



Methodenbericht zur Erstellung der Wohnlagenkarte für das Stadtgebiet Duisburg

Stand 19.12.2023

1. Einleitung	1
2. Konzept	1
3. Methodik der Wohnlagenermittlung	2
4. Indikatorenauswahl und Datenerhebung	3
4.1. Indikatorenauswahl	3
4.2. Datenerhebung	4
5. Auswahl und Training der modellgestützten Verfahren	6
6. Homogenisierung und kartographische Darstellung	11
7. Qualitäts- und Plausibilitätsprüfung	14
8. Literaturverzeichnis	16
9. Anlage 1 :	16
10. Anhang.....	18

1. Einleitung

Eine Wohnlagenkarte erfüllt vielfältige Funktionen in der Öffentlichkeit, der Wirtschaft sowie innerhalb der Verwaltung. Insbesondere für den Gutachterausschuss stellt sie ein wesentliches Instrument dar, um fundierte und nachvollziehbare Bewertungen von Immobilien vorzunehmen. Darüber hinaus dient sie als zentrale Grundlage für die Erstellung eines qualifizierten Mietspiegels. Für das Stadtgebiet Duisburg wird erstmalig eine umfassende und offizielle Wohnlagenkarte auf Grundlage einer datenbasierten und standardisierten Unterstützung, sowie nachvollziehbarer und dokumentierter Methodik erstellt. Gleichzeitig trägt die Wohnlagenkarte zur Verbesserung der Transparenz auf dem Grundstücksmarkt bei und unterstützt eine einheitliche Einschätzung von Wohnlagen im gesamten Stadtgebiet.

2. Konzept

Die Bewertung von Wohnlagen ist mit einer Vielzahl unterschiedlicher Einflussfaktoren verbunden. Neben objektiv messbaren Größen, wie etwa baulichen Strukturen, Dichtekennzahlen oder infrastrukturellen Entfernungen, spielen auch qualitativ geprägte Merkmale („weiche Faktoren“), wie beispielsweise das Image eines Gebietes, eine wichtige Rolle. Diese lassen sich

häufig nur eingeschränkt quantifizieren, sind jedoch für die Wahrnehmung und Bewertung von Wohnlagen von erheblicher Bedeutung.

Die zentrale Herausforderung besteht daher darin, eine Methodik zu entwickeln, die sowohl quantifizierbare als auch schwer messbare Einflussgrößen berücksichtigt und daraus eine konsistente und nachvollziehbare Gesamtbewertung ableitet. Zur Unterstützung dieses Prozesses werden modellgestützte Verfahren des maschinellen Lernens eingesetzt. Diese ermöglichen die Analyse komplexer und nichtlinearer Zusammenhänge zwischen den einzelnen Einflussgrößen und erlauben es, Muster in umfangreichen Datensätzen zu erkennen, die in einer rein manuellen Auswertung nur eingeschränkt zugänglich wären.

Der Einsatz dieser Verfahren dient dabei nicht der automatisierten „Entscheidungsfindung“, sondern der datenbasierten und standardisierten Unterstützung der Wohnlagenbewertung. Die Bewertung selbst erfolgt stets auf Grundlage definierter Kriterien und in Verbindung mit fachlicher Expertise. Im Rahmen der Wohnlagenermittlung kommen unterschiedliche Modellansätze zum Einsatz, darunter insbesondere baumbasierte Verfahren, neuronale Netzwerke sowie kernelbasierte Methoden. Diese bilden verschiedene methodische Zugänge zur Klassifikation von Wohnlagen ab und werden im weiteren Verlauf hinsichtlich ihrer Eignung bewertet.

3. Methodik der Wohnlagenermittlung

Die in die Wohnlagenbewertung einfließenden Daten zeichnen sich durch eine hohe Komplexität sowie teilweise nichtlineare Zusammenhänge aus. Klassische statistische Verfahren, insbesondere Regressionsmodelle als Teil der multivariaten Analyse, sind jedoch primär auf die Abbildung linearer Beziehungen und vergleichsweise einfach strukturierter Datensätze ausgerichtet und können in komplexeren Anwendungskontexten an Grenzen stoßen.

Unabhängig vom gewählten Modellansatz besteht bei der Modellbildung grundsätzlich die Gefahr einer Überanpassung („Overfitting“), insbesondere bei der Verwendung sehr flexibler oder komplexer Modelle sowie bei hochdimensionalen Datensätzen. In solchen Fällen kann ein Modell die Trainingsdaten zwar sehr genau abbilden, verliert jedoch an Generalisierungsfähigkeit für neue Daten.

Vor diesem Hintergrund erfordert eine datengestützte Ermittlung von Wohnlagen Verfahren, die in der Lage sind, sowohl nichtlineare Zusammenhänge als auch Wechselwirkungen zwischen verschiedenen Einflussgrößen abzubilden und gleichzeitig eine hinreichende Robustheit gegenüber Überanpassung aufweisen.

Ausgehend von diesen Anforderungen wurde für das Stadtgebiet Duisburg eine spezifische Methodik entwickelt, die modellgestützte Verfahren mit fachlichen Bewertungen kombiniert. Dabei wird insbesondere berücksichtigt, dass bestimmte Einflussgrößen – wie etwa das Image – nicht vollständig objektivierbar sind und daher ergänzend durch Experteneinschätzungen abgebildet werden müssen.

Die Ermittlung der Wohnlagen erfolgt in mehreren aufeinander aufbauenden Arbeitsschritten, die im weiteren Verlauf detailliert dargestellt werden. Diese lassen sich wie folgt zusammenfassen:

- datengestützte Ermittlung der Wohnlagen auf Grundlage eines Indikatorensystems
- anschließende räumliche Generalisierung der Ergebnisse

- abschließende fachliche Plausibilisierung

Das zugrunde liegende Indikatorensystem wurde auf Basis fachwissenschaftlicher Erkenntnisse sowie unter Berücksichtigung der Vorgaben der Mietspiegelverordnung entwickelt.

Die Wohnlagen werden im Ergebnis in vier Kategorien eingeteilt:

- sehr gute Lage
- gute Lage
- mittlere Lage
- einfache Lage

4. Indikatorenauswahl und Datenerhebung

4.1. Indikatorenauswahl

Die Qualität der Wohnlage wird im vorliegenden Ansatz im Wesentlichen durch vier Indikatorengruppen beschrieben:

- Indikatoren zur Kennzeichnung der Bebauungs- und Bevölkerungsdichte
- Indikatoren zur Charakterisierung der physischen Umweltbedingungen (z. B. Lärmimmissionen, umgebende Nutzungen, Begrünung)
- Indikatoren zur Kennzeichnung von Zentralität und Infrastruktur
- Indikator „Image“

Entsprechend den Vorgaben der Mietspiegelverordnung können zur Abbildung der Wohnlage ergänzend auch wertbeeinflussende Größen wie Bodenrichtwerte oder die allgemeine Beliebtheit von Wohngebieten herangezogen werden, sofern die Lagequalität durch strukturelle Indikatoren allein nicht hinreichend differenziert, beschrieben werden kann.

Eine Einbeziehung von Bodenrichtwerten als ergänzender Indikator erfolgt im vorliegenden Fall jedoch nicht. Bodenrichtwerte stellen als marktbasierter Kenngröße bereits das Ergebnis normierter und am Immobilienmarkt orientierter Wertermittlungsverfahren dar und bilden lageprägende Einflussfaktoren in aggregierter Form ab.

In der immobilienökonomischen und ökonometrischen Methodik wird grundsätzlich zwischen erklärenden Merkmalen und Marktergebnissen unterschieden. Die Verwendung von Bodenrichtwerten als erklärender Indikator innerhalb der Wohnlagenklassifikation würde daher zu einer Vermischung von Ursache und Wirkung führen und ist im Sinne der Vermeidung endogener Zusammenhänge methodisch kritisch zu beurteilen. Vor diesem Hintergrund wurde im Rahmen der methodischen Konzeption bewusst entschieden, auf die Einbeziehung von Bodenrichtwerten zu verzichten.

Stattdessen wird der Aspekt der allgemeinen Beliebtheit („Image“) über eine strukturierte und dokumentierte Experteneinschätzung des Gutachterausschusses berücksichtigt. Dieses Vorgehen ermöglicht es, lageprägende, qualitativ geprägte Faktoren zu erfassen, die sich nicht oder nur eingeschränkt durch objektivierbare Indikatoren abbilden lassen.

Der Indikator „Image“ wird dabei nicht ausschließlich mit zentralen Lagen gleichgesetzt. Vielmehr können sowohl zentrale als auch periphere Wohnlagen aufgrund unterschiedlicher Lagequalitäten und Umfeldmerkmale eine hohe oder geringe Attraktivität aufweisen.

Für die erstmalige Ermittlung der Wohnlagen wurden die Indikatoren auf Basis fachwissenschaftlicher Erkenntnisse, der verfügbaren Datenbasis sowie unter Berücksichtigung der Mietspiegelverordnung ausgewählt. Dabei wurden insbesondere folgende Aspekte berücksichtigt:

- Aussagekraft der Indikatoren (u. a. Datenaktualität, räumliche Differenzierbarkeit, Skalenniveau)
- Verfügbarkeit und Qualität der Daten
- Vermeidung inhaltlicher Redundanzen zwischen einzelnen Indikatoren

4.2. Datenerhebung

Im Rahmen des angestrebten, datengestützten KI-Modells wurden die vorliegenden verfügbaren Daten der im vorherigen Kapitel genannten potenziellen Indikatoren für die Bestimmung der Wohnlage statistisch untersucht und entsprechend in Ansatz gebracht. Grundlage der weiteren Auswertung ist eine flurstücksgenaue Abbildung der Indikatoren, daher wurden die in Anlage 1 aufgeführten Indikatoren mit allen Flurstücken verschnitten.

Der besseren Lesbarkeit halber sind hier nur einige stellvertretende Indikatoren und deren Erläuterung aufgeführt

- ANZAHL_FLAECHEN
- AmtlicheFlaeche_FLSTK
- BEBAUTE_FLAECHEN,
Die Gebäudeflächen (BEBAUTE_FLAECHEN) und die amtlichen Flurstücksflächen (AmtlicheFlaeche_FLSTK) entstammen dem amtlichen Liegenschaftskatasterinformationssystem (ALKIS) bzw. dem vorliegenden 3D-Stadtmodell der Stadt Duisburg.
- Ewolnsgesamt, EwolnsgesamtMaennlich, EwolnsgesamtWeiblich
- NAMGMK.baerl
- Ortsteilname.duissern
Die Einwohnerzahlen sind ortsteilscharf mit den Flurstücken verschnitten worden, die Attribute NAMGMK und Ortsteilnamen geben an, ob das Flurstück in der angegebenen Gemarkung oder dem angegebenen Ortsteil liegt.
- GFZ_BERECHNET_Median
- GRZ_BERECHNET_Median
- storeysAboveAvg (gemittelte Anzahl der Geschosse oberhalb des Erdgeschosses sämtlicher Gebäude auf dem Flurstück), storeysAboveByArea

Die Grundflächenzahlen (GRZ) sowie die Geschossflächenzahlen (GFZ) wurden aus den Einzelangaben, den amtlichen Flurstücksflächen, bebauter Fläche sowie nach Flächenanteilen der Bebauungsfläche der einzelnen Gebäude auf dem Flurstück (storeysAboveAvg, storeysAboveByArea) ermittelt. Unter storeysAboveAvg versteht man die gemittelte Anzahl der Geschosse oberhalb des Erdgeschosses sämtlicher

Gebäude auf dem Flurstück. Sie ist „0“, sofern es kein Geschoss oberhalb des Erdgeschosses gibt.

- NEAR_Baeume_DIST, NEAR_Baeume_FID
- NEAR_Bildungseinrichtungen_DIST, NEAR_Bildungseinrichtungen_FID, Attribute, welche mit dem Präfix „NEAR_“ beginnen, stammen aus den Gebäudeteilen und -flächen des Amtlichen Liegenschaftskatasterinformationssystems (ALKIS). Anhand der Gebädefunktion bzw. -bezeichnung wurde gefiltert, ob es sich um einen entsprechenden Point of Interest (POI) handelt. So zählen beispielsweise Gebäude mit der Funktion bzw. Bezeichnung „Allgemeinbildende Schule“ zur Kategorie der „Bildungseinrichtungen“. Diese werden über die Attribute „NEAR_Bildungseinrichtungen_DIST“ (Luftlinienentfernung zur nächstgelegenen Einrichtung dieser Kategorie) sowie „NEAR_Bildungseinrichtungen_FID“ (OBJECTID der entsprechenden ALKIS-Geometrie) referenziert (vgl. Anhang). Dieses Verfahren ermöglicht eine einheitliche und flächendeckend vergleichbare Ermittlung von Distanzbeziehungen. Es ist jedoch zu berücksichtigen, dass die Luftlinienentfernung in Einzelfällen von tatsächlich zurückzulegenden Wegenetzdistanzen abweichen kann, da infrastrukturelle Gegebenheiten wie Straßenführungen, Barrieren oder Zugänglichkeiten nicht berücksichtigt werden. Vor dem Hintergrund einer stadtweiten, generalisierenden Bewertung wird die Luftlinienentfernung dennoch als sachgerechte und methodisch konsistente Näherung verwendet.

- NEAR_Kinderbetreuungseinrichtungen_DIST, NEAR_Kinderbetreuungseinrichtungen_FID

Diese Einrichtungen entstammen der Points-of-Interest-Datenbank „XErleben“, deren Name die folgenden Worte bzw. Wortbruchstücke enthält:

- Kindergarten
- Kindertageseinrichtung
- Kindertagesstätte
- Kita
- NEAR_Nahversorgungsbetriebe_DIST, NEAR_Nahversorgungsbetriebe_FID
Diese Daten entstammen der Einzelhandelsermittlung (EHE) für das Einzelhandels- und Zentrenkonzept 2019 [EHZK].

- NEAR_Parkraeume_DIST, NEAR_Parkraeume_FID

Diese Daten basieren ebenfalls auf der „XErleben“-Datenbank mit den folgenden Worten bzw. Wortbruchstücken:

- Parkhaus
- Parkplatz
- NEAR_Spielplaetze_DIST, NEAR_Spielplaetze_FID, NEAR_Sportanlage_oder_Kinderspielplatz_DIST, NEAR_Sportanlage_oder_Kinderspielplatz_FID

Die Sportanlagen und Spielplätze stammen von „XErleben“, wobei folgende Wortbruchstücke erlaubt sind:

- Kinderspiel
- Spielplatz
- Schwimm
- Sport mit folgenden Ausnahmen (dürfen nicht enthalten sein):

Denkmal, Dezernat, Fashion, Hotel, Sportcenter, Sportmanagement, Sportswear, Sportwetten, Mit Beschränkter Haftung, GmbH, KG, Klinik, Parkplatz, Transport, UG, Verlag

- **NEAR_VIA_Strassen_DIST**
Die Daten der Infrastruktur (Straßen) nebst der dazugehörigen Straßenklasse (z.B. Weg, Erschließungsstraße, Hauptstraße, Autobahn, etc.) stammen aus der „VIAVIS“-Datenbank.
- **Image:**
Für den Indikator „Image“ wurden durch den Gutachterausschuss insgesamt 515 Flurstücke ausgewählt, um ein möglichst repräsentatives Abbild der unterschiedlichen Wohnlagen im Stadtgebiet zu erhalten.
Die Bewertung dieser Flurstücke erfolgte durch die ehrenamtlichen Mitglieder des Gutachterausschusses unabhängig und anonym. Ziel war es, eine fachlich fundierte Einschätzung der allgemeinen Beliebtheit unterschiedlicher Wohnlagen zu erhalten. Zur Überprüfung der Konsistenz dieser Bewertungen wurde ein Reliabilitätstest durchgeführt. Die Reliabilität, verstanden als die Zuverlässigkeit und Wiederholbarkeit der Bewertungsergebnisse, wurde mittels Krippendorffs Alpha bestimmt.
Der ermittelte Wert liegt im akzeptablen Bereich und zeigt eine hinreichende Übereinstimmung zwischen den Bewertungen der beteiligten Sachverständigen. Unter Berücksichtigung der naturgemäß vorhandenen Interpretationsspielräume bei der Bewertung qualitativer Faktoren ist dieses Ergebnis als sachgerecht einzustufen.
Auf Grundlage dieser Bewertung konnte für jedes der 515 Flurstücke eine einheitliche Lageklassifikation als Referenzdatensatz für die weitere Modellierung abgeleitet werden.

5. Auswahl und Training der modellgestützten Verfahren

Künstliche Intelligenz umfasst eine Vielzahl von Methoden und Techniken, die darauf abzielen, aus Daten Muster zu erkennen und auf dieser Grundlage Vorhersagen oder Klassifikationen vorzunehmen. Die im vorliegenden Kontext eingesetzten Verfahren basieren dabei insbesondere auf statistischen Lernmethoden, die Zusammenhänge in den Daten modellieren und für die Bewertung nutzbar machen.

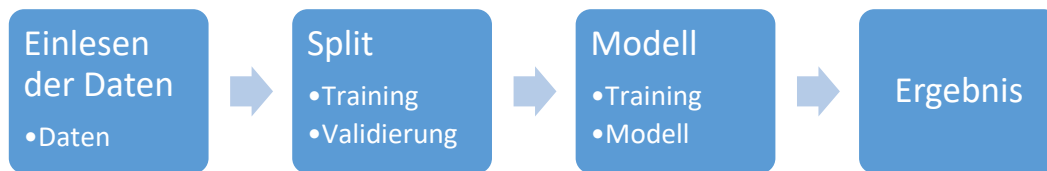
Die Auswahl und das nachfolgende Training wurden für die folgenden modellgestützten Verfahren durchgeführt, die im Folgenden kurz erläutert werden:

- **Neuronale Netzwerke** sind eine Klasse von Algorithmen, die von der Funktionsweise des menschlichen Gehirns inspiriert sind. Sie bestehen aus künstlichen Neuronen, die in Schichten angeordnet sind. Jedes Neuron ist mit anderen Neuronen verbunden und hat eine Gewichtung. Durch das Training des Netzwerks werden diese Gewichtungen angepasst, um Muster in den Daten zu erkennen.
- **Random Forest** ist eine Ensemble-Methode, die mehrere Entscheidungsbäume kombiniert, um genauere und robustere Vorhersagen zu treffen. Jeder Baum im Wald wird auf einem zufälligen Unterdatensatz des Gesamtdatensatzes trainiert. Die Vorhersage wird dann durch die Aggregation der Ergebnisse aller Bäume getroffen, was zu einer verbesserten Leistung führt.
- **Support Vector Machine (SVM)** ist ein Machine Learning-Algorithmus, der für Klassifikations- und Regressionsaufgaben verwendet wird. Das Hauptziel der SVM ist es, eine Trennlinie oder Hyperebene zu finden, die den Datensatz in Klassen oder

Gruppen aufteilt. Die SVM kann sowohl für lineare als auch nichtlineare Probleme eingesetzt werden.

Zur Erstellung der Wohnlagenkarte wurden flurstücksscharf sämtliche der in Anlage 1 genannten Attribute gesammelt und auch weitere Attribute ebenfalls mit den Flurstücken verschnitten.

Der Ablauf der einzelnen Schritte stellt sich grob dabei wie folgt dar:



Die eingelesenen Daten werden in Trainings- und Validierungsdaten unterteilt. Die Trainingsdaten dienen dem Aufbau der jeweiligen Modelle, während die Validierungsdaten zur Bewertung der Modellgüte herangezogen werden. Die Validierungsdaten wurden während des Trainings nicht zur Anpassung der Modellparameter verwendet. Die Trennung der Datensätze ermöglicht es, die Modellleistung auf Daten zu überprüfen, die im Trainingsprozess nicht unmittelbar verwendet wurden. Ein separater Testdatensatz wird im vorliegenden Fall nicht verwendet, sodass die Validierungsdaten als maßgebliche Grundlage für die vergleichende Beurteilung der Modelle dienen.

Im Rahmen der Modellentwicklung wurden systematische Variationen der Eingangsattribute durchgeführt. Dabei wurde jeweils untersucht, ob einzelne Attribute oder Attributgruppen die Modellgüte verbessern oder verschlechtern, während die übrigen Variablen unverändert beibehalten wurden. So wurden beispielsweise die Klassen „Bäume“, „Wälder“ und „Wiesen“ gegenüber aggregierten Grünvolumenkennwerten bzw. Isotopen der Stadt Duisburg bevorzugt, da diese in den Modellvergleichen zu stabileren und genaueren Ergebnissen führten.

Darüber hinaus wurden verschiedene Modellansätze unter Verwendung identischer Eingabeattribute miteinander verglichen, um die Leistungsfähigkeit unterschiedlicher Verfahren beurteilen zu können.

Für jedes Modell wurden mehrere Trainingsdurchläufe mit identischen Modellparametern durchgeführt, um zufallsbedingte Einflüsse im Trainingsprozess zu reduzieren und die Stabilität der Ergebnisse zu prüfen.

Die Bewertung der Modelle erfolgte auf Grundlage der in den Validierungsdaten erzielten Modellgüte. Modelle mit höheren Genauigkeitswerten wurden gegenüber solchen mit geringerer Genauigkeit bevorzugt. Die Genauigkeit dient dabei als zentrale Kennzahl zur vergleichenden Bewertung der verschiedenen Modellansätze.

Die Wohnlagen wurden methodisch in folgenden Schritten ermittelt:

- 1. Datenaufbereitung und Zusammenführung**
Verschneidung der 515 Referenzdatensätze des Gutachterausschusses (mit Wohnlagenklassifikation) mit den flurstücksscharfen Indikatoren zur Erstellung eines analysierbaren Datensatzes.
- 2. Aufteilung in Trainings- und Validierungsdaten**
Unterteilung des Datensatzes in Trainings- und Validierungsanteile (70 % / 30 %) zur getrennten Modellentwicklung und Bewertung.

3. Modellentwicklung und Bewertung

Training der Modelle auf Basis der Trainingsdaten sowie Bewertung der Modellgüte anhand der Validierungsdaten.

Die modellgestützten Verfahren wurden jeweils auf Basis der Trainingsdaten, einschließlich der zugehörigen Wohnlagenklassifikation, trainiert. Die Modellgüte wurde anschließend durch Vergleich der modellierten Wohnlagen mit den tatsächlichen Klassifikationen in den Validierungsdaten bestimmt.

Zur Berücksichtigung zufallsbedingter Einflüsse im Trainingsprozess wurden für jedes Modell mehrere Trainingsdurchläufe mit identischen Parametereinstellungen durchgeführt. Für die nachfolgende Auswertung wurde jeweils der Modelllauf mit der höchsten Genauigkeit in den Validierungsdaten herangezogen.

Die nachfolgend dargestellten Ergebnisse beziehen sich somit auf die jeweils besten erzielten Modellgüten je Modelltyp.

Im Rahmen der Modellentwicklung wurden bewusst unterschiedliche Modellklassen berücksichtigt, um sowohl lineare als auch nichtlineare Zusammenhänge sowie unterschiedliche Modellansätze abzudecken. Hierzu zählen insbesondere baumbasierte Verfahren, kernelbasierte Methoden sowie neuronale Netzwerke.

Die Auswahl der Modelle erfolgte dabei nicht mit dem Ziel der Vollständigkeit aller verfügbaren Verfahren, sondern orientierte sich an der praktischen Eignung für die vorliegenden Datenstrukturen sowie an der Reproduzierbarkeit und Interpretierbarkeit der Ergebnisse.

- Es wurden folgende Modelle mit den jeweils besten Ergebnissen (Maximum über 25 Trainingsdurchläufe) berücksichtigt: Support Vector Machine (SVC),
- Random Forest (RF) mit 245 Bäumen,
- die übrigen Modelle als neuronale Netzwerke (NN) sind nach der Optimierungsfunktion des NN benannt (vergleiche jeweils https://www.tensorflow.org/api_docs/python/tf/keras/optimizers):

Ada-delta	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion Ada-delta
Ada-factor	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion Ada-factor
Ada-grad	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion Ada-grad
Adam	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion Adam
Adamax	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion Adamax
AdamW	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion AdamW
Nadam	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion NAdam
RF	scikit-learn: Random Forest Modell auf Basis von mehreren unkorrelierten Entscheidungsbäumen (RandomForestClassifier)
RMS-prop	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion RMS-prop
SGD	tensorflow: Neuronales Netzwerk mit Optimierungsfunktion SGD
SVC	scikit-learn: Support Vector Machine (SupportVectorClassification)

Für die eingesetzten Modelle wurden jeweils modelltypische Hyperparameter definiert, die den Trainingsprozess und die Modellkomplexität steuern. Die Auswahl dieser Hyperparameter erfolgte auf Grundlage etablierter Standardwerte der jeweiligen Bibliotheken (insbesondere *scikit-learn* und *TensorFlow*) sowie durch gezielte Variation einzelner Parameter im Rahmen der Modellentwicklung.

Ziel war es, eine ausgewogene Balance zwischen Modellkomplexität, Generalisierungsfähigkeit und Rechenaufwand zu erreichen. Eine vollständige

systematische Durchsuchung aller möglichen Parameterkombinationen wurde nicht durchgeführt, da dies angesichts der Vielzahl möglicher Konfigurationen zu unverhältnismäßig hohem Rechenaufwand führen würde und erfahrungsgemäß nur begrenzte zusätzliche Erkenntnisse liefert.

Stattdessen wurden für die einzelnen Modellklassen exemplarisch relevante Hyperparameter systematisch variiert und deren Einfluss auf die Modellgüte analysiert. Die Auswahl der finalen Parameter erfolgte anhand der erreichten Modellgüte sowie der Stabilität der Ergebnisse über mehrere Trainingsdurchläufe hinweg.

Die Genauigkeit der Berechnungsergebnisse basiert auf dem Anteil korrekt klassifizierter Wohnlagen (Accuracy). Diese Kennzahl gibt an, in welchem Umfang die vom Modell prognostizierten Wohnlagen mit den tatsächlich zugrunde gelegten Klassifikationen übereinstimmen.

Es ist zu beachten, dass die Accuracy insbesondere bei ungleich verteilten Klassen nur eine eingeschränkte Aussagekraft besitzen kann. Daher wurde sie im vorliegenden Zusammenhang primär als vergleichende Kennzahl zur Gegenüberstellung der verschiedenen Modelle herangezogen.

Die Interpretation der Accuracy erfolgt unter Berücksichtigung der Verteilung der Wohnlagenklassen.

Die Ergebnisse der Modellvergleiche sind in der nachfolgenden Tabelle dargestellt:

Modell	Accuracy
Random Forest (RF)	69,48%
SGD	64,94 %
Adamax (NN)	64,29 %
Adam (NN)	62,34 %
Nadam (NN)	61,69 %
AdamW (NN)	61,04 %
RMSprop (NN)	60,39 %
AdaGrad (NN)	59,09 %
AdaDelta (NN)	37,66 %
Support Vector Machine (SVC)	35,45 %
AdaFactor (NN)	24,67 %

Für die Erstellung der Duisburger Wohnlagenkarte wurde untersucht, welche der vorgenannten Methoden die für diesen Anwendungszweck höchstmögliche Aussagefähigkeit besitzt:

Neuronales Netzwerk

Die untersuchten neuronalen Netzwerke erzielten im Vergleich der betrachteten Modelle die zweithöchsten Werte hinsichtlich der Modellgüte.

Die unterschiedlichen Varianten der neuronalen Netzwerke unterscheiden sich primär in dem verwendeten Optimierungsverfahren (z. B. Adam, SGD, RMSprop), welches zur Minimierung der Fehlerfunktion (Loss Function) während des Trainings eingesetzt wird. Diese Optimierungsverfahren beeinflussen vor allem die Konvergenzgeschwindigkeit und Stabilität des Trainingsprozesses, nicht jedoch die grundsätzliche Modellstruktur.

Trotz vergleichsweise guter Ergebnisse sind neuronale Netzwerke in ihrer internen Funktionsweise nur eingeschränkt nachvollziehbar („Blackbox“), da die Modellentscheidungen auf einer Vielzahl nichtlinearer Transformationen beruhen und sich nur begrenzt direkt interpretieren lassen.

Random Forest

Das Random-Forest-Modell wurde auf Grundlage der im Vergleich höchsten erzielten Modellgüte für die Version 1.1 der Wohnlagenkarte ausgewählt.

Random Forest gehört zu den baumbasierten Ensembleverfahren und basiert auf einer Vielzahl einzelner Entscheidungsbäume, die unter Einsatz von Zufallsmechanismen (z. B. Stichprobenziehung und zufällige Merkmalsauswahl) unabhängig voneinander trainiert werden. Jeder Entscheidungsbaum wird dabei auf einer unterschiedlichen Teilmenge der Daten aufgebaut.

Die Vorhersage erfolgt durch Aggregation der Einzelvorhersagen aller Bäume (Mehrheitsentscheidung), wodurch robuste und gegenüber Ausreißern weniger anfällige Ergebnisse erzielt werden.

Hinsichtlich der Nachvollziehbarkeit ist festzuhalten, dass Random-Forest-Modelle aufgrund ihrer Ensemble-Struktur nicht als vollständig transparent einzustufen sind. Allerdings bieten sie im Vergleich zu komplexeren Modellansätzen wie neuronalen Netzwerken erweiterte Möglichkeiten zur Interpretation, beispielsweise durch die Auswertung von Merkmalsbedeutungen oder die Analyse einzelner Entscheidungsbäume.

Vor diesem Hintergrund stellt das Random-Forest-Modell im vorliegenden Kontext einen geeigneten Kompromiss zwischen hoher Modellgüte und begrenzter, jedoch vorhandener Interpretierbarkeit dar.

Support Vector Machine

Das Support-Vector-Machine-Modell erzielte im Vergleich der untersuchten Verfahren eine geringe Modellgüte und wurde daher im weiteren Verlauf nicht weiter berücksichtigt.

Das Random-Forest-Modell mit 245 Bäumen erreichte im Rahmen der Cross-Validation die höchste Modellgüte (Accuracy: 69,48 %) und wurde daher für die weitere Anwendung ausgewählt.

Bei Anwendung dieses Modells auf den gesamten zugrunde liegenden Datensatz ergibt sich eine Übereinstimmung zwischen modellierten und referenzierten Wohnlagen von ca. 91 %. Hierbei ist zu berücksichtigen, dass es sich nicht um eine unabhängige Modellbewertung handelt, da diese Daten bereits in den Trainings- und Validierungsprozess eingeflossen sind. Die Kennzahl dient daher ausschließlich der ergänzenden Beschreibung der Modellanpassung an die vorhandenen Daten.

Im Zuge der Anwendung zeigte sich, dass das Modell in einzelnen Bereichen zu einer sehr kleinteiligen Differenzierung der Wohnlagen führt. Insbesondere konnten Wohnlagenwechsel

auf Ebene einzelner Flurstücke („Inselwohnlagen“) entstehen, deren Klassifikation von der umgebenden Wohnlage abweicht.

Diese Effekte sind typisch für datengetriebene Modellansätze mit hoher Auflösung und resultieren aus lokalen Unterschieden in den zugrunde liegenden Eingangsmerkmalen. Im weiteren Verfahren wurde daher eine fachliche Bewertung dieser Ergebnisse vorgenommen.

Unter Abwägung von Modellgüte, Stabilität der Ergebnisse und begrenzter Interpretierbarkeit wurde das Random-Forest-Modell als geeignetes Verfahren zur Ermittlung der Wohnlagen im Stadtgebiet Duisburg ausgewählt.

6. Homogenisierung und kartographische Darstellung

Die durch das eingesetzte Modell ermittelte Wohnlagenklassifikation erfolgte zunächst auf flurstücksscharfer Ebene. Aufgrund der detaillierten Datenlage und der hohen Modellauflösung zeigte sich im Ergebnis eine sehr feingliedrige Differenzierung der Wohnlagen.

Insbesondere traten lokal begrenzte Abweichungen („Insellagen“) sowie kleinräumige Lageklassensprünge auf, bei denen einzelne Flurstücke eine andere Wohnlagenklassifikation aufwiesen als die unmittelbar angrenzenden Bereiche.

Diese Effekte sind typisch für datengetriebene Modellansätze mit hoher räumlicher Auflösung und resultieren aus lokalen Unterschieden in den zugrunde liegenden Eingangsmerkmalen.

Über sämtliche Datensätze, die voreingeschätzt waren, konnte man eine Genauigkeit von rd. 76 % ermitteln (Abbildung 1). Es zeigten sich dabei jedoch teilweise Lageklassensprünge von bis zu 3 Lageklassen zwischen Imageeinschätzung und KI-Einschätzung.

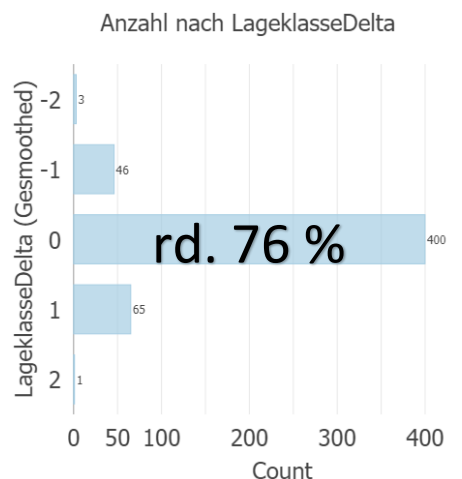


Abbildung 1

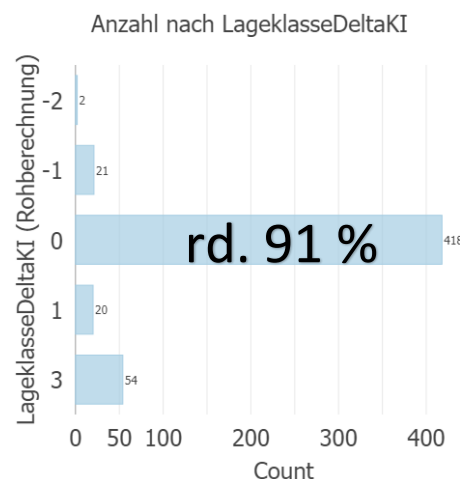


Abbildung 2

Zur Erhöhung der räumlichen Konsistenz und der kartographischen Interpretierbarkeit wurde daher eine nachgelagerte Generalisierung der Ergebnisse vorgenommen.

Die zunächst flurstücksscharf ermittelte Wohnlagenklassifikation wurde in einem zweiten Schritt auf Baublockebene aggregiert und vereinheitlicht.

Hierzu wurde für jeden Baublock der Mittelwert der zugehörigen flurstücksbezogenen Wohnlagen bestimmt und anschließend allen Flurstücken innerhalb des jeweiligen

Baublocks zugewiesen. Die kleinste räumliche Einheit der Wohnlagenkarte ist somit der Baublock.

Bei diesem Vorgehen handelt es sich methodisch nicht um ein klassisches Smoothing im Sinne gleitender Mittelwerte oder kernelbasierter Glättungsverfahren, sondern um eine räumliche Aggregation mit anschließender Homogenisierung innerhalb definierter Gebietseinheiten.

Durch die Aggregation werden kleinräumige Ausreißer („Insellagen“) sowie abrupte Wohnlagensprünge reduziert, wodurch eine stärker generalisierte und räumlich konsistente Abbildung der Wohnlagen erreicht wird.

Nach Durchführung der Aggregation ergibt sich eine Modellgüte (Accuracy) von rd. 78 % (vgl. Abbildung 2). Die geringere Genauigkeit nach der Aggregation ist eine Folge der angestrebten räumlichen Generalisierung

Die resultierende Übereinstimmung zwischen modellierten und referenzierten Wohnlagen liegt in einem Bereich, der als fachlich angemessen und konsistent bewertet werden kann.

Die erstellte Wohnlagenkarte stellt sich wie folgt dar:

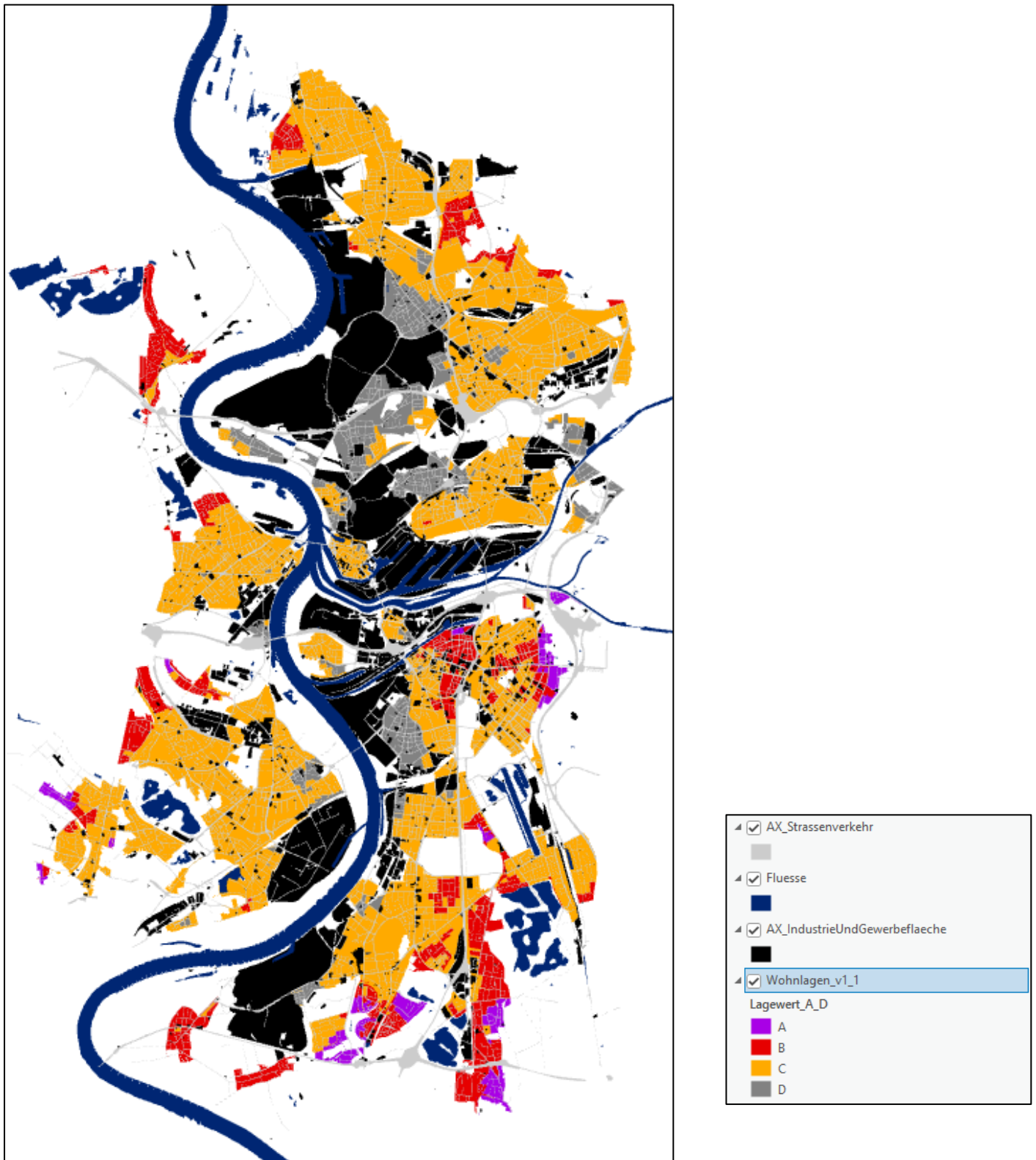


Abbildung 3

Die Wohnlagen wurden bei der erstmaligen Erstellung in die vier Klassen A (sehr gute Lage), B (gute Lage), C (mittlere Lage) und D (einfache Lage) unterteilt. In den weiteren Veröffentlichungen und Darstellungen wird die Einteilung in sehr gute bis einfache Lage verwendet.

7. Qualitäts- und Plausibilitätsprüfung

Im Anschluss an die modellgestützte Ermittlung und die räumliche Generalisierung der Wohnlagen wurde eine umfassende Qualitäts- und Plausibilitätsprüfung durchgeführt.

Hierzu wurden sämtliche Baublöcke im Stadtgebiet Duisburg auf auffällige Abweichungen untersucht. Der Fokus lag insbesondere auf der Identifikation von:

- abrupten Wohnlagensprüngen
- systematischen Abweichungen gegenüber den zugrunde liegenden Referenzdaten (Imagebewertungen)
- räumlich nicht plausiblen Klassifikationen

Die Überprüfung erfolgte insbesondere auf Grundlage der kartographischen Darstellung der Wohnlagen sowie im Vergleich mit den durch den Gutachterausschuss vorgenommenen Referenzbewertungen.

Im Rahmen dieser Prüfung konnten keine systematischen oder fachlich nicht erklärbaren Abweichungen festgestellt werden. Die verbleibenden Unterschiede zwischen modellierten und referenzierten Wohnlagen sind vielmehr auf lokale Variationen in den zugrunde liegenden Eingangsmerkmalen zurückzuführen und im Rahmen der Modellierung plausibel erklärbar.

Diese Abweichungen können als Residuen verstanden werden, also als Differenzen zwischen modellierter Klassifikation und Referenzbewertung. Eine systematische Verzerrung oder strukturelle Fehlklassifikation konnte nicht festgestellt werden.

Die Ergebnisse der Qualitätsprüfung bestätigen somit die insgesamt hohe Konsistenz und Plausibilität der ermittelten Wohnlagen.

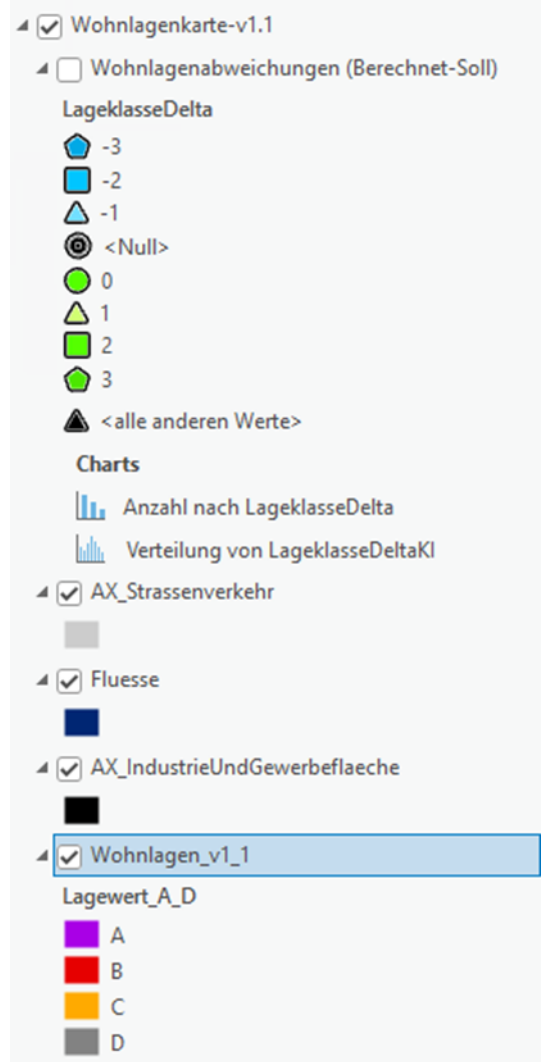
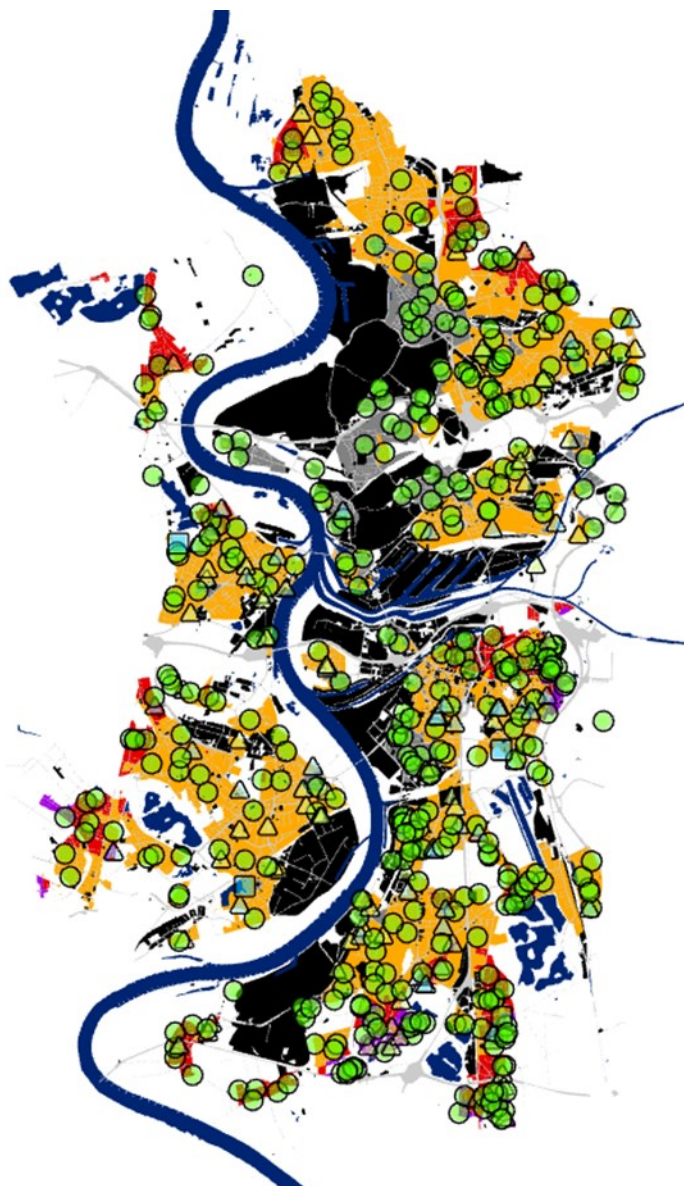


Abbildung 4

Es wurden keine auffälligen Differenzen in der Wohnlagenkarte gefunden.

Die Anzahl der unterschiedlichen Lageneinstufungen bezogen auf die Anzahl der Flurstücke über das Stadtgebiet stellt sich wie folgt dar:

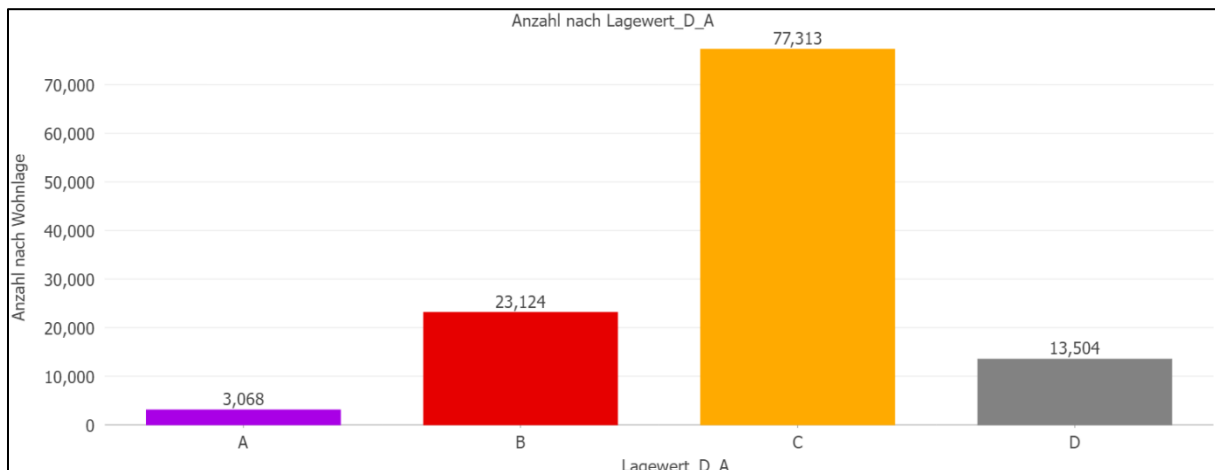


Abbildung 5

8. Literaturverzeichnis

- [EHZK]: Einzelhandels- und Zentrenkonzept der Stadt Duisburg, 2019
- [Entscheidungsbaum] <https://de.wikipedia.org/wiki/Entscheidungsbaum>

9. Anlage 1:

- ANZAHL_FLAECHEN,
- AmtlicheFlaeche_FLSTK,
- BEBAUTE_FLAEICHE,
- Ewolnsgesamt, EwolnsgesamtMaennlich, EwolnsgesamtWeiblich,
- GFZ_BERECHNET_Median,
- GRZ_BERECHNET_Median,
- Hochwert,
- NAMGMK.baerl, NAMGMK.beeck, NAMGMK.duisburg, NAMGMK.hamborn, NAMGMK.homberg, NAMGMK.huckingen, NAMGMK.kaldenhausen, NAMGMK.meiderich, NAMGMK.muendelheim, NAMGMK.rheinhausen, NAMGMK.ruhrort, NAMGMK.rumeln, NAMGMK.walsum,
- NEAR_Baeume_DIST, NEAR_Baeume_FID,
- NEAR_Bahnhoefe_DIST, NEAR_Bahnhoefe_FID,
- NEAR_Bildungseinrichtungen_DIST, NEAR_Bildungseinrichtungen_FID,
- NEAR_Einkaufszentren_DIST, NEAR_Einkaufszentren_FID,
- NEAR_Feuerwehr_DIST, NEAR_Feuerwehr_FID,
- NEAR_Firmen_DIST, NEAR_Firmen_FID,
- NEAR_Freizeiteinrichtungen_DIST, NEAR_Freizeiteinrichtungen_FID,
- NEAR_Handel_DIST, NEAR_Handel_FID,
- NEAR_Industrie_100m_DIST, NEAR_Industrie_100m_FID, NEAR_Industrie_100m_JN, NEAR_Industrie_150m_DIST, NEAR_Industrie_150m_FID, NEAR_Industrie_150m_JN, NEAR_Industrie_50m_DIST, NEAR_Industrie_50m_FID, NEAR_Industrie_50m_JN, NEAR_Industrie_DIST, NEAR_Industrie_FID, NEAR_Industrie_Klasse,

- NEAR_Infrastruktur_DIST, NEAR_Infrastruktur_FID,
- NEAR_Inseln_DIST, NEAR_Inseln_FID,
- NEAR_Kinderbetreuungseinrichtungen_DIST,
NEAR_Kinderbetreuungseinrichtungen_FID,
- NEAR_Krankenhaeuser_DIST, NEAR_Krankenhaeuser_FID,
- NEAR_Medizinische_Einrichtungen_DIST, NEAR_Medizinische_Einrichtungen_FID,
- NEAR_Nahversorgungsbetriebe_DIST, NEAR_Nahversorgungsbetriebe_FID,
- NEAR_Parken_DIST,
- NEAR_Parkraeume_DIST, NEAR_Parkraeume_FID,
- NEAR_Parks_DIST, NEAR_Parks_FID,
- NEAR_Religioese_Einrichtungen_DIST, NEAR_Religioese_Einrichtungen_FID,
- NEAR_Reparaturgebiet_DIST, NEAR_Reparaturgebiet_FID,
- NEAR_Schulen_DIST, NEAR_Schulen_FID,
- NEAR_Sicherheit_und_Ordnung_DIST, NEAR_Sicherheit_und_Ordnung_FID,
- NEAR_Spielplaetze_DIST, NEAR_Spielplaetze_FID,
- NEAR_Sportanlage_oder_Kinderspielplatz_DIST,
NEAR_Sportanlage_oder_Kinderspielplatz_FID,
- NEAR_Strecken_Bahn_DIST, NEAR_Strecken_Bahn_FID,
- NEAR_Tankstellen_DIST, NEAR_Tankstellen_FID,
- NEAR_VIA_Strassen_DIST, NEAR_VIA_Strassen_FID,
- NEAR_VIA_Strassen_VIA_TYP_KLASSE,
- NEAR_Verwaltungsgebäude_DIST, NEAR_Verwaltungsgebäude_FID,
- NEAR_Waelder_DIST, NEAR_Waelder_FID,
- NEAR_Wiesen_DIST, NEAR_Wiesen_FID,
- NEAR_Wirtschaftsgebäude_DIST, NEAR_Wirtschaftsgebäude_FID,
- NEAR_Wochenmaerkte_DIST, NEAR_Wochenmaerkte_FID,
- NORDWERT,
- OrtsteilHOCHWERT, OrtsteilRECHTSWERT,
- OrtsteilNummer,
- OrtsteilShape_Area, OrtsteilShape_Length,
- Ortsteilname.aldenrade, Ortsteilname.althamborn, Ortsteilname.althomburg,
Ortsteilname.altstadt, Ortsteilname.altwalsum, Ortsteilname.baerl, Ortsteilname.beeck,
Ortsteilname.beeckerwerth, Ortsteilname.bergheim, Ortsteilname.bissingheim,
Ortsteilname.bruckhausen, Ortsteilname.buchholz, Ortsteilname.dellviertel,
Ortsteilname.duissern, Ortsteilname.fahrn, Ortsteilname.friemersheim,
Ortsteilname.grossenbaum, Ortsteilname.hochemmerich, Ortsteilname.hochfeld,
Ortsteilname.hochheide, Ortsteilname.huckingen, Ortsteilname.huettenheim,
Ortsteilname.kasslerfeld, Ortsteilname.laar, Ortsteilname.marxloh,
Ortsteilname.mittelmeiderich, Ortsteilname.muendelheim, Ortsteilname.neudorfnord,
Ortsteilname.neudorfsued, Ortsteilname.neuenkamp, Ortsteilname.neumuehl,
Ortsteilname.obermarxloh, Ortsteilname.obermeiderich, Ortsteilname.overbruch,
Ortsteilname.rahm, Ortsteilname.rheinhausenmitte, Ortsteilname.roettgersbach,
Ortsteilname.ruhrort, Ortsteilname.rumelnkaldenhausen, Ortsteilname.ungelsheim,
Ortsteilname.untermeiderich, Ortsteilname.vierlinden,
Ortsteilname.wanheimangerhausen, Ortsteilname.wanheimerort, Ortsteilname.wedau,
Ortsteilname.wehofen,
- SHAPE_Area, SHAPE_Length,
- Strasse_Nacht_DBA,
- Strasse_Nacht_LAERM_DBA_KZ, Strasse_Nacht_LAERM_area,
Strasse_Nacht_LAERM_count,

- Strasse_Tag_DBA,
- Strasse_Tag_LAERM_DBA_KZ, Strasse_Tag_LAERM_area,
Strasse_Tag_LAERM_count,
- storeysAboveAvg (gemittelte Anzahl der Geschosse oberhalb des Erdgeschosses
sämtlicher Gebäude auf dem Flurstück), storeysAboveByArea

10. Anhang

- [Gebaeudebezeichnungen_TableToExcel.xlsx]: Zuordnung Gebäudebezeichnungen
ALKIS <-> POIs für NEAR-Klassen