



Ingeniería en Computación

IC-6200 Inteligencia Artificial

Juan Manuel Esquivel Rodriguez

Proyecto Corto 1 - Simulador Votantes

Fabio Mora Cubillo - 2013012801

Sergio Moya Valerin - 2013015682

Gabriel Venegas Castro - 2013115967

I Semestre, 2018

Naturaleza del proyecto

El proyecto consiste en generar datos sintéticos que sigan la distribución teórica de los datos de las elecciones nacionales de Costa Rica del año 2018. Consta de dos funciones principales:

- `generar_muestra_pais(n)`: retornará una muestra simulada con n votantes con etiquetas y atributos proporcionales a la distribución de resultados nacionales.
- `generar_muestra_provincia(n, nombre_provincia)`: retornará una muestra simulada siguiendo la distribución de datos para una provincia específica (donde `nombre_provincia` es uno de "SAN JOSE", "ALAJUELA", "CARTAGO", "HEREDIA", "GUANACASTE", "PUNTARENAS", "LIMON")

Datos a utilizar

Los datos a utilizar en el proyecto son archivos csv generados a partir de los siguientes:

- Actas de escrutinio: Votos por junta receptora:
http://www.tse.go.cr/elecciones2018/actas_escrutinio.htm
- Mapeo juntas a cantones: Información sobre mapeo juntas-cantones:
<http://www.tse.go.cr/pdf/nacional2018/JRV.pdf>
- Indicadores cantonales: Atributos e indicadores cantonales:
https://www.estadonacion.or.cr/files/biblioteca_virtual/otras_publicaciones/Indicadores_cantonales_Censos2000y2011.xlsx

Decisiones técnicas sobre los datos

1. Se toman en cuenta los indicadores de ruralidad/urbanidad.
2. La relación de cantidad entre hombres y mujeres se modifica mediante software para ser usado como porcentaje. $(\text{Hombres} * 100 / (\text{Hombres} + \text{Mujeres}) = \text{Porcentaje de hombres})$
3. El dato "relación de dependencia" se toma como sólo las personas mayores de 65, dejando fuera a las menores de 15, ya que no forman parte del padrón electoral.
4. Se incluyen en la generación de la población todos los posibles indicadores.
5. Se agrega un indicador "PROVINCIA" en la penúltima posición de la lista de indicadores
6. Se agrega un indicador "VOTO" en la última posición de la lista de indicadores
7. El tamaño de la lista interna que retornan las funciones (la que contiene todos los indicadores) es $[\text{numero_atributos} + 2]$, ya que se añaden PROVINCIA y VOTO.

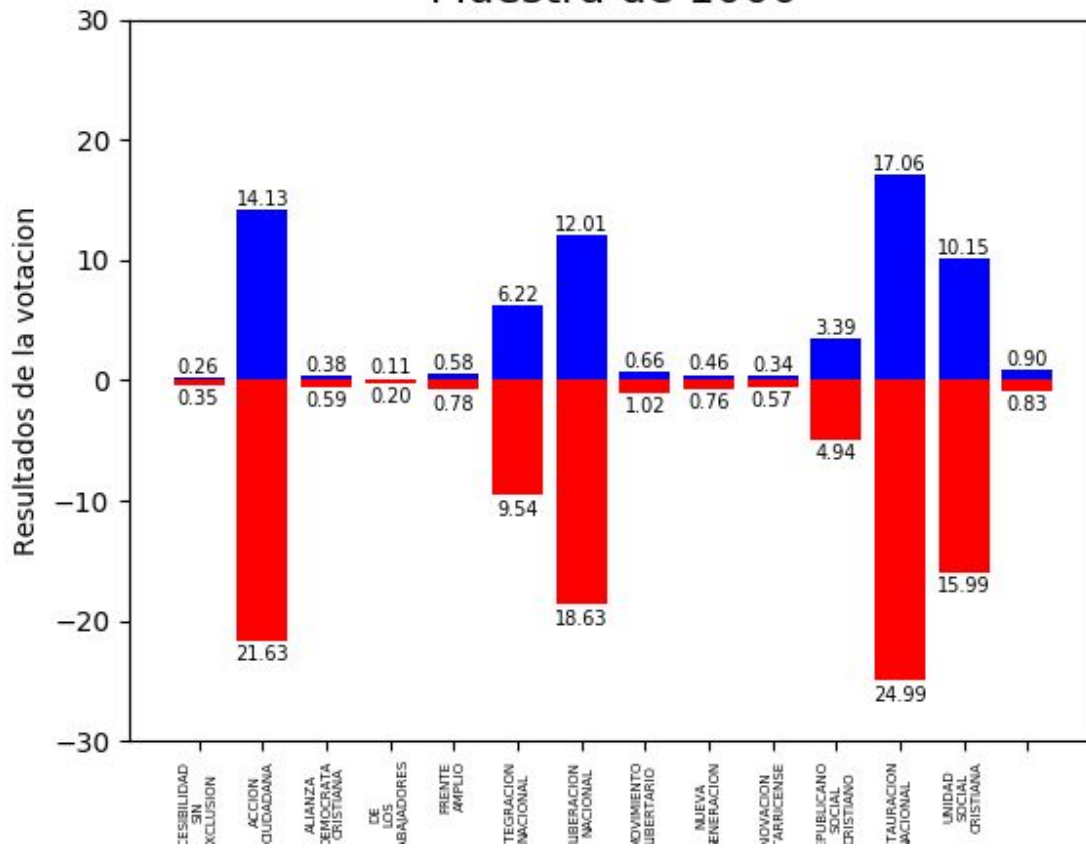
Otras consideraciones

1. Para las pruebas pytest, se toman los porcentajes como un valor flotante entre 0.0 y 1.0
2. Para las pruebas pytest, se hacen pruebas de aleatoriedad para verificar que el generador de números aleatorios tenga el comportamiento esperado. Se puede especificar el tamaño de la muestra y el porcentaje de tolerancia con respecto al porcentaje fijo establecido.
3. Para muestras pequeñas, no tiene sentido mantener un valor de tolerancia bajo; para muestras medianas y grandes, la distribución converge, permitiendo tener valores de tolerancia más bajos.

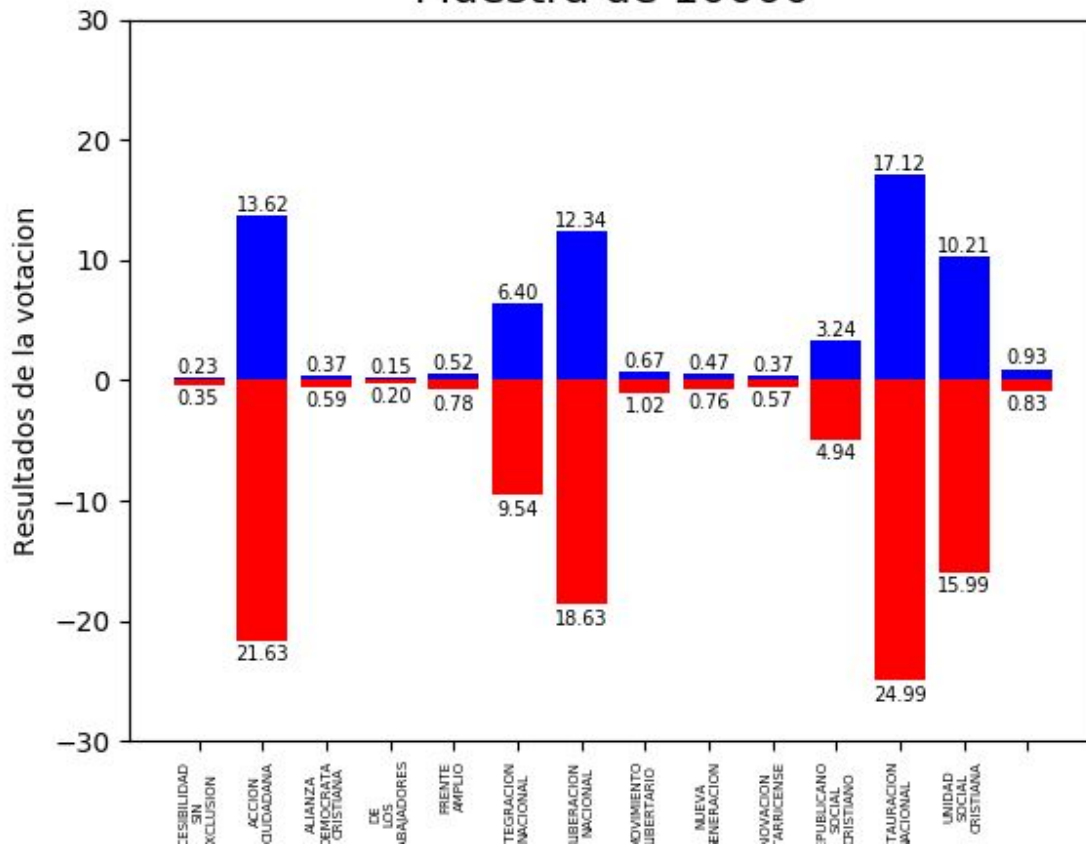
Análisis de resultados

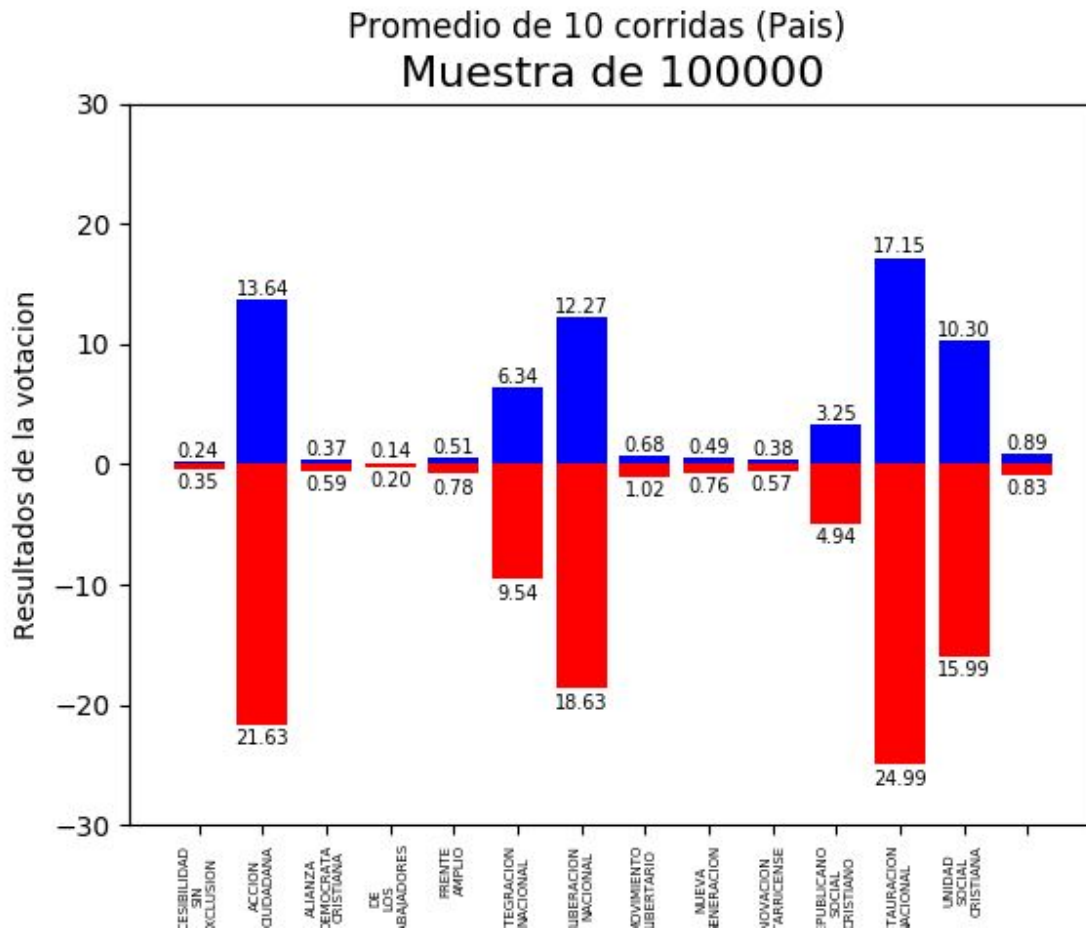
Para el análisis de los resultados se determinó un estándar que consiste en promediar 10 corridas de cada caso. Los casos son por país y por provincias, con 1000 y 10000 muestras. Se mostrarán una serie de gráficos y su descripción relativa. En dichos gráficos las barras azules equivalen al valor promediado de las 10 corridas y el rojo al valor real provisto por el TSE, además el último valor a la derecha son los votos vacíos y blancos. Se utilizó matplotlib para automatizar la generación de los gráficos.

Promedio de 10 corridas (Pais) Muestra de 1000



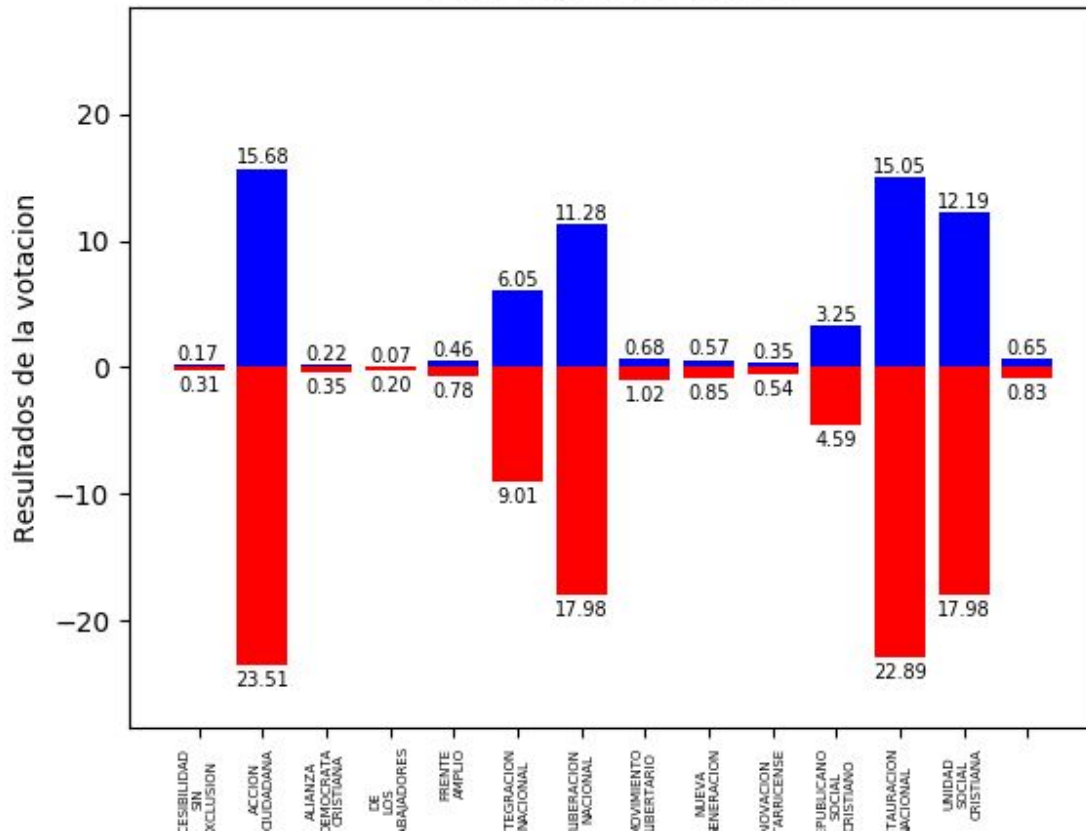
Promedio de 10 corridas (Pais) Muestra de 10000



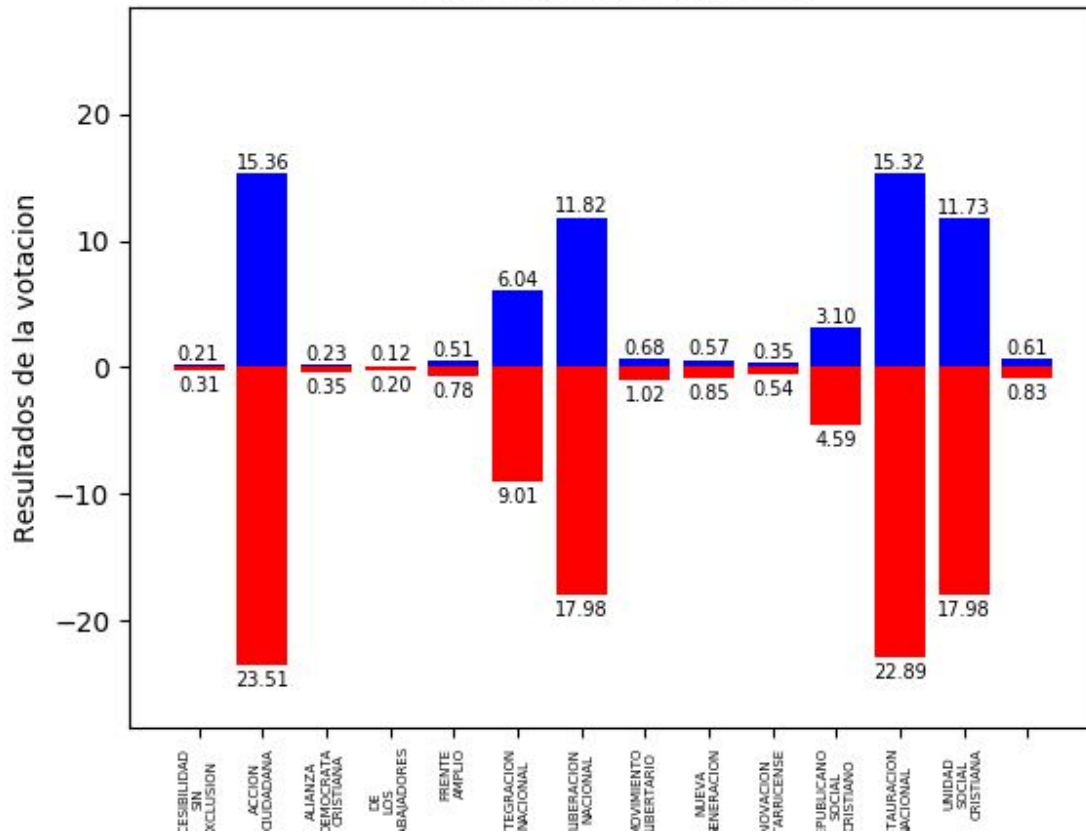


Al generar la muestra de cien mil para el país se determinó que la generación de muestras es lo suficientemente sólida e independientemente del tamaño de la muestra, mientras esta sea lo suficientemente grande, siempre lanzará resultados casi iguales.

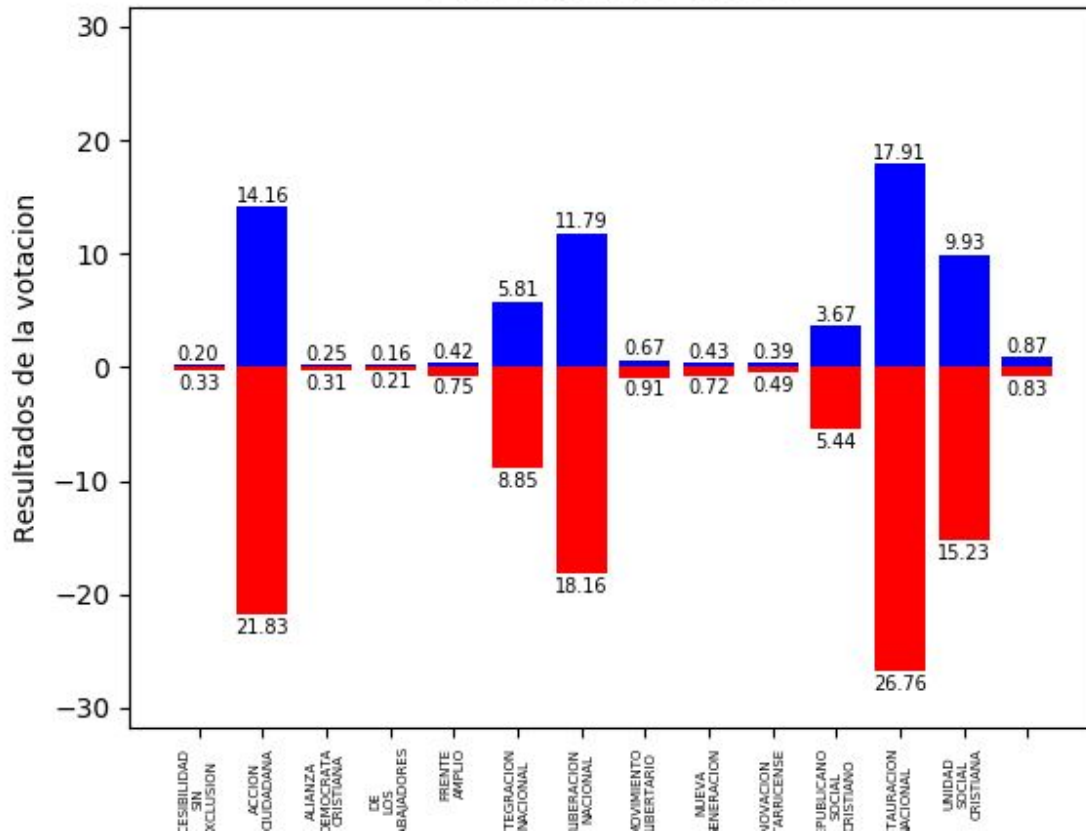
Promedio de 10 corridas (San Jose) Muestra de 1000



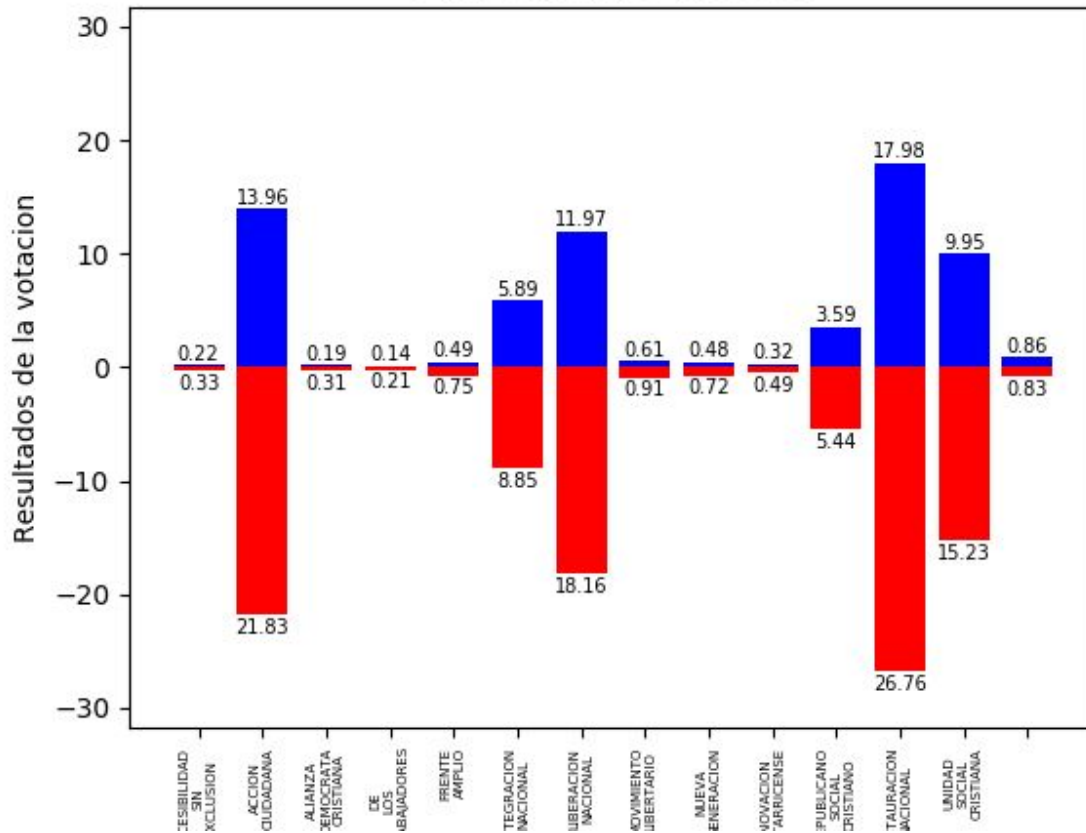
Promedio de 10 corridas (San Jose) Muestra de 10000



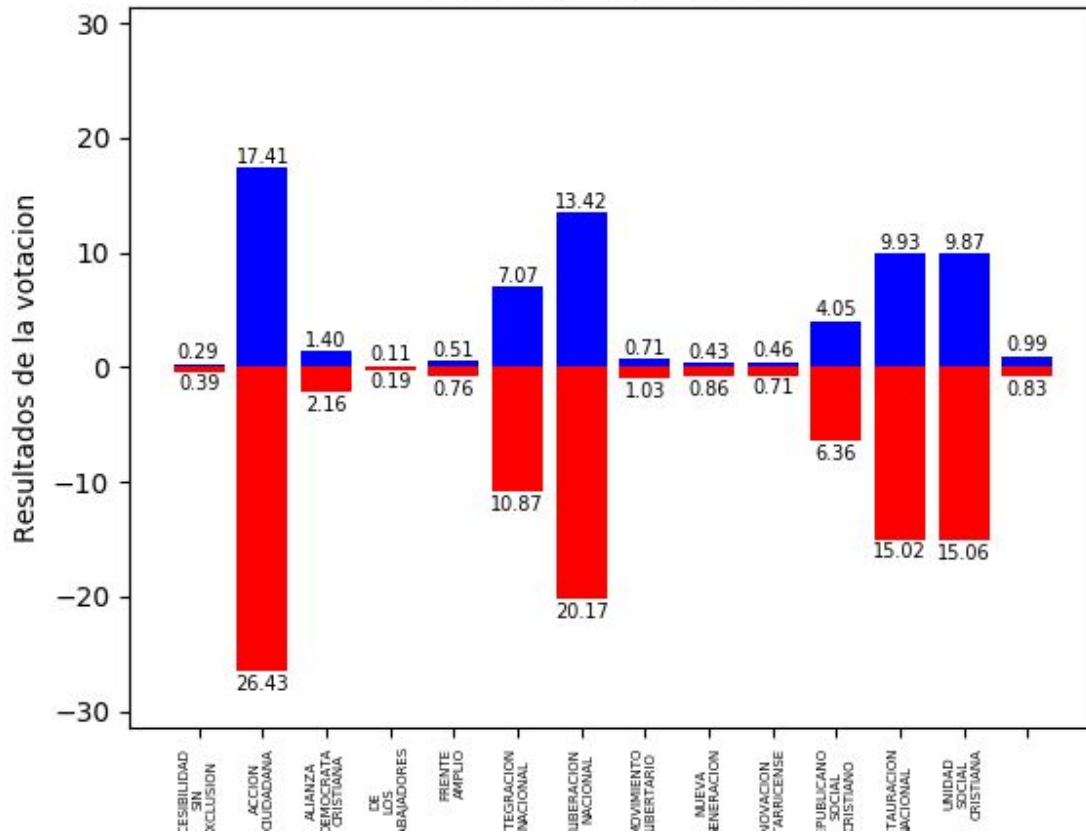
Promedio de 10 corridas (Alajuela) Muestra de 1000



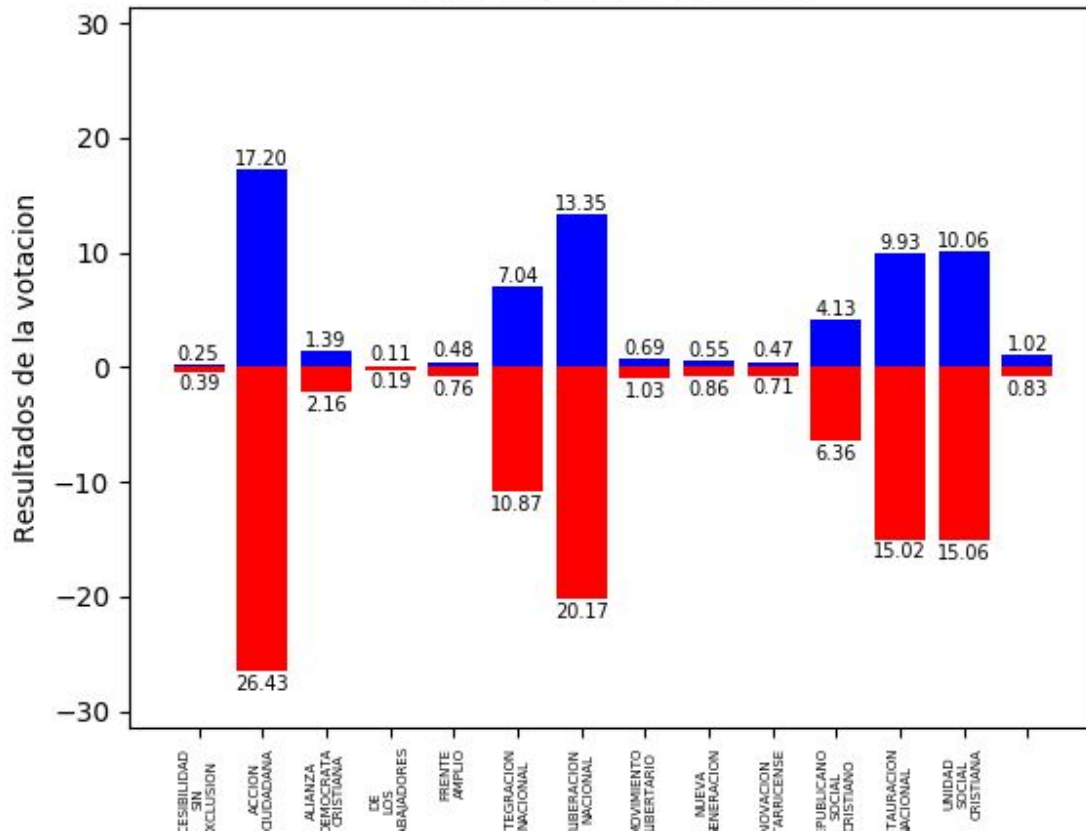
Promedio de 10 corridas (Alajuela) Muestra de 10000



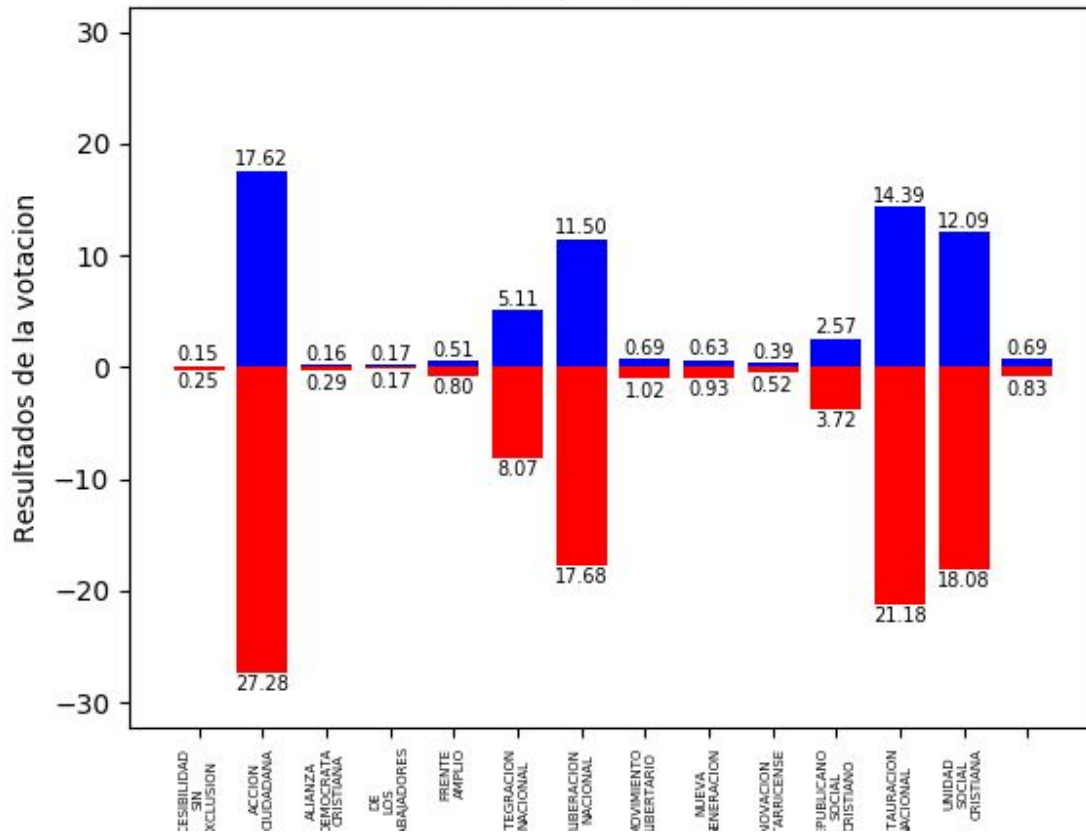
Promedio de 10 corridas (Cartago) Muestra de 1000



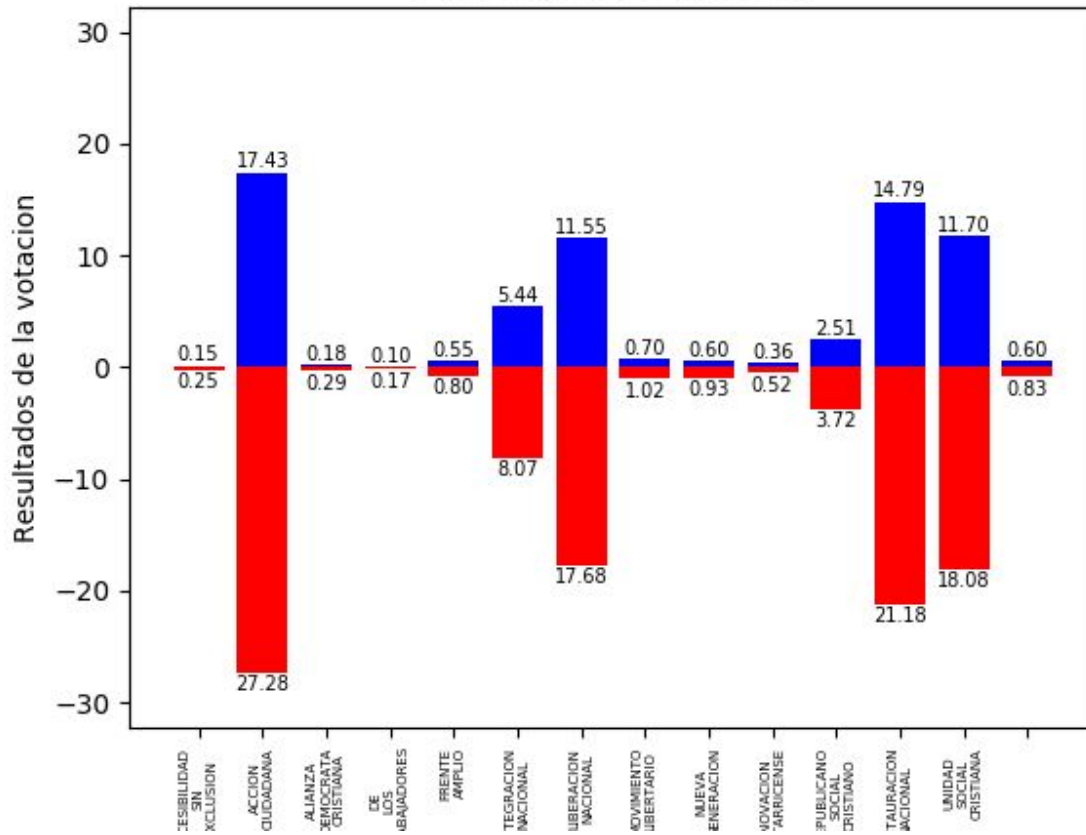
Promedio de 10 corridas (Cartago)
Muestra de 10000



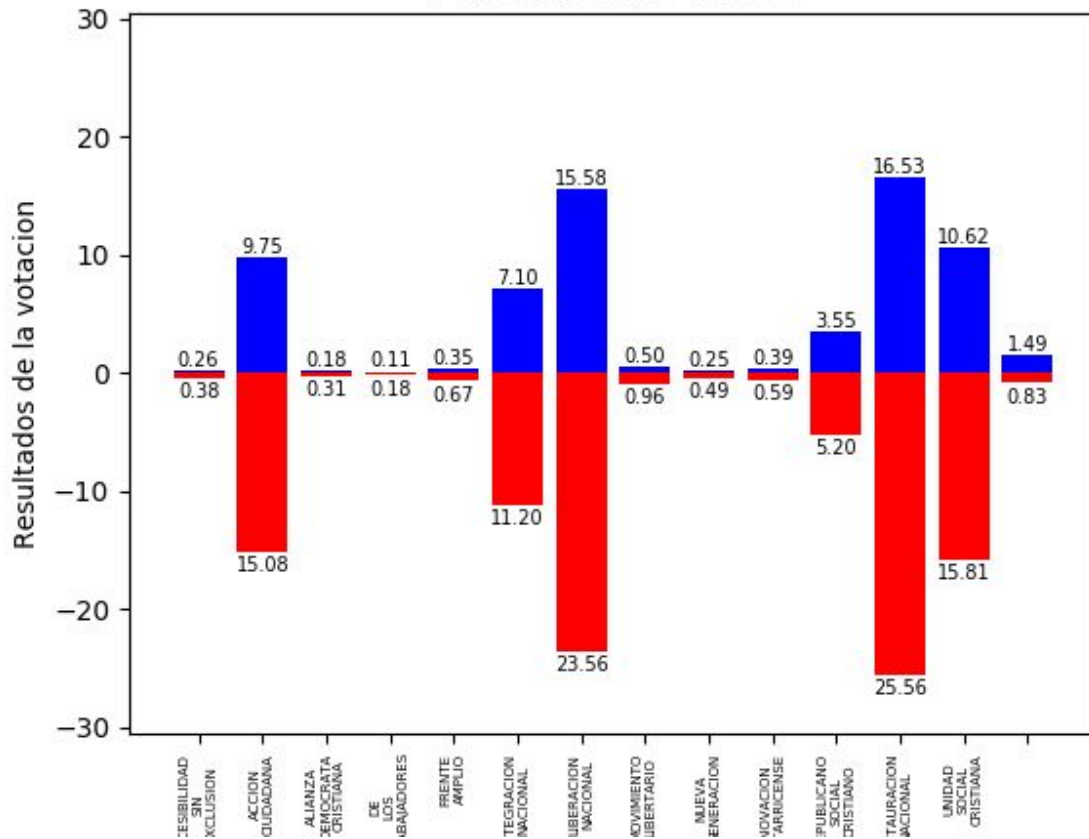
Promedio de 10 corridas (Heredia) Muestra de 1000



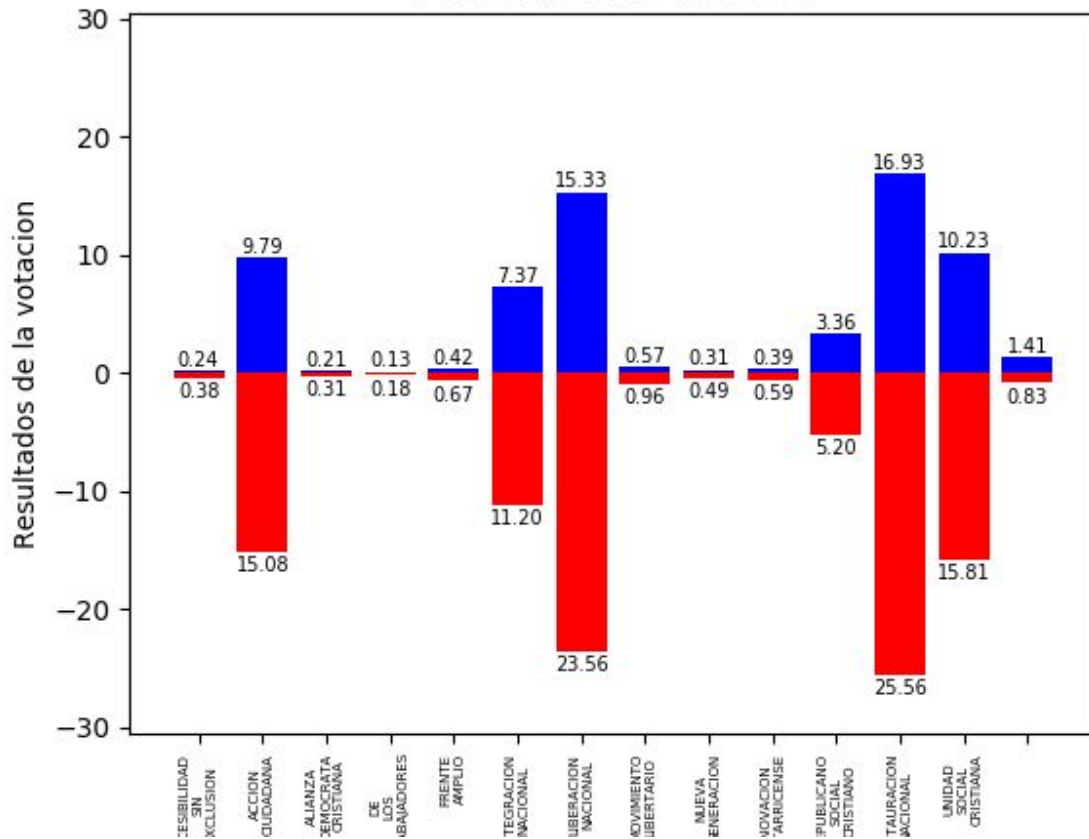
Promedio de 10 corridas (Heredia) Muestra de 10000



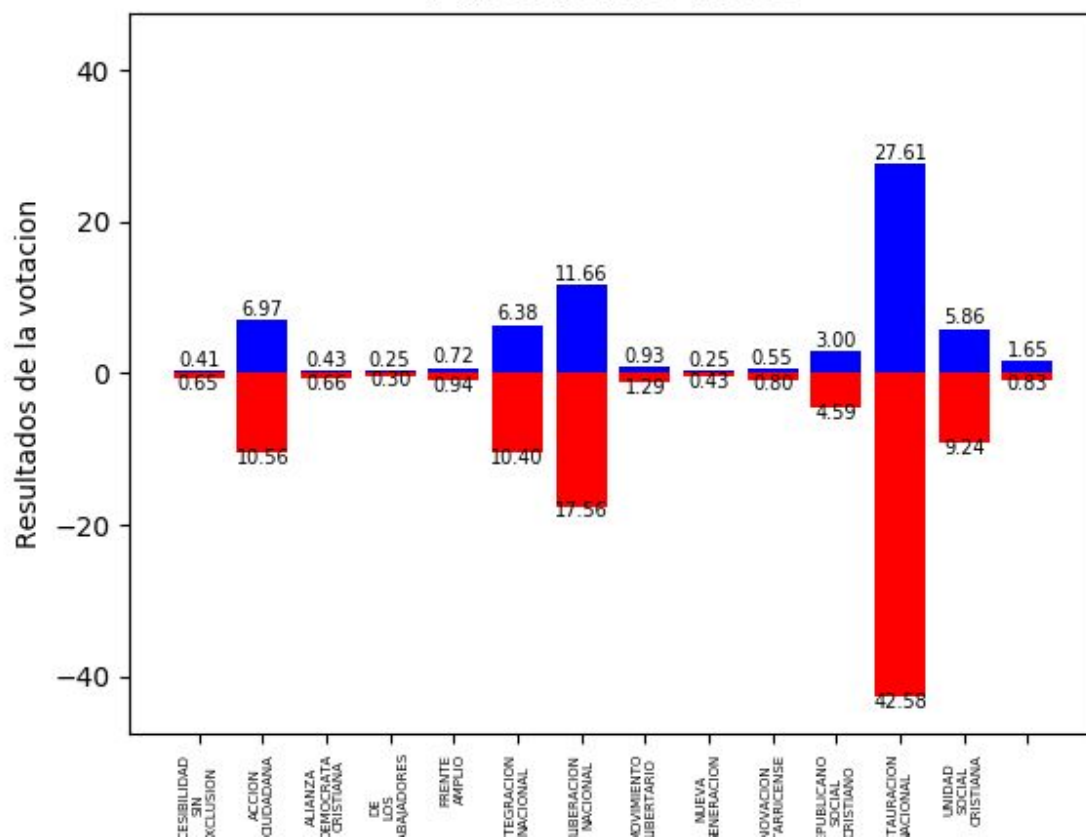
Promedio de 10 corridas (Guanacaste) Muestra de 1000



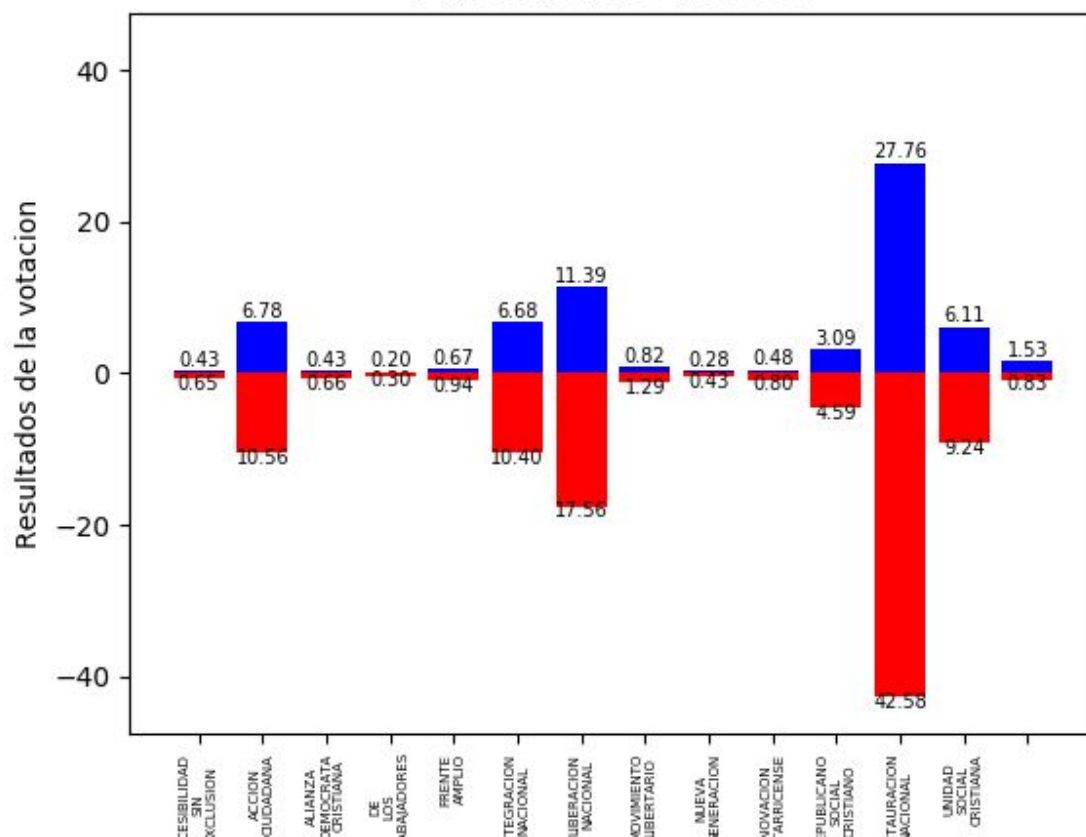
Promedio de 10 corridas (Guanacaste)
Muestra de 10000



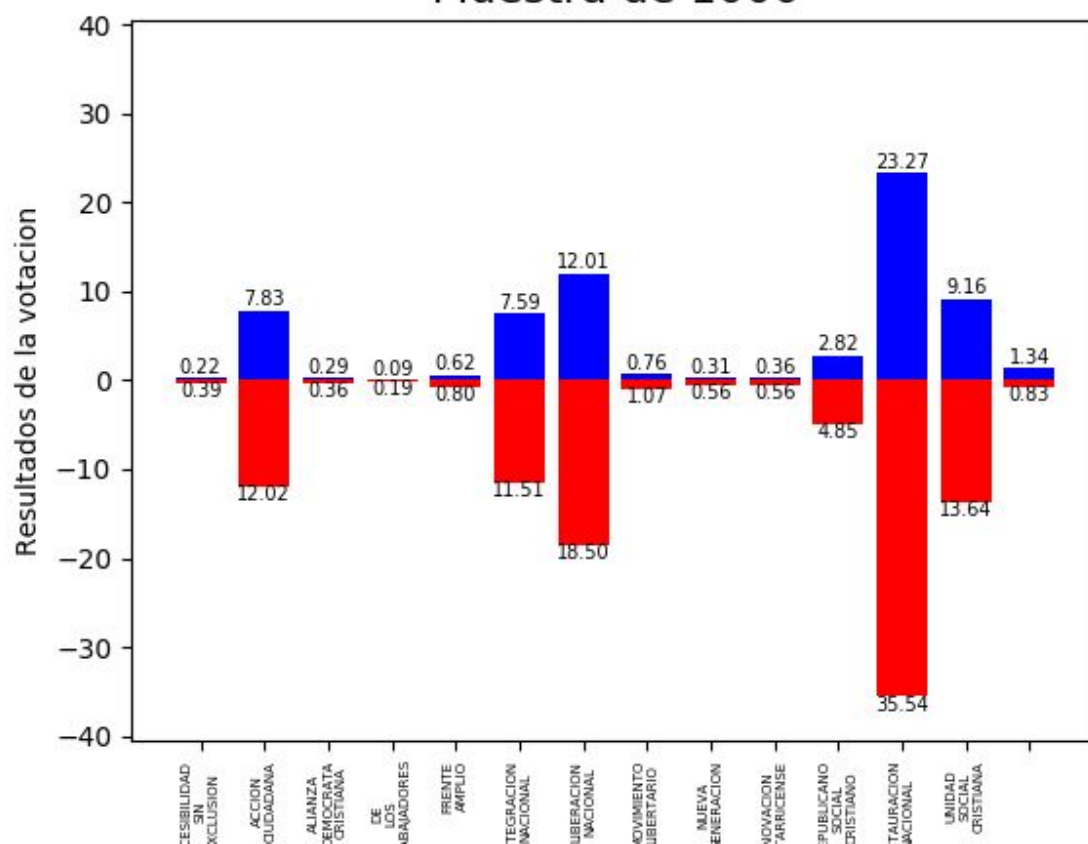
Promedio de 10 corridas (Limon) Muestra de 1000

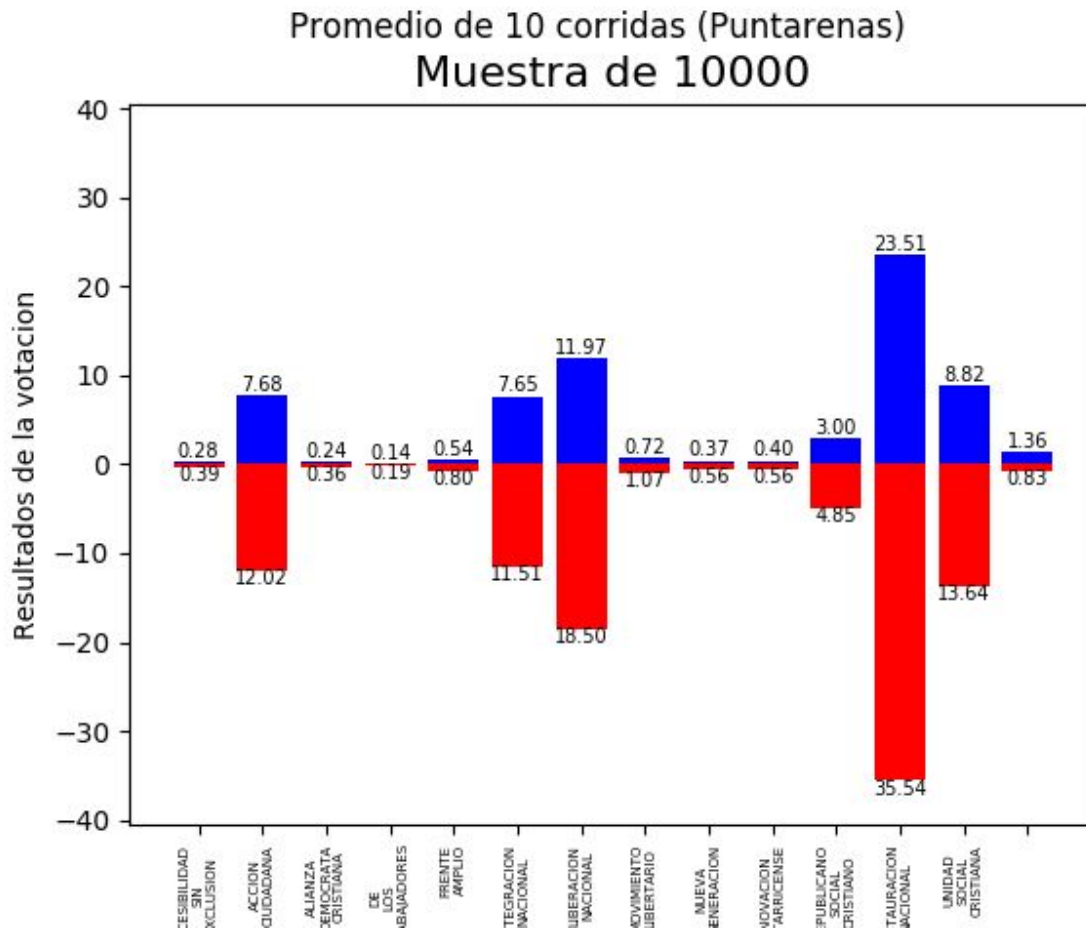


Promedio de 10 corridas (Limon) Muestra de 10000



Promedio de 10 corridas (Puntarenas) Muestra de 1000





Finalmente se tiene que, pese a que el porcentaje de votos para cada partido baja con respecto al porcentaje real, es esperable ya que el valor correspondiente a los porcentajes de la muestra es un promedio. Lo importante es visualizar que la distribución de votos se mantiene y el partido ganador de cada provincia en los datos reales también gana en los datos de la muestra, por lo que se determina que el generador de población es válido para representar la población de votantes en el contexto dado.