

PLS-PLR

El modelo predictivo PLS-PLR de BLSD se construyó utilizando dataset de entrenamiento y se evaluó su desempeño utilizando un dataset de prueba. En primer lugar, se seleccionó el parámetro de penalización Ridge (λ), luego el modelo fue ajustado utilizando el parámetro de penalización que ofreció el mejor ajuste y finalmente fue validado usando validación cruzada y validación externa.

Algoritmo PLS-PLR

Sea \mathbf{Y} la variable respuesta binaria y $(\mathbf{X}_1, \dots, \mathbf{X}_p)$ un conjunto de predictores. Para una muestra de tamaño n los datos pueden ser organizados en un vector respuesta $\mathbf{y} = (y_1, \dots, y_n)^T$ y una matriz de predictores $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p) = (x_{ij})$ ($i = 1, \dots, n; j = 1, \dots, p$), donde y_i es 0 o 1 para presencia o ausencia de la característica principal y x_{ij} es el valor del i^{th} individuo en el j^{th} predictor. Las columnas de \mathbf{X} se suponen centradas y escaladas.

La ecuación de regresión de la variable dependiente $\mathbf{y} = \mathbf{T}\mathbf{c} + \mathbf{F}$ en el algoritmo PLS que explica \mathbf{y} desde las componentes \mathbf{T} es adaptada para tomar en cuenta la respuesta binaria. Bastien et al. (2005) generaliza el método cuando la respuesta desde una familia exponencial usando $g(\hat{\mathbf{y}}) = \mathbf{T}\mathbf{c}$ y propone un algoritmo para estimar los parámetros. Nosotros incluimos una constante en el modelo porque la variable binaria no puede ser centrada.

$$g(\hat{\mathbf{y}}) = c_0 + \mathbf{T}\mathbf{c} \quad (1)$$

El modelo Regresión Logística PLS es escrito como:

$$E(\mathbf{y}) = \hat{\mathbf{y}} = \frac{1}{\mathbf{1} + e^{-(c_0 \mathbf{1} + \sum_{h=1}^m c_h \mathbf{t}_h)}} \quad (2)$$

o

$$\text{logit}(\hat{\mathbf{y}}) = \log\left(\frac{\hat{\mathbf{y}}}{\mathbf{1} - \hat{\mathbf{y}}}\right) = c_0 \mathbf{1} + \sum_{h=1}^m c_h \mathbf{t}_h \quad (3)$$

O en forma matricial $\text{logit}(\hat{\mathbf{y}}) = c_0 \mathbf{1} + \mathbf{T}\mathbf{c}$, donde $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^T$ es el vector de probabilidades estimadas de la presencia en cada individuo y $\mathbf{c} = (c_1, \dots, c_m)^T$ son los coeficientes de la regresión sobre las componentes.

En términos de las variables originales,

$$\text{logit}(\hat{\mathbf{y}}) = c_0 \mathbf{1} + \mathbf{T}\mathbf{c} = c_0 \mathbf{1} + \mathbf{X}\mathbf{W}\mathbf{c} = c_0 \mathbf{1} + \mathbf{X}\mathbf{b} \quad (4)$$

Donde $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m)$ and $\mathbf{b} = (b_1, \dots, b_p)^T$ son los coeficientes sobre las variables observadas.

Para estimar $\mathbf{T}, \mathbf{W}, \mathbf{c}, c_0$ and \mathbf{b} nosotros usamos el algoritmo desarrollado por Bastien et al., (2005) con las modificaciones que se detallaron en los primeros párrafos de este apartado.

1. Calculo de \mathbf{t}_1 , la primera componente PLS.
 - a. Para cada predictor ($j = 1, \dots, p$), calcule el coeficiente de regresión w_{1j} de x_j , en la regresión logística de \mathbf{y} sobre \mathbf{x}_j , para obtener $\mathbf{w}_1 = (w_{11}, \dots, w_{1p})^T$
 - b. Normalice el vector $\mathbf{w}_1 := \mathbf{w}_1 / \|\mathbf{w}_1\|$.
 - c. Calcule las puntuaciones de la componente (scores) $\mathbf{t}_1 = \mathbf{X}\mathbf{w}_1 / \mathbf{w}_1^T \mathbf{w}_1$
2. Calculo de \mathbf{t}_h , la h^{th} componente PLS. Las componentes $\mathbf{t}_1, \dots, \mathbf{t}_{h-1}$ han sido obtenidas.

- a. Para cada predictor ($j = 1, \dots, p$), calcule el coeficiente de regresión w_{hj} de \mathbf{x}_j , en la regresión logística de \mathbf{y} sobre $\mathbf{t}_1, \dots, \mathbf{t}_{h-1}$ y \mathbf{x}_j , para obtener $\mathbf{w}_h = (w_{h1}, \dots, w_{hp})^T$
 - b. Normalice el vector $\mathbf{w}_h := \mathbf{w}_h / \|\mathbf{w}_h\|$.
 - c. Calcule la matriz residual \mathbf{X}_{h-1} de la regresión lineal de \mathbf{X} sobre $\mathbf{t}_1, \dots, \mathbf{t}_{h-1}$.
 - d. Calcule las puntuaciones de la componente (scores) $\mathbf{t}_h = \mathbf{X}\mathbf{w}_h / \mathbf{w}_h^T \mathbf{w}_h$.
3. \mathbf{X} es factorizado como $\mathbf{X} = \mathbf{TP}$
 4. Regresión Logística de \mathbf{y} sobre las componentes PLS retenidas

$$\text{logit}(\hat{y}) = c_0 \mathbf{1} + \sum_{h=1}^m c_h \mathbf{t}_h$$

5. Expresión del modelo en términos de las predictoras originales $\mathbf{b} = \mathbf{W}$.

$$\text{logit}(\hat{\mathbf{y}}) = c_0 \mathbf{1} + \mathbf{T}\mathbf{c} = c_0 \mathbf{1} + \mathbf{X}\mathbf{W}\mathbf{c} = c_0 \mathbf{1} + \mathbf{X}\mathbf{b}$$

Penalización Ridge

La función log-verosimilitud con penalización Ridge es,

$$L_{\text{ridge}}(\boldsymbol{\beta}) = L(\boldsymbol{\beta}) - \lambda \|\boldsymbol{\beta}\|^2 \quad (5)$$

El estimador Ridge iterativo logístico $\hat{\boldsymbol{\beta}}_{\text{Ridge}}$ se lo obtiene usando el método de Newton-Raphson.

(6)

$$\hat{\boldsymbol{\beta}}_{\text{Ridge}} = \hat{\boldsymbol{\beta}}_{\text{Ridge}} + \{ \mathbf{X}^T \hat{\mathbf{V}}(\hat{\boldsymbol{\beta}}_{\text{Ridge}}) \mathbf{X} + 2\lambda \mathbf{I} \}^{-1} \{ \mathbf{U}(\hat{\boldsymbol{\beta}}_{\text{Ridge}}) - 2\lambda \hat{\boldsymbol{\beta}}_{\text{Ridge}} \}$$

Donde,

$U(\hat{\beta}_{Ridge}) = X^T[y - \pi(\hat{\beta}_{Ridge})]$ es el vector de primeras derivadas parciales de $L(\beta)$

y

$$\hat{V} = diag\{\hat{\pi}_1(1 - \hat{\pi}_1), \dots, \hat{\pi}_n(1 - \hat{\pi}_n)\}.$$

Bondad de ajuste del modelo

Para evaluar la bondad de ajuste del modelo se calculó la diferencia de devianza y los Pseudo R^2 de McFadden, Cox&Snell, Nagelkerke.

Diferencia de Devianza

$$D_0 - D_M = \left(-2 \ln \frac{L_0(\beta)}{L_F(\beta)}\right) - \left(-2 \ln \frac{L_M(\beta)}{L_F(\beta)}\right) \quad (7)$$

$$D_0 - D_M = -2 \left(\ln \frac{L_0(\beta)}{L_F(\beta)} - \ln \frac{L_M(\beta)}{L_F(\beta)} \right) \quad (8)$$

$$D_0 - D_M = -2 \ln \left(\frac{L_0(\beta)}{L_M(\beta)} \right) \quad (9)$$

$$DiffDeviance = -(2LL_0(\beta) - 2LL_M(\beta)) \quad (10)$$

Donde:

$LL_M(\beta)$ Log verosimilitud del modelo.

$LL_0(\beta)$ Log verosimilitud del modelo nulo (null model)

La diferencia de devianza (*DiffDeviance*) es interpretada como una medida de la variación de los datos explicada por el modelo con predictores, pero sin constantes (modelo nulo). Este estadístico tiene una distribución χ^2 ($\chi^2_{s_M - s_0}$), con grados de

libertad igual a la diferencia entre los números de parámetros de los modelos. De esta manera, la hipótesis nula será rechazada para el nivel de significancia α cuando $DiffDeviance > \chi^2$, lo que es equivalente a que el valor_p de contraste sea menor que el nivel de α fijado (Hosmer et al., 1998).

Pseudo R^2

$R^2_{McFadden}$

$$R^2_{McFadden} = 1 - \left(\frac{LL_M}{LL_0} \right) \quad (11)$$

Donde:

LL_M es log verosimilitud del modelo

LL_0 es log verosimilitud de modelo nulo (null model)

$R^2_{Cox\&Snell}$

$$R^2_{Cox\&Snell} = 1 - \left(\frac{L_0}{L_M} \right)^{\frac{2}{n}} \quad (12)$$

Donde

L_M es el valor de verosimilitud del modelo.

L_0 es el valor de verosimilitud del modelo nulo (null model).

$R^2_{Nagelkerke}$

$$R^2_{Nagelkerke} = \frac{R^2_{Cox\&Snell}}{1 - (L_0)^{\frac{2}{n}}} \quad (13)$$

Donde,

L_M es el valor de verosimilitud del modelo.

L_0 es el valor de verosimilitud del modelo nulo (null model).

Selección de parámetro de penalización Ridge (λ)

Tabla 1 Métricas de bondad de ajuste para modelos PLS-PLR con valores incrementales de penalización Ridge (λ)

λ	<i>Diff-Deviance</i> ^a	<i>R</i> ² <i>CoxSnell</i> ^b	<i>R</i> ² <i>Nagelkerke</i> ^c	<i>R</i> ² <i>MacFadden</i> ^d
0.1	88.488	0.573	0.994	0.991
0.2	88.005	0.571	0.991	0.986
0.3	87.668	0.570	0.988	0.982
0.4	87.405	0.568	0.986	0.979
0.5	87.187	0.568	0.985	0.976
0.6	86.999	0.567	0.984	0.974
0.7	86.832	0.566	0.982	0.972
0.8	86.682	0.565	0.981	0.971
0.9	86.543	0.565	0.980	0.969

^a difference of Deviance (Hosmer et al.,1998).

^b *R*²*Cox&Snell*, ^c *R*²*Nagelkerke* and ^d *R*²*MacFadden* son índices pseudo *R*² para modelos de regresión logística binaria (Allison, 2014; Walker & Smith, 2016).

Con un valor de λ igual a 0.1, se obtuvieron las mejores medidas de bondad de ajuste.

El modelo generado con las dos primeras componentes PLS está representado en la siguiente función:

$$P_y = \frac{e^{(27.226+1.546t_1+1.318t_2)}}{1 + e^{(27.226+1.546t_1+1.318t_2)}}$$

Donde,

P_y es la probabilidad de la presencia de la enfermedad

t_1 es la primera componente PLS

t_2 es la segunda componente PLS

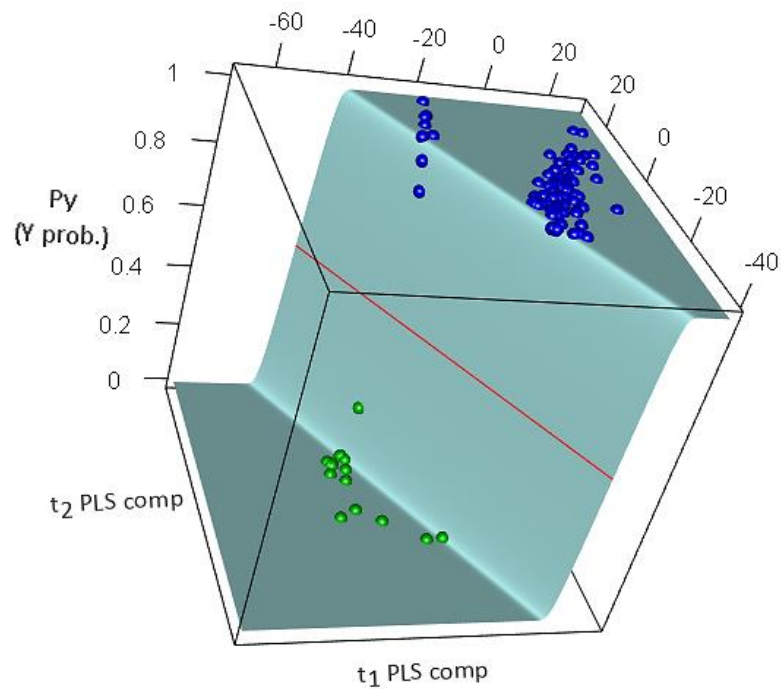


Figura 1 Gráfico 3D de la respuesta del modelo PLS-PLR.

Predicción y validación del modelo PLS-PLR

La capacidad predictiva del modelo fue valorada utilizando el método de validación cruzada LOOCV (Leave-One-Out-Cross-Validation).

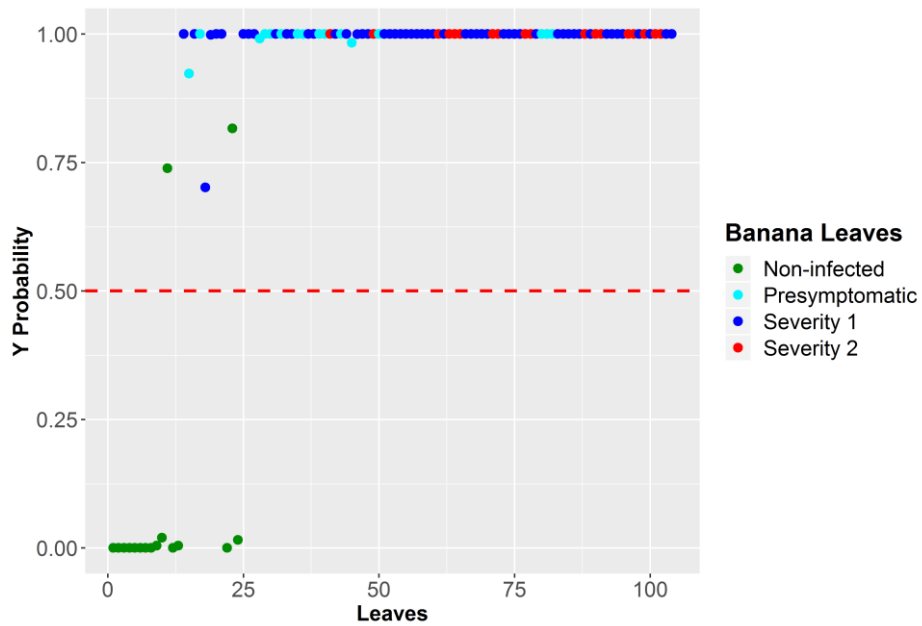


Figura 2 Probabilidad estimada por PLS-PLR con validación cruzada.

Matriz de confusión

Tabla 2 Matriz de confusión del modelo PLS-PLR con validación cruzada.

	Hojas Infectadas	Hojas no-infectadas	
Resultado Test	TP 88	FP 2	Precisión 0.98
	FN 0	TN 14	Valor Pred. Negativo 1
	Sensibilidad 1	Especificidad 0.88	Exactitud 0.98

Validación Externa del modelo PLS-PLR

El modelo PLS-PLR ajustado al conjunto de datos de entrenamiento se usó para predecir la presencia de la enfermedad en nuevas hojas y evaluar la eficacia del modelo.

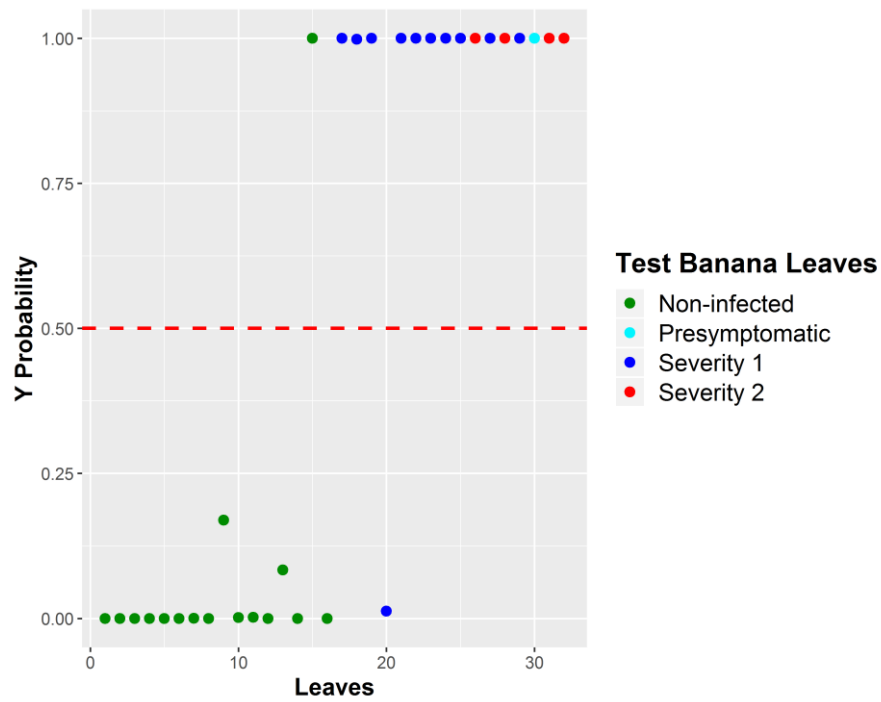


Figura 3 Probabilidad estimada por el modelo PLS-PLR en prueba de validación.

Tabla **¡Error! No hay texto con el estilo especificado en el documento.** Matriz de confusión del modelo PLS-PLR en prueba de validación

	Hojas Infectadas	Hojas no-infectadas	
Resultado Test	TP	FP	Precisión
	15	1	0.94
	FN	TN	Valor Pred. Negativo
	1	15	0.94
	Sensibilidad	Especificidad	Exactitud
	0.94	0.94	0.94

Área bajo la curva ROC (AUC): para calcular el área bajo la curva ROC se utilizó la función *auc()* del lenguaje R, el resultado fue 0.94.

HS-BILOT

Hyperspectral Biplot (HS-Biplot) es una representación gráfica de las hojas, las longitudes de onda y las regiones de predicción y permite una inspección visual de las relaciones entre ellos. Las longitudes de onda fueron representadas por líneas coloreadas de acuerdo con la banda espectral a la que pertenecen.

Longitudes de onda de las regiones visible y cercana al infrarrojo

Espectro visible	
Color	Long. De onda
Violet	380 - 427 nm
Blue	427 - 476 nm
Cyan	476 - 497 nm
Green	497 - 570 nm
Yellow	570 - 581 nm
Orange	581 - 618 nm
Red	618 - 780 nm
Infrarrojo cercano	
Color	Long. De onda
Gray	780 - 1350 nm

Las dos primeras componentes del modelo PLS-PLR fueron utilizadas para graficar el HS-Biplot. Las puntuaciones filas ***T*** (scores) y las puntuaciones columna ***P*** (loadings) fueron proyectadas en un plano que tiene como ejes las componentes del modelo PLS-PLR. Las puntuaciones filas se muestran como puntos y representan las hojas de banano. Las puntuaciones columnas proveen la dirección de las líneas que representan las longitudes de onda.

$$X \cong TP + E$$

El porcentaje de variabilidad capturada por la aproximación es,

$$\rho^2 = \frac{tr(\hat{X}^T \hat{X})}{tr(X^T X)} \times 100 \quad (14)$$

Donde \hat{X} es la matriz de predictores aproximada y X es la matriz original.

Es posible identificar las variables relacionadas a las componentes PLS calculando bondad de ajuste por columna

$$\rho_j^2 = \frac{tr(\hat{x}_{[j]}^T \hat{x}_{[j]})}{tr(x_{[j]}^T x_{[j]})} \times 100 \quad (j = 1, \dots, p) \quad (15)$$

Donde $\hat{x}_{[j]}$ y $x_{[j]}$ son el j^{th} columnas de la matriz ajustada y de la matriz original respectivamente.

La bondad de ajuste de cada fila es,

$$\rho_i^2 = \frac{tr(\hat{x}_{[i]} \hat{x}_{[i]}^T)}{tr(x_{[i]} x_{[i]}^T)} \times 100 \quad (i = 1, \dots, n) \quad (16)$$

Donde $\hat{x}_{[i]}$ y $x_{[i]}$ son el i^{th} filas de la matriz ajustada y de la matriz original respectivamente.

Estas medidas también se llaman calidad de la representación o predictividad.

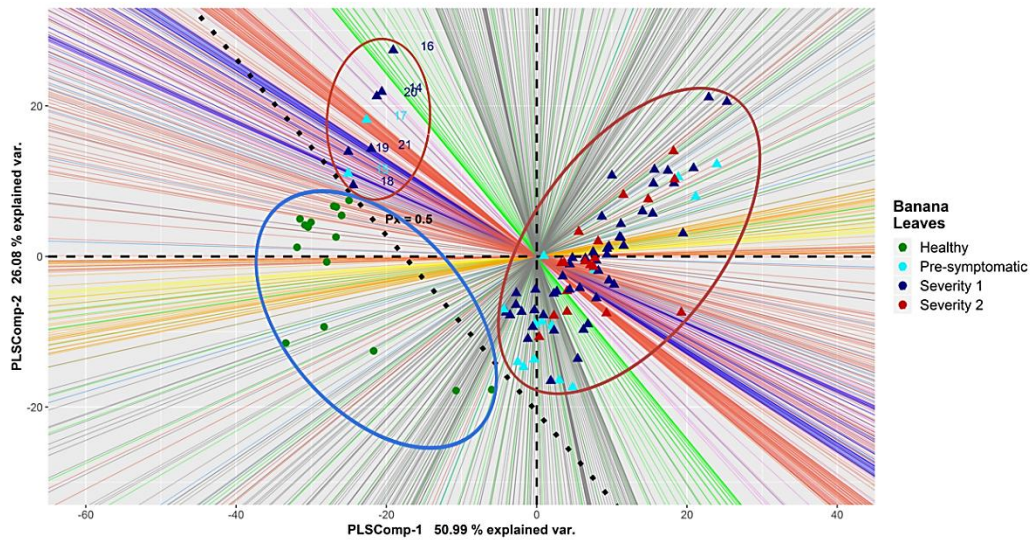


Figura 4 HS-Biplot del dataset de entrenamiento.

Las filas de la matriz reducida (hojas de banano) están representadas por puntos y las columnas (longitudes de onda) están representadas por líneas rectas.

El espacio cartesiano está dividido por una línea puntuada oblicua formando dos regiones separadas que predicen la presencia o ausencia de enfermedad. La línea corta el plano en el valor de predicción de 0.5 y corresponde al umbral de clasificación, lo que implica que los puntos ubicados sobre dicha línea son hojas clasificadas como infectadas y, por el contrario, si están abajo de la línea son hojas no infectadas.

Se observó exactamente el mismo comportamiento en el conjunto de datos de validación externa.

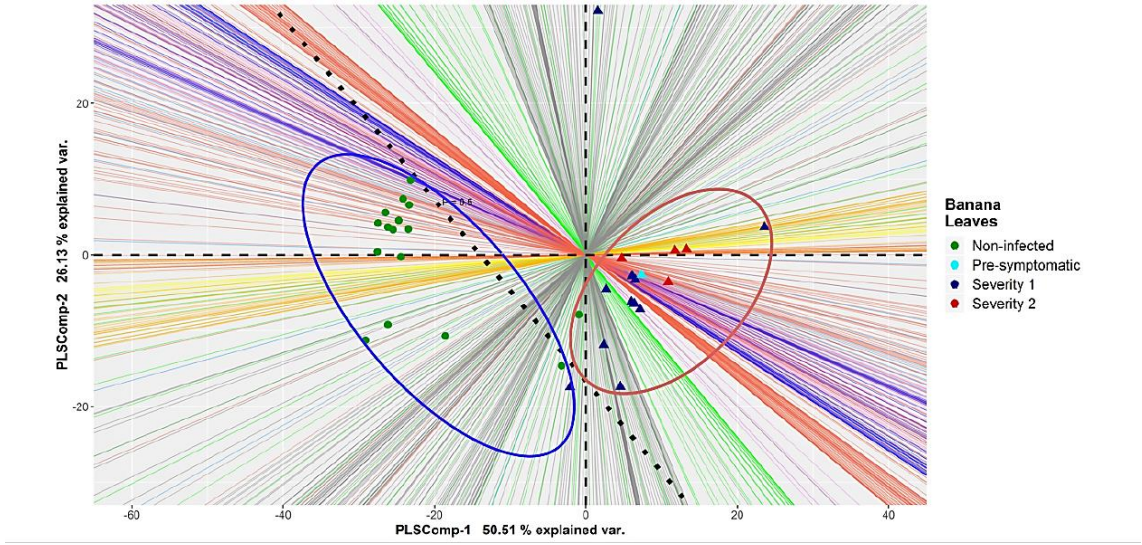


Figura 5 HS-Biplot del dataset de validación.

La bondad de ajuste global del HS-Biplot (el coeficiente de correlación al cuadrado entre los valores ajustados y observados) fue del 77.07%.

$$\rho^2 = \frac{\text{tr}(\hat{\mathbf{X}}^T \hat{\mathbf{X}})}{\text{tr}(\mathbf{X}^T \mathbf{X})} \times 100 = 77.07$$

Donde

$\hat{\mathbf{X}} = \mathbf{TP}^T$ es la matriz de datos aproximada.

\mathbf{X} es la matriz de datos original

La contribución de las componentes a cada variable

$$\rho_j^2 = \frac{\text{tr}(\hat{\mathbf{x}}_{[j]}^T \hat{\mathbf{x}}_{[j]})}{\text{tr}(\mathbf{x}_{[j]}^T \mathbf{x}_{[j]})} \times 100 \quad (j = 1, \dots, p)$$

ρ_j^2 es la contribución de cada columna j

$\hat{\mathbf{x}}_{[j]}$ la columna j de la matriz aproximada $\hat{\mathbf{X}}$

$\mathbf{x}_{[j]}$ es la columna j de la matriz original \mathbf{X}

La contribución de las componentes a cada fila o individuo

$$\rho_i^2 = \frac{\text{tr}(\hat{\mathbf{x}}_{[i]}\hat{\mathbf{x}}_{[i]}^T)}{\text{tr}(\mathbf{x}_{[i]}\mathbf{x}_{[i]}^T)} \times 100 \quad (i = 1, \dots, n)$$

ρ_i^2 es la contribución de cada columna i

$\hat{\mathbf{x}}_{[i]}$ es la fila i de la matriz aproximada $\hat{\mathbf{X}}$

$\mathbf{x}_{[i]}$ es la fila i de la matriz original \mathbf{X}