

Curso-Taller de R para investigadores

UMSA, La Paz, Bolivia

23 - 25 Feb 2023

Saras Windecker & David Uribe

Código de Conducta

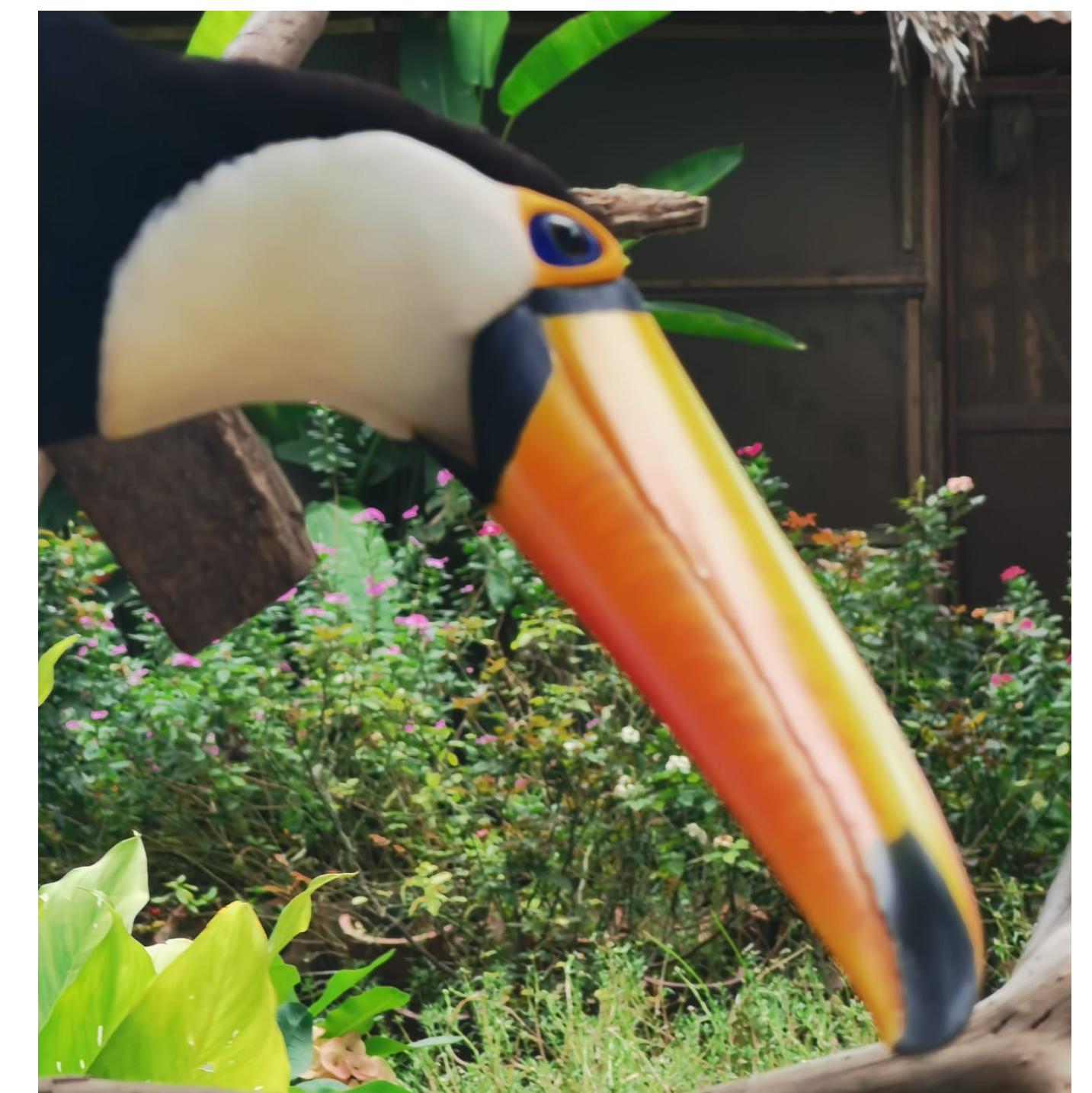
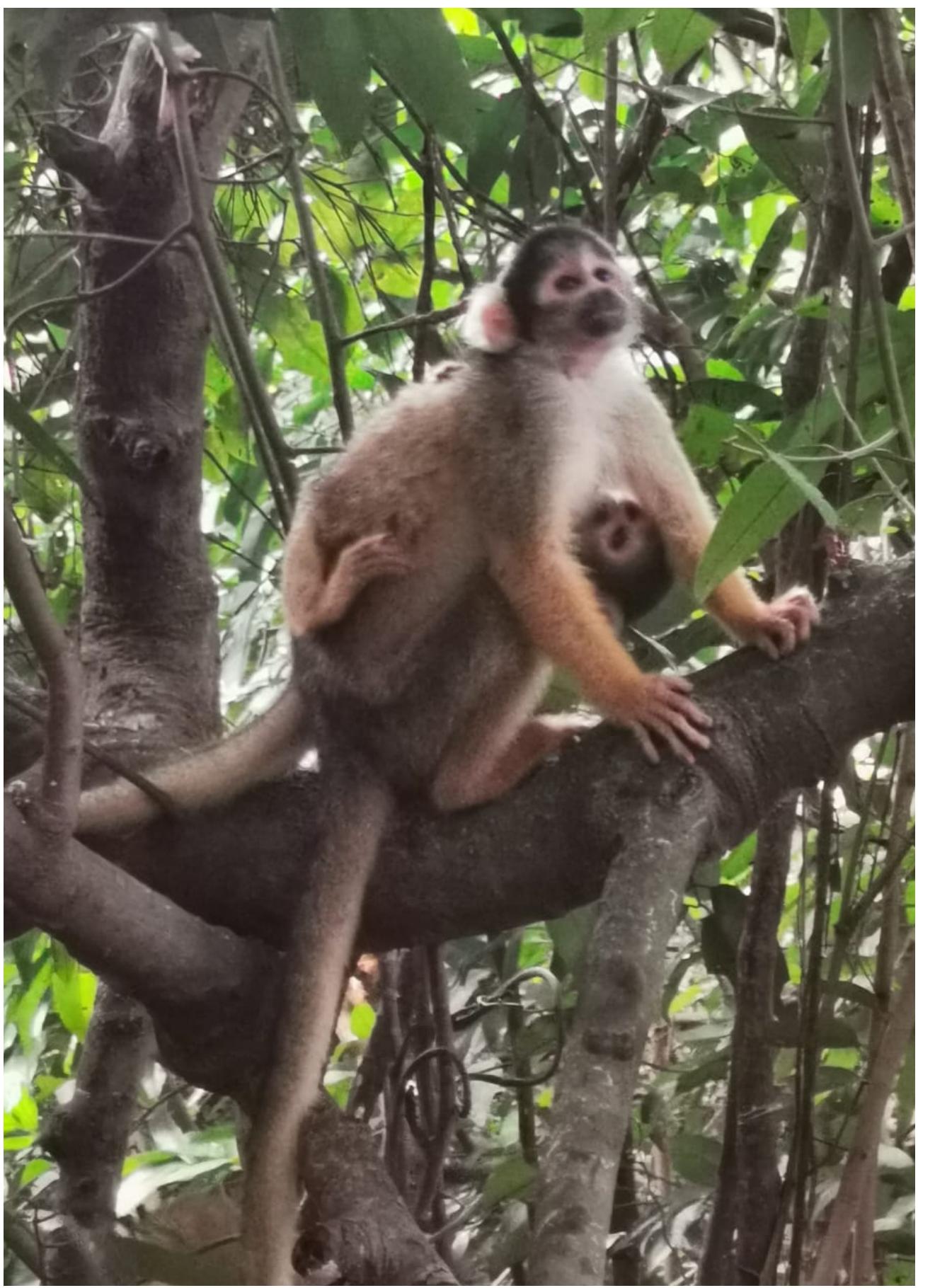
<https://github.com/smwindecker/R-para-ecologia>

Financia:



Invita:







Qué esperan aprender en este curso?

Con qué tipos de datos trabajan?

Introducción

aprender nueva terminología
identificar modelos apropiados para sus datos
entender supuestos de modelos comunes
saber dónde buscar ayuda

Introducción

- Como funciona R y Rstudio (uso de funciones, paquetes, ambiente)

Introducción

- Como funciona R y Rstudio (uso de funciones, paquetes, ambiente)
- Análisis de varianza y regresiones lineales

Introducción

- Como funciona R y Rstudio (uso de funciones, paquetes, ambiente)
- Análisis de varianza y regresiones lineales
- Supuestos y generalizaciones de los modelos lineales

Introducción

- Como funciona R y Rstudio (uso de funciones, paquetes, ambiente)
- Análisis de varianza y regresiones lineales
- Supuestos y generalizaciones de los modelos lineales
- Distribuciones de datos y funciones de enlace

- Como funciona R y Rstudio (uso de funciones, paquetes, ambiente)
- Análisis de varianza y regresiones lineales
- Supuestos y generalizaciones de los modelos lineales
- Distribuciones de datos y funciones de enlace
- Efectos fijos y efectos aleatorios

- Como funciona R y Rstudio (uso de funciones, paquetes, ambiente)
- Análisis de varianza y regresiones lineales
- Supuestos y generalizaciones de los modelos lineales
- Distribuciones de datos y funciones de enlace
- Efectos fijos y efectos aleatorios
- Diagnosticando y evaluando modelos lineales

- Como funciona R y Rstudio (uso de funciones, paquetes, ambiente)
- Análisis de varianza y regresiones lineales
- Supuestos y generalizaciones de los modelos lineales
- Distribuciones de datos y funciones de enlace
- Efectos fijos y efectos aleatorios
- Diagnosticando y evaluando modelos lineales
- Interpretación y presentación de resultados

- Como funciona R y Rstudio (uso de funciones, paquetes, ambiente)
- Análisis de varianza y regresiones lineales
- Supuestos y generalizaciones de los modelos lineales
- Distribuciones de datos y funciones de enlace
- Efectos fijos y efectos aleatorios
- Diagnosticando y evaluando modelos lineales
- Interpretación y presentación de resultados

* Análisis estadístico reproducible

I. Entendiendo los modelos

Cómo empezamos un proyecto de investigación?

Modelo conceptual

Formular la pregunta

Diseño experimental

Recolección de datos

Escribir y ajustar el modelo

Presentar resultados

Modelo conceptual

Formular la pregunta

Diseño experimental

Recolección de datos

Escribir y ajustar el modelo

Presentar resultados

Modelo conceptual

Formular la pregunta / escribir el modelo

Diseño experimental

Recolección de datos

Ajustar el modelo

Presentar resultados

Entendiendo los modelos

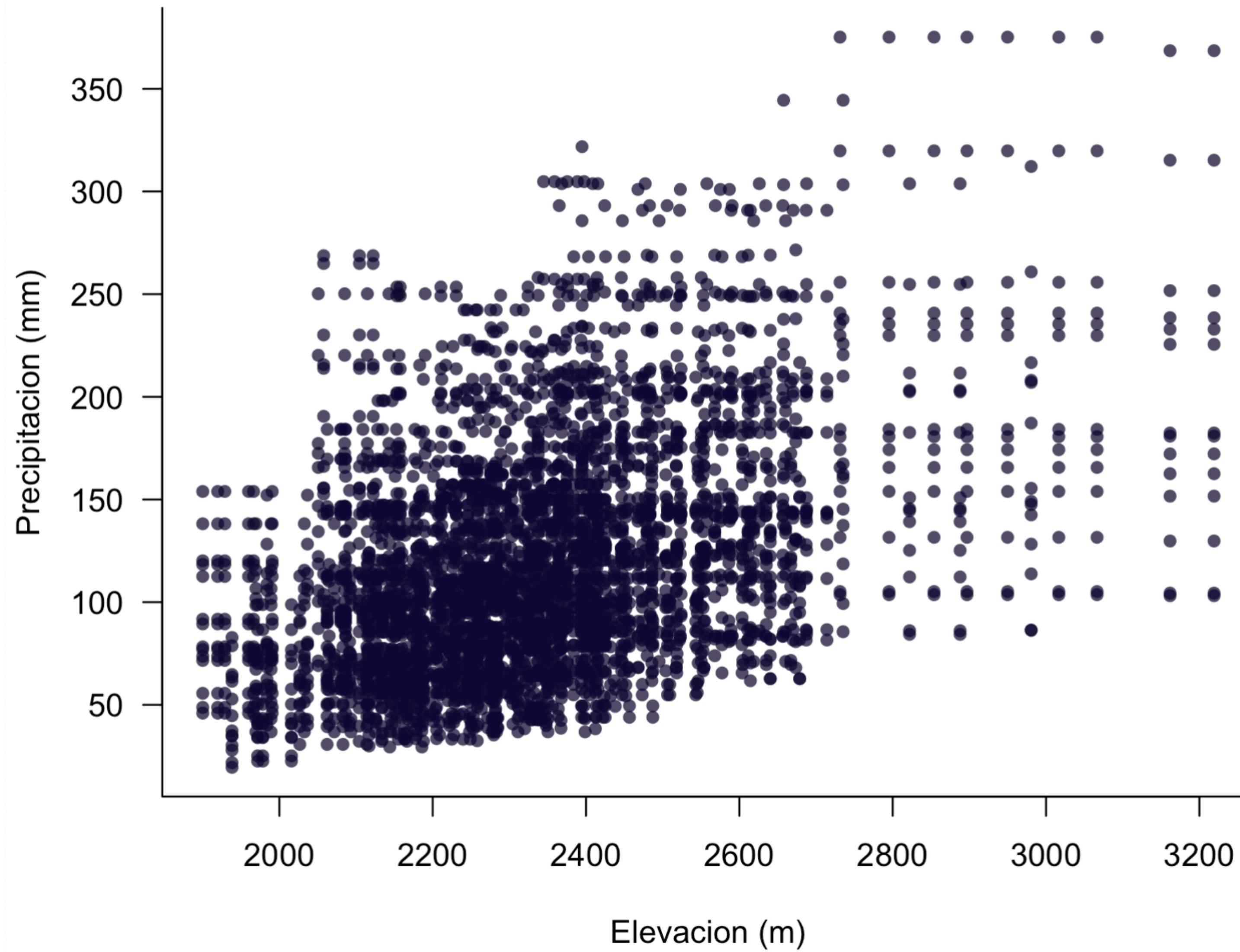
Una pregunta de investigación bien formulada debería poder representarse como modelo estadístico.

Entendiendo los modelos

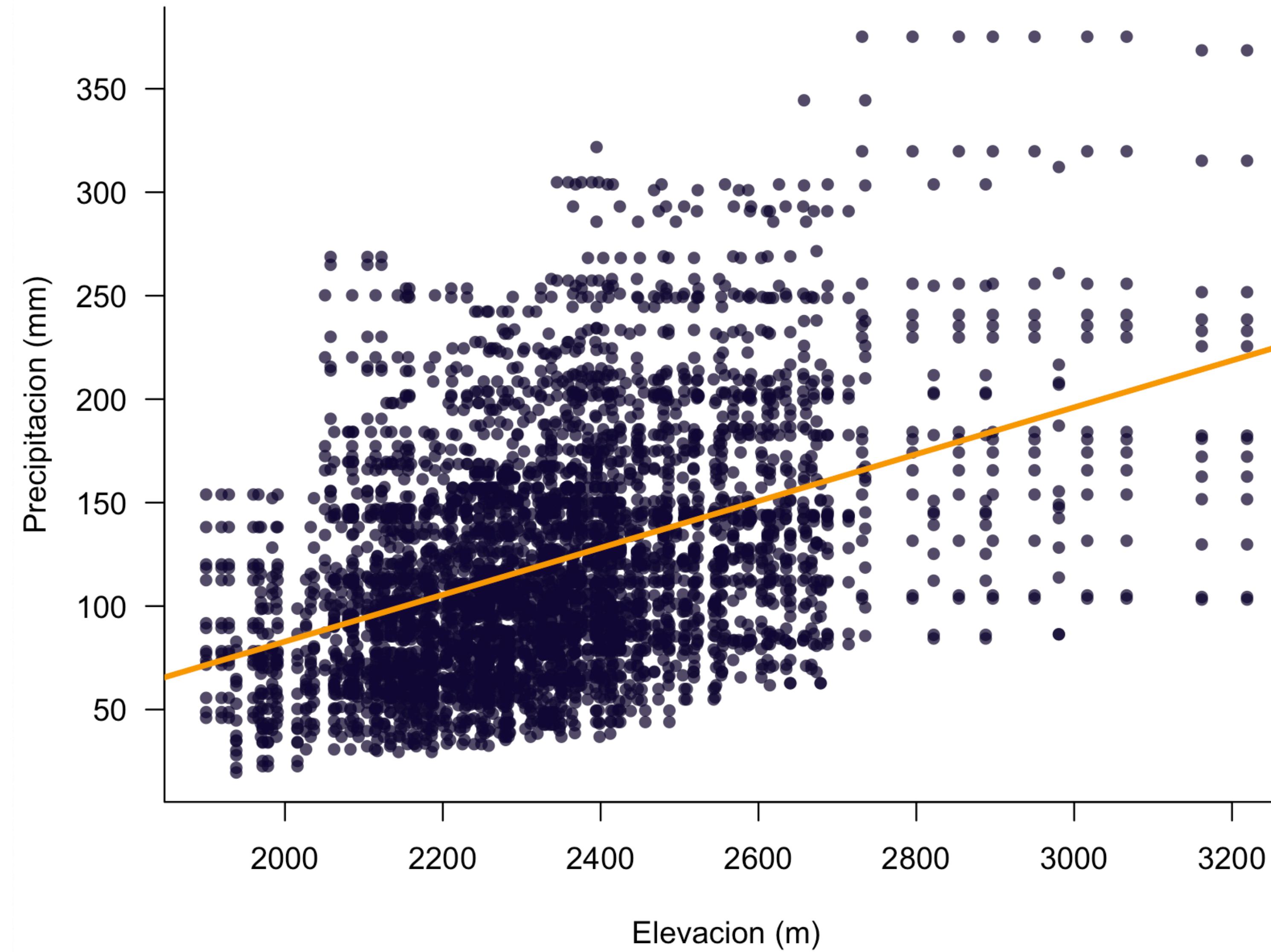
example with anova?

II. Modelos lineales

Modelos lineales



Modelos lineales

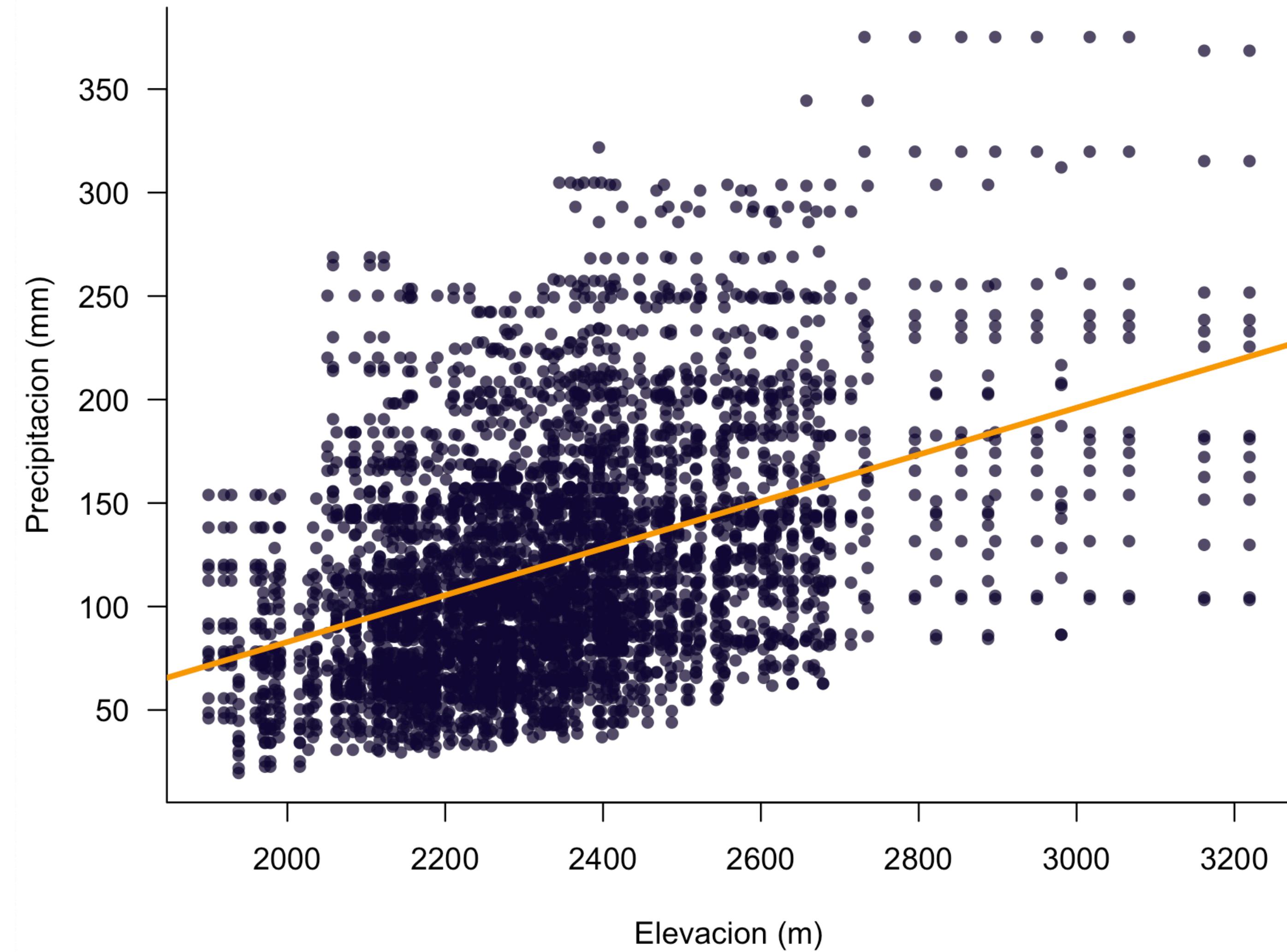


Modelos lineales

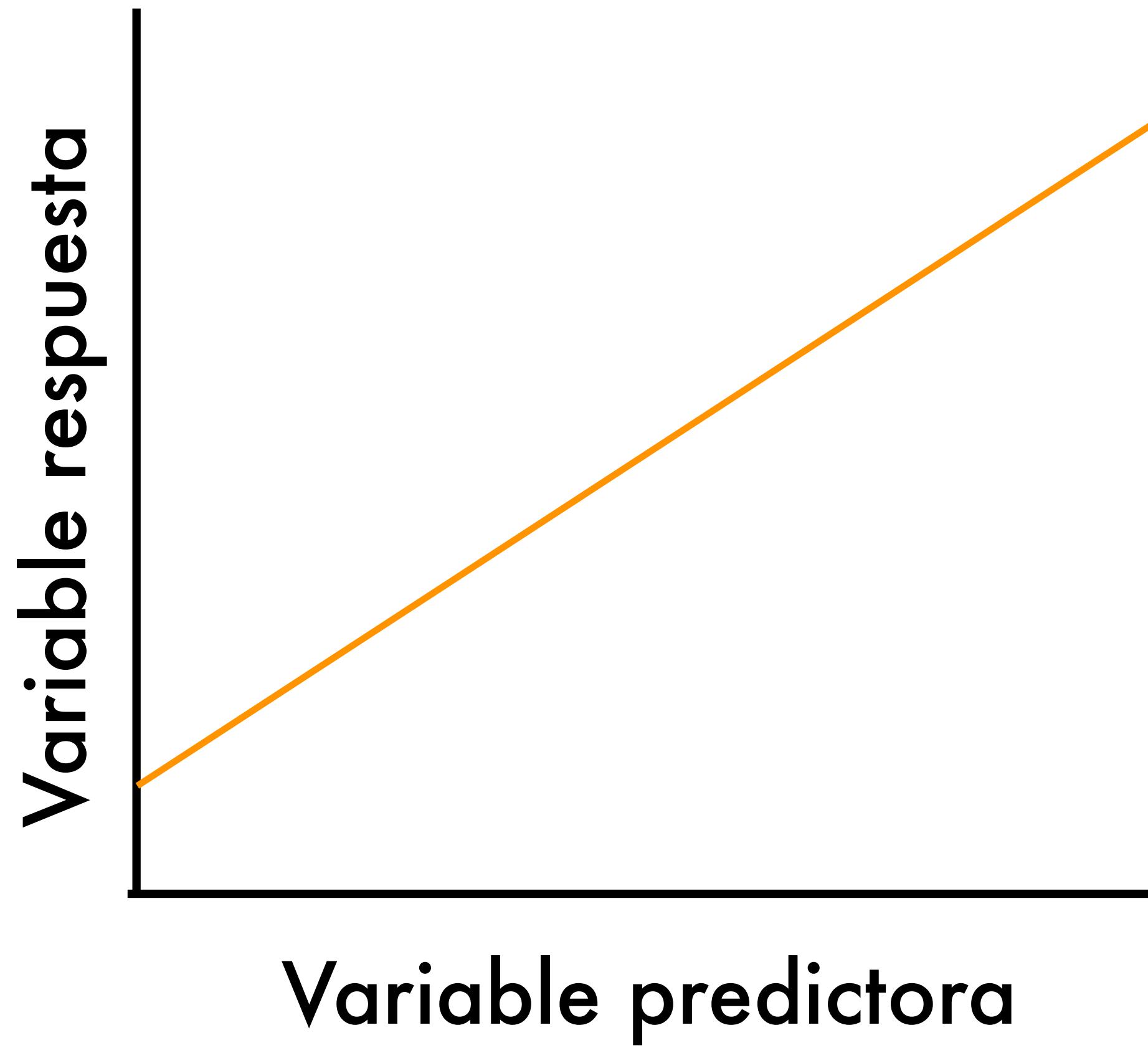
Qué caracteriza este ejemplo?

variable respuesta continua

variable predictoria continua

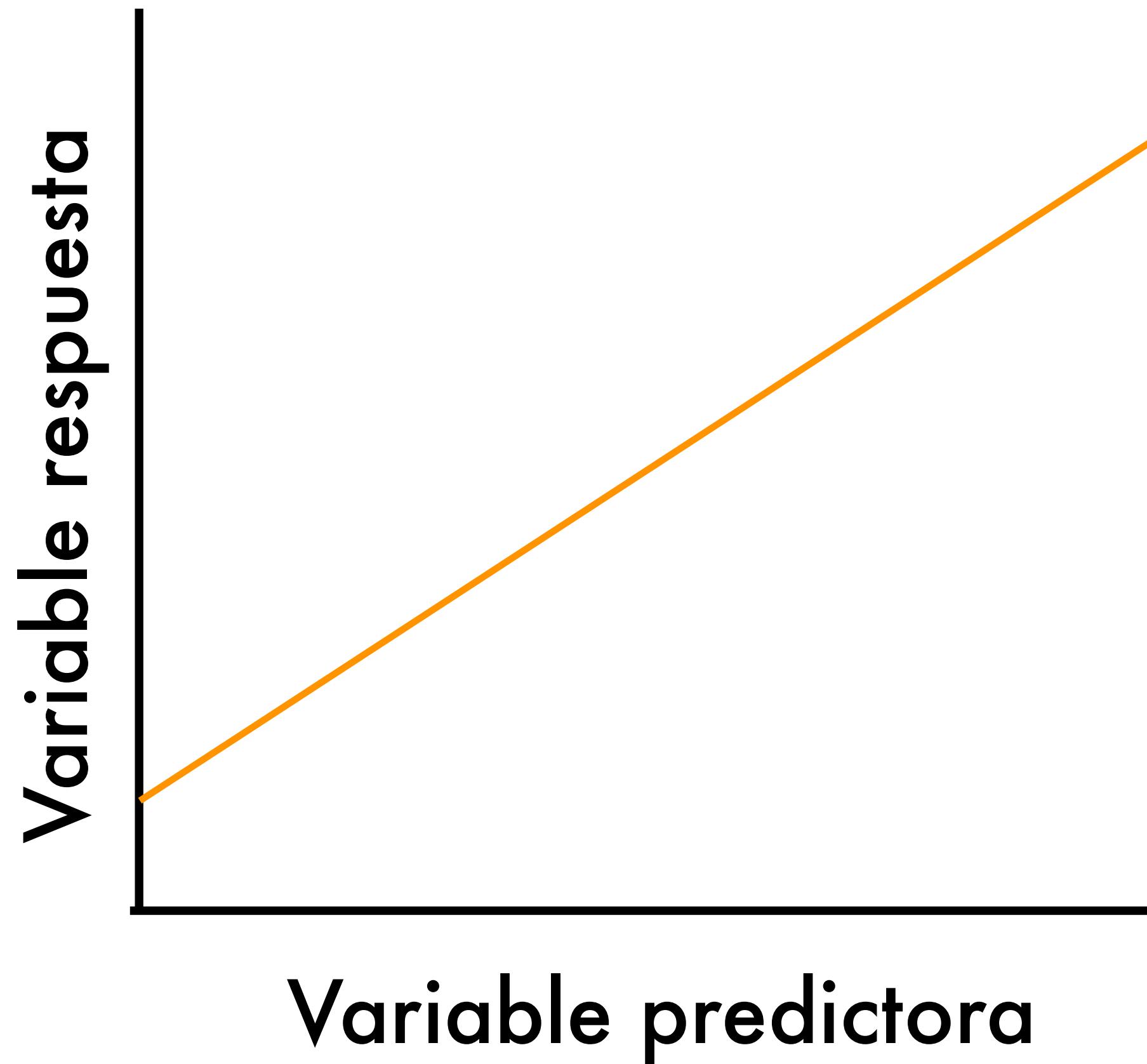


Modelos lineales



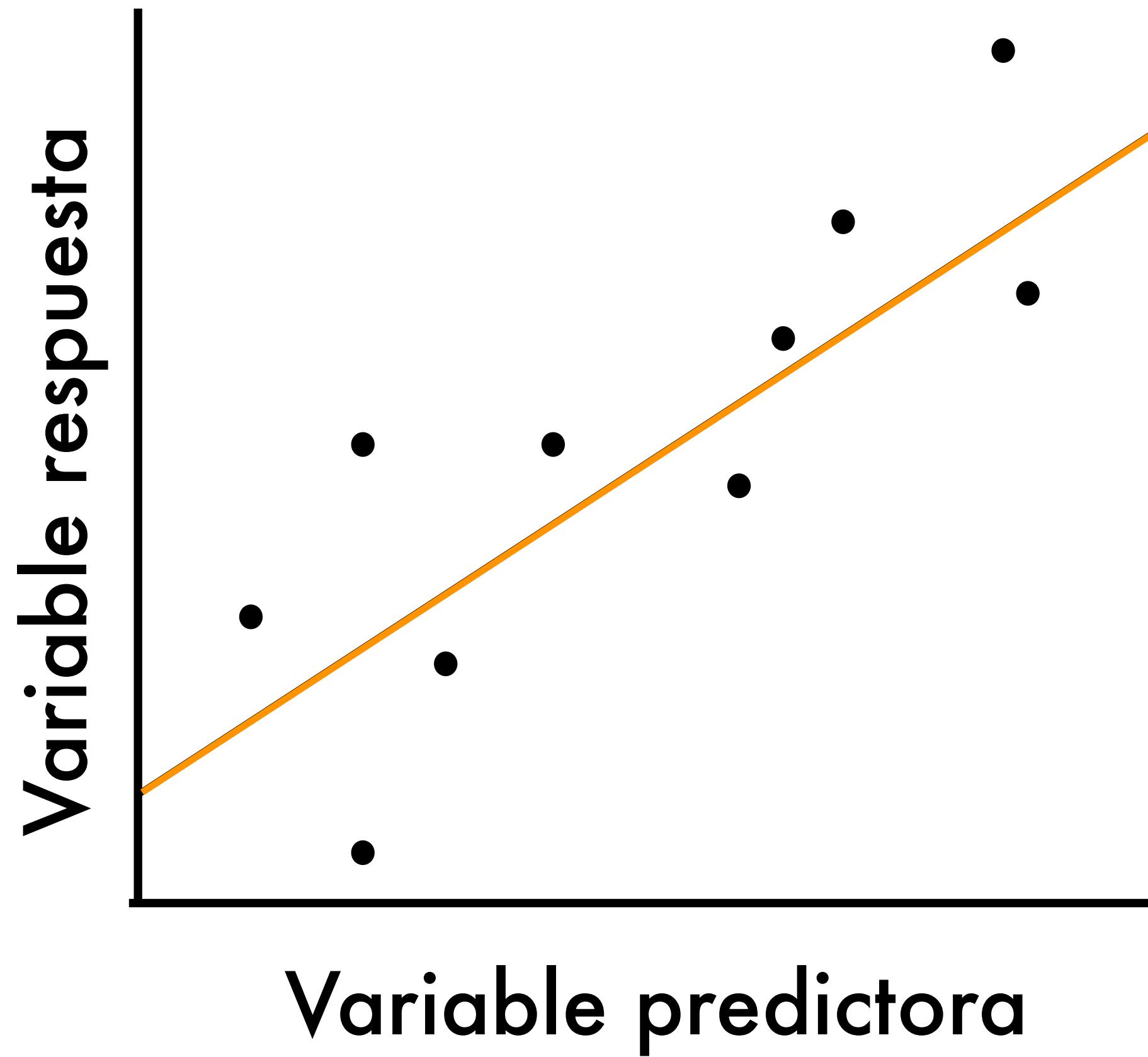
Modelos lineales

$$y_i = \alpha + \beta * x_i$$



Queremos estimar los
valores mas probables de
 α y β

Modelos lineales



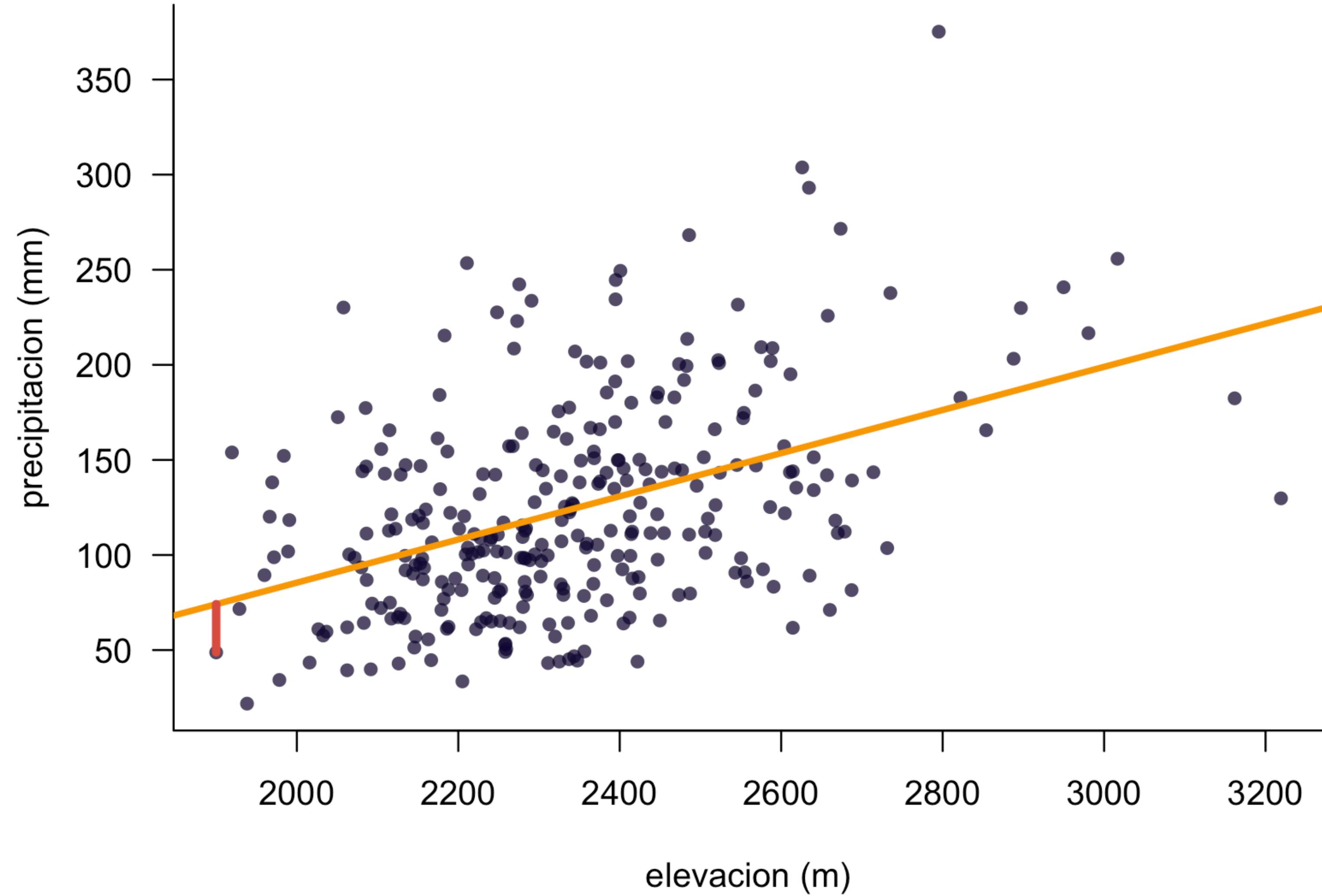
$$y_i = \alpha + \beta * x_i + \epsilon_i$$

Las diferencias entre las observaciones y la linea es el error residual ϵ_i

Como podemos
entender el ϵ ?

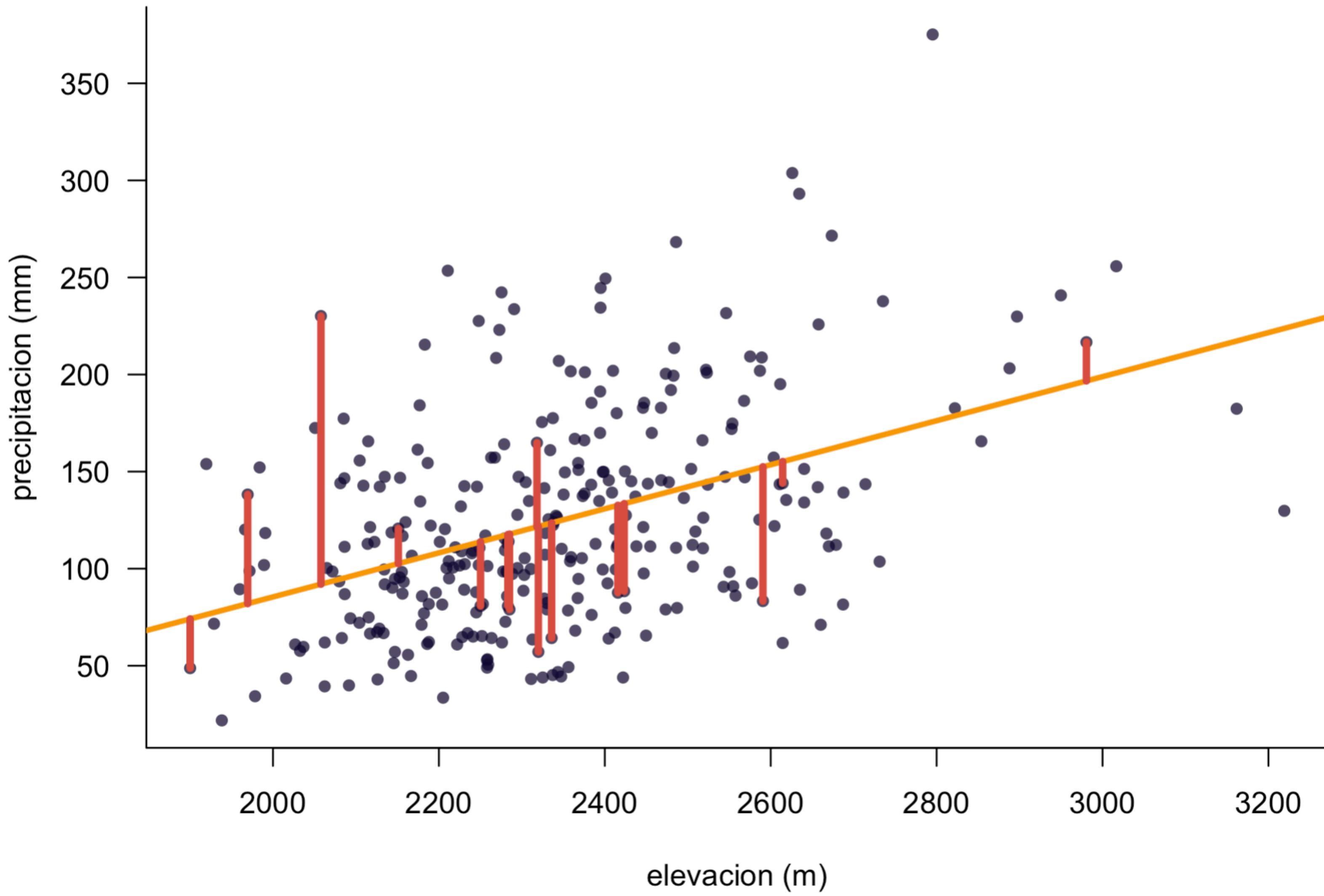
Modelos lineales

Como podemos entender el ϵ ?



Modelos lineales

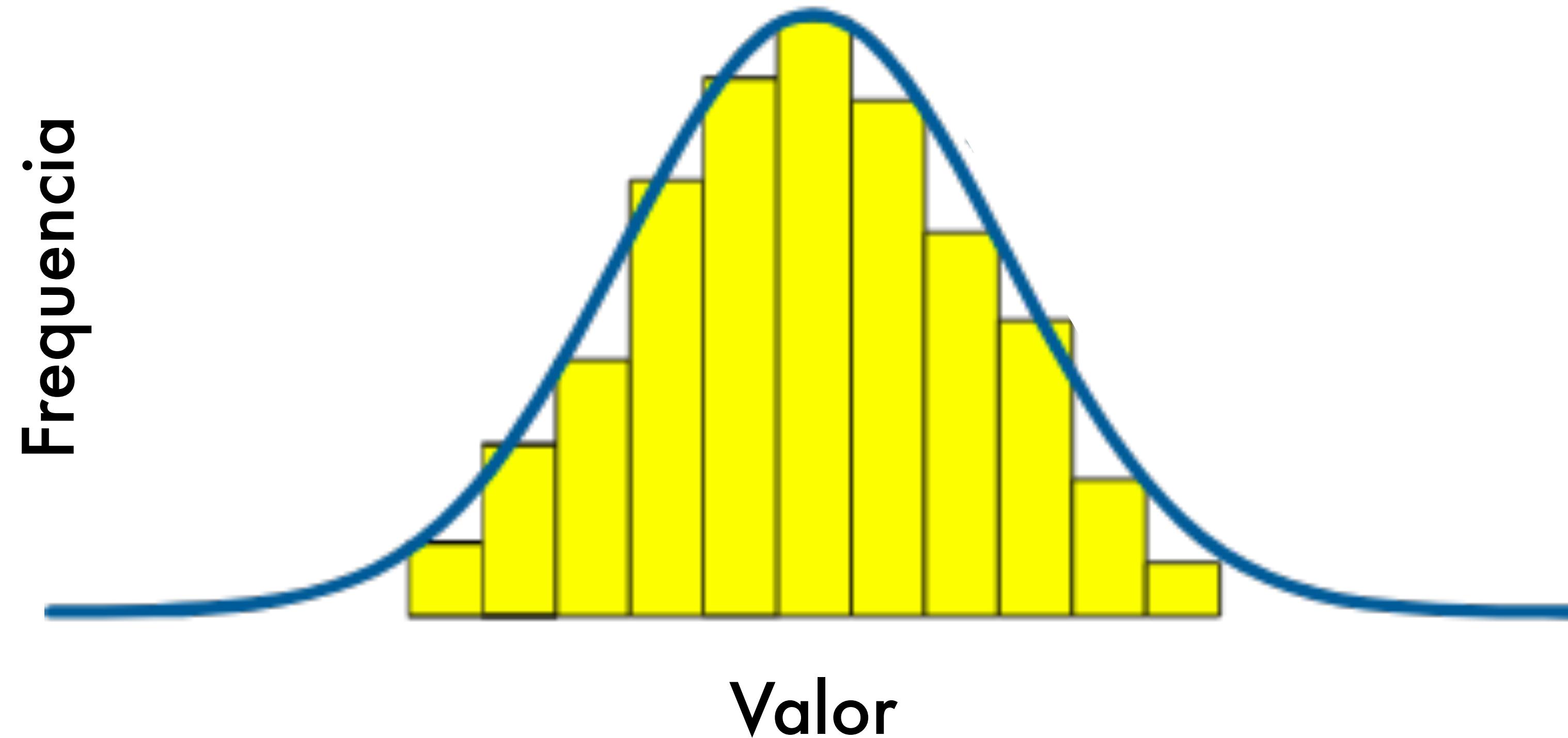
Como podemos entender el ϵ ?



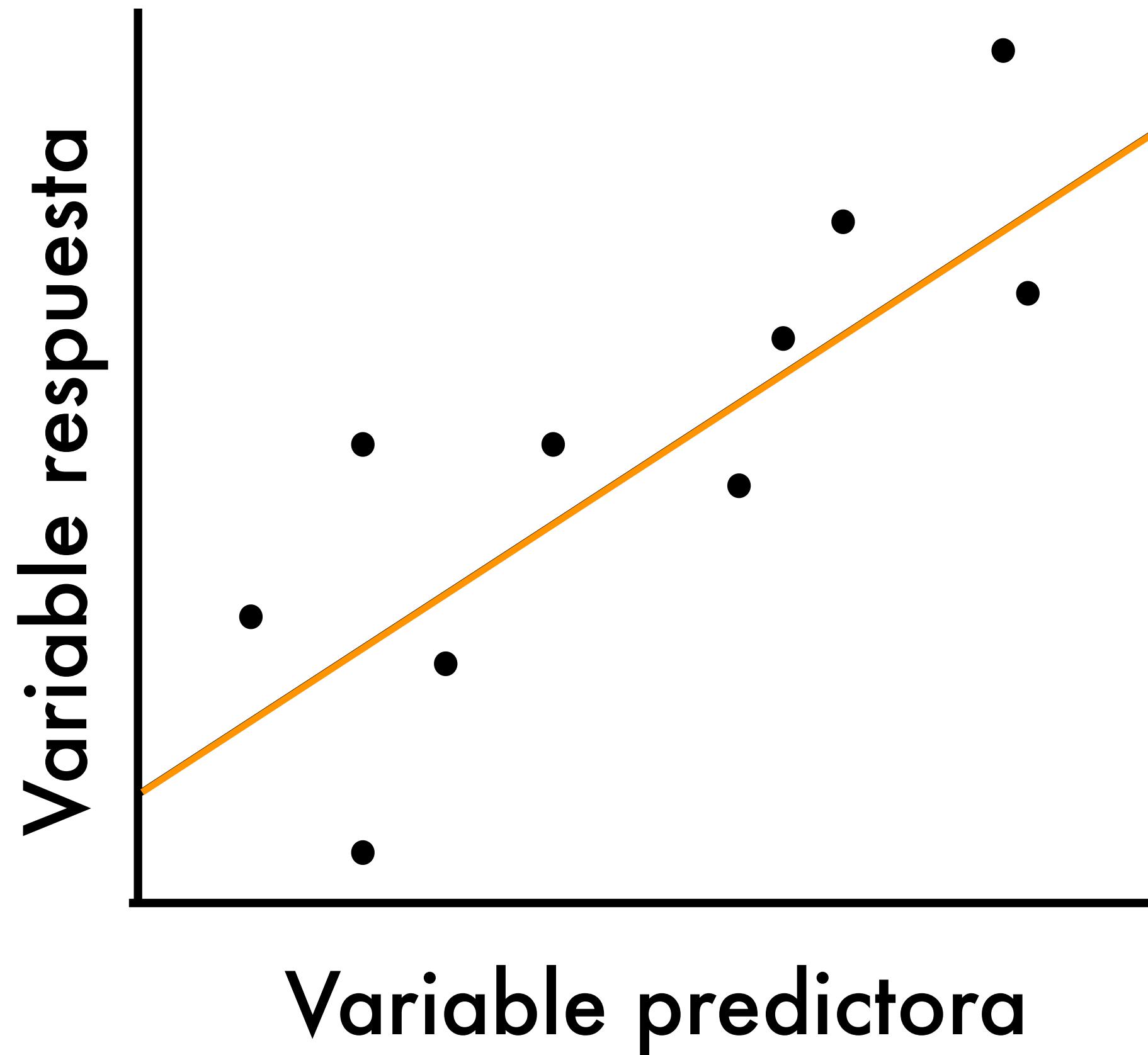
Supuesto de regresión lineal 1: Observaciones son independientes

- cada observación provee nueva información
- observaciones no-independientes dan menos información
- buen diseño de estudio puede eliminar problemas de independencia

Supuesto 2: Residuales se ajustan a una distribución normal

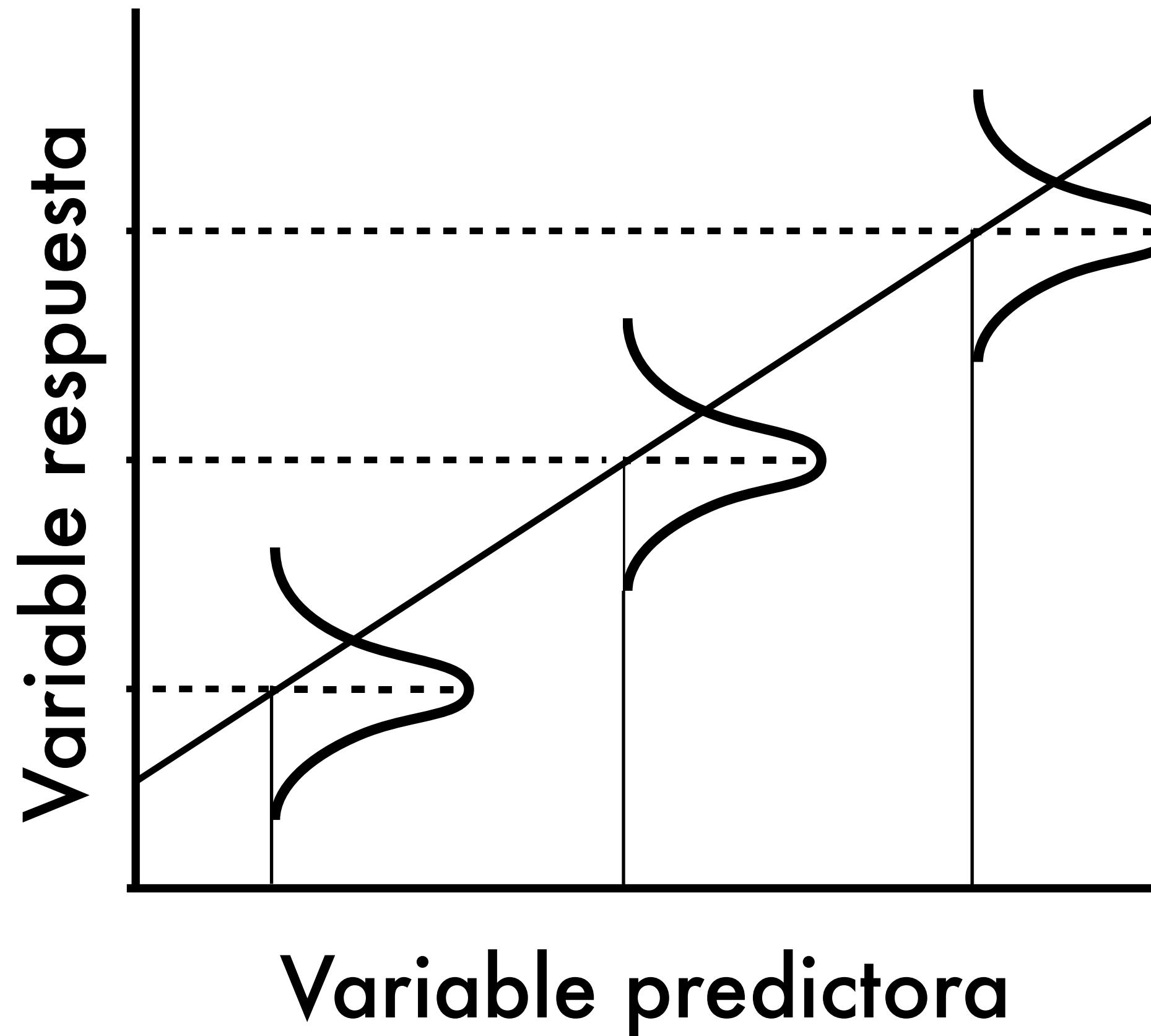


Supuesto 2: Residuales se ajustan a una distribución normal



$$y_i = \alpha + \beta * x_i + \epsilon_i$$

Supuesto 2: Residuales se ajustan a una distribución normal



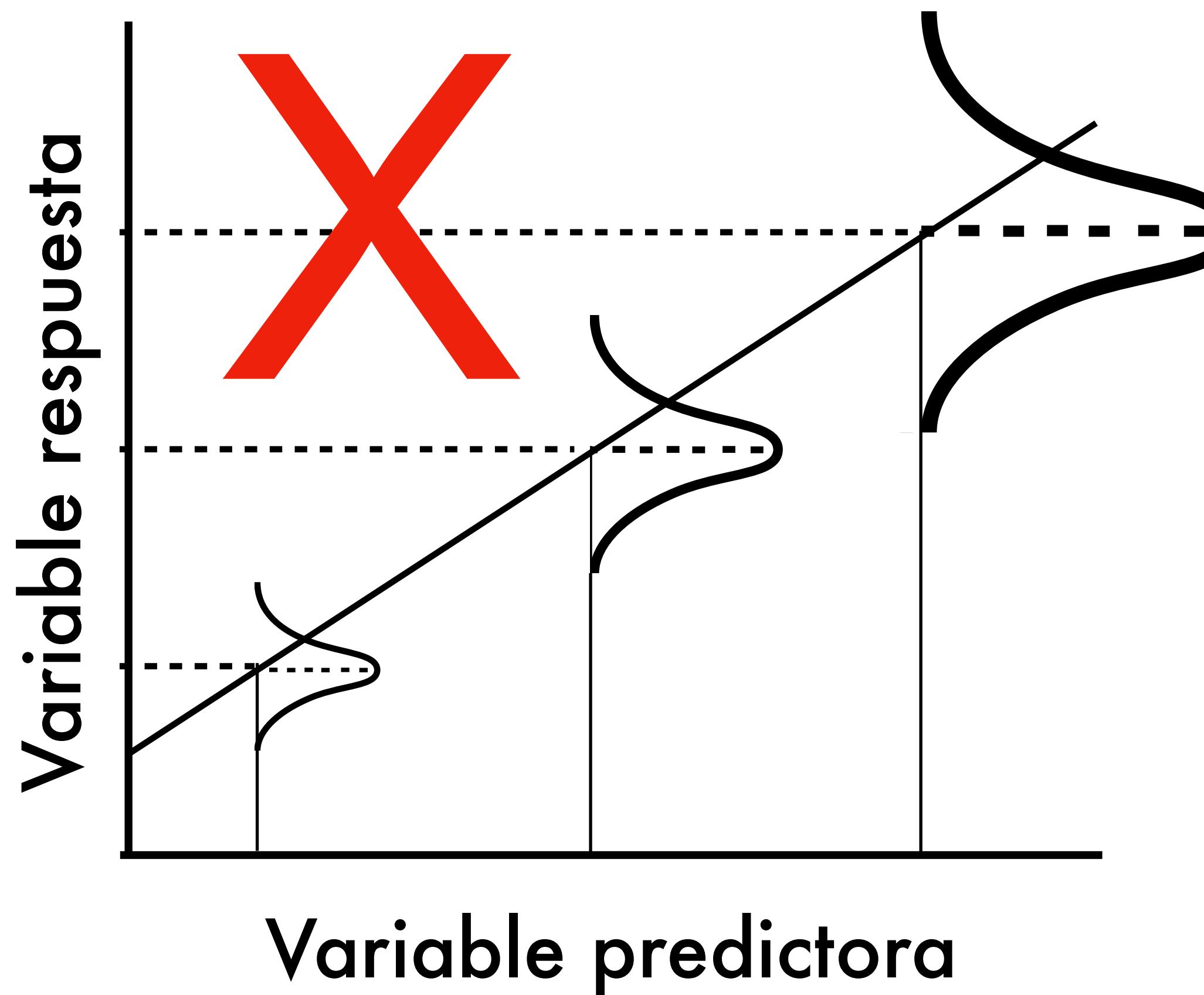
$$y_i = \alpha + \beta * x_i + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2)$$

Supuesto 2: Residuales se ajustan a una distribución normal

- Siempre hay que poner a prueba el supuesto, ya sea visualmente o con una prueba de ajuste (por ej. Kolmogorov-Smirnov)
- Si el supuesto no se cumple, se puede usar transformaciones, o modelos lineales generalizados

Supuesto 3: Homocedasticidad



- La varianza de los residuales debe ser constante a través de la variable predictora

Modelos lineales

practica en R

El manera de interpretar/ilustrar las estimaciones de un modelo depende de la pregunta.

El manera de interpretar/ilustrar las estimaciones de un modelo depende de tu pregunta.

Qué tan bien se ajusta el modelo a los datos?

usa R²

Modelos lineales

El manera de interpretar/ilustrar las estimaciones de un modelo depende de tu pregunta.

Qué tan bien se ajusta el modelo a los datos?

usa R²

Hay apoyo estadística para una asociación?

usa p-values

Modelos lineales

El manera de interpretar/ilustrar las estimaciones de un modelo depende de tu pregunta.

Qué tan bien se ajusta el modelo a los datos?

usa R²

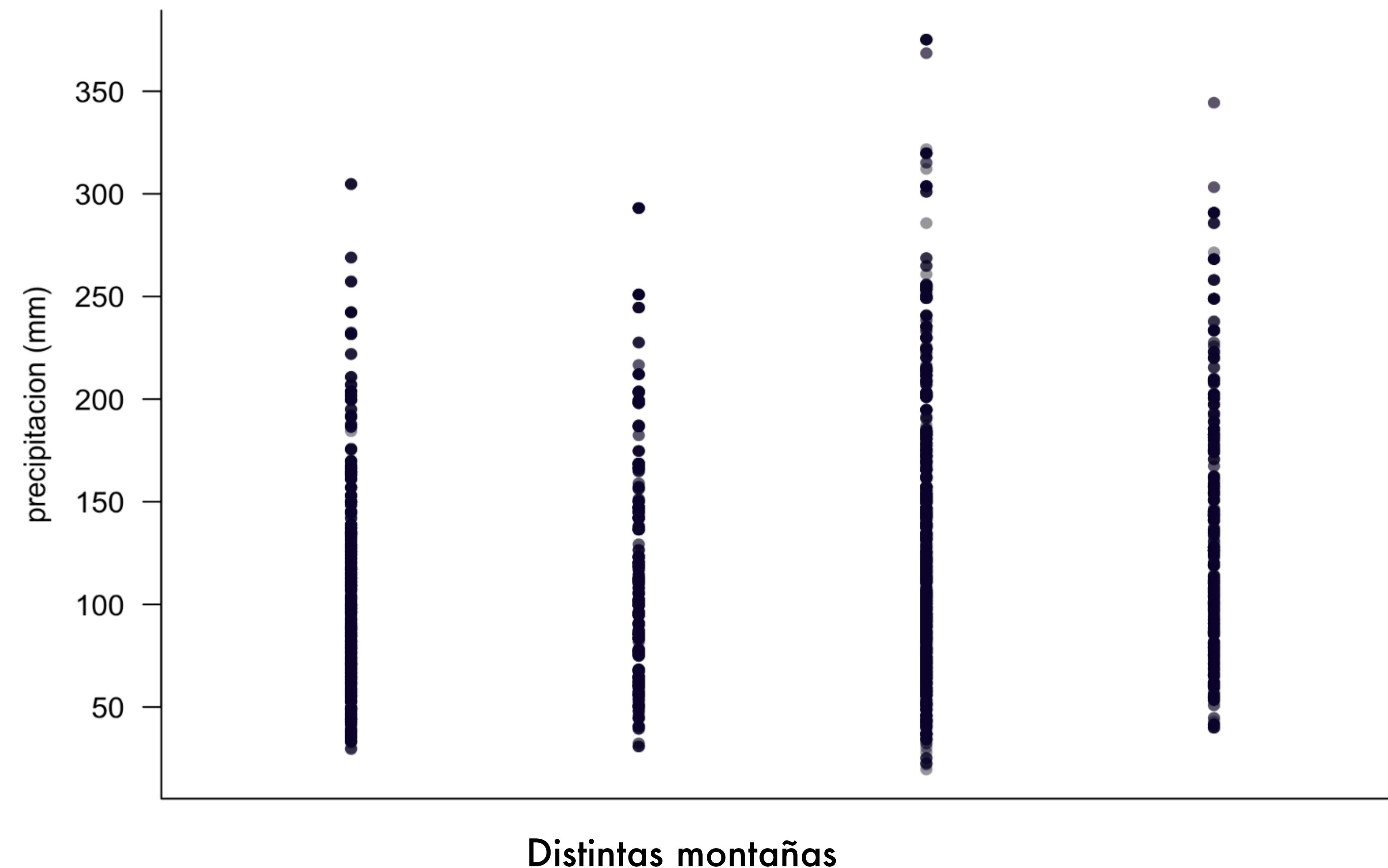
Hay apoyo estadística para una asociación?

usa p-values

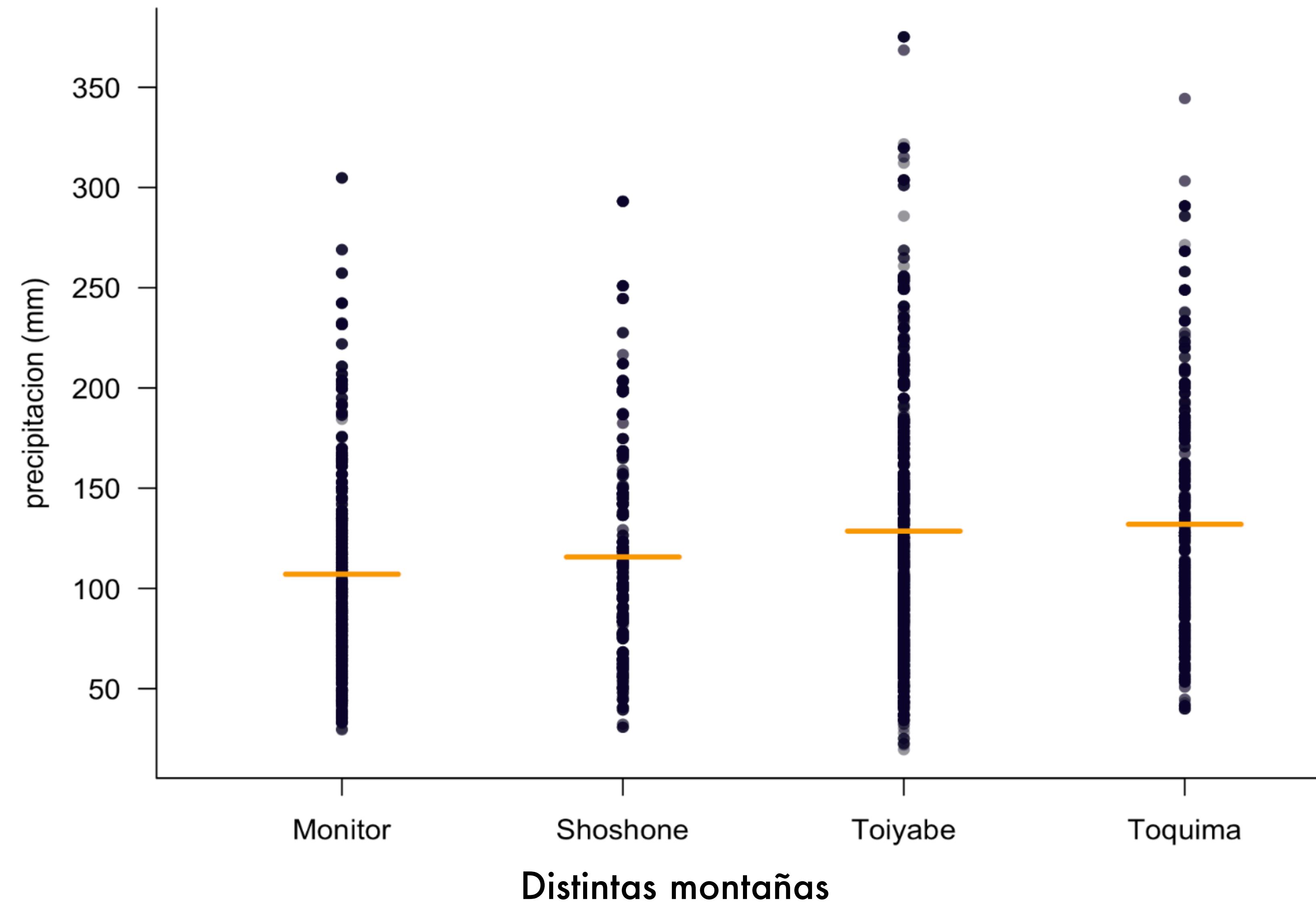
Una asociación estadísticamente significativa tiene sentido?

mira los coeficientes
(tamaño y signo)

ANOVA



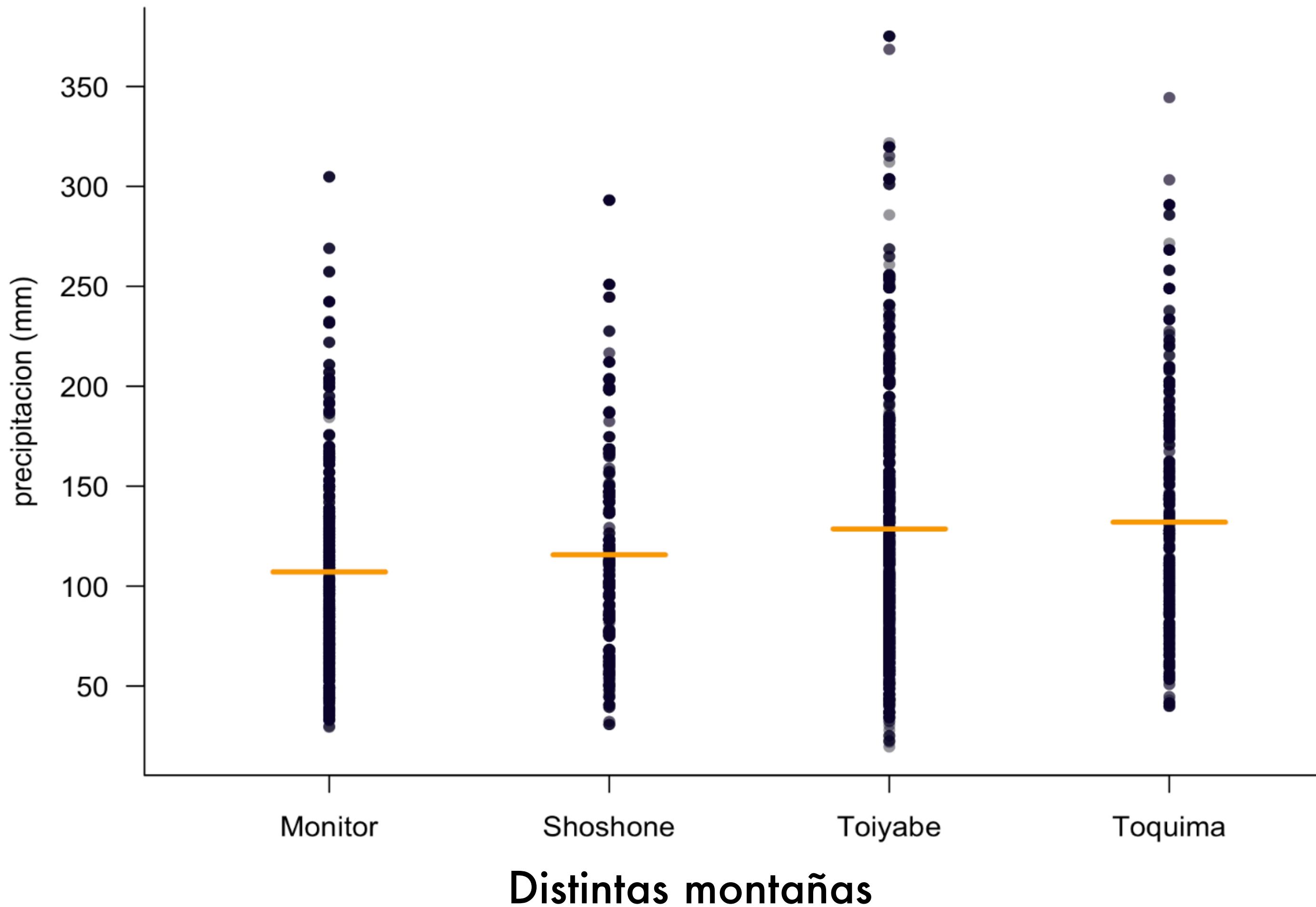
ANOVA



Qué caracteriza esta
ejemplo?

variable respuesta continua

variable predictoria
categorica



Supuestos: identicos a los de regresiones lineales

En este caso x es una categoría, entonces se puede entender como el intercepto que cambia dado al variable.

En otras palabras, diferencias en la media entre distintas categorías

$$y_i = \alpha + \beta * x_i + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2)$$

t-test

t-test: un ANOVA caso especial, con solo dos niveles

ANOVA

practica en R

Model matrices

Para usar variables categóricas o multiples variables continuas, es necesario entender la matemática.

$$y_i = \alpha + \beta^T x_i + \epsilon_i$$

$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}; \quad x_i = \begin{pmatrix} x_{i,1} \\ \vdots \\ x_{i,k} \end{pmatrix}$$

$$\beta^T = (\beta_1 \quad \dots \quad \beta_k)$$

Model matrices

Como programamos esto?

$$\beta^T = (\beta_1 \quad \beta_2 \quad \beta_3)$$

Transformamos las observaciones
para clasificarlas en 0 o 1:

$$x_i = \begin{pmatrix} \\ 0 \\ 1 \end{pmatrix}$$

$$\beta^T x_i = 0 + \beta_2 + 0$$

Model matrices

R te lo hace:

```
model.matrix(~ discrete_predictor)
```

```
##      (Intercept) mountain_rangeShoshone mountain_rangeToiyabe
## [1,]          1                  0                  0
## [2,]          1                  1                  0
## [3,]          1                  0                  0
## [4,]          1                  0                  0
## [5,]          1                  0                  1
##      mountain_rangeToquima
## [1,]          0
## [2,]          0
## [3,]          0
## [4,]          1
## [5,]
```

Model matrices

Con mas variables:

$$\beta^T = (\beta_1 \ \ \beta_2 \ \ \beta_3)$$

$$\begin{array}{c} \begin{array}{|c|c|c|} \hline & \text{Red} & \text{Green} & \text{Blue} \\ \hline & \text{Red} & \text{Green} & \text{Blue} \\ \hline & \text{Red} & \text{Green} & \text{Blue} \\ \hline \end{array} & \times & \begin{array}{|c|} \hline \text{a} \\ \hline \text{b} \\ \hline \text{c} \\ \hline \end{array} & = & \begin{array}{|c|} \hline \text{a} \\ \hline \text{Red} \\ \hline \end{array} & + & \begin{array}{|c|} \hline \text{b} \\ \hline \text{Green} \\ \hline \end{array} & + & \begin{array}{|c|} \hline \text{c} \\ \hline \text{Blue} \\ \hline \end{array} & = & \begin{array}{|c|} \hline \text{Pink} \\ \hline \text{Pink} \\ \hline \end{array} \end{array}$$

$$x_i = \begin{pmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{pmatrix}$$

$$\beta^T x_i = \beta_1 x_{i,1} + \beta_2 x_{i,2} + \beta_3 x_{i,3}$$

Model matrices

Con mas variables:

```
model.matrix( ~ predictor1 + predictor2 + predictor3)
```

```
##          (Intercept) predictor1 predictor2 predictor3
## [1,]           1       2285.1      78.99     20.8
## [2,]           1       2304.6      67.31     17.1
## [3,]           1       2330.1      64.27     16.7
## [4,]           1       2589.2     144.02     15.8
## [5,]           1       3016.6     235.46      6.2
```

Multiples variables predictoras

Es necesaria pensar en la escala de los variables para que puedas comparar el tamaño del efecto.

```
# standardise continuous predictors  
predictors_std <- scale(predictors)
```

```
##      predictor1 predictor2 predictor3  
## [1,] -0.7065051 -0.5328227  1.00678496  
## [2,] -0.6438887 -0.6923145  0.32702139  
## [3,] -0.5620058 -0.7338261  0.25353344  
## [4,]  0.2699889  0.3551696  0.08818554  
## [5,]  1.6424108  1.6037937 -1.67552534
```

Multiples variables predictoras

Se puede integrar variables predictoras continuas y categóricas en el mismo modelo:

```
mod <- lm(response ~ continuous1 + continuous2 + discrete)
```

```
##          (Intercept) continuous1 continuous2 discrete1 discrete2 discrete3
## [1,]           2285.1       78.99         0         0         0
## [2,]           2304.6       67.31         1         0         0
## [3,]           2330.1       64.27         0         0         0
## [4,]           2589.2      144.02         0         0         1
## [5,]           3016.6      235.46         0         1         0
```

Multiples variables predictoras

Suposiciones son lo mismo, ademas suponemos que los variables son relativamente independiente.

Si no lo son, lo llamamos "multicollinearity".

Es difícil para el modelo de distinguir entre los efectos de cada variable.

Multiples variables predictoras

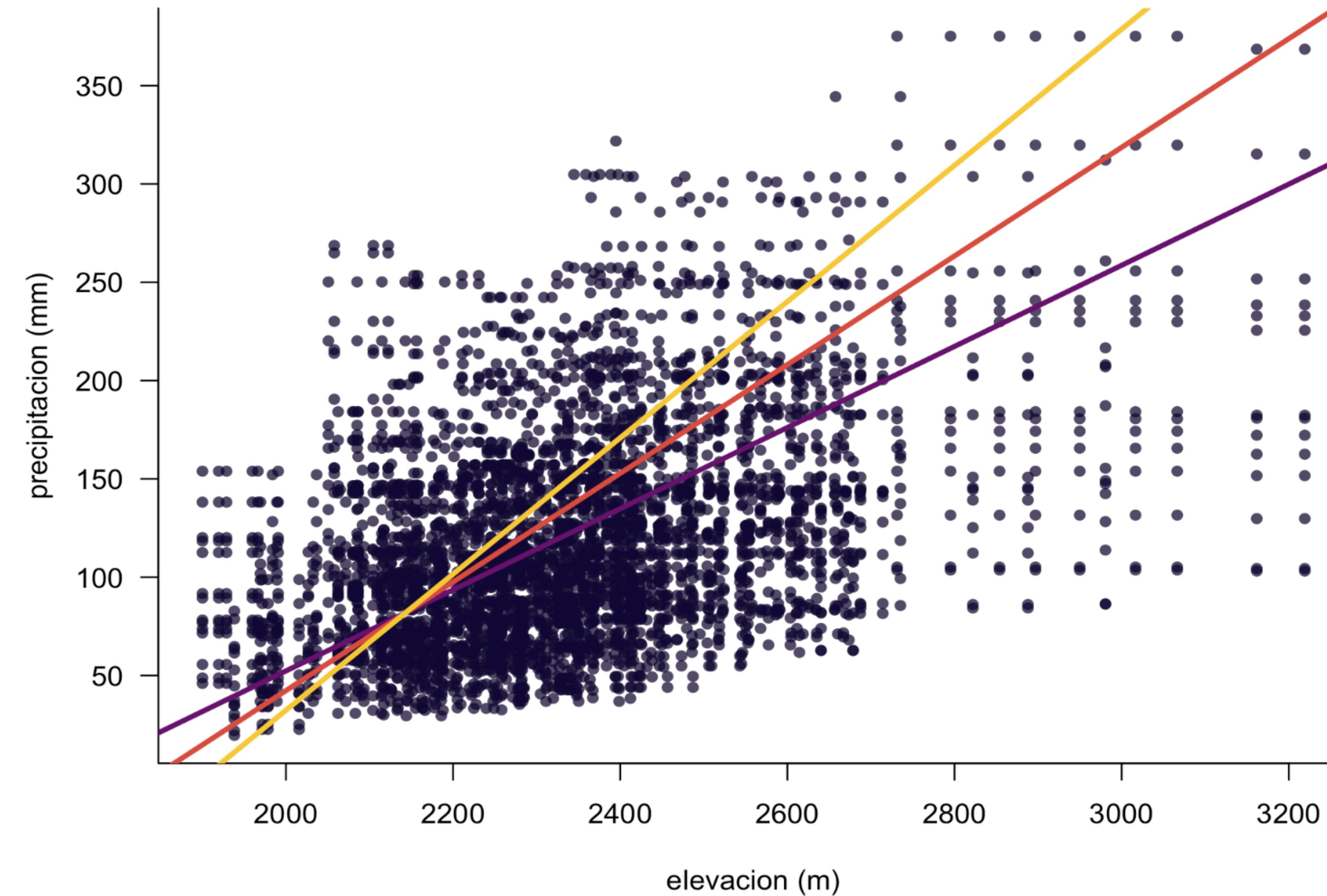
Suposiciones son lo mismo, ademas suponemos que los variables son relativamente independiente.

Si no lo son, lo llamamos "multicollinearity".

Es difícil para el modelo de distinguir entre los efectos de cada variable.

Se arregla por elegir variables con cuidado, y sacar variables colineales.

Interacciones



Interacciones

practica en R

Interacciones

- Es difícil interpretar coeficientes

```
mod <- lm(precipitation ~ elevation * mountain_range)
summary(mod)
```

```
##  
## Call:  
## lm(formula = precipitation ~ elevation * mountain_range)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -120.75  -35.97   -9.48   30.03  202.73  
##  
## Coefficients:  
##  
## (Intercept)                         Estimate Std. Error t value Pr(>|t|)  
## elevation                                1.364e-01  1.366e-02  9.981  < 2e-16  
## mountain_rangeShoshone                  1.842e+01  4.067e+01  0.453  0.65053  
## mountain_rangeToiyabe                   9.473e+01  3.346e+01  2.832  0.00465  
## mountain_rangeToquima                  1.798e+01  4.192e+01  0.429  0.66795  
## elevation:mountain_rangeShoshone -3.521e-03  1.747e-02 -0.202  0.84028  
## elevation:mountain_rangeToiyabe -3.089e-02  1.435e-02 -2.152  0.03145  
## elevation:mountain_rangeToquima -2.771e-03  1.767e-02 -0.157  0.87538  
##
```

Interacciones

- Es difícil interpretar coeficientes
- El efecto de cada parameter depende del valor estimado de otros parameters

```
mod <- lm(precipitation ~ elevation * mountain_range)
summary(mod)
```

```
##
## Call:
## lm(formula = precipitation ~ elevation * mountain_range)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -120.75  -35.97   -9.48   30.03  202.73
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)                 -2.104e+02  3.185e+01 -6.607  4.4e-11
## elevation                      1.364e-01  1.366e-02  9.981 < 2e-16
## mountain_rangeShoshone        1.842e+01  4.067e+01  0.453  0.65053
## mountain_rangeToiyabe         9.473e+01  3.346e+01  2.832  0.00465
## mountain_rangeToquima         1.798e+01  4.192e+01  0.429  0.66795
## elevation:mountain_rangeShoshone -3.521e-03  1.747e-02 -0.202  0.84028
## elevation:mountain_rangeToiyabe -3.089e-02  1.435e-02 -2.152  0.03145
## elevation:mountain_rangeToquima -2.771e-03  1.767e-02 -0.157  0.87538
##
```

Interacciones

- Es difícil interpretar coeficientes
- El efecto de cada parameter depende del valor estimado de otros parameters
- Particularmente difícil si los dos variables son continuos

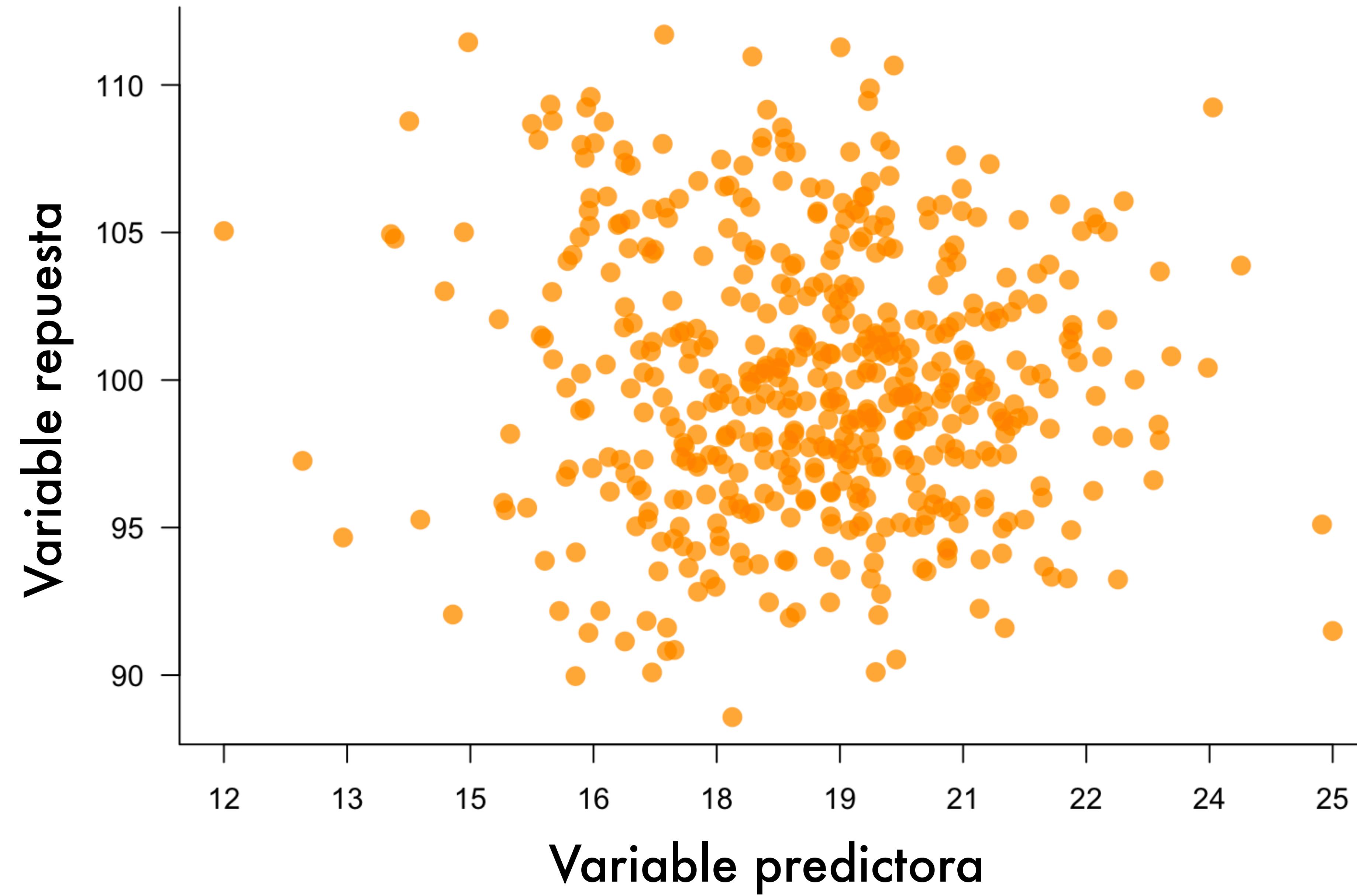
```
mod <- lm(precipitation ~ elevation * mountain_range)
summary(mod)
```

```
##  
## Call:  
## lm(formula = precipitation ~ elevation * mountain_range)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -120.75  -35.97   -9.48   30.03  202.73  
##  
## Coefficients:  
##  
## (Intercept)                         Estimate Std. Error t value Pr(>|t|)  
## elevation                                1.364e-01  1.366e-02  9.981  < 2e-16  
## mountain_rangeShoshone                  1.842e+01  4.067e+01  0.453  0.65053  
## mountain_rangeToiyabe                   9.473e+01  3.346e+01  2.832  0.00465  
## mountain_rangeToquima                  1.798e+01  4.192e+01  0.429  0.66795  
## elevation:mountain_rangeShoshone    -3.521e-03  1.747e-02 -0.202  0.84028  
## elevation:mountain_rangeToiyabe     -3.089e-02  1.435e-02 -2.152  0.03145  
## elevation:mountain_rangeToquima   -2.771e-03  1.767e-02 -0.157  0.87538  
##
```

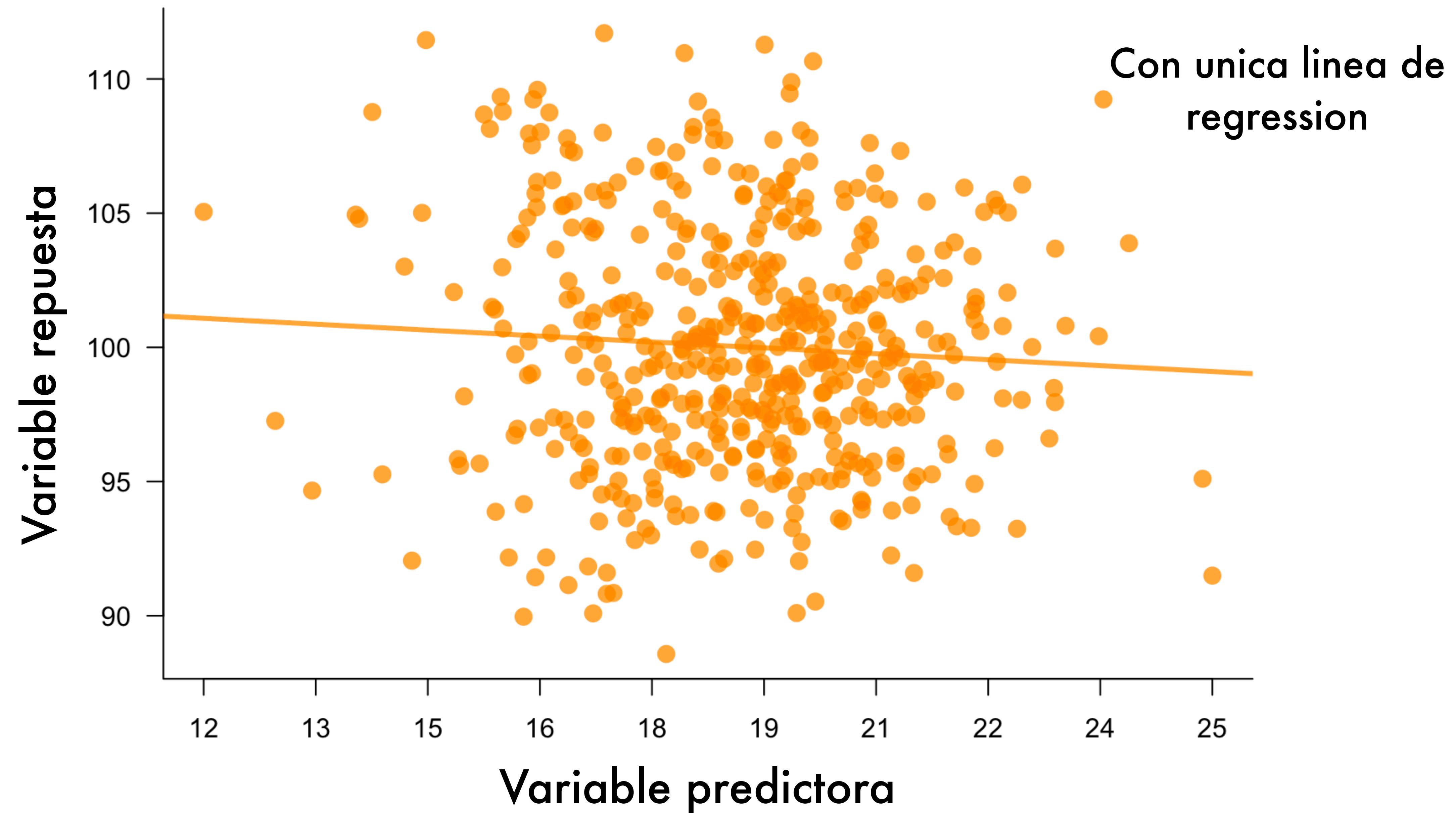
III. Modelos mixtos



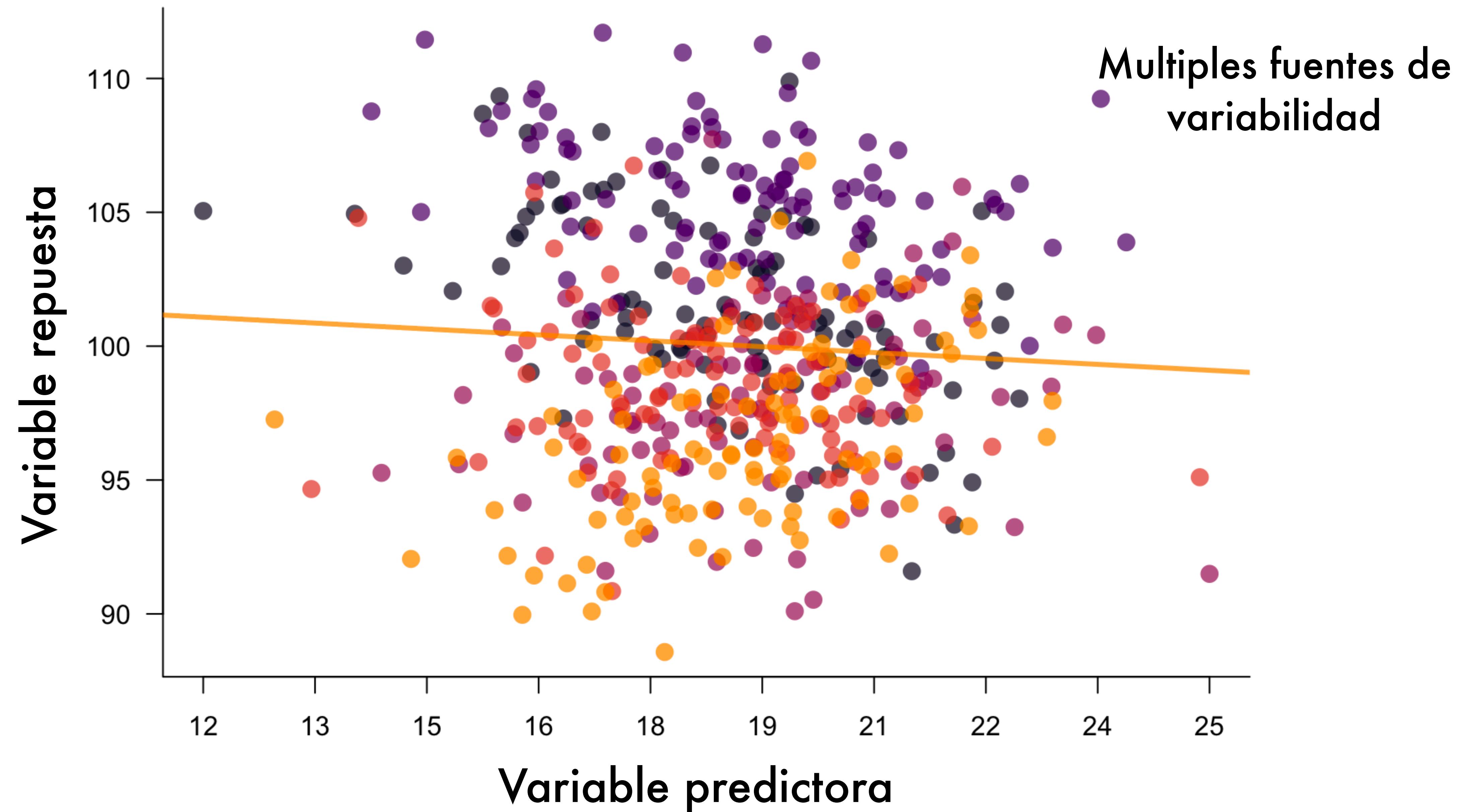
Modelos mixtos



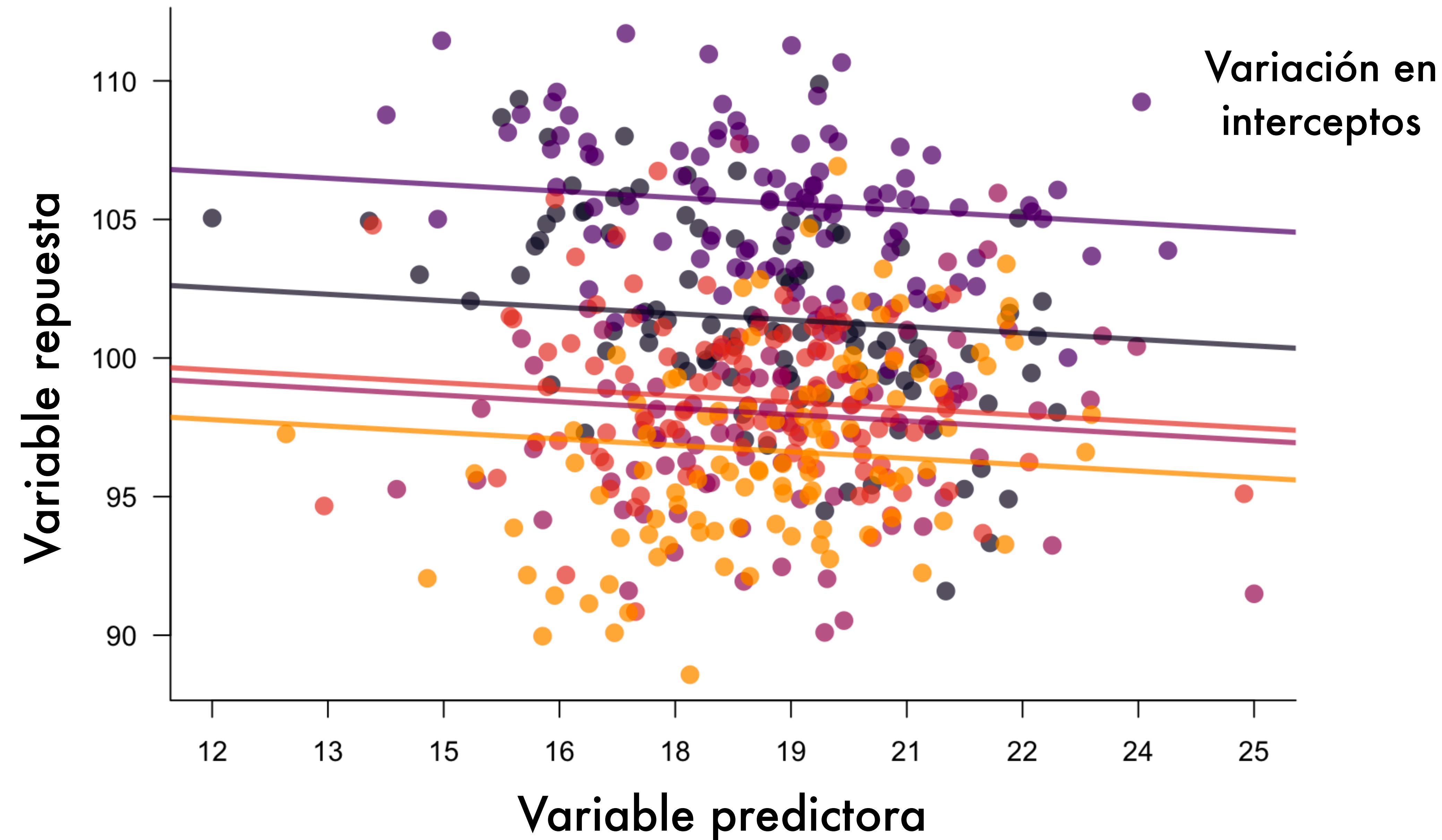
Modelos mixtos



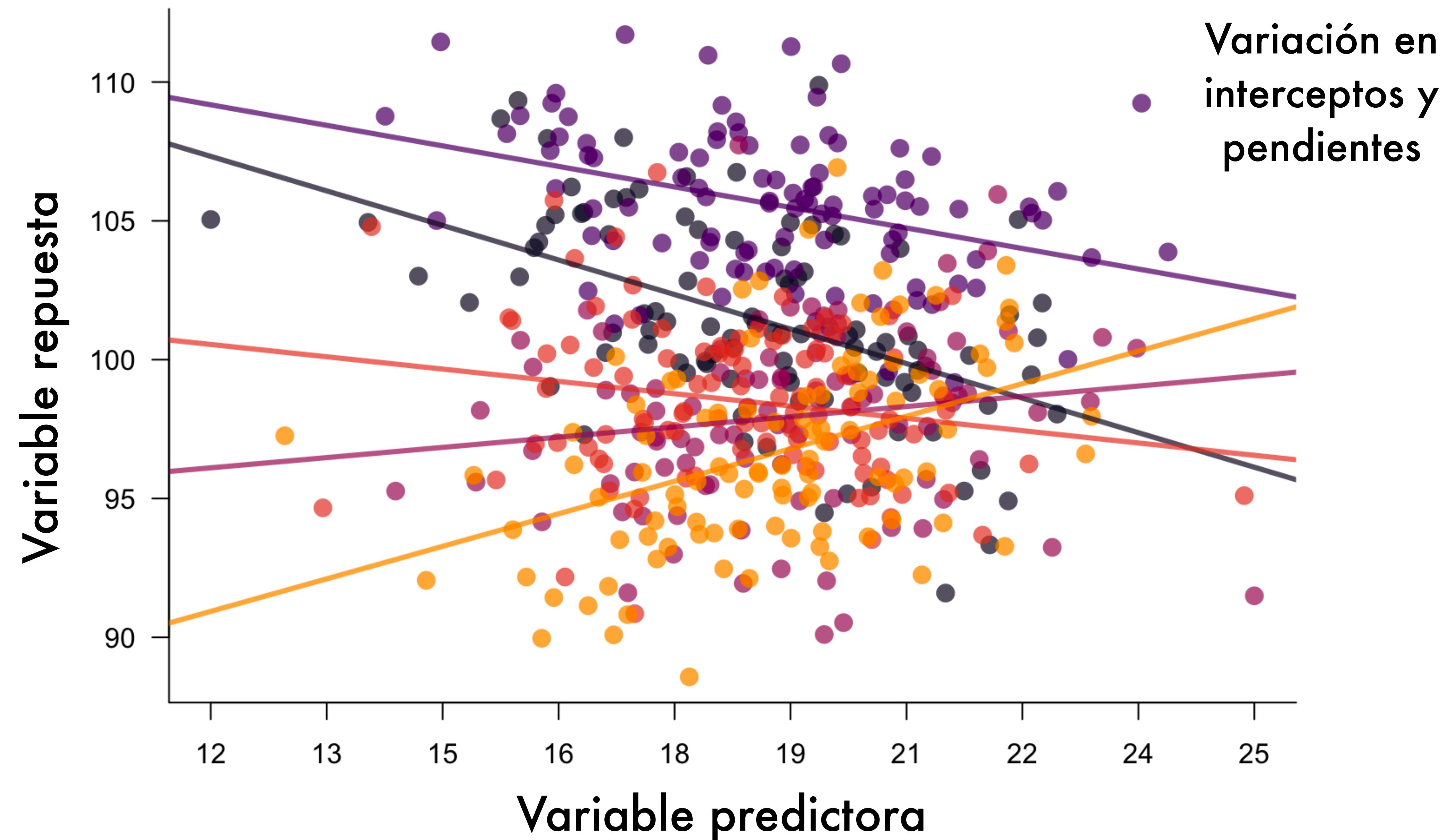
Modelos mixtos



Modelos mixtos

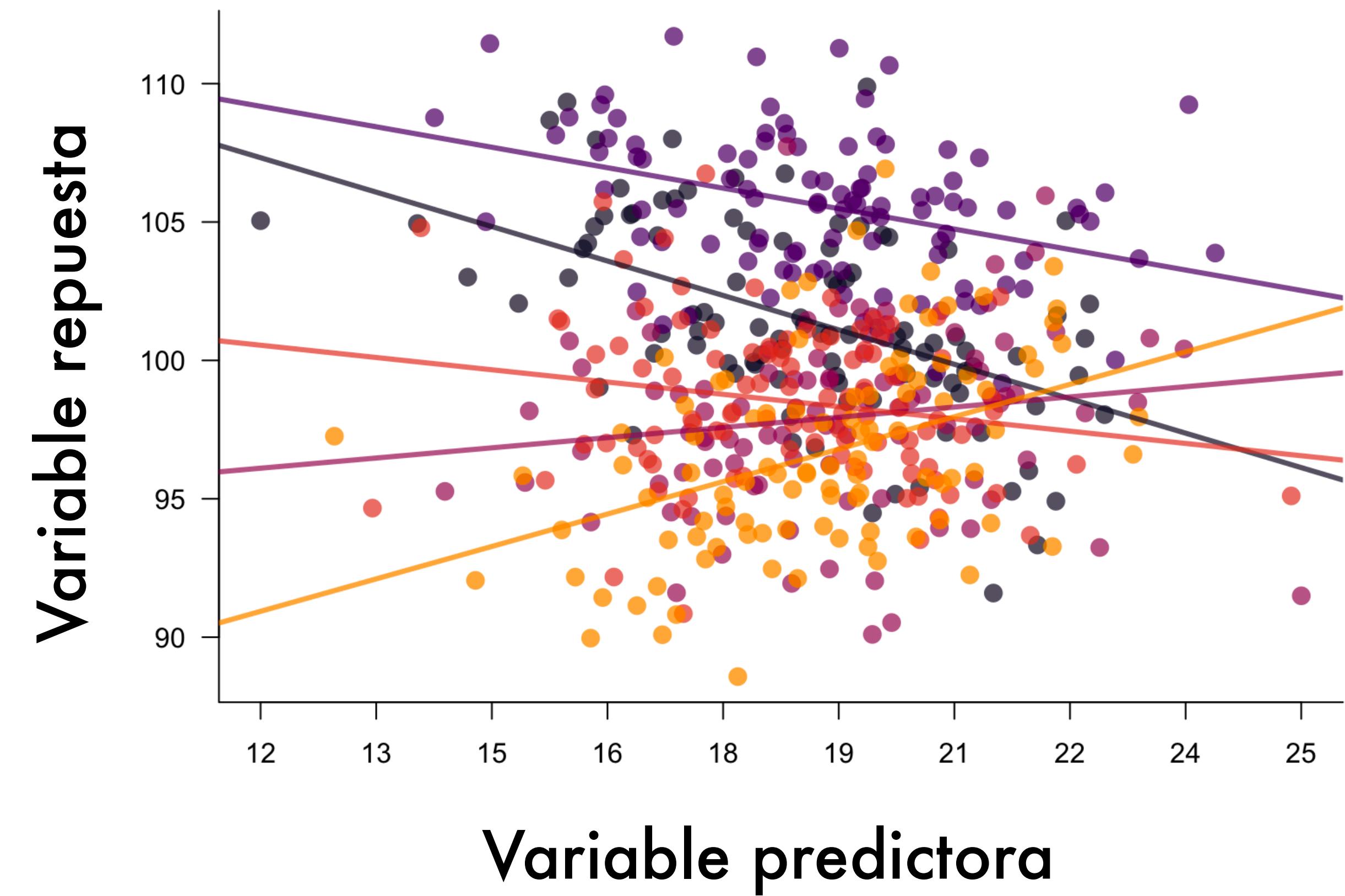


Modelos mixtos



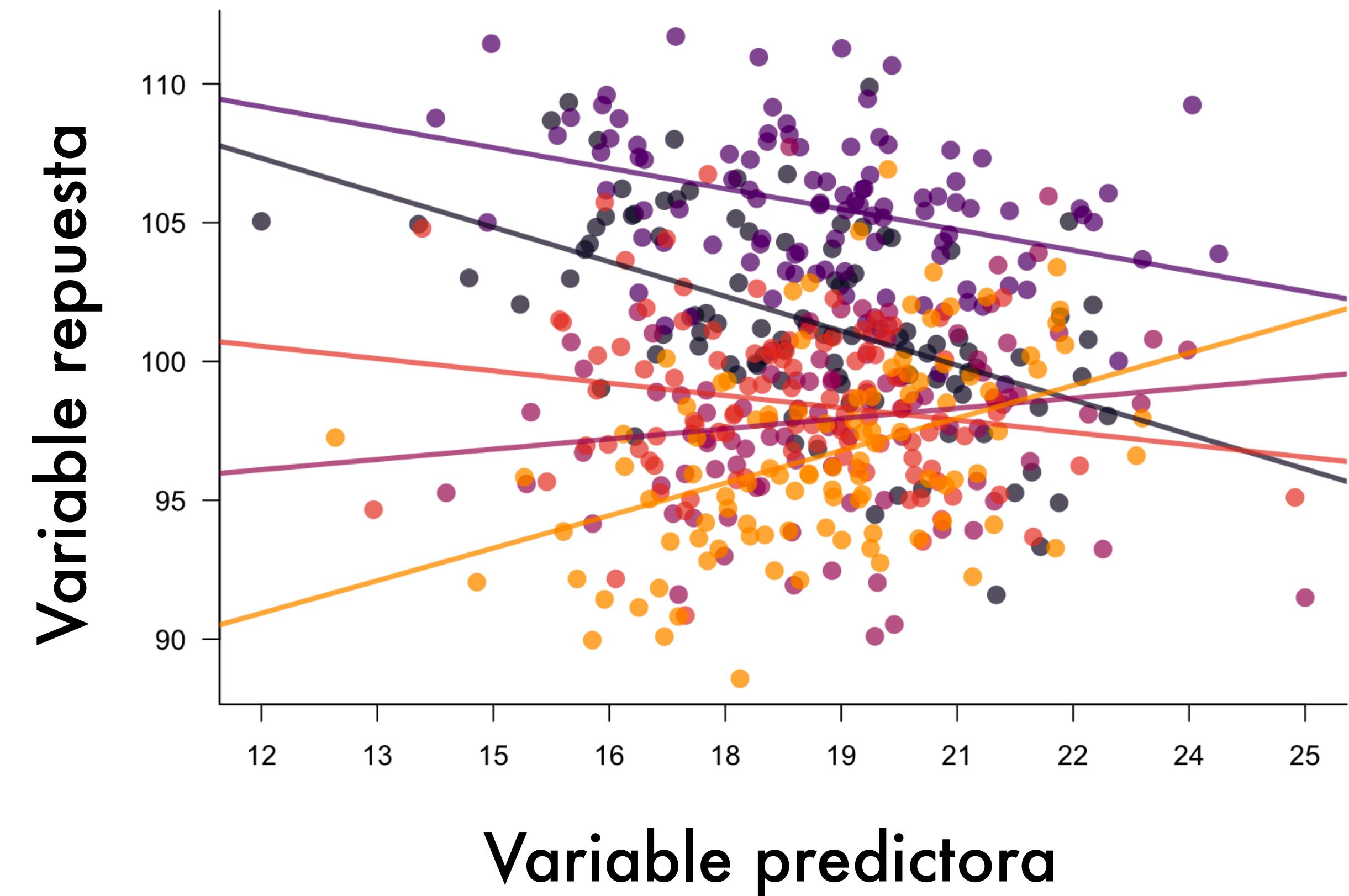
Modelos mixtos

- Como vimos con variables categóricas, grupos pueden variar en su relación con la variable respuesta.



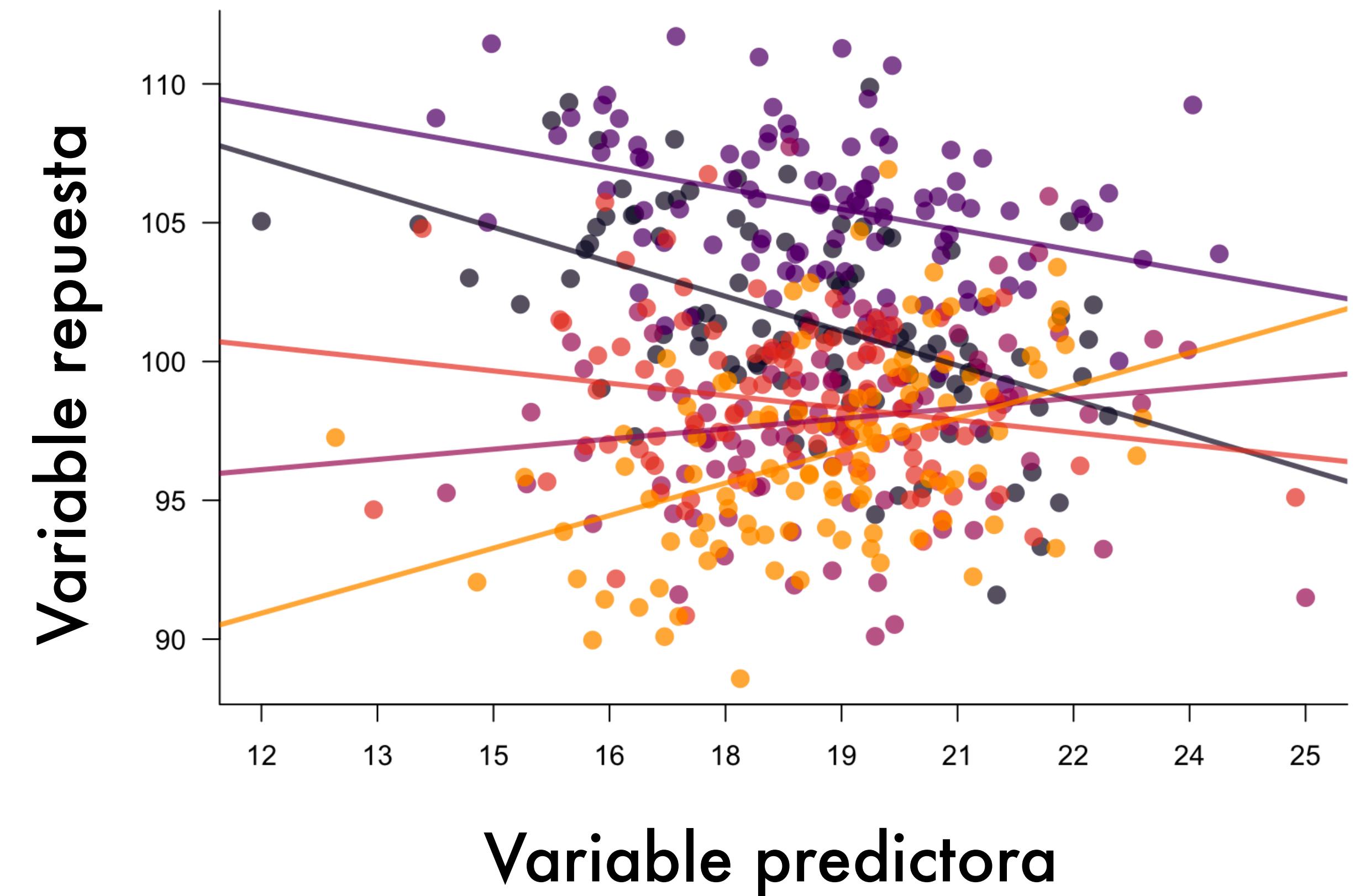
Modelos mixtos

- Como vimos con variables categóricas, grupos pueden variar en su relación con la variable respuesta.
- Con variables categóricas, los grupos varian en su relación con la variable respuesta en una manera que nos interesa y que sean independientes.



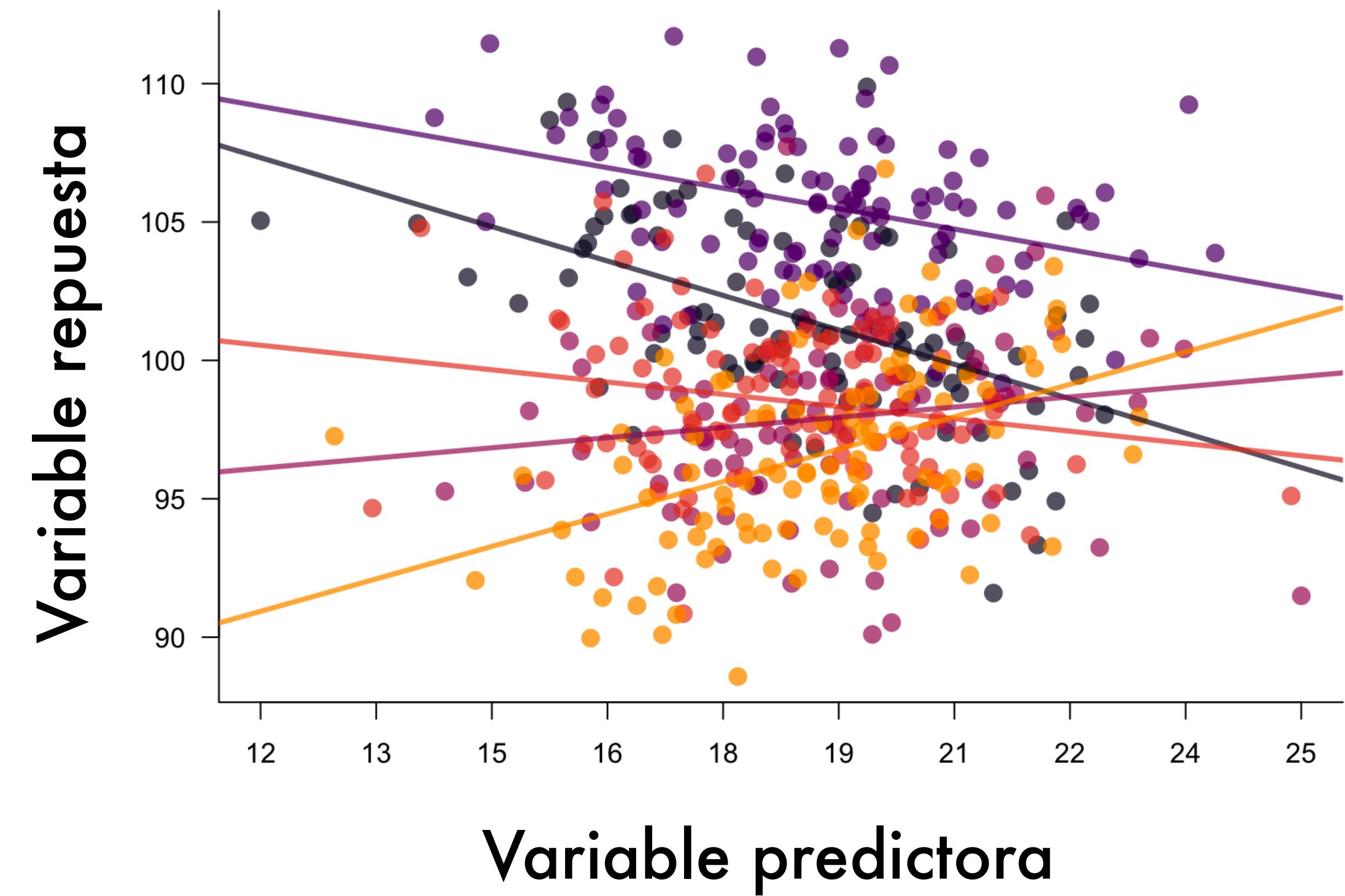
Modelos mixtos

- Como vimos con variables categóricas, grupos pueden variar en su relación con la variable respuesta.
- Con variables categóricas, los grupos varian en su relación con la variable respuesta en una manera que nos interesa y que sean independientes.
- Pero a veces la variabilidad entre grupos no es informativa, y no son independiente.



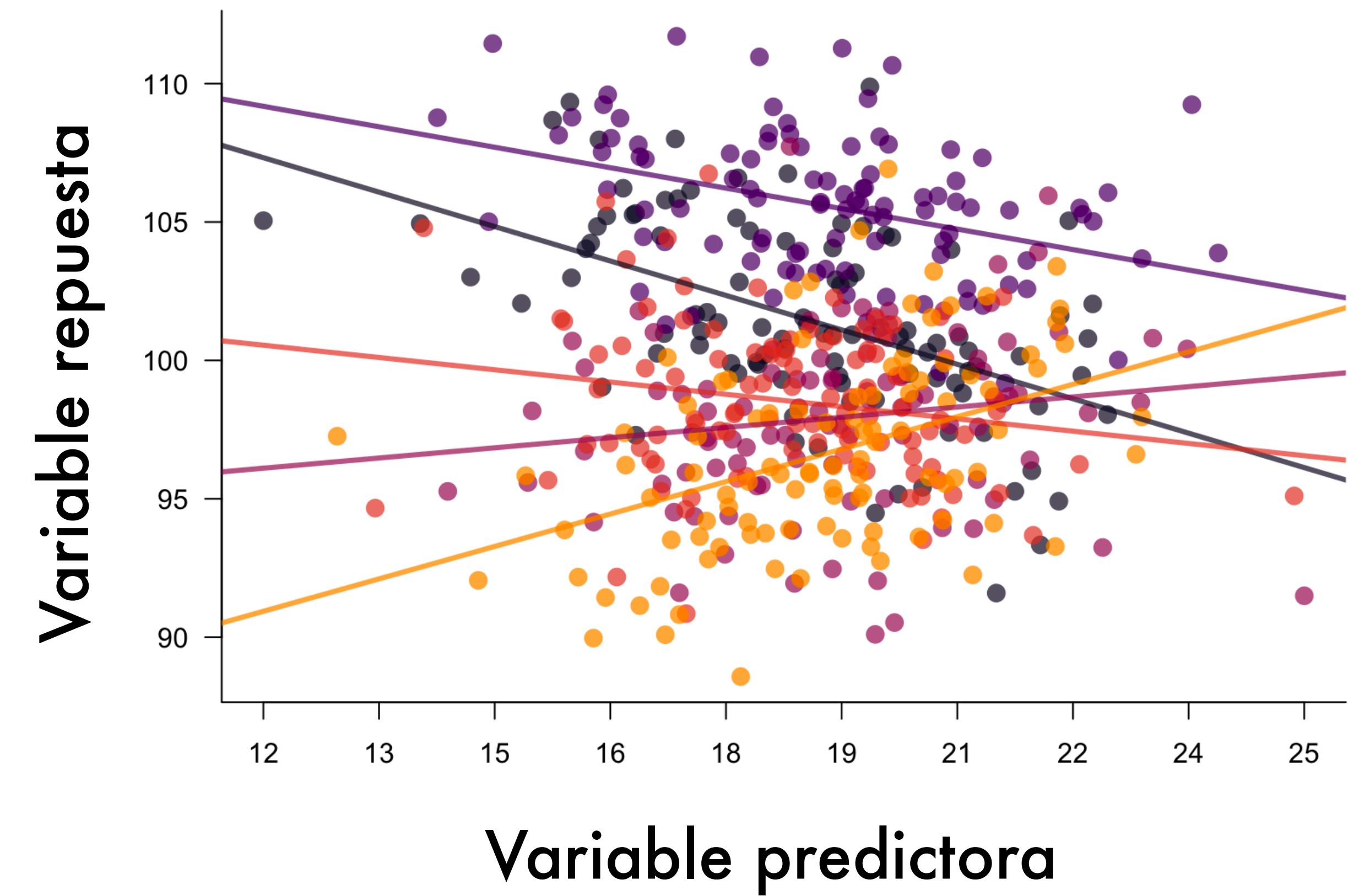
Modelos mixtos

- Este tipo de variabilidad presente un problema porque puede introducir ruido al modelo o romper la suposición de independencia.



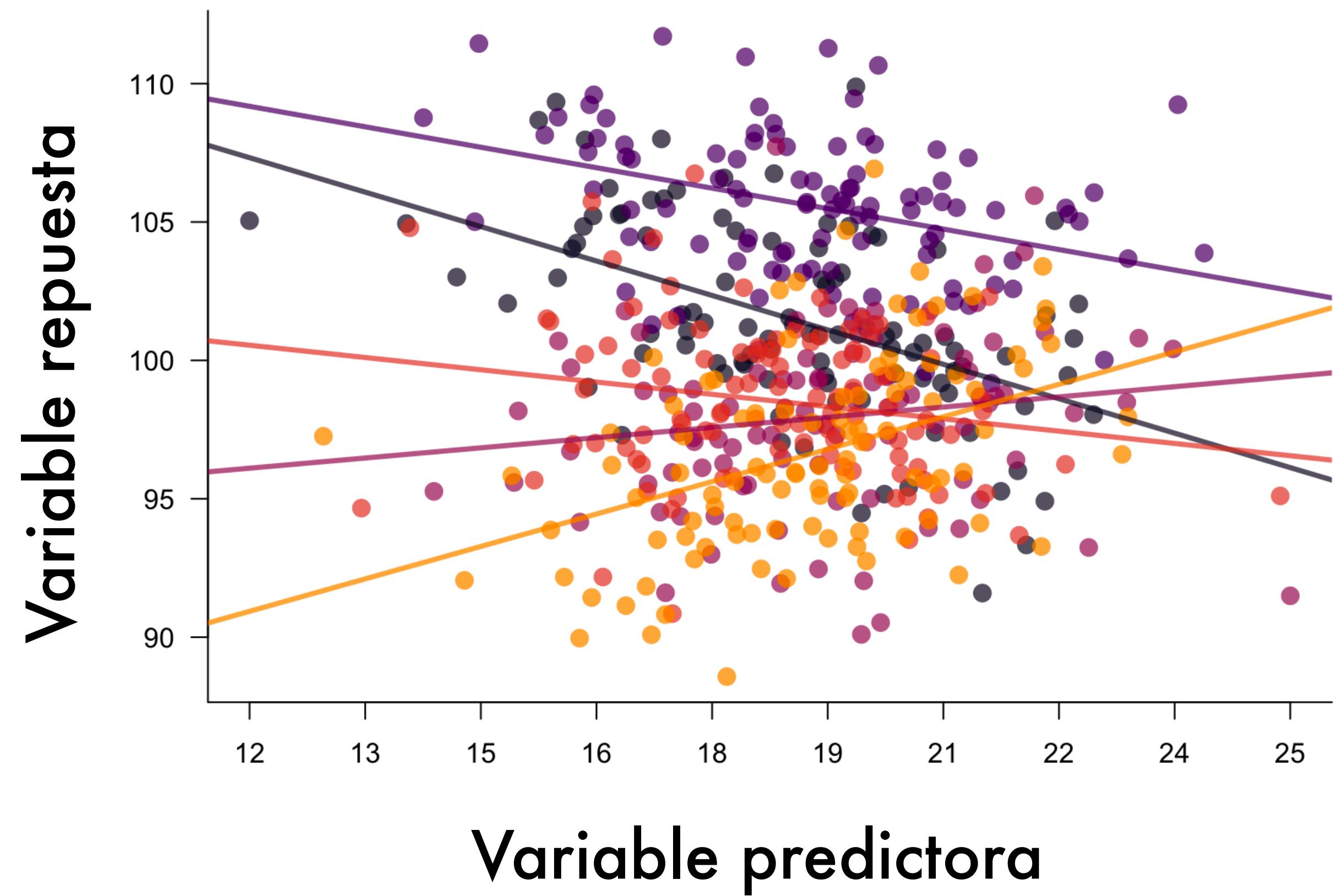
Modelos mixtos

- Este tipo de variabilidad presente un problema porque puede introducir ruido al modelo o romper la suposición de independencia.
- En contraste al efectos fijos, este tipo de variabilidad se llama efectos aleatorios.



Modelos mixtos

- Este tipo de variabilidad presente un problema porque puede introducir ruido al modelo o romper la suposición de independencia.
- En contraste al efectos fijos, este tipo de variabilidad se llama efectos aleatorios.
- Por ejemplo, muestreros replicados o bloques espacial.



Modelos mixtos

En R usando paquete *lme4*

Modelos mixtos

En R usando paquete *lme4*

```
# Cargar paquete lme4
library(lme4)

# Ajustar modelo con intercepto unico y un predictor
mod_lm <- lm(response ~ predictor)
```

Modelos mixtos

En R usando paquete *lme4*

```
# Cargar paquete lme4
library(lme4)

# Ajustar modelo con intercepto unico y un predictor
mod_lm <- lm(response ~ predictor)

# Ajustar modelo con interceptos aleatorios
mod_int <- lmer(response ~ predictor + (1 | block))
```

Modelos mixtos

En R usando paquete *lme4*

```
# Cargar paquete lme4
library(lme4)

# Ajustar modelo con intercepto unico y un predictor
mod_lm <- lm(response ~ predictor)

# Ajustar modelo con interceptos aleatorios
mod_int <- lmer(response ~ predictor + (1 | block))

# Ajustar modelo con interceptos y pendientes aleatorios
mod_slope <- lmer(response ~ predictor + (1 + predictor | block))
```

Modelos mixtos

Modelos mixtos

- Normalmente no interpretamos los efectos aleatorios. Si quieres interpretarlo, puede que sea mejor como efecto fijo.
- Recomendamos decidir sus efectos aleatorios a priori por dependen el diseño experimental y no de la hipótesis.
- Aunque normalmente no interpretamos los valores del efectos aleatorios, puede ser útil examinar componentes de variabilidad.

Modelo conceptual

Formar la pregunta / escribir el modelo

Diseño experimental

Colección de datos

Armar el modelo

Creer resultados

Modelo conceptual

Formar la pregunta

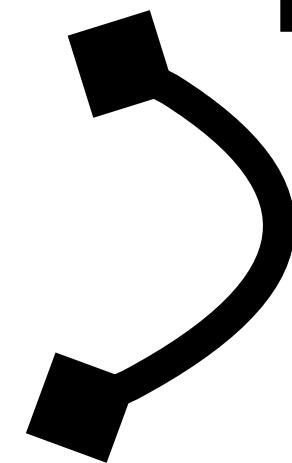
Diseño experimental

Colección de datos

Armar el modelo

Creer resultados

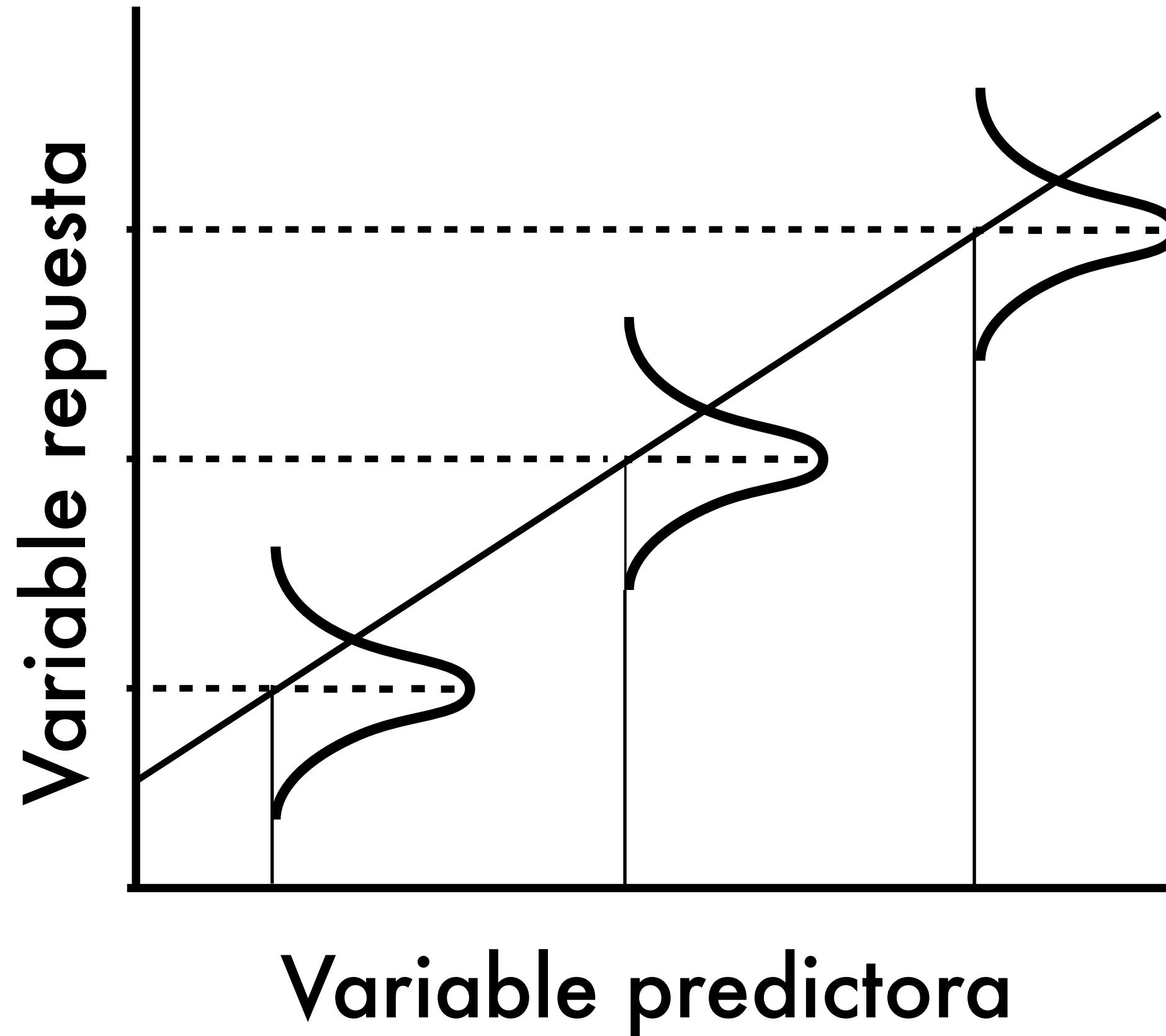
Escribir el modelo



IV. Modelos lineales generalizados

Recuerden suposición 2 de modelos lineales:

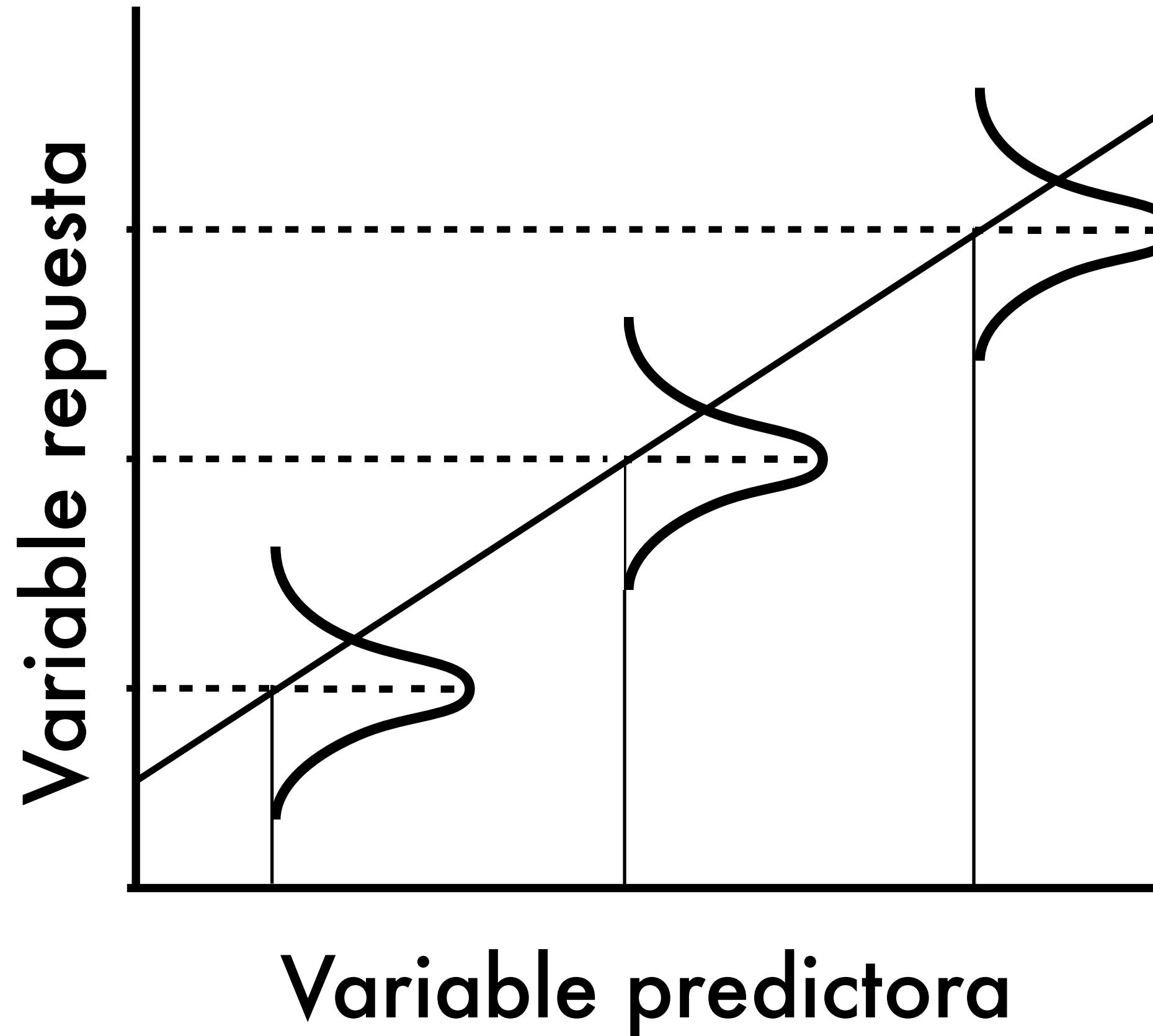
Residuales son con distribución normal



$$y_i = \alpha + \beta * x_i + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2)$$

El matemática también puede escribir así:



$$y_i = \alpha + \beta * x_i + \epsilon_i$$

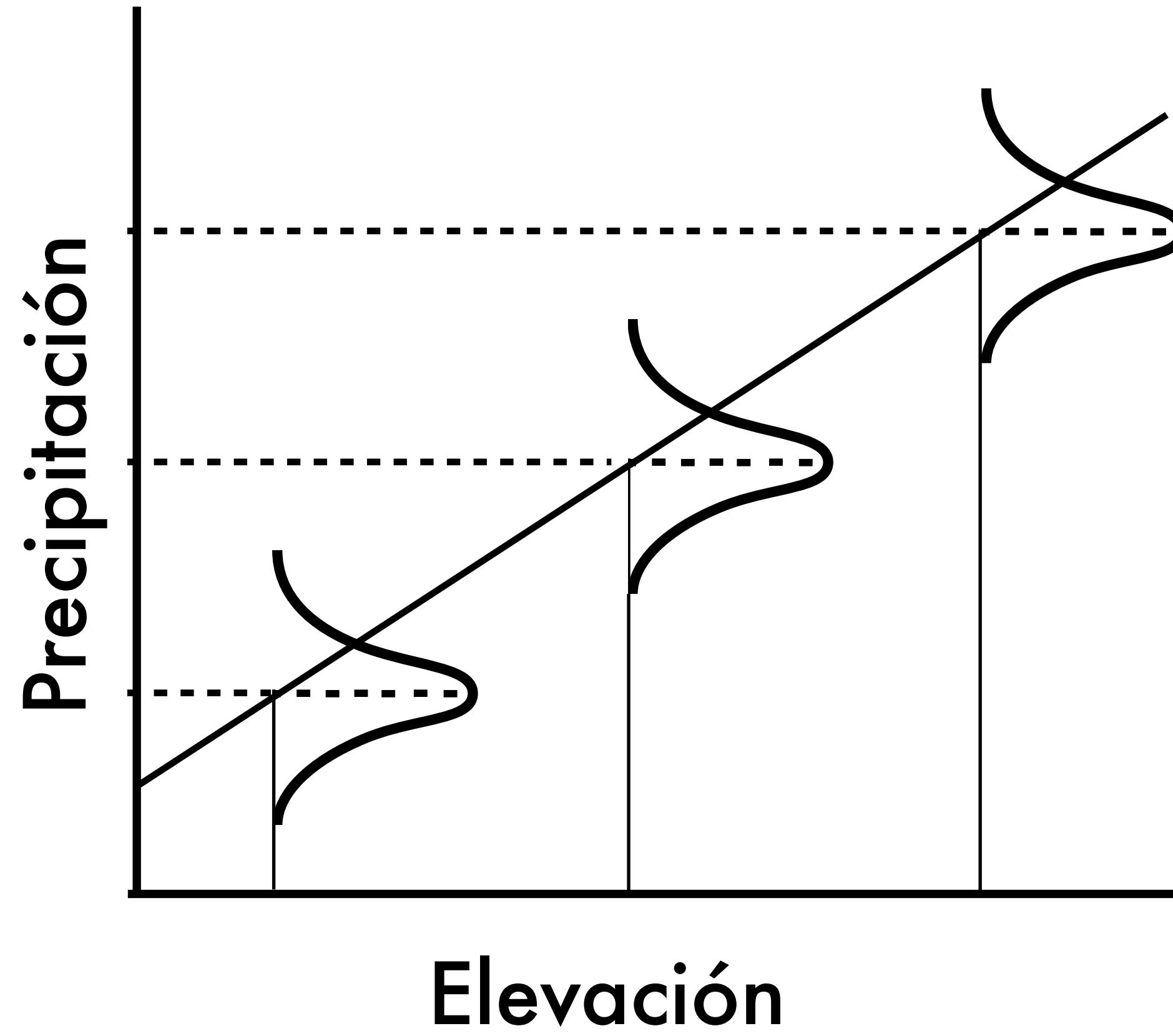
$$\epsilon_i \sim N(0, \sigma^2)$$

==

$$y_i \sim N(\mu_i, \sigma^2)$$

$$\mu_i = \alpha + \beta * x_i$$

Modelos mixtos



$$Y = \alpha + \beta * X + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

$$Y \sim N(\mu_t, \sigma^2)$$

$$\mu_t = \alpha + \beta * X$$

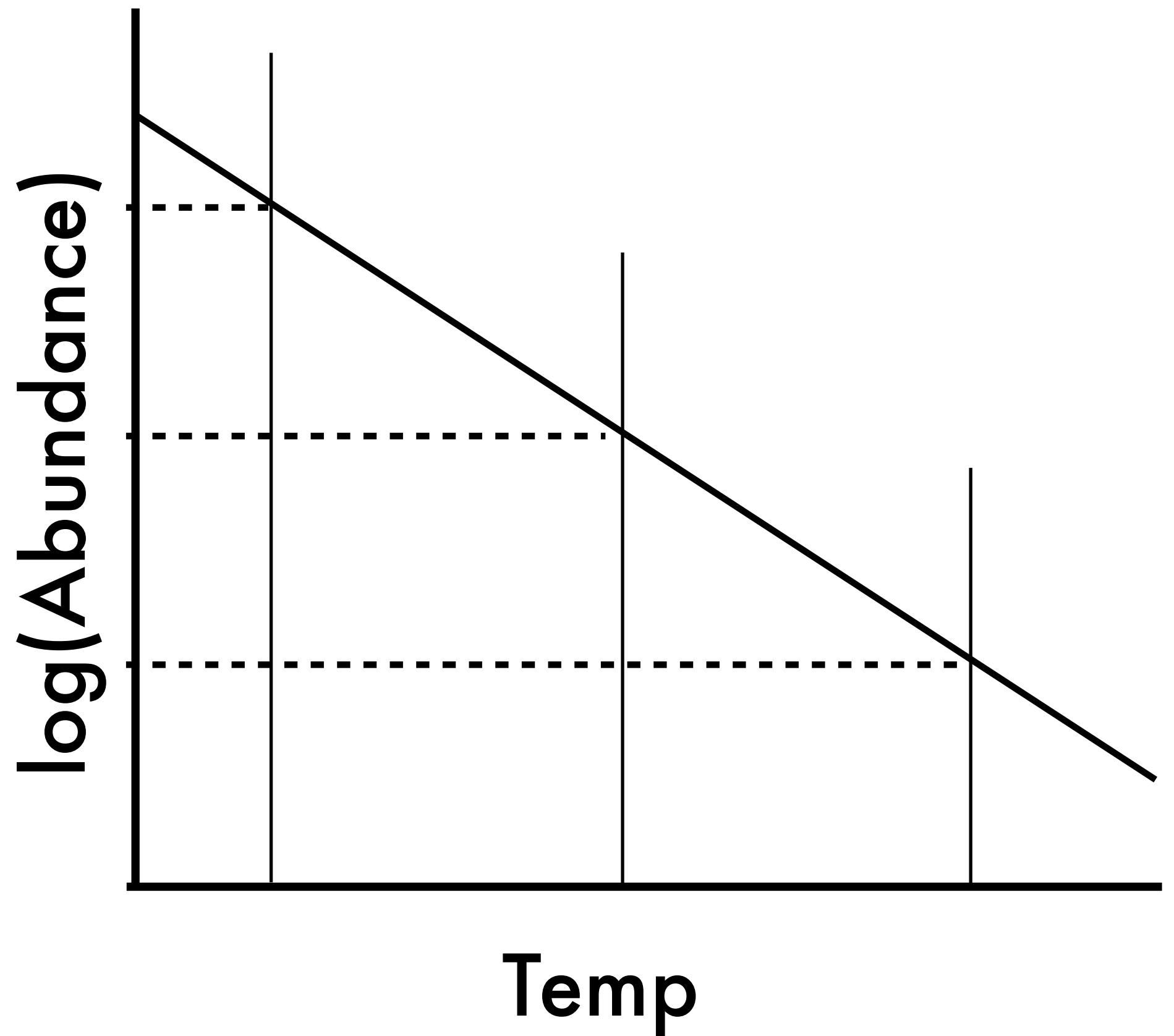
Writing data generative models

How does temperature of hollow
affect ringtail possum abundance?



$$A \sim \text{Poisson}(\lambda)$$

$$\log(\lambda) = \alpha + \beta * \text{temp}$$



Exercise 2: Write a data generative linear model

Simulating data

Model



Model-defined
data

Simulating data

Model



Model-defined
data

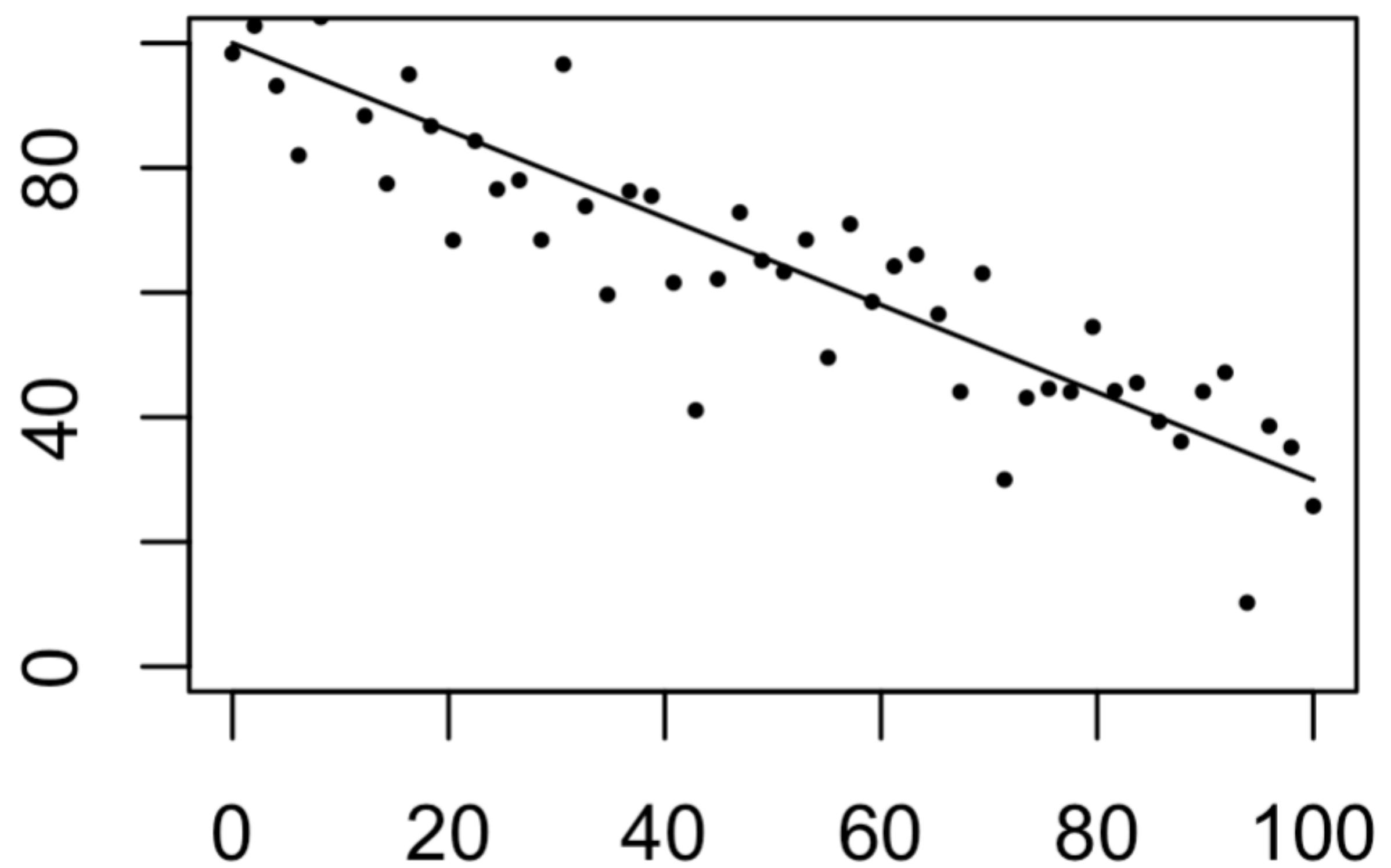
$$M_t \sim N(\mu_t, \sigma^2)$$

$$\mu_t = \alpha + \beta * time$$

R practical: simulating data

$$y = a + bx + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$



What's it all about?

$$y = a + bx + \epsilon$$

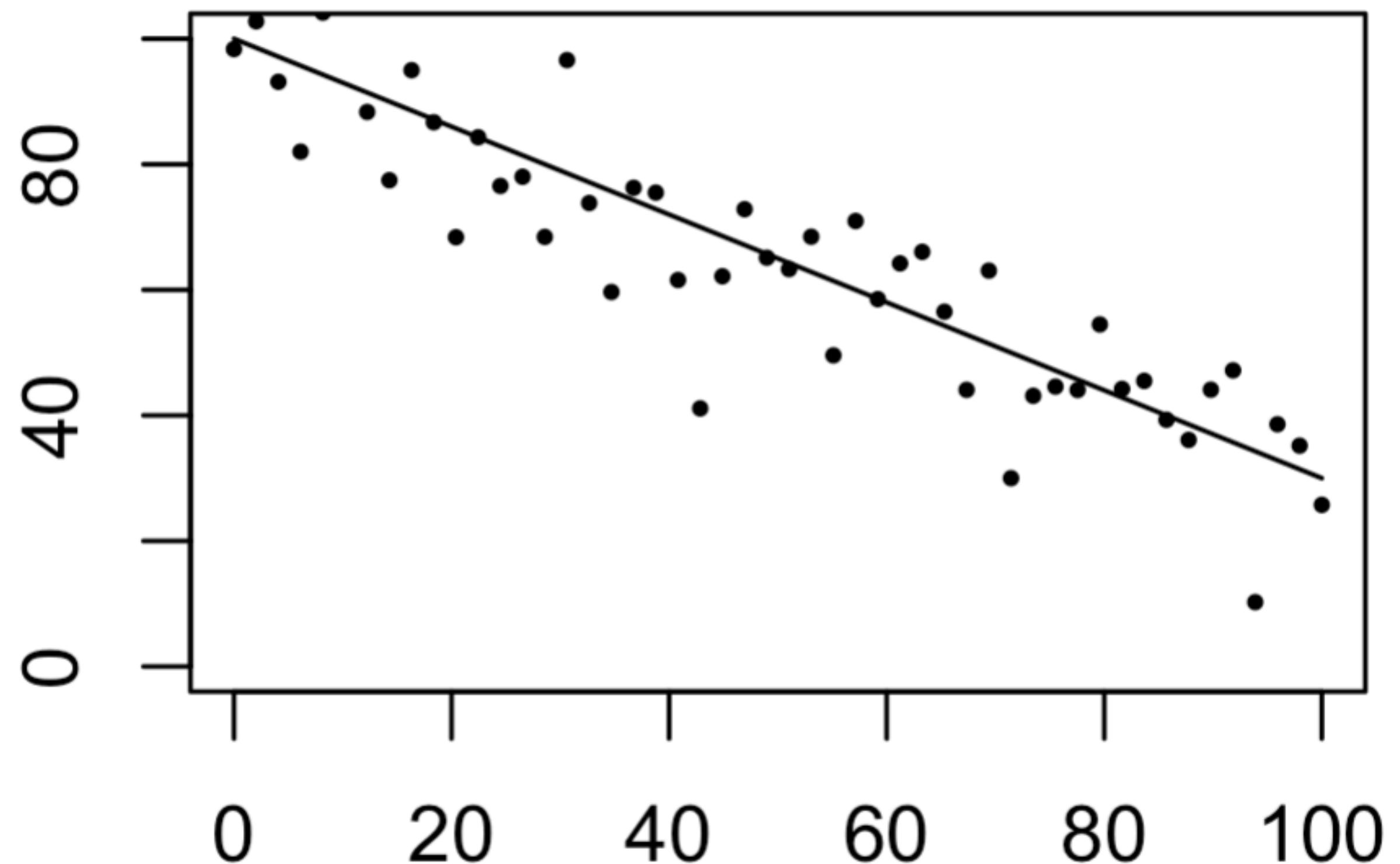
$$\epsilon \sim N(0, \sigma^2)$$

MLE

$$b = -0.67$$

Posterior mean

$$b = -0.67$$



What's it all about?

Why Bayes

Flexibility in building models > can test any hypothesis, not just the null

Principled way of building your model > think about what you know & the implications of what you think you know

Your inferences will fail often, giving you a way to diagnose the issue

Scales up, same tools and workflow for range of model processes > from t test to bayesian network analysis

More subjective > you make more choices.

What's it all about?

Statistical models exist independent of the method of estimating parameters

What's it all about?

Statistical models exist independent of the method of estimating parameters

There are no **Frequentist models** or
Bayesian models

What's it all about?

Statistical models exist independent of the method of estimating parameters

There are no Frequentist models or
Bayesian models

May choose to analyse a model in a Bayesian way

I. Thinking about models

- WHAT DISTRIBUTION?
 - LINK FUNCTION NEEDED?
- HIERARCHICAL STRUCTURE?
- INTERACTIONS?

WRITE MODEL

WRITE QUESTION

DETERMINE AIM

- ESTIMATE VALUE OF PARAMETER
- HYPOTHESIS TESTING
- PREDICTION

FIT MODEL

MODEL CHECKING

MODEL INTERPRETATION

PREDICTION

HYPOTHESIS TESTING

ESTIMATE VALUE OF PARAMETER

► RESIDUALS

► RESPONSE PLOTS

► P-VALUE

► SUMMARIES

Exercise 1: Research workflow

Conceptual model

Prepare question

Experimental design

Data collection

Write and build model

Create model outputs

Thinking about models

Conceptual model

Prepare question

Experimental design

Data collection

Write and build model

Create model outputs

Thinking about models

== Using a model to understand our data.

Conceptual model

Prepare question

Experimental design

Data collection

Write and build model

Create model outputs

Thinking about models

Conceptual model

Prepare question / Write model

Experimental design

Data collection

Build model

Create model outputs

Thinking about models

Conceptual model

Prepare question / Write model

Experimental design

We want to
collect data to
test our model.

Data collection

Build model

Create model outputs

Thinking about models

Using a model to understand our data is asking:

“What is the probability of this data, given a certain model?”

Thinking about models

Using a model to understand our data is asking:

“What is the probability of this data, given a certain model?”

Collecting data to test our model is asking:

“What is the probability of this model, given observed data?”

Thinking about models

A well-specified research question defines a statistical model.

AND

A statistical model underpins a research question.

Thinking about models

Model



**Model-defined
data**

Thinking about models

Model



Model-defined
data

?=

Experiment



Collected data

VI. Model checking

Addressing your question

Estimating value of parameter:

Forest plot, summarise, mean and sd of posterior of parameter.

Hypothesis testing:

Bayesian p-value. "what is the posterior probability that X has an effect on Y."

Prediction:

calculate, effect/response plot to new data. Cross-validation - testing on new data.

Review

Special thanks to

Nick Golding

Marc Kery

Gerry Ryan

Richard McElreath

III. Thinking about probability

Thinking about probability

Using a model to understand our data is asking:

“What is the probability of this data, given a certain model?”

Thinking about probability

Using a model to understand our data is asking:

“What is the probability of this data, given a certain model?”

$p(\text{data} \mid \text{parameter})$

Thinking about probability

Using a model to understand our data is asking:

“What is the probability of this data, given a certain model?”

$p(\text{data} \mid \text{parameter})$

Collecting data to test our model is asking:

“What is the probability of this model, given observed data?”

Thinking about probability

Using a model to understand our data is asking:

“What is the probability of this data, given a certain model?”

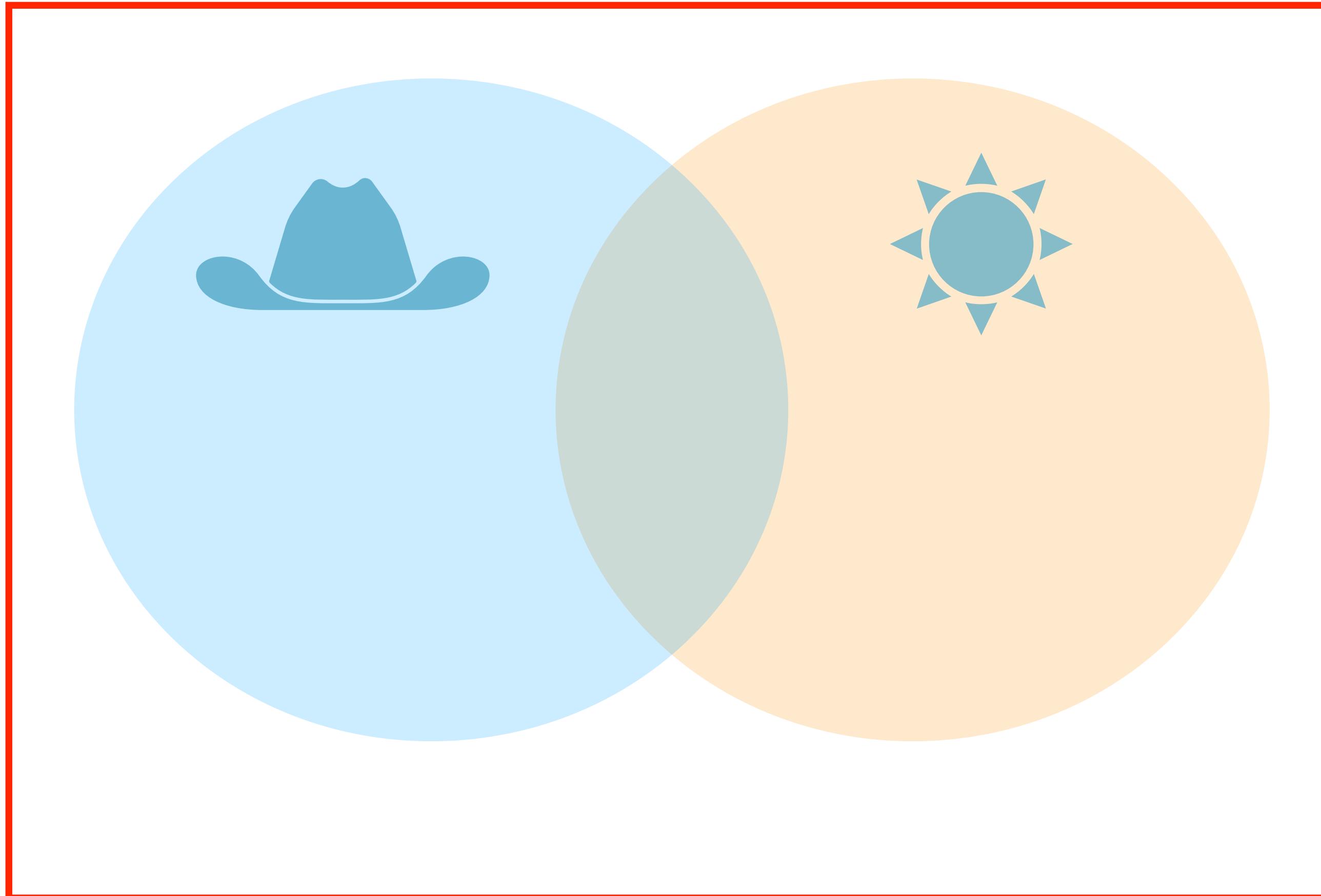
$p(\text{data} \mid \text{parameter})$

Collecting data to test our model is asking:

“What is the probability of this model, given observed data?”

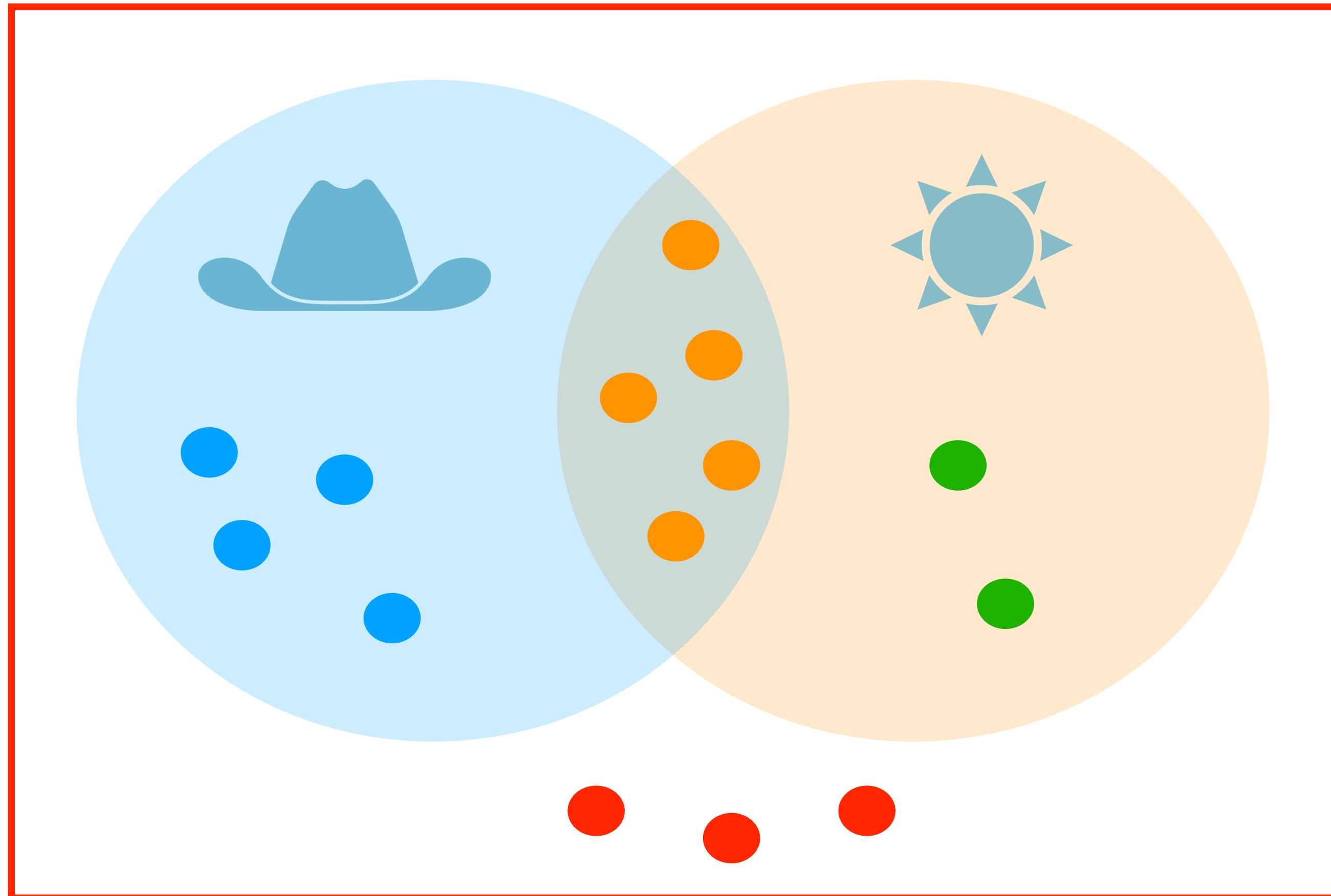
$p(\text{parameter} \mid \text{data})$

Thinking about probability

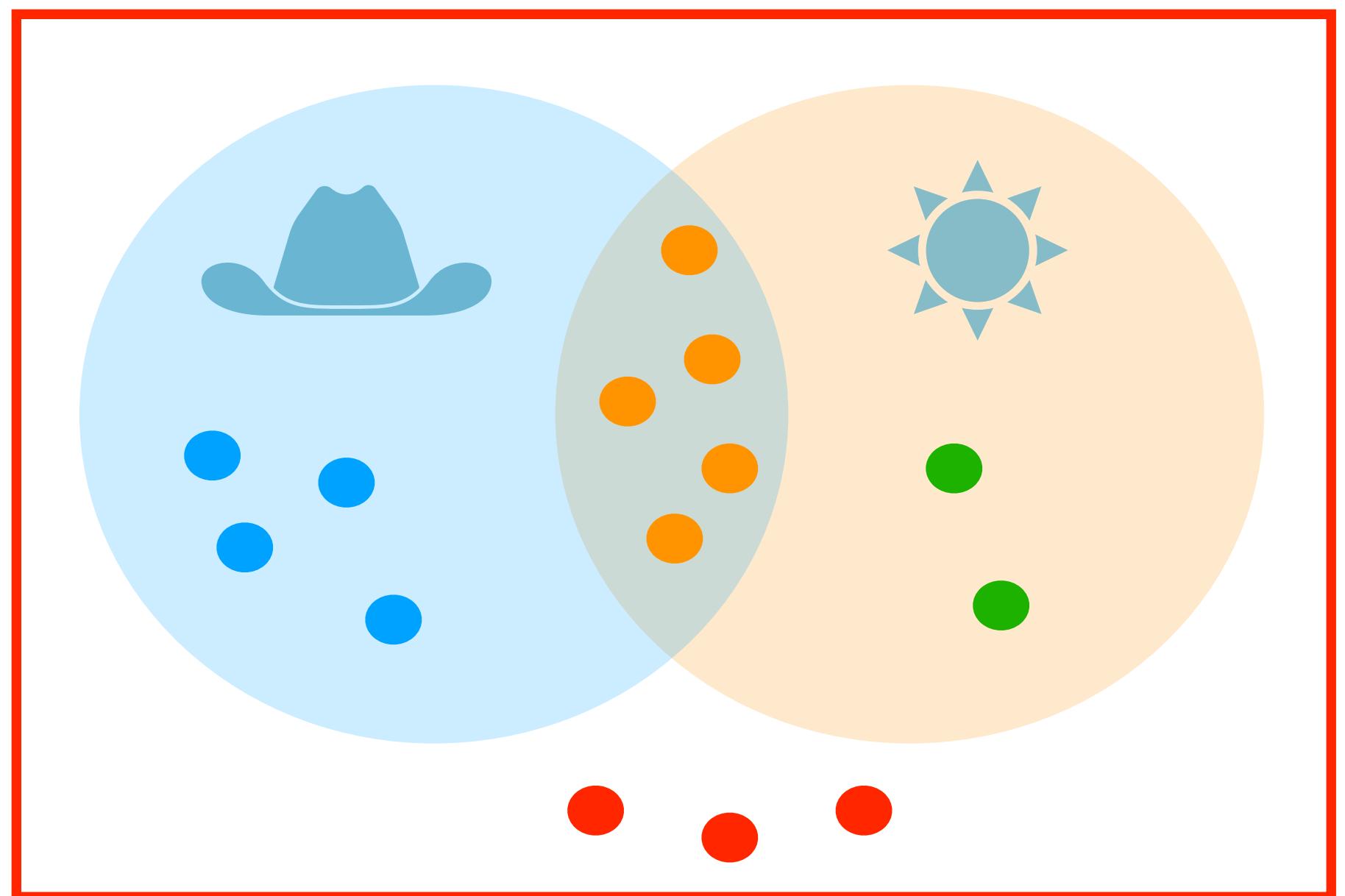


Check out: https://www.youtube.com/watch?v=9wCnvr7Xw4E&ab_channel=StatQuestwithJoshStarmer

Thinking about probability

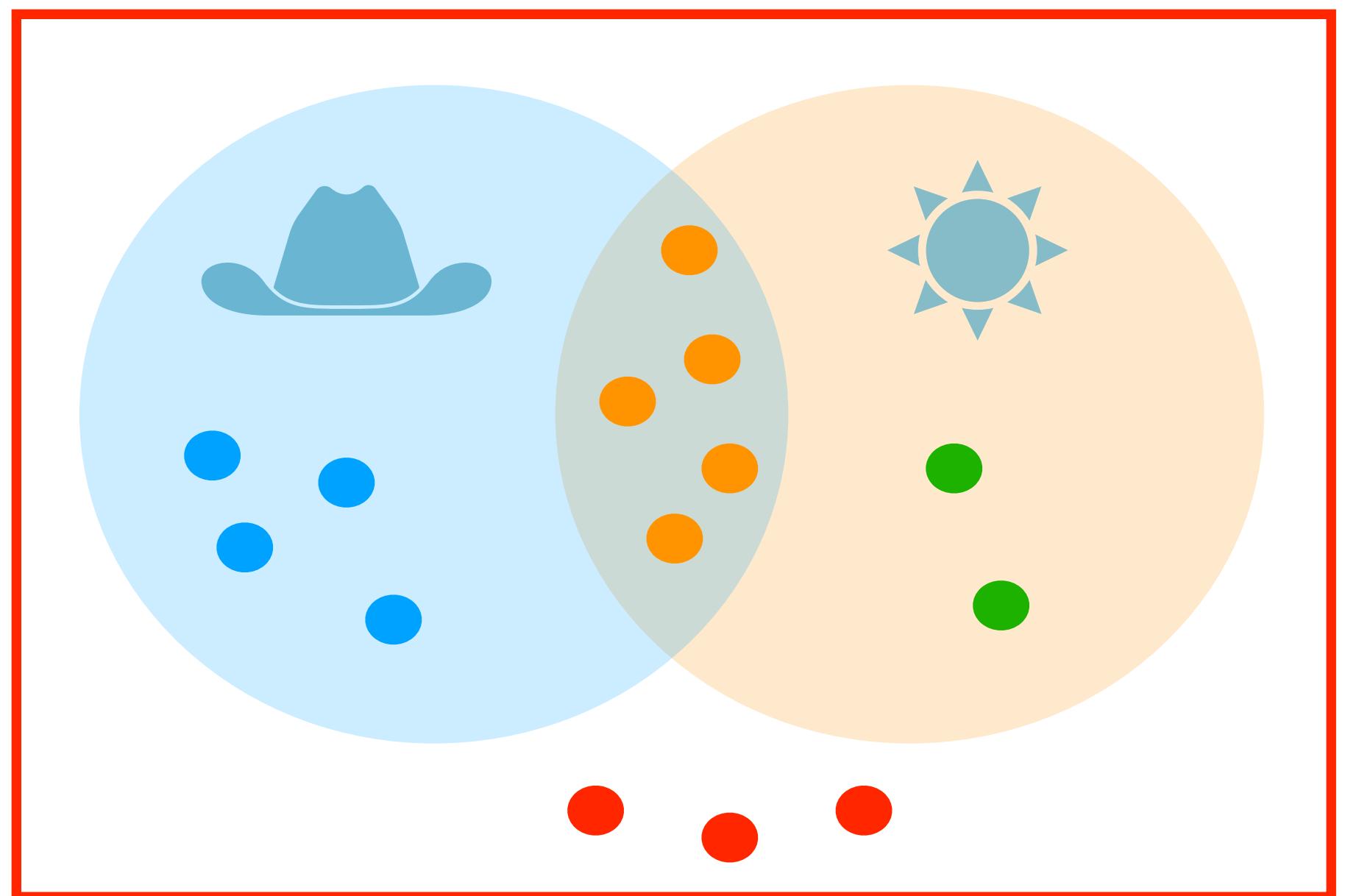


Thinking about probability



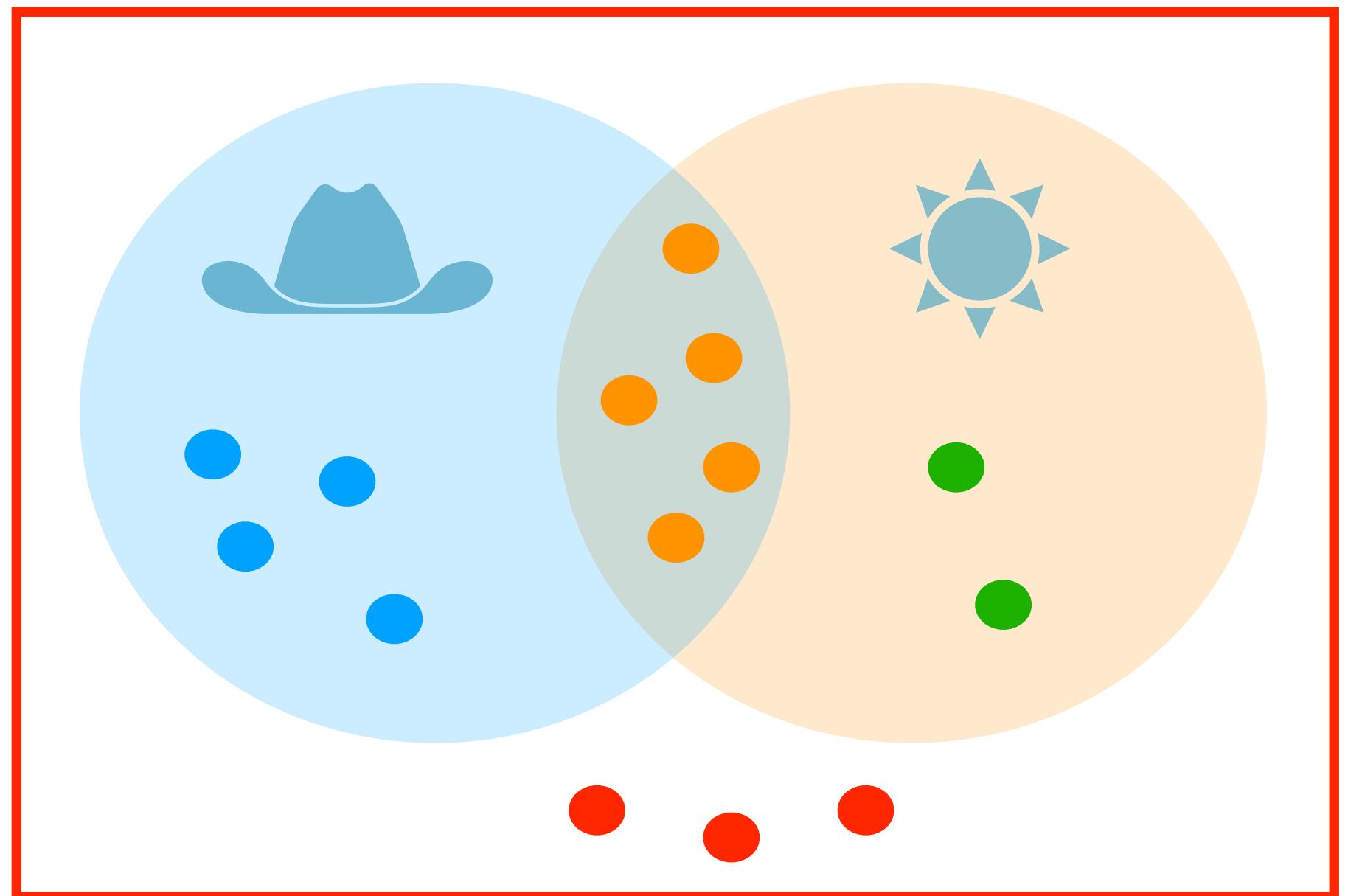
	No sun	Sun	Total
Wore hat	4/14		
Did not wear hat			
Total			

Thinking about probability



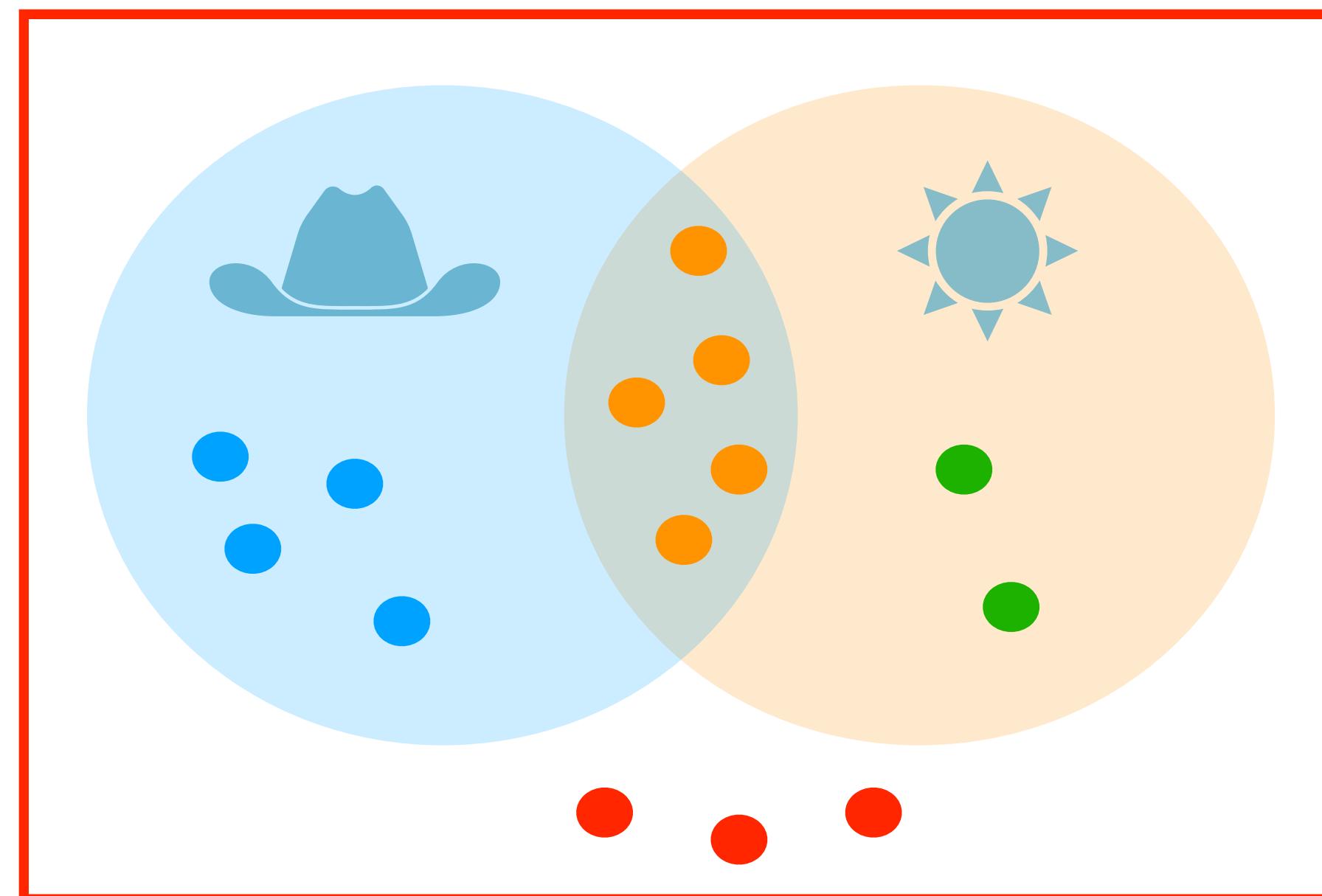
	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat			
Total			

Thinking about probability



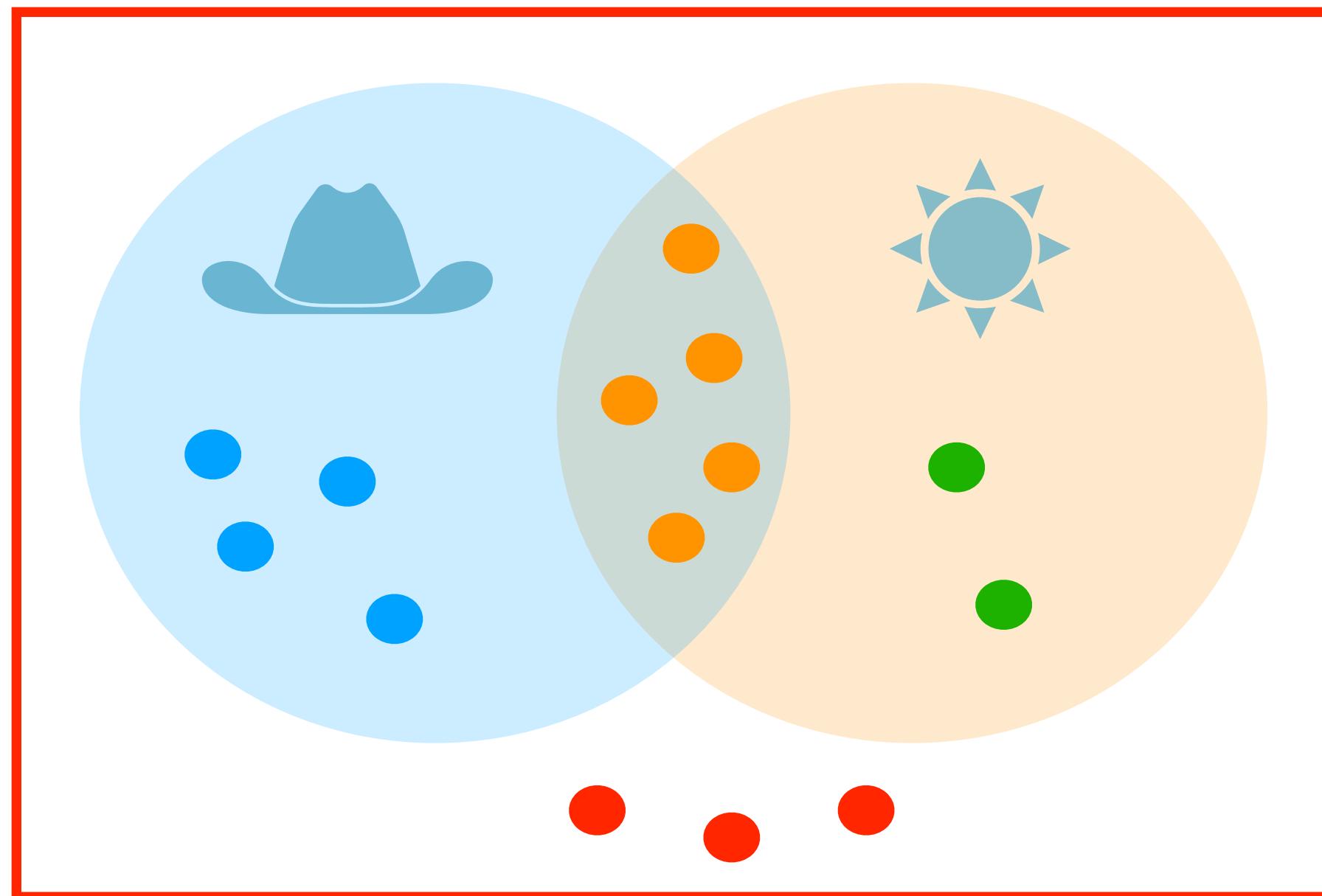
	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat	3/14		
Total			

Thinking about probability



	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat	3/14	2/14	5/14
Total	7/14	7/14	

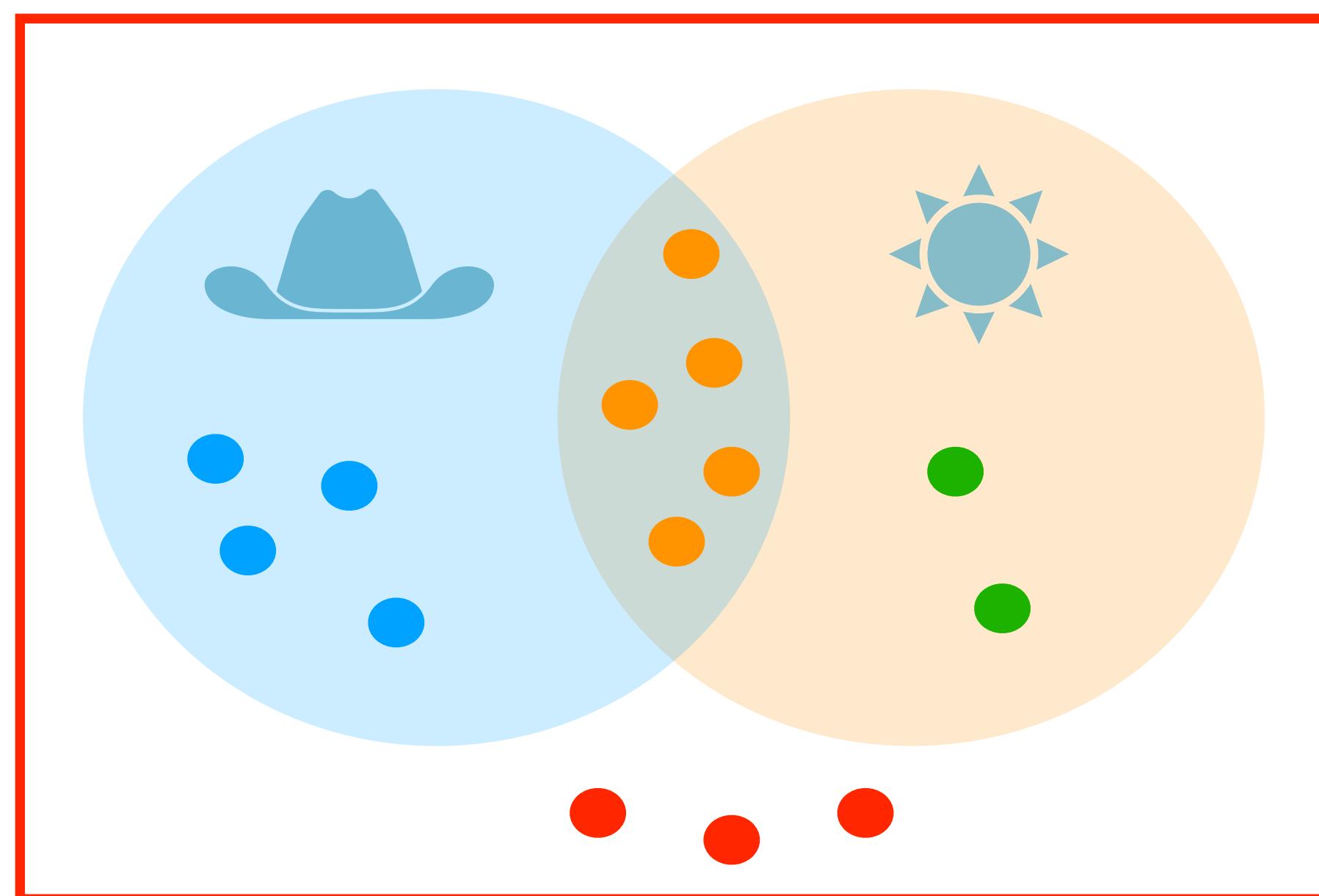
Thinking about probability



	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat	3/14	2/14	5/14
Total	7/14	7/14	

$$p(\text{sun} \& \text{hat} \mid \text{hat}) =$$

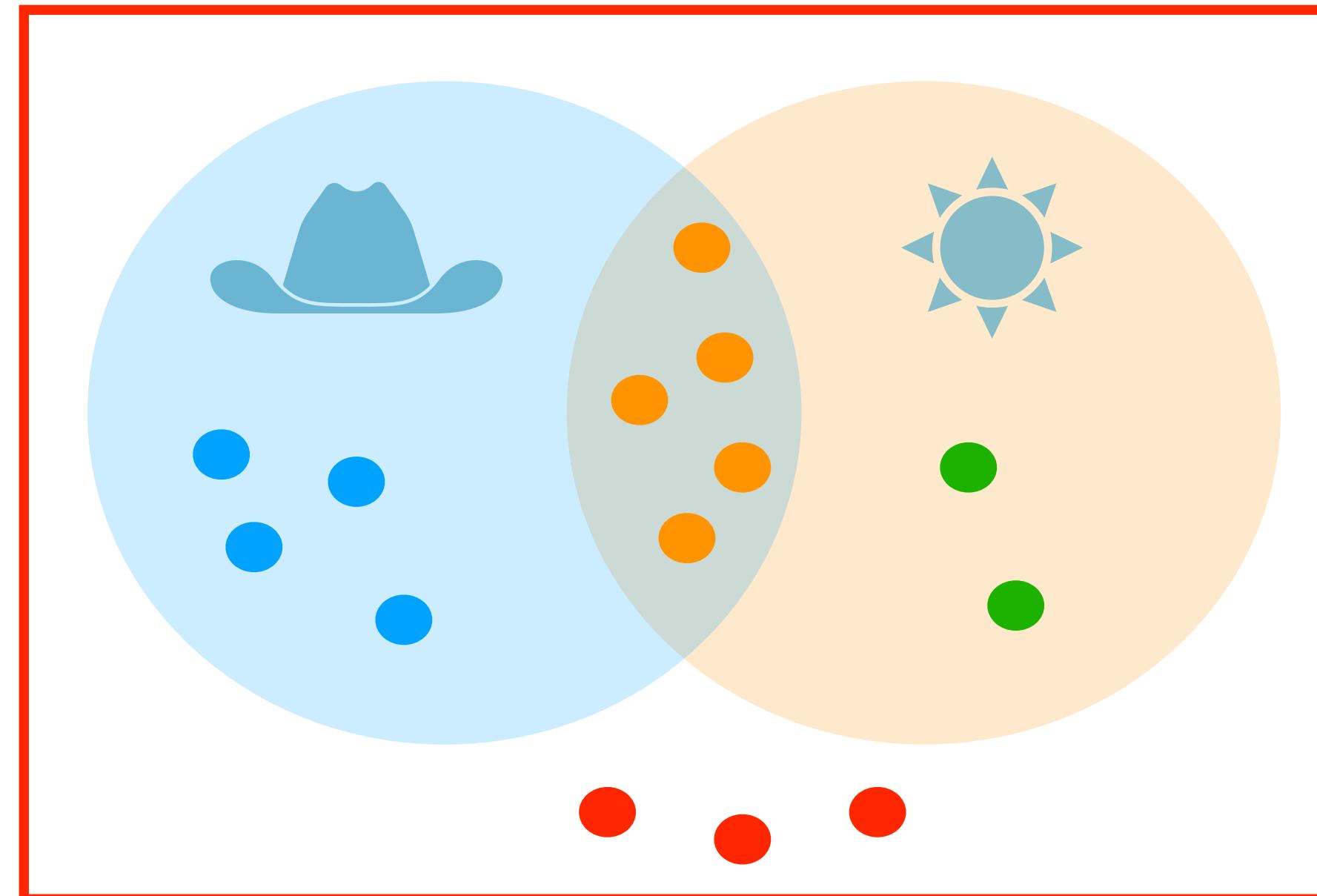
Thinking about probability



	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat	3/14	2/14	5/14
Total	7/14	7/14	

$$p(\text{sun} \& \text{hat} \mid \text{hat}) = \frac{\text{ }}{p(\text{hat})}$$

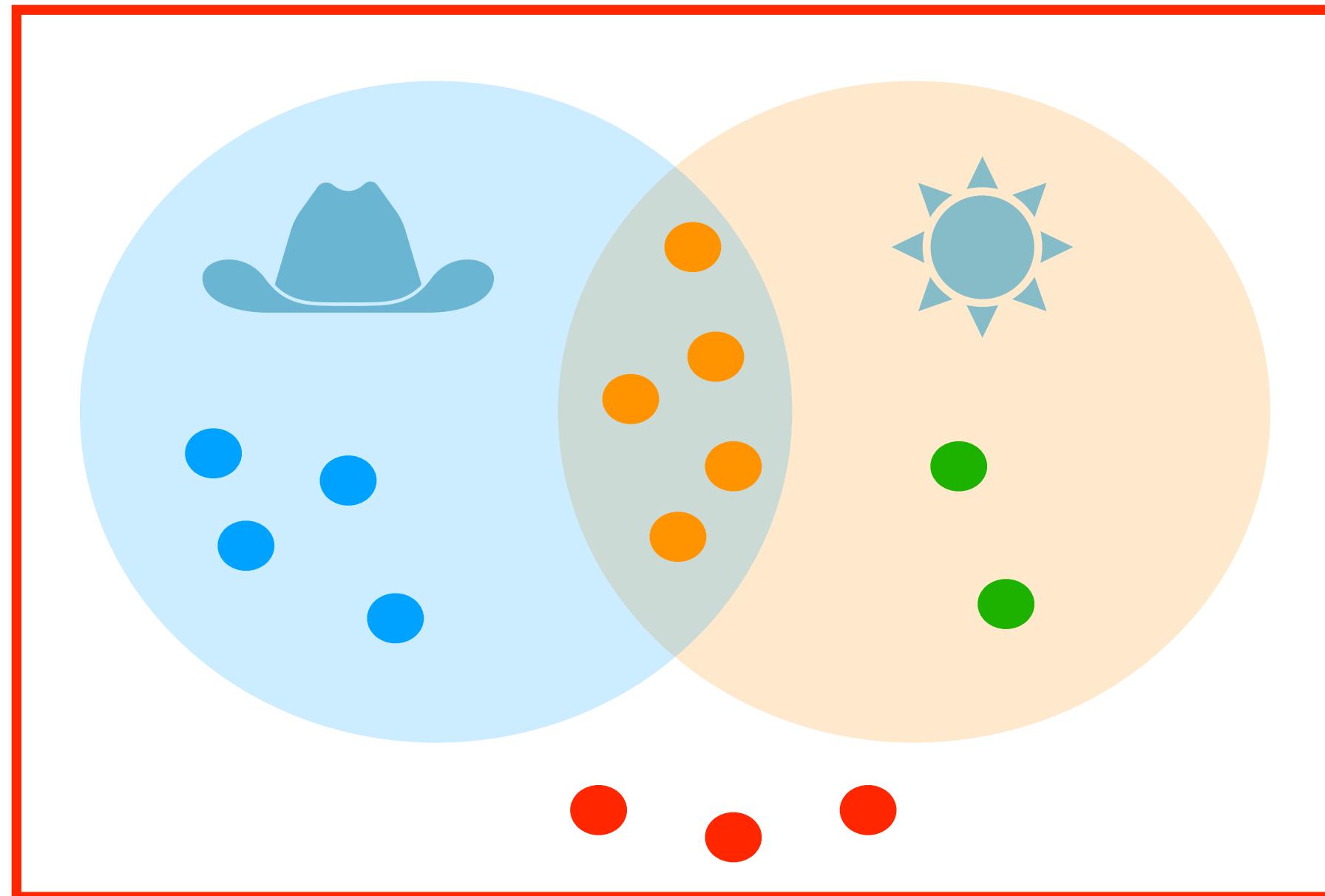
Thinking about probability



	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat	3/14	2/14	5/14
Total	7/14	7/14	

$$p(\text{sun} \& \text{hat} \mid \text{hat}) = \frac{p(\text{sun} \& \text{hat})}{p(\text{hat})}$$

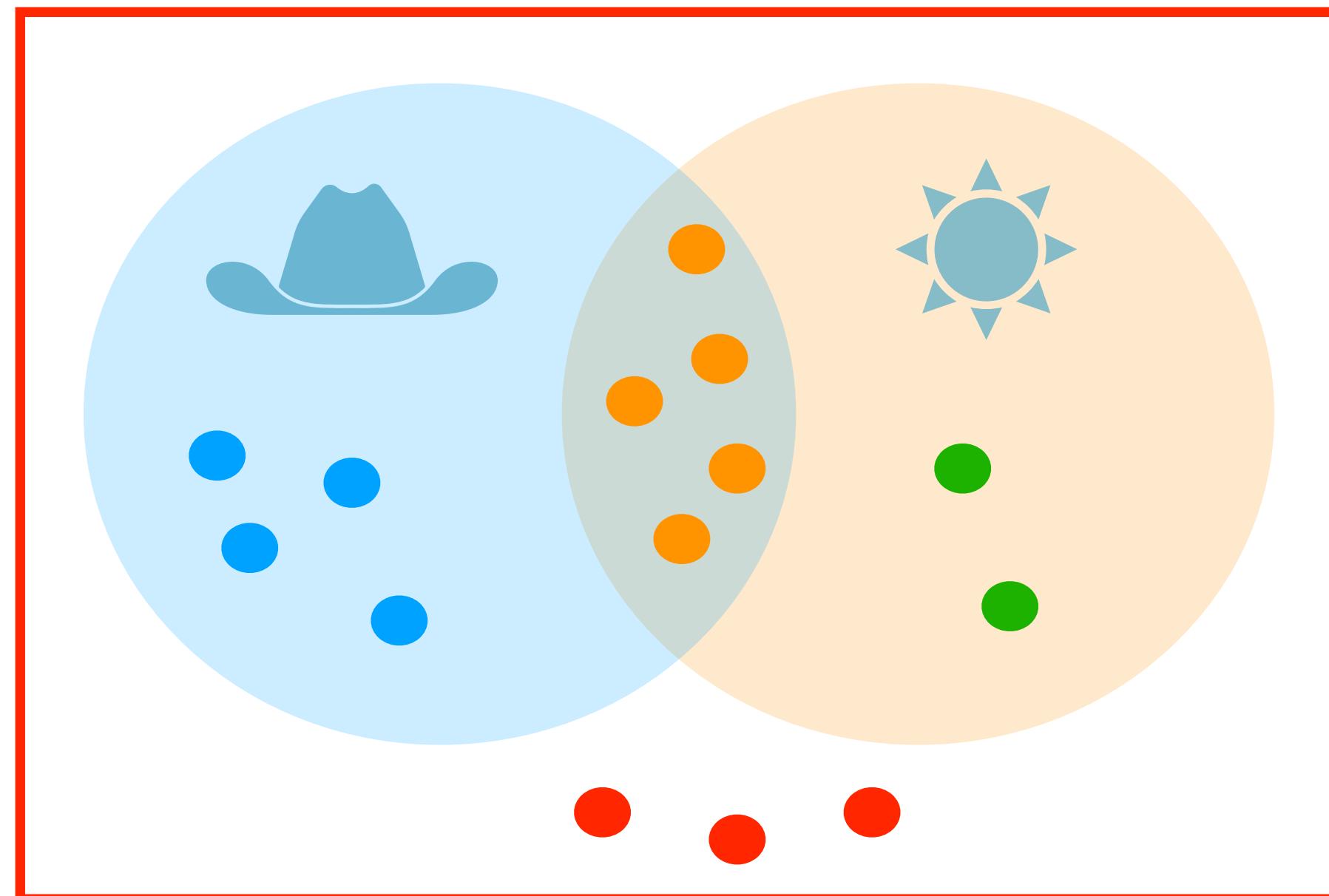
Thinking about probability



	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat	3/14	2/14	5/14
Total	7/14	7/14	

$$p(\text{sun} \& \text{hat} \mid \text{hat}) = \frac{5/14}{9/14} = .55$$

Thinking about probability



Wore hat

Did not
wear hat

Total

No sun

4/14

3/14

7/14

Sun

5/14

2/14

7/14

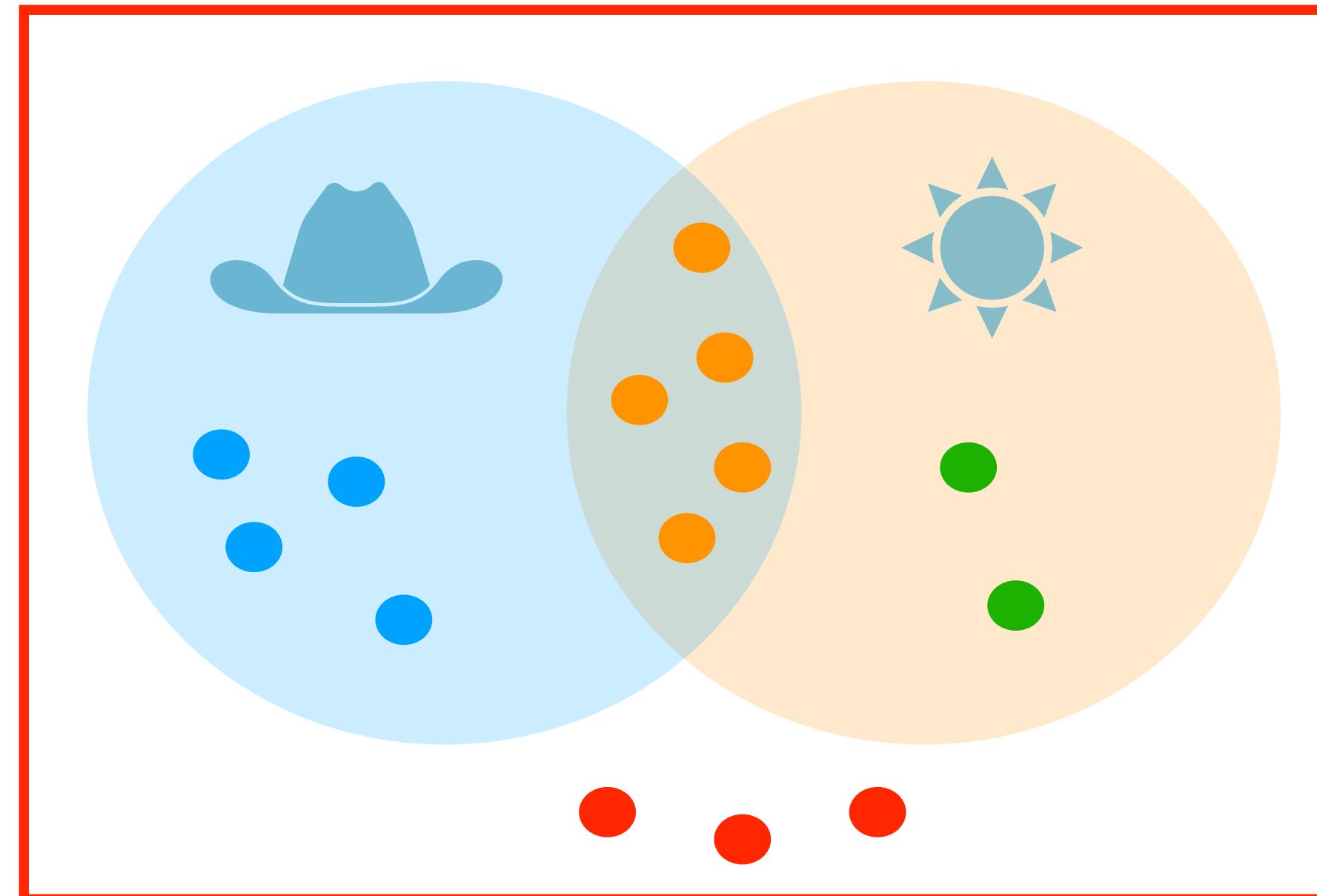
Total

9/14

5/14

$$p(\text{sun} \& \text{hat} \mid \text{sun}) = \frac{\text{P(Sun and Hat)}}{\text{P(Sun)}}$$

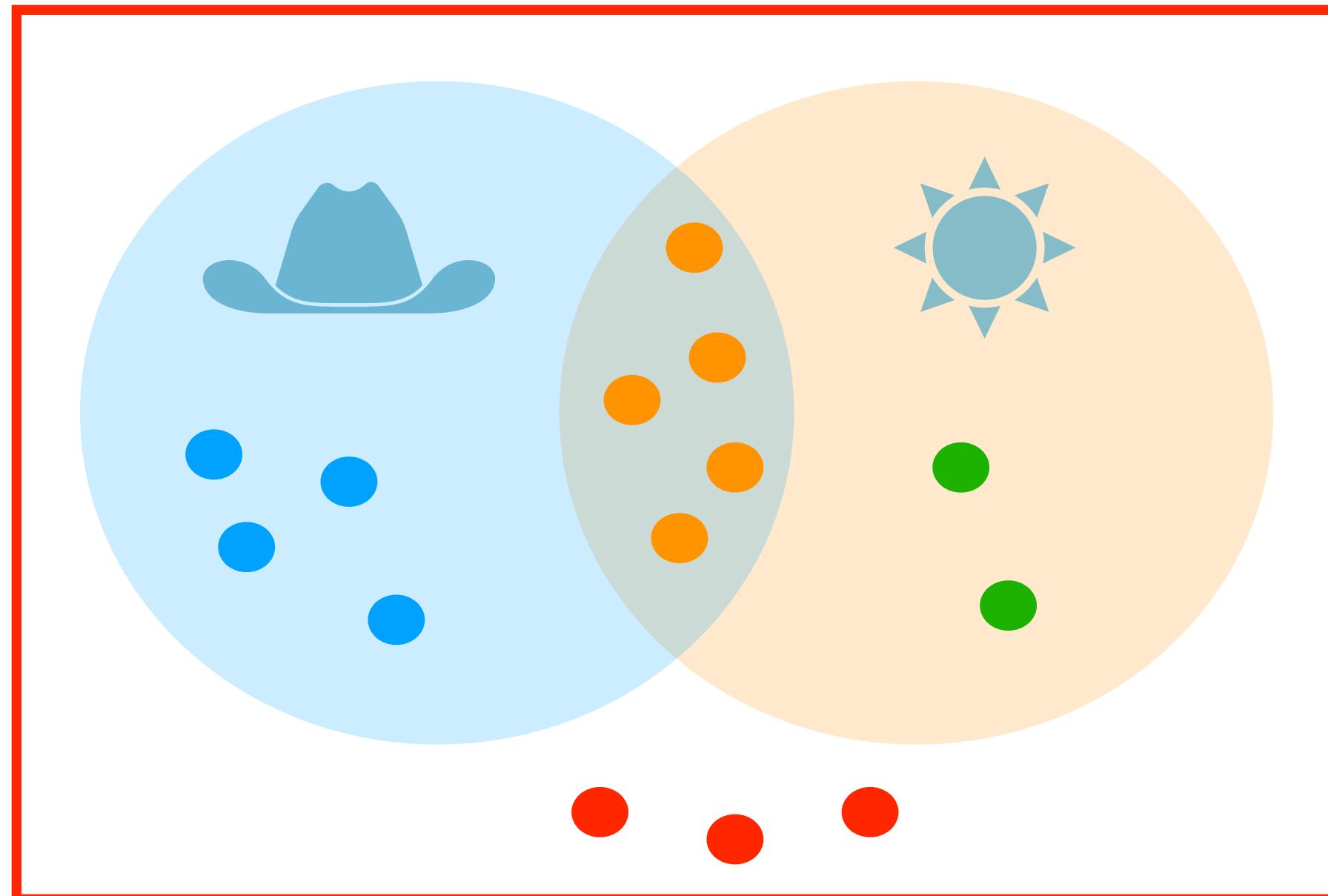
Thinking about probability



	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat	3/14	2/14	5/14
Total	7/14	7/14	

$$p(\text{sun} \& \text{hat} \mid \text{sun}) = \frac{p(\text{sun} \& \text{hat})}{p(\text{sun})}$$

Thinking about probability



	No sun	Sun	Total
Wore hat	4/14	5/14	9/14
Did not wear hat	3/14	2/14	5/14
Total	7/14	7/14	

$$p(\text{sun \& hat} \mid \text{sun}) = \frac{5/14}{7/14} = .71$$

Thinking about probability

$$p(\text{sun} \& \text{hat} \mid \text{hat}) =$$

$$\frac{p(\text{sun} \& \text{hat})}{p(\text{hat})}$$

$$p(\text{sun} \& \text{hat} \mid \text{sun}) =$$

$$\frac{p(\text{sun} \& \text{hat})}{p(\text{sun})}$$

Thinking about probability

$$p(\text{sun} \& \text{hat} \mid \text{hat}) =$$

$$\frac{p(\text{sun} \& \text{hat})}{p(\text{hat})}$$

$$p(\text{sun} \& \text{hat} \mid \text{sun}) =$$

$$\frac{p(\text{sun} \& \text{hat})}{p(\text{sun})}$$

We don't usually know the probabilities of both events, so we ask, is it possible to estimate the conditional p, without that data?

Thinking about probability

$$p(\text{sun} \& \text{hat} \mid \text{hat}) * p(\text{hat}) = \text{p(sun \& hat)}$$

$$p(\text{sun} \& \text{hat} \mid \text{sun}) * p(\text{sun}) = \text{p(sun \& hat)}$$

Thinking about probability

$$p(\text{sun} \ \& \ \text{hat} \mid \text{hat}) * p(\text{hat}) = p(\text{sun} \ \& \ \text{hat} \mid \text{sun}) * p(\text{sun})$$

Thinking about probability

$$p(\text{sun} \& \text{hat} \mid \text{hat}) * p(\text{hat}) = p(\text{sun} \& \text{hat} \mid \text{sun}) * p(\text{sun})$$

$$p(\text{sun} \& \text{hat} \mid \text{hat}) = \frac{p(\text{sun} \& \text{hat} \mid \text{sun}) * p(\text{sun})}{p(\text{hat})}$$

Thinking about probability

$$p(\text{sun} \& \text{hat} \mid \text{hat}) * p(\text{hat}) = p(\text{sun} \& \text{hat} \mid \text{sun}) * p(\text{sun})$$



$$p(\text{sun} \& \text{hat} \mid \text{hat}) = \frac{p(\text{sun} \& \text{hat} \mid \text{sun}) * p(\text{sun})}{p(\text{hat})}$$

Thinking about probability

$$p(\text{sun \& hat} \mid \text{hat}) = \frac{p(\text{sun \& hat} \mid \text{sun}) * p(\text{sun})}{p(\text{hat})}$$

== Bayes Theorem! 🎉🎉🎉

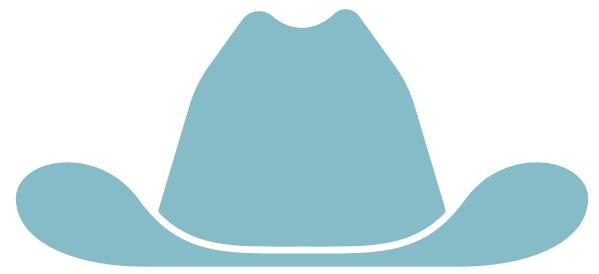
Thinking about probability

$$p(\text{sun} \& \cancel{\text{hat}} \mid \text{hat}) = \frac{p(\cancel{\text{sun}} \& \text{hat} \mid \text{sun}) * p(\text{sun})}{p(\text{hat})}$$

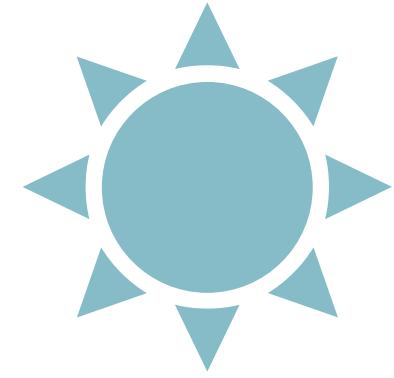
Thinking about probability

$$p(\text{sun} \mid \text{hat}) = \frac{p(\text{hat} \mid \text{sun}) * p(\text{sun})}{p(\text{hat})}$$

Thinking about probability



evidence/
data



hypothesis/
parameter

Thinking about probability

$$p(\text{parameter} \mid \text{data}) = \frac{p(\text{data} \mid \text{parameter}) * p(\text{parameter})}{p(\text{data})}$$

Thinking about probability

$$p(\text{parameter} \mid \text{data}) = \frac{p(\text{data} \mid \text{parameter}) * p(\text{parameter})}{p(\text{data})}$$

IV. Bayesian analysis

Bayesian analysis

$$p(\text{parameter} \mid \text{data}) = \frac{p(\text{data} \mid \text{parameter}) * p(\text{parameter})}{p(\text{data})}$$

Likelihood

.....

$$p(\text{parameter} \mid \text{data}) = \frac{p(\text{data} \mid \text{parameter}) * p(\text{parameter})}{p(\text{data})}$$

What is the probability of observing this data, given
this parameter value?

Bayesian analysis

$$p(\text{parameter} \mid \text{data}) = \frac{p(\text{data} \mid \text{parameter}) * p(\text{parameter})}{p(\text{data})}$$

The equation is displayed with two dotted arrows above it. The left arrow, colored blue, points from the term $p(\text{parameter})$ to the term $p(\text{parameter} \mid \text{data})$. The right arrow, colored red, points from the term $p(\text{data} \mid \text{parameter})$ to the term $p(\text{data})$. The word "Posterior" is written in blue above the first arrow, and the word "Likelihood" is written in red above the second arrow.

What is the probability of a given parameter value, given the data we have observed?

Bayesian analysis

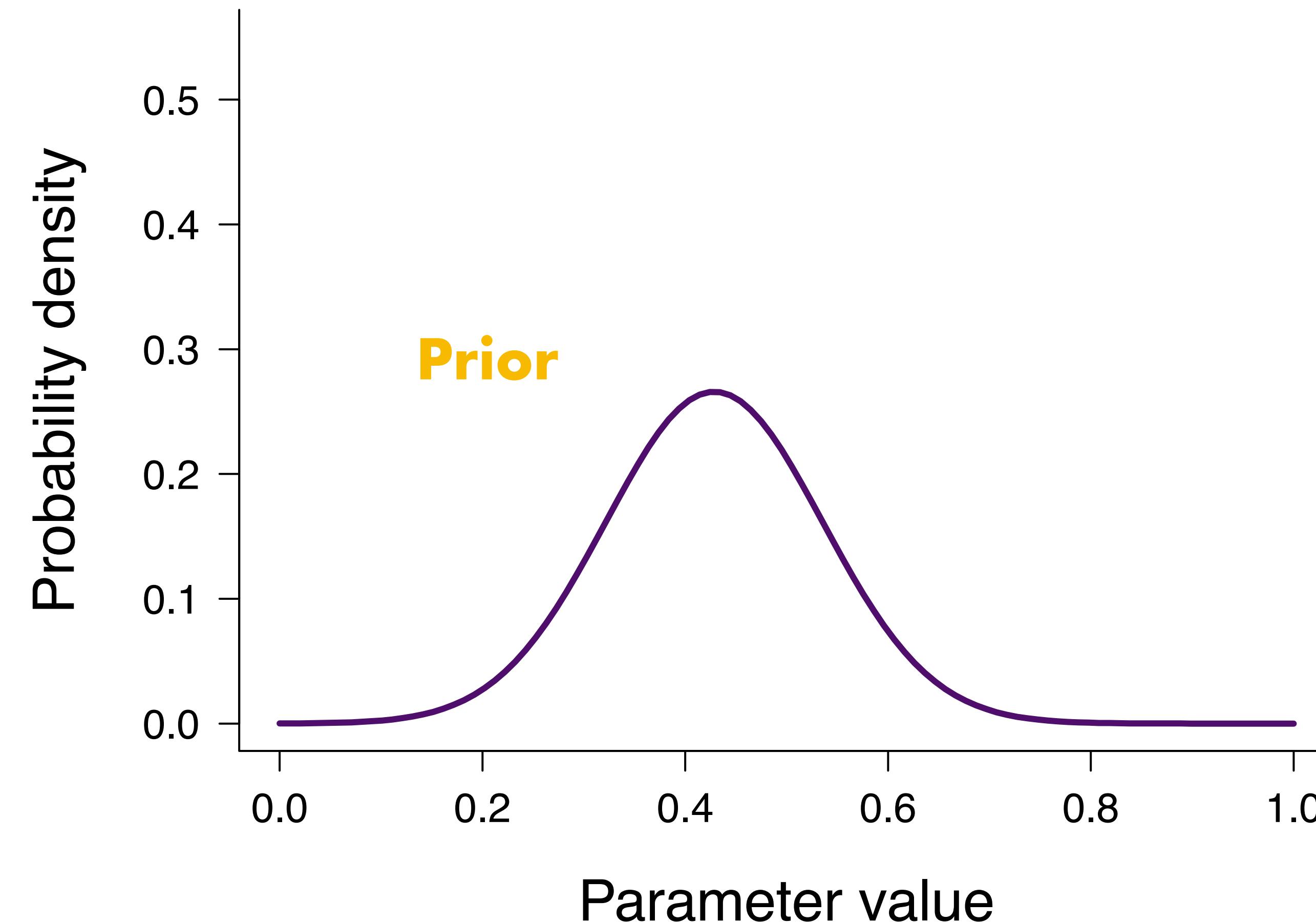
$$p(\text{parameter} \mid \text{data}) = \frac{p(\text{data} \mid \text{parameter}) * p(\text{parameter})}{p(\text{data})}$$

Posterior **Likelihood** **Prior**

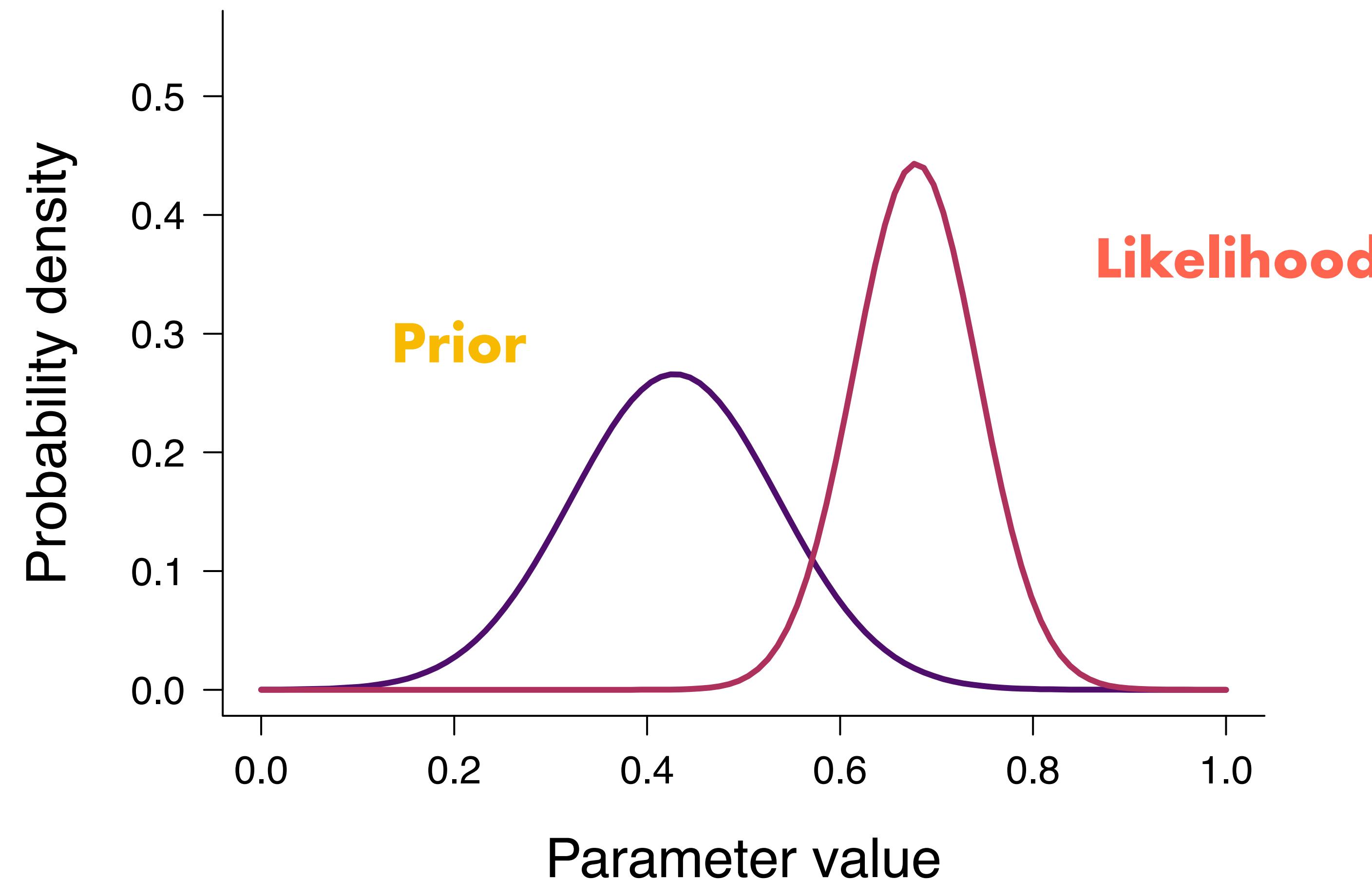
The diagram shows three bell-shaped curves side-by-side. The first curve on the left is blue and labeled 'Posterior'. The middle curve is red and labeled 'Likelihood'. The third curve on the right is yellow and labeled 'Prior'.

What is our expectation for the probability of the parameter?

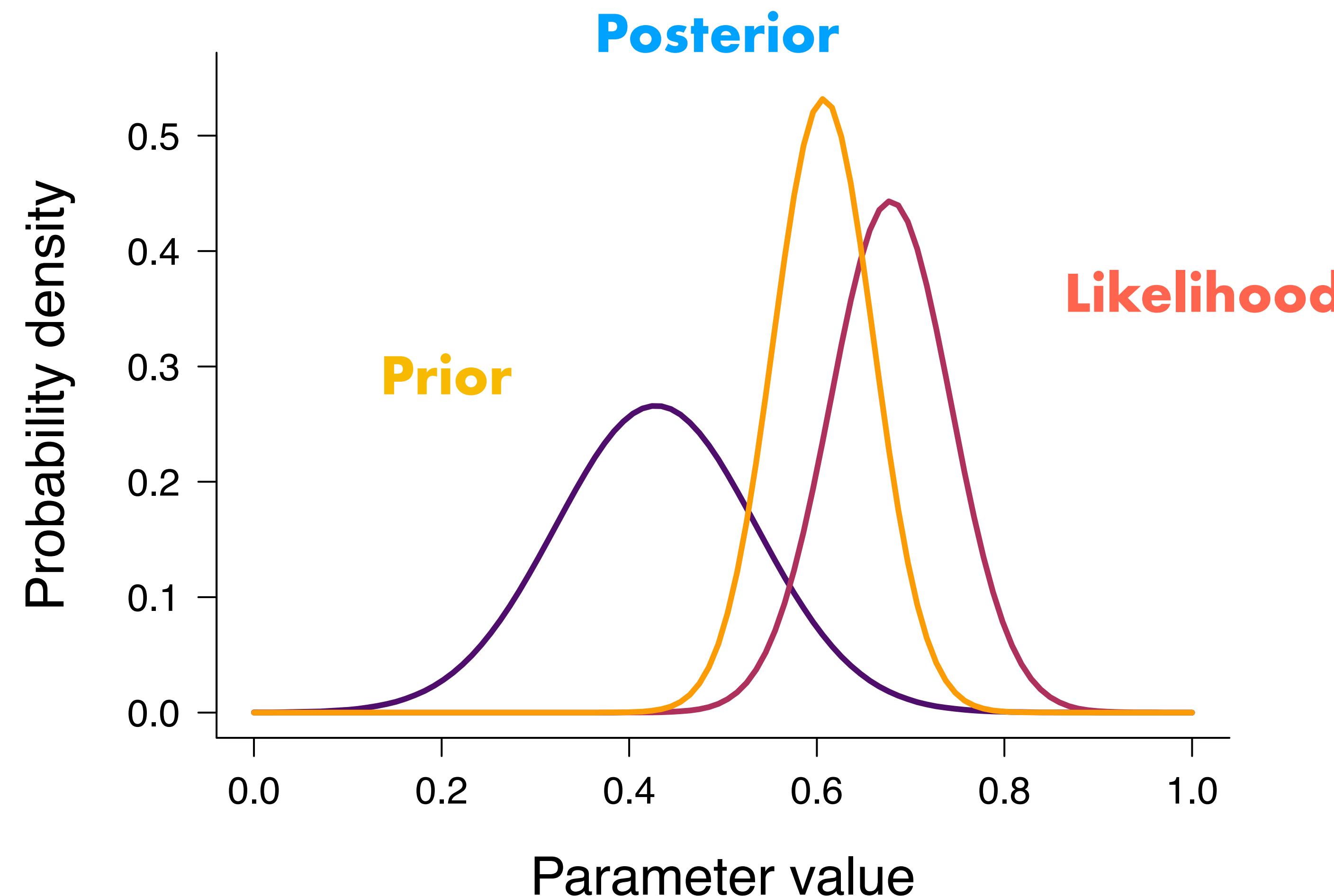
We have three quantities of interest



Bayesian analysis



Bayesian analysis



Use prior knowledge to define *prior distributions*

Observe data and calculate the *likelihood*

Apply Bayes' theorem to estimate *posterior distributions*

**Exercise 3: Explore impact of likelihood and prior
on posterior probability with online tool**

Defining a prior

Defining a prior

- vague/minimally informative

Defining a prior

- vague/minimally informative
- subjective/expert opinion

Defining a prior

- vague/minimally informative
- subjective/expert opinion
- estimate from previous data

Defining a prior

- vague/minimally informative
- subjective/expert opinion
- estimate from previous data

DON'T define your prior based on examining your data

Exercise 4: Defining a prior

Instead of a single value estimate for a parameter, we have an entire probability distribution across all unknowns.

Bayesian analysis

$$y = a + bx + \epsilon$$

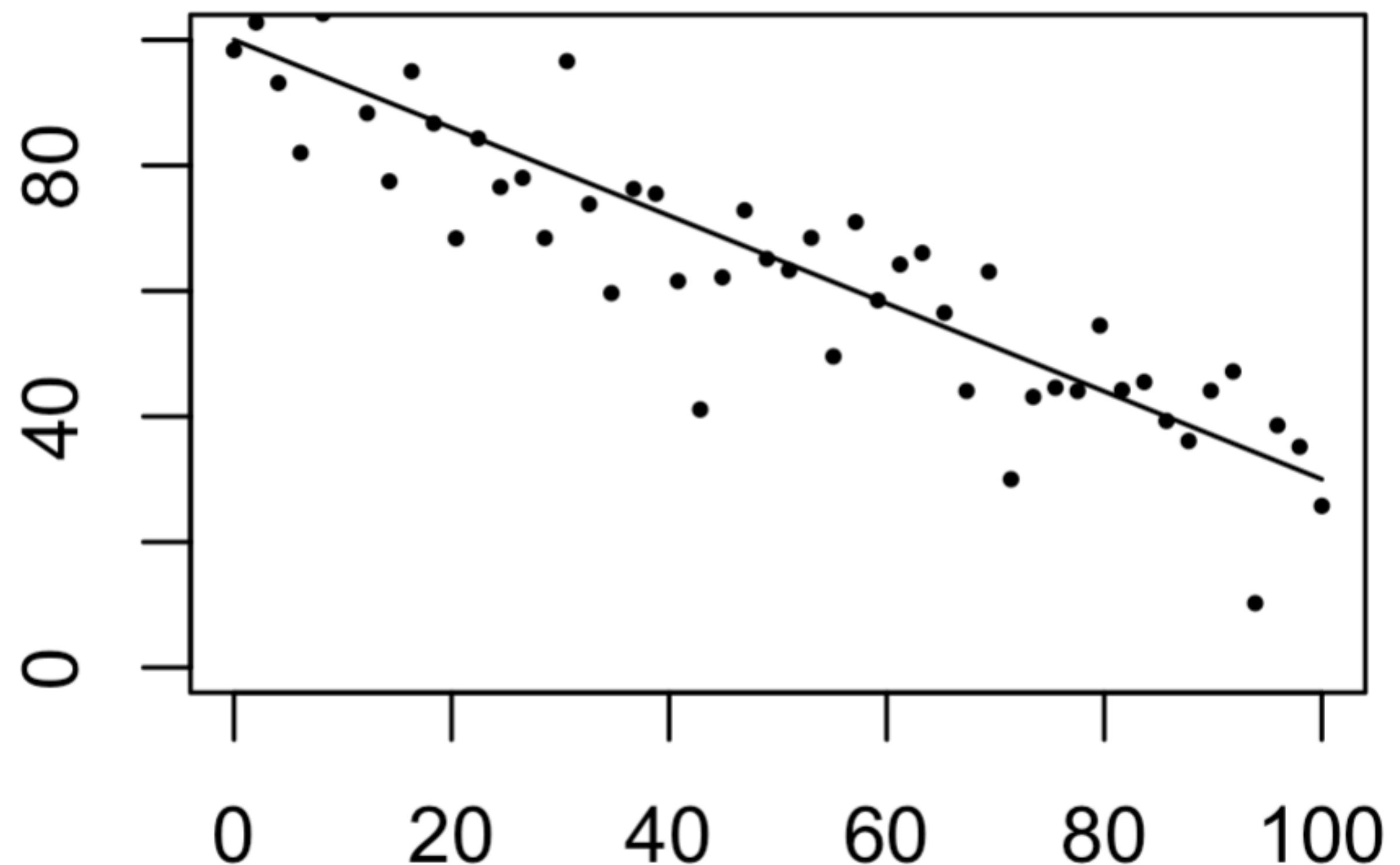
$$\epsilon \sim N(0, \sigma^2)$$

MLE

$$b = -0.67$$

Posterior mean

$$b = -0.67$$



Bayesian inference is intuitive

Frequentist inference is confusing

Bayesian inference is intuitive

Frequentist inference is confusing

Exercise 5: Confidence intervals

Frequentist confidence interval

“Were this procedure to be repeated on numerous samples, the fraction of calculated confidence intervals (which would differ for each sample) that encompass the true population parameter would tend toward 90%.”

Frequentist confidence interval

“Were this procedure to be repeated on numerous samples, the fraction of calculated confidence intervals (which would differ for each sample) that encompass the true population parameter would tend toward 90%.”

Bayesian credible interval

“An interval within which an unobserved parameter value falls with a particular probability.”

Frequentist approach:

- No place for “prior beliefs”
- Inference should only depend on the data (likelihood)
- Probability is the same as frequency
- Point estimate

Frequentist approach:

- No place for “prior beliefs”
- Inference should only depend on the data (likelihood)
- Probability is the same as frequency
- Point estimate

Bayesian approach:

- Inference depends on prior knowledge and available data
- Probability is subjective; it is a degree of belief
- It is more intuitive! “I am 95% certain that...”

V. Introduction to brms

The diagram illustrates the components of Bayesian posterior probability. It features three horizontal dotted lines: a blue line labeled "Posterior", a red line labeled "Likelihood", and a yellow line labeled "Prior". Below these lines is a mathematical equation:

$$p(\text{parameter} \mid \text{data}) = \frac{p(\text{data} \mid \text{parameter}) * p(\text{parameter})}{p(\text{data})}$$

$$p(\text{parameter} \mid \text{data}) = \frac{p(\text{data} \mid \text{parameter}) * p(\text{parameter})}{\text{Likelihood}} \cdot \text{Posterior} \cdot \text{Prior}$$

The equation illustrates the Bayesian formula for updating prior beliefs. The posterior distribution is proportional to the product of the likelihood (the probability of the data given the parameter) and the prior (the initial belief about the parameter). The terms "Posterior", "Likelihood", and "Prior" are represented by dotted lines of increasing length and color (blue, red, yellow) from left to right.

The diagram illustrates the formula for the posterior distribution. It features three colored dots at the top: blue for 'Posterior', red for 'Likelihood', and yellow for 'Prior'. Below each dot is a dotted line of the same color, representing a probability density function. The blue dotted line is flat, the red is bell-shaped, and the yellow is a sharp peak. Below the dots is the mathematical equation:

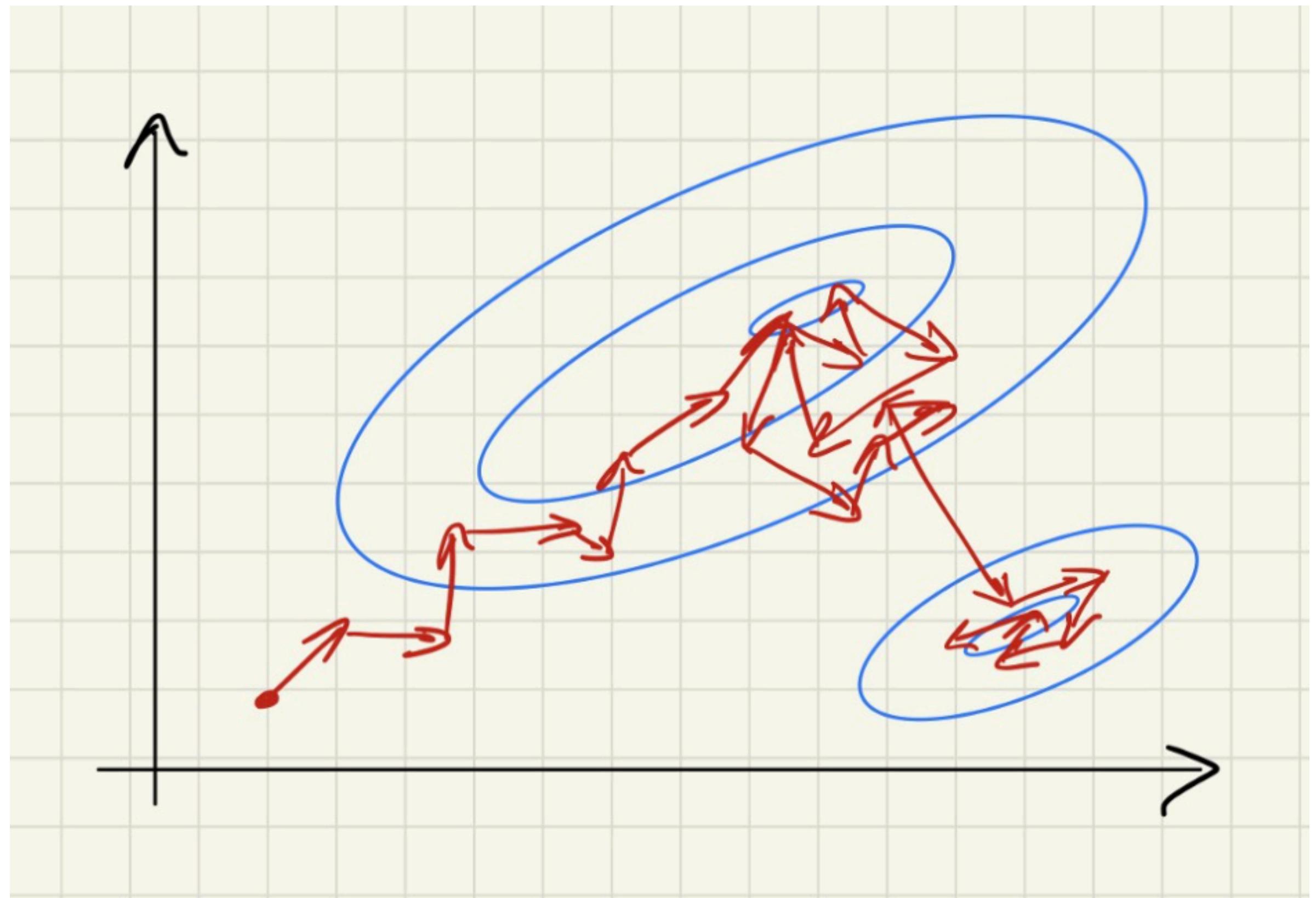
$$p(\text{parameter} \mid \text{data}) \propto p(\text{data} \mid \text{parameter}) * p(\text{parameter})$$

We can approximate the posterior by drawing a large random sample from the distribution using Markov chain Monte Carlo (MCMC)

Making a posterior

Stan and brms

1. Sample a starting point uniformly from prior distribution
2. Calculate probability density function (pdf) at that point.
3. Propose a new sample by stepping away according to some state transition function.
4. Calculate new pdf.
5. If new pdf is higher than old pdf, accept the step.
6. If less, ... reject add current position as new sample and take new step.



Animated MCMC

Stan and brms

'brms()' is a software package that leverages lme4-like syntax to make implementation of Stan functionality more accessible.

Stan and brms

'Stan' is a software package that comes with a programming language to implement MCMC.

Stan and brms

'Stan' is a software package that comes with a programming language to implement MCMC.

Stan uses Hamiltonian Monte Carlo and No-U-Turn Sampler (NUTS) algorithms to implement the MCMC sampling.

Stan and brms

Chains:

Iterations:

Warmup:

Thin:

Draws:

Stan and brms

Chains: Number of Markov chains

Iterations:

Warmup:

Thin:

Draws:

Stan and brms

Chains: Number of Markov chains

Iterations: Number of steps per chain

Warmup:

Thin:

Draws:

Stan and brms

Chains: Number of Markov chains

Iterations: Number of steps per chain

Warmup: First walks around parameter space that you throw away as the chain searches for the right area.

Thin:

Draws:

- Chains:** Number of Markov chains
- Iterations:** Number of steps per chain
- Warmup:** First walks around parameter space that you throw away as the chain searches for the right area.
- Thin:** Prevents correlation between steps by removing steps at this rate.
- Draws:**

Stan and brms

- Chains:** Number of Markov chains
- Iterations:** Number of steps per chain
- Warmup:** First walks around parameter space that you throw away as the chain searches for the right area.
- Thin:** Prevents correlation between steps by removing steps at this rate.
- Draws:** $(\text{Iterations}-\text{warmup}) \times \text{chains}$

Convergence: Do the samples in the chains converge in to the same maxima of the posterior distribution

1. Whether each chain converges on an estimate
2. Whether all chains converge on the same estimate

Convergence: Do the samples in the chains converge in to the same maxima of the posterior distribution

- Traceplot
- Rhat
- Effective sample size

Convergence: Do the samples in the chains converge in to the same maxima of the posterior distribution

- Traceplot
- Rhat
- Effective sample size

R practical: introduction to brms() & diagnosing convergence