## Algorithmic Paradigms

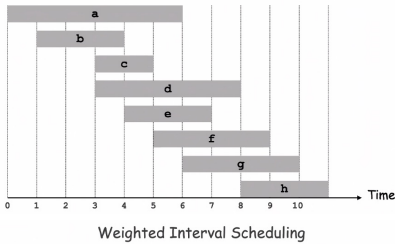**Greedy.** Build up a solution incrementally, myopically optimizing some local criterion.

**Divide-and-conquer.** Break up a problem into sub-problems, solve each sub-problem independently, and combine solution to sub-problems to form solution to original problem.

**Dynamic programming.** Break up a problem into a series of overlapping sub-problems, and build up solutions to larger and larger sub-problems.

## Dynamic Programming Applications

**Areas.**
- Bioinformatics.
- Control theory.
- Information theory.
- Operations research.
- Computer science: theory, graphics, AI, compilers, systems, ....

**Some famous dynamic programming algorithms.**
- Unix diff for comparing two files.
- Viterbi for hidden Markov models.
- Smith-Waterman for genetic sequence alignment.
- Bellman-Ford for shortest path routing in networks.
- Cocke-Kasami-Younger for parsing context free grammars.

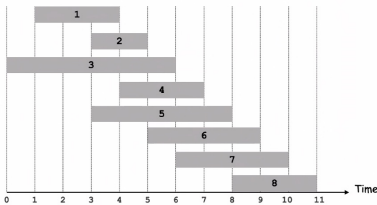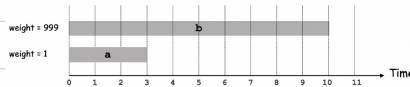## Weighted Interval Scheduling

**Weighted interval scheduling problem.**
- Job j starts at $s_j$, finishes at $f_j$, and has weight or value $v_j$.
- Two jobs compatible if they don't overlap.
- Goal: find maximum weight subset of mutually compatible jobs.



## Unweighted Interval Scheduling Review

**Recall.** Greedy algorithm works if all weights are 1.
- Consider jobs in ascending order of finish time.
- Add job to subset if it is compatible with previously chosen jobs.

**Observation.** Greedy algorithm can fail spectacularly if arbitrary weights are allowed.



## Weighted Interval Scheduling

**Notation.** Label jobs by finishing time: $f_1 \le f_2 \le \ldots \le f_n$.
**Def.** p(j) = largest index i < j such that job i is compatible with j.

**Ex:** p(8) = 5, p(7) = 3, p(2) = 0.



## Dynamic Programming: Binary Choice

**Notation.** OPT(j) = value of optimal solution to the problem consisting of job requests 1, 2, ..., j.

- Case 1: OPT selects job j.
  - collect profit $v_j$
  - can't use incompatible jobs { p(j) + 1, p(j) + 2, ..., j - 1 }
  - must include optimal solution to problem consisting of remaining compatible jobs 1, 2, ..., p(j)    *optimal substructure*

- Case 2: OPT does not select job j.
  - must include optimal solution to problem consisting of remaining compatible jobs 1, 2, ..., j-1

## Weighted Interval Scheduling: Brute Force

Brute force algorithm.

```
Input: n, s₁,…,sₙ, f₁,…,fₙ, v₁,…,vₙ

Sort jobs by finish times so that f₁ ≤ f₂ ≤ ... ≤ fₙ.

Compute p(1), p(2), …, p(n)

Compute-Opt(j) {
    if (j = 0)
        return 0
    else
        return max(vⱼ + Compute-Opt(p(j)), Compute-Opt(j-1))
}
```

## Weighted Interval Scheduling: Memoization

Memoization. Store results of each sub-problem in a cache; lookup as needed.

```
Input: n, s₁,…,sₙ, f₁,…,fₙ, v₁,…,vₙ

Sort jobs by finish times so that f₁ ≤ f₂ ≤ ... ≤ fₙ.
Compute p(1), p(2), …, p(n)

for j = 1 to n
    M[j] = empty          ← global array
M[0] = 0

M-Compute-Opt(j) {
    if (M[j] is empty)
        M[j] = max(vⱼ + M-Compute-Opt(p(j)), M-Compute-Opt(j-1))
    return M[j]
}
```

## Weighted Interval Scheduling: Running Time

Claim. Memoized version of algorithm takes $O(n \log n)$ time.
- Sort by finish time: $O(n \log n)$.
- Computing $p(\cdot)$: $O(n \log n)$ via sorting by start time.

- `M-Compute-Opt(j)`: each invocation takes $O(1)$ time and either
  - (i) returns an existing value `M[j]`
  - (ii) fills in one new entry `M[j]` and makes two recursive calls

- Progress measure $\Phi$ = # nonempty entries of `M[]`.
  - initially $\Phi = 0$, throughout $\Phi \le n$.
  - (ii) increases $\Phi$ by 1 $\Rightarrow$ at most $2n$ recursive calls.

- Overall running time of `M-Compute-Opt(n)` is $O(n)$. ∎

Remark. $O(n)$ if jobs are pre-sorted by start and finish times.

## Weighted Interval Scheduling: Finding a Solution

Q. Dynamic programming algorithms computes optimal value. What if we want the solution itself?
A. Do some post-processing.

```
Run M-Compute-Opt(n)
Run Find-Solution(n)

Find-Solution(j) {
    if (j = 0)
        output nothing
    else if (vⱼ + M[p(j)] > M[j-1])
        print j
        Find-Solution(p(j))
    else
        Find-Solution(j-1)
}
```

- # of recursive calls $\le n \Rightarrow O(n)$.
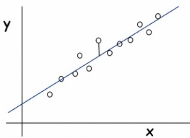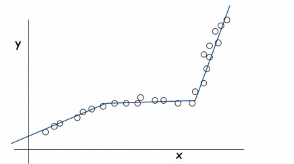
## Segmented Least Squares

Least squares.
- Foundational problem in statistic and numerical analysis.
- Given n points in the plane: $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$.
- Find a line $y = ax + b$ that minimizes the sum of the squared error:



Solution. Calculus $\Rightarrow$ min error is achieved when

## Segmented Least Squares

Segmented least squares.
- Points lie roughly on a sequence of several line segments.
- Given n points in the plane $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$ with
- $x_1 < x_2 < \ldots < x_n$, find a sequence of lines that minimizes $f(x)$.

Q. What's a reasonable choice for $f(x)$ to balance accuracy and parsimony?
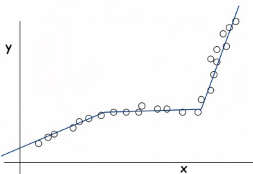        ↑                                    ↑
  number of lines                      goodness of fit



## Segmented Least Squares

Segmented least squares.
- Points lie roughly on a sequence of several line segments.
- Given n points in the plane $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$ with
- $x_1 < x_2 < \ldots < x_n$, find a sequence of lines that minimizes:
  - the sum of the sums of the squared errors E in each segment
  - the number of lines L
- Tradeoff function: $E + c\,L$, for some constant $c > 0$.



## Dynamic Programming: Multiway Choice

Notation.
- $OPT(j)$ = minimum cost for points $p_1, p_{i+1}, \ldots, p_j$.
- $e(i, j)$ = minimum sum of squares for points $p_i, p_{i+1}, \ldots, p_j$.

To compute $OPT(j)$:
- Last segment uses points $p_i, p_{i+1}, \ldots, p_j$ for some i.
- Cost = $e(i, j) + c + OPT(i-1)$.

## Segmented Least Squares:  Algorithm

```
INPUT: n, p₁,...,pₙ , c

Segmented-Least-Squares() {
   M[0] = 0
   for j = 1 to n
       for i = 1 to j
           compute the least square error eᵢⱼ for
           the segment pᵢ,..., pⱼ

   for j = 1 to n
       M[j] = min 1 ≤ i ≤ j (eᵢⱼ + c + M[i-1])

   return M[n]
}
```

**Running time.** $O(n^3)$.  ← can be improved to $O(n^2)$ by pre-computing various statistics

- Bottleneck = computing $e(i, j)$ for $O(n^2)$ pairs, $O(n)$ per pair using previous formula.

## Dynamic Programming:  False Start

**Def.**  OPT(i) = max profit subset of items 1, ..., i.

- Case 1:  OPT does not select item i.
  - OPT selects best of { 1, 2, ..., i-1 }

- Case 2:  OPT selects item i.
  - accepting item i does not immediately imply that we will have to reject other items
  - without knowing what other items were selected before i, we don't even know if we have enough room for i

**Conclusion.**  Need more sub-problems!

## Knapsack Problem:  Bottom-Up

**Knapsack.**  Fill up an n-by-W array.

```
Input: n, W, w₁,...,wₙ, v₁,...,vₙ

for w = 0 to W
    M[0, w] = 0

for i = 1 to n
    for w = 1 to W
        if (wᵢ > w)
            M[i, w] = M[i-1, w]
        else
            M[i, w] = max {M[i-1, w], vᵢ + M[i-1, w-wᵢ]}

return M[n, W]
```

## Knapsack Problem:  Running Time

**Running time.**  $\Theta(n\,W)$.
- Not polynomial in input size!
- "Pseudo-polynomial."
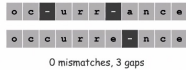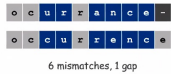- Decision version of Knapsack is NP-complete.  [Chapter 8]

**Knapsack approximation algorithm.**  There exists a poly-time algorithm that produces a feasible solution that has value within 0.01% of optimum. [Section 11.8]

## Knapsack Problem

**Knapsack problem.**
- Given n objects and a "knapsack."
- Item i weighs $w_i > 0$ kilograms and has value $v_i > 0$.
- Knapsack has capacity of W kilograms.
- Goal:  fill knapsack so as to maximize total value.

**Ex:** { 3, 4 } has value 40.

W = 11

| # | value | weight |
|---|-------|--------|
| 1 | 1 | 1 |
| 2 | 6 | 2 |
| 3 | 18 | 5 |
| 4 | 22 | 6 |
| 5 | 28 | 7 |

**Greedy:**  repeatedly add item with maximum ratio $v_i / w_i$.
**Ex:** { 5, 2, 1 } achieves only value = 35  $\Rightarrow$  greedy not optimal.

## Dynamic Programming:  Adding a New Variable

**Def.**  OPT(i, w) = max profit subset of items 1, ..., i with weight limit w.

- Case 1:  OPT does not select item i.
  - OPT selects best of { 1, 2, ..., i-1 } using weight limit w

- Case 2:  OPT selects item i.
  - new weight limit = w − wᵢ
  - OPT selects best of { 1, 2, ..., i-1 } using this new weight limit

## Knapsack Algorithm

W + 1

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|----|----|
| φ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| { 1 } | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| { 1, 2 } | 0 | 1 | 6 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| { 1, 2, 3 } | 0 | 1 | 6 | 7 | 7 | 18 | 19 | 24 | 25 | 25 | 25 | 25 |
| { 1, 2, 3, 4 } | 0 | 1 | 6 | 7 | 7 | 18 | 22 | 24 | 28 | 29 | 29 | 40 |
| { 1, 2, 3, 4, 5 } | 0 | 1 | 6 | 7 | 7 | 18 | 22 | 28 | 29 | 34 | 34 | 40 |

n + 1

OPT: { 4, 3 }
value = 22 + 18 = 40

W = 11

| Item | Value | Weight |
|------|-------|--------|
| 1 | 1 | 1 |
| 2 | 6 | 2 |
| 3 | 18 | 5 |
| 4 | 22 | 6 |
| 5 | 28 | 7 |

## Dynamic Programming Summary

**Recipe.**
- Characterize structure of problem.
- Recursively define value of optimal solution.
- Compute value of optimal solution.
- Construct optimal solution from computed information.

**Dynamic programming techniques.**
- Binary choice:  weighted interval scheduling.
- Multi-way choice:  segmented least squares. ← Viterbi algorithm for HMM also uses DP to optimize a maximum likelihood tradeoff between parsimony and accuracy
- Adding a new variable:  knapsack.
- Dynamic programming over intervals:  RNA secondary structure. ← CKY parsing algorithm for context-free grammar has similar structure

**Top-down vs. bottom-up:**  different people have different intuitions.

# Sequence Alignment 6.5

## String Similarity

**How similar are two strings?**
- ocurrance
- occurrence


6 mismatches, 1 gap


1 mismatch, 1 gap


0 mismatches, 3 gaps

## Edit Distance

**Applications.**
- Basis for Unix diff.
- Speech recognition.
- Computational biology.

**Edit distance.** [Levenshtein 1966, Needleman-Wunsch 1970]
- Gap penalty $\delta$; mismatch penalty $\alpha_{pq}$.
- Cost = sum of gap and mismatch penalties.


$\alpha_{TC} + \alpha_{GT} + \alpha_{AG} + 2\alpha_{CA}$


$2\delta + \alpha_{CA}$

## Sequence Alignment

**Goal:** Given two strings $X = x_1 x_2 \ldots x_m$ and $Y = y_1 y_2 \ldots y_n$ find alignment of minimum cost.
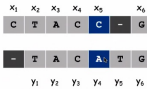
**Def.** An *alignment* M is a set of ordered pairs $x_i$-$y_j$ such that each item occurs in at most one pair and no crossings.

**Def.** The pair $x_i$-$y_j$ and $x_{i'}$-$y_{j'}$ *cross* if $i < i'$, but $j > j'$.

$$cost(M) = \underbrace{\sum \alpha_{x_i y_j}}_{\text{mismatch}} + \underbrace{\sum \delta + \sum \delta}_{\text{gap}}$$

**Ex:** CTACCG vs. TACATG.
**Sol:** $M = x_2$-$y_1$, $x_3$-$y_2$, $x_4$-$y_3$, $x_5$-$y_4$, $x_6$-$y_6$.



## Sequence Alignment: Problem Structure

**Def.** OPT(i, j) = min cost of aligning strings $x_1 x_2 \ldots x_i$ and $y_1 y_2 \ldots y_j$.
- Case 1: OPT matches $x_i$-$y_j$.
  - pay mismatch for $x_i$-$y_j$ + min cost of aligning two strings $x_1 x_2 \ldots x_{i-1}$ and $y_1 y_2 \ldots y_{j-1}$
- Case 2a: OPT leaves $x_i$ unmatched.
  - pay gap for $x_i$ and min cost of aligning $x_1 x_2 \ldots x_{i-1}$ and $y_1 y_2 \ldots y_j$
- Case 2b: OPT leaves $y_j$ unmatched.
  - pay gap for $y_j$ and min cost of aligning $x_1 x_2 \ldots x_i$ and $y_1 y_2 \ldots y_{j-1}$

$$OPT(i,j) = \begin{cases} j\delta & \text{if } i = 0 \\ \min \begin{cases} \alpha_{x_i y_j} + OPT(i-1, j-1) \\ \delta + OPT(i-1, j) \\ \delta + OPT(i, j-1) \end{cases} & \text{otherwise} \\ i\delta & \text{if } j = 0 \end{cases}$$
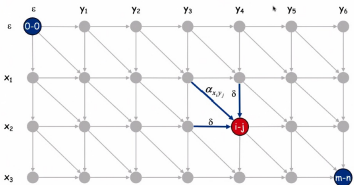
## Sequence Alignment: Algorithm

```
Sequence-Alignment(m, n, x₁x₂...xₘ, y₁y₂...yₙ, δ, α) {
    for i = 0 to m
        M[0, i] = iδ
    for j = 0 to n
        M[j, 0] = jδ

    for i = 1 to m
        for j = 1 to n
            M[i, j] = min(α[xᵢ, yⱼ] + M[i-1, j-1],
                          δ + M[i-1, j],
                          δ + M[i, j-1])
    return M[m, n]
}
```

**Analysis.** $\Theta(mn)$ time and space.
**English words or sentences:** $m, n \leq 10$.
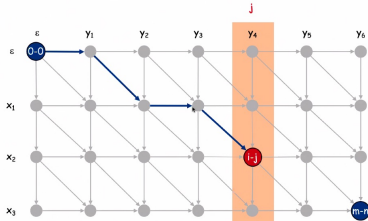**Computational biology:** $m = n = 100{,}000$. 10 billions ops OK, but 10GB array?

## Sequence Alignment: Linear Space

**Q.** Can we avoid using quadratic **space**?

**Easy.** Optimal **value** in O(m + n) space and O(mn) time.
- Compute OPT(i, $\cdot$) from OPT(i-1, $\cdot$).
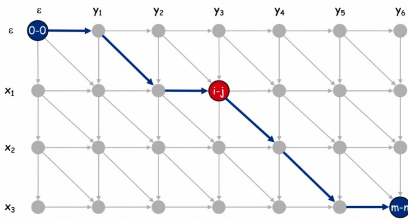- No longer a simple way to recover alignment itself.

**Theorem.** [Hirschberg 1975] Optimal **alignment** in O(m + n) space and O(mn) time.
- Clever combination of divide-and-conquer and dynamic programming.
- Inspired by idea of Savitch from complexity theory.

### Sequence Alignment: Linear Space

**Edit distance graph.**
- Let f(i, j) be shortest path from (0,0) to (i, j).
- Observation: f(i, j) = OPT(i, j).



### Sequence Alignment: Linear Space

**Edit distance graph.**
- Let f(i, j) be shortest path from (0,0) to (i, j).
- Can compute f ($\cdot$, j) for any j in O(mn) time and O(m + n) space.

## Sequence Alignment: Linear Space

**Observation 1.** The cost of the shortest path that uses (i, j) is
f(i, j) + g(i, j).



7, 7, 3

## Sequence Alignment: Running Time Analysis

**Theorem.** Let T(m, n) = max running time of algorithm on strings of length m and n. T(m, n) = O(mn).

**Pf.** (by induction on n)
- O(mn) time to compute f( · , n/2) and g ( · , n/2) and find index q.
- T(q, n/2) + T(m - q, n/2) time for two recursive calls.
- Choose constant c so that:

$$
\begin{aligned}
T(m,\ 2) &\leq cm \\
T(2,\ n) &\leq cn \\
T(m,\ n) &\leq cmn + T(q,\ n/2) + T(m-q,\ n/2)
\end{aligned}
$$

- Base cases: m = 2 or n = 2.
- Inductive hypothesis: T(m, n) ≤ 2cmn.

$$
\begin{aligned}
T(m,n) &\leq T(q,n/2) + T(m-q,n/2) + cmn \\
&\leq 2cqn/2 + 2c(m-q)n/2 + cmn \\
&= cqn + cmn - cqn + cmn \\
&= 2cmn
\end{aligned}
$$