

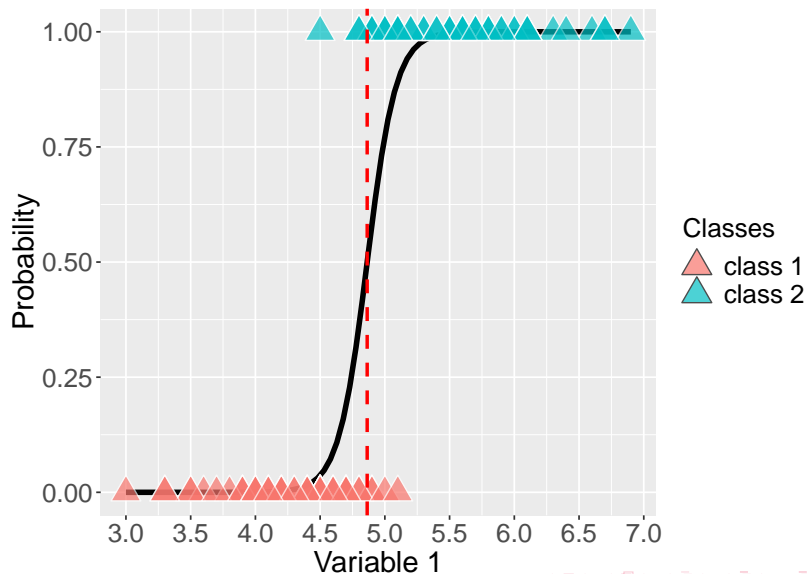
A Dive into Logistic Regression Modeling

Crista Moreno

June 7, 2019

- 1 Goal: Build a Good Model
- 2 Structure of Biomedical Data
- 3 Which Variables to Include in the Model?
- 4 Logistic Regression
- 5 Overfitting the Model
- 6 Cross Validation
- 7 Software

Goal Build a Good Logistic Regression Model



Structure of Biomedical Data

Data

Patient	Variable 1	Variable 2	...	Variable M	Class
1	$m_{1,1}$	$m_{1,2}$...	$m_{1,M}$	normal
2	$m_{2,1}$	$m_{2,2}$...	$m_{2,M}$	sick
3	$m_{3,1}$	$m_{3,2}$...	$m_{3,M}$	normal
4	$m_{4,1}$	$m_{4,2}$...	$m_{4,M}$	normal
\vdots	\vdots	\vdots	...	\vdots	\vdots
N	$m_{N,1}$	$m_{N,2}$...	$m_{N,M}$	sick

Data

Patient	Variable 1	Variable 2	...	Variable M	Class
1	$m_{1,1}$	$m_{1,2}$...	$m_{1,M}$	normal
2	$m_{2,1}$	$m_{2,2}$...	$m_{2,M}$	sick
3	$m_{3,1}$	$m_{3,2}$...	$m_{3,M}$	normal
4	$m_{4,1}$	$m_{4,2}$...	$m_{4,M}$	normal
\vdots	\vdots	\vdots	...	\vdots	\vdots
N	$m_{N,1}$	$m_{N,2}$...	$m_{N,M}$	sick

Important Numbers

$m_{3,4}$ - 4th variable measurement for the 3rd patient

Structure of Biomedical Data

Data

Patient	Variable 1	Variable 2	...	Variable M	Class
1	$m_{1,1}$	$m_{1,2}$...	$m_{1,M}$	normal
2	$m_{2,1}$	$m_{2,2}$...	$m_{2,M}$	sick
3	$m_{3,1}$	$m_{3,2}$...	$m_{3,M}$	normal
4	$m_{4,1}$	$m_{4,2}$...	$m_{4,M}$	normal
\vdots	\vdots	\vdots	...	\vdots	\vdots
N	$m_{N,1}$	$m_{N,2}$...	$m_{N,M}$	sick

Important Numbers

$m_{3,4}$ - 4th variable measurement for the 3rd patient

N - Total Number of Data Points (number of rows, patients etc.)

Structure of Biomedical Data

Data

Patient	Variable 1	Variable 2	...	Variable M	Class
1	$m_{1,1}$	$m_{1,2}$...	$m_{1,M}$	normal
2	$m_{2,1}$	$m_{2,2}$...	$m_{2,M}$	sick
3	$m_{3,1}$	$m_{3,2}$...	$m_{3,M}$	normal
4	$m_{4,1}$	$m_{4,2}$...	$m_{4,M}$	normal
\vdots	\vdots	\vdots	...	\vdots	\vdots
N	$m_{N,1}$	$m_{N,2}$...	$m_{N,M}$	sick

Important Numbers

$m_{3,4}$ - 4th variable measurement for the 3rd patient

N - Total Number of Data Points (number of rows, patients etc.)

M - Total Number of Variables (measurements, parameters etc.)

Structure of Biomedical Data

Data

Patient	Variable 1	Variable 2	...	Variable M	Class
1	$m_{1,1}$	$m_{1,2}$...	$m_{1,M}$	normal
2	$m_{2,1}$	$m_{2,2}$...	$m_{2,M}$	sick
3	$m_{3,1}$	$m_{3,2}$...	$m_{3,M}$	normal
4	$m_{4,1}$	$m_{4,2}$...	$m_{4,M}$	normal
\vdots	\vdots	\vdots	...	\vdots	\vdots
N	$m_{N,1}$	$m_{N,2}$...	$m_{N,M}$	sick

Important Numbers

$m_{3,4}$ - 4th variable measurement for the 3rd patient

N - Total Number of Data Points (number of rows, patients etc.)

M - Total Number of Variables (measurements, parameters etc.)

C - Total Number of Classes (labels i.e. normal vs diseased)

Which Variables to Include in the Model?

- $N = |\{X^1, X^2, \dots, X^N\}|$ (data)
- $C = |\{c_1, c_2\}| = 2$ e.g. {normal, diseased}
- $M = |\{v_1, v_2, \dots, v_M\}|$ e.g. {sex, weight, blood type, etc.}

Which Variables to Include in the Model?

- $N = |\{X^1, X^2, \dots, X^N\}|$ (data)
- $C = |\{c_1, c_2\}| = 2$ e.g. {normal, diseased}
- $M = |\{v_1, v_2, \dots, v_M\}|$ e.g. {sex, weight, blood type, etc.}

Question 1

What is the **probability** that patient belongs to class c , given that the data X is equal to x ?

$$P(Y = c | X = x)$$

Which Variables to Include in the Model?

- $N = |\{X^1, X^2, \dots, X^N\}|$ (data)
- $C = |\{c_1, c_2\}| = 2$ e.g. {normal, diseased}
- $M = |\{v_1, v_2, \dots, v_M\}|$ e.g. {sex, weight, blood type, etc.}

Question 1

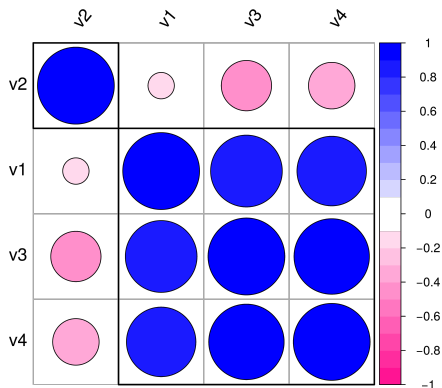
What is the **probability** that patient belongs to class c , given that the data X is equal to x ?

$$P(Y = c | X = x)$$

Question 2

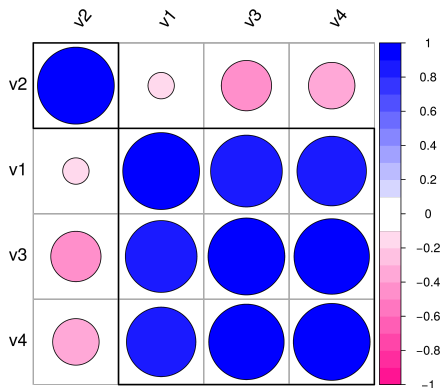
How to determine which of the M variables to use for the model?

Answer Question 2: Correlation Matrix



$$\text{Cor}(x, y) = \frac{\langle x^*, y^* \rangle}{\|x^*\| \cdot \|y^*\|}, \quad x^* = x - \bar{x}, y^* = y - \bar{y}$$

Answer Question 2: Correlation Matrix



$$\text{Cor}(x, y) = \frac{\langle x^*, y^* \rangle}{\|x^*\| \cdot \|y^*\|}, \quad x^* = x - \bar{x}, y^* = y - \bar{y}$$

v1 and v2 look interesting

Answer Question 1: Logistic Regression Model

$$p(Y = c|X = x)$$

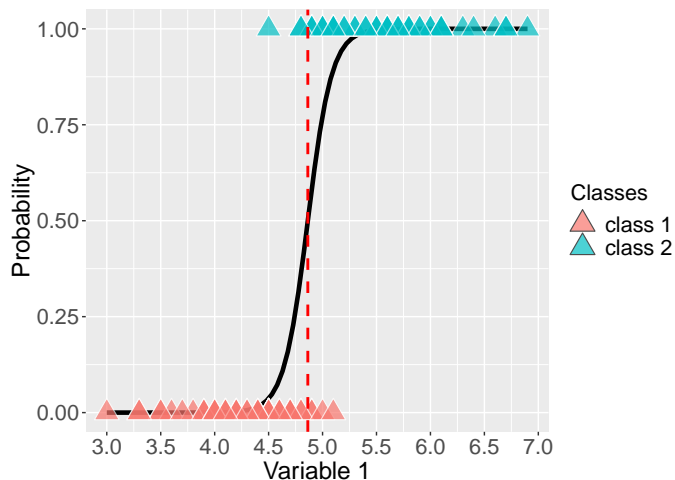
Answer Question 1: Logistic Regression Model

$$\begin{aligned} p(Y = c|X = x) \\ = p(X) = \frac{e^{\beta \cdot \mathbf{x}}}{1 + e^{\beta \cdot \mathbf{x}}} \end{aligned}$$

Answer Question 1: Logistic Regression Model

$$\begin{aligned} p(Y = c|X = x) \\ &= p(X) = \frac{e^{\beta \cdot \mathbf{x}}}{1 + e^{\beta \cdot \mathbf{x}}} \\ p(-\infty) &= 0, \quad p(+\infty) = 1 \end{aligned}$$

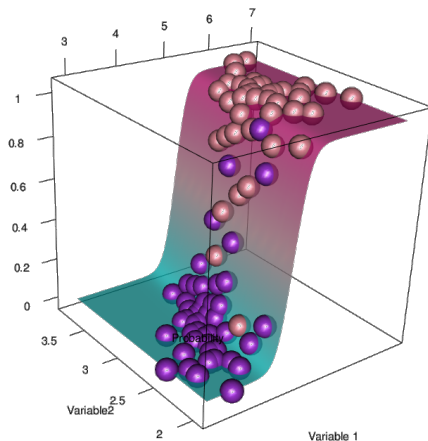
Logistic Regression Model with a Single Variable



$$p(X) = \frac{e^{\beta_0 + \beta_1 X_1}}{1 + e^{\beta_0 + \beta_1 X_1}}$$

Red dashed line is the decision boundary.

Logistic Regression Model with Multiple Variables



$$p(\mathbf{X}) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_\gamma X_\gamma}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_\gamma X_\gamma}} = \frac{e^{\beta \cdot \mathbf{X}}}{1 + e^{\beta \cdot \mathbf{X}}}$$

Great! We have a model!
Are we done?

Great! We have a model!
Are we done?
Not so fast.

- Adding too many variables \Rightarrow Overfitting & \uparrow Accuracy

Dangers with Overfitting the Model

- Adding too many variables \Rightarrow Overfitting & \uparrow Accuracy
- \uparrow Accuracy \Rightarrow \uparrow Sensitivity and \uparrow Specificity

Dangers with Overfitting the Model

- Adding too many variables \Rightarrow Overfitting & \uparrow Accuracy
- \uparrow Accuracy \Rightarrow \uparrow Sensitivity and \uparrow Specificity
- \uparrow Sensitivity and \uparrow Specificity \Rightarrow Misleading Model

Dangers with Overfitting the Model

- Adding too many variables \Rightarrow Overfitting & \uparrow Accuracy
- \uparrow Accuracy \Rightarrow \uparrow Sensitivity and \uparrow Specificity
- \uparrow Sensitivity and \uparrow Specificity \Rightarrow Misleading Model

The simplest model that fits the data is also the most plausible. (Occam's Razor)

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}, \quad \text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad \text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

How do we validate our model?

Validating a Model: Back to Data

Patient	Variable 1	Variable 2	...	Variable M	Class
1	$m_{1,1}$	$m_{1,2}$...	$m_{1,M}$	normal
2	$m_{2,1}$	$m_{2,2}$...	$m_{2,M}$	sick
3	$m_{3,1}$	$m_{3,2}$...	$m_{3,M}$	normal
4	$m_{4,1}$	$m_{4,2}$...	$m_{4,M}$	normal
⋮	⋮	⋮	...	⋮	⋮
N	$m_{N,1}$	$m_{N,2}$...	$m_{N,M}$	sick

Validating a Model: Back to Data

Patient	Variable 1	Variable 2	...	Variable M	Class
1	$m_{1,1}$	$m_{1,2}$...	$m_{1,M}$	normal
2	$m_{2,1}$	$m_{2,2}$...	$m_{2,M}$	sick
3	$m_{3,1}$	$m_{3,2}$...	$m_{3,M}$	normal
4	$m_{4,1}$	$m_{4,2}$...	$m_{4,M}$	normal
⋮	⋮	⋮	...	⋮	⋮
N	$m_{N,1}$	$m_{N,2}$...	$m_{N,M}$	sick

Data \mathcal{D} for Model

Data	Observation	Input (Variables)	Output (Class)
\mathcal{D}_1	1	X^1	Y^1
\mathcal{D}_2	2	X^2	Y^2
\mathcal{D}_3	3	X^3	Y^3
\mathcal{D}_4	4	X^4	Y^4
⋮	⋮	⋮	⋮
\mathcal{D}_N	N	X^N	Y^N

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \dots, \mathcal{D}_N\}$$

- Given the data for modeling

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}\}$$

Classic Approach: Partition Data for Validating the Model

- Given the data for modeling

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}\}$$

- Partition the data \mathcal{D} into a **Training set** and **Test set**

$$\underbrace{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7}_{\text{Training Set}}, \underbrace{\mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}}_{\text{Test Set}}$$

Classic Approach: Partition Data for Validating the Model

- Given the data for modeling

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}\}$$

- Partition the data \mathcal{D} into a **Training set** and **Test set**

$$\underbrace{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7}_{\text{Training Set}}, \underbrace{\mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}}_{\text{Test Set}}$$

- Use **Training set** to build a model

Classic Approach: Partition Data for Validating the Model

- Given the data for modeling

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}\}$$

- Partition the data \mathcal{D} into a **Training set** and **Test set**

$$\underbrace{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7}_{\text{Training Set}}, \underbrace{\mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}}_{\text{Test Set}}$$

- Use **Training set** to build a model
- Test model on **Test set** to get the model's performance.

Classic Approach: Partition Data for Validating the Model

- Given the data for modeling

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}\}$$

- Partition the data \mathcal{D} into a **Training set** and **Test set**

$$\underbrace{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7}_{\text{Training Set}}, \underbrace{\mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}}_{\text{Test Set}}$$

- Use **Training set** to build a model
- Test model on **Test set** to get the model's performance.

Now we have the model's performance, are we done?

Classic Approach: Partition Data for Validating the Model

- Given the data for modeling

$$\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}\}$$

- Partition the data \mathcal{D} into a **Training set** and **Test set**

$$\underbrace{\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7}_{\text{Training Set}}, \underbrace{\mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}}_{\text{Test Set}}$$

- Use **Training set** to build a model
- Test model on **Test set** to get the model's performance.

Now we have the model's performance, are we done?

Classic approach works well for large data sets.

k-fold Cross Validation Works for Smaller Data Sets

k-fold Cross Validation Works for Smaller Data Sets

- Select k e.g. $k = 5$ folds

k-fold Cross Validation Works for Smaller Data Sets

- Select k e.g. $k = 5$ folds
- Partition data into a **Training Set** and **Validation Set** in the following fashion:

k-fold Cross Validation Works for Smaller Data Sets

- Select k e.g. $k = 5$ folds
- Partition data into a **Training Set** and **Validation Set** in the following fashion:

fold 1 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

fold 2 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

fold 3 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

fold 4 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

fold 5 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

k-fold Cross Validation Works for Smaller Data Sets

- Select k e.g. $k = 5$ folds
- Partition data into a **Training Set** and **Validation Set** in the following fashion:

fold 1 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

fold 2 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

fold 3 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

fold 4 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

fold 5 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

- For each fold train model on the **Training Set** and get performance on **Validation Set**

k-fold Cross Validation Works for Smaller Data Sets

- Select k e.g. $k = 5$ folds
- Partition data into a **Training Set** and **Validation Set** in the following fashion:

fold 1 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$
fold 2 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$
fold 3 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$
fold 4 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$
fold 5 $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3, \mathcal{D}_4, \mathcal{D}_5, \mathcal{D}_6, \mathcal{D}_7, \mathcal{D}_8, \mathcal{D}_9, \mathcal{D}_{10}$

- For each fold train model on the **Training Set** and get performance on **Validation Set**
- The variance and average of the performance helps indicate how well this model can make predictions on future data.

Software

- **R** - A **free software** (GNU Affero GPL) environment for statistical computing and graphics.
- **RStudio** is a free and open-source integrated development environment (IDE) for R.

