

9-3-2013

# Assessing Uncertainty in Volunteered Geographic Information for Emergency Response

Michael Camponovo

Follow this and additional works at: [https://digitalrepository.unm.edu/geog\\_etds](https://digitalrepository.unm.edu/geog_etds)

---

## Recommended Citation

Camponovo, Michael. "Assessing Uncertainty in Volunteered Geographic Information for Emergency Response." (2013).  
[https://digitalrepository.unm.edu/geog\\_etds/18](https://digitalrepository.unm.edu/geog_etds/18)

This Thesis is brought to you for free and open access by the Electronic Theses and Dissertations at UNM Digital Repository. It has been accepted for inclusion in Geography ETDs by an authorized administrator of UNM Digital Repository. For more information, please contact [disc@unm.edu](mailto:disc@unm.edu).

Michael E. Camponovo

*Candidate*

---

Geography and Environmental Studies

*Department*

---

This thesis is approved, and it is acceptable in quality and form for publication:

*Approved by the Thesis Committee:*

Dr. Scott M. Freundsuh, Chairperson

---

Dr. Karl Benedict

---

Dr. Christopher Lippitt

---

---

---

---

---

---

---

---

---

**ASSESSING UNCERTAINTY IN VOLUNTEERED GEOGRAPHIC  
INFORMATION FOR EMERGENCY RESPONSE**

**BY**

**MICHAEL E. CAMPONOVO**

B.S., Agriculture, Tennessee Technological University, 2004  
M.A., Curriculum and Instruction, Tennessee Technological University, 2008

THESIS

Submitted in Partial Fulfillment of the  
Requirements for the Degree of

**Master of Science  
Geography**

The University of New Mexico  
Albuquerque, New Mexico

**July 2013**

## **DEDICATION**

This research project is dedicated to the volunteers who worked tirelessly to help the victims of the devastating earthquake that struck Haiti on January 12, 2010.

## ACKNOWLEDGEMENTS

This research project would not have been possible without the combined effort and encouragement of numerous individuals. I would like to thank my committee chair, Dr. Scott Freundsuh, for his guidance, advice, encouragement, and motivation. His tough questions and unwavering support are responsible for all that is good in this research project.

I would also like to thank my other committee members, Dr. Karl Benedict and Dr. Chris Lippitt. Dr. Benedict has provided numerous technical contributions to this project as well as encouraged me to pursue my passion in GIS as a career. Dr. Lippitt helped keep me grounded when struggling with early ideas of this project and introduced me to other geospatial professionals for advice.

The other faculty and staff members within the department have helped prepare me for this project in so many ways. Dr. Paul Zandbergen provided me with numerous opportunities to learn first-hand how to conduct quality geospatial research. Dr. Maria Lane worked tirelessly to improve early versions of several sections of this paper. Dr. John Carr has always been available to offer guidance and advice in all manners of challenges, both personal and academic. Dr. Bradley Cullen has been instrumental in helping me learn the statistics necessary to complete this project. Mary Thomas' organizational ability ensured that deadlines and paperwork were always turned in on time.

Prior to my time at UNM, I took courses in GIS at several institutions. Dr. Reed Cripps, Jason Duke, Pat Wurth, Rich Winterfield, and Mark Young helped give me the geospatial background to make it as far as I have.

Susan Finger and Joshua Johnson were instrumental in encouraging me to continue pursuing my career goals when the economy seemed most unwilling to cooperate.

My fellow graduate students have served as sounding boards and cheerleaders throughout this entire process. Thank you all for helping me keep my sanity through the last several years.

My family has been incredibly patient with me as visits home have always involved a great deal of time sitting in front of a computer while working on classwork and research projects. I can't wait to catch up on lost time.

My wife Sarah has borne the greatest burden from my seeming inability to get out of school. So many evenings and weekends were spent working on homework, studying for tests, preparing for presentations, and going to conferences. Thank you for always supporting me, never losing your patience over the time that I was away, and for always being my best friend and partner. I am so thankful for the time we have shared together.

**ASSESSING UNCERTAINTY IN VOLUNTEERED GEOGRAPHIC  
INFORMATION FOR EMERGENCY RESPONSE**

BY

**MICHAEL E. CAMPONOVO**

B.S., Agriculture, Tennessee Technological University, 2004  
M.A., Curriculum and Instruction, Tennessee Technological University, 2008  
M.S., Geography, University of New Mexico, 2013

**ABSTRACT**

This research project examines data produced by volunteers through the Ushahidi web platform in response to the earthquake that struck Haiti in January 2010. Volunteers translated messages submitted by victims in Haiti, categorized each message based on its content, and georeferenced each message on a dynamic web based map. When categorizing the data, volunteers were able to assign up to 8 main and 42 subcategories to each message. Initial inspection of the attribute data produced by the volunteers indicated a strong discrepancy between the contents of the messages submitted by the victims and the corresponding attributes assigned to those messages by the volunteers. By comparing the attributes of the data originally produced by the volunteers to data that I re-categorized, I was able to examine the degree of inconsistency among the attribute data produced by the volunteers. I found that only 26.59% of the messages submitted by the victims were consistently categorized compared to the data set that I re-categorized. However, when aggregating the subcategories up to their appropriate main category, I

found 49.88% of messages were consistently categorized indicating that approximately half of the messages were conveying the main idea or ideas of the victims' messages. These numbers are significantly lower than the estimate of 64% correct categorization produced by an independent review of the Ushahidi platform. Despite these low indicators of consistent categorization, the volunteer response to the Haitian earthquake represents a paradigm shift in emergency response and victim empowerment that has been repeated in numerous natural and man-made disasters around the world.

## TABLE OF CONTENTS

<b>LIST OF FIGURES .....</b>	<b>viii</b>
<b>LIST OF TABLES .....</b>	<b>ix</b>
<b>Chapter One – Introduction.....</b>	<b>1</b>
Project Description .....	2
Background .....	2
Significance.....	4
<b>Chapter Two – Literature Review .....</b>	<b>6</b>
GIS and Hazards .....	6
Crowdsourcing and Social Media .....	11
Data Quality.....	18
<b>Chapter Three – Research Design .....</b>	<b>23</b>
Methodology.....	27
Initial Data Collection .....	27
Database Creation .....	31
Re-categorization .....	35
Category and Subcategory Comparisons .....	38
Statistics .....	40
<b>Chapter Four – Results.....</b>	<b>41</b>
<b>Aggregate Results .....</b>	<b>41</b>
Main Categories .....	41
Subcategories.....	44
<b>Daily Results.....</b>	<b>50</b>
Main Categories .....	51
Subcategories.....	64
<b>Statistical Results .....</b>	<b>81</b>
Main Categories and Subcategories, Independent Ushahidi Review versus Researcher Findings.....	81
<b>Chapter Five – Discussion .....</b>	<b>83</b>
<b>Chapter Six – Conclusions .....</b>	<b>87</b>
Summary of Findings.....	87
Limitations.....	88
Future Research.....	90
<b>Appendices.....</b>	<b>94</b>
<b>Appendix A: Ushahidi Haiti Category Rules and Definitions .....</b>	<b>94</b>
<b>Appendix B: Main Category Python Script .....</b>	<b>97</b>
<b>Appendix C: Subcategory Python Script.....</b>	<b>99</b>
<b>Appendix D: Subcategory Entries by Date.....</b>	<b>104</b>
<b>Appendix E: Haiti Map.....</b>	<b>151</b>
<b>References.....</b>	<b>152</b>



## LIST OF FIGURES

Figure 1 - Emergency Management Cycle .....	7
Figure 2 - Ushahidi Haiti Flow Chart .....	25
Figure 3 - Sample Ushahidi Form.....	26
Figure 4 - Ushahidi Haiti Web Interface (Edublog, 2011) .....	26
Figure 5 - Breakdown of how many records remain after each step in removing records not in English and duplicate records.....	34
Figure 6 - Instances of Each Main Category in Total.....	42
Figure 7 – Instances of Each Subcategory, Total .....	45
Figure 8 - Total Reports per Day .....	51
Figure 9 - Original Reports by Category per Day.....	52
Figure 10 – Revised Reports by Category per Day .....	53
Figure 11 - Cumulative Entries per Day, Original.....	55
Figure 12 - Cumulative Entries by Day, Revised .....	56
Figure 13 - Category 1, Original versus Revised, by Day .....	57
Figure 14 - Category 2, Original versus Revised, by Day .....	58
Figure 15 - Category 3, Original versus Revised, by Day .....	59
Figure 16 - Category 4, Original versus Revised, by Day .....	60
Figure 17 - Category 5, Original versus Revised, by Day .....	61
Figure 18 - Category 6, Original versus Revised, by Day .....	62
Figure 19 - Category 7, Original versus Revised, by Day .....	63
Figure 20 - Category 8, Original versus Revised, by Day .....	64
Figure 21 - Events for Subcategories 1a-1d by Day .....	66
Figure 22 - Events for Subcategories 2a-2f by Day.....	68
Figure 23 - Events for Subcategories 3a-3g by Day .....	70
Figure 24 - Events for Subcategories 4a-4e by Day .....	72
Figure 25 - Events for Subcategories 5a-5e by Day .....	74
Figure 26 - Events for Subcategories 6a-6c by Day .....	76
Figure 27 - Events for Subcategories 7a-7h by Day.....	78
Figure 28 - Events for Subcategories 8a-8h by Day.....	80

## LIST OF TABLES

Table 1 - Sample Ushahidi Haiti CSV .....	30
Table 2 - Sample Original Date to Julian Date Values .....	32
Table 3 - Main Categories in Ushahidi Haiti Deployment .....	35
Table 4 - Subcategories for Ushahidi Haiti Deployment.....	35
Table 5 - Subcategories with more than one definition .....	38
Table 6 - Subcategories Moved to Different Main Categories .....	38
Table 7 – Sample CSV generated by Python script comparing original and re-categorized results per record.....	40
Table 8 - Instances of Each Main Category in Total .....	41
Table 9 - Subcategory 1, Total Entries .....	46
Table 10 - Subcategory 2, Total Entries .....	46
Table 11 - Subcategory 3, Total Entries .....	47
Table 12 - Subcategory 4, Total Entries .....	48
Table 13 - Subcategory 5, Total Entries .....	48
Table 14 - Subcategory 6, Total Entries .....	49
Table 15 - Subcategory 7, Total Entries .....	49
Table 16 - Subcategory 8, Total Entries .....	50
Table 17 - Category Assessment by Discrepancy Type by Independent Ushahidi Evaluators .....	81
Table 18 - Category Assessment by Discrepancy Type, by Researcher at Main Category Level .....	81
Table 19 - Category Assessment by Discrepancy Type, by Researcher at Subcategory Level .....	82

## Chapter One – Introduction

January 12, 2010 was just another Tuesday for Jens Kristensen at the United Nations Headquarters in Port-au-Prince until a magnitude 7.0 earthquake struck the capital at approximately 4:53p.m. In a matter of seconds Mr. Kristensen went from being someone who provided relief to those in Haiti to a victim. Mr. Kristensen's third floor office was now at ground level surrounded by piles of debris. Fortunately for Mr. Kristensen, as the debris fell around him it left a pocket where he was able to take shelter mostly unscathed. Within a few short hours of the earthquake, a small group of volunteers over sixteen hundred miles away in Boston, Massachusetts set up a website to collect tweets, text messages, emails, and news reports about the disaster and place those reports on a dynamic web based map. At exactly 11:01 the next day the volunteers received and published the following message for anyone, including relief agencies in Haiti, to view, "Over 100 #UN personnel trapped in collapsed headquarters in #Haiti earthquake..."(Ushahidi, 2011). Meanwhile in Haiti Mr. Kristensen did his best to keep his situation as positive as possible by collecting anything that would help him survive and thinking about his family. Mr. Kristensen then waited for the next five days. On the afternoon of Sunday January 17 the Fairfax County Search and Rescue Team pulled Mr. Kristensen from the rubble of his former office (LaFranchi, 2010).

Situations like this have become common occurrences in the wake of natural disasters since the mid 2000's because the technology that we carry continues to advance at an astounding pace. In disasters that have spanned six continents, everyday citizens are reporting geo-located information through social media applications like Twitter that traditional aid agencies are absorbing into their protocols and acting upon, oftentimes

saving lives in the process. The convergence of technology, social media, and geography have resulted in a new area of research referred to as volunteered geographic information (VGI) that enables anyone, regardless of geographic knowledge, to produce geographic data with very significant results (Goodchild, 2007). This research is designed to assess the quality of data submitted by “citizen” geographers in the emergency response phase of disaster situations.

### **Project Description**

Despite VGI facilitating the rapid accumulation of data from numerous human sensors, it is often criticized in the literature for a lack of quality (Goodchild and Glennon, 2010). This research project examines one component of geospatial data quality for a data set produced in response to a time sensitive emergency. A thorough examination of the attributes of the Ushahidi database was conducted in order to assess the consistency of the data produced by volunteers during the disaster. This thesis will address the question, “What is the nature of uncertainties in the attribute data distributed via the Ushahidi geospatial platform in response to the Haiti earthquake of 2010?”

### **Background**

Historically, geospatial data has been created by and shared from a select group of organizations commonly referred to as authoritative data sources. In the United States, agencies like the United States Geologic Survey (USGS) have been responsible for collecting geographic information and disseminating that information in both paper and digital formats. The relative paucity of data collectors was due to the high cost associated with collecting geographic data. This high cost included the technical training necessary

to map objects in the field, the cost of the mapping equipment, and the time necessary to complete mapping projects.

In recent years, the high cost associated with mapping has significantly decreased. Handheld GPS units are relatively inexpensive and they are now included on most cellular phones. Open source and free versions of mapping software have been produced that provide various levels of the functionality of commercial geographic information systems. The increase in availability of broadband Internet access has allowed for data to be shared quickly and easily. By lowering the cost of entry for creating and sharing geographic information, mapping has become another addition to the Web 2.0 movement that allows not only the use of data shared over the Internet, but also the creation of data. This melding of Web 2.0 and geography has resulted in what Goodchild refers to as “volunteered geographic information” (VGI) (2007). VGI has many potential benefits for the geospatial community. For instance, VGI can provide free access to data instead of relying only on commercial options. For example, OpenStreetMap provides free road network data that can be used in place of costly commercial data. VGI facilitates faster data updates, when in-car navigation companies accept user generated corrections and additions to their road networks instead of relying only on their own data collection processes. In addition, data that was previously too insignificant or costly to collect is now map-able. An example of this is during the Super Bowl, maps depicting where fans for each team were located in the United States were generated by collecting the locations of Twitter feeds and parsing their content (Bloch and Carter, 2012).

While this democratization of geospatial data has many potential benefits, it is not without problems. In the past, quality of geospatial data was assumed to be relatively

high. This was due to the nature of the provider, oftentimes government agencies and reputable commercial providers, and the time and effort that went into collecting and producing data. The notion of assumed quality has declined since almost anyone can now create geographic data, regardless of training or expertise. Recent disasters in Haiti and Libya have shown that agencies with a geographic component are willing to make decisions based on volunteered data that may be of questionable accuracy and credibility (Standby Volunteer Task Force, 2011). As more agencies begin to embrace social media, more questions will surely arise over the quality of the information that is being shared by individuals and the impacts that data quality will have on the decisions that the agencies make. Agencies that work in fields that are not greatly affected by time have the luxury to assess volunteered data and determine its worth (Haklay, 2010; Haklay *et al.*, 2010; Haklay and Ellul, 2010; Girres and Touya, 2010; Zielstra and Zipf, 2010). But what about agencies that regularly have to make rapid decisions that can have a profound impact on life or death? For example, agencies like the American Red Cross, firefighters, and 911 dispatchers may not have the time to assess the volunteered data that is being shared with them to determine its quality and credibility. If these agencies act on data of questionable quality, what impact will the quality of that data have on their operations? This thesis seeks to begin to answer these questions by studying the quality of VGI produced during a disaster. The goal of this research is to help relief agencies more accurately assess what data they may want from volunteers and in what capacity those data will be employed.

### **Significance**

There is a pre-existing and substantial body of published research on the individual components of this research project. The field of spatial data quality has been

researched for decades with some of the most prominent names in geographic information science contributing to the field. Research has focused on defining quality as well as creating methods to measure its different components. The use of geographic information systems in the disaster management cycle has also been well documented. Because hazards have a strong spatial component, GIS has been shown to be beneficial in all stages of emergency preparation and response. The role of social media and crowdsourcing (when an undefined group of individuals are tasked with solving a problem rather than designating a specific person or entity to solve it (Howe, 2008)) in society has been extensively studied with research focusing on topics as varied as why people volunteer, to the asserted reliability of the data. A recent emphasis on a specific type of crowdsourcing, VGI, is also gaining more attention as more researchers contribute to this specific field. Some researchers have even combined two of the above components in their research by studying the quality of volunteered geographic information (Haklay, 2010; Haklay *et al.*, 2010; Haklay and Ellul, 2010; Girres and Touya, 2010; Zielstra and Zipf, 2010) or the role VGI plays in disaster response (Standby Task Force, 2011; Norheim-Hagtun and Meier, 2010; Pitzer, 2011). However, there seems to be a lack of research that investigates the intersection of data quality, crowdsourcing/social media, and GIS/hazards. This unique combination will provide the research community and emergency responders with new insight as to how VGI can be used as part of their arsenal of response tools. This research will attempt to aid researchers and responders in determining the fitness of use of VGI data.

## **Chapter Two – Literature Review**

This literature review explores three topics related to the field of volunteered geographic information and disaster response. The first topic, GIS and Hazards, discusses the role of GIS in emergency response as well as the data needs and current limitations within the field. The second topic pertains to Crowdsourcing and Social Media and provides a general history of the field along with specific applications to geography and strengths and weaknesses within the field. The last topic is GIS Data Quality and presents methods for measuring the quality of spatial data as well as specific case studies pertaining to VGI data quality.

### **GIS and Hazards**

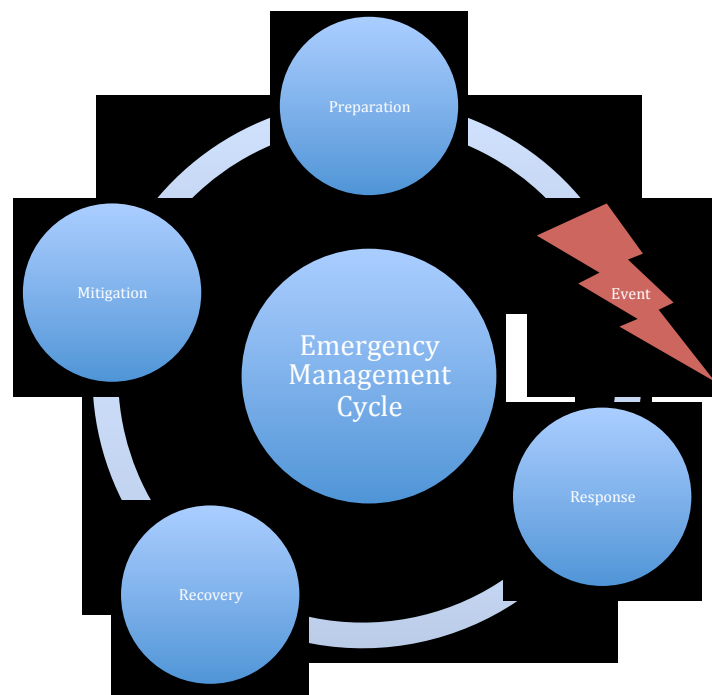
Natural hazards exhibit strong spatial patterns and are therefore a suitable topic of study for the field of geography. Within the field of geography, geographic information systems help emergency responders make empirical decisions related to spatial questions. This section provides an overview of the emergency management cycle and the role that GIS takes within it while focusing on the response phase of the emergency cycle and the specific needs for geospatial analysis within that phase.

Natural and manmade hazards continue to plague people as evidenced by news and media sources around the world. Because all of these hazards occur at or near the earth's surface, they all have a spatial component and are therefore well suited to analysis by geographers. For instance, some regions are known by their corresponding natural hazard like "Tornado Alley" and the "Pacific Ring of Fire," while other regions are known for cyclical disasters like flooding and tropical storms and hurricanes. Due to the geography of certain regions, some disasters can typically only occur in particular regions



like tsunamis along the coast and earthquakes along plate boundaries (Keller and Blodgett 2006). For all of these reasons, hazards are well suited to analysis through a specific branch of geography, geographic information science (GIS).

It is common practice for planning and academic purposes to refer to the actions taken as a result of a disaster as the emergency management cycle. This cycle is typically categorized by four phases: preparedness, response, recovery, and mitigation. Refer to Figure 1.



**Figure 1 - Emergency Management Cycle**

Each phase is defined by the actions taken within it and when those actions take place. The preparedness phase occurs prior to the onset of the emergency and consists of activities that help prepare the community for the upcoming disaster. An example of a preparedness procedure is the evacuation of at-risk areas due to wildfire or hurricane risk. Because some disasters are slow-onset and others are fast-onset, every disaster may not incorporate all of the actions typical with this phase. The response phase takes place

immediately after the disaster and is typically associated with actions taken to reduce the loss of life and property. Response procedures might include search and rescue operations and the distribution of water and blankets. The response phase transitions into the recovery phase. The recovery phase is typified by the actions necessary to return life in the affected community to normal. Recovery activities might include the reconstruction of homes and businesses in an affected community. The mitigation phase strives to limit the effect of future emergencies on the affected community. Examples of typical actions that might result during the mitigation phase include legislation that prohibits construction in a flood plain or requires wind-resistant construction practices in hurricane prone areas (Cova, 1999; Radke *et al.*, 2000; National Research Council, 2007a). While given four distinct names and comprising specific activities, it is not uncommon for the different phases to overlap depending on the emergency.

Geographic Information Systems can and do play a key role in each phase of the emergency management cycle (Cova, 1999; Radke *et al.*, 2000; National Research Council, 2007a). In the preparedness phase GIS can be used to help model a hurricane to predict where it will make landfall and identify which areas should be evacuated. GIS can also help determine the routes needed once the evacuation order is in place. During the response phase GIS can be used to produce maps to guide search and rescue teams as well as determine where to allocate resources like emergency shelters. GIS can be used during the recovery phase to determine the extent of damage to determine where improved construction techniques may be needed. During the mitigation phase GIS can be used to help politicians and lawmakers delineate at-risk areas or make sure future construction is not located within at-risk areas (Maliszewski and Horner, 2010).

In order to make the GIS most effective, emergency responders need very specific pieces of spatial data for the affected region. Prior to the disaster, responders need to accumulate baseline vector and raster data (National Research Council, 2007a; van Westen and Georgiadou, 2001). The raster data may be in the form of satellite or aerial imagery in the visible spectrum or in bandwidths that indicate heat. Vector data should include information like roads and transportation networks, water and sewer lines, gas and electric utilities, communication infrastructures for landlines, Internet, and cellular services, hospitals, fire stations, law enforcement, hazardous materials, and emergency resources like supplies and pre-established shelters. Responders will also need census and demographic data about the affected region (National Research Council, 2007b). Once the disaster has taken place, responders will need access to new imagery to delineate the affected region and make damage estimates (Kelmelis *et al.*, 2006; National Research Council, 2007a). They will also need data pertaining to where affected citizens are moving to and congregating (Kaiser *et al.*, 2003). In order to make search and rescue teams most effective, the responders will also need to know where people are trapped and in what condition they are in.

All of the above data will eventually make its way into a spatial decision support system (SDSS). This system is designed to help emergency responders make empirical decisions in time critical situations. These systems may be standalone (Tomaszewski, 2011) or incorporate add-ons to existing GIS software (Nguyen, 2005). The SDSS will conduct analysis of the data that has been collected to help the decision makers. For instance, the SDSS may be used to locate emergency shelters based on demographic and census data along with delineations of the affected region, or evacuation routes may be

updated based on infrastructure damage data collected in the field. The results of the analysis and the decisions made from those analyses will eventually make their way to responders via maps, either in print or electronic forms. These maps may consist of grids for search and rescue teams or the routes to be taken by convoys delivering food, water, and medical supplies.

While the use of GIS in responding to emergencies has increased in frequency and effectiveness, it is not without limitations. For instance, the best imagery in the world is of no use if it is saved in a file format that is incompatible with the hardware and software that the emergency responders are using (Heinzelman *et al.*, 2010). Nor will data be useful if it comes attached to licensing restrictions that forbid its dissemination (National Research Council, 2007a; van Westen and Georgiadou, 2001). An SDSS that accurately determines where trapped victims of an earthquake are is of no use if the geospatial analyst does not have a good communication protocol established with the search and rescue teams (Piotrowski, 2010). Nor does it help to have a powerful GIS that no one is trained to use (Zerger and Smith, 2003). An overreliance on Internet access can also be a problem because communications infrastructure are often affected during emergencies (Frassl *et al.*, 2010). In many areas outside of the developed world, there are limited geospatial resources and much of the baseline data that responders would like access to are unavailable (Kelmelis *et al.*, 2006; Cutter, 2003). Another key component of accurate use of GIS in emergencies is the temporal quality, or timeliness, of the data that is being used within the SDSS (Kelmelis *et al.*, 2006; Cutter, 2003).

The next section of this literature review describes the tools and technology that have had a significant impact on response methods, particularly how advances in

technology and communications can compensate for the dearth of baseline and real time data that emergency responders are often faced with.

### **Crowdsourcing and Social Media**

Advances in technology have greatly increased the ease with which creation and sharing of digital content in what is referred to as Web 2.0 (O'Reilly, 2005). This technological revolution has played a part in geography as well as leading to the relatively new field of VGI where amateurs can submit and create some types of geographic content just as easily as professionals. Despite the increase in the amount of geospatial data that is now available, it may not be as useful as many users would like. This section provides a brief history of the tools and advances that allowed for the development of Web 2.0 and the ensuing developments of crowdsourcing and social media. The section focuses on the role that these two developments have played in geography.

The Internet began as a tool for consuming information, but as more material was provided online and more organizations created a presence on the web, there was a shift of purpose on the Internet. The term Web 2.0 describes this shift from consuming data to producing data on the Internet (O'Reilly, 2005). Examples of websites that use Web 2.0 technology include YouTube (<http://www.youtube.com>) that allows users to submit their own videos and Flickr (<http://www.flickr.com>) that allows users to share their own photos. E-commerce sites like Amazon (<http://www.amazon.com>) also embrace the technology by allowing users to submit ratings for products. The transformation of the Internet is due to the lowering cost of digital data creation tools (Howe, 2008) in combination with advances in web technology to enable interactive web tools. Many households now have digital still and video cameras that are relatively inexpensive and

allow for the production of quality media. Many households also have the prerequisite computer hardware and software that allows for the editing and manipulating of digital content. Finally, many households have high-speed Internet access that allows for the easy sharing and distribution of digital content.

The transformative ability of technology did not just stop with videos and pictures. The ability to create digital content now also applies to geography (Goodchild, 2007). The cost of global positioning satellite receivers has decreased over time while their accuracy and ubiquity have increased. Organizations like Google (<http://maps.google.com>) and Microsoft (<http://www.bing.com/maps>) provide free imagery online as well. Many pieces of software for the creation of maps are now available for free (<http://earth.google.com> and <http://www.google.com/mapmaker>) while sophisticated GIS analysis can be performed without purchasing expensive proprietary software (<http://grass.fbk.eu> and <http://www.qgis.org>). These free software packages often utilize best practices and advanced algorithms allowing for amateurs to produce aesthetically pleasing maps and sophisticated geospatial analysis without any formal geographical training (Crampton, 2009) in what has become known as neo-geography (Turner, 2006). Online users have created social networks based on their common interests in utilizing these tools.

These social networks with common interests have spawned what are referred to as crowdsourcing (Howe, 2008). Crowdsourcing allows for anyone who is interested to participate in solving a problem or reaching a goal. Common examples of crowdsourcing are the free online encyclopedia Wikipedia (<http://www.wikipedia.org>) and the free, open source Linux operating system (<http://www.linux.com>). The benefits of crowdsourcing

can be many. For instance, a group of individuals who are interested in a problem but lack technical expertise may find a novel solution to that problem because they are not encumbered by the dogma of the discipline. For example, chemistry technicians at Colgate trying to solve a problem related to toothpaste manufacturing were stymied until a crowdsourced physicist applied his knowledge of electrically charged particles in what seemed an obvious solution to him (Howe, 2008:150). Or, the crowd may be able to solve a problem that is beyond the scope of any agency to solve alone because of the many participants in the crowd. Amazon (<http://www.amazon.com>), an online merchant, for instance, would have a difficult time rating all of its various products, so it allows its users to do so. This benefits Amazon because it does not have to pay someone to review thousands of products and the users benefit because they get various opinions rather than just one from an Amazon employee.

Crowdsourcing has several applications within the field of emergency management and disaster response and therefore, geography. During the Indian Ocean Tsunami of 2004, many affected countries did not have access to sophisticated and oftentimes expensive emergency response software. Programmers in Sri Lanka pooled their coding skills and within three weeks of the emergency created a free and open source emergency response software package that was modifiable, scalable, operated with minimal hardware and software, and protected the privacy of users and contributors to the software (Currion, De Silva, and van De Walle, 2007). Their software, named Sahana (<http://sahanafoundation.org>), has since been used in numerous emergency situations around the world. Another application designed for smart phones, Outbreaks Near Me (<http://healthmap.org/outbreaksnearme>) also incorporates crowdsourcing. Users

of the app can report locations where people are infected with various diseases, like swine flu, allowing health professionals and the community to monitor the spread of disease (Freifeld *et al.*, 2010).

Businesses and organizations can also participate and benefit from crowdsourcing. Businesses in South Florida are susceptible to hurricanes just like people are and they operate more efficiently when in sync with their suppliers and distributors. Crowdsourcing applications that mechanically parse reports issued by participating companies can be used to gather data about when businesses will be open and operational so that their partners can make better decisions about their own operations. In this situation, businesses participating in the crowd gain the ability to operate more effectively and efficiently (Zheng *et al.*, 2010).

Government and nongovernmental (NGO) agencies that participate in disaster response can also benefit from crowdsourcing. Oftentimes agencies are duplicating efforts related to data collection or services without realizing it or acting on a need that has already been met. By participating in crowdsourcing applications, participating agencies can reduce the duplications of efforts and make more efficient use of limited resources (Gao *et al.*, 2011).

Another important component of crowdsourcing is online social media. Examples of social media applications include Facebook (<http://www.facebook.com>) and Twitter (<http://twitter.com>) where users can share information with their friends and followers. Shared information can be in the form of pictures and videos as well as hyperlinks and plain text. A significant majority of homes in America have computers with Internet access which allows for the use of social media from home (Gutnick *et al.*, 2011), while



almost half of all mobile subscribers in the United States have a smart phone which allows for the use of social media at all times (Nielsen, 2012). As the popularity of social media increases, not only are individuals using the services but also so are organizations and corporations (American Red Cross, 2010).

Social media has several applications within the field of emergency management and disaster response. Because many young people are accustomed to sharing their personal problems on social websites, they also share their physical problems and pleas for help instead of dialing 9-1-1 (Benko, 2011). Among older users of social media, there is also increasing use of applications like Twitter to report emergencies. Many users of social media feel that relief agencies like the American Red Cross should monitor Twitter for requests for assistance (American Red Cross, 2010). Social media does not have to provide only for the physical needs of affected users either. Some forms of social media, like EagleVox (<http://www.cersi.it/projects.html?view=project&task=show&id=4>) encourage users to communicate about the emotional aspects of a disaster because, for the survivors, these can be just as traumatizing as the physical ones (Banzato *et al.*, 2010).

Combining social media and crowdsourcing with geography has led researchers to this new field of VGI (Goodchild, 2007). An example of an application of VGI is the Ushahidi platform (<http://ushahidi.com>) that was developed in Kenya. Kenya was suffering from post-election violence during 2007 and 2008, yet citizens could not learn much about the situation because of a media blackout. The designer of the software operated a blog and asked her followers to send her emails pertaining to the violence they witnessed. She was quickly overrun with reports and asked software developers and

programmers to develop a system that would help automate the reporting process. The resulting software collects citizen reports that are categorized and geo-tagged and then placed on an online map. The software is open to anyone to report and anyone can access the reports (Okolloh, 2008).

A combination of Ushahidi as well as other VGI platforms played a key role a few years later during the Haiti earthquake of 2010. At the onset of the disaster in Port-au-Prince there was little geospatial data for the affected region. Therefore, the first task was to create the base data that first responders needed to make decisions. The Humanitarian OpenStreetMap Team (HOT) (<http://hot.openstreetmap.org/>) began by heads-up digitizing donated aerial and satellite imagery of the affected region with streets and buildings. However, the volunteers did not know the names of the streets or buildings in the maps that they were making because the maps lacked attribute information. The volunteers enlisted the help of Haitian expatriates to label the streets and buildings in order to make the maps more usable. The HOT team relied on Web 2.0 technology and crowdsourcing to quickly create geospatial data where none existed previously (Nelson, Sigal, and Zambrano, 2010). Ushahidi was used extensively to collect real-time information from affected Haitians on the ground (Norheim-Hagtun and Meier, 2010). A special short message system (SMS) was set up that allowed affected citizens to report conditions, and request help via texts (Pitzer, 2011). For instance, some text messages were sent that stated where a person was trapped in a building or that there was not enough water at a shelter. An initial problem with the texts, however, was that few of the responders and Ushahidi volunteers spoke Haitian Creole. A crowdsourcing system was developed to allow Haitian expatriates to translate the text messages to English so that

they could be mapped, coded, and then acted upon by relief agencies (Munro, 2010). Other volunteers mined Twitter and Facebook for updates related to the emergency and mapped those on the Ushahidi map as well (Nelson, Sigal, and Zambrano, 2010). The Sahana emergency response software platform was used by various NGO's to organize their response to the emergency (Nelson, Sigal, and Zambrano, 2010). Traditional Aid agencies like the US military used a form of wiki software to help alleviate bottlenecks in their relief efforts by allowing their members to contribute information about the disaster that responders could use to help make better and faster decisions (Yates and Paquette, 2011).

The various VGI platforms employed during the Haiti earthquake illustrate how the technology can be used to meet the needs of emergency responders through real-time data access and baseline data creation. For instance, in areas where there are few resources for critical needs like water and shelter, VGI can create the necessary baseline data that first responders need for allocating those resources. VGI can provide real-time data on the conditions of victims and infrastructure that may take days or weeks for authoritative sources to develop (Goodchild and Glennon, 2010). The data created by VGI can be open source and free, resulting in datasets that are stored in formats that are nonproprietary and can be used by anyone. VGI also helps meet the emotional needs of affected people by transforming them from powerless victims to empowered citizens (Elwood, 2008).

The Internet was originally used to consume information, but as the technology matured, it developed into a tool to create information as well. As users gained easier access to the tools necessary to create digital data, they also began to easily share that

data on the World Wide Web through social media. In addition, crowdsourcing allows Internet users to quickly and easily work collaboratively to solve problems. This transformation has had a significant impact on the collection of geographic data. As new tools are developed to easily collect and share geospatial and georeferenced data, geospatial analysts and researchers are provided with an abundance of new data resources, whether individually geotagged “tweets” or entire street networks from OpenStreetMap. These advances in technology are being implemented in new ways to respond to emergency situations around the globe. While some researchers see VGI as a benefit because of its ability to democratize access to and the creation of geospatial data (Elwood, 2008), others worry about the quality of the data that are being produced by non-professionals (Goodchild and Glennon, 2010). This concern is valid because of the important role that geospatial data provides in decision support systems, especially related to emergency response when lives are at risk. The last section of this literature review will provide an overview of techniques that are used to assess the quality of spatial data and models and specific techniques that can be used to assess VGI.

### **Data Quality**

Traditional measures of spatial data quality were appropriate when data was collected and distributed by a few agencies, but those same techniques may not be well suited when anyone is capable of producing spatial content (Goodchild, 2008). This section provides a brief background on spatial data quality followed by specific techniques used to measure the most common aspects of spatial data quality. Examples of techniques used to measure the quality components of VGI follow.

The quality of geospatial data has been an important consideration since the birth of maps and GIS, in particular due to the role of maps and spatial data in the decision

making process (Chrisman, 1991; Foody, 2003). Though most people initially think of positional accuracy when concerned with geospatial data quality, it is actually comprised of many different components. These components are lineage, positional accuracy, attribute accuracy, logical consistency, completeness, semantic accuracy, usage, purpose and constraints, and temporal quality (van Oort, 2006). GIScience researchers have developed several techniques to measure these different quality components of geospatial data.

Most research in this field focuses on positional accuracy. One technique for measuring the positional accuracy of a geospatial data set uses two data sets to compare to each other. This technique assumes one dataset — the reference data set — is of higher quality than the test data set, so this is a relative measure of accuracy. A buffer is created around the features in the reference data set and a percentage of the test set that falls within the buffer is calculated. This technique is best used when comparing datasets with few linear features like interstates or streams (Goodchild and Hunter, 1997). A modification of this technique looks for corresponding intrinsic nodes within the reference and test data sets, like street intersections, and compares the Euclidean distance between them (Tveite and Langass, 1999; van Niel and McVicar, 2002). Another method is to geocode addresses with the two data sets, and compare the locations of corresponding geocoded results using Euclidean distance tools (Lee, 2009).

Other components of spatial data quality are also measured, though less often. Completeness is an assessment of the absence of data or the presence of non-existent data. For point data this assessment is accomplished by summing the number of features of a certain type in a prescribed area in both a reference data set and a test data set and

comparing the results. Linear features are assessed in a similar manner by summing the length of features instead of their quantity (Haklay *et al.*, 2010; Zielstra and Zipf, 2010). The component of quality most relevant to this research is attribute accuracy, however. Girres and Touya (2010) list three separate components of attribute accuracy. One component, quantitative accuracy, can be assessed using statistical methods while another component, non-quantitative attributes, can be assessed using the Levenstein method for string comparison. The final component, and the most important for this research, relates to the correct classification of features.

Variations of these same techniques for measuring completeness, positional accuracy, and attribute accuracy are being used on VGI as well as techniques to measure other components of spatial data quality. The most common type of VGI analyzed is OpenStreetMap (OSM) data. These data are easy to access and download and are easy to add to existing GIS software for analysis. The reference data set is typically an authoritative data set produced by a commercial mapping company or national mapping agency. OpenStreetMap data from England, Germany, and France have been analyzed for positional accuracy, completeness, and attribute accuracy (Girres and Touya, 2010; Haklay, 2010; Haklay *et al.*, 2010; Haklay and Ellul 2010; Zielstra and Zipf, 2010). Through these measures certain trends in the OSM dataset were discovered. The quality of the OSM data set improves with population and socioeconomic status and decreases where there is low population density and low socioeconomic status (Girres and Touya, 2010; Haklay *et al.*, 2010; Haklay and Ellul, 2010; Zielstra and Zipf, 2010). In urban areas with relatively high socioeconomic status the positional accuracy, completeness, and attribute accuracy rivals that of more authoritative sources but changes to a more

heterogeneous quality as the distance from urban areas increases (Haklay *et al.*, 2010; Haklay and Ellul, 2010; Zielstra and Zipf, 2010).

While this research is beneficial because it illustrates that data created through crowdsourcing and VGI can be nearly as accurate as authoritative data sets, it does not address the unique situation of data created for time critical emergencies like disasters. OSM data has the benefit of being collected without a time constraint and with little risk to life and property if the data collector spends days or weeks collecting his or her data. In the response phase of emergencies, decision makers do not have the luxury of that kind of time. No research appears to have yet been done that uses these same techniques for testing relative data quality on data collected during an emergency. This may be due to the data itself because there are few authoritative sources of locations of trapped victims or low levels of supplies at emergency shelters. This does not mean, however, that researchers are not trying to determine VGI quality.

One technique would use Tobler's First Law of Geography and the crowd itself to evaluate asserted content. Tobler's First Law states that objects that are closer together are more similar than objects that are farther away (Tobler, 1970). Therefore a crowd of editors, similar to the Wikipedia model, would quickly be able to tell if asserted content was similar to its surroundings or not (Goodchild, 2008). Similarly, using computer algorithms, it may be possible to use spatial autocorrelation techniques to accomplish the same goal without human intervention (Sui, 2004). Another option would be to treat emergency data sources just like other social media and crowdsourcing data. Users could be ranked or rated by their peers based on the quality of their submissions or judged based on their number of followers (Flanagin and Metzger, 2008). While these techniques

may prove useful in the future, there do not appear to be any practical applications of these methods currently in use.

There are numerous methods for assessing the quality of geospatial data. The choice of methodology depends on the purpose of the research and the types of data available to analyze. Positional accuracy and completeness are easily compared as long as two data sets exist for the same phenomenon and one of those data sets is considered more accurate than the other. Analysis methods become more complicated, however, when only one data set exists. Using a combination of these approaches, this research project assesses the consistency of attribute data produced by volunteers in response to the Haiti earthquake of January 2010.



## Chapter Three – Research Design

Prior to a description of specific methods, a brief explanation of the original Ushahidi Haiti data set and how it was created is necessary. In broad terms, the data set was created by three sets of interconnected participants. Victims of the disaster comprise the majority of the first group. The victims consist of people in Haiti who were affected by the earthquake and its resulting damage. The victims were responsible for producing much of the raw data that was incorporated into the Ushahidi platform. Off-site volunteers comprise the majority of the second group. These volunteers were responsible for turning the raw data produced by the victims into information to be used by the third group. Relief agencies comprise the majority of the third group. This group was primarily responsible for providing aid to the victims and took advantage of the information provided on the Ushahidi Haiti website. These groups are not, however, mutually exclusive. For instance, members of relief agencies often provided information to the Ushahidi Haiti website in regards to their relief efforts. A more detailed accounting of the steps and roles of the three groups is described below. Refer to Figure 2. On January 12, 2010 a magnitude 7.0 earthquake struck the country of Haiti. Within a matter of hours the Ushahidi Haiti website had been established by volunteers in Boston, Massachusetts. These volunteers began collecting data from Twitter, email, and traditional media sources and georeferencing this data on a web-based map of Haiti. On January 16 the short message (SMS) code 4636 was established and advertised on local Haitian radio as a way to report your needs. Once the 4636 SMS was established, text messages comprised the majority of the incoming data for the web site [see Morrow *et al.* (2011) for a detailed timeline of events and Meier (2012) for a first hand account of the creation and

motivation for deploying the Ushahidi website]. The following steps describe the process of turning raw data from the victims into actionable information for the relief agencies. A victim submits a text message to the SMS 4636 from a cell phone. At this point the volunteers take the raw message and conduct three primary tasks. The first step is to translate the message from Haitian Creole to English. The next step is to read the message and categorize the contents with regard to the type of emergency. The last step is to georeference the message to a location in Haiti. If a message contains all of the necessary information to complete these three tasks it is entered into the Ushahidi Haiti database using a form and shows up on the map on the website where the information is available to anyone with an Internet connection. (See Figure 3 for a sample Ushahidi form, see Figure 4 for a screen capture of the Ushahidi Haiti website). At this point relief agencies on the ground can respond to individual messages, or reports, that are on the Ushahidi Haiti website by clicking on individual dots on the online map. When clicking on an individual dot, the web site presents the title of the report, the contents of the translated message, a list of categories and subcategories that the contents of the message relate to, the time the message was added to the web site, and the latitude and longitude assigned to the report. Reports can be filtered by date using the slider at the bottom of the website or by category by selecting the appropriate category on the right side of the website.

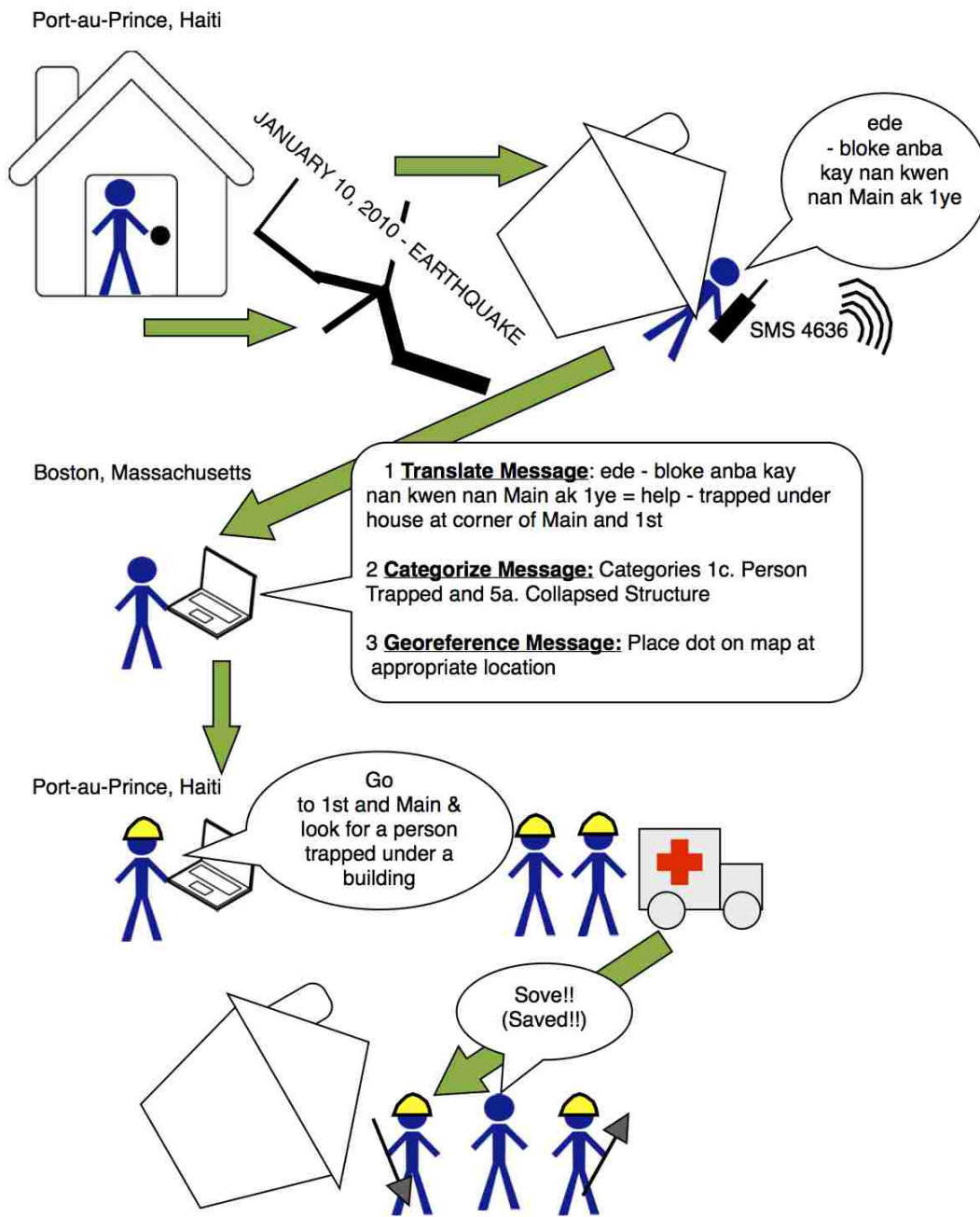


Figure 2 - Ushahidi Haiti Flow Chart

### New Report

SAVE REPORT SAVE & CLOSE SAVE & ADD NEW CANCEL

**Form** (Select A Form Type)

**Title \***

**Description** Include as much detail as possible.  
 Allowed HTML tags: "a, p, img, br, b, u, strong, em, i".  
 Iframes are only allowed from: www.youtube.com/embed/,  
 player.vimeo.com/video/, w.soundcloud.com/player

**Modify Date:** Today at 10:57 pm Modify Date

**Categories** (Select as many as needed) New Category

- Damaged Sidewalk
- Skateboard Damage
- Graffiti

**Incident Location** Latitude:  Longitude:  Wider Map

FIND LOCATION

\* Search for your location using a location name OR latitude,longitude coordinates (format: 38.19, 85.61), OR click on the map to pinpoint the correct location..

**Location Name \***  
 Example: Corner of City Market, 5th Street & 4th Avenue, Johannesburg

**News Source Link**

Figure 3 - Sample Ushahidi Form

## Haiti

The 2010 Earthquake in Haiti

Ushahidi-Haiti @ Tufts UNIVERSITY

**Transitioning to Noula.ht**  
 Our local partner, Haiti-company Solutions, has the most up-to-date reports for first responders — please visit their site for the most current data.

+ SUBMIT INCIDENT

Search Reports Here:

DOWNLOAD REPORTS (3605) REPORTS RSS

HOME   REPORTS   SUBMIT INCIDENT   GET ALERTS   CONTACT US   HOW TO HELP   ABOUT

**FILTERS** → REPORTS   NEWS   PICTURES   VIDEO   TODO

**VIEWS** → CLUSTERS

Terms of Use

Scale = 1 : 433K   -72.28568, 18.56383   EPSG:900913   © OpenStreetMap contributors, CC-BY-SA

**↓ TIMELINE OF EVENTS**

From:  To:  PLAY

Jan 12 2010	May 6 2010	Aug 28 2010	Apr 15 2010
-------------	------------	-------------	-------------

**↓ CATEGORY FILTER**

ALL CATEGORIES

- ! 1. URGENCES | EMERGENCY
- ⚡ 2. URGENCES LOGISTIQUES | VITAL LINES
- ⚕ 3. PUBLIC HEALTH
- 🛡 4. MENACES | SECURITY THREATS
- 🏠 5. INFRASTRUCTURE DAMAGE
- ⚡ 6. NATURAL HAZARDS
- + 7. SECOURS | SERVICES AVAILABLE
- ◇ 8. AUTRE | OTHER

Figure 4 - Ushahidi Haiti Web Interface (Edublog, 2011)

Initially I intended to conduct a variety of geospatial data quality comparisons between the Ushahidi data set and data produced by traditional aid agencies like the United Nations. However, upon examining the original Ushahidi data I realized that this would not be possible because there appeared to be significant discrepancies between the contents of the messages submitted by the victims and the categories assigned to those messages by the volunteers. This research project, therefore, was designed to evaluate the consistency of attribute values in the dataset produced by volunteers using the Ushahidi platform in response to the earthquake that struck Haiti on January 12, 2010. In particular, this project assesses the consistency of how the categories and subcategories were applied by the volunteers to the raw data from the victims. Unfortunately, there is no other appropriate data set with which to compare the Ushahidi data. As a result, I created a data set for comparison by re-categorizing the raw data from the victims. The specific methods outlined below rely on a quantitative analysis of the category attributes produced by the volunteers and those produced by myself. This technique is a consensus classification approach where agreement indicates an increased likelihood of correct entry.

## **Methodology**

### **Initial Data Collection**

The original database was not available to download so I relied on a comma separated value (CSV) file that I exported from the Ushahidi Haiti website. The original CSV formatted file was downloaded from the website <http://haiti.ushahidi.com> in Spring 2012. The CSV consists of 10 columns and 3,606 rows. The column headings are: “#”, INCIDENT TITLE, INCIDENT DATE, LOCATION, DESCRIPTION, CATEGORY, LATITUDE, LONGITUDE, APPROVED, and VERIFIED. The “#” symbol column

contains the unique identifier for each record in the data as a one, two, three, or four digit integer. The INCIDENT TITLE column contains the title given to each record in the database by the volunteer as determined from the original victim's message. The INCIDENT DATE column contains the date and time stamp for when the record was added to the database in the format YYYY-MM-DD HH:MM:SS. During the original export from the Ushahidi database to the CSV, however, the values were converted from a date/time stamp to a simple string. The LOCATION column contains a written description of the location referenced in the original message from the victim. This is a written description of the location, like the name of a specific community within a town or the nearest landmark, rather than a latitude and longitude. The DESCRIPTION column contains the contents of the message submitted by the victim. The DESCRIPTION column also occasionally contains notes from the volunteers about the particular message primarily intended for other volunteers. The CATEGORY column contains a list of the categories and subcategories that the volunteers determined were related to the original message submitted by the victim. The volunteers generated this list by selecting a checkbox next to each category and subcategory that he or she felt applied to the message. The format of the category information includes the number and/or number letter combination for each category and subcategory that pertained to the message as well as a written description of the category and/or subcategory name in both English and Haitian Creole. The checkboxes were not mutually exclusive so the CATEGORY field might contain a combination of categories and sub-categories. The next two columns, LATITUDE and LONGITUDE provide the geographic coordinates of the record in WGS84 decimal degrees. The Ushahidi software generated these attributes automatically

when the volunteer added a dot to the map where they thought the victims' reports were located. The Ushahidi platform has the ability to control which records become public through the use of information contained in the last two columns. Moderators have the ability to approve a message based on whatever criteria they establish and the results are found in the APPROVED column as either YES or NO. The VERIFIED column operates in the same manner and gives the moderators a chance to document if a message has been corroborated in some way. See Table 1 for a sample of the original Ushahidi CSV.

#	INCIDENT TITLE	INCIDENT DATE	LOCATION	DESCRIPTION	CATEGORY	LATITUDE	LONGITUDE	APPROVED	VERIFIED
638	Anesthesiologists needed	2010-01-17 02:08:00	Port-au-Prince	in the Hospital de la Paix are anesthesiologists and other doctors needed	1. Urgences   Emergency, 1b. Urgence medicale   Medical Emergency,	18.556439	-72.298248	YES	NO
4054	Pharmacy open	2010-07-27 20:20:00	lamare	Pharmacy at Lammare Street in Petionville open and selling medications	7. Secours   Services Available,	49.295769	-0.892294	NO	NO
4051	Food-Aid sent to Fondwa, Haiti	2010-06-28 23:06:00	fondwa	Please help food-aid.org deliver more food to Haiti through your financial gift and donations.	1. Urgences   Emergency, 2. Urgences logistiques   Vital Lines,	50.226029	5.729886	NO	NO

Table 1 - Sample Ushahidi Haiti CSV



### **Database Creation**

For the purpose of this research project, the data needed to be searchable, sortable, and georeferenced, which was not possible using a traditional text editor or spreadsheet program. Before the data could be converted to a database, however, a number of formatting inconsistencies were cleaned up. The steps involved were:

1. Using a text editor, remove the empty first row and remove all spaces from the field names
2. Perform a search and replace to remove all tabs and replace with \*tab\*
3. Open the CSV in a spreadsheet software package and export as a tab separated value (TSV)
4. Open the newly created TSV file in a text editor
5. Perform another search and replace to replace the \x0A hex character (upside down  $\iota$ ) with a space
6. Open the TSV file in a GIS package and export it as a SpatiaLite database

The free, open source Quantum GIS (QGIS) package was used to create the SpatiaLite database. SpatiaLite was chosen because of limitations to the field length in more traditional GIS file types like shapefiles and keyhole markup language (KML) files. The database generated through QGIS resulted in a database with 3,604 records ranging from January 12, 2010 to September 7, 2011. In order to perform the necessary analysis, the data had to be sorted by date; however, no true date field existed in the database. The database was exported as a CSV file and opened in a spreadsheet software package. A formula was used to convert the original, non-sortable date into a Julian date starting on January 1, 2010. Refer to Table 2 for sample dates in their original and Julian format.

Table 2 - Sample Original Date to Julian Date Values

INCIDENT DATE	Julian Date
2012-02-11	2012042
2011-09-07	2011250
2011-05-10	2011130

All fields except the unique identifier and the Julian date were removed from the CSV and were joined back to the database. At this point in the project I removed those records that were submitted after January 27, 2010. There were multiple reasons for this removal. The research project is primarily concerned with the response phase of the disaster management cycle. Going beyond January 27, 2010 begins reaching past the response phase into the rebuilding phase. The January 27 cutoff date also coincides with data sets that are available from authoritative relief agencies that can be used in future research. The removal of records resulted in a database with 2,608 values ranging from January 12, 2010 to January 27, 2010. The original database contained 8 main categories and 42 subcategories at the time that it was downloaded. See Table 3 and Table 4 for category and subcategory definitions (Note that these Tables do not represent the exact original subcategories as will be explained below). In order to more easily sort the database by category, a binary field was created for each category (a “1” indicating that a category relates to the record and a “0” indicating that a category does not relate to the record), adding 8 new fields to the database. Structured Query Language (SQL) was used to identify each record containing a category and the appropriate matching binary field was populated with a 1. This process was repeated for the remaining seven categories. Forty-two binary fields were created for each subcategory as well and they were populated using the same technique described above.

Any categories not in English were removed resulting in 28 records being removed and 2,580 remaining (99 percent of the total). Next, records that were obvious repeats were removed from the database. Criteria for determining if a record was a repeat were based on temporal proximity to each other and identical language. This resulted in 90 records being removed and 2,490 records remaining in the database (95.5 percent of the total). See Figure 5 for a breakdown of the number of records in the database at each stage described above.

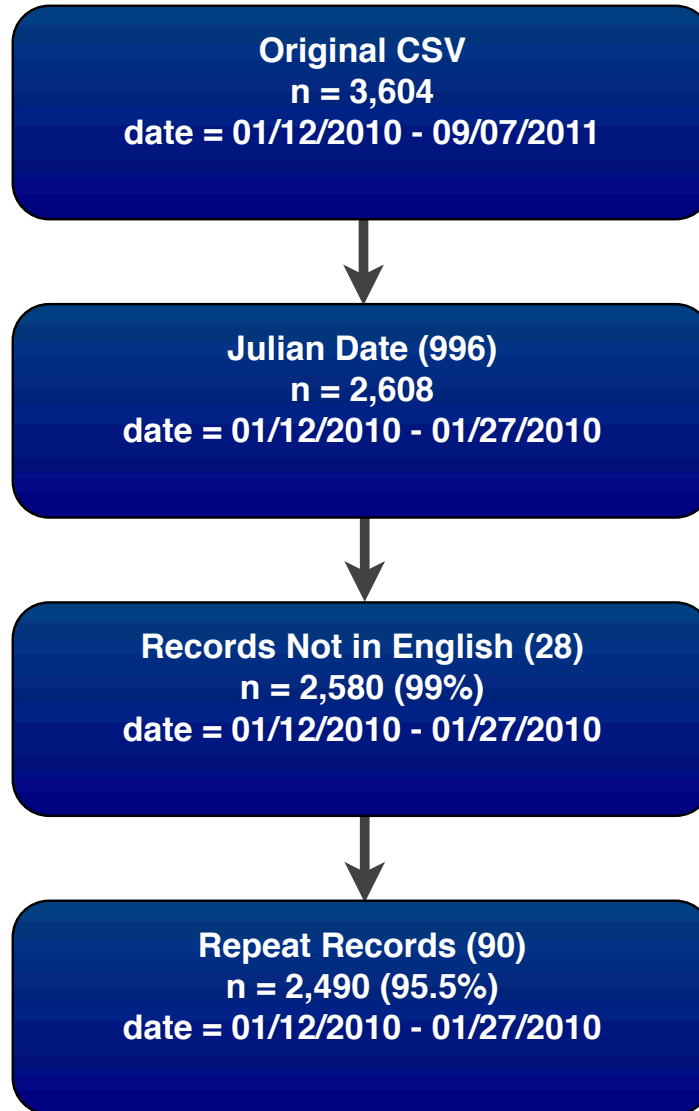


Figure 5 - Breakdown of how many records remain after each step in removing records not in English and duplicate records

Table 3 - Main Categories in Ushahidi Haiti Deployment

Category Number	Title
1	Emergency
2	Vital Lines
3	Public Health
4	Security Threats
5	Infrastructure Damage
6	Natural Hazards
7	Services Available
8	Other

Table 4 - Subcategories for Ushahidi Haiti Deployment

Category Number	Title	Category Number	Title
1a	Highly Vulnerable	5a	Collapsed Structure
1b	Medical Emergency	5b	Unstable Structure
1c	People Trapped	5c	Road Blocked
1d	Fire	5d	Compromised Bridge
2a	Food Shortage	5e	Communication Lines Down
2b	Water Shortage	6a	Floods
2c	Contaminated Water	6b	Landslides
2d	Shelter Needed	6c	Earthquakes and Aftershocks
2e	Fuel Shortage	7a	Food Distribution Point
2f	Power Outage	7b	Water Distribution Point
3a	Infectious Human Disease	7c	Non-Food Distribution Point
3b	Chronic Care Needs	7d	Hospital/Clinics Operating
3c	Medical Equipment and Supply Needs	7e	Feeding Centers Available
3d	OBGYN/Women's Health	7f	Shelter Offered
3e	Psychiatric Need	7g	Human Remains Management
3f	Water Sanitation and Hygiene Promotion	7h	Rubble Removal
3g	Deaths	8a	IDP Concentration
4a	Looting	8b	Aid Manipulation
4b	Theft of Aid	8c	Price Gouging
4c	Group Violence	8d	Search and Rescue
4d	Riot	8e	Person News
4e	Security Concern	8f	Other
		8g	Missing Persons
		8h	Asking to Forward a Message

### Re-categorization

In order to check for consistency among the categories and subcategories assigned to the victims' messages, I needed to know the number of times each category and

subcategory were used in the original database. I also had to have a data set to compare the original number of entries per category and subcategory with. To accomplish this goal I needed to re-categorize each message in the database for main and subcategories. I wanted my re-categorization to be as consistent with the original volunteers' as possible, so I reached out to the volunteer community to ask for any categorization guides, definitions, or training material that they were provided with. I quickly discovered that no such documents were created or used during the emergency (J. Valuch, personal communication, January 15, 2013). This led me to develop my own definitions and guidelines for each category and subcategory in order to consistently re-categorize the records within the database (See Appendix A for my category and subcategory definitions). All references to the original categories and subcategories were removed in order to reduce the likelihood of bias prior to my re-categorization. Each of the 2,490 remaining records in the database was re-categorized based on the rules in Appendix A. Re-categorization was repeated to check for consistency. Records whose re-categorization did not match were checked again. The two narratives below are examples of how I conducted my re-categorization.

#### Re-categorization Narrative:

Record # 516

Contents of Message: Carrefour, Fontamura, Bizoton, Thor: Hopital Adventiste de Diquini (Haitian Adventist Hospital) is treating and receiving patients in and around Carrefour (Fontamara, Bizoton, Thor, etc).

Original Categories from Volunteers: 7d. Hospital/clinics operating

My Categories: 7d. Hospital/clinics operating

Category Definition (from researcher): reports that medical services are being provided at this location



**Table 5 - Subcategories with more than one definition**

Subcategory	Definition 1	Definition 2
2c	Contaminated water	Security Concern
6a	Floods	Deaths
6b	Landslides	Missing Persons
6c	Earthquakes and Aftershocks	Asking to forward a message

In other instances, subcategories were moved to more closely match their more appropriate main category. See Table 6 below for examples.

**Table 6 - Subcategories Moved to Different Main Categories**

Original Subcategory & Definition	Original Category	New Category	New Subcategory
2c – Security Concern	2 - Vital Lines	4 – Security Threats	4e
4e - Water Sanitation	4 - Security Threats	3 – Public Health	3f
6a – Deaths	Natural Hazards	3 – Public Health	3g
6b – Missing Persons	Natural Hazards	8 – Other	8g
6c – Asking to Forward a Message	Natural Hazards	8 – Other	8h

### **Category and Subcategory Comparisons**

My main goal for this project was to compare the instances of original categories and subcategories in the original database to the instances of re-categorized categories and subcategories that I created. In order to conduct this comparison, I needed to compare which categories and subcategories the original volunteers associated with each record in the database to my own re-categorization for each record. In order to automate this process, a Python script (see Appendix B and Appendix C) was used to compare the volunteer produced categories and subcategories for each record to those generated by



me. This Python script resulted in a new CSV file that contained the unique identifier, a column containing the original, volunteer produced categories and subcategories, and a column containing the categories and subcategories I generated per record. A spreadsheet software program was used to generate descriptive statistics to determine how many records were a perfect match between the two data sets at both the category level and subcategory level. See Table 7 below as an example for what the output of the Python script looks like as well as a description of what the columns represent. The original volunteers determined that record 67 (Column 1) in the database related to categories 1, 7, and 8 (Column 2). When I re-categorized the record, I found that the message only related to category 8 (Column 3). The Python script then determined that for record 67, the original categorization and my conducted re-categorization both related to category 8 (Column 4). The last two columns indicate which categories appeared in only the original volunteer categorization and my categorization. The table indicates that for record 68, volunteers and myself agreed that the message related to category 6 while the volunteers thought the message also applied to category 7 and I thought the message also applied to category 8. The table indicates that for record 69, there was no agreement between the original volunteers and myself

**Table 7 – Sample CSV generated by Python script comparing original and re-categorized results per record.**

1	2	3	4	5	6
Unique ID	Original Categories (each digit represents a category)	Revised Categories (each digit represents a category)	Both Only (each digit represents a category)	Original Only (each digit represents a category)	Revised Only (each digit represents a category)
67	178	8	8	17	
68	67	68	6	7	8
69	125	8		125	8

### **Statistics**

Descriptive statistics were generated for each original and revised category and subcategory in aggregate and by day. These descriptive statistics included total numbers of instances for each category and subcategory as well as percentages of totals used to compare the original and re-categorized values for the main and subcategories.

## Chapter Four – Results

The results section of this thesis is separated into three subsections. The first section focuses on the aggregate of the results at the main category and subcategory level. The second section focuses on results aggregated by day at the main category and subcategory level. The third section focuses on statistical comparisons at various levels for the main and subcategories between the original number of entries and the re-categorized number of entries.

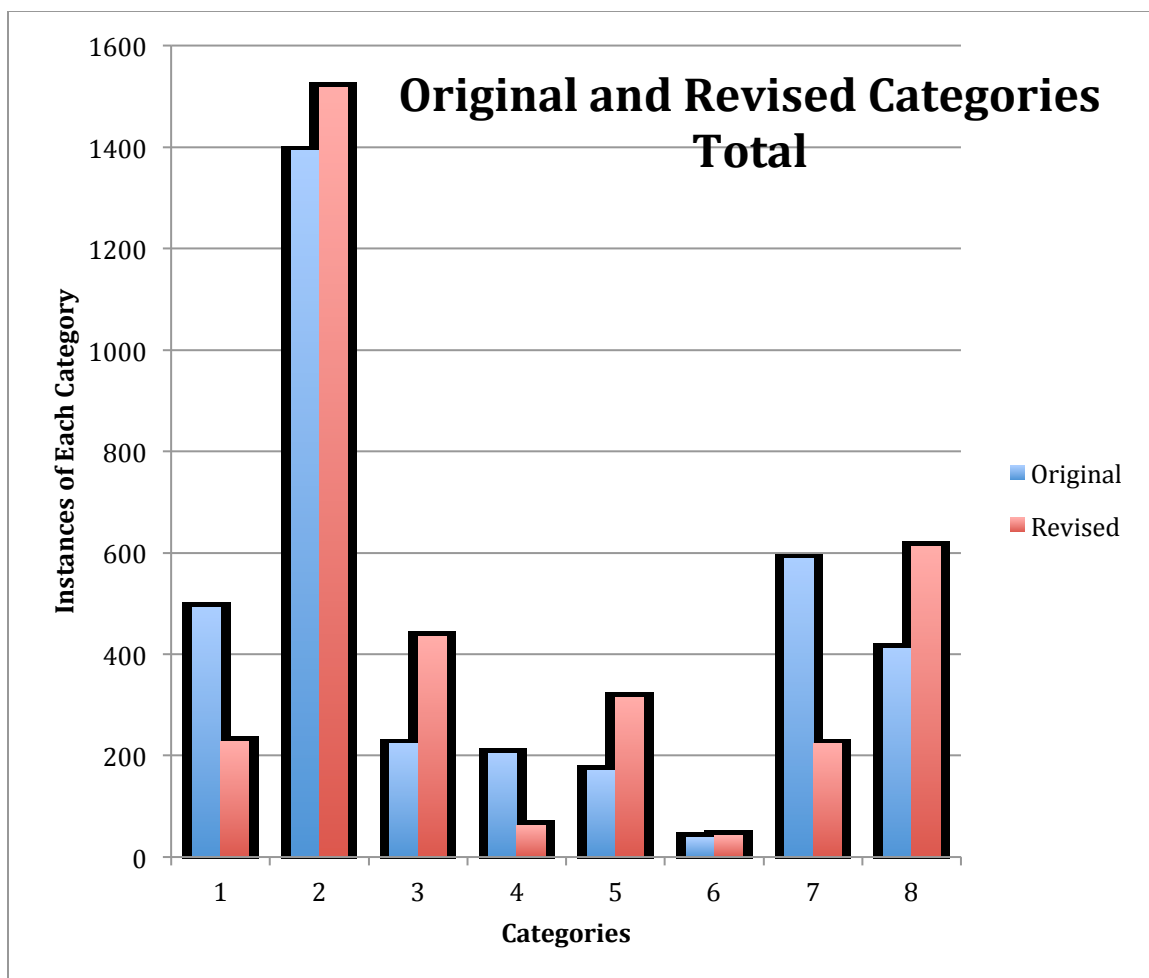
### Aggregate Results

#### Main Categories

The aggregate results for the main categories show that the number of differences in entries for each category and the total number of entries for each category varied widely across categories and subcategories. Refer to Table 8 and Figure 6.

**Table 8 - Instances of Each Main Category in Total**

Category	Number of Entries in Each Category for the Original Data Set	Number of Entries in Each Category After Re-categorization
1	492	228
2	1393	1518
3	223	435
4	204	62
5	171	315
6	38	42
7	588	223
8	411	613
total	3520	3436



**Figure 6 - Instances of Each Main Category in Total**

Category 1, Emergency, had significantly fewer entries after re-categorization than before. Category 1 comprised 13.98 percent of the total original entries and 6.64 percent of the total revised entries. Category 2, Vital Lines, contained the highest percentage of the total records in both the original data set and the revised data set. Category 2 comprised 39.57 percent of the total original entries and 44.18 of the total revised entries. Category 3, Public Health, had a significant increase in the number of entries after re-categorization. Category 3 comprised 6.34 percent of the total original entries and 12.66 percent of the total revised entries. Category 4, Security Threats, appears to have significantly fewer entries after re-categorization than before. Category 4

comprised 5.80 percent of the total original entries and 1.80 percent of the total revised entries. Category 5, Infrastructure Damage, had significantly more entries after re-categorization than before. Category 5 comprised 4.86 percent of the total original entries and 9.17 percent of the total revised entries. Category 6, Natural Hazards, did not have a significant difference in the number of entries when comparing the original categorization to the re-categorization. Category 6 contained the fewest number of entries in the original and revised data sets with the original comprising 1.08 percent of the total and the revised comprising 1.22 percent of the total. Category 7, Services Available, had significantly fewer entries after re-categorization than before. Category 7 comprised 16.70 percent of the original entries and 6.49 percent of the total revised entries. This resulted in Category 7 having the second most entries in the original data set to the third least entries in the revised data set. Category 8, Other, had significantly more entries after re-categorization than before. Category 8 comprised 11.68 percent of the total original entries and 17.84 percent of the total revised entries.

There are several reasons that the number of instances for the categories may be inconsistent. Many of the discrepancies highlighted above can be attributed to omission, commission, or both when comparing the original data set to my revised data set. If the revised data set is considered to be more consistent, and therefore higher quality, then inconsistencies of omission are when I assigned a category to a record that the original volunteers did not. Inconsistencies of commission are when the volunteers assigned a category to a record that I did not. However, these types of inconsistencies do not fully explain the discrepancies I observed. For instance, when examining the original categories and subcategories prior to re-categorization, I realized that some subcategories

appeared to be in the wrong main categories. The next subsection will highlight these differences.

**Subcategories**

The number of entries within each subcategory both for the original data set and the revised data set are highly variable. Even within Category 2, the largest category both before and after re-categorization, there were some individual subcategories with very few entries. Some of this is attributed to omission and commission inconsistencies between the original and revised data set while others are due to the rearranging of some subcategories. Rearranging of subcategories will be highlighted below as appropriate. Refer to Figure 7 below.

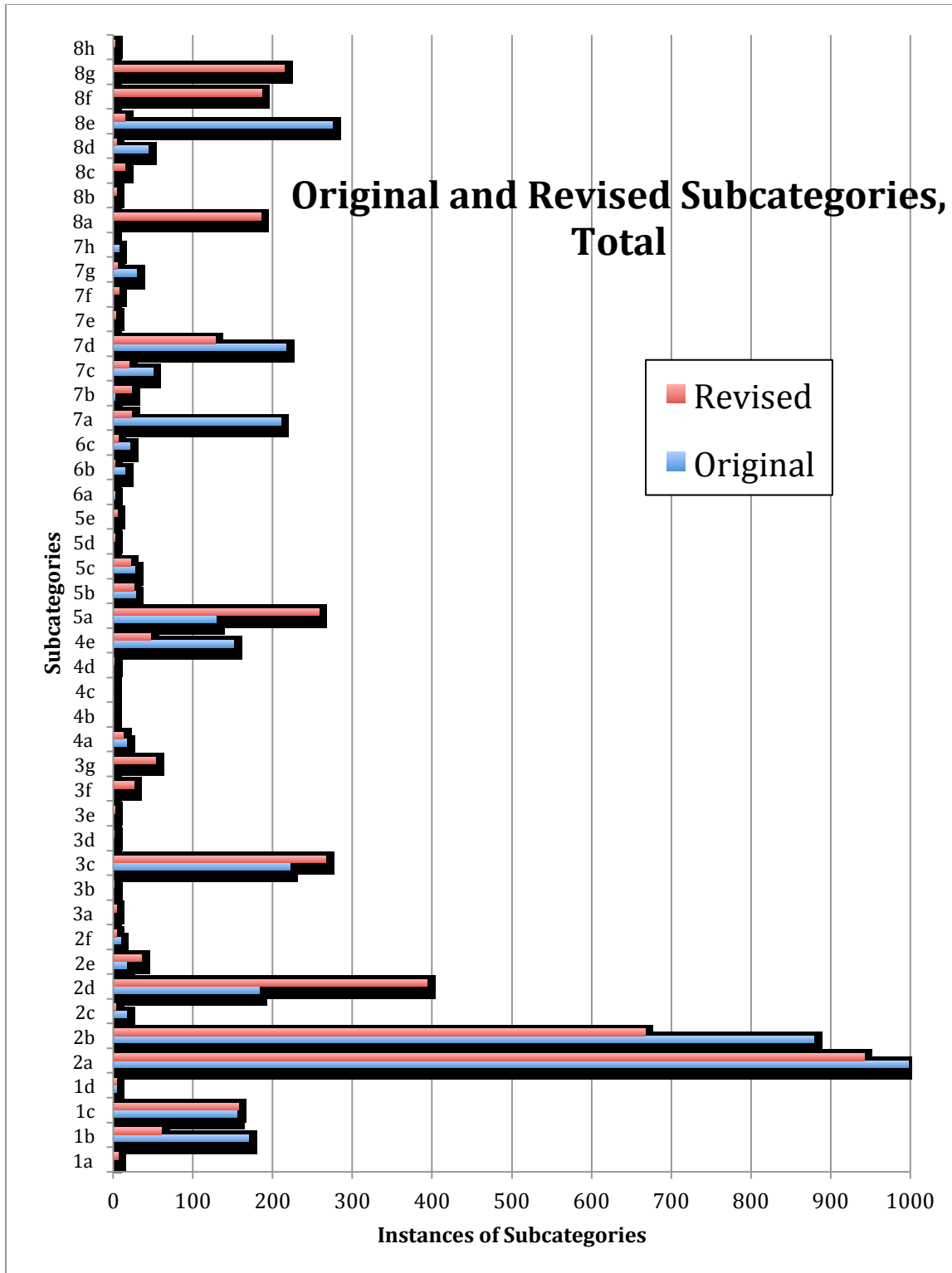


Figure 7 – Instances of Each Subcategory, Total

Category 1 consisted of four subcategories related to emergencies. Refer to Table 9. Subcategory 1d had very few entries both before and after re-categorization but had the exact same number of entries in both cases. Subcategory 1a had no original entries but 6 entries in the revised data set. Subcategories 1b, Medical Emergencies, and 1c, People Trapped, comprised the majority of the entries both before and after re-categorization. There were no discrepancies between the original and revised subcategories for Category 1.

**Table 9 - Subcategory 1, Total Entries**

Subcategory	Number of Entries in Each Subcategory for the Original Data Set	Number of Entries in Each Subcategory after Re-categorization
1a	0	6
1b	170	61
1c	155	157
1d	4	4

Category 2 consisted of 6 subcategories related to needs. Refer to Table 10. Subcategories 2a, Food Shortage, and 2b, Water Shortage, comprised the overwhelming majority of the entries both before and after re-categorization. There was one discrepancy between the original and revised subcategories for Category 2.

**Table 10 - Subcategory 2, Total Entries**

Subcategory	Number of Entries in Each Subcategory for the Original Data Set	Number of Entries in Each Subcategory after Re-categorization
2a	998	942
2b	879	667
2c	17	3
2d	183	394
2e	17	36
2f	9	4



Category 3 was originally comprised of 5 subcategories related to public health, however, after examining all of the data I moved two subcategories into this category. Refer to Table 11. Original category 4e, Water Sanitation and Hygiene Promotion, became 3f and 6a, Deaths, became 3g. Subcategories 3a, 3b, 3d, and 3e contained no entries in the original data set and a very limited number of entries in the revised data set. Category 3c, Medical Equipment and Supply Needs, comprised the overwhelming majority of entries in both the original and revised data sets. Due to subcategories 3f and 3g being moved into Category 3 during re-categorization, they have nothing to be compared with.

**Table 11 - Subcategory 3, Total Entries**

Subcategory	Number of Entries in Each Subcategory for the Original Data Set	Number of Entries in Each Subcategory after Re-categorization
3a	0	4
3b	0	1
3c	222	267
3d	0	1
3e	0	2
3f	0	26
3g	0	53

Category 4 originally consisted of 5 subcategories related to security threats. Refer to Table 12. Subcategory 4e was moved to Category 3 and replaced with subcategory 2c (In the original data set there were two 2c's. One referenced Contaminated Water while the other referenced Security Concern.). Original subcategories 4b, 4c, and 4d did not contain any entries in the original dataset and only 4d contained an entry in the revised data set. The majority of entries in the original and revised data sets related to original category 4a and revised category 4e.

Table 12 - Subcategory 4, Total Entries

Subcategory	Number of Entries in Each Subcategory for the Original Data Set	Number of Entries in Each Subcategory after Re-categorization
4a	17	13
4b	0	0
4c	0	0
4d	0	1
4e	151	47

Category 5 consisted of 5 subcategories related to infrastructure damage. Refer to Table 13. The majority of entries in both the original and revised data sets related to subcategory 5a followed by 5b and 5c. There were no original entries in subcategories 5d and 5e and only a few in the revised data set. There were no discrepancies between the original and revised subcategories for Category 5.

Table 13 - Subcategory 5, Total Entries

Subcategory	Number of Entries in Each Subcategory for the Original Data Set	Number of Entries in Each Subcategory after Re-categorization
5a	129	258
5b	28	26
5c	27	22
5d	0	2
5e	0	5

Category 6 was a challenging category. Refer to Table 14. The original data set consisted of two definitions each for Category 6's three subcategories, which were supposed to relate to natural hazards. As a result, the second 6a, Deaths, was moved to Category 3, the second 6b, Missing Persons, was moved to Category 8, and the second 6c, Asking to Forward a Message, was moved to Category 8. The remaining three subcategories related to natural hazards. The overall number of entries for Category 6 in

both the original and revised data sets was very small with no instances of subcategory 6a in the revised data set.

**Table 14 - Subcategory 6, Total Entries**

Subcategory	Number of Entries in Each Subcategory for the Original Data Set	Number of Entries in Each Subcategory after Re-categorization
6a	2	0
6b	15	2
6c	21	6

Category 7 consisted of 8 subcategories related to services available. Refer to Table 15. The two largest subcategories in both the original and revised data sets were 7a, Food Distribution Points, and 7d, Hospital/Clinics Operating. There were no original entries for subcategories 7e and 7f and there were no revised entries for subcategory 7h. There were no discrepancies between the original and revised subcategories for Category 7.

**Table 15 - Subcategory 7, Total Entries**

Subcategory	Number of Entries in Each Subcategory for the Original Data Set	Number of Entries in Each Subcategory after Re-categorization
7a	210	23
7b	2	23
7c	50	20
7d	217	128
7e	0	3
7f	0	7
7g	29	5
7h	7	0

Category 8 originally consisted of 6 subcategories that did not fit in any other main category. Refer to Table 16. After I moved the second subcategories 6b and 6c, however, Category 8 consisted of 8 subcategories. Four of the original subcategories did not contain any entries with most of the original entries going to subcategory 8e, Person

News. Somewhat surprisingly, there were no original entries for subcategory 8a, Internally Displaced Persons (IDP) Concentrations.

**Table 16 - Subcategory 8, Total Entries**

Subcategory	Number of Entries in Each Subcategory for the Original Data Set	Number of Entries in Each Subcategory after Re-categorization
8a	0	185
8b	0	4
8c	0	15
8d	44	4
8e	275	15
8f	0	186
8g	0	215
8h	0	2

### **Daily Results**

The next section of results relate to how the original and re-categorized data sets changed over the days following the earthquake. Figure 8 reflects the total number of reports that were filed with the Ushahidi platform following the earthquake without making any distinctions between classes. Note that the number of reports is declining starting on the 14<sup>th</sup> but begins to increase starting on the 16<sup>th</sup>, coinciding with the start of the 4636 short message (SMS) code program. The peak of the reports occurs on the 23<sup>rd</sup>, which is the date that the Haitian government called an end to the response phase of the disaster (Batty, 2010).

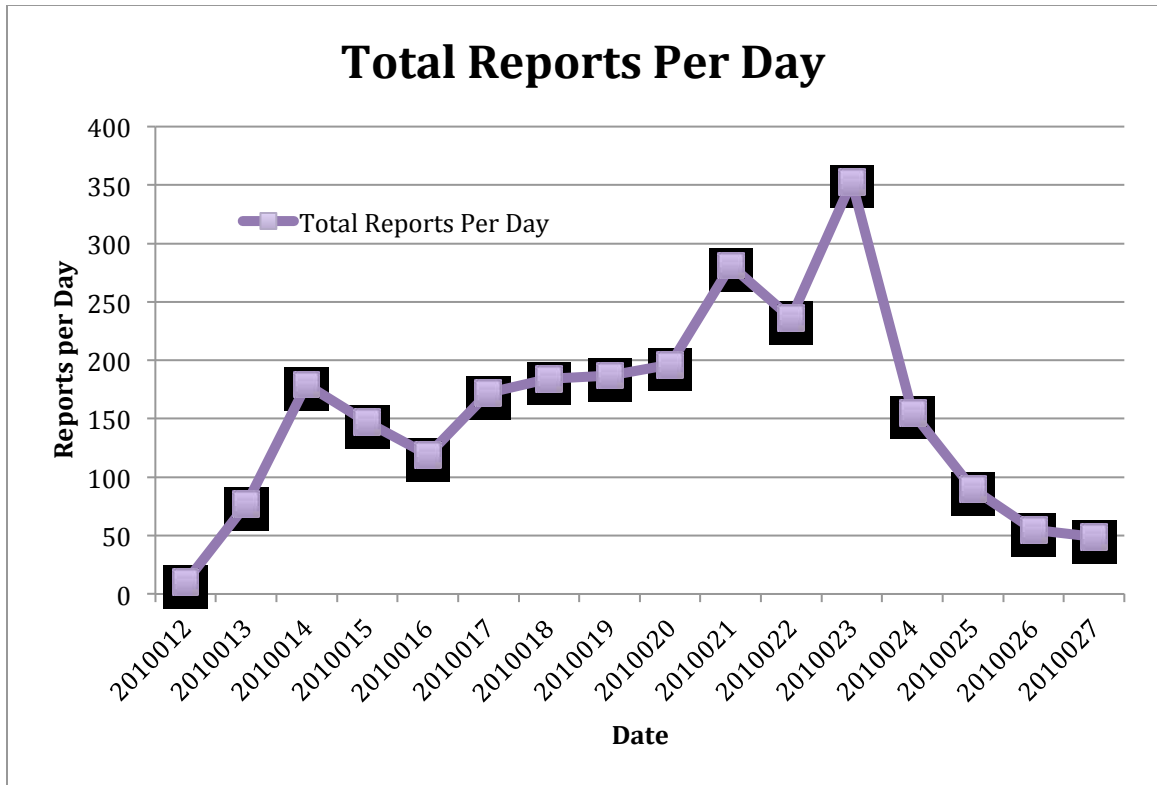


Figure 8 - Total Reports per Day

### Main Categories

Just as the overall number of entries in the original and revised data sets varied by category, they also varied by date. The following results pertain to how the number of entries in the original and revised main category data sets changed over time. Figures 9 and 10 compare the total entries in the original data set and the total entries in the revised data set by main category. In both cases notice the large number of reports related to needs (Category 2) compared to the other categories.

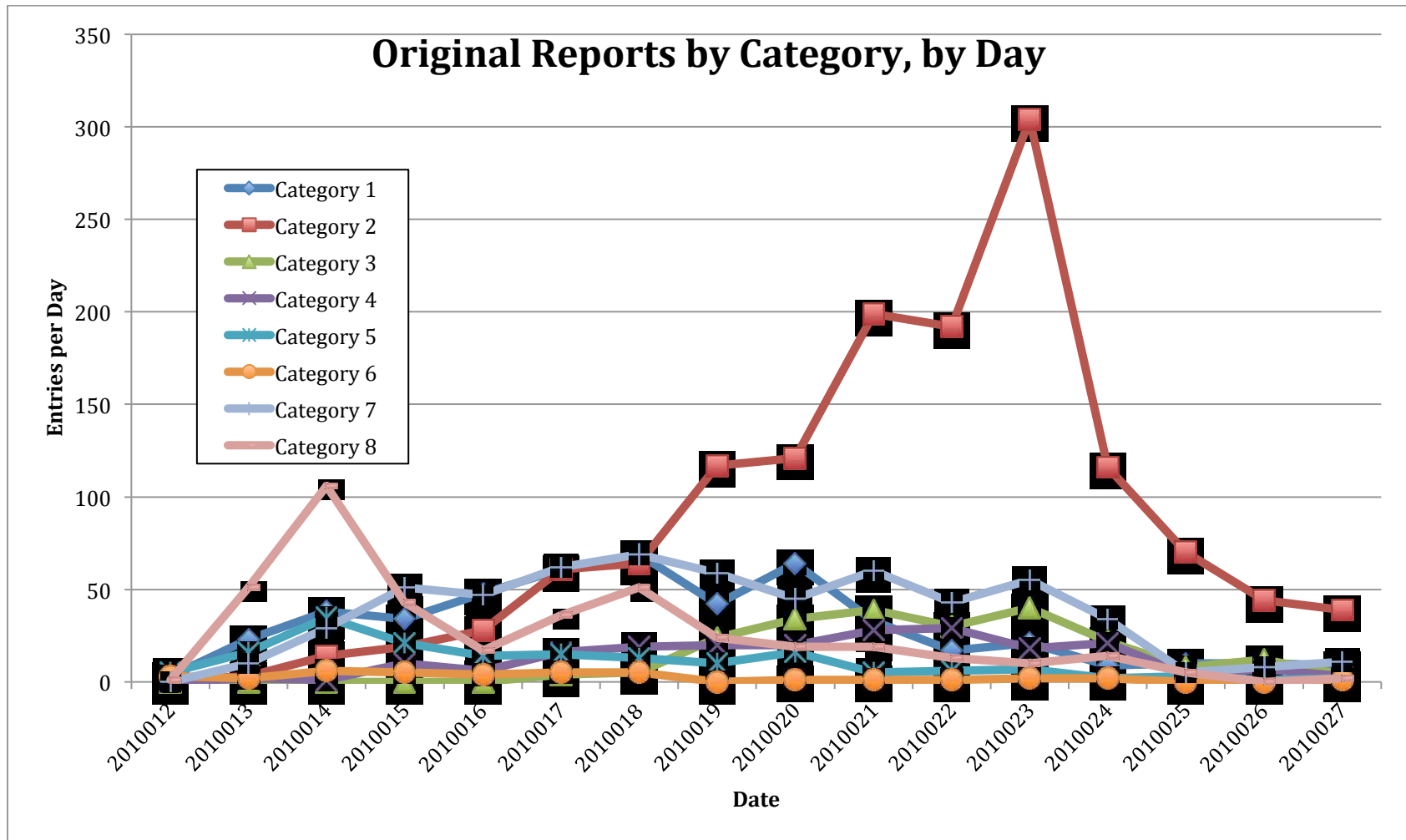


Figure 9 - Original Reports by Category per Day

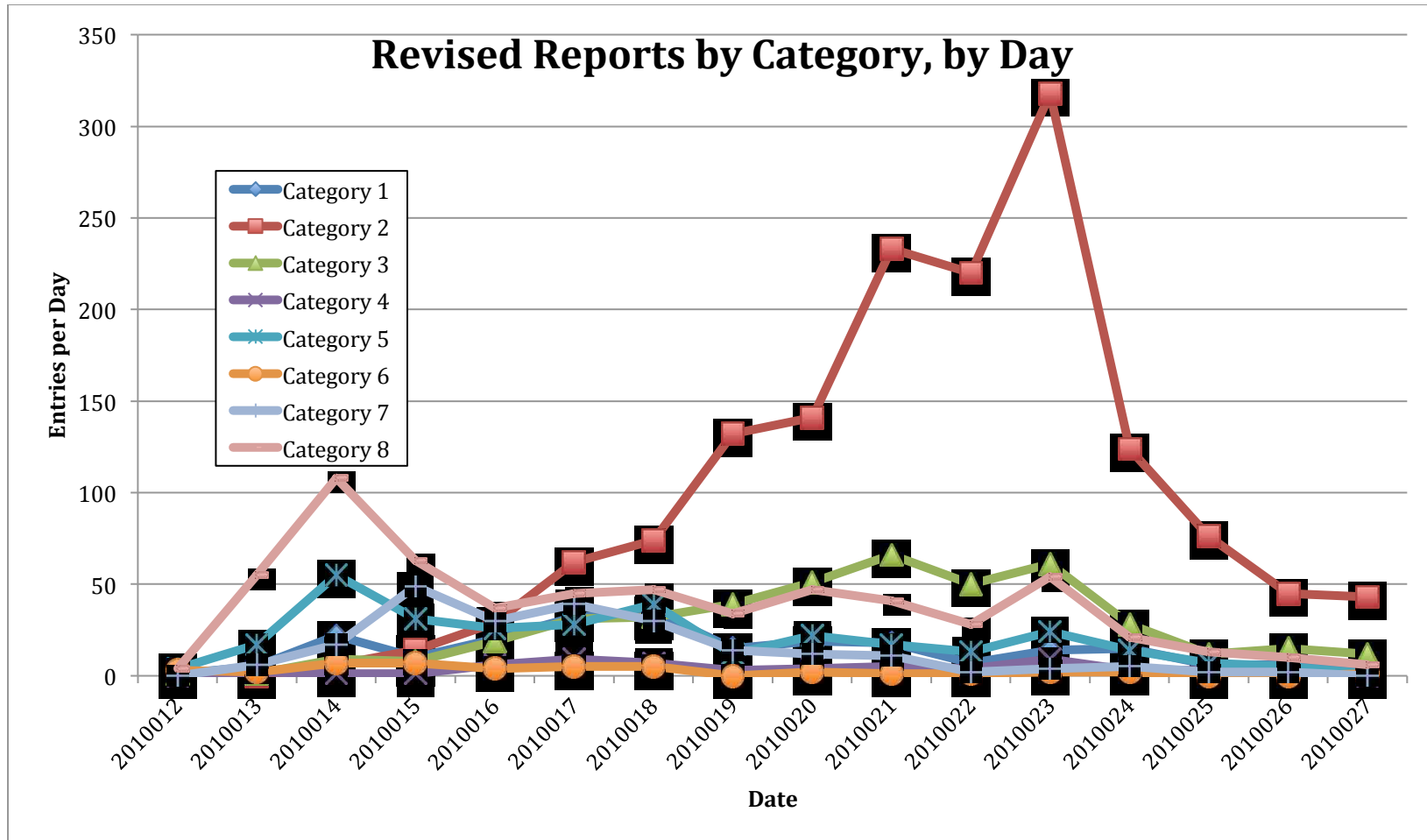


Figure 10 – Revised Reports by Category per Day

Another way to view these two sets of values is using a stacked area chart that highlights the relative number of entries in each category by day and provides a cumulative number of entries. Refer to Figures 11 and 12 below.



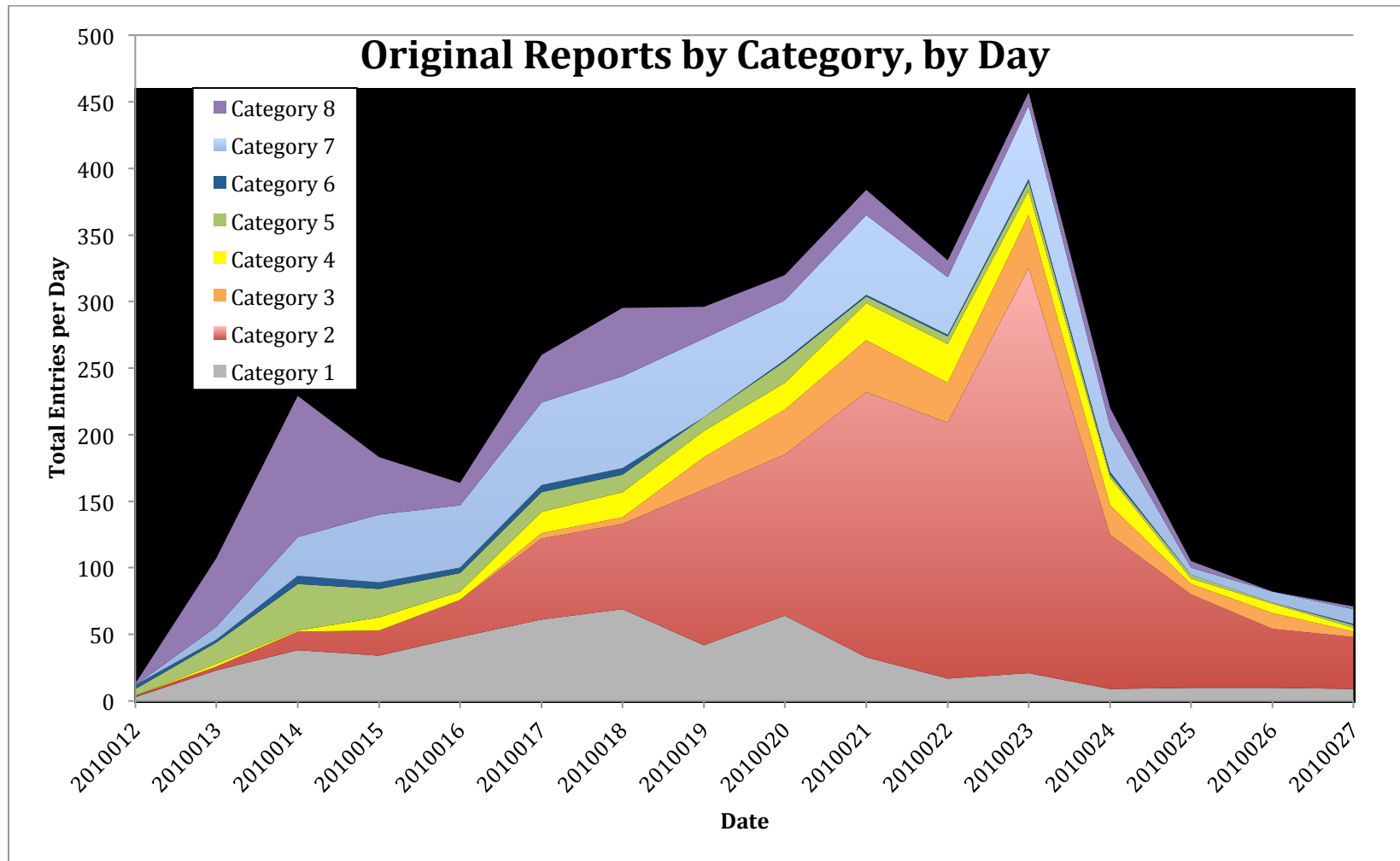


Figure 11 - Cumulative Entries per Day, Original

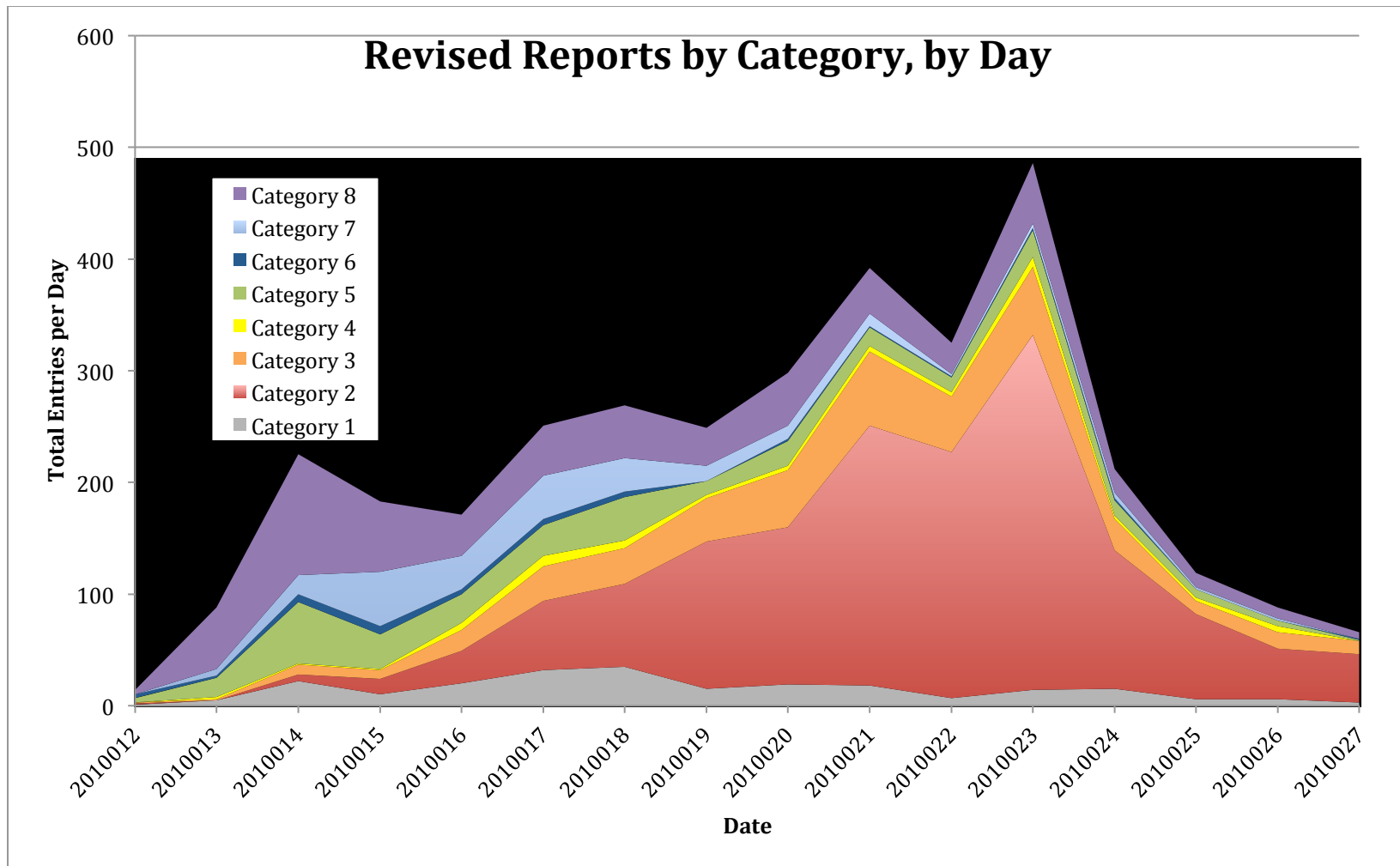


Figure 12 - Cumulative Entries by Day, Revised

The number of entries in the original and revised Category 1 data sets indicate discrepancies related to commission for every day in the study period except one. Refer to Figure 13. The discrepancy between the two data sets indicates a difference between the categorization by the researcher and the volunteers.

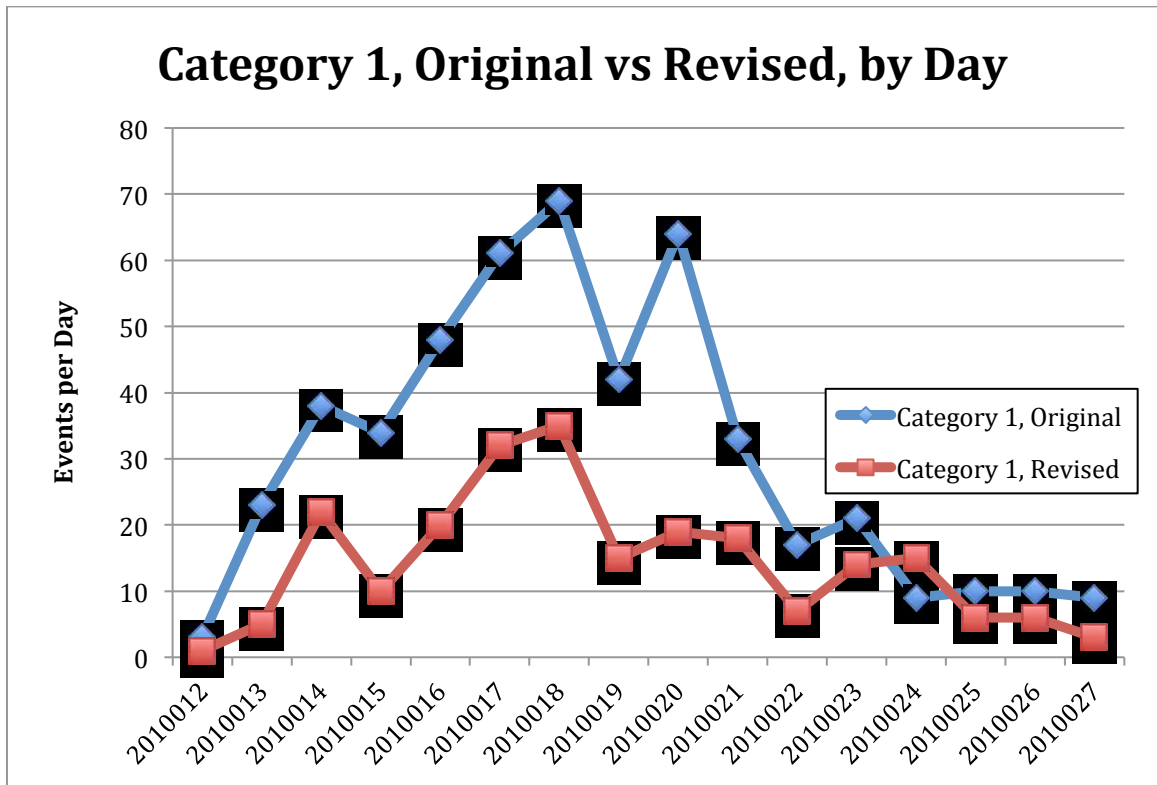


Figure 13 - Category 1, Original versus Revised, by Day

The number of entries in the original and revised Category 2 data sets rarely fluctuate, and when they do, by a relatively small amount. This consistency between the two data sets indicates a high degree of agreement between the categorization by the volunteers and the researcher. Refer to Figure 14.

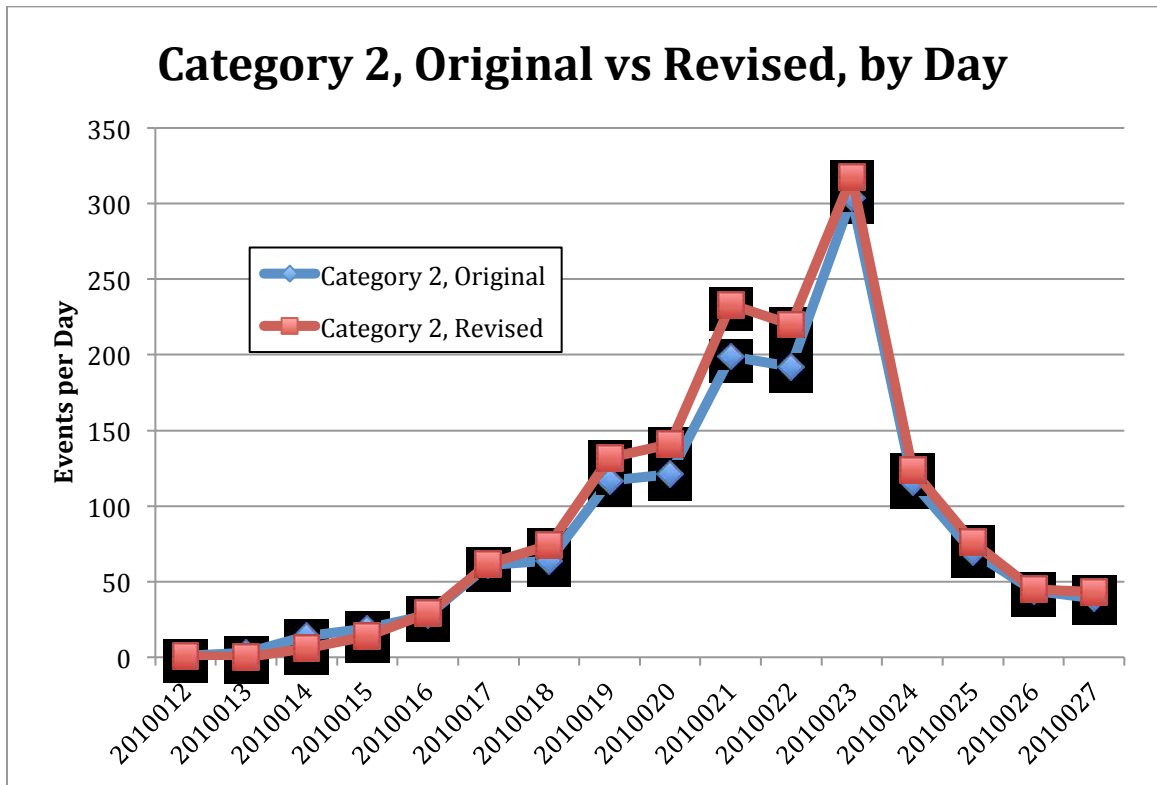


Figure 14 - Category 2, Original versus Revised, by Day

The number of entries in the original and revised Category 3 data sets indicates discrepancies related to omission for every day in the study period. Refer to Figure 15.

The discrepancy between the two data sets indicates a difference between the categorization by the researcher and the volunteers.

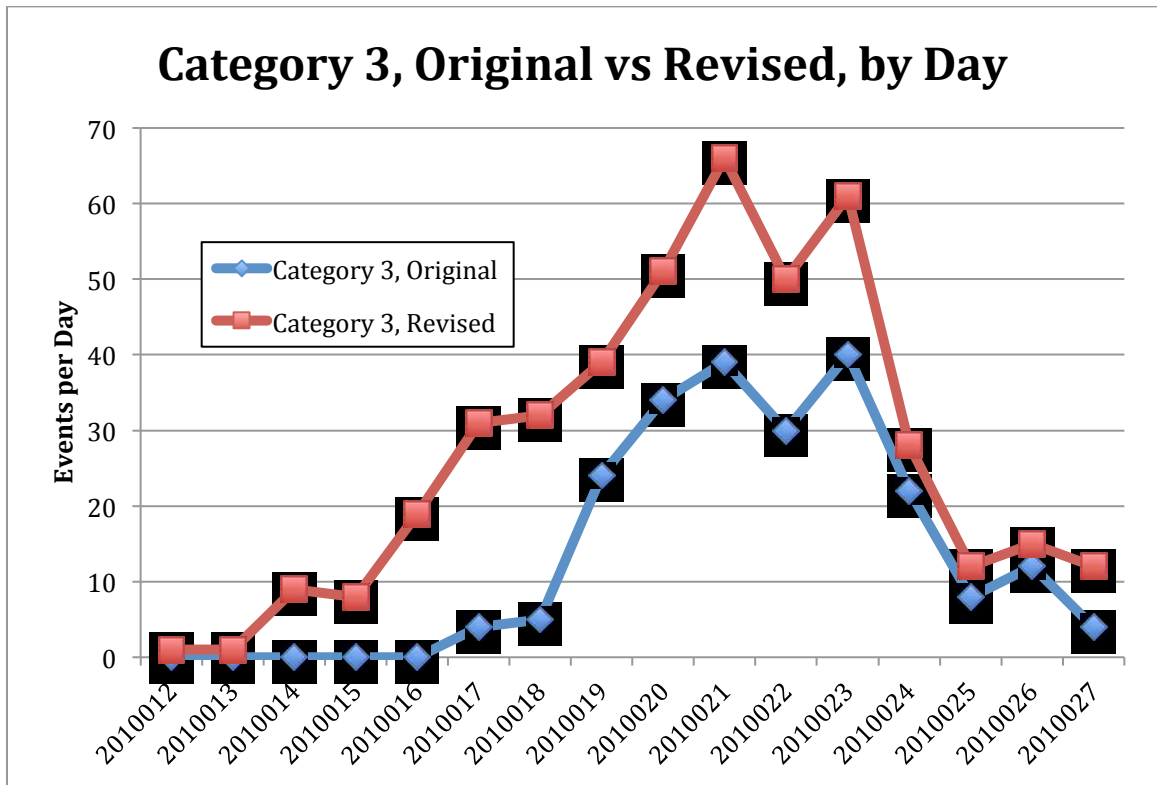


Figure 15 - Category 3, Original versus Revised, by Day

The number of entries in the original and revised Category 4 data sets indicates discrepancies related to commission for a majority of the days in the study period. Refer to Figure 16. The discrepancy between the two data sets indicates a difference between the categorization by the researcher and the volunteers.

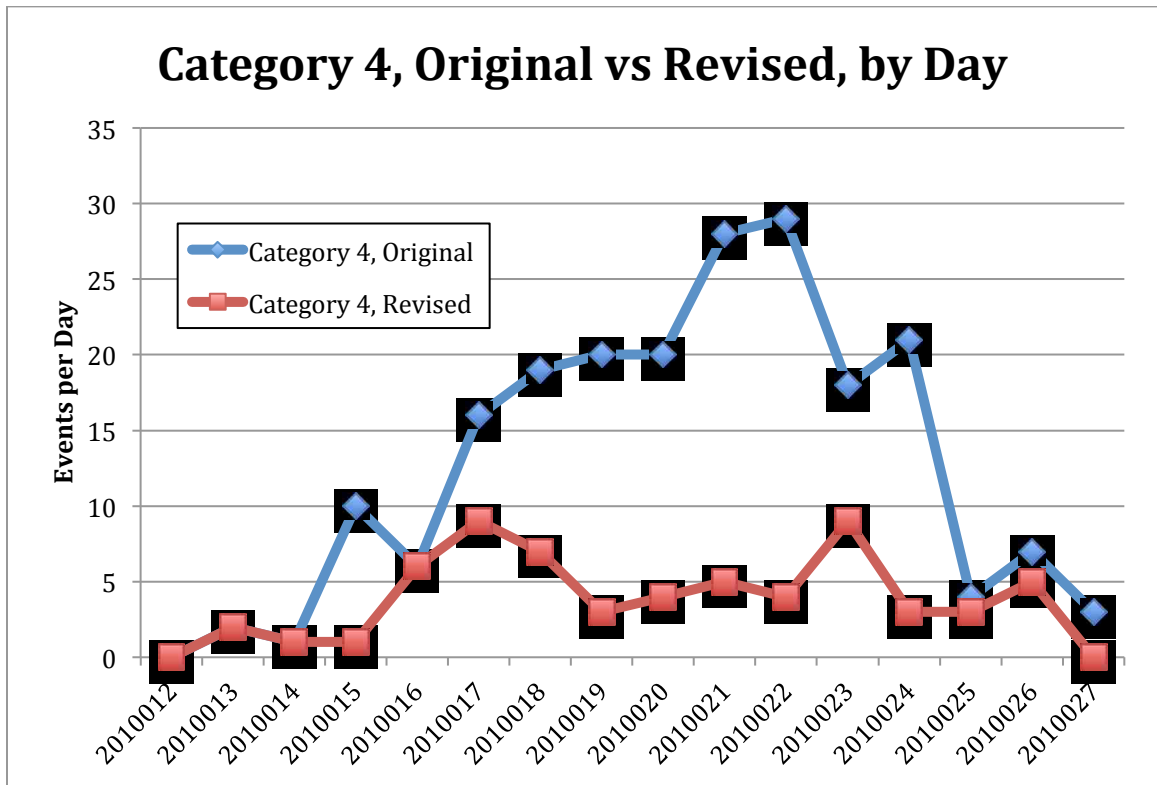


Figure 16 - Category 4, Original versus Revised, by Day

The number of entries in the original and revised Category 5 data sets indicates discrepancies related to omission for a majority of the days in the study period. Refer to Figure 17. The discrepancy between the two data sets indicates a difference between the categorization by the researcher and the volunteers.

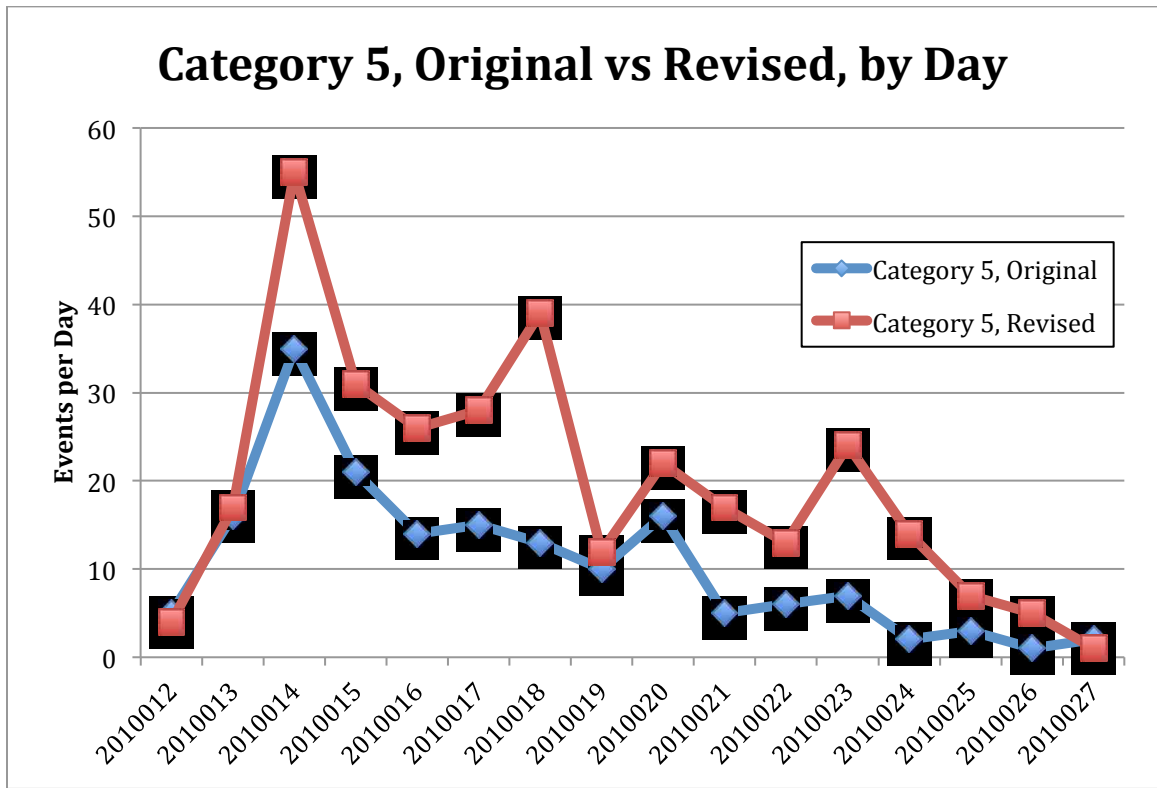


Figure 17 - Category 5, Original versus Revised, by Day

The results of comparing the number of entries in the original and revised Category 6 by day defy expectation based on what we now know about the inconsistencies in the original subcategories. This may be a result of the category's small size and should not be viewed as agreement between the categorization between the researcher and the volunteers. Refer to Figure 18.

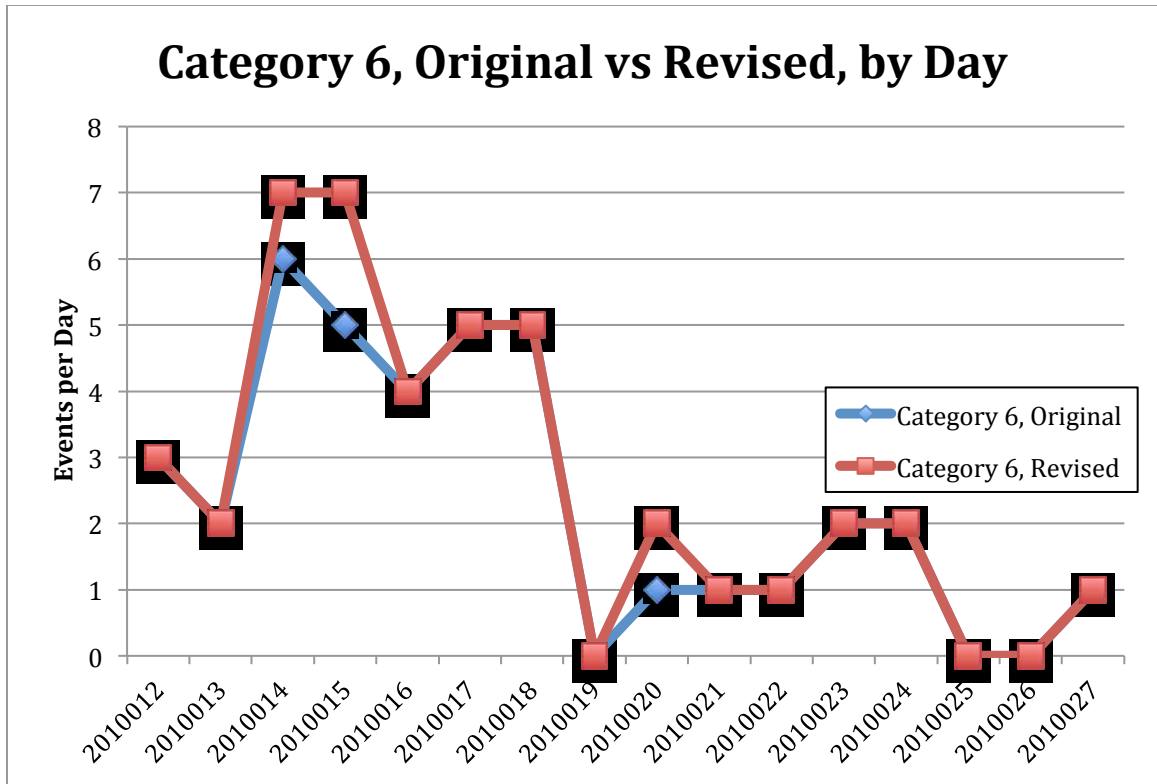


Figure 18 - Category 6, Original versus Revised, by Day



The number of entries in the original and revised Category 7 data sets indicates discrepancies related to commission for all but one day in the study period. Refer to Figure 19. The discrepancy between the two data sets indicates a difference between the categorization by the researcher and the volunteers.

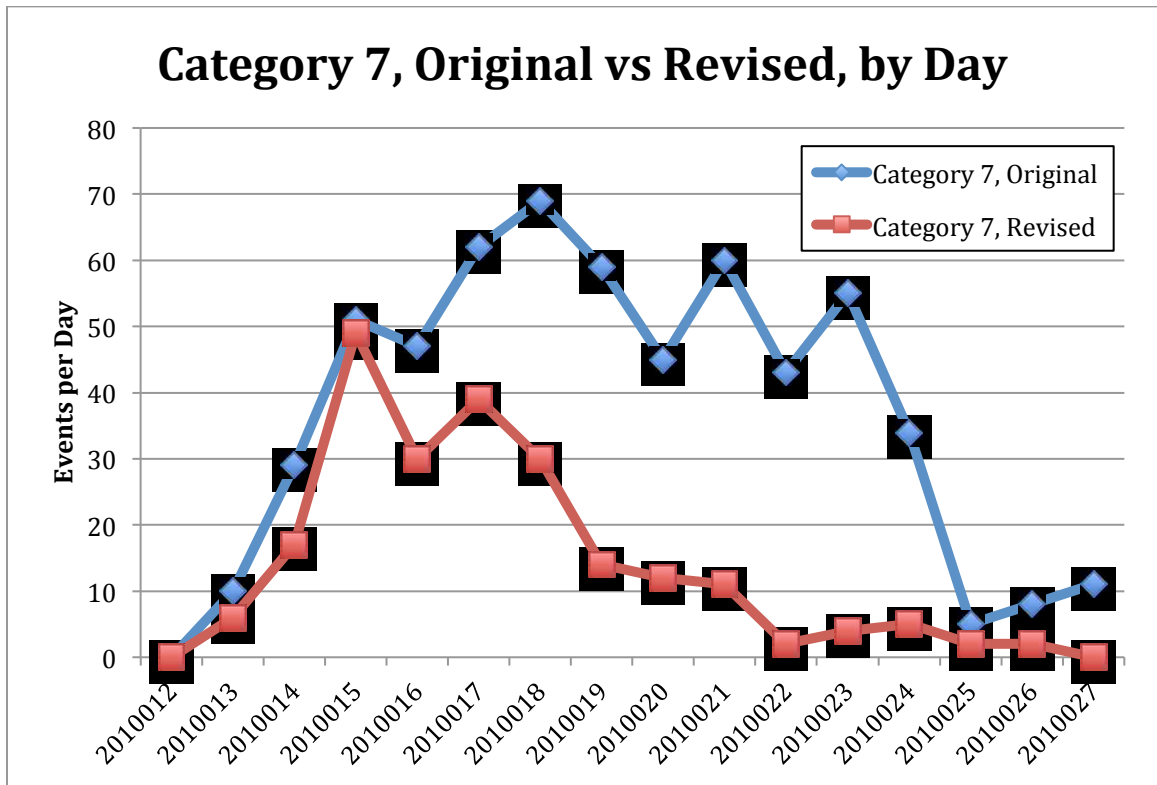


Figure 19 - Category 7, Original versus Revised, by Day

The number of entries in the original and revised Category 8 data sets indicates discrepancies related to omission for a majority of the days in the study period. Refer to Figure 20. The discrepancy between the two data sets indicates a difference between the categorization by the researcher and the volunteers.

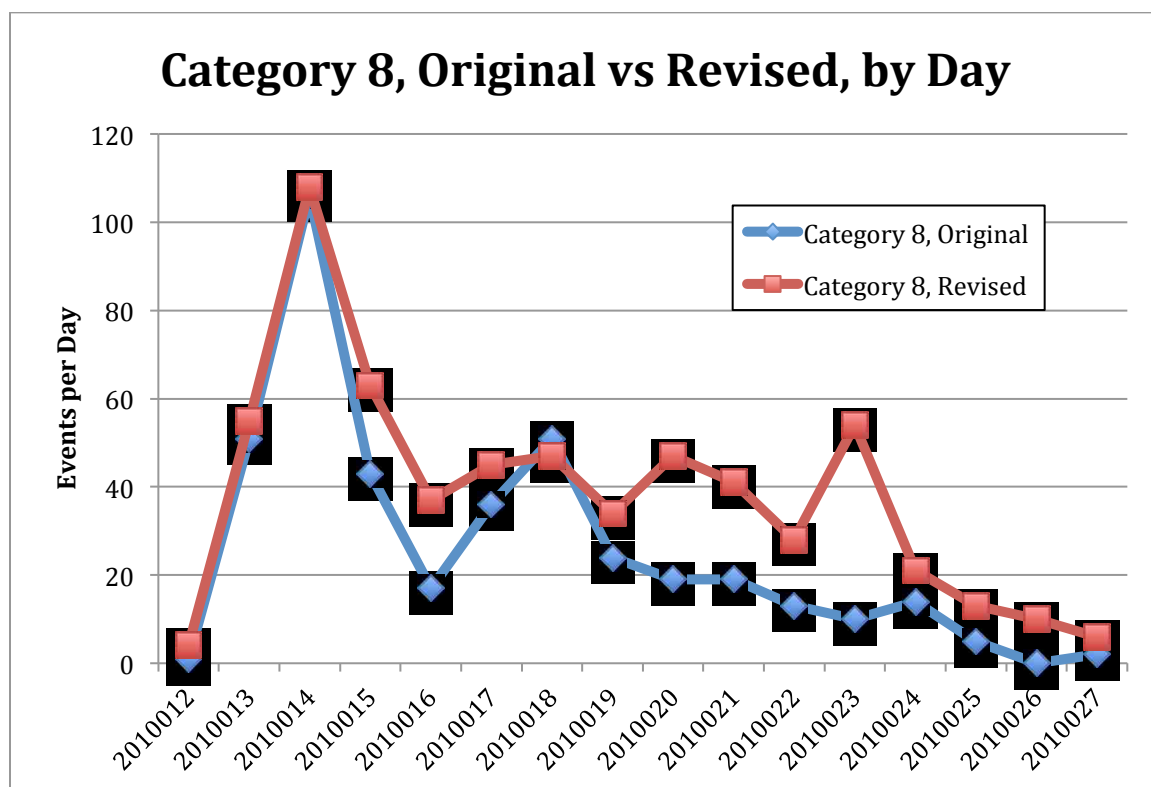


Figure 20 - Category 8, Original versus Revised, by Day

### Subcategories

The following results pertain to how the number of entries in the original and revised subcategory data sets changed over time. There are too many subcategories to display on one graph so a separate graph is displayed for the subcategories that pertain to each main category. The x and y axes for the graphs within a single main category are the same but the axes may change between categories. Full sized versions of each graph can be found in Appendix D.

The subcategories that comprise Category 1 highlight the differences in the number of entries among each of the subcategories as well as how the number of reports

fluctuates per day. Subcategory 1b shows a primary discrepancy by commission while Subcategory 1c shows a general agreement between the original and revised entries. Subcategories 1a and 1d consist of few events in both data sets. Refer to Figure 21.

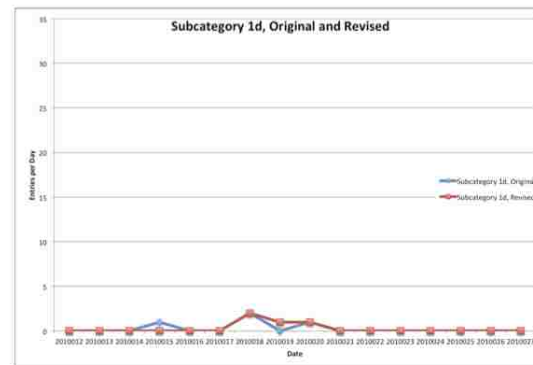
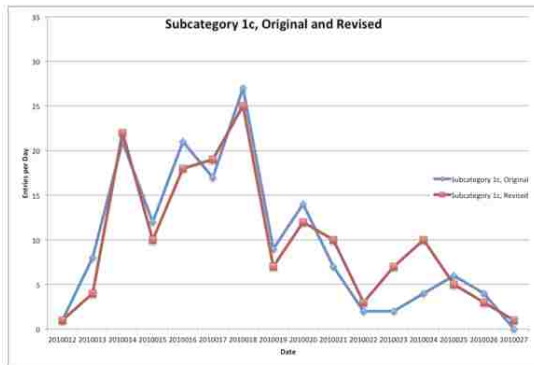
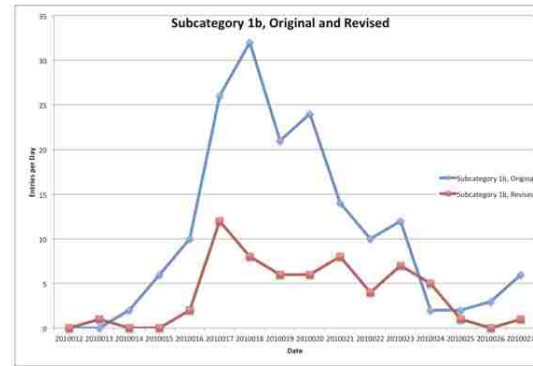
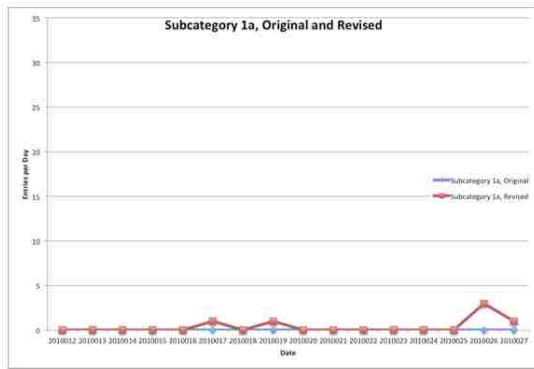


Figure 21 - Events for Subcategories 1a-1d by Day

The subcategories that comprise Category 2 highlight the differences in the number of entries among each of the subcategories as well as how the number of reports fluctuates per day. Subcategory 2b shows a primary discrepancy by commission, Subcategory 2d shows a primary discrepancy by omission, while Subcategory 2a shows a general agreement between the original and revised entries. Subcategories 2c, 2e, and 2f consist of few events in both data sets. Refer to Figure 22.

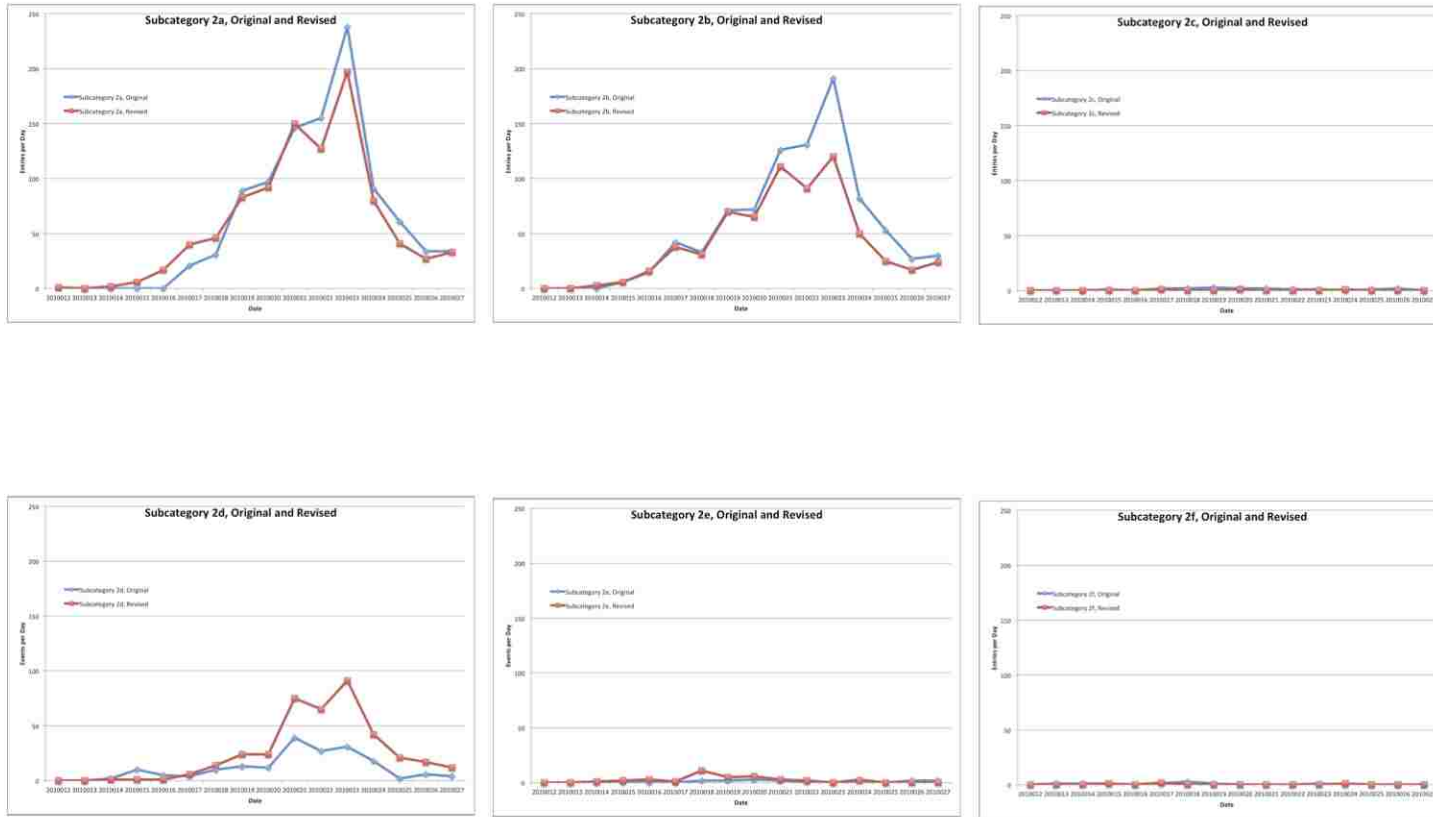


Figure 22 - Events for Subcategories 2a-2f by Day

The subcategories that comprise Category 3 highlight the differences in the number of entries among each of the subcategories as well as how the number of reports fluctuates per day. Subcategory 3c shows a primary discrepancy by commission. Subcategories 3a, 3b, 3d, and 3e consist of few events in both data sets. Subcategories 3f and 3g did not exist in the original data set so cannot be compared. Refer to Figure 23.

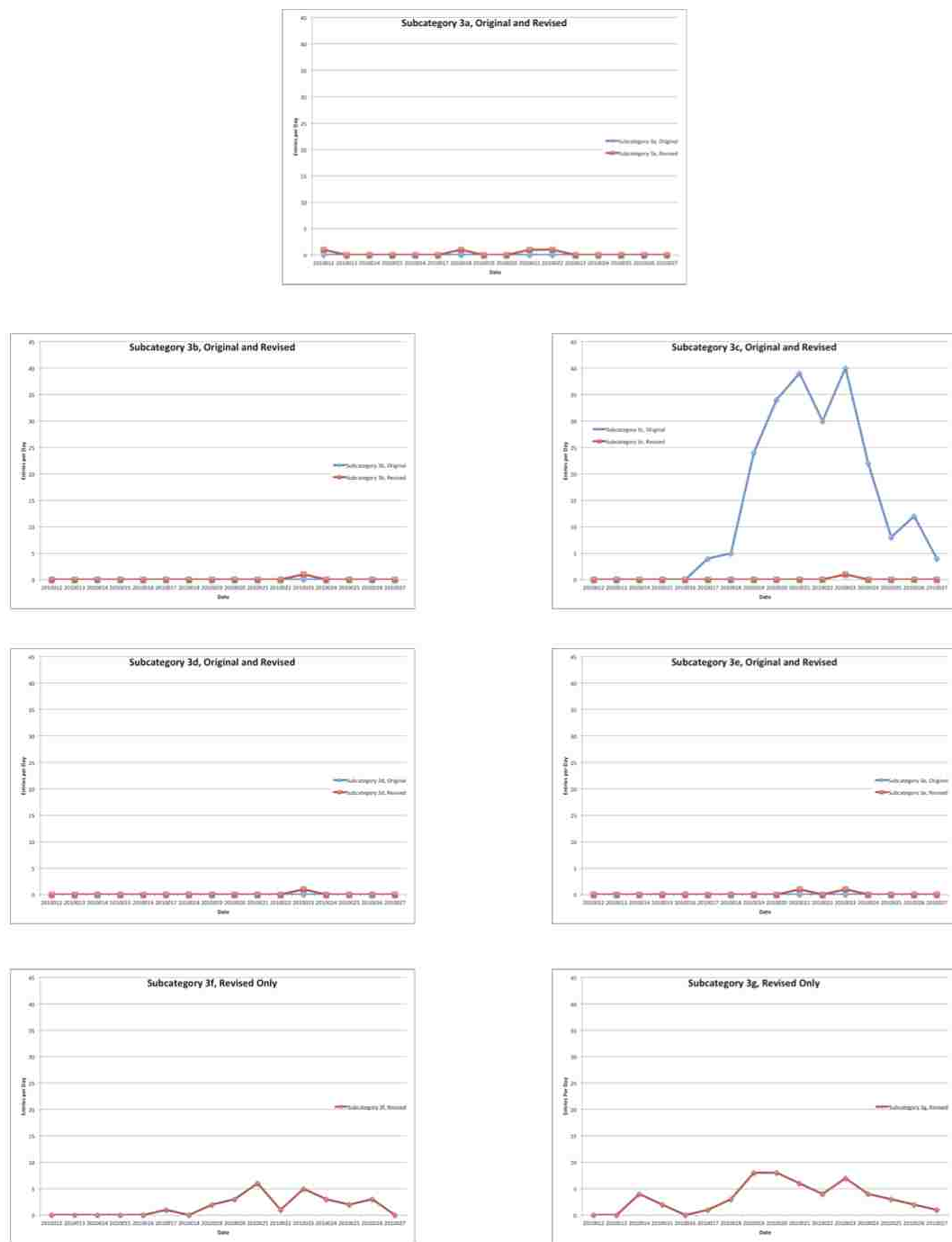


Figure 23 - Events for Subcategories 3a-3g by Day



The subcategories that comprise Category 4 highlight the differences in the number of entries among each of the subcategories as well as how the number of reports fluctuates per day. Subcategory 4e shows a primary discrepancy by commission. Subcategories 4b, 4c, and 4d consist of few events in both data sets. The number of events in Subcategory 4a generally matches between the two data sets. Refer to Figure 24.

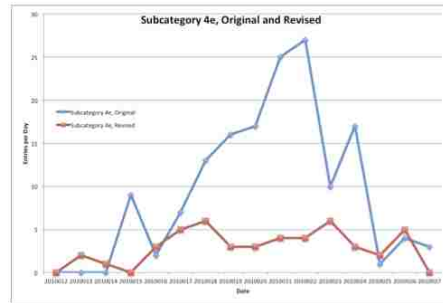
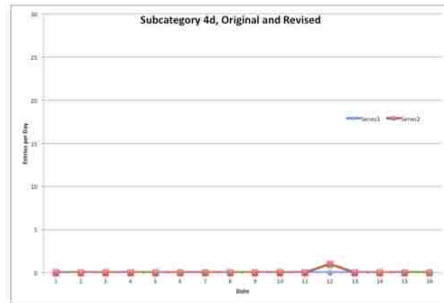
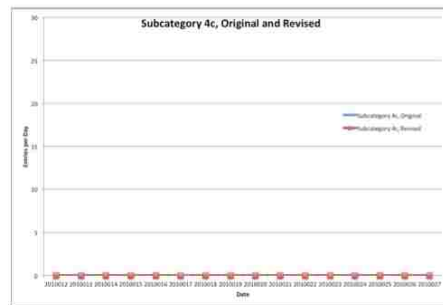
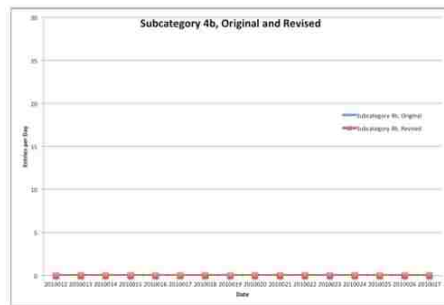
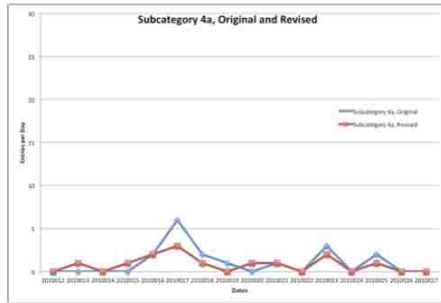


Figure 24 - Events for Subcategories 4a-4e by Day

The subcategories that comprise Category 5 highlight the differences in the number of entries among each of the subcategories as well as how the number of reports fluctuates per day. Subcategory 5a shows a primary discrepancy by omission. Subcategories 5d and 5e consist of few events in both data sets. The number of events in Subcategories 5b and 5c generally match between the two data sets. Refer to Figure 25.

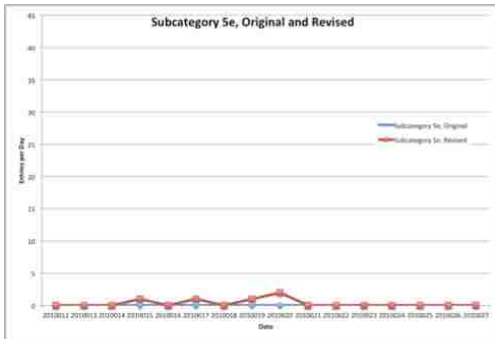
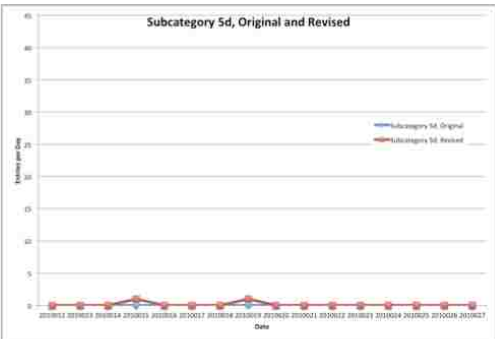
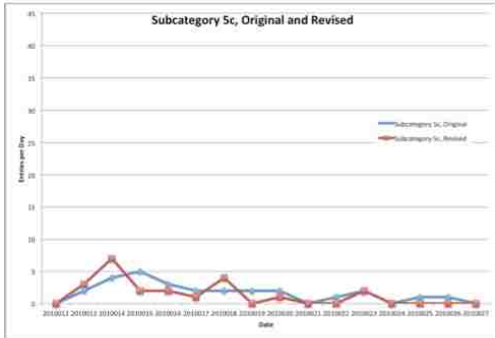
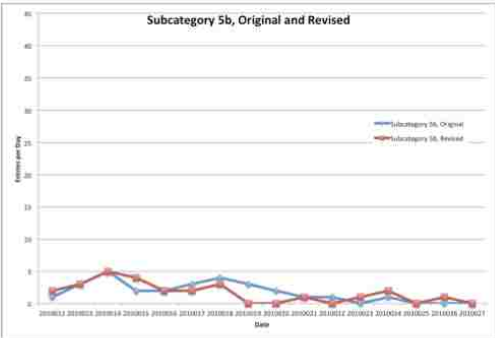
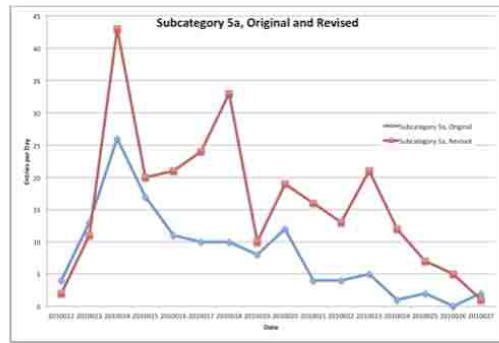


Figure 25 - Events for Subcategories 5a-5e by Day

The subcategories that comprise Category 6 are too few to draw any meaningful conclusions. While the values fluctuate per day and by subcategory, there were at most 5 events on any given day in either data set. Refer to Figure 26.

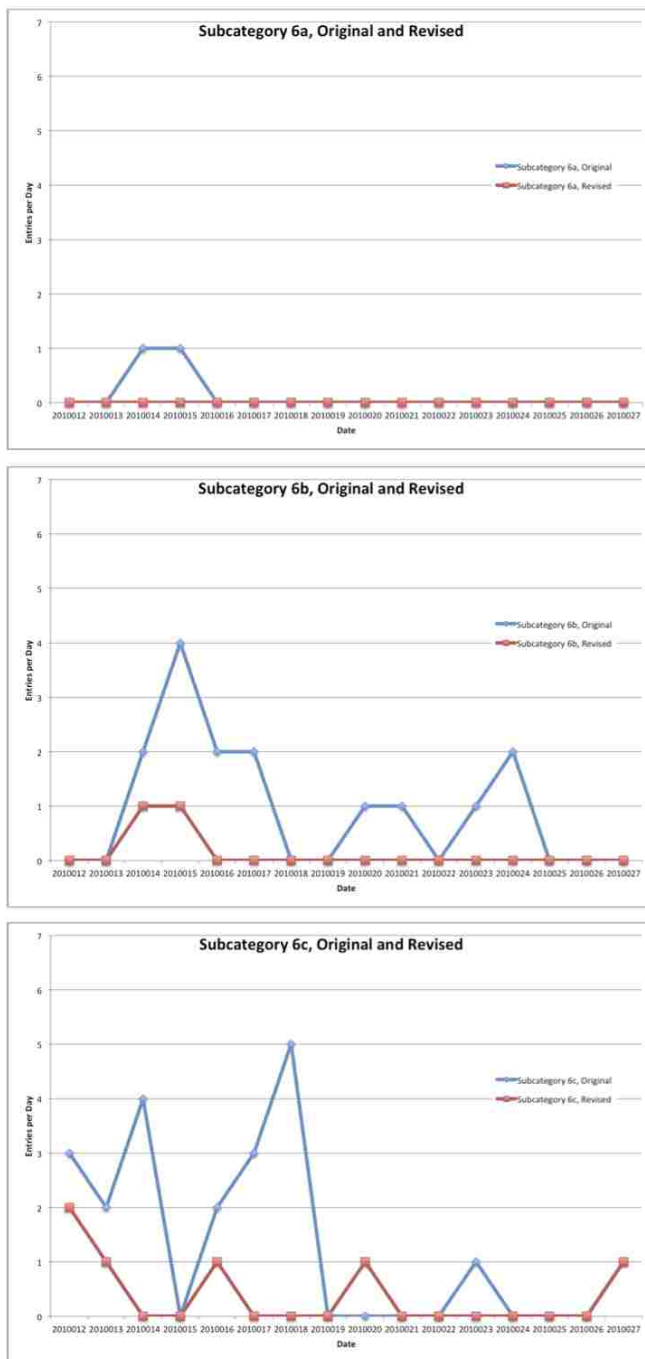


Figure 26 - Events for Subcategories 6a-6c by Day

The subcategories that comprise Category 7 highlight the differences in the number of entries among each of the subcategories as well as how the number of reports fluctuates per day. Subcategories 7a, 7d, and 7g show a primary discrepancy by commission. Subcategory 7b shows a slight discrepancy by omission while 7c exhibits both commission and omission depending on the day. The number of events in Subcategories 7e, 7f, and 7h generally match between the two data sets. Refer to Figure 27.

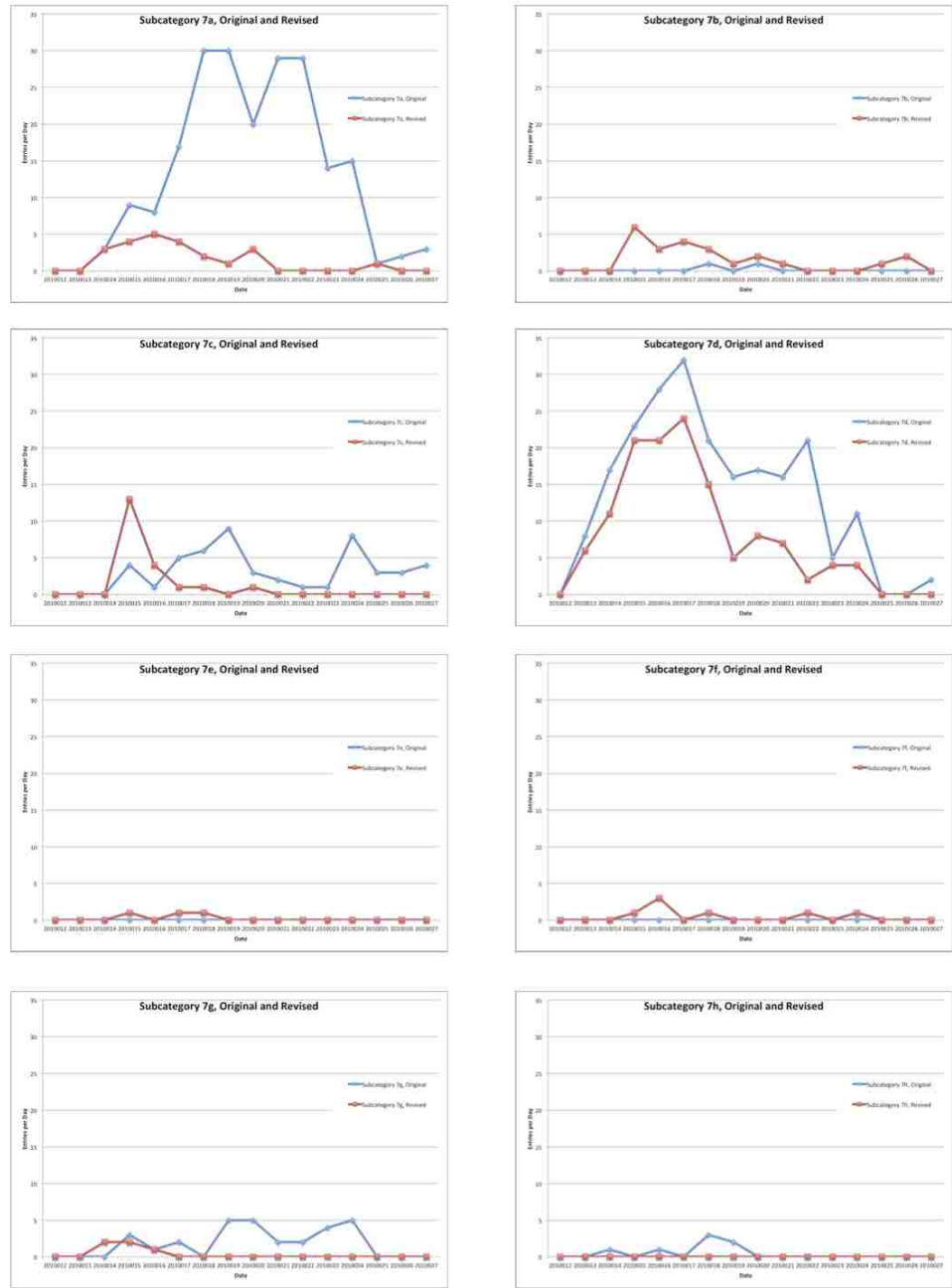


Figure 27 - Events for Subcategories 7a-7h by Day



The subcategories that comprise Category 8 highlight the differences in the number of entries among each of the subcategories as well as how the number of reports fluctuates per day. Subcategories 8a and 8f show a primary discrepancy by omission. Subcategories 8b and 8c consist of few events in both data sets. Subcategories 8d and 8e show a primary discrepancy by commission. Subcategories 8g and 8h did not exist in the original data set so there is nothing to compare them with. Refer to Figure 28.



Figure 28 - Events for Subcategories 8a-8h by Day

## Statistical Results

### Main Categories and Subcategories, Independent Ushahidi Review versus Researcher Findings

Table 17 was created by an independent review of the Ushahidi deployment in Haiti approximately a year and a half after the earthquake (Morrow *et al.*, 2011). The table was created by randomly selecting 50 entries from the total of 3,584 and assessing their original categories for errors of omission, commission, or a combination of omission and commission. Using these randomly selected entries, the reviewers estimated an overall error rate of 36% with a perfect match rate of 64%.

**Table 17 - Category Assessment by Discrepancy Type by Independent Ushahidi Evaluators**

Error Type	% of all Reports
Reports with incorrect category tag (Commission)	18
Reports missing a critical category tag (omission)	30
Both incorrect and missing tags	6
Missing or incorrect category tag (overall error rate)	36
Reports with neither missing nor incorrect tags	64

Table 18 contains a summary of the commission, omission, and combined commission and omission discrepancies for my study period at the main category level. This is the most forgiving level of assessment because it focuses on the main category level while ignoring specific subcategories. Note that the overall rate of discrepancy was determined to be 50.12%, which resulted in a perfect agreement rate of 49.88%.

**Table 18 - Category Assessment by Discrepancy Type, by Researcher at Main Category Level**

Error Type	% of all Reports
Reports with incorrect category tag (Commission)	12.65
Reports missing a critical category tag (omission)	14.50
Both incorrect and missing tags	22.97
Missing or incorrect category tag (overall error rate)	50.12
Reports with neither missing nor incorrect tags	49.88

Table 19 contains my results for my study period at the subcategory level that most closely matches the analysis done by the independent reviewers. This is a much more stringent level of assessment because it checks for agreement among all of the subcategories. Note that the overall rate of discrepancy was determined to be 73.41%, which resulted in a perfect agreement rate of 26.59%.

**Table 19 - Category Assessment by Discrepancy Type, by Researcher at Subcategory Level**

Error Type	% of all Reports
Reports with incorrect category tag (Commission)	16.83
Reports missing a critical category tag (omission)	23.90
Both incorrect and missing tags	16.83
Missing or incorrect category tag (overall error rate)	73.41
Reports with neither missing nor incorrect tags	26.59

## Chapter Five – Discussion

Based on the results from comparing the number of events for the categories and subcategories in the original data set produced by the volunteers with the re-categorized data that I produced, it would be easy to dismiss VGI as a viable data source during disasters. However, the simple percentages that I produced do not tell the entire story of this data set and its significance. Prior to the earthquake that struck Haiti in early 2010, almost all data produced during disasters came from authoritative sources. This data had a tendency to be slow to produce and was often not available to anyone outside the disaster response community. In addition, the collected data were often derived from remote sensing products or a limited number of data collectors on the ground. This greatly limited the amount of data produced and the types of data that could be collected. In contrast to this historical approach to data collection, the volunteer response to the Haiti disaster produced a major paradigm shift in emergency response.

In a way, the Ushahidi deployment in Haiti served as a national 9-1-1 system that was completely operated by volunteers. As a result of the Ushahidi deployment in Haiti, the amounts and types of data available to responders was greatly increased. Volunteers were able to collect data at a very fine scale from individual victims. Volunteers were also able to collect data that is not readily visible like hunger, thirst, and price gouging that may be missing from remotely sensed data. The data that were collected by volunteers was immediately available rather than having to wait to be incorporated into the next map update. The data were also available in an open format to anyone interested in the disaster. This helped reduce many of the barriers present in traditional emergency

response situations when dealing with incompatible software or data or when trying to navigate licensing agreements.

As was presented in the results section, however, the use of VGI data produced using the Ushahidi platform has room to improve. The following ideas are meant to stimulate discussion about ways to improve the quality of data produced as part of the Ushahidi platform but should not be considered recommendations. As this research project revolved around attribute consistency in the data, my ideas are focused on techniques that may improve that consistency.

First, organizers of the online crisis mapping community should consider developing documentation and training materials that can be provided to volunteers to help them more consistently categorize data produced by victims of disasters. This documentation might include categories and subcategories with definitions and examples. As part of the process of training volunteers, organizers of volunteers should consider hosting training exercises that can be used to assess the quality of the work conducted by the volunteers that could lead to improved training materials. Members of the online crisis mapping community should examine other hierarchical systems for classification, like the USGS Land Cover System, to determine if there are techniques or methods that can be incorporated into their own categorization schemes. Any changes that are introduced could be assessed during the training exercises discussed above. In order to assist researchers, future volunteers, and reviewers, volunteers should consider creating documentation for the procedures and tools that are used during an emergency. For instance, based on anecdotal evidence in the original CSV it appears that the categories and subcategories changed at least once during the first few weeks that the Haiti site was

online. When I asked the volunteer community if anyone knew what those original categories were or when they were changed, no one knew the answers to my questions (J. Valuch, personal communication, January 15, 2013).

I am also presenting some ideas that are specific to this scenario but may be applicable in future deployments. Volunteers should consider using a character encoding that is appropriate for the language of the victims. Numerous reports in the original CSV were difficult to read because of what appeared to be missing or inappropriate characters. Numerous reports described a need or activity at more than one location. Volunteers should determine if a system for splitting messages might be useful in the future if a message contains references to more than one location. There were several instances of roads being closed due to debris or landslides and then follow up reports that said those same features were now open. These different messages provided conflicting information. Volunteers should consider whether adding an ability to change the status of a report is needed. Many reports in the original CSV were meant as ways to collect data directly from victims. However, numerous reports were also published that related to information sharing between relief agencies. It may be worth investigating whether a multiple tiered system could be appropriate so that data can be collected from victims at one level and so that relief agencies can coordinate with each other at another level. Many messages from victims voiced a great deal of frustration at having reported a need but not receiving help. In the future it may be beneficial to help manage the expectations of the victims so that they do not feel taken advantage of or neglected when they are already in a vulnerable situation. While some messages were quite specific about the types of aid that were needed and where the victims were located, many other messages

were quite generic with people just asking for help. It may also be worthwhile to investigate ways to better communicate what information is needed from the victims in order to act on their reports. There were many reports that matched a main category but did not suite any of the appropriate subcategories. It might be worth considering adding an “other” subcategory to each main category rather than having a single main category labeled “other”. For instance, there were several instances where people requested help with corpse removal. This type of request may best fit in the main category related to public health, but because there was not a specific subcategory for corpse removal, I placed these reports in the other main category.



## Chapter Six – Conclusions

### Summary of Findings

This research project set out to better understand the consistency of categorization by volunteers in a time critical emergency. An independent review of the Ushahidi Haiti deployment estimated that 64 percent of the reports in the original database were correctly categorized (Morrow *et al.*, 2011). After re-categorizing the data, I found that during my study period the data was correctly categorized nearly 27 percent of the time at the subcategory level. My estimates are less than half of what the independent reviewers estimated. The process of comparing the number and distribution of subcategories between the two data sets is a very strict comparison between the two data sets. For instance, a volunteer may have categorized a message as 5a. Collapsed Structure and I may have categorized the same message as 5b. Unstable Structure. While both subcategories imply that a structure may be structurally unsound, they do not exactly match so in my assessment these would be inconsistent. In order to better understand if the main ideas of the messages were consistently identified, I aggregated the subcategories up to their appropriate main categories and compared each record across the two data sets. Using the example scenario from above, in this case the 5a and 5b would both be aggregated as main category 5 and would therefore be a match. I found that when comparing the two data sets at the main category level almost 50 percent of the messages were consistently categorized. While this is a marked improvement over 27 percent at the subcategory level, it is still considerably less than the 64 percent estimated by the independent reviewers.

## **Limitations**

There are several limitations related to the methods and results of this research project that need to be addressed. While every effort was made to be as consistent as possible during the re-categorization process, only one person (me) was responsible for this step of the project. Therefore the results of the re-categorization process may be biased based on my own background and understanding of the original work of the volunteers. It is also possible that as I was manually re-categorizing each entry in the database that I mistyped the categories or subcategories as I was entering them in the computer. In addition, when I attempted to obtain the training materials that were provided to the original volunteers to help them determine how to categorize the contents of the data that was submitted to the Ushahidi platform, I discovered that there were not any materials. As a result, I had to develop my own definitions for each category. Some categories and subcategories were self-explanatory while others were ambiguous or seemed redundant. If a different person were responsible for writing their own category and subcategory definitions, then they may develop very different results.

Per more traditional methods of comparing attributes, it is better to compare a dataset of unknown quality to one of known quality (Goodchild and Hunter, 1997), however, no data set exists to my knowledge that is suitable for comparison to the original Ushahidi data set. Due to the reasons listed above, the attributes of the re-categorized data set that I created are but one example and should not be considered “THE” ultimate categorization.

Though there are as many as eight components of geospatial data quality (van Oort, 2006), this research project only addresses one of those components. I discovered the difficulties of comparing specific categories and subcategories as I wrestled with how

to categorize the data. For instance, Category 6 was described as “Natural Hazards”. But as I reviewed the original data I discovered two separate definitions for subcategories 6a, 6b, and 6c. One set of original subcategories clearly related to the main category heading (Floods, Landslides, Earthquakes and Aftershocks), while the second subcategories seemed to belong in different main categories (Deaths, Missing Persons, Asking to Forward a Message). Because the second set of subcategories did not match the original main category, I separated those and moved them to what I felt was a more appropriate category. How best, then, to compare the number of entries in the original data with the number of entries in the revised data? My method was to compare categories and subcategories with matching number and letter designations no matter the definitions. The reason for this is that I wanted to create the “best” version of the original data that I could where categories and their subcategories were consistent and related. While this may not be the most appropriate method, in most cases the number of entries in the affected subcategories was quite small and statistically insignificant.

There are also potential limitations that affect the findings of my research. While this research project addressed the quantity of inconsistencies in the categorization process, it does not investigate the cause of those discrepancies. In addition, this research project only investigates a single Ushahidi deployment. The results of Ushahidi deployments that took place after this disaster may have significantly different results. This research project does not address ways in which to better incorporate VGI with more traditional relief agencies. In addition, this research project does not address the effects of error in the attribute data and whether or not they have any impact on relief activities.

## **Future Research**

Some future research should be conducted as a way to investigate the limitations of this research project as discussed above. For instance, future research could look for and examine any trends in inconsistencies between the original data and the revised data over time, and as the number of events increases or decreases. While it is beneficial to know the magnitude of the discrepancies between the two data sets, it would also be useful to study the nature or cause of the omission and commission discrepancies.

Attempts should be made to study as many different components of geospatial data quality as possible related to this disaster. Longitudinal studies could also be conducted to determine if the discrepancies are consistent across multiple deployments of the Ushahidi platform.

Further research could also provide direct benefits to individuals or organizations that utilize the Ushahidi platform. For instance, what are some ways to incorporate the strengths of VGI in disaster response to supplement data produced by traditional aid agencies? GIS models could be developed to help predict the accuracy of information collected by volunteers to help responders assess the appropriateness of a data source. Sensitivity analyses could also be conducted to determine which aspects of geospatial quality have the greatest impact on activities undertaken by emergency responders.

Beyond the specific research ideas mentioned above that relate specifically to this thesis, I also feel that there are several broad categories that deserve further attention. An important component of future VGI and disaster response research should focus on the equality of access to the tools needed to create and share the data. Haklay and Ellul (2010) present evidence to suggest that the producers of VGI do not span all socioeconomic segments of society, but rather skew away from people who are at the low

end of the socioeconomic spectrum. This could prove dangerous because these same people may be most at risk of being unable to evacuate prior to an emergency and may not have access to the resources necessary to survive and rebuild as quickly as people from higher up on the socioeconomic ladder. Future research may need to focus on making the tools necessary to create and share VGI more equally available.

Another important component of future VGI research is the degree that data creation is affected by the disaster itself. Nelson, Sigal, and Zambrano (2010) report that, despite the unprecedented role of crowdsourcing and VGI in response to the Haiti earthquake, the use of cell phones was reliant on damaged and overtaxed cellular networks and access to electricity to charge and use the devices, while inexpensive, battery powered low-tech FM radio was available throughout the duration of the earthquake. It may prove that, despite the allure of new technology like location aware cell phones and web applications, a combination of low and high tech efforts may be the best way to collect and provide information following a devastating disaster that disrupts critical infrastructure. Future research could focus on combining technology that is more resistant to damage from disasters with newer technology that may or may not be available due to damage.

Future research could also focus on how people actually interact with social media tools, especially in high stress situations (American Red Cross, 2010). Will people be more likely to use social media tools during an emergency if they use them in their normal life? Will the data they provide be more useful if they are already familiar with how the tools work rather than trying to learn them in a high-stress situation? The same research could be applied to those institutions and agencies that may want to incorporate

VGI into their decision support system. If these agencies are able to test and experiment with VGI during normal operations, will they be more likely to use VGI during disasters? Also, will they make better use of the data if they already have a system in place to take advantage of it rather than waiting until an emergency to learn?

A combination of law, politics, and geography may also present a new research agenda related to VGI. Are relief agencies required to act on data that are submitted through social media like Twitter? What if the data are submitted anonymously? Will they be held liable if they fail to act even if responding uses limited resources? Will society provide the funding and resources to provide relief agencies with the new tools and training to take advantage of social media and VGI?

As the cost of entry continues to lower, the role of VGI in society will continue to grow. More and more people have access to cell phones and the Internet. The capabilities of those cell phones continue to expand, and even when they do not provide smart phone capabilities, they can be used to provide useful information (Munro, 2010). The adoption of social media applications like Facebook and Twitter show no signs of slowing, and as they continue to gain new users and features, their role in society will continue to grow and provide researchers with ever more data to mine. Relief agencies like the United Nations have begun to recognize the importance of VGI in their operations and are conducting their own research into VGI's usefulness and applications (Standby Volunteer Task Force, 2011). As it becomes easier and more common to produce and share georeferenced information, geography may see an increase in attention as new researchers and existing disciplines seek to combine the potential of VGI with their own

discipline. If this increased attention to geography proves true, then countless new variations of the research proposed in this study will be added to our existing literature.

## Appendices

### Appendix A: Ushahidi Haiti Category Rules and Definitions

**Category 1** – Emergency = a time critical response is necessary in order to preserve life and/or property

*1a* – Highly Vulnerable = reports of victims who are especially vulnerable like children or the disabled but not necessarily in formalized settings like an orphanage or nursing home

*1b* – Medical emergency = reports of immediate, life threatening illnesses or injuries that involve heavy bleeding, head trauma, etc

*1c* – People trapped = reports of people trapped who are still alive, must be specific about location

*1d* – Fire = reports of fire

**Category 2** – Vital Lines = requests for services or goods that are necessary to sustain life, if a generic request that does not mention a specific subcategory below use 2

*2a* – Food shortage = reports requesting food or mentioning “hungry”, “starving” or other words that indicate hunger

*2b* – Water shortage = reports requesting water or mentioning “thirst”, “dehydration”, or other words that indicate thirst

*2c* – Contaminated water = reports indicating that a water supply is not safe to drink

*2d* – Shelter needed = reports requesting tents, sleeping bags, tarps, clothing, or any other words that indicate cold, wet, damp, sleeping on street, etc

*2e* – Fuel shortage = reports requesting fuel for generators, vehicles, or any other liquid fossil fuel purpose

*2f* – Power outage = reports indicating a lack of electricity

**Category 3** – Public Health = requests for services or goods that are necessary to prevent or treat illness or injury, also generic medical care – requests for medical care that are not life threatening, for instance, when someone says they are sick, ill, need a doctor; also requests for corpse removal

*3a* – Infectious human disease = reports indicating a specific illness not simply feeling ill or being sick

*3b* – Chronic care needs = reports indicating assistance needed for victims with life-long health needs

*3c* – Medical equipment and supply needs = requests for any items used for treating illnesses or injuries, including medically trained staff (doctors, nurses, specialists, etc)

*3d* – OBGYN/Women’s health = requests for medical care that are specific to women’s health or delivery

*3e* – Psychiatric need = requests for help with psychiatric problems

*3f* – Water sanitation and hygiene promotion (move to category 3 public health) = reports requesting items necessary for personal sanitation or hygiene like soap, water purification methods, toilets/latrines, etc (formerly 4e)



**3g – Deaths (move to category 3 public health) = reports that contain specific mentions of death at a specific location (should be a building or address) (former 6a duplicate)**

**Category 4 – Security Threats = requests for security or reports of threats**

*4a – Looting = reports indicating looting or theft, but not of aid*

*4b – Theft of aid = reports indicating theft of aid*

*4c – Group Violence = reports indicating riots or large groups intent on harm*

*4d – Riot = reports indicating riots or large groups*

*4e – Security Concern (former 2c duplicate)*

**Category 5 – Infrastructure Damage = reports of damage to infrastructure**

*5a – Collapsed structure = reports indicating collapsed buildings that are specific with a building/company name or address*

*5b – Unstable structure = reports indicating collapsed buildings that are specific with a building/company name or address*

*5c – Road blocked = reports indicating that a road is blocked by natural or manmade means*

*5d – Compromised bridge = reports indicating that a bridge has sustained damage or has collapsed*

*5e – Communication Lines down = reports indicating that any form of communication has stopped functioning including radio, television, cellular telephone, landline, etc.*

**Category 6 – Natural Hazards = reports of natural hazards**

*6a – Floods = reports of flooding*

*6b – Landslides = reports of landslides*

*6c – Earthquake and aftershocks = specific reports of earthquake or aftershocks*

**Category 7 – Services Available = services that are available to victims or services or supplies that are available to other aid agencies for the ultimate purpose of helping victims**

*7a – Food distribution point = reports indicating that food is being distributed from this location*

*7b – Water distribution point = reports that water is being distributed from this location*

*7c – Non-food distribution point = reports that non-food/water items are being distributed from this location*

*7d – Hospital/clinics operating = reports that medical services are being provided at this location*

*7e – Feeding centers available = reports that meals (not food) are being provided at this location*

*7f – Shelter offered = reports that shelter is being offered at this location*

*7g – Human remains management = reports that human remains are being managed at this location*

*7h – Rubble removal = reports that rubble removal is taking place at this location*

**Category 8** – Other = any reports or requests that do not fit the above criteria

*8a* – IDP Concentration = reports that people are concentrating at this location, must be more specific than a city or region, greater than 20 people

*8b* – Aid manipulation = reports that aid is not being distributed fairly or is being manipulated in any other way

*8c* – Price gouging = reports of price gouging

*8d* – Search and rescue = specific requests for search and rescue

*8e* – Person news = reports that people are safe if they do not already fit another category

*8f* – Other = requests for transportation, requests for tools to remove rubble without mentioning trapped people, requests for money, requests for jobs, reports of infrastructure repair, ambiguous message

*8g* – Missing Persons (move to category 8 other) = reports indicating that someone is missing, can't be found, is not answering their phone, or has not been heard from (former 6b duplicate)

*8h* – Asking to forward a message (move to category 8 other) = reports asking for someone to pass on a message to someone else, but not related to requests for aid/help (former 6c duplicate)

## Appendix B: Main Category Python Script

```

## this script was generated in order to work with data created for my thesis
## the purpose of this script is to open a csv file, create two lists for each
row in the csv, then to compare each of those lists
## author: michael camponovo
## email: mecampo@unm.edu
## date: 20130207

## import required modules
import csv
import os

## define original csv
original_csv = ##replace this comment with the file path and file name, with
extension to your original csv

## create empty list to store csv
master_csv_list = []
output_list = []

##open original csv
with open(original_csv, 'rU') as original_csv_file:
    original_csv_reader = csv.reader(original_csv_file, delimiter=',')
    for row in original_csv_reader:
        master_csv_list.append(row)
    ##print master_csv_list[0] ##this currently print the first line about 2000
times
for row in master_csv_list[1:]:
    original_list = [] #create empty list for original category values
    revised_list = [] #create empty list for revised category values
    if row[1] == '1':# populate each list with the values 1-8 based on
whether that value is present in the original csv
        original_list.append('1')
    if row[2] == '1':
        original_list.append('2')
    if row[3] == '1':
        original_list.append('3')
    if row[4] == '1':
        original_list.append('4')
    if row[5] == '1':
        original_list.append('5')
    if row[6] == '1':
        original_list.append('6')
    if row[7] == '1':
        original_list.append('7')
    if row[8] == '1':
        original_list.append('8')
    if row[9] == '1':
        revised_list.append('1')
    if row[10] == '1':
        revised_list.append('2')
    if row[11] == '1':
        revised_list.append('3')

```

```

    if row[12] == '1':
        revised_list.append('4')
    if row[13] == '1':
        revised_list.append('5')
    if row[14] == '1':
        revised_list.append('6')
    if row[15] == '1':
        revised_list.append('7')
    if row[16] == '1':
        revised_list.append('8')
    both = list(set(original_list) & set(revised_list))# create a list
of values that are in both lists
    original_only = list(set(original_list) - set(revised_list))#
create a list of values that are only in the original list
    revised_only = list(set(revised_list) - set(original_list))#
create a list of values that are only in the revised list
    # '|'.join(original_list)
# '|'.join(revised_list)
# '|'.join(both)
# '|'.join(original_only)
# '|'.join(revised_only)
    my_string = ','.join([row[0], '|'.join(original_list),
'|'.join(revised_list), '|'.join(both), '|'.join(original_only),
'|'.join(revised_only)])
    output_list.append(my_string)
with open('new.csv', 'w') as new_csv:
    new_csv.write('\n'.join(output_list))

```

## Appendix C: Subcategory Python Script

```

## this script was generated in order to work with data created for my thesis
## the purpose of this script is to open a csv file, create two lists for each
row in the csv, then to compare each of those lists
## author: michael camponovo
## email: mecampo@unm.edu
## date: 20130207

## import required modules
import csv
import os

## define original csv
original_csv = ##replace this comment with the file path and file name, with
extension to your original csv

## create empty list to store csv
master_csv_list = []
output_list = []

##open original csv
with open(original_csv, 'rU') as original_csv_file:
    original_csv_reader = csv.reader(original_csv_file, delimiter=',')
    for row in original_csv_reader:
        master_csv_list.append(row)
    ##print master_csv_list[0] ##this currently print the first line about 2000
times
for row in master_csv_list[1:]:
    original_list = [] #create empty list for original category values
    revised_list = [] #create empty list for revised category values
    if row[1] == '1':# populate each list with the values 1a-8h based
on whether that value is present in the original csv
        original_list.append('1a')
    if row[2] == '1':
        original_list.append('1b')
    if row[3] == '1':
        original_list.append('1c')
    if row[4] == '1':
        original_list.append('1d')
    if row[5] == '1':
        original_list.append('2a')
    if row[6] == '1':
        original_list.append('2b')
    if row[7] == '1':
        original_list.append('2c')
    if row[8] == '1':
        original_list.append('2d')
    if row[9] == '1':
        original_list.append('2e')
    if row[10] == '1':
        original_list.append('2f')
    if row[11] == '1':
        original_list.append('3a')

```

```
if row[12] == '1':
    original_list.append('3b')
if row[13] == '1':
    original_list.append('3c')
if row[14] == '1':
    original_list.append('3d')
if row[15] == '1':
    original_list.append('3e')
if row[16] == '1':
    original_list.append('4a')
if row[17] == '1':
    original_list.append('4b')
if row[18] == '1':
    original_list.append('4c')
if row[19] == '1':
    original_list.append('4d')
if row[20] == '1':
    original_list.append('4e')
if row[21] == '1':
    original_list.append('5a')
if row[22] == '1':
    original_list.append('5b')
if row[23] == '1':
    original_list.append('5c')
if row[24] == '1':
    original_list.append('5d')
if row[25] == '1':
    original_list.append('5e')
if row[26] == '1':
    original_list.append('6a')
if row[27] == '1':
    original_list.append('6b')
if row[28] == '1':
    original_list.append('6c')
if row[29] == '1':
    original_list.append('7a')
if row[30] == '1':
    original_list.append('7b')
if row[31] == '1':
    original_list.append('7c')
if row[32] == '1':
    original_list.append('7d')
if row[33] == '1':
    original_list.append('7e')
if row[34] == '1':
    original_list.append('7f')
if row[35] == '1':
    original_list.append('7g')
if row[36] == '1':
    original_list.append('7h')
if row[37] == '1':
    original_list.append('8a')
if row[38] == '1':
```

```
        original_list.append('8b')
if row[39] == '1':
    original_list.append('8c')
if row[40] == '1':
    original_list.append('8d')
if row[41] == '1':
    original_list.append('8e')
if row[42] == '1':
    original_list.append('8f')
if row[43] == '1':
    revised_list.append('1a')
if row[44] == '1':
    revised_list.append('1b')
if row[45] == '1':
    revised_list.append('1c')
if row[46] == '1':
    revised_list.append('1d')
if row[47] == '1':
    revised_list.append('2a')
if row[48] == '1':
    revised_list.append('2b')
if row[49] == '1':
    revised_list.append('2c')
if row[50] == '1':
    revised_list.append('2d')
if row[51] == '1':
    revised_list.append('2e')
if row[52] == '1':
    revised_list.append('2f')
if row[53] == '1':
    revised_list.append('3a')
if row[54] == '1':
    revised_list.append('3b')
if row[55] == '1':
    revised_list.append('3c')
if row[56] == '1':
    revised_list.append('3d')
if row[57] == '1':
    revised_list.append('3e')
if row[58] == '1':
    revised_list.append('3f')
if row[59] == '1':
    revised_list.append('3g')
if row[60] == '1':
    revised_list.append('4a')
if row[61] == '1':
    revised_list.append('4b')
if row[62] == '1':
    revised_list.append('4c')
if row[63] == '1':
    revised_list.append('4d')
if row[64] == '1':
    revised_list.append('4e')
```

```

if row[65] == '1':
    revised_list.append('5a')
if row[66] == '1':
    revised_list.append('5b')
if row[67] == '1':
    revised_list.append('5c')
if row[68] == '1':
    revised_list.append('5d')
if row[69] == '1':
    revised_list.append('5e')
if row[70] == '1':
    revised_list.append('6a')
if row[71] == '1':
    revised_list.append('6b')
if row[72] == '1':
    revised_list.append('6c')
if row[73] == '1':
    revised_list.append('7a')
if row[74] == '1':
    revised_list.append('7b')
if row[75] == '1':
    revised_list.append('7c')
if row[76] == '1':
    revised_list.append('7d')
if row[77] == '1':
    revised_list.append('7e')
if row[78] == '1':
    revised_list.append('7f')
if row[79] == '1':
    revised_list.append('7g')
if row[80] == '1':
    revised_list.append('7h')
if row[81] == '1':
    revised_list.append('8a')
if row[82] == '1':
    revised_list.append('8b')
if row[83] == '1':
    revised_list.append('8c')
if row[84] == '1':
    revised_list.append('8d')
if row[85] == '1':
    revised_list.append('8e')
if row[86] == '1':
    revised_list.append('8f')
if row[87] == '1':
    revised_list.append('8g')
if row[88] == '1':
    revised_list.append('8h')
both = list(set(original_list) & set(revised_list))# create a list
of values that are in both lists
original_only = list(set(original_list) - set(revised_list))#
create a list of values that are only in the original list

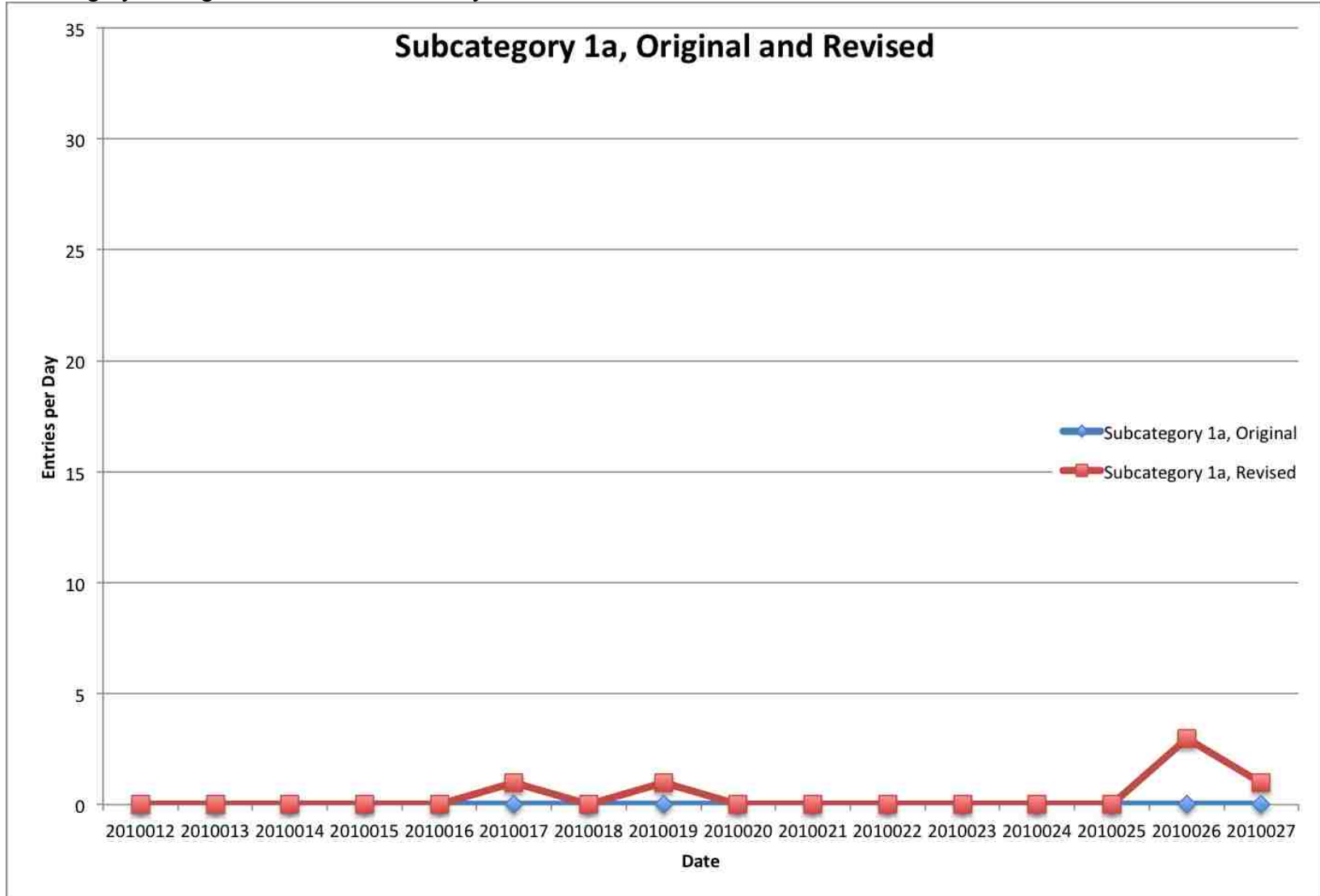
```



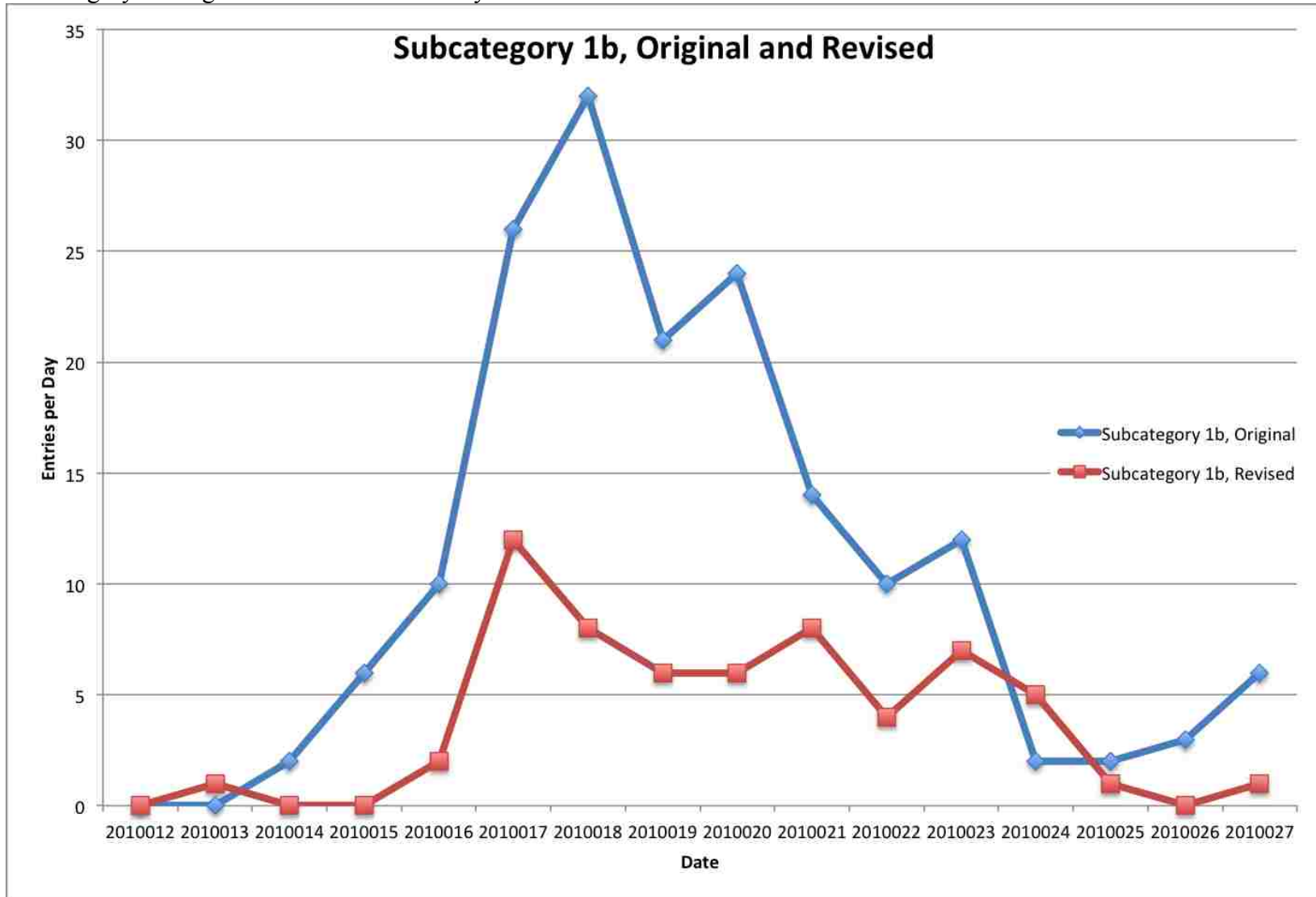
```
        revised_only = list(set(revised_list) - set(original_list))#
create a list of values that are only in the revised list
        # '|'.join(original_list)
#        '|'.join(revised_list)
#        '|'.join(both)
#        '|'.join(original_only)
#        '|'.join(revised_only)
        my_string = ','.join([row[0], '|'.join(original_list),
'|'.join(revised_list), '|'.join(both), '|'.join(original_only),
'|'.join(revised_only)])
        output_list.append(my_string)
with open('newsub.csv', 'w') as new_subcsv:
        new_subcsv.write('\n'.join(output_list))
```

**Appendix D: Subcategory Entries by Date**

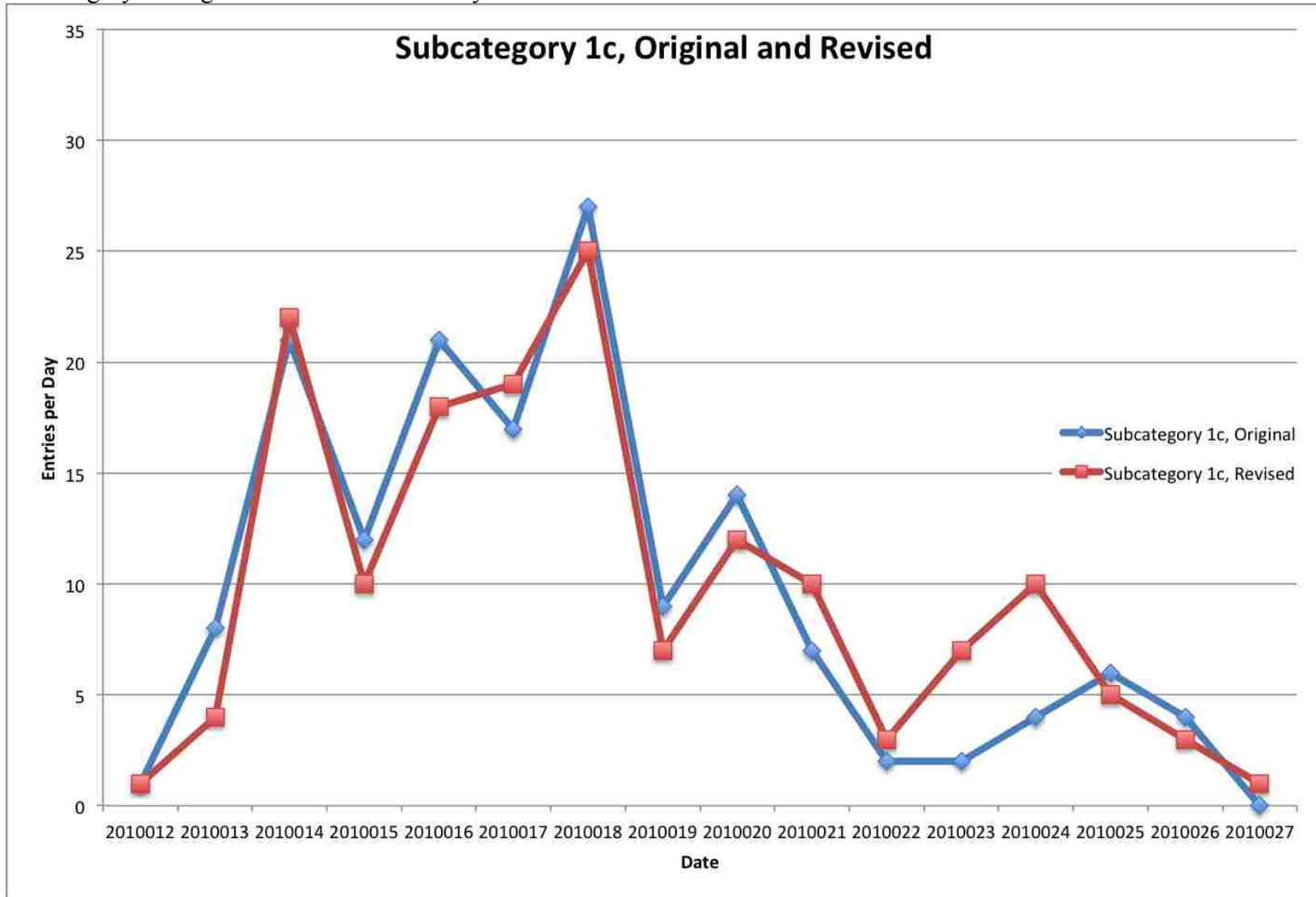
Subcategory 1a original and revised entries by date



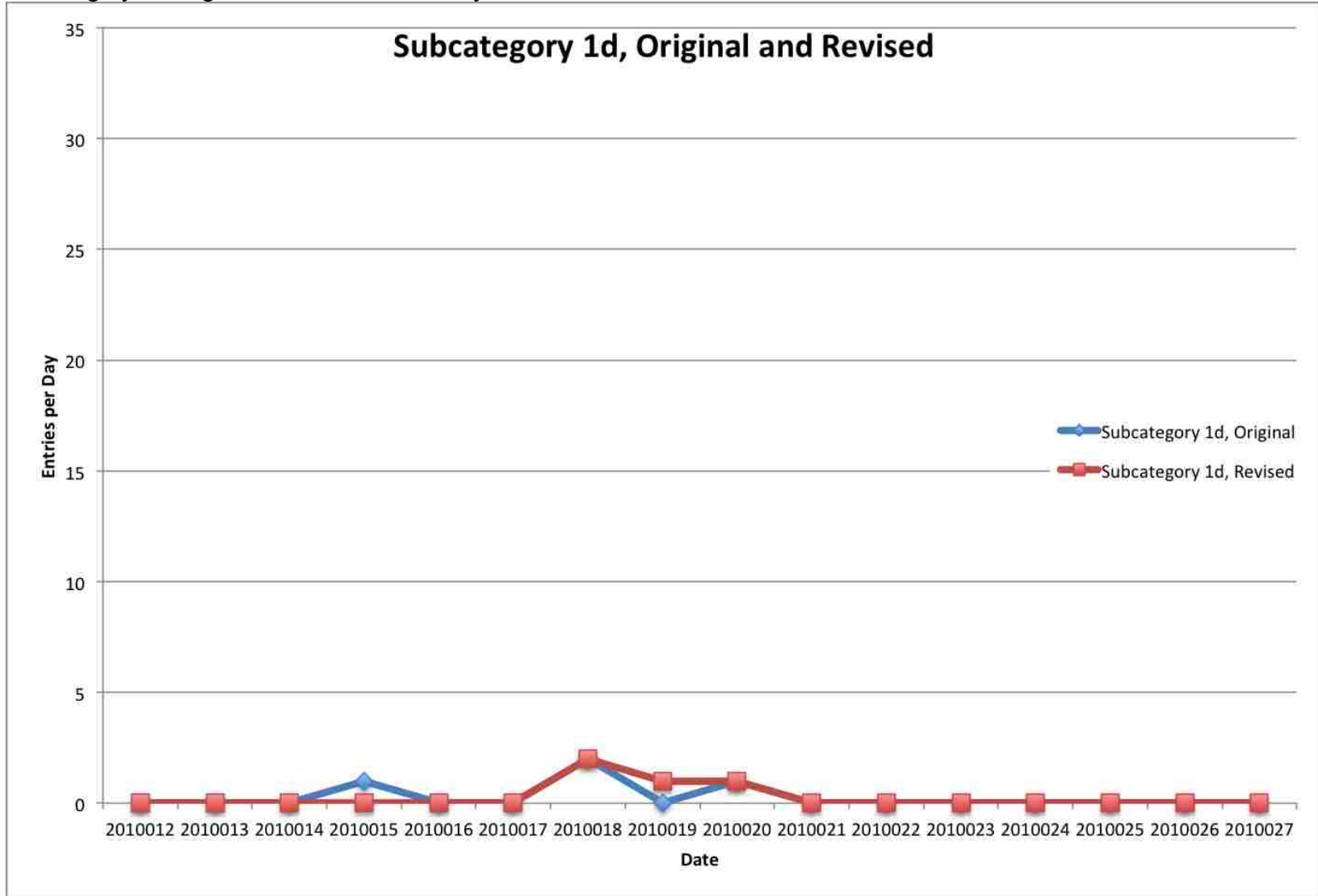
Subcategory 1b original and revised entries by date



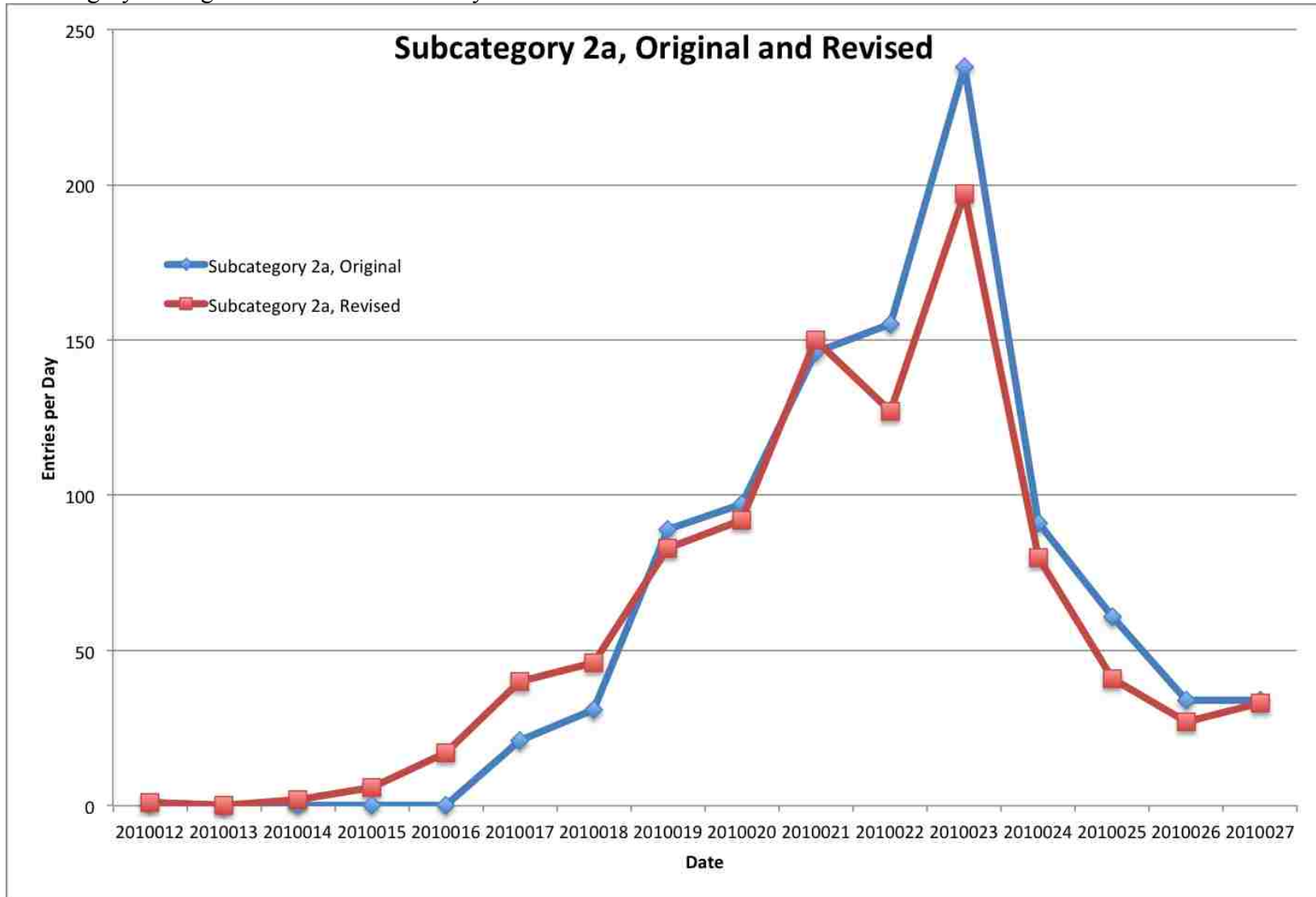
Subcategory 1c original and revised entries by date



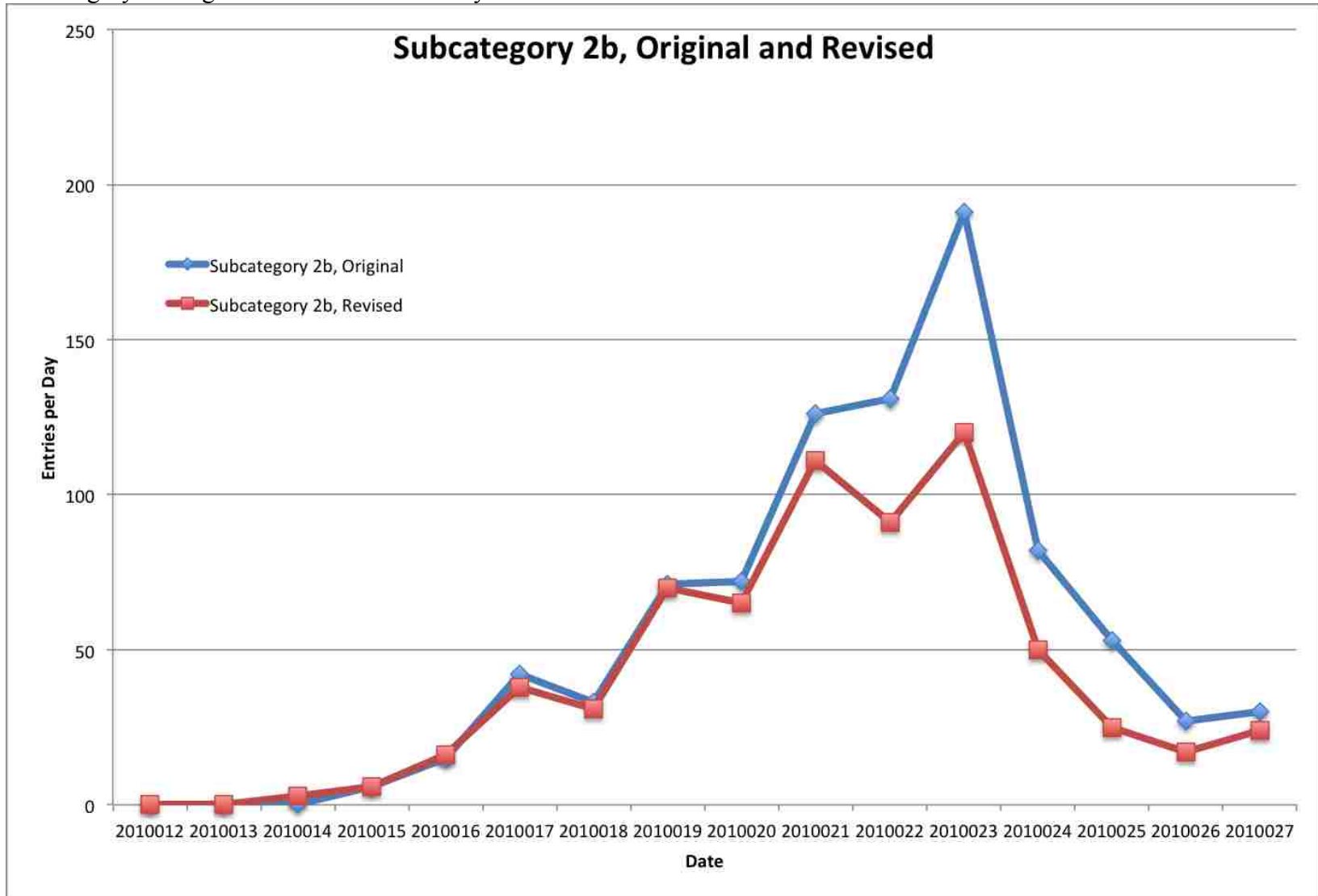
Subcategory 1d original and revised entries by date



Subcategory 2a original and revised entries by date

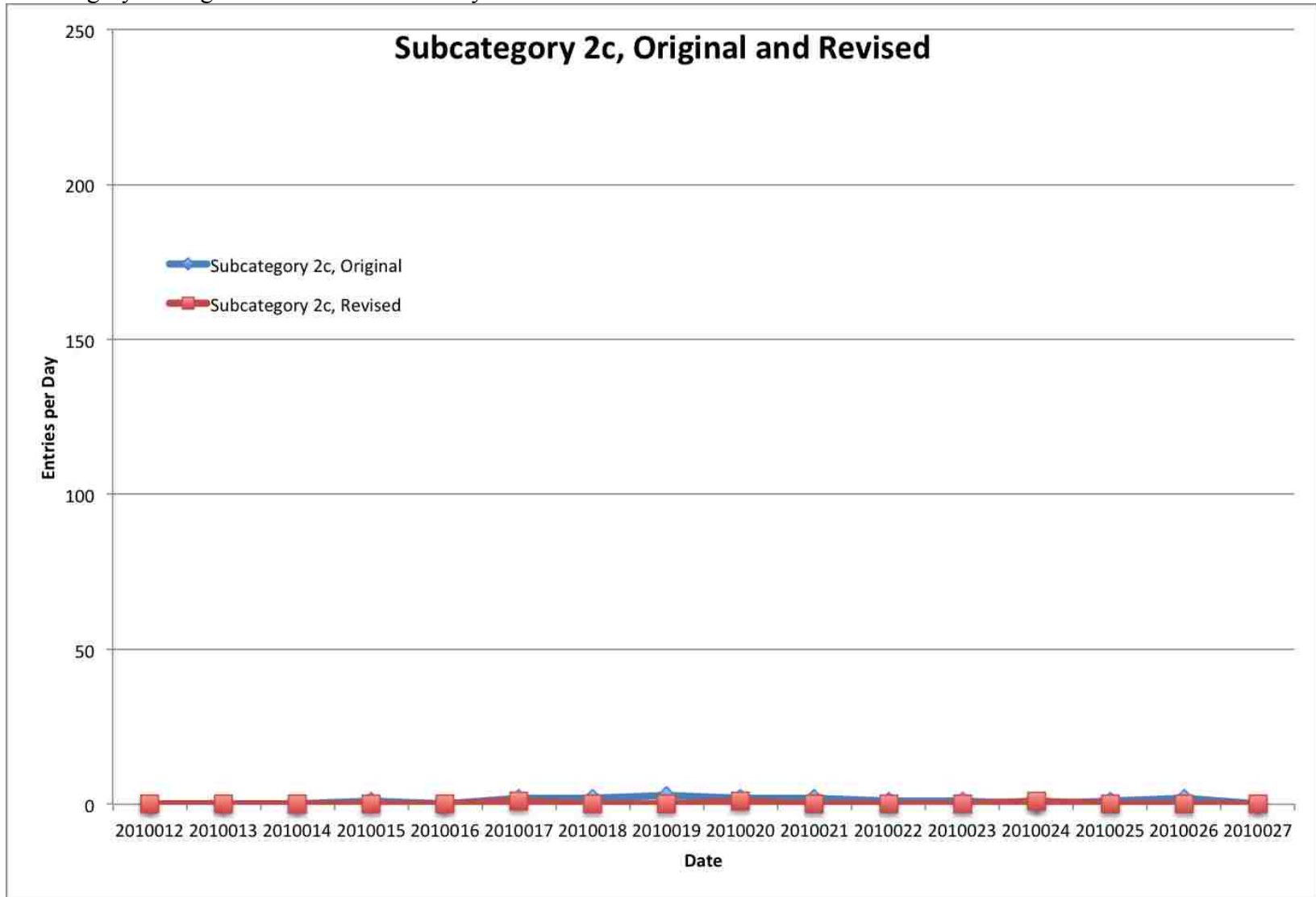


Subcategory 2b original and revised entries by date

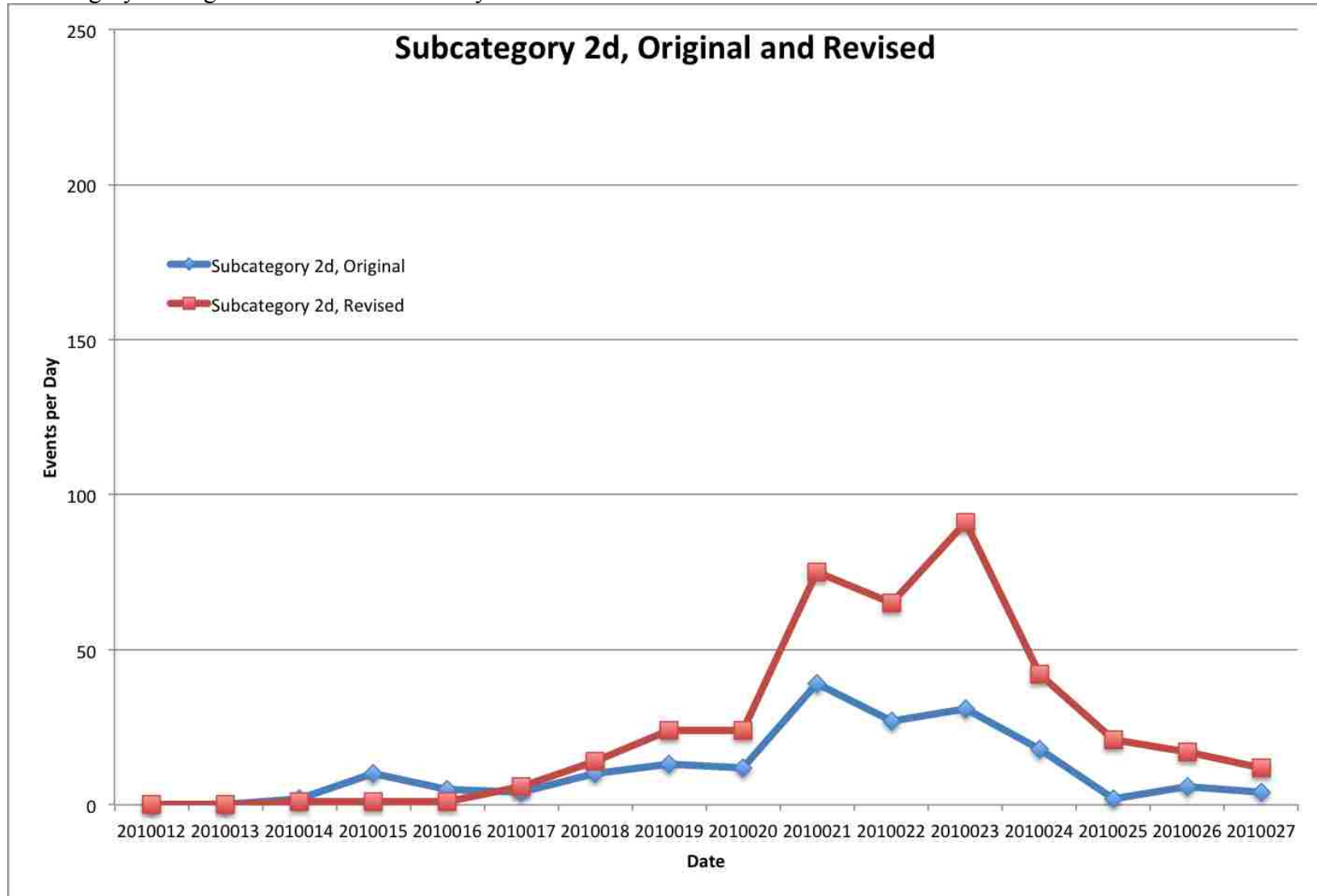




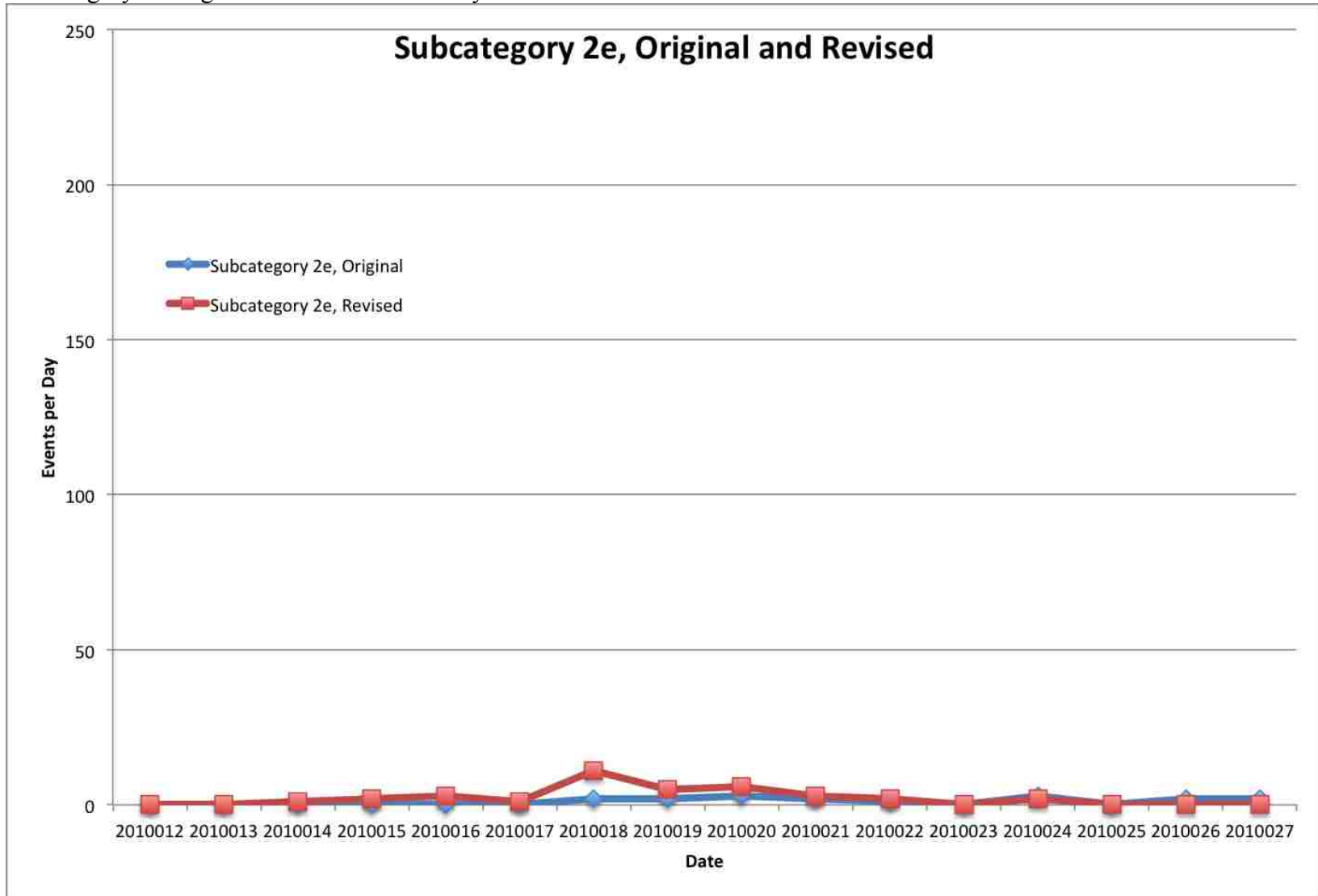
Subcategory 2c original and revised entries by date



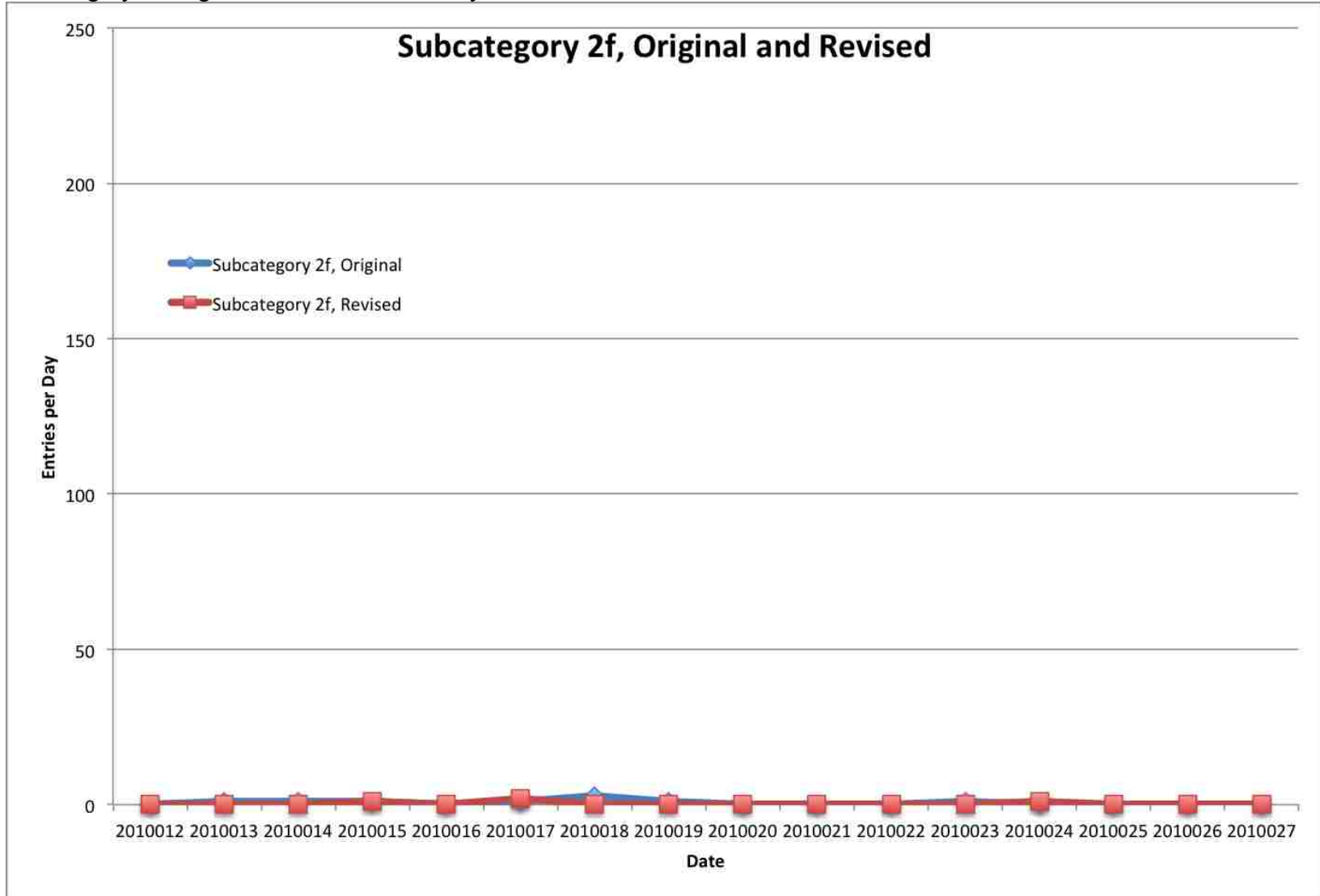
Subcategory 2d original and revised entries by date



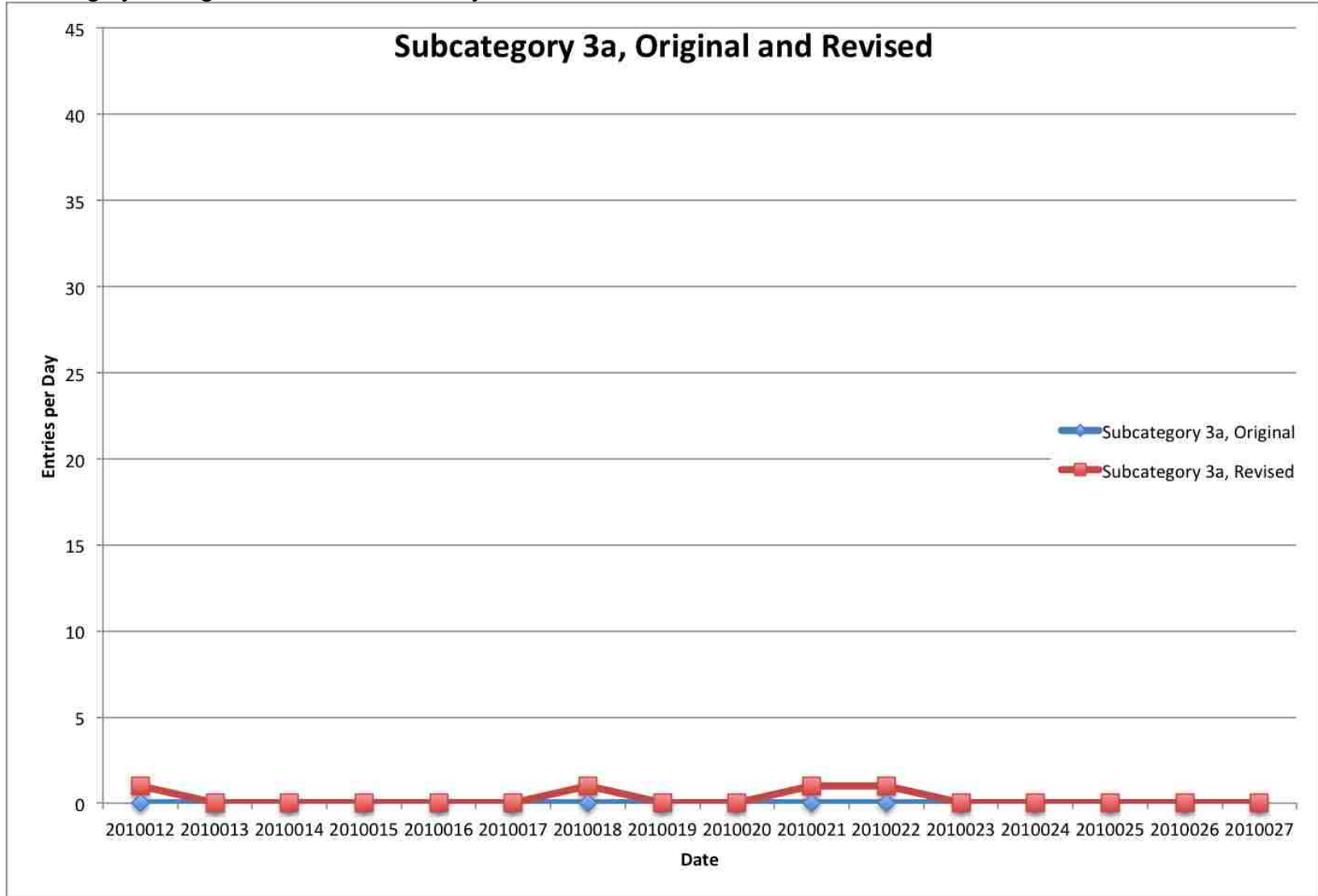
Subcategory 2e original and revised entries by date



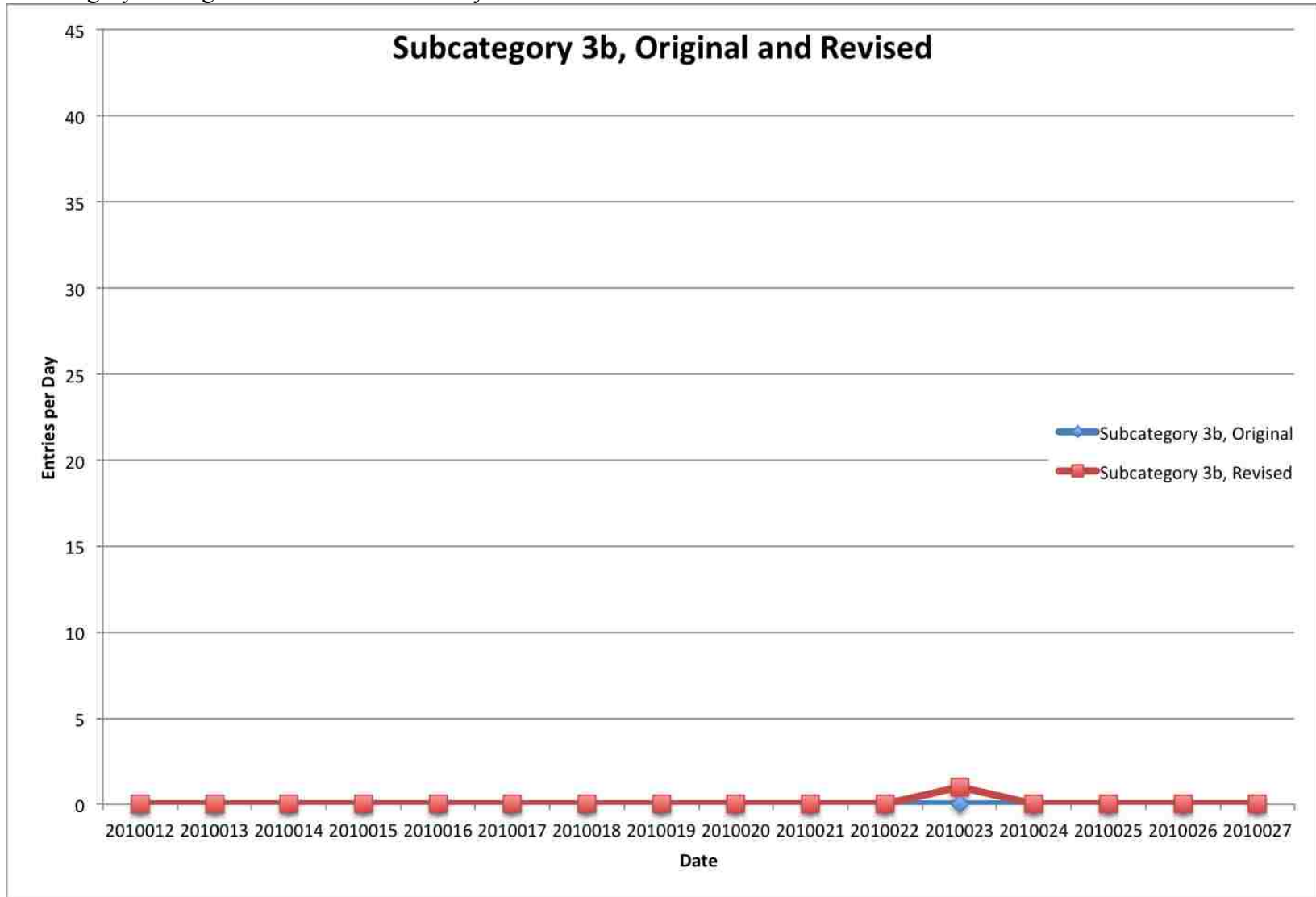
Subcategory 2f original and revised entries by date



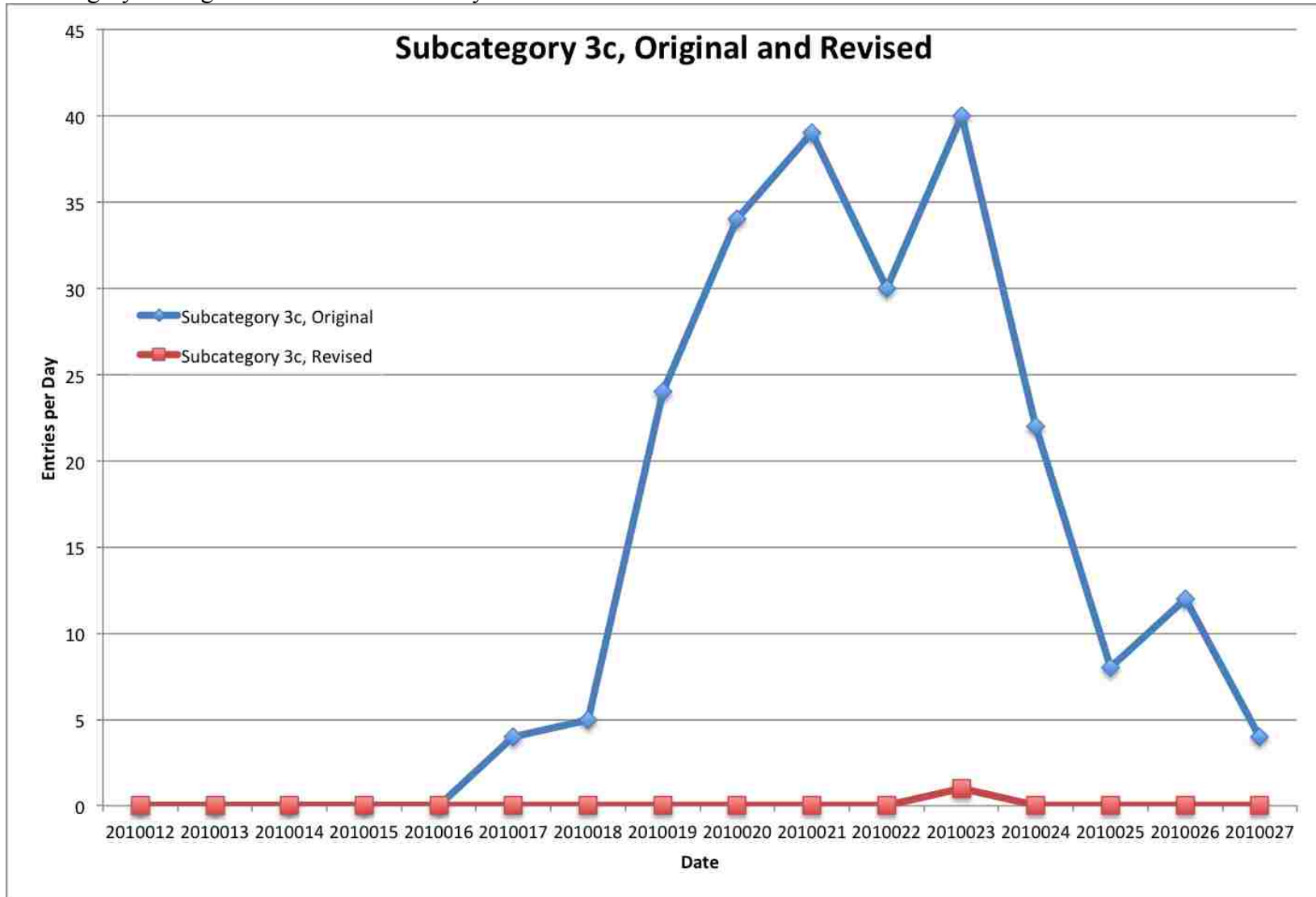
Subcategory 3a original and revised entries by date



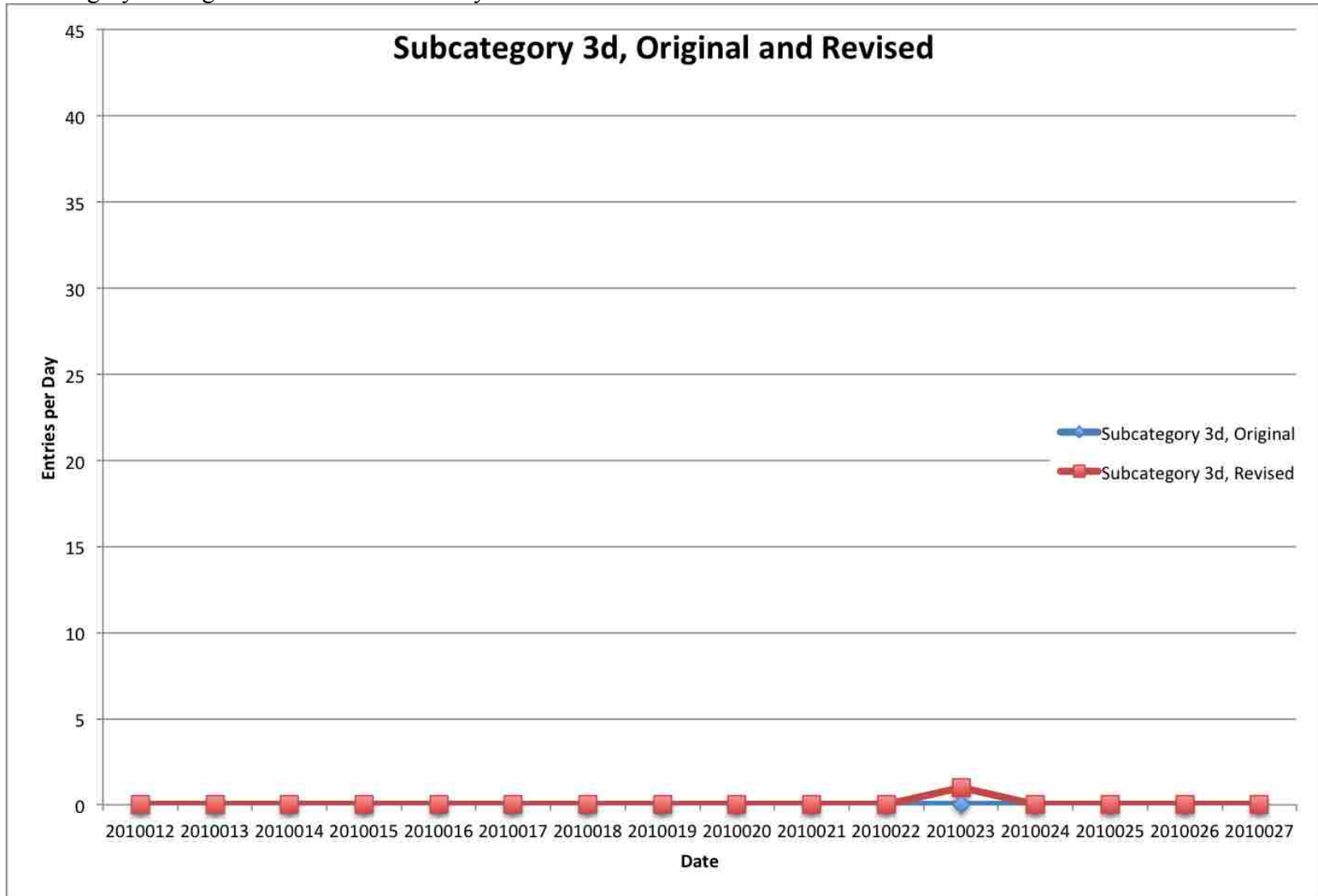
Subcategory 3b original and revised entries by date



Subcategory 3c original and revised entries by date

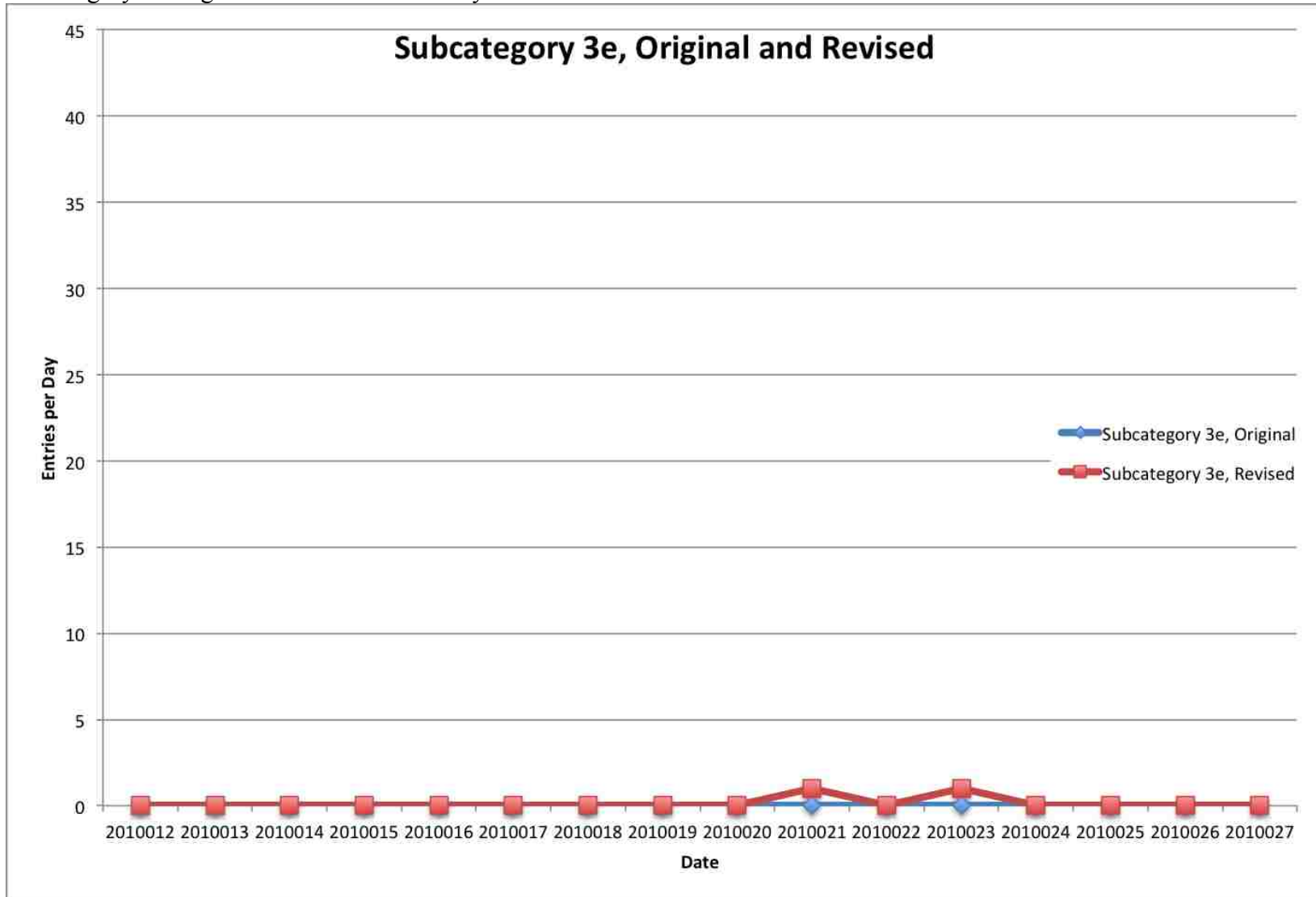


Subcategory 3d original and revised entries by date

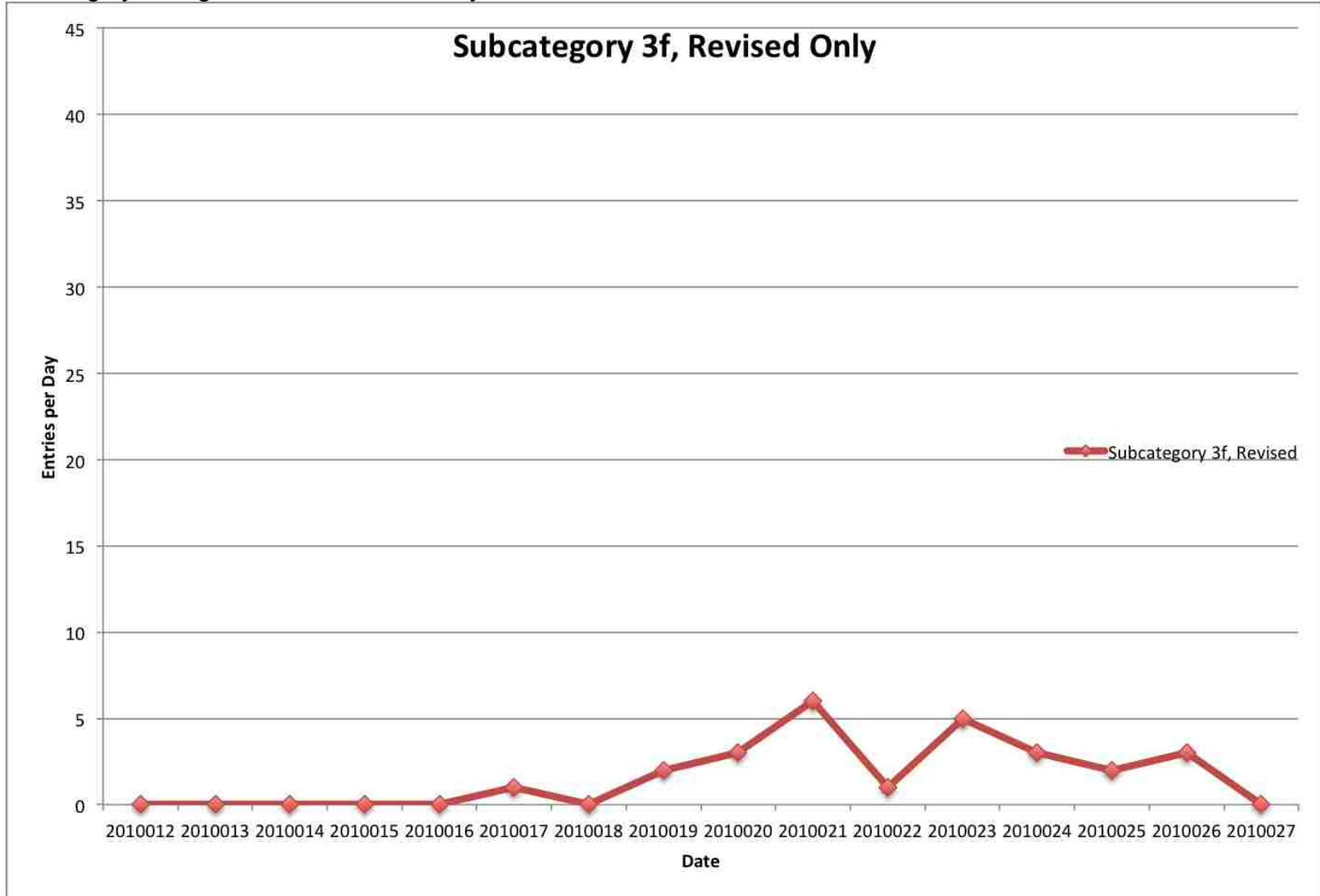




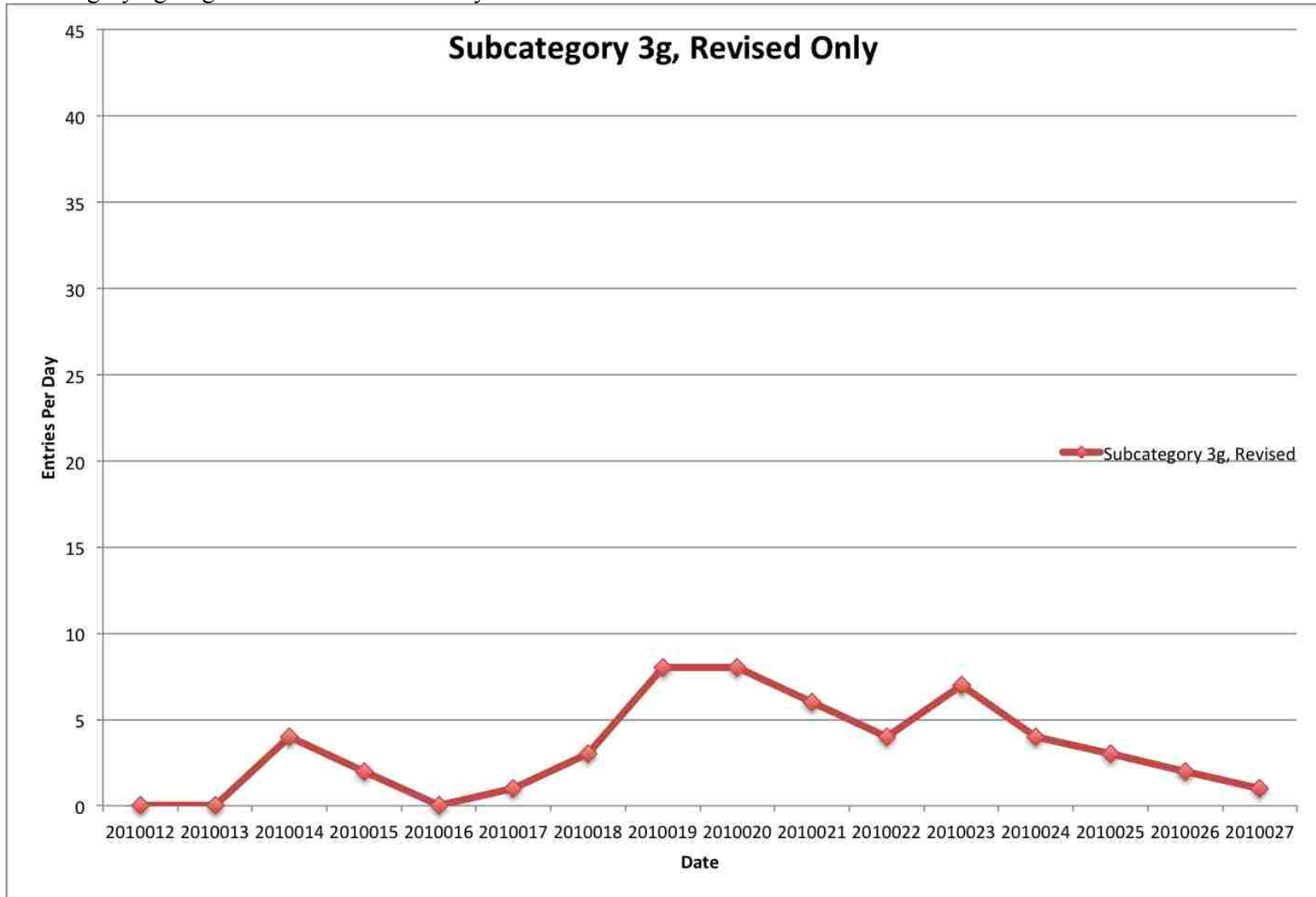
Subcategory 3e original and revised entries by date



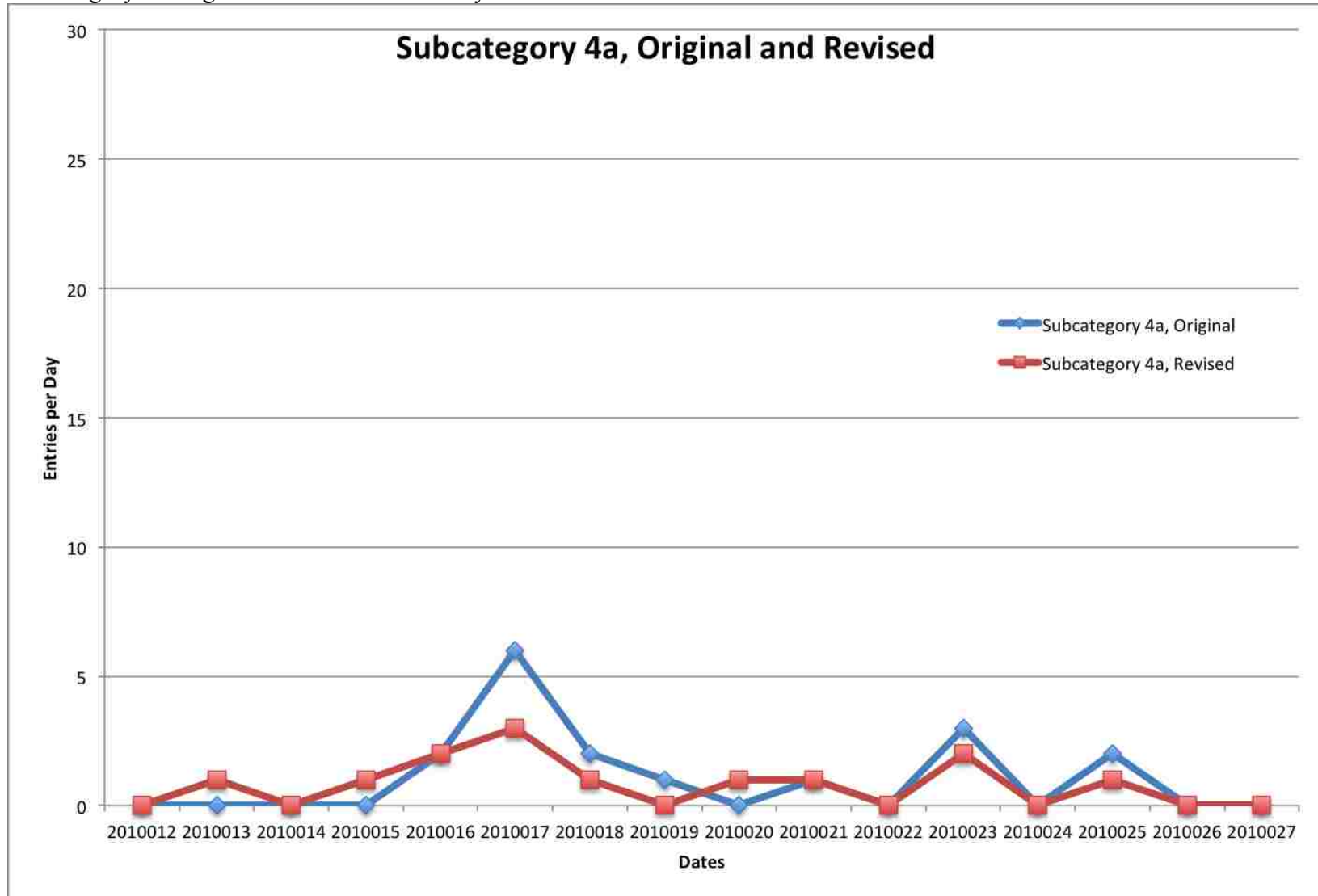
Subcategory 3f original and revised entries by date



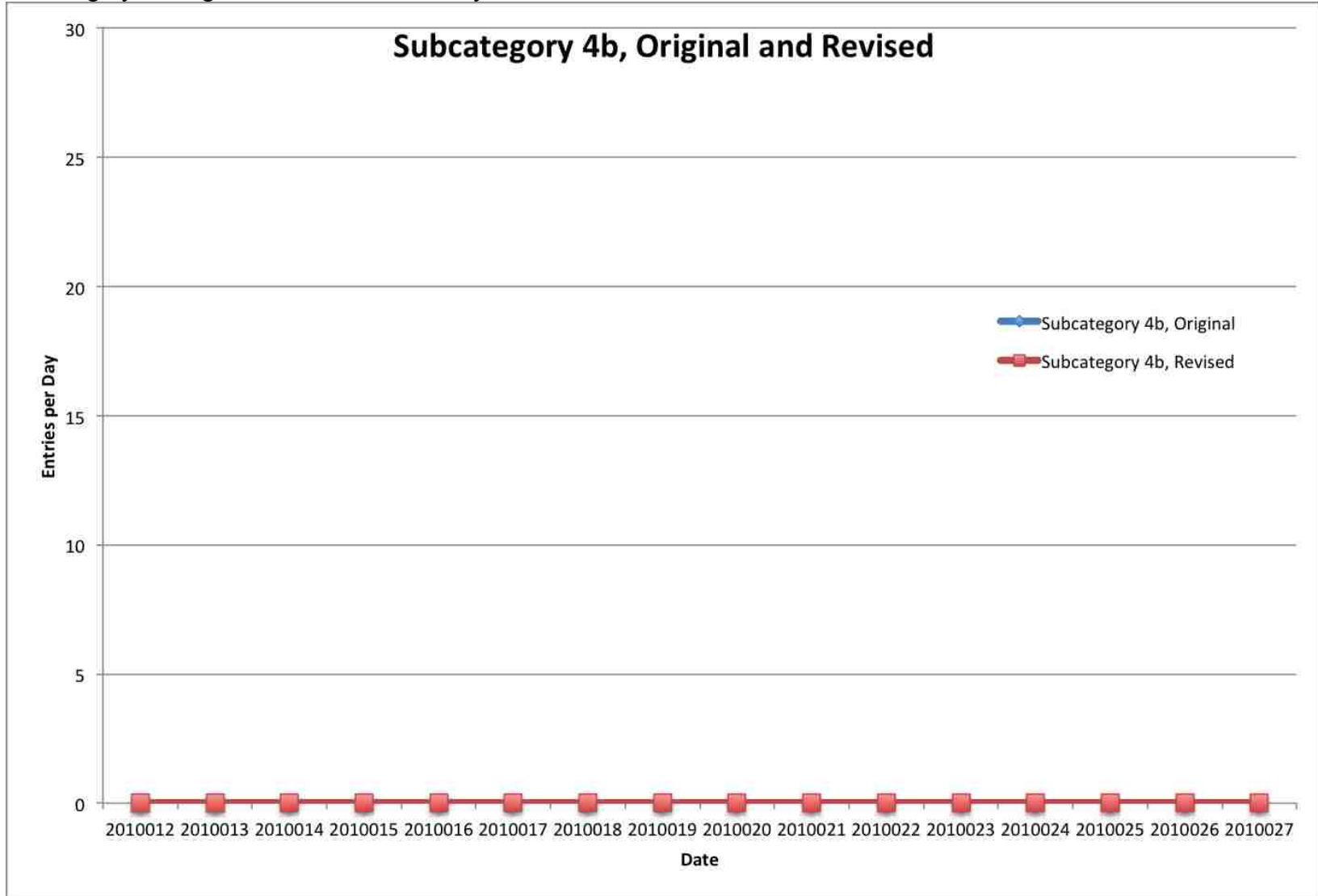
Subcategory 3g original and revised entries by date



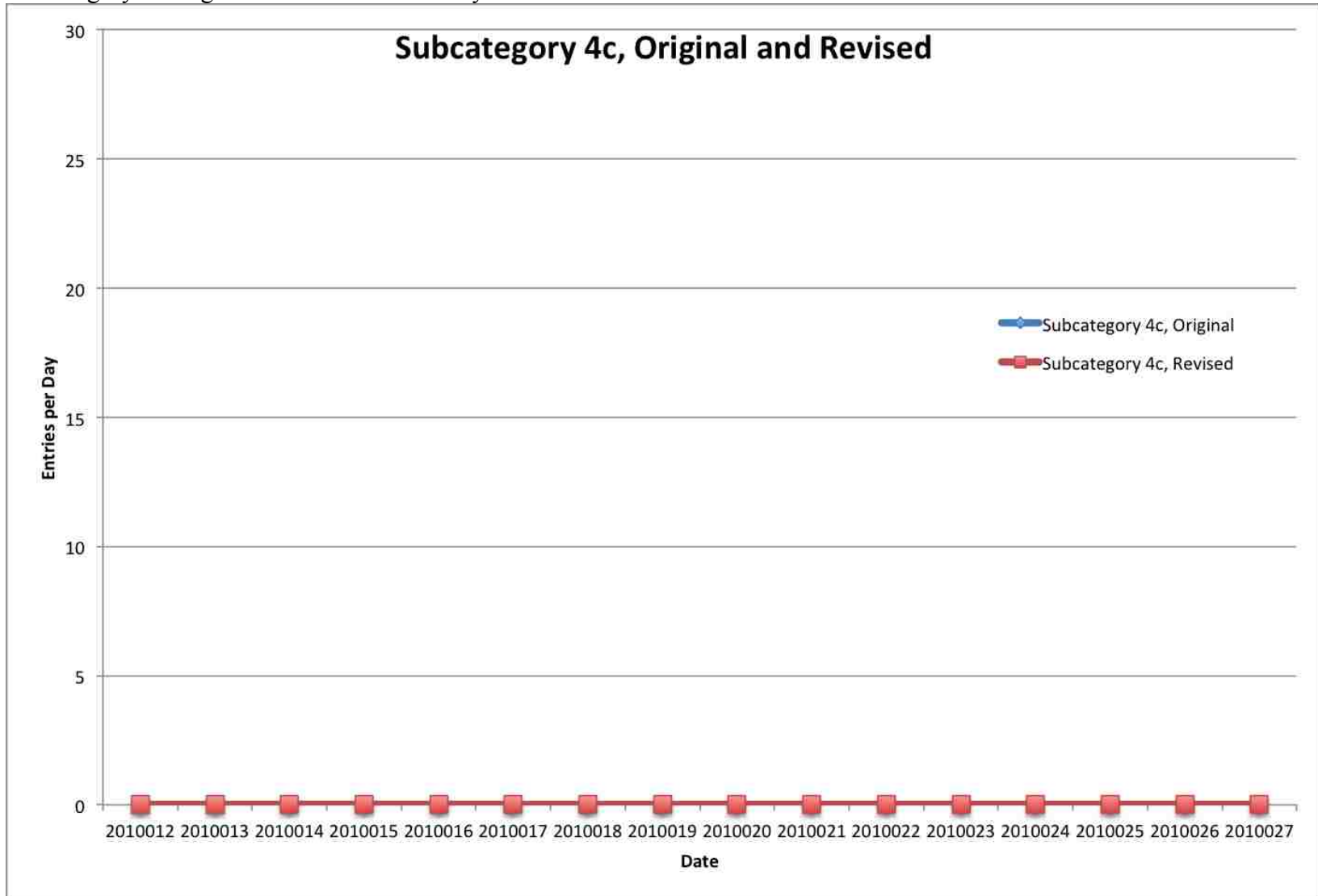
Subcategory 4a original and revised entries by date



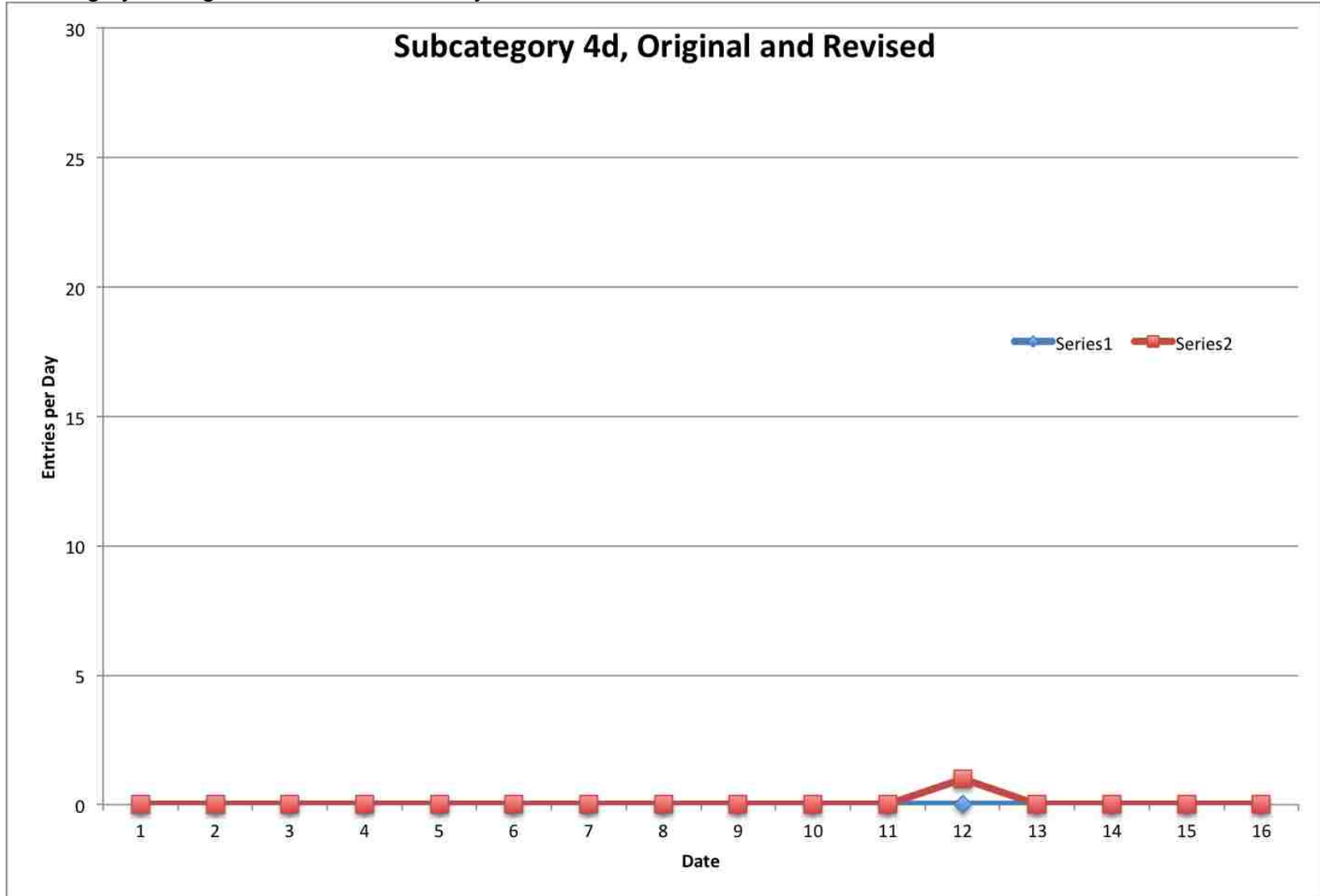
Subcategory 4b original and revised entries by date



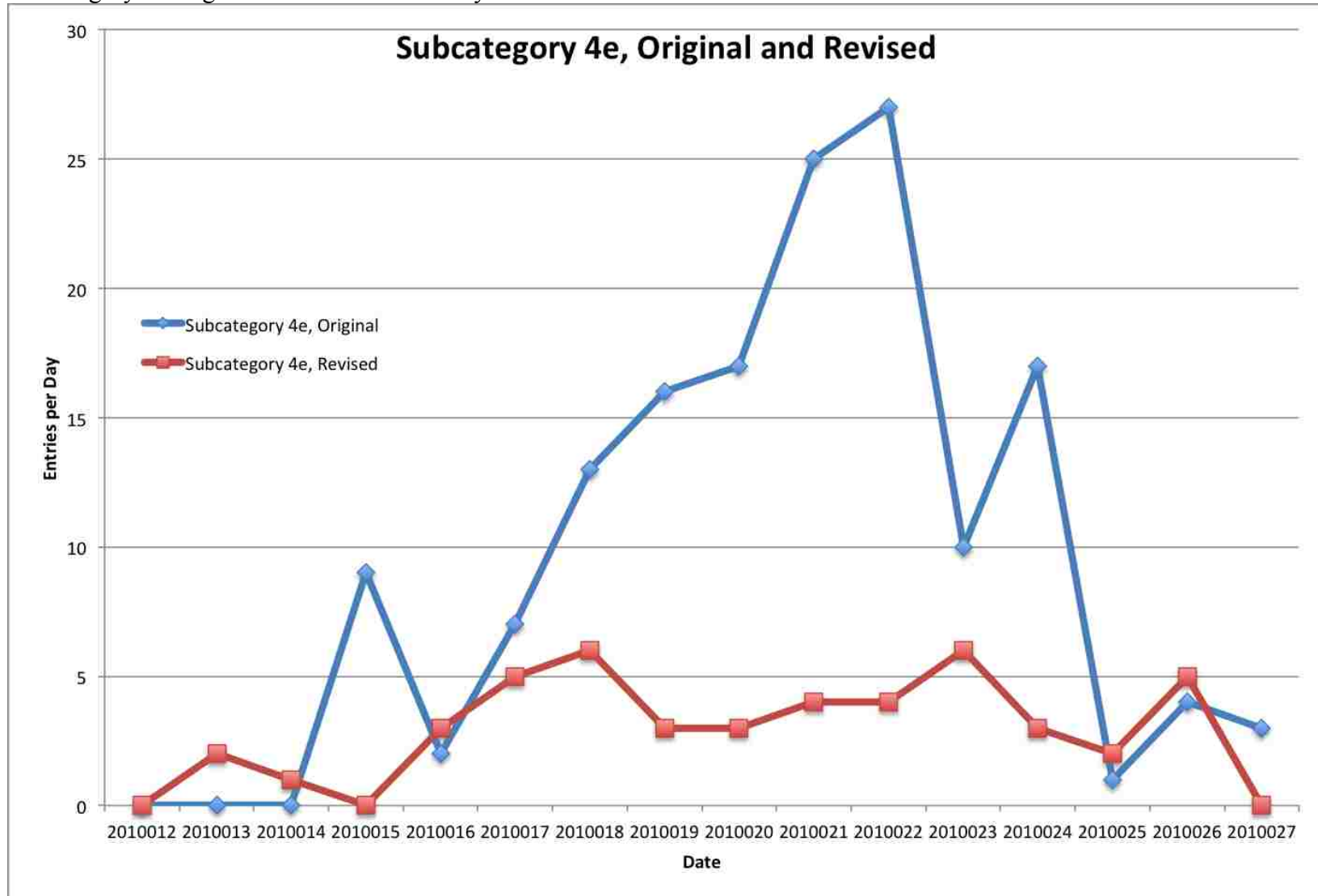
Subcategory 4c original and revised entries by date



Subcategory 4d original and revised entries by date

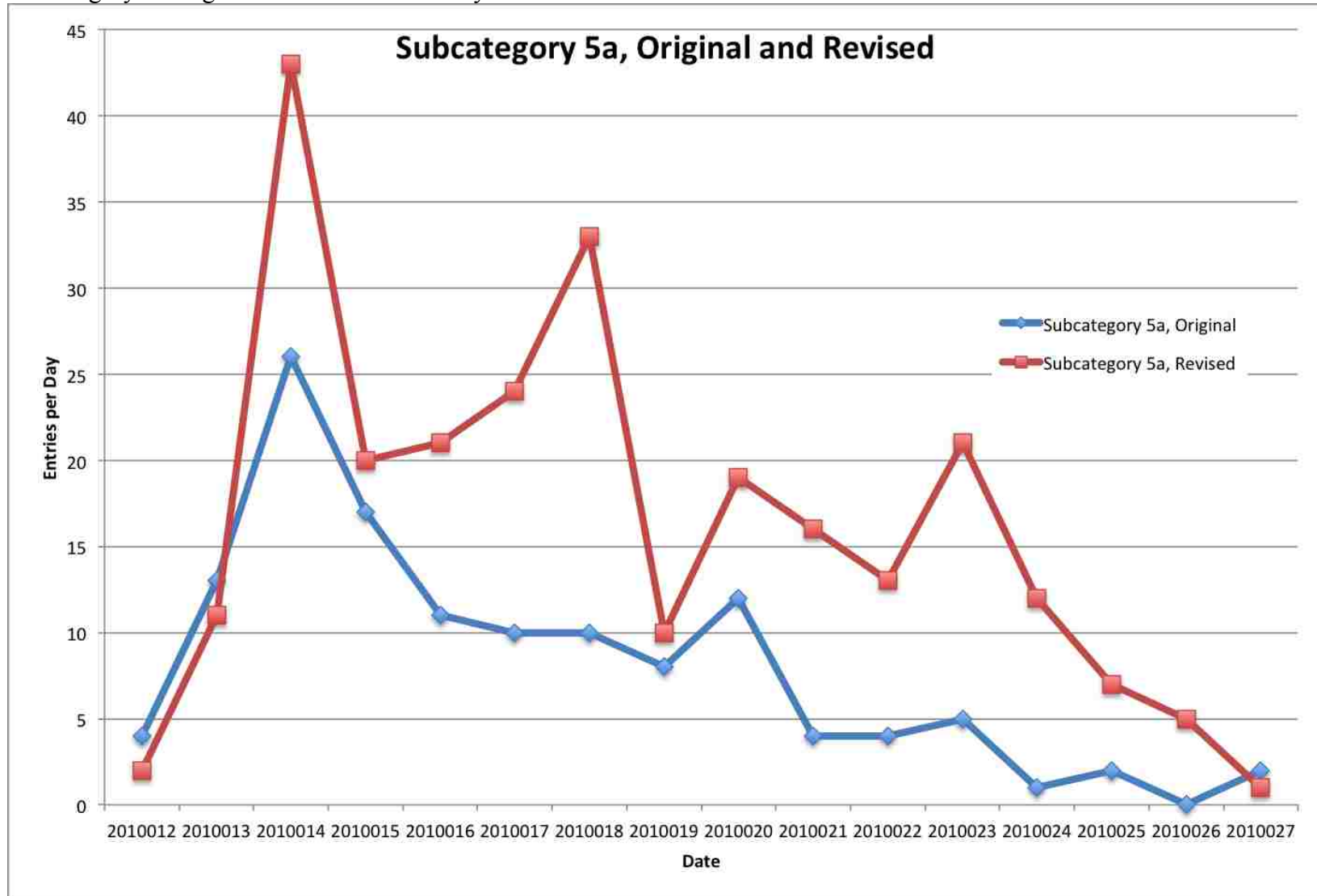


Subcategory 4e original and revised entries by date

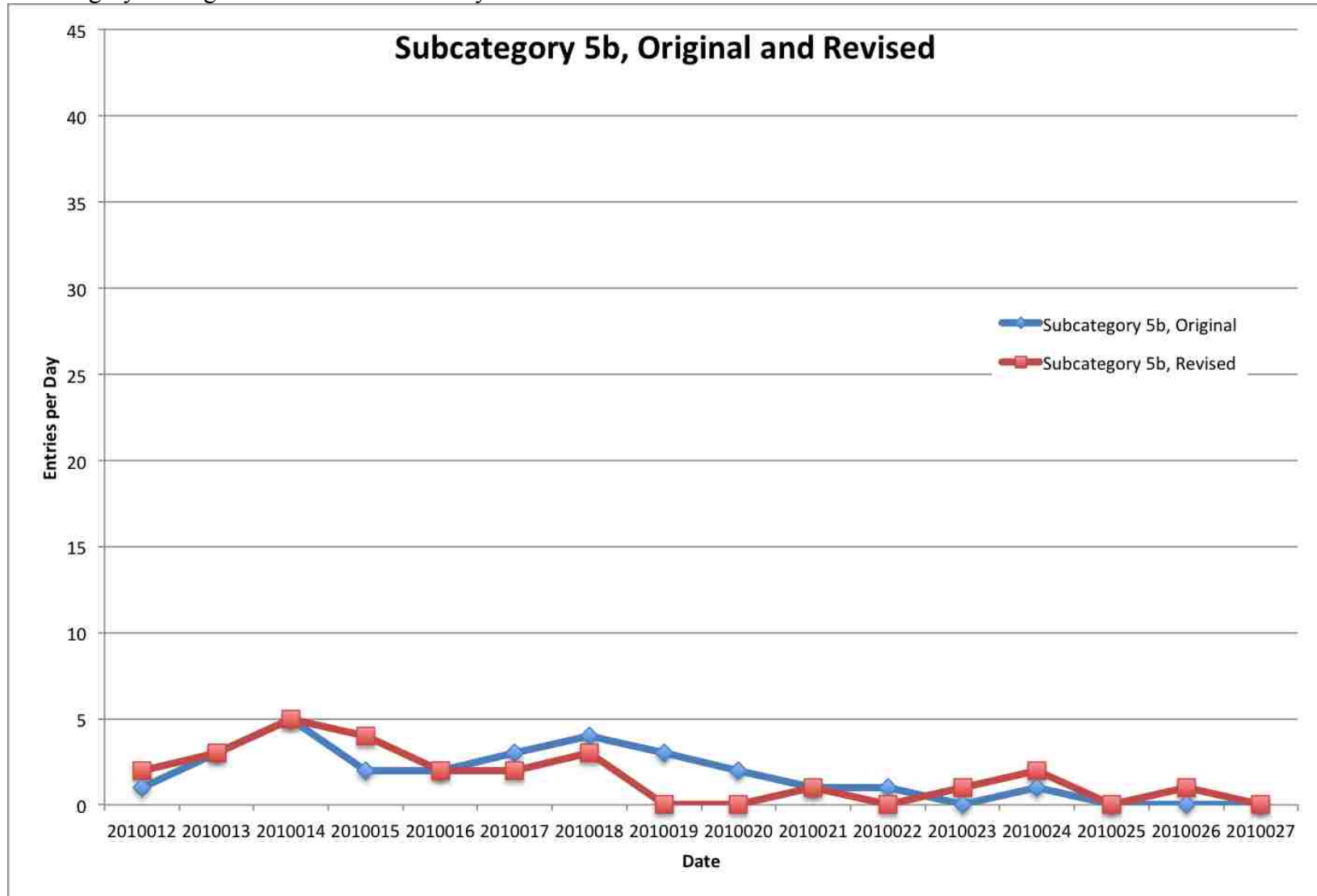




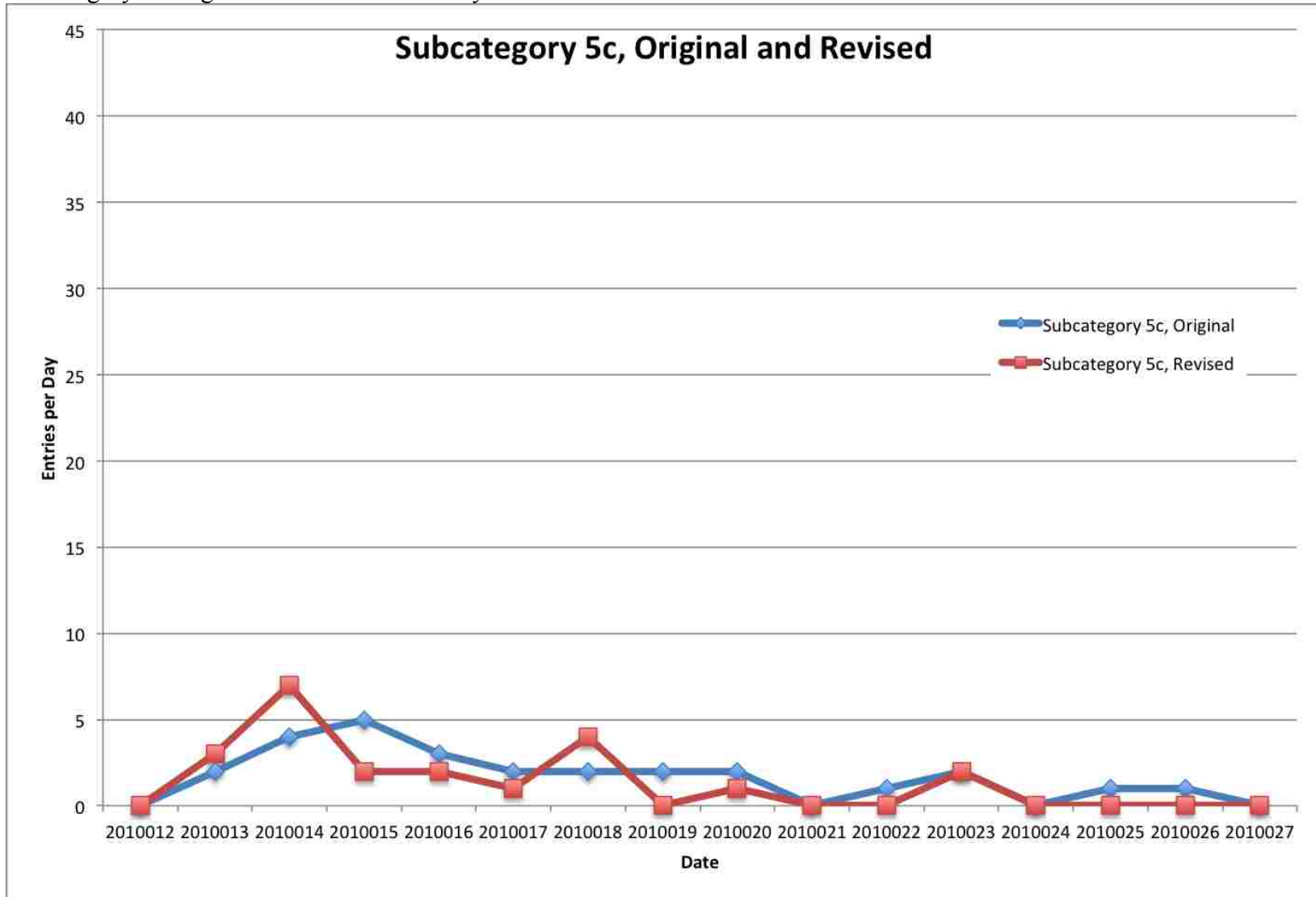
Subcategory 5a original and revised entries by date



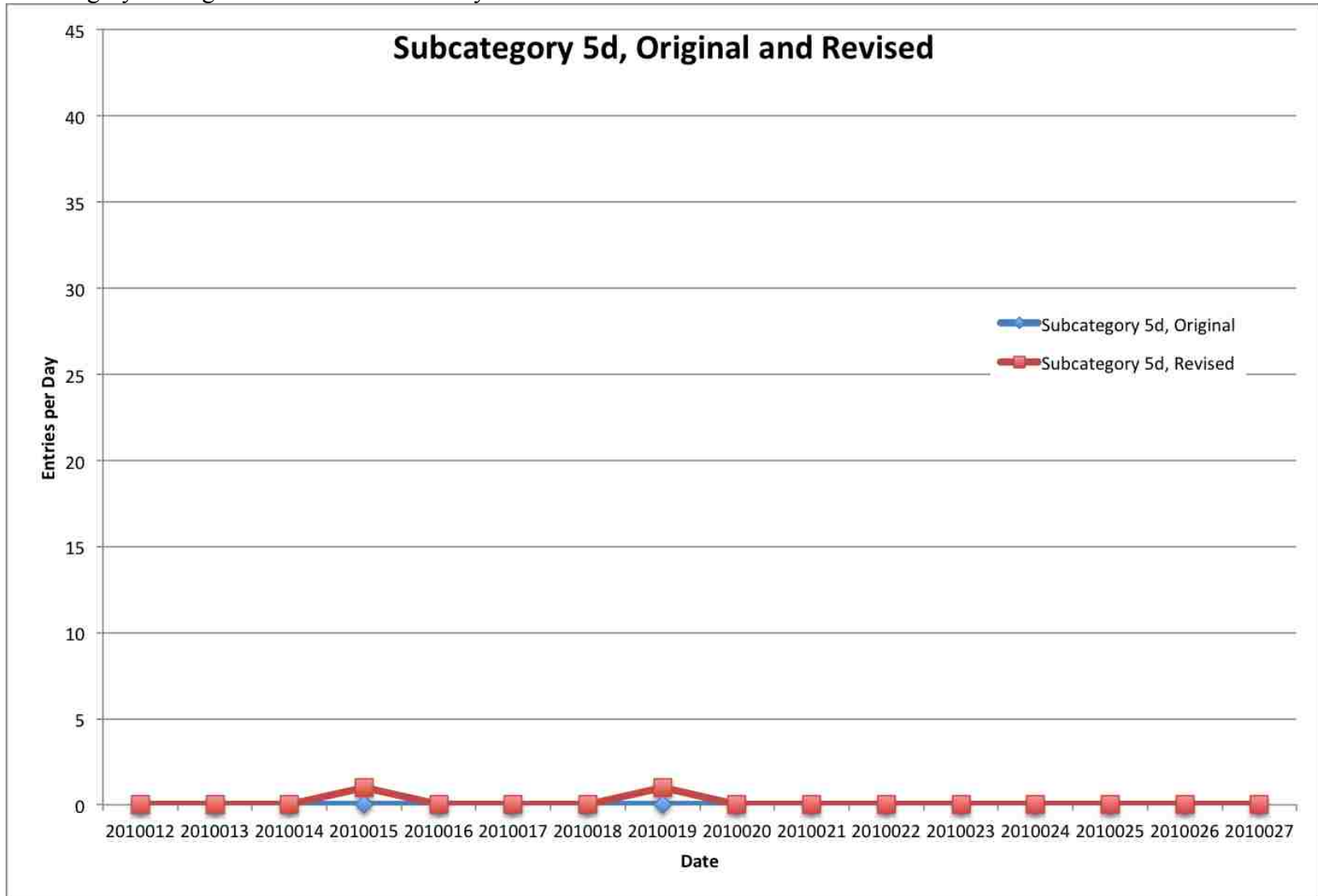
Subcategory 5b original and revised entries by date



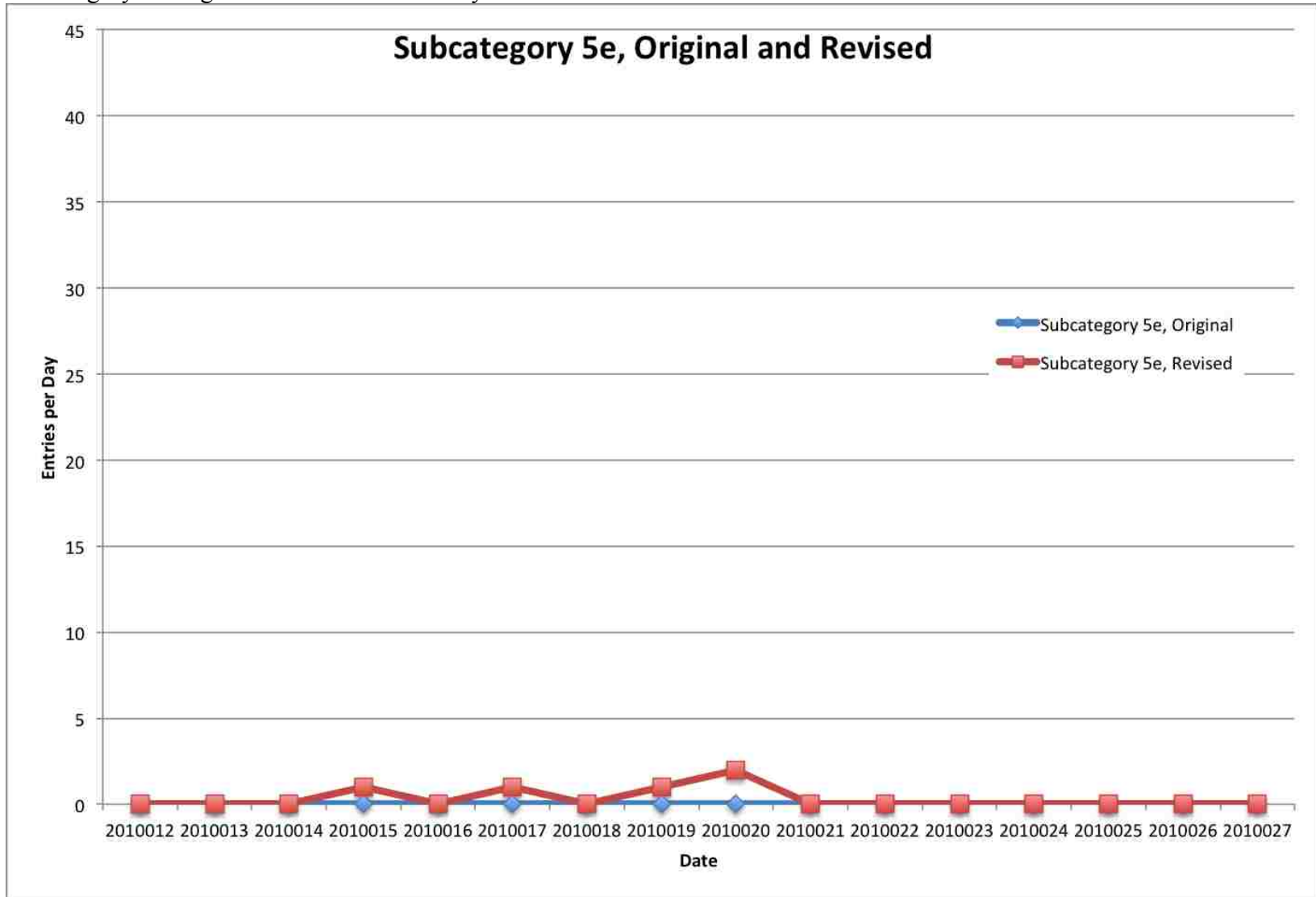
Subcategory 5c original and revised entries by date



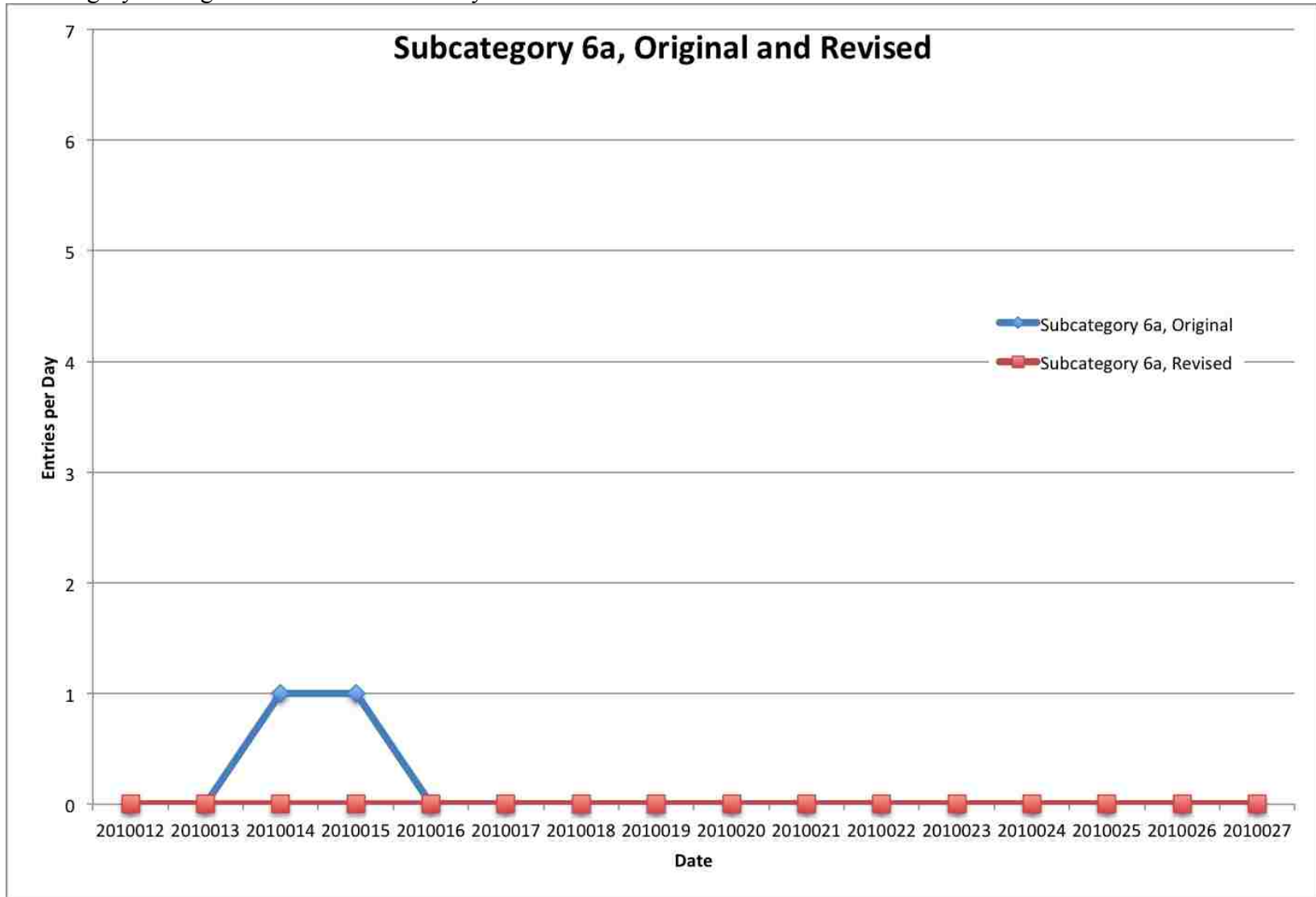
Subcategory 5d original and revised entries by date



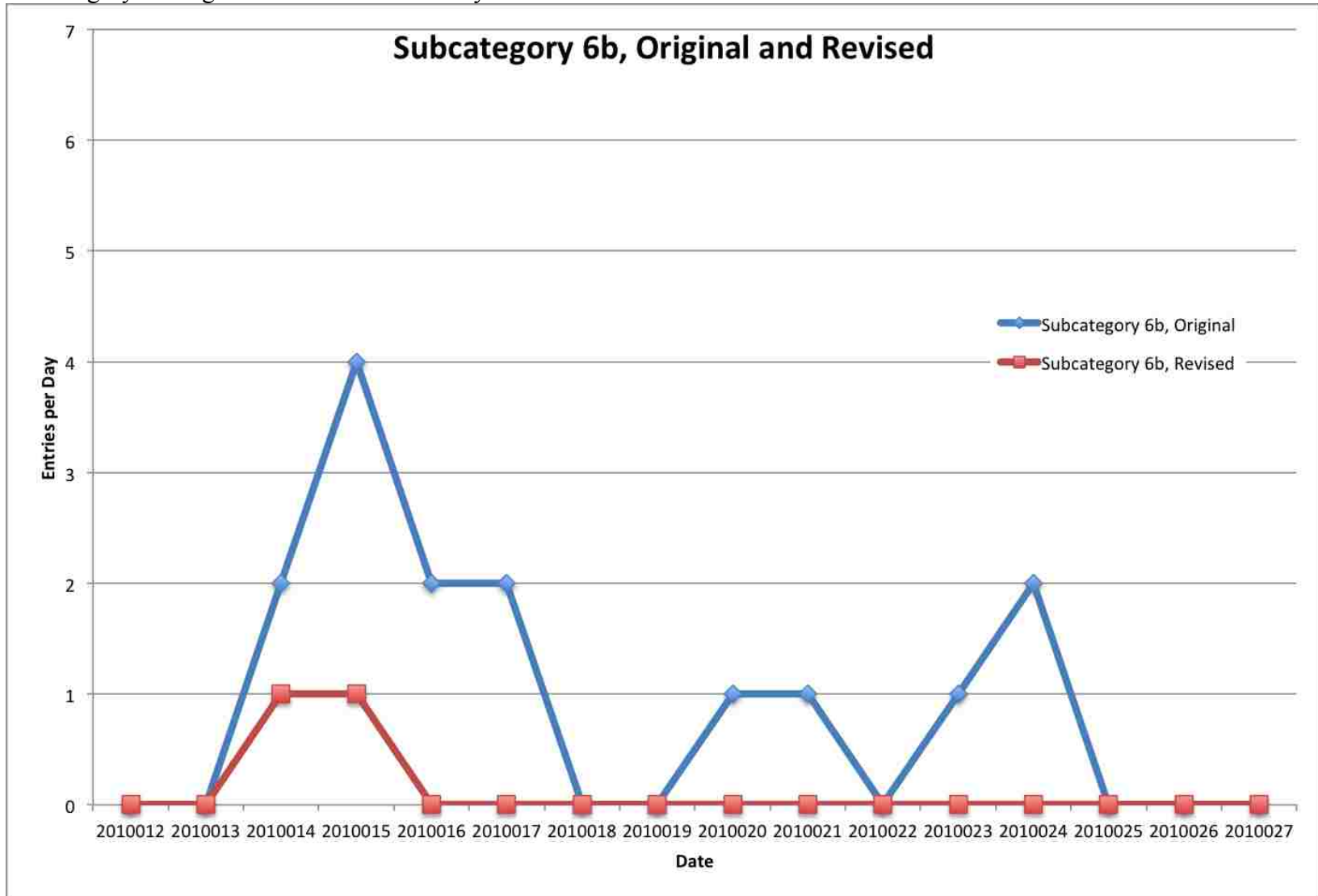
Subcategory 5e original and revised entries by date



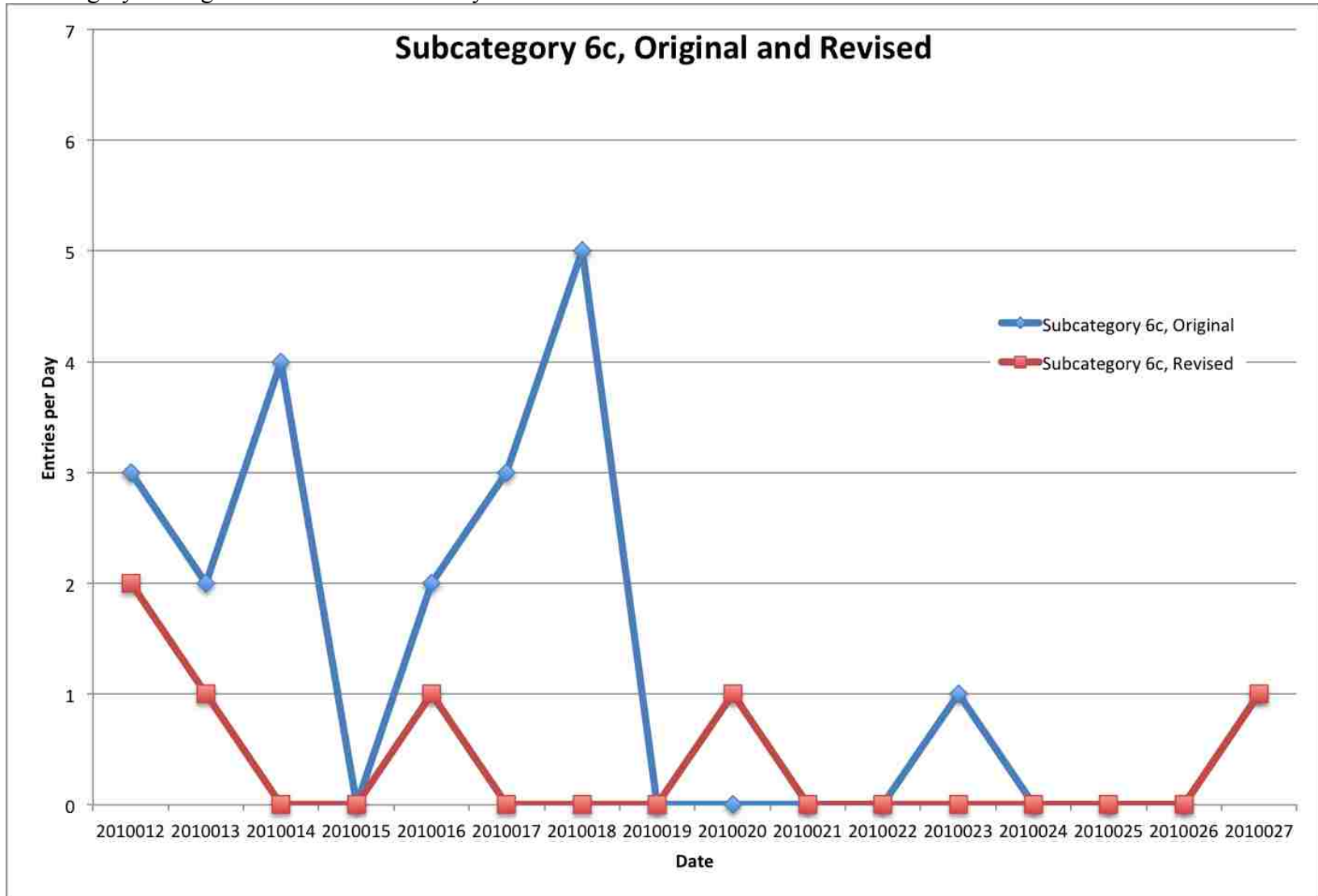
Subcategory 6a original and revised entries by date



Subcategory 6b original and revised entries by date

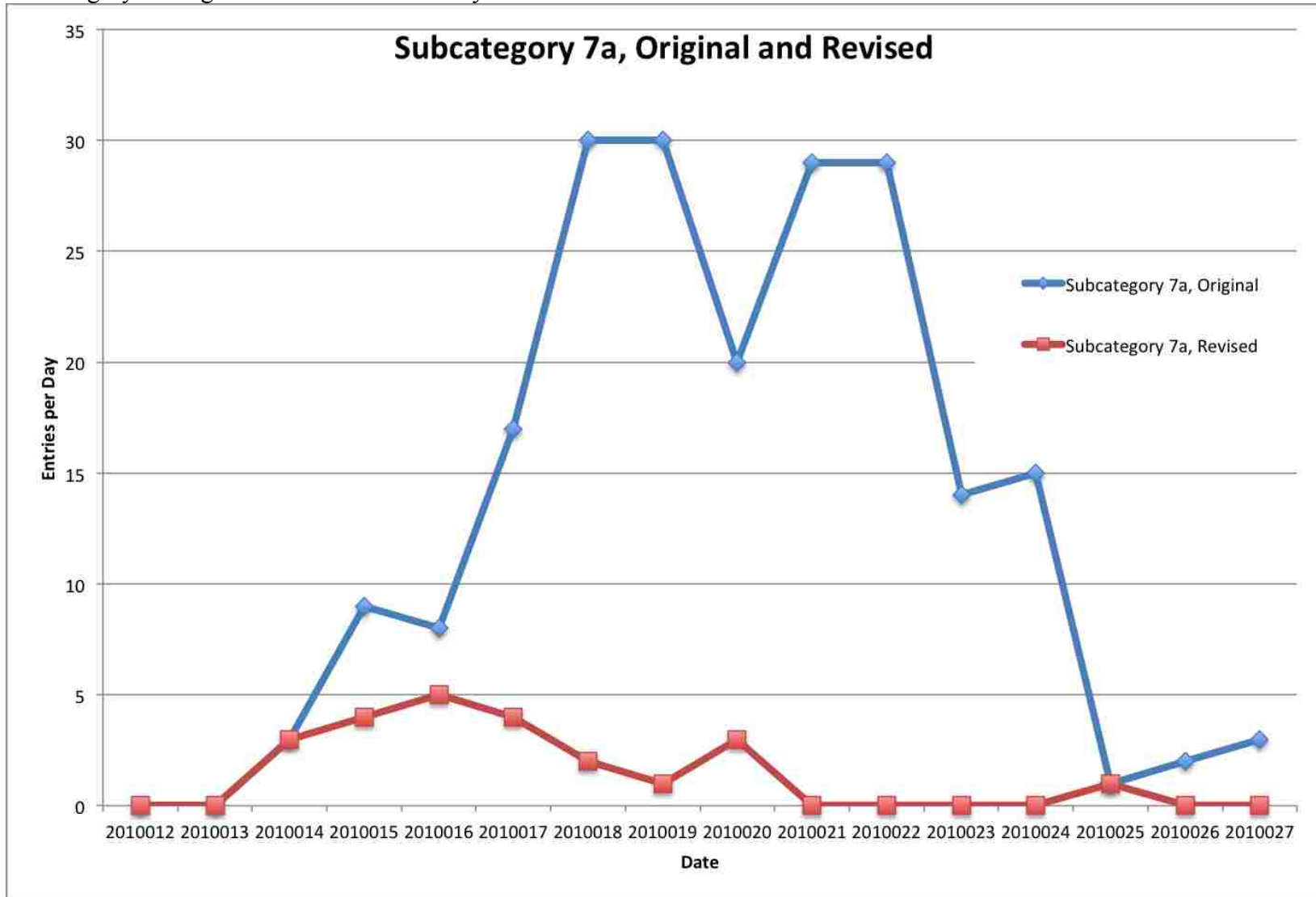


Subcategory 6c original and revised entries by date

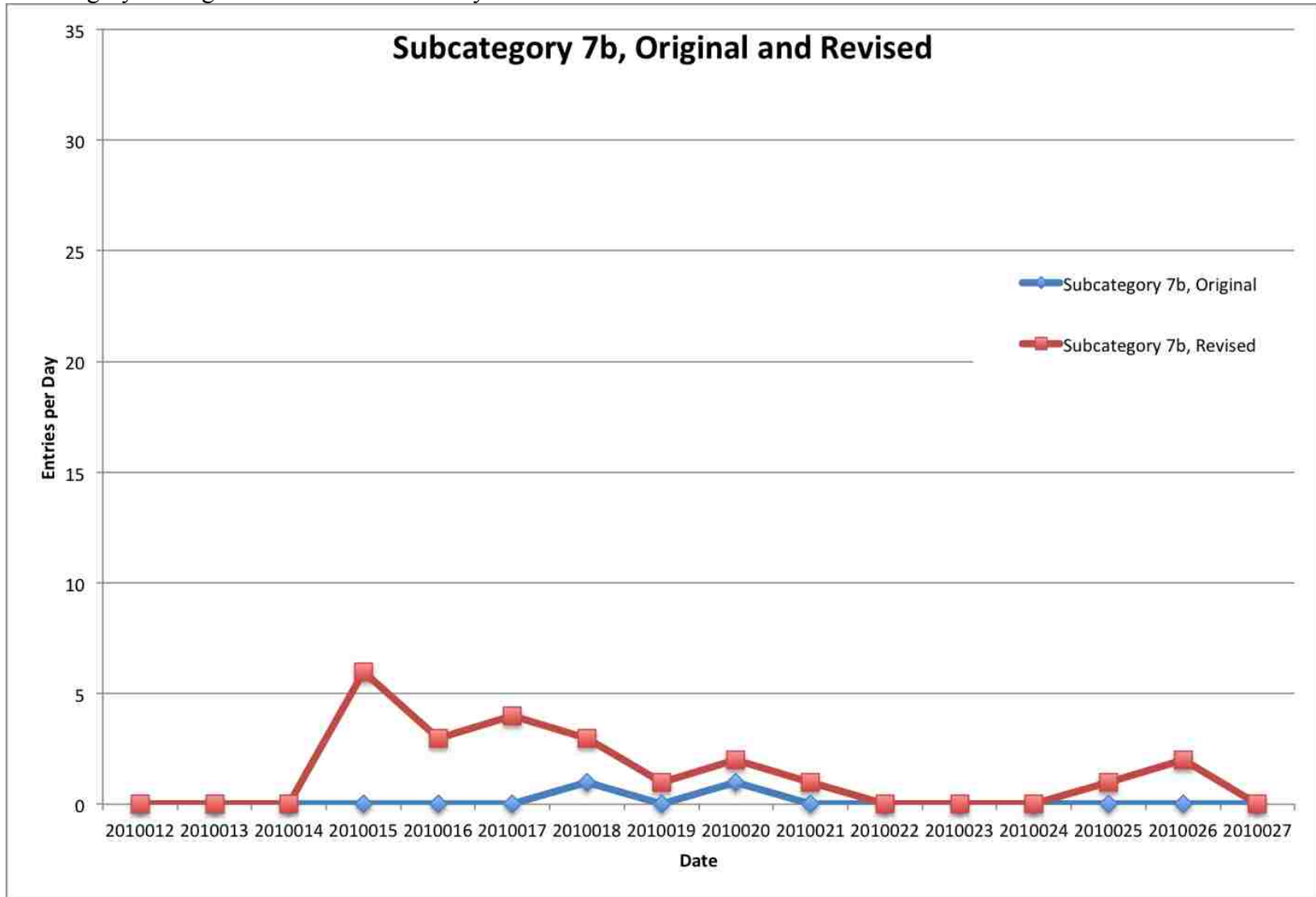




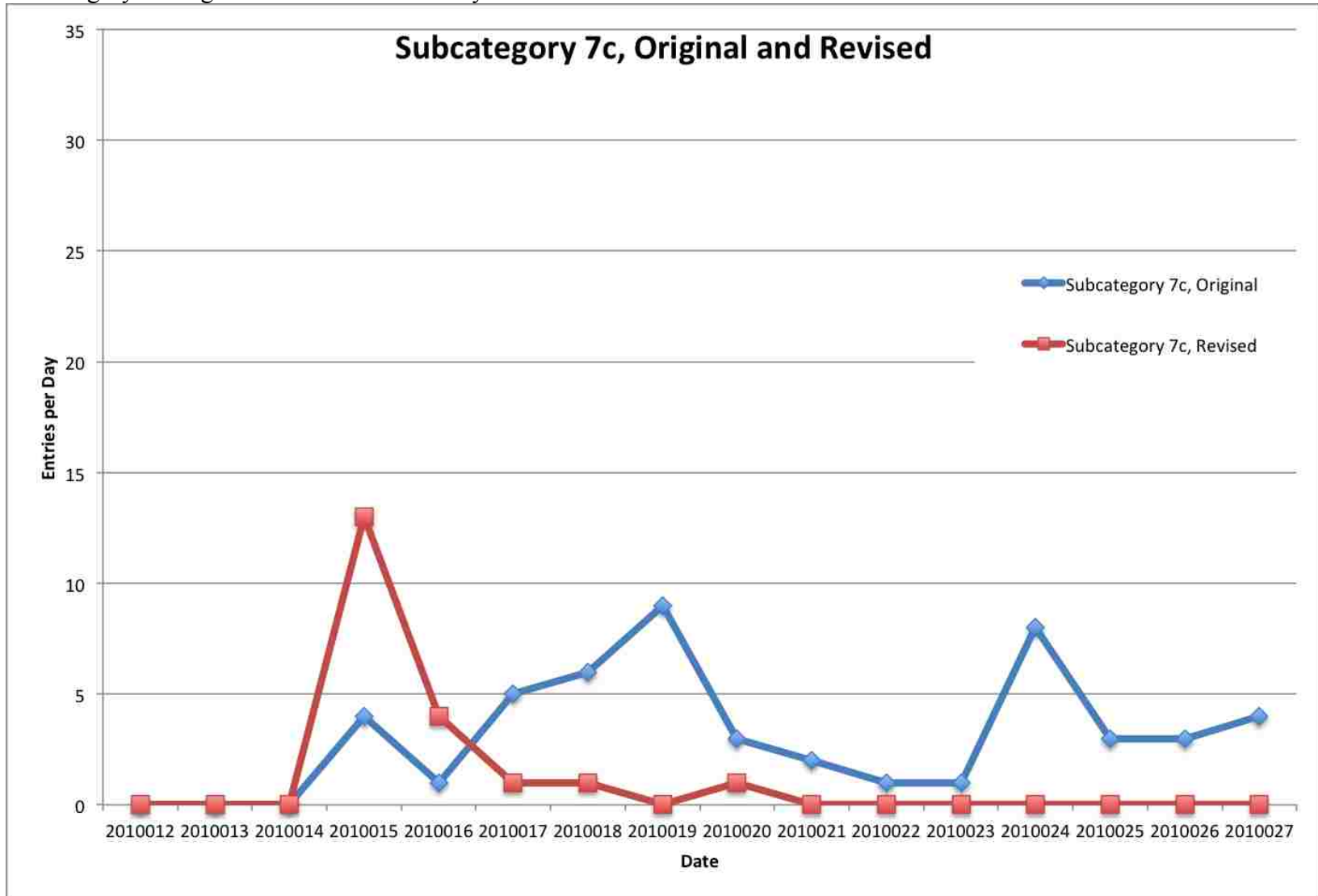
Subcategory 7a original and revised entries by date



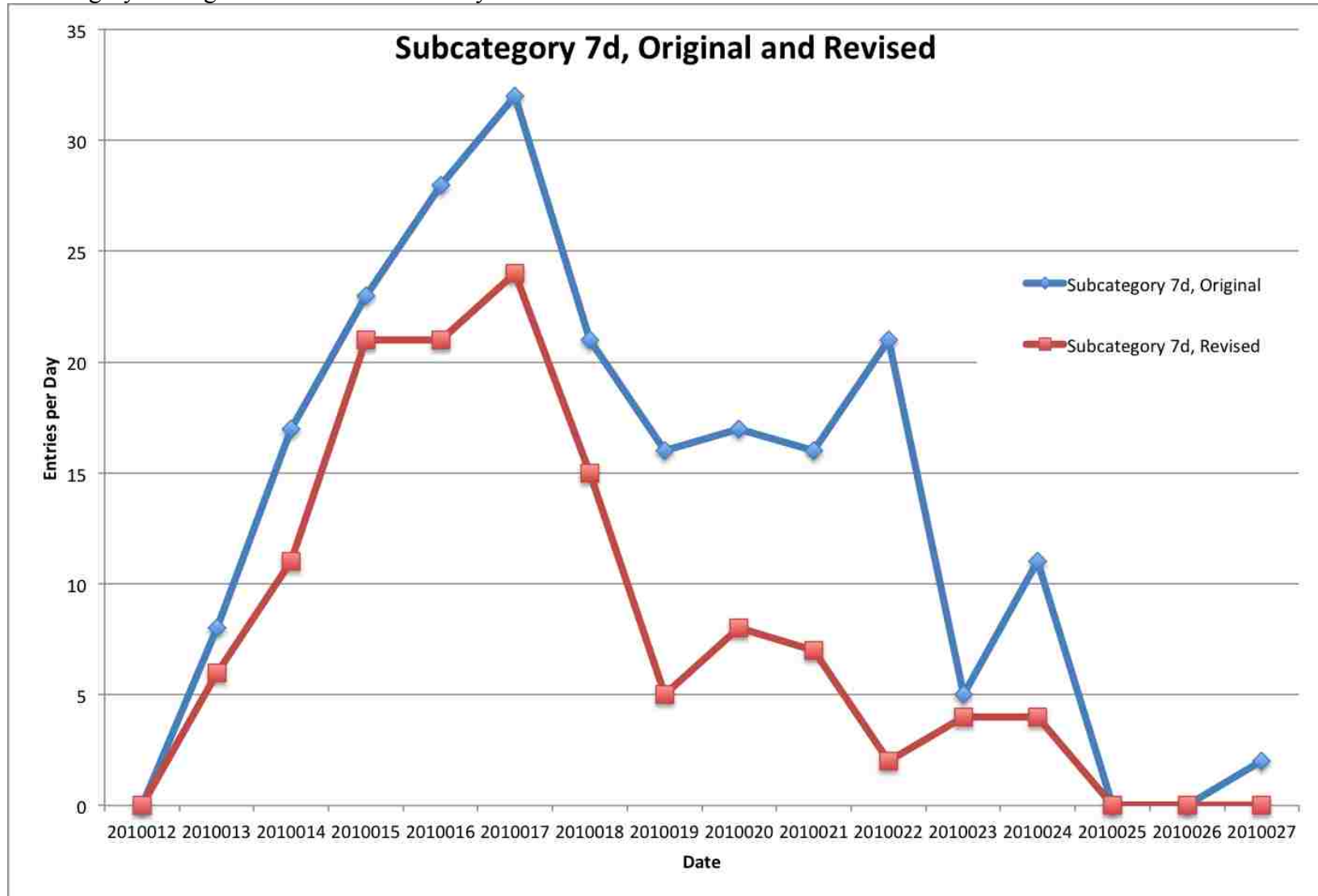
Subcategory 7b original and revised entries by date



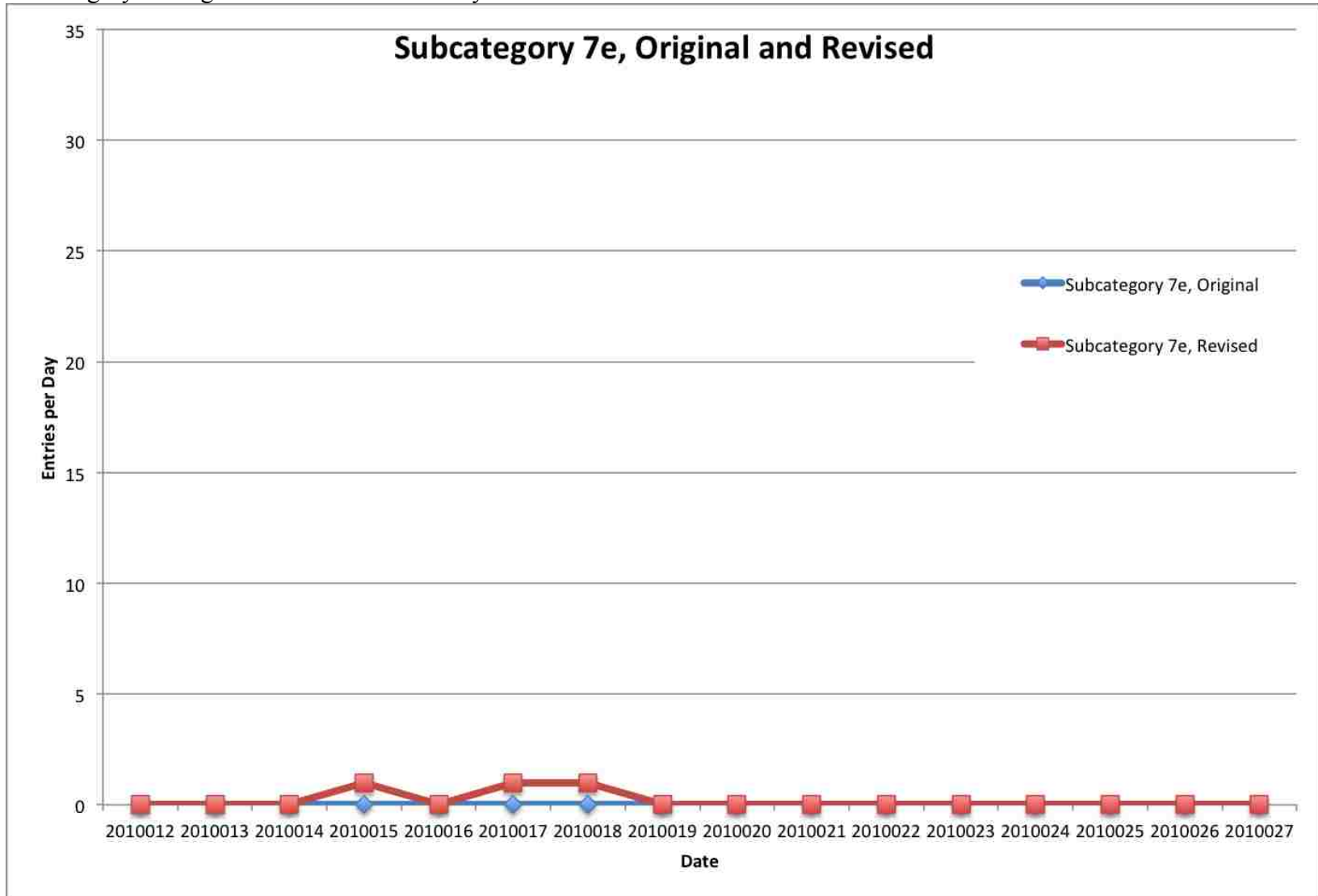
Subcategory 7c original and revised entries by date



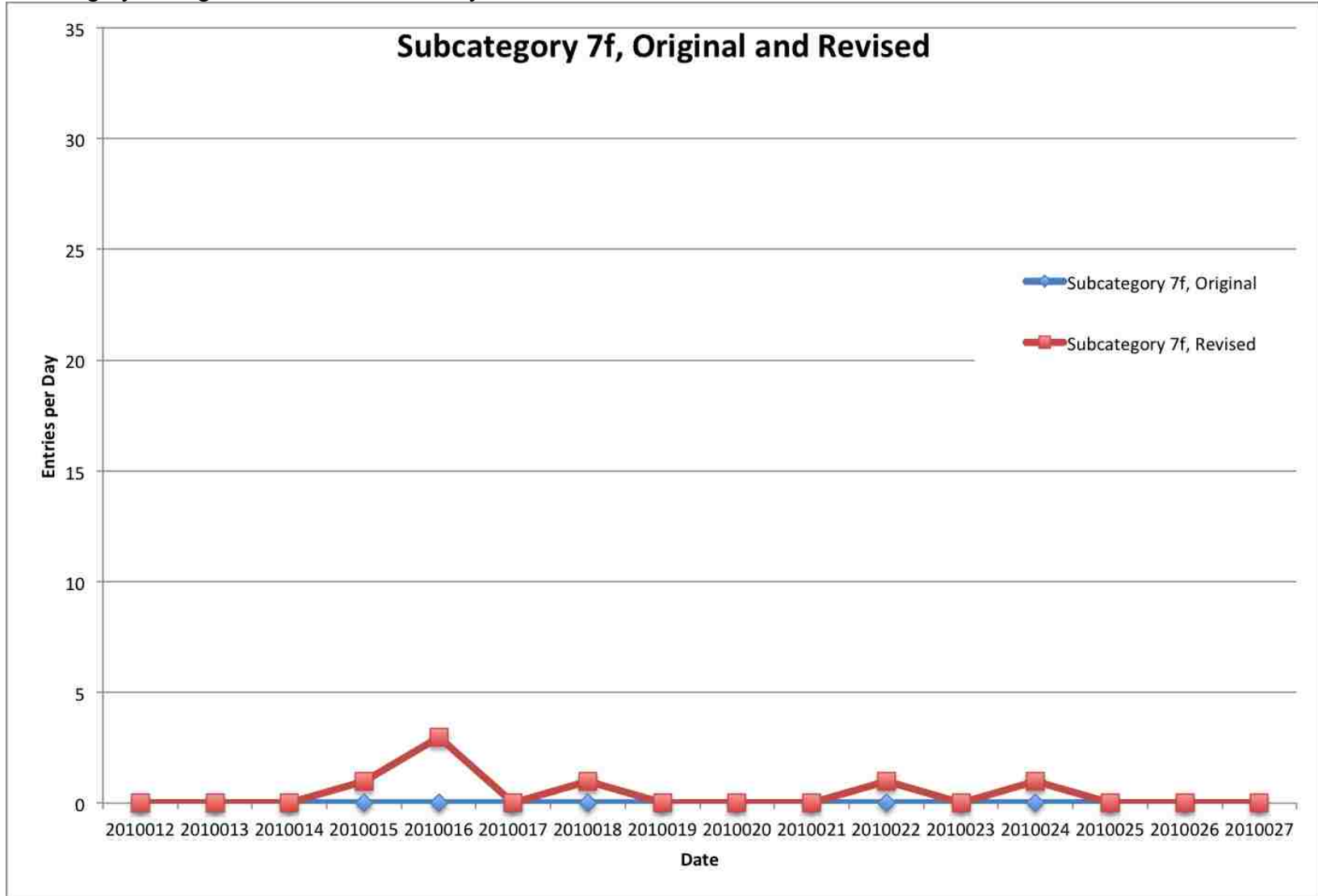
Subcategory 7d original and revised entries by date



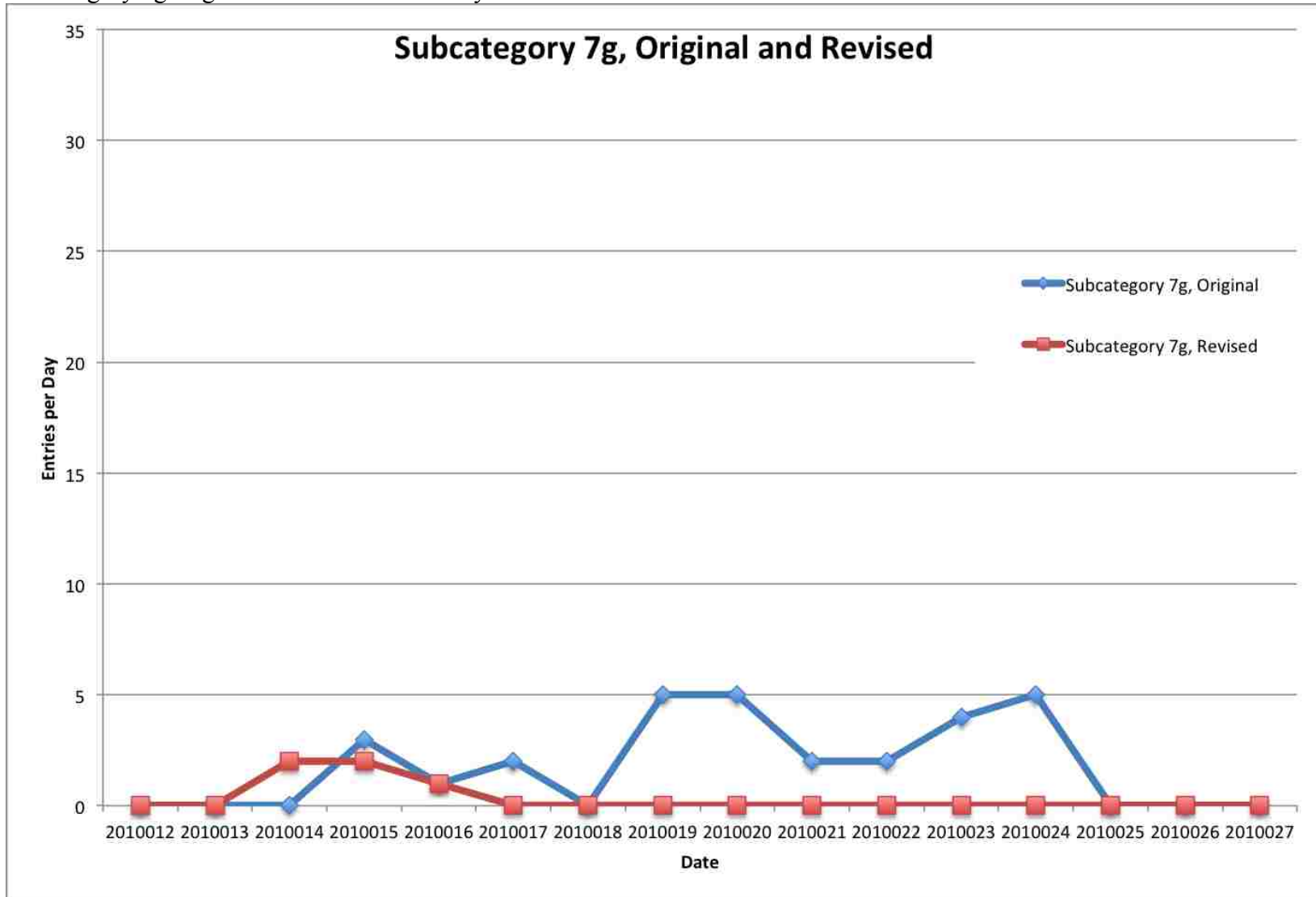
Subcategory 7e original and revised entries by date



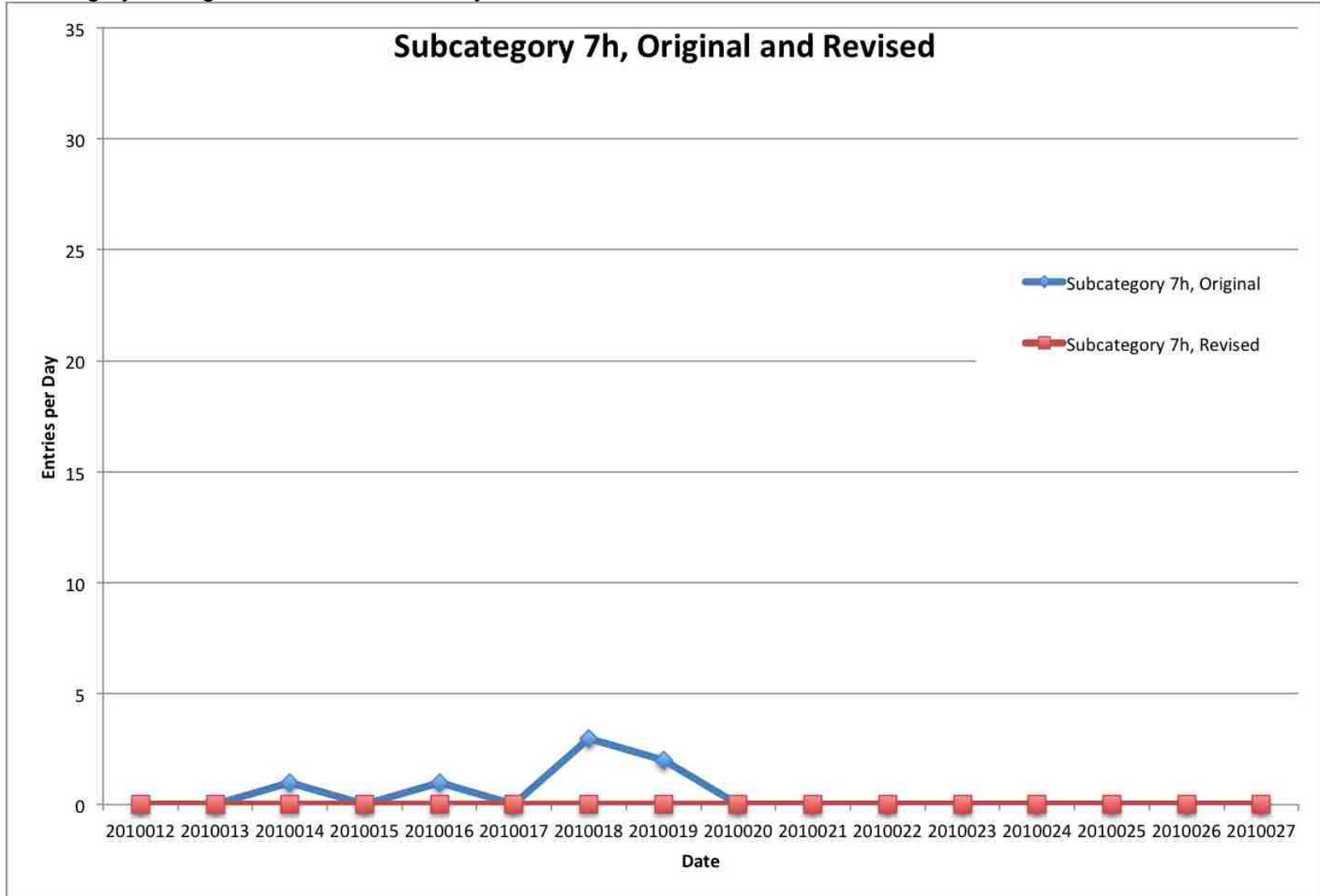
Subcategory 7f original and revised entries by date



Subcategory 7g original and revised entries by date

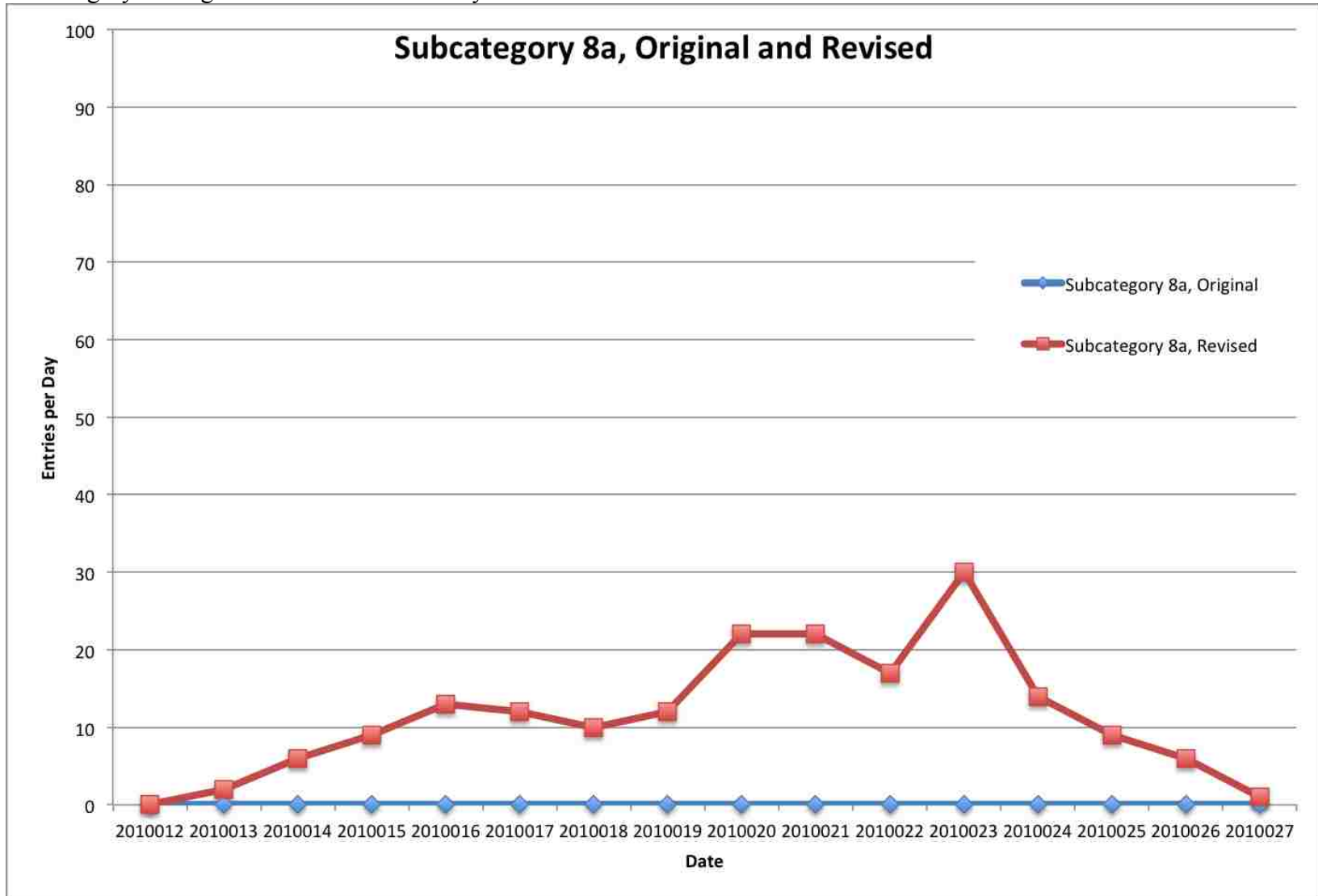


Subcategory 7h original and revised entries by date

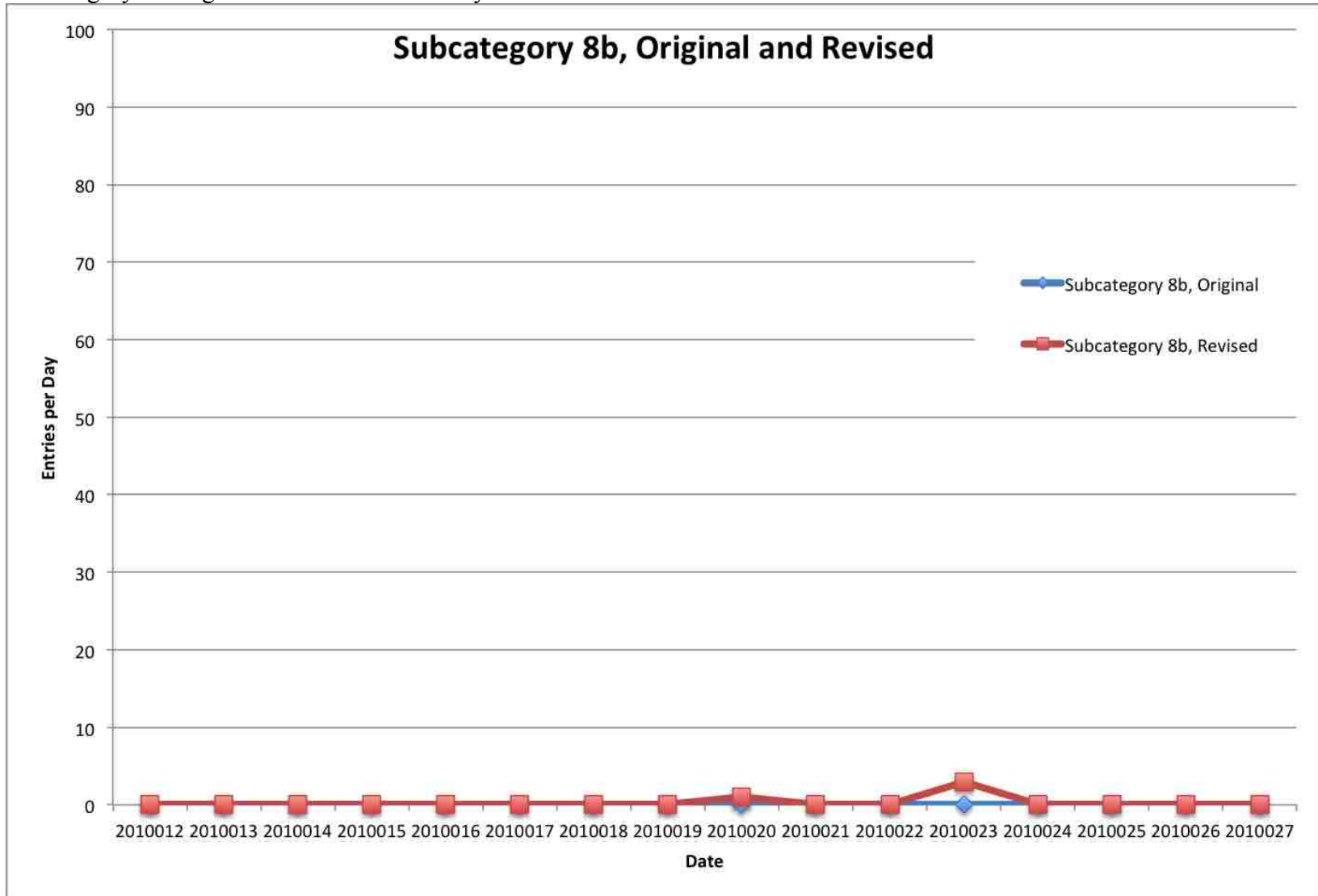




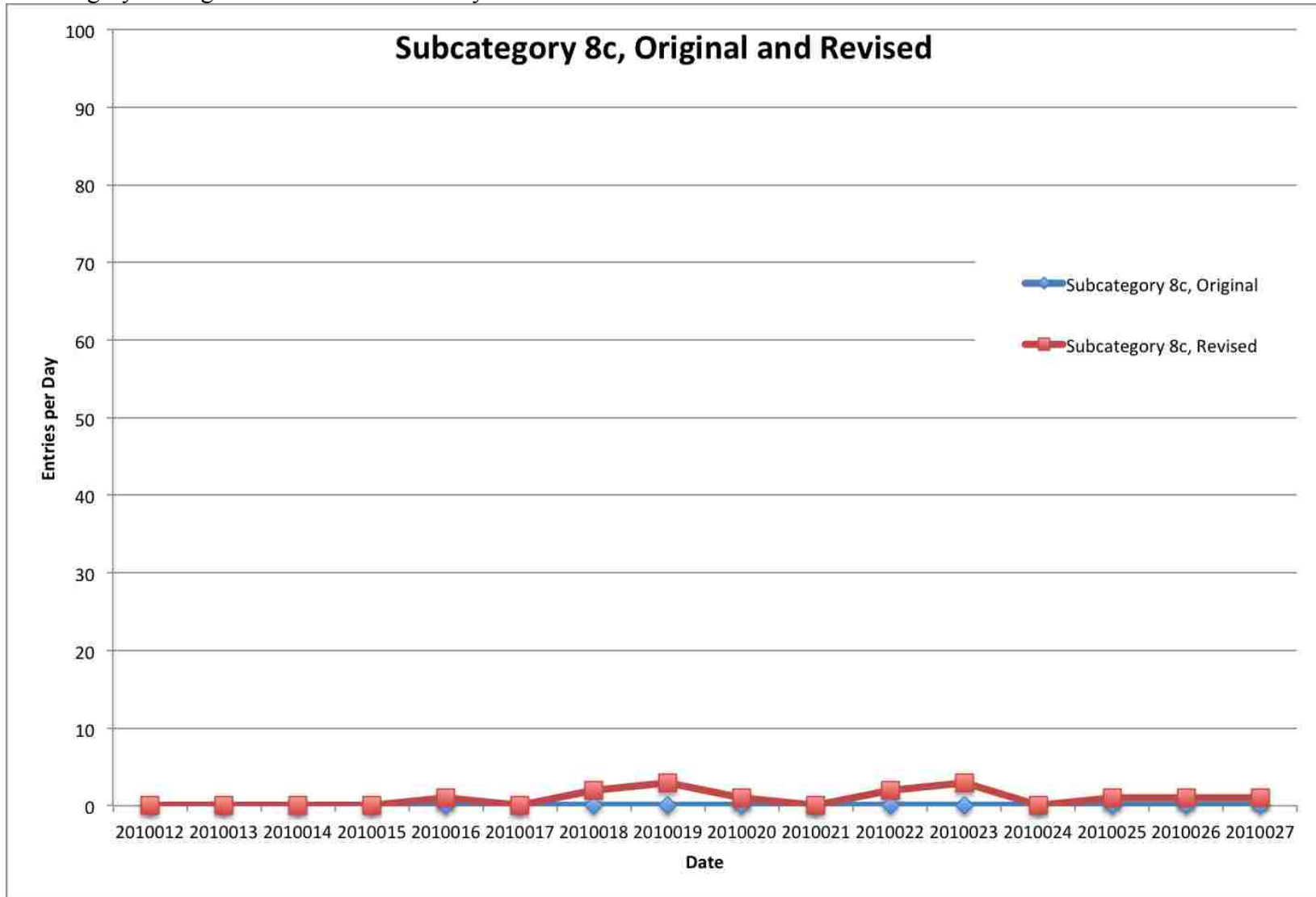
Subcategory 8a original and revised entries by date



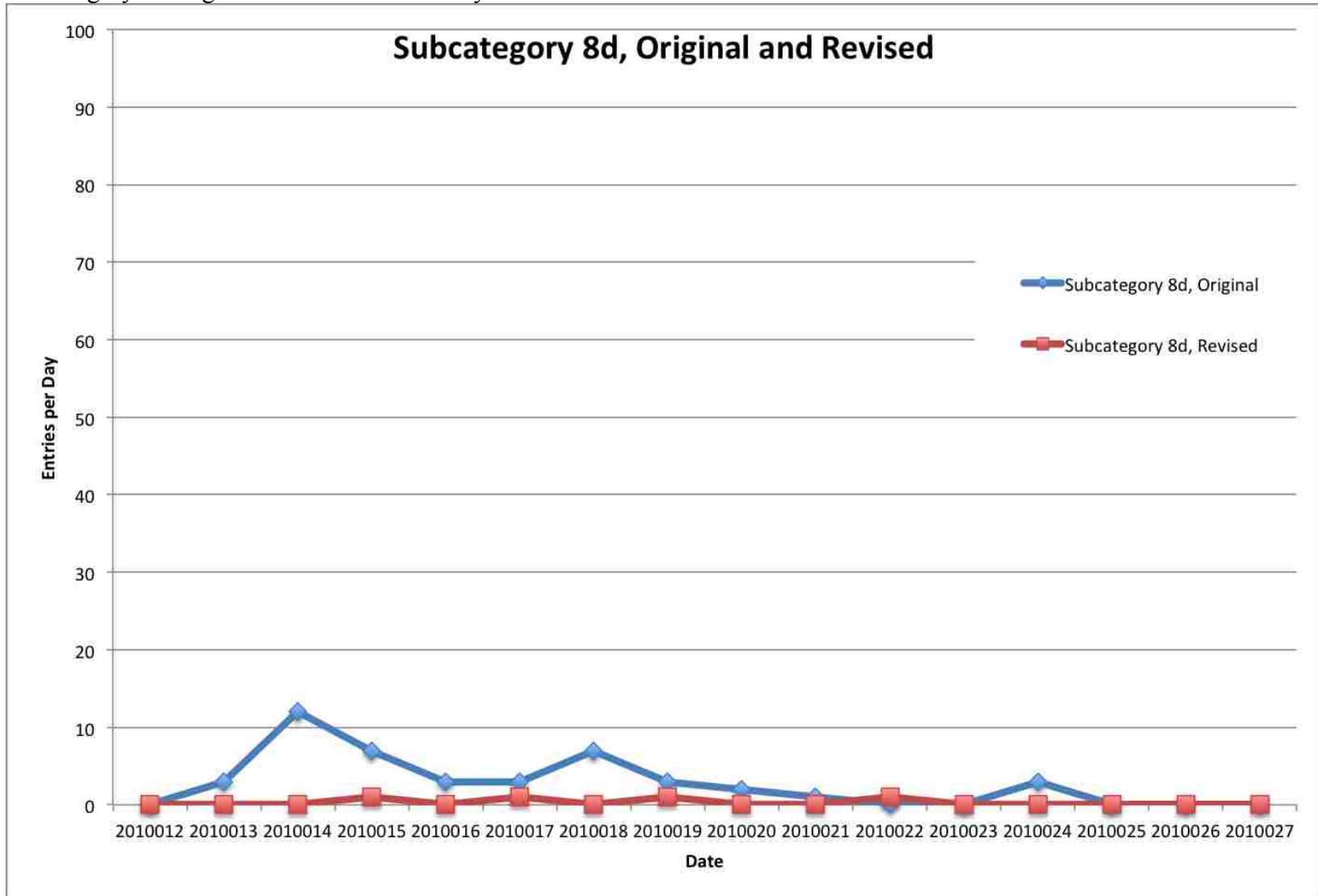
Subcategory 8b original and revised entries by date



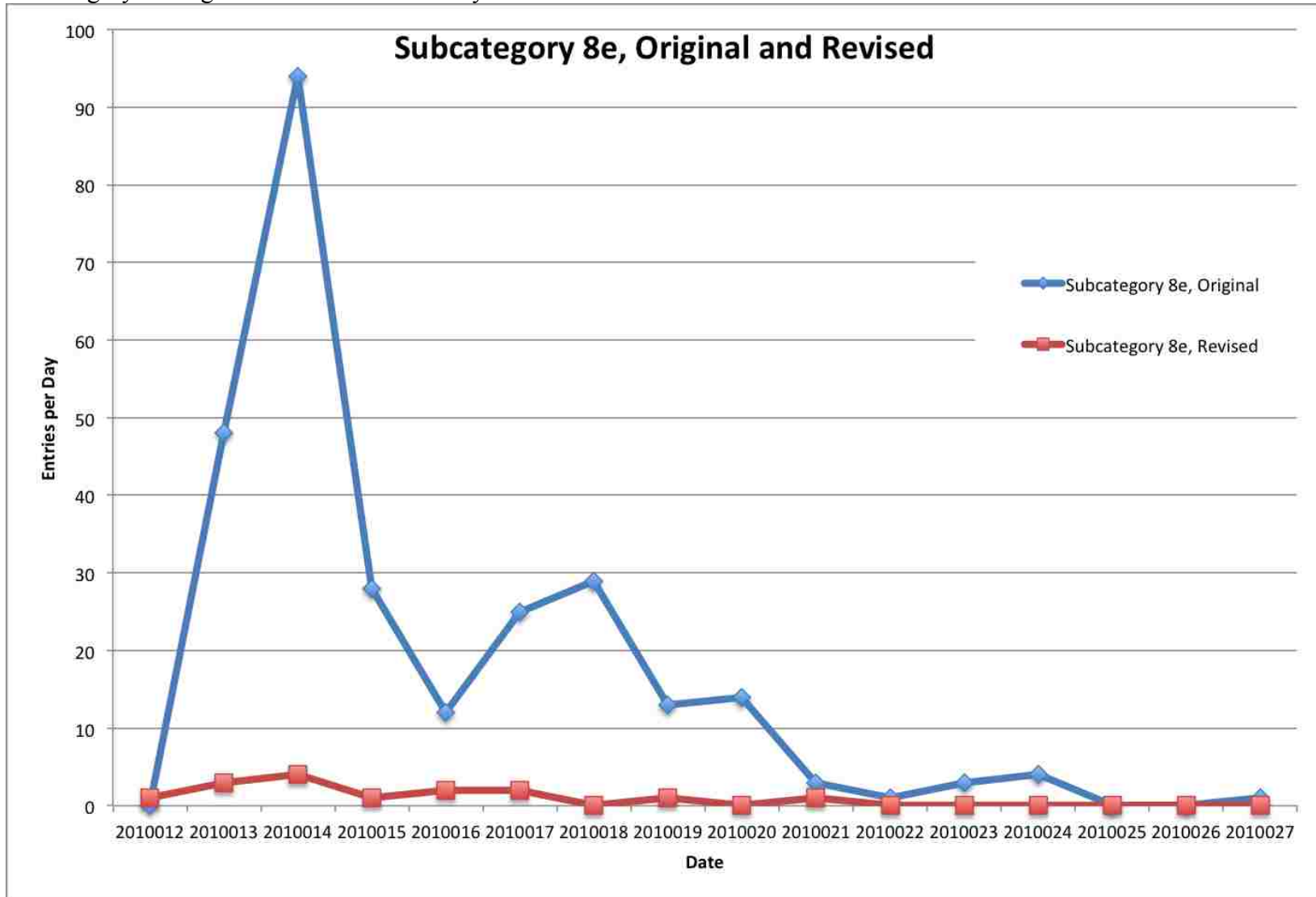
Subcategory 8c original and revised entries by date



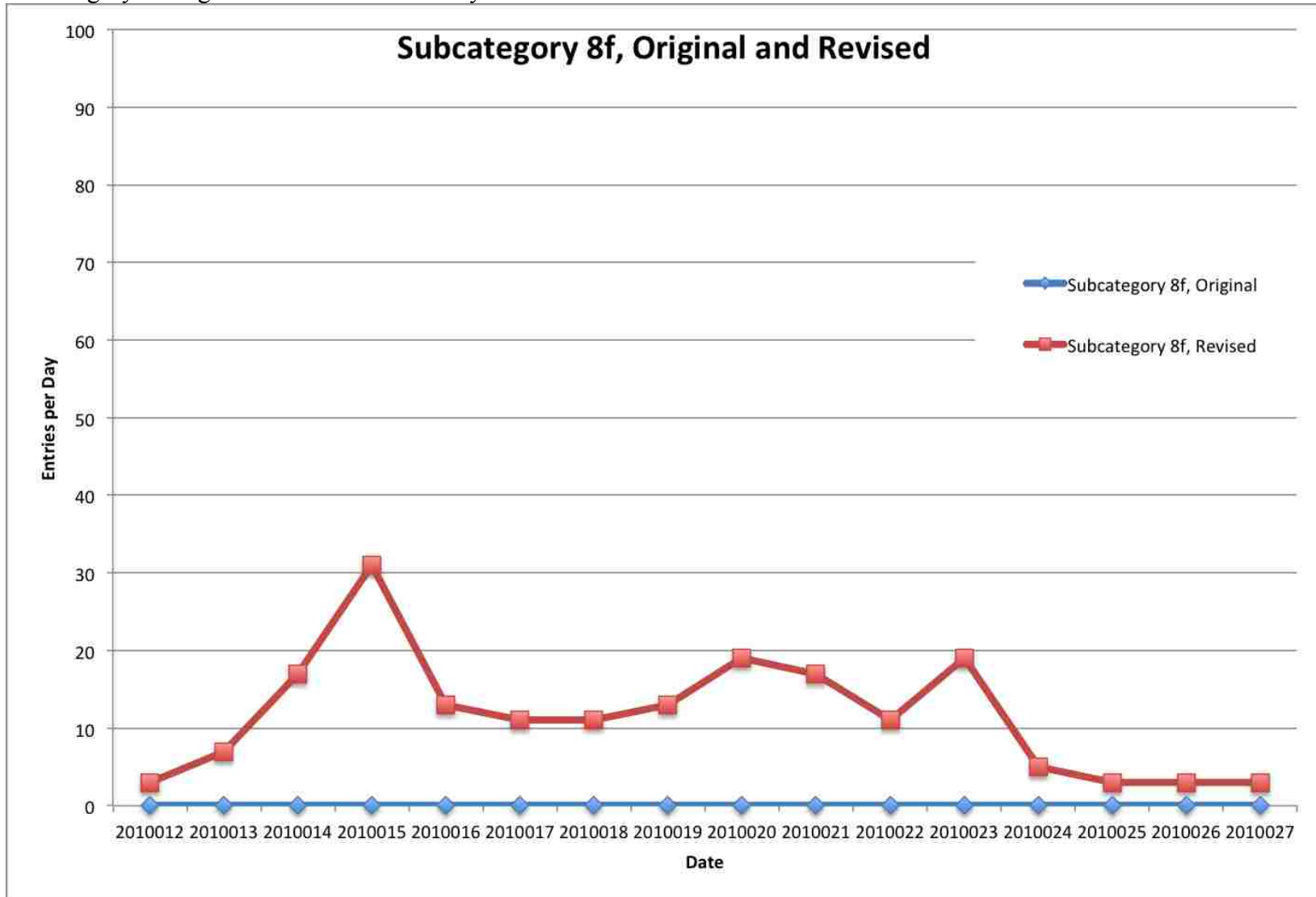
Subcategory 8d original and revised entries by date



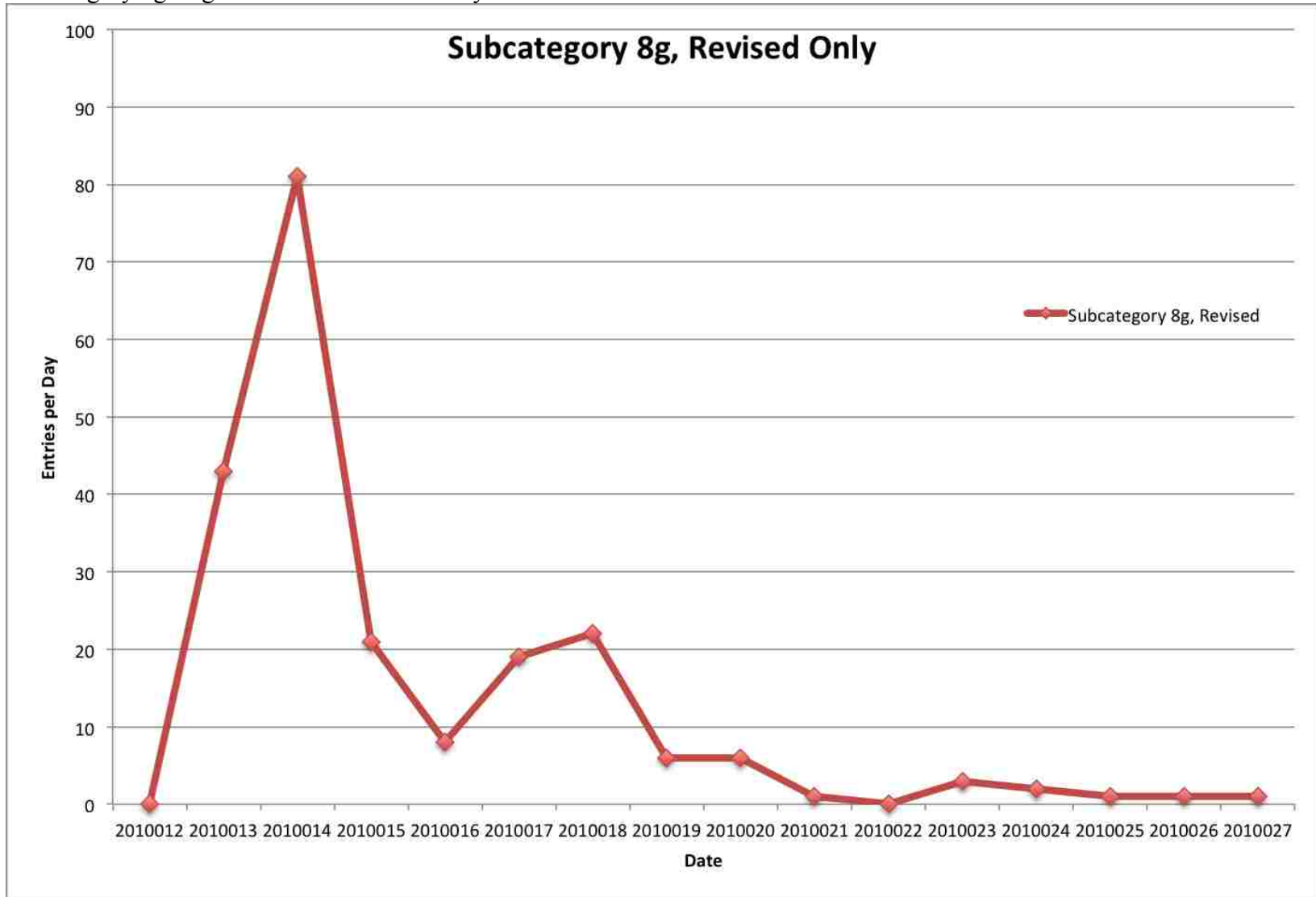
Subcategory 8e original and revised entries by date



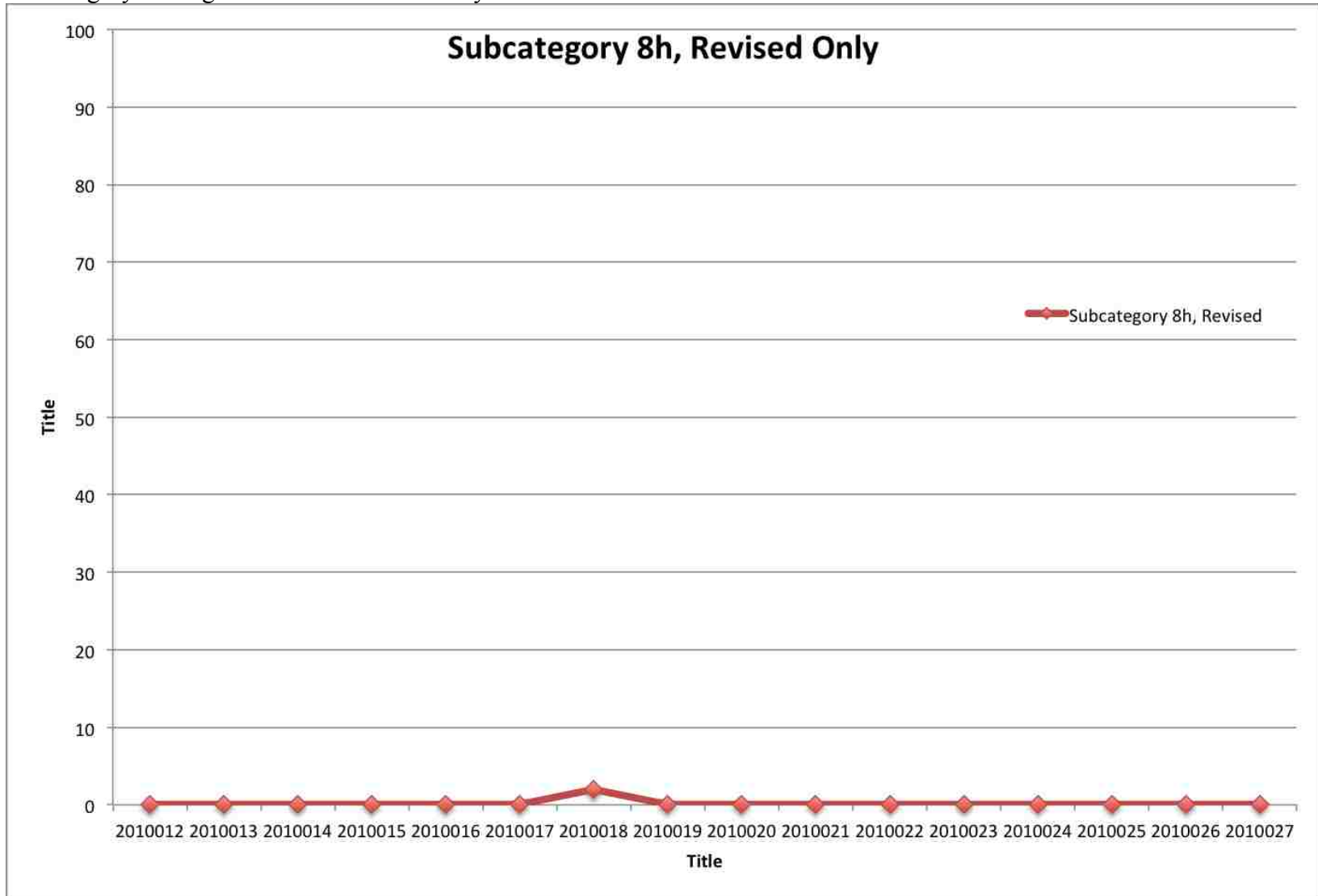
Subcategory 8f original and revised entries by date



Subcategory 8g original and revised entries by date



Subcategory 8h original and revised entries by date





## Appendix E: Haiti Map



## References

- American Red Cross. 2010. Social Media in Disasters and Emergencies. Accessed December 3, 2011. [www.americanredcross.org](http://www.americanredcross.org).
- Banzato, A., F. Barbini, A. D'Atri, E. D'Atri, and S. Za. 2010. Social Networks and Information Systems to Handle Emergency and Reconstruction in Natural Disasters: the L'Aquila Earthquake Case Study Paper read at ALPIS, at Italy.
- Batty, David. 2010. Haiti Ends Quake Rescue Operations. [www.guardian.co.uk](http://www.guardian.co.uk). Accessed on June 4, 2013. <http://www.guardian.co.uk/world/2010/jan/23/haiti-ends-quake-rescue-operations>
- Benko, Boris. 2011. Social Media and Emergency Response. [www.Mashable.com](http://www.Mashable.com). Accessed on December 3, 2011. <http://mashable.com/2011/02/11/social-media-in-emergencies/>.
- Bloch, Matthew, and Shan Carter. 2012. Twitter Chatter During the Super Bowl. The New York Times 2009 [cited April 15, 2012 2012]. Available from [http://www.nytimes.com/interactive/2009/02/02/sports/20090202\\_superbowl\\_twitter.html](http://www.nytimes.com/interactive/2009/02/02/sports/20090202_superbowl_twitter.html).
- Chrisman, N. R. 1991. The Error Component in Spatial Data. In *Geographic Information Systems*, edited by P A Longley, M F Goodchild, D J Maguire and D. W. Rhind: Longman Scientific and Technical.
- Cova, Thomas J. 1999. GIS in emergency management. In *Geographical Information Systems: Principles, Techniques, Applications, and Management*, edited by M. F. G. P.A. Longley, D.J. Maguire, D.W. Rhind. New York: John Wiley & Sons.
- Crampton, Jeremy W. 2009. Cartography: maps 2.0. *Progress in Human Geography* 33 (1):91-100.
- Currión, P, C De Silva, and B van De Walle. 2007. Open Source Software for Emergency Management. *Communications of the ACM* 50 (3):61-65.
- Cutter, Susan. 2003. GI Science, Disasters, and Emergency Management. *Transactions in GIS* 7 (4):439-445.
- Edublogs. 2011. Ushahidi-Haiti. Last modified 4 11, 2011. Accessed June 5, 2013. <http://kristenm051.edublogs.org/files/2011/04/Screen-shot-2011-04-11-at-8.42.46-AM-uyf7we.png>.
- Elwood, Sarah. 2008a. Volunteered Geographic Information: Future Research Directions motivated by critical participatory, and feminist GIS. *GeoJournal* 72:173-183.

- Elwood, Sarah. 2008b. Volunteered geographic information: key questions, concepts and methods to guide emerging research and practice. *GeoJournal* 72:133-135.
- Flanagin, Andrew J, and Miriam J Metzger. 2008. The Credibility of Volunteered Geographic Information. *GeoJournal* 72:137-148.
- Foody, G. M. 2003. Uncertainty, knowledge discovery and data mining in GIS. *Progress in Physical Geography* 27 (1):113-121.
- Frassl, M, M Lichtenstein, M Khider, and M Angermann. 2010. Developing a system for Information Management in Disaster Relief - Methodology and Requirements. Paper read at Proceedings of the 7th International ISCRAM Conference, May 2010, at Seattle, USA.
- Freifeld, Clark C., Rumi Chunara, Sumiko R. Mearu, Emily H. Chan, Taha Kass-Hout, Anahi Ayala Iacucci, and John S. Brownstein. 2010. Participatory Epidemiology: Use of Mobile Phones for Community-Based Health Reporting *PLOS Medicine* 7 (12):1-5.
- Gao, Huiji, Geoffrey Barbier, and Rebecca Goolsby. 2011. Harnessing the Crowdsourcing Power of Social Media for Disaster Relief. *IEEE Intelligent Systems* 26 (3):10-14.
- Gao, Huiji, Xufei Wang, Geoffrey Barbier, and Huan Liu. 2011. Promoting Coordination for Disaster Relief - From Crowdsourcing to Coordination. Paper read at International Conference on Social Computing, Behavioral- Cultural Modeling, and Prediction (SBP 2011), at Maryland.
- Girres, J. F., and G. Touya. 2010. Quality Assessment of the French OpenStreetMap Dataset. *Transactions in GIS* 14 (4):435-459.
- Goodchild, M. F. and G. J. Hunter. 1997. Dealing with Error in Spatial Databases: A Simple Case Study. *Photogrammetric Engineering and Remote Sensing* 61 (5):529-537.
- Goodchild, M. F., and J. A. Glennon. 2010. Crowdsourcing geographic information for disaster response: a research frontier. *International Journal of Digital Earth* 3 (3):231-241.
- Goodchild, Michael F. 2007. Citizens as sensors: the world of volunteered geography. *GeoJournal* 69 (4):211-221.
- Goodchild, Michael F. 2008. Assertion and authority: the science of user-generated geographic content. Paper read at *Colloquium for Andrew U. Frank's 60th Birthday*, at Vienna University of Technology.

- Goodchild, Michael F. 2008. Spatial Accuracy 2.0. Paper read at 8th international symposium on spatial accuracy assessment in natural resources and environmental sciences, at Shanghai.
- Gutnick, Aviva Lucas, Michael Robb, Lori Takeuchi, and Jennifer Kotler. 2011. *Always Connected: The new digital media habits of young children*. New York: The Joan Ganz Cooney Center at Sesame Workshop.
- Haklay, M. 2010. How good is volunteered geographical information? A comparative study of OpenStreetMap and Ordnance Survey datasets. *Environment and Planning B-Planning & Design* 37 (4):682-703.
- Haklay, M., S. Basiouka, V. Antoniou, and A. Ather. 2010. How Many Volunteers Does it Take to Map an Area Well? The Validity of Linus' Law to Volunteered Geographic Information. *Cartographic Journal* 47 (4):315-322.
- Haklay, Mordechai, and Claire Ellul. 2010. Completeness in volunteered geographical information the evolution of OpenStreetMap coverage in England (2008--2009) *Journal Of Spatial Information Science*.
- Heinzelman, Jessica, D. Roz Sewell, Jen Ziemke, and Patrick Meier. 2011. *Lessons from Haiti and Beyond : Report from the 2010 International Conference on Crisis*. edited by U. S. I. o. Peace. Washington D.C.
- HOT. Haiti. Accessed June 30, 2013. <http://hot.openstreetmap.org/projects/haiti-2>.
- Howe, Jeff. 2008. *Crowdsourcing: Why the power of the crowd is driving the future of business*. New York: Three Rivers Press.
- Kaiser, R, P Spiegel, A Henderson, and M Gerber. 2003. The Application of Geographic Information Systems and Global Positioning Systems in Humanitarian Emergencies: Lessons Learned, Programme Implications and Future Research, Disasters. *Disasters* 27 (2):127-140.
- Keller, Edward A., and Robert H. Blodgett. 2006. *Natural Hazards*. New Jersey: Pearson Prentice Hall.
- Kelmelis, J, L Schwartz, C Christian, M Crawford, and D King. 2006. Use of Geographic Information in Response to the Sumatra-Andaman Earthquake and Indian Ocean Tsunami of December 26, 2004. *Photogrammetric Engineering and Remote Sensing* 72:862-876.
- LaFranchi, Howard. 2010. The Christian Science Monitor, "Haiti earthquake: How a top UN official was plucked from the rubble." Last modified 01 26, 2010. Accessed June 4, 2013. <http://www.csmonitor.com/World/Americas/2010/0126/Haiti-earthquake-How-a-top-UN-official-was-plucked-from-the-rubble>.

- Lee, B. 2009. Spatial patterns of uncertainties: An accuracy assessment of the TIGER files. *Journal of Geography and Geology* 1(2): 2–12
- Maliszewski, Paul, and Mark Horner. 2010. A spatial modeling framework for siting critical supply infrastructures. *The Professional Geographer* 62 (3):426-441.
- Meier, Patrick. 2012. How Crisis Mapping Saved Lives in Haiti. Last modified July 2, 2012. <http://newswatch.nationalgeographic.com/2012/07/02/crisis-mapping-haiti/>. Last Accessed June 4, 2013.
- Morrow, Nathan, Nancy Mock, Adam Papendieck, and Nicholas Kocmich. 2011. Independent Evaluation of the Ushahidi Haiti Project. <https://sites.google.com/site/haitiushahidieval/news/finalreportindependentevaluationoftheushahidihaitiproject>. Last Accessed June 4, 2013.
- Munro, Robert. 2010. Crowdsourced translation for emergency response in Haiti: the global collaboration of local knowledge. Paper read at AMTA Workshop on Collaborative Crowdsourcing for Translation, at Denver, Colorado.
- National Research Council. 2007a. Successful Response Starts with a Map: Improving Geospatial Support for Disaster Management. Washington D.C.: National Academies Press.
- National Research Council. 2007b. Tools and Methods for Estimating Populations at Risk from Natural Disasters and Complex Humanitarian Crises. Washington D.C.: National Academies Press. <http://www.nap.edu/catalog/11895.html>.
- Nelson, Anne, Ivan Sigal, and Dean Zambrano. 2010. Media, Information Systems and Communities: Lessons from Haiti. Knight Foundation.
- Nguyen, Hong Phuong. 2005. Development of a Decision Support System for Earthquake Risk Assessment and Loss Mitigation: The Hanoi Case Study. *International Journal of Geoinformatics* 1 (1):191-196.
- Nielsen. 2012. Smartphones account for half of all mobile phones, Dominate new phone purchases in the US. In Nielsen Wire. [blog.nielsen.com/nielsenwire/](http://blog.nielsen.com/nielsenwire/): Nielsen.
- Norheim-Hagtun, Ida, and Patrick Meier. 2010. Crowdsourcing for Crisis Mapping in Haiti. *Innovations* 5 (4):81-89.
- O'Reilly, Tim. 2005. What is Web 2.0: Design Patterns and Business Models for the Next Generation of Software. Accessed June 30, 2013. <http://oreilly.com/web2/archive/what-is-web20.html>
- Okolloh, O. 2008. Ushahidi, or “testimony”: Web 2.0 tools for crowdsourcing crisis information. *Participatory Learning and Action* 59:65-70.

- Ortmann, Jens, Minu Limbu, Dong Wang, and Tom Kauppinen. 2011. Crowdsourcing Linked Open Data for Disaster Management. *Framework*:12.
- Piotrowski, Chris. 2010. EARTHQUAKE IN HAITI: The Failure of Crisis Management. *Organization Development Journal* 28 (1):107-112.
- Pitzer, William. 2011. Media and Communications Ecosystem. Accessed December 3, 2011. <https://knight.box.net/shared/static/y4m3iq9a9k43kyfkelm8.pdf>.
- Radke, John, Tom Cova, Micahel Sheridan, Austin Troy, Mu Lan, and Russ Johnson. 2000. Application Challenges for Geographic Information Science: Implications for Research, Education, and Policy for Emergency Preparedness and Response. *URISA Journal* 12 (2):15-30.
- Standby Volunteer Task Force. 2011. Libya Crisis Map Deployment. Accessed December 3, 2011. <http://blog.standbytaskforce.com/libya-crisis-map-report/>.
- Sui, D. Z. 2004. Tobler's first law of geography: A big idea for a small world? *Annals of the Association of American Geographers* 94 (2):269-277.
- Tveitte, H. and Langass S. 1999. An accuracy assessment method for geographical line data sets based on buffering. *International Journal of Geographical Information Science*. 13: 27 – 47.
- Tobler, Waldo. 1970. A computer movie simulating urban growth in the Detroit region. *Economic Geography* 46 (2):234-240.
- Tomaszewski, B. 2011. Situation awareness and virtual globes: Applications for disaster management. *Computers and Geosciences* 37:86-92.
- Turner, A. 2006. Introduction to neogeography. Sebastopol, CA: O'Reilly.
- Ushahidi. 2011. Ushahidi-Haiti." Last modified 9 7, 2011. Last Accessed March 2012. <http://haiti.ushahidi.com>.
- Van Niel T. and T McVicar. 2002. Experimental evaluation of positional accuracy estimates from a linear network using point and line based testing methods. *International journal of Geographical Information Science*. 16: 455-73.
- Van Oort, P.A.J. 2006. Spatial data quality: from description to application. PhD Thesis, Wageningen: Wagenigen Universiteit. 125.
- Van Westen, C.J., and P.Y. Georgiadou. 2001. Spatial data requirements and infrastructure for geologic risk assessment. Paper read at workshop on natural disaster management, ISPRS Technical Committee VII, at Ahmedabad, India.

- Yates, Dave, and Scott Paquette. 2011. Emergency knowledge management and social media technologies: A case study of the 2010 Haitian earthquake. *International Journal of Information Management* 31 (1):6-13.
- Zerger, A, and A Smith. 2003. Impediments to Using GIS for real-time disaster decision support. *Computers, Environment, and Urban Systems* 27:123-141.
- Zheng, Li, Chao Shen, Liang Tang, Tao Li, Steve Luis, Shu-Ching Chen, and Vagelis Hristidis. 2010. Using data mining techniques to address critical information exchange needs in disaster affected public-private networks. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*. Washington, DC, USA: ACM.
- Zielstra, Dennis, and Alexander Zipf. 2010. A Comparative Study of Proprietary Geodata and Volunteered Geographic Information for Germany In *13th AGILE International Conference on Geographic Information Science* Portugal.