

Grundlagen_der_Visualisierung_czlabinger

Zlabinger Christof

2023-11-13

Grundlagen der Visualisierung

Merkmale und Daten

Qualitativ

Alles was nur in Kategorien eingeteilt werden kann. Es gibt ein Maximum wie viele es gibt. Es kann keine "Zwischen Kategorien geben" z.B.: 1, 1.5, 2

nominal Keine Vordefinierte Anordnung. z.B.: Geschlecht, Farben

ordinal Klare Ordnung z.B.: Noten Es kann nicht mit ihnen gerechnet werden. Selbe Differenzen können verschieden gewichtet sein.

Quantitativ

Können gemessen und verglichen werden. Es kann mit ihnen gerechnet werden.

intervallskalierte Differenzen können ausgerechnet werden aber keine Vielfachheiten interpretieren kann. z.B.: Das Jahr 2000 ist später als das Jahr 1000 aber das Jahr 1000 ist nicht das Doppelte vom Jahr 500. Mit dem Jahr -500 kann keine sinnvolle Berechnung durchgeführt werden.

Man kann nicht sagen -5°C ist das -1 fache von 5°C . Wenn die Werte 0 oder negative Zahlen beinhalten können kann nicht von Vielfachen gesprochen werden sondern nur von Differenzen. 5°C ist um 10°C mehr als -5°C

rationalskaliert Daten die nur positiv sein können. z.B.: Kelvin da es einen absoluten 0-Punkt hat der in der Natur nicht erreicht werden kann. 1000K sind das Doppelte von 500K

diskret z.B.: Zähler Variablen: 0,1,2,...

stetig Beliebige Werte innerhalb eines Intervalls, eingeschränkt durch die Messgenauigkeit.

Datenmatrix

In den Spalten stehen die Merkmale aka Variablen In den Zeilen stehen die Untersuchungseinheiten

Schätzung und Darstellung von qualitativer Messungen

absolute Häufigkeit

relative Häufigkeit = absolute/gesamt

Können als Balken- oder Tortendiagramm dargestellt werden

Nichkummulierte Darstellung so wie sie sind gegenüber

Längen können gut geschätzt werden. Winkel nicht. Vergleich von % kein Tortendiagramm Balkendiagramm besser geeignet. Tortendiagramm gut um zu sehen wann es eine Mehrheit gibt und um wie viel.

Kumulative Addieren alle Gruppen bis zu einer bestimmten auf. Wenn es eine zu Grunde liegende Anordnung gibt. Mehr Sinn für ordinal als nominal.

Eigenschaften von quantitativen Daten:

Lage - Modus Unimodalität ist die Voraussetzung für alle anderen Lageschätzer wie den Mittelwert. Bimodalität oder Multimodalität müssen die Daten in verschiedene "Lager" aufgeteilt werden, um sie besser zu beschreiben.

Streuung - Varianz & Standardabweichung Varianz: $\sigma^2 = \frac{1}{n} \sum (i, n, (\text{Datenpunkt} - \text{Mittelwert})^2)$ (nicht stabil)

Standardabweichung: $\sigma = \sqrt{\sigma^2}$ (nicht stabil)

Interquartildistanz: Distanz zwischen dem 25%-Quantil (1. Quartil) und dem 75%-Quantil (3. Quartil) (stabil)

Spannweite: Maximalwert - Minimalwert (nicht stabil)

Symmetrie und Schiefe

- Normalverteilung: Symmetrisch mit einem Gipfel.
- Gleichverteilung: Symmetrisch ohne Peak - alle Werte sind gleich hoch.
- Bimodal symmetrisch: Zwei Peaks; Mittelwert und Median sind nicht sinnvoll, wenn sie bei 0 sind.
- Linksschief: Die Daten steigen rechts im Graphen stark an und fallen links langsam ab.
- Rechtsschief: Die Daten steigen links im Graphen stark an und fallen rechts langsam ab.
- Gewicht an den Rändern (Kurtosis)
 - Gibt an, wie stark Daten von ihrem Zentrum abweichen; erkennbar am Verschieben des Medians bei linksschiefen oder rechtsschiefen Daten.

Quantilschätzer

Gibt einen Zahlenwert für einen beliebigen Wert σ zwischen 0 und 1 aus, der dem Anteil aller Messwerte im Datensatz entspricht. Der Median ist beispielsweise das 50%-Quantil. Das Minimum (0%-Quantil) und Maximum (100%-Quantil) sind ebenfalls enthalten. Quantilschätzer geben an, wie viele Daten unbrauchbar sein können, aber dennoch funktionieren. Zum Beispiel können beim Median 50% der Daten unbrauchbar sein, und das Ergebnis ist dennoch sinnvoll. Der Anteil wird durch σ angegeben.

Dichtefunktion

Eine Dichtefunktion ist eine Funktion, bei der die Summe der n -ten Potenz der Datenwerte die Fläche unter dem Graphen ist.

2. Video

Arithmetisches Mittel

Alle Werte aufaddieren und durch die Gesamtanzahl dividieren.

`mean()`

Median

Der Wert in der Mitte bei Geordneten Daten.

`median()`

Modus

Wert mit der grössten Häufigkeit

`mode()`

Schätzung und Darstellung von quantitativer Messungen

Zählvariablen

Nur in ganzen Messeinheiten Zählbar. Wie oft kommt welche Grundanzahl vor. Anhand dessen absolute und relative Häufigkeit. Auch aufsummierbar.

Stetige numerische Variablen

Nicht wahrscheinlich dass 2 mal der gleiche Wert vorkommt. Somit keine Kategorien bilden. z.B: Messung von Dichten. Zusammenfassung in Teilbereiche. (Intervalle) Mit Intervallen wieder als Balkendiagramm darstellbar. (Histogramm um zu sehen ob Daten Unimodal oder Multimodal sind)

Kumulative Verteilungsfunktion Hilft mögliche Ausreisser zu erkennen wenn sie weit Weg von den anderen Werten sind.

Boxplot Zusammenfassung der Quartile. Box läuft vom 1. bis zum 3. Quantil. (25%-75%) In der Box ist der Median 50% Quantil. Robust gegenüber Ausreissern. Ausserhalb der Box sind die "Whiskers". Interquartilsdistanz ist die Distanz von 1. Quantil bis zum 3. Quantil. Laufen bis zur 1.5 fachen interquartilsdistanz. Schneidet aber beim letzten Wert innerhalb von 1.5 Box Längen ab.

Normal Q-Q Plot

Vergleicht Messwerte mit den Theoretischen Quantilen Wie kann ich feststellen ob es eine Normalverteilung gibt. Normalverteilt wenn die Punkte auf der Referenz liegen (Mit gewisser Toleranz). Wenn sehr viele Punkte gibt die nicht auf der Referenz liegen sind es keine Normalverteilte Daten.

Ca um den Median zentriert. Leichte Kante.

Bei schwerer Kante Boxplot nicht mehr gut um Ausreisser zu erkennen.