

Algoritmos en Bioinformática

Resumen

En bioinformática se necesita aplicar conocimientos para resolver problemas en nuevos contextos. Además, se necesita capacidad de elaborar proyectos de investigación o aplicaciones en bioinformática, incorporando soluciones innovadoras, anticipando dificultades y valorando estrategias alternativas de contingencia, así como consideraciones en cuanto a responsabilidad social, ética y legal. Para ello, es imprescindible conocer y manejar los principales métodos de algoritmia y su aplicación en bioinformática. Concretamente, nos centraremos en las estructuras de datos, la notación O y órdenes de ejecución, la búsqueda y ordenación, la programación dinámica y aplicaciones algorítmicas en bioinformática para la búsqueda de perfiles y alineamientos.

Índice general

I	Introducción a los algoritmos y estructuras de datos	2
I.1	Algoritmos y estructura de datos	2
I.1.1	Algoritmos	2
I.1.2	Ejemplo: Algoritmo de Euclid	2
I.1.3	Estructura de datos	3
I.2	Diseño de algoritmos	4
I.2.1	El problema del cambio - algoritmo codicioso	4
I.2.2	Las torres de Hanoi - algoritmo recursivo	5
I.3	Eficiencia de algoritmos	6
I.3.1	Estimar tiempos de ejecución	6
I.3.2	Multiplicación de matrices	7
I.3.3	Búsqueda lineal	7

Capítulo I

Introducción a los algoritmos y estructuras de datos

I.1. Algoritmos y estructura de datos

Un programa es el resultado de la **ecuación de Wirth**, es decir, la suma de algoritmos y estructura de datos.

I.1.1. Algoritmos

Los algoritmos tienen muchas definiciones, pero ninguna es muy precisa. Wikipedia define los algoritmos como "un conjunto de reglas que definen con precisión una secuencia de operaciones para realizar alguna tarea y que finalmente se detienen". Normalmente, están escritos en **pseudocódigo**, algo intermedio entre lenguaje natural y código de ordenador. Los tres bloques principales de un algoritmo son:

- **Bloque secuencial:** bloques de sentencias (ordinarias) que se ejecutan secuencialmente en su totalidad. Las sentencias pueden tener sólo cálculos directos o varias llamadas a funciones. El orden de ejecución es según la ley de la gravedad. En Python, se definen como bloques formados por sentencias con la misma sangría.
- **Selecciones:** sentencias en las que la ejecución se bifurca a diferentes bloques según alguna condición. En Python se reconoce en bloques if, elif y else.
- **Repeticiones o loops:** un bloque de sentencias se repite mientras se cumpla alguna condición. Puede haber un bucle for para un cierto número de repeticiones o un bucle while si hay una condición.

I.1.2. Ejemplo: Algoritmo de Euclid

El algoritmo de Euclid calcula el máximo común divisor de dos números positivos a , b calculando repetidamente $r = a \% b$ y sustituyendo a por b y b por r mientras $r > 0$ (siendo r el resto). En Python:

```
def euclid_gcd(a, b):
    while b > 0:
        r = a % b
        a = b
        b = r
        #Alternativa pitónica: a, b = b, a % b

    return a
```

La ecuación de Wirth es bonita y correcta en general, pero los programas también necesitan bastante **manejo de excepciones**, es decir, detectar situaciones excepcionales y decir al programa qué hacer cuando se producen. Los errores de argumentos son fáciles de prevenir y de manejar. Las excepciones de ejecución, cosas que van mal durante la ejecución, son más difíciles de detectar y prevenir. La programación debe ser muy **defensiva**.

I.1.3. Estructura de datos

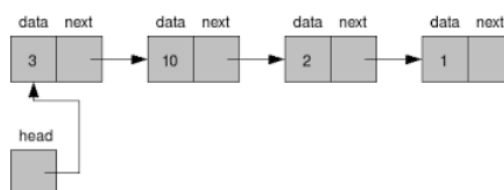
Los algoritmos trabajan con datos. Variables individuales están bien para algoritmos simples. Las estructuras de datos son formas de organizar datos complejos para algoritmos avanzados. Las estructuras de datos más simples son strings, listas y arrays, las avanzadas son diccionarios y sets, y las más avanzadas listas enlazadas, árboles y grafos, aunque estos últimos están disponibles por la importación de módulos.

I.1.3.1. Strings, listas, arrays y diccionarios

Los elementos de strings, listas y arrays son accesibles mediante índices, mientras que los diccionarios están compuestos por parejas clave:valor. Todos son objetos de Python con atributos, variables con información del objeto, y métodos, funciones que actúan sobre el contenido del objeto. `dir(objeto)` lista todos los atributos y métodos del objeto.

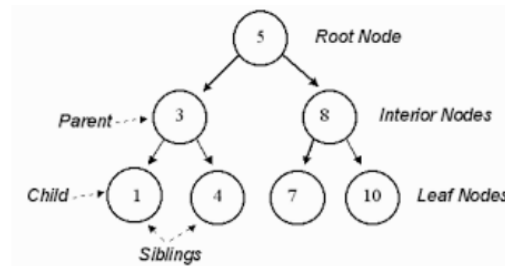
I.1.3.2. Listas enlazadas

Las listas enlazadas están compuestas por nodos con campos `data` que contienen la información del nodo y `next` que apunta al siguiente nodo. Son una versión dinámica de los arrays, y son útiles cuando el número de nodos y/o su localización no se conoce previamente.



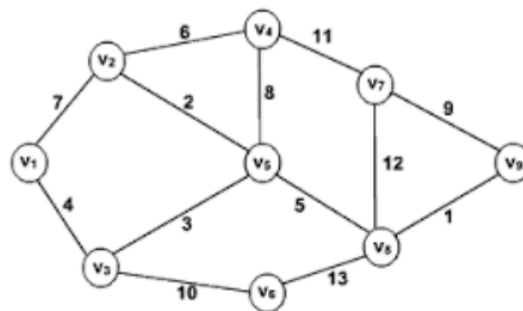
I.1.3.3. Árboles

Los árboles contienen nodos de datos organizados de forma jerárquica con un único nodo raíz y los demás tienen un padre y quizás hijos.



I.1.3.4. Grafos

Los grafos están compuestos por nodos o vértices conectados por edges. Posiblemente es la estructura de datos más general: puede representar mapas de carreteras, redes sociales, interacciones de proteínas, etc.



I.2. Diseño de algoritmos

La escritura de algoritmos (y la programación en general) suele hacerse ad hoc. Es un acto creativo: debe seguir las reglas de programación pero también requiere imaginación, creatividad y experiencia. Lo mismo ocurre con la escritura ordinaria, ya que no podemos llenar una página vacía sólo con reglas gramaticales. La programación también requiere trabajo duro, mucha práctica y, además, bastante lectura de algoritmos. A veces podemos aprovechar las técnicas generales de diseño derivadas de una larga experiencia en resolución de problemas y análisis de algoritmos. No pueden aplicarse como reglas empíricas automáticas, pero pueden tener un amplio rango de aplicabilidad. Consideraremos tres: **algoritmos codiciosos**, **algoritmos de divide y vencerás** (también conocidos como **recursivos**) y **programación dinámica**.

I.2.1. El problema del cambio - algoritmo codicioso

Supongamos que tenemos trabajo como cajero y nuestros clientes quieren cambio en el menor número de monedas posible. ¿Cómo podemos proceder? La idea más

sencilla: dar a cada paso la moneda más grande y más pequeña que la cantidad que queda por cambiar. Ejemplo: ¿cómo dar cambio de 3,44 euros? Fácil: una moneda de 2 euros, una moneda de 1 euro, dos monedas de 20 céntimos, dos monedas de 2 céntimos. Hay que escribir el algoritmo pero la idea general es codiciosa: Intentamos minimizar **globalmente** el número total de monedas, pero lo hacemos **localmente** usando en cada paso la moneda más grande posible para minimizar la cantidad que queda por cambiar.

Suponiendo que trabajamos con monedas/billetes de 1, 2, 5, 10, 20, 50, 100 y 200, queremos guardar el número de monedas/billetes de cada tipo a devolver en un diccionario:

```
def change(c):
    assert c >= 0, "change for positive amounts only"
    l_coin_values = [1, 2, 5, 10, 20, 50, 100, 200]
    d_change = {}

    for coin in sorted(l_coin_values) [::-1]:
        d_change[coin] = c//coin
        c = c%coin

    return d_change
```

Aparentemente, esto funciona. Pero si debemos dar un cambio de 7 con monedas 1, 3, 4 o 5, la respuesta más eficiente solo requiere dos monedas, una de 4 y una de 3, pero este algoritmo cogería una de 5 y tendrá que dar dos monedas de 1. Esto ocurre bastante con algoritmos codiciosos: son muy naturales, pero pueden dar una respuesta equivocada. La forma de resolver esto sería con programación dinámica.

1.2.2. Las torres de Hanoi - algoritmo recursivo

Se nos da un conjunto de 64 discos de oro de diferentes tamaños apilados en la pila A en tamaños crecientes, y otras dos pilas vacías B, C. Queremos mover la primera pila a B un disco cada vez usando C como clavija auxiliar obedeciendo la regla de que ningún disco puede colocarse encima de otro disco más pequeño. Esto es fácil para 2 discos, no muy difícil para 3, pero para 4 la dificultad aumenta.

Se puede obtener una solución recursiva sencilla para N discos. Primero se mueven los primeros N-1 discos de la pila A a la C utilizando B como pila auxiliar. El disco restante se mueve de A a B. Los N-1 discos restantes se mueven de C a B usando A como pila auxiliar.

```
def hanoi(n_disks, a=1, b=2, c=3):
    assert n_disks > 0, "n_disks at least 1"

    if n_disks == 1:
        print("Move disk from %d to %d" % (a,b))
    else:
        hanoi(n_disks - 1, a, c, b)
        print("move disk from %d to %d" % (a,b))
```

```
hanoi(n_disks - 1, c, b, a)
```

Con esto, hay que tener cuidado con los tiempos de ejecución incluso para `n_disks` pequeños. De hecho, el problema general de Hanoi es extremadamente costoso incluso para un número moderado de discos.

Los algoritmos recursivos suelen derivar de una estrategia de «divide y vencerás»: Dividir un problema P en M subproblemas P_m , resolverlos por separado obteniendo soluciones S_m y combinar estas soluciones en una solución S de P .

En el caso de las torres de Hanoi se pueden dividir dos subproblemas: P_1 es el subproblema de mover $N - 1$ discos de A a C usando B , y P_2 el subproblema de mover $N - 1$ discos de C a B usando A . Se pueden combinar los movimientos según el código de Python. Los algoritmos son eficientes si los subproblemas son sustancialmente más pequeños - pero esto no es el caso de Hanoi.

1.3. Eficiencia de algoritmos

En primer lugar, los algoritmos deben ser correctos, ya que un algoritmo rápido, pero erróneo, es inútil. También es deseable que no requieran (mucho) memoria extra. La función `hanoi` cumple esto: sólo se usan sus parámetros. Algo a tener en cuenta en bioinformática, ya que los datos pueden ser muy grandes, es que también es muy deseable que los algoritmos sean lo más rápidos posible. Pero un algoritmo debe leer sus entradas, y si hay muchas y grandes, esto ralentizará el algoritmo. No obstante, los tiempos de ejecución deseables no deberían estar muy por encima del **mismo orden de magnitud que el tamaño de sus entradas**.

1.3.1. Estimar tiempos de ejecución

En primer lugar, no se miden solo los tiempos reales, ya que dependen del lenguaje, la máquina, el programador y, por supuesto, las entradas. Por tanto, dependen demasiado del contexto para permitir generalizaciones significativas. En su lugar, hay que centrarse en **tiempos abstractos** medidos contando las **operaciones clave** que el algoritmo realiza en una entrada dada. Para los algoritmos iterativos, normalmente se busca la operación clave en el bucle más interno. Contando cuántas veces se realizan estas operaciones clave se obtiene una buena estimación del tiempo que tardarán los algoritmos. De esta forma, el coste del algoritmo de cambio viene dado por la longitud de la lista de monedas.

El análisis de algoritmos recursivos es (mucho) más difícil. Para Hanoi, la operación clave puede ser `print("move disk from %d to %d" % (a, b))`, pero aunque aparece explícitamente en el código, también tiene lugar dentro de las llamadas recursivas. Esto da lugar a estimaciones recurrentes del coste de los algoritmos recursivos que a menudo son difíciles de escribir y resolver. Se pueden desarrollar algunas estrategias generales en algoritmos mucho más sencillos basados en bucles.

I.3.2. Multiplicación de matrices

Un algoritmo muy conocido y relativamente costoso es $c_{i,j} = \sum_{k=1}^n a_{i,k} b_{k,j}$. Un código de Python simple y malo que describe esto es el siguiente:

```
def matrix_multiplication(m_1, m_2):
    n_rows, n_interm, n_columns = m_1.shape[0], m_2.shape[0], m_2.shape[1]
    m_product = np.zeros( (n_rows, n_columns) )

    for p in range(n_rows):
        for q in range(n_columns):
            for r in range(n_interm):
                m_product[p, q] += m_1[p, r] * m_2[r, q]

    return m_product
```

Aquí, la operación clave es $m_1[p, r] * m_2[r, q]$. Asumiendo matrices cuadradas con N filas y columnas, esta operación clave se repite $N \times N \times N = N^3$, lo cual es sustancialmente más grande que el tamaño del problema $N^2 + N^2 = 2N^2$.

I.3.3. Búsqueda lineal