

Анализ методик интерпретации данных в применении к эксперименту АНИ

А.А.Чилингарян (ЕрФИ)

Проводится анализ существующих методик классификации многомерных данных, методик определения истинной размерности данных, линейных и нелинейных преобразований признакового пространства с целью выделения наиболее эффективных признаков. Приводится схема процедуры принятия статистических решений для эксперимента АНИ.

ВВЕДЕНИЕ

Изучать сильные взаимодействия в области энергий $> 10^{15}$ эВ в настоящее время возможно только с помощью экспериментов с космическими лучами.

Экспериментальные установки регистрируют потоки вторичных частиц, образующих ядерно-электромагнитный каскад (ЯЭК), или, иначе говоря, широкий атмосферный ливень (ШАЛ). Замена измерения характеристик первичного взаимодействия измерением параметров ШАЛ ведет к значительной reduktion информации о характере сильного взаимодействия. Прямо не измеряются три самых важных параметра - энергия и тип первичной частицы, высота первого взаимодействия. Неопределенность этих параметров приводит к очень большим флуктуациям экспериментально измеряемых величин. Поэтому моделирование ЯЭК является единственной возможностью сравнения экспериментальных данных о ШАЛ с различными моделями сильного взаимодействия. Вычислительные мощности современных ЭВМ позволяют создавать весьма детальные имитационные программы развития ЯЭК в пространстве и времени [1,2]. Однако принятие решения о выборе правильной модели из списка альтернатив затруднено из-за статистического характера и неизбежных упрощающих предположений моделирования, а также экспериментальных ошибок.

Обычно не удается установить взаимно-однозначного соответствия между моделью и полученной экспериментальной информацией.

Целью эксперимента АНИ [3] является получение данных в области энергий 10^{15} - 10^{16} эВ, достаточных для выбора модели сильного взаимодействия.

Для планирования и интерпретации эксперимента АНИ предлагается использовать методы и процедуры, развитые в рамках подхода распознавания образов.

Разработка методики проведения статистического анализа позволит дать рекомендации об оптимальном составе измерительного комплекса, оценить необходимое время измерений.

Всякое распознавающее устройство состоит из измерительного комплекса и системы принятия решений.

Задается список альтернатив - возможных состояний природы, отвечающих определенным модельным представлениям. Каждый тип (образ, класс) характеризуется некоторой областью в признаковом пространстве - совокупностью реализаций моделирующей программы. Задачей распознавания является выбор признакового пространства и решающего правила, позволяющего принять точные статистические решения о принадлежности результатов эксперимента одному из классов, не превышая допустимых затрат.

Признаковое пространство - N -мерное пространство измеряемых величин. Каждая координата соответствует экспериментально измеряемому параметру ШАЛ. Решающее правило - алгоритм принятия решения о сходстве или различии вектора измерения с имеющимися классами.

Стоимостные характеристики - время эксперимента и стоимость оборудования, задают ограничения на точность и количество измеряемых параметров и достижимую статистику.

Совокупность классов, заданных в признаковом пространстве, назовем обучающей выборкой, ее элементы A_{ij}^k - векторами, точками, представителями.

$$A_{ij}^k \quad \begin{array}{l} k = 1, \dots, N - \text{размерность пространства,} \\ i = 1, \dots, M - \text{число векторов в классе,} \\ j = 1, \dots, L - \text{число классов,} \end{array}$$

или (\vec{A}_i, ω_j) $\omega_i - i$ - i класс.

Если не делаем различия в классах, то

$$A_i^k, \quad i = 1, \dots, M; \quad ML = M \cdot L.$$

Векторы экспериментальных величин, поступающие на распознавание, обозначим

$$\vec{x}, \quad x^k \quad k = 1, \dots, N.$$

Будем считать, что априорные вероятности появления представителей различных классов одинаковы и равны $1/L$. Для вычисления степени "близости" векторов в признаковом пространстве будем использовать функцию расстояния с обычными свойствами [4]. Из большого числа предложенных метрик в N -мерном пространстве [5] выберем евклидову

$$D^k(\vec{A}_i, \vec{x}) = A_i^k - x^k; \quad D(\vec{A}_i, \vec{x}) = \left[\sum_{k=1}^N D^k(\vec{A}_i, \vec{x})^2 \right]^{1/2}.$$

Для примера используем простую модель классификации - случай $L = 2$. Общий случай можно привести к нему последовательным разбиением.

Через $P(\vec{x}), P(\vec{x}/\omega_1), P(\vec{x}/\omega_2)$ обозначим соответственно функцию распределения вероятности случайной величины и ее плотность и условные функции распределения и ее плотность.

Под ранжированием выборки A_i будем понимать такое ее упорядочение, что $A_i \leq A_{i+1}$ $i = 1, \dots, M-1$ или $A_i \geq A_{i+1}$ $i = 1, \dots, M-1$. В дальнейшем знак \leq обозначает математическое ожидание, $-T$ - транспонирование, λ - оценку, ε - дисперсию, tr - штур.

I. КЛАССИФИКАТОРЫ

Классификацией назовем процесс установления соответствия между случайнм вектором \vec{x} и классами ω_j , $j = 1, \dots, L$.

В N -мерном пространстве предъявляется вектор \vec{x} . Надо определить, какому из распределений, описанных в этом пространстве, он принадлежит.

Допустим, плотности вероятности (ПВ) этих распределений заданы явно, тогда решающее правило (РН) можно записать в виде [6]

$$\ell(\vec{x}) = P(\vec{x}/\omega_1)/P(\vec{x}/\omega_2) \geq 1 \rightarrow \vec{x} \in \{\omega_1\}, \quad (I.1)$$

$\ell(\vec{x})$ - отношение правдоподобия.

Вектор \vec{x} относится к тому классу, чья ПВ в точке \vec{x} наибольшая. РН (I.1) называют байесовским критерием, минимизирующим ошибку решения. Это РН оптимально в том смысле, что его использование обеспечивает наименьшую вероятность совершения ошибки классификации.

Можно сформировать РН, используя понятие штрафных функций. Пусть C_{ij} - штраф за решение $\vec{x} \in \omega_j$, если на самом деле $\vec{x} \in \omega_i$. Тогда байесовский критерий, минимизирующий средние потери, имеет вид

$$P(\vec{x}/\omega_1)/P(\vec{x}/\omega_2) \geq \frac{C_{21} - C_{22}}{C_{12} - C_{11}} \rightarrow \vec{x} \in \{\omega_1\}. \quad (I.2)$$

Если положить $C_{11} = 0$; $C_{ij} = 1 - \delta$, $\delta = \text{const}$, то есть если в случае правильной классификации штраф равен нулю, а в случае любой ошибки одинаков, то (I.2) совпадает с (I.1).

Для конструирования оптимального классификатора необходимо вычисление ПВ для всех классов. Значение $P(\vec{x}/\omega_j)$ в явном виде редко бывает известно. Для их оценки используется обучающая выборка. Оценивание параметров ПВ известного вида и дальнейшую классификацию объектов на основании этих оценок называют параметрическим методом классификации. В случае отсутствия априорных знаний мы должны как-то оценить ПВ, не зная ее структуры. Этот случай называется непараметрическим оцениванием.

Широкое распространение в этом подходе получил метод потенциальных функций [7], в котором предполагается, что каждая точка обучающей выборки является источником "поля". На основании этой аналогии можно допустить существование эквипотенциальных поверхностей, которые описываются потенциальной функцией $\psi(\vec{A}_i, \vec{x})$. Значение ПВ в точке \vec{x} оценивается как суперпозиция

потенциальных функций

$$\hat{P}(\vec{x}/\omega_j) = \frac{1}{M} \sum_{i=1}^M \delta(\vec{x}, \vec{A}_{ij}). \quad (I.3)$$

Желательно выбирать потенциальные функции гладкими и монотонными [8] с выполнением условий:

- 1) $\delta(\vec{A}_i, \vec{x})$ — максимальна при $\vec{x} = \vec{A}_i$;
- 2) $\delta(\vec{A}_i, \vec{x}) \rightarrow 0$ при $D(\vec{A}_i, \vec{x}) \rightarrow \infty$.

Возможен и другой тип разложения ПВ [9]

$$\hat{P}(\vec{x}/\omega_j) = \vec{B}_j^T \vec{\phi}(\vec{x}) \quad j = 1, \dots, L, \quad (I.4)$$

где $\vec{\phi}^T = (\varphi_1, \varphi_2, \dots, \varphi_q)$,

$$\vec{B}_j^T = (b_{j1}, b_{j2}, \dots, b_{jq}).$$

Здесь $\varphi_i(\vec{x})$ — линейно независимый базис,

B_j — вектор неизвестных параметров, который должен быть определен так, чтобы ошибка аппроксимации $P(\vec{x}/\omega_j)$ была минимальна.

Обычно минимизируют критерий вида

$$J(B) = \sum_{i=1}^{M_L} [\Pi(\vec{A}_i) - B^T \phi(\vec{A}_i)]^T [\Pi(\vec{A}_i) - B^T \phi(\vec{A}_i)], \quad (I.5)$$

где $\Pi(\vec{A}_i) = [P(\vec{A}_i/\omega_1), P(\vec{A}_i/\omega_2), \dots, P(\vec{A}_i/\omega_L)]$,

$$\vec{B}^T = [B_1, B_2, \dots, B_L].$$

Для минимизации функционала $J(B)$ вместо неизвестной функции $\Pi(\vec{A}_i)$ можно использовать векторы $[\vec{d}_1, \vec{d}_2, \dots, \vec{d}_L]$:

$$\vec{d}_{ij} = \begin{cases} 1 & \vec{A}_i \in \omega_j \\ 0 & \vec{A}_i \notin \omega_j \end{cases} \quad j = 1, 2, \dots, L.$$

Решая уравнение $\nabla J(B) = 0$, получаем решение в виде

$$B = \Delta^{-1} C; \quad \Delta = \sum_{i=1}^{M_L} \Phi(\vec{A}_i) \Phi^T(\vec{A}_i); \quad C = \sum_{i=1}^{M_L} \Phi(\vec{A}_i) \vec{d}_i^T. \quad (I.6)$$

Следующий тип оценки определяет ПВ как математическое ожидание δ -функции [10]

$$P(x) = \int \delta(x-A) P(A) dA = E\{\delta(x-A)\}. \quad (I.7)$$

Вычислим математическое ожидание по обучающей выборке и заменим δ -функцию ее аппроксимацией:

$$K\left(\frac{x-A}{h}\right)/h \rightarrow \delta(x-A) \quad \text{при } h \rightarrow 0,$$

где $K(x)$ — ядро оценки,

$$\hat{P}_M(x) = E_A \left\{ \frac{1}{h} K\left(\frac{x-A}{h}\right) \right\} \approx \frac{1}{Mh} \sum_{i=1}^M K\left(\frac{x-A_i}{h}\right), \quad (I.8)$$

$\hat{P}_M(x) \rightarrow P(x)$, если $h \rightarrow 0$ при $M \rightarrow \infty$ и $Mh \rightarrow \infty$ при $M \rightarrow \infty$

В работе [11] (I.8) рассматривается как некорректная задача численного дифференцирования и $\hat{P}(x)$ определяется как решение интегрального уравнения Фредгольма I рода.

В работе [12] для классификации предлагается использовать концепцию статистически эквивалентных блоков. Обучающая выборка (A_i, ω_i) ранжируется и разбивается на блоки по K членов $K \sim \sqrt{M}$ и, если

$$\min_{i \in \omega_j} \{C_{ij} \times N\omega_i\} < C_0 \times K, \quad (I.9)$$

то блоку присваивается уровень ω_i ,

$N\omega_i$ — число членов, принадлежащих ω_i в блоке,

C_0 — штраф за отказ принять решение,

C_{ij} — штраф за ошибочное решение.

Далее, объединяя блоки с одинаковой классификацией и распределяя неопределенные блоки по соседним, получаем дробление признакового пространства на области с известной классификацией.

Наиболее распространенный классификатор основан на априорном представлении о "близости" объектов, принадлежащих одному классу [13]. Вычисляются расстояния от точки \vec{x} до всех представителей обучающей выборки и ранжируются в порядке возрастания.

Оценку ПВ в точке \vec{x} получим в виде

$$P_M(\vec{x}) = \frac{K}{M} \frac{1}{V(K, N, \vec{x})}, \quad (I.10)$$

где $V(K, N, \vec{x})$ - объем гипершара радиуса R_K с центром в точке \vec{x} . R_K - расстояние от точки \vec{x} до K -го ближайшего соседа (КБС). $P_K(\vec{x})$ является несмещенной и состоятельной оценкой ПВ, если

$$\lim_{M \rightarrow \infty} K = \infty \quad \text{и} \quad \lim_{M \rightarrow \infty} K/M = 0 \quad (I.11)$$

условие (I.11) не является необходимым. В работе [14] показано, что даже для малых ($M \sim 10$) обучающих выборок КБС оценки близки к оптимальным.

Подставляя КБС оценку ПВ в (I.1), получаем

$$\text{где } K_1 \geq K_2 \rightarrow \vec{x} \in \{\omega_1, \omega_2\}, \quad K_1 + K_2 = K. \quad (I.12)$$

Таким образом, решение о принадлежности \vec{x} можно принять, определив K ближайших соседей, и выбрав класс, представленный наибольшим числом векторов.

Развитие КБС метода идет в направлении увеличения его точности и уменьшения требований к памяти ЭВМ.

В работе [15] предложено КБС правило с предварительным редактированием:

- для каждого члена обучающей выборки (\vec{A}_i, ω_j) находят K ближайших соседей и получают КБС оценку $\hat{\omega}_j$, и, если $\hat{\omega}_j \neq \omega_j$, \vec{A}_i удаляют из обучающей выборки.

В работе [16] показано, что использование отредактированной обучающей выборки, наряду с экономией памяти ЭВМ, приводит к значительному улучшению асимптотических свойств оценок.

Если в обучающей выборке каждому классу соответствует неодинаковое количество представителей, то вводится взвешенное расстояние [17]

$$D(\vec{A}_{ij}, \vec{x}) = \left[\frac{M_j}{\lambda_j \sum_{j=1}^M M_j} \right]^{1/N} D(\vec{A}_{ij}, \vec{x}), \quad (I.13)$$

где M_j - число представителей в классе ω_j ;

λ_j - априорная вероятность j -го класса.

Для того, чтобы избежать неправильной классификации отдельных, далеко отстоящих точек, предложен [18] (K, K') БС алгоритм, в котором требуется, чтобы K ближайших соседей вектора \vec{x} находились на расстоянии, не большем R_p .

R_p - пороговое значение, при превышении которого \vec{x} не может считаться близким к классу ω_j . При выполнении этого условия \vec{x} относится к классу ω_j , если не менее K' из K соседей \vec{x} принадлежат ω_j ($K' < K$).

КБС правило зависит от выбранной метрики и, следовательно, от выбранных единиц измерения признаков. Если измеряются качественно разные признаки (энергия, координаты и т.д.) трудно выбрать согласованный масштаб переменных. В работе [19] предложено эмпирическое расстояние, инвариантное относительно монотонных преобразований координатных осей. Сначала множества

$$\vec{A}_i^K, \vec{x}^K \quad i = 1, \dots, M; \\ K = 1, \dots, N$$

ранжируются в порядке возрастания. Обозначим ранги \vec{A}_i^K и \vec{x}^K через τ_i^K и τ_x^K . Расстояние между \vec{x} и \vec{A} определяется так:

$$D(\vec{A}_i^K, \vec{x}^K) = \max_{1 \leq k \leq N} |\tau_i^K - \tau_x^K|.$$

При сборе информации с крупной экспериментальной установки трудно добиться постоянной 100% эффективности всех каналов регистрации. Поэтому векторы измерений будут содержать пропуски, и расстояния между ними и векторами обучающей выборки не будут определены. В работе [20] исследованы различные алгоритмы "заполнения" пропусков:

I. Алгоритм нормализации расстояния

$$D^k(A_i^k, x^k) = \begin{cases} 0, & \text{если } x^k \text{ отсутствует} \\ A_i^k - x^k, & \end{cases} \quad (I.14)$$

$$D^2(\vec{A}_i, \vec{x}) = \frac{N}{N-NB} \sum_{k=1}^N D_k^2,$$

NB - число пропусков.

2. Усреднение

$$D^k(A_i^k, x^k) = \begin{cases} A_i^k - \bar{x}^{*k}, & \text{если } x^k \text{ отсутствует} \\ A_i^k - x^k, & \end{cases} \quad (I.15)$$

$$\bar{x}^{*k} = \frac{2}{M_2(M_2-1)} \sum_{i=1}^{M_2} |x_i^k - \bar{x}^k|,$$

M₂ - число экспериментальных векторов.

3. Зануление

$$D^k(A_i^k, x^k) = \begin{cases} 0, & \text{если } x^k \text{ отсутствует} \\ A_i^k - x^k, & \end{cases} \quad (I.16)$$

Эти алгоритмы проверялись для широкого диапазона данных и разного процента пропусков. Лучшие результаты получены при использовании формулы (I.15).

2. ПРЕОБРАЗОВАНИЕ ПРИЗНАКОВОГО ПРОСТРАНСТВА И ИСТИННАЯ РАЗМЕРНОСТЬ ДАННЫХ

Выбор множества признаков, обеспечивающих необходимое качество классификации, представляет собой одну из наиболее трудных задач построения распознающих систем.

Очевидно, что наиболее важными являются признаки, различающиеся для разных классов, а одинаковые для всех классов признаки - **малоценны**. Выбор признаков можно рассматривать как процесс преобразования исходных измерений в более эффективные. Если это преобразование является линейным, то задача сводится к нахождению коэффициентов линейной функции, максимизирующей или минимизирующей некоторый критерий, экстремум которого можно найти с помощью методов линейной алгебры, или, в случае сложного критерия, с помощью методов оптимизации.

Следует учесть, что разделимость классов зависит не только от распределения векторов, но и от используемого классификатора. Лучшим критерием эффективности признаков является малая вероятность ошибки классификации. Из-за того, что вероятность ошибки часто не имеет явного математического выражения, ее подсчитывают экспериментально, выбрав набор признаков и построив классификатор.

В работе [21] предложено несколько критериев разделимости классов. Они строятся исходя из следующих требований:

- 1) монотонная связь с вероятностью ошибки или с верхней и нижней границей вероятности;
- 2) инвариантность относительно взаимно-однозначных отображений;
- 3) аддитивность по отношению к независимым признакам.

Критерии разделимости классов формулируются с использованием матрицы рассеяния внутри классов

$$G_\omega = \sum_{j=1}^L E \{ (\vec{A} - \vec{M}_j) (\vec{A} - \vec{M}_j)^T \} = \sum_{j=1}^L \Delta_j, \quad (2.1)$$

Δ_j - выборочная ковариационная матрица и матрица рассеяния между классами

$$G_g = \sum_{j=1}^L (\vec{M}_j - \vec{M}_o) (\vec{M}_j - \vec{M}_o)^T, \quad (2.2)$$

\vec{M}_j - математическое ожидание векторов класса,

\vec{M}_o - математическое ожидание смеси распределений.

Матрицы G_ω и G_g инвариантны относительно трансляций системы координат. Критерий разделимости должен увеличиваться при увеличении рассеяния между классами и при уменьшении рассеяния внутри классов. Поэтому можно использовать критерий вида

$$J(\vec{A}_i) = t_2(G_\omega^{-1} G_g). \quad (2.3)$$

Преобразование исходного пространства, максимизирующее J , получаем следующим образом. Нахо-

дим собственные функции и собственные значения матрицы $\{G_{\omega}^{-1} \times G_{\theta}\}$. Выбираем К "существенных" собственных функций, соответствующих наибольшим собственным значениям, $K < N$, и строим из них матрицу С, которую используем для преобразования N -мерных векторов в новое К-мерное пространство:

$$\vec{A}'_i = C^T \vec{A}_i \quad C^T = (\vec{\Psi}_1, \vec{\Psi}_2, \dots, \vec{\Psi}_K).$$

Отметим, что ранг матрицы G_{θ} равен L , поэтому матрица $G_{\omega}^{-1} \times G_{\theta}$ имеет только L ненулевых собственных значений. Остальные $N-L$ признаков не вносят вклад в J .

Алгоритмы линейного преобразования прости и не требуют больших вычислительных мощностей, однако в информации о ШАЛ имеются важные признаки, являющиеся существенно нелинейными функциями исходных измерений.

В настоящее время не существует сколько-нибудь общей теории нелинейных преобразований, поэтому процедуры выделения признаков в сильной степени зависят от конкретной задачи и обычно носят итеративный характер.

Размерность вектора исходной информации с эксперимента АИ может достигать $\sim 10^3$. Поэтому важно выбрать преобразование, проецирующее эти векторы в пространство наименьшей возможной размерности.

Минимальное число параметров, необходимое для генерации многомерных данных называют истинной или топологической размерностью. Кривая в N -мерном евклидовом пространстве обладает истинной размерностью $N_t = 1$, поверхность любой формы $N_t = 2$ и т.д... Если N -мерные данные лежат на гиперповерхности размерности $K < N$, то их истинная размерность равна K . K задает нижнюю границу возможного понижения размерности и определяет необходимое и достаточное число эффективных признаков.

Идея, используемая для определения истинной размерности многомерных данных, основана на изучении расстояний соседства в локальных областях N -мерного пространства. Рассмотрим векторы, равномерно распределенные в гипершаре радиуса r . Нормализуем евклидово расстояние между двумя векторами из этого шара,

$$R_N = |\vec{x}_i - \vec{x}_j| / 2r.$$

Исследование распределения расстояний R_N показало [22], что дисперсия величины R_N обратно пропорциональна размерности пространства N ,

$$N^6 (R_N) \approx \text{const}, \quad (2.6)$$

то есть с уменьшением размерности пространства происходит "расширение" локальных сгустков. Следовательно, ключевая процедура уменьшения размерности – локальное увеличение дисперсии межточечных расстояний. Этого можно добиться, увеличивая расстояния, большие среднего, и уменьшая расстояния, меньшие среднего. Такая деформация, однако, может сильно искажить исходное распределение. Поэтому надо определить самые существенные свойства исходного распределения и потребовать их инвариантность относительно преобразований, понижающих размерность. Предложенные к настоящему времени инварианты строятся как функции расстояния между точками исходного пространства.

В первых работах по нелинейному проецированию и многомерному шкалированию [23-25] за инвариант был принят ранговый порядок межточечных расстояний в исходном пространстве.

Сначала вычисляются ММ расстояний $MM = M \cdot (M-1)/2$, M – общее число векторов обучающей выборки, множество расстояний ранжируется и ранги всех расстояний запоминаются. Затем выбирается произвольная конфигурация точек в пространстве минимальной размерности, вычисляются ММ расстояний в этом пространстве и рассматривается критерий сходства двух конфигураций:

$$S_{KR} = \sum_{i,j} S_{ij} (D_{ij} - D_{ij}^*) / D_{ij}, \quad (2.7)$$

или

$$S_{SP} = \sum_{i,j} f_{ij} (R(D_{ij}) - R(D_{ij}^*)), \quad (2.8)$$

где D_{ij}^* – межточечные расстояния в новом пространстве, сохраняющие ранговый порядок исходной выборки; D_{ij} – текущее значение межточечных расстояний; $R(D_{ij})$ – ранг расстояния между точками выборки \vec{A}_i и \vec{A}_j . f_{ij} – весовой множитель. В частности, f_{ij} можно использовать для определения локальной области соседства, положив $f_{ij} = 0$, если $D_{ij} > D_{ij}^*$, или, если \vec{A}_i и \vec{A}_j , принадлежат разным классам обучающей выборки. Критерий минимизируется изменением положения всех точек. Итеративный процесс продолжается до тех пор, пока S не перестает уменьшаться. Если значение S все еще достаточно велико, то размерность исходного пространства увеличивается.

ется на 1 и процесс минимизации продолжается. Когда S достигнет величины $5 - 10\%$, можно считать, что размерность пространства совпадает с топологической.

Определить истинную размерность можно и другим путем [22]. В исходном пространстве увеличиваем дисперсии локальных расстояний, следя за тем, чтобы ранговые соотношения в локальной области не изменялись. Итеративный процесс происходит до тех пор, пока дальнейшее увеличение дисперсии становится невозможным. Затем вычисляется ковариационная матрица полученной конфигурации точек, и истинная размерность определяется количеством "существенных" собственных векторов этой матрицы.

Но как определить размеры локальных областей? Ответ на этот вопрос можно получить, определив структурный инвариант, автоматически определяющий локальную близость [26].

Минимальным графом (МГ) M точек в метрическом пространстве назовем древовидную структуру, соединяющую все точки и имеющую минимальную длину.

Основные свойства МГ:

- 1) Ближайшие соседи (БС) соединены между собой в МГ;
- 2) МГ инвариантен относительно любых преобразований, сохраняющих ранговый порядок межточечных расстояний;
- 3) МГ легко вычисляется и его чертеж ясно отражает структуру выборки.

Локальная область определяется соединенными за МГ точками.

Понижение размерности осуществляется специальным преобразованием, сохраняющим структуру МГ. Концепция МГ была использована для построения алгоритма геометрической реконструкции сложных событий в дрейфовых и пропорциональных камерах [27]. Быстродействие алгоритма увеличилось на порядок в сравнении с традиционными методами.

Отрицательной стороной изложенных методик определения истинной размерности является итеративный характер алгоритмов и эвристические правила остановки процесса минимизации.

В работе [28] предложен неитеративный алгоритм, основанный на информации ближнего соседства. Используя оценку плотности вероятности расстояния до K -го ближайшего соседа $\bar{\tau}_k$, удается получить линейное уравнение, угловой коэффициент которого можно использовать как оценку истинной размерности:

$$\log \bar{\tau}_k = (1/\hat{N}) \log K + \text{const}. \quad (2.9)$$

$\bar{\tau}_k$ - усредненное по выборке расстояние до K -го соседа,

\hat{N} - оценка истинной размерности.

Вычислив $\bar{\tau}_k$ для нескольких значений K , можно получить значение \hat{N} методом линейной регрессии. Практически, за истинную размерность берется ближайшее к \hat{N} целое число.

Алгоритм проверялся с использованием данных выборок различного типа: многомерного нормального распределения, равномерного и т.д. Чем больше значений K и $\bar{\tau}_k$ используется, тем точнее будет линейная аппроксимация, но даже для прямых, построенных по трем значениям $\bar{\tau}_k$ ($K = 1, 2, 3$), получена хорошая сходимость алгоритма. Алгоритм вычисления прост и включает вычисление БС информации, используемой при классификации, так что объем добавочных вычислений невелик. Для получения проекций M векторов исходного пространства в пространство "истинной" размерности можно использовать алгоритмы случайного поиска с возвратом при неудачном шаге [29].

Итеративный процесс начинается с M произвольных векторов в пространстве "истинной" размерности. На каждом шаге минимизации критерий сходства двух выборок генерируется случайный вектор \vec{u} и новые координаты выборки вычисляются по формуле

$$\vec{A}_i = \vec{A}_i + \vec{u}. \quad (2.10)$$

Если критерий, вычисленный с новыми координатами, уменьшился, то новые значения принимаются, в противном случае \vec{u} вычитается и генерируется новый случайный вектор. Процесс продолжается до тех пор, пока дальнейшего улучшения получить не удается. Минимизируется критерий вида

$$S_{\text{cal}} = \sum_{i < j} f_{ij} D_{ij}^2 / D_{ij}^2, \quad (2.11)$$

В алгоритме используются ограничения, запрещающие сближение точек, принадлежащих разным классам.

Таким образом, можно получить обучавшую выборку компактных, хорошо разделенных классов. Но как использовать эту выборку для классификации? Ведь у нас нет явного математического выражения преобразования признакового пространства, а применять итеративное преобразование к векторам, поступающим на распознавание, бессмысленно.

Имея результаты действия неизвестного нелинейного преобразования выборки $\vec{A}_i' - \vec{A}_i$, можно поставить задачу его аппроксимации, предположив

$$\vec{A}' = \sum_{j=1}^q b_j \varphi_j(\vec{A}_i),$$

где φ_j – известные функции.

Коэффициенты b_j находятся из условия минимизации среднеквадратичной ошибки аппроксимации по всей обучающей выборке:

$$J(\vec{b}, \vec{\varphi}) = \sum_{i=1}^M |\vec{A}_i' - \vec{A}_i|^2. \quad (2.12)$$

Если удастся получить хорошую точность для не очень больших значений q , то (2.11) можно использовать для преобразования экспериментальных векторов и их последующей классификации.

Рассмотрим теперь специальный случай проецирования на признаковое пространство размерности $N=2$. Это преобразование дает возможность отображать данные на экране индикатора и анализировать данные в режиме диалога с ЭВМ.

В работе [30] для целей двухмерного проектирования предложено минимизировать критерий

$$S_{SAM} = \sum_{i,j} f_{ij} (D_{ij} - D_{ij}^*) , \quad (2.13)$$

то есть ищется конфигурация в двухмерном пространстве с теми же межточечными расстояниями, что и в исходном пространстве. Итеративный процесс начинается с произвольной конфигурации, для достижения минимума S используется метод наискорейшего спуска – минимизирующий шаг делается в антиградиентном направлении.

В работе [31] предложен алгоритм последовательного проецирования N -мерных точек на плоскость. Первые три точки можно спроектировать так, чтобы все 3 межточечных расстояния сохранились. Последующие точки можно спроектировать с сохранением лишь двух расстояний. Всего сохраняется $3 + 2 \times (M-3) = 2M-3$ расстояний из M . Как уже указывалось, минимальный граф МГ содержит $M-1$ расстояние. Поэтому при проецировании можно сохранить структуру МГ. Оставшиеся $M-2$ расстояния можно использовать для сохранения расстояний до второго ЕС, или можно потребовать сохранения всех расстояний до специально выделенной точки. Такой подход позволяет получить проекции, каждая из которых отвечает определенной точке зрения на структуру данных.

ЗАКЛЮЧЕНИЕ

Назначение статистических методов анализа – интерпретация данных с ощущимой случайной изменчивостью. Традиционный способ статистического анализа – выбор семейства вероятностных моделей с последующим оцениванием параметров этих моделей. Если, из-за сложности явления, трудно описать физическую ситуацию адекватной вероятностной моделью, используются непараметрические методы, которые позволяют формировать статистики с распределениями, одинаковыми для более общего семейства плотностей вероятности, чем параметрические семейства.

Для задач классификации достаточно оценить, какое из распределений обеспечивает наибольшую плотность вероятности в исследуемой точке признакового пространства.

Результаты моделирования ядерно-электромагнитного каскада в атмосфере трудно представить вероятностной моделью с небольшим числом параметров. Кроме того, моделирование ЯЭК с энергией $> 10^{15}$ эВ требует очень больших машинных мощностей и, следовательно, количество реализаций каскадов с такими энергиями будет невелико.

Поэтому для интерпретации данных с эксперимента АНИ предлагается использовать методы, развитые в рамках концепций распознавания образов – ведущего подхода в попытках применения математических методов обработки информации, относящейся к сложным и плохо formalizedанным задачам [32].

Схема процедуры принятия статистических решений может быть представлена следующим образом:

1. Получение обучающей выборки – реализация имитационной программы для разных моделей сильного взаимодействия.
2. Определение истинной размерности векторов обучающей выборки.
3. Нелинейное преобразование векторов обучающей выборки и выделения эффективных признаков.
4. Редактирование обучающей выборки.
5. Заполнение "пропусков" в векторах экспериментальной информации.
6. Исследование качества классификации и выбор оптимального K в правиле K ближайших соседей.

7. Аппроксимация нелинейного преобразования и преобразование экспериментальных векторов.
8. Проведение классификации и выбор модели сильного взаимодействия.

Автор благодарен Т.Л.Асатиани и Э.А. Мамиджаняну за интерес к работе и поддержку, А.А.Аглинцеву и А.М.Дунаевскому за полезные обсуждения.

СПИСОК ЛИТЕРАТУРЫ

1. Grieder P.K.F., The relation between extensive air shower data and high energy particle production models, Revista del N. Cim., 1977, vol. 7, p. 1.
2. Dunaevskii A.M., Emelyanov Yu.A., Shorin B.P., Urysson A.V., The calculation of nuclear-electromagnetic cascades, Preprint FIAN, N 149, 1980.
3. Никольский С.И., Туким Е.И., Файнберг Е.Л. и др. Исследование взаимодействия адронов и ядер космического излучения при энергиях $10^3 - 10^5$ ТэВ (проект эксперимента АНИ). -НСЕИИ, 358(16), 1974, Ереван, 1974.
4. Елисеева И.И., Рукавишников В.О. Группировка, корреляция, распознавание образов. М.:Статистика, 1977, с.35.
5. Доран Б., Оделл П. Кластерный анализ ,М.:Статистика, 1977, с.17.
6. Фукунага К. Введение в статистическую теорию распознавания образов, М.:Наука, 1979, с.57, 176, 188, 267.
7. Айзerman M.A., Браверман Э.М., Розеноэр Л.И. Теоретические основы метода потенциальных функций в задачах об обучении автоматов разделению входных ситуаций на классы. -Автоматика и телемеханика, 1964, т.25, с.91%.
8. Meisel W.S., Potential function in mathematical pattern recognition, IEEE Trans. comput. 1969, vol. C 18, p. 911.
9. Kittler J., Classification of incomplete pattern vectors using modified discriminant functions.-IEEE Trans. comput. 1978, vol. C 27, p. 367.
10. Банник В.И., Стефанюк А.Р. Ненараметрические методы восстановления плотности вероятности. -Автоматика и телемеханика, 1978, т.39, с.38.
- II.Фукунага К. Введение в статистическую теорию распознавания образов, М.: Наука, 1979, с.57, 176, 188, 267.
12. Henrichson E.G., Fu K.S., A nonparametric Partitioning procedure for pattern classification, IEEE Trans. comput. 1969, vol. C 18, p. 614.
13. Фукунага К. Введение в статистическую теорию распознавания образов. М.: Наука, 1979, с.57, 176, 188, 267.
14. Levine A., Lustick L., Saltzberg B., The nearest-neighbor rule for small samples drawn from uniform distributions.-IEEE Trans. inform theory, 1973, vol. IT-19, p. 697.
15. Wilson D.C., Asymptotic properties of nearest neighbor rules using edited data. IEEE Trans. syst., man. and cybernet., 1972, vol. SMC 2, p. 408.
16. Wagner T.J., Convergence of the edited nearest neighbor.-IEEE Trans. inform theory, 1973, vol. IT 19, p. 696.
17. Brown T.A., Koplowitz T., The weighted nearest neighbor rule for class dependent sample size.-IEEE Trans. comput., 1978, vol. C 27, p. 453.
18. Ben-Bassat M., Newhouse M., Bailis E.I., Object recognition and learning of known and unknown types,-IEEE Trans. comput., 1978, vol. C 27, p. 66.
19. Devroye L.P., A universal k-nearest neighbor procedure in discrimination.-IEEE Trans. comput., 1978, vol. C 27, p. 142.
20. Dixon J.K., Pattern recognition partly missing data, -IEEE Trans. syst., man. and cybernet., 1979, vol. SMC9, p. 617.
21. Фукунага К. Введение в статистическую теорию распознавания образов. М.: Наука 1979, с.57, 176, 188, 267.
22. Bennet R.S., The intrinsic dimensionality of signal collections, IEEE Trans. inform. theory, 1969, vol. IT 15, p. 517.
23. Shepard R.N., The analysis of proximities: multidimensional scaling with an unknown distance function.-Psychometrika, 1962, vol. 27, p. 12, 5, 219.
24. Kruskal J.B., Multidimensional scaling by optimizing goodness of fit to a nonparametric hypothesis.-Psychometrika, 1964, vol. 29, p. 115.

25. Kruskal J.B., Nonlinear multidimensional scaling: a nuclear method, *Psychometrika*, 1964, vol. 29, p. 1.
26. Schwartzmann D.H., Vidal J.J., An algorithm for determining the topological dimensionality of point clusters.-*IEEE Trans. comput.*, 1975, vol. C 24, p. 1175.
27. Kowalski H., Pattern recognition in track chamber spectrometers, DESY - 80/72, 1980.
28. Pettis K.W., Baiby T.A., Tain A.K., An intrinsic dimensionality estimator from near-neighbor information.-*IEEE Trans. Pattern Analysis Machine Intell.*, 1979, vol. PAMI 1, p. 25
29. Calvert W., Young T.Y., Randomly generated nonlinear transformations for pattern recognition.-*IEEE Trans. syst. science cybern.*, 1969, vol. SSC5, p. 266.
30. Sammon T.W., A nonlinear mapping for data structure analysis.-*IEEE Trans. comput.*, 1969, vol. C 18, p. 401.
31. Lee R.C.T., Slagle J.R., Blum H., A triangulation method for the sequential mapping of points from N-space to two-space.-*IEEE Trans. comput.*, 1977, vol. C 26, p. 288.
32. Сборник "Современные проблемы кибернетики", М.: Знание, № II, 1979.

Статья поступила в редакцию 15 декабря 1980 года.