

1b. Intrinsic Analysis of OSM Data

This notebook analyzes the quality of OSM bicycle infrastructure data for a given area. The quality assessment is *intrinsic*, i.e. based only on the one input data set without making use of external information. For an extrinsic quality assessment that compares the OSM data 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 data set as the ground truth, no universal claims of data quality can be made. The idea is rather to enable those working with OSM-based bicycle networks to assess whether the data are 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 [Ferster et al. \(2020\)](#), [Hochmair et al. \(2015\)](#), [Barron et al. \(2014\)](#), and [Neis et al. \(2012\)](#).

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](#)
 - [Under/overshoots](#)
 - [Missing intersection nodes](#)
- [Network components](#)
 - [Disconnected components](#)
 - [Components per grid cell](#)
 - [Component size distribution](#)
 - [Largest connected component](#)
 - [Missing links](#)

- Component connectivity
- Summary

Data completeness

Network density

In this setting, network density refers to the length of edges or number of nodes per km². This is the usual definition of network density in spatial (road) networks, which is distinct from the *structural* network density known more generally in network science. Without comparing to a reference data set, network density 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.

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 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:

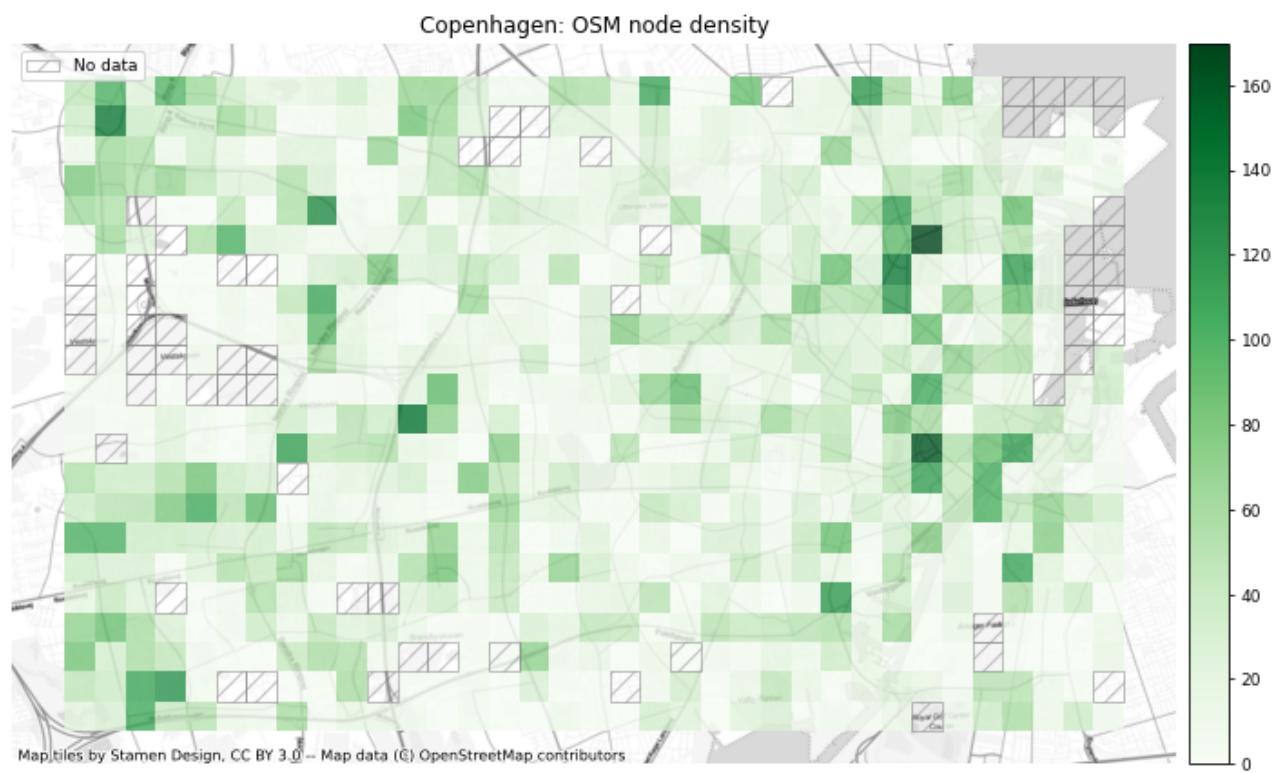
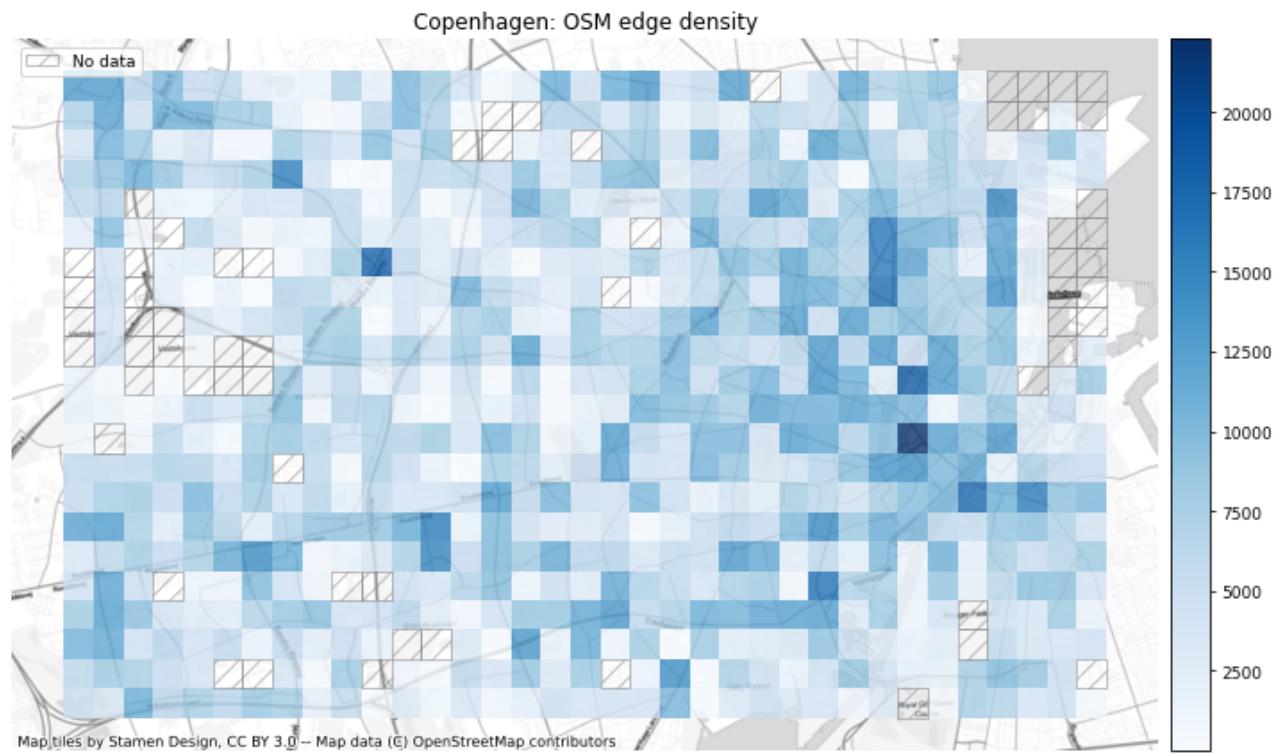
- lower-than-average infrastructure density indicates a locally sparser network
- higher-than-average node density indicates that there are relatively many intersections in a grid cell
- higher-than-average dangling node density indicates that there are relatively many dead ends in a grid cell

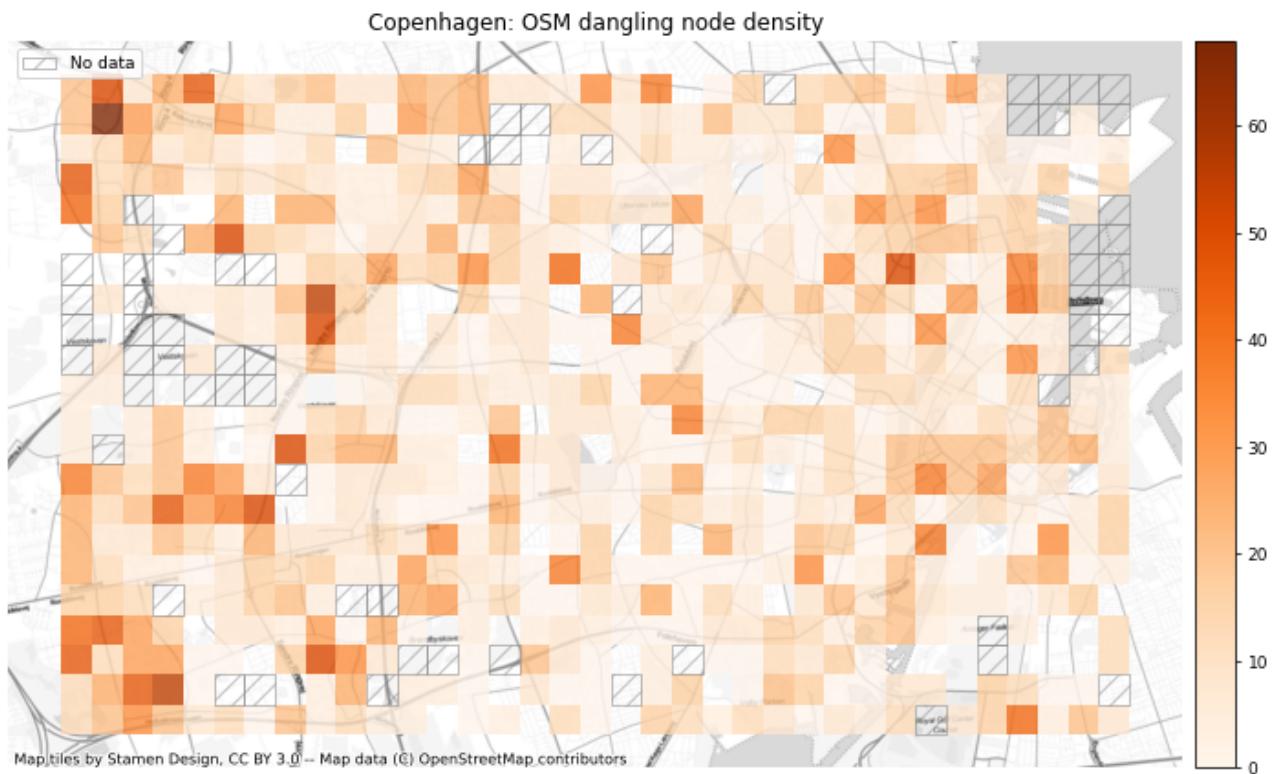
Global network density

For the entire study area, there are:

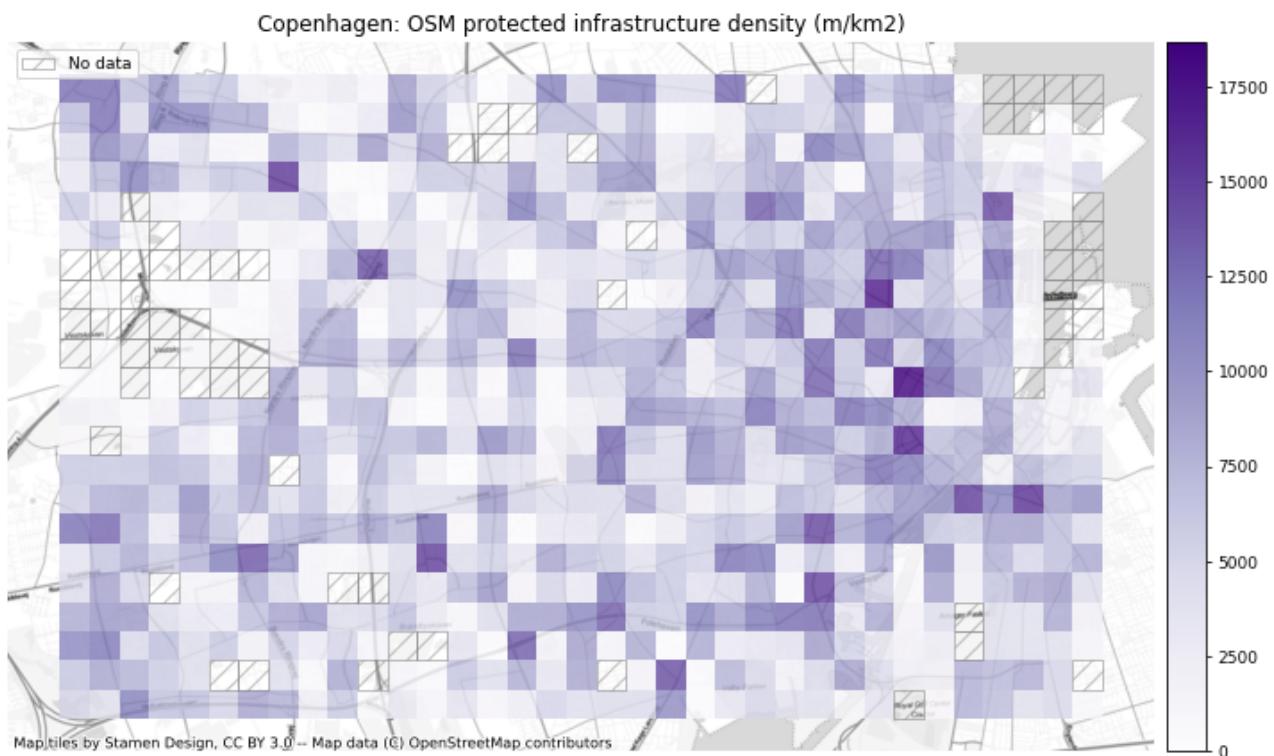
- 5861.46 meters of bicycle infrastructure per km².
- 27.15 nodes in the bicycle network per km².
- 10.02 dangling nodes in the bicycle network per km².
- 5302.84 meters of protected bicycle infrastructure per km².
- 499.41 meters of unprotected bicycle infrastructure per km².
- 59.21 meters of mixed protection bicycle infrastructure per km².

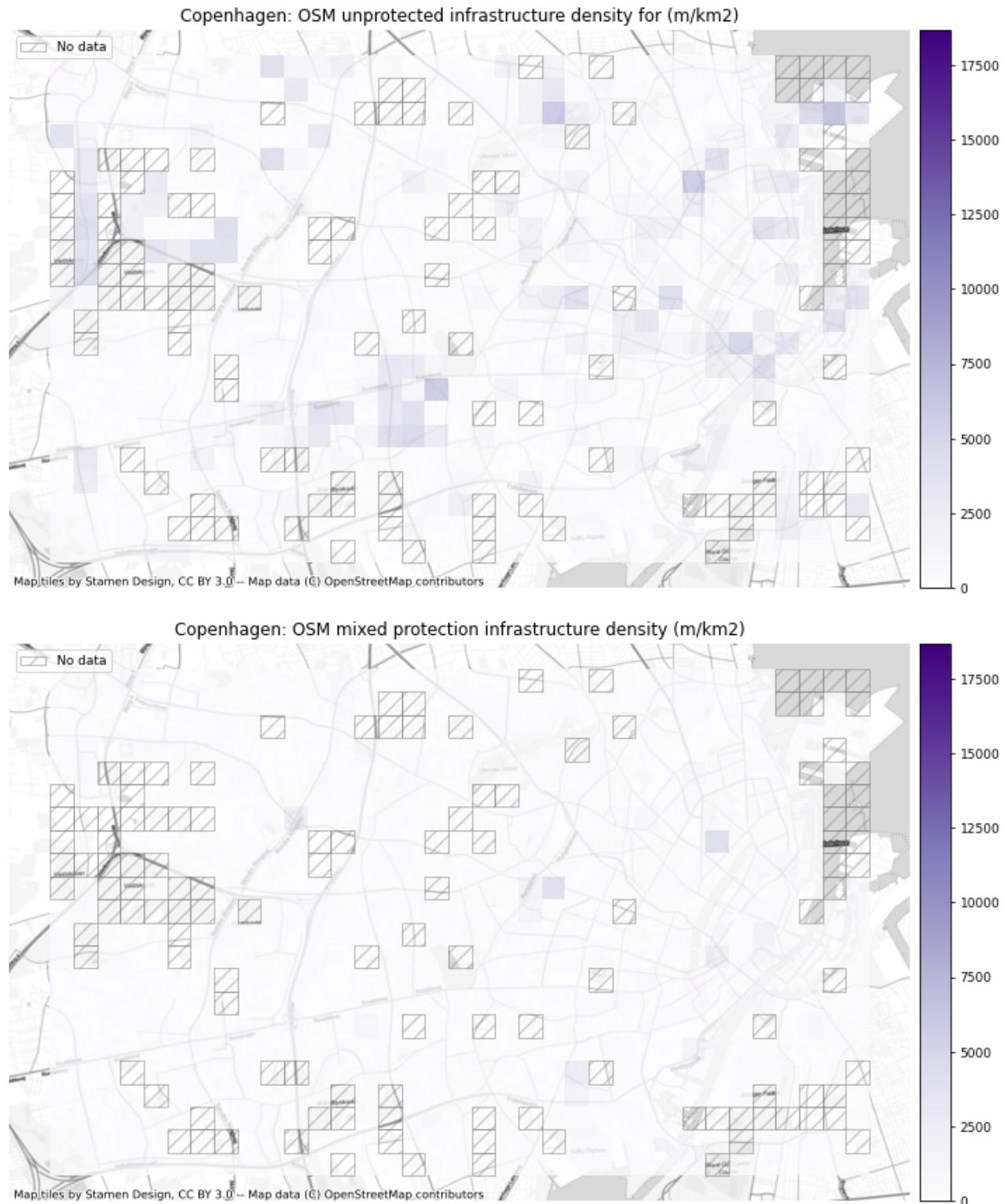
Local network density





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 an 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 sometimes result in conflicting tags which can undermine the evaluation of cycling conditions.

This section includes analyzes of missing tags (edges with tags that lack information), incompatible tags (edges with tags labelled with two or more contradictory tags), and tagging patterns (the spatial variation of which tags are being used to describe bicycle infrastructure).

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 depends on the use case - for example, the tag `lit` for a project that studies light conditions on cycle paths. The workflow below allows to quickly analyze the percentage of network edges that have a value available 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 the study area level, a higher percentage of existing tag values indicates in principle a higher quality of the data set. However, this is different from an estimation of whether the existing tag values are truthful. On the grid cell level, lower-than-average percentages for existing tag values can indicate a more poorly mapped area. However, 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 that contains 200 edges.

Global missing tags

Analysing tags describing:

surface – width – speedlimit – lit –

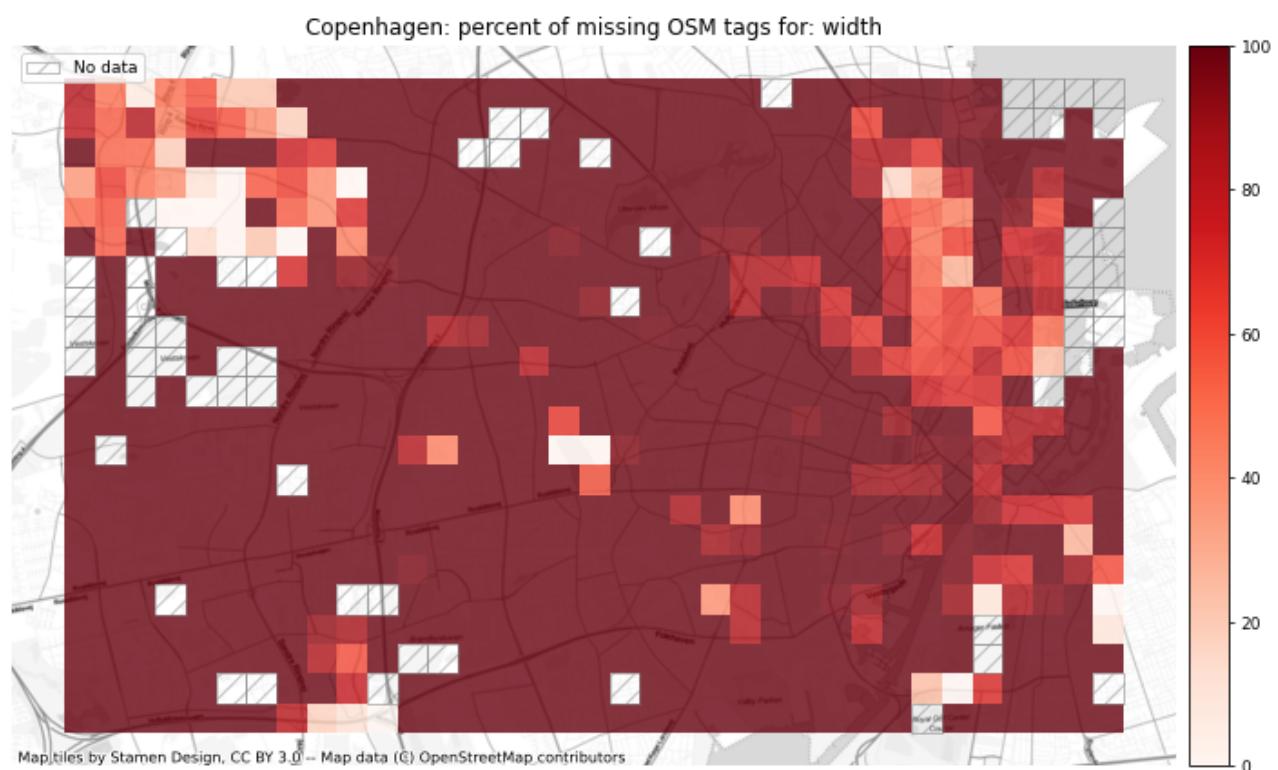
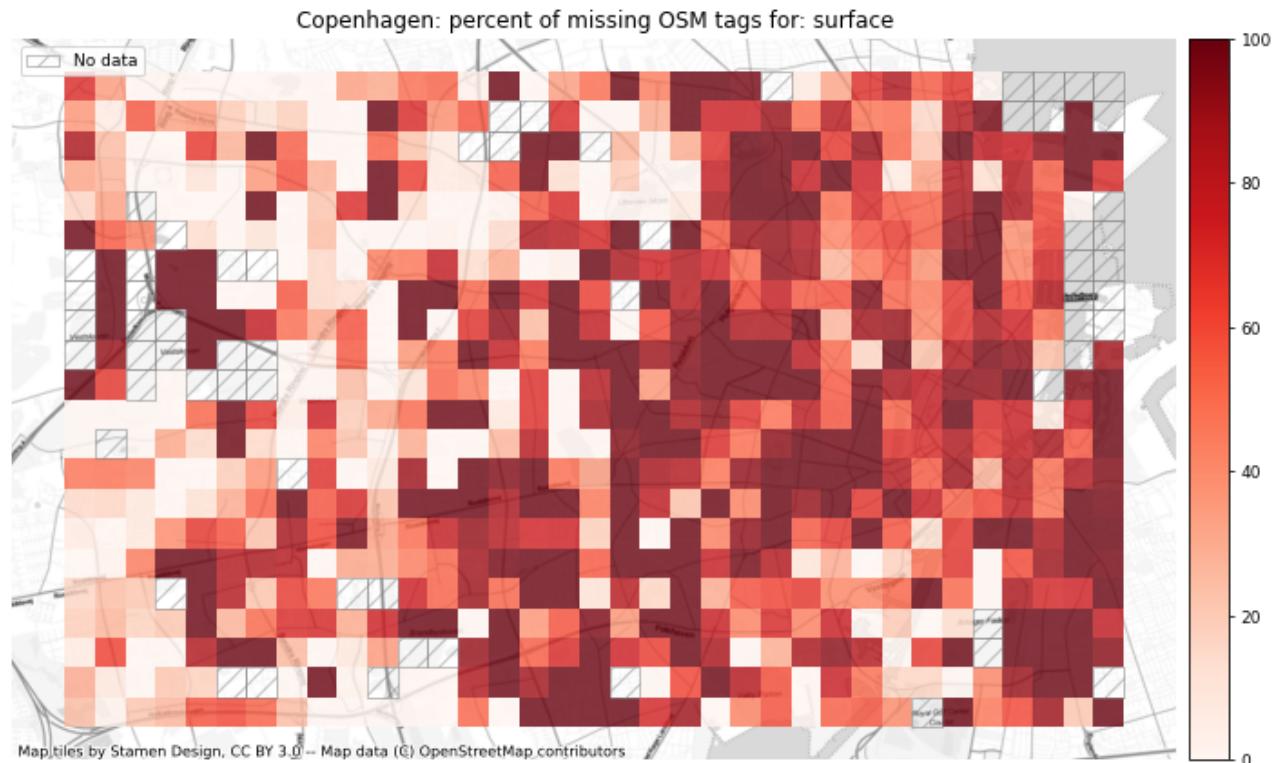
surface: 23325 out of 50959 edges (45.77%) have information.
surface: 552 out of 1285 km (42.95%) have information.

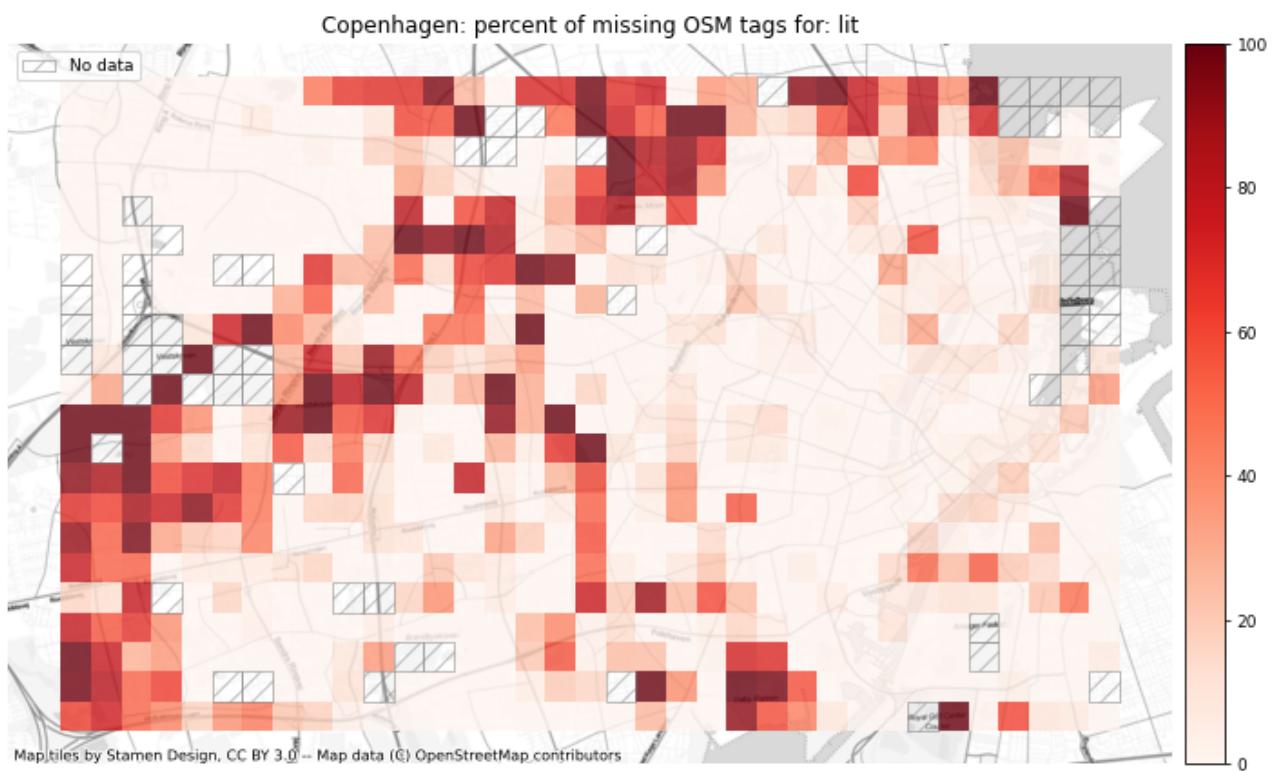
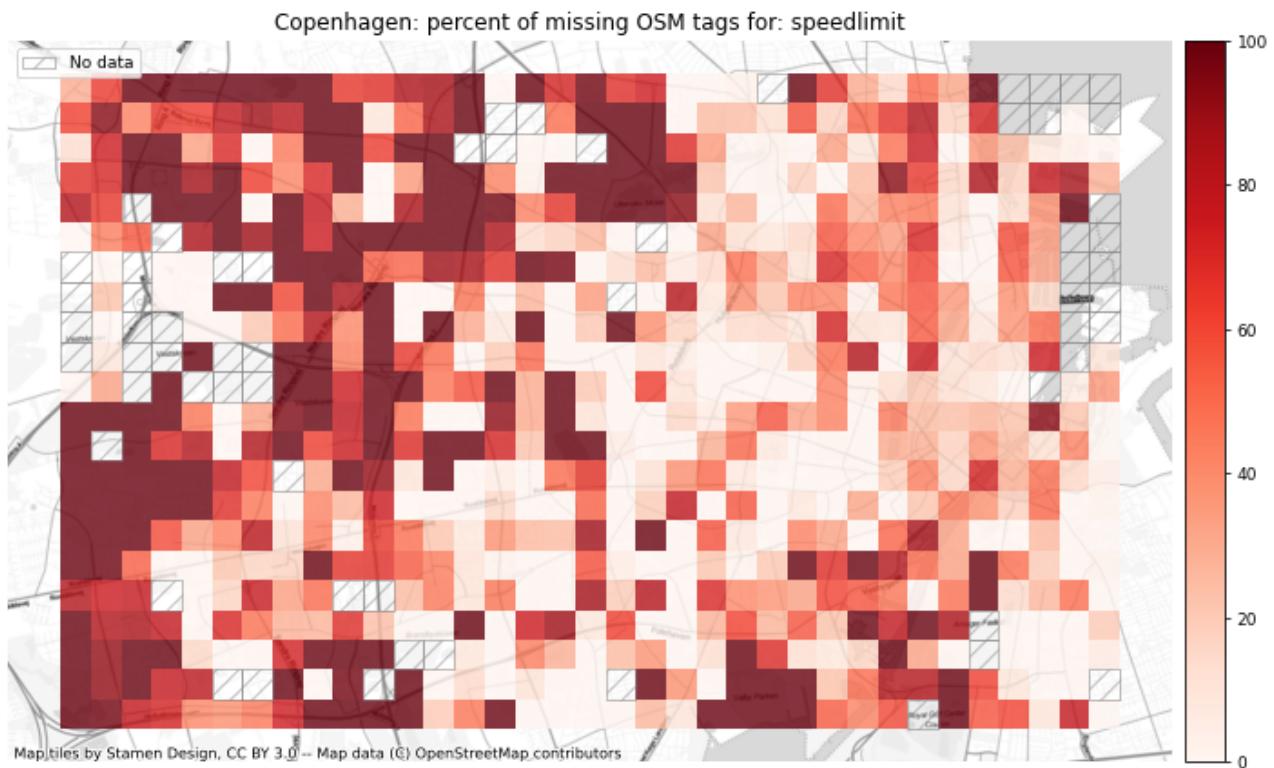
width: 5015 out of 50959 edges (9.84%) have information.
width: 97 out of 1285 km (7.56%) have information.

speedlimit: 25348 out of 50959 edges (49.74%) have information.
speedlimit: 684 out of 1285 km (53.21%) have information.

lit: 39318 out of 50959 edges (77.16%) have information.
lit: 1008 out of 1285 km (78.45%) 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

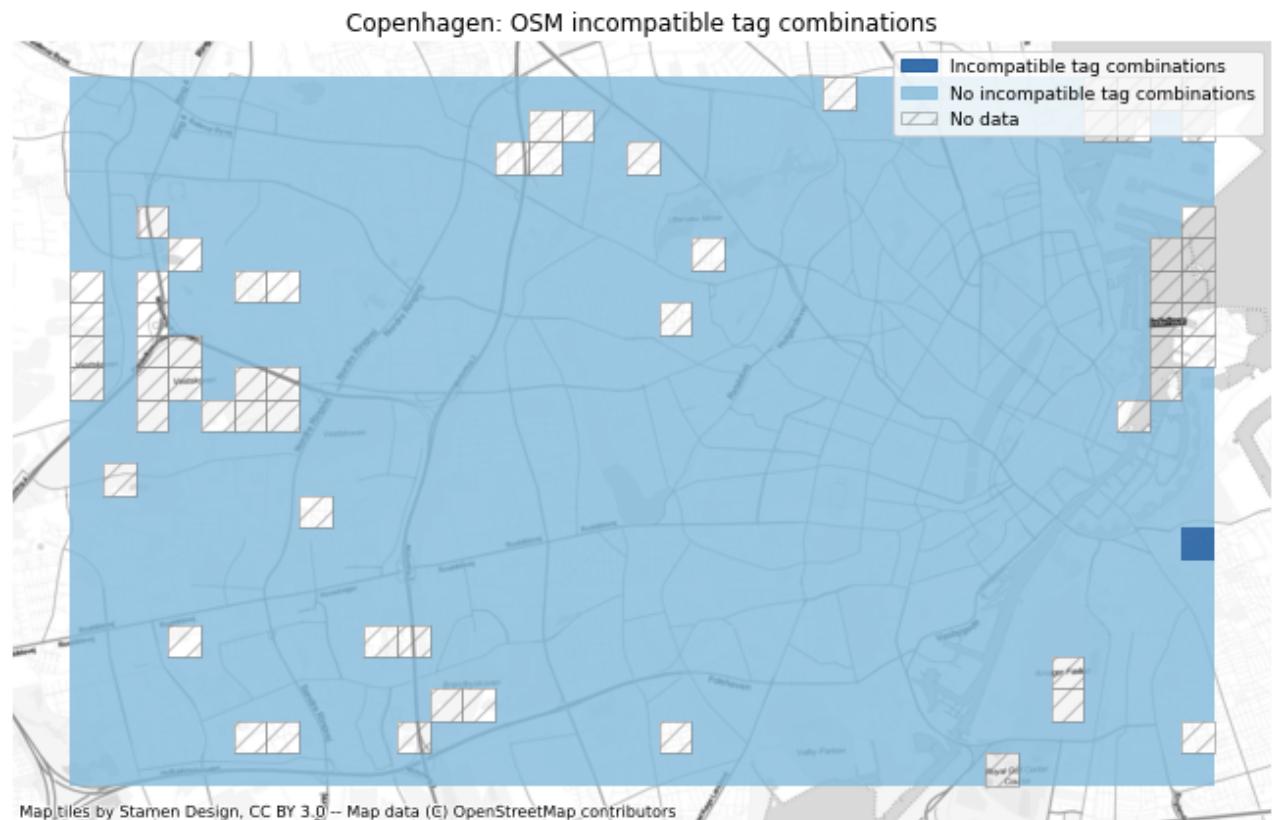
In the `config.yml` file, a list of incompatible key-value pairs for tags in the `incompatible_tags_analysis` is defined. Since there is no limitation to which tags a data set could potentially contain, the list is, by definition, non-exhaustive, and can be adjusted by the user. In the section below, `check_incompatible_tags` is run, which identifies all incompatibility instances for a given area, first on the study area level and then on the 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 data set, there are 2 incompatible tag combinations (of those defined in the configuration file).

Local incompatible tags (per grid cell)



Plotting incompatible tag geometries

Interactive map saved at results/OSM/cph_geodk/maps_interactive/tagsincompatible_osm.html

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 can 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

We first plot 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

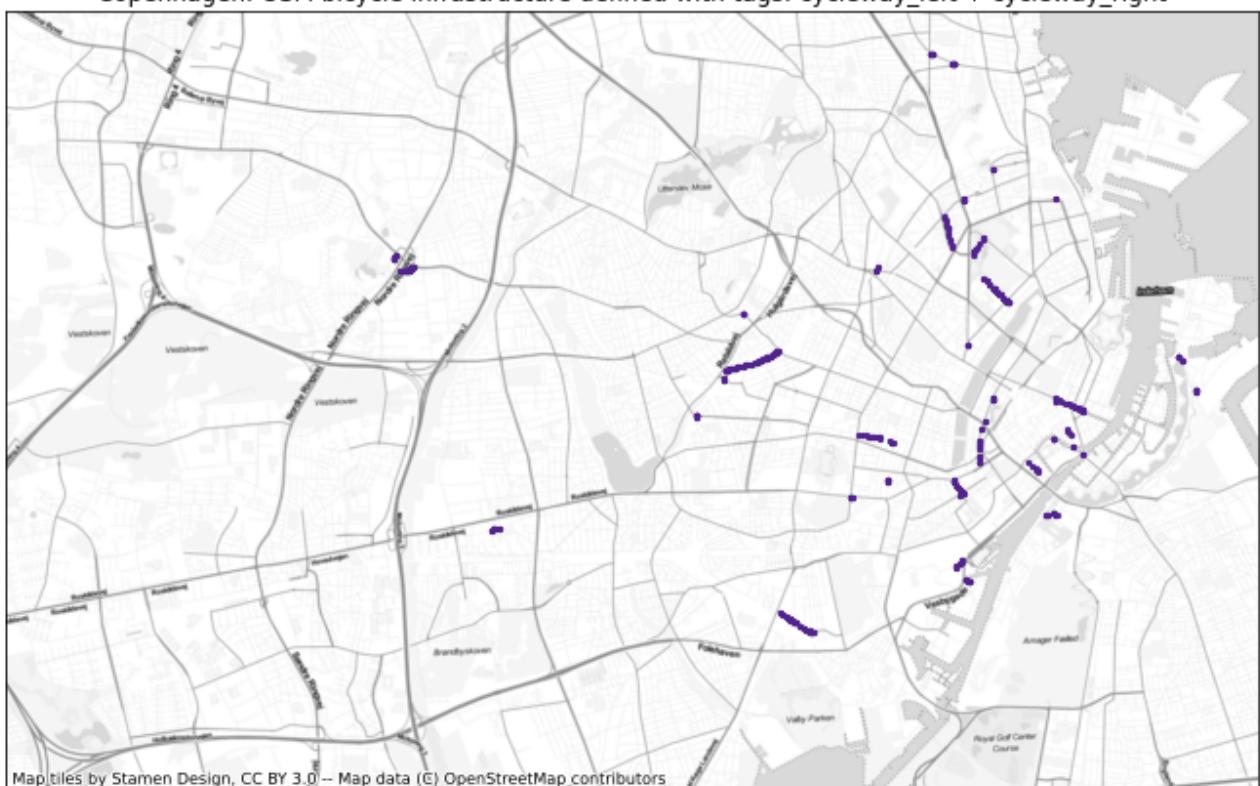
The plots for each tagging type allow for a quick visual overview of different tagging patterns present in the area. Based on local knowledge, the user may estimate whether the differences in tagging types are due to actual physical differences in the infrastructure or rather an artefact of the OSM data. 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

Interactive map saved at [results/OSM/cph_geodk/maps_interactive/taggingtypes_osm.html](#)

Multiple tagging

Copenhagen: OSM bicycle infrastructure defined with tags: `cycleway_left + cycleway_right`



Copenhagen: OSM bicycle infrastructure defined with tags: `cycleway + cycleway_both`



Copenhagen: OSM bicycle infrastructure defined with tags: highway + cycleway



Copenhagen: OSM bicycle infrastructure defined with tags: `cycleway_left` + `cycleway_both`



Copenhagen: OSM bicycle infrastructure defined with tags: cycleway + cycleway_right



Copenhagen: OSM bicycle infrastructure defined with tags: cycleway_right + cycleway_both



Interactive map saved at results/OSM/cph_geodk/maps_interactive/taggingcombinations_osm.html

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 bicycle networks, many metrics, such as missing links or network gaps, can simply reflect the true extent of the infrastructure ([Natera Orozco et al., 2020](#)). This is different for road 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.

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 1a 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 vast 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

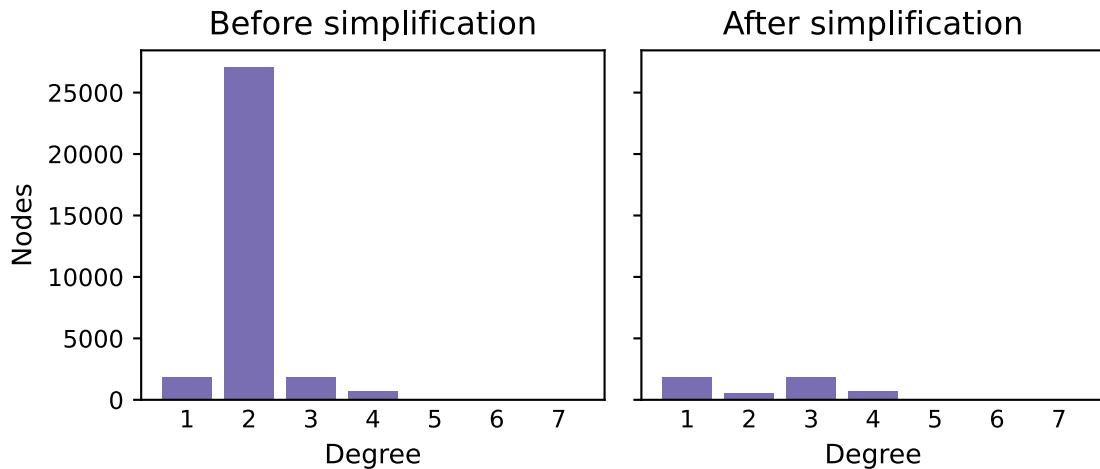
The degree distributions before and after simplification are plotted below.

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

Copenhagen: OSM degree distributions



Dangling nodes

Dangling nodes are nodes of degree one, i.e. they 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 (see 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 containing many dead-ends can indicate digitization errors and problems with edge over/undershoots.



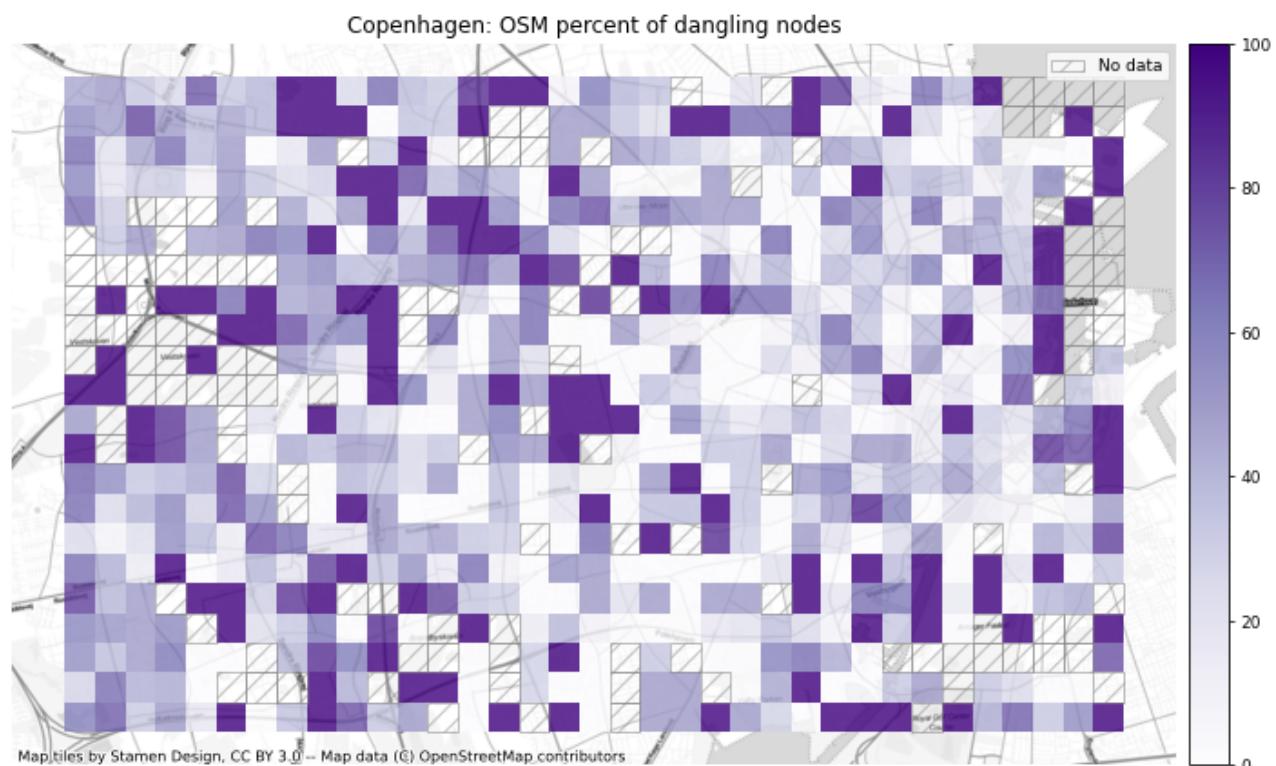
Left: Dangling nodes occur where road features end. Right: However, when separate features are joined at the end, there will be no dangling nodes.

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 color, all other nodes are shown in black.

Interpretation

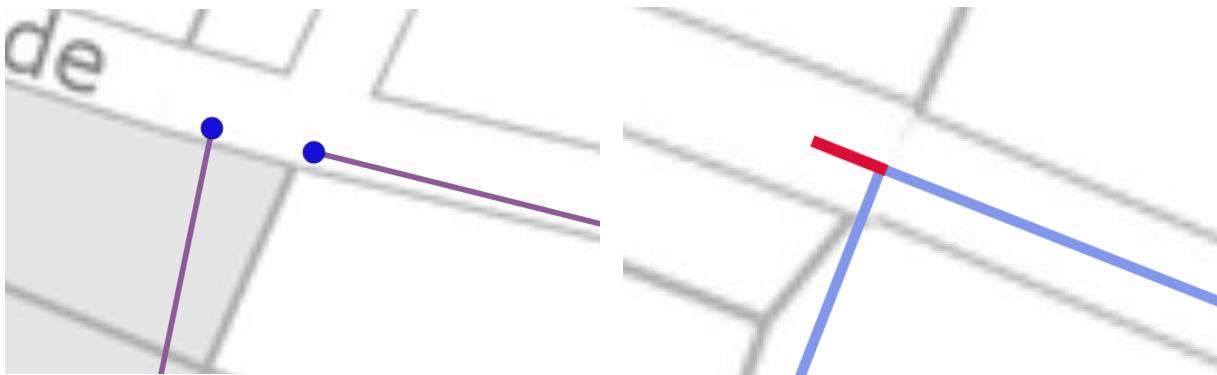
We recommend a visual analysis in order to interpret the spatial distribution of dangling nodes, with particular attention to areas of high dangling node density. It is important to understand where dangling nodes come from: are they actual dead-ends or digitization errors (e.g., over/undershoots)? A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



Interactive map saved at results/OSM/cph_geodk/maps_interactive/danglingmap_osm.html

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 undershoot occurs when two features are supposed to meet, but instead are just in close proximity to each other. An overshoot occurs when two features meet and one of them extends beyond the other. See the image below for an illustration of an undershoot (left) and an overshoot (right). For a more detailed explanation of over/undershoots, see the [GIS Lounge website](#).



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

Method

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.

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.

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

Interpretation

Under/overshoots 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

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

Interactive map saved at [results/OSM/cph_geodk/maps_interactive/underovershoots_3_3.osm.html](#)

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 the area, in order to determine ~~the origin of intersection issues and confirm/correct/reject them~~.

1 place(s) appear to be missing an intersection node or a bridge/tunnel tag.

Interactive map saved at `results/OSM/cph_geodk/maps_interactive/intersection_issues_osm.html`

Network 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

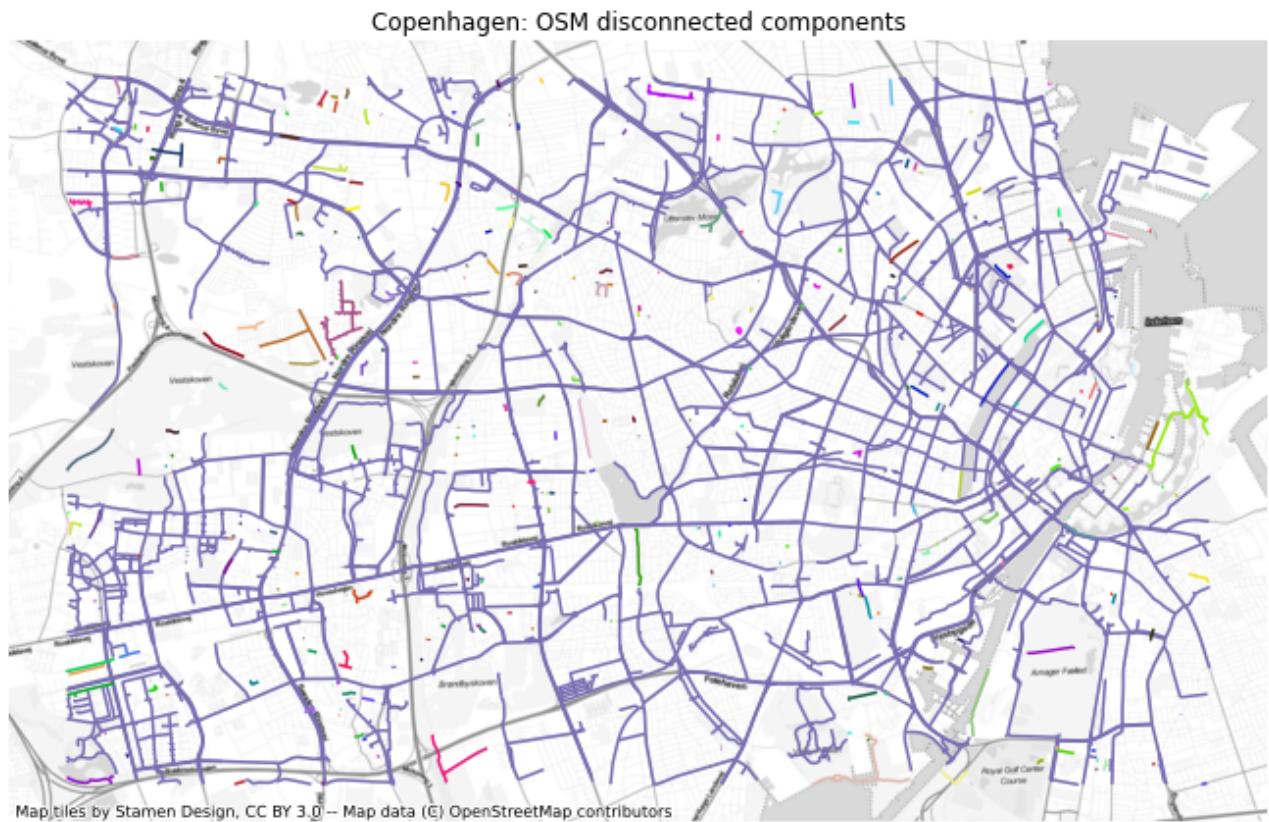
First, with the help of `return_components`, a list of all (disconnected) components of the network is obtained. The total number of components is printed and all components are plotted in different colors for visual analysis. Next, the component size distribution (with components ordered by the network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

Interpretation

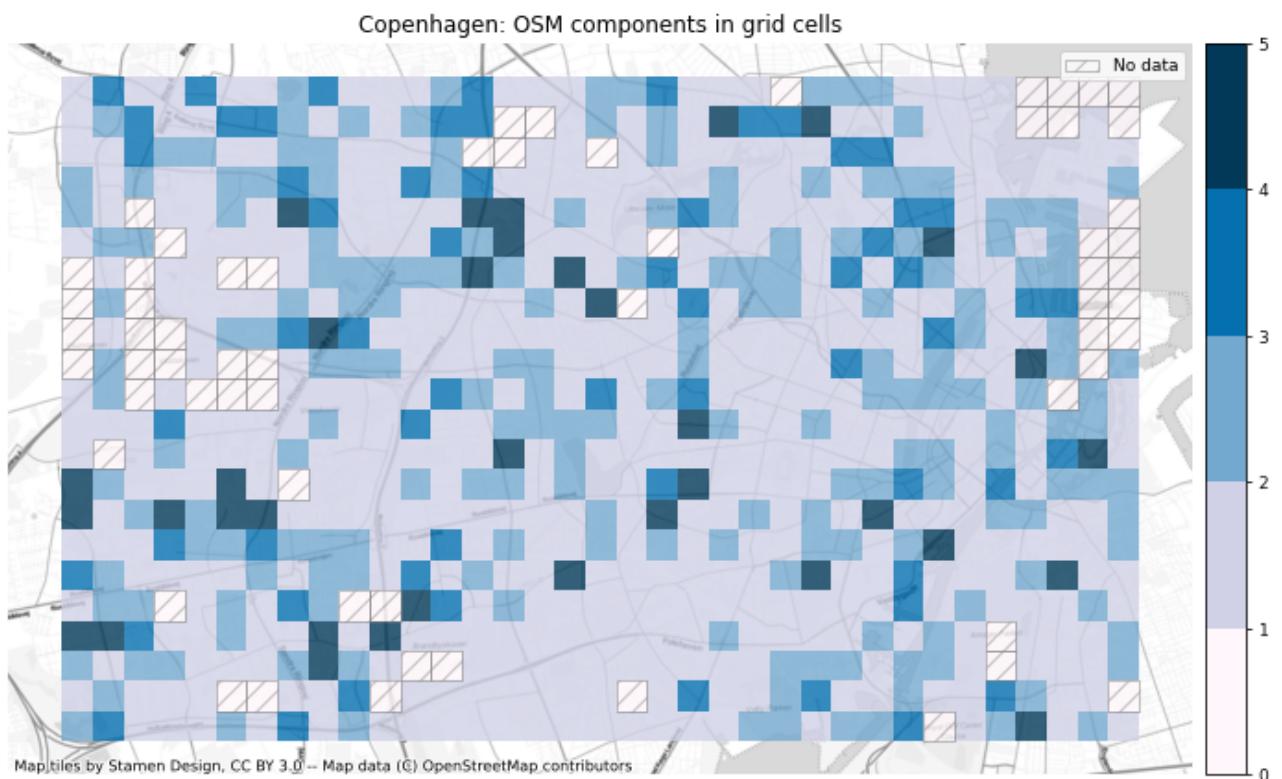
As with many of the previous analysis steps, knowledge of the area is crucial for a correct interpretation of component analysis. Given that the data represents the actual infrastructure accurately, bigger components indicate coherent network parts, while smaller components indicate scattered infrastructure (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.

Disconnected components

The network in the study area has 352 disconnected components.



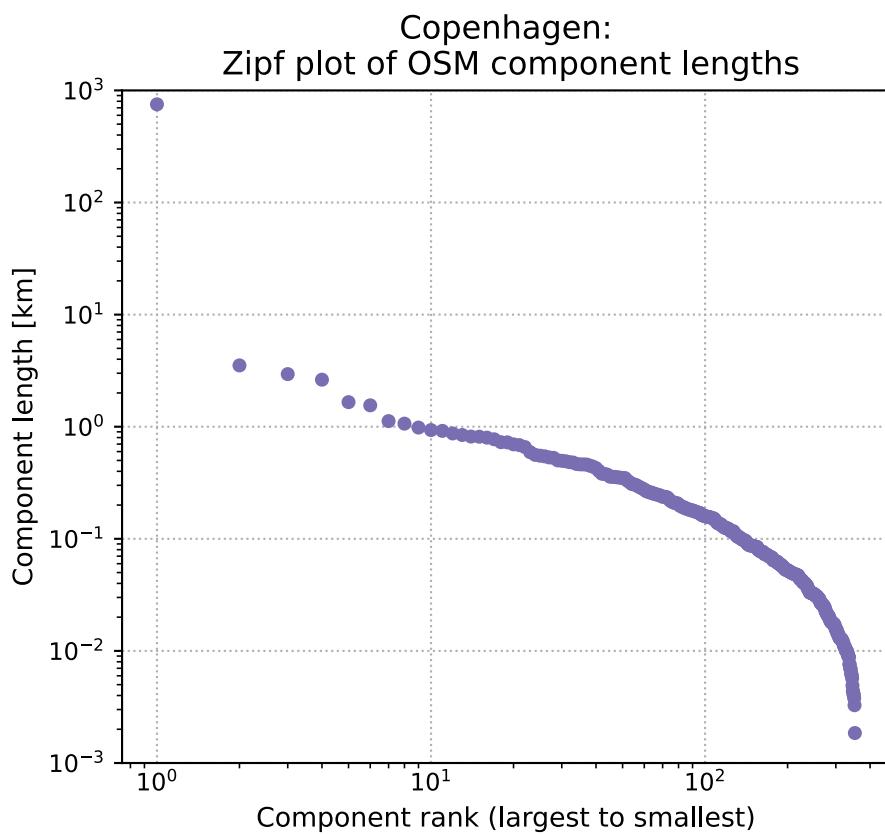
Components per grid cell



Component size distribution

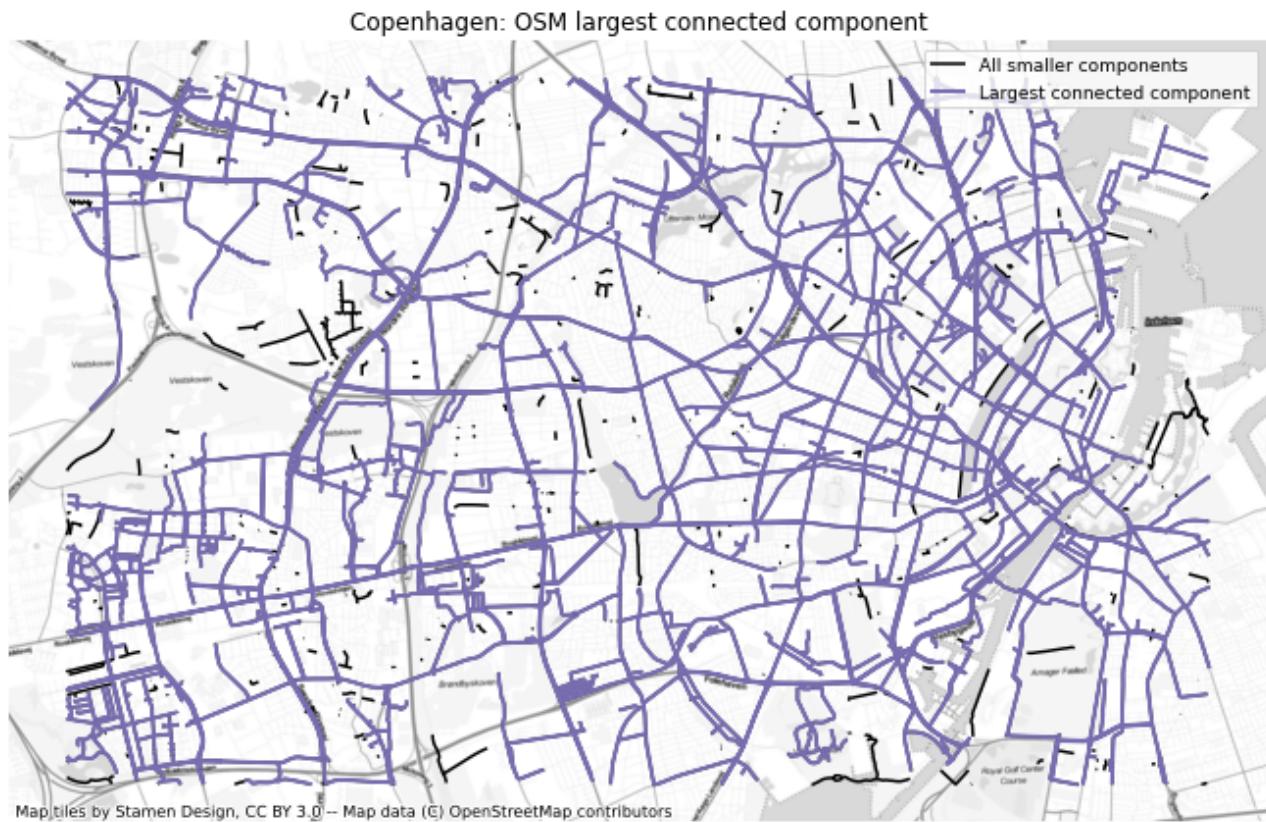
Many empirical distributions are skewed and often follow a power law, i.e. a straight line in a log-log plot, due to natural processes such as multiplicative network growth ([Clauset et al., 2009](#)). The network component size distribution (where size is length) can be visualized with a so-called Zipf plot, which plots the frequency of a component versus its rank (from largest to smallest). When a Zipf plot follows a straight line in log-log scale, it means that there is much higher chance to find small disconnected components than expected by a distribution from an exponential family (like a normal distribution). This can mean that there has been no consolidation of the network, only piece-wise or random additions ([Szell et al., 2022](#)).

However, it can also happen that the largest connected component (the leftmost marker in the plot at rank $1 = 10^0$) is a clear outlier, while the rest of the plot follows a different shape. This can mean that a consolidation *has* taken place, and that either a central planner has deliberately targeted to connect the network, or that the data are of high enough quality to have overcome many gaps.



Largest connected component

The largest connected component contains 92.30% of the network length.



Missing links

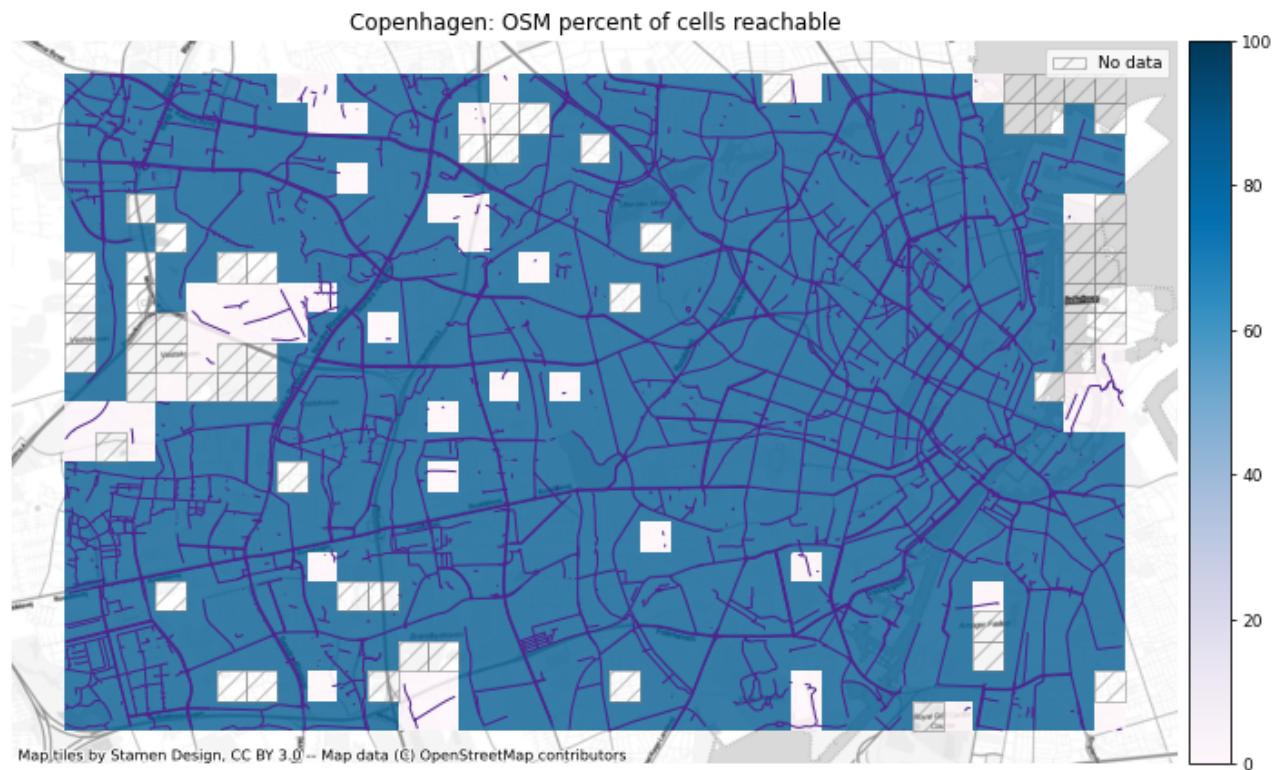
In the plot of potential missing links between components, all edges that are within the specified distance of an edge on another component are plotted. The gaps between disconnected edges are highlighted with a marker. The map thus highlights edges which, despite being in close proximity of each other,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

Analysis with component distance threshold of 10 meters:

Interactive map saved at results/OSM/cph_geodk/maps_interactive/component_gaps_10_osm.html

Component connectivity

Here we visualize 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 quickly highlight stark differences in network connectivity.



Summary

Intrinsic Quality Metrics - OSM data

Total infrastructure length (km)	1,063
Protected bicycle infrastructure density (m/km ²)	5,303
Unprotected bicycle infrastructure density (m/km ²)	499
Mixed protection bicycle infrastructure density (m/km ²)	59
Bicycle infrastructure density (m/km ²)	5,861
Nodes	4,925
Dangling nodes	1,818
Nodes per km ²	27
Dangling nodes per km ²	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