

Исследование важности входных признаков при прогнозировании геомагнитного индекса алгоритмами машинного обучения

Р.Д. Владимиров
Московский государственный
университет
имени М.В. Ломоносова
Москва, Россия
vladimirov.rd16@physics.msu.ru

И.Н. Мягкова
Московский государственный
университет
имени М.В. Ломоносова
НИИ ядерной физики
имени Д.В. Скобельцына
Москва, Россия
irina@srd.sinp.msu.ru

В.Р. Широкий
Московский государственный
университет
имени М.В. Ломоносова
НИИ ядерной физики
имени Д.В. Скобельцына
Москва, Россия
shirokiy@srd.sinp.msu.ru

С.А. Доленко
Московский государственный
университет
имени М.В. Ломоносова
НИИ ядерной физики
имени Д.В. Скобельцына
Москва, Россия
dolenko@srd.sinp.msu.ru

О.Г. Баринов
Московский государственный
университет
имени М.В. Ломоносова
НИИ ядерной физики
имени Д.В. Скобельцына
Москва, Россия
obar@snp.msu.ru

Аннотация—Одним из эффективных инструментов для прогнозирования временных рядов является использование методов машинного обучения, в частности, искусственных нейронных сетей. Однако при этом необходимым этапом исследования является понижение размерности входных данных. В данной работе рассматриваются результаты понижения размерности данных на основе ранжирования входных признаков по их существенности при решении задачи прогнозирования геомагнитного индекса Dst. Для оценки относительной существенности признаков используется итеративный подход, связанный с перебором моделей-кандидатов путём отбрасывания признаков по одному.

Ключевые слова— многомерный временной ряд, прогнозирование, отбор существенных признаков, искусственные нейронные сети, магнитосфера Земли, геомагнитный индекс Dst.

1. ВВЕДЕНИЕ

Геомагнитные возмущения представляют собой один из существенных факторов космической погоды, который с развитием космической отрасли становится все более важным. Прогнозирование геомагнитных возмущений, представляет интерес, так как сильные возмущения (магнитные бури) могут стать причиной нарушений в работе телеграфных линий и радиосвязи, трубопроводов, линий электропередач и энергосетей [1]. Также они опосредованно оказывают влияние на радиационные условия в космическом пространстве, поскольку после примерно половины магнитных бурь возрастает поток релятивистских электронов внешнего радиационного пояса Земли, что сбоям в электронных микросхемах спутниковой аппаратуры.

Состояние магнитосферы Земли характеризуется геомагнитными индексами. Одним из наиболее используемых является индекс Dst. Магнитосфера Земли и представляет собой динамическую систему, будущее состояние которой зависит не только от её текущего состояния и от текущего воздействия со стороны

солнечного ветра, но и от предыстории. Это приводит к тому, что размерность используемых для прогнозирования данных, описывающих текущее состояние магнитосферы и предысторию на необходимую глубину, оказывается достаточно высока.

Одним из эффективных инструментов для прогнозирования временных рядов является использование методов машинного обучения (МО), в частности, искусственных нейронных сетей [2, 3]. Однако при этом высокая размерность входных данных может приводить к нежелательным последствиям, таким, как высокая вычислительная стоимость обучения и переучивание. Также отбор существенных входных признаков (ВП) может позволить сделать некоторые выводы о взаимосвязях различных физических величин, значения которых используются в качестве ВП.

Целью настоящей работы являлось применение адаптивного метода ранжирования ВП по их существенности на основе их отбрасывания по одному, а также оценка полученных результатов.

2. ВХОДНЫЕ ДАННЫЕ И ОПИСАНИЕ АЛГОРИТМА

В качестве входных признаков использовались исторические значения характеристик солнечного ветра (скорости SW_spd и плотности H_den) и межпланетного магнитного поля (компонент V_y и V_z в системе GSM и модуля вектора поля $|V|$), измеренные в точке Лагранжа L1 в система Солнце-Земля, и значения самого прогнозируемого индекса Dst, а также временные характеристики, описывающие фазу вращения Земли вокруг Солнца и вокруг своей оси. Каждый параметр, кроме временных характеристик, описывался 24 среднечасовыми значениями: текущим и 23 историческими значениями за последние сутки (всего 6 признаков * 24 часа + 4 = 148 ВП). Использовались данные за период с октября 1997 по конец 2016 года (тренировочный набор, из которого 20% случайно

