

# CONSTRUCCIÓN BODEGA DE DATOS

## SABERPRO

### GUIA 3

## Limpieza de Datos

Crear la tabla *saberpro\_limpio* a partir de la tabla *saberpro\_2012\_2014*

Realizar los siguientes procesos de limpieza de datos:

1. Reemplazar los valores 'EK20123', EK20133', 'EK20142', 'EK20143 y 'EKO2014' del atributo *estu\_cod\_aplicacion* por 'EK20122', 'EK20132', 'EK20141', EK20142' y EK20141 respectivamente:
2. Reemplazar los valores '20123', 20133', '20142', '20143 y 'O2014' del atributo *prueba* por '20122', '20132', '20141', 20142' y 20141 respectivamente:
3. Reemplazar los valores nulos del atributo *estu\_genero* por la moda de este atributo
4. Reemplazar los valores nulos del atributo *estu\_nacimiento\_dia* por la moda de este atributo en la tabla *saberpro\_2012\_2014*.
5. Reemplazar el valor '5' del atributo *estu\_nacimiento\_dia* por '05'
6. Reemplazar los valores nulos del atributo *estu\_nacimiento\_mes* por la moda de este atributo en la tabla *saberpro\_2012\_2014*.
7. Reemplazar los valores nulos del atributo *estu\_nacimiento\_anno* por la moda de este atributo en los respectivos años de las pruebas 2012, 2013 y 2014.
8. Reemplazar los valores nulos del atributo *estu\_pais\_reside* por 'COLOMBIA'
9. Reemplazar los valores 'CO' del atributo *estu\_pais\_reside* por 'COLOMBIA'
10. Reemplazar los valores nulos del atributo *estu\_estado\_civil* por la moda de este atributo.
11. Reemplazar los valores nulos del atributo *estu\_reside\_codmpio* por los valores del atributo *estu\_exam\_codmpio\_presentacion*.
12. Estandarizar el valor del atributo *estu\_exam\_mpio\_presentacion* de "BOGOTÁ D.C." por el valor "BOGOTA D.C".
13. Estandarizar el valor del atributo *estu\_exam\_dpto\_presentacion* de "BOGOTÁ" por el valor BOGOTA, "BOYACÁ" por BOYACA, CAQUETÁ por CAQUETA,

"ATLÁNTICO" por ATLANTICO, "CHOCÓ" por CHOCO, "CÓRDOBA" por CORDOBA," SAN ANDRÉS" por SAN ANDRES, "VAUPÉS" por VAUPES, "QUINDÍO" por QUINDIO, "GUAINÍA" por GUANIA," BOLÍVAR" por BOLIVAR,

14. Reemplazar los valores nulos del atributo *estu\_exam\_cod* por 141 para la prueba 20141 y por 142 para la prueba 20142.
15. Reemplazar los valores nulos del atributo *estu\_exam\_nombre* por EXAMEN SABER PRO 2014-1 para la prueba 20141 y por EXAMEN SABER PRO 2014-2 para la prueba 20142.
16. Estandarizar los nombres de las instituciones del atributo *inst\_nombre\_institucion*, el origen *inst\_origen* el caracter academico *inst\_caracter\_academico* a partir de los nombres, origen y caracter academico de las instituciones de la tabla IES teniendo en cuenta los códigos de las instituciones
17. Reemplazar los valores nulos del atributo *estu\_nivel\_prgm\_academico* por el valor no nulo del atributo *estu\_nivel\_prgm\_academico* de los programas iguales de la misma institución
18. Reemplazar los valores nulos del atributo *estu\_nivel\_prgm\_academico* por "UNIVERSITARIA" para aquellas instituciones cuyo nombre *inst\_nombre\_institucion* inicie con "UNIVERSIDAD"
19. Reemplazar los valores nulos del atributo *estu\_nivel\_prgm\_academico* por "UNIVERSITARIA" para aquellos programas cuyo nombre *estu\_prgm\_academico* contengan las palabras INGENIERIA, DERECHO, MEDICINA, ODONTOLOGIA, CONTADURIA, LICENCIATURA, PSICOLOGIA, ARQUITECTURA, ECOLOGIA, FISIOTERAPIA, ZOOTECNIA, ADMINISTRACIÓN, DISEÑO, SOCIAL, PUBLICIDAD, INTERNACIONAL, ARTES, PERIODISMO, GASTRONOMÍA, CIENCIA, TERAPIA, SALUD
20. Reemplazar los valores del atributo *estu\_nivel\_prgm\_academico* por "TECNOLOGIA" para aquellos programas cuyo nombre *estu\_prgm\_academico* inicie con la palabra 'TECNOLOGÍA' o 'TECNOLOGIA'
21. Reemplazar los valores del atributo *estu\_nivel\_prgm\_academico* por "TECNICO" para aquellos programas cuyo nombre *estu\_prgm\_academico* inicie con la palabra 'TECNICO', 'TÉCNICO', 'TÉCNICA', 'TECNICA'

22. Reemplazar los valores nulos del atributo *estu\_metodo\_prgm* por el valor no nulo del atributo *estu\_metodo\_prgm* de la institución con igual código (*inst\_cod\_institucion*) o nombre (*inst\_nombre\_institucion*), y de los programas del mismo código(*estu\_prac\_id\_prgrm\_academico*) o nombre (*estu\_prgm\_academico*)
23. Reemplazar los valores nulos del atributo *estu\_metodo\_prgm* por los valores del atributo *metodología* de la tabla *programas\_ies* teniendo en cuenta los códigos de las instituciones y de los programas
24. Reemplazar los 2 restantes valores nulos del atributo *estu\_metodo\_prgm* por la moda ‘PRESENCIAL’
25. Estandarizar los valores del atributo *estu\_metodo\_prgm* en mayúsculas;
26. Reemplazar los valores nulos del atributo *dipo\_codigomunicipio* por los valores del atributo *estu\_exam\_codmpio\_presentacion*:
27. Reemplazar los valores nulos del atributo *inst\_cod\_jornada* por la moda “1” para los nombres de instituciones que contengan la palabra ‘NORMAL’.
28. Reemplazar los valores nulos del atributo *inst\_cod\_jornada* por la moda del resto de instituciones que es ‘12’
29. Reemplazar los valores nulos del atributo *estu\_area\_conoc de saberpro\_limpio*, por el valor no nulo del atributo *estu\_area\_conoc* de los programas del mismo nombre (*estu\_prgm\_academico*) de la tabla *saberpro\_2012\_2014*
30. Reemplazar los valores nulos del atributo *estu\_nucleo\_pregrado de saberpro\_limpio*, por el valor no nulo del atributo *estu\_nucleo\_pregrado* de los programas del mismo nombre (*estu\_prgm\_academico*) de la tabla *saberpro\_2012\_2014*
31. Reemplazar los valores nulos del atributo *estu\_cod\_grupo\_ref de saberpro\_limpio*, por el valor no nulo del atributo *estu\_cod\_grupo\_ref* de los grupos de referencia con el nombre (*estu\_grupo\_referencia*) de la tabla *saberpro\_2012\_2014*
32. Actualizar en la tabla *saberpro\_limpio*, el atributo *estu\_cod\_grupo\_ref* utilizando en algunos casos caracteres comodines “%” y “\_” así: (utilice select case when)

BELLAS ARTES Y DISEÑO por el código 1,

CIENCIAS NATURALES Y EXACTAS por el código 2,

CIENCIAS SOCIALES por el código 3,

HUMANIDADES por el código 4,

DERECHO por el código 5,

COMUNICACI\_N, PERIODISMO Y PUBLICIDAD por el código 6,  
 CIENCIAS MILITARES Y NAVALES por el código 7,  
 %CIENCIAS AGROPECUARIAS% por el código 8,  
 %ADMINISTRACI\_N% por el código 9,  
 EDUCACI\_N por el código 10,  
 ARQUITECTURA Y URBANISMO por el código 11,  
 INGENIER\_A por el código 12,  
 %SALUD% por el código 13,  
 MEDICINA por el código 14,  
 %INGENIER\_A, INDUSTRIA Y MINAS% por el código 15  
 %TIC% por el código 17,  
 %ARTES - DISEÑO – COMUNICACI\_N% por el código 19,  
 %CIENCIAS AGROPECUARIAS% por el código 20,  
 NORMALES SUPERIORES por el código 27,  
 %JUDICIAL% por el código 28;  
 %MILITAR Y POLICIAL% por el código 29  
 RECREACI\_N Y DEPORTES por el código 30,  
 %ECONOM% por el código 31,  
 %CONTADUR\_A% por el código 32,  
 PSICOLOG\_A por el código 33,  
 ENFERMER\_A por el código 34;  
 GRUPO REFERENCIA NACIONAL% por el código 42,

33. Reemplazar los nulos del atributo *estu\_semestre\_cursa* por la moda (de las NORMALES) del atributo *estu\_semestre\_cursa* de la tabla *saberpro\_2012\_2014*, para los nombres de instituciones que contengan la palabra 'NORMAL'.
34. Reemplazar los valores nulos del atributo *estu\_semestre\_cursa* de *saberpro\_limpio*, por el valor no nulo del atributo *estu\_semestre\_cursa* de la tabla *saberpro\_2012\_2014* de igual instituciones con los mismos programas
35. Reemplazar los nulos del atributo *estu\_pje\_creditos* por el código '0' para aquellas instituciones cuyo valor del atributo *estu\_nivel\_prgm\_academico* es 'NORMALISTA' o el nombre de la institución del atributo *inst\_nombre\_institucion* contiene la sigla SENA.
36. Reemplazar los valores nulos del atributo *estu\_pje\_creditos* por el valor no nulo del

- atributo *estu\_pje\_creditos* de la institución con igual código (*inst\_cod\_institucion*) o nombre (*inst\_nombre\_institucion*), y de los programas del mismo código(*estu\_prac\_id\_prgrm\_academico*) o nombre (*estu\_prgrm\_academico*)
37. Reemplazar los valores nulos del atributo *inst\_vlr\_matricula\_ant* por la moda (de las NORMALES) del atributo *inst\_vlr\_matricula\_ant* de la tabla *saberpro\_2012\_2014*, para los nombres de instituciones que contengan la palabra ‘NORMAL’.
  38. Reemplazar los valores nulos del atributo *inst\_vlr\_matricula\_ant* de *saberpro\_limpio*, por el valor no nulo del atributo *inst\_vlr\_matricula\_ant* de la tabla *saberpro\_2012\_2014* de igual instituciones con los mismos programas
  39. Reemplazar los valores nulos del atributo *estu\_titulo\_bto* por la moda (de las NORMALES) del atributo *estu\_titulo\_bto* de la tabla *saberpro\_2012\_2014*, para los nombres de instituciones que contengan la palabra ‘NORMAL’.
  40. Reemplazar los valores nulos del atributo *estu\_titulo\_bto* al resto de instituciones por la moda del atributo *estu\_titulo\_bto* de la tabla *saberpro\_2012\_2014*
  41. Reemplazar los valores nulos del atributo *estu\_hogar\_actual* por la moda del atributo *estu\_hogar\_actual* de la tabla *saberpro\_2012\_2014*
  42. Reemplazar los valores nulos del atributo *fami\_num\_pers\_grup\_fam* por la media (redondeada a entera) de los valores no nulos de *fami\_num\_pers\_grup\_fam* de la tabla *saberpro\_2012\_2014*;
  43. Reemplazar los valores nulos del atributo *estu\_sn\_cabeza\_fmilia* por la moda del atributo *estu\_sn\_cabeza\_fmilia* de la tabla *saberpro\_2012\_2014*.
  44. Reemplazar los valores nulos del atributo *fami\_num\_pers\_cargo* por la media (redondeada a entera) de los valores no nulos de *fami\_num\_pers\_cargo* de la tabla *saberpro\_2012\_2014*;
  45. Reemplazar los valores nulos del atributo *fami\_cod\_educa\_padre* por el código ‘99’
  46. Reemplazar los valores 1,2,3,4,5,6,7,8 del atributo *fami\_cod\_educa\_padre* por el código ‘9’
  47. Reemplazar los valores nulos del atributo *fami\_cod\_educa\_madre* por el código ‘99’
  48. Reemplazar los valores 1,2,3,4,5,6,7,8 del atributo *fami\_cod\_educa\_madre* por el código ‘9’
  49. Reemplazar los valores nulos del atributo *fami\_cod\_ocup\_madre* por el código ‘99’
  55. Reemplazar los valores 01,02,03,04,05,06,07,08,09,10,11,12 del atributo *fami\_cod\_ocup\_madre* por la moda del atributo *fami\_cod\_ocup\_madre* de la tabla

*saberpro\_2012\_2014*

56. Reemplazar los valores nulos del atributo *estu\_estrato* por la moda del atributo *estu\_estrato* de la tabla *saberpro\_2012\_2014*
57. Reemplazar los valores nulos del atributo *fami\_nivel\_sisben* por la moda del atributo *fami\_nivel\_sisben* de la tabla *saberpro\_2012\_2014*
58. Reemplazar los valores nulos del atributo *econ\_material\_pisos* por la moda del atributo *econ\_material\_pisos* de la tabla *saberpro\_2012\_2014*
59. Reemplazar los valores nulos del atributo *econ\_sn\_telefonia* por la moda del atributo *econ\_sn\_telefonia* de la tabla *saberpro\_2012\_2014*
60. Reemplazar los valores nulos del atributo *econ\_sn\_internet* por la moda del atributo *econ\_sn\_internet* de la tabla *saberpro\_2012\_2014*
61. Reemplazar los valores nulos del atributo *econ\_sn\_servicio\_tv* por la moda del atributo *econ\_sn\_servicio\_tv* de la tabla *saberpro\_2012\_2014*
62. Reemplazar los valores nulos del atributo *econ\_sn\_computador* por la moda del atributo *econ\_sn\_computador* de la tabla *saberpro\_2012\_2014*
63. Actualizar el atributo *econ\_sn\_computador* por 1 para aquellos valores del atributo *econ\_sn\_computador* > 0.
64. Reemplazar los valores nulos del atributo *econ\_sn\_celular* y los de *econ\_sn\_celular* código 2 por la moda del atributo *econ\_sn\_celular* de la tabla *saberpro\_2012\_2014*
65. Reemplazar los valores nulos del atributo *econ\_sn\_dvd* y los de *econ\_sn\_dvd* código 2 por la moda del atributo *econ\_sn\_dvd* de la tabla *saberpro\_2012\_2014*
66. Reemplazar los valores nulos del atributo *econ\_sn\_lavadora* por la moda del atributo *econ\_sn\_lavadora* de la tabla *saberpro\_2012\_2014*
67. Reemplazar los valores nulos del atributo *econ\_sn\_microondas* por la moda del atributo *econ\_sn\_microondas* de la tabla *saberpro\_2012\_2014*
68. Reemplazar los valores nulos del atributo *econ\_sn\_automovil* por la moda del atributo *econ\_sn\_automovil* de la tabla *saberpro\_2012\_2014*
69. Reemplazar los valores nulos del atributo *econ\_sn\_horno* por la moda del atributo *econ\_sn\_horno* de la tabla *saberpro\_2012\_2014*
70. Reemplazar los valores nulos del atributo *econ\_sn\_nevera* por la moda del atributo *econ\_sn\_nevera* de la tabla *saberpro\_2012\_2014*
71. Reemplazar los valores nulos del atributo *infa\_dormitorios* por la moda del atributo *infa\_dormitorios* de la tabla *saberpro\_2012\_2014*

72. Reemplazar los valores nulos del atributo *fami\_ing\_fmliar\_mensual* por la moda del atributo *fami\_ing\_fmliar\_mensual* de la tabla *saberpro\_2012\_2014*
73. Reemplazar los valores nulos del atributo *estu\_trabaja* y los valores de *estu\_trabaja* códigos 1,6,7 por la moda del atributo *estu\_trabaja* de la tabla *saberpro\_2012\_2014*
74. Reemplazar los valores nulos del atributo *estu\_horas\_trabajo* por 0 para los estudiantes que no trabajan.
75. Reemplazar los valores nulos del atributo *estu\_horas\_trabajo* por la media entera de los valores del atributo *estu\_horas\_trabajo* de los estudiantes cuyo atributo *estu\_trabaja* es de código '3' de la tabla *saberpro\_2012\_2014*. Actualizar únicamente para los estudiantes que trabajan con código '3'
76. Reemplazar los valores nulos del atributo *estu\_horas\_trabajo* por la media entera de los valores del atributo *estu\_horas\_trabajo* de los estudiantes cuyo atributo *estu\_trabaja* es de código '4' de la tabla *saberpro\_2012\_2014*. Actualizar únicamente para los estudiantes que trabajan con código '4'
77. Reemplazar los valores nulos del atributo *estu\_horas\_trabajo* por la media entera de los valores del atributo *estu\_horas\_trabajo* de los estudiantes cuyo atributo *estu\_trabaja* es de código '5' de la tabla *saberpro\_2012\_2014*. Actualizar únicamente para los estudiantes que trabajan con código '5'
78. Reemplazar la coma por punto decimal a los valores del atributo *mod\_lectura\_critica*
79. Reemplazar los valores nulos del atributo *mod\_lectura\_critica* por 'NP'
80. Reemplazar la coma por punto decimal a los valores del atributo *mod\_comunica\_escrita\_punt*
81. Reemplazar los valores nulos del atributo *mod\_comunica\_escrita\_punt* por 'NP'
82. Estandarizar a punto decimal los valores del atributo *mod\_razona\_cuantitativo\_punt*
83. Reemplazar los valores nulos del atributo *mod\_razona\_cuantitativo\_punt* por 0
84. Estandarizar a punto decimal los valores del atributo *mod\_ingles\_punt*.
85. Reemplazar los valores nulos del atributo *mod\_ingles\_punt* por 0
86. Estandarizar a punto decimal los valores del atributo *mod\_comp\_ciudadanas\_punt*.
87. Reemplazar los valores nulos del atributo *mod\_comp\_ciudadanas\_punt* por 0