

Project Motivation

Over the last year I've been using a relatively new (2009 announce, first "stable" release 28 Mar 2012) programming language called [Go](#), for various projects. As it's a new language it lacks a lot the niche type libraries common in older languages. Go has an existing Porter stemming algorithm implementation and while it is old and well tested, it isn't very flexible or customizable. I decided to create a [Go implementation](#) of the [Paice/Husk](#) rule based stemmer to help fill a gap in the languages libraries, and make it freely available under a MIT style license. [A simple website](#) hosted on Google App Engine was also created to demonstrate use of the library, as well as a separate [evaluation program](#).

Architecture and Implementation

The project was split into 3 separate parts:

1. [The paicehusk stemming library](#)
2. [A demonstration website](#)
3. [The evaluation program](#)

The Paice/Husk library:

The library was developed using only functionality provided by the Go standard Library. It consists of around 450 lines of Go. The stemmer follows the algorithm defined at <http://www.comp.lancs.ac.uk/computing/research/stemming/Links/paice.htm>. The exposed API consists of the following functionality:

NewRuleTable

Takes a slice of strings in rule format, validates them and returns a new RuleTable.

```
1 // NewRuleTable returns a new RuleTable instance
2 func NewRuleTable(f []string) (table *RuleTable)
3
4 // Example usage:
5 var DefaultRules = NewRuleTable(strings.Split(defaultRules, "\n"))
```

DefaultRules

Provides a default set of rules for a user.

```
1 // DefaultRules is a default ruleset for the english language.
2 var DefaultRules = NewRuleTable(strings.Split(defaultRules, "\n"))
3
4 // Example Usage:
5 for i := range testFile {
6     fmt.Println(paicehusk.DefaultRules.Stem(testFile[i]))
7 }
```

ValidRule

Extracts a rule from the provided string using the regexp `[a-zA-Z]**?[0-9][a-zA-z]*[.>]gim`

```
1 // Validates a rule
2 func ValidRule(s string) (rule string, ok bool)
3
4 // Example Usage:
5 s, ok := ValidRule(s)
6 if !ok {
7     return nil, false
8 }
```

ParseRule

Converts the textual representation of a rule to a rule struct

```
1 // ParseRule parses a rule in the form:
2 // |suffix|intact flag|number to strip|Append|Continue flag
3 func ParseRule(s string) (r *rule, ok bool)
4
5 // Example Usage:
6 table = &RuleTable{Table: make(map[string][]*rule)}
7 for _, s := range f {
8     if r, ok := ParseRule(s); ok {
9         table.Table[r.suf[:1]] = append(table.Table[r.suf[:1]], r)
10    }
11 }
```

Stem

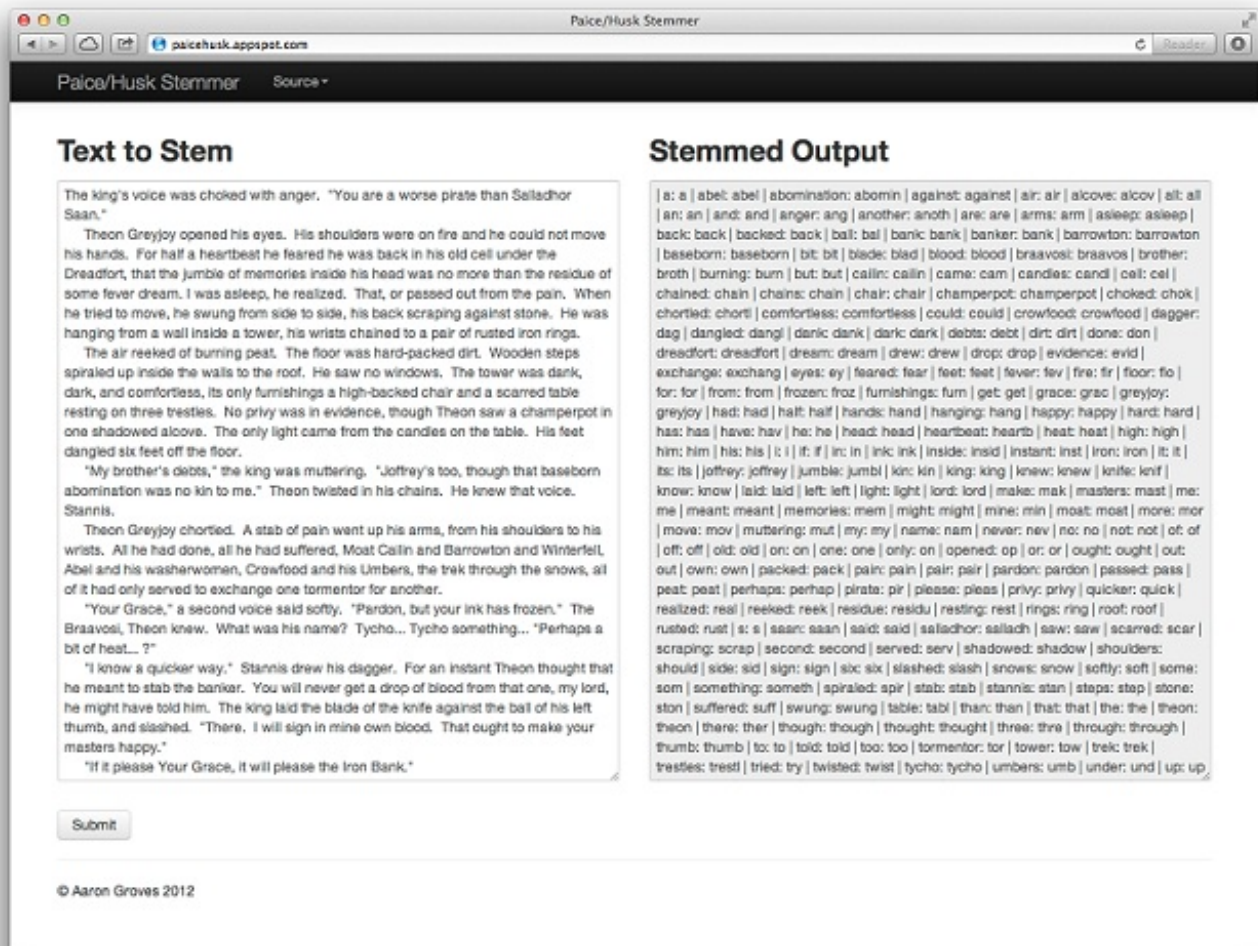
Method on RuleTable that returns the stem of the param word using the stored rules.

```
1 // Stem a string, word, based on the rules in *RuleTable r, by following
2 // the algorithm described at:
3 // http://www.comp.lancs.ac.uk/computing/research/stemming/Links/paice.htm
4 func (r *RuleTable) Stem(word string) string
5
6 // Example Usage:
7 for i := range testFile {
8     fmt.Println(paicehusk.DefaultRules.Stem(testFile[i]))
9 }
```

Demonstration Site:

To demonstrate the use of the library and to provide an easy to use stemming service, I created a [simple website](#) hosted on [Google App Engine](#). Go is one of the 3 run-times supported by App Engine, and as such, allowed me to integrate my library very quickly and easily.

I used the [Bootstrap](#) front-end framework so I didn't have to employ my terrible design skills.



The page communicates with the server via Ajax requests.

```

1 $.ajax({
2     type: "POST",
3     url: "/stem",
4     data: JSON.stringify(input),
5 }).done(function(data) {
6     var res = "|";
7     var raw = "\n\nRaw output:\n";
8     $.each(data, function(key, value){
9         res += key + ": " + value + " | ";
10        raw += value + "\n";
11    });
12    $("#output").val(res + raw);
13 });

```

The Contents of the Input text box is converted to a JSON object, then posted to the server. The server un-marshals the JSON and stems the input word by word, storing the resulting stems in a map, de-duplicating the input. This map is then converted to JSON and returned to the page, where it is written to the out put text box.

```

1 if err := decoder.Decode(text); err != nil {
2     http.Error(w, err.Error(), http.StatusInternalServerError)
3 }
4
5 words := wordreg.FindAllString(text.Text, -1)
6 for _, word := range words {
7     out[strings.ToLower(word)] = paicehusk.DefaultRules.Stem(word)
8 }
9
10 encoder := json.NewEncoder(w)

```

```
11 if err := encoder.Encode(out); err != nil {  
12     http.Error(w, err.Error(), http.StatusInternalServerError)  
13 }
```

Evaluation

Astemmers strength is based on how likely it is to merge words together. Aweak stemmer will only merge clearly related words, eg plurals by dropping the 's' or verb variants such as dropping the 'ing' of 'fighting'. Astrong stemmer will be much more aggressive, however this increases the likely hood of [stemming errors](#). The Paice/Husk stemmer is generally considered a "strong" stemming algorithm.

The Go implementation was evaluated against 3 word lists. It was also compared to an older Pascal implementation to test implementation accuracy, as well as the Porter Stemmer to provide another comparison.

The following metrics were calculated:

1. Mean number of words per conflation class
2. Index Compression Factor
3. Number of words and stems that differ
4. Mean number of characters removed
5. Median number of characters removed
6. The mode number of characters removed
7. The mean Hamming distance
8. The median Hamming distance
9. The mode Hamming distance
10. Fox and Frake Similarity Metric
11. Chris O'Neill Similarity Metric
12. Character removal distribution table
13. Hamming distance distribution table

Implementation Accuracy

I believe my implementation of the Paice/Husk algorithm to be accurate, however it does have some differences from the older Pascal version, and another C implementation that I tested against (in which I discovered several bugs). I've went over the algorithm and rule files many times, and have come to the conclusion that the older versions have errors, or have modifications to the algorithm.