

Integration of Traditional and Telematics data for Efficient Insurance Claims Prediction

Hashan Peiris

© Hashan Peiris 2023

SIMON FRASER UNIVERSITY

Spring 2023

Declaration of Committee

Name: Hashan Peiris
Degree: Master of Science
Thesis title: Integration of Traditional and Telematics data for Efficient Insurance Claims Prediction

Committee:

Chair:



Himchan Jeong



Gary Parker



Joan Hu



Abstract

The first part of the abstract discusses the background and objectives of the study. It highlights the importance of understanding the relationship between the variables being investigated and the need for a comprehensive analysis.

The second part of the abstract describes the methodology used in the study. It details the data collection process, the statistical models employed, and the steps taken to ensure the accuracy and reliability of the results.

Keywords: This section lists the key terms and concepts used throughout the study, providing a clear reference for researchers and readers interested in the topic.

Dedication

For my family and friends who have supported me throughout my life.

Acknowledgements

Table of Contents

Declaration of Committee	ii
Abstract	iii
Dedication	iv
Acknowledgements	v
Table of Contents	vi
List of Tables	viii
List of Figures	ix
1 Introduction	1

1.4




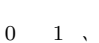
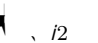




















































18.995 (elematics)-3322997 (in)-3P14 (.)-496 (Data)Descrip0.909-661TD85 (.)-500.004 (.)-

5 Data analysis		24
6 Conclusions		34
Bibliography		35
Appendix A Results		39
Appendix B Code		42
Appendix C Basic Setup of Proposed Method		49

List of Tables

Table 1	1
Table 2	2
Table 3	3
Table 4	4
Table 5	5
Table 6	6
Table 7	7
Table 8	8
Table 9	9
Table 10	10
Table 11	11
Table 12	12
Table 13	13
Table 14	14
Table 15	15
Table 16	16
Table 17	17
Table 18	18
Table 19	19
Table 20	20
Table 21	21
Table 22	22
Table 23	23
Table 24	24
Table 25	25
Table 26	26
Table 27	27
Table 28	28
Table 29	29
Table 30	30
Table 31	31
Table 32	32
Table 33	33
Table 34	34
Table 35	35
Table 36	36
Table 37	37
Table 38	38
Table 39	39
Table 40	40
Table 41	41
Table 42	42
Table 43	43
Table 44	44
Table 45	45
Table 46	46
Table 47	47
Table 48	48
Table 49	49
Table 50	50
Table 51	51
Table 52	52
Table 53	53
Table 54	54
Table 55	55
Table 56	56
Table 57	57
Table 58	58
Table 59	59
Table 60	60
Table 61	61
Table 62	62
Table 63	63
Table 64	64
Table 65	65
Table 66	66
Table 67	67
Table 68	68
Table 69	69
Table 70	70
Table 71	71
Table 72	72
Table 73	73
Table 74	74
Table 75	75
Table 76	76
Table 77	77
Table 78	78
Table 79	79
Table 80	80
Table 81	81
Table 82	82
Table 83	83
Table 84	84
Table 85	85
Table 86	86
Table 87	87
Table 88	88
Table 89	89
Table 90	90
Table 91	91
Table 92	92
Table 93	93
Table 94	94
Table 95	95
Table 96	96
Table 97	97
Table 98	98
Table 99	99
Table 100	100

List of Figures

Chapter 1

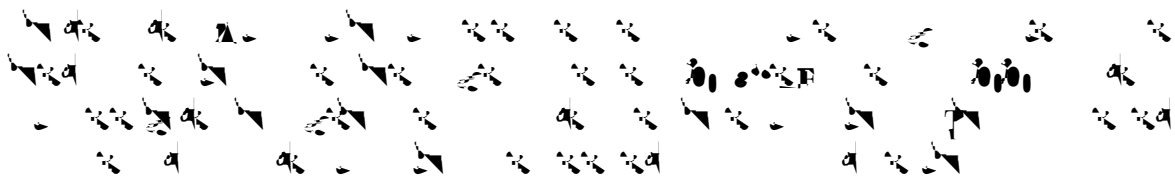
Introduction



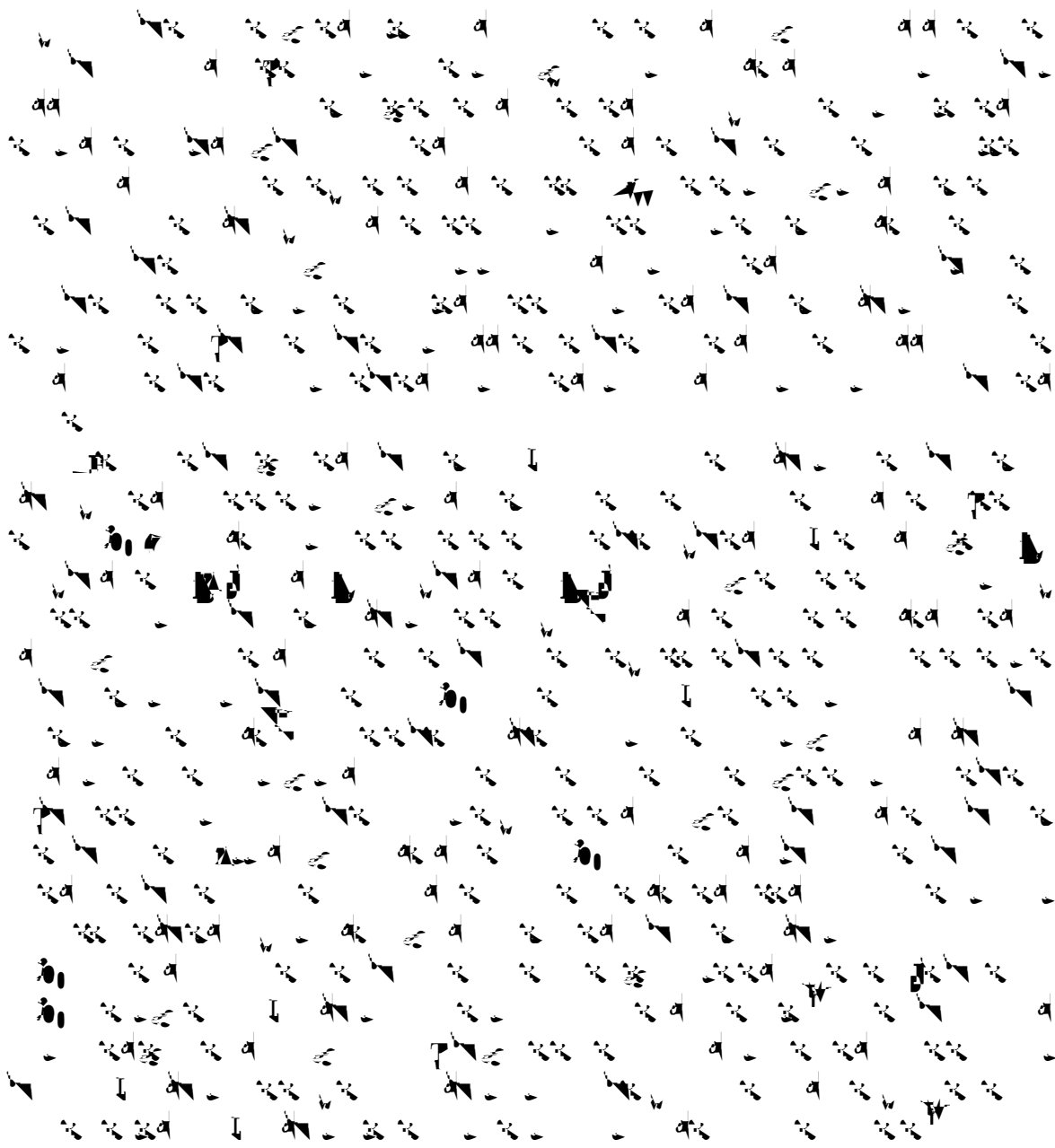
1.1 Modeling Claim Counts



1. 2019年12月31日，甲公司“应付账款”科目所属各明细科目期末贷方余额如下表所示。甲公司2019年12月31日资产负债表中“应付账款”项目期末余额为多少万元？

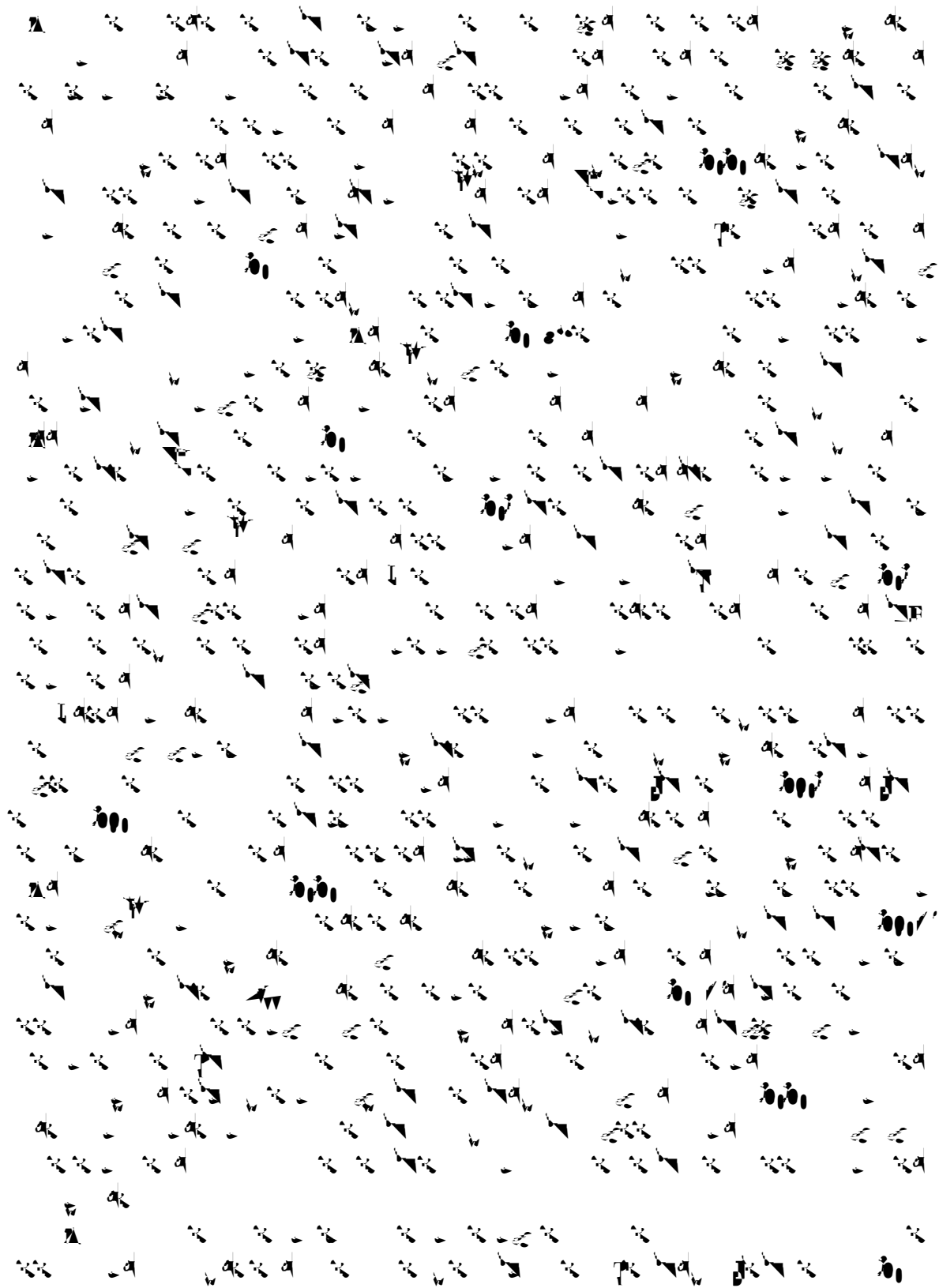


1.2 Telematics in Insurance



1.2.1 Uses of Telematics Data

1.2.2 Challenges



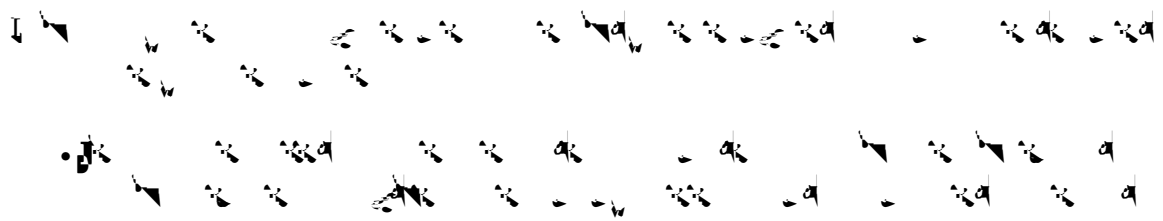


1.3 Motivation





1.4 Summary



Chapter 2

Data structure and problem description

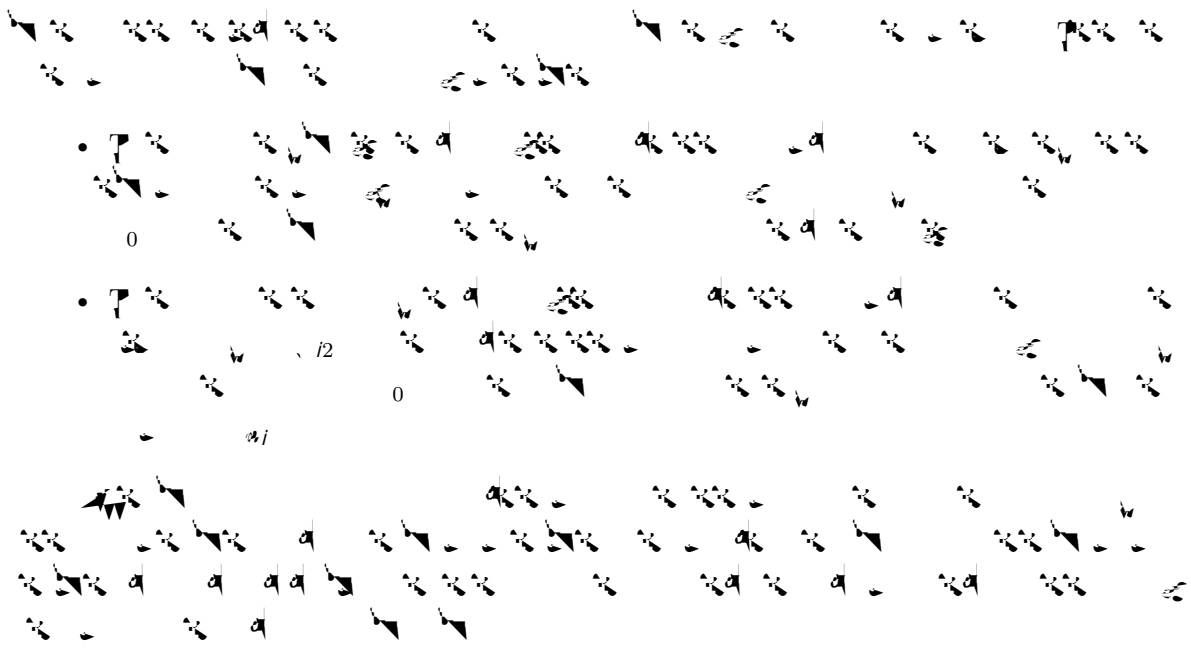
- i_1 is the i_1 -th element of a in the i -th iteration, $i = 0, 1, \dots$
- i_2 is the i_2 -th element of a in the i -th iteration, $i = 0, 1, \dots$
- $i = (i_1, i_2)$



$0, 1, i_1, a, i_2$

2.2 Problem Description

The diagram illustrates a sequence of operations on a data structure. It features a large light blue rectangle labeled 'M' at the top, a smaller blue rectangle below it, and a white rectangle with a red border at the bottom. A small blue square labeled 'S' is positioned to the right of the white rectangle. A pink 'M' is visible in the bottom right corner of the diagram area.



Chapter 3

Methodology





$$\sum_{i=1}^M a_i - \phi(\cdot) \nabla(\cdot) = 0$$

(.)

$\mu_i = (\mu_{i1}, \mu_{i2})$

3.2 Proposed Method

$\mu_i = (\mu_{i1}, \mu_{i2})$

3.2.1 Estimation of Parameters

$\mu_i = (\mu_{i1}, \mu_{i2})$

$$f(x; \mu) = \frac{1}{\sigma} \exp\left(-\frac{x - \mu}{\sigma}\right)$$

$\mu_i = (\mu_{i1}, \mu_{i2})$

$$i_{S_0} [1^i 1^i \cdots 1^L] = \sum_{i=1}^M [1^i 1^i \cdots 1^L]$$

$i_{S_0} [1^i 1^i \cdots 1^L] = [1^i x_{i1} \cdots 1^L] = [1^i x_{i1} \cdots 1^L] = 2^* + 1$

$$i_{S_0} (i; i^i) = \sum_{i=1}^M i (i; i^i) + \sum_{i=1}^M (1 - i) \binom{L}{k=0} i^k$$

$$= \sum_{i=1}^M i (i; i^i) + (1 - i) \binom{L}{k=0} i^k$$

$$+ \sum_{i=1}^M i(i-1) \binom{L}{k=0} i^k$$

$$w = (w_0, w_1, \dots, w_L)$$

$$\sum_{i=1}^M i(i-1) \binom{L}{k=0} i^k = 0$$

$$i_{S_0} (i; i^i) = \sum_{i=1}^M i (i; i^i) + (1 - i) \binom{L}{k=0} i^k$$

$$\theta = (\theta_0, \theta_1, \dots, \theta_L)$$

$$\hat{\theta}_i = 1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1i} + \dots + \hat{\theta}_L x_{Li})$$

$$\hat{\theta}_i = \frac{1}{\sum_{j=1}^S \hat{\theta}_j} \left(\sum_{j=1}^S \hat{\theta}_j \right)$$

$$\hat{\theta}_i(\hat{\theta}) = \hat{\theta}_i(\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_L) = 0$$

$$\hat{\theta}_i(\hat{\theta}) = \frac{1}{\sum_{j=1}^S \hat{\theta}_j} \left(\sum_{j=1}^S \hat{\theta}_j \right) = \frac{1}{\sum_{j=1}^S \left(1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1j} + \dots + \hat{\theta}_L x_{Lj}) \right)}$$

$$\hat{\theta}_i(\hat{\theta}) = \frac{1}{\sum_{j=1}^S \left(1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1j} + \dots + \hat{\theta}_L x_{Lj}) \right)}$$

3.2.2 Standard Errors of Estimates

$$\hat{\theta}_i(\hat{\theta}) = \frac{1}{\sum_{j=1}^S \left(1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1j} + \dots + \hat{\theta}_L x_{Lj}) \right)}$$

$$\hat{\theta}_i(\hat{\theta}) = \frac{1}{\sum_{j=1}^S \left(1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1j} + \dots + \hat{\theta}_L x_{Lj}) \right)}$$

$$\hat{\theta}_1(\hat{\theta}) = \frac{\hat{\theta}_1}{\sum_{j=1}^S \left(1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1j} + \dots + \hat{\theta}_L x_{Lj}) \right)}$$

$$\hat{\theta}_2(\hat{\theta}) = \frac{\hat{\theta}_2}{\sum_{j=1}^S \left(1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1j} + \dots + \hat{\theta}_L x_{Lj}) \right)}$$

$$\hat{\theta}_i(\hat{\theta}) = \frac{\hat{\theta}_i}{\sum_{j=1}^S \left(1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1j} + \dots + \hat{\theta}_L x_{Lj}) \right)}$$

$$\hat{\theta}_i(\hat{\theta}) = 1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1i} + \dots + \hat{\theta}_L x_{Li})$$

$$\hat{\theta}_i(\hat{\theta}) = \frac{1}{\sum_{j=1}^S \left(1 + \frac{1}{0} \exp(\hat{\theta}_0 + \hat{\theta}_1 x_{1j} + \dots + \hat{\theta}_L x_{Lj}) \right)}$$

$$\hat{\theta}_1(\hat{\theta}) = 0 \quad \hat{\theta}_2(\hat{\theta}) = 0$$

$$\hat{\gamma} = (\hat{\gamma}_1, \hat{\gamma}_2)$$

$$\hat{\gamma}(\hat{\theta}) = \begin{pmatrix} \hat{\gamma}_1(\hat{\theta}) \\ \hat{\gamma}_2(\hat{\theta}) \end{pmatrix}$$

$$\hat{\gamma}(\hat{\theta}) = \hat{\gamma}^{-1}(\hat{\theta}) \hat{\gamma}^{-1'}$$

$$\hat{\gamma}(\hat{\theta}) = \hat{\gamma}^{-1}(\hat{\theta}) \hat{\gamma}^{-1'}$$

$$\hat{\gamma}(\hat{\theta}) = \hat{\gamma}^{-1}(\hat{\theta}) \hat{\gamma}^{-1'}$$

$$\hat{\gamma}(\hat{\theta}) = \hat{\gamma}^{-1}(\hat{\theta}) \hat{\gamma}^{-1'}$$

$$\hat{\gamma}(\hat{\theta}) = \hat{\gamma}^{-1}(\hat{\theta}) \hat{\gamma}^{-1'}$$

$$\hat{\gamma}_i = \frac{1}{i} \sum_{j=1}^i \hat{\gamma}_j(\hat{\theta}) \quad \hat{\gamma}_i = \frac{1}{i} \sum_{j=1}^i \hat{\gamma}_j(\hat{\theta})$$

3.3 Estimation scheme

$$\mathcal{H} = \{ \hat{\gamma}_1(\hat{\theta}_1), \dots, \hat{\gamma}_L(\hat{\theta}_L) \}$$

$$\mathcal{H} = \{ \hat{\gamma}_1(\hat{\theta}_1), \dots, \hat{\gamma}_L(\hat{\theta}_L) \}$$

Chapter 4

Simulation study

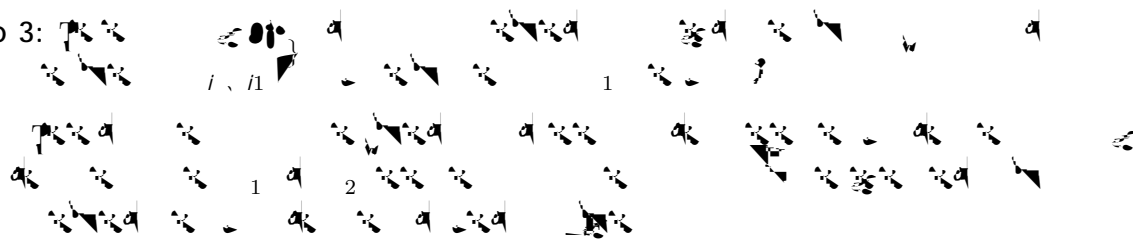


$\frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV + \frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV$

- $\mathcal{I}A$ $\frac{d}{dt} \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV + \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV$
- $\mathcal{I}G$ $\frac{d}{dt} \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV + \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV$
- $\mathcal{I}T$ $\frac{d}{dt} \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV + \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV$

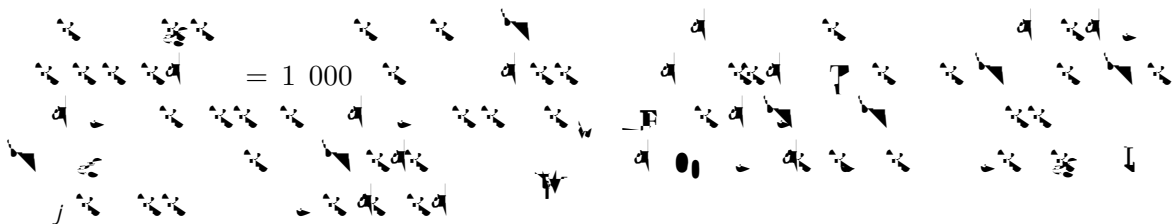
$\frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV + \frac{1}{2} \frac{d}{dt} \int_{\Omega} \rho \mathbf{v} \cdot \mathbf{v} \, dV$

Step 3:



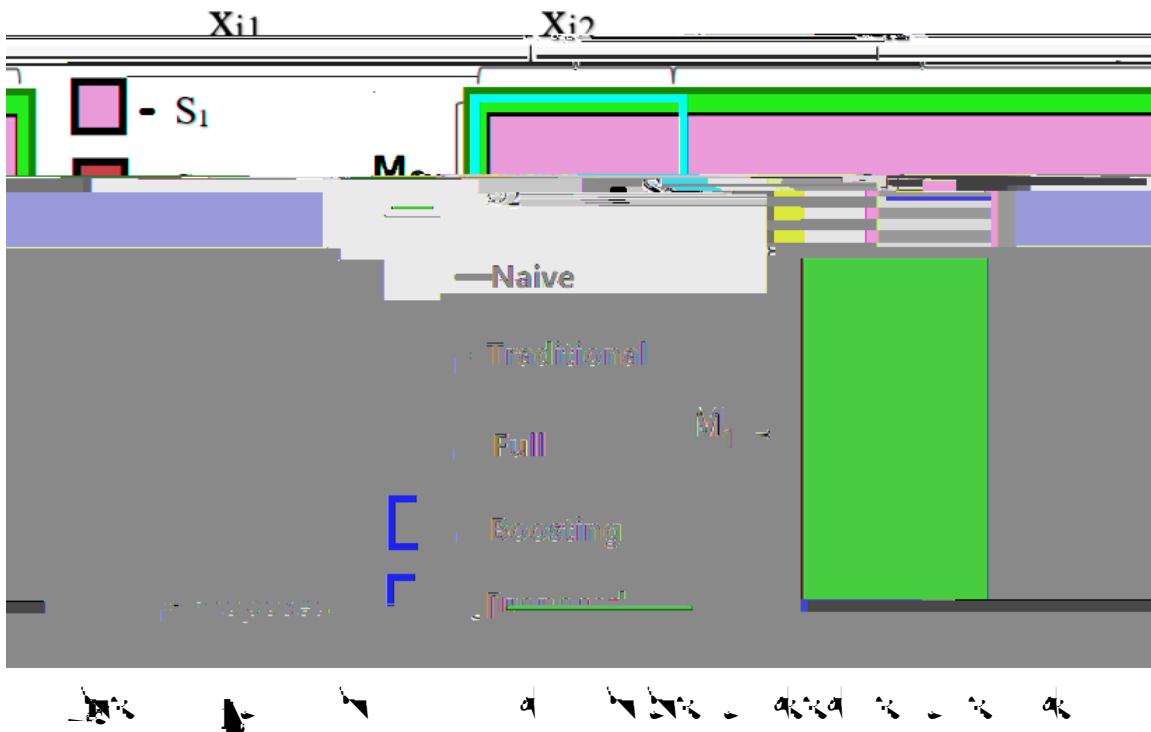
- Naive model $\mathbb{P}(i, i_1)$
- Traditional model $\mathbb{P}(i, i_1) \rightarrow \theta H_1$
- Full model $\mathbb{P}(i, i_1, j, w, a)$
- Boosting model $\eta_j = \exp(\hat{h}_j)$
 $\log \eta_j = \hat{h}_j$
- Proposed model \hat{h}_j *ith*

4.3 Evaluation Procedure



$$j = \frac{1}{r-1} \sum_{r=1}^R (j - \hat{j}^{(r)})$$

$$j = \frac{1}{r-1} \sum_{r=1}^R (j - \hat{j}^{(r)})^2$$



$$I_j = \frac{1}{R} \mathbb{1}_{\left\{ \left| \hat{j}_j^{(r)} - \hat{j}_j^{(r)} \right| < 1 \cdot \text{SE}(\hat{j}_j^{(r)}) \right\}}$$

$\hat{j}_j^{(r)}$ is the j th r th $(\hat{j}_j^{(r)})$

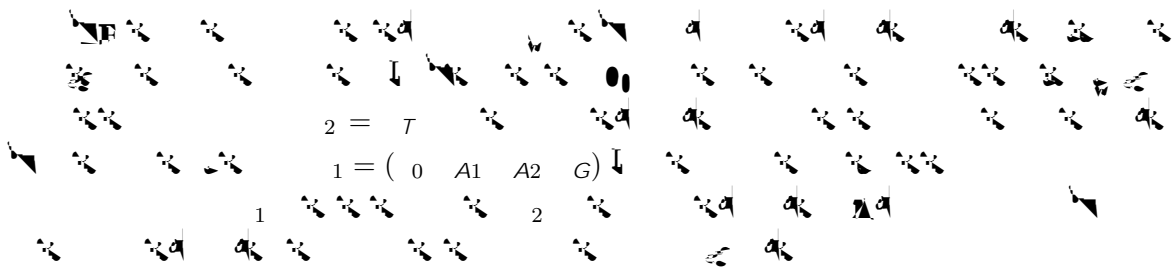
$\hat{j}_j^{(r)}$

$r = 1, \dots, R$

	$i S_0 (; i a_i) = 0$	$\hat{i} = \exp(, i1 \hat{1} +, i2 \hat{2})$
	$i S (1; i1 a_i) = 0$	$\hat{i} = \exp(, i1 \hat{1})$
	$i S^* (; i a_i) = 0$	$\hat{i} = \exp(, i1 \hat{1} +, i2 \hat{2})$
	$i S (1; i1 a_i) = 0$	$\hat{i} = \eta \hat{i} \exp(, i2 \hat{2})$
	$M_0 \left. \begin{aligned} a_{i-1} - \eta \hat{i} \exp(, i2 \hat{2}) \end{aligned} \right\} i2 = 0$	$\eta \hat{i} = \exp(, i1 \hat{1})$
	$\eta \hat{i} = \exp(, i1 \hat{1}) \quad b = 0$	$a = b$
	$i S i \hat{i} () (; i a_i) = 0$	$\hat{i} = \exp(, i1 \hat{1} +, i2 \hat{2})$
	$\hat{i} ()$	
	$a_i = \mathbb{I}(b \in 0)$	

4.4 Results





						\mathbb{W}					\mathbb{I}						
	N	T	B	F	P	N	T	B	F	P	N	T	B	F	P		
Random selection																	
0	0 0 0 1	0 1 1	0 1 1	0 0 0 1	0 0 0 1	0 1 0 1 z ⁰ 0 1 0 1 0 0 0 0 0 0 1 z ⁰	0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1 z ⁰	
A1	0 0 1	0 0 1	0 0 1	0 0 1	0 0 1	0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1	
A2	0 0 1	0 0 1	0 0 1	0 0 1	0 0 1	0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1	
G	0 0 0 1	0 0 0 1	0 0 0 1	0 0 0 1	0 0 0 1	0 1 0 1 0 1 0 1 0 1 0 1	0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1
T	0 0 0 0		0 0 1	0 0 0 0	0 0 0 0	0 1 0 1		0 1	0 1 0 1	0 1 0 1	0 1 z ⁰		0 1 z ⁰	0 1 z ⁰	0 1 z ⁰		
Age selection																	
a ₀	0 0 1	0 1 1	0 1 1	0 0 0 1	0 0 0 1	0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1	0 1	0 1	0 1 z ⁰	0 1 z ⁰		
a ₁	0 0 1 0	0 0 1	0 0 1	0 0 1	0 0 1	0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1	
A2	0 0 1	0 0 1	0 0 1	0 0 1	0 0 1	0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1	
G	0 0 0 1	0 0 0 1	0 0 0 1	0 0 0 1	0 0 0 1	0 1 0 1 0 1 0 1 0 1 0 1	0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1
T	0 0 0 1		0 0 1 z ⁰	0 0 0 0	0 0 0 1	0 1 0 1 z ⁰		0 1	0 1 0 1	0 1 0 1	0 1 0 1		0 1 0 1	0 1 z ⁰	0 1 0 1		
Adverse selection																	
0		0 1 1	0 1 1	0 0 0 1	0 0 0 1 z ⁰	z ⁰ 0 1 0 1 0 1 0 1 0 1 0 1	0 1 z ⁰	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1	0 1	0 1	0 1 z ⁰	0 1 z ⁰	
A1	0 1 z ⁰	0 0 1	0 0 1	0 0 1	0 0 1 z ⁰	0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1	
A2	0 1 z ⁰	0 0 1	0 0 1	0 0 1	0 0 1 z ⁰	0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1	
G	0 0 0 1 z ⁰	0 0 0 1	0 0 0 1	0 0 0 1	0 0 0 1 z ⁰	0 1 z ⁰ 0 1 0 1 0 1 0 1 0 1	0 1	0 1	0 1	0 1	0 1	0 1 z ⁰ 0 1 z ⁰ 0 1 z ⁰ 0 1	0 1	0 1	0 1	0 1	0 1
T	0 0 0 1		0 1	0 0 0 0	0 0 1	0 1		0 1	0 1 0 1	0 1 0 1	0 1 0 1		0 0 0 0	0 1 z ⁰	0 1 z ⁰		



	Durati on	
	Insured. age	
	Insured. sex	
	Car. age	
	Mari tal	
	Car. use	
	Credi t. score	
	Regi on	
	Annual . mi les. dri ve	
	Years. nocl ai ms	
	Terri toryEmb	
	Annual . pct. dri ven	
	Total . mi les. dri ven	
	Pct. dri ve. xxx	
	Pct. dri ve. rush. am	
	Pct. dri ve. rush. pm	
	Avgdays. week	
	Accel . 06mi les	
	Brake. 06mi les	
	Acbr. others	
	Left. turns	
	Right. turns	
	NB_Cl ai m	

5.2 Estimation and Evaluation



5.2.1 Estimation

- Random selection
- Age selection $1/(1 + \exp(0.031 \text{ insured. age}_i))$
- Adverse selection $1/(1 + \exp(\text{NB_CI ai m}_i))$

5.2.2 Evaluation

$$j = \frac{1}{r-1} \sum_{r=1}^R (\hat{\theta}_r - \hat{\theta}_j)$$

r56j ET qv7 7. TD5.5

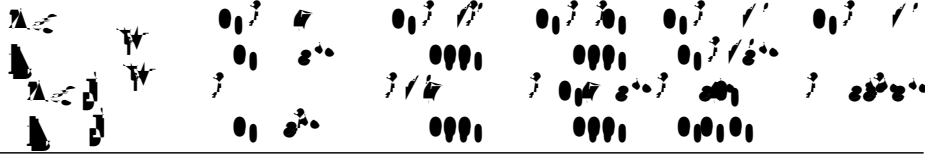
Handwritten musical notation on a page with five systems of staves. Each system begins with a treble clef and a common time signature (C). The notation consists of rhythmic stems, beams, and note heads, with some systems including a single note head. The notation is dense and covers most of the page.

0 1 2 3 4 5 6 7 8 9

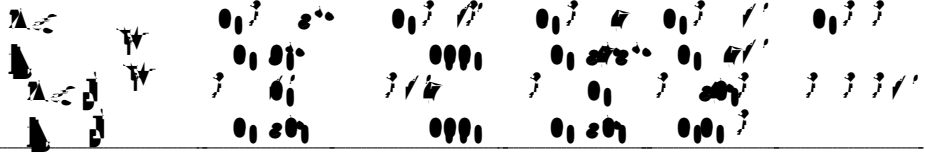
Random selection

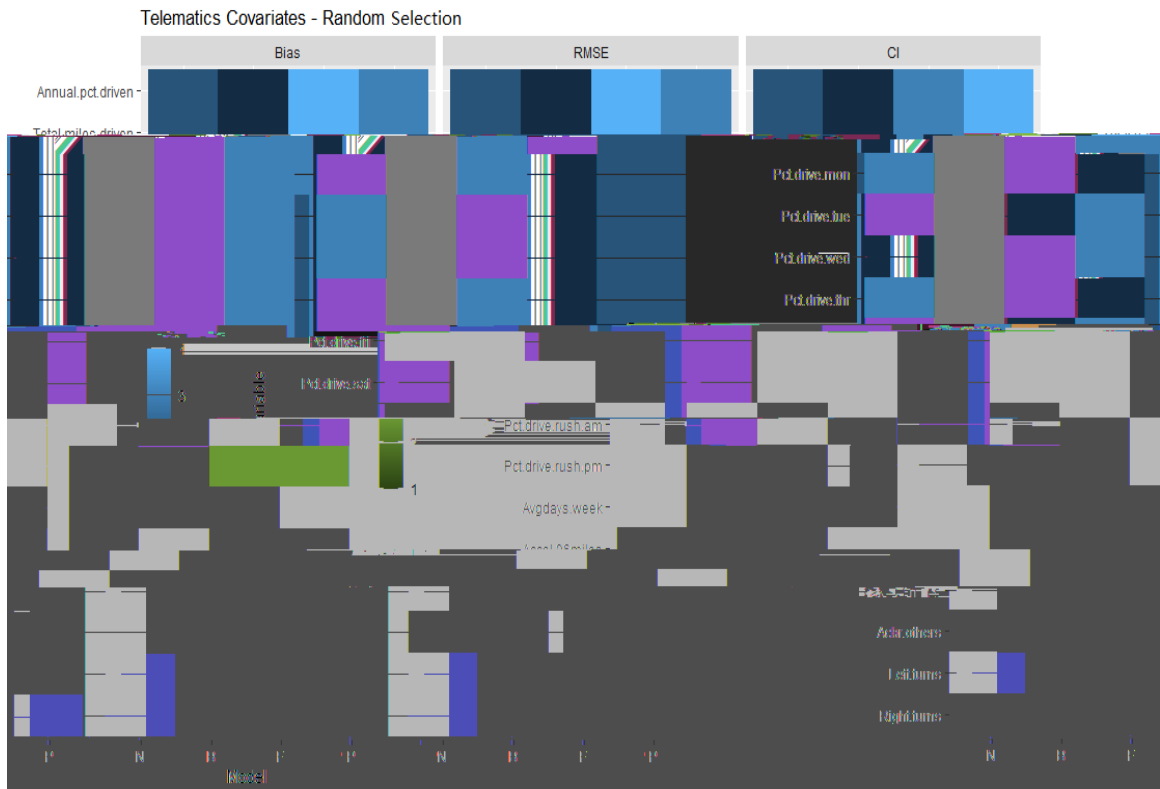


Age selection

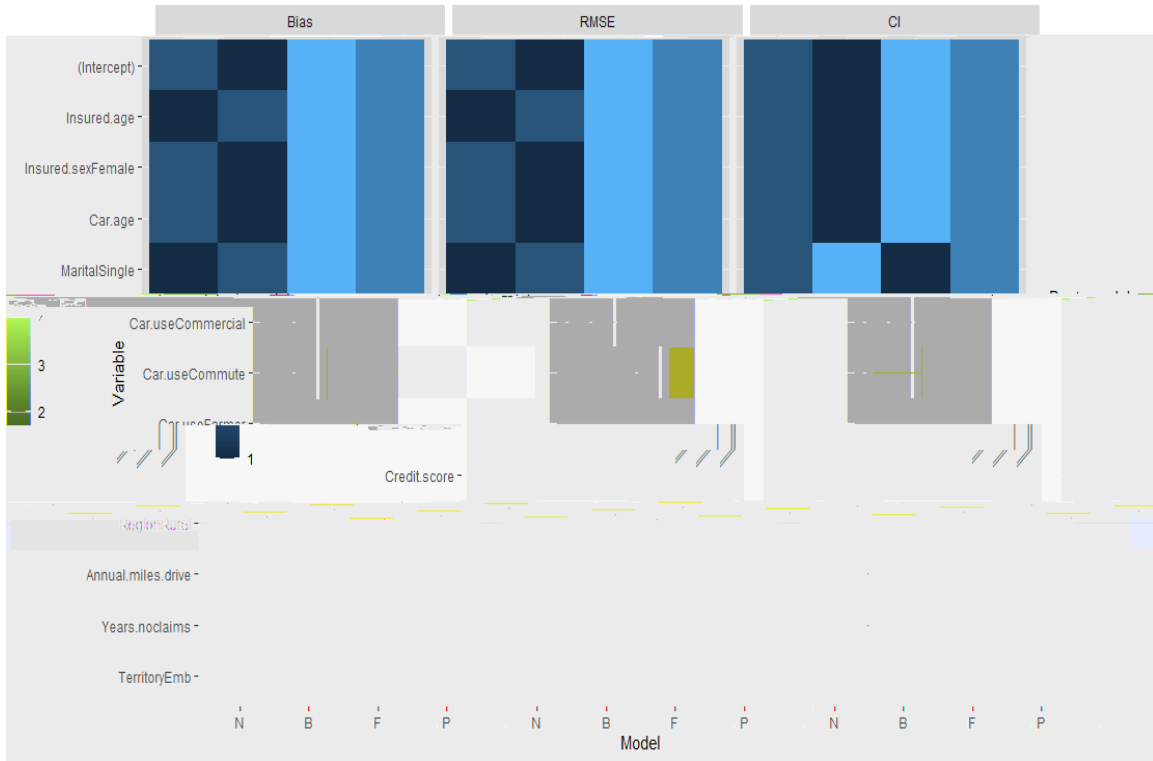


Adverse selection

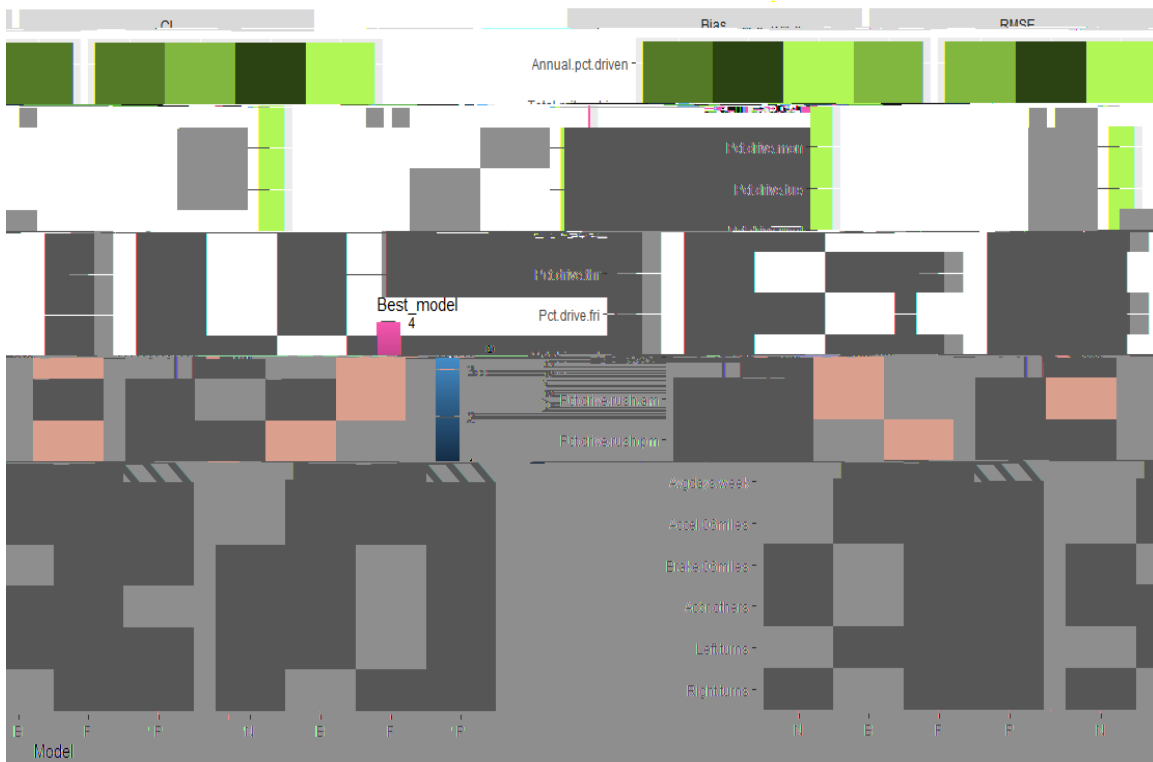




Traditional Covariates - Age Selection



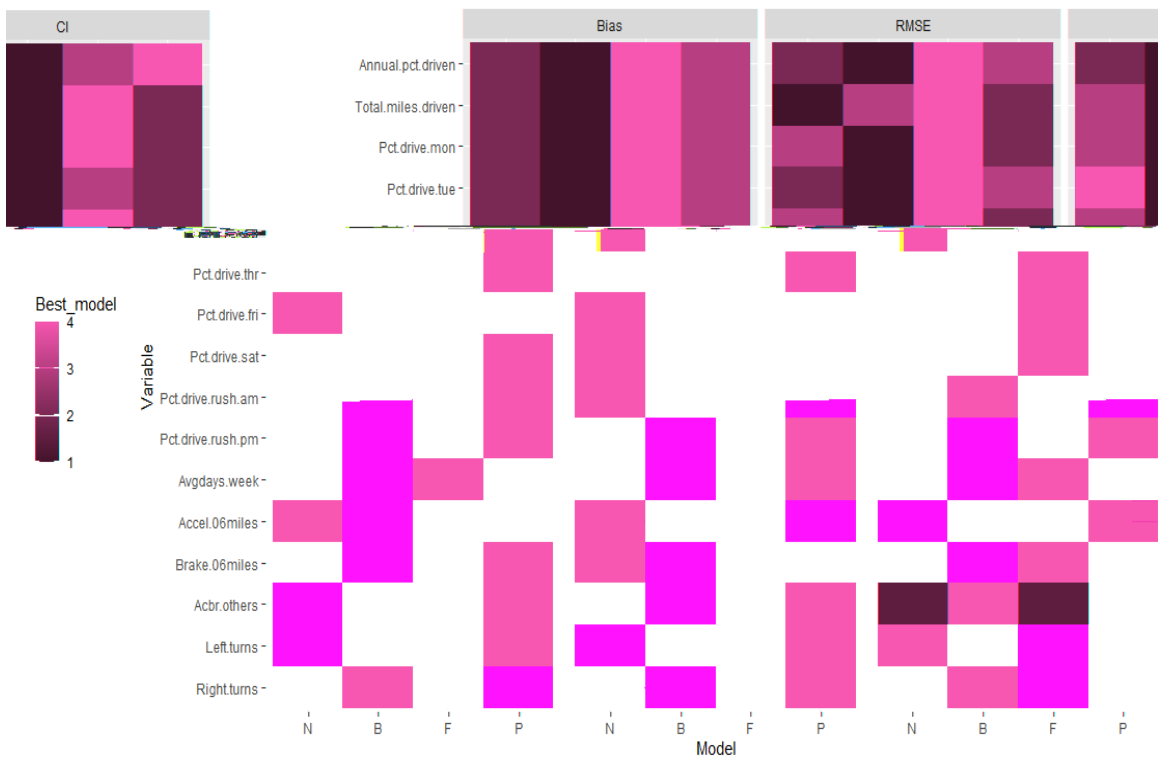
Telematics Covariates - Age Selection



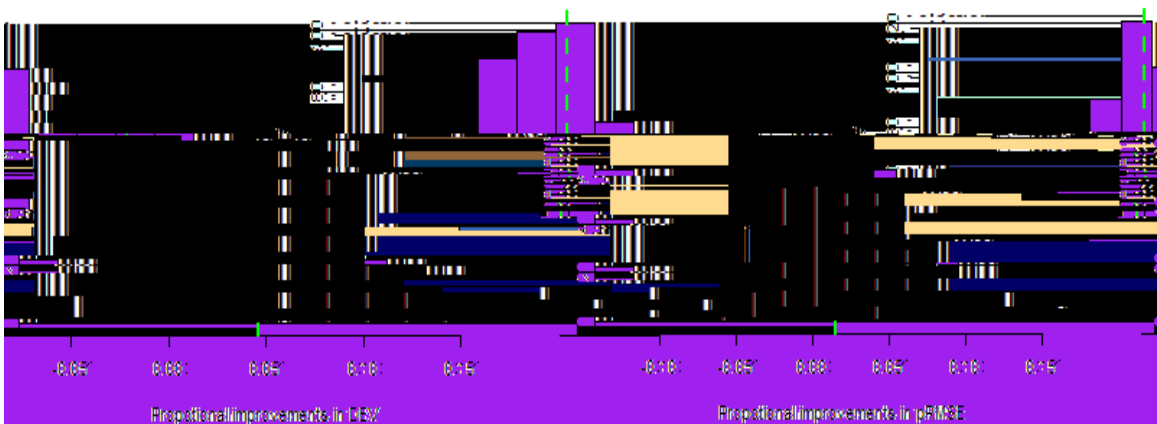
Traditional Covariates - Adverse Selection



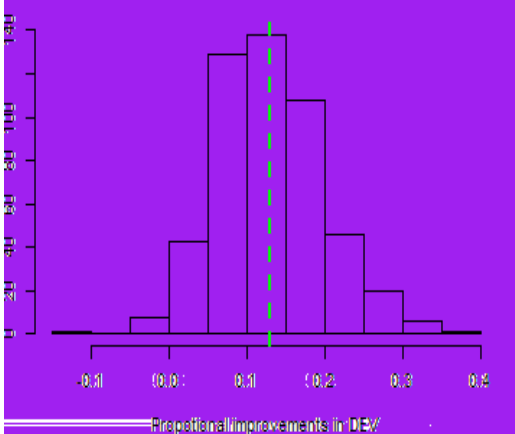
Telomatics Covariates - Adverse Selection



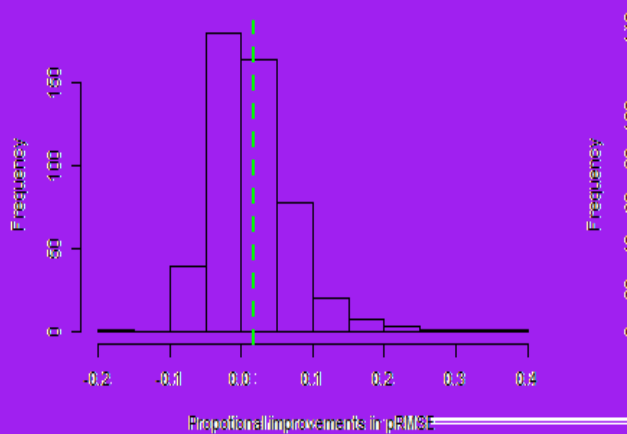
Proposed vs Naive (in percentage) from Random Selection



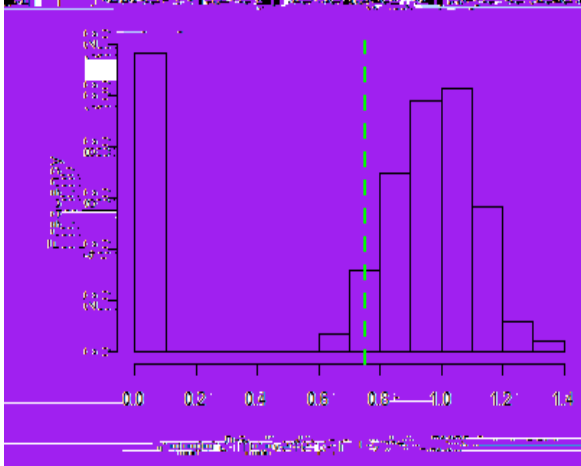
Proposed vs Naive (in percentage) from Age Selection



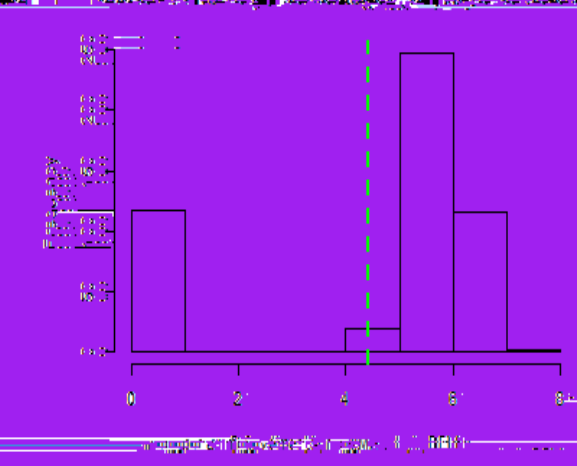
Proposed vs Naive (in percentage) from Age Selection



Proposed vs Naive (in percentage) from Age Selection



Proposed vs Naive (in percentage) from Age Selection



Bibliography

Categorical data analysis

Accident Analysis and Prevention

Risks

Transportation

Decision Support Systems

Insurance: Mathematics and Economics

Accident Analysis & Prevention

IAA *ASTIN Bulletin: The Journal of the*

Risk analysis

Journal of Risk and Insurance

Journal of the Royal Statistical Society: Series D (The Statistician)

Scandinavian Actuarial Journal

Procedia Engineering

Journal of the Royal Statistical Society: Series B (Statistical Methodology)

Proceedings of the 6th ACM Symposium on Development and Analysis of Intelligent Vehicular Networks and Applications

Management

Journal of Risk

Available at SSRN 3251623

Variance

Mathematics and Economics

Insurance:

Winter 2011 Volume 2

Casualty Actuarial Society E-Forum,

Canadian Electronic Library

<https://policycommons.net/artifacts/1189025/distance-based-vehicli e-insurance/1742147/>

Transportation Research Part A: Policy and Practice

Big Data & Society

Sustainability

Risks

Biometrika

International Journal of Statistics and Systems

Journal of the operational research society

Risks

Accident Analysis and Prevention

Principles and Applications of Narrowband Internet of Things (NBIoT)

Journal of the Royal Statistical Society: Series C (Applied Statistics)

Big Data for Insurance Companies

arXiv e-prints

Appendix A

Results

W

↓

Appendix B

Code

ak   ak

```

x4      <- rnorm(1)
lambda <- exp(-1.3-4*x1 + 3.4*x2 + 0.1*x3 + 0.5*x4)
NB_Claim <- rpois(1, lambda)
Duration <- rep(1, 1)
fdata   <- as.data.frame(cbind(x1, x2, x3, x4, Duration, NB_Claim))

#for testing
set.seed(j+1000)
test_ind <- sample(1:nrow(fdata), 10000)
forr.data<-fdata[ test_ind,]
trtt.data<-fdata[-test_ind,]
#sampling- when using a specific sampling method, comment other two sampling sections.
#####RS
set.seed(j+2000)
tele_ind <- sample(1:90000, nrow(fdata)*0.1)
ntr <- length(tele_ind)
#####NIS(advsel)
#set.seed(j+2000)
#dz <- 1/(1+exp(2*trtt.data$NB_Claim))
#dz <- dz/mean(dz)/9
#dzz <- rbinom(90000, 1, dz)
#tele_ind <- (1:90000)*(dzz==1)
#rm(dz, dzz)
#tele_ind <- tele_ind[tele_ind!=0]
#ntr <- length(tele_ind)
#####MAR(agesel)
#set.seed(j+2000)
#dz <- 1/(1+exp(3*trtt.data$x1))
#dz <- dz/mean(dz)/9
#dzz <- rbinom(90000, 1, dz)
#tele_ind <- (1:90000)*(dzz==1)
#rm(dz, dzz)
#tele_ind <- tele_ind[tele_ind!=0]
#ntr <- length(tele_ind)
#####
S0 <- trtt.data[ tele_ind, ]
# A small dataset that contains both telematics and traditional features
S1 <- trtt.data[-tele_ind, -tele_ind, -tele_ind, ]
T*Td(<-)Tj 0g14. 1220Td[(trtt.data[-525.004(tele_ind,)]Tj150220Td[(trt. 12214_in9D4]Tj5atetel emati cs)-

```

```

b_S0 <- as.matrix(cbind(S0[,c(5, 1:3)], S0[,6]*S0[,c(5, 1:3)]))
b_S  <- as.matrix(cbind(S[,c(5, 1:3)], S[,6]*S[,c(5, 1:3)]))

#function for optimize using nleqslv()
cal_eqn <- function(parm) {
  result <- colSums(as.vector(1+nrow(S1)/nrow(S0)*exp(parm%%t(b_S0)))*b_S0)-colSums(b_S)
  return(result) }

#find for parameters of basis functions
fit2 <- nleqslv(rep(0,8), cal_eqn)

#calculate weights from information projection
w3 <- 1+nrow(S1)/nrow(S0)*exp(b_S0 %% fit2$x)
#####
#combine weights to S0
SS6<-cbind(S0,w3)

#fitted the model with ws
glm.freq.S3 <- glm(NB_C1aim ~ .-Duration-w3, offset=log(Duration),
  weights= w3, data=SS6, family=poisson())
x_S0 <- model.matrix(glm.freq.S3)

#coef of proposed model
prop2_coef[, ] <- summary(glm.freq.S3)$coefficients[, 1]

# sandwich formula for variance estimation
Ui <- cbind(c(as.vector(w3)-1, rep(-1, nrow(S1))) * b_S,
  c(w3*(SS6$NB_C1aim-fitted(glm.freq.S3)), rep(0, nrow(S1))) * as.matrix(S[,c(5, 1:4)]))

Ui <- Ui - rep(colMeans(Ui), each=nrow(S))
V_U <- ( t(Ui) %% Ui)
tau <- rbind(cbind(t(b_S0) %% (as.vector(w3-1)*b_S0),
  matrix(0, ncol=ncol(x_S0), nrow=ncol(b_S0)),
  cbind(t(x_S0) %% (as.vector(w3-1)*(SS6$NB_C1aim-fitted(glm.freq.S3))*b_S0),
  -t(x_S0) %% (as.vector(w3*fitted(glm.freq.S3))*x_S0) ))
invtau <- solve(tau)
prop2_stde[, ] <- sqrt(diag(invtau %% V_U %% invtau))[-(1:ncol(b_S0))]
#####boosting model
glm.freq.boost <- glm(NB_C1aim ~ x4-1, data=S0,
  offset=log(Duration)+predict(glm.freq.trad, S0), family=poisson())
#coefficients and SE
boost_coef[, ] <- c(summary(glm.freq.trad)$coefficients[, 1],
  summary(glm.freq.boost)$coefficients[, 1])
boost_stde[, ] <- c(summary(glm.freq.trad)$coefficients[, 2],
  summary(glm.freq.boost)$coefficients[, 2])
#####try forecast
pred.nai ve <- predict(glm.freq.nai ve, newdata = forr.data, type="response")
pred.full <- predict(glm.freq.full , newdata = forr.data, type="response")
pred.S3 <- exp(as.matrix(forr.data[,c(5, 1:4)]) %% prop2_coef[, ])
pred.trad <- predict(glm.freq.trad , newdata = forr.data, type="response")
pred.boost <- pred.trad * exp(coef(glm.freq.boost)*forr.data$x4)

#remove datasets for this split
rm(tele_ind, test_ind)

#RMSE
RMSEs[, -4] <- sqrt(c(
  mean((forr.data$NB_C1aim-pred.nai ve)^2),

```

```

    mean((forr.data$NB_Claim-pred.trad )^2),
    mean((forr.data$NB_Claim-pred.full )^2),
    mean((forr.data$NB_Claim-pred.S3 )^2),
    mean((forr.data$NB_Claim-pred.boost)^2)))

#MAE
MAEs[j, -4] <- c(
  mean(abs(forr.data$NB_Claim-pred.nai ve)),
  mean(abs(forr.data$NB_Claim-pred.trad )),
  mean(abs(forr.data$NB_Claim-pred.full )),
  mean(abs(forr.data$NB_Claim-pred.S3 )),
  mean(abs(forr.data$NB_Claim-pred.boost)))

#DEV
DEVs[j, -4] <- c(
  Poi sson.Devi ance(pred.nai ve, forr.data$NB_Claim),
  Poi sson.Devi ance(pred.trad , forr.data$NB_Claim),
  Poi sson.Devi ance(pred.full , forr.data$NB_Claim),
  Poi sson.Devi ance(pred.S3 , forr.data$NB_Claim),
  Poi sson.Devi ance(pred.boost, forr.data$NB_Claim))

})

#summarizing the outputs
col Means(RMSEs)
col Means(MAEs)
col Means(DEVs)

#true coefficients used for data generation
true_coef <- c(-1.3, -4, 3.4, 0.1, 0.5)

#bias of each estimator
bi as_nai ve <- true_coef - col Means(nai ve_coef)
bi as_trad <- true_coef - col Means(trad_coef)
bi as_prop2 <- true_coef - col Means(prop2_coef) #proposed
bi as_boost <- true_coef - col Means(boost_coef)
bi as_full <- true_coef - col Means(full _coef)

#RMSE of estimates
rmse_nai ve <- sqrt(col Means((nai ve_coef-rep(true_coef, each=J))^2))
rmse_trad <- sqrt(col Means((trad_coef -rep(true_coef, each=J))^2))
rmse_prop2 <- sqrt(col Means((prop2_coef-rep(true_coef, each=J))^2)) #proposed
rmse_boost <- sqrt(col Means((boost_coef-rep(true_coef, each=J))^2))
rmse_full <- sqrt(col Means((full _coef -rep(true_coef, each=J))^2))

#CI of estimator
nai ve_90CI <- col Means((nai ve_coef-1.645*nai ve_stde<rep(true_coef, each=J))^*
  (nai ve_coef+1.645*nai ve_stde>rep(true_coef, each=J))^*1)
trad_90CI <- col Means((trad_coef -1.645*trad_stde<rep(true_coef, each=J))^*
  (trad_coef +1.645*trad_stde>rep(true_coef, each=J))^*1)
prop2_90CI <- col Means((prop2_coef-1.645*prop2_stde<rep(true_coef, each=J))^*
  (prop2_coef+1.645*prop2_stde>rep(true_coef, each=J))^*1,
  na.rm=TRUE) #proposed
boost_90CI <- col Means((boost_coef-1.645*boost_stde<rep(true_coef, each=J))^*
  (boost_coef+1.645*boost_stde>rep(true_coef, each=J))^*1)
full _90CI <- col Means((full _coef -1.645*full _stde<rep(true_coef, each=J))^*
  (full _coef +1.645*full _stde>rep(true_coef, each=J))^*1)

```



```
#### Preliminary analysis ####  
#naive model  
glm.freq.nai ve <- glm(NB_Claim ~ . -Duration, offset=log(Duration), data=S0, family=poisson())  
#full model  
glm.freq.full <- glm(NB_Claim ~ . -Duration, offset=log(Duration), data=S , family=poisson())  
#traditional model  
glm.freq.trad <- glm(NB_Claim ~ . -Duration, offset=log(Duration),  
                    data=S[, c(1:13, 30)], family=poisson())
```



```

data=S0, offset=log(Duration)+predict(glm.freq.trad, S0)
, family=poisson())

#coefficients and SE
boost_coef[j,] <- c(summary(glm.freq.trad)$coefficients[,1],
                    summary(glm.freq.boost)$coefficients[,1])
boost_stde[j,] <- c(summary(glm.freq.trad)$coefficients[,2],
                    summary(glm.freq.boost)$coefficients[,2])

#####try forecast
pred.nai ve <- predict(glm.freq.nai ve, newdata = forr.data, type="response")
pred.full <- predict(glm.freq.full, newdata = forr.data, type="response")
pred.S3 <- exp(as.matrix(forr.data[,c(2:29)]) %*%
               propd_coef[j, 2:29]+propd_coef[j, 1]+ log(forr.data[, 1]))
pred.trad <- predict(glm.freq.trad, newdata = forr.data, type="response")
pred.boost <- pred.trad * exp(as.matrix(forr.data[14:29])%*%coef(glm.freq.boost))

#remove datasets for this split
rm(tele_ind, test_ind)

#RMSE
RMSEs[j,] <- sqrt(c(
  mean((forr.data$NB_Cli m-pred.nai ve)^2),
  mean((forr.data$NB_Cli m-pred.trad)^2),
  mean((forr.data$NB_Cli m-pred.full)^2),
  mean((forr.data$NB_Cli m-pred.S3)^2),
  mean((forr.data$NB_Cli m-pred.boost)^2)))

#DEV
DEVs[j,] <- c(
  Poisson.Deviance(pred.nai ve, forr.data$NB_Cli m),
  Poisson.Deviance(pred.trad, forr.data$NB_Cli m),
  Poisson.Deviance(pred.full, forr.data$NB_Cli m),
  Poisson.Deviance(pred.S3, forr.data$NB_Cli m),
  Poisson.Deviance(pred.boost, forr.data$NB_Cli m))
})

```

Appendix C

Basic Setup of Proposed Method

$$E \left(\sum_{i=1}^N X_1 X_2 \right) = 0$$

$(a_i X_{i1})$

$$\perp \mathcal{X}^2(\mathcal{X}_1) \quad ;$$

$$E^c(\mathcal{X}_1) = \{1(\mathcal{X}_1) \dots L(\mathcal{X}_1)\} \\ \perp \mathcal{X}^2(\mathcal{X}_1)$$

$$E^c(\mathcal{X}_1) = E^c(\mathcal{X}_1) \quad = 1$$

$$\frac{1}{0} \sum_{i=1}^M (\mathcal{X}_i \mathcal{X}_1) (\mathcal{X}_i \mathcal{X}_1) = \frac{1}{1} \sum_{i=1}^M (1 - i) (\mathcal{X}_i \mathcal{X}_1)$$

$$E^c(\mathcal{X}_1)$$

$$\hat{SPS}(\) = \frac{1}{i-1} \sum_{i=1}^M (1 + \frac{1}{0}) (\mathcal{X}_i \mathcal{X}_1) \quad (; \mathcal{X}_i \mathcal{X}_1 \mathcal{X}_2)$$

$$E^c(\mathcal{X}_1) \mathcal{X}_2 = 1 = E^c(\mathcal{X}_1) = 1$$