# UC Berkeley
## Research Reports

**Title**
Evaluating Research on Data Linkage to Assess Underreporting of Pedestrian and Bicyclist Injury in Police Crash Data

**Permalink**
https://escholarship.org/uc/item/0jq5h6f5

**Authors**
Doggett, Sarah
Ragland, David R.
Felschundneff, Grace

Peer reviewed

1  **Evaluating Research on Data Linkage to Assess Underreporting of Pedestrian**
2  **and Bicyclist Injury in Police Crash Data**
3
4  **Sarah Doggett, Corresponding Author**
5  Graduate Student, Safe Transportation Research and Education Center (SafeTREC)
6  University of California, Berkeley
7  2614 Dwight Way #7374
8  Berkeley, CA 94720-7374
9  Phone: 510-642-0655, Fax: 510-643-9922
10  doggett_sarah@berkeley.edu
11
12  **David R. Ragland, PhD, MPH**
13  Adjunct Professor Emeritus, School of Public Health
14  Director, Safe Transportation Research and Education Center (SafeTREC)
15  University of California, Berkeley
16  2614 Dwight Way #7374
17  Berkeley, CA 94720-7374
18  Phone: 510-642-0655, Fax: 510-643-9922
19  davidr@berkeley.edu
20
21  **Grace Felschundneff**
22  Senior Editor, Safe Transportation Research and Education Center (SafeTREC)
23  University of California, Berkeley
24  2614 Dwight Way #7374
25  Berkeley, CA 94720-7374
26  Phone: 510-642-0655, Fax: 510-643-9922
27  gracefelschundneff@gmail.com
28
29
30  **Word Count:** 4,426 words + 3 tables/figures (250 words each) = 5,176 words
31
32  **Submission Date:** July 30, 2018

1   **ABSTRACT**
2   Traffic safety decisions are based predominantly on information from police collision reports.
3   However, a number of studies suggest that such reports tend to underrepresent bicycle and
4   pedestrian collisions. Underreporting could lead to inaccurate evaluation of crash rates and may
5   under- or overestimate the effects of road safety countermeasures. This review examined ten
6   studies that used data linkage to explore potential underreporting of pedestrian and/or bicyclist
7   injury in police collision reports. Due to variations in definitions of reporting level, periods of
8   study, and study locations, it was difficult to directly compare the studies. Even among the six
9   studies using the hospital link definition, estimates of reporting levels ranged from 44 to 75
10  percent for pedestrian crashes, and from 7 to 46 percent for bicycle crashes, suggesting a severe
11  underreporting problem. However, few of the studies provided estimates of the error around their
12  reporting level estimates, and as a result, it is difficult to determine the true level of
13  underreporting. It may be that bicycle and pedestrian crashes appear in both police and hospital
14  datasets but are less likely to be linked. Due to linkage error, link rate can only be used to
15  *estimate* reporting level. Without the *variance* of that estimate, the effect of underreporting on
16  traffic safety analyses cannot be accurately determined. Future studies should include estimates
17  of the error present in their data linkage process for greater accuracy of the underreporting in
18  police data. Datasets should be designed for easier linkage with hospital data and other datasets.
19
20  *Keywords*: Data, Linkage, Pedestrian, Bicycle, Crash, Underreporting
21
22

1   **INTRODUCTION**
2   Traffic safety decisions are almost universally based on information from police collision
3   reports. However, many researchers believe that police collision reports have limitations, and fail
4   to include all crashes that occur on the road. For example, a number of studies suggest that police
5   collision reports underrepresent bicycle and pedestrian collisions.
6        According to the Federal Highway Administration (FHWA), traditional crash data
7   sources are insufficient because they exclude both crashes that take place in non-roadway
8   locations (e.g., parking lots, driveways, and sidewalks) and bicycle crashes and pedestrian
9   injuries that do not involve motor vehicles *(1)*.
10       Ideally, the degree of bias in police collision reports could be measured and accounted for
11  in analyses *(2)*. This process requires both the reporting level and the uncertainty in the estimate
12  of the reporting level to be known *(3)*. Although researchers have estimated reporting level by
13  linking police collision reports with hospital data, few have determined the level of uncertainty
14  surrounding their estimates.
15       This paper summarizes existing research using data linkage to study pedestrians and
16  bicyclists, explores potential problems concerning linkage, and offers suggestions on how to
17  improve the data linkage process.
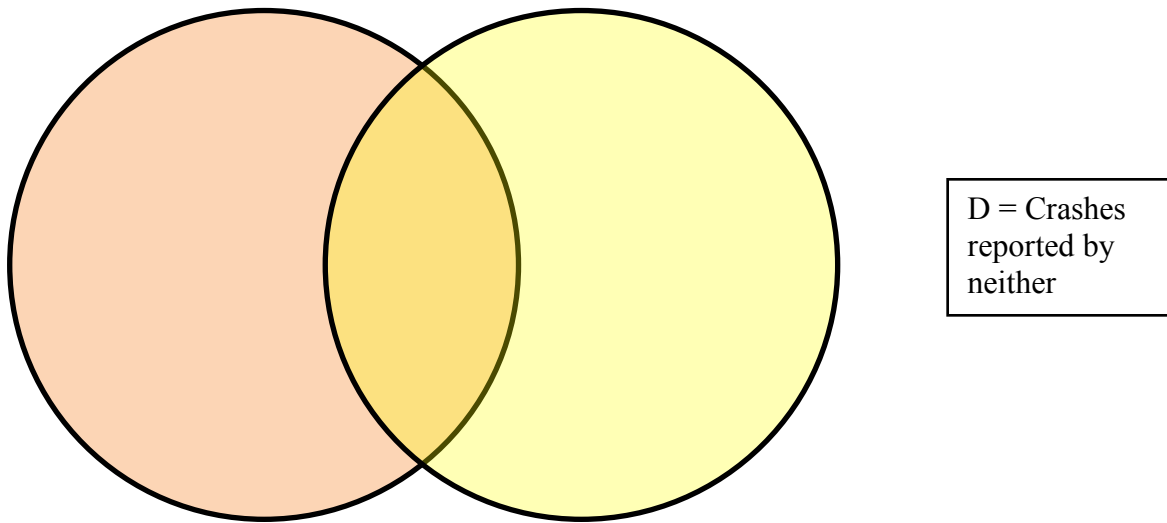18
19  **METHODOLOGY AND DEFINITIONS**
20  Searches using Google Scholar and Science Direct were conducted to identify potentially
21  relevant articles, published between 1999 and 2017, for the literature review. Search terms
22  included record linkage, data linkage, crash data linkage, pedestrian underreporting,
23  underreporting, crash injury assessment, and police crash data limitations.
24       The abstracts of the articles returned through this search were evaluated to assess their
25  relevancy to the general topic of crash reporting and heath data linkage. Articles with relevant
26  abstracts were read in full to determine whether they reported findings related to pedestrians
27  and/or bicyclists.
28     Five of the articles explored data linkage in general but mentioned findings specific to
29  pedestrians and/or bicyclists. Two focused exclusively on bicyclists, two exclusively on
30  pedestrians, and one on both pedestrians and bicyclists.
31       Most of the studies included in the literature review defined underreporting as hospital
32  records without counterparts in police crash reports. However, Janstrup et al. found that there are
33  also police crash reports that do not have counterparts in hospital records *(4)*. This expands the
34  definition of underreporting to crashes that exist only in one dataset. This concept is illustrated in
35  Figure 1.
36

D = Crashes reported by neither

A = Crashes reported by police

Figure 1  **Patterns of Reporting Pedestrian and Bicyclist Crashes by Source**

Some crashes may be reported in both datasets but cannot be linked due to errors in data recording—these crashes are *misreported* but cannot not be distinguished from underreported crashes using current linkage techniques. There are also crashes that are *not reported* to either the police or to hospitals (D in Figure 1), due to lack of serious injury or an unwillingness on the part of those involved in the crash to contact the authorities. To identify these unreported crashes, other datasets such as those collected by insurance companies can be analyzed. One of the studies addressed this, and found that relying solely on police and hospital injury data can result in a significant underestimation in the actual number and severity of crashes *(5)*.

Several different definitions of reporting level exist. Using the terminology described in Table 1, these definitions are explained as follows:

- A: Cases only appearing in police crash reports
- B: Cases appearing in both police crash reports and hospital records
- C: Cases only appearing in hospital records
- D: Cases that do not appear in either police crash reports or hospital records

**Table 1  Definitions of Reporting Level**

| Definition of Reporting Level | Formula | Question Addressed |
|---|---|---|
| 1: Hospital Link Rate | B/(B+C) | Proportion of cases reported in hospital records for which a police crash report is found |
| 2: Police Link Rate | B/(A+B) | Proportion of cases reported in police crash reports for which a hospital record is found |
| 3: Police Ascertainment Rate | (A+B)/(A+B+C) | Proportion of identified cases accounted for in police crash reports |
| 4: Hospital Ascertainment Rate | (B+C)/(A+B+C) | Proportion of identified cases accounted for in hospital records |

Several studies used the capture-recapture method to determine reporting level, which is not derived from those shown in Table 1. Originally, this method was used to estimate animal populations based on how many animals were captured and then later recaptured in a least two

1   random samples *(6)*. However, researchers have noted that this method may not be appropriate
2   for estimating reporting levels because some of its underlying assumptions, such as a closed
3   study population and an equal probability for each individual to be captured in each sample, may
4   be violated *(6,4)*.
5      The definition used can significantly affect the estimated reporting level *(2)*. As an
6   example, assume that 50 cases only appear in police reports (A), 20 cases appear in both police
7   and hospital reports (B), and 30 cases only appear in hospital reports (C). Under this scenario,
8   hospital link rate would be 0.4, police link rate would be 0.29, police ascertainment rate would
9   be 0.7, and hospital ascertainment would be 0.5. While definitions 3 and 4 are theoretically the
10   most correct definitions of reporting level, few studies apply these when estimating reporting
11   level.
12
13   **LITERATURE REVIEW**
14   Due to the different definitions of reporting level, different periods of study, and different study
15   locations, it is difficult to directly compare the ten studies. As shown in Table 2, six of the ten
16   studies used hospital link rate as their definition of reporting level, while two used capture-
17   recapture, two used police link rate, two used police ascertainment rate, and one used hospital
18   ascertainment rate (this adds up to more than ten because the article by Short and Caufield
19   reported on four different definitions). Even among the six studies using the hospital link
20   definition, estimates of reporting levels ranged from 44 to 75 percent for pedestrian crashes, and
21   from 7 to 46 percent for bicycle crashes.
22      However, the articles all agreed that police collision reports have limitations and
23   underreport certain types of crashes, especially those involving pedestrians and bicyclists. For
24   example, Stutts and Hunter *(9)* found that police were unlikely to be contacted about pedestrian
25   incidents not involving motor vehicles, although more than two-thirds of serious injuries fell into
26   this category. Sciortino *(7)* found that pedestrian injury victims who were African American,
27   male, or sustained minor injuries were less likely to be included in police crash records. The
28   author suggested that such bias was likely due to the reluctance on the part of pedestrians to
29   summon the police if the police were not initially present at the crash scene. Tarko and Azam *(8)*
30   found that pedestrians were less likely to be included in the linked database if they were struck
31   by vehicles on state roads, at Y intersections, or on divided roadways. However, they were more
32   likely to be included in the database if they were struck while crossing a road instead of walking
33   along its edge or standing outside of the roadway.
34      Drivers are more likely to be included in crash reports than are cyclists *(9)*. Langley *(10)*
35   found that only 22 percent of bicyclist crashes on public roads could be linked to a crash
36   report—the percentage increased to 54 percent when it included crashes that also involved motor
37   vehicles. Because crash reports are mainly focused on incidents that involve motor vehicles and
38   which occur on public roadways, they likely capture fewer than one-third of bicyclist injury
39   cases serious enough to require medical treatment *(1)*. According to Janstrup et al. *(4)*, the
40   underreporting rate for crashes involving cyclists in Denmark was 14 percent for those resulting
41   in serious injuries, and 7 percent for those resulting in slight injuries.
42      The literature identified several possible reasons for underreporting. According to
43   Langley *(10)*, motor vehicle crashes listed in hospital data are assumed to have occurred on
44   roadways unless another location is specified, thus hospital derived estimates may overstate the
45   number of motor vehicle crashes on public roads and understate crashes that occur in other
46   locations or those that do not involve motor vehicles. Another possible reason for underreporting

1    is that pedestrian and bicycle crashes are less likely to result in insurance compensation than are
2    motor vehicle crashes. Lujic et al. *(11)* found that those entitled to insurance payouts had higher
3    linkage rates than those who were not, presumably because police reports were required as a
4    condition for receiving compensation. According to Watson et al. *(12)*, it is possible that the
5    severity of injuries resulting from a collision with another vehicle is likely to be more serious,
6    and therefore more likely to be reported.
7           Several studies have shown that compared with hospital data, police crash reports do not
8    accurately report injury severity. Injury classifications in crash reports are usually based on the
9    KABCO scale, which is less nuanced than the Injury Severity Score that many hospitals use *(13)*.
10   Additionally, KABCO classifications are made at the scene of the crash by officers who typically
11   lack medical training—therefore, less visible but life-threatening injuries such as internal
12   bleeding may be misclassified as non-severe, while more obvious minor injuries such as minor
13   lacerations with profuse bleeding may be misclassified as severe *(13)*. Another problem related
14   to estimation of injury severity, is that the injury classification in police reports is static and does
15   not necessarily reflect subsequent developments *(14)*.
16
17   **Table 2  Summary of Data Linkage Articles Relevant to Pedestrian and Bicyclist Safety**
18

| Source | Study Period | Study Location | Focus | Definition Used | Findings/Conclusions |
|---|---|---|---|---|---|
| Conderino, Fung, Sedlar & Norton, 2017 | 2009-2015 | New York City | General | Hospital Link Rate | • 50% of hospital reports involving a pedestrian crash linked to a police report<br>• 45% of hospital reports involving a bicyclist crash linked to a police report<br>• Sensitivity - 74%<br>• Specificity - 93% |
| Janstrup et al., 2016 | 2003-2007 | Funen, Denmark | General | Capture-Recapture | • Compared with car occupants, pedestrians are more likely to appear in both police and hospital databases; bicyclists are more likely to appear in either<br>• Only 7% of bicycle crashes resulting in slight injury and only 15% of bicycle crashes resulting in severe injury are reported by the police |
| Langley, 2003 | 1995-1999 | New Zealand | Bicyclist | Hospital Link Rate | • Only 22% of bicycle crashes on public roads could be linked to a crash report<br>• When limited to bicycle crashes on public roads involving motor vehicles, 54% could be linked to a crash report |
| Lujic, Finch, Boufous, Hayen & Dunsmuir, 2008 | 2000-2001 | New South Wales | General | Hospital Link Rate | • 69% of road traffic crashes were linked to police records<br>• Drivers were most likely to have their hospital records linked to police records (83%)<br>• 46% of bicyclist crashes and 75% of pedestrian crashes were linked to police records<br>• Authors hypothesized that underreporting for cyclists is due to "ambiguity of…laws and regulations" and the fact that cyclists are "less likely to cause property damage" |
| Sciortino, Vassar, Radetsky & Knudson, 2005 | 2000-2001 | San Francisco | Pedestrian | Police Ascertainment Rate | • Police reports underestimate the number of pedestrian injuries by 21% (e.g., reporting level is 79%)<br>• African-American pedestrians were less likely than white pedestrians to be linked to a police report<br>• Women were more likely than men to be linked to a police report |

| Source | Study Period | Study Location | Focus | Definition Used | Findings/Conclusions |
|---|---|---|---|---|---|
| Short & Caulfield, 2016 | 2005-2011 (Police and Hospital Data) 2010-2011 (Injuries Board Data) | Ireland | General | Hospital Link Rate, Police Link Rate, Hospital Ascertainment Rate, Police Ascertainment Rate | • For pedestrian injuries, 28.9% of police records were matched with a hospital record; 44.3% of hospital records were matched with a police record<br>• For bicyclist injuries, 24.8% of police records were matched with a hospital record; 8.2% of hospital records matched with a police record<br>•Police Ascertainment Rate was 73.4% for pedestrians and 26.4% for bicyclists<br>•Hospital Ascertainment Rate was 47.7% for pedestrians and 80.2% for bicyclists<br>• False Positive Rate – 13%<br>• False Negative Rate – 13% |
| Stutts & Hunter, 1999 | 1995-1996 | Various locations in California, New York, and North Carolina | Pedestrian and Bicyclist | Hospital Link Rate | • 70% of bicyclist injuries reported by hospitals did not involve a motor vehicle<br>• 64% of pedestrian injuries reported by hospitals did not involve a motor vehicle<br>• 31% of bicyclist and 53% of pedestrian injuries occurred in non-roadway locations (e.g., sidewalks, parking lots, trails)<br>• Police crash reports capture less than 33% of serious bicyclist injuries |
| Tarko & Azam, 2011 | 2003-2008 | Indiana | Pedestrian | Police Link Rate | • Pedestrians struck on a state road, at a Y intersection, or on a divided roadway were less likely to be included in both police and hospital databases<br>• Pedestrians struck while crossing a road, as opposed to walking along or standing outside of the roadway were much more likely to be included in both databases<br>• Authors hypothesized that more severe injuries were more likely to appear in both databases |
| Tin Tin, Woodward & Ameratunga, 2013 | 2006-2011 | New Zealand | Bicyclist | Capture-Recapture | • Police reports were linked to insurance, hospital, and mortality records<br>• 13% of hospital reported crashes and 64% of hospital reported collisions were linked to police records<br>• 39% of police reported crashes and 43% of police reported collisions were linked to hospital records<br>• When compared with self-reported data from the cyclists, the entire linked dataset had a sensitivity of 63.1% and specificity of 93.5%<br>• The collision-only dataset showed a 40.0% sensitivity and a 99.9% specificity |
| Watson, Watson & Vallmuur, 2015 | 2009 | Queensland | General | Hospital Link Rate | • The study used discordance rates between police data and hospital records to measure underreporting of crashes to/by the police<br>• The discordance rate was 44% for pedestrians and 93% for bicyclists (e.g., a reporting level of 56% and 7% respectively)<br>• Authors hypothesized that bicyclist injures are not reported to the police because injuries are generally less serious, are less likely to have insurance implications, and are more likely to involve young people, who generally have high discordance rates |

1

2 **POTENTIAL IMPLICATIONS OF STUDY RESULTS**
3 Findings from these ten studies indicate that there is a severe underreporting problem in datasets
4 based only on police collision reports. Underreporting can result in the forecasting of incorrect
5 estimates of crash and fatality rates, and identification of erroneous factors responsible for crash
6 occurrence, thereby making the entire road safety exercise ineffective, according to Ahmed,

1    Sadullah, & Yahya *(15)*. When road safety is evaluated based on data other than the actual
2    number of crashes that occurred, there is a tendency to mistake trends in crash reporting with
3    trends in traffic safety *(3)*. In addition, the authors found that inaccurate crash data can result in
4    improper prioritization of funding and resources, and under- or overestimation of injury severity
5    risk. For example, one study found that estimates based on underreported police-reported crash
6    data minimize the effectiveness of seat belt use in injury severity risk and could have serious
7    policy implications *(16)*.
8        In cases of crash underreporting, analysis relying on police data may be biased, according
9    to Janstrup et al. *(4)*. Reliance on hospital data may also be problematic, as Watson et al. *(12)*
10   reported that the level of underreporting varied depending on the data with which the police data
11   was linked. Watson found that when hospital data was examined, approximately two thirds of the
12   data were not linked to police data. Similarly, when Short and Caulfield *(5)* added injury claims
13   data to their analysis, the total number of identifiable injuries was found to be more than three
14   times greater than the number identified by police reports, and five times greater than the number
15   identified in hospital records.
16       Because of the inaccuracy of injury classifications in police collision reports, injury
17   severity cannot be used to match datasets. In addition, by only using injury data from police
18   reports, financial costs to crash victims or to the public for health care associated with a crash
19   cannot be easily and accurately determined *(14)*.
20
21   **LIMITATIONS OF EXISTING STUDIES**
22   Hauer and Hakkert *(3)* proposed methods to account for underreporting in police crash reports.
23   They argue that the variance of the estimated number of crashes that occur can be calculated if
24   the following factors are known: reported number of crashes, reporting level, and the variance of
25   the reporting level. While the studies included in this literature review have attempted to
26   establish the reporting level associated with various types of crashes, few have reported on the
27   error surrounding their estimates.
28       In most real-world cases, true match status is unknown and link status is used as a proxy.
29   Under the presence of a perfect matching process, link status and match status would be identical
30   and link rate would be equivalent to the reporting level. However, there are two sources of error
31   in the process of data linkage—false positives and false negatives. False positives are linked
32   records that do not belong to the same person/event *(17)*. False positives are more likely to occur
33   when identifiers are not discriminative and when files are large *(18)*. False negatives, also known
34   as *missed matches*, are records that belong to the same person/event but that are not linked *(17)*.
35   This occurs when records have inaccurate or missing data *(18)*.
36       Because linkage and analysis processes are often separated due to privacy concerns,
37   researchers using linked datasets may be unable to determine the extent of bias that linkage
38   errors have introduced into their study *(19)*.
39       It is difficult to predict the direction and magnitude of bias resulting from linkage errors
40   due to the "distribution of errors with respect to variables of interest" which is usually unknown
41   *(19)*. Missed matches reduce sample size and statistical power and can lead to under-
42   underestimates of exposures or outcomes if the linkage is informative *(19)*. Because missed
43   matches do not necessarily occur at random, the linked data may not accurately represent the
44   study population, reducing the viability of the research effort and may introduce bias *(19)*.
45   Selection bias can occur if an individual's presence in the linked dataset is related both to
46   exposure and an outcome of interest *(19)*.

1    When data for key variables are missing, cases that are linked are less likely to be
2    representative of the study population because they have fewer common values such as unusual
3    zip codes *(17)*. Unfortunately, missing data is common in data linkage. The individual datasets
4    used for linking—police, hospital, insurance, traffic—are generally not developed for research
5    purposes. The intended use of a dataset largely determines its collection method *(20)*. Therefore,
6    the contents and details of their attributes, and the application of various datasets for purposes
7    other than those for which they were designed may result in decreased data quality. Additionally,
8    many variables that would aid in correctly linking the datasets are often redacted for privacy
9    purposes (i.e., name, address, social security number). Data may be inaccurate due to
10   misspellings or lack of information (i.e., staff might guess the age of an unconscious patient and
11   set the birth month and day to '01') *(21)*. Finally, the method of data collection can influence
12   data quality because errors are more likely to occur when data is copied from handwritten forms
13   or transcribed from conversations than when the information is entered directly into a database
14   *(18)*.
15        While there are statistical methods for managing data that are completely or partially
16   missing at random, it is likely that missing data in crash and hospital data linkages are usually
17   missing *not* at random, which introduces systemic bias into the analysis of the linked dataset
18   *(17)*. Studies of linked databases may be significantly impacted by linkage errors, especially if
19   certain types of people or events are more or less likely to have the outcome of interest *(17)*.
20   Often, the bias introduced by these errors cannot be determined because the errors themselves
21   are unknown.
22        Linkage error can be estimated by using gold standard data, which is data with a known
23   true match status. These data are rarely available in the real world, although sometimes synthetic
24   datasets are used instead. To estimate linkage error, gold standard data is matched using the same
25   process as the study datasets. This allows link status to be compared with match status and for
26   the calculation of metrics such as sensitivity and specificity.
27        Sensitivity, or the true positive rate, is the proportion of matches that are correctly
28   identified as links. Specificity, or the true negative rate, is the proportion of non-matches that are
29   correctly identified as non-links. High sensitivity reduces the number of false negatives, while
30   high specificity reduces the number of false positives. Often, there is a tradeoff between
31   sensitivity and specificity.
32        Although sensitivity and specificity can be affected by the method of data linkage used,
33   they are much more dependent on the quality of the data that is being linked. Data quality is
34   affected by the accuracy, completeness, consistency, timeliness, accessibility, and believability
35   of the variables used to link the databases *(21)*.
36        Of the ten studies included in the present review, only two report the sensitivity and
37   specificity of the data linkage process. Both showed a higher specificity than sensitivity, which
38   indicates that there are more false negatives than false positives in their datasets. One reported
39   the false positive and false negative rates, which are similar though not interchangeable metrics.
40        Few of the studies provide an estimate of the error around their reporting level estimates,
41   thus it is possible that bicycle and pedestrian crashes are not as underreported as these studies
42   suggest. It may be that bicycle and pedestrian crashes appear in both police and hospital datasets
43   but are less likely to be linked. Because of linkage error, link rate can only be used to *estimate*
44   reporting level. Without knowledge of the *variance* of that estimate, the effect of underreporting
45   on traffic safety analyses cannot be accurately determined.
46

1  **SUGGESTIONS FOR IMPROVEMENT**
2  Prior to linking data, researchers should carefully examine the individual datasets to ensure data
3  integrity and completeness. Bohensky *(17)* recommends that direct, unique identifiers be
4  included within datasets to reduce bias from incorrectly entered data or missing variables.
5  Unfortunately, privacy concerns may prevent this method from being implemented in the United
6  States. In the absence of direct identifiers, Bohensky *(17)* suggests the use of financial incentives
7  to encourage data custodians to improve data quality and consistency. Definitions of variables
8  used for linking, such as injury severity, should be standardized at the national or international
9  level so that datasets from different sources can be reliably linked. Data linkage studies need to
10  develop a clear and systematic method of reporting their methodology so that they can be easily
11  repeatable by other researchers *(17)*.
12      Additionally, data linkers must develop indicators to describe the linkage errors present
13  in a linked dataset *(17)*. Harron suggests three approaches to evaluating linkage quality *(19)*. The
14  first approach is use of a gold standard dataset to directly measure missed and false matches—
15  while this is the most accurate approach, it is difficult to apply this methodology in practice as
16  gold standard data are rarely available. The second approach is to compare the characteristics of
17  linked and unlinked data to determine potential sources of bias—this approach requires a linkage
18  design in which all records in at least one file are expected to link, as well as provision of
19  characteristics of the unlinked records to the researchers. The third approach is to conduct a
20  sensitivity analysis to evaluate how sensitive the results are in response to changes in the linkage
21  method—this approach requires that researchers have access to match weights for each linked
22  record. Match weights do not reveal any sensitive information and thus are more likely to be
23  shared with researchers. However, the sensitivity analysis may be difficult to interpret as the
24  effects of missed matches and false matches cannot be distinguished from each other *(19)*.
25      Further research must be conducted to identify the populations most likely to suffer from
26  selection bias during data linkage, in addition to determining the effects that different linkage
27  processes have on reducing such bias *(17)*.
28      Additionally, researchers must use a consistent definition of reporting level so that results
29  can be compared. While hospital and police ascertainment rates are the most theoretically
30  accurate definitions of reporting level, researchers may consider use of hospital link rate instead
31  because it is the most common definition. Finally, linkage with other datasets should be
32  explored. For example, the addition of linked traffic data could help account for exposure in
33  crash safety analysis.
34
35  **CONCLUSIONS**
36  Research has argued that police collision reports tend to underrepresent bicycle and pedestrian
37  crashes, especially when motor vehicles are not involved. To account for this bias when making
38  traffic safety decisions based on police data, the estimated reporting level and the uncertainty of
39  the estimated reporting level must be known. Ten studies using data linkage to explore
40  pedestrian and/or bicyclist safety were evaluated and summarized. There may be other relevant
41  studies that could be reviewed in a more extensive literature review. Due to different definitions
42  of reporting level, periods of study, and study locations, it was difficult to directly compare the
43  studies. Even among the six studies using the hospital link definition, estimates of reporting
44  levels ranged from 44 to 75 percent for pedestrian crashes, and from 7 to 46 percent for bicycle
45  crashes.

1         These results indicate a severe underreporting problem in police collision reports, which
2    could lead to inaccurate estimates of crash rates and could under- or overestimate the effects of
3    road safety countermeasures.
4         However, few of the studies provided an estimate of the error around their reporting level
5    estimates. Therefore, it is possible that bicycle and pedestrian crashes are not as underreported as
6    these studies suggest. It may be that bicycle and pedestrian crashes appear in both police and
7    hospital datasets but are less likely to be linked. Because of linkage error, link rate can only be
8    used to *estimate* reporting level. Without knowledge of the *variance* of that estimate, the effect
9    of underreporting on traffic safety analyses cannot be accurately determined.
10        Future studies must include estimates of the error present in their data linkage process so
11   that the level of underreporting in police data can accurately be measured and accounted for.
12   Additionally, datasets should be designed so that they can more easily be linked. This could
13   involve standardizing the definition of common fields in each dataset, but could also involve the
14   introduction of some type of individual identifier so that records can be linked automatically.
15   Finally, linkage with other datasets should be explored.
16
17   **The authors confirm contribution to the paper as follows: study conception and design: S.**
18   **Doggett, D. Ragland; data collection: S. Doggett; analysis and interpretation of results: S.**
19   **Doggett, D. Ragland; draft manuscript preparation: S. Doggett. G. Felschundneff. All**
20   **authors reviewed the results and approved the final version of the manuscript.**
21

22   **REFERENCES**

23   1.  Stutts, Jane C., and William W. Hunter. 1999. *Injuries to Pedestrians and Bicyclists: An*
24       *Analysis Based on Hospital Emergency Department Data*. Final FHWA-RD-99-078.
25       https://www.fhwa.dot.gov/publications/research/safety/pedbike/99078/index.cfm. Accessed
26       March 18, 2018.
27   2.  Elvik, Rune, and Anne Mysen. 1999. Incomplete Accident Reporting: Meta-Analysis of
28       Studies Made in 13 Countries. *Transportation Research Record: Journal of the*
29       *Transportation Research Board* 1665 (January): 133–40. https://doi.org/10.3141/1665-18.
30   3.  Hauer, E, and A S Hakkert. 1988. Extent and Some Implications of Incomplete Accident
31       Reporting. *Transportation Research Record* Methods for Evaluating Highway Improvements
32       (1185): 1–10.
33   4.  Janstrup, Kira H., Sigal Kaplan, Tove Hels, Jens Lauritsen, and Carlo G. Prato. 2016.
34       Understanding Traffic Crash Under-Reporting: Linking Police and Medical Records to
35       Individual and Crash Characteristics. *Traffic Injury Prevention* 17 (6): 580–84.
36       https://doi.org/10.1080/15389588.2015.1128533.
37   5.  Short, Jack, and Brian Caulfield. 2016. Record Linkage for Road Traffic Injuries in Ireland
38       Using Police Hospital and Injury Claims Data. *Journal of Safety Research* 58 (Supplement
39       C): 1–14. https://doi.org/10.1016/j.jsr.2016.05.002.
40   6.  Tin Tin, Sandar, Alistair Woodward, and Shanthi Ameratunga. 2013. Completeness and
41       Accuracy of Crash Outcome Data in a Cohort of Cyclists: A Validation Study. *BMC Public*
42       *Health* 13 (1). https://doi.org/10.1186/1471-2458-13-420.
43   7.  Sciortino, Stanley, Mary Vassar, Michael Radetsky, and M. Margaret Knudson. 2005. San
44       Francisco Pedestrian Injury Surveillance: Mapping, under-Reporting, and Injury Severity in
45       Police and Hospital Records. *Accident Analysis & Prevention* 37 (6): 1102–13.
46       https://doi.org/10.1016/j.aap.2005.06.010.

1    8.  Tarko, Andrew, and Md. Shafiul Azam. 2011. Pedestrian Injury Analysis with Consideration
2        of the Selectivity Bias in Linked Police-Hospital Data. *Accident Analysis & Prevention* 43
3        (5): 1689–95. https://doi.org/10.1016/j.aap.2011.03.027.
4    9.  Conderino, Sarah, Lawrence Fung, Slavenka Sedlar, and Jennifer M. Norton. 2017. Linkage
5        of Traffic Crash and Hospitalization Records with Limited Identifiers for Enhanced Public
6        Health Surveillance. *Accident Analysis & Prevention* 101 (April): 117–23.
7        https://doi.org/10.1016/j.aap.2017.02.011.
8    10. Langley, J D. 2003. Missing Cyclists. *Injury Prevention* 9 (4): 376–79.
9        https://doi.org/10.1136/ip.9.4.376.
10   11. Lujic, Sanja, Caroline Finch, Soufiane Boufous, Andrew Hayen, and William Dunsmuir.
11       2008. How Comparable Are Road Traffic Crash Cases in Hospital Admissions Data and
12       Police Records? An Examination of Data Linkage Rates. *Australian and New Zealand*
13       *Journal of Public Health* 32 (1): 28–33. https://doi.org/10.1111/j.1753-6405.2008.00162.x.
14   12. Watson, Angela, Barry Watson, and Kirsten Vallmuur. 2015. Estimating Under-Reporting of
15       Road Crash Injuries to Police Using Multiple Linked Data Collections. *Accident Analysis &*
16       *Prevention* 83 (October): 18–25. https://doi.org/10.1016/j.aap.2015.06.011.
17   13. California Department of Public Health. 2015. *Exploratory Analysis of Injury Classification*
18       *of Crash Victims Using Crash-Medical Linked Data in California*. California Department of
19       Public Health.
20       https://archive.cdph.ca.gov/programs/Documents/Exploratory%20Analysis%20of%20Injury
21       %20Classification%20of%20Crash%20Victims%20Using%20Crash-
22       Medical%20Linked%20Data_Final.pdf. Accessed March 18, 2018.
23   14. Washington Area Bicyclist Association. 2015. *Modernizing the Collection, Integration and*
24       *Disclosure of Crash Data: Policy Recommendations for the District of Columbia.*
25       Washington Area Bicyclist Association. http://www.waba.org/wp-content/uploads/
26       2016/01/DC-Crash-Data-Policy-Paper-July-2015.pdf. Accessed March 24, 2018.
27   15. Ahmed, Ashar, Ahmad Farhan Mohd Sadullah, and Ahmad Shukri Yahya. 2017. Errors in
28       Accident Data, Its Types, Causes and Methods of Rectification-Analysis of the Literature.
29       *Accident Analysis & Prevention*, July. https://doi.org/10.1016/j.aap.2017.07.018.
30   16. Abay, Kibrom A. 2015. Investigating the Nature and Impact of Reporting Bias in Road Crash
31       Data. *Transportation Research Part A: Policy and Practice*, 71 (January): 31-45.
32       https://doi.org/10.1016/j.tra.2014.11.002
33   17. Bohensky, Megan. 2016. Bias in Data Linkage Studies. In *Methodological Developments in*
34       *Data Linkage*, First, 63–82. John Wiley and Sons.
35   18. Harron, Katie, Harvey Goldstein, and Chris Dibben. 2016. Introduction. In *Methodological*
36       *Developments in Data Linkage*, First, 63–82. John Wiley and Sons.
37   19. Harron, Katie, James C Doidge, Hannah E Knight, Ruth E Gilbert, Harvey Goldstein, David
38       A Cromwell, and Jan H van der Meulen. 2017. A Guide to Evaluating Linkage Quality for
39       the Analysis of Linked Data. *International Journal of Epidemiology* 46 (5): 1699–1710.
40       https://doi.org/10.1093/ije/dyx177.
41   20. Imprialou, Marianna, and Mohammed Quddus. 2017. Crash Data Quality for Road Safety
42       Research: Current State and Future Directions. *Accident Analysis & Prevention*.
43       https://doi.org/10.1016/j.aap.2017.02.022.
44   21. Christen, Peter. 2012. Data Pre-Processing. In *Data Matching*, by Peter Christen, 39–67.
45       Berlin, Heidelberg: Springer Berlin Heidelberg. https://doi.org/10.1007/978-3-642-31164-
46       2_3.