PharmaSUG 2019 - Paper DS-070

Patient Narrative Generation, PROC Report/Stream vs R - Markdown

Yuanyuan Gu and Mengmeng Zhao, dMed Biopharmaceutical Co., Ltd.

ABSTRACT

Patient narrative is essential in both drug safety and patient safety monitoring during clinical trials. As part of clinical research result or CSR section, it costs much time and effort to generate patient narratives. Meanwhile, it is often cumbersome to manage, which may take up resources and affect timelines. Hence relevant automated forms or tools are urgently needed to increase efficiency and to reduce error rate. This presentation will give an overview of the pros and cons of SAS proc steps (Proc REPORT/STREAM) vs R - Markdown for the automation to solve this problem, showing how to use SAS to generate narrative RTF files separately and automatically for each patient and to batch convert .rtf to .docx files. Based on the current steps, it is possible to explore further improvement processes, such as which way is more suitable for statistical programmers to generate patient narratives and to manage a certain frequency of output results, and whether there exists more reasonable ways to tabulate the effects finally demonstrated in CSR.

INTRODUCTION

Patient narrative is an important part of the appendix (i.e. appendix 11 or 12.) to clinical study reporting, submissions of which would require demographic information, concomitant medications, medical history, exposures, procedures, laboratory test results, serious adverse events and adverse events that the principal investigators hope to report, and so on. The medical writer usually generates narratives by manually writing in plain text, which has several shortcomings. First, the number of required medical writers is large, and management and arrangement of workload is a problem. Second, dealing with above data requires medical writers to have certain qualifications and experiences. Third, it is difficult to control delivery time, accuracy, and quality control procedures. Hence, we need to find out a new method to do this. Developing an automated narrative report generation tool with robust SAS or R methods would bring benefits in saving a large amount of human resource, where automation can be repeated in various studies. Proc REPORT uses macro variables and macros; Proc STREAM uses macros combined with HTML statements and R - Markdown uses Markdown syntaxes and chunks of R code, all of which could automatically generate patient narratives. We will introduce how these three methods automatically generate patient narratives in batches, and comprehensively explore the advantages and disadvantages of the three methods from the aspects of editability, efficiency, format customization, and output format.

THE PROCEDURE SYNTAX

PROC REPORT

The COLUMN statement is used to list variables to specify columns. Each column has a DEFINE statement that sets the value of an attribute.

```
PROC REPORT DATA=datasetname <options>;

COLUMN variable list and column specifications;

DEFINE column / define type and column attributes;

DEFINE column / define type and column attributes;

...

RUN;
```

PROC STREAM

Specifies an external file with the "OUTFILE=" keyword, as well as optional arguments. Any arbitrary text to be output should begin with "BEGIN" and ends with four semicolons ";;;;". There is no "RUN" or "QUIT" statement.

```
PROC STREAM OUTFILE= fileref <option(s)>; BEGIN
    text-1
    <text-n>
;;;;
```

R-MARKDOWN

The file contains three sections of contents: An (optional) YAML header embedded by ---s, R code chunks embedded by ```s, and texts mixed with inline codes which output simple text formatting.

```
title: "Untitled"
author: "author"
date: "mm dd,yyyy"
output: word_document
```

```
---
```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = TRUE)

Text `r `
```

### DATA PROCESSING

The narrative report is generated using SDTM data sets, which is standardized data sets of CDISC standard. The data sets involve DM, EX, AE, CM, PR, LB, etc. The following are main steps in data processing.

### STEP 1: CATEGORIZING OF NARRATIVE REASON

```
/* Category of narrative reasons*/
data ae_all;
 set ae_all;
 length AEREASON $200;
 then AEREASON="SAE"; /*AE severity is yes*/
 if kindexc(AESER, 'Y是')
 if kindexc(DCSREAS,'不良事件') then AEREASON='Withdrawal';/*AE action taken with permanent withdrawal*/if kindexc(AESDTH,'Y是') then AEREASON="Death";/*AE results in death*/
 if not missing (AESICAT) then AEREASON="AESI";/*AE specified for each study*/
 if not missing (AEREASON);
run;
∃data ae;
 set ae_all;
 AEDECOD_REASON=strip(AEDECOD)||' ('||strip(AEREASON)||')'; /*combine AEDECOD and AEREASON with "()"*/
 proc sort;
 by USUBJID AESEQ;
run;
```

#### STEP 2: COUNTING AND LISTING UNIQUE SUBJECT

Proc SQL is used to create macro variables SUBJ LST and SUBJ CNT.

#### STEP 3: OBTAINING DATA STRUCTURE FOR TABLE

In order to make COL1 display the label, COL2 display the value and ORD display the order of desired information, we need to transpose the subject's ("ABC-xxx-10011008") multiple records to one record.

```
/*Data structure for table display;*/
data part1(where=(not missing(COL1)));
 set mg_all(where=(usubjid="ABC-xxx-10011008"));
 if n =1;
 length COL1-COL2 $300;/*define length two colum named as COL1 and COL2*/
 output;
 ord=1; col1='方案编号'; col2=studyid; output;
 ord=2; col1='受试者编号'; col2=subjid; output;
 ord=3; col1='病例描述原因'; col2=aedecod_reason; output;
 ord=4; col1='国家'; col2=country; output;
 ord=5; col1='随机化治疗组'; col2=arm; output;
 ord=6; col1='实际治疗组'; col2=actarm; output;
 ord=7; col1='随机化日期'; col2=dsstdtc; output;
 ord=8; col1='首次用药日期'; col2=rfxstdtc; output;
 ord=9; col1='用药时间(天)'; col2=exdur; output;
 ord=10; col1='性别'; col2=sex; output;
 ord=11; col1='年龄(岁)'; col2=strip(put(age,??best.)); output;
 ord=12; coll='种族'; col2=race; output; ord=13; coll='既往病史'; col2=mhdecodg; output; ord=14; coll='既往用药史'; col2=cmtrts; output;
 keep ORD COL1 COL2;
run;
data part1;
 if missing(COL2) then COL2="NA";/*display missing value as "NA"*/
run;
```

## STEP 4: OBTAINING DATA STRUCTURE FOR DESCRIBE SENTENCE

The logic of free text processing is to generate a variable (SENTCx) for each sentence and combine all SENTCx variables. Macro calls, macro variables and IF statements could automatically achieve statement judgment, while loops and macros could also be used to batch run. The %SYSFUNC function is used to check if variables exist. Here is part of the processing code.

```
%macro out subi; /*define macro out subi*/
%macro out_subj:/*define macro out_subj*/
%do i=1 %to &subj_cnt.; /*use macro variable subj_cnt to do loop*/
%let subj=%scan(&subj_lst,&i.,'*'); /*use macro variable subj_lst to get current usubjid*/
%let dsid=%sysfunc(open(mg_all,i));
%let AESEV=%sysfunc(varnum(&dsid,AESEV)); /*check if AESEV is exist*/
 %let AETOXGR=%sysfunc(varnum(&dsid,AETOXGR));/*check if AETOXGR is exist*/
%let dsid=%sysfunc(close(&dsid));
 a ae_standard;
length SENTENCE $10000 SENTC1-SENTC3 TEMP1-TEMP3 $200;
set mg_all (where=(USUB3ID="ssubj."));/* read in one person records*/
SENTC1=catx('','这位',AGE,'岁的',SEX,'性受试者',"sPT_TYPE.",',入组本研究后于',DSSTDTC,'被随机化后接受',ACTARM,'用药。');
if not missing(AESTDTC) and not missing(AESTDY) then TEMP1=catx('','在研究的第',AESTDY,'天(',AESTDTC,')');
else if not missing(AESTDTC) and missing(AESTDY) then TEMP1=catx('','在','(',AESTDTC,')');
else tEMP1='未报告日期';/*generate the first sentence*/
%if &AESEV ne 0 %then %do;/*if contain AESEV then do the following code */
 if kindexc(AESER,'Y是') then TEMP2='严重不良事件';
 else if kindexc(AESER,'N否') then TEMP2='不良事件';
 if not missing(AESEV) then do; /*check if AESEV missing*/
 SENTC2=catx('',TEMP1,', 患者发生',TEMP2,'严重程度为',AESEV,'逐字报告为',AEDECOD);
end;
 data ae standard;
 end:
 if missing(AESEV) then do; /*check if AESEV missing*/
SENTC2=catx('',TEMP1,', 患者发生',TEMP2,'(严重程度未报告),逐字报告为',AEDECOD);
 %end:
 %if &AESEV eq 0 %then %do;/*if not contain AESEV then do the following code */
SENTC2=catx('',TEMP1,',患者发生',TEMP2,'逐字报告为',AEDECOD);
 %end;/*generate the second sentence*/
%if &AETOXGR ne 0 %then %do;
 if not missing(AETOXGR) then sentc3=cats(', 当时的CTCAE为',AETOXGR,'级,'); else if missing(AETOXGR) then sentc3=', CTCAE未报告,';
 %end;
%if &AETOXGR eq 0 %then %do;
 SENTC3=', CTCAE未报告,';
%end;/*generate the third sentence*/
SENTENCE='^R" "'||compress((catx('',of SENTC2-SENTC3))," ");/*Indent 4 spaces before the paragraph*/
 run;
 %end;/*end loop*/
 %mend; /*end macro*/
```

### STEP 5: FINALIZING DATA STRUCTURE

### 1. The dataset structure for table

COL1	COL2	ord
方案编号	ABC-xxx	1
受试者编号	10011008	2
病例描述原因	贫血 (SAE)	3
国家	CHN	4
随机化治疗组	ABC 300mg	5
实际治疗组	ABC 300mg	6
随机化日期	2018-03-13	7
首次用药日期	2018-03-14	8
用药时间(天)	119	9
性别	女	10
年龄(岁)	52	11
种族	亚洲人	12
既往病史	卵巢癌、胆结石手术	13
既往用药史	紫杉醇+卡铂	14

## 2. The dataset for sentence

usubjid	sentence	title
ABC-xxx-10011008	这位52岁的女性受试者被诊断为xxx患病者,入组本研究后于2018-03-13被随机化后接受ABC300mg用药。	病例描述
ABC-xxx-10011008	^R"	病例描述

## **APPLICATION TO PATIENT NARRATIVES**

### PROC REPORT

To allow for various forms of layout, we need to carry out at least two Proc REPORT processes. Proc REPORT is often used to output reports in .rtf format. This format enables customized headers and footers and is convenient for third-party editing.

## **Report Code**

```
options orientation=portrait papersize=A4 topmargin=lin bottommargin=lin leftmargin=lin lin lect=32767;/*define margins*
title1 j=l h=10.5 pt font = 'Times New Roman' '獨例描述' j=r "产品编号: &productid.: 研究编号: &studyid.";/*header*/
title2 j=r h=10.5 pt font = 'Times New Roman' "版本号: 1.0";/*header*/
footnote j=l h=10.5 pt font = 'Times New Roman' "严格保密" j=c h=9pt "<草稿版>" j=r "第 ^{thispage} 页";/*footnote*/
ods noresults;
ods oresults;
ods scapechar='n';
ods listing close;
ods rtf file="c;\usersname\Narrative.rtf";/* output path ;startpage=now*/

proc report nowd da a=part1 style=[background=white] list;
 column col1 col2;
 define col1 / display " " style(column)=[just=left fontsize=12pt width=35%];
 define col2 / display " " style(column)=[just=left fontsize=12pt width=64%];

run;

proc report data=part2 style=narrative list;
 column title sentence;
 define title/group noprint " " style(column)=[pretext='^R"" fontsize=12pt just=left fontweight=bold width=99%]; /*indent 4 space
 define sentence/group display " " style(column)=[pretext='^R"" fontsize=12pt just=left width=99%]; /*indent 4 spaces*/
 compute before title/style=[asis=on fontweight=bold fontsize=12pt];
 line @1 title $40.;
 endcomp;
 ods rtf startpage=no;/* not force a page break*/
run;

ods rtf close;
ods listing;
ods results on;
```

## Output (.rtf)

It is easy to convert .rtf to .docx document by tools.

病例描述

产品编号: ABC; 研究编号: ABC-xxx 版本号: 1.0

### 患者信息及用药情况

方案编号	ABC-xxx
受试者编号	10011008
病例描述原因	贫血 (SAE)
国家	CHN
随机化治疗组	ABC 300mg
实际治疗组	ABC 300mg
随机化日期	2018-03-13
首次用药日期	2018-03-14
用药时间(天)	119
性别	女
年齢(岁)	52
种族	<b>亚洲</b> 人
既往病史	卵巢癌、胆结石手术
既往用药史	紫杉醇+卡铂

#### 』 病例描述

这位52岁的女性受试者被诊断为XXX患病者,入组本研究后于2018-03-13被随机化后接受ABC300mg用药。

在研究的第71天(2018-05-23),患者发生逐字报告为贫血,当时的CTCAE为3级,不良事件需要住院治疗/需要延长住院治疗。对患者给予以下矫正治疗措施研究药物因此不良事件暂停使用。患者在2018-06-23从不良事件中恢复。研究者对于不良事件与研究药物相关性的判断如下:与研究药物相关。

## **PROC STREAM**

Proc STREAM uses HTML tags, macros, and macro variables to incorporate flexible inclusion of texts, tables, charts, and lists in a patient narrative document. If a patient narrative document contains only simple text, Proc STREAM could

output in rtf and HTML format. But if the report contains not only texts but also tables or figures, the output in HTML format is the most appropriate choice.

Tags	Function
 	carriage return
•	bullet symbol
	non-breaking space

Table 1 HTML Tags

## **Output Code**

%gp macro calls for individual patient characteristics to get the parameter value for parameters, e.g. calling %gp(SEX) returns the value " $\checkmark$ ".

## Output(.html)



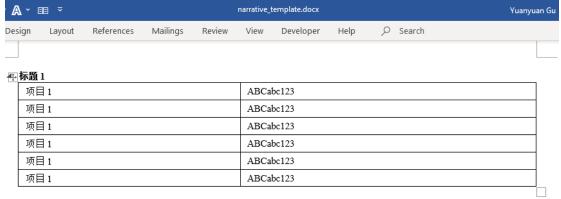
### R-MARKDOWN

It is a great way to create dynamic documents with embedded chunks of R codes. The document is then 'knit' using knitr to create a HTML file. Markdown documents can be converted to many other file types including .html, .pdf and .docx.

## **Report Code**

## Part 1 YAML Header

If the final document needs no title, author or date information, we could delete those keywords. R-Markdown could output both html. and docx. If document style editing is needed, save a style you wanted as your style reference document and use it as the 'reference\_docx:' argument in the front matter.



#### 标题 2

测试测试 test test测试测试 test test

#### Reference Document

```
header-includes:
output:
html_document: #output in html
word_document: #output in docx.
reference_docx: "narrative_template.docx" # format refrence document
```

### Part 2 Code Chunk

Option	Function
include = FALSE	prevents code and results from appearing in the finished file.
echo = FALSE	prevents code, but not the results from appearing in the finished file.
message = FALSE	prevents messages that are generated by code from appearing in the finished file.
warning = FALSE	prevents warnings that are generated by code from appearing in the finished.
fig.cap = ""	adds a caption to graphical results.

## Table 2 Chunk Option

```
<!-- ## read rawdata -->
```{r setup, include=FALSE }
knitr::opts_knitiset(root.dir = "C:/Users/yuanyuan.gu/Desktop/narrative") #Set path
knitr::opts_chunkSset(echo = TRUE)
sessionInfo()
Sys.setloade("LC_CTYPE", "Chinese")
options(Encoding="UTF-8") #Set to Chinese environment
library(knitr) #read library
library(mintr) #read library
library(markdown)
library(markd
```

Part3 Text

Using r makes it easy to update the report to refer other functions of R. r ae_t\$AGE is similar to macro variable in SAS, which returns '65' in this example. & emsp use to insert spaces in .html.

```
## 病例描述 <!-- title -->
&emsp;&emsp;这位'n ac_t$AGE '岁的'n ac_t$SEX '性受诱者被诊断为XXX患病者,入组本研究后于'n ac_t$OSSTDTC '被随机化后接受'n ac_t$AGM'用药。

&emsp;&emsp;在研究第 n ac_t$AESTDV '天 ('n ac_t$AESTDTC'). 患者报告 'n ifelse(ac_t$AESER =="N", "'凡良事件")" 是严重的CTCAE等级为'n ac_t$AESTDXGR
'级,逐字报告为'n ac_t$AEDCDO'. 编码为'n ac_t$AEBDDSYS'.'n ifelse((is.na(AERES)),"",paste("该产重不良事件",AERES,sep = ": ")) '.' r ifelse(ac_t$AEACN =="",",paste("该究药物因此不良事件",ac_t$AEACN,sep = ": ")) '.' r AEOUT1 '. 研究者对于不良事件与研究药物相关性的判断如下:'n ac_t$AEREL'.
```

Batch Run Code

Output-1(.docx)

患者信息及用药情况

70 F 10 10 10 10 10 10 10 10 10 10 10 10 10		
方案编号	ABC-xxx	
受试者编号	ABC-xxx-10011008	
病例描述原因	贫血 (SAE)	
国家	CHN	
随机化治疗组	ABC 300mg	
实际治疗组	ABC 300mg	
随机化日期	2018-03-13	
研究数据收集截至日期	2018-07-10	
实际用药时间(天)	92	
性别	女	
年齢 (岁)	52	
种族	人帐亚	
既往用药史	紫杉醇+卡铂	
既往骨髓抑制史	中性粒细胞减少症	

病例描述

这位 52 岁的女性受试者被诊断为 XXX 患病者,入组本研究后于 2018-03-13 被随机化后接受 ABC 300mg 用药。

在研究第 71 天(2018-0523),患者报告严重不良事件最严重的 CTCAE 等级为 3 级,逐字报告为贫血,编码为血液及淋巴系统疾病。该严重不良事件:需要住院治疗/需要延长住院治疗。研究药物因此不良事件:研究药物暂停。患者在研究第 102 天(2018-06-23)痊愈/已恢复。研究者对于不良事件与研究药物相关性的判断如下:有关。

Output-2 (.html):



患者信息及用药情况

ABC-xxx ABC-xxx-10011008 方案编号 受试者编号 病例描述原因 贫血 (SAE) CHN ABC 300mg 随机化治疗组 实际治疗组 ABC 300mg 随机化日期 2018-03-13 研究数据收集截至日期2018-07-10 实际用药时间(天) 92 性別 年齢 (岁) 种族 **亚洲人** 既往用药史 既往骨髓抑制史 紫杉醇+卡铂 中性粒细胞减少症

病例描述

这位52岁的女性受试者被诊断为XXX患病者,入组本研究后于2018-03-13被随机化后接受ABC 300mg用药。

在研究第71天(2018-05-23),患者报告严重不良事件最严重的CTCAE等级为3级,逐字报告为兹血,编码为血液及淋巴系统疾病。该严重不良事件:需要住院治疗需要延长住院治疗。研究药物因此不良事件:研究药物暂停。患者在研究第102天(2018-06-23)痊愈/已恢复。研究者对于不良事件与研究药物相关性的判断如下:有关。

ADVANTAGE & DISADVANTAGE

The advantages and disadvantages of Proc REPORT, Proc STREAM and R-Markdown are as following:

- 1. Output file: The most appropriate output formats for Proc REPORT and Proc STREAM are .rtf, and .html. R-Markdown supports a variety of file formats, including .docx, .html and .pdf format. However, in practice, medical writers prefer to use .docx as they hope to edit in narrative reports while reviewing.
- 2. Dataset: All of three ways support the types of .sas7bdat, of course, also support .excel if needed.
- 3. Styles or templates: Proc STREAM does not support to adjusted or customized layout or format. Proc REPORT supports customized layout or format by options and keywords. R-Markdown could also be implemented by reference documents or VBA macro, but it is more complicated than Proc REPORT.
- 4. Efficiency: Proc REPORT, Proc STREAM and R-Markdown can automatically generate multiple patient narratives in a matter of minutes.
- 5. Cost of software: R-Markdown is free but Proc REPORT and Proc STREAM are paid.
- Debug processing: Proc REPORT, Proc STREAM and R-Markdown use consoles or logs to locate problems which will greatly facilitate debugging.

CONCLUSION

The methods presented in this paper are very useful to generate narrative report. By using the Proc REPORT, Proc STREAM or R - Markdown, this paper provides the possibilities in producing consistent, high-quality patient narratives without unnecessary increases in costs and threats to timelines. Of course, this is only the initial configuration and debugging phase before the final output. The following service steps include: The medical writers manually edit and supplement the information as required which cannot be done automatically. The combined output is then provided to physicians for review. The draft is finalized according to the opinions of all parties. Need translation(optional). In general, the first step of development of patient narrative generation brings great convenience to medical writers: it provides efficient, accurate, repeatable and editable results for the subsequent final report.

REFERENCES

1. The REPORT Procedure (SAS Institute Inc)

https://documentation.sas.com/?docsetId=proc&docsetTarget=p0bqogcics9o4xn17yvt2qjbgdpi.htm

2. The STREAM Procedure (SAS Institute Inc)

https://documentation.sas.com/?docsetId=proc&docsetTarget=n12zrkr08eiacmn17lcv4fmt79tb.htm

3.R-Markdown (RStudio Inc.)

https://rmarkdown.rstudio.com/

- 4. Joseph Hinson. "Proc STREAM: The Perfect Tool For Creating Patient Narratives", Proceedings of the SAS Global Forum 2015 Conference, Paper 1738-2014. https://pharmasug.org/proceedings/2015/AD/PharmaSUG-2015-AD03.pdf
- 5. Henderson, Don. "PROC STREAM and SAS® Server Pages: Generating Custom

HTML Reports" Proceedings of the SAS Global Forum 2014 Conference, Paper 1738-2014. http://support.sas.com/resources/papers/proceedings14/1738-2014.pdf

6. Renuka Tammisetti and Karthika Bhavadas. "A Guide to Programming Patient Narratives" Proceedings of the SAS Global Forum 2017 Conference, Paper PO24. https://www.pharmasug.org/proceedings/2017/PO/PharmaSUG-2017-PO24.pdf