

Manual Práctico de easystats en R

Análisis Estadístico Aplicado a Datos Abiertos del Gobierno del
Perú

Encuesta de Intenciones de Siembra 2023–2024
Región Cusco

Material didáctico para usuarios principiantes

December 15, 2025

Contents

1	Introducción	7
1.1	¿Qué es R y por qué usarlo en análisis estadístico?	7
1.1.1	Ventajas de R para análisis estadístico	7
1.2	¿Qué es el ecosistema easystats?	7
1.3	¿Para qué sirve easystats dentro de un flujo de análisis de datos?	7
1.4	Ventajas de easystats frente a enfoques tradicionales	8
1.4.1	Comparación con métodos tradicionales	8
1.4.2	Ventajas clave de easystats	8
2	El ecosistema easystats	9
2.1	Paquetes principales	9
2.1.1	insight: Acceso unificado a modelos	9
2.1.2	parameters: Extracción y descripción de parámetros	9
2.1.3	performance: Evaluación de modelos	10
2.1.4	effectsize: Tamaños de efecto	10
2.1.5	modelbased: Estimaciones basadas en modelos	10
2.1.6	see: Visualización	11
2.1.7	report: Reportes automáticos	11
2.2	Flujo de trabajo integrado	11
3	Descripción de los datasets utilizados	12
3.1	Procedencia y contexto	12
3.1.1	Datos Abiertos del Gobierno del Perú	12
3.2	Encuesta de Intenciones de Siembra 2023–2024	12
3.2.1	Descripción general	12
3.2.2	Características técnicas	13
3.3	Estructura del dataset	13
3.3.1	Variables geográficas	13
3.3.2	Variables agrícolas	13
3.3.3	Tipos de variables	14
3.4	Justificación del dataset para análisis estadístico	14
4	Preparación de los datos	15
4.1	Configuración inicial	15
4.2	Carga de datos	15
4.2.1	Lectura desde archivo CSV	15
4.3	Exploración inicial	16
4.3.1	Diagnóstico de datos con datawizard	16
4.4	Limpieza de datos	17
4.4.1	Manejo de nombres de columnas	17
4.4.2	Conversión de tipos de datos	17
4.4.3	Manejo de valores faltantes	18
4.5	Creación de variables derivadas	19
4.6	Transformación de datos	19
4.6.1	De formato ancho a largo	20

4.7	Verificación final	20
4.8	Resumen de la preparación	20
5	Caso 1: Regresión Lineal	21
5.1	Objetivo del análisis	21
5.2	Pregunta estadística	21
5.3	Preparación de datos específica	21
5.4	Modelo estadístico propuesto	22
5.5	Ajuste del modelo en R	22
5.6	Análisis con easystats	23
5.6.1	Extracción de parámetros	23
5.6.2	Evaluación del desempeño del modelo	23
5.6.3	Verificación de supuestos	24
5.7	Cálculo de tamaños de efecto	25
5.8	Visualización de resultados	25
5.8.1	Gráfico de coeficientes	25
5.8.2	Gráficos de efectos predichos	26
5.9	Reporte automático	26
5.10	Interpretación de resultados	26
5.11	Comparación de modelos	27
5.12	Predicción con el modelo	28
5.13	Resumen del flujo de regresión lineal	28
6	Caso 2: Regresión Logística	29
6.1	Objetivo del análisis	29
6.2	Pregunta estadística	29
6.3	Preparación de datos específica	29
6.4	Modelo estadístico propuesto	30
6.5	Ajuste del modelo en R	30
6.6	Análisis con easystats	31
6.6.1	Extracción de parámetros	31
6.6.2	Métricas de desempeño	31
6.6.3	Matriz de confusión y métricas de clasificación	32
6.6.4	Verificación de supuestos	33
6.7	Interpretación de Odds Ratios	33
6.8	Predicción de probabilidades	34
6.9	Efectos marginales	34
6.10	Visualización de resultados	35
6.11	Reporte automático	35
6.12	Comparación de modelos logísticos	36
6.13	Validación cruzada	36
6.14	Resumen del flujo de regresión logística	37
7	Caso 3: Comparación de Grupos	38
7.1	Objetivo del análisis	38
7.2	Preguntas estadísticas	38
7.3	Preparación de datos	38

7.4	Caso 3.1: Prueba t para dos grupos	39
7.4.1	Prueba t independiente	39
7.4.2	Prueba t pareada	40
7.5	Caso 3.2: ANOVA de una vía	40
7.5.1	Tamaño de efecto para ANOVA	41
7.6	Caso 3.3: ANOVA de dos vías (factorial)	41
7.6.1	Interpretación de interacciones	42
7.7	Verificación de supuestos de ANOVA	43
7.8	Alternativas no paramétricas	43
7.9	Comparaciones múltiples y corrección por múltiples pruebas	44
7.10	Visualizaciones avanzadas	44
7.11	Reporte automático	45
7.12	Ejemplo de interpretación integrada	45
7.13	Resumen del flujo de comparación de grupos	46
8	Caso 4: Modelos con Interacciones	46
8.1	Objetivo del análisis	46
8.2	Pregunta estadística	46
8.3	Modelo estadístico con interacción	47
8.4	Preparación de datos	47
8.5	Modelos con y sin interacción	48
8.6	Análisis del modelo con interacción	48
8.7	Interpretación de interacciones	48
8.7.1	Efectos simples (Simple Effects)	48
8.8	Visualización de interacciones	49
8.9	Descomposición de la interacción	50
8.10	Interacciones continuas × categóricas	51
8.11	Interacciones continuas × continuas	52
8.12	Cuándo incluir interacciones	53
8.13	Verificación de supuestos	53
8.14	Centrado de variables	54
8.15	Reporte de modelos con interacción	54
8.16	Ejemplo de interpretación integrada	55
8.17	Resumen del análisis de interacciones	55
9	Caso 5: Modelos de Efectos Mixtos	55
9.1	Objetivo del análisis	56
9.2	¿Cuándo usar modelos mixtos?	56
9.3	Estructura de datos jerárquicos	56
9.4	Modelo estadístico de efectos mixtos	57
9.4.1	Modelo con intercepto aleatorio	57
9.4.2	Modelo con pendiente aleatoria	57
9.5	Ajuste de modelos mixtos en R	57
9.5.1	Modelo nulo (solo intercepto aleatorio)	58
9.5.2	Modelo con predictores de nivel 1	59
9.5.3	Modelo con predictores de nivel 2	59
9.5.4	Modelo con pendiente aleatoria	60

9.5.5	Estructura de efectos aleatorios anidados	60
9.6	Análisis de resultados	61
9.6.1	Extracción de efectos fijos	61
9.6.2	Extracción de efectos aleatorios	61
9.6.3	Componentes de varianza	61
9.7	Evaluación del modelo	62
9.8	Predicciones con modelos mixtos	63
9.9	Visualización de modelos mixtos	63
9.10	Comparación de modelos mixtos	64
9.11	Modelos mixtos generalizados (GLMM)	65
9.12	Reporte de modelos mixtos	65
9.13	Ejemplo de interpretación integrada	66
9.14	Resumen del flujo de modelos mixtos	66
10	Visualización y Comunicación de Resultados	67
10.1	Filosofía de visualización en easystats	67
10.2	Gráficos de parámetros	67
10.3	Gráficos de diagnóstico	68
10.4	Gráficos de efectos y predicciones	68
10.5	Gráficos de comparación de modelos	69
10.6	Gráficos para modelos logísticos	69
10.7	Paletas de colores y temas	70
10.8	Exportación de gráficos	70
10.9	Gráficos combinados (multipanel)	71
10.10	Tablas para publicación	72
10.11	Reportes reproducibles con R Markdown	72
11	Buenas Prácticas y Errores Comunes	73
11.1	Buenas prácticas generales	73
11.1.1	Documentación y reproducibilidad	73
11.1.2	Comentarios y estilo de código	73
11.1.3	Flujo de análisis sistemático	74
11.2	Errores comunes en análisis estadístico	74
11.2.1	Error 1: P-hacking y HARKing	74
11.2.2	Error 2: Ignorar supuestos del modelo	75
11.2.3	Error 3: Confundir significancia estadística con importancia práctica	75
11.2.4	Error 4: Sobreajuste del modelo	76
11.2.5	Error 5: Ignorar datos faltantes	76
11.2.6	Error 6: Múltiples comparaciones sin ajuste	77
11.3	Buenas prácticas específicas con easystats	78
11.3.1	Usar funciones consistentes	78
11.3.2	Combinar paquetes del ecosistema	78
11.4	Checklist de análisis	79
11.5	Recomendaciones para adaptar a nuevos datasets	79
11.5.1	Lista de verificación al usar nuevos datos	79
11.5.2	Template de análisis reutilizable	80

12 Conclusiones	82
12.1 Resumen del manual	82
12.2 Ventajas de easystats	83
12.3 Próximos pasos	83
12.3.1 Práctica	83
12.3.2 Profundización	83
12.4 Recursos adicionales	84
12.4.1 Documentación oficial de easystats	84
12.4.2 Libros recomendados	84
12.4.3 Cursos en línea	84
12.4.4 Comunidades y foros	84
12.4.5 Datos abiertos en Perú	84
12.5 Palabras finales	85
13 Referencias	85
13.1 Publicaciones de easystats	85
13.2 Libros y artículos de referencia	86
13.3 Recursos en línea	86
A Instalación y configuración	86
A.1 Instalación de R y RStudio	86
A.2 Instalación de easystats	86
A.3 Paquetes complementarios recomendados	87
B Glosario de términos estadísticos	87
C Atajos de teclado útiles en RStudio	89

1 Introducción

1.1 ¿Qué es R y por qué usarlo en análisis estadístico?

R es un lenguaje de programación y entorno de software libre especializado en análisis estadístico y visualización de datos. Desde su creación en 1993 por Ross Ihaka y Robert Gentleman en la Universidad de Auckland, R se ha convertido en una herramienta fundamental en ciencia de datos, investigación académica, y análisis aplicado.

1.1.1 Ventajas de R para análisis estadístico

R ofrece múltiples ventajas que lo hacen ideal para el análisis de datos:

- **Gratis y de código abierto:** No requiere licencias costosas y su código es auditable por la comunidad.
- **Amplia comunidad:** Miles de paquetes desarrollados por expertos en diversas áreas.
- **Reproducibilidad:** Los análisis pueden ser compartidos y replicados fácilmente.
- **Flexibilidad:** Desde estadística básica hasta machine learning avanzado.
- **Visualización de alta calidad:** Gráficos publicables para revistas académicas.
- **Integración:** Funciona con bases de datos, APIs, y otros lenguajes de programación.

1.2 ¿Qué es el ecosistema easystats?

El ecosistema `easystats` es una colección integrada de paquetes de R diseñada para facilitar el análisis estadístico. Desarrollado por Dominique Makowski, Daniel Lüdecke, Mattan S. Ben-Shachar y colaboradores, `easystats` busca hacer la estadística más accesible, intuitiva y reproducible.

A diferencia de enfoques tradicionales que requieren conocer múltiples funciones dispersas en diferentes paquetes, `easystats` proporciona una interfaz coherente y consistente para todo el flujo de análisis estadístico.

1.3 ¿Para qué sirve easystats dentro de un flujo de análisis de datos?

El ecosistema `easystats` cubre las cinco etapas fundamentales del análisis estadístico:

1. **Modelamiento:** Ajustar modelos estadísticos con sintaxis estándar de R
2. **Evaluación:** Verificar supuestos y desempeño del modelo
3. **Interpretación:** Extraer y comprender parámetros y efectos
4. **Visualización:** Crear gráficos diagnósticos y de resultados

5. Reporte:

Generar descripciones textuales automáticas de resultados

Easystats actúa como un “traductor” que convierte objetos complejos de modelos estadísticos en tablas, gráficos y texto comprensibles, facilitando tanto el análisis como la comunicación de resultados.

1.4 Ventajas de easystats frente a enfoques tradicionales

1.4.1 Comparación con métodos tradicionales

En R tradicional, obtener un análisis completo requiere conocer funciones de múltiples paquetes:

Listing 1: Enfoque tradicional (fragmentado)

```

1 # Ajustar modelo
2 modelo <- lm(y ~ x1 + x2, data = datos)
3
4 # Ver resultados (base R)
5 summary(modelo)
6
7 # Calcular R^2 ajustado manualmente
8 # Extraer residuos con resid()
9 # Calcular IC con confint()
10 # Hacer diagnósticos con plot()
11 # Calcular tamaños de efecto manualmente
12 # Escribir interpretación manualmente

```

Con easystats, el mismo análisis es más intuitivo:

Listing 2: Enfoque easystats (integrado)

```

1 library(easystats)
2
3 # Ajustar modelo
4 modelo <- lm(y ~ x1 + x2, data = datos)
5
6 # Análisis completo
7 parameters(modelo)    # Parámetros con IC
8 performance(modelo)   # Métricas de bondad de ajuste
9 check_model(modelo)    # Diagnósticos visuales
10 effectsize(modelo)    # Tamaños de efecto
11 report(modelo)        # Reporte en lenguaje natural

```

1.4.2 Ventajas clave de easystats

- **Consistencia:** La misma sintaxis funciona para regresión lineal, logística, modelos mixtos, etc.
- **Automatización:** Los reportes se generan automáticamente en formato APA.
- **Visualización integrada:** Gráficos diagnósticos con una sola función.

- **Interpretabilidad:** Resultados en lenguaje claro, no solo números.
- **Mejores prácticas:** Incorpora recomendaciones estadísticas modernas.
- **Interoperabilidad:** Funciona con cientos de tipos de modelos existentes.

2 El ecosistema easystats

El ecosistema easystats está compuesto por varios paquetes especializados que trabajan conjuntamente. Cada paquete tiene un propósito específico pero comparte una filosofía común de diseño.

2.1 Paquetes principales

2.1.1 insight: Acceso unificado a modelos

El paquete `insight` es la base del ecosistema. Proporciona una interfaz consistente para extraer información de más de 200 tipos de modelos estadísticos diferentes.

Funciones clave:

- `get_data()`: Extrae los datos utilizados en el modelo
- `find_variables()`: Identifica variables del modelo
- `get_parameters()`: Obtiene coeficientes del modelo
- `model_info()`: Información sobre el tipo de modelo

Utilidad: Permite que los demás paquetes de easystats funcionen con cualquier tipo de modelo sin necesidad de código específico para cada uno.

2.1.2 parameters: Extracción y descripción de parámetros

El paquete `parameters` extrae, calcula y presenta los parámetros de modelos estadísticos de forma estandarizada.

Funciones clave:

- `model_parameters()`: Tabla completa de parámetros con intervalos de confianza
- `standardize_parameters()`: Coeficientes estandarizados
- `p_value()`: Valores p con diferentes métodos
- `bootstrap_parameters()`: Parámetros mediante bootstrap

Utilidad: Proporciona estimaciones, errores estándar, intervalos de confianza y estadísticos de prueba en un formato consistente y legible.

2.1.3 performance: Evaluación de modelos

El paquete `performance` evalúa la calidad y adecuación de modelos estadísticos.

Funciones clave:

- `model_performance()`: Índices de bondad de ajuste (R^2 , AIC, RMSE, etc.)
- `check_model()`: Diagnósticos visuales completos
- `check_normality()`: Prueba de normalidad de residuos
- `check_heteroscedasticity()`: Prueba de homogeneidad de varianzas
- `compare_performance()`: Comparación entre múltiples modelos

Utilidad: Verifica supuestos estadísticos y cuantifica qué tan bien el modelo se ajusta a los datos.

2.1.4 effectsize: Tamaños de efecto

El paquete `effectsize` calcula medidas estandarizadas de magnitud de efectos.

Funciones clave:

- `standardize()`: Estandariza variables o coeficientes
- `cohens_d()`: d de Cohen para diferencias entre grupos
- `eta_squared()`: Eta cuadrado para ANOVA
- `effectsize()`: Detecta automáticamente el tamaño de efecto apropiado
- `interpret_*`(): Interpreta la magnitud (pequeño, mediano, grande)

Utilidad: Los tamaños de efecto complementan los valores p al indicar la magnitud práctica de las diferencias o relaciones, no solo su significancia estadística.

2.1.5 modelbased: Estimaciones basadas en modelos

El paquete `modelbased` genera predicciones y estimaciones marginales a partir de modelos.

Funciones clave:

- `estimate_means()`: Medias estimadas por grupo
- `estimate_contrasts()`: Contrastes entre niveles
- `estimate_slopes()`: Pendientes (efectos simples)
- `estimate_response()`: Predicciones del modelo

Utilidad: Facilita la interpretación de modelos complejos (con interacciones o efectos no lineales) mediante estimaciones concretas en valores específicos de las variables.

2.1.6 see: Visualización

El paquete `see` extiende `ggplot2` para crear visualizaciones especializadas de resultados estadísticos.

Funciones clave:

- `plot()`: Método genérico que crea gráficos apropiados para cada objeto easystats
- Gráficos de diagnóstico automáticos
- Visualizaciones de efectos e interacciones
- Gráficos de comparación de modelos

Utilidad: Produce gráficos de calidad publicable con configuración mínima, manteniendo la flexibilidad de `ggplot2`.

2.1.7 report: Reportes automáticos

El paquete `report` genera descripciones textuales de análisis estadísticos en formato académico.

Funciones clave:

- `report()`: Genera un reporte completo del análisis
- `report_text()`: Solo el texto narrativo
- `report_table()`: Solo la tabla de resultados
- `report_performance()`: Reporte de bondad de ajuste

Utilidad: Automatiza la redacción de resultados estadísticos siguiendo convenciones académicas (APA), ahorrando tiempo y reduciendo errores de transcripción.

2.2 Flujo de trabajo integrado

El poder de easystats reside en la integración de estos paquetes. El flujo típico es:

1. **Ajustar modelo** con funciones estándar de R (`lm()`, `glm()`, `lmer()`, etc.)
2. **Extraer parámetros** con `parameters()`
3. **Evaluar desempeño** con `performance()` y `check_model()`
4. **Calcular efectos** con `effectsize()` y `modelbased()`
5. **Visualizar** con `plot()` del paquete `see`
6. **Reportar** con `report()`

Este flujo es consistente independientemente del tipo de modelo, lo que facilita el aprendizaje y la aplicación a nuevos problemas.

3 Descripción de los datasets utilizados

3.1 Procedencia y contexto

Los datos utilizados en este manual provienen de la plataforma de Datos Abiertos del Gobierno del Perú, específicamente de la Gerencia Regional de Agricultura del Cusco (GERAGRI).

3.1.1 Datos Abiertos del Gobierno del Perú

El Gobierno del Perú, mediante la Autoridad Nacional de Transparencia y Acceso a la Información Pública, promueve la apertura de datos gubernamentales. Estos datos son:

- De acceso público y gratuito
- Legalmente reutilizables para cualquier propósito
- Actualizados periódicamente
- Disponibles en formatos procesables (CSV, JSON, etc.)

Acceso: <https://www.datosabiertos.gob.pe/>

3.2 Encuesta de Intenciones de Siembra 2023–2024

3.2.1 Descripción general

La Encuesta de Intenciones de Siembra es un instrumento estadístico aplicado por las Direcciones y Gerencias Regionales de Agricultura para estimar las áreas que los agricultores planean sembrar en la próxima campaña agrícola.

Objetivo del dataset: Proporcionar información oportuna sobre las intenciones de siembra de cultivos por distrito, cultivo y mes de siembra, para:

- Planificación agrícola regional
- Estimación de producción futura
- Política de seguridad alimentaria
- Asignación de recursos (semillas, fertilizantes, créditos)

3.2.2 Características técnicas

Table 1: Características del dataset

Característica	Descripción
Institución	GERAGRI Cusco
Período	Campaña agrícola 2023–2024
Ámbito geográfico	Región Cusco (13 provincias)
Formato	CSV
Separador	Punto y coma (;
Codificación	ISO-8859-1 (Latin-1)
Número de registros	Variable (miles de observaciones)
Fecha de corte	10 de noviembre de 2025

3.3 Estructura del dataset

El dataset contiene información desagregada geográficamente y por cultivo:

3.3.1 Variables geográficas

- **DEPARTAMENTO**: Nombre del departamento (Cusco)
- **PROVINCIA**: Nombre de la provincia (13 provincias del Cusco)
- **COD_UBIGEO**: Código único de ubicación geográfica (INEI)
- **DISTRITO**: Nombre del distrito
- **SECTOR_ESTADISTICO**: Subdivisión del distrito para fines de encuesta

3.3.2 Variables agrícolas

- **CULTIVO_AGRICOLA**: Nombre del cultivo (papa, maíz, quinua, etc.)
- **CAMPANIA_AGRICOLA**: Período agrícola (20232024)
- **SIEMBRAS_[MES]**: Área de siembra en hectáreas por mes
 - AGO, SET, OCT, NOV, DIC (agosto–diciembre 2023)
 - ENE, FEB, MAR, ABR, MAY, JUN, JUL (enero–julio 2024)
- **FECHA_CORTE**: Fecha de corte de la información

3.3.3 Tipos de variables

Table 2: Tipología de variables

Variable	Tipo	Escala
DEPARTAMENTO	Categórica nominal	—
PROVINCIA	Categórica nominal	—
COD_UBIGEO	Identificador	—
DISTRITO	Categórica nominal	—
SECTOR_ESTADISTICO	Categórica nominal	—
CULTIVO_AGRICOLA	Categórica nominal	—
CAMPANIA_AGRICOLA	Categórica ordinal	—
SIEMBRAS_*	Numérica continua	Razón (ha)
FECHA_CORTE	Fecha	—

3.4 Justificación del dataset para análisis estadístico

Este dataset es ideal para aprender análisis estadístico por varias razones:

1. **Realidad aplicada:** Datos reales con implicaciones en política pública y economía agrícola.
2. **Estructura jerárquica:** Los datos están naturalmente anidados (sectores dentro de distritos, distritos dentro de provincias), lo que permite explorar modelos multinivel.
3. **Variables diversas:** Combina variables categóricas (cultivo, ubicación) y continuas (área), permitiendo múltiples tipos de análisis.
4. **Potencial de análisis:**
 - Modelos de regresión (¿qué factores explican el área de siembra?)
 - Comparaciones entre grupos (¿difiere la siembra entre provincias?)
 - Series temporales (¿cuáles son los patrones estacionales?)
 - Modelos espaciales (¿existe autocorrelación geográfica?)
5. **Interpretabilidad:** Los resultados son comprensibles para no expertos (hectáreas de papa, no variables abstractas).
6. **Relevancia local:** Para usuarios peruanos, trabajar con datos nacionales es más motivador y relevante que datasets extranjeros.
7. **Accesibilidad:** Datos abiertos, legales y gratuitos que pueden ser utilizados sin restricciones.

4 Preparación de los datos

La preparación de datos es un paso crucial que puede representar hasta el 80% del tiempo en un proyecto de análisis. En esta sección aprenderemos a cargar, explorar, limpiar y transformar los datos de forma reproducible.

4.1 Configuración inicial

Antes de comenzar, debemos cargar los paquetes necesarios y configurar nuestro entorno de trabajo.

Listing 3: Configuración inicial del entorno

```

1 # Instalar easystats (solo la primera vez)
2 # install.packages("easystats")
3
4 # Cargar paquetes
5 library(easystats)      # Ecosistema completo
6 library(tidyverse)       # Manipulación de datos
7 library(here)            # Gestión de rutas de archivos
8
9 # Configurar opciones
10 options(scipen = 999)   # Desactivar notación científica
11 options(digits = 3)     # Limitar decimales en salida
12
13 # Verificar versiones
14 sessionInfo()
```

Nota sobre tidyverse: Aunque easystats puede trabajar con R base, utilizaremos funciones de `tidyverse` (particularmente `dplyr` y `tidyr`) para manipulación de datos, ya que su sintaxis es más intuitiva para principiantes.

4.2 Carga de datos

4.2.1 Lectura desde archivo CSV

El dataset está en formato CSV con separador de punto y coma y codificación ISO-8859-1 (Latin-1), común en datos del gobierno peruano.

Listing 4: Carga del dataset desde archivo local

```

1 # Cargar datos con encoding correcto
2 datos_raw <- read.csv2(
3   file = "intencionesSIEMBRA.csv",
4   encoding = "Latin1",
5   stringsAsFactors = FALSE
6 )
7
8 # Alternativa con readr (tidyverse)
9 library(readr)
10 datos_raw <- read_delim(
```

```

11   file = "intencionesSIEMBRA.csv",
12   delim = ";",
13   locale = locale(encoding = "Latin1")
14 )
15
16 # Verificar carga exitosa
17 dim(datos_raw) # Dimensiones: filas y columnas

```

Parámetros importantes:

- `sep = ";"`: El separador es punto y coma, no coma
- `encoding = "Latin1"`: Preserva caracteres especiales (tildes, eñes)
- `stringsAsFactors = FALSE`: Mantiene texto como texto, no factor

4.3 Exploración inicial

Una vez cargados los datos, debemos explorarlos antes de cualquier análisis.

Listing 5: Exploración inicial del dataset

```

1 # Ver primeras filas
2 head(datos_raw, n = 10)
3
4 # Estructura del dataset
5 str(datos_raw)
6
7 # Resumen estadístico básico
8 summary(datos_raw)
9
10 # Nombres de columnas
11 names(datos_raw)
12
13 # Número de observaciones y variables
14 nrow(datos_raw) # Filas (observaciones)
15 ncol(datos_raw) # Columnas (variables)
16
17 # Ver valores únicos de variables categóricas
18 unique(datos_raw$PROVINCIA)
19 unique(datos_raw$CULTIVO_AGRICOLA)
20
21 # Contar observaciones por categoría
22 table(datos_raw$PROVINCIA)
23 table(datos_raw$CULTIVO_AGRICOLA)

```

4.3.1 Diagnóstico de datos con datawizard

El paquete `datawizard` (parte de `easystats`) ofrece funciones especializadas para exploración:

Listing 6: Diagnóstico avanzado con datawizard

```

1 library(datawizard)
2
3 # Resumen descriptivo completo
4 describe_distribution(datos_raw)
5
6 # Identificar valores faltantes
7 data_missing(datos_raw)
8
9 # Identificar variables constantes
10 data_findconstants(datos_raw)
11
12 # Ver tipos de datos
13 data_properties(datos_raw)

```

4.4 Limpieza de datos

4.4.1 Manejo de nombres de columnas

Los nombres de columnas con caracteres especiales pueden causar problemas:

Listing 7: Estandarización de nombres de columnas

```

1 # Ver nombres actuales
2 names(datos_raw)
3
4 # Opción 1: Limpiar nombres con janitor
5 library(janitor)
6 datos <- clean_names(datos_raw)
7
8 # Verificar cambios
9 names(datos)
10 # "sector_estadístico" se convierte en "sector_estadistico"
11
12 # Opción 2: Renombrar manualmente si es necesario
13 datos <- datos %>%
14   rename(
15     departamento = DEPARTAMENTO,
16     provincia = PROVINCIA,
17     cod_ubigeo = COD_UBIGEO,
18     distrito = DISTRITO,
19     sector = SECTOR_ESTADISTICO,
20     cultivo = CULTIVO_AGRICOLA,
21     campania = CAMPANIA_AGRICOLA
22   )

```

4.4.2 Conversión de tipos de datos

Listing 8: Conversión de tipos de datos

```

1 # Convertir variables categóricas a factor
2 datos <- datos %>%
3   mutate(
4     departamento = as.factor(departamento),
5     provincia = as.factor(provincia),
6     distrito = as.factor(distrito),
7     cultivo = as.factor(cultivo),
8     sector = as.factor(sector)
9   )
10
11 # Verificar conversión
12 str(datos)
13
14 # Convertir fecha
15 datos <- datos %>%
16   mutate(
17     fecha_corte = as.Date(fecha_corte, format = "%Y%m%d")
18   )

```

4.4.3 Manejo de valores faltantes

Listing 9: Identificación y tratamiento de valores faltantes

```

1 # Contar valores faltantes por columna
2 colSums(is.na(datos))
3
4 # Visualizar patrón de valores faltantes
5 library(naniar)
6 vis_miss(datos)
7
8 # Decidir estrategia según el caso:
9
10 # 1. Eliminar filas con NA (si son pocas)
11 datos_completos <- na.omit(datos)
12
13 # 2. Reemplazar NA con 0 en columnas de siembra
14 # (asumiendo que NA = no hay siembra)
15 meses_siembra <- grep("^\$siembras_\$", names(datos), value = TRUE)
16 datos <- datos %>%
17   mutate(across(all_of(meses_siembra), ~replace_na(., 0)))
18
19 # 3. Imputación más sofisticada (si es necesario)
20 # library(mice)
21 # datos_imputados <- mice(datos, m = 5, method = 'pmm')

```

4.5 Creación de variables derivadas

Crear nuevas variables puede facilitar el análisis y la interpretación.

Listing 10: Creación de variables derivadas

```

1 # 1. Total de área de siembra por observación
2 datos <- datos %>%
3   mutate(
4     area_total = rowSums(select(., starts_with("siembras_"))),
5     na.rm = TRUE)
6 )
7
8 # 2. Indicador binario: hubo siembra?
9 datos <- datos %>%
10  mutate(
11    hubo_siembra = ifelse(area_total > 0, 1, 0),
12    hubo_siembra = as.factor(hubo_siembra)
13  )
14
15 # 3. Categorizar cultivos por tipo
16 datos <- datos %>%
17  mutate(
18    tipo_cultivo = case_when(
19      grepl("PAPA|OLLUCO|OCA|MASHUA", cultivo) ~ "Tubérculos",
20      grepl("MAIZ|TRIGO|CEBADA|QUINUA|AVENA", cultivo) ~ "Cereales",
21      grepl("HABA|ARVEJA|TARWI", cultivo) ~ "Leguminosas",
22      TRUE ~ "Otros"
23    ),
24    tipo_cultivo = as.factor(tipo_cultivo)
25  )
26
27 # 4. Mes con mayor siembra
28 datos <- datos %>%
29  rowwise() %>%
30  mutate(
31    mes_max_siembra = names(.)[which.max(c_across(starts_with("siembras_")))
32      ])
33 ) %>%
34 ungroup()
35
36 # 5. Número de meses con siembra
37 datos <- datos %>%
38  mutate(
39    n_meses_siembra = rowSums(select(., starts_with("siembras_")) > 0)
40  )

```

4.6 Transformación de datos

Para ciertos análisis, necesitamos transformar la estructura de los datos.

4.6.1 De formato ancho a largo

El formato actual es "ancho" (una columna por mes). Para algunos análisis es útil el formato "largo":

Listing 11: Transformación a formato largo

```

1 # Convertir de ancho a largo
2 datos_largo <- datos %>%
3   pivot_longer(
4     cols = starts_with("siembras_"),
5     names_to = "mes",
6     values_to = "area_ha",
7     names_prefix = "siembras_"
8   )
9
10 # Ver resultado
11 head(datos_largo)
12
13 # Ahora cada combinación cultivo-distrito-mes es una fila
14 # Esto facilita análisis por mes o gráficos temporales

```

4.7 Verificación final

Antes de proceder al análisis, verificamos que los datos estén listos:

Listing 12: Verificación final de datos preparados

```

1 # 1. Dimensiones correctas
2 dim(datos)
3
4 # 2. No hay valores faltantes en variables clave
5 colSums(is.na(datos))
6
7 # 3. Tipos de datos correctos
8 str(datos)
9
10 # 4. Rangos razonables en variables numéricas
11 summary(datos$area_total)
12
13 # 5. Niveles apropiados en factores
14 levels(datos$cultivo)
15 levels(datos$provincia)
16
17 # 6. Guardar datos limpios para análisis
18 saveRDS(datos, file = "datos_limpios.rds")
19 write.csv(datos, file = "datos_limpios.csv", row.names = FALSE)

```

4.8 Resumen de la preparación

El proceso de preparación de datos incluye:

1. **Carga:** Leer el archivo con encoding correcto
2. **Exploración:** Entender estructura y contenido
3. **Limpieza:** Estandarizar nombres, tipos, valores faltantes
4. **Transformación:** Crear variables derivadas útiles
5. **Verificación:** Confirmar que los datos están listos

Con los datos preparados, podemos proceder al análisis estadístico con confianza.

5 Caso 1: Regresión Lineal

La regresión lineal es uno de los métodos estadísticos más utilizados y fundamentales. Permite modelar la relación entre una variable respuesta continua y una o más variables predictoras.

5.1 Objetivo del análisis

Queremos responder: **¿Qué factores explican el área total de siembra de papa en la región Cusco?**

5.2 Pregunta estadística

Específicamente, investigaremos:

- ¿El área de siembra varía significativamente entre provincias?
- ¿Existe relación entre el número de meses de siembra y el área total?
- ¿Ciertas variedades de papa (color vs. nativa) se siembran en mayor extensión?

5.3 Preparación de datos específica

Para este análisis, filtraremos solo las observaciones de papa:

Listing 13: Preparación de datos para regresión lineal

```

1 # Filtrar solo cultivos de papa
2 datos_papa <- datos %>%
3   filter(grepl("PAPA", cultivo)) %>%
4   filter(area_total > 0) %>% # Solo donde hay siembra
5   droplevels() # Eliminar niveles no usados de factores
6
7 # Crear variable de tipo de papa
8 datos_papa <- datos_papa %>%
9   mutate(
10     tipo_papa = case_when(
11       grepl("COLOR", cultivo) ~ "Papa Color",

```

```

12     grepl("NATIVA", cultivo) ~ "Papa Nativa",
13     grepl("BLANCA", cultivo) ~ "Papa Blanca",
14     TRUE ~ "Otra Papa"
15   ),
16   tipo_papa = as.factor(tipo_papa)
17 )
18
19 # Crear variable log de área (para normalizar distribución)
20 datos_papa <- datos_papa %>%
21   mutate(
22     log_area = log(area_total + 1) # +1 para evitar log(0)
23   )
24
25 # Resumen descriptivo
26 summary(datos_papa$area_total)
27 table(datos_papa$provincia)
28 table(datos_papa$tipo_papa)

```

5.4 Modelo estadístico propuesto

Ajustaremos un modelo de regresión lineal múltiple:

$$\log_{-}\text{area}_i = \beta_0 + \beta_1 \text{provincia}_i + \beta_2 \text{tipo_papa}_i + \beta_3 \text{n_meses}_i + \epsilon_i \quad (1)$$

Donde:

- $\log_{-}\text{area}_i$: Logaritmo del área total de siembra (variable respuesta)
- provincia_i : Provincia (variable categórica)
- tipo_papa_i : Tipo de papa (variable categórica)
- n_meses_i : Número de meses con siembra (variable continua)
- ϵ_i : Error aleatorio, $\epsilon_i \sim N(0, \sigma^2)$

Nota: Usamos transformación logarítmica porque el área de siembra suele tener distribución asimétrica (muchos valores pequeños, pocos valores grandes).

5.5 Ajuste del modelo en R

Listing 14: Ajuste de regresión lineal múltiple

```

1 # Ajustar modelo
2 modelo_lm <- lm(
3   log_area ~ provincia + tipo_papa + n_meses_siembra,
4   data = datos_papa
5 )
6
7 # Ver resumen básico (base R)
8 summary(modelo_lm)

```

5.6 Análisis con easystats

5.6.1 Extracción de parámetros

Listing 15: Parámetros del modelo con parameters()

```

1 # Tabla de parámetros con easystats
2 params <- model_parameters(modelo_lm)
3 print(params)
4
5 # Parámetros estandarizados (coeficientes beta)
6 params_std <- standardize_parameters(modelo_lm)
7 print(params_std)
8
9 # Interpretación: Los coeficientes estandarizados permiten comparar
10 # la magnitud relativa del efecto de cada predictor

```

Interpretación de coeficientes:

- **Intercepto:** Valor esperado de log_area cuando todas las demás variables están en su categoría de referencia
- **provincia[nivel]:** Diferencia en log_area respecto a la provincia de referencia
- **tipo_papa[nivel]:** Diferencia en log_area respecto al tipo de papa de referencia
- **n_meses_siembra:** Por cada mes adicional de siembra, log_area aumenta en β_3 unidades

5.6.2 Evaluación del desempeño del modelo

Listing 16: Métricas de bondad de ajuste con performance()

```

1 # Métricas de desempeño
2 perf <- model_performance(modelo_lm)
3 print(perf)
4
5 # Métricas importantes:
6 # - R2: Proporción de varianza explicada (0-1)
7 # - R2_adjusted: R2 ajustado por número de predictores
8 # - RMSE: Raíz del error cuadrático medio
9 # - Sigma: Desviación estándar residual
10 # - AIC/BIC: Criterios de información (menor es mejor)

```

Interpretación de métricas:

Table 3: Interpretación de métricas de desempeño

Métrica	Rango	Interpretación
R^2	0–1	% de varianza explicada
R^2 ajustado	0–1	R^2 penalizado por predictores
RMSE	$0 - \infty$	Error promedio (mismas unidades que Y)
MAE	$0 - \infty$	Error absoluto promedio
AIC	$-\infty - \infty$	Comparar modelos (menor mejor)
BIC	$-\infty - \infty$	AIC con penalización mayor

5.6.3 Verificación de supuestos

La regresión lineal asume:

1. **Linealidad:** Relación lineal entre X e Y
2. **Normalidad:** Residuos siguen distribución normal
3. **Homoscedasticidad:** Varianza constante de residuos
4. **Independencia:** Observaciones independientes

Listing 17: Diagnóstico visual completo con check_model()

```

1 # Gráficos diagnósticos automáticos
2 check_model(modelo_lm)
3
4 # Este comando genera múltiples gráficos:
5 # 1. Posterior Predictive Check: Ajuste general
6 # 2. Linearity: Linealidad de la relación
7 # 3. Homogeneity of Variance: Homoscedasticidad
8 # 4. Influential Observations: Valores influyentes
9 # 5. Normality of Residuals: Normalidad de residuos
10 # 6. Multicollinearity: Correlación entre predictores

```

Listing 18: Pruebas específicas de supuestos

```

1 # Prueba de normalidad de residuos
2 check_normality(modelo_lm)
3
4 # Prueba de homogeneidad de varianza
5 check_heteroscedasticity(modelo_lm)
6
7 # Detección de multicolinealidad
8 check_collinearity(modelo_lm)
9
10 # Detección de valores influyentes
11 check_outliers(modelo_lm)

```

Interpretación de diagnósticos:

- Si normalidad se viola: Considerar transformación de Y o modelo robusto
- Si homoscedasticidad se viola: Usar errores estándar robustos
- Si multicolinealidad es alta ($VIF > 5$): Eliminar variables redundantes
- Si hay outliers influyentes: Investigar y potencialmente excluir

5.7 Cálculo de tamaños de efecto

Los valores p indican si un efecto es estadísticamente significativo, pero no su magnitud práctica. Los tamaños de efecto cuantifican la importancia práctica.

Listing 19: Tamaños de efecto con effectsize()

```

1 # Tamaño de efecto para el modelo completo
2 effectsize(modelo_lm)

3

4 #  $f^2$  de Cohen para cada predictor
5 f2 <- cohens_f_squared(modelo_lm)
6 print(f2)

7

8 # Interpretación de  $f^2$ :
9 # Pequeño: 0.02
10 # Mediano: 0.15
11 # Grande: 0.35
12
13 # Interpretar automáticamente
14 interpret_f_squared(f2)
```

5.8 Visualización de resultados

5.8.1 Gráfico de coeficientes

Listing 20: Visualización de coeficientes estimados

```

1 library(see)
2
3 # Gráfico de coeficientes con IC al 95%
4 plot(params)

5

6 # Alternativa: solo coeficientes significativos
7 plot(params, show_intercept = FALSE)

8

9 # Los intervalos que no cruzan cero son estadísticamente
10 # significativos al nivel alfa = 0.05
```

5.8.2 Gráficos de efectos predichos

Listing 21: Predicciones y efectos marginales

```

1 library(modelbased)
2
3 # Estimar medias por provincia
4 pred_provincia <- estimate_means(modelo_lm, at = "provincia")
5 plot(pred_provincia)
6
7 # Estimar medias por tipo de papa
8 pred_tipo <- estimate_means(modelo_lm, at = "tipo_papa")
9 plot(pred_tipo)
10
11 # Efecto del número de meses (manteniendo otras variables constantes)
12 pred_meses <- estimate_relation(modelo_lm, at = "n_meses_siembra")
13 plot(pred_meses)

```

5.9 Reporte automático

Listing 22: Generación de reporte con report()

```

1 # Reporte completo
2 report_modelo <- report(modelo_lm)
3 print(report_modelo)
4
5 # Solo texto narrativo
6 report_text(modelo_lm)
7
8 # Solo tabla de resultados
9 report_table(modelo_lm)
10
11 # Reporte de desempeño
12 report_performance(modelo_lm)

```

El reporte automático genera texto como:

“We fitted a linear model (estimated using OLS) to predict log_area with provincia, tipo_papa and n_meses_siembra (formula: log_area ~ provincia + tipo_papa + n_meses_siembra). The model explains a statistically significant and moderate proportion of variance ($R^2 = 0.XX$, $F(XX, XXX) = XX.XX$, $p < .001$, adj. $R^2 = 0.XX$). The model’s intercept, corresponding to provincia = [referencia], is at X.XX (95% CI [X.XX, X.XX], $t(XXX) = XX.XX$, $p < .001$). Within this model: ...”

5.10 Interpretación de resultados

Para interpretar correctamente los resultados:

1. **Coeficientes:** Cuantifican el cambio en log_area por unidad de cambio en el predictor. Para interpretar en escala original: area = $e^{\text{log-area}}$
2. **Significancia estadística:** Un valor $p < 0.05$ indica que el efecto es improbable bajo la hipótesis nula (no hay efecto).
3. **Intervalos de confianza:** Si el IC al 95% no incluye cero, el efecto es significativo al $\alpha = 0.05$.
4. R^2 : Indica qué porcentaje de la variabilidad en área de siembra es explicado por el modelo. Un $R^2 = 0.40$ significa que el modelo explica el 40% de la varianza.
5. **Tamaños de efecto:** Complementan los valores p indicando la magnitud práctica del efecto.

Ejemplo de interpretación integrada:

El modelo de regresión lineal explica el 45% de la varianza en el área de siembra de papa ($R^2 = 0.45$, $p \downarrow 0.001$). La provincia tiene un efecto significativo ($F = 12.5$, $p \downarrow 0.001$, $f^2 = 0.18$, efecto mediano). Específicamente, la provincia de Canchis siembra en promedio 35% más área que la provincia de referencia (Acomayo), manteniendo constantes el tipo de papa y meses de siembra ($b = 0.30$, IC 95% [0.15, 0.45], $p \downarrow 0.001$). El número de meses de siembra también predice positivamente el área ($b = 0.12$ por mes adicional, IC 95% [0.08, 0.16], $p \downarrow 0.001$).

5.11 Comparación de modelos

A menudo queremos comparar modelos con diferente complejidad:

Listing 23: Comparación de modelos anidados

```

1 # Modelo simple (solo provincia)
2 modelo_1 <- lm(log_area ~ provincia, data = datos_papa)
3
4 # Modelo con tipo de papa
5 modelo_2 <- lm(log_area ~ provincia + tipo_papa, data = datos_papa)
6
7 # Modelo completo (añade meses)
8 modelo_3 <- lm(log_area ~ provincia + tipo_papa + n_meses_siembra,
9                 data = datos_papa)
10
11 # Comparar desempeño
12 comparacion <- compare_performance(modelo_1, modelo_2, modelo_3)
13 print(comparacion)
14
15 # Visualizar comparación
16 plot(comparacion)
17
18 # Test de razón de verosimilitud (likelihood ratio test)
19 test_likelihoodratio(modelo_1, modelo_2, modelo_3)

```

Criterios de selección:

- **AIC/BIC más bajo:** Mejor balance entre ajuste y complejidad
- **R^2 ajustado más alto:** Mayor varianza explicada ajustada por predictores
- **Test de razón de verosimilitud significativo:** El modelo más complejo mejora significativamente el ajuste

5.12 Predicción con el modelo

Una vez ajustado y validado el modelo, podemos hacer predicciones:

Listing 24: Predicciones con el modelo ajustado

```

1 # Crear datos nuevos para predicción
2 nuevos_datos <- expand.grid(
3   provincia = c("CUSCO", "CANCHIS", "QUISPICANCHI"),
4   tipo_papa = c("Papa Color", "Papa Nativa"),
5   n_meses_siembra = c(1, 2, 3)
6 )
7
8 # Predicciones puntuales
9 predicciones <- predict(modelo_lm, newdata = nuevos_datos)
10
11 # Predicciones con intervalos de confianza
12 predicciones_ic <- predict(modelo_lm, newdata = nuevos_datos,
13                           interval = "confidence", level = 0.95)
14
15 # Predicciones con intervalos de predicción (más amplios)
16 predicciones_pred <- predict(modelo_lm, newdata = nuevos_datos,
17                               interval = "prediction", level = 0.95)
18
19 # Combinar con datos originales
20 resultado <- cbind(nuevos_datos, predicciones_ic)
21
22 # Volver a escala original (des-logaritmizar)
23 resultado$area_pred <- exp(resultado$fit)
24 resultado$area_lwr <- exp(resultado$lwr)
25 resultado$area_upr <- exp(resultado$upr)
26
27 print(resultado)

```

5.13 Resumen del flujo de regresión lineal

El análisis completo de regresión lineal con easystats sigue estos pasos:

1. Preparar datos (filtrado, transformaciones)
2. Ajustar modelo con `lm()`

3. Extraer parámetros con `parameters()`
4. Evaluar desempeño con `performance()`
5. Verificar supuestos con `check_model()`
6. Calcular tamaños de efecto con `effectsize()`
7. Visualizar con `plot()` y `see`
8. Generar reporte con `report()`
9. Interpretar resultados en contexto
10. Comparar modelos alternativos si es relevante
11. Hacer predicciones si es necesario

Este flujo es sistemático, reproducible y aplicable a cualquier análisis de regresión lineal.

6 Caso 2: Regresión Logística

La regresión logística se utiliza cuando la variable respuesta es binaria (0/1, sí/no, presencia/ausencia). Es fundamental en clasificación y predicción de probabilidades.

6.1 Objetivo del análisis

Queremos responder: **¿Qué factores predicen la presencia o ausencia de siembra de quinua en un sector?**

6.2 Pregunta estadística

Específicamente:

- ¿La probabilidad de sembrar quinua varía entre provincias?
- ¿Existe relación entre la altitud (aproximada por distrito) y la siembra de quinua?
- ¿Podemos predecir qué sectores sembrarán quinua basándonos en características geográficas?

6.3 Preparación de datos específica

Listing 25: Preparación de datos para regresión logística

```
1 # Crear variable binaria: se sembró quinua?
2 datos_quinua <- datos %>%
3   mutate(
4     siembra_quinua = ifelse(cultivo == "QUINUA" & area_total > 0, 1, 0)
```

```

5 ) %>%
6 # Agregar por sector (un sector puede tener múltiples cultivos)
7 group_by(departamento, provincia, distrito, sector) %>%
8 summarise(
9   siembra_quinua = max(siembra_quinua), # 1 si hubo quinua
10  n_cultivos = n(),                      # Número de cultivos
11  area_total_sector = sum(area_total), # Área total del sector
12  .groups = "drop"
13 ) %>%
14 mutate(
15   siembra_quinua = as.factor(siembra_quinua)
16 )
17
18 # Ver distribución de la variable respuesta
19 table(datos_quinua$siembra_quinua)
20 prop.table(table(datos_quinua$siembra_quinua))
21
22 # Idealmente, la distribución no debería ser extremadamente
23 # desbalanceada (ej. 99% en una categoría)

```

6.4 Modelo estadístico propuesto

Ajustaremos un modelo de regresión logística:

$$\text{logit}(p_i) = \log \left(\frac{p_i}{1 - p_i} \right) = \beta_0 + \beta_1 \text{provincia}_i + \beta_2 \text{n_cultivos}_i + \beta_3 \log(\text{area}_i) \quad (2)$$

Donde:

- p_i : Probabilidad de sembrar quinua en el sector i
- $\text{logit}(p_i)$: Transformación logit (log-odds)
- Variables predictoras: provincia, número de cultivos, log(área)

La relación entre probabilidad y log-odds es:

$$p_i = \frac{e^{\beta_0 + \beta_1 x_1 + \dots}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots)}} \quad (3)$$

6.5 Ajuste del modelo en R

Listing 26: Ajuste de regresión logística

```

1 # Ajustar modelo logístico
2 modelo_logit <- glm(
3   siembra_quinua ~ provincia + n_cultivos + log(area_total_sector + 1),
4   data = datos_quinua,
5   family = binomial(link = "logit")

```

```

6  )
7
8 # Resumen básico
9 summary(modelo_logit)

```

Parámetros importantes de `glm()`:

- `family = binomial`: Especifica distribución binomial
- `link = "logit"`: Función de enlace logit (por defecto para binomial)
- Alternativas: `link = "probit"` (distribución normal acumulativa)

6.6 Análisis con easystats

6.6.1 Extracción de parámetros

Listing 27: Parámetros del modelo logístico

```

1 # Parámetros en escala log-odds
2 params_logit <- model_parameters(modelo_logit)
3 print(params_logit)
4
5 # Parámetros como Odds Ratios (OR)
6 params_or <- model_parameters(modelo_logit, exponentiate = TRUE)
7 print(params_or)

```

Interpretación de Odds Ratios:

- $OR = 1$: No hay efecto
- $OR > 1$: Aumenta las odds de $Y = 1$ (efecto positivo)
- $OR < 1$: Disminuye las odds de $Y = 1$ (efecto negativo)
- $OR = 2$: Duplica las odds
- $OR = 0.5$: Reduce las odds a la mitad

Ejemplo: Si el OR para n_cultivos es 1.5, significa que por cada cultivo adicional, las odds de sembrar quinua aumentan en 50%.

6.6.2 Métricas de desempeño

Listing 28: Evaluación del modelo logístico

```

1 # Métricas de bondad de ajuste
2 perf_logit <- model_performance(modelo_logit)
3 print(perf_logit)
4
5 # Métricas importantes:
6 # - AIC/BIC: Criterios de información

```

```

7 # - Tjur's R2: Pseudo-R2 apropiado para logística
8 # - Log-likelihood: Verosimilitud logarítmica
9 # - Score_log: Score logarítmico (menor es mejor)

```

Pseudo-R² en regresión logística:

A diferencia de la regresión lineal, no existe un R² único. Existen múltiples pseudo-R²:

Table 4: Tipos de Pseudo-R² para regresión logística

Tipo	Rango	Interpretación
Tjur's R ²	0–1	Diferencia de prob. entre Y=0 y Y=1
McFadden's R ²	0–1	Basado en log-verosimilitud
Nagelkerke's R ²	0–1	Versión ajustada de Cox-Snell

6.6.3 Matriz de confusión y métricas de clasificación

Listing 29: Evaluación de capacidad predictiva

```

1 # Obtener probabilidades predichas
2 datos_quinua$prob_pred <- predict(modelo_logit, type = "response")
3
4 # Clasificar con umbral 0.5
5 datos_quinua$pred_clase <- ifelse(datos_quinua$prob_pred > 0.5, 1, 0)
6 datos_quinua$pred_clase <- as.factor(datos_quinua$pred_clase)
7
8 # Matriz de confusión
9 library(caret)
10 confusionMatrix(datos_quinua$pred_clase, datos_quinua$siembra_quinua)
11
12 # Alternativamente con performance
13 performance_accuracy(modelo_logit)
14
15 # Curva ROC y AUC
16 library(pROC)
17 roc_obj <- roc(datos_quinua$siembra_quinua, datos_quinua$prob_pred)
18 auc(roc_obj)
19 plot(roc_obj, main = "Curva ROC - Predicción Siembra Quinua")

```

Métricas de clasificación:

- **Accuracy (Precisión)**: % de predicciones correctas
- **Sensitivity (Sensibilidad/Recall)**: % de 1's correctamente identificados
- **Specificity (Especificidad)**: % de 0's correctamente identificados
- **Precision (Precisión positiva)**: % de predicciones "1" que son correctas
- **AUC (Area Under Curve)**: Capacidad discriminatoria general (0.5–1)

6.6.4 Verificación de supuestos

Listing 30: Diagnósticos del modelo logístico

```

1 # Diagnóstico visual
2 check_model(modelo_logit)
3
4 # Multicolinealidad
5 check_collinearity(modelo_logit)
6
7 # Valores influyentes
8 check_outliers(modelo_logit)
9
10 # Overdispersion (para modelos binomiales)
11 check_overdispersion(modelo_logit)

```

Nota: La regresión logística no asume normalidad de residuos (el outcome es binario), pero sí asume:

1. Independencia de observaciones
2. Linealidad entre log-odds y predictores continuos
3. Ausencia de multicolinealidad perfecta

6.7 Interpretación de Odds Ratios

Los Odds Ratios son la medida de efecto natural en regresión logística:

Listing 31: Cálculo e interpretación de Odds Ratios

```

1 # Extraer Odds Ratios con IC
2 or_tabla <- model_parameters(modelo_logit, exponentiate = TRUE)
3 print(or_tabla)
4
5 # Visualizar Odds Ratios
6 plot(or_tabla)
7
8 # Interpretar automáticamente
9 interpret_oddsratio(or_tabla$Coefficient)

```

Guía de interpretación:

Table 5: Interpretación práctica de Odds Ratios

OR	Cambio en odds	Interpretación
0.5	-50%	Reduce odds a la mitad
0.8	-20%	Reducción moderada
1.0	0%	Sin efecto
1.2	+20%	Aumento moderado
2.0	+100%	Duplica las odds
3.0	+200%	Triuplica las odds

Ejemplo de interpretación completa:

El modelo logístico muestra que el número de cultivos en un sector predice significativamente la siembra de quinua (OR = 1.35, IC 95% [1.18, 1.54], p < 0.001). Específicamente, por cada cultivo adicional, las odds de sembrar quinua aumentan en 35%, manteniendo constantes la provincia y el área. La provincia de Chumbivilcas tiene 2.5 veces más odds de sembrar quinua que la provincia de referencia (OR = 2.50, IC 95% [1.80, 3.47], p < 0.001).

6.8 Predicción de probabilidades

Listing 32: Estimación de probabilidades con modelbased

```

1 # Probabilidad media por provincia
2 prob_provincia <- estimate_expectation(
3   modelo_logit,
4   at = "provincia"
5 )
6 plot(prob_provincia)
7
8 # Probabilidades según número de cultivos
9 prob_cultivos <- estimate_expectation(
10   modelo_logit,
11   at = "n_cultivos = c(1, 3, 5, 7, 9)"
12 )
13 plot(prob_cultivos)
14
15 # Superficie de respuesta (2 predictores)
16 prob_superficie <- estimate_expectation(
17   modelo_logit,
18   at = c("n_cultivos", "area_total_sector")
19 )
20 plot(prob_superficie)

```

6.9 Efectos marginales

Los efectos marginales expresan el cambio en probabilidad (no en odds):

Listing 33: Cálculo de efectos marginales

```

1 library(marginaleffects)
2
3 # Efecto marginal promedio (AME)
4 ame <- marginaleffects(modelo_logit)
5 summary(ame)
6
7 # Efecto de un cambio de 1 unidad en cada variable
8 # sobre la probabilidad predicha
9

```

```

10 # Efectos marginales en valores específicos
11 mfx_cultivos <- marginaleffects(
12   modelo_logit,
13   variables = "n_cultivos",
14   newdata = datagrid(n_cultivos = 1:10)
15 )
16 plot(mfx_cultivos)

```

6.10 Visualización de resultados

Listing 34: Visualizaciones para regresión logística

```

1 # 1. Gráfico de Odds Ratios
2 plot(model_parameters(modelo_logit, exponentiate = TRUE))
3
4 # 2. Probabilidades predichas por grupo
5 library(ggplot2)
6 ggplot(prob_provincia, aes(x = provincia, y = Predicted)) +
7   geom_point(size = 3) +
8   geom_errorbar(aes(ymin = CI_low, ymax = CI_high), width = 0.2) +
9   coord_flip() +
10  labs(
11    title = "Probabilidad de Siembra de Quinua por Provincia",
12    x = "Provincia",
13    y = "Probabilidad Predicha"
14  ) +
15  theme_minimal()
16
17 # 3. Curva de probabilidad según predictor continuo
18 ggplot(prob_cultivos, aes(x = n_cultivos, y = Predicted)) +
19   geom_line(size = 1) +
20   geom_ribbon(aes(ymin = CI_low, ymax = CI_high), alpha = 0.2) +
21   labs(
22    title = "Efecto del Número de Cultivos en Prob. de Quinua",
23    x = "Número de Cultivos",
24    y = "Probabilidad de Siembra de Quinua"
25  ) +
26  theme_minimal()
27
28 # 4. Curva ROC
29 plot(roc_obj,
30   main = "Curva ROC",
31   print.auc = TRUE,
32   col = "blue",
33   lwd = 2)

```

6.11 Reporte automático

Listing 35: Reporte de regresión logística

```

1 # Reporte completo
2 report(modelo_logit)
3
4 # Reporte de desempeño
5 report_performance(modelo_logit)
6
7 # Reporte de efectos (Odds Ratios)
8 report_parameters(modelo_logit)

```

6.12 Comparación de modelos logísticos

Listing 36: Comparación de modelos logísticos anidados

```

1 # Modelo 1: Solo provincia
2 m1 <- glm(siembra_quinua ~ provincia,
            data = datos_quinua, family = binomial)
3
4 # Modelo 2: + número de cultivos
5 m2 <- glm(siembra_quinua ~ provincia + n_cultivos,
            data = datos_quinua, family = binomial)
6
7 # Modelo 3: + área
8 m3 <- glm(siembra_quinua ~ provincia + n_cultivos +
            log(area_total_sector + 1),
            data = datos_quinua, family = binomial)
9
10 # Comparar
11 comp <- compare_performance(m1, m2, m3, rank = TRUE)
12 print(comp)
13 plot(comp)
14
15 # Test de razón de verosimilitud
16 test_likelihoodratio(m1, m2, m3)

```

6.13 Validación cruzada

Para evaluar la capacidad predictiva en datos no vistos:

Listing 37: Validación cruzada del modelo logístico

```

1 library(caret)
2
3 # Configurar validación cruzada de 10 pliegues
4 control <- trainControl(
5   method = "cv",
6   number = 10,
7   classProbs = TRUE,
8   summaryFunction = twoClassSummary

```

```
9 )
10
11 # Preparar datos (caret requiere factores con nombres)
12 datos_cv <- datos_quinua %>%
13   mutate(
14     siembra_quinua = factor(
15       siembra_quinua,
16       levels = c(0, 1),
17       labels = c("No", "Si")
18     )
19   )
20
21 # Entrenar con CV
22 modelo_cv <- train(
23   siembra_quinua ~ provincia + n_cultivos +
24     log(area_total_sector + 1),
25   data = datos_cv,
26   method = "glm",
27   family = binomial,
28   trControl = control,
29   metric = "ROC"
30 )
31
32 # Ver resultados
33 print(modelo_cv)
34 modelo_cv$results
```

6.14 Resumen del flujo de regresión logística

El análisis completo de regresión logística incluye:

1. Crear variable binaria (0/1)
2. Ajustar modelo con `glm(family = binomial)`
3. Extraer coeficientes en log-odds y OR con `parameters()`
4. Evaluar con pseudo-R² y métricas de clasificación
5. Verificar supuestos con `check_model()`
6. Interpretar Odds Ratios
7. Estimar probabilidades con `modelbased`
8. Visualizar OR, probabilidades y curva ROC
9. Generar reporte con `report()`
10. Validar con datos independientes o CV

7 Caso 3: Comparación de Grupos

La comparación de grupos es fundamental cuando queremos determinar si existen diferencias significativas entre dos o más grupos. Utilizaremos pruebas t y ANOVA, complementadas con tamaños de efecto.

7.1 Objetivo del análisis

Queremos responder: *¿Existen diferencias significativas en el área de siembra entre diferentes provincias y tipos de cultivos?*

7.2 Preguntas estadísticas

Investigaremos:

- ¿Difiere el área promedio de siembra entre dos provincias específicas? (prueba t)
- ¿Existen diferencias en área entre múltiples provincias? (ANOVA)
- ¿El tipo de cultivo (tubérculos vs cereales vs leguminosas) afecta el área sembrada?
- ¿Qué tan grandes son estas diferencias? (tamaños de efecto)

7.3 Preparación de datos

Listing 38: Preparación de datos para comparación de grupos

```

1 # Datos agregados por sector y cultivo
2 datos_grupos <- datos %>%
3   filter(area_total > 0) %>%
4   mutate(
5     # Clasificar cultivos por tipo
6     tipo_cultivo = case_when(
7       grepl("PAPA|OLLUCO|OCA|MASHUA", cultivo) ~ "Tubérculos",
8       grepl("MAIZ|TRIGO|CEBADA|QUINUA|AVENA|KIWICHA", cultivo) ~ "Cereales"
9       ,
10      grepl("HABA|ARVEJA|TARWI|FRIJOL", cultivo) ~ "Leguminosas",
11      TRUE ~ "Otros"
12    ),
13    tipo_cultivo = as.factor(tipo_cultivo),
14    # Log-transformación
15    log_area = log(area_total + 1)
16  )
17
18 # Resumen descriptivo por grupo
19 datos_grupos %>%
20   group_by(provincia) %>%
21   summarise(
22     n = n(),
23     media = mean(area_total),

```

```

23     sd = sd(area_total),
24     mediana = median(area_total)
25   )
26
27 datos_grupos %>%
28   group_by(tipo_cultivo) %>%
29   summarise(
30     n = n(),
31     media = mean(area_total),
32     sd = sd(area_total),
33     mediana = median(area_total)
34   )

```

7.4 Caso 3.1: Prueba t para dos grupos

7.4.1 Prueba t independiente

Comparemos el área de siembra entre dos provincias específicas: Cusco y Canchis.

Listing 39: Prueba t independiente con easystats

```

1 # Filtrar dos provincias
2 datos_2prov <- datos_grupos %>%
3   filter(provincia %in% c("CUSCO", "CANCHIS"))
4
5 # Prueba t tradicional (para referencia)
6 t.test(area_total ~ provincia, data = datos_2prov, var.equal = TRUE)
7
8 # Con easystats: información más completa
9 # Primero ajustamos un modelo lineal simple (equivalente a t-test)
10 modelo_t <- lm(log_area ~ provincia, data = datos_2prov)
11
12 # Parámetros del modelo
13 parameters(modelo_t)
14
15 # El coeficiente de provincia es la diferencia entre grupos
16
17 # Tamaño de efecto: d de Cohen
18 d_cohen <- cohens_d(log_area ~ provincia, data = datos_2prov)
19 print(d_cohen)
20 interpret_cohens_d(d_cohen$Cohens_d)
21
22 # Visualización
23 library(ggplot2)
24 ggplot(datos_2prov, aes(x = provincia, y = log_area, fill = provincia)) +
25   geom_boxplot() +
26   geom_jitter(alpha = 0.2, width = 0.2) +
27   labs(
28     title = "Comparación de Área de Siembra entre Cusco y Canchis",
29     x = "Provincia",

```

```

30     y = "Log(Área en hectáreas)"
31   ) +
32   theme_minimal() +
33   theme(legend.position = "none")

```

Interpretación de d de Cohen:

Table 6: Interpretación de d de Cohen

Valor —d—	Magnitud
< 0.2	Muy pequeño
0.2 – 0.5	Pequeño
0.5 – 0.8	Mediano
> 0.8	Grande

7.4.2 Prueba t pareada

Si tuviéramos mediciones repetidas (ej. mismo sector en dos años), usaríamos:

Listing 40: Ejemplo de prueba t pareada

```

1 # Ejemplo hipotético con datos pareados
2 # t.test(antes, despues, paired = TRUE)
3
4 # Con lmer para diseños más complejos:
5 # library(lme4)
6 # modelo_pareado <- lmer(area ~ tiempo + (1/sector), data = datos_long)

```

7.5 Caso 3.2: ANOVA de una vía

Cuando comparamos más de dos grupos, utilizamos ANOVA.

Listing 41: ANOVA de una vía con easystats

```

1 # Modelo ANOVA: área según tipo de cultivo
2 modelo_anova <- lm(log_area ~ tipo_cultivo, data = datos_grupos)
3
4 # Tabla ANOVA
5 anova_tabla <- anova(modelo_anova)
6 print(anova_tabla)
7
8 # Con easystats
9 params_anova <- model_parameters(modelo_anova)
10 print(params_anova)
11
12 # Pruebas post-hoc (comparaciones múltiples)
13 library(emmeans)
14 medias_est <- emmeans(modelo_anova, specs = "tipo_cultivo")
15 contrastes <- contrast(medias_est, method = "pairwise", adjust = "tukey")
16 print(contrastes)

```

```

17
18 # Con modelbased
19 medias_tipo <- estimate_means(modelo_anova, at = "tipo_cultivo")
20 print(medias_tipo)
21 plot(medias_tipo)
22
23 # Contrastos
24 contrasts <- estimate_contrasts(modelo_anova, contrast = "tipo_cultivo")
25 print(contrasts)
26 plot(contrasts)

```

7.5.1 Tamaño de efecto para ANOVA

Listing 42: Eta cuadrado y omega cuadrado

```

1 # Eta cuadrado ( )
2 eta_sq <- eta_squared(modelo_anova)
3 print(eta_sq)
4 interpret_eta_squared(eta_sq$Eta2)
5
6 # Omega cuadrado ( ) - menos sesgado
7 omega_sq <- omega_squared(modelo_anova)
8 print(omega_sq)
9
10 # Epsilon cuadrado ( ) - el menos sesgado
11 epsilon_sq <- epsilon_squared(modelo_anova)
12 print(epsilon_sq)

```

Interpretación de η^2 :

Table 7: Interpretación de eta cuadrado

η^2	Magnitud
< 0.01	Muy pequeño
0.01 – 0.06	Pequeño
0.06 – 0.14	Mediano
> 0.14	Grande

7.6 Caso 3.3: ANOVA de dos vías (factorial)

Cuando tenemos dos factores, podemos evaluar sus efectos principales e interacción.

Listing 43: ANOVA factorial con easystats

```

1 # Seleccionar provincias principales
2 datos_factorial <- datos_grupos %>%
3   filter(provincia %in% c("CUSCO", "CANCHIS", "QUISPICANCHI")) %>%
4   filter(tipo_cultivo %in% c("Tubérculos", "Cereales", "Leguminosas")) %>%
5   droplevels()

```

```

6
7 # Modelo ANOVA 2x3 (provincia x tipo_cultivo)
8 modelo_2way <- lm(log_area ~ provincia * tipo_cultivo,
9                      data = datos_factorial)
10
11 # Tabla ANOVA tipo II (para diseños desbalanceados)
12 library(car)
13 Anova(modelo_2way, type = "II")
14
15 # Parámetros del modelo
16 parameters(modelo_2way)
17
18 # Medias estimadas
19 medias_2way <- estimate_means(
20   modelo_2way,
21   at = c("provincia", "tipo_cultivo")
22 )
23 plot(medias_2way)
24
25 # Efectos simples (efecto de provincia para cada tipo de cultivo)
26 efectos_simples <- estimate_slopes(
27   modelo_2way,
28   trend = "provincia",
29   at = "tipo_cultivo"
30 )
31 print(efectos_simples)
32 plot(efectos_simples)
33
34 # Tamaño de efecto para cada factor
35 eta_sq_2way <- eta_squared(modelo_2way, partial = TRUE)
36 print(eta_sq_2way)

```

7.6.1 Interpretación de interacciones

Una interacción significativa indica que el efecto de un factor depende del nivel del otro factor.

Listing 44: Visualización de interacciones

```

1 # Gráfico de interacción
2 library(ggplot2)
3 ggplot(medias_2way, aes(x = provincia, y = Mean,
4                           color = tipo_cultivo, group = tipo_cultivo)) +
5   geom_line(size = 1) +
6   geom_point(size = 3) +
7   geom_errorbar(aes(ymin = CI_low, ymax = CI_high), width = 0.2) +
8   labs(
9     title = "Interacción entre Provincia y Tipo de Cultivo",
10    subtitle = "Efecto en Log(Área de Siembra)",
11    x = "Provincia",

```

```

12     y = "Media Estimada (Log-escala)",
13     color = "Tipo de Cultivo"
14 ) +
15 theme_minimal()
16
17 # Las líneas paralelas indican ausencia de interacción
18 # Las líneas que se cruzan indican interacción

```

7.7 Verificación de supuestos de ANOVA

Listing 45: Diagnósticos para ANOVA

```

1 # Diagnósticos visuales
2 check_model(modelo_anova)
3
4 # Supuestos específicos:
5
6 # 1. Normalidad de residuos
7 check_normality(modelo_anova)
8 plot(check_normality(modelo_anova))
9
10 # 2. Homogeneidad de varianzas (homocedasticidad)
11 check_heteroscedasticity(modelo_anova)
12
13 # Test de Levene para homogeneidad de varianzas
14 library(car)
15 leveneTest(log_area ~ tipo_cultivo, data = datos_grupos)
16
17 # Si se viola homogeneidad: considerar Welch ANOVA
18 oneway.test(log_area ~ tipo_cultivo, data = datos_grupos,
19             var.equal = FALSE)

```

Si los supuestos se violan:

- **No normalidad:** Transformar datos (log, sqrt) o usar pruebas no paramétricas (Kruskal-Wallis)
- **No homogeneidad:** Usar correcciones (Welch ANOVA) o modelos robustos

7.8 Alternativas no paramétricas

Si los supuestos de ANOVA se violan severamente:

Listing 46: Pruebas no paramétricas

```

1 # Kruskal-Wallis (alternativa no paramétrica a ANOVA)
2 kruskal.test(area_total ~ tipo_cultivo, data = datos_grupos)
3
4 # Pruebas post-hoc para Kruskal-Wallis
5 library(FSA)

```

```

6 dunnTest(area_total ~ tipo_cultivo, data = datos_grupos,
7           method = "bonferroni")
8
9 # Mann-Whitney U (alternativa no paramétrica a prueba t)
10 wilcox.test(area_total ~ provincia, data = datos_2prov)
11
12 # Tamaño de efecto para Mann-Whitney
13 rank_biserial(area_total ~ provincia, data = datos_2prov)

```

7.9 Comparaciones múltiples y corrección por múltiples pruebas

Cuando realizamos múltiples comparaciones, aumenta el riesgo de error tipo I (falsos positivos).

Listing 47: Corrección para comparaciones múltiples

```

1 # Comparaciones pareadas con ajuste
2 library(emmeans)
3 emm <- emmeans(modelo_anova, specs = "tipo_cultivo")
4
5 # Diferentes métodos de ajuste:
6 # 1. Tukey HSD (para todas las comparaciones pareadas)
7 pairs(emm, adjust = "tukey")
8
9 # 2. Bonferroni (más conservador)
10 pairs(emm, adjust = "bonferroni")
11
12 # 3. Holm (menos conservador que Bonferroni)
13 pairs(emm, adjust = "holm")
14
15 # 4. FDR - False Discovery Rate (control de tasa de descubrimientos falsos)
16 pairs(emm, adjust = "fdr")
17
18 # Con modelbased
19 contrasts_ajustados <- estimate_contrasts(
20   modelo_anova,
21   contrast = "tipo_cultivo",
22   adjust = "tukey"
23 )
24 print(contrasts_ajustados)

```

7.10 Visualizaciones avanzadas

Listing 48: Visualizaciones para comparación de grupos

```

1 # 1. Boxplot con significancia
2 library(ggpubr)
3 ggboxplot(datos_grupos, x = "tipo_cultivo", y = "log_area",
4            fill = "tipo_cultivo", palette = "jco") +

```

```

5   stat_compare_means(method = "anova", label.y = max(datos_grupos$log_area)
6     + 0.5) +
7   stat_compare_means(comparisons = list(
8     c("Tubérculos", "Cereales"),
9     c("Cereales", "Leguminosas"),
10    c("Tubérculos", "Leguminosas")
11  ), method = "t.test")
12
13 # 2. Violin plot con puntos
14 ggplot(datos_grupos, aes(x = tipo_cultivo, y = log_area,
15   fill = tipo_cultivo)) +
16   geom_violin(alpha = 0.7) +
17   geom_boxplot(width = 0.2, fill = "white", outlier.shape = NA) +
18   geom_jitter(alpha = 0.1, width = 0.2) +
19   labs(
20     title = "Distribución de Área de Siembra por Tipo de Cultivo",
21     x = "Tipo de Cultivo",
22     y = "Log(Área en hectáreas)"
23   ) +
24   theme_minimal() +
25   theme(legend.position = "none")
26
27 # 3. Gráfico de medias con IC
28 plot(estimate_means(modelo_anova, at = "tipo_cultivo"))
29
30 # 4. Gráfico de contrastes
31 plot(estimate_contrasts(modelo_anova, contrast = "tipo_cultivo"))

```

7.11 Reporte automático

Listing 49: Reporte de comparación de grupos

```

1 # Reporte de ANOVA
2 report(modelo_anova)
3
4 # Reporte de prueba t
5 report(modelo_t)
6
7 # Reporte de tamaños de efecto
8 report(cohens_d(log_area ~ provincia, data = datos_2prov))
9 report(eta_squared(modelo_anova))

```

7.12 Ejemplo de interpretación integrada

Se realizó un ANOVA de una vía para comparar el área de siembra entre tipos de cultivo (Tubérculos, Cereales, Leguminosas). El análisis reveló un efecto significativo del tipo de cultivo sobre el área sembrada, $F(2, 1247) = 45.3$, $p < 0.001$, $\eta^2 = 0.068$ (efecto mediano).

Las comparaciones post-hoc (Tukey HSD) mostraron que los Tubérculos se siembran en mayor extensión promedio ($M = 2.8$ ha, $SD = 1.5$) que los Cereales ($M = 2.1$ ha, $SD = 1.2$, $p < 0.001$, $d = 0.51$) y las Leguminosas ($M = 1.6$ ha, $SD = 0.9$, $p < 0.001$, $d = 0.95$). Los Cereales también difieren significativamente de las Leguminosas ($p < 0.01$, $d = 0.43$). Los tamaños de efecto indican diferencias de magnitud práctica moderada.

7.13 Resumen del flujo de comparación de grupos

Para comparar grupos con easystats:

1. Explorar datos y verificar distribución de grupos
2. Dos grupos: usar `lm(y ~ grupo)` o `t.test()`
3. Múltiples grupos: usar `lm(y ~ grupo)` y `anova()`
4. Calcular tamaños de efecto: `cohens_d()`, `eta_squared()`
5. Verificar supuestos con `check_model()`
6. Comparaciones post-hoc con `emmeans` o `estimate_contrasts()`
7. Visualizar con boxplots, violin plots, gráficos de medias
8. Reportar con `report()`
9. Interpretar significancia estadística Y práctica (tamaño de efecto)

8 Caso 4: Modelos con Interacciones

Las interacciones ocurren cuando el efecto de una variable sobre la respuesta depende del nivel de otra variable. Son comunes en ciencias agrícolas, donde factores ambientales y de manejo interactúan.

8.1 Objetivo del análisis

Queremos responder: **¿El efecto de la provincia sobre el área de siembra difiere según el tipo de cultivo?**

8.2 Pregunta estadística

Específicamente:

- ¿Existe interacción significativa entre provincia y tipo de cultivo?
- ¿Cómo varía el efecto de la provincia para cada tipo de cultivo?
- ¿Cómo se interpretan y visualizan estas interacciones?

8.3 Modelo estadístico con interacción

$$y_{ij} = \beta_0 + \beta_1 X_i + \beta_2 Z_j + \beta_3 (X_i \times Z_j) + \epsilon_{ij} \quad (4)$$

Donde:

- X_i : Provincia (factor A)
- Z_j : Tipo de cultivo (factor B)
- $X_i \times Z_j$: Término de interacción
- β_3 : Coeficiente de interacción

Interpretación del término de interacción:

- Si $\beta_3 = 0$: Los efectos son aditivos (sin interacción)
- Si $\beta_3 \neq 0$: Los efectos no son aditivos (hay interacción)

8.4 Preparación de datos

Listing 50: Preparación de datos para modelo con interacciones

```

1 # Seleccionar provincias principales y tipos de cultivo
2 datos_interaccion <- datos %>%
3   filter(
4     provincia %in% c("CUSCO", "CANCHIS", "QUISPICANCHI", "URUBAMBA"),
5     area_total > 0
6   ) %>%
7   mutate(
8     tipo_cultivo = case_when(
9       grepl("PAPA|OLLUKO|OCA", cultivo) ~ "Tubérculos",
10      grepl("MAIZ|TRIGO|CEBADA|QUINUA", cultivo) ~ "Cereales",
11      grepl("HABA|ARVEJA|TARWI", cultivo) ~ "Leguminosas",
12      TRUE ~ "Otros"
13    ),
14    tipo_cultivo = as.factor(tipo_cultivo),
15    log_area = log(area_total + 1)
16  ) %>%
17  filter(tipo_cultivo != "Otros") %>%
18  droplevels()
19
20 # Resumen por combinación
21 datos_interaccion %>%
22   group_by(provincia, tipo_cultivo) %>%
23   summarise(
24     n = n(),
25     media = mean(area_total),
26     sd = sd(area_total),
27     .groups = "drop"
28   )

```

8.5 Modelos con y sin interacción

Listing 51: Ajuste de modelos con y sin interacción

```

1 # Modelo sin interacción (efectos principales)
2 modelo_aditivo <- lm(
3   log_area ~ provincia + tipo_cultivo,
4   data = datos_interaccion
5 )
6
7 # Modelo con interacción
8 modelo_interaccion <- lm(
9   log_area ~ provincia * tipo_cultivo,
10  data = datos_interaccion
11 )
12
13 # El símbolo * indica: efectos principales + interacción
14 # Equivalente a: provincia + tipo_cultivo + provincia:tipo_cultivo
15
16 # Comparar modelos
17 compare_performance(modelo_aditivo, modelo_interaccion, rank = TRUE)
18
19 # Test de significancia de la interacción
20 test_likelihoodratio(modelo_aditivo, modelo_interaccion)
21 anova(modelo_aditivo, modelo_interaccion)

```

8.6 Análisis del modelo con interacción

Listing 52: Parámetros del modelo con interacción

```

1 # Parámetros del modelo
2 params_int <- model_parameters(modelo_interaccion)
3 print(params_int)
4
5 # Los coeficientes de interacción representan la diferencia
6 # en el efecto de provincia para cada nivel de tipo_cultivo
7 # (respecto a los niveles de referencia)
8
9 # Desempeño del modelo
10 performance(modelo_interaccion)
11
12 # Tamaño de efecto
13 eta_squared(modelo_interaccion, partial = TRUE)

```

8.7 Interpretación de interacciones

8.7.1 Efectos simples (Simple Effects)

Los efectos simples muestran el efecto de un factor en cada nivel del otro factor.

Listing 53: Análisis de efectos simples

```

1 library(emmeans)
2
3 # Medias estimadas para cada combinación
4 emm_int <- emmeans(modelo_interaccion,
5                      specs = ~ provincia | tipo_cultivo)
6 print(emm_int)
7
8 # Efecto de provincia para cada tipo de cultivo
9 contrast(emm_int, method = "pairwise", by = "tipo_cultivo")
10
11 # Alternativamente, efecto de tipo de cultivo para cada provincia
12 emm_int2 <- emmeans(modelo_interaccion,
13                      specs = ~ tipo_cultivo | provincia)
14 contrast(emm_int2, method = "pairwise", by = "provincia")
15
16 # Con modelbased
17 efectos_prov <- estimate_slopes(
18   modelo_interaccion,
19   trend = "provincia",
20   at = "tipo_cultivo"
21 )
22 print(efectos_prov)
23 plot(efectos_prov)

```

8.8 Visualización de interacciones

La visualización es crucial para entender interacciones.

Listing 54: Visualización de interacciones

```

1 # 1. Gráfico de perfiles (líneas)
2 medias_int <- estimate_means(
3   modelo_interaccion,
4   at = c("provincia", "tipo_cultivo")
5 )
6
7 library(ggplot2)
8 ggplot(medias_int, aes(x = provincia, y = Mean,
9                         color = tipo_cultivo, group = tipo_cultivo)) +
10   geom_line(size = 1.2) +
11   geom_point(size = 3) +
12   geom_errorbar(aes(ymin = CI_low, ymax = CI_high),
13                 width = 0.2, alpha = 0.7) +
14   labs(
15     title = "Interacción: Provincia Tipo de Cultivo",
16     subtitle = "Efecto sobre Log(Área de Siembra)",
17     x = "Provincia",
18     y = "Media Estimada (Log-escala)",

```

```

19     color = "Tipo de Cultivo"
20 ) +
21 theme_minimal() +
22 theme(
23   axis.text.x = element_text(angle = 45, hjust = 1),
24   legend.position = "bottom"
25 )
26
27 # Interpretación:
28 # - Líneas paralelas No hay interacción
29 # - Líneas que se cruzan Interacción presente
30 # - Mayor diferencia entre líneas Interacción más fuerte
31
32 # 2. Gráfico de barras con facetas
33 ggplot(medias_int, aes(x = provincia, y = Mean, fill = provincia)) +
34   geom_col() +
35   geom_errorbar(aes(ymax = CI_high, ymin = CI_low), width = 0.3) +
36   facet_wrap(~ tipo_cultivo, ncol = 3) +
37   labs(
38     title = "Área de Siembra por Provincia y Tipo de Cultivo",
39     x = "Provincia",
40     y = "Media Estimada (Log-escala)"
41   ) +
42   theme_minimal() +
43   theme(
44     axis.text.x = element_text(angle = 45, hjust = 1),
45     legend.position = "none"
46   )
47
48 # 3. Heatmap de interacción
49 ggplot(medias_int, aes(x = provincia, y = tipo_cultivo, fill = Mean)) +
50   geom_tile() +
51   geom_text(aes(label = round(Mean, 2)), color = "white", size = 5) +
52   scale_fill_viridis_c() +
53   labs(
54     title = "Heatmap de Interacción",
55     x = "Provincia",
56     y = "Tipo de Cultivo",
57     fill = "Log(Área)"
58   ) +
59   theme_minimal() +
60   theme(axis.text.x = element_text(angle = 45, hjust = 1))

```

8.9 Descomposición de la interacción

Listing 55: Análisis detallado de la interacción

```

1 # Contrastes de interacción
2 contrast(emm_int, interaction = c("pairwise", "pairwise"))

```

```

3
4 # Con modelbased: contrastes específicos
5 contrasts_int <- estimate_contrasts(
6   modelo_interaccion,
7   contrast = c("provincia", "tipo_cultivo")
8 )
9 print(contrasts_int)
10
11 # Visualizar contrastes
12 plot(contrasts_int)

```

8.10 Interacciones continuas × categóricas

Cuando una variable es continua y otra categórica:

Listing 56: Interacción continua × categórica

```

1 # Modelo: área según número de meses y tipo de cultivo
2 modelo_cont_cat <- lm(
3   log_area ~ n_meses_siembra * tipo_cultivo,
4   data = datos_interaccion
5 )
6
7 # Parámetros
8 parameters(modelo_cont_cat)
9
10 # La interacción indica que la pendiente de n_meses difiere
11 # entre tipos de cultivo
12
13 # Visualizar pendientes
14 pendientes <- estimate_slopes(
15   modelo_cont_cat,
16   trend = "n_meses_siembra",
17   at = "tipo_cultivo"
18 )
19 print(pendientes)
20
21 # Gráfico de pendientes
22 ggplot(datos_interaccion,
23   aes(x = n_meses_siembra, y = log_area, color = tipo_cultivo)) +
24   geom_point(alpha = 0.3) +
25   geom_smooth(method = "lm", se = TRUE, size = 1.5) +
26   labs(
27     title = "Interacción entre Meses de Siembra y Tipo de Cultivo",
28     x = "Número de Meses con Siembra",
29     y = "Log(Área en hectáreas)",
30     color = "Tipo de Cultivo"
31   ) +
32   theme_minimal()

```

```

34 # Si las pendientes son similares poca interacción
35 # Si las pendientes difieren sustancialmente interacción fuerte

```

8.11 Interacciones continuas × continuas

Para dos variables continuas:

Listing 57: Interacción entre dos variables continuas

```

1 # Ejemplo: área según n_meses y log(area_total_sector)
2 # (requiere datos agregados por sector)
3
4 datos_continuas <- datos %>%
5   group_by(distrito, sector) %>%
6   summarise(
7     n_cultivos = n(),
8     area_total = sum(area_total),
9     n_meses_promedio = mean(n_meses_siembra),
10    .groups = "drop"
11  ) %>%
12  filter(area_total > 0) %>%
13  mutate(log_area = log(area_total))
14
15 modelo_cont_cont <- lm(
16   log_area ~ n_cultivos * n_meses_promedio,
17   data = datos_continuas
18 )
19
20 parameters(modelo_cont_cont)
21
22 # Superficie de respuesta 3D
23 library(modelbased)
24 pred_superficie <- estimate_response(
25   modelo_cont_cont,
26   data = expand.grid(
27     n_cultivos = seq(min(datos_continuas$n_cultivos),
28                      max(datos_continuas$n_cultivos),
29                      length.out = 30),
30     n_meses_promedio = seq(min(datos_continuas$n_meses_promedio),
31                           max(datos_continuas$n_meses_promedio),
32                           length.out = 30)
33   )
34 )
35
36 # Gráfico de contorno
37 ggplot(pred_superficie,
38         aes(x = n_cultivos, y = n_meses_promedio, z = Predicted)) +
39         geom_contour_filled(alpha = 0.8) +
40         geom_point(data = datos_continuas,
41                    aes(x = n_cultivos, y = n_meses_promedio)),

```

```

42         inherit.aes = FALSE, alpha = 0.3) +
43     labs(
44       title = "Superficie de Respuesta: Interacción Continua",
45       x = "Número de Cultivos",
46       y = "Número de Meses Promedio",
47       fill = "Log(Área)\nPredicha"
48     ) +
49     theme_minimal()

```

8.12 Cuándo incluir interacciones

Razones teóricas para incluir interacciones:

- Existe evidencia previa o teoría que sugiere interacción
- Los factores están relacionados conceptualmente
- Se busca modelar efectos contextuales o de moderación

Razones empíricas:

- La interacción mejora significativamente el ajuste del modelo
- Los gráficos exploratorios sugieren no-aditividad
- El test de razón de verosimilitud es significativo

Precauciones:

- Las interacciones de orden superior (≥ 3 factores) son difíciles de interpretar
- Aumentan la complejidad del modelo y reducen grados de libertad
- Requieren más datos para estimarse con precisión
- Pueden llevar a sobreajuste si no están justificadas

8.13 Verificación de supuestos

Listing 58: Diagnósticos del modelo con interacción

```

1 # Diagnósticos visuales completos
2 check_model(modelo_interaccion)
3
4 # Normalidad
5 check_normality(modelo_interaccion)
6
7 # Homoscedasticidad
8 check_heteroscedasticity(modelo_interaccion)
9
10 # Multicolinealidad (importante en modelos con interacciones)

```

```

11 check_collinearity(modelo_interaccion)
12
13 # Nota: La interacción X*Z está naturalmente correlacionada con X y Z
14 # VIF altos para términos de interacción son esperables
15 # Lo importante es centrar las variables antes si VIF > 10

```

8.14 Centrado de variables

Para reducir multicolinealidad en interacciones con variables continuas:

Listing 59: Centrado de variables para reducir multicolinealidad

```

1 # Centrar variables continuas (restar la media)
2 datos_centrados <- datos_continuas %>%
3   mutate(
4     n_cultivos_c = n_cultivos - mean(n_cultivos),
5     n_meses_promedio = n_meses_promedio - mean(n_meses_promedio)
6   )
7
8 # Modelo con variables centradas
9 modelo_centrado <- lm(
10   log_area ~ n_cultivos_c * n_meses_c,
11   data = datos_centrados
12 )
13
14 # Comparar multicolinealidad
15 compare_performance(modelo_cont_cont, modelo_centrado)
16 check_collinearity(modelo_cont_cont)
17 check_collinearity(modelo_centrado)
18
19 # Las estimaciones son idénticas, pero la multicolinealidad se reduce
20 # La interpretación del intercepto cambia (ahora es el valor en las medias)

```

8.15 Reporte de modelos con interacción

Listing 60: Reporte automático de interacciones

```

1 # Reporte completo
2 report(modelo_interaccion)
3
4 # Reporte enfocado en la interacción
5 report_text(modelo_interaccion)
6
7 # Tabla de parámetros
8 report_table(modelo_interaccion)

```

8.16 Ejemplo de interpretación integrada

Se ajustó un modelo de regresión lineal con interacción para evaluar si el efecto de la provincia sobre el área de siembra varía según el tipo de cultivo. El modelo con interacción ($R^2 = 0.52$) tuvo un ajuste significativamente mejor que el modelo aditivo ($R^2 = 0.48$), $\chi^2(6) = 45.2$, $p < 0.001$, indicando que la interacción es significativa.

El análisis de efectos simples reveló que el efecto de la provincia difiere sustancialmente entre tipos de cultivo. Para Tubérculos, la provincia de Canchis siembra 45% más área que Cusco ($b = 0.37$, IC 95% [0.25, 0.49], $p < 0.001$), mientras que para Cereales esta diferencia es solo del 15% ($b = 0.14$, IC 95% [0.05, 0.23], $p < 0.01$). Para Leguminosas, no se observó diferencia significativa entre provincias ($b = 0.08$, IC 95% [-0.05, 0.21], $p = 0.23$).

Estos resultados sugieren que las políticas agrícolas deben considerar no solo las diferencias provinciales, sino también cómo estas varían según el tipo de cultivo cultivado.

8.17 Resumen del análisis de interacciones

Pasos para analizar interacciones con easystats:

1. Justificar teórica o empíricamente la inclusión de interacciones
2. Ajustar modelo con término de interacción (`X * Z`)
3. Comparar con modelo aditivo usando `compare_performance()`
4. Extraer parámetros con `parameters()`
5. Analizar efectos simples con `emmeans` o `estimate_slopes()`
6. Visualizar interacciones con gráficos de perfiles o superficies
7. Verificar supuestos (especialmente multicolinealidad)
8. Centrar variables continuas si es necesario
9. Interpretar en contexto del problema
10. Reportar con `report()`

9 Caso 5: Modelos de Efectos Mixtos

Los modelos mixtos (también llamados multinivel o jerárquicos) son apropiados cuando los datos tienen estructura anidada o agrupada, como observaciones dentro de distritos, distritos dentro de provincias, etc.

9.1 Objetivo del análisis

Queremos responder: **¿Cómo modelar el área de siembra considerando la estructura jerárquica geográfica de los datos?**

9.2 ¿Cuándo usar modelos mixtos?

Los modelos mixtos son necesarios cuando:

- **Datos agrupados:** Observaciones anidadas en grupos (ej. sectores dentro de distritos)
- **Medidas repetidas:** Múltiples mediciones del mismo sujeto
- **Estructura jerárquica:** Niveles naturales de organización (estudiantes en escuelas, pacientes en hospitales)
- **Heterogeneidad entre grupos:** Los grupos difieren sistemáticamente
- **Variabilidad intra e intergrupo:** Interesa modelar ambas fuentes de variación

Ventajas sobre regresión estándar:

- Controla por la no independencia de observaciones
- Proporciona estimaciones específicas por grupo (BLUPs)
- Permite modelar variabilidad en interceptos y pendientes
- Más potencia estadística y estimaciones precisas

9.3 Estructura de datos jerárquicos

En nuestro dataset:

- **Nivel 1:** Observaciones individuales (cultivo en sector)
- **Nivel 2:** Sectores estadísticos
- **Nivel 3:** Distritos
- **Nivel 4:** Provincias

Listing 61: Exploración de estructura jerárquica

```

1 # Contar observaciones por nivel
2 datos %>%
3   summarise(
4     n_obs = n(),
5     n_sectores = n_distinct(sector),
6     n_distritos = n_distinct(distrito),
7     n_provincias = n_distinct(provincia)

```

```

8   )
9
10 # Observaciones por grupo
11 datos %>%
12   group_by(distrito) %>%
13   summarise(n = n()) %>%
14   summary()
15
16 # Idealmente, cada grupo debería tener múltiples observaciones

```

9.4 Modelo estadístico de efectos mixtos

9.4.1 Modelo con intercepto aleatorio

El modelo más simple incluye solo interceptos aleatorios:

$$y_{ij} = (\beta_0 + u_{0j}) + \beta_1 X_{ij} + \epsilon_{ij} \quad (5)$$

Donde:

- y_{ij} : Observación i en grupo j
- β_0 : Intercepto fijo (media poblacional)
- u_{0j} : Desviación del grupo j respecto al intercepto ($u_{0j} \sim N(0, \tau^2)$)
- β_1 : Pendiente fija (igual para todos los grupos)
- ϵ_{ij} : Residuo ($\epsilon_{ij} \sim N(0, \sigma^2)$)

Interpretación: Cada grupo tiene su propio intercepto, pero la pendiente es común.

9.4.2 Modelo con pendiente aleatoria

Podemos permitir que la pendiente también varíe entre grupos:

$$y_{ij} = (\beta_0 + u_{0j}) + (\beta_1 + u_{1j}) X_{ij} + \epsilon_{ij} \quad (6)$$

Donde u_{1j} es la desviación de la pendiente del grupo j respecto a la pendiente media.

9.5 Ajuste de modelos mixtos en R

Listing 62: Preparación de datos para modelos mixtos

```

1 # Filtrar y preparar datos
2 datos_mixtos <- datos %>%
3   filter(
4     area_total > 0,
5     provincia %in% c("CUSCO", "CANCHIS", "QUISPICANCHI", "URUBAMBA")
6   ) %>%
7   mutate(

```

```

8     tipo_cultivo = case_when(
9       grepl("PAPA|OLLUCO|OCA", cultivo) ~ "Tubérculos",
10      grepl("MAIZ|TRIGO|CEBADA|QUINUA", cultivo) ~ "Cereales",
11      grepl("HABA|ARVEJA|TARWI", cultivo) ~ "Leguminosas",
12      TRUE ~ "Otros"
13    ),
14    tipo_cultivo = as.factor(tipo_cultivo),
15    log_area = log(area_total + 1)
16  ) %>%
17  filter(tipo_cultivo != "Otros") %>%
18  droplevels()
19
20 # Verificar número de observaciones por distrito
21 datos_mixtos %>%
22   group_by(distrito) %>%
23   summarise(n = n()) %>%
24   arrange(n)
25
26 # Filtrar distritos con suficientes observaciones (>= 5)
27 distritos_suficientes <- datos_mixtos %>%
28   group_by(distrito) %>%
29   summarise(n = n()) %>%
30   filter(n >= 5) %>%
31   pull(distrito)
32
33 datos_mixtos <- datos_mixtos %>%
34   filter(distrito %in% distritos_suficientes) %>%
35   droplevels()

```

9.5.1 Modelo nulo (solo intercepto aleatorio)

Listing 63: Modelo nulo con lme4

```

1 library(lme4)
2
3 # Modelo nulo (sin predictores, solo intercepto aleatorio por distrito)
4 modelo_nulo <- lmer(
5   log_area ~ 1 + (1 | distrito),
6   data = datos_mixtos,
7   REML = TRUE
8 )
9
10 summary(modelo_nulo)
11
12 # Con easystats
13 parameters(modelo_nulo)
14 performance(modelo_nulo)

```

Cálculo del ICC (Coeficiente de Correlación Intraclass):

Listing 64: Coeficiente de Correlación Intraclass

```

1 # ICC: proporción de varianza entre grupos
2 library(performance)
3 icc_resultado <- icc(modelo_nulo)
4 print(icc_resultado)
5
6 # Interpretación del ICC:
7 # ICC = 0: No hay agrupación (usar regresión estándar)
8 # ICC = 0.05-0.10: Agrupación pequeña pero presente
9 # ICC > 0.10: Agrupación sustancial (modelos mixtos recomendados)
10 # ICC cercano a 1: Observaciones dentro de grupos son muy similares

```

9.5.2 Modelo con predictores de nivel 1

Listing 65: Modelo con predictores fijos e intercepto aleatorio

```

1 # Modelo con tipo de cultivo (predictor de nivel 1)
2 modelo_nivel1 <- lmer(
3   log_area ~ tipo_cultivo + (1 | distrito),
4   data = datos_mixtos,
5   REML = TRUE
6 )
7
8 # Parámetros
9 params_m1 <- model_parameters(modelo_nivel1)
10 print(params_m1)
11
12 # Desempeño
13 perf_m1 <- model_performance(modelo_nivel1)
14 print(perf_m1)
15
16 # Comparar con modelo nulo
17 compare_performance(modelo_nulo, modelo_nivel1)
18 test_likelihoodratio(modelo_nulo, modelo_nivel1)

```

9.5.3 Modelo con predictores de nivel 2

Listing 66: Modelo con predictor de nivel 2 (provincia)

```

1 # Provincia es nivel 2 (contiene distritos)
2 # Primero agregar provincia al nivel de distrito
3 datos_mixtos <- datos_mixtos %>%
4   mutate(provincia = as.factor(provincia))
5
6 modelo_nivel2 <- lmer(
7   log_area ~ tipo_cultivo + provincia + (1 | distrito),
8   data = datos_mixtos,
9   REML = TRUE

```

```

10  )
11
12 # Nota: provincia es un predictor de nivel 2, pero distrito está
13 # anidado en provincia, así que el efecto aleatorio controla
14 # por variabilidad dentro de provincia
15
16 parameters(modelo_nivel2)
17 performance(modelo_nivel2)

```

9.5.4 Modelo con pendiente aleatoria

Listing 67: Modelo con intercepto y pendiente aleatorios

```

1 # Permitir que el efecto de tipo_cultivo varíe por distrito
2 modelo_pendiente <- lmer(
3   log_area ~ tipo_cultivo + (1 + tipo_cultivo | distrito),
4   data = datos_mixtos,
5   REML = TRUE
6 )
7
8 # Este modelo permite que cada distrito tenga:
9 # - Su propio intercepto
10 # - Su propia pendiente para tipo_cultivo
11
12 summary(modelo_pendiente)
13 parameters(modelo_pendiente)
14
15 # Comparar modelos
16 compare_performance(modelo_nivel1, modelo_pendiente)
17
18 # Nota: Modelos con pendientes aleatorias requieren más datos
19 # y pueden no converger si hay pocos grupos o datos por grupo

```

9.5.5 Estructura de efectos aleatorios anidados

Listing 68: Efectos aleatorios anidados en múltiples niveles

```

1 # Efectos aleatorios en distrito Y provincia
2 modelo_anidado <- lmer(
3   log_area ~ tipo_cultivo + n_meses_siembra +
4     (1 | provincia/distrito),
5   data = datos_mixtos,
6   REML = TRUE
7 )
8
9 # La sintaxis (1 / provincia/distrito) es equivalente a:
10 # (1 / provincia) + (1 / provincia:distrito)
11 # Es decir, efectos aleatorios en provincia y en distrito dentro de
  provincia

```

```

12
13 summary(modelo_anidado)
14 parameters(modelo_anidado)
15
16 # Extrae componentes de varianza
17 get_variance(modelo_anidado)

```

9.6 Análisis de resultados

9.6.1 Extracción de efectos fijos

Listing 69: Parámetros de efectos fijos

```

1 # Efectos fijos (parámetros poblacionales)
2 efectos_fijos <- model_parameters(
3   modelo_nivel2,
4   effects = "fixed"
5 )
6 print(efectos_fijos)
7
8 # Visualizar efectos fijos
9 plot(efectos_fijos)

```

9.6.2 Extracción de efectos aleatorios

Listing 70: Efectos aleatorios (BLUPs)

```

1 # Efectos aleatorios (desviaciones por grupo)
2 efectos_aleatorios <- ranef(modelo_nivel2)
3 print(efectos_aleatorios)
4
5 # BLUPs: Best Linear Unbiased Predictors
6 # Son las estimaciones específicas de cada grupo
7
8 # Visualizar con lattice
9 library(lattice)
10 dotplot(ranef(modelo_nivel2, condVar = TRUE))
11
12 # Distritos con efectos aleatorios más extremos
13 efectos_df <- as.data.frame(efectos_aleatorios)
14 efectos_df %>% arrange(desc(condval))
15
16 # Estos representan distritos con mayor/menor área
17 # de siembra después de controlar por tipo de cultivo

```

9.6.3 Componentes de varianza

Listing 71: Descomposición de varianza

```

1 # Extraer componentes de varianza
2 var_components <- get_variance(modelo_nivel2)
3 print(var_components)

4

5 # Varianza explicada por efectos fijos
6 var_fixed <- var_components$var.fixed

7

8 # Varianza entre grupos (distrito)
9 var_random <- var_components$var.random

10

11 # Varianza residual (dentro de grupos)
12 var_residual <- var_components$var.residual

13

14 # R condicional: varianza explicada por efectos fijos + aleatorios
15 # R marginal: varianza explicada solo por efectos fijos
16 r2_modelo <- r2(modelo_nivel2)
17 print(r2_modelo)

18

19 # R marginal: proporción explicada por predictores fijos
20 # R condicional: proporción explicada por modelo completo (fijos +
aleatorios)

```

9.7 Evaluación del modelo

Listing 72: Diagnósticos de modelos mixtos

```

1 # Diagnósticos visuales
2 check_model(modelo_nivel2)

3

4 # Normalidad de residuos
5 check_normality(modelo_nivel2)
6 plot(check_normality(modelo_nivel2))

7

8 # Normalidad de efectos aleatorios
9 qqnorm(ranef(modelo_nivel2)$distrito[,1])
10 qqline(ranef(modelo_nivel2)$distrito[,1])

11

12 # Homoscedasticidad
13 check_heteroscedasticity(modelo_nivel2)

14

15 # Valores influyentes
16 check_outliers(modelo_nivel2)

17

18 # Convergencia (importante en modelos mixtos complejos)
19 check_convergence(modelo_nivel2)

```

9.8 Predicciones con modelos mixtos

Listing 73: Predicciones de modelos mixtos

```

1 # Predicciones de nivel poblacional (usando solo efectos fijos)
2 pred_poblacion <- predict(modelo_nivel2, re.form = NA)
3
4 # Predicciones de nivel de grupo (usando efectos fijos + aleatorios)
5 pred_grupo <- predict(modelo_nivel2)
6
7 # Con modelbased
8 # Medias estimadas por tipo de cultivo (nivel poblacional)
9 medias_poblacion <- estimate_means(
10   modelo_nivel2,
11   at = "tipo_cultivo"
12 )
13 print(medias_poblacion)
14 plot(medias_poblacion)
15
16 # Predicciones específicas por distrito
17 # (usando el efecto aleatorio de cada distrito)
18 nuevos_datos <- expand.grid(
19   tipo_cultivo = levels(datos_mixtos$tipo_cultivo),
20   distrito = unique(datos_mixtos$distrito)[1:5], # Primeros 5 distritos
21   provincia = unique(datos_mixtos$provincia)[1]
22 )
23
24 pred_distrito <- predict(modelo_nivel2, newdata = nuevos_datos)

```

9.9 Visualización de modelos mixtos

Listing 74: Visualizaciones para modelos mixtos

```

1 # 1. Gráfico de efectos aleatorios (caterpillar plot)
2 library(sjPlot)
3 plot_model(modelo_nivel2, type = "re")
4
5 # 2. Gráfico de efectos fijos
6 plot_model(modelo_nivel2, type = "est")
7
8 # 3. Predicciones por grupo
9 library(ggeffects)
10 pred_eff <- ggpredict(modelo_nivel2, terms = c("tipo_cultivo", "distrito [",
11   sample=10]"))
11 plot(pred_eff)
12
13 # 4. Comparar interceptos entre distritos
14 ggplot(data.frame(
15   distrito = rownames(ranef(modelo_nivel2)$distrito),

```

```

16     intercepto = ranef(modelo_nivel2)$distrito[,1]
17   )) +
18   aes(x = reorder(distrito, intercepto), y = intercepto) +
19   geom_point(size = 3) +
20   geom_hline(yintercept = 0, linetype = "dashed", color = "red") +
21   coord_flip() +
22   labs(
23     title = "Efectos Aleatorios por Distrito",
24     subtitle = "Desviaciones del intercepto poblacional",
25     x = "Distrito",
26     y = "Efecto Aleatorio (Intercepto)"
27   ) +
28   theme_minimal()
29
30 # 5. Gráfico de líneas por grupo (si hay pendiente aleatoria)
31 # Mostrar regresión específica de cada grupo
32 if("tipo_cultivo" %in% names(ranef(modelo_pendiente)$distrito)) {
33   ggplot(datos_mixtos, aes(x = tipo_cultivo, y = log_area, group = distrito
34     )) +
35     geom_line(alpha = 0.3) +
36     geom_line(stat = "smooth", method = "lm", se = FALSE, aes(group = NULL)
37       ,
38       size = 2, color = "red") +
39   labs(
40     title = "Pendientes Específicas por Distrito",
41     subtitle = "Línea roja = efecto fijo poblacional",
42     x = "Tipo de Cultivo",
43     y = "Log(Área)"
44   ) +
45   theme_minimal()
46 }

```

9.10 Comparación de modelos mixtos

Listing 75: Selección del mejor modelo mixto

```

1 # Comparar múltiples modelos
2 comp_mixtos <- compare_performance(
3   modelo_nulo,
4   modelo_nivel1,
5   modelo_nivel2,
6   modelo_pendiente,
7   rank = TRUE
8 )
9 print(comp_mixtos)
10 plot(comp_mixtos)
11
12 # Criterios de selección:
13 # - AIC/BIC más bajo

```

```

14 # - R condicional más alto
15 # - Parsimonia (modelo más simple que explica bien)
16
17 # Test de razón de verosimilitud (para modelos anidados)
18 # IMPORTANTE: Usar REML = FALSE para comparar modelos con diferentes
   efectos fijos
19 modelo_ml1 <- lmer(log_area ~ tipo_cultivo + (1 | distrito),
   data = datos_mixtos, REML = FALSE)
20 modelo_ml2 <- lmer(log_area ~ tipo_cultivo + provincia + (1 | distrito),
   data = datos_mixtos, REML = FALSE)
21
22 test_likelihoodratio(modelo_ml1, modelo_ml2)
23
24

```

9.11 Modelos mixtos generalizados (GLMM)

Para respuestas no gaussianas (ej. binarias, conteos):

Listing 76: Modelo logístico mixto

```

1 # Ejemplo: Presencia/ausencia de quinua por sector
2 datos_glmm <- datos %>%
3   mutate(
4     siembra_quinua = ifelse(cultivo == "QUINUA" & area_total > 0, 1, 0)
5   ) %>%
6   group_by(distrito, sector) %>%
7   summarise(
8     siembra_quinua = max(siembra_quinua),
9     n_cultivos = n(),
10    .groups = "drop"
11  ) %>%
12  filter(distrito %in% distritos_suficientes)
13
14 # Modelo logístico mixto
15 modelo_logit_mixto <- glmer(
16   siembra_quinua ~ n_cultivos + (1 | distrito),
17   data = datos_glmm,
18   family = binomial(link = "logit")
19 )
20
21 summary(modelo_logit_mixto)
22 parameters(modelo_logit_mixto, exponentiate = TRUE) # Odds Ratios
23 performance(modelo_logit_mixto)
24
25 # ICC para modelos binomiales
26 icc(modelo_logit_mixto)

```

9.12 Reporte de modelos mixtos

Listing 77: Reporte automático de modelos mixtos

```

1 # Reporte completo
2 report(modelo_nivel2)
3
4 # Reporte de efectos fijos
5 report_parameters(modelo_nivel2, effects = "fixed")
6
7 # Reporte de efectos aleatorios
8 report_parameters(modelo_nivel2, effects = "random")
9
10 # Reporte de desempeño
11 report_performance(modelo_nivel2)

```

9.13 Ejemplo de interpretación integrada

Se ajustó un modelo lineal mixto para predecir el área de siembra (log-transformada), incluyendo tipo de cultivo y provincia como efectos fijos, y distrito como efecto aleatorio. El modelo mostró un buen ajuste (R^2 condicional = 0.58, R^2 marginal = 0.35), indicando que los efectos fijos explican el 35% de la varianza y el modelo completo (fijos + aleatorios) explica el 58%.

El coeficiente de correlación intraclass fue $ICC = 0.28$, indicando que el 28% de la varianza en área de siembra se debe a diferencias entre distritos, justificando el uso de modelos mixtos. El tipo de cultivo tuvo un efecto significativo: Los Tubérculos se siembran en mayor área que los Cereales ($b = 0.42$, $SE = 0.08$, $t = 5.25$, $p < 0.001$) y las Leguminosas ($b = 0.68$, $SE = 0.09$, $t = 7.56$, $p < 0.001$).

Los efectos aleatorios mostraron heterogeneidad sustancial entre distritos ($SD = 0.52$), con algunos distritos sembrado consistentemente más área que otros, incluso después de controlar por tipo de cultivo y provincia. Esta variabilidad podría deberse a factores locales no medidos como altitud, acceso a agua, o tradiciones agrícolas específicas.

9.14 Resumen del flujo de modelos mixtos

Pasos para analizar con modelos mixtos:

1. Identificar estructura jerárquica o agrupación en los datos
2. Calcular ICC con modelo nulo para justificar modelo mixto
3. Ajustar modelos con `lmer()` (continuos) o `glmer()` (otros)
4. Especificar efectos aleatorios apropiados: `(1 | grupo)`
5. Extraer efectos fijos y aleatorios con `parameters()`
6. Evaluar con R^2 marginal y condicional

7. Verificar supuestos con `check_model()`
8. Comparar modelos con `compare_performance()`
9. Visualizar efectos con `plot_model()` o `ggeffects`
10. Interpretar en términos de efectos poblacionales (fijos) y variabilidad entre grupos (aleatorios)
11. Reportar con `report()`

10 Visualización y Comunicación de Resultados

La visualización efectiva es fundamental para comunicar resultados estadísticos. El paquete `see` de easystats facilita la creación de gráficos de calidad publicable.

10.1 Filosofía de visualización en easystats

El paquete `see` extiende `ggplot2` con métodos especializados para objetos de easystats:

- **Automatización inteligente:** Detecta el tipo de objeto y crea el gráfico apropiado
- **Estética profesional:** Paletas de colores y temas diseñados para publicación
- **Personalizable:** Todos los gráficos son objetos `ggplot2` que pueden modificarse
- **Consistencia:** Misma sintaxis para diferentes tipos de análisis

10.2 Gráficos de parámetros

Listing 78: Visualización de coeficientes de regresión

```

1 library(see)
2
3 # Ajustar modelo
4 modelo <- lm(log_area ~ provincia + tipo_cultivo + n_meses_siembra,
5               data = datos_papa)
6
7 # Extraer parámetros
8 params <- model_parameters(modelo)
9
10 # Gráfico básico
11 plot(params)
12
13 # Personalizar
14 plot(params,
15       show_intercept = FALSE, # Omitir intercepto
16       sort = TRUE) +          # Ordenar por magnitud
17       labs(title = "Coeficientes del Modelo de Regresión",
18             subtitle = "Intervalos de confianza al 95%") +

```

```

19   theme_modern()
20
21 # Parámetros estandarizados
22 params_std <- standardize_parameters(modelo)
23 plot(params_std) +
24   labs(title = "Coeficientes Estandarizados (Betas)")

```

10.3 Gráficos de diagnóstico

Listing 79: Diagnósticos visuales completos

```

1 # Diagnósticos automáticos (múltiples paneles)
2 check_model(modelo)
3
4 # Diagnósticos individuales
5 check_normality(modelo)
6 plot(check_normality(modelo))
7
8 check_heteroscedasticity(modelo)
9 plot(check_heteroscedasticity(modelo))
10
11 check_collinearity(modelo)
12 plot(check_collinearity(modelo))
13
14 check_outliers(modelo)
15 plot(check_outliers(modelo))
16
17 # Personalizar diagnósticos
18 check_model(
19   check = c("linearity", "homogeneity", "qq", "outliers"),
20   colors = c("blue", "red"))

```

10.4 Gráficos de efectos y predicciones

Listing 80: Visualización de efectos marginales

```

1 library(modelbased)
2
3 # Medias estimadas
4 medias <- estimate_means(modelo, at = "tipo_cultivo")
5 plot(medias) +
6   labs(title = "Medias Estimadas por Tipo de Cultivo",
7       y = "Log(Área de Siembra)")
8
9 # Efectos según variable continua
10 efecto_meses <- estimate_relation(modelo, at = "n_meses_siembra")
11 plot(efecto_meses) +
12   labs(title = "Efecto del Número de Meses de Siembra",

```

```

13     x = "Meses de Siembra",
14     y = "Predicción de Log(Área)"
15
16 # Contrastes entre grupos
17 contrasts <- estimate_contrasts(modelo, contrast = "tipo_cultivo")
18 plot(contrasts) +
19   labs(title = "Contrastes entre Tipos de Cultivo")

```

10.5 Gráficos de comparación de modelos

Listing 81: Comparación visual de modelos

```

1 # Comparar desempeño
2 m1 <- lm(log_area ~ provincia, data = datos_papa)
3 m2 <- lm(log_area ~ provincia + tipo_papa, data = datos_papa)
4 m3 <- lm(log_area ~ provincia + tipo_papa + n_meses_siembra, data = datos_
  papa)
5
6 comp <- compare_performance(m1, m2, m3)
7 plot(comp) +
8   labs(title = "Comparación de Modelos")
9
10 # Visualizar múltiples métricas
11 plot(comp, type = "radar")

```

10.6 Gráficos para modelos logísticos

Listing 82: Visualizaciones específicas para regresión logística

```

1 # Modelo logístico
2 modelo_logit <- glm(siembra_quinua ~ provincia + n_cultivos,
3                      data = datos_quinua, family = binomial)
4
5 # Odds Ratios
6 or_params <- model_parameters(modelo_logit, exponentiate = TRUE)
7 plot(or_params) +
8   geom_vline(xintercept = 1, linetype = "dashed", color = "red") +
9   labs(title = "Odds Ratios",
10        subtitle = "OR > 1: aumenta probabilidad | OR < 1: disminuye
11                  probabilidad")
12
13 # Probabilidades predichas
14 probs <- estimate_expectation(modelo_logit, at = "provincia")
15 plot(probs) +
16   labs(title = "Probabilidad de Siembra de Quinua por Provincia",
17        y = "Probabilidad")

```

10.7 Paletas de colores y temas

Listing 83: Personalización de estética

```

1 library(see)
2
3 # Temas disponibles
4 # theme_modern() - Limpio y moderno
5 # theme_lucid() - Alto contraste
6 # theme_abyss() - Fondo oscuro
7 # theme_blackboard() - Estilo pizarra
8
9 # Paletas de colores
10 # Paletas discretas
11 plot(params) +
12   scale_color_material_d() + # Material Design
13   theme_modern()
14
15 # Paletas continuas
16 plot(params) +
17   scale_color_viridis_c() +
18   theme_lucid()
19
20 # Paletas personalizadas
21 colores_cusco <- c("#E63946", "#F1FAEE", "#A8DADC", "#457B9D", "#1D3557")
22 plot(medias) +
23   scale_fill_manual(values = colores_cusco) +
24   theme_abyss()

```

10.8 Exportación de gráficos

Listing 84: Guardar gráficos para publicación

```

1 # Crear gráfico
2 p <- plot(params, show_intercept = FALSE) +
3   labs(title = "Coeficientes del Modelo de Papa",
4       subtitle = "Región Cusco, Campaña 2023-2024") +
5   theme_modern()
6
7 # Guardar en alta resolución
8 ggsave(
9   filename = "grafico_coeficientes.png",
10  plot = p,
11  width = 10,
12  height = 6,
13  dpi = 300,
14  units = "in"
15 )
16

```

```

17 # Formato PDF (vectorial, ideal para publicaciones)
18 ggsave(
19   filename = "grafico_coeficientes.pdf",
20   plot = p,
21   width = 10,
22   height = 6
23 )
24
25 # Formato TIFF (requerido por algunas revistas)
26 ggsave(
27   filename = "grafico_coeficientes.tiff",
28   plot = p,
29   width = 10,
30   height = 6,
31   dpi = 600,
32   compression = "lzw"
33 )

```

10.9 Gráficos combinados (multipanel)

Listing 85: Combinación de múltiples gráficos

```

1 library(patchwork)
2
3 # Crear gráficos individuales
4 g1 <- plot(params, show_intercept = FALSE) +
5   labs(title = "A) Coeficientes")
6
7 g2 <- plot(check_normality(modelo)) +
8   labs(title = "B) Normalidad")
9
10 g3 <- plot(medias) +
11   labs(title = "C) Medias Estimadas")
12
13 g4 <- plot(check_collinearity(modelo)) +
14   labs(title = "D) Multicolinealidad")
15
16 # Combinar en panel 2x2
17 (g1 | g2) / (g3 | g4) +
18   plot_annotation(
19     title = "Análisis Completo del Modelo de Siembra de Papa",
20     subtitle = "Región Cusco",
21     caption = "Fuente: GERAGRI Cusco, Encuesta de Intenciones de Siembra
22       2023-2024"
23   )
24
25 # Guardar panel
26 ggsave("panel_analisis.png", width = 14, height = 10, dpi = 300)

```

10.10 Tablas para publicación

Listing 86: Exportación de tablas en formato publicable

```

1 # Tabla de parámetros
2 params_table <- model_parameters(modelo)
3
4 # Exportar como dataframe
5 as.data.frame(params_table)
6
7 # Formato APA con kableExtra
8 library(kableExtra)
9 kable(params_table,
10       format = "latex",
11       booktabs = TRUE,
12       caption = "Parámetros del modelo de regresión lineal") %>%
13       kable_styling(latex_options = c("striped", "hold_position"))
14
15 # Formato para Word
16 library(flextable)
17 flextable(as.data.frame(params_table)) %>%
18   theme_vanilla() %>%
19   autofit()
20
21 # Exportar a Excel
22 library(writexl)
23 write_xlsx(as.data.frame(params_table), "parametros_modelo.xlsx")

```

10.11 Reportes reproducibles con R Markdown

Listing 87: Integración con R Markdown

```

1 # En un documento .Rmd:
2 # `'{r analisis, echo=TRUE, message=FALSE, warning=FALSE}`
3 # library(easystats)
4 #
5 # # Cargar datos
6 # datos <- readRDS("datos_limpios.rds")
7 #
8 # # Ajustar modelo
9 # modelo <- lm(log_area ~ provincia + tipo_papa, data = datos)
10 #
11 # # Reportar automáticamente
12 # report(modelo)
13 # ``
14 #
15 # `'{r diagnosticos, fig.width=12, fig.height=8}`
16 # check_model(modelo)
17 # ``

```

```

18 #
19 # El reporte se genera automáticamente en HTML, PDF o Word

```

11 Buenas Prácticas y Errores Comunes

Esta sección recopila recomendaciones para un análisis estadístico robusto y reproducible, así como errores frecuentes a evitar.

11.1 Buenas prácticas generales

11.1.1 Documentación y reproducibilidad

Listing 88: Estructura de proyecto reproducible

```

1 # Estructura recomendada de directorios:
2 # proyecto/
3 #         datos/
4 #             raw/          # Datos originales (nunca modificar)
5 #                 processed/    # Datos limpios
6 #                     scripts/
7 #                         01_cargar_datos.R
8 #                         02_limpiar_datos.R
9 #                         03_analisis_exploratorio.R
10 #                         04_modelamiento.R
11 #                     resultados/
12 #                         figuras/
13 #                         tablas/
14 #                     reportes/
15 #                     README.md
16
17 # Usar proyectos de RStudio (.Rproj)
18 # Usar control de versiones (Git)
19
20 # Documentar sesión
21 sessionInfo()
22 writeLines(capture.output(sessionInfo()), "session_info.txt")
23
24 # Usar here() para rutas relativas
25 library(here)
26 datos <- read.csv(here("datos", "processed", "datos_limpios.csv"))

```

11.1.2 Comentarios y estilo de código

Listing 89: Código bien documentado

```

1 # BUENA PRÁCTICA: Comentarios claros y código legible
2
3 # Cargar paquetes necesarios

```

```

4 library(easystats)
5 library(tidyverse)
6
7 # Cargar y preparar datos de siembra
8 datos <- read.csv("intencionesSIEMBRA.csv", encoding = "Latin1") %>%
9   # Filtrar solo observaciones con siembra
10  filter(area_total > 0) %>%
11  # Crear variable log-transformada para normalizar
12  mutate(log_area = log(area_total + 1)) %>%
13  # Convertir variables categóricas a factor
14  mutate(across(c(provincia, cultivo), as.factor))
15
16 # Ajustar modelo de regresión lineal
17 # Variable respuesta: log del área
18 # Predictores: provincia y tipo de cultivo
19 modelo <- lm(log_area ~ provincia + tipo_cultivo, data = datos)
20
21 # MALA PRÁCTICA: Código sin comentarios
22 d<-read.csv("i.csv",encoding="L1")%>%filter(at>0)%>%mutate(la=log(at+1))%>%
  mutate(across(c(p,c),as.factor))
23 m<-lm(la~p+tc,data=d)

```

11.1.3 Flujo de análisis sistemático

1. **Planificación:** Definir preguntas de investigación y análisis antes de ver los datos
2. **Exploración:** Entender los datos antes de modelar
3. **Limpieza:** Documentar todas las decisiones de limpieza
4. **Modelamiento:** Ajustar modelo, verificar supuestos, interpretar
5. **Validación:** Verificar robustez con análisis de sensibilidad
6. **Comunicación:** Reportar con claridad, incluir limitaciones

11.2 Errores comunes en análisis estadístico

11.2.1 Error 1: P-hacking y HARKing

P-hacking: Probar múltiples análisis hasta encontrar un valor $p \leq 0.05$.

HARKing: Hypothesizing After Results are Known (formular hipótesis después de ver resultados).

Listing 90: Evitar P-hacking

```

1 # MALO: Probar muchos modelos hasta encontrar significancia
2 modelo1 <- lm(y ~ x1, data = datos) # p = 0.08, no significativo
3 modelo2 <- lm(y ~ x2, data = datos) # p = 0.06, no significativo
4 modelo3 <- lm(y ~ x3, data = datos) # p = 0.04, significativo!
5 # ... reportar solo modelo3

```

```

6
7 # BUENO: Pre-registrar análisis, reportar todos los modelos
8 # 1. Definir hipótesis a priori
9 # 2. Especificar análisis planeados
10 # 3. Reportar todos los resultados, no solo los "significativos"
11 # 4. Distinguir análisis confirmatorios de exploratorios
12
13 # Si múltiples modelos son plausibles, reportar todos
14 modelos <- list(
15   "Modelo 1" = lm(y ~ x1, data = datos),
16   "Modelo 2" = lm(y ~ x2, data = datos),
17   "Modelo 3" = lm(y ~ x1 + x2, data = datos)
18 )
19
20 compare_performance(modelos)

```

11.2.2 Error 2: Ignorar supuestos del modelo

Listing 91: Verificar siempre los supuestos

```

1 # MALO: Ajustar modelo y reportar sin verificar supuestos
2 modelo <- lm(area_total ~ provincia, data = datos)
3 summary(modelo) # Asumir que está bien
4
5 # BUENO: Verificar supuestos sistemáticamente
6 modelo <- lm(log_area ~ provincia, data = datos) # Log para normalizar
7
8 # Verificación completa
9 check_model(modelo)
10
11 # Verificaciones específicas
12 check_normality(modelo) # Residuos normales?
13 check_heteroscedasticity(modelo) # Varianza homogénea?
14 check_collinearity(modelo) # Multicolinealidad?
15 check_outliers(modelo) # Valores influyentes?
16
17 # Si supuestos se violan, tomar medidas correctivas:
18 # - Transformar variables
19 # - Usar modelos robustos
20 # - Usar métodos no paramétricos
21 # - Incluir más predictores

```

11.2.3 Error 3: Confundir significancia estadística con importancia práctica

Listing 92: Reportar tamaños de efecto

```

1 # MALO: Enfocarse solo en valores p
2 modelo <- lm(log_area ~ provincia, data = datos)

```

```

3  summary(modelo) # "p < 0.001, es significativo!"
4
5  # BUENO: Reportar significancia estadística Y tamaño de efecto
6  parameters(modelo)      # Coeficientes e IC
7  effectsize(modelo)       # Tamaño de efecto
8  cohens_f_squared(modelo) # f de Cohen
9
10 # Interpretar
11 # p < 0.001: Poco probable que no haya efecto (significancia)
12 # f = 0.02: Pero el efecto es muy pequeño (importancia práctica)
13
14 # Reportar: "El efecto es estadísticamente significativo (p < 0.001)
15 # pero de magnitud práctica pequeña (f = 0.02)."

```

11.2.4 Error 4: Sobreajuste del modelo

Listing 93: Evitar sobreajuste

```

1  # MALO: Incluir demasiados predictores
2  # Regla general: Al menos 10-20 observaciones por predictor
3  n <- nrow(datos) # 100 observaciones
4  modelo_sobreajustado <- lm(
5    y ~ x1 + x2 + x3 + x4 + x5 + x6 + x7 + x8 + x9 + x10, # 10 predictores
6    data = datos
7  ) # Ratio: 10 obs/predictor (límite)
8
9  # BUENO: Usar selección de variables justificada
10 # 1. Basar en teoría/conocimiento previo
11 # 2. Usar validación cruzada
12 # 3. Comparar modelos con criterios de información
13
14 # Selección paso a paso (con precaución)
15 library(MASS)
16 modelo_completo <- lm(y ~ ., data = datos)
17 modelo_step <- stepAIC(modelo_completo, direction = "both")
18
19 # Validación cruzada
20 library(caret)
21 train_control <- trainControl(method = "cv", number = 10)
22 modelo_cv <- train(y ~ ., data = datos, method = "lm",
23                      trControl = train_control)

```

11.2.5 Error 5: Ignorar datos faltantes

Listing 94: Manejo apropiado de valores faltantes

```

1  # MALO: Eliminar filas con NA sin pensar
2  datos_completos <- na.omit(datos) # Sesgo? Pérdida de información?

```

```

3
4 # BUENO: Analizar patrón de valores faltantes
5 library(naniar)
6 vis_miss(datos) # Visualizar NA
7 gg_miss_var(datos) # NA por variable
8
9 # Son MCAR (completamente aleatorios)?
10 # Son MAR (aleatorios condicionalmente)?
11 # Son MNAR (no aleatorios)?
12
13 # Opciones según el caso:
14 # 1. Si pocos NA y MCAR: na.omit() es aceptable
15 # 2. Si muchos NA o MAR: imputación múltiple
16 library(mice)
17 imp <- mice(datos, m = 5, method = "pmm")
18 modelo_imp <- with(imp, lm(y ~ x1 + x2))
19 pool(modelo_imp)
20
21 # 3. Si MNAR: modelos de selección o pattern-mixture

```

11.2.6 Error 6: Múltiples comparaciones sin ajuste

Listing 95: Ajustar por comparaciones múltiples

```

1 # MALO: Hacer muchas comparaciones sin ajuste
2 # Probabilidad de al menos un falso positivo aumenta con cada prueba
3 pvalues <- c()
4 for(i in 1:20) {
5   test <- t.test(grupo1[[i]], grupo2[[i]])
6   pvalues[i] <- test$p.value
7 }
8 sum(pvalues < 0.05) # Probablemente algunos falsos positivos!
9
10 # BUENO: Ajustar valores p
11 library(emmeans)
12 emm <- emmeans(modelo_anova, specs = "grupo")
13
14 # Método conservador
15 pairs(emm, adjust = "bonferroni")
16
17 # Método balanceado
18 pairs(emm, adjust = "holm")
19
20 # Control de False Discovery Rate
21 pairs(emm, adjust = "fdr")
22
23 # Para familias de hipótesis
24 p_adjust(pvalues, method = "BH") # Benjamini-Hochberg

```

11.3 Buenas prácticas específicas con easystats

11.3.1 Usar funciones consistentes

Listing 96: Aprovechar la consistencia de easystats

```

1 # La misma sintaxis funciona para múltiples tipos de modelos
2
3 # Regresión lineal
4 m1 <- lm(y ~ x, data = datos)
5
6 # Regresión logística
7 m2 <- glm(y ~ x, data = datos, family = binomial)
8
9 # Modelo mixto
10 library(lme4)
11 m3 <- lmer(y ~ x + (1 | grupo), data = datos)
12
13 # Para todos:
14 parameters(m1)
15 parameters(m2)
16 parameters(m3)
17
18 performance(m1)
19 performance(m2)
20 performance(m3)
21
22 check_model(m1)
23 check_model(m2)
24 check_model(m3)
25
26 report(m1)
27 report(m2)
28 report(m3)
```

11.3.2 Combinar paquetes del ecosistema

Listing 97: Flujo integrado completo

```

1 library(easystats) # Carga todos los paquetes
2
3 # 1. Preparar datos (datawizard)
4 datos <- datos_raw %>%
5   standardize() # Estandarizar variables
6
7 # 2. Ajustar modelo
8 modelo <- lm(y ~ x1 + x2, data = datos)
9
10 # 3. Parámetros (parameters)
11 params <- model_parameters(modelo)
```

```
12  
13 # 4. Desempeño (performance)  
14 perf <- model_performance(modelo)  
15  
16 # 5. Verificar supuestos (performance)  
17 check_model(modelo)  
18  
19 # 6. Tamaños de efecto (effectsize)  
20 efectos <- effectsize(modelo)  
21  
22 # 7. Estimaciones (modelbased)  
23 medias <- estimate_means(modelo, at = "x1")  
24  
25 # 8. Visualizar (see)  
26 plot(params)  
27 plot(medias)  
28  
29 # 9. Reportar (report)  
30 report(modelo)
```

11.4 Checklist de análisis

Antes de reportar resultados, verificar:

- ¿Los datos están limpios y bien preparados?
- ¿La pregunta de investigación está clara?
- ¿El modelo es apropiado para el tipo de datos?
- ¿Los supuestos del modelo se cumplen?
- ¿Se reportan intervalos de confianza, no solo valores p?
- ¿Se calculan y reportan tamaños de efecto?
- ¿Se consideran explicaciones alternativas?
- ¿El análisis es reproducible?
- ¿Se documentan limitaciones?
- ¿Las visualizaciones son claras e informativas?

11.5 Recomendaciones para adaptar a nuevos datasets

11.5.1 Lista de verificación al usar nuevos datos

1. Entender el contexto

- ¿Cuál es la fuente de los datos?

- ¿Cómo se recopilaron?
- ¿Cuál es la unidad de análisis?
- ¿Hay documentación o libro de códigos?

2. Exploración inicial

- Dimensiones (filas, columnas)
- Tipos de variables
- Valores faltantes
- Valores atípicos
- Distribuciones

3. Definir preguntas

- ¿Qué se quiere predecir o explicar?
- ¿Qué variables son relevantes?
- ¿Hay hipótesis a priori?

4. Seleccionar método apropiado

- Tipo de variable respuesta (continua, binaria, conteo, etc.)
- Estructura de datos (independientes, agrupados, temporales)
- Supuestos del método

5. Implementar con easystats

- Usar funciones consistentes
- Verificar supuestos
- Interpretar en contexto
- Validar resultados

11.5.2 Template de análisis reutilizable

Listing 98: Script template para nuevos análisis

```
1 # =====
2 # ANÁLISIS ESTADÍSTICO CON EASYSTATS
3 # Dataset: [NOMBRE]
4 # Fecha: [FECHA]
5 # Analista: [NOMBRE]
6 # =====
7
8 # 1. CONFIGURACIÓN -----
9 library(easystats)
10 library(tidyverse)
```

```
12 # Configurar opciones
13 options(scipen = 999, digits = 3)
14
15 # 2. CARGAR DATOS -----
16 datos_raw <- read.csv("ruta/al/archivo.csv")
17
18 # 3. EXPLORACIÓN INICIAL -----
19 glimpse(datos_raw)
20 summary(datos_raw)
21 vis_miss(datos_raw)
22
23 # 4. LIMPIEZA Y PREPARACIÓN -----
24 datos <- datos_raw %>%
25   # Filtros
26   filter(...) %>%
27   # Transformaciones
28   mutate(...) %>%
29   # Conversiones de tipo
30   mutate(across(c(...), as.factor))
31
32 # 5. ANÁLISIS DESCRIPTIVO -----
33 datos %>%
34   group_by(...) %>%
35   summarise(...)
36
37 # 6. MODELAMIENTO -----
38 modelo <- lm(y ~ x1 + x2, data = datos)
39
40 # 7. EVALUACIÓN -----
41 # Parámetros
42 parameters(modelo)
43
44 # Desempeño
45 performance(modelo)
46
47 # Supuestos
48 check_model(modelo)
49
50 # Tamaños de efecto
51 effectsize(modelo)
52
53 # 8. VISUALIZACIÓN -----
54 plot(parameters(modelo))
55 plot(estimate_means(modelo, at = "x1"))
56
57 # 9. REPORTE -----
58 report(modelo)
59
60 # 10. GUARDAR RESULTADOS -----
```

```
61 saveRDS(modelo, "resultados/modelo.rds")
62 ggsave("resultados/figuras/grafico.png", width = 10, height = 6, dpi = 300)
63
64 # 11. INFORMACIÓN DE SESIÓN -----
65 sessionInfo()
```

12 Conclusiones

12.1 Resumen del manual

Este manual ha cubierto el uso completo del ecosistema easystats para análisis estadístico aplicado, utilizando datos reales de la Encuesta de Intenciones de Siembra de la Región Cusco.

Conceptos clave aprendidos:

1. **Ecosistema easystats:** Una colección integrada de paquetes que facilita todo el flujo de análisis estadístico, desde el modelamiento hasta el reporte.
2. **Flujo de análisis sistemático:**

- Preparar datos (limpieza, transformación)
- Ajustar modelos apropiados
- Extraer parámetros (`parameters`)
- Evaluar desempeño (`performance`)
- Verificar supuestos (`check_model`)
- Calcular tamaños de efecto (`effectsize`)
- Visualizar (`see`)
- Reportar (`report`)

3. **Métodos estadísticos aplicados:**

- Regresión lineal (variables continuas)
- Regresión logística (variables binarias)
- Comparación de grupos (t-tests, ANOVA)
- Modelos con interacciones
- Modelos de efectos mixtos (datos jerárquicos)

4. **Buenas prácticas:**

- Verificar supuestos siempre
- Reportar tamaños de efecto, no solo p-valores
- Documentar y hacer análisis reproducibles
- Visualizar resultados efectivamente
- Interpretar en contexto del problema

12.2 Ventajas de easystats

Para principiantes:

- Sintaxis consistente y predecible
- Automatización de tareas complejas
- Reportes en lenguaje natural
- Visualizaciones profesionales con mínimo código
- Documentación extensa y ejemplos

Para usuarios avanzados:

- Flexibilidad y personalización
- Integración con otros paquetes de R
- Soporte para cientos de tipos de modelos
- Implementación de métodos estadísticos modernos
- Código abierto y extensible

12.3 Próximos pasos

Para continuar desarrollando habilidades en análisis estadístico con R y easystats:

12.3.1 Práctica

- Aplicar los métodos a sus propios datos
- Explorar otros datasets de Datos Abiertos del Gobierno del Perú
- Replicar análisis de artículos científicos
- Participar en comunidades de R (R-Ladies, useR!, etc.)

12.3.2 Profundización

- Modelos más avanzados (GAM, machine learning, etc.)
- Análisis bayesianos con `bayestestR`
- Series temporales y análisis espacial
- Meta-análisis y síntesis de evidencia

12.4 Recursos adicionales

12.4.1 Documentación oficial de easystats

- Sitio web: <https://easystats.github.io/easystats/>
- GitHub: <https://github.com/easystats>
- Viñetas de cada paquete (muy detalladas)
- Blog con tutoriales: <https://easystats.github.io/blog/>

12.4.2 Libros recomendados

- *R for Data Science* - Hadley Wickham & Garrett Grolemund
- *Statistical Rethinking* - Richard McElreath
- *An Introduction to Statistical Learning* - James et al.
- *Regression and Other Stories* - Gelman, Hill & Vehtari
- *Data Analysis Using Regression and Multilevel/Hierarchical Models* - Gelman & Hill

12.4.3 Cursos en línea

- DataCamp: Cursos de R y estadística
- Coursera: Especialización en Data Science (Johns Hopkins)
- edX: Estadística y R (Harvard)
- RStudio Education: <https://education.rstudio.com/>

12.4.4 Comunidades y foros

- RStudio Community: <https://community.rstudio.com/>
- Stack Overflow (etiqueta [r])
- Cross Validated (Stack Exchange para estadística)
- Twitter: #rstats, #easystats
- R-Ladies: <https://rladies.org/>

12.4.5 Datos abiertos en Perú

- Portal de Datos Abiertos: <https://www.datosabiertos.gob.pe/>
- INEI (Instituto Nacional de Estadística): <https://www.inei.gob.pe/>
- MIDAGRI (Ministerio de Desarrollo Agrario): <https://www.midagri.gob.pe/>
- SENAMHI (Servicio Nacional de Meteorología): <https://www.senamhi.gob.pe/>

12.5 Palabras finales

El análisis estadístico es una habilidad fundamental en la era de los datos. El ecosistema easystats democratiza el acceso a métodos estadísticos sofisticados, permitiendo que investigadores, estudiantes y profesionales realicen análisis rigurosos y reproducibles sin necesidad de ser expertos en programación o estadística avanzada.

Principios para recordar:

1. **Entender antes de aplicar:** Los métodos estadísticos son herramientas, no recetas. Comprender los supuestos y limitaciones es crucial.
2. **Simplicidad y parsimonia:** Comenzar con modelos simples y añadir complejidad solo cuando esté justificado.
3. **Transparencia y reproducibilidad:** Documentar todo el proceso analítico. La ciencia se basa en la replicabilidad.
4. **Significancia estadística ≠ importancia práctica:** Los valores p indican evidencia contra la hipótesis nula, pero los tamaños de efecto indican magnitud práctica.
5. **Comunicación clara:** Los resultados estadísticos deben ser comprensibles para la audiencia objetivo, no solo para estadísticos.
6. **Ética en análisis:** Reportar todos los análisis realizados, no solo los "significativos". Reconocer limitaciones y sesgos potenciales.
7. **Aprendizaje continuo:** La estadística y sus métodos evolucionan constantemente. Mantenerse actualizado es esencial.

Este manual es solo el comienzo de su viaje en el análisis estadístico con R y easystats. La práctica constante, la curiosidad y la disposición para aprender de los errores son las claves para desarrollar experticia. Los datos abiertos del Gobierno del Perú ofrecen innumerables oportunidades para aplicar estos conocimientos a problemas reales que afectan a millones de personas.

¡Buena suerte en sus análisis futuros!

13 Referencias

13.1 Publicaciones de easystats

- Lüdecke, D., Ben-Shachar, M.S., Patil, I., Waggoner, P., & Makowski, D. (2021). *performance: An R Package for Assessment, Comparison and Testing of Statistical Models*. Journal of Open Source Software, 6(60), 3139. <https://doi.org/10.21105/joss.03139>
- Makowski, D., Ben-Shachar, M.S., & Lüdecke, D. (2019). *bayestestR: Describing Effects and their Uncertainty, Existence and Significance within the Bayesian Framework*. Journal of Open Source Software, 4(40), 1541. <https://doi.org/10.21105/joss.01541>

- Ben-Shachar, M.S., Lüdecke, D., & Makowski, D. (2020). *effectsize: Estimation of Effect Size Indices and Standardized Parameters*. Journal of Open Source Software, 5(56), 2815. <https://doi.org/10.21105/joss.02815>
- Lüdecke, D., Waggoner, P., & Makowski, D. (2019). *insight: A Unified Interface to Access Information from Model Objects in R*. Journal of Open Source Software, 4(38), 1412. <https://doi.org/10.21105/joss.01412>

13.2 Libros y artículos de referencia

- Gelman, A., & Hill, J. (2006). *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge University Press.
- Field, A., Miles, J., & Field, Z. (2012). *Discovering Statistics Using R*. Sage Publications.
- Wickham, H., & Grolemund, G. (2017). *R for Data Science*. O'Reilly Media. <https://r4ds.had.co.nz/>
- Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences* (2nd ed.). Lawrence Erlbaum Associates.

13.3 Recursos en línea

- Easystats website: <https://easystats.github.io/easystats/>
- The R Graph Gallery: <https://r-graph-gallery.com/>
- STHDA (Statistical tools for high-throughput data analysis): <http://www.sthda.com/>
- Datos Abiertos del Gobierno del Perú: <https://www.datosabiertos.gob.pe/>

A Instalación y configuración

A.1 Instalación de R y RStudio

1. Descargar R desde CRAN: <https://cran.r-project.org/>
2. Descargar RStudio Desktop: <https://www.rstudio.com/products/rstudio/download/>
3. Instalar ambos programas siguiendo las instrucciones

A.2 Instalación de easystats

Listing 99: Instalación completa del ecosistema

```

1 # Instalar desde CRAN
2 install.packages("easystats")
3
4 # Alternativamente, instalar versión de desarrollo desde GitHub
5 # install.packages("remotes")
6 # remotes::install_github("easystats/easystats")
7
8 # Cargar el ecosistema
9 library(easystats)
10
11 # Verificar paquetes cargados
12 easystats::easystats_packages()

```

A.3 Paquetes complementarios recomendados

Listing 100: Otros paquetes útiles

```

1 # Manipulación de datos
2 install.packages("tidyverse")
3
4 # Visualización avanzada
5 install.packages("ggpubr")
6 install.packages("patchwork")
7
8 # Modelos mixtos
9 install.packages("lme4")
10 install.packages("lmerTest")
11
12 # Análisis de valores faltantes
13 install.packages("naniar")
14 install.packages("mice")
15
16 # Tablas
17 install.packages("kableExtra")
18 install.packages("flextable")
19
20 # Reportes
21 install.packages("rmarkdown")
22 install.packages("bookdown")

```

B Glosario de términos estadísticos

AIC (Akaike Information Criterion) Criterio de información para comparación de modelos. Menor es mejor.

ANOVA (Analysis of Variance) Análisis de varianza, prueba para comparar medias de múltiples grupos.

BIC (Bayesian Information Criterion) Similar a AIC pero con mayor penalización por complejidad.

BLUE (Best Linear Unbiased Estimator) Mejor estimador lineal insesgado.

BLUP (Best Linear Unbiased Predictor) Predictor lineal óptimo en modelos mixtos.

Coeficiente de determinación (R^2) Proporción de varianza explicada por el modelo.

d de Cohen Medida de tamaño de efecto para diferencias entre grupos.

Efecto aleatorio Variación entre grupos que se modela como una distribución.

Efecto fijo Parámetro del modelo que se estima directamente.

Efecto simple Efecto de un factor en un nivel específico de otro factor.

Eta cuadrado (η^2) Tamaño de efecto para ANOVA, proporción de varianza explicada.

Heteroscedasticidad Varianza no constante de los residuos.

Homoscedasticidad Varianza constante de los residuos (supuesto de regresión).

ICC (Intraclass Correlation) Proporción de varianza entre grupos en modelos mixtos.

Intervalo de confianza Rango de valores plausibles para un parámetro.

Multicolinealidad Correlación alta entre predictores.

Odds Ratio (OR) Razón de momios, medida de asociación en regresión logística.

Pseudo- R^2 Análogos al R^2 para modelos no lineales.

REML (Restricted Maximum Likelihood) Método de estimación para modelos mixtos.

RMSE (Root Mean Square Error) Raíz del error cuadrático medio.

Tamaño de efecto Magnitud estandarizada de un efecto estadístico.

Valor p Probabilidad de observar los datos (o más extremos) si la hipótesis nula es cierta.

VIF (Variance Inflation Factor) Medida de multicolinealidad.

C Atajos de teclado útiles en RStudio

Table 8: Atajos de teclado en RStudio

Atajo	Acción
Ctrl + Enter	Ejecutar línea/selección
Ctrl + Shift + M	Insertar pipe (%>%)
Ctrl + Shift + C	Comentar/descomentar
Alt + -	Insertar operador de asignación (-)
Ctrl + Shift + K	Compilar documento R Markdown
Ctrl + Shift + F10	Reiniciar sesión de R
Tab	Autocompletar código
Ctrl + 1	Ir a editor de código
Ctrl + 2	Ir a consola
Ctrl + L	Limpiar consola

— *Fin del Manual* —