

REGRESIÓN SIMBÓLICA

- Búsqueda de funciones que se ajusten a una serie de puntos dado
- Búsqueda en el espacio de todas las fórmulas matemáticas posibles las que mejor predicen la variable de salida tomando como entrada las variables de entrada
 - Usando un conjunto de funciones base como las funciones aritméticas, trigonométricas y/o exponenciales
 - Espacio de búsqueda enorme
- Al contrario que las técnicas de regresión clásicas como RR.NN.AA. o SVR (SVM para Regresión), el resultado es una ecuación matemática explícita

REGRESIÓN SIMBÓLICA

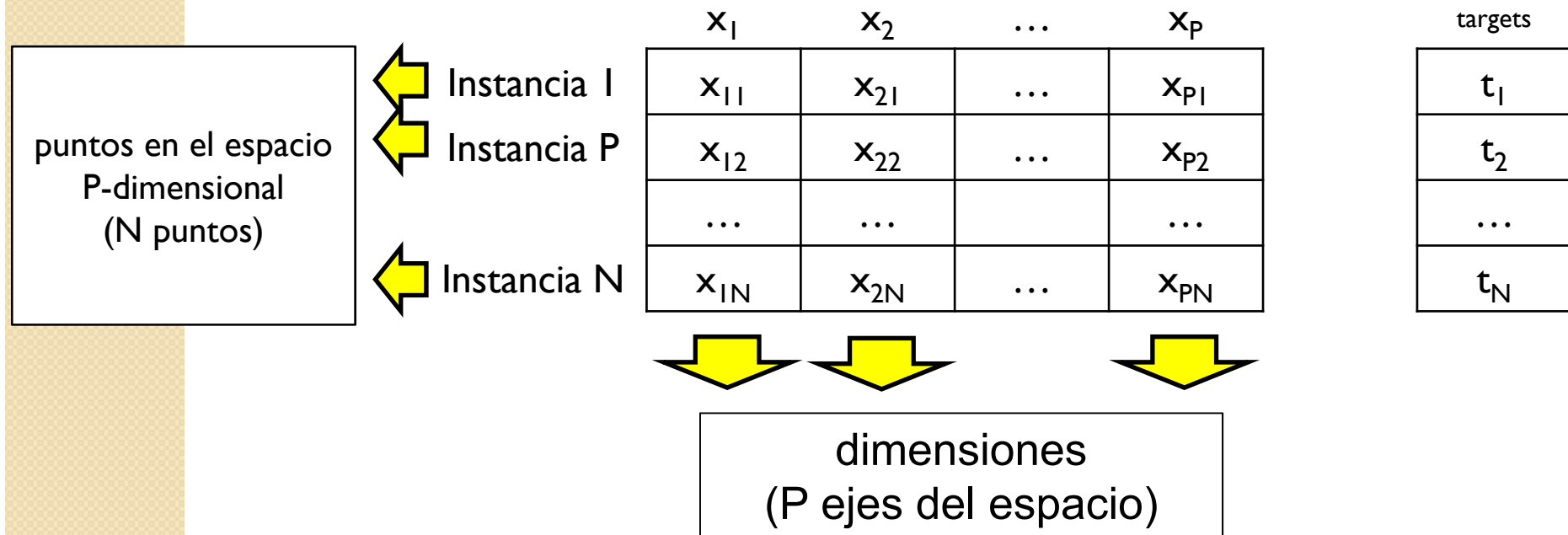
- Principales técnicas, en dos grandes grupos:
 - Basadas en Programación Genética
 - Genetic Programming (Koza, 1992)
 - Age-fitness Pareto Optimization (Schmidt and Lipson, 2011)
 - ϵ -Lexicase selection (La Cava et al., 2016)
 - Geometric Semantic Genetic Programming (Moraglio et al., 2012)
 - Multiple Regression Genetic Programming (Arnaldo et al., 2014)

REGRESIÓN SIMBÓLICA

- Principales técnicas, en dos grandes grupos:
 - Técnicas de Aprendizaje Automático (incluye técnicas de regresión no simbólica):
 - Adaptive Boosting (AdaBoost) Regression (Drucker, 1997)
 - Gradient Boosting Regression (Friedman, 2000)
 - Kernel Ridge (Murphy, 2012)
 - Least-Angle Regression with Lasso (Tibshirani, 1994)
 - Linear Regression (Efron et al., 2004)
 - Linear Support Vector Regression (Smola and Schölkopf, 2004)
 - Multilayer Perceptrons (MLPs) Regressor (Kingma and Ba, 2014)
 - Random Forests Regression (Breiman, 2001)
 - Stochastic Gradient Descent Regression (Pedregosa et al., 2011)
 - Extreme Gradient Boosting (Chen and Guestrin, 2016)
 - **Development of Mathematical Expressions**

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Cambio en la forma de interpretar los datos
 - Forma «clásica» en Aprendizaje Automático
 - N patrones, P variables:

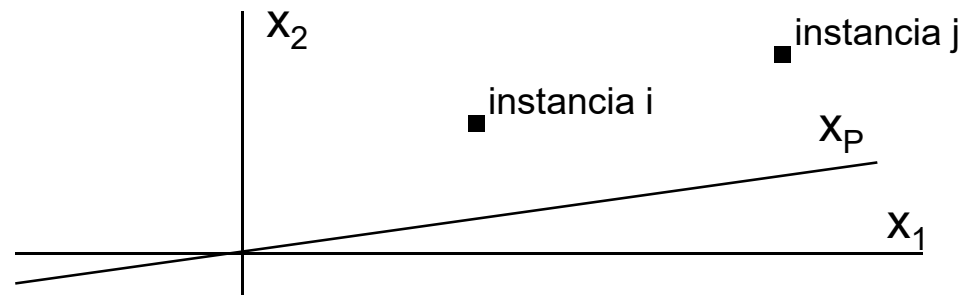


REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Cambio en la forma de interpretar los datos
 - Forma «clásica» en Aprendizaje Automático
 - N patrones, P variables:

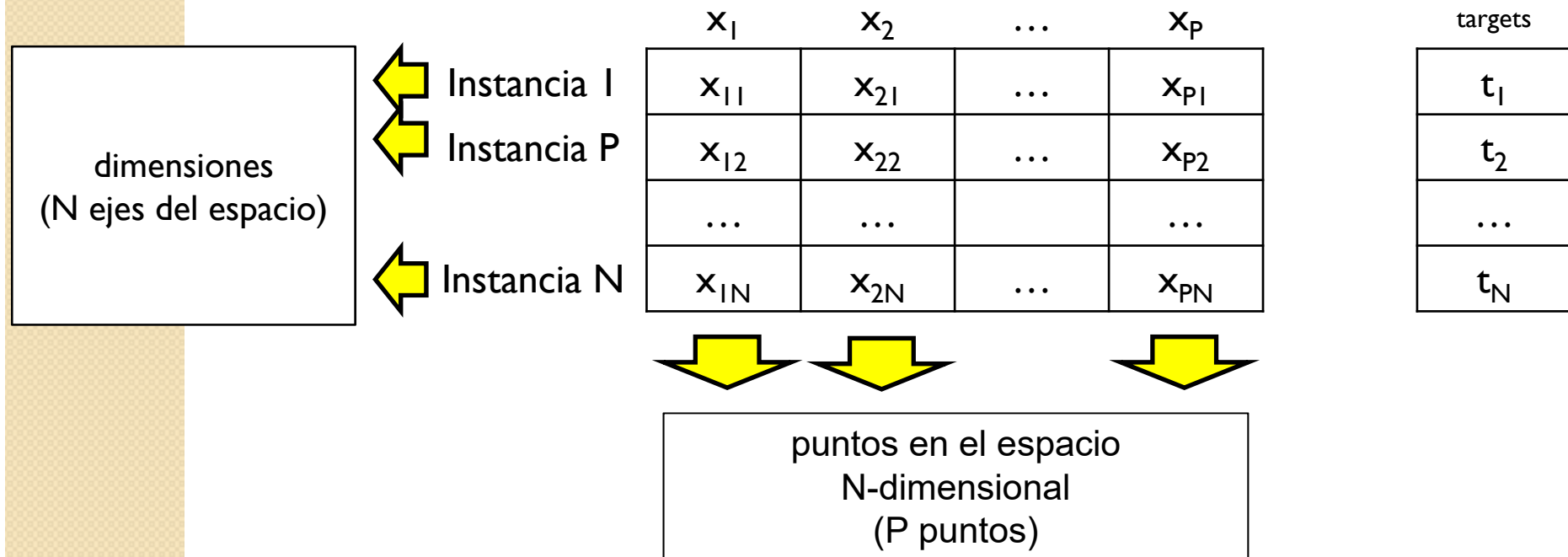
	x_1	x_2	...	x_p	targets
Instancia I	x_{11}	x_{21}	...	x_{p1}	t_1
Instancia P	x_{12}	x_{22}	...	x_{p2}	t_2

Instancia N	x_{1N}	x_{2N}	...	x_{pN}	t_N



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Cambio en la forma de interpretar los datos
 - Espacio semántico
 - N patrones, P variables:

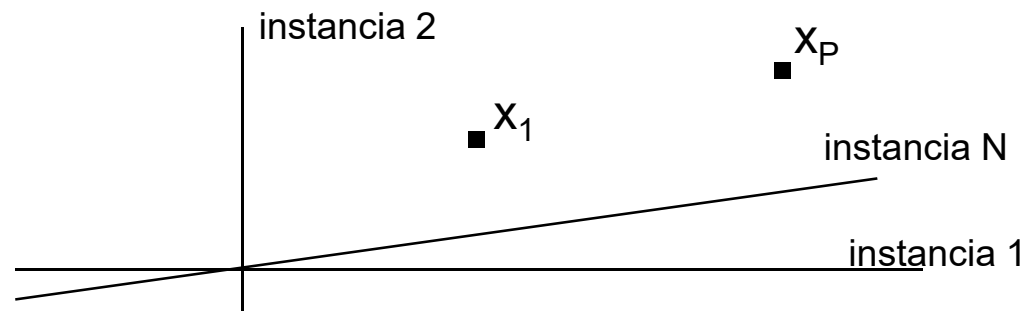


REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Cambio en la forma de interpretar los datos
 - Espacio semántico
 - N patrones, P variables:

	x_1	x_2	...	x_p	targets
Instancia 1	x_{11}	x_{21}	...	x_{p1}	t_1
Instancia P	x_{12}	x_{22}	...	x_{p2}	t_2

Instancia N	x_{1N}	x_{2N}	...	x_{pN}	t_N



REGRESIÓN SIMBÓLICA

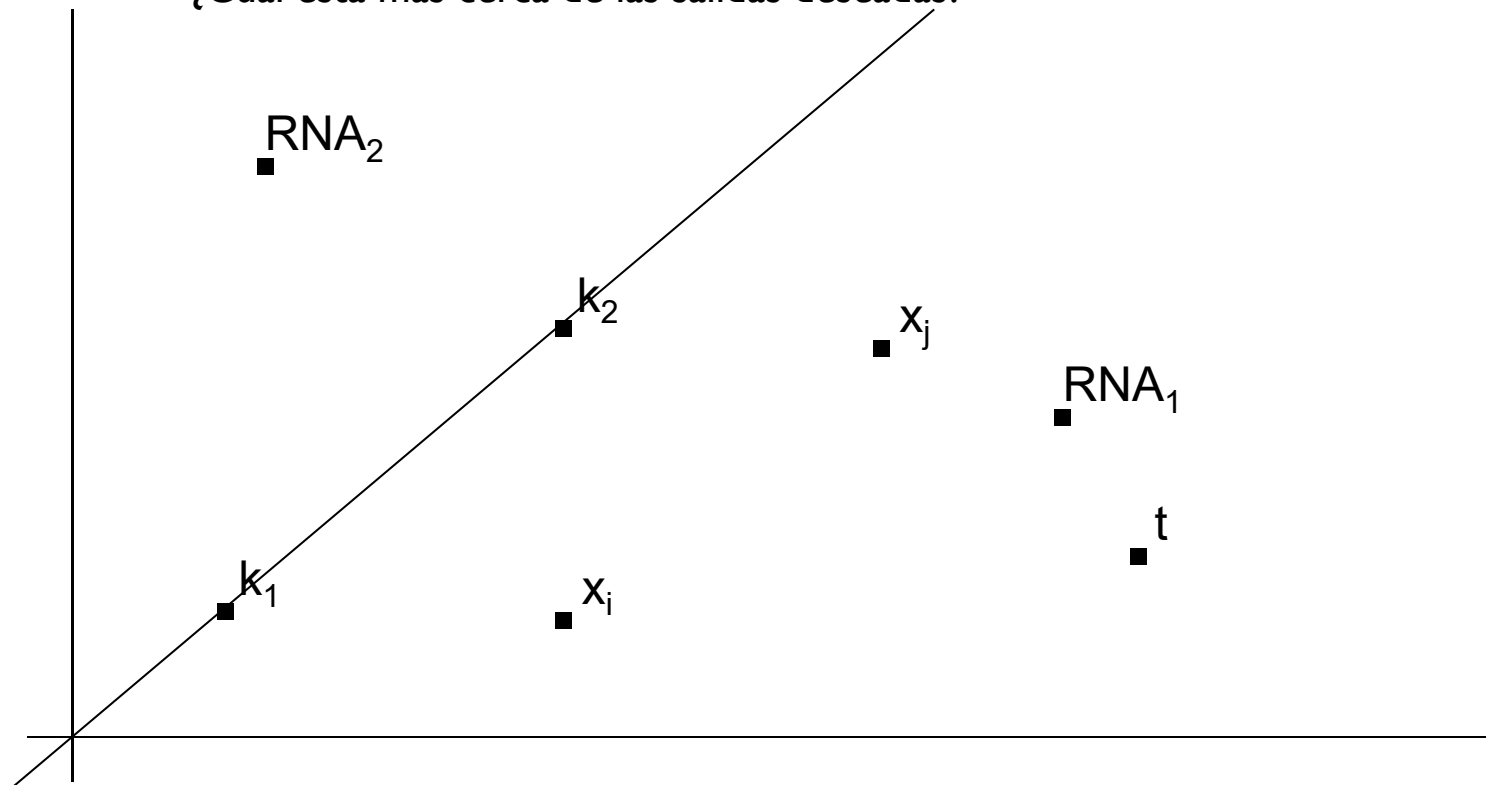
- Development of Mathematical Expressions (DoME)
 - Espacio semántico
 - Cada punto del espacio N-dimensional está definido por un valor para cada patrón del conjunto de entrenamiento
 - Por lo tanto, cada punto puede ser:
 - Las salidas de un modelo, para cada patrón
 - Sea el modelo que sea (RNA, SVR, etc.)
 - Una constante
 - Las constantes se pueden interpretar como un modelo:
 - $f(x_1, x_2, \dots, x_p) = k$
 - Modelo que devuelve el mismo valor para cada patrón
 - Las constantes se sitúan en la línea $(k, k, \dots, k) = k^*(1, 1, \dots, 1)$
 - Los distintos valores de cada variable para todos los patrones
 - Las variables se pueden interpretar como un modelo:
 - $f(x_1, x_2, \dots, x_p) = x_i$
 - Las salidas deseadas (targets)
 - El modelo que se busca

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

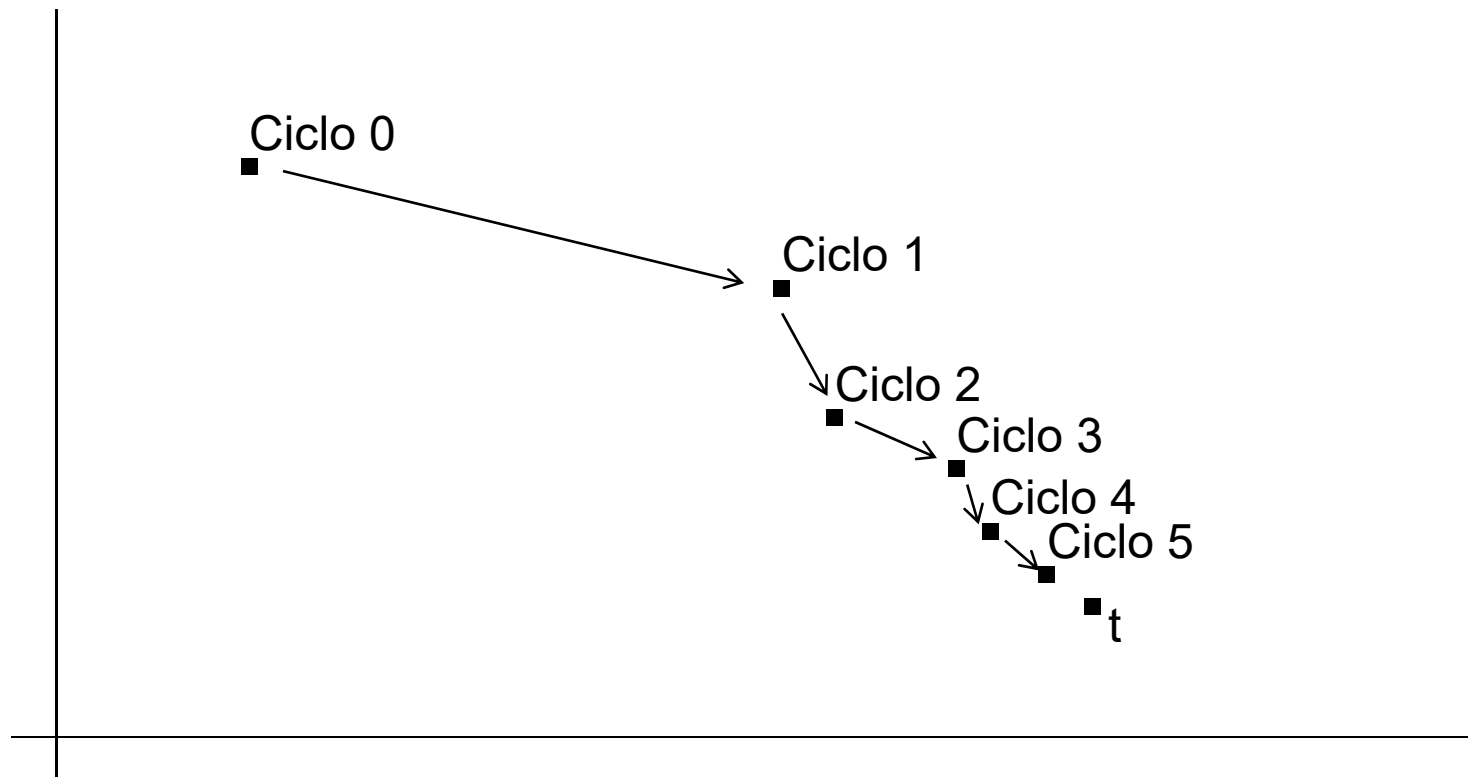
- Espacio semántico

- ¿Qué modelo de los siguientes es mejor?
 - ¿Cuál está más cerca de las salidas deseadas?



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Espacio semántico
 - Entrenamiento de una RNA



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Espacio semántico
 - Medida del error: distancia a las salidas deseadas:

$$\sqrt{\sum_{i=0}^N (o_i - t_i)^2} = \sqrt{SSE} = RSSE$$

- o_i : salidas del modelo
 - SSE: *Sum Squared Error*
 - RSSE: *Root Sum Squared Error*
- Ecuación de una esfera N-dimensional centrada en t de radio RSSE
 - En dos dimensiones, la ecuación de un círculo centrado en (x_0, y_0) es

$$(x - x_0)^2 + (y - y_0)^2 = R^2$$

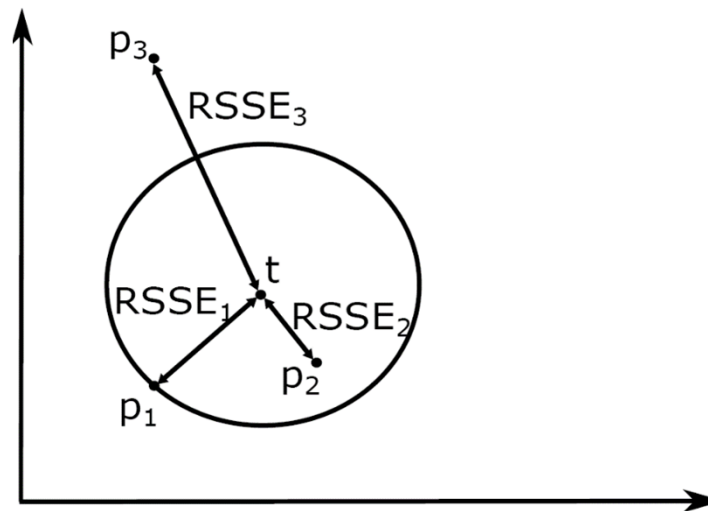
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

- Espacio semántico

- Las salidas deseadas t junto con el modelo p_i definen una esfera N-dimensional

- Radio: $RSSE_i$



- Un modelo que esté en el interior de la esfera tendrá menor distancia a t , y por tanto menor RSSE
 - Modelo p_2 , con $RSSE_2$
 - Un modelo en el exterior tendrá mayor RSSE
 - Modelo p_3 , con $RSSE_3$

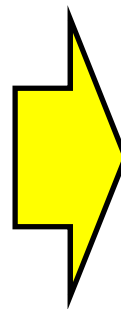
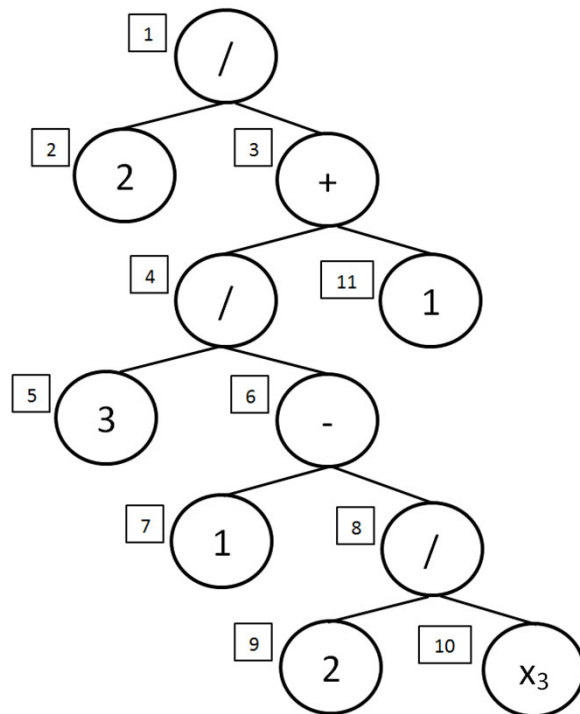
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Espacio semántico
 - Para reducir el impacto del número de patrones, se introduce este en la ecuación
 - La medida del error es el MSE (*Mean Squared Error*):

$$MSE = \frac{1}{N} \sum_{i=0}^N (o_i - t_i)^2$$

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Representación en forma de árbol
 - Similar a la Programación Genética clásica
 - Ejemplo:



$$f(x_1, x_2, x_3, \dots) = \frac{2}{\frac{3}{1 - \frac{2}{x_3}} + 1}$$

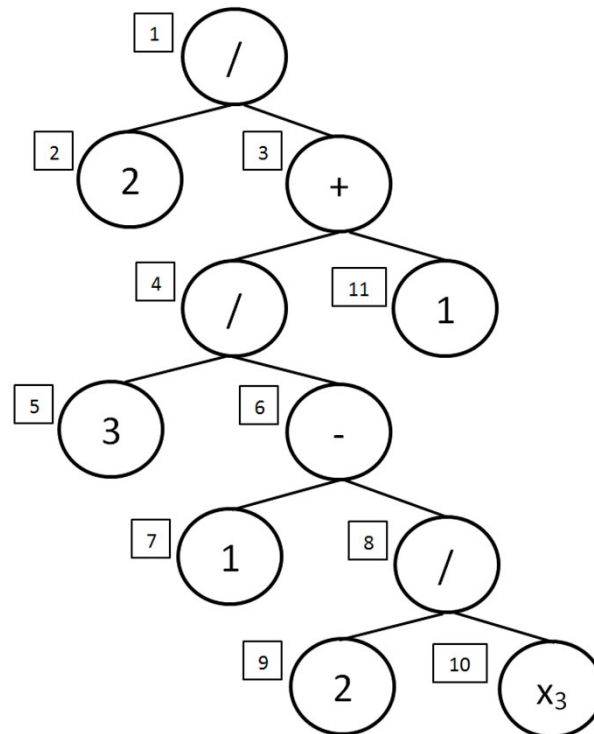
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Representación en forma de árbol
 - Dos tipos de nodos:
 - Nodos terminales:
 - Variables
 - Constantes
 - Nodos no terminales:
 - Funciones aritméticas: $+$, $-$, $*$, $/$
 - La semántica del árbol se define como las salidas del árbol para cada patrón del conjunto de entrenamiento
 - Cada nodo del árbol define un subárbol, que se puede interpretar como un nuevo modelo
 - Una nueva ecuación
 - Un «trozo» de la ecuación final
 - Por tanto, **cada nodo tiene su propia semántica**
 - Resultado de evaluar ese subárbol con el conjunto de entrenamiento
 - **Cada nodo es un punto en el espacio N-dimensional**

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Representación en forma de árbol
 - Ejemplo:

- Valores de x_3 : (1, 4, -2)



Node	Semantic
1	$(-1, 0.286, 0.8)$
2	$(2, 2, 2)$
3	$(-2, 7, 2.5)$
4	$(-3, 6, 1.5)$
5	$(3, 3, 3)$
6	$(-1, 0.5, 2)$
7	$(1, 1, 1)$
8	$(2, 0.5, -1)$
9	$(2, 2, 2)$
10	$(1, 4, -2)$
11	$(1, 1, 1)$

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Idea clave: calcular el error cometido **a partir de la semántica de cada nodo del árbol**
 - Cada nodo tiene asociados 4 vectores: a, b, c, d
 - La ecuación para calcular el MSE para un nodo es:

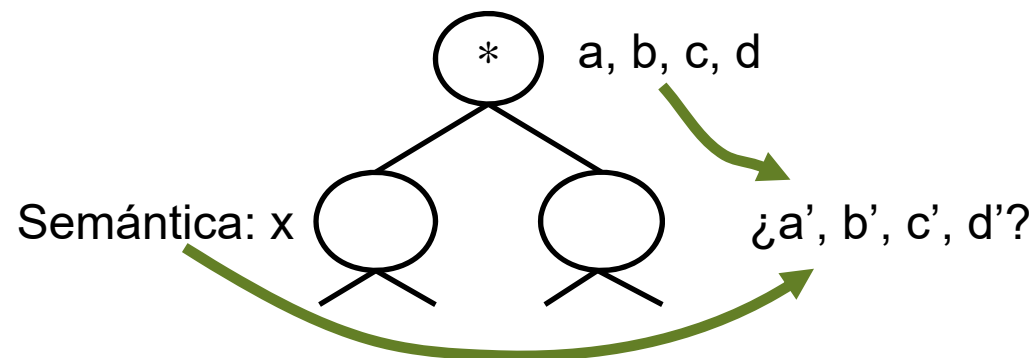
$$MSE = \frac{1}{N} \sum_{i=0}^N \left(\frac{a_i \cdot o_i - b_i}{c_i \cdot o_i - d_i} \right)^2$$

- o_i : salidas de ese nodo (semántica)
- Para el nodo raíz: $a_i=1$, $b_i=t_i$, $c_i=0$, $d_i=-1$
 - Lleva a

$$MSE = \frac{1}{N} \sum_{i=0}^N (o_i - t_i)^2$$

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - ¿Cómo se calculan los vectores a, b, c, d para cada nodo?
 - Cada nodo no terminal calcula los vectores a', b', c', d' de cada hijo en función de sus propios vectores a, b, c, d y la salida del otro hijo
 - Si x define la semántica del hijo izq., y define la semántica del hijo dcho
 - Hijo izq: a', b', c', d' se definen en función de a, b, c, d, y
 - Hijo dcho: a', b', c', d' se definen en función de a, b, c, d, x
 - Por ejemplo, nodo $*$, para el hijo 2:



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - ¿Cómo se calculan los vectores a,b,c,d para cada nodo?
 - Cada nodo no terminal calcula los vectores a',b',c',d' de cada hijo en función de sus propios vectores a,b,c,d y la salida del otro hijo
 - Si x define la semántica del hijo izq., e y define la semántica del hijo dcho
 - Hijo izq: a',b',c',d' se definen en función de a,b,c,d,y
 - Hijo dcho: a',b',c',d' se definen en función de a,b,c,d,x
 - Por ejemplo, nodo *, para el hijo 2:

$$MSE = \frac{1}{N} \sum_{i=0}^N \left(\frac{a_i \cdot o_i - b_i}{c_i \cdot o_i - d_i} \right)^2 = \frac{1}{N} \sum_{i=0}^N \left(\frac{a_i \cdot (x_i \cdot y_i) - b_i}{c_i \cdot (x_i \cdot y_i) - d_i} \right)^2 = \frac{1}{N} \sum_{i=0}^N \left(\frac{(a_i \cdot x_i) \cdot y_i - b_i}{(c_i \cdot x_i) \cdot y_i - d_i} \right)^2 = \frac{1}{N} \sum_{i=0}^N \left(\frac{a'_i \cdot y_i - b'_i}{c'_i \cdot y_i - d'_i} \right)^2$$

que es la ecuación $MSE = \frac{1}{N} \sum_{i=0}^N \left(\frac{a'_i \cdot o_i - b'_i}{c'_i \cdot o_i - d'_i} \right)^2$ para el segundo hijo con

$$a'_i = a_i \cdot x_i \quad b'_i = b_i \quad c'_i = c_i \cdot x_i \quad d'_i = d_i$$

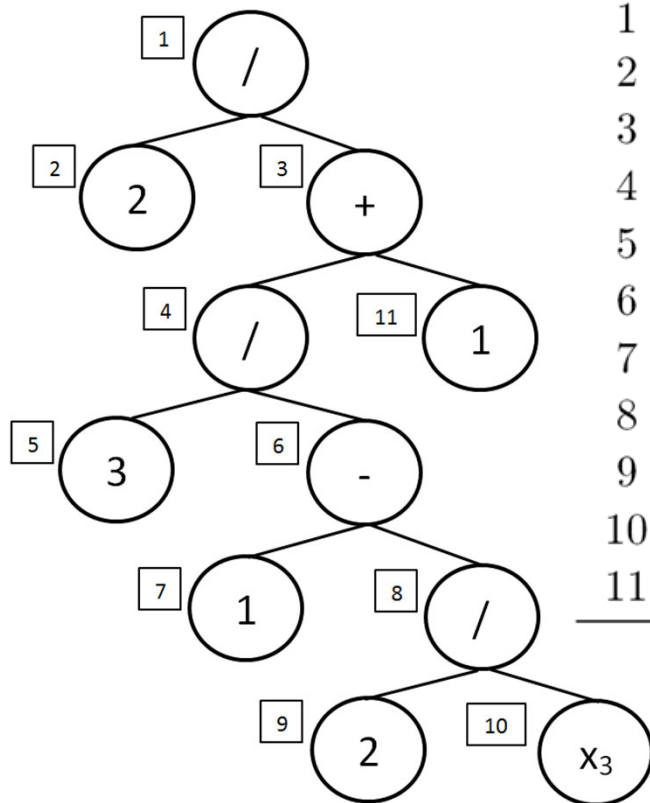
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - ¿Cómo se calculan los vectores a, b, c, d para cada nodo?
 - Haciendo un desarrollo similar con el resto de operaciones e hijos:

Operation	Left Child				Right Child			
	a'_i	b'_i	c'_i	d'_i	a'_i	b'_i	c'_i	d'_i
+	a_i	$b_i - a_i \cdot y_i$	c_i	$d_i - c_i \cdot y_i$	a_i	$b_i - a_i \cdot x_i$	c_i	$d_i - c_i \cdot x_i$
-	a_i	$b_i + a_i \cdot y_i$	c_i	$d_i + c_i \cdot y_i$	a_i	$a_i \cdot x_i - b_i$	c_i	$c_i \cdot x_i - d_i$
*	$a_i \cdot y_i$	b_i	$c_i \cdot y_i$	d_i	$a_i \cdot x_i$	b_i	$c_i \cdot x_i$	d_i
/	a_i	$b_i \cdot y_i$	c_i	$d_i \cdot y_i$	b_i	$a_i \cdot x_i$	d_i	$c_i \cdot x_i$

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Ejemplo:



Node	Semantic	Equation
1	$(-1, 0.286, 0.8)$	$\frac{1}{3}((o_i - 5)^2 + (o_i - 4)^2 + (o_i - 1)^2)$
2	$(2, 2, 2)$	$\frac{1}{3}((\frac{o_i+10}{2})^2 + (\frac{o_i-28}{-7})^2 + (\frac{o_i-2.5}{-2.5})^2)$
3	$(-2, 7, 2.5)$	$\frac{1}{3}((\frac{-5 \cdot o_i + 2}{o_i})^2 + (\frac{-4 \cdot o_i + 2}{o_i})^2 + (\frac{-o_i + 2}{o_i})^2)$
4	$(-3, 6, 1.5)$	$\frac{1}{3}((\frac{-5 \cdot o_i - 3}{o_i + 1})^2 + (\frac{-4 \cdot o_i - 2}{o_i + 1})^2 + (\frac{-o_i + 1}{o_i + 1})^2)$
5	$(3, 3, 3)$	$\frac{1}{3}((\frac{-5 \cdot o_i + 3}{o_i - 1})^2 + (\frac{-4 \cdot o_i - 1}{o_i + 0.5})^2 + (\frac{-o_i + 2}{o_i + 2})^2)$
6	$(-1, 0.5, 2)$	$\frac{1}{3}((\frac{-3 \cdot o_i - 15}{o_i + 3})^2 + (\frac{-2 \cdot o_i - 12}{o_i + 3})^2 + (\frac{o_i - 3}{o_i + 3})^2)$
7	$(1, 1, 1)$	$\frac{1}{3}((\frac{-3 \cdot o_i - 9}{o_i + 1})^2 + (\frac{-2 \cdot o_i - 11}{o_i + 2.5})^2 + (\frac{o_i - 2}{o_i + 4})^2)$
8	$(2, 0.5, -1)$	$\frac{1}{3}((\frac{3 \cdot o_i - 18}{-o_i + 4})^2 + (\frac{2 \cdot o_i - 14}{-o_i + 4})^2 + (\frac{-o_i - 2}{-o_i + 4})^2)$
9	$(2, 2, 2)$	$\frac{1}{3}((\frac{3 \cdot o_i - 18}{-o_i + 4})^2 + (\frac{2 \cdot o_i - 56}{-o_i + 16})^2 + (\frac{-o_i + 4}{-o_i - 8})^2)$
10	$(1, 4, -2)$	$\frac{1}{3}((\frac{-18 \cdot o_i + 6}{4 \cdot o_i - 2})^2 + (\frac{-14 \cdot o_i + 4}{4 \cdot o_i - 2})^2 + (\frac{-2 \cdot o_i - 2}{4 \cdot o_i - 2})^2)$
11	$(1, 1, 1)$	$\frac{1}{3}((\frac{-5 \cdot o_i + 17}{o_i - 3})^2 + (\frac{-4 \cdot o_i - 22}{o_i + 6})^2 + (\frac{-o_i + 0.5}{o_i + 1.5})^2)$

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Observaciones (1/5):
 - Cada nodo tiene una ecuación que permite calcular el MSE global
 - Para cada nodo va a dar el mismo valor
 - Esta ecuación depende de:
 - Los vectores a, b, c, d de ese nodo
 - La semántica de ese nodo
 - Las salidas de ese nodo para cada patrón del conjunto de entrenamiento

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Observaciones (2/5):
 - Para calcular los vectores a, b, c, d de ese nodo no es necesario conocer la semántica de ese nodo
 - Es decir, estos vectores son independientes de ese subárbol
 - Dependen del resto del árbol que no cuelga de ese nodo
 - Para calcular la semántica de ese nodo no es necesario conocer el resto del árbol
 - Sólo es necesario conocer ese subárbol
 - Es independiente del resto del árbol que no cuelga de ese nodo

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

- Observaciones (3/5):

- Ejemplo: Nodo 6:

- Para calcular la semántica se necesitan las semánticas de los nodos 6, 7, 8, 9, 10

- De todo el subárbol que cuelga de ese nodo, expresión « $1-(2/x_3)$ »

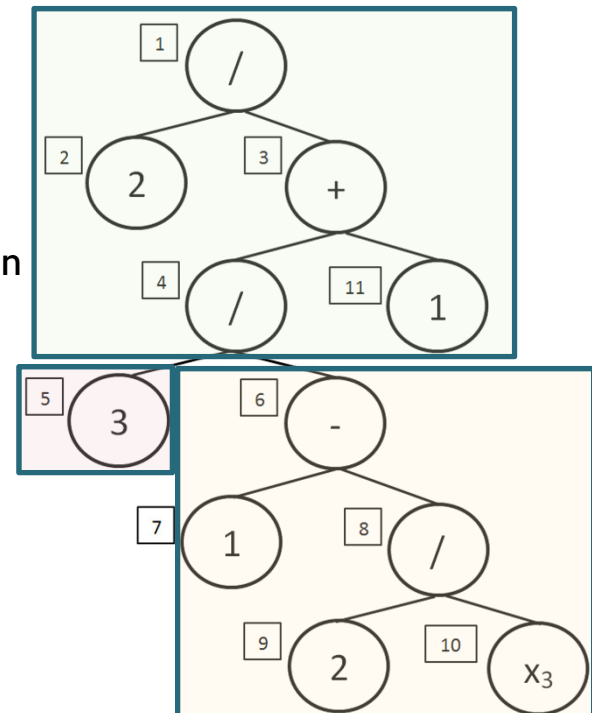
- Para calcular los vectores a,b,c,d se necesita:

- La semántica de su nodo «hermano» 3

- Todo el subárbol que cuelga de ese nodo

- Los vectores a,b,c,d del nodo 4

- Para calcularlos, es necesaria la información de toda la parte superior del árbol

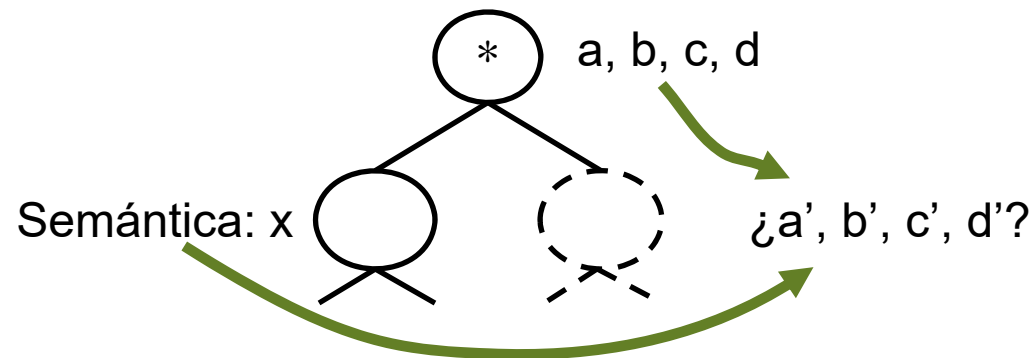


REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

- Observaciones (4/5):

- Los vectores a, b, c, d de un nodo son independientes de este y de todo el subárbol que cuelga de él
 - No dependen de la semántica de ese nodo



- Se pueden calcular para el hijo de un nodo padre, aunque todavía no haya subárbol
 - Una vez se haya creado un subárbol «candidato», se pueden usar estos 4 vectores a', b', c', d' junto con la semántica de ese subárbol para calcular cuál sería el MSE al insertar el subárbol en ese sitio

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Observaciones (5/5):
 - Los vectores a, b, c, d de un nodo son independientes de este y de todo el subárbol que cuelga de él
 - No dependen de la semántica de ese nodo
 - Por tanto, si se calculan para un subárbol y se cambia ese subárbol por otro, los vectores a', b', c', d' no cambian
 - Esto permite calcular cuál sería el valor del MSE global como resultado de cambiar ese subárbol por otro
 - Se calcula la semántica de ese nuevo subárbol
 - Se aplica la ecuación de cálculo de MSE con los vectores a, b, c, d de ese nodo y esta nueva semántica
 - Si el valor del nuevo MSE es inferior al anterior, se puede realizar el cambio del nodo por el nuevo subárbol

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

- Interpretación geométrica

- Operadores suma y resta

- Ejemplo: Árbol p, con hijos p₁ y p₂:

- Semántica de p₁: x = (1,2) Semántica de p₂: y = (2,1)

- Ecuaciones de ambos hijos:

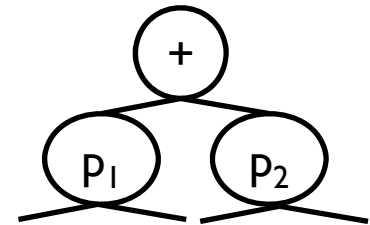
- Hijo 1:

$$MSE = \frac{1}{N} \sum_{i=0}^N (o_i - t_i)^2 = \frac{1}{N} \sum_{i=0}^N ((x_i + y_i) - t_i)^2 = \frac{1}{N} \sum_{i=0}^N (x_i - (t_i - y_i))^2$$

- Hijo 2:

$$MSE = \frac{1}{N} \sum_{i=0}^N (o_i - t_i)^2 = \frac{1}{N} \sum_{i=0}^N ((x_i + y_i) - t_i)^2 = \frac{1}{N} \sum_{i=0}^N (y_i - (t_i - x_i))^2$$

- Se crean dos nuevas esferas centradas en t-p₂ y t-p₁, con el mismo radio



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

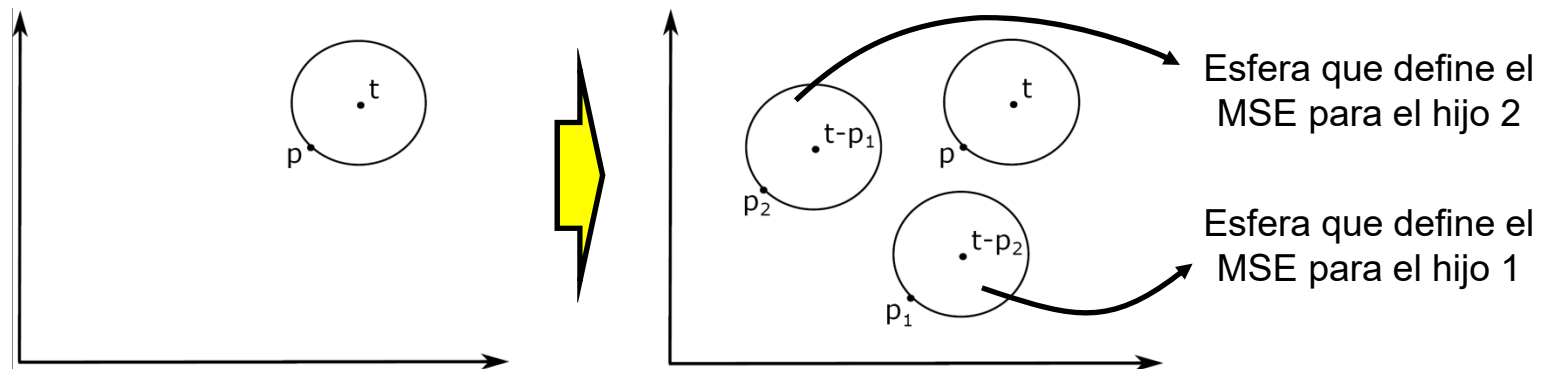
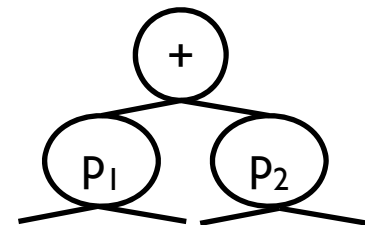
- Interpretación geométrica

- Operadores suma y resta

- Ejemplo: Árbol p , con hijos p_1 y p_2 :

- Semántica de p_1 : $x = (1,2)$ Semántica de p_2 : $y = (2,1)$

- Se crean dos nuevas esferas centradas en $t-p_2$ y $t-p_1$, con el mismo radio



Esfera que define el MSE para el hijo 2

Esfera que define el MSE para el hijo 1

- En lugar de buscar un árbol nuevo dentro de la esfera original, se puede buscar un árbol dentro de las dos nuevas esferas
 - Búsqueda local en cada nodo (nuevos sitios de búsqueda)
 - Si se encuentra, sustituir al hijo correspondiente por ese árbol
 - Como consecuencia, el árbol global p se acercará al punto t

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

- Interpretación geométrica

- Operador multiplicación

- Ejemplo: Árbol p, con hijos p₁ y p₂:

- Semántica de p₁: x = (1,2) Semántica de p₂: y = (2,1)

- Ecuaciones de ambos hijos:

- Hijo 1:

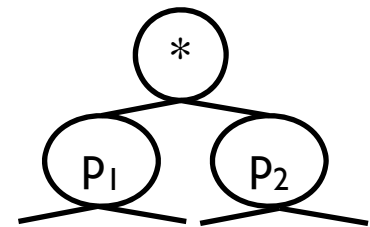
$$MSE = \frac{1}{N} \sum_{i=0}^N (x_i \cdot y_i - t_i)^2 \quad \Rightarrow \quad \sum_{i=0}^N \frac{\left(x_i - \frac{t_i}{y_i}\right)^2}{\left(\frac{SSE}{y_i}\right)^2} = 1$$

- Hijo 2:

$$MSE = \frac{1}{N} \sum_{i=0}^N (y_i \cdot x_i - t_i)^2 \quad \Rightarrow \quad \sum_{i=0}^N \frac{\left(y_i - \frac{t_i}{x_i}\right)^2}{\left(\frac{SSE}{x_i}\right)^2} = 1$$

- 2 nuevas esferas, centradas en t/p₂ y t/p₁, con los radios modificados cada eje

- Elipses



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

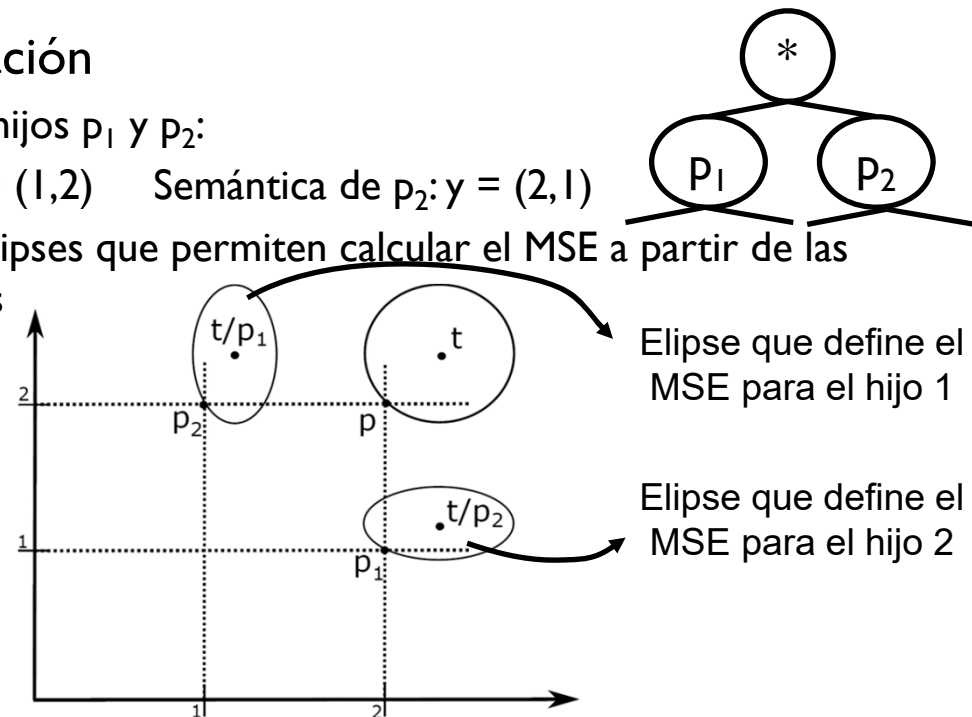
- Interpretación geométrica

- Operador multiplicación

- Ejemplo: Árbol p , con hijos p_1 y p_2 :

- Semántica de p_1 : $x = (1,2)$ Semántica de p_2 : $y = (2,1)$

- Se crean dos nuevas elipses que permiten calcular el MSE a partir de las semánticas de los hijos



- Al igual que antes, se puede realizar búsqueda local en cada nodo
 - Buscar un árbol dentro de cada elipse
 - Si se encuentra, se sustituye ese hijo por ese árbol y el árbol global p se habrá acercado a t

REGRESIÓN SIMBÓLICA

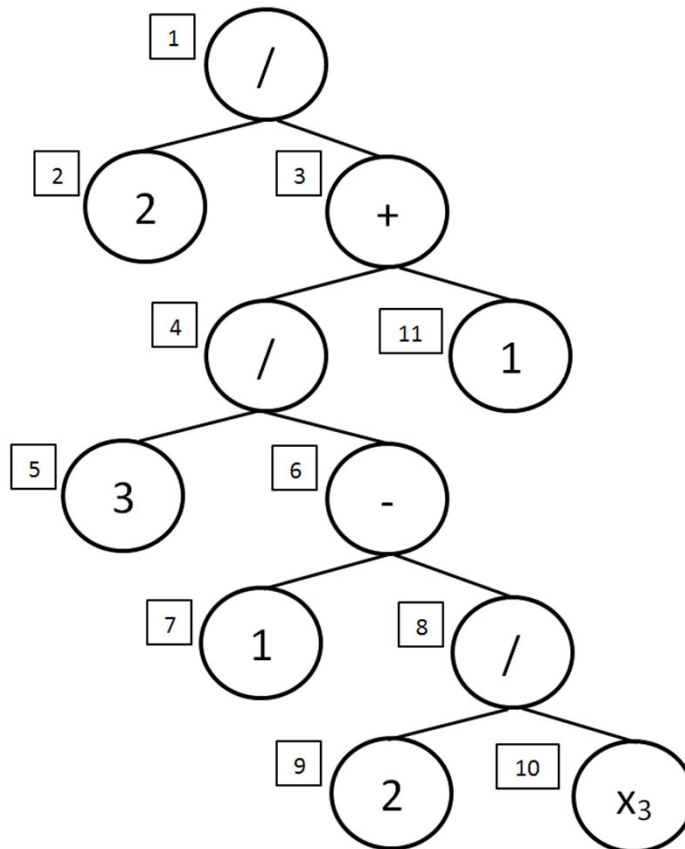
- Development of Mathematical Expressions (DoME)
 - Interpretación geométrica
 - Operador división
 - El lugar de crear esferas o elipses (elipsoides), crea formas totalmente distintas
 - Al igual que antes, se puede realizar búsqueda local en cada nodo
 - Buscar un árbol dentro de cada forma
 - Si se encuentra, se sustituye ese hijo por ese árbol y el árbol global p se habrá acercado a t

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Además, a la vez que se calculan los vectores a, b, c, d para cada nodo, es necesario calcular un conjunto S para cada nodo
 - Contiene qué valores están «prohibidos» en las semánticas de ese nodo
 - Porque llevarían a realizar una división por 0 en un nodo no terminal en alguna parte superior del árbol
 - El conjunto S del segundo hijo de una división (el divisor) siempre contiene la semántica $(0, 0, \dots, 0)$
 - Nodo raíz: $S = \emptyset$

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Ejemplo:

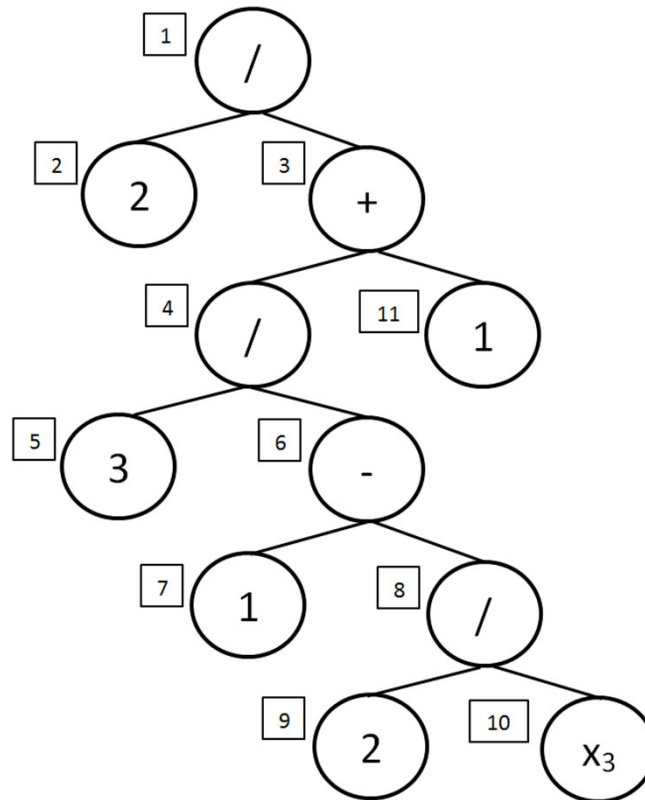
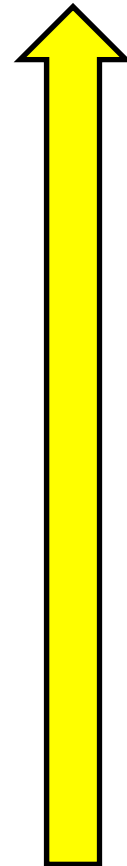


Node	S
1	\emptyset
2	\emptyset
3	$\{(0, 0, 0)\}$
4	$\{(-1, -1, -1)\}$
5	$\{(1, -0.5, -2)\}$
6	$\{(0, 0, 0), (-3, -3, -3)\}$
7	$\{(2, 0.5, -1), (-1, -2.5, -4)\}$
8	$\{(1, 1, 1), (4, 4, 4)\}$
9	$\{(1, 4, -2), (4, 16, -8)\}$
10	$\{(0, 0, 0), (2, 2, 2), (0.5, 0.5, 0.5)\}$
11	$\{(3, -6, -1.5)\}$

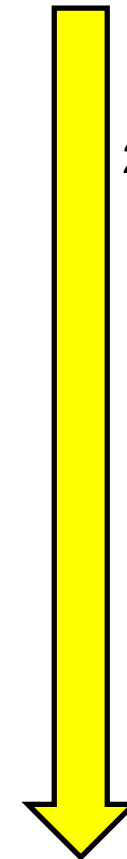
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Por lo tanto, la evaluación de un árbol se hace en 2 pases:

1. Evaluación del árbol desde los nodos terminales hasta la raíz (cálculo de las semánticas de los nodos)



2. Cálculo de los vectores a, b, c, d y conjunto S de cada nodo desde la raíz hasta los nodos terminales



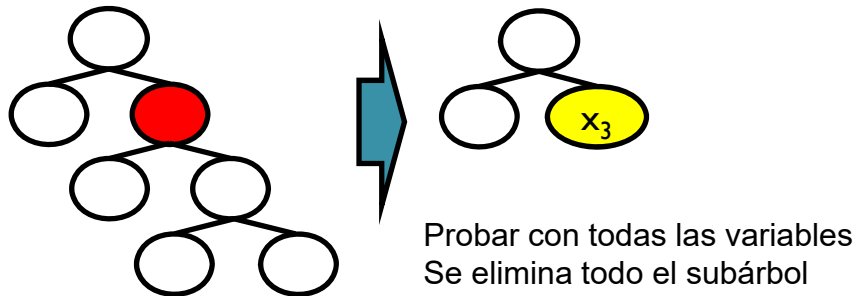
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo
 - Dado un nuevo subárbol con una semántica como candidato para reemplazar un nodo
 - Se evalúa la ecuación de MSE de ese nodo con la nueva semántica
 - Si se disminuye el MSE, se puede sustituir el nodo por ese subárbol
 - Eliminar el nodo y todo el subárbol que cuelga de él
 - ¿Cómo hallar un subárbol para un nodo concreto?
 - *Constant search*
 - Sustituir un subárbol por un nodo terminal con una constante
 - *Variable search*
 - Sustituir un subárbol por un nodo terminal con una variable
 - *Constant-variable search*
 - Sustituir un subárbol por un nuevo subárbol con 3 nodos:
 - <constante> <operación> <variable>
 - *Constant-expression search*
 - Desplazar un subárbol «hacia abajo» añadiendo dos nodos
 - <constante> <operación> <subárbol anterior>

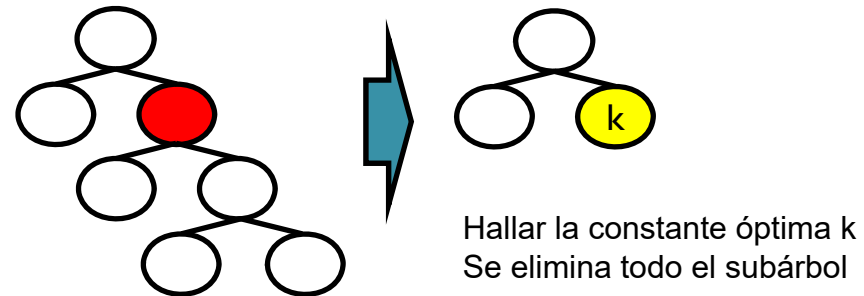
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo
 - Dado un nodo (marcado en rojo)

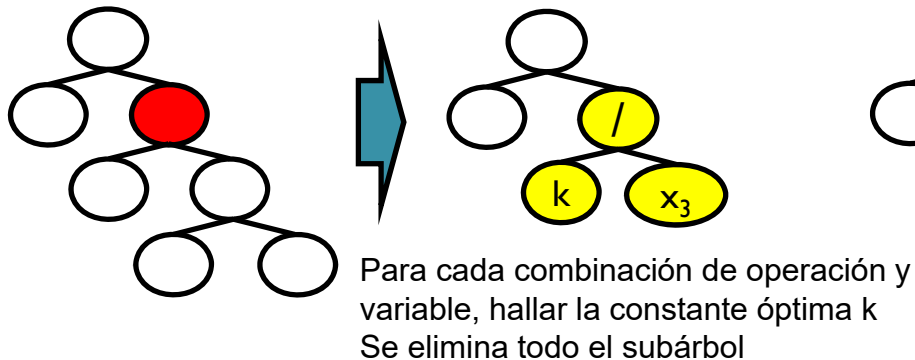
Variable search



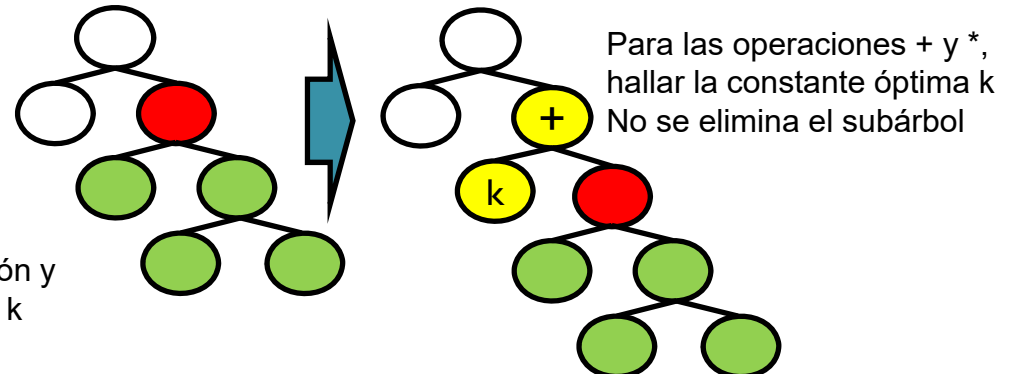
Constant search



Constant-variable search

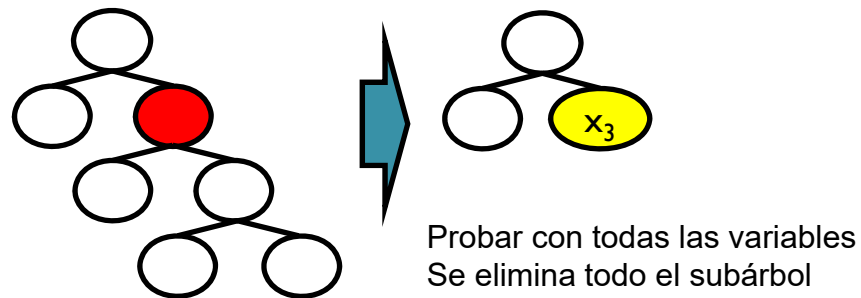


Constant-expression search



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *variable search*
 - Se busca sustituir el nodo por un nodo terminal con una variable

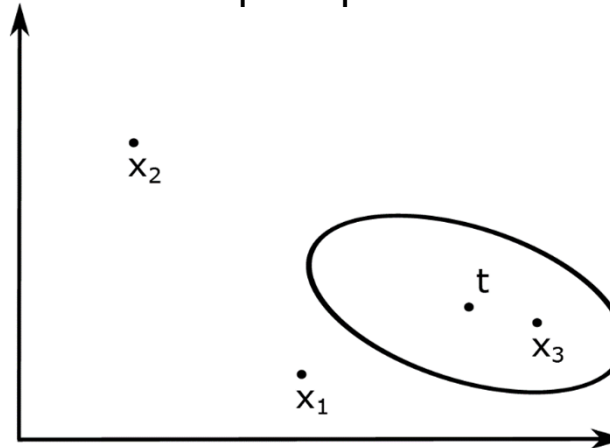


REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

- Búsqueda local en un nodo: *variable search*

- Se busca sustituir el nodo por un nodo terminal con una variable
 - Cada variable tiene una semántica que define un punto en el espacio
 - Si una está en el interior de la forma, se puede crear un nodo terminal con esa variable para sustituir el nodo de búsqueda por ese nodo terminal
 - Ejemplo: variable x_3



- Para saber esto, se usa la ecuación de MSE de ese nodo

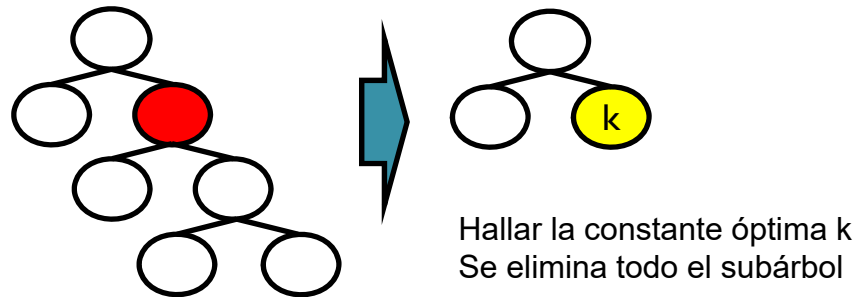
$$MSE = \frac{1}{N} \sum_{i=0}^N \left(\frac{a_i \cdot o_i - b_i}{c_i \cdot o_i - d_i} \right)^2$$

o_i : salida de esa variable para el patrón i

- Si el MSE es menor, se puede sustituir

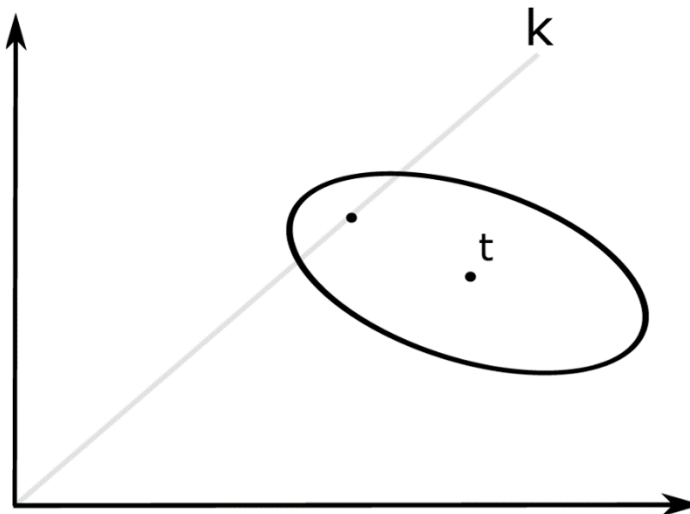
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant search*
 - Se busca un subárbol formado por un nodo terminal que contenga una constante k



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant search*
 - Se busca un subárbol formado por un nodo terminal que contenga una constante k
 - La semántica será (k, k, \dots, k) , es decir, $o_i = k$
 - Situada en la línea definida por el origen y el vector $(1, 1, \dots, 1)$
 - **El valor de k será el punto de esa línea más cercano a t**
 - Si este punto está en el interior de la forma correspondiente, el nodo terminal k puede sustituir al nodo donde se está produciendo la búsqueda



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant search*
 - Para hallar el mejor valor de k para ese nodo:

- Ecuación para calcular el MSE en ese nodo:

$$MSE = \frac{1}{N} \sum_{i=0}^N \left(\frac{a_i \cdot o_i - b_i}{c_i \cdot o_i - d_i} \right)^2 = \frac{1}{N} \sum_{i=0}^N \left(\frac{a_i \cdot k - b_i}{c_i \cdot k - d_i} \right)^2$$

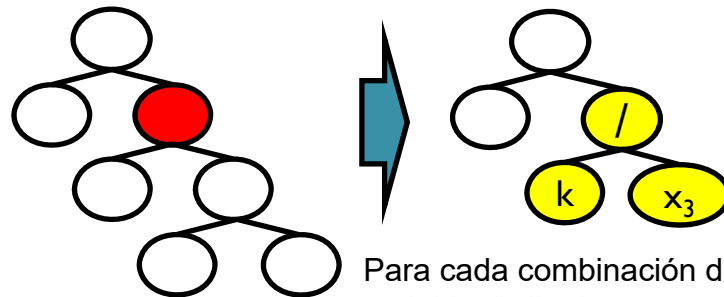
- Valor de k que minimice el MSE:
 - Se deriva esta expresión y se iguala a 0:

$$\frac{2}{N} \sum_{i=0}^N (b_i \cdot c_i - a_i \cdot d_i) \frac{a_i \cdot k - b_i}{(c_i \cdot k - d_i)^3} = 0$$

- Se despeja el valor de k
 - Existen expresiones analíticas para los casos más comunes

REGRESIÓN SIMBÓLICA

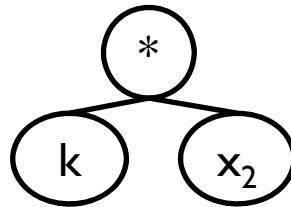
- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-variable search*
 - Se busca sustituir el nodo por un subárbol que consta de 3 nodos:
 - $\langle \text{constante (nodo terminal)} \rangle \langle \text{operación} \rangle \langle \text{variable (nodo terminal)} \rangle$



Para cada combinación de operación y variable, hallar la constante óptima k
Se elimina todo el subárbol

REGRESIÓN SIMBÓLICA

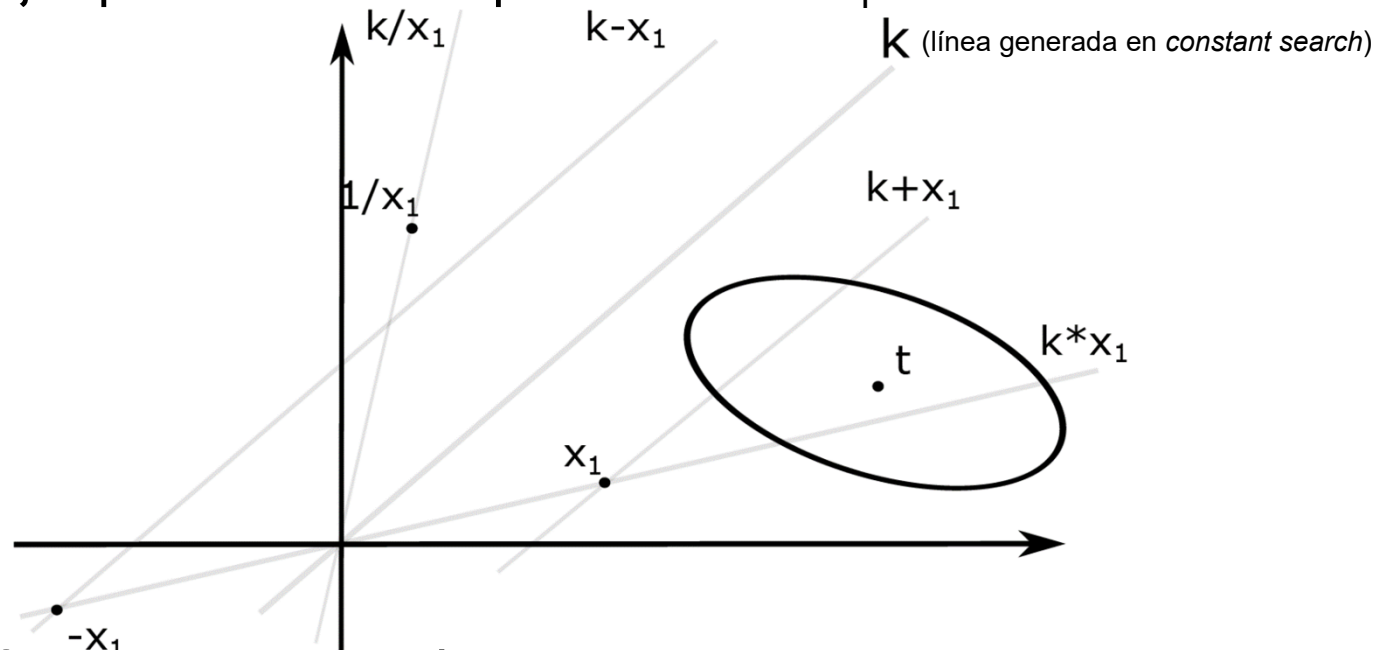
- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-variable search*
 - Se busca sustituir el nodo por un subárbol que consta de 3 nodos:
 - $\langle \text{constante (nodo terminal)} \rangle \langle \text{operación} \rangle \langle \text{variable (nodo terminal)} \rangle$
 - Por ejemplo:



- Se crean 4 líneas (una por operador) **para cada variable**, y para cada línea se escoge el punto que minimice el MSE
 - Si algún punto mejora el MSE, el nodo puede ser sustituido por el subárbol formado por estos 3 nodos
 - *Constant search* creaba una línea y escogía un punto de esa línea

REGRESIÓN SIMBÓLICA

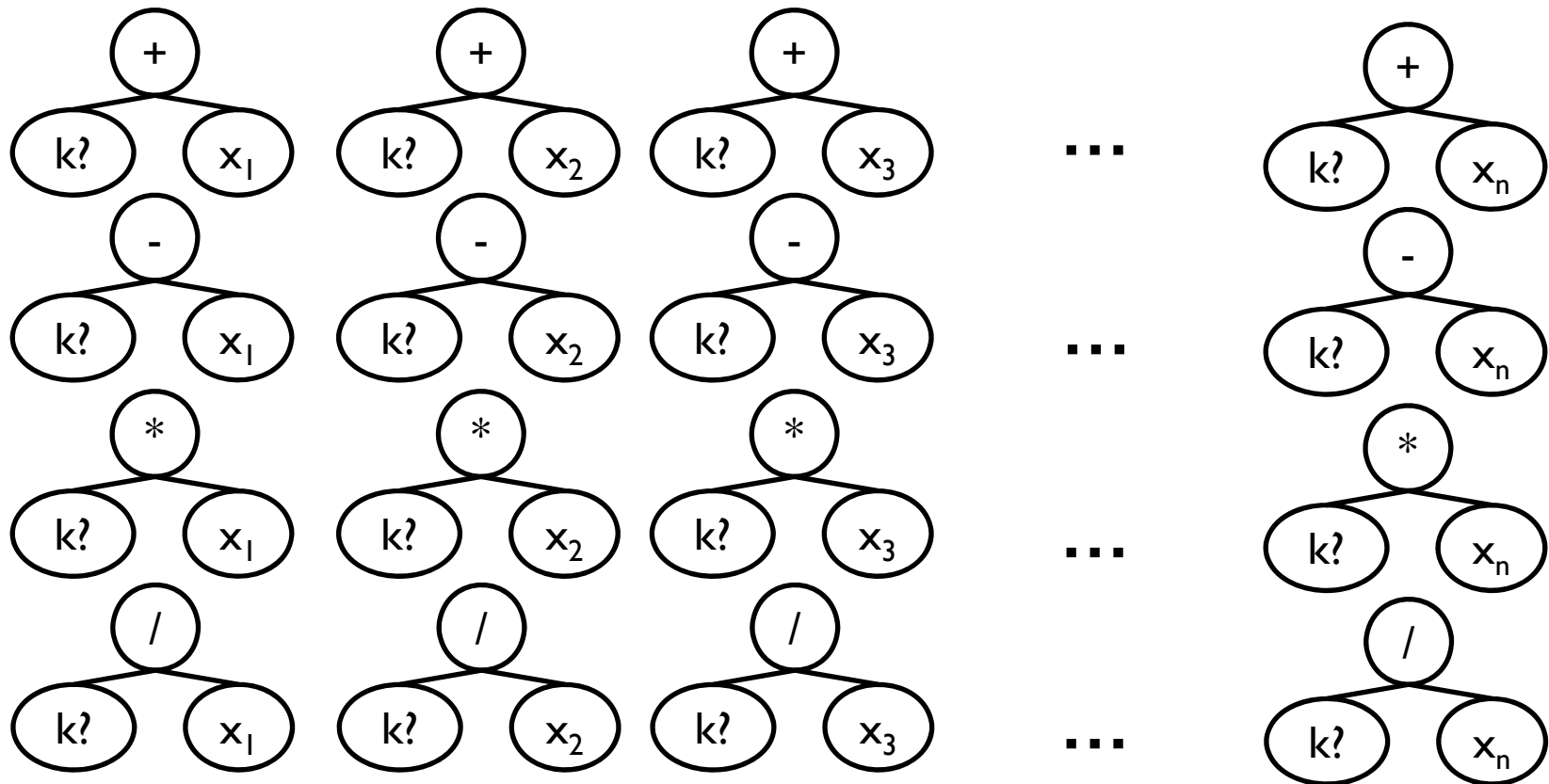
- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-variable search*
 - Ejemplo de las 4 líneas para la variable x_1 :



- Se busca el punto más cercano a t
 - El que minimice el MSE
- Se realiza este proceso para el resto de variables

REGRESIÓN SIMBÓLICA

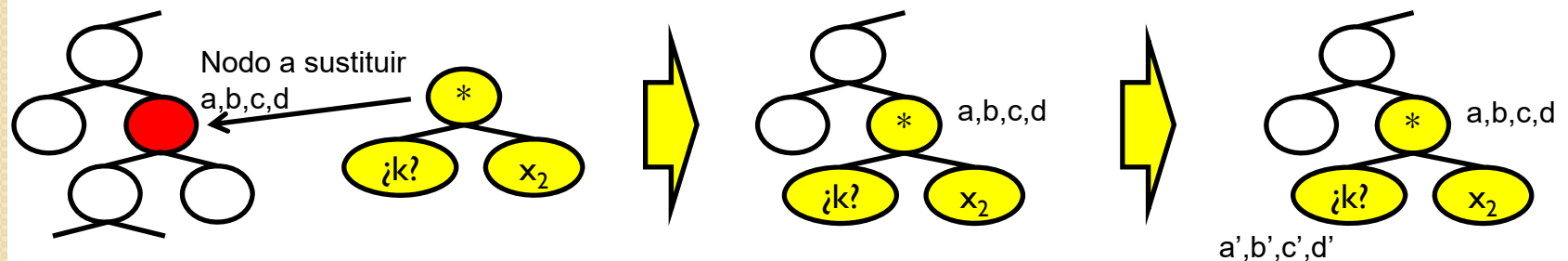
- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-variable search*
 - Para cada operación (+, -, *, /) y cada variable x_i , se construyen:



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-variable search*

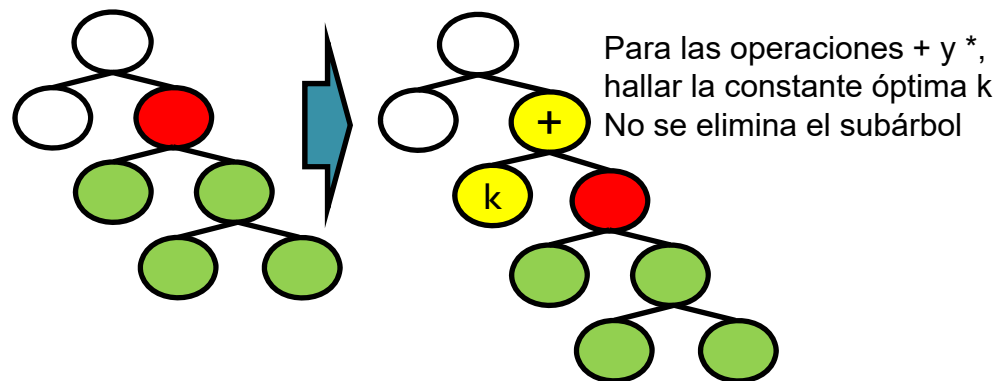
- Para cada uno de estos subárboles:
 - Se reemplaza temporalmente el nodo seleccionado por el nuevo subárbol
 - Se calcula la ecuación del nodo con la constante (es decir, se calcula a', b', c', d')



- Se calcula el mejor valor de k (el que minimiza el ECM)
 - Se realiza una *constant search*
- Con este valor de k , se calcula el ECM con el subárbol
- Si este nuevo ECM mejora el ECM anterior, el cambio se puede dejar permanente
 - Si no, se deshace el cambio
- Repetir el este proceso para cada uno de los nuevos subárboles creados

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-expression search*
 - Se busca sustituir el nodo por un subárbol que consta de:
 - Nodo padre: Operación
 - Sólo suma y multiplicación
 - Primer hijo: constante a hallar
 - Segundo hijo: el subárbol designado por el nodo seleccionado a ser sustituido
 - Tiene una semántica

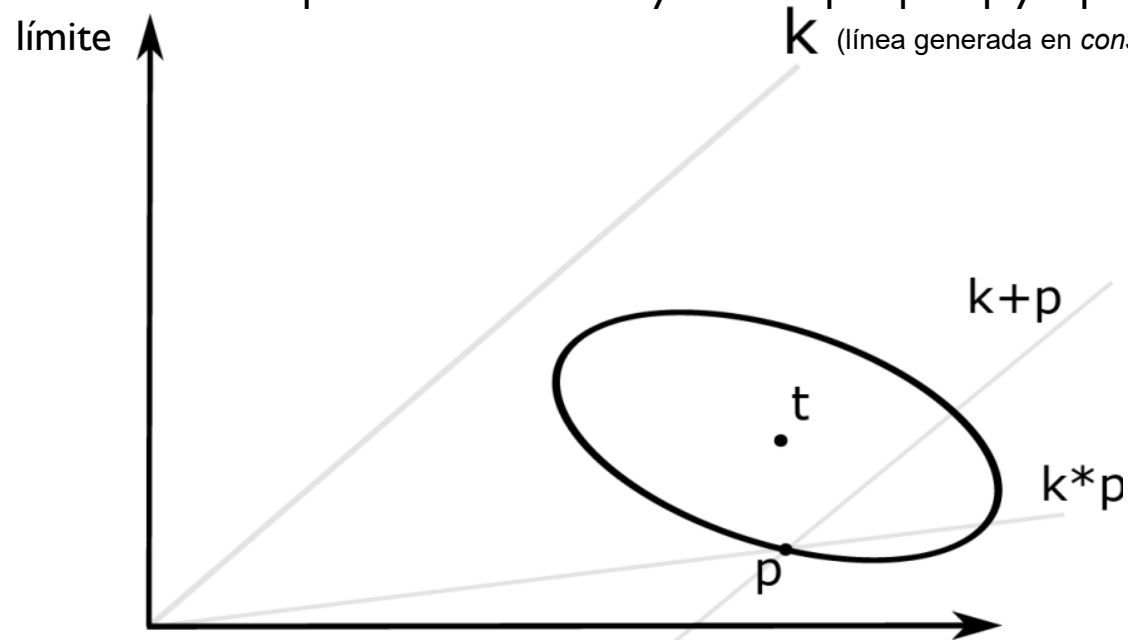


REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-expression search*
 - Se busca sustituir el nodo por un subárbol que consta de:
 - Nodo padre: Operación
 - Sólo suma y multiplicación
 - Primer hijo: constante a hallar
 - Segundo hijo: el subárbol designado por el nodo seleccionado a ser sustituido
 - Tiene una semántica
 - Se aprovecha la semántica del nodo para acercar el árbol al objetivo
 - **Esa semántica está en el límite de la forma**
 - Parece fácil desplazarla hacia el interior con una línea como las anteriores
 - No se elimina el subárbol, sino que se desplaza «hacia abajo» dentro del árbol

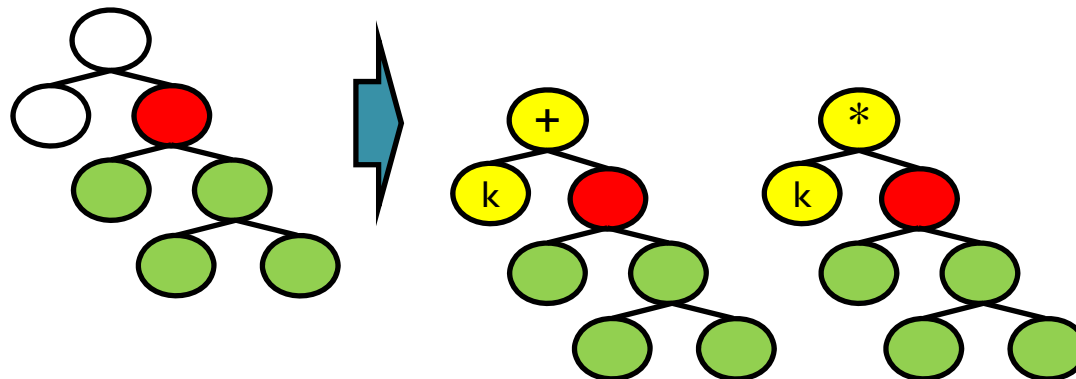
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-expression search*
 - Ejemplo:
 - Nodo a sustituir es el punto p (en el límite de la forma)
 - Se puede desplazar al interior mediante las operaciones de suma y multiplicación
 - Usando las constantes adecuadas
 - No se usan las operaciones de resta y división porque $-p$ y $1/p$ no están en el límite



REGRESIÓN SIMBÓLICA

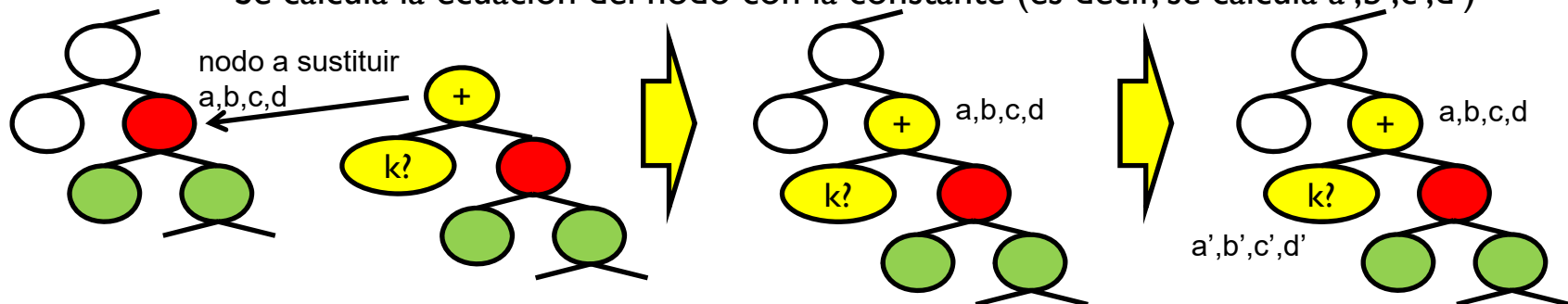
- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-expression search*
 - De forma similar a *constant-variable search*:
 - En primer lugar, construir estos dos subárboles:
 - Raíz: nodo no terminal con operación + o *
 - Primer hijo: constant k que determinar
 - Segundo hijo: nodo a reemplazar
 - y todo el subárbol que cuelga de él



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Búsqueda local en un nodo: *constant-expression search*

- Para cada uno de estos dos subárboles:
 - Se reemplaza temporalmente el nodo seleccionado por el nuevo subárbol
 - Se calcula la ecuación del nodo con la constante (es decir, se calcula a', b', c', d')



- Se calcula el mejor valor de k (el que minimiza el ECM)
 - Se realiza una *constant search*
- Con este valor de k , se calcula el ECM con el subárbol
- Si este nuevo ECM mejora el ECM anterior, el cambio se puede dejar permanente
 - Si no, se deshace el cambio
- Repetir el este proceso para cada uno de los nuevos subárboles creados

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Restricciones:
 - Número máximo de nodos
 - Limita la complejidad del modelo desarrollado
 - Valores muy altos:
 - MSE en entrenamiento muy bajo, pero posible *overfitting*
 - Valores muy bajos:
 - Buena generalización, pero posible *underfitting*
 - Limita la posibilidad de usar *constant-variable search* y *constant-expression search*
 - Pueden aumentar el número de nodos en 2 si se aplican en un nodo terminal
 - *Constant search* y *variable search* no aumentan el número de nodos

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Algoritmo:
 - Árbol inicial:
 - Formado por un nodo terminal con una constante
 - *Constant search* usando como vectores los propios de la raíz del árbol:
 - $a_i=1, b_i=t_i, c_i=0, d_i=-1$
 - A partir de ese árbol, se inicia un proceso iterativo
 - En cada iteración, se sustituye un nodo por otro nodo o subárbol si se mejora el MSE
 - Resultado de una de las 4 búsquedas locales en algún nodo
 - ¿Cómo saber en qué nodos buscar y qué búsquedas hacer en cada iteración?
 - Recorrer el árbol aplicando las búsquedas siguiente una **estrategia**
 - Ejemplos:
 - Exhaustiva
 - Selectiva

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Estrategias: Exhaustiva:
 - Recorrer todos los nodos, aplicando las 4 búsquedas en cada nodo
 - Reemplazar el nodo con el resultado de la estrategia que disminuya más el MSE
 - Computacionalmente costoso
 - Muchos cálculos que no van a llevar a mejoras
 - Se pueden ahorrar muchos cálculos realizando las búsquedas locales en los nodos donde suelen tener éxito
 - Realizar las búsquedas en el resto de los nodos sólo si las anteriores no tuvieron éxito
 - Estrategia selectiva

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Estrategias: Selectiva:
 - En cada iteración, realizar los siguientes pasos. Si alguno tiene éxito, aplicar el cambio correspondiente y comenzar la siguiente iteración:
 1. Realizar *constant search* en los nodos terminales
 2. Realizar *variable search* en los nodos terminales
 3. Realizar *constant-expression search* en los nodos no terminales
 4. Realizar *constant-variable search* en los nodos terminales
 5. Realizar de forma conjunta:
 - *Constant search* en nodos terminales con variables y nodos no terminales
 - *Variable search* en nodos terminales con variables y nodos no terminales
 - *Constant-variable search* en nodos no terminales

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Algoritmo:
 - Criterios de parada:
 - Se ha alcanzado el valor de MSE objetivo
 - En la iteración actual la estrategia no ha tenido éxito
 - No se han encontrado cambios en el árbol que mejoren el MSE en una cantidad determinada
 - Hiperparámetros:
 - Número máximo de nodos
 - Estrategia a utilizar
 - Valor de MSE a alcanzar. Por defecto: 0
 - Mejora mínima en el MSE para que una búsqueda tenga éxito
 - Multiplicada por el valor del MSE actual
 - Por ejemplo, un valor de 10^{-2} implica que es necesaria una mejorar el MSE como mínimo en un 1%
 - Valores típicos: 10^{-5} , 10^{-6} , 10^{-7}

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Ejemplo: Ley de la Gravitación Universal de Newton:

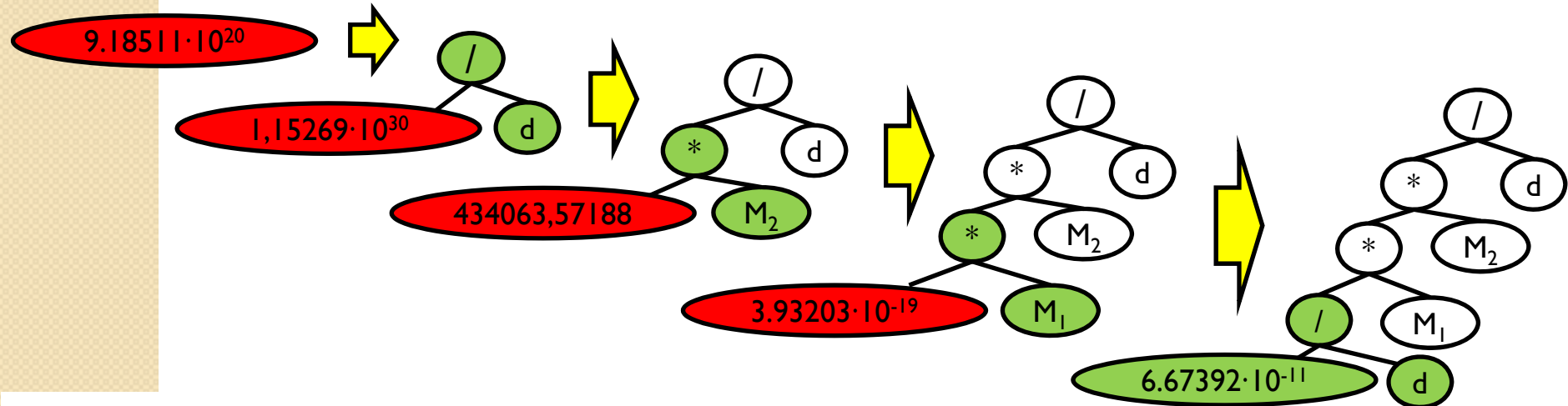
$$F = G \cdot \frac{M_1 \cdot M_2}{d^2}$$

- $G = 6.67392 \cdot 10^{-11}$
- Base de datos artificial
 - 1000 instancias, cada una con:
 - M_1 : masa del primer planeta
 - Valor aleatorio entre 10^{23} y 10^{25}
 - M_2 : masa del segundo planeta
 - Valor aleatorio entre 10^{23} y 10^{25}
 - d : distancia entre ambos planetas
 - Valor aleatorio
 - Target: resultado de la ecuación anterior

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Ejemplo: Ley de la Gravitación Universal de Newton:

• Resultados:



Iteration	Node selected for substitution	Resulting expression	MSE
0	-	$9.185106275464827 \cdot 10^{20}$	$2.7645 \cdot 10^{43}$
1	$9.185106275464827 \cdot 10^{20}$	$(1.1526861538137104 \cdot 10^{30}/d)$	$2.1717 \cdot 10^{43}$
2	$1.1526861538137104 \cdot 10^{30}$	$((434063.57187533064 \cdot M_2)/d)$	$1.9183 \cdot 10^{43}$
3	434063.57187533064	$((((3.93202929320367 \cdot 10^{-19} \cdot M_1) \cdot M_2)/d)$	$5.1733 \cdot 10^{42}$
4	$3.93202929320367 \cdot 10^{19}$	$(((((6.6739200000000007 \cdot 10^{-11}/d) \cdot M_1) \cdot M_2)/d)$	$3.1828 \cdot 10^{13}$

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:
 - 32 datasets del repositorio PMLB
 - *Penn Machine Learning Benchmark*
 - Resultados en la tabla:
 - Mejores configuraciones para cada problema
 - Estrategia
 - Num. nodos
 - Reducción mínima de MSE
 - Tiempos de ejecución
 - Estabilidad

REGRESIÓN SIMBÓLICA

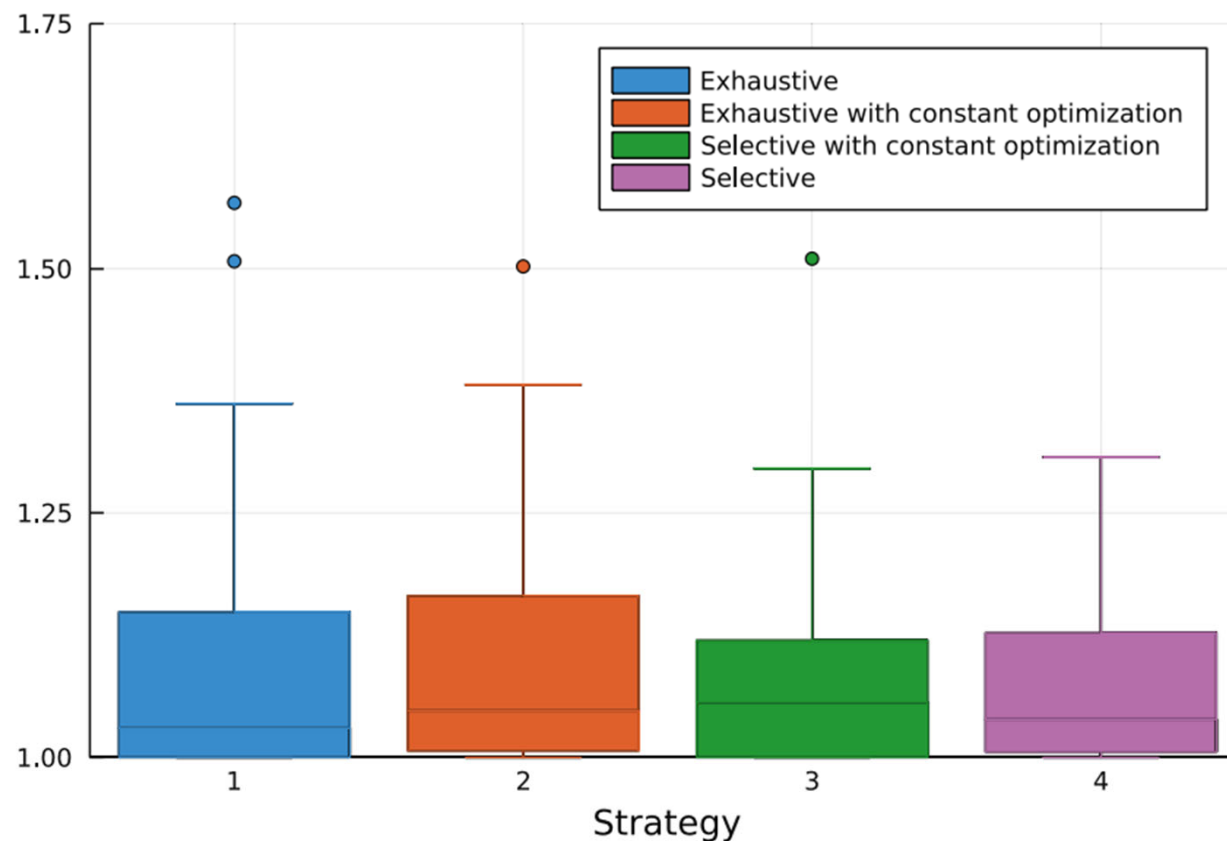
- Development of Mathematical Expressions (DoME)
 - Resultados:

Results obtained. For each problem, this table shows the best median MSE obtained in test, the hyperparameter configuration (strategy used to go through the nodes, minimum improvement in MSE for a search to be considered successful and maximum number of nodes in the tree) that returned this value, the median training time of this configuration (measured in seconds), and a stability measure of this configuration.

Problem	Best test <i>MSE</i>	Strategy	Min MSE reduction	Maximum num nodes	Average time (s)	Stability (100%)
ESL	0.254	2	10^{-5}	100	8.3	4.35%
SWD	0.368	1	10^{-5}	60	4.15	1.49%
LEV	0.35	3	10^{-4}	45	0.0356	0.89%
ERA	2.21	2	10^{-7}	100	1.49	2.03%
USCrime	193	3	10^{-3}	20	0.0124	2.48%
FacultySalaries	0.728	1	10^{-7}	155	165	52.44%
vineyard	4.1	4	10^{-3}	15	0.00136	16.33%
auto price	3.43e6	1	10^{-7}	65	0.978	19.50%
autoPrice	3.41e6	2	10^{-4}	160	7.94	17.61%
cloud	0.0676	4	10^{-7}	55	0.0903	27.87%
elusage	89.8	1	10^{-3}	15	0.00706	2.81%
machine cpu	1.27e3	4	10^{-7}	70	8.41	5.84%
analcata data vehicle	5.32e3	4	10^{-5}	145	2.13	51.20%
vinnie	2.36	3	10^{-4}	5	0.000948	1.41%
pm10	0.545	2	10^{-4}	30	0.677	1.57%
analcata data neavote	0.814	3	10^{-2}	10	0.00081	12.57%
analcata data election2000	1.88e7	3	10^{-5}	100	1.91	99.83%
pollution	1.26e3	4	10^{-4}	30	0.0186	5.40%
no2	0.242	1	10^{-4}	50	3.4	3.23%
analcata data apnea2	5.58e5	4	10^{-4}	115	0.749	9.53%
analcata data apnea1	6.97e5	1	10^{-6}	125	5.99	8.98%
cpu	15.2	3	10^{-7}	140	140	79.49%
sleuth ex1714	7.66e5	3	10^{-2}	40	0.0691	13.25%
rabe 266	2.65	1	10^{-7}	185	262	34.92%
sleuth case2002	52.7	3	10^{-6}	30	0.109	2.81%
rmftsa ladata	2.81	1	10^{-3}	70	0.211	8.50%
visualizing environmental	6.82	3	10^{-5}	30	0.0472	6.63%
sleuth ex1605	85.2	2	10^{-7}	30	0.104	8.01%
visualizing galaxy	384	1	10^{-7}	75	19.3	8.13%
chatfield 4	212	3	10^{-6}	30	0.319	5.03%
sleuth case1202	1.56e3	3	10^{-3}	70	0.0875	11.89%
chscase geyser1	33.3	1	10^{-3}	45	0.0219	2.11%

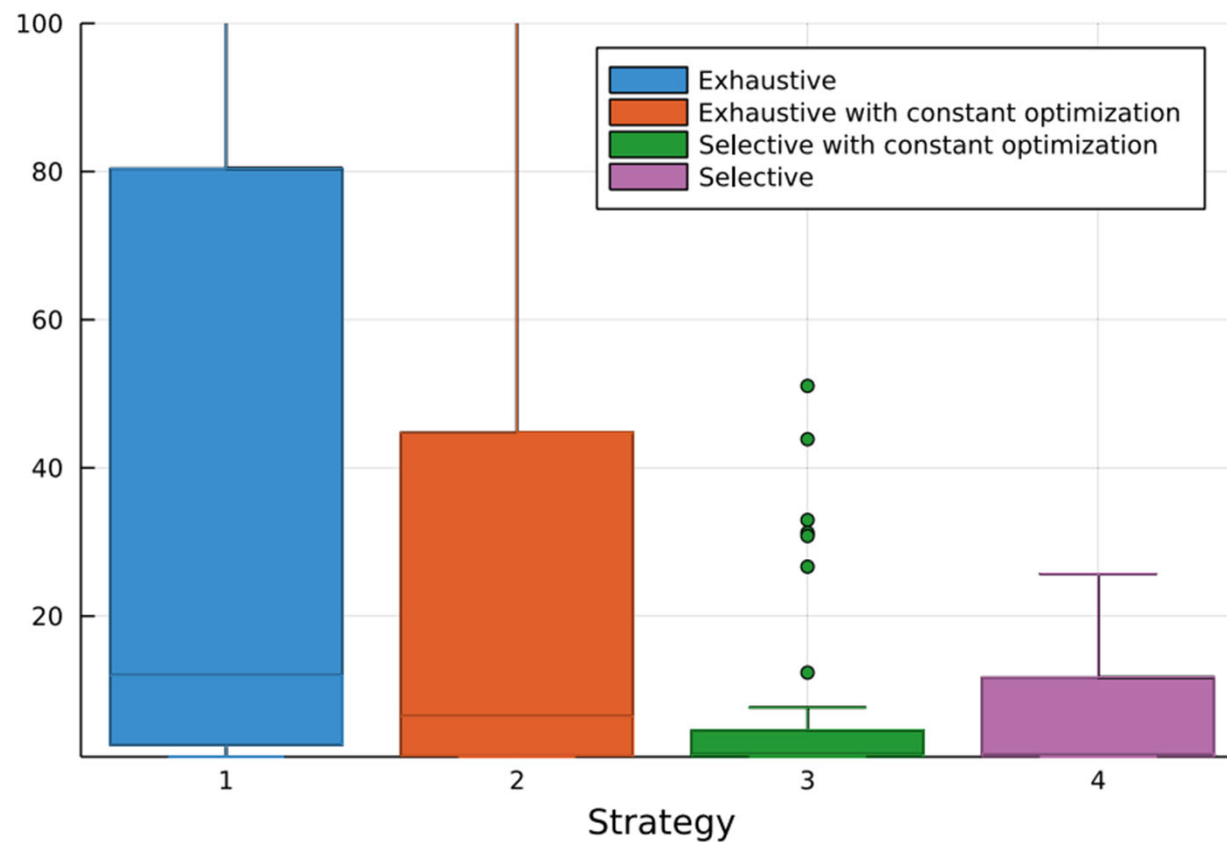
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:
 - ECM relativo entre las 4 estrategias:



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:
 - Tiempo de entrenamiento relativo entre las 4 estrategias:



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:
 - Comparación con el trabajo
 - Moscato, P., Sun, H., & Haque, M. N. (2021). Analytic continued fractions for regression: A memetic algorithm approach. Expert Systems with Applications, 179, Article 115018
 - <http://dx.doi.org/10.1016/j.eswa.2021.115018>
 - <https://www.sciencedirect.com/science/article/pii/S0957417421004590>
 - Comparativa de distintas técnicas de regression: Basados en PG:
 - afp: Age-fitness Pareto Optimization
 - eplex: ϵ -Lexicase selection
 - eplex-Im: variación de eplex con criterio de parada de un millón de evaluaciones (tamaño de población x generaciones)
 - gsgp: Geometric Semantic Genetic Programming
 - mrgp: Multiple Regression Genetic Programming

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:
 - Comparación con el trabajo
 - Moscato, P., Sun, H., & Haque, M. N. (2021). Analytic continued fractions for regression: A memetic algorithm approach. Expert Systems with Applications, 179, Article 115018
 - Comparativa de distintas técnicas de regression: Algoritmos de AA:
 - ada-b: Adaptive Boosting (AdaBoost) Regression
 - grad-b: Gradient Boosting Regression
 - krnl-r: Kernel Ridge
 - lasso-l: Least-Angle Regression con Lasso
 - l-regr: linear-regression: Linear Regression
 - l-svr: linear-svr: Linear Support Vector Regression
 - mlp: Multilayer Perceptrons (MLPs) Regressor
 - rf: Random Forests Regression
 - sgd-r: Stochastic Gradient Descent Regression
 - xg-b: Extreme Gradient Boosting (xgboost)

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)

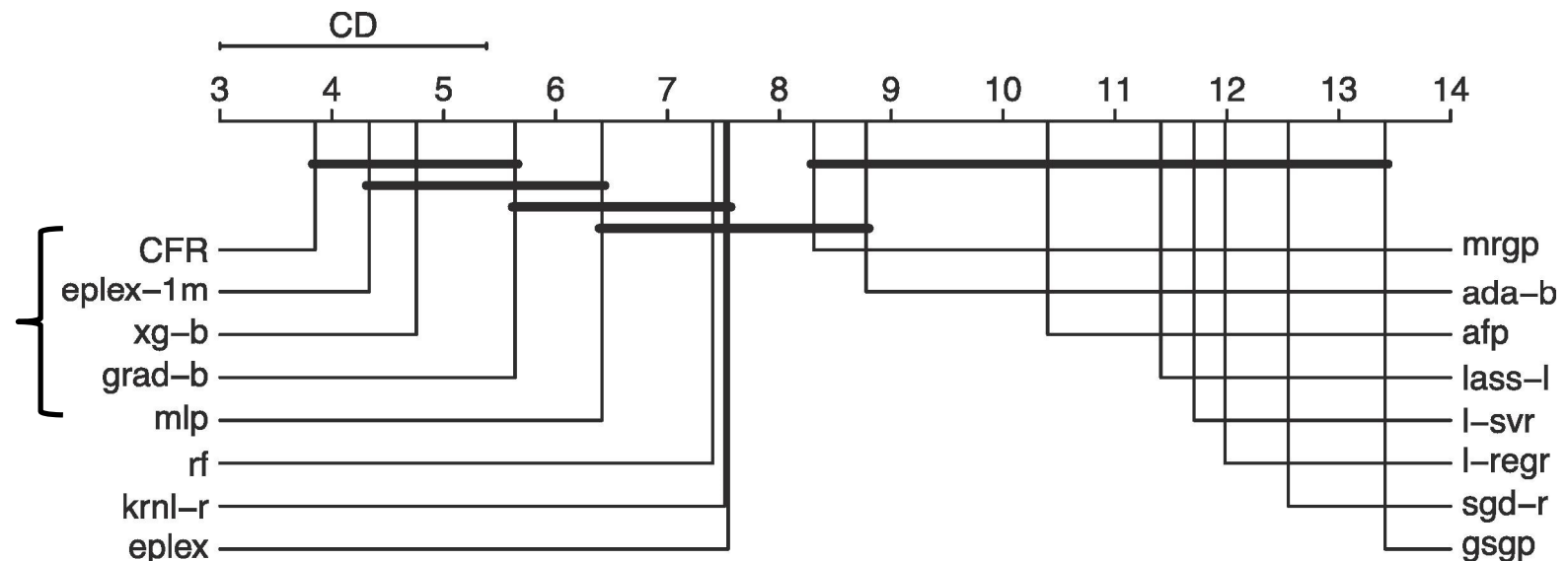
- Resultados:

- Comparación con el trabajo

- Moscato, P., Sun, H., & Haque, M. N. (2021). Analytic continued fractions for regression: A memetic algorithm approach. Expert Systems with Applications, 179, Article 115018

- <http://dx.doi.org/10.1016/j.eswa.2021.115018>

- <https://www.sciencedirect.com/science/article/pii/S0957417421004590>



REGRESIÓN SIMBÓLICA

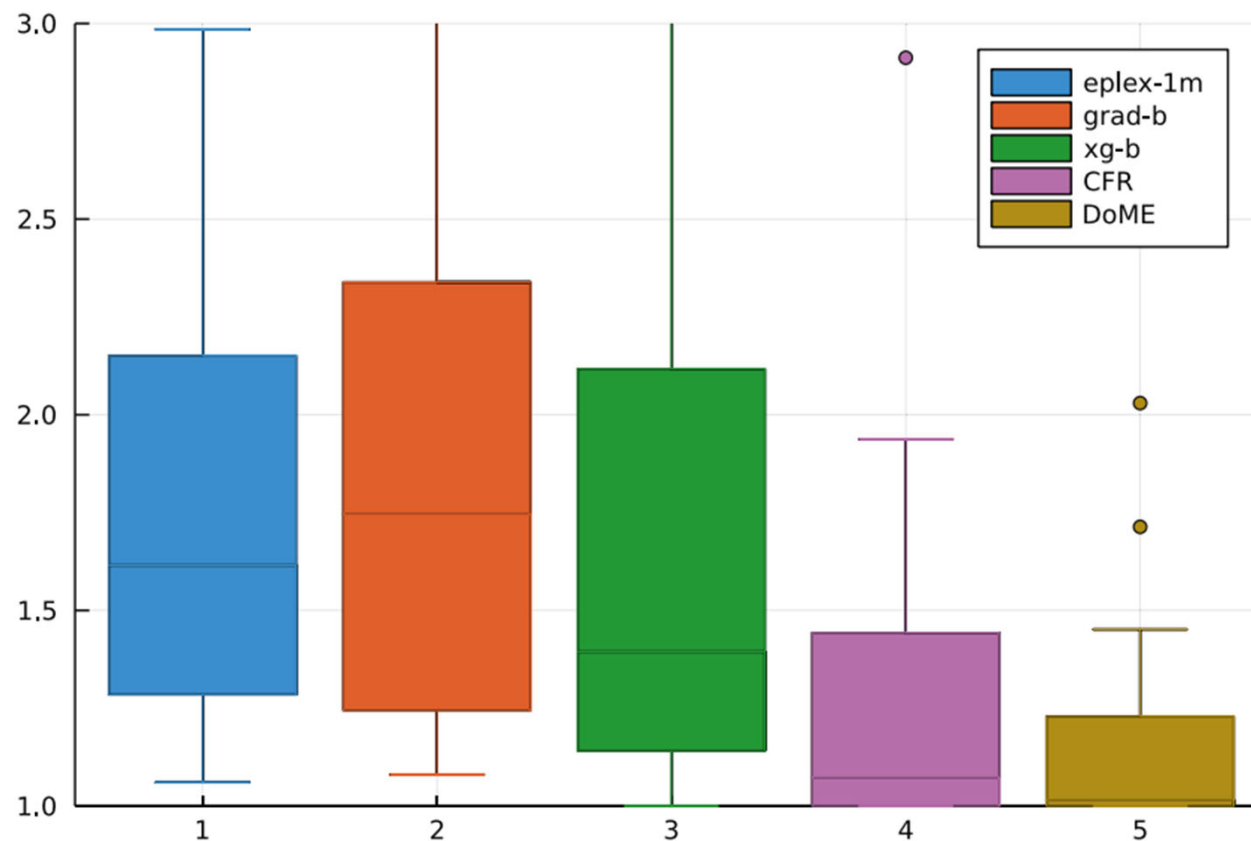
- DoME
 - Resultados:
 - Comparando DoME con las 4 mejores técnicas en esos mismos datasets:

Comparison between different methods (*MSE*)

Problem	eplex-1 m	grad-b	xg-b	CFR	DoME
ESL	0.274	0.319	0.272	0.268	0.254
SWD	0.39	0.405	0.408	0.432	0.368
LEV	0.425	0.424	0.422	0.353	0.35
ERA	2.51	2.58	2.57	2.45	2.21
USCrime	393	257	378	220	193
FacultySalaries	4.04	8.07	4.11	1.28	0.728
vineyard	6.01	8.22	7.83	4.22	4.1
auto price	5.89e6	3.89e6	4.03e6	6.01e6	3.43e6
autoPrice	4.17e6	5.29e6	2.87e6	4.83e6	3.41e6
cloud	0.11	0.208	0.144	0.095	0.0676
elusage	135	199	137	65.8	89.8
machine cpu	3.8e3	2.23e3	2.69e3	2.1e3	1.27e3
analcata data vehicle	4.14e4	2.41e4	4.2e4	1.55e4	5.32e3
vinnie	2.29	2.86	2.66	1.93	2.36
pm10	0.64	0.431	0.399	0.621	0.545
analcata data neavote	1.18	0.818	0.917	0.401	0.814
analcata data election2000	4.33e7	3.4e8	7.72e8	5.09e5	1.88e7
pollution	1.87e3	2.19e3	1.67e3	1.42e3	1.26e3
no2	0.272	0.227	0.21	0.295	0.242
analcata data apnea2	1.12e6	9.42e5	7.86e5	6.09e5	5.58e5
analcata data apnea1	8.16e5	9.98e5	5.28e5	6.96e5	6.97e5
cpu	175	2.36e3	883	164	15.2
sleuth ex1714	1.42e6	1.57e6	2.31e6	6.83e5	7.66e5
rabe 266	7.11	7.32	3.04	2.64	2.65
sleuth case2002	75.8	56.2	72.4	41.7	52.7
rmftsa ladata	3.01	3.51	3.21	2.76	2.81
visualizing environmental	9.62	9.8	9.54	4.7	6.82
sleuth ex1605	102	98.4	92	84	85.2
visualizing galaxy	313	268	224	434	384
chatfield 4	282	384	288	189	212
sleuth case1202	3.29e3	3.42e3	3.2e3	1.39e3	1.56e3
chscase geyser1	38.9	39.9	42.9	31.1	33.3

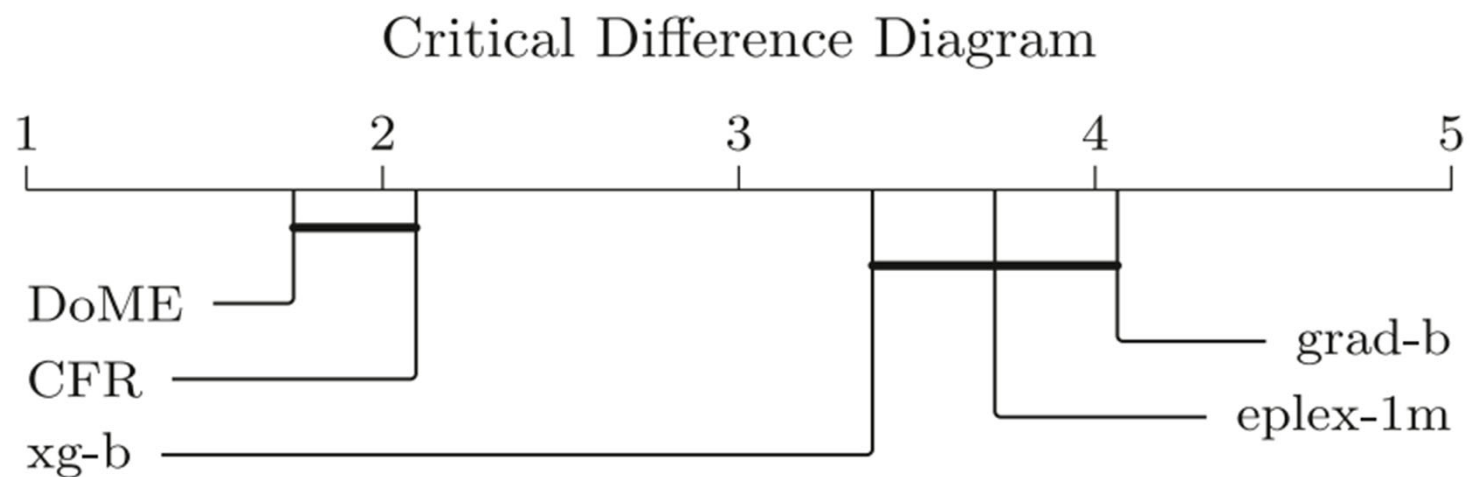
REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:
 - ECM promedio relativo entre las distintas técnicas



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:
 - Diagrama de diferencia crítica con la comparativa entre las distintas técnicas



REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:

Summary of the expressions obtained.

Problem	Nodes	Height	Num. inputs	Problem	Nodes	Height	Num. inputs
ESL	99	17	4	analcata election2000	99	16	14
SWD	59	14	10	pollution	29	12	15
LEV	45	13.5	4	no2	49	12.5	7
ERA	78	18.5	4	analcata apnea2	115	23	3
USCrime	19	7	13	analcata apnea1	125	15.5	3
FacultySalaries	155	20	4	cpu	139	24	7
vineyard	9	4	2	sleuth ex1714	39	8.5	7
auto price	65	16.5	15	rabe 266	185	46.5	2
autoPrice	159	25.5	15	sleuth case2002	29	11	6
cloud	32	9	5	rmftsa ladata	15	7	10
elusage	15	5	2	visualizing environmental	24	9	3
machine cpu	69	15.5	6	sleuth ex1605	29	11	5
analcata vehicle	145	24	4	visualizing galaxy	75	13.5	4
vinnie	5	3	2	chatfield 4	29	10	12
pm10	15	6	7	sleuth case1202	69	13.5	6
analcata neavote	9	5	2	chscase geyser1	23	8	2

REGRESIÓN SIMBÓLICA

- Development of Mathematical Expressions (DoME)
 - Resultados:

Examples of the final expressions obtained.

Problem	Expression
LEV	$(-0.122) + (((((0.265 \cdot X_1) + (0.416 \cdot X_2)) + (0.152 \cdot X_3)) + (0.155 \cdot X_4))$
USCrime	$\left(\left(\frac{(-0.822)}{X_{10}} + \frac{\left(X_9 + \frac{648854.722}{X_2} \right)}{X_6} \right) \cdot X_1 \right) + \frac{72690.856}{X_{13}}$
vineyard	$0.968 \cdot ((5.663 + X_1) + X_2)$
cloud	$((-0.147) \cdot X_5) + ((0.479 \cdot X_3) + ((-0.047) + (0.733 \cdot X_4)))$
elusage	$(((-0.002) \cdot X_1) \cdot X_1) + \frac{((2362.397 - X_1) - X_1)}{((-2.965) + X_1)}$
vinnie	$(-1.584) + (0.582 \cdot X_2)$
pm10	$0.351 \cdot \left(\frac{13.246}{((0.442 \cdot X_1) + (1.155 \cdot X_3))} + (1.021 \cdot X_1) \right)$
analcata neavote	$3.732 \cdot (2.602 - ((3.022 - X_1) \cdot X_1))$
visualizing environmental	$4.346 \cdot \left(1.335 + \frac{(112.678 - (1.017 \cdot X_1))}{X_3} \right)$
sleuth ex1605	$0.761 \cdot \left(\left(\left(\frac{31.767}{X_5} \cdot X_4 \right) \cdot \left(\frac{(((\frac{327.570 \cdot X_5}{X_2} - X_2) - X_2) - X_2)}{((-95.548) + X_4)} + X_2 \right) \right) + X_5 \right)$
chscase geyser1	$((((31.425 - (0.923 \cdot (0.984 \cdot X_2))) - X_2) - X_2) \cdot \left(\frac{(2.508 - (((0.002 \cdot X_2) \cdot X_2) \cdot X_1) + X_2))}{X_2} + X_2 \right)$

REGRESIÓN SIMBÓLICA

- Ventajas de la Regresión Simbólica
 - Al contrario que las técnicas de regresión clásicas como RR.NN.AA. o SVR (SVM para Regresión), el resultado es una ecuación matemática explícita. Esto conlleva una serie de ventajas importantes:
 - Modelo explicable
 - **IA explicable y confiable**
 - Ofrece información sobre los datos
 - Permite ser analizada para obtener nuevo conocimiento
 - Permite descubrir relaciones ocultas entre los datos
 - Modelo muy portátil: no es necesario implementarlo en ningún lenguaje ni cargar librerías para utilizarlo
 - Caso muy habitual: se puede usar fácilmente en una hoja de cálculo

REGRESIÓN SIMBÓLICA

- Ventajas de la Regresión Simbólica
 - Muchas técnicas intentan reducir el tamaño de las expresiones que devuelve
 - Modelos más sencillos → menos sobreajuste
 - Los resultados no son peores a los obtenidos por el resto de las técnicas de Regresión como RR.NN.AA. o SVR (SVM para Regresión)
 - Ninguna técnica de Regresión va a ser mejor que el resto en promedio
 - Teoremas No Free Lunch (NFL)
 - Los resultados parecen indicar que las técnicas de Regresión Simbólica ofrecen mejores resultados que las técnicas de regresión clásicas
 - Pero muchas veces llevan a sobreajustar → Importante limitar la complejidad

REGRESIÓN SIMBÓLICA

- Desventajas de la Regresión Simbólica
 - Espacio de búsqueda mucho mayor que el resto de técnicas de Regresión
 - Formado por todas las posibles ecuaciones matemáticas que aproximen una relación entrada/salida
 - En otras técnicas como RR.NN.AA. o SVR (SVM para Regresión) se parte de una arquitectura fijada de antemano, con lo que el número de parámetros es fijo
 - Este espacio de búsqueda mucho mayor hace que el proceso de búsqueda sea generalmente más lento
 - Sobre todo en las técnicas basadas en poblaciones (como Programación Genética y técnicas derivadas de ella)
 - Se desperdicia mucho tiempo evaluando individuos que no van a contribuir en la solución final

REGRESIÓN SIMBÓLICA

- Ventajas de DoME
 - Al contrario que las técnicas de Regresión basadas en poblaciones, esta tiene una base matemática que explica su funcionamiento y convergencia
 - Más rápida que las técnicas basadas en poblaciones
 - Se puede limitar explícitamente la complejidad de las expresiones que devuelve
 - Controlar el sobreajuste
 - Número muy bajo de hiperparámetros que controlan el proceso
 - Fácil experimentar

REGRESIÓN SIMBÓLICA

- Desventajas de DoME
 - En cada iteración se realizan búsquedas en muchos nodos pero sólo se modifica uno
 - Se realizan muchos cálculos que no llevan a ninguna modificación
 - Esto se acentúa conforme el árbol se va haciendo más grande
 - Aún así, es más eficiente que las técnicas basadas en poblaciones
 - En general, que el resto de técnicas de Regresión Simbólica
 - Estos cálculos innecesarios se pueden controlar escogiendo una estrategia adecuada para recorrer los nodos del árbol
 - Por ahora la formulación sólo permite usar operadores aritméticos
 - Apropriados para problemas del mundo real, pero pueden tener limitaciones en otros tipos de problemas donde los datos tengan una naturaleza donde claramente se necesite alguna función determinada (por ejemplo, con datos periódicos)