

Fachtagung Bodenindikatoren (17.04.2024)

# Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

Viola Rädle

Data Scientist

KI-Lab Leipzig, Umweltbundesamt

[ki-anwendungslabor@uba.de](mailto:ki-anwendungslabor@uba.de)

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Das KI-Lab am Umweltbundesamt

Anwendungslabor für Künstliche Intelligenz und Big Data

Informieren • Involvieren • Begeistern



**Ziele:** Forschung, Entwicklung, Kompetenzaufbau, Beratung, Wissenstransfer

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Welche Arten von Datenlücken gibt es?

Bsp: Umfrage zu Einkommen

**Umfrage**

Geschlecht:

m    w    d

Haarfarbe: \_\_\_\_\_

Alter: \_\_\_\_\_

Schuhgröße: \_\_\_\_\_

Einkommen: \_\_\_\_\_



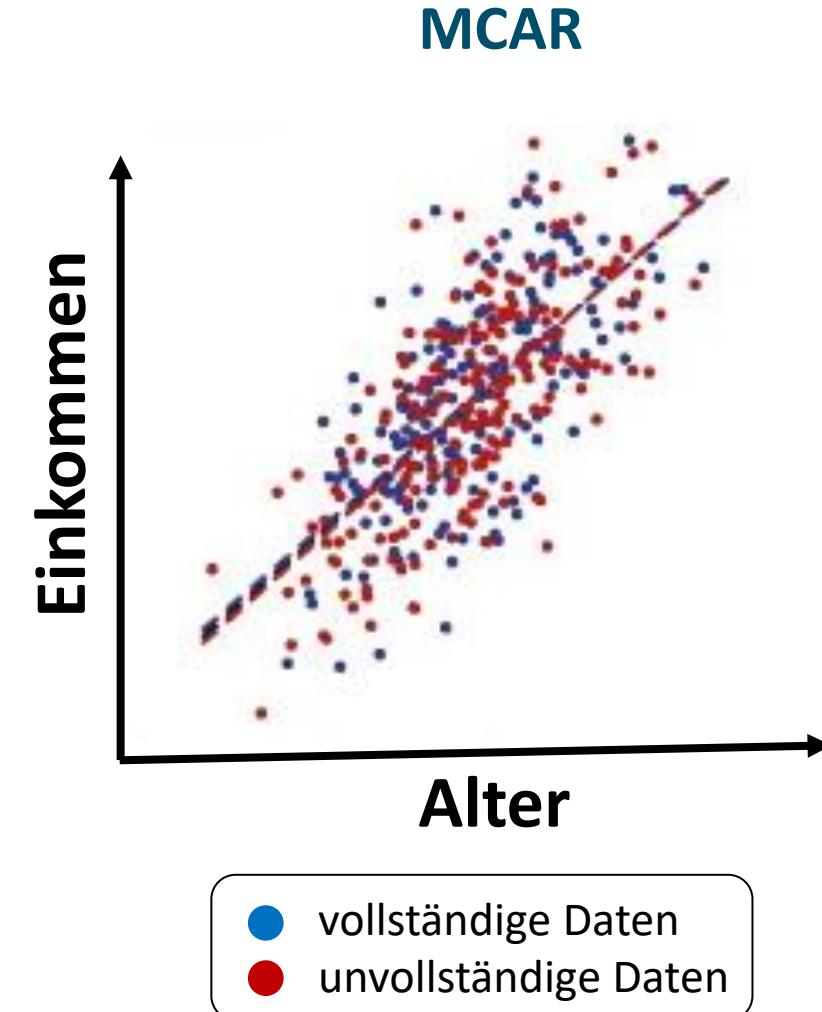
Quelle: <https://de.vecteezy.com/vektorkunst/554935-lust-auf-kugelschreiber-vektor-icon>

# Welche Arten von Datenlücken gibt es?

Bsp: Umfrage zu Einkommen

## Missing completely at random (MCAR):

- **kein Zusammenhang** mit erfassten Daten
- Bsp: technische Fehler, Übersehen eines Feldes



Quelle: <https://hushuli.github.io/Metabox-Blog.github.io/posts/2018-11-08-missing-value/>

# Welche Arten von Datenlücken gibt es?

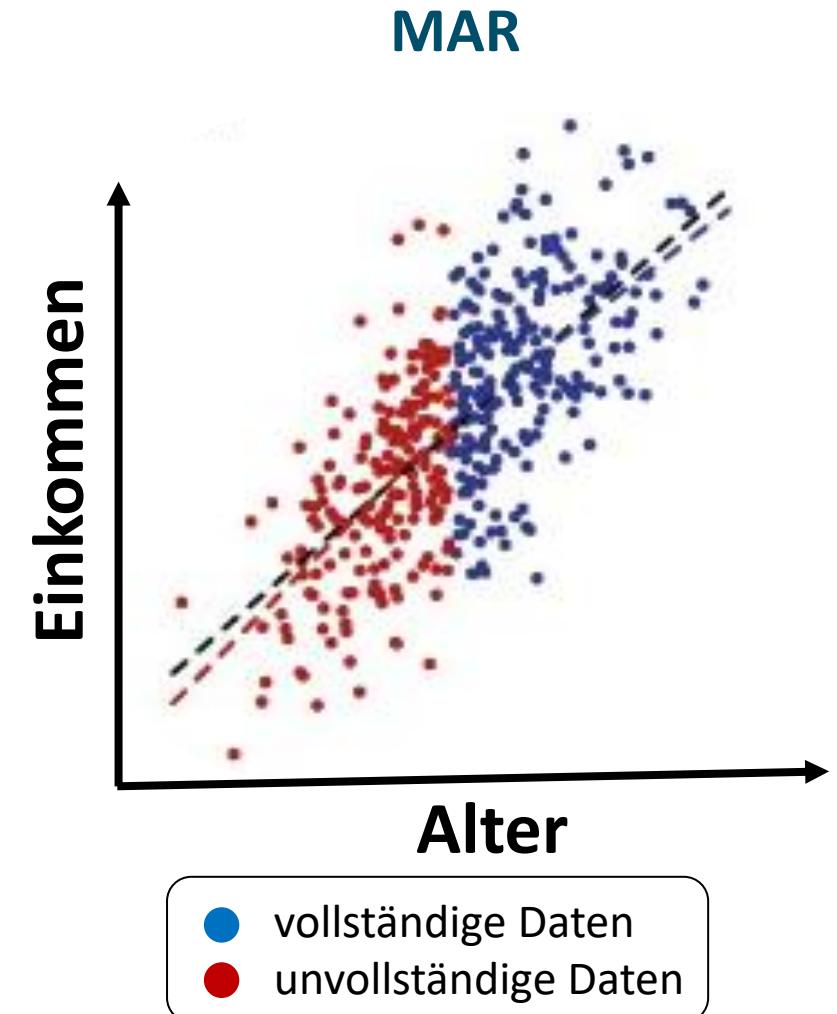
Bsp: Umfrage zu Einkommen

## Missing completely at random (MCAR):

- **kein Zusammenhang** mit erfassten Daten
- Bsp: technische Fehler, Übersehen eines Feldes

## Missing at random (MAR):

- Zusammenhang mit **anderer Variable**
- Bsp: Jüngere Menschen wollen keine Auskunft über ihr Einkommen geben



Quelle: <https://hushuli.github.io/Metabox-Blog.github.io/posts/2018-11-08-missing-value/>

# Welche Arten von Datenlücken gibt es?

Bsp: Umfrage zu Einkommen

## Missing completely at random (MCAR):

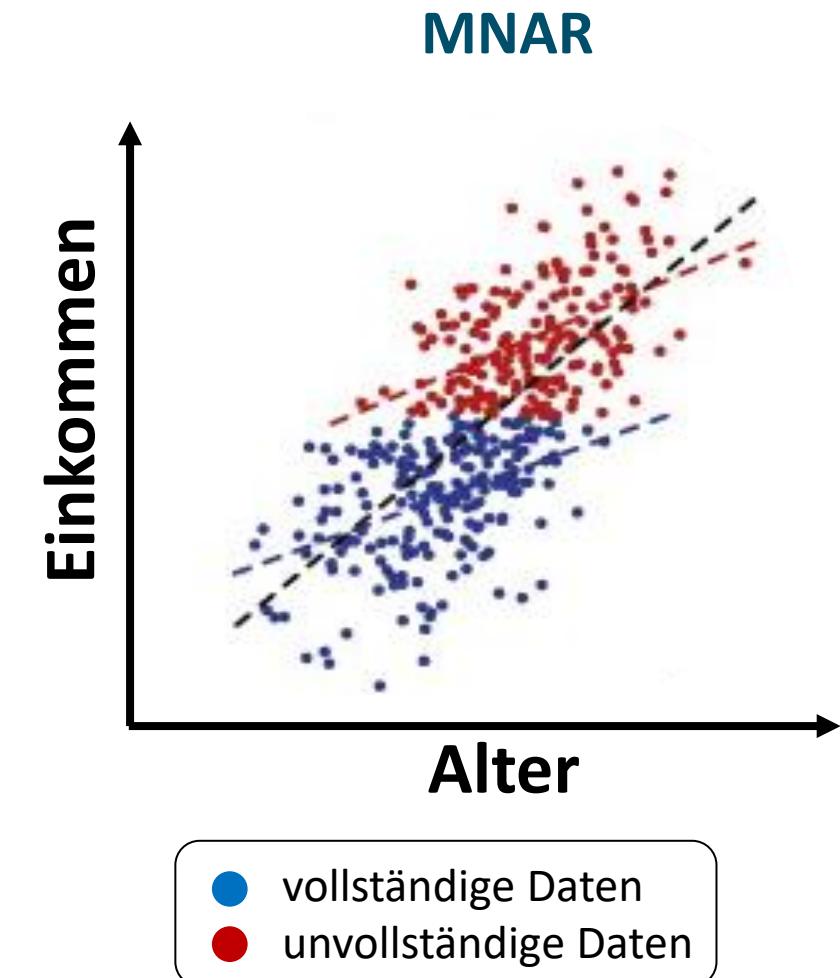
- **kein Zusammenhang** mit erfassten Daten
- Bsp: technische Fehler, Übersehen eines Feldes

## Missing at random (MAR):

- Zusammenhang mit **anderer Variable**
- Bsp: Jüngere Menschen wollen keine Auskunft über ihr Einkommen geben

## Missing not at random (MNAR):

- Zusammenhang mit **lückenhafter Variable**
- Bsp: Menschen mit höherem Einkommen wollen keine Auskunft geben



# Welche Arten von Datenlücken gibt es?

## Beispiele aus der Umweltforschung

### **Missing completely at random (MCAR):**

- technische Probleme / Stromausfall
- finanzieller Aufwand der Messung
- Messpersonal erkrankt



### **Missing at random (MAR):**

- Monitoring der Wasserqualität in Industrienähe nicht möglich (Gelände abgesperrt)
- weniger Messstationen in schwer zugänglichen Gebieten als in Zivilisationsnähe
- Messfahrten aufgrund schlechter Witterungsbedingungen abgebrochen

### **Missing not at random (MNAR):**

- Monitoring der Artenvielfalt: bedrohte Arten sind schwieriger zu finden
- Sensor kann Werte oberhalb des Messbereichs nicht mehr erfassen
- Radioaktivitätsmessungen: Probenahme in kontaminierten Gebieten zu gefährlich

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# data imputation



## data dropping

lösche lückenhafte  
Messungen

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Der einfachste Weg: Data dropping

Lösche alle Messungen mit fehlenden Werten

Datum	Temperatur in °C	Wassergehalt in %	Leitfähigkeit in dS/m	pH	Stickstoff in %
20.03.2023		22.6	1.2	4.9	0.09
01.04.2023	12.5	20.4	1.8	5.2	0.07
13.04.2023	14.2	25.3		6.3	0.10
02.05.2023	15.6	33.7	1.9	5.0	
24.05.2023	16.4	30.1	2.1	5.7	0.12
	19.0	28.3	2.3	6.5	0.14
02.09.2023	18.7	35.8	2.6	6.8	0.11
17.09.2023	17.2		2.0	6.2	0.06

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Der einfachste Weg: Data dropping

Lösche alle Messungen mit fehlenden Werten

Datum	Temperatur in °C	Wassergehalt in %	Leitfähigkeit in dS/m	pH	Stickstoff in %
20.03.2023		22.6	1.2	4.9	0.09
01.04.2023	12.5	20.4	1.8	5.2	0.07
13.04.2023	14.2	25.3		6.3	0.10
02.05.2023	15.6	33.7	1.9	5.0	
24.05.2023	16.4	30.1	2.1	5.7	0.12
	19.0	28.3	2.3	6.5	0.14
02.09.2023	18.7	35.8	2.6	6.8	0.11
17.09.2023	17.2		2.0	6.2	0.06

# data imputation

## single imputation

fülle die Lücken mit  
anderem Wert

## data dropping

lösche lückenhafte  
Messungen

## univariat

- Hot-Deck
- Mittelwert/Median
- Regression

## multivariat

- k-nearest-neighbors
- Random Forest
- Support Vector Machine
- multiple Regression
- Hauptkomponentenanalyse

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Hot-Deck-Imputation

Fülle die Datenlücke mit einem zufälligen Wert aus der Spalte

Datum	Temperatur in °C	Wassergehalt in %	Leitfähigkeit in dS/m	pH	Stickstoff in %
20.03.2023		22.6	1.2	4.9	0.09
01.04.2023	12.5	20.4	1.8	5.2	0.07
13.04.2023	14.2	25.3		6.3	0.10
02.05.2023	15.6	33.7	1.9	5.0	
24.05.2023	16.4	30.1	2.1	5.7	0.12
	19.0	28.3	2.3	6.5	0.14
02.09.2023	18.7	35.8	2.6	6.8	0.11
17.09.2023	17.2		2.0	6.2	0.06

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

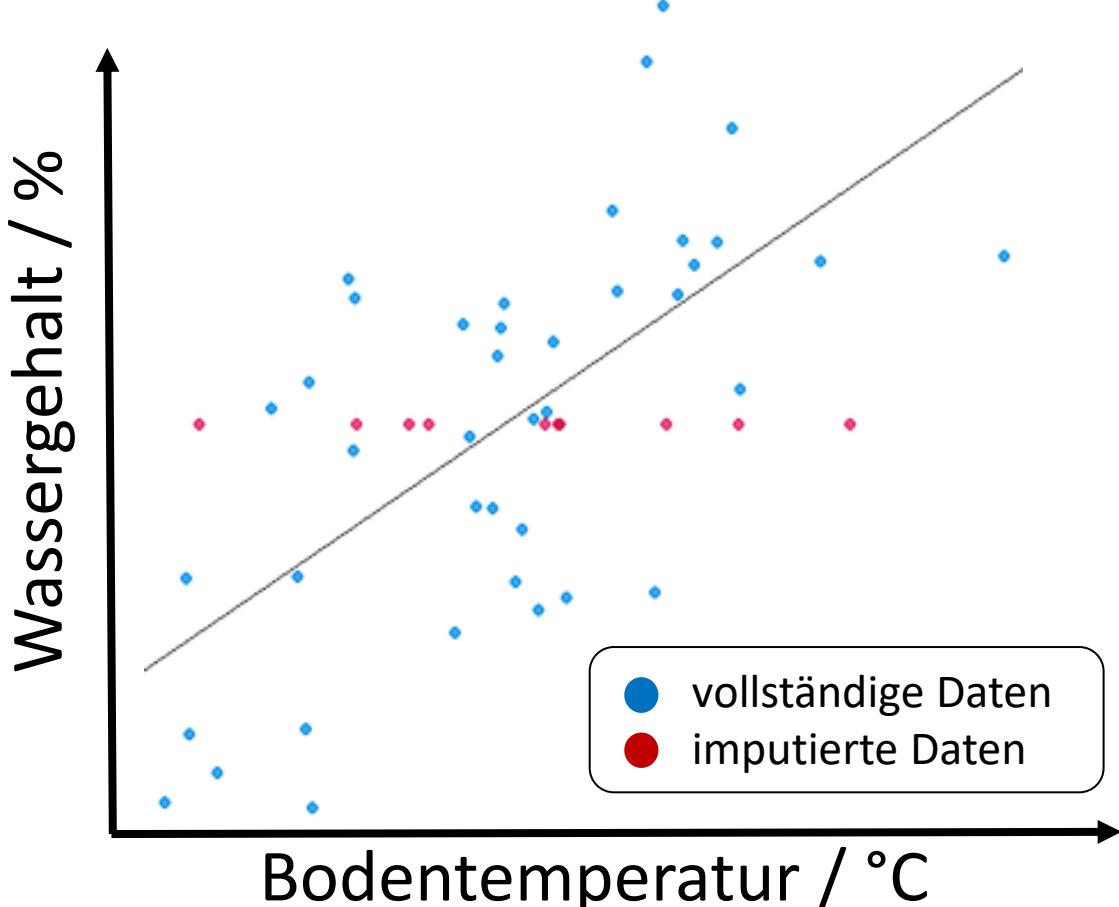
# Mean/Median Imputation

	Datum	Temperatur in °C	Wassergehalt in %	Leitfähigkeit in dS/m	pH	Stickstoff in %
	20.03.2023		22.6	1.2	4.9	0.09
	01.04.2023	12.5	20.4	1.8	5.2	0.07
	13.04.2023	14.2	25.3		6.3	0.10
	02.05.2023	15.6	33.7	1.9	5.0	
	24.05.2023	16.4	30.1	2.1	5.7	0.12
		19.0	28.3	2.3	6.5	0.14
	02.09.2023	18.7	35.8	2.6	6.8	0.11
	17.09.2023	17.2		2.0	6.2	0.06
<b>Mean</b>	<b>29.05.2023</b>	<b>16.2</b>	<b>28.0</b>	<b>2.0</b>		<b>0.1</b>
<b>Median</b>	<b>02.05.2023</b>	<b>16.4</b>	<b>28.3</b>	<b>2.0</b>		<b>0.1</b>

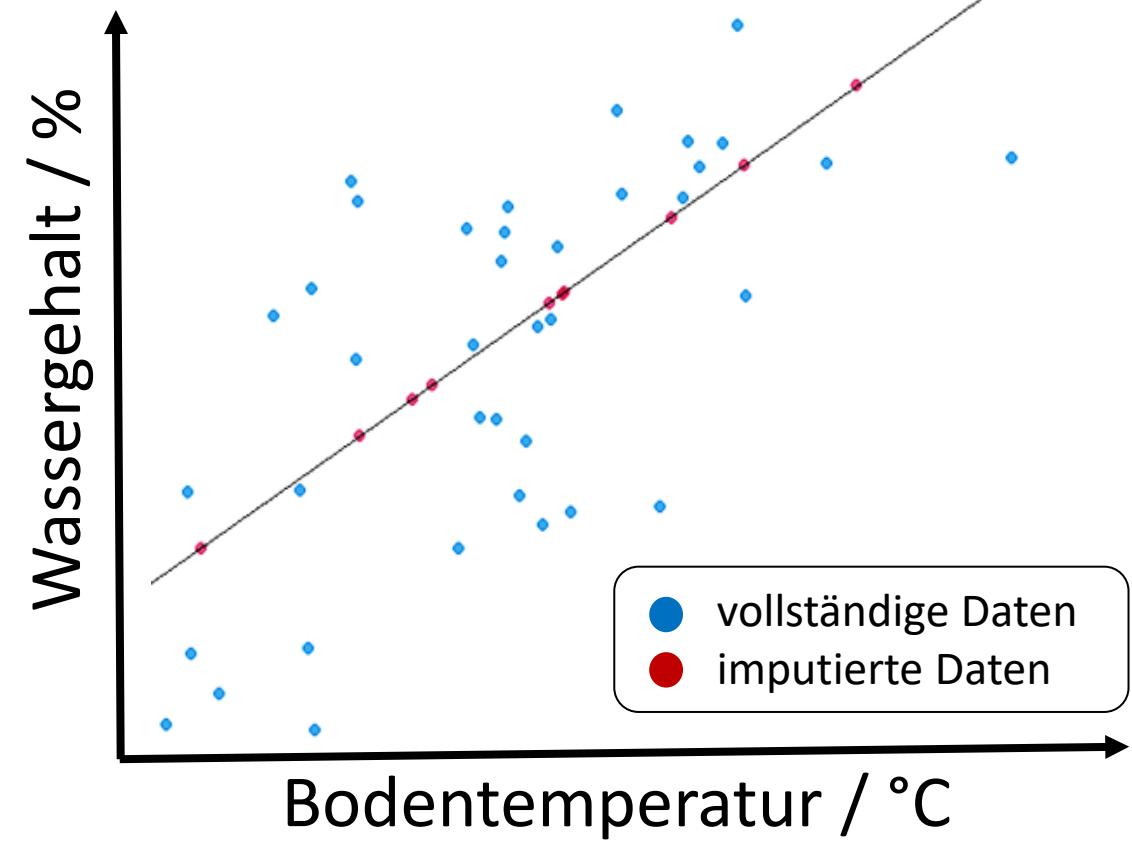
Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Imputation durch...

## ...den Mittelwert



## ...Regression



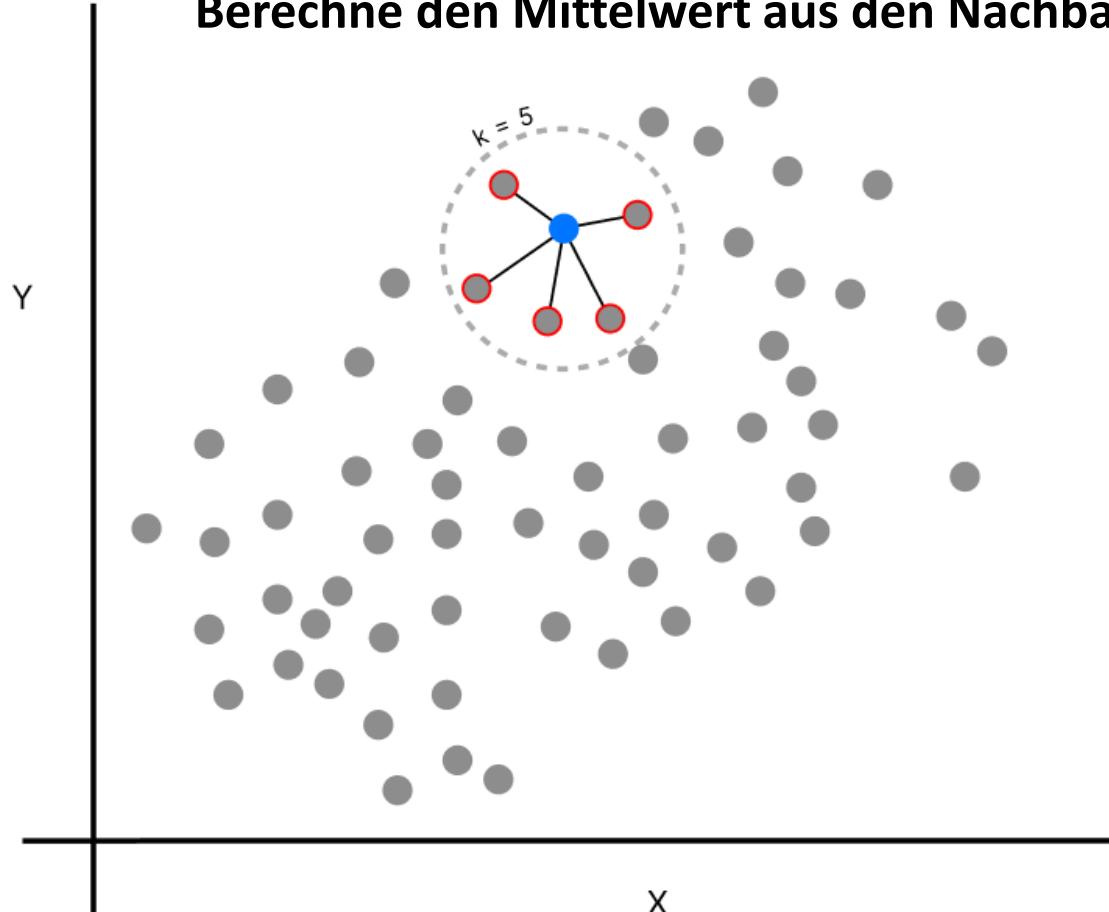
Quelle: <https://www.missingdata.nl/missing-data/missing-data-methods/imputation-methods/>

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Imputation durch...

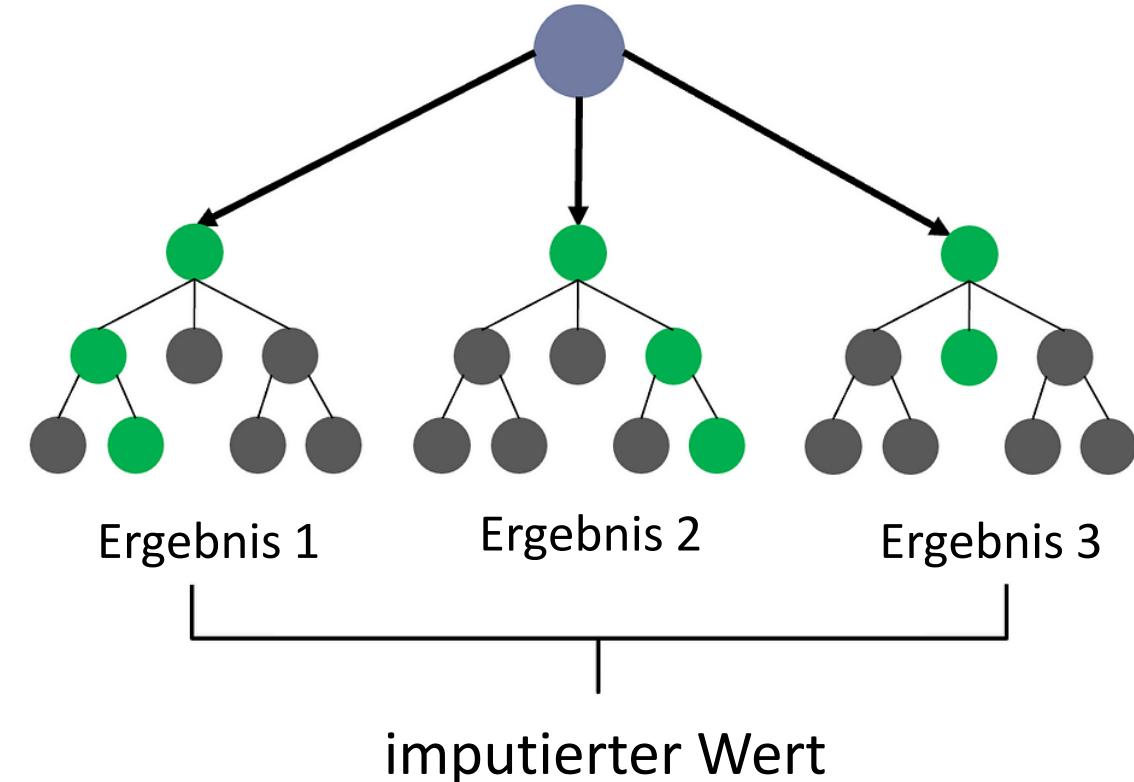
## k-nearest-neighbours

Berechne den Mittelwert aus den Nachbarn



## random forest

Mittle über mehrere Entscheidungsbäume



Quellen: <https://pub.towardsai.net/understanding-k-nearest-neighbors-a-simple-approach-to-classification-and-regression-e4b30b37f151>, <https://medium.com/@roiyeho/random-forests-98892261dc49>

# data imputation

## single imputation

fülle die Lücken mit  
anderem Wert

## data dropping

lösche lückenhafte  
Messungen

## multiple imputation

fülle Lücken +  
beachte Unsicherheiten

### univariat

- Hot-Deck
- Mittelwert/Median
- Regression

### multivariat

- k-nearest-neighbors
- Random Forest
- Support Vector Machine
- multiple Regression
- Hauptkomponentenanalyse

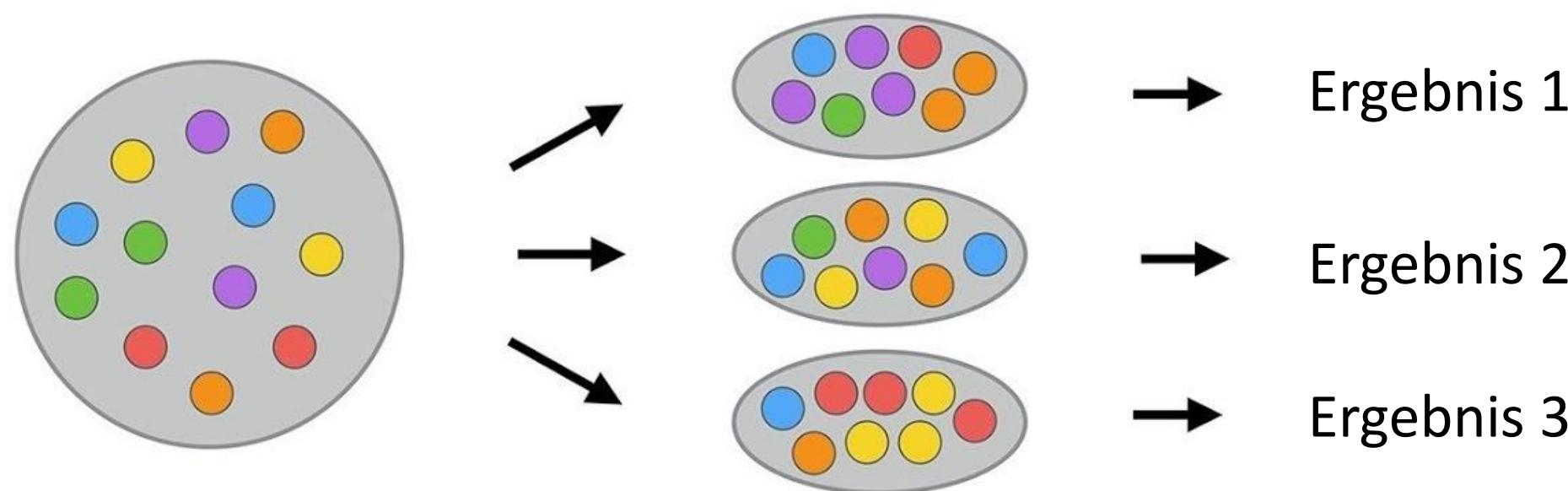
- Bootstrap
- Bayesianische Statistik
- Markov Chain Monte Carlo
- Expectation Maximization
- MICE (Multiple Imputation by Chained Equations)

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Imputation durch...

## Bootstrapping

Resampling zur Unsicherheitsabschätzung



Quelle: <https://www.youtube.com/watch?app=desktop&v=d3mcuJycJfl>

# data imputation

## single imputation

fülle die Lücken mit  
anderem Wert

## data dropping

lösche lückenhafte  
Messungen

## multiple imputation

fülle Lücken +  
beachte Unsicherheiten

### univariat

- Hot-Deck
- Mittelwert/Median
- Regression

### multivariat

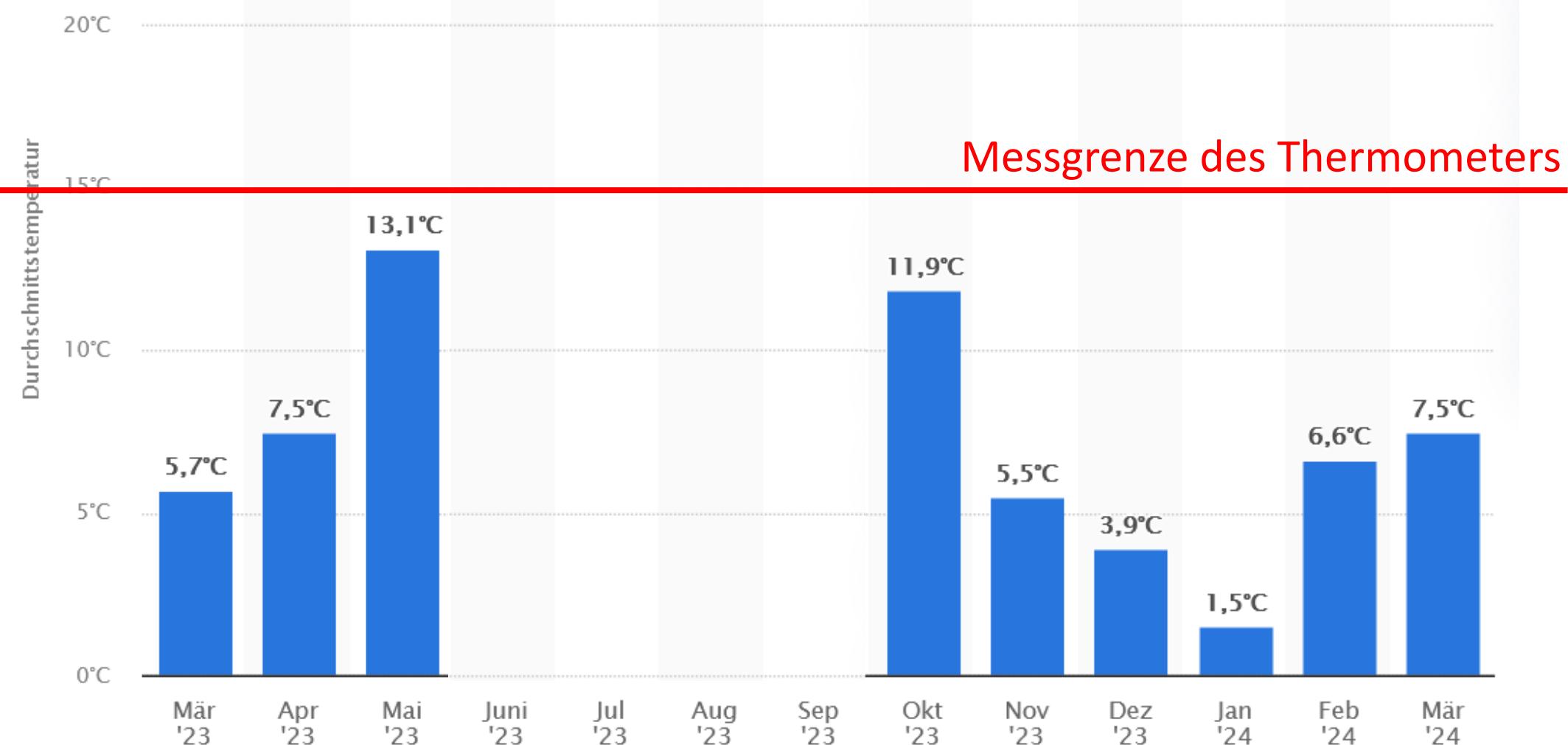
- k-nearest-neighbors
- Random Forest
- Support Vector Machine
- multiple Regression
- Hauptkomponentenanalyse

- Bootstrap
- Bayesianische Statistik
- Markov Chain Monte Carlo
- Expectation Maximization
- MICE (Multiple Imputation by Chained Equations)

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Limitierungen von Imputation

## Entzauberung von Künstlicher Intelligenz

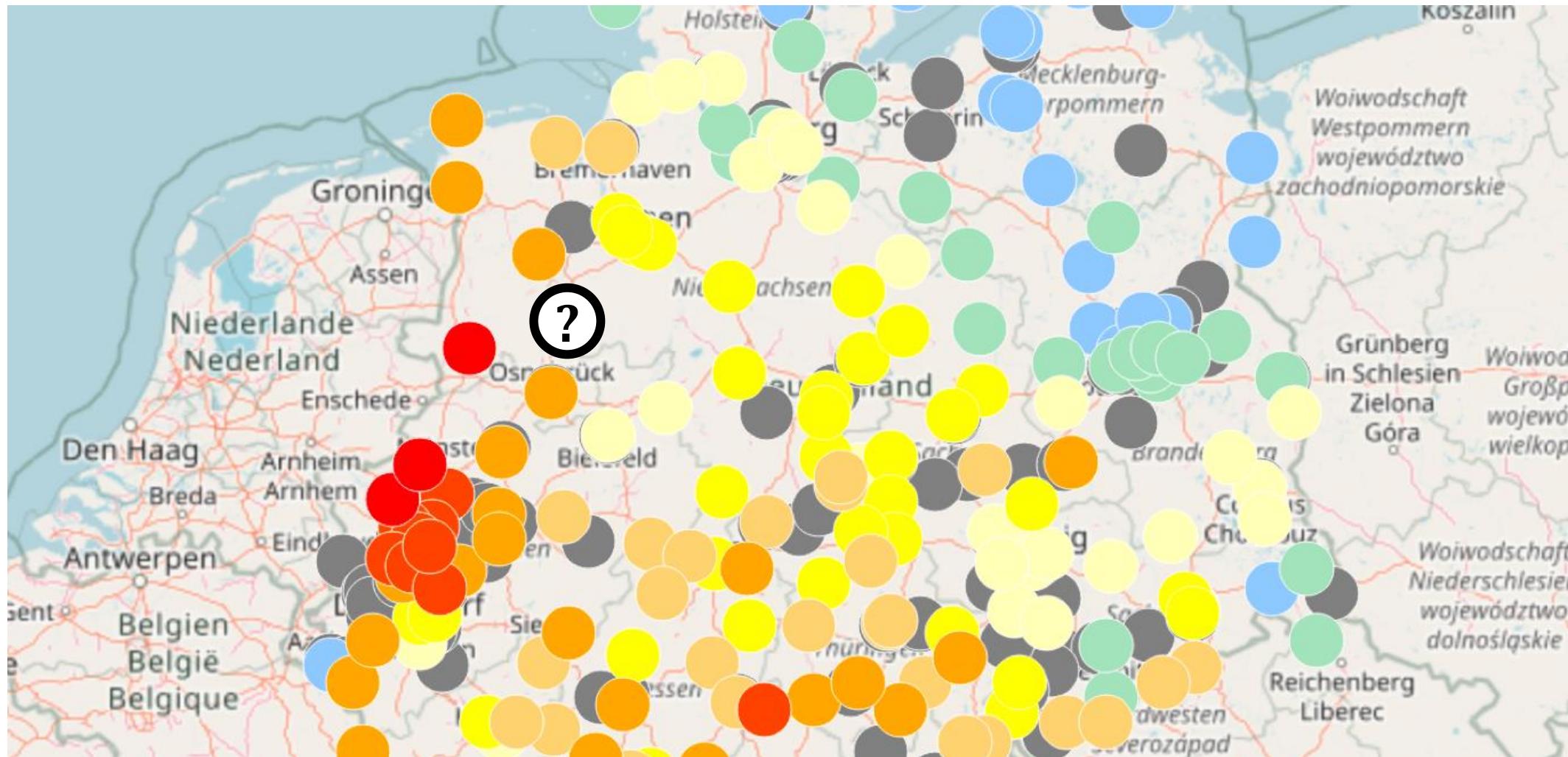


Quelle: <https://de.statista.com/statistik/daten/studie/5564/umfrage/monatliche-durchschnittstemperatur-in-deutschland/>

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Räumliche Datenlücken

Bsp: Messstationen für Luftqualität

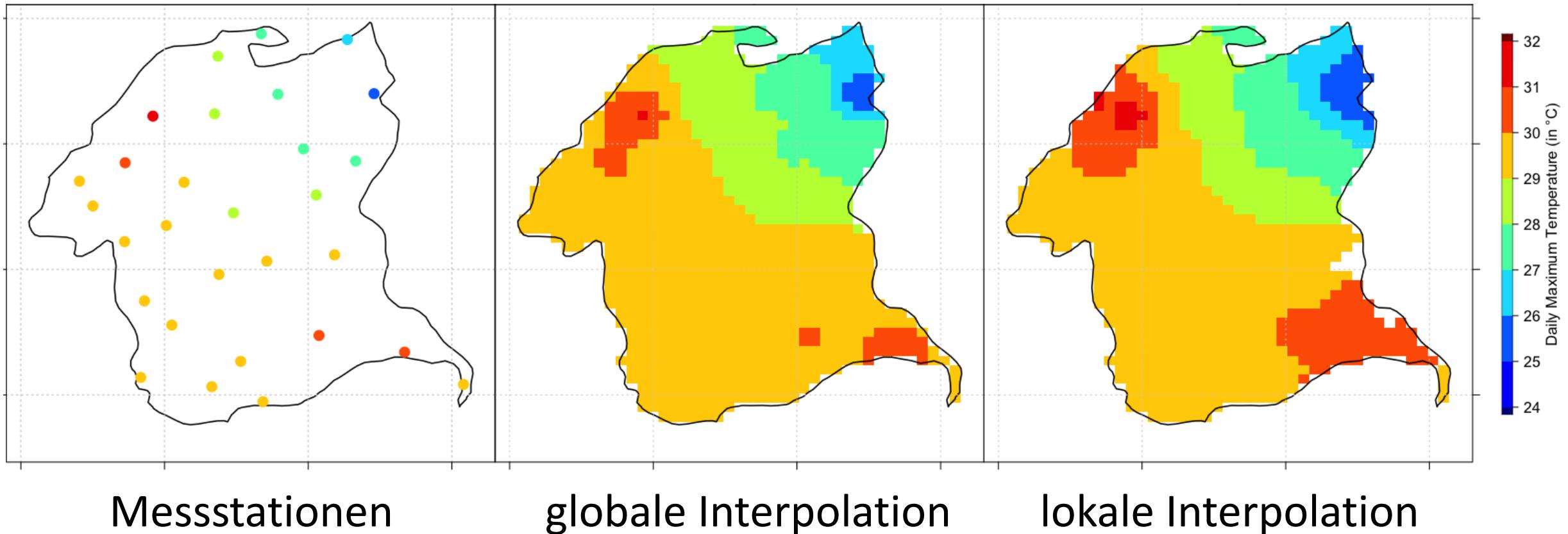


Quelle: <https://www.umweltbundesamt.de/daten/luft/luftdaten>

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Räumliche Interpolation

## global vs. lokal (Methode: *inverse distance weight*)

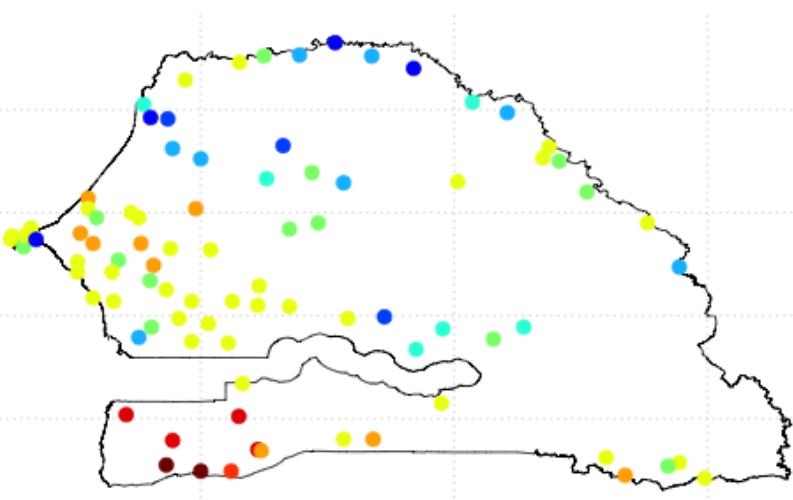


Quelle: [https://iri.columbia.edu/~rijaf/CDTUserGuide/html/interpolation\\_methods.html](https://iri.columbia.edu/~rijaf/CDTUserGuide/html/interpolation_methods.html)

Wie kann KI dabei helfen, Datenlücken bei der Entwicklung bodenbezogener Indikatoren zu schließen?

# Räumliche Interpolation

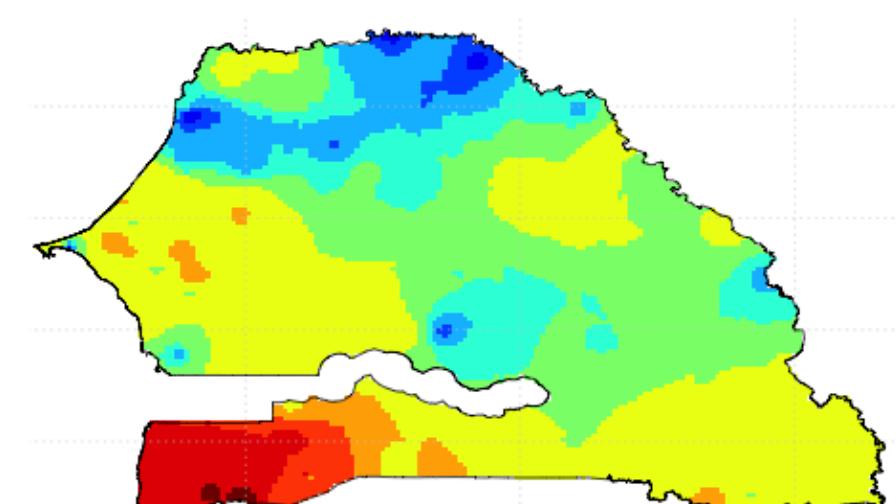
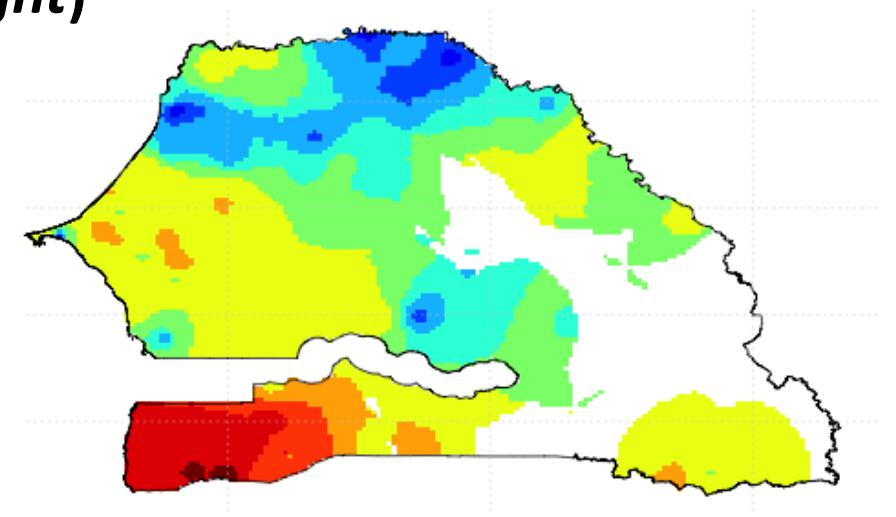
Lokale Interpolation (Methode: *inverse distance weight*)



Messstationen

fester Radius  $r_{max}$

variabler Radius  $r_{max}$

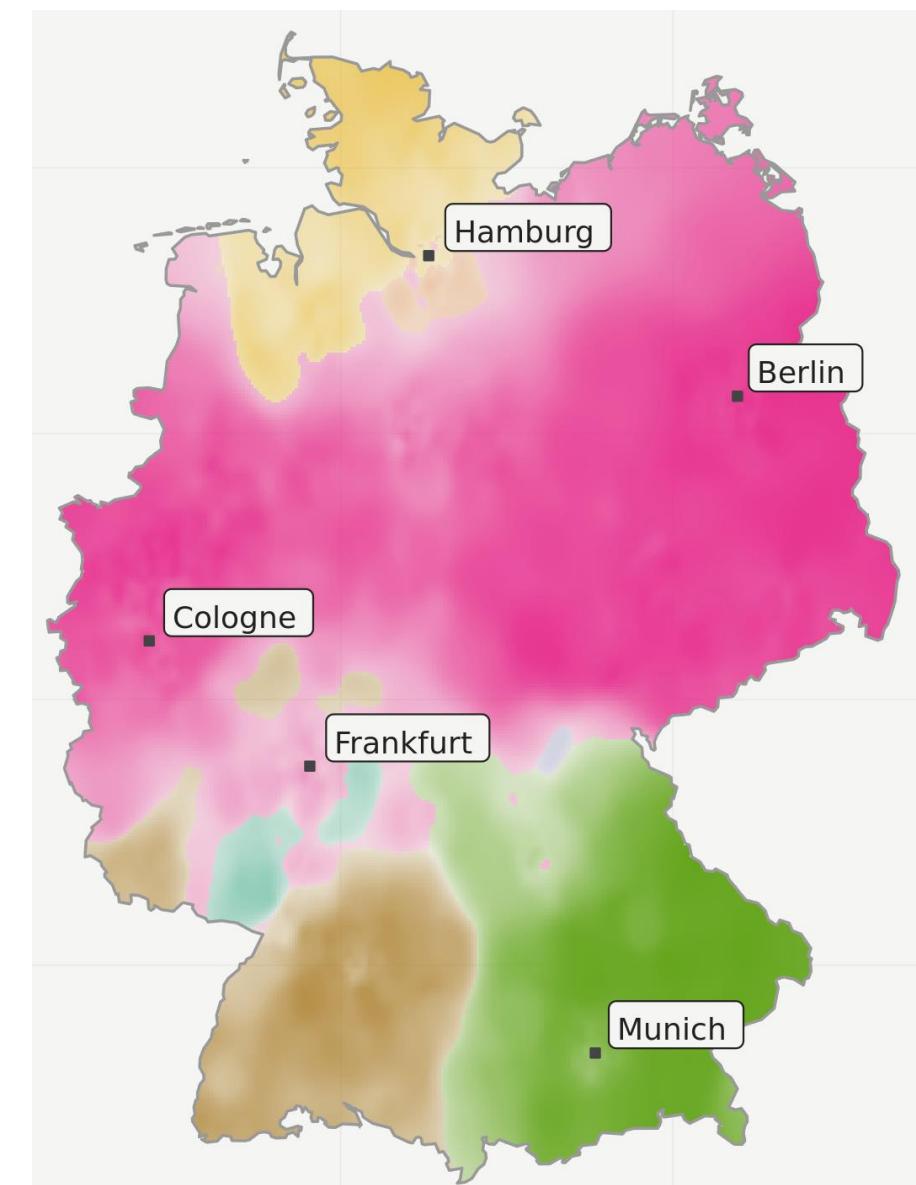
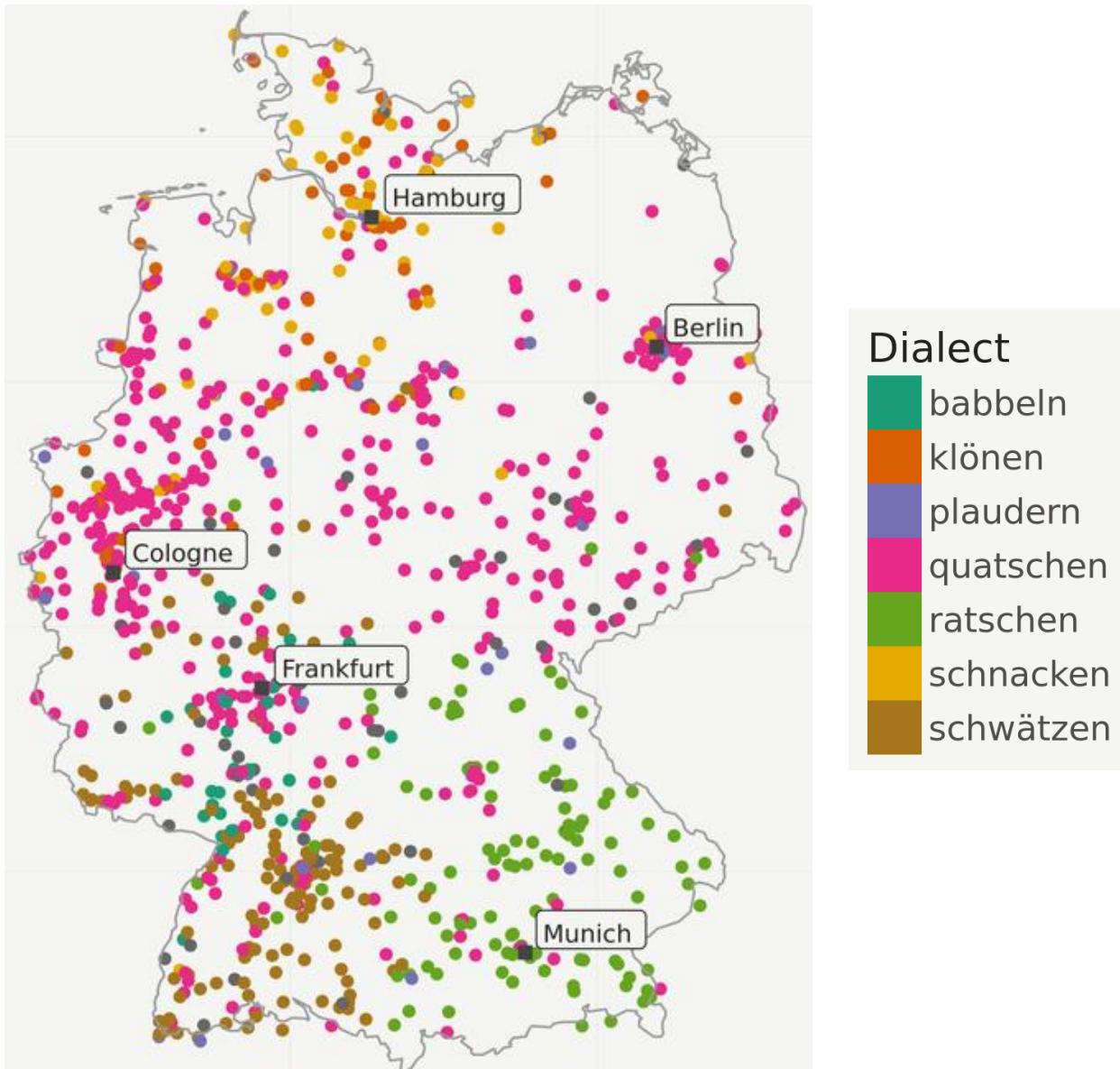


Quelle: [https://iri.columbia.edu/~rijaf/CDTUserGuide/html/interpolation\\_methods.html](https://iri.columbia.edu/~rijaf/CDTUserGuide/html/interpolation_methods.html)

Danke fürs Zuhören ☺

Jetzt: Eure Fragen!





Quelle: <https://timogrossenbacher.ch/categorical-spatial-interpolation-with-r/>