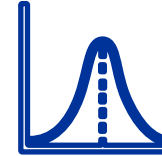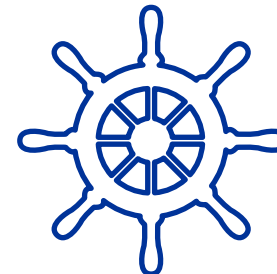# Feature Engineering in GLM

**Titanic Exercise**

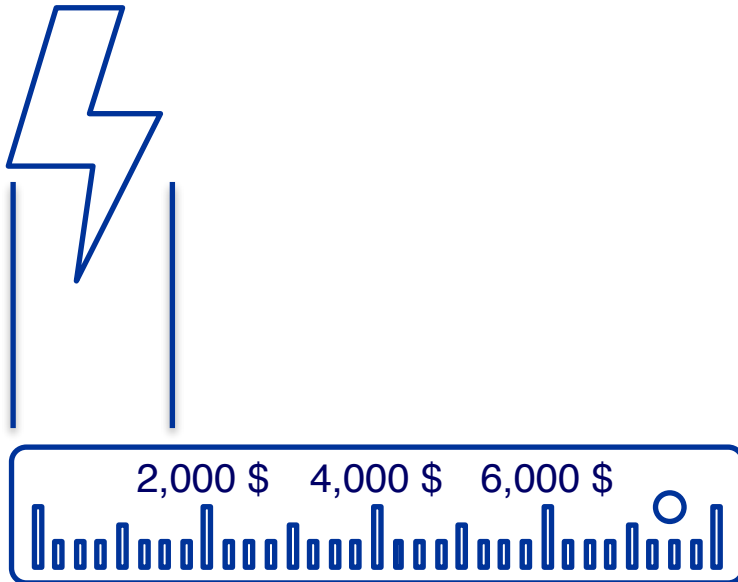*Peter Cvacho, Data Scientist*
*EUBA - Bratislava*
2018-03-26

# Risk in our everyday life

# How to measure and predict?

2,000 $    4,000 $    6,000 $

# Solution
## Generalized Linear Models

George E. P. Box:
"*Essentially, all models are wrong, but some of them are useful*"
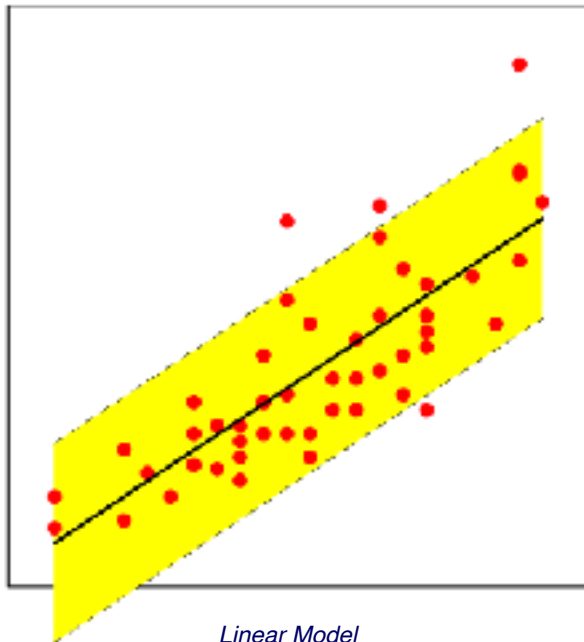
*So, Why GLM?*

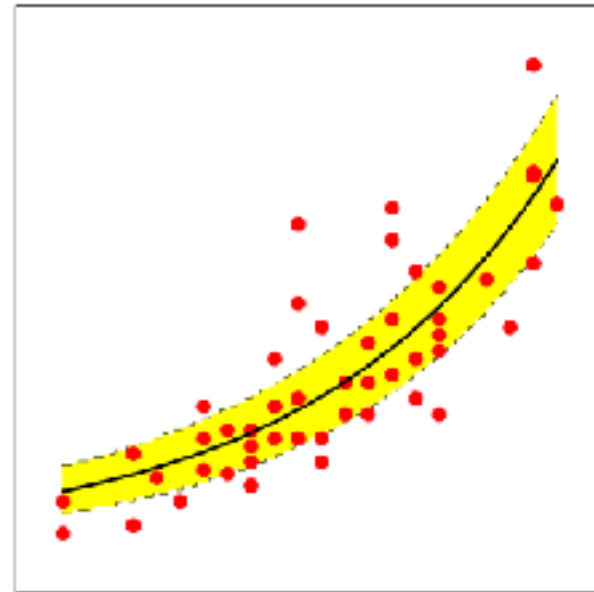*Because we can*.

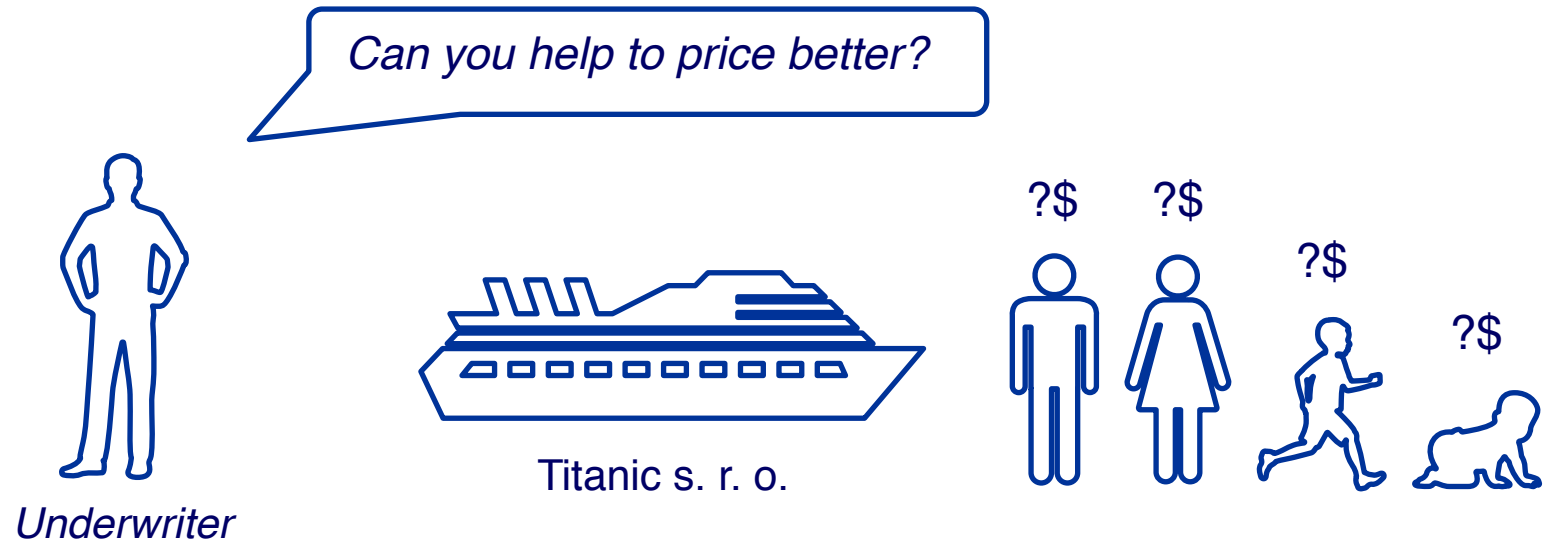*if assumptions of GLM are satisfied

# LM vs. GLM



Linear Model

Generalized Linear Model (Poisson)

*Credit: https://freakonometrics.hypotheses.org/9593*

*simplified visualisation*

# Always care about assumptions!

# Case Intro

UW Department: Our pricing is not accurate



*Can you help to price better?*

*Underwriter*

Titanic s. r. o.

?$  ?$  ?$  ?$

# Case Exploration

Statistics might help!

## Titanic s. r. o - Pricing Tool

**Sum Insured**

**Sex**
male

**Age**

Probability of having claim:
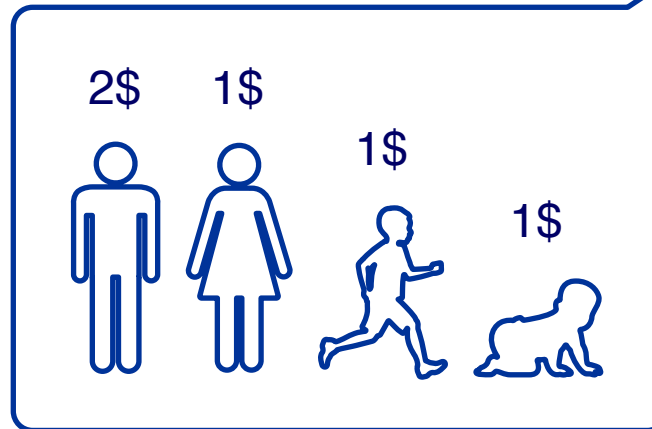
**0.146**

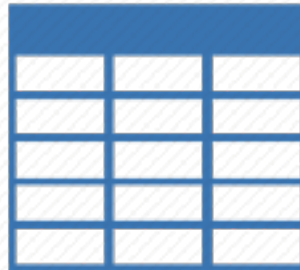Recommended price:

**GBP2.19**

Titanic s. r. o.

2$ 1$ 1$ 1$

# Solution
## Generalized Linear Models
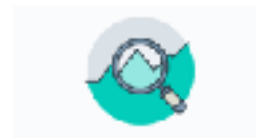
9891 passengers

Having (1) / Not having (0) Claim

Features: Sex, Age, …

*Binomial Distribution* *in Regression*

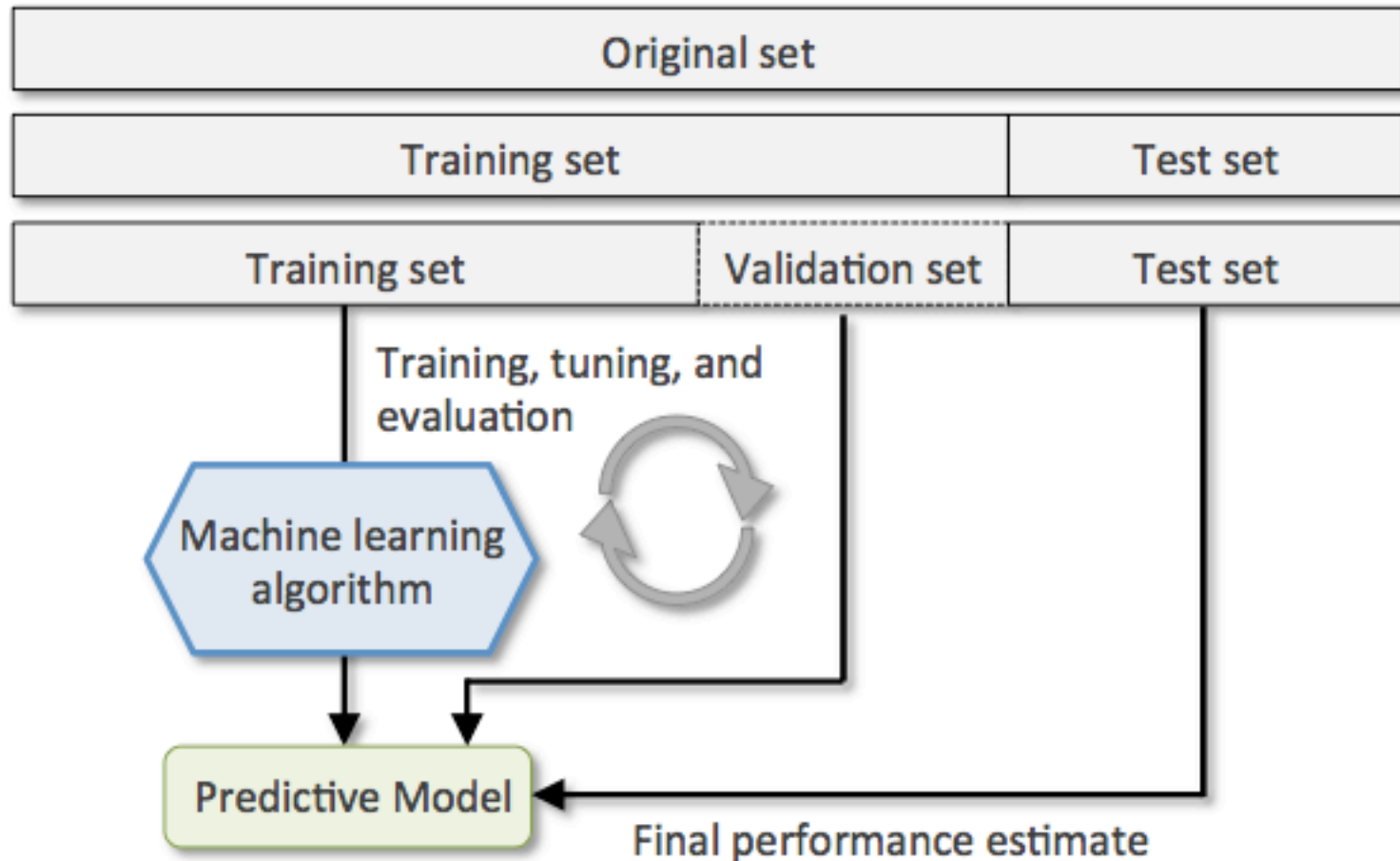*Def.:*
$$\mathbf{E(Y)} = \mu = g^{-1}(\mathbf{X}\beta)$$

*Link Function:*
$$\mathbf{X}\beta = \ln\left(\frac{\mu}{1-\mu}\right)$$

*Mean Function:*
$$\mu = \frac{\exp(\mathbf{X}\beta)}{1 + \exp(\mathbf{X}\beta)} = \frac{1}{1 + \exp(-\mathbf{X}\beta)}$$

# Modeling
Training vs. Validation



We will use rule of 70:20:10

# Modeling

Exercise

*Our Goal:*
We want to know, what is the risk, passengers are facing during
the trip and what should be sustainable price in case of accident.

## Improve Current Model

```
glm(data = train_70,
        formula = Claim ~ sex + age,
        family = binomial())
```

## How?

- Add New Features
- Capping
- Grouping Continuous Features
- Normalisation
- Interactions
- Elimination of correlation
- Etc.

# Cheatsheet

Base R

*http://github.com/rstudio/cheatsheets/raw/master/base-r.pdf*

# Modeling
## Outliers

# Summary

What have you learnt today?

- Training vs. Validation

- Checking of missing values and its imputation

- Identifying outliers and capping them

# Materials

- Many Cheatsheets
  https://www.rstudio.com/resources/cheatsheets/

- More about Shiny (gallery, tutorials, articles, …)
  https://shiny.rstudio.com/

- R Programming
  http://www.cookbook-r.com/

- GLM Paper
  https://www.jstor.org/stable/2344614?seq=1#page_scan_tab_contents

https://www.surveymonkey.com/r/X6CTS3D