

CS 519 Applied Machine Learning I

Open machine learning project

Team members:

Troy McMillan, William Baker, Aaron Hudson

Topic: Outcome Prediction for Shelter Animals

Source:

The following website presents this data set as a challenge:

<https://www.kaggle.com/c/shelter-animal-outcomes/overview/description>

Problem:

Succinctly stated, Animal rights advocates from the Austin Animal Center and the ASPCA have put together data for shelter animals. It is the desire of these organizations, and our goal, to use this data to try and accurately predict the likely outcomes for animals when they leave the shelter.

Solution:

Machine Learning can be implemented to determine if a shelter animal has a high chance of being adopted or not. This will require a high amount of data preprocessing, as the initial data is highly categorical. Fortunately, it's also distinct without too many levels, so a labeler can be used to get numerical values for use. The key traits that will be used for training will be AnimalType, AgeuponOutcome, SexuponOutcome, and OutcomeType. OutcomeType will be the training target.

Since this outcome relies on so many traits, a simple classifier will likely not give good results. Thus, the group will use ensemble methods for predictions, though the inclusion of a more basic learner such as a support vector machine may be worth trying simply for comparison.

Solution: Walkthrough

Preparing the data:

Current data can be found in: https://github.com/ghost8472/CS519_2020S_PS/tree/master/scratch/data

The source data consists of 10 features, 6 of which are labeled classes (including our target feature of OutcomeType), two of which are identification features, and two time-related features. The data centers around the Outcome faced by the animal at the end of their stay at the shelter. The features are:

- A. AnimalID – a unique identifier for each animal
- B. Name – an identifier, with no guarantee of uniqueness
- C. DateTime – the date/time of when the Outcome occurred
- D. OutcomeType – the action that occurred for the Outcome

- E. OutcomeSubType – extra information about the Outcome action
- F. AnimalType – label class of either Dog or Cat
- G. SexuponOutcome – label class identifying the gender/sterilization of the animal
- H. AgeuponOutcome – the age (in days, weeks, months, or years) of the animal
- I. Breed – label class of the breed(s) of the animal, possibly including “Mix”
- J. Color – label class of the primary/secondary color(s) of the animal

Each type of feature (identification, time, and label classes) require some special attention, in order to make them useful for feeding to the machine. Used heavily are binary features, which are features where it is either existent (1) or non-existent (0), which are meant to be easy for classifiers to partition. For example, Name is converted to a binary feature by: if they have a Name it is a 1, and if they didn't it is a 0. This is important for splitting out the label classes of Breed and Color, as each may represent two classes, and important information becomes lost if each combination is merely treated with the standard labeling technique of assigning a single number to each combination.

- A. AnimalID – this unique identifier was removed, as it had no discernable numeric significance.
- B. Name – this identifier was converted into a binary feature, of whether they have a name or not.
- C. DateTime – in addition to converting the date/time into a timestamp, it also spawns into new binary features, one for each month, one for each day of the week, and ones for morning, midday, and night, for a total of 22 new features. The times of day features were based on a manual pre-analysis that their business hours are from 7am to 7pm, thereby making morning be from 7am to noon, afternoon from noon to 7pm, and night from 7pm to 7am.
- D. OutcomeType – this was moved to the end of the feature list, and left in its original form of a label class. This is the target we are using, so the standard label handling is acceptable.
- E. OutcomeSubType – this feature was removed, as it is part of the outcome, and not included in the “test” data from the source's challenge.
- F. AnimalType – this feature was turned into two binary features: one for Cat, and one for Dog. This handling keeps it consistent with the handling of the other label class features. Standard labeling technique is also maintained.
- G. SexuponOutcome – this feature was turned into four binary features: fertile/sterilized, males, females, and unknown. Standard labeling technique is also maintained.
- H. AgeuponOutcome – this feature is initially written with different time units (days, weeks, months, years), so this was standardized into a numeric value of years, for example “3 months” becomes 0.25
- I. Breed – this label class contains many instances of an animal belonging to two breeds, or a breed with a “Mix” afterward. To represent this as binary features, each breed becomes a feature, and a “Mix” feature was added as well. When an animal belongs to two breeds, the features for both breeds will be 1, and “Mix” will be 1. When an animal belongs to one breed, but is a “Mix,” then the feature for that breed will be 1, and “Mix” will be 1. If of an animal is of a single breed, only the one breed will be set to 1. Standard labeling technique is also maintained.
- J. Color – this label class contains many instances of an animal possessing two colors. Each color became a binary feature, as well as a “Mix” feature. If an animal has two colors, then the features for both colors will be 1, as will the “Mix” feature. Standard labeling technique is also retained.

Initial comparison of classifiers

Initial comparisons of the classifiers were conducted using `train_test_split` on the converted data. The data, without the extra binary features described above, was also tested, to see the impact of adding the binary features.

	Input file: traditional label encoding train_conversion_lowfeat_sansheader.csv				Input file: with added features train_conversion_extfeat_sansheader.csv			
	Training		Testing		Training		Testing	
	Time	F1	Time	F1	Time	F1	Time	F1
Perceptron	5.261	.320	.426	.324	70.166	.363	.575	.354
SVM	12.759	.365	9.390	.365	215.047	.307	191.495	.308
Decision Tree	.039	.407	.440	.414	.385	.407	.546	.414
KNN	.074	.528	1.900	.400	3.062	.456	219.731	.323
Random Forrest	.162	.295	.530	.297	.405	.269	.570	.271
Ada Boost	1.232	.388	.691	.383	10.240	.387	1.621	.381
Bagging	1.633	.407	.523	.414	5.289	.407	2.929	.414

Observations:

- The best performance was the K Nearest Neighbors classifier
- The extra features did help some classifiers, such as the Perceptron, but not the KNN classifier
- Decision Tree and Bagging were about equal, but neither was as good as KNN

Seeing that KNN is the best classifier, and the dataset utilizing standard labeling gets better results, next is comparing this combination with different parameters for the KNN classifier.

	Input file: traditional label encoding train_conversion_lowfeat_sansheader.csv											
	Testing											
	F1											
neighbors	1	2	3	4	5	6	7	8	9	10	20	30
KNN	.376	.372	.385	.400	.400	.402	.402	.408	.395	.400	.401	.400

Observations:

- The highest F1 for testing data occurred when neighbors was set to 8.
- Increasing the number of neighbors does not keep increasing the F1.
- The accuracy for the test data when neighbors is 8: 0.623

Next steps

Attempt using newer classifiers, and the parameters for them.

Appendix A: Command sequence “readme.txt”

The following commands were used to generate the data in the report:

```
python main.py perceptron data\train_conversion_lowfeat_sansheader.csv
python main.py svm data\train_conversion_lowfeat_sansheader.csv
python main.py dtree data\train_conversion_lowfeat_sansheader.csv
python main.py knn data\train_conversion_lowfeat_sansheader.csv
python main.py randforest data\train_conversion_lowfeat_sansheader.csv
python main.py adaboost data\train_conversion_lowfeat_sansheader.csv
python main.py bagging data\train_conversion_lowfeat_sansheader.csv
python main.py perceptron data\train_conversion_extfeat_sansheader.csv
python main.py svm data\train_conversion_extfeat_sansheader.csv
python main.py dtree data\train_conversion_extfeat_sansheader.csv
python main.py knn data\train_conversion_extfeat_sansheader.csv
python main.py randforest data\train_conversion_extfeat_sansheader.csv
python main.py adaboost data\train_conversion_extfeat_sansheader.csv
python main.py bagging data\train_conversion_extfeat_sansheader.csv
```

```
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=1
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=2
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=3
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=4
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=5
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=6
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=7
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=8
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=9
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=10
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=20
python main.py knn data\train_conversion_lowfeat_sansheader.csv -neighbors=30
```

Appendix B: Machine learning script

This script was used for the data processing discussed in this report. This script was run with the following command sequence:

Appendix C: Inputs conversion script

This script will take the downloaded “train.csv” and convert it into the various formats discussed in this report.

The results of this script can be downloaded from:

https://github.com/ghost8472/CS519_2020S_PS/tree/master/scratch/data

converter.php:

```
<?php

echo "Started at ".date("H:m:s");

$in_raw = file_get_contents('train.csv');
$in = explode("\n",$in_raw);

//helper function to get the Default Save Value
function DefSV($def,$obj,$key) {
    if (isset($obj[$key])) {
        return $obj[$key];
    } else {
        return $def;
    }
}

//luckily, the data is clean enough for doing dirty csv processing
$keyed = [];
$colmap = explode(",",$in[0]);
for($i = 1; $i < count($in); $i++) {
    $line = explode(",",$in[$i]);
    if (count($line) == 1) continue;

    $row = [];
    foreach ($colmap as $c=>$col) {
        $row[$col] = $line[$c];
    }
    $keyed[] = $row;
}
//echo "<pre>".var_export($keyed,true)."</pre>";

//Now, find all values, for enumeration, and/or column layout
$alls = [];
$splits = [];
foreach ($keyed as $i=>$row) { foreach($row as $col=>$val) {
    $alls[$col][$val] = $val;

    if ($col == "Breed") {
        $ex = explode("/",str_replace(" Mix","", $val));
        foreach($ex as $exv) { $splits[$col.'-'.$exv] = []; }
    } else if ($col == "Color") {
        $ex = explode("/", $val);
        foreach($ex as $exv) { $splits[$col.'-'.$exv] = []; }
    } else if ($col == "AnimalType") {
        $splits[$col.'-'.$val] = [];
    } else if ($col == "SexuponOutcome") {
        // ignore, done manually below
    } else {
        $splits[$col][$val] = $val;
    }
}}
//Add a few more columns to the extended feature version
$splits["SexuponOutcome-Female"] = [];
$splits["SexuponOutcome-Male"] = [];
$splits["SexuponOutcome-Sterilized"] = [];
$splits["SexuponOutcome-Unknown"] = [];
$splits["Breed-Mix"] = [];
$splits["Color-Mix"] = [];
for($i = 1; $i <= 12; $i++) {
    $splits["DateTime-Month{$i}"] = [];
}
for ($i = 1; $i <= 7; $i++) {
    $splits["DateTime-DOW{$i}"] = [];
```

```
}
$plits["DateTime-Morning"] = [];
$plits["DateTime-Afternoon"] = [];
$plits["DateTime-Night"] = [];

//Get rid of columns we won't be keeping
unset($alls['AnimalID']);
unset($plits['AnimalID']);
unset($alls['OutcomeSubtype']);
unset($plits['OutcomeSubtype']);

ksort($alls);
ksort($plits);

//echo "<pre>".var_export($alls,true)."</pre>";
//echo "<pre>".var_export($plits,true)."</pre>";

//Go through each row, and each data cell, and transform it.
foreach ($keyed as $i=>$row) { foreach($row as $col=>$val) {

    if ($col == "AnimalID") {
        unset($row[$col]); //this doesn't get used as an input, so would require
    }
    if ($col == "Name") {
        $row[$col] = ($val ? 1:0);
    }
    if ($col == "DateTime") {
        $row[$col] = strtotime($val);
        $ex = explode(" ", $val);
        $exh = explode(":", $ex[1]); $hour = intval($exh[0]);
        $exm = explode("-", $ex[0]); $mon = intval($exm[1]);
        $row[$col.'-Month'.'.mon'] = 1;
        $row[$col.'-DOW'.date('N')] = 1;
        if ($hour >= 7 && $hour < 12) {
            $row[$col.'-Morning'] = 1;
        } else if ($hour >= 12 && $hour <= 19) {
            $row[$col.'-Afternoon'] = 1;
        } else { // $hour < 7 || $hour > 19
            $row[$col.'-Night'] = 1;
        }
    }
    if ($col == "AgeuponOutcome") {
        $ex = explode(" ", $val);
        $ex[0] = intval($ex[0]);
        if ($val == "") {
            $row[$col] = 7.77; //unknown value, distinct though
        } else if ($ex[1] == "years" || $ex[1] == "year") {
            $row[$col] = $ex[0];
        } else if ($ex[1] == "months" || $ex[1] == "month") {
            $row[$col] = $ex[0]/12;
        } else if ($ex[1] == "weeks" || $ex[1] == "week") {
            $row[$col] = $ex[0]/52;
        } else if ($ex[1] == "days" || $ex[1] == "day") {
            $row[$col] = $ex[0]/365;
        }
    }
    if ($col == "OutcomeType") {
        //no modification, final combiner will put this last
    }
    if ($col == "OutcomeSubtype") {
        unset($row[$col]); //this is not in the "test.csv" so we can't use it for the challenge
    }
    if ($col == "AnimalType") {
        $row[$col] = array_search($val, array_keys($alls[$col]));
        $row[$col.'-'.$val] = 1;
    }
    if ($col == "SexuponOutcome") {
        if ($val == "") { $val = "Unknown"; } //one blank to fill in
        $row[$col] = array_search($val, array_keys($alls[$col]));
        $row[$col.'-Unknown'] = ($val=="Unknown"?1:0);
        $row[$col.'-Sterilized'] = (strpos($val, "Neutered") !== strpos($val, "Spayed") ? 1:0); //both
        "false" if Intact and not Unknown
        $row[$col.'-Male'] = (strpos($val, "Male") !== false ? 1:0);
        $row[$col.'-Female'] = (strpos($val, "Female") !== false ? 1:0);
    }
}
```

```
    if ($col == "Breed") {
        $row[$col] = array_search($val,array_keys($alls[$col]));
        $mix = (strpos($val," Mix") != false ? 1:0);
        $ex = explode("/",str_replace(" Mix","", $val));
        $mix = (count($ex) > 1 ? 1:$mix);
        foreach($ex as $exv) { $row[$col.'-'.$exv] = 1; }
        $row[$col.'-Mix'] = $mix;
    }
    if ($col == "Color") {
        $row[$col] = array_search($val,array_keys($alls[$col]));
        $ex = explode("/", $val);
        $mix = (count($ex) > 1 ? 1:0);
        foreach($ex as $exv) { $row[$col.'-'.$exv] = 1; }
        $row[$col.'-Mix'] = $mix;
    }
} unset($row);

//echo "<pre>".var_export($keyed,true)."</pre>";

//function to convert the definition array and data array into CSV ready lines
function Transform($define, $keyed, $header) {
    $conv = "";
    if ($header) {
        $head = [];
        foreach($define as $col=>$vals) {
            if ($col == "OutcomeType") continue;
            $head[] = $col;
        }
        $head[] = "OutcomeType";
        $conv .= implode(",",$head)."\n";
    }
    foreach ($keyed as $i=>$row) {
        $line = ""; $d = "";
        foreach($define as $col=>$vals) {
            if ($col == "OutcomeType") continue;

            $line .= $d.DefSV(0,$row,$col);
            $d = ",";
        }
        $line .= ",".DefSV("", $row, "OutcomeType");

        $conv .= $line."\n";
    }
    return $conv;
}

//write the various versions to the files
file_put_contents("train_conversion_lowfeat_withheader.csv",Transform($alls, $keyed, true));
file_put_contents("train_conversion_lowfeat_sansheader.csv",Transform($alls, $keyed, false));
file_put_contents("train_conversion_extfeat_withheader.csv",Transform($splits,$keyed, true));
file_put_contents("train_conversion_extfeat_sansheader.csv",Transform($splits,$keyed, false));

echo "<br/>Finished at ".date("H:m:s");
```