# A Theory of Segregation Measurement

Sofía Correa      Daniel Hojman[*]

This version: March 2021

Latest Version Here

## Abstract

This paper proposes a theory of segregation measurement based on the intensity and social diversity of pairwise interactions. In our framework, societies are described by a space of locations, a space of social groups, and agents' distribution across locations and groups. Locations can be schools in a district, residences in a city, or platforms such as media outlets where individuals interact or meet. Social groups can be defined by race, socioeconomic status, political ideology, or any other source of social identity. We axiomatize measures that can be expressed as a weighted sum across pairs of an interaction *intensity* that depends on locations and an interaction *value* that depends on social identities. We prove that the index is proportional to the covariance between spatial and social distances, so that high segregation is associated with a high correlation between location and social proximity. We use our framework to study two segregation phenomena. The first one measures socioeconomic segregation in Chilean schools, showing variation across cities in line with residential segregation and across grades in line with differences in elementary and high school supply. The second one measures ideological segregation in media outlets' consumption, for different media platforms -newspapers, radio, TV- for 28 European countries. We find systematic differences in segregation across countries and platforms.

JEL Classification: C43, C81, D63, L82.

1

# 1 Introduction

Segregation in different domains remains a pervasive fact in contemporary societies. The lack of socioeconomic and racial diversity of interactions in schools and neighborhoods, and the exposure to like-minded ideological content can hinder a society's ability to embrace the value of diversity. Social inequality is often times both a cause and a consequence of residential, school and cultural segregation. Further, drawing on a long tradition in the social sciences, recent work in sociology has renewed attention on how barriers across social groups and segregation remains a fundamental barrier to equal opportunities.[1] The resurgence of inequality in local and global politics and the loss of trust in political and financial elites perceived as detached from "main street" also points in this direction.[2] Some of the most relevant social and technological recent changes may also contribute to a renewed interest in segregation. Specifically, recent mass migration waves across the world are often times associated with new segregated communities in the places of destination. In a different domain, the radical changes in the media landscape and new forms in which people access to news and information have also been associated with ideological segregation, a potential driver of the recent political polarization.[3] To the extent that different forms of segregation persist or arise, affecting social cohesion and the ability to construct a common ground for life in a democratic community, improving our understanding and measurement of segregation remains essential.

This paper provides a general framework to study segregation in different domains. We aim to contribute to a long tradition in the measurement of segregation. Our theory's starting point is shared by many conceptualizations of this phenomenon: segregation is the lack of interactions between individuals belonging to different social groups. Our framework considers a society of individuals that can differ in two dimensions, a *social type* that defines their social group, and a *location* (or a set of locations) they occupy. Depending on the application, social types can be race, socioeconomic status, ethnicity, nationality, religion, ideology, or any mix of social characteristics that define the groups we are interested in. An individual's location could be a

---

[1] For recent work on this front, see Lamont et al. (2014).

[2] See Atkinson et al. (2011) and related work for recent evidence on the evolution of inequality in a large number of countries. Sandel (2020) analyses the resentment against governing elites and anti-elite movements in a context of income polarization. Political scientists Hacker and Pierson (2010) and Bartels (2018) analyze the role unequal political representation on policy outcomes and inequality in the United States.

[3] See, for example, Gentzkow & Shapiro (2011) for evidence on ideological segregation in different media platforms. Campante & Hojman (2013) provide a theory and historical evidence on how changes in the media environment may cause political polarization.

home address in the case of residential segregation, a school, or the media outlets she visits to acquire information. For example, if the analyst is interested in racial segregation in schools, social types are races and, locations, schools. Suppose instead, she is interested in ideological segregation in media consumption. In that case, social types are individuals' political ideologies, and the location corresponds to the set of media outlets she consumes.

The theory presented in this paper has two building blocks. First, we take pairwise interactions as the basic unit of our measure. In practice, the measures we propose aggregate the contribution of each pairwise interaction to overall segregation. Second, the contribution of each interaction depends on two dimensions: the *intensity* of the interaction and the *diversity* of the social characteristics of the pair involved. On the one hand, the more two individuals have access to each other, the more they meet or encounter, the larger the intensity of that particular interaction. Intuitively, individuals who occupy the same location interact more than those occupying different locations. On the other hand, the value of an interaction depends on the social types of the pair interacting. In principle, an interaction between two individuals from the same social group contributes more to segregation than one between those belonging to different groups. In sum, both the intensity and the diversity of the interactions matter.

Our axiomatization allows us to obtain a simple formula in terms of these two dimensions characterizing each interaction. Building on axioms that parallel those of Expected Utility theory, Theorem 1 provides a representation of segregation that is precisely the sum over pairwise interactions of a product between the intensity of an interaction—analogous to the probability of two agents meeting—, and a social value of the interaction that depends on the social types of the pair.

The second general result of the paper shows that the measure can be interpreted in terms of a covariance of distances: a distance in the space of social types—a social distance—and a distance in the space of locations—a spatial distance. Specifically, we identify a large class of segregation measurement problems, *distance-based* problems, in which the intensity and the diversity of interactions are linear functions of distances in the space of locations and the space of social types, respectively. In the space of social types, it is always possible to define a distance between social groups. The choice of the distance depends on the problem under study and the researcher's objective. For instance, if we are only interested in distinguishing whether two individuals belong to same group or not—such as it is normally done in the study of racial segregation—two individuals with the same race are assigned distance 0 while two individuals of different races are assigned distance 1. The distance can be richer if the social types' space has an

order, such as unidimensional ideology space or income levels. In the first case, someone who is to the far right of the spectrum is further away from a left-winger than a moderate. Similarly, it is also natural to define a distance in the space of locations. Again, if the analyst is simply interested in distinguishing whether two individuals coincide at one location (e.g. students attending the same school), a discrete distance can be used. In other applications, more sophisticated notions of a distance could be useful (e.g. distance between residences, ideological distance between media outlets, etc.). Theorem 2 shows that in this context segregation is proportional to the covariance between the social distance—distance across social groups—and the spatial or location distance. Intuitively, segregation is large when two individuals of the same or proximate social groups are spatially proximate to each other, so that most interactions or encounters occur between individuals who share social characteristics.

We distinguish and obtain characterizations for two cases: with and without an individual resource constraints. The resource constraint represents a limit to the interactions an agent can have. Intuitively, if individuals have a time constraint, the more interactions they have, the smaller the time spent on each of them. This constraint reflects on the interaction intensity measure. We show that if individuals have access to a single location (e.g., school, residence) and have resource constraints, minimal segregation configurations are associated with a homogeneous distribution of social groups across areas. This need not be the case in the absence of resource constraint, as in this case, scale effects may play a role.

We use our framework to study two empirical applications. The first one measures socioeconomic segregation in the Chilean school system. We use micro-data containing administrative information on the socioeconomic status of the parents of each child in fourth and tenth grade. The social type of each student is defined by the education level of the parents and each student's location is simply the school she attends. We compute the school segregation index for the 22 largest cities in the country. We find that the segregation measures obtained for schools are highly correlated with the most recent measures of socioeconomic residential segregation. The variation across cities seems to be explained by sensible variables such as differences in the structure of local school supply. We also show that segregation depends on the grade, with more segregation in elementary school than high school, likely due to differences in mobility.

Our second application measures segregation in media consumption for 28 European countries, using survey data from Eurobarometer. Agents are characterized by an ideology—their social type—and a set of outlets where they get information from—their location. Agents are allowed to get information from many outlets, so that an individual location is really a vector for each

media platform—radio, TV, and newspapers—rather than a single location. Each component of the vector is associated with a particular outlet. In this case individuals can *meet* in more than one location. We find a strong correlation between segregation levels across media environments for each country, suggesting that there are some fundamental features, possibly related to the idiosyncratic political environment. In addition, we find that, for each media platform there is some correlation between segregation and the number of existing outlets, but this is not a general rule.

The paper contributes to the literature on three different margins. First, our framework is relatively general and can be applied broadly. In most previous theories of segregation measurement, both the social types or groups space and the space of locations is given. For example, in a classical paper such as Duncan & Duncan (1955) social types are race and locations residences (see also Massey & Denton (1988)). By considering a more general set up, most known applications can be accommodated by the framework. Second, an important consequence of a more general framework is that is it allows to explore more general propositions. In concrete, by introducing the notion of distances in the interaction space and the social types' space allows to provide a natural interpretation of segregation as a covariance between spatial and social distances of pairs of individuals. Another issue illuminated by our framework is that the segregation order induced by the measure may coincide with widely used measures such as the Duncan or Atkinson measures under some assumptions but not others. Specifically, for these and other measures, minimal segregation is achieved by a configuration in which interactions in each local community reproduces the distribution of social types in the general population. This is true in our framework under some assumptions such as the existence of a resource or budget constraint for each individual, but may not hold if this assumption is relaxed.[4] Finally, the theoretical flexibility in our framework allows to tackle problems in which individuals may encounter at multiple locations such as the consumption of multiple media outlets, as illustrated by our second application. The rest of the paper is organized as follows. Section 1.1 reviews the literature and its relationship with this work. The basic framework and the axiomatization of our measures is presented in Section 2. The general characterization of our measures is

---

[4]We explore in detail this issue in a companion paper (Correa & Hojman (2021)) where we show that, in the absence of capacity constraints, minimal segregation could be achieved by configurations in which some locations are associated with a relatively equal distribution of agents while other locations involve segregated individuals, that is, a share of the most numerous groups in locations with a socially homogeneous population. This configuration can sometimes maximize interactions between a social majority and a social minority group, due to scale effects.

summarized by theorems 1 and 2 in Section 3. Section 6.2 presents our two applications, one on socioeconomic segregation in Chilean schools and the other on the ideological segregation of media consumption in European countries. New questions and extensions of the framework are discussed in the conclusion section.

## 1.1 Related Literature

A detailed review of the most commonly used segregation indexes and their properties in social sciences can be found in James & Taeuber (1985) and Massey & Denton (1988). Most traditional segregation measures focus on residential segregation by race can be classified in two categories, *evenness* or *exposure*. Evenness measures focus on the differential distribution of groups across the city and include the Dissimilarity, Gini and Atkinson indexes. Instead, exposure measures refer to the potential contact between members of different groups, and the most used of this class is the Isolation index. A concise critique of the limitations of both types of measures can be found in Echenique & Fryer Jr (2007).

The economics literature on segregation measurement remains small, with relatively recent relevant contributions such as Echenique & Fryer Jr (2007) and Frankel & Volij (2011). These papers propose axiomatic foundations that identify the desirable properties of a segregation ranking, as we do. Alonso-Villar & Del Río (2010) and Frankel & Volij (2011) study multi-group indexes. Frankel & Volij (2011) propose two multi-group indexes for school segregation: the Atkinson Index for cases with a fixed number of social groups and the general case's Mutual Information index.

Echenique & Fryer Jr (2007) propose the Spectral Index to tackle two-race residential segregation measurement. Their measure is based on the aggregation of individual segregation measures in a residential network, where nodes are residences linked to neighboring residences. The value of individual segregation is implicitly defined as the weighted sum of direct neighbors' values in the network, where the weight is zero for neighbors of a different race and a positive number for same-race neighbors. This measure captures the idea agents are more segregated the more segregated are the same-race agents with whom they interact. Agents are budget-constrained in their interactions so that the positive weights just mentioned are inversely proportional to the total number of direct neighbors.

Our framework is more flexible than the ones provided by these papers. It can be used for any number of social groups. It is also suitable to study different environments, such as the school and residential segregation problems in each of these papers, and others, such as ideological

segregation in media consumption. Also, we consider the possibility of agents who interact in multiple locations and may or may not have budget constraints (which may be natural in some applications but not others).[5] In addition to a more general framework, some of the principles in our measures differ from these papers. For example, while the Spectral Index proposed by Echenique & Fryer Jr (2007) is based on the aggregation of individual segregation measures, our building blocks are pairwise interactions, so that our measures can be expressed as a weighted sum of values of each pairwise interactions. Each interaction's weights and values depend on distances in the landscape and across the social group of each pair. The characterization that we propose, in which both the diversity and intensity of an interaction between two agents are a function of distances, is suitable for different types of segregation problems and permits a relatively general characterization in terms of a covariance of the distances in the interaction and social space. Our paper is the first work proposing a segregation measure based on these concepts to the best of our knowledge.[6]

The applications presented provide empirical evidence of socioeconomic segregation in Chilean schools and ideological segregation in media consumption in Europe. The literature on socioeconomic segregation among schools is surprisingly undeveloped, most likely because of both the student level's lack of income data and the focus on segregation by race (see Reardon & Owens (2014) for a detailed review). In a novel study, Owens et al. (2016) study the evolution of income segregation between schools and school districts in the US. To overcome the lack of data, they use the count of enrolled students who are eligible for free lunch as a proxy for income. In this paper, instead, we use the parents' educational level as a proxy for students' income, which allows us to have a finer measure of income levels. We find that income segregation between schools is closely related to measures of residential segregation.

There is a burgeoning literature on segregation in media consumption, primarily driven by scholars' and policy-makers' concerns regarding the impact of new technologies and platforms on the supply of tailor-made content and increased selective exposure to like-minded views. For

---

[5]Using the terminology introduced later, we allow for general landscapes.

[6]To force a comparison, we can express our measure as a sum across individuals of the contribution of each of the pairs that include this agent (times one half to avoid double counting) and define the distance in the residential landscape as the distance in the residential network as defined by the the shortest path. As in Echenique & Fryer Jr (2007) the value of segregation increases with same-group interactions (a social distance equal to zero). However, in contrast to their index, the value of an interaction between individuals $i$ and $j$ does not depend on the individuals' races in the shortest path (or any path) between these individuals, just on the distance. In our framework, larger distances, regardless of intermediaries' race, are interpreted as having a lower intensity of interaction.

instance, Gentzkow & Shapiro (2011) study ideological segregation for media outlets online and offline using a very rich dataset. They compute segregation by using a dissimilarity index. The main difference between our approach and theirs is that our index considers that agents can interact on different locations simultaneously. Moreover, it allows for a complete ideological scale. Given that nowadays, the span of various news outlets' ideologies and specific content is more prominent, using measures that enable the analyst to use a finer grid of ideological positions makes a difference. Although the effects of segregation over political outcomes are out of the scope of this paper, we refer the reader to related works, such as Stroud (2008), DellaVigna & Kaplan (2007), and Campante & Hojman (2013).

## 2    A Segregation Model Based on Pairwise Interactions

In our framework, agents are characterized by a social type and a location in the space of interactions. This pair of characteristics is flexible enough allows us to include a large number of applications. An agent's social type could be race, income, education, ideology, ethnicity, religion, or any other characteristic (or combination of them) of interest to the analyst. If we are interested in racial segregation, the social type is race. If we are interested in socioeconomic segregation, the social type can be any measure of socioeconomic status.

On the other hand, implicit or explicitly, any segregation measurement takes into account an environment in which people interact. For residential segregation, this environment is a city where residences locate; for racial segregation in schools, it is a set of schools in a school district or city. In the case of ideological segregation in media consumption, the media outlets or platforms available to consumers are locations where agents can encounter. In these and other examples, we can identify a set of locations where people can coincide or not. In addition to locations, in each application, the environment has a configuration that can affect agents' distribution across locations. For example, in the case of schools, each school has a capacity that limits students who can use that location. In contrast, capacity constraints may not be relevant in the media outlets example, as the number of agents with access to a given outlet is unlimited in most cases. We refer to the set of locations and their capacities as a landscape. Both the locations and their capacities are features of the interaction structure determined mainly by either market forces or public policies.

While social types identify social differentiation, locations identify the space in which agents may or may not encounter and interact with each other. The measures of segregation we propose

aim to capture the extent to which agents with different social types—i.e., belonging to different social groups—are distantly located in the space of interactions.

We introduce the framework's essential elements, the definition of a distance-based measure, and the axioms postulated for our segregation measures.

## 2.1   The Basic Framework

The three main elements in our environment are: (i) a set of agents, $N$, (ii) a space of social types, $\Sigma$; (iii) a landscape, $\Lambda$. The set of agents $N$, with generic element $i$, is finite and, with some abuse of notation, $N$ also denotes its cardinality.

The space of social types $\Sigma$ is a finite set. Each element identifies a social characteristic defining the social group each agent belongs to. This space sets the stage for the specific problem of interest and remains fixed throughout the analysis.

**Definition 1** *A landscape $\Lambda$ is a pair $< L, Q >$, where*

*(i) $L$ is a set of locations, and*

*(ii) $Q = \{Q_l\}_{l \in L}$ is a collection of location capacities, with $Q_l \in \mathbb{R}_+$ possibly unbounded.*

*The space of admissible landscapes is denoted by $\mathcal{L}$.*

Note that given a landscape $\Lambda$, we define a space of individual locations $X(\Lambda)$, which determines the way agents can be distributed across the landscape. More precisely, agent $i$'s location is a vector $x^i = (x_1^i, ..., x_L^i) \in X(\Lambda)$, which describes the location(s) in the landscape the agent visits. For instance, in the case of schools, students can only attend one school at a time, and then the space of individual locations is described as

$$X(\Lambda) = \{x^i \in \{0,1\}^L | \sum_{l \in L} x_l = 1\}. \tag{2.1}$$

If schools have some capacity constraints, the space of individual assignments would be a constrained space, i.e.,

$$X(\Lambda) = \left\{ x^i \in \{0,1\}^L | \sum_{l \in L} x_l = 1 \wedge \sum_{i \in N} x_l^i \leq Q_l \right\}. \tag{2.2}$$

We assume that the set of admissible landscapes includes the possibility of complete segregation, i.e., a profile of location assignments such that for any two agents $i$ and $j$ with different social types, if $x_l^j > 0$ then $x_l^j = 0$.

Given the space of individual locations $X(\Lambda)$, define a profile of locations $\mathbf{x} \in \mathbb{R}^{L \times N}$ as the set of all agents' locations. In particular, $\mathbf{x}$ is a matrix whose $i$th row is given by $x^i \in X(\Lambda)$. Analogously, given the space of social types $\Sigma$, define a profile of social types $\mathbf{s} \in \mathbb{R}^{|\Sigma|}$, as a vector with generic element $s^i \in \Sigma$. We put these profiles together in the following definition.

**Definition 2** *A profile of agents' characteristics is given by a pair of profiles* $(\mathbf{x}, \mathbf{s}) \in \mathbb{R}^{L \times N} \times \mathbb{R}^{|\Sigma|}$, *with generic element* $(x^i, s^i) \in X(\Lambda) \times \Sigma$.

The unit on which we measure segregation is called a *community*. A community might be, for instance, a city, a set of schools, or a media platform, in which a set of agents with some social types interact. As we mention above—and as it is usually the case in the study of segregation— we keep the space of social types fixed across communities and allow all the other features of the problem to varying. To illustrate this, consider the problem of residential segregation by race. In that case, we fix $\Sigma$ (a set of races) and compare the segregation levels observed in different cities. Each city is characterized by a landscape, a set of agents, and a profile describing the distribution of individual races and locations in the city. Formally:

**Definition 3** *Fix* $\Sigma$. *A community is given by a tuple* $(N, \Lambda, X(\Lambda), (\mathbf{x}, \mathbf{s}))$, *where:*

(i) $N$ *is a finite set of agents;*

(ii) $\Lambda$ *is a* landscape, *with associated space of individual locations* $X(\Lambda)$; *and,*

(iii) $(\mathbf{x}, \mathbf{s}) \in \mathbb{R}^{L \times N} \times \mathbb{R}^{|\Sigma|}$ *is a profile of agents' characteristics.*

To fix ideas, consider the case of racial segregation in schools for a given city. The set of $N$ agents corresponds to a set of students in the city. The social space is a set of racial groups $\Sigma = \{\text{Black, White, Asian,...}\}$, which remains fixed across communities. The landscape $\Lambda$ is composed of a set of schools $L = \{\text{school } 1, \text{school } 2, ...\}$, and a set of capacities for each school $Q = \{Q_1, Q_2, ...\}$, where the latter corresponds to the maximum number of students that each school can admit. Each student's social type is a race $s^i \in \Sigma$. An individual location is a vector $x^i = (x_1^i, ..., x_L^i)$ such that $x_l^i = 1$ if student $i$ attends school $l$, and zero otherwise.

The building blocks of our segregation measure are pairwise interactions. More specifically, our measures are based on the aggregation of each pairwise interaction's value, weighted by a measure of the interaction's intensity. Denote by $\Pi(N)$ to the set of possible pairwise interactions between agents in $N$, with generic element $\pi = (i, j)$.[7] We omit $N$ when it is clear from the

---

[7]More precisely, $\Pi(N) = \{\pi = (i, j) \in N \times N \mid i \neq j\}$.

context. Let $\Delta$ denote the space of probability distributions over $\Pi$. Each pairwise interaction $\pi = (i, j)$ contributes to the level of segregation through two components: (i) an *intensity*, describing how likely it is that a pair of agents meet in the landscape, and (ii) a *social value*, describing the value of the interaction. Throughout the paper, we focus on location-based intensity functions—i.e., the intensity of an interaction between two agents depends on their locations—, and social values that are type-based—i.e., the value of an interaction between two agents depends on their social types.

**Definition 4** *Let $\Pi$ be the set of possible pairwise interactions.*

(i) *$\mu : \Pi \to \Delta$ is a location-based intensity function if for some function $m : X(\Lambda) \to \Delta$, $\mu(i, j) = m(x^i, x^j)$ for all $(i, j) \in \Pi$.*

(ii) *$\rho : \Pi \to \mathbb{R}$ is a type-based social value if for some function $r : \Sigma \to \mathbb{R}$, $\rho(i, j) = r(s^i, s^j)$ for all $(i, j) \in \Pi$.*

We interpret the intensity function as a descriptive measure of the likelihood that any two people in a community meet or have access to each other. We note, however, that a normative baseline could define the relative intensity of a particular interaction. To illustrate this issue, consider racial segregation in schools. Implicitly, most segregation measures used for this problem assume that any agents in the same school (more generally, two agents in the same location) are accounted equally by the measure. This assumption does not consider that within a school, students may endogenously sort based on homophily, that is, with a tendency to interact more with students of the same race. For example, take two schools $A$ and $B$, each one with 20 students (2 black, 2 Asian, 2 Latino, and 14 white). Suppose that in school $A$, students interact randomly, disregarding race, while in school $B$, students only mingle with those of the same race. Existing measures would not distinguish two situations. There might be at least two reasons for this. The first one might be a practical limitation of the data: it seems pointless to distinguish between these two schools if the information regarding detailed interactions or social networks within a school is not available. The second reason normative: from the perspective of social cohesion and empathy, it may be valuable for students to have access or contact with a diverse population of students.

The group-based nature of the social value represents the fact that the value of an interaction between two agents depends on their social types. The term "social value" emphasizes the idea that value in a theory of segregation is not associated with a measure of personal productivity but with the social diversity of interactions.

As seen shortly, it will prove convenient to use an aggregate version of the intensity function, $\tilde{\mu} : \Sigma \times \Sigma \to \Delta$, defined by

$$\tilde{\mu}(s, s') := \sum_{i:s^i=s} \sum_{j:s^j=s'} \mu(i, j). \tag{2.3}$$

This corresponds to the aggregate fraction of interactions between agents in groups $s$ and $s'$. In addition to this, we denote by $N_l$ to the number of agents in location in $l \in L$, and $N_s$ to the number of agents with social type $s \in \Sigma$.

## 2.2 Distance-based Measures

In this section, we consider the case in which both the space of social types and locations can be associated with natural notions of distances. The motivation for this is both conceptual and computational. As shown in Section 3, our measures conceptualize segregation as the covariance between a social distance and location distance. Intuitively high segregation can be associated with individuals with the same or proximate social types located in the same or near locations.

Let us consider first the space of locations. The notion of distance is natural in this case. Consider, for instance, residential segregation. People naturally live closer to some neighbors than to others, and there are many distances between residences that can capture that. One might consider a discrete distance indicating whether individuals live in the same block, a geodesic distance, or even walking times between two residences. This case differs from school segregation. When the space of interactions is a set of schools, it is usually assumed that students only interact with other children in the same school. In that case, it is direct to endow the space of interactions with a discrete metric taking value 0 for children in the same school and 1 for children in different schools.

In the space of social types, the notion of a distance might be more subtle. Social types might take the form of groups clearly identified, an ordered set of categories, or a more continuous scale. In the case of race or ethnicity, a trivial notion of distance is given by the discrete metric that assigns distance 0 to those in the same group and distance 1 to those in different groups. Indeed, this notion of distance can always be defined and associated with space of social types $\Sigma$. There are other cases, however, in which the structure of $\Sigma$ induces a natural order and distance. For example, if the socioeconomic status is defined by a scale such as income or education years, types can be ordered, and there is a naturally defined metric space. This is also the case for social types defined on a uni-dimensional ideology line.

We say that a problem is *distance-based* if it is possible to characterize both the intensity

and the social value of an interaction by some notion of distance. More precisely, a problem is distance-based if both the space of individual assignments in the landscape, and the space of social types, are endowed with a metric. Denote these metrics by $d_x : X \times X \to \mathbb{R}$ and $d_s : \Sigma \times \Sigma \to \mathbb{R}$. If this is the case, the intensity and social value functions from Definition 1 can be defined in terms of distances.

Definition 5 formalizes the notion of linear distance-based problems, that will be our relevant framework to obtain Theorems 2 and 3. For simplicity of notation, we define the distance in the landscape by $d_\Lambda(i,j) = d_x(x^i, x^j)$, and the distance in the space of social types by $d_\Sigma(i,j) = d_s(s^i, s^j)$.

**Definition 5** *A distance-based problem is a problem in which $(\Sigma, d_s)$ and $(X, d_x)$ are metric spaces. A distance-based problem is linear if the functions $\mu : \Pi \to \Delta$ and $\rho : \Pi \to \mathbb{R}$ satisfy (i) $\mu(i,j) = m_0 - m_1 d_\Lambda(i,j)$ for some $m_0, m_1 > 0$, and (ii) $\rho(i,j) = r_0 - r_1 d_\Sigma(i,j)$, for some $r_0, r_1 > 0$.*

## 2.3 Axioms

In this section we state the main axioms. The first one, *Anonimity*, reflects the fact that a segregation order does not depend on agents' identity but on their social characteristics.

**Axiom 1 (Anonimity)** *Any permutation of agents that preserves the original social groups does not change segregation.*

This axiom allows us to focus on the functions $\tilde{\rho}$ and $\tilde{\mu}$ instead of $\rho$ and $\mu$. The next two axioms refer to how changes in communities, and combinations of them, affect segregation. We combine communities by combining the aggregate intensity functions that characterize them (see equation 2.3). We formalize this idea in the following definition.

**Definition 6** *Let $C_1$, $C_2$ be two communities, with associated aggregate intensity functions $\tilde{\mu}^1$ and $\tilde{\mu}^2$, respectively. Then, we define the community $C = \alpha C_1 + (1 - \alpha)C_2$ as one with aggregate intensity $\tilde{\mu} = \alpha \tilde{\mu}^1 + (1 - \alpha)\tilde{\mu}^2$.*

In Appendix C.1, we illustrate these operations through an example for the case of socioeconomic segregation in schools.

**Axiom 2 (Continuity)** *Let $C_1, C_2, C_3$ be three communities such that $C_1 \succeq_s C_2 \succeq_s C_3$. Then, there exist $\alpha$ such that $\alpha C_1 + (1 - \alpha)C_3 \sim_s C_2$.*

**Axiom 3 (Independence)** *Let $C_1, C_2$ be two communities such that $C_1 \succeq_s C_2$. Then, for any $C_3$ and $\alpha \in [0,1]$, $\alpha C_1 + (1-\alpha) C_3 \succeq_s \alpha C_2 + (1-\alpha) C_3$.*

In addition to these technical axioms, we introduce axioms that allow us to define the distance-based notion of segregation.

**Axiom 4 (Spatial Proximity)** *Let $\pi = (i,j), \pi' = (i',j')$ to be two pairs such that $d_\Lambda(i,j) > d_\Lambda(i',j')$. Then, $\mu(i,j) < \mu(i',j')$.*

Spatial Proximity refers to the idea that the probability of two agents meeting in the landscape is decreasing in the distance between them.[8]

**Axiom 5 (Social Proximity)** *Let $\pi = (i,j), \pi' = (i',j')$ to be two pairs such that $d_\Sigma(i,j) > d_\Sigma(i',j')$. Then, $\rho(i,j) < \rho(i',j')$.*

The Social Proximity axiom corresponds to the idea that segregation decreases with more diverse interactions. In other words, the more similar two people are, the more does their interaction contributes to segregation. Conversely, diverse interactions contribute to lower segregation. We observe that Integrating different groups of people might affect economic outcomes through changes in stereotypes, beliefs about others, and changes in social interactions. There is vast literature on the effects of social and ethnic diversity over several outcomes.[9] For instance, Alesina & La Ferrara (2000) and Easterly & Levine (1997) show a negative relation between social diversity and public good provision and GDP. However, there is also evidence that goes in the opposite direction. Hong & Page (2001) and Hong & Page (2004) develop a theory of problem-solving, in which social diversity is beneficial since it brings with it a different interpretation to a problem, increasing innovation and the probability of solving it. Alesina & La Ferrara (2005) develop a model in which preference diversity may imply public goods under-provision, which might be socially costly.

Still, there might be benefits in terms of innovation and creativity, which might help overcome the costs depending on the development of the group under study. Also, there might be differences in terms of how the benefits of diversity present over time. Putnam (2007) shows that social diversity costs in terms of trust in society are observed only in the short term, but social diversity is welfare improving in the long term. Beyond the evidence on the value of diversity on different types of productivity, the social value's monotonicity concerning social diversity could

---

[8]This is the same intuition present in the axiom *School Division Property* in Frankel and Volij (2011): dividing a school is analogous to increasing the distance between agents.

[9]For a detailed literature review on the topic see Alesina & La Ferrara (2005).

be defended from a normative perspective. We aim to extend the analysis to measures that do not impose this monotonicity condition in future work.

So far, axioms 4 and 5 ensure that the intensity of an interaction is a decreases with the distance in the landscape, and that the value of each interaction decreases with the distance in the types' space. These axioms do not restrict the curvature of these functions.

The following axioms, impose structure on these functions that allow for a simple interpretation and a representation convenient for applications. They do not seem essential to us but yield simpler formulas, and simplicity is useful for applications.

**Axiom 6 (Linear Intensity)** *The effect of an additive increase in the spatial distance between two agents over their interaction intensity does not depend on their original distance. More precisely, take a pair $(i,j)$ with two possible assignments $(x^i, x^j)$ and $(x_0^i, x_0^j)$ with intensities $\mu(i,j)$ and $\mu_0(i,j)$, respectively. Then, $\mu(i,j) - \mu_0(i,j) = K \cdot |d_x(x^i, x^j) - d_x(x_0^i, x_0^j)|$ for some constant $K < 0$.*

**Axiom 7 (Linear Value)** *The effect of an additive increase in the social distance between two agents over their value to segregation does not depend on their original distance. More precisely, take a pair $(i,j)$ with two possible social types, $(s^i, s^j)$ and $(s_0^i, s_0^j)$ with intensities $\rho(i,j)$ and $\rho_0(i,j)$, respectively. Then, $\rho(i,j) - \rho_0(i,j) = K \cdot |d_s(s^i, s^) - d_\Lambda(s_0^i, s_0^j)|$ for some constant $K < 0$.*

# 3 Representation Results

In this section we state our main results. For the ease of exposition all the proofs are relegated to Appendix A.

## 3.1 General Representation Result

We first prove that a segregation order can be represented by a function valuing pairwise interactions, which is the baseline for the measures obtained in sections 3.2 and 4.

**Theorem 1** *The preference order $\succeq_s$ satisfies axioms 1-3 if and only if such preferences are represented by:*

$$S = \sum_{(s,s') \in \Sigma \times \Sigma} \tilde{\mu}(s, s')\rho(s, s'), \tag{3.1}$$

*where $\tilde{\mu}(s, s')$ is defined by equation 2.3.*

Although this representation might seem general, it gives us an intuition of understanding segregation in a society. The basic idea is intuitive: the distribution of agents in the space determines the probability of observing each type of interaction. Each city is like a lottery of interaction values. To measure segregation, we measure the value of these *lotteries* for a given segregation order. In the next section, we analyze the case in which interactions can be characterized in terms of distances.

## 3.2  A Distance-based Representation

In linear distance-based problems, both the intensity and the social value of an interaction are linear functions of distances in the corresponding spaces (see Definition 5). These functions' linearity allows us to prove a very intuitive result: segregation is proportional to the covariance between social and spatial distances. The idea is that a community is more segregated if similar people are more likely to meet, i.e., similar social types are closer in the space of interactions. Analogously, a community is more segregated if different people are unlikely to meet, i.e., different social types are far from each other in the space of interactions. The more the social and spatial distances covary, the more segregated the society is. In the extreme, a completely segregated society would be one in which there is a one-to-one mapping from the space of social types to the space of interactions.

**Theorem 2** *A segregation index satisfies Axioms 1-7 if and only if it is proportional to $S = cov(d_\Lambda, d_\Sigma)$.*

Note that axioms 4 to 7 are only consistent with linear distance-based problems. In particular, we can show that a problem is linear distance-based if and only if it satisfies axioms 4-7. An alternative interpretation of Theorem 2 is that segregation is proportional to the coefficient of a regression of the distance in the space of locations on the distance in social types. Thus, segregation is a measure of the linear association between both.

# 4  Index Normalization

In its more general form, our index is defined by:

$$S = \frac{1}{\Pi} \sum_{(i,j) \in \Pi} \rho(i,j)\mu(i,j) \tag{4.1}$$

In most of the segregation literature, measures are normalized between 0 and 1, where the lower bound represents a configuration that achieves the minimal segregation and the upper

16

bound the maximal. For example, for the Duncan index, which is defined for two social types, minimal segregation is achieved when each location's population reproduces the global distribution of types in the community. The maximal segregation is obtained when individuals of the same social group occupy all locations.

Recent work shows that normalizing a segregation index is not trivial.[10] In our case, in line with Expected Utility Theory, it is always possible to consider a linear transformation of our measure that respects the segregation ranking and is associated with a normalization:

$$\hat{S} = \frac{S - S^{min}}{S^{max} - S^{min}} \tag{4.2}$$

which by construction is 0 for the minimum value of the index, $S^{min}$, and 1 for its maximum value, $S^{max}$. However, computing these numbers for a given application is not necessarily immediate. More precisely, a uniform distribution of social types across locations does not necessarily achieve minimal segregation, in contrast to the Duncan index.[11]

The underlying normalization criteria of segregation measures are often down-played in the literature. In the residential segregation literature, normalization typically considers the landscape as fixed, i.e., for a fixed $\Lambda$. In practice, this assumes a given built (or physical) environment so that the minimal and maximal segregation are found by varying the individual's assignment profile across locations. The underlying thought experiment is that people of different social groups can move within the same city. With a few exceptions, this is usually the case in the school segregation literature as well. From a descriptive perspective, it is reasonable to take the landscape as fixed in the short run. From an analytical perspective, this is reasonable if the segregation measures provided are invariant to landscape changes. However, neither of these assumptions is obvious.

Market supply and public policy that affect the landscape can have a substantial impact on segregation. For example, new urban developments influence segregation in a city. In recent work, Agostini et al. (2016) show that the evolution of segregation in Santiago—one of Latin America's most populated towns—is determined mainly by housing policies that resulted in

---

[10]Echenique and Fryer (2006) do not impose normalization for the maximum segregation.

[11]This issue is treated in detail in our companion paper Correa and Hojman (2021). There we show, for the case of schools, that the configuration that achieves the minimal segregation is not obvious and is sensitive to assumptions that may seem innocuous at first. This shows that our measure captures a notion of segregation that, depending on the assumptions, may coincide with traditional measures. However, it may also differ in a meaningful way. We also show that the associate minimization problem is well behaved, i.e., it associates with the optimization of a quadratic form with linear constraints and can always be solved numerically.

the allocation of low-income families in the new peripheral neighborhoods of the city. Large government-funded projects, new secluded high-income developments in the city's borders, and central areas' densification. On the other hand, changes in media technologies and markets have drastically affected the supply of content in the media market, affecting ideological segregation and polarization (see Campante & Hojman (2013) and Levy & Razin (2019)). The landscape can also be affected by regulations and government initiatives—for instance, publicly-funded school supply and capacities. In an urban setting, construction regulations, location of social housing are only some examples. In media markets, the regulation of media concentration ownership, entry, and policies to foster balance in media supply ideologies provide additional illustrations.

Simultaneously, from a normative stance, it is far from evident whether segregation measures should be invariant to changes the landscape. In a companion paper, we show for the case of school segregation that the nature of assignments that minimize segregation can vary substantially across different landscapes. Perhaps more importantly. If changes in the landscape affect segregation or the minimum and maximal segregation, why should segregation normalization take them as fixed? To illustrate this issue, consider two societies, A and B, with the same population and distribution of social groups. Suppose that in A there are 1000 small schools and in B ten large schools. Why should segregation normalization take the number of schools as fixed if this variable changes over time? Similarly, if A and B are two countries, one with 1000 news websites and the other with 10, why should we normalize segregation taking the media landscape as given?

In principle, our framework does not take a stance on whether the landscape should vary in the calculation of maximizing segregation. In fact, in the short run the landscape is not flexible. We may consider two variants of the optimization problem leading to the maximal and minimum segregation levels. The first one takes the landscape $\Lambda$ as given and considers the location assignment profile as the optimization variable:

$$\min_{\mathbf{x} \in \mathbf{X}(\mathbf{\Lambda})} S(\mathbf{x}, s; \Lambda). \tag{4.3}$$

The solution of this problem is some profile $\mathbf{x}^*(\mathbf{\Lambda})$. The value of the problem is $\Phi(\Lambda) = S(\mathbf{x}^*(\mathbf{\Lambda}), s; \Lambda)$ (analogous for maximization). If minimal and maximal segregation allow to vary the landscape, we can consider a second stage optimization across admissible landscapes,

$$\min_{\Lambda \in \mathcal{L}} \Phi(\Lambda), \tag{4.4}$$

yielding a solution $\Lambda^*$. In the Appendix, we describe the normalization optimization problem for the application to ideological segregation in media outlets, which is solved using numerical

18

methods.

In the following section we specialize the measures to the case of unit location assignments and individual capacity constraints. In this case, large class of problems, we provide an explicit normalization and formula.

# 5 Individual Capacity Constraints

So far, we have considered the case in which agents do not have capacity constraints in their interactions with others. For example, given two schools, $A$ and $B$, one with 20 students and the other with 100, and assuming that agents interact only with agents in their school, the intensity of interaction functions give equal weight to any interaction in school $A$ and school $B$. The unconstrained capacity could be a reasonable assumption for situations in which what matters is whether two students share the same space but not the amount of time spend with each other. If empathy mainly were associated with sharing a common location with people of a different race —access to a diverse social group— but not with the time spent with people of other groups, this makes sense. This is also a sensible assumption if the most relevant activities in school involve the whole population. However, in many cases, it makes sense to assume that each agent has a capacity constraint, a time o resource budget that she allocates to each interaction (see, for example, Echenique & Fryer Jr (2007)).

In this section, we extend the framework to consider agents having a time budget or capacity constraint. In the school case, a capacity constraint implies that the probability of meeting other students decreases with the size of the school, or, in other words, the time spent with each classmate decreases as the number of classmates increases. Thus, the weight of a particular interaction in school $A$, the smaller school, is larger than the weight of an interaction in school $B$, the larger school.

We introduce a new axiom that formalizes capacity constraints in our framework. The intuition is that, as agents have a limited time to interact with each other, to increase the interaction with one agent, they have to reallocate their resources across space, decreasing the intensity of interaction with others.

**Axiom 8 (Resource Constraint)** *Agents have limited resources to interact with each other, which is represented by the following resource constraint:*

$$\sum_j \mu(i,j) = T \qquad \forall i \tag{5.1}$$

We restrict the analysis to a framework in which agents can only be assigned to a unique location. This imposes a restriction over the space of individual assignments $X$, for any given landscape $\Lambda$. The following definition formalizes this notion.

**Definition 7** *An individual assignment is a Unit Location Assignment (ULA), if (i) $x_i \in \{0,1\}^L$, and (ii) $x_i^T e_L = 1$.*

Let $\Pi_l = \{(i,j)|x_i = x_j = e_l\}$ to be the set of possible pairwise interactions within location $l \in L$, and denote by $\bar{d}_\Sigma^l = \frac{1}{\Pi_l} \sum_{(i,j) \in \Pi_l} d_\Sigma(i,j)$ to the average social distance between pairs of agents within location $l$.

**Proposition 1** *Consider an assignment satisfying (ULA), and suppose axioms 1-8 hold. Then, the normalized segregation index is given by*

$$S = 1 - \frac{1}{\bar{d}_\Sigma} \sum_{l=1}^{L} w_l \bar{d}_\Sigma^l, \tag{5.2}$$

*where $w_l = \frac{N_l}{N}$ is the share of the population in location $l$.*

The proof can be found in the Appendix. Observe that the index can be expressed as

$$S = \sum_{l=1}^{L} w_l \left(1 - \frac{\bar{d}_l^\Sigma}{\bar{d}^\Sigma}\right), \tag{5.3}$$

that is, a weighted sum across locations of an expression that compares the local social distance with the global social distance. The weights are given by the share of the population in each location. This is parallel to the dissimilarity index, which can also be expressed as a weighed sum across locations of a quantity that compares the local share of a minority relative to the aggregate share of the minority.

# 6 Applications

## 6.1 Socioeconomic Segregation in Schools

In this section, we use distance-based segregation measures to study school segregation. Given the poor availability of income data at the student level, the empirical work on schools' economic segregation is not very developed. In an attempt to fill this gap, in this paper, we use administrative individual-level data of Chilean students provided by the Ministry of Education to explore socioeconomic segregation in schools across different cities. For each student in 4th grade and

10th grade in 2014, the dataset identifies the school they attend, in addition to sociodemographic and family background variables.[12] In the absence of reliable family income, we use the parents' educational level —average years of education of mother and father—as a measure of each student's socioeconomic status.

Segregation measures are computed for all regional capitals, including the three Chilean metropolitan areas, Santiago, Valparaíso and Concepción. In our main terminology, each of these regional capitals constitutes a different community. A set of students inhabits each city, and each student is characterized by a socioeconomic level and a single school to which she attends.

We denote each school by $l \in \{1, ..., L\}$, and $l_i$ corresponds to the school attended by student $i$. We use $N$ to denote the total number of students in a district, $N_l$ to the number of students at school $l$, and $N_{g,l}$ to the number of students in quintile $g$ attending school $l$. Also, denote by $\Pi$ and $\Pi_l$ the number of pairs in the population, and in school $l$, respectively.

In this context, the landscape is a partition, as in Unit Location Assignments (see Definition 7). As students only interact with other students attending their same schools and assuming that interaction is uniform, the corresponding metric in the space of interactions is a discrete metric.

We divide the income distribution into income quintiles, and associate to each student the average parents' educational level of the corresponding quintile. Thus, $d^{\Sigma}(i, j) = |y_i - y_j|$, where $y_i$ is average years of education of the respective quintile. Then, from proposition 1 the normalized index can be computed using the following equation:

$$S = 1 - \frac{1}{\overline{d}_{\Sigma}} \sum_{l \in L} \frac{N_l}{N} \overline{d}_{\Sigma}^{l} \tag{6.1}$$

where $\overline{d}_{\Sigma}$ is the average social distance in the population, and $\overline{d}_{\Sigma}^{l}$ is the average social distance at school $l$.

The results are shown in Table 1. A few remarks are in place. First, for both grades, segregation is highly correlated with the city's size (in terms of students and schools). In particular, Santiago is the most segregated city in both grades, followed by Temuco and Valparaiso. Secondly, in each city, segregation is higher in 4th grade than in 10th grade. These facts are consistent with residential socioeconomic segregation patterns in Chilean metropolitan areas

---

[12]In Chile, students of these grades in all schools —with a few exceptions— are required to take a standardized test in math and language, the SIMCE. The Ministry uses a complementary questionnaire to gather sociodemographic and family background information in addition to the individual scores. The dataset covers roughly 95 percent of all Chilean schools, excluding new and special education schools.

and the role of distance in school choice. Indeed, we compare our school segregation index with residential segregation measures for Chilean cities obtained by Agostini et al. (2016), who use household income as an SES measure and Census data. We obtain a correlation coefficient between school and residential segregation of 0.79 for 4th-grade students, vs. a 0.67 for 10th grade. Even when both are high, the correlation for 4th grade is slightly higher, probably reflecting the fact that younger children have more mobility constraints than students in 10th grade and tend to attend schools that are, on average, closer to their residences than those in 10th grade. Both comparisons are plotted in Figure 1.

Table 1: School socioeconomic segregation by city

| Region | 4th Grade | | | 10th Grade | | |
| --- | --- | --- | --- | --- | --- | --- |
| | S | Schools | Students | S | Schools | Students |
| Santiago | 0.395 | 1,244 | 50,535 | 0.346 | 753 | 41,020 |
| Temuco | 0.385 | 109 | 3,031 | 0.308 | 189 | 7,320 |
| Valparaiso | 0.353 | 340 | 9,079 | 0.308 | 108 | 6,145 |
| La Serena | 0.343 | 87 | 2,773 | 0.276 | 31 | 1,626 |
| Concepcion | 0.340 | 279 | 9,351 | 0.302 | 39 | 1,522 |
| Puerto Montt | 0.339 | 100 | 2,886 | 0.286 | 44 | 2,434 |
| Valdivia | 0.336 | 60 | 1,552 | 0.302 | 18 | 1,119 |
| Talca | 0.331 | 67 | 2,658 | 0.273 | 55 | 2,146 |
| Antofagasta | 0.331 | 74 | 4,214 | 0.253 | 44 | 3,107 |
| Rancagua | 0.310 | 76 | 2,930 | 0.255 | 42 | 2,612 |
| Coyhaique | 0.287 | 25 | 725 | 0.175 | 45 | 2,906 |
| Punta Arenas | 0.287 | 37 | 1,409 | 0.246 | 35 | 1,355 |
| Copiapo | 0.271 | 35 | 1,965 | 0.193 | 30 | 1,553 |
| Iquique | 0.262 | 55 | 2,176 | 0.208 | 11 | 470 |
| Arica | 0.241 | 65 | 2,422 | 0.215 | 22 | 946 |
| Mean | 0.321 | 177 | 6514 | 0.263 | 98 | 5085 |
| Max | 0.395 | 1244 | 50535 | 0.346 | 753 | 41020 |
| Min | 0.241 | 25 | 725 | 0.175 | 11 | 470 |

Figure 1: School vs Residential Segregation.

## 6.2 Segregation in Media Consumption

In many situations, agents can interact with others in multiple locations. In particular, in media consumption, agents can obtain information from different outlets, and on each of them, they *meet* or coincide with other agents. This meeting can be a direct interaction, as may happen in an online news forum, or indirectly, from obtaining the same information when they read the same newspaper. This section measures ideological segregation in media consumption where each location is a specific media outlet, and individuals may consume more than one outlet. That is, location is a vector with one component per outlet.

Measures of segregation such as the Isolation index (see Gentzkow & Shapiro (2011)) may not fully consider two media consumption features. First, agents are characterized by a rich set of social types such as ideology, and this richness may be partially lost by imposing a discrete metric over the social space. And second, agents can obtain information from (and henceforth, interact through) more than one media source. Our framework can accommodate these issues.

We assume that agents have a budget of time that can allocate across news sources. We consider the individual *location* vector that proxy the share of time spent on each of them. We analyze three markets separately, i.e., TV, newspapers, and radio stations' consumption. The social characteristic we are interested in is a measure of political ideology. Then, our segregation index measures the extent to which individuals with different ideologies are sharing media consumption. We use survey data from Eurobarometer 82.4 (2014) that covers 28 European

countries. The survey has a series of questions that ask individuals about the TV stations, radio stations, newspapers, and websites they use.[13][14]. The survey also asks individuals to self-identify in an ideology uni-dimensional ten-point scale, where 1 is the extreme left and 10 is the extreme right.[15] In Table 2 we show some descriptive statistics of the dataset.

Let $\Lambda = 1, ..., L$ to be the set of media outlets on each market, and $x_l^i \in \{0, 1\}$ an indicator variable taking the value 1 if $i$ consumes outlet $l \in L$ and zero otherwise. As defined in Section 2.1, the vector $x^i = (x_1^i, ..., x_L^i)$ summarizes individual $i$'s consumption bundle, and we focus on the normalized version $\hat{x}_i$ with generic element $\hat{x}_l^i = \frac{x_l^i}{\sum_{l=1}^{L} x_l^i}$. Individuals interact to the extent the share media consumption, that is, if they coincide in their locations. The effective number of location coincidences between two individuals $i$ and $j$ is then $\sum_{l=1}^{L} \min\{\hat{x}_l^i, \hat{x}_l^j\}$. This defines our distance in the landscape:

$$d_\Lambda(i, j) = 1 - \sum_{l=1}^{L} \min\{\hat{x}_l^i, \hat{x}_l^j\}. \tag{6.2}$$

The following example illustrates the intuition for this distance. There are two outlets, $\Lambda = \{1, 2\}$, and two agents $i$ and $j$. Agent $i$ gets information from both outlets, spending half of the time on each, but agent $j$ only acquires information from outlet 1, spending all the time there. Formally, $\hat{x}^i = (0.5, 0.5)$ and $\hat{x}^j = (1, 0)$. Then, they both share only half of the total time together: they *coincide* $\min\{0.5, 1\} = 0.5$ on outlet 1, and $\min\{0.5, 0\} = 0$ on outlet 2. Their distance is $d_\Lambda(i, j) = 1 - 0.5 = 0.5$. This distance is a natural but necessary approximation of the actual time they spend on the same outlet. As in our data we only have information about which outlets the agent visits, but not the time she spends there.

Let $y_i \in \{1, 2..., 10\}$ stand for the answer of individual $i$ in the ideological self-identification question, a number between 1 to 10, where 1 is extreme left and 10 is extreme right. We define the social distance between $i$ and $j$ as $d_\Sigma(i, j) = |y_i - y_j|$. Let $\bar{d}_\Sigma$ be the average social distance of a uniform matching pairing, which is computed as in the previous application.

Our index is defined as:

$$S = \sum_{(i,j) \in \Pi} \mu(i, j) \rho(i, j) \tag{6.3}$$

for some linear functions $\mu$ and $\rho$ consistent with Definition 2 (see Section 2.2). Using theorem

---

[13]Questions QP17a,b,c y d.

[14]We omit websites from the analysis, as the consumption of websites is likely to include the consumption of online versions of radio and newspapers

[15]Media consumption is captured by questions QP17a,b,c and d, and ideology by question D1.

2, the measure is proportional to the covariance between social and spatial distances, i.e.

$$S = \frac{1}{\Pi} \sum_{(i,j) \in \Pi} d_\Sigma(i,j) d_\Lambda(i,j) - \overline{d}_\Sigma \overline{d}_\Lambda \tag{6.4}$$

We compute a normalized version of the index, which requires finding the expression above's minimal and maximal values. Maximal segregation is achieved by a media environment completely segregated, one in which all agents with the same ideology visit the same outlet. To compute this, we compute equation 6.4 for this configuration. As anticipated in Section 4, minimal segregation is calculated by solving the corresponding optimization problem numerically. In most countries, the latter's solution involves having some fully segregated groups and others uniformly distributed across media outlets. Details of the solution to the optimization problem can be found in Appendix B.

Table 3 shows the segregation index values for each country and media environment. The country with the highest segregation is Malta, which has the maximum value for the index in radio, TV, and newspaper markets. Cyprus has the lowest index value for TV and newspapers, while the Netherlands has the lowest segregation in radio stations. In half of the countries, the newspapers' market is the most segregated, while in other eight countries, the TV market is the most segregated.

In general, the correlation between segregation and the number of outlets is low but positive, ranging from 0.04 for the radio market to 0.28 for $TV$. This result is in line with theories predicting that increasing competition in the market for news, understood as increasing competitors, could exacerbate segregation. The main channel is that competition brings outlet differentiation and allows consumers to self-select more effectively into news outlets with like-minded opinions (see Mullainathan & Shleifer (2005)).

We find a high correlation of segregation across media markets. The correlation between segregation in radio and newspapers is 0.68, while for TV and newspapers, the number is 0.80. This might suggest that there could be structural political conditions associated with segregation. To explore this idea, we compare segregation indices for each environment with an index of polarization obtained from the Varieties of Democracy (V-Dem) Project (see Coppedge (2020)) As shown in Figure 2, segregation and polarization are positively correlated in all media outlets.
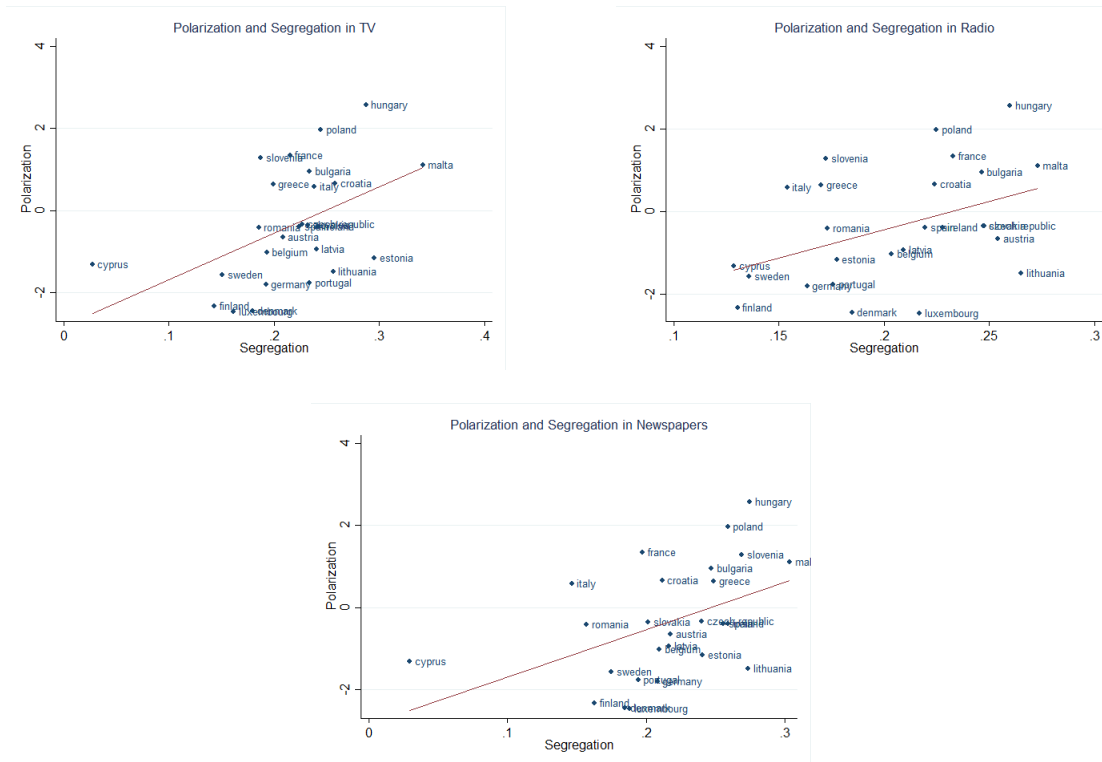
Figure 2: Polarization and Segregation

Table 2: Ideological self-identification in european countries, Eurobarometer 2014

| Country | Mean | Median | Std. Deviation | Sample size |
|---|---|---|---|---|
| Austria | 4.83 | 5 | 1.91 | 929 |
| Belgium | 5.08 | 5 | 2.00 | 924 |
| Bulgaria | 5.48 | 5 | 2.63 | 784 |
| Cyprus | 5.44 | 5 | 3.00 | 288 |
| Czech Republic | 5.39 | 5 | 2.36 | 949 |
| Germany | 4.89 | 5 | 1.75 | 1381 |
| Denmark | 5.47 | 5 | 2.35 | 955 |
| Estonia | 5.71 | 5 | 2.21 | 679 |
| Spain | 4.34 | 5 | 1.91 | 838 |
| Finland | 5.52 | 5 | 1.93 | 893 |
| France | 5.01 | 5 | 2.16 | 829 |
| Great Britain | 5.06 | 5 | 1.84 | 1092 |
| Greece | 5.06 | 5 | 2.08 | 729 |
| Croatia | 5.38 | 5 | 2.47 | 786 |
| Hungary | 5.84 | 5 | 2.32 | 842 |
| Ireland | 5.28 | 5 | 1.84 | 838 |
| Italy | 5.05 | 5 | 2.18 | 686 |
| Lithuania | 5.17 | 5 | 2.63 | 708 |
| Luxembourg | 5.25 | 5 | 1.94 | 411 |
| Latvia | 5.88 | 5 | 2.15 | 781 |
| Malta | 5.10 | 5 | 2.34 | 303 |
| The Netherlands | 5.08 | 5 | 1.85 | 969 |
| Poland | 5.94 | 5 | 2.42 | 799 |
| Portugal | 4.68 | 5 | 1.90 | 668 |
| Romania | 6.01 | 5 | 2.82 | 691 |
| Sweden | 5.22 | 5 | 2.35 | 1008 |
| Slovenia | 4.60 | 5 | 2.61 | 684 |
| Slovakia | 5.12 | 5 | 2.37 | 869 |

Table 3: Ideological Segregation in european countries, Eurobarometer 2014

| Country | TV $\hat{S}$ | TV L | Radio $\hat{S}$ | Radio L | Newspaper $\hat{S}$ | Newspaper L |
|---|---|---|---|---|---|---|
| Austria | 0.208 | 26 | 0.254 | 24 | 0.218 | 29 |
| Belgium | 0.193 | 26 | 0.204 | 29 | 0.210 | 27 |
| Bulgaria | 0.234 | 19 | 0.246 | 11 | 0.247 | 28 |
| Cyprus | 0.028 | 18 | 0.129 | 22 | 0.030 | 19 |
| Czech Republic | 0.227 | 23 | 0.247 | 18 | 0.240 | 29 |
| Germany | 0.192 | 28 | 0.163 | 27 | 0.208 | 28 |
| Denmark | 0.179 | 13 | 0.185 | 9 | 0.185 | 29 |
| Estonia | 0.295 | 22 | 0.177 | 25 | 0.241 | 29 |
| Spain | 0.224 | 23 | 0.219 | 22 | 0.256 | 28 |
| Finland | 0.143 | 15 | 0.130 | 22 | 0.163 | 28 |
| France | 0.215 | 27 | 0.233 | 22 | 0.197 | 27 |
| Great Britain | 0.247 | 26 | 0.235 | 24 | 0.293 | 28 |
| Greece | 0.199 | 25 | 0.170 | 26 | 0.249 | 25 |
| Croatia | 0.258 | 24 | 0.224 | 26 | 0.211 | 26 |
| Hungary | 0.287 | 19 | 0.260 | 21 | 0.275 | 27 |
| Ireland | 0.242 | 21 | 0.228 | 26 | 0.259 | 26 |
| Italy | 0.238 | 19 | 0.154 | 21 | 0.146 | 27 |
| Lithuania | 0.256 | 23 | 0.265 | 26 | 0.274 | 17 |
| Luxembourg | 0.162 | 5 | 0.217 | 12 | 0.188 | 25 |
| Latvia | 0.240 | 16 | 0.209 | 17 | 0.217 | 26 |
| Malta | 0.341 | 18 | 0.273 | 25 | 0.303 | 21 |
| Netherlands | 0.130 | 19 | 0.103 | 28 | 0.160 | 29 |
| Poland | 0.244 | 28 | 0.225 | 30 | 0.259 | 29 |
| Portugal | 0.234 | 24 | 0.176 | 23 | 0.194 | 29 |
| Romania | 0.185 | 29 | 0.173 | 27 | 0.157 | 19 |
| Sweden | 0.151 | 9 | 0.136 | 7 | 0.175 | 25 |
| Slovenia | 0.188 | 18 | 0.172 | 22 | 0.269 | 19 |
| Slovakia | 0.232 | 22 | 0.247 | 24 | 0.202 | 26 |
| Min | 0.028 | 5 | 0.103 | 7 | 0.030 | 17 |
| Max | 0.341 | 29 | 0.273 | 30 | 0.303 | 29 |
| Average | 0.213 | 21 | 0.202 | 22 | 0.215 | 26 |

# 7 Conclusion

This paper proposes a theory of segregation measurement based on the intensity and social diversity of pairwise interactions. In our framework, societies are described by a space of locations, a space of social groups, and agents' distribution across locations and groups. Both the space of locations and the space of social groups are flexible enough to include many different segregation problems. Locations can be schools in a district, residences in a city, or platforms such as media outlets where individuals interact. Social groups can be race, socioeconomic status, political ideology, or any other social identity. We axiomatize measures that can be expressed as a weighted sum across pairs of an interaction intensity that depends on locations and the value of pairwise interactions that relies on social identities. We prove that the index is proportional to the covariance between spatial and social distances.

We then use our segregation measures to study two segregation problems. First, we measure socioeconomic segregation in Chilean schools using Chilean micro-data, which includes information on each student's parents' socioeconomic status. There is variation across cities and grades, and school segregation highly correlates with residential segregation. As our index allows for multiple simultaneous interactions, in a second application, we use it to measure ideological segregation in media outlets' consumption for different media platforms -newspapers, radio, TV- for 28 European countries. There are systematic differences in segregation across countries and platforms, suggesting that some fundamental features, probably related to the political environment, explain these segregation levels. The stark correlation between our segregation measures and political polarization indexes suggests this might be the case. Further inquiry is required to understand this relationship better.

There are numerous possibilities for future research. First, the framework can be extended to consider other segregation problems, such as segregation in consumption patterns (e.g., cultural consumption). In contrast to the media consumption problem analyzed in this paper, this requires considering the consumption of goods in a continuum that differs in more than one dimension. Another analytical extension that could be explored is the structure of the social types' space. Specifically, in some segregation problems, we might be interested in multi-dimensional social types. This extension may allow us to understand the extent to which smaller groups and intersectionality drive racial or socioeconomic segregation. In concrete, it seems crucial to have a framework to understand if racial segregation is associated with the isolation of low-income people (as opposed to more affluent individuals that could be less segregated). Relaxing the social proximity axiom to allow for not monotonic measures in the social distance

29

and achieve a maximum at an intermediate level of social diversity is also left for future research. The framework presented in this framework is also well-suited for a richer understanding of the dynamics of segregation. By explicitly accounting for the space of social types and the landscape of interactions, it may allow decomposing time-changes in segregation associated with the social, market, or regulation forces that affect these primitives of the framework.

# A    Appendix: Main Proofs

**Proof of Theorem 1** From von Neumann Morgenstern utility representation theorem, a complete and transitive preference relation satisfies continuity and independence if and only if admits a expected utility representation:

$$S(\mu) = \sum_{(i,j)\in\Pi} \mu(i,j)\rho(i,j).$$

Using anonimity axiom we get the result. □

**Proof of Theorem 2** First, let's prove that axioms 1-7 imply $S = cov(d_\Lambda, d_\Sigma)$.

It is direct to see that axioms $1-3$ imply: $S = \sum_{(i,j)\in\Pi} \mu(i,j)\rho(i,j)$. By axioms 4 and 5,

$$S = \sum_{(i,j)\in\Pi} f(d_\Lambda(\lambda_i,\lambda_j))g(d_\Sigma(s_i,s_j)) \tag{A.1}$$

for some decreasing functions $f, g$. Moreover, by axioms 6 and 7, the functions $f(\cdot)$ and $g(\cdot)$ are linear. Then, there exit constants $m_0, m_1$ and $r_0, r_1$ such that:

$$S = \sum_{(i,j)\in\Pi} (m_0 - m_1 d_\Lambda(\lambda_i,\lambda_j))(r_0 - r_1 d_\Sigma(s_i,s_j)) \tag{A.2}$$

Note that $\sum_{(i,j)\in\Pi} \mu(i,j) = 1$ imposes a constraint over parameters $m_0, m_1$, such that:

$$m_0 = \frac{1}{\Pi} + m_1 \overline{d}_\Lambda \tag{A.3}$$

Plugging this in equation A.2 and with a little algebra we obtain the result:

$$
\begin{aligned}
S &= m_0 r_0 - r_0 m_1 \overline{d}_\Lambda(i,j) - r_1 m_0 \overline{d}_\Sigma + r_1 m_1 \frac{\sum_{(i,j)} d_\Lambda(\lambda_i,\lambda_j) d_\Sigma(s_i,s_j)}{\Pi} \tag{A.4}\\
&= \frac{r_0 - r_1 \overline{d}_\Sigma}{\Pi} + r_1 m_1 \left[ \frac{\sum_{(i,j)} d_\Lambda(\lambda_i,\lambda_j) d_\Sigma(s_i,s_j)}{\Pi} - \overline{d}_\Lambda \overline{d}_\Sigma \right] \tag{A.5}\\
&= \overline{\rho} + r_1 m_1 cov(d_\Lambda, d_\Sigma) \tag{A.6}
\end{aligned}
$$

which completes the proof in this direction.

Now suppose a segregation measure proportional to $S = cov(d_\Lambda, d_\Sigma)$. Note that the contribution of each interaction takes the form: $(\bar{d}_\Lambda - d_\Lambda(\lambda_i, \lambda_j))(\bar{d}_\Sigma - d_\Sigma(s_i, s_j))$, which is a decreasing function of the spatial and social distances. Moreover, any additive variation in $d_\Lambda(\lambda_j, \lambda_j)$ only changes the first component in an additive way, and same for additive variations in $d_\Sigma(s_i, s_j)$. Thus, axioms 4-7 hold. The proof of axioms 1-3 is analogous to the previous theorem. This completes the proof. $\square$

**Proof of Proposition 1** In order to prove the proposition, we first prove the following auxiliary lemma.

**Lemma 1** *Consider an assignment satisfying (ULA), and suppose axioms 1-8 hold. Then,* $\mu(i, j) = \frac{2}{N(N_l - 1)}$ *for any* $(i, j) \in \Pi_l$.

**Proof of Lemma 1** We already know that $\mu(i, j) = 0$ if $x_i \neq x_j$. Since $\mu(i, j) = \mu(d_\Lambda(i, j))$, we have that $\mu(i, j) = \mu_l$ for each $(i, j) \in \Pi_l$. With and individual resource constraint, the latter implies that $\mu_l(N_l - 1) = T$, from which $\mu_l = \frac{T}{N_l - 1}$. Now, since $\mu$ is a probability distribution,

$$\sum_{i,j} \mu(i, j) = \sum_{l=1}^{L} \Pi_l \mu_l,$$

where, with some abuse of notation, $\Pi_l = N_l(N_l - 1)/2$ is the number of pairs in $l$. Using the previous expression for $\mu_l$, we must have that $T = N/2$, which yields that $\mu(i, j) = \frac{2}{N(N_l - 1)}$ for any $(i, j) \in \Pi_l$. $\blacksquare$

From theorem 2 and lemma 1, the expression for $S$ becomes

$$S = \frac{2}{N} \sum_{l=1}^{L} \left(\frac{1}{N_l - 1}\right) \sum_{(i,j) \in \Pi_l} \rho(i, j) = \sum_{l=1}^{L} w_l \bar{\rho}_l,$$

where $w_l = \frac{N_l}{N}$ is the share of the population in location $l$, $\bar{\rho}_l = \frac{1}{\Pi_l} \sum_{(i,j) \in \Pi_l} \rho(i, j)$, and $\Pi_l = N_l(N_l - 1)$.

Now, with (DB), if $\rho(i, j) = r_0 - r_1 d^\Sigma(i, j)$, with $r_1 > 0$, then $\bar{\rho}_l = r_0 - r_1 \bar{d}_l^\Sigma$. Combining this with nn, we have

$$S = r_0 - r_1 \sum_{l=1}^{L} w_l \bar{d}_l^\Sigma.$$

Note that $S \leq r_0$, which is achieved with equality if and only if $\bar{d}_l^\Sigma = 0$ for all $l$. This is indeed the case if all agents in each location have the same social type. A sufficient condition to achieve

this maximal segregation is thus that $L = G$ is an admissible landscape. Normalizing maximal segregation to 1, it follows that $r_0 = 1$.

On the other hand, minimal segregation is obtained by maximizing the expression $\sum_{l=1}^{L} w_l \bar{d}_l^{\Sigma}$, as a function of the number of agents of each social type in each location. It can be shown that this es achieved by the population distributed uniformly across locations or by having all agents in the same location (which is equivalent to setting $L = 1$). In this case , for all $l$, $\bar{d}_l^{\Sigma} = \bar{d}^{\Sigma}$. Minimal segregation is thus $S^{min} = r_0 - r_1 \bar{d}^{\Sigma}$ and normalizing minimal segregation to 0, yields $r_1 = 1/\bar{d}^{\Sigma}$. $\square$

# B   Appendix: Minimal Segregation in Media Consumption

Let $G$ be the number of social types. For any two social types $s, s' \in \Sigma$, define $A_{s,s'} = d_\Sigma(s, s') - \bar{d}_\Sigma$.

Suppose assumption ULA holds. Let $N_s$ be the number of agents with social type $s$, and $N_{ls}$ the number of agents type $s$ in location $l$. We define $x_l = (x_{1,l}, ..., x_{G,l})' \in \mathbb{N}^G$ as a vector such that $x_{l,s} = N_{l,s}$. In this context, the minimization problem reduces to:

$$\min_{x_l \in \mathbb{N}^G} \quad \sum_{s \in \Sigma} \sum_{s' \in \Sigma} A_{s,s'} \cdot \sum_{l=1}^{L} N_{l,s} N_{l,s'}$$
$$\text{s.t.} \quad \sum_{l=1}^{L} N_s$$

This is the general problem to be solved. Note that the problem corresponds to minimizing a quadratic function over a simplex.

# C   Appendix: Examples and Additional results

## C.1   Community Operations: An Illustrative Example

When comparing segregation levels across communities, we keep fixed $(N, \Sigma, \rho)$, and study how changes in the distribution of agents over the space of interactions affect segregation.

Consider the study of school segregation by income. Students are characterized by their family income, which can be *low*, *middle* or *high*. We denote this social space by $\Sigma = \{l, m, h\}$. The proportion of each social group in the population are given by $\left(\frac{N_l}{N}, \frac{N_m}{N}, \frac{N_h}{N}\right) = (0.25, 0.5, 0.25)$. The space of locations $\Lambda$ is composed by eight schools, all of which have the same capacity.

Each agent $i$ has associated a social group $s^i \in \Sigma$, and a an assignment $x^i \in X$. We assume that students only interact with other students in the same school.

Consider three possible distributions of students across schools. Each of these distributions will generate a different community, which we denote by $C_1, C_2, C_3$. In Table 4 we illustrate the distribution of social groups across the space for each of them. Each column represents a different community, and then we compute segregation on each of them separately. Each row corresponds to a location in the space (i.e. a school). The triplets on each cell correspond to the share of low, middle and high income students on each school, for a given community. For instance, in the first community, schools 1 and 2 have only low income students, schools 3 to 6 only middle income, and schools 7 and 8 only high income. In community $C_2$ there is some mixing, so for instance half of the students in schools 1 and 2 have low income, and half have middle income.

Table 4: Distribution of students (As a share of school capacity)

|  | $C_1$ | $C_2$ | $C_3$ |
|---|---|---|---|
| School 1 | (1,0,0) | $(\frac{1}{2},\frac{1}{2},0)$ | (1, 0,0) |
| School 2 | (1,0,0) | $(\frac{1}{2},\frac{1}{2},0)$ | $(\frac{1}{2},\frac{1}{2},0)$ |
| School 3 | (0,1,0) | $(\frac{1}{2},0,\frac{1}{2})$ | $(\frac{1}{2},0,\frac{1}{2})$ |
| School 4 | (0,1,0) | $(\frac{1}{2},0,\frac{1}{2})$ | (0,1,0) |
| School 5 | (0,1,0) | (0,1,0) | (0,1,0) |
| School 6 | (0,1,0) | (0,1,0) | (0,1,0) |
| School 7 | (0,0,1) | $(0,\frac{1}{2},\frac{1}{2})$ | $(0,\frac{1}{2},\frac{1}{2})$ |
| School 8 | (0,0,1) | $(0,\frac{1}{2},\frac{1}{2})$ | (0,0,1) |

The distribution of students across schools generates probabilities of observing each type interaction. For instance, the probability of observing an interaction between two low income children in district $C_1$ is 1/4 (they can meet in two out of eight schools). This is the aggregate intensity function $\tilde{\mu}_\Sigma$ defined in equation 2.3, when $\mu$ is correctly normalized. Let $\tilde{\mu}^1_\Sigma$ be the aggregate intensity function of community $C_1$. This function satisfies $\tilde{\mu}^1_\Sigma(l,l) = \tilde{\mu}^1_\Sigma(h,h) = 1/4$; $\tilde{\mu}^1(m,m) = 1/2$, and $\tilde{\mu}^1_\Sigma(s,s') = 0$ for $s \neq s'$.

Following the same reasoning for the second community, we obtain $\tilde{\mu}^2_\Sigma(l,m) = \tilde{\mu}^2(l,h) = \tilde{\mu}^2_\Sigma(h,m) = \tilde{\mu}^2_\Sigma(m,m) = 1/4$, and $\tilde{\mu}^2_\Sigma(s,s) = 0$ for $s \in \{l,h\}$. It is clear that community 1 is more segregated than community 2: in the first one students with different income levels do not

interact, while in the second one there is some mixing.

Now suppose we are interested in combining communities 1 and 2, to obtain a community $C = \alpha C_1 + (1 - \alpha)C_2$ with $\alpha = 1/2$. This is analogous to combining the interaction intensities generated by communities $C_1$ and $C_2$: the new community $C$ is consistent with a new intensity $\tilde{\mu}$ such that $\tilde{\mu}(s, s') = \alpha\tilde{\mu}^1(s, s') + (1 - \alpha)\tilde{\mu}^2(s, s')$, for all $s, s' \in \Sigma$.

Note that this combination generates community $C_3$, represented in the third column in Table 4. The interaction intensities are as follows:

$$\tilde{\mu}^3(l, l) = \tilde{\mu}^3(l, m) = \tilde{\mu}^3(l, h) = \tilde{\mu}^3(m, h) = \tilde{\mu}^3(h, h) = \frac{1}{8}; \ \tilde{\mu}^3(m, m) = \frac{1}{4}. \tag{C.1}$$

## C.2 Equivalent Metrics

**Lemma 2** *Let $f, g : \mathbb{R} \to \mathbb{R}$ be increasing, continuous and subadditive functions. Then, the order of segregation is preserved for any metrics $d'_\Gamma = f(d_\Gamma)$, $d'_\Lambda = g(d_\Lambda)$.*

**Proof.** For any function $f$ we can do a first-order approximation around $E(d)$, so that

$$g(d) = g(E(d)) + (d - E(d)) \left. \frac{\partial g(d)}{\partial d} \right|_{E(d)}$$

Thus, without loss of generality fix $d_\Gamma$, and take $g(d_\Lambda)$. Then,

$$
\begin{aligned}
cov(g(d_\Lambda), d_\Gamma) &= cov\left( g(E(d_\Lambda)) + (d - E(d_\Lambda)) \left. \frac{\partial g(d_\Lambda)}{\partial d_\Lambda} \right|_{E(d_\Lambda)}, d_\Gamma \right) \\
&= cov\left( g(E(d_\Lambda)), d_\Gamma \right) + cov\left( (d - E(d_\Lambda)) \left. \frac{\partial g(d_\Lambda)}{\partial d_\Lambda} \right|_{E(d_\Lambda)}, d_\Gamma \right) \\
&= \left. \frac{\partial g(d_\Lambda)}{\partial d_\Lambda} \right|_{E(d_\Lambda)} cov\left( d_\Gamma, d_\Gamma \right)
\end{aligned}
$$

Thus, the order is preserved. Moreover, $g(d_\Gamma)$ and $d_\Gamma$ are equivalent metrics. Following the same reasoning for $d_\Lambda$, we get the result. $\qquad \square$

# References

Agostini, C., Hojman, D., Román, A. & Valenzuela, L. (2016), 'Segregación residencial de ingresos en el gran santiago, 1992-2002: Una estimación robusta', *EURE* **42**(127), 159–184.

Alesina, A. & La Ferrara, E. (2000), 'Participation in heterogeneous communities', *The Quarterly Journal of Economics* **115**(3), 847–904.

Alesina, A. & La Ferrara, E. (2005), 'Preferences for redistribution in the land of opportunities', *Journal of public Economics* **89**(5), 897–931.

Alonso-Villar, O. & Del Río, C. (2010), 'Local versus overall segregation measures', *Mathematical Social Sciences* **60**(1), 30–38.

Atkinson, A. B., Piketty, T. & Saez, E. (2011), 'Top incomes in the long run of history', *Journal of economic literature* **49**(1), 3–71.

Bartels, L. M. (2018), *Unequal democracy: The political economy of the new gilded age*, Princeton University Press.

Campante, F. R. & Hojman, D. A. (2013), 'Media and polarization: Evidence from the introduction of broadcast tv in the united states', *Journal of Public Economics* **100**, 79–92.

Coppedge, Michael, e. a. (2020), 'Varieties of democracy (v-dem) project'.

Correa, S. & Hojman, D. (2021), 'A characterization of minimal segregation', *Working Paper* .

DellaVigna, S. & Kaplan, E. (2007), 'The fox news effect: Media bias and voting', *The Quarterly Journal of Economics* **122**(3), 1187–1234.

Duncan, O. D. & Duncan, B. (1955), 'A methodological analysis of segregation indexes', *American sociological review* **20**(2), 210–217.

Easterly, W. & Levine, R. (1997), 'Africa's growth tragedy: policies and ethnic divisions', *The Quarterly Journal of Economics* **112**(4), 1203–1250.

Echenique, F. & Fryer Jr, R. G. (2007), 'A measure of segregation based on social interactions', *The Quarterly Journal of Economics* **122**(2), 441–485.

Frankel, D. M. & Volij, O. (2011), 'Measuring school segregation', *Journal of Economic Theory* **146**(1), 1–38.

Gentzkow, M. & Shapiro, J. M. (2011), 'Ideological segregation online and offline', *The Quarterly Journal of Economics* **126**(4), 1799–1839.

Hong, L. & Page, S. E. (2001), 'Problem solving by heterogeneous agents', *Journal of Economic Theory* **97**(1), 123–163.

Hong, L. & Page, S. E. (2004), 'Groups of diverse problem solvers can outperform groups of high-ability problem solvers', *Proceedings of the National Academy of Sciences of the United States of America* **101**(46), 16385–16389.

James, D. R. & Taeuber, K. E. (1985), 'Measures of segregation', *Sociological Methodology* **15**, 1–32.

Lamont, M., Beljean, S. & Clair, M. (2014), 'What is missing? cultural processes and causal pathways to inequality', *Socio-Economic Review* **12**(3), 573–608.

Levy, G. & Razin, R. (2019), 'Echo chambers and their effects on economic and political outcomes', *Annual Review of Economics* **11**, 303–328.

Massey, D. S. & Denton, N. A. (1988), 'The dimensions of residential segregation', *Social Forces* **67**(2), 281–315.

Mullainathan, S. & Shleifer, A. (2005), 'The market for news', *American Economic Review* **95**(4), 1031–1053.

Owens, A., Reardon, S. F. & Jencks, C. (2016), 'Income segregation between schools and school districts', *American Educational Research Journal* **53**(4), 1159–1197.

Putnam, R. D. (2007), 'E pluribus unum: Diversity and community in the twenty-first century the 2006 johan skytte prize lecture', *Scandinavian political studies* **30**(2), 137–174.

Reardon, S. F. & Owens, A. (2014), '60 years after brown: Trends and consequences of school segregation', *Annual Review of Sociology, 40:1, 199-218* .

Sandel, M. J. (2020), *The Tyranny of Merit: What's Become of the Common Good?*, Allen Lane London.

Stroud, N. J. (2008), 'Media use and political predispositions: Revisiting the concept of selective exposure', *Political Behavior* **30**(3), 341–366.