Le meilleur défi de Friends of SAS

(probablement)

Par Mathieu Gaouette

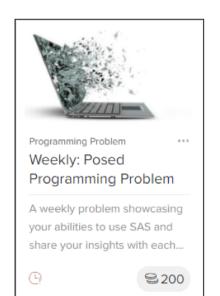


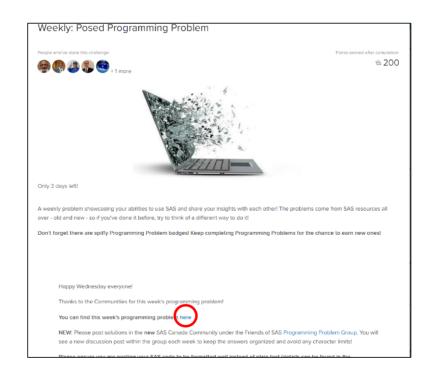
Plan

- Le problème
- Solution avec vecteurs
- Solution avec table hash
- Conclusion



Le problème







Le problème

Clinic Visits



Suppose you want to calculate one patient's 'maximum visit' of clinic, which changes over time.

Below is the raw data that is medical billing data. 'id 1' visited 3 different clinics and 'id 2' visits 2 different clinics.

In'20030124'(third row) 'id 1"s most visiting clinic is '12135' and maximum value is '2'.

However, in '20030607'(6th row), it is '75423' and maximum value is '3'. and in '20030815'(7th row) most visiting clinic is '12135' and '75423', maximum value is still '3'.

'id 2"s most visiting clinic is '25875', and maximum value changes over time. it is 1, 2, 2, 3, 4, 5 in order.

Source

Destination

Count: compte du nombre de visite à cette clinique Max: Nombre maximum de visite à une Clinique pour ce client

In addition, because it is national billing data, the type of clinic is more than 10 thousand, so it is impossible to make virtual array by clinic. SAS can't comprehend that much column.



Le problème

Source

Analyse de la source

- Champ « id » identifiant le patient
- Champ « date » identifiant la date de visite
- Champ « clinic » identifiant la clinique

Ordre des données

• Trié par id et date



Solution avec vecteurs

```
data destination ;
    set source ;
    by id;
    length count 8 max 8;
    retáin Max .
               clin1- clin9999 .;
                                                              Identification source, ordre source, variables
                                                              agrégées et « feuille de notes »
    /* Création vecteurs */
    array cliniques (*) clin1- clin9999;
    /* On vide le vecteur pour un nouveau patient */
    if first.id then do ;
    do _i=1 to 9999;
                                                                       Initialisation « feuille de notes » pour nouveau
             -cliniques(i) = 0;
                                                                       patient.
         max = 0;
     end;
    /* Incrémentation du nombre de visite */
cliniques(clinic) = cliniques(clinic) + 1;
count = cliniques(clinic);
                                                                Incrémentation visite clinique
    max = max(max, count)
                                      Mise à jour maximum visites pour une clinique
run ;
```



Solution avec table hash

```
data destination ;
    set source ;
    by id;
    length count 8 max 8;
    retain max . ;
                                              Identification source, ordre source, variables
    /* Création de la table hash */
    if n = 1 then do;
                                              agrégées et « feuille de notes »
        declare hash h();
h.defineKey('clinic');
h.defineData('count');
        h.defineDone();
    end ;
    /* On vide la table hash pour un nouveau patient */
                                                                Initialisation « feuille de notes » pour nouveau
    if first.id then do;
        h.CLEAR();
                                                                patient.
        max = 0;
    end:
    /* Incrémentation du nombre de visite */
    if h.find() EQ 0 then do;
        count = count + 1;
    end ;
                                                                Incrémentation visite clinique
    else do ;
        count = 1;
    end ;
    h.replace();
                                   Mise à jour maximum visites pour une clinique
   max = max(max, count);
```



Conclusion

Table hash

- Très performant
- Utilisation optimale de la mémoire
- Syntaxe peu intuitive

Vecteurs

- Syntaxe intuitive
- Performant (généralement)
- Peu optimale dans certains cas (comme le notre)



