

PHUSE US Connect 2021

14-18 June 2021 Paper DH12 – Data Handling (DH)

The Encoding Dilemma

Angelo Tinazzi, Cytel Inc. Geneva – Switzerland angelo.tinazzi@cytel.com



Introduction

Encoding Dilemma?

Encoding in a Nutshell

What is Encoding Encoding Methods

Encoding and SAS

Encoding vs Transcoding

Dealing with Transcoding in SAS

wlatin1 → UTF-8, UTF-8 → wlatin1, Failing transcode to wlatin1

Transcode datasets to your local encoding

Summary of Options

Encoding and Data Submission

FDA, PMDA and NMPA Requirements
The SAS XPT Challenge with Multiple Byte Character Set (MBCS)

Conclusions

Summary

Other aspects to consider → References

AGENDA



Introduction

Encoding Dilemma.....what's that "weird" messages are telling me?

Introduction

What's that "weird" messages are telling me?





NOTE: Data file SJISRAW.ADVEVE_1_0.DATA is in a format that is native to another host, or the file encoding does not match the session encoding.

Cross Environment Data Access

will be used, which might require additional CPU resources and might reduce performance.

Encoding?
Cross Environment Data Access?



WARNING: Some character data was lost during transcoding in the dataset SJISRAW.ADVEVE_1_0. Either the data contains characters that are not representable in the new encoding or truncation occurred during transcoding.

Transcoding? Truncation?



ERROR: Some character data was lost during transcoding in the dataset TOCREATE.TESTTOUTF8.

Either the data contains characters

that are not representable in the new encoding or truncation occurred during transcoding.

NOTE: The DATA step has been abnormally terminated.

NOTE: The SAS System stopped processing this step because of errors.

NOTE: There were 2 observations read from the data set TOCREATE.TESTTOUTF8.

WARNING: The data set WORK.TESTTOUTF8 may be incomplete. When this step was stopp

NOTE: DATA statement used (Total process time).

ERROR!!!



Encoding in a Nutshell

What is Encoding

Encoding Methods

Encoding (and Transcoding) in a Nutshell What is Encoding?



Characters are stored as a series of bytes, where 1 byte = 8 bits e.g. "Character" is intended a letter, a number, a symbol, etc.

Encoding is the way a computer interprets and represents the "Characters" within the data by assigning an integer number to each "Character" (code page)

"Sequence" of bits can represent an integer number i.e. with 1 one byte, numbers between 0 and 255 (2^8-1=256)

A **Character Set** is the set of characters used by a language or a group of language

e.g. English or more in general western languages, Chinese, Japanese, etc.

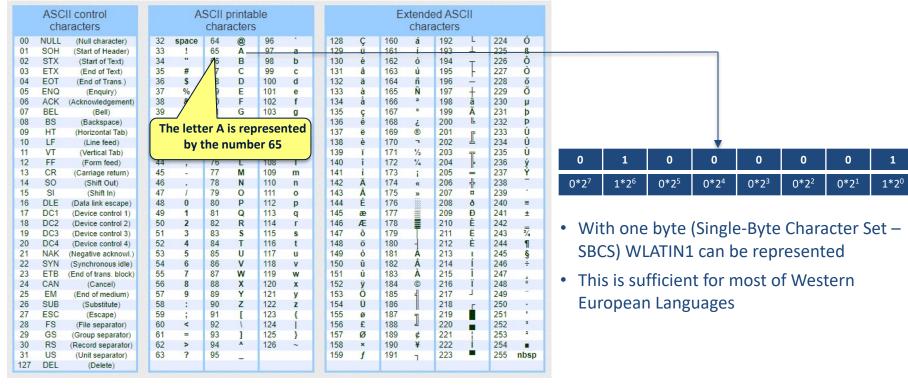
An **Encoding Method** is a set of rules that assign the numeric representations to the set of characters

SAS® 9.4 National Language Support (NLS): Reference Guide, Fifth Edition, SAS https://documentation.sas.com/api/collections/pgmsascdc/9.4_3.5/docsets/nlsref/content/nlsref.pdf?locale=en#nameddest=n04275wmuchnqjn1hfx0yzpopw5a

Encoding (and Transcoding) in a Nutshell

Encoding Methods Example WLATIN1 (Windows Latin 1, US and Western European)





https://theasciicode.com.ar/

Encoding (and Transcoding) in a Nutshell Encoding Methods SBCS vs DBCS vs MBCS



A set of rules that assign the numeric representations to the set of characters

- Size of the encoding (type and number of characters represented)
- Number of bytes (bits) used
 - Single-Byte Character Set (SBCS)
 - Double-Byte Character Set (DBCS)
 - Multiple-Byte Character Set (MBCS)

Single-Byte Character Set (SBCS)

- wlatin1 or wlatin9 for Western European
- wlatin2 for Eastern European

Double-Byte Character Set (DBCS)

- > Shift-Jis for Japanese
- **Big5** or **EUC-CN** for Chinese

Encoding (and Transcoding) in a Nutshell

Encoding Methods SBCS vs DBCS vs MBCS



A set of rules that assign the numeric representations to the set of characters

- Size of the encoding (type and number of characters represented)
- Number of bytes (bits) used
 - Single-Byte Character Set (SBCS)
 - Double-Byte Character Set (DBCS)
 - Multiple-Byte Character Set (MBCS)

Multiple-Byte Character Set (DBCS)

- > UTF-8 Unicode Universal Character set Transformation
- An encoding that attempt to represent all characters in all languages
- ➤ MBCS (it uses 4 bytes)
- >120000 characters. For some 1 byte is still sufficient for many others not



Encoding vs Transcoding

Dealing with Transcoding in SAS

wlatin1 → UTF-8, UTF-8 → wlatin1, Failing transcode to wlatin1

Transcode datasets to your local encoding

Summary of "Options"

Encoding vs Transcoding – Detecting your SAS Encoding



Encoding in SAS establishes the default working environment for your SAS session i.e. WLATIN1 in a Windows Os for US and Western European Languages

The encoding applied by default in your SAS environment determines the default encoding of generated permanent datasets by your SAS programs i.e. SDTM

Checking the encoding used by default by your SAS configuration

ENCODING=WLATIN1 Spe

proc options option=encoding;
run;

```
ENCODING=WLATIN1 Specifies the default character-set encoding for the SAS session.

NOTE: PROCEDURE OPTIONS used (Total process time):
real time 0.01 seconds
cpu time 0.01 seconds
```

Encoding and SAS Encoding vs Transcoding



Transcoding is the process of converting data from one encoding to another

The same character could have different numeric representations in two different encodings

SAS Cross-Environment Data Access (CEDA) in many cases does the transcoding for you

Dealing with Transcoding in SAS



Checking the dataset encoding

```
/*Checking encoding*/
%macro encod(libin=,dsin=);
 %put;
 %put ----CHECKING FOR ENCODING-----;
 %put &libin..&dsin;
 %let dsid=%sysfunc(open(&libin..&dsin,i));
 %put ENCODING is %sysfunc(attrc(&dsid,encoding));
 %put ------
 %put;
 %put;
 %let dsclose=%sysfunc(close(&dsid));
%mend;
```

Dealing with Transcoding in SAS



Checking the dataset encoding, 4 examples

- UTF8 a SAS library with SDTM datasets generated in an UTF-8 Unicode encoding environment
- > SJID a SAS library with SDTM datasets generated in a shift-Jis Japanese encoding environment
- > WLATIN a SAS library with SDTM datasets generated in WLatin1 western encoding environment
- SJISRAW a SAS library with legacy datasets generated in shift-Jis Japanese encoding environment.
 Some datasets contain some variables with Japanese characters

My SAS session is using a WLatin1 western encoding

Dealing with Transcoding in SAS



Checking the dataset encoding

```
%encod(libin=utf8,dsin=ae);
                                                                    Transcoding Successful
----CHECKING FOR ENCODING-----
utf8.ae
NOTE: Data file UTF8.AE.DATA is in a format that is native to another host, or the file encoding does not match the session encoding.
     performance.
ENCODING is utf-8 Unicode (UTF-8)
    %encod(libin=SJIS,dsin=ae);
                                                                    Transcoding Successful
----CHECKING FOR ENCODING------
SJIS.ae
NOTE: Data file SJIS.AE.DATA is in a format that is native to another host, or the file encoding does not match the session encoding.
     performance.
ENCODING is shift-iis Japanese (SJIS)
    %encod(libin=wlatin,dsin=Adveve_1_0);
----CHECKING FOR ENCODING-----
                                                             Same encoding, no need for transcoding
wlatin.Adveve_1_0
ENCODING is wlatin1 Western (Windows)
```

Dealing with Transcoding in SAS



Opening a dataset with **shift-Jis Japanese** (**SJIS**) encoding with Japanese Characters in a **WLatin1** SAS encoding environment

VIEWTABLE: Sjisraw.Adveve_1_0 (Pretreatment/Adverse Events)								
	AESEQ	AEVTJ	AESDT	AEEDT				
1	1	+++	11APR2012	30MAY2012				
2	1	++++	21MAR2012	26MAR2012				
3	1	++	21MAR2012	22MAR2012				
4	1		21FEB2012	06MAR2012				
5	2	++++	26FEB2012	29NOV2012				
6	2	++++	17SEP2012	08OCT2012				
7	1	AST++	21SFP2012	19NOV2012				

CEDA in action

NOTE: Data file SJISRAW.ADVEVE_1_0.DATA is in a format that is native to another host, or the file encoding does not match the session encoding.

Cross Environment Data Access will be used, which might require additional CPU resources and might reduce performance.

WARNING: Some character data was lost during transcoding in the dataset SJISRAW.ADVEVE_1_0. Either the data contains characters that are not representable in the new encoding or truncation occurred during transcoding.

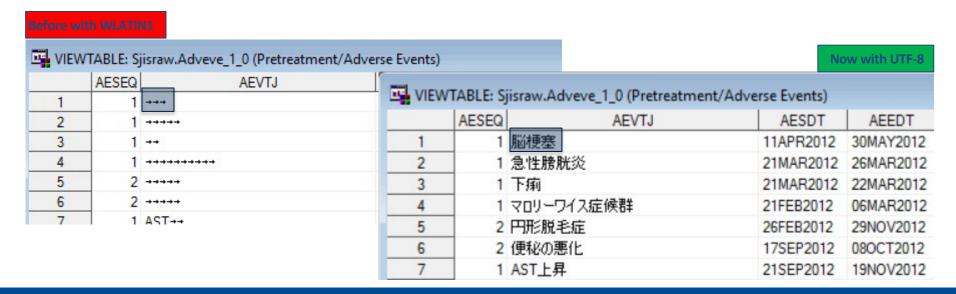
Transcoding Failed

Dealing with Transcoding in SAS



When transcoding fails, the only solution is to have your SAS session using the same Encoding of the original SAS session (or use UTF-8), this is determined at SAS startup (AUTOEXEC)

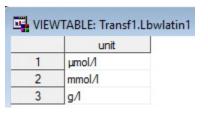
-ENCODING UTF-8

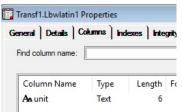


Encoding and SAS wlatin1 → UTF-8



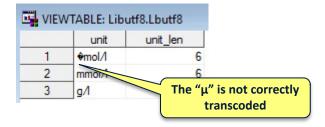
LButf8 a dataset generated with UTF-8 encoding





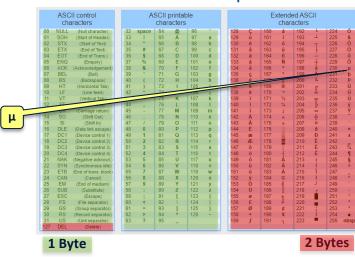
Reading LButf8 in an wlatin1 encoding SAS session

data libutf8.LBUTF8;
 set libwlat.LBwlatin1;
run;





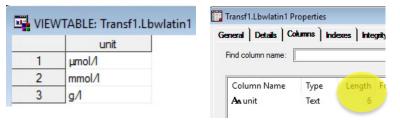
UTF-8 Extended ASCII Representation



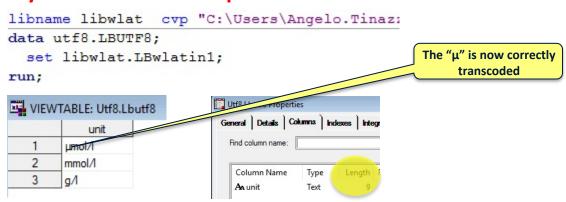
Encoding and SAS wlatin1 → UTF-8



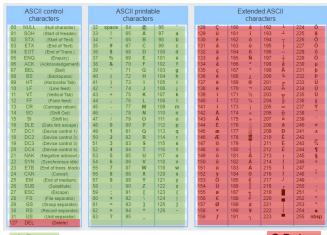
LBwlatin1 a dataset generated with WLATIN1 encoding



Reading LBwlatin1 in an UTF-8 encoding SAS session Try the CVP libname option



UTF-8 Extended ASCII Representation

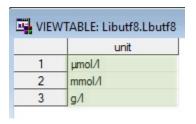


1 Byte

2 Bytes



LButf8 a dataset generated with UTF-8 encoding



Reading LButf8 in an wlatin1 encoding SAS session

VIEW	TABLE: Wor	k.Lbwlatin1
	unit	
1	µmol/l	
2	mmol/l	
3	g/l	

NOTE: Data file LIBUTF8.LBUTF8.DATA is in a format that Cross Environment Data Access will be used.......

The transcoding was successful, no errors

Failing transcode to wlatin1



SJISRaw a dataset generated with SJIS encoding containing Japanese characters

```
data Adveve_1_0;

set SJISRaw.Adveve_1_0;

run;

ERROR: Some character data was lost during transcoding in the dataset SJISRAW.ADVEVE_1_0. Either the NOTE: The DATA step has been abnormally terminated.

NOTE: The SAS System stopped processing this step because of errors.

WARNING: The data set WORK.ADVEVE_1_0 may be incomplete. When this step was stopped there were 0 obs WARNING: Data set WORK.ADVEVE_1_0 was not replaced because this step was stopped.

data Adveve_1_0;

set SJISRaw.Adveve_1_0 (encoding=ANY);

run;

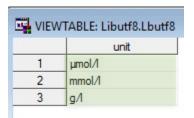
NOTE: There were 212 observations read from the data set SJISRAW.ADVEVE_1_0.

NOTE: The data set WORK.ADVEVE_1_0 has 212 observations and 50 variables.
```

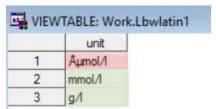
Failing transcode to wlatin1



LButf8 a dataset generated with UTF-8 encoding



Reading LButf8 in an wlatin1 encoding SAS session with encoding=ANY option

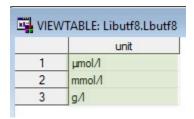


```
data lbwlatin1;
  set libutf8.lbutf8(encoding=any);
run;
```

Failing transcode to wlatin1



LButf8 a dataset generated with UTF-8 encoding



Reading LButf8 in an wlatin1 encoding SAS session with encoding=ANY option

```
VIEWTABLE: Work.Lbwlatin1

unit

1 

µmol/1

2 

mmol/1

3 

g/1
```

```
data lbwlatin1;
  set libutf8.lbutf8(encoding=ANY);
  array allchar (*) _character_;
  do i=1 to dim(allchar);
    txt=KCVT(allchar(i), "UTF-8", "WLATIN1");
    varname=vname(allchar(i));
    if txt^=allchar(i) then do;
        put ".... Instring for variable " varname " : " allchar(i) "converted to: " txt;
        allchar(i)=txt;
    end;
end;
drop txt varname;
run;
```

.... Instring for variable unit : ΑμμοΙ/Ι converted to: μμοΙ/Ι

Transcode datasets to your local encoding



SAS CEDA transcode for you

```
NOTE: Cross Environment Data Access will be used, which might require additional CPU resource and might reduce performance
```

But it requires additional CPU resources, Time!!!!

proc copy inlib=sjis outlib=transf2 noclone;

run;

```
/*Test use of CPU: Using original datasets */
proc sql noprint;
  create table testSQLOriginalEncoding as
    select a.*, b.arm, b.race, b.sex
    from SJIS.AE a left join SJIS.DM b
    on a.usubjid=b.usubjid;
quit;
```

```
NOTE: PROCEDURE SQL used (Total process time):
real time 0.13 seconds
cou time 0.06 seconds
```

```
/*Test use of CPU: Using transcoded datasets */
proc sql noprint;
  create table testSQLTanscodedEncoding as
    select a.*, b.arm, b.race, b.sex
    from TRANSF2.AE a left join TRANSF2.DM b
    on a.usubjid=b.usubjid;
quit;
```

```
NOTE: PROCEDURE SQL used (Total process time):
real time 0.06 seconds
cpu time 0.04 seconds
```



Encoding and Data Submission

FDA, PMDA and NMPA Requirements

The SAS XPT Challenge with Multiple Byte Character Set (MBCS)

Encoding and Data Submission FDA, PMDA and NMPA Requirements



FDA

From the FDA Study Data Technical Conformance Guide section 3.3.5

Variable names, as well as variable and dataset labels should include American Standard Code for Information Interchange (ASCII) text codes only. Variable values are the most broadly compatible with software and operating systems when they are restricted to ASCII text codes (printable values below 128). Use UTF-8 for extending character sets; however, the use of extended mappings is not recommended. Transcoding errors, variable length errors, and lack of software support for multi byte UTF-8 encodings can result in incorrect character display and variable value truncations. Ensure that LBSTRESC and controlled terminology extensions in LBTEST do not contain byte values 160-191 as some character mappings in that range may interfere with agency processes.

FDA, PMDA and NMPA Requirements



FDA

restricted to ASCII text codes (printable values below 128)

4	C 3	LIIC	eu lo Asc	Ш	LE	AL CC	ues	pi	1111	ubie
ASCII control					ASCII printable					
characters				characters						
Н	00	NULL	(Null character)		32	space	64	@	96	,
	01	SOH	(Start of Header)		33	!	65	A	97	а
н	02	STX	(Start of Text)		34	"	66	В	98	b
п	03	ETX	(End of Text)		35	#	67	С	99	С
н	04	EOT	(End of Trans.)		36	\$	68	D	100	d
п	05	ENQ	(Enquiry)		37	%	69	E	101	е
ш	06	ACK	(Acknowledgement)		38	&	70	F	102	f
П	07	BEL	(Bell)		39		71	G	103	g
ш	80	BS	(Backspace)		40	(72	H	104	h
Ш	09	HT	(Horizontal Tab)		41)	73	- 1	105	i
ш	10	LF	(Line feed)		42	*	74	J	106	j
Ш	11	VT	(Vertical Tab)		43	+	75	K	107	k
ш	12	FF	(Form feed)		44	,	76	L	108	- 1
Ш	13	CR	(Carriage return)		45	-	77	M	109	m
Ш	14	SO	(Shift Out)		46		78	N	110	n
П	15	SI	(Shift In)		47	1	79	0	111	0
ш	16	DLE	(Data link escape)		48	0	80	Р	112	р
Ш	17	DC1	(Device control 1)		49	1	81	Q	113	q
Ш	18	DC2	(Device control 2)		50	2	82	R	114	r
ш	19	DC3	(Device control 3)		51	3	83	S	115	S
	20	DC4	(Device control 4)		52	4	84	Т	116	t
	21	NAK	(Negative acknowl.)		53	5	85	U	117	u
	22	SYN	(Synchronous idle)		54	6	86	V	118	٧
	23	ETB	(End of trans. block)		55	7	87	W	119	W
	24	CAN	(Cancel)		56	8	88	X	120	X
	25	EM	(End of medium)		57	9	89	Y	121	У
	26	SUB	(Substitute)		58	- :	90	Z	122	Z
	27	ESC	(Escape)		59	;	91	[123	{
	28	FS	(File separator)		60	<	92	1	124	
	29	GS	(Group separator)		61	=	93]	125	}
	30	RS	(Record separator)		62	>	94	٨	126	~
	31	US	(Unit separator)		63	?	95	_		
L	127	DEL	(Delete)							

Removing "control characters"

```
mychar=compress(mychar,,'kw');
```

Other "control characters"

Encoding and Data Submission FDA, PMDA and NMPA Requirements



FDA

LBTEST do not contain byte values 160-191

		uo i	101	COITE	1111	Dyte	Vario		
		E	vton	And AS	CII				
Extended ASCII									
characters									
128	Ç	160	á	192	L	224	Ó		
129	ü	161	í	193	\perp	225	ß		
130	é	162	Ó	194	Т	226	Ô		
131	â	163	ú	195	Ŧ	227	Ò		
132	ä	164	ñ	196	-	228	õ		
133	à	165	Ñ	197	+ ã	229	Õ		
134	å	166	3	198	ã	230	μ		
135	ç	167	۰	199	Ã	231	þ		
136	ê	168	ż	200	L	232	Þ		
137	ë	169	8	201	1	233	Ú		
138	è	170	7	202	ᆚ	234	Û		
139	Ï	171	1/2	203	ī	235	Ù		
140	î	172	1/4	204	Ţ	236	Ý		
141	ì	173	i i	205	=	237	Ý		
142	A	174	« «	206	#	238	_		
143	Å	175	>>	207		239	•		
144	É	176		208	ð	240	=		
145	æ	177		209	Đ	241	± .		
146	Æ	178		210	Ê	242	=		
147	ô	179		211	Ë	243	3/4		
148	Ö	180		212	È	244	¶		
149	Ò	181	A	213	ļ	245	§		
150	û	182	Â	214	ĺ	246	÷		
151	ù	183	À	215	Î	247	0		
152	ÿ	184	0	216	Ï	248	٥		
153	Ö	185	=	217		249			
154	Ü	186		218		250			
155	Ø	187]	219		251	1		
156	£	188		220		252	3		
157	Ø	189	¢	221	Ī	253	2		
158	×	190	¥	222		254			
159	f	191	٦	223	•	255	nbsp		

Use of Extended ASCII characters require careful investigation → no automatic replacement



Non-Printable and Special Characters? ... BYTE me!

L. Sims, PHUSE 2016

Encoding and Data Submission FDA, PMDA and NMPA Requirements



PMDA (Japan) and NMPA (China)

- Some variables need to be provided in both English and Japanese/Chinese i.e. AE verbatim
- NMPA have additional requirements
 - Define-xml and Reviewer Guide in Chinese
 - > Encoding used to be specified in the reviewer guide
- SAS Datasets and Variables Labels in Chinese
- SAS XPT is still the data transfer format (not necessarily XPT 5 for NMPA)

Encoding and Data Submission

The SAS XPT Challenge with Multiple Byte Character Set (MBCS)



- Some variables need to be provided in both English and Japanese/Chinese i.e. AE verbatim
- NMPA have additional requirements
 - Define-xml and Reviewer Guide in Chinese
 - SAS Datasets Labels in Chinese
 - > SAS XPT is SBCS

Some SDTM Labels won't fit!

CMTRT: Reported Name of Drug, Med, or Therapy

Chinese: 报告药品报告药品名称,药物治疗

requires 42 bytes, available are only 40 - last character

will be lost

HODECOD: Dictionary-Derived Term for the Healthcare

Encounter

Chinese: 词典中针对医疗保健遇到的术语

requires 42 bytes, available are only 40 - last character

will be lost

http://xml4pharma.com/publications/Poster Jozef Aerts Chinese characters XPT.pdf



Conclusions

Summary

Other aspects to consider → References

Conclusions Summary



Make sure the encoding of the data you receive is compatible with your encoding and apply the needed transformations



- Length of variables when encoding is not SBCS
- Working with MBCS Encoding (Use of "K" functions)



- CDISC-CT, US-FDA requirements
- Other Agencies requirements



PHUSE US CONNECT 2021

Paper DH12

The encoding dilemma

Angelo Tinazzi, Cytel Inc., Geneva, Switzerland

ABSTRACT

It is not uncommon to see in SAS log files messages such as "Some character date was not during francocting in the dataset." Must are we risking if we don't take care of such a message PWhat is the message leding use for care up revent the message loc occur or what actions do we need to take to eventually correctly reading the data we Moreover, regulatory agencies such as the Japanese (PMDA) or more recently the Chinage (NMPA) (I.j. are now

asking to submit datasets, and not surprisingly for surprisingly for someone), they are requesting not only datasets in CDISC format but they are also requesting or suggesting the use of SAS XPT format as a fise data format, furthermore, the NMPA has specific requirements with regards to the use of Chinese language for things tillo label or datasets and valenties and contact of character variables or g, adverse event terms.

with an paper I would need to stress the importance of rendoming in Social flow in Could affect the Context handling of text. Options will be provided to make sure your character data are correctly retrieved when 'special' characters are handled. The risk of data lose when using SAS XPT and possible solutions will be also discussed.

Conclusions

Other aspects to consider \rightarrow References



- ➤ Tips and Fixes for Cross-Environment Batch Transfer of SAS® Data Y. Zhuo; PharmaSUG 2018
- ➤ Data Encoding: All Characters for All Countries D. Dutton; PHUSE 2015
- > The impact of Change from wlatin1 to UTF-8 in SAS Environment H. Song, A. Koster; PharmaSUG 2016
- UTF What? A Guide for Handling SAS Transcoding Errors with UTF-8 Encoded Data M. Stackhouse, L. Pogula, PharmaSUG 2018
- > SAS 9.3 UTF-8 Encoding Support and Related Issue Troubleshooting J. Liang, Edmonton User Group 2016
- Non-Printable and Special Characters? ... BYTE me! L. Sims, PHUSE 2016
- SAS® 9.4 National Language Support (NLS): Reference Guide, Fifth Edition, SAS https://documentation.sas.com/api/collections/pgmsascdc/9.4-3.5/docsets/nlsref/content/nlsref.pdf?locale=en#nameddest=n04275wmuchnqjn1hfx0yzpopw5a
 - Chapter 3 "Encoding for NLS"
 - Chapter 4 "Transcoding for NLS"
 - Chapter 5 "Double-Byte Characters Sets (DBCS)"
- ➤ Guideline on the Submission of Clinical Trial, NMPA Center for Drug Evaluation Data (https://www.nmpa.gov.cn/directory/web/nmpa/images/obbSqc7vwdm0ssrU0enKb7dtd29u9a4tbzUrdTyo6jK1NDQo6mhty5wZGY=.pdf)
- Chinese and Asian characters in SAS Transport 5 datasets: Why that is the worst possible choice, J. Aerts, 2020 (unpublished CDISC-US Interchange poster) http://xml4pharma.com/publications/Poster Jozef Aerts Chinese characters XPT.pdf





Angelo Tinazzi, Senior Director
Cytel Inc.
Standards, Systems, CDISC Consulting, Statistical Programming
Clinical Research Services
Route de Pré-Bois 20
C.P 1839, 1215 Geneva, SWITZERLAND

Email: angelo.tinazzi@cytel.com



Backup Slides

Dealing with Transcoding in SAS



Getting more details with PROC CONTENTS (dataset with UTF-8 encoding)

```
proc contents data=utf8.ae;
run;
```

The CONTENTS Procedure

Data Set Name
Member Type
Engine
Created
Last Modified
Protection
Data Set Type
Label
Data Representation
Encoding

UTF8.AE DATA U9 05/22/2018 23:39:22 05/22/2018 23:39:22

Adverse Events

HP_UX_64, RS_6000_AIX_64, SOLARIS_64, HP_IA64 utf-8 Unicode (UTF-8)

Different encoding Different Host/Os

Observations	200
Jariables	31
Indexes	0
Observation Length	648
Deleted Observations	0
Compressed	NO
Sorted	NO

run;

Dealing with Transcoding in SAS



Getting more details with PROC CONTENTS (dataset with SJIS encoding) proc contents data=SJIS.ae;

Encoding Details for SJIS.ae . The CONTENTS Procedure

Data Set Name
Member Type
Engine
Created
Last Modified
Protection
Data Set Type
Label
Data Representation
Encoding

SJIS.AE DATA V9 10/24/2019 04:10:06 10/24/2019 04:10:06

Adverse Events WINDOWS_64 shift-jis Japanese (SJIS)

Different encoding Same Host/Os

Observations	1137
Jariables	33
Indexes	0
Observation Length	672
Deleted Observations	0
Compressed	NO
Sorted	NO

Dealing with Transcoding in SAS



Getting more details with PROC CONTENTS (dataset with wlatin1 encoding)

```
proc contents data=wlatin.Adveve_1_0;
run;
```

The CONTENTS Procedure

Data Set Name
Member Type
Engine
Created
Last Modified
Protection
Data Set Type
Label
Data Representation
Encoding

```
WLATIN.ADVEVE_1_0
DATA
V9
05/20/2020 06:42:11
05/20/2020 06:42:11
```

```
WINDOWS_64
wlatin1 Western (Windows)
```

Same encoding Same Host/Os

```
Observations 314
Variables 50
Indexes 0
Observation Length 4568
Deleted Observations 0
Compressed NO
Sorted NO
```

Encoding (and Transcoding) in a Nutshell

Encoding Methods SBCS vs DBCS vs MBCS



A set of rules that assign the numeric representations to the set of characters

- Size of the encoding (type and number of characters represented)
- Number of bytes (bits) used
 - Single-Byte Character Set (SBCS)
 - Double-Byte Character Set (DBCS)
 - Multiple-Byte Character Set (MBCS)

Encoding	Character Set	Size (Nr. of Bytes)
WLATIN1 (SBC)	ASCII (and extended ASCII)	1
SHIFT-JIS	ASCII, Katakana, other Japanese	1 or 2 bytes
UTF-8	ASCII, foreign languages, special symbols and more	1 to 4 bytes

Summary of "Options"



Checking your SAS session encoding

```
proc options option=encoding;
run;
```

Checking dataset encoding

```
%let dsid=%sysfunc(open(<libname>.<dataset name>,i));
%let encodDS=%sysfunc(%attrc(&dsid,encoding));
```

Changing your SAS session encoding (AUTOEXEC)

-encoding UTF-8

If you need to switch to a different session

Avoiding data truncation from wlatin1 to UTF-8

```
libname libor cvp "<original datasets folder>";
proc copy inlib=libor outlib=libnew noclone;
run;
```

Avoiding errors when reading UTF-8 datasets in wlatin1 session

```
set <UTF-8 dataset>(encoding=ANY);
libname libor "<original datasets folder>"
    inencoding=ANY;
```

Assessing and trying to correct encoding=ANY