

УДК 621.391:681.3.06.

АНАЛИЗ И СИНТЕЗ РЕЧИ НА ОСНОВЕ Z-ОПИСАНИЯ

В.С. Лозовский

1. А п п р о к с и м а ц и я

В работе [1] предложено для звуков речи с тональным возбуждением аппроксимировать отклик речеобразующего тракта на участке смыкания голосовых связок дробно-рациональной функцией вида:

$$\hat{H}(z) = \frac{N(z)}{D(z)} = \frac{a_0 + a_1 z^{-1} + \dots + a_n z^{-n}}{1 + b_1 z^{-1} + \dots + b_m z^{-m}}, \quad (1)$$

где $z = e^{j\omega}$. Выражение (1) представляет собой передаточную функцию дискретной линейной динамической системы, полюса которой соответствуют формантам исследуемого речевого сигнала. При этом не выдвигается требование наблюдаемости всех состояний речеобразующей системы, поскольку для слухового восприятия важны не собственные частоты реальной физической системы, а лишь те колебания, которые присутствуют в конкретном, воспринимаемом ухом сигнале.

Для сегментов речи с шумовым возбуждением предлагается воспользоваться методами идентификации линейных динамических систем в стохастическом случае [2]. Итак, в каждый дискретный момент времени t сигнал на выходе аппроксимирующей системы, описываемой выражением (1), имеет вид:

$$y_t = -b_1 y_{t-1} - \dots - b_m y_{t-m} + a_0 x_t + \dots + a_n x_{t-n}, \quad (2)$$

где x_t и y_t - отсчеты соответственно входного и выходного сигналов.

Предполагается, что при произнесении фриктивных звуков на вход резонансной системы речеобразующего тракта подан белый шум. В модели ему соответствует дискретный белый шум - сигнал с некоррелированными соседними отсчетами, значения которых распределены равномерно, например, в интервале $(1, -1)$.

Рассмотрим систему, описываемую уравнением $1/D(z)$, полагая, что на её вход поступает сигнал:

$$u_t = \sum_{j=0}^t a_j x_{t-j}. \quad (3)$$

Множество $\{u_t\}$ представляет собой так называемый линейный шум, в котором соседние отсчеты коррелированы. Из (3) следует, что с увеличением интервала между рассматриваемыми значениями u_t корреляция между ними уменьшается; отсчеты, отстоящие друг от друга на интервал, не меньший $n+1$, некоррелированы.

Перепишем (2) в виде:

$$-\sum_{j=1}^m b_j y_{t-j} = y_t - u_t \triangleq f_t. \quad (4)$$

Итак, f_t - это значения функции, аппроксимации которой в наилучшем смысле мы будем добиваться, выбирая b_j . Функция $\{f_t\}$ включает в себя наблюдаемую составляющую y_t и линейно-шумовую составляющую $-u_t$, нам не известную. Выбрав отсчеты f_t с шагом $n+1$, мы переходим от линейного к белому шуму.

Составим систему из m уравнений (4), полагая $f_t = y_t$. Восстановив по полученным b_j аппроксимирующую функцию, обнаруживаем, что на каждом уравнении мы допустили ошибку u_t , которая есть несмещенный белый шум. При большем числе исходных уравнений к ошибкам $|u_t|$ будут добавляться ошибки средне-квадратической аппроксимации, которые также носят характер белого шума. Таким образом, можно полагать, что полученная описанным путем оценка b_j будет несмещенной.

Подчеркнем, что в случае стохастической аппроксимации определяются лишь коэффициенты $D(z)$.

При решении задач x -идентификации, о которых шла речь выше, возникает необходимость в решении переопределенных систем вида:

$$-\sum_{j=1}^n b_j y_{t-j} = y_t, \quad t=1, \dots, m; \quad m \geq n. \quad (5)$$

Можно показать, что подобные системы являются плохо обусловленными; поэтому корректный выбор вычислительной процедуры для их решения является весьма важной задачей. Использование известной трансформации Гаусса позволяет преобразовать матрицу системы (5) к квадратной симметричной положительно-определенной форме; для решения полученной системы весьма эффективен метод квадратных корней. Однако точность полученного решения оказывается недостаточной. Дело в том, что уже на этапе составления нормальной системы при вычислении скалярных произведений допускаются ошибки округления, которые в условиях слабой определенности (5) иногда приводят даже к отрицательно определенным системам. Конечно, можно вычислять скалярное произведение и даже полностью решать полученную систему с удвоенным числом знаков мантиссы, но, как показал эксперимент, в рассматриваемом случае этого недостаточно. Хорошие результаты дает использование матриц отражения Хаусхолдера [3] для триангуляризации исходной системы (5); этот алгоритм с некоторыми модификациями и был принят в качестве рабочего.

Обсуждая методы решения переопределенных систем и делая определенные выводы относительно физической стороны исследуемых процессов, нельзя обойти вопрос о корректности решаемой задачи. В работе [4] показано, что при задании исходной информации (значения y_t, y_{t-j} в (5)) с некоторым приближением задача определения нормального решения системы (5) некорректна в смысле Адамара. Некорректность проявляется в больших вариациях решения b_j при малых вариациях исходных данных. Прежде чем приступить к регуляризации исходной системы, которая позволяет до некоторой степени "выправить" решение, остановимся на физическом смысле некорректности нашей задачи.

Итак, нас интересует построение модельной системы, отклик которой с нужной точностью аппроксимирует анализируемый речевой сигнал. Ясно, что существует множество систем, отвечающих этому требованию. Разница между откликами этих систем и анализируемым сигналом, поскольку мы используем методы среднеквадратической аппроксимации, носит характер белого шума. Некорректность нашей задачи проявляется в том, что близкие в указанном смысле во временной области системы могут заметно отличаться по положению полюсов, т.е. формант. Это приводит к выводу, что, с одной стороны, при вычислении положения формант для со-

седних сегментов речи возможны их резкие скачки, а, с другой стороны, эти нестационарности отражают физику явления в рамках рассматриваемой модели и на качестве синтеза по полученным параметрам сказываться не должны. Здесь уместно снова подчеркнуть, что наша задача — построить модельную систему, близкую исследуемой во временной области, а не определить передаточную функцию речеобразующего тракта. По-видимому, вторая задача не может быть строго решена без получения более подробной информации об источнике возбуждения ^{*)}, чем мы можем почерпнуть из исследования лишь осциллограммы речи.

Теперь о регуляризации. Потребность в ней может возникнуть, например, из соображений уменьшения ширины полосы сигнала при передаче значений коэффициентов аппроксимации для синтеза речи на приемном конце. Суть известных методов регуляризации задачи линейной алгебры [4] сводится к тому, что минимизируемый функционал (в нашем случае сумма квадратов ошибок аппроксимации) дополняется функционалом, зависящим от разности между искомым решением и некоторым фиксированным вектором, например, нулевым, или же, в нашем случае, равным вектору неизвестных коэффициентов, определенному для предыдущего сегмента.

Выведем необходимые соотношения. Пусть исходная система уравнений имеет вид:

$$a_{i:m, 1:n} \times z_{1:n} = b_{i:m}, \quad m \geq n. \quad (6)$$

Здесь одна индексная пара соответствует вектору-столбцу, а две — матрице. Минимизируемый функционал имеет вид:

$$F = \| a_{i:m, 1:n} \times z_{1:n} - b_{i:m} \|^2 + \alpha \| z_{1:n} - z_{1:n}^0 \|^2 \quad (7)$$

$$= \sum_{i=1}^m \left(\sum_{j=1}^n a_{ij} z_j - b_i \right)^2 + \alpha \sum_{j=1}^n (z_j - z_j^0)^2,$$

где $z_{1:n}^0$ — некоторый фиксированный вектор, а α — параметр регуляризации.

$$\frac{\partial F}{\partial z_k} = 2 \sum_{i=1}^m \left(\sum_{j=1}^n a_{ij} z_j - b_i \right) a_{ik} + 2\alpha (z_k - z_k^0),$$

$$k = 1, \dots, n.$$

Полученная квадратичная форма имеет лишь один экстремум:

^{*)} Строго говоря, при этом надо ещё потребовать наблюдаемости всех состояний тракта.

это - минимум, поскольку форма существенно положительна.

$$\sum_{j=1}^n \sum_{i=1}^m a_{ki}^T a_{ij} x_j + \alpha x_k = \sum_{i=1}^m a_{ki}^T b_i + \alpha x_k^0, \quad (8)$$

$k = 1, \dots, n.$

Выражение (8) - регуляризованная система нормальных уравнений для решения задачи среднеквадратической аппроксимации, большей наглядностью обладает матричная форма записи:

$$(a_{1:n, 1:m}^T a_{1:m, 1:n} + \alpha J_{1:n, 1:n}) x_{1:n} = a_{1:n, 1:m}^T b_{1:m} + \alpha x_{1:n}^0 \quad (9)$$

здесь J - единичная матрица.

В связи с тем, что принятый нами алгоритм позволяет обойти составление нормальных уравнений, перейдем к записи исходной переопределенной системы. Предварительно перепишем (9) в блочном виде:

$$(a_{1:n, 1:m}^T \sqrt{\alpha} J_{1:n, 1:n}) \begin{pmatrix} a_{1:m, 1:n} \\ \sqrt{\alpha} J_{1:n, 1:n} \end{pmatrix} x_{1:n} = (a_{1:n, 1:m}^T \sqrt{\alpha} J_{1:n, 1:n}) \begin{pmatrix} b_{1:m} \\ \sqrt{\alpha} x_{1:n}^0 \end{pmatrix} \quad (9')$$

Отсюда легко получить исходную систему:

$$\begin{pmatrix} a_{1:m, 1:n} \\ \sqrt{\alpha} J_{1:n, 1:n} \end{pmatrix} x_{1:n} = \begin{pmatrix} b_{1:m} \\ \sqrt{\alpha} x_{1:n}^0 \end{pmatrix}, \quad (10)$$

которая решается с использованием триангуляризации Хаусхолдера.

Дальнейшего повышения точности решения задачи x -идентификации можно добиться с помощью весовых функций. Весовая функция - это массив коэффициентов $w_{1:m+n}$, на которые должны умножаться уравнения исходной системы в соответствии с достоверностью каждого. Значения последних n весов уже определены - $\sqrt{\alpha}$; на выборе же первых m коэффициентов мы сейчас остановимся.

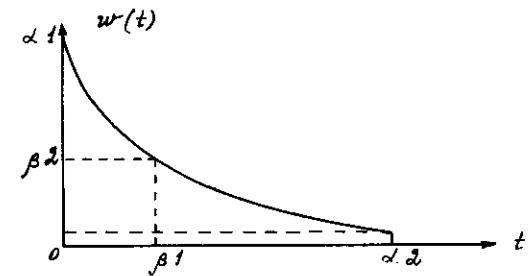
Прежде всего, трудно рассчитывать на то, что в реальном речевом сигнале при анализе озвученных сегментов удастся четко фиксировать границы интервала смыкания голосовых связок. Поэтому первым и последним уравнениям (10) (без учета регуляризу-

ющей добавки) следует придавать меньший вес. Подобное подавление граничных уравнений с помощью гладкой весовой функции будет также способствовать сглаживанию вектора решений при выполнении аппроксимации в скользящем режиме при поиске интервала аппроксимации. Специфика нашей задачи заключается в том, что аппроксимируемые функции вычисляются рекуррентно. При этом в связи с конечной точностью вычислений происходит накопление ошибок. Поэтому к концу интервала исходным уравнениям следует придавать меньший вес.

Практически использовались функции четырех типов:

Тип 1. Значение функции принимается равным постоянной величине на всем интервале аппроксимации.

Тип 2. Экспонента.



При этом задаются:

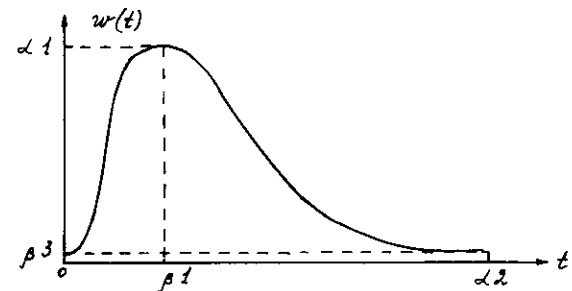
- α_1 - значение ординаты в начале,
- β_1 и β_2 - абсцисса и ордината опорной точки,
- β_3 - ордината на конце интервала аппроксимации: α_2 .

Аналитически функция записывается в виде:

$$w(t) = \alpha_1 \cdot e^{-\lambda_1 \left(\frac{t}{\alpha_2}\right)^{\lambda_2}}, \quad (II)$$

где необходимые коэффициенты λ_1 и λ_2 легко определяются.

Тип 3. Сдвинутая косинусоида.



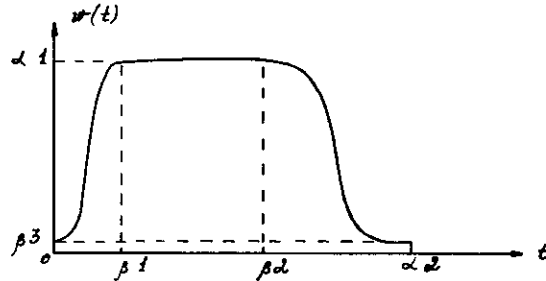
Функция имеет вид:

$$w(t) = \frac{1 + \beta_3 - \cos\left(2\pi\left(\frac{t}{\alpha_2}\right)^{\lambda_1}\right)(1 - \beta_3)}{\alpha_2} \cdot \alpha_1, \quad (12)$$

$$\lambda_1 = \ln 0.5 / \ln 31,$$

β_1 и β_3 нормированы к единице.

Тип.4. Комбинированная функция.



Здесь передний фронт строится по формуле:

$$w_{np}(t) = \frac{1 + \beta_3 - \cos\left(\frac{\pi t}{\beta_1 \cdot \alpha_2}\right)(1 - \beta_3)}{\alpha_2} \cdot \alpha_1, \quad (13)$$

а задний:

$$w_{zp}(t) = \frac{1 + \beta_3 + \cos\left(\frac{\pi(t - \alpha_2 \cdot \alpha_2)}{(1 - \beta_3) \cdot \alpha_2}\right)(1 - \beta_3)}{\alpha_2} \cdot \alpha_1, \quad (14)$$

на интервале (β_1, β_2) значение весовой функции постоянно и равно α_1 .

Варьируя параметр α_1 весовых функции, удобно масштабировать исходную систему, что способствует дополнительному повышению точности вычислений.

2. Поиск интервала

Для озвученных сегментов речи не меньше, чем сама аппроксимация, значение имеет поиск интервала, соответствующего смыканию голосовой щели, т.е. свободным колебаниям изучаемой си-

стемы. при этом попутно с поиском интервала можно было бы определять и период основного тона (OT), но проще это сделать отдельно.

Итак, исследуемый отрезок речи разбивается на сегменты по 20 мсек. Если энергия сигнала на сегменте меньше пороговой, аппроксимация не выполняется и сегменту присваивается код: П - пауза.

Для поиска OT на сегментах используется модификация автокорреляционного метода [5], в результате применения которого вычисляется индекс аperiodичности сегмента (ИА). Для строго периодического сигнала ИА = 0. Устанавливается величина порога (~0,5). Считается, что при больших значениях ИА мы имеем дело с шумовым возбуждением (код типа сегмента: Ш), при меньших - с тональным (код: Т). В последнем случае вычисляется частота OT. Для Ш - сегментов сразу выполняется аппроксимация сигнала, пропущенного через фильтр высоких частот, по способу, описанному в разделе I. В Т - сегментах предварительно ищется оптимальный интервал аппроксимации.

Поиск выполняется методом пробной аппроксимации в скользящем режиме. Предварительно сигнал пропускается через фильтр низких частот κ для понижения частоты квантования вдвое. Этот прием облегчает работу алгоритма аппроксимации и сокращает машинное время. В качестве начальных точек для пробной аппроксимации выбираются не все точки, а лишь соответствующие моментам перехода сигнала через нуль, а также дополнительные точки с тем, чтобы расстояние между соседними не превышало задаваемого фиксированного порога.

При исследовании методов поиска было испробовано несколько алгоритмов.

Интервал постоянства коэффициентов. Длина интервала при пробной аппроксимации выбирается заведомо меньше ожидаемой. Производится аппроксимация для выбранной последовательности начальных точек и сравниваются значения коэффициентов b_j (1), полученные для каждого участ-

*) Параметры рекурсивных ФНЧ и ФВЧ определялись по программе [6]; на частоте квантования 16 кгц их данные: ФНЧ 0 - 3700 гц, аппроксимация по Чебышеву пятого порядка, неравномерность на вершине: 0.04. ФВЧ моделируется полосовым фильтром: $f_{cp} = 4000$ гц, $(2af)_{0.707} = 6000$ гц, аппроксимация по Чебышеву, порядок НЧ прототипа: 5, неравномерность: 0.04.

ка. При закрытой голосовой щели система совершает свободные колебания, параметры её должны изменяться мало; с открытием щели меняются резонансные характеристики (главным образом, добротность), в систему поступает энергия извне, и параметры аппроксимации должны проявить заметную нестационарность. На практике же ситуация осложняется. Как уже говорилось в разделе I, решаемая задача анализа является некорректной, т.е. даже в условиях фактического постоянства физических параметров исследуемой системы им будут соответствовать заметно различающиеся значения параметров аппроксимации. Конечно, можно регуляризовать задачу, но и в этом случае нужно алгоритмически сформулировать определение, что считать стационарным участком, а что нет. Был исследован метод поиска интервала постоянства с использованием таксономии на потенциальных функциях, но сложность алгоритма и время вычислений не оправдывались качеством полученных результатов.

Обратная фильтрация. Сигнал ОТ можно выделить в чистом виде, если пропустить речь через фильтр, нули которого в точности соответствуют полюсам речевого тракта. Здесь получается заколдованный круг: ведь мы как раз ищем участок аппроксимации, чтобы по нему определить резонансные параметры тракта. Применение рассматриваемого метода упирается в выбор функционала оценки качества обратной фильтрации. А это не просто и требует дополнительных затрат машинного времени при вычислениях.

Взвешенная ошибка аппроксимации и экстраполяции. Этот метод требует небольших дополнительных затрат машинного времени, обладает достаточной надежностью, в связи с чем он и был принят в качестве рабочего ^{*}). Суть его в следующем.

Для аппроксимации реального сигнала выбрана линейная динамическая модель. При этом качество аппроксимации будет тем выше, чем точнее выполняются требования линейности для исследуемого сигнала. Кроме этого, требуется, чтобы аппроксимирующая

^{*}) Первоначально [7] в качестве критерия использовался минимум ошибки аппроксимации, являющийся частным случаем описываемого здесь алгоритма.

система была пассивной. Таким образом, выбирая участок, на котором ошибка аппроксимации минимальна, мы автоматически получаем участок, который с наибольшим основанием может быть принят в качестве отклика линейной динамической системы. Точность определения параметров системы, а также неизменность их во времени должны подтверждаться небольшим значением ошибки экстраполяции в течение нескольких отсчетов сигнала, примыкающих к концу исследуемого интервала.

Практически описываемый алгоритм реализуется следующим образом. Выбирается для взвешивания исходной системы уравнений функция типа 3 (раздел I) с малым значением β . Производится определение коэффициентов b_j и уточнение a_j по всей длине рассматриваемого интервала [I]. В связи с малой величиной β значение весовой функции $w(t)$ к концу интервала весьма невелико, т.е. последние уравнения исходной переопределенной системы практически уже не участвуют в определении коэффициентов. Таким образом, к концу интервала $\alpha/2$ уже имеет смысл говорить не об аппроксимации, а об экстраполяции исследуемого сигнала. Ошибка вычисляется как разность между исходной и аппроксимирующей функциями на всем интервале $0 - \alpha/2$ и без взвешивания нормируется по интенсивности сигнала. Запоминается значение этой ошибки для всех исследованных на данном сегменте интервалов и находится её минимум. Найденная точка считается началом оптимального интервала для точной аппроксимации. Степень аппроксимирующей функции определяется естественным образом: исходная система уравнений строится для максимальной степени, представляющей интерес, например, восьмой. В процессе триангуляризации системы контролируется норма ведущего столбца. Когда она станет меньше заданного порога (т.е. преобразующая матрица отражений вырождается в единичную), процесс обрывается и действительная степень аппроксимации полагается равной числу преобразованных столбцов.

В самом алгоритме аппроксимации не заложено требование пассивности. Поэтому после того как будет выполнена аппроксимация сигнала на выбранном участке, вычисляется [8] положение полюсов аппроксимирующего полинома. Если хотя бы один полюс выходит за пределы единичной окружности в z -плоскости, выбранный интервал считается активным, а аппроксимация повторяется для следующего по значению ошибки интервала аппроксимации.

3.0 некоторых аспектах Z - анализа

Результатом решения задачи идентификации параметров речевого сигнала методами рациональной z -аппроксимации является представление вида $N(z)/D(z)$ для озвученных 20 мсек - сегментов и $a_n/D(z)$ - для шумовых. Вопрос об использовании $N(z)$ в первом случае не вполне ясен, в частности, потому, что не всегда в передаточной функции тракта нули отсутствуют, а отличить нули источника от нулей тракта пока не всегда удается; с другой стороны, даже выявив нули источника, мы таким образом получим информацию лишь о конечном участке импульса возбуждения. Таким образом, остановимся пока на том, что сегменты обоих типов будут характеризоваться лишь полиномами $D(z)$.

Корни полинома $D(z)$ соответствуют формантам [1], т.е. в принципе отсюда можно определить и частоты, и ширину полос найденных формант. Возникает вопрос, удобно ли пользоваться формантным описанием, в частности, для передачи по каналам связи. В этой связи следует заметить, что до сих пор не известен строгий и однозначный алгоритм соотношения спектральных максимумов (полусов) и формант [9]. Если в определенной ситуации и удастся найти первые четыре форманты, то через несколько десятков миллисекунд какая-либо из них может исчезнуть под влиянием нуля источника возбуждения, форманты могут слиться, могут появиться субформанты, назальные, ложные и т.п. При идентификации формант, например, по картинкам "видимой речи" рекомендуется [9] иметь в виду "критерий непрерывности". Однако этот критерий весьма проблематичен, поскольку некоторые полости при артикуляции открываются, закрываются, делятся на две и т.д. Большая часть имеющихся в литературе критериев рассчитана на то, что почти все фонемы уже распознаны, расчленены, изучены переходы; только после этого, опираясь на физическую модель тракта, можно пытаться соотнести найденные максимумы с номерами соответствующих формант.

В инженерной практике при построении формантных систем частотный диапазон делится на фиксированные, реже перестраиваемые области, которые и соотносятся с формантами. Возможны выходы формант за границы выделенных областей. Ясно, что первый путь очень неформален, а второй излишне формализован, и построить корректную автоматическую схему формантного анализа по

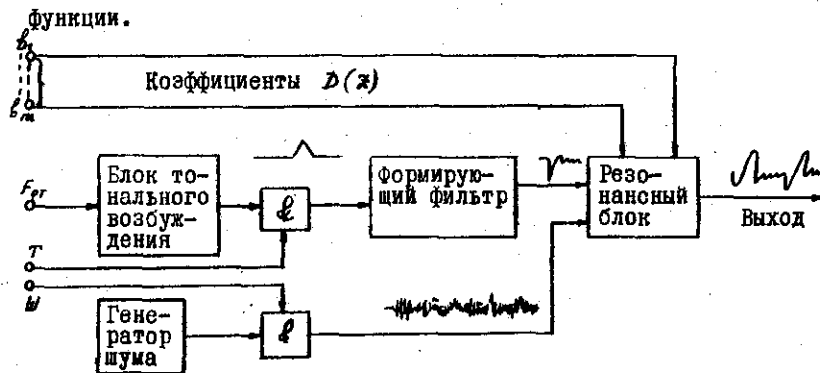
сигналу от микрофона на их основе очень непросто.

С другой стороны, методы z -анализа, описываемые в настоящей работе, будучи по существу описанием, близким к формантному, позволяют избежать нетривиальной задачи индексации формант при анализе путем передачи по каналу либо коэффициентов $D(z)$ непосредственно, либо подвергнув их предварительно ортогональной кодировке. Это дает возможность избежать определенных затруднений при конструировании синтезатора (параллельная или последовательная схема, управление шириной полос, амплитуд формант, частотная коррекция и т.п.). В нашем случае синтезатор по коэффициентам $D(z)$ непосредственно во временной области восстанавливает исходный сигнал традиционными методами рекурсивной фильтрации.

4. z - синтез

Ниже приводится блок-схема предлагаемого синтезатора. В качестве управляющих используются следующие сигналы:

- тональное (T) или шумовое (Ш) возбуждение.
- значение частоты основного тона F_{0T} .
- значения m коэффициентов $D(z)$ аппроксимирующей функции.



Блок тонального возбуждения генерирует треугольные импульсы. Назначение формирующего фильтра - придать возбуждающему импульсу форму, близкую к реальной. Критерий настройки - естественность синтезируемой речи. Фильтр включает от нуля до пяти резонаторов, которые могут использоваться в качестве полюсов и нулей. Резонансный блок реализует алгоритм рекурсивной фильтра-

ции :

$$y_t = x_t - b_1 y_{t-1} - \dots - b_m y_{t-m}, \quad (15)$$

где x_t и y_t - отсчеты функций на его входе и выходе соответственно, а b_1, \dots, b_m - коэффициенты аппроксимации $D(x)$.

Генератор шума выдает последовательность псевдослучайных чисел, равномерно распределенных в диапазоне $(-0,5, 0,5)$.

Простота схемы синтезатора избавляет от необходимости описывать принцип его работы.

5. Результаты эксперимента

Работа анализатора и синтезатора моделировалась на ЭВМ БЭСМ-6. С помощью специальной программы речевой сигнал вводился в машину через девятиразрядный аналого-цифровой преобразователь. Программа фиксировала начало и конец слова или фразы и переписывала сигнал в цифровом виде на ленту БЭСМ-6 для последующей обработки.

На рис. 1 иллюстрируется работа алгоритма x -аппроксимации. Основные параметры программы анализа:

Порядок $D(x)$ при поиске интервала аппроксимации: $m = 6$.

При поиске коэффициенты b_j не уточнялись, а a_j определялись по всему интервалу.

Длина пробного интервала аппроксимации при поиске оптимального интервала $l_{as} = 0,9 T_{ор}$.

При поиске интервала исследуется участок сигнала длиной $l_s = 2,1 T_{ор}$.

Порог по нормам столбцов матрицы при триангуляризации $\epsilon_s = 10^{-6}$.

Порядок $D(x)$ при точной аппроксимации $m = 8$.

Осуществлялись две итерации уточнения b_j .

Длина интервала аппроксимации $l_a = 0,6 T_{ор}$.

Применена регуляризация: в систему включен модуль разности между вектором неизвестных коэффициентов и вектором коэффициентов $D(x)$, полученных на предыдущем сегменте. Параметр регуляризации $\alpha = 0,06$.

Весовая функция при поиске начала аппроксимации: Тип 3;

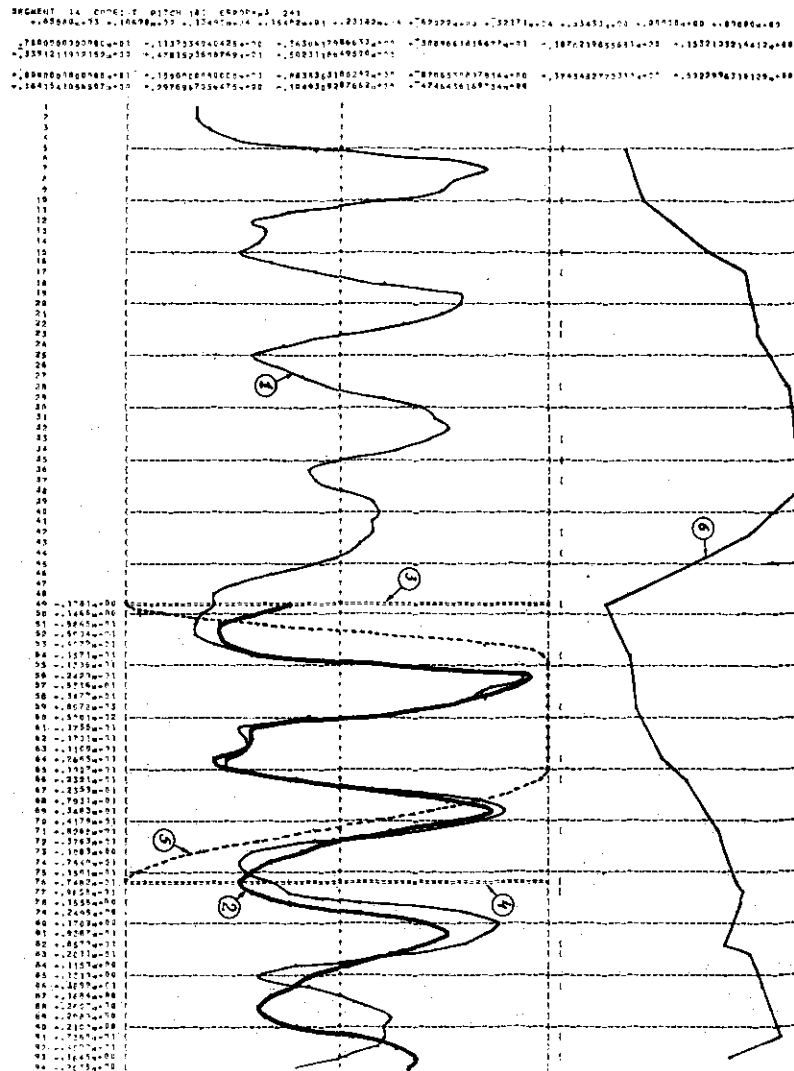


Рис. 1

положение максимума : $\beta 1s = 0.2$; значение в конце интервала : $\beta 3s = 0.002$.

Весовая функция при точной аппроксимации : Тип 4, $\beta 1=0.2$, $\beta 2 = 0.6$, $\beta 3 = 0.006$.

Кривая I-1*) - речевой сигнал после фильтрации до 3700гц с частотой квантования 8 кгц. Это "а" в слове "вада". Минимум кривой I-6, представляющей график взвешенной ошибки аппроксимации / экстраполяции, принят в качестве начала (I-3) оптимального интервала аппроксимации. I-5 - график взвешивающей функции при аппроксимации, I-2 - аппроксимирующий сигнал. Регуляризация существенного влияния на качество аппроксимации во временной области не оказывает. I-4 - конец интервала аппроксимации. Видно, что интервал закрытия голосовой щели заметно превышает интервал аппроксимации - открытие начинается в районе 85-89 отсчетов. Частота основного тона: 180 гц.

На рис. 2 фигурирует отрезок сигнала, представляющего конец второй фонемы "а" в том же слове "вада". Параметры анализа совпадают с приведенными выше, за исключением двух: порядок

$D(x)$ как при поиске интервала, так и при самой аппроксимации принят равным четырем.

После завершения анализа печатается картинка видимой речи обработанного объекта, полученная с помощью программы "FAST-60" [10]. Частоты каналов анализатора указаны в таблице I. С целью контроля для каждого сегмента, в котором была осуществлена аппроксимация, вычислялось положение полюсов передаточной функции. Эти полюса нумеровались в порядке возрастания частот, после чего были проведены линии, соединяющие соответствующие полюса для соседних сегментов. Следует подчеркнуть, что эта процедура выполнялась чисто интуитивно и, так же как и отпечатанные номера полюсов, не соответствует тем "истинным" формантам, которые определяются конфигурацией тракта диктора. Однако сходство полученных картинок и того, что принято считать формантным рисунком, особенно на рис. 4, очевидна. Число коэффициентов

$D(x)$ при аппроксимации было положено равным шести. На рис. 3 представлена нерегуляризованная картина; на рис. 4 - результаты анализа с параметром регуляризации 0.1. В нижней части рис. 3 и 4 расположены шесть колонок. Это :

- код типа сегмента (П, Т или Ш),
- номер сегмента (20 мсек - участка),

*) Первая цифра - номер рисунка, вторая - номер кривой.

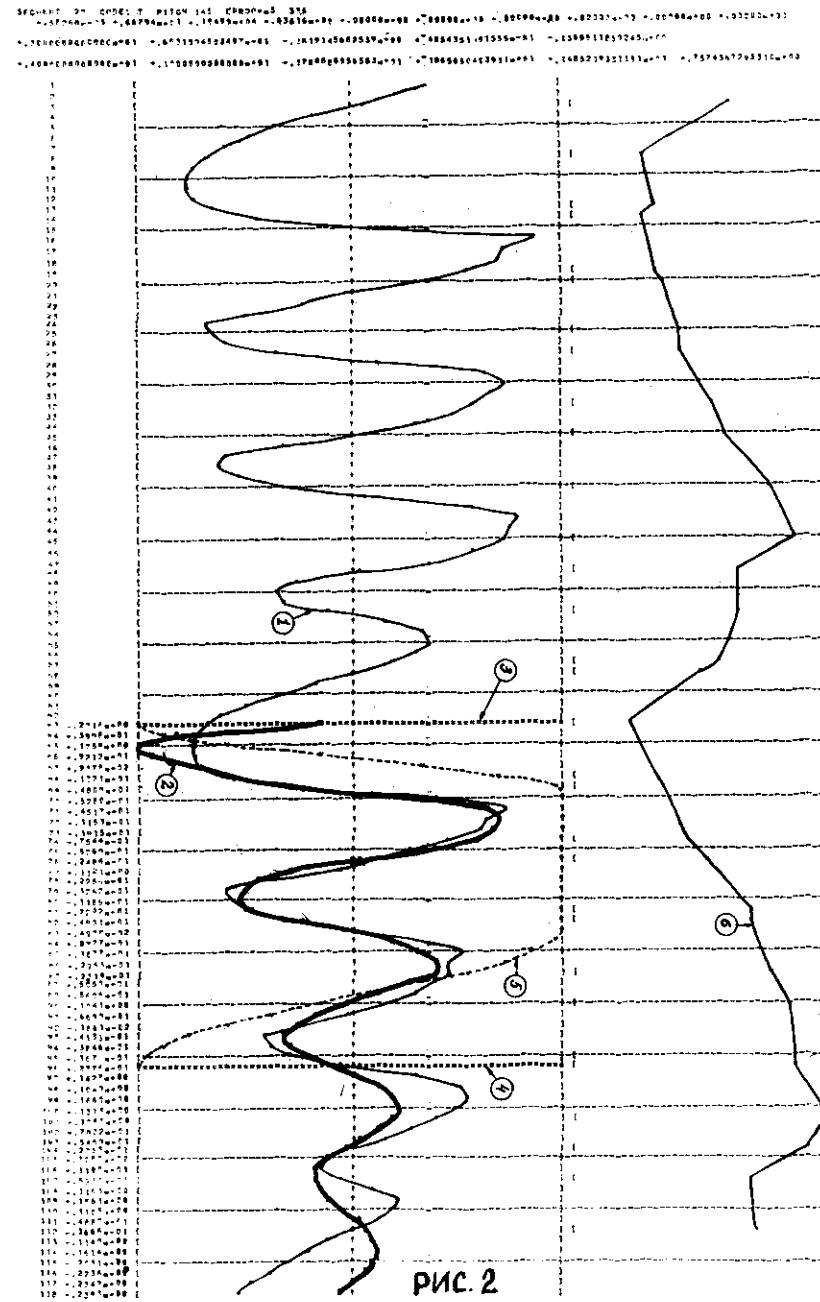


РИС. 2

Таблица I

Центральные частоты каналов (гц) на картинках "видимой речи"

номера каналов	0	1	2	3	4	5	6	7	8	9
0	-	75	98	121	145	170	195	220	246	273
10	300	328	356	385	415	445	476	507	540	572
20	606	640	675	711	748	785	823	862	902	942
30	984	1026	1069	1114	1159	1205	1252	1300	1349	1399
40	1451	1503	1556	1611	1667	1724	1782	1841	1902	1964
50	2027	2092	2158	2225	2294	2364	2436	2509	2584	2661
60	2739	2819	2900	2984	3069	3156	3244	3335	3428	3522
70	3619	3717	3818	3921	4026	4133	4243	4355	4469	4586
80	4705	4827	4952	5079	5208	5341	5476	5615	5756	5900
90	6048	6198	6352	6509	6669	6833	7000	-	-	-

Пример: центральная частота канала номер 63 : 2984 гц.

- индекс аperiodичности, умноженный для удобства на сто,
- относительный уровень сигнала на сегментах,
- взвешенная ошибка аппроксимации,
- частота ОТ в герцах.

На рисунках хорошо заметно (сегменты II2-III2), что когда первые и вторые полюса близки, то последние располагаются выше спектральных максимумов, которые человек принял бы за проявление второй форманты. Это соответствует физике явления, поскольку для близко расположенных резонансных кривых положение максимумов смещается.

При синтезе речевого сигнала были выбраны следующие параметры : длительность импульса возбуждения $\tau_n = 0,2 T_{от}$, форма импульса - треугольная, передний фронт $\tau_{нф} = 0,9 \tau_n$.

Характеристики формирующего фильтра подбирались экспериментально - это нуль на частоте 0 гц шириной 10 гц, нуль на частоте 8 гц шириной 1000 гц и полюс на частоте 3 кгц шириной 1000 гц.

Форма сигнала возбуждения после формирующего фильтра приведена на графике 5 - 1. 5-2 - отрезок синтезированного речевого сигнала. Это "а" в слове "убывала".

На качество синтезированной речи в большой степени влияет длительность импульса возбуждения и характеристики формирующего фильтра. Соотношение между фронтами возбуждающего импульса отражается на тембре речи, на осциллограмме, но незначительно влияет на качество.

Качество синтеза оценивалось путем неформального прослушивания и признано удовлетворительным. Случайное изменение в пределах нескольких процентов соседних периодов ОТ или длительностей импульса возбуждения не оказывало заметного влияния на качество. Особый интерес представляло сравнение качества синтеза по данным анализа с регуляризацией и без неё. Ощущается несколько большая плавность регуляризованной речи, но при этом ухудшается разборчивость, звук делается более глухим, машиноподобным.

Для корректной оценки объема информации, необходимого для передачи \mathcal{Z} -описания по линиям связи, необходимо было бы подвергнуть коэффициенты $D(\alpha)$ одному из видов оптимального кодирования. На данном этапе проверено лишь усечение длины мантиссы при передаче коэффициентов δ_j . Процедура эта выполня-

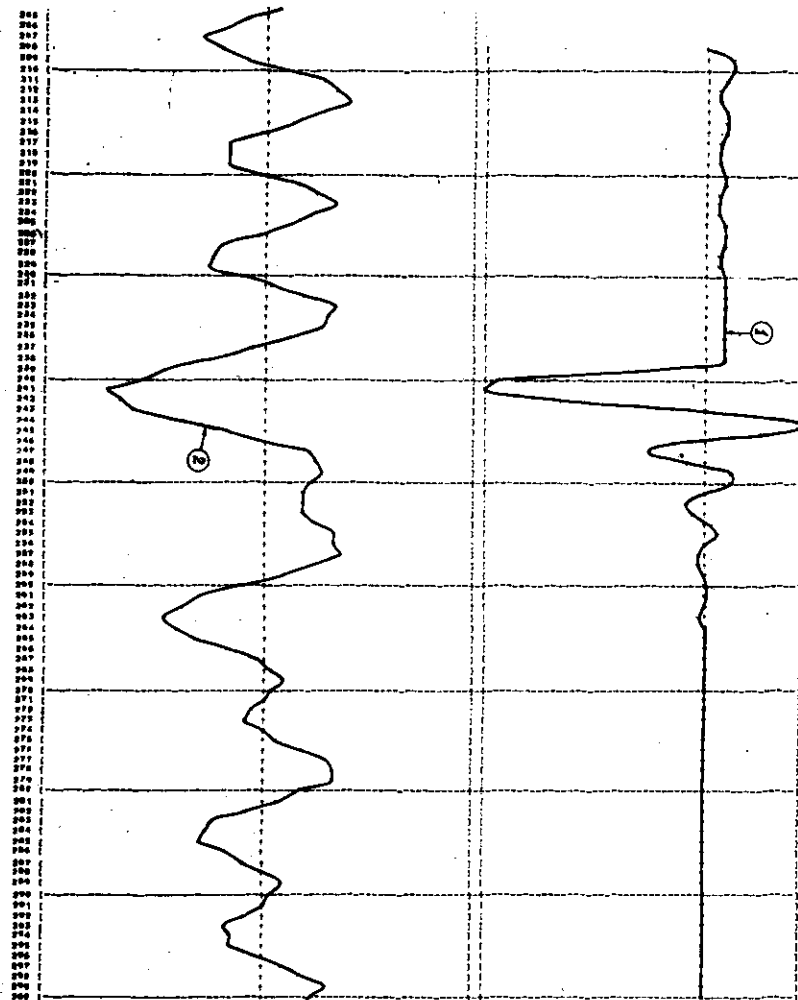


рис. 5

лась следующим образом. Находился динамический диапазон изменения коэффициента на длительности фразы; значение каждого коэффициента для каждого сегмента с учетом найденных диапазонов приводилось к интервалу 0-.9999, переводилось в представление с фиксированной запятой, после чего сохранялось заданное в программе число старших разрядов мантиссы. Эксперимент показал, что при уменьшении числа бит на каждый коэффициент до 15 ухудшения качества синтезированной речи на слух не заметно. Небольшие искажения появляются при числе бит 10. При числе бит 8 и менее система не работает. Допустимое количество бит (10) не зависит от степени используемой аппроксимации. Воспроизводился синтез при аппроксимации 8,6 и 4 порядка (соответственно 4,3 или 2 форманты). При $m = 4$ речь разборчива, но делается глуше. Возможно, в какой-то степени это поддастся корректировке при более тщательном подборе параметров источника возбуждения и формирующих фильтров.

Для более качественного звучания синтезированной речи следует ввести ещё одну кодировку типа сегмента: "K" - комбинированное (голос + шум) возбуждение. Это сделать, по-видимому, несложно, оценивая на стадии аппроксимации голосовых сегментов точность: если она ниже обычной, значит, в сигнале возбуждения присутствует шум.

6. Выводы

1. При работе на ЦВМ или специальной цифровой аппаратуре есть смысл пользоваться \mathcal{Z} - описанием дискретных систем, вместо того, чтобы допускать всевозможные неточности из-за попыток втиснуть в рамки дискретной техники непрерывные методы и описания.

Предлагается отказаться от попыток получить физиологически корректную формантную картину, учитывая трудности формализации процесса индексации спектральных максимумов. Предлагается синтезировать некую модельную (линейную динамическую) систему, добиться достаточного сходства сигналов во временной области, и \mathcal{Z} - параметры полученной передаточной функции модели использовать в качестве сокращенного описания речи.

2. Выяснено, что задача анализа, базирующаяся на средне-квадратической аппроксимации сигнала во временной области, является некорректной по Адамару с точки зрения вычисления поло-

жения полюсов (формант), т.е. близким сигналам во временной области могут соответствовать заметно различающиеся и в общем случае нестационарные формантные картиннки. Отмечено, что для корректного в смысле реального физического речевого тракта формантного анализа недостаточно осциллограммы речи; кроме более подробных знаний об источнике, необходимо быть уверенным в наблюдаемости всех состояний системы. Этот путь для практических применений, например, в связи, неконструктивен.

3. Рассматриваются алгоритмы анализа речевого сигнала на основе \mathcal{Z} -описания [1]. Описан алгоритм поиска оптимального интервала аппроксимации, дающий хорошие результаты. Указан способ получения \mathcal{Z} -описания речи для сегментов с шумовым возбуждением.

4. Для контроля качества анализа была смоделирована на ЦВМ работа синтезатора. К достоинствам предлагаемого \mathcal{Z} -описания следует отнести очевидную простоту синтезирующей части. Получены первые экспериментальные результаты по синтезу. Как и ожидалось, регуляризация анализа, приводящая к сглаживанию формантных картинок, качества звучания не улучшила. Без дополнительных перекодировок при передаче \mathcal{Z} -параметров достаточно сохранять по 10 бит на каждый коэффициент. Рекомендуется использовать в описании передаточные функции шестой степени, однако, допустимо понижение степени до четвертой.

Л и т е р а т у р а

1. Лозовский В.С. Аппроксимация отклика системы в \mathcal{Z} -плоскости и формантный анализ речи. Сб. "Вычислительные системы", Новосибирск, № 37, 1969.
2. Роберт Ли. Оптимальные оценки, определение характеристик и управление. Москва, Физматгиз "Наука", 1966.
3. G.Golub, Numerical Methods for Solving Linear Least Squares Problems, Numerische Mathematik, В.7, Н.3, 1965.
4. Тихонов А.Н. О некорректных задачах линейной алгебры и устойчивом методе их решения. ДАН СССР, т.163, № 3, 1968.
5. Лозовский В.С. Модифицированный разностный метод определения основного тона речи. "Труды акустического института", Москва, вып. XII, 1970.
6. Лозовский В.С. Программа синтеза рекурсивных цифровых фильтров.

тров. Сб. "Вычислительные системы", Новосибирск, № 35, 1969.

7. Лозовский В.С. Метод получения полюсно-нулевого описания речевого сигнала (Аннотация доклада на Всесоюзной школе-семинаре АРСО-5, Сухуми, 1969 г.), Труды Акустического института, М., вып. XII, 1970.
8. Лозовский В.С. Процедура решения полиномиальных уравнений методом Миллера. (наст. сборник).
9. Фант Г. Анализ и синтез речи, пер. с англ. под редакцией Н.Г. Загоруйко, Новосибирск, "Наука", 1970.
10. Лозовский В.С. Программа спектрального анализа в частотной и временной областях. Сб. "Вычислительные системы", Новосибирск, 1971, (в печати).

Поступила в редакцию
13.1.1971 г.