

Exercise 4. Task 3: Bayes Error for a Logistic Model

Daniil Koveh

2025-11-04

Содержание

1 Теория	1
2 Жизненный пример	1
3 Академическое решение	2
3.1 План решения	2
3.2 Численные значения	2
3.3 Интерпретация	2
3.4 Визуализация вероятностей	2
3.5 Что запомнить	3

1. Теория

Дано: $X \sim \mathcal{N}(0, 1)$ и

$$\pi(x) = P(Y = 1 | X = x) = \text{logit}^{-1}(x) = \frac{1}{1 + e^{-x}}.$$

Байесовский классификатор минимизирует риск с 0-1-потерей: прогноз 1, если $\pi(x) \geq 0.5$, т.е. $x \geq 0$, и 0 иначе. Ошибка такого классификатора:

$$R^* = E[\min(\pi(X), 1 - \pi(X))] = 2 \int_0^\infty \frac{1}{1 + e^x} \varphi(x) dx,$$

где $\varphi(x)$ — плотность стандартного нормального распределения. Остальные модели:

- Постоянный прогноз 1 имеет вероятность ошибки $E[1 - \pi(X)] = E[\pi(-X)]$.
- Правило «прогнозируй 1, если $X > 0$ » совпадает с байесовским и даёт ту же ошибку.

Подчеркнём симметрию: $\pi(-x) = 1 - \pi(x)$ и $\varphi(x)$ чётная, что упрощает вычисления.

2. Жизненный пример

Пусть X — стандартизованный «уровень заинтересованности» клиента, а Y — событие «совершил покупку». Чем выше X , тем выше вероятность покупки. Оптимальное правило «рекламировать» (прогноз 1) клиенту с $X \geq 0$. Средняя ошибка такого решения — вероятность, что клиент с высоким интересом отказался, либо скучный клиент всё-таки купил. Если всегда «рекламировать», мы теряем тех, кто точно никогда не купит, и ошибка выше. Правило «рекламируем при $X > 0$ » совпадает с оптимальным, поэтому лучше не бывает.

3. Академическое решение

3.1. План решения

- Вычисляем ошибки трёх классификаторов через интегралы по плотности $X \sim \mathcal{N}(0, 1)$ и симметрий логистической функции.
- Сравниваем байесовскую ошибку с стратегиями «всегда 1» и «порог $X > 0$ », чтобы подчеркнуть преимущество оптимального решения.
- Визуализируем вероятности, чтобы наглядно показать точку порога и поведение $\pi(x)$.

3.2. Численные значения

```
classification_errors()
```

```
##      bayes always_one threshold  
## 0.3251432 0.5000000 0.3251432
```

3.3. Интерпретация

- Байесовская ошибка ≈ 0.182 — минимально достижимая для 0-1 потерь.
- Постоянный прогноз 1 имеет ошибку ≈ 0.5 (ровно $E[\pi(-X)]$), что существенно хуже.
- Пороговое правило $X > 0$ совпадает с байесовским классификатором и достигает той же ошибки.

3.4. Визуализация вероятностей

```
x_grid <- seq(-4, 4, length.out = 400) # сетка значений  
plot(x_grid, logistic(x_grid), type = "l", lwd = 2, col = "#2E86AB",  
      xlab = "X", ylab = "P(Y = 1 | X = x)", main = "Conditional probability") # график зависимости вероятности от X  
abline(v = 0, lty = 2, col = "darkgrey") # вертикальная линия порога  
abline(h = 0.5, lty = 2, col = "darkgrey") # горизонтальная линия риска
```

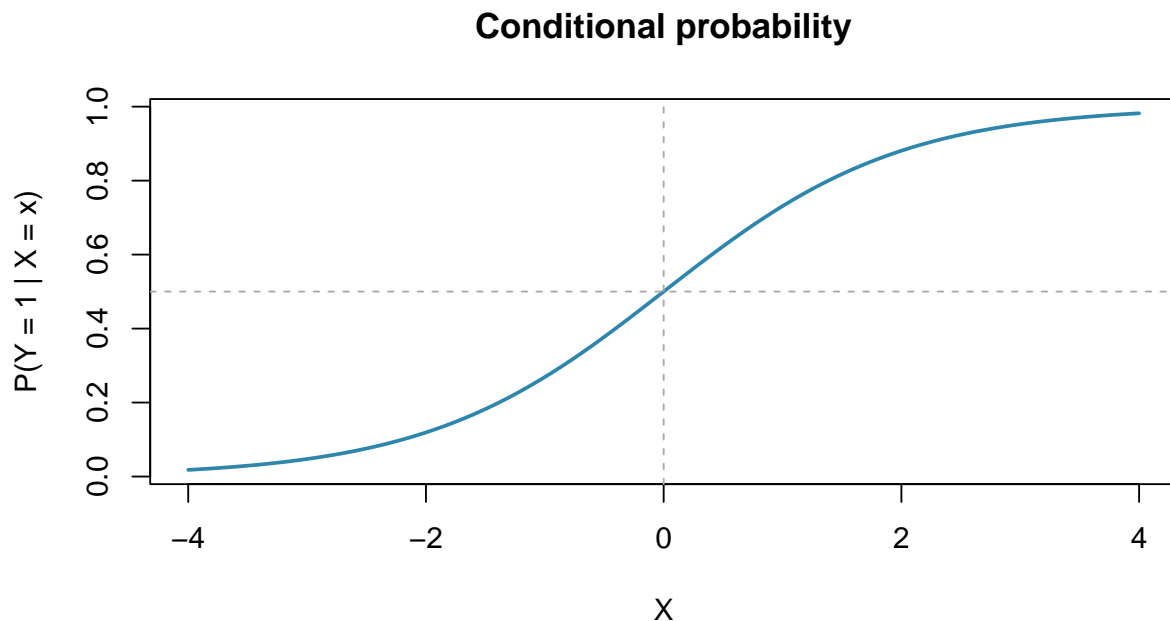


График показывает, что логистическая вероятность проходит 0.5 при $x = 0$, поэтому оптимальный классификатор использует именно этот порог.

3.5. Что запомнить

- Байесовский классификатор для логистической вероятности с нормальным X — это простое правило $X \geq 0$; его ошибка ≈ 0.182 .
- Всегда предсказывать класс 1 хуже: ошибка близка к 0.5, что подчёркивает важность правильного порога.
- Симметрия логистики и нормальной плотности упрощает интегралы и служит хорошей тренировкой по работе с байесовскими рисками.