

ИЗМЕНЕНИЕ ОСНОВНОГО ТОНА, ОСНОВАННОЕ НА
МАРКОВСКОЙ МОДЕЛИ

Л.С. БЕРНЕР, К.С. КУЗЬМИЧ

Институт электротехники им. В.М. Гугенца АН УССР
Киев, СССР 250007

АННОТАЦИЯ

Учет основных свойств траекторий основного тона на базе марковской модели позволяет существенно повысить надежность определения периода основного тона и предсказания его длины.

РЕЗЮМЕ

Учет основных свойств траекторий основного тона // // позволяет существенно повысить надежность определения периода основного тона и предсказания его длины.

- 1) Модель, описывающая зависимость длины основного периода от значения параметра s , позволяет существенно повысить надежность определения периода основного тона и предсказания его длины.
- 2) Модель, описывающая зависимость длины траекторий основного тона от значения параметра s , позволяет существенно повысить надежность определения периода основного тона и предсказания его длины.

Для построения модели необходимо знать значения параметров s и v в каждом из элементов траектории. Для этого на каждом элементе траектории необходимо определить значение параметра s и значения v в каждом из элементов траектории. Для этого на каждом элементе траектории необходимо определить значение параметра s и значения v в каждом из элементов траектории.

Учет основных свойств траекторий основного тона на базе марковской модели позволяет существенно повысить надежность определения периода основного тона и предсказания его длины.

спектра (например, путем обратной фильтрации системы линейного предсказания) и пропускать через фильтр низких частот, после чего вычислять нормированную автокорреляционную функцию $R(i)$, $0 \leq i \leq L$.

Для каждого из 12 сегментов найдем опорное значение T_s :

$$T_s = \arg \max_{1 \leq i \leq 12} R(i)$$

и в дальнейшем вместо автокорреляционной функции $R(i)$ будем рассматривать ее "дробную" $\bar{R}(s) = R(T_s)$, $1 \leq s \leq 12$. Кроме того, введем, говоря о периоде OT , шаг s , в который попадает этот период. Номер сегмента s является округленным значением OT ; уточненное значение, соответствующее округленному значению s , равно опорному значению T_s .

МОДЕЛЬ ОСНОВНОГО ТОНА НА
ОПРЕДЕЛЕННОМ ИНТЕРВАЛЕ АНАЛИЗА

Построение такой модели основано на предположении зависимости плотности совместных двумерных вероятностей состояний траектории автокорреляционной функции (АКФ) от значения признака s и периода OT $s = s(\bar{R}/s, v)$, где v соответствует признаку "тон", а s - признаку "длина".

Значение периода OT равномерно распределено, независимо от того же от параметров s и v . В этом случае допустимо считать, что s , v и OT независимы.

нать, например, к такому упрощению. Найдем максимум ПФ \bar{R} :

$$s^* = \arg \max_{1 \leq s \leq 12} \bar{R}(s)$$

Объявление s^* округленным значением периода OT соответствует корреляционному методу выделения OT ; мы, однако, сделаем некоторый шаг вперед по сравнению с этим методом, если учтем вероятности получения различных значений s^* при фактическом значении s : $P_1(s^*/s)$. Причем здесь можно ограничиться заданием вероятности правильной оценки $P_1(s/s)$ и вероятности ошибки, соответствующей удвоенному периоду OT , $P_1(s+4/s)$, полагая остальные ошибки равновероятными с вероятностью $P_1(s/s) = (1 - P_1(s/s) - P_1(s+4/s))/10$, $s \neq s, s \neq s+4$.

Величина $R^* = \bar{R}(s^*)$ в корреляционном методе служит для определения признака тон/шум; при этом решение принимается в результате сравнения этой величины с некоторым порогом. Мы, однако, определим вероятности значений R^* в зависимости от значения признака тон/шум: $p_2(R^*/v)$, сохранив больше информации, чем используется в корреляционном методе.

Теперь статистические модели вокализованного и невокализованного интервалов могут быть определены, например, следующим образом:

$$p(\bar{R}/s, v=1) = p(R^*, s^*/s, v=1) = P_1(s^*/s) \cdot p_2(R^*/v=1),$$

$$p(\bar{R}/s, v=0) = p(R^*, s^*/s, v=0) = C \cdot p_2(R^*/v=0),$$

где $C = 1/12$ - вероятность любого конкретного значения s^* на невокализованном интервале.

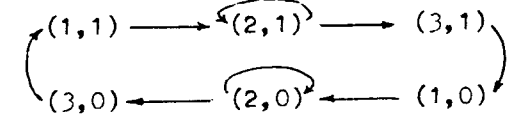
Более точная статистическая модель получится, если учесть зависимость между $\bar{R}(s)$ и $\bar{R}(s+4)$ при фактическом значении OT s . Возможны и другие варианты.

МАРКОВСКАЯ МОДЕЛЬ ТРАЕКТОРИЙ ОТ

Модель траекторий OT опирается на следующие их свойства. Во-первых, продолжительность вокализованных и невокализованных отрезков речевых сигналов ограничена снизу. Во-вторых, имеет место сильная зависимость между значениями периода OT на соседних вокализованных интервалах анализа. Аналогичная зависимость, хотя и в меньшей степени, наблюдается между значениями периода OT на интервалах, разделенных невокализованным участком сигнала, причем эта зависимость убывает с увеличением длительности разделяющего невокализованного участка.

Первое из этих свойств хорошо отражается следующей марковской моделью, генерирующей траектории признака тон/шум. Допустим, что минимальная длина вокализованных и невокализованных участков равна 3. Введем по три состояния $(1, v)$, $(2, v)$ и $(3, v)$ для каждого значения признака тон/шум $v=0, 1$.

Граф переходов модели имеет вид:



Состояния $(1, 1)$ и $(1, 0)$ являются начальными при порождении вокализованного и невокализованного участков соответственно, а $(3, 1)$ и $(3, 0)$ - конечными. Повторение состояний $(2, 1)$ и $(2, 0)$ позволяет генерировать участки вокализованной и невокализованной речи любой длины, большей или равной трем интервалам анализа. Для полного описания этой модели необходимо задать $P_v(2/2)$ - вероятность перехода по петле в каждое из этих повторяемых состояний, при этом $P_v(3/2) = 1 - P_v(2/2)$.

Второе из сформулированных выше свойств также может быть описано с помощью некоторой марковской модели; эта модель будет генерировать траектории значе-

ний периода ОТ. Модель содержит 12 состояний, соответствующих огрубленным значениям ОТ s . С достаточно хорошим приближением можно считать, что огрубленные значения периода ОТ s и s' на соседних интервалах анализа могут отличаться не более, чем на единицу: $|s - s'| \leq 1$.
 Задав вероятности переходов $P_3(s'/s)$, определяющие среднюю скорость изменения периода в огрубленной шкале, получаем марковскую модель траекторий периода ОТ.

Объединение модели траекторий признака тон/шум и модели траекторий периода ОТ приводит к совокупной марковской модели траекторий ОТ. Возможные переходы и их вероятности представлены ниже:

- $(s, 1, 1) \rightarrow (s', 2, 1), P_3(s'/s),$
- $(s, 2, 1) \rightarrow (s', 2, 1), P_3(s'/s) \cdot P_v(2/2),$
- $(s, 2, 1) \rightarrow (s', 3, 1), P_3(s'/s) \cdot P_v(3/2),$
- $(s, 3, 1) \rightarrow (s', 1, 0), P_3(s'/s),$
- $(s, 1, 0) \rightarrow (s', 2, 0), P_3(s'/s),$
- $(s, 2, 0) \rightarrow (s', 2, 0), P_3(s'/s) \cdot P_v(2/2),$
- $(s, 2, 0) \rightarrow (s', 3, 0), P_3(s'/s) \cdot P_v(3/2),$
- $(s, 3, 0) \rightarrow (s', 1, 1), P_3(s'/s),$

причем во всех случаях $|s - s'| \leq 1$.

Состояния с различными значениями s для невокализованных интервалов вводятся, чтобы отразить зависимость значений периода ОТ на вокализованных интервалах, примыкающих с двух сторон к невокализованному участку.

Введем следующие обозначения:

- $S_n = (s_n, i_n, v_n)$ - состояние модели на n -м интервале,
- R_n - ПАФ сигнала на n -м интервале;
- $p(R_n/S_n)$ - условная плотность вероятности параметров ПАФ в зависимости от значения признака тон/шум и периода ОТ;
- $P(S_{n+1}/S_n)$ - вероятность перехода из состояния S_n в состояние S_{n+1} ;
- $P(S_1/S_0)$ - распределение вероятностей начального состояния.

В качестве $P(S_1/S_0) = P(S_1)$ естественно взять распределение с равными вероятностями для всех невокализованных состояний и нулевыми вероятностями для вокализованных.

Теперь задача выделения ОТ сводится к задаче определения состояний построенной марковской модели, наилучшим образом соответствующих наблюдаемому речевому сигналу. Возможны различные постановки этой задачи.

Можно, например, поставить задачу отыскания наиболее правдоподобной траектории ОТ по всей реализации речевого сигнала. Для этого следует максимизировать апостериорную вероятность траектории:

$$P(S_1, \dots, S_N / R_1, \dots, R_N) = \prod_{n=1}^N P(S_n / S_{n-1}) \cdot p(R_n / S_n) / p(R_1, \dots, R_N). \quad (1)$$

Эта задача решается следующим алгоритмом динамического программирования:

$$I(S', n) = \operatorname{argmax}_{S \in Q(S')} (F(S, n-1) + g(R_n, S, S')) \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} 2 \leq n \leq N, \quad (2)$$

$$F(S', n) = \max_{S \in Q(S')} (F(S, n-1) + g(R_n, S, S'))$$

$$S_N^* = \operatorname{argmax}_S F(S, N),$$

$$S_{n-1}^* = I(S_n^*, n), \quad n = N, \dots, 2.$$

В этих формулах:

- $g(R_n, S, S') = \ln P(S'/S) + \ln p(R_n/S')$,
- $F(S, 1) = \ln P(S) + \ln p(R_1/S)$,
- $Q(S')$ - множество состояний, из которых возможен переход в состояние S' .

Другим важным вариантом постановки задачи является минимизация вероятности ошибки выделения ОТ на m -м интервале анализа по фрагменту реализации из первых $m+k$ интервалов. Для решения этой задачи, как известно, следует определить, какое из интересующих нас событий имеет максимальную апостериорную вероятность. Мы будем рассматривать 13 событий: одно, означающее невокализованность m -го ин-

тервала, и еще 12, соответствующих двенадцати возможным значениям s_m периода на m -м интервале, $1 \leq s_m \leq 12$.

Для вычисления вероятностей этих событий следует просуммировать апостериорные вероятности состояний модели на m -м интервале, составляющих эти события. Апостериорная вероятность состояния $S_m = (s_m, i_m, v_m)$ определяется следующей формулой:

$$P(s_m, i_m, v_m / R_1, \dots, R_{m+k}) = P(S_m / R_1, \dots, R_{m+k}) = \sum_{S_1, \dots, S_{m-1}} (1/A) \cdot \prod_{n=1}^{m+k} P(S_n / S_{n-1}) \cdot p(R_n / S_n),$$

где $A = p(R_1, \dots, R_{m+k})$. Суммирование производится по всем состояниям на первых $m+k$ интервалах, кроме m -го.

Теперь можно вычислить апостериорные вероятности огрубленных значений периода:

$$P(s_m, v_m = 1 / R_1, \dots, R_{m+k}) = \sum_{i_m} P(s_m, i_m, v_m = 1 / R_1, \dots, R_{m+k}), \quad 1 \leq s_m \leq 12,$$

и апостериорную вероятность того, что m -й интервал невокализован:

$$P(v_m = 0 / R_1, \dots, R_{m+k}) = \sum_{i_m, s_m} P(s_m, i_m, v_m = 0 / R_1, \dots, R_{m+k})$$

Выбор наибольшего из этих 13 чисел определяет оптимальное значение признака тон/шум и периода ОТ.

Итак, предлагаемый метод применим и в случае, когда траекторию ОТ можно определять после завершения ввода всей реализации, и в случае оперативного оценивания.

В первом случае ошибки в значениях ОТ практически исключаются. Во втором случае вероятность ошибки после 4-5 вокализованных интервалов также довольно мала.

УПРОЩЕННАЯ СХЕМА: МОДИФИКАЦИЯ КОРРЕЛЯЦИОННОГО МЕТОДА

Ниже предлагается упрощенный вариант,

являющийся фактически усовершенствованной модификацией корреляционного метода, поскольку в нем принятие решений осуществляется на основании сумм коэффициентов автокорреляции на траектории ОТ.

Сначала для каждого состояния модели $S_n = (s_n, i_n, v_n)$ определяется вес: $d(R_n, S_n) = \begin{cases} R_n(s_n), & \text{если } v_n = 1, \\ A - \max_s R_n(s), & \text{если } v_n = 0. \end{cases}$

где A - эмпирическая константа. Затем для каждой пары состояний, между которыми возможен переход определяется величина:

$$g(R_n, S_{n-1}, S_n) = \begin{cases} d(R_n, S_n), & \text{если } s_n = s_{n-1}, \\ d(R_n, S_{n-1}) - D, & \text{если } s_n = s_{n-1} \pm 1, \end{cases}$$

где D - также эмпирическая константа.

Теперь ставится задача поиска такой последовательности состояний $S_n, 1 \leq n \leq N$, связанных допустимыми переходами, которая имела бы наибольшую сумму:

$$F = \sum_{n=1}^N g(R_n, S_{n-1}, S_n).$$

Эта задача решается алгоритмом (2).

Описанный упрощенный вариант проверялся на тестовом материале из 9 фраз, произнесенных двумя дикторами-женщинами и одним диктором-мужчиной; частота дискретизации 16 кГц, порядок обратного фильтра 14, частота среза НЧ-фильтра 1,25 кГц, длина интервала анализа 30 мс, шаг - 20 мс. Общий объем речевого сигнала около 800 интервалов анализа. Эксперименты показали высокую надежность алгоритма: не было ни одного ошибочного значения периода ОТ. Были только некоторые сомнения относительно признака тон/шум на стыках вокализованных и невокализованных участков.

ЗАКЛЮЧЕНИЕ

Предварительные эксперименты говорят о перспективности применения марковской модели для надежного выделения ОТ. При этом роль модели будет тем больше, чем менее информативным является используемое первичное описание, что имеет место, например, в случае зашумленной речи.