

Programacion Genetica

PG a Vuelo de Pajaro

- Desarrollado: EEUU en los 90s
- Pioneros: J. Koza pero...
- Generalmente aplicado a:
 - prediccion, clasificacion
- Propiedades generales:
 - compite con NN y similares
 - necesita poblaciones bastante grandes (miles?)
 - lenta
- Propiedades particulares:
 - cromosomas no lineales: arboles, grafos
 - mutacion posible pero no necesaria (disputado!)

Tableau Tecnico

Representacion	Arboles
Recombinacion	Intercambio de sub-arboles
Mutacion	Modificaciones al azar en los (sub)arboles
Seleccion de padres	Proporcional al fitness
Seleccion de sobrevivientes	Reemplazamiento generacional

Ejemplo: evaluación de créditos

- Los bancos necesitan distinguir entre buenos aplicantes (a créditos) y malos
- Se necesita un modelo que se ajuste a datos históricos

ID	Nro de hijos	Salario	E.civil	OK?
ID-1	2	45000	Casado	0
ID-2	0	30000	Soltero	1
ID-3	1	40000	Divorciado	1
...				

- Un modelo posible:

IF (NDH = 2) AND (S > 80000) THEN bueno ELSE malo

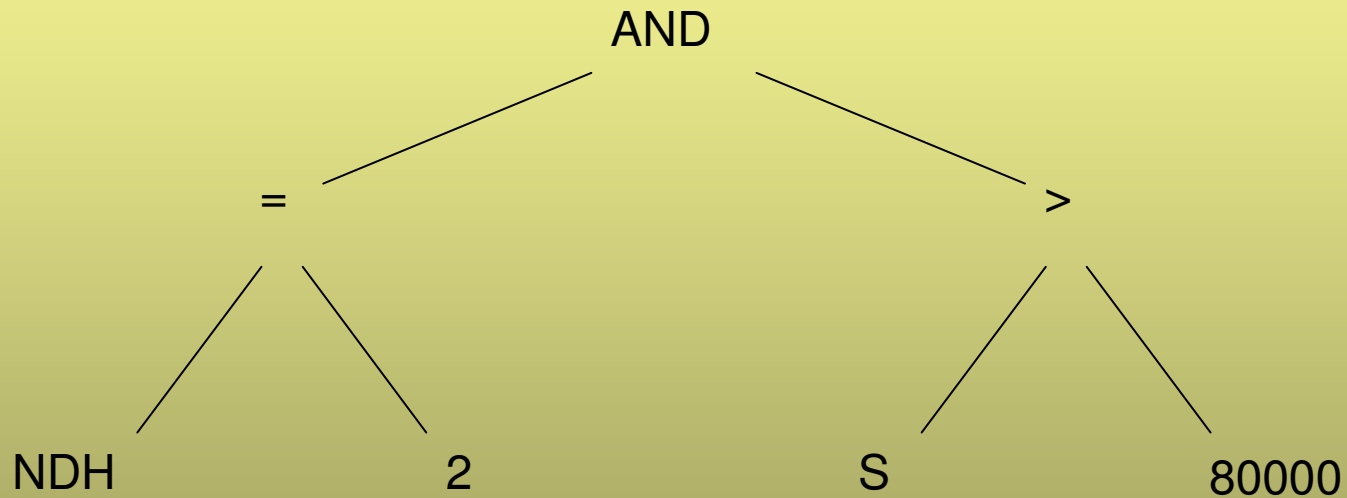
- En general:

IF formula THEN bueno ELSE malo

- No se conoce la estructura de la formula, entonces
- El espacio de busqueda es el espacio de formulas (fenotipo)
- Fitness de una formula designoide: porcentaje de casos bien clasificados
- La representacion natural de las formulas (genotipo) es un “parse tree” (arbol de clasificacion)

IF (NDH = 2) AND (S > 80000) THEN bueno ELSE malo

se representa con:



Representacion basada en Arboles

Representacion muy general, Ejemplo:

- Formula aritmetica

$$2 \cdot \pi + \left((x + 3) - \frac{y}{5 + 1} \right)$$

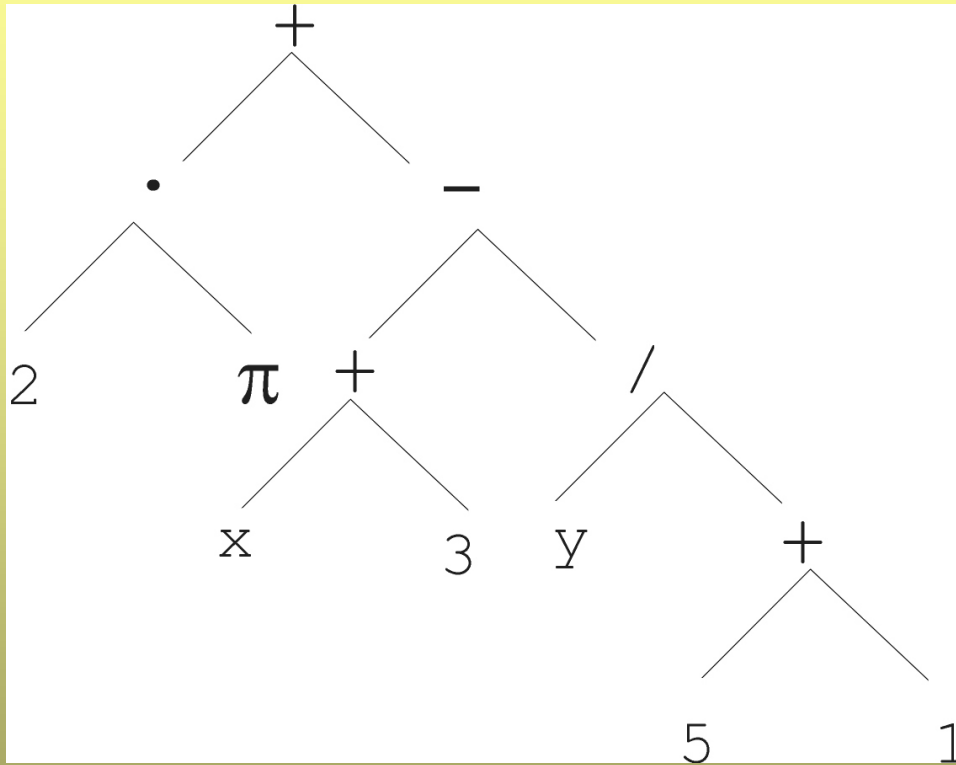
- Formula logica

$$(x \wedge \text{true}) \rightarrow ((x \vee y) \vee (z \leftrightarrow (x \wedge y)))$$

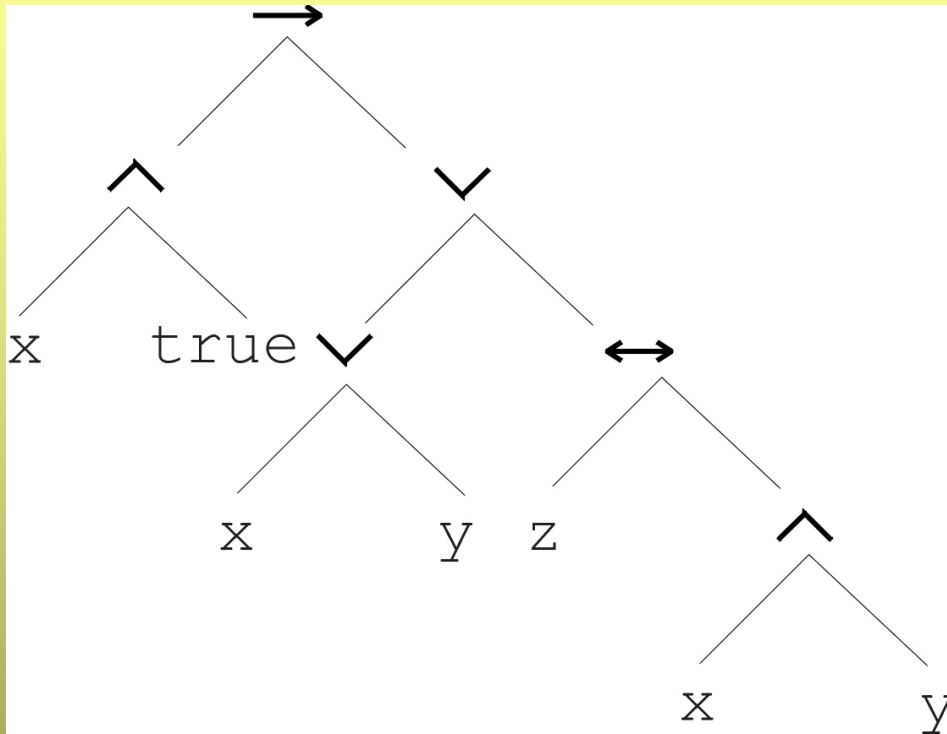
- Programa

```
i = 1;
while (i < 20)
{
    i = i + 1
}
```

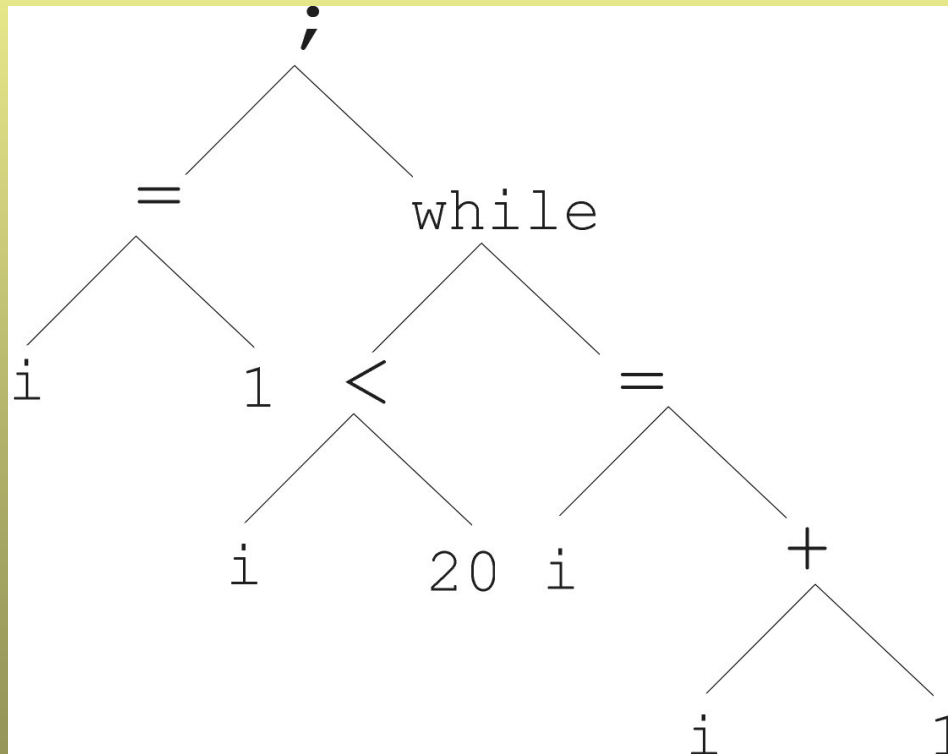
$$2 \cdot \pi + \left((x+3) - \frac{y}{5+1} \right)$$



$$(x \wedge \text{true}) \rightarrow ((x \vee y) \vee (z \leftrightarrow (x \wedge y)))$$



```
i = 1;
while (i < 20)
{
    i = i + 1
}
```



- En AGs, PE y EE la representación se basa en strings, por ejemplo, de bits, enteros, reales (vectores, permutaciones, etc)
- Estos árboles son estructuras no lineales
- En AGs, PE y EE los cromosomas (en general) son de una longitud fija
- En PG los árboles pueden variar de estructura → varían de longitud (profundidad, densidad, etc)

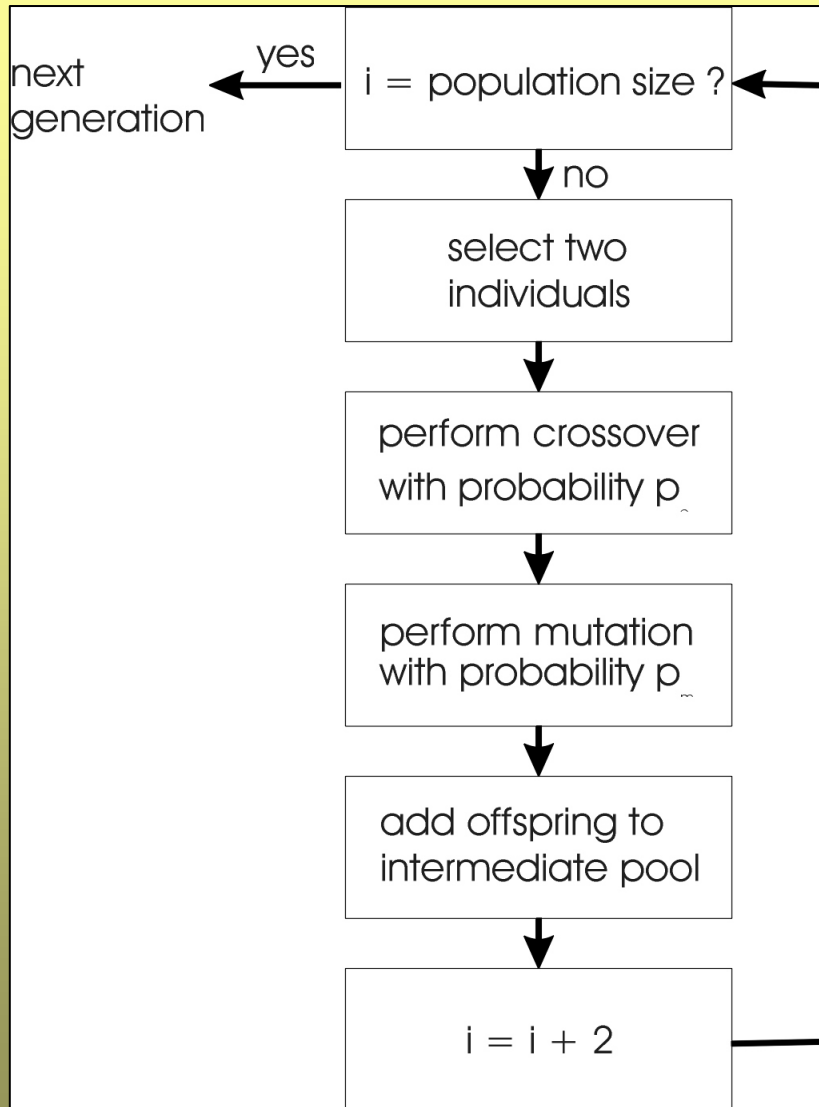
- Las expresiones representadas por arboles constan de:
 - un conjunto terminal T
 - un conjunto de funciones F (con sus aridades)
- Se adopta la siguiente definicion recursiva:
 1. Cada $t \in T$ es una expresion correcta
 2. $f(e_1, \dots, e_n)$ es una expresion correcta if $f \in F$, $\text{arity}(f)=n$ y e_1, \dots, e_n son expresiones correctas
 3. No hay otras formas correctas
- En general, GP no es fuertemente tipado (propiedad de clausura: cualquier $f \in F$ puede tomar $g \in F$ como argumento)

Creacion de Hijos

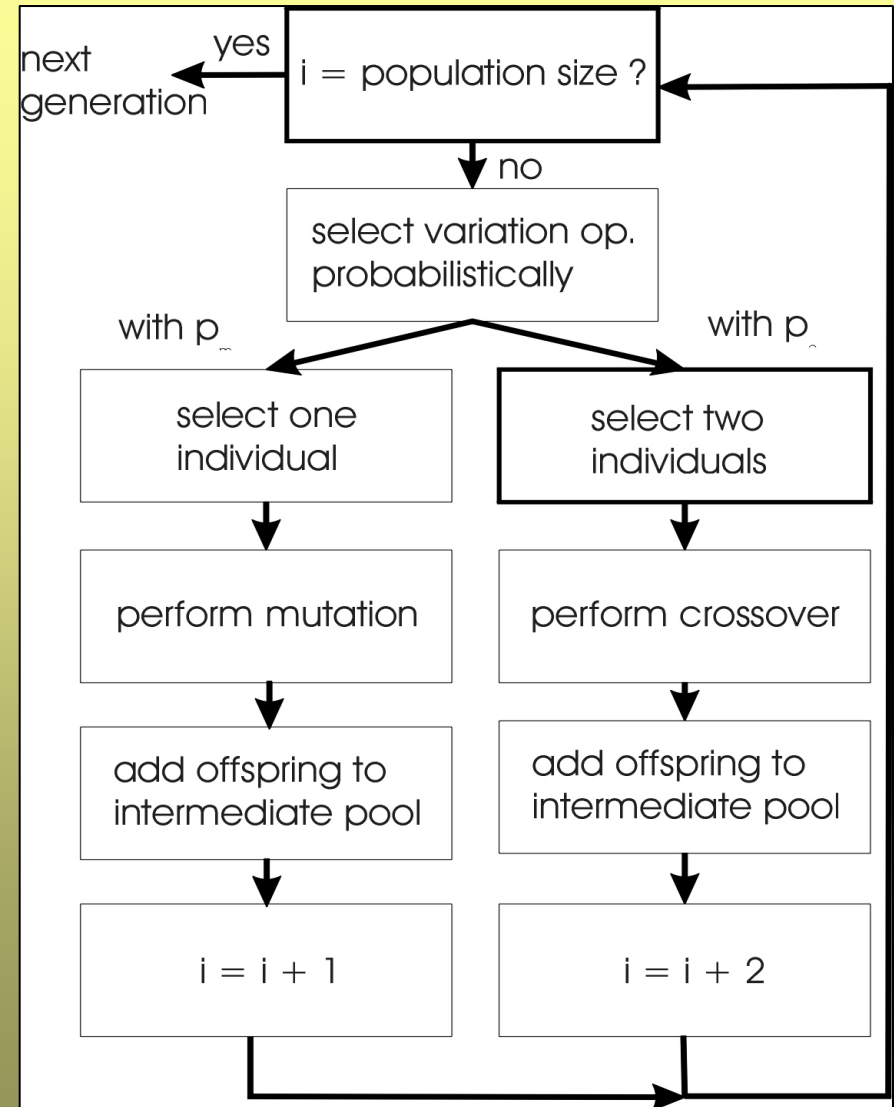
Comparen:

- Un AG usa secuencialmente mutacion y crossover (con probs)
- Un PG usa alternativamente mutacion o crossover (con probs)

GA flowchart

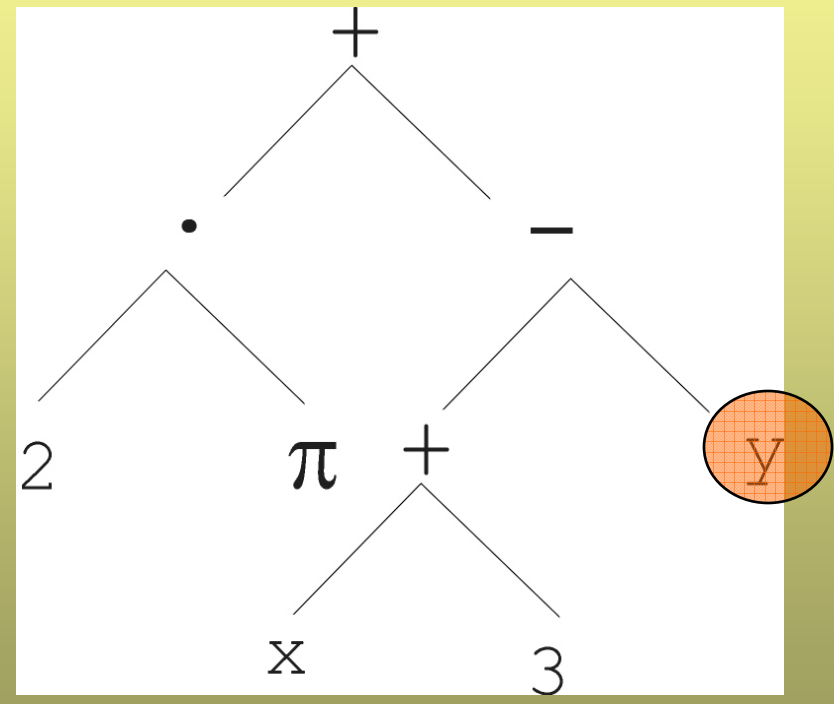
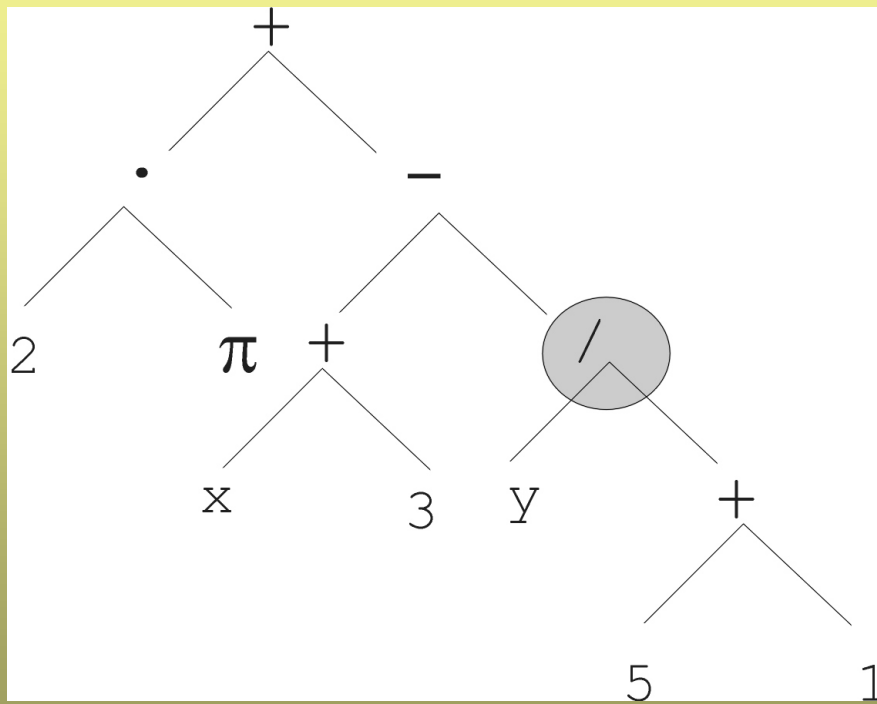


GP flowchart



Mutación

- La mutación mas comun: reemplazar al azar un subtree con uno nuevo creado al azar

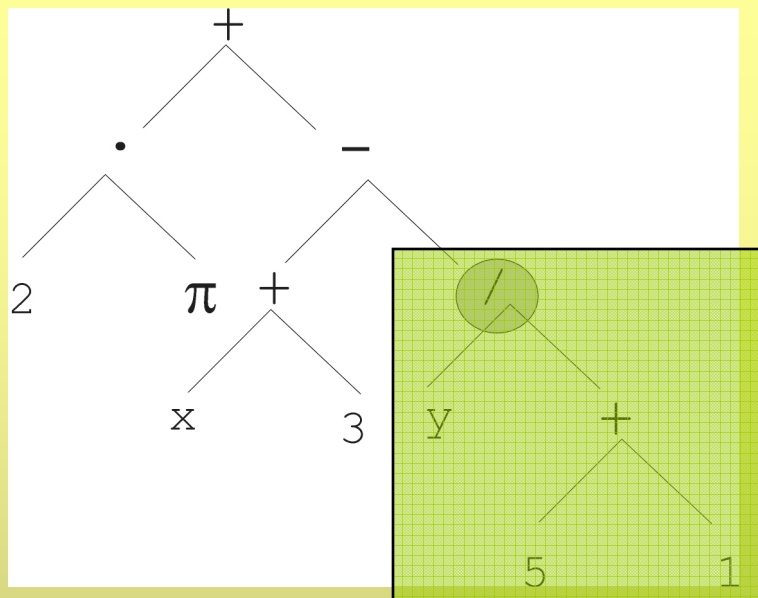


- La mutacion tiene dos parametros:
 - Probabilidad p_m de seleccionar mutacion en vez de crossover
 - Probabilidad de seleccionar un nodo interno como la raiz del subarbol a mutar
- Notablemente Koza recomienda $p_m = 0$ (Koza'92) o muy chica, como por ejemplo 0.05 (Banzhaf et al. '98)
- Notar que el tamaño del hijo puede exceder el tamaño del padre

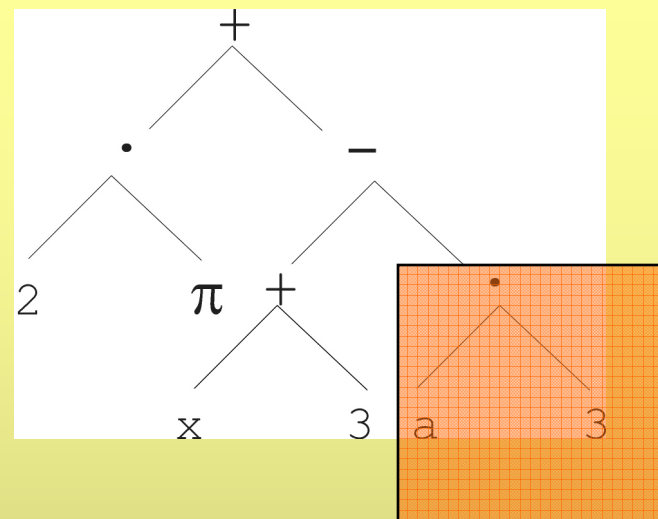
Recombinacion

- La recombinacion mas comun: intercambiar dos subarboles entre los padres
- Recombinacion usa dos parametros:
 - Probabilidad p_c de elegir recombinacion y no mutacion
 - La probabilidad de elegir un nodo interno para recombinar
- Aqui tambien el resultado puede ser un hijo mas grande que los padres

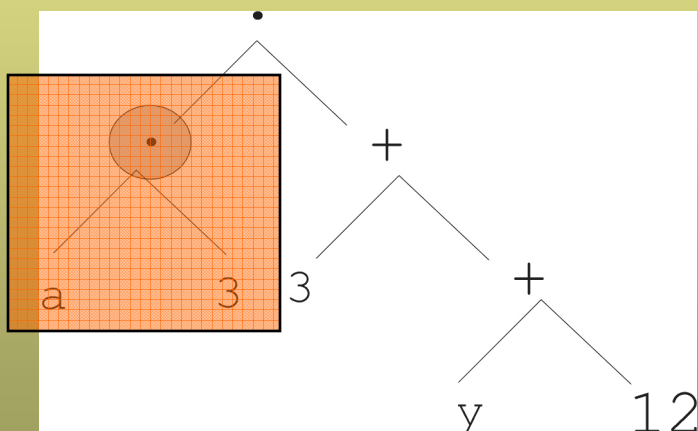
P1



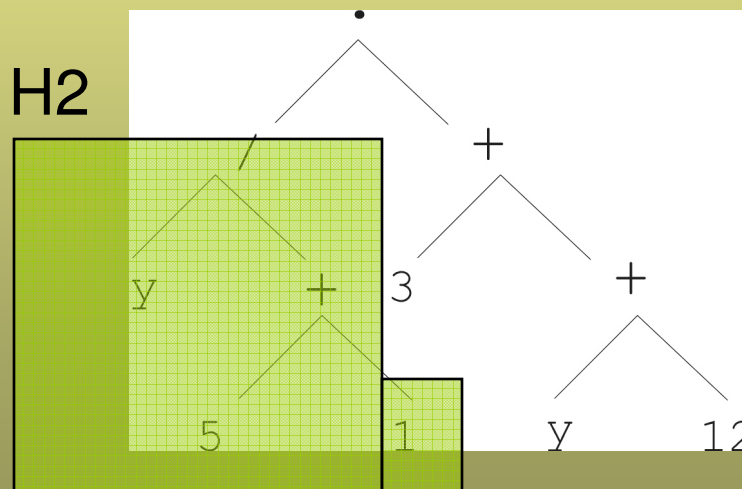
H1



P2



H2



Selección

- Para la selección de padres se usa generalmente proporcional al fitness
- Sobre selección en poblaciones grandes:
 - se rankea y divide la población en dos grupos:
 - grupo 1: los mejores $x\%$ de la población, grupo 2 el resto
 - El 80% de los padres viene del grupo 1 (para la reproducción) y solo un 20% del segundo
 - para $|poblacion| = 1000, 2000, 4000, 8000$ $x = 32\%, 16\%, 8\%, 4\%$
 - motivación: incrementar la eficiencia, el $\%$ es un “rule of thumb”
- Selección de sobrevivientes
 - Típicamente: generacional, esto es hijos reemplazan padres
 - Sin embargo recientemente la gente utiliza steady state por el elitismo

Initialisacion

- Se define una profundidad maxima D_{\max} para los arboles
- Metodo completo (cada rama realiza D_{\max}):
 - nodos con profundidad $d < D_{\max}$ elegidos al azar del conjunto F
 - nodes con profundidad $d = D_{\max}$ elegidos al azar del conjunto T
- Metodo constructivo (cada rama con prof $\leq D_{\max}$):
 - nodos con profundidad $d < D_{\max}$ elegidos al azar de la union de F y T
 - nodes con profundidad $d = D_{\max}$ elegidos al azar del conjunto T
- Inicialisaciones: “ramped half-and-half”, donde cada uno de los anteriores produce la mitad de la poblacion

“Bloat”

- Bloat = “Selección del más gordo”, i.e., los árboles en la población tienden a crecer
- Hay un debate sobre las razones de bloat
- Tomar contra-medidas, e.g.
 - Operadores que crean chicos más bien “deportivos” y no tan “gordos”
 - Parsimonia: penalizar los designoides grandes sin cambiar los operadores
 - Híbrido de ambos métodos

PG para controladores

- Árboles para “data fitting” vs. árboles (programas) ejecutables
- La ejecución de un programa podría cambiar el entorno donde el designoide está embebido → cambio de la función de fitness (si es que la evolución acontece on-line)
- Ejemplo: robot controller
- Los designoides (controladores) se evalúan off-line por simulaciones que suelen ser o caras o MUY caras
- Sin embargo tienden a ser bastante buenos

Problema mas tipico: regresion simbolica

- Dados puntos ejemplos en \mathbf{R}^2 , $(x_1, y_1), \dots, (x_n, y_n)$
- Encontrar una funcion $f(x)$ s.t. $\forall i = 1, \dots, n : f(x_i) = y_i$
- PG al rescate:
 - Representacion con $F = \{+, -, /, \sin, \cos\}$, $T = \mathbf{R} \cup \{x\}$
 - El inverso del error es el fitness
 - Operadores estandard $err(f) = \sum_{i=1}^n (f(x_i) - y_i)^2$
 - $|poblacion| = 1000$, inicializacion ramped half-half
 - Terminacion: n “hits” o 50000 evaluaciones de fitness (donde “hit” es $|f(x_i) - y_i| < 0.0001$, precision)

Discussion

Que es PG:

El arte de evolucionar programas computacionales ?

El medio para programacion automatica?

Un AG con otra representacion?