

Revealing the determinants of gender inequality in urban cycling with large-scale data

Alice Battiston,¹ Ludovico Napoli,² Paolo Bajardi,³ André Panisson,³
 Alan Perotti,³ Michael Szell,^{4,3,5} and Rossano Schifanella^{1,3}

¹*University of Turin, Via Giuseppe Verdi, 8, 10124 Torino TO, Italy*

²*Central European University, Quellenstraße 51, 1100 Wien, Austria*

³*ISI Foundation, Via Chisola 5, 10126 Torino TO, Italy*

⁴*IT University of Copenhagen, Rued Langgaards Vej 7, 2300 København, Denmark*

⁵*Complexity Science Hub, Josefstädter Str. 39, 1080 Wien, Austria*

Cycling is an outdoor activity with massive health benefits, and an effective solution towards sustainable urban transport. Despite these benefits and the recent rising popularity of cycling, most countries still have a negligible uptake. This uptake is especially low for women: there is a largely unexplained, persistent gender gap in cycling. To understand the determinants of this gender gap in cycling at scale, here we use massive, automatically-collected data from the tracking application Strava on outdoor cycling for 61 cities across the United States, the United Kingdom, Italy and the Benelux area. Leveraging the associated gender and usage information, we first quantify the emerging gender gap in recreational cycling at city-level. A comparison of cycling rates of women across cities within similar geographical areas unveils a broad range of gender gaps. On a macroscopic level, we link this heterogeneity to a variety of urban indicators and provide evidence for traditional hypotheses on the determinants of the gender-cycling-gap. We find a positive association between female cycling rate and urban road safety. On a microscopic level, we identify female preferences for street-specific features in the city of New York. Enhancing the quality of the dedicated cycling infrastructure may be a way to make urban environments more accessible for women, thereby making urban transport more sustainable for everyone.

Cycling is an outdoor activity associated with many individual and societal benefits. From the individual perspective, cycling has a positive impact on both physical and mental health, with a strong link to improved cardio-respiratory fitness, decreased cardiovascular mortality risk, and reduced stress-levels [1–3]. From a societal viewpoint, cycling is an environmentally-friendly and highly economic commuting option, especially for typical urban trips [4]. Recently, the United Nation (UN) Sustainable Development Goals (SDG) identified it as a pivotal component of a sustainable urban-mobility system [5]. Interventions targeted at increasing the number of cyclists are recommended as one of the solutions against traffic congestion, increased emissions, poor air quality and road safety.

Despite these wide-ranging benefits, cycling is mostly a male-dominated activity with a large gap in participation rates between men and women. Cycling research and policy making that is mostly focused on improving mobility for the existing, dominant group, risks to ignore half of the population and sustainable mobility solutions for everybody [6, 7]. Data on the use of bike-sharing services in three large US cities (New York, Boston and Chicago) show that only one in four bicycle trips in the 4-year period between 2014 and 2018 was made by a woman; other modes of transport, however, do not display comparable trip-share gaps [8]. Similarly, in San Francisco, only 29% of cyclists are women [9]. Recent data for England show that on average, not only do men take more bicycle trips per week than women, but they also cover longer distances [10]. A few European countries however, such as Denmark, Germany and the Netherlands represent the

main exception to this pattern, with women making up for more than 45% of all cyclists in these areas already in 2005 [11]. The evidence from this group of countries demonstrates that the reason for any kind of gender gap is not intrinsic but comes from place-specific barriers that need to be identified and, whenever possible, removed, if cycling should become a universal mode of transport.

The academic literature aimed at understanding the determinants of the gender gap in cycling links it on one hand to behavioral and psychological hypotheses. Women perceive cycling as a riskier activity compared to men, which would directly translate into a stronger preference for cycling infrastructure that is physically separated from motorized traffic [12–15]. On the other hand, physical route characteristics can play a role, for example in San Francisco, where women disfavor steep slopes, particularly for commuting [16]. In low-cycling contexts, women also report other deterring factors, such as an aversion for long distances and poor weather conditions, and a generally lower confidence in their cycling skills [17, 18]. Differences in preferences are typically stronger among occasional or non-cyclists than among regular cyclists [12], thus suggesting that policies targeting women are particularly needed to increase cycling uptake. The main limitation of these studies is that they are mostly conducted via surveys or experiments with typically low sample sizes and/or a limited geographical breadth, and therefore low statistical explanatory power – especially for the large number of possible confounders.

Recently, the emergence of new technologies for cycle-tracking and online-based services (e.g. bike-sharing) generated an unprecedented stream of automatically collected

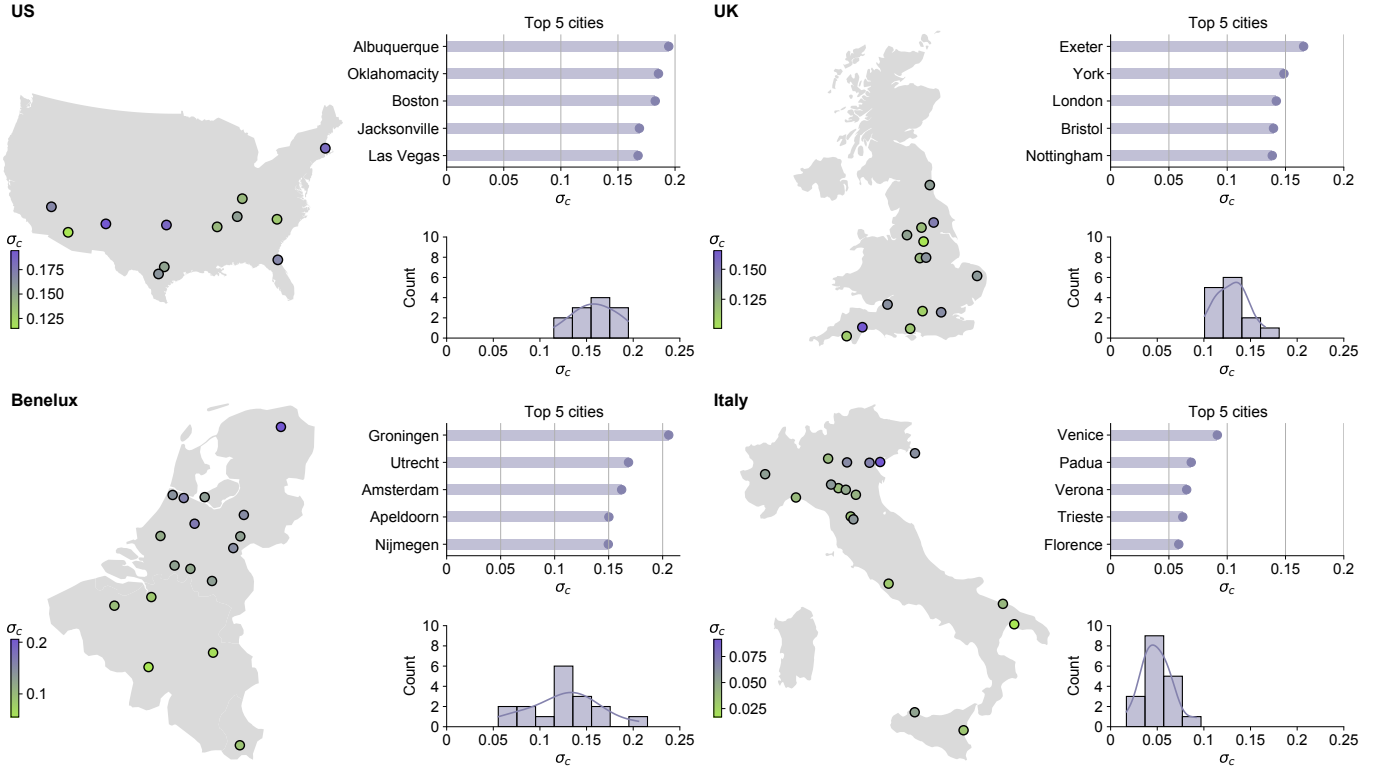


FIG. 1. **Gender gap in recreational cycling in Strava: overview of cities included in the study.** For each of the four geographical areas covered by the study, the figure depicts: the location of cities included in the analysis, the value of the female cycling rate σ_c for the five cities displaying the highest σ_c and the distribution of σ_c in the geographical area.

data on cycling behavior, which enlarge the potential for research in this area. In this context, data from bike-sharing services have been used to study whether interventions to the bike-sharing facilities impacted men and women differently in the city of New York [19] and, more generally, to study factors affecting the demand for these types of services [20]. Data from Strava Metro, a service provided by the sport-tracking application Strava, have been used to study exposure to air pollution for different groups of cyclists in the city of Glasgow [21] and cycling patterns and trends for the city of Johannesburg [22]. A few pioneering works have used GPS-based data to study route choices for different demographic groups in the city of San Francisco and Atlanta [14, 16].

In this study, we contribute to this strand of literature and use data from Strava to investigate the determinants of the gender gap in recreational cycling at a larger scale. With about 36 million users (2018 data) over 195 countries, Strava represents an unique data source on cycling-related behaviour [23], both in terms of the number of cyclists involved and the extent of the geographical coverage with a methodologically homogeneous data collection. For this study, we collect and use data for over 60 cities in four geographical areas across the United States and Europe, to explore the gender-cycling-gap at two different levels. First, we exploit the heterogeneity in the gender gap across

the various cities in our dataset to challenge traditional hypotheses from the literature on the determinants of the gender gap in cycling. In particular, we study the strength of association between the gender gap in cycling measured at the level of urban centers and a set of urban indicators, spanning from morphological characteristics of the cities to safety indicators capturing the prevalence of cycleways and streets with low-speed limit in the road network. In the second part of the study, we move the analysis from a macro to a micro level. Here, focusing on the city of New York, we model the gender-cycling-gap measured at street-level in terms of specific urban features. By using logistic regression analysis, we investigate the association between the presence of dedicated cycling infrastructure and the volume of female cyclists on the street relative to men. The results indicate that streets with cycling infrastructures, particularly those ensuring the presence of physical separation for motorized traffic, are associated to a more balanced gender ratio, suggesting a way for policy makers to intervene to make urban environments more accessible for women.

RESULTS

Using Strava data to measure the gender gap in recreational cycling

We use Strava data to measure the gender gap in recreational cycling in 61 urban centers across four geographical areas: United States, United Kingdom, Italy and Benelux. Strava is an Internet service for tracking human exercise that relies on GPS data. The service supports up to 33 different activities, but it is mostly used for cycling and running. At the time of the data collection in 2018, Strava counted around 36 million users worldwide, corresponding to 0.6 billion recorded activities [23]. Of these, 284 millions were cycling-activities (47%), and approximately one in five of cycling-uploads were by women (50 million)[23]. Tracking of commuting is growing in popularity on Strava [23], however the majority of uploads refers to recreational and athletic cycling. The raw data consist of a collection of Strava segments, with information on users training on these from the associated leader boards. The data were processed to map gender and usage information from Strava segments to a network-based definition of streets and then aggregated for the entire city, following the pipeline described in the SI.

For each city c , we define the gender-cycling-gap as the ratio σ_c between the total kilometers travelled by female cyclists and the overall kilometers travelled by cyclists of both genders. This measure accounts both for gaps in trip-shares among men and women and for differences in travelled distances. By construction, σ_c varies between 0 (no female cyclists) and 1 (no male cyclists): a value below 0.5 indicates the presence of a positive gender-cycling-gap (i.e. men cycling more than women). The closer the value to 0 the stronger the gap. For each geographical area covered by the study, Fig. 1 provides an overview of the cities included in the study, showing the five urban centers associated with the highest σ_c for each area, as well as the location and the distribution of σ_c of all covered cities, highlighting a persistent gender gap in recreational cycling (the full ranking is provided in the SI). In our sample, the largest value for σ_c is 0.21 in the municipality of Groningen, Netherlands, indicating the presence of a substantial gender gap in recreational cycling for all cities under consideration.

Even within homogeneous geographical areas, we observe a substantial heterogeneity in σ_c across cities. In the area of Benelux, in particular, σ_c ranges between 0.06 (Charleroi, Belgium) and 0.21 (Groningen, Netherlands). Dutch cities (particularly those in the northern regions) generally outperform cities in Belgium and Luxembourg. Among Italian cities, we observe a characteristic geographical pattern, with urban centers in the north-east displaying a lower gender gap than cities in the south and north-west. This north-south dichotomy is likely to be linked to the morphological characteristics of the country

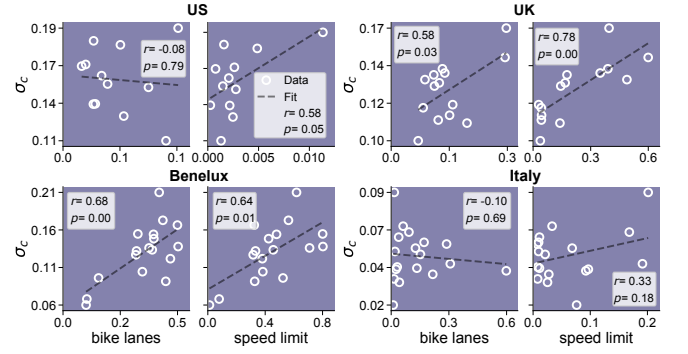


FIG. 2. **Correlations between gender ratio and urban road safety indicators.** The scatter plots show the correlations between two urban road safety indicators and the gender ratio σ_c , for cities in the four geographical areas separately. Each data point represents a city. The black line is the linear fit. The two urban road safety indicators capture the density of streets with a cycle lane in the street-network (indicator: *bike lanes*) and the density of streets with a speed limit up to 20 mi/h or 30 km/h in the street network (indicator: *speed limit*). A formal definition of the two indicators is provided in the SI.

and the presence of a large flat land with a well-established cycling tradition. Differences in economic development might partially explain this structure as well. No geographical patterns are instead observable for cities in the United States and in the United Kingdom included in our sample. Interestingly, there is no evident link between the gender ratio and the size of a city. For instance, large cities such as Boston, Amsterdam and London perform high in the corresponding ranking, while top-ranking positions in Italy are dominated by relatively smaller urban areas. When comparing different countries, we stress that in general the penetration and typical use of Strava might differ across geographical areas. For instance, data on the use of Strava indicates different usage patterns and adoption rates in the United States compared to other countries [24]. Therefore, we limit the comparison to cities within the same geographical area to ensure homogeneity in the Strava adoption by the general population, and we recommend to interpret a worldwide ranking with caution.

Cross-city analysis of the gender gap

The survey-based literature on the gender gap in cycling suggests that women are more-risk averse, which would result in a lower cycling rate than men in environments perceived as risky [12]. Following this hypothesis, we investigate the association between the gender ratio σ_c and two indicators of urban road safety, constructed using OpenStreetMap (OSM) data [25]. The first indicator (hereafter *bike lanes*) measures the proportion of streets

with cycleways (either protected or unprotected) in the street network. The second metric (hereafter *speedlimit*) provides the proportion of streets with a speed-limit equal or lower than 20 mi/h or 30 km/h. Both metrics are weighted using the length of each street. Figure 2 reports the scatter plots between the gap σ_c and the two urban road safety metrics, for the four main geographical areas separately. Each marker corresponds to a city, the black line is the linear fit. Both measures of road safety display a positive correlation with the observed gender ratio for the UK and for the area of Benelux. For cities in the United States, a positive (but weaker) correlation is only observable for the *speedlimit* indicator. For Italian cities, in contrast, both correlations are not statistically different from 0. This lack of significant correlations may be due to the presence of a large number of Italian cities with a very low degree of development of dedicated cycling infrastructure compared to cities in the UK and Benelux. Further, although some of these positive correlations appear to be driven by the presence of extreme values (Fig. S2 in the SI provides a similar picture without outliers), these are legitimate observations as the data were collected similarly for all cities. Although limited to specific geographical areas, the positive correlations suggest an association between the degree of road safety and σ_c , thus supporting the hypothesis that low levels of women engagement with cycling may be explained by a greater concern for safety compared to men.

To untangle the effect of confounding factors, we explore the relationship between σ_c and the two indicators of urban road safety controlling for a range of city-level indicators. To provide a thorough characterization of each city, the indicators are chosen from four domains: 1) *E: Environment*, such as share of population in green areas, 2) *BEI: Built-Environment & Industrialization*, such as concentration of PM 2.5, 3) *SED: Socio-Economics & Demographics*, such as GDP per person, and 4) *M: Street Morphology*, such as average street grade. A full list of indicators is provided in Table I. The correlation matrix of the indicators across the entire sample is provided in Figure S3 in the SI. We include geographical dummies for the macro areas to account for different penetration levels of Strava worldwide.

Coefficients (and 95% confidence intervals) of a linear regression model estimated via Ordinary Least Squares (OLS) are shown in Fig. 3, with statistically significant coefficients at 0.05 level (two-tailed test) pictured in purple. Information on the model selection is provided in the Materials and Methods. The regression analysis confirms the positive association between the gender ratio of cyclists and the *speedlimit* indicator. This association means that urban centers with a relatively wider low-speed zone typically present a more balanced cycling uptake between men and women, after controlling for other confounding factors. Under the assumption that a wider low-speed zone indicates a less risky environment, this result also

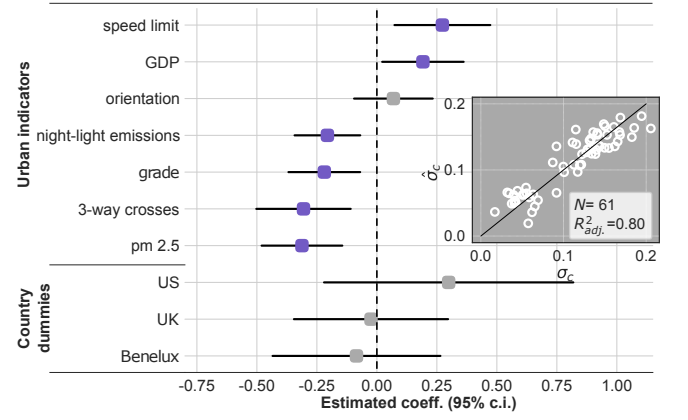


FIG. 3. **Results of regression analysis.** The main plot shows the estimated coefficients (square markers) with 95% confidence intervals (black lines) for the final set of regressors included in the model. Model estimated via Ordinary Least Square, selection performed via exhaustive search. Selection criterion: Akaike Information Criterion (AIC). Statistically significant coefficient at 0.05 significance level are pictured in purple. The scatter plot displays the observed σ_c vs the fitted σ_c .

confirms that women are more susceptible than men to the perceived level of risk of the cycling environment. Other insights emerge from the analysis of the control variables. First, we observe a negative association between σ_c and the proportion of 3-way crosses. From a topological view-point, cities with a high proportion of 3-way intersections deviate from grid-like street networks, that, by contrast, present a large prevalence of (mostly orthogonal) 4-way intersections [26]. This result can be interpreted again under the lens of the degree of safety of the urban environment for cycling. Indeed, the literature has shown that not only are crashes involving cyclists more likely to happen at non-orthogonal crosses than at right intersections, but the former are more likely to lead to severe injuries [27]. Another key urban feature relates to the morphology of the street-network. The negative association between σ_c and the *grade* indicator shows that hillier cities display a larger gender gap in recreational cycling, controlling for all other factors. This result aligns with previous findings that women would have a preference for flatter routes [16] which may indicate a structural limit in the potential for cycling uptake by women in particular urban environments. Interestingly, the analysis also indicates a lower gender ratio in cities with worse air quality (higher concentration of PM 2.5). In absence of a (quasi-)experimental setting, however, we are unable to determine whether the air quality is a relevant feature per se or if it acts as a proxy for other city-level characteristics such as motorized traffic. Finally, the results indicate a more balanced cycling uptake between men and women in relatively wealthier cities (with a larger GDP per person) and cities with a lower degree

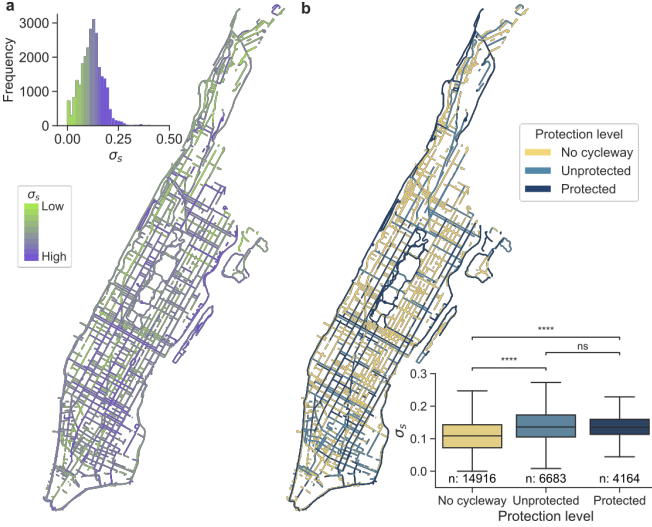


FIG. 4. Streets characteristics and σ_s : The case of New York City. a) The map displays streets in the borough of Manhattan included in the final sample. A 10-quantile color scheme has been used for the value of σ_s . The inset is the distribution of σ_s for streets included in the final sample (computed over the entire city of New York). b) The map displays the protection level of streets in the borough of Manhattan included in our sample. Yellow: no cycleway, light-blue: unprotected cycleway, dark blue: protected cycleway. The inset displays the box plots of σ_s for streets with different levels of protection (computed over the entire city of New York). 'ns', '*', '**', '***', '****' indicate the significance level of a Mann-Whitney-Wilcoxon test two-sided with Bonferroni correction, with the following p-values thresholds: $1e-4$: "****", $1e-3$: "***", $1e-2$: "**", 0.05 : "*", 1 : "ns".

of night-light emissions (which can be a proxy for the size of the city). To test the robustness of this analysis, we estimate three additional models where we adopt different strategies to account for the different levels of penetration of Strava worldwide. These strategies differ in terms of: geographical coverage, specification of the geographical dummies and standardization of the input and target variables (Fig. S4 in SI). The results are largely consistent with the preferred specification provided here.

Street-level analysis of the gender gap

The results in the previous section show that aggregated urban features model well the heterogeneity of the gender gap in cycling observed across different cities. They also confirm and provide quantitative support to traditional hypotheses from the literature, which are typically grounded on small-sample survey-based analyses. Though informative and affirmative, the previous analysis leaves open the question: Where exactly do women prefer to cycle? Also, which concrete interventions could policy makers implement to enhance cycling for women?

To answer these questions, we shift the focus from a

macro-level comparison across cities, to a micro-level setting where the unit of analysis are streets within one city as opposed to the entire city itself. This shift in perspective allows us to examine the preferences of women for street-level characteristics in greater detail, thus identifying potential targets for interventions by policy makers. Among the available cities, we select as a case study the city of New York, whose large collection of administrative datasets represents an opportunity to enrich the analysis with data not otherwise available from OSM only. In particular, using OSM data, we are able to characterize each street in our sample with information on: the presence (or absence) of a protected (or unprotected) cycleway, the presence of public lighting, the type of surface (paved vs unpaved), whether the streets are close to a park or to a coastline. The administrative data are instead used to measure the number of crashes (any type of vehicles or bicycle-related only) on the street (normalized by the street length) and to associate each street to a neighborhood. Finally, to proxy for traffic flow, we compute the normalized edge betweenness [28] of each street in the street network. The edge betweenness is a network centrality measure capturing the number of the shortest paths that go through an edge in the network. A summary of all features is provided in Table S3. As for the city-level analysis, we use Strava data on cycling to quantify female preferences for a street s . We measure the proportion of female cyclists out of all cyclists travelling via street s , and call this metric σ_s . The indicator σ_s is a direct street-level extension of σ_c – indeed σ_c can be constructed averaging over σ_s with weights equal to the product between the length of each street and the total number of cyclists on it. The larger σ_s the greater are womens' preferences to cycle on street s . Compared to a simple count of female cyclists, this relative measure has the advantage of quantifying female-specific preferences towards a street s , irrespective of the total level of *popularity* of the street. Therefore, the metrics will not be distorted towards streets that are very popular for cyclists in general (for instance for their position in the street network), but that may not present features that are particularly appreciated by our target group. In addition, we adopt a data-driven approach to filter streets with a low number of cyclists (described in the SI). This filtering ensures that the observed σ_s is computed on a sufficiently large cyclist base. The distribution of σ_s is bell-shaped with a mean around 0.12 and a range between 0.00 and 0.41 (FIG.4(a)). Stratifying the distribution by protection level of the street ('No cycleway', 'Unprotected cycleway', 'Protected cycleway') -Fig. 4(d)- we see that streets with no forms of dedicated-infrastructure are typically associated with lower σ_s than streets with either protected or unprotected cycleway: the median value of σ_s for streets with no cycleway roughly corresponding to the 25th percentile of both the distributions of streets with protected or unprotected cycleways. This

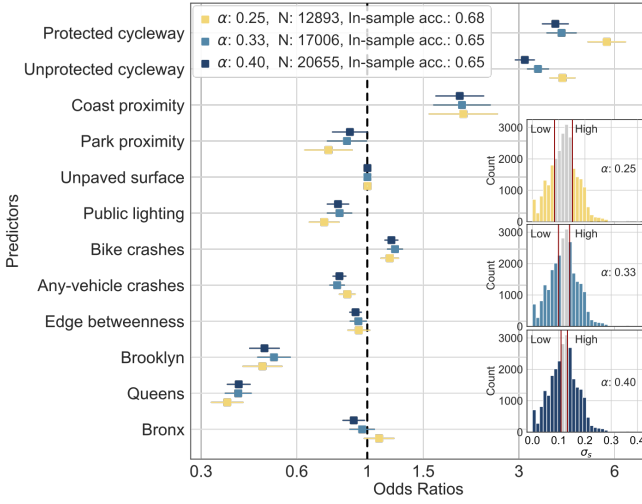


FIG. 5. **Odds Ratios of multivariate logistic regressions, for several level of the threshold α** a) the chart presents estimated ORs for a multivariate logistic regression where the target variable is the binarized σ_s and the predictors are listed in Table S3. The squared dots are the point estimates. The straight lines are the estimated 95% confidence interval for the corresponding OR. The model was estimated on three different sample selections, with the threshold α corresponding to 0.25 (yellow), 0.33 (light-blue) and 0.40 (dark blue). The ORs are computed exponentiating the corresponding estimated coefficients. For each estimated model, the legend reports the value of the threshold α , the number of observations and the in-sample accuracy. b) The histograms show the mapping between the σ_s and the binarized σ_s for the three values of the threshold α : 0.25 (yellow), 0.33 (light blue) and 0.40 (dark blue).

descriptive evidence provides a first indication that streets with some form of cycling infrastructure are typically used more intensively by women than streets with no dedicated infrastructure at all. To delve deeper into women’s preferences for dedicated cycling infrastructure, we study the degree of association between the presence of protected and unprotected cycleways and σ_s by means of a logistic regression analysis. We classify streets into two classes, *Low* and *High*, corresponding to the bottom, and top 33% of the distribution of σ_s and estimate the Odds Ratios (OR) via multivariate logistic regression. To check the robustness of the results, the analysis is repeated choosing different thresholds α (0.25 and 0.40, instead of 0.33) for the classification. The results (presented in Figure 5) are consistent across sample specifications, with generally slightly larger estimates on more extreme samples (lower values of the threshold α).

The main result pertains to the role of dedicated cycling infrastructure. With an estimated OR of around 4.08 (95% confidence interval: [3.67, 4.54]), the analysis indicates that the odds to be classified *High* are more than four times greater for protected cycleways than for streets with no cycleway (used as baseline). This result

is largely in line with the survey-based literature on the gender-cycling-gap, according to which women would favor physical separation more than men [12–14]. Though smaller in magnitude, we estimate a similarly positive association between the presence of an unprotected cycleway and σ_s . This analysis suggests that, whenever protected cycleways are not feasible due to either budget or physical constraints, the use of shared unprotected cycleways would still be a way to make the urban environment more accessible for female cyclists. In light of recent findings [29, 30] which suggest that unprotected cycleways would not enhance the degree of safety of the road-network for cyclists, our results suggest that subjective safety may matter more than objective safety. In terms of other control variables, in line with the assumption that women favor more quiet streets, we estimate an OR below 1 for our proxy for traffic-flow (*Edge-betweenness*) and for the volume of accidents (by any type of vehicle). The positive association with the volume of bicycle crashes, on the other hand, is likely to be the effect of reverse causality: a more balanced gender ratio is typically associated with a larger volume of cyclists, with an increased likelihood of bicycle crashes. The two dummies on coast and park proximity, here inserted as a proxy for the natural environment within which the street is located, appear to have opposite effects, with an estimated OR above 1 for *Coast proximity* and below 1 for *Park proximity* (with the latter being only statistically significant at 0.05 level for $\alpha = 0.25$). On one hand, the reason for the high coast proximity value is due to the morphology of the city of New York and the presence of a long protected cycleway along the coastline of Manhattan acting as an attractive infrastructure and impacting nearby streets too, with many cyclists riding through to reach it. On the other hand, the negative association with the *Park proximity* dummy can be traced back to the location of the green areas under consideration, often in non-central locations (note that streets within the Central Parks largely fall into the excluded part of the distribution around the median value). Information on the presence of public lighting is generally very sparse in OSM and particularly for New York City (we assume public lighting to be absent only whenever explicitly stated, with less than 100 streets classified as without public lighting), therefore the negative estimated OR requires further analysis with more complete data. Finally, although hard to generalize to other urban contexts, we observe strong negative neighbourhood effects, particularly for the boroughs of Brooklyn and Queens (compared to the baseline borough of Manhattan).

The multivariate logistic regression presented in this paragraph allowed us to investigate the average impact of specific street-level features on the probability that a street belongs to the *High* or *Low* σ_s group. The overall in-sample accuracy of the model oscillates between 0.68 and 0.65, depending on the threshold α . To check the

robustness of our analysis, we additionally compared the results obtained using the multivariate logistic regression to the results of a Random Forest classifier. The forest-based classifier provides a higher accuracy (average out-of-sample accuracy of 0.83 for $\alpha = 0.33$ compared to 0.65), at the expense of a lower degree of explainability. Nevertheless, using the game-theory inspired concept of shapely values [31], for each street, we can quantify the impact on the prediction of the considered features. In Fig. S6 in the SI, we present shapely values computed on a random selection of 500 data points, which largely confirm the results of the logistic regression in terms of the central role of the dedicated cycling infrastructure.

DISCUSSION

In this study, we investigated the determinants of the gender-cycling-gap using data for over 60 cities in Europe and the United States. Unlike the vast majority of previous analyses that used survey-based data, we leveraged large automatically collected data from the online sport-tracking application Strava. We first related female cycling rates in different European and American cities to city-level characteristics and found evidence for traditional hypotheses which link the observed gender gap in cycling to gender specific preferences on road safety. Additionally, we found higher female cycling rates in flatter than in hillier cities, also in line with the literature [16]. This is an interesting result as there may be structural, morphological or cultural [32] constraints for specific places where the cycling uptake is harder to increase for women. For urban planning this result suggests that ad-hoc infrastructural interventions such as the provision of cycleways or the enlargement of the low-speed limit zones could have limited efficacy in these contexts and may require concurrent behavioral incentives, for instance to expand the adoption of e-bikes. A novel result concerns the strong association between the gender-cycling-gap and the air quality of a city, which however requires further research within a (quasi-)experimental setting.

In the second part of the study, we shifted the focus from a macro comparison across cities to a micro-level analysis, at the level of single streets. If the first analysis successfully provided evidence for and expanded existing hypotheses (further validating our data as a reliable source on cycling behavior), the second aims at capturing the role of urban features modelled at a higher resolution and delving deeper into the association between the gender-cycling-gap and the presence of dedicated cycling infrastructure. We selected the city of New York as case study for this component of the study. Using multivariate logistic regression analysis, we have shown the existence of a positive association between the volume of female cyclists (relative to men) and the presence of dedicated cycling infrastructure. The positive association between

σ_s and the presence of a protected cycleway was expected and well-documented in the literature, which highlights the strong preference of women for physical separation from motorized traffic [12–14]. More novel and interesting is the observed association with the presence of an unprotected cycleway. In light of recent studies showing that unprotected cycleways may not enhance the degree of objective road safety [29, 30], our result suggests that the perceived degree of safety may induce women to cycle more than the actual degree of safety. Therefore, in contexts where no physical separation is possible (for instance for space or budget constraints), the provision of shared cycleway may still act as a way to make to urban environment perceived as more accessible by women. However, given that the increase in the perception of safety induced by this type of infrastructure may not always translate into a lower risk cycling environment, the planning of this type of infrastructure should be evaluated carefully by city planners, for instance favoring specific solutions associated to greater safety levels.

Overall our study validated survey-based results quantitatively using unprecedentedly large-scale automatically collected data. With around 36 million users worldwide in 2018, Strava was among the major applications for sport tracking and as such, a reliable information on cycling-behavior for regular cyclists. The main limitation of our study pertains to the representativeness of Strava users and the purposes of Strava trips. For example, having a considerable gender gap in the Netherlands (Fig. 1), contrary to expectations [11], the Strava data are clearly not representative, and neither users nor purposes of use can be clearly inferred. We therefore stress that Strava does not necessarily reflect recreational cycling only, and that such assumptions should be challenged and explored with richer data sets or qualitative methods. Nevertheless, we did our best to account for the representative challenge of this data set, first by comparing only cities in the same region, and second by comparing streets only in the same city, aiming to minimize user and trip purpose variation. A second limitation of the Strava data set is the inability to extract the potentially useful information of cyclist volumes [33], as the raw data are not individual cycling traces but Strava segments with only aggregated statistics. This aggregation also implies that the same cyclists may cycle on many segments in one or multiple sessions and we would not be able to identify them.

It is unclear to which extent our results can be generalized to cycling for purposes other than recreational, such as transport, and to less-skilled cyclists (occasional cyclists and not cyclists). It is therefore important to find data sources that are able to reliably distinguish between such purposes and users, since gender-based constraints can differ between these categories [34]. However, since the survey-based academic literature on gender-cycling-gap indicates that cycling preferences differ less among regular cyclists than among occasional ones [12, 35], the

results of our analysis could be interpreted as a lower-bound and it is likely that the identified factors play an even larger role in explaining the gender-cycling-gap in the general population. Another limitation concerns the cross-sectional nature of the available cycling data. The absence of a longitudinal dimension limited the extent to which temporal variations could be analysed in the data, hindering the use of policy-evaluation statistical tools such as diff-in-diff techniques to evaluate casual effects along with correlations.

Finally, there is a variety of gender-specific constraints apart from street safety that future studies should take into account, from cultural and psychological reasons [32, 33], to other environmental factors and harassment by motorists [34, 36]. Gender inequality and gendered transport habits may also play a large role, such as more frequent trip chaining by women due to childcare and other errands [37, 38]. Therefore, while street safety and urban design are undoubtedly important ingredients, there is no universal, simple fix for getting rid of the gender gap in cycling towards more sustainable mobility. It remains a complex societal issue that needs to be tackled from multiple angles [7].

MATERIALS AND METHODS

Data sources The study used data from multiple sources summarised below. More detailed information on the collection, cleansing and processing of each of these strands of data is provided in the SI.

- Strava data on biking extracted from Strava heat maps using the public API in November 2018.

Cross-cities analysis

- City-indicators from the Global Human Settlement - Urban Centre Database 2015 (GHS-UCDB R2019A) [39].
- Street-network indicators for the urban centers of interests extracted from [26].
- OpenStreetMap [25] data extracted through the Python library OSMnx [40] to compute urban safety indicators measuring: 1) the proportion of the street network with bike-lanes, and 2) the proportion of streets with low speed limit. Detailed information on the extraction pipeline is provided in the SI.

Case study on the City of New York

- OpenStreetMap [25] data on street-level characteristics extracted during the process of remapping of Strava data via the python library OSMnx [40]. For each street, we retained information on: the presence of public lighting, the presence of protected or unprotected bike-lanes, proximity with a park

or with the coastline and whether the surface is paved. A full list of the OSM tags used is provided in Table S2. In addition, for streets in the largest component of the street network, we computed the edge-betweenness via the python library *graph-tool*. Streets outside the largest component of the network (i.e. excluding streets in the borough of Staten Island) were excluded from the sample.

- Administrative data from the OpenData Portal of the city of New York on location of all (any-vehicle) accidents and bike accidents [41].
- Shapefiles of the administrative boundaries of boroughs in the city of New York. Available at [41].

Ordinary Least Squares regression We estimate a linear regression model via Ordinary Least Squares (OLS) of the form:

$$\sigma_c = \sum_{j=1}^N \beta_j z_{j,c} + \epsilon_c \quad c = 1, \dots, 61 \quad (1)$$

where the list of regressors z_j includes: *speed limit*, *orientation*, *GDP*, *3-way crosses*, *night-light emissions*, *grade*, *pm2.5* plus three dummy variables for the macro area to which the city belong (US, UK, Benelux, *baseline*: Italy). Prior to undertake any analysis, we perform a z-score-transformation of all regressors. As such, the results should be interpreted in terms of standard deviations. Out of the initial 15 city-level indicators collected (provided in Table 1), the final subset of seven indicators (plus the three country-level dummies) included in the regression were selected via exhaustive search to minimize the Akaike Information Criterion (AIC) of the model. For estimation of the Linear regression, we use the OLS function of the Python library *statsmodel* [42].

Multivariate logistic regression To assess the degree of association between σ_s and the presence of cycling dedicated infrastructure, we use a multivariate logistic regression. We restrict the sample to streets belonging to the bottom and top 33% of the distribution of σ_s and classify streets in *Low* and *High* σ_s , respectively. As a robustness check, the analysis is repeated for alternative values of this threshold (25% and 40%). We use features described in Table S3 as predictors and the binarized σ_s as the target variable. Moreover, we scale continuous predictors (*Any-vehicle crashes*, *Bike crashes* and *Edge-betweenness*) using a z-score-transformation to normalize the magnitude of the estimated coefficients. For estimation of the Multivariate Logistic regression, we use the Logit function of the Python library *statsmodel* [42]. The OR for each predictor is computed exponentiating the corresponding estimated coefficient.

Classification task: models comparison We evaluate the accuracy of the multivariate Logistic Regression and the Random Forest classifier using a 80-20 train-validation

split of the sample selected using $\alpha = 0.33$. Hyperparameters setting via grid-search optimization is performed for the Random Forest classifier (Selected parameters: `'max_features': 'auto', 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 150`). A stratified 10-fold approach is then adopted to evaluate the variability in the accuracy of the two models. For this strand of the analysis, we used the scikit-learn library [43] for Python.

CODE AVAILABILITY

The Python code developed for the data analysis is available at: <https://github.com/alibatti/GenderCyclingGapUsingStrava>.

ACKNOWLEDGEMENT

We thank Ane Rahbek Vierø for helpful discussions. PB, APa, APe acknowledge partial support from Intesa Sanpaolo Innovation Center. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

-
- [1] P. Oja, S. Titze, A. Bauman, B. De Geus, P. Krenn, B. Reger-Nash, and T. Kohlberger, *Scandinavian journal of medicine & science in sports* **21**, 496 (2011).
 - [2] L.-A. Leyland, B. Spencer, N. Beale, T. Jones, and C. M. Van Reekum, *PloS one* **14**, e0211779 (2019).
 - [3] I. Avila-Palencia, A. de Nazelle, T. Cole-Hunter, D. Donaire-Gonzalez, M. Jerrett, D. A. Rodriguez, and M. J. Nieuwenhuijsen, *BMJ open* **7**, e013542 (2017).
 - [4] S. Gössling, A. Choi, K. Dekker, and D. Metzler, *Ecological Economics* **158**, 65 (2019).
 - [5] U. N. G. Assembly, a/res/70/1. Technical report, United Nations General Assembly (2015).
 - [6] J. Monk and S. Hanson, *The Professional Geographer* **34**, 11 (1982).
 - [7] S. Hanson, *Gender, Place & Culture* **17**, 5 (2010).
 - [8] K. Hosford and M. Winters, *Findings*, 10802 (2019).
 - [9] D. Funaki, *Why Don't Women Cycle?: A Case Study of Women's Perceptions of Cycling in the SOMA District of San Francisco*, Ph.D. thesis, University of California, Berkeley (2019).
 - [10] Cycling UK, Technical Report (2021).
 - [11] J. Pucher and R. Buehler, *Transport reviews* **28**, 495 (2008).
 - [12] R. Aldred, B. Elliott, J. Woodcock, and A. Goodman, *Transport reviews* **37**, 29 (2017).
 - [13] J. Garrard, G. Rose, and S. K. Lo, *Preventive medicine* **46**, 55 (2008).
 - [14] A. Misra and K. Watkins, *Transportation Research Record* **2672**, 145 (2018).
 - [15] J. Dill, T. Goddard, C. Monsere, and N. McNeil, in *Transportation Research Board 94th Annual Meeting* (2014).
 - [16] J. Hood, E. Sall, and B. Charlton, *Transportation letters* **3**, 63 (2011).
 - [17] G. Akar, N. Fischer, and M. Namgung, *International Journal of Sustainable Transportation* **7**, 347 (2013).
 - [18] E. Heinen, K. Maat, and B. van Wee, *Transportation* **40**, 23 (2013).
 - [19] K. Wang and G. Akar, *Journal of Transport Geography* **76**, 1 (2019).
 - [20] E. Eren and V. E. Uz, *Sustainable Cities and Society* **54**, 101882 (2020).
 - [21] Y. Sun and A. Mobasheri, *International journal of environmental research and public health* **14**, 274 (2017).
 - [22] W. Musakwa and K. M. Selala, *Data in brief* **9**, 898 (2016).
 - [23] Strava, "Year in sport 2018," Available at <https://blog.strava.com/press/2018-year-in-sport/> (2018).
 - [24] Strava, "Year in sport 2019," Available at <https://blog.strava.com/press/strava-releases-2019-year-in-sport-data-report/> (2019).
 - [25] OpenStreetMap contributors, "Planet dump,".
 - [26] G. Boeing, *Geographical Analysis* (2021).
 - [27] M. Asgarzadeh, S. Verma, R. A. Mekary, T. K. Courtney, and D. C. Christiani, *Injury prevention* **23**, 179 (2017).
 - [28] V. Latora, V. Nicosia, and G. Russo, *Complex Networks: Principles, Methods and Applications*, Complex Networks: Principles, Methods and Applications (Cambridge University Press, 2017).
 - [29] L. Pearson, J. Dipnall, B. Gabbe, S. Braaf, S. White, M. Backhouse, and B. Beck, *Journal of Transport & Health* **24**, 101290 (2022).
 - [30] W. E. Marshall and N. N. Ferenchak, *Journal of Transport & Health* **13**, 100539 (2019).
 - [31] S. M. Lundberg and S.-I. Lee, in *Proceedings of the 31st international conference on neural information processing systems* (2017) pp. 4768–4777.
 - [32] T. Goddard and J. Dill, in *Transportation Research Board 93rd Annual Meeting. Washington, DC* (2014).
 - [33] R. Aldred, J. Woodcock, and A. Goodman, *Transport reviews* **36**, 28 (2016).
 - [34] K. C. Heesch, S. Sahlqvist, and J. Garrard, *International Journal of Behavioral Nutrition and Physical Activity* **9**, 1 (2012).
 - [35] G. Prati, F. Fraboni, M. De Angelis, L. Pietrantonio, D. Johnson, and J. Shires, *Journal of transport geography* **78**, 1 (2019).
 - [36] M. Graystone, R. Mitra, and P. M. Hess, *Transportation Research Part D: Transport and Environment* **105**, 103237 (2022).
 - [37] G. Prati, *Journal of transport geography* **66**, 369 (2018).
 - [38] J. Garrard, S. Handy, and J. Dill, *City cycling* **2012**, 211 (2012).
 - [39] A. Florczyk, M. Melchiorri, C. Corbane, M. Schiavina, M. Maffeni, M. Pesaresi, P. Politis, S. Sabo, S. Freire, D. Ehrlich, *et al.*, Public Release (2019).
 - [40] G. Boeing, *Computers, Environment and Urban Systems* **65**, 126 (2017).
 - [41] "Administrative data from the new york city council," <https://opendata.cityofnewyork.us/> (2020), [Online; accessed July 2021].
 - [42] S. Seabold and J. Perktold (2010).
 - [43] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cour-

- napeau, M. Brucher, M. Perrot, and E. Duchesnay, Journal of Machine Learning Research **12**, 2825 (2011).
- [44] W. E. Forum, Insight Report (2020).
- [45] T. P. Peixoto, “The graph-tool python library,” (2014).

TABLE I. Description and source of indicators at the level of urban center included in the study, by category. ^a

Category	Variable name	Description	Data source
E	σ_c	Proportion of kilometers rode by female cyclists to the overall kilometers rode by any cyclist within the urban area	Strava*
	share green	Share of population living in the high green area in 2015 in the Urban Centre of 2015. Ranging between 0-1	[39]
	open space	Percentage of open-spaces within the spatial domain of the Urban Centre. Ranging between 0-100	[39]
BEI	built area	Amount of the built-up area per person in 2015 calculated within the spatial domain of the Urban Centre. Expressed in square meters per person	[39]
	light emissions	Average night time night-light emission calculated within the Urban Centre spatial domain. Expressed in nano-watt per steradian per square centimetre	[39]
	pm2.5	Total concentration of PM2.5 for reference epoch 2014, calculated over the Urban Centre. Expressed in $\mu g/m^3$	[39]
SED	area	Area of the spatial domain of the Urban Centre. Expressed in square meters	[39]
	population	Population density within the spatial domain of the Urban Centre	[39] *
	GDP	GDP per capita for year 2015 within the Urban Centre. Expressed in US dollars	[39] *
M	degree	Average node degree of street network within the spatial domain of the Urban Centre	[26]
	grade	Average absolute inclination of streets within the spatial domain of the Urban Centre. Expressed in percentage	[26]
	orientation	Orientation order of street network bearings within the spatial domain of the Urban Centre.	[26]
RS	3-way crosses	Proportion of nodes that represent a 3-ways street intersection in the street network within the spatial domain of the Urban Area. Ranging between 0-1	[26]
	straightness	Ratio of straightline distances to street lengths for streets in the street network within the spatial domain of the Urban Area	[26]
	bike lanes	Proportion of streets with cycleways (either protected or unprotected) computed on streets within the spatial domain of Urban Centre	OSM*
	speed limit	Proportion of streets with a speed-limit equal or lower than 20 <i>mi/h</i> or 30 <i>km/h</i> computed on streets within the spatial domain of Urban Centre	OSM*

^a Categories: E: Environment, BEI: Built-Environment and Industrialization, SED: Socio-Economics and Demographics, M: Street Network Morphology, RS: Road Safety. *Indicates that the data from the original data sources required specific processing described in the SI.

TABLE II. List of street-level urban features included in the New York City case study.^a

Variable name	Description	Data source
Unprotected cycleway	Dummy for presence of a shared or unprotected bike-lane	OSM
Protected cycleway	Dummy for either the presence of a protected bike-lane or streets with no vehicles	OSM
Public lighting ¹	Dummy for the presence of public lighting	OSM
Unpaved surface	Dummy for unpaved surface	OSM
Park proximity	Dummy for streets next to a park (within 15 meters)	OSM*
Any-vehicle crashes	Number of crashes involving any type of vehicles per 10m of street length	NYC*
Bike crashes	Number of bicycle crashes per 10m of street length	NYC*
{borough name}	Dummy for boroughs (baseline: Manhattan)	NYC
Coast proximity	Dummy for street next to the river coast	OSM*
Edge betweenness	Edge betweenness of the streets computed for streets within the largest component of the street network	

^a ¹ Information on the presence of public lighting is very sparse in OSM. In case of missing information, we assumed that the public lighting is available. *Indicates that the data from the original data sources required specific processing (e.g. for normalization) described in the SI.

Supplemental Material:

Revealing the determinants of gender inequality in urban cycling with large-scale data

I. STRAVA DATA ON RECREATIONAL CYCLING: DATA COLLECTION AND PROCESSING

A. Data collection

Raw Strava data consist of a collection of Strava segments for 62 cities located in four geographical areas: the United States, United Kingdom, Benelux (Belgium, Netherlands and Luxembourg) and Italy. For the sensitivity analysis only, the dataset was extended to include 8 additional cities across other European countries. A Strava segment is a single portion of a road or a trail upon which Strava users compete by recording their times. The performance of a user exercising upon a segment is automatically recorded into its leader board, which in turn provides a picture of the characteristics of users exercising on a specific trail. Each raw data record consists of geographic information in the form of a linestring of lat-long coordinates plus two statistics extracted from the associated leaderboard:

1. the total number of unique cyclists training on the segment. This information corresponds to the sum of the length of the female and male leader boards, (it should be noted that -irrespective of the number of training performed on the segment- each cyclist is only included once in the corresponding leader board, according to their best performance on the segment);
2. the gender split of users training on the segment, in terms of the length of the female and male leader boards respectively (corresponding to the number of unique female and male cyclists training on the segment).

The data collection comprised of two phases, both undertaken in November 2018. In the first phase, we collected the whole corpus of segments (16.4 million) available at the time through the Strava API. This step provided us with a series of information related to each segment, in particular its unique ID and the linestring geometry. The second phase, consisted in the collection of data from the female and male leader boards (with two separate queries for the same ID) associated with each segment. In particular, given a city, we made queries for the leader boards of all the segments whose geometry is contained for at least 75% within the city boundary (extracted from OpenStreetMap). The information from the leader board was then processed to extract the statistics described above.

B. Characteristics of raw Strava segments

Strava segments are not pre-defined by the app developers but they are directly generated by the users of the app. The result of this user-driven generating process is a set of segments highly heterogeneous in length, both across cities and within the same city. Indeed, some segments may correspond to portions of a street, while others define long trails spanning multiple streets. Furthermore, segments can partially or completely overlap each others. An illustration of the extent of this characteristic is provided in FIG. S1, where we present a visualisation of raw Strava segments for nine cities. As segments are plotted with the same color intensity, darker areas on the map indicate a series of overlapping segments. Furthermore, TABLE S4 provides summary information on the raw data for the cities included in the final sample for the study. In terms of length of segments, the table depicts well the discussed heterogeneity both across cities and within the same urban center. An illustration of the large heterogeneity for segments within the same city is provided by the city of London, for which the length of segments (expressed in kilometers) spans from 0.00 to 104.97, with a mean of 2.61. By contrast, the range is much smaller (less than 7 km) for the cities of Apeldoorn and Luxembourg in the macro area of Benelux. It is noteworthy that the information on the two cities of Rotterdam and The Hague is here presented separately. However, in the study these urban centers are analysed together to match the definition of the Global Human Settlement - Urban Centers Database [39].

C. Remapping of information from Strava segments to streets in the street-network of each city

To identify the gender of cyclists travelling upon the streets network of each city, the collection of Strava segments for each city c was re-projected on the corresponding street-network following the pipeline described below.

1. Load the Strava data for city c .

TABLE S1. OpenStreetMap key-value pairs for the classification of cycleways in protected and unprotected.

Type	OpenStreetMap key	OpenStreetMap value
Protected cycleway	<i>highway</i>	<i>['cycleway', 'path', 'footway', 'bridleway', 'track']</i>
Protected cycleway	<i>cycleway</i>	<i>['track', 'opposite_track']</i>
Protected cycleway	<i>cycleway : left</i>	<i>['track', 'opposite_track']</i>
Protected cycleway	<i>cycleway : right</i>	<i>['track', 'opposite_track']</i>
Protected cycleway	<i>cycleway : both</i>	<i>['track', 'opposite_track']</i>
Protected cycleway	<i>bicycle</i>	<i>['designated']</i>
Unprotected cycleway	<i>cycleway</i>	<i>['lane', 'opposite_lane', 'share_busway', 'shared_lane', 'designated', 'yes']</i>
Unprotected cycleway	<i>cycleway : left</i>	<i>['lane', 'opposite_lane', 'share_busway', 'shared_lane', 'designated', 'yes']</i>
Unprotected cycleway	<i>cycleway : right</i>	<i>['lane', 'opposite_lane', 'share_busway', 'shared_lane', 'designated', 'yes']</i>
Unprotected cycleway	<i>cycleway : both</i>	<i>['lane', 'opposite_lane', 'share_busway', 'shared_lane', 'designated', 'yes']</i>
Unprotected cycleway	<i>bicycle</i>	<i>['yes', 'permissive', 'destination', 'private']</i>

2. Extract the bounding box of city c from the GHS-UCDB R2019A [39].
3. From OpenStreetMap, extract the street-network within the polygon defined in the bounding box using the OSMnx library [40]. Set: *network_type = 'bike', retain_all = True, simplify = True*.
4. Classify streets in the street network based on OpenStreetMap attributes in: ‘street with protected bike lane’, ‘street with unprotected bike lane’ and ‘street with no cycleway’. The (*key, value*) pairs for the classification are provided in the TABLE S1. All other bikable streets are classified as ‘no cycleway’.
5. Proceed with the *preferential assignment* of Strava segments as follows. Buffer with a 10 meters radius the geometries of the street network. Select all streets categorized as ‘protected cycleways’ and intersect each Strava segment with the network. Re-project each segment (or portion(s) of a segment) on all streets with an intersection of at least 30 meters. Finally compute the geometries of Strava segments left unassigned - that could be either a full segment or portion(s) of a segment- and repeat the procedure selecting ‘unprotected cycleways’ first and subsequently streets with ‘no cycleways’.
6. Compute the gender ratio of each street in the street network using statistics from the re-projected Strava segments. In particular, letting I be the set of segments re-projected to street s , $Females_i$ ($Males_i$) the number of unique female (male) cyclists on segment i , the total number of female cyclists on streets s (and correspondingly for male cyclists) is defined as:

$$Females_s = \sum_{i \in I} Females_i \quad (S1)$$

The gender ratio (σ_s) of street s is then computed as:

$$\sigma_s = \frac{Females_s}{Males_s + Females_s} = \frac{\sum_{i \in I} Females_i}{\sum_{i \in I} Females_i + Males_i} \quad (S2)$$

The rationale for the *preferential assignment* is that if a cycleway runs parallel to a street with no cycleway and the linestring geometry for the Strava segment is compatible with both streets (i.e. it falls within the buffered geometry of both streets), we assume that the cyclists rode on the cycleway as opposed to the street with no cycling-dedicated infrastructure. This approach prevents us from remapping the same portion of a Strava segment to multiple parallel streets with different characteristics.

D. Construction of city-level index of gender gap in recreational cycling

Strava data remapped on the street-network of each urban area were then used to construct an index of the gender-cycling-gap for all cities included in our sample. The gender-cycling-gap of city c is measured by σ_c defined as the ratio between the total kilometers travelled by female cyclists and the overall kilometers travelled by cyclists of both gender within the urban area. The rationale for the use of this metric is its ability to capture two forms of gender gaps described in the literature on cycling and gender, i.e. the propensity of women to make less trips than men and

the propensity to cycle shorter distance. This measure is equivalent to the weighted sum of the gender ratio on streets (σ_s) within the urban area, with weights equal to the product of the length and the total popularity (total number of cyclists) of the street. I.e., letting S be the set of streets in the street network of city C , $Females_s$ ($Males_s$) the number of female (male) cyclists on s and l_s the length of street s expressed in kilometers:

$$\sigma_c = \frac{\sum_{s \in S} Females_s * l_s}{\sum_{s \in S} (Females_s + Males_s) * l_s} = \frac{\sum_{s \in S} \sigma_s * l_s * (Females_s + Males_s)}{\sum_{s \in S} (Females_s + Males_s) * l_s} \quad (S3)$$

A value of σ_s between 0 and 0.5 indicates the presence of a positive gender gap (with more men cycling then women), while a value above 0.5 indicates a negative gender gap. In our sample, the maximum value of σ_s stands at 0.2, indicating a positive gender gap for all cities and a monotonic relationship between the gender gap and σ_c . A full ranking of cities for the four geographical areas is provided in TABLE S5.

II. UNDERSTANDING THE DETERMINANTS OF GENDER-CYCLING-GAP - A CROSS-CITIES ANALYSIS: DATA, METHODOLOGY AND SENSITIVITY ANALYSIS

A. Data sources

A full list of data sources used for this strand of the study is provided below.

- Data on recreational cycling at city-level from Strava. The data were processed following the steps described in the previous section.
- City-indicators from the Global Human Settlement - Urban Centre Database 2015 (GHS-UCDB R2019A) [39]. The following information was extracted:
 1. the share of population living in green areas: data field *SDG_A2G14*;
 2. the percentage of open space: data field *SDG_OS15MX*;
 3. the built-up area per capita, data field *BUCAP15*;
 4. the average night-light emission: data field *NTL_AV*;
 5. the concentration levels of PM2.5: data field *E_CPM2_T14*,
 6. the city area: data field *AREA*;
 7. the population density: computed as $\frac{P15}{AREA}$,
 8. the GDP per person, computed as $\frac{GDP15_SM}{P15}$.
- Street-network indicators from [26]. Out of the list of available indicators, we extract the average absolute street grade, the average degree, orientation order, the proportion of three-way intersections and the average street straightness.
- Urban safety indicators measuring the proportion of the street network with cycleways and the proportion of streets with low speed limit. These data were directly constructed from OpenStreetMap [25] information following the the pipeline in Section II B below.
- The Global Gender Gap Index (country-level) from the World Economic Forum [44]. Included in the sensitivity analysis only.

It should be noted that the final sample for this component of the analysis consists of 61 cities. The city of New York was excluded from this component of the study due to the large discrepancy between the administrative area of this city and the bounding box of the GHS, which would have made the indicators based on this definition of urban center not representative for the area actually covered by the cycling data.

TABLE S2. Definition of bike lanes for construction of city-level indicators of urban road safety. The table provides the OpenStreetMap (*key, value*) pairs used for the identification of streets with some form of cycling dedicated infrastructure, simply indicated as *bike lane* in the main text. This information was then used to measure the size of the cycling-dedicated-infrastructure in the city and construct the corresponding city-level indicator.

OpenStreetMap key	OpenStreetMap value
<i>highway</i>	<code>['cycleway']</code>
<i>cycleway</i>	<code>['track','opposite_track','lane','opposite_lane','opposite','share_busway','shared_lane','designated','yes']</code>
<i>cycleway : left</i>	<code>['track','opposite_track','lane','opposite_lane','opposite','share_busway','shared_lane','designated','yes']</code>
<i>cycleway : right</i>	<code>['track','opposite_track','lane','opposite_lane','opposite','share_busway','shared_lane','designated','yes']</code>
<i>cycleway : both</i>	<code>['track','opposite_track','lane','opposite_lane','opposite','share_busway','shared_lane','designated','yes']</code>

B. Construction of urban road safety indicators

OpenStreetMap information accessed via the Python library OSMnx [40] was used to construct the two indicators on urban road safety. The indicator on the proportion of streets with max-speed limit equal or below 20mi/h or 30km/h (simply referred to as *speed limit* in the Main) was constructed according to the pipeline described below.

1. For each city c , extract the bounding box of city c from the GHS-UCDB R2019A [39].
2. Extract the street-network from the polygon defined in the bounding box using the OSMnx library [40]. Set: `network_type = 'drive', retain_all = True`.
3. Compute the proportion of streets satisfying the condition on the speed limit. Weight each street with its length.

The indicator on the proportion of streets with cycling-dedicating infrastructure (simply referred to as *bike lanes*) was constructed according to the pipeline below.

1. For each city c , extract the bounding box of city c from the GHS-UCDB R2019A [39].
2. Extract the street-network from the polygon defined in the bounding box using the OSMnx library [40]. Set: `network_type = 'bike', retain_all = True`. Call this graph G_0 .
3. From OpenStreetMap [25], extract the street-network from the polygon defined in the bounding box using the OSMnx library. Set: `network_type = 'drive', retain_all = True`. Call this graph G_1 .
4. Define as cycleways all streets in G_0 with the pairs of OSM attribute described in TABLE S2.
5. Sum over the length of all 'bike-lanes' in G_0 .
6. Sum over the length of all streets in G_1 .
7. Define the index as the ratio between the metric computed at point 5 and the metric computed at point 6.

C. Regression analysis

We estimated a linear regression model via Ordinary Least Squares (OLS) of the form:

$$\sigma_c = \sum_{j=1}^N \beta_j z_{j,c} + \epsilon_c \quad c = 1, \dots, 61 \quad (\text{S4})$$

where the list of regressors z_j in the preferred model includes: *speed limit, orientation, GDP, 3-way crosses, night-light emissions, grade, pm2.5* plus three dummy variables for the macro area to which the city belong (US, UK, Benelux, *baseline*: Italy). All continuous regressors were normalised using a z-score transformation. Out of the initial 15 city-level indicators collected (provided in Table 1), the final subset of seven indicators (plus the three country-level dummies) included in the regression were selected via exhaustive search to minimize the Akaike Information Criterion (AIC) of the model. The model is estimated using the OLS function of the Python library *statsmodel* [42].

D. Sensitivity analysis

As sensitivity analysis, we estimated three additional regression models, in which we adopted different strategies to account for the heterogeneous levels of penetration of Strava across countries covered in the study. The three additional models share the linear formulation and the estimation technique (OLS) with the preferred one, but differ in (at least one) of the following characteristics:

1. Standardization of target and input variables on either the full sample or by geographical area. For model 2, the standardization was still performed on the full sample (as for the preferred model) for both target and input variables. For model 3, the target variable was standardized by continent ('US' vs 'Europe') while the input variables were standardized at the level of the entire sample. For model 4, both target and input variables were standardized by continent ('US' vs 'Europe').
2. Inclusion and exclusion of geographical dummies. Model 2 includes a geographical dummy for the US (with baseline: Europe). Model 3 and 4 do not include geographical dummies, but the standardization of the target and/or inputs is performed by geographical area ('US' vs 'Europe').
3. Enlargement of the sample to include 8 additional cities in Europe, located outside of the four main geographical areas covered by the study.

A characterization of each model is provided in FIG. S4-a. Model 1 corresponds to the preferred model described in the main analysis. Variables selection for each additional model was performed via extensive search to minimize the Akaike Information Criterion (AIC) of the model. Scatter plots of the actual values vs the fitted values of σ_c for each model are presented in FIG. S4-a, suggesting that overall all four models explain well the variability in the observed σ_c . Figure S4-b shows the estimated coefficients of the preferred model and the three additional models. Matrix cells outlined in black indicate that the corresponding coefficient is statistically significant at 0.05 level. The sensitivity analysis highlights that the results discussed in the main text concerning the role of the speed limit, average hilliness and complexity of the street network are robust across the different specifications, both in terms of the sign of the effect of each regressor and its significance.

III. CASE STUDY ON THE CITY OF NEW YORK: DATA AND METHODOLOGY

A. Data sources

A full list of data sources used for this component of the study is provided below.

- Data on recreational cycling at street-level for the city of New York from Strava. The raw Strava data were processed and remapped to the street-network of each city extracted from OpenStreetMap following the steps described in Section I. A network definition of streets was used, which does not reflect a the toponymy of streets.
- OpenStreetMap data on street-level characteristics extracted during the process of remapping of Strava data via the python library OSMnx [40]. In particular, for each street, we retained information on: the presence of public lighting, the presence of protected or unprotected cycleways, proximity with a park or with the coastline and whether the surface is paved. A list of OSM key-value pairs is provided in TABLE S3. In addition, for streets in the largest component of the street network, we computed the edge-betweenness [28] via the Python library *graph-tool* [45]. Streets outside the largest component of the network (i.e. excluding streets in the borough of Staten Island) were excluded from the sample.
- Administrative data from the OpenData Portal of the city of New York on location of all (any-vehicle) accidents and bike accidents only. The data of car and bike accidents were processed to compute the number of car and bike accidents per 10 meters, for each street. The raw data are available at [41].
- Shapefiles of the administrative boundaries of boroughs in the city of New York. Available at [41].

B. Methodology

Data filtering The vast majority of Strava users are men. Because of this, the probability of observing $\sigma_s = 0$ is a decreasing function of the number of cyclists on the street. As such, observing no women on a road may be due to two factors. First, it may be the case that the street is not (overall) sufficiently popular: being women underrepresented among Strava users, there will not be female cyclists on it. On the other hand, the segment might be overall sufficiently popular, but it's attractiveness being low among women compared to men. These effects are hard to disentangle, limiting our ability to interpret an extreme values of σ_s , on segments with low overall popularity. To mitigate this issue, we filter the data to exclude those segments with a small number of cyclists. In particular, for a city c , the probability of observing $\sigma_s = 0$ on a segment s conditional to observing N_s cyclists on s is given by (assuming replacement):

$$P_N = P(\sigma_s = 0|N_s) = \left(\frac{Males_c}{Males_c + Females_c} \right)^{N_s} \quad (S5)$$

which is decreasing in N_s . To ensure that there is a weak dependence between the popularity of the segment and the observed gender ratio, we choose a filtering threshold on N_s such that we retain only those segments s for which $P_{N_s} - P_{N_s-1}$ is low(we fixed $P_{N_s} - P_{N_s-1} < 0.015$). The impact of this filtering-rule on the distribution of gender ratio for the city of New York is depicted in FIG. S5. The top panel shows P_{N_s} as a function of N_s for the City of New York, together with the vertical line which denotes the selected threshold according to the aforementioned criterion. The bottom panels show the distribution of gender ratio before (left) and after (right) the filtering. As expected, the observed zero inflation of the distribution is much less marked after the filtering, confirming the effectiveness of the adopted mitigation strategy. In addition, the new distribution appears much less right-skewed.

Multivariate logistic regression To assess the degree of association between σ_s and the presence of cycling-dedicated infrastructure, we use a multivariate logistic regression. We restrict the sample to streets belonging to the bottom and top 33% of the distribution fo σ_s and classify streets in *Low* and *High* σ_s respectively. As a robustness check, the analysis is repeated for alternative values of this threshold (0.25 and 0.40, instead of 0.33). We use features described in Table S3 as predictors and the binarized σ_s as the target variable. Moreover, we scale continuous predictors (*Any-vehicle crashes*, *Bike crashes* and *Edge-betweenness*) using a z-score-transformation to normalize the magnitude of the estimated coefficients. For estimation of the Multivariate Logistic regression, we use the Logit function of the Python library *statsmodel* [42].



FIG. S1. Visualisation of raw Strava segments for nine cities. Graphical visualisation of raw Strava segments for nine cities for the four main geographical areas covered by the study (United Kingdom: Nottingham, Bristol and Manchester; Benelux: Amsterdam and Liege, Italy: Rome and Turin, United States of America: Boston and Memphis). All segments are plotted with the same color intensity. Darker lines indicate the presence of overlapping segments.

TABLE S3. List of street-level urban features included in the New York City case study.^a

Variable name	Description	Data source	OSM key-values pairs
Unprotected cycleway	Dummy for presence of a shared or unprotected bike-lane	OSM	{'cycleway:right': ['lane', 'opposite_lane', 'share_busway', 'shared_lane', 'designated', 'yes'], 'cycleway:left': ['lane', 'opposite_lane', 'share_busway', 'shared_lane', 'designated', 'yes'], 'cycleway:both': ['lane', 'opposite_lane', 'share_busway', 'shared_lane', 'designated', 'yes'], 'cycleway': ['lane', 'opposite_lane', 'opposite', 'share_busway', 'shared_lane', 'designated', 'yes'], 'bicycle': ['yes', 'permissive', 'destination', 'private']}
Protected cycleway	Dummy for either the presence of a protected bike-lane or streets with no vehicles	OSM	{'highway': ['cycleway', 'path', 'footway', 'bridleway', 'track'], 'cycleway': ['track', 'opposite_track'], 'cycleway:left': ['track', 'opposite_track'], 'cycleway:right': ['track', 'opposite_track'], 'cycleway:both': ['track', 'opposite_track'], 'bicycle': ['designated']}
Public lighting	Dummy for the presence of public lighting	OSM	{'lit': ['yes', '24/7']} or missing tag
Unpaved surface	Dummy for unpaved surface	OSM	{'surface': ['unpaved', 'compacted', 'fine_gravel', 'gravel', 'pebblestone', 'ground', 'earth', 'dirt', 'grass', 'grass_paver', 'sand', 'mud']}
Park proximity	Dummy for streets next to a park (within 15 meters)	OSM*	{'leisure': 'park': } + 15m proximity from the geometry
Any-vehicle crashes	Number of crashes involving any type of vehicles per 10m of street length	NYC*	Not applicable
Bike crashes	Number of bike crashes per 10m of street length	NYC*	Not applicable
{borough.name}	Dummy for boroughs (baseline: Manhattan)	NYC	Not applicable
Coast proximity	Dummy for street next to the river coast	OSM*	'natural': 'coastline' + 150m proximity from the geometry
Edge betweenness	Edge betweenness of the streets computed for streets within the largest component of the street network	OSM*	OSM street-network used to compute edge-betweenness on the largest component of the network. Network extracted using OSMNX, network_type='bike'.

^a *Indicates that the data from the original data sources required specific processing (e.g. for normalisation).

TABLE S4. Geographic characteristics of raw Strava segments by city. For each city, the table reports the total number of Strava segments, the length of the shortest segment in the data collection in kilometers (Min (km)), the length of the longest segment in the data collection in kilometers (Max (km)), the mean value of the distribution of the lengths of segments (Mean (km)) and its standard deviation (Std (km)), expressed in kilometers. All figures are rounded to the nearest 0.01.

	City name	Country	N seg.	Min (km)	Max (km)	Mean (km)	Std (km)
0	Albuquerque	USA	1997	0.25	126.51	6.88	8.97
1	Almere	Benelux	818	0.08	24.03	2.31	3.13
2	Amsterdam	Benelux	2091	0.00	23.39	1.64	2.13
3	Antwerp	Benelux	1014	0.04	23.96	1.67	2.68
4	Apeldoorn	Benelux	357	0.03	7.12	1.14	1.04
5	Arnhem	Benelux	1608	0.00	19.75	1.24	1.79
6	Austin	USA	5310	0.07	143.23	7.16	12.00
7	Bari	Italy	129	0.03	14.50	2.25	2.63
8	Bologna	Italy	653	0.10	44.52	2.85	4.68
9	Boston	USA	1400	0.19	53.87	4.45	5.96
10	Breda	Benelux	457	0.07	18.42	1.64	2.24
11	Brescia	Italy	665	0.03	21.19	2.20	2.76
12	Bristol	Uk	3221	0.00	40.34	1.61	2.87
13	Catania	Italy	271	0.14	49.78	3.95	6.40
14	Charleroi	Benelux	355	0.07	13.51	1.12	1.53
15	Charlotte	USA	1664	0.05	136.31	7.50	13.18
16	Derby	UK	1381	0.04	18.27	1.34	2.04
17	Eindhoven	Benelux	779	0.03	18.11	1.27	1.57
18	Exeter	UK	1674	0.00	24.83	1.70	2.94
19	Florence	Italy	689	0.00	24.50	1.53	1.82
20	Genoa	Italy	2107	0.00	35.92	2.35	2.88
21	Ghent	Benelux	1170	0.04	20.71	1.26	1.97
22	Groningen	Benelux	1279	0.00	11.34	1.45	1.51
23	Haarlem	Benelux	220	0.01	10.78	1.11	1.29
24	Jacksonville	USA	1032	0.36	160.42	10.54	16.88
25	Las Vegas	USA	1948	0.01	102.45	7.63	9.40
26	Leeds	UK	9991	0.00	82.89	2.90	5.87
27	Liege	Benelux	770	0.08	14.62	1.19	1.20
28	London	UK	18232	0.00	104.97	2.61	4.69
29	Louisville	USA	1833	0.19	208.30	9.15	14.50
30	Luxembourg	Benelux	606	0.06	7.36	1.14	0.93
31	Manchester	UK	2812	0.02	23.17	1.65	2.47
32	Memphis	USA	709	0.06	102.86	9.68	14.41
33	Modena	Italy	191	0.13	17.24	3.27	3.52
34	Nashville	USA	2186	0.13	168.36	7.78	12.32
35	Newcastle upon Tyne	UK	2085	0.03	17.23	1.51	2.21
36	New York	USA	7122	0.00	256.20	11.13	19.44
37	Nijmegen	Benelux	551	0.10	5.75	0.93	0.86
38	Norwich	UK	1065	0.00	20.73	1.14	1.82
39	Nottingham	UK	2168	0.02	22.73	1.43	2.30
40	Oklahomacity	USA	1502	0.23	277.57	16.75	28.65
41	Padua	Italy	130	0.20	10.47	1.58	1.40
42	Palermo	Italy	929	0.10	57.93	3.15	4.23
43	Parma	Italy	155	0.18	17.93	3.77	3.23
44	Phoenix	USA	6564	0.05	219.64	11.22	17.39
45	Plymouth	UK	3503	0.03	47.29	1.63	2.70
46	Prato	Italy	288	0.17	15.19	1.94	2.02
47	Reading	UK	830	0.08	16.75	1.09	1.51
48	Reggio Emilia	Italy	178	0.06	28.69	3.55	3.91
49	Rome	Italy	3257	0.00	52.60	3.00	5.11
50	Rotterdam	Benelux	1871	0.00	20.35	1.76	2.42
51	San Antonio	USA	3277	0.11	167.88	8.68	14.84
52	Sheffield	UK	7733	0.00	88.06	2.50	5.11
53	Southampton	UK	1495	0.07	17.03	1.34	1.68
54	Taranto	Italy	152	0.07	11.87	1.95	2.25
55	The Hague	Benelux	1018	0.02	23.95	1.27	1.90
56	Tilburg	Benelux	525	0.05	24.37	1.62	2.41
57	Trieste	Italy	1092	0.00	19.19	1.92	2.22
58	Turin	Italy	1418	0.00	19.27	2.06	2.10
59	Utrecht	Benelux	1221	0.06	50.08	2.05	4.15
60	Venice	Italy	341	0.15	75.69	3.31	7.39
61	Verona	Italy	1154	0.07	29.3	2.19	2.43
62	York	UK	1711	0.03	52.88	2.12	3.45

TABLE S5. Ranking of cities by σ_c , by geographical area. The table reports the ranking of cities for the four geographical areas. All figures are rounded to the nearest 0.01.

Ranking	City name	Geographical area	Country	σ_c
1	Groningen	Benelux	Netherlands	0.21
2	Utrecht	Benelux	Netherlands	0.17
3	Amsterdam	Benelux	Netherlands	0.16
4	Apeldoorn	Benelux	Netherlands	0.15
5	Nijmegen	Benelux	Netherlands	0.15
6	Haarlem	Benelux	Netherlands	0.14
7	Almere	Benelux	Netherlands	0.13
8	Eindhoven	Benelux	Netherlands	0.13
9	Arnhem	Benelux	Netherlands	0.13
10	Breda	Benelux	Netherlands	0.13
11	Tilburg	Benelux	Netherlands	0.12
12	Rotterdam/The Hague	Benelux	Netherlands	0.12
13	Ghent	Benelux	Belgium	0.10
14	Luxembourg city	Benelux	Luxembourg	0.09
15	Antwerp	Benelux	Belgium	0.09
16	Liege	Benelux	Belgium	0.06
17	Charleroi	Benelux	Belgium	0.06
1	Venice	Italy	Italy	0.09
2	Padua	Italy	Italy	0.07
3	Verona	Italy	Italy	0.07
4	Trieste	Italy	Italy	0.06
5	Florence	Italy	Italy	0.06
6	Parma	Italy	Italy	0.06
7	Turin	Italy	Italy	0.05
8	Palermo	Italy	Italy	0.05
9	Modena	Italy	Italy	0.05
10	Bologna	Italy	Italy	0.04
11	Bari	Italy	Italy	0.04
12	Brescia	Italy	Italy	0.04
13	Genoa	Italy	Italy	0.04
14	Reggio Emilia	Italy	Italy	0.04
15	Prato	Italy	Italy	0.04
16	Catania	Italy	Italy	0.03
17	Rome	Italy	Italy	0.03
18	Taranto	Italy	Italy	0.02
1	Exeter	UK	UK	0.17
2	York	UK	UK	0.15
3	London	UK	UK	0.14
4	Bristol	UK	UK	0.14
5	Nottingham	UK	UK	0.14
6	Norwich	UK	UK	0.14
7	Newcastle upon Tyne	UK	UK	0.13
8	Manchester	UK	UK	0.13
9	Derby	UK	UK	0.12
10	Leeds	UK	UK	0.12
11	Southampton	UK	UK	0.12
12	Plymouth	UK	UK	0.11
13	Reading	UK	UK	0.11
14	Sheffield	UK	UK	0.10
1	Albuquerque	US	US	0.19
2	Oklahomacity	US	US	0.19
3	Boston	US	US	0.18
4	Jacksonville	US	US	0.17
5	Las Vegas	US	US	0.17
6	San Antonio	US	US	0.16
7	Nashville	US	US	0.15
8	Austin	US	US	0.15
9	Memphis	US	US	0.14
10	Louisville	US	US	0.14
11	Charlotte	US	US	0.13
12	Phoenix	US	US	0.11

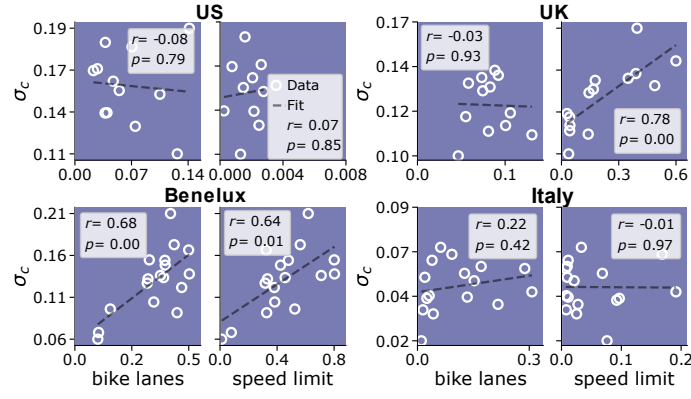


FIG. S2. **Correlations between gender ratio and urban road safety indicators, excluding outliers.** The scatter plots show the correlations between two urban road safety indicators and the gender ratio σ_c , for cities in the four geographical areas separately. For each area, outliers to the three distributions of σ_c , *bike lanes* and *speed limit* were identified using the IQR Score method and excluded. Each data point represents a city. The black line is the linear fit. The two urban road safety indicators capture the density of streets with a bike-lane in the street-network (indicator: *bike lanes*) and the density of streets with a speed limit up to 20 mi/h or 30 km/h in the street network (indicator: *speed limit*).

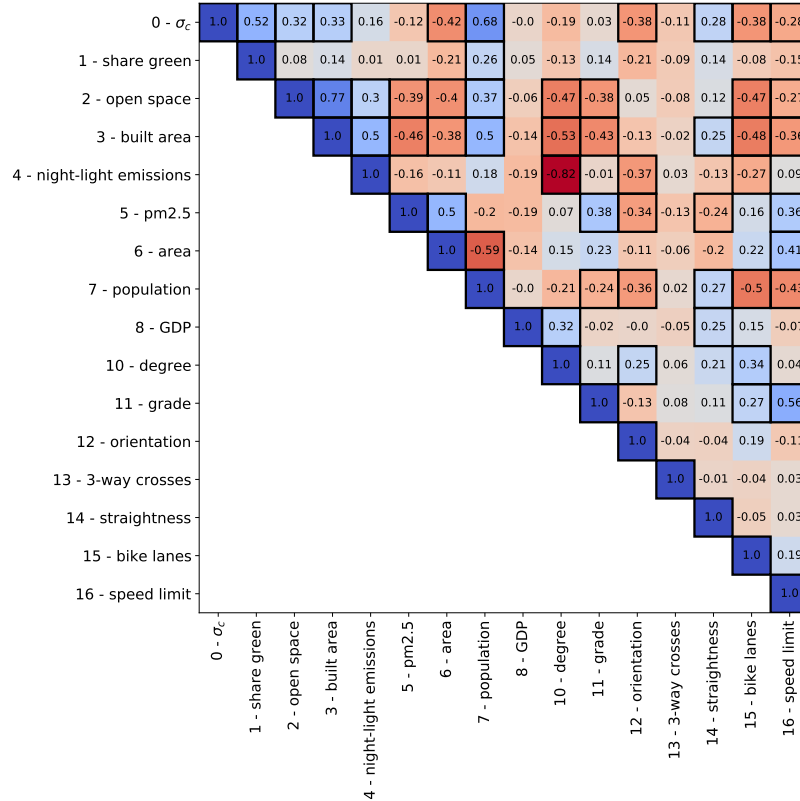
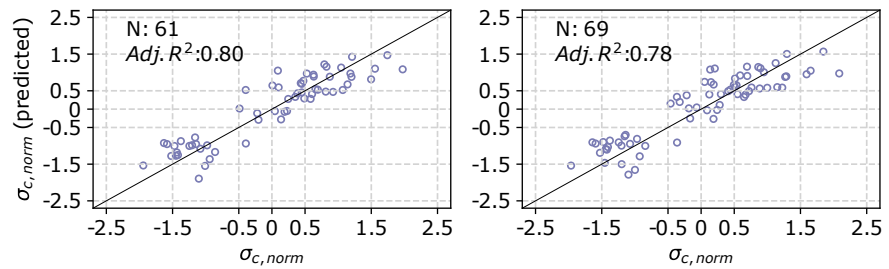


FIG. S3. **Correlations between gender ratio and city-level features.** Correlation matrix among city-level characteristics. The correlations are computed for the entire sample of 61 cities on z-scored transformed variables. Cells outlined in black indicate that the correlation is statistically different from 0 at a significance level of $\alpha = 0.05$.

a

**Model :****Geo coverage:****Stand. target, geo:****Stand. inputs, geo:****Geo dummies:**

Main
 ['US', 'UK', 'Italy', 'Benelux']
 entire sample
 entire sample
 Yes

2
 ['US', 'Europe']
 entire sample
 entire sample
 Yes

Model :**Geo coverage:****Stand. target, geo:****Stand. inputs, geo:****Geo dummies:**

3
 ['US', 'Europe']
 ['US', 'Europe']
 entire sample
 No

4
 ['US', 'Europe']
 ['US', 'Europe']
 ['US', 'Europe']
 No

b

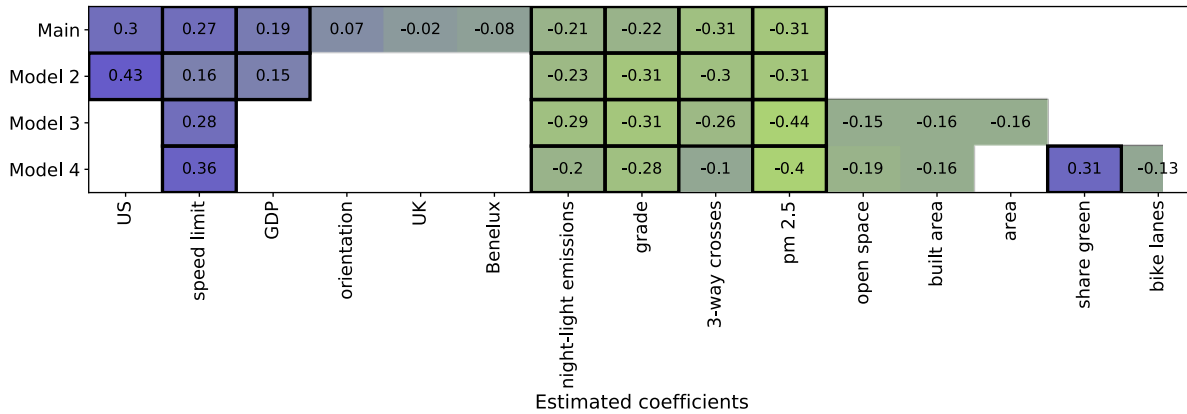


FIG. S4. **Understanding the determinants of gender gap in recreational cycling: cross-cities comparison: sensitivity analysis.** The panel provides a comparison between the preferred model (Main model) discussed in the main text and three additional models, included as sensitivity analysis. a) A summary of the characteristics of each model and scatter plots of (normalized) actual vs fitted values, for each model separately. All models share the linear formulation and the estimation technique (OLS) with the main one, but differ in terms of geographical coverage, standardization sample for target and input variables and the presence/absence of geographical dummies (Section III D). b) The heatmap displaying the estimated coefficient for each model. Cells outlined in black correspond to coefficients statistically significant at 0.05 significance level.

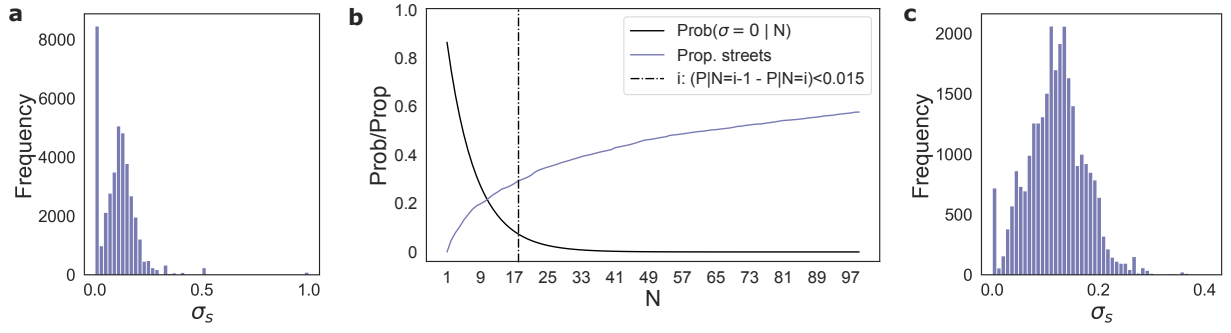


FIG. S5. **Impact of streets filtering.** a) The distribution of σ_s for streets in the City of New York, before filtering. b) The black line depicts the probability of observing $\sigma_s = 0$ conditional on the number of cyclists. The light blue line depicts the proportion of streets in the street network as a function of the number of cyclists. The vertical dashed line depicts the selected threshold for the data filtering. c) The distribution of σ_s for streets in the City of New York, after filtering of segments with number of cyclists below the selected threshold.

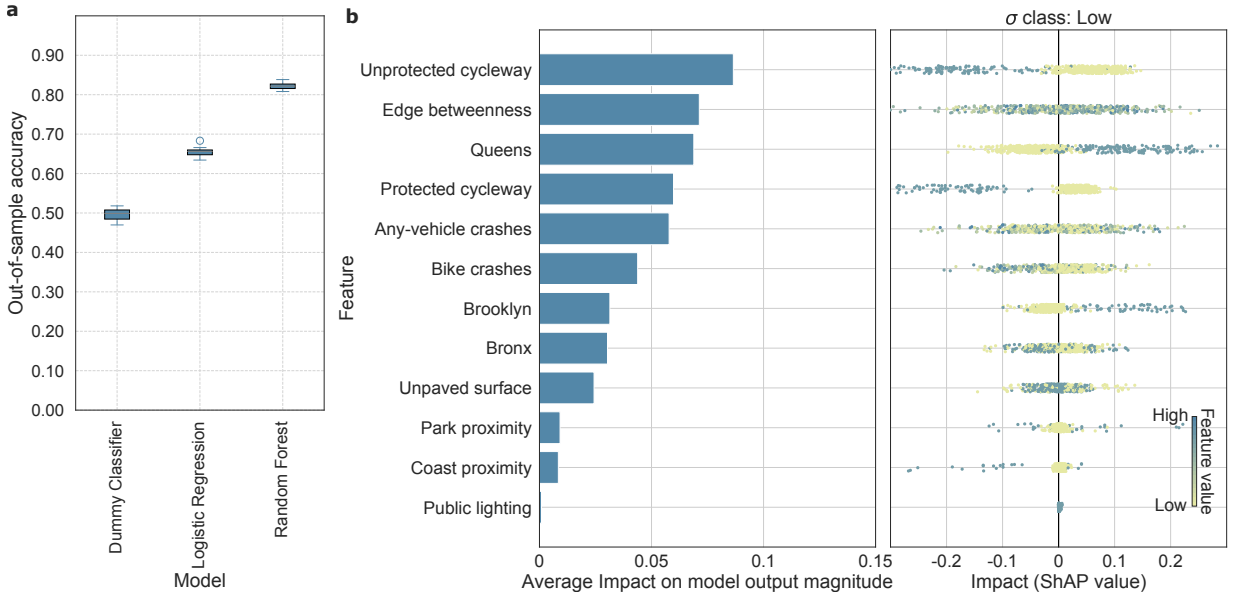


FIG. S6. **Classification task, city of New York - models comparison.** a) The boxplot present the out-of-sample accuracy of the logistic regression vs random-forest, for a value of the threshold $\alpha = 0.33$, computing using a stratified 10-fold approach. b) The bar plot presents average shapely values computed on a random selection of 500 data point. The adjacent strip plot presents the impact (in terms of shapely values) of a given feature on the probability that the street belong to the *Low* class, for 500 randomly selected data points [31]. Each point on the summary plot is a shapely value for a feature and an observation. The color represents the value of the feature from low to high. The negative values associated with the mass of blue dots for unprotected and protected cycleways indicate that, in the random selection of points here presented, the presence of a cycleway decrease the probability that the streets belong to the *Low* σ_s class. Given the symmetry of shape values in a binary classification task, this additionally indicates that for the data points in random selection, the presence of a cycleway increases the probability that the streets belongs to the *high* σ_s class.