

# Construção da Coorte de 100 milhões de brasileiros

Adriana Viriato Ribeiro, David Gorender, Luan Menezes





Para a construção da Coorte de 100 milhões de brasileiros foi necessária a utilização dos *backups* do sistema do Cadastro Único (versões 6, 7 e 7.1) cedidos pelo Ministério de Desenvolvimento Social. Esses *backups* contém cópias das bases em diferentes momentos entre os anos de 2006 a 2015. Para cada ano de backup existem tabelas associadas às informações de família e do indivíduo. Na versão 6 existem três tabelas, as tabelas A, B e C. Enquanto na versão 7 existem as tabelas de 1 a 17. Nem todas as tabelas foram utilizadas na construção da Coorte. Na versão 6 foram utilizadas as tabelas A e B e na versão 7 foram usadas as tabelas 1 a 9, 10, 11, 13 e 15.

A lista com todas as bases utilizadas pode ser observada na Tabela x.

Ano referente	Data da extração	Nome diretório	Versão
2006	Março de 2007	MARCO_2007	6
2007	Dezembro de 2007	DEZEMBRO_2 007	
2008	Dezembro de 2008	DEZEMBRO_2 008	
2009	Dezembro de 2009	DEZEMBRO_2 009	
2010	Setembro de 2010	SETEMBRO_2 010	
2011	Janeiro de 2012		7
2012	Janeiro de 2013		
2013	Dezembro de 2013		
2014	Dezembro de 2014		





2015	Novembro de 2015		
------	------------------	--	--

Essas bases são consideradas as bases "originais" e representam as cópias cedidas sem nenhum tipo de modificação ou tratamento. No entanto, para a construção da Coorte e para preparar o banco final para análise por parte dos pesquisadores, é necessário realizar algumas tarefas em relação à limpeza e tratamento dos dados. Desta forma, algumas etapas foram definidas para a criação da Coorte baseada no Cadastro Único, a saber:

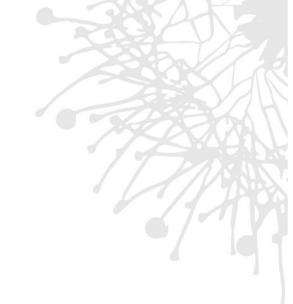
- 1. Conversão para csv
- 2. Criação de um dicionário do CIDACS
- 3. Padronização das variáveis
- 4. Harmonização das variáveis
- 5. Split por data de atualização
- 6. Bind por ano

Cada uma dessas etapas será explicada nas próximas seções, considerando algumas subtarefas, as diferentes características entre as versões e as avaliações de integridade e qualidade dos dados.

#### 1 Conversão para csv

A conversão para csv é uma etapa realizada para tratamento inicial dos dados. A ideia é que os dados sejam convertidos para csv para que as





variáveis categóricas, inteiros, datas, entre outras, sejam tratadas antes de ser aplicado o *schema* em parquet. Desta forma, quando a cópia da base cedida ao Cidacs está em formato csv, é feita uma verificação de inconsistências e ela é salva novamente no mesmo formato adequando-se aos padrões de codificação e outras premissas. Quando a cópia da base chega em outros formatos, é realizada a conversão de acordo com as mesmas verificações de inconsistências. Em relação ao Cadastro Único, as bases relacionadas à versão 6 tem formato *fixed width*, enquanto as bases derivadas das versões 7 e 7.1 tem formato csv.

#### 1.1 Conversão da Versão 6 para csv

A versão 6 do Cadastro Único é composta por duas tabelas: tabela A e tabela B. A tabela A refere-se às informações de família, enquanto a tabela B mantém informações dos indivíduos. Ambas tabelas são utilizadas na construção da Coorte, portanto, é preciso fazer a conversão de todas as tabelas em todos os anos (de 2006 a 2010). Desta forma, fez-se a conversão de *fixed width* para csv dos seguintes arquivos:

```
tbas = ['MARCO_2007/CNT.ICO.RJ.BHA1.ICOD884A.D070402', 
 'DEZEMBRO_2007/CNT.ICO.RJ.BHA1.ICOD884A.D080102', 
 'DEZEMBRO_2008/CNT.ICO.RJ.BHA1.ICOD884A.D090102', 
 'DEZEMBRO_2009/CNT.ICO.RJ.BHA1.ICOD884A.D100101', 
 'SETEMBRO_2010/CNT.ICO.RJ.BHA1.ICOD884A.D100901']
```

```
tbbs = ['MARCO_2007/CNT.ICO.RJ.BHA1.ICOD884B.D070402', 
'DEZEMBRO_2007/CNT.ICO.RJ.BHA1.ICOD884B.D080102', 
'DEZEMBRO_2008/CNT.ICO.RJ.BHA1.ICOD884B.D090102', 
'DEZEMBRO_2009/CNT.ICO.RJ.BHA1.ICOD884B.D100101', 
'SETEMBRO_2010/CNT.ICO.RJ.BHA1.ICOD884B.D100901']
```





Antes de realizar a conversão, é preciso substituir os caracteres que podem interferir na quantidade de colunas em um arquivo csv. A saber: contra barra, vírgula e aspas duplas. Para isso, antes da conversão é realizada uma alteração na cópia do arquivo original (nomedoarquivooriginal\_sed). Essa alteração é feita no próprio terminal do linux utilizando o seguinte comando:

#### sed -i 'y/,\"\\/.../' nomedoarquivo

Depois dessa etapa, inicia-se o processo de conversão para csv. Esse processo consiste na utilização das chaves cedidas pelo MDS para mapeamento das variáveis. Existem duas chaves, uma para a tabela A e uma para tabela B, e as chaves consistem do nome da variável, sua posição inicial e sua extensão. Por exemplo:

Desta forma, CD\_DOMICILIAR é o nome da variável iniciada na posição 1, cuja extensão é de 9 caracteres. Assim, basta aplicar as chaves das tabelas A e B para os arquivos anuais que representam cada tabela. Como a intenção é realizar a conversão para csv, o procedimento feito é o mapeamento de cada coluna, que é armazenada em um arquivo de texto que utiliza vírgula como separador (csv), todo esse procedimento pode ser observado no código 01 v6-to-csv.

Depois de salvar todos os arquivos de todas as tabelas, é preciso adicionar o cabeçalho com o nome das variáveis. Para isso, basta executar os seguintes procedimentos:





1 - Criar arquivos com os cabeçalhos das tabelas A e B:

vim header tbla lower.csv

#Digite:

ds\_logradouro,cd\_destino\_lixo,cd\_tipo\_domicilio,cd\_construcao, cd\_eas\_ms,cd\_familiar,nm\_bairro\_logradouro,cd\_ddd\_logradouro,cd\_didentificacao\_domicilio,cd\_escoamento\_sanitario,nm\_entrevist ador,nm\_estabelecimento\_saude,dt\_cadastro\_domicilio,cd\_ilumina cao,cd\_situacao\_domicilio,cd\_identificacao\_caixa,cd\_abastecime nto\_agua,cd\_tipo\_cobertura,dt\_pesquisa\_cadastro,qt\_pessoas\_inf ormada,cd\_tipo\_localidade,nu\_telefone\_logradouro,qt\_mulheres\_g ravidas,nu\_comodos,in\_domicilio\_excluido,in\_domicilio\_ativo,nu\_nis\_entrevistador,cd\_ibge\_logradouro,qt\_maes\_amamentando,nu\_c npj\_pref\_orgao,dt\_alteracao\_domicilio,nm\_logradouro,in\_complem entado\_bes,cd\_tratamento\_agua,nu\_logradouro,dt\_inclusao\_domici lio,qt\_deficientes,cd\_cep\_logradouro,sg\_uf\_logradouro,cd\_orige m\_cadastro,cd\_domiciliar,ds\_complemento\_logradouro

vim header\_tblb\_lower.csv

#Digite:

cd\_folha\_doc\_certidao,nu\_ordem\_rl,dt\_chegada\_brasil,dt\_nascime
nto,cd\_ocupacao,nm\_outro\_benef\_pessoa,tp\_documento\_certidao\_ci
vil,vl\_despesa\_prest\_habitacional,vl\_remuneracao\_emprego,vl\_re
nda\_pensao\_alimenticia,cd\_parentesco,cd\_censo\_inep,cd\_benefici





o\_peti,in\_participa\_proger,in\_deficiencia\_mental,rf\_anos\_morad ia,cd\_nacionalidade,ds\_orgao\_emissor\_identidade,dt\_ultima\_alte racao\_domic,vl\_despesa\_luz,in\_outro\_programa,in\_nenhum\_program a,sg\_uf\_emissao\_ctps,dt\_ultima\_alteracao\_pessoa,nu\_ordem pesso a,dt\_inclusao\_peti,in\_participa\_pronaf,dt\_emissao\_doc\_identida de, nu natural pessoa, cd grau instrucao, vl renda seguro desempr ego,nu\_termo\_doc\_certidao,vl\_beneficio\_peti,in\_outra\_deficienc ia,in\_complementado\_bes,vl\_despesa\_gas,nm\_pai,in\_beneficiario\_ bal,rf\_mes\_gravidez,sg\_uf\_doc\_certidao,in\_pessoa\_duplicado,in\_ deficiencia\_fisica,nu\_nis\_pessoa,dt\_emissao\_ctps,in\_pessoa\_exc luida,vl\_despesa\_transporte,cd\_pais\_origem,in\_deficiencia\_cegu eira,dt\_inclusao\_agente\_jovem,sg\_uf\_emissao\_identidade,dt admi ssao\_empresa,nm\_empresa,cd\_sexo,vl\_despesa\_aluguel,cd\_ocupacao \_exercida\_peti,in\_participa\_agente\_jovem,in\_ocupacao\_peti,nu\_o rdem\_esposo,cd\_escola,nu\_cnpj\_empresa,vl\_renda\_aposentadoria,c d\_livro\_doc\_certidao,nu\_ctps,in\_participa\_juv\_cidada,nu\_ordem\_ mae,in\_domicilio\_excluido,in\_participa\_prev\_rural,nu\_documento \_identidade,nm\_pessoa,in\_participa\_loas\_bpc,dt\_emissao\_doc\_cer tidao,in\_participa\_peti,nu\_ordem\_pai,cd\_ibge\_nascimento,vl\_out ras\_despesas,vl\_despesa\_agua,cd\_estado\_civil,nu\_serie\_ctps,in\_ beneficiario\_bes,nu\_cpf,vl\_despesa\_medicamentos,cd\_mercado\_tra balho, ds complemento doc identidade, cd familiar, in amamentando ,in\_participa\_bolsa\_alimentacao,cd\_raca\_cor,vl\_despesa\_aliment acao, nu\_secao\_eleitor, cd\_crianca\_quem\_fica, in\_deficiencia\_surd ez,nu\_nis\_original,in\_deficiencia\_mudez,vl\_outras\_rendas,nu\_zo na eleitor,nm mae,nm escola,in domicilio ativo,cd tipo ocupaca o\_peti,dt\_inclusao\_proger,in\_nenhuma\_deficiencia,nu\_pessoa\_viv em renda, in liberto trab escravo, in participa bolsa escola, nu titulo eleitor, rf mes moradia, cd serie escolar





2 - Adicionar o cabeçalho aos arquivos das tabelas (exemplo para tabela A de 2011):

```
cat header_tbla_lower.csv > tbla2011_csv/tba.csv
cat tbla2011_csv/part* >> tba.csv
```

Depois de realizado o procedimento de conversão para csv, foi feita uma descritiva para avaliar a integridade e consistência dos dados de ambas as tabelas. Nesse processo, foi identificada a existência de 12 registros inconsistentes nas tabelas B dos diversos anos:

2006	2
2007	2
2008	4
2009	2
2010	2

A inconsistência refere-se ao número de colunas desses registros. Enquanto os registros normais da tabela B tem 107 variáveis, esses 12 registros tinham apenas 42. Como 42 é a quantidade de atributos encontrados na tabela A, tentou-se aplicar a chave da tabela A nesse registros para recuperá-los, no entanto, ainda assim os registros permaneceram inconsistentes. Desta forma, esses registros foram removidos da base.





#### 1.1 Conversão da Versão 7 para csv

As bases da versão 7 estão no formato *csv.* Algumas operações foram realizadas nas tabelas para evitar erros no tratamento das bases. No cabeçário de cada uma foi necessário realizar a substituição do caractere "\_" pelo "." e todas as letras foram convertidas em minúsculas. Foi necessário padronizar a codificação de caracteres de cada tabela, convertendo de UTF-8 para ASCII.

As bases originais possuem inconsistências com a quantidade de colunas, e quantidade total de registros, o que traz problemas para os procedimentos futuros. Foi identificado que o caracter "\" levava a quebra de linhas resultando nessas inconsistências citadas. Então foi necessário a remoção do mesmo para que a bases da versão 7 pudessem ser devidamente utilizadas.





## 2 Merge entre as tabelas

As bases do Cadastro Único são compostas por tabelas anuais que podem representar informações de família, domicílio e do indivíduo, por exemplo. Para a construção dos bancos anuais, é necessário realizar um join entre as tabelas de cada ano.

## 2.1 Merge da Versão 6

Como foi citado anteriormente, as tabelas da versão 6 que foram utilizadas na construção da Coorte foram as tabelas A e B. A tabela A corresponde às informações de família, enquanto a tabela B corresponde às informações individuais. Assim, foi realizado um merge<sup>1</sup> entre as tabelas A e B a partir da variável CD\_FAMILIAR, resultando em bancos compostos por todos os indivíduos de cada ano, na qual cada linha representa as informações individuais e familiares.

O merge realizado foi do tipo full

Na versão 7 do Cadastro Único, foram utilizadas as tabelas de 1 a 15 (com exceção das tabelas 9,10 e 14) para construção dos bancos anuais. As tabelas 1 e 3 contém informações da família, a tabela 2 mantém informações do domicílio e as tabelas de 4 a 8 são compostas por informações do indivíduo. As tabelas de 11 a 13 e 15 referem-se a informações suplementares de pendência, referente à família e indivíduo. Assim, o merge entre as tabelas de 1 a 4, 11 e 13f foi realizado apenas com a variável cod\_familiar\_fam, enquanto o merge com as outras tabelas foi realizado com as variáveis cod\_familiar\_fam e num\_membro\_fmla.

TABELA 1: REPRESENTAÇÃO DA TRANSFORMAÇÃO ENTRE TABELAS PARA BANCOS ANUAIS NA VERSÃO 6

<sup>&</sup>lt;sup>1</sup> Merge ou join é uma função que associa dois bancos de acordo com uma ou mais variáveis em comum.





Tabela A_2006	0000
Tabela B_2006	2006

Tabela 2: Representação da transformação entre tabelas para bancos anuais na versão 7

Tabela 1_2015	
Tabela 2_2015	
Tabela 3_2015	
Tabela 4_2015	
Tabela 5_2015	
Tabela 6_2015	0045
Tabela 7_2015	2015
Tabela 8_2015	
Tabela 11_2015	
Tabela 12_2015	
Tabela 13_2015	
Tabela 15_2015	





#### 2.1 Inconsistência no merge da versão 7

Após a realização do merge entre as tabelas da versão 7 foi identificado mudanças na quantidade de registros, das variáveis categóricas quando comparados, tabela original e base pós merge. Aumento esse na quantidade de registros nulos e também em algumas categorias.

Foi realizado um merge full entre as tabelas, o que significa que são combinados registros iguais entre elas e também os divergentes. A combinação dos registros divergentes ( $cod_familiar_fam e num_membro_fmla$ ) explica o acréscimo de registros nulos. As tabelas de 1 a 13 (com exceção das tabelas 9 e 10) não possuem registros duplicados, porém foi identificado registros duplicados na tabela 15 explicando o acréscimo de registros nas categorias. A tabela 15 possui dois tipos duplicação, onde foi realizado deduplicação para os registros que possuem valores iguais nas variavéis ( $cod_familiar_fam,num_membro_fmla,cod_prog_prohab_memb$ ) porém quando há duplicação apenas nas variáveis ( $cod_familiar_fam,num_membro_fmla$ ) foi realizado a criação de uma nova coluna( $cod_prog_prohab_memb_comp$ ) para armazenar os duplicados.





# 3 Compatibilização da versão 7

As versões 7.1 e 7.2 do Cadastro Único são muito similares, com exceção de algumas variáveis que foram adicionadas/removidas na versão 7.2. Todas as variáveis que foram adicionadas/removidas na versão 7.2 foram excluídas dos registros em ambas as versões. Enquanto as variáveis com nomenclatura diferente tiveram seus nomes harmonizados. Desta forma, foi preciso mapear e tratar essas variações.

As diferenças mapeadas foram:

#### Tabela 1:

Removida	vazio
	num_reg_arquivo
Adicionada	qtde_pessoas
Nomenclatura diferente	dat_atual_fam (v 7.1) e dat_alteracao_fam (v 7.2)

Tabela 2:



Tabela 7:



Removida	vazio	
	num_reg_arquivo	
Tabela 3:		
Removida	vazio	
	num_reg_arquivo	
	vazio2	
Nomenclatura diferente	nu_estbo_saude (v 7.1) e cod_eas_fam (v 7.2)	
Tabela 4:		
Removida	vazio	
	num_reg_arquivo	
	ind_transferencia_pessoa	
Adicionada	matchcode	
	matchcode_mae	
Tabela 5:		
Removida	vazio	
	num_reg_arquivo	
abela 6:		
Removida	vazio	
	num_reg_arquivo	





Removida	vazio
	num_reg_arquivo

#### Tabela 8:

Removida	vazio
	num_reg_arquivo

# 4 Padronização

## 4.1 Criação dos dicionários das versões 6 e 7

A criação dos dicionários é necessária para exposição das variáveis para os pesquisadores e para a automatização do processo de padronização. O dicionário é caracterizado pelos seguintes campos:

- 1) Nome das variáveis.
- 2) Descrição original, incluindo as categorias e formatos.
- 3) Descrição das categorias e formatos padronizados.
- 4) Uma coluna utilizada para a marcação das categorias que será utilizada no script para geração do código de padronização (esse campo não precisa estar presente na versão pública do dicionário).
- 5) Tipo da variável, a saber: integer, long, byte e string.
- 6) Uma coluna de transformação para os casos em que as variáveis mudam de categoria.





- 7) Um campo de comentários para identificar possíveis problemas, características e verificações.
- 8) Campos com a quantidade de missings de cada variável por ano. Esses campos devem ser preenchidos.

Desta forma, foram criados os dicionários para as versões 6 e 7 do Cadastro Único. Esses arquivos podem ser vistos em.

### 4.2 Padronização da versão 6 e 7

Uma vez definido o dicionário, é inicializado o processo de padronização das variáveis. Como foram criados campos nos dicionários para automatizar o processo de padronização, esses campos são lidos e utilizados como insumo para o script de geração automática das funções de padronização.

As funções de padronização dependem dos tipos de variáveis utilizadas:

- 1) Date: são verificadas inconsistências do ponto de vista sintático, a depender de algumas regras: tamanho do campo, verificação se as partes do campo correspondem a dia, mês e ano válidos, entre outras verificações.
- 2) Integer e Long: garantia de que existem apenas números nesses campos.
- Byte: verificação e transformação das categorias de acordo com o que foi definido no dicionário.

Desta forma, verificou-se o tipo da variável de acordo com o campo "tipo" do dicionário e para cada tipo foram aplicadas as transformações/verificações de acordo com as funções de padronização. Depois disso, foram realizadas as chamadas das funções.

Além disso, na versão 6, as variáveis que indicam valor (renda e despesa), além de serem transformadas em inteiro, tiveram seus valores divididos por 100. Visto que na versão 6 o valor cem, por exemplo, está no formato 10000, com duas casas decimais. Esse processo foi realizado para que as variáveis de valor da versão 6 fossem compatíveis com a versão 7.

É importante ressaltar, que nenhuma variável original foi substituída, para cada uma dessas alterações foi criada uma variável com a tag "\_original", que foi mantida sem alterações. Por





exemplo, foi criada uma variável sexo\_original que não sofre alterações, enquanto a variável sexo foi padronizada.