

Understanding Concepts of NLTK: The Natural Language Toolkit

Thanushree V, Pruthvi B, V Nagasarshan

Course name: Information Retrieval and organization

Abstract

NLTK, the Natural Language Toolkit, is a suite of open-source program modules, tutorials, and ready-to-use problem sets, providing computational linguistics courseware. NLTK covers statistical natural symbolic and language processing and is interfaced to annotated corpora. Students augment and replace existing components, learn structured programming by example, and manipulate sophisticated models from the outset.

1 Introduction

Educators starting seminars on computational etymology are often confronted with the challenge of setting up a commonsense programming part for understudy tasks and projects. This is a challenging task on the grounds that different computational etymology areas require an assortment of different information structures furthermore, abilities, and since a different scope of

subjects might be remembered for the schedule. A far-reaching practice is to use differently

programming dialects, where every language gives local information designs and abilities

that is a decent ⁻t for the main job. For models, a course could involve Prolog for standards, Perl for corpus handling, and a definite-state toolbox for morphological investigation. By depending on the underlying highlights of different dialects, the educator tries not to need to foster a great deal of programming foundation.

A lamentable result is just a significant piece of such courses should be dedicated to showing programming dialects. Further, many fascinating ventures length an assortment of

areas, and would require that different dialects be spanned. For instance, an understudy project that elaborates syntactic parsing of corpus information from a morphologically rich language may include every one of the three of the dialects referenced above: Perl for string handling; a definite state

toolbox for morphological examination; and Prolog for parsing. Obviously these extensive

overheads and weaknesses call for an innovative approach.

Aside from the functional part, computational etymology courses may likewise rely upon programming for in-class showings. This context calls for an exceptionally intuitive graphical client interfaces, making it conceivable to see program state (for example the graph of an outline parser), notice

program execution bit by bit (for instance execution of a definite-state machine), and even make

minor modifications to programs because of imagine a scenario where" inquiries from the class. Since

of these difficulties, it is normal to stay away from life shows and save classes for theoretical introductions as it were. Aside from being dull, this approach passes on understudies to address significant useful issues all alone or to manage them less efficiently in one hour.

A streamlined and versatile means of arranging the practical part of a basic computational linguistics course in this paper we introduce a fresh way of the abovementioned challenges. We describe NLTK, the Natural Language Toolkit, which we've developed in conjunction with a program.

The Natural Language Toolkit is available under a source that is available from

http://nltk.sf.net/. NLTK runs on all platforms supported by Python, including Windows, OS X, Linux, and Unix.

2 Choice of Programming Language

The most basic step in installing a practical component is selecting a programming language that works. Several considerations affected our option. First, the language will need to have a learning that is superficial, in order that novice coders get instant benefits because of their efforts. Second, the language must support quick prototyping and a short develop/test cycle; an obligatory compilation step is just a detraction that is serious. Third, the code should be self-documenting, having a syntax that is clear semantics. Fourth, it ought to be an easy task to write organized programs, ideally object-oriented but without the burden associated with languages like C++. Finally, the language should have a graphics which can be

easy-to-use to support the growth of graphical user interfaces. In surveying the available languages, we believe that Python has a fit that is especially good the above-mentioned demands. Python is an object-oriented language that is scripting by Guido van Rossum and available platforms on all (www.python.org). Python provides a learning that is shallow; it absolutely was built to be effortlessly learnt by children (van Rossum, 1999). As a language that is interpreted Python would work for rapid prototyping. Python rule is exceptionally readable and possesses been praised as "executable pseudocode." Python is a language that is object-oriented not punitively so, which is easy to encapsulate data and methods inside Python classes. Finally, Python comes with a user interface towards the Tk pictures toolkit (Lundh, 1999), and writing interfaces being graphical simple.

3 Design Criteria

Several criteria were considered in the design and implementation of the toolkit. These design criteria are listed in the order of their importance. It was also important to decide what goals the toolkit would *not* attempt to accomplish; we therefore include an explicit set of nonrequirements, which the toolkit is not expected to satisfy.



3.1 Requirements

Easy to use. The main goal of the toolkit is to allow students to focus on building natural language processing (NLP) systems. The more time students must spend learning how to use the toolkit, the less useful it will be. consistent data structures and interfaces. Expandability. The toolkit should easily accommodate new components, regardless of whether those components replicate or extend the toolkit's existing functionality. The toolkit should be structured in such a way that it is obvious where new extensions would fit into the toolkit infrastructure. Documentation. Data structures and their implementation must be carefully and comprehensively documented. All nomenclature must be carefully selected and used consistently. Simplicity. The toolkit should structure, not hide, the complexity of building NLP systems. Therefore, each class defined by the toolkit must be simple enough for a student to implement when finished. an introductory course in computational linguistics. modularity. Interaction between the different components of the toolset should be kept to a minimum using simple and well-defined interfaces. Specifically, it should be possible to complete individual projects using small pieces of the toolset without worrying about how they interact with the rest of the toolset. toolbox. This allows students to gradually learn how to use the toolkit during a course. The modularity also makes it easy to change and expand the toolkit.

3.2 non-Requirements

Completeness. The toolset is certainly not designed to provide a complete set of tools. In fact, there should be a variety of ways students can extend the toolkit. Efficiency: The toolset does not need to be heavily optimized for runtime performance. However, it must be efficient enough for students to use their NLP systems to perform real tasks Intelligence: Clean designs and implementations are preferable to clever but undecipherable ones.

4 Modules

The toolkit is enforced as a series of freelancing modules, each of which outlines a specific arrangement or task.

Throughout the toolkit, a collection of core modules defines the basic information varieties and process systems. Using the token module, you can process individual elements of text, words, or sentences. Tree modules represent tree structures over text using data structures such as syntax trees and morphological trees. There are classes in the chance module that write frequency distributions and probability distributions, as well as a range of applied mathematics smoothing techniques.

Data structures and interfaces for playacting specific NLP tasks are also defined in the remaining modules. Adding new tasks and algorithms to the toolkit can result in this list of modules growing over time. Parsing Modules

In the parser module, trees that represent text structures are represented as high-level interfaces. As part of the chunk parser module, parsers that recognize nonoverlapping linguistic groups (such as noun phrases) are offered a sub-interface. Implementations of these abstract interfaces are provided by four modules. A shift-reduce parser is implemented by the srparser module. Using a chart to record hypotheses about syntactic constituents, the chart parser module defines a flexible parser. PCFGparser offers several different probabilistic parsers grammars. And the rechunk parser module defines a transformational regular expressionbased implementation of the chunk parser interface.

Tagging Modules

In the tagger module, supplementary information such as the token's part of speech and WordNet synset tags are added to each token and several different implementations are provided for this interface. Finite State Automata

The fsa module defines a data type for encoding finite state automata; and an interface for creating automata from regular expressions.

Type Checking

An important factor in the toolkit's ease of use is the debugging time. A type-checking module reduces the number of time students must spend debugging their code by ensuring that functions are given valid arguments. All the basic types of data and processing classes utilize the type-checking module.

The toolkit may run slower since type checking is explicitly done.

Visualization

The draw. Plot chart component could possibly be utilized to graph functionality being numerical. The draw's component creates an appliance this is certainly graphical displaying and simulating declare that was limited . The draw. Chart component has a software that will be entertaining are visual tinkering with chart parsers.

The visualization segments give connects for communication and testing; they simply try not to right apply NLP information frameworks or activities. Ease of execution was thus a reduced amount of a pressing problem when considering visualization segments than it truly is for many of the different toolkit.

Visualization segments seeing that will be establish information frameworks, and visual apparatus for testing out NLP jobs. The draw tree component offers a program this is certainly smooth are visual exhibiting forest structures. The draw tree modify component offers a software for building and modifying forest buildings.

Text Classification

The classifier component describes a screen this is certainly regular classifying texts into groups. This screen happens to be applied by two segments. The classifier naive Bayes component describes a book classifier according to the Bayes this is certainly naive presumption. The classifier maxent module describes the entropy this is certainly greatest for book category and implements two formulas for exercises the unit: Generalized Iterative Scaling and Improved Iterative Scaling. The classifier features election module describes a regular user interface for selecting featuring is pertinent to get a category project this is certainly specific. Close element option can fix category show notably

The classifier feature component offers a common encoding for your info which is used to create choices regarding category job this is certainly specific. These expectations encoding allows students to try out the distinctions between various book category formulas, making use of function this is certainly similar.

5 Documentation

The toolkit is coupled with extensive paperwork which explains the toolkit and represent ways to use and expand it. This record is divided in to three classes being main

Training instructs youngsters strategies for the toolkit, within the perspective of executing jobs that include certain. Each guide focuses primarily on a site that will be unmarried these as marking, probabilistic techniques, or book category. The training adds a conversation that will be high-level details and inspires the site, followed closely by an in-depth walk-through that makes use of instances to display just how NLTK enables you to do activities.

Specialized states describe and validate the toolkit's execution and concept. The builders make use of them regarding the toolkit to steer and record the toolkit's development. College students may also seek advice from this research it's designed this way should they need more information on how the toolkit was created, and just why. Resource records produces descriptions which can be accurate every component, user interface, lessons, strategy, features, and changeable during the toolkit. It's instantly taken from docstring commentary within the Python provider laws, utilizing Epydoc (Loper, 2002).

6 Uses of NLTK

6.1 Assignments

The chunking motivates that were tutorial parsing, describes each guideline type, and offers almost all the signal that is essential the project. The code that is given in charge of loading the chunked, part-of-speech text that is marked a pre-existing tokenizer, creating an unchunked type of the writing, using the amount principles in to the unchunked text, and scoring the impact. College students pay attention to the NLP chore only – providing a guideline set aided by the coverage that is most readily useful.

Example: Amount Parsing

For instance, ChunkRule('<NN.*>') builds chunks from sequences of successive nouns; ChinkRule('<VB.>') excises verbs from present pieces; Split Rule('<NN>', '<DT>') splits any chunk that is actually existing includes a noun that will be singular closely by determiner into two pieces; and Merge Rule('<JJ>', '<JJ>') combines two adjacent chunks whenever in actuality the earliest amount ends while the second chunk begins with adjectives.

For instance, of an assignment that is moderately difficult we asked students to generate a chunk parser that correctly identifies base noun expression chunks in confirmed text, by identifying a cascade of transformational rules that are chunking. The NLTK rechunk parser module supplies many regular expressions this is certainly various rule type, that the students can instantiate to produce rules that are complete.



NLTK can help create student assignments of varying difficulty and scope. When you glance at the simplest assignments, children experiment a module that is existing. The range that will be wide of modules provide most opportunities for generating these tasks that are straightforward. As soon as students are more acquainted with the toolkit, they were able to be questioned to make adjustment that are minor extensions in a module that is present. A much more task that is challenging be produce a module that was new. Right here, NLTK provides some useful points being starting predefined connects and data and structures. existing modules that implement the interface that is same.

Within the remainder with this section, we replicate a true number of the cascades created by the youngsters. The example that is first a mixture of a few rule types:

cascade = [

ChunkRule('<DT><NN.*><VB.><NN.*>'), ChunkRule('<DT><VB.><NN.*>'), ChunkRule('<.*>'), UnChunkRule('<IN|VB.*|CC|MD|RB.*>'), UnChunkRule("<,|\\.|"|">"), MergeRule('<NN.*|DT|JJ.*|CD>', '<NN.*|DT|JJ.*|CD>'),

SplitRule('<NN.*>', '<DT|JJ>')]

The second example illustrates a brute-force approach that is statistical. The student calculated how many times each part-ofspeech tag was contained in a noun phrase. They then constructed chunks from any sequence of tags that took place a noun phrase a lot more than 50% of that time.

cascade = [

ChunkRule('<\\\$|CD|DT|EX|PDT

|PRP.*|WP.*|\\#|FW |JJ.*|NN.*|POS|RBS|WDT>*') cascade = [
 ChunkRule('<.*>+')
 ChinkRule('<VB.*|IN|CC|R.*|MD|WRB|TO|.|,>
 +')
]

6.2 Class demonstrations

NLTK supplies tools that are graphical can feel used in class demonstrations to simply help explain basic NLP concepts and algorithms. These tools which can be interactive be used to show related data structures and to demonstrate the step-bystep execution of algorithms. Both information frameworks and regulation flow tend to be easily modified during the demonstration, in reaction to questions from the class.

Since these tools that include graphical included with the toolkit, they can be used by college students. This allows pupils to experiment at home making use of algorithms that they have seen presented in class.

Example: The Chart Parsing Tool

The chart tool this is certainly parsing the process of parsing a single sentence, with a given grammar and lexicon. Its display is divided into three sections: the bottom section displays the chart; the middle section displays the sentence; and the top section displays the syntax that will be partial corresponding towards the selected edge. Buttons along the base of the windows is used to control the execution of the algorithm. The display that is main for the chart parsing appliance is shown in Figure 1.

This tool could be used to explain several different facets of chart parsing. First, it is generally used to explain chart that will be basic design, and to display how edges can represent hypotheses about syntactic constituents. It could then be used to demonstrate and describe the individual rules that the chart parser uses to generate edges that are new. Finally, it enables you to show just how whenever edges which are brand-new be part of the data. This ready of rules handles the habits that are total of parser (e.g., whether or not it parses bottom-up or top-down).

These specific rules combine to locate a parse this is certainly complete a given sentence.

The data tool that are parsing for flexible command over the parsing algorithm. The consumer can select which guideline or method they want to apply at each step of the formula. This permits the user to experiment with mixing strategies that are differente.g., topdown and bottom-up). The user can exercising controls this is certainly fine-grained the formula by selecting which edge they wish to make use of a tip to. This flexibility permits lecturers to utilize the tool to respond to a variety that is wide of

The user can determine a listing of predetermined charts to reduce the cost of setting up demonstrations during lecture. The device can be reset to then any one of them charts at anytime.

The chart tool that is parsing a typical example of a tool that is graphical by NLTK. This tool may be employed to explain the principles that were chart that will be fundamental, and to show how the algorithm works. Chart parsing try a parsing that is flexible that utilizes a information framework also known as a given information to register hypotheses about syntactic constituents. A edge represents each theory that try single the information and knowledge. A group of principles determine



Figure 1: Chart Parsing Tool

The consumer can establish a list of present charts to lessen the expense of installing demonstrations during lecture. The tool can be reset to then any one of those maps at any time. The chart software that will be parsing for versatile controls of the parsing algorithm. The user can select which rule or strategy they desire to apply at each step regarding the formula. This permits an individual to experiment with blending strategies that are different., top-down and bottom-up). The user can exercise control that is fine-grained the algorithm by selecting which edge they wish to apply a rule to. This flexibility allows lecturers to make use of the device to react to a variety that will be wide of

7.2 Conclusions and Future Work

NLTK supplies an easy, extensible, consistent platform for assignments, works, and class

demonstrations. It is well noted, easy to master, and simple to make utilize of. We hope that NLTK will enable computational linguistics classes to include more experience that is using that is hands-on building NLP components and systems.

NLTK is unique for the mixture off three issues. First, it was purposely developed as courseware and provides needs which can be pedagogical that is primary. Second, the target readers consist of both linguists and scientists that include pc plus it really is accessible and challenging at many levels of earlier skill that include computational. Sooner or later, it is reliant on a scripting that is object-oriented supporting quick prototyping and programming that is literate. We choose to continue extending the breadth of products covered from the toolkit. We become currently operating on NLTK modules for concealed Markov types, language modelling, and tree grammars that comprise adjoining. We additionally propose to boost the actual number of formulas applied by some modules which can be established for instance the book classification module. Discovering corpora that is correct a necessity for scholar tasks which were many and projects. We are therefore placing collectively a class that is matched of data that become containing for every module defined from the toolkit.NLTK is a source that will be open, and we welcome any contributions. Readers who are interested in adding to NLTK, or exactly who have ideas for modifications, add encouraged to contact the writers.

10 Acknowledgments

We are indebted to our students for feedback on the toolkit, and to anonymous reviewers, Jee Bang, and the workshop organizers for comments on an earlier version of this paper. We are grateful to Nasar Uddin and the Department of and Information Science at the Presidency university for sponsoring the work reported here.

References

- Jason Baldridge, John Dowding, and Susana Early. 2002a. Leo: an architecture for sharing resources for unification-based grammars. In *Proceedings of the Third Language Resources and Evaluation Conference*. Paris: European Language Resources Association. http://www.iccs.informatics.ed.ac.uk/ ~jmb/leo-lrec.ps.gz.
- Jason Baldridge, Thomas Morton, and Gann Bierner. 2002b. The MaxEnt project. http://maxent.sourceforge.net/.
- Kenneth R. Beesley and Lauri Karttunen. 2002. *Finite-State Morphology: Xerox Tools and Techniques*. Studies in Natural Language Processing. Cambridge University Press.

View publication stats

- Kalina Bontcheva, Hamish Cunningham, Valentin Tablan, Diana Maynard, and Oana Hamza. 2002. Using GATE as an environment for teaching NLP. In *Proceedings of the ACL Workshop on Effective Tools and Methodologies for Teaching NLP and CL*. Somerset, NJ: Association for Computational Linguistics.
- Philip R. Clarkson and Ronald Rosenfeld. 1997. Statistical language modelling using the CMU-Cambridge Toolkit. In *Proceedings of the 5th European Conference on Speech Communication and Technology (EUROSPEECH '97)*. http://svrwww.eng.cam.ac.uk/~prc14/ eurospeech97.ps.
- Ann Copestake. 2000. The (new) LKB system. http://www-csli.stanford.edu/~aac/doc5-2. pdf.

- Michael Hammond. 2002. *Programming for Linguistics: Java Technology for Language Researchers*. Oxford: Blackwell. In press.
- Jonathan Harrington and Steve Cassidy. 1999. *Techniques in Speech Acoustics*. Kluwer.
- John M. Lawler and Helen Aristar Dry, editors. 1998. *Using Computers in Linguistics*. London: Routledge.
- Edward Loper. 2002. Epydoc. http://epydoc.sourceforge.net/.
- Fredrik Lundh. 1999. An introduction to tkinter. http://www.pythonware.com/library/ tkinter/introduction/index.htm.
- Kazuaki Maeda, Steven Bird, Xiaoyi Ma, and Haejoong Lee. 2002. Creating annotation tools with the annotation graph toolkit. In *Proceedings of the Third International Conference on Language Resources and Evaluation*. http://arXiv.org/ abs/cs/0204005.
- Fernando C. N. Pereira and David H. D. Warren. 1980. Definite clause grammars for language analysis – a survey of the formalism and a comparison with augmented transition grammars. *Artificial Intelligence*, 13:231–78.
- Carl Pollard and Ivan A. Sag. 1994. *Head-Driven Phrase Structure Grammar*. Chicago University Press.
- Guido van Rossum. 1999. Computer programming for everybody. Technical report, Corporation for National Research Initiatives. http: //www.python.org/doc/essays/cp4e.html.