

# Cortana for Dummies

Marvin Meeng, Wouter Duivesteijn

November 20, 2012

## 1 What is Cortana, and what does it do?

Cortana is a tool which, given a dataset, performs Supervised Descriptive Local Pattern Mining:

**Local Pattern Mining:** a run of Cortana strives to find *subsets* of the dataset where *something interesting* is going on. We try to pinpoint local exceptionalities in your dataset.

**Descriptive:** subsets delivered by Cortana are not just any subsets of the data, but *coherent* subsets, defined on attributes. Such subsets are called *subgroups*.

**Supervised:** the interestingness is measured with respect to a user-defined *target concept*: a subset of attributes that we are particularly interested in. Cortana finds subgroups for which the target concept is substantially different than the target concept over the entire dataset.

As an example of the things Cortana could find, suppose a dataset about people, and as target concept the distribution of just one binary attribute: whether the person develops lung cancer or not. Cortana will find subgroups like *smoking = true* (the smokers, which have a substantially higher incidence of lung cancer), and *weekly walking kilometers  $\geq 50$*  (athletes, which have a substantially lower incidence of lung cancer) [TODO: might need a better example here].

## 2 Launching Cortana

After obtaining and unpacking a copy of Cortana from <http://datamining.liacs.nl/cortana.html>, navigate to the Cortana directory. There you will find a `cortana.jar` file. One could start Cortana by double clicking the jar, however, it is recommended to start Cortana from the command line. This way, a lot more information about the Subgroup Discovery process is fed back to the user. To start Cortana this way, open a terminal or command window, navigate to the Cortana directory containing `cortana.jar`, and type: `java -jar cortana.jar`. Alternatively one can use either `cortana.bat` (for Windows) or `cortana.sh` (Bash shell script). Be sure to read Appendix C if you do, and experience memory issues.

After starting Cortana one is asked to select the file that contains the dataset to be analysed, after which Cortana's main screen is shown.

Currently, two types of files can be loaded into Cortana, ARFF files and plain (comma separated) text files. The first is very common in data mining and describes the data and additionally defines the attribute type of each attribute. As such, ARFF files are preferred over plain text files. For text files, Cortana will try to infer the correct attribute type of each attribute. Unfortunately, this may not always deliver the desired outcome. We therefore recommend to check whether Cortana was able to infer the correct attribute type, via the *Meta Data...* button on the main screen. One can find more about this in Section 3.1.3.

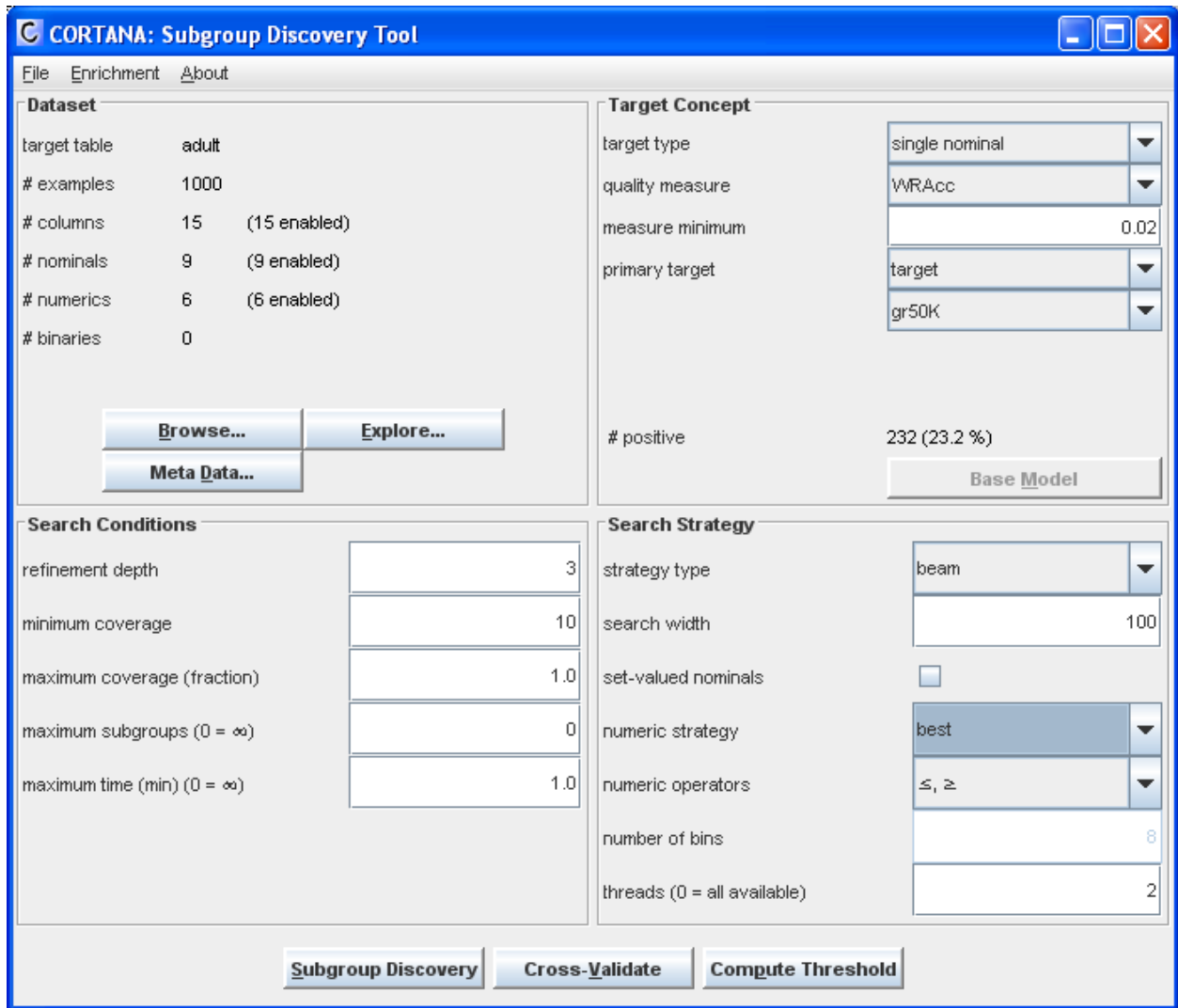


Figure 1: Cortana’s main window.

### 3 Main Screen

The main screen of Cortana is divided into four major panels, *Dataset*, *Target Concept*, *Search Conditions*, and *Search Strategy*.

#### 3.1 Dataset Panel

The dataset panel gives some information about the dataset that is currently loaded:

**target table** shows the name of the data file used, which for text files is just the filename, and for ARFF files is the name defined in the '@relation' field;

**# examples** shows the number of examples in the dataset;

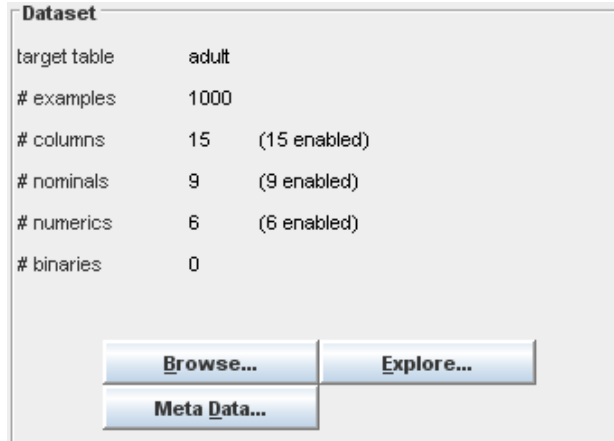


Figure 2: Dataset details in the main window.

**# columns** shows the number of columns in the dataset. Remember that columns are also referred to as attributes.

Finally, there is a number of fields that indicate the number of attributes from each data type. The type of an attribute determines the sort of mining algorithm or quality measure that is applicable to it. More about this can be found in section ??.

In addition to the fields described above, there are three buttons present on the **Dataset** panel, **Browse...**, **Meta Data...** and **Explore...**.

### 3.1.1 Browse Window

By clicking the **Browse...** button, a Browse Window is presented, showing a table with the data in the state that it is currently in. Additionally, in the table header, it shows the number of distinct values for each attribute. Note that the data may not be in the same state as when it was loaded, as it can be modified using functionality of the Meta Data Window, presented after pressing the **Meta Data...** button, which is also present on the **Dataset** panel. Section 3.1.3 describes the Meta Data Window, its components, and the data manipulation functionalities, in more detail.

### 3.1.2 Explore Window

**Explore...** button.

### 3.1.3 Meta Data Window

The **Meta Data...** button gives access to a new window that, next to displaying some additional information about the dataset loaded, also allows changing some of the (characteristics of) the data. The upper part of the *Meta Data Window* shows a table with six columns. The lower part contains a number of panels that allow modification of the data as it is in memory, note that no modifications are made to the original data file. First the properties shown in the table in the upper part will be described, the purpose of the various data manipulations will be explained after that.

Data for: adult									
age (66 distinct)	workclass (7 distinct)	fnlwgt (987 distinct)	education (16 distinct)	education-num (16 distinct)	marital-status (7 distinct)	occupation (15 distinct)	relationship (6 distinct)	race (5 distinct)	(
39	State-gov	77,516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Mal
50	Self-emp-not-i...	83,311	Bachelors	13	Married-civ-sp...	Exec-manage...	Husband	White	Mal
38	Private	215,646	HS-grad	9	Divorced	Handlers-clea...	Not-in-family	White	Mal
53	Private	234,721	11th	7	Married-civ-sp...	Handlers-clea...	Husband	Black	Mal
28	Private	338,409	Bachelors	13	Married-civ-sp...	Prof-specialty	Wife	Black	Fen
37	Private	284,582	Masters	14	Married-civ-sp...	Exec-manage...	Wife	White	Fen
49	Private	160,187	9th	5	Married-spou...	Other-service	Not-in-family	Black	Fen
52	Self-emp-not-i...	209,642	HS-grad	9	Married-civ-sp...	Exec-manage...	Husband	White	Mal
31	Private	45,781	Masters	14	Never-married	Prof-specialty	Not-in-family	White	Fen
42	Private	159,449	Bachelors	13	Married-civ-sp...	Exec-manage...	Husband	White	Mal
37	Private	280,464	Some-college	10	Married-civ-sp...	Exec-manage...	Husband	Black	Mal
30	State-gov	141,297	Bachelors	13	Married-civ-sp...	Prof-specialty	Husband	Asian-Pac-Isl...	Mal
23	Private	122,272	Bachelors	13	Never-married	Adm-clerical	Own-child	White	Fen
32	Private	205,019	Assoc-acdm	12	Never-married	Sales	Not-in-family	Black	Mal
40	Private	121,772	Assoc-voc	11	Married-civ-sp...	Craft-repair	Husband	Asian-Pac-Isl...	Mal
34	Private	245,487	7th-8th	4	Married-civ-sp...	Transport-mo...	Husband	Amer-Indian-E...	Mal
25	Self-emp-not-i...	176,756	HS-grad	9	Never-married	Farming-fishing	Own-child	White	Mal
32	Private	186,824	HS-grad	9	Never-married	Machine-op-in...	Unmarried	White	Mal
38	Private	28,887	11th	7	Married-civ-sp...	Sales	Husband	White	Mal
43	Self-emp-not-i...	292,175	Masters	14	Divorced	Exec-manage...	Unmarried	White	Fen
40	Private	193,524	Doctorate	16	Married-civ-sp...	Prof-specialty	Husband	White	Mal
54	Private	302,146	HS-grad	9	Separated	Other-service	Unmarried	Black	Fen
35	Federal-gov	76,845	9th	5	Married-civ-sp...	Farming-fishing	Husband	Black	Mal
43	Private	117,037	11th	7	Married-civ-sp...	Transport-mo...	Husband	White	Mal
59	Private	109,015	HS-grad	9	Divorced	Tech-support	Unmarried	White	Fen
56	Local-gov	216,851	Bachelors	13	Married-civ-sp...	Tech-support	Husband	White	Mal
19	Private	168,294	HS-grad	9	Never-married	Craft-repair	Own-child	White	Mal
54	?	180,211	Some-college	10	Married-civ-sp...	?	Husband	Asian-Pac-Isl...	Mal
39	Private	367,260	HS-grad	9	Divorced	Exec-manage...	Not-in-family	White	Mal
49	Private	193,366	HS-grad	9	Married-civ-sp...	Craft-repair	Husband	White	Mal

Figure 3: The browse window.

C Meta Data for: adult						
Attribute	Cardinality	Type	Enabled	Values Missing	Value for Missing	
age	66	numeric	yes	no		
workclass	7	nominal	yes	no		
fnlwgt	987	numeric	yes	no		
education	16	nominal	yes	no		
education-num	16	numeric	yes	no		
marital-status	7	nominal	yes	no		
occupation	15	nominal	yes	no		
relationship	6	nominal	yes	no		
race	5	nominal	yes	no		
sex	2	nominal	yes	no		
capital-gain	36	numeric	yes	no		
capital-loss	30	numeric	yes	no		
hours-per-week	56	numeric	yes	no		
native-country	29	nominal	yes	no		
target	2	nominal	yes	no		

Select

All
All Nominal
All Numeric
All Binary
Clear Selection

Set Type

☒ nominal
☐ numeric
☐ binary

Change Type

Set Disabled/Enabled

Disable Selected
Enable Selected
Toggle Selected

Set Value for Missing

?
Change Value

Close

Last Action: Meta Data loaded for adult

Figure 4: Meta data window.

**Meta Data Table** The first column in this table, *Attribute*, lists all *attribute names* of the attributes in the dataset. The remaining columns show some information for each of these attributes. *Cardinality* gives the number of distinct values for an attribute. *Type* shows its attribute type. *Enabled* indicates whether the attribute is *enabled* or *disabled*. *Values Missing* indicates whether the attribute contains missing values or not. And finally, *Missing Value* shows the value that is currently used for missing values in the data. Obviously, this field is blank if there are no values missing for the attribute.

**Meta Data Functions** The panels in the lower part of the *Meta Data Window* all allow selecting or changing the data. The first, *Select*, allows selecting all attributes of a certain type. This is a convenience method to be used in combination with the functionalities available in other panels.

**Set Type** allows changing the attribute type. This can be useful for various reasons. The first is that, after loading a plain text file, it is observed that Cortana was not capable to infer the correct type for an attribute. Related to this is the possibility to change the type of an attribute to allow other quality measures to be used in the Subgroup Discovery process. An example of this would be an attribute that describes the number of doors in a car dataset. If this value is used as a target value, one might treat it as a *nominal* property, forcing the Subgroup Discovery process to only perform equality tests on this *attribute value* for the creation of the conditions used to form subgroups. *Instances* in the dataset are then either in the target set if they have the same value for the ‘doors’ attribute as the selected target value, or are in the complement of the set formed by those instances. If the ‘doors’ attribute is treated as *numeric*, any, combination, of the ‘<=’, ‘>=’ and ‘=’ tests can be used to create conditions to perform on the *attribute values*. This means that the size of the set of instances selected using an *attribute value* might be bigger than in the *nominal* case. A condition using *doors* >= 2 will select all cars having two or more doors. In the *nominal* case it would not be possible to select this group using only one condition (assuming the set of cars having more than

two doors is not empty). Obviously, it would be possible to select the same group using a set of conditions like *doorsequals*‘2’  $\vee$  *doorsequals*‘3’..., but, among other negative characteristics, creating such conditions would be computationally more demanding, and less intuitive. Note that if there are missing values for an attribute, the *missing value* value for this attribute might be automatically changed to a value that is relevant to the attribute type. See *Set Value for Missing* below for more on the *missing value* values for the different attribute types.

**Set Disabled/Enabled** allows to disable or enable an attribute. When an attribute is disabled, it will not be considered by the mining algorithm to form conditions with to create subgroups. Note that disabling an attribute does not affect the possibility to select it as a target concept (see section ?? for more on target concepts).

**Set Value for Missing** can be used to change the value that is currently used for values that were missing in the data. The value that is used for missing values depends on the type of the attribute. If, in an ARFF file, values are declared missing, using the ‘?’ directive, Cortana’s file loader might replace this value with one that makes more sense in its Subgroup Discovery setting. For *nominal* types it will leave this value as is. This will result in ‘?’ being one of the possible target values one can select for the corresponding attribute. However, one can assign a different value to the *missing values*. One then has two options, either assign the *missing values* a value that is an existing one for the attribute, or a non-existing one. In the first case one effectively assigns all instances that have a missing value for the corresponding attribute to one of the other *attribute values*. In the latter case, one just changes the value. When changing the *missing value* value of an attribute, the *Cardinality* column is updated accordingly. For *numeric* and *binary* attribute types Cortana’s file loader will replace ‘?’ values with 0.0 and *false*, respectively. Again, if this is incorrect, or one wishes to assign the *missing values* another *attribute value*, either existing or non-existing, this can be done analogously to the *nominal* case.

## 3.2 Target Concept Panel

## 3.3 Search Conditions Panel

## 3.4 Search Strategy Panel

# 4 Result Window

# Appendices

## A Glossary

## B Autorun

## C Memory Issues

Although Cortana is written in Java, and therefore platform independent, it will behave slightly different on different operation systems and/or platforms. These differences arise from small variations in the Java Virtual Machines, used in different situations. The main issue is with 32-bit operating systems (OS). On such systems the maximum amount of memory the Java Virtual Machine (JVM) can use is around 1600 MegaBytes. However, the actual amount depends on the amount of RAM available. The `cortana.bat` and `cortana.sh` file included in the `Cortana.zip` set the maximum amount of memory the JVM can use to 1600 MegaBytes, through the `-Xmx` option. The value should, at most, be set to half the amount of available RAM, meaning eg. for a 2GB machine to `-Xmx1000m`. For 64-bit OSes no such limit exists, and it should be save to remove the `-Xmx`. Note that the above means that, especially for 32-bit OSes, not all datasets will fit into memory.