

These notes analyze algorithms for optimization problems involving matchings in bipartite graphs. Matching algorithms are not only useful in their own right (e.g., for matching clients to servers in a network, or buyers to sellers in a market) but also furnish a concrete starting point for learning many of the recurring themes in the theory of graph algorithms and algorithms in general. Examples of such themes are augmenting paths, linear programming relaxations, and primal-dual algorithm design.

1 Bipartite maximum matching

In this section we introduce the bipartite maximum matching problem, present a naïve algorithm with $O(mn)$ running time, and then present and analyze an algorithm due to Hopcroft and Karp that improves the running time to $O(m\sqrt{n})$.

1.1 Definitions

Definition 1. A *bipartite graph* is a graph whose vertex set is partitioned into two disjoint sets L, R such that each edge has one endpoint in L and the other endpoint in R . When we write a bipartite graph G as an ordered triple $G = (L, R, E)$, it means that L and R are the two vertex sets (called the *left set* and *right set*, respectively) and E is the edge set.

Definition 2. A *matching* in an undirected graph is a set of edges such that no vertex belongs to more than one element of the set.

The *bipartite maximum matching problem* is the problem of computing a matching of maximum cardinality in a bipartite graph.

We will assume that the input to the bipartite maximum matching problem, $G = (L, R, E)$, is given in its adjacency list representation, and that the bipartition of G —that is, the partition of the vertex set into L and R —is given as part of the input to the problem.

Exercise 1. Prove that if the bipartition is not given as part of the input, it can be constructed from the adjacency list representation of G in linear time.

(Here and elsewhere in the lecture notes for CS 6820, we will present exercises that may improve your understanding. You are encouraged to attempt to solve these exercises, but they are not homework problems and we will make no effort to check if you have solved them, much less grade your solutions.)

1.2 Alternating paths and cycles; augmenting paths

The following sequence of definitions builds up to the notion of an *augmenting path*, which plays a central role in the design of algorithms for the bipartite maximum matching problem.

Definition 3. If G is a graph and M is a matching in G , a vertex is called *matched* if it belongs to one of the edges in M , and *free* otherwise.

An *alternating component with respect to M* (also called an *M -alternating component*) is an edge set that forms a connected subgraph of G of maximum degree 2 (i.e., a path or cycle), in which every degree-2 vertex belongs to exactly one edge of M . An *augmenting path with respect to M* is an M -alternating component which is a path both of whose endpoints are free vertices.

In the following lemma, and throughout these notes, we use the notation $A \oplus B$ to denote the *symmetric difference* of two sets A and B , i.e. the set of all elements that belong to one of the sets but not the other.

Lemma 1. *If M is a matching and P is an augmenting path with respect to M , then $M \oplus P$ is a matching containing one more edge than M .*

Proof. P has an odd number of edges, and its edges alternate between belonging to M and its complement, starting and ending with the latter. Therefore, $M \oplus P$ has one more edge than M . To see that it is a matching, note that vertices in the complement of P have the same set of neighbors in M as in $M \oplus P$, and vertices in P have exactly one neighbor in $M \oplus P$. \square

Lemma 2. *A matching M in a graph G is a maximum cardinality matching if and only if it has no augmenting path.*

Proof. We have seen in Lemma 1 that if M has an augmenting path, then it does not have maximum cardinality, so we need only prove the converse. Suppose that M^* is a matching of maximum cardinality and that $|M| < |M^*|$. The edge set $M \oplus M^*$ has maximum degree 2, and each vertex of degree 2 in $M \oplus M^*$ belongs to exactly one edge of M . Therefore each connected component of $M \oplus M^*$ is an M -alternating component. At least one such component must contain more edges of M^* than of M . It cannot be an alternating cycle or an even-length alternating path; these have an equal number of edges of M^* and M . It also cannot be an odd-length alternating path that starts and ends in M . Therefore it must be an odd-length alternating path that starts and ends in M^* . Since both endpoints of this path are free with respect to M , it is an M -augmenting path as desired. \square

1.3 Bipartite maximum matching: Naïve algorithm

The foregoing discussion suggests the following general scheme for designing a bipartite maximum matching algorithm.

Algorithm 1 Naïve iterative scheme for computing a maximum matching

- 1: Initialize $M = \emptyset$.
 - 2: **repeat**
 - 3: Find an augmenting path P with respect to M .
 - 4: $M \leftarrow M \oplus P$
 - 5: **until** there is no augmenting with respect to M .
-

By Lemma 1, the invariant that M is a matching is preserved at the end of each loop iteration. Furthermore, each loop iteration increases the cardinality of M by 1, and the cardinality cannot exceed $n/2$, where n is the number of vertices of G . Therefore, the algorithm terminates after at most $n/2$ iterations. When it terminates, M is guaranteed to be a maximum matching by Lemma 2.

The algorithm is not yet fully specified because we have not indicated the procedure for finding an augmenting path with respect to M . When G is a bipartite graph, there is a simple linear-time procedure that we now describe.

Definition 4. If $G = (L, R, E)$ is a bipartite graph and M is a matching, the graph $D(G, M)$ is the directed graph formed from G by orienting each edge from L to R if it does not belong to M , and from R to L otherwise.

Lemma 3. *Suppose M is a matching in a bipartite graph G , and let F denote the set of free vertices. M -augmenting paths are in one-to-one correspondence with directed paths from $L \cap F$ to $R \cap F$ in $D(G, M)$.*

Proof. If P is a directed path from $L \cap F$ to $R \cap F$ in $D(G, M)$ then P starts and ends at free vertices, and its edges alternate between those that are directed from L to R (which are in the complement of M) and those that are directed from R to L (which are in M), so the undirected edge set corresponding to P is an augmenting path.

Conversely, if P is an augmenting path, then each vertex in the interior of P belongs to exactly one edge of M , so when we orient the edges of P as in $D(G, M)$ each vertex in the interior of P has exactly one incoming and one outgoing edge, i.e. P becomes a directed path. This path has an odd number of edges so it has one endpoint in L and the other endpoint in R . Both of these endpoints belong to F , by the definition of augmenting paths. Thus, the directed edge set corresponding to P is a path in $D(G, M)$ from $L \cap F$ to $R \cap F$. \square

Lemma 3 implies that in each loop iteration of Algorithm 1, the step that requires finding an augmenting path (if one exists) can be implemented by building the auxiliary graph $D(G, M)$ and running a graph search algorithm such as BFS or DFS to search for a path from $L \cap F$ to $R \cap F$. Building $D(G, M)$ takes $O(m + n)$ time, where m is the number of edges in G , as does searching $D(G, M)$ using BFS or DFS. For convenience, assume $m \geq n/2$; otherwise G contains isolated vertices which may be eliminated in a preprocessing step requiring only $O(n)$ time. Then Algorithm 1 runs for at most $n/2$ iterations, each requiring $O(m)$ time, so its running time is $O(mn)$.

Remark 1. When G is not bipartite, our analysis of Algorithm 1 still proves that it finds a maximum matching after at most $n/2$ iterations. However, the task of finding an augmenting path, if one exists, is much more subtle. The first polynomial-time algorithm for finding an augmenting path was discovered by Jack Edmonds in a 1965 paper entitled “Paths, Trees, and Flowers” that is one of the most influential papers in the history of combinatorial optimization. Edmonds’ algorithm finds an augmenting path in $O(mn)$ time, leading to a running time of $O(mn^2)$ for finding a maximum matching in a non-bipartite graph. Faster algorithms have subsequently been discovered.

1.4 The Hopcroft-Karp algorithm

One potentially wasteful aspect of the naïve algorithm for bipartite maximum matching is that it chooses one augmenting path in each iteration, even if it finds many augmenting paths in the process of searching the auxiliary graph $D(G, M)$. The Hopcroft-Karp algorithm improves the running time of the naïve algorithm by correcting this wasteful aspect; in each iteration it attempts to find many disjoint augmenting paths, and it uses all of them to increase the size of M .

The following definition specifies the type of structure that the algorithm searches for in each iteration.

Definition 5. If G is a graph and M is a maximum matching, a *blocking set of augmenting paths* with respect to M is a set $\{P_1, \dots, P_k\}$ of augmenting paths such that:

1. the paths P_1, \dots, P_k are vertex disjoint;
2. they all have the same length, ℓ ;
3. ℓ is the minimum length of an M -augmenting path;
4. every augmenting path of length ℓ has at least one vertex in common with $P_1 \cup \dots \cup P_k$.

In other words, a blocking set of augmenting paths is a (setwise) maximal collection of vertex-disjoint minimum-length augmenting paths.

The following lemma generalizes Lemma 1 and its proof is a direct generalization of the proof of that lemma.

Lemma 4. *If M is a matching and $\{P_1, \dots, P_k\}$ is any set of vertex-disjoint M -augmenting paths then $M \oplus P_1 \oplus P_2 \oplus \dots \oplus P_k$ is a matching of cardinality $|M| + k$.*

Generalizing Lemma 2 we have the following.

Lemma 5. *Suppose G is a graph, M is a matching in G , and M^* is a maximum matching; let $k = |M^*| - |M|$. The edge set $M \oplus M^*$ contains at least k vertex-disjoint M -augmenting paths. Consequently, G has at least one M -augmenting path of length less than n/k , where n denotes the number of vertices of G .*

Proof. The edge set $M \oplus M^*$ has maximum degree 2, and each vertex of degree 2 in $M \oplus M^*$ belongs to exactly one edge of M . Therefore each connected component of $M \oplus M^*$ is an M -alternating component. Each M -alternating component which is *not* an augmenting path has at least as many edges in M as in M^* . Each M -augmenting path has exactly one fewer edge in M as in M^* . Therefore, at least k of the connected components of $M \oplus M^*$ must be M -augmenting paths, and they are all vertex-disjoint. To prove the final sentence of the lemma, note that G has only n vertices, so it cannot have k disjoint subgraphs each with more than n/k vertices. \square

These lemmas suggest the following method for finding a maximum matching in a graph, which constitutes the outer loop of the Hopcroft-Karp algorithm.

Algorithm 2 Hopcroft-Karp algorithm, outer loop

- 1: $M = \emptyset$
 - 2: **repeat**
 - 3: Let $\{P_1, \dots, P_k\}$ be a blocking set of augmenting paths with respect to M .
 - 4: $M \leftarrow M \oplus P_1 \oplus P_2 \oplus \dots \oplus P_k$
 - 5: **until** there is no augmenting path with respect to M
-

The key to the improved running-time guarantee is the following pair of lemmas which culminate in an improved bound on the number of outer-loop iterations.

Lemma 6. *The minimum length of an M -augmenting path strictly increases after each iteration of the Hopcroft-Karp outer loop in which a non-empty blocking set of augmenting paths is found.*

Proof. We will use the following notation.

$$\begin{aligned} M &= \text{matching at the start of one loop iteration} \\ P_1, \dots, P_k &= \text{blocking set of augmenting paths found} \\ Q &= P_1 \cup \dots \cup P_k \\ R &= E \setminus Q \\ M' &= M \oplus Q = \text{matching at the end of the iteration} \\ F &= \{\text{vertices that are free with respect to } M\} \\ F' &= \{\text{vertices that are free with respect to } M'\} \\ d(v) &= \text{length of shortest path in } D(G, M) \text{ from } L \cap F \text{ to } v \\ &\quad (\text{If no such path exists, } d(v) = \infty.) \end{aligned}$$

If (x, y) is any edge of $D(G, M)$ then $d(y) \leq d(x) + 1$. Edges of $D(G, M)$ that satisfy $d(y) = d(x) + 1$ will be called *advancing* edges, and all other edges will be called *retreating* edges. Note that a shortest path in $D(G, M)$ from $L \cap F$ to any vertex v must be formed entirely from advancing edges. In particular, Q is contained in the set of advancing edges.

In the edge set of $D(G, M')$, the orientation of every edge in Q is reversed and the orientation of every edge in R is preserved. Therefore, $D(G, M')$ has three types of directed edges (x, y) :

1. reversed edges of Q , which satisfy $d(y) = d(x) - 1$;
2. advancing edges of R , which satisfy $d(y) = d(x) + 1$;
3. retreating edges of R , with satisfy $d(y) \leq d(x)$.

Note that in all three cases, the inequality $d(y) \leq d(x) + 1$ is satisfied.

Now let ℓ denote the minimum length of an augmenting path with respect to M , i.e. $\ell = \min\{d(v) \mid v \in R \cap F\}$. Let P be any path in $D(G, M')$ from $L \cap F'$ to $R \cap F'$. The lemma asserts that P has at least ℓ edges. The endpoints of P are free in M' , hence also in M . As w ranges over the vertices of P , the value $d(w)$ increases from 0 to at least ℓ , and each edge of P increases the value of $d(w)$ by at most 1. Therefore P has at least ℓ

edges, and the only way that it can have ℓ edges is if $d(y) = d(x) + 1$ for each edge (x, y) of P . We have seen that this implies that P is contained in the set of advancing edges of R , and in particular P is edge-disjoint from Q . It cannot be vertex-disjoint from Q because then $\{P_1, \dots, P_k, P\}$ would be a set of $k + 1$ vertex-disjoint minimum-length M -augmenting paths, violating our assumption that $\{P_1, \dots, P_k\}$ is a blocking set. Therefore P has at least one vertex in common with P_1, \dots, P_k , i.e. $P \cap Q \neq \emptyset$. The endpoints of P cannot belong to Q , because they are free in M' whereas every vertex in Q is matched in M' . Let w be a vertex in the interior of P which belongs to Q . The edge of M' containing w belongs to P , but it also belongs to Q . This violates our earlier conclusion that P is edge-disjoint from Q , yielding the desired contradiction. \square

Lemma 7. *The Hopcroft-Karp algorithm terminates after fewer than $2\sqrt{n}$ iterations of its outer loop.*

Proof. After the first \sqrt{n} iterations of the outer loop are complete, the minimum length of an M -augmenting path is greater than \sqrt{n} . This implies, by Lemma 5, that $|M^*| - |M| < \sqrt{n}$, where M^* denotes a maximum cardinality matching. Each remaining iteration strictly increases $|M|$, hence there are fewer than \sqrt{n} iterations remaining. \square

The inner loop of the Hopcroft-Karp algorithm must compute a blocking set of augmenting paths with respect to M . We now describe how to do this in linear time.

Recalling the distance labels $d(v)$ defined in the proof of Lemma 6; $d(v)$ is the length of the shortest alternating path from a free vertex in L to v ; if no such path exists $d(v) = \infty$. Recall also that an *advancing* edge in $D(G, M)$ is an edge (x, y) such that $d(y) = d(x) + 1$, and that every minimum-length M -augmenting path is composed exclusively of advancing edges. The Hopcroft-Karp inner loop begins by performing a breadth-first search to compute the distance labels $d(v)$, along with the set A of advancing edges and a counter $c(v)$ for each vertex that counts the number of incoming advancing edges at v , i.e. advancing edges of the form (u, v) for some vertex u . It sets ℓ to be the minimum length of an M -augmenting path (equivalently, the minimum of $d(v)$ over all $v \in R \cap F$), marks every vertex as unexplored, and repeatedly finds augmenting paths using the following procedure. Start at an unexplored vertex v in $R \cap F$ such that $d(v) = \ell$, and trace backward along incoming edges in A until a vertex u with $d(u) = 0$ is reached. Add this path P to the blocking set and add its vertices to a “garbage collection” queue. While the garbage collection queue is non-empty, remove the vertex v at the head of the queue, mark it as explored, and delete its incident edges (both outgoing and incoming) from A . When deleting an outgoing edge (v, w) , decrement the counter $c(w)$, and if $c(w)$ is now equal to 0, then add u to the garbage collection queue.

The inner loop performs only a constant number of operations per edge — traversing it during the BFS that creates the set A , traversing it while creating the blocking set of paths, deleting it from A during garbage collection, and decrementing its tail’s counter during garbage collection — and a constant number of operations per vertex: visiting it during the BFS that creates the set A , initializing $d(v)$ and $c(v)$, visiting it during the search for the blocking set of paths, marking it as explored, inserting it into the garbage collection queue, and removing it from that queue. Therefore, the entire inner loop runs in linear time.

By design, the algorithm discovers a set of minimum-length M -augmenting paths that are vertex disjoint, so we need only prove that this set is maximal. By induction on the number of augmenting paths the algorithm has discovered, the following invariants hold whenever the garbage collection queue is empty.

1. For every vertex v , $c(v)$ counts the number of advancing edges (u, v) that have not yet been deleted from A .
2. Whenever an edge e is deleted or a vertex v is placed into the garbage collection queue, any path made up of advancing edges that starts in $L \cap F$ and includes edge e or vertex v must have a vertex in common with the selected set of paths.
3. For every unmarked vertex v , $c(v) > 0$ and there exists a path in A from $L \cap F$ to v . (The existence of such a path follows by tracing backwards along edges of A from v to a vertex u such that $d(u) = 0$.)

The third invariant ensures that whenever the algorithm starts searching for an augmenting path at an unmarked free vertex, it is guaranteed to find such a path. The second invariant ensures that when there are no longer any unmarked free vertices v with $d(v) = \ell$, the set of advancing edges no longer contains a path from $L \cap F$ to $R \cap F$ that is vertex-disjoint from the selected ones; thus, the selected set forms a blocking set of augmenting paths as desired.

2 Non-bipartite matching

When the graph G is not bipartite, Lemma 2 is still valid: a matching has maximum cardinality if and only if it has no augmenting path. Hence, as before, the problem of finding a maximum matching reduces to the problem of finding an augmenting path with respect to a given matching, or else certifying that there is none. However, whereas in the bipartite case the problem of finding an augmenting path reduced to searching for a path in the directed graph $D(G, M)$, in the non-bipartite case there is no correspondingly simple reduction.

To see why, it's useful to consider what goes wrong with the naïve idea of searching for an augmenting path using “breadth-first search over the set of alternating paths”. Here's one way of making this information idea precise. For a graph G and matching M , define $H(G, M)$ to be a directed graph with the same set of vertices as G , and with a directed edge (u, v) for every pair of vertices such that G contains a path made up of two edges (u, u') and (u', v) such that $(u, u') \notin M$ and $(u', v) \in M$. Note that if G contains an M -augmenting path made up of $2k + 1$ edges, then the first $2k$ of those edges correspond to a k -edge path in $H(G, M)$ that starts in F , the set of free vertices, and ends in $\Gamma(F)$, the set of vertices that are adjacent to a free vertex.

If the converse were true, i.e. if finding a path in $H(G, M)$ from F to $\Gamma(F)$ were equivalent to finding an M -augmenting path in G , then we could design a maximum non-bipartite matching algorithm along exactly the same lines as in the bipartite case. Instead, there is a second alternative represented by the diagram in Figure 1: a simple path in $H(G, M)$ from F to $\Gamma(F)$ might correspond to a *non-simple* alternating walk from F to $N(F)$ in G , i.e. an alternating walk that repeats some vertices.

Definition 6. If G is a graph and M is a matching in G , a *flower* with respect to M is an M -alternating walk u_0, u_1, \dots, u_s such that:

1. u_0 is a free vertex with respect to M ;
2. the vertices u_0, \dots, u_{s-1} are distinct, whereas $u_s = u_r$ for some number $r < s$
3. r is even and s is odd.

The *stem* of the flower is the path u_0, u_1, \dots, u_r . (Note that it is possible that $r = 0$, in which case the stem is an empty path.) The *blossom* of the flower is the cycle u_r, \dots, u_s .

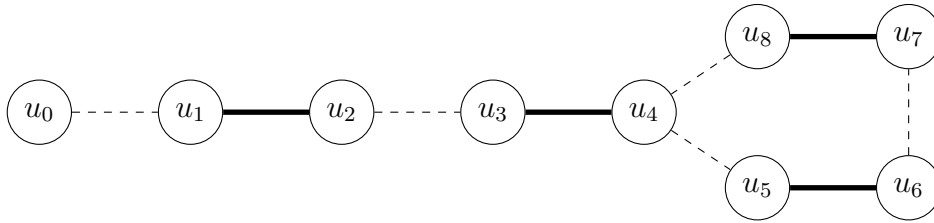


Figure 1: A flower

Lemma 8. *If the graph $H(G, M)$ contains a path P from F to $\Gamma(F)$, then G contains either an M -augmenting path or a flower.*

Proof. Suppose that $P = v_0, v_1, \dots, v_k$ is a simple path in $H(G, M)$ from F to $\Gamma(F)$. For $i = 1, 2, \dots, k$, the edge (v_{i-1}, v_i) in $H(G, M)$ corresponds to a sequence of two edges in G , the first lying outside M and the second belonging to M . Denote these two edges by $(u_{2i-2}, u_{2i-1}) \notin M$ and $(u_{2i-1}, u_{2i}) \in M$. Since $v_k = u_{2k}$ belongs to $\Gamma(F)$, we may choose a free vertex u_{2k+1} adjacent to u_{2k} . Now consider the alternating walk $u_0, u_1, \dots, u_{2k+1}$ in G . If all of its vertices are distinct, then it is an M -augmenting path. Otherwise, let u_s be the earliest instance of a repeated vertex in the alternating walk, and let u_r denote the earlier occurrence of this same vertex.

If s is even, then the edge (u_{s-1}, u_s) belongs to M . Note that this means s is not a free vertex, so $r > 0$. This means that either (u_{r-1}, u_r) or (u_r, u_{r+1}) belongs to M , hence u_{s-1} is equal to either u_{r-1} or u_{r+1} . This contradicts our choice of s unless $r + 1 = s - 1$, which is impossible because the case $(u_r, u_{r+1}) \in M$ only occurs when r is odd, in which case $r + 1 \neq s - 1$ because the left side is even and the right side is odd. Hence, the assumption that s is even leads to a contradiction.

If s and r are both odd, then $(u_r, u_{r+1}) \in M$ so u_s is not a free vertex. In particular this means that $s < 2k + 1$. The edge (u_s, u_{s+1}) belongs to M , which implies that $u_{s+1} = u_{r+1}$. However, since $s + 1$ and $r + 1$ are even, the vertices u_{s+1} and u_{r+1} both belong to P , contradicting the assumption that P is a simple path.

By process of elimination, we have deduced that s is odd and r is even, in which case the sequence u_0, \dots, u_s constitutes a flower. \square

If G contains a flower with blossom B , our algorithm for finding an M -augmenting path in G will depend on an operation called *blossom shrinking* which forms a new graph G/B with matching M/B , as follows. The vertices of B are replaced with a single vertex $\{v_b\}$. Edges having both endpoints in B are removed. For those having exactly one endpoint in B , that endpoint is changed to v_b and the other endpoint is preserved. Edges having no endpoints in B are unchanged. Note that this operation may produce a multigraph (i.e., there may be multiple edges between the same two vertices) in the case that there is a vertex having more than one neighbor in B . In the event that G/B contains multiple edges between the same two vertices, we can discard all but one of those edges without affecting the algorithm's correctness; however, in our analysis we prefer to treat G/B as a multigraph because it means that every edge of G/B has one unambiguous corresponding edge in G , which simplifies the analysis.

Let M/B denote the set of edges in G/B whose corresponding edge in G belongs to M . Note that M/B is a matching: for all vertices other than v_b it is clear that they belong to at most one edge in M/B , while for v_b this holds because if u_r, u_{r+1}, \dots, u_s denotes the list of vertices in B , in the order that they occur in the flower, then u_r (also known as u_s) is the only vertex in the blossom that potentially belongs to an edge of M whose other endpoint lies outside of M .

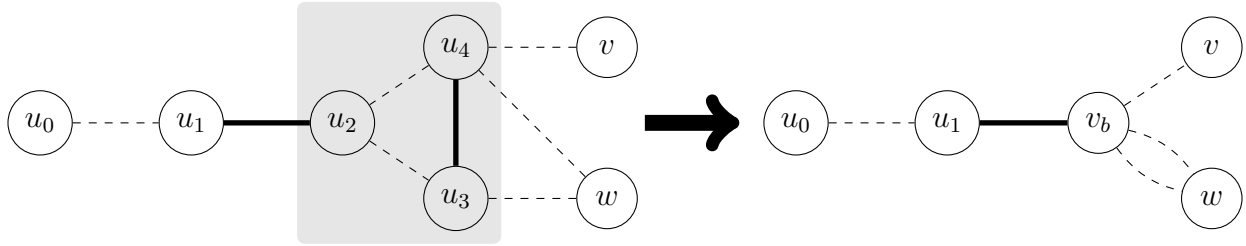


Figure 2: Shrinking a blossom

The following lemma on the relationship between augmenting paths in G and those in G/B accounts for the importance of the blossom shrinking operation.

Lemma 9. *If G is a graph, M is a matching, and B is the blossom of a flower with respect to M , then G/B contains an (M/B) -augmenting path if and only if G contains an M -augmenting path. Furthermore, any (M/B) -augmenting path in G/B can be modified into an M -augmenting path in G in time $|B|$.*

Proof. Denote the vertices of the flower containing B by u_0, \dots, u_s , numbered as in Definition 6. If P is an (M/B) -augmenting path in G/B and P does not contain v_b , then it is already an M -augmenting path in G . Otherwise, P contains an edge (w, v_b) that does not belong to M/B . Let (w, u_t) be the corresponding edge of G , where $r \leq t < s$. An M -augmenting path in G can be constructed by replacing edge (w, v_b) with an M -alternating path from w to u_r in G whose first and last edges do not belong to M . If t is even, then replace edge (w, u_t) with path $w, u_t, u_{t-1}, \dots, u_r$; if t is odd, then replace (w, u_t) with $w, u_t, u_{t+1}, \dots, u_s$. Notice that the path segment that replaces (w, u_t) has fewer than $|B|$

edges, and the replacement can be done in $O(|B|)$ time if we use suitable data structures, e.g. representing the path P as a doubly linked list of edges. This justifies the running time bound in the last sentence of the lemma statement.

It remains for us to prove that if G contains an M -augmenting path, then G/B contains an (M/B) -augmenting path. One might expect this to be a simple matter of reversing the operation defined in the preceding paragraph, but in fact it's a little trickier. To see why, consider the augmenting path $\langle v, u_4, u_3, w \rangle$ in Figure 2. The corresponding path in G/B is $\langle v, v_b, w \rangle$ which is not an alternating path with respect to M/B .

Instead, we first let $S = \{(u_{i-1}, u_i) \mid i = 1, \dots, r\}$ denote the set of edges belonging to the stem of the flower, and we modify M to $M' = M \oplus S$. Note that $|M'| = |M|$ and u_r, u_{r+1}, u_s is a flower with respect to M' (having an empty stem). Hence, B is still a blossom with respect to M' , and M'/B is still a matching in G/B , with the same number of edges as M/B . We will now apply the following chain of reasoning to deduce the existence of an (M/B) -augmenting path in G/B .

$$\begin{aligned}
G \text{ has an } M\text{-augmenting path} &\Rightarrow M \text{ is not a maximum matching in } G && \text{(Lemma 2)} \\
&\Rightarrow M' \text{ is not a maximum matching in } G && (|M| = |M'|) \\
&\Rightarrow G \text{ has an } M'\text{-augmenting path} && \text{(Lemma 2)} \\
&\Rightarrow G/B \text{ has an } (M'/B)\text{-augmenting path} && \text{(proven below)} \\
&\Rightarrow M'/B \text{ is not a maximum matching in } G/B && \text{(Lemma 2)} \\
&\Rightarrow M/B \text{ is not a maximum matching in } G/B && (|M/B| = |M'/B|) \\
&\Rightarrow G/B \text{ has an } (M/B)\text{-augmenting path} && \text{(Lemma 2)}
\end{aligned}$$

The only step that remains to be justified is that the existence of an M' -augmenting path in G implies the existence of an (M'/B) -augmenting path in G/B . Suppose that P is an M' -augmenting path in G . If P does not intersect B then it is already an augmenting path in G/B . Otherwise, since B contains only one free vertex (namely u_r), we know that at least one endpoint of P does not belong to B . Number the vertices of P as v_0, v_1, \dots, v_t with $v_0 \notin B$, and suppose that v_k is the lowest-numbered vertex of P that belongs to B . Then the path $\langle v_0, v_1, \dots, v_{k-1}, v_b \rangle$ is an (M'/B) -augmenting path in G/B , as desired. \square

Lemma 9 inspires the following algorithm for solving the maximum perfect matching problem in non-bipartite graphs.

Algorithm 3 Edmonds' non-bipartite matching algorithm

```
1: Initialize  $M = \emptyset$ .
2: repeat
3:    $P = \text{SEARCH}(G, M)$            // Return augmenting path, or empty set if none exists.
4:    $M \leftarrow M \oplus P$ 
5: until  $P = \emptyset$ 

6: procedure  $\text{SEARCH}(G, M)$ 
7:   Build the graph  $H(G, M)$ .
8:   Search for a path  $\hat{P}$  from  $F$  to  $\Gamma(F)$  in  $H(G, M)$ .
9:   if no path found then
10:    Return  $P = \emptyset$ 
11:  end if
12:  Post-process  $\hat{P}$ , as in the proof of Lemma 8, to extract an augmenting path or flower.
13:  if augmenting path  $P$  is found then
14:    Return  $P$ 
15:  else
16:    Let  $B$  be the blossom of the flower.
17:    Let  $P' = \text{SEARCH}(G/B, M/B)$ 
18:    if  $P' = \emptyset$  then
19:      Return  $P = \emptyset$ 
20:    else
21:      Transform  $P'$  to an  $M$ -augmenting path  $P$  as in the proof of Lemma 9.
22:    end if
23:  end if
24: end procedure
```

The algorithm's outer loop (Lines 2 and 5) iterates at most $n/2$ times. In each iteration, we make a sequence of recursive calls to the `SEARCH` procedure, which searches for an augmenting path. Each recursive call involves shrinking a blossom, which reduces the number of vertices in the graph by at least 2. Hence, we call `SEARCH` at most $n/2$ times within each iteration of the outer loop. To assess the amount of work done in each call to `SEARCH` (excluding work done in recursive sub-calls to the same procedure) we start by observing that the graph $H(G, M)$ has n vertices and at most $2m$ edges, since the number of outgoing edges from a vertex in $H(G, M)$ is bounded above by the degree of that vertex in G . Hence, building and searching the graph $H(G, M)$ takes $O(m)$ time. (Assuming, as always, that G has no isolated vertices so that $n = O(m)$.) The remaining steps of `SEARCH` also take $O(m)$ steps, if not fewer, as can be seen by reviewing the proofs of Lemma 8 and Lemma 9. Hence, the overall running time of Edmonds' algorithm is bounded above by $O(mn^2)$.

3 Bipartite min-cost perfect matching

In the bipartite minimum-cost perfect matching problem, we are given an undirected bipartite graph $G = (L, R, E)$ as before, together with a (non-negative, real-valued) cost c_e for each edge $e \in E$. Let $c(u, v) = c_e$ if $e = (u, v)$ is an edge of G , and $c(u, v) = \infty$ otherwise. As always, let n denote the number of vertices and m the number of edges of G .

3.1 LP relaxation

A perfect matching M can be described by a matrix (x_{uv}) of 0's and 1's, where $x_{uv} = 1$ if and only if $(u, v) \in M$. The sum of the entries in each row and column of this matrix equals 1, since each vertex belongs to exactly one element of M . Conversely, for any matrix with $\{0, 1\}$ -valued entries, if each row sum and column sum is equal to 1, then the corresponding set of edges is a perfect matching. Thus, the bipartite minimum-cost matching problem can be expressed as follows.

$$\begin{aligned} \min \quad & \sum_{u,v} c(u, v)x_{uv} \\ \text{s.t.} \quad & \sum_v x_{uv} = 1 \quad \forall u \\ & \sum_u x_{uv} = 1 \quad \forall v \\ & x_{uv} \in \{0, 1\} \quad \forall u, v \end{aligned}$$

This is a discrete optimization problem because of the constraint that $x_{uv} \in \{0, 1\}$. Although we already know how to solve this discrete optimization problem in polynomial time, many other such problems are not known to have any polynomial-time solution. It's often both interesting and useful to consider what happens when we relax the constraint $x_{uv} \in \{0, 1\}$ to $x_{uv} \geq 0$, allowing the variables to take any non-negative real value. This turns the problem into a continuous optimization problem, in fact a *linear program*.

$$\begin{aligned} \min \quad & \sum_{u,v} c(u, v)x_{uv} \\ \text{s.t.} \quad & \sum_v x_{uv} = 1 \quad \forall u \\ & \sum_u x_{uv} = 1 \quad \forall v \\ & x_{uv} \geq 0 \quad \forall u, v \end{aligned}$$

How should we think about a matrix of values x_{uv} satisfying the constraints of this linear program? We've seen that if the values are integers, then it represents a perfect matching. A general solution of this constraint set can be regarded as a *fractional perfect matching*. What does a fractional perfect matching look like? An example is illustrated in Figure 3. Is it possible that this fractional perfect matching achieves a lower cost than any perfect matching? No, because it can be expressed as a convex combination of perfect matchings (again, see Figure 3) and consequently its cost is the weighted average of the costs of those perfect matchings. In particular, at least one of those perfect matchings costs no more than the fractional perfect matching illustrated on the left side of the figure. This state of affairs is not a coincidence. The *Birkhoff-von Neumann Theorem* asserts that every fractional perfect matching can be decomposed as a convex combination of perfect matchings. (Despite the eminence of its namesakes, the theorem is actually quite easy to prove. You should try finding a proof yourself, if you've never seen one.)

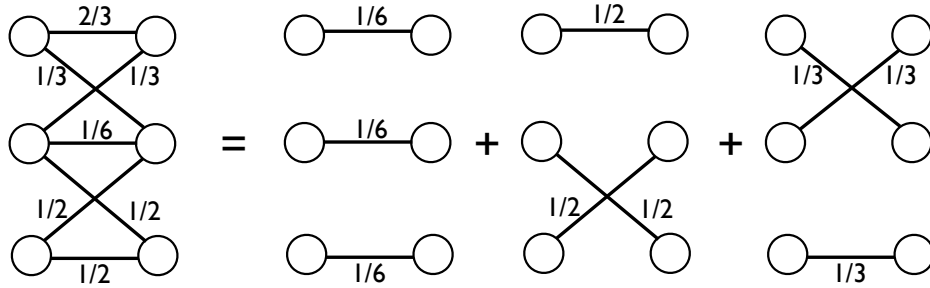


Figure 3: A fractional perfect matching.

Now suppose we have an instance of bipartite minimum-cost perfect matching, and we want to prove a *lower bound* on the optimum: we want to prove that every fractional perfect matching has to cost at least a certain amount. How might we prove this? One way is to run a minimum-cost perfect matching algorithm, look at its output, and declare this to be a lower bound on the cost of any fractional perfect matching. (There exist polynomial-time algorithms for minimum-cost perfect matching, as we will see later in this lecture.) By the Birkhoff-von Neumann Theorem, this produces a valid lower bound, but it's not very satisfying. There's another, much more direct, way to prove lower bounds on the cost of every fractional perfect matching, by directly combining constraints of the linear program. To illustrate this, consider the graph with edge costs as shown in Figure 4. Clearly, the

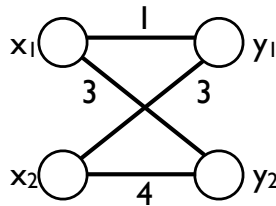


Figure 4: An instance of bipartite minimum cost perfect matching.

minimum cost perfect matching has cost 5. To prove that no fractional perfect matching can cost less than 5, we combine some constraints of the linear program as follows.

$$\begin{aligned} 2x_{11} + 2x_{21} &= 2 \\ -x_{11} - x_{12} &= -1 \\ 4x_{12} + 4x_{22} &= 4 \end{aligned}$$

Adding these constraints, we find that

$$x_{11} + 3x_{12} + 2x_{21} + 4x_{22} = 5 \tag{1}$$

$$x_{11} + 3x_{12} + 3x_{21} + 4x_{22} \geq 5 \tag{2}$$

Inequality (2) is derived from (1) because the only change we made on the left side was to increase the coefficient of x_{21} from 2 to 3, and we know that $x_{21} \geq 0$. The left side of (2)

is the cost of the fractional perfect matching \vec{m} . We may conclude that the cost of every fractional perfect matching is at least 5.

What's the most general form of this technique? For every vertex $w \in L \cup R$, the linear program contains a “degree constraint” asserting that the degree of w in the fractional perfect matching is equal to 1. For each degree constraint, we multiply its left and right sides by some coefficient to obtain

$$\sum_v p_u x_{uv} = p_u$$

for some $u \in L$, or

$$\sum_u q_v x_{uv} = q_v$$

for some $v \in R$. Then we sum all of these equations, obtaining

$$\sum_{u,v} (p_u + q_v) x_{uv} = \sum_u p_u + \sum_v q_v. \quad (3)$$

If the inequality $p_u + q_v \leq c(u, v)$ holds for every $(u, v) \in L \times R$, then in the final step of the proof we (possibly) increase some of the coefficients on the left side of (3) to obtain

$$\sum_{u,v} c(u, v) x_{uv} \geq \sum_u p_u + \sum_v q_v,$$

thus obtaining a lower bound on the cost of every fractional perfect matching. This technique works whenever the coefficients p_u, q_v satisfy $p_u + q_v \leq c(x, y)$ for every edge (x, y) , regardless of whether the values p_u, q_v are positive or negative. To obtain the strongest possible lower bound using this technique, we would set the coefficients p_u, q_v by solving the following linear program.

$$\begin{aligned} \max \quad & \sum_u p_u + \sum_v q_v \\ \text{s.t.} \quad & p_u + q_v \leq c(u, v) \quad \forall u, v \end{aligned}$$

This linear program is called the *dual* of the min-cost-fractional-matching linear program. We've seen that its optimum constitutes a lower bound on the optimum of the min-cost-fractional-matching LP. For any linear program, one can follow the same train of thought to develop a dual linear program. (There's also a formal way of specifying the procedure; it involves taking the transpose of the constraint matrix of the LP.) The dual of a minimization problem is a maximization problem, and its optimum constitutes a lower bound on the optimum of the minimization problem. This fact is called **weak duality**; as you've seen, weak duality is nothing more than an assertion that we can obtain valid inequalities by taking linear combinations of other valid inequalities, and that this sometimes allows us to bound the value of an LP solution from above or below. But actually, the optimum value of an LP is always *exactly equal* to the value of its dual LP! This fact is called **strong duality** (or sometimes simply “duality”), it is far from obvious, and it has important ramifications for algorithm design. In the special case of fractional perfect matching problems, strong duality says that the simple proof technique exemplified above is actually powerful enough to prove the *best possible* lower bound on the cost of fractional perfect matchings, for *every* instance of the bipartite min-cost perfect matching problem.

It turns out that there is a polynomial-time algorithm to solve linear programs. As you can imagine, this fact also has extremely important ramifications for algorithm design, but that's the topic of another lecture.

3.2 Primal-dual algorithm

In this section we will construct a fast algorithm for the bipartite minimum-cost perfect matching algorithm, exploiting insights gained from the preceding section. The basic plan of attack is as follows: we will design an algorithm that simultaneously computes two things: a minimum-cost perfect matching, and a dual solution (vector of p_u and q_v values) whose value (sum of p_u 's and q_v 's) equals the cost of the perfect matching. As the algorithm runs, it maintains the a dual solution \vec{p}, \vec{q} and a matching M , and it preserves the following invariants:

1. Every edge (u, v) satisfies $p_u + q_v \leq c(u, v)$. If $p_u + q_v = c(u, v)$ we say that edge $e = (u, v)$ is *tight*.
2. The elements of M are a subset of the tight edges.
3. The cardinality of M increases by 1 in each phase of the algorithm, until it reaches n .

Assuming the algorithm can maintain these invariants until termination, its correctness will follow automatically. This is because the matching M at termination time will be a perfect matching satisfying

$$\sum_{(u,v) \in M} c(u, v) = \sum_{(u,v) \in M} p_u + q_v = \sum_{u \in L} p_u + \sum_{v \in R} q_v,$$

where the final equation holds because M is a perfect matching. The first invariant of the algorithm implies that \vec{p}, \vec{q} is a feasible dual solution, hence the right side is a lower bound on the cost of any fractional perfect matching. The left side is the cost of the perfect matching M , hence M has the minimum cost of any fractional perfect matching.

So, how do we maintain the three invariants listed above while growing M to be a perfect matching? We initialize $M = \emptyset$ and $\vec{p} = \vec{q} = 0$. Note that the three invariants are trivially satisfied at initialization time. Now, as long as $|M| < n$, we want to find a way to either increase the value of the dual solution or enlarge M without violating any of the invariants. The easiest way to do this is to find an M -augmenting path P consisting of tight edges: in that case, we can update M to $M \oplus P$ without violating any invariants, and we reach the end of a phase. However, sometimes it's not possible to find an M -augmenting path consisting of tight edges: in that case, we must adjust some of the dual variables to make additional edges tight.

The process of adjusting dual variables is best described as follows. The easiest thing would be if we could find a vertex $u \in L$ that doesn't belong to any tight edges. Then we could raise p_u by some amount $\delta > 0$ until an edge containing u became tight. However, maybe every $u \in L$ belongs to a tight edge. In that case, we need to raise p_u by δ while lowering some other q_v by the same amount δ . This is best described in terms of a vertex

set T which will have the property that if one endpoint of an edge $e \in M$ belongs to T , then both endpoints of e belong to T . Whenever T has this property, we can set

$$\delta = \min\{c(u, v) - p_u - q_v \mid u \in L \cap T, v \in R \setminus T\} \quad (4)$$

and adjust the dual variables by setting $p_u \leftarrow p_u + \delta, q_v \leftarrow q_v - \delta$ for all $u \in L \cap T, v \in R \cap T$. This preserves the feasibility of our dual solution \vec{p}, \vec{q} (by the choice of δ) and it preserves the tightness of each edge $e \in M$ because every such edge has either both or neither of its endpoints in T .

Let F be the set of free vertices, i.e. those that don't belong to any element of M . T will be constructed by a sort of breadth-first search along tight edges, starting from the set $L \cap F$ of free vertices in L . We initialize $T = L \cap F$. Since $|M| < n$, T is nonempty. Define δ as in (4); if $\delta > 0$ then adjust dual variables as explained above. Call this a *dual adjustment step*. If $\delta = 0$ then there is at least one tight edge $e = (u, v)$ from $L \cap T$ to $R \setminus T$. If v is a free vertex, then we have discovered an augmenting path P consisting of tight edges (namely, P consists of a path in T that starts at a free vertex in L , walks to u , then crosses edge e to get to v) and we update M to $M \oplus P$ and finish the phase. Call this an *augmentation step*. Finally, if v is not a free vertex then we identify an edge $e = (u', v) \in M$ and we add both v and u' to T and call this a *T -growing step*. Notice that the left endpoint of an edge of M is always added to T at the same time as the right endpoint, which is why T never contains one endpoint of an edge of M unless it contains both.

A phase can contain at most n T -growing steps and at most one augmentation step. Also, there can never be two consecutive dual adjustment steps (since the value of δ drops to zero after the first such step) so the total number of steps in a phase is $O(n)$. Let's figure out the running time of one phase of the algorithm by breaking it down into its component parts.

1. There is only one augmentation step and it costs $O(n)$.
2. There are $O(n)$ T -growing steps and each costs $O(1)$.
3. There are $O(n)$ dual adjustment steps and each costs $O(n)$.
4. Finally, every step starts by computing the value δ using (4). Thus, the value of δ needs to be computed $O(n)$ times. Naïvely it costs $O(m)$ work each time we need to compute δ .

Thus, a naïve implementation of the primal-dual algorithm takes $O(mn^2)$.

However, we can do better using some clever book-keeping combined with efficient data structures. For a vertex $w \in T$, let $s(w)$ denote the number of the step in which w was added to T . Let δ_s denote the value of δ in step s of the phase, and let Δ_s denote the sum $\delta_1 + \dots + \delta_s$. Let $p_{u,s}, q_{v,s}$ denote the values of the dual variables associated to vertices u, v

at the end of step s . Note that

$$p_{u,s} = \begin{cases} p_{u,0} + \Delta_s - \Delta_{s(u)} & \text{if } u \in L \cap T \\ p_{u,0} & \text{if } u \in L \setminus T \end{cases} \quad (5)$$

$$q_{v,s} = \begin{cases} q_{v,0} - \Delta_s + \Delta_{s(v)} & \text{if } v \in R \cap T \\ q_{v,0} & \text{if } v \in R \setminus T \end{cases} \quad (6)$$

Consequently, if $e = (u, v)$ is any edge from $L \cap T$ to $R \setminus T$ at the end of step s , then

$$c(u, v) - p_{u,s} - q_{v,s} = c(u, v) - p_{u,0} - \Delta_s + \Delta_{s(u)} - q_{v,0}$$

The only term on the right side that depends on s is $-\Delta_s$, which is a global value that is common to all edges. Thus, choosing the edge that minimizes $c(u, v) - p_{u,s} - q_{v,s}$ is equivalent to choosing the edge that minimizes $c(u, v) - p_{u,0} + \Delta_{s(u)} - q_{v,0}$. Let us maintain a priority queue containing all the edges from $L \cap T$ to $R \setminus T$. An edge $e = (u, v)$ is inserted into this priority queue at the time its left endpoint u is inserted into T . The value associated to e in the priority queue is $c(u, v) - p_{u,0} + \Delta_{s(u)} - q_{v,0}$, and this value never changes as the phase proceeds. Whenever the algorithm needs to choose the edge that minimizes $c(u, v) - p_{u,s} - q_{v,s}$, it simply extracts the minimum element of this priority queue, repeating as necessary until it finds an edge whose right endpoint does not belong to T . The total amount of work expended on maintaining the priority queue throughout a phase is $O(m \log n)$.

Finally, our gimmick with the priority queue eliminates the need to actually update the values p_u, q_v during a dual adjustment step. These values are only needed for computing the value of δ_s , and for updating the dual solution at the end of the phase. However, if we store the values $s(u), s(v)$ for all u, v as well as the values Δ_s for all s , then one can compute any specific value of $p_{u,s}$ or $q_{v,s}$ in constant time using (5)-(6). In particular, it takes $O(n)$ time to compute all the values p_u, q_v at the end of the phase, and it only takes $O(1)$ time to compute the value $\delta_s = c(u, v) - p_u - q_v$ once we have identified the edge $e = (u, v)$ using the priority queue. Thus, all the work to maintain the values p_u, q_v amounts to only $O(n)$ per phase.

In total, the amount of work in any phase is bounded by $O(m \log n)$ and consequently the algorithm's running time is $O(mn \log n)$.