

BikeDNA Report

Copenhagen: GeoDanmark



Report generated on: 2022-12-10 15:07:25

1a. Load and Process OSM data

This notebook:

- Loads the polygon defining the study area and then creates a grid overlay for the study area.
- Downloads street network data for the study area using OSMnx.
- Creates a network only with bicycle infrastructure (with queries defined in `config.yml`).
- Creates additional attributes in the data to be used in the analysis.

Sections

- [Load data for study area and define analysis grid](#)
 - [Load and process OSM data](#)
-

Load data for study area and define analysis grid

This step:

- Loads settings for the analysis from the config file.
- Reads data for the study area.
- Creates a grid overlay of the study area, with grid cell size as defined in configuration file `config.yml`.

Read in data for study area

The study area is defined by the user-provided polygon. It will be used for the computation of *global* results, i.e. for the entire study area.

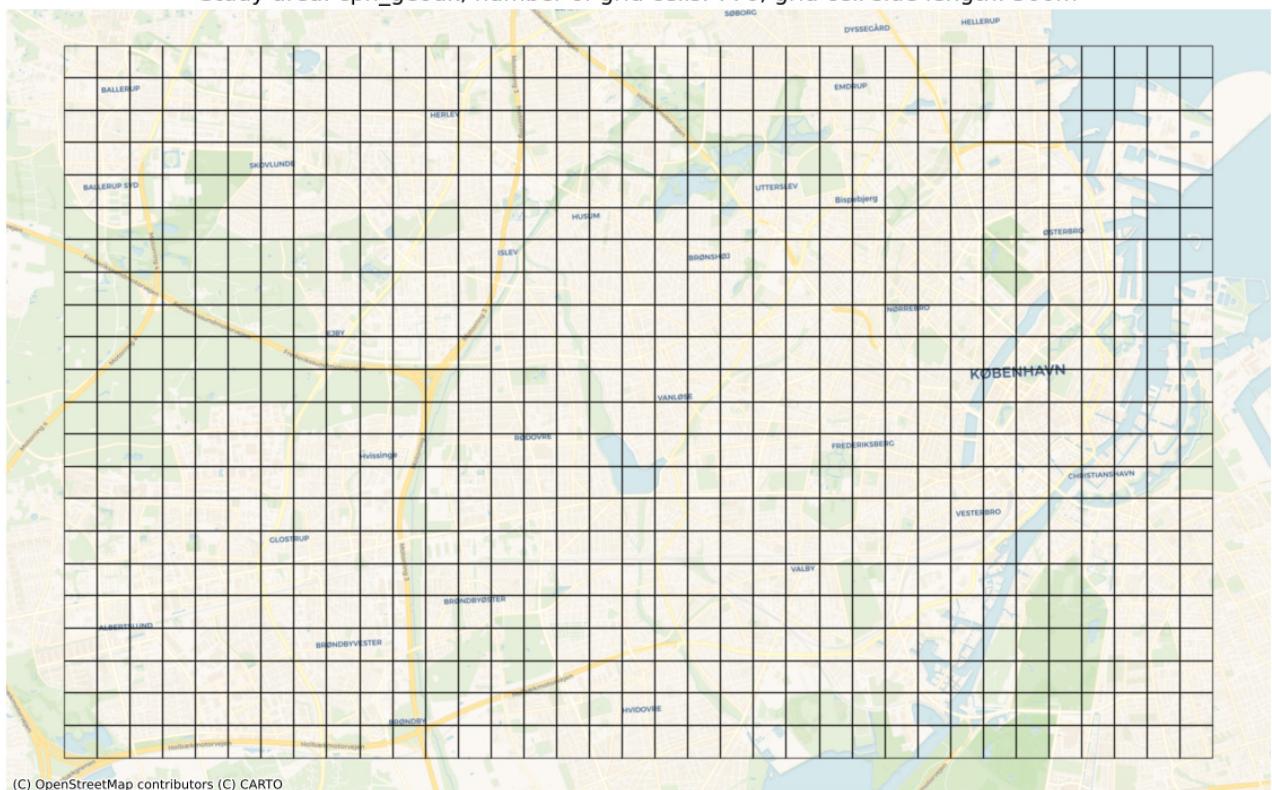
The size of the study area is 181.38 km².

Create analysis grid

The grid contains 770 square cells with a side length of 500 m and an area of 0.25 km².

This grid will be used for local (grid cell level) analysis:

Study area: cph_geodk, number of grid cells: 770, grid cell side length: 500m



Load and process OSM data

This step:

- Downloads data from OpenStreetMap using OSMnx.
- Projects the data to the chosen CRS.
- Creates a subnetwork consisting only of bicycle infrastructure.
- Classifies all edges in the bicycle network based on whether they are protected or unprotected bicycle infrastructure, how they have been digitized, and whether they allow for bidirectional travel or not
- Simplifies the network (*to read more about the modified OSMnx simplification (Boeing, 2017) used here, we refer to this [GitHub repository](#) which contains both the simplification functions, explanation of the logic and a demonstration*).
- Creates copies of all edge and node datasets indexed by their intersecting grid cell.

Successfully created network only with bicycle infrastructure!

Number of edges where 'bicycle_bidirectional' is False: 33552 out of 50973 (65.82%)

Number of edges where 'bicycle_bidirectional' is True: 17421 out of 50973 (34.18%)

Number of edges where the geometry type is 'centerline': 25738 out of 50973 (50.49%)

Number of edges where the geometry type is 'true_geometries': 25235 out of 50973 (49.51%)

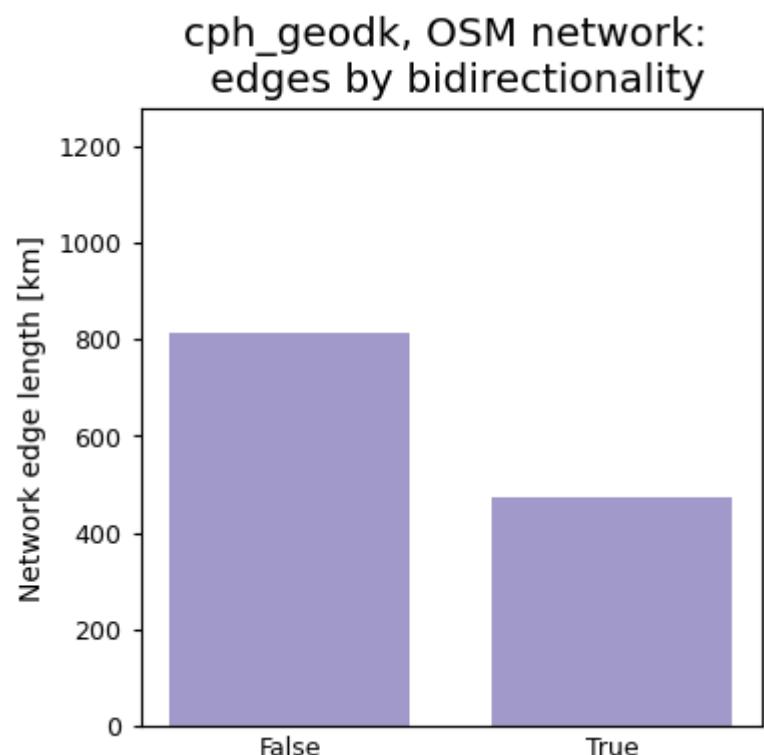
Number of edges where the protection level is 'protected': 46585 out of 50973 (91.39%)

Number of edges where the protection level is 'unprotected': 3725 out of 50973 (7.31%)

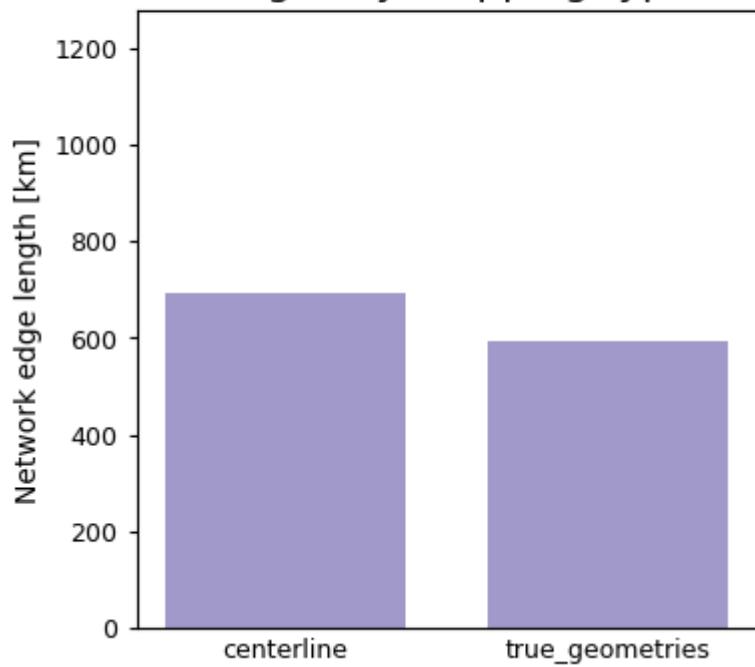
Number of edges where the protection level is 'mixed': 663 out of 50973 (1.3%)

The graph covers an area of 179.70 km².

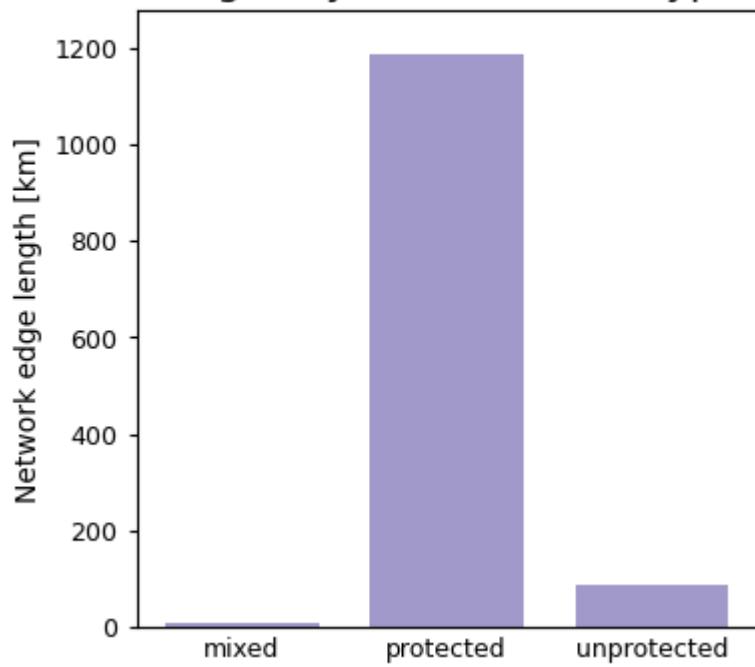
The length of the OSM network with bicycle infrastructure is 1063.46 km.

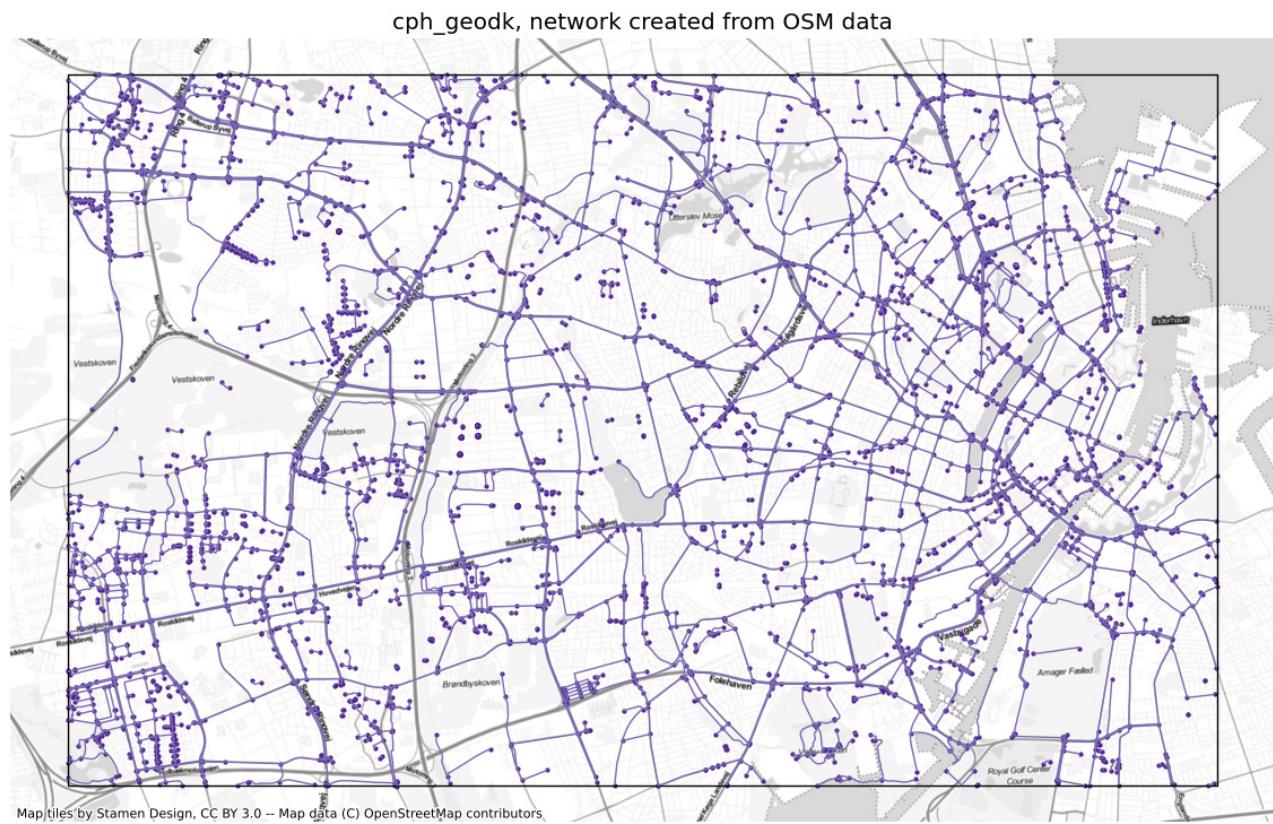


cph_geodk, OSM network:
edges by mapping type



cph_geodk, OSM network:
edges by infrastructure type





Time of analysis: Sat, 10 Dec 2022 00:15:49

1b. Intrinsic Analysis of OSM Bicycle Network Data

This notebook analyses the quality of OSM data on bicycle infrastructure for a given area. The quality assessment is *intrinsic*, i.e. based only on the input data set, and making no use of information external to the data set. For an extrinsic quality assessment that compares the OSM data set to a user-provided reference data set, see the notebooks 3a and 3b.

The analysis assesses the *fitness for purpose* ([Barron et al., 2014](#)) of OSM data for a given area. Outcomes of the analysis can be relevant for bicycle planning and research - especially for projects that include a network analysis of bicycle infrastructure, in which case the topology of the geometries is of particular importance.

Since the assessment does not make use of an external reference dataset as the ground truth, no universal claims of data quality can be made. The idea is rather to enable those working with bicycle networks based on OSM to assess whether the data is good enough for their particular use case. The analysis assists in finding potential data quality issues, but leaves the final interpretation of the results to the user.

The notebook makes use of quality metrics from a range of previous projects investigating OSM/VGI data quality, such as [Antoniou & Skopeliti \(2015\)](#), [Fonte et al. \(2017\)](#) and [Fester et al. \(2020\)](#).

Familiarity required

For a correct interpretation of some of the metrics for spatial data quality, some familiarity with the area is necessary.

Sections

- [Data completeness](#)
 - Network density
- [OSM tag analysis](#)
 - Missing tags
 - Incompatible Tags
 - Tagging patterns
- [Network topology](#)
 - Simplification outcome
 - Dangling nodes
 - Over/undershoots

- Missing intersection nodes
- Network components
 - Disconnected components
 - Potential missing links
- Summary
- Save results

Data completeness

Network Density

In this setting, network density refers to the length of edges or number of nodes per square kilometer (i.e., the definition of network density usually used when looking at street networks, which is distinct from the definition usually found in graph theory). Network density without comparing to a reference dataset does not in itself indicate spatial data quality. For anyone familiar with the study area, network density can however indicate whether parts of the area appear to be under- or over-mapped and is thus included here.

Method

The density here is not based on the geometric length of edges, but instead on the computed length of the infrastructure. For example, a 100-meter-long bidirectional path contributes with 200 meters of bicycle infrastructure. With `compute_network_density`, the number of elements (nodes; dangling nodes; and total infrastructure length) per unit area is calculated. The density is computed twice: first for the study area for both the entire network ('global density'), then for each of the grid cells ('local density'). Both global and local densities are computed for the entire network and for respectively protected and unprotected infrastructure.

Interpretation

Since the analysis conducted here is intrinsic, i.e., it makes no use of external information, it cannot be known whether a low-density value is due to incomplete mapping, or due to actual lack of infrastructure in the area. However, a comparison of the grid cell density values can provide some insights, for example:

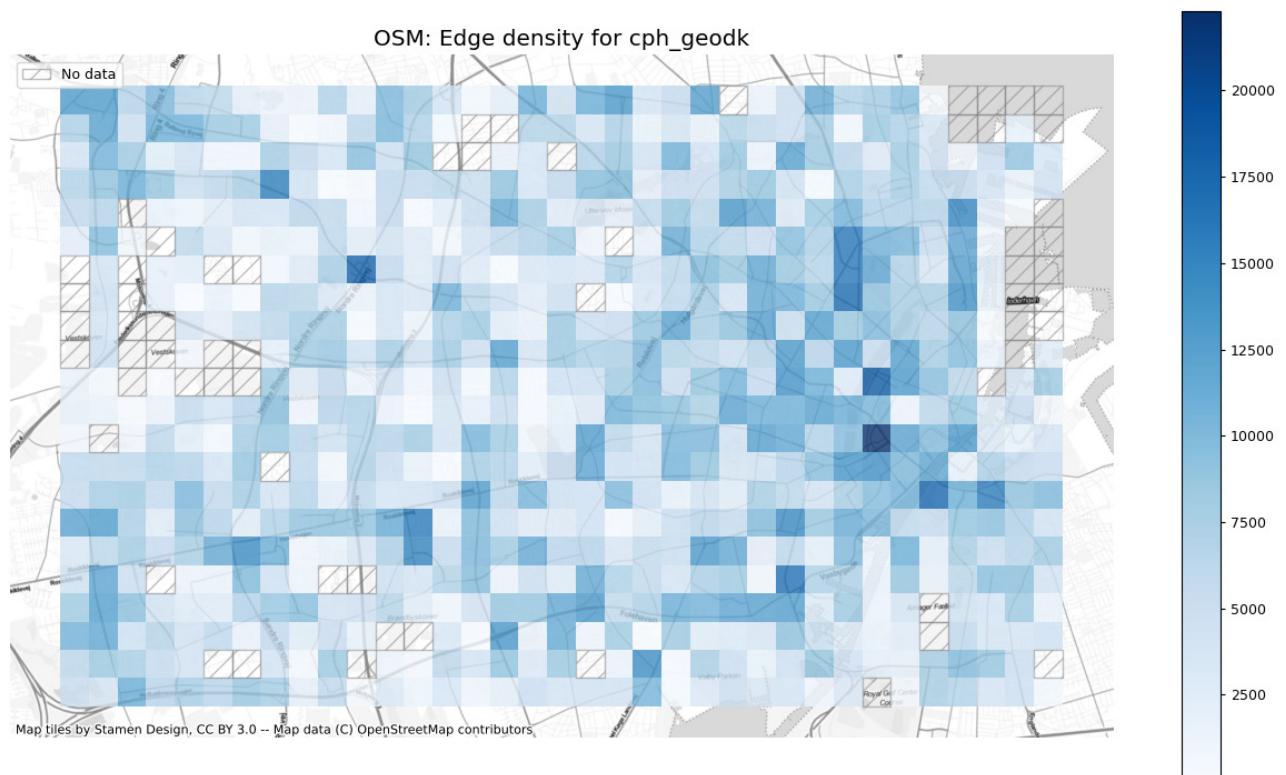
- lower-than-average infrastructure density indicates a locally sparser network
- higher-than-average node density indicates that there are relatively many intersections in the grid cell
- higher-than-average dangling node density indicates that there are relatively many dead ends in the grid cell

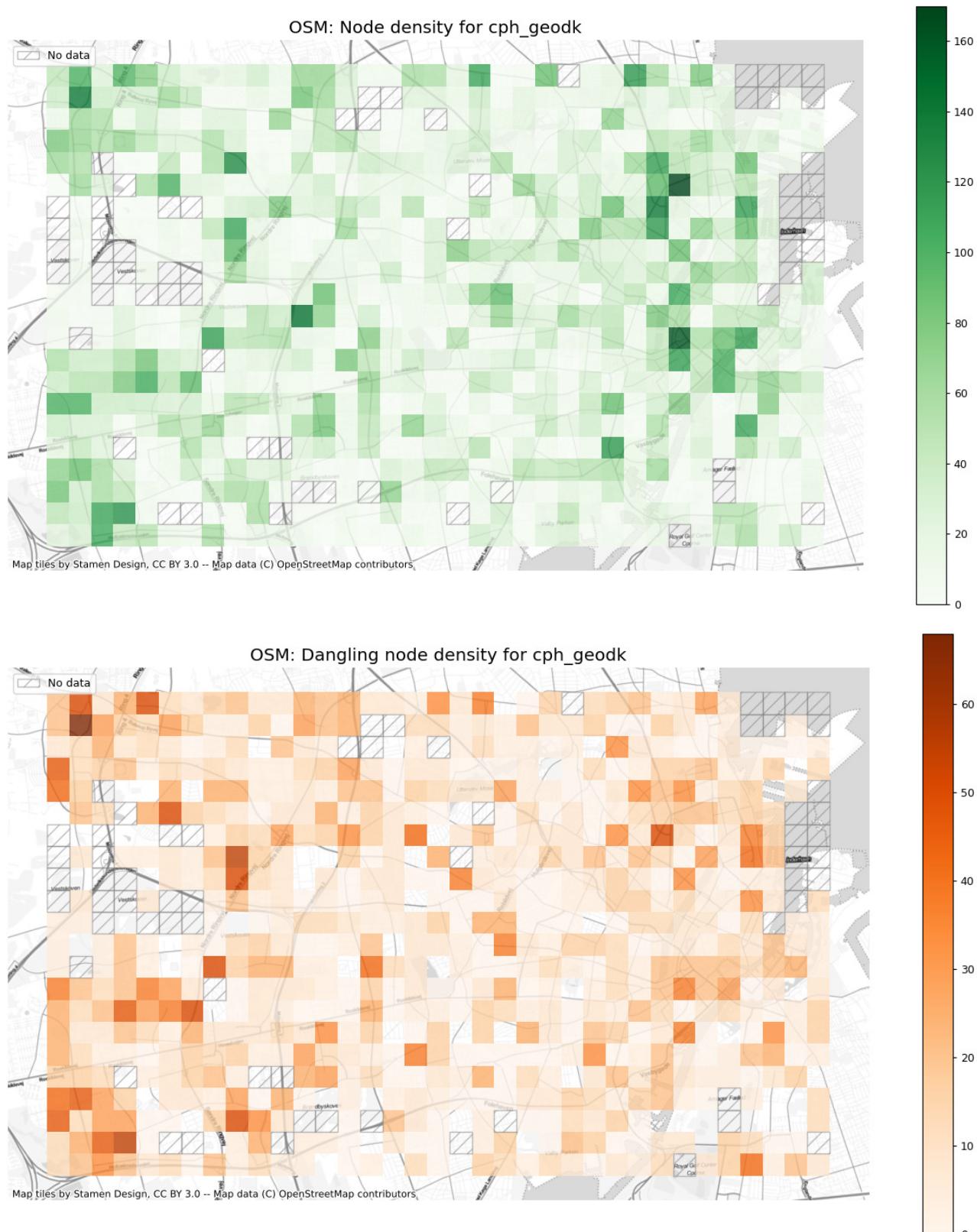
Global network density

For the entire study area, there are:

- 5862.22 meters of bicycle infrastructure per km².
- 27.16 nodes in the bicycle network per km².
- 10.03 dangling nodes in the bicycle network per km².
- 5300.30 meters of protected bicycle infrastructure per km².
- 502.71 meters of unprotected bicycle infrastructure per km².
- 59.21 meters of mixed protection bicycle infrastructure per km².

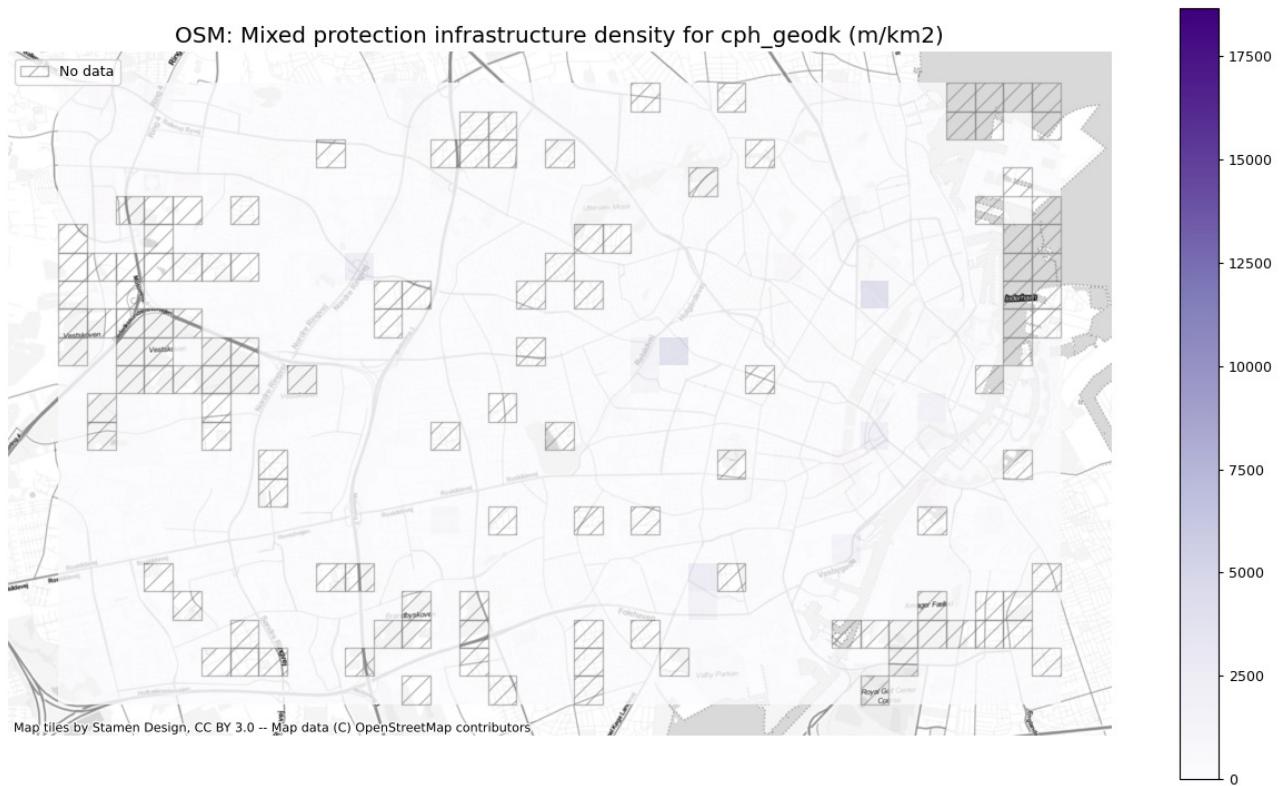
Local network density





Densities of protected and unprotected infrastructure





OSM tag analysis

For many practical and research purposes, more information than just the presence/absence of bicycle infrastructure is of interest. Information about e.g., the width of the infrastructure, speed limits, streetlights etc. can be of high relevance, for example when evaluating the bike friendliness of an area or individual network segment. The presence of these tags (describing attributes of the bicycle infrastructure) is however highly unevenly distributed in OSM, which poses a barrier to evaluations of bikeability and traffic stress. Likewise, the lack of restrictions on how OSM features can be tagged often result in conflicting tags, which undermines the evaluation of cycling conditions.

The section includes analyses of a. missing tags (edges with tags that lack information), b. incompatible tags (edges with tags labelled with two or more contradictory tags), and c. tagging patterns (the spatial variation of what tags are being used to describe bicycle infrastructure).

Note that for the evaluation of tags, the non-simplified edges should be used to avoid issues with tags that have been aggregated in the simplification process.

Missing tags

The information that is required or desirable to obtain from the OSM tags will depend on the use case - for example, the tag `lit` for a project looking at light conditions on cycle paths. The workflow below allows to quickly analyze the percentage of network edges that have a value for the tag of interest.

Method

We analyze all tags of interest as defined in the `existing_tag_analysis` section of `config.yml`. For each of these tags, `analyze_existing_tags` is used to compute the total number and the percentage of edges that have a corresponding tag value.

Interpretation

On study area level, a higher percentage of existing tag values in principle indicates a higher quality of the data set; however, note that this is different from an estimation of whether the existing tag values are truthful. On grid cell level, lower-than-average percentages for existing tag values can indicate a more poorly mapped area. However, note that the percentages are less informative for grid cells with a low number of edges: for example, if a cell contains one single edge that has a tag value for `lit`, the percentage of existing tag values is 100% - but given that there is only 1 data point this is less informative than say a value of 80% for a cell with many edges.

Global missing tags

Analysing tags describing:

surface – width – speedlimit – lit –

surface: 23339 out of 50964 edges (45.80%) have information.

surface: 552 out of 1286 km (42.95%) have information.

width: 5035 out of 50964 edges (9.88%) have information.

width: 97 out of 1286 km (7.58%) have information.

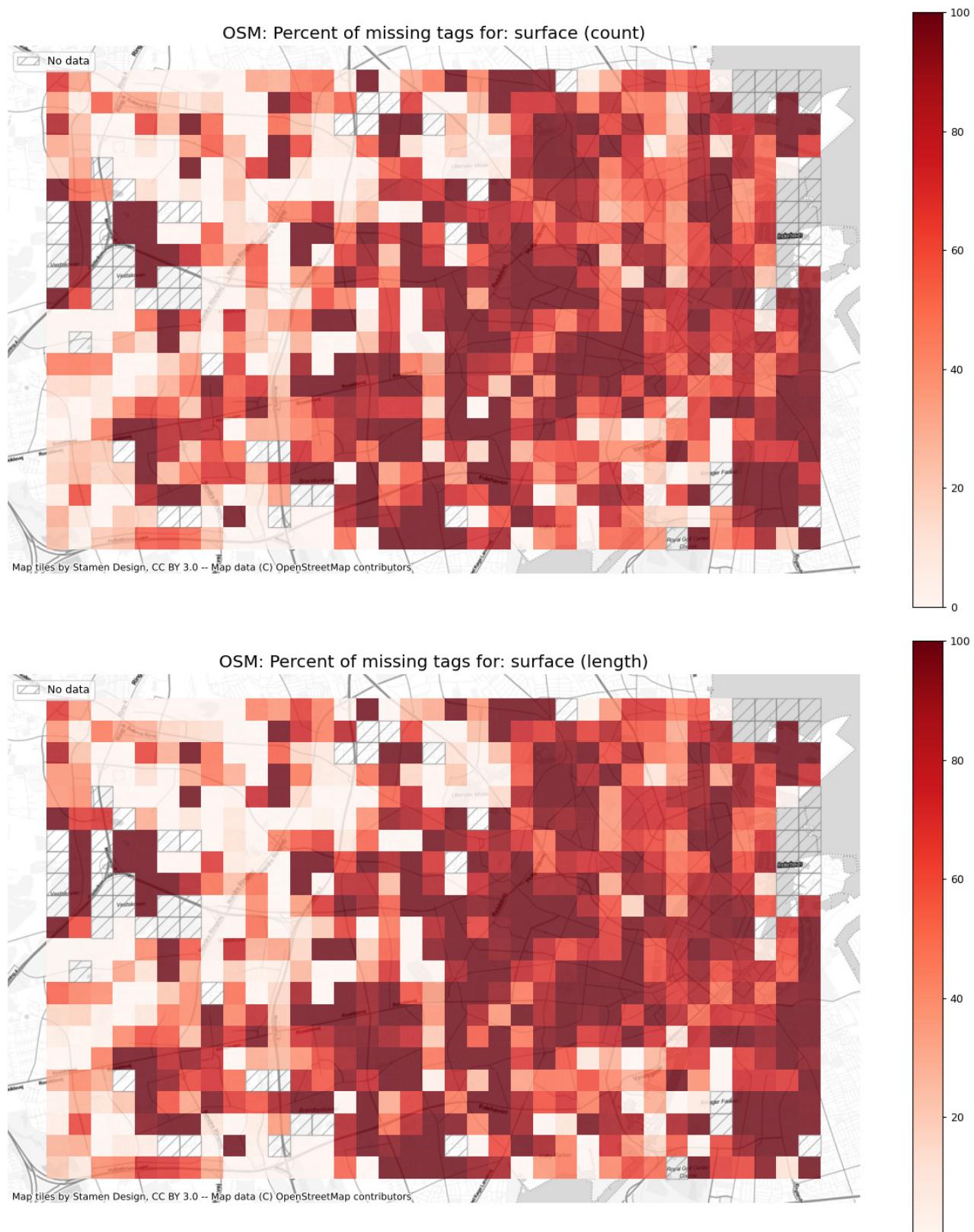
speedlimit: 25343 out of 50964 edges (49.73%) have information.

speedlimit: 684 out of 1286 km (53.22%) have information.

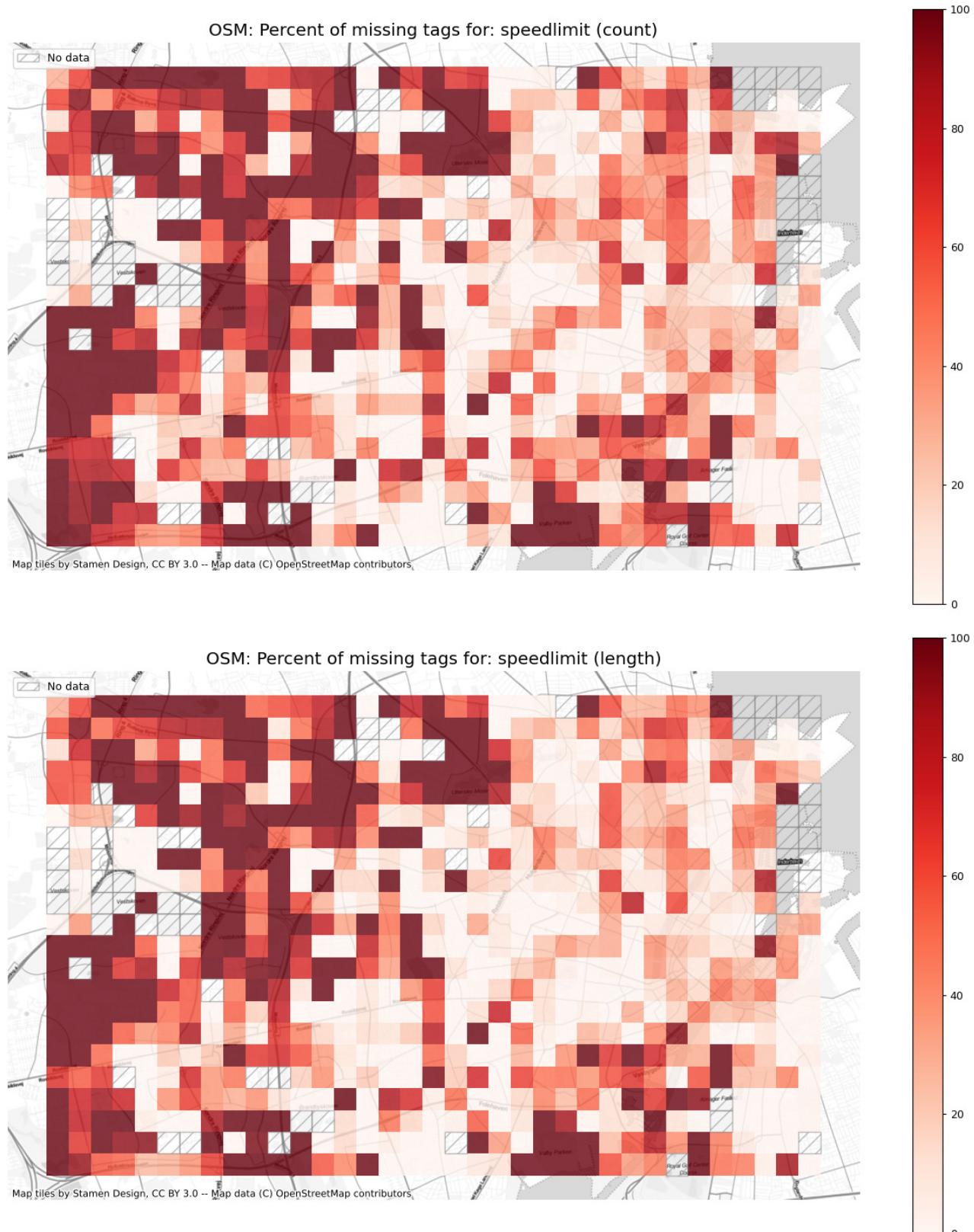
lit: 39313 out of 50964 edges (77.14%) have information.

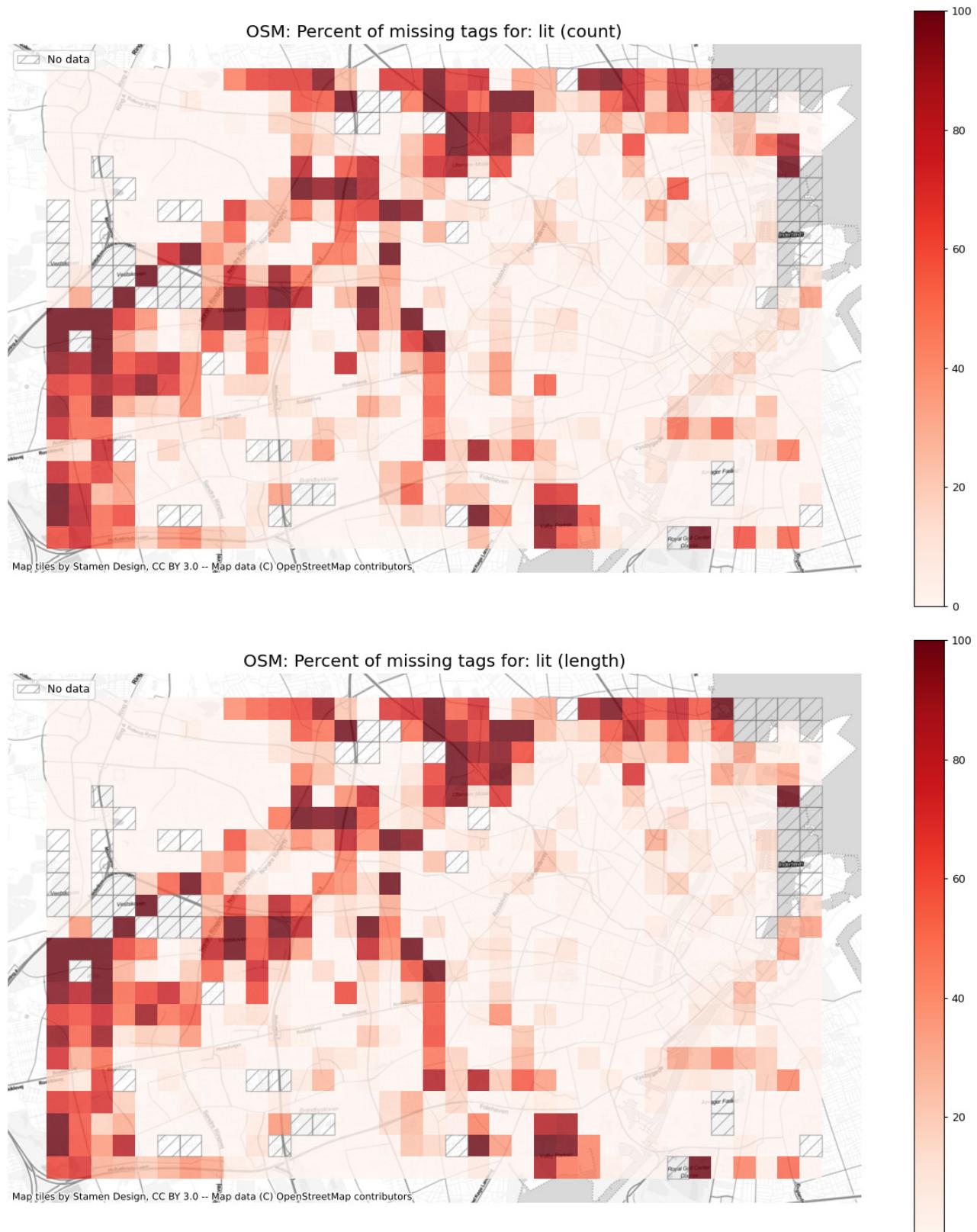
lit: 1008 out of 1286 km (78.42%) have information.

Local missing tags









Incompatible tags

Given that the tags in OSM data lack coherency at times and there are no restrictions in the tagging process (cf. [Barron et al., 2014](#)), incompatible tags might be present in the data set. For example, an edge might be tagged with the following two contradicting key-value pairs: `bicycle_infrastructure = yes` and `bicycle = no`.

Method

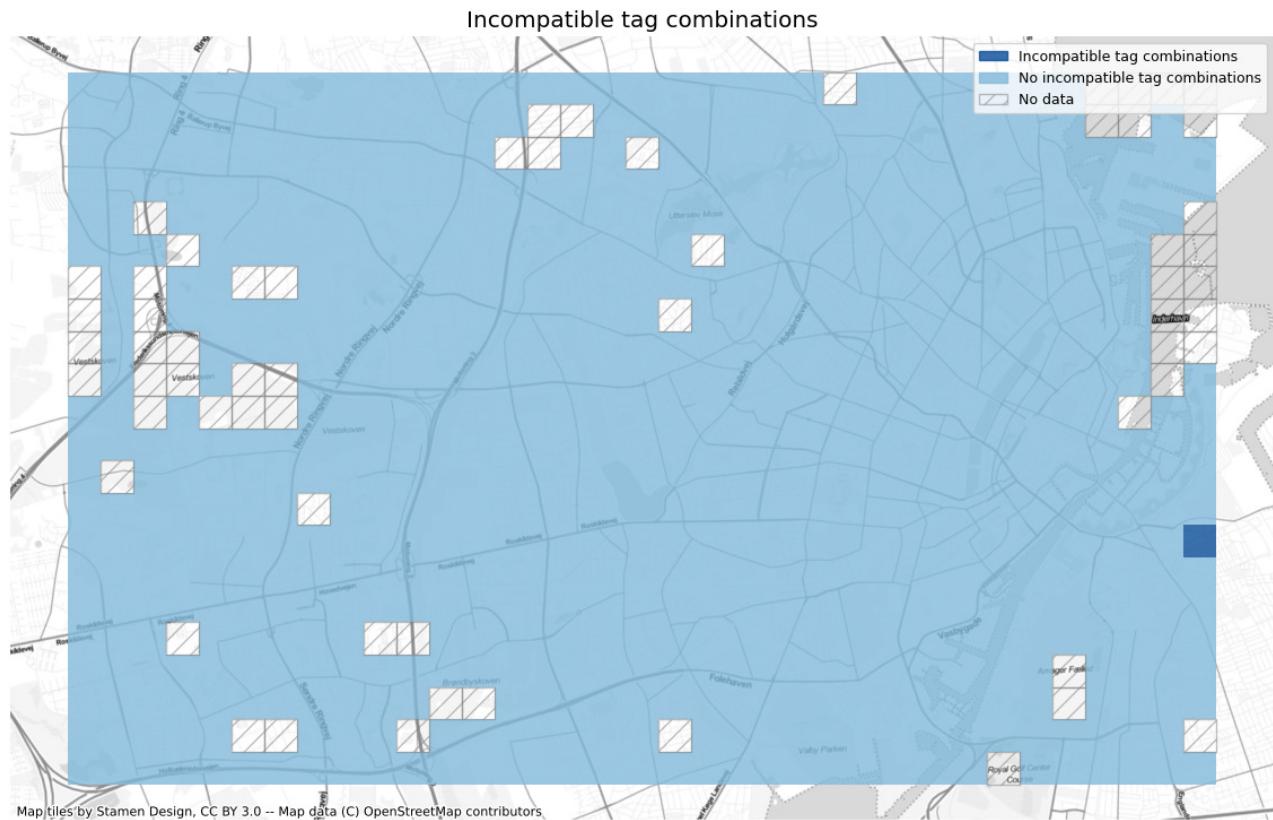
In the `config.yml` file, a list of incompatible key-value pairs for tags in the `incompatible_tags_analysis` is defined. Since there is no limitation as to which tags a data set could potentially contain, the list is, by definition, non-exhaustive, and can be adjusted by the user. In the section below, `check_incompatible_tags` is run, which identifies all incompatibility instances for a given area, first for study area level and then for grid cell level.

Interpretation

Incompatible tags are an undesired feature of the data set and render the corresponding data points invalid; there is no straightforward way to resolve the arising issues automatically, making it necessary to either correct the tag manually or to exclude the data point from the data set. A higher-than-average number of incompatible tags in a grid cell suggests local mapping issues.

In the entire dataset, there are 2 incompatible tag combinations (of those defined in the configuration file).

Local incompatible tags (per grid cell)



Plotting incompatible tag geometries



Tagging patterns

Identifying bicycle infrastructure in OSM can be tricky due to the many different ways in which the presence of bicycle infrastructure can be indicated. The [OSM Wiki](#) is a great resource for recommendations for how OSM features should be tagged, but some inconsistencies and local variations do remain. The analysis of tagging patterns allows to visually explore some of the potential inconsistencies.

Regardless of how the bicycle infrastructure is defined, examining which tags contribute to which parts of the bicycle network allows to visually examine patterns in tagging methods. It also allows to estimate whether some elements of the query will lead to the inclusion of too many or too few features.

Likewise, 'double tagging' where several different tags have been used to indicate bicycle infrastructure can lead to misclassifications of the data. For this reason, identifying features that are included in more than one of the queries defining bicycle infrastructure can indicate issues with the tagging quality.

Method

The section below first plots individual subsets of the OSM data set for each of the queries listed in `bicycle_infrastructure_queries`, as defined in the `config.yml` file. The subset defined by a query is the set of edges for which this query is True. Since several queries can be True for the same edge, the subsets can overlap. In the second step below, all overlaps between 2 or more queries are plotted (i.e. all edges that have been assigned several, potentially competing, tags).

Interpretation

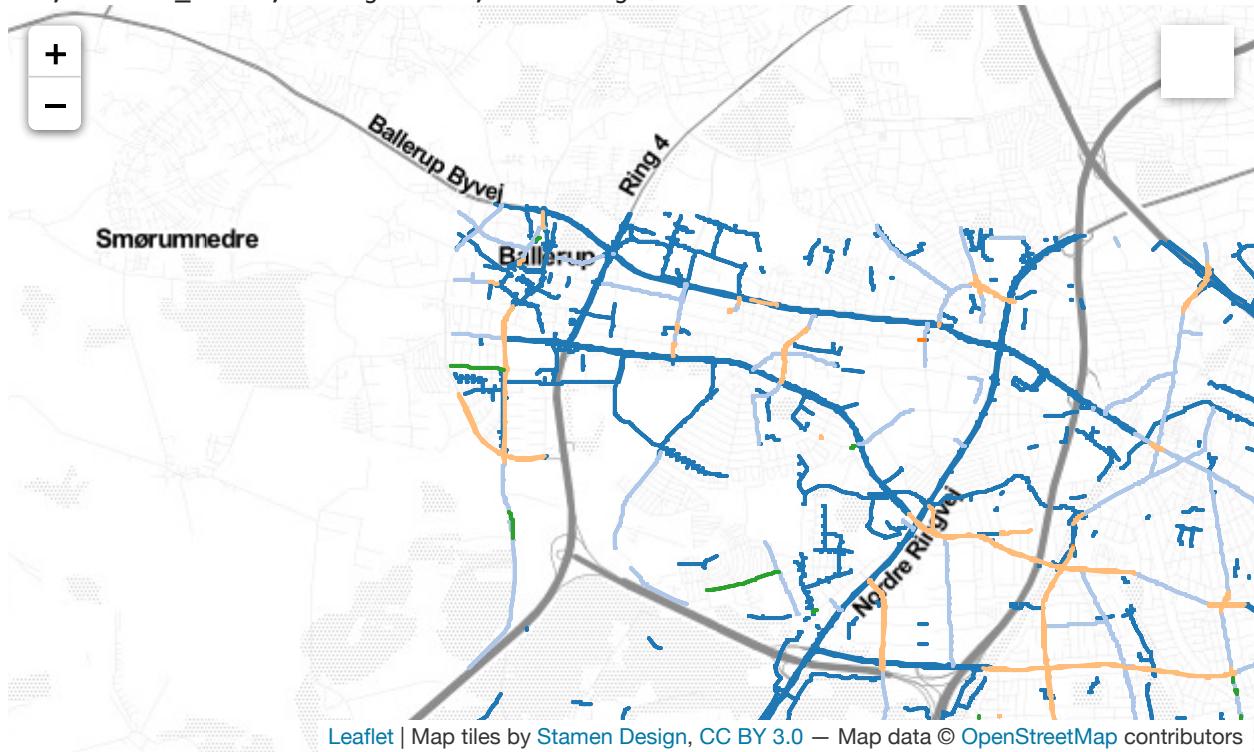
The plots for each tagging type allow for a quick visual overview of different tagging patterns present in the area. Based on local knowledge, the user may estimate whether the differences in tagging types are due to actual physical differences in the infrastructure or rather an artefact of the OSM data. Next, the user can access overlaps between different tags; depending on the specific tags, this may or may not be a data quality issue. For example, in case of '`'cycleway:right'`' and '`'cycleway:left'`', having data for both tags is valid, but other combinations such as '`'cycleway='track'`' and '`'cycleway:left=lane'`' gives an ambiguous picture of what type of bicycle infrastructure is present.

Tagging types

```
Tagging type A: highway == 'cycleway'  
Tagging type B: cycleway in ['lane', 'track', 'opposite_lane', 'opposite_track', 'shared_lane', 'designated', 'crossing']  
Tagging type C: cycleway_left in ['lane', 'track', 'opposite_lane', 'opposite_track', 'shared_lane', 'designated', 'crossing']
```

Tagging type D: cycleway_right in ['lane', 'track', 'opposite_lane', 'opposite_track', 'shared_lane', 'designated', 'crossing']

Tagging type E: cycleway_both in ['lane', 'track', 'opposite_lane', 'opposite_track', 'shared_lane', 'designated', 'crossing']

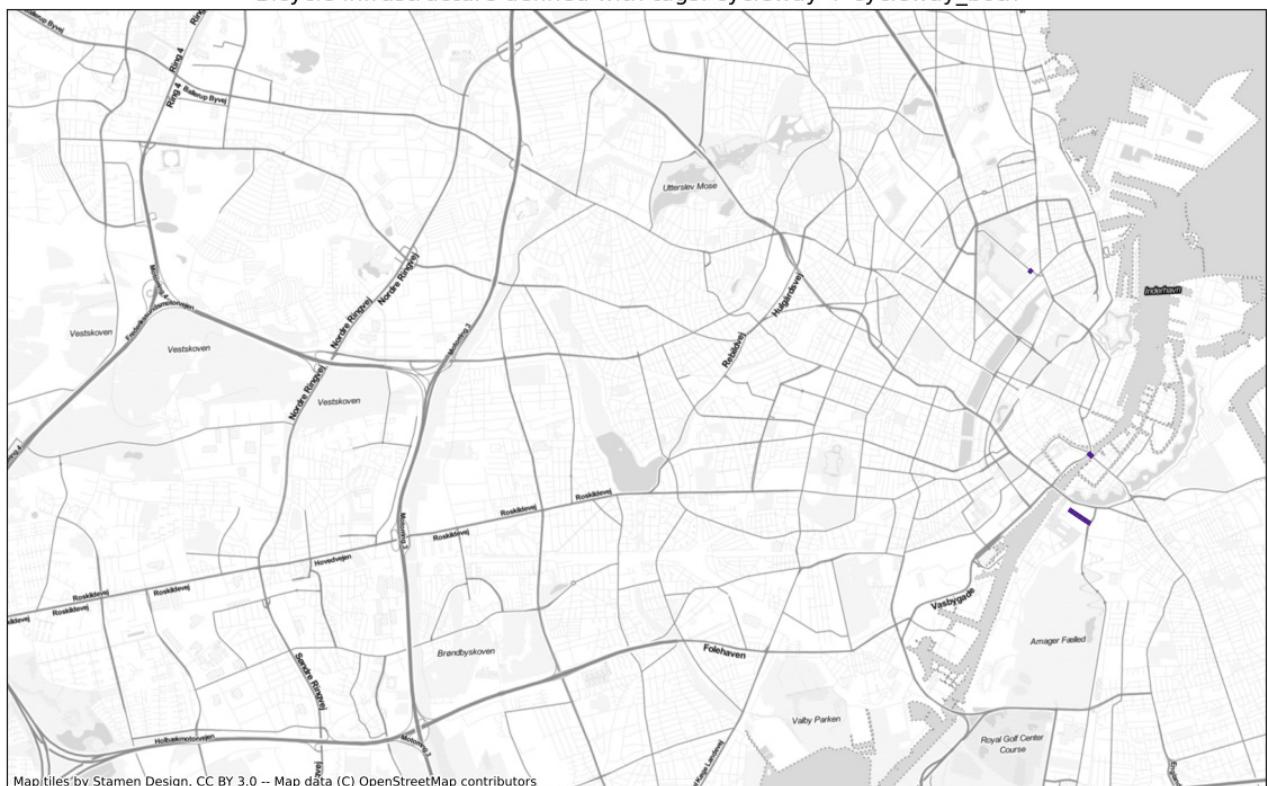


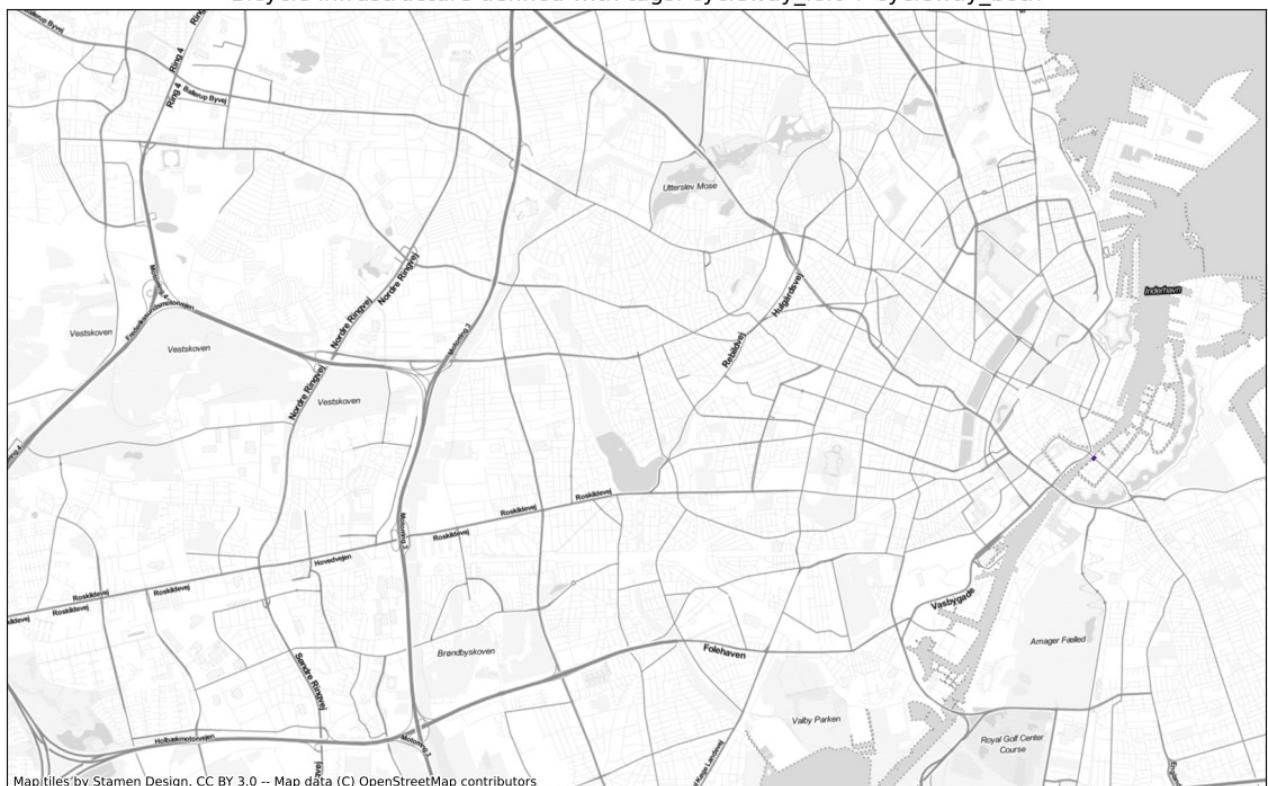
Multiple tagging

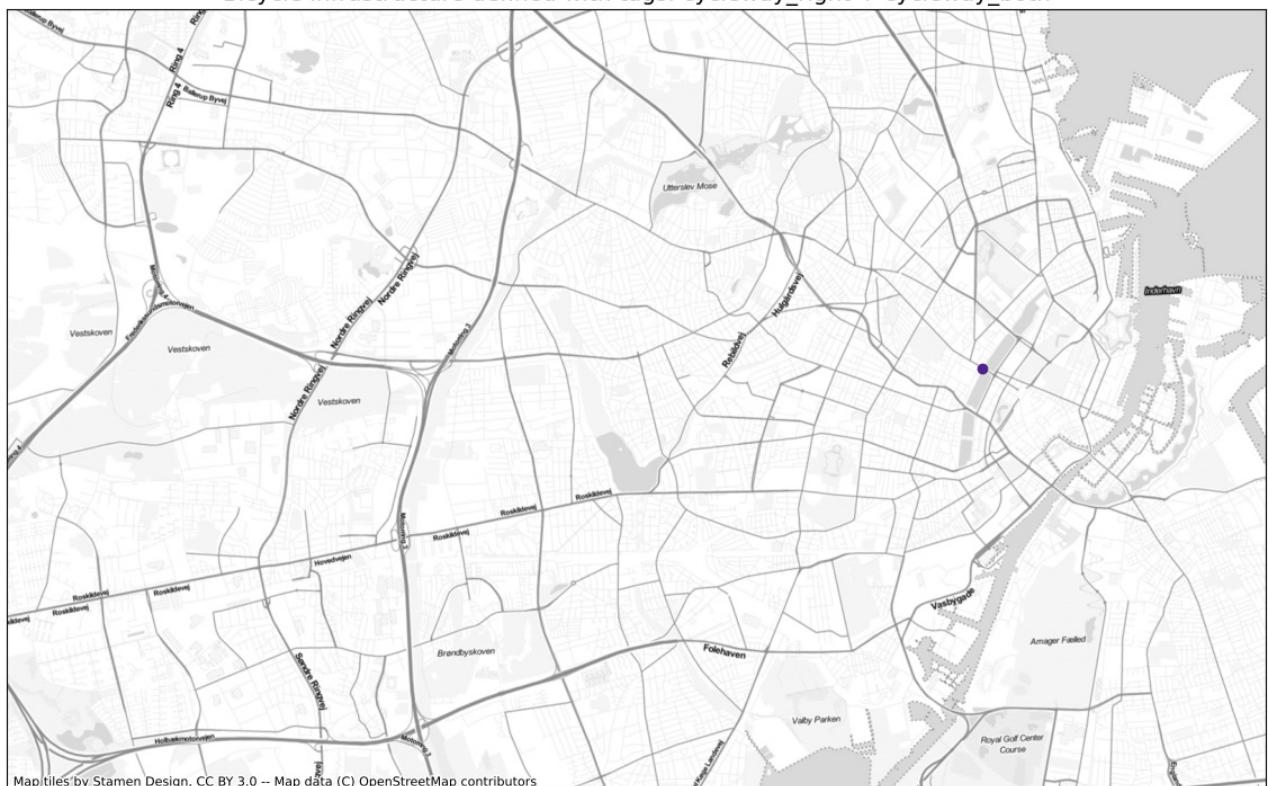
Bicycle infrastructure defined with tags: `cycleway_left + cycleway_right`



Bicycle infrastructure defined with tags: `cycleway + cycleway_both`



Bicycle infrastructure defined with tags: highway + cycleway**Bicycle infrastructure defined with tags: cycleway_left + cycleway_both**

Bicycle infrastructure defined with tags: `cycleway + cycleway_right`Bicycle infrastructure defined with tags: `cycleway_right + cycleway_both`

Network topology

This section explores the geometric and topological features of the data.

These are, for example, network density, disconnected components, dangling (degree one) nodes; it also includes exploring whether there are nodes that are very close to each other but do not share an edge - a potential sign of edge undershoots - or if there are intersecting edges without a node at the intersection, which might indicate a digitizing error that will distort routing attempts on the network.

Due to the fragmented nature of most networks of bicycle infrastructure, many metrics, such as missing links or network gaps, simply reflect the true extent of the infrastructure ([Natera Orozco et al., 2020](#)). This is different for car networks, where e.g., disconnected components could more readily be interpreted as a data quality issue.

Therefore, the analysis only takes very small network gaps into account as potential data quality issues.

Subsections:

- [Simplification outcome](#)
- [Dangling nodes](#)
- [Under/Overshoots](#)
- [Missing intersection nodes](#)

Simplification outcome

When converting a set of geocoded linestrings (polygonal chains) to graph format, not all vertices (nodes) are of equal meaning. For geometry of the infrastructural element, all nodes are needed as an ordered list. For the topology of the network, however, only those nodes that are endpoints or intersection points with other edges are needed, while all other (so-called 'interstitial') nodes do not add any information. To compare the structure and true ratio between nodes and edges in a network, a simplified network representation which only includes nodes at endpoints and intersections, or where the value of important attributes changes, is required. Therefore, in notebook 1 the bicycle network was simplified by removing all interstitial nodes from the graph object (retaining, however, the complete node lists in the geometry attribute of each edge). An additional advantage of simplifying the network is the resulting substantial reduction of the number of nodes and edges, which makes computational routines much faster.

Comparing the degree distribution for the networks before and after simplification is a quick sanity check for the simplification routine. Typically, the big majority of nodes in the non-simplified network will be of degree two; in the simplified network, however, most nodes will have degrees other than two. Degree two nodes are retained in only two cases: if they represent a connection point between two different types of infrastructure; or if they are needed in order to avoid self-loops (edges whose start and end points are identical) or multiple edges between the same pair of nodes.

As part of the simplification routine, in cases where there are several edges between the same pair of nodes ('parallel edges' or 'multiedges'), only one of the edges is retained. Within the routine, the number edges removed in this way are counted.

Method

The degree distributions before and after simplification are plotted below.

Interpretation

Typically, the degree distribution will go from high (before simplification) to low (after simplification) counts of degree two nodes, while it will not change for all other degrees (1, or 3 and higher). Further, the total number of nodes will see a strong decline. If the simplified graph still maintains a relatively high number of degree two nodes, or if the number of nodes with other degrees changes after the simplification, this might point to issues either with the graph conversion or with the simplification process.

Simplifying the network decreased the number of edges by 88.9% and the number of nodes by 84.3%.

Before the network simplification the OSM graph had:

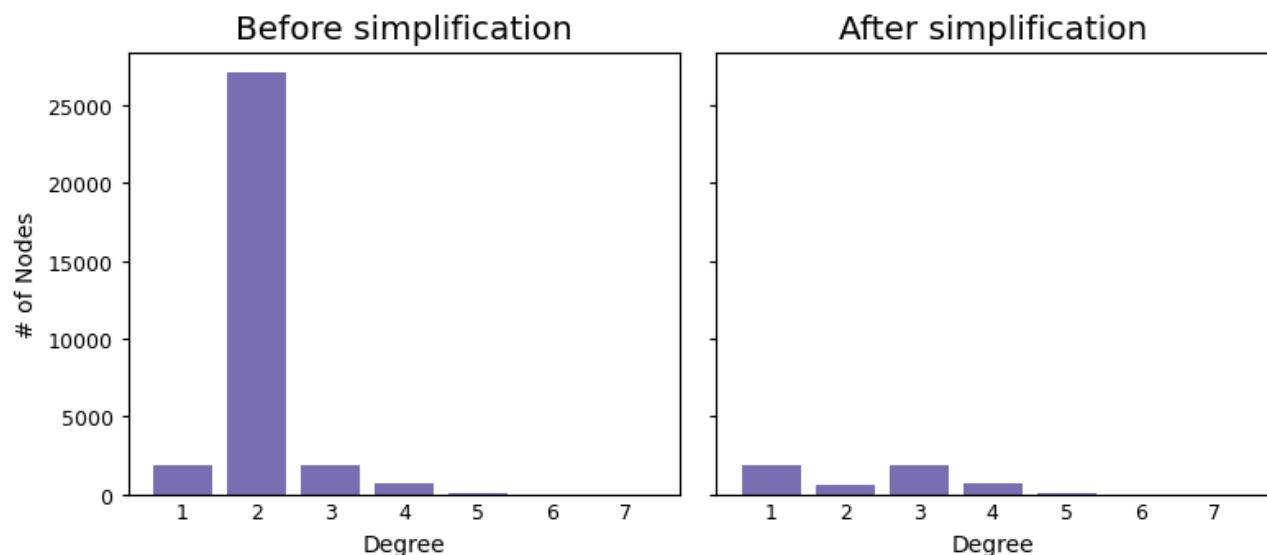
- 1 node(s) with degree 7
- 3 node(s) with degree 6
- 29 node(s) with degree 5

- 702 node(s) with degree 4
- 1807 node(s) with degree 3
- 27073 node(s) with degree 2
- 1819 node(s) with degree 1

After the network simplification the OSM graph had:

- 1 node(s) with degree 7
- 3 node(s) with degree 6
- 29 node(s) with degree 5
- 702 node(s) with degree 4
- 1807 node(s) with degree 3
- 566 node(s) with degree 2
- 1819 node(s) with degree 1

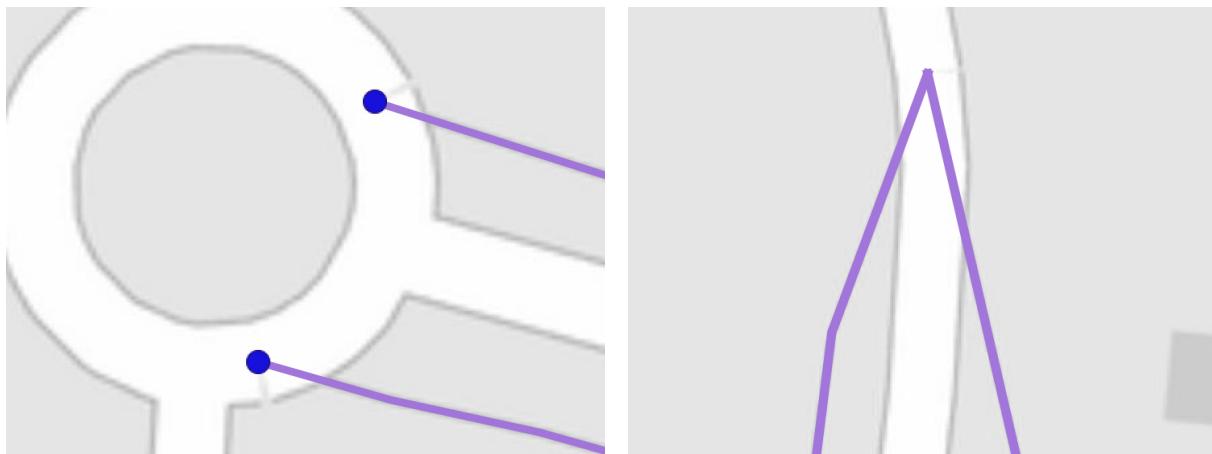
OSM: Degree distributions for the network in cph_geodk



Dangling nodes

Dangling nodes are nodes of degree one - in other words, nodes that have only one single edge attached to them. Most networks will naturally contain a number of dangling nodes. Dangling nodes can occur at actual dead-ends (representing a cul-de-sac) or at the endpoints of certain features (e.g., when a bicycle path ends in the middle of a street). However, dangling nodes can also occur as a data quality issue in case of over/undershoots (as described in detail in the next section). The number of dangling nodes in a network does to some extent also depend on the digitization method, as shown in the illustration below.

Therefore, the presence of dangling nodes is in itself not a sign of low data quality. However, a high number of dangling nodes in an area that is not known for suffering from many dead-ends can indicate digitization errors and problems with edge over/undershoots.



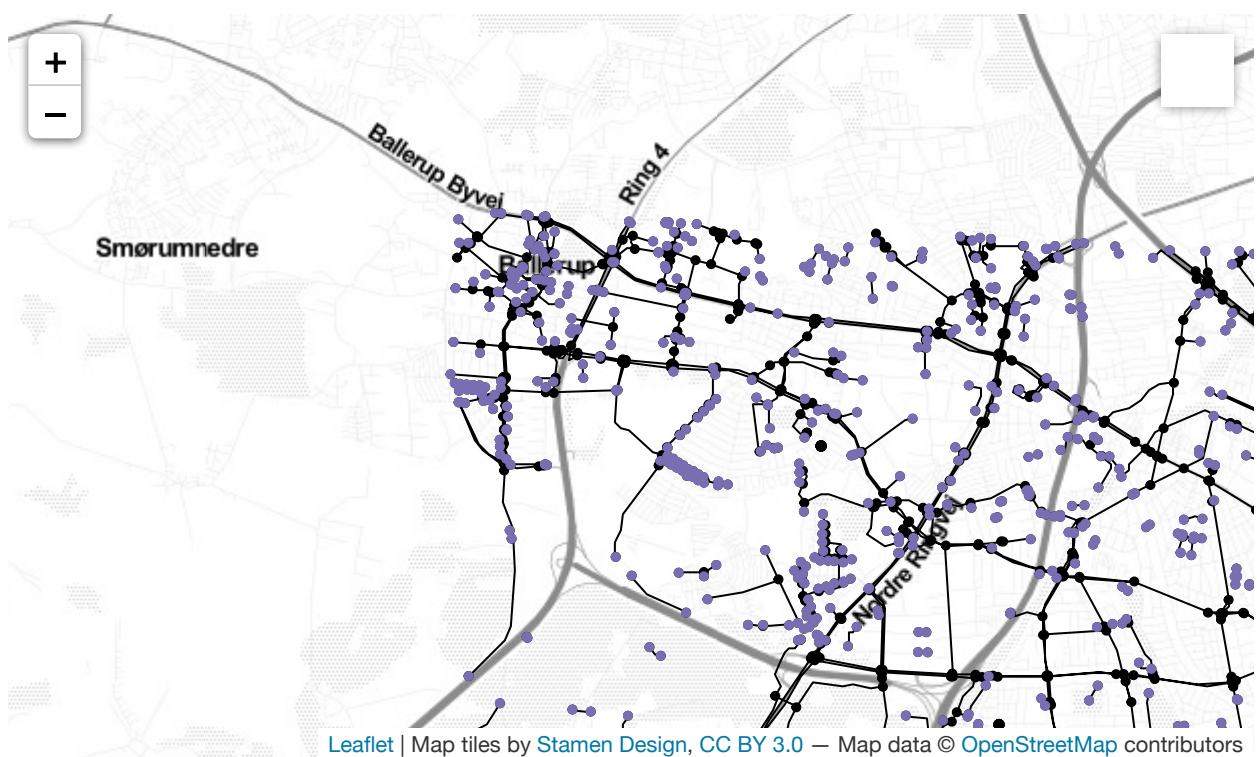
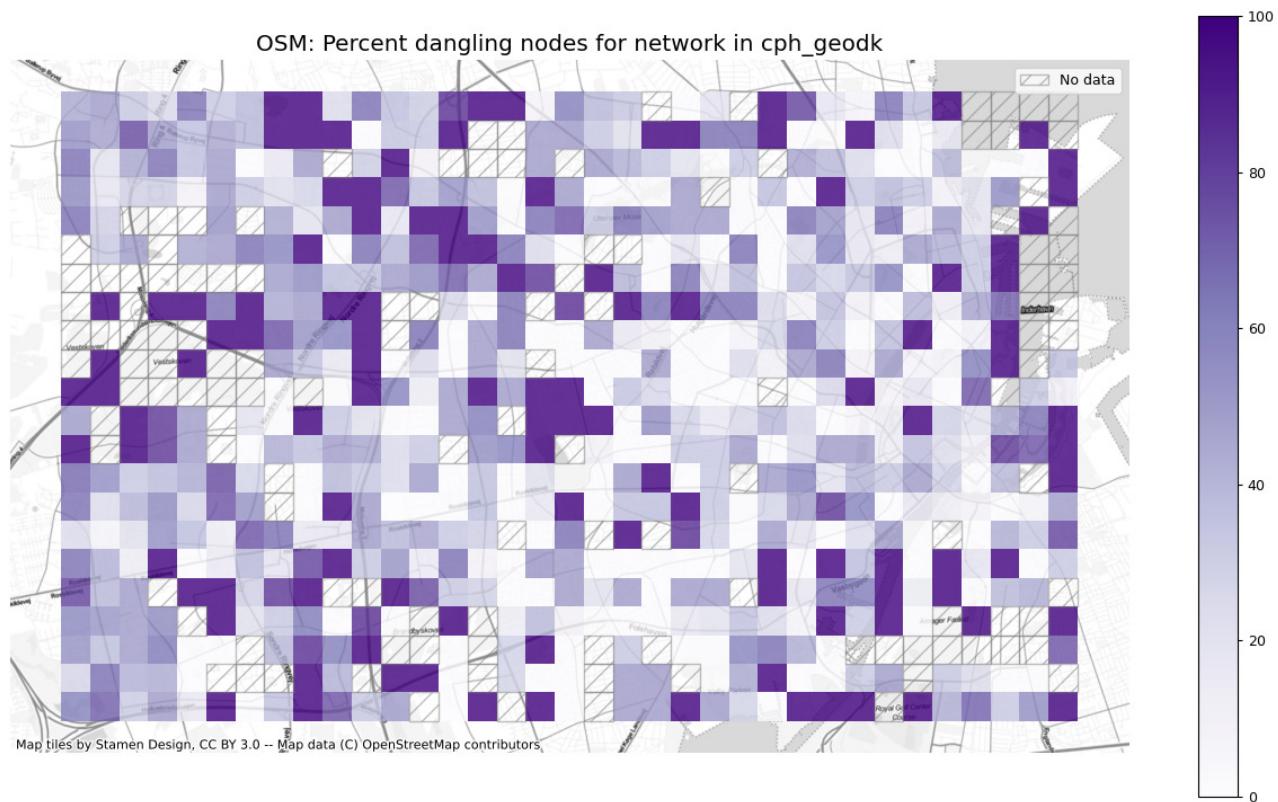
Dangling nodes occur where road features end (left), but when separate features are joined at the end (right), there will be no dangling nodes

Method

Below, a list of all dangling nodes is obtained with the help of `get_dangling_nodes`. Then, the network with all its nodes is plotted. The dangling nodes are shown in purple; all other nodes are shown in black.

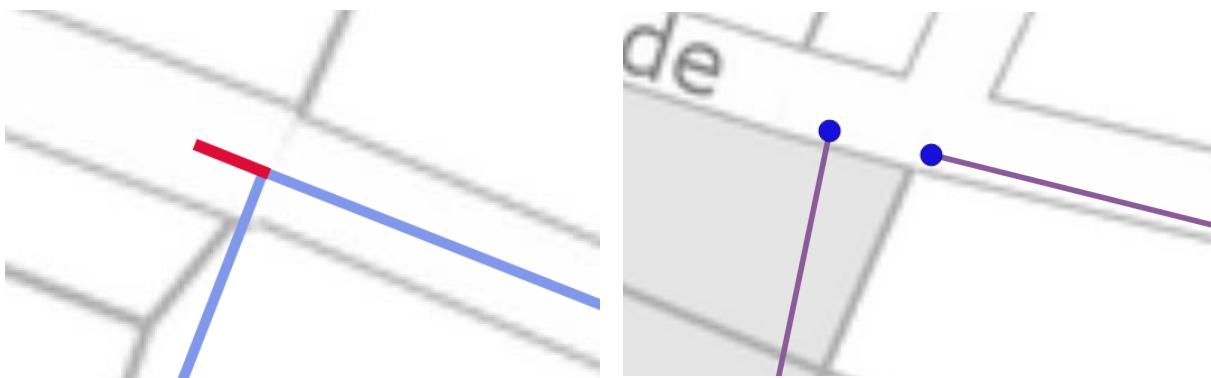
Interpretation

We recommend a visual analysis in order to interpret the spatial distribution of dangling nodes, with particular attention to areas of high dangling node density. It is important to understand where dangling nodes come from: are they actual dead-ends or digitization errors (e.g., over/undershoots)? A higher number of digitization errors points to a lower quality of the data.



Under/overshoots

When two nodes in a simplified network are placed within a distance of a few meters, but do not share a common edge, it is often due to an edge over/undershoot or another digitizing error. An overshoot occurs when two features meet and one of them extends beyond the other. An undershoot occurs when two features are supposed to meet, but instead are just in close proximity to each other. See the image below for an illustration of an overshoot (left) and an undershoot (right). For a more detailed explanation of over/undershoots, see the [GIS Lounge website](#).



Overshoots refer to situations where a line feature extends too far beyond at intersecting line, rather than ending at the intersection (left). Undershoots happen when two line features are not properly joined, for example at intersection (right)

Method

Overshoots: First, the `length_tolerance` (in meters) is defined in the cell below. Then, with `find_overshoots`, all network edges that have a dangling node attached to them and that have a maximum length of `length_tolerance` are identified as overshoots, and the results are plotted.

Undershoots: First, the `length_tolerance` (in meters) is defined in the cell below. Then, with `find_undershoots`, all pairs of dangling nodes that have a maximum of `length_tolerance` distance between them, are identified as undershoots, and the results are plotted.

The workflow for over/undershoot detection below is inspired by [Neis et al. 2012](#).

Interpretation

Note that over/undershoots are not necessarily always a data quality issue - they might be instead an accurate representation of the network conditions or of the digitization strategy (for example, a cycle path might end abruptly soon after a turn, which results in an overshoot; protected cycle paths are often digitized in OSM as interrupted at intersections, which results in intersection 'undershoots').

The interpretation of the impact of over/undershoots on data quality is context dependent. For certain applications, such as routing, overshoots do not present a particular challenge; they can, however, pose an issue for other applications such as network analysis, given that they skew the network structure. Undershoots, on the contrary, are a serious problem for routing applications, especially if only bicycle infrastructure is considered; they also pose a problem for network analysis if they result in disconnected components.

9 potential overshoots were identified using a length tolerance of 3 meters.
 14 potential undershoots were identified using a length tolerance of 3 meters.



Missing intersection nodes

When two edges intersect without having a node at the intersection – and if neither edges are tagged as a bridge or a tunnel – there is a clear indication of a topology error.

Method

The workflow below is inspired by [Neis et al. 2012](#). First, with the help of `check_intersection`, each edge which is not tagged as either tunnel or bridge is checked for any crossing with another edge of the network. If this is the case, the edge is marked as having an intersection issue. The number of intersection issues found is printed and the results are plotted for visual analysis.

Interpretation

A higher number of intersection issues points to a lower data quality. However, it is recommended with a manual visual check of all intersection issues with a certain knowledge of 1 place(s) appear to be missing an intersection node or a bridge/tunnel tag.



Network components

Disconnected components

Disconnected components do not share any elements (nodes/edges). In other words, there is no network path that could lead from one disconnected component to the other. As mentioned above, most real-world networks of bicycle infrastructure do consist of many disconnected components ([Natera Orozco et al., 2020](#)). However, when two disconnected components are very close to each other, it might be a sign of a missing edge or another digitizing error.

Method

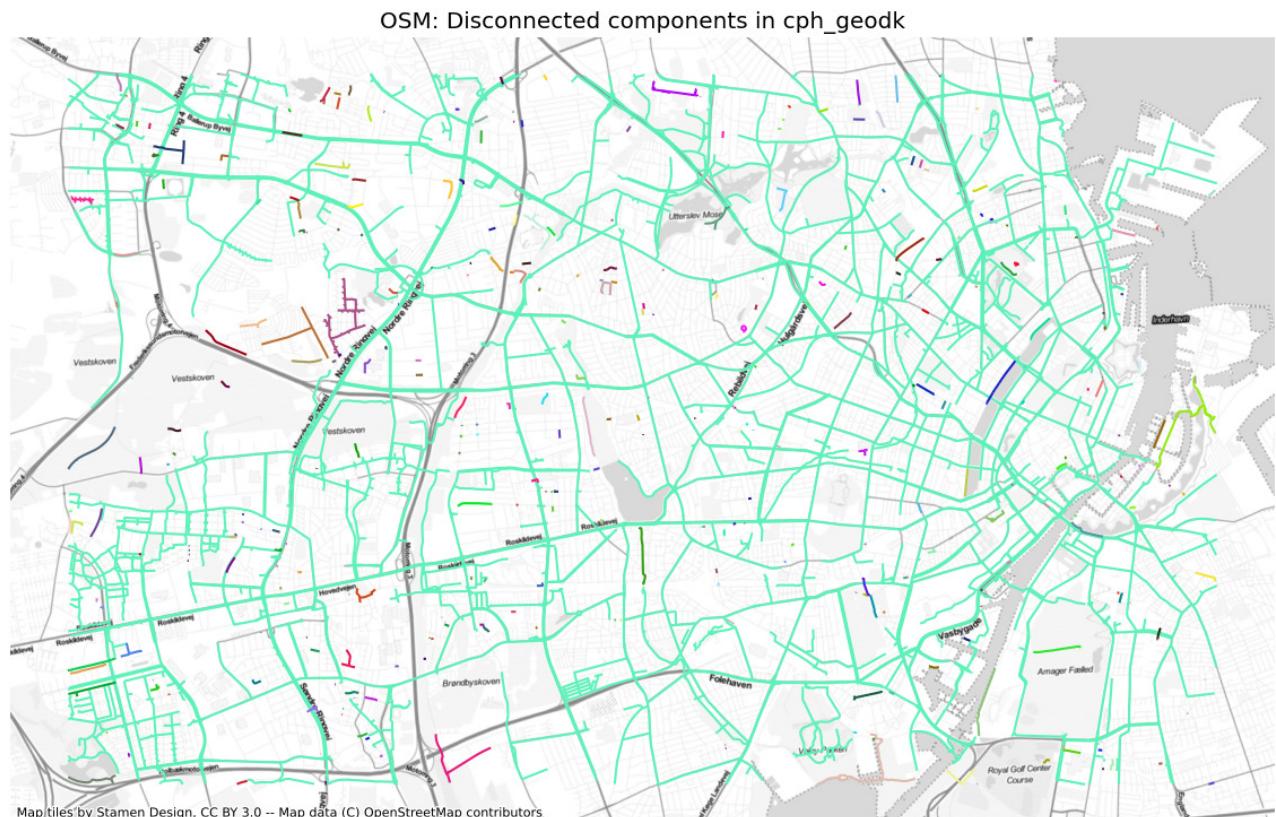
First, with the help of `return_components`, a list of all (disconnected) components of the network is obtained. The total number of components is printed and all components are plotted in different colors for visual analysis. Next, the component size distribution (with components ordered by the network length they contain) is plotted, followed by a plot of the largest connected component.

Interpretation

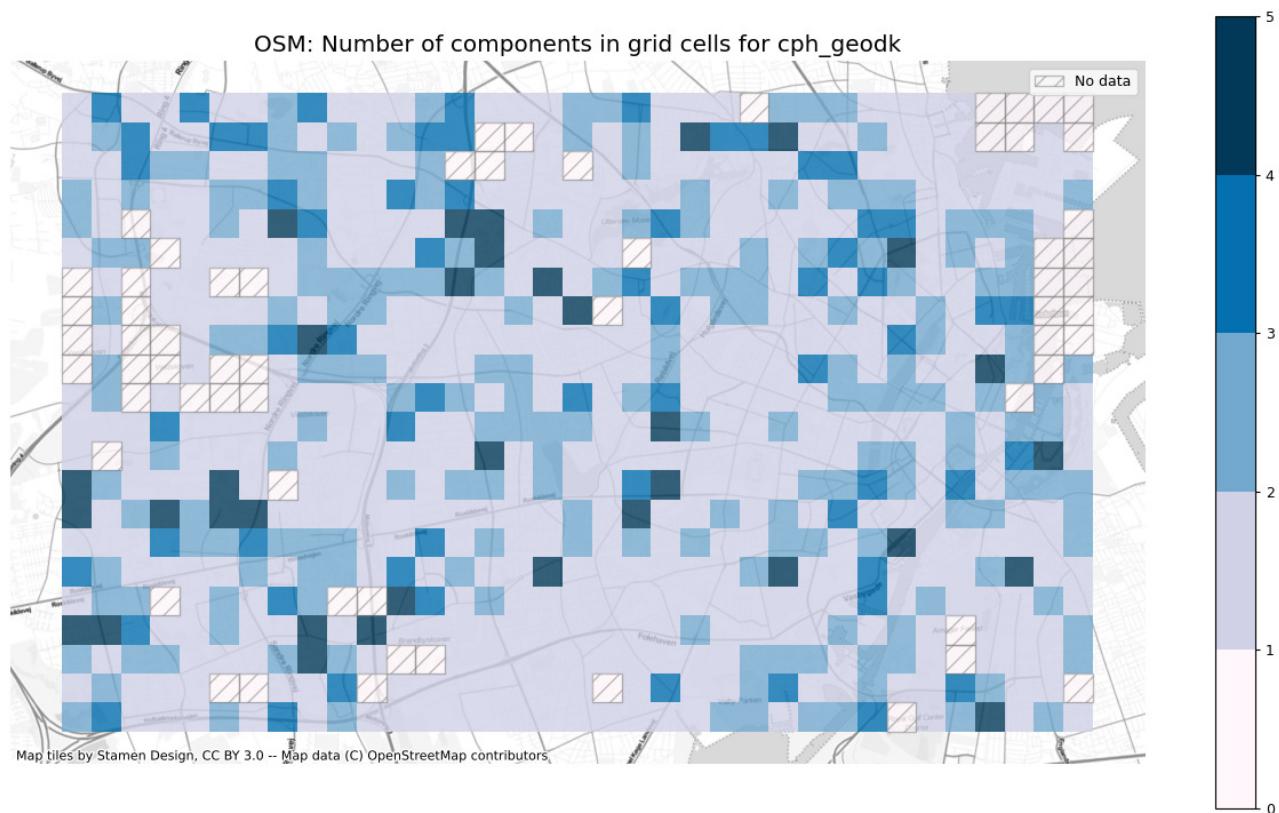
As with many of the previous analysis steps, knowledge of the area is crucial for a correct interpretation of component analysis. Given that the data represents the actual infrastructure accurately, bigger components indicate coherent network parts, while smaller components indicate scattered infrastructure (e.g., one single bicycle path along a street that does not connect to any other bicycle infrastructure). A high number of disconnected components in near vicinity of each other could indicate digitization errors or missing data.

Number of components

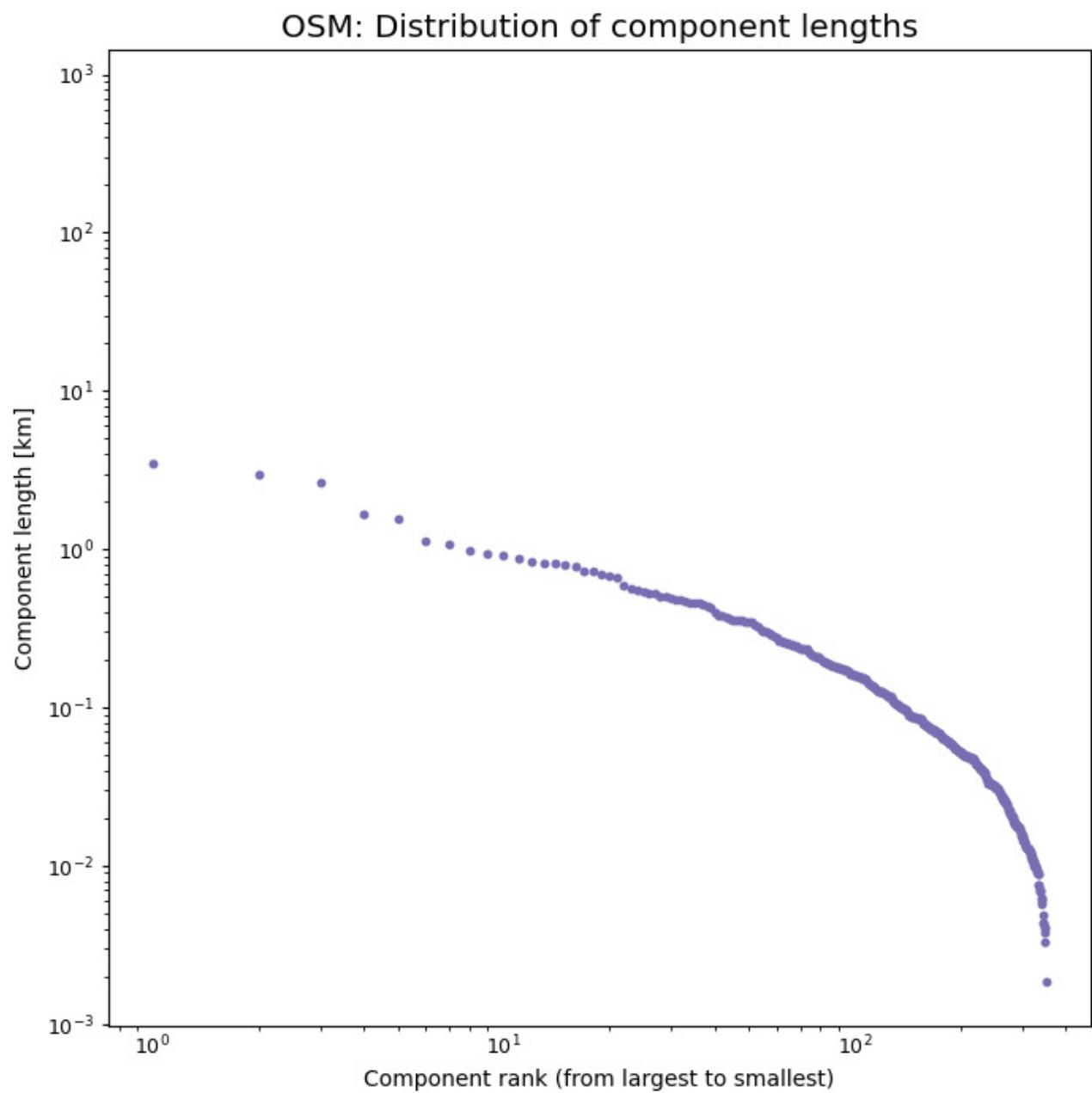
The network in the study area consists of 352 disconnected components.

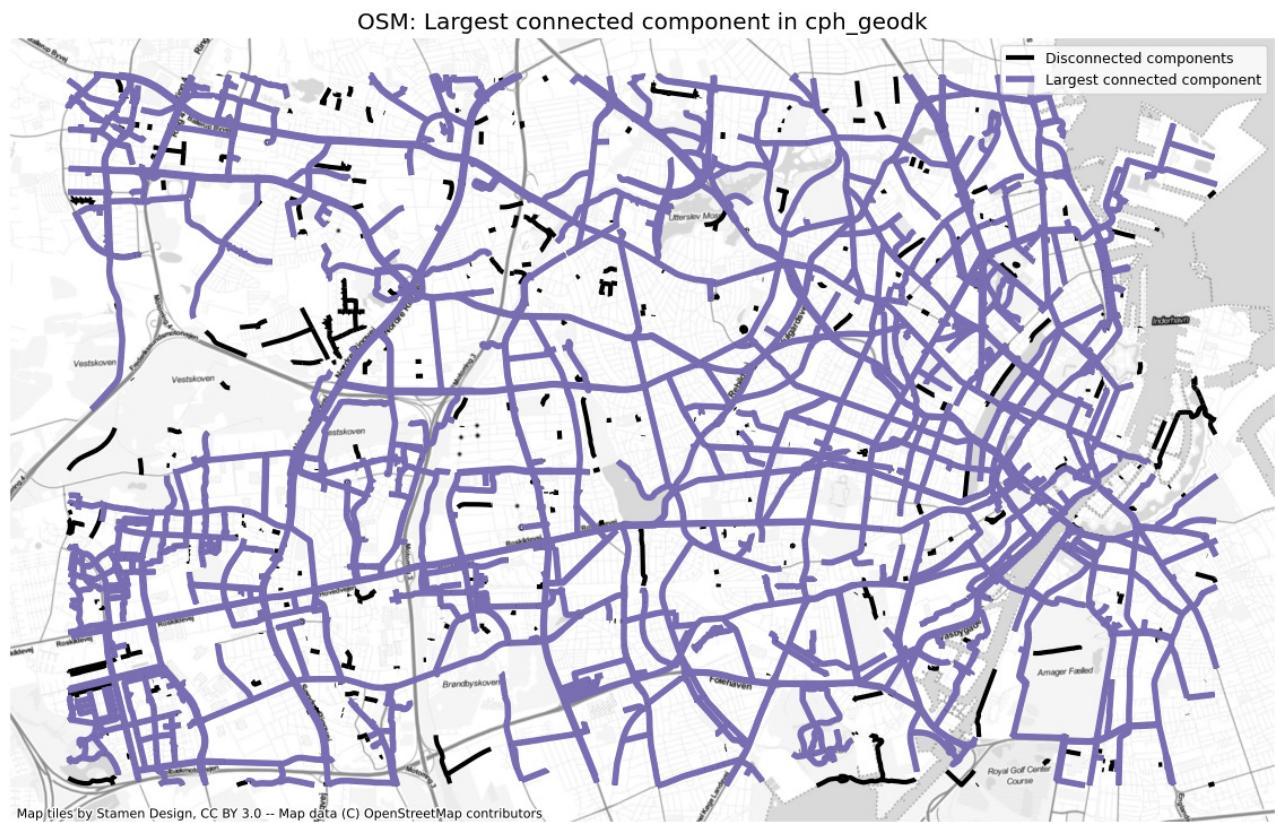


Number of components per grid cell



Distribution of network length per component

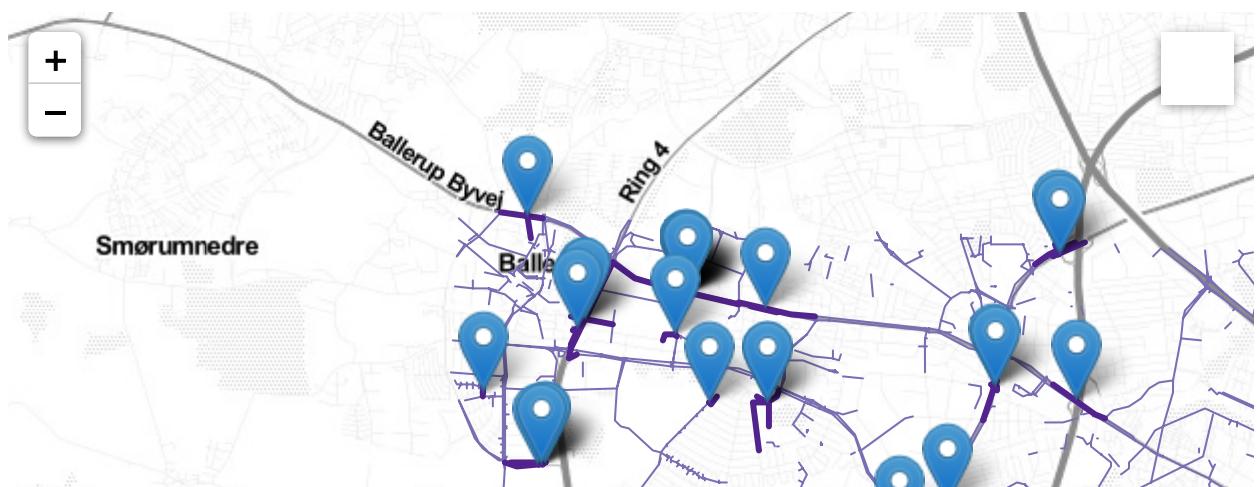




Potential missing links between disconnected components

In the plot of potential missing links between components, all edges that are within the specified distance of an edge on another component are plotted. The gaps between disconnected edges are highlighted with a marker. The map thus highlights edges which, despite being in close proximity of each other, are disconnected and where it thus would not be possible to bike on cycling infrastructure between the edges.

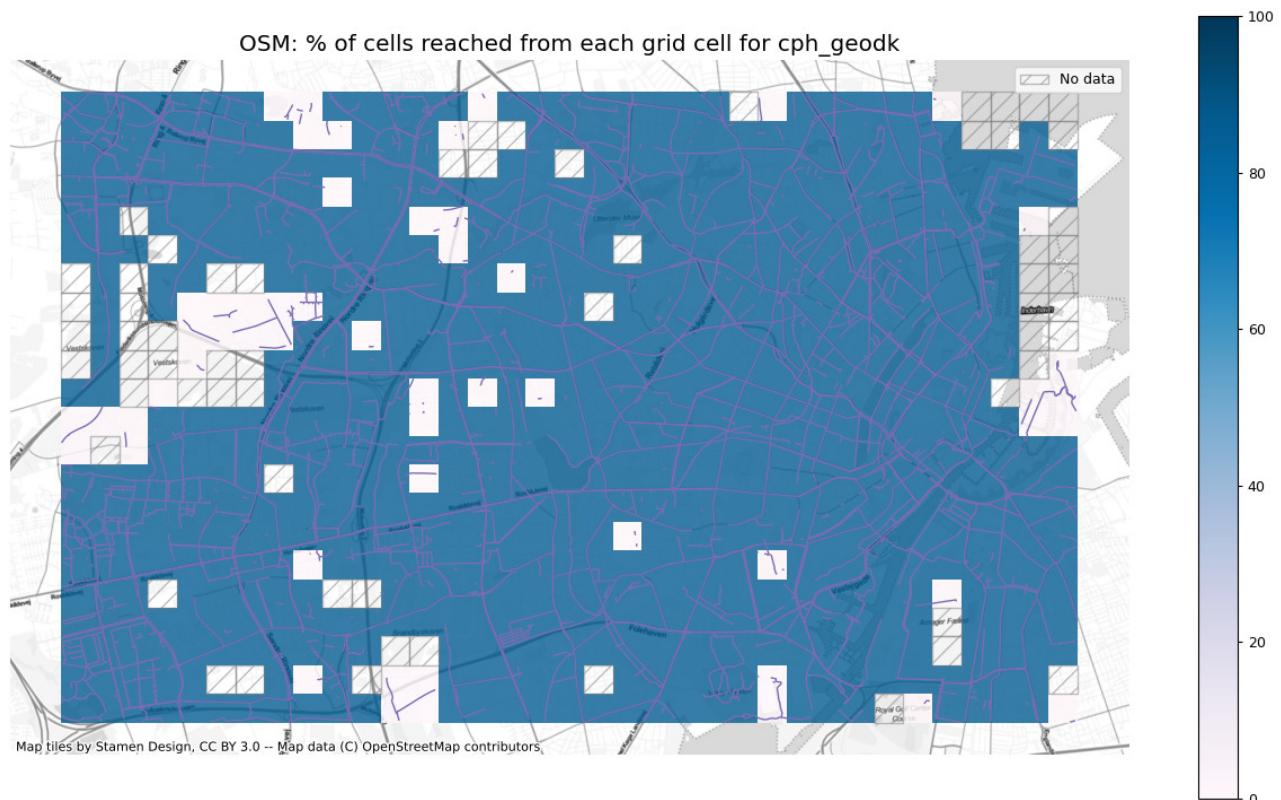
Analysis with component distance threshold of 10 meters:



Component connectivity

Visualizing differences between how many cells can be reached from each cell.

This is a crude measure for network connectivity but has the benefit of being computationally cheap and thus able to highlight stark differences in network connectivity in very little time.



Summary

Intrinsic Quality Metrics - OSM data

Total infrastructure length (km) 1,063

Protected bicycle infrastructure density (m/km2)	5,300
Unprotected bicycle infrastructure density (m/km2)	503
Mixed protection bicycle infrastructure density (m/km2)	59
Bicycle infrastructure density (m/km2)	5,862
Total number of nodes (count)	4,927
Dangling nodes	1,819
Nodes per km2	27
Dangling nodes per km2	10
Incompatible tag combinations	2
Overshoots	9
Undershoots	14
Missing intersection nodes	1
Components	352
Length of largest component (km)	752
Largest component's share of network length	92%
Component gaps	91

Time of analysis: Sat, 10 Dec 2022 00:19:47

2a. Load and Process Reference Data

This notebook:

- Loads the polygon defining the study area and then creates a grid overlay for the study area.
- Loads the reference data.
- Processes the reference data to create the network structure and attributes needed in the analysis.

Sections

- [Load data for study area and define analysis grid](#)
 - [Load and process reference data](#)
-

Load data for study area and define analysis grid

This step:

- Loads settings for the analysis from the configuration file.
- Reads data for the study area.
- Creates a grid overlay of the study area, with grid cell size as defined in configuration file `config.yml`.

Read in data for study area

The study area is defined by the user-provided polygon. It will be used for the computation of **global** results, i.e. for the entire study area.

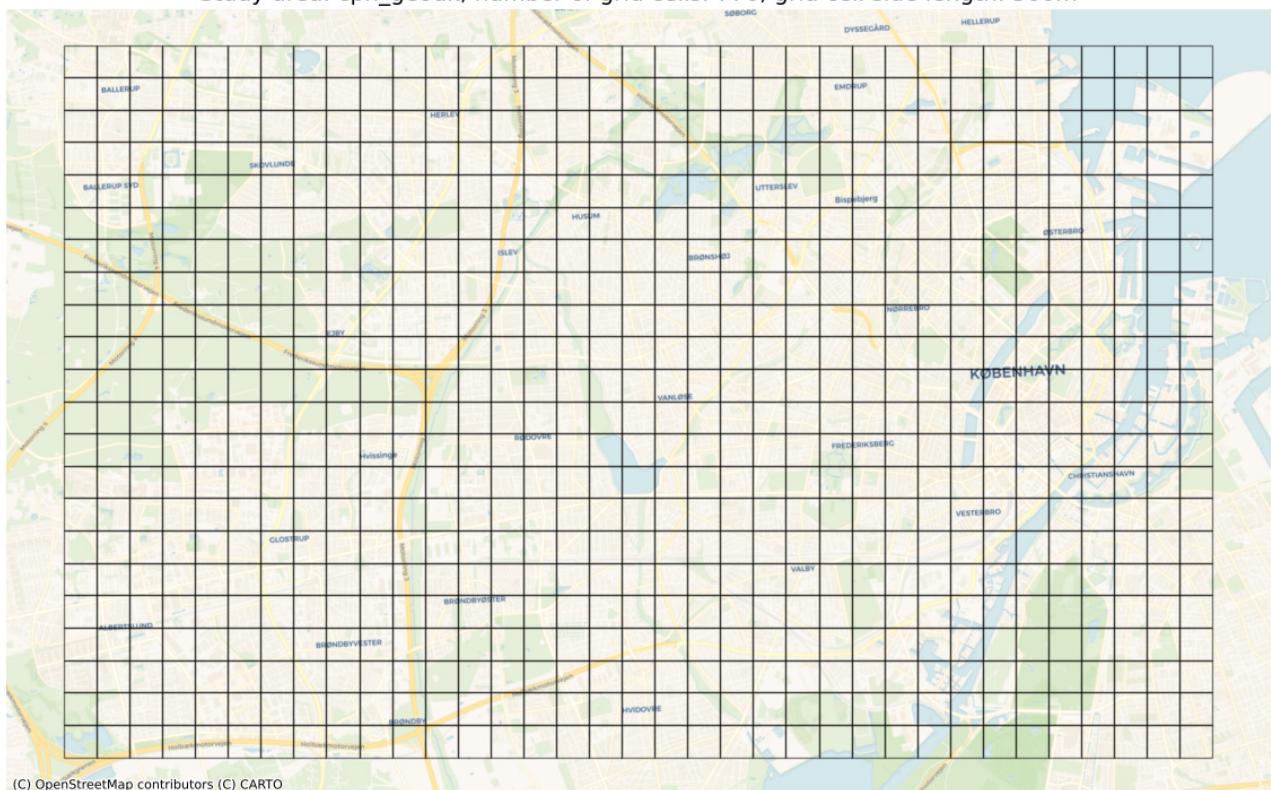
The size of the study area is 181.38 km².

Create analysis grid

The grid contains 770 square cells with a side length of 500 m and an area of 0.25 km².

This grid will be used for local (grid cell level) analysis:

Study area: cph_geodk, number of grid cells: 770, grid cell side length: 500m



Load and process reference data

This step:

- Creates a network from the reference data.
- Projects it to the chosen CRS.
- Clips the data to the polygon defining the study area.
- Measures the infrastructure length of the edges based on the geometry type and whether they allow for bidirectional travel or not.
- Simplifies the network.
- Creates copies of all edge and node datasets indexed by their intersecting grid cell.

The reference data covers an area of 169.76 km².

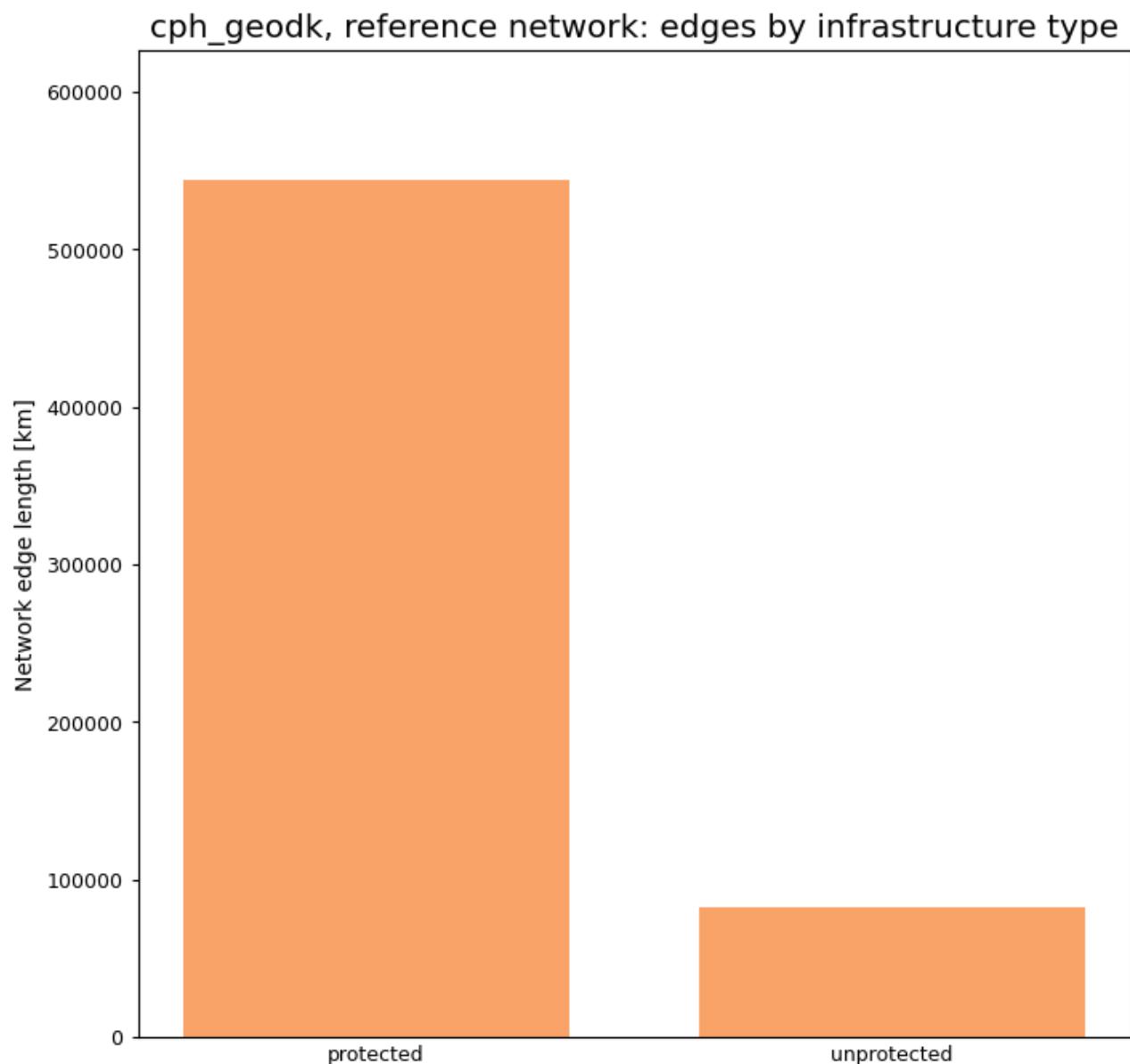
Number of edges where the protection level is 'protected': 46097 out of 53580 (86.03%)

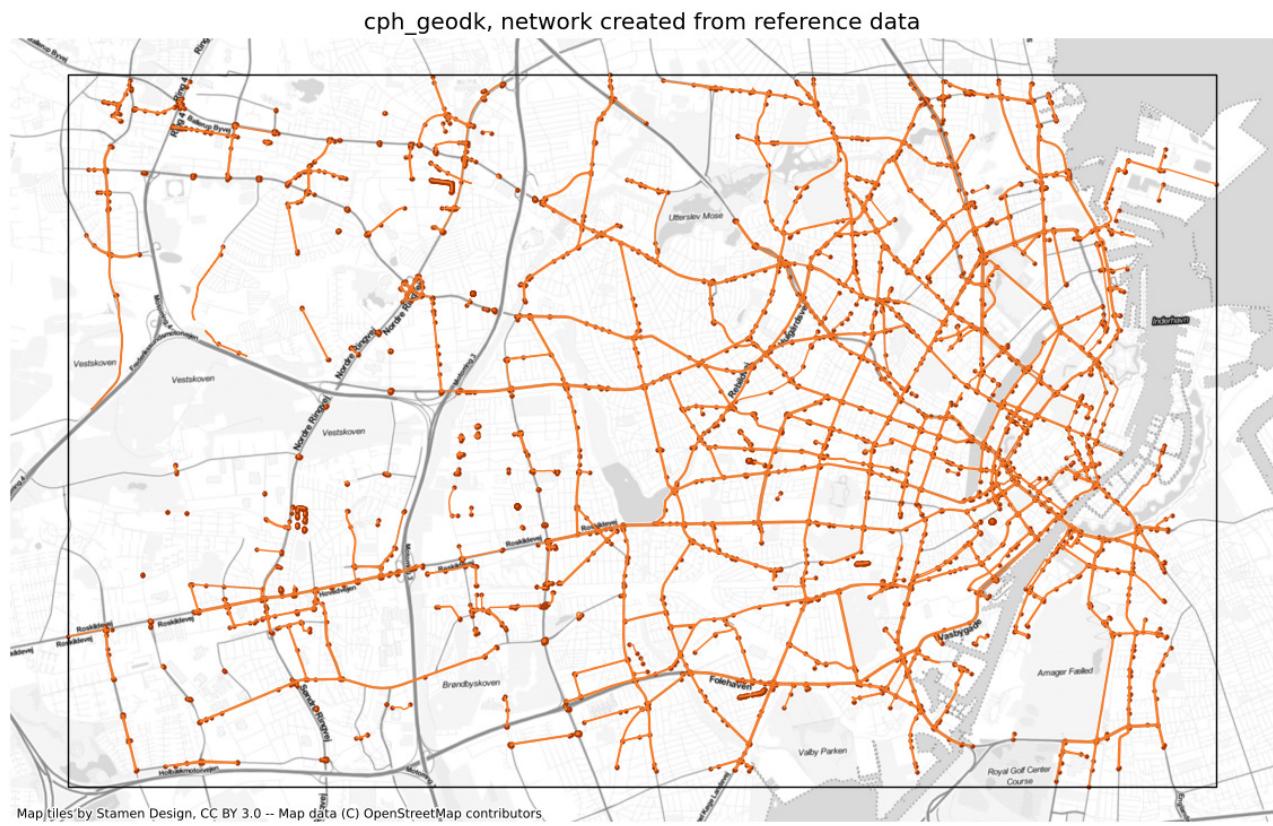
Number of edges where the protection level is 'unprotected': 7483 out of 53580 (13.97%)

No column on bidirectional edges has been found in the data set.

No column on geometry types has been found in the data set.

The length of the reference network is 626.48 km.





Time of analysis: Sat, 10 Dec 2022 13:09:57

2b. Intrinsic Analysis of Reference Bicycle Network Data

This notebook analyses the quality of a user-provided reference bicycle infrastructure data set for a given area. The quality assessment is *intrinsic*, i.e. based only on the input data set, and making no use of information external to the data set. For an extrinsic quality assessment that compares the reference data set to corresponding OSM data, see the notebooks 3a and 3b.

The analysis assesses the *fitness for purpose* ([Barron et al., 2014](#)) of the reference data for a given area. Outcomes of the analysis can be relevant for bicycle planning and research - especially for projects that include a network analysis of bicycle infrastructure, in which case the topology of the geometries is of particular importance.

Since the assessment does not make use of an external reference dataset as the ground truth, no universal claims of data quality can be made. The idea is rather to enable those working with bicycle networks to assess whether their data is good enough for their particular use case. The analysis assists in finding potential data quality issues, but leaves the final interpretation of the results to the user.

The notebook makes use of quality metrics from a range of previous projects investigating OSM/VGI data quality, such as [Antoniou & Skopeliti \(2015\)](#), [Fonte et al. \(2017\)](#) and [Fester et al. \(2020\)](#).

Familiarity required

For a correct interpretation of some of the metrics for spatial data quality, some familiarity with the area is necessary.

Sections

- [Data completeness](#)
 - Network density
- [Network topology](#)
 - Simplification outcome
 - Dangling nodes
 - Under/overshoots
- [Network components](#)
 - Disconnected components
 - Potential missing links
- [Summary](#)

- Save results

Data completeness

Network Density

In this setting, network density refers to the length of edges or number of nodes per square kilometer (i.e., the definition of network density usually used when looking at street networks, which is distinct from the definition usually found in graph theory). Network density without comparing to a reference dataset does not in itself indicate spatial data quality. For anyone familiar with the study area, network density can however indicate whether parts of the area appear to be under- or over-mapped and is thus included here.

Method

The density here is not based on the geometric length of edges, but instead on the computed length of the infrastructure. For example, a 100-meter-long bidirectional path contributes with 200 meters of bicycle infrastructure. With `compute_network_density`, the number of elements (nodes; dangling nodes; and total infrastructure length) per unit area is calculated. The density is computed twice: first for the study area for both the entire network ('global density'), then for each of the grid cells ('local density'). Both global and local densities are computed for the entire network and for respectively protected and unprotected infrastructure.

Interpretation

Since the analysis conducted here is intrinsic, i.e., it makes no use of external information, it cannot be known whether a low-density value is due to incomplete mapping, or due to actual lack of infrastructure in the area. However, a comparison of the grid cell density values can provide some insights, for example:

- lower-than-average infrastructure density indicates a locally sparser network
- higher-than-average node density indicates that there are relatively many intersections in the grid cell
- higher-than-average dangling node density indicates that there are relatively many dead ends in the grid cell

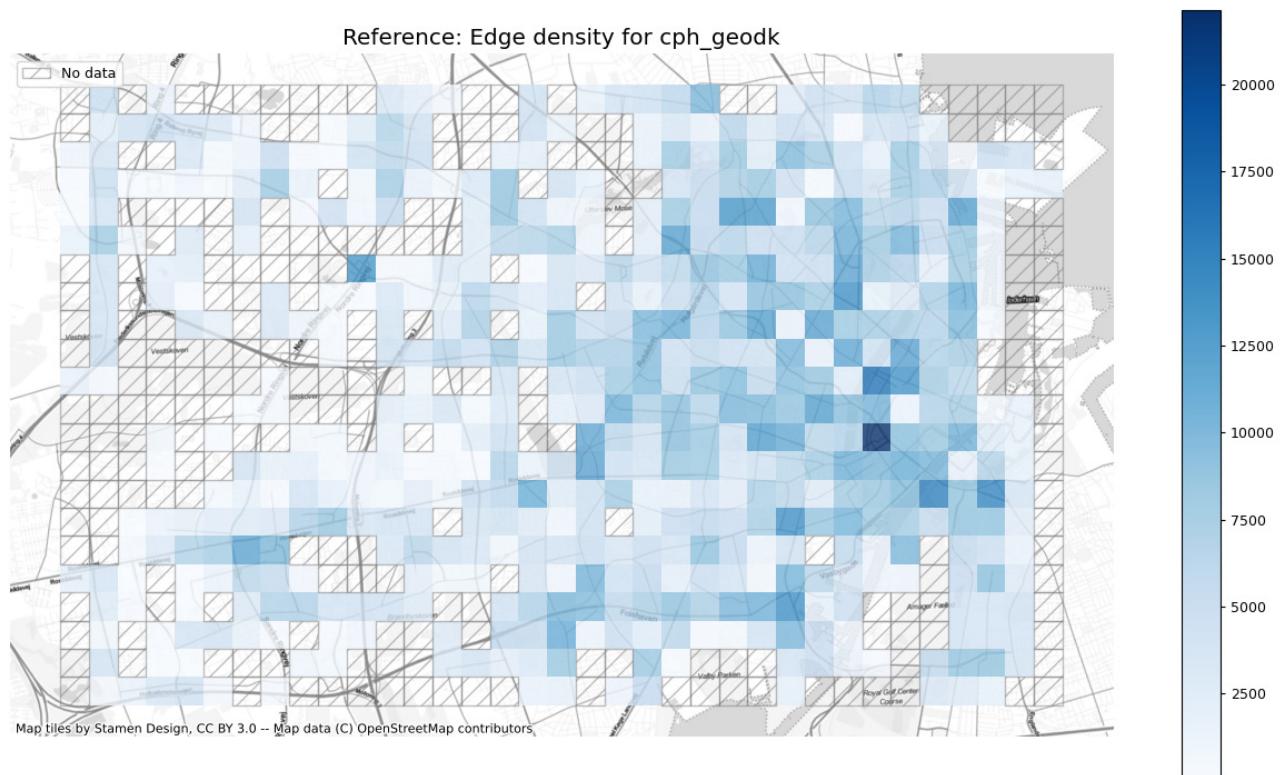
Global network density

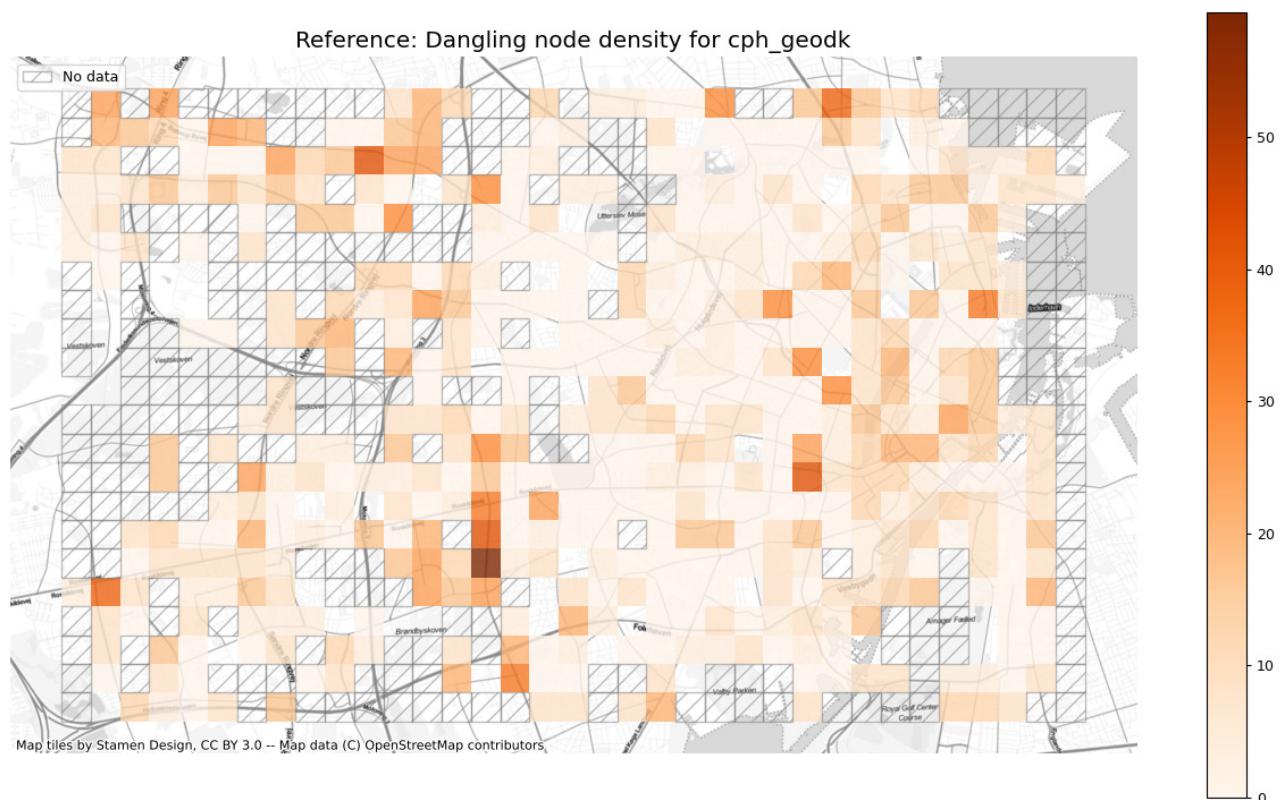
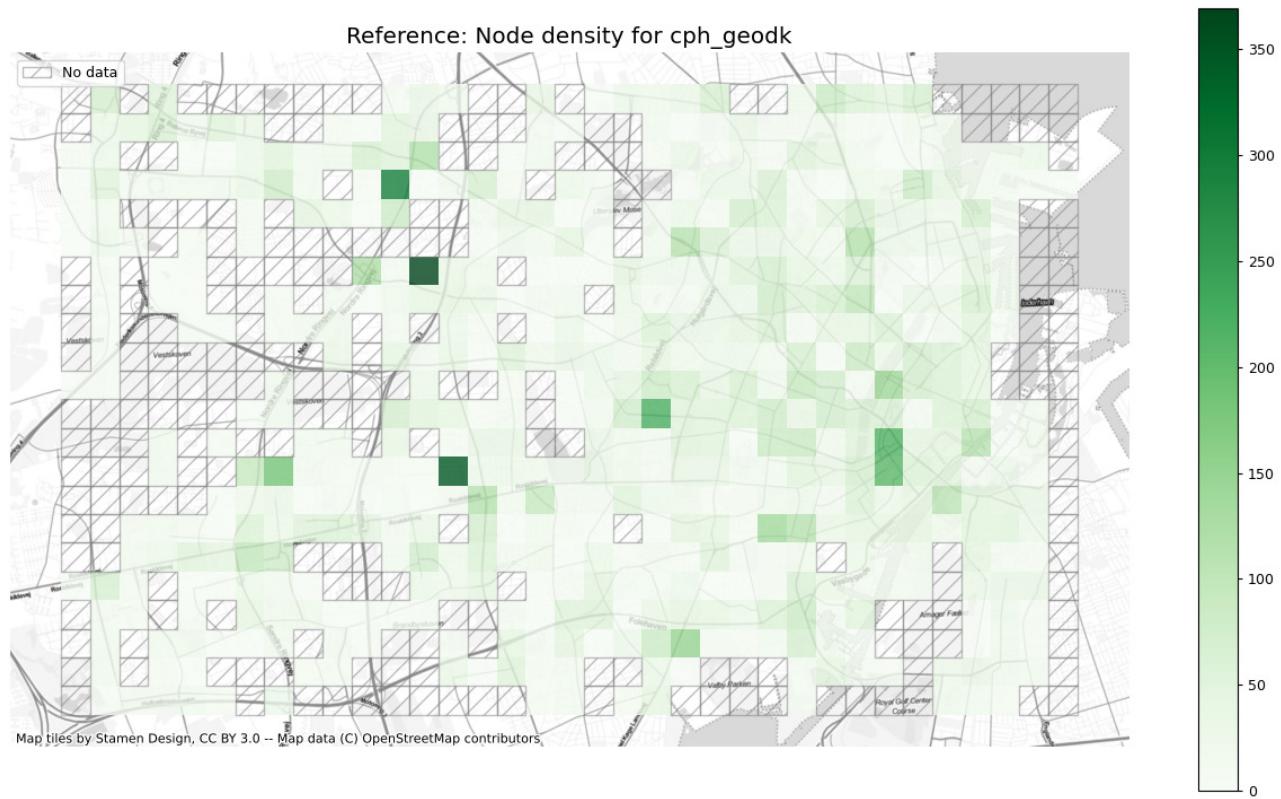
For the entire study area, there are:

- 3453.85 meters of bicycle infrastructure per km².
- 22.74 nodes in the bicycle network per km².
- 4.80 dangling nodes in the bicycle network per km².
- 2998.80 meters of protected bicycle infrastructure per km².

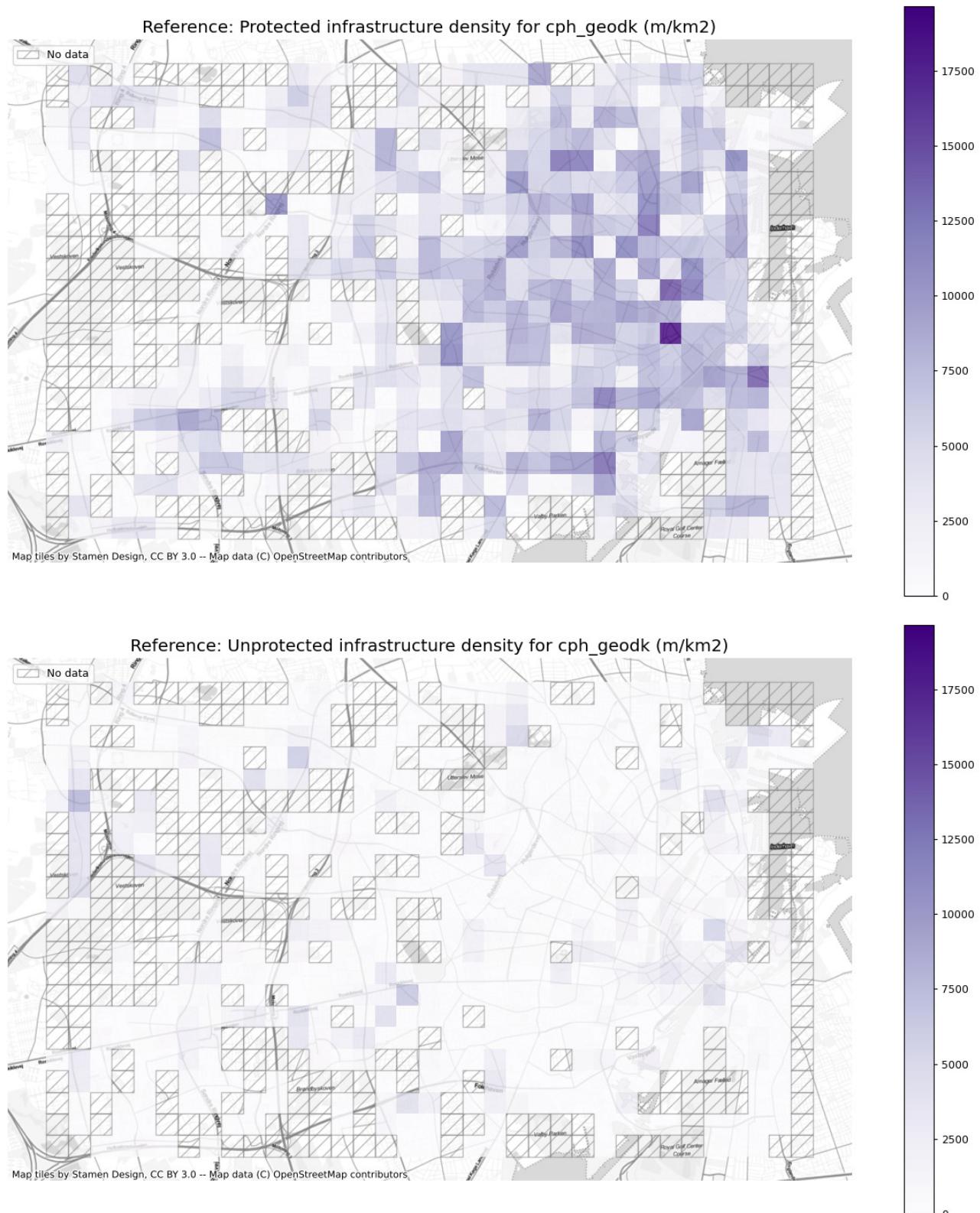
- 455.05 meters of unprotected bicycle infrastructure per km².
- 0.00 meters of mixed protection bicycle infrastructure per km².

Local network density





Densities of protected and unprotected infrastructure:



Network topology

This section explores the geometric and topological features of the data.

These are, for example, network density, disconnected components, dangling (degree one) nodes; it also includes exploring whether there are nodes in close proximity, that do not share an edge - a potential sign of edge undershoots - or if there are intersecting edges without a node at the intersection, which might indicate a digitizing error that will distort any attempts at routing on the network.

Due to the fragmented nature of most networks of bicycle infrastructure, many metrics, such as missing links or network gaps, simply reflect the true extent of the infrastructure ([Natera Orozco et al., 2020](#)). This is different for car networks, where e.g., disconnected components could more readily be interpreted as a data quality issue.

Therefore, the analysis only takes very small network gaps into account as potential data quality issues.

Subsections:

- [Simplification outcome](#)
- [Dangling nodes](#)
- [Under/Overshoots](#)

Simplification outcome

When converting a set of geocoded linestrings (polygonal chains) to graph format, not all vertices (nodes) are of equal meaning. For geometry of the infrastructural element, all nodes are needed as an ordered list. For the topology of the network, however, only those nodes that are endpoints or intersection points with other edges are needed, while all other (so-called 'interstitial') nodes do not add any information. To compare the structure and true ratio between nodes and edges in a network, a simplified network representation which only includes nodes at endpoints and intersections, or where the value of important attributes changes, is required. Therefore, in the notebook `01_load_data` the bicycle network was simplified by removing all interstitial nodes from the graph object (retaining, however, the complete node lists in the geometry attribute of each edge). An additional advantage of simplifying the network is the resulting substantial reduction of the number of nodes and edges, which makes computational routines much faster.

Comparing the node degree distribution for the networks before and after simplification is a quick sanity check for the simplification routine. Typically, the big majority of nodes in the non-simplified network will be of degree two; in the simplified network, however, most nodes will have degrees other than two. Degree two nodes are retained in only two cases: if they represent a connection point between two different types of infrastructure; or if they are needed in order to avoid self-loops (edges whose start and end points are identical) or multiple edges between the same pair of nodes.

As part of the simplification routine, in cases where there are several edges between the same pair of nodes ('parallel edges' or 'multiedges'), only one of the edges is retained. Within the routine, the number edges removed in this way are counted.

Method

The node degree distributions before and after simplification are plotted below.

Interpretation

Typically, the node degree distribution will go from high (before simplification) to low (after simplification) counts of degree two nodes, while it will not change for all other degrees (1, or 3 and higher). Further, the total number of nodes will see a strong decline. If the simplified graph still maintains a relatively high number of degree two nodes, or if the number of nodes with other degrees changes after the simplification, this might point to issues either with the graph conversion or with the simplification process.

Simplifying the network decreased the number of edges with 91.2% and the number of nodes with 92.2%.

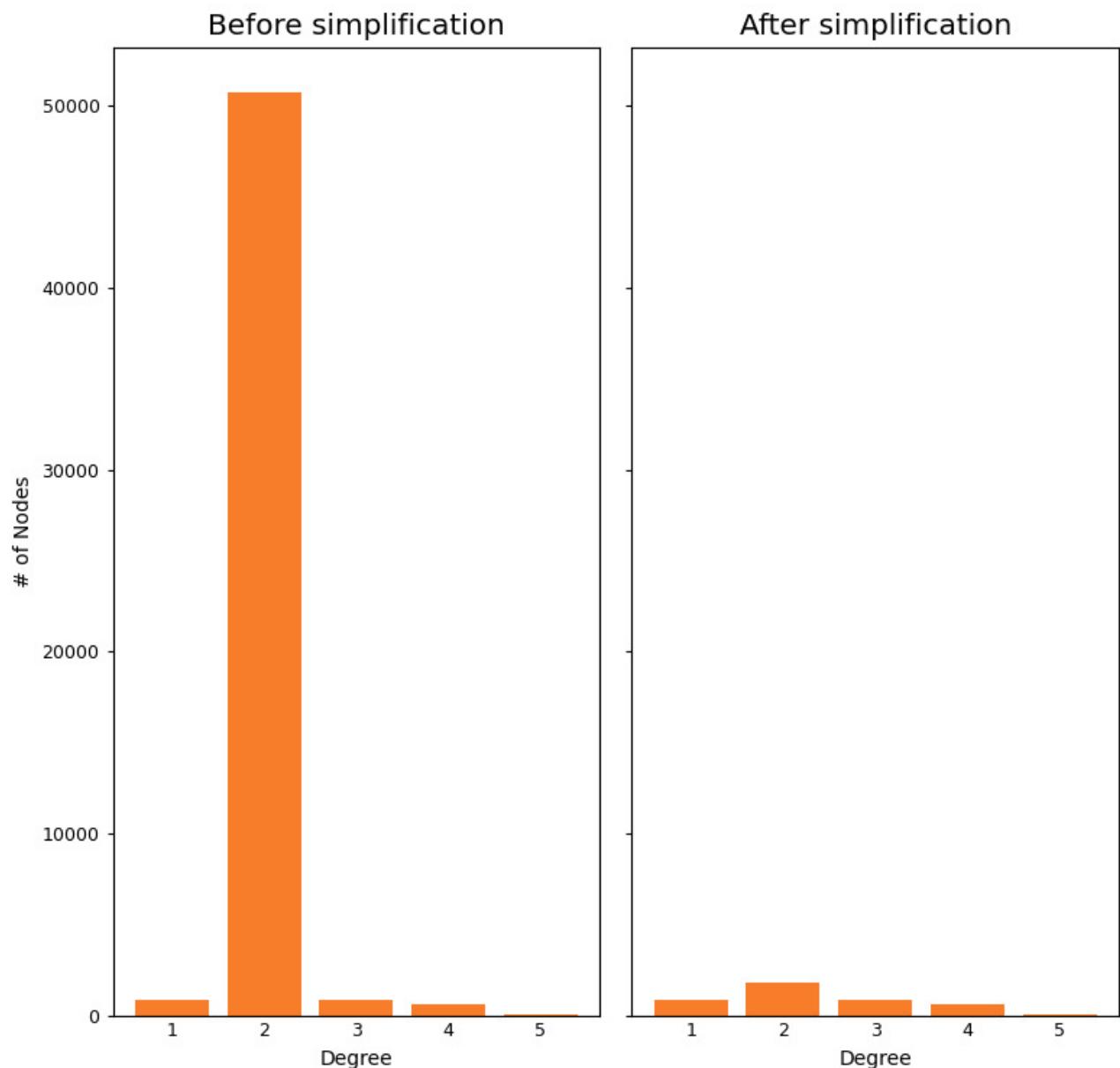
Before the network simplification the Reference graph had:

- 7 node(s) with degree 5
- 591 node(s) with degree 4
- 827 node(s) with degree 3
- 50705 node(s) with degree 2
- 870 node(s) with degree 1

After the network simplification the Reference graph had:

- 7 node(s) with degree 5
- 591 node(s) with degree 4
- 827 node(s) with degree 3
- 1830 node(s) with degree 2
- 870 node(s) with degree 1

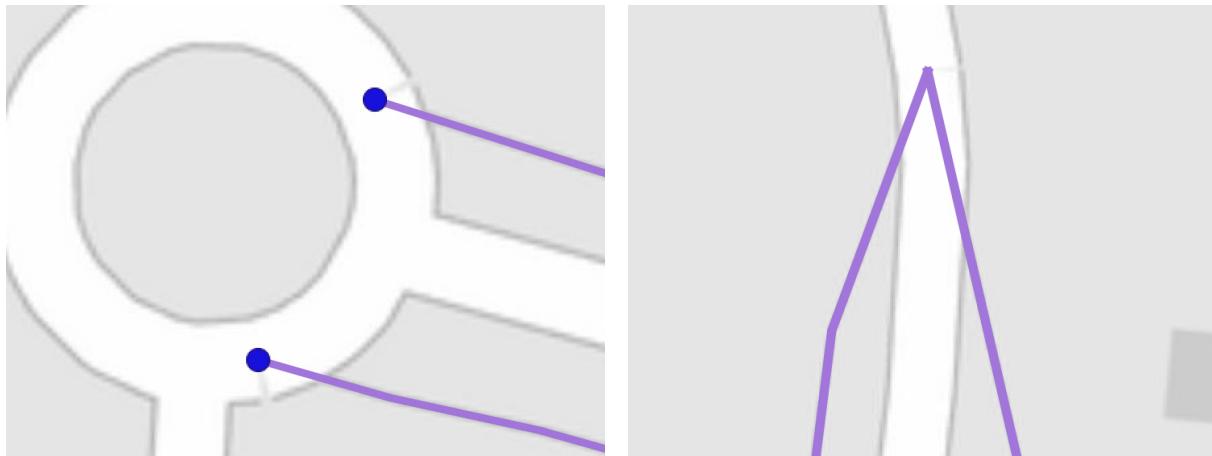
Reference: Node degree histograms for the network in cph_geodk



Dangling nodes

Dangling nodes are nodes of degree one - in other words, nodes that have only one single edge attached to them. Most networks will naturally contain a number of dangling nodes. Dangling nodes can occur at actual dead-ends (representing a cul-de-sac) or at the endpoints of certain features (e.g., when a bicycle path ends in the middle of a street). However, dangling nodes can also occur as a data quality issue in case of over/undershoots (as described in detail in the next section). The number of dangling nodes in a network does to some extent also depend on the digitization method, as shown in the illustration below.

Therefore, the presence of dangling nodes is in itself not a sign of low data quality. However, a high number of dangling nodes in an area that is not known for suffering from many dead-ends can indicate digitization errors and problems with edge over/undershoots.



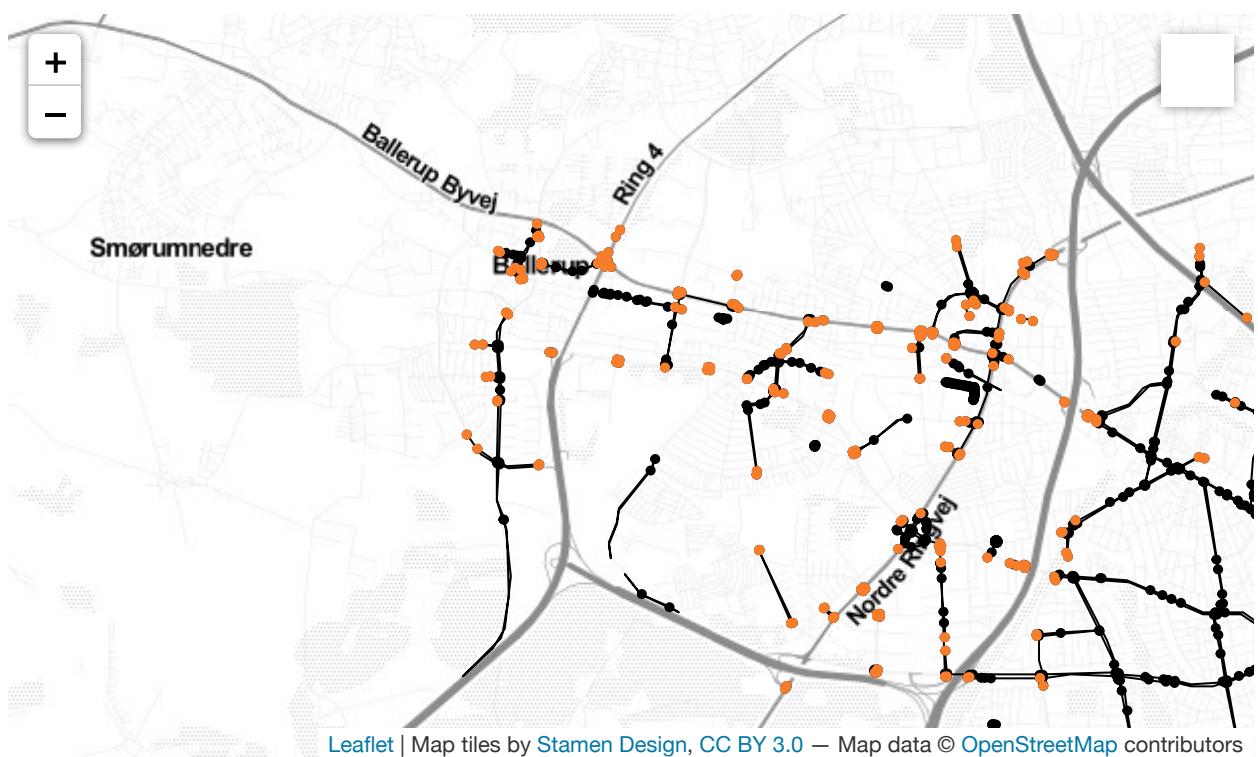
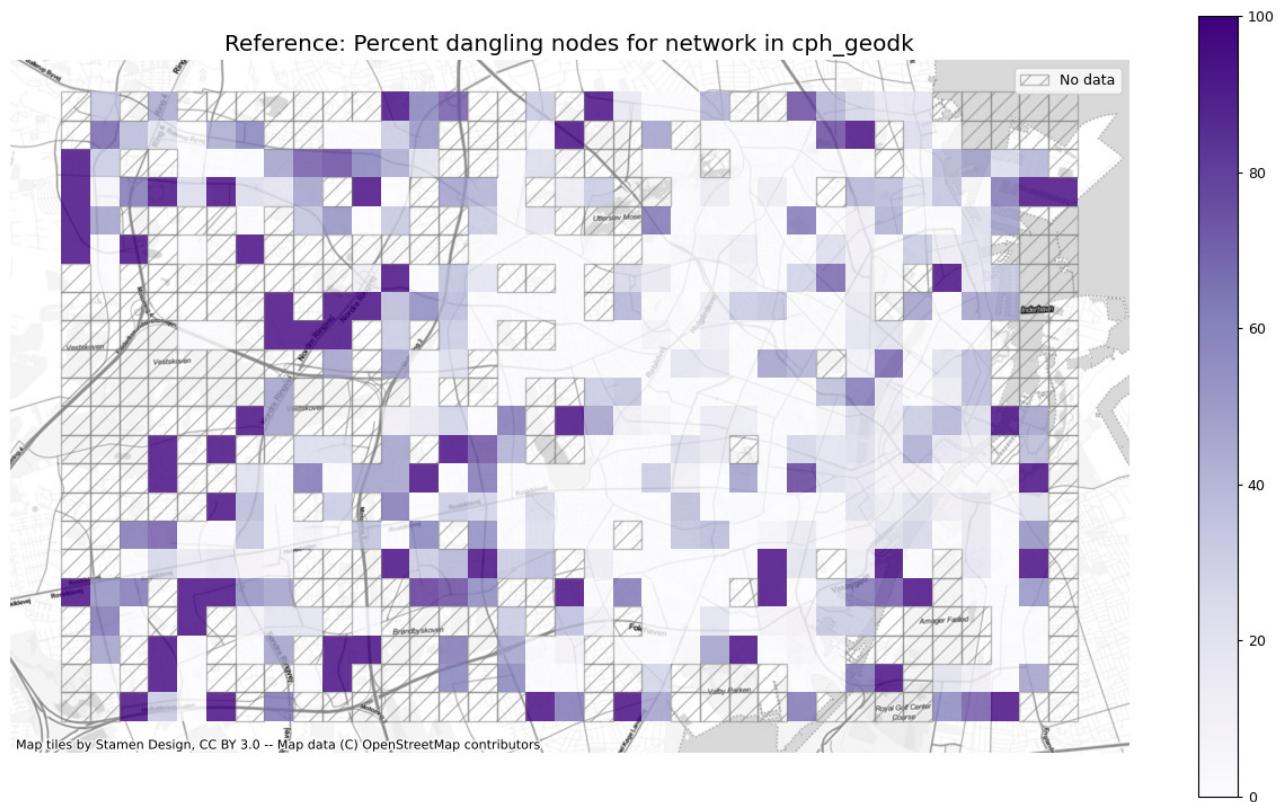
Dangling nodes occur where road features end (left), but when separate features are joined at the end (right), there will be no dangling nodes

Method

Below, a list of all dangling nodes is obtained with the help of `get_dangling_nodes`. Then, the network with all its nodes is plotted. The dangling nodes are shown in orange; all other nodes are shown in black.

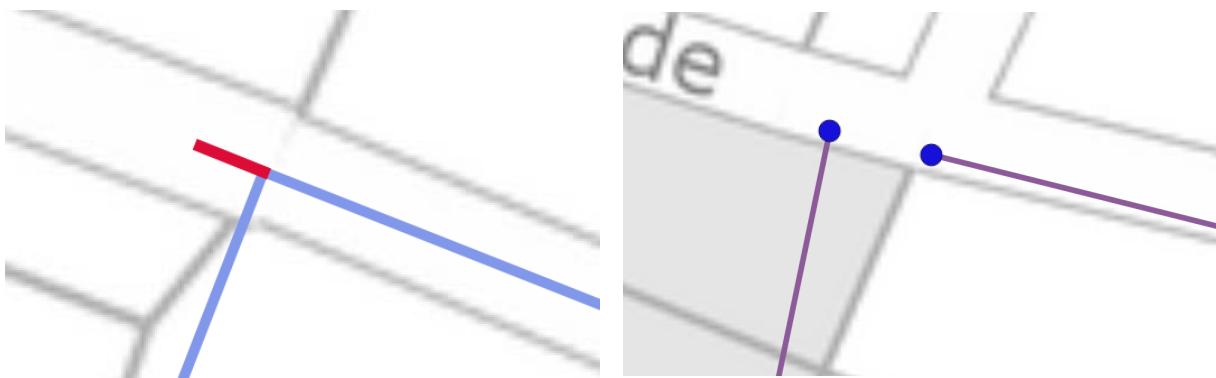
Interpretation

We recommend a visual analysis in order to interpret the spatial distribution of dangling nodes, with particular attention to areas of high dangling node density. It is important to understand where dangling nodes come from: are they actual dead-ends or digitization errors (e.g., over/undershoots)? A higher number of digitization errors points to a lower quality of the data.



Under/overshoots

When two nodes in a simplified network are placed within a distance of a few meters, but do not share a common edge, it is often due to an edge over/undershoot or another digitizing error. An overshoot occurs when two features meet and one of them extends beyond the other. An undershoot occurs when two features are supposed to meet, but instead are just in close proximity to each other. See the image below for an illustration of an overshoot (left) and an undershoot (right). For a more detailed explanation of over/undershoots, see the [GIS Lounge website](#).



Overshoots refer to situations where a line feature extends too far beyond at intersecting line, rather than ending at the intersection (left). Undershoots happen when two line features are not properly joined, for example at intersection (right)

Method

Overshoots: First, the `length_tolerance` (in meters) is defined in the cell below. Then, with `find_overshoots`, all network edges that have a dangling node attached to them and that have a maximum length of `length_tolerance` are identified as overshoots, and the results are plotted.

Undershoots: First, the `length_tolerance` (in meters) is defined in the cell below. Then, with `find_undershoots`, all pairs of dangling nodes that have a maximum of `length_tolerance` distance between them, are identified as undershoots, and the results are plotted.

The workflow for over/undershoot detection below is inspired by [Neis et al. 2012](#).

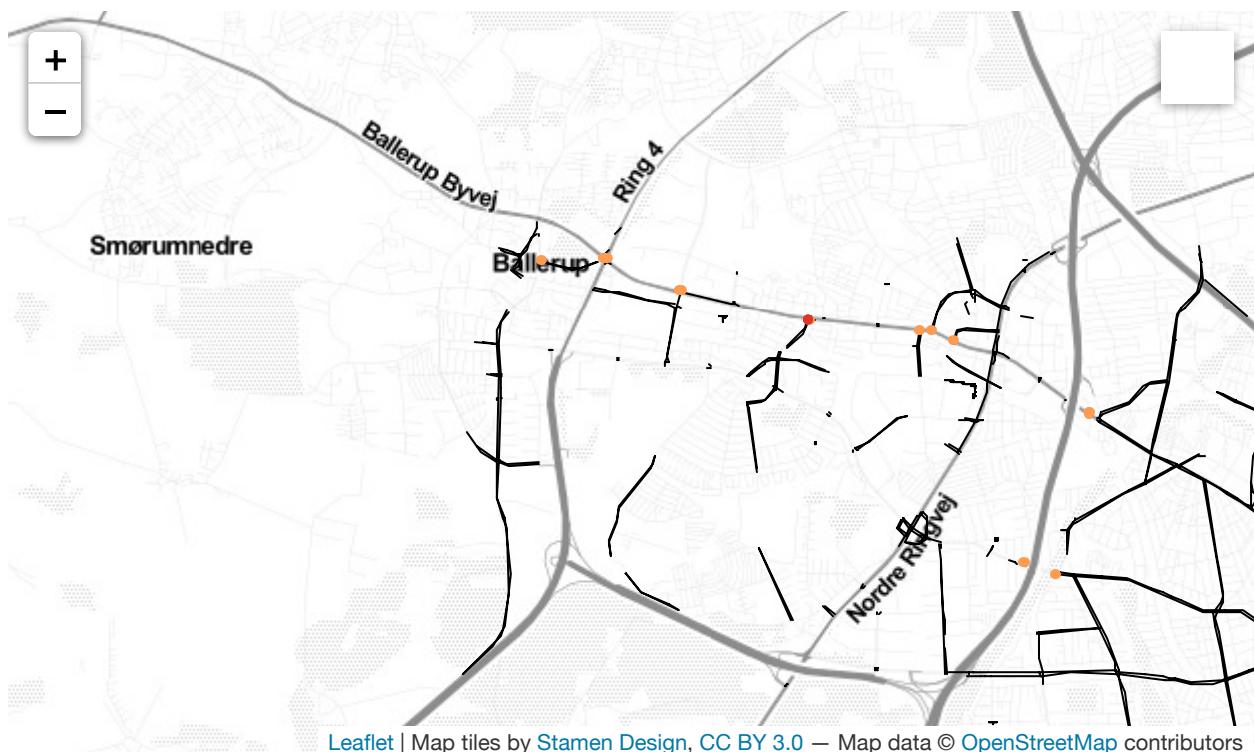
Interpretation

Note that over/undershoots are not necessarily always a data quality issue - they might be instead an accurate representation of the network conditions or of the digitization strategy (for example, a cycle path might end abruptly soon after a turn, which results in an overshoot; protected cycle paths are often digitized in OSM as interrupted at intersections, which results in intersection 'undershoots').

The interpretation of the impact of over/undershoots on data quality is context dependent. For certain applications, such as routing, overshoots do not present a particular challenge; they can, however, pose an issue for other applications such as network analysis, given that they skew the network structure. Undershoots, on the contrary, are a serious problem for routing applications, especially if only bicycle infrastructure is considered; they also pose a problem for network analysis if one wants to identify connected components.

Under/overshoots

21 potential overshoots were identified with a length tolerance of 3 meters.
11 potential undershoots were identified with a length tolerance of 3 meters.



Network components

Disconnected components

Disconnected components do not share any elements (nodes/edges). In other words, there is no network path that could lead from one disconnected component to the other. As mentioned above, most real-world networks of bicycle infrastructure do consist of many disconnected components ([Natera Orozco et al., 2020](#)). However, when two disconnected components are very close to each other, it might be a sign of a missing edge or another digitizing error.

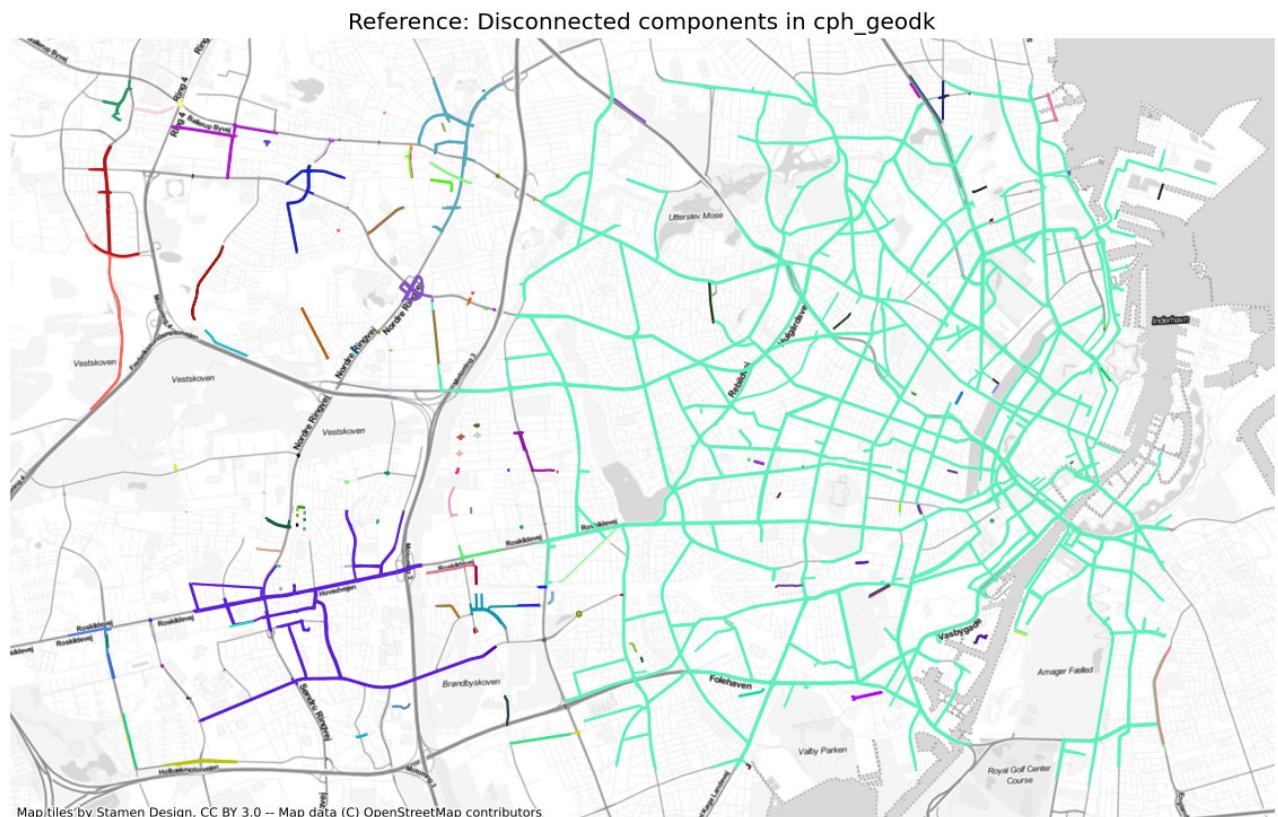
Method

First, with the help of `return_components`, a list of all (disconnected) components of the network is obtained. The total number of components is printed and all components are plotted in different colors for visual analysis. Next, the component size distribution (with components ordered by the network length they contain) is plotted, followed by a plot of the largest connected component.

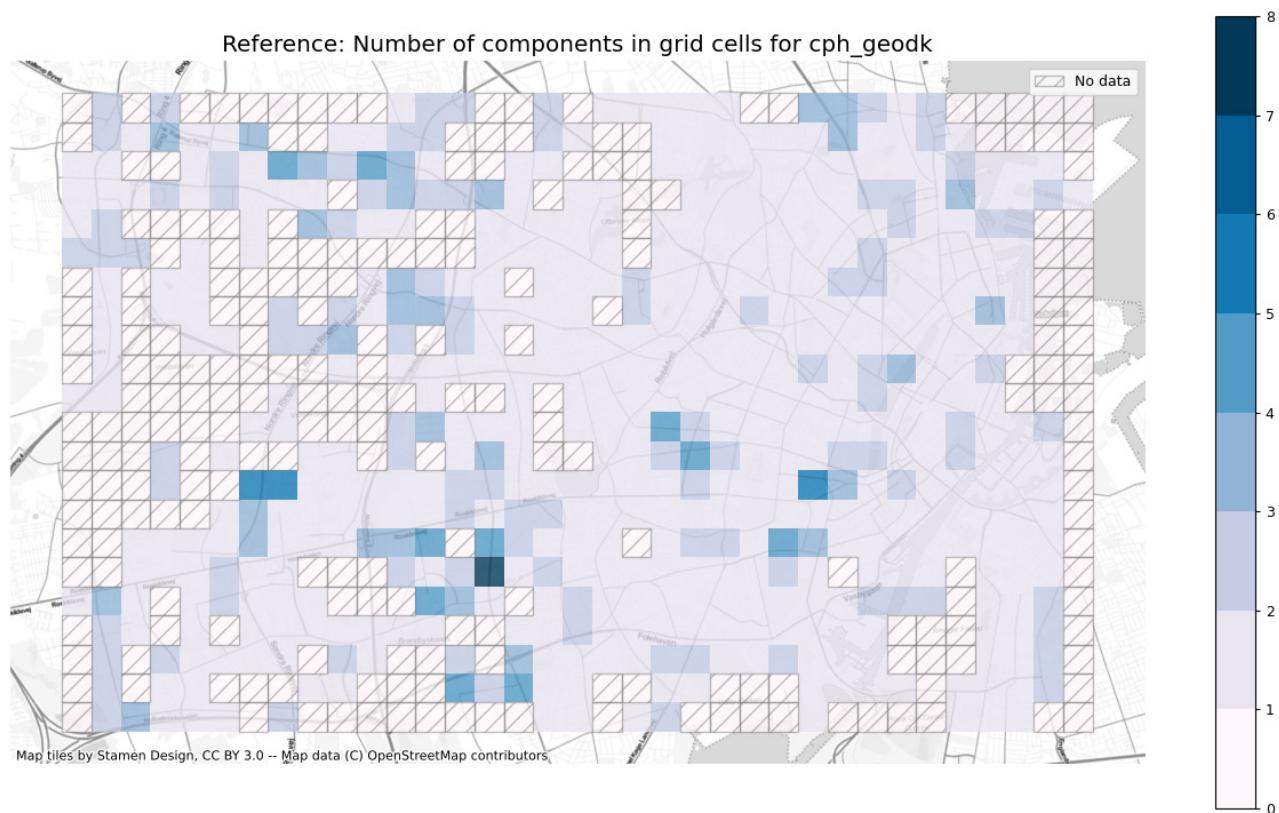
Interpretation

As with many of the previous analysis steps, knowledge of the area is crucial for a correct interpretation of component analysis. Given that the data represents the actual infrastructure accurately, bigger components indicate coherent network parts, while smaller components indicate scattered infrastructure (e.g., one single bicycle path along a street that does not connect to any other bicycle infrastructure). A high number of disconnected components in near vicinity of each other could indicate digitization errors or missing data.

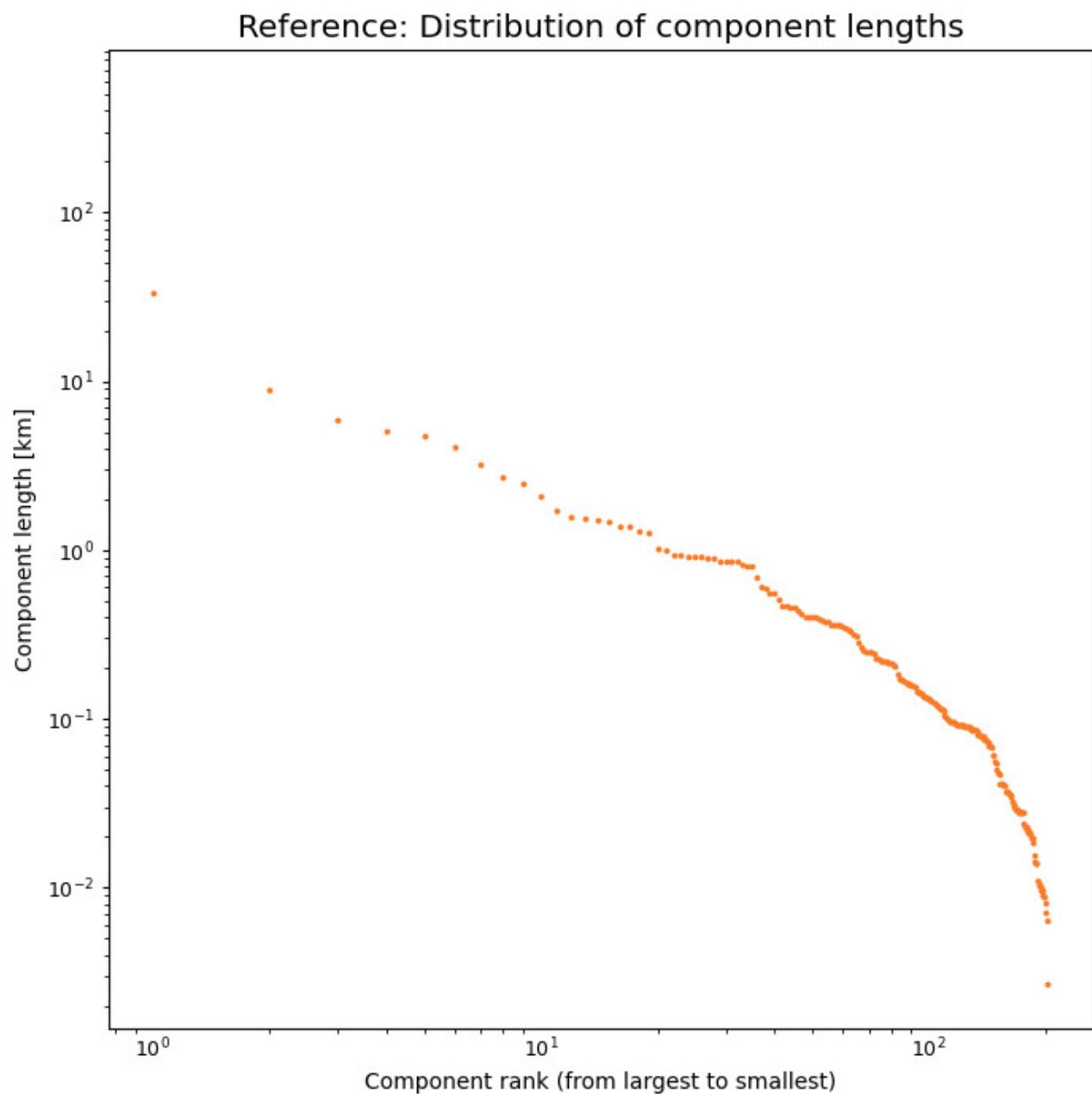
The network in the study area consists of 204 disconnected components.



Number of components per grid cell

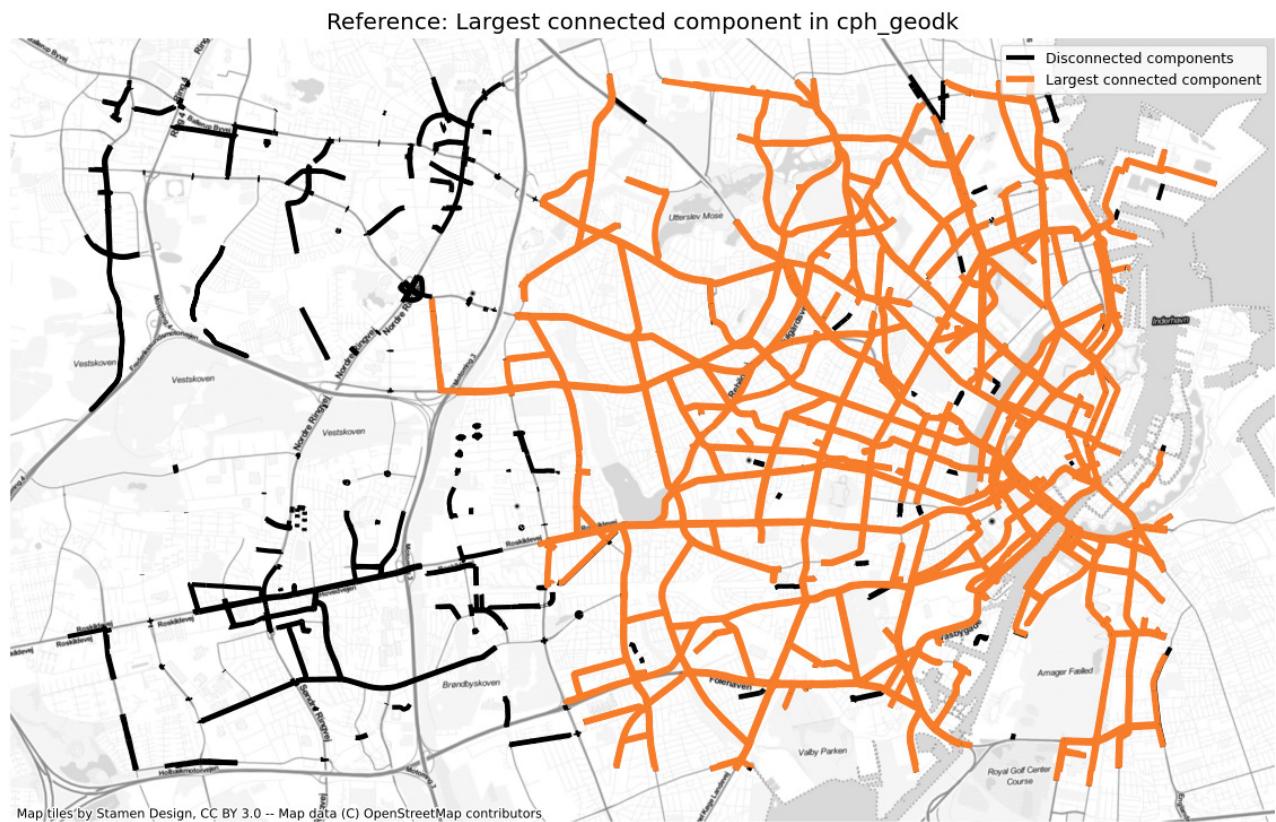


Distribution of network length per component



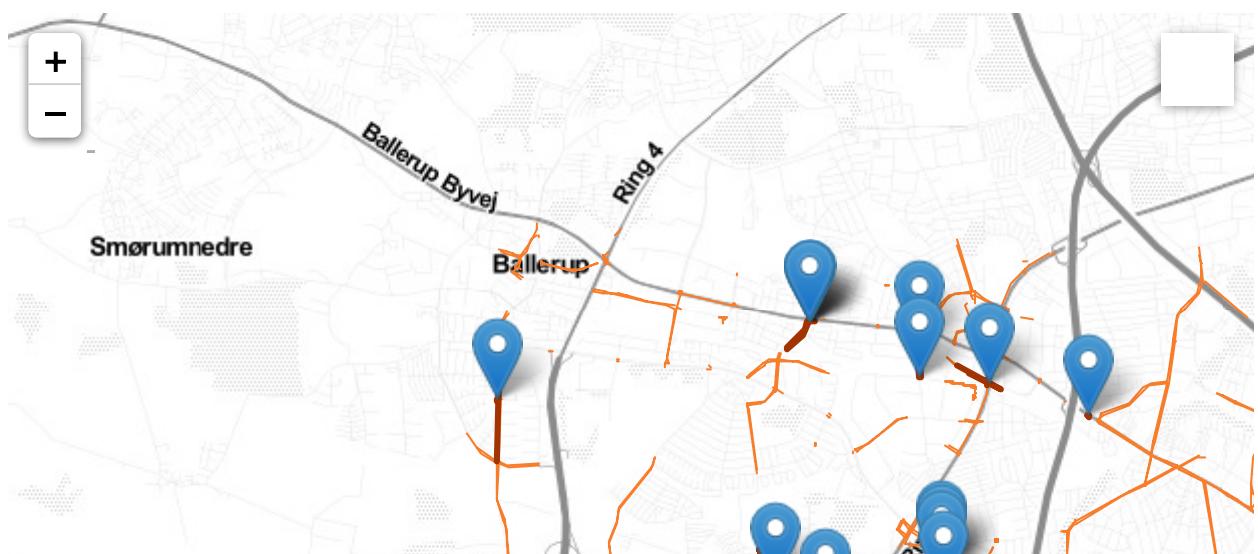
The largest connected component contains 80.04% of the network length.

Largest connected component



Potential missing links between disconnected components

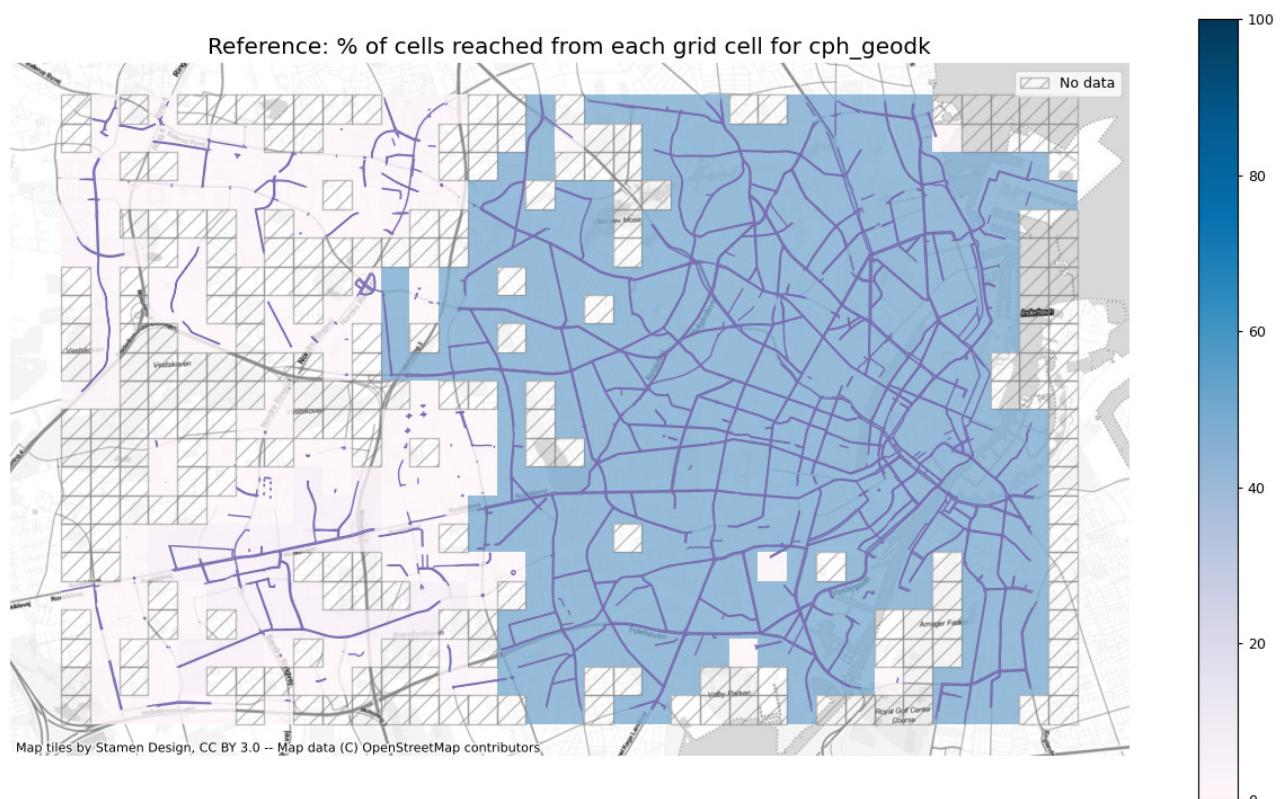
In the plot of potential missing links between components, all edges that are within the specified distance of an edge on another component are plotted. The gaps between disconnected edges are highlighted with a marker. The map thus highlights edges which, despite being in close proximity of each other, are disconnected and where it thus would not be possible to bike on cycling infrastructure between the edges.



Component connectivity

Visualizing differences between how many cells can be reached from each cell.

This is a crude measure for network connectivity but has the benefit of being computationally cheap and thus able to highlight stark differences in network connectivity in very little time.



Summary

Intrinsic Quality Metrics - Reference data

Total infrastructure length (km)

626

Protected bicycle infrastructure density (m/km²)

2,999

Unprotected bicycle infrastructure density (m/km2)	455
Mixed protection bicycle infrastructure density (m/km2)	0
Bicycle infrastructure density (m/km2)	3,454
Total number of nodes (count)	4,125
Dangling nodes	870
Nodes per km2	23
Dangling nodes per km2	5
Overshoots	21
Undershoots	11
Components	204
Length of largest component (km)	501
Largest component's share of network length	80%
Component gaps	52

Time of analysis: Sat, 10 Dec 2022 13:10:38

3a. Extrinsic Analysis: Comparison of OSM & Reference Data

This notebook compares the OSM data on bicycle infrastructure for a given area with a provided reference data set in a so-called extrinsic quality assessment. To run this part of the analysis, a reference data set must be available for comparison.

The analysis is based on comparing OSM data to the reference dataset and highlighting how and where they differ, both in terms of *how much* bicycle infrastructure is mapped in the two datasets, but also when it comes to *how* the infrastructure is mapped (e.g., looking at differing network topology and structure).

All differences are computed with OSM values as the baseline. For example, the difference in network density is computed by subtracting the reference density from the OSM density: $OSM - reference = difference$. Hence, positive difference values (over 0) indicate how much higher the OSM value is; negative difference values (below 0) indicate how much lower the OSM value is. Accordingly, if differences are given in percent, the OSM value is taken to be the total value (100%).

While the analysis is based on a comparison, the analysis makes no a priori assumptions about which dataset is better. The same goes for the identified differences: the workflow does not lead to an automatic conclusion as to which data set is of better quality, but instead requires the user to interpret the meaning of the differences found, e.g., whether differing features are results of errors of omission or commission, and which dataset is more correct.

The goal is that the identified differences can be used by the user to both assess the quality of the OSM and the reference datasets, and to support the decision of which dataset should be used for further analysis.

Familiarity required

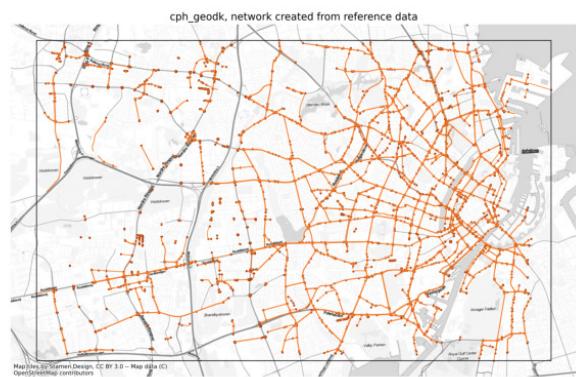
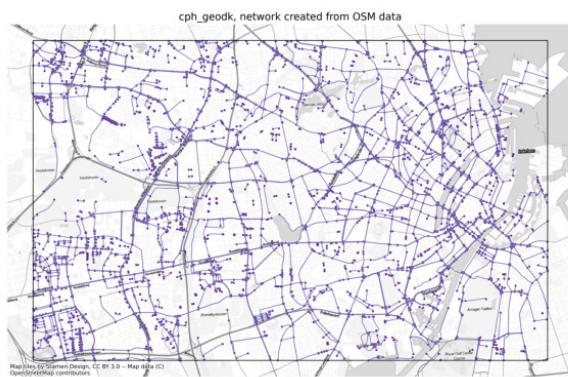
For a correct interpretation of some of the metrics for spatial data quality, some familiarity with the area is necessary.

Sections

- [Data completeness](#)
 - Network length
 - Network density
- [Network topology](#)

- Simplification outcome
- Alpha, beta, and gamma indices
- Dangling nodes
- Under & overshoots
- Network components
 - Disconnected components
 - Potential missing links
- Summary
- Save results

OSM & reference networks



Data completeness

This section compares the OSM and reference datasets in terms of data completeness.

The goal is to identify whether one dataset has more bicycle infrastructure mapped than the other, and if so, whether those differences are concentrated in some areas.

The section starts with a comparison of the total length of the infrastructure in both data sets. Then, infrastructure, node and dangling node densities (i.e., the length/number of infrastructure/nodes per square kilometer) is compared first at a global (study area) and at local (grid cell) level. Finally, density differences for protected and unprotected bicycle infrastructure are compared separately.

Method

To account for differences in how bicycle infrastructure has been mapped, the computation of network length and density is based on the infrastructure length, not the geometric length of the network edges. For example, a 100 meter long **bidirectional** path (geometric length: 100m) contributes with 200 meters of bicycle infrastructure (infrastructure length: 200m).

Interpretation

Density differences can point to incomplete data. For instance, if a grid cell has a significantly higher edge density in the OSM than in the reference data set, this indicates unmapped grid cell features in the reference data set (or potentially that a street mistakenly has been tagged as bicycle infrastructure.)

Network length

Length of the OSM data set: 1063.46 km

Length of the reference data set: 626.48 km

The OSM data set is 436.99 km longer than the reference data set.

The OSM data set is 41.09% longer than the reference data set.

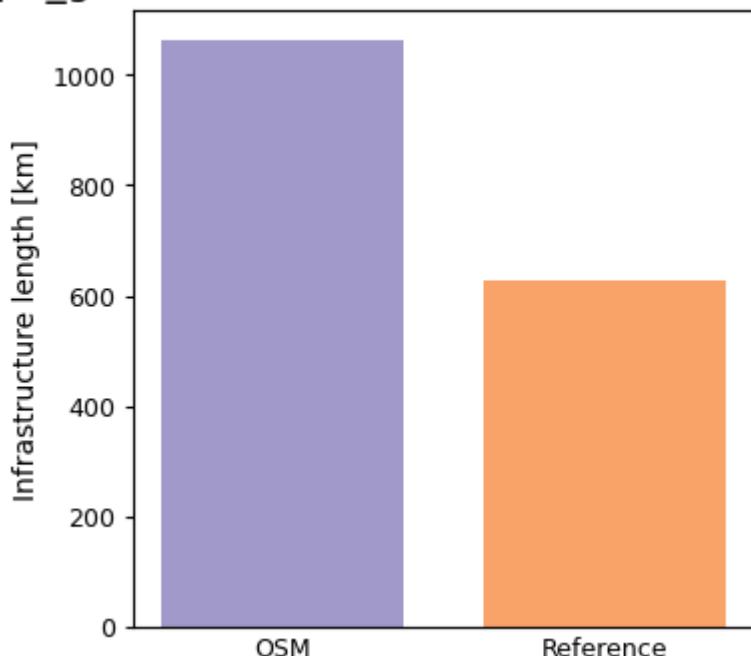
Length of the OSM data set: 1063.46 km

Length of the reference data set: 626.48 km

The OSM data set is 436.99 km longer than the reference data set.

The OSM data set is 41.09% longer than the reference data set.

cph_geodk: Total network infrastructure length



1b. Network Densities

Global network densities

In the OSM data, there are:

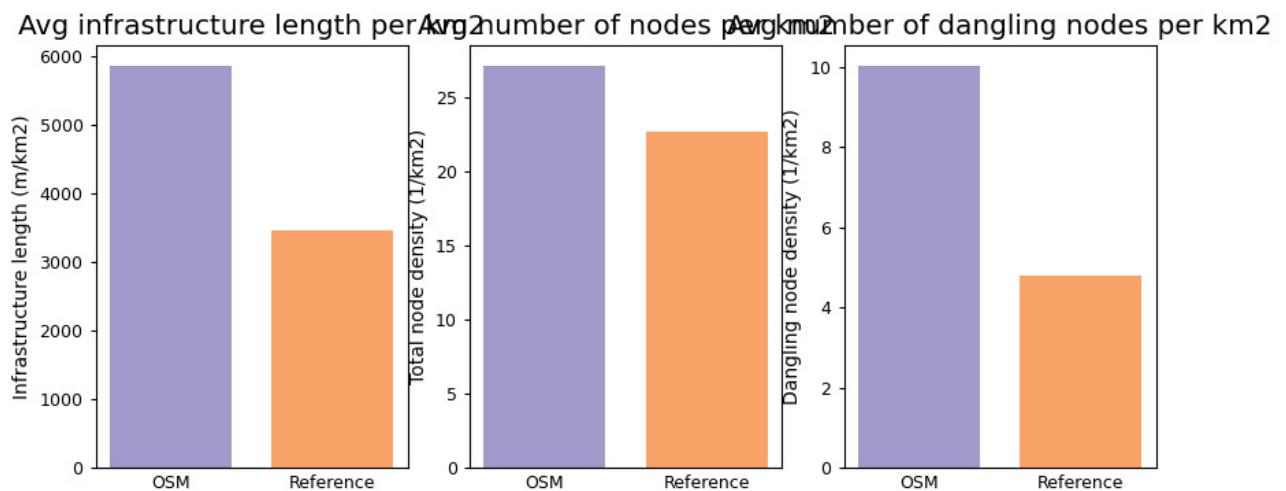
- 5862.22 meters of cycling infrastructure per km².

- 27.16 nodes in the cycling network per km2.
- 10.03 dangling nodes in the cycling network per km2.

In the reference data, there are:

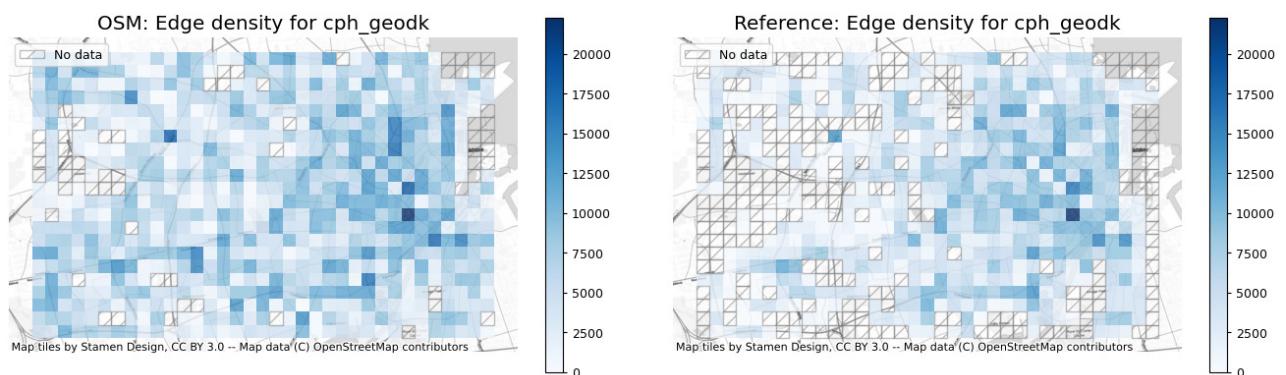
- 3453.85 meters of cycling infrastructure per km2.
- 22.74 nodes in the cycling network per km2.
- 4.80 dangling nodes in the cycling network per km2.

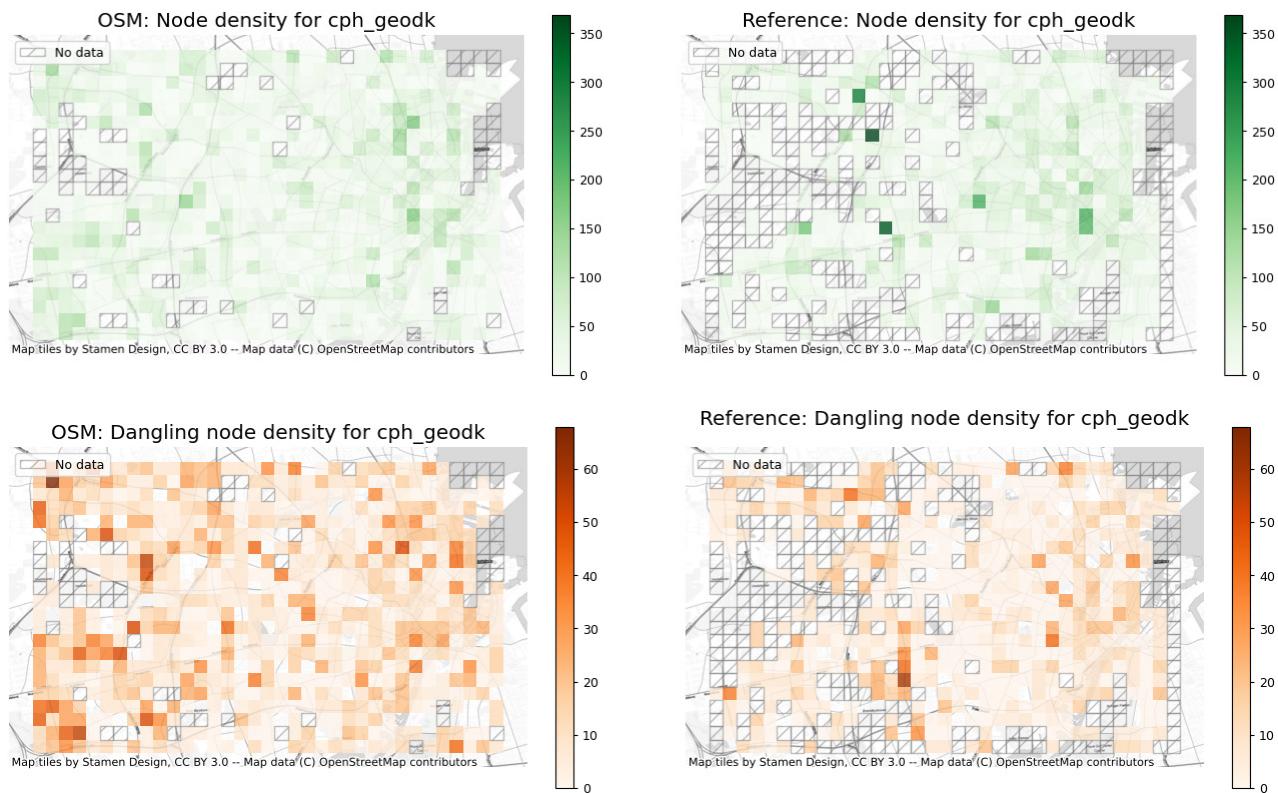
Global average network densities (per km2):



Local network densities

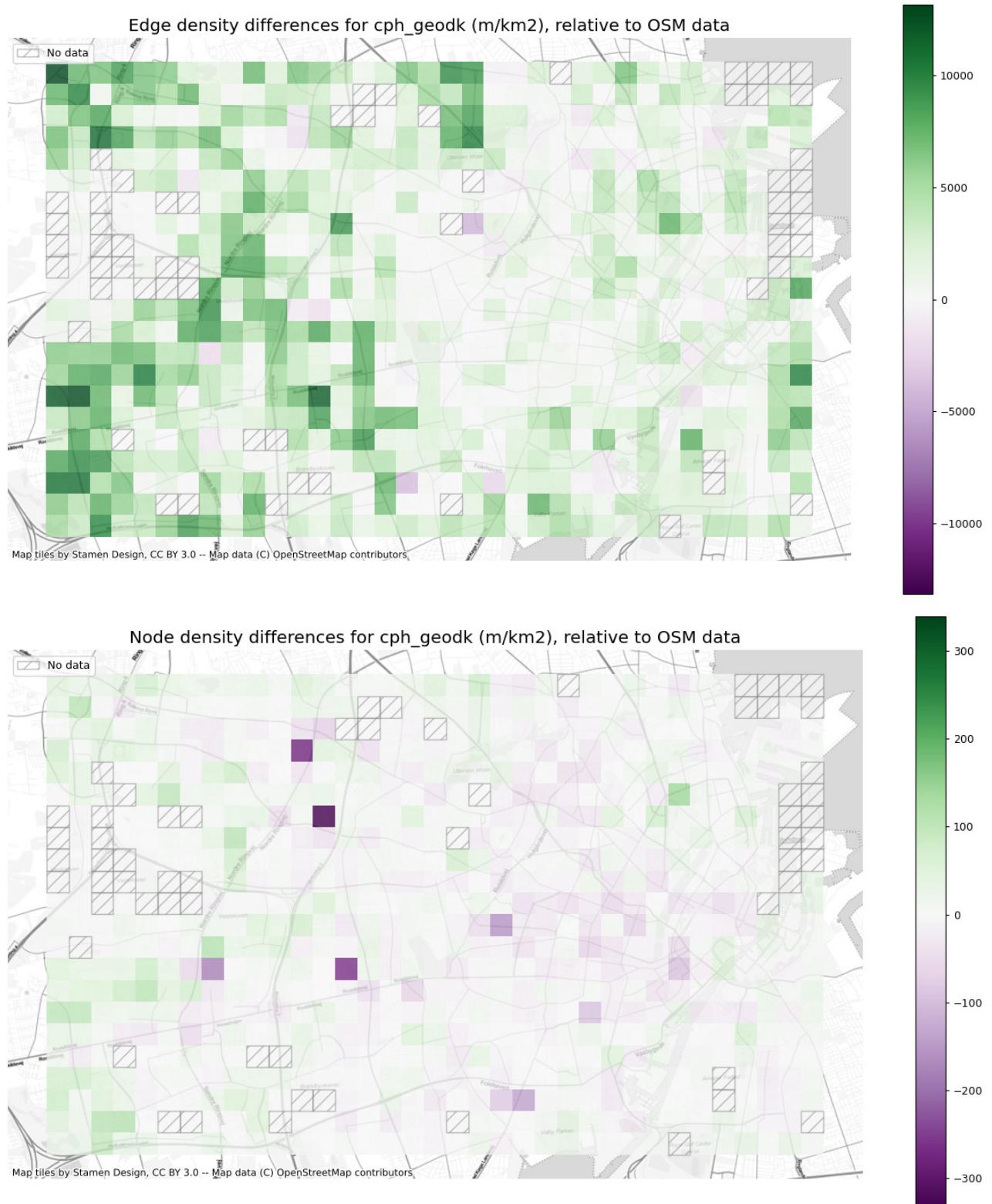
Please note that the plots of densities for respectively the OSM and reference data are produced in the intrinsic notebooks. The plots will thus not necessarily have the same value range for the color bars.

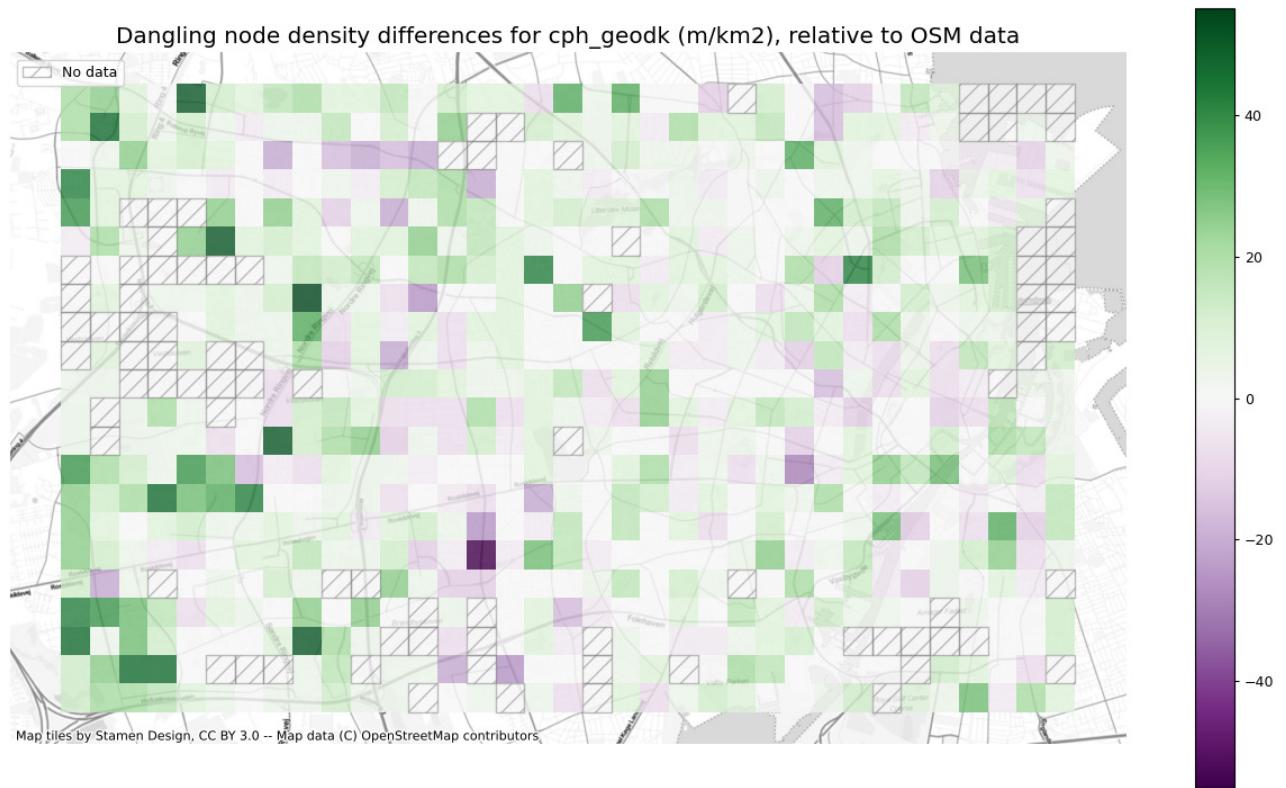




Local differences in network densities

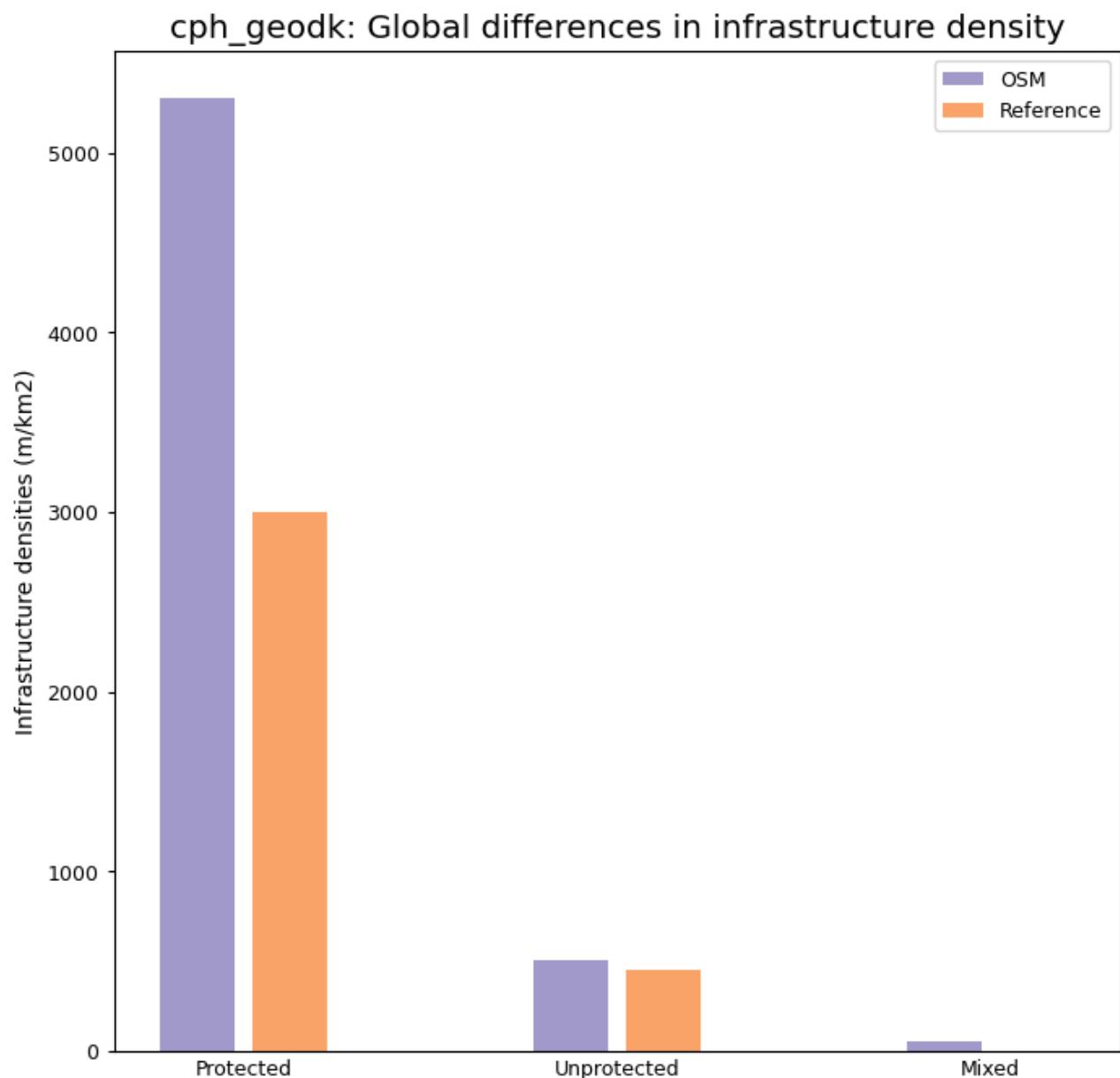
The densities in the OSM data set are taken as base line for comparison. Hence, positive values indicate that the OSM density of the infrastructure type is higher than the reference density; negative values indicate that the OSM density is lower than the reference density.





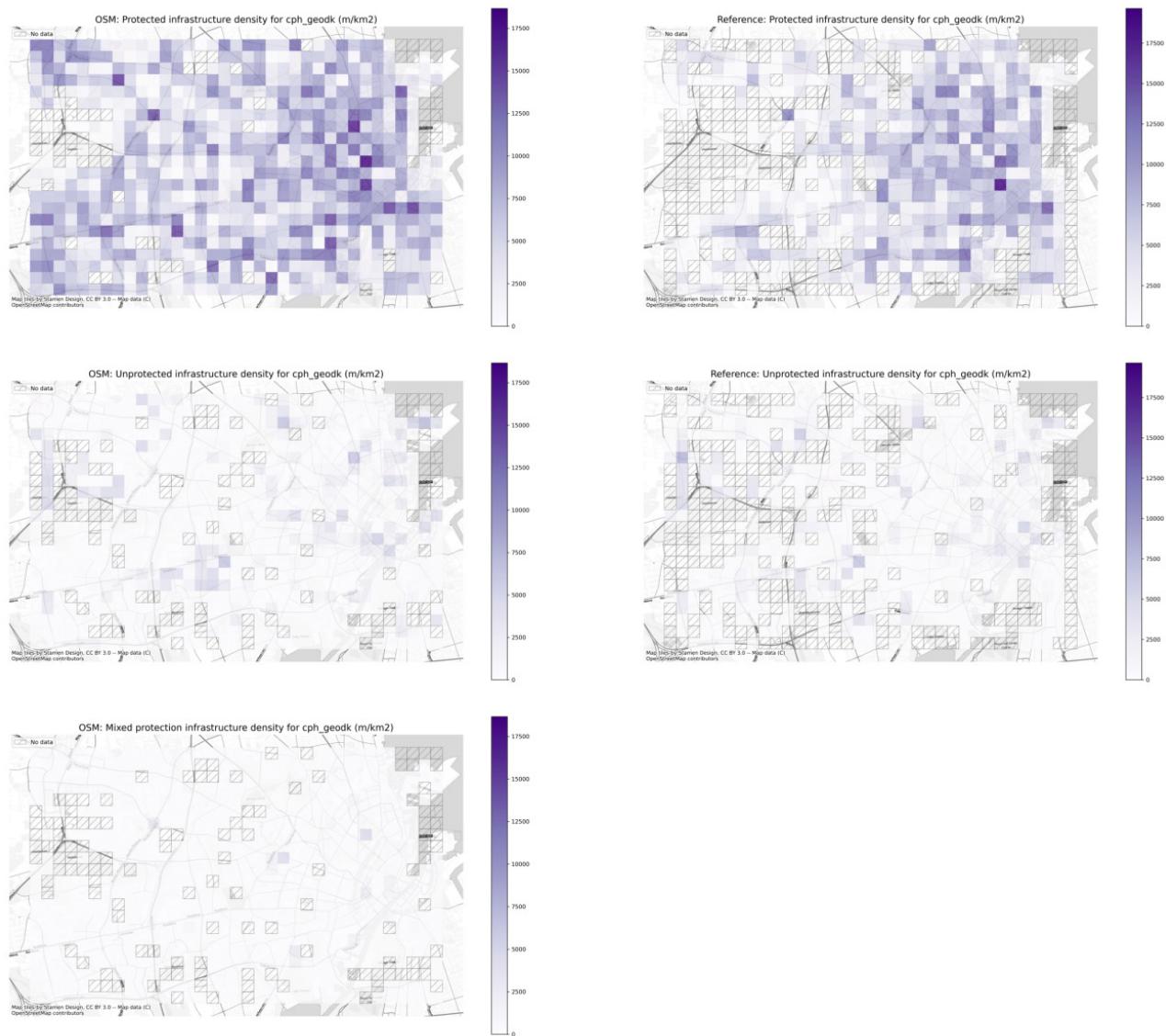
Densities of protected and unprotected bicycle infrastructure

Global network densities for protected/unprotected infrastructure



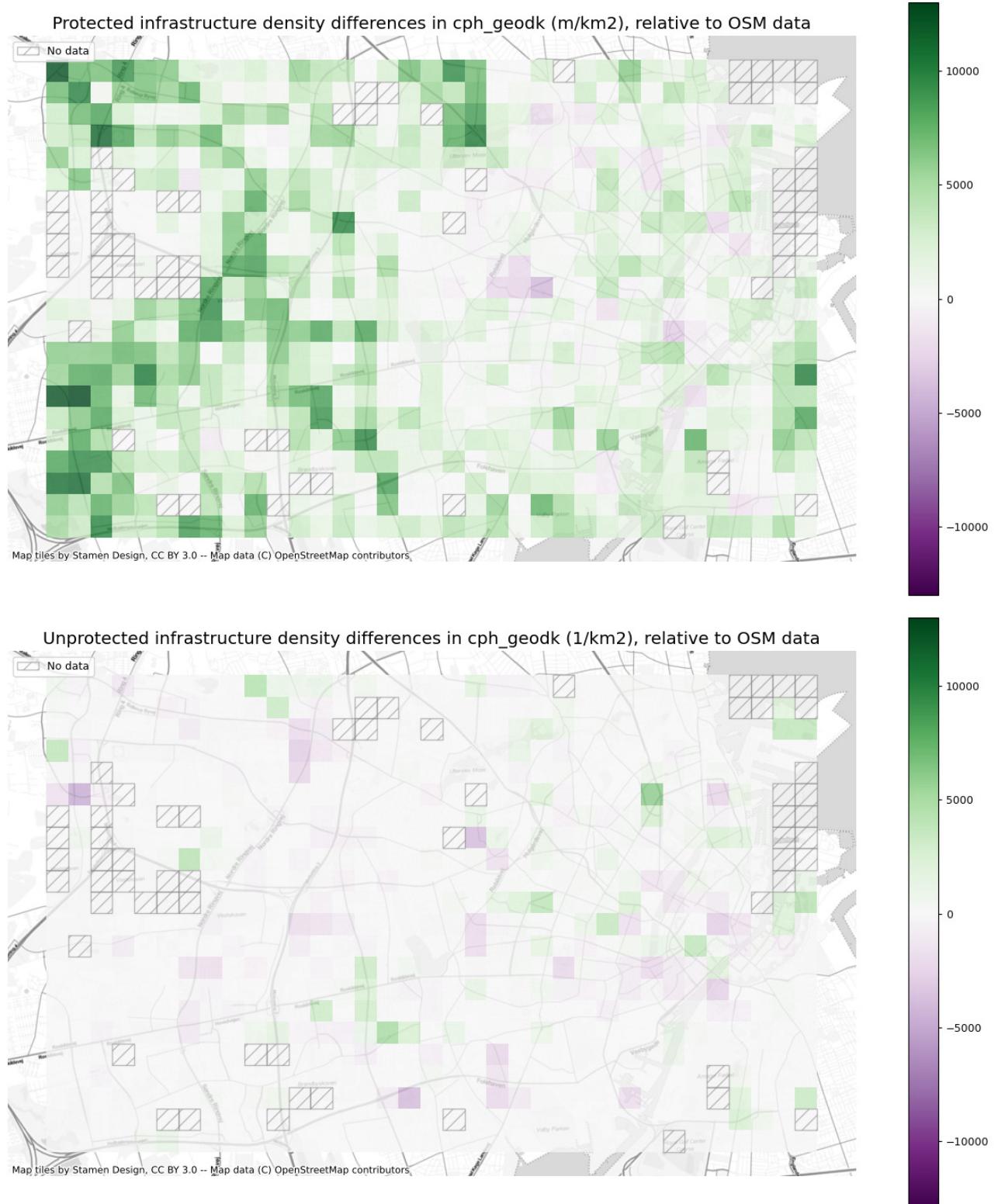
Local network densities for protected/unprotected infrastructure

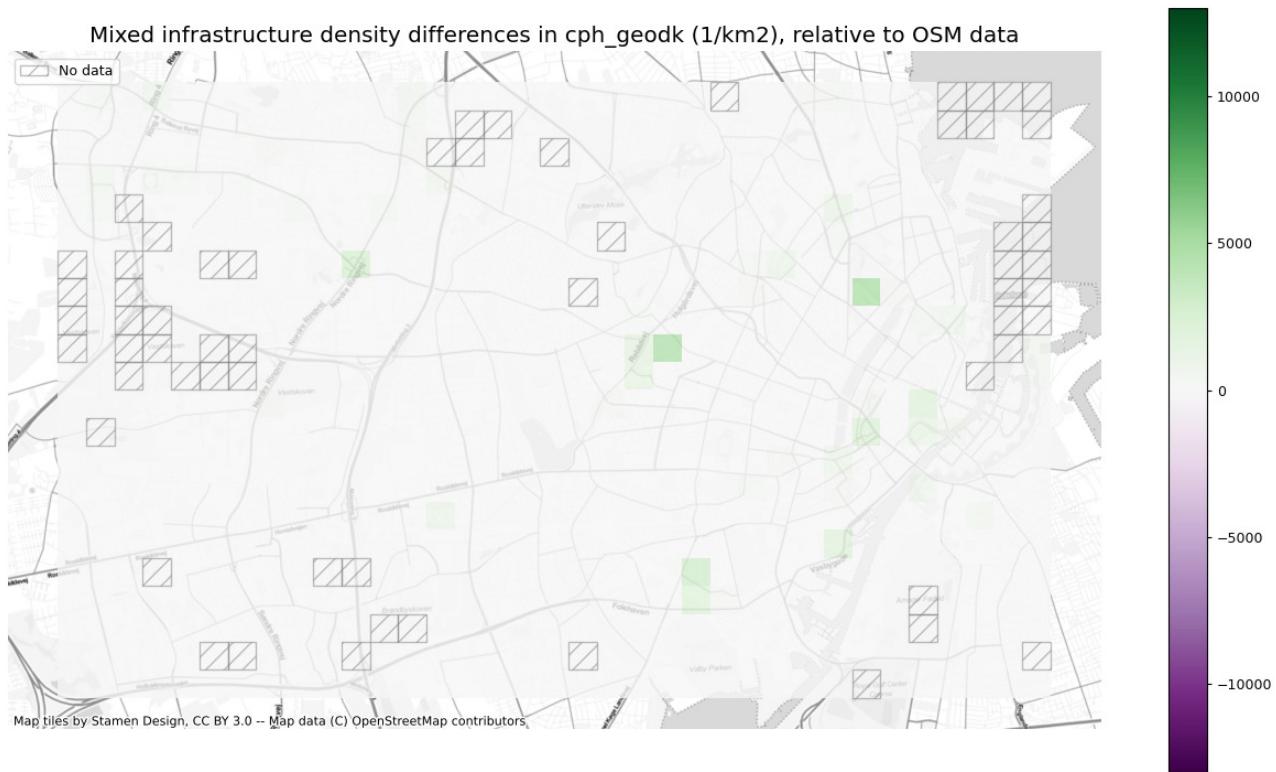
The densities in the OSM data set are taken as base line for comparison. Hence, positive values indicate that the OSM density of the infrastructure type is higher than the reference density; negative values indicate that the OSM density is lower than the reference density.



Differences in infrastructure type density

No infrastructure is mapped as mixed protected/unprotected in the reference data.





Network topology

After having compared data completeness, i.e., *how much* infrastructure is mapped, in the section above, we now will focus on differences in network *topology*, which can give some information about *how* the infrastructure is mapped in both datasets. Here we also analyze the extent to which network edges are connected to one or more other edges, or if they on the other hand end in a dangling node. The extent to which edges are properly connected to adjacent edges are important for, for example, analyses of accessibility and routing.

When working with data on bicycle networks, a dataset without gaps between actually connected network elements is preferred - while of course reflecting the real conditions. Identifying the dangling nodes in a network are a quick and easy way to identify edges that end in a 'dead end'. Under and overshoots offer a more precise picture of respectively network gaps and overextended edges, that give a misleading count of dangling nodes.

Method

To identify potential gaps or missing links in the data, first the dangling nodes in both datasets are plotted. Then, the local percentage of dangling nodes out of all nodes in each dataset is plotted separately followed by a plot showing the local difference in the percent of nodes categorized as dangling.

Under and overshoots in both OSM and reference data are finally plotted together in an interactive plot for further inspection.

Interpretation

If an edge ends in a dangling node in one dataset but not the other, this indicates a problem with the data quality and that there either is a missing connection in the data, or that two edges have been connected erroneously. Similarly, different local rates in the share of dangling nodes indicates differences in how the bicycle networks have been mapped - although differences in data completeness of course should be considered in the interpretation.

Undershoots are clear indications of misleading gaps in network data - although they might also represent actual gaps in bicycle infrastructure. Comparing undershoots in one dataset with another dataset can help identify whether it is a question of data quality or the quality of the actual infrastructure. Vast differences in the presence of undershoots or gaps across intersections might be an indication in differing digitizing strategies, since some approaches will map a bike lane crossing a street as a connected stretch, while others will introduce a gap in the width of the crossing street. While both approaches can be considered correct, datasets created with the former method are more suited for network-based analysis.

Overshoots will often be less consequential for analysis, but a high number of overshoots will introduce false dangling nodes and distort measures for network structure based on e.g., node degree or the ratio between nodes and edges.

Simplifying the OSM network decreased the number of edges by 88.9%.
Simplifying the OSM network decreased the number of nodes by 84.3%.

Simplifying the reference network decreased the number of edges by 91.2%.
Simplifying the reference network decreased the number of nodes by 92.2%.

Node degree distribution

Before the network simplification the OSM graph had:

- 1 node(s) with degree 7
- 3 node(s) with degree 6
- 28 node(s) with degree 5
- 703 node(s) with degree 4
- 1808 node(s) with degree 3
- 27077 node(s) with degree 2
- 1821 node(s) with degree 1

After the network simplification the OSM graph had:

- 1 node(s) with degree 7
- 3 node(s) with degree 6

- 28 node(s) with degree 5
- 703 node(s) with degree 4
- 1808 node(s) with degree 3
- 566 node(s) with degree 2
- 1821 node(s) with degree 1

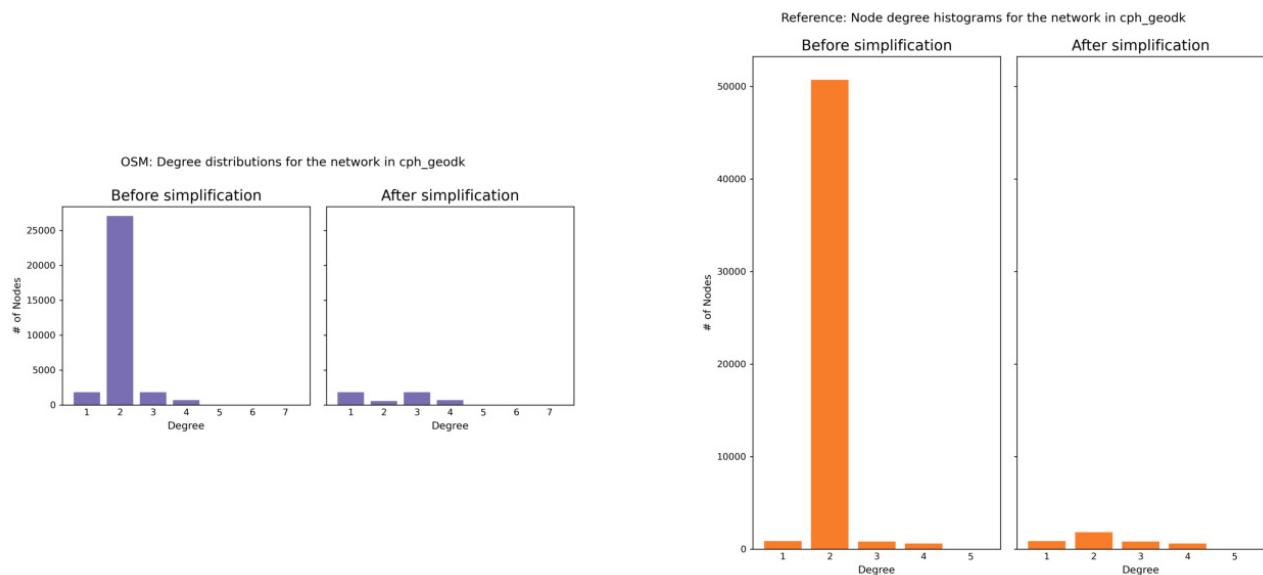
Before the network simplification the reference graph had:

- 7 node(s) with degree 5
- 591 node(s) with degree 4
- 827 node(s) with degree 3
- 50705 node(s) with degree 2
- 870 node(s) with degree 1

After the network simplification the reference graph had:

- 7 node(s) with degree 5
- 591 node(s) with degree 4
- 827 node(s) with degree 3
- 1830 node(s) with degree 2
- 870 node(s) with degree 1

Note that the two figures below have differing y-axis scales.



Alpha, beta & gamma indices

In this subsection, we compute and contrast the three aggregated network metrics alpha, beta, and gamma. These metrics are often used to describe network structure, but as measures of data quality, they are only meaningful when compared to the values of a corresponding dataset. For this reason, alpha, beta, and gamma are only part of the extrinsic analysis and not included in the intrinsic notebooks.

While no conclusion can be made about data quality based on any of the three metrics by itself, a comparison of the metrics for the two data sets can indicate differences in network topology, and hence differences in how the infrastructure has been mapped.

Method

All three indices are computed with `eval_func.compute_alpha_beta_gamma`.

The **alpha** value is the ratio of actual to possible *cycles* in the network. A network cycle is defined as a closed loop - i.e. a path that ends on the same node that it started from. The value of alpha ranges from 0 to 1. An alpha value of 0 means that the network has no cycles at all (e.g. a tree-like structure); an alpha value of 1 means that the network is fully connected, which is very rarely the case.

The **beta** value is the ratio of existing edges to existing nodes in the network. The value of beta ranges from 0 to N-1, where N is the number of existing nodes. A beta value of 0 means that the network has no edges; a beta value of N-1 means that the network is fully connected (see also gamma value of 1). The higher the beta value, the more different paths (on average) can be chosen between any pair of nodes.

The **gamma** value is the ratio of existing to *possible* edges in the network. Any edge that connects two of the existing network nodes is defined as "possible". Hence, the value of gamma ranges from 0 to 1. A gamma value of 0 means that the network has no edges; a gamma value of 1 means that every node of the network is connected to every other node.

For all three indices, see [Ducruet and Rodrigue, 2020](#). All three indices can be interpreted in respect to network connectivity: The higher the alpha value, the more cycles are present in the network; the higher the beta value, the higher the number of paths and thus the higher the complexity of the network; and the higher the gamma value, the fewer edges lie between any pair of nodes.

Interpretation

These metrics do not say much about the data quality itself, nor are they useful for a topological comparison of networks of similar size. However, some conclusions can be made through a comparison. For example, if the indices are very similar for the two networks, despite the networks e.g., having very different geometric lengths, this suggests that the data sets have been mapped in roughly the same way, but than one simply includes more features than the other. On the other hand, if the networks have roughly the same total geometric length, but the values from alpha, beta and gamma differ, this can be an indication that the structure and topology of the two datasets are fundamentally different.

Alpha for the simplified OSM network: 0.11

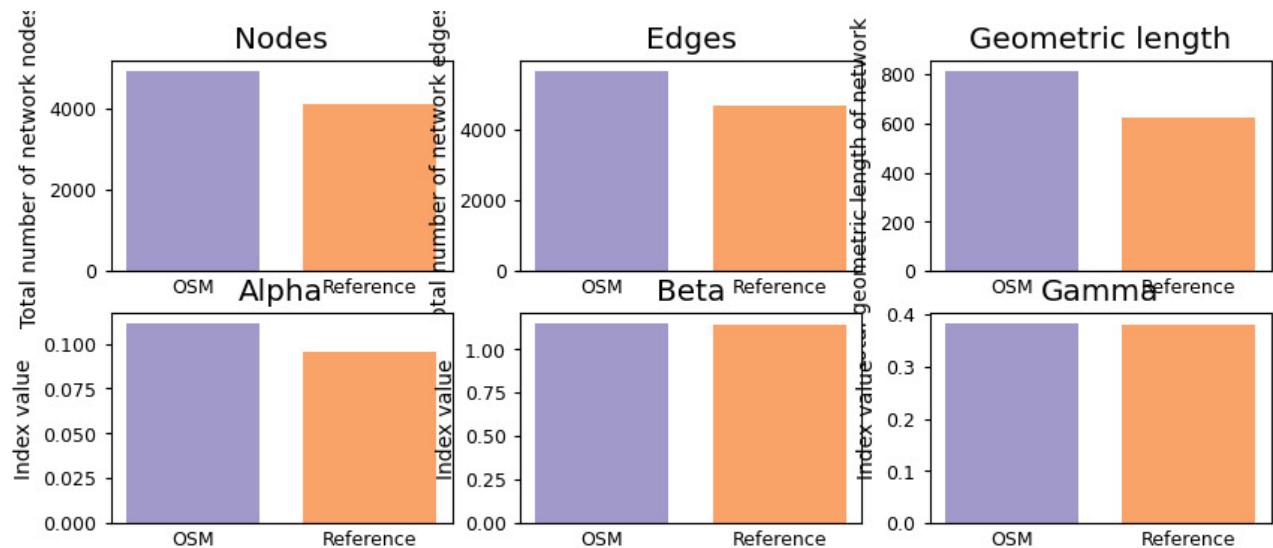
Beta for the simplified OSM network: 1.15

Gamma for the simplified OSM network: 0.38

Alpha for the simplified reference network: 0.10

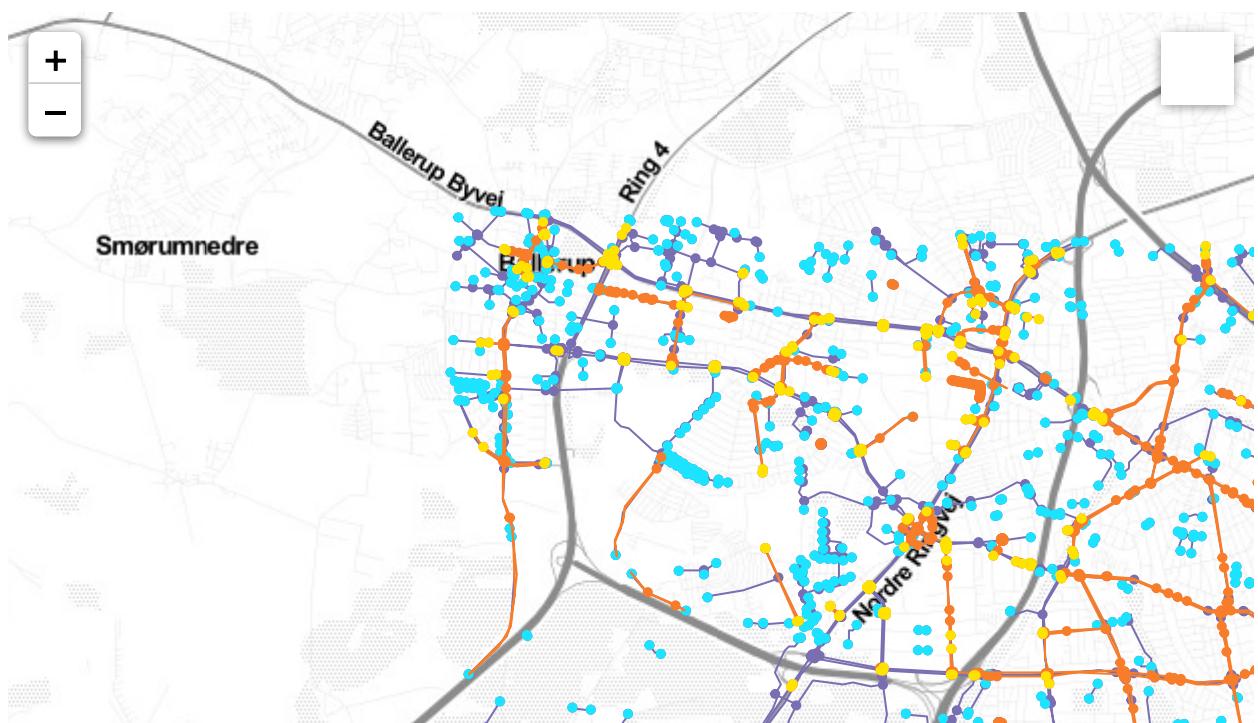
Beta for the simplified reference network: 1.14

Gamma for the simplified reference network: 0.38



Dangling nodes

Dangling nodes in OSM & reference networks

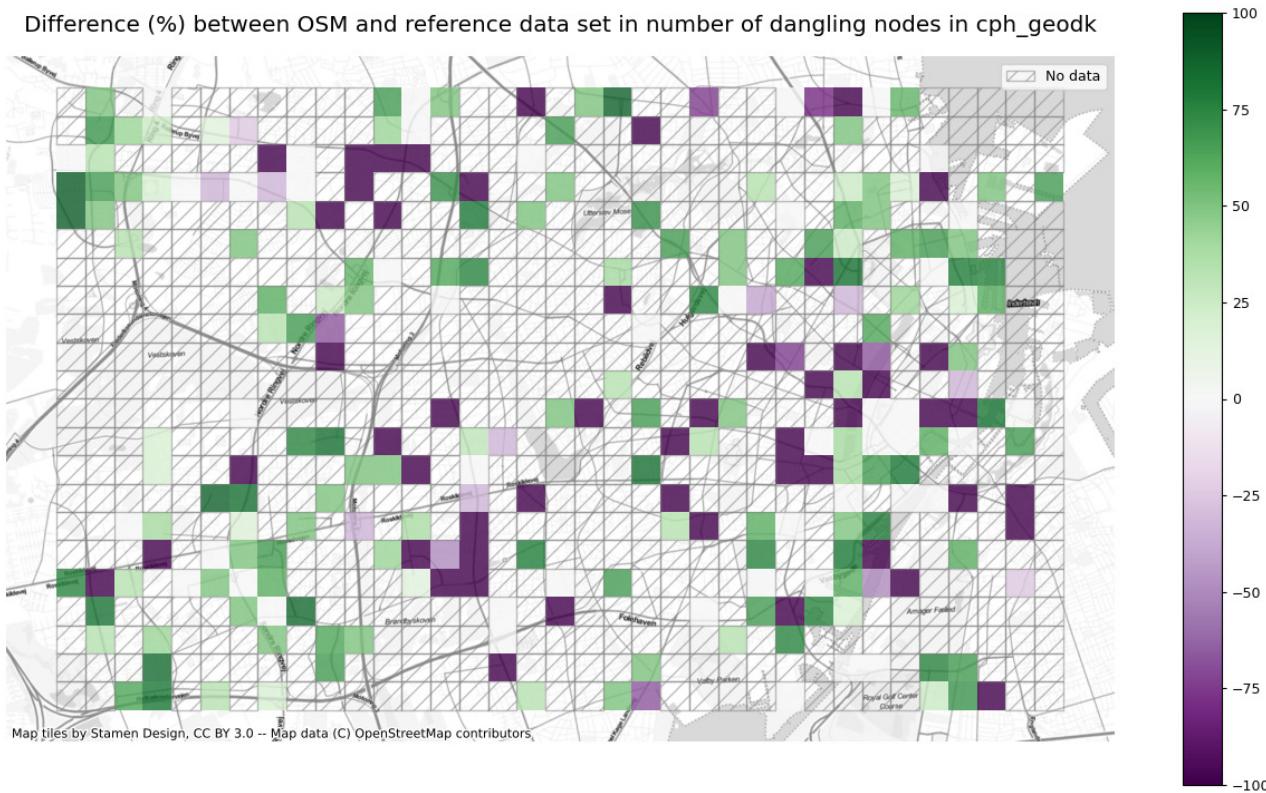


Local values for dangling nodes

Dangling nodes as percentage of all nodes:

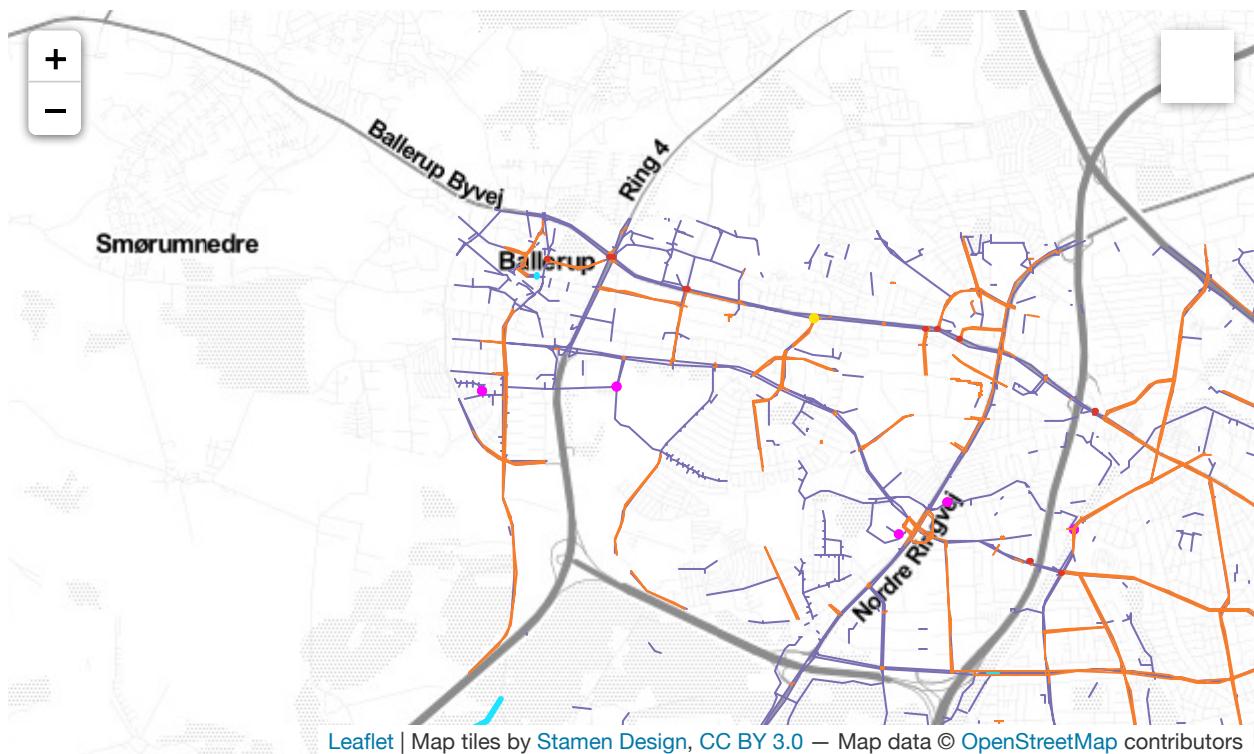


Local differences in percent dangling nodes:



Under/overshoots

Over and undershoots in OSM and reference networks



Network components

This section takes a close look at the network component characteristics for each of the two data sets.

Disconnected components do not share any elements (nodes/edges). In other words, there is no network path that could lead from one disconnected component to the other. As mentioned above, most real-world networks of bicycle infrastructure do consist of many disconnected components ([Natera Orozco et al., 2020](#)). However, when two disconnected components are very close to each other, it might be a sign of a missing edge or another digitizing error.

Method

To compare the number and pattern of disconnected components in OSM and reference data, all component results from the intrinsic analyses are juxtaposed and two new plots showing respectively components gaps for OSM and reference data and the difference in component connectivity are produced.

Interpretation

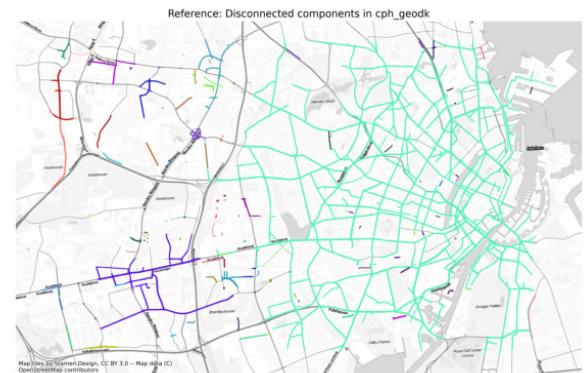
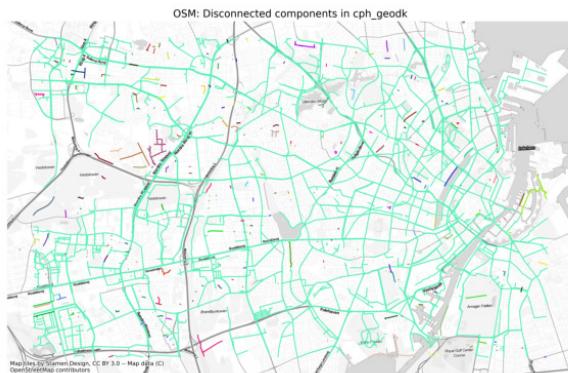
The fragmented nature of many bicycle networks make it hard to assess whether disconnected components are a question of a lack of data quality or a lack of properly connected bicycle infrastructure. Comparing disconnected components in two datasets enables a more accurate assessment of whether a disconnected component is a data or a planning issue.

Subsections

- Disconnected components
- Component size distribution
- Largest connected component
- Potentially missing links
- Number of components per grid cell
- Component connectivity

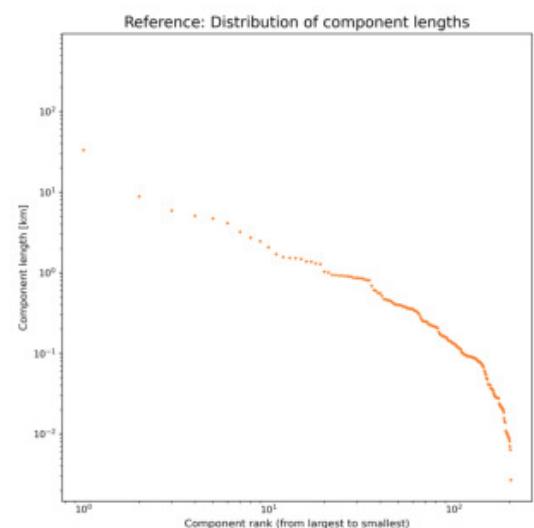
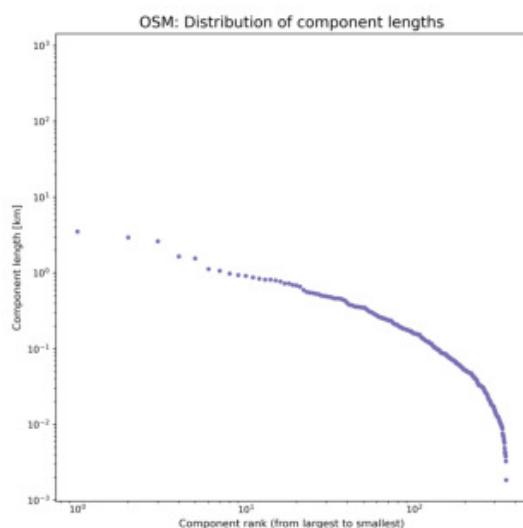
Disconnected components

The OSM network in the study area consists of 352 disconnected components. The reference network in the study area consists of 204 disconnected components.



Component size distribution

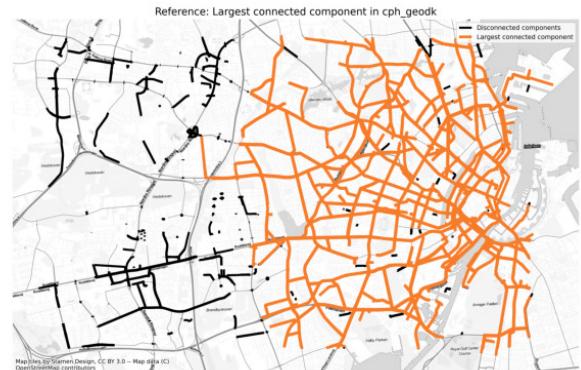
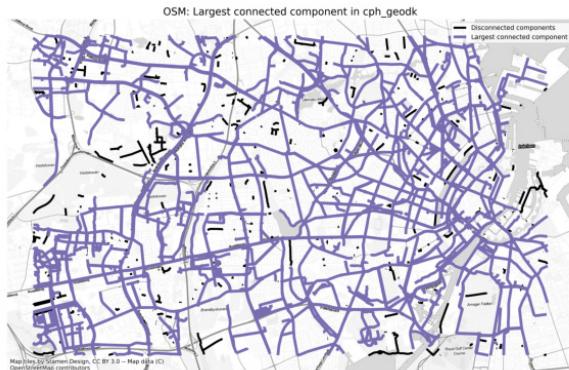
TODO: add explanation of plot (after plot has been redone)



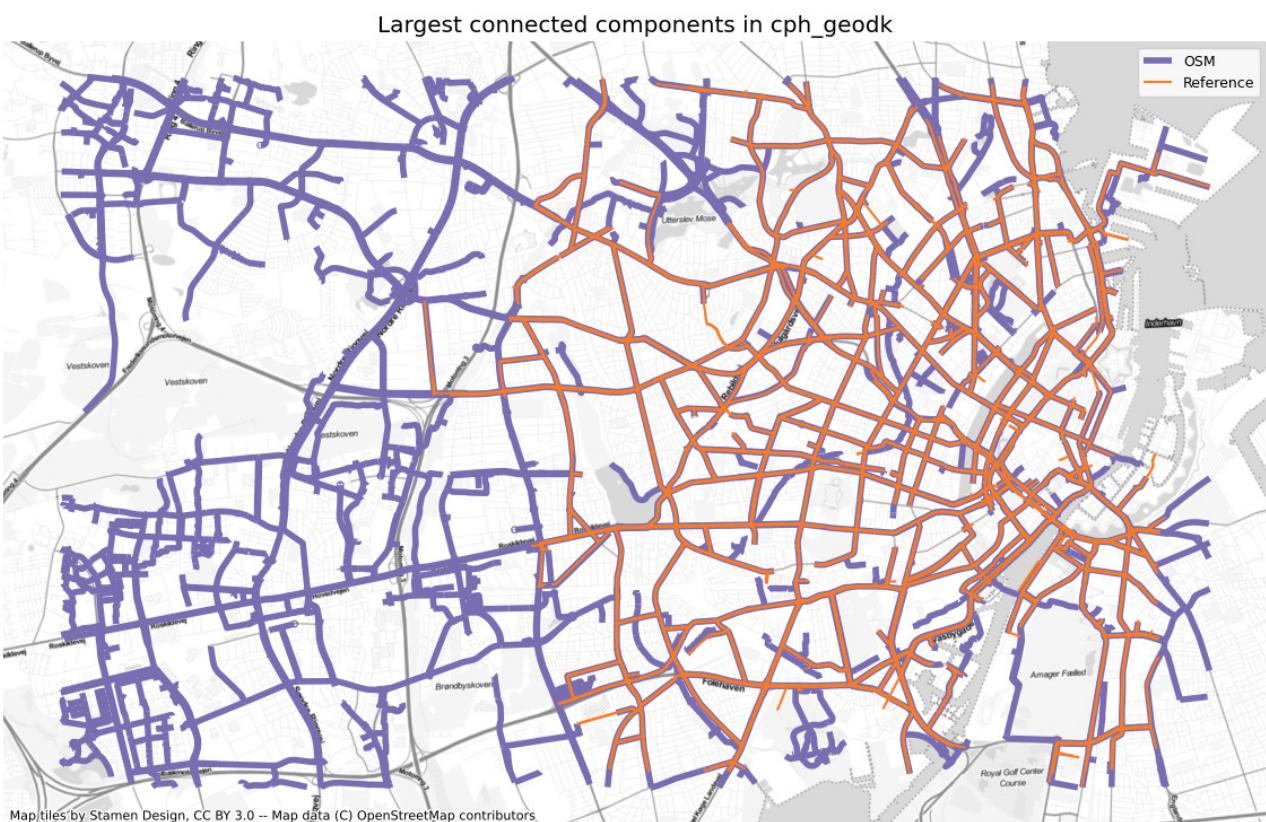
Largest connected component

The largest connected component in the OSM network contains 92.27% of the network length.

The largest connected component in the reference network contains 80.04% of the network length.

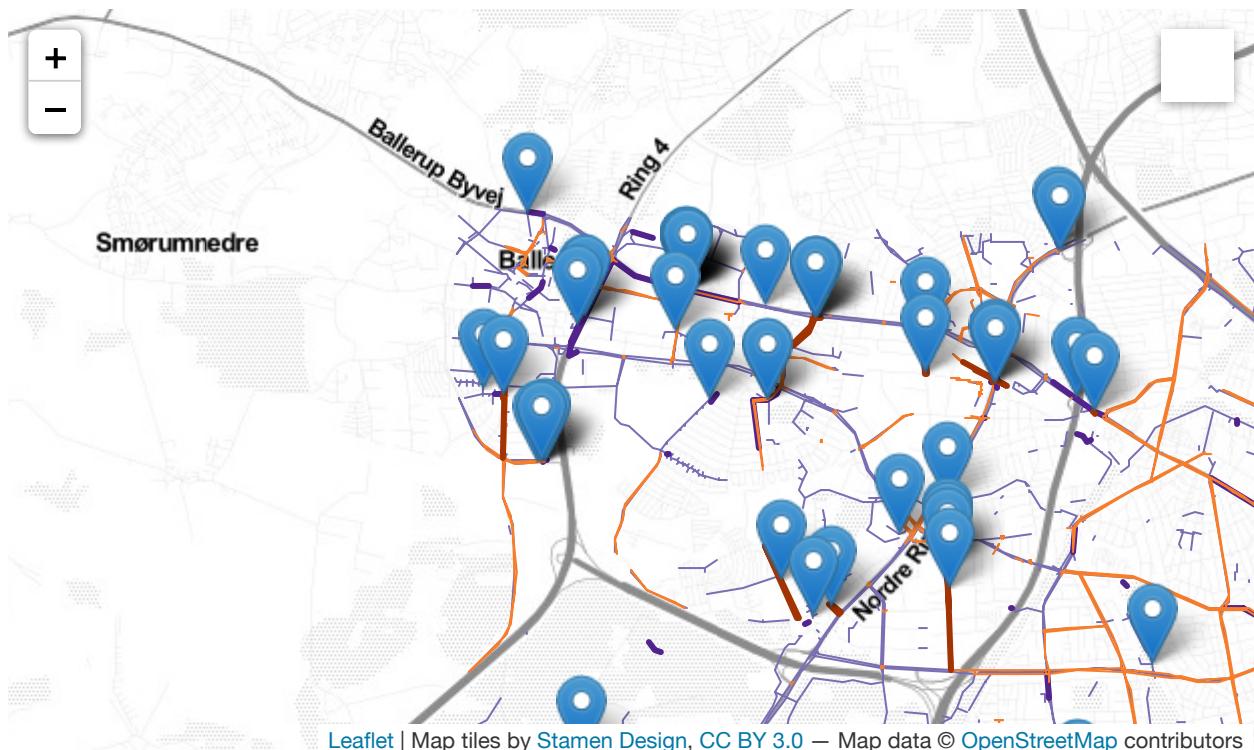


Overlay of largest connected component in OSM and reference networks



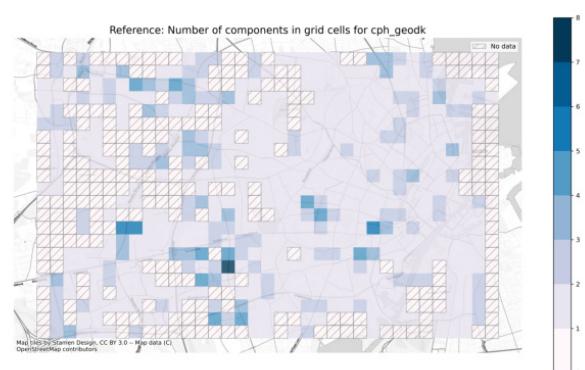
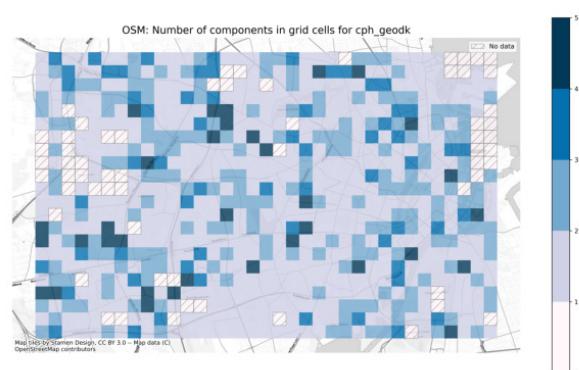
Potentially missing links

In the plot of potential missing links between components, all edges that are within the specified distance of an edge on another component are plotted. The gaps between disconnected edges are highlighted with a marker. The map thus highlights edges which, despite being in close proximity of each other are disconnected and where it thus would not be



Number of components per grid cell

The plots below show the number of components intersecting a grid cell. A high number of components in a grid cell is generally an indication of poor network connectivity - either due to fragmented infrastructure or because of problems with the data quality.



Component connectivity

Visualizing differences between how many cells can be reached from each cell. In the plot showing the relative difference in percent cells reached, positive values means that more cells can be reached from this particular cell using the OSM data, while negative values indicate a higher connectivity using the reference data set.

The metric is a crude measure for network connectivity but has the benefit of being computationally cheap and thus able to highlight stark differences in network connectivity in very little time.



Summary

```

KeyError                                     Traceback (most recent call last)
Cell In [93], line 17
  9 ref_intrinsic_df = pd.read_csv(
10     ref_results_data_fp + "intrinsic_summary_results.csv",
11     index_col=0,
12     names=["Reference"],
13     header=0,
14 )
15 # Drop rows from OSM results not available for reference
--> 17 osm_intrinsic_df.drop([
18     "Count of incompatible tag combinations", "Count of missing intersection nodes"],
19     axis=0,
20     inplace=True,
21 )
22 # Save new results
24 osm_intrinsic_df.at["Alpha", "OSM"] = osm_alpha

File ~/anaconda3/anaconda3/envs/bikedna/lib/python3.11/site-packages/pandas/util/_decorators.py:331, in deprecate_nonkeyword_arguments.<locals>.decorate.<locals>.wrapper(*args, **kwargs)
  325 if len(args) > num_allow_args:
  326     warnings.warn(
  327         msg.format(arguments=_format_argument_list(allow_args)),
  328         FutureWarning,
  329         stacklevel=find_stack_level(),
  330     )
--> 331 return func(*args, **kwargs)

File ~/anaconda3/anaconda3/envs/bikedna/lib/python3.11/site-packages/pandas/core/frame.py:5396, in DataFrame.drop(self, labels, axis, index, columns, level, inplace, errors)
  5248 @deprecate_nonkeyword_arguments(version=None, allowed_args=["self", "labels"])
  5249 def drop(  # type: ignore[override]
  5250     self,
  5251     ...,
  5252     errors: IgnoreRaise = "raise",
  5253 ) -> DataFrame | None:
  5254     """
  5255     Drop specified labels from rows or columns.
  5256
  5257     Parameters
  5258     ----------
  5259     weight : float, optional
  5260         weight 1.0      0.8
  5261
  5262     Returns
  5263     -------
  5264     DataFrame
  5265
  5266     See Also
  5267     --------
  5268     DataFrame.dropna
  5269
  5270     Examples
  5271     -----
  5272     drop(['a', 'c'])
  5273
  5274     drop(['a', 'c'], axis=1)
  5275
  5276     drop(['a', 'c'], axis=0)
  5277
  5278     drop(['a', 'c'], axis=1, labels=['a', 'c'])
  5279
  5280     drop(['a', 'c'], axis=0, labels=['a', 'c'])
  5281
  5282     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8])
  5283
  5284     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8])
  5285
  5286     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise')
  5287
  5288     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise')
  5289
  5290     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='ignore')
  5291
  5292     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='ignore')
  5293
  5294     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5295
  5296     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5297
  5298     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='ignore', errors='raise')
  5299
  5300     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='ignore', errors='raise')
  5301
  5302     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5303
  5304     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5305
  5306     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5307
  5308     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5309
  5310     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5311
  5312     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5313
  5314     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5315
  5316     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5317
  5318     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5319
  5320     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5321
  5322     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5323
  5324     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5325
  5326     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5327
  5328     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5329
  5330     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5331
  5332     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5333
  5334     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5335
  5336     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5337
  5338     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5339
  5340     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5341
  5342     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5343
  5344     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5345
  5346     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5347
  5348     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5349
  5350     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5351
  5352     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5353
  5354     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5355
  5356     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5357
  5358     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5359
  5360     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5361
  5362     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5363
  5364     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5365
  5366     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5367
  5368     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5369
  5370     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5371
  5372     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5373
  5374     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5375
  5376     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5377
  5378     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5379
  5380     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5381
  5382     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5383
  5384     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5385
  5386     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5387
  5388     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5389
  5390     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5391
  5392     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5393
  5394     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5395
  5396     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5397
  5398     drop(['a', 'c'], axis=1, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5399
  5400     drop(['a', 'c'], axis=0, labels=['a', 'c'], weights=[1.0, 0.8], errors='raise', errors='raise')
  5401
  5402
  5403
  5404
  5405
  5406
  5407
  5408
  5409
  5410
  5411
  5412
  5413
  5414
  5415
  5416
  5417
  5418
  5419
  5420
  5421
  5422
  5423
  5424
  5425
  5426
  5427
  5428
  5429
  5430
  5431
  5432
  5433
  5434
  5435
  5436
  5437
  5438
  5439
  5440
  5441
  5442
  5443
  5444
  5445
  5446
  5447
  5448
  5449
  5450
  5451
  5452
  5453
  5454
  5455
  5456
  5457
  5458
  5459
  5460
  5461
  5462
  5463
  5464
  5465
  5466
  5467
  5468
  5469
  5470
  5471
  5472
  5473
  5474
  5475
  5476
  5477
  5478
  5479
  5480
  5481
  5482
  5483
  5484
  5485
  5486
  5487
  5488
  5489
  5490
  5491
  5492
  5493
  5494
  5495
  5496
  5497
  5498
  5499
  5500
  5501
  5502
  5503
  5504
  5505
  5506
  5507
  5508
  5509
  5510
  5511
  5512
  5513
  5514
  5515
  5516
  5517
  5518
  5519
  5520
  5521
  5522
  5523
  5524
  5525
  5526
  5527
  5528
  5529
  5530
  5531
  5532
  5533
  5534
  5535
  5536
  5537
  5538
  5539
  5540
  5541
  5542
  5543
  5544
  5545
  5546
  5547
  5548
  5549
  5550
  5551
  5552
  5553
  5554
  5555
  5556
  5557
  5558
  5559
  5560
  5561
  5562
  5563
  5564
  5565
  5566
  5567
  5568
  5569
  5570
  5571
  5572
  5573
  5574
  5575
  5576
  5577
  5578
  5579
  5580
  5581
  5582
  5583
  5584
  5585
  5586
  5587
  5588
  5589
  5590
  5591
  5592
  5593
  5594
  5595
  5596
  5597
  5598
  5599
  5600
  5601
  5602
  5603
  5604
  5605
  5606
  5607
  5608
  5609
  5610
  5611
  5612
  5613
  5614
  5615
  5616
  5617
  5618
  5619
  5620
  5621
  5622
  5623
  5624
  5625
  5626
  5627
  5628
  5629
  5630
  5631
  5632
  5633
  5634
  5635
  5636
  5637
  5638
  5639
  5640
  5641
  5642
  5643
  5644
  5645
  5646
  5647
  5648
  5649
  5650
  5651
  5652
  5653
  5654
  5655
  5656
  5657
  5658
  5659
  5660
  5661
  5662
  5663
  5664
  5665
  5666
  5667
  5668
  5669
  5670
  5671
  5672
  5673
  5674
  5675
  5676
  5677
  5678
  5679
  5680
  5681
  5682
  5683
  5684
  5685
  5686
  5687
  5688
  5689
  5690
  5691
  5692
  5693
  5694
  5695
  5696
  5697
  5698
  5699
  5700
  5701
  5702
  5703
  5704
  5705
  5706
  5707
  5708
  5709
  5710
  5711
  5712
  5713
  5714
  5715
  5716
  5717
  5718
  5719
  5720
  5721
  5722
  5723
  5724
  5725
  5726
  5727
  5728
  5729
  5730
  5731
  5732
  5733
  5734
  5735
  5736
  5737
  5738
  5739
  5740
  5741
  5742
  5743
  5744
  5745
  5746
  5747
  5748
  5749
  5750
  5751
  5752
  5753
  5754
  5755
  5756
  5757
  5758
  5759
  5760
  5761
  5762
  5763
  5764
  5765
  5766
  5767
  5768
  5769
  5770
  5771
  5772
  5773
  5774
  5775
  5776
  5777
  5778
  5779
  5780
  5781
  5782
  5783
  5784
  5785
  5786
  5787
  5788
  5789
  5790
  5791
  5792
  5793
  5794
  5795
  5796
  5797
  5798
  5799
  5800
  5801
  5802
  5803
  5804
  5805
  5806
  5807
  5808
  5809
  5810
  5811
  5812
  5813
  5814
  5815
  5816
  5817
  5818
  5819
  5820
  5821
  5822
  5823
  5824
  5825
  5826
  5827
  5828
  5829
  5830
  5831
  5832
  5833
  5834
  5835
  5836
  5837
  5838
  5839
  5840
  5841
  5842
  5843
  5844
  5845
  5846
  5847
  5848
  5849
  5850
  5851
  5852
  5853
  5854
  5855
  5856
  5857
  5858
  5859
  5860
  5861
  5862
  5863
  5864
  5865
  5866
  5867
  5868
  5869
  5870
  5871
  5872
  5873
  5874
  5875
  5876
  5877
  5878
  5879
  5880
  5881
  5882
  5883
  5884
  5885
  5886
  5887
  5888
  5889
  5890
  5891
  5892
  5893
  5894
  5895
  5896
  5897
  5898
  5899
  5900
  5901
  5902
  5903
  5904
  5905
  5906
  5907
  5908
  5909
  5910
  5911
  5912
  5913
  5914
  5915
  5916
  5917
  5918
  5919
  5920
  5921
  5922
  5923
  5924
  5925
  5926
  5927
  5928
  5929
  5930
  5931
  5932
  5933
  5934
  5935
  5936
  5937
  5938
  5939
  5940
  5941
  5942
  5943
  5944
  5945
  5946
  5947
  5948
  5949
  5950
  5951
  5952
  5953
  5954
  5955
  5956
  5957
  5958
  5959
  5960
  5961
  5962
  5963
  5964
  5965
  5966
  5967
  5968
  5969
  5970
  5971
  5972
  5973
  5974
  5975
  5976
  5977
  5978
  5979
  5980
  5981
  5982
  5983
  5984
  5985
  5986
  5987
  5988
  5989
  5990
  5991
  5992
  5993
  5994
  5995
  5996
  5997
  5998
  5999
  6000
  6001
  6002
  6003
  6004
  6005
  6006
  6007
  6008
  6009
  6010
  6011
  6012
  6013
  6014
  6015
  6016
  6017
  6018
  6019
  6020
  6021
  6022
  6023
  6024
  6025
  6026
  6027
  6028
  6029
  6030
  6031
  6032
  6033
  6034
  6035
  6036
  6037
  6038
  6039
  6040
  6041
  6042
  6043
  6044
  6045
  6046
  6047
  6048
  6049
  6050
  6051
  6052
  6053
  6054
  6055
  6056
  6057
  6058
  6059
  6060
  6061
  6062
  6063
  6064
  6065
  6066
  6067
  6068
  6069
  6070
  6071
  6072
  6073
  6074
  6075
  6076
  6077
  6078
  6079
  6080
  6081
  6082
  6083
  6084
  6085
  6086
  6087
  6088
  6089
  6090
  6091
  6092
  6093
  6094
  6095
  6096
  6097
  6098
  6099
  6100
  6101
  6102
  6103
  6104
  6105
  6106
  6107
  6108
  6109
  6110
  6111
  6112
  6113
  6114
  6115
  6116
  6117
  6118
  6119
  6120
  6121
  6122
  6123
  6124
  6125
  6126
  6127
  6128
  6129
  6130
  6131
  6132
  6133
  6134
  6135
  6136
  6137
  6138
  6139
  6140
  6141
  6142
  6143
  6144
  6145
  6146
  6147
  6148
  6149
  6150
  6151
  6152
  6153
  6154
  6155
  6156
  6157
  6158
  6159
  6160
  6161
  6162
  6163
  6164
  6165
  6166
  6167
  6168
  6169
  6170
  6171
  6172
  6173
  6174
  6175
  6176
  6177
  6178
  6179
  6180
  6181
  6182
  6183
  6184
  6185
  6186
  6187
  6188
  6189
  6190
  6191
  6192
  6193
  6194
  6195
  6196
  6197
  6198
  6199
  6200
  6201
  6202
  6203
  6204
  6205
  6206
  6207
  6208
  6209
  6210
  6211
  6212
  6213
  6214
  6215
  6216
  6217
  6218
  6219
  6220
  6221
  6222
  6223
  6224
  6225
  6226
  6227
  6228
  6229
  6230
  6231
  6232
  6233
  6234
  6235
  6236
  6237
  6238
  6239
  6240
  6241
  6242
  6243
  6244
  6245
  6246
  6247
  6248
  6249
  6250
  6251
  6252
  6253
  6254
  6255
  6256
  6257
  6258
  6259
  6260
  6261
  6262
  6263
  6264
  6265
  6266
  6267
  6268
  6269
  6270
  6271
  6272
  6273
  6274
  6275
  6276
  6277
  6278
  6279
  6280
  6281
  6282
  6283
  6284
  6285
  6286
  6287
  6288
  6289
  6290
  6291
  6292
  6293
  6294
  6295
  6296
  6297
  6298
  6299
  6300
  6301
  6302
  6303
  6304
  6305
  6306
  6307
  6308
  6309
  6310
  6311
  6312
  6313
  6314
  6315
  6316
  6317
  6318
  6319
  6320
  6321
  6322
  6323
  6324
  6325
  6326
  6327
  6328
  6329
  6330
  6331
  6332
  6333
  6334
  6335
  6336
  6337
  6338
  6339
  6340
  6341
  6342
  6343
  6344
  6345
  6346
  6347
  6348
  6349
  6350
  6351
  6352
  6353
  6354
  6355
  6356
  6357
  6358
  6359
  6360
  6361
  6362
  6363
  6364
  6365
  6366
  6367
  6368
  6369
  6370
  6371
  6372
  6373
  6374
  6375
  6376
  6377
  6378
  6379
  6380
  6381
  6382
  6383
  6384
  6385
  6386
  6387
  6388
  6389
  6390
  6391
  6392
  6393
  6394
  6395
  6396
  6397
  6398
  6399
  6400
  6401
  6402
  6403
  6404
  6405
  6406
  6407
  6408
  6409
  6410
  6411
  6412
  6413
  6414
  6415
  6416
  6417
  6418
  6419
  6420
  6421
  6422
  6423
  6424
  6425
  6426
  6427
  6428
  6429
  6430
  6431
  6432
  6433
  6434
  6435
  6436
  6437
  6438
  6439
  6440
  6441
  6442
  6443
  6444

```

```

5402         inplace=inplace,
5403         errors=errors,
5404     )
File ~/anaconda3/anaconda3/envs/bikedna/lib/python3.11/site-packages/pandas/util/
_decorators.py:331, in deprecate_nonkeyword_arguments.<locals>.decorate.<locals>.
wrapper(*args, **kwargs)
325 if len(args) > num_allow_args:
326     warnings.warn(
327         msg.format(arguments=_format_argument_list(allow_args)),
328         FutureWarning,
329         stacklevel=find_stack_level(),
330     )
--> 331 return func(*args, **kwargs)

File ~/anaconda3/anaconda3/envs/bikedna/lib/python3.11/site-packages/pandas/core/
generic.py:4505, in NDFrame.drop(self, labels, axis, index, columns, level, inpla
ce, errors)
4503 for axis, labels in axes.items():
4504     if labels is not None:
-> 4505         obj = obj._drop_axis(labels, axis, level=level, errors=errors)
4507 if inplace:
4508     self._update_inplace(obj)

File ~/anaconda3/anaconda3/envs/bikedna/lib/python3.11/site-packages/pandas/core/
generic.py:4546, in NDFrame._drop_axis(self, labels, axis, level, errors, only_sl
ice)
4544     new_axis = axis.drop(labels, level=level, errors=errors)
4545 else:
-> 4546     new_axis = axis.drop(labels, errors=errors)
4547     indexer = axis.get_indexer(new_axis)
4549 # Case for non-unique axis
4550 else:

File ~/anaconda3/anaconda3/envs/bikedna/lib/python3.11/site-packages/pandas/core/
indexes/base.py:6977, in Index.drop(self, labels, errors)
6975 if mask.any():
6976     if errors != "ignore":
-> 6977         raise KeyError(f"{list(labels[mask])} not found in axis")
6978     indexer = indexer[~mask]
6979 return self.delete(indexer)

KeyError: "['Count of incompatible tag combinations', 'Count of missing intersect
ion nodes'] not found in axis"

```

3b. Feature Matching

The feature matching takes its starting point in the reference data and attempts to identify corresponding features in the OSM data set. Feature matching is a necessary precondition to compare single features rather than feature characteristics on study area a grid cell level, as well as for merging two data sets.

Method

Matching features in two road datasets with each their way of digitizing features and a potential one-to-many relationship between edges (for example in the case where one data set only maps road center lines, while the other map the geometries of each bike lane) is not a trivial task.

The method used here converts all network edges to smaller segments of a uniform length before looking for a potential match between the reference and the OSM data. The matching is done on the basis of the buffered distance between objects, the angle, and the undirected Hausdorff distance, and is based on the work of [Koukoletsos et al. \(2012\)](#) and [Will \(2014\)](#).

Based on the matching results, the following is computed:

- The number and length of matched and unmatched edges, in total and per grid cell
- A comparison of the attributes of the matched edges (is their classification of cycling infrastructure as protected or unprotected the same?)

Interpretation

Once the feature matching is complete, it is important to visually explore the results, since the success rate of the matching influences how the analysis of number of matches should be interpreted.

If the features in the two data sets have been digitized in very different ways - e.g. if one data set has digitized bike tracks as mostly straight lines, while the other includes more winding tracks, the matching will fail. This is also the case if they are placed too far from each other. If it visually can be confirmed that the same features do exist in both data sets, a lack of matches indicates that the geometries in the two data sets are too different. If it on the other hand can be confirmed that most real corresponding features have been identified, a lack of matches in an area indicates errors of commission or omission.

Sections

- [Match features](#)

- Run feature matching
 - Plot results: matched vs. unmatched features
 - [Analyze feature matching results](#)
 - Matched features by infrastructure type
 - Local success rate of feature matching
 - [Summary](#)
 - Save results
-

Match Features

Run feature matching

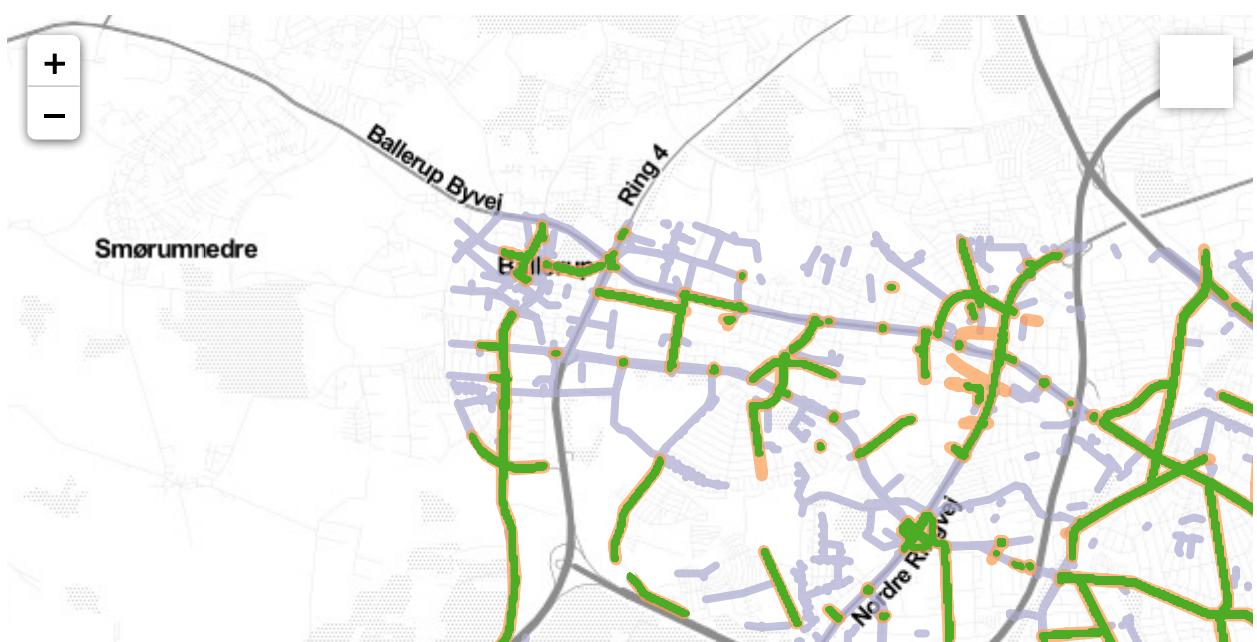
Segments created successfully!

Starting matching process using a buffer distance of 15 meters, a Hausdorff distance of 17 meters, and a max angle of 30 degrees.

Buffer matches found! Continuing with final matching process...

61040 reference segments were matched to OSM edges
2655 reference segments were not matched
Feature matching completed!

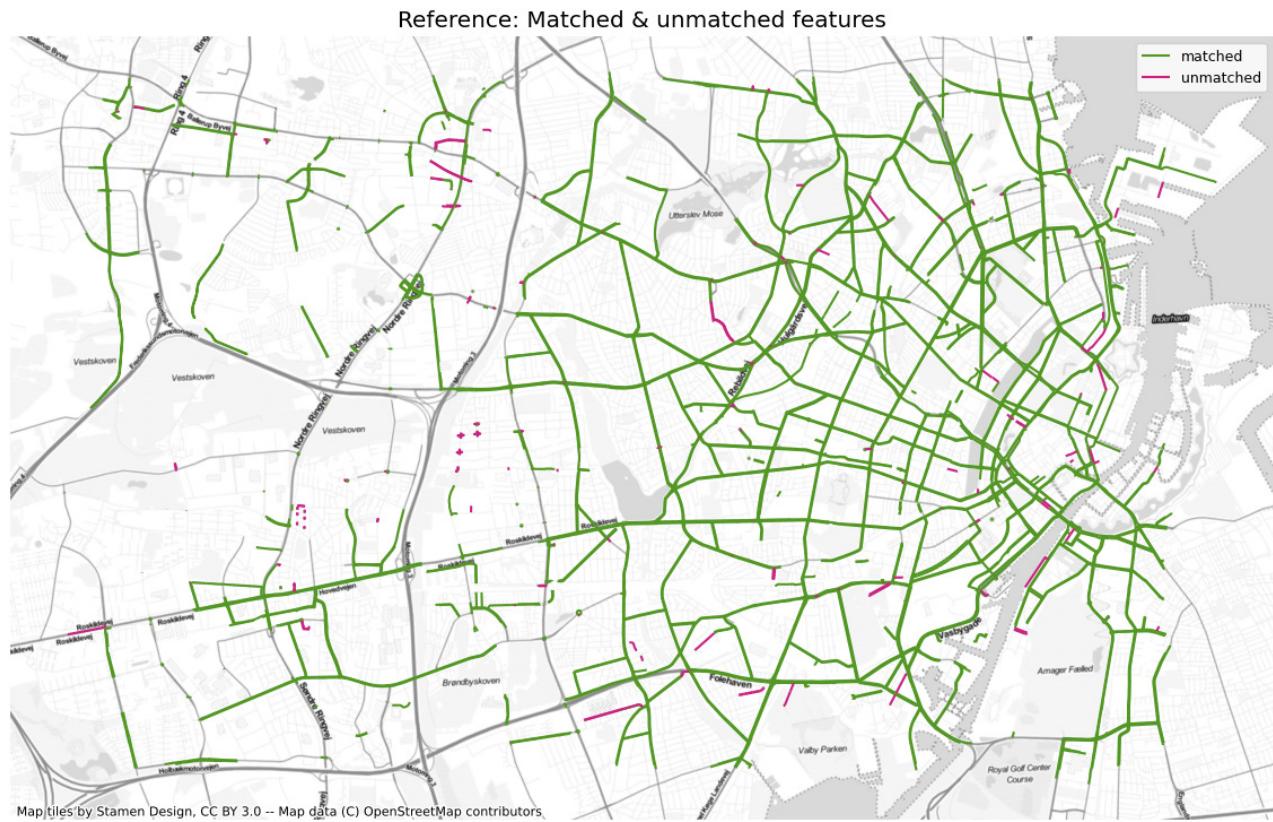
Plot feature matching results



Matched and unmatched features in OSM and reference data sets

OSM: Matched & unmatched features





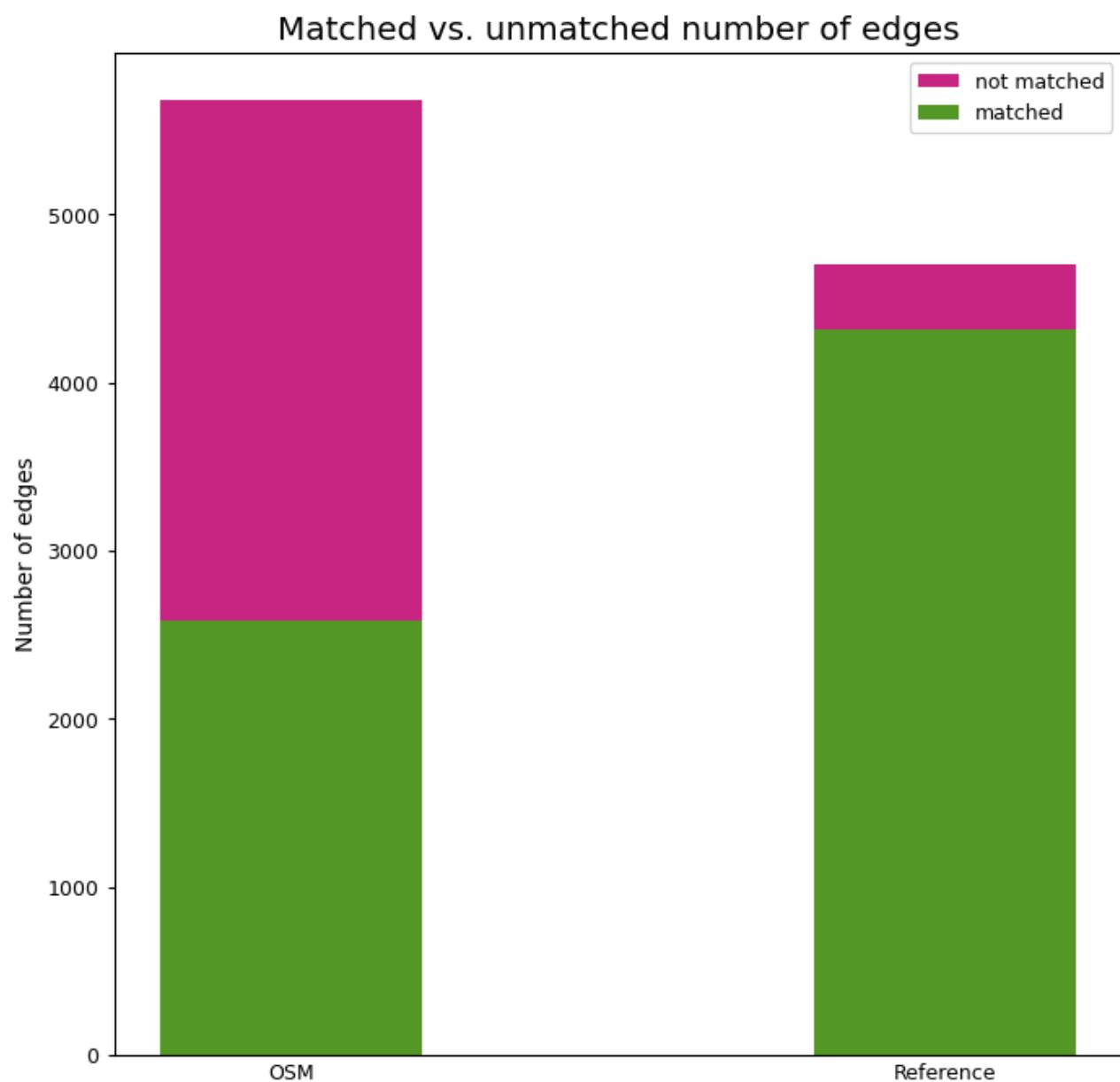
Summarized results

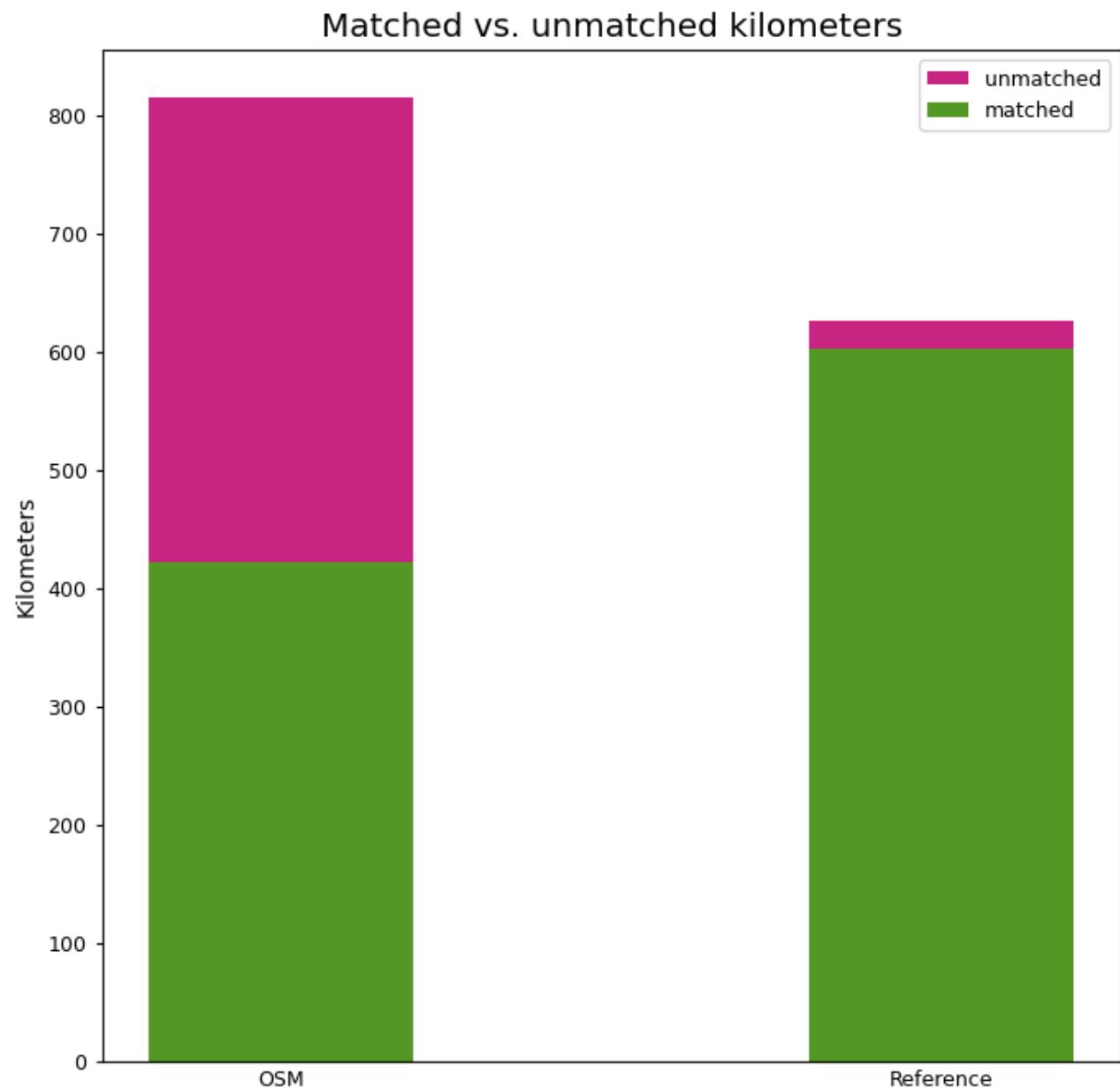
Edge count: 2588 of 5677 OSM edges (45.59%) were matched with a reference edge.

Edge count: 4313 out of 4705 reference edges (91.67%) were matched with an OSM edge.

Length: 422.50 km out of 814.72 km of OSM edges (51.86%) were matched with a reference edge.

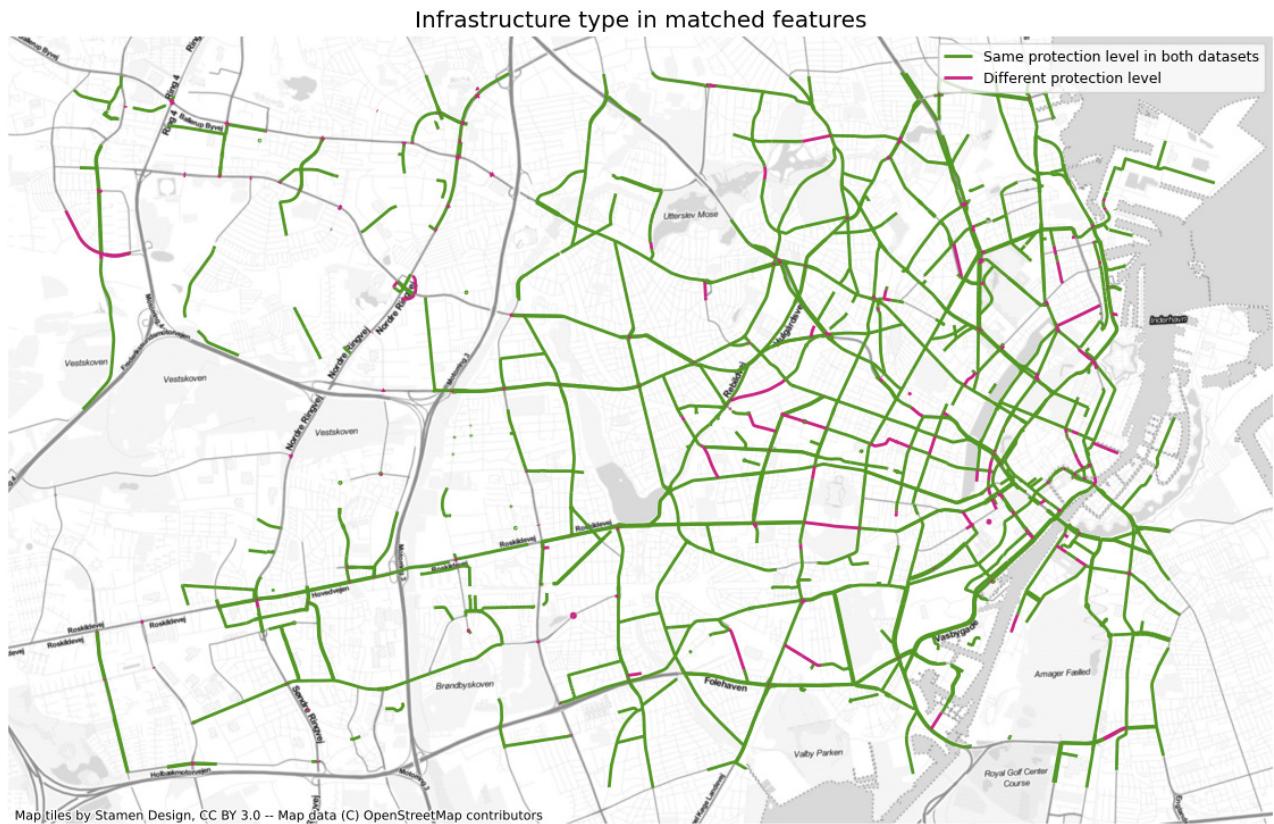
Length: 603.15 km out of 626.48 km of reference edges (96.28%) were matched with an OSM edge.





Analyze feature matching results

Comparison of infrastructure type in matched features

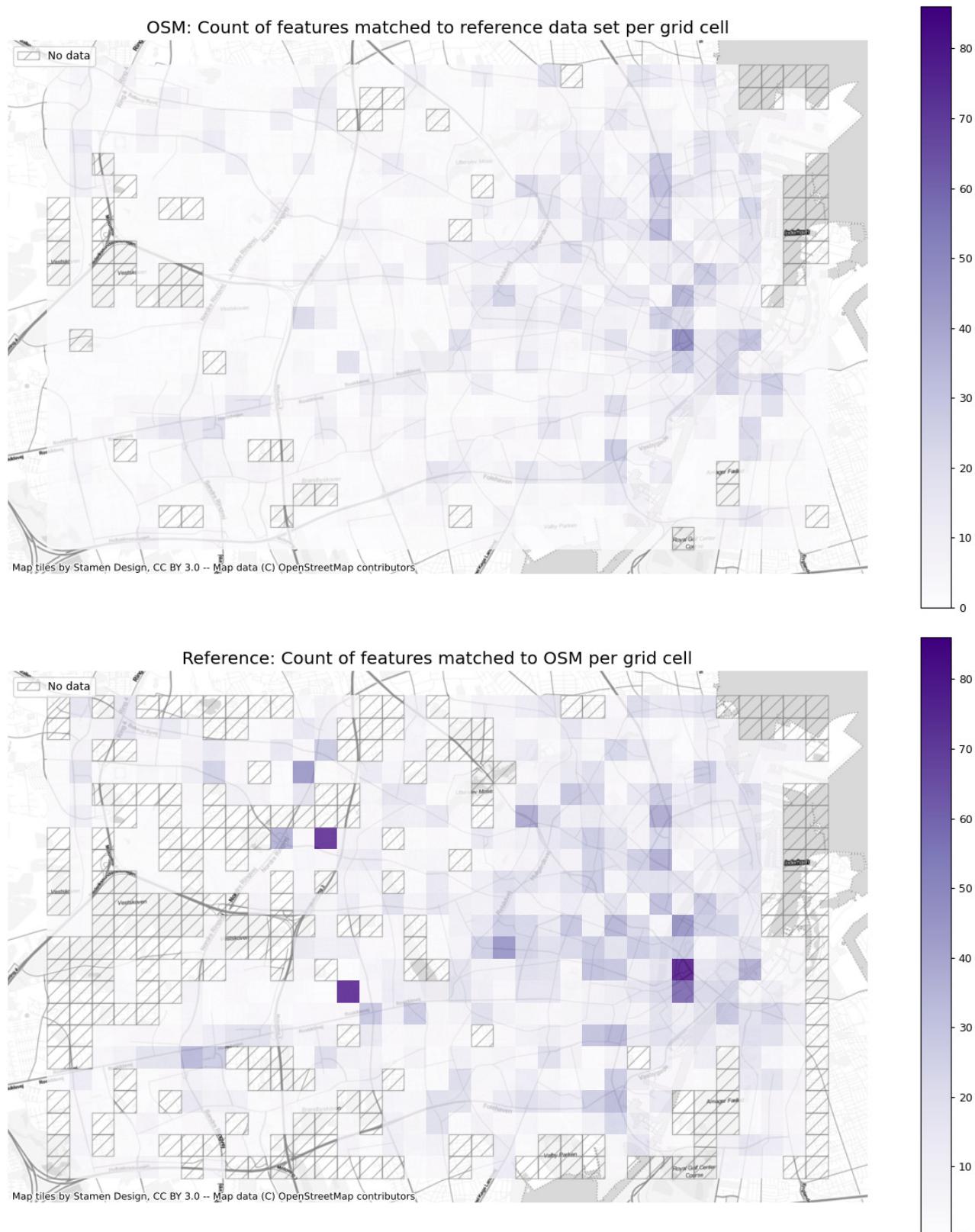


Local success rate of feature matching

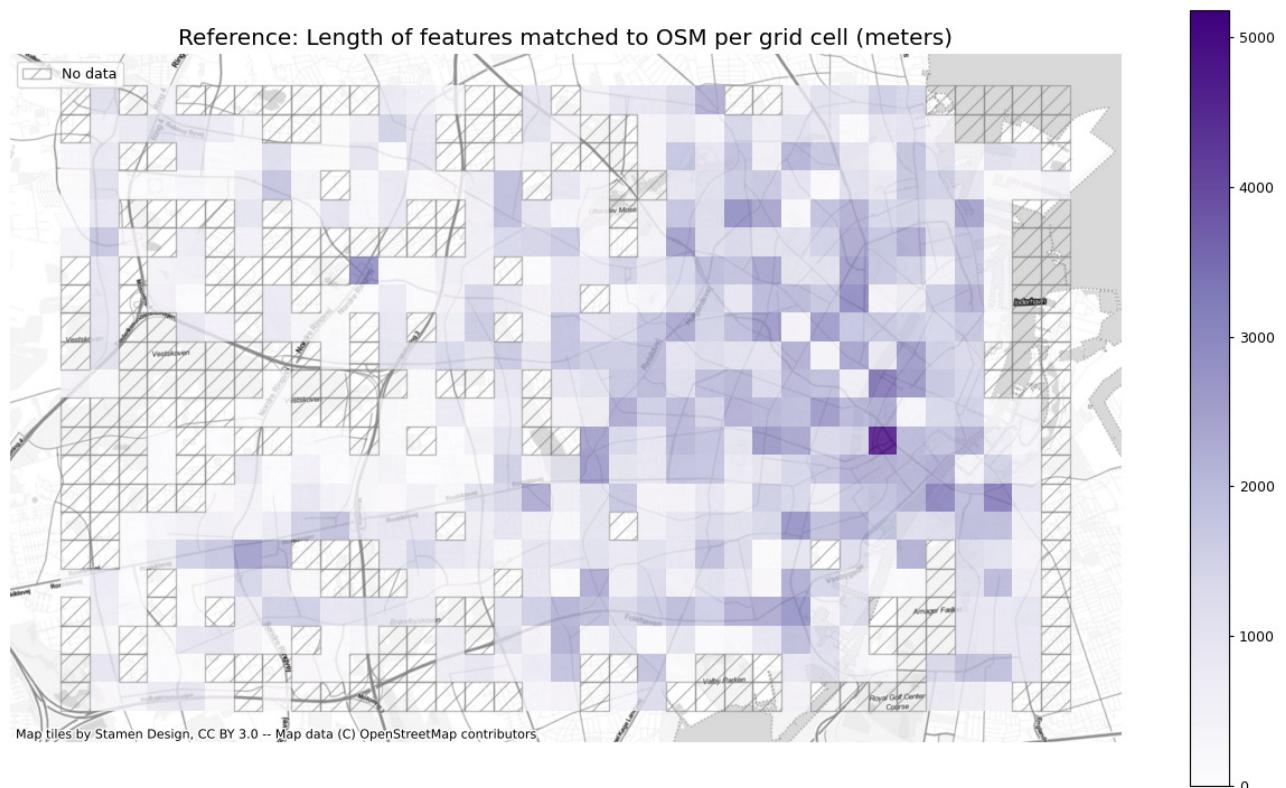
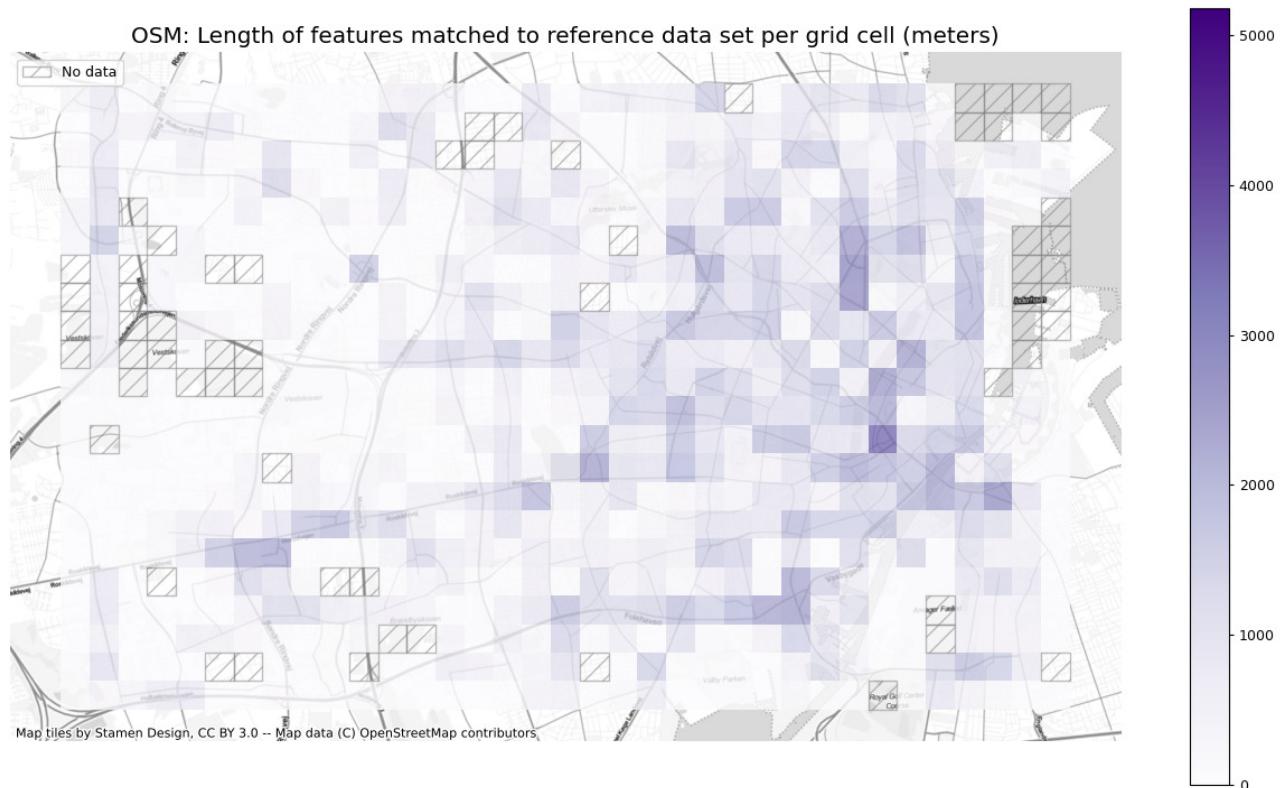
In the plots below, the count, percent, and length of matched and unmatched features in each dataset are summarized.

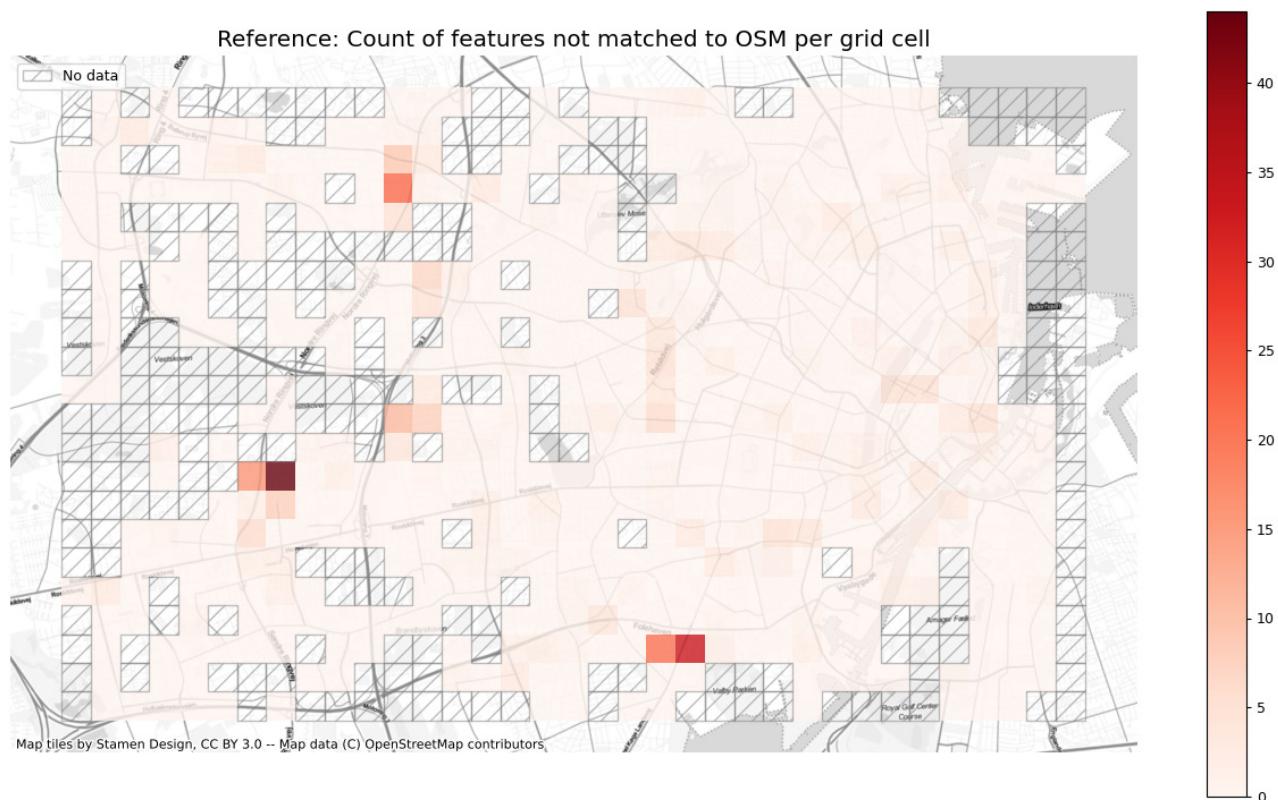
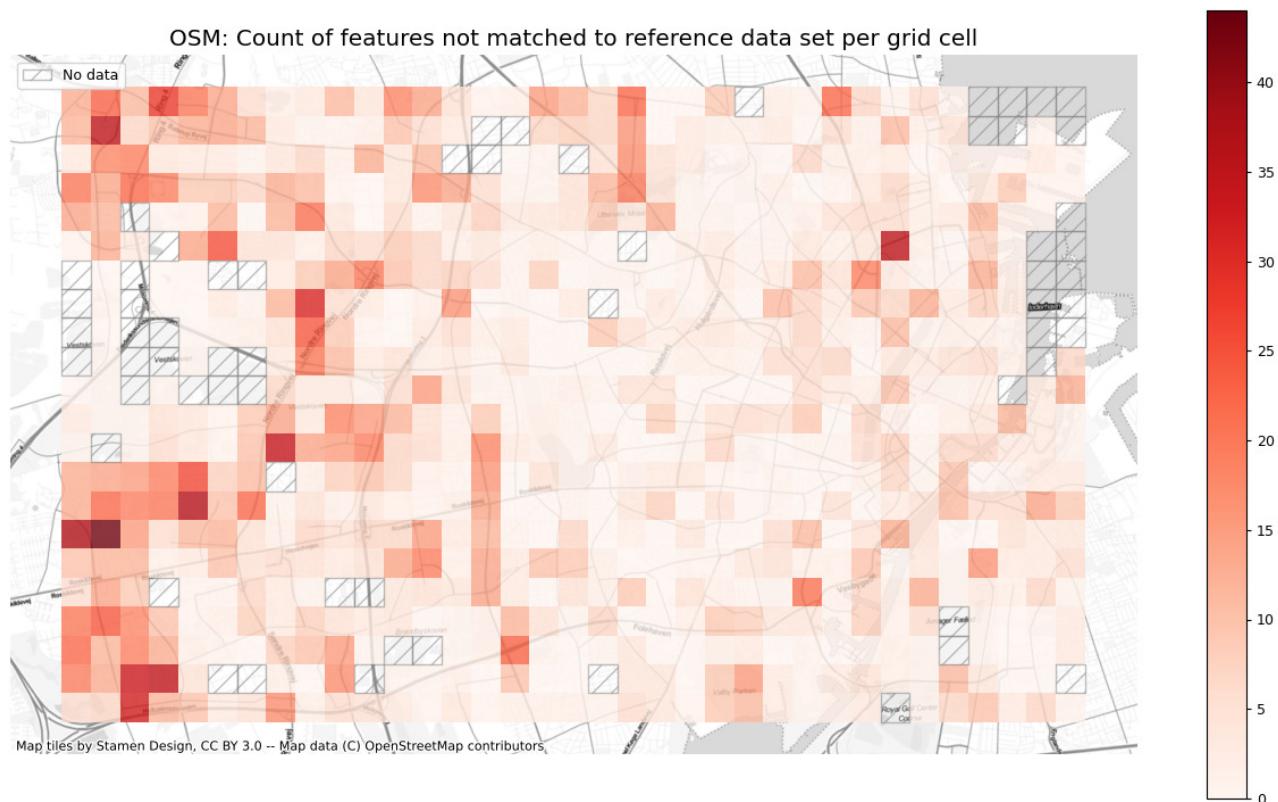
Warning

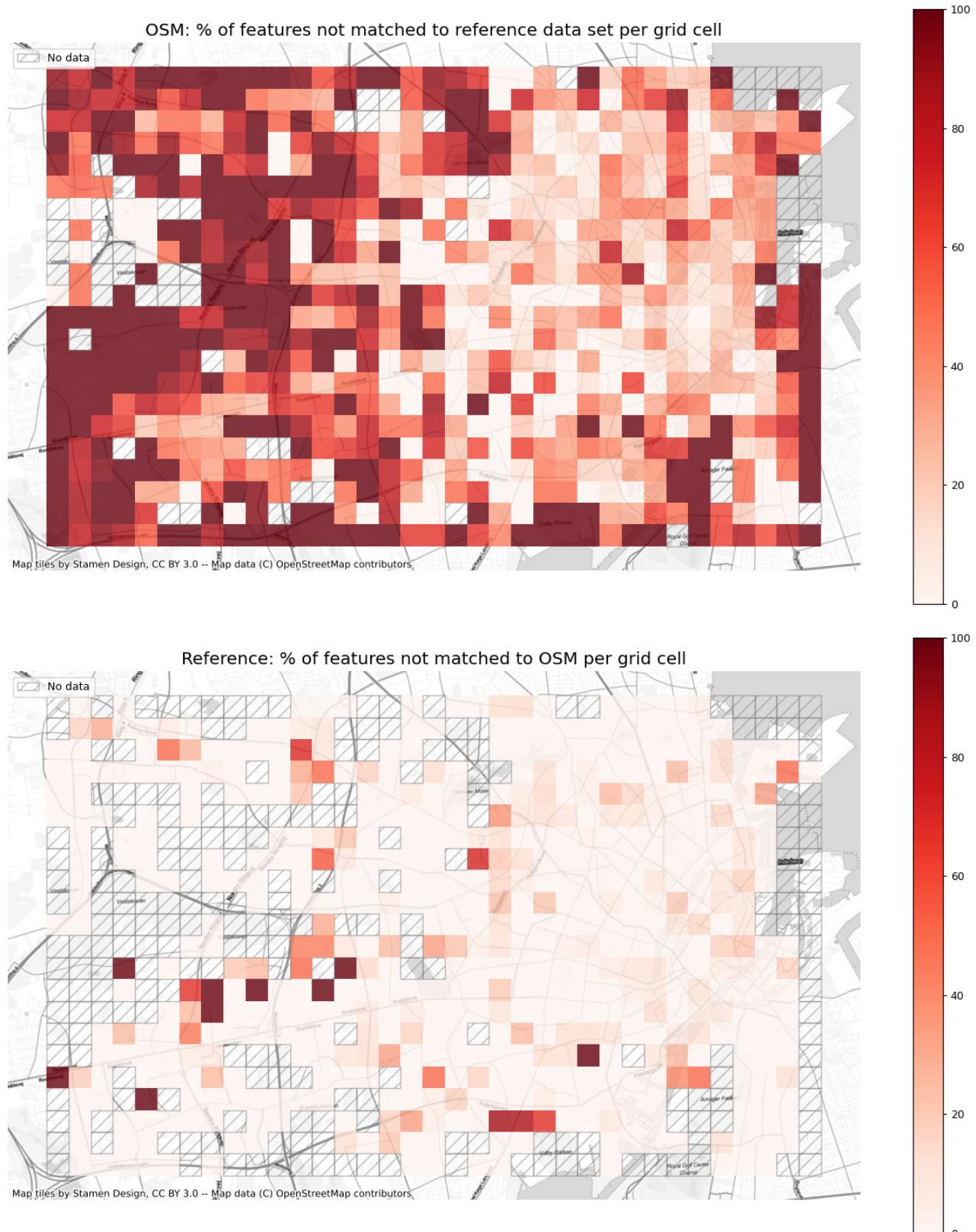
The number of matched features in one dataset in a grid cell does not necessarily reflect the number of matched features in the other dataset, since an edge can be matched to a corresponding edge in another cell. Moreover, the local count refers to edges intersected with the grid cell. For example, a long bike lane crossing 3 cells will thus be counted as matched in 3 different cells. This does not change the relative distribution of matched/unmatched features, but it does entail that the overall summary of matched/unmatched features above uses a different total count of edges than the plots below.

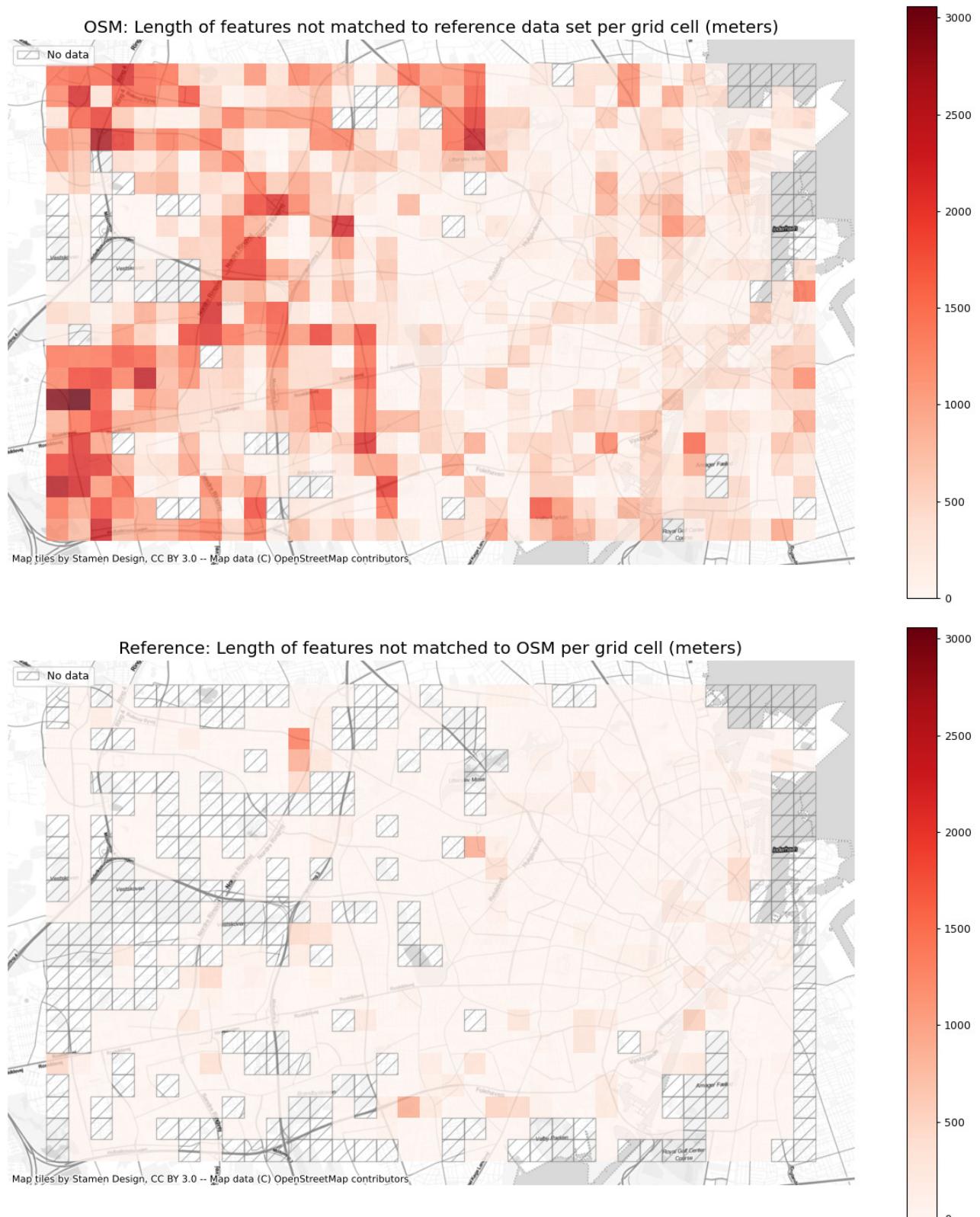












Summary

Results of Feature Matching

	OSM	Reference
Count of matched edges	2,588	4,313
Percent matched edges	46%	92%
Length of matched edges (km)	423	603
Percent of matched network length	52	96
Local min of % matched edges	4%	nan%
Local max of % matched edges	100%	nan%
Local average of % matched edges	62%	nan%
Local (cell) min of % matched edges	-	12
Local (cell) max of % matched edges	-	100
Local (cell) average of % matched edges	-	95

Time of analysis: Sat, 10 Dec 2022 13:27:14