, Climate, and Society*, 2015.
- [20] J. Garrido, C. Genest, and J. Schulz, "Generalized linear models for dependent frequency and severity of insurance claims," *Insurance: Mathematics and Economics*, 2016.
- [21] E. W. Frees, R. A. Derrig, and G. Meyers, *Predictive Modeling in Actuarial Science*, ser. International Series on Actuarial Science. Cambridge University Press, 2014, p. 1–12.
- [22] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*, 2021.
- [23] The CAS Institute, "Study note: Model validation and holdout data," Online; accessed 2025-06-19, Casualty Actuarial Society Institute, Tech. Rep., 2019. [Online]. Available: <https://thecasinstitute.org/wp-content/uploads/2019/01/Exam-3-Study-Note-Model-Validation-01162019.pdf>
- [24] C. Ye, L. Zhang, M. Han, Y. Yu, B. Zhao, and Y. Yang, "Combining predictions of auto insurance claims," *Econometrics*, 2018.
- [25] T. Pijl, "A framework to forecast insurance claims," M.Sc. thesis, Erasmus University Rotterdam, Erasmus School of Economics, Rotterdam, Netherlands, Aug. 2017.
- [26] A. I. Company and J. Moser, "Allstate claim prediction challenge," 2011, kaggle. [Online]. Available: <https://kaggle.com/competitions/ClaimPredictionChallenge>