



MINISTRY OF ECONOMY
DEPARTMENT OF STATISTICS MALAYSIA

Healthcare Insurance Fraud Detection using Benford's Law

Wei Han YAP, Kee Huong LAI

School of Accounting and Finance,
Taylor's Business School, Taylor's University

**11th MALAYSIA
STATISTICS CONFERENCE**
"Data and Artificial Intelligence: Empowering the Future"

**19th September
2024**

Organized by:



OVERVIEW

1. INTRODUCTION
2. METHODOLOGY
3. RESULT
4. DISCUSSION
5. CONCLUSION
6. REFERENCES

INTRODUCTION

11th MALAYSIA STATISTICS CONFERENCE
"Data and Artificial Intelligence: Empowering the Future"

INTRODUCTION – Losses due to Fraud

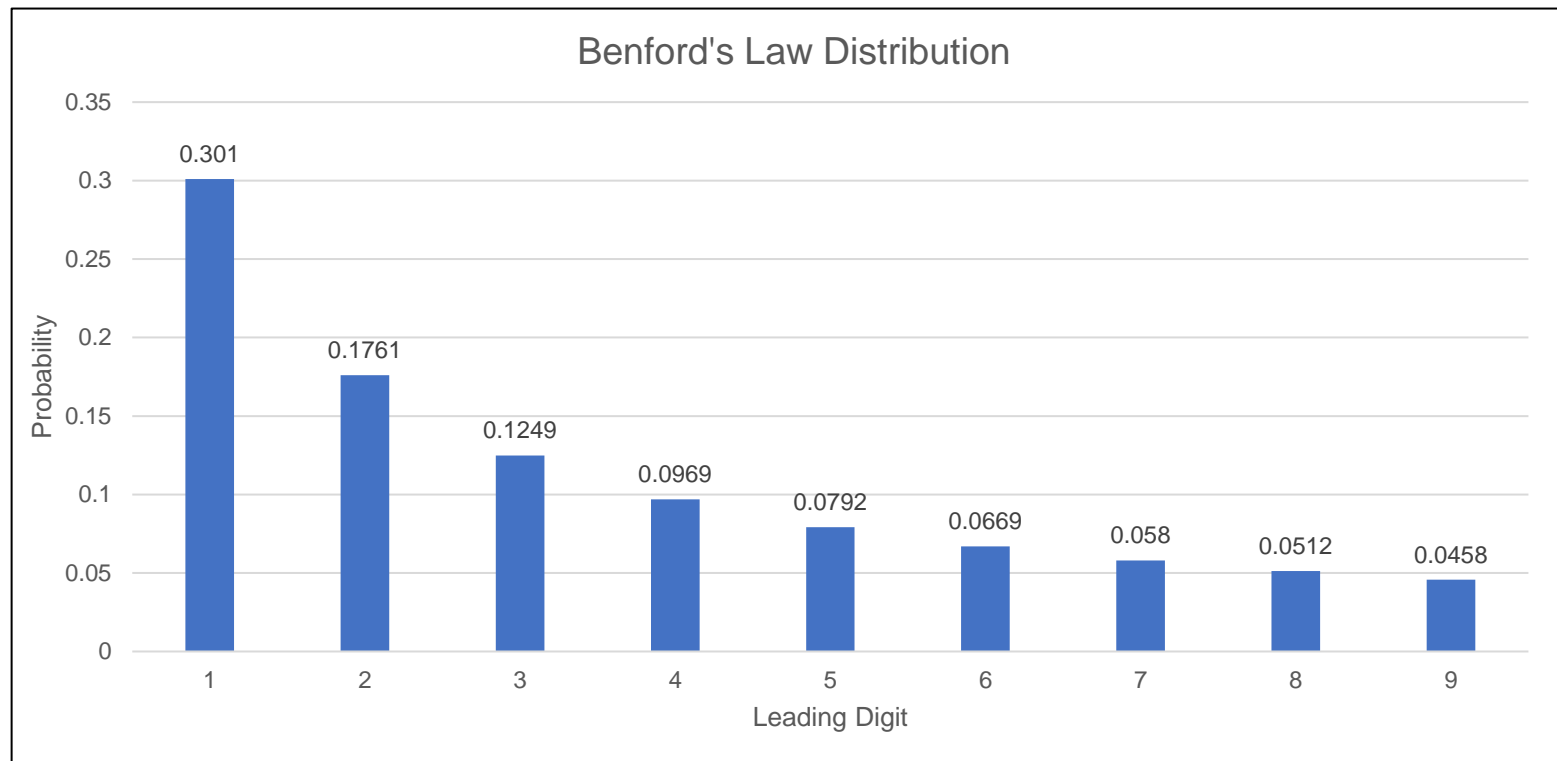
- **An alarming annual loss of USD 308.6 billion due to insurance fraud (Kilroy, 2024)**
- **Conservative pricing – transfers financial burden to policyholders – higher premiums (Chen et al., 2020)**
- **Automated fraud detection technologies – machine learning models (Nabrawi & Alanazi, 2023) – sophisticated and computationally expensive**

INTRODUCTION – Benford's Law

- A straightforward statistical tool
- Applied successfully in forensic accounting (Druica et al., 2018) and electoral fraud detection (Gueron & Pellegrini, 2022)
- Applies exclusively to naturally occurring numbers (insurance claim amounts, stock prices, etc)
- Not applicable to manipulated or pre-assigned numbers (phone numbers, aggregate claim amounts after the policy limit is applied)

INTRODUCTION – Benford's Law – Formula and Distribution

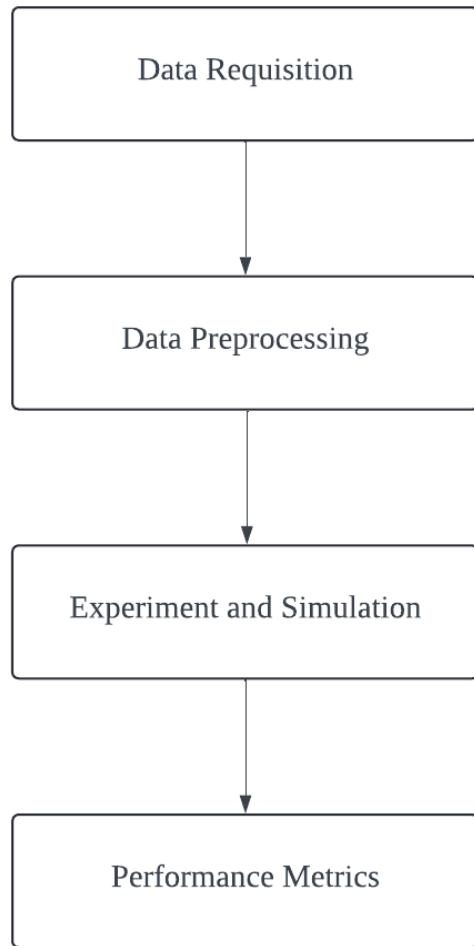
$$P(d) = \log_{10} \left(1 + \frac{1}{d} \right), \text{ where } d = \{1, 2, \dots, 9\}.$$



METHODOLOGY

11th MALAYSIA STATISTICS CONFERENCE
"Data and Artificial Intelligence: Empowering the Future"

METHODOLOGY - Flowchart



1. Data Requisition

- Synthetic dataset obtained from Kaggle
- A dataset with 63,968 observations focusing on annual reimbursement amounts for Medicare
- Research focus on IPAnnualReimbursementAmt, IPGrossClaim, OPAnnualReimbursementAmt, and OPGrossClaim

2. Data Preprocessing

- Utilised both Microsoft Excel and R-programming
- The first digit of the respective inpatient and outpatient reimbursement amount and gross claim were extracted

3. Experiment and Simulation

- The distribution of each category is fitted to Benford's Law distribution.
- Combo charts were generated for data visualization.

4) Performance Metrics

$$Z = \frac{|AP - EP| - \frac{1}{2N}}{\sqrt{\frac{EP(1-EP)}{N}}}$$

- **Z-test** is used to measure how many standard deviations the observed distribution of each respective digit are from Benford's Law distribution.

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

- **Chi-squared goodness-of-fit test** compares the overall conformity of the observed distribution to the expected distribution, is affected by the number of observations.

$$MAD = \frac{\sum_{i=1}^k |AP - EP|}{k}$$

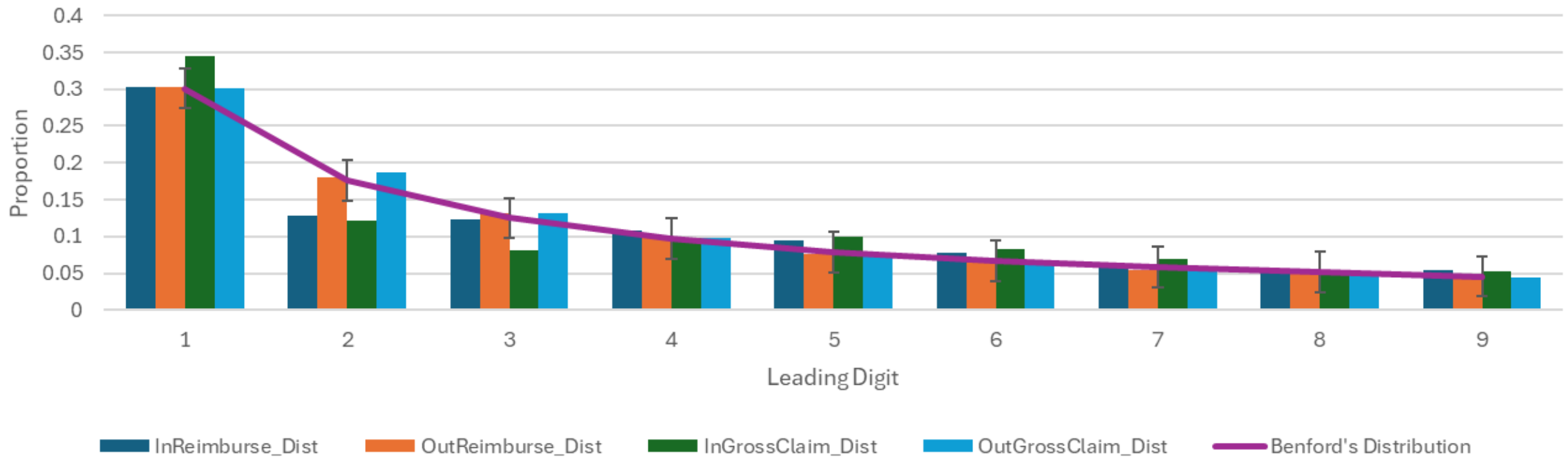
- **Mean Absolute Deviation (MAD)** compares the overall conformity of the observed distribution to Benford's Law distribution, not affected by the number of observations

RESULT

11th MALAYSIA STATISTICS CONFERENCE
"Data and Artificial Intelligence: Empowering the Future"

RESULT – Comparison of All Categories

Comparison of all categories actual distribution with Benford's Law distribution



RESULT – Inpatient Categories

Digit	Benford's law	Inpatient Reimbursement	z-stat	Inpatient Gross Claim	z-stat
1	0.3010	0.3024	0.4212	0.3455	13.7297*
2	0.1761	0.1280	17.7637*	0.1214	20.3182*
3	0.1249	0.1232	0.7501	0.0803	19.0966*
4	0.0969	0.1073	4.9091*	0.0980	0.5256
5	0.0792	0.0948	8.1215*	0.0990	10.3752*
6	0.0669	0.0771	5.7159*	0.0823	8.6848*
7	0.0580	0.0577	0.1366	0.0686	6.4413*
8	0.0512	0.0559	3.0317*	0.0517	0.3400
9	0.0458	0.0536	5.2667*	0.0531	4.9286*
χ^2 p-value		1.494e-83	χ^2 p-value	6.662e-216	* Z > 1.96
MAD		0.0111	MAD	0.0221	
Interpretation		Acceptable conformity	Interpretation	Non-conformity	

RESULT – Outpatient Categories

Digit	Benford's law	Outpatient Reimbursement	z-stat	Outpatient Gross Claim	z-stat	
1	0.3010	0.3029	0.9958	0.3011	0.0496	
2	0.1761	0.1796	2.3193*	0.1870	7.1578*	
3	0.1249	0.1317	5.1056*	0.1322	5.5429*	
4	0.0969	0.0978	0.7309	0.0984	1.2497	
5	0.0792	0.0754	3.4714*	0.0749	3.9737*	
6	0.0669	0.0653	1.6151	0.0619	5.0904*	
7	0.0580	0.0542	4.0552*	0.0533	5.0172*	
8	0.0512	0.0487	2.7934*	0.0475	4.1746*	
9	0.0458	0.0444	1.6193	0.0437	2.4504*	
χ^2 p-value		1.442e-11	χ^2 p-value		1.287e-29	* Z > 1.96
MAD		0.00288	MAD		0.00439	
Interpretation		Close conformity	Interpretation		Close conformity	

DISCUSSION

11th MALAYSIA STATISTICS CONFERENCE
"Data and Artificial Intelligence: Empowering the Future"

DISCUSSION – Contradictory Results

- **Inpatient gross claim – raw loss data**
- **Inpatient reimbursement amount – revised amount after deductible**
- **Intuitively speaking, the raw data would conform to Benford's law**
- **Contradicting results**

DISCUSSION – Excess Power Problem

- Large sample size – p -values obtained are small, while the test statistics have large values
- Excess power problem encountered by the chi-squared test (Kossovsky, 2021).
- Other statistical tests, such as the MAD, should be used to complement the chi-squared test to provide additional insights

CONCLUSION

11th MALAYSIA STATISTICS CONFERENCE
"Data and Artificial Intelligence: Empowering the Future"

CONCLUSION – Concluding Remarks

- **Inpatient reimbursement amount conforms to Benford's Law more than inpatient gross claim (raw data), which contradicts Benford's Law.**
- **Excess power problem faced by chi-square test and z-test**

CONCLUSION – Limitations

- **Lack of accessibility of real-world insurance datasets due to privacy matters.**
- **Deviations from Benford's Law do not necessarily imply fraudulent cases as Benford's Law only serves as a preliminary statistical tool.**
- **More advanced tests and algorithms should be employed for further investigation.**

CONCLUSION –Future Works

- **Using simulated datasets (Campo & Antonio, 2023) that are realistic and representative of actual insurance datasets.**
- **Complementing the results from Benford's Law with more advanced machine learning models to accurately identify fraudulent cases**
- **To examine insurance datasets that include co-payments**

REFERENCES

11th MALAYSIA STATISTICS CONFERENCE
"Data and Artificial Intelligence: Empowering the Future"

REFERENCES

1. Chen, Z. X., Hohmann, L., Banjara, B., Zhao, Y., Diggs, K., & Westrick, S. C. (2020). Recommendations to protect patients and health care practices from medicare and medicaid fraud. *Journal of the American Pharmacists Association*, 60(6), e60-e65.
2. Nabrawi, E., & Alanazi, A. (2023). Fraud detection in healthcare insurance claims using machine learning. *Risks*, 11(9), 160.
3. Gueron, E., & Pellegrini, J. (2022). Application of Benford–Newcomb law with base change to electoral fraud detection. *Physica A: Statistical Mechanics and its Applications*, 607, 128208.
4. Druică, E., Oancea, B., & Vâlsan, C. (2018). Benford's law and the limits of digit analysis. *International Journal of Accounting Information Systems*, 31, 75-82.
5. Gupta, R. A. (2019). Healthcare Provider Fraud Detection Analysis. Retrieved May 1, 2024, from <https://www.kaggle.com/datasets/rohitrox/healthcare-provider-fraud-detection-analysis/code>
6. Nigrini, M. J. (2012). *Benford's Law: Applications for forensic accounting, auditing, and fraud detection* (Vol. 586). John Wiley & Sons.
7. Kossovsky, A. E. (2021). On the mistaken use of the chi-square test in Benford's law. *Stats*, 4(2), 419-453.
8. Campo, B. & Antonio, K. (2023). Insurance fraud network data simulation machine: Generating synthetic fraud network data sets to develop and to evaluate insurance fraud detection strategies. In *Insurance Data Science Conference, Location: London, United Kingdom*.

Thank you

**11th MALAYSIA
STATISTICS CONFERENCE**
"Data and Artificial Intelligence: Empowering the Future"

**19th September
2024**

Organized by:

