

Geostatistik und räumliche Vorhersage

am Beispiel von Angebotsmietdaten in Magdeburg

Georg Wiegler

Amt für Statistik, Wahlen und Digitalisierung Magdeburg

Erfurt, 16.09.2022

Einleitung

"... everything is related to everything else, but near things are more related than distant things"

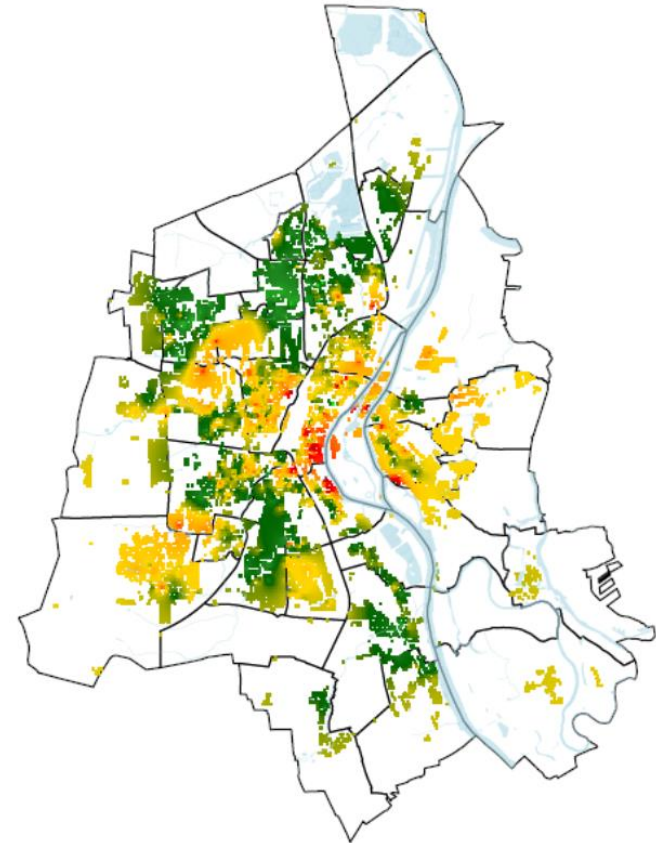
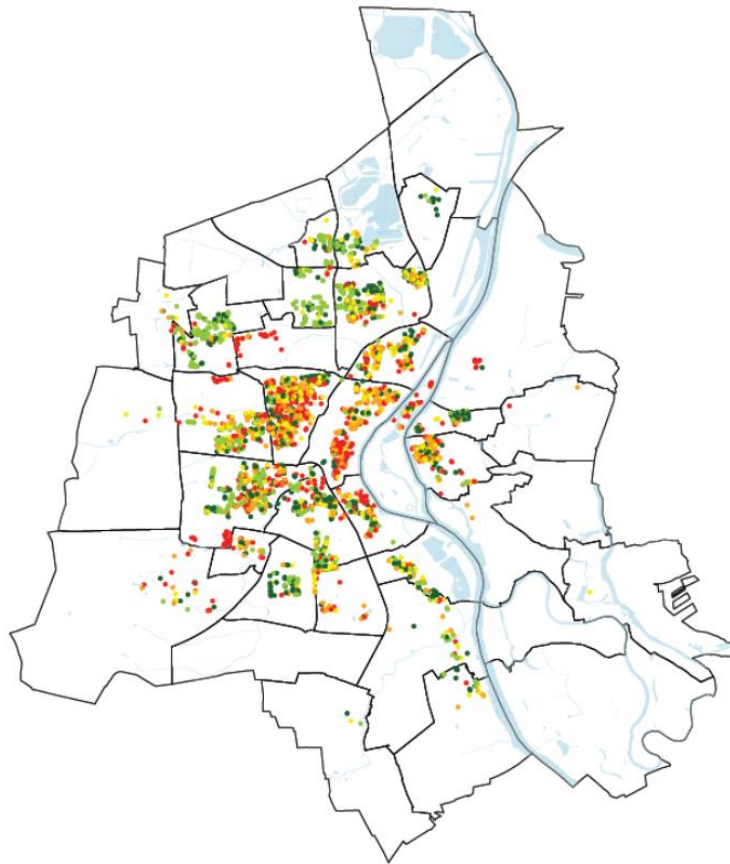
— Tobler's first law of geography

- Wohnungsmietpreise als **prominenter Untersuchungsgegenstand**
- Zentrales Anliegen: **Vorhersage/Prädiktion** von Mietpreisen anhand anderer Merkmale
- Einfluss der **Lage** allgemein anerkannt, aber **problematisch**
- In georeferenzierten Daten liegt großes (häufig ungenutztes) Potenzial Vorhersage/Prädiktion von Mietpreisen zu verbessern

Geostatistik

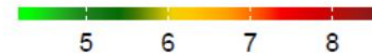
- Geostatistik Sammlung von stochastischen Methoden zur Modellierung georeferenzierter Daten unter Annahme räumlicher Autokorrelation.
- **Ziel:** Schätzung und Vorhersage/Prädiktion auf Basis der Positionen von beobachteten Untersuchungsobjekten (Kriging)

Angebotsmietpreise in Magdeburg



Kaltmiete pro Quadratmeter

● (2.75,5.1]	● (5.1,5.5]	● (5.5,5.94]	● (5.94,6.3]	● (6.3,12.5]
--------------	-------------	--------------	--------------	--------------



(Mathematische) Grundlagen

Allg. Z Stationarität Isotropie

- **Traditionelle Statistik:** Merkmalsträger/ Objekte voneinander unabhängig (iid etc.)
- **Geostatistik:** Räumliche Abhängigkeit, Phänomen ist stochastischer Prozess (Zufallsfeld), Abhängigkeit durch Abstand
- $Z(\mathbf{x})$ ist Zufallsvariable in Abhängigkeit von $\mathbf{x} \in \mathbb{R}^2$ als Koordinate, z.B. $Z((52.12, 11.65))$
- Zusammenhang über Kovarianzfunktion $C(\mathbf{x}, \mathbf{y})$ oder Variogramm $\gamma(\mathbf{x}, \mathbf{y})$ mit \mathbf{x} und \mathbf{y} als 2 Punkten im relevanten Raum $D \subseteq \mathbb{R}^2$
- Annahmen: (intrinsische) Stationarität und Isotropie $\rightarrow \gamma$ und C nur vom (euklidischen) Abstand $|\mathbf{h}| = \|\mathbf{y} - \mathbf{x}\|$ abhängig

(Mathematische) Grundlagen

Allg. Z Stationarität Isotropie

Zufallsfeld

Sei $D \subseteq \mathbb{R}^2$ und $Z(x): \Omega \rightarrow \mathbb{R}$ für alle $x \in D$ messbar.

$$Z := (Z(x))_{x \in D}$$

heißt Zufallsfeld Z .

Kovarianzfunktion

$$C(x, y) = \text{Cov}[Z(x), Z(y)]$$

Variogramm

$$\gamma(x, y) = \frac{1}{2} \text{Var}[Z(y) - Z(x)]$$

(Mathematische) Grundlagen

Allg. Z Stationarität Isotropie

Stationarität

Z heißt **stationär**, wenn $\mathbb{E}[Z(x)] \equiv \mu$
und wenn C verschiebungsinvariant ist.
Für den Abstandsvektor $\mathbf{h} = \mathbf{y} - \mathbf{x}$ gilt:
 $R(\mathbf{h}) := C(\mathbf{x}, \mathbf{x} + \mathbf{h}) = C(\mathbf{x}, \mathbf{y})$.

Intrinsische Stationarität

Z heißt **intrinsisch stationär**, wenn für
alle $\mathbf{h} \in D$ $(Z(\mathbf{x} + \mathbf{h}) - Z(\mathbf{x}))_{\mathbf{x} \in D}$
Stationarität und Zentriertheit ($\mathbb{E} \equiv 0$)
vorliegt. In diesem Fall ist das
Variogramm eine Funktion von \mathbf{h} : $\gamma(\mathbf{h})$.

- Stationäre Zufallsfelder sind auch intrinsisch stationär (gilt nicht umgekehrt)

(Mathematische) Grundlagen

Allg. Z Stationarität Isotropie

Isotropie

Z heißt **isotrop**, wenn es stationär ist und die Kovarianzfunktion invariant gegenüber der Ausrichtung des Abstandsvektors \mathbf{h} ist. In diesem Fall schreibt man $R_0(|\mathbf{h}|) = C(\mathbf{x}, \mathbf{y})$

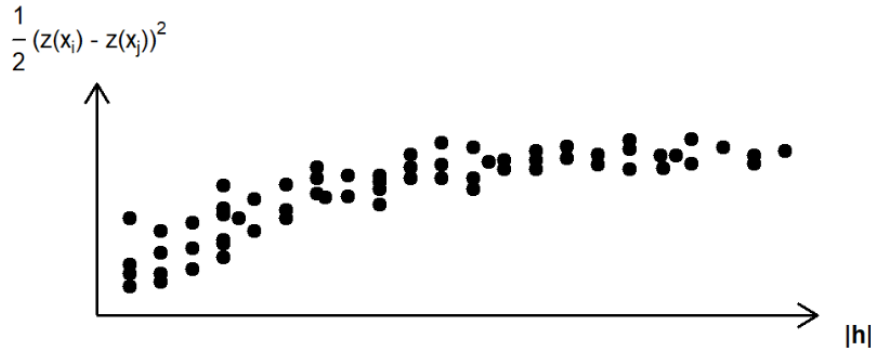
- Im isotropen Fall ist nur die Länge des Abstandsvektors bzw. der (euklidische) Abstand zwischen beiden Punkten relevanten
- Analog gilt im isotropen dass sich das Variogramm als eine Funktion des Abstands zwischen zwei Punkten darstellen lässt

Variographie

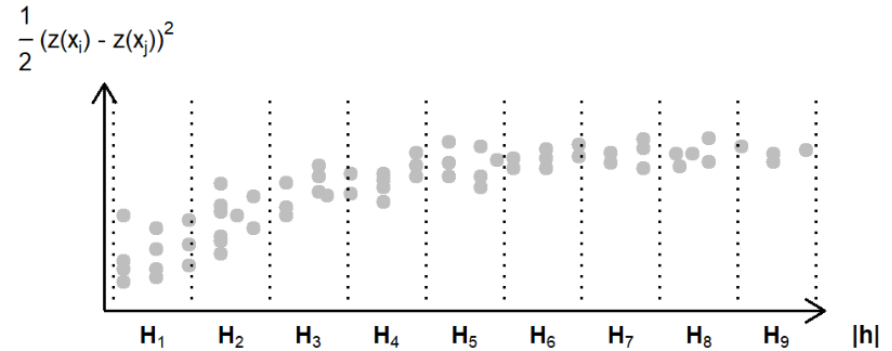
Illustration

Beispiel Magdeburg

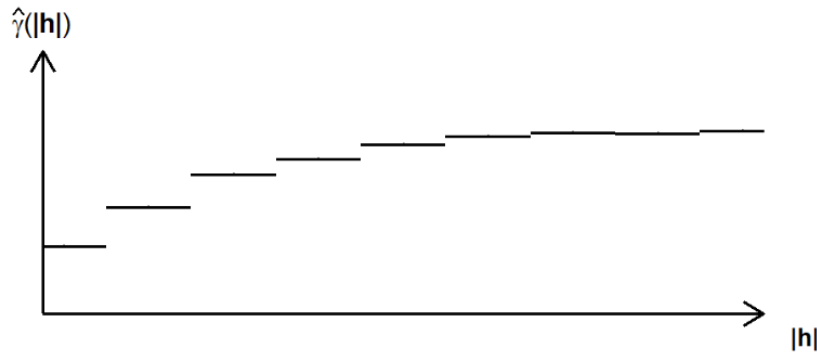
Sinn und Zweck



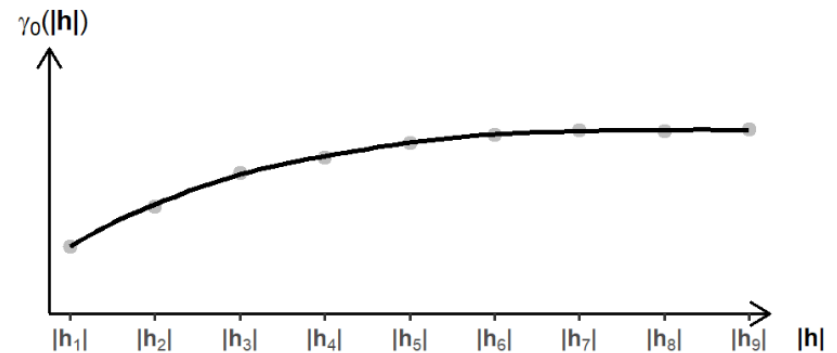
(a) Erstellung der Variogrammwolke



(b) Unterteilung in Distanzgruppen



(c) Bestimmung des empirischen Variogramms



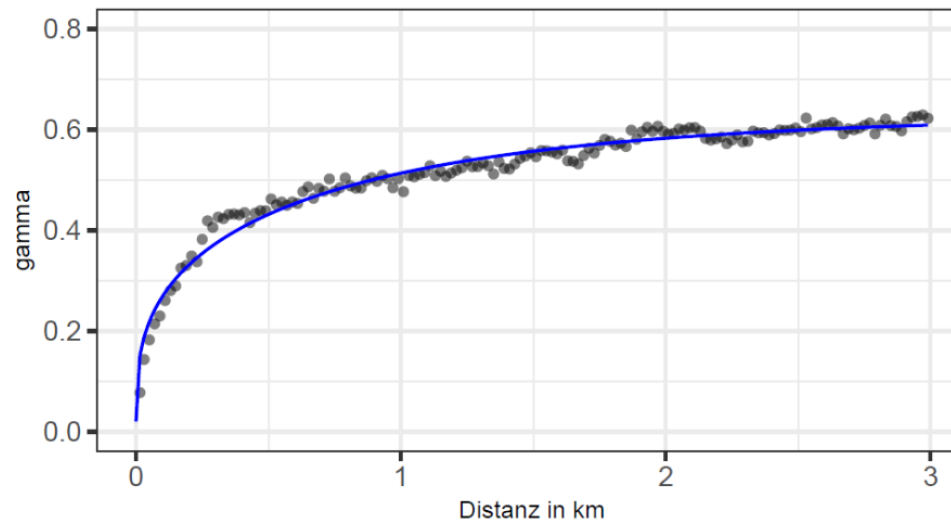
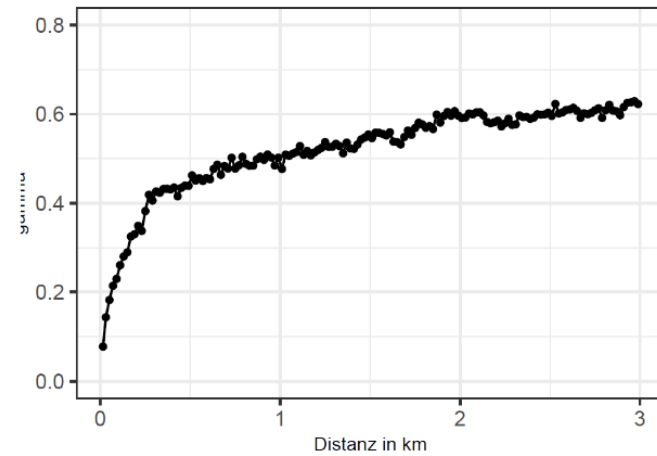
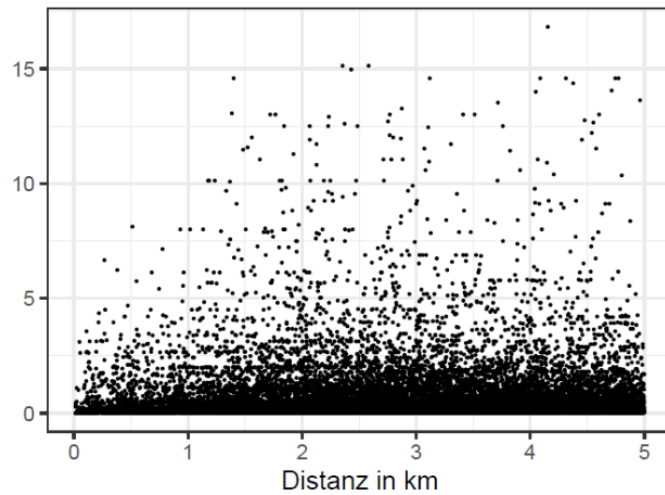
(d) Anpassung der Variogrammfunktion

Variographie

Illustration

Beispiel Magdeburg

Sinn und Zweck

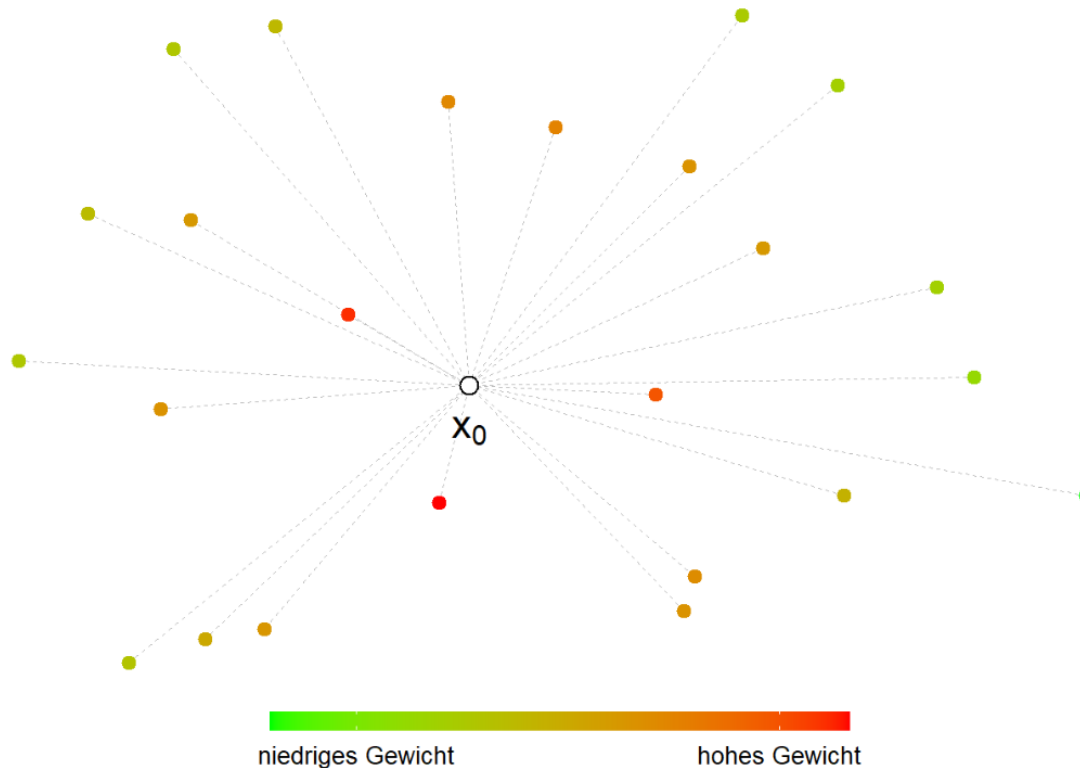


Räumliche Vorhersage

Beobachtungen $z(x_1), \dots, z(x_n)$ aus Z . Was ist eine **beste Vorhersage/Prädiktion** für $z(x_0)$ für $x_0 \in D$?

Gewichtung und Abstand

Linearer Prädiktor



Räumliche Vorhersage

Beobachtungen $z(\mathbf{x}_1), \dots, z(\mathbf{x}_n)$ aus Z . Was ist eine **beste Vorhersage/Prädiktion** für $z(\mathbf{x}_0)$ für $\mathbf{x}_0 \in D$?

Gewichtung und Abstand

Linearer Prädiktor

Ein **linearer Prädiktor** $\hat{Z}(\mathbf{x}_0)$ für $Z(\mathbf{x}_0)$ an der Position \mathbf{x}_0 ist definiert als:

$$\hat{Z}(\mathbf{x}_0) := \sum_{i=1}^n w_i Z(\mathbf{x}_i) = \mathbf{w}^T \mathbf{Z}(\mathbf{x})$$

mit $\mathbf{w} = (w_1, \dots, w_n)^T$ und $\mathbf{Z}(\mathbf{x}) = (Z(\mathbf{x}_1), \dots, Z(\mathbf{x}_n))^T$.

$\hat{Z}^*(\mathbf{x}_0)$ ist **bester, linearer, unverzerrter Prädiktor (BLUP)** für $Z(\mathbf{x}_0)$, wenn gilt:

$$\mathbf{Unverzerrtheit} : \text{Bias}[\hat{Z}^*(\mathbf{x}_0)] := \mathbb{E}[\hat{Z}^*(\mathbf{x}_0) - Z(\mathbf{x}_0)] = 0$$

$$\mathbf{Minimale Varianz} : \text{Var}[Z^*(\mathbf{x}_0) - \hat{Z}(\mathbf{x}_0)] \leq \text{Var}[\hat{Z}(\mathbf{x}_0) - Z(\mathbf{x}_0)]$$

für alle linear unverzerrten Prädiktoren $\hat{Z}(\mathbf{x}_0)$.

Ordinary Kriging

Idee für Mietpreise

Beispiel Magdeburg

Modell

Gleichungssystem

Annahmen

- Mietpreis wird durch einen durchschnittlichen Wert und ein Zufallsfeld bestimmt
- Der Durchschnittswert ist nicht bekannt
- Das Zufallsfeld hat eine gewisse Stabilität (intrinsische Stationarität + Isotropie)
- Das Variogramm ist bekannt bzw. wurde aus den Daten geschätzt

Dann ist der Ordinary-Kriging-Vorhersage die statistisch beste Prädiktion (BLUP)

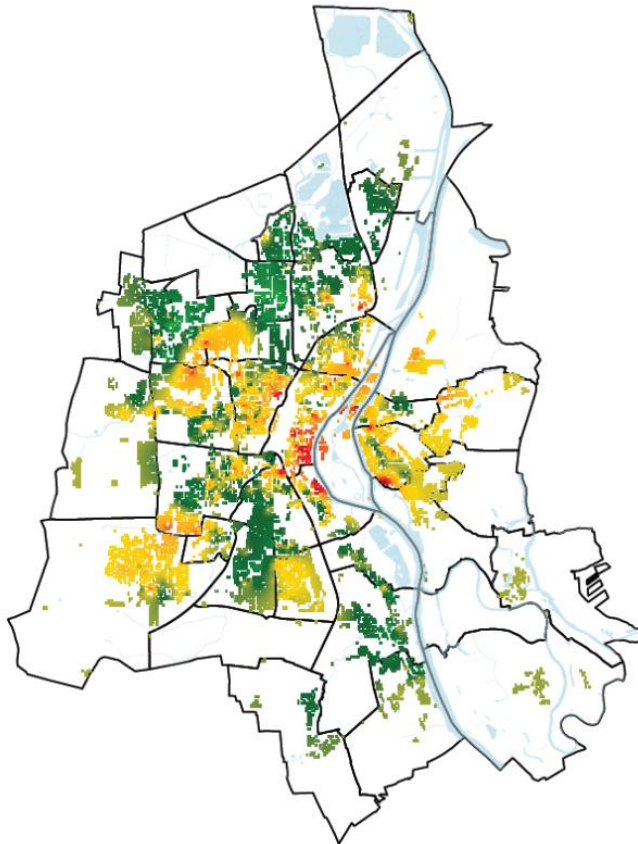
Ordinary Kriging

Idee für Mietpreise

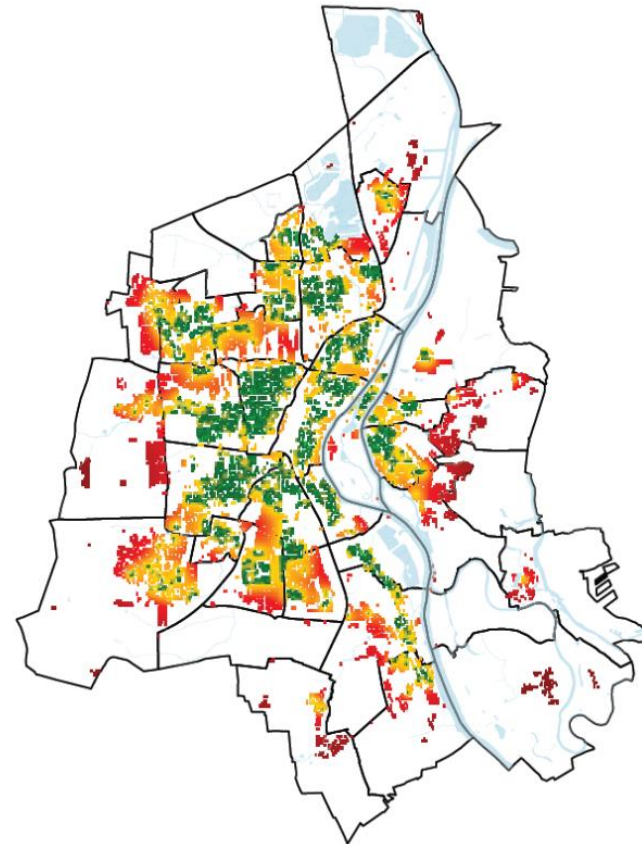
Beispiel Magdeburg

Modell

Gleichungssystem



(a) Prädiktion



(b) Varianz

Ordinary Kriging

Idee für Mietpreise

Beispiel Magdeburg

Modell

Gleichungssystem

- Beobachtungen aus intrinsisch stationären Zufallsfeld
- konstanter, aber unbekannter Erwartungswert μ

Gegeben Sei ein intrinsisch stationäres Zufallsfeld Z mit $\mathbb{E} \equiv \mu \in \mathbb{R}$ unbekannt und Variogramm γ . Für beobachtete Positionen $x_1, \dots, x_n \in D$ sei $\mathbf{V} := [\gamma(x_i - x_j)]_{i,j=1}^n$ invertierbar. Dann ist

$$\hat{Z}(x_0) := \mathbf{w}_{OK}^T Z(\mathbf{x})$$

BLUP für $Z(x_0)$, wenn $\mathbf{w}_{OK} := \mathbf{V}^{-1} \left(\mathbf{v}_0 - \mathbf{1}_n \left(\frac{\mathbf{1}_n^T \mathbf{V}^{-1} \mathbf{v}_0 - 1}{\mathbf{1}_n^T \mathbf{V}^{-1} \mathbf{1}_n} \right) \right)$ mit

$$\mathbf{v}_0 := (\gamma(x_1 - x_0), \dots, \gamma(x_n - x_0))^T.$$

Ordinary Kriging

Idee für Mietpreise

Beispiel Magdeburg

Modell

Gleichungssystem

$$\begin{pmatrix} 0 & \gamma(x_1 - x_2) & \cdots & \gamma(x_1 - x_n) & 1 \\ \gamma(x_2 - x_1) & 0 & \cdots & \gamma(x_2 - x_n) & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \gamma(x_n - x_1) & \gamma(x_n - x_2) & \cdots & 0 & 1 \\ 1 & 1 & \cdots & 1 & 0 \end{pmatrix} \begin{pmatrix} w_1^{OK} \\ w_2^{OK} \\ \vdots \\ w_n^{OK} \\ \lambda \end{pmatrix} = \begin{pmatrix} \gamma(x_1 - x_0) \\ \gamma(x_2 - x_0) \\ \vdots \\ \gamma(x_n - x_0) \\ 1 \end{pmatrix}$$

Universal Kriging

Idee für Mietpreise

Trendoberfläche

Beispiel Magdeburg

Modell

Gleichungssystem

Annahmen

- Mietpreis wird durch eine Trendoberfläche (Durchschnittswert verändert sich) und ein Zufallsfeld bestimmt
- Die Trendoberfläche ist nicht bekannt
- Das Zufallsfeld hat eine gewisse Stabilität (intrinsische Stationarität + Isotropie)
- Das Variogramm ist bekannt bzw. wurde aus den Daten geschätzt

Dann ist der Universal-Kriging-Vorhersage die statistisch beste Prädiktion (BLUP)

Universal Kriging

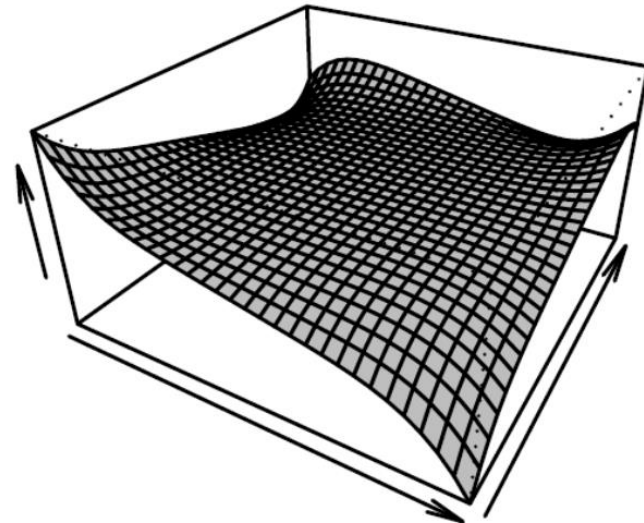
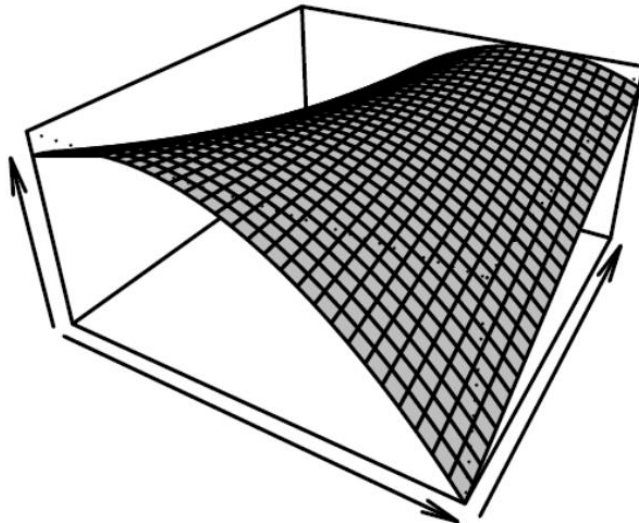
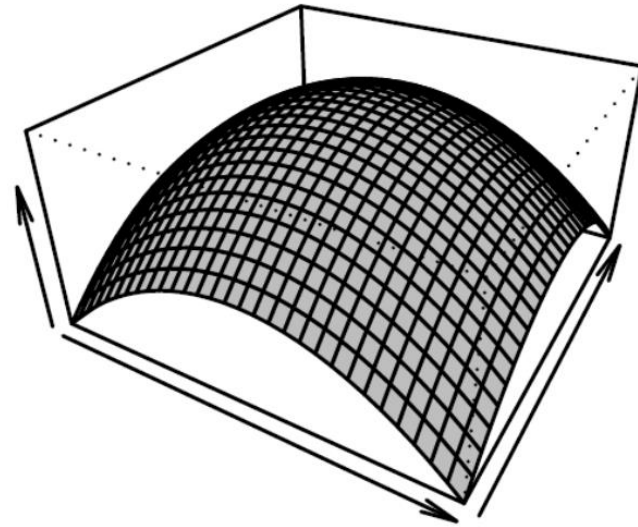
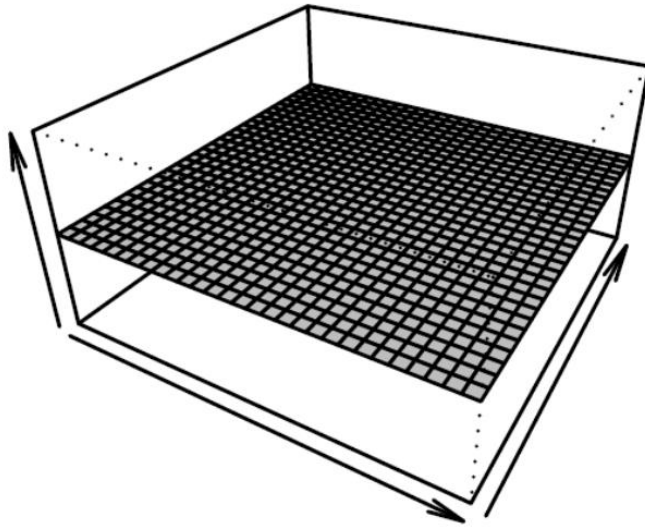
Idee für Mietpreise

Trendoberfläche

Beispiel Magdeburg

Modell

Gleichungssystem



Universal Kriging

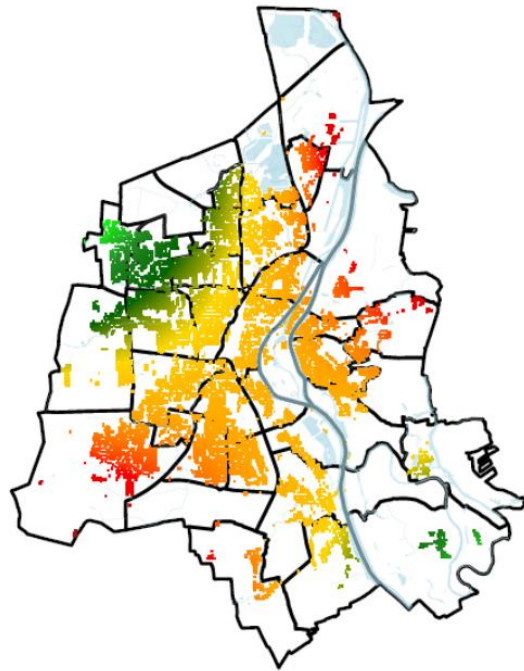
Idee für Mietpreise

Trendoberfläche

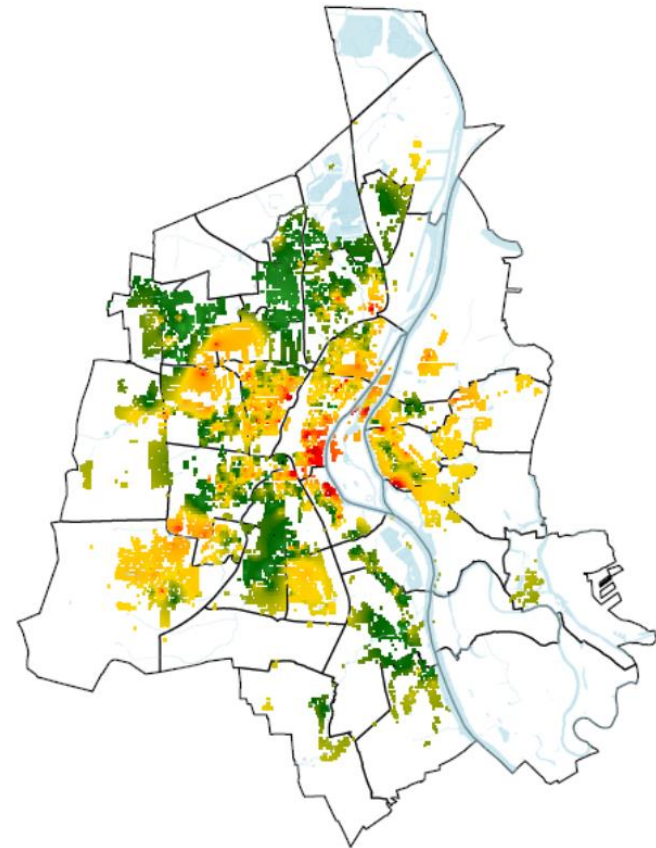
Beispiel Magdeburg

Modell

Gleichungssystem



5.4 5.6 5.8 6.0 6.2



5 6 7 8

Universal Kriging

Idee für Mietpreise

Trendoberfläche

Beispiel Magdeburg

Modell

Gleichungssystem

- Zufallsfeld $Z(\mathbf{x}) = m(\mathbf{x}) + \epsilon(\mathbf{x})$ mit $m(\mathbf{x})$ deterministisch und $\epsilon(\mathbf{x})$ zufällig (universelles Modell spatialer Variation)
 - $m(\mathbf{x})$ wird Drift genannt und beschreibt eine (unbekannte) polynomiale Trendoberfläche, die sich durch ein lineares Modell schätzen lässt (Koordinaten als erklärende Variablen)
 - $\epsilon(\mathbf{x})$ ist ein intrinsisch stationäres, zentriertes Zufallsfeld

Sei Z ein Zufallsfeld mit den oben beschriebenen Annahmen. Sei \mathbf{F} eine Designmatrix für ein allgemeines lineares Modell und $\boldsymbol{\alpha}$ der zugehörige Parametervektor für die polynomiale Gestalt der Trendoberfläche. Für beobachtete Positionen $\mathbf{x}_1, \dots, \mathbf{x}_n \in D$ sei $\mathbf{V}_\epsilon := [\gamma_\epsilon(\mathbf{x}_i - \mathbf{x}_j)]_{i,j=1}^n$ invertierbar. Dann ist der Prädiktor $\hat{Z}(\mathbf{x}_0) := \sum_{i=1}^n w_i Z(\mathbf{x}_i) = \mathbf{w}^T \mathbf{Z}(\mathbf{x})$ BLUP, wenn

$$\mathbf{w} = \mathbf{w}_{UK} := \mathbf{V}_\epsilon^{-1} \left(\mathbf{v}_{\epsilon,0} - \mathbf{F} \left(\mathbf{F}^T \mathbf{V}_\epsilon^{-1} \mathbf{F} \right)^{-1} \left(\mathbf{F}^T \mathbf{V}_\epsilon^{-1} \mathbf{v}_{\epsilon,0} - \mathbf{f}_0 \right) \right).$$

Universal Kriging

Idee für Mietpreise

Trendoberfläche

Beispiel Magdeburg

Modell

Gleichungssystem

$$\begin{pmatrix}
 \gamma_\epsilon(x_1 - x_1) & \cdots & \gamma_\epsilon(x_1 - x_n) & 1 & f_1(x_1) & \cdots & f_L(x_1) \\
 \gamma_\epsilon(x_2 - x_1) & \cdots & \gamma_\epsilon(x_2 - x_n) & 1 & f_1(x_2) & \cdots & f_L(x_2) \\
 \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 \gamma_\epsilon(x_n - x_1) & \cdots & \gamma_\epsilon(x_n - x_n) & 1 & f_1(x_n) & \cdots & f_L(x_n) \\
 \hline
 1 & \cdots & 1 & 0 & 0 & \cdots & 0 \\
 f_1(x_1) & \cdots & f_1(x_n) & 0 & 0 & \cdots & 0 \\
 \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 f_L(x_1) & \cdots & f_L(x_n) & 0 & 0 & \cdots & 0
 \end{pmatrix}
 \begin{pmatrix}
 w_1^{UK} \\
 w_2^{UK} \\
 \vdots \\
 w_n^{UK} \\
 \alpha_0^{UK} \\
 \alpha_1 \\
 \vdots \\
 \alpha_L^{UK}
 \end{pmatrix}
 =
 \begin{pmatrix}
 \gamma_\epsilon(x_1 - x_0) \\
 \gamma_\epsilon(x_2 - x_0) \\
 \vdots \\
 \gamma_\epsilon(x_n - x_0) \\
 f_0(x_0) \\
 f_1(x_0) \\
 \vdots \\
 f_L(x_0)
 \end{pmatrix}$$

$$\begin{pmatrix}
 \mathbf{V}_\epsilon & \mathbf{F} \\
 \mathbf{F}^T & \mathbf{0}
 \end{pmatrix}
 \begin{pmatrix}
 \mathbf{w}^{UK} \\
 \boldsymbol{\lambda}
 \end{pmatrix}
 =
 \begin{pmatrix}
 \mathbf{v}_{\epsilon,0} \\
 \mathbf{f}_0
 \end{pmatrix}$$

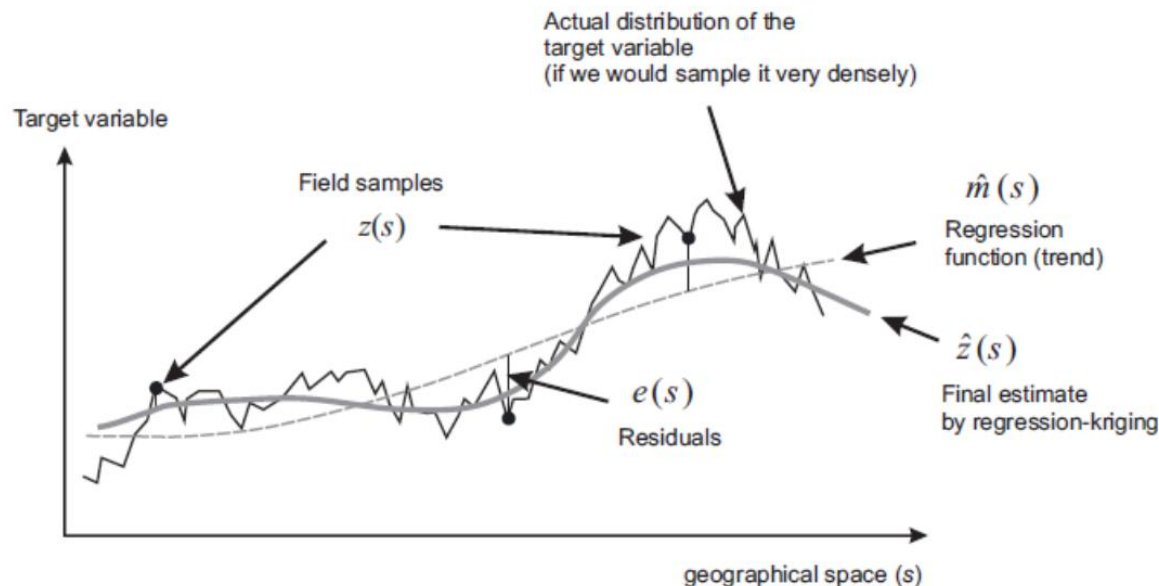
Weitere Spielarten von Kriging

Regression-Kriging

Block-Kriging

Cokriging

- auch: Gauss-Prozess-Regression
- Methodisch/ Mathematisch äquivalent zu Universal Kriging.
- Strukturelle Eigenschaften statt Geokoordinaten für "Trend" (z.B. Wohnungsgröße, Ausstattung, Baujahr, ...)
- $Z(x, s) = m(s) + \epsilon(x)$
 - $m(s)$ als deterministisch, Resultat der Regressionsfunktion mit Merkmalskombination s
 - $\epsilon(x)$ als Zufallsfeld (intrinisch stationär, zentriert, isotrop)



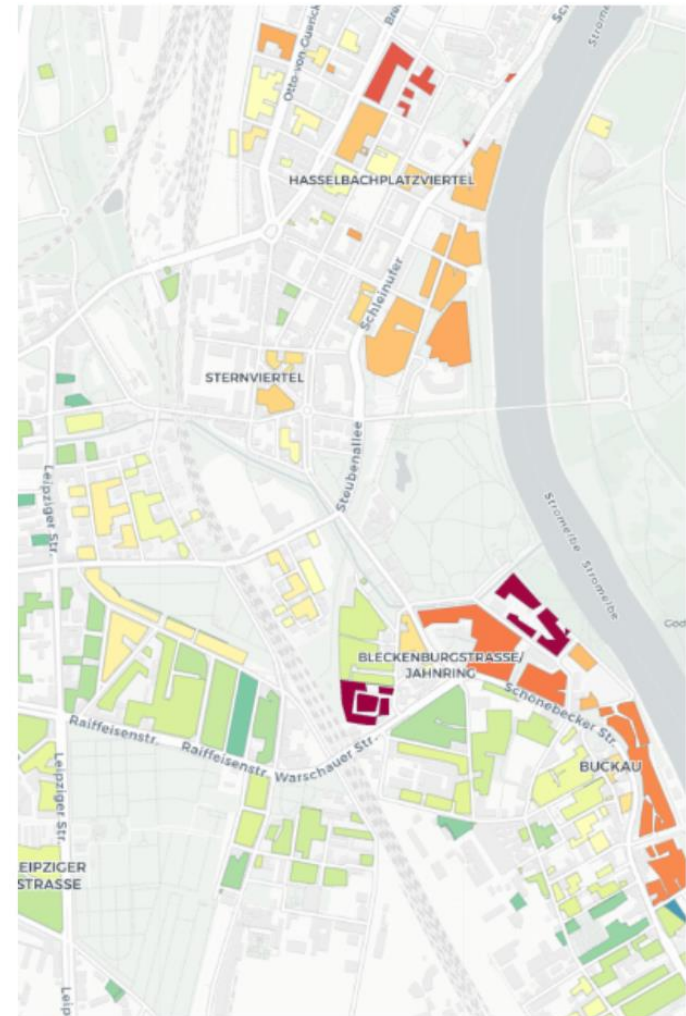
Weitere Spielarten von Kriging

Regression-Kriging

Block-Kriging

Cokriging

- Statt einer einzelnen Position (Point-Kriging) wird der erwartete Wert über ein gesamtes Polygon $P \subseteq D$ interpoliert
 - Stadtteile, Statistische Bezirke, Blöcke, Flur-Stücke,...
- Definiert als ein Integral des Zufallsfeldes über P , Annäherung über Durchschnitte mehrerer Prädiktionen
- Besonderheit:
 - Auch Werte außerhalb der durch die Polygon definierten Bereiche werden berücksichtigt
 - keine Restriktion bzgl. Anzahl der Merkmalsträger
 - Bei Autokorrelation bessere Vorhersage (Standardfehler geringer)



Weitere Spielarten von Kriging

Regression-Kriging

Block-Kriging

Cokriging

- Natürliche Erweiterung von Kriging für den multivariaten Fall.
- Viele Anwendungsszenarien:
 - Simuliertes Kriging mehrerer abhängiger Variablen
 - Berücksichtigung mehrerer Stichproben
 - Verbesserung der Prädiktion einer räumlich weniger gut verteilten Variable durch Hilfsvariablen, die regelmäßiger beobachtet wurden (Collocated Cokriging)
 - Berücksichtigung der Korrelation von exogenen Merkmalen
- Denkbare Anwendungsfälle:
 - Verbesserung der Vorhersage von Bestandsmietdaten durch Hinzunahme von Angebotsmietdaten
 - Hinzunahme exogener Variablen, wie Bodenrichtwerte, Abstände zu Grünflächen etc., die an jeder Position ermittelt werden können

Regression-Block-Kriging Magdeburg

Regressionskoeffizienten

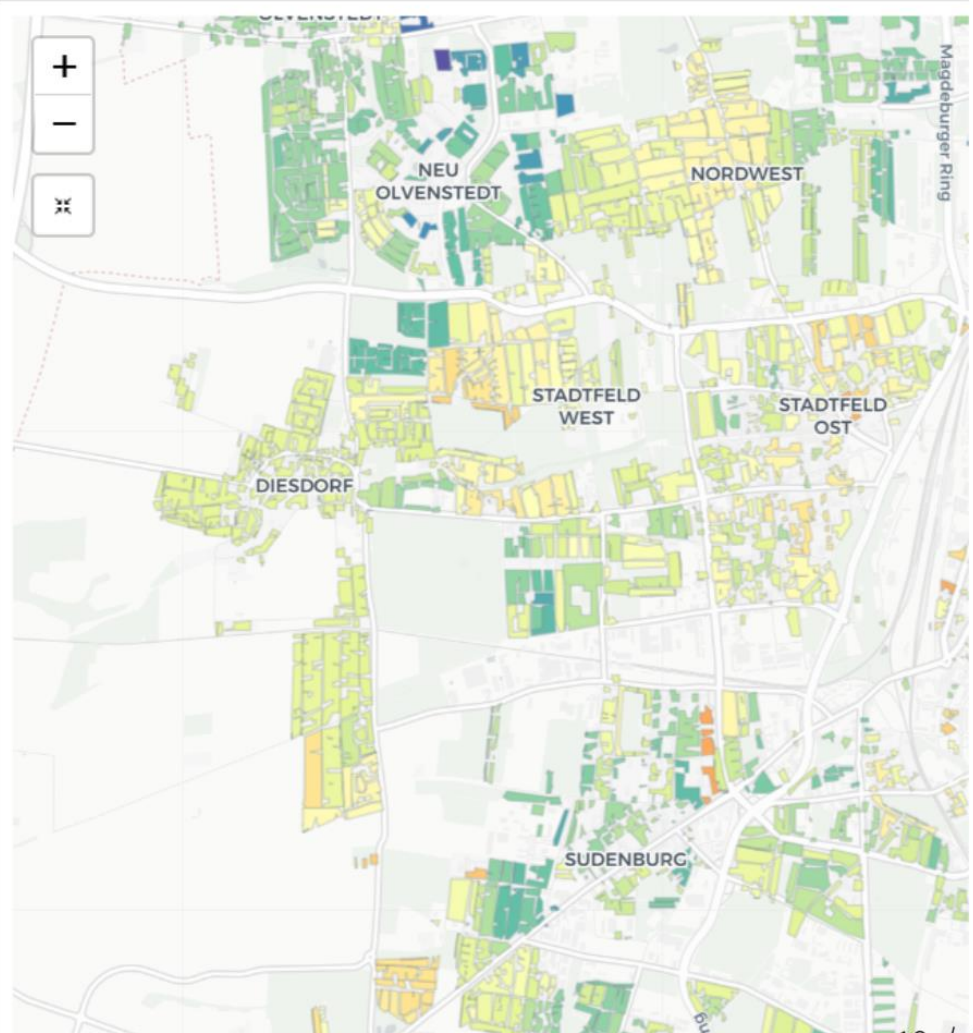
Download ▾

Filter

Schätzung ⚡ P-Wert ⚡

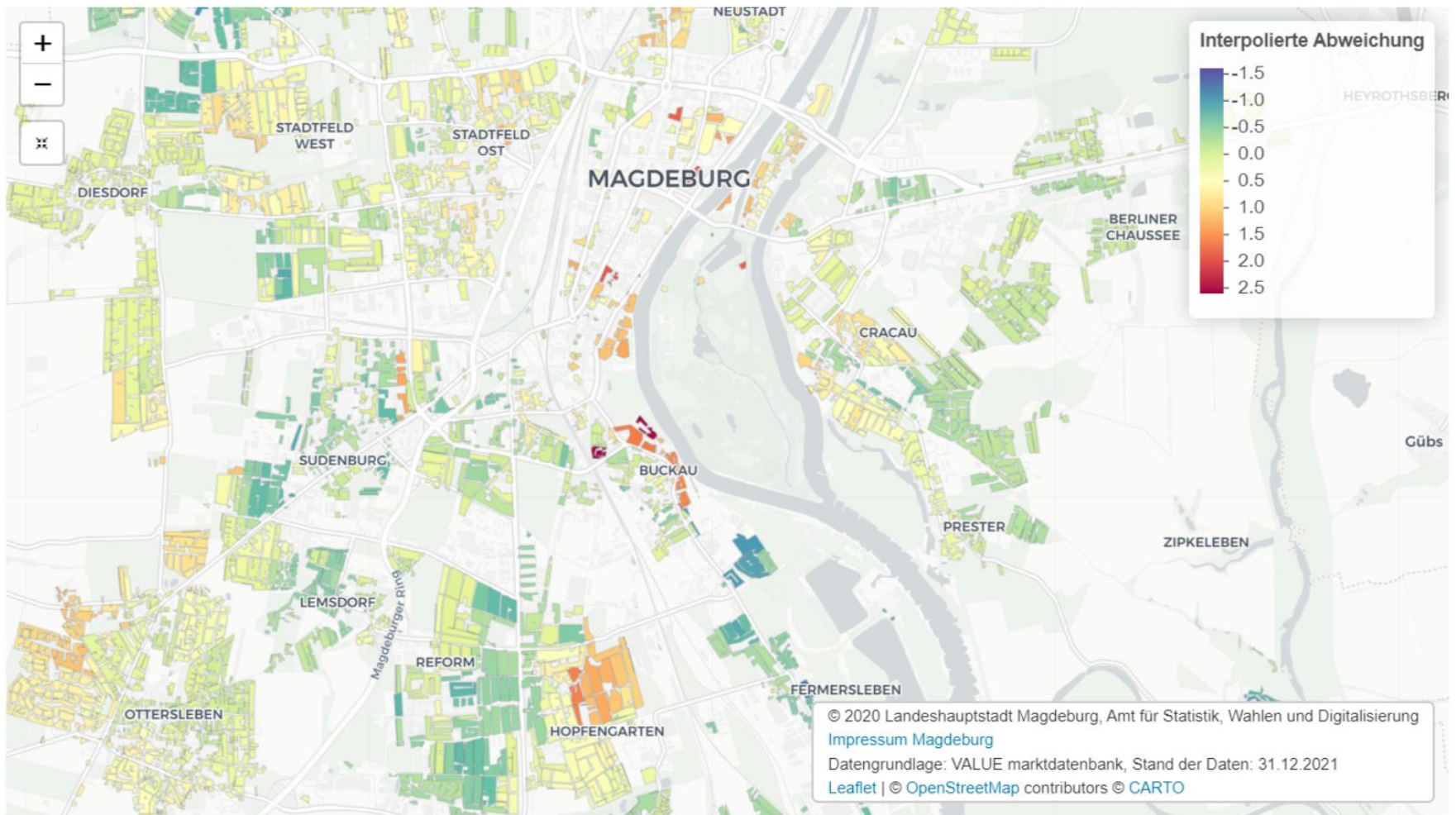
Ausgangskaltmiete pro qm	6.473	0
40 bis 59 qm	-0.299	0
60 bis 79 qm	-0.333	0
80 bis 99 qm	-0.15	0.003
110 bis 119 qm	0.17	0.034
120 und mehr qm	0.224	0.012
Lift vorhanden	0.317	0
Einbauküche	0.449	0
Fußbodenheizung	1.057	0
Balkon oder Terasse	0.27	0
Neuwertig	0.17	0.03
Etage (1,2,3,...)	-0.039	0

Block-Kriging-Resultate



Regression-Block-Kriging Magdeburg

Block-Kriging-Resultate



Verwendete und weiterführende Literatur

Chica-Olmo, J., Cano-Guervos, R. \& Chica-Rivas, M. (2019). Estimation of Housing Price Variations Using Spatio-Temporal Data. *Sustainability*, 11(6), 1551.

Chilès, Jean-Paul; Delfiner, Pierre (1999): *Geostatistics. Modeling spatial uncertainty*. New York, NY: Wiley (A Wiley-Interscience publication).

Cressie, Noel A. C. (1993): *Statistics for spatial data*. Rev. ed. New York: Wiley (Wiley series in probability and mathematical statistics Applied probability and statistics).

Del Giudice, V. & de Paola, P. (2017). Spatial Analysis of Residential Real Estate Rental Market with Geoaddivitive Models. In M. D'AMATO \& T. KAUKO (Hrsg.), *Advances in Automated Valuation Modeling* (Bd. 86, S. 155–162). *Studies in Systems, Decision and Control*.

Gaetan, Carlo; Guyon, Xavier (2010): *Spatial Statistics and Modeling*. New York, NY: Springer Science+Business Media LLC (Springer Series in Statistics).

Hengl, Tomislav (2009): *A practical guide to geostatistical mapping*. 2nd extended ed. Amsterdam: Hengl.

Kuntz, M. & Helbich, M. (2014). Geostatistical mapping of real estate prices: an empirical comparison of kriging and cokriging. *International Journal of Geographical Information Science*, 28(9), 1904–1921.

Lerbs, Oliver; Sebastian, Steffen (2015): *Mietspiegel aus ökonomischer Sicht – Vorschläge für eine Neuregulierung*.

Verwendete und weiterführende Literatur

Lichtenstern, Andreas (2013): Kriging methods in spartial statistics. Technische Universität München, München. Department of Mathematics.

Matheron, Georges (1971): The Theory of Regionalized Variables and Its Applications. Les Cahiers du Centre de Morphologie Mathematique in Fontainebleu, Paris.

Olea, R. A. (1999). Geostatistics for Engineers and Earth Scientists.

Schernthanner, Harald (2017): Räumliche Analyse und Visualisierung von Mietpreisdaten. Dissertation. Universität Potsdam.

Stein, Michael L. (1999): Interpolation of Spatial Data. Some Theory for Kriging. New York, NY: Springer (Springer Series in Statistics).

Tobler, W. R. (1970): A Computer Movie Simulating Urban Growth in the Detroit Region. In: Economic Geography 46, S. 234.

Wackernagel, Hans (2003): Multivariate geostatistics. An introduction with applications : 7 tables. 3., completely revised ed. Berlin: Springer.