

Предсказание качества для процедуры выбора признаков*

Аминов Т.В.

aminov.tv@phystech.edu

¹Московский Физико-технический институт

В случае избыточного признакового пространства предсказательная модель машинного обучения является неустойчивой. Для повышения устойчивости модели применяются методы выбора признаков. Задача выбора признаков состоит в поиске решения среди экспоненциального числа возможных решений, поэтому она является NP-задачей. Для решения применяются эвристические, субоптимальные алгоритмы. В данной статье предлагается свести дискретную задачу выбора признаков к задаче непрерывной оптимизации. Строится процедуры предсказания качества модели на тестовой выборке для подмножества признаков. Решение задачи выбора признаков восстанавливается из решения непрерывной задачи. Проводится вычислительный эксперимент, чтобы сравнить результаты предложенного метода с существующими алгоритмами.

Ключевые слова: Выбор признаков, многомерные пространства, отображение булева куба.

1 Введение

При решении многих задач машинного обучения, таких как классификация и регрессия, исходное пространство признаков оказывается избыточным. Не все признаки коррелируют с целевой переменной, нерелевантные признаки делают модель неустойчивой. Данная работа посвящена проблеме выбора оптимального подмножества признаков. Эта задача является NP-задачей, так как точное решение ищется среди всех возможных $2^n - 1$ вариантов, где n - количество признаков. Поэтому предполагается отыскание приближенного решения.

Задача выбора признаков возникает при распознавании лиц [1]. В статье [2] данная проблема в пространствах малой размерности (порядка 1000). За счет огромного

*Работа выполнена при финансовой поддержке РФФИ, проект №00-00-00000. Научный руководитель: Стрижов В. В. Задачу поставил: Роман Р. В. Консультант: Роман Р. В.

прироста вычислительной мощности за последние 20 лет задача выбора признаков может решаться точно за разумное время. В статье [3] представлена идея построения оптимальной архитектуры нейронной сети. Существующие методы, независимо от того, основаны ли они на обучении с подкреплением [4] или на эволюционных алгоритмах (ЕА) [5], выполняют поиск архитектуры в дискретном пространстве, который крайне неэффективен.

Авторы предлагают простой и эффективный метод автоматического выбора признаков непрерывной оптимизации. Весь алгоритм разбивается на три этапа: (1) энкодер - отображает пространство подмножества признаков в линейное пространство; (2) предиктор - принимает непрерывное представление подмножества признаков в качестве входных данных и прогнозирует ошибку на тестовой выборке; (3) декодер - отображает непрерывное представление обратно в набор признаков. Существенным отличием нашего метода от уже существующих состоит в том что размерность булева куба признаков очень велика, а значит для прироста скорости нам придется отображать его в пространство гораздо меньшей размерности.

В качестве показателя эффективности предложенного алгоритма предполагается сравнение с уже существующими: CSO (PSO) [6], CFS [7], QPFS [8], на реальных и синтетических данных.

2 Постановка задачи

Пусть $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$ заданная матрица, где $\mathbf{x}_j \in \mathbb{R}^m$ j -ый признак. Пусть $\mathbf{y} \in \mathbb{R}^m$. Множество $\mathcal{A} \subseteq \{1, \dots, n\}$ признаков. Существует соответствие между множеством \mathcal{A} и двоичным вектором $\mathbf{a} \in \mathbb{B}^n$:

$$\mathcal{A} = \{j : a_j = 1\}.$$

Функция $f(\mathbf{x}, \mathbf{w}, \mathcal{A})$ предсказывает y по заданному объекту \mathbf{x} и используя только признаки из множества \mathcal{A} . Данные разбиваются на две части: train ($\mathbf{X}_{\text{tr}}, \mathbf{y}_{\text{tr}}$) и test ($\mathbf{X}_{\text{te}}, \mathbf{y}_{\text{te}}$). Чтобы измерить качество модели f , вводится функция ошибки $s(\mathbf{w}, \mathbf{X}, \mathbf{y}, \mathcal{A})$. Для обучения подбираются наилучшие параметры \mathbf{w}^* , исходя из условия:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} s(\mathbf{w}, \mathbf{X}_{\text{tr}}, \mathbf{y}_{\text{tr}}, \mathcal{A}).$$

58 Целью является нахождение такого множества \mathcal{A} на котором функция ошибки на те-
 59 стовой выборке будет минимальна. Предлагается следовать следующему методу. Со-
 60 здается набор данных $\{(\mathbf{a}_i, s_i)\}_{i=1}^N$. Выбираются N векторов \mathbf{a} из двоичного куба \mathbb{B}^n .
 61 Затем оценивается ошибка s_i для каждого набора \mathbf{a}_i . Идея в том чтобы отобра-
 62 зить дискретную область \mathbb{B}^n в непрерывную с меньшей размерностью $h < n$. Та-
 63 ким образом вектор признаков будет обладать непрерывным представлением $\mathbf{u} \in \mathbb{R}^h$.
 64 Это представление создается энкодером $e(\mathbf{a}, \mathbf{w}_e)$. После этого модель $p(\mathbf{u}, \mathbf{w}_p)$ пыта-
 65 ется предсказать ошибку на тестовой выборке s . Тогда модель принимает вид $s =$
 66 $= p(\mathbf{u}, \mathbf{w}_p) = p(e(\mathbf{a}, \mathbf{w}_e), \mathbf{w}_p)$. Возможно использование любой функции потерь для
 67 оценки параметров $\mathbf{w}_e, \mathbf{w}_p$. Например это может быть квадратичная функция ошибки

$$68 \quad L_{error}(\mathbf{w}_e, \mathbf{w}_p, \mathbf{a}, s) = \|p(e(\mathbf{a}, \mathbf{w}_e), \mathbf{w}_p) - s\|^2 \rightarrow \min_{\mathbf{w}_e, \mathbf{w}_p}.$$

69 Так же требуется построить модель обратного отображения вектора \mathbf{a} по его непре-
 70 рывному представлению \mathbf{u} . Эта модель будет называться декодер $\sigma(\mathbf{u}, \mathbf{w}_\sigma)$. Вводится
 71 функция ошибки отображения - кросс-энтропия между начальным вектором \mathbf{a} и тем
 72 что получается на выходе модели декодер σ

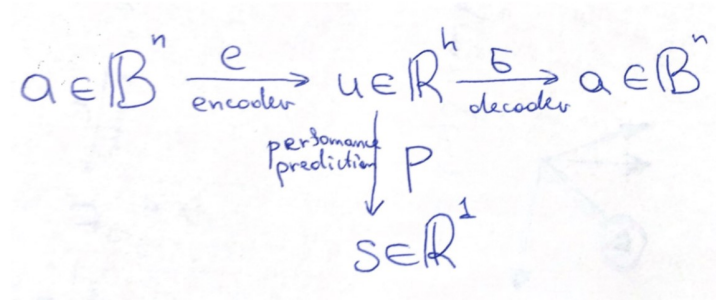
$$73 \quad L_{rec}(\mathbf{w}_e, \mathbf{w}_\sigma, \mathbf{a}) = \sum_{i=1}^n a_i \log(\sigma_i) + (1 - a_i) \log(1 - \sigma_i) \rightarrow \min_{\mathbf{w}_e, \mathbf{w}_\sigma},$$

74 где $\sigma_i = \sigma(e(\mathbf{a}, \mathbf{w}_e), \mathbf{w}_\sigma)_i$. Общая функция потерь будет линейной комбинацией ошибки
 75 предсказания и ошибки отображения

$$76 \quad L = L_{error} + \alpha L_{rec}.$$

77 Совокупностью параметров $\mathbf{w}_e, \mathbf{w}_p, \mathbf{w}_\sigma$ задается наиболее подходящее подмножество
 78 признаков. Максимизируется предсказательное качество модели

$$79 \quad \mathbf{u}^* = \arg \max_{\mathbf{u}} p(\mathbf{u}, \mathbf{w}_p).$$



Литература

- [1] Dong Chen, Xudong Cao, Fang Wen, and Jian Sun. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. pages 3025–3032, 06 2013.
- [2] Dimitris Bertsimas, Angela King, and Rahul Mazumder. Best subset selection via a modern optimization lens. *Ann. Statist.*, 44(2):813–852, 04 2016.
- [3] Shengcao Cao, Xiaofang Wang, and Kris M. Kitani. Learnable embedding space for efficient neural architecture compression. In *International Conference on Learning Representations*, 2019.
- [4] Leslie Pack Kaelbling, Michael Littman, and Andrew P Moore. Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, 4:237–285, 04 1996.
- [5] Q. Song, J. Ni, and G. Wang. A fast clustering-based feature subset selection algorithm for high-dimensional data. *IEEE Transactions on Knowledge and Data Engineering*, 25(1):1–14, Jan 2013.
- [6] Zhaleh Manbari, Fardin Akhlaghian Tab, and Chiman Salavati. Hybrid fast unsupervised feature selection for high-dimensional data. *Expert Systems with Applications*, 124, 06 2019.
- [7] Shenkai Gu, Ran Cheng, and Yaochu Jin. Feature selection for high-dimensional classification using a competitive swarm optimizer. *Soft Computing*, 22(3):811–822, Feb 2018.
- [8] Irene Rodriguez-Lujan, Ramon Huerta, Charles Elkan, and Carlos Santa Cruz. Quadratic programming feature selection. *J. Mach. Learn. Res.*, 11:1491–1516, August 2010.