Пошаговая дискриминация, кросс-валидация и бутстрап в задаче классификации пострадавших с сочетанной травмой груди.

Феофанов Василий

Санкт-Петербургский государственный университет, факультет прикладной математики - процессов управления, кафедра МТИСР

5 апреля, 2016

Оглавление

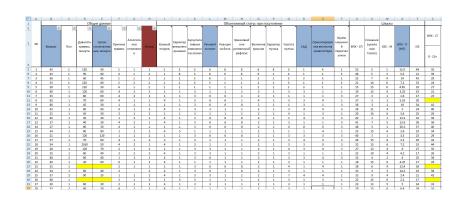
- Постановка задачи
- 2 Линейный дискриминатный анализ
- Пошаговая дискриминация
- Оценка вероятности ошибочной классификации
- \delta Классификация пострадавших с сочетанной травмой груди
- Выводы
- 🕡 Список литературы



Постановка задачи

В исследовании рассматриваются 52 пациента с сочетанной травмой грудной клетки, госпитализированные в экстренном порядке. База содержит общую информацию о пострадавших, результаты лабораторных и инструментальных исследований, проведенных в течение первых 12 часов с момента поступления в больницу пострадавшего, а также исход получения травмы. По имеющимся 160 признакам, описывающих каждого пациента, необходимо построить классификационное правило, которое позволит спрогнозировать летальность исхода для будущих пациентов с сочетанной травмой груди.

Фрагмент базы данных



Линейный дискриминатный анализ (ЛДА)

Один из классических методов классификации. Если μ_1, μ_2 — математические ожидания классов W_1 и W_2 соответственно, а Σ — общая ковариационная матрица, тогда для произвольного наблюдения x классификационное правило будет выглядеть следующим образом:

$$W_1: (\mu_1 - \mu_2)^T \Sigma^{-1} x \geqslant \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 + \mu_2) + \ln\left(\frac{p_2}{p_1}\right)$$

$$W_2: (\mu_1 - \mu_2)^T \Sigma^{-1} x < \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 + \mu_2) + \ln \left(\frac{p_2}{p_1}\right),$$

где p_1, p_2 — априорные вероятности принадлежности к классам W_1 и W_2 .

Проблемы применения ЛДА

- наблюдений меньше, чем оцениваемых параметров.
- риск получения плохо обусловленной ковариационной матрицы
- ошибка на обучении: $\overline{err} = \frac{1}{n} \sum_{i=1}^n |y_i m(x_i)|$ не является корректной оценкой вероятности ошибочной классификации

Пошаговая дискриминация

В основе пошагового дискриминатного анализа лежит тест на добавочную информацию, который определяет значимость вклада новых включенных переменных по отношению к старым при проверки гипотезы $H_0: \mu_1 = \mu_2$. Для этого вычисляется частная лямбда Уилкса:

$$\Lambda(x|t) = \frac{\Lambda(t,x)}{\Lambda(t)},$$

где t — изначальный набор признаков, x — набор добавленных переменных, а $\Lambda(z)$:

$$\Lambda = \frac{|W|}{|W+B|},$$

где B и W меж- и внутри- групповые ковариационные матрицы.

Алгоритм Forward Selection

- 1. Вначале из модели удаляются все рассматриваемые признаки, число которых p.
- 2. Задаемся значением Л-включения.
- 3. Для каждого x_j высчитывается $\Lambda(x_j)$ и затем включается в модель та переменная, значение соответствующей статистики которой наименьшее среди рассматриваемых. Включенный в модель признак обозначим за t_1 .

Алгоритм Forward Selection

4. Среди остальных p-1 переменных ищется признак с наименьшей частной лямбдой Вилкса:

$$\Lambda(x_j|t_1) = \frac{\Lambda(t_1, x_j)}{t_1},$$

Вместе с тем, статистика должна удовлетворять условию: $\Lambda \leqslant \Lambda$ -включения.

5. Процесс продолжается аналогичным образом до тех пор, пока ни одна из переменных не будет удовлетворять условию или все переменные не войдут в модель.

Оценка вероятности ошибочной классификации

1. Cross-validation leave-one-out:

$$\widehat{Err}^{(CV)} = \frac{1}{n} \sum_{i=1}^{n} |y^{(i)} - m^{-(i)}(x^{(i)})|$$

2. Bootstrap leave-one-out:

$$\widehat{Err}^{(LOOB)} = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} |y^{(i)} - m^{(b)}(x^{(i)})|,$$

где C^{-i} — набор индексов, идентифицирущие те бустрап-выборки, которые не содержат объект i.

Оценка вероятности ошибочной классификации

3. Bootstrap 0.632:

$$\widehat{Err}^{(0.632)} = 0.368 \cdot \overline{err} + 0.632 \cdot \widehat{Err}^{(LOOB)}$$

4. Bootstrap 0.632+:

$$\widehat{Err}^{(0.632+)} = (1 - \alpha) \cdot \overline{err} + \alpha \cdot \widehat{Err}^{(LOOB)},$$

где $\alpha = \frac{0.632}{1-0.368\widehat{R}}$, а \widehat{R} — относительная частота переобучения.

Классификация пострадавших с сочетанной травмой груди. Этап I.

Все 160 признаков были разбиты на следующие группы:

Группа 1: Общие данные

Группа 2: Объективный статус при поступлении + шкалы

Группа 3: Шкалы

Группа 4: Структура повреждений внутренних органов

Группа 5: Параметры ИВ Π

 Γ руппа 6: BCP + BCAД + BДАД

Группа 7: ВД + общие данные спироартериокардиоритмографа (САКР)

Группа 8: Электрокардиография (ЭКГ)

Группа 9: Общий анализ крови (ОАК)

Группа 10: Биохимия

Группа 11: Маркеры повреждения сердца + газы крови

Результаты пошагового дискриминантного анализа

| Группа | Признак 1 | Λ | Признак 2 | Λ | Признак 3 | Λ | Признак 4 | Λ | Признак 5 | Λ |
|-----------|------------------------|------|-------------------------|------|-------------------------|------|------------------------|------|--------------------------------------|------|
| Группа 1 | Возраст | 0.6 | Сроки госпитализации | 0.73 | | | | | | |
| Группа 2 | ВПХ - П (МТ) | 0.59 | Речевой контакт | 0.55 | Величина кровопотери | 0.51 | САД при поступлении | 0.44 | Частота пульса При поступлении | 0.4 |
| Группа 3 | ВПХ - голова | 0.76 | ВПХ - грудь | 0.63 | ВПХ - таз | 0.53 | AIS - грудь | 0.45 | AIS - таз | 0.41 |
| Группа 4 | Повреждение ЦНС САК | 0.85 | | | | | | | | |
| Группа 5 | ивл адд | 0.92 | | | | | | | | |
| Группа 6 | BCP LF (n. u.) | 0.8 | Вар. САД ТР | 0.65 | Вар. САД НЕ | 0.57 | Bap. САД LF (n. u.) | 0.5 | Bap. САД VLF | 0.43 |
| Группа 7 | CAKP PQ | 0.84 | CAKP 4CCcp | 0.73 | САКР АДСмакс | 0.57 | САКР АДДср | 0.46 | САКР Вар. Дых. НF n. | 0.36 |
| Группа 8 | ЭКГ RR | 0.92 | | | | ' | | | | |
| Группа 9 | ОАК Гемоглобин | 0.77 | | | | | | | | |
| Группа 10 | Биохимия Натрий | 0.63 | | | | | | | | |
| Группа 11 | Газы крови FiO2 | 0.79 | Газы крови RI | 0.67 | | | | | | |

Оценка вероятности ошибочной классификации.

| Группа | Ошибка | Cross- | Bootstrap | Bootstrap | Bootstrap 0.632+ | |
|-----------|-------------|-----------------------------|---------------|-----------|---------------------|--|
| | на обучении | validation Leave-one-out | Leave-one-out | 0.632 | | |
| | | | | | | |
| Группа 1 | 0.289 | 0.333 | 0.399 | 0.358 | 0.367 | |
| Группа 2 | 0.118 | 0.157 | 0.177 | 0.155 | 0.157 | |
| Группа 3 | 0.178 | 0.222 | 0.242 | 0.219 | 0.221 | |
| Группа 4 | 0.28 | 0.28 | 0.305 | 0.296 | 0.296 | |
| Группа 5 | 0.306 | 0.347 | 0.345 | 0.331 | 0.332 | |
| Группа 6 | 0.114 | 0.114 | 0.164 | 0.146 | 0.147 | |
| Группа 7 | 0.074 | 0.185 | 0.222 | 0.168 | 0.179 | |
| Группа 8 | 0.304 | 0.326 | 0.375 | 0.349 | 0.353 | |
| Группа 9 | 0.362 | 0.426 | 0.385 | 0.376 | 0.377 | |
| Группа 10 | 0.267 | 0.267 | 0.291 | 0.282 | 0.282 | |
| Группа 11 | 0.25 | 0.296 | 0.309 | 0.287 | 0.289 | |
| | | | | | | |

Классификация пострадавших с сочетанной травмой груди. Этап II.

Группа А: Отобранные признаки из групп 1-3 Группа В: Отобранные признаки из групп 4-6 Группа С: Отобранные признаки из групп 7-11 Группа F: Отобранные признаки из групп А-С

| Группа | Признак 1 | Λ | Признак 2 | Λ | Признак 3 | Λ | Признак 4 | Λ | Признак 5 | Λ | Признак 6 | Λ | Признак 7 | Λ |
|----------|--------------------|------|--------------------|------|------------------------|------|-------------------------|------|-------------------------|------|-------------------------|------|------------------------|------|
| Группа А | ВПХ - П (MT) | 0.61 | Речевой контакт | 0.53 | Возраст | 0.45 | Величина кровопотери | 0.39 | САД при поступлении | 0.33 | Сроки госпитализации | 0.3 | | |
| Группа В | BCPLF (n. u.) | 0.83 | Вар. САД ТР | 0.67 | Вар. САД НЕ | 0.58 | Bap. САД LF (n. u.) | 0.51 | ИВЛАДД | 0.42 | Bap. CAД VLF | 0.38 | | |
| Группа С | Газы крови FiO2 | 0.85 | Газы крови RI | 0.77 | | | | | | | | | | |
| Группа F | ВПХ - П (MT) | 0.63 | Возраст | 0.53 | Bap. САД LF (n. u.) | 0.44 | Речевой контакт | 0.4 | Величина кровопотери | 0.34 | Сроки госпитализации | 0.3 | САД при поступлении | 0.26 |

| Группа | Ошибка | Cross- | Bootstrap | Bootstrap | Bootstrap 0.632+ | |
|----------|-------------|---------------|---------------|-----------|---------------------|--|
| | на обучении | validation | Leave-one-out | 0.632 | | |
| | | Leave-one-out | | | | |
| Группа А | 0.044 | 0.044 | 0.099 | 0.079 | 0.08 | |
| Группа В | 0.116 | 0.163 | 0.161 | 0.144 | 0.145 | |
| Группа С | 0.25 | 0.296 | 0.312 | 0.289 | 0.291 | |
| Группа F | 0.075 | 0.1 | 0.124 | 0.106 | 0.107 | |
| | | | | | | |

Выводы

В задаче прогнозирования летальности исхода постравдшего с сочетанной травмой груди стоит принимать во внимание такие параметры, как:

- Шкала военной-полевой хирургии $\Pi(MT)$ (Π -повреждение, MT-механическая травма)
- Систолическое артериальное давление при поступлении в больницу
- Величина кровопотери и сроки госпитализации
- Вариабельность сердечного ритма (мощность низкочастотного компонента)
- Вариабельность систолического артериального давления
- Диастолическое артериальное давление
- Речевой контакт



Выводы

- Наилучшая точность была достигнута на группе признаков А, но не на группе F. Это подтверждает, что пошаговый дискриминантный анализ не выдает гарантированно оптимальный набор признаков
- Несмотря на наличие большого числа признаков, превышающее количество наблюдений, удалось добиться высокой точности классификации (90.1–95.6%)
- Ошибка на обучении в среднем давала оптимистически заниженную оценку в сравнении с методами кросс-валидации и бутстрап
- 0.632 и 0.632+ показали примерно одинаковые результаты, что может говорить о том, что в исследовании наблюдался незначимый эффект переобучения



Список литературы

- [1] Буре В. М., Парилина Е. М., Рубша А. И., Свиркина Л. А. Анализ выживаемости по медицинской базе данных больных раком предстательной железы // Вестн. С.-Петерб. ун-та. Сер. 10: Прикладная математика, информатика и процессы управления. 2014. N. 2. C. 27–35.
- [2] Буре В. М., Щербакова А. А. Применение дискриминантного анализа и метода деревьев принятия решений для диагностики офтальмологических заболеваний // Вестн. С.-Петерб. ун-та. Сер. 10: Прикладная математика, информатика и процессы управления. 2013. N. 1. С. 70–76.

- [3] Rencher A. C. Methods of Multivariate Analysis. 2nd Ed. New York: John Wiley & Sons, Inc. 2002. P. 738.
- [4] Рао С. Р. Линейные статистические методы и их применения. / науч ред. Линник Ю. В.; пер. с англ. Калинина В. М. и др. М.: Наука, 1968. 548 с.
- [5] Fu W. J., Carroll R. J., Wang S. Estimating misclassification error with small samples via bootstrap cross-validation. // Bioinformatics. 2005. Vol. 21(9). P. 1979–1986.
- [6] Hastie T., Tibshirani R., Friedman J. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. 2nd Ed. New York: Springer-Verlag. 2009. P. 745.

- [7] Molinaro A. M., Simon R., Pfeiffer R. M. Prediction error estimation: a comparison of resampling methods. // Bioinformatics. 2005. Vol. 21(15). P. 3301–3307.
- [8] Zavorka S., Perrett J. J. Minimum sample size considerations for two-group linear and quadratic discriminant analysis with rare populations // Communications in Statistics - Simulation and Computation. 2014. Vol. 43(7). P. 1726–1739.