

Modelos para la implementación de monitoreo de homicidios

Felipe González

Diciembre, 2018

Contents

1	Introducción	1
2	Datos	1
3	Metodología	2
4	Métodos de monitoreo: bayesiano simple	2
5	Métodos de monitoreo: Farrington	7
6	Otros ejemplos	10
7	Propuesta: Método mixto	14
7.1	Método mixto para todos los municipios.	15
	References	18

1 Introducción

En este documento mostramos fundamentos y ejemplos de implementación de un sistema de monitoreo de violencia en México basado en los datos de homicidios a nivel municipal provistos por el Secretariado Ejecutivo del Sistema Nacional de Seguridad Pública (“Datos Abiertos de Incidencia Delictiva” 2018a)).

El objetivo es producir un detector de anomalías (datos inusualmente altos en incidencia) de los datos reportados, con el fin de enfocar el análisis a acciones a municipios de interés. Esos municipios de interés deberán ser identificados por reportes inusualmente altos de incidencia de homicidio, en contraste a aquellos municipios donde los datos reportados tengan variación consistente con la que se ha observado históricamente.

Aunque el interés general es medición de violencia, la decisión de utilizar solamente los datos de homicidios (homicidios más feminicidios según la nueva metodología de la SESNSP) se debe a que otros tipos de crímenes violentos tienden a tener niveles altos de cifra negra – alrededor de 90% o más según la Encuesta de Victimización del INEGI, mientras que esto sucede en mucho menor grado con homicidios. Otras indicadores de crimen con cifra negra relativamente baja son el robo de automóvil, por ejemplo, pero consideramos que este no necesariamente refleja el interés en medir la violencia que experimenta la población en cada municipio.

2 Datos

Usamos los datos de homicidios y feminicidios a nivel municipal provistos por el Secretariado Ejecutivo del Sistema Nacional de Seguridad Pública (“Datos Abiertos de Incidencia Delictiva” 2018a)). Estos datos son publicados cada mes.

Otra fuente de este tipo de datos son los reportados por INEGI, que se consideran más completos y precisos. Los datos de INEGI, sin embargo, son publicados con retraso de un año, lo que los descalifica para un sistema de monitoreo que pretende evaluar la situación actual de los municipios.

```
library(tidyverse)
library(lubridate)
library(surveillance)
```

El archivo de datos que utilizamos es preprocesado de los datos crudos según el código de <https://github.com/diegovalle/new.crimenmexico>, a partir de los datos del Secretariado Ejecutivo del Sistema Nacional de Seguridad Pública.

```
nm <- read_csv("./data/nm-fuero-comun-municipios.csv", progress = FALSE)
nm <- nm %>% mutate(date = ymd(paste0(date, '-01')))
tipo_homicidio <- c("FEMINICIDIO", "HOMICIDIO")
homicidios <- nm %>% filter(tipo %in% tipo_homicidio) %>%
  select(state_code, state, mun_code, municipio, date, count, population) %>%
  group_by(state_code, state, mun_code, municipio, date) %>%
  summarise(population = first(population), count = sum(count))
```

3 Metodología

Consideramos que el enfoque apropiado para construir monitores de violencia es el adoptado en la epidemiología para detección de brotes de enfermedades infecciosas ((Farrington and Andrews 2003)). La detección temprana de estos brotes en un aspecto importante de sistemas de monitoreo epidemiológicos, y están implementados en muchos países (Farrington and Andrews 2003) con el fin de proveerlos de respuesta epidemiológica oportuna.

Existen varias aproximaciones para la construcción de estos sistemas de detección temprana. Por ejemplo:

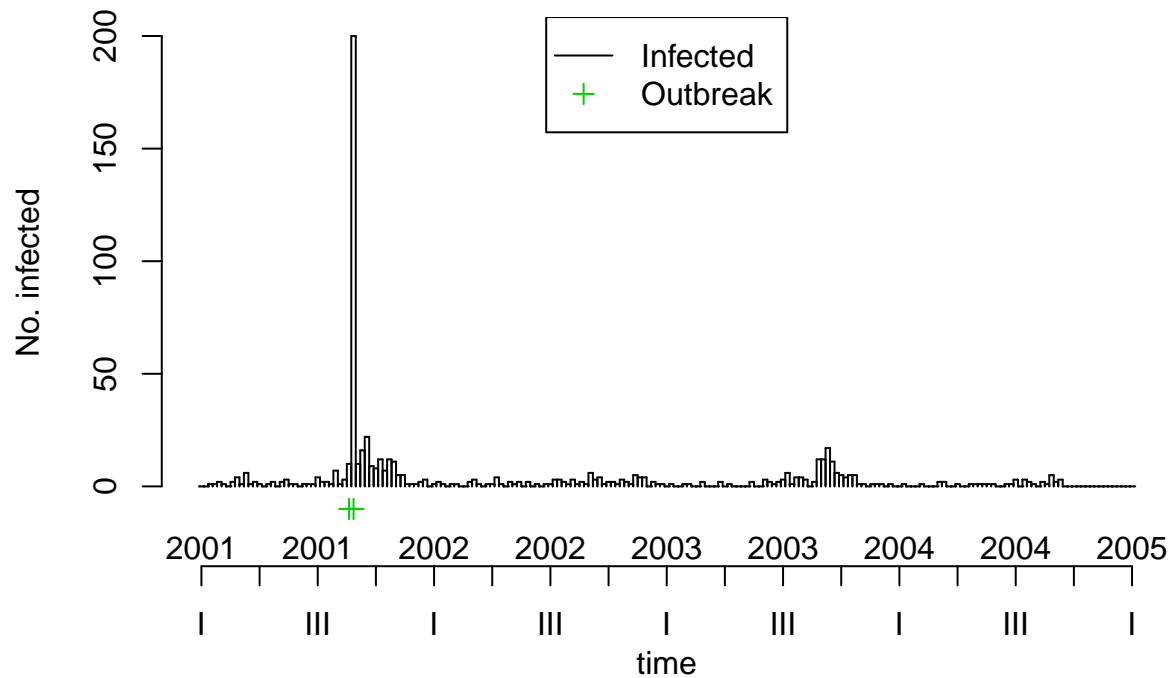
- CDC y otros: utilizan media, varianza y utilizan una aproximación normal simple. Algunos modelos simples bayesianos son similares a este método.
- Agencias europeas de monitoreo: método de Farrington y variaciones. Similar al modelo bayesiano, pero con un modelo para la media que puede incluir tendencia y estacionalidad.

En este documento, examinaremos algunos de los que más ampliamente implementados, como algunos métodos bayesianos y el método de Farrington (Meyer, Held, and Höhle 2017). Utilizamos las implementaciones del paquete de R surveillance.

4 Métodos de monitoreo: bayesiano simple

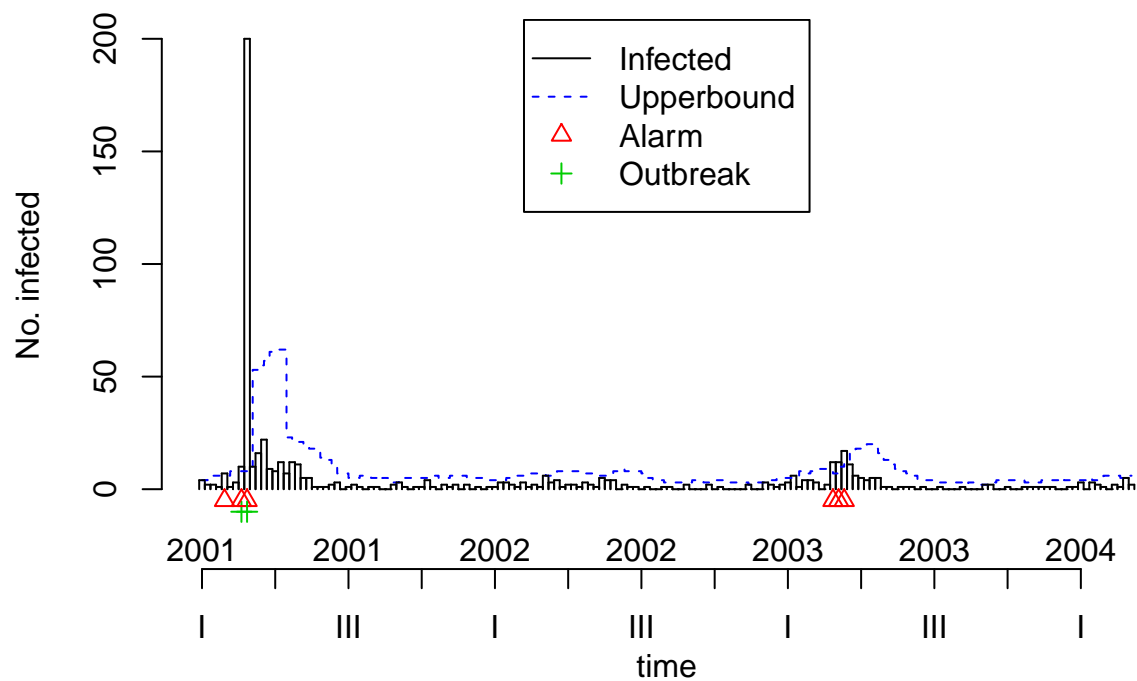
Uno de los métodos más simples, de desempeño competitivo ((Yang et al. 2018)) utiliza un enfoque parcialmente bayesiano construido a partir de una regresión Poisson con sobredispersión. Este es un ejemplo típico utilizado en la detección de brotes de *Kryptosporidiose* en Alemania (Höhle 2007):

```
data(k1)
plot(k1)
```



```
# Evaluar de 27 a 192
k1_surveil <- algo.bayes(k1,
  control = list(range = 27:192, b = 0, w = 6, alpha = 0.01))
plot(k1_surveil)
```

Analysis of k1 using bayes(6,6,0)



Para aplicar este método, primero definimos:

- Valores de referencia para estimar parámetros del modelo, que son las observaciones pasadas de conteos de homicidios y_i .

- Región donde queremos crear predicciones, límites superiores, y detectar anomalías.

Una vez definida la ventana de valores de referencia y horizonte de predicción, el modelo para nuevas observaciones (ver referencia) y_k es

- $P(y_k|\lambda)$ es Poisson con parámetro λ
- λ es Gamma con parámetros (α, β) .

De forma que la posterior para λ es Gamma con parámetros $(\alpha + \sum_i y_i, \beta + n)$, y la posterior predictiva es Poisson-Gamma con los mismos parámetros.

Una vez que calculamos la posterior predictiva, establecemos una alarma si la observación y_k satisface $y_k > y_\alpha$, donde y_α es el cuantil $1 - \alpha$ de la posterior predictiva.

En los ejemplos siguientes, pondremos

- $\alpha = 0.02$
- Observaciones de 6 meses anteriores para estimar parámetros, solo del año corriente.

4.0.1 Ejemplos: método bayesiano.

Consideramos tres municipios de Puebla:

```
#puebla_edo <- filter(homicidios, state == "PUEBLA")
puebla_edo <- readRDS("data/puebla.rds")
puebla_edo %>% group_by(municipio) %>% summarise(count = sum(count)) %>%
  arrange(desc(count)) %>%
  filter(municipio %in% c("PUEBLA", "CHILCHOTLA", "HUAUCHINANGO"))

## # A tibble: 3 x 2
##   municipio    count
##   <chr>        <dbl>
## 1 PUEBLA         982
## 2 HUAUCHINANGO   152
## 3 CHILCHOTLA      4

huauchinango <- puebla_edo %>% filter(state == "PUEBLA", municipio == "HUAUCHINANGO")

## Warning: Detecting old grouped_df format, replacing `vars` attribute by
## `groups`

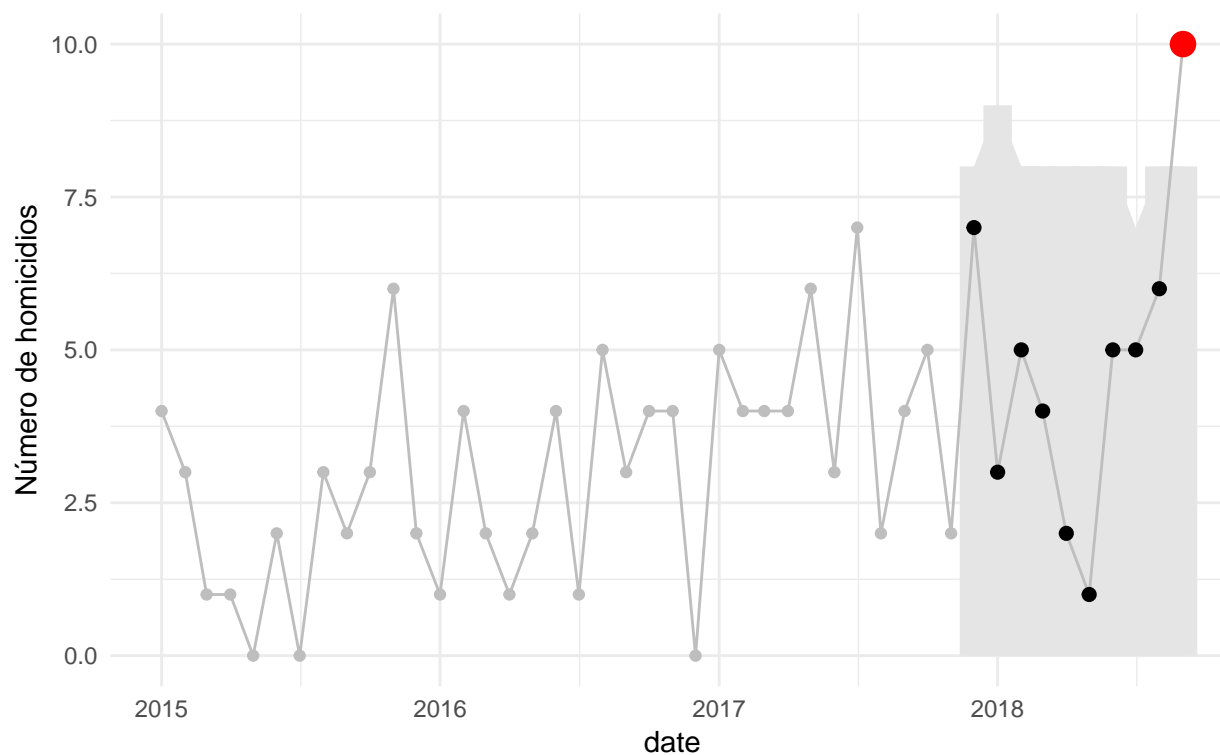
puebla <- puebla_edo %>% filter(state == "PUEBLA", municipio == "PUEBLA")
chilchotla <- puebla_edo %>% filter(state == "PUEBLA", municipio == "CHILCHOTLA")

source("./deteccion-homicidios.R")

monitor_huau <- monitor_bayes(huauchinango, alpha = 0.05)
graf_monitor(monitor_huau)
```

Huauchinango, Puebla

Método: bayes



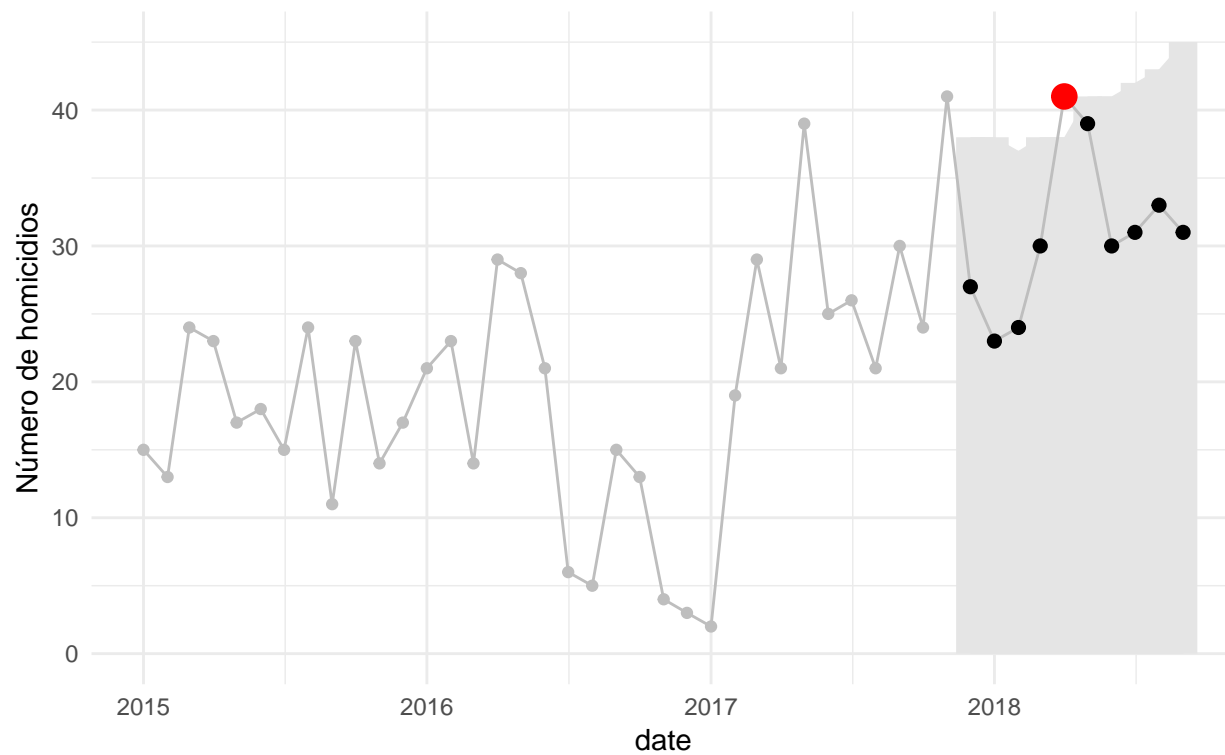
```
monitor_huau$data %>% tail
```

```
## # A tibble: 6 x 10
## # Groups:   state_code, state, mun_code, municipio [1]
##   state_code state mun_code municipio date      population count alerta
##   <int> <chr>    <int> <chr>    <date>         <int> <dbl> <lgl>
## 1      21 PUEB~      71 HUAUCHIN~ 2018-04-01     107811      2 FALSE
## 2      21 PUEB~      71 HUAUCHIN~ 2018-05-01     107898      1 FALSE
## 3      21 PUEB~      71 HUAUCHIN~ 2018-06-01     107986      5 FALSE
## 4      21 PUEB~      71 HUAUCHIN~ 2018-07-01     108072      5 FALSE
## 5      21 PUEB~      71 HUAUCHIN~ 2018-08-01     108158      6 FALSE
## 6      21 PUEB~      71 HUAUCHIN~ 2018-09-01     108244     10 TRUE
## # ... with 2 more variables: alerta_nivel <dbl>, metadata <chr>
```

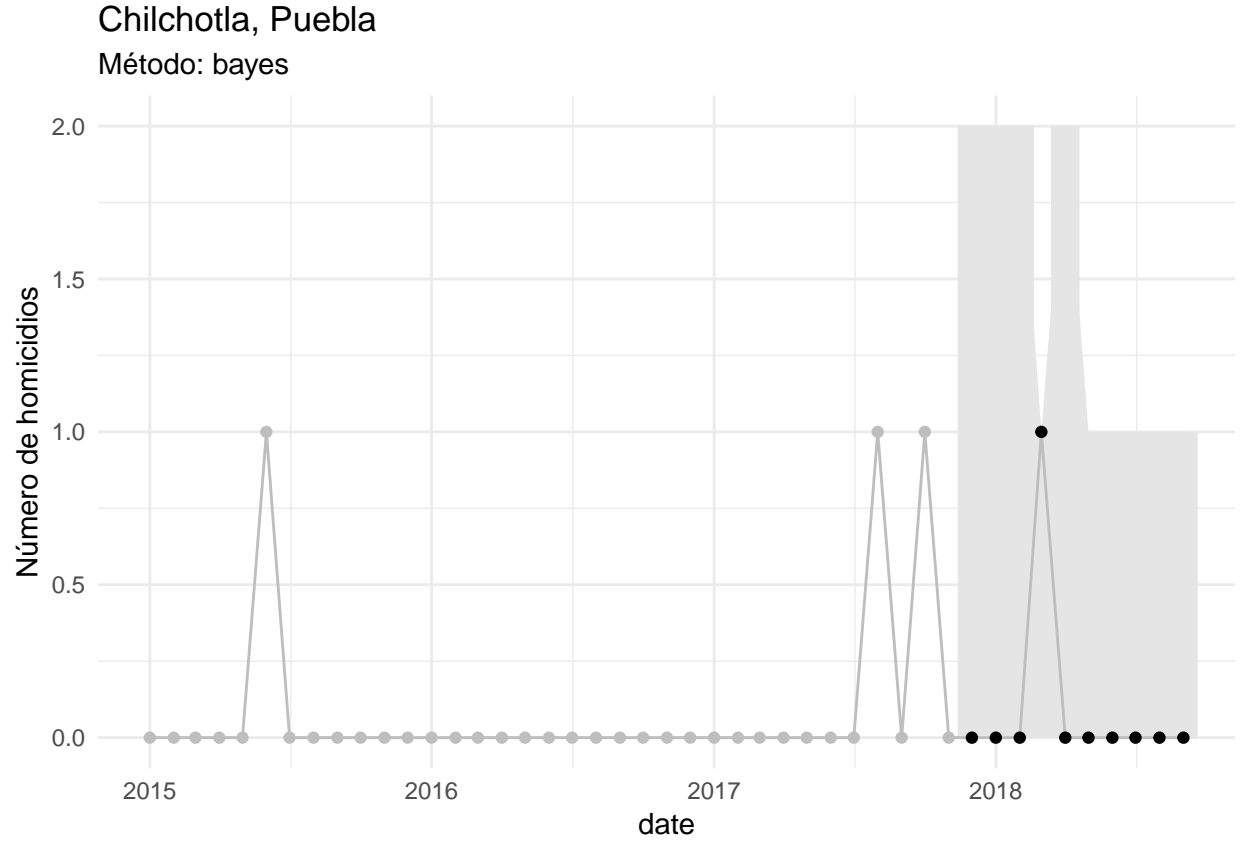
```
monitor_pue <- monitor_bayes(puebla, alpha = 0.05)
graf_monitor(monitor_pue)
```

Puebla, Puebla

Método: bayes



```
monitor_chi <- monitor_bayes(chilchotla, alpha = 0.05)
graf_monitor(monitor_chi)
```



5 Métodos de monitoreo: Farrington

Consideramos ahora el método de Farrington (mejorado), que utiliza un modelo lineal generalizado Poisson:

1. Supongamos que los valores de referencia son y_t
2. Se ajusta un modelo Poisson sobredisperso a los datos de referencia, con

$$E(y_t) = \mu_t, Var(y_t) = \phi \mu_t,$$

donde ϕ es el parámetro de sobredispersión, y donde

$$\log(\mu_k) = \alpha + \beta t$$

sirve para modelar tendencia.

3. Se calculan pesos para los datos, reduciendo el peso de *outliers* según la distribución obtenida de los datos de referencia.
4. Se reajusta el modelo con los pesos calculados en el paso anterior
5. Para el paso de detección, usamos el método de Noufaily (2012) (ver referencia de paquete), donde la distribución de referencia para el nuevo valor y_{t_0} es

$$y_{t_0} \sim NB(\mu_{t_0}, \eta)$$

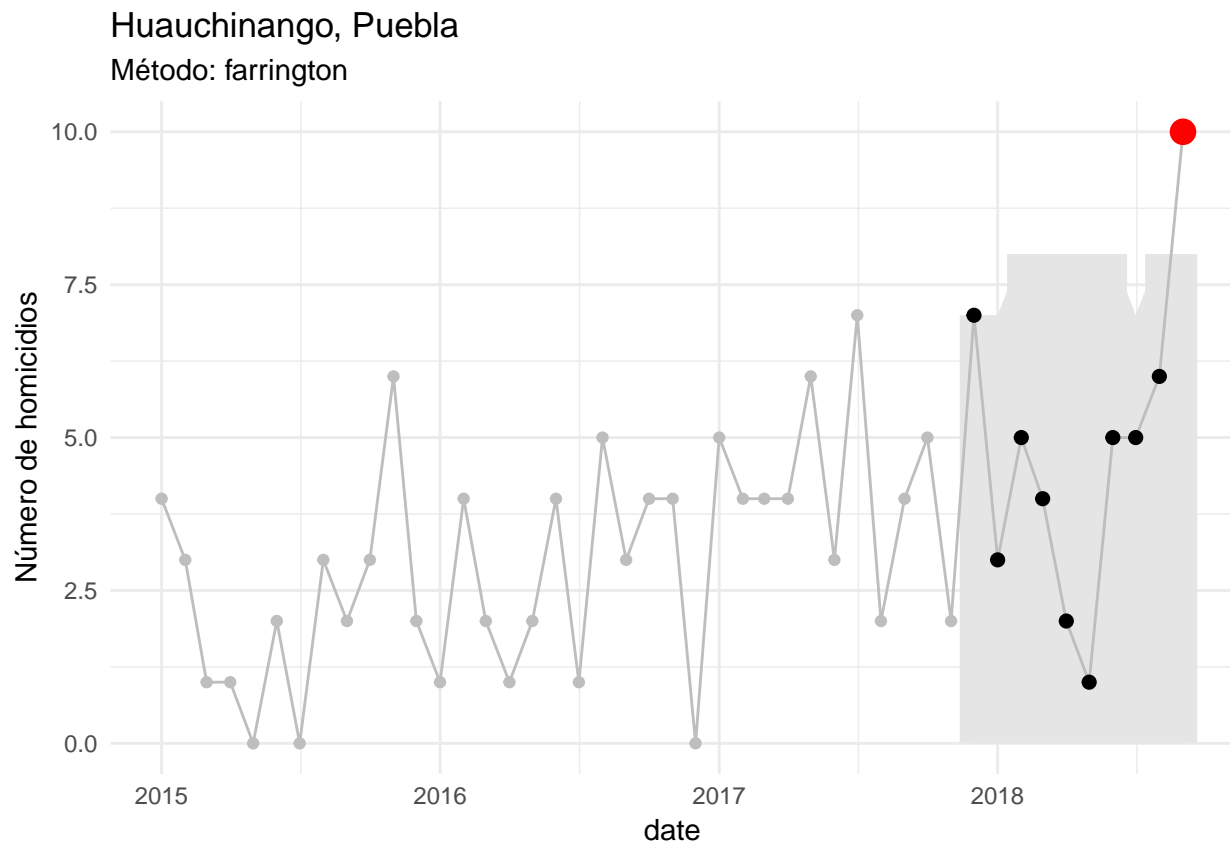
donde $\eta = \frac{\mu_{t_0}}{\phi - 1}$.

Notas:

- Este enfoque ignora la incertidumbre en la estimación de los parámetros μ_{t_0} y ϕ .
 - Para declarar una anomalía, se requieren también al menos 5 casos en los últimos 3 meses (incluyendo el mes que se está evaluando).
6. Utilizamos 4 meses alrededor de del valor que queremos predecir, considerando dos años atrás (para un total de $2 + 2(4) + 1 = 11$ datos de referencia).

5.0.1 Ejemplos: método de Farrington mejorado

```
monitor_huau <- monitor_farrington(huachinango, alpha=0.05, periods = 10)
graf_monitor(monitor_huau)
```



```
monitor_huau$data %>% tail
```

```
## # A tibble: 6 x 10
## # Groups:   state_code, state, mun_code, municipio [1]
##   state_code state mun_code municipio date      population count alerta
##   <int> <chr> <int> <chr> <date>      <int> <dbl> <lgl>
## 1      21 PUEB~      71 HUAUCHIN~ 2018-04-01    107811      2 FALSE
## 2      21 PUEB~      71 HUAUCHIN~ 2018-05-01    107898      1 FALSE
## 3      21 PUEB~      71 HUAUCHIN~ 2018-06-01    107986      5 FALSE
## 4      21 PUEB~      71 HUAUCHIN~ 2018-07-01    108072      5 FALSE
## 5      21 PUEB~      71 HUAUCHIN~ 2018-08-01    108158      6 FALSE
## 6      21 PUEB~      71 HUAUCHIN~ 2018-09-01    108244     10 TRUE
## # ... with 2 more variables: alerta_nivel <dbl>, metadata <chr>
```



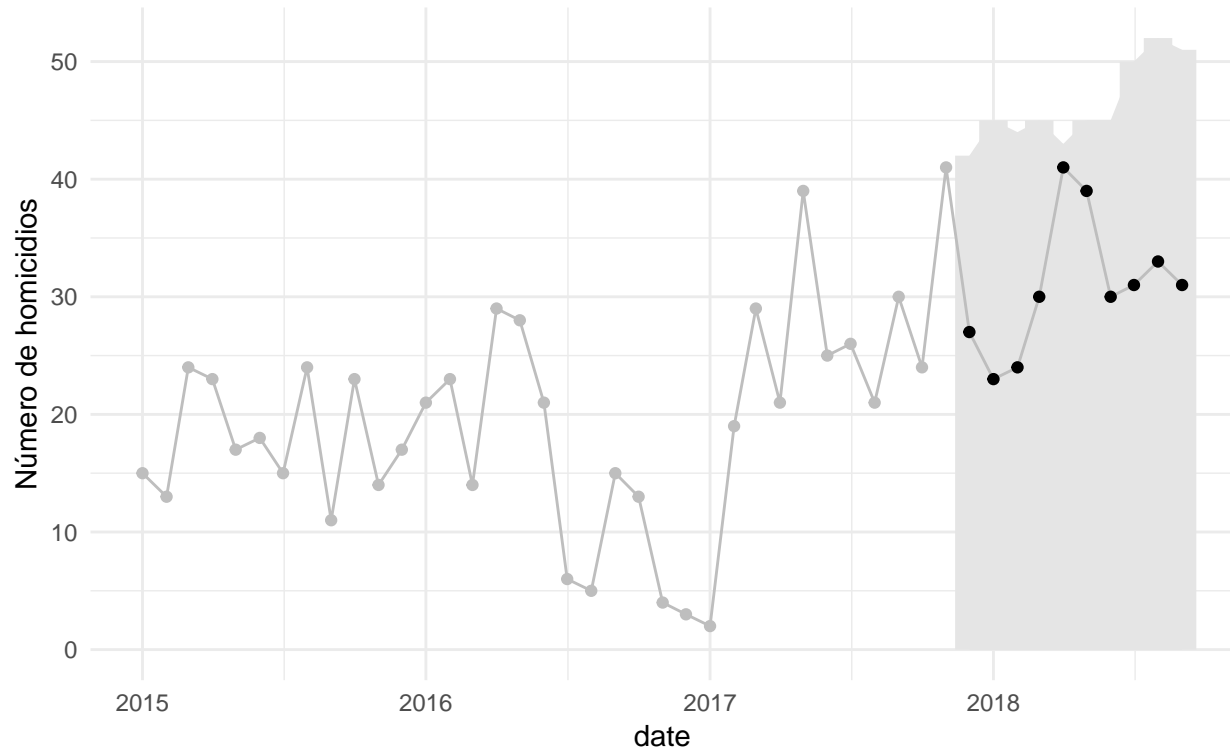
```
monitor_pue <- monitor_farrington(puebla, alpha=0.05, periods = 10)
graf_monitor(monitor_pue)
```

```
## Warning: Removed 35 rows containing missing values (geom_linerange).
```

```
## Warning: Removed 35 rows containing missing values (geom_point).
```

Puebla, Puebla

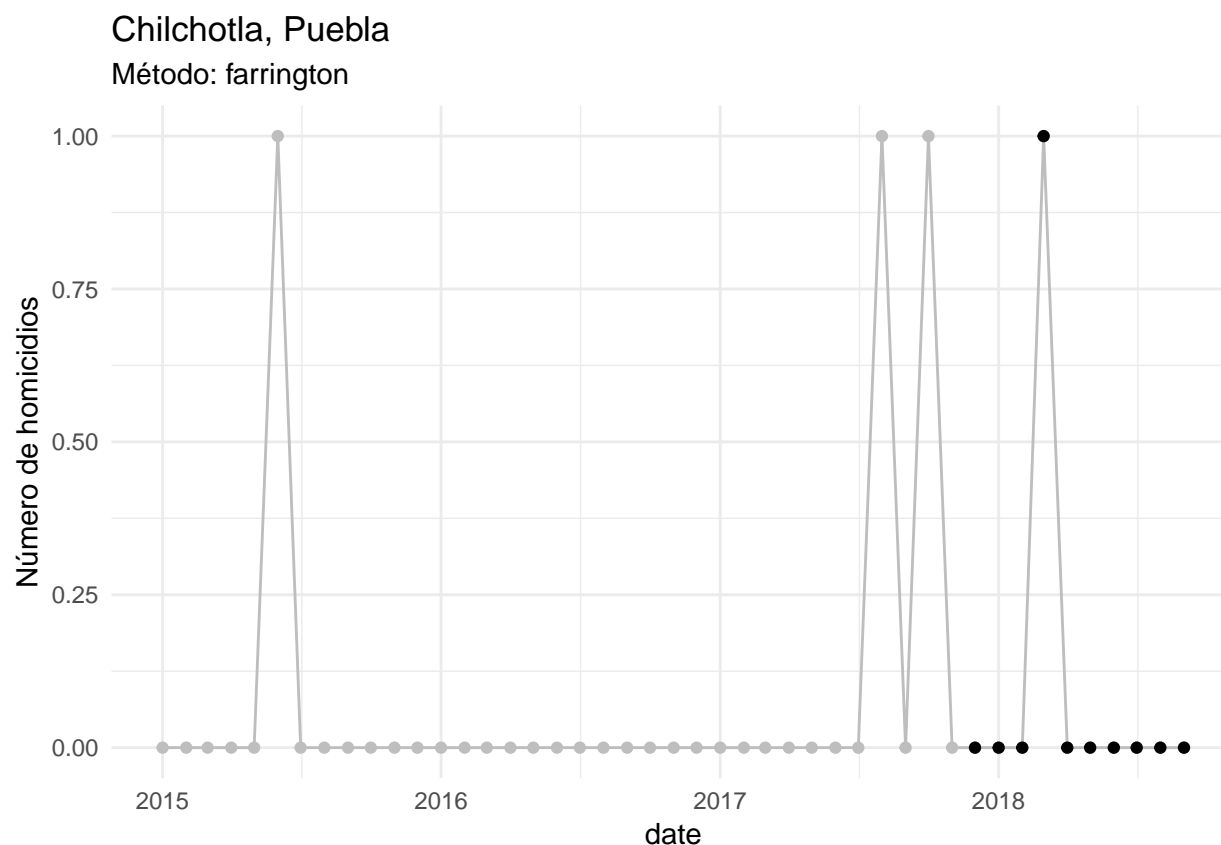
Método: farrington



```
monitor_puebla <- monitor_farrington(puebla, periods = 10)
```

Nótese que las bandas en la estimación de Farrington puede colapsarse para municipios con niveles bajos de homicidios:

```
monitor_chi <- monitor_farrington(chilchotla, alpha = 0.05, periods = 10)
graf_monitor(monitor_chi)
```



6 Otros ejemplos

En Sinaloa, por ejemplo, no observamos ninguna alarma reciente. Aún cuando algunas tasas de homicidio son muy altas en todo el periodo de referencia/evaluación, no se detectan brotes recientes:

```
#sinaloa_edo <- filter(homicidios, state == "SINALOA")
sinaloa_edo <- readRDS("data/sinaloa.rds")
sinaloa_edo %>% group_by(municipio) %>% summarise(count = sum(count)) %>%
  arrange(desc(count))
```

```
## # A tibble: 18 x 2
##   municipio      count
##   <chr>         <dbl>
## 1 CULIACÁN      2449
## 2 MAZATLÁN      763
## 3 AHOME         544
## 4 GUASAVE       484
## 5 NAVOLATO      404
## 6 EL FUERTE     215
## 7 MOCORITO      180
## 8 SALVADOR ALVARADO 150
## 9 ROSARIO       132
## 10 CONCORDIA     128
## 11 SINALOA       127
## 12 BADIRAGUATO   115
## 13 ELOTA         106
```

```

## 14 ESCUINAPA          97
## 15 SAN IGNACIO        87
## 16 ANGOSTURA          83
## 17 CHOIX              65
## 18 COSALÁ            26

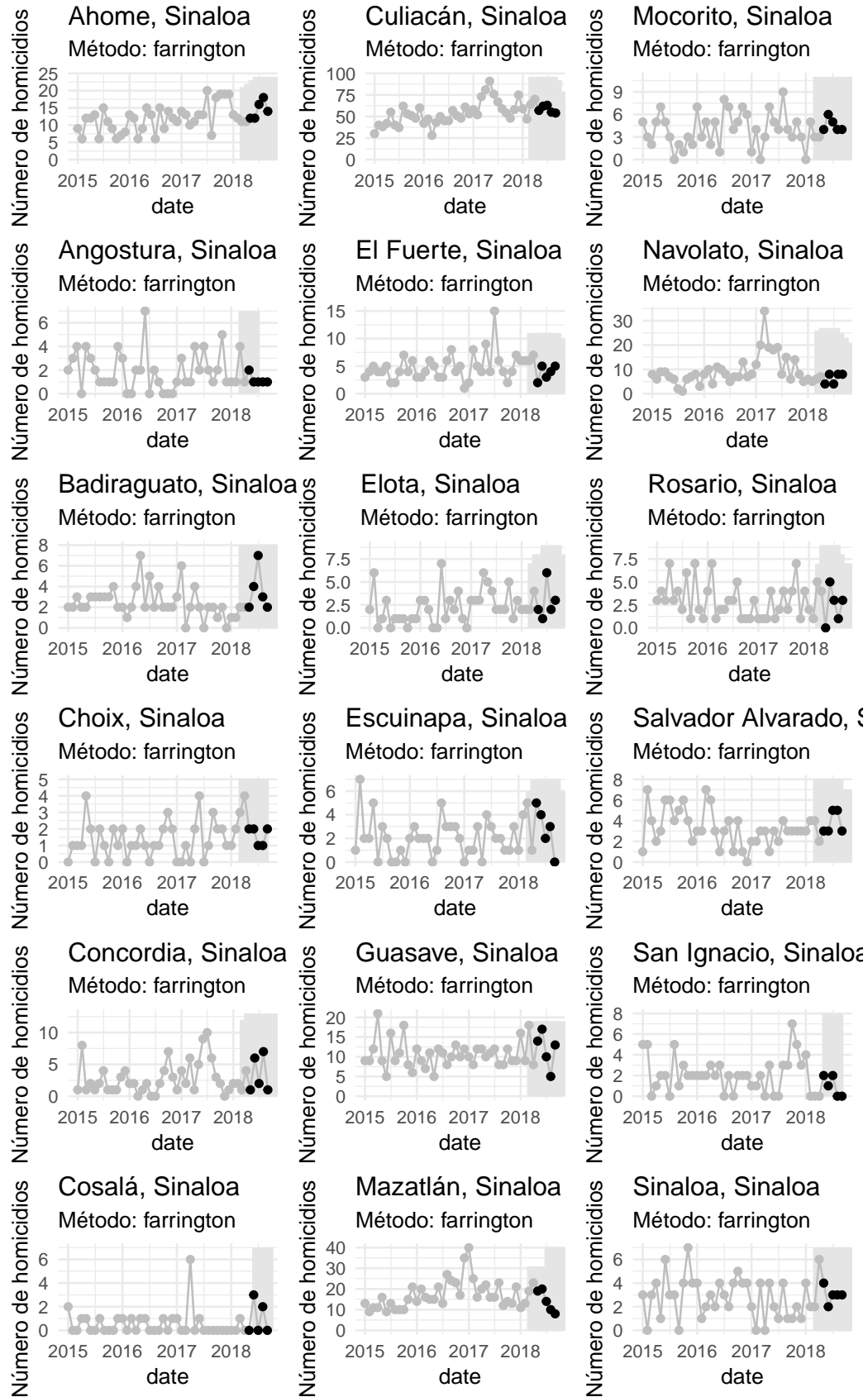
dat_sinaloa <- sinaloa_edo %>% split(.$municipio)
monitores_sinaloa <-
  lapply(dat_sinaloa, function(df){
    monitor_mun <- monitor_farrington(df, alpha = 0.02, periods = 5)
    g <- graf_monitor(monitor_mun)
    g
  })

library(gridExtra)

##
## Attaching package: 'gridExtra'
## The following object is masked from 'package:dplyr':
##
##      combine

marrangeGrob(monitores_sinaloa, nrow = 6, ncol = 3)

```



En los siguientes ejemplos sólo mostramos los municipios con detección de anomalías:

```
mex_edo <- filter(homicidios, state == "MICHOACÁN")
#sinaloa_edo <- readRDS("data/sinaloa.rds")
mex_edo %>% group_by(municipio) %>% summarise(count = sum(count)) %>%
  arrange(desc(count))
```

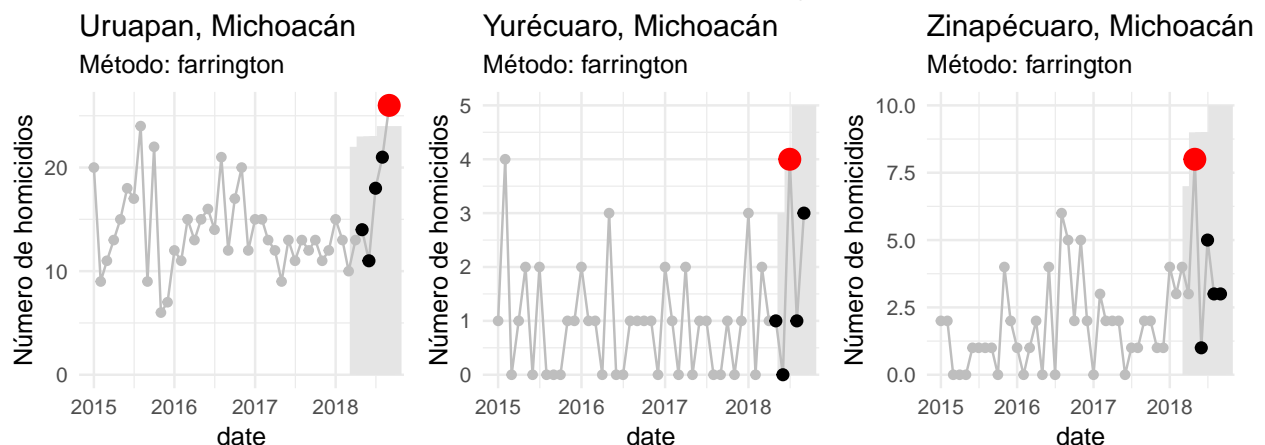
```
## # A tibble: 113 x 2
##   municipio      count
##   <chr>         <dbl>
## 1 MORELIA       1490
## 2 URUAPAN       639
## 3 ZAMORA        504
## 4 APATZINGÁN    471
## 5 LÁZARO CÁRDENAS 437
## 6 LA PIEDAD     266
## 7 SAHUAYO       236
## 8 MÚGICA        205
## 9 ZITÁCUARO     179
## 10 BUENAVISTA   165
## # ... with 103 more rows
```

```
dat_mex <- mex_edo %>% split(.$municipio)
monitores_mex <-
  lapply(dat_mex, function(df){
    monitor_mun <- monitor_farrington(df, alpha = 0.02, periods = 5)
    g <- graf_monitor(monitor_mun)
    list(graf = g, monitor = monitor_mun)
  })
```

```
tiene_alerta <- function(monitor){
  any(monitor$monitor$data$alerta == 1, na.rm = TRUE)
}
monitores_alarma <- keep(monitores_mex, tiene_alerta)
```

```
graficas <- monitores_alarma %>% map(~.$graf)
marrangeGrob(graficas, nrow = 1, ncol = 4)
```

page 1 of 1



```
mex_edo <- filter(homicidios, state == "VERACRUZ")
#sinaloa_edo <- readRDS("data/sinaloa.rds")
```

```
mex_edo %>% group_by(municipio) %>% summarise(count = sum(count)) %>%
  arrange(desc(count))
```

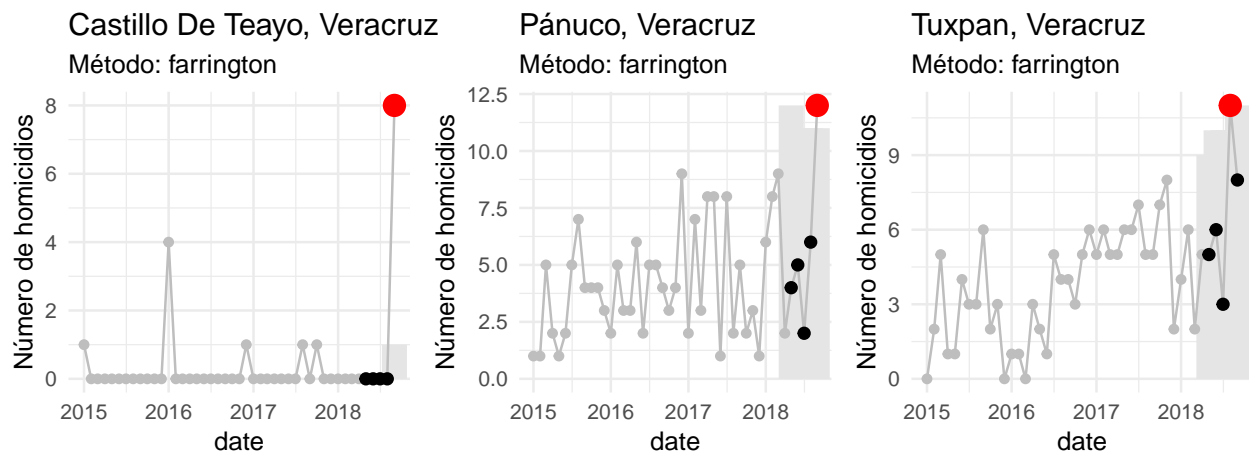
```
## # A tibble: 214 x 2
##   municipio      count
##   <chr>         <dbl>
## 1 VERACRUZ      471
## 2 CORDOBA       362
## 3 XALAPA        326
## 4 COATZACOALCOS 320
## 5 POZA RICA DE HIDALGO 286
## 6 PAPANTLA      227
## 7 ACAYUCAN      195
## 8 PÁNUCO        194
## 9 TUXPAN        182
## 10 TIERRA BLANCA 177
## # ... with 204 more rows
```

```
dat_mex <- mex_edo %>% split(.$municipio)
monitores_mex <-
  lapply(dat_mex, function(df){
    monitor_mun <- monitor_farrington(df, alpha = 0.02, periods = 5)
    g <- graf_monitor(monitor_mun)
    list(graf = g, monitor = monitor_mun)
  })
```

```
monitores_alarma <- keep(monitores_mex, tiene_alerta)
```

```
graficas <- monitores_alarma %>% map(~.$graf)
marrangeGrob(graficas, nrow = 1, ncol = 4)
```

page 1 of 1



7 Propuesta: Método mixto

Proponemos un método mixto para construir rango y alertas para todos los municipios:

- Aplicamos primero el método Farrington flexible. Si el algoritmo converge, y hay suficientes datos (dos años atrás), utilizamos su salida.
- En otro caso, utilizamos el método bayesiano simple aplicado a los últimos 6 meses de datos.

Si hay suficientes datos hacia atrás, Farrington puede fallar cuando observamos conteos muy bajos (proporción alta de ceros).

Típicamente, estos métodos se evalúan con etiquetado de brotes según expertos en epidemiología ((Yang et al. 2018)). En nuestro caso, consideramos propiedades básicas de la estimación y comportamiento con los datos observados, pero los métodos deberán ser evaluados por expertos en el tema que puedan etiquetar cuáles saltos indican brotes de violencia o no.

7.1 Método mixto para todos los municipios.

Al aplicar para detección del último mes con los datos actuales obtenemos la siguiente proporción de uso de Farrington vs Bayesiano simple:

```
dat_nal <- homicidios %>% split(list(. $municipio, . $state), drop = TRUE)
length(dat_nal)
```

```
## [1] 2469
```

```
monitores_nal <-
  lapply(dat_nal, function(df){
    # farrington por default, a menos que haya menos de 2 años de datos
    # o no converja.
    monitor_mun <- monitor(df, alpha = 0.05, periods = 1)
    #g <- graf_monitor(monitor_mun)
    #list(graf = g, monitor = monitor_mun)
    monitor_mun
  })
tipos <- sapply(monitores_nal, function(monitor){ monitor $tipo })
table(tipos)
```

```
## tipos
##      bayes farrington
##      561      1908
```

Y en la siguiente tabla mostramos que cuando hay menos de un año de datos, se utiliza siempre bayes simple. El de Farrington se utiliza en la mayor parte de los casos cuando hay más de un año de datos, excepto en algunos casos de conteo bajo donde no converge apropiadamente.

```
conteos_bajo <- sapply(monitores_nal, function(monitor){ sum(tail(monitor $data $count,6)) < 5})
menos_1_año <- sapply(monitores_nal, function(monitor){ length(monitor $data $count) < 24})
data_frame(conteo_bajo = conteos_bajo, menos_1_año = menos_1_año, tipo = tipos) %>%
  group_by(menos_1_año, conteo_bajo, tipo) %>% tally() %>% spread(tipo, n)
```

```
## Warning: `data_frame()` is deprecated, use `tibble()`.
## This warning is displayed once per session.
```

```
## # A tibble: 4 x 4
## # Groups:   menos_1_año, conteo_bajo [4]
##   menos_1_año conteo_bajo bayes farrington
##   <lg1>      <lg1>      <int>      <int>
## 1 FALSE      FALSE         NA         675
## 2 FALSE      TRUE         23        1233
## 3 TRUE       FALSE         14          NA
## 4 TRUE       TRUE         524          NA
```

El número de alertas generadas para este último mes, los municipios afectados:

```

alertas <- sapply(monitores_nal, function(monitor){
  any(monitor$data$alerta, na.rm = TRUE) }
)
table(alertas)

```

```

## alertas
## FALSE TRUE
## 2440 29

```

```

alertas_mun <- which(alertas)
alertas_mun

```

```

##          MEOQUI.CHIHUAHUA          CELAYA.GUANAJUATO
##                200                332
##      IRAPUATO.GUANAJUATO      LEÓN.GUANAJUATO
##                341                344
##      PÉNJAMO.GUANAJUATO      SALAMANCA.GUANAJUATO
##                348                352
##    TLAPA DE COMONFORT.GUERRERO      EL SALTO.JALISCO
##                445                572
##      GUADALAJARA.JALISCO      TALA.JALISCO
##                577                623
##    IXTAPAN DE LA SAL.MÉXICO      TENANGO DEL VALLE.MÉXICO
##                705                755
##      TEXCOCO.MÉXICO      URUAPAN.MICHOACÁN
##                764                891
##      TEMIXCO.MORELOS      TECUALA.NAYARIT
##                919                951
##      GARCÍA.NUEVO LEÓN      JUÁREZ.NUEVO LEÓN
##                971                983
##    MONTERREY.NUEVO LEÓN SANTO DOMINGO TEHUANTEPEC.OAXACA
##                994                1521
##    TLALIXTAC DE CABRERA.OAXACA      HUAUCHINANGO.PUEBLA
##                1554                1647
##    SAN MARTÍN TEXMELUCAN.PUEBLA      CAJEME.SONORA
##                1707                1916
##    CASTILLO DE TEAYO.VERACRUZ      EMILIANO ZAPATA.VERACRUZ
##                2123                2157
##    LA ANTIGUA.VERACRUZ      PÁNUCO.VERACRUZ
##                2187                2223
##    LORETO.ZACATECAS
##                2435

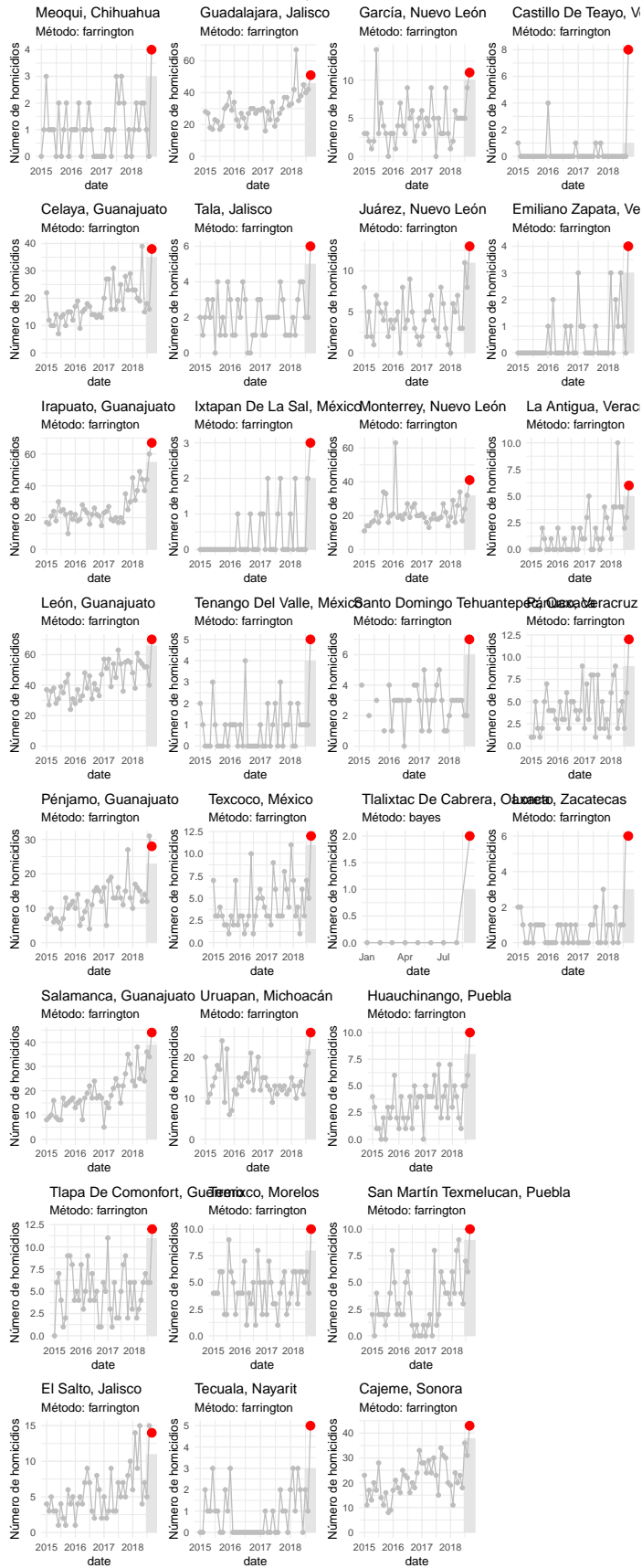
```

Finalmente, mostramos los datos detectados para todos los municipios con alerta:

```

monitores_alerta <- monitores_nal[alertas_mun]
graficas <- monitores_alerta %>% map(graf_monitor)
marrangeGrob(graficas, nrow = 8, ncol = 4)

```

References

- “Datos Abiertos de Incidencia Delictiva.” 2018a. <https://www.gob.mx/sesnsp/acciones-y-programas/datos-abiertos-de-incidencia-delictiva?state=published>.
- “Datos Abiertos de Incidencia Delictiva.” 2018b. <https://github.com/diegovalle/new.crimenmexico>.
- “Encuesta Nacional de Victimización y Percepción sobre Seguridad Pública (ENVIPE) 2018.” n.d.
- Farrington, Paddy, and Nick Andrews. 2003. “Outbreak Detection: Application to Infectious Disease Surveillance.” In *Monitoring the Health of Populations: Statistical Principles and Methods for Public Health Surveillance*, edited by Ron Brookmeyer and Donna F. Stroup, 203–31. New York, NY, USA: OUP USA. <http://oro.open.ac.uk/22646/>.
- Höhle, Michael. 2007. “Surveillance: An R Package for the Monitoring of Infectious Diseases.” *Computational Statistics* 22 (4): 571–82. <https://doi.org/10.1007/s00180-007-0074-8>.
- Meyer, Sebastian, Leonhard Held, and Michael Höhle. 2017. “Spatio-Temporal Analysis of Epidemic Phenomena Using the R Package surveillance.” *Journal of Statistical Software* 77 (11): 1–55. <https://doi.org/10.18637/jss.v077.i11>.
- Yang, Eunjoo, Hyun Woo Park, Yeon Hwa Choi, Jusim Kim, Lkhagvadorj Munkhdalai, Ibrahim Hussein Musa, and Keun Ho Ryu. 2018. “A Simulation-Based Study on the Comparison of Statistical and Time Series Forecasting Methods for Early Detection of Infectious Disease Outbreaks.” In *International Journal of Environmental Research and Public Health*.