# CS2309 Project Presentation: Web Crawling

Lim Jia Yee

November 11, 2016

Problem
Web Crawling
Decision Making
Parsing
Experiments

Motivation
Relevance
Solution

## Problem: All thanks to you.

- It is easy to create content on the web.
- It is easier for the web to expand now thanks to you.
- It is easiest to take web search for granted.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Motivation
Relevance
Solution

## Motivation

- Ensure web search remains optimised; prevent world destruction.
- Vested interest in learning from data (i.e. machine learning).
- Brainchild: Exploring the possibility of enhancing web crawling with decisions based on data.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Motivation
Relevance
Solution

## Relevance

- Educate on the considerations of designing a web crawler, or equivalent systems.
- Increasing the scope of machine learning as a solution.
- Target Audience: People who maintain focused web crawlers

Problem
Web Crawling
Decision Making
Parsing
Experiments

Motivation
Relevance
Solution

## Three-Part Solution

1. Requesting
2. Deciding
3. Parsing

Problem
**Web Crawling**
Decision Making
Parsing
Experiments

Context
Graph Algorithms
Algorithm Analysis

## Single Web Crawler Algorithm

1. Decide on a good hyperlink to begin crawling from.
2. Fetch the corresponding web page of the hyperlink in (1).
3. Parse for all hyperlinks and store them.
4. Process the contents of the web page.
5. From the storage of hyperlinks, extract an unvisited one.
6. Repeat from (2).

Problem
Web Crawling
Decision Making
Parsing
Experiments

Context
Graph Algorithms
Algorithm Analysis

## Graph Problem

- The World Wide Web is the graph.
- Directed, unweighted, unknown.
- Vertices: Web pages and their contents
- Edges: Hyperlinks

Problem
Web Crawling
Decision Making
Parsing
Experiments

Context
Graph Algorithms
Algorithm Analysis

## Graph Exploration

- Breadth-first search
- Depth-first search
- Iterative deepening depth-first search
- Beam search (i.e. enhanced best-first search)

Problem
Web Crawling
Decision Making
Parsing
Experiments

Context
Graph Algorithms
Algorithm Analysis

## Complexity Analysis

- Breadth-first search: Memory
- Depth-first search: Narrow scope
- Iterative deepening depth-first search: Data structure
- Beam search: Accuracy of the heuristic

Problem
Web Crawling
Decision Making
Parsing
Experiments

Context
Graph Algorithms
Algorithm Analysis

## Solution Part I: Graph Exploration

- Modified depth-first search
- Does not always push to the stack

Problem
Web Crawling
Decision Making
Parsing
Experiments

Reinforcement Learning and MDP
Components of MDP
Finite State MDP

# Alert: Another Performance Bottleneck

- Fetching the web page and parsing it.
- Why not decide beforehand whether we should even do it?

# Solution Part II: Reinforcement Learning

- Why?
- No training data.
- Unknown until experienced.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Reinforcement Learning and MDP
Components of MDP
Finite State MDP

# Markov Decision Process (MDP)

1. State
2. Action: To parse or not to parse, that is the question.
3. Reward
4. Policy

Problem
Web Crawling
Decision Making
Parsing
Experiments

Reinforcement Learning and MDP
Components of MDP
Finite State MDP

# Reinforcement Learning: State

1. How relevant the current host is.
2. How relevant the previous host was.
3. How many web pages belonging to the current host were actually parsed, and not skipped.
4. How relevant the URL is.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Reinforcement Learning and MDP
Components of MDP
Finite State MDP

# Reinforcement Learning: State

- Real-valued states $\Rightarrow$ infinite states.
- Reduce dimension via intervals.
- Result: $2^4 = 16$ (high and low features) states.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Reinforcement Learning and MDP
Components of MDP
Finite State MDP

# Reinforcement Learning: Action

- Parse or not.

# Reinforcement Learning: Reward

- Number of "high" features in the state.
- No additional penalisation of "low" features in the state.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Reinforcement Learning and MDP
Components of MDP
Finite State MDP

# Reinforcement Learning: Policy

- Value Iteration Algorithm
- Policy Iteration Algorithm

Problem
Web Crawling
Decision Making
Parsing
Experiments

Reinforcement Learning and MDP
Components of MDP
Finite State MDP

## Value Iteration Algorithm

1. Assign random **true values** to each of the states.
2. For every state, calculate a new true value based on its neighbours' current true values.
3. Terminate if any of the true values in (2) changes by more than a user-specified $\delta$. Else, repeat from (2).

# Value Iteration Algorithm: Limitations

- Slow convergence.
- Do we want true value or policy?

Problem
Web Crawling
**Decision Making**
Parsing
Experiments

Reinforcement Learning and MDP
Components of MDP
Finite State MDP

## Policy Iteration Algorithm

1. Create a random **policy** (i.e. assign a random action for each state).
2. Calculate the true value of each state given the policy in (1).
3. Based on these new true values, choose the optimal action for each state.
4. Terminate if none of the actions in (3) is changed. Else, repeat from (2).

Problem
Web Crawling
Decision Making
**Parsing**
Experiments

RAKE
word2vec

## Solution Part III: Parsing

- Not the main focus, but needed for testing.
- Define "relevant": Keywords in the web page match the list of "search words" prepared beforehand.

Problem
Web Crawling
Decision Making
**Parsing**
Experiments

RAKE
word2vec

# Rapid Automatic Keyword Extraction (RAKE)

1. Remove punctuation and special characters.
2. Remove stop words.
3. Stem the remaining words or phrases.
4. Find the degree of each word or phrase in (3).
5. Count the frequency of each word or phrase in (3).
6. Compute $score = \frac{degree}{frequency}$ for each word or phrase in (3).

Problem
Web Crawling
Decision Making
**Parsing**
Experiments

RAKE
word2vec

## RAKE: Specifications

- Assumes input is in standard English.
- Numbers are also extracted.

Problem
Web Crawling
Decision Making
**Parsing**
Experiments

RAKE
word2vec

## Similarity Measure: word2vec

- Words as vectors.
- Similarity $\Rightarrow$ distance between words.

Problem
Web Crawling
Decision Making
**Parsing**
Experiments

RAKE
word2vec

## word2vec: Limitations

- Bias towards exact words.
- Resolve by classifying a range as "similar".

Problem
Web Crawling
Decision Making
Parsing
Experiments

Modified DFS
Policy Iteration
Phrase Similarity
Future Work

## Three-Part Experiment

- Modified DFS
- Policy Iteration
- RAKE and word2vec

Problem
Web Crawling
Decision Making
Parsing
Experiments

Modified DFS
Policy Iteration
Phrase Similarity
Future Work

## Modified DFS versus BFS

- Reduced execution time (approx. 40%)
- Reduced memory usage (approx. 40%)
- Not verified: "Quality" of the visited web pages
- Can be verified via logic, but the algorithm will have to be even more accurate in targeting particular kinds of URL.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Modified DFS
Policy Iteration
Phrase Similarity
Future Work

## Final Policy

- Policy was "False" for every state $\Rightarrow$ nothing will ever be parsed.
- Policy remains the same despite increase of $\gamma$, number of iterations, and increase in penalty for skipping.
- Conclusion:
    - Insufficient domain knowledge,
    - Should not define parameters by hand, and/or
    - Unsuitable library
- Try model-free learning instead.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Modified DFS
Policy Iteration
Phrase Similarity
Future Work

## Phrase Similarity: Generality

- Keywords which are more general than search words can still score very high in similarity.
- Discovered: "Relevance" and "similarity" are not the same.

Problem
Web Crawling
Decision Making
Parsing
Experiments

Modified DFS
Policy Iteration
Phrase Similarity
Future Work

# Future Work: Beam Search and URL Analysis

1. Rank vertices by their edges (URLs).
2. Add only the first $N$ ranked edges to the heap.

## Conclusion

- Failure: Nope. Instead, you have *succeeded* in proving that this failed. Try something else.
- Failure: Is when you give up.