

Biometria II

TP N° 6

Modelos Lineales Generalizados Regresion Poisson

Hasta ahora teniamos para los modelos lineales, lo siguiente:

$$Y_i \sim N(\mu, \sigma^2)$$
$$Y_i = g(x_i) + \epsilon_i$$

Con lo cual:

$$E[Y_i] = \mu_i = g(x_i)$$

Sin embargo ahora supondremos:

$$\mu_i = e^{g(x_i)}$$

o

$$\log(\mu_i) = g(x_i)$$

Entonces:

$$Y_i \sim P(\mu_i)$$

y

$$E[Y_i / \text{continuo}_i] = \mu_i = e^{g(x_i)} = e^{\beta_0 + \beta_1 * X_{i1} + \dots + \beta_p * X_{ip}}$$

Distribucion de Poisson

La distribucion de *Poisson* es muy utilizada para modelar el numero de ocurrencias de un evento en un intervalo de tiempo dado. Por ejemplo, la espera de un colectivo, cantidad de presas encontradas por un predador por mes, o el numero de individuos / m² tipicamente se pueden modelar con una distribucion de *Poisson*. Uno de los supuestos basicos sobre los que esta distribucion se construye es, que para intervalos de tiempo cortos, la probabilidad de ocurrencia del evento es proporcional a la medida del tiempo esperado. Esta distribucion posee solo un parametro λ , que es la esperanza. Entonces una *v.a* tendra una distribucion de *Poisson* si:

$$P(X = x | \lambda) = (e^{-\lambda} * \lambda^x) / x!$$

$$E(X) = \lambda$$

$$Var[X] = \lambda$$

Antes de arrancar vamos a inspeccionar un poco como se ve la distribucion de Poisson conforme cambia el valor de λ . Abra una explorador de internet y vaya al siguiente <https://homepage.divms.uiowa.edu/~mbognar/>.

Explore la *applet* para la distribucion de Poisson y responda:

-Cual es la probabilidad de que $X=2$ cuando $\lambda=7$?

-Cual es la probabilidad de que $X=2$ cuando $\lambda=7$?

Problema 1. Reservas Urbanas y conservación de aves

Las reservas urbanas pueden, entre otros servicios ecosistémicos, brindar refugio para las aves nativas. Esto representa un factor importante en la conservación de la biodiversidad. Con el fin de evaluar este servicio, se llevo a cabo un estudio para cuantificar la diferencia en abundancia de aves nativas entre ambientes antropizados y reservas urbanas con planes de manejo. Para ello, se seleccionaron 35 puntos de muestreo en la Ciudad de Buenos Aires de los cuales 21 correspondieron a ambientes antropizados y 14 a la Reserva Ecologica Costanera Sur. En cada uno de los puntos se realizaron observaciones de aves.

El protocolo original consistio en realizar conteos de aves nativas en cada punto durante 5 minutos, pero dependiendo de las condiciones en las que se daba el registro, la duracion final de las observaciones vario entre 2 y 7 minutos. Por ello para cada conteo se cuenta tambien con el registro del tiempo de observacion (en minutos). Base de datos “*aves.txt*”.

- Indique cual es la variable dependiente o respuesta. ¿Cual es su potencial distribucion de probabilidades? ¿Cual es la variable explicatoria? ¿De que tipo es?
- Plantee y escriba el modelo en terminos de regresion.
- Describa grafica y estadisticamente los datos.
- Pruebe los supuestos del modelo. ¿Se verifican?

Chequeamos residuos vs. predichos.

- Valide el modelo.
- En que medida las reservas favorecen la conservacion de aves nativas?. Informe la magnitud del efecto en escala de la variable respuesta.

Notar que el summary nos devuelve: a) Intercept y b) Ambiente Reserva.

- Intercept (tasa para el area metropolitana)

$$\hat{\beta}_0 = \log(\widetilde{aves}_{metropolitana})$$

-Ambiente Reserva (efecto reserva sobre metropolitana)

$$\hat{\beta}_1 = \log(\widetilde{aves}_{reserva}) - \log(\widetilde{aves}_{metropolitana})$$

Es decir, el incremento respecto del ambiente metropolitano.

Problema 2. Atropellamiento de anfibios en una carretera en las cercanias de un parque natural (Modificado a partir de base de datos de Zuur 2009)

Los datos presentados provienen de un estudio de dos anios sobre vertebrados atropellados en una ruta nacional del sur de Portugal, pavimentada y con trafico moderado. En las cercanias se encuentran ambientes boscosos, tierras abiertas, incluyendo pastos, prados y barbecho.

La ruta fue inspeccionada cada dos semanas por dos anios. Se identifico a cada animal encontrado muerto a nivel de especie, siempre que hubiese sido posible, y se registro su ubicacion geografica en coordenadas. Para fines del analisis de datos, la ruta se dividio en 52 segmentos de 500 mts y se presenta el numero total de anfibios muertos por segmento (“TOT.N”). Los datos se encuentran en la base *roadkills.txt*.

En cada punto se registraron variables ambientales (ver tabla). En particular se desea conocer si el numero de anifbios atropellados (TOT.N) se encuentra relacionado con las variables OPEN.L, MONT.S, POLIC, SHRUB, WAT.RES, L.WAT.C, L.P.ROAD, D.WAT.COUR y D.PARK. No se tienen evidencias de que haya interaccion entre las variables.

- Realice un analisis exploratorio de las variables involucradas. ¿Detecta datos atipicos? ¿Como son las relaciones entre las variables explicatorias? ¿Y entre las explicatorias y la respuesta?
- Se decide aplicar raiz cuadrada sobre las variables POLIC, SHRUB, WAT.RES, L.P.ROAD, D.WAT.COUR. ¿Por que cree que se realizo este procedimiento?

Se aplica la raiz cuadrada debido a los altos valores que toman las variables.

- Plantee el modelo adecuado.

- Calcule la sobredispersión del modelo aditivo completo.
- ¿Es correcto modelar suponiendo distribución de Poisson? ¿Por qué? Si considera que no es correcto utilizar la distribución *Poisson*, ¿qué alternativas de análisis propone?.

A que llamamos sobredispersión o subdispersion? Como se ve?

Supongamos que estamos modelando la disposición en un determinado lugar de individuos de una especie. Si bien esta disposición de los individuos puede ser modelada usando una distribución de Poisson, algunos individuos podrían **agruparse** para ocupar el terreno (dispersión por contagio) o **alejarse** entre sí lo más posible (dispersión por rechazo). Estos casos no pueden ser adecuadamente modelados por la distribución de *Poisson* porque se viola el supuesto de la independencia entre los eventos. En estos casos la variabilidad se aleja de la media. El coeficiente de dispersión ($CD = \text{varianza}/\text{media}$) permite identificar estas situaciones. Si los individuos se agrupan siguiendo un patrón de contagio, la variabilidad aumenta más allá de la explicada por la distribución de *Poisson* (sobredispersión). En cambio, si los individuos se agrupan siguiendo un patrón de rechazo, la variabilidad estará por debajo de la explicada por la distribución de *Poisson* (subdispersión).

Time out para un ejemplo

- Se sabe que la cantidad de hojas de una plántula a los 30 días post-germinación tiene una media de 8.56 hojas. Corra el script.
- ¿Qué media y qué varianza esperaría obtener? ¿Es razonable suponer una distribución *Poisson*?

Una forma efectiva de modelar datos de conteo con sobredispersión es usando la distribución binomial negativa.

Que ventaja tiene la distribución Binomial Negativa?

La distribución *Binomial* cuenta el número de éxitos en un número prefijado de ensayos de *Bernoulli*. Supongamos, en cambio, que contamos el número de ensayos de *Bernoulli* requeridos para conseguir un número prefijado de éxitos. Esta formulación nos anticipa la distribución *Binomial Negativa*.

Entonces, en una secuencia de ensayos independientes de *Bernoulli* (p), sea la v.a. X , que denota el ensayo para el cual el r -ésimo éxito ocurre, donde r es un entero prefijado. Entonces:

$$P(X = r | r, p) = \binom{x-1}{r-1} * p^r * (1-p)^{x-r} \quad (1)$$

donde (1) tiene distribución *Binomial Negativa* con parámetros r y p .

La obtención de esta fórmula se entiende rápidamente de la distribución binomial. El evento $\{X=x\}$ puede ocurrir solamente si hay exactamente $r-1$ éxitos en los primeros $x-1$ ensayos y un éxito en el ensayo x . La probabilidad de $r-1$ éxitos en $x-1$ ensayos es la probabilidad binomial:

$$\binom{x-1}{r-1} * p^{r-1} * (1-p)^{x-r}$$

y con probabilidad p hay un éxito en el ensayo x . Multiplicando estas probabilidades se llega a (1).

La distribución *Binomial Negativa* está relacionada con la distribución de *Poisson* y puede parametrizarse usando μ y α ($\alpha > 0$). La media de la distribución *Binomial Negativa* es μ y su varianza $\mu + \alpha * \mu^2$. De esta manera si $\alpha = 0$, la distribución Binomial Negativa es igual a la Poisson.

Ahora vamos a caracterizar a la distribución *Binomial Negativa* para entender cómo influyen los diferentes parámetros sobre ella.

En R tenemos que $\text{mu} = \text{media}$ y $\text{varianza} = \text{mu} + (\text{mu}^2)/\text{size}$

Por otra parte, es conveniente visualizar qué relación guarda la *esperanza* con la *varianza*.

Seguimos con el ejercicio

- Realice el modelo utilizando la distribución Binomial negativa. Aplique un método de selección de modelos de manera de incorporar variables de a una. Utilice AIC para decidir la inclusión o no de una variable en el modelo.

Se puede chequear que el modelo elegido tiene todos los factores importantes mediante *drop1*.

- Evalúe los supuestos del modelo seleccionado.

Chequeamos colinealidad entre las variables.

Chequeamos residuos vs. predichos.

- Concluya en relación a las condiciones que favorecen el atropellamiento de anfibios.
- Para saber más: Zuur (2009), plantea este caso y realiza diversos análisis con distinto nivel de complejidad, debido a la naturaleza de los datos y las relaciones entre las variables; además compara distintas formas de modelar. Es un libro de cabecera, aproveche este ejercicio para conseguirlo, enriquecer el contexto y aprender sobre otras estrategias de análisis.

Problema 3. Riqueza de especies de anfibios en la superficie central de América del Sur.

En Base a los datos correspondientes al gradiente latitudinal de diversidad de aves anfibios y las variables ambientales relacionadas a las hipótesis planteadas (18 variables totales, ver enunciados de guías anteriores), se quiere responder:

- ¿La riqueza de especies de anfibios (*Amphibian.Richness*) responde a la temperatura y a la precipitación en la porción central de América del Sur? (*Annual.Mean.Temperature* y *Mean.Annual.Precipitation*) ¿Cómo es esa relación?
- Indique cuál es la variable dependiente o respuesta. ¿Cuál es su potencial distribución de probabilidades? ¿Cuáles son las variables explicatorias? ¿De qué tipo son?
- Explore cómo es la relación entre las variables.
- Plantee los modelos aditivos y con interacción.
- ¿Qué modelo le parece que describe mejor a la Riqueza de especies?
- Ajuste el modelo.
- Valide el modelo.

Problema 4. Efecto del herbicida glifosato sobre la fecundidad de arañas

Las arañas son depredadores importantes de varias plagas agrícolas y desempeñan un papel importante como indicadores de disturbio del ecosistema. En Argentina, el cultivo de soja ha aumentado desde la introducción de la soja transgénica resistente al glifosato. Esta expansión produjo un aumento en el uso del glifosato, un herbicida de amplio espectro, cuyo efecto sobre la fisiología de los artrópodos es muy poco conocida. En los cultivos transgénicos de soja de la provincia de Buenos Aires, *Alpaida veniliae* (Araneae, Araneidae) es una de las arañas tejedoras más abundantes. El propósito de este trabajo fue estudiar los efectos del glifosato sobre algunos atributos reproductivos de *A. veniliae*, en laboratorio. Para ello, hembras fecundadas fueron criadas en frascos de vidrio individuales y alimentadas con moscas, que previamente habían sido tratadas con dosis distintas de glifosato (2) o con solvente. Se utilizaron 10 hembras por tratamiento. Para cada hembra se registró la fecundidad (número de huevos) y la fertilidad (número de crías). No se observó un efecto letal del glifosato. Los resultados se encuentran en el archivo *arana.csv*.

- Indique cuál es la variable dependiente o respuesta. ¿Cuál es su potencial distribución de probabilidades? ¿Cuál es la variable explicatoria? ¿De qué tipo es?
- Describa gráficamente y estadísticamente los datos.
- Plantee el modelo.
- Pruebe los supuestos del modelo. ¿Se verifican?
- Ajuste el modelo. Calcule las predicciones en escala de la variable respuesta para los tres tratamientos.
- Valide el modelo.
- Compare los tres tratamientos y concluya en relación al efecto del glifosato sobre la fecundidad de *A. veniliae*.

Basado en Benamú, M. A., Schneider, M. I., & Sánchez, N. E. (2010). Effects of the herbicide glyphosate on biological attributes of Alpaida veniliae (Araneae, Araneidae), in laboratory. Chemosphere, 78(7), 871-876.

Problema 5.

En base al problema del atropellamiento de anfibios presentado, se quiere modelar la Riqueza de anfibios (“S.RICH”) en función del porcentaje de espacio abierto (“OPEN.L”).

- Plantee el modelo teórico en parámetros y en el contexto de la experiencia.
- Realice un procedimiento adecuado para modelar los datos.
- ¿Cómo se modifica la riqueza de anfibios en función del espacio abierto?