

Министерство образования Республики Беларусь

Учреждение образования

БЕЛОРУССКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ

ИНФОРМАТИКИ И РАДИОЭЛЕКТРОНИКИ

Факультет Информационных технологий и управления

Кафедра Интеллектуальных информационных технологий

Индивидуальная практическая работа №1

По дисциплине статические основы индукционного вывода

на тему

Построение бинарного классификатора средствами MS EXCEL

Выполнил Мулярчик Д.С.

Проверил Ефремов А.А

Минск 2023

Задание:

1. Подобрать в открытых источниках data set, состоящий из результативного признака (заданного бинарной переменной) и нескольких факторных признаков (не менее 3).
2. Построить бинарный классификатор, пользуясь методическими указаниями из примера ниже.
3. В отчёте представить: постановку задачи с описанием переменных (А), фрагмент таблицы с исходными данными (Б), уравнение логистической регрессии (В), значение Zгр (Г), оценку надёжности классификатора через расчёт процента ошибок (Д).

Оценка наличия ухудшения ментального здоровья:

В качестве dataseta использовался <https://www.kaggle.com/datasets/csafr12/maternal-health-risk-data/>

1. Увеличение риска
2. Нормальная частота пульса в состоянии покоя в ударах в минуту. HeartRate - x1
3. Верхнее значение артериального давления в мм рт. ст. SystolicBP - x2
4. Нижнее значение артериального давления в мм рт.ст. DiastolicBP - x3
5. Уровень глюкозы в крови выражен в молярной концентрации, ммоль/л. BS - x4

Коэффициенты потенциального x1,x2,x3,x4 человека являются следующими 102.5, 78.7, 10.6, 70.9.

Требуется:

- 1) построить линейную регрессионную модель для оценки кредитного риска и с ее помощью оценить вероятность дефолта для потенциального заемщика;
- 2) построить регрессионную дискриминантную модель, найти граничное значение $Z_{гр}$ и отнести потенциального заемщика к группе с высоким либо низким кредитным риском.

Информация по кредитам приведена в следующей таблице.

Age	SystolicBP	DiastolicBP	BS	HeartRate	RiskLevel
35	85	60	11	86	1
42	130	80	18	70	1
23	90	60	7	76	0
25	110	89	7	77	0
15	120	80	7	70	0
50	140	90	15	90	1
25	140	100	7	80	1
10	70	50	6	70	0
40	140	100	18	90	1

Пункт 1

Введем таблицу с данными в Excel. Линейная регрессионная модель для оценки кредитного риска в данном случае имеет вид: $Z = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + b_4x_4 + e$.

Для оценки коэффициентов β_k , $k=0,4$, будем использовать модуль «Анализ данных», который вызывается из «Сервиса» в главном меню. В «Анализе данных» найдем инструмент «Регрессия» и вызовем его. В появившемся окне укажем входные интервалы Y и X.

Входной интервал Y – это массив ячеек (в таблице исходных данных), содержащих значения объясняемой переменной Z. Входной интервал X – это массив ячеек, содержащих значения объясняющих переменных x_1, x_2, x_3, x_4 .

После ввода входных интервалов, нажмем на кнопку «ОК». В результате появится новый лист с параметрами регрессионной модели. Оценка β_0 коэффициента β_0 равна значению коэффициента для «Y-пересечения», а оценки $\beta_1, \beta_2, \beta_3, \beta_4$ коэффициентов $\beta_1, \beta_2, \beta_3, \beta_4$ равны значениям коэффициентов для переменных X_1, X_2, X_3, X_4 .

Вероятность ухудшения ментального здоровья оценивается по формуле $\hat{\beta}_2 x_2 + \hat{\beta}_3 x_3 + \hat{\beta}_4 x_4$, где x_1, x_2, x_3, x_4 заданные значения для потенциального пациента.

В результате вычислений мы получили следующие значения:

1. $\hat{\beta}_0 = 14.0374285$.
2. $\hat{\beta}_1 = 0.4658391$.
3. $\hat{\beta}_2 = -2.0573952$.
4. $\hat{\beta}_3 = -0.7652419$.
5. $\hat{\beta}_4 = 0.9274294$.

Пункт 2

В качестве регрессионной дискриминантной модели можно взять модель из п.1.

Для каждого наблюдения вычисляются прогнозные значения показателя z по формуле:

$$\hat{Z}^i = \hat{\beta}_0 + \hat{\beta}_1 x_1^i + \hat{\beta}_2 x_2^i + \hat{\beta}_3 x_3^i = \overline{1, N}$$

Затем с помощью функций СРЗНАЧ и СТАНДОТКЛОН нужно найти средние значения z_1 и z_2 , и стандартные отклонения σ_1 и σ_2 для наблюдений с дефолтом (1-й массив) и для наблюдений без дефолта. (Для этого предварительно следует упорядочить таблицу соответствующим образом.)

1. $z_1 = 3.5723852$
2. $z_2 = 0.2473618$
3. $\sigma_1 = 2.5063951$
4. $\sigma_2 = 1.8294629$

$$Z_{гр} = \frac{\sigma_1 \overline{Z_2} + \sigma_2 \overline{Z_1}}{\sigma_1 + \sigma_2}$$

Граничное значение $z_{гр}$ вычисляется по формуле: равное
1.6498662053056.

Поскольку $z_1 > z_2$, вероятность заболевания оценивается как низкая, если $z < z_{гр}$, и высокая если $z > z_{гр}$.

Для потенциального пациента вычислим $Z_{пот} = 1.442$.

Следовательно, у данного пациента низкая вероятность ухудшения ментального здоровья.