

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ
«МОСКОВСКИЙ ИНСТИТУТ ЭЛЕКТРОННОЙ ТЕХНИКИ»

ОТЧЕТ ПО ПРОИЗВОДСТВЕННОЙ ПРАКТИКЕ

Направление подготовки — 01.03.04 «Прикладная математика»

Профиль — «Применение математических методов к решению инженерных и
экономических задач»

Выполнил студент Димаков Владислав Сергеевич
Оценка руководителя практики от кафедры ВМ-1
Козлитин Иван Алексеевич, доцент, к.ф.-м.н.

Группа: МП-40

(оценка)

(подпись)

Москва

2017

Оглавление.

Введение.	2
Задача классификации.	2
Обзор существующих методов классификации.	3
Деревья решений.	4
Алгоритм CART.....	5
Оптимальное расщепление вершин.....	5
Алгоритм построения случайного леса.....	6
Признаки Хаара.	6
Создание обучающей выборки.	7
Тестирование классификатора.	7
Заключение.....	8
Список использованных источников.....	9

Введение.

В настоящее время от систем видеонаблюдения требуется не только предоставление возможности воспроизведения и записи видеопотока с камеры, но и возможности решения в автоматическом режиме множество задач без участия человека, начиная от простого детектирования движения в области наблюдения, заканчивая высокоточным подсчётом проехавших машин или прошедших людей. Одной из таких задач является задача классификации объектов, присутствующих на цифровой видеозаписи.

Задача классификации.

Вероятностная постановка задачи классификации выглядит следующим образом. Пусть множество пар «объект, метка класса» $X \times Y$ является вероятностным пространством с неизвестной вероятностной мерой P . Имеется конечная обучающая выборка наблюдений $D = \{x_i, y_i\}_{i=1}^l$, сгенерированная согласно вероятностной мере P . Требуется построить алгоритм $\alpha: X \rightarrow Y$, способный классифицировать произвольный объект $x \in X$

Далее под обучающей выборкой будем понимать независимую выборку $D = \{\bar{x}_i, y_i\}_{i=1}^l$ из некоторого неизвестного распределения $P(\bar{x}, y) = P(\bar{x})P(y|\bar{x})$. Здесь $\bar{x}_i = \{x_i^s\}_{s=1}^n$, $i = 1, 2, \dots, l$ – векторы признаков, координаты которых

представляют собой значения n признаков, измеряемых на некотором объекте; y_i – метки классов, $y_i \in \{\omega_1, \omega_2, \dots, \omega_c\}$, $c \geq 2$.

Обзор существующих методов классификации.

Существуют различные методы решения задачи классификации, такие как наивный байесовский классификатор, метод k ближайших соседей, искусственные нейронные сети с прямой связью и случайные леса.

Наивный байесовский классификатор – вероятностный классификатор, основанный на применении Теоремы Байеса со строгими предположениями о том, что объекты описываются n статистически независимыми признаками. Предположение о независимости существенно упрощает задачу, так как оценить n одномерных плотностей гораздо легче, чем одну n -мерную плотность. К сожалению, данное предположение крайне редко выполняется на практике, что влечет за собой относительно низкое качество классификации в большинстве реальных задач.

Метод k ближайших соседей – метод метрической классификации, основанный на оценивании сходства объектов. Классифицируемый объект относится к тому классу, к которому принадлежат k ближайших к нему объектов обучающей выборки. Недостатками метода является сложность выбора параметра k и высокая вычислительная трудоемкость, которая увеличивается квадратично с ростом числа записей в наборе данных.

Искусственная нейронная сеть с прямой связью является универсальным средством аппроксимации функций, что позволяет использовать ее в решении задач классификации. Как правило, нейронные сети оказываются наиболее эффективным методом классификации, так как способны генерировать большое число моделей, которые используются в решении задач классификации статистическими методами. Главными недостатками нейронных сетей являются сложность реализации их структур и подбор параметров моделей.

Случайный лес – статистический метод, предназначенных для решения задач классификации и регрессии. Метод основан на построении большого числа деревьев решений, каждое из которых строится по выборке, получаемой из обучающей с помощью бутстрепа (т.е. выборки с возвращением). В отличие от классических алгоритмов построения деревьев решений, в методе случайных лесов при построении каждого дерева на стадии расщепления вершин используется

только фиксированное число случайно отбираемых признаков и строится полное дерево, т.е. дерево без усечения.

Данный метод обладает следующими достоинствами:

- Способность эффективно обрабатывать данные с большим числом признаков и классов;
- Нечувствительность к масштабированию значений признаков;
- Одинаково хорошая обработка как непрерывных, так и дискретных признаков;
- Для построения случайного леса по обучающей выборке требуется задание всего двух параметров – количества деревьев в ансамбле и числа случайно отбираемых признаков на стадии расщепления вершин.

Несмотря на то, что данный метод также как и нейронные сети требует построения модели большого размера, он наилучшим образом подходит для решения поставленной задачи классификации объектов, присутствующих на цифровой видеозаписи.

Деревья решений.

Обозначим χ – пространство образов, т.е. множество всех возможных значений векторов признаков. Тогда деревом решений будет называться дерево, с каждой вершиной t которого связаны:

- Некоторое подмножество $\chi_t \subset \chi$. С корневой вершиной связывается все пространство образов χ .
- Подвыборка $D_t \subset D$ обучающей выборки D , такая, что $D_t = \{(\bar{x}, y) \in D \mid \bar{x} \in \chi_t\}$. С корневой вершиной связывается вся выборка D .
- Некоторая функция (правило) $f_t: \chi \rightarrow \{0, 1, \dots, k_t - 1\}$ (здесь $k_t \geq 2$ – количество потомков вершины t), определяющая разбиение множества χ на k_t непересекающихся подмножеств. С листьями дерева не связывается никакая функция.

Цель построения дерева решений состоит в классификации векторов \bar{x} из распределения $P(\bar{x})$. Процесс принятия решений начинается с корневой вершины и состоит в последовательном применении правил, связанных с вершинами дерева. Результатом этого процесса является определение листа t такого, что $\bar{x} \in \chi_t$. В

этом случае вектор \bar{x} относится к классу, являющемуся мажорантным (наиболее часто встречающимся) в подвыборке D_t , соответствующей данному листу.

Алгоритм CART.

Алгоритм CART (Classification and Regression Tree) предназначен для решения задач классификации и регрессии построением бинарного дерева решений. На каждом шаге построения дерева правило f_t , формируемое в узле t , делит обучающую выборку на две более однородные подвыборки:

Обычно вместо меры однородности используется противоположная по смыслу мера загрязненности. Пусть t – некоторая вершина дерева решений, D_t – подвыборка, связанная с этой вершиной, $i(t)$ – загрязненность вершины. Необходимо потребовать, чтобы загрязненность вершины была равна 0, если D_t содержит прецеденты только одного класса и была максимальной в случае, если D_t содержит одинаковое число прецедентов каждого класса.

Одной из наиболее часто используемых является мера загрязненности вершины, формализованная в индексе Gini:

$$i(t) = Gini(t) = 1 - \sum_{m=1}^c P^2(\omega_m),$$

где $P(\omega_m)$ – доля примеров класса ω_m в подвыборке D_t .

Оптимальное расщепление вершин.

Правило разбиения множества χ , связанное с каждой вершиной дерева решений, называется расщеплением. Бинарное расщепление вершины t можно рассматривать как функцию $f_t: \chi \rightarrow \{0,1\}$, $\bar{x} \in \chi$, где в случае $f_t(\bar{x}) = 0$ вектор \bar{x} относится к первому (левому) потомку, а в случае $f_t(\bar{x}) = 1$ – ко второму (правому).

Расщепление подвыборки естественно осуществлять таким образом, чтобы максимально уменьшить загрязненность. Уменьшение загрязненности вершины t для бинарных деревьев определяется как

$$\Delta i(t) = i(t) - P_L i(t_L) - P_R i(t_R),$$

где P_L и P_R – доли примеров подвыборки D_t , соответствующие левому и правому потомкам (t_L и t_R). Наилучшим расщеплением вершины t считается

разбиение, которое максимизирует величину $\Delta i(t)$, т.е. расщепление выполняется таким образом, чтобы $\Delta i(t) \rightarrow \max$.

Алгоритм построения случайного леса.

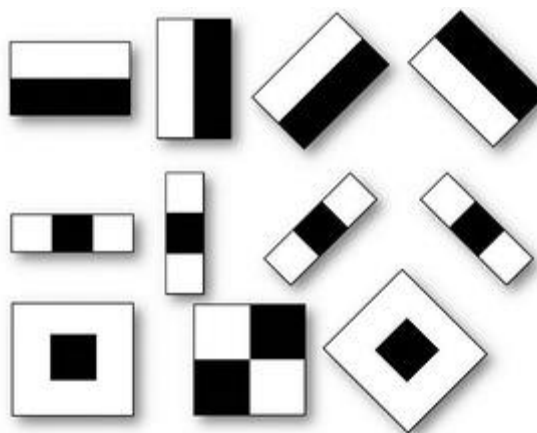
Алгоритм построения случайного леса может быть представлен в следующем виде:

Для $i = 1, 2, \dots, B$ (здесь B – количество деревьев в ансамбле) выполнить:

1. Сформировать бутстреп выборку S размера l по исходной обучающей выборке $D = \{\bar{x}_i, y_i\}_{i=1}^l$.
2. По бутстреп выборке S построить неусеченное дерево решений T_i с минимальным количеством наблюдений в терминальных вершинах равным n_{min} , рекурсивно следуя следующему подалгоритму:
 - a. из исходного набора n признаков случайно выбрать p признаков (в задачах классификации обычно $p \approx \sqrt{n}$);
 - b. из p признаков выбрать признак, который обеспечивает наилучшее расщепление;
 - c. расщепить выборку, соответствующую обрабатываемой вершине, на две подвыборки;

Признаки Хаара.

Признаки Хаара – это признаки цифрового изображения, представляющие собой прямоугольные области, состоящие из смежных прямоугольных подобластей. Своим названием они обязаны сходством с вейвлетами Хаара.



Значение двух-прямоугольного признака вычисляется, как разность между интегральными суммами пикселей в двух смежных прямоугольных подобластях. Для трех-прямоугольного признака значение вычисляется, как интегральная сумма

двух внешних подобластей, вычитаемая из суммы в центральной подобласти. Значение четырех-прямоугольного признака вычисляется, как разность между суммами диагональных пар подобластей.

Рассмотрим вектор признаков $\bar{x}_i = \{x_i^s\}_{s=1}^n$, координаты которого представляют значения n признаков Хаара, измеряемых в некоторой области, полученной в результате сегментации изображения. Значение признака x_i^s примем равным единице в случае, если разность интегральных сумм подобластей неотрицательна, и нулю – в противоположном случае.

Создание обучающей выборки.

Цифровая видеозапись представляет собой последовательность двумерных матриц растровых полутоновых изображений $\{I_t\}$, где t – номер кадра в последовательности, а каждый элемент матрицы $I(x, y)$ характеризует цвет соответствующего пикселя.

В качестве обучающей видеозаписи была выбрана видеозапись движения цели – радиоуправляемой машинки. Данная видеозапись включала в себя временные промежутки с изменением освещенности области наблюдения и временные промежутки с изменением размеров цели вследствие ее удаления от видеокамеры.

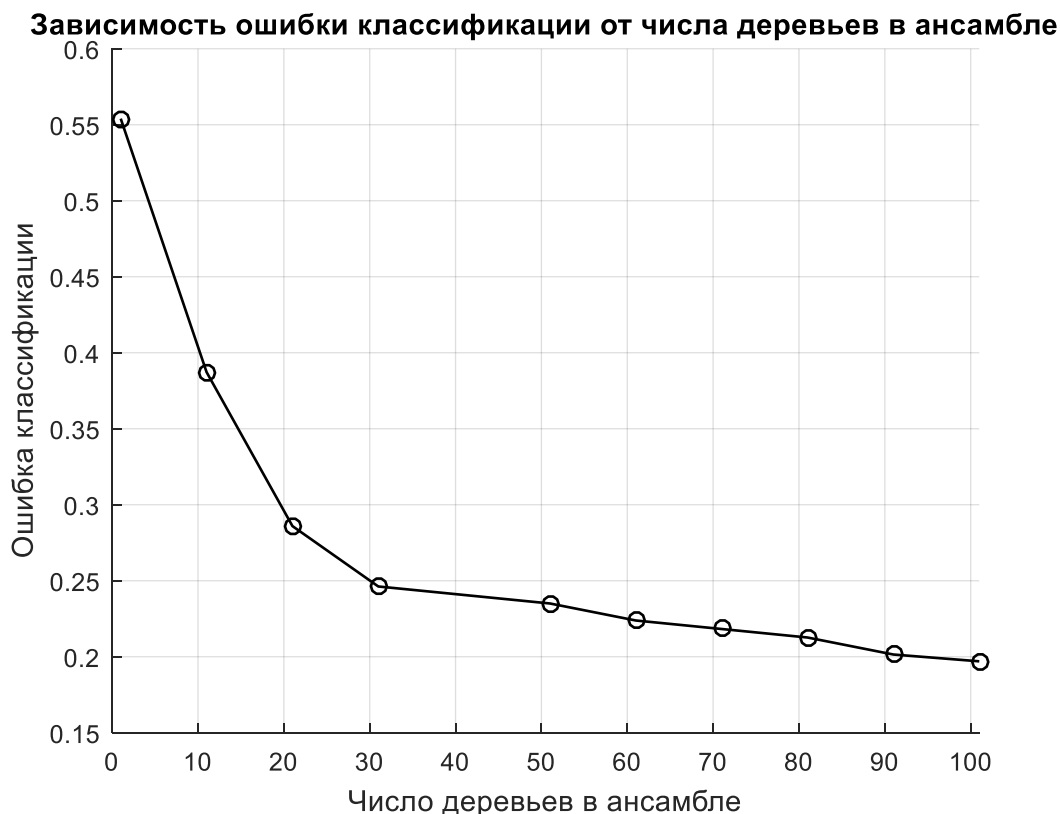
Для создания обучающей выборки $D = \{\bar{x}_i, y_i\}_{i=1}^l$ на основе видеозаписи был выполнен следующий алгоритм:

1. Для каждого кадра видеозаписи была выполнена сегментация изображения;
2. Для всех полученных сегментов были вычислены значения случайно сгенерированных векторов признаков Хаара \bar{x}_i ;
3. Значение метки класса y_i принималось равным единице, если \bar{x}_i – вектор признаков, соответствующий сегменту, охватывающему цель, нулю – в противоположном случае.

Тестирование классификатора.

Для создания тестовой выборки была использована видеозапись, сходная с обучающей. Данная видеозапись имела отличную от обучающей видеозаписи область наблюдения и включала в себя иные траектории движения цели.

По окончании тестирования был проведен анализ полученных результатов классификации и построен график ошибки.



Заключение.

Применение классификатора на основе случайного леса для решения задачи обнаружения объектов, присутствующих на цифровой видеозаписи, показало высокие результаты, стойкие к масштабированию объектов, изменению освещенности в области наблюдения, исчезновению объектов из области наблюдения. При этом точность классификации оказалась выше при использовании большого числа деревьев в ансамбле.

Существенным недостатком метода является то, что фазы обучения и тестирования разделены. Так как на практике обучающие данные чаще всего поступают последовательно (например, в приложениях трекинга), то в таких ситуациях алгоритм обучения классификатора должен работать on-line, т.е. «на лету». On-line обучение имеет большое число преимуществ перед off-line методом, например, требует меньший объем выделяемой памяти. Также on-line метод показывает высокие результаты в случае, если распределение, лежащее в основе обучающей выборки, меняется с течением времени.

Список использованных источников.

- [1] С.П. Чистяков, «Случайные леса: обзор». Труды Карельского научного центра РАН №1. 2013. С. 117-136.
- [2] L. Breiman, J.H. Friedman, R.A. Olshen, and C.J. Stone. Classification and Regression Trees. Wadsworth, Belmont, Ca, 1983.
- [3] A. Saffari, C. Leistner, J. Santner, M. Godec, and H. Bischof, «On-line Random Forests», 3rd IEEE ICCV Workshop on On-line Computer Vision, 2009.