

## The Effects of Oversampling Non-Independent Data on Concussion Injury Risk Functions

Gunter P. Siegmund, Benjamin S. Elkin, Stephanie J. Bonin, Allen W. Yu, Adam J. Bartsch

### I. INTRODUCTION

Different injury risk functions for concussion have relied on different methods of sampling players and impacts. Based on 25 concussions in 58 individuals, Pellman [1] produced a concussion risk function that yielded 1% and 50% risks of concussion at peak resultant linear accelerations (PLA) of 5 g and 81 g, respectively. A concern about undersampling non-injurious head impacts led Funk [2] to add 27,315 non-injurious impacts and four concussions in 64 players to yield 1% and 50% risks of concussion at 109 g and 358 g, respectively. A later concern about oversampling, i.e., using multiple head impacts from a single player, led Freeman [3] to recalculate 1% and 50% risks of concussion at about 12 g and 63 g, respectively. These vastly different thresholds for concussion undermine the utility of these risk functions. Therefore, the goal of this study was to examine how oversampling non-injurious impacts and individual players affects concussion risk functions.

### II. METHODS

We used non-parametric bootstrapping to resample PLA data for injurious and non-injurious impacts from the Head Impact Telemetry System (HITS) [2][4-6] and the Prevent Biometrics Impact Monitoring Mouthguard (IMM) [7] (Fig. 1a). For HITS, we used 143 concussions in 1,417 unique players (10% injured), and for IMM we used 15 video-confirmed impacts that caused obvious performance decrements (OPDs) – treated here as injuries – in 8 of 71 unique players (11% injured). To examine the effect of oversampling non-injurious impacts, we calculated risk functions for three datasets: 50 concussions randomly sampled from the HITS concussions combined with 500, 5,000 and 50,000 non-injurious impacts (~10%, 1% and 0.1% injured, Fig. 1b-d, respectively) randomly sampled from the Funk distribution [2] (Fig. 1a). These data were considered non-independent because players could appear multiple times in the datasets.

We then used the IMM data to examine the effect of oversampling individual players. With the IMM data, we first calculated a risk function using the minimum PLAs that caused OPDs in eight unique players and the maximum PLAs that did not cause OPDs in all 71 players (including those with OPDs from other impacts) (Fig. 1e). Since each player appeared only once in the non-OPD data, we considered the data in this analysis to be independent. We then repeated this analysis but oversampled the maximum non-OPD data 10 times (Fig. 1f). To then compare the HITS and IMM datasets, we calculated risk functions for two more IMM datasets: (i) all 15 OPD and all 1,432 non-OPD impacts from the available data set (Fig. 1g); and (ii) all 15 OPD events and 10 times the number of non-OPD events (Fig. 1h). These latter two analyses had injury/non-injury ratios (C/N ratios in Fig. 1) of ~1% and ~0.1%, similar to the latter two HITS analyses.

For all bootstrap analyses, the number of injuries and non-injuries remained constant within a given analysis. Two risk functions (logistic regression, LR; and consistent threshold, CT) were calculated for each of the 500 iterations of each dataset, and the median and 95<sup>th</sup> percentile confidence interval (CI) of the risk functions were then extracted (Fig. 1b-h). The 1% and 50% injury risks based on the median risk function were also calculated.

### III. RESULTS

The LR and CT risk functions for the HITS data shifted to the right as oversampling of the non-injurious data increased (C/N ratio decreased) (Fig. 1b-d). For LR, the 50% risk increased from 90 g to 207 g, which exceeded the maximum PLA value in both the injury and non-injury data. In contrast, the 50% risk for CT increased from 79 g to 157 g and, by definition, remained within the data. Both risk functions for the independent IMM data indicated a 50% risk of OPD at 59 g (Fig. 1e). When non-OPD data were oversampled, the 50% risks increased to 72 g and 62 g for the LR and CT functions, respectively (Fig. 1f). Using all of the IMM data, i.e., over-sampling the

G. P. Siegmund (e-mail: gunter.siegmund@meaforensic.com; tel: +1 604 277 3659) is the Director of Research at MEA Forensic Engineers & Scientists in Richmond, BC, Canada and an Adjunct Professor at the UBC School of Kinesiology. B. S. Elkin is a Biomechanical Engineer at MEA Forensic in Toronto, ON, Canada. S. J. Bonin is a Senior Engineer and A. W. Yu is a Biomechanist at MEA Forensic in Los Angeles, CA, USA. A. J. Bartsch is the Chief Science Officer at Prevent Biometrics in Minneapolis, MN, USA.

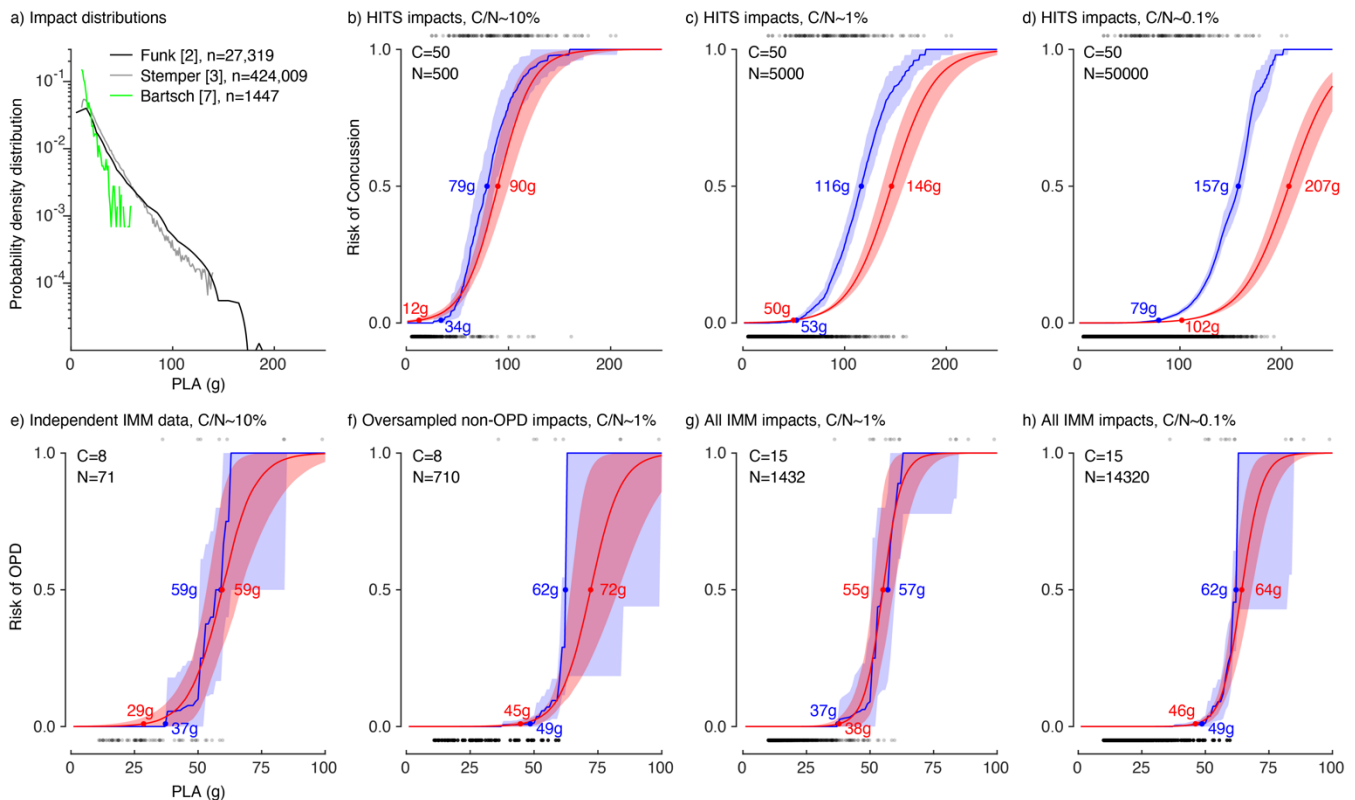


Fig. 1. a) Three impact distributions for peak linear acceleration (PLA) [2-3][7]; b), c) and d) injury risk functions based on resampled HITS data with 50 concussive impacts (C) and 10, 100 and 1,000 times more non-concussive impacts (N); e) injury risk function based on the independent IMM data, i.e., the minimum PLA of eight unique players with obvious performance decrements (OPDs) and the maximum PLA for all players with no OPDs; f) same as e) but with 10 times as many non-OPD impacts resampled from the maximum PLA data; g) and h) injury risk functions based on 15 OPD events resampled from all 15 OPD events and about 100 and 1,000 times as many non-OPD events resampled from all 1,432 actual non-OPD events. Red=logistic regression, blue=consistent threshold; solid line= median, shaded region=95<sup>th</sup> percentile confidence interval; 50% and 1% PLA risks shown; gray dots show sample data from one bootstrap iteration.

injury and non-injury data, the 50% risks reduced to 55 g and 57 g for the LR and CT functions, respectively, the confidence intervals narrowed and the 1% risks for both functions aligned. Using 10 times the non-OPD data increased the 50% risks to 64 g and 62 g for the LR and CT functions, respectively.

#### IV. DISCUSSION

Oversampling non-independent, non-injurious impacts when calculating injury risk functions artificially shifts the risk curve and underestimates injury risk for a given exposure level. This phenomenon was larger in the HITS data than in the IMM data, possibly because the separation between the injured and non-injured data was poorer for HITS than for IMM, which in turn may be related to HITS' poorer measurement precision [7-8]. Even with better separation of the injury and non-injury data in the IMM data, oversampling non-independent data shifts the risk functions and/or narrows the confidence intervals compared to the independent IMM dataset (Fig. 1e), though to a lesser extent than for HITS. Both effects should be considered when applying risk functions to real-world exposures. Our analysis focused on oversampling and ignored other possible issues, like injury underreporting, diagnostic accuracy, data censoring and repetitive impacts, all of which need to be explored in more depth. It is also unclear why the range of PLA values for HITS and IMM were so different, though it may be related to different injury metrics: diagnosed concussions for HITS versus video observations of altered player behaviour for IMM. Despite these limitations, our results show that concussion injury risk functions should avoid using non-independent data, particularly with poorly separated injury and non-injury data.

#### V. REFERENCES

- [1] Pellman, *Neurosurg*, 2003.
- [2] Funk, AAAM, 2007.
- [3] Freeman, *J Forensic Biomed*, 2018.
- [4] Greenwald, *Neurosurg*, 2008.
- [5] Rowson, *Ann Biomed Eng*, 2012.
- [6] Stemper, *Ann Biomed Eng*, 2018.
- [7] Bartsch, *Ann Biomed Eng*, 2020.
- [8] Siegmund, *Ann Biomed Eng*, 2016.