

Coling 2008

**22nd International Conference on
Computational Linguistics**

**Advanced Dynamic Programming in
Computational Linguistics:
Theory, Algorithms and Applications
Tutorial notes**

Liang Huang

Department of Computer and Information Science
University of Pennsylvania

©2008, Liang Huang

Licensed under the *Creative Commons Attribution-Noncommercial-Share Alike 3.0 Nonported* license
<http://creativecommons.org/licenses/by-nc-sa/3.0/>
Some rights reserved

Order copies of this and other Coling proceedings from:

Association for Computational Linguistics (ACL)
209 N. Eighth Street
Stroudsburg, PA 18360
USA
Tel: +1-570-476-8006
Fax: +1-570-476-0860
acl@aclweb.org

Design by Chimney Design, Brighton, UK
Production and manufacture by One Digital, Brighton, UK

Outline

Dynamic Programming (DP) is an important class of algorithms widely used in many areas of speech and language processing. It provides efficient solutions to seemingly intractable inference over exponentially-large spaces by sharing overlapping subproblems. Notable instances of DP include the well-known Viterbi and Forward-Backward Algorithms for finite-state models, and the CKY Algorithm for context-free parsing and syntax-based machine translation. The Dijkstra and A* Algorithms, although less obvious, can also be viewed as DP instances. In fact, almost all inference algorithms in the NLP/CL literature involve some sort of DP.

With this overwhelming popularity, a unified view of various DP algorithms would not only provide NLP researchers a better understanding but help them design new DP algorithms in practice. This tutorial, therefore, surveys two such theoretical frameworks: the semiring framework in the context of finite-state methods, and the hypergraph framework in the context of parsing and machine translation. Under each of these two paradigms, we review two most important types of DP algorithms, namely the Viterbi-style topological algorithms and the Dijkstra-style best-first algorithms. Wherever relevant, we will discuss typical instances of these algorithms in practice, which include applications in tagging and chunking, word alignment, phrase-based translation, syntactic parsing, and syntax-based translation.

Structure

The tutorial will be structured as follows:

1. Dynamic Programming on Lattices/Graphs under the Semiring Framework
 - (a) Motivations and Examples
 - (b) Semirings
 - (c) Viterbi Algorithm
 - (d) Dijkstra and A* Algorithms
 - (e) Comparison between Viterbi and Dijkstra/A* Algorithms
2. Dynamic Programming on Packed Forests under the Hypergraph Framework
 - (a) Hypergraphs and Related Formalisms
 - (b) Examples in Parsing and Machine Translation
 - (c) Generalized Viterbi Algorithm; CKY Parsing
 - (d) Knuth and A* Algorithms
3. Extensions: Non-Optimization Problems and k-best Inference
 - (a) Forward-backward and Inside-Outside Algorithms
 - (b) k-best Inference in Graphs and Hypergraphs

Prerequisites

None!

Instructor

Liang Huang// Department of Computer & Information Science (CIS)// 3330 Walnut Street// Levine Hall// University of Pennsylvania// Philadelphia, PA 19104, USA// Email: lhuang3 at cis.upenn.edu

Liang Huang is finishing his PhD study at the University of Pennsylvania, co-supervised by Aravind Joshi and Kevin Knight (USC/ISI). He is mainly interested in the theoretical aspects of computational linguistics, in particular, efficient algorithms in parsing and machine translation, generic dynamic programming, and formal properties of synchronous grammars. His thesis work develops a series of "forest-based methods" that have been applied to many problems in NLP including k-best parsing, forest rescoring and reranking, and forest-based translation. This work received the outstanding paper award at ACL 2008 and a best paper award nomination at ACL 2007. He was also a recipient of the University Prize for Excellence in Teaching at Penn, and had designed and lectured one of the first Python Programming courses for CS undergraduates in the US.

Advanced Dynamic Programming in Semiring and Hypergraph Frameworks *

Liang Huang

Department of Computer and Information Science
University of Pennsylvania
lhuang3@cis.upenn.edu

July 15, 2008

Abstract

Dynamic Programming (DP) is an important class of algorithms widely used in many areas of speech and language processing. Recently there have been a series of work trying to formalize many instances of DP algorithms under algebraic and graph-theoretic frameworks. This tutorial surveys two such frameworks, namely semirings and directed hypergraphs, and draws connections between them. We formalize two particular types of DP algorithms under each of these frameworks: the Viterbi-style topological algorithms and the Dijkstra-style best-first algorithms. Wherever relevant, we also discuss typical applications of these algorithms in Natural Language Processing.

1 Introduction

Many algorithms in speech and language processing can be viewed as instances of *dynamic programming* (DP) (Bellman, 1957). The basic idea of DP is to solve a bigger problem by divide-and-conquer, but also reuses the solutions of *overlapping subproblems* to avoid recalculation. The simplest such example is a Fibonacci series, where each $F(n)$ is used twice (if cached). The correctness of a DP algorithm is ensured by the *optimal substructure property*, which informally says that an optimal solution must contain optimal subsolutions for subproblems. We will formalize this property as an algebraic concept of *monotonicity* in Section 2.

*Survey paper to accompany the COLING 2008 tutorial on dynamic programming. The material presented here is based on the author's candidacy exam report at the University of Pennsylvania. I would like to thank Fernando Pereira for detailed comments on an earlier version of this survey. This work was supported by NSF ITR EIA-0205456.

search space \ ordering	topological-order	best-first
graph + semirings (2)	Viterbi (3.1)	Dijkstra/A* (3.2)
hypergraph + weight functions (4)	Gen. Viterbi (5.1)	Knuth/A* (5.2)

Table 1: The structure of this paper: a two dimensional classification of dynamic programming algorithms, based on search space (rows) and propagation ordering (columns). Corresponding section numbers are in parentheses.

This report surveys a two-dimensional classification of DP algorithms (see Table 1): we first study two types of search spaces (rows): the semiring framework (Mohri, 2002) when the underlying representation is a directed graph as in finite-state machines, and the hypergraph framework (Gallo et al., 1993) when the search space is hierarchically branching as in context-free grammars; then, under each of these frameworks, we study two important types of DP algorithms (columns) with contrasting order of visiting nodes: the Viterbi style topological-order algorithms (Viterbi, 1967), and the Dijkstra-Knuth style best-first algorithms (Dijkstra, 1959; Knuth, 1977). This survey focuses on *optimization problems* where one aims to find the best solution of a problem (e.g. shortest path or highest probability derivation) but other problems will also be discussed.

2 Semirings

The definitions in this section follow Kuich and Salomaa (1986) and Mohri (2002).

Definition 1. A *monoid* is a triple $(A, \otimes, \bar{1})$ where \otimes is a closed *associative binary operator* on the set A , and $\bar{1}$ is the identity element for \otimes , i.e., for all $a \in A$, $a \otimes \bar{1} = \bar{1} \otimes a = a$. A monoid is *commutative* if \otimes is commutative.

Definition 2. A *semiring* is a 5-tuple $R = (A, \oplus, \otimes, \bar{0}, \bar{1})$ such that

1. $(A, \oplus, \bar{0})$ is a commutative monoid.
2. $(A, \otimes, \bar{1})$ is a monoid.
3. \otimes distributes over \oplus : for all a, b, c in A ,

$$(a \oplus b) \otimes c = (a \otimes c) \oplus (b \otimes c),$$

$$c \otimes (a \oplus b) = (c \otimes a) \oplus (c \otimes b).$$

4. $\bar{0}$ is an *annihilator* for \otimes : for all a in A , $\bar{0} \otimes a = a \otimes \bar{0} = \bar{0}$.

Semiring	Set	\oplus	\otimes	$\bar{0}$	$\bar{1}$	intuition/application
Boolean	$\{0, 1\}$	\vee	\wedge	0	1	logical deduction, recognition
Viterbi	$[0, 1]$	max	\times	0	1	prob. of the best derivation
Inside	$\mathbb{R}^+ \cup \{+\infty\}$	+	\times	0	1	prob. of a string
Real	$\mathbb{R} \cup \{+\infty\}$	min	+	$+\infty$	0	shortest-distance
Tropical	$\mathbb{R}^+ \cup \{+\infty\}$	min	+	$+\infty$	0	with non-negative weights
Counting	\mathbb{N}	+	\times	0	1	number of paths

Table 2: Examples of semirings

Table 2 shows some widely used examples of semirings and their applications.

Definition 3. A semiring $(A, \oplus, \otimes, \bar{0}, \bar{1})$ is *commutative* if its multiplicative operator \otimes is commutative.

For example, all the semirings in Table 2 are commutative.

Definition 4. A semiring $(A, \oplus, \otimes, \bar{0}, \bar{1})$ is *idempotent* if for all a in A , $a \oplus a = a$.

Idempotence leads to a comparison between elements of the semiring.

Lemma 1. Let $(A, \oplus, \otimes, \bar{0}, \bar{1})$ be an idempotent semiring, then the relation \leq defined by

$$(a \leq b) \Leftrightarrow (a \oplus b = a)$$

is a partial ordering over A , called the natural order over A .

However, for optimization problems, a partial order is often not enough since we need to compare arbitrary pair of values, which requires a total ordering over A .

Definition 5. An idempotent semiring $(A, \oplus, \otimes, \bar{0}, \bar{1})$ is *totally-ordered* if its natural order is a total ordering.

An important property of semirings when dealing with optimization problems is *monotonicity*, which justifies the *optimal subproblem property* in dynamic programming (Cormen et al., 2001) that the computation can be factored (into smaller problems).

Definition 6. Let $K = (A, \oplus, \otimes, \bar{0}, \bar{1})$ be a semiring, and \leq a partial ordering over A . We say K is *monotonic* if for all $a, b, c \in A$

$$(a \leq b) \Rightarrow (a \otimes c \leq b \otimes c)$$

$$(a \leq b) \Rightarrow (c \otimes a \leq c \otimes b)$$

Lemma 2. Let $(A, \oplus, \otimes, \bar{0}, \bar{1})$ be an idempotent semiring, then its natural order is monotonic.

In the following section, we mainly focus on totally-ordered semirings (whose natural order is monotonic).

Another (optional) property is *superiority* which corresponds to the *non-negative weights* restriction in shortest-path problems. When superiority holds, we can explore the vertices in a best-first order as in the Dijkstra algorithm (see Section 3.2).

Definition 7. Let $K = (A, \oplus, \otimes, \bar{0}, \bar{1})$ be a semiring, and \leq a partial ordering over A . We say K is *superior* if for all $a, b \in A$

$$a \leq a \otimes b, \quad b \leq a \otimes b.$$

Intuitively speaking, superiority means the combination of two elements always gets worse (than each of the two inputs). In shortest-path problems, if you traverse an edge, you always get worse cost (longer path). In Table 2, the Boolean, Viterbi, and Tropical semirings are superior while the Real semiring is not.

Lemma 3. Let $(A, \oplus, \otimes, \bar{0}, \bar{1})$ be a superior semiring with a partial order \leq over A , then for all $a \in A$

$$\bar{1} \leq a \leq \bar{0}.$$

Proof. For all $a \in A$, we have $\bar{1} \leq \bar{1} \otimes a = a$ by superiority and $\bar{1}$ being the identity of \otimes ; on the other hand, we have $a \leq \bar{0} \otimes a = \bar{0}$ by superiority and $\bar{0}$ being the annihilator of \otimes . \square

This property, called *negative boundedness* in (Mohri, 2002), intuitively illustrates the direction of optimization from $\bar{0}$, the initial value, towards as close as possible to $\bar{1}$, the best possible value.

3 Dynamic Programming on Graphs

Following Mohri (2002), we next identify the common part shared between these two algorithms as the generic shortest-path problem in graphs.

Definition 8. A (*directed*) *graph* is a pair $G = (V, E)$ where V is the set of vertices and E the set of edges. A *weighted (directed) graph* is a graph $G = (V, E)$ with a mapping $w : E \mapsto A$ that assigns each edge a weight from the semiring $(A, \oplus, \otimes, \bar{0}, \bar{1})$.

Definition 9. The *backward-star* $BS(v)$ of a vertex v is the set of incoming edges and the *forward-star* $FS(v)$ the set of outgoing edges.

Definition 10. A *path* π in a graph G is a sequence of consecutive edges, i.e. $\pi = e_1 e_2 \cdots e_k$ where e_i and e_{i+1} are connected with a vertex. We define the *weight* (or *cost*) of path π to be

$$w(\pi) = \bigotimes_{i=1}^k w(e_i) \quad (1)$$

We denote $P(v)$ to be the set of all paths from a given source vertex s to vertex v . In the remainder of the section we only consider *single-source* shortest-path problems.

Definition 11. The *best weight* $\delta(v)$ of a vertex v is the weight of the best path from the source s to v .¹

$$\delta(v) = \begin{cases} \bar{1} & v = s \\ \bigoplus_{\pi \in P(v)} w(\pi) & v \neq s \end{cases} \quad (2)$$

For each vertex v , the current estimate of the best weight is denoted by $d(v)$, which is initialized in the following procedure:

procedure INITIALIZE(G, s)
 for each vertex $v \neq s$ **do**
 $d(v) \leftarrow \bar{0}$
 $d(s) \leftarrow \bar{1}$

The goal of a shortest-path algorithm is to repeatedly update $d(v)$ for each vertex v to some better value (based on the comparison \oplus) so that eventually $d(v)$ will converge to $\delta(v)$, a state we call *fixed*. For example, the generic update along an incoming edge $e = (u, v)$ for vertex v is²

$$d(v) \oplus = d(u) \otimes w(e) \quad (3)$$

Notice that we are using the current estimate of u to update v , so if later on $d(u)$ is updated we have to update $d(v)$ as well. This introduces the problem of *cyclic updates*, which might cause great inefficiency. To alleviate this problem, in the algorithms presented below, we will *not* trigger the update until u is fixed, so that the $u \rightarrow v$ update happens at most once.

3.1 Viterbi Algorithm for DAGs

In many NLP applications, the underlying graph exhibits some special structural properties which lead to faster algorithms. Perhaps the most common

¹By convention, if $P(v) = \emptyset$, we have $\delta(v) = \bar{0}$.

²Here we adopt the C notation where $a \oplus = b$ means the assignment $a \leftarrow a \oplus b$.

of such properties is *acyclicity*, as in Hidden Markov Models (HMMs). For acyclic graphs, we can use the Viterbi (1967) Algorithm ³ which simply consists of two steps:

1. topological sort
2. visit each vertex in the topological ordering and do updates

The pseudo-code of the Viterbi algorithm is presented in Algorithm 1.

Algorithm 1 Viterbi Algorithm.

```

1: procedure VITERBI( $G, w, s$ )
2:   topologically sort the vertices of  $G$ 
3:   INITIALIZE( $G, s$ )
4:   for each vertex  $v$  in topological order do
5:     for each edge  $e = (u, v)$  in  $BS(v)$  do
6:        $d(v) \oplus = d(u) \otimes w(e)$ 

```

The correctness of this algorithm (that $d(v) = \delta(v)$ for all v after execution) can be easily proved by an induction on the topologically sorted sequence of vertices. Basically, at the end of the outer-loop, $d(v)$ is fixed to be $\delta(v)$.

This algorithm is widely used in the literature and there have been some alternative implementations.

Variante 1. If we replace the backward-star $BS(v)$ in line 5 by the forward-star $FS(v)$ and modify the update accordingly, this procedure still works (see Algorithm 2 for pseudo-code). We refer to this variante the *forward-update* version of Algorithm 1.⁴ The correctness can be proved by a similar induction (that at the beginning of the outer-loop, $d(v)$ is fixed to be $\delta(v)$).

Algorithm 2 Forward update version of Algorithm 1.

```

1: procedure VITERBI-FORWARD( $G, w, s$ )
2:   topologically sort the vertices of  $G$ 
3:   INITIALIZE( $G, s$ )
4:   for each vertex  $v$  in topological order do
5:     for each edge  $e = (v, u)$  in  $FS(v)$  do
6:        $d(u) \oplus = d(v) \otimes w(e)$ 

```

³Also known as the Lawler (1976) algorithm in the theory community, but he considers it as part of the folklore.

⁴This is *not* to be confused with the *forward-backward* algorithm (Baum, 1972). In fact both forward and backward updates here are instances of the forward phase of a forward-backward algorithm.

Variante 2. Another popular implementation is *memoized recursion* (Cormen et al., 2001), which starts from a target vertex t and invokes recursion on sub-problems in a top-down fashion. Solved sub-problems are memoized to avoid duplicate calculation.

The running time of the Viterbi algorithm, regardless of which implementation, is $O(V + E)$ because each edge is visited exactly once.

It is important to notice that this algorithm works for all semirings as long as the graph is a DAG, although for non-total-order semirings the semantics of $\delta(v)$ is no longer “best” weight since there is no comparison. See Mohri (2002) for details.

Example 1 (Counting). Count the number of paths between the source vertex s and the target vertex t in a DAG.

Solution Use the counting semiring (Table 2).

Example 2 (Longest Path). Compute the longest (worst cost) paths from the source vertex s in a DAG.

Solution Use the semiring $(\mathbb{R} \cup \{-\infty\}, \max, +, -\infty, 0)$.

Example 3 (HMM Tagging). See Manning and Schütze (1999, Chap. 10).

3.2 Dijkstra Algorithm

The well-known Dijkstra (1959) algorithm can also be viewed as dynamic programming, since it is based on optimal substructure property, and also utilizes the overlapping of sub-problems. Unlike Viterbi, this algorithm does not require the structural property of acyclicity; instead, it requires the algebraic property of superiority of the semiring to ensure the correctness of best-first exploration.

Algorithm 3 Dijkstra Algorithm.

```

1: procedure DIJKSTRA( $G, w, s$ )
2:   INITIALIZE( $G, s$ )
3:    $Q \leftarrow V[G]$  ▷ prioritized by  $d$ -values
4:   while  $Q \neq \emptyset$  do
5:      $v \leftarrow$  EXTRACT-MIN( $Q$ )
6:     for each edge  $e = (v, u)$  in  $FS(v)$  do
7:        $d(u) \oplus = d(v) \otimes w(e)$ 
8:       DECREASE-KEY( $Q, u$ )

```

The time complexity of Dijkstra Algorithm is $O((E + V) \log V)$ with a binary heap, or $O(E + V \log V)$ with a Fibonacci heap (Cormen et al., 2001).

Since Fibonacci heap has an excessively high constant overhead, it is rarely used in real applications and we will focus on the more popular binary heap case below.

For problems that satisfy both acyclicity and superiority, which include many applications in NLP such as HMM tagging, both Dijkstra and Viterbi can apply (Nederhof, 2003). So which one is better in this case?

From the above analysis, the complexity $O((V + E) \log V)$ of Dijkstra look inferior to Viterbi's $O(V + E)$ (due to the overhead for maintaining the priority queue), but keep in mind that we can quit as long as the solution for the target vertex t is found, at which time we can ensure the current solution for the target vertex is already optimal. So the real running time of Dijkstra depends on how early the target vertex is popped from the queue, or how good is the solution of the target vertex compared to those of other vertices, and whether this early termination is worthwhile with respect to the priority queue overhead. More formally, suppose the complete solution is ranked r th among V vertices, and we prefer Dijkstra to be faster, i.e.,

$$\frac{r}{V}(V + E) \log r < (V + E),$$

then we have

$$r \log r < V \tag{4}$$

as the condition to favor Dijkstra to Viterbi. However, in many real-world applications (especially AI search, NLP parsing, etc.), often times the complete solution (a full parse tree, or a source-sink path) ranks very low among all vertices (Eq. 4 does not hold), so normally the direct use of Dijkstra does not bring speed up as opposed to Viterbi. To alleviate this problem, there is a popular technique named A* (Hart et al., 1968) described below.

3.2.1 A* Algorithm for State-Space Search

We prioritize the queue using a combination

$$d(v) \otimes \hat{h}(v)$$

of the known cost $d(v)$ from the source vertex, and an estimate $\hat{h}(v)$ of the (future) cost from v to the target t :

$$h(v) = \begin{cases} \bar{1} & v = t \\ \bigoplus_{\pi \in P(v,t)} w(\pi) & v \neq t \end{cases} \tag{5}$$

where $P(v, t)$ is the set of paths from v to t . In case where the estimate $\hat{h}(v)$ is *admissible*, namely, no worse than the true future cost $h(v)$,

$$\hat{h}(v) \leq h(v) \text{ for all } v,$$

we can prove that the optimality of $d(t)$ when t is extracted still holds. Our hope is that

$$d(t) \otimes \hat{h}(t) = d(t) \otimes \bar{1} = d(t)$$

ranks higher among $d(v) \otimes \hat{h}(v)$ and can be popped sooner. The Dijkstra Algorithm is a special case of the A* Algorithm where $\hat{h}(v) = \bar{1}$ for all v .

4 Hypergraphs

Hypergraphs, as a generalization of graphs, have been extensively studied since 1970s as a powerful tool for modeling many problems in Discrete Mathematics. In this report, we use *directed hypergraphs* (Gallo et al., 1993) to abstract a hierarchically branching search space for dynamic programming, where we solve a big problem by dividing it into (more than one) sub-problems. Classical examples of these problems include matrix-chain multiplication, optimal polygon triangulation, and optimal binary search tree (Cormen et al., 2001).

Definition 12. A (*directed*) *hypergraph* is a pair $H = \langle V, E \rangle$ with a set \mathbf{R} , where V is the set of *vertices*, E is the set of *hyperedges*, and \mathbf{R} is the set of *weights*. Each hyperedge $e \in E$ is a triple $e = \langle T(e), h(e), f_e \rangle$, where $h(e) \in V$ is its *head vertex* and $T(e) \in V^*$ is an ordered list of *tail* vertices. f_e is a *weight function* from $\mathbf{R}^{|T(e)|}$ to \mathbf{R} .

Note that our definition differs slightly from the classical definitions of Gallo et al. (1993) and Nielsen et al. (2005) where the tails are *sets* rather than *ordered lists*. In other words, we allow multiple occurrences of the same vertex in a tail and there is an ordering among the components. We also allow the head vertex to appear in the tail creating a self-loop which is ruled out in (Nielsen et al., 2005).

Definition 13. We denote $|e| = |T(e)|$ to be the *arity* of the hyperedge⁵. If $|e| = 0$, then $f_e() \in \mathbf{R}$ is a constant (f_e is a *nullary* function) and we call $h(e)$ a *source vertex*. We define the *arity* of a hypergraph to be the maximum arity of its hyperedges.

A hyperedge of arity one degenerates into an edge, and a hypergraph of arity one is standard graph.

Similar to the case of graphs, in many applications presented below, there is also a distinguished vertex $t \in V$ called *target vertex*.

We can adapt the notions of backward- and forward-star to hypergraphs.

⁵The arity of e is different from its *cardinality* defined in (Gallo et al., 1993; Nielsen et al., 2005) which is $|T(e)| + 1$.

Definition 14. The *backward-star* $BS(v)$ of a vertex v is the set of incoming hyperedges $\{e \in E \mid h(e) = v\}$. The *in-degree* of v is $|BS(v)|$. The *forward-star* $FS(v)$ of a vertex v is the set of outgoing hyperedges $\{e \in E \mid v \in T(e)\}$. The *out-degree* of v is $|FS(v)|$.

Definition 15. The *graph projection* of a hypergraph $H = \langle V, E, t, \mathbf{R} \rangle$ is a directed graph $G = \langle V, E' \rangle$ where

$$E' = \{(u, v) \mid \exists e \in BS(v), \text{s.t. } u \in T(e)\}.$$

A hypergraph H is *acyclic* if its graph projection G is acyclic; then a *topological ordering* of H is an ordering of V that is a topological ordering in G .

4.1 Weight Functions and Semirings

We also extend the concepts of monotonicity and superiority from semirings to hypergraphs.

Definition 16. A function $f : \mathbf{R}^m \mapsto \mathbf{R}$ is *monotonic* with regarding to \preceq , if for all $i \in 1..m$

$$(a_i \preceq a'_i) \Rightarrow f(a_1, \dots, a_i, \dots, a_m) \preceq f(a_1, \dots, a'_i, \dots, a_m).$$

Definition 17. A hypergraph H is *monotonic* if there is a total ordering \preceq on \mathbf{R} such that every weight function f in H is monotonic with regarding to \preceq . We can borrow the additive operator \oplus from semiring to define a comparison operator

$$a \oplus b = \begin{cases} a & a \preceq b, \\ b & \text{otherwise.} \end{cases}$$

In this paper we will assume this monotonicity, which corresponds to the optimal substructure property in dynamic programming (Cormen et al., 2001).

Definition 18. A function $f : \mathbf{R}^m \mapsto \mathbf{R}$ is *superior* if the result of function application is worse than each of its argument:

$$\forall i \in 1..m, a_i \preceq f(a_1, \dots, a_i, \dots, a_m).$$

A hypergraph H is superior if every weight function f in H is superior.

4.2 Derivations

To do optimization we need to extend the notion of paths in graphs to hypergraphs. This is, however, not straightforward due to the asymmetry of the head and the tail in a hyperedge and there have been multiple proposals in the literature. Here we follow the recursive definition of *derivations* in (Huang and Chiang, 2005). See Section 6 for the alternative notion of hyperpaths.

Definition 19. A *derivation* D of a vertex v in a hypergraph H , its size $|D|$ and its weight $w(D)$ are recursively defined as follows:

- If $e \in BS(v)$ with $|e| = 0$, then $D = \langle e, \epsilon \rangle$ is a derivation of v , its size $|D| = 1$, and its weight $w(D) = f_e()$.
- If $e \in BS(v)$ where $|e| > 0$ and D_i is a derivation of $T_i(e)$ for $1 \leq i \leq |e|$, then $D = \langle e, D_1 \cdots D_{|e|} \rangle$ is a derivation of v , its size $|D| = 1 + \sum_{i=1}^{|e|} |D_i|$ and its weight $w(D) = f_e(w(D_1), \dots, w(D_{|e|}))$.

The ordering on weights in \mathbf{R} induces an ordering on derivations: $D \preceq D'$ iff $w(D) \preceq w(D')$.

We denote $\mathcal{D}(v)$ to be the set of derivations of v and extend the *best weight* in definition 11 to hypergraph:

Definition 20. The *best weight* $\delta(v)$ of a vertex v is the weight of the best derivation of v :

$$\delta(v) = \begin{cases} \bar{1} & v \text{ is a source vertex} \\ \bigoplus_{D \in \mathcal{D}(v)} w(D) & \text{otherwise} \end{cases} \quad (6)$$

4.3 Related Formalisms

Hypergraphs are closely related to other formalisms like AND/OR graphs, context-free grammars, and deductive systems (Shieber et al., 1995; Nederhof, 2003).

In an AND/OR graph, the OR-nodes correspond to vertices in a hypergraph and the AND-nodes, which links several OR-nodes to another OR-node, correspond to a hyperedge. Similarly, in context-free grammars, nonterminals are vertices and productions are hyperedges; in deductive systems, items are vertices and instantiated deductions are hyperedges. Table 3 summarizes these correspondences. Obviously one can construct a corresponding hypergraph for any given AND/OR graph, context-free grammar, or deductive system. However, the hypergraph formulation provides greater

hypergraph	AND/OR graph	context-free grammar	deductive system
vertex	OR-node	symbol	item
source-vertex	leaf OR-node	terminal	axiom
target-vertex	root OR-node	start symbol	goal item
hyperedge	AND-node	production	instantiated deduction
$(\{u_1, u_2\}, v, f)$		$v \xrightarrow{f} u_1 u_2$	$\frac{u_1 : a \quad u_2 : b}{v : f(a, b)}$

Table 3: Correspondence between hypergraphs and related formalisms.

modeling flexibility than the weighted deductive systems of Nederhof (2003): in the former we can have a separate weight function for each hyperedge, where as in the latter, the weight function is defined for a deductive (template) rule which corresponds to many hyperedges.

5 Dynamic Programming on Hypergraphs

Since hypergraphs with weight functions are generalizations of graphs with semirings, we can extend the algorithms in Section 3 to the hypergraph case.

5.1 Generalized Viterbi Algorithm

The Viterbi Algorithm (Section 3.1) can be adapted to acyclic hypergraphs almost without modification (see Algorithm 4 for pseudo-code).

Algorithm 4 Generalized Viterbi Algorithm.

- 1: **procedure** GENERAL-VITERBI(H)
 - 2: topologically sort the vertices of H
 - 3: INITIALIZE(H)
 - 4: **for** each vertex v in topological order **do**
 - 5: **for** each hyperedge e in $BS(v)$ **do**
 - 6: e is $(\{u_1, u_2, \dots, u_{|e|}\}, v, f_e)$
 - 7: $d(v) \oplus = f_e(d(u_1), d(u_2), \dots, d(u_{|e|}))$
-

The correctness of this algorithm can be proved by a similar induction. Its time complexity is $O(V + E)$ since every hyperedge is visited exactly once (assuming the arity of the hypergraph is a constant).

The forward-update version of this algorithm, however, is not as trivial as the graph case. This is because the tail of a hyperedge now contains several vertices and thus the forward- and backward-stars are no longer symmetric. The naive adaption would end up visiting a hyperedge many

times. To ensure that a hyperedge e is fired only when all of its tail vertices have been fixed to their best weights, we maintain a counter $r[e]$ of the remaining vertices yet to be fixed (line 5) and fires the update rule for e when $r[e] = 0$ (line 9). This method is also used in the Knuth algorithm (Section 5.2).

Algorithm 5 Forward update version of Algorithm 4.

```

1: procedure GENERAL-VITERBI-FORWARD( $H$ )
2:   topologically sort the vertices of  $H$ 
3:   INITIALIZE( $H$ )
4:   for each hyperedge  $e$  do
5:      $r[e] \leftarrow |e|$  ▷ counter of remaining tails to be fixed
6:   for each vertex  $v$  in topological order do
7:     for each hyperedge  $e$  in  $FS(v)$  do
8:        $r[e] \leftarrow r[e] - 1$ 
9:       if  $r[e] == 0$  then ▷ all tails have been fixed
10:          $e$  is  $(\{u_1, u_2, \dots, u_{|e|}\}, h(e), f_e)$ 
11:          $d(h(e)) \oplus = f_e(d(u_1), d(u_2), \dots, d(u_{|e|}))$ 

```

5.1.1 CKY Algorithm

The most widely used algorithm for parsing in NLP, the CKY algorithm (Kasami, 1965), is a specific instance of the Viterbi algorithm for hypergraphs. The CKY algorithm takes a context-free grammar G in Chomsky Normal Form (CNF) and essentially intersects G with a DFA D representing the input sentence to be parsed. The resulting search space by this intersection is an acyclic hypergraph whose vertices are items like (X, i, j) and whose hyperedges are instantiated deductive steps like $(\{(Y, i, k)(Z, k, j)\}, (X, i, j), f)$ for all $i < k < j$ if there is a production $X \rightarrow YZ$. The weight function f is simply

$$f(a, b) = a \otimes b \otimes w(X \rightarrow YZ).$$

The Chomsky Normal Form ensures acyclicity of the hypergraph but there are multiple topological orderings which result in different variants of the CKY algorithm, e.g., bottom-up CKY, left-to-right CKY, and right-to-left CKY, etc.

5.2 Knuth Algorithm

Knuth (1977) generalizes the Dijkstra algorithm to what he calls the *grammar problem*, which essentially corresponds to the search problem in a monotonic superior hypergraph (see Table 3). However, he does not provide

an efficient implementation nor analysis of complexity. Graehl and Knight (2004) present an implementation that runs in time $O(V \log V + E)$ using the method described in Algorithm 5 to ensure that every hyperedge is visited only once (assuming the priority queue is implemented as a Fibonacci heap; for binary heap, it runs in $O((V + E) \log V)$).

Algorithm 6 Knuth Algorithm.

```

1: procedure KNUTH( $H$ )
2:   INITIALIZE( $H$ )
3:    $Q \leftarrow V[H]$  ▷ prioritized by  $d$ -values
4:   for each hyperedge  $e$  do
5:      $r[e] \leftarrow |e|$ 
6:   while  $Q \neq \emptyset$  do
7:      $v \leftarrow \text{EXTRACT-MIN}(Q)$ 
8:     for each edge  $e$  in  $FS(v)$  do
9:        $e$  is  $(\{u_1, u_2, \dots, u_{|e|}\}, h(e), f_e)$ 
10:       $r[e] \leftarrow r[e] - 1$ 
11:      if  $r[e] == 0$  then
12:         $d(h(e)) \oplus = f_e(d(u_1), d(u_2), \dots, d(u_{|e|}))$ 
13:        DECREASE-KEY( $Q, h(e)$ )

```

5.2.1 A* Algorithm on Hypergraphs

We can also extend the A* idea to hypergraphs to speed up the Knuth Algorithm. A specific case of this algorithm is the A* parsing of Klein and Manning (2003) where they achieve significant speed up using carefully designed heuristic functions. More formally, we first need to extend the concept of (exact) *outside cost* from Eq. 5:

$$\alpha(v) = \begin{cases} \bar{1} & v = t \\ \bigoplus_{D \in \mathcal{D}(v,t)} w(D) & v \neq t \end{cases} \quad (7)$$

where $\mathcal{D}(v, t)$ is the set of (partial) derivations using v as a leaf node. This outside cost can be computed from top-down following the inverse topological order: for each vertex v , for each incoming hyperedge $e = (\{u_1, \dots, u_{|e|}\}, v, f_e) \in BS(v)$, we update

$$\alpha(u_i) \oplus = f_e(d(u_1) \dots d(u_{i-1}), \alpha(v), d(u_{i+1}) \dots d(u_{|e|})) \text{ for each } i.$$

Basically we replace $d(u_i)$ by $\alpha(v)$ for each i . In case weight functions are composed of semiring operations, as in shortest paths (+) or probabilistic

grammars (\times), this definition makes sense, but for general weight functions there should be some formal requirements to make the definition sound. However, this topic is beyond the scope of this paper.

6 Extensions and Discussions

In most of the above we focus on optimization problems where one aims to find the best solution. Here we consider two extensions of this scheme: *non-optimization problems* where the goal is often to compute the summation or closure, and *k-best problems* where one also searches for the 2nd, 3rd, through *k*th-best solutions. Both extensions have many applications in NLP. For the former, algorithms based on the Inside semiring (Table 1), including the forward-backward algorithm (Baum, 1972) and Inside-Outside algorithm (Baker, 1979; Lari and Young, 1990) are widely used for unsupervised training with the EM algorithm (Dempster et al., 1977). For the latter, since NLP is often a pipeline of several modules, where the 1-best solution from one module might *not* be the best input for the next module, and one prefers to postpone disambiguation by propagating a *k*-best list of candidates (Collins, 2000; Gildea and Jurafsky, 2002; Charniak and Johnson, 2005; Huang and Chiang, 2005). The *k*-best list is also frequently used in discriminative learning to approximate the whole set of candidates which is usually exponentially large (Och, 2003; McDonald et al., 2005).

6.1 Beyond Optimization Problems

We know that in optimization problems, the criteria for using dynamic programming is *monotonicity* (definitions 6 and 16). But in non-optimization problems, since there is no comparison, this criteria is no longer applicable. Then when can we apply dynamic programming to a non-optimization problem?

Cormen et al. (1990) develop a more general criteria of *closed semiring* where \oplus is idempotent and infinite sums are well-defined and present a more sophisticated algorithm that can be proved to work for all closed semirings. This definition is still not general enough since many non-optimization semirings including the Inside semiring are not even idempotent. Mohri (2002) solves this problem by a slightly different definition of *closedness* which does not assume idempotence. His generic single-source algorithm subsumes many classical algorithms like Dijkstra, Bellman-Ford (Bellman, 1958), and Viterbi as specific instances.

It remains an open problem how to extend the *closedness* definition to the case of weight functions in hypergraphs.

6.2 k -best Extensions

The straightforward extension from 1-best to k -best is to simply replace the old semiring $(A, \oplus, \otimes, \bar{0}, \bar{1})$ by its k -best version $(A^k, \oplus^k, \otimes^k, \bar{0}^k, \bar{1}^k)$ where each element is now a vector of length k , with the i th component represent the i th-best value. Since \oplus is a comparison, we can define \oplus^k to be the top- k elements of the $2k$ elements from the two vectors, and \otimes^k the top- k elements of the k^2 elements from the cross-product of two vectors:

$$a \oplus^k b = \oplus'_k(\{a_i \mid 1 \leq i \leq k\} \cup \{b_j \mid 1 \leq j \leq k\})$$

$$a \otimes^k b = \oplus'_k\{a_i \otimes b_j \mid 1 \leq i, j \leq k\}$$

where \oplus'_k returns the *ordered list* of the top- k elements in a set. A similar construction is obvious for the weight functions in hypergraphs.

Now we can re-use the 1-best Viterbi Algorithm to solve the k -best problem in a generic way, as is done in (Mohri, 2002). However, some more sophisticated techniques that breaks the modularity of semirings results in much faster k -best algorithms. For example, the Recursive Enumeration Algorithm (REA) (Jiménez and Marzal, 1999) uses a lazy computation method on top of the Viterbi algorithm to efficiently compute the i th-best solution based on the 1st, 2nd, ..., $(i - 1)$ th solutions. A simple k -best Dijkstra algorithm is described in (Mohri and Riley, 2002).

For the hypergraph case, the REA algorithm has been adapted for k -best derivations (Jiménez and Marzal, 2000; Huang and Chiang, 2005). Applications of this algorithm include k -best parsing (McDonald et al., 2005; Mohri and Roark, 2006) and machine translation (Chiang, 2007). It is also implemented as part of Dyna (Eisner et al., 2005), a generic language for dynamic programming. The k -best extension of the Knuth Algorithm is studied by Huang (2005). A separate problem, k -shortest hyperpaths, has been studied by Nielsen et al. (2005).

Eppstein (2001) compiles an annotated bibliography for k -shortest-path and other related k -best problems.

7 Conclusion

This report surveys two frameworks for formalizing dynamic programming and presents two important classes of DP algorithms under these frameworks. We focused on 1-best optimization problems but also discussed other scenarios like non-optimization problems and k -best solutions. We believe that a better understanding of the theoretical foundations of DP is beneficial for NLP researchers.

References

- Baker, James K. 1979. Trainable grammars for speech recognition. In *Proceedings of the Spring Conference of the Acoustical Society of America*, pages 547–550.
- Baum, L. E. 1972. An inequality and associated maximization technique in statistical estimation of probabilistic functions of a markov process. *Inequalities*, (3).
- Bellman, Richard. 1957. *Dynamic Programming*. Princeton University Press.
- Bellman, Richard. 1958. On a routing problem. *Quarterly of Applied Mathematics*, (16).
- Charniak, Eugene and Mark Johnson. 2005. Coarse-to-fine-grained n -best parsing and discriminative reranking. In *Proceedings of the 43rd ACL*.
- Chiang, David. 2007. Hierarchical phrase-based translation. In *Computational Linguistics*. To appear.
- Collins, Michael. 2000. Discriminative reranking for natural language parsing. In *Proceedings of ICML*, pages 175–182.
- Cormen, Thomas H., Charles E. Leiserson, and Ronald L. Rivest. 1990. *Introduction to Algorithms*. MIT Press, first edition.
- Cormen, Thomas H., Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. 2001. *Introduction to Algorithms*. MIT Press, second edition.
- Dempster, A. P., N. M. Laird, and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39:1–38.
- Dijkstra, Edsger W. 1959. A note on two problems in connexion with graphs. *Numerische Mathematik*, (1):267–271.
- Eisner, Jason, Eric Goldlust, and Noah A. Smith. 2005. Compiling comp ling: Weighted dynamic programming and the dyna language. In *Proceedings of HLT-EMNLP*.
- Eppstein, David. 2001. Bibliography on k shortest paths and other “ k best solutions” problems. <http://www.ics.uci.edu/~eppstein/bibs/kpath.bib>.
- Gallo, Giorgio, Giustino Longo, and Stefano Pallottino. 1993. Directed hypergraphs and applications. *Discrete Applied Mathematics*, 42(2):177–201.
- Gildea, Daniel and Daniel Jurafsky. 2002. Automatic labeling of semantic roles. *Computational Linguistics*, 28(3):245–288.
- Graehl, Jonathan and Kevin Knight. 2004. Training tree transducers. In *HLT-NAACL*, pages 105–112.
- Hart, P. E., N. J. Nilsson, and B. Raphael. 1968. A formal basis for the heuristic determination of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics*, 4(2):100–107.
- Huang, Liang. 2005. k -best Knuth algorithm and k -best A* parsing. Unpublished manuscript.
- Huang, Liang and David Chiang. 2005. Better k -best Parsing. In *Proceedings of the Ninth International Workshop on Parsing Technologies (IWPT-2005)*.

- Jiménez, Víctor and Andrés Marzal. 1999. Computing the k shortest paths: A new algorithm and an experimental comparison. In *Algorithm Engineering*, pages 15–29.
- Jiménez, Víctor M. and Andrés Marzal. 2000. Computation of the n best parse trees for weighted and stochastic context-free grammars. In *Proc. of the Joint IAPR International Workshops on Advances in Pattern Recognition*.
- Kasami, T. 1965. An efficient recognition and syntax analysis algorithm for context-free languages. Technical Report AFCRL-65-758, Air Force Cambridge Research Laboratory, Bedford, MA†.
- Klein, Dan and Chris Manning. 2003. A* parsing: Fast exact Viterbi parse selection. In *Proceedings of HLT-NAACL*.
- Knuth, Donald. 1977. A generalization of Dijkstra’s algorithm. *Information Processing Letters*, 6(1).
- Kuich, W. and A. Salomaa. 1986. *Semirings, Automata, Languages*. Number 5 in EATCS Monographs on Theoretical Computer Science. Springer-Verlag, Berlin, Germany.
- Lari, K. and S. J. Young. 1990. The estimation of stochastic context-free grammars using the inside-outside algorithm. *Computer Speech and Language*, 4:35–56.
- Lawler, E. L. 1976. *Combinatorial Optimization: Networks and Matroids*. Holt, Rinehart, and Winston.
- Manning, Chris and Hinrich Schütze. 1999. *Foundations of Statistical Natural Language Processing*. MIT Press.
- McDonald, Ryan, Koby Crammer, and Fernando Pereira. 2005. Online large-margin training of dependency parsers. In *Proceedings of the 43rd ACL*.
- Mohri, Mehryar. 2002. Semiring frameworks and algorithms for shortest-distance problems. *Journal of Automata, Languages and Combinatorics*, 7(3):321–350.
- Mohri, Mehryar and Michael Riley. 2002. An efficient algorithm for the n -best-strings problem. In *Proceedings of the International Conference on Spoken Language Processing 2002 (ICSLP ’02)*, Denver, Colorado, September.
- Mohri, Mehryar and Brian Roark. 2006. Probabilistic context-free grammar induction based on structural zeros. In *Proceedings of HLT-NAACL*.
- Nederhof, Mark-Jan. 2003. Weighted deductive parsing and Knuth’s algorithm. 29(1):135–143.
- Nielsen, Lars Relund, Kim Allan Andersen, and Daniele Pretolani. 2005. Finding the k shortest hyperpaths. *Computers and Operations Research*.
- Och, Franz Joseph. 2003. Minimum error rate training in statistical machine translation. In *Proceedings of ACL*, pages 160–167.
- Shieber, Stuart, Yves Schabes, and Fernando Pereira. 1995. Principles and implementation of deductive parsing. *Journal of Logic Programming*, 24:3–36.
- Viterbi, Andrew J. 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory*, IT-13(2):260–269, April.