

Lenguajes de Programación, 2019-I

Nota de clase 8: Tipos de Datos Estructurados^{*}

Favio E. Miranda Perea Lourdes del Carmen González Huesca
Facultad de Ciencias UNAM

1 de octubre de 2018

Resulta conveniente agregar a nuestro lenguaje nucleo funcional `MinHs` algunos tipos adicionales como tuplas, listas ó árboles. Estas estructuras o tipos de datos se pueden clasificar en tres categorías.

- Tipos producto: son aquellos cuyos valores son tuplas de valores donde cada componente es de un tipo específico. El caso base es el de tuplas con 2 elementos o pares ordenados.
- Tipos suma: también conocidos como uniones (disjuntas), son aquellos cuyos valores son valores de alguno de n tipos específicos con una etiqueta explícita indicando de cual de estos n tipos se trata. El caso base es el de una unión disjunta de dos tipos.
- Tipos recursivos: estos son tipos auto-referenciales en el sentido de que sus valores pueden definirse mediante otros valores del mismo tipo recursivo. Por ejemplo las listas y los árboles.

En nuestro curso nos dedicaremos a las sumas y los productos que proporcionan estructuras fundamentales básicas para definir tipos registro (records) y tipos variante (opciones). Ambas nociones guardan una dualidad encontrada en diversas áreas. Por ejemplo:

- Lógica: conjunción y disyunción.
- Lenguajes formales: concatenación y unión.
- Programación: secuencias y alternativas en instrucciones.
- Definiciones inductivas: varias premisas y varias reglas.

Iniciamos con los tipos producto.

1. Productos

Con frecuencia un lenguaje de programación ofrece una manera de fusionar elementos de dos o más tipos mediante el uso de un nuevo tipos. Por ejemplo dos valores flotantes representan un punto en el espacio, etc.

En Haskell por ejemplo tenemos:

^{*}Estas notas se basan en los libros de Harper y Pierce y en material de Frank Pfenning, Ralf Hinze y Gerwin Klein.

- Tuplas:
 - Declaración: `type Punto = (Float,Float)`
 - Uso: `let x = (1.75,1.21) in fst x`
- Tuplas con campos nombrados:
 - Declaración: `data Punto = Punto x::Float, y::Float`
 - Uso: `let p=Point 5.0 10.1 in...`, `let p = Point x=5.0 y=10.1 in...`
 se generan funciones selectoras `x,y:: Punto ->Float` de manera automática, por ejemplo `xp` devuelve el valor de `x` en el punto `p`.

Similarmente en C podemos definir:

```
struct point {
    float x;
    float y;
};
```

```
struct point p;
p.x = 10.1;
```

Mientras que en Java podemos definir clases degeneradas que no contienen métodos:

```
class Point {
    public float x;
    public float y;
}
```

O, como es más usual, proteger los campos de datos y dando métodos para generar, acceder y alterar estructuras:

```
class Point {
    private float x;
    private float y;

    public Point (float x, float y) {
        this.x = x;
        this.y = y;
    }

    public float getX () {return x;}
    public float getY () {return y;}

    public void setX (float x) {this.x = x;}
    public void setY (float y) {this.y = y;}
}
```

1.1. Tipos producto y sus constructores

Los tipos producto son nuestro primer ejemplo de un tipo estructurado. Un tipo producto $T \times S$ representa al tipo de parejas de valores cuyo primer componente es de tipo T y el segundo de tipo S . Como ya es costumbre presentamos aquí únicamente los nuevos constructores del lenguaje. Los dos niveles de sintaxis se resumen como sigue:

Sintaxis concreta	Sintaxis abstracta
<code>Void, 1, ()</code>	<code>Void</code>
<code>void,()</code>	<code>void</code>
<code>$T \times S$</code>	<code>prod(T, S)</code>
<code>(e_1, e_2)</code>	<code>pair(t_1, t_2)</code>
<code>fst e</code>	<code>fst(t)</code>
<code>snd e</code>	<code>snd(t)</code>

Distinguimos tres categorías:

- Tipos: `Void`, `1` o `()` para tuplas sin elementos y $T \times S$ para pares ordenados.
- Constructores: `()` o `void` es la única tupla sin elementos y (e_1, e_2) es un par ordenado.
- Selectores o destructores: las proyecciones `fst e` , `snd e` devuelven el primer y segundo elemento del par ordenado e respectivamente.

Es importante aclarar que el nombre `void` es usual en lenguajes de programación aunque `void` significa vacío y este tipo no es vacío sino que tiene un único elemento.

La semántica estática está dada por las siguientes reglas:

$$\overline{\Gamma \vdash () : \text{Void}}$$

$$\frac{\Gamma \vdash e_1 : T \quad \Gamma \vdash e_2 : S}{\Gamma \vdash \text{pair}(e_1, e_2) : \text{prod}(T, S)}$$

$$\frac{\Gamma \vdash e : \text{prod}(T, S)}{\Gamma \vdash \text{fst}(e) : T} \quad \frac{\Gamma \vdash e : \text{prod}(T, S)}{\Gamma \vdash \text{snd}(e) : S}$$

Como ya es usual usaremos con frecuencia la sintaxis concreta, por ejemplo usando la regla para tipar pares como sigue:

$$\frac{\Gamma \vdash e_1 : T \quad \Gamma \vdash e_2 : S}{\Gamma \vdash (e_1, e_2) : T \times S}$$

Para la estrategia de evaluación de la semántica dinámica definimos los siguientes nuevos valores:

$$v ::= \dots \mid () \mid (v_1, v_2)$$

Es decir, la tupla vacía `()` es un valor y cualquier par ordenado de valores es un valor.

La semántica dinámica se especifica mediante las siguientes reglas, nuevamente en sintaxis concreta, obsérvese que estas reglas modelan la estrategia de evaluación ansiosa:

$$\begin{array}{c}
\frac{e_1 \rightarrow e'_1}{(e_1, e_2) \rightarrow (e'_1, e_2)} \qquad \frac{e_2 \rightarrow e'_2}{(v, e_2) \rightarrow (v, e'_2)} \\
\\
\frac{e \rightarrow e'}{\text{fst } e \rightarrow \text{fst } e'} \qquad \frac{e \rightarrow e'}{\text{snd } e \rightarrow \text{snd } e'} \\
\\
\frac{}{\text{fst}(v_1, v_2) \rightarrow v_1} \qquad \frac{}{\text{snd}(v_1, v_2) \rightarrow v_2}
\end{array}$$

El sistema es seguro, es decir se cumple lo siguiente:

Proposición 1 (Seguridad) *Se cumplen las siguientes propiedades:*

- Si $\Gamma \vdash e : \mathsf{T}$ y $e \rightarrow e'$ entonces $\Gamma \vdash e' : \mathsf{T}$.
- Si $\vdash e : \mathsf{T}$ entonces e es un valor o existe e' tal que $e \rightarrow e'$.

1.2. Las operaciones *Curry* y *Uncurry*

La relación entre las funciones de orden superior y las funciones de primer orden con múltiples argumentos es una pregunta de interés práctico en el cálculo lambda. Por lo general estamos acostumbrados a usar y definir funciones de múltiples argumentos, sin embargo en el cálculo lambda con tipos simples todas las funciones son de un argumento. Si tenemos una función de dos argumentos, por ejemplo la suma $\text{suma}(x : \mathsf{Nat}, y : \mathsf{Nat}) = x + y$ generalmente la pensamos como una función cuyo dominio son los pares ordenados de $\mathsf{Nat} \times \mathsf{Nat}$, es decir el tipo en que estamos pensando es $\text{suma} : \mathsf{Nat} \times \mathsf{Nat} \rightarrow \mathsf{Nat}$. Al tener disponibles ahora los tipos producto podemos escribir una versión de esta función como sigue:

$$\text{suma} = \lambda p : \mathsf{Nat} \times \mathsf{Nat} . \text{fst } p + \text{snd } p$$

En el cálculo lambda sin productos esta función no puede escribirse, pero sí la siguiente función relacionada a la anterior.

$$\text{sumaCurry} = \lambda x : \mathsf{Nat} \lambda y : \mathsf{Nat} . x + y$$

la cual devuelve también la suma pero recibe a sus argumentos uno por uno, por ejemplo la función $\text{sumaDos } x = 2 + x$ puede definirse como $\text{sumaDos } x = \text{suma} \langle 2, x \rangle$ o bien como $\text{sumaDos} = \text{sumaCurry } 2$. Obsérvese que la segunda definición no requiere de un argumento x para definirse, lo cual es obligatorio en la primera. Funciones como sumaCurry se llaman funciones *Curryficadas* debido a que fue el lógico Haskell Curry quién las popularizo aunque realmente fue Schönfinkel quien descubrió su definición y relación con las funciones de múltiples argumentos. De hecho, en el lenguaje HASKELL existe la función de orden superior `curry` que transforma una función de dos argumentos en una función que toma a sus argumentos de uno en uno, definida como $\text{curry } f \ x \ y = f \ \langle x, y \rangle$. Análogamente existe la función inversa `uncurry` que transforma una función de orden superior de un argumento en una función que toma a sus dos argumentos a la vez, definida como $\text{uncurry } f \ \langle x, y \rangle = f \ x \ y$.

1.3. El tipo unitario

El tipo unitario () o Void aparentemente no es de gran utilidad puesto que sólo tiene un valor, a saber (). Sin embargo su utilidad surge cuando se necesita devolver un resultado irrelevante como veremos más adelante. Dada su simplicidad no hay ninguna operación interesante sobre este tipo por lo cual no figura en la semántica dinámica.

1.4. Tuplas

Podemos generalizar el constructor producto a un número arbitrario de tipos generando tuplas. La sintaxis es:

Sintaxis concreta	Sintaxis abstracta
$T_1 \times T_2 \times \dots \times T_k$	$\mathbf{tpl}(T_1, \dots, T_k)$
(e_1, e_2, \dots, e_k)	$\mathbf{tpl}(t_1, t_2, \dots, t_k)$
$t.i$	$\mathbf{prj}[i](t)$

Obsérvese que las funciones selectoras se indican por posición usando una notación de punto $t.i$ que indica a la i -ésima posición de la tupla t . Esta notación no debe confundirse con la notación punto para ligado utilizada en la sintaxis abstracta.

No daremos más detalles de esta extensión, en su lugar tratamos una generalización más útil, los tipos registro.

1.5. Tipos Registro

Los tipos producto pueden usarse para definir tipos de tuplas con nombre llamados tipos registro o *records* cuyas componentes se accesan mediante nombres simbólicos en lugar de por posición. De esta forma un *record* modela a una colección de tuplas de diversos tipos con una operación al nivel de valores para extraer los elementos de esta familia por su nombre.

La sintaxis para los tipos registro es:

Concreta	Abstracta
$(\ell_1 : T_1, \ell_2 : T_2, \dots, \ell_k : T_k)$	$\mathbf{rcd}(\ell_1 : T_1, \ell_2 : T_2, \dots, \ell_k : T_k)$
$\{\ell_1 = e_1, \ell_2 = e_2, \dots, \ell_k = e_k\}$	$\mathbf{rcd}(\ell_1 = t_1, \ell_2 = t_2, \dots, \ell_k = t_k)$
$t.\ell$	$\mathbf{prj}[\ell](t)$

Aquí se introduce una nueva categoría de etiquetas ℓ los cuales son identificadores pero no variables. En un registro $\{\ell_1 = e_1, \ell_2 = e_2, \dots, \ell_k = e_k\}$ todas las etiquetas deben ser distintas para evitar confusiones entre los campos del registro.

La semántica estática se rige por las siguientes reglas:

$$\frac{\Gamma \vdash e_1 : T_1 \dots \Gamma \vdash e_k : T_k}{\Gamma \vdash \mathbf{rcd}(\ell_1 = e_1, \dots, \ell_k = e_k) : \mathbf{rcd}(\ell_1 : T_1, \dots, \ell_k : T_k)}$$

$$\frac{\Gamma \vdash e : \mathbf{rcd}(\ell_1 : \mathsf{T}_1, \dots, \ell_j : \mathsf{T}_j, \dots, \ell_k : \mathsf{T}_k))}{\Gamma \vdash \mathbf{prj}[\ell_j](e) : \mathsf{T}_j}$$

A la categoría de valores agregamos un registro de valores:

$$v ::= \dots \mid \mathbf{rcd}(\ell_1 = v_1, \dots, \ell_k = v_k)$$

Finalmente la semántica dinámica para la estrategia de evaluación ansiosa es:

$$\frac{e_j \rightarrow e'_j}{\mathbf{rcd}((\ell_i = v_i)_{i=1}^{j-1}, \ell_j = e_j, (\ell_k = e_k)_{k=j+1}^n) \rightarrow \mathbf{rcd}((\ell_i = v_i)_{i=1}^{j-1}, \ell_j = e'_j, (\ell_k = e_k)_{k=j+1}^n)}$$

La notación $(\ell_i = v_i)_{i=1}^{j-1}$ abrevia a $(\ell_1 = v_1, \dots, \ell_{j-1} = v_{j-1})$.

$$\frac{e \rightarrow e'}{\mathbf{prj}[\ell](e) \rightarrow \mathbf{prj}[\ell](e')}$$

$$\overline{\mathbf{prj}[\ell_j](\mathbf{rcd}(\ell_1 = v_1, \dots, \ell_j = v_j, \dots, \ell_k = v_k)) \rightarrow v_j}$$

2. Sumas

Los tipos suma proporcionan una manera para fusionar en un mismo tipo elementos de tipos distintos. Un elemento del tipo suma $\mathsf{T} + \mathsf{S}$ puede pensarse como un elemento de T o bien como un elemento de S los cuales se encuentran etiquetados para denotar su origen.

Sintaxis concreta	Sintaxis abstracta
$\mathbf{Empty}, \perp, 0$	\mathbf{empty}
$\mathbf{abort}_{\mathsf{T}}(e)$	$\mathbf{abort}[\mathsf{T}](e)$
$\mathsf{T} + \mathsf{S}$	$\mathbf{suma}(\mathsf{T}, \mathsf{S})$
$\mathbf{inl}_{\mathsf{T}} e$	$\mathbf{inl}[\mathsf{T}](e)$
$\mathbf{inr}_{\mathsf{T}} e$	$\mathbf{inr}[\mathsf{T}](e)$
$\mathbf{case } e \text{ of } \{\mathbf{inl } x \Rightarrow e_1 \mid \mathbf{inr } y \Rightarrow e_2\}$	$\mathbf{case}[\mathsf{T}, \mathsf{S}](e, x.e_1, y.e_2)$

Nuevamente distinguimos tres categorías:

- Tipos: **Empty** es el tipo vacío, que representa una elección entre cero alternativas. Por supuesto no hay valores que habiten este tipo. El tipo $\mathsf{T} + \mathsf{S}$ es la suma o unión disjunta de dos tipos.

- Constructores: no existen constructores para **Empty**. Los constructores para una suma binaria se conocen como inyecciones. Se denotan con **inl**, **inr** e indican a cual de las dos posibles elecciones para el tipo suma nos referimos, al tipo de la izquierda o al de la derecha.
- Selectores o destructores: el destructor **abort_T** aborta el cómputo de un programa que intentó evaluar una expresión de tipo **empty** lo cual es imposible. El selector **case** hace un análisis de casos de acuerdo a la etiqueta de un valor de tipo suma.

La semántica estática es:

$$\frac{\Gamma \vdash e : \text{Empty}}{\Gamma \vdash \text{abort}[\mathbf{T}](e) : \mathbf{T}}$$

$$\frac{\Gamma \vdash e : \mathbf{T}}{\Gamma \vdash \text{inl}[\mathbf{S}](e) : \mathbf{T} + \mathbf{S}} \quad \frac{\Gamma \vdash e : \mathbf{S}}{\Gamma \vdash \text{inr}[\mathbf{T}](e) : \mathbf{T} + \mathbf{S}}$$

$$\frac{\Gamma \vdash e : \mathbf{T} + \mathbf{S} \quad \Gamma, x : \mathbf{T} \vdash e_1 : \mathbf{R} \quad \Gamma, y : \mathbf{S} \vdash e_2 : \mathbf{R}}{\Gamma \vdash \text{case}[\mathbf{T}, \mathbf{S}](e, x.e_1, y.e_2) : \mathbf{R}}$$

Los nuevos valores son:

$$v ::= \dots \mid \text{inl}[\mathbf{T}](v) \mid \text{inr}[\mathbf{T}](v)$$

La semántica dinámica para la evaluación ansiosa es:

$$\frac{\frac{e \rightarrow e'}{\text{inl}[\mathbf{T}](e) \rightarrow \text{inl}[\mathbf{T}](e')} \quad \frac{e \rightarrow e'}{\text{inr}[\mathbf{T}](e) \rightarrow \text{inr}[\mathbf{T}](e')}}{e \rightarrow e'} \quad \frac{}{\text{case}[\mathbf{T}, \mathbf{S}](e, x.e_1, y.e_2) \rightarrow \text{case}[\mathbf{T}, \mathbf{S}](e', x.e_1, y.e_2)}$$

$$\frac{}{\text{case}[\mathbf{T}, \mathbf{S}](\text{inl}[\mathbf{S}](v), x.e_1, y.e_2) \rightarrow e_1[x := v]}$$

$$\frac{}{\text{case}[\mathbf{T}, \mathbf{S}](\text{inr}[\mathbf{T}](v), x.e_1, y.e_2) \rightarrow e_2[y := v]}$$

Nuevamente esta extensión con tipos suma es segura. Los detalles se dejan al lector.

2.1. Ejemplos del uso del tipo suma

2.1.1. Booleanos

El tipo booleano puede definirse mediante azucar sintáctica como una suma, en la presencia adicional del tipo unitario como sigue:

$$\text{Bool} =_{\text{def}} \text{Void} + \text{Void}$$

$$\text{true} =_{\text{def}} \text{inl}_{\text{Void}}()$$

$$\text{false} =_{\text{def}} \text{inr}_{\text{Void}}()$$

$$\text{if } e_{\text{Void}} \text{ then } e_2 \text{ else } e_3 =_{\text{def}} \text{case } e_{\text{Void}} \text{ of } \{\text{inl } x \Rightarrow e_2 \mid \text{inr } y \Rightarrow e_3\} \text{ donde } x \notin FV(e_2), y \notin FV(e_3)$$

2.1.2. Tipos Enumerados

Los tipos enumerados son aquellos con un número finito de valores los cuales se pueden enumerar directamente. Estos tipos son definibles mediante sumas del tipo unitario, por ejemplo en Haskell el tipo

```
data Color = Verde | Blanco | Rojo
```

puede definirse como

$$\begin{aligned}\text{Color} &=_{def} \text{Void} + (\text{Void} + \text{Void}) \\ \text{Verde} &=_{def} \text{inl}_{\text{Void} + \text{Void}}() \\ \text{Blanco} &=_{def} \text{inr}_{\text{Void}}(\text{inl}_{\text{Void}}()) \\ \text{Rojo} &=_{def} \text{inr}_{\text{Void}}(\text{inr}_{\text{Void}}())\end{aligned}$$

2.1.3. Tipos Opción

Un tipo de gran interés y utilidad es el tipo opción, llamado `Maybe` en Haskell, definido como

```
data Maybe a = Just a | Nothing
```

el cual puede definirse como

$$\begin{aligned}\text{MaybeT} &=_{def} \text{Void} + \text{T} \\ \text{Just } e &=_{def} \text{inr}_{\text{Void}} e \\ \text{Nothing} &=_{def} \text{inl}_{\text{T}}()\end{aligned}$$

Este tipo es de gran importancia para un manejo de errores en ciertos casos, por ejemplo usando `Maybe Float` podemos devolver el valor `Nothing` al tratar de dividir entre cero.

Existen otros tipos de datos con elecciones como

```
data Nat = Cero | Suc Nat
```

Este no es un tipo suma propiamente dicho puesto que, visto como suma, su definición es $\text{Nat} = \text{Void} + \text{Nat}$ pero esta no es una suma válida pues implica recursión en tipos. De manera que `Nat` es un tipo recursivo, cuya definición formal no trataremos.

2.2. Tipos Variante

Considérese el siguiente tipo en Haskell

```
data Expr = LitN Int | LitF Float | LitA Char | Ident String | Error
```


Este tipo puede representarse como una suma $Int + (Float + (Char + (String + Void)))$ sin embargo el acceso a cada componente particular resulta complicado al tener que componer inyecciones, por ejemplo un elemento de tipo **Char** debe etiquetarse mediante $\text{inr}(\text{inr}(\text{inl}))$. Es más fácil manejar directamente las etiquetas **LitN**, **LitF**, **LitA**, **Ident**, **Error** que utilizar la referencias a su posición en la suma mediante combinaciones de etiquetas inl , inr . Este mecanismo genera los tipos variantes o sumas etiquetadas.

La sintaxis es:

Sintaxis concreta	Sintaxis abstracta
$[\ell_1 : T_1, \ell_2 : T_2, \dots, \ell_k : T_k]$	$\text{var}(\ell_1 : T_1, \ell_2 : T_2, \dots, \ell_k : T_k)$
ℓe	$\text{in}[\ell](e)$
$\text{case } e \text{ of } \{\ell_1 x_1 \Rightarrow e_1 \mid \dots \mid \ell_n x_n \Rightarrow e_n\}$	$\text{case}[\ell_1 : T_1, \ell_2 : T_2, \dots, \ell_k : T_k](e, x_1.e_1, \dots, x_n.e_n)$

La semántica estática se define como sigue:

$$\frac{\Gamma \vdash e : T_i}{\Gamma \vdash \text{in}[\ell_i](e) : \text{var}(\ell_1 : T_1, \dots, \ell_i : T_i, \dots, \ell_k : T_k)}$$

Por otro lado la regla para tipar una expresión **case** es:

$$\frac{\begin{array}{c} \Gamma \vdash e : \text{var}(\ell_1 : T_1, \dots, \ell_k : T_k) \\ \Gamma, x_1 : T_1 \vdash e_1 : T \\ \vdots \\ \Gamma, x_n : T_n \vdash e_n : T \end{array}}{\Gamma \vdash \text{case}[\ell_1 : T_1, \dots, \ell_k : T_k](e, x_1.e_1, \dots, x_n.e_n) : T}$$

Se le deja al lector la labor de definir la semántica dinámica de esta extensión.

Veamos cómo se ve el tipo de HASKELL del ejemplo inicial en nuestra extensión.

```
data Expr = LitN Int | LitF Float | LitA Char | Ident String | Error
```

Las etiquetas son **LitN**, **LitF**, **LitA**, **Ident**, **Error**.

El tipo **Expr** es entonces una variante con cinco alternativas, definido como

```
var(LitN : Int, LitF : Float, LitA : Char, Ident : String, Error : Void)
```

Obsérvese que la etiqueta **Error** que no tiene un tipo asociado en HASKELL debe declararse aquí con el tipo unitario o bien podemos convenir, mediante azucar sintáctica, que una entrada que sólo tenga una etiqueta ℓ en un tipo variante, significa realmente $\ell : \text{void}$.