

EE-613: Machine Learning for Engineers

Practical session 10

(Decision Trees)

1 Download and test

```
1 | wget http://www.idiap.ch/~fleuret/files/EE613/EE613-pw10.tgz
2 | tar zxvf EE613-pw10.tgz
3 | cd EE613/pw10
4 | wget http://www.idiap.ch/~fleuret/files/EE613/mnist.tgz
5 | tar zxvf mnist.tgz
6 | ./do.sh -v demo graph
```

after a couple of minutes, you should have obtained a printout of the data set check, followed by the compilation and the computation itself, ending with.

```
1 | Depth 19 ... done (0.283333% / 14.47%).
```

and a display of an error graph.

2 Programming

2.1 Compilation

You can compile the source code with `make -k`.

It is suggested to implement each question in the corresponding function

```
1 | void computation_question1(VignetteSet *train_image_set,  
2 |                             VignetteSet *test_image_set) {  
3 |     ...  
4 | }
```

and to test it by calling the `./do.sh question1`, `./do.sh question2`, etc.

Note that during your development, if you compile the code with

```
1 | make -k DEBUG=yes
```

the `VignetteSet::pixel()` method will check for out-of-bound errors.

2.2 Classes

A **VignetteSet** contains small gray-scale images, each with an integer label (note that with the MNIST data set `width = height = 24`).

```
1 | class VignetteSet {  
2 | public:  
3 |     inline int nb_vignettes();  
4 |     inline int nb_classes();  
5 |     inline int width();  
6 |     inline int height();  
7 |  
8 |     inline unsigned char pixel(int p, int x, int y);  
9 |     inline unsigned char label(int p);  
10 |  
11 | void load_mnist_format(char *picture_file_name,  
12 |                       char *label_file_name);
```

```

13 |
14 |     void bootstrap(int nb, VignetteSet *vs);
15 | };

```

A **Classifier** is a virtual object that is trainable and can predict the label of a vignette.

```

1 | class Classifier {
2 | public:
3 |     virtual void train(VignetteSet *train_set) = 0;
4 |     virtual int predict(VignetteSet *vs, int n_vignette) = 0;
5 |     int nb_errors(VignetteSet *vs);
6 | };

```

A **DecisionTree** is a subclass of **Classifier** implementing a simple decision tree

```

1 | class DecisionTree : public Classifier {
2 | public:
3 |     DecisionTree(QuestionSet *question_set,
4 |                 int depth_max,
5 |                 int nb_min_samples_for_split,
6 |                 int nb_questions_for_optimization);
7 |     virtual ~DecisionTree();
8 |     virtual void train(VignetteSet *train_set);
9 |     virtual int predict(VignetteSet *vs, int n_vignette);
10 | };

```

A **QuestionSet** is a virtual object standing for a set of Boolean functions of vignettes.

```

1 | class QuestionSet {
2 | public:
3 |     virtual int nb_questions() = 0;
4 |     virtual bool response(int n_question,
5 |                          VignetteSet *vs, int n_vignette) = 0;
6 | };

```

The **PixelQuestionSet** is a set of $28 \times 28 = 784$ questions, each testing if a certain pixel of the vignette is greater than 128.

3 Questions

Question 1: Performance vs. training set size

Train one tree with 16, 32, 64, etc. up to $2^{15} = 32,768$ samples, and compute the test error each time. Use depth max 10, number of samples to split 10, and look at 100 questions at every node for training.

Help: Use `VignetteSet::bootstrap` to extract sub-sets from the full training sets.

Question 2: New feature set

Write a sub-class `MinQuestionSet` of `QuestionSet` which implements questions which test that the minimum of the pixel values over rectangles of size $p \times q$ is greater than 128. Compute the same test errors vs. depths as in the demo program. Try with $p = 4, q = 1$ and $p = 4, q = 4$.

Help: Both the `.h` and `.cc` files of that new class will be extremely similar to those of the `PixelQuestionSet` class.

Question 3: Forest

Write a subclass `Forest` of the class `Classifier` which implements a voting procedure with trees. Compute the same errors as in the demo program with 5 trees.

Help: Here is the declaration of that class

```
1 | class Forest : public Classifier {
2 |     int _nb_trees;
3 |     DecisionTree **_trees;
4 | public:
5 |
6 |     Forest(int nb_trees,
7 |           QuestionSet *question_set,
8 |           int depth_max,
9 |           int nb_min_samples_for_split,
```

```
10         int nb_questions_for_optimization);
11     ~Forest();
12     void train(VignetteSet *train_set);
13     int predict(VignetteSet *vs, int n_vignette);
14 };
```

Question 4: Nearest-Neighbors

Implement a sub-class `NearestNeighbors` of `Classifier` which implements the k -NN classifier. Compute the test error vs. the parameter k . You may reduce the training set size to run at reasonable speed.

Help: The classifier will have to keep a pointer to a `VignetteSet`.