

Knowledge-base for Word Prediction

Submit your questions to [Moshe](#) by mail.

Abstract

In this assignment you will generate a knowledge-base for English word-prediction system, based on Google 3-Gram English dataset, using Amazon Elastic Map-Reduce (EMR). The produced knowledge-base indicates the probability of each word trigram found in the corpus. In addition, you should examine the quality of your algorithm according to statistic measures and manual analysis.

The Assignment

Probability Function

In class, we presented two methods for estimating the probability of a given word sequence: the Maximum Likelihood Estimation (MLE), and back-off (bo):

$$P_{MLE}(w_1 \dots w_n) = \frac{C(w_1 \dots w_n)}{N}$$

$$P_{bo}(w_1 \dots w_n) = \alpha_1 \cdot P_{MLE}(w_1) + \alpha_2 \cdot P_{MLE}(w_1 w_2) + \dots + \alpha_n \cdot P_{MLE}(w_1 \dots w_n), \sum_{\alpha_i} = 1$$

In this assignment, we will implement a *held out* method, named *deleted estimation*.

Held out estimators divide the training data (the corpus) into two parts, build initial estimates by doing counts on one part, and then use the other pool of held out data to refine those estimates.

The **deleted estimation** method, for instance, uses a form of two-way cross validation, as follows:

$$P_{del}(w_1 \dots w_n) = \frac{T_r^{01} + T_r^{10}}{N(N_r^0 + N_r^1)}, \text{ where } C(w_1 \dots w_n) = r$$

Where:

- N is the number of n -gram instances in the whole corpus.
- N_r^0 is the number of n -gram types occurring r times in the first part of the corpus.
- T_r^{01} is the total number the n -grams of the first part (of N_r^0) appear the second part of the corpus (instances).
- N_r^1 is the number of n -gram types occurring r times in the second part of the corpus.

- T_r^{10} is the total number the n-grams of the second part (of N_r^1) appear in the first part of the corpus (instance).

An example (for 2-grams):

Part 0 of the corpus: ילד מותר ילד אסור ילד מותר

Part 1 of the corpus: ילד אסור ילד מוזר ילד אסור ילד מותר

$N = 12$ [ילד מותר, מותר ילד, ילד אסור, אסור ילד, ילד מותר, ילד אסור, אסור ילד, ילד מוזר, מוזר ילד, ילד מותר]

[מוזר ילד, ילד אסור, אסור ילד, ילד מותר]

$N_1^0 = 3$ [מוזר ילד, ילד אסור, אסור ילד]

$T_1^0 = 4$ [ילד אסור, ילד אסור, אסור ילד, אסור ילד]

$N_1^1 = 3$ [ילד מוזר, מוזר ילד, ילד מותר]

$T_1^1 = 2$ [ילד מותר, ילד מותר]

You can find some explanation on the deleted-estimation formula, given in assignment 2, in sections 6.2.3 and 6.2.4 [here](#). Make sure you understand this formula well.

Your Task

You are asked to build a map-reduce system for calculating the probability of each *trigram* (w_1, w_2, w_3) found in a given corpus, to run it on the Amazon Elastic MapReduce service, and to generate the output knowledge base with the resulted probabilities.

The input corpus is the English 3-Gram dataset of [Google Books Ngrams](#).

The output of the system is a list of word trigrams (w_1, w_2, w_3) and their probabilities ($P(w_1 w_2 w_3)$). The list should be ordered: (1) by $w_1 w_2$, ascending; (2) by the probability for $w_1 w_2 w_3$, descending.

For example (for the case of Hebrew):

קפה נמס עלית 0.6

קפה נמס מגורען 0.4

קפה שחור חזק 0.6

קפה שחור טעים 0.3

קפה שחור חם 0.1

...

שולחן עבודה ירוק 0.7

שולחן עבודה מעץ 0.3

...

Scalability, Memory Assumptions

Your code must be scalable, *i.e.*, should successfully run on much larger input. You CANNOT assume that a table which maps each r to its list of ngrams can be stored in the memory. For any case, you should justify your memory assumption.

Reports

Statistics

You are required to provide the number of key-value pairs that were sent from the mappers to the reducers in your map-reduce runs, and their size (Hint: take a look at the log file of Hadoop), **with and without local aggregation**.

Analysis

Choose 10 'interesting' word pairs and show their top-5 next words. Judge whether the system got to a reasonable decision for these cases.

Stop Words

Stop words are words which appear very frequently in the corpus. In many natural language processing algorithms, the stop words are filtered (think why). You are required to remove all bigram that contain stop words and not include them in your counts (a list of stop-words for English is provided in the assignment archive, feel free to find better lists in order to improve results).

Technical Stuff

Amazon Abstraction of Map Reduce

Amazon has introduced two abstractions for its Elastic MapReduce framework and they are: Job Flow, Job Flow Step.

Job Flow

A Job Flow is a collection of processing steps that Amazon Elastic MapReduce runs on a specified dataset using a set of Amazon EC2 instances. A Job Flow consists of one or more steps, each of which must complete in sequence successfully, for the Job Flow to finish.

Job Flow Step

A Job Flow Step is a user-defined unit of processing, mapping roughly to one algorithm that manipulates the data. A step is a Hadoop MapReduce application implemented as a Java jar or a streaming program written in Java, Ruby, Perl, Python, PHP, R, or C++. For example, to count the frequency with which words appear in a document, and output them sorted by the count, the first step would be a MapReduce application which counts the occurrences of each word, and the second step would be a MapReduce application which sorts the output from the first step based on the calculated frequencies.

Example Code

Here is a small piece of code to help you get started:

```
AWSCredentials credentials = new PropertiesCredentials(...);
AmazonElasticMapReduce mapReduce = new
AmazonElasticMapReduceClient(credentials);

HadoopJarStepConfig hadoopJarStep = new HadoopJarStepConfig()
    .withJar("s3n://yourbucket/yourfile.jar") // This should be a full map
reduce application.
    .withMainClass("some.pack.MainClass")
    .withArgs("s3n://yourbucket/input/", "s3n://yourbucket/output/");

StepConfig stepConfig = new StepConfig()
    .withName("stepname")
    .withHadoopJarStep(hadoopJarStep)
    .withActionOnFailure("TERMINATE_JOB_FLOW");

JobFlowInstancesConfig instances = new JobFlowInstancesConfig()
    .withInstanceCount(2)
    .withMasterInstanceType(InstanceType.M4Large.toString())
    .withSlaveInstanceType(InstanceType.M4Large.toString())
    .withHadoopVersion("2.6.0").withEc2KeyName("yourkey")
    .withKeepJobFlowAliveWhenNoSteps(false)
    .withPlacement(new PlacementType("us-east-1a"));

RunJobFlowRequest runFlowRequest = new RunJobFlowRequest()
    .withName("jobname")
    .withInstances(instances)
    .withSteps(stepConfig)
    .withLogUri("s3n://yourbucket/logs/");

RunJobFlowResult runJobFlowResult = mapReduce.runJobFlow(runFlowRequest);
String jobFlowId = runJobFlowResult.getJobFlowId();
System.out.println("Ran job flow with id: " + jobFlowId);
```

Notice that order of commands matters.

Reading the n-grams File

The n-grams file is in sequence file format with block level LZ0 compression. In order to read it, use the code:

```
Configuration conf = new Configuration();
Job job = new Job(conf, "...");
...
job.setInputFormatClass(SequenceFileInputFormat.class);
```

Passing Parameters from the Main to the Mappers or Reducers

You can use the "Configuration" object to pass parameters from the main to the mapper/reducer:

- In order to set the value of the parameter:

```
Configuration jobconf = new Configuration();
jobconf.set("threshold", args[3]);
//threshold is the name of the parameter - you can write whatever you
like here, and arg[3] is its value in this case.
```

- To get the value in the mapper/reducer we use the context to get the configuration and then take the parameter value from it:

```
context.getConfiguration().get("threshold", "1")
// "1" is the returned value if threshold has not been set.
```

Local single-node

During the development, you can run your code locally on a single-node cluster of Hadoop, installed on your computer.

The installation is quite simple: [Linux](#), [Windows](#).

Additional Notes

- Notice that some parts of the AWS SDK are deprecated due to SDK updates. It is okay to use these parts here.
- If you choose not to install Hadoop locally, you must pay attention to the fees, especially since the instances that are required in EMR are not included in the ["Free Usage Tier"](#). For debugging purpose, it is recommended to choose the weakest instance possible (M4Large), and the lowest number instances to complete the job. In addition, start by testing your system on a small file, and only after you make sure all of the parts work properly, move to the big corpus.
Consider every code you run, since each run is directly associated with money coming out of your budget!
- The version of Hadoop we are going to use is 2.6.0, however if you encounter any problem you may choose any other version.
- Notice that EMR uses other different products of AWS to complete its job, particularly: ec2 (management and computation) and S3 (storing logs, results, etc.). Make sure that every resource you have used is released /deleted after you're done.

Grading

- The assignment will be graded in a frontal setting.
- All information mentioned in the assignment description, or learnt in class is mandatory for the assignment.
- You will be reduced points for not reading the relevant reading material, not implementing the recommendations mentioned there, and not understanding them.
- Students belonging to the same group will not necessarily receive the same grade.
- All the requirements in the assignment will be checked, and any missing functionality will cause a point reduction. Any additional functionality will compensate on lost points. Whatever is not precisely defined in the assignment, you have the freedom to choose how to implement it.
- You should strive to make your implementation as scalable and efficient as possible, in terms of: time, memory, and money.

Submission

Submit a zip file that contains: (1) all your sources (no need to submit the jars); (2) the output mentioned above, **OR** a link to the output directory on S3 in the README file; (3) The required reports: statistics and analysis. Include a file called README that contains your usernames, names, ids, how to run the project, and any notes regarding your implementation or/and your map-reduce steps.