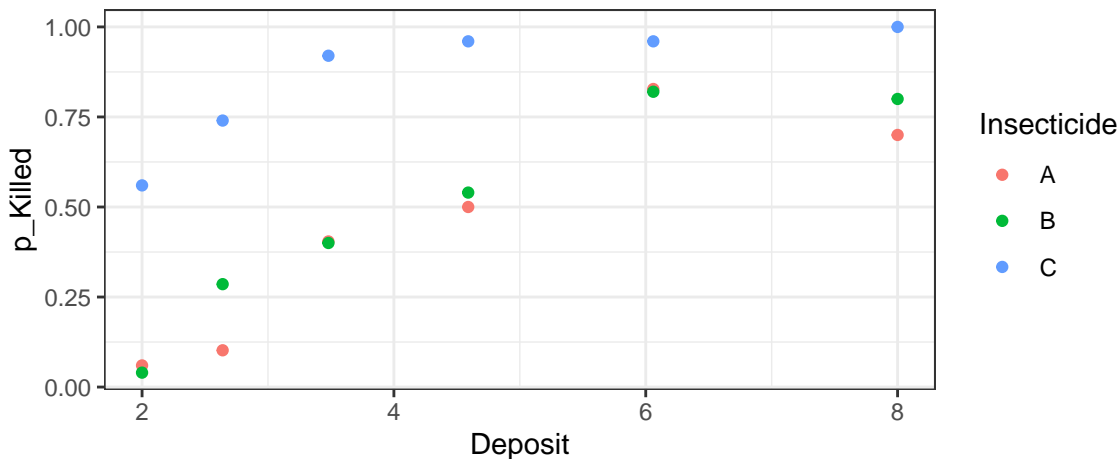


3. Modelos lineales generalizados para datos binarios

La base de datos Preg3B.csv contiene información sobre 862 insectos que fueron expuestos a diferentes dosis en mg (Deposit) de tres insecticidas (Insecticide). La asignación a una dosis y al tipo de insecticida se realizó de forma aleatoria. Después de seis días se analizó si los insectos se habían muerto, de manera que la base de datos contiene también el número de insectos muertos (Killed) y el número total de insectos expuestos (Number) por cada dosis e insecticida. Dado que se asume que el costo de los insecticidas es el mismo, el objetivo del análisis es identificar para cada insecticida qué dosis es la mínima con la que se puede indicar que el 70% de los insectos se muere, así como si considerando la menor de esas tres dosis se puede afirmar que un insecticida es el mejor comparado con el resto. El evento de interés es si el insecto muere o no (died).

i) Gráfica de dispersión de dosis del insecticida y la proporción de insectos muertos.

Se presenta una gráfica de dispersión en donde en el eje x se incluye la dosis del insecticida (Deposit) y en el eje y la proporción de insectos muertos observados (p_Killed) para cada combinación dosis-insecticida (Deposit-Insecticide), distinguiendo con un color el insecticida asociado. Se puede observar que el insecticida C tiene una mayor tasa de mortalidad para todas las seis dosis consideradas (solamente la primera dosis es menor a 70%). Para el caso de los insecticidas A y B, los resultados son muy parecidos, aunque marginalmente parece que el insecticida A tiene menor tasa de mortalidad, al menos de manera evidente en tres dosis distintas.



ii) Ajuste modelos para datos binarios 1

Ajustaremos modelos para datos binarios (ligas: logit, probit, y cloglog) en donde se incluyen como covariables a Insecticide y $\ln D$ ($\ln D = \ln(\text{Deposit})$), así como su interacción. Se calcularon los tres modelos con interacciones y se muestran en el siguiente Cuadro. De acuerdo con el criterio AIC el modelo más adecuado es el de la liga probit, cuyo AIC fue de 789.28 (el del logit de 789.44 y cloglog de 800.46). Los términos de las interacciones no son significativas para los tres modelos (no se rechaza la hipótesis nula de que los coeficientes son cero), mientras que para el intercepto, InsecticideC y $\ln D$ sí se rechaza la hipótesis nula. Esto sugiere que podría ser más adecuado el modelo reducido.

Adicionalmente, se calcularon los tres modelos (ligas logit, probit y cloglog) reducidos, sin las interacciones Insecticide- $\ln D$. Todos tienen un menor AIC, en particular el modelo probit. Se puede observar que en estos casos incluso InsecticideB podría ser estadísticamente significativo si consideramos un nivel de significancia estadística del 10%. Si consideramos el modelo reducido, el modelo probit tiene un mejor desempeño por su AIC y por ser más parsimonioso, con componente lineal o sistemático $\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 = -2.623 + 0.209\text{InsecticideB} + 1.672\text{InsecticideC} + 1.690\ln D$.

La prueba de hipótesis global con la chi-cuadrada del modelo **probit reducido** muestra un valor Chisq de 264.5619 y un p-value muy pequeño ($\Pr(>\text{Chisq})$: 4.633875e-57), mucho menor a 0.05, es decir se rechaza la hipótesis nula, por lo que podríamos proceder con el análisis de los supuestos del modelo. Antes de continuar, revisaremos en el siguiente inciso algunos modelos que incluyan $(\ln D)^2$, y veremos si tienen menor AIC.

Table 1:

	<i>Dependent variable:</i>					
	died					
	<i>logistic</i>	<i>probit</i>	<i>glm: binomial link = cloglog</i>	<i>logistic</i>	<i>probit</i>	<i>glm: binomial link = cloglog</i>
	(1)	(2)	(3)	(4)	(5)	(6)
InsecticideB	0.188 (0.722)	0.105 (0.400)	0.260 (0.530)	0.349* (0.206)	0.209* (0.120)	0.249* (0.135)
InsecticideC	2.110*** (0.790)	1.505*** (0.433)	2.350*** (0.485)	2.840*** (0.254)	1.672*** (0.141)	1.706*** (0.151)
lnD	2.727*** (0.349)	1.634*** (0.194)	1.861*** (0.234)	2.887*** (0.224)	1.690*** (0.122)	1.714*** (0.134)
InsecticideB:lnD	0.111 (0.487)	0.072 (0.270)	-0.004 (0.319)			
InsecticideC:lnD	0.661 (0.671)	0.137 (0.347)	-0.486 (0.327)			
Constant	-4.231*** (0.524)	-2.543*** (0.289)	-3.377*** (0.392)	-4.461*** (0.356)	-2.623*** (0.194)	-3.138*** (0.238)
Observations	862	862	862	862	862	862
Log Likelihood	-388.721	-388.640	-394.229	-389.246	-388.727	-395.786
Akaike Inf. Crit.	789.443	789.280	800.458	786.491	785.454	799.571

Note:

*p<0.1; **p<0.05; ***p<0.01

iii) Ajuste modelos para datos binarios 2

A continuación incluiremos, adicional a los términos de las covariables anteriores, a la interacción de Insecticide con el término cuadrático $(\ln D)^2$. Para las ligas logit y probit, ninguna de las intersecciones con $\ln D$ y $(\ln D)^2$ rechazan la hipótesis nula, es decir ninguna aparece estadísticamente significativa porque el p-value asociado es mayor a 0.05. Para el caso del cloglog, la única intersección estadísticamente significativa al 5% de significancia estadística es InsecticideC:lnD. En los tres modelos se rechaza la hipótesis nula para el intercepto, InsecticideC, $\ln D$ y $\ln D^2$. Los AIC son 786.61, 786.92 y 786.06 para los modelos con liga logit, probit y cloglog, respectivamente, lo que indica que el mejor modelo por el criterio AIC es el de la liga cloglog.

Adicionalmente, se procedió a hacer un modelo reducido con sólo efectos principales, sin estas interacciones y el resultado es que hay menores AIC para los tres modelos considerando las variables explicativas Insecticide, $\ln D$ y $\ln D^2$, sin las interacciones. Por ejemplo, el menor AIC es de 780.01 para el caso de la liga probit, con componente lineal o sistemático $\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 = -3.920 + 0.195 \text{InsecticideB} + 1.701 \text{InsecticideC} + 3.775 \ln D - 0.750 (\ln D)^2$.

La prueba de hipótesis global con la chi-cuadrada del modelo **probit reducido** muestra un valor Chisq de 254.2325 y un p-value muy pequeño ($\Pr(>\text{Chisq})$: 7.974736e-54), mucho menor a 0.05, es decir se rechaza la hipótesis nula, por lo que podemos proceder con el análisis de los supuestos de este modelo reducido más sencillo con el menor AIC.

Se puede notar que una ventaja de introducir el componente $(\ln D)^2$ es que los AIC disminuyeron, por lo que nos quedamos con este modelo probit reducido, para los análisis subsecuentes.

En la prueba de normalidad Lilliefors (Kolmogorov-Smirnov) **normality test** tenemos que el p-value es de 0.88982146208239, por lo que no se rechaza la hipótesis nula de normalidad. Por otra parte, pa la prueba de normalidad de Shapiro-Wilk **normality test** el p-value es de 0.974134219816826, lo que también no rechaza la hipótesis nula de normalidad. Esto se observa en la siguiente Gráfica.

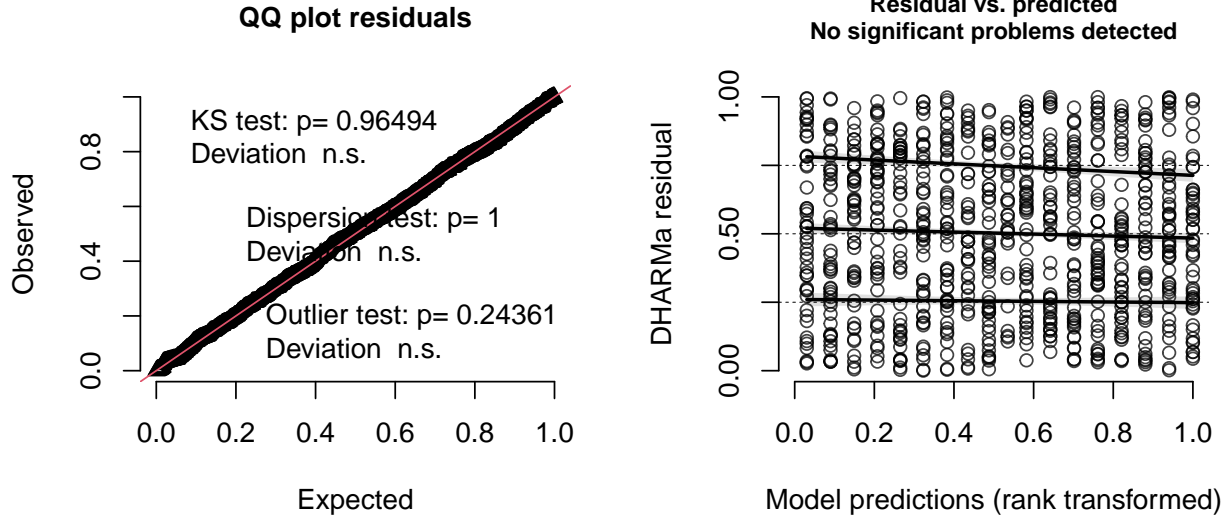
Table 2:

	<i>Dependent variable:</i>					
	died					
	<i>logistic</i>	<i>probit</i>	<i>glm: binomial link = cloglog</i>	<i>logistic</i>	<i>probit</i>	<i>glm: binomial link = cloglog</i>
	(1)	(2)	(3)	(4)	(5)	(6)
InsecticideB	2.013 (2.589)	0.679 (1.361)	1.973 (2.085)	0.325 (0.204)	0.195 (0.120)	0.221* (0.133)
InsecticideC	6.150** (2.684)	2.934** (1.388)	6.139*** (1.825)	2.976*** (0.271)	1.701*** (0.145)	1.663*** (0.152)
lnD	9.085*** (2.778)	4.717*** (1.474)	8.599*** (2.173)	6.813*** (1.408)	3.775*** (0.782)	4.117*** (0.857)
lnD2	-2.167** (0.918)	-1.066** (0.499)	-2.198*** (0.691)	-1.407*** (0.491)	-0.750*** (0.276)	-0.844*** (0.295)
InsecticideB:lnD	-2.479 (3.663)	-0.773 (1.982)	-2.376 (2.790)			
InsecticideC:lnD	-5.238 (4.300)	-1.872 (2.198)	-5.572** (2.530)			
InsecticideB:lnD2	0.839 (1.231)	0.277 (0.678)	0.748 (0.895)			
InsecticideC:lnD2	1.971 (1.656)	0.628 (0.817)	1.558* (0.843)			
Constant	-8.512*** (1.993)	-4.560*** (1.022)	-8.123*** (1.637)	-6.946*** (0.967)	-3.920*** (0.524)	-4.647*** (0.594)
Observations	862	862	862	862	862	862
Log Likelihood	-384.307	-384.460	-384.028	-385.098	-385.006	-391.604
Akaike Inf. Crit.	786.613	786.919	786.055	780.196	780.011	793.208

Note:

*p<0.1; **p<0.05; ***p<0.01

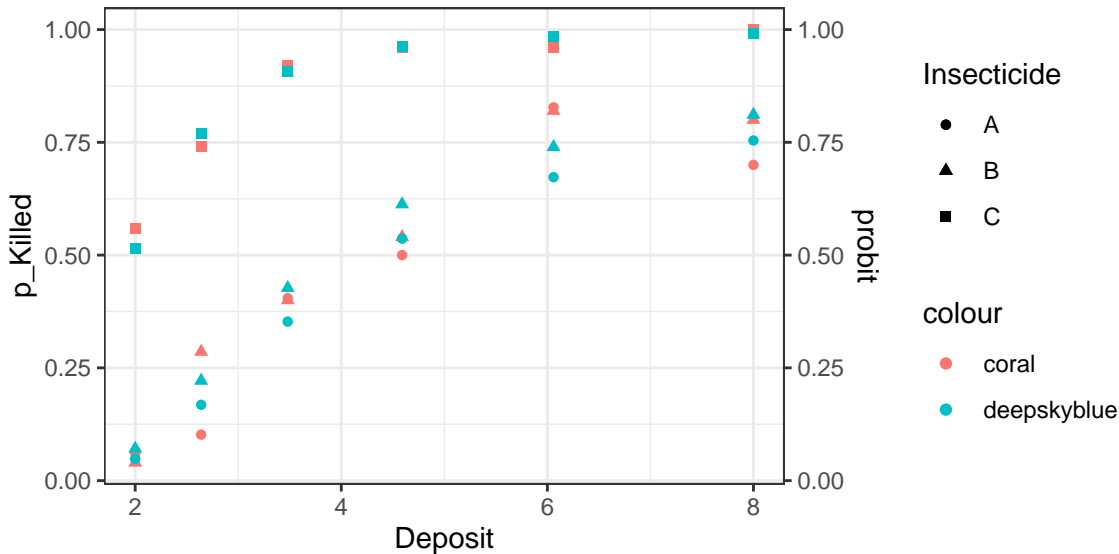
DHARMA residual



La regla de dedo para verificar el parámetro de dispersión de 1, con la devianza de residuales entre grados de libertad, muestra un valor de 0.8984962, lo cual se acerca a 1.

iv) Modelo adecuado. Comparaciones, probabilidades y prueba de hipótesis.

La siguiente gráfica muestra en el eje x la dosis del insecticida (Deposit) y en el eje y izquierdo la proporción de insectos muertos observados (se generó la variable p_Killed) para cada combinación dosis-insecticida (Deposit-Insecticide), distinguiendo con una figura distinta el insecticida asociado. Adicionalmente se agrega un eje y derecho, con las probabilidades obtenidas con el modelo probit para cada dosis e insecticida. El color coral (●) indica los datos de la tasa de mortalidad (p_Killed) y el color deepskyblue (●) indica los valores de las probabilidades calculadas con el modelo (probit). Con el modelo se obtuvieron probabilidades muy cercanas a las proporciones o tasas de mortalidad observadas, especialmente para el insecticida C (■).



A continuación se muestra un cuadro de la dosis mínima para cada insecticida con la que se puede indicar que el 70% de los insectos se muere. Para ello recordemos que $P(died; A) = \beta_0 + \beta_3 \ln D + \beta_4 (\ln D)^2 = 0.7$, $P(died; B) = \beta_0 + \beta_1 InsecticidaB + \beta_3 \ln D + \beta_4 (\ln D)^2 = 0.7$ y $P(died; C) = \beta_0 + \beta_2 InsecticidaC + \beta_3 \ln D + \beta_4 (\ln D)^2 = 0.7$, por lo que resolviendo para cada insecticida, se obtienen los respectivos valores de $\ln D$ y por lo tanto de D que es la dosis en

mg (Deposit). Es decir, para A, resolveremos $\beta_4(\ln D)^2 + \beta_3 \ln D + (\beta_0 - 0.7)$, para B $\beta_4(\ln D)^2 + \beta_3 \ln D + (\beta_0 + \beta_1 - 0.7)$ y para C $\beta_4(\ln D)^2 + \beta_3 \ln D + (\beta_0 + \beta_2 - 0.7)$.

Insecticida	A	B	C
Dosis	8.1600976	6.410879	2.5961564

Table 3:

Como se observa en la Gráfica anterior, el insecticida C es mejor, pues con menores dosis se tienen mayor probabilidad de muerte que A y B según el modelo probit. Además, como se mostró en el cuadro anterior, se encontró que la menor dosis mínima con la que el 70% se muere es para el insecticida C. A continuación mostramos una prueba de hipótesis que comprueba esto. Planteamos entonces que $\beta_0 + \beta_2 \text{InsecticidaC} + \beta_3 \ln D + \beta_4 (\ln D)^2 > \beta_0 + \beta_3 \ln D + \beta_4 (\ln D)^2$ y $\beta_0 + \beta_2 \text{InsecticidaC} + \beta_3 \ln D + \beta_4 (\ln D)^2 > \beta_0 + \beta_1 \text{InsecticidaB} + \beta_3 \ln D + \beta_4 (\ln D)^2$, de donde obtenemos la hipótesis nula $H_0 : \beta_2 \text{InsecticidaC} > 0$ y $\beta_2 \text{InsecticidaC} > \beta_1 \text{InsecticidaB}$ y la hipótesis alternativa $H_a : \beta_2 \text{InsecticidaC} \leq 0$ o $\beta_2 \text{InsecticidaC} \leq \beta_1 \text{InsecticidaB}$.

Resultado: Chisq: 152.8355 y p-value: 0.00000000000000000000000000000000006489137.

A continuación se muestra la prueba de hipótesis que muestra si A y B tienen un desempeño similar. En este caso planteamos que $\beta_0 + \beta_3 \ln D + \beta_4 (\ln D)^2 = \beta_0 + \beta_1 \text{InsecticidaB} + \beta_3 \ln D + \beta_4 (\ln D)^2$ de donde tenemos la prueba de hipótesis $H_0 : \beta_1 \text{InsecticidaB} = 0$ y la alternativa $H_a : \beta_1 \text{InsecticidaB} \neq 0$.

Resultado: Chisq: 2.652954 y p-value: 0.1033576.