

**Big Data und Machine Learning**  
**bei Schadenzahlprozessen:**  
**Möglichkeiten und Grenzen**

**Dr. Olaf Kruse**

**KSFE 2022 - 15.09.2022**

# 1 Agenda

## Einführung

1. Einführung

2. Modellbewertung

3. Analyse der Anomalien

4. Zusammenfassung

*... keine Agenda*  
*... lediglich in Reisebericht durch die*  
*Untiefen von Big Data und ML*  
*... mit ungewissen Ausgang und einer*  
*offenen Frage*

# 1 Schadenzahlprozesse

## Einführung

### Zweistufiger Zufalls-Prozess

- ⇒ 1. Stufe: Es tritt ein (Schaden-) Ereignis ein (oder nicht) ⇒ z.B. Poisson-Prozess
- ⇒ 2. Wie „teuer“ ist das Ereignis, wenn es eintritt ⇒ z.B. Gamma-Verteilung

### Nicht untypisch

- ⇒ Online-Shopping: Kaufentscheidung ja/nein, Größe des Warenkorbbs
- ⇒ Kreditrisiko: Kreditausfall ja/nein, Ausfallvolumen
- ⇒ Versicherungen: Schadenereignisse



# 1 Big Data Einführung

*Big Data is the frontier of a firm's ability to store, process, and access all the data it needs to operate effectively, make decisions, reduce risks, and serve customers.*  
[Forrester Research]



# 1 Machine Learning

## Einführung



Mithilfe des **maschinellen Lernens** werden **IT-Systeme** in die Lage versetzt, auf Basis vorhandener Datenbestände und Algorithmen Muster und Gesetzmäßigkeiten zu erkennen und Lösungen zu entwickeln.

Die aus den Daten gewonnenen Erkenntnisse lassen sich verallgemeinern und für neue Problemlösungen oder für die Analyse von bisher unbekanntem Daten verwenden.



⇒ **Machine Learning:** Summe aller statistischen Optimierungsalgorithmen

... nur ohne Statistik

# 2 Daten- und Methoden-Setting

## Motivation

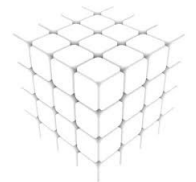
### Tarifikalkulation in der Kraftfahrtversicherung

- ⇒ Wie hoch ist der Schadenbedarf (Prämie) für ein zu versicherndes Risiko
- ⇒ Datenbestand von ca. 10 Mio. Risiken p.a.
- ⇒ ⇒ ⇒ Viele Daten und große / komplexe Modelle



### Methodischer Ansatz

- ⇒ Ordinalskaliert Merkmale, die einen n-Dimensionalen Tarifwürfel aufspannen
- ⇒ Verallgemeinerte Lineare Modelle in allen Facetten



$$P_{\lambda}(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

**„Erstes“ Tarif-Modell** (nach Tarif-Freigabe 1992)      **100.000 Tarifzellen**  
 (SF-Klasse – Tarif-Gruppe – Typ-Klasse – Regionalklasse – Garage – Fahrleistung)

**„Normales“ Tarif-Modell**      **0,35 Mrd. Tarifzellen**  
 (+ Nutzer-Alter + Kfz.-Alter bei Erwerb + Nutzerkreis + Wohngebäude + Halter ≠ VN)

**„Fortgeschrittenes“ Tarif-Modell**      **45 Mrd. Tarifzellen**  
 (+ Rabattschutz + Zahlweise + Mahnverfahren + Antriebsart)

**Unser XXL-Modell**      **9.000 Mrd. Tarifzellen**  
 (+ Dauer der Kundenbeziehung + Erst/Zweitwagen + Fahrzeugtyp (Cabrio, SUV, Sport-Coupe...))

**VST Kraftfahrt-Bestand**      **0,01 Mrd. Risiken**  
**Anzahl Schäden p.a.**      **0,001 Mrd. Schäden**

Masterarbeit:

Foundations of modern regression analysis and  
application to the **analysis of telematics data**

**Hauptkomp 2 (Itzehoer)**



**Hauptkomp 1 (Itzehoer)**



**Hauptkomp 3 (Itzehoer)**



**Hauptkomp 4 (Itzehoer)**



**Ergebnis:** relativ stabiles Modell auf Basis von 6 Telematik-Parametern  
„Vielfahrer“ – „Kurzstreckenpendler“ – „Gelegenheitsfahrer“ .....

**Wie geht das mit wenigen 1.000 Telematik-Fahrern, die keine 100 Schäden produzieren?**

- ⇒ Es bestand gar **kein** Zugriff auf Schadendaten
- ⇒ Ersatz durch die Versicherungsprämie aus dem „normalen“ Tarifmodell
- ⇒ „Selbsterfüllende Prophezeiung“

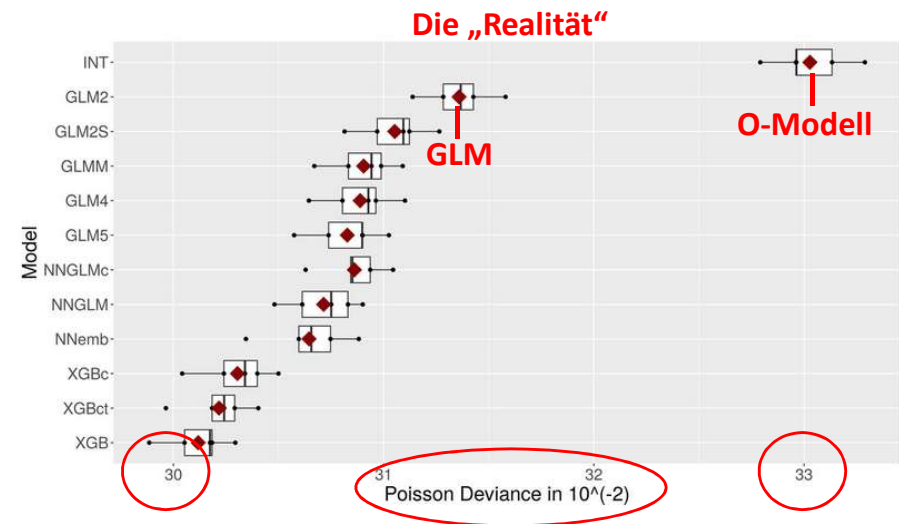
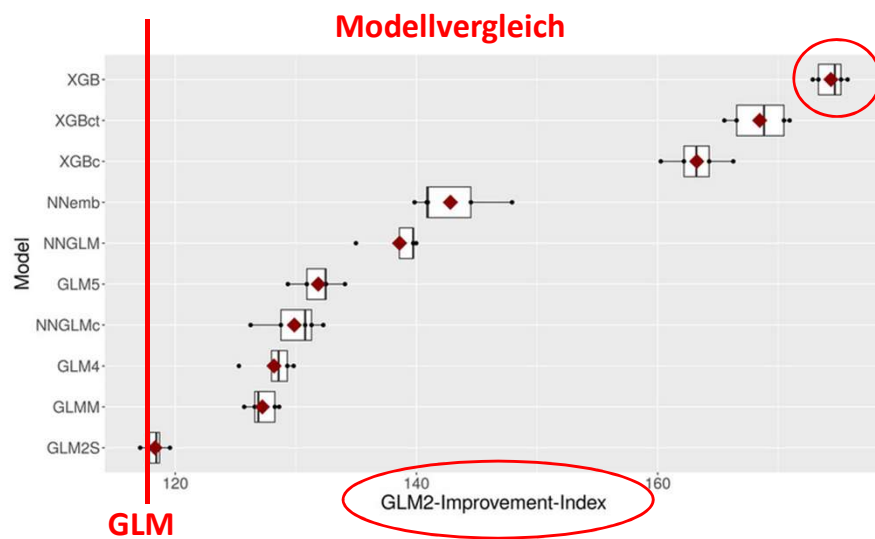
**VST**



### GLM, Neural Nets and XGBoost for Insurance Pricing

#### Verfahrensvergleich an einem frei zugänglichen kleinen KFZ-Datensatz

⇒ Klassisches Big-Data Setting: Trainings- und Test-Daten für Modellbildung und -bewertung



⇒ Moderne Verfahren passen die Daten bis zu 80% besser an, als klassische GLMs

⇒ Aber: Alle Modelle performen sehr schlecht (Erklärung von ca. 10% der Gesamt-Devianz)

# 2 Können wir Big-Data?!

## Motivation

### Entwicklung eines **maximalen Risikomodells**

- ⇒ Keine Rücksichtnahme auf **Umsetzbarkeit** in einen Tarif  
(weder rechtlich, noch technisch, noch Marktgängigkeit oder Organik)
- ⇒ Zusammenfassung **mehrerer Jahre** als Grundlage der **Modellentwicklung**  
(systematische Effekte als zeit- und zufallsstabil identifizieren)
- ⇒ Entwicklung effizienter **Modellbewertungsstrategien**



### Sukzessive Aufnahme neuer Merkmale

Basis-Modell	Modell 1	Modell 2	Modell 3	Modell 4
Erweitertes VST-Modell	+ Kfz-Alter	+ Mahnverf.	+ Erst-/Zweit-FZ	+ Antrieb/Sitzplätze

### Modellierung über verschiedene Statistikjahre

Modell 4	Modell 4	Modell 4	Modell 4	Modell 4	Modell 4
Stat.-Jahr 2018	Stat.-Jahr 2017	Stat.-Jahr 2016	Stat.-Jahr 2015	Stat.-Jahr 2014	Stat.-Jahr 2013

⇒ Unterschiedliche Sichtweise, aber vergleichbare Problemstellung

**VST**

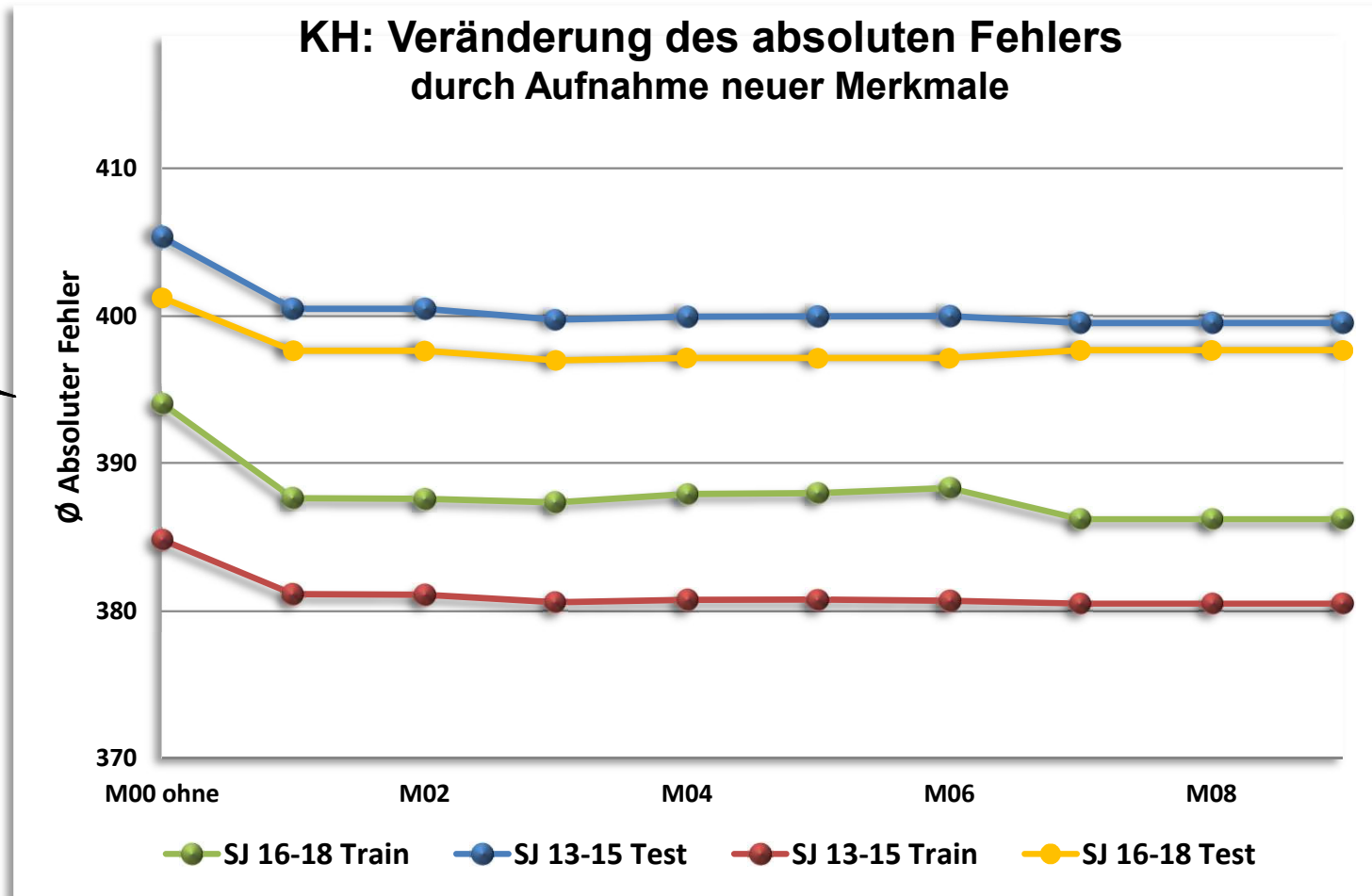
### Typische Anpassungsmaße

Anpasungsmaß	Trainingsdaten	Testdaten
-2 Log Likelihood	80,549,463	90,746,279
AIC (kleiner ist besser)	80,549,597	90,746,413
AICC (kleiner ist besser)	80,549,597	90,746,413
BIC (kleiner ist besser)	80,550,630	90,747,442
Pearson Chi-Quadrat	3,517,758,153	3,693,232,867
Pearson Chi-Quadrat / DF	1,000.9	1,020.8
Durchschnittlicher absoluter Fehler	382.7	396.4

- ⇒ Berechnung über alle Zellen jeweils für Trainings- und Test-Datensatz
- ⇒ Gleiche Aufteilung Statistikjahre 2013-2015 und 2016-2018

# 3 Modellanpassung

## Modellbewertung



⇒ Geringe Unterschiede zwischen Trainings- und Test-Daten (++)

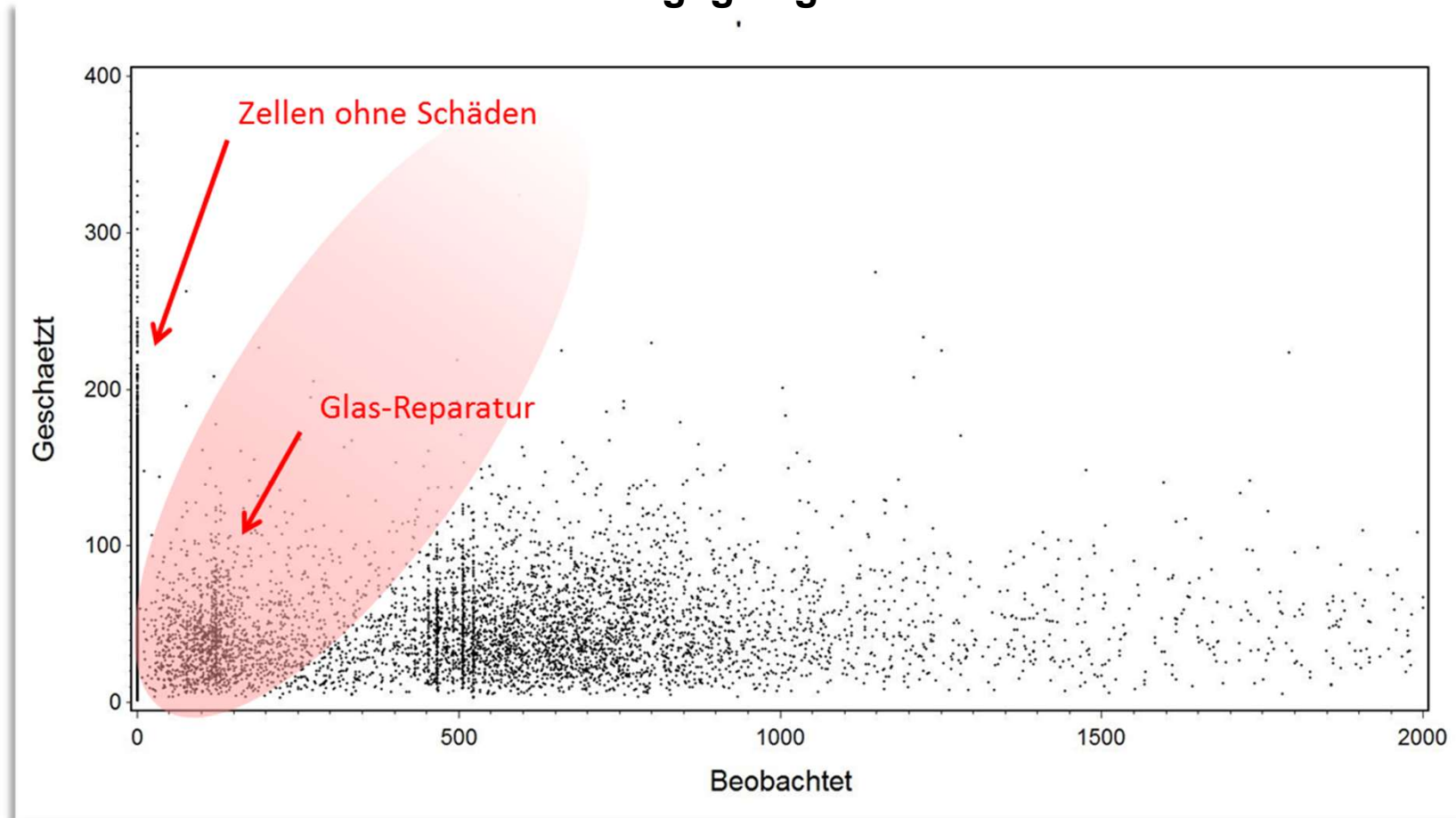
⇒ Marginale Modellverbesserung durch Aufnahme weiterer Merkmale (--)

**VST**

# 4

## Prognosevergleich Ergebnisanalyse

### TK: Beobachteter gegen geschätzter S-Bedarf



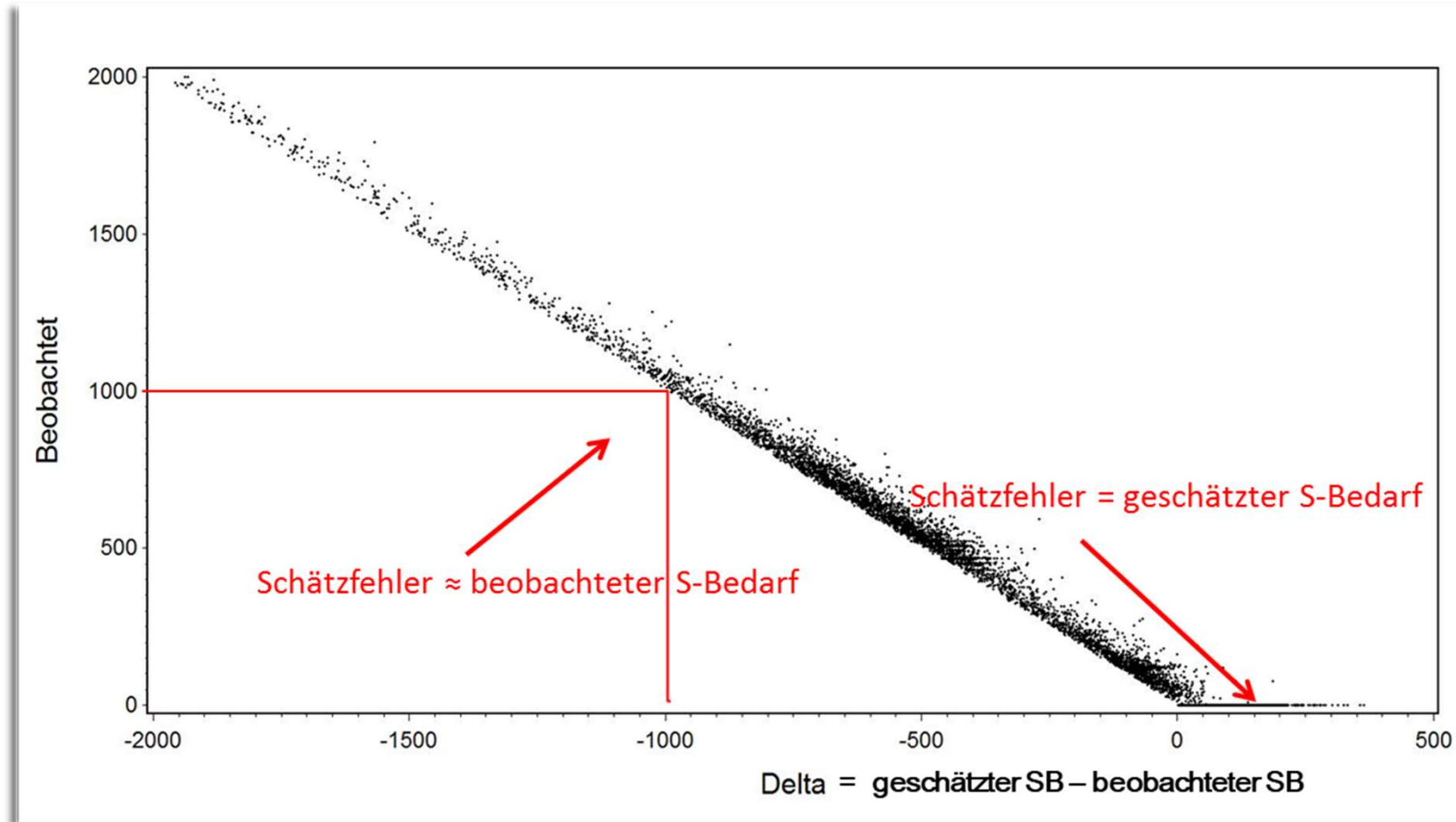
⇒ Keine positive Korrelation zwischen Schätzung und Beobachtung

# 4

## Struktur der Schätzfehler

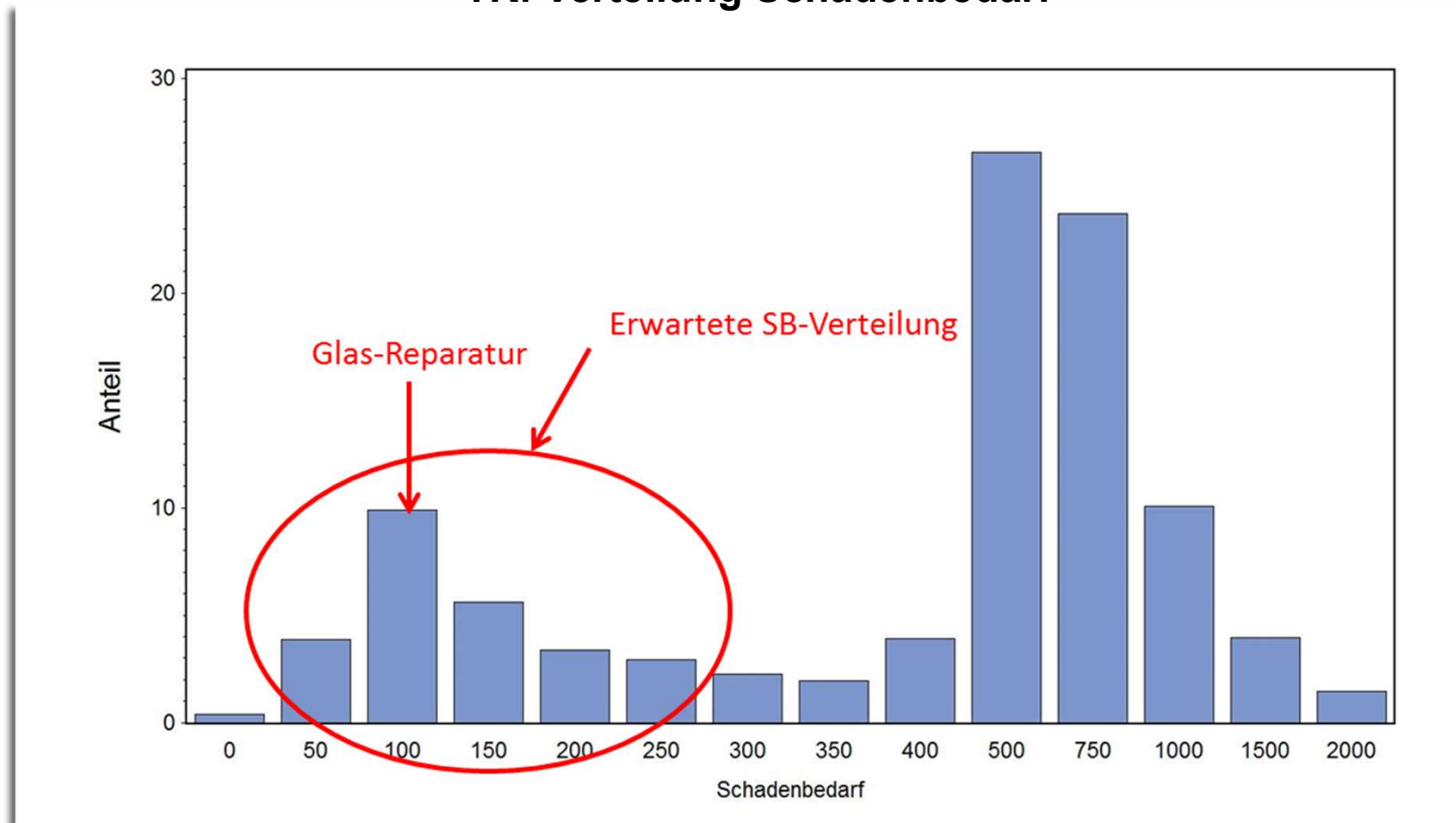
### Ergebnisanalyse

#### TK: Schadenbedarf gegen Schätzfehler



⇒ Kaum Beobachtungen mit geringem Schätzfehler

### TK: Verteilung Schadenbedarf



⇒ S-Bedarfs-Verteilung entspricht der S-Höhenverteilung

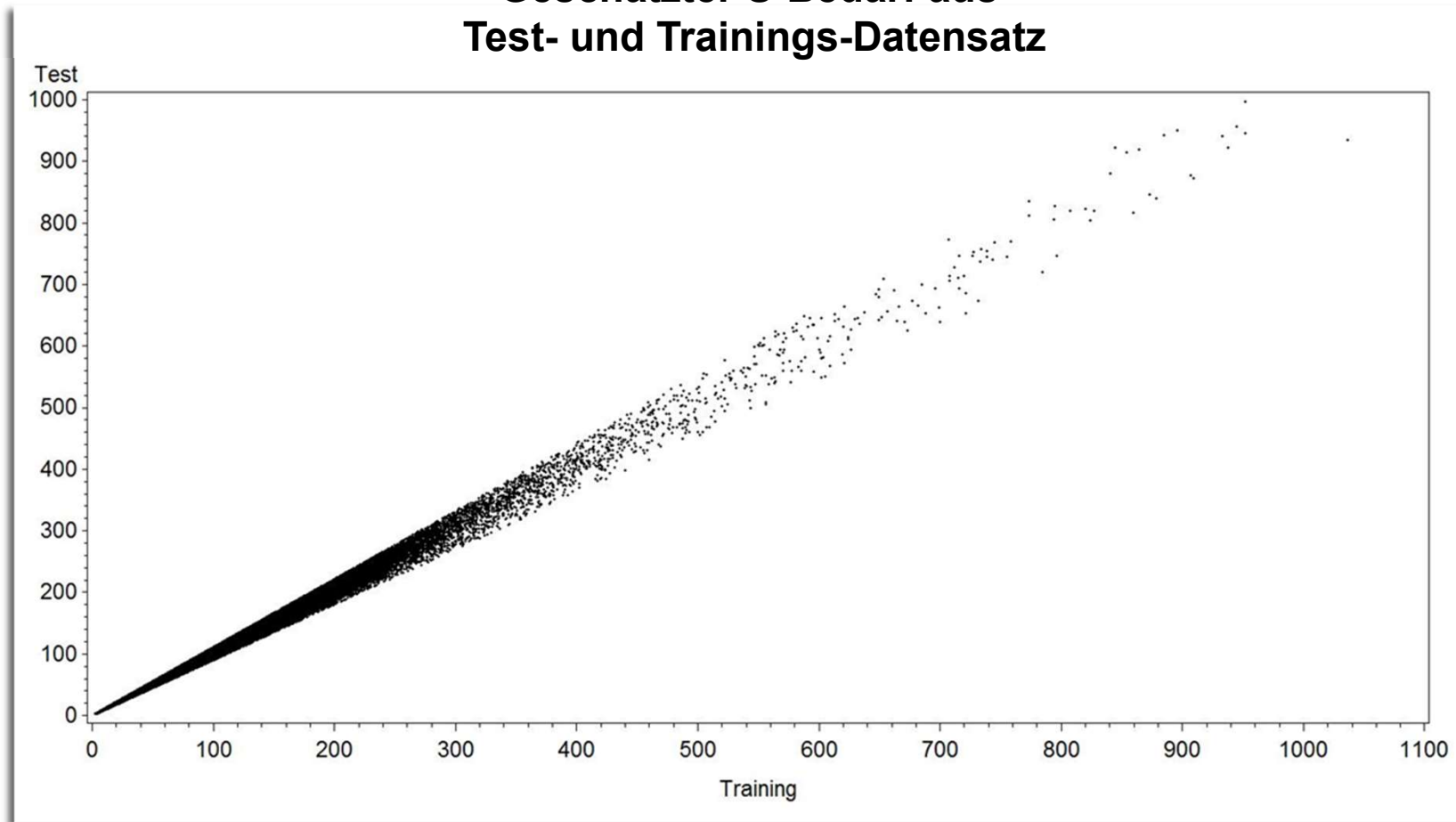


# 4

## Test- und Trainingsdaten

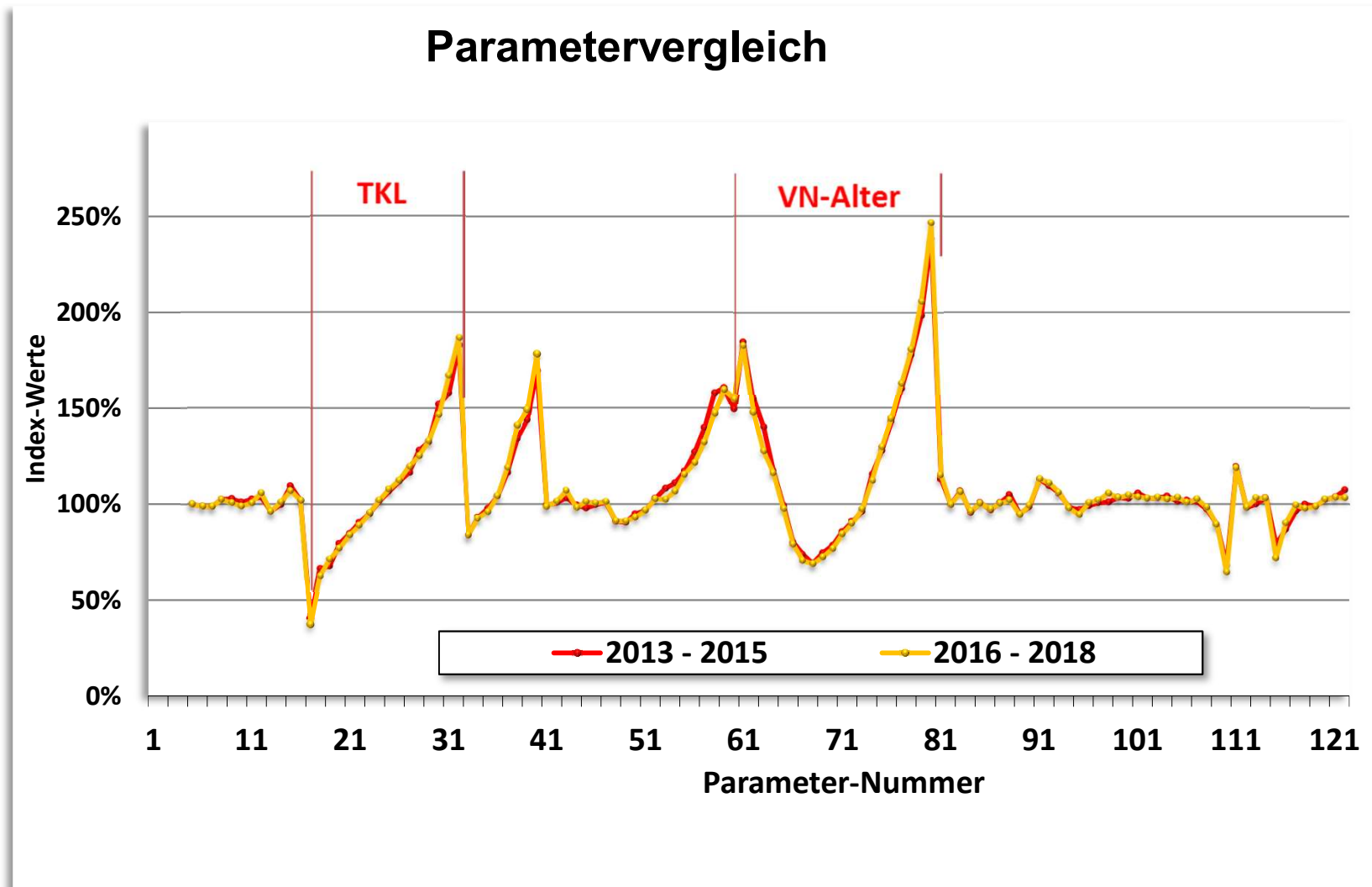
### Ergebnisanalyse

**Geschätzter S-Bedarf aus  
Test- und Trainings-Datensatz**

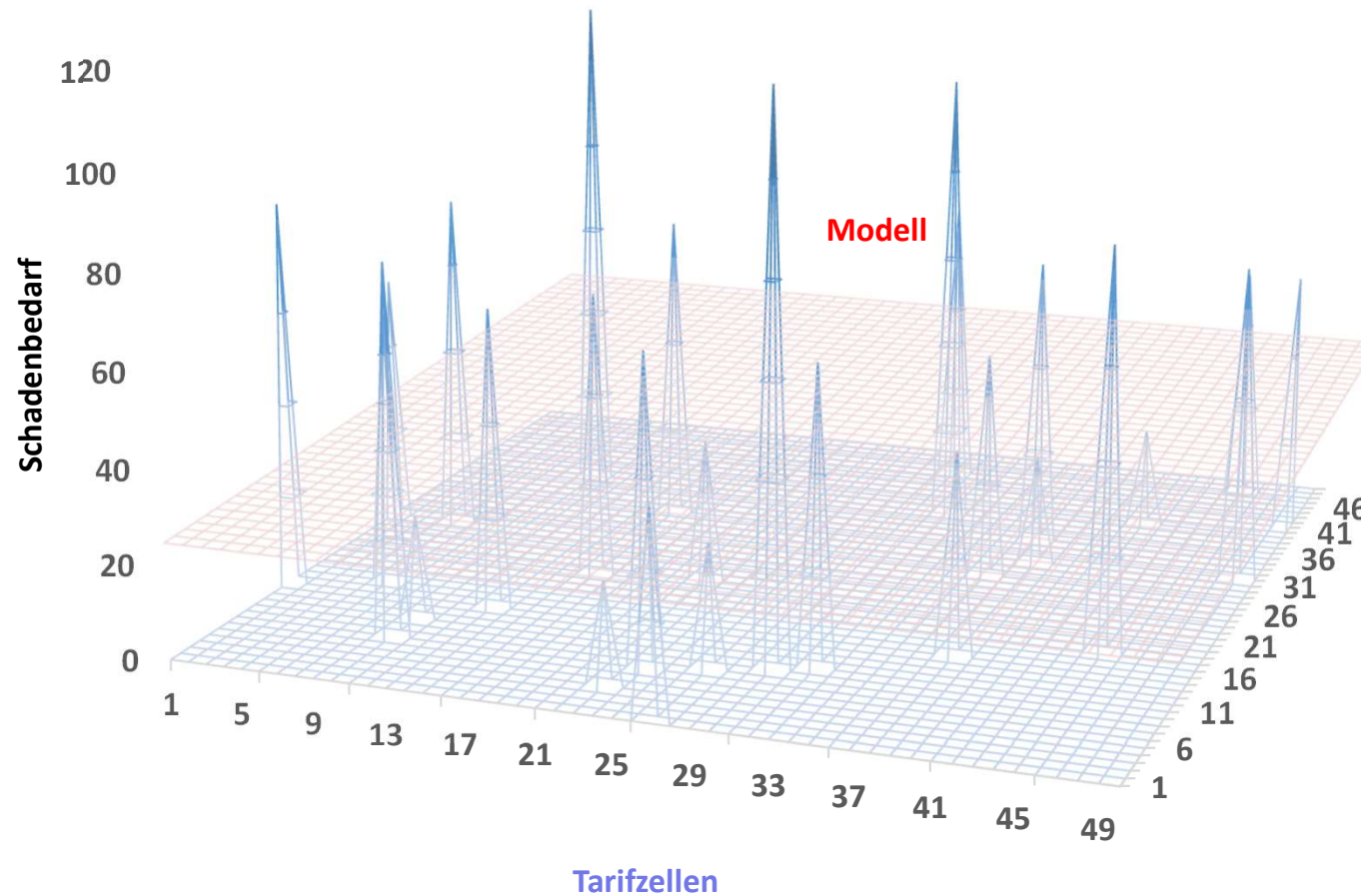


⇒ Stabile Schätzergebnisse zeigen sich auch im Parameter-Vergleich

**VST**



⇒ Trotz alledem: „Gute“ Modelle mit stabilen Modellparametern



⇒ Schadenzahlprozesse sind mit „modernen“ ML-/Big-Data-Methoden schwer greifbar

### Datensituation

- ⇒ Anzahl Tarifzellen >> Anzahl Beobachtungen >> Anzahl Schäden
- ⇒ pro Zelle **nicht aussagekräftig**/repräsentativ
- ⇒ Perfekte Anpassung an diese Daten kann **nicht Ziel** der Modellierung sein

### Modellbewertung

- ⇒ Klassische ~ bewerten Anpassung an o.g. nichtrepräsentative Datenstruktur
- ⇒ Aufteilung in Test- und Trainingsdaten hilft **nicht**
- ⇒ Gefahr von Fehlinterpretationen

### Schlussfolgerung

- ⇒ Dennoch sind die Modelle aussagekräftig und zeitstabil
- ⇒⇒⇒ Weitere Forschungsarbeit notwendig

# Fragen bitte !!

Ende

## Big Data und Machine Learning bei Schadenzahlprozessen: Möglichkeiten und Grenzen

Dr. Olaf Kruse

VST Gesellschaft für Versicherungsstatistik mbH

Roscherstr. 10

30161 Hannover

[olaf.kruse@vst-gmbh.de](mailto:olaf.kruse@vst-gmbh.de)