

EFFICIENCY AND FAIRNESS OF SYSTEM-OPTIMAL ROUTING WITH USER CONSTRAINTS*

ANDREAS S. SCHULZ¹ AND NICOLÁS E. STIER-MOSES²

August 2004

ABSTRACT. We study the route-guidance system proposed by Jahn, Möhring, Schulz, and Stier-Moses (2004) from a theoretical perspective. This approach computes a traffic pattern that minimizes the total travel time subject to user constraints, which ensure that routes suggested to users are not much longer than shortest paths. We show that when distances are measured with respect to travel times at equilibrium, the resulting traffic assignment is efficient and fair.

1. INTRODUCTION

Transportation authorities and users alike hope that route-guidance systems can help to mitigate the congestion generated by the ever-increasing amount of vehicular traffic. In particular, in-car navigation devices might be used not only to provide drivers with map information and current traffic conditions, but also to optimize the entire network. With this application in mind, Jahn, Möhring, Schulz, and Stier-Moses (2004) introduced a route-guidance system which computes a traffic pattern that minimizes the total travel time subject to certain user constraints. These constraints are designed to overcome an inherent problem of system-optimal guidance. Indeed, it is well known that in a system-optimal flow, some users may be assigned to considerably longer routes for the benefit of others. User constraints are intended to guarantee that no recommended route is significantly longer than that suggested to any other user with the same origin and destination. For the sake of algorithmic efficiency, Jahn et al. (2004) proposed to perform this comparison based on *normal path lengths* instead of actual travel times. The normal length of a path is defined via some a priori estimate of the real impedance. Based on extensive computational studies on real-world

Key words and phrases. Selfish Routing, Price of Anarchy, Computational Game Theory, Multicommodity Flows, Route Guidance, Traffic Assignment.

* An extended abstract of a preliminary version appeared in the Proceedings of the 14th Annual ACM-SIAM Symposium on Discrete Algorithms, Baltimore, MD, January 12–14, 2003.

¹ Andreas S. Schulz, Sloan School of Management, Massachusetts Institute of Technology, Office E53-361, 77 Massachusetts Avenue, Cambridge, MA 02139-4307, Email: schulz@mit.edu

² Nicolás E. Stier-Moses, Columbia Business School, Uris Hall, Room 418, 3022 Broadway, New York, NY 10027-6902, Email: nicolas.stier@columbia.edu.

instances, Jahn et al. concluded that the resulting *constrained* system-optimal flow has two desirable properties when the normal length is properly chosen: the total travel time in the network is close to that of the unconstrained system optimum, and individual users do not experience a notably larger travel time than others. While normal lengths have to be set in advance, it is important to point out that users do not need to know the normal lengths; they are merely an artifact to select solutions without detours and are controlled by the traffic authority. Jahn et al. considered three possible definitions: *geographic* lengths, *free-flow* travel times (travel times in the uncongested network), and *user equilibrium* travel times (travel times at user equilibrium). As geographic distances and free-flow travel times are highly correlated, we confine our study to free-flow travel times and user equilibrium travel times.

In this note, we provide a theoretical framework for the work of Jahn et al. For our analysis, we rely on the *price of anarchy* concept, discussed in this context by Koutsoupias and Papadimitriou (1999), Roughgarden and Tardos (2002), Roughgarden (2003b), and Correa, Schulz, and Stier-Moses (2004b), among others. While the price of anarchy is originally defined as the ratio of the total travel time of a user equilibrium to that of an ordinary system optimum, we adopt a more pragmatic perspective for the application considered here. Instead of comparing the user equilibrium to a system-optimal traffic flow, which can realistically not be implemented because of its unfairness, we measure the price of anarchy with respect to a traffic pattern that can potentially be used in practice. In other words, we evaluate the efficiency of the route-guidance system with the help of the worst-case ratio of the total travel time of an equilibrium to that of a constrained system optimum. In addition, we compare the travel times experienced by different users between the same origin-destination pair. A primary goal of any route-guidance system is to offer routes with similar travel times; otherwise, route recommendations would most likely be dismissed. Moreover, as the system assigns users to routes randomly, it is desirable that routes offered under similar circumstances (e.g., the same user performs a trip every day) have similar latencies to reduce the variance of latencies experienced by individual users.

Section 2 formally introduces constrained system optima and the price of anarchy. In addition, we study instances for which the corresponding equilibria have high total cost. In Section 3, we concentrate on free-flow travel times and argue that, under these normal lengths, user equilibria might be preferable over constrained system optima. Section 4 analyzes the case in which normal lengths are defined as user equilibrium travel times. In contrast to the previous case, constrained

system optima turn out to be provably efficient and fair. Section 5 shows that the unfairness of constrained system optima is bounded from above by a small constant. All results established here corroborate the conclusions drawn by Jahn et al. (2004)

2. PRELIMINARIES

The road network is represented by a directed multigraph $G = (N, A)$ with two attributes on each arc $a \in A$: the normal length $\tau_a \geq 0$ gives an a priori estimate of the actual traversal time in the solution we seek; the latency function $\ell_a : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ maps the traffic flow x_a on arc a to its traversal time $\ell_a(x_a)$. The normal arc lengths can be any metric that is fixed in advance. However, their proper choice will allow us to compute solutions that users of the route-guidance device are likely to accept. The latency functions are assumed to be continuous and nondecreasing. These assumptions are naturally met by common latency functions (see, e.g., Sheffi 1985; Cohen 1991). We only consider latencies that belong to a specified set \mathcal{L} . In practice, latency functions take a specific form, and the bounds one obtains without this restriction are unnecessarily pessimistic (Roughgarden and Tardos 2002). For some results, we will additionally assume that latencies are linear functions; i.e., they belong to $\mathcal{L}_{\text{lin}} := \{\ell : \ell(x) = qx + r \text{ for some } q, r \geq 0\}$. Although this may appear restrictive at first, congestion effects and counterintuitive phenomena can already occur (e.g., Braess 1968).

Vehicles are grouped according to their origin-destination (OD) pairs $K \subseteq N \times N$. For each OD pair $k = (s_k, t_k) \in K$, let \mathcal{P}_k be the set of directed (simple) paths in G from s_k to t_k , and let $d_k > 0$ be the demand rate associated with commodity k . Let $\mathcal{P} := \bigcup_{k \in K} \mathcal{P}_k$ be the set of paths between all OD pairs. Because route-guidance systems eventually have to propose paths to the drivers, our formulation is path-based: A feasible flow x assigns a nonnegative value x_P to every path $P \in \mathcal{P}$ such that $\sum_{P \in \mathcal{P}_k} x_P = d_k$ for all $k \in K$. Note that flows are not required to be integral; they describe average rates. Furthermore, we define the latency of a path $P \in \mathcal{P}$ under a given flow x as $\ell_P(x) := \sum_{a \in P} \ell_a(x_a)$, where $x_a := \sum_{Q \in \mathcal{P}} x_Q$. We refer to the maximum latency of a flow-carrying path in \mathcal{P}_k as $L_k(x) := \max\{\ell_P(x) : P \in \mathcal{P}_k, x_P > 0\}$, and to the shortest normal path length for OD pair k as $T_k := \min_{P \in \mathcal{P}_k} \tau_P$. Here, $\tau_P := \sum_{a \in P} \tau_a$ is the normal length of path P .

There are two aspects that define the quality of a flow. The fairness of the route assignment is of importance to the users while the total travel time in the system is of importance to the

traffic authority. Without any control (but with perfect information), different users with the same OD pair should experience the same travel time. If that was not the case, users would have an incentive to switch routes. We say that a flow with this property is *fair*. In a seminal contribution, Wardrop (1952) stated a principle that formalizes this notion: “The journey times on all the routes actually used are equal, and less than those which would be experienced by a single vehicle on any unused route.” Traffic patterns satisfying this principle are called *user equilibria* (Dafermos and Sparrow 1969), and will be denoted by f . Although there may be multiple equilibria, the travel time that users experience is invariant across different equilibria (Beckmann, McGuire, and Winsten 1956). In particular, all equilibria share the same total cost. We define the *unfairness* of a given flow x as the maximum ratio of the experienced travel times of two users sharing the same OD pair, i.e., as $\max\{\ell_Q(x)/\ell_R(x) : Q, R \in \mathcal{P}_k, x_Q, x_R > 0, k \in K\}$.

It is well known that user equilibria can be inefficient (Dupuit 1849) and that system optima can be unfair (Ben-Akiva et al. 1997). The route-guidance system proposed by Jahn et al. (2004) is designed to compute the most efficient traffic assignment among those that are not too unfair. As it is difficult to control the unfairness directly, the authors introduced an upper bound on the ratio of the normal length of different users traveling between the same OD pair. Let us describe this approach in more detail. Let $\varphi \geq 1$ quantify the notional tolerance of users to suboptimal paths. This factor is used to prevent users from being assigned to paths that are longer than φ times the length of a shortest path between their OD pair. In other words, users of OD pair $k \in K$ may only be assigned to paths in $\mathcal{P}_k^\varphi := \{P \in \mathcal{P}_k : \tau_P \leq \varphi T_k\}$. We call a path feasible when it belongs to $\mathcal{P}^\varphi := \bigcup_{k \in K} \mathcal{P}_k^\varphi$. Selecting the solution with minimum total travel time from all assignments of users to paths in \mathcal{P}^φ is equivalent to solving the following minimum cost multicommodity flow problem with path constraints and a separable nonlinear objective function:

$$\begin{aligned}
 \text{(CSO)} \quad & \min \quad C(x) := \sum_{a \in A} \ell_a(x_a) x_a \\
 & \text{s.t.} \quad \sum_{P \in \mathcal{P}_k^\varphi} x_P = d_k \quad \text{for all } k \in K, \\
 & \quad \quad \sum_{P: a \in P \in \mathcal{P}^\varphi} x_P = x_a \quad \text{for all } a \in A, \\
 & \quad \quad x_P \geq 0 \quad \text{for all } P \in \mathcal{P}^\varphi.
 \end{aligned}$$

A *constrained system optimum* with tolerance factor φ , denoted by f^φ , is an optimal solution to this problem. A *system optimum* is an optimal solution to a similar problem without path constraints, i.e., with $\varphi = \infty$. We use f^* to denote a system optimum. It is evident that the larger is the factor φ , the larger is the feasible region. Consequently, $C(f^\varphi)$ is a nonincreasing function of φ and $C(f^*) \leq C(f^\varphi)$ for all $\varphi \geq 1$. When $C(x)$ is convex, Beckmann, McGuire, and Winsten (1956) proved that a flow f^* is a system optimum with respect to latencies $\ell_a(x)$ if and only if it is a user equilibrium with respect to the modified latencies $\ell_a^*(x) := \ell_a(x) + x\ell'_a(x)$. The latency functions ℓ_a^* include an extra term that accounts for the service degradation caused to the other users of arc a . As f^* is at equilibrium with respect to ℓ_a^* , we denote the common travel time for OD pair $k \in K$ by $L_k^*(f^*)$.

2.1. The Price of Anarchy. Koutsoupias and Papadimitriou (1999) proposed to measure the inefficiency of equilibria by comparing the worst equilibrium to the system optimum. Roughgarden and Tardos (2002), Roughgarden (2003b), and Correa, Schulz, and Stier-Moses (2004b) studied the efficiency of user equilibria in traffic networks, under less and less restrictive assumptions. In particular, the cost $C(f)$ of a user equilibrium is bounded from above by $\alpha(\mathcal{L})$ times the cost $C(f^*)$ of a system optimum (Roughgarden 2003b). The constant $\alpha(\mathcal{L})$, defined in Section 2.2, depends only on the set of allowed latency functions; for example, it is $4/3$ for linear, 1.626 for quadratic, and 1.896 for cubic functions, respectively. For polynomials of degree p , it grows like $p/\ln p$.

To measure the potential benefits of the route-guidance system proposed by Jahn et al. (2004), we compare the cost of user equilibria to that of constrained system optima. Although the original definition of the price of anarchy relies on ordinary system optima, our notion is arguably more realistic in our context because ordinary system optima cannot be implemented in practice because of their unfairness. In the following definition, $f_{\mathcal{I}}$ and $f_{\mathcal{I}}^\varphi$ denote the user equilibrium and constrained system optimum of an instance \mathcal{I} , respectively. To simplify notation, we drop the subindex \mathcal{I} afterwards. The price of anarchy for a given tolerance factor φ and a given set \mathcal{L} of allowed latency functions is defined as follows:

$$\alpha^\varphi(\mathcal{L}) := \sup_{\mathcal{I} \in \text{inst}(\mathcal{L})} \frac{C(f_{\mathcal{I}})}{C(f_{\mathcal{I}}^\varphi)}. \quad (2.1)$$

Here, $\text{inst}(\mathcal{L})$ is the set of instances with latency functions drawn from \mathcal{L} . It is immediately clear that $\alpha^1(\mathcal{L}) \geq 1$ and that $\alpha^\varphi(\mathcal{L})$ is a nondecreasing function of φ . In addition, the bounds previously

mentioned imply that

$$C(f) \leq \alpha(\mathcal{L}) C(f^*) \leq \alpha(\mathcal{L}) C(f^\varphi). \quad (2.2)$$

Equivalently, $\alpha^\varphi(\mathcal{L}) \leq \alpha(\mathcal{L})$ for all $\varphi \geq 1$. Moreover, for instances with positive minimum normal length T_k for all OD pairs $k \in K$, a constrained system optimum with large tolerance is optimal in the unconstrained sense; i.e., $C(f^\varphi) = C(f^*)$ when φ is sufficiently large.

2.2. Tight Instances. For a specific instance, we refer to the ratio of the cost of a user equilibrium to that of a constrained system optimum as the *coordination ratio* of the instance. Aiming at an understanding of the cause of inefficiency, we now characterize instances with high coordination ratio. While we work with ordinary system optima in this section, we will concentrate on arbitrary tolerance factors later. We call an instance *tight* when its coordination ratio $C(f)/C(f^*)$ matches the upper bound $\alpha(\mathcal{L})$. Here, f and f^* are a user equilibrium and a system optimum of the corresponding instance, respectively. Roughgarden and Tardos (2002) presented an example, originally suggested by Pigou (1920), that is tight for linear latency functions. Subsequently, Roughgarden (2003b) proved that the maximum coordination ratio among all instances with two parallel arcs and latencies drawn from \mathcal{L} matches the price of anarchy. By a careful analysis of Correa, Schulz, and Stier-Moses' proof of that result, we can establish conditions that characterize instances that are tight in the unconstrained case. For completeness, we include the proof below. Let us first define parameters that help us to express $\alpha(\mathcal{L})$ easily: setting $\beta(v, \ell) := \max \left\{ \frac{x}{v} (1 - \ell(x)/\ell(v)) : 0 \leq x \leq v \right\}$, and $\beta(\mathcal{L}) := \sup \{ \beta(v, \ell) : v \geq 0 \text{ and } \ell \in \mathcal{L} \}$, it can be proved that $\alpha(\mathcal{L}) = (1 - \beta(\mathcal{L}))^{-1}$.

Before the proof, let us present the mentioned tight example, which will also be helpful for subsequent results. Figure 1 illustrates the network, which consists of two nodes connected by two parallel arcs and a demand rate equal to d . For a given function $\ell \in \mathcal{L}$, the latency of the top arc is constant at $\ell(d)$, while the bottom arc has a flow-dependent latency of $\ell(x)$. The user equilibrium f assigns the entire demand d to the lower arc, whereas the system optimum f^* routes $d - x^*$ on the top arc. Here, $x^* := \arg \max x(\ell(d) - \ell(x))$. The costs are $C(f) = \ell(d)d$ and $C(f^*) = \ell(d)(d - x^*) + \ell(x^*)x^*$, respectively. Taking the supremum of the ratio leads precisely to $\alpha(\mathcal{L})$.

Theorem 2.1 (Correa, Schulz, and Stier-Moses 2004b, Theorem 3.6). *Consider an instance with latency functions drawn from a set of continuous and nondecreasing latency functions \mathcal{L} . Then, $C(f) \leq (1 - \beta(\mathcal{L}))^{-1} C(f^*)$.*

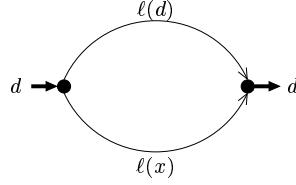


FIGURE 1. Simple tight instance

Proof. The claim follows from

$$C(f) \leq \sum_{a \in A} \ell_a(f_a) f_a^* \leq \sum_{a \in A} \beta(f_a, \ell_a) \ell_a(f_a) f_a + \sum_{a \in A} \ell_a(f_a^*) f_a^* \leq \beta(\mathcal{L}) C(f) + C(f^*). \quad (2.3)$$

The first inequality is valid because f routes flow along shortest paths with respect to $\ell_a(f_a)$; the other two follow directly from the definitions of β . \square

Observation 2.2. *Let \mathcal{L} be a family of continuous and nondecreasing latency functions. An unconstrained instance with latency functions drawn from the set \mathcal{L} is tight if and only if the following three conditions are satisfied:*

$$\text{for all } k \in K \text{ and } P \in \mathcal{P}_k : f_P^* > 0 \Rightarrow \ell_P(f) = L_k(f), \quad (2.4a)$$

$$\text{for all } a \in A : f_a^* = \arg \max_{x \geq 0} x(\ell_a(f_a) - \ell_a(x)), \quad (2.4b)$$

$$\text{for all } a \in A : \ell_a(f_a) f_a > 0 \Rightarrow \beta(f_a, \ell_a) = \beta(\mathcal{L}). \quad (2.4c)$$

Proof. An instance is tight if and only if all inequalities in the proof of the previous theorem are equalities. The conditions correspond to the three inequalities in (2.3), in the order in which they appear. \square

Let us make a few remarks related to Observation 2.2:

- (i) When latency functions are differentiable, setting the derivative of the right-hand side of Condition (2.4b) to zero, we see that $\ell_a^*(f_a^*) = \ell_a(f_a)$. This implies that $\ell_P^*(f^*) = \ell_P(f)$ for all $P \in \mathcal{P}$, and that $L_k^*(f^*) = L_k(f)$ for all $k \in K$. For example, when latencies are linear, the user equilibrium and system optimum of a tight instance must satisfy $f_a^* = f_a/2$ for all arcs $a \in A$ with the exception of those that have constant latency.
- (ii) Under the additional assumptions in (i), Condition (2.4a) is implied by the optimality of f^* and the prior remark.

- (iii) For arcs with strictly increasing latency functions, Condition (2.4b) implies that either $f_a^* < f_a$ or $f_a^* = f_a = 0$.
- (iv) Assume that $\beta(\mathcal{L}) > 0$ (otherwise the price of anarchy is 1 and all instances are tight). If an arc a carries flow in the user equilibrium, it must satisfy $\ell_a(0) = 0$.

Proof. Suppose $\ell_a(0) > 0$. Let \hat{x} be the argument that maximizes $x(\ell_a(f_a) - \ell_a(x))$ in the definition of $\beta(f_a, \ell_a)$. Because $0 < \ell_a(0) \leq \ell_a(f_a)$, Condition (2.4c) implies that $\beta(f_a, \ell_a) > 0$ meaning that $\ell_a(\hat{x}) < \ell_a(f_a)$. Therefore,

$$\frac{\hat{x}}{f_a} \left(1 - \frac{\ell_a(\hat{x})}{\ell_a(f_a)} \right) < \frac{\hat{x}}{f_a} \left(1 - \frac{\ell_a(\hat{x}) - \ell_a(0)}{\ell_a(f_a) - \ell_a(0)} \right).$$

This shows that $\beta(f_a, \ell_a) < \beta(f_a, \ell_a - \ell_a(0))$, which is a contradiction to Condition (2.4c). \square

The following lemma gives a necessary condition for an instance to be tight.

Lemma 2.3. *Let \mathcal{L} be a family of continuous and nondecreasing latency functions that are either constant or strictly increasing. If $\ell_a(0) = 0$ for all $a \in A$, then the instance cannot be tight.*

Proof. Given the assumption, there are two classes of arcs: those with latencies identically equal to 0 (we refer to them as 0-latency arcs), and those with strictly increasing latencies. Consider a user equilibrium, a system-optimal flow, and an OD pair $k \in K$. Let us define a set $C_k := \{i \in N : \text{there is a path from } s_k \text{ to } i \text{ using 0-latency arcs}\}$. If $t_k \in C_k$, both flows route the demand of OD pair k along a 0-latency path. If this happens for all OD pairs, the instance cannot be tight. Therefore, consider an OD pair $k \in K$ for which $t_k \notin C_k$. Thus, C_k defines an s_k - t_k -cut. By construction of C_k , a forward arc in the cut cannot be a 0-latency arc, and a backward arc cannot carry (user or system optimal) flow for OD pair k . The former is obvious; to see the latter, if a backward arc a carried flow, a would be a 0-latency arc because all flow reaching C_k must use paths with latency equal to zero. Thus, its tail would belong to C_k , too. Hence, there is no flow of OD pair k entering C_k , and all flow exits C_k along non 0-latency arcs. This is a contradiction to Remark (iii) because the sum of the flow on forward arcs of the cut C_k must equal the demand d_k . \square

Under the assumptions of the previous lemma, if an instance is tight, it must have an arc satisfying $\ell_a(0) > 0$. As Remark (iv) prevents that arc from carrying flow in a user equilibrium, either it carries flow in a system optimum or it is not used by either flow. (In the latter case, we can remove it and get a smaller tight instance.)

3. FREE-FLOW TRAVEL TIMES AS NORMAL LENGTHS

In this section, we assume that normal lengths are defined as travel times in the uncongested network; i.e., $\tau_a = \ell_a(0)$ for all $a \in A$. The theoretical results presented next help to explain the conclusion derived from the computational study of real-world instances by Jahn et al. (2004): user equilibria have improved performance guarantees when compared to constrained system optima rather than system optima. In fact, for small values of φ , constrained system optima can even be worse than user equilibria.

Let us study the function $\alpha^\varphi(\mathcal{L})$ to understand how the price of anarchy depends on the users' tolerance to unfairness. We start by proving a structural property that implies that the price of anarchy is subadditive. For this purpose, we introduce a construction that enables us to modify the tolerance factor of an instance without altering its coordination ratio too much. For a fixed tolerance factor φ , consider an instance \mathcal{I} with large coordination ratio; i.e., the instance satisfies

$$\frac{C(f)}{C(f^\varphi)} \geq \alpha^\varphi(\mathcal{L}) - \varepsilon, \quad (3.1)$$

for some $\varepsilon > 0$. We construct a new instance $\tilde{\mathcal{I}}$, which is equal to \mathcal{I} except for the following modifications. The origins in $\tilde{\mathcal{I}}$ are new vertices \tilde{s}_k for $k \in K$ (instead of s_k), which are connected to s_k with arcs of constant latency M_k , specified below. The natural extension \tilde{x} of a flow x to the new instance is defined as $\tilde{x}_{\tilde{P}} := x_P$ for $P \in \mathcal{P}_k$ and $k \in K$, where \tilde{P} starts at \tilde{s}_k and then continues with the original path P . The extensions \tilde{f} and \tilde{f}^* of a user equilibrium f and a system optimum f^* of \mathcal{I} are a user equilibrium and a system optimum for $\tilde{\mathcal{I}}$, respectively. The next lemma establishes a relation between the constrained system optima of the two instances.

Lemma 3.1. *Consider a fixed $\tilde{\varphi}$ such that $1 < \tilde{\varphi} < \varphi$, and set $M_k := \frac{\varphi - \tilde{\varphi}}{\tilde{\varphi} - 1} T_k$. If f^φ is a φ -constrained system optimum of \mathcal{I} , then its natural extension \tilde{f}^φ is a $\tilde{\varphi}$ -constrained system optimum of $\tilde{\mathcal{I}}$.*

Proof. All paths in \mathcal{P}_k that carry flow under f^φ have a normal length between T_k and φT_k . After adding M_k to each of them, their lengths are between $M_k + T_k$ and $M_k + \varphi T_k = \tilde{\varphi}(M_k + T_k)$. It follows that \tilde{f}^φ is a $\tilde{\varphi}$ -constrained system optimum. \square

Observe that extending a flow x of \mathcal{I} to a flow \tilde{x} of $\tilde{\mathcal{I}}$ changes its cost by a fixed amount M ; that is, $C(\tilde{x}) = M + C(x)$ with $M := \sum_{k \in K} M_k d_k$. Moreover, for this choice of normal lengths,

$T_k \leq \ell_P(x)$ for all $P \in \mathcal{P}_k$, which implies that

$$M = \frac{\varphi - \tilde{\varphi}}{\tilde{\varphi} - 1} \sum_k T_k d_k \leq \frac{\varphi - \tilde{\varphi}}{\tilde{\varphi} - 1} C(x). \quad (3.2)$$

We can now prove that the price of anarchy cannot increase too fast.

Theorem 3.2. *The function $\alpha^\varphi(\mathcal{L})/(\varphi - 1)$ is nonincreasing in φ .*

Proof. Consider the instance \mathcal{I} with large coordination ratio that we selected in (3.1), and let f be a user equilibrium and f^φ be a φ -constrained system optimum. Furthermore, their natural extensions to $\tilde{\mathcal{I}}$ are referred to as \tilde{f} and \tilde{f}^φ , respectively. We bound the price of anarchy of the new instance $\tilde{\mathcal{I}}$ with that of the original instance \mathcal{I} :

$$\alpha^{\tilde{\varphi}}(\mathcal{L}) \geq \frac{C(\tilde{f})}{C(\tilde{f}^\varphi)} = \frac{M + C(f)}{M + C(f^\varphi)} \geq \frac{C(f)}{\frac{\varphi - 1}{\tilde{\varphi} - 1} C(f^\varphi)} \geq \frac{\tilde{\varphi} - 1}{\varphi - 1} (\alpha^\varphi(\mathcal{L}) - \varepsilon) \quad \text{for all } \tilde{\varphi} < \varphi.$$

The inequalities follow from (2.1), (3.2) and (3.1), respectively. As ε can be made arbitrarily small, $\alpha^{\tilde{\varphi}}(\mathcal{L}) \geq \frac{\tilde{\varphi} - 1}{\varphi - 1} \alpha^\varphi(\mathcal{L})$ for all $\tilde{\varphi} < \varphi$. \square

The last theorem implies that the price of anarchy is subadditive as a function of δ , where $\delta \geq 0$ is a modified tolerance factor defined as $\varphi - 1$.

Corollary 3.3. *The function $\alpha^{1+\delta}(\mathcal{L})$ is subadditive in δ .*

3.1. Bad Instances. In this section, we extend the results presented in Section 2.2. We call an instance *tight* if $C(f)/C(f^\varphi)$ matches the upper bound $\alpha(\mathcal{L})$, where f and f^φ are a user equilibrium and a constrained system optimum of the corresponding instance, respectively. We will use Observation 2.2 and Lemma 2.3 to show that, under mild assumptions, there cannot exist tight instances for the constrained case. Note that this does not prevent $\alpha^\varphi(\mathcal{L})$ from being equal to $\alpha(\mathcal{L})$ for some φ .

Theorem 3.4. *Consider an instance with latency functions drawn from a set \mathcal{L} of continuous and nondecreasing latency functions that are either strictly increasing or constant. Then, the coordination ratio $C(f)/C(f^\varphi) < \alpha(\mathcal{L})$ for all $\varphi \geq 1$.*

Proof. Suppose that the coordination ratio equals $\alpha(\mathcal{L})$. In this case, f^φ is a system optimum because the cost of the system optimum is a lower bound for that of f^φ , and the coordination ratio cannot be larger than $\alpha(\mathcal{L})$. From Remark (iv) in Section 2.2, we know that $\ell_a(0) = 0$ for all arcs a with $f_a > 0$. Hence, there is a path joining each OD pair with free-flow travel time equal to zero.

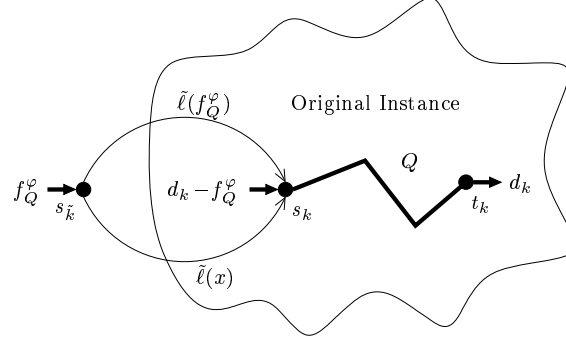


FIGURE 2. Modified instance used in the proof of Theorem 3.5

In other words, the normal length T_k has to be 0 for all $k \in K$, which implies that a path belongs to \mathcal{P}_k^φ only if its normal length is zero. Therefore, $\ell_a(0) = 0$ for all arcs a with flow in f or f^φ , contradicting Lemma 2.3. \square

We turn our attention to characterizing instances with large coordination ratio for fixed φ . We say that a path $P \in \mathcal{P}_k$ is *longest* if its normal length τ_P equals the maximum possible value φT_k .

Theorem 3.5. *Consider a family \mathcal{L} of differentiable latency functions that is closed under multiplication by constants. Furthermore, assume that \mathcal{L} is closed under scaling, i.e., if $\ell \in \mathcal{L}$, $\bar{\ell} : x \rightarrow \ell(\rho x)$ belongs to \mathcal{L} for all $\rho \geq 0$. Let $\varphi \geq 1$ and f^φ be a φ -constrained system optimum of a given instance with latencies drawn from \mathcal{L} . If f^φ routes flow along a path that is not longest, then the instance can be modified to increase the coordination ratio $C(f)/C(f^\varphi)$.*

Proof. Assume that for some $k \in K$ there is path $Q \in \mathcal{P}_k^\varphi$ that is not *longest* such that $f_Q^\varphi > 0$. We insert the network shown in Figure 1 at the source s_k , the origin of path Q . (This modification is illustrated in Figure 2.) After the modification, two parallel arcs connect a new origin $s_{\tilde{k}}$ to s_k . Furthermore, f_Q^φ units of demand are reassigned from OD pair k to a new OD pair \tilde{k} with terminals $s_{\tilde{k}}$ and t_k . We will now construct latency functions that will make the added network tight. Let $v > 0$ and $\ell \in \mathcal{L}$ be such that $\beta(v, \ell) = \beta(\mathcal{L})$. Consider the latency function $\tilde{\ell}$ defined as $\tilde{\ell}(x) := \ell(xv/f_Q^\varphi)/M$, where $M > 0$ is a constant to be specified later. Note that the assumptions imply that $\tilde{\ell} \in \mathcal{L}$. It can be seen that $\beta(f_Q^\varphi, \tilde{\ell}) = \beta(\mathcal{L})$ as well. We let the latencies of the new arcs be the constant $\tilde{\ell}(f_Q^\varphi)$ and the function $\tilde{\ell}(x)$. After the modification, there are two possible extensions of Q : the path Q_\uparrow (resp. Q_\downarrow) starts with the constant (resp. nonconstant) arc just added and continues along Q . Denoting the set of paths in the new instance by $\tilde{\mathcal{P}}$ and setting M such

that $\tilde{\ell}(f_Q^\varphi) + \tau_Q < \varphi T_k \leq \varphi T_{\tilde{k}}$, τ_{Q_\uparrow} and τ_{Q_\downarrow} are bounded from above by $\varphi T_{\tilde{k}}$. Thus, Q_\uparrow and Q_\downarrow belong to $\tilde{\mathcal{P}}_{\tilde{k}}$.

The user equilibrium \tilde{f} of the new instance can be constructed as an extension of f . We reassign (any) f_Q^φ units of flow from OD pair k to OD pair \tilde{k} , route them along the new nonconstant arc, and let them follow their original paths. It is straightforward to see that \tilde{f} is feasible and at equilibrium. Similarly, the constrained system optimum \tilde{f}^φ of the new instance can be obtained easily from f^φ . Indeed, we distribute the flow originally in Q along the paths Q_\downarrow and Q_\uparrow in a way that satisfies $\tilde{\ell}_{Q_\downarrow}^*(\tilde{f}^\varphi) = \tilde{\ell}_{Q_\uparrow}^*(\tilde{f}^\varphi)$. As $\tilde{f}_a^\varphi = f_a^\varphi$ for the original arcs a , $\tilde{\ell}_P^*(\tilde{f}^\varphi) = \tilde{\ell}_P^*(f^\varphi)$ for all $P \in \tilde{\mathcal{P}}_j^\varphi$ with $j \neq \tilde{k}$. Therefore, this extension is a constrained system optimum. Computing the total travel times of both flows and using that the subgraph is tight, the coordination ratio $C(\tilde{f})/C(\tilde{f}^\varphi)$ equals

$$\frac{C(f) + \tilde{\ell}(f_Q^\varphi)f_Q^\varphi}{C(f^\varphi) + \tilde{\ell}(f_Q^\varphi)f_Q^\varphi/\alpha(\mathcal{L})}.$$

Note that the last expression is a convex combination of $C(f)/C(f^\varphi)$ and $\alpha(\mathcal{L})$. As the former is smaller than the latter, the new instance has worse performance. \square

The previous theorem is still valid if latencies are not differentiable. To establish this result, we cannot use the modified latencies ℓ^* , but the extension can be done using the system optimum of the tight subinstance. The theorem also remains valid for user equilibrium normal lengths. We will state this explicitly in Section 4.

3.2. Bounds for the Price of Anarchy. In this section, we present upper and lower bounds for the function $\alpha^\varphi(\mathcal{L}_{\text{lin}})$. We start with an upper bound that improves on $\alpha^\varphi(\mathcal{L}_{\text{lin}}) \leq \alpha(\mathcal{L}_{\text{lin}}) = 4/3$.

Theorem 3.6. *The price of anarchy $\alpha^\varphi(\mathcal{L}_{\text{lin}}) \leq (2 - \varphi)^{-1}$ for all $1 \leq \varphi < 2$. In particular, $\alpha^1(\mathcal{L}_{\text{lin}}) = 1$ and $\alpha^\varphi(\mathcal{L}_{\text{lin}}) < 4/3$ for $\varphi < 5/4$.*

Proof. Consider a factor $1 \leq \varphi < 2$, and let f^φ and f be a φ -constrained system optimum and a user equilibrium, respectively. We define the function $h : [0, 1] \rightarrow \mathbb{R}$ by $h(z) := C(f + z(f^\varphi - f))$. Due to the convexity of $C(\cdot)$, $h(1) \geq h(0) + h'(0)$. To prove the claim we verify that $h(0) + h'(0) \geq (2 - \varphi)h(0)$

because then $C(f^\varphi) = h(1) \geq (2 - \varphi)h(0) = (2 - \varphi)C(f)$, as required. Now,

$$\begin{aligned} h'(0) &= \sum_a \ell_a^*(f_a)(f_a^\varphi - f_a) = \sum_a [2\ell_a(f_a) - \ell_a(0)](f_a^\varphi - f_a) \\ &\geq 2\left(\sum_k L_k(f)d_k - \sum_k L_k(f)d_k\right) + \sum_k T_k d_k - \varphi \sum_k T_k d_k \\ &= (1 - \varphi) \sum_k T_k d_k \geq (1 - \varphi)C(f) = (1 - \varphi)h(0). \end{aligned}$$

The first inequality follows from the fact that $\ell_P(f) = L_k(f)$ for every $P \in \mathcal{P}_k$ such that $f_P > 0$, and $\ell_P(f) \geq L_k(f)$ in general. Moreover, $\tau_P \leq \varphi T_k$ for every P such that $f_P^\varphi > 0$, and $T_k \leq \tau_P$ in general. \square

We now give lower bounds for $\alpha^\varphi(\mathcal{L}_{\text{lin}})$ by providing corresponding instances. Although the tight instance in Figure 1 can be used, a stronger bound can be given with a collection of instances based on the Braess Paradox network (Braess 1968).

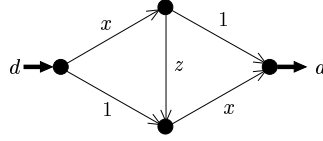


FIGURE 3. Instance used in Lemma 3.7

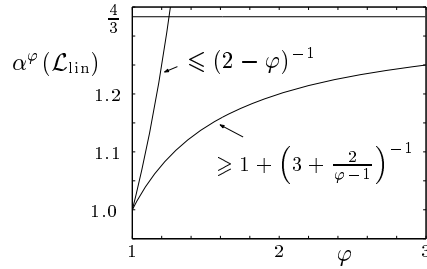
Lemma 3.7. *The price of anarchy $\alpha^\varphi(\mathcal{L}_{\text{lin}}) \geq 1 + \left(3 + \frac{2}{\varphi-1}\right)^{-1}$.*

Proof. Consider the network depicted in Figure 3, where $z \geq 0$ is a constant and the demand between the single OD pair is $d \geq 0$. Maximizing the coordination ratio over z and d , we obtain the claim. \square

Figure 4 summarizes the bounds for $\alpha^\varphi(\mathcal{L}_{\text{lin}})$. The main conclusion is that $\alpha^\varphi(\mathcal{L}_{\text{lin}})$ is relatively close to 1 when φ is close to 1. Therefore, we do not need to compute the exact value of the price of anarchy to establish that user equilibria are almost as good as constrained system optima. We conclude that, if free-flow travel times are used, central coordination is not beneficial: the route-guidance system does not deliver solutions of significantly improved quality.

4. USER EQUILIBRIUM TRAVEL TIMES AS NORMAL LENGTHS

In this section, we assume that normal lengths are set equal to the travel times experienced in a user equilibrium. Based on superior empirical performance compared to free-flow travel times,

FIGURE 4. Bounds for $\alpha^\varphi(\mathcal{L}_{\text{lin}})$

Jahn et al. (2004) concluded that these normal lengths are the “correct” choice. We now provide a theoretical foundation for their findings. The improvement with respect to free-flow normal lengths comes from the fact that user equilibrium normal lengths depend on the actual demand.

A user equilibrium f is a feasible solution to (CSO) because all paths used in f are feasible. Therefore, the constrained system optimum f^φ satisfies

$$C(f^\varphi) \leq C(f) \quad \text{for all } \varphi \geq 1. \quad (4.1)$$

As in the previous section, we obtain a lower bound for the function $\alpha^\varphi(\mathcal{L})$ by providing an appropriate instance. However, in this case, the lower bound matches the upper bound.

Lemma 4.1. *The price of anarchy $\alpha^\varphi(\mathcal{L}) = \alpha(\mathcal{L})$ for all $\varphi \geq 1$.*

Proof. In the equilibrium of the instance depicted in Figure 1, both paths have latency $\ell(d)$. Hence, regardless of the value of φ , the system optimum is a φ -constrained system optimum. The claim follows by taking the supremum over $\ell \in \mathcal{L}$. \square

The lemma implies that the price of anarchy is the same regardless of whether it is defined with respect to the system optimum or the constrained system optimum. The two bounds presented in (2.2) and (4.1) can be tight. For example, the proof of Lemma 4.1 describes an instance for $\varphi = 1$ satisfying $C(f^*) = C(f^1) = C(f)/\alpha(\mathcal{L})$. On the other hand, if we add a small constant $\varepsilon > 0$ to the latency of the first arc, the constrained system optimum coincides with the user equilibrium. Therefore, $C(f) = C(f^1) \approx \alpha(\mathcal{L})C(f^*)$.

Finally, let us point out that Theorem 3.5, proved before for free-flow normal lengths, is also valid when normal lengths are set to user equilibrium travel times. It is enough to note that at equilibrium the lengths of the two new arcs equal $\tilde{\ell}(f_Q^\varphi)$; therefore, Q_\downarrow and Q_\uparrow belong to $\tilde{\mathcal{P}}^\varphi$.

Observation 4.2. *Consider a family \mathcal{L} of differentiable latency functions that is closed under multiplication by constants. Furthermore, assume that \mathcal{L} is closed under scaling. Let $\varphi \geq 1$ and f^φ be a φ -constrained system optimum of a given instance with latencies drawn from \mathcal{L} . If f^φ routes flow along a path that is not longest, then the instance can be modified to increase the coordination ratio $C(f)/C(f^\varphi)$.*

5. FAIRNESS

As we mentioned earlier, a common argument against using the system optimum in the design of route-guidance devices for traffic assignment is that it generally assigns some drivers to unacceptably long paths in order to use shorter paths for most other drivers. This section presents results related to the unfairness of system optima and constrained system optima. In this section, we work with arbitrary normal lengths, unless explicitly stated otherwise.

The following theorem quantifies the severity of this effect by characterizing the unfairness of the system optimum. It turns out that there is a relation to earlier work by Roughgarden (2002), who compared the maximum latency of a system optimum in a single-sink single-source network to the latency of a user equilibrium. He showed that for a given class of latency functions \mathcal{L} , this ratio is bounded from above by $\gamma(\mathcal{L})$. Here, $\gamma(\mathcal{L})$ is defined to be the smallest value that satisfies $\ell^*(x) \leq \gamma(\mathcal{L})\ell(x)$ for all $\ell \in \mathcal{L}$ and $x \geq 0$. For example, $\gamma(\{\text{polynomials of degree } p \text{ with positive coefficients}\}) = p + 1$. Roughgarden (2003a) and Correa, Schulz, and Stier-Moses (2004a) independently proved that the unfairness of a system optimum is in fact bounded by the same constant, even for general instances with multiple commodities. This bound implies Roughgarden's earlier result: for the single-source single-sink case, $L(f^*) \leq \gamma(\mathcal{L})L(f)$.

It is not difficult to extend the bound on the unfairness of system optima to constrained system optima. Notice that the following theorem does not assume any particular definition of normal lengths.

Theorem 5.1. *Let f^φ be a constrained system optimum in a multicommodity flow network with latency functions drawn from a family \mathcal{L} of differentiable and nondecreasing latency functions. Then, the unfairness of f^φ is bounded from above by $\gamma(\mathcal{L})$.*

Proof. Using the definitions of ℓ^* and $\gamma(\mathcal{L})$, it is clear that $\ell_a(x) \leq \ell_a^*(x) \leq \gamma(\mathcal{L})\ell_a(x)$ for all $x \geq 0$. The first-order optimality conditions of (CSO) imply that for a constant $L_k^*(f^\varphi)$, $\ell_P^*(f^\varphi) = L_k^*(f^\varphi)$

for all $P \in \mathcal{P}_k^\varphi$ such that $f_P^\varphi > 0$. Therefore, for all paths $P \in \mathcal{P}_k^\varphi$ carrying flow,

$$\frac{L_k^*(f^\varphi)}{\gamma(\mathcal{L})} \leq \ell_P(f^\varphi) \leq L_k^*(f^\varphi).$$

Consequently, $\ell_Q(f^\varphi)/\ell_R(f^\varphi) \leq \gamma(\mathcal{L})$ for all $Q, R \in \mathcal{P}_k^\varphi$ with positive flow. \square

Correa, Schulz, and Stier-Moses (2004a) presented an instance that can be used to show that the bound given in the last theorem is tight. Consider the instance shown in Figure 1 with $d = 1$ and $\ell(x) := x$. User equilibrium normal lengths are equal to 1 for both arcs; therefore, both paths are feasible regardless of the value of φ . This means that any constrained system optimum is an unconstrained system optimum and its unfairness is $\gamma(\mathcal{L}_{\text{lin}}) = 2$. Nevertheless, in practice these bounds are loose, as the extensive experiments performed by Jahn et al. (2004) show. Note that Theorem 5.1 does not imply that the unfairness of constrained system optima with factor φ is nondecreasing as a function of φ . We now present two examples corresponding to the two definitions of normal lengths that we studied. The example using free-flow travel times is depicted in Figure 5. The instance has unit demand, and two nodes are connected by three arcs with latencies $1 + \varepsilon$, $1 + x$ and $1 + x^2$, respectively. A constrained system optimum with $\varphi = 1$ can only route flow in the last two arcs and therefore has an unfairness strictly larger than 1. For $\varphi \geq 1 + \varepsilon$, all arcs can be used, and it is easy to see that the value of unfairness approaches 1 when $\varepsilon \rightarrow 0$.

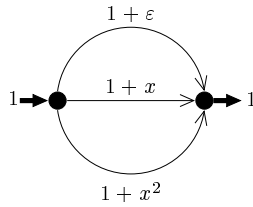


FIGURE 5. The unfairness may decrease when φ increases (normal length = free-flow travel time)

For the case when normal lengths are user equilibrium travel times, consider the instance shown in Figure 6. There are five arcs a , b , c , d , and e with latencies $x + 2\varepsilon$, 1 , 2ε , $1 + 5\varepsilon$, and x , respectively. The user equilibrium routes flow only along paths ab and ace ; at equilibrium the path de is too long to carry flow. Therefore, the constrained system optimum f^1 can only use paths ab and ace , and its unfairness is $\frac{4+4\varepsilon}{3+6\varepsilon}$. If $\varphi \geq \frac{2+3\varepsilon}{2+2\varepsilon}$, the constrained system optimum f^φ can use all three paths. In that case, it routes the flow along ab and cd , and its unfairness is $\frac{6+17\varepsilon}{6+11\varepsilon}$. For small

enough values of ε , the unfairness of the constrained system optimum with $\varphi = 1$ is arbitrarily close to $4/3$ while the unfairness for a large enough tolerance factor φ is arbitrarily close to 1.

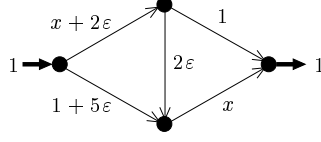


FIGURE 6. The unfairness may decrease when φ increases (normal length = user equilibrium travel times)

Finally, we show that if a system-optimal solution uses paths that are short with respect to free-flow normal lengths, they must also be short with respect to experienced travel times. We let the average latency of a flow x between OD pair $k \in K$ be defined as $\bar{C}_k(x) := \sum_{P \in \mathcal{P}_k} \ell_P(x) x_P / d_k$.

Theorem 5.2. *Consider an instance with linear latencies. Let $P \in \mathcal{P}_k$ be a path satisfying $f_P^* > 0$ and $\ell_P(0) \leq \varepsilon \bar{C}_k(f^*)$ for some $\varepsilon \geq 0$. Then, the experienced travel time $\ell_P(f^*)$ is bounded from above by $(1 + \varepsilon/2) \bar{C}_k(f^*)$.*

Proof. As latencies are linear functions and a system optimum is an equilibrium with respect to latencies ℓ^* , we have

$$\begin{aligned} \ell_P(f^*) &= \ell_P^*(f^*) - \sum_{a \in P} q_a f_a^* = \left(\sum_{Q \in \mathcal{P}_k} \frac{f_Q^* \ell_Q^*(f^*)}{d_k} \right) - (\ell_P(f^*) - \ell_P(0)) \\ &\leq \left(\sum_{Q \in \mathcal{P}_k} 2 \frac{f_Q^* \ell_Q^*(f^*)}{d_k} \right) - \ell_P(f^*) + \varepsilon \bar{C}_k(f^*). \end{aligned}$$

Therefore, $2\ell_P(f^*) \leq (2 + \varepsilon) \bar{C}_k(f^*)$. □

REFERENCES

- Beckmann, M. J., C. B. McGuire, and C. B. Winsten (1956). *Studies in the Economics of Transportation*. Yale University Press, New Haven, CT.
- Ben-Akiva, M. E., M. Bierlaire, J. Bottom, H. N. Koutsopoulos, and R. G. Mishalani (1997). Development of a route guidance generation system for real-time application. In M. Papanaghiou and A. Pouliezios (Eds.), *Proceedings of the 8th IFAC Symposium on Transportation Systems*, Chania, Greece, pp. 405–410. Elsevier Science, Oxford.

- Braess, D. (1968). Über ein Paradoxon aus der Verkehrsplanung. *Unternehmensforschung* 12, 258–268.
- Cohen, S. (1991). Flow variables. In M. Papageorgiou (Ed.), *Concise Encyclopedia of Traffic & Transportation Systems*, pp. 139–143. Pergamon Press, Oxford.
- Correa, J. R., A. S. Schulz, and N. E. Stier-Moses (2004a). Computational complexity, fairness, and the price of anarchy of the maximum latency problem. In D. Bienstock and G. Nemhauser (Eds.), *Proceedings of the 10th Conference on Integer Programming and Combinatorial Optimization (IPCO)*, New York, NY, Volume 3064 of *Lecture Notes in Computer Science*, pp. 59–73. Springer, Berlin.
- Correa, J. R., A. S. Schulz, and N. E. Stier-Moses (2004b). Selfish routing in capacitated networks. *Mathematics of Operations Research*. To appear. A preliminary version appeared as Working Paper No. 4319-03, Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA, June 2003.
- Dafermos, S. C. and F. T. Sparrow (1969). The traffic assignment problem for a general network. *Journal of Research of the U.S. National Bureau of Standards* 73B, 91–118.
- Dupuit, J. (1849). On tolls and transport charges. *Annales des Ponts et Chaussées*. Reprinted in *International Economic Papers* 11 (1962), 7–31.
- Jahn, O., R. H. Möhring, A. S. Schulz, and N. E. Stier-Moses (2004). System-optimal routing of traffic flows with user constraints in networks with congestion. *Operations Research*. To Appear. A preliminary version appeared as Working Paper No. 4394-02, Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA, November 2002.
- Koutsoupias, E. and C. H. Papadimitriou (1999). Worst-case equilibria. In C. Meinel and S. Tiison (Eds.), *Proceedings of the 16th Annual Symposium on Theoretical Aspects of Computer Science (STACS)*, Trier, Germany, Volume 1563 of *Lecture Notes in Computer Science*, pp. 404–413. Springer, Berlin.
- Pigou, A. C. (1920). *The Economics of Welfare*. Macmillan, London.
- Roughgarden, T. (2002). How unfair is optimal routing? In *Proceedings of the 13th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, San Francisco, CA, pp. 203–204. SIAM, Philadelphia, PA.
- Roughgarden, T. (2003a). Personal communication.

- Roughgarden, T. (2003b). The price of anarchy is independent of the network topology. *Journal of Computer and System Sciences* 67, 341–364.
- Roughgarden, T. and É. Tardos (2002). How bad is selfish routing? *Journal of the ACM* 49, 236–259.
- Sheffi, Y. (1985). *Urban Transportation Networks*. Prentice-Hall, Englewood, NJ.
- Wardrop, J. G. (1952). Some theoretical aspects of road traffic research. *Proceedings of the Institution of Civil Engineers, Part II, Vol. 1*, 325–378.