# Revealing the Autonomous System Taxonomy:
# The Machine Learning Approach

Xenofontas Dimitropoulos[1], Dmitri Krioukov[2], George Riley[1], and KC Claffy[2]

[1] Georgia Institute of Technology
`fontas@ece.gatech.edu`
`riley@ece.gatech.edu`
[2] Cooperative Association for Internet Data Analysis (CAIDA)
`dima@caida.org`
`kc@caida.org`

**Abstract.** Although the Internet AS-level topology has been extensively studied over the past few years, little is known about the details of the AS taxonomy. An AS "node" can represent a wide variety of organizations, e.g., large ISP, or small private business, university, with vastly different network characteristics, external connectivity patterns, network growth tendencies, and other properties that we can hardly neglect while working on veracious Internet representations in simulation environments. In this paper, we introduce a radically new approach based on machine learning techniques to map all the ASes in the Internet into a natural AS taxonomy. We successfully classify 95.3% of ASes with expected accuracy of 78.1%. We release to the community the AS-level topology dataset augmented with: 1) the AS taxonomy information and 2) the set of AS attributes we used to classify ASes. We believe that this dataset will serve as an invaluable addition to further understanding of the structure and evolution of the Internet.

## 1 Introduction

The rapid expansion of the Internet in the last two decades has produced a large-scale system of thousands of diverse, independently managed networks that collectively provide global connectivity across a wide spectrum of geopolitical environments. From 1997 to 2005 the number of globally routable AS identifiers has increased from less than 2,000 to more than 20,000, exerting significant pressure on interdomain routing as well as other functional and structural parts of the Internet. This impressive growth has resulted in a heterogenous and highly complex system that challenges accurate and realistic modeling of the Internet infrastructure. In particular, the AS-level topology is an intermix of networks owned and operated by many different organizations, e.g., backbone providers, regional providers, access providers, universities and private companies. Statistical information that faithfully characterizes different AS types is on the critical path toward understanding the structure of the Internet, as well as for modeling its topology and growth.

In topology modeling, knowledge of AS types is mandatory for augmenting synthetically constructed or measured AS topologies with realistic intra-AS and inter-AS router-level topologies. For example, we expect the network of a dual-homed

university to be drastically different from that of a dual-homed small company. The university will likely contain dozens of internal routers, thousands of hosts, and many other network elements (switches, servers, firewalls). On the other hand, the small company will most probably have a single router and a simple network topology. Since there is such a diversity among different network types, we cannot accurately augment the AS-level topology with appropriate router-level topologies if we cannot characterize the composing ASes.

Moreover, annotating the ASes in the AS topology with their types is a prerequisite for modeling the evolution of the Internet, since different types of ASes exhibit different growth patterns. For example, Internet Service Providers (ISP) grow by attracting new customers and by engaging in business agreements with other ISPs. On the other hand, small companies that connect to the Internet through one or few ISPs do not grow significantly over time. Thus, categorizing different types of ASes in the Internet is necessary to identify network evolution patterns and develop accurate evolution models.

An AS taxonomy is also necessary for mapping IP addresses to different types of users. For example, in traffic analysis studies its often required to distinguish between packets that come from home and business users. Given an AS taxonomy, its possible to realize this goal by checking the type of AS that originates the prefix in which an IP address lies.

In this work, we introduce a radically new approach based on machine learning to construct a representative AS taxonomy. We develop an algorithm to classify ASes based on empirically observed differences between AS characteristics. We use a large set of data from the Internet Routing Registries (IRR) [12] and from Route-Views [9] to identify intrinsic differences between ASes of different types. Then, we employ a novel machine learning technique to build a classification algorithm that exploits these differences to classify ASes into six representative classes that reflect ASes with different network properties and infrastructures. We derive macroscopic statistics on the different types of ASes in the Internet and validate our results using a sample of 1200 manually identified AS types. Our validation demonstrates that our classification algorithm achieves high accuracy: 78.1% of the examined classifications were correct. Finally, we make our results and our classifier publicly available to promote further research and understanding of the Internet's structure and evolution.

In Section 2 we start with a brief discussion of related work. Section 3 describes the data we used, and in Section 4 we specify the set of AS classes we use in our experiments. Section 5 introduces our classification approach and results. We validate them in Section 6 and conclude in Section 7.

## 2   Related Work

Several works have developed techniques decomposing the AS topology into different levels or tiers based on connectivity properties of BGP-derived AS graphs. Govindan and Reddy [8] propose a classification of ASes into four levels based on their AS degree. Ge *et al.* [7] classify ASes into seven tiers based on inferred customer-

to-provider relationships. Their classification exploits the idea that provider ASes should be in higher tiers than their customers. Subramanian *et al.* [14] classify ASes into five tiers based on inferred customer-to-provider as well as peer-to-peer relationships.

Our work differs from previous approaches in the following ways:

1. We do not employ heuristics and ad-hoc thresholds to define the AS levels. Instead, we use a novel machine learning algorithm to identify intrinsic features distinguishing different AS types.
2. We do not rely exclusively on AS graphs, which often miss a substantial fraction of the true AS links in the Internet, resulting in incomplete AS topologies. Instead, we use an extensive set of diverse data including IRR records, inferred AS relationships, IP prefixes, and AS graphs.
3. We do not classify ASes into hierarchies of levels or tiers extracted from AS graphs using degree-based or more sophisticated mechanisms; these methods tend to mix ASes with substantially different network properties into a single AS group.[3] Instead, we specify a representative set of AS classes characterized by unique signatures of network properties.

## 3   Data Sources and AS Attributes

To construct the set of AS attributes that we use in our AS classification, we collect data from the following databases and measurement projects:

1) **IRRs** [12]. The IRRs constitute a distributed database containing records on ASes' routing policies, assigned IP prefixes, contact information, etc. A natural approach to identifying the type of an AS, given its AS number, is to lookup the AS number in the IRRs and examine its organization description record. In the RPSL [1] terminology, this record is the `descr` attribute of the RPSL class `aut-num`. It contains the name or a short description of the organization that owns the AS number. For example, the following are entries for the `descr` attribute found in the IRRs: "Intervivo Networks, a broadband Internet access provider" and "Auckland Peering Exchange". The `descr` attribute does not have a standard representation. It usually consists of a short description as in the examples above, but in some cases it only contains an acronym, e.g., "KPMG LLP", "LTI". For the purposes of this work, the `descr` record is our first AS attribute, from which we extract useful information by means of text classification techniques. We downloaded the mapping of AS numbers to organization description records on 04/08/2005 from the CIDR Report [2], which provides on a daily basis mappings of AS numbers to organization description records extracted from ARIN, RIPE NCC, LAPNIC, APNIC, KRNIC, TWNIC, and JPNIC databases. We preprocess the organization records by removing stop words, i.e., words with little semantic meaning, such as "of" or "the", using

---

[3] According to our analysis, small regional providers often have small AS degrees, as low as 1 or 2. The previous heuristics thus tend to assign these ASes to the lowest levels, where small companies and multihomed customers naturally reside.

the stop word list [10]. Then, using the Porter stemming algorithm [11], we replace words with their stem.

We note that IRRs contain significant portions of incomplete or obsolete records, which is not a serious problem for this study since we are only concerned with the `descr` attribute, which changes relatively rarely.

2) **RouteViews** [9]. RouteViews is a measurement project that collects and archives a union of BGP tables from a large number of ISPs. We download all 12 BGP table snapshots archived from the collector `route-views2.oregon-ix.net` on 07/18/2005. For each table snapshot we extract AS paths and remove AS sets and private AS numbers. Then, we merge the extracted AS paths into an AS topology and use the AS relationship inference heuristics of [5, 4] to annotate the AS links with customer-to-provider and peer-to-peer relationships. Having the AS relationships inferred, we calculate the following three attributes for every AS: the number of providers, the number of customers, and the number of peers a given AS is connected to. Large ISPs typically have a large number of customers, zero providers, and a small number of peers, while small ISPs typically have few providers, a small number of customers and a large number of peers. Stub university or company networks typically have zero customers, zero or few peers and a small number of providers.

From the RouteViews data, we also extract information on IP prefixes. We use the chronologically first table snapshot from the same snapshot set to construct a mapping of AS numbers to the IP prefixes they advertise. Then, for each AS we count the number of advertised prefixes and use this number as another AS attribute. Small ASes, with tiny portions of IP address space allocated to them, as well as older ASes with large IP blocks, tend to advertise few prefixes. On the contrary, large ASes with relatively high IP address utilization and diversified routing policies tend to advertise many prefixes of various lengths.

One problem with this attribute is that IP prefixes are of drastically different sizes. Advertised IPv4 prefixes range in size from a /8, covering $2^{24}$ IP addresses, down to a /32, covering just one address. The prefix length of 24 bits is generally the smallest IPv4 prefix size that is globally routed, which suggests our last AS attribute to be the number of unique /24 prefixes found within the union of address space advertised by the AS. This attribute is likely to have small values for smaller ASes that use and advertise smaller portions of IP address space, while it is likely to be at its maximum for large or old ASes (those that appeared in the Internet early, e.g., some military and academic networks) since they often have huge chunks of assigned IP address space that they utilize scarcely.

In summary, we collect data from the IRRs and RouteViews. We find 19,537 ASes. Using the collected data, we annotate every AS with the following six attributes: 1) the organization description record (the *description* attribute), 2) the number of inferred customers (the *customer* attribute), 3) the number of inferred providers (the *provider* attribute), 4) the number of inferred peers (the *peer* attribute), 5) the number of advertised IP prefixes (the *prefix* attribute), and 6) the equivalent number of /24 prefixes covering all the advertised IP space (the *space* attribute).

## 4 The AS Class Set

In this work, we focus on network properties of an AS as the main criterion determining the set of AS classes. In other words, to construct the AS class set, we use the rule that ASes in the same class should have similar network properties, while ASes in different classes should have different network properties. ASes in the same class may still have significant network differences, however these differences should be small compared with differences between networks of different classes. For example, a small university and a large university may have quite different networks, however these differences are less significant compared to the differences between a typical university network and a typical ISP network.

Besides employing our *de facto* empirical knowledge, we perform the following experiment to specify the set of AS classes. We *randomly* select 1200 ASes and then, for each AS, we examine its attributes, visit its website (if possible), search for references to its organization name and determine its business profile. After examining the spectrum of these 1200 ASes, we construct our AS class set:

1. *Large ISPs*: Large backbone providers, tier-1 ISPs, with intercontinental networks.
2. *Small ISPs*: Regional and access providers with small metropolitan or larger regional networks.
3. *Customer ASes*: Companies or organizations that run their own networks but as opposed to members of the previous two classes do not provide Internet connectivity services. We find a wide range of ASes in this class, like web hosting companies, technology companies, consulting companies, hospitals, banks, military networks, government networks, etc.
4. *Universities*: University or college networks. We distinguish these networks from members of the Customer AS class, since they typically have substantially larger networks that serve thousands of end hosts.
5. *Internet exchange points (IXPs)*: Small networks serving as interconnection points for the members of the first two classes.
6. *Network information centers (NICs)*: Networks that host important network infrastructure, such as root or TLD servers.

## 5 AS Classification: Algorithms and Results

We build our classification algorithm using *AdaBoost* [6], a very powerful machine learning technique. The main idea behind AdaBoost is to combine multiple simple classification rules into a efficient composite classifier. These simple classification rules are called *weak hypotheses* and by definition are only required to perform slightly better than random guessing. Intuitively, weak hypotheses reflect simple "rules of thumb" that are usually much easier to construct than a complex classifier. AdaBoost works iteratively over a set of training examples; at each iteration it finds a weak hypothesis that performs well on the examples which the weak hypotheses of previous iterations erroneously classified. One constructs a weak hypothesis by means of a *weak learning algorithm* or simply *weak learner*. The power of AdaBoost

**Fig. 1**: *AdaBoost.MH* pseudocode

---

**Input**: $S = \{(x_1, y_1), \ldots, (x_m, y_m)\}$
**Initialize**: $D_1(x, y) = 1/mk$;      // $k$ is the total number of classes
**for** $t = 1$ **to** $T$ **do**

    Pass distribution $D_t$ and examples $S$ to weak learner
    Get weak hypothesis $h_t : X \times Y \to \mathbb{R}$
    Update distribution

$$D_{t+1}(x, y) = \frac{D_t(x, y)exp(-P(x, y)h_t(x, y))}{Z_t}$$

    // where $Z_t$ is a normalization coefficient
    // chosen so that $D_{t+1}$ will be a distribution
**end**
**Output**: $f(x, y) = \sum_{t=1}^{T} h_t(x, y)$

---

lies in a well-developed theoretical framework that intelligently combines the weak hypotheses into a composite classifier.

Let $X$ denote the set of ASes that we want to classify and $Y$ be the set of possible classes, such that each AS $x \in X$ belongs to a unique class in $Y$. If $y \in Y$ is the correct class for AS $x$, than let $P(x, y) = 1$, otherwise $P(x, y) = -1$. The goal is to find a classifier that for each AS produces a *ranking* of all the possible classes. More formally, we will compute a ranking function $f : X \times Y \to \mathbb{R}$: for each AS $x \in X$, the classes in $Y$ will be ordered according to $f(x, \cdot)$—the higher the value of $f(x, y)$, the more likely $x$ belongs to $y$.

In Figure 1, we illustrate the AdaBoost.MH algorithm [13], a special member of the AdaBoost algorithm family that is suited for solving multiclass problems. Let $S = \{(x_1, y_1), \ldots, (x_m, y_m)\}$ be the set of training examples. In our case, we construct $S$ by manually determining the correct class of 1220 ASes: 1200 are randomly selected; and 20 are well-known large ISPs and IXPs, which we use to increase the number of these two types of ASes in the initial random sample. Let $T$ be the total number of iterations. For each iteration $t = 1 \ldots T$, we maintain a distribution $D_t$ of weights over the set of examples and classes $D_t : X \times Y \to \mathbb{R}$. At the first iteration, we initialize $D_t$ to the uniform distribution, i.e., $D_1$ is constant for all $(x_i, y_i)$, $1 \leqslant i \leqslant m$ . At each subsequent iteration, we pass the distribution $D_t$ and the training examples $S$ to a weak learner that computes a weak hypothesis $h_t : X \times Y \to \mathbb{R}$. A positive (negative) sign of the weak hypothesis $h_t(x, y)$ corresponds to the prediction that AS $x$ is (is not) a member of class $y$. The value $|h_t(x, y)|$ of the weak hypothesis reflects the confidence level of the prediction. Then, we update the distribution $D_t$ so that the $x$, $y$ pairs that were erroneously predicted, i.e., the signs of $h(x, y)$ and $P(x, y)$ differ, receive a exponentially higher weight. By assigning higher weight to incorrect predictions, we force the algorithm to focus on these difficult examples in the next round. The final classifier $f$ is the sum of votes of the weak hypotheses in all rounds $h_t$, $t = 1 \ldots T$.

A weak hypothesis is equivalent to a one-level decision tree that checks a single AS attribute. For the description AS attribute, a weak hypothesis searches for the presence of a term or a sequence of terms in a given record, and if a match occurs, it outputs a confidence value for each of the classes. For example, upon finding the term "university" in the record "Seoul National University of Education" the weak hypothesis will likely output a high positive confidence value for the University AS class and a negative confidence value for the other AS classes. For scalar attributes, a weak hypothesis asks if a given attribute value is above or below a certain threshold. Depending on the outcome, the hypothesis outputs different confidence values.

The weak learner builds a weak hypothesis by exhaustively evaluating the attributes in the given weighted training examples. For a text attribute, it builds a candidate weak hypothesis by evaluating all possible terms and sequences of terms. For each term or sequence of terms, it calculates the appropriate confidence values by minimizing the Hamming loss, which is the fraction of examples $x$ and classes $y$, for which the sign of the final classifier $f(x, y)$ differs from $P(x, y)$. Similarly, for each scalar attribute, the weak learner builds a candidate weak hypothesis by exhaustively searching the threshold and confidence values minimizing the Hamming loss. On its output, the learner returns the weak hypothesis attaining the minimum Hamming loss.[4]

To realize our classification algorithm, we use BoosTexter [13], a publicly available implementation of AdaBoost. In Table 1 we depict the weak hypotheses that our algorithm discovered during its first six iterations. For each weak hypothesis, we illustrate the selected AS attribute, the term or threshold that is looked for, and the computed confidence values. The first weak hypothesis deals with the space attribute. If the IP address space advertised by an AS is less than 8.5 (equivalent) /24 prefixes then, the hypothesis assigns a positive confidence value to the Customer AS class and negative confidence values to the other classes. If the value of the space attribute is above 8.5, the hypothesis assigns negative or very close to zero confidence values to all the classes, which means that in this case it cannot make a confident positive prediction. The second weak hypothesis checks the description AS attribute for the presence of term "network inform".[5] If it finds one, it assigns a high positive confidence to the NIC AS class and negative confidence values to the other classes. Note, that in some cases the weak hypothesis assigns zero confidence values, meaning that it abstains from making any prediction.

We experimentally fix the number of rounds $T$ to 28, since subsequent iterations lead to overfitting. Overfitting is a common problem in machine learning. It is a consequence of too extensive training of an algorithm on one dataset. The undesired effect is that the algorithm is memorizing the training examples instead of extracting concepts from them. Fortunately, we can easily detect if the algorithm tends to overfit by examining the produced classification rules.

Having the number of iterations fixed, we finally apply our classifier to the set of 19,537 ASes and calculate the ranking of the AS classes for them. For each AS we consider the highest ranked class as its predicted class. If the class with the

---

[4] See [13] for the analytic expression of the Hamming loss.
[5] Recall that words have been stemmed.

**Table 1.** The list of weak hypotheses computed by AdaBoost.MH in the first six iterations. The first column is the iteration number; the second is the AS attribute that the weak hypothesis is checking; the third is the term or threshold of the weak hypothesis; and the remaining columns depict the computed negative or positive confidence values for each of the AS classes.

| Round | Attribute | Term/Threshold | L.ISP | S.ISP | Cusmr | Uni | IXP | NIC |
|-------|-----------|----------------|-------|-------|-------|-----|-----|-----|
| 1 | space | < 8.5 | | | | | | |
|   |       | > 8.5 | | | | | | |
| 2 | description | "network inform" | | | | | | |
| 3 | customer | < 1.5 | | | | | | |
|   |          | > 1.5 | | | | | | |
| 4 | description | "exchang" | | | | | | |
| 5 | description | "univers" | | | | | | |
| 6 | customer | < 97 | | | | | | |
|   |          | > 97 | | | | | | |

highest rank value for an AS has the confidence value less than or equal to zero, then the classifier abstains from making a prediction since the given information is not sufficient to produce a reliable assignment. Overall, the classifier abstains from making a prediction for 923 ASes, which accounts for 4.7% of the total number of ASes in our dataset. In Table 2 we show the per category classification statistics. Among the classified ASes, 63.0% are Customers, 30.1% are small ISPs, 4.7% are Universities, 1.8% are NICs, 0.2% are ISPs, and 0.2% are large IXPs.

**Table 2.** Numbers of ASes in each AS class.

|      | Large ISPs | Small ISPs | Customer ASes | Universities | IXPs | NICs |
|------|-----------|-----------|---------------|--------------|------|------|
| ASes | 44 | 5,599 | 11,729 | 877 | 33 | 332 |
| %    | 0.2 | 30.1 | 63.0 | 4.7 | 0.2 | 1.8 |

## 6 Validation

To validate our results, we employ the standard machine learning methodology called *cross-validation*. Cross-validation is the process of splitting the training examples into two subsets. One then uses the first subset to train a new classifier and the second subset to validate the results of this new classifier.

From our main set of 1200 training examples, we randomly extract 1100 ASes and use varying size subsets of these 1100 ASes to train new classifiers. We validate the predictions of the new classifiers against the remaining 100 examples. We repeat the random selection process 400 times and for each iteration we compute the following evaluation metrics: 1) *accuracy*, which we define as the percentage of ASes for which the AS class with the highest rank value is their correct class. The disadvantage of this metric is that it checks only the top of the ranking, ignoring the remaining positions. To address this problem, we use 2) *coverage*, which we define as the average position number of the correct class of an AS. For each AS, we number AS class positions incrementally starting from zero for the class with the highest positive confidence value. Thus, if all the predictions are correct the coverage is zero.

In Figures 2(a) and 2(b) we plot the average accuracy and coverage versus the size of the training set. As the size of the training set increases, the accuracy increases and the coverage decreases. For $|S| = 1100$ the accuracy reaches 0.781, e.g., 78.1%, and the coverage 0.251. The increasing accuracy trend suggests that for the training size of 1200 that we use in our final classification, the expected accuracy must be even higher. The low value of the coverage indicates that when the correct class is not of the top rank value, it is close to the top. More specifically, we find that for 97.7% of the predictions the correct class is in the top two positions of the ranking.

We next analyze the per class percentage of correct predictions. We find that for $|S| = 1100$ the percentage of correct predictions is on average: 100% for large ISPs, 100% for NICs, 100% for IXPs, 92.8% for Universities, 79.2% for Customer ASes and 72.1% for small ISPs. The actual distribution of ASes among the classes in our training set is: 684 Customer ASes, 401 small ISPs, 66 Universities, 36 NICs, 11 IXPs, and 2 large ISPs. The lower accuracy for customer ASes and small ISPs illustrates that these are the hardest classes to identify. We explain this effect by similarities between the characteristics of these two AS classes: 1) more than a half of the small ISPs appear to have the AS degree of 1 or 2, which is also typical for customer ASes; 2) some customer ASes, especially web hosting companies, advertise large numbers of different IP prefixes or large chunks of address space, which is also typical for ISPs.

In summary, we find that in the examined examples our classifier almost perfectly identifies large ISPs, NICs, IXPs and universities, while it also produces accurate predictions for customer ASes and small ISPs, which are the hardest to classify.

## 7 Conclusion

In this work, we establish a classification of ASes required for expanding our understanding of the Internet infrastructure and for creating realistic models of its topology and evolution. We develop a novel classification methodology that we apply to an exhaustive set of AS data to obtain the first statistics on the different AS classes in the Internet. We validate our results and demonstrate that our classifier achieves accuracy of 78.1% in the examined data. To promote further analysis
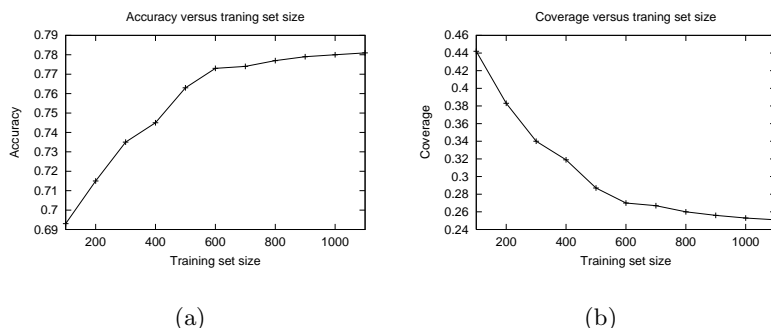
**Fig. 2.** Accuracy and coverage of computed predictions versus the size of the training set.

and to inspire development of better topology models, we release to the community our classification dataset along with the AS class predictions [3]. To the best of our knowledge, our dataset is the most comprehensive collection of AS macroscopic characteristics. In addition to AS topology and taxonomy information, it includes organization description records, AS business relationship information, and information on advertised IP prefixes and space.

# References

1. C. Alaettinoglu, C. Villamizar, E. Gerich, D. Kessens, D. Meyer, T. Bates, D. Karrenberg, and M. Terpstra. *Routing Policy Specification Language (RPSL)*. IETF, RFC 2622, 1999.
2. Tony Bates, Philip Smith, and Geoff Huston. The CIDR report. `http://bgp.potaroo.net/cidr/`.
3. Xenofontas Dimitropoulos. Revealing the Autonomous System Taxonomy: The Machine Learning Approach. Data Page. `http://www.ece.gatech.edu/research/labs/MANIACS/as_taxonomy/`.
4. Xenofontas Dimitropoulos, Dmitri Krioukov, Marina Fomenkova, Bradley Huffaker, kc claffy, and George Riley. AS relationships: Inference and validation. In *under submission*, 2005.
5. Xenofontas Dimitropoulos, Dmitri Krioukov, Bradley Huffaker, kc claffy, and George Riley. Inferring AS relationships: Dead end or lively beginning? In *Proceedings of 4th Workshop on Efficient and Experimental Algorithms (WEA' 05)*, May 2005.
6. Y. Freund and R. E. Schapire. A decision theoretic generalization of on-line learning and an application to boosting. In *Second European Conference on Computational Learning Theory (EuroCOLT-95)*, pages pages 23–37, 1995.
7. Z. Ge, D. Figueiredo, S. Jaiwal, , and L. Gao. On the hierarchical structure of the logical Internet graph. In *SPIE ITCOM*, 2001.
8. R. Govindan and A. Reddy. An analysis of Internet inter-domain topology and route stability. In *IEEE INFOCOM*, 1997.
9. David Meyer. University of Oregon Route Views Project. `http://www.routeviews.org/`.
10. Ted Pedersen. WordNet stop list. `http://www.d.umn.edu/~tpederse/Group01/wordnet.html`.
11. Martin Porter. An algorithm for suffix stripping. *Program*, 14:130–137, 1980.
12. Internet Routing Registries. `http://www.irr.net/`.
13. R. E. Schapire and Y. Singer. BoosTexter: A boosting-based system for text categorization. *Machine Learning*, 39:135–168, 2000.
14. L. Subramanian, S. Agarwal, J. Rexford, and R. H. Katz. Characterizing the Internet hierarchy from multiple vantage points. In *IEEE INFOCOM*, 2002.