

Bericht der Gruppe E

Identifikation von Nutzergruppen - Reddit \r/place

Christian Funk, Sebastian Heuer, Greta Wetzel, Nicole Derr, Jannik Littmann

Martin-Luther-Universität Halle-Wittenberg

christian.funk@student.uni-halle.de, sebastian.heuer@student.uni-halle.de, greta.wetzel@student.uni-halle.de,
nicole.derr@student.uni-halle.de, jannik.littmann@student.uni-halle.de

Figure 1: Some large rubber duck.

ABSTRACT

Some general hints on what to mention in an abstract: What are the questions you address? Why are they interesting? What approaches did you use? What answers did you find?

As for how to structure the abstract: Give some motivation and context on the general topic you address (1–2 sentences). Then state the specific questions you address (1–2 sentences) and describe how you approach them (2–3 sentences). Finally, results and some conclusion (1–3 sentences).

KEYWORDS

Übung „Big Data Analytics“, Sommersemester 20XX

1 INTRODUCTION

Jedes Jahr zum 1. April führt Reddit ein soziales Experiment durch. Im Jahr 2017 wurde dabei zum ersten Mal r/place vollkommen unangekündigt durchgeführt. Fünf Jahre später wurde das gleiche Experiment wiederholt, jedoch mit Vorankündigung. Bei r/place können Nutzer alle 5 Minuten einen Pixel auf einem Canvas setzen. In den Jahren 2017 und 2022 war der Canvas 1000 x 1000 Pixel bzw. 2000 x 2000 Pixel groß. Das Experiment dauerte 4 Tage.

Nach erster Verwirrung der Nutzer zu Beginn des Experiments von 2017, bildeten sich einzelne Communities, welche gemeinsam Projekte verfolgten. Es entstanden erste Artworks und kreative Explosionen bis hin zu Kämpfen um bestimmte Regionen auf dem Canvas.

Vor allem bei der Wiederholung 2022 kamen vermehrt Scripts für Bots zum Einsatz, was vor allem der Vorankündigung geschuldet war, sodass Nutzer die Möglichkeit bekamen sich gezielt vorzubereiten. Anders als 2017, wurde die Wiederholung 2022 von Moderatoren begleitet, welche unerwünschte Artworks zensurierten.

Da sich viele Communities vor allem 2022 zusammenschlossen, wurde die Frage, "ob es möglich ist, Nutzer anhand von ähnlicher Verhaltensweisen automatisch Nutzergruppen zuzuordnen?" zum Ziel dieses Projekts.

Weiterhin wurde auch betrachtet, ob Bots und Moderatoren identifizierbar sind und ob entsprechende Statistiken zu diesen entwickelt und ausgewertet werden können.

2 DATA

Der offizielle Datensatz von r/place aus der Wiederholung 2022 beinhaltet 72 Millionen Pixel, von 6 Millionen Nutzern und besitzt eine Größe von ca. 22GB. Dieser ist in einer CSV-Datei strukturiert, welche den Zeitpunkt des gesetzten Pixels (timestamp), eine gehashte Nutzeridentifikation (*userid*), die verwendete Farbe

(*pixel_color*) und die x- und y-Koordinaten auf dem Canvas (coordinates) beinhaltet. Eine Besonderheit stellen die Moderatoren im Datensatz dar. Diese besitzen in der Spalte coordinates nicht nur eine x- und y-Koordinate, wie alle anderen Nutzer, sondern zwei x- und y-Koordinaten. Das erste Koordinatenpaar bildet dabei die obere linke Ecke des Zensur-Quadrates eines Moderators und das zweite Koordinatenpaar die rechte untere Ecke.

Zu finden sind die Daten online unter: https://www.reddit.com/r/place/comments/txvk2d/rplace_datasets_april_fools_2022/ und https://www.reddit.com/r/redditdata/comments/6640ru/place_datasets_april_fools_2017/ für 2017.

Der Datensatz der ersten Durchführung 2017 ist mit einer Größe von ca. 1GB um einiges kleiner als der von 2022. Dieser hat eine ähnliche Gliederung wie der von 2022, jedoch wurden hier die x- und y-Koordinaten direkt getrennt aufgelistet statt unter einer gemeinsamen Spalte, wie in 2022. Zudem ist hier die Restriktion zu beachten, dass Nutzer nur zwischen 16 Farben wählen konnten und diese im Datensatz mit Integer-Values zwischen 0 bis 15 codiert wurden. Unter dem angegebenen Link für den Datensatz von 2017 findet sich eine entsprechende Tabelle, die wiedergibt, welcher Integer für welchen Farb-Hashwert steht.

Nachdem die Datensätze heruntergeladen waren, mussten diese zunächst minimal bereinigt bzw. angepasst werden. Im Datensatz für 2017 mussten die Integer-Values der Farben in den jeweiligen Hashcode umgewandelt werden und für jeweils beide Datensätze wurden die Werte der timestamp-Spalte in nutzbare Werte umgewandelt. Dafür wurde das UTC-Timestampformat zunächst in das Unix-Epoch-Seconds-Format umgewandelt und anschließend der kleinste Wert von allen timestamps abgezogen. Damit waren die timestamps so normalisiert, dass der erste Eintrag mit einem timestamp gleich null beginnt. Eine weitere Bereinigung der Daten beinhaltet, dass Moderatoren und "normale" Nutzer gespalten wurden, sodass diese in separaten DataFrames betrachtet werden können. Im Datensatz von 2022 wurde damit die Spalte coordinates jeweils in die x- und y-Koordinaten für "normale" Nutzer aufgespalten und für Moderatoren in x1-, y1-, x2- und y2-Koordinaten. Eine solche Spaltung musste für den 2017-Datensatz entsprechend nicht durchgeführt werden. - Beispieldaten

3 IDENTIFIKATION VON NUTZERGRUPPEN

Fragestellung 1: Ist es möglich anhand von ähnlichen Verhaltensweisen Nutzer ihren Nutzergruppen zuzuordnen?

26.04.

Betrachtung verschiedener *UserIDs*, die mehr als x-Mal zur selben Zeit und recht nah beieinander (Koordinaten) Pixel platziert/verändert haben

Zusätzlich anhand der gewählten Farben der User bzw. Nutzergruppen betrachten, ob zwei oder mehr Gruppen miteinander

konkurrierten

03.05.

Definition Nutzergruppe:

- Nutzer ist zu verschiedenen Zeitintervallen aktiv
- Nutzer kann mehrere Interessen vertreten
- In jedem aktiven Zeitintervall vertritt ein Nutzer pro Raumzeit-gebiet aber nur genau ein Interesse
- Überschneiden sich Interessen mehrerer Nutzer räumlich und zeitlich, so stehen diese im Konflikt zueinander
- Beispiel Deutschland - Frankreich Flaggen

10.05.

- Aggregation Nutzerdaten
- Raumzeitmetrik (SVM)
- 3-dimensionalität der Daten, x-, y-Koordinate, Zeit
- Zusammenfassen von Raumzeiten (Bounding Boxes)
- Infrastruktur?

17.05.

- Bilder Bounding Boxes und zusammenhängende Pixel (Zusammenhangskomponenten) mit Beschreibung

14.06. - Methodik weiteres Vorgehen

Abschluss

Was ist unser Abschluss, Ergebnis?

4 BOTS UND MODERATOREN

Statistische Auswertungen des Datensatzes:

In diesem Abschnitt soll nun auf die Auswertung hinsichtlich der Bots und Moderatoren eingegangen werden.

Ist es möglich Bots und Moderatoren anhand des Datensatzes zu identifizieren?

Können wir dazu entsprechende Statistiken entwickeln und auswerten?

17.05.

Vermutungen:

- viele Nutzer setzen nur wenige Pixel, große Häufigkeit bei kleiner Pixelzahl
- Bots setzen konstant viele Pixel, leicht erhöhte Häufigkeit bei großer Pixelzahl
- Testdatensatz umfasst eine Stunde, maximal 12 Pixel pro Nutzer möglich

Vorgehen:

- Nach Nutzern aggregierte Pixeldaten nochmal nach Pixelanzahl aggregieren
- Grafik

30.05.

Statistiken:

- meiste umkämpfter Pixel: **Grafik** ev. in 3-Dimensionen
- meiste verwendete Farben: **Grafik**

21.06.

- Filterkriterien Bots: ca. 14.000 Bots im Testdatensatz
- Filter Moderatoren, Welche Farben wurden verwendet?

28.06.

- Grafik Bots 2022

05.07.

- Präzisierung Filter Bots
- Bearbeitung Datensatz 2017, **Grafik Bots 2017**
- Grafik Mods
- Fragen zur Rekonstruktion

4.1 Moderatoren

Hier Input von Nicole! - Die Struktur der CSV-Datei erleichterte die Identifikation von Moderatoren im Datensatz 2022

- Der Datensatz von 2017 blieb dabei unbetrachtet, da hier unklar ist, ob Moderatoren aktiv waren und wenn ja, ob diese mit in den Datensatz aufgenommen wurden

- Wie bereits eingangs erwähnt, besitzen Moderatoren in der Spalte coordinates vier Werte, welche die x- und y-Koordinate der linken oberen Ecke und der rechten unteren Ecke des Zensur-Quadrates darstellen

- Anhand dieser Struktur wurde für die Moderatoren ein von den "normalen" Nutzern getrenntes DataFrame generiert, sodass keine Identifikation in engeren Sinne durchgeführt werden musste

- Zunächst legte sich der Fokus auf einfach auszuwertende Statistiken:

- Während des gesamten Experiments 2022 waren 16 Moderatoren aktiv, von denen jeder im Durchschnitt ca. 6262 Pixel zensierte und insgesamt 100.197 Pixel zensiert wurden.

- Im Vergleich dazu wurden 159.024.375 Pixel von "normalen" Nutzer gesetzt.

- Teilt man nun die Anzahl der Zensuren durch die Anzahl an gesetzten Pixeln, so erhält man mit 0,063% das Verhältnis, in welchem zensiert wurde.

- Aus diesem prozentualen Anteil lässt sich schließen, dass Moderatoren kaum Pixel zensierten während des gesamten Verlaufes.

- Das bedeutet wiederum, dass entweder kaum unerwünschte Artworks von Nutzern gezeichnet wurden oder Artworks so schnell bezeichnet wurden, dass die Moderatoren kaum aktiv werden mussten.

- HIER BIDL EINFÄGEN MIT DEM CANVAS WO DIE RAHMEN DER ZENSUR-QUADRATE ZU SEHEN SIND

- Die Abbildung Nr. 7 zeigt ein Canvas auf dem die Rahmen der Zensur-Quadrat der Moderatoren nachgezeichnet wurden.

- Damit erhält man einen Überblick darüber, in welchen Regionen des gesamten Canvas Moderatoren überhaupt aktiv werden mussten

- Die Abbildung verdeutlicht noch einmal, dass Moderatoren kaum aktiv werden mussten und das nur in der unteren Region zensiert werden musste

- Da das Canvas während des gesamten Verlaufes 2022 zwei mal verdoppelt wurde und der untere Teil dabei als letztes hinzukam, lässt sich daraus schließen, dass Moderatoren erst zu einem späteren Zeitpunkt des Experimentes aktiv werden mussten

- Eine aufwendigere Statistik für die Moderatoren, die betrachtet werden sollte:

- Wie viele Pixel wurden von den Bots bzw. "normalen" Nutzern im Verhältnis zensiert? Lässt sich daraus ableiten, dass eher von Bots oder von Nutzern gesetzte Pixel zensiert werden mussten?

- Um dies zu beantworten, musste zunächst eine Funktion geschrieben werden, welche rückwirkend die Pixel identifiziert, die mit dem jeweiligen Zensur-Quadrat eines Moderators zensiert wurden.

- Beschreibung wie die Funktion vorgeht:

- Vor dem Funktionsaufruf werden zunächst alle Pixel bestimmt, die von den Zensuren der Moderatoren betroffen sind

- Dies dient als einer der Inputs der Funktion

- Den zweiten Input bietet ein DataFrame mit allen gesetzten

Table 1: Some example table.

Some entries	Some numbers
Entry A	400 million
Entry B	300 million
Entry C	200 million

Pixeln im Verlaufe des Experimentes

- Dafür wird innerhalb der Funktion jeder Pixel schrittweise betrachtet, der von einer Zensur betroffen ist
- Mittels einer For-Schleife wird dann von dem Zeitpunkt der Zensur aus runtergezählt, bis der Zeitpunkt 0 erreicht wurde oder der zensierte Pixel gefunden wurde
- wurde ein zensierter Pixel identifiziert, wird dieser in ein neues DataFrame geschrieben, welches alle zensierten Pixel mit ihren Daten beinhaltet
- Am Ende liefert die Funktion das entsprechende DataFrame mit allen zensierten Pixel und deren Daten, wie UserID, x- und y-Koordinate, etc.
- Ein Zwischenergebnis dieser Funktion ist, dass die zensierten Pixel nachgezeichnet werden können.
- Eine Erwartung war, dass die zensierten Artworks evtl "interessante" Themen darstellen, wie z.B. politisch Themen die unerwünscht waren
- Jedoch wurde diese Erwartung nicht erfüllt, da es sich bei den zensierten Pixeln nur um unangemessene, kleinere Zeichnungen handelte, die in diesem Rahmen auch nicht weiter betrachtet werden sollen.
- Aufgrund von Problemen mit der Infrastruktur, war es nicht möglich den gesamten Datensatz auf einmal auszuwerten, so dass die Frage "Wie viele Pixel wurden von den Bots bzw. "normalen" Nutzern im Verhältnis zensiert? Lässt sich daraus ableiten, dass eher von Bots oder von Nutzern gesetzte Pixel zensiert werden mussten?" entsprechend nicht beantwortet werden konnte
- (Ein Problem ist, dass die Datensätze voneinander für diese Auswertung zu stark zusammenhängen, als das man diese hätte einzeln auswerten und zu einem Gesamtergebnis zusammentragen können) -> entsprechend im Fazit weiter drauf eingehen bzw. sagen dass es einen möglichen Ausblick bieten könnte!
- damit bleibt diese Frage weiterhin offen

Vergleiche

- Statistiken 2022 und 2017 im Vergleich - Grafiken
- Vergleich Arbeit und Menge Bots
- Vergleich der Zensur von Mods

5 EVALUATION

Some evaluation section if appropriate. You might want to refer to some table with results in this section (e.g., to Table ??).

6 CONCLUSION

The introduction in less detail. Summarize story in retrospective, give outlook on possible next steps. Semi-technical ...