

Bericht über die Datenanalyse des SHP-Datensatzes

Dimitriev, Erismann, Hosseini

2022-04-28

Dokumentation Projekt 2

Einleitung

Dieser Bericht behandelt folgende Fragestellung:

“Welche Faktoren stehen mit einer Veränderung der allgemeinen Lebenszufriedenheit seit Beginn der Verbreitung von Covid-19 in Beziehung?”

Die Zielgrösse '*Änderung der allgemeinen Lebenszufriedenheit von 2019 auf 2020*' soll untersucht werden. Dazu kommen Variablen, die mögliche Einflussfaktoren darstellen. Des weiteren bestimmten wir drei zusätzliche Variablen, die möglicherweise einen Einfluss auf die Zielvariable haben. Die verwendeten Daten stammen vom Swiss Household Panel (SHP). Es sind nur volljährige Personen berücksichtigt worden.

Organisation

| Woche | Auftrag | Zuständigkeit |
|------------------|--|---|
| 28.03 - 01.04 | <ul style="list-style-type: none">• Einarbeiten in den Auftrag• Datenstruktur planen und festlegen: git.zhaw• Daten herunterladen• Einlesen der ersten Daten in R• Uns mit den erhaltenen Daten vertraut machen• Definieren der eigenen Variablen• Erstellen der Dokumentation | Team Patrick Ahmad Patrick Team Team Team |
| 04.04 - 08.04 | <ul style="list-style-type: none">• Alle Daten einlesen• Daten gemäss Vorgaben aufbereiten• Vorgegebene Variablen einlesen und wenn nötig definieren• Eigene Variablen erstellen | Patrick Ahmad Team Patrick |
| 11.04 - 15.04 | <ul style="list-style-type: none">• Re-Check der definierten Variablen• Prüfen wo es NA's hat, nötigenfalls Annahmen treffen• Abgleich der übriggebliebenen Daten mit Daten vom BFS | Team Patrick Stojche |

| Woche | Auftrag | Zuständigkeit |
|------------------|---|-----------------------------|
| 18.04 - 22.04 | <ul style="list-style-type: none"> • EXPD • Auswertung der Daten • R-Markdown erstellen für Abgabe | Ahmad Stojche Patrick |

Variablen

Vordefinierte Variablen:

- 1 Änderung der allgemeinen Lebenszufriedenheit von 2019 auf 2020 (schlechter, gleich, besser)
- 2 Geschlecht
- 3 Alter
- 4 Kanton
- 5 Gemeindetyp
- 6 Höchstes Ausbildungszertifikat
- 7 Beziehungsstatus (in Beziehung, nicht in Beziehung)
- 8 Änderung des Beziehungsstatus von 2019 auf 2020 (Trennung, keine Änderung, Neu in Beziehung)
- 9 Kurzfristige Änderung der Arbeitssituation nach Pandemiebeginn (ja = wurde Arbeitslos, Kurzarbeit, arbeitet weniger z.B. wegen Kinderbetreuung, eigenes Geschäft ist direkt von Pandemie betroffen, sonst nein)
- 10 Corona Infektion (nein, ja)
- 11 Haushaltstyp
- 12 Jährliches Netto-Haushaltseinkommen pro Haushaltsmitglied
- 13 Finanzielle Probleme in der Jugend
- 14 War Person mindestens einmal Arbeitslos zwischen 2015 und 2019?

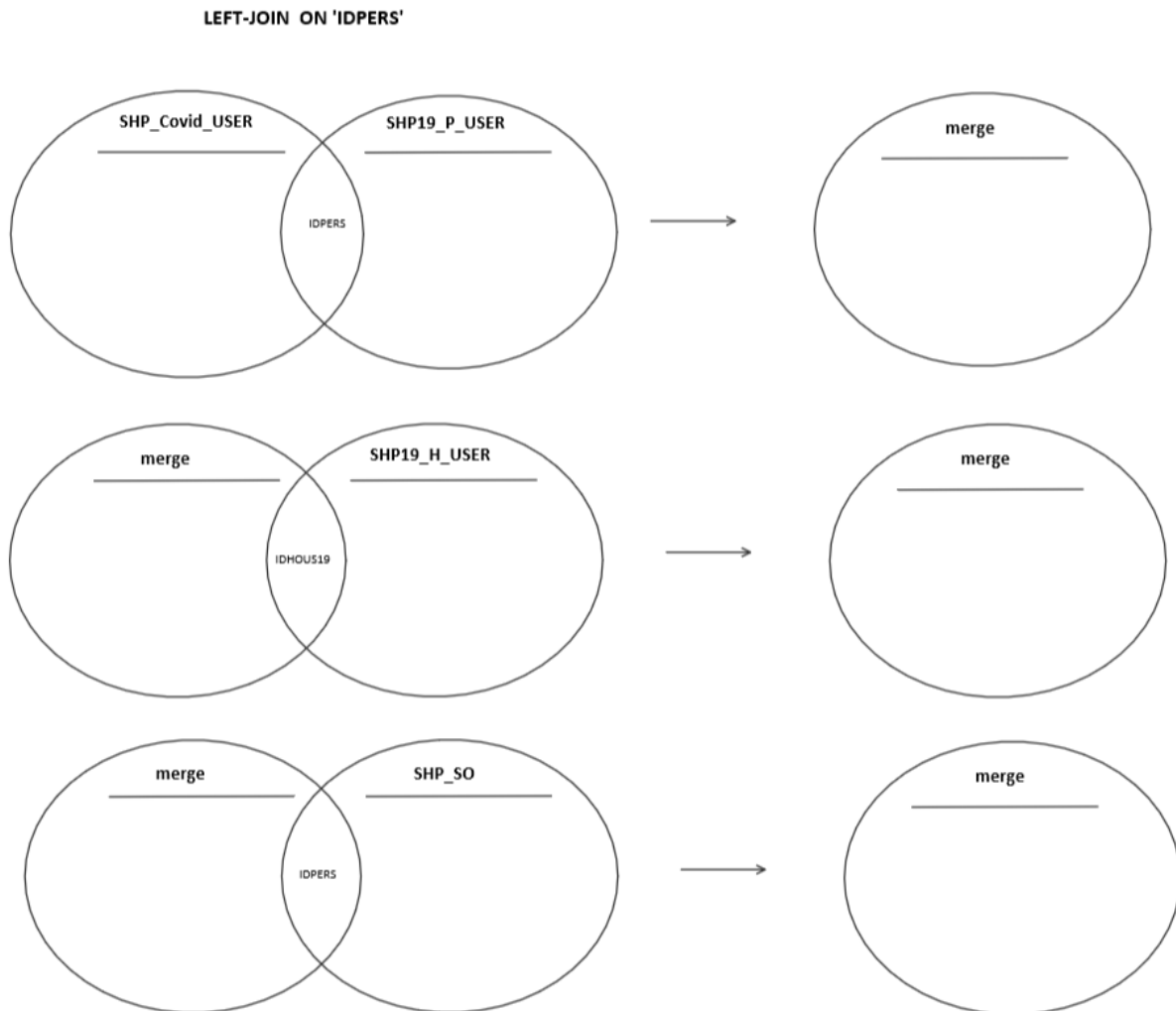
Selbstdefinierte Einflussfaktoren

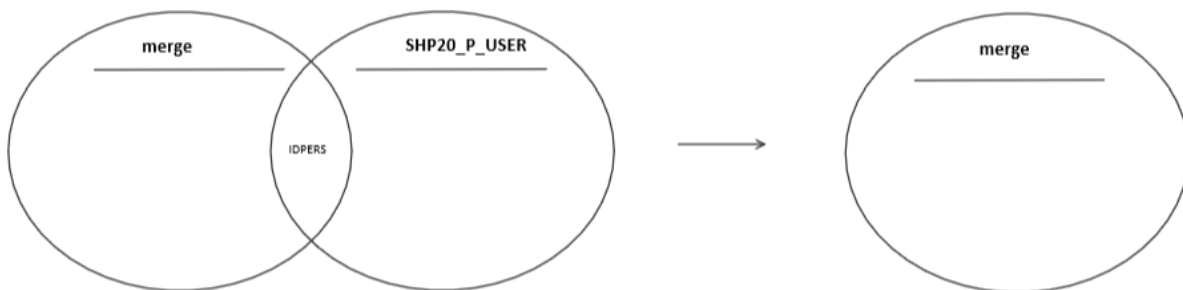
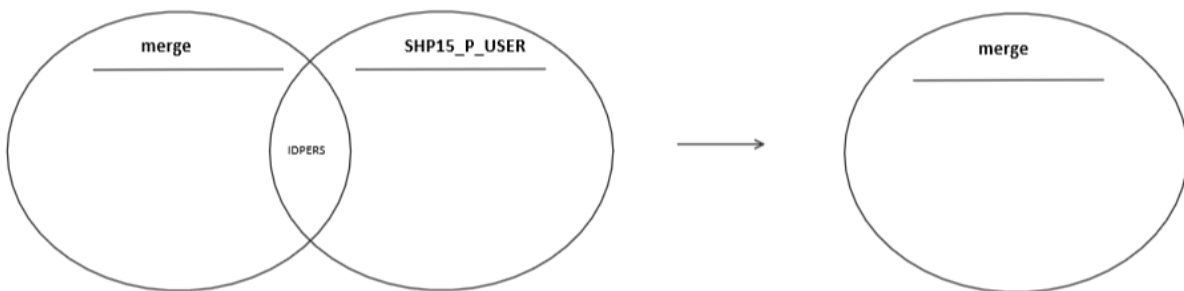
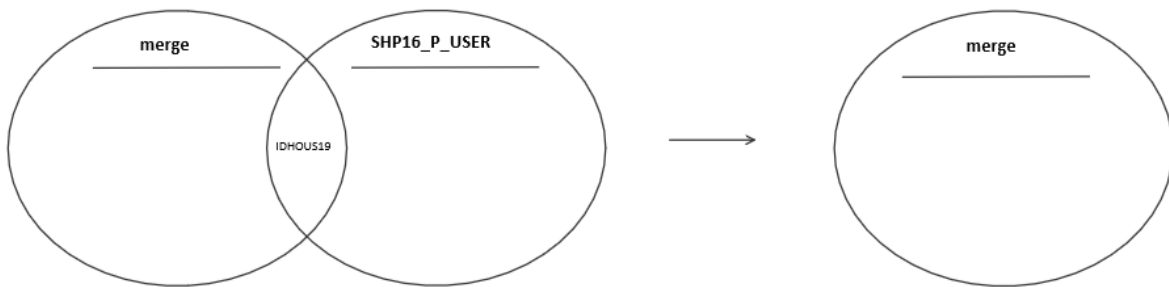
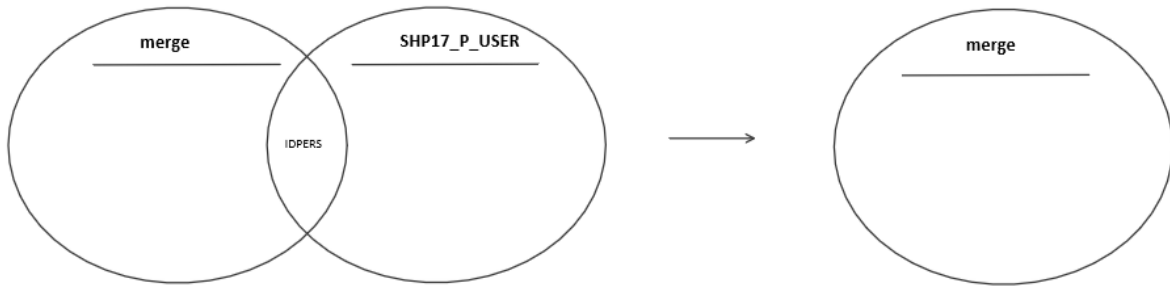
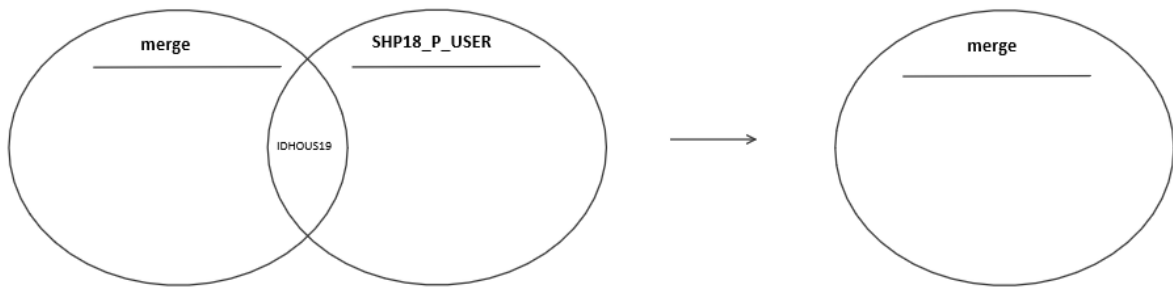
1. Illness / Accident (P20L01)
Hypothese: In der Pandemie kam es zu weniger Unfällen, wegen den strikten Einschränkungen in der Mobilität, deshalb sind die Menschen glücklicher.
2. Happy with Partner(P20F54)
Hypothese: Zielgrösse ist abhängig von der Zufriedenheit mit dem Partner.
3. Trust in Federal Government(P17P04)
Hypothese: Wer durch die Pandemie das Vertrauen in die Regierung verlor, ist unglücklicher.
Die Variable ist den Jahren 2019 und 2018 nicht verfügbar, deshalb griffen wir auf das Jahr 2017 zurück.

Datenaufbereitung

Die Datenaufbereitung erwies sich als keine leichte Kost. Zuerst selektierten wir die benötigten Variablen, um die Zusammenführung der Daten nicht unnötig gross zu gestalten. Danach wurden die Datensätze wie folgt zusammengeführt (Alle sind mit Left-Joins zusammengeführt).

Zusammenführen der Daten gemäss Skizze





Unser Datensatz ‘merge’, beinhaltet somit alle benötigten Variablen, um die Bereinigung der Daten zu starten.

Resultate des Parallel-Codings

Stojche und Ahmed führten gleichzeitig eine explorative Datenanalyse der Daten durch und verglichen schliesslich die Analyse. Ahmad erstellte Diagramme, die sich auf die Verteilung der einzelnen Kategorien konzentrieren. Stojche erstellte Diagramme, die sich sowohl auf die Verteilung jeder Kategorie als auch zusätzlich auf die Anzahl der Datenpunkte, die zu dieser Kategorie gehören, konzentrieren. Seine Variante deutet darauf hin, ob die Verteilung in jeder Kategorie im Verhältnis zur Anzahl der Datenpunkte in dieser Kategorie sinnvoll ist.

Datenqualität

NA-Werte:

| Variable | Anzahl (% - Anteil der Gesamtdaten) | Grund |
|---------------------|-------------------------------------|---|
| Change_satisfaction | 63 (~1%) | Der Grossteil der NA's ist aus des Covid-Daten(C20PC44), es handelt sich unserer Meinung nach um eine geringe Anzahl NA' bezogen auf die Gesamtdaten. |
| Relation_status | 72 (~1,5%) | Coviddatensatz beinhaltet fehlende Daten. |
| Change_in_relation | 170 (~3,5%) | Die Variable (C20D29) beinhaltet 72 NA's, die sich durch die neue Definition der Variable auf 170 erhöhten. |
| Workingsituation | 151 (~3%) | Entstehung der NA's durch die Kombination verschiedener Variablen. Mögliche Gründe für diese NA's könnten sein, dass Leute ihre Arbeitssituation preisgeben wollen. |
| Corona_infection | 99 (~2%) | Fehlende Werte sind gering. |
| Household_type | 3 (~0.05%) | Werte kommen von der Variable. |
| Income | 383 (~7%) | Grosse Anzahl fehlende Werte. NA's sind bereits von der Variable (I19HTYN). Es könnte durch Schweizer-Spiessbürgertum erklärt werden. Man redet nicht über Geld. |
| Financial_problem | 309 (~6%) | Grosse Anzahl fehlende Werte. Variable ist keine Kombination aus mehreren Variablen! Man sieht ähnlich Situation wie Income. |
| Illness_accident | 368 (~6.5%) | Die allermeisten Werte entstammen der Variable (P20L01), die Variable (P19L01) enthält nur drei NA's. |
| Happy_with_partner | 1585 (~28%) | Die beiden Grundvariablen (P19F54,P20F54) beinhalten beide über 1000 fehlende Werte. |

| Variable | Anzahl (% - Anteil der Gesamtdaten) | Grund |
|---------------------|-------------------------------------|--|
| Trust_in_government | 819 (~14%) | Die Variable (P17P04) enthält 533 fehlende Werte und die Variable (P20P04) enthält 390 NA's. |
| Total | 4022 (~71%) | Wenn wir alle NA's entfernen würden, gingen drei-viertel der Beobachtungen verloren. |

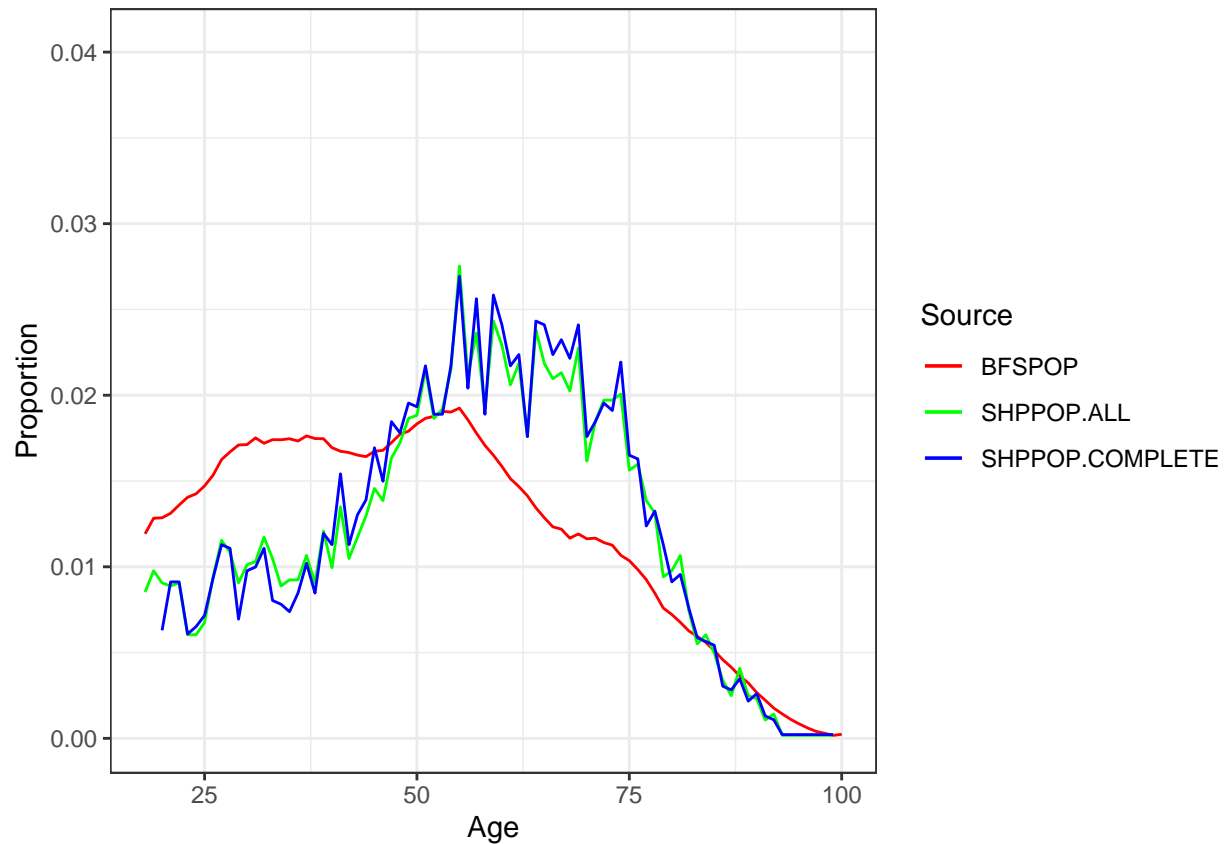
Fazit:

Variablen, die im Bezug auf die Gesamtbeobachtungen über 10% fehlende Werte haben, sollten für die Auswertung nicht verwendet werden, da die Stichprobe stark an ihrer Repräsentativität einbüsst, wenn alle fehlenden Werte entfernt werden würden.

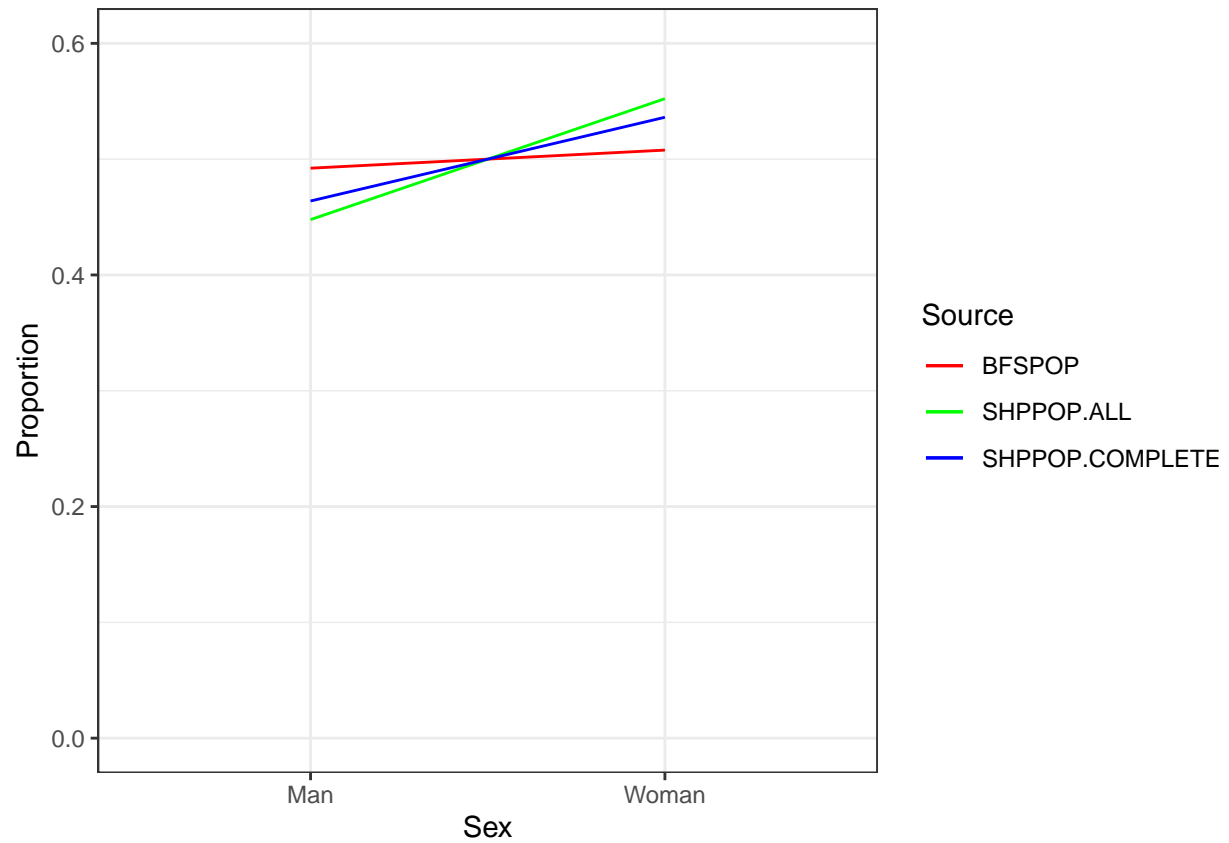
Deshalb änderten wir bei unseren selbstdefinierten Variablen alle NA's zum Zustand 'unknown', damit keine unnötigen Verluste der Beobachtungen in Kauf genommen werden müssen. Dadurch konnten wir die Verluste der NA's auf 1025 Beobachtungen reduzieren.

Plausibilität der Daten

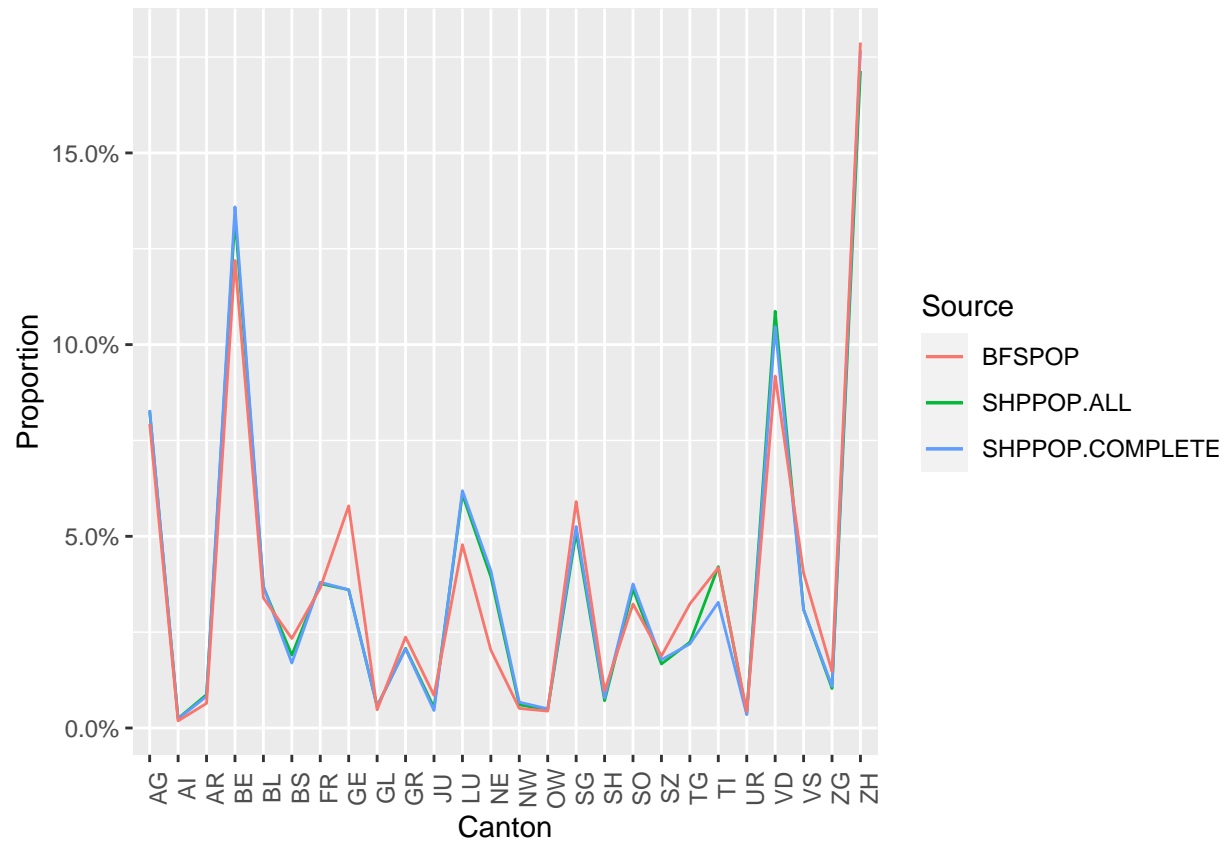
Um zu prüfen, dass die Stichprobe die Bevölkerung der Schweiz repräsentiert, vergleichen wir die Stichprobe mit Daten des BFS. Die drei Grafiken sollen darüber aufschliessen, wie repräsentativ unser Datensatz ist:



Wie sich auf der Grafik erkennen lässt, sind jüngere Altersgruppen, sprich die 25 bis 35-Jährigen Menschen in der Stichprobe unterrepräsentiert. Hingegen sind ältere Altersgruppen überrepräsentiert.



Im Bezug auf die Geschlechterverteilung der Stichprobe sind keine nennenswerte Unterschiede zu erläutern.

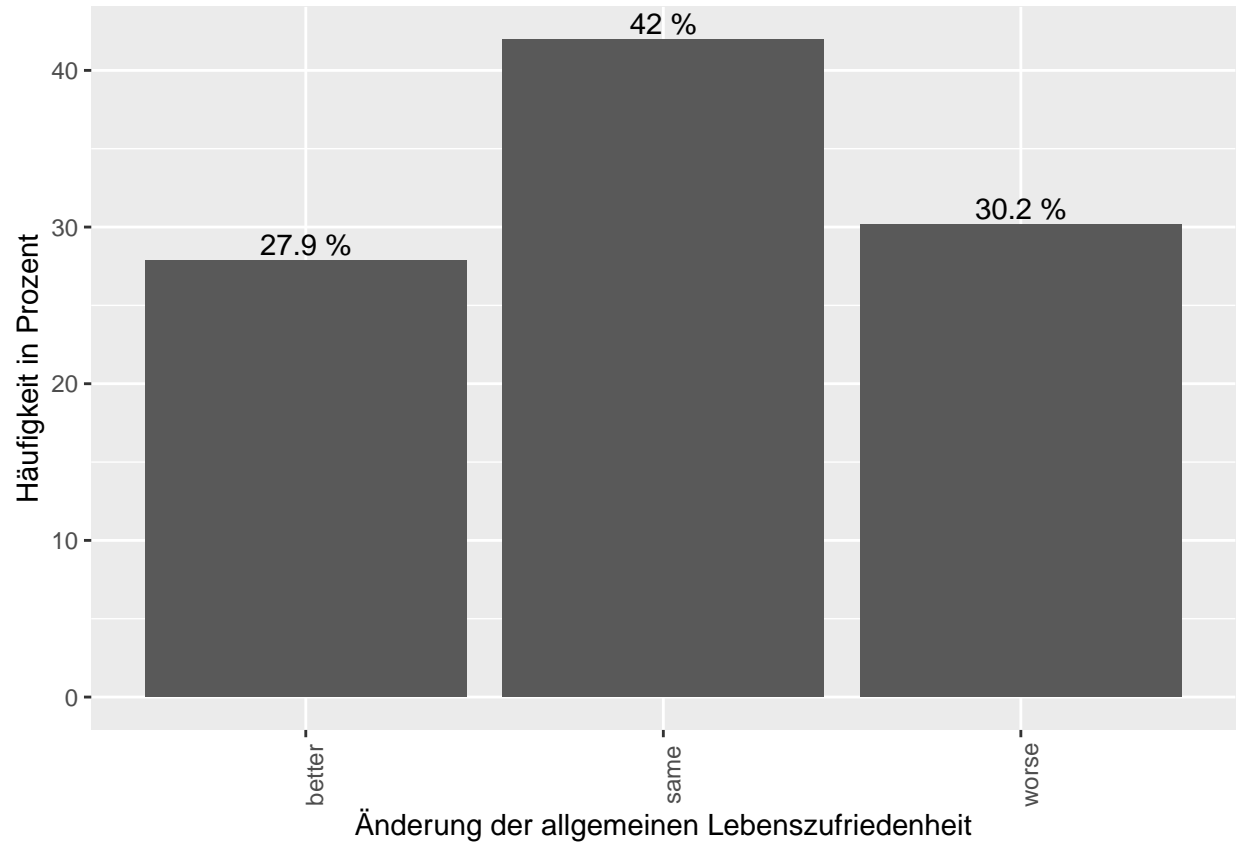


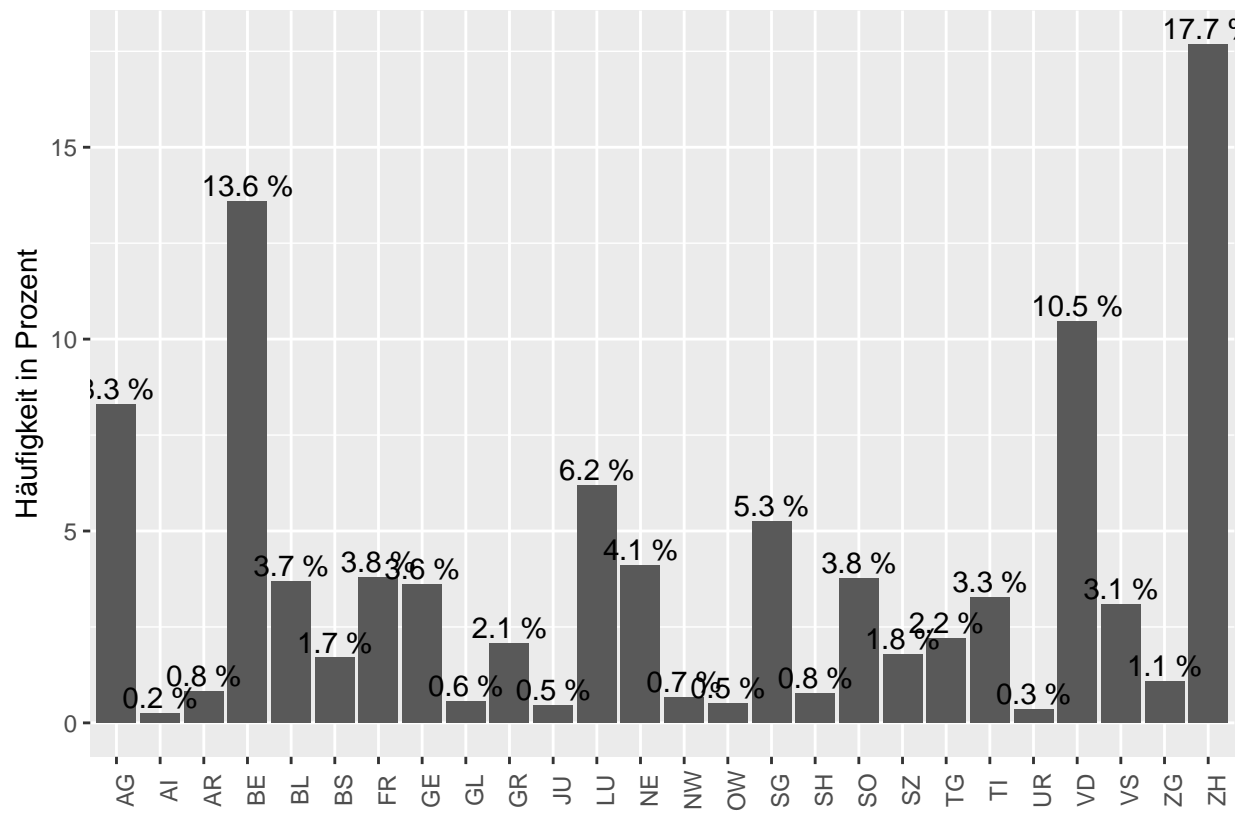
Die Vertretung der Kantone in der Stichprobe ist für die Schweizerbevölkerung repräsentativ.

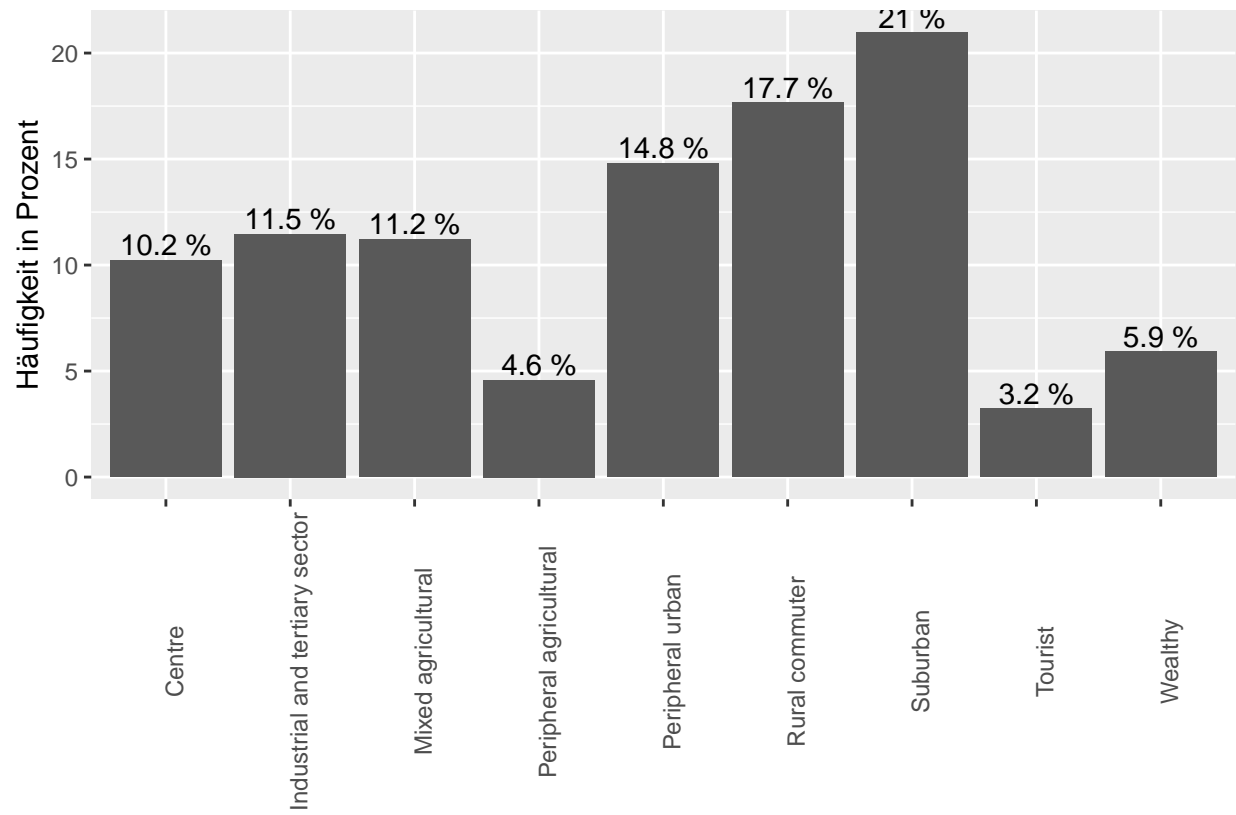
Auswertung

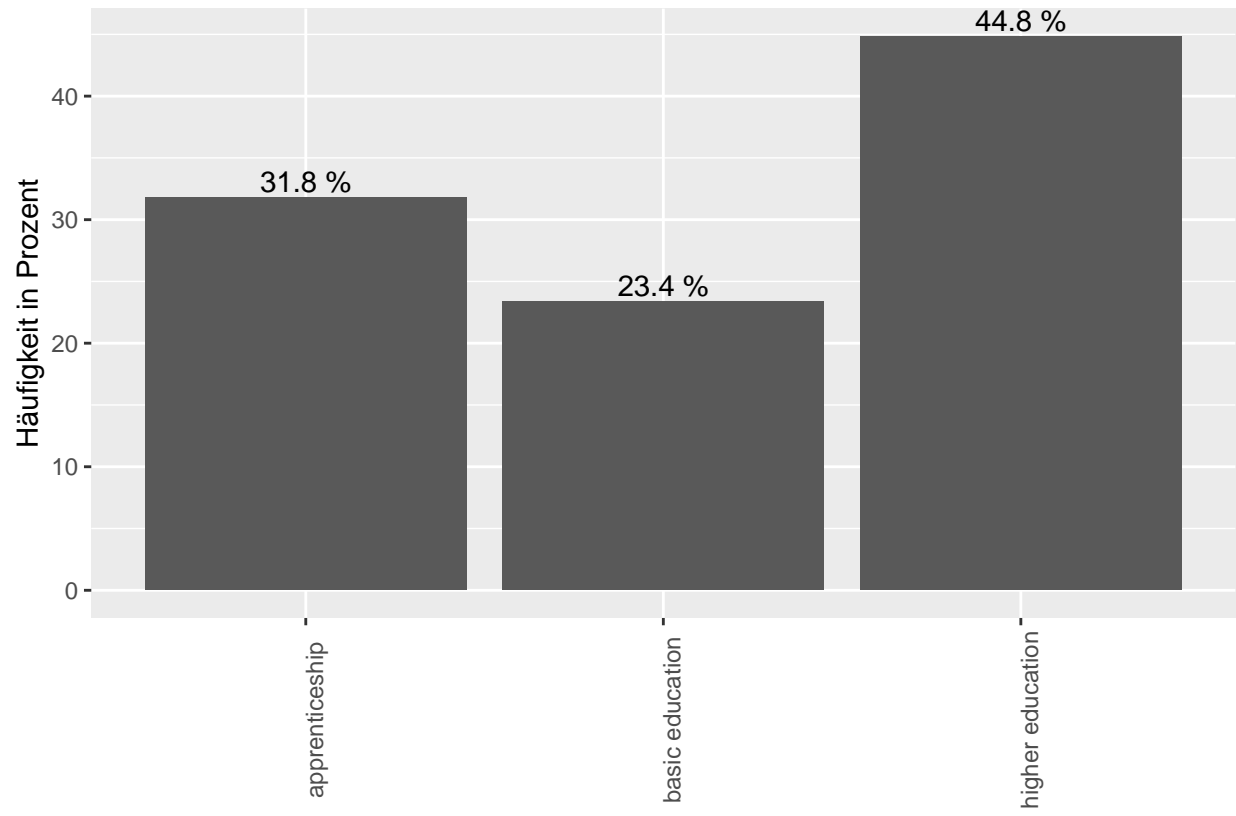
Univariate Grafiken

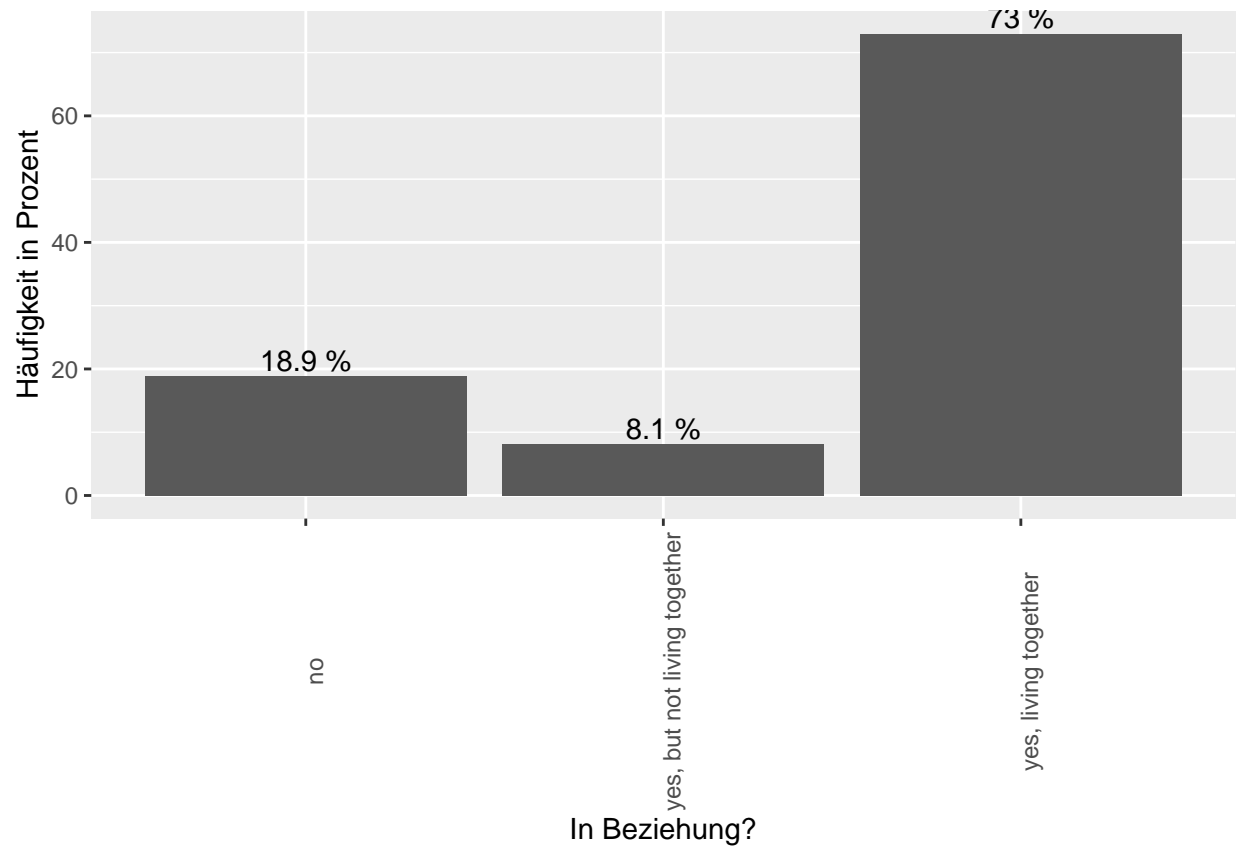
Die folgenden Abbildungen zeigen die explorative Datenanalyse mit jeweils einer Variable.

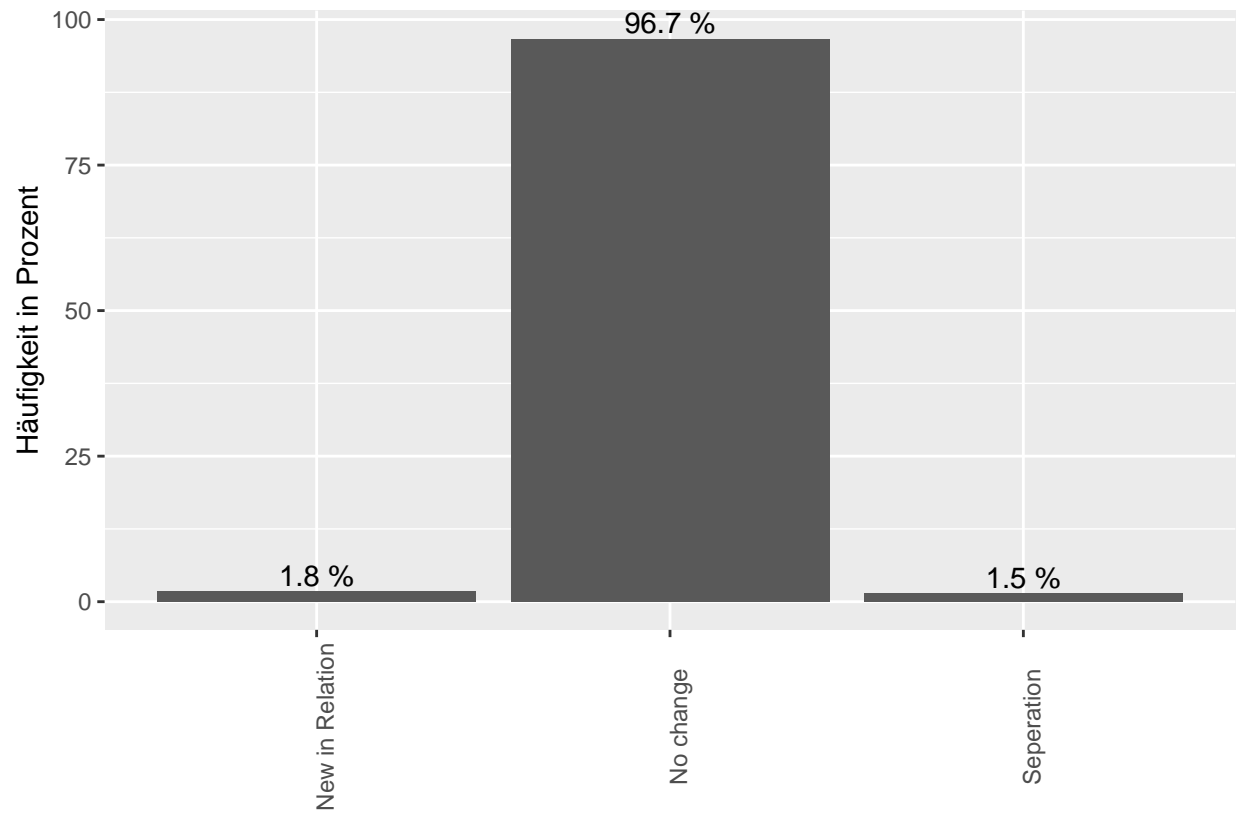


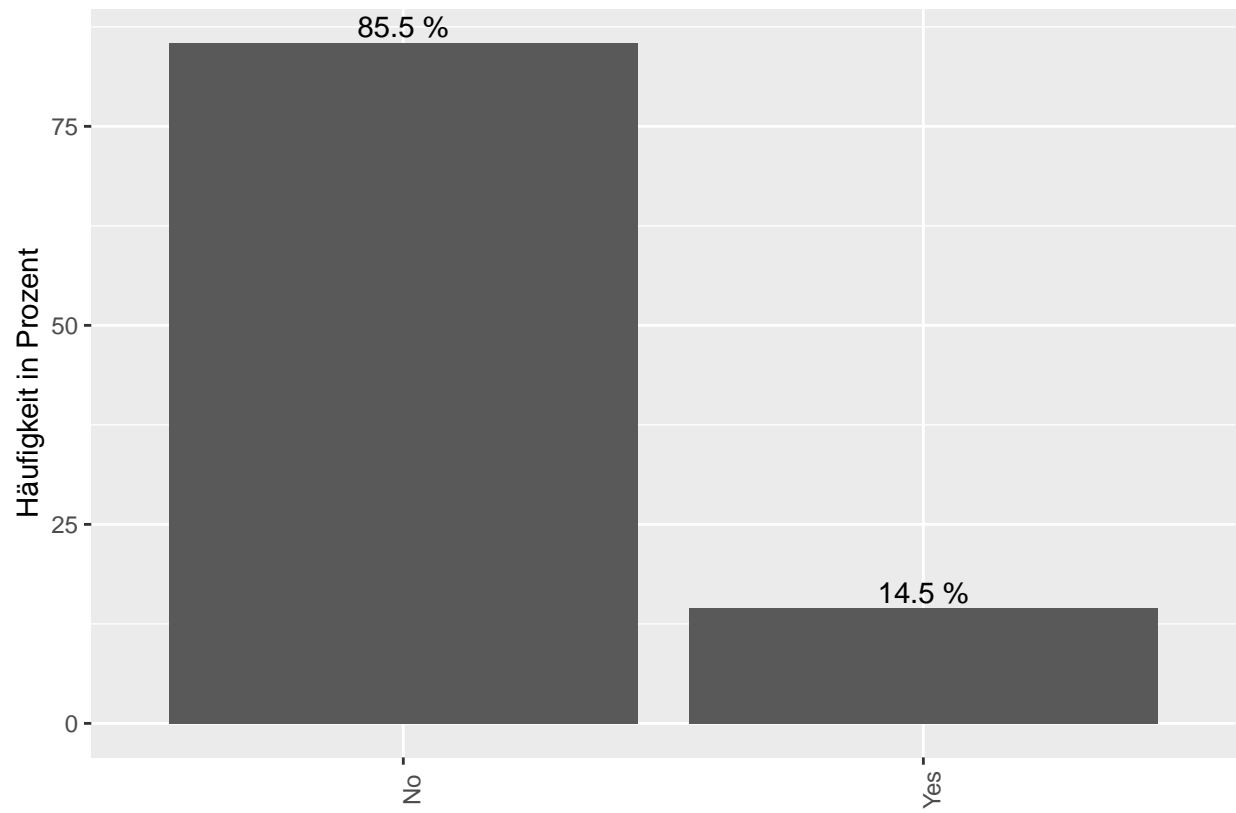


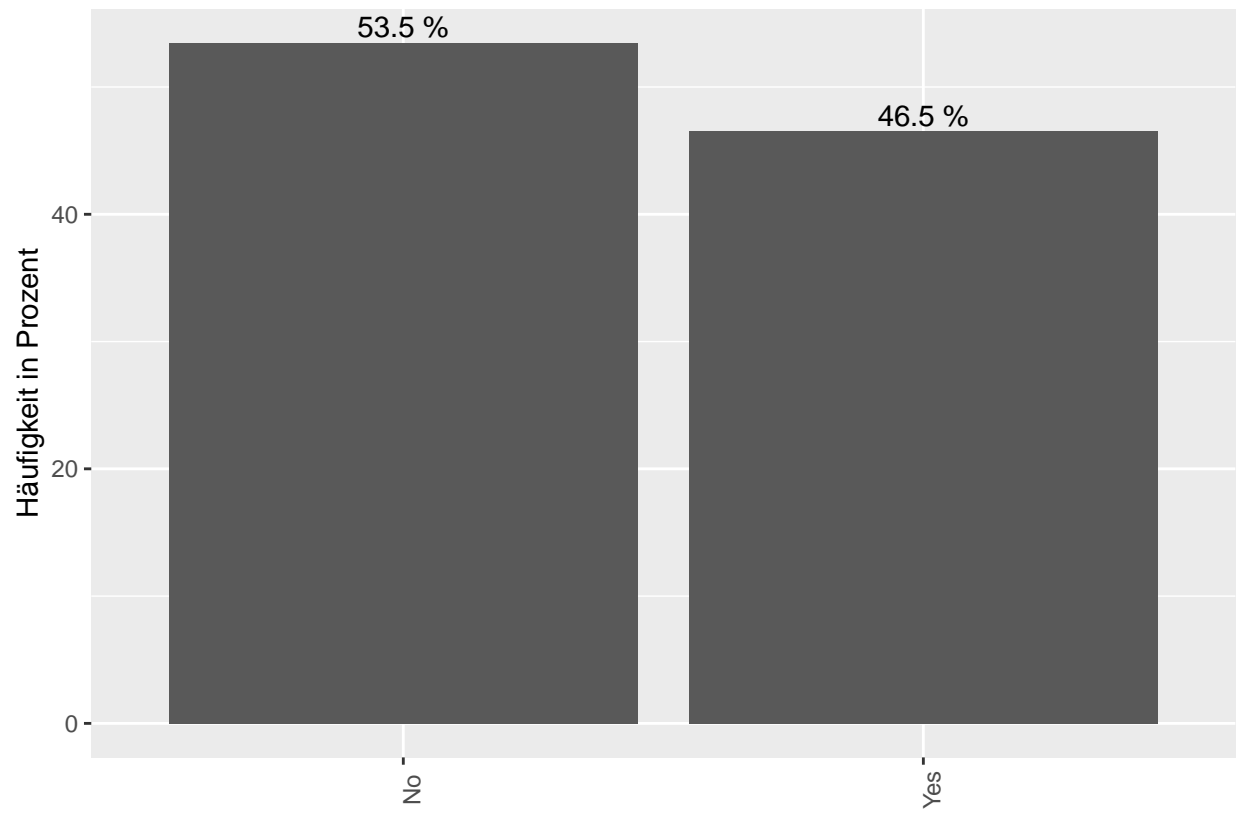


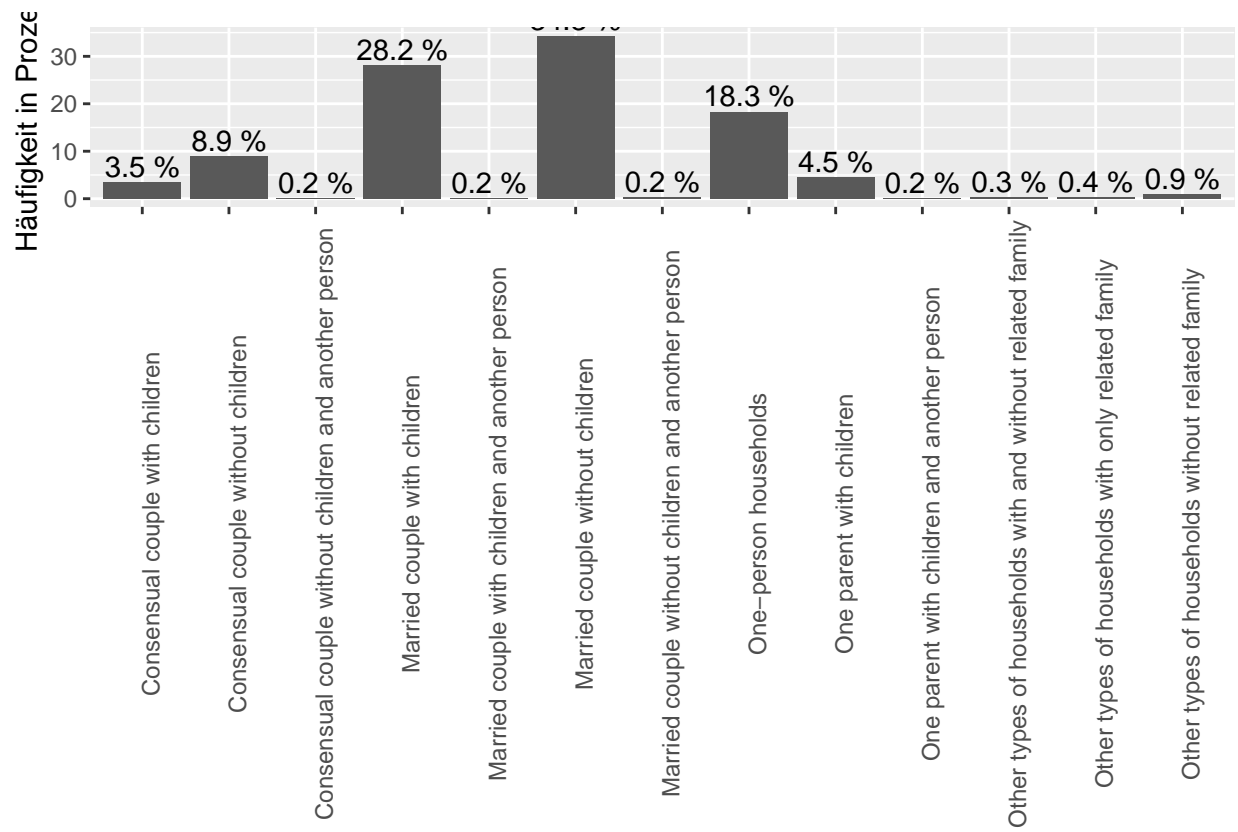


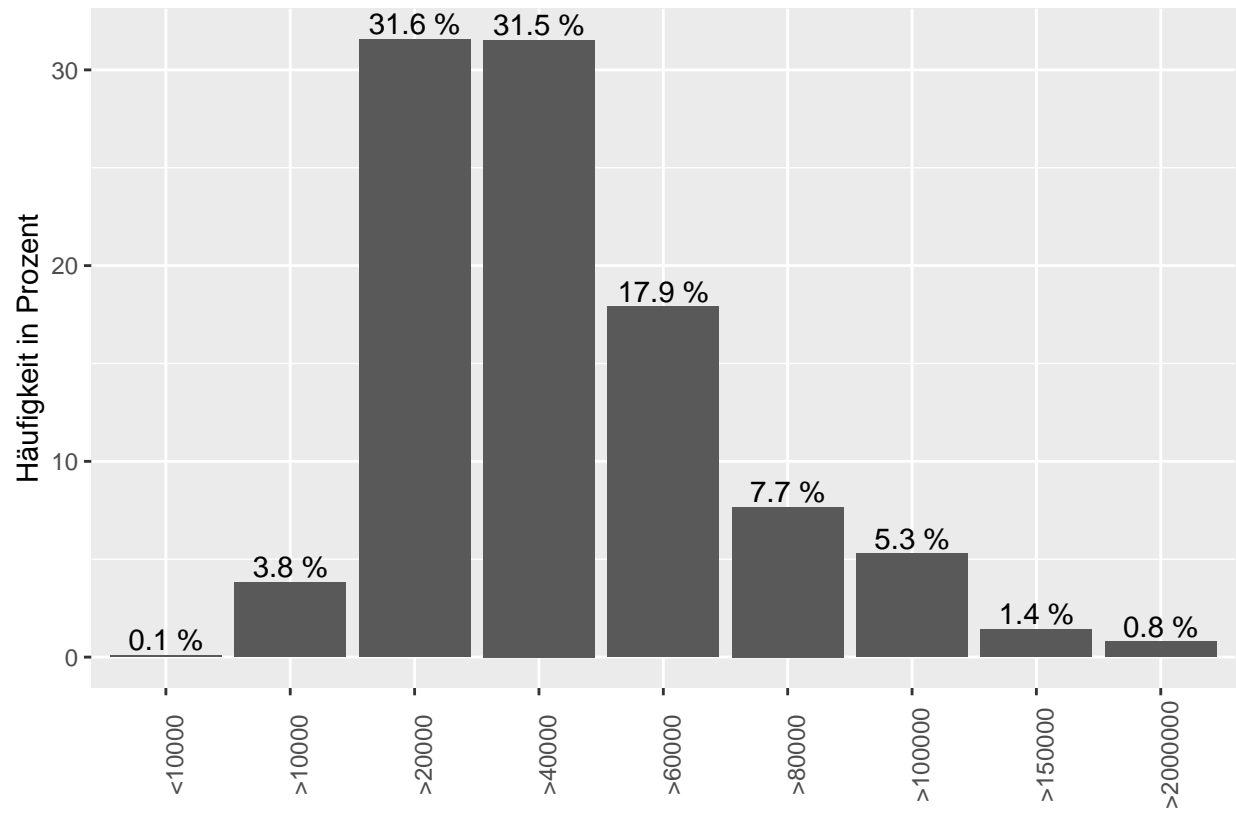


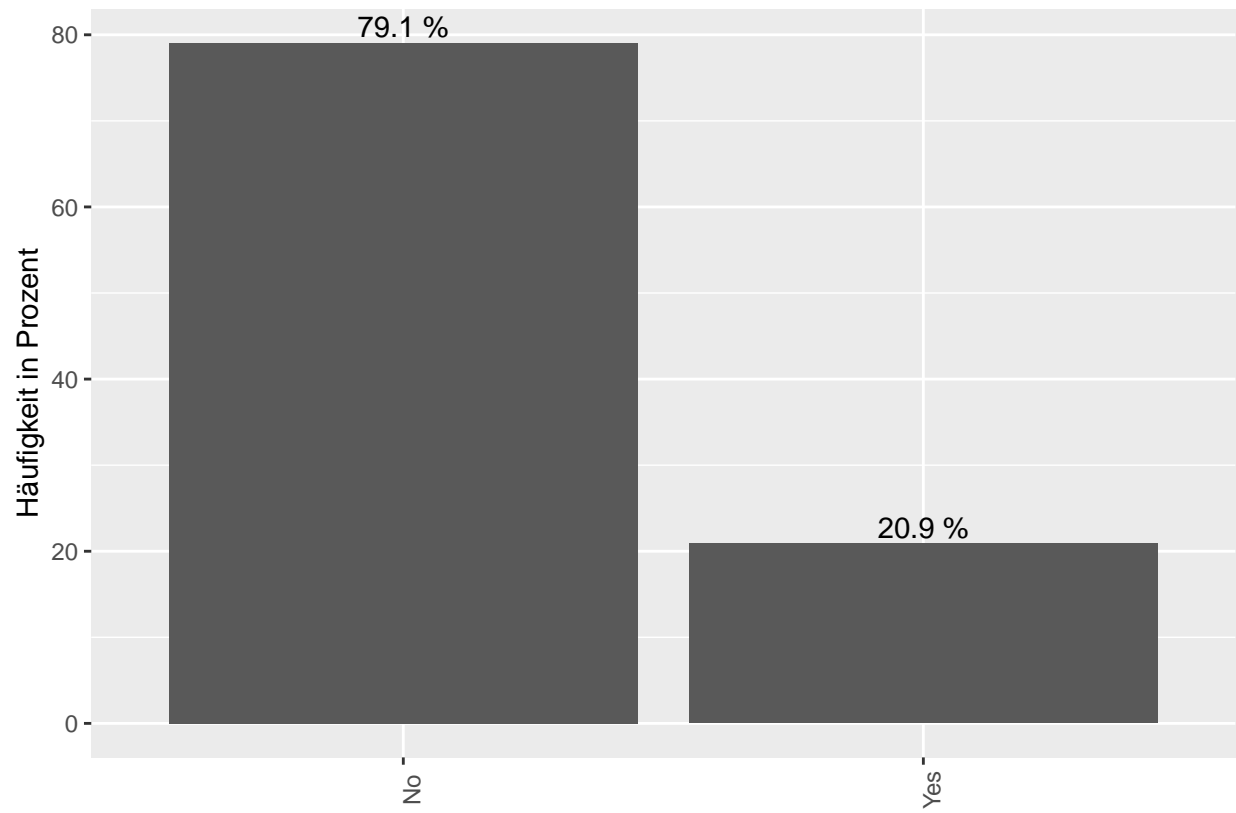


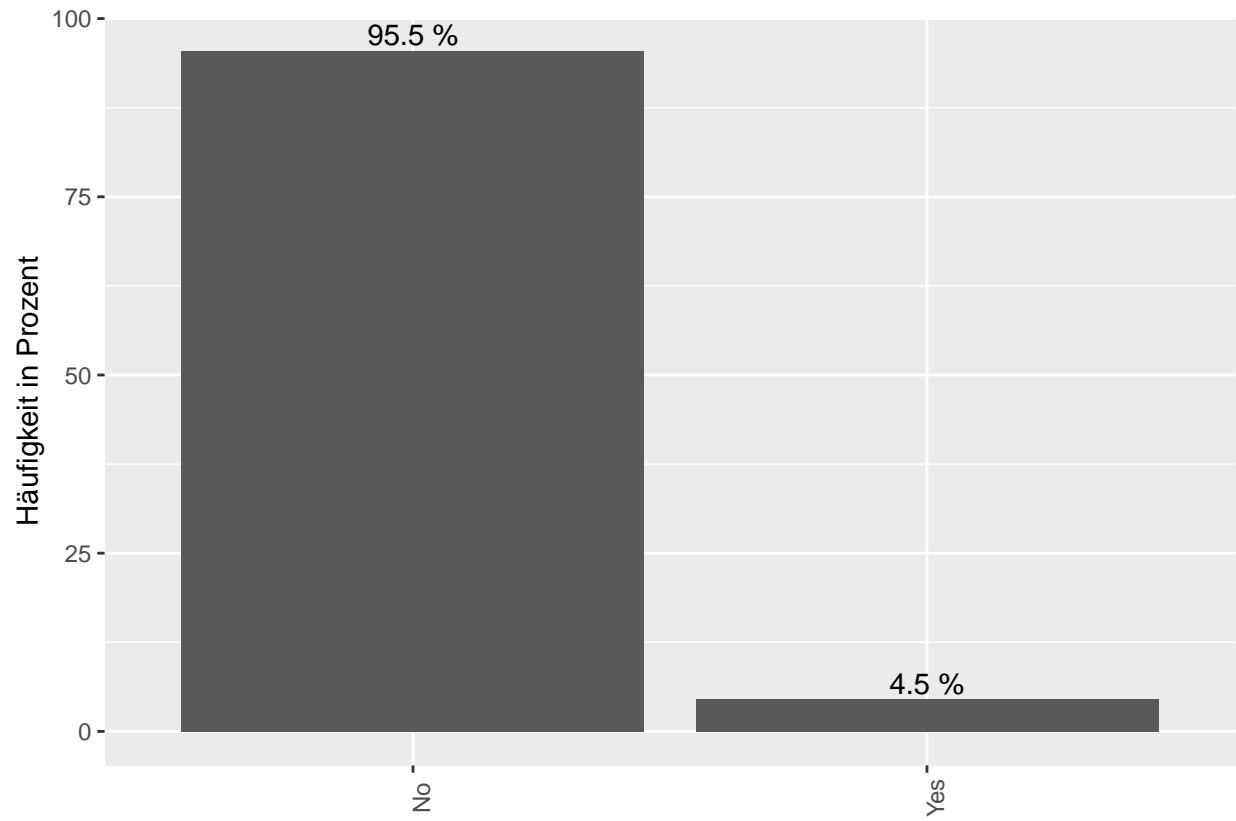


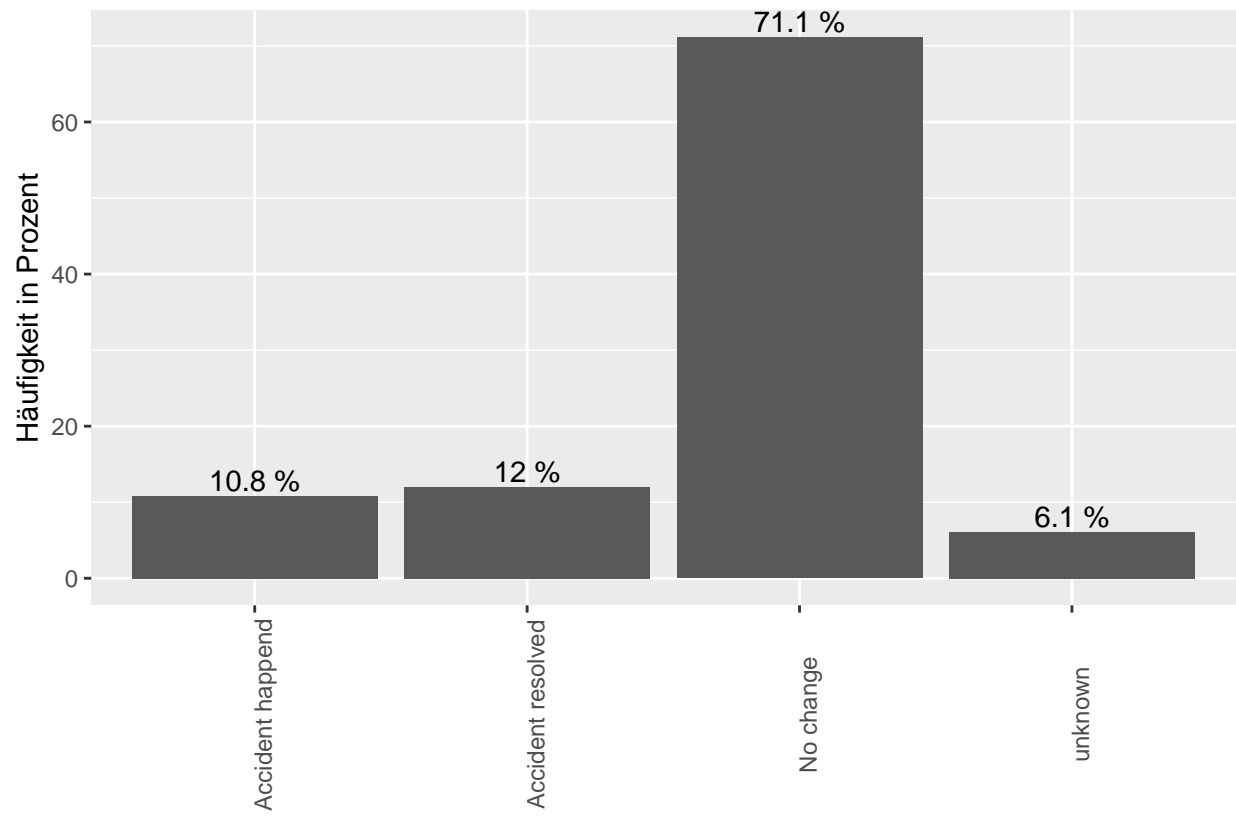


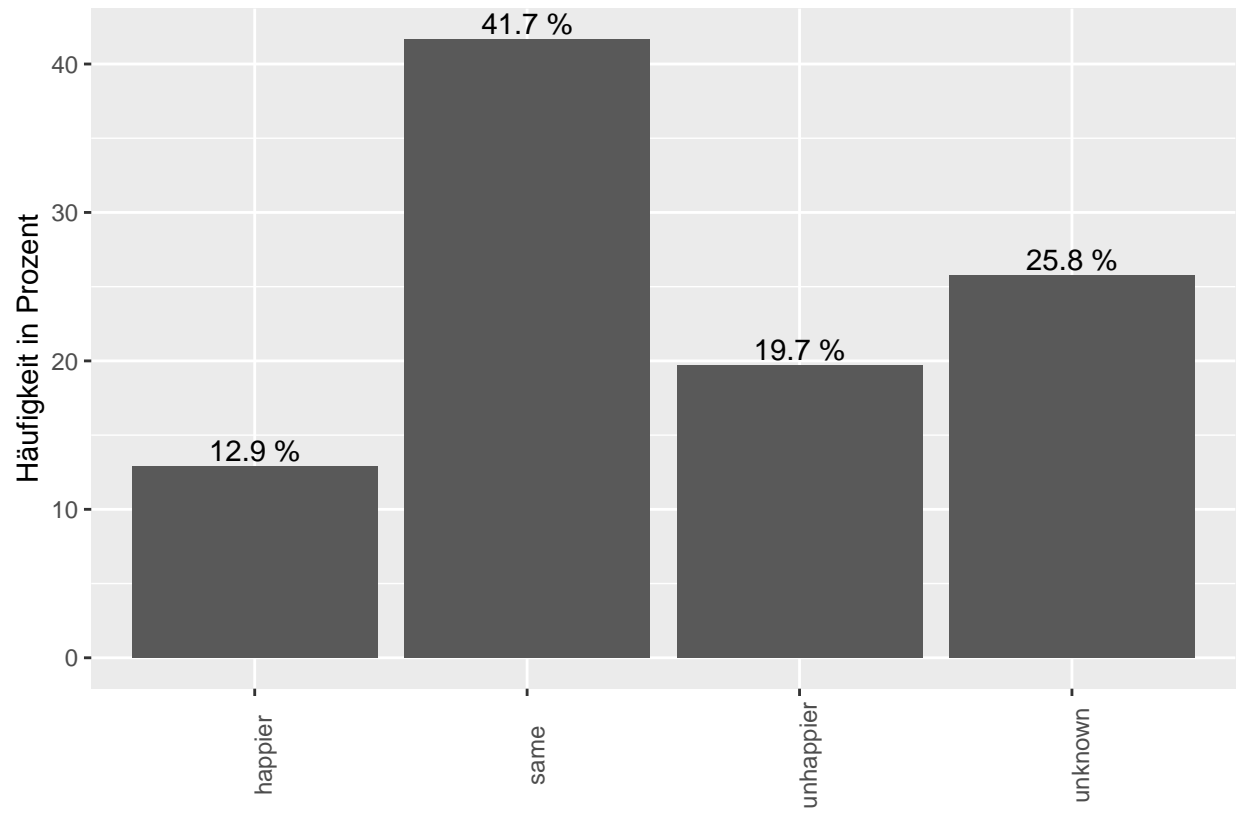


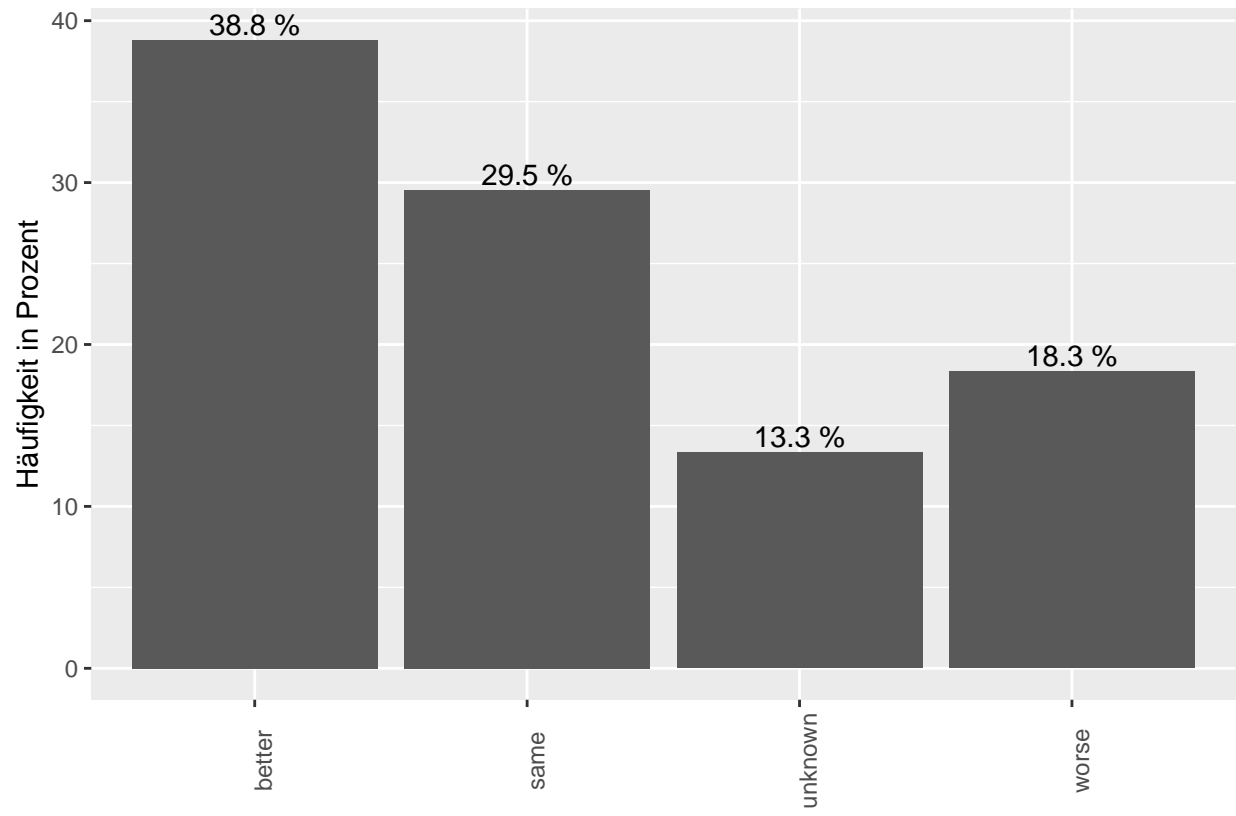








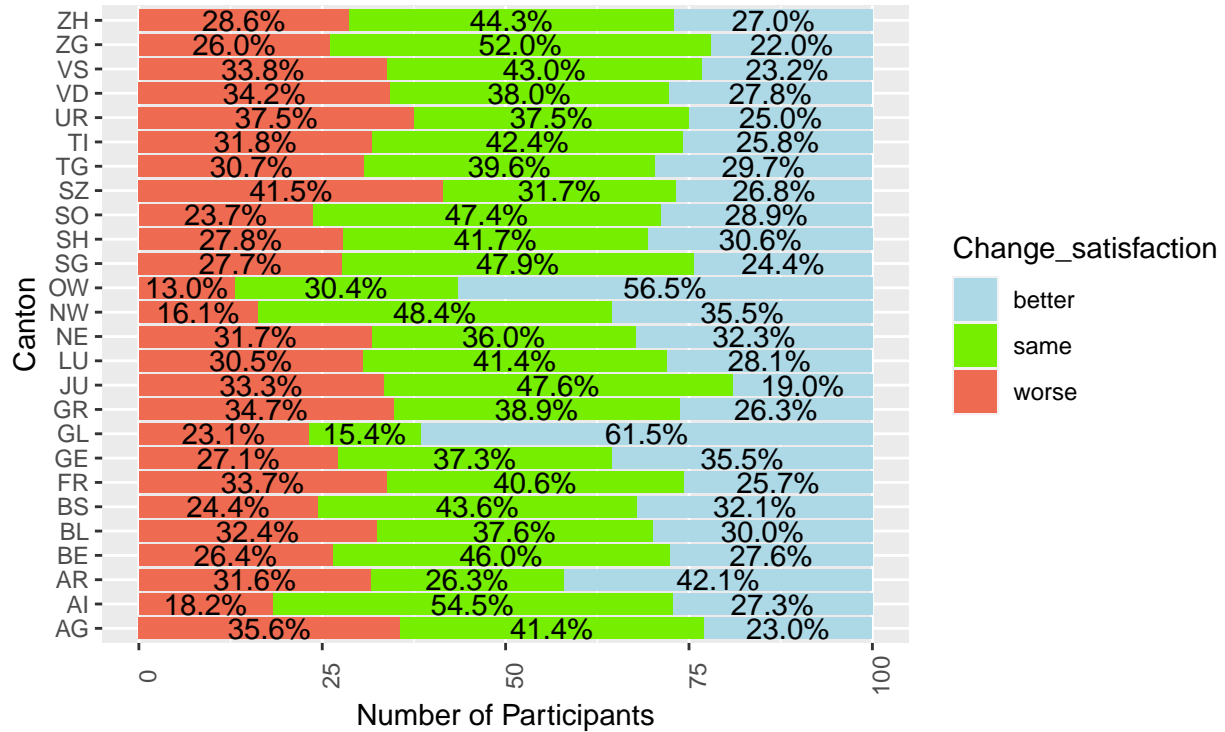




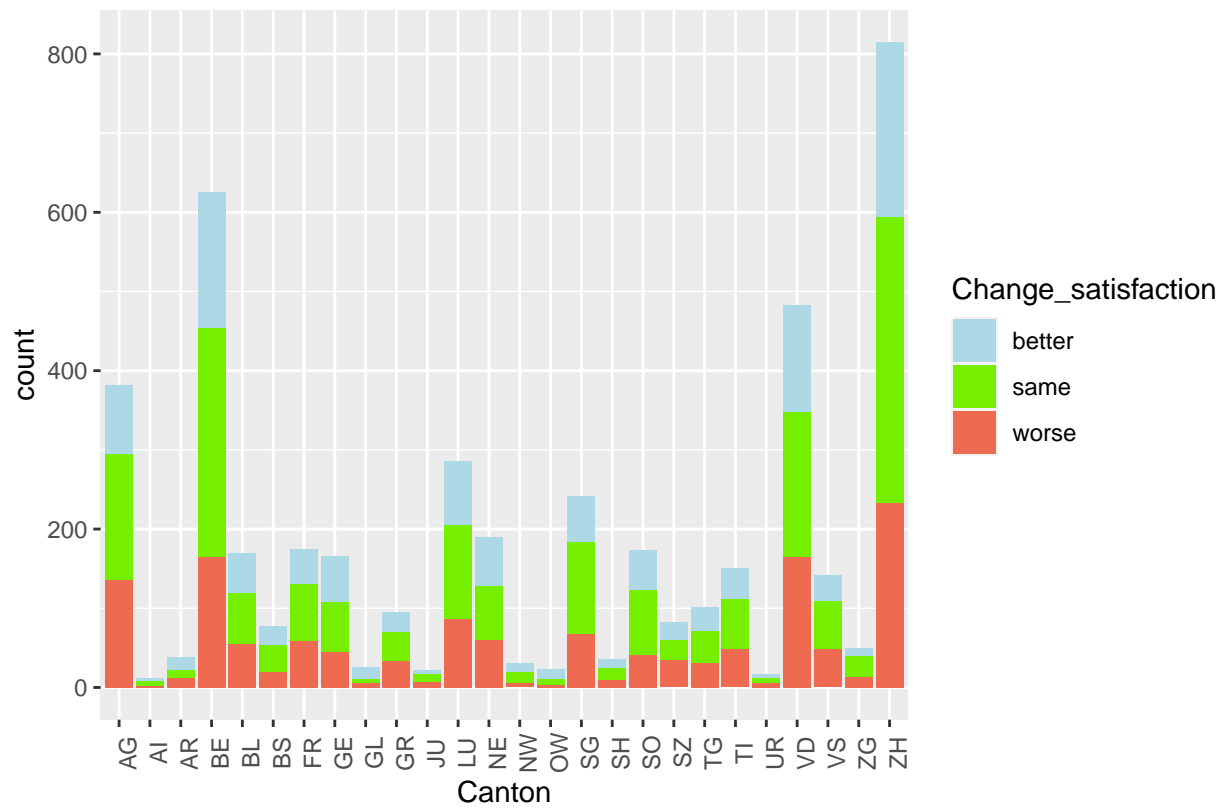
Bivariate Grafiken

Die nachstehenden Abbildungen sollen mögliche Zusammenhänge zwischen den Variablen und der Zielgrösse veranschaulichen.

Untersuchung des Einflusses auf die Zielgrösse

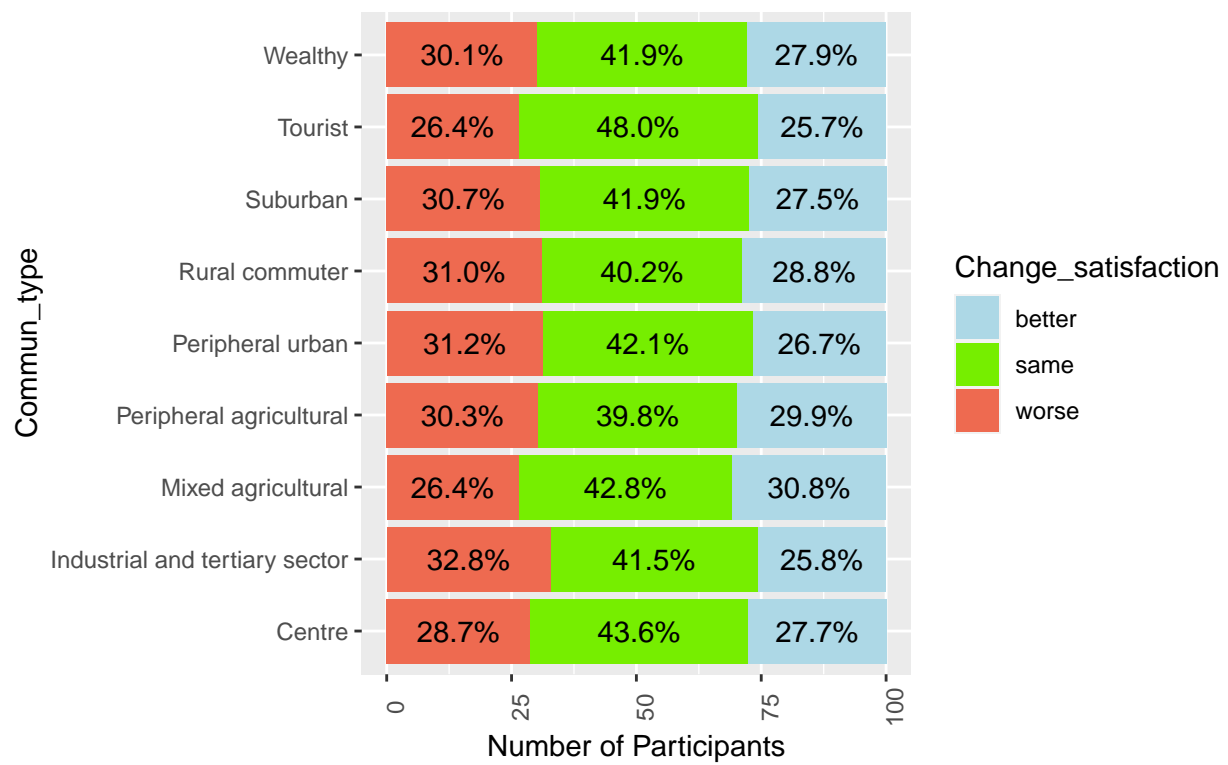


Im Kanton Glarus ist die stärkste positive Veränderung sichtbar,
im Kanton Schwyz trat genau das Gegenteil ein.

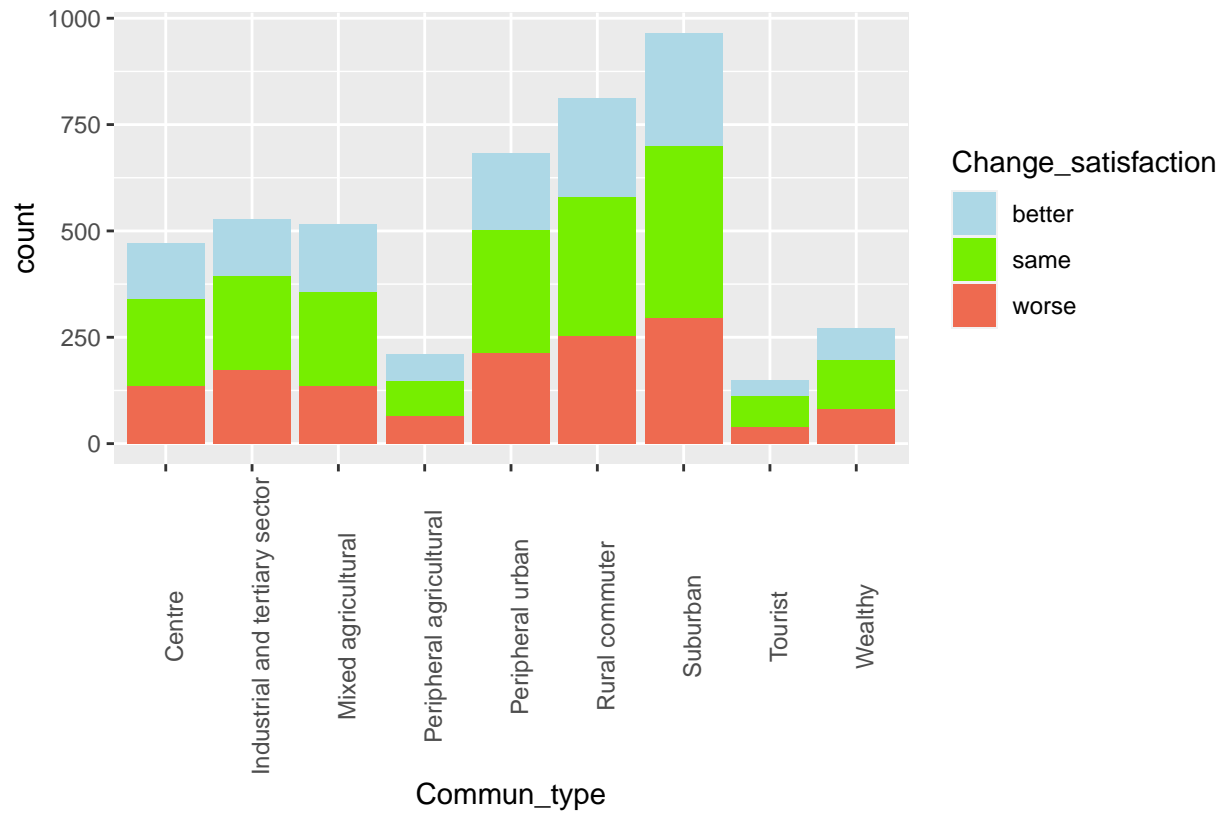


Absolute Häufigkeiten

Untersuchung des Einflusses auf die Zielgrösse

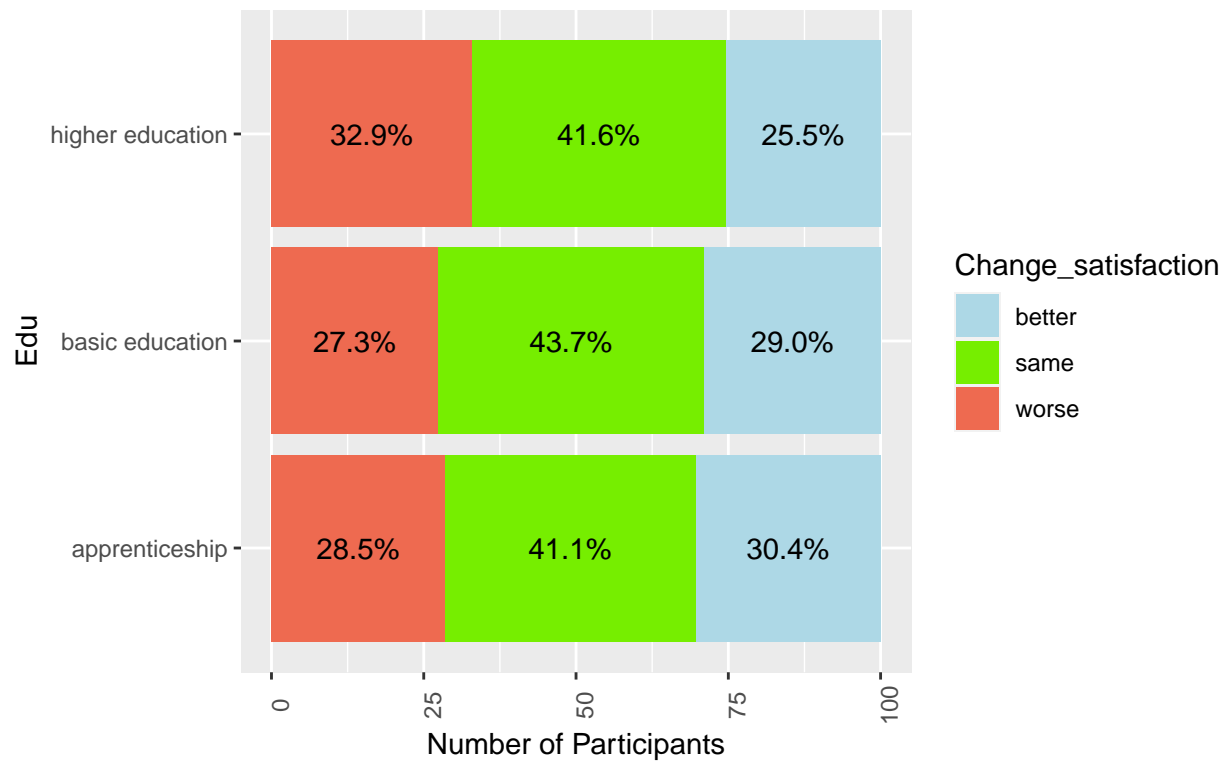


Der Gemeindetyp 'Industrie' hat die grösste negative Veränderung im Bereich der Lebenszufriedenheit.

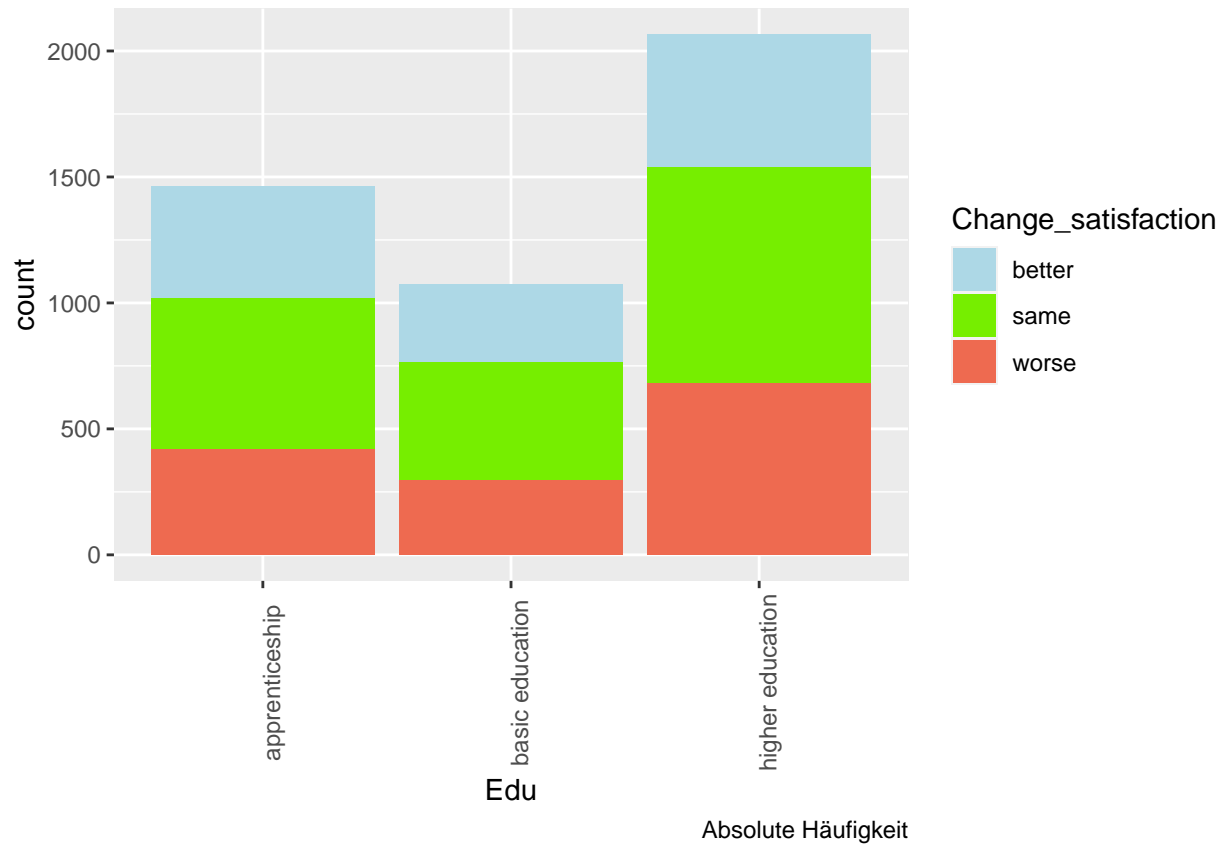


Absolute Häufigkeiten

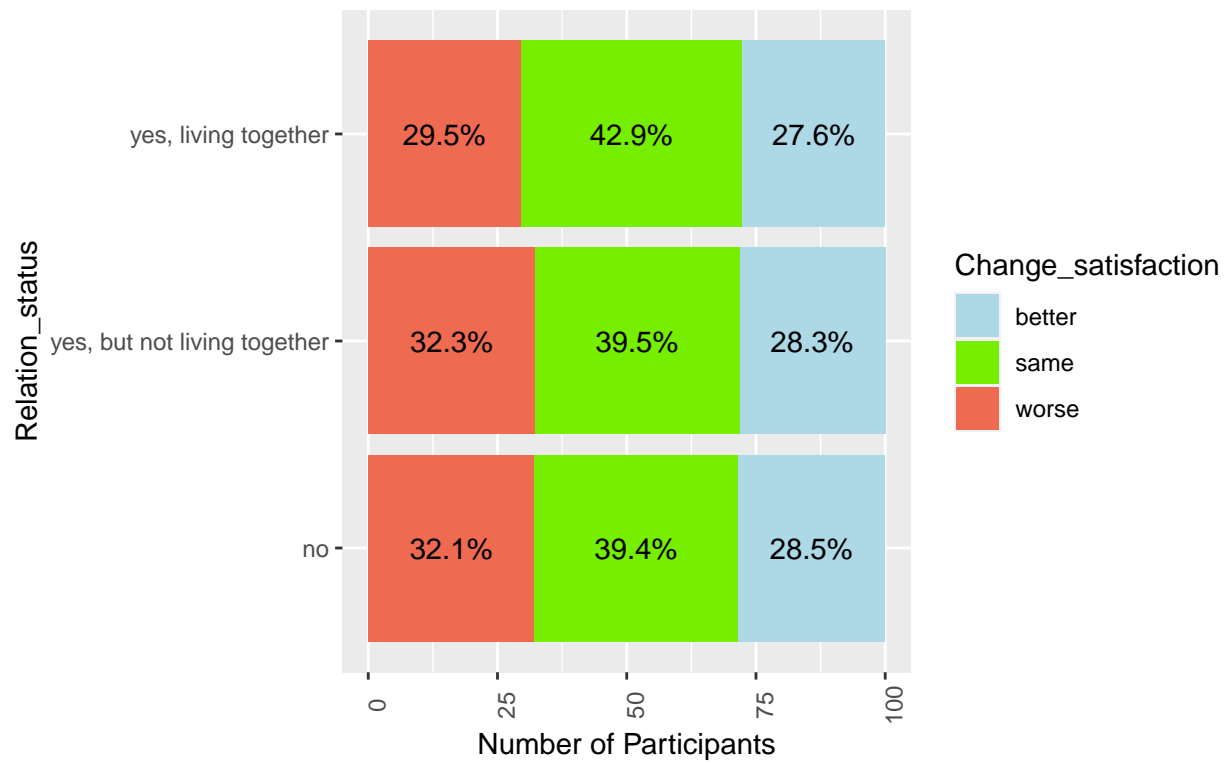
Untersuchung des Einflusses auf die Zielgrösse



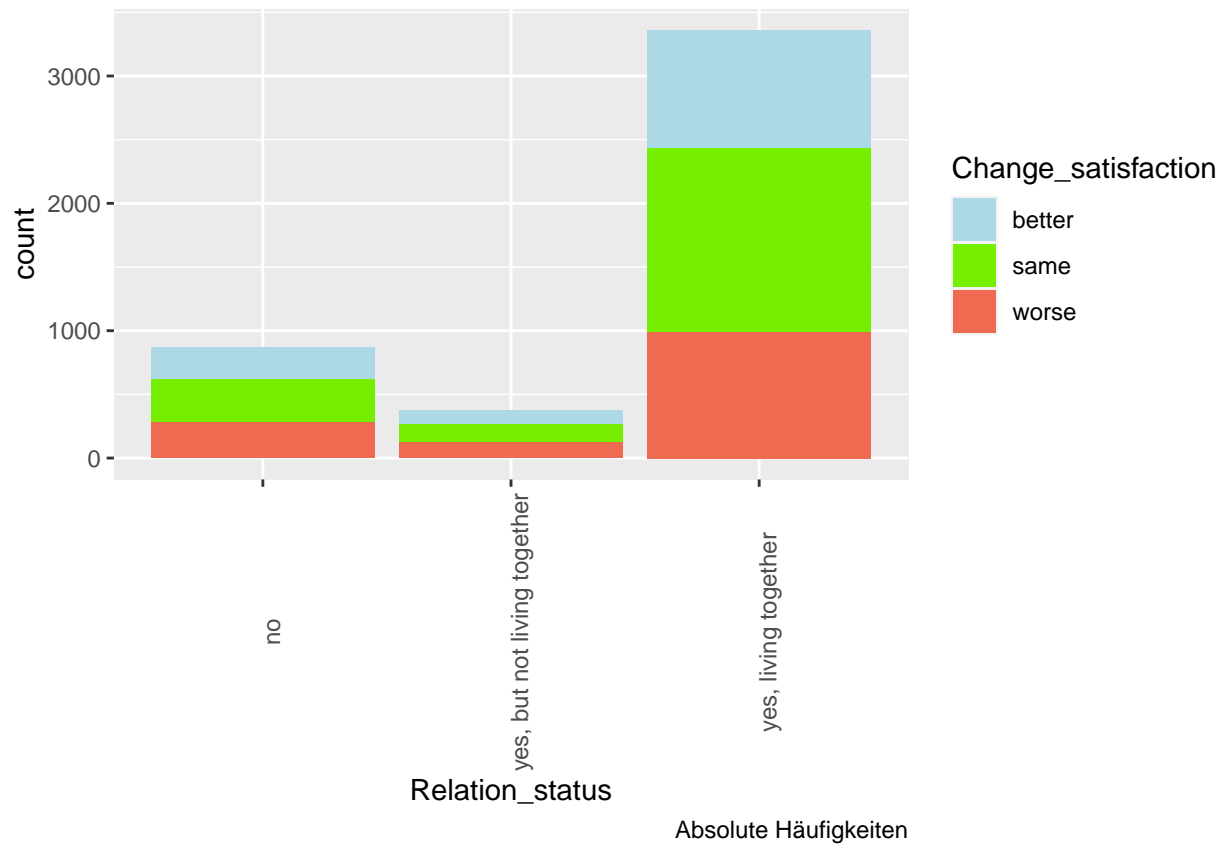
Leute mit tiefen Bildungsstand laut Graphik wurden in der Krise unzufriedener.



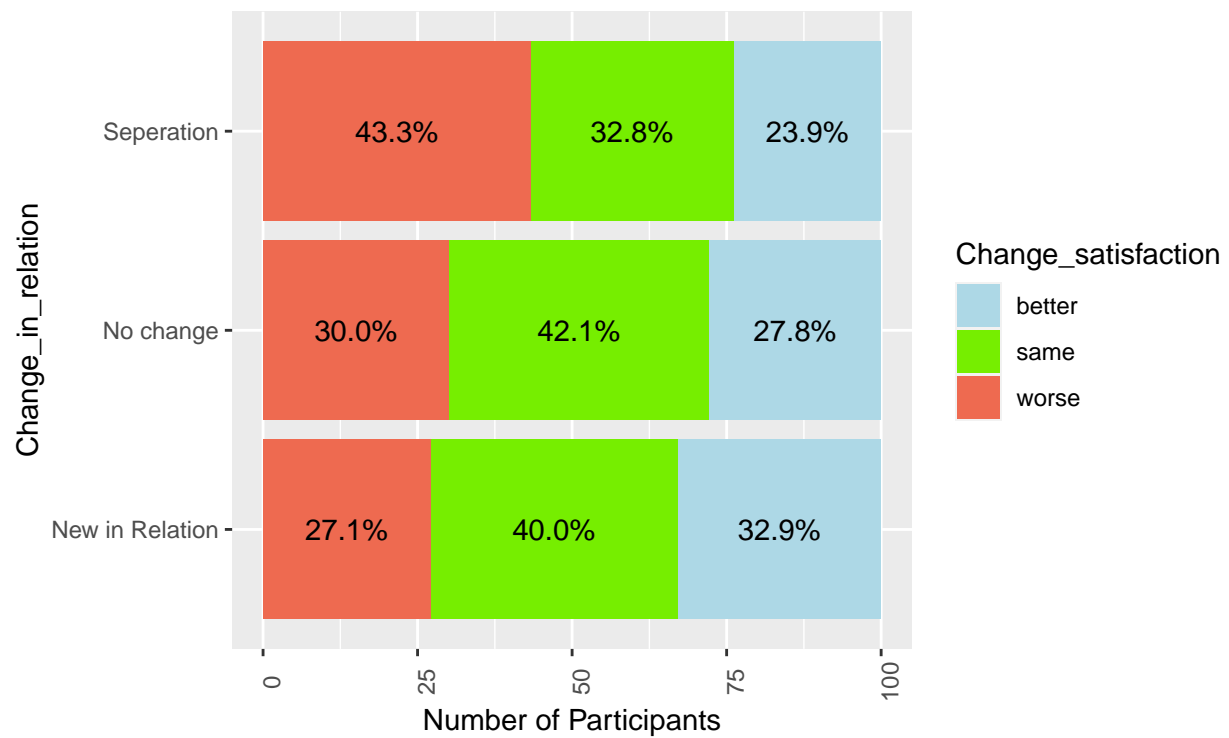
Untersuchung des Einflusses auf die Zielgrösse



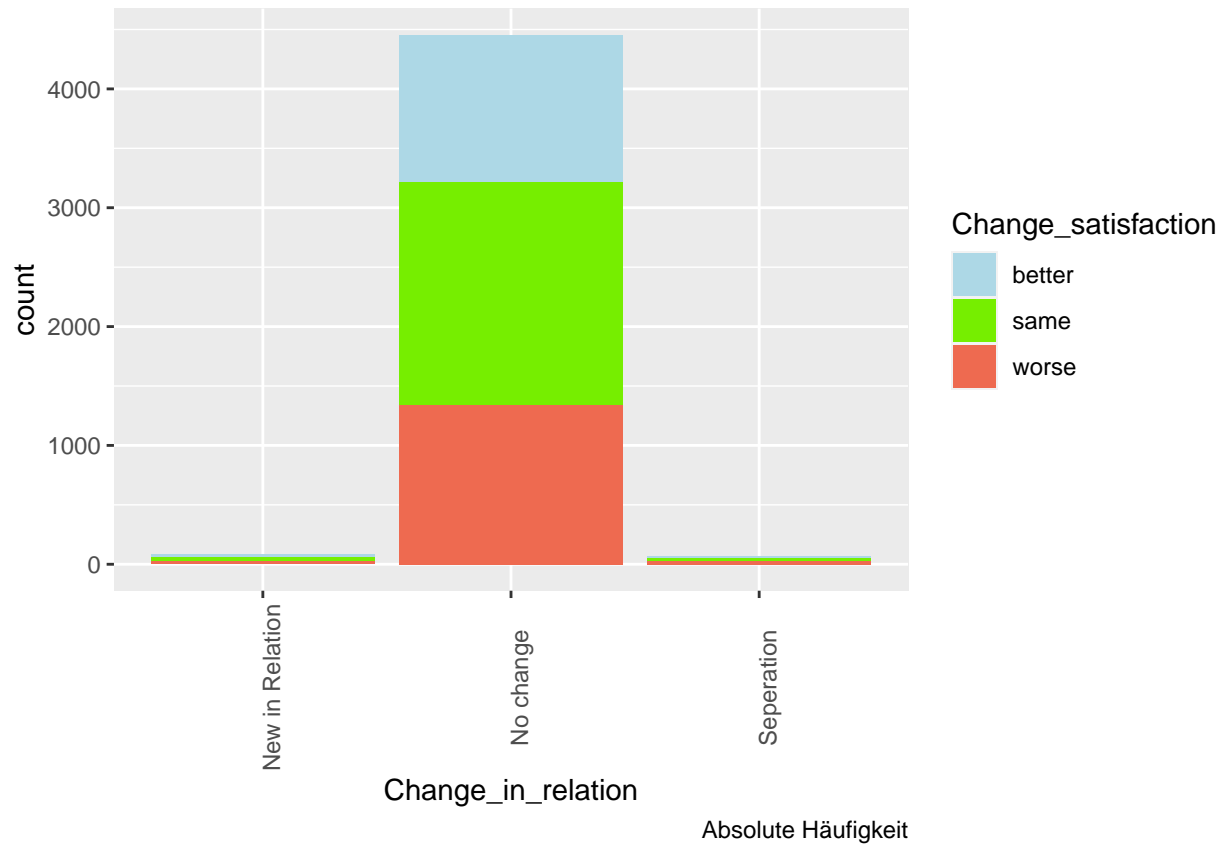
Der Beziehungsstatus hat keinen markanten Einflusses auf die Zufriedenheit



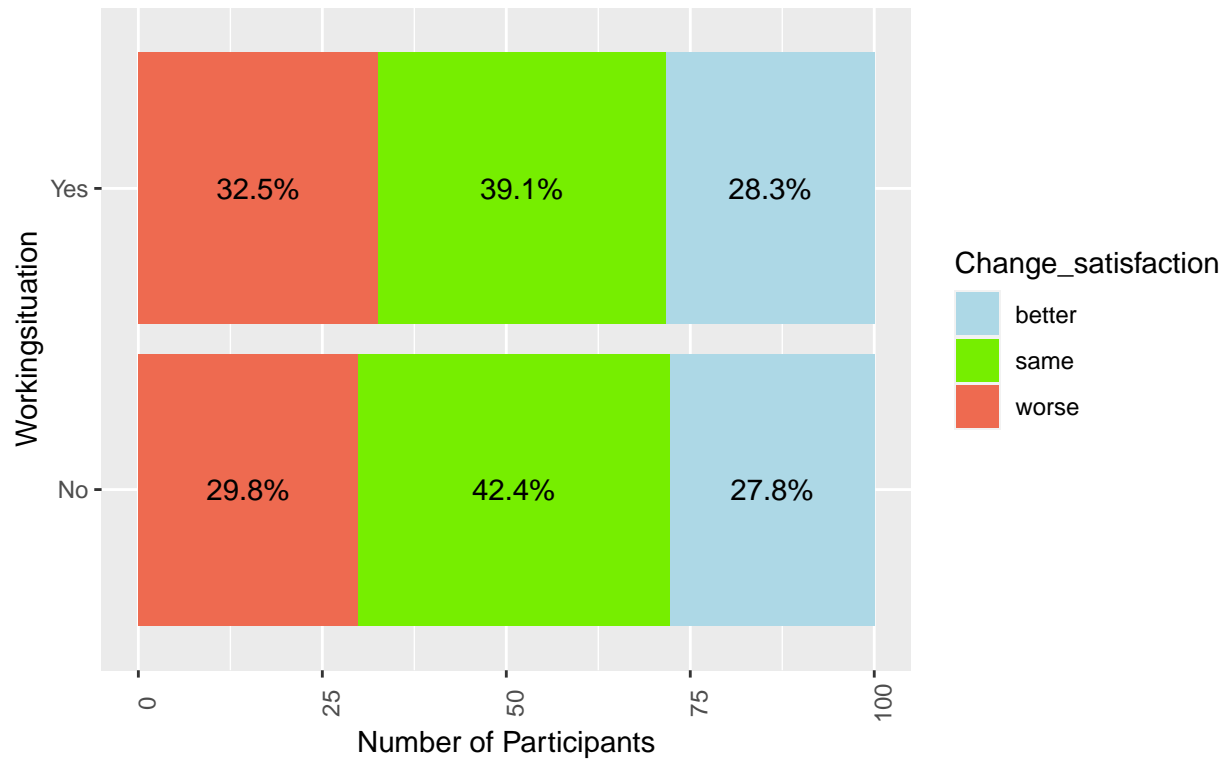
Untersuchung des Einflusses auf die Zielgrösse



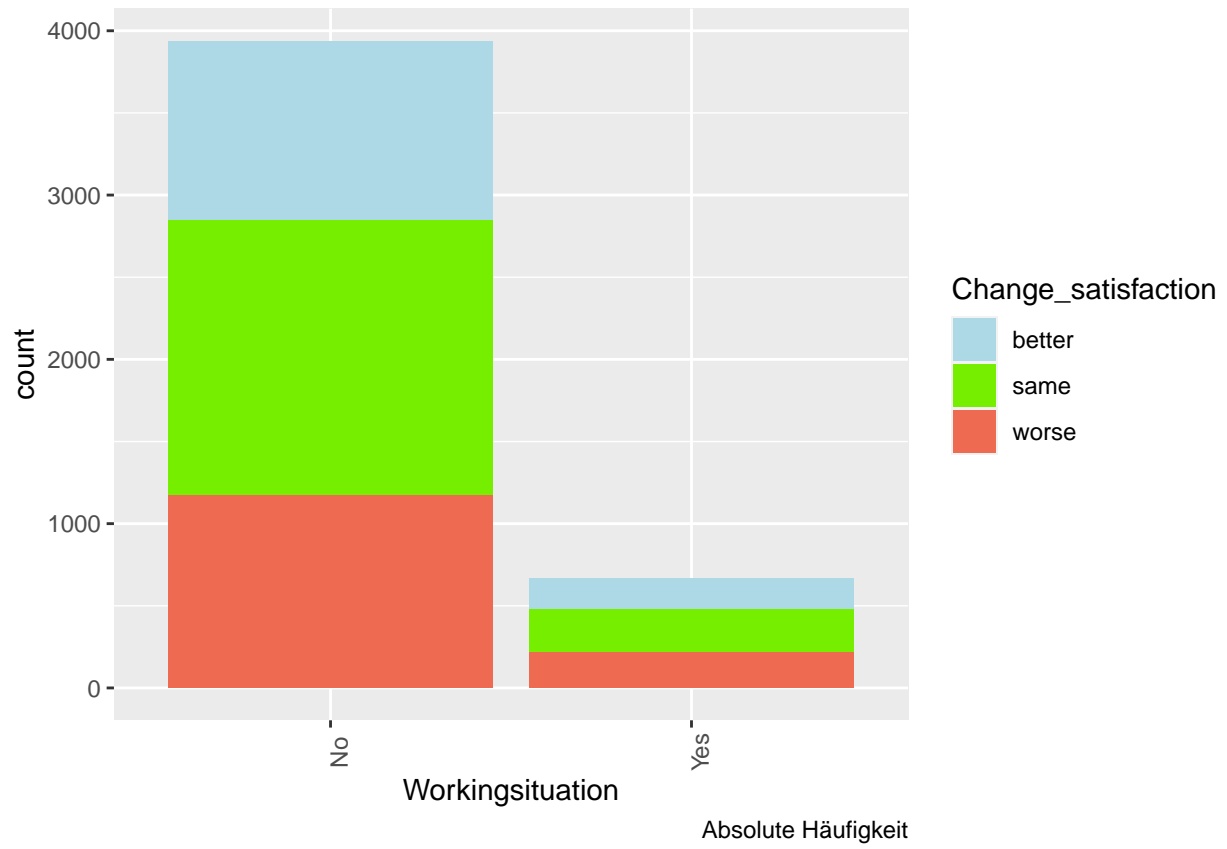
Die Grafik zeigt, dass es den Menschen mit einer neuen Beziehung besser geht als den Menschen, die sich von ihrem Partner getrennt haben.



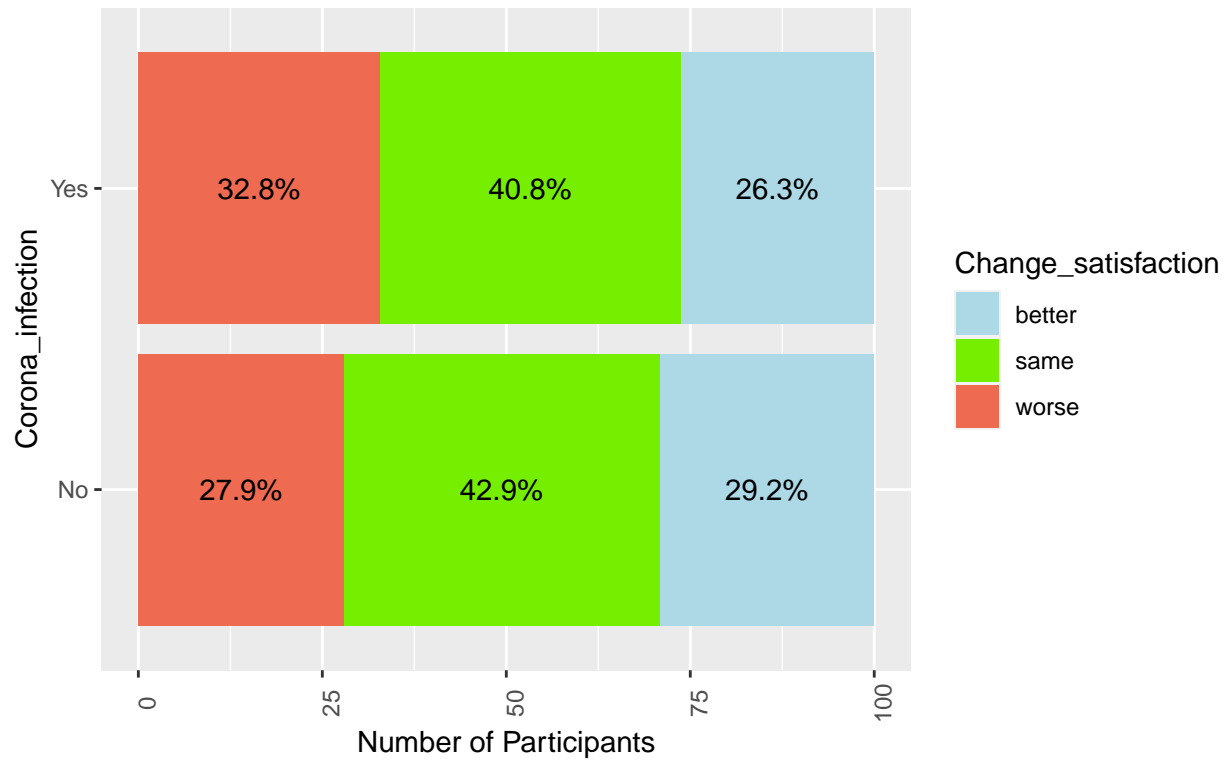
Untersuchung des Einflusses auf die Zielgrösse



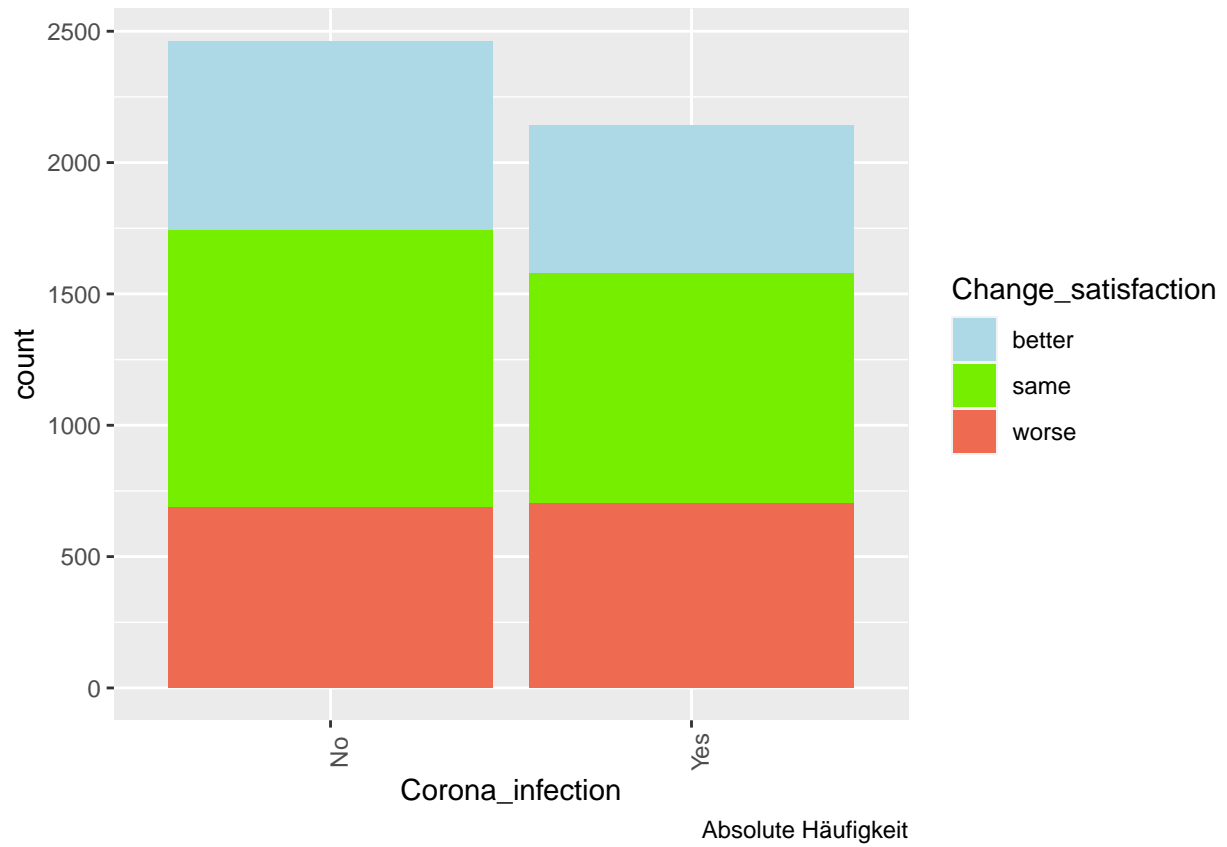
Die Änderung der Arbeitsituation hat keinen (prozentual gesehen) Einfluss



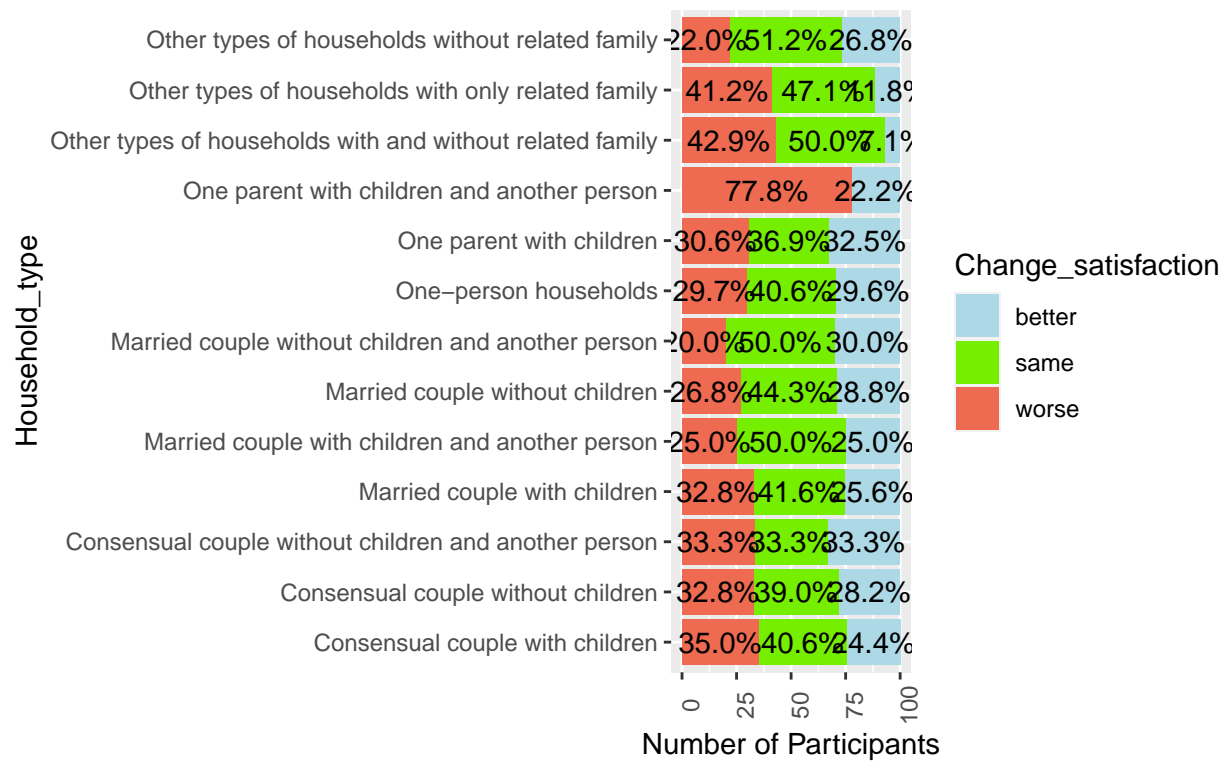
Untersuchung des Einflusses auf die Zielgrösse



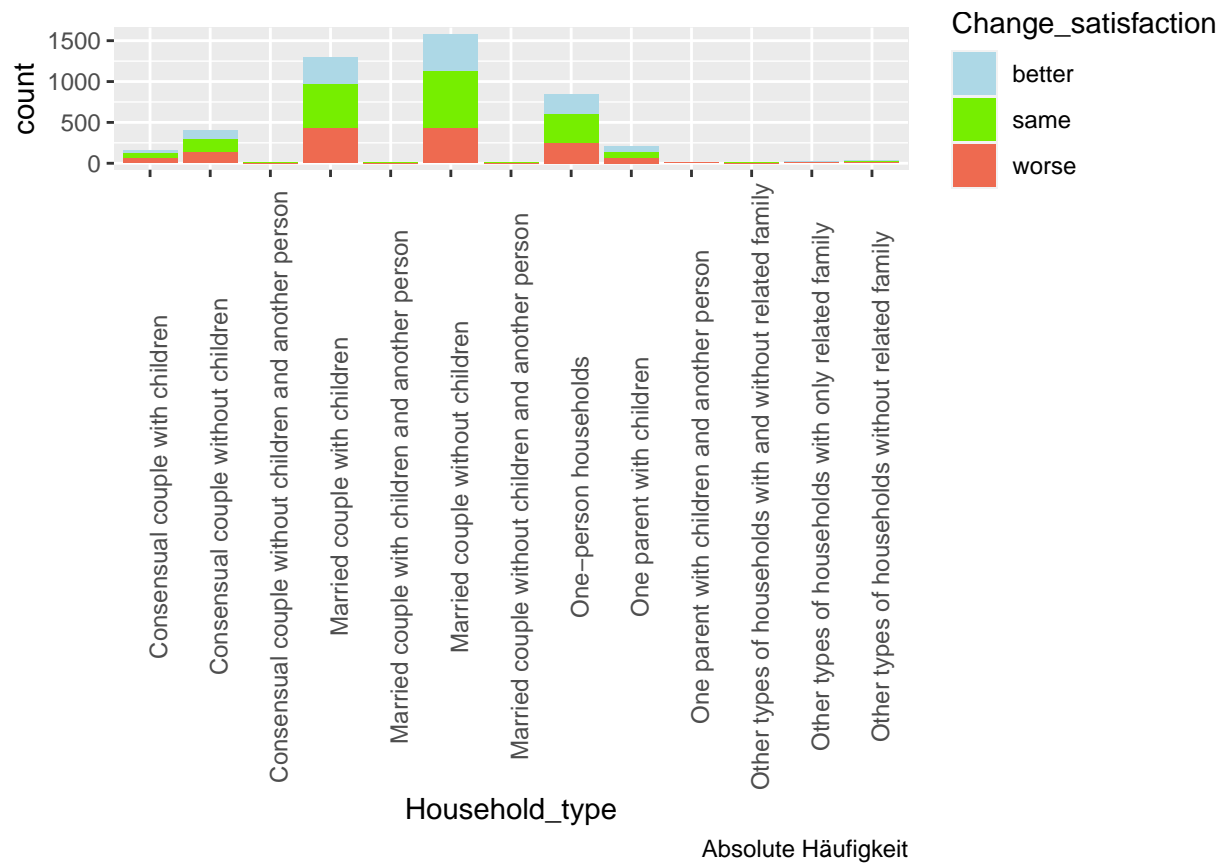
Menschen mit einer Corona-Infektion büssten in Lebenszufriedenheit ein



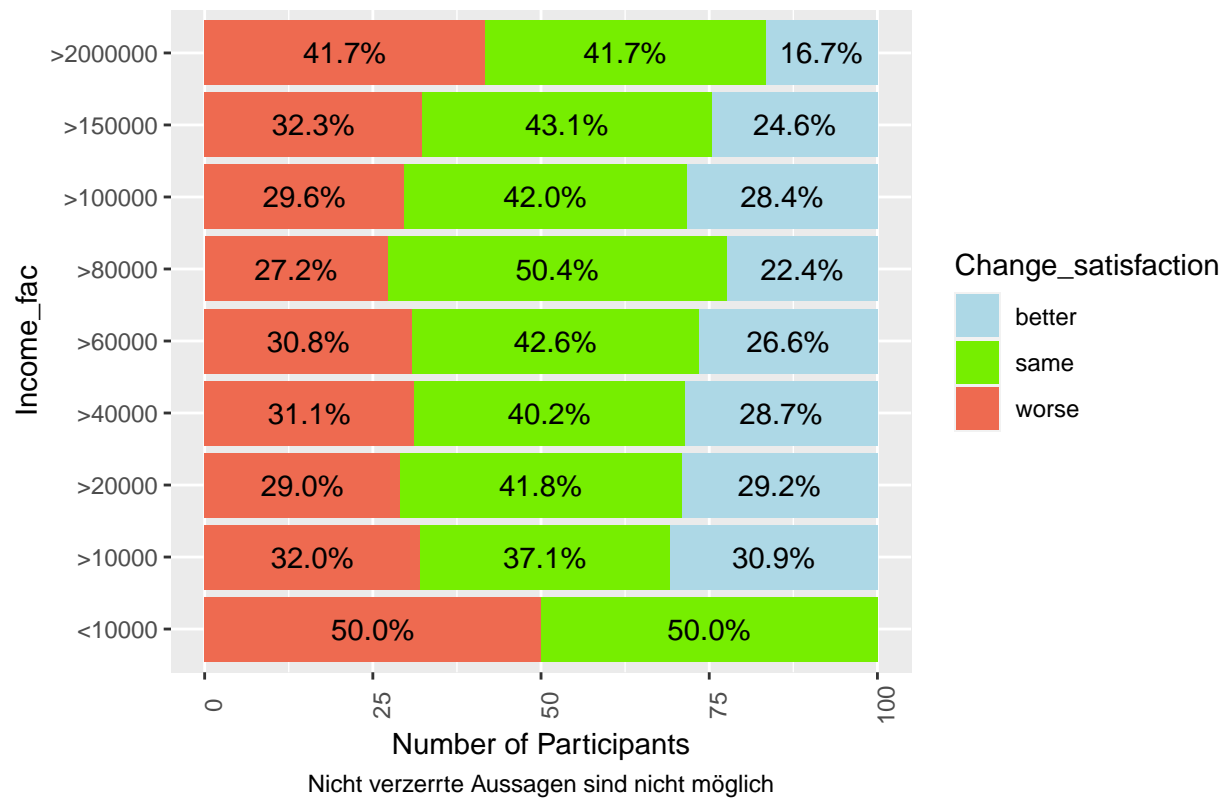
Untersuchung des Einflusses auf die :

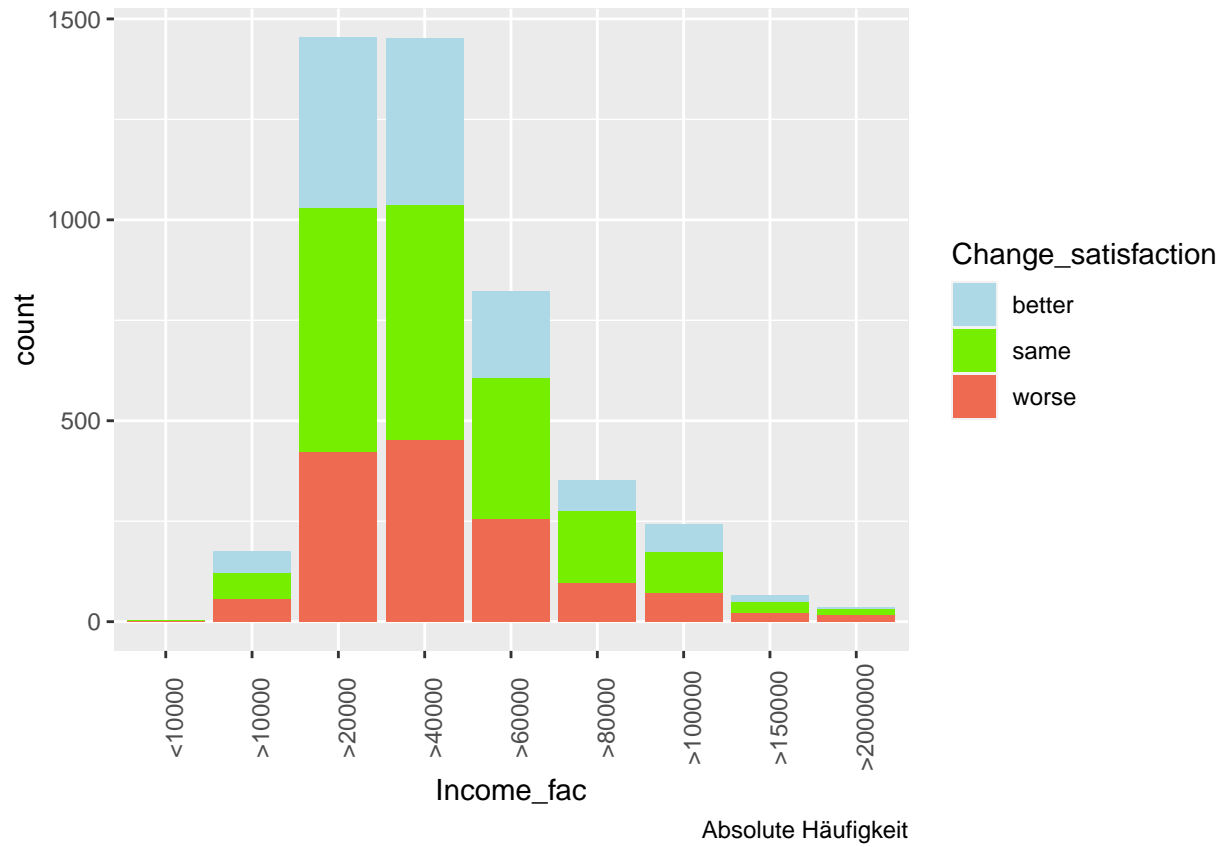


Veränderung in diesem Vergleich zu erkennen sind sinnlos

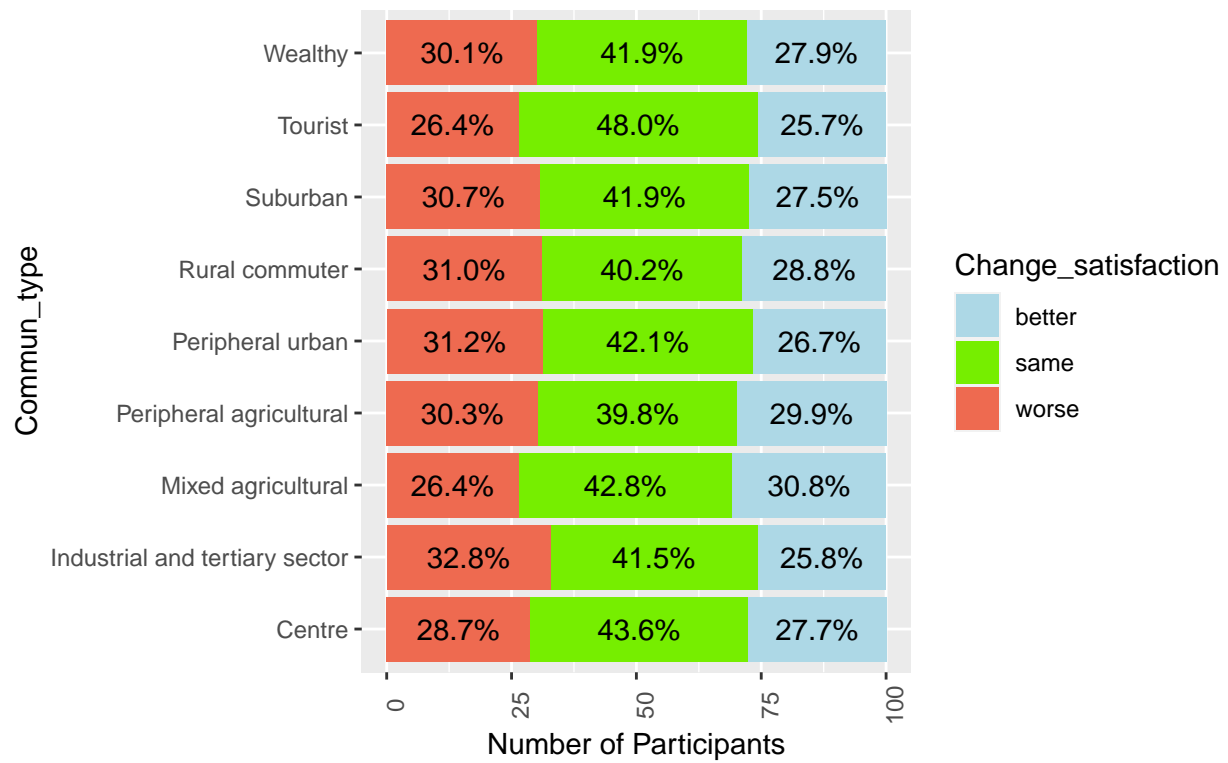


Untersuchung des Einflusses auf die Zielgrösse

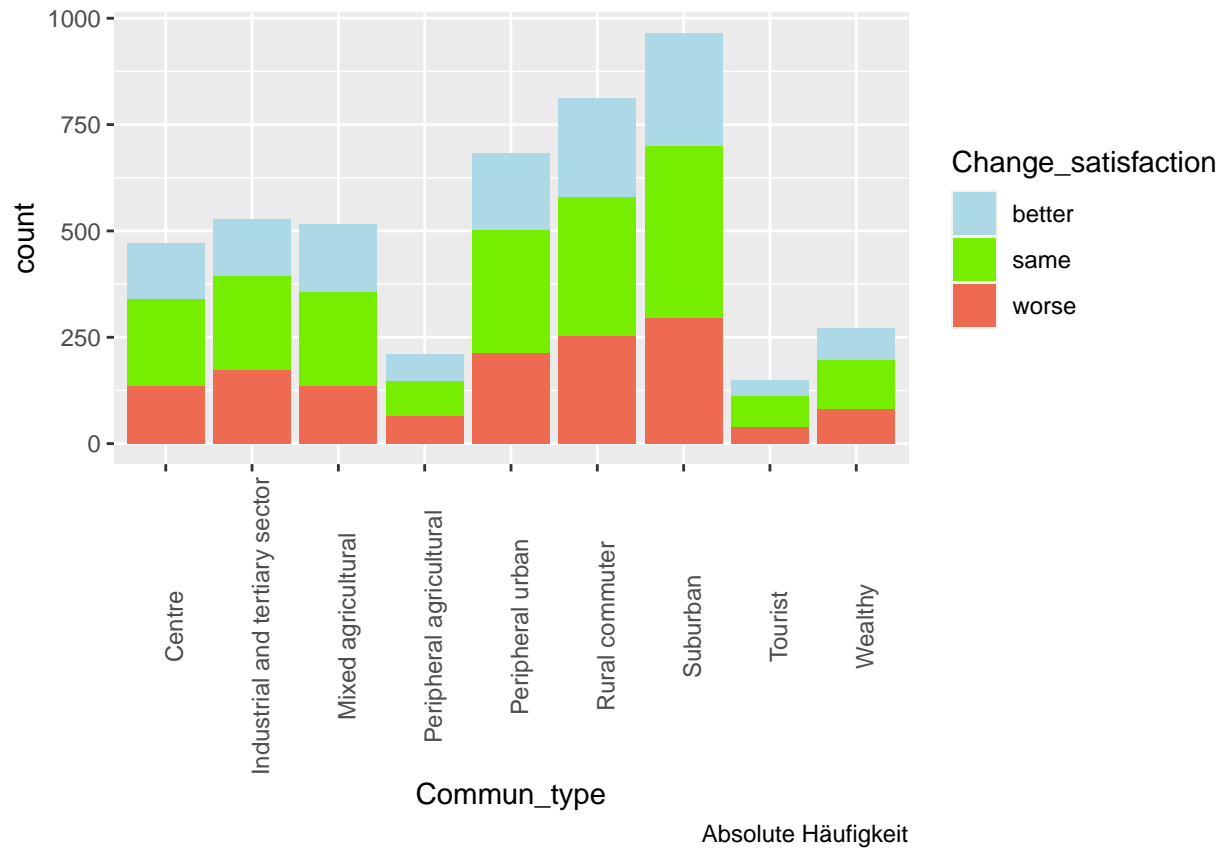




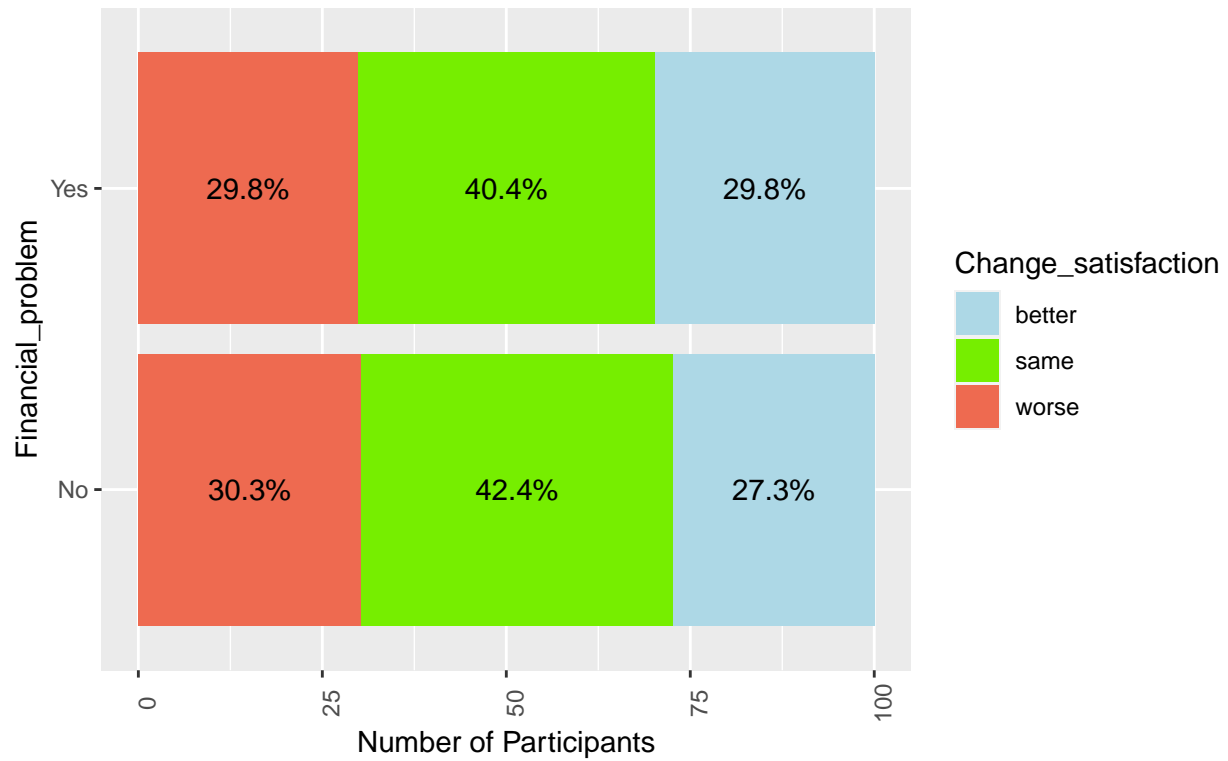
Untersuchung des Einflusses auf die Zielgrösse



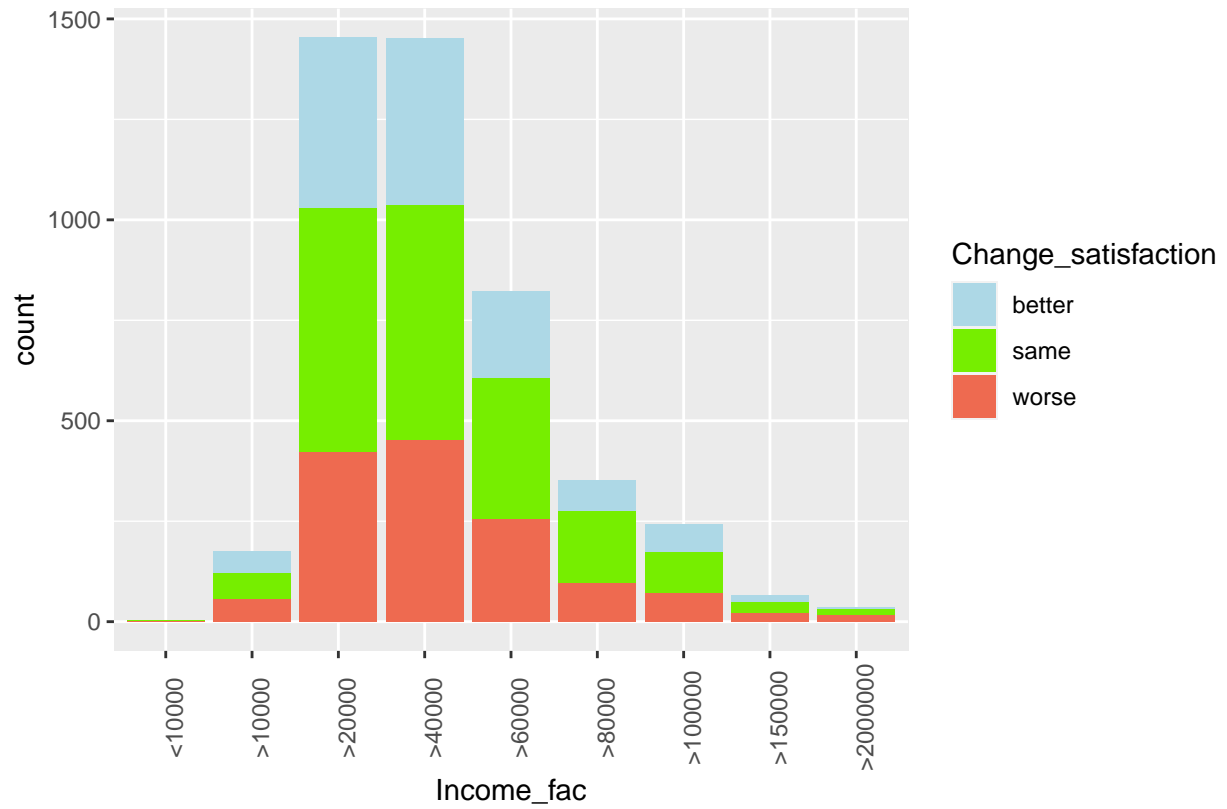
Der Gemeindetyp 'Industrie' hat die grösste negative Veränderung im Bereich der Lebenszufriedenheit



Untersuchung des Einflusses auf die Zielgrösse

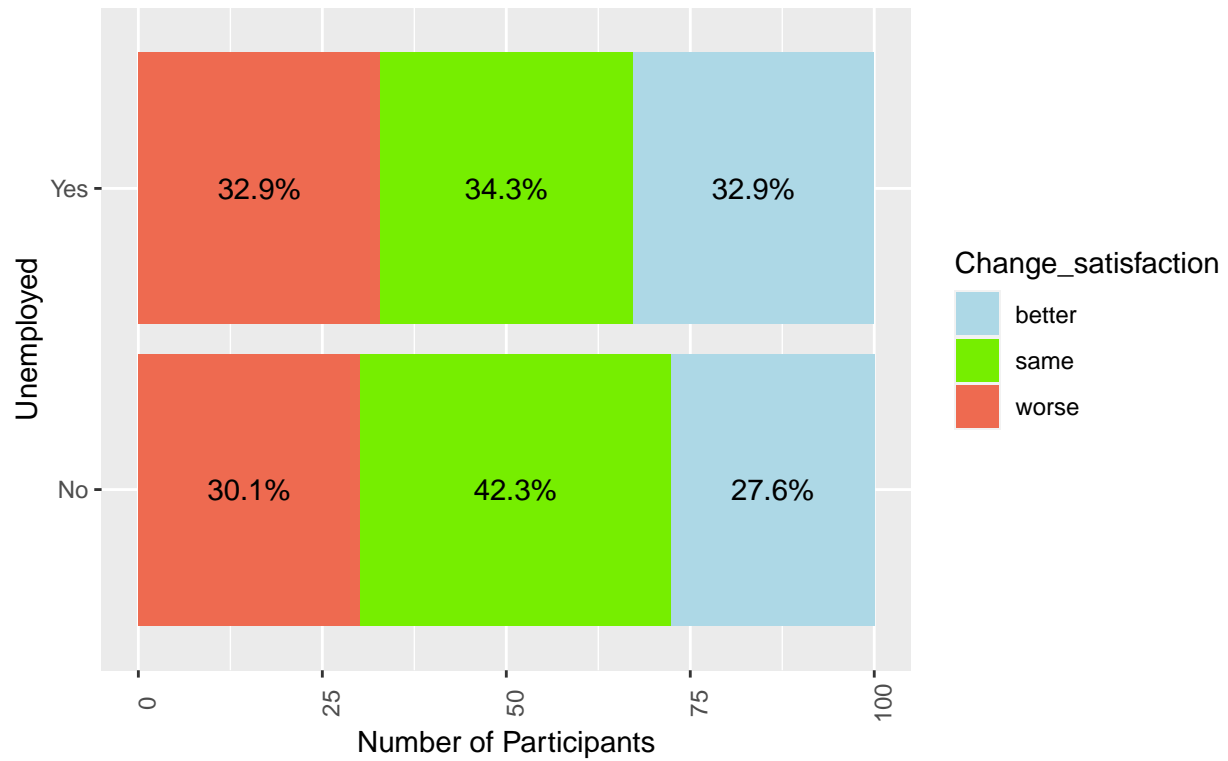


:irstaunlicherweise erging es laut Graphik Leuten mit ehemaligen Finanziellen Problemen besser.

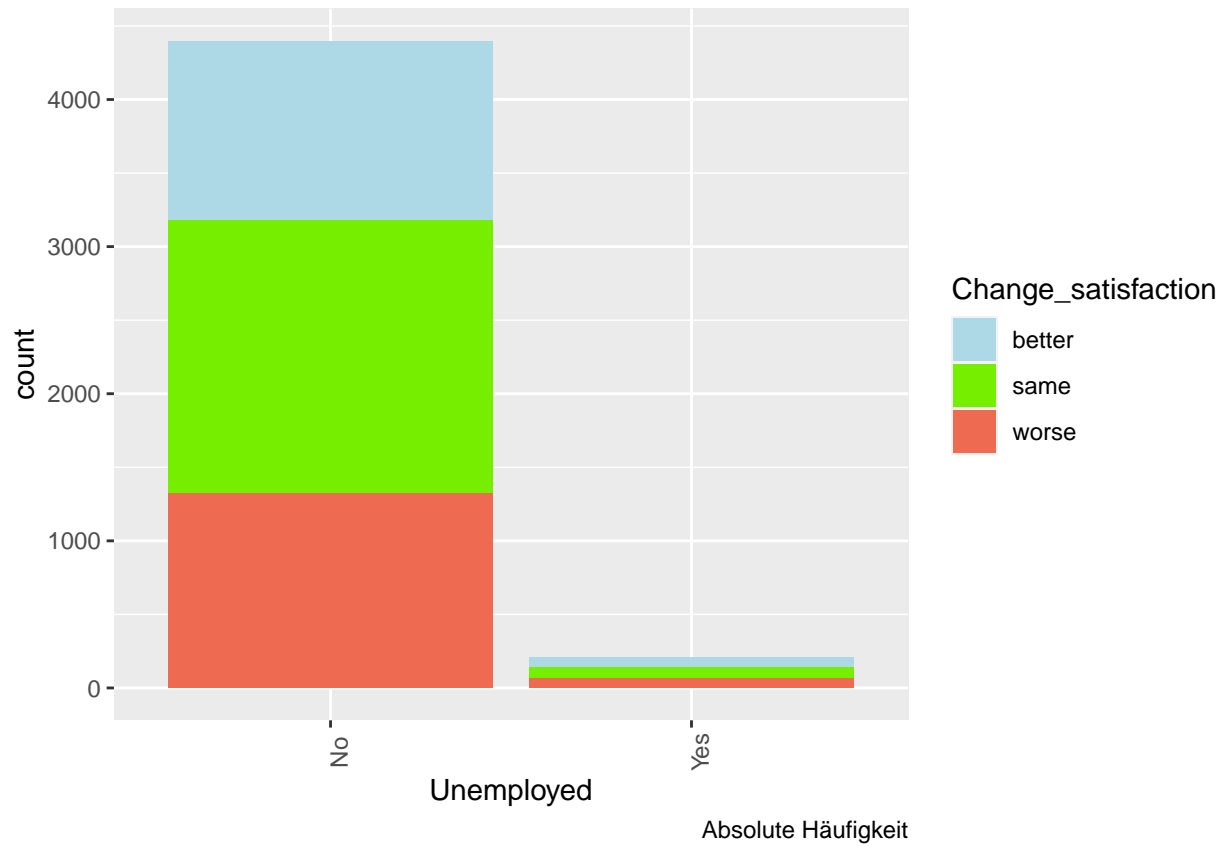


Absolute Häufigkeit

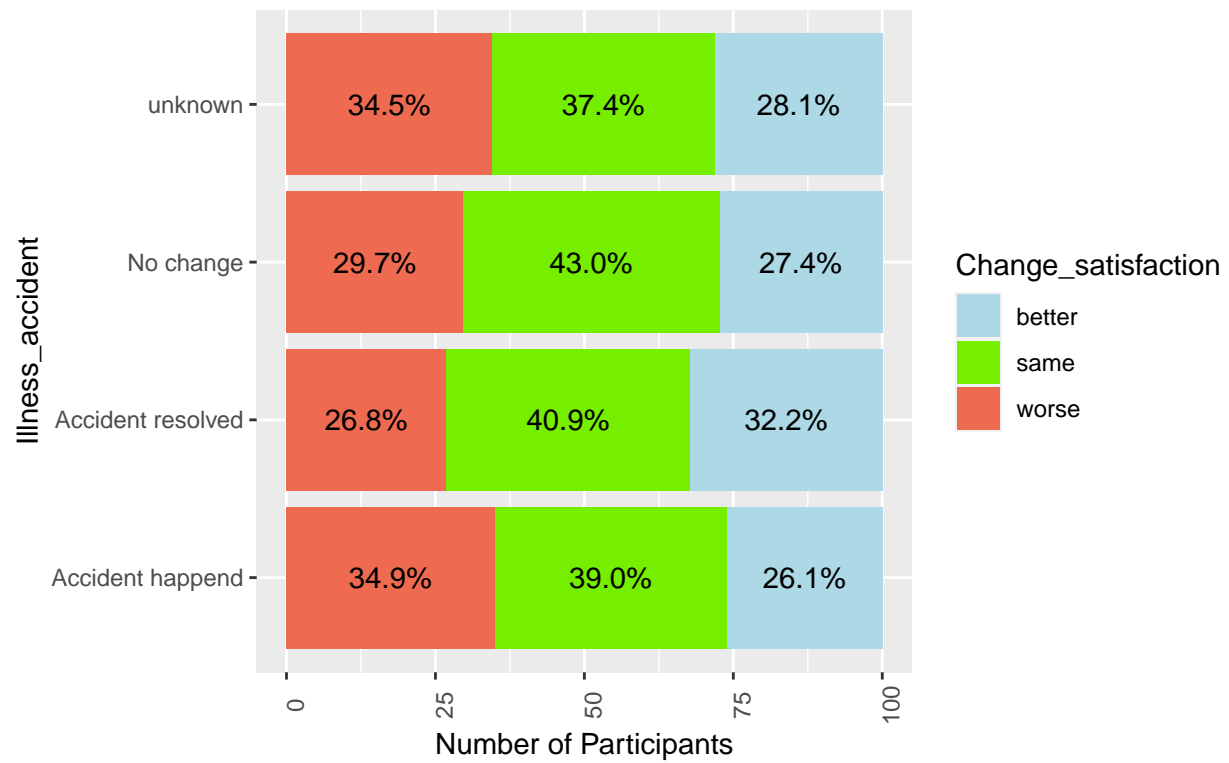
Untersuchung des Einflusses auf die Zielgrösse



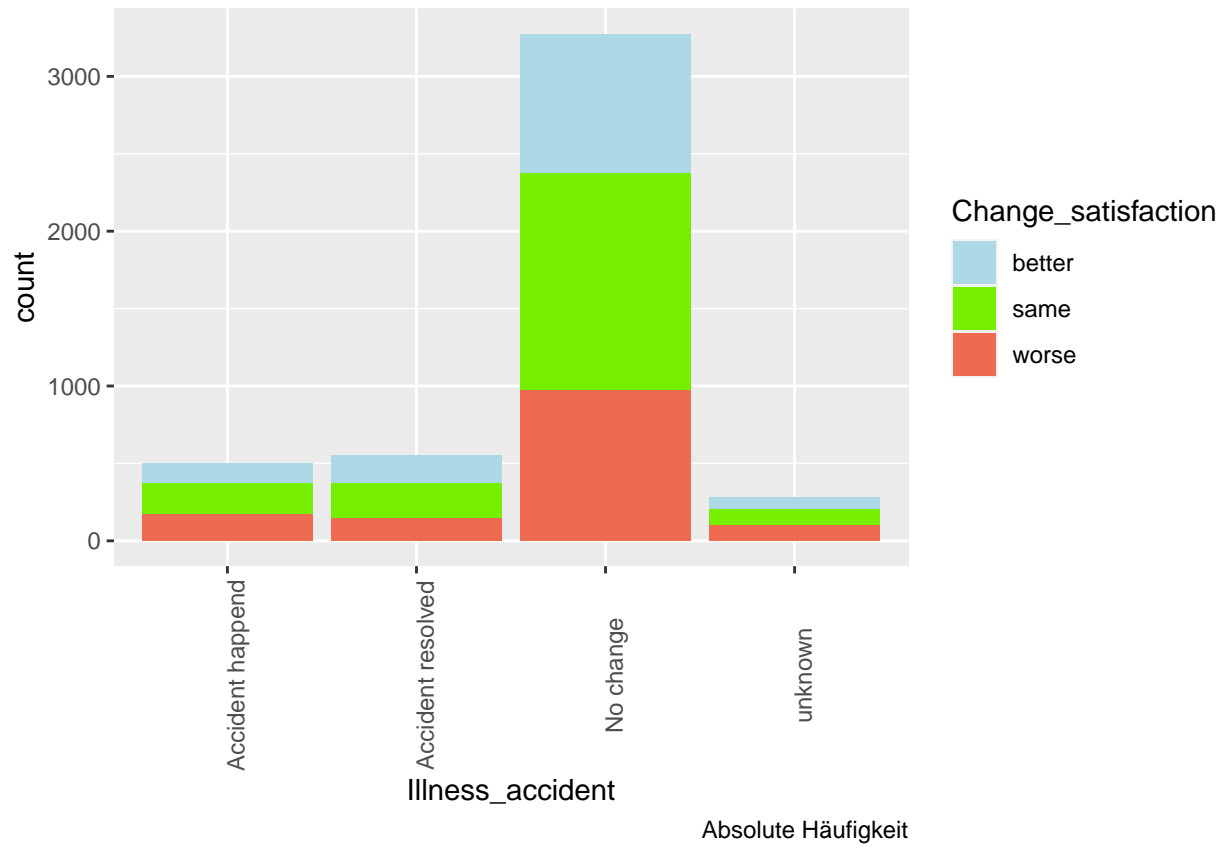
Laut Graphik hat die Arbeitslosigkeit keinen merklichen Einflusses auf die Zielgrösse



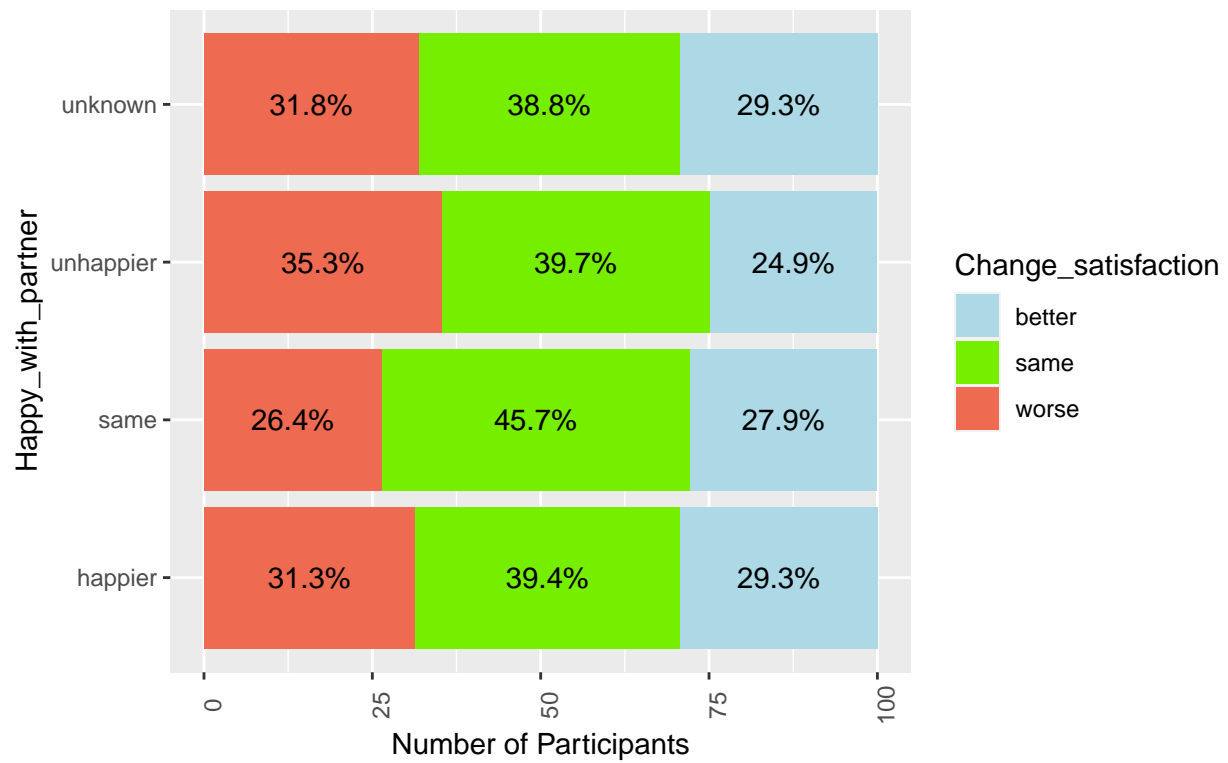
Untersuchung des Einflusses auf die Zielgrösse



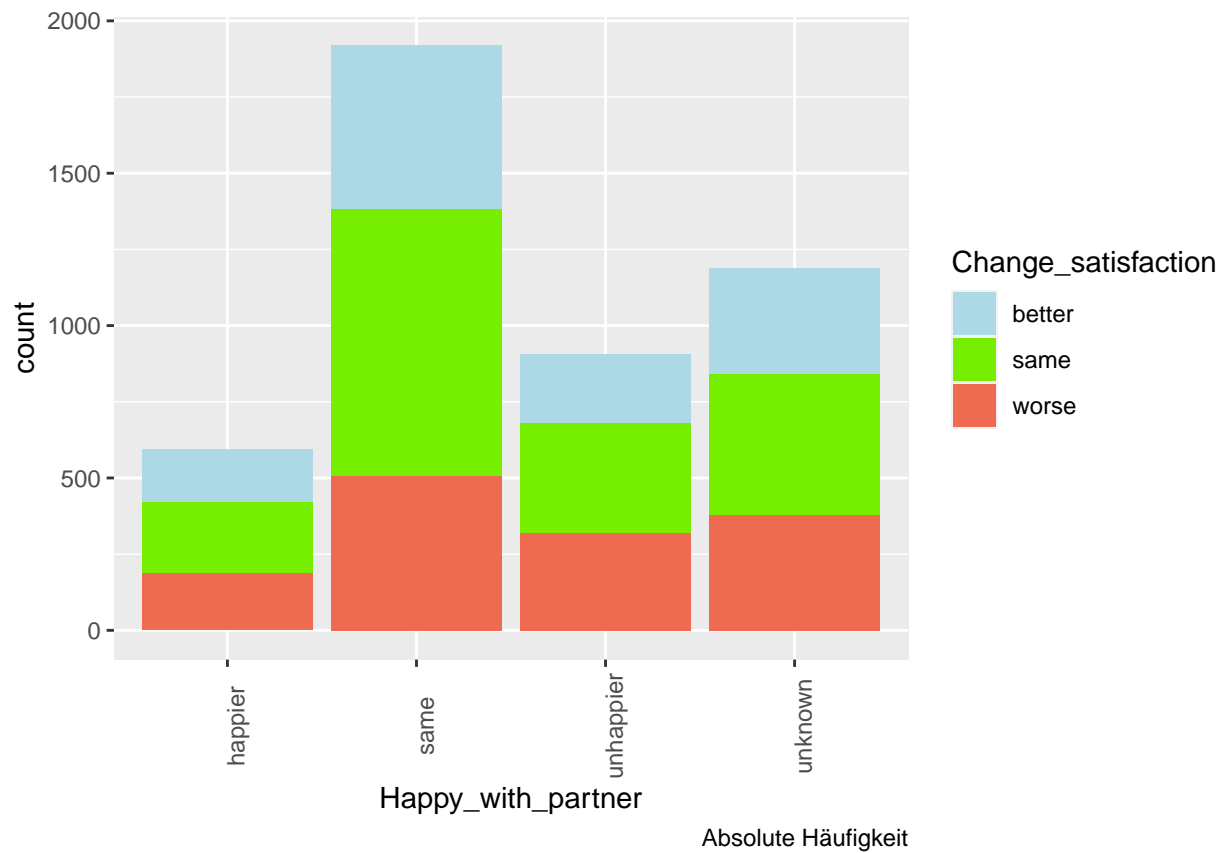
Personen, die einen Unfall oder Krankheit erlitten, erfuhren laut Graphik negative Auswirkungen in der Lebensqualität

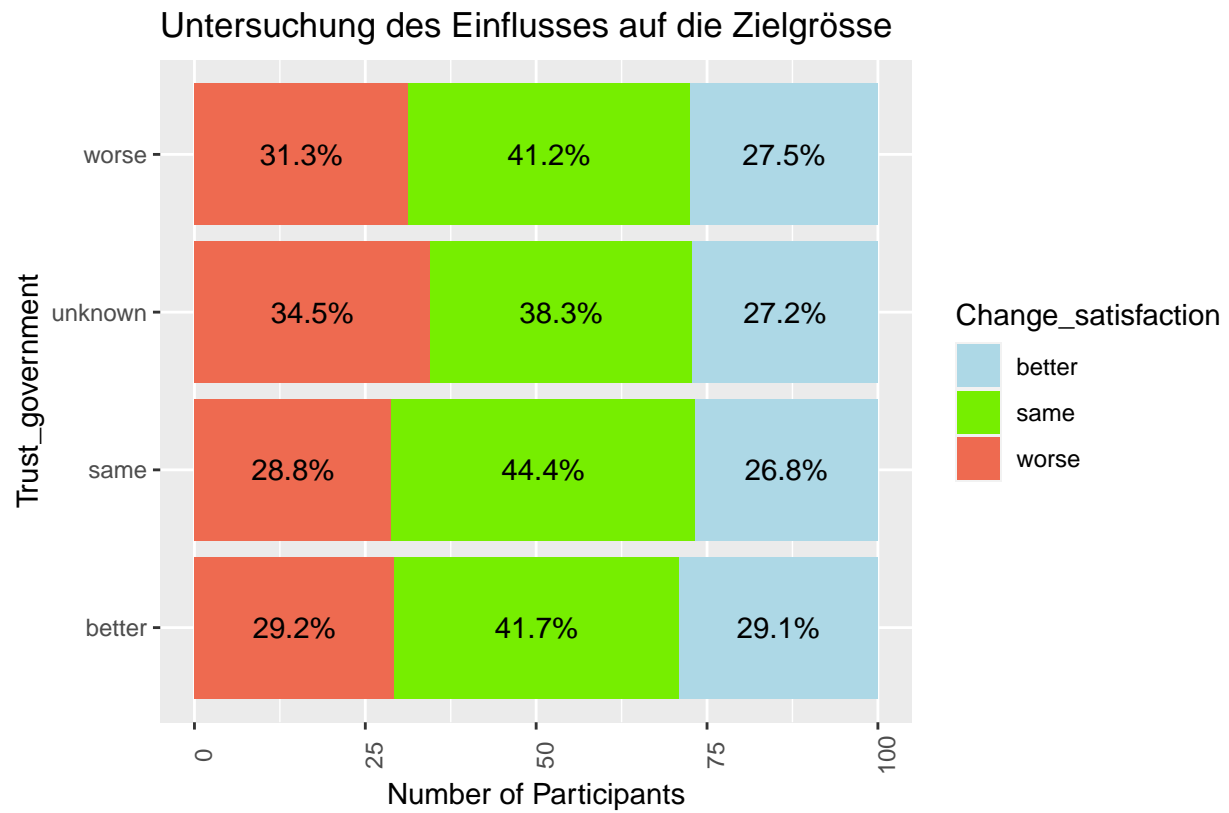


Untersuchung des Einflusses auf die Zielgrösse

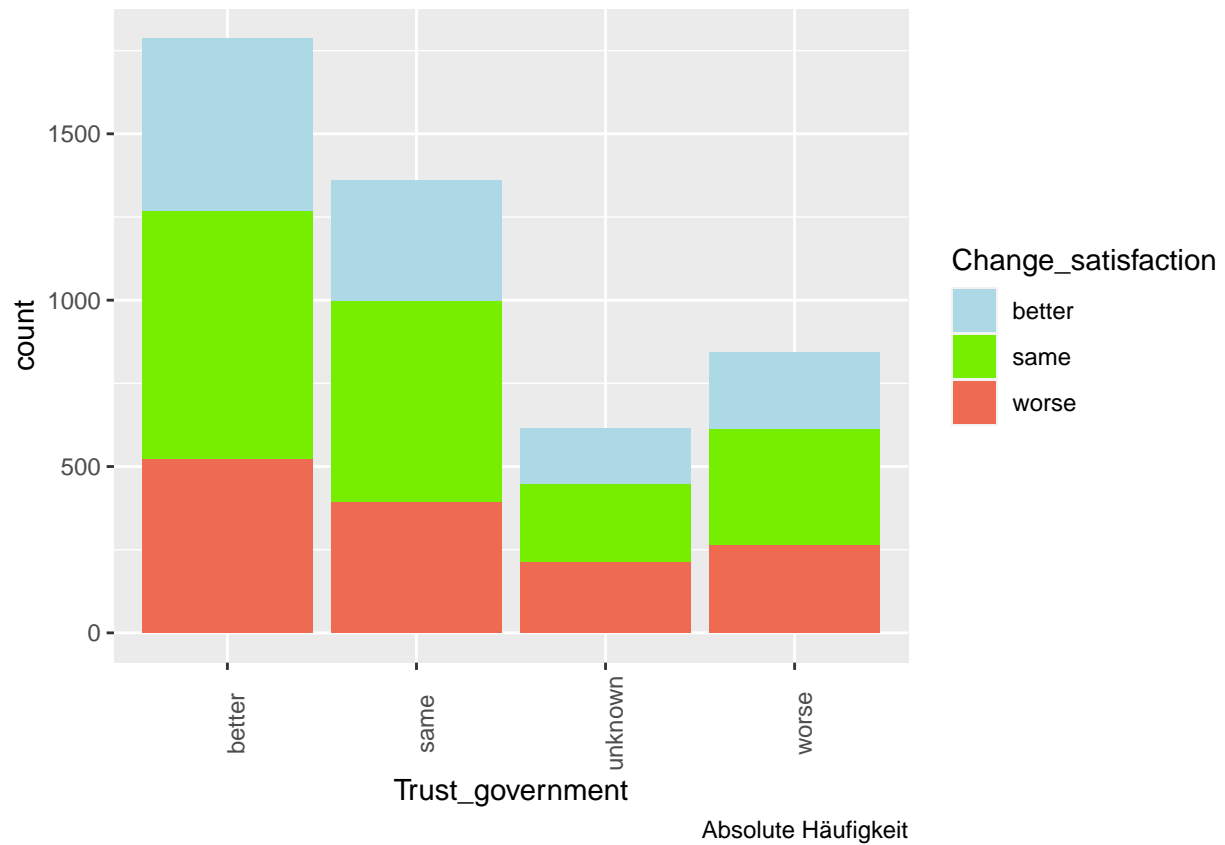


Leute, die unglücklich in der Beziehung waren, verloren laut Graphik auch an Lebensqualität.





schen mit Vertrauensverlust in die Regierung erlitten erhöhte negative Auswirkungen in der Lebensqualität



Schlusswort

Die Zielvariable hat sich im Verlauf der Pandemie nicht merklich verändert. Die Lebenszufriedenheit hat sich für etwa $\frac{2}{5}$ der Bevölkerung nicht verändert. Für die restlichen $\frac{3}{5}$ kam es entweder zu einer Verbesserung oder Verschlechterung der Lebenszufriedenheit. Man kann aber nicht den Schluss ziehen, dass es an der Pandemie lag.