

Modification in ‘KNN’ Clustering Algorithm for Distributed Data

Eena Gilhotra¹, Priyanka Trikha²

¹Department of Computer Science
Sri Ganganagar Engineering College
Sri Ganganagar (Raj.), India

²SBCET –Jaipur, India

¹eena.er@gmail.com

Abstract - Clustering has become an increasingly important task in modern application domains such as marketing and purchasing assistance, multimedia, molecular biology etc. The goal of clustering is to decompose or partition a data set into groups such that both the intra-group similarity and the inter-group dissimilarity are maximized. In many applications, the size of the data that needs to be clustered is much more than what can be processed at a single site. Further, the data to be clustered could be inherently distributed. The increasing demand to scale up to these massive data sets which are inherently distributed over networks with limited bandwidth and computational

resources has led to methods for parallel and distributed data clustering. In this thesis, we present a cohesive framework for cluster identification and outlier detection for distributed data. The core idea is to generate independent local models and combine the local models at a central server to obtain global clusters. A feedback loop is then provided from the central site to the local sites to complete and refine the global clusters obtained. Our experimental results show the efficiency and accuracy of our approach.

Keywords- Cluster, Data mining, Data warehousing

I. INTRODUCTION

Knowledge Discovery in Databases (KDD) is the process of searching large volumes of data for the non trivial extraction of implicit, novel, and potentially useful information. Traditional KDD applications require complete access to the data which is going to be analyzed. Nowadays, a huge amount of heterogeneous, complex data resides on different computers which are connected to each other via local or wide area networks (LANs or WANs). One of the most common approaches for business applications to perform data mining on such massive datasets is to centralize distributed data in a data warehouse on which the data mining techniques are applied. Data warehousing is a widely used technology which integrates data from multiple data sources into a single repository in order to efficiently execute complex analysis queries. However, despite its commercial success, this approach may be impractical or even impossible for certain business settings, for instance:

- When a huge amount of data is frequently produced at different sites and the cost of its centralization cannot scale in terms of communication, storage and computation. For example, call data records of telephonic conversations.
- Whenever data owners cannot or do not want to release information. This could be in order to

maintain privacy or because disclosing such information may result in a competitive disadvantage or a considerable loss in commercial value. For example, data mining across banks by the Reserve Bank of India.

Such scenarios call for distributed data mining. Distributed data mining deals with pattern extraction problem. One of the most studied data mining techniques is clustering. The goal of this technique is to decompose or partition a data set into groups such that both intra-group similarity and inter-group dissimilarity are maximized. In particular, clustering is fundamental in knowledge acquisition. It is applied in various fields including data mining, statistical data analysis, compression and vector quantization.

In this paper, we present a cohesive framework for cluster identification and outlier detection for distributed data. The data is either distributed originally because of its production at different locations or is distributed in order to gain a computational speed up. We use a parameter free clustering algorithm to cluster the data at local sites. These clusters are actually the partial results which are communicated in the form of local models to the central server. The central server aggregates these partial results to give a global solution. A feedback loop is then provided to purify and enhance this global solution.

Below we give a brief statement of the problem addressed in the paper:

Given a dataset D of n points distributed across s sites and a central server, find the clusters of the dataset D by communicating minimum information to the central server such that the accuracy of the obtained clusters is comparable to the results obtained using a centralized clustering approach...

II. RELATED WORK

A. Clustering

Clustering partitions a dataset into highly dissimilar groups of similar points. The definition of clusters and outliers depends very much on the domain of the dataset. For the sake of clarity, we provide here general definitions quoted in the literature. A cluster is a set of similar points that are highly dissimilar with other points in the dataset. An outlier or a noise point is an observation which appears to be inconsistent with the remainder of the data. Next, we discuss the various clustering algorithms briefly.

Partitioning techniques K-means, PAM, CLARA and CLARANS are good examples for clustering based on partitioning techniques. K-means is in fact a very popular clustering algorithm. These clustering algorithms work under the assumption that there is no noise in the dataset, the clusters are spherical-shaped and the points in the cluster follow uniform distribution. They perform well on datasets containing spherical clusters but also do include noise in the cluster results. They cannot identify clusters of irregular shapes and sizes. CURE is another partitioning based clustering technique developed to reduce the effect of noise and identify elliptical shaped clusters. It uses multiple representative points for each cluster to reduce the effect of noise but cannot capture clusters of different densities.

1) Hierarchical techniques

- Complete link or MAX or CLIQUE
- Average

Density-based techniques: Examples of this technique are DBSCAN and OPTICS. They consider the density around each point to identify the cluster boundaries and the core cluster points. The close cluster points in a single neighborhood are then merged. OPTICS does not generate a clustering solution; instead it generates an augmented ordering of the points. Given a reachability distance, it generates reachability plots that capture the local density settings around each point. It can be effective for noisy datasets having arbitrary shaped clusters of different sizes and densities. The clustering results for both these algorithms depend largely on the provided parameters. Finding these parameters for both these algorithms is a challenge for the user.

III. A FRAMEWORK FOR DISTRIBUTED CLUSTERING

The proposed framework uses object distributed data and is a centralized ensemble based method. This is a cohesive framework for cluster identification and outlier detection for object distributed data. From now onwards by distributed data, we refer to object distributed/homogeneous data. The core idea of the framework is to generate independent local models and combine the local models at a central server to obtain global clusters. A feedback loop is then provided from the central site to the local sites to complete and refine the obtained global clusters. We illustrate the this framework in Figure 3.1. The steps involved in the process of distributed clustering for the framework are:

A. DD: Data Distribution

As stated earlier, the need for distributed clustering might occur not only when the data is inherently distributed but also when the data cannot be processed using a single processor and thus to gain a computational speed up, it is distributed across multiple processors. The step DD in this process involves the distribution of data across the local sites. There could be different data partitioning strategies which could be employed for distributed data mining. The one we use is cover based partitioning, i.e. In our approach the partitions are overlapping. Here two partitions D_i and D_j are overlapping if $|D_i \cap D_j| > 0$.

B. Local Clustering Algorithm

The clustering algorithm that we use for clustering the data at local sites is the Stability based RECORD Algorithm (SRA) with minor modifications.

The SRA algorithm uses the notion of k-Reverse Nearest Neighbor (kRNN) and Strongly Connected Components (SCC) to:

- Derive clearly distinguishable clusters, and
- To identify and remove outliers during the process of clustering.

SRA presents a cohesive framework for both cluster analysis and outlier detection that improves the quality of the clustering solutions.

1) Definitions

kRNN (p_j): The set of points which consider p_j as one of their k th nearest neighbors.

$kRNN(p_j) = \cup_{i=1}^k p_{rnn_i}(p_j)$

Strongly Connected Components (SCC): The SCC of a digraph partitions the vertices's into subsets such that all the points of each subset are mutually reachable.

Core point: In a cluster, a core point is a point that lies amidst the dense set of points of the cluster. The concept of kRNN is used to capture the nature of this neighborhood of the point. For a given value of k , a

point having $|\text{kRNN}| \geq k$ is flagged as a core point.

Boundary point and Outlier: A boundary point in a cluster is a point which lies in the transition region of dense points to noise points in the cluster. An outlier is a point far from most others in a set of data. Based on large number of experiments done over various datasets, it was found in [11] that an outlier has $|\text{kRNNs}| < k$ and most of its kRNNs are outliers. Also for a boundary point, its $|\text{kRNNs}| < k$ and most of its kRNNs are core points.

C. Basic RECORD Algorithm

Stability based RECORD algorithm (SRA) is a hierarchical clustering algorithm which uses the basic RECORD algorithm presented in [11]. Therefore we first describe the basic RECORD algorithm and then move on to the SRA in the later parts of the section. The major steps involved in the basic RECORD algorithm are as follows:

Step 1-kRNN Computation: Generate the distance matrix based on the distance function $d(i, j)$. Here distance function is a measure to define the distance between two data points. Calculate the kRNN sets of each point as follows:

- For each point p , identify the k -nearest points $\text{knn}(p)$. For every point $q \in \text{knn}(p)$, increase the count of its reverse nearest neighbors by 1.
- Add the directed edge ($q \in \text{knn}(p), p$) to the kRNN graph.

Step 2-Outlier detection: For each data point, the numbers of points in kRNN (p) are checked. For any data point p if $|\text{kRNN}| < k$, it is flagged as an outlier. The point p and all out-going edges from p are removed from the kRNN graph. After removing all outlier points from the graph, the modified graph is represented as $\text{kRNN} > k$. Also we have the graph $\text{kRNN} < k$ which is given by $\text{kRNN} - \text{kRNN} > k$. This $\text{kRNN} > k$ includes only outlier points and their corresponding out-going edges and is used at later stages to get a stable clustering solution.

Step 3-Cluster Identification: Eliminating outliers from kRNN results in $\text{kRNN} > k$. The clusters are now computed based on this $\text{kRNN} > k$. Each SCC in the $\text{kRNN} > k$ graph becomes a cluster as it follows two rules:

- Each member of the SCC is accepted as k -nearest by at least k points and
- Every possible pair of members in the SCC is mutually reachable depicting the cohesiveness of the graph.

SCCs are also computed on $\text{kRNN} < k$ (sub-graph kRNN containing only outliers). The SCCs obtained from this graph are clusters of noise points and are very sparse. The SCCs so obtained are used to identify stable clustering solution (Stability based RECORD). An efficient approach to compute these

SCCs incrementally is provided in [11].

Step 4-Local outlier incorporation: Since even the boundary points could have their $|\text{kRNN}| < k$, the outlier detection technique discussed above could be very stringent. Because of this, some of the boundary points of the clusters are also identified as outliers. Due to this, the clusters generated are highly dense and incomplete. To avoid such eliminations, after the generation of SCCs, an effort is made to draw each of the outlier points to its nearest cluster. Nearest cluster is the closest majority cluster among kRNN points of the outlier under study. At the end of this step, an outlier point such that majority of its neighbors lie in one cluster, gets flagged as a cluster point of that cluster.

D. Stability based RECORD algorithm

Having discussed the basic RECORD algorithm, we now give a brief explanation of the Stability-based RECORD Algorithm (SRA) algorithm. The SRA utilizes the basic RECORD algorithm, and generates the most stable clustering result as its output. Starting from $k = 1$, for each value of k the kRNN is computed. The kRNN is split into two subgraphs $\text{kRNN} < k$ and $\text{kRNN} > k$. SCCs are computed for both the subgraphs. As k increases, the number of SCCs in $\text{kRNN} < k$ and $\text{kRNN} > k$ are checked. After a certain value of k these numbers do not change. This would imply that no new clusters are formed and no variation in noise behavior is detected. So in [11], the number of SCCs in both the subgraphs is observed for some l consecutive values of k and if they remain constant then the clusters obtained are considered to be stable. The value of l is determined experimentally. For the distributed scenario, the stability constraint has to be less stringent because at any given site only a subset of data is clustered. Because of this, the stability of the clustering algorithm is expected to go down.

For our experiments, we chose l to be 1 and terminate the local clustering algorithm when the number of SCCs remains constant for both the subgraphs for at least two successive values of k .

The first step of the SRA requires the computation of the complete distance matrix to compute the kNN of a data point. This becomes a bottleneck in the performance of the SRA algorithm. In the next subsection, we present an efficient approach to compute these k nearest neighbors of a data point.

E. Efficient Computation of kNN

For computing the kNN of a point one needs to compute the distance/similarity matrix. A similarity matrix is a $n \times n$ matrix of distance measures which expresses the similarity between any two data points. The $O(n^2)$ complexity becomes a bottleneck in the performance of the SRA algorithm. The experiments

carried out in [11] show that stable clustering solution is found by exploring within 100 neighbors of every data point. Hence we propose an efficient method to compute the kNN of every data point where k has some known value which is much less than n. The basic idea here is to convert the kNN search of a data point from the entire search space to a search in some neighborhood of the point using the property explained in Theorem 1.

Consider a data point A and its kNN points ordered on distance: N_1, N_2, \dots, N_k . Let N_k be the kth nearest neighbor of A and $d(A, N_k) = x$, where $d(A, N_k)$ is the distance between A and N_k . We define a neighborhood value of a data point with respect to a value k, as the radius of the circle that encompasses its k nearest neighbors. Thus x is the neighborhood value of A.

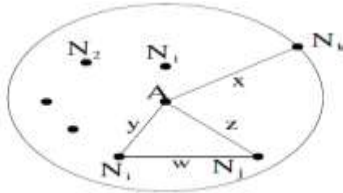


Figure 1: Efficient Computation of kNN

Theorem 1 Consider any one kNN, N_i of A and let $d(A, N_i) = y$. Then the circle with center as N_i and radius $x+y$, encompasses all the points $N_1, N_2, \dots, N_{i-1}, N_{i+1}, \dots, N_k$ and A (refer Figure 4.1).

Proof: Since radius of the circle is $x + y$ and point A is at distance y from N_i , A lies in the circle.

Now consider any kNN of A other than N_i say N_j . Let $d(N_i, N_j) = w$. Then in the $\Delta(N_i, N_j, A)$, by triangle inequality $w < y + z$. Since N_k is the kth nearest neighbor, $z < x$. Therefore, $w < x + y$. Thus N_j lies inside the circle. Hence proved.

Based on theorem 1, the working of the kNN computation algorithm is given below.

Initialize P S and HP S to be empty sets. Here P S is the list of points which are to be processed and HP S is the set of points which are already processed.

Let d_i be any data point whose neighbors are not yet found. Its kNN set $\{N_1, \dots, N_k\}$ is computed by determining the distance of d_i with every data point. Let $NBHD_i$ be the neighborhood value of d_i . The kNN set is then added to a processing list P S which is a set of ordered pairs $(N_i, NBHD_{ii} + d(N_i, d_i))$ where $d(N_i, d_i)$ denotes the distance between N_i and d_i .

For each pair $(T, V) \in P S$, determine all points $\{P_1, \dots, P_l\}$ that lie inside the circle that is drawn with T as center and V as radius. Note that $l \geq k$ by theorem 1. The distance of T from each $P_i, \forall i = 1, \dots, l$ is computed and the kNN set of T, $\{O_1, \dots, O_k\}$ is determined as the k closest points out of P_1, P_2, \dots, P_l .

Let NBHT be the neighborhood value of T. Then the set of pairs $(O_i, NBHT + d(T, O_i))$ for all $i = 1, \dots, k$, such that O_i is not present in any of the ordered pairs in P S and in HP S, is added to P S. The pair (T, V) is removed from P S and T is added to HP S.

Step 3 and 4 are repeated until P S becomes empty.

Repeat step 2, 3, 4 and 5 till the k neighbors for all the data points are found.

The above approach helps us in overcoming the earlier $O(n^2 \log n)$ time complexity of computing the k nearest neighbors for each point. Next we present an analysis of the time complexity of the above approach. We also provide experimental results in which support the obtained worst case bounds on the time spent.

The total complexity due to all points is:

$$O(s*n*\log n + (n-s)*[\log n + m\log m])$$

The overall cost for computing the k nearest neighbors can be given as

$O(c*n*\log n + k^2 * n*\log(n/c))$. The experimental results presented in Table 1 show that as the difference between the k value and the size of the dataset i.e. n increases, the efficiency gained by employing the proposed approach also increases.

In the next section, we discuss the information that is communicated to the central server, using which the global clusters can be acquired.

F. Acquiring the Local Model

After having clustered the data locally at each local site, we would like to communicate to the central server adequate amount of information that will describe the local clustering results. These local models have to be such that they give an accurate description of the local clustering as far as possible.

The information that can be communicated from the local clusters in the form of local model is:

- Cluster statistics such as density information and number of data points in the cluster.
- Representative points from the cluster. These could be chosen from the core/dense regions of the cluster and also from the boundary regions of the cluster.

G. Acquiring the Global Model

Once the local model from each site is communicated to the central server, a global merging algorithm has to be employed on the local models to acquire the global model. Let there be any two representative SCCs of clusters C_{ij} (jth cluster from i th site) and C_{mn} (mth cluster from nth site), which are candidates to be merged. Consider the points $p, t \in C_{ij}$ and $q, r \in C_{mn}$. The clusters C_{ij} and C_{mn} would merge if for any value k, $p \in kRNN(q)$ and $r \in kRNN(t)$ as in Figure 5.1. Here p, q, r, t are the representatives that have been communicated to the

central server in the respective local models.

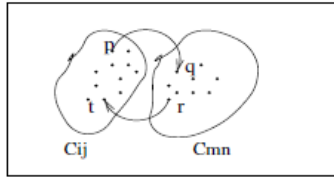


Figure 2: Merging of Clusters

As can be seen, the global merging algorithm is required to cluster the data that is communicated in the form of representatives of the local clusters. At the central server, we again apply the Stability based RECORD algorithm to cluster these local representatives effectively. As stated earlier, the termination condition for the SRA is as follows: the number of SCCs in both the subgraphs ($kRNG < k$ for the outliers and $kRNG \geq k$ for the core points) are observed for some l consecutive values of k and if they remain constant then the clusters obtained are considered to be stable. The global clusters obtained at this stage might not be complete

H. Updation of the global and the local models

After finishing the global clustering, we send the complete global model to all the client sites. Using this information, the client sites can assign each of the objects a label that corresponds to the global cluster to which it belongs.

IV. EXPERIMENTS AND RESULTS

We conducted the experiments on synthetic datasets and real life datasets to evaluate the efficacy of our approach.

A. Results for efficient computation of nearest neighbors

In Table 1 we present the results obtained using the approach proposed by us for computing the k nearest neighbors of a data point efficiently. Since the results in [11] show that computing just 100 nearest neighbors is sufficient to get the clusters present in a dataset we take the k value as 100 and show the results for it.

TABLE 1:- SPEED UP GAINED IN COMPUTING THE K NEAREST NEIGHBORS

Dataset size	Ratio T1/T2
500	1.089
1000	2.280
3000	3.373
5000	5.735
10000	9.654
100000	34.173

There are two columns in the table1. Column 1 shows the size of the dataset and column 2 shows the ratio $T1 /T2$. $T2$ is the time taken by the approach proposed by us in the Section 2.3. $T 1$ is the time taken

by the naive approach where to get the k nearest neighbors the entire similarity matrix is computed. As is demonstrated by the results, with the increase in the difference between the k value i.e. 100 and the size of dataset, the efficiency gained by the proposed technique also increases.

B. Efficacy of our approach

In this section we discuss the various efficiency and accuracy results obtained for our framework.

1)Results on synthetic datasets

Syndeca tool set was used to generate synthetic datasets with clusters for experimentation. In order to compute the total time taken by our approach, we carried out the local clustering for all sites and denoted the time taken by these as a set $\{TLCMi\}$ for $i = 1, \dots, s$. Here the time taken by the individual sites for this updation phase was denoted by a set $\{TUM\}$ for $i = 1, \dots, s$. The overall runtime T for the approach was calculated using the formula given below:

$$T = T_g + \max \{TLCi\}, +\max \{TUMi\}$$

The variable i spans from 0 to s taking into considerations all the local sites. Here s is the total number of sites, T_g is the time taken by the global merging algorithm GM, $T LCMi$ is the time taken for LCM on the i th site and $T UMi$ is the time taken for the updation phase UM on the i th site.

2)Experiments with real life datasets:

In this section we provide the results for the proposed framework on three real life datasets. We also provide the performance of the DBDC approach on the same datasets and compare the accuracy results obtained by both. To get the performance results for DBDC, we implemented it in C++. All the experiments were carried out on an AMD dual 64 bit machine with 2 GB of RAM. The evaluation metric that we use for the real life datasets is cluster homogeneity. We define accuracy as:

$$\text{accuracy} = (\sum_{ki=0} a_i)/n$$

Here n is the total number of points to be clustered and a_i is the number of points that belong to the dominance class of the k th cluster.

3)Iris Dataset:

The Iris flower data set or Fisher's Iris data set is a multivariate data set introduced by Sir Ronald Aylmer Fisher (1936). The dataset consists of 3 classes, 50 instances each, where each class refers to a type of iris plant namely Iris Setosa, Iris Versicolor, and Iris Verginica. The first class is linearly separable from others while the latter two are not linearly separable. There are four continuous attributes and a target attribute which determines the class of the tuple. The attribute measurement consists of the sepal and petal lengths and

widths in cms.

TABLE 2:- RESULTS ON IRIS DATASET

	Kmeans	DBSCAN	DBDC	SRA	Our approach
No. of sites	1	1	2	1	2
Dataset size	150	150	75	150	75
Parameters	K=3	Mpts=3 Eps=0.5	LMpts=3 LEps=1 GMPT S=2 GEps=2	None	None
No.of clusters	3	10	5	23	6
Accuracy	0.667	0.5	0.634	0.663	UM:0.812

4) Comparison of the existing work with our approach

As discussed in the earlier sections of the thesis, a distributed clustering algorithm is expected to exhibit several properties.

TABLE 3: COMPARISON OF VARIOUS DISTRIBUTED CLUSTERING APPROACHES

Property	Our approach	SRA	DBDC	SDBC	DMBC
Efficiency	Yes	Yes	yes	yes	Yes
Clusters of mixed densities	Yes	Yes	No	no	-
Clusters of varied shapes	Yes	Yes	yes	yes	-
Privacy	No	No	No	no	Yes
Needs parameters for local clustering	No	No	yes	yes	Yes
Amount of info communicated for obtaining global model is tunable	Yes	Yes	No	yes	No
Detects clusters lost due to Distribution	Yes	Yes	No	no	No
No of rounds of communication	2	2	1	1	1

V. CONCLUSIONS

In this thesis, we present a cohesive framework for the identification of clusters and outliers in the distributed environment. The presented framework uses homogeneously distributed data and is a centralized ensemble based method. It requires two rounds of communication between the central server and the local sites. In the first round a global model is obtained from the local models generated as a result of the clustering at local sites. The second round is employed to complete and purify the global clusters obtained in the first round. Our approach uses the

parameter free clustering algorithm SRA for the clustering of data at local sites. It also proposes an efficient way to compute the kNN. Unlike previous centralized ensemble based approaches for homogeneous distributed data, using CIODD, we are able to detect clusters of mixed densities and varied shapes placed in close vicinity of each other. Earlier approaches also fail to handle the local outliers, where as our approach is able to do so because of the presence of the grid based feedback loop. Moreover, the increase in accuracy obtained due to the introduction of the feedback loop is much more than the increase in the overhead caused by its computation. We also show that without compromising much with the accuracy, the time taken by our approach is much less than the classical centralized clustering.

For example we were able to cluster a set of 20,000 points by distributing the data on six sites with an accuracy of 98.54% with the amount of information communicated being equal to 9% of the data. Also the time taken by the entire process was 75% less than what was taken by a centralized solution.

REFERENCES

- [1] K. Jain, M. N. Murthy, and P. J. Flynn., "Data clustering: A review".
- [2] F. E. Grubbs., "Procedures for detecting outlying observations in samples," in In Technometrics, pp. 2-21, 1969.
- [3] P. Berkhin., "Survey of clustering data mining techniques" in Tech. Report, Accrue Software.
- [4] D. Fasulo, "An analysis of recent work on clustering algorithms: a technical report".
- [5] S. Guha, R. Rastogi, and K. Shim., "Cure: An efficient clustering algorithm for large databases" in SIGMOD, 1998.
- [6] T. Zhang, R. Ramakrishnan, and M. Livny, "Birch: A efficient data clustering method for very large databases" in SIGMOD, 1996.
- [7] M. Ester, H. P. Kriegel, J. Sander, and X. Xu., "A density based algorithm for discovering clusters in large spatial databases with noise" 1996.
- [8] G. Karypis, E. Han and V. Kumar, "Chameleon: Hierarchical clustering using dynamic modeling" Computer, vol. 32, no. 8, pp. 68-75, 1999.
- [9] S. Bandyopadhyay, C. Gianella, U. Maulik, H. Kargupta, K. Liu, and S. Datta, "Clustering distributed data streams in peer-to-peer environments", Information Science Journal, 2004.
- [10] M. Klusch, S. Lodi, and G. L. Moro, "Distributed clustering based on sampling local density estimates", in Proceedings of International Joint Conference on Artificial Intelligence (IJCAI 2003), (Mexico), pp. 485-490, August 2003.
- [11] S. Vadapalli, S. R. Valluri, and K. Karlapalem, "simple yet effective data clustering algorithm" ICDM, 2006.
- [12] E. Januzaj, H. P. Kriegel, and M. Pfeifle, "Towards effective and efficient distributed clustering" Workshop on Clustering Large Data Sets (ICDM2003), 2003.
- [13] E. Januzaj, H. P. Kriegel, and M. Pfeifle, "Scalable density-based distributed clustering" Proc. 8th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD), 2004.
- [14] D. K. Tasoulis and M. N. Vrahatis, "Unsupervised distributed clustering," in In Proceedings of the IASTED International Conference on Parallel and Distributed Computing and Networks. Innsbruck, Austria, 2004.