

Multiset Discrimination

R.Paige, R.Tarjan e.a.

Literatur.:

J.Cai, R.Paige: Using multiset discrimination to solve language processing problems without hashing. TCS 145 (1995), 189-228

F.Henglein: Multiset Discrimination. <http://www.Plan-X.org>

Problem:

Gegeben eine Multimenge (als Liste von Objekten, bei der die Anordnung irrelevant ist). Unterscheide voneinander verschiedene Elemente und identifiziere gleiche Elemente.

Verallgemeinerung 1:

Gegeben eine Multimenge und eine Äquivalenzrelationen auf den Elementen. Unterscheide inäquivalente Elemente und fasse äquivalente in Klassen zusammen.

Verallgemeinerung 2:

Gegeben eine Multimenge von (Objekt,Info)-Paaren und eine Äquivalenzrelation auf den Objekten. Fasse die Informationen zu den Objekten derselben Äquivalenzklassen zusammen.

Ziel:

- ohne Hashing,
- Linearzeit-Algorithmus für azyklische Objekte,
- $n \log n$ -Algorithmus für Graphen (z.B. Automaten)

Anwendungen

- (i) Bei 2 Elementen: Gleichheits- bzw. Äquivalenztest
 - (a) Isomorphietest bei Bäumen
 - (b) Dagifizierung von Bäumen (Trie \rightarrow DAWG)
 - (c) Minimierung von det. Automaten
 - (d) Termgleichheit bei ass./komm./idempot.Funktionen
- (ii) Bei größeren Multimengen:
 - (a) Sind zwei Listen von Ganzzahlen dieselbe Menge?
 - (b) Finde Duplikate in einer Liste von Strings
 - (c) Verbesserung lexikographischen Sortierens

Vorteil der Diskriminierung: in vielen Fällen

- (i) kann Hashing und Sortieren ersetzen
- (ii) ist asymptotisch und praktisch linear bzw. $n \log n$
- (iii) ist relativ einfach zu implementieren (Zeigermaschine)
- (iv) ist auch für rekursive Datentypen möglich

Trennfunktionen

Sei τ ein Typ und \sim eine Äquivalenzrelation auf V_τ . Ein *Trenner* $d : \tau^* \rightarrow (\tau^*)^*$ ordnet jeder Liste $\vec{v} = [v_1, \dots, v_n] : \tau^*$ eine Liste $d(\vec{v}) = [\vec{u}_1, \dots, \vec{u}_k] : (\tau^*)^*$ von τ -Listen zu, so daß

- (i) die Verkettung von $d(\vec{v})$ eine Permutation von \vec{v} ist und
- (ii) jedes \vec{u}_i , $1 \leq i \leq k$, eine Permutation der Teilfolge aller $v \sim v_i$ von \vec{v} ist, für ein v_i in \vec{v} .

Der Trenner d ist *ordnungserhaltend*, wenn die Elemente in jedem u_i in derselben Reihenfolge wie in \vec{v} vorkommen.

Verallgemeinerte Form mit $d_\tau : (\tau \times \iota)^* \rightarrow (\iota^*)^*$:

Eingabeliste $[(v_1, w_1), \dots, (v_n, w_n)]$, wobei die Ausgabe bei \vec{u}_i statt $\vec{u}_i = [v_{i,1}, \dots, v_{i,k}]$ die zugehörige Liste $[w_{i,1}, \dots, w_{i,k}] : \iota^*$ der $w_{i,j}$ mit $(v_{i,j}, w_{i,j})$ in der Eingabe hat.

Die $w_i : \iota$ sind Informationen eines festen Typs ι und gehen nicht in die Äquivalenzrelation \sim ein. (Einfache Form: $w_i = v_i$)

Unterscheidungsproblem

Das *Unterscheidungsproblem* für τ modulo \sim besteht darin, eine Trennfunktion d_τ zu finden, die $d_\tau(\vec{v})$ in der Zeit $O(|\vec{v}|)$ berechnet. Das wird induktiv über den Aufbau von τ und unabhängig von den w_i erfolgen.

Typen τ und Werte $v : \tau$

τ	$:=$	v	$:=$
Grundtyp	a	c	
Einheitstyp	1	$()$	
Paare	$\tau_1 \times \tau_2$	(v_1, v_2)	
Disj. Vereinigung	$(\tau_1 + \tau_2)$	$in_l(v) \mid in_r(v)$	
Adresse	$ref(\tau)$	n	
Typvariable	α		
Induktiver Typ	$\mu\alpha.\tau$		
	a $:=$	$\text{char} \mid \text{bool}.$	
	τ^* $:=$	$\mu\alpha(1 + \tau \times \alpha)$	
	string $:=$	char [*] .	

$V_a :=$ der Bereich der Werte des Grundtyps a .

$V[N] :=$ Bereich aller aus $N \subseteq \mathbb{N}$ und den V_a gebildeten Werte.

Typisierung von Werten

Werte enthalten i.a. Adressen n . Sie werden in einem Kontext Γ von Typannahmen $n : \tau$ nach folgenden Regeln getypt:

$$\begin{array}{c}
 \overline{\Gamma \vdash 1 : ()} \quad \overline{\Gamma \vdash c : a} \\
 \\
 \frac{\Gamma \vdash v_1 : \tau_1 \quad \Gamma \vdash v_2 : \tau_2}{\Gamma \vdash (v_1, v_2) : \tau_1 \times \tau_2} \\
 \\
 \frac{\Gamma \vdash v : \tau_1}{\Gamma \vdash in_l(v) : (\tau_1 + \tau_2)} \quad \frac{\Gamma \vdash v : \tau_2}{\Gamma \vdash in_r(v) : (\tau_1 + \tau_2)} \\
 \\
 \frac{n : ref(\tau) \in \Gamma}{\Gamma \vdash n : ref(\tau)} \\
 \\
 \frac{\Gamma \vdash v : \tau[\mu\alpha.\tau/\alpha]}{\Gamma \vdash v : \mu\alpha.\tau}
 \end{array}$$

V_τ sei der Bereich der Werte vom Typ τ .

Größe von Objekten

Auf den Grundtypen a sei eine *Größe* $|\cdot|_a : V_a \rightarrow \mathbb{N}$ und eine Äquivalenzrelation $=_a$ gegeben.

Sie wird wie folgt auf beliebige Werte fortgesetzt:

$$\begin{aligned}
 |()| &:= 0 \\
 |(v_1, v_2)| &:= |v_1| + |v_2| \\
 |in_l(v)| &:= 1 + |v| \\
 |in_r(v)| &:= 1 + |v| \\
 |v| &:= |c|_a \quad \text{für } v \in V_a \\
 |n| &:= 1
 \end{aligned}$$

Beachte:

- Alle Adressen haben dieselbe Größe 1.
- Bei Paaren addiert sich die Größe der Komponenten, d.h. es wird angenommen, daß sich die Werte gleiche Teile *nicht* teilen. (Structure sharing nur via Adressen!)
- Damit Informationen $w : \iota$ die Größe 1 haben, nehmen wir für ι nicht `string`, sondern eher `ref(string)`; wir brauchen aber, daß ι unter $\tau \times \iota$ abgeschlossen ist und nehmen daher $\iota := \bigcup_{\tau} \text{ref}(\tau)$,

Graphendarstellung von Objekten

Ein Wert v wird als Graph dargestellt, mit v als Marke an der Wurzel des Graphen und den Teilen an den Nachfolgerknoten; Adressen treten nur an den Blättern auf. Die Blätter *sind* keine Adressen $n \in N$, sondern nur mit Adressen markiert.

Ein „Sharing“ erfolgt noch nicht im Graphen für v , sondern erst, wenn man !-Kanten von den Blättern mit Marke n zum Graphen für $!n$ zieht.

Typisierter Speicher

Trenner $d_\tau : (\tau \times \iota)^* \rightarrow \iota^{**}$ verändern den Speicher S , d.h. sie benutzen S als implizites Argument.

Ein *Speicher* $S = (N, !)$ ist eine endliche Menge N von Adressen mit einer Auswertung $! : N \rightarrow V[N]$.

Der in $n \in S$ gespeicherte Wert ist $!n \in V[N]$. Die *Größe* des Speichers ist

$$|(N, !)| := \sum_{n \in N} |!n|.$$

Eine *Typisierung des Speichers* ist eine bezüglich \subseteq maximale Menge Γ von Typannahmen $n : \text{ref}(\tau)$, so daß $\Gamma \vdash !n : \tau$ für alle $n : \text{ref}(\tau) \in \Gamma$.

Jede Äquivalenzrelation \sim_N auf S induziert

(i) die Kongruenzrelation $\simeq_{V[N]}$ auf $V[N]$,

$$\simeq_{V[N]} := \bigcap \{ \simeq \mid \simeq \text{ ist Kongruenzrelation, } \sim_N \subseteq \simeq, =_a \subseteq \simeq \}$$

(ii) die Äquivalenzrelation $\sim_{(N,!)}$ auf N ,

$$\sim_{(N,!)} := \{ (n, n') \in N \times N \mid !n \simeq_{V[N]} !(n') \}.$$

Die durch $S = (N, !)$ definierte *Isomorphie* \equiv_S ist die größte Äquivalenzrelation \sim auf N , die sich selbst induziert, d.h. für die $\sim = \sim_S$ gilt.

Rekursive Definition der Trenner

Man muß i.a. den Trenner d_τ mit Hilfe einer (zu einer Typumgebung $[\alpha_i/\tau_i, \dots]$ passenden) *Trennerumgebung* $\Gamma = [\alpha_i/d_{\tau_i}, \dots]$ definieren, die für jede freie Typvariable α_i von τ einen Trenner

$$d_{\tau_i} : \forall \beta ((\tau_i \times \beta)^* \rightarrow \beta^{**}) \quad (8)$$

vorgibt. Die Trenner werden rekursiv über den Typaufbau durch *Trennerkombinatoren* D_c für die Typkonstruktoren c aus Trennern für die Teiltypen zusammengesetzt:

$$\begin{aligned} d_\alpha^\Gamma &:= \Gamma(\alpha), \\ d_{\text{char}}^\Gamma &:= D_{\text{char}}, \\ d_1^\Gamma &:= D_{\text{unit}}, \\ d_{\tau \text{ ref}}^\Gamma &:= D_{\text{ref}}(d_\tau^\Gamma), \\ d_{\sigma \times \tau}^\Gamma &:= D_\times(d_\sigma^\Gamma, d_\tau^\Gamma) =: d_\sigma^\Gamma \times d_\tau^\Gamma, \\ d_{\sigma + \tau}^\Gamma &:= D_+(d_\sigma^\Gamma, d_\tau^\Gamma) =: d_\sigma^\Gamma + d_\tau^\Gamma, \\ d_{\mu\alpha.\tau}^\Gamma &:= D_\mu(\lambda \tilde{d}. d_\tau^{\Gamma[\alpha/\tilde{d}]}) =: \mu \tilde{d}. d_\tau^{\Gamma[\alpha/\tilde{d}]} \end{aligned}$$

Die einzelnen D_c folgen unten, von D_{ref} und D_μ nur Spezialfälle. Die Trenner für rekursive Typen sind auch rekursiv zu definieren, z.B. der Trenner für τ -Listen:

$$d_{\tau^*}^\Gamma := d_{\mu\alpha.(1+\tau \times \alpha)}^\Gamma = \mu \tilde{d}. (d_1^\Gamma + d_\tau^\Gamma \times \tilde{d}) =: D_*(d_\tau^\Gamma).$$

Das kann man nur in Programmiersprachen ausführen, die Typen wie in (8) erlauben (da zwei Rekursionsaufrufe eines Trenners zwei Ausgabetypen erfordern können), und wo man an die jeweilige Typumgebung herankommt.

Konstruktion von Trennern $d_\tau : (\tau \times \iota)^* \rightarrow \iota^{**}$
 $\tau := \text{bool} \mid \text{char} \mid (\tau \times \tau) \mid \tau^* \mid \text{ref}(\tau).$

Informationen $w : \iota$ seien oBdA stets in Adressen gehalten, also $V_\iota \subseteq N$ und $\iota := \bigcup_\tau \text{ref}(\tau).$

• *Fall* $\tau = \text{bool}$:

- (i) Initialisiere¹ drei Adressen $n_{eq} : \text{ref}(\text{bool}^*), n_t : \text{ref}(\iota^*)$ und $n_f : \text{ref}(\iota^*)$ mit leeren Listen

$$n_{eq} := [], \quad n_t := [], \quad n_f := [].$$

- (ii) Durchlaufe die Eingabe $[(v_1, w_1), \dots, (v_n, w_n)]$ und aktualisiere die Werte gemäß

$$\begin{aligned} n_t &:= \begin{cases} w_i :: !n_t, & \text{falls } v_i = t, \\ !n_t, & \text{sonst} \end{cases} \\ n_f &:= \begin{cases} w_i :: !n_f, & \text{falls } v_i = f, \\ !n_f, & \text{sonst} \end{cases} \end{aligned}$$

und falls dabei für $b : \text{bool}$ die Liste $!n_b$ nicht leer wird, setze $n_{eq} := b :: !n_{eq}.$

- (iii) Gib als Wert $d_{\text{bool}}([(v_1, w_1), \dots, (v_n, w_n)])$ die Liste der $!n_b$ mit b auf der Liste $!n_{eq}$ aus, setze für diese b dann $n_b := []$ und schließlich $n_{eq} := [].$

Das kann in $O(n)$ Schritten gemacht werden, wenn das Ablesen/Speichern eines Werts bei einer Adresse nur $O(1)$ Schritte kosten.

Anschließend sind unter n_{eq}, n_t, n_f wieder $[]$ gespeichert.

¹Notation: $n := v$ meint, unter Adresse n den Wert v zu speichern.

- *Fall char*: Das geht analog, mit Adressen $n_{\text{eq}} : \text{ref}(\text{char}^*)$ und $n_b : \text{ref}(\iota^*)$ für jedes $b \in V_{\text{char}}$.

(Das wird unten in der SML-Implementierung mit einem Array der Länge $|V_{\text{char}}| = 256$ gemacht.)

Implementierung

Wir bilden einen universellen Typ

$$U \simeq V_{\text{char}} + V_{\text{string}} + U^* + (U \times U) + \text{ref}(U)$$

und setzen $\iota = \text{ref}(U)$. Damit haben alle $w : \iota$ die Größe 1, aber man braucht lästige Einbettungen von U in ι und umgekehrt.

$\langle \text{discriminate.sml} \rangle \equiv$

```
structure Discriminate =
```

```
  struct
```

```
    datatype U = C of char
```

```
              | S of string
```

```
              | L of U list
```

```
              | P of U * U
```

```
              | R of (U ref) (* info as value *)
```

```
              | Bag of U list
```

```
              | Set of U list
```

```
              | N of (U ref) * (U ref list ref)
```

```
  type info = U ref
```

```
  type class = info list (* infos as eq-class of nodes
```

```
  type discriminator = (U * info) list -> class list
```

```
  exception Undefined
```

Für die Trennung nach Adressen modulo $\sim_{\text{ref}(U)}$ wird ein Summand mit $N : \text{ref}(U) \times \text{ref}(\text{ref}(U)^*) \rightarrow U$ in U eingebettet; eine Trennung nach Adressengleichheit $=_{\text{ref}(U)}$ (ohne Angabe der Äquivalenzrelation) ist nicht implementiert, wäre aber mit einer Typänderung von N machbar (s.u.).

Beispiel: Trenner für char

$\langle CharDiscriminator \rangle \equiv$

```
structure CharDiscr =
  struct
    val size = 256
    exception Size
```

```
val dict = Array.array(size, []:class) (* value restr
val labels = ref [] : int list ref
```

```
local
```

```
  fun add (v:char,w) =
```

```
    let
```

```
      val i = Char.ord v
```

```
      val ts = Array.sub(dict,i)
```

```
    in
```

```
      if List.exists (fn j => j=i) (!labels)
```

```
      then () else labels := i::(!labels);
```

```
      Array.update(dict,i,ts @ [w])
```

```
    end
```

```
  fun getLists () =
```

```
    let
```

```
      val classes = map (fn i => Array.sub(dict,
                                          (rev (!labels)))
```

```
    in
```

```
      List.app (fn i => Array.update(dict,i,[]
                                     (!labels));
```

```
      labels := [];
```

```
      classes
```

```
    end
```

```
in
fun discriminate vs =
  (List.app add vs ; getLists())
  handle Subscript => raise Size
end
end
```

Beispiel

$\langle \text{discriminate.expl.sml} \rangle \equiv$

```
use "discriminate.sml";
open Discriminate;
```

```
fun discriminateToStrings l =
  map (map toString) (discriminate l);
```

(* Example 1: Discriminate according to Char-labels *)

```
discriminateToStrings [(C #"a",toInfo "Affe"),
                       (C #"b",toInfo "Baer"),
                       (C #"a",toInfo "Albatross"),
                       (C #"b",int_toInfo 3)];
```

$\langle \text{Ergebnis} \rangle \equiv$

```
val it = [["Affe","Albatross"],["Baer","3"]] : string list
```

Trenner für Paare

$$d_{\sigma \times \tau} : ((\sigma \times \tau) \times \iota)^* \rightarrow \iota^{**}$$

- Fall $\sigma \times \tau$:

In der Eingabe $[(v_1, w_1), \dots, (v_n, w_n)]$ sei $v_i = (v_{i,1}, v_{i,2})$.

- Für $i = 1, \dots, n$ speichere die zweite Komponente $v_{i,2}$ des Werts unter einer frischen Adresse $\overline{v_{i,2}}$, bilde das Paar $(\overline{v_{i,2}}, w_i) : \iota \times \iota$ und speichere diesen Wert in einer ebenfalls neuen Adresse $\overline{(\overline{v_{i,2}}, w_i)} : \iota$.
- Mit dem induktiv schon definierten Trenner d_σ trenne

$$\begin{aligned} d_\sigma([(v_{1,1}, \overline{(\overline{v_{1,2}}, w_1)})], \dots, (v_{n,1}, \overline{(\overline{v_{n,2}}, w_n)})) \\ = [[\overline{(\overline{v_{1,2}}, w_1)}, \dots], \dots, [\overline{(\overline{v_{k,2}}, w_k)}, \dots]] \end{aligned}$$

nach gleichen ersten Komponenten $v_{1,i}$ der v_i .

- Ersetze die Adressen durch die dort gespeicherten Werte, und trenne jede Liste mit dem Trenner d_τ :

$$[d_\tau([(v_{1,2}, w_1), \dots]), \dots, d_\tau([(v_{k,2}, w_k), \dots])]$$

- Verkette die Ergebnislisten: das ergibt eine Trennung der w_j nach gleichen Wertpaaren $(v_{i,1}, v_{i,2})$,

$$[[w_{1,1} \dots, w_{1,k_1}], \dots, [w_{r,1}, \dots, w_{r,k_r}]].$$

Die Speicherung in neuen Adressen ist nur wegen der Wahl des Typs ι notwendig.

Trennung nach \sim_N -äquivalenten Adressen

$$d_{ref(\tau)} : (ref(\tau) \times \iota)^* \rightarrow \iota^{**}$$

Die Eingabeliste $[(n_1, w_1), \dots, (n_m, w_m)]$ soll bzgl. einer Äquivalenz \sim_N auf dem Speicher N getrennt werden:

$$[[w_{1,1}, \dots, w_{1,k_1}], \dots, [w_{r,1}, \dots, w_{r,k_r}]]$$

wobei die $w_{i,j}$ die Informationen aller $n \sim_N n_i$ seien.

Stelle den Speicher $S = (N, !)$ und die Äquivalenz \sim_N durch *discrimination records* $r = \{\mathbf{cont} : \mathbf{U} \mathbf{ref}, \mathbf{eq} : \mathbf{N} \mathbf{ref}\}$ dar:

Jeder Adresse n entspricht ein $r = (!n, n')$, wobei $n' \in N$ die Klasse $[n]_{\sim_N}$ darstellt, d.h. es ist

$$n_1 \sim_N n_2 \iff n'_1 = n'_2$$

(aber $!n'$ muß nicht eine Liste der Elemente von $[n]_{\sim_N}$ sein).

Die Äquivalenz zweier Adressen n_1, n_2 kann dann in $O(1)$ Schritten festgestellt werden.

Änderungen von E und \sim_N werden durch Änderungen bestehender $r = (!n, n')$ erreicht, und Änderungen an N durch Allokation neuer r oder Löschung vorhandener.

Die Äquivalenz \sim_N sei also unter n mit übergeben, d.h. $!n = (!n, n')$, oBdA mit $n' : ref(\iota^*)$.

Reduziere Trennung modulo \sim_N auf Trennung modulo $=_N$:

$$d_{ref(\tau)}([(n_1, w_1), \dots]) := d_{ref(\iota^*)}([(n'_1, w_1), \dots])$$

Trennung nach gleichen Adressen

$$d_{\text{ref}(\iota^*)} : (\text{ref}(\iota^*) \times \iota)^* \rightarrow \iota^{**}$$

- *Fall* $\text{ref}(\iota^*)$: Zur Eingabeliste $[(n_1, w_1), \dots, (n_m, w_m)]$ sollen die w_i nach $=_N$ -gleichen Adressen n_i getrennt werden, unabhängig von den gespeicherten Werten.

(i) Wähle frische Adresse $n_{\text{eq}} : \text{ref}(\iota^*)$, setze $n_{\text{eq}} := []$ und $n_i := []$ für $i = 1, \dots, m$.

(ii) Für $i = 1, \dots, m$ setze

$$n_{\text{eq}} := \begin{cases} n_i :: !n_{\text{eq}}, & \text{falls } !n_i = [], \\ !n_{\text{eq}}, & \text{sonst} \end{cases}$$

und anschließend (für ordnungserhaltendes $d_{\text{ref} \iota^*}$)

$$n_i := !n_i @ [w_i]$$

Danach steht unter n_{eq} eine Anordnung der Menge $\{n_1, \dots, n_m\}$, mit jedem n_i genau einmal, und bei n_i die Liste $[w_{i,1}, \dots, w_{i,k_i}]$ aller $w_{i,j}$, für die $(n_i, w_{i,j})$ in der Eingabeliste vorkommt.

(iii) Gib für jedes n auf der Liste $!n_{\text{eq}}$ die Liste $!n$ aus.

Achtung:

Die anfangs unter n_i gespeicherten Werte sind verloren.

Das ist harmlos, falls die n_i nur als Darstellung von \sim_N -Klassen benutzt werden: dazu braucht man nur Gleichheiten $n_i =_N n_j$, und die sind von den $E(n_i)$ unabhängig.

Voraussetzung: unter Adressen ι kann man Informationen vom Typ ι^* speichern, also $\text{ref}(\iota^*) \subseteq \iota$.

Trenner nach gleichen Adressen

$$d_{ref(\tau)} : (ref(\tau) \times \iota)^* \rightarrow \iota^{**}$$

- *Fall* $ref(\tau)$: Die Eingabe $[(n_1, w_1), \dots, (n_m, w_m)]$ soll nach $=_N$ -gleichen Adressen getrennt werden, wobei $n_i : ref(\tau)$.²

Mit einem universellen Typ U kann man das auf den Fall $ref(\iota^* \times \tau)$ zurückführen:

- Ändere die gespeicherten Werte zu $n_i := ([], !n_i)$ für $i = 1, \dots, m$.
- Trenne $[(n_1, w_1), \dots, (n_m, w_m)]$ wie im Fall $ref(\iota^*)$, ohne die zweiten Komponenten zu verändern.
- Setze die ursprünglichen Werte mit $n_i := second(!n_i)$ wieder her, für $i = 1, \dots, m$.
- Gib das Ergebnis von (ii) aus.

Da die Gleichheit der Adressen beibehalten bleibt, erhält man die korrekte Antwort. Damit im ersten Schritt kein Typfehler auftritt, muß man $ref(\tau)$ und $ref(\iota^* \times \tau)$ als Spezialfälle von $\iota = ref(U)$ behandeln können.

²Es ist zu teuer, zuerst die Gleichheiten unter den n_i zu bestimmen; das geht auch nicht mit

$$\begin{aligned} & d_{set(ref(\tau))}([([n_1], (n_1, w_1)), \dots, ([n_m], (n_m, w_m))]) \\ &= [[(n_{j_1}, w_{j_1}), \dots, (n_{j_{k_1}}, w_{j_{k_1}})], \dots, [(n_{j_r}, w_{j_r}), \dots, (n_{j_{k_r}}, w_{j_{k_r}})]], \end{aligned}$$

denn dabei würde die schwache Sortierung der Verkettung von Einermengen $[(n_i, \tilde{w}_i)]$ benötigt, die rekursiv die Trennung nach den n_i erfordert.

Trennung nach gleichen Adressen

$\langle \textit{Discriminator for references} \rangle \equiv$

```

structure RefDiscr =
  struct
    type discrRef = {content: U ref, eq: class ref}

    fun eqClass (v:discrRef) = #eq (!v)
    fun isEmpty (v:discrRef) = (!(eqClass v) = [])
    fun add (v:discrRef) (w:info) =
      let val c = (eqClass v) in c := (!c)@[w] end
    fun getThenResetClass (v:discrRef) =
      let val els = !(eqClass v)
      in (eqClass v) := [] ; els end

    val eqClasses = ref [] : discrRef list ref

    fun addClass (v:discrRef,w:info) =
      (if isEmpty v then
        eqClasses := v::(!eqClasses)
      else ();
      (add v w); ())

    fun discriminate (vws : (discrRef * info) list)
      let val _ = List.app addClass vws
        val classes =
          map (fn v => getThenResetClass v)
            (!eqClasses)
        in
          eqClasses := []; classes
        end
      end
  end
end

```

Trenner für Listen

$$d_{\tau^*} : (\tau^* \times \iota)^* \rightarrow \iota^{**}$$

- Fall $\tau^* = \mu\alpha(1 + (\tau \times \alpha))$:
 - (i) Zerlege $[(v_1, w_1), \dots, (v_n, w_n)]$ in die Paare mit $v_i : 1$ und die mit $v_i : \tau \times \tau^*$, d.h. die mit leerer Liste bzw. nicht-leerer Liste v_i .
 - (ii) Alle w_i mit $v_i : 1$ bilden eine Klasse A .
 - (iii) Trenne die Teilliste der übrigen (v_i, w_i) mit $d_{\tau \times \tau^*}$ in

$$\begin{aligned} & d_{\tau \times \tau^*}([(v_1, w_1), \dots, (v_r, w_r)]) \\ &= [[w_{1,1}, \dots, w_{1,k_1}], \dots, [w_{r,1}, \dots, w_{r,k_r}]] \end{aligned}$$

und füge A hinzu, falls A nicht leer war.

Das ist keine Induktion über die Typausdrücke, terminiert aber, da bei der rekursiven Verwendung von d_{τ^*} nur kürzere Listen in den Eingabetupeln auftreten.

Analog werden andere Datentypen $\mu\alpha\tau$ behandelt.

Trennung modulo $f : \sigma \rightarrow \tau$

Trennung einer Eingabe $[(v_1, w_1), \dots, (v_n, w_n)]$ nicht nach gleichen Werten $v_i =_\sigma v_j$, sondern nach f -gleichen Werten,

$$v_i =_\sigma^f v_j : \iff f(v_i) =_\tau f(v_j).$$

$$\begin{aligned} d_\sigma^f([(v_1, w_1), \dots, (v_n, w_n)]) \\ := d_\tau([(f(v_1), w_1), \dots, (f(v_n), w_n)]). \end{aligned}$$

Trennung nach *bag*-Gleichheit

$$d_{\mathbf{bag}(\tau)} : (\tau^* \times \iota)^* \rightarrow \iota^{**}$$

Unter $[(\vec{v}_1, w_1), \dots, (\vec{v}_n, w_n)]$ soll man diejenigen w_i in eine Liste zusammenfassen, die in einem (\vec{v}_j, w_j) mit $\vec{v}_j \simeq_{\mathbf{bag}(\tau)} \vec{v}_i$ als w_j vorkommen.

1.Idee: Zuerst aus \vec{v} eine neue Liste $\mathbf{bag}(\vec{v}) : \tau^*$ konstruieren, in der alle gleichen Elemente aus \vec{v} nebeneinander stehen, und dann $[(\mathbf{bag}(\vec{v}_1), w_1), \dots, (\mathbf{bag}(\vec{v}_n), w_n)]$ mit d_{τ^*} trennen.

Aber die Konstruktion von $\mathbf{bag}(\vec{v})$ erfordert schon zu viele Vergleiche, und man muß die entstandenen Listen noch nach einer Ordnung auf τ sortieren, damit die Trennung die für $=_{\mathbf{bag}(\tau)}$ richtigen Ergebnisse liefert.

2.Effiziente Lösung (B.Paige): durch *schwache Sortierung* (erfordert keine Ordnung auf τ !)

- (i) Definieren Sie eine Funktion $weaksort : \tau^{**} \rightarrow \tau^{**}$, die jede Liste \vec{v}_i der Eingabe $[\vec{v}_1, \dots, \vec{v}_n]$ in die Liste \vec{v}'_i umordnet, bei der die Elemente wie ihre ersten Vorkommen in der Verkettung von $[\vec{v}_1, \dots, \vec{v}_n]$ vorkommen.
- (ii) Berechnen Sie aus $[(\vec{v}_1, w_1), \dots, (\vec{v}_n, w_n)]$ zuerst $[\vec{v}'_1, \dots, \vec{v}'_n] = weaksort([\vec{v}_1, \dots, \vec{v}_n])$ und wenden den Trenner d_{τ^*} auf $[(\vec{v}'_1, w_1), \dots, (\vec{v}'_n, w_n)]$ an. Das gibt das korrekte Ergebnis, da $\vec{v}'_i =_{\tau^*} \vec{v}'_j$ genau dann gilt, wenn $\vec{v}_i =_{\mathbf{bag}(\tau)} \vec{v}_j$, weil die Komponenten in \vec{v}'_i, \vec{v}'_j in schwacher Sortierung auftreten.

Schwache Sortierung von $[\vec{v}_1, \dots, \vec{v}_n]$
weaksort : $\tau^{**} \rightarrow \tau^{**}$

- (i) Speichere für spätere Trennung jedes $\vec{v}_i = [c_{i,1}, \dots, c_{i,k_i}]$ in einer frischen Adresse \overline{v}_i und dann für jedes j das Paar $(\overline{v}_i, c_{i,j})$ in einer frischen Adresse $w_{i,j} := \overline{(v_i, c_{i,j})}$. Bilde für jedes i die Liste

$$\tilde{v}_i := [(c_{i,1}, (\overline{v}_i, w_{i,1})), \dots, (c_{i,k_i}, (\overline{v}_i, w_{i,k_i}))].$$

- (ii) Verkette die Ergebnisse $[\tilde{v}_1, \dots, \tilde{v}_n]$ zu

$$[(c_{1,1}, (\overline{v}_1, w_{1,1})), \dots, (c_{n,k_n}, (\overline{v}_n, w_{n,k_n}))]$$

und trenne diese Liste mit d_τ nach den Elementen $c_{i,j}$ (entsprechend einer bestimmten Anordnung $<_{[\vec{v}_1, \dots, \vec{v}_n]}$ der Elemente $c_{i,j}$). Die Ergebnisliste

$$[[(\overline{v}_i, w_{i,j_1}), \dots,], \dots, [\dots]]$$

faßt alle $(\overline{v}_i, w_{i,j})$ in ein Element zusammen, deren Liste v_i ein bestimmtes Element c aus der Verkettung der Argumentlisten $[\vec{v}_1, \dots, \vec{v}_n]$ enthält.

- (iii) Trenne die Verkettung dieser Ergebnisse (ordnungserhaltend) mit $d_{ref(\tau^*)}$ bezüglich der Gleichheit $=_N$ als Äquivalenz \sim_N : das Ergebnis

$$d_{ref(\tau^*)}([(v_i, w_{i,j_1}), \dots,]) = [[w_{i,j_1}, \dots, w_{i,j_{k_i}}], \dots, [\dots]]$$

enthält in den zweiten Komponenten der $[w_{i,j_1}, \dots, w_{i,j_{k_i}}]$ genau die Elemente $c_{i,j}$ aus \vec{v}_i in der Ordnung $<_{[\vec{v}_1, \dots, \vec{v}_n]}$ aufgeführt, in den ersten Komponenten jeweils \overline{v}_i .

Daraus kann eine Zuordnung $\bar{v}_i \mapsto \bar{v}'_i := \mathbf{bag}(\bar{v}_i)$ bestimmt werden, mit der die Eingabe $[\bar{v}_1, \dots, \bar{v}_n]$ in $[\bar{v}'_1, \dots, \bar{v}'_n]$ umgewandelt werden kann (in $O(n)$, wenn man sich in a) die Liste $[\bar{v}_1, \dots, \bar{v}_n]$ merkt.)

Trennung nach *set*-Gleichheit

$$d_{\mathbf{set}(\tau)}(\tau^* \times \iota)^* \rightarrow \iota^{**}$$

- (i) Definiere *weaksort* : $\tau^{**} \rightarrow \tau^{**}$ wie bei der Trennung nach *bag*-Gleichheit, nur daß im Schritt (c) aus den Listen $[w_{i,j_1}, \dots, w_{i,j_{k_i}}]$ keine Multimenge, sondern eine (schwach sortierte) Menge gebildet wird, d.h. Duplikate ignoriert werden.
- (ii) Berechne aus $[(\bar{v}_1, w_1), \dots, (\bar{v}_n, w_n)]$ zuerst $[\bar{v}'_1, \dots, \bar{v}'_n] = \mathit{weaksort}([\bar{v}_1, \dots, \bar{v}_n])$ und wende den Trenner d_{τ^*} auf $[(\bar{v}'_1, w_1), \dots, (\bar{v}'_n, w_n)]$ an. Das gibt das korrekte Ergebnis, da $\bar{v}'_i =_{\tau^*} \bar{v}'_j$ genau dann gilt, wenn $\bar{v}_i =_{\mathbf{set}(\tau)} \bar{v}_j$, weil die Komponenten in \bar{v}'_i, \bar{v}'_j in schwacher Sortierung und jeweils nur einmal auftreten.

Beispiel: Schwache Sortierung für Multimengen

Sortiere $[aba, bca, aab, ba, bcb]$ mit $d_{\text{bag}(\text{char})}$.

(i) Paare die Elemente mit Adressen ihrer Listen und sich:

$$aba \mapsto [(a, (aba, a)), (b, (aba, b)), (a, (aba, a))]$$

$$\vdots$$

$$bcb \mapsto [(b, (bcb, b)), (c, (bcb, b)), (b, (bcb, b))]$$

(ii) Trenne die Verkettung der erhaltenen Listen mit d_{char} :

$$\begin{aligned} & [[(aba, a), (aba, a), (bca, a), (aab, a), (aab, a), (ba, a)], \\ & \quad [(aba, b), (bca, b), (aab, b), (ba, b), (bcb, b), (bcb, b)], \\ & \quad [(bca, c), (bcb, c)]] \end{aligned}$$

(iii) Trenne deren Verkettung mit $d_{\text{ref}(\text{char}^*)}$ nach Gleichheit der String-Adressen:

$$\begin{aligned} & [[a, a, b], [a, b, c], [a, a, b], [a, b], [b, b, c]] \\ = & [aab, abc, aab, ab, bbc] \end{aligned}$$

Im Programm wird eine andere Anordnung erzeugt:

$\langle \text{discriminate.expl.sml} \rangle + \equiv$

```
bagsort [[C #"a",C #"b",C #"a"], [C #"b",C #"c",C #"a"],
         [C #"a",C #"a",C #"b"], [C #"b",C #"a"],
         [C #"b",C #"c",C #"b"]];
```

$\langle \text{Ergebnis} \rangle + \equiv$

```
val it =
  [[C #"b",C #"a",C #"a"], [C #"c",C #"b",C #"a"],
   [C #"b",C #"a",C #"a"],
   [C #"b",C #"a"], [C #"c",C #"b",C #"b"] ] : U list list
```

Anwendung: Dokumentklassifizierung

Klassifiziere Dokumente nach darin vorkommenden Begriffen.

Vereinfacht: Dokument := Zeichenreihe, Begriff := Buchstabe

$\langle \text{discriminate.expl.sml} \rangle + \equiv$

```
fun discrChar (s : string list) =
  let
    fun string2Uref str = R (ref (S str))
    fun cPair (u as (R(Ref as (ref(S str)))))) =
      map (fn c => (C c, Ref)) (String.explode str)
    val Lists = map cPair (map string2Uref s)
  in
    discriminate (List.concat Lists)
  end;
```

```
map (map toString) (discrChar ["abc","aa","b","ac"]);
```

$\langle \text{Ergebnis} \rangle + \equiv$

```
val it = [["abc","aa","aa","ac"],["abc","b"],["abc","ac"]
```

Ignorieren wir die Einbettungen in den Typ U, so bildet man zu jeder Zeichenreihe die Paare aus den Zeichen und der Reihe,

abc \mapsto [(a, abc), (b, abc), (c, abc)]

aa \mapsto [(a, aa), (a, aa)]

b \mapsto [(b, b)]

ac \mapsto [(a, ac), (c, ac)]

und trennt die Verkettung der Listen nach den Zeichen.

Anwendung: Index aller Wörter eines Texts

Der Index des Texts gibt zu den Wörtern ihre Vorkommen an.

Man paart jedes Wort mit seiner Position und trennt die Paare dann nach den Wörtern:

```

⟨discriminate.expl.sml⟩+≡
  fun discrString text =
    let
      fun addPosition (s::ss) n =
          (S s, int_toInfo n)::(addPosition ss (n+1))
        | addPosition [] n = []
    in
      (discriminateToStrings (addPosition text 1))
    end;

  discrString ["aber", "er", "aber", "so", "und", "so"];

```

```

⟨Ergebnis⟩+≡
  val it = [["1","3"],["2"],["4","6"],["5"]]

```

Um die Wörter in der Ausgabe zu haben, paart man jedes w im Text mit seiner Position i und w selbst, und trennt die entstandene Liste der $(w, (w, i))$ nach den Wörtern.

Lokale Hilfsfunktionen

$\langle \text{Hilfsfunktionen} \rangle \equiv$

```
local
```

```
  fun unC (C c,w) = (c,w)
    | unC _ = raise Inhomogeneity
```

```
  fun reassoc [] = [] : (U * info) list
    | reassoc ((P(u,v),w) :: uvws) =
        (u,ref(P(v,R w))) :: (reassoc uvws)
    | reassoc (_ :: uvws) = raise Inhomogeneity
```

```
  fun unP (ref(P(v,R w))) = (v,w)
    | unP _ = raise Inhomogeneity
```

```
  fun splitList ((L v,w)::vws) =
    let val (A,B) = splitList vws
    in case v of [] => (w::A,B)
        | (h::r) => (A,(P(h,L r),w)::B)
    end
    | splitList (_::vws) = raise Inhomogeneity
    | splitList [] = ([],[])
```

```
  fun unN (N(v,e),w) = (ref{content=v,eq=e},w)
    | unN _ = raise Inhomogeneity
```

```
  fun addClass cs =
    let
      val (i,eqi) = (ref (L cs), ref ([]:class))
      fun pack c = (c, ref (P(N(i,eqi),P(R i,c))))
    in
      map pack cs
    end
```

```

    end

    fun toRefs (ref (P (N(i,eqi), u)))
      = (N(i,eqi), ref u)
      | toRefs _ = raise Undefined

    fun collect f (ws as ((ref (P (R i,c))))::_)
      = (i, f [] ws)
      | collect f _ = raise Undefined

    fun subst ((lst,us:U list)::ps) (l:U list) =
      if (L l) = !lst then us else subst ps l
      | subst [] l = raise Undefined
  in

```

⟨Trennung nach dem Typ des ersten Listenelements⟩≡

```

  fun discriminate ([] : (U * info) list) = [] : class li
  | discriminate [(v,w)] = [[w]]
  | discriminate (vws as ((v,w) :: _)) =
    case v of
      (C _) => CharDiscr.discriminate (map unC vws)
    | (L _) =>
      let
        val (A,B) = splitList vws
        val Bs = discriminate B
      in
        if A = [] then Bs else A :: Bs
      end
    | P(_,_) =>
      let
        val Ss = discriminate (reassoc vws)
        val Ts = map (map unP) Ss

```

```

    in
      List.concat (map discriminate Ts)
    end
  | (S _) =>
    let
      fun toCs (S s) = L (map C (String.explode s)
        | toCs _ = raise Inhomogeneity
      in
        discriminate (map (fn(v,w) => (toCs v,w)) vws)
      end
    | (N _) => (* modulo equivalence on adresses *)
      RefDiscr.discriminate (map unN vws)
    | (R _) => (* modulo identiy on adresses *)
      raise Undefined
      (* If we had defined N : U ref * U ref -> U ref
        N : U ref * U ref list ref -> U, using L : U ref
        we could use the references as the equivalence relation *)
    (R _) =>
      let
        val us = map (fn (R v,w) => !v) vws
        fun toN (R v,w) = (v := ([] : class) ; (N(v,w)
          | toN _ = raise Undefined
        val ws = discriminate (map toN vws)
      in
        "restore the values assigned to vs in vws using
        ws
      end
      *)
  | (Bag _) =>
    let
      val (vs,ws) = ListPair.unzip vws
      fun unBag (Bag v) = v | unBag _ = raise Undefined
    end

```

```

        val svcs = map L (bagsort (map unBag vs))
        val newvws = ListPair.zip (svcs,ws)
    in
        discriminate newvws
    end
| (Set _) =>
    let
        val (vs,ws) = ListPair.unzip vws
        fun unSet (Set v) = v | unSet _ = raise Und
        val svcs = map L (setsort (map unSet vs))
        val newvws = ListPair.zip (svcs,ws)
    in
        discriminate newvws
    end

and bagsort us = (* Standard weak sorter for bags *)
    let
        val byElements =
            discriminate (List.concat (map addClass us))
        val byLists =
            discriminate (map toRefs (List.concat byEle
fun asBag acc [] = acc : U list
    | asBag acc ((ref (P (R i,c))))::ws)
    = asBag (c::acc) ws
    | asBag _ _ = raise Undefined
        val Subst = map (collect asBag) byLists
    in
        map (subst Subst) us
    end
end

```

\langle *Schwache Sortierung von Mengen* $\rangle \equiv$

```
and setsort us = (* Standard weak sorter for sets *)
  let
    val byElements =
      discriminate (List.concat (map addClass us))
    val byLists =
      discriminate (map toRefs (List.concat byEle
fun asSet acc [] = acc : U list
  | asSet [] ((ref (P (R i,d))))::ws) =
    asSet [d] ws
  | asSet (c::acc) ((ref (P (R i,d))))::ws) =
    if d = c then asSet (c::acc) ws
      else asSet (d::c::acc) ws
  | asSet _ _ = raise Undefined
    val Subst = map (collect asSet) byLists
  in
    map (subst Subst) us
  end

fun discriminateToStrings l =
  map (map toString) (discriminate l)
end
```


Kombinatoren zum Aufbau von Trennfunktionen

Wir hatten einen universellen Typ U eigener Typen definiert, um die Trenner d_τ rekursiv über den Typaufbau von τ zu definieren. In d_τ werden ständig Typumwandlungen zwischen U und seinen Summanden durchgeführt.

Wenn man, um diese Umwandlungen zu vermeiden, die in SML eingebauten Typen verwenden will, kann man noch einen Teil der Definition von d_τ programmieren, sofern man polymorphe Rekursion hat (Emms/Leiß, SML+, TCS 1997)

Man definiert dazu *Kombinatoren*, die aus Trennern für bestimmte Typen einen Trenner für einen daraus aufgebauten Typen bilden, z.B. aus Trennern d_σ und d_τ einen Trenner $d_{\sigma \times \tau}$.

$\langle \text{discrpr.sml} \rangle \equiv$

(* Hilfsfunktionen *)

```
fun reassoc (((u,v),w) :: rs) =
  (u,(v,w)) :: (reassoc rs)
```

```
fun concat (w::ws) = w @ (concat ws)
  | concat [] = []
```

```
fun split ((inl x,y)::vs) =
  let val (ls,rs) = split vs in ((x,y)::ls,rs) end
  | split ((inr x,y)::vs) =
  let val (ls,rs) = split vs in (ls,(x,y)::rs) end
  | split [] = ([],[])
```

$\langle \text{discrpr.sml} \rangle + \equiv$

```
(* Kombinatoren zum Aufbau von Discriminatoren,  
    mit Fritz Henglein, 29.4.2005 *)
```

```
fun ProdDisc (da,db) xs =  
  let val ys = reassoc xs  
      val ss = da ys  
      val ws = map db ss  
  in concat ws end
```

```
fun PlusDisc (da,db) xs =  
  let val (ys,zs) = split xs  
      val (w1,w2) = (da ys, db zs)  
  in w1 @ w2 end
```

```
fun UnitDisc (xs: (unit * 'a) list) =  
  [map #2 xs]
```

```
(* Ein Rekursionskombinator ist nicht allgemein,  
    aber im Einzelfall programmierbar: *)
```

```
datatype ('a,'b) sum = inl of 'a | inr of 'b  
datatype 'a tree =  
  Tree of (unit, 'a tree * 'a tree) sum
```

```
fun untree (Tree x,y) = (x,y)
```

```
fun treeDisc xs =  
  (PlusDisc (UnitDisc,  
            ProdDisc(treeDisc, treeDisc)))  
  (map untree xs)
```