

Proposition d'interface pour les wrappers de modèle dans CROPTIMIZR

S. Buis

13/05/2020

Problématique

L'objectif de cette note technique est de proposer une solution concernant le design de l'interface des wrappers de modèles de culture dans CROPTIMIZR en prenant en compte les éléments de contraintes et de contexte des différents cas d'étude envisagés pour l'instant.

Les principaux objectifs sont les suivants :

- Faciliter au maximum le développement des wrappers pour une utilisation « standard » de CROPTIMIZR (i.e. hors modes avancés), afin de permettre le développement de nombreux wrappers, y compris de façon occasionnelle pour des modèles jouets très simples.
- Donner un maximum de souplesse dans l'interface des wrappers pour permettre des utilisations avancées et des cas particuliers.
- Faciliter l'utilisation des wrappers et de CROPTIMIZR en ayant des arguments faciles à manipuler
- Optimiser consommation mémoire et temps de calcul sachant que des cas d'étude peuvent être gourmands.

Les cas d'études considérés sont les suivants :

- Estimation de paramètre :
 - Cas standard : on estime des paramètres du modèle en minimisant un critère calculant un écart entre des observations d'une sélection de ses variables de sortie à certaines dates sur plusieurs situations et les valeurs simulées correspondantes.
 - Cas multi-variétal : idem mais les valeurs des paramètres dépendent des situations à simuler.
 - Cas simulations ou observations transformées dynamiquement : Il est possible que les observations ne correspondent pas aux variables simulées par le modèle, ou que la fonction objectif à minimiser soit composée de termes calculés à partir des simulations et observations. Dans ce cas, des template de fonctions de transformations prenant en entrée les résultats des simulations et les observations peuvent être utilisées pour redéfinir les valeurs simulées et/ou observées. Elles sont appelées systématiquement après l'exécution du wrapper et avant le calcul du critère.
 - Méthodes d'estimation : soit purement itératives (une valeur par paramètre est fournie au wrapper pour réaliser les simulations, puis le critère est évalué, le tout répété x fois), soit non itératives (un plan d'expérience numérique composé d'un jeu de valeurs pour les paramètres est fourni au wrapper en une seule fois puis le critère est calculé pour chaque individu du plan d'expérience), soit mixte.

=> Dans ce cas, le wrapper doit retourner les données simulées pour le jeu de valeurs de paramètres fourni, **a minima** pour les situations, dates et variables observées (ou celles nécessaires pour pouvoir effectuer les transformations de variables). En effet, ces données sont ensuite intersectées avec les observations avant le calcul du critère.

- Analyse d'incertitude et de sensibilité (AI / AS):
 - Un plan d'expérience numérique composé d'un jeu de valeurs pour les paramètres est fourni au wrapper en une seule fois. Ces valeurs peuvent éventuellement varier selon les situations à simuler. Le wrapper retourne les résultats simulés pour une sélection de situations, variables et dates ou stades.
- Multi-simulations :
 - Cas assez similaire au précédent mais en général composé d'un grand nombre de situations à simuler et d'un seul jeu de valeurs de paramètres par situation.

Trois questions sont traitées ici : comment designer les structures de données pour 1) le passage des valeurs des paramètres au wrapper, 2) la spécification des situations, variables et dates à simuler et pour lesquelles le wrapper doit renvoyer les résultats et 3) le stockage des sorties simulées.

Valeurs des paramètres

Sous quelle forme transmettre au wrapper les valeurs des paramètres à utiliser pour effectuer les simulations ?

Eléments de contraintes et de contexte :

- Les valeurs des paramètres peuvent varier selon les situations (e.g. estimation multi-variétale, multi-simulations) ... ou pas (estimation de paramètres standard) !
- Il peut y avoir un très grand nombre de valeurs de paramètres => ne pas les répliquer inutilement s'ils sont identiques pour plusieurs situations

Solution proposée

Un tibble appelé `param_values` contenant une colonne par paramètre + 1 colonne "situation" optionnelle.

Si la colonne situation n'est pas fournie les simulations des différentes situations doivent être réalisées pour l'ensemble des situations. Ce sera typiquement le cas en analyse de sensibilité / incertitude et estimation de paramètres standard.

```
## # A tibble: 4 x 3
##   p1    p2    p3
##   <dbl> <dbl> <dbl>
## 1     1     1     1
## 2     2     1     1
## 3     3     2     1
## 4     4     2     1
```

La colonne situation sera fournie si les valeurs des paramètres dépendent des situations. Par exemple:

```
## # A tibble: 4 x 4
##   situation    p1    p2    p3
##   <chr>      <dbl> <dbl> <dbl>
## 1 sit1         1     1     1
## 2 sit2         1     1     1
## 3 sit3         2     2     1
## 4 sit4         2     2     1
```

dans le cas d'une estimation de paramètre multi-variétale, avec `p1` et `p2` paramètres variétaux, `p3` paramètre spécifique, et ("sit1", "sit2") et ("sit3",sit4") des situations observées sur 2 variétés différentes.

Un autre exemple typique sera le cas multi-simulation pour lequel des paramètres varient en fonction de la situation (par exemple condition initiale, stade forcé, ...):

```
## # A tibble: 4 x 4
##   situation    p1    p2    p3
##   <chr>      <dbl> <dbl> <dbl>
## 1 sit1         1     1     1
## 2 sit2         2     1     1
## 3 sit3         3     2     1
## 4 sit4         4     2     1
```

Dans le cas d'un plan d'expérience numérique qui dépend des situations, il y aura des réplifications de valeur dans la colonne situation :

```
## # A tibble: 4 x 4
##   situation    p1    p2    p3
##   <chr>      <dbl> <dbl> <dbl>
## 1 sit1         1     1     1
## 2 sit1         2     1     1
## 3 sit2         3     2     1
## 4 sit2         4     2     1
```

Eléments de Discussion

Cette solution a l'avantage :

- de gérer la plupart des cas sans duplication d'information,
- d'être simple à manipuler pour l'utilisateur qui souhaite exécuter des simulations directement avec le wrapper en forçant des valeurs de paramètres.

Mais a l'inconvénient de devoir traiter 2 formes dans les entrées du wrapper, donc complique un peu (mais pas beaucoup) son développement.

Le plan d'expérience numérique à appliquer par situation sera soit le tibble si pas de colonne "situation", soit une extraction du tibble en fonction de la colonne "situation":

```
# boucle sur les situations à simuler
for (sit in situation_list) {

  # extraction du plan d'expérience numérique en fonction de la situation à simuler
  if (is.null(param_values$situation)) {
    params <- param_values
  } else {
    params <- filter(param_values, situation==sit)
  }

  # boucle sur le plan d'experience numérique par situation
  for (iparams in 1:nrows(params)) {

    # run the model with parameters values defined by params[iparams,]

  }
}
```

Gestion des situations, variables et dates

Comment spécifier au wrapper les situations, variables et dates à simuler et pour lesquelles le wrapper doit renvoyer les résultats ?

Eléments de contraintes et de contexte :

- Dans la plupart des cas d'estimation de paramètres la liste des situations, variables et dates peut être fournie automatiquement puisqu'elle est déduite des observations. Ce n'est cependant pas le cas s'il est nécessaire d'effectuer des transformations (observations qui ne correspondent pas à des variables simulées par exemple).
- Le temps de calcul des wrappers peut parfois être diminué si un nombre réduit de variables est simulé. Cela dépend du modèle (e.g. requête dans une base de donnée pour lire les résultats pour APSIM, écriture et lecture +/- lourde dans les fichiers de sortie des modèles si on peut leur spécifier les variables à écrire, ...).

Solution proposée

Trois arguments :

- `sit_names` : vecteur de nom des situations à simuler
- `var_names` : vecteur de nom des variables à simuler
- `dates` : vecteur des dates à simuler
- `stages` : vecteur des stades à simuler
- `sit_var_date_mask` : masque des situations, variables et dates à simuler (homogène au format des observations)

Tous ces arguments seraient optionnels : cela permet de construire rapidement un wrapper qui fonctionne. Par exemple:

```
simple_wrapper <- function(param_values, model_options, ...) {  
  
  workspace_path <- model_options$workspace_path  
  
  # lecture de la liste des situations à simuler dans le workspace  
  
  # simulation de toutes les variables du modèle pour toutes dates sur toutes les situations  
  # pour les valeurs prescrites des paramètres  
  
}
```

permettra facilement de faire de l'estimation de paramètre, de la multi-simulation ou même AI/AS. Dans ces deux derniers cas, l'utilisateur devra par contre lui-même extraire les résultats qui l'intéressent en terme de dates.

Un cas de développement de wrapper un peu plus avancé mais a priori pas très compliqué à réaliser serait de prendre en compte les arguments `var_names`, `dates` et `stages` : l'utilisateur définit les situations à simuler via son workspace, comme précédemment, mais le wrapper ne renvoie les résultats que pour les variables, dates et stades sélectionnés, ce qui évite des manipulations pour le post-traitement des résultats. Le caractère optionnel (au moins de `dates` et `stages`, puisqu'il ne sont pas fournis en estimation de paramètres) doit cependant être géré dans le wrapper.

L'argument `sit_names` est également assez facile à intégrer dans le wrapper pour sélectionner une sous-liste de situations à simuler au sein de celles définies dans le workspace, bien qu'une vérification de la cohérence des deux s'impose.

L'utilisation de `sit_var_date_mask` n'est peut-être à conseiller que pour les modèles pour lesquels la gestion des sorties simulées est très coûteuse. Cet argument permet en effet de régler plus finement la liste des dates et variables à retourner en fonction des situations à simuler, mais la multiplicité des argument rend le wrapper plus complexe à implémenter.

Si tous les arguments sont fournis et utilisés dans le wrapper, `sit_var_date_mask` doit être prioritaire.

Les noms des arguments sont fixés pour faciliter le passage d'un modèle à l'autre et pour fournir automatiquement les informations dans le cas estimation de paramètres (on aurait pu sinon laisser les développeurs de wrapper gérer cela via les options du wrapper).

Les arguments `sit_names`, `var_names` et `sit_var_date_mask` sont automatiquement fournis d'après les observations dans le cas estimation de paramètre. Si l'utilisateur souhaite faire une transformation des sorties simulées ou des observations et a besoin pour cela de simuler des variables non observées, alors il peut définir ces arguments en entrée d'`estim_param`, ils seront passés tels quels au wrapper (peut-être via une solution du type ... pour `estim_param` : tous les arguments non définis dans l'interface de `estim_param` pourraient être passés à l'ensemble des fonctions utilisateurs (wrapper, fonction de transformation, ...)).

Eléments de Discussion

Cette solution a l'avantage :

- de fonctionner dans les cas standards même si ces arguments ne sont pas gérés dans les wrappers à partir du moment où il simule au moins les situations, variables et dates nécessaires,
- d'être assez souple pour le développement des wrappers (seuls certains de ces 3 arguments peuvent être pris en compte dans le wrapper en fonction des cas d'études traités, méthodes utilisées, optimisation du temps de calcul du wrapper ...),
- de ne pas demander à l'utilisateur de spécifier les situations, variables et dates à simuler dans le cas estimation de paramètre sans transformation.

Mais a l'inconvénient :

- de complexifier le développement des wrappers si on veut qu'ils fonctionnent de façon optimale dans tous les cas (il faut traiter les différents arguments, leur caractère optionnel, compatibilité avec ce qui peut être réellement fourni par le modèle, ...),
- de multiplier la liste des arguments des wrappers

Stockage des sorties simulées

Comment stocker les valeurs de différentes variables simulées pour différentes situations, dates / stades et valeurs de paramètres ?

Eléments de contraintes et de contexte :

- la liste des variables peut différer selon les situations, et la liste des dates / stades selon les situations et variables.
- Comme il peut y avoir beaucoup de situations et/ou variables et/ou dates et/ou valeurs de paramètres, l'optimisation de la quantité de mémoire utilisée et du coût traitement de ces données est important
- Les tests réalisés par Thomas Robine (cf. https://github.com/SticsRPacks/documentation/tree/master/tests_data_structure/Output_tests/output_tests.html) montrent qu'un stockage sous forme d'un seul tibble est efficace en temps de traitement et ergonomique mais peut poser des problèmes mémoires pour des tailles importantes, mais réalistes, des différents cas d'étude. Par ailleurs, l'implémentation actuelle (liste de liste de `data.frame`) n'est pas satisfaisante en raison de son manque d'ergonomie.

Solution proposée

Une solution, intermédiaire entre la solution actuelle et la solution « un seul tibble », est une liste de tibble par situation. Le tibble contiendrait les résultats des simulations des différentes variables (en colonne), dates et stades pour la situation concernée et pour l'ensemble des valeurs du plan d'expérience. Une colonne `Doe` contiendrait le numéro de ligne du plan d'expérience numérique, une colonne `Date` la date et une colonne 'Stage' le (ou les) noms des stades pour lesquels la variable est demandée.

Dans un cas multi-simulation simple de 2 situations, pour lesquelles on demande 2 dates et 1 stade “rec” on aurait par exemple :

```
## $sit1
## # A tibble: 3 x 4
##   Doe Date                var1 Stage
##   <dbl> <dtm>                <dbl> <chr>
## 1     1 2020-01-01 00:00:00    1.1 <NA>
## 2     1 2020-03-01 00:00:00    1.2 <NA>
## 3     1 2020-07-12 00:00:00    1.3 rec
##
## $sit2
## # A tibble: 3 x 4
##   Doe Date                var1 Stage
##   <dbl> <dtm>                <dbl> <chr>
## 1     1 2020-01-01 00:00:00    1.1 <NA>
## 2     1 2020-03-01 00:00:00    1.2 <NA>
## 3     1 2020-05-01 00:00:00    1.3 rec
```

Dans un cas AI /AS avec une situation et deux stades, on aurait :

```
## $sit1
## # A tibble: 6 x 4
##   Doe Date                var1 Stage
##   <dbl> <dtm>                <dbl> <chr>
## 1     1 2020-03-11 00:00:00    1.1 lev
## 2     1 2020-07-12 00:00:00    1.3 rec
## 3     2 2020-03-10 00:00:00    1.4 lev
## 4     2 2020-07-10 00:00:00    1.6 rec
## 5     3 2020-03-09 00:00:00    1.6 lev
## 6     3 2020-07-09 00:00:00    1.8 rec
```

Eléments de Discussion

Cette solution a l’avantage :

- de résoudre les problèmes mémoires rencontrés avec la solution « un seul tibble » en répartissant les données simulées par variable en différents vecteurs (1 par situation)
- d’être plus ergonomique que la solution actuellement implémentée (liste de liste) :
 - extraction de simulations d’une variable pour une date ou un stade et différents individus du DOE (cas typique AI/AS) : `filter(data,Date=...)$var` ou `filter(data,Stage=...)$var`
 - extraction des simulations d’une variable à toutes (ou une liste de) dates (cas typique estimation de paramètres ou multi-simulation, 1 seul DOE) : `data$var` ou `filter(data,Date=...)$var`
 - L’extraction de simulations d’une variable pour différentes situations pourra se faire soit via la transformation de la liste en un seul tibble si la mémoire le permet (`bind_rows`) soit via une instruction type `apply`.

Exemple d’un wrapper simple qui implémente ces solutions

CODAGE DU WRAPPER (sans `sit_var_date`)

wrapper : équation bonzai, on définit les situations et T° en dur dans le code ...

Dans les différents cas (avec ou sans prise en compte des arguments)

Illustrer uTILISATION mode AI/AS, multisimulation, estimation de paramètres

Remarques

Les cas intercropping et successions culturales ne sont pas (encore) pris en compte dans cette réflexion.