

# Sandbox Project: Prediction of Building Claims Costs in SMEs



**Queenie  
Huang**



**Melantha  
Wang**



**Michael  
Jacinto**



**Kevin  
Li**



**Harry  
Peng**

# Background to the project



Project

The (challenging) task of predicting building claims costs in SMEs (Small and Medium Enterprises) across Australia



Issues

Insufficient claims experience per occupation to understand occupation-level risks

High heterogeneity between policies → aggregate modelling would be inefficient

← More granular risk modelling



→ Less differentiating risk modelling



Goals

1) Develop occupational rating scheme both accurate and consistent with domain knowledge

2) Build and test models for predicting working claims cost of SME building insurance



Impact

- Fast growing SME sector → market potential
- More accurate model → sustainable and competitive pricing → edge over other insurers

# Agenda

---

1

## Data Preparation

Impute missing values and handle invalid data.

2

## Exploratory Data Analysis

Identify patterns in the data, group high cardinal variables.

3

## Occupation Grouping

Cluster occupations by claim size and proportion of fire claims, rank them by claim size.

4

## Claims Cost Model

Propose a GLM-XGBoost hybrid model, demonstrate how it can be applied in practice.

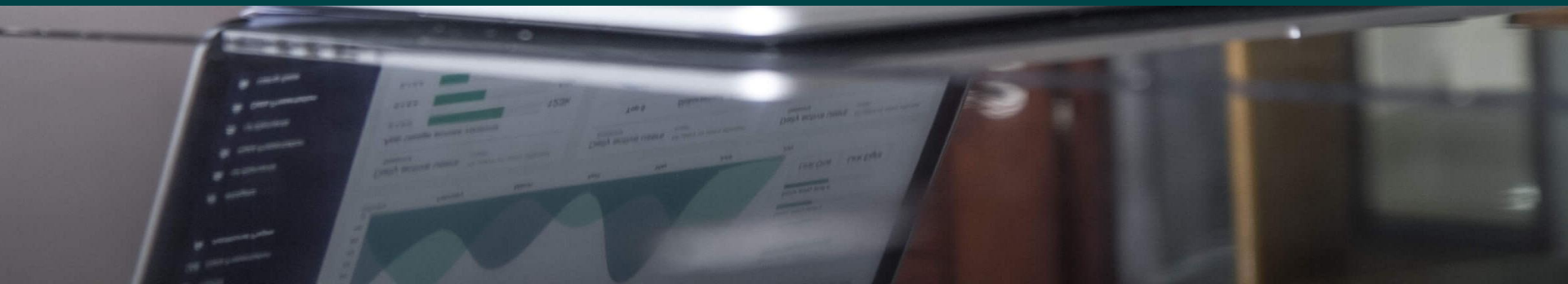
5

## Limitations and Next Steps (including Appendix)

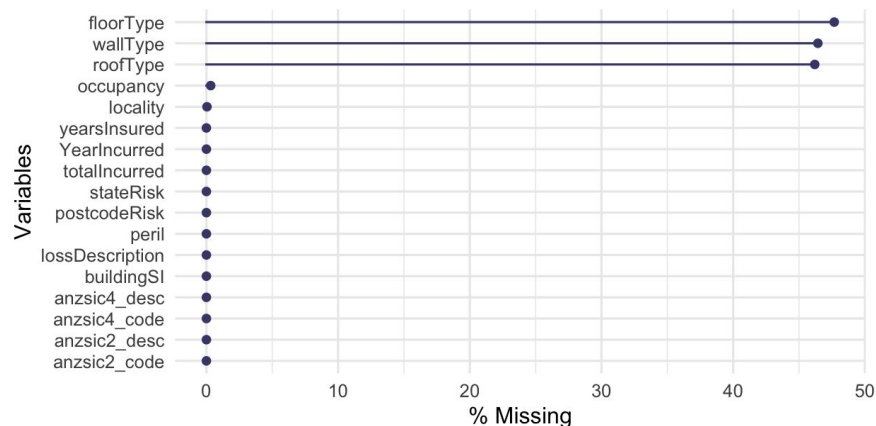
Considered limitations and ways of improvements, further findings



# Data Preparation



## Building materials was plagued with **missing data**



floorType, wallType and roofType are 45-47% missing

Imputation not feasible as too many values missing

Data preparation

EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

## Building materials was plagued with **missing data**

locality	roofType	wallType	floorType	fireprotecAlarm	fire
Officefirst	Alloysteel	Brick	Concrete	0	
Officefirst	Alloysteel	Brick	Concrete	0	
Indust	Alloysteel	Brick	Concrete	0	
Officefirst	Alloysteel	Brick	Concrete	0	
Conversion	Asbestos	Brick	Z-Missing	0	
Conversion	Alloysteel	Brick	Concrete	0	
Conversion	Alloysteel	Concretsla	Concrete	0	
Substfirst	Othmixed	Brick	Othmixed	0	
Majorrd	Alloysteel	Brick	Concrete	0	

floorType, wallType and roofType are  
45-47% missing

Imputation not feasible as too many  
values missing

**Solution: code the missing values as  
a separate factor**

Data preparation

EDA

Occupation  
grouping

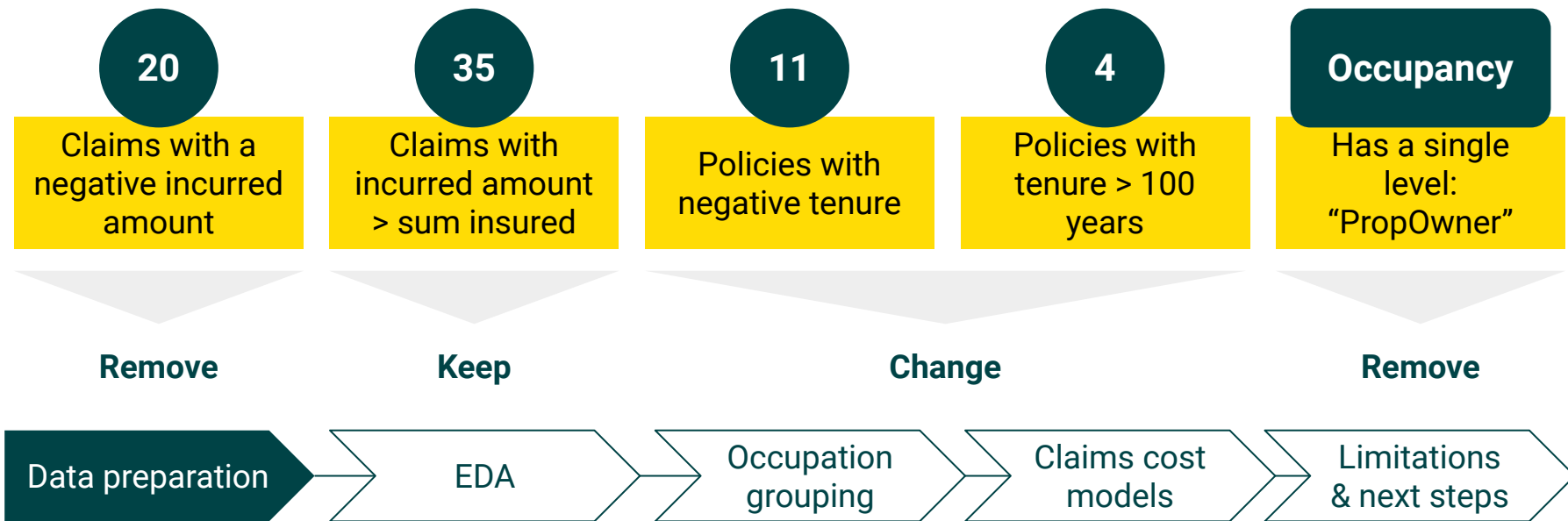
Claims cost  
models

Limitations  
& next steps

## Invalid data was treated depending on the variable

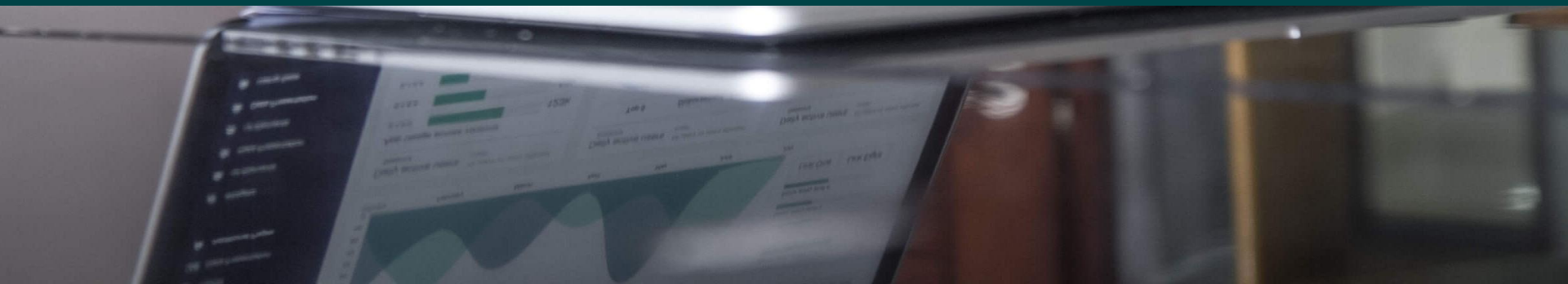
### Anomalies (invalid data)

Out of the **6,746** claims from the data,



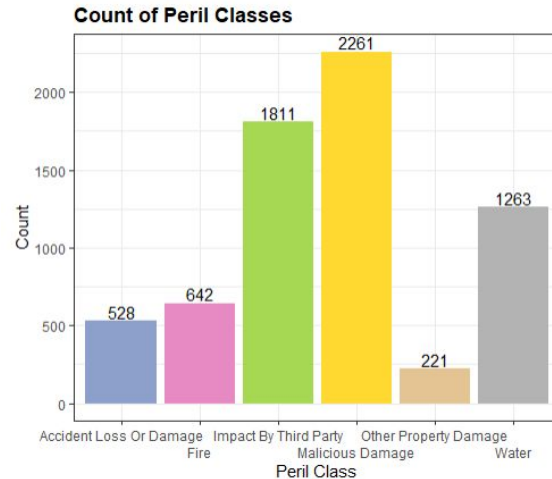
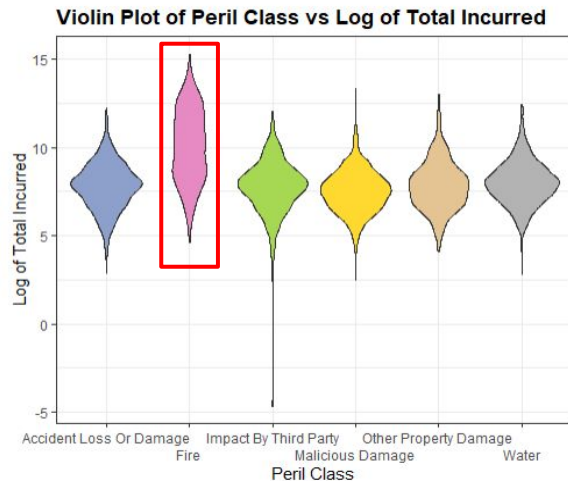


# Exploratory Data Analysis





## Fire claims show different behaviours to other claim types



**Fire claims, although less frequent, tend to be more severe than all other perils**

Data preparation

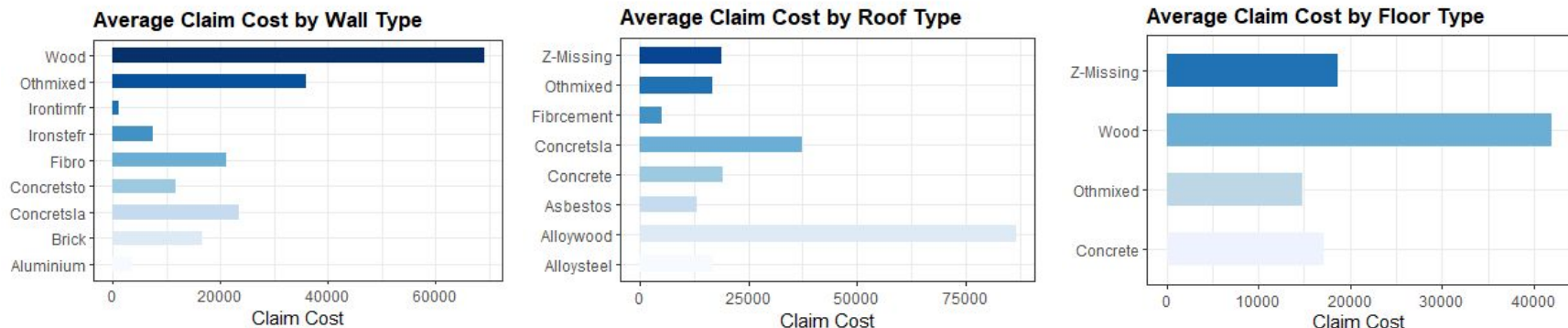
EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

## Building materials have a significant impact on claim size



**Alloywood and wood have higher average claims costs than other materials**

Data preparation

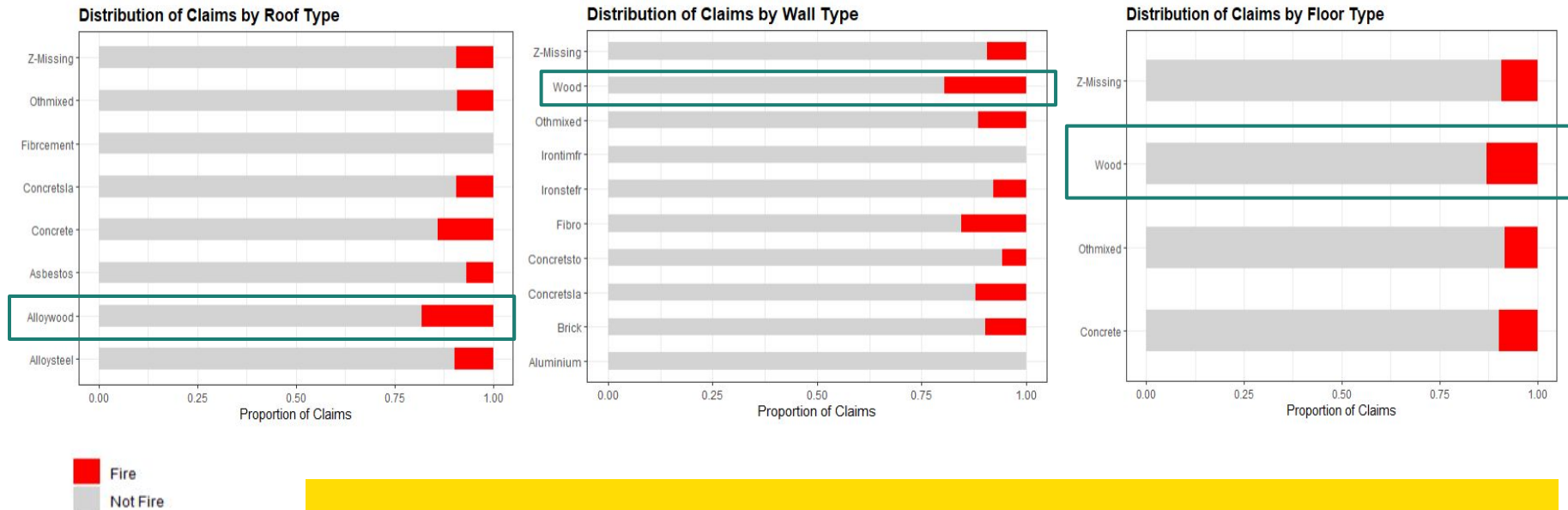
EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

# Building materials have a significant impact on claim size



**Wooden roofs, walls, and floors are more susceptible to fire than other materials**

Data preparation

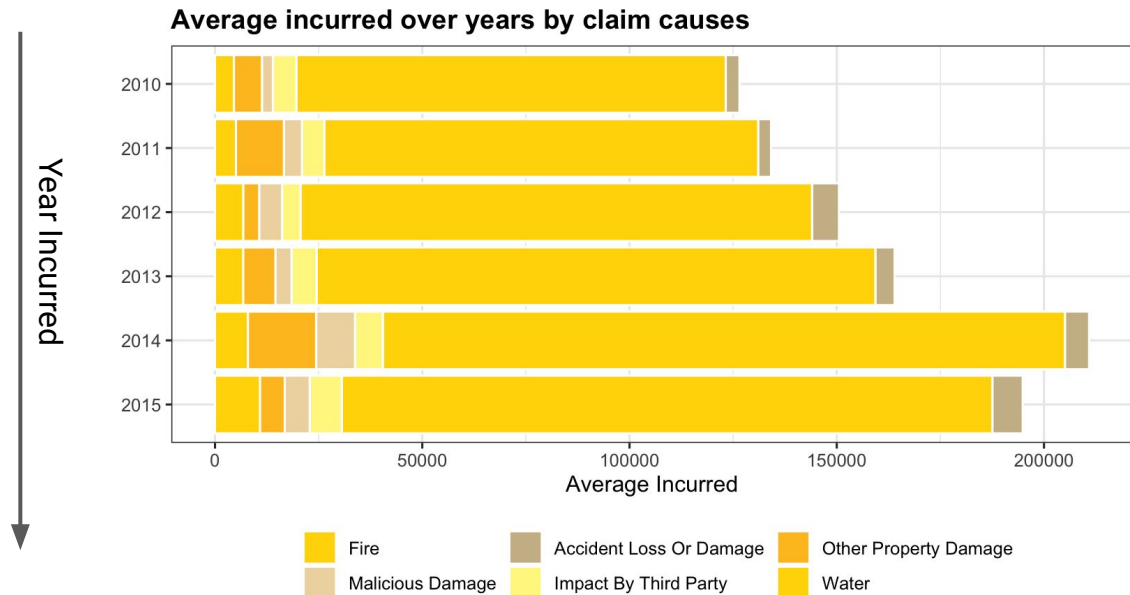
EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

## Later modelling must account for **superimposed inflation**



Steady rise in average incurred across years

Potential inflation impact

Data preparation

EDA

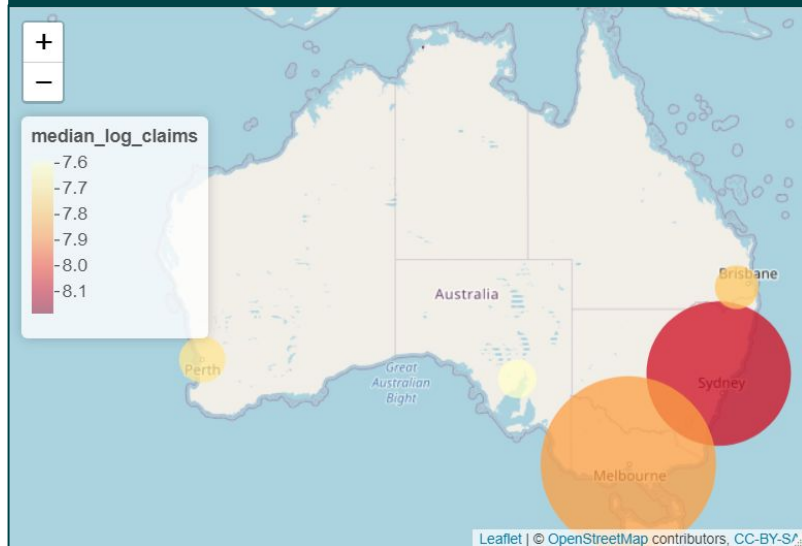
Occupation grouping

Claims cost models

Limitations & next steps

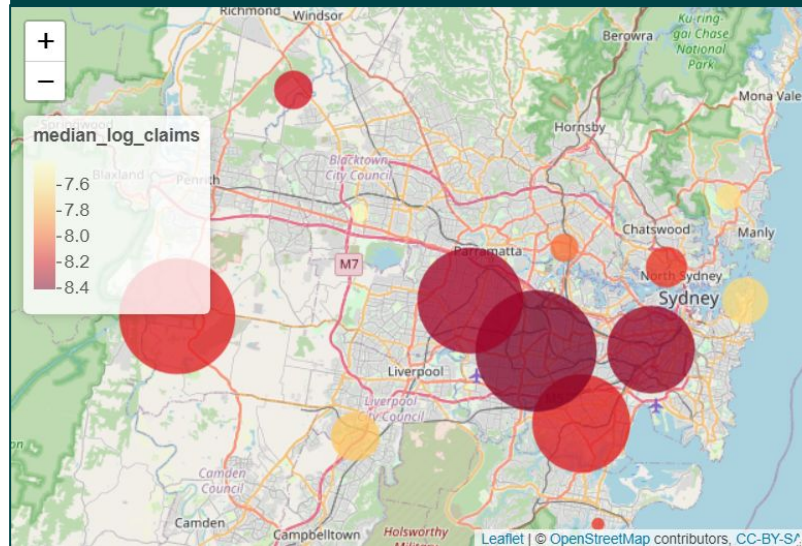
# Postcode risk was regrouped into Statistical Area (level 4) to reduce feature cardinality

Median total incurred\* by states is too broad  
(while postcode is too granular)



*\*size of circles indicates number of claims*

Median total incurred by  
Statistical Area Level 4



Data preparation

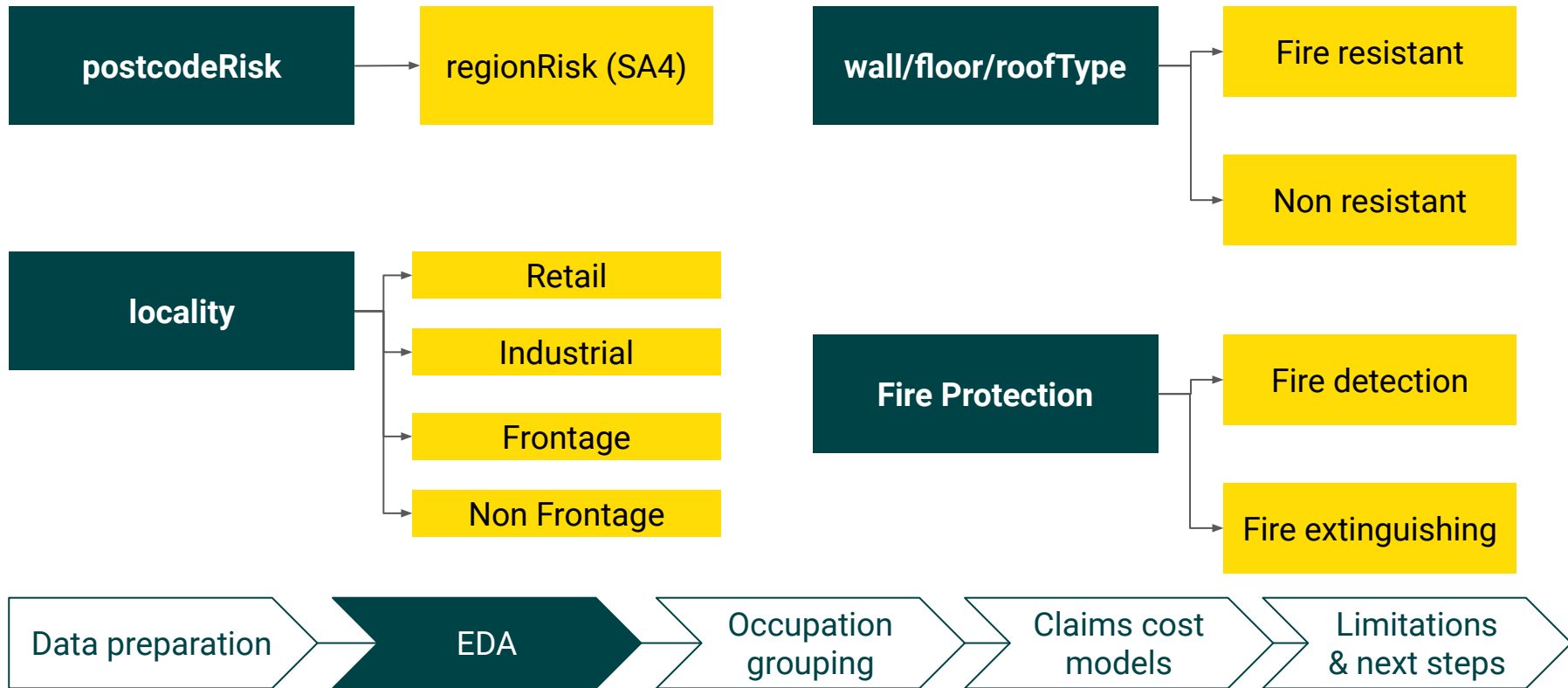
EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

## Features were **grouped or combined** to simplify the dataset

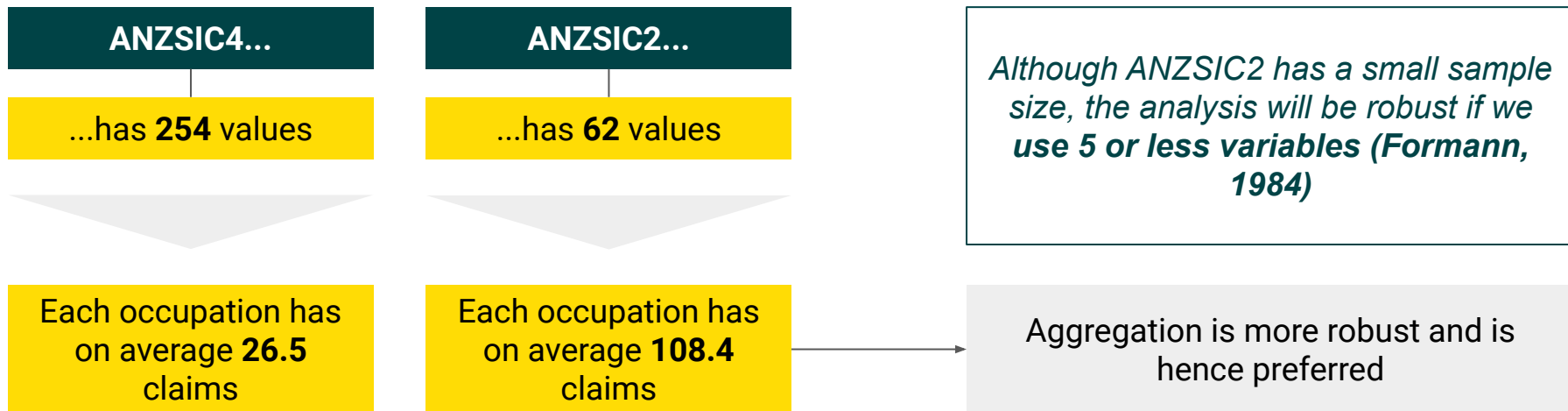




# Occupation Grouping

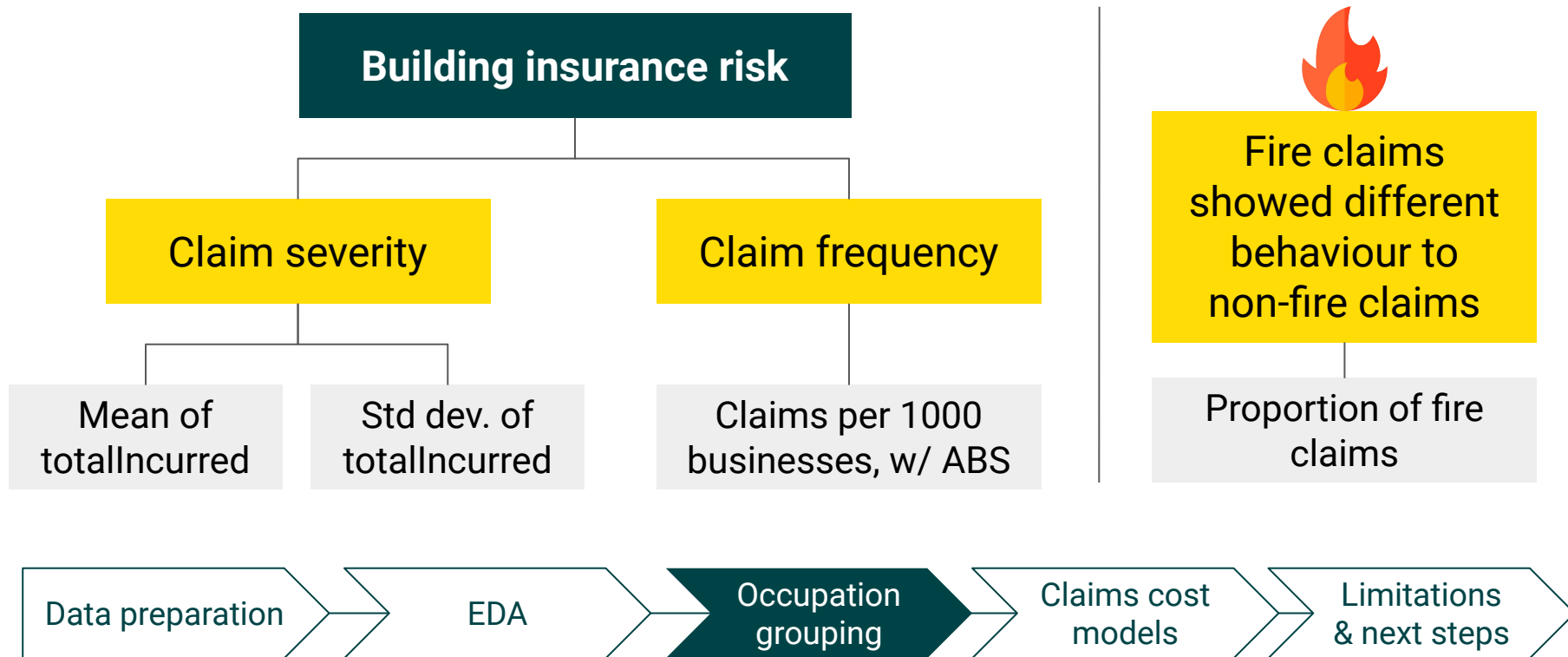


## The data was **aggregated by ANZSIC2** code

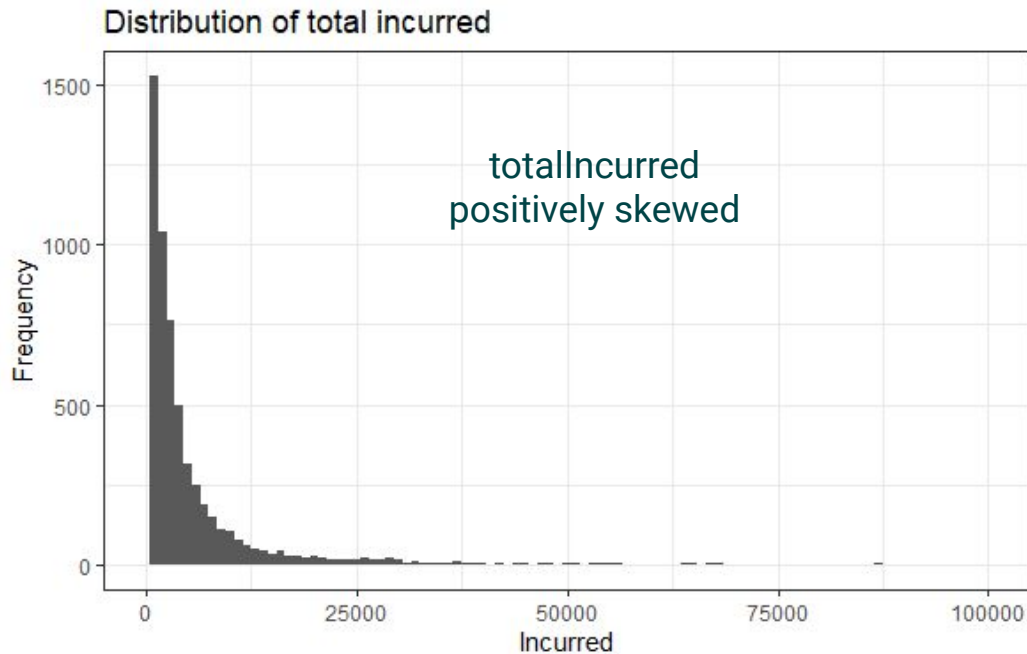




Claim **severity**, **frequency** and **fire proportion** were used to quantify risk



## Variables were **normalised** to minimise the impact of outliers



**Take logarithm of  
totalIncurred**

Data preparation

EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

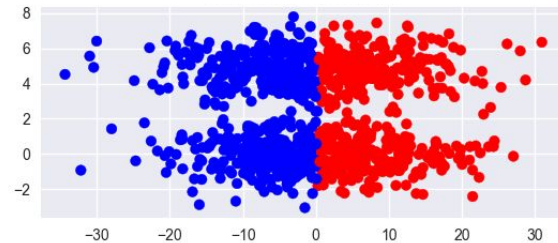
# Variables were **normalised** to minimise the impact of outliers

Clustering methods use distance-based measures

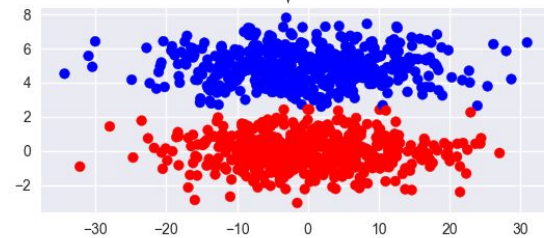
Algorithms become sensitive to large values

Min-Max scaler to map values to  $[0,1]$

Before normalising



After normalising



Data preparation

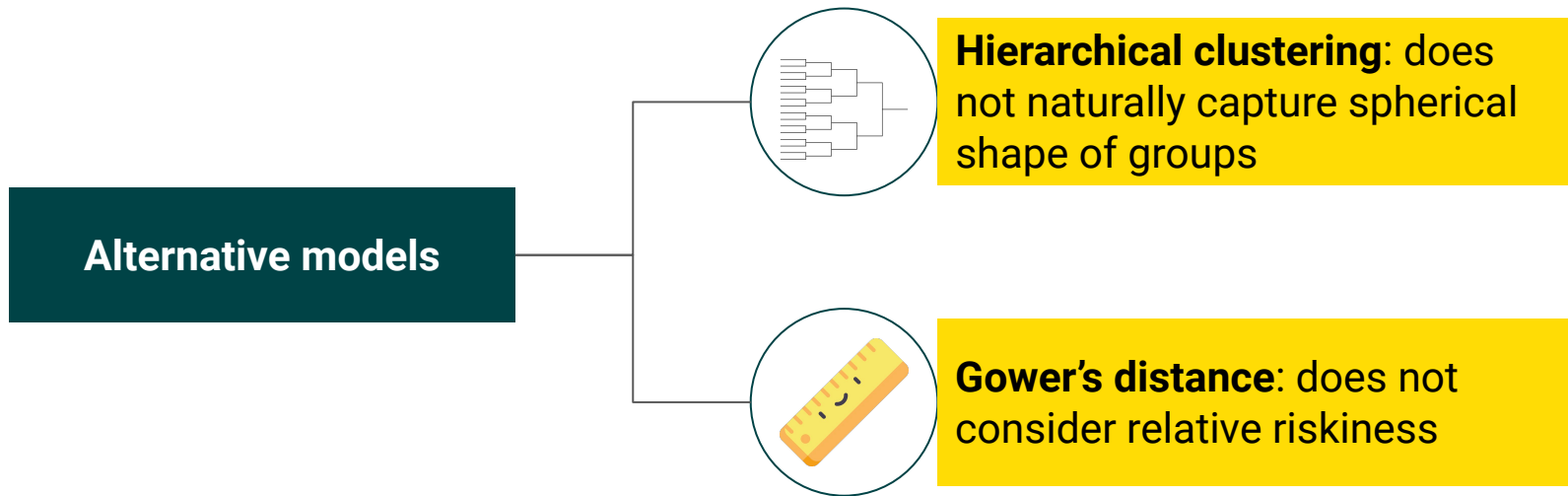
EDA

Occupation  
grouping

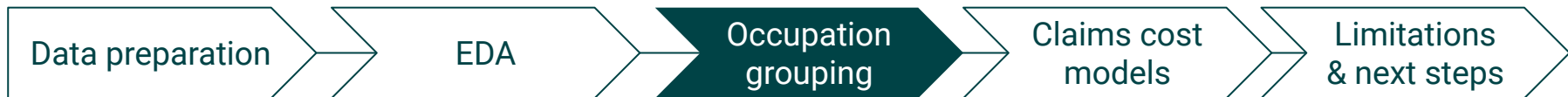
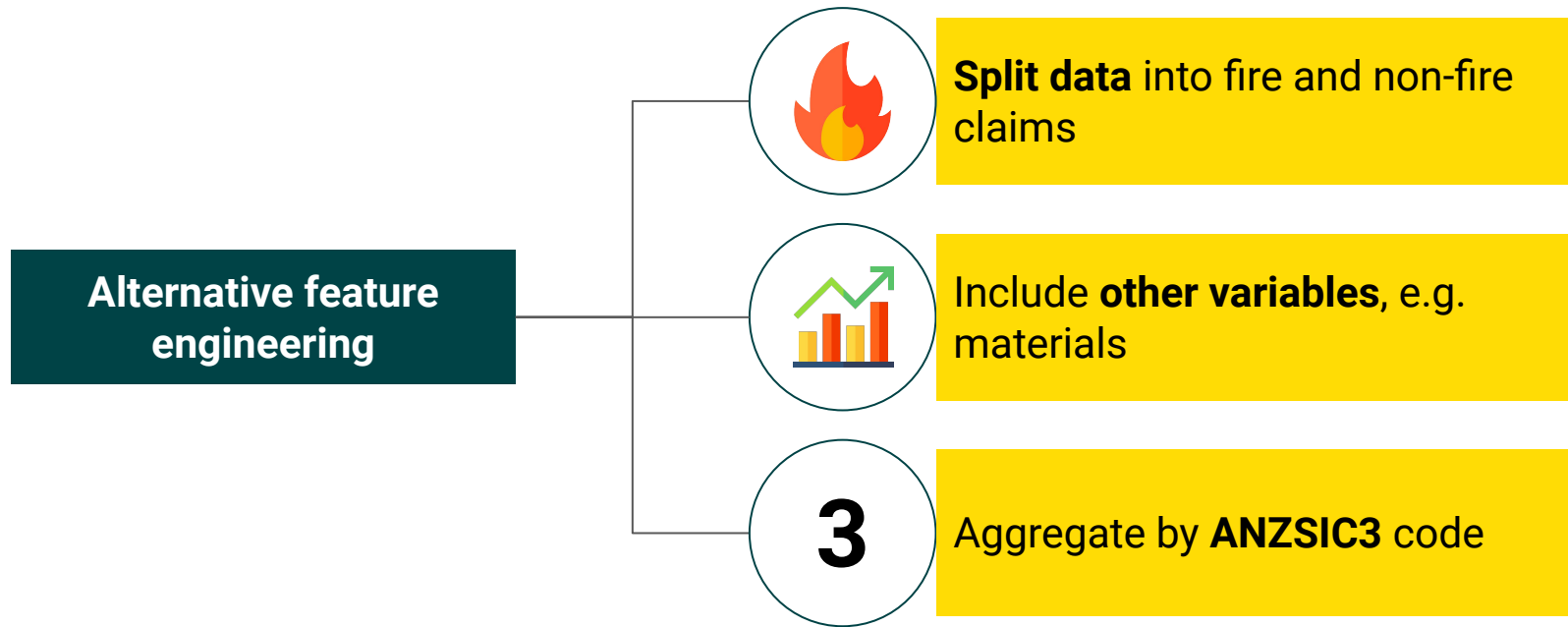
Claims cost  
models

Limitations  
& next steps

## Other models, settings, and variables were considered

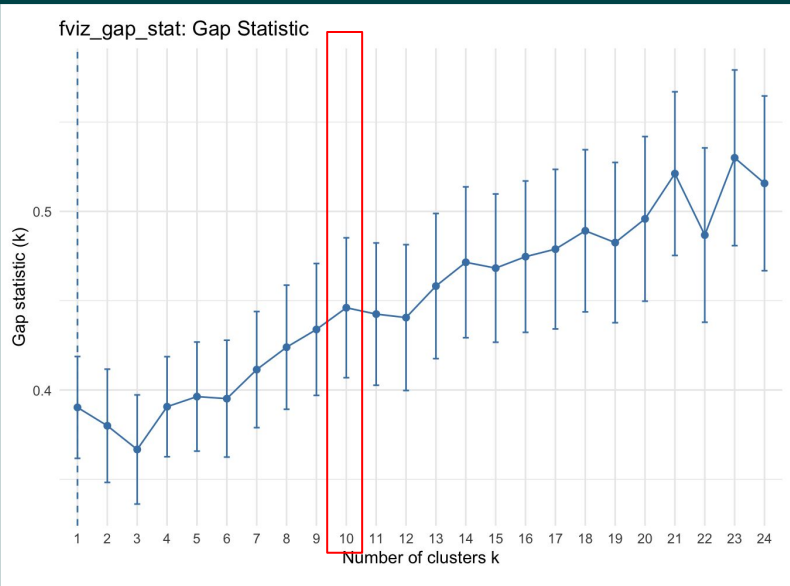


## Other models, settings, and variables were considered



# In K-means clustering, choosing the right number of clusters is important

Gap plot shows how gap between clusters increases with k



Elbows at 10, 14, 18, 21, 23

Too many clusters lead to overfitting, too many categorical variables

10 clusters chosen

Data preparation

EDA

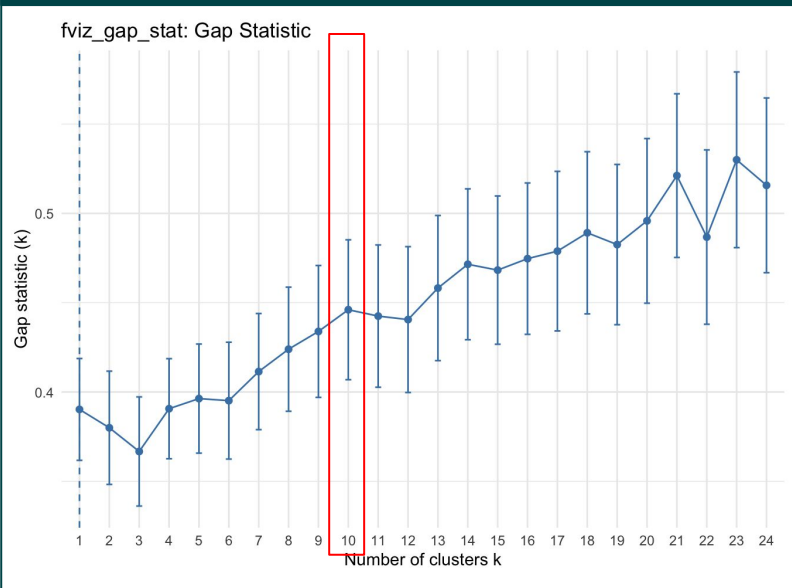
Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

# Choosing the right number of clusters is important

Gap plot shows how gap between clusters increases with  $k$



Elbows at 10, 14, 18, 21, 23

Too many clusters lead to overfitting, too many categorical variables

10 clusters chosen

Combine groups 9 and 10

9 distinct groups

Data preparation

EDA

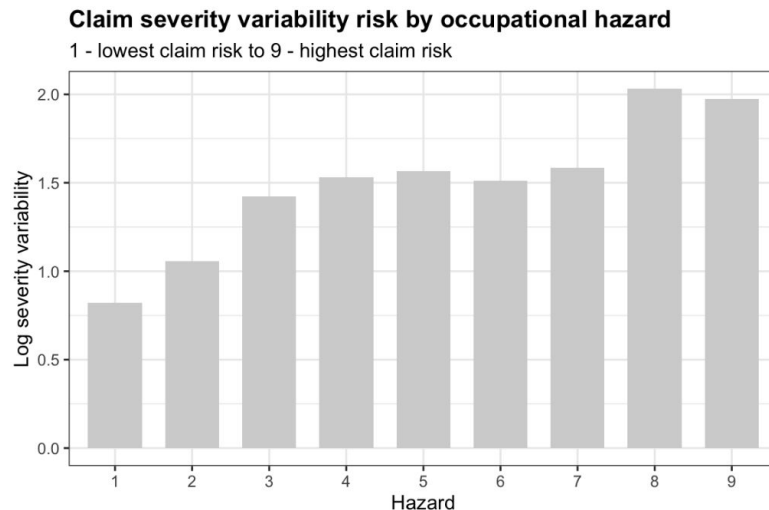
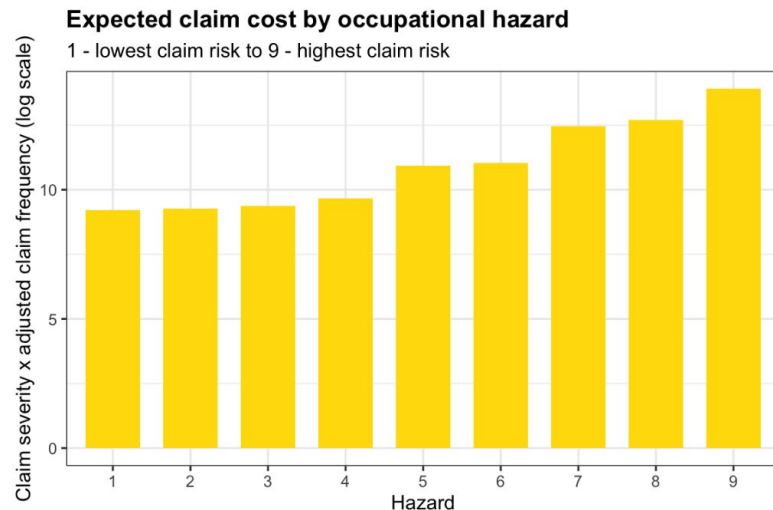
Occupation grouping

Claims cost models

Limitations & next steps

# Clusters were **ranked** based on **expected claim size**

Expected claim size  $\approx$  Claim frequency  $\times$  Claim severity



Data preparation

EDA

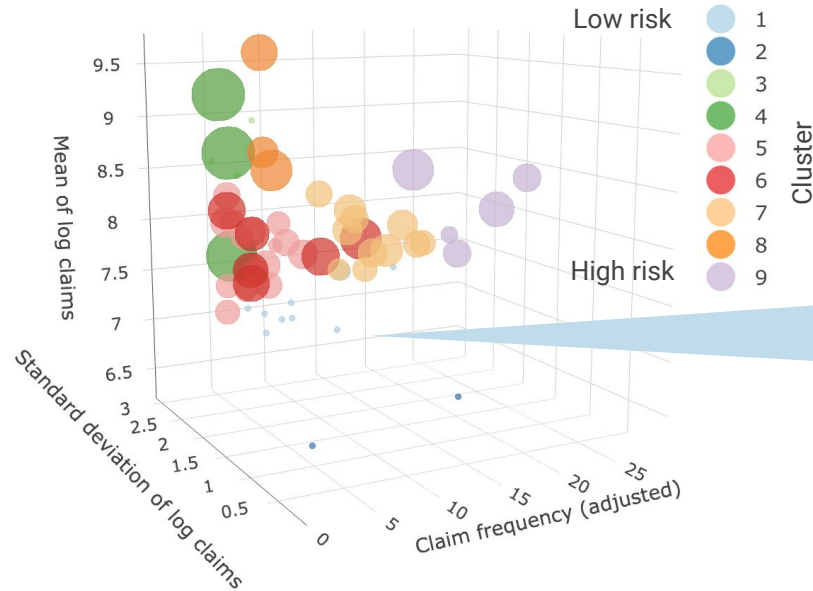
Occupation  
grouping

Claims cost  
models

Limitations  
& next steps



# Education and media businesses are prevalent in the **low-risk** group



Size of bubble indicates proportion of fire claims

## Group 1 includes:

- Education (adult, tertiary)
- Media and arts (broadcasting, motion picture, sound recording, performing arts)
- Construction
- Telecommunications

Data preparation

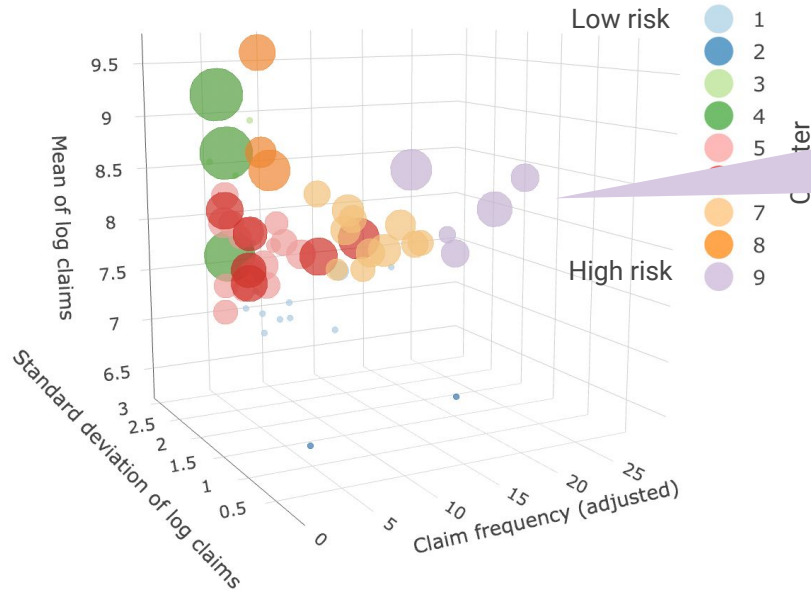
EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

# Retail and hospitality makes up the riskiest occupations



## Group 9 includes:

- Fuel Retailing
- Food Retailing
- Food and Beverage Services
- Heritage Activities
- Library and Information Services

Size of bubble indicates proportion of fire claims

Data preparation

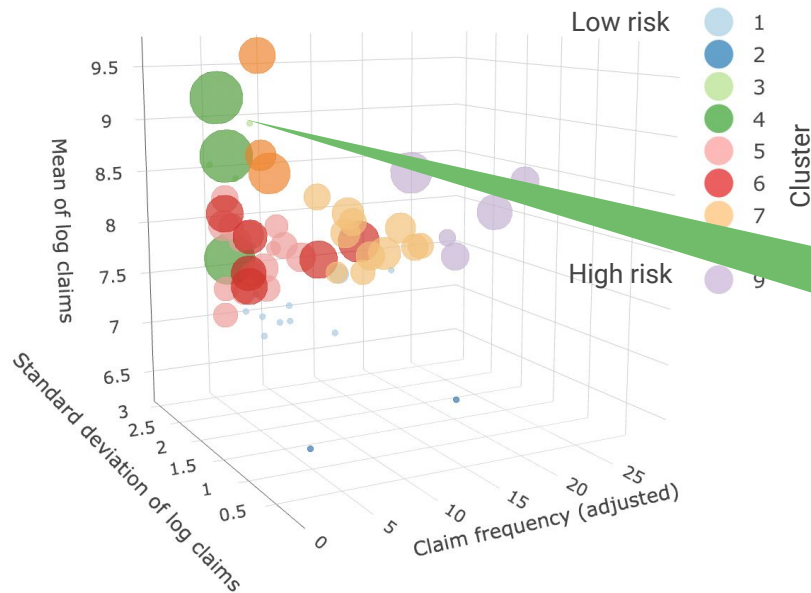
EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

## Group 4 suffers from large claims, but benefits from a low claim frequency



### Group 4 includes:

- Agriculture
- Transport support services
- Internet service providers

Size of bubble indicates proportion of fire claims

Data preparation

EDA

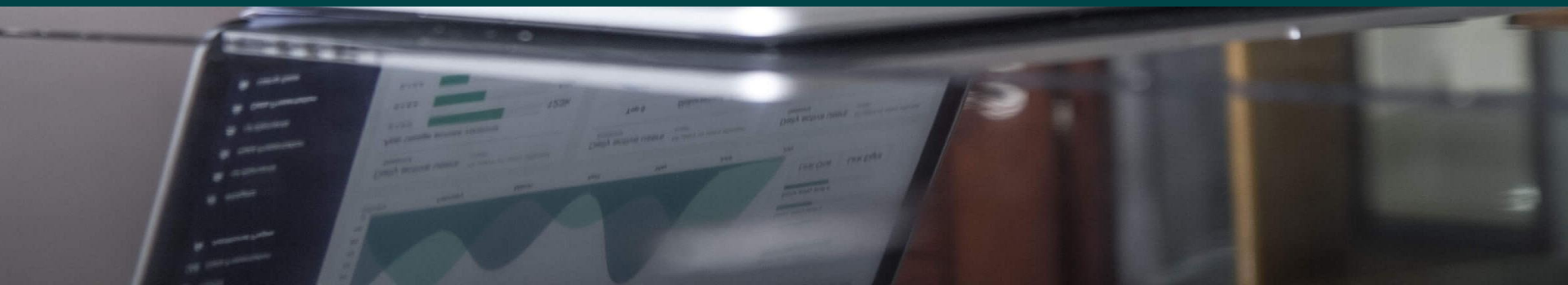
Occupation  
grouping

Claims cost  
models

Limitations  
& next steps



# Claims Cost Models



## The data was **split into a training and testing set**

**Objective: To accurately model working claims cost for building insurance**

The data was randomly split into...

**Training set**

80% of data

To aid in model  
tuning and fitting  
processes

**Testing set**

20% of data

To evaluate  
**performance**  
across models

**Root mean square error (RMSE)**

- Learning objective (for training models)
- Performance metric (for model selection)

$$RMSE = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$

Data preparation

EDA

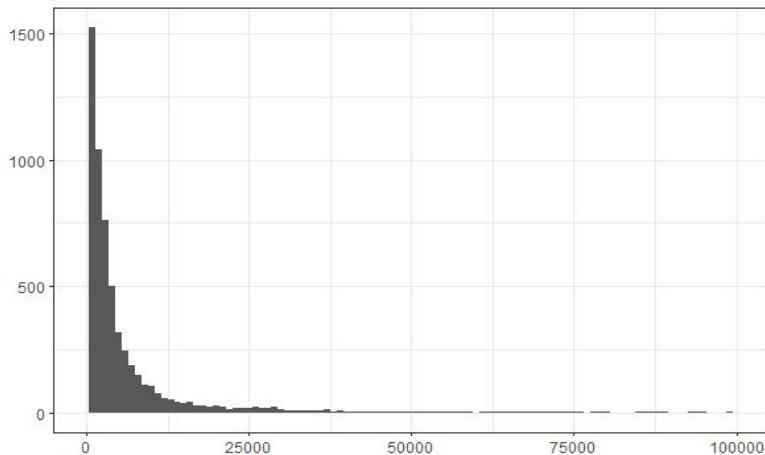
Occupation  
grouping

**Claims cost  
models**

Limitations  
& next steps

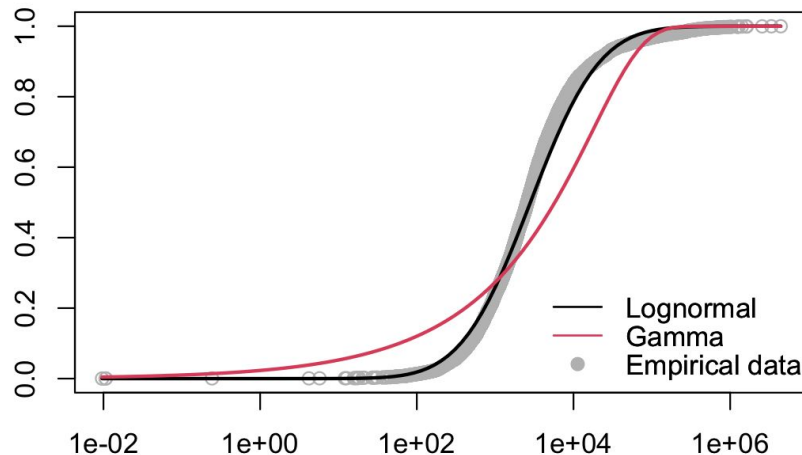
# We fit a normal distribution to the **log transformed** data to optimise the GLM

Distribution of total incurred



Incurred

Empirical and fitted parametric CDFs (without predictors)



Incurred

Data preparation

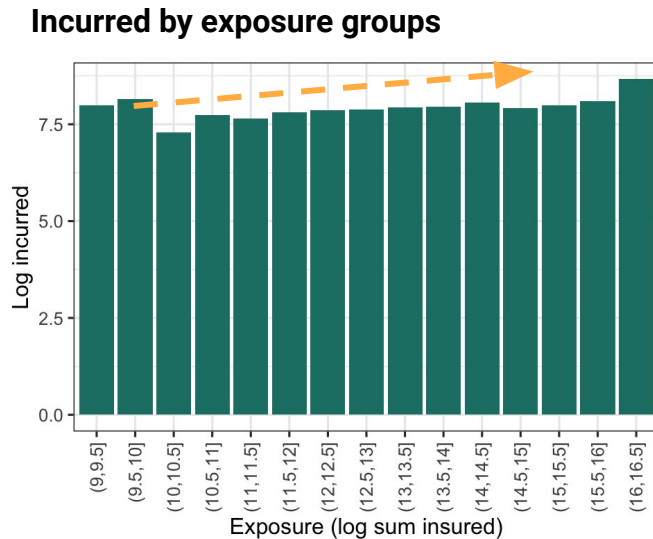
EDA

Occupation  
grouping

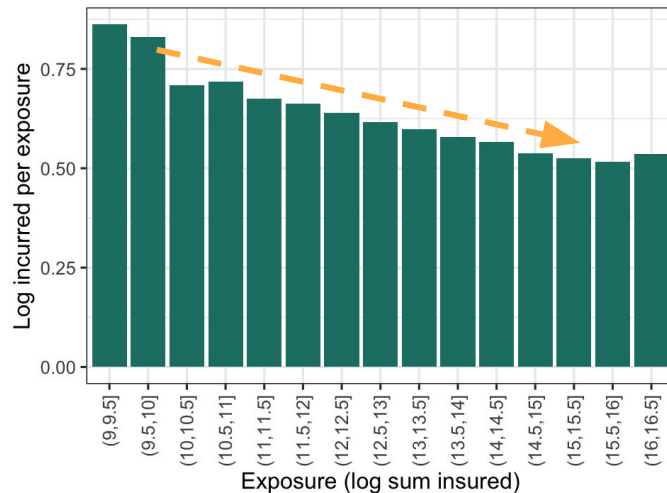
Claims cost  
models

Limitations  
& next steps

## Sum insured was treated as a predictor



**Incurred per exposure by exposure groups**



Relationship between incurred and sum insured may not be strictly pro-rata

Data preparation

EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

## Other **model considerations** were assessed and included to the model

---

**Material types, fire protection,  
occupations, postcodes**

In line with EDA results

**Superimposed inflation**

yearIncurred as a categorical predictor

**Fire/non-fire vs aggregate model**

Balanced the lack of data and the significantly  
different distributions

Data preparation

EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps



# The models need to be both **accurate and interpretable** to meet business objectives

## Business objectives



### Accuracy

**High** - Accurate prediction for claim costs is the basis for competitive and sustainable pricing schemes.



### Interpretability

**High** - It is crucial to understand the risk drivers in order to inform product design and risk management, and communicate the insights to other teams and/or management.



### Efficiency

**Low** - The dataset is relatively small (thousands of observations x 10-20 features). The models need not to be re-trained frequently.

Data preparation

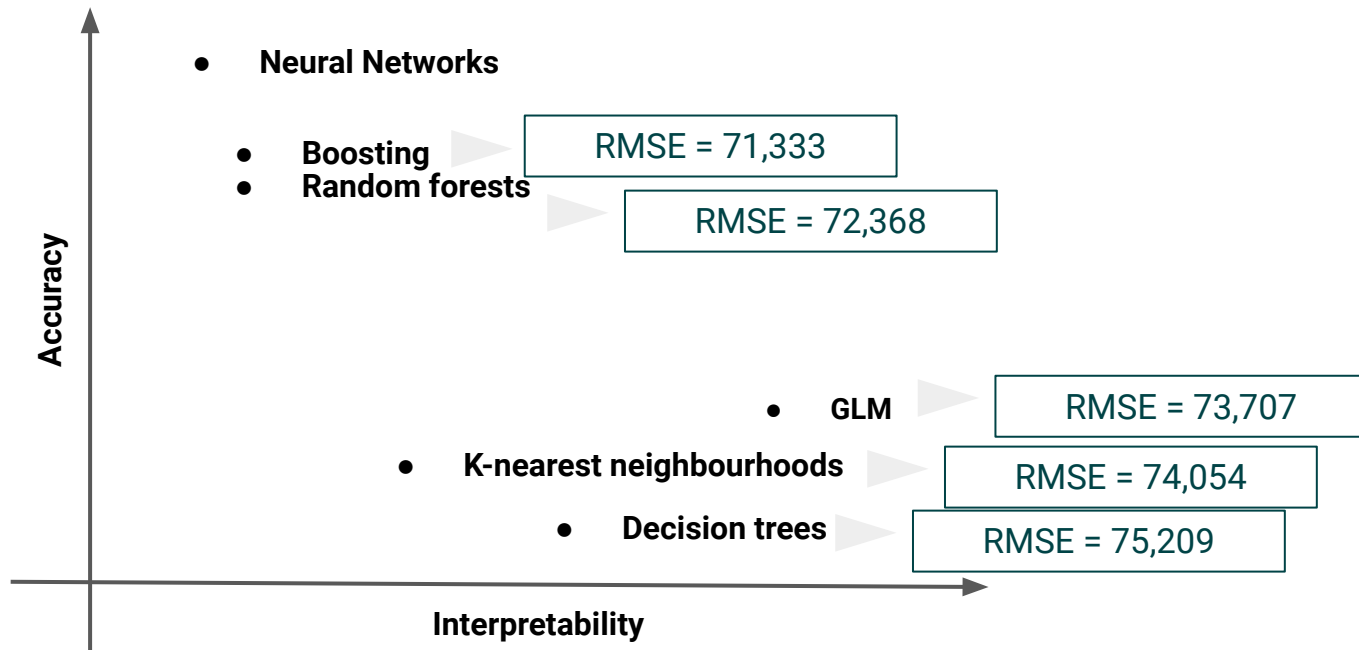
EDA

Occupation  
grouping

Claims cost  
models

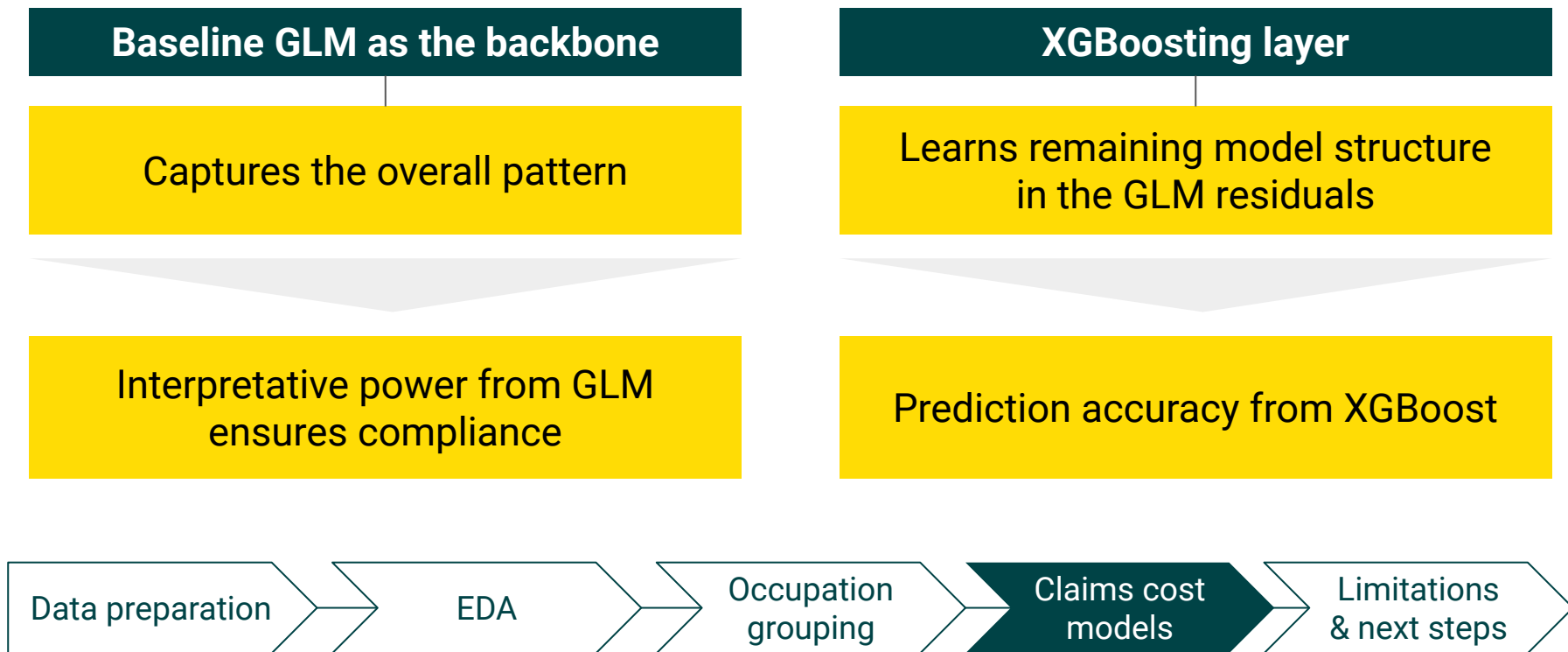
Limitations  
& next steps

But it is well known that a **compromise exists** between accuracy and interpretability



To balance these objectives, our team proposes an innovative model:  
**XGBoosted GLM!**

---



## XGBoosted GLM maintains XGBoost's **accuracy** while retaining **interpretability** of GLM

Error comparison of GLM, XGBoost, and XGBoosted GLM

RMSE	GLM	XGBoost	XGBoosted GLM	XGBoosted GLM: Improvement over GLM baseline
Fire claims	525,465	505,244	508,709	3.2%
Non-fire claims	22,841	22,482	22,456	1.7%
Aggregate	73,707	71,333	71,660	2.8%



## GLM identified **different significant predictors** for fire and non-fire claims

Fire			
Variable	Estimated effect	F-statistics	Significance
<b>buildingSI</b>	<b>1.17e-07</b>	<b>8.008</b>	★★☆
YearIncurred	Multiple factors	1.126	
yearsInsured	-7.38e-03	0.424	
<b>regionRisk</b>	<b>Multiple factors</b>	<b>1.650</b>	★★☆
locality	Multiple factors	1.123	
roofType_resist	-6.25e-01	1.429	
wallType_resist	-1.25	1.017	
floorType_resist	4.53e-01	0.307	
<b>fire_detection</b>	<b>6.50e-01</b>	<b>5.695</b>	★★☆
fire_extinguishing	3.23e-01	1.327	
<b>hazard</b>	<b>Multiple factors</b>	<b>1.988</b>	☆☆☆

Significance: ★★★ <0.001; ★★☆ <0.01; ★☆☆ <0.05; ☆☆☆ <0.10

Non-fire			
Variable	Estimated effect	F-statistics	Significance
<b>buildingSI</b>	<b>2.28e-08</b>	<b>9.433</b>	★★☆
<b>YearIncurred</b>	<b>Multiple factors</b>	<b>12.950</b>	★★★
<b>yearsInsured</b>	<b>-8.19e-03</b>	<b>6.382</b>	★★☆
<b>regionRisk</b>	<b>Multiple factors</b>	<b>2.245</b>	★★★
locality	Multiple factors	1.413	
roofType_resist	-7.76e-02	0.438	
wallType_resist	4.94e-01	0.107	
<b>floorType_resist</b>	<b>-6.03e-01</b>	<b>6.928</b>	★★★
fire_detection	8.93e-02	0.363	
fire_extinguishing	-5.75e-02	0.868	
<b>hazard</b>	<b>Multiple factors</b>	<b>3.813</b>	★★★
<b>peril</b>	<b>Multiple factors</b>	<b>18.985</b>	★★★

The split models together have an **RMSE of 73,726**, a vast improvement over the previous two models!

Data preparation

EDA

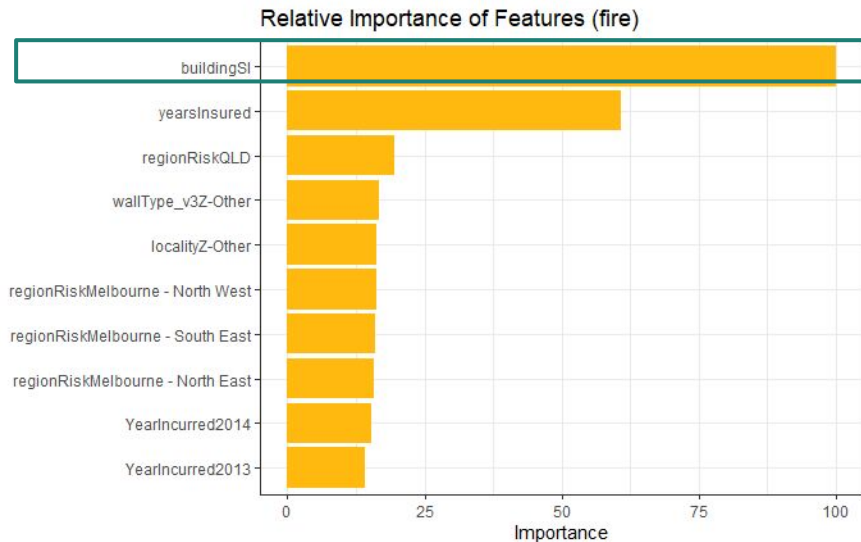
Occupation  
grouping

Claims cost  
models

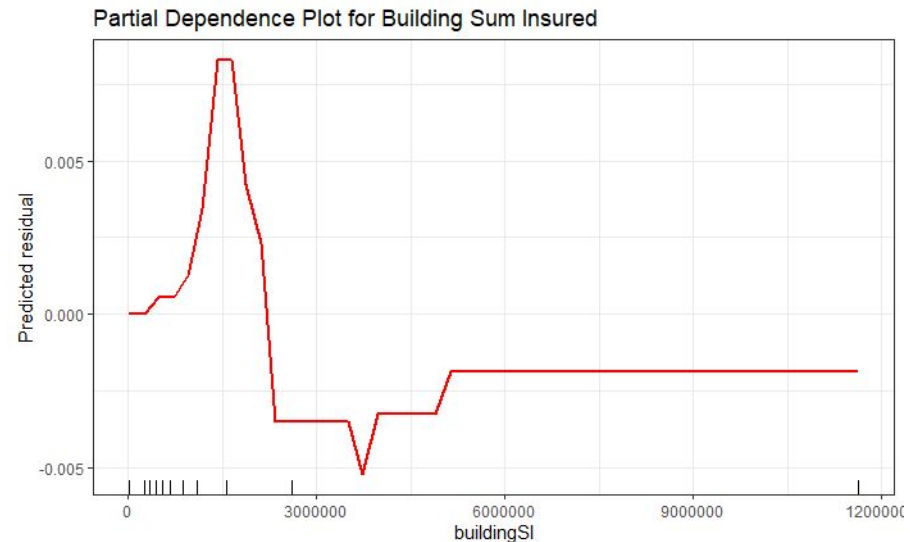
Limitations  
& next steps

## XGBoost effectively picked up the **remaining patterns**: fire example

### Variable importance



### Partial dependence



Data preparation

EDA

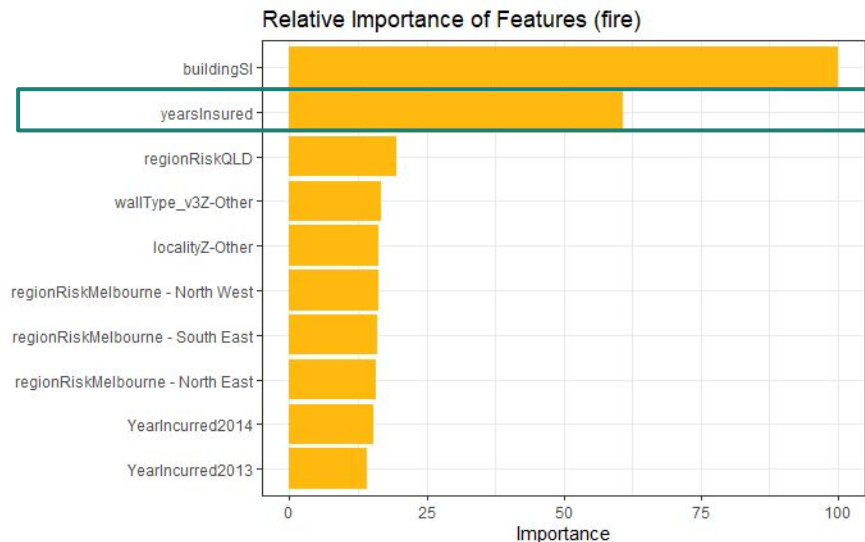
Occupation  
grouping

Claims cost  
models

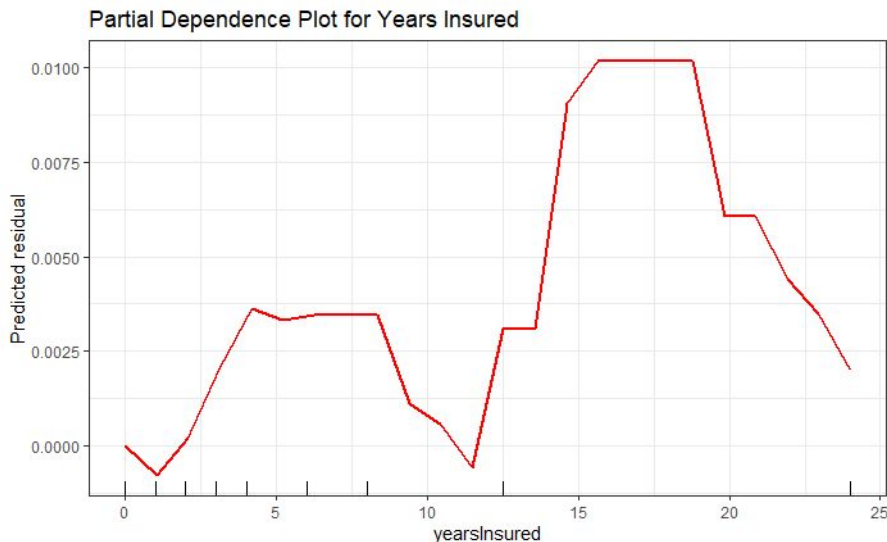
Limitations  
& next steps

## XGBoost effectively picked up the **remaining patterns**: fire example

### Variable importance



### Partial dependence



Data preparation

EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps

## We make three recommendations to make our model **commercially viable**

### Converting models into rating tables

Discriminate customers based on risk drivers identified by GLM (e.g., buildingSI, yearIncurred, occupation)

### Further experimentation with the hybrid model

Consider other parametric models for baseline and/or deep learning techniques for learning residuals

### Integration with peril classifier (either model or underwriter)

Instead of assuming known perils, we can integrate our model with peril classifier that estimates  $\Pr(\text{fire claim})$

Data preparation

EDA

Occupation  
grouping

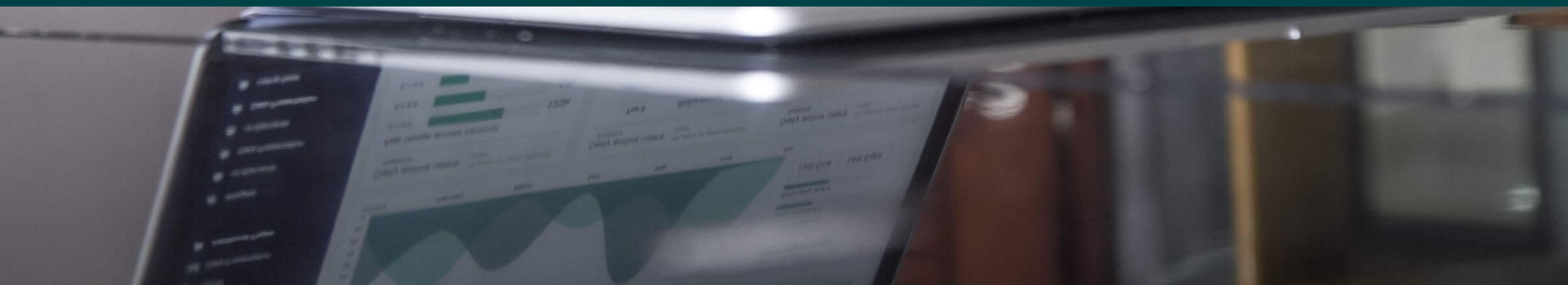
Claims cost  
models

Limitations  
& next steps





# Limitations & Next Steps



## Our results inherited any limitations in the data sources



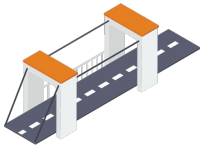
### Limitations

1

Relatively small size of data may constrain the performance of ML models

2

Lack of policy information/exposure data



### Next Steps

1

Augment with industry and external datasets (e.g. population density)

2

Explore additional models (e.g. credibility models)

3

Linking fire and non-fire models trained on separate datasets

Data preparation

EDA

Occupation  
grouping

Claims cost  
models

Limitations  
& next steps