

# Cátedra INEGI

Estado de ánimo de los tuiteros en México.



INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA



Universidad Veracruzana



Big Data??

Big data is like teenage sex:  
everyone talks about it,  
nobody really knows how to do it,  
everyone thinks everyone else is  
doing it, so everyone claims they  
are doing it...

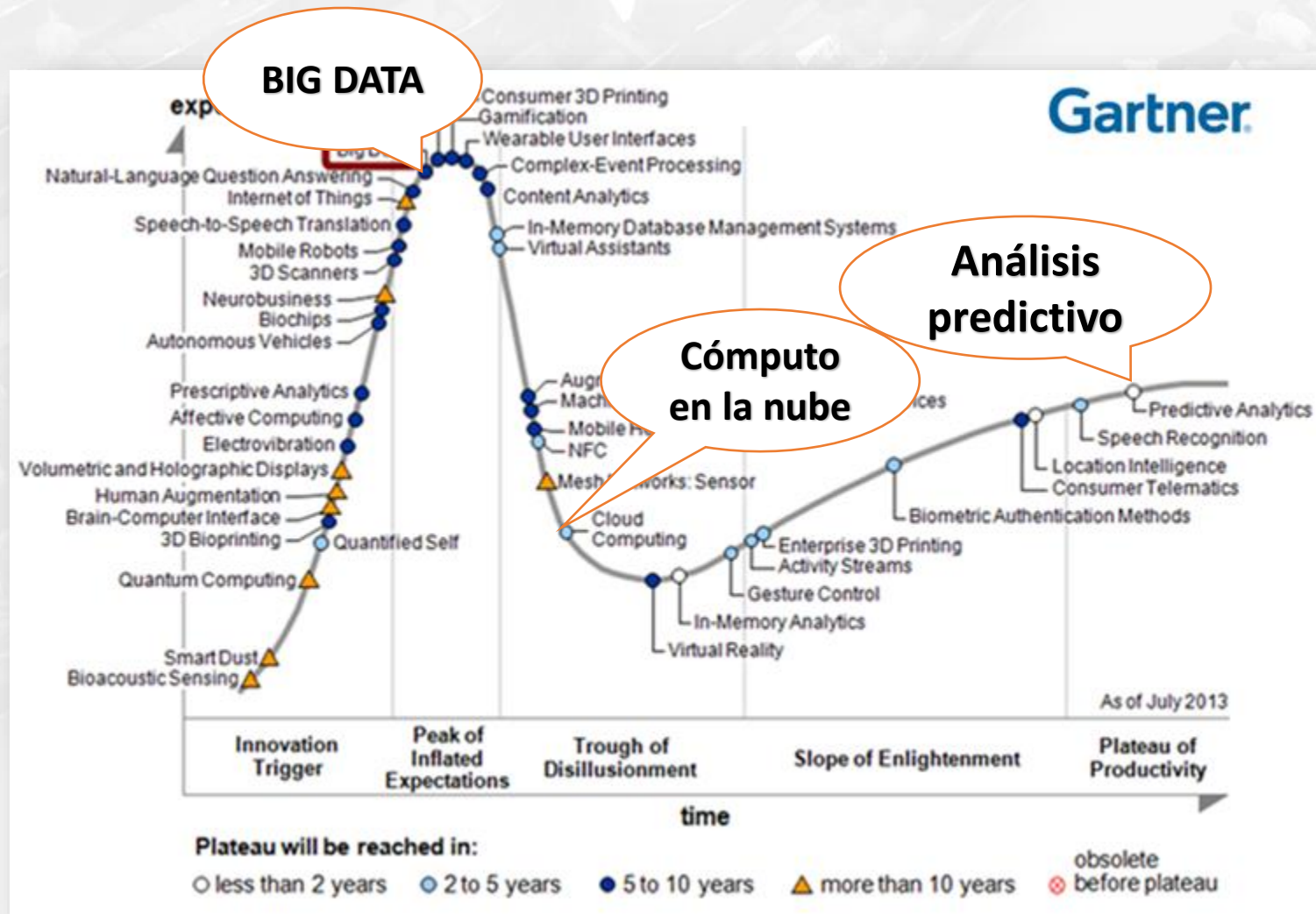
(Dun & Amy)

# EL PARADIGMA BIG DATA

- Dejar que los datos hablen.
- N=todo (¿Mucho?).
- Granularidad constante.
- No hay diseño, ni pregunta, ni muestreo pues los datos no son escasos.
- Correlación más que causalidad.
- Menos exactitud y mucha paja o ruido.

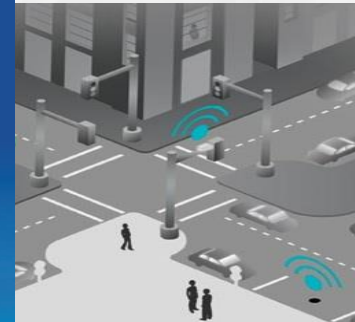


# Posicionamiento de Big Data en el Ciclo de las Expectativas de Gartner





# Las fuentes de información están creciendo y evolucionando...

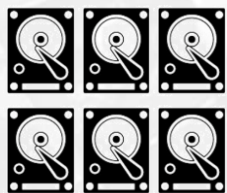


# CONOCIMIENTO

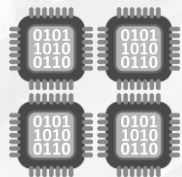


Expertos marcan el rumbo.

## Infraestructura Informática



Sistemas Distribuidos



Computo Paralelo  
y Concurrente

## Estadística Matemática

Análisis de Datos

Statistical Learning

Minería de Datos

Estratificaciones

**Y mucho más...**

Muestreo

Análisis de Redes (Grafos)

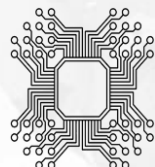
Machine Learning

Análisis de Regresión

## Internet de las cosas



Sensores



# BIG DATA:

## Internet de las Personas



Redes Sociales





# QUÉ SIGNIFICA “BIG DATA” PARA LA INFORMACIÓN OFICIAL?

Cátedra  
**INEGI**



# BIG DATA Y ESTADÍSTICA OFICIAL





# Posibles fuentes de “Big Data”

- Administrativas (gubernamentales o no);
- Comerciales o transaccionales;
- Sensores;
- De desplazamiento;
- De comportamiento;
- Opinión;

- Legislativos,
- Privacidad,
- Financieros,
- Administrativos,
- Metodológicos,
- Tecnológicos.



## 3 ámbitos abiertos a experimentación:

- Corto plazo: Combinación de Big data con estadísticas oficiales.
- Mediano plazo: Producción de información relativa a temas emergentes basada en Big Data (en particular, cuando no hay encuesta ni registro).
- Largo plazo: Reemplazo de parte de la estadística oficial producida mediante Big Data.

Qué estamos haciendo en  
INEGI?

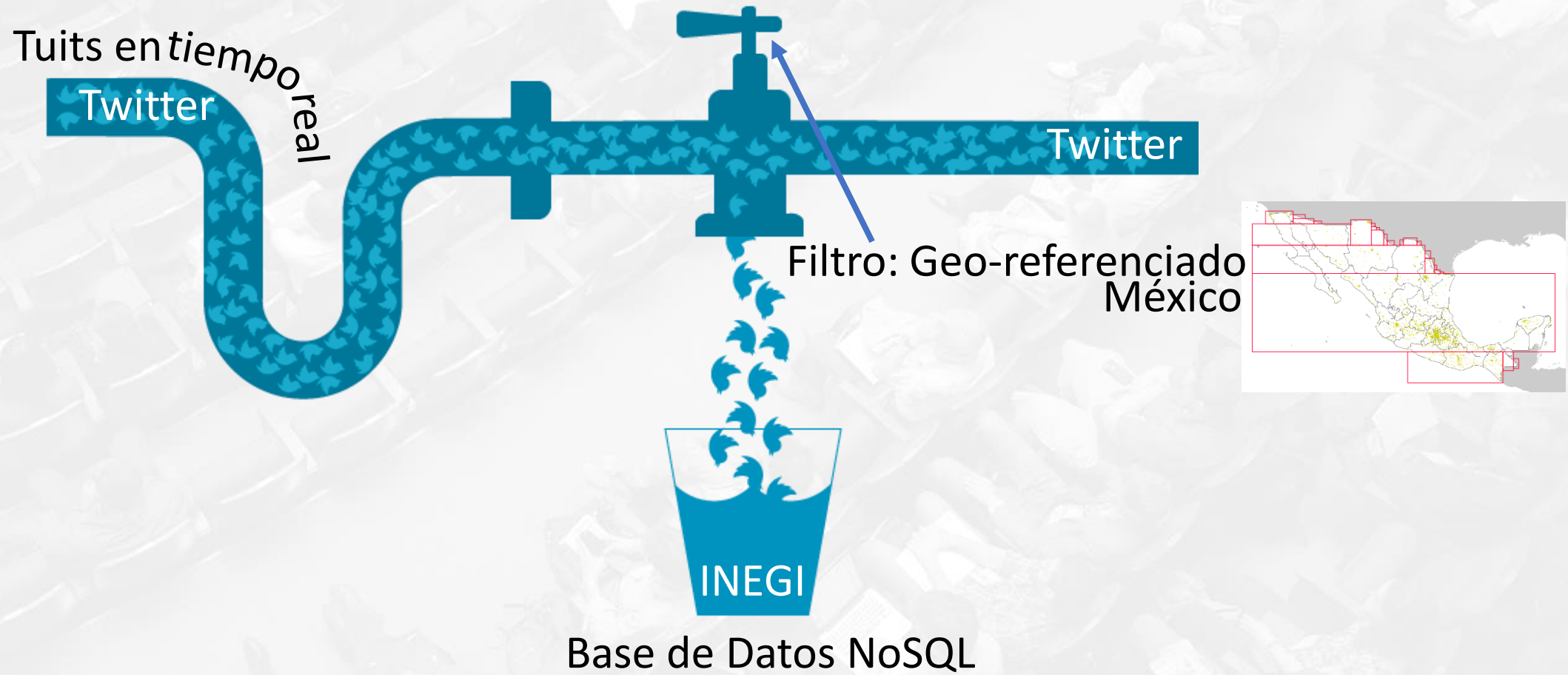
**PROYECTO PILOTO:  
EXPLOTACIÓN DE TWITTER**

Cátedra  
**INEGI**





# Twitter como fuente de datos



# Por qué Twitter?



- Disponible fácilmente
- Hasta el 1% de los tweets mundiales sin costo
- Alrededor de 12 M de cuentas en México
- 700 mil cuentas con tuits georreferenciados en MX
- 100 millones de tuits descargados desde enero de 2014



# Infraestructura para recolectar tuits

Prueba de concepto



Cluster (Hydra)



NoSql Database "Elasticsearch"

Unix "Red Hat"

Big Data Layers







# TEMAS DE ESTUDIO

- Bienestar Subjetivo a partir de Redes Sociales
- Flujos Turísticos
- Movilidad
- Elecciones 2018

# TWITTER PARA BIENESTAR SUBJETIVO

Cátedra  
**INEGI**





## Proceso de análisis de tuits:

- Se utiliza un **Método supervisado de clasificación**:
  - Humanos califican el sentimiento y clasifican el tema de un conjunto de tuits (conjunto de entrenamiento)
  - El conjunto de entrenamiento se utiliza para que la máquina “aprenda” a calificar y clasificar por similitudes nuevos tuits





INSTITUTO NACIONAL  
DE ESTADÍSTICA Y GEOGRAFÍA

<http://cienciadedatos.inegi.org.mx/pioanalysis>

API para  
desarrollar  
conjunto de  
entrenamiento



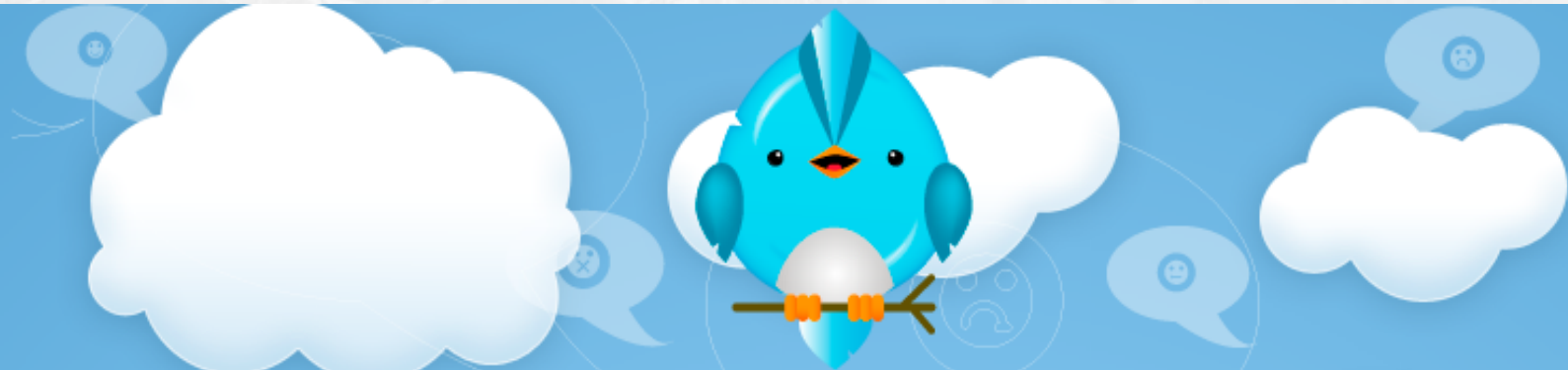
Califica sentimiento de  
tuits en positivo, negativo o  
neutro, y clasifícalos en  
varios temas.

**PIOANALISIS**

Universidad Tec Milenio  
apoya con estudiantes en  
muchos estados del país para  
producir el conjunto de  
entrenamiento

**Bienestar Subjetivo y las redes sociales como fuente de datos**





# Que pague con diamantes su pecado

¿El tema del tuit  es?

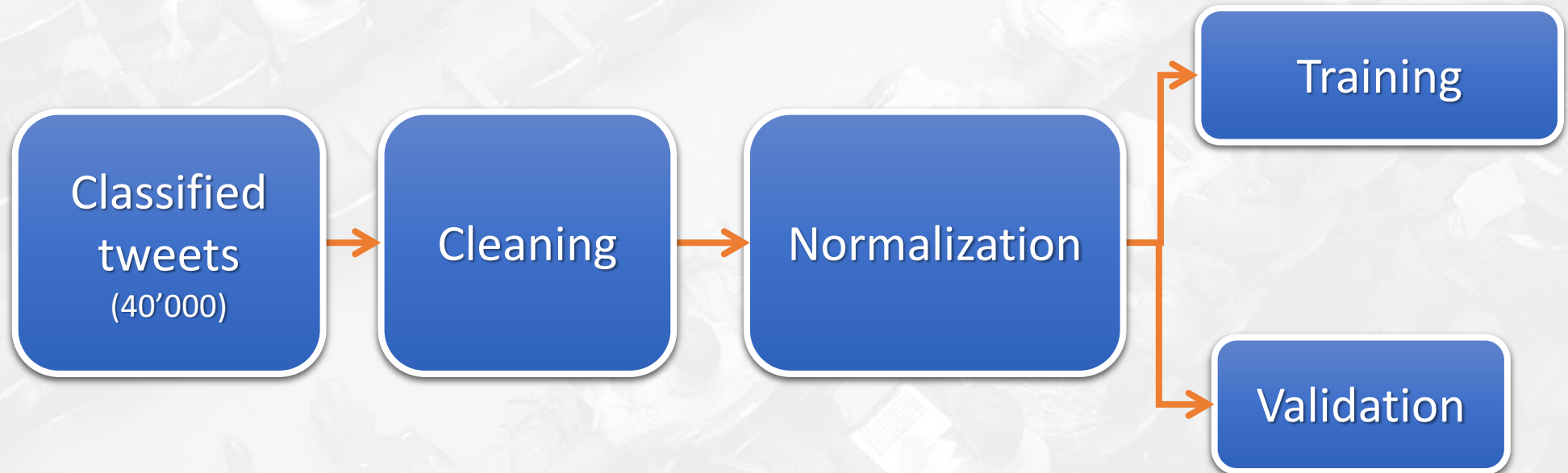
- Política
- Cultural / Entretenimiento
- Deporte
- Escolar / Laboral
- Personal
- Ni idea

¿El tuitero se sentia?

Four icons in a row: a yellow smiley face, a grey neutral face, a red sad face with a white box labeled "Negativo" above it, and a black question mark.

# Colaboración con la academia para clasificar

La información clasificada se divide en dos conjuntos como sigue:

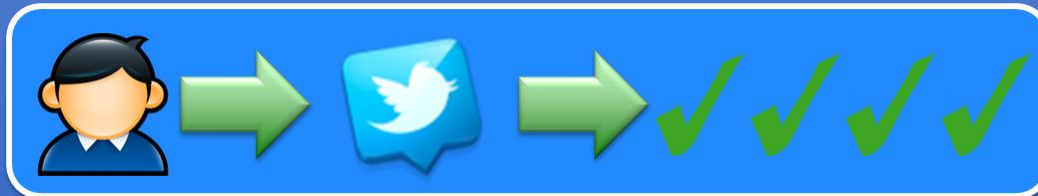




# Depuración de cada conjunto clasificado (Limpieza)

Classified tweets  
(40,136)

Cleaning  
Contradictions and repetitions



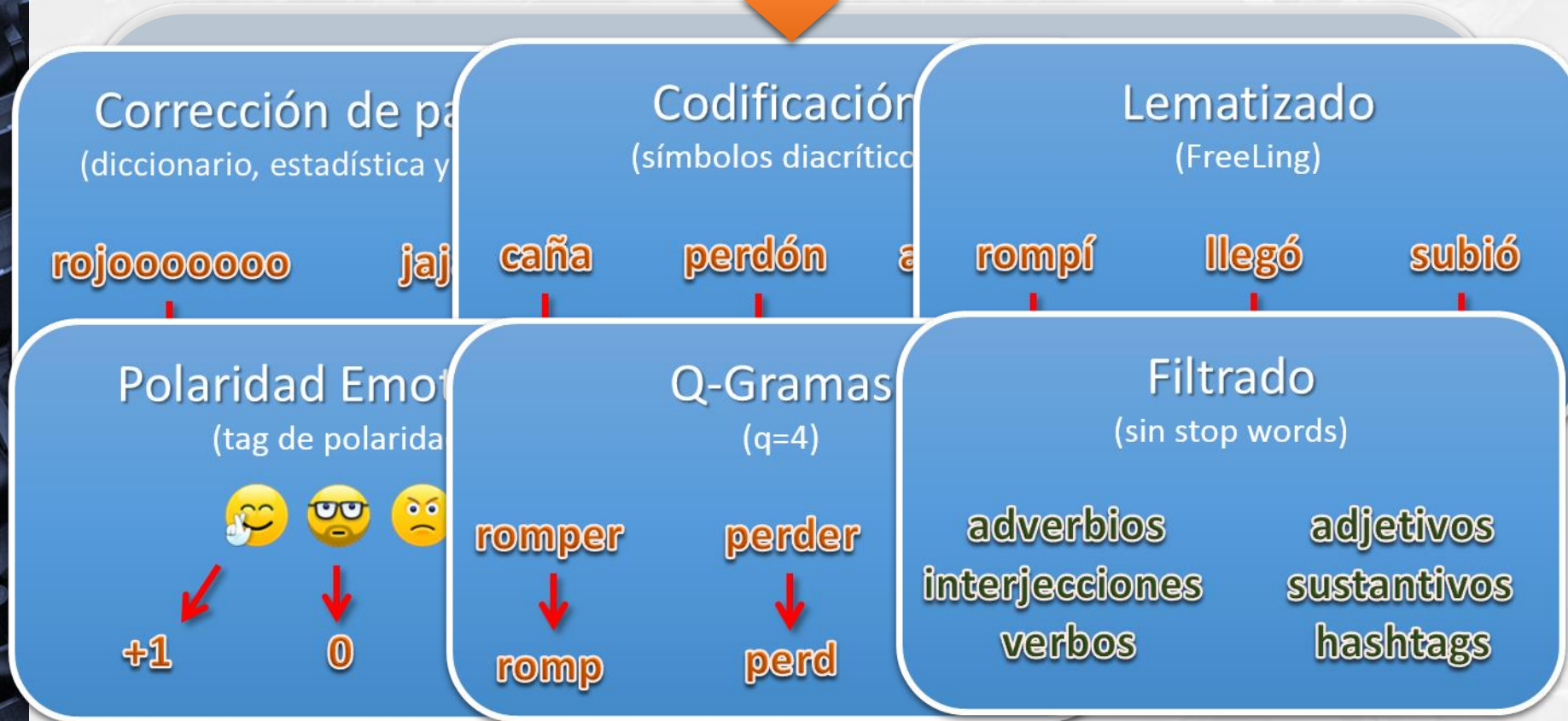
Limpieza  
entropía



# Pre-procesamiento de cada conjunto clasificado (Normalización)



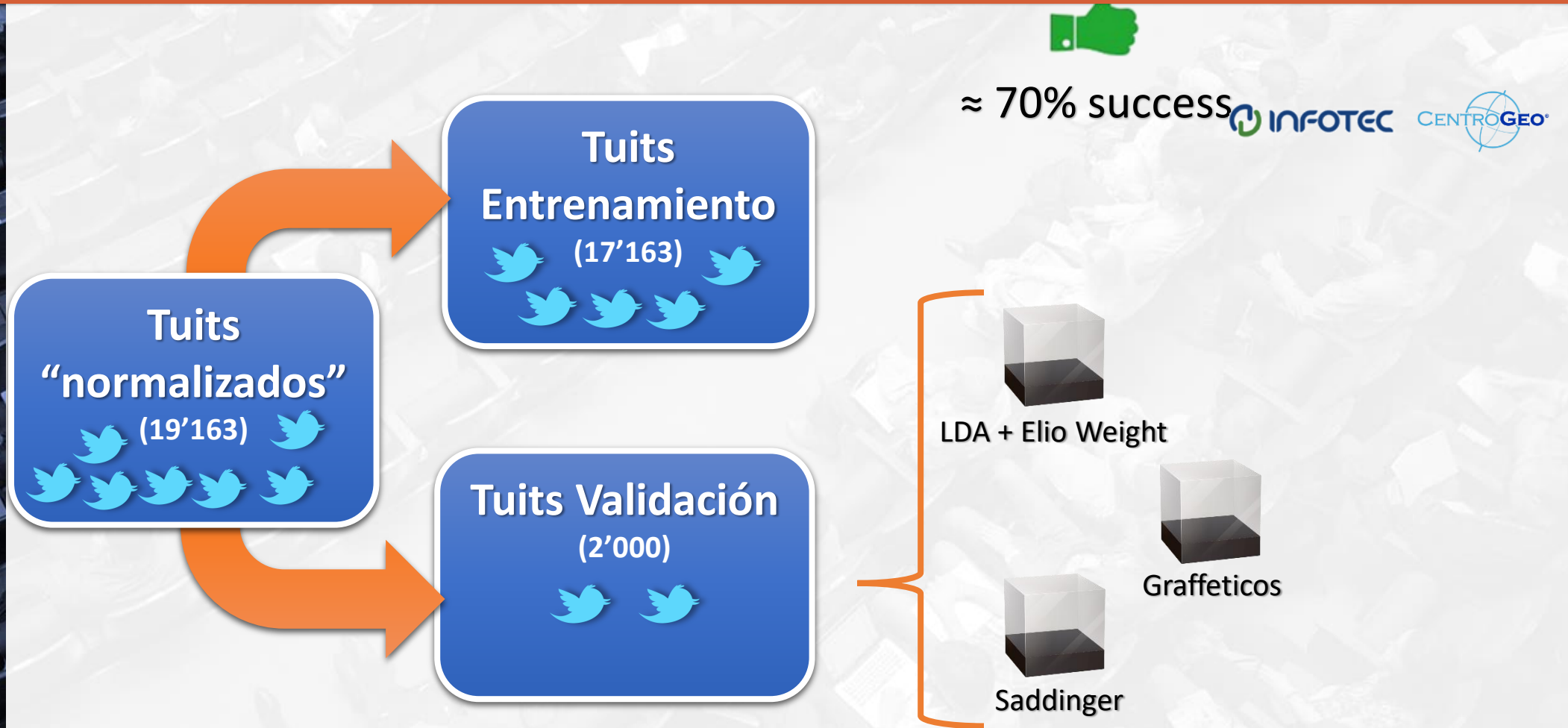
Tuits Limpios



Tuits Normalizados



# Resultados de clasificación



# Mejorar precisión ensamble de clasificadores innovadores



≈ 80% de  
acierto

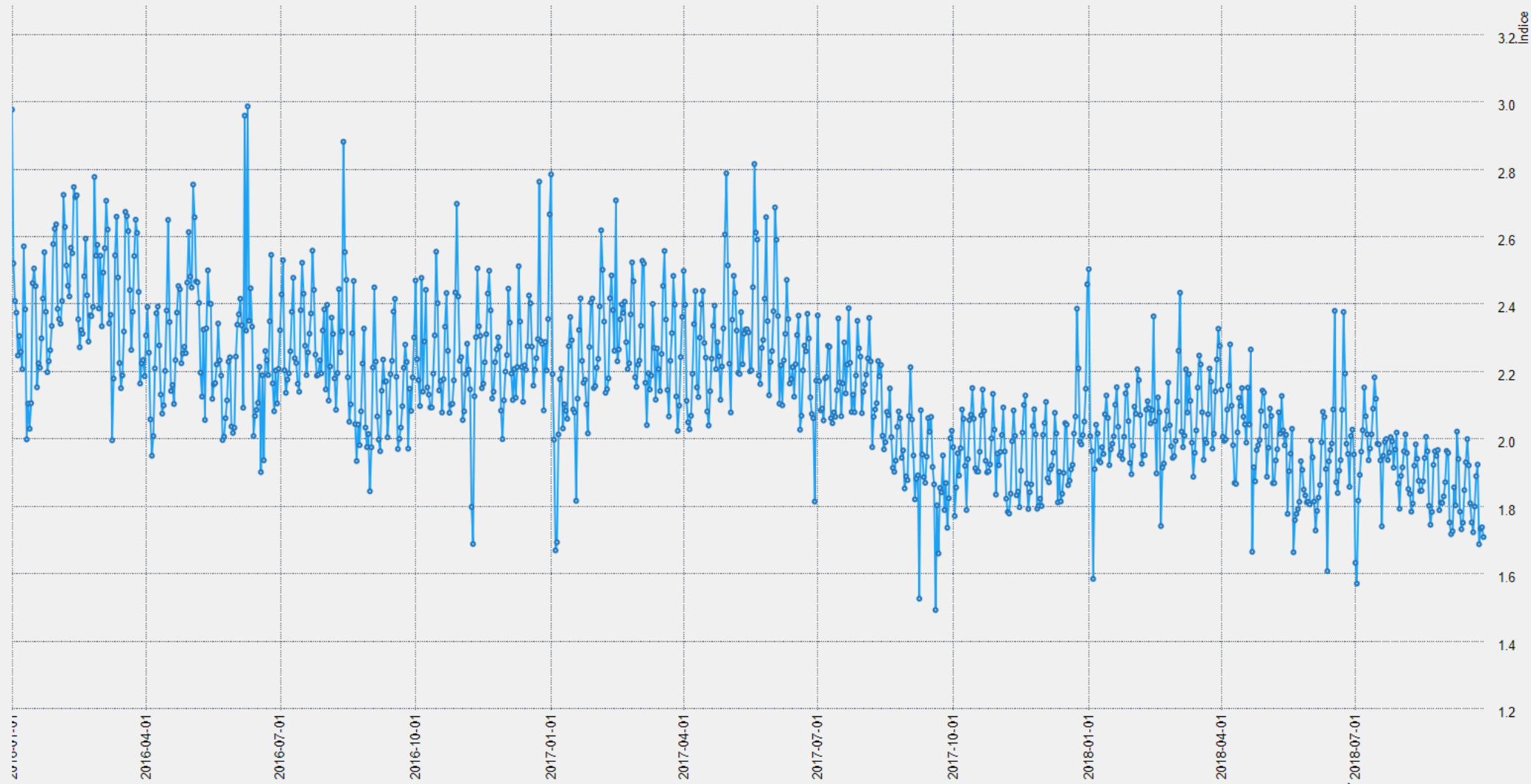
Positivo

Negativo



# Estado de ánimo de los tuiteros en México

1 de enero de 2016 - 26 de septiembre de 2018

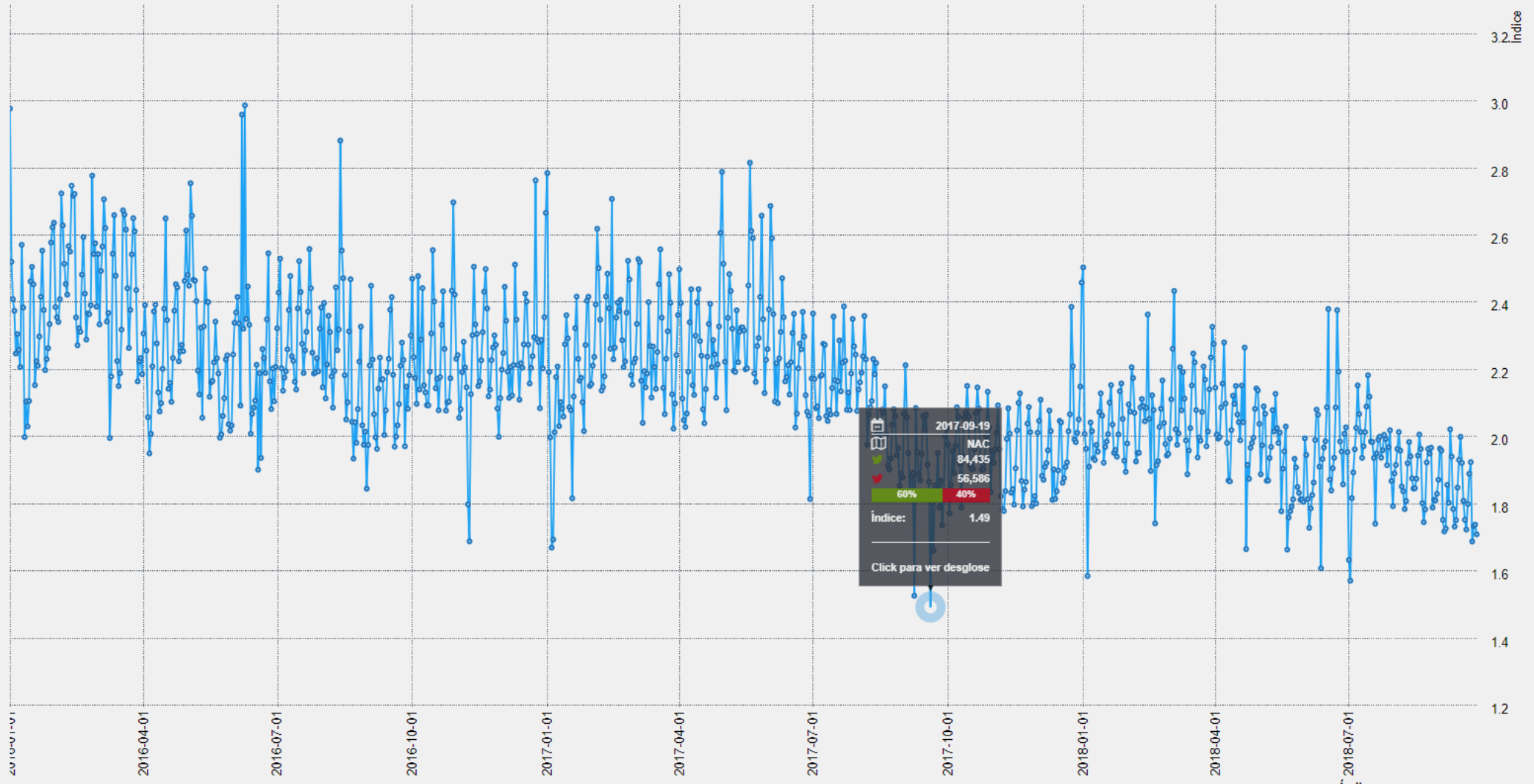


?  
  
  
CSV 

Índice = Positivos (👍) / Negativos (👎)

# Estado de ánimo de los tuitos en México

1 de enero de 2016 - 26 de septiembre de 2018



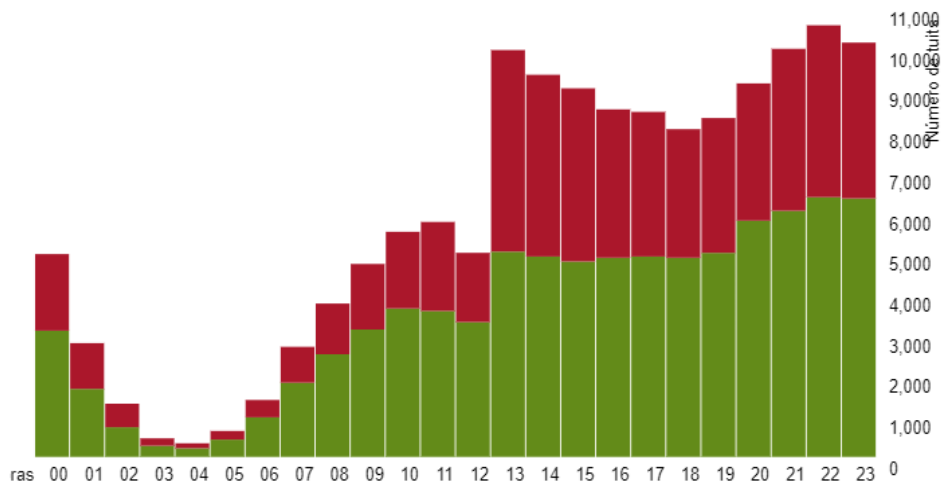
Índice = Positivos (👍) / Negativos (👎)



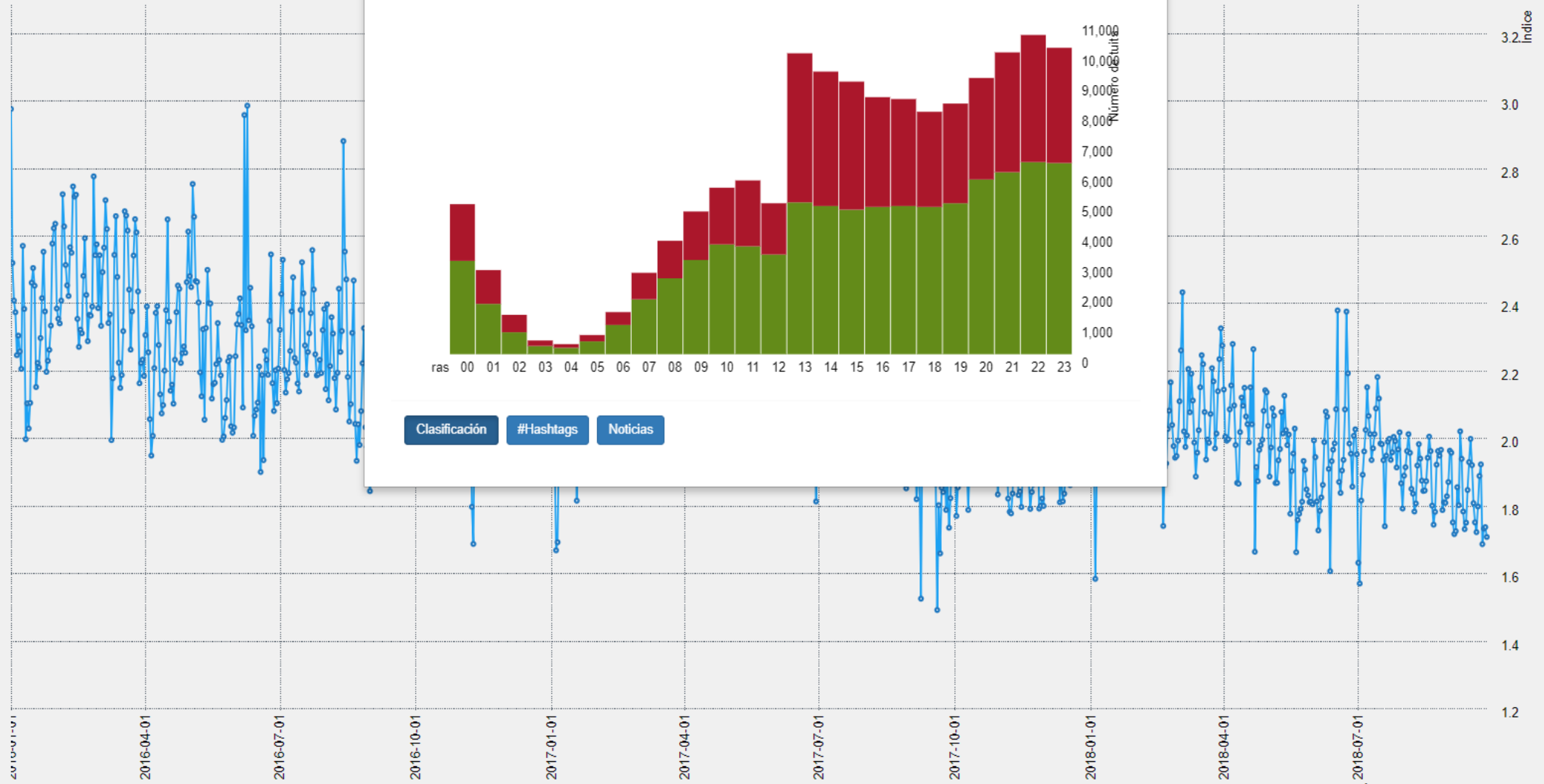


### Clasificación por hora del día (NAC)

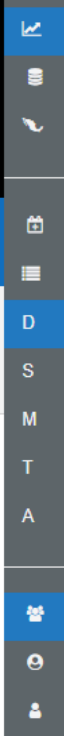
19 de septiembre de 2017



Clasificación #Hashtags Noticias

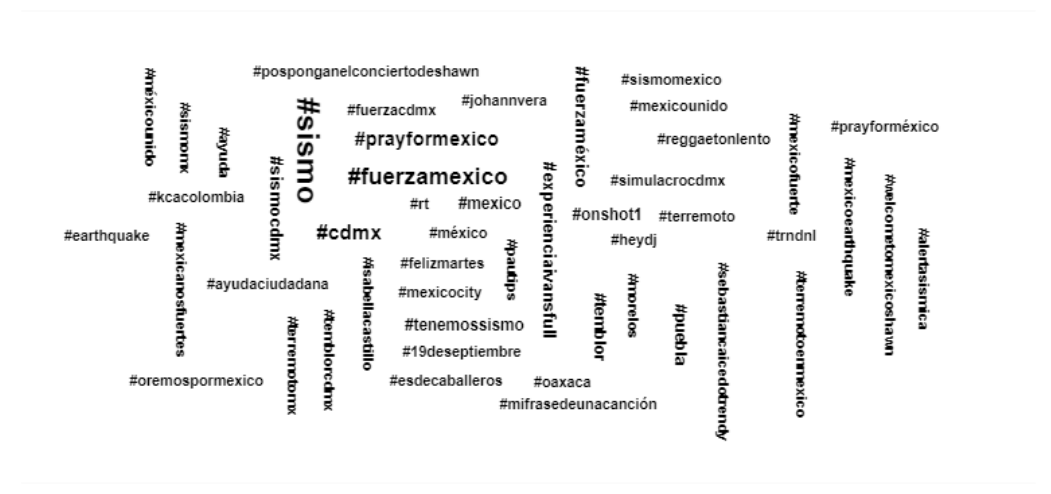


Índice = Positivos (🟢) / Negativos (🔴)

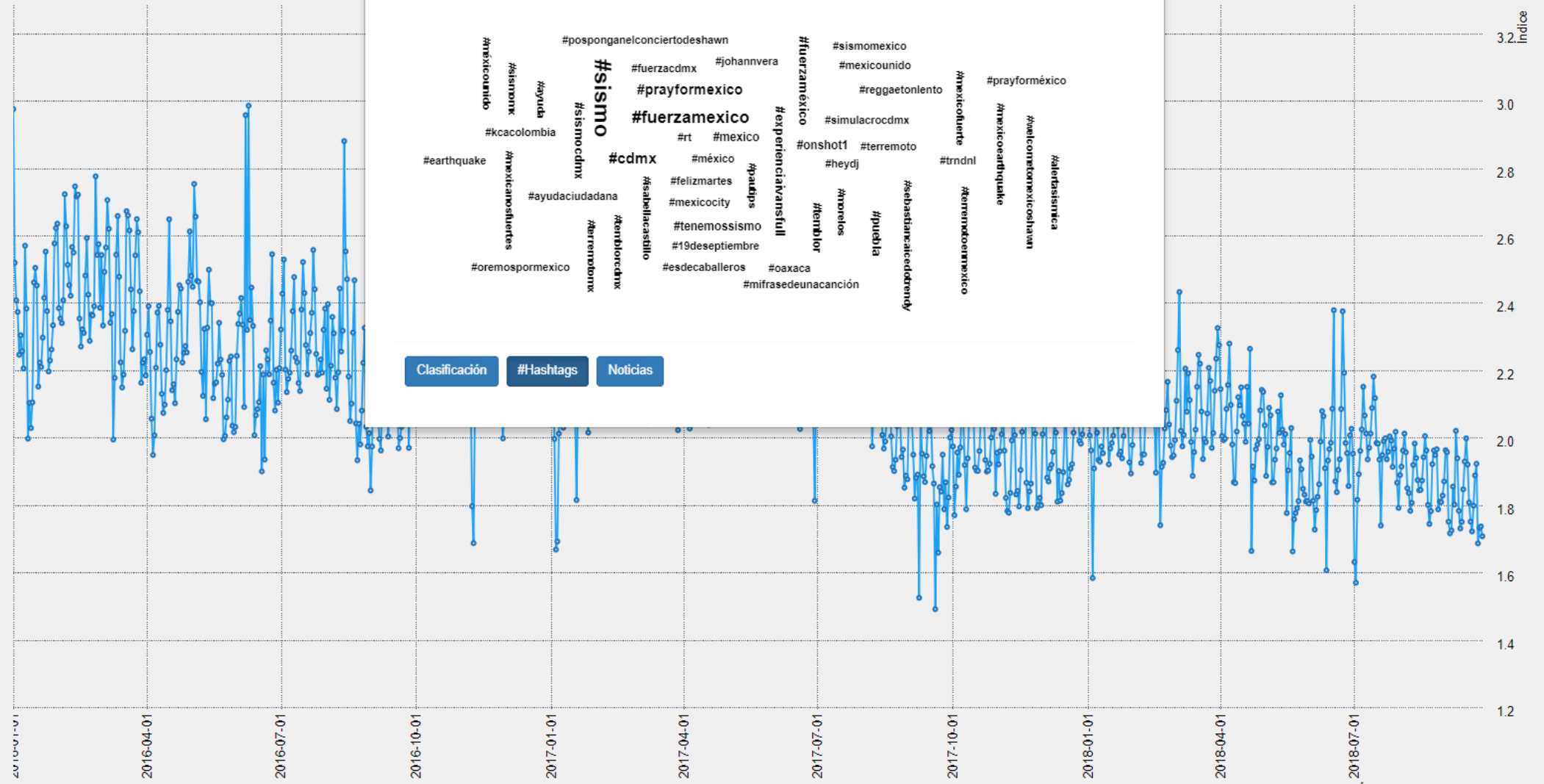


### Nube de #hashtags del día (NAC)

19 de septiembre de 2017

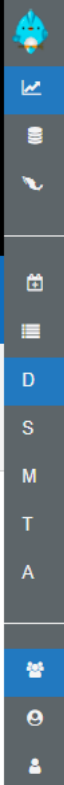


Clasificación #Hashtags Noticias



Índice = Positivos (👍) / Negativos (👎)





# Noticias del día (NAC)

20 de septiembre de 2017



Para ayudar a comprender el comportamiento de un día se ofrece una nube de palabras con los #Hashtags del día y ligas a sitios externos de noticias del día, que no son responsabilidad de INEGI. Se seleccionaron como buscadores de medios online a Google por ser el más utilizado con 3.5+ millones de peticiones por minuto, y a Kiosko.net porque reúne 3,600+ medios informativos online de todo el mundo, de los cuales 197 corresponden a México, aún y cuando para Colima, Durango y Nayarit no exista información (datos correspondientes a diciembre 2017).

Al dar clic en cualquier medio online, saldrá del sitio.

- Clasificación
- #Hashtags
- Noticias



Índice = Positivos (👍) / Negativos (👎)



[Otro terremoto en 19 de septiembre](#)[El Universal](#) - 19 sep 2017

mario.dorantes@eluniversal.com.mx. Ni un día más, ni un día menos. Justo hoy 19 de septiembre de 2017 la Ciudad de México volvió a vivir las mismas ...



Noticieros Televisa

[Aumenta a 369 víctimas por sismo en México](#)[El Financiero](#) - 19 sep 2017

Al menos 369 personas fallecieron tras el sismo de 7.1 de magnitud que el 19 de septiembre sacudió el centro y sur del país. El coordinador nacional de ...

[Fotos del sismo de la Ciudad de México que no pudimos transmitir en ...](#)[Noticieros Televisa](#) - 19 sep 2017

La impotencia ante la tragedia #Sismo #Puebla estas imágenes son muy fuertes pero hablan de la solidaridad pic.twitter.com/YBPTeXY2j3. — Maricela Luna ...

[19 de septiembre, otra vez; coincidencias y diferencias entre sismos ...](#)[El Universal](#) - 19 sep 2017

19 de septiembre, otra vez. Fue el 19 de septiembre pero de 1985 cuando un terremoto de 8.1 grados sacudió a la Ciudad de México. Otro sismo este 19 de ...

[Así se vivieron las primeras horas del terremoto de 1985](#)[Milenio.com](#) - 19 sep 2017

El 19 de septiembre de 1985 está marcado como una de las fechas más trágicas en la historia de México, pues ese día un terremoto de 8.1 grados Richter ...

[Aeropuerto de México reanuda operaciones](#)[El Diario de Juárez](#) - 19 sep 2017

Ciudad de México– Tras el sismo de 7.1 grados, el Aeropuerto Internacional de la Ciudad de México (AICM) reactivó sus operaciones a las 16:00 horas.

Kiosko.NET

Periódicos de México  
Toda la prensa de hoy

Hasta **20** Megas por solo \$435 al mes\*  
Con telefonía ilimitada

Hemeroteca ▾ 20/Sep/2017

Inicio África Asia-Pacífico Europa Latinoamérica USA Canadá

Últimas ediciones ▾

México

- Aguascalientes
- Baja California
- Campeche
- Chiapas
- Chihuahua
- Coahuila
- Colima
- Durango
- Est. México
- Guanajuato
- Guerrero
- Hidalgo
- Jalisco
- México D.F.
- Michoacán
- Morelos
- Nayarit
- Nuevo León
- Oaxaca
- Puebla
- Querétaro
- Quintana Roo
- San Luis Potosí
- Sinaloa
- Sonora

Prensa de Información General ▾



publicidad

Prensa Deportiva ▾



Prensa Económica ▾



# Mapa del estado de ánimo de los tuiteros en México

1 de enero de 2016 - 26 de septiembre de 2018

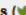
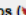


Índice promedio nacional: 2.15

Rango del periodo: [-0.1 , 3.43]

Rango seleccionado: [0 , 3.43]



Índice = Positivos  / Negativos 

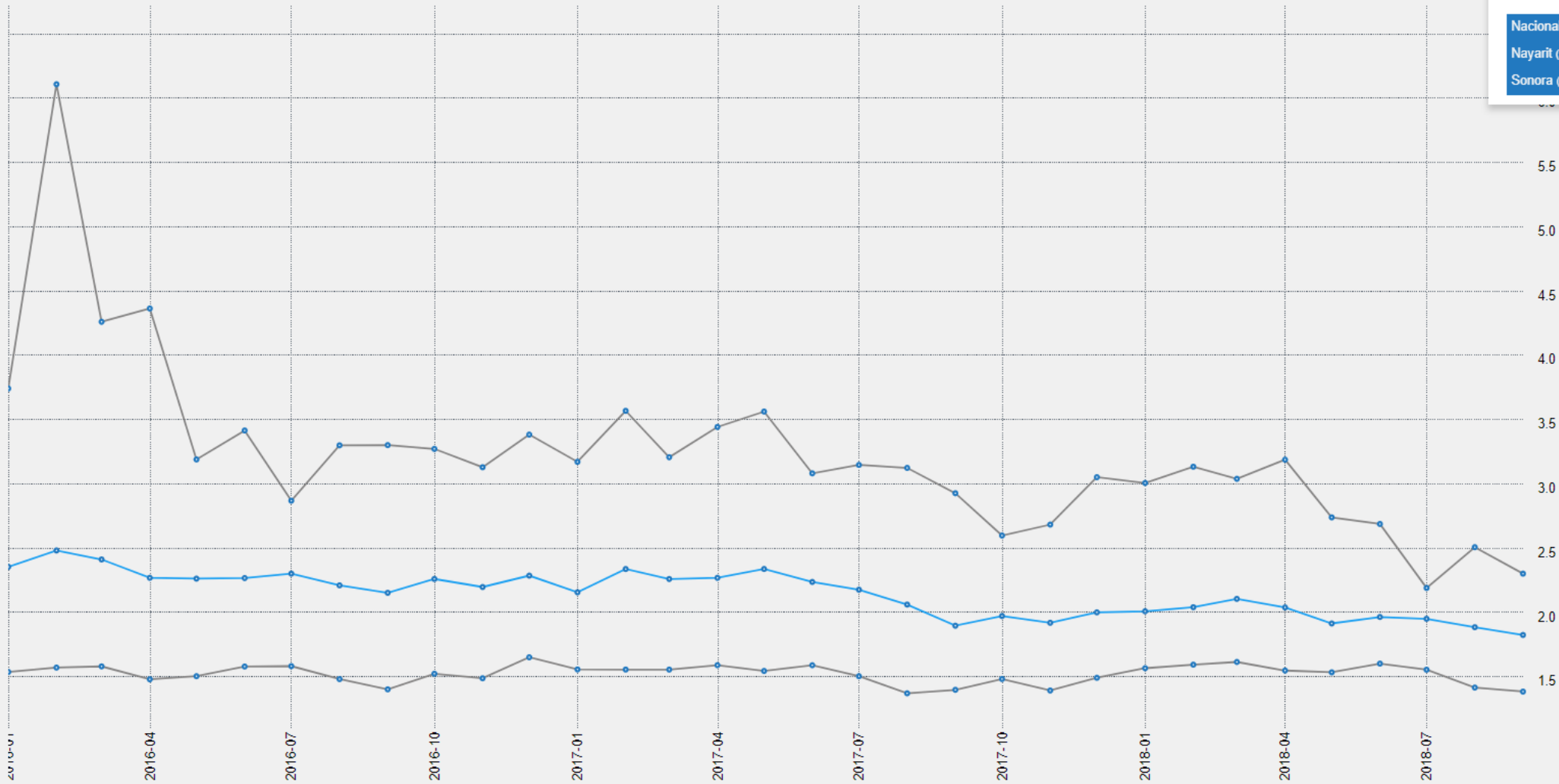


# Estado de ánimo de los tuitos en México

1 de enero de 2016 - 26 de septiembre de 2018

Entidad federativa

- Nacional (NAC) ✕
- Nayarit (NAY) ✕
- Sonora (SON) ✕



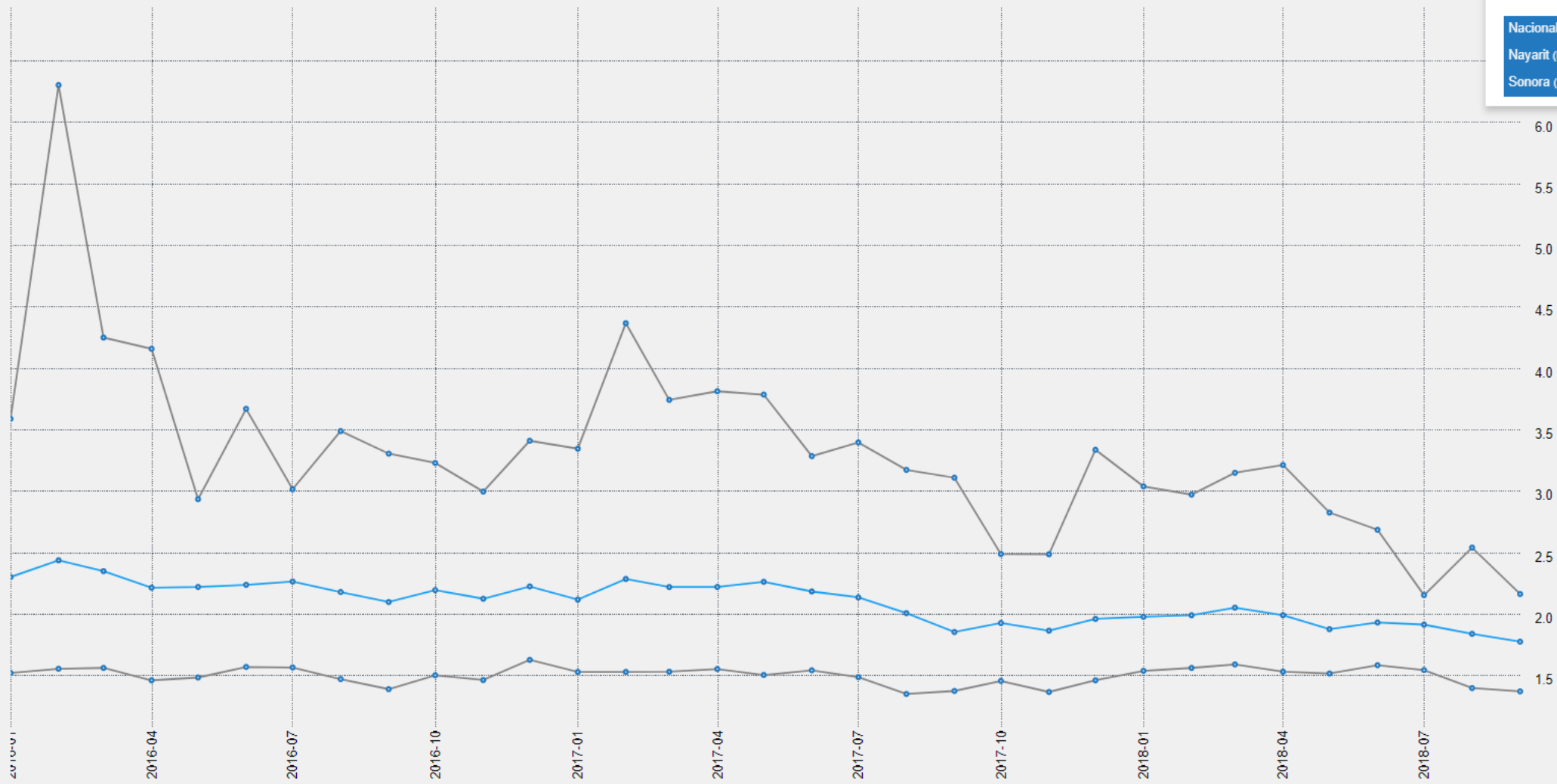
Índice = Positivos (😊) / Negativos (😡)

# Estado de ánimo de los tuiteros locales en México

1 de enero de 2016 - 26 de septiembre de 2018

Entidad federativa ✕

- Nacional (NAC) ✕
- Nayarit (NAY) ✕
- Sonora (SON) ✕



Índice = Positivos (😊) / Negativos (😞)

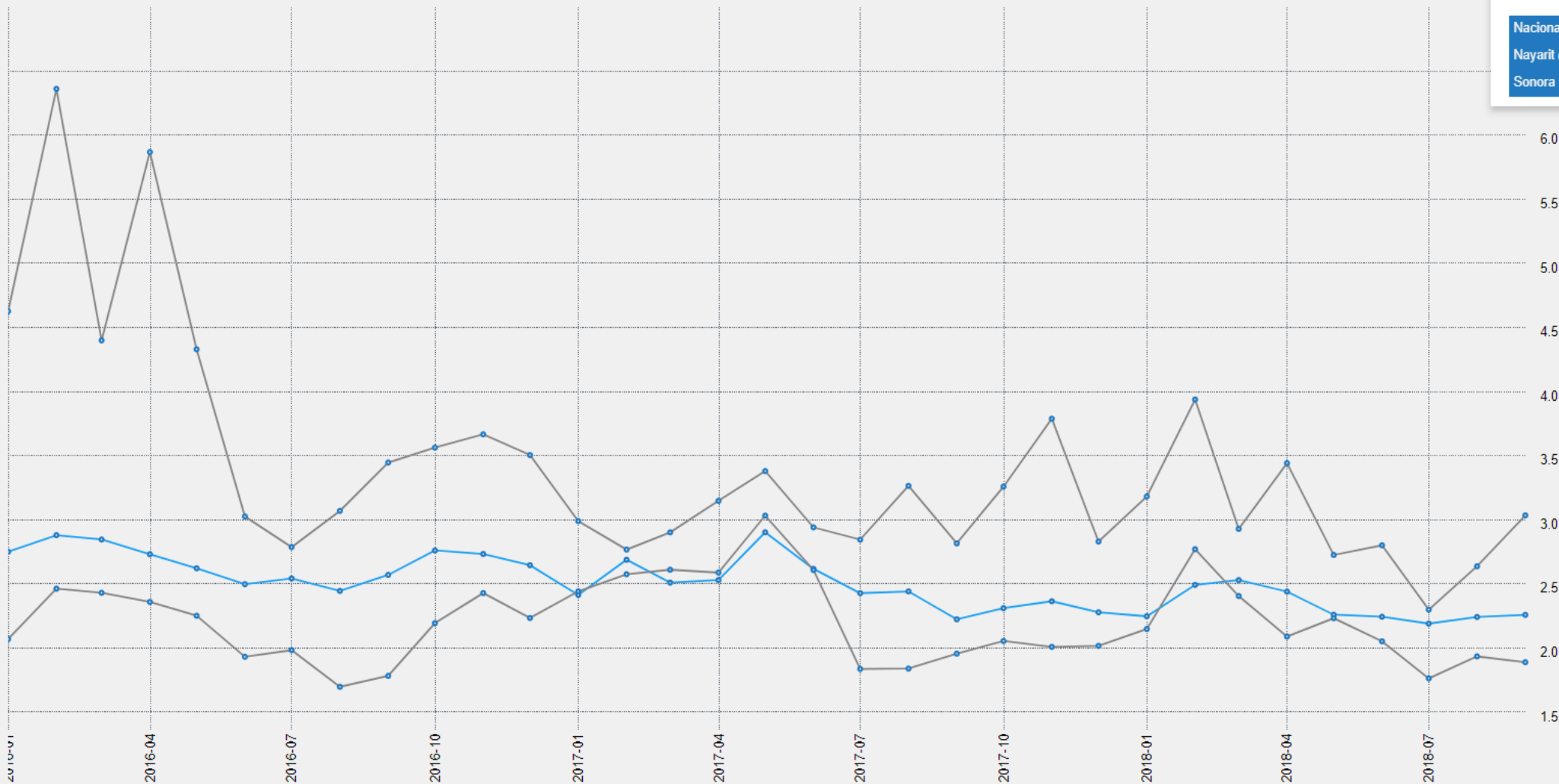


# Estado de ánimo de los tuitos visitantes en México

1 de enero de 2016 - 26 de septiembre de 2018

Entidad federativa ✕

- Nacional (NAC) ✕
- Nayarit (NAY) ✕
- Sonora (SON) ✕



Índice = Positivos (😊) / Negativos (😡)





# ELECCIONES, 2018

Cátedra  
**INEGI**

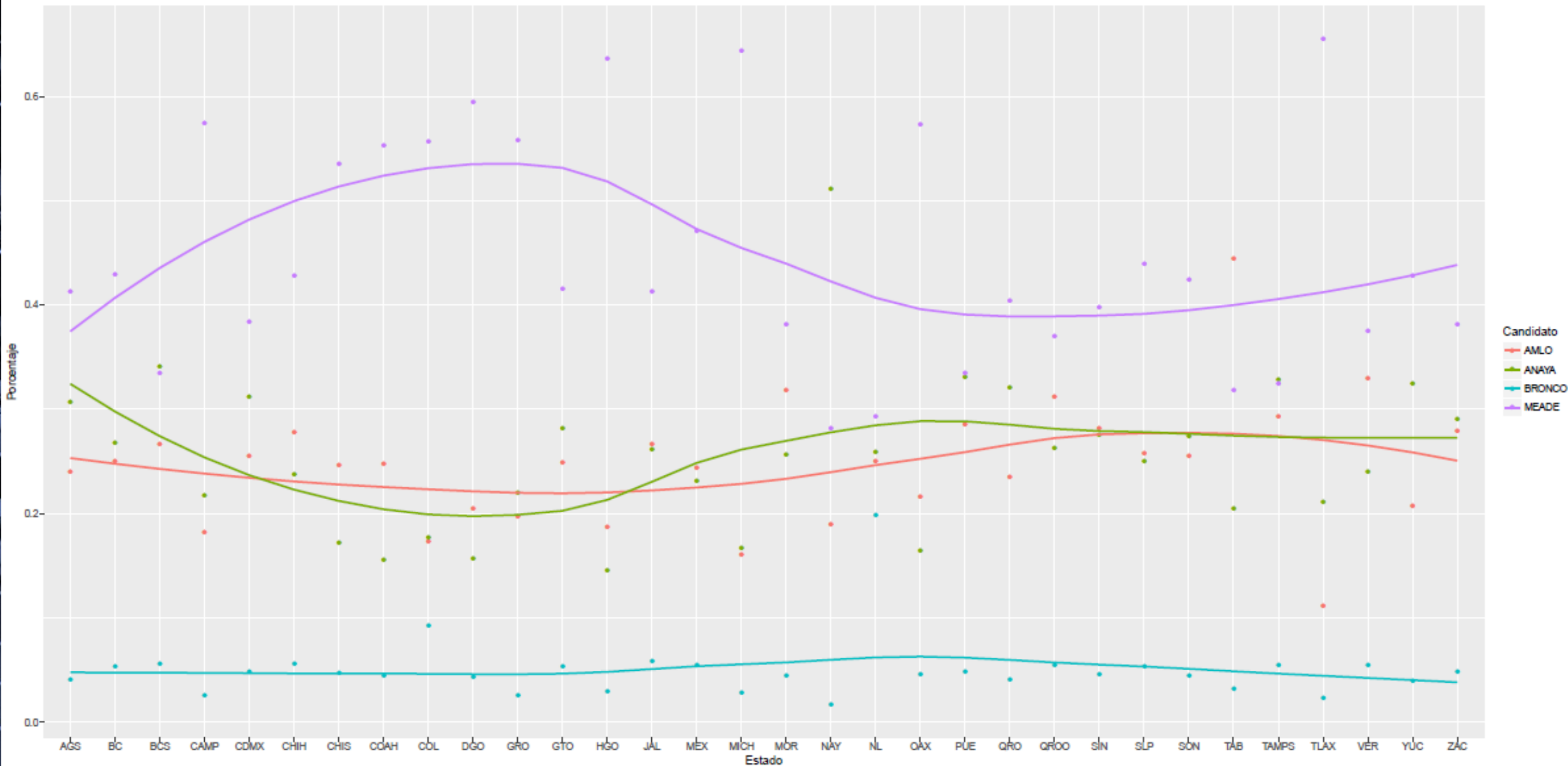


- Selección de tweets clasificados como “positivos”, filtrando por fecha y por texto.
- Texto es relacionado con cada uno de los cuatro candidatos, para generar una consulta que contenga información solo de ese candidato en un periodo de tiempo determinado.
- Las consultas fueron semanales a partir de Abril hasta Junio del presente año.

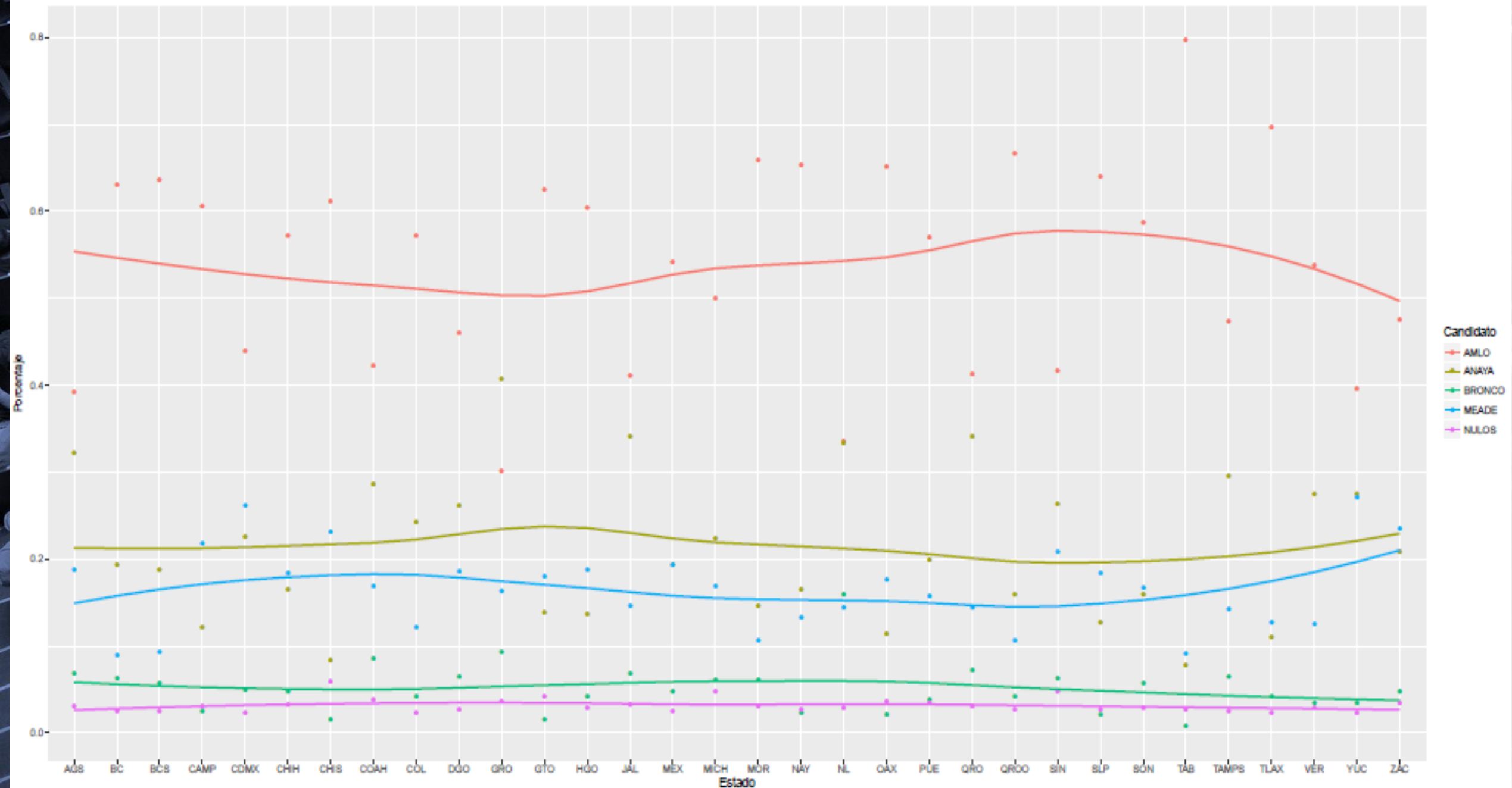
- Además, base de datos que contiene el conteo de votos por candidato, proporcionada por el INE.
- Con la base del INE se calculó el promedio de votos a favor de cada candidato por estado.
- En el caso de los tweets, se calculó el promedio de datos semanales y por estado, a favor de cada candidato.
- Dando como resultado las siguientes gráficas.

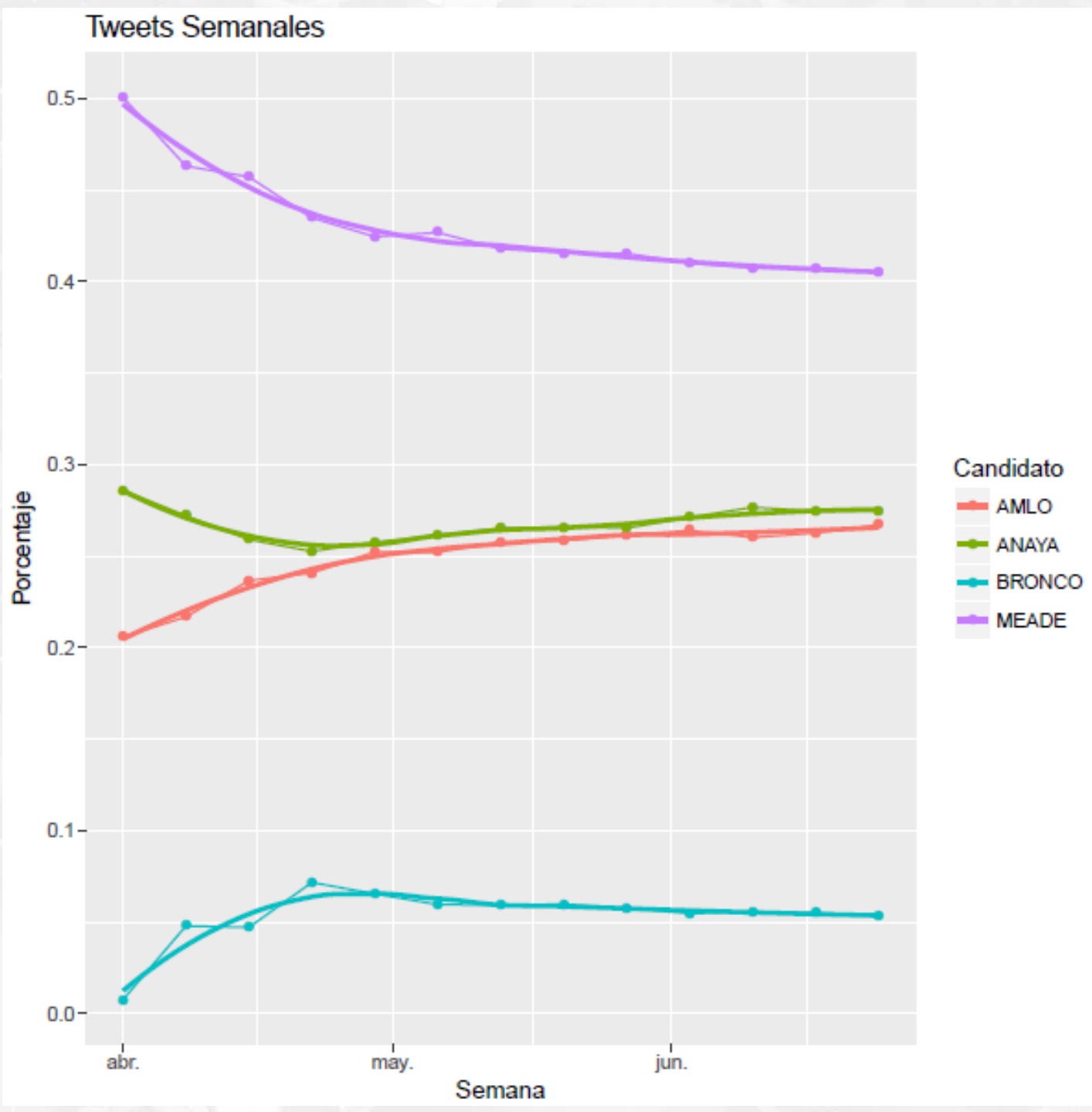


Tweets por Estado



# Resultados INE



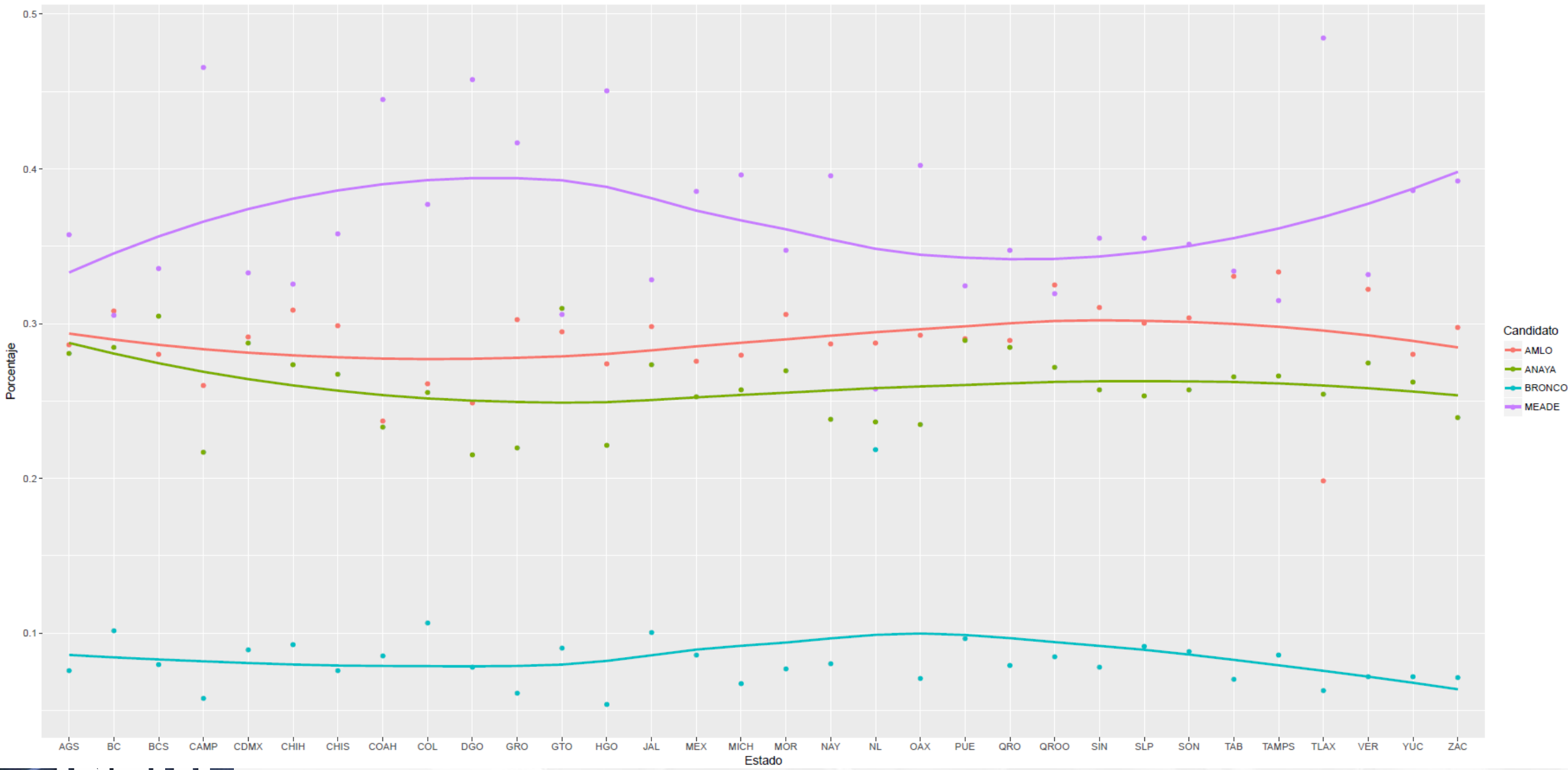




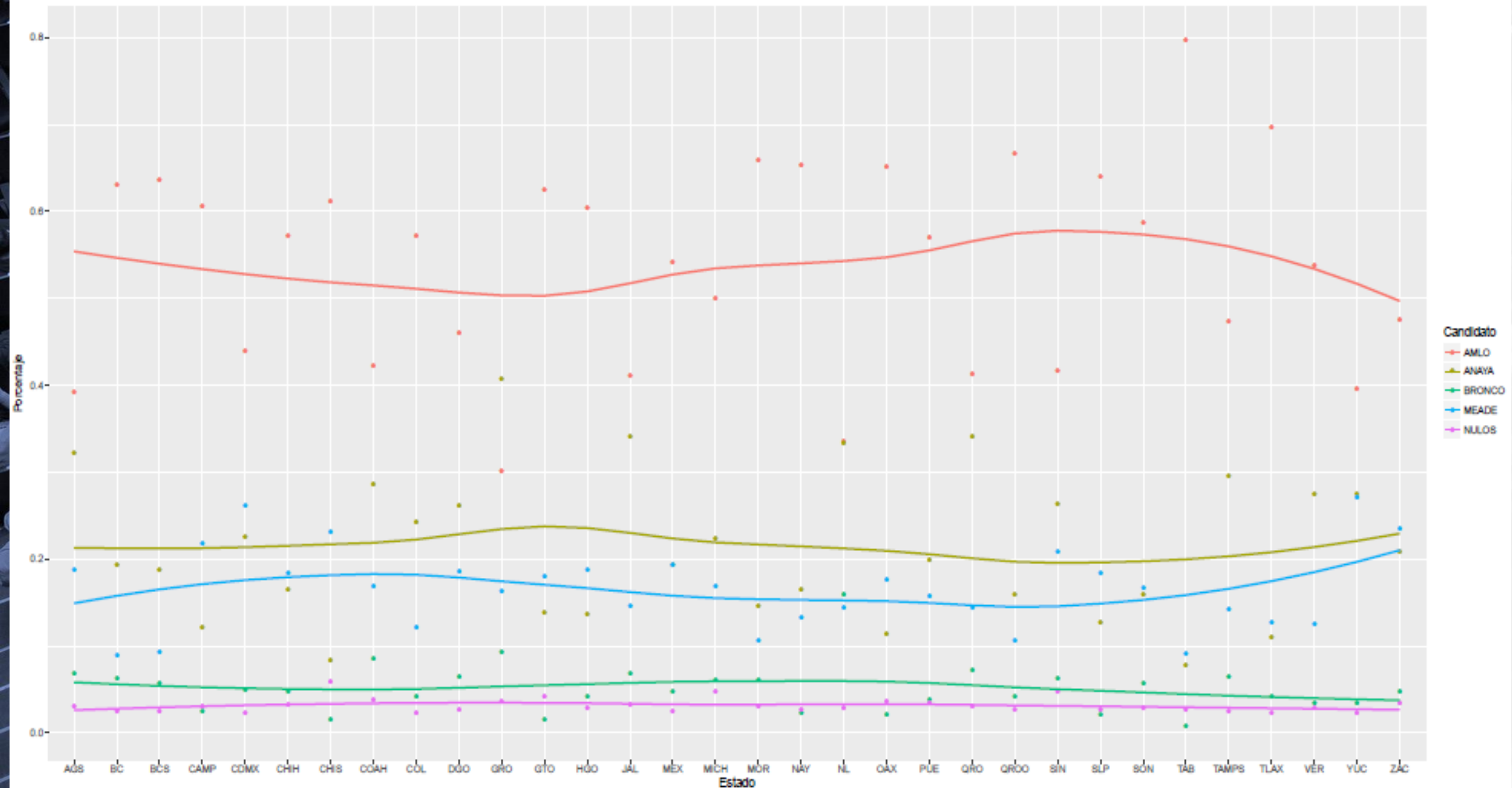
# Segundo intento

- La contabilización en el ejercicio anterior corresponde al número de tuits.
- Es decir, un tuitero puede votar más de una vez.
- Durante la elección ocurre algo diferente.
- En consecuencia, se decidió eliminar repeticiones del tuitero.
- En otras palabras, ahora contabilizamos tuiteros y no tuits
- Los resultados cambian.

Tweets por Estado



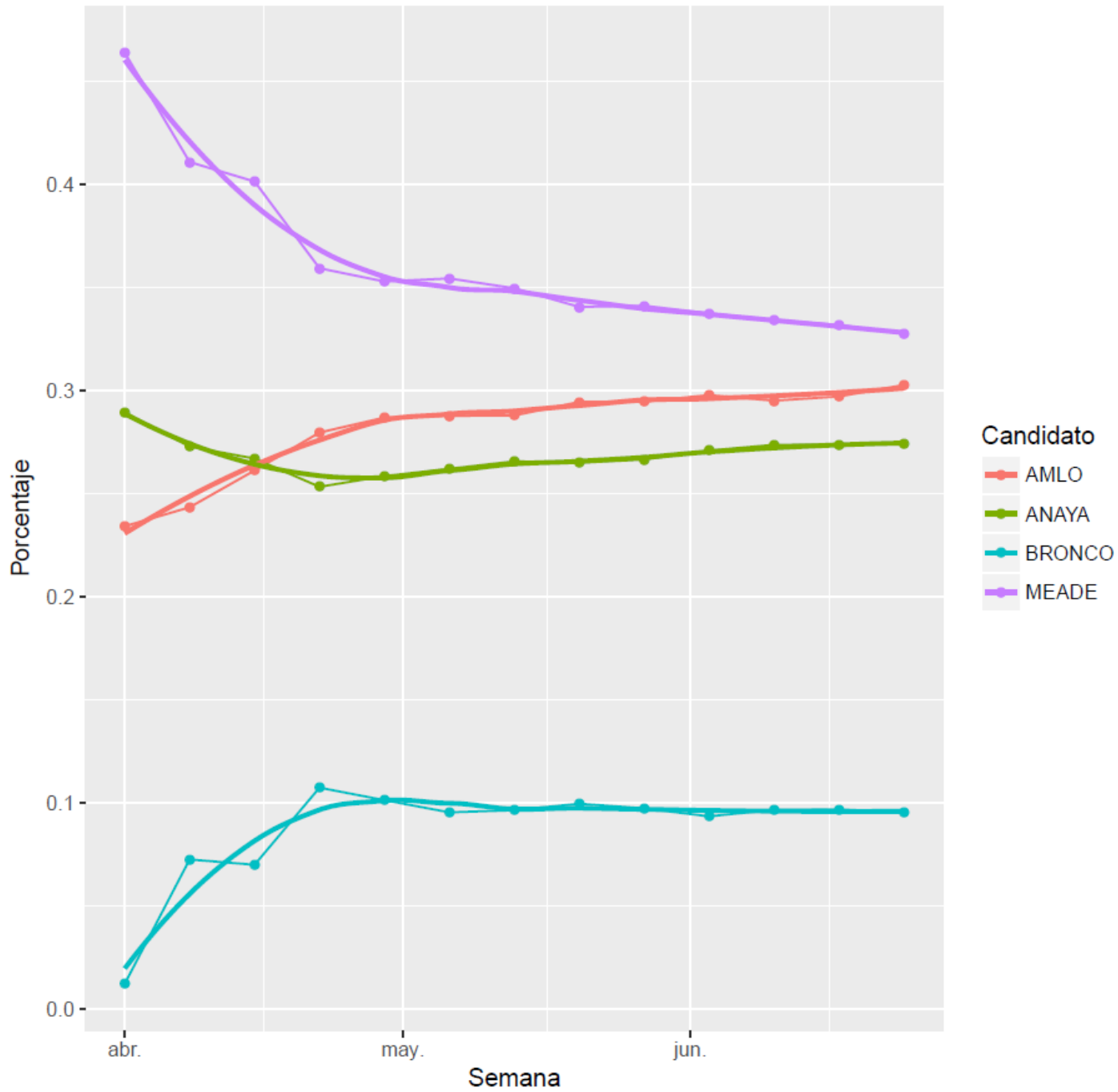
# Resultados INE



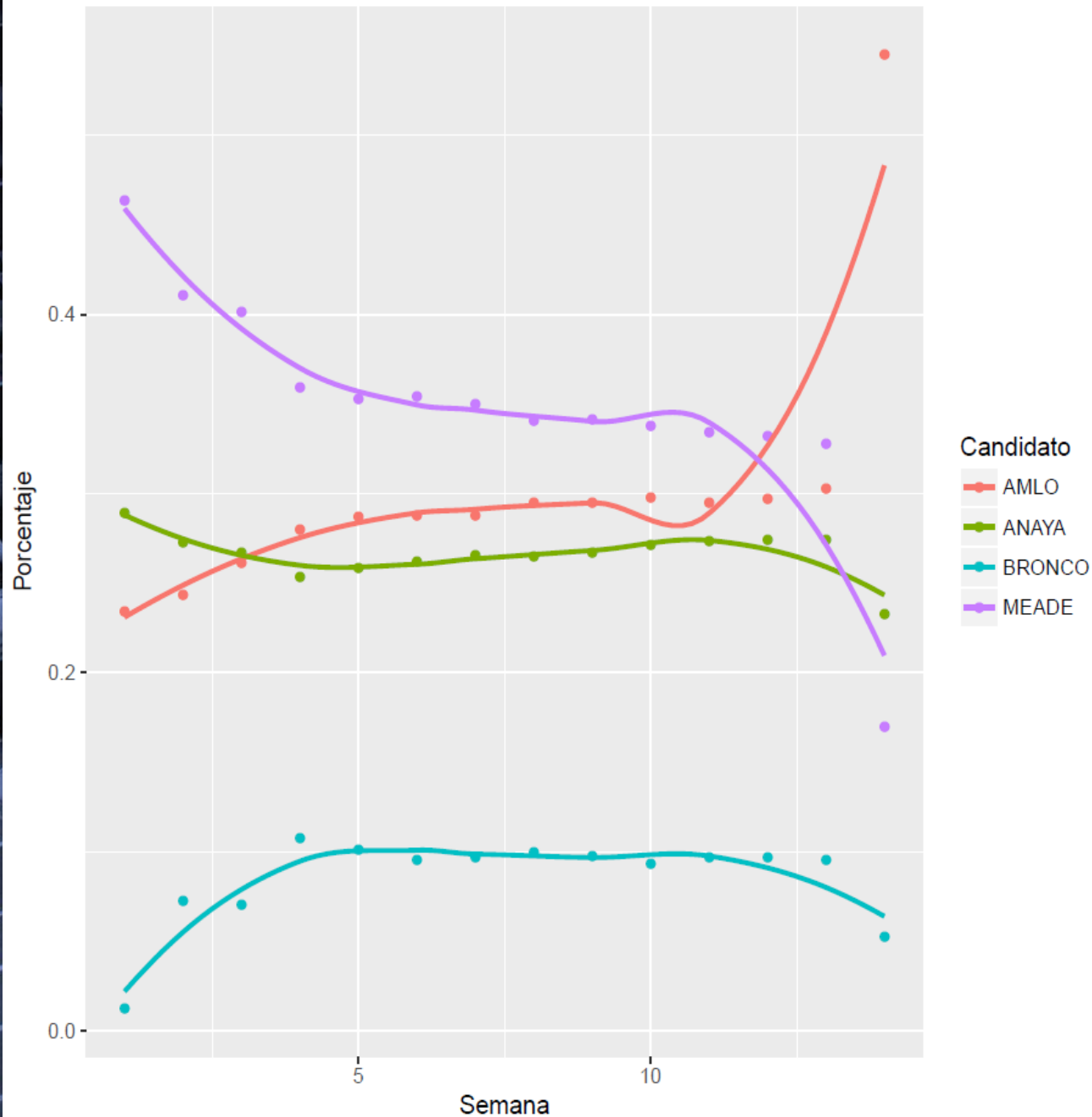




Tweets Semanales



Tweets vs Resultado Elecciones



Agregando a la gráfica de los Twitteros semanales el resultado de las elecciones, se ve un salto no trivial en la línea.

Indica que el resultado de las elecciones no es predecible a partir de los datos encontrados en el análisis twittero.

# El futuro de Big Data en la producción de información

- ¿Estamos dispuestos a integrarnos en un proceso de modernización acelerada?
- ¿Cuál será la aportación de valor de la información que produce el INEGI en el futuro?
- ¿Cómo nos adaptaremos a las nuevas necesidades de información de los ciudadanos?







# Conociendo México

01 800 111 46 34

[www.inegi.org.mx](http://www.inegi.org.mx)

[atencion.usuarios@inegi.org.mx](mailto:atencion.usuarios@inegi.org.mx)



**INEGI** Informa