USER MANUAL OF HAZY-CLASSIFY

M. Levent Koc University of Wisconsin-Madison koc@cs.wisc.edu Christopher Ré University of Wisconsin-Madison koc@cs.wisc.edu

February 5, 2011

Contents

1 Overview		2					
2	Installation						
	1 Installing Prerequisites	2					
	2 Configuring HAZY-CLASSIFY	4					
3	sing Hazy-Classify	4					
3	sing Hazy-Classify 1 Loading demo examples	4 5					
	 Loading demo examples Creating Classification 	5					
	1 Loading demo examples	5					

1 Overview

Classification is a ubiquitous and well-studied statistical task. Many applications such as Facebook and Twitter use classification. There are great tools to classify the data. Some examples are: SVM Light, Svm Sgd and Libsvm. However, the downside of these tools is that they are not able to incrementally maintain the output of classification results. Also, they are not well integrated with open source databases such as PostgreSQL.

HAZY-CLASSIFY is a system that implements classification inside an RDBMS for dynamic data as part of the ongoing Hazy Project. It supports classification techniques such as SVMs. One of the main contributions of Hazy-Classify is that it can efficiently maintain the output of classification results by using novel techniques. It achieves scalability and orders of magnitude speedup compared to prior implementations on many datasets including Forest Cover, DBLife and Citeseer. Hazy-Classify is implemented in C++ on top of PostgreSQL.

This manual contains an installation and usage guide for HAZY-CLASSIFY. For more information, please visit our project website.

2 Installation

HAZY-CLASSIFY is implemented in C++ on top of PostgreSQL. Therefore, it requires gcc/g++ compiler and PostgreSQL (8.4 or higher). In addition, it requires ocaml in order to parse Hazy queries. Figure 1 lists the URLs where you can get them.

gcc/g++	http://gcc.gnu.org/releases.html			
PostgreSQL	http://www.postgresql.org/download/			
ocaml	http://caml.inria.fr/download.en.html			
HAZY-CLASSIFY	http://www.cs.wisc.edu/hazy/hazy-classify/download.php			

Figure 1: Download links of the softwares

2.1 Installing Prerequisites

Install gcc/g++ If you don't have gcc/g++ compiler, you can download and install it from this link.

Install and Configure PostgreSQL You can download PostgreSQL 8.4(or higher) from here. Some systems can have PostgreSQL by default, but since HAZY-CLASSIFY uses libpq library of PostgreSQL, you need to download the source code and install it from scratch.

1. Let PG_DIST be the location where you unpacked the source distribution. Let PG_PATH be the location where you want to install PostgreSQL. Run the following commands:

```
cd PG_DIST
./configure --prefix=PG_PATH
gmake; gmake install
cd PG_PATH/bin
mkdir PG_PATH/data
```

initdb -D PG_PATH/data
postmaster -D PG_PATH/data &
createdb hazy_demo

initdb command initializes the database and postmaster runs PostgreSQL server. At the end, a database called hazy_demo is created.

2. You also need to add PG_PATH to the PATH variable of your system in order to run PostgreSQL commands without navigating to PG_PATH/bin directory.

PATH=\$PATH:PG_PATH/bin

3. You need to set LD_LIBRARY_PATH to include PG_PATH/lib directory in order to run Hazy.

```
LD_LIBRARY_PATH=$LD_LIBRARY_PATH:PG_PATH/lib
```

4. (Optional step) Finally, you can tune configuration parameters to obtain better performance results with HAZY-CLASSIFY. Assuming 4GB RAM (you can change parameters proportionally to RAM), the parameters in postgresql.conf (under PG_PATH/data) can be changed as follows:

shared buffers=896MB
work mem=256MB
maintenance work mem=512MB
effective cache size=1792MB
checkpoint segments=65536

The system might not allow shared buffer parameter to be increased since it uses more System V shared memory. In order to change this, you can edit /etc/sysctl.conf file and add these two lines, then run the command "sysctl -p". Note that, it requires root access.

```
kernel.shmall=671088640
kernel.shmmax=671088640
```

Install ocaml You can download ocaml from this link.

1. Let OCAML_DIST be the location where you unpacked the source distribution. Let OCAML_PATH be the location where you want to install ocaml. Run the following commands:

```
./configure -prefix OCAML_PATH
make world
make opt
make install
```

2. You also need to add OCAML_PATH to the PATH variable of your system in order to run ocaml commands without navigating to OCAML_PATH/bin directory.

PATH=\$PATH:OCAML_PATH/bin

2.2 Configuring Hazy-Classify

Download HAZY-CLASSIFY from the project website. Let the directory where you unpack is HAZY_HOME. After unpacking, you should see two directories: "doc" and "system_src". When you navigate to "system_src", you can see HAZY-CLASSIFY configuration file named "hazy.conf" (Figure 2).

- POSTGRES_PATH: PG_PATH specified during the installation of PostgreSQL.
- Database: Name of the database you want to use HAZY-CLASSIFY
- Storage_Manager: Type of the storage manager that HAZY-CLASSIFY uses. Options: Main_Memory_SM, Ondisk_SM, Hybrid_SM
- Use_Eps_Map: HAZY-CLASSIFY uses eps map when Storage_Manager is Hybrid_SM and this variable is set to 1.
- Hazy_Strategy: The strategy that is used by Hazy-CLASSIFY during incremental updates. Options: EAGER_HAZY, EAGER_NAIVE, LAZY_HAZY, LAZY_NAIVE
- Reservoir_Size: If you want HAZY-CLASSIFY to use reservoir during the training, you can set this variable to a value > 0 (e.g. 1000)
- Buffer_Size: Size of the buffer that is used by HAZY-CLASSIFY when Storage_Manager is Hybrid_SM. (Usually, 1% of the # of entities).

```
POSTGRES_PATH=/hazy/Programs/postgresql
Database=hazy_demo
Storage_Manager=Ondisk_SM
Use_Eps_Map=0
Hazy_Strategy=EAGER_HAZY
Reservoir_Size=0
Buffer_Size=0
```

Figure 2: Default Configuration File

After setting the parameters in the configuration file, you can install HAZY-CLASSIFY by running the following commands. The sample output of this installation is shown in Figure 3.

```
cd HAZY_HOME
cd system_src/Hazy
sh scripts/install.sh
```

Note that, you need to call install.sh script for each database that you want to use HAZY-CLASSIFY.

3 Using Hazy-Classify

In this section, we describe how to use HAZY-CLASSIFY. In order to explain how HAZY-CLASSIFY can be used, we will use the demo examples that you can find under HAZY_HOME/system_src/Hazy/demo. We first show how to load demo examples and then how to create, update and query classifications by using HAZY-CLASSIFY.

```
CREATE FUNCTION
NOTICE: function getentitylabelwithsgl(float8,pair[]) does not exist, skipping
CREATE FUNCTION
CREATE FUNCTION
CREATE FUNCTION
NOTICE: function getentitylabelwithsqlnoeps(pair[]) does not exist, skipping
CREATE FUNCTION
ERROR: type "sm_test_external_table" does not exist
CREATE FUNCTION
ERROR: relation "hazy_cls_view_catalog" already exists
make -C installation install
make[1]: Entering directory `/scratch/koc/hazy_svn/HAZY_VIEWS/trunk/hazy_system/Hazy_Syst
em/system_src/src/installation'
Trying to install in /scratch/koc/Programs/postgresql/lib
cp s_libs/functions/frontend_functions.so s_libs/functions/backend_functions.so s_libs/fu
nctions/common_functions.so s_libs/functions/inplace.so s_libs/triggers/update_trig.so s_
libs/triggers/trig_noop.so s_libs/triggers/trig_noop_server.so /scratch/koc/Programs/post
gresql/lib
make[1]: Leaving directory `/scratch/koc/hazy svn/HAZY_VIEWS/trunk/hazy_system/Hazy_Syste
m/system_src/src/installation'
[koc@francisco] (43)$
```

Figure 3: Terminal Output of Hazy installation

3.1 Loading demo examples

Demo examples contain 50 computer science publication names. Our purpose is to classify these papers as database and non-database papers and maintain them efficiently in the presence of updates. You can load these examples to the database by running the following commands:

```
cd HAZY_HOME/system_src/Hazy
psql hazy_demo < demo/demo_papers.sql</pre>
```

This will create two tables in hazy_demo database: demo_papers that contains 50 publication names and demo_training_examples that will be used to give training examples to HAZY-CLASSIFY. Sample terminal output of these commands is shown in Figure 4.

3.2 Creating Classification

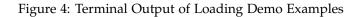
In order to use HAZY-CLASSIFY, first you need to run Hazy server. To start the Hazy server, run the following commands as shown in Figure 5.

```
cd HAZY_HOME/system_src/Hazy
./exec/hazy_main
```

After starting Hazy server, now we are ready to create the view. In order to create a view, you need to write a query in Hazy-CLASSIFY syntax. HAZY_HOME/system_src/Hazy/demo/create_demo_view_query contains create view query for the demo example.

INSERT 0	1
INSERT 0	1
INSERT 0	-
INSERT 0	
INSERT 0	-
INSERT 0	
LKOCGITA	ncisco] (45)\$

٠



```
[koc@francisco] (14)$ cd ~/hazy/system_src/Hazy
[koc@francisco] (14)$ ./exec/hazy_main &
[1] 3398
[koc@francisco] (15)$ server started
[koc@francisco] (15)$ []
```

Figure 5: Terminal Output of Starting Hazy Server

```
CONTEXT: SQL statement "DROP TABLE IF EXISTS entity_features_temp5;"
PL/pgSQL function "create_entity_features" line 2 at EXECUTE statement
NOTICE: table "hazy_entity_table5" does not exist, skipping
CONTEXT: SQL statement "DROP TABLE IF EXISTS hazy_entity_table5;"
PL/pgSQL function "create_entity_features" line 3 at EXECUTE statement
NOTICE: sequence "serial5" does not exist, skipping
CONTEXT: SQL statement "DROP SEQUENCE IF EXISTS serial5;"
PL/pgSQL function "create_entity_features" line 4 at EXECUTE statement
NOTICE: table "featureids5" does not exist, skipping
CONTEXT: SQL statement "DROP TABLE IF EXISTS featureIds5;"
PL/pgSQL function "create_entity_features" line 5 at EXECUTE statement
NOTICE: index "fids5" does not exist, skipping
CONTEXT: SQL statement "DROP INDEX IF EXISTS fids5;"
PL/pgSQL function "create_entity_features" line 6 at EXECUTE statement
NOTICE: table "training5" does not exist, skipping
CONTEXT: SQL statement "DROP TABLE IF EXISTS TRAINING5;"
PL/pgSQL function "prepare_train" line 2 at EXECUTE statement
/scratch/koc/hazy_svn/HAZY_VIEWS/trunk/hazy_system/Hazy_System/system_src/src/exec/../ins
tallation/functions/inplace5.c: In function 'update_hack5':
/scratch/koc/hazy_svn/HAZY_VIEWS/trunk/hazy_system/Hazy_System/system_src/src/exec/../ins
tallation/functions/inplace5.c:69: warning: assignment makes pointer from integer without
a cast
/scratch/koc/hazy_svn/HAZY_VIEWS/trunk/hazy_system/Hazy_System/system_src/src/exec/../ins
tallation/functions/inplace5.c:119: warning: initialization makes pointer from integer wi
thout a cast
NOTICE: function update_trig5() does not exist, skipping
NOTICE: function update_trig_bin5() does not exist, skipping
NOTICE: trigger "training_trig5" for table "demo_training_examples" does not exist, skip
ping
view is successfully created with view id = 5
[koc@francisco] (21)$
```

Figure 6: Terminal Output of Create View

```
CREATE CLASSIFICATION VIEW demo_labeled_papers KEY id
ENTITIES FROM demo_papers KEY id
LABELS FROM Paper_Area LABEL 1
EXAMPLES FROM demo_training_examples KEY id LABEL 1
FEATURE FUNCTION tf_bag_of_words;
```

This query declares a view demo_labeled_papers that contains each paper in demo_papers where papers are labeled by HAZY-CLASSIFY as database or non-database papers. Name of the entity table that contains id and title of the publication is demo_papers. demo_training_examples table contains training examples to train a model by HAZY-CLASSIFY. Also, to extract the features of the publications, HAZY-CLASSIFY uses tf_bag_of_words. Paper_Area table specifies possible labels (in our case, database paper or not).

In order to create the view, execute the following commands:

```
cd HAZY_HOME/system_src/Hazy
cat demo/create_demo_view_query | ./exec/hazy_parser
```

As you can see from Fig 6, the view is successfully created with view id 5.

3.3 Querying the View

Since our view demo_labeled_papers is created in hazy_demo database of PostGreSQL, we can query it like a regular database table (Actually, since our strategy is eager in hazy.conf file, the view is already a regular database table). First, we need to login PostGreSQL:

```
[koc@francisco] (28)$ psql hazy_demo
psql (8.4.1)
Type "help" for help.
hazy_demo=# select id, class from demo_labeled_papers limit 10;
 id | class
 ---+----
 1 |
          Θ
  2 |
          Θ
  3 |
          Θ
  4 |
          Θ
  5 |
          Θ
  6
          Θ
  7
    8 |
          Θ
  9 |
          Θ
 10 |
          Θ
(10 rows)
hazy_demo=#
```

Figure 7: Terminal Output of Query View

psql hazy_demo

Then, we can run the following SQL query to retrieve publication id and their classes:

hazy_demo=# select id, class from demo_labeled_papers limit 10;

Output of this query is shown in Fig 7.

As we did not provide any training example, the SVM model initially classifies every entity as nondatabase paper (class is 0).

3.4 Updating the View

We can provide training examples to HAZY-CLASSIFY in order to obtain more successful classification results. For example, when we query demo_papers, we see that the title of the paper with id = 7 is "Query Processing in Spatial Network Databases" which is a database paper.

```
hazy_demo=# select * from demo_papers where id = 7;
id | name
----+
7 | Query Processing in Spatial Network Databases.
```

However, this paper is classified as non-database paper as shown in Fig 7. We can insert a training example to HAZY-CLASSIFY that the paper with id 7 is a database paper (class = 1) and it handles re-classification of entities and updating the labels of entities automatically. However, we need to have paper's feature vector since demo_training_examples table requires feature vector field. One of HAZY-CLASSIFY's internal tables (hazy_entity_table5) contains all feature vectors. Note that, 5 at the end of hazy_entity_table is id of our view. In order to insert this training example, we run the following SQL query:

hazy_demo=# insert into demo_training_examples values(7, 1, (select feature_vector FROM hazy_entity_ta

```
hazy_demo=# insert into demo_training_examples values(7, 1, (select feature_vector FROM h ▲
azy_entity_table5 where id = 7));
INSERT 0 1
hazy_demo=# insert into demo_training_examples values(32, -1, (select feature_vector FROM
hazy_entity_table5 where id = 32));
INSERT 0 1
hazy_entity_table5 where id = 22));
INSERT 0 1
hazy_demo=# []
```

Figure 8: Terminal Output of Inserting Examples

We can feed HAZY-CLASSIFY with more training examples:

```
hazy_demo=# insert into demo_training_examples values(32, -1, (select feature_vector FROM hazy_entity_hazy_demo=# insert into demo_training_examples values(22, 1, (select feature_vector FROM hazy_entity_t
```

Figure 8 shows the terminal output of inserting examples.

Now, when we re-query the view, we see that the classes of the examples change as shown in Figure 9.

hazy_demo=# select * from demo_labeled_papers;

id	eps	I	class
	+	+	
1	0	L	1
2	0	Ì.	1
3	j 0	Ĺ	1
4	0	Í.	1
5	0	Ì.	1
6	0	Ì.	1
7	j 0	Ĺ	1
8	0	Ì.	1
9	j O	Ĺ	1
10	0	I.	1
11	0	Í.	1
12	0	Ì.	1
13	0	I.	1
14	0	L	1
15	0	Ì.	1
16	0	I.	1
17	0	L	1
18	0	L	1
19	0	Ì.	Θ
20	0	I.	1
21	0	L	1
22	O	Í.	1
23	0	I.	1
24	0	T	1
25	0	Í.	1
26	Θ	I.	1
27	j 0	Í.	1
28	0	Í.	1
Mo	re		

