

УДК 519.226:534.784

ОЦЕНИВАНИЕ ПАРАМЕТРОВ РЕЧЕВОГО ТРАКТА В КЛАССЕ МОДЕЛЕЙ
АВТОРЕГРЕССИИ СО СТАЦИОНАРНОЙ СЕЗОННОЙ РАЗНОСТЬЮ
ПЕРВОГО ПОРЯДКА И СТАЦИОНАРНОЙ РАЗНОСТЬЮ ПЕРВОГО ПОРЯДКА

А.В. Кальманов

Цель данной работы является исследование возможности оценивания параметров речевого тракта по акустическому сигналу для вокализованных звуков речи в классе моделей авторегрессии со стационарной сезонной разностью первого порядка и стационарной разностью первого порядка.

Как известно, при анализе речевых сигналов возникает задача нахождения параметров речевого тракта. При ее решении предполагают, что процесс речеобразования описывается некоторой физической моделью, известной приближенно, а исследователь наблюдает сигнал только на выходе этой модели. Естественный способ решения такой задачи - построение математической модели по известному сигналу с последующим установлением соответствия между построенной и физической моделями. После установления такого соответствия выводы о физической модели можно строить в терминах математической.

В [1] показано, что модель авторегрессии неадекватно описывает голосовые звуки речи при асинхронном анализе, что говорит о невозможности получения точного описания речеобразования в рамках такой модели и необходимости построения другой модели. Данную работу следует рассматривать как попытку построения новой модели и установления соответствия между ее параметрами и параметрами физической модели.

§ I. Постановка задачи

Основные предположения (являющиеся физической моделью) о процессе речеобразования голосовых звуков речи не выходят за рамки общепринятой теории [2] и состоят в следующем.

1. Речевой сигнал является случайным процессом со спектром, ограниченным диапазоном частот $F_n/2$ так, что его можно квантовать с частотой F_n и задавать в виде временного ряда $\{a_n\}$.

2. Речь образуется в результате свертки функции возбуждения u_n с импульсной реакцией h_n речевого тракта. Амплитудно-частотная характеристика излучателя ($20 \log |H_{\text{изл.}}|$) имеет наклон +6 дБ/октаву.

3. В первом приближении речевой тракт можно рассматривать как линейную динамическую многорезонансную систему конечного порядка P_T с сосредоточенными параметрами. Передаточная функция речевого тракта содержит только полюсы

$$H_{PT}(z^{-1}) = \frac{1}{B(z^{-1})} = \frac{1}{1 - \sum_{i=1}^{P_T} e_i z^{-i}} \quad (I)$$

а каждый отсчет h_n может быть предсказан по P_T предшествующим отсчетам с ошибкой $\epsilon \sim N[0, \sigma_\epsilon^2]$ так, что

$$h_n = \sum_{i=1}^{P_T} e_i h_{n-i} + \epsilon_n \quad (2)$$

Каждой паре комплексно-сопряженных полюсов $H_{PT}(z^{-1})$ $z_j = \alpha_j \pm i\beta_j$, $j=1, 2, \dots, P_T/2$, соответствует формант F_j , частота которой F_j и полоса ΔF_j задаются в виде:

$$F_j = \frac{1}{2\pi T} \text{Im}[\ln z_j], \quad \Delta F_j = \frac{1}{\pi T} \text{Re}[\ln z_j], \quad (3)$$

где $T = 1/F_n$ - интервал квантования сигнала.

4. Функция возбуждения u_n для вокализованных участков речи представляет собой импульсный случайный процесс с детерминированным тактовым интервалом T_0 (периодом основного тона) так, что момент возникновения n -го импульса $t_n = nT_0 + v$, где v случайная величина с матожиданием $M[v] = 0$, $|v| \leq T_0/c$, $c \geq 2$, а $T_{\min} \leq t_n - t_{n-1} \leq T_{\max}$, где $T_{\min} = \frac{c-1}{c} T_0$, $T_{\max} = \frac{c+1}{c} T_0$. Форма импульсов

возбуждения известна приближенно и изменяется достаточно мало на коротких участках речи (интервалах анализа) длительностью $T_A = \lambda T_0$ при $2 \leq \lambda \leq 4$, а наклон спектра источника ($20 \log |N_{ГМ}|$) может изменяться от одного интервала анализа к другому в пределах от -5 дБ/октаву до -15 дБ/октаву.

Учитывая предположения, задачу нахождения параметров речевого тракта можно сформулировать следующим образом: по выборке $\{a_n\}$ объема $N = T_A F_s + 1$ определить параметры передаточной функции речевого тракта e_i , $i = 1, 2, \dots, P_T$, или формантные параметры $\langle F_j, \Delta F_j \rangle$, $j = 1, 2, \dots, P_T/2$.

§ 2. Решение задачи

Нетрудно убедиться в том, что задача относится к классу некорректно поставленных задач [3] и не может быть строго решена без учета информации об источнике возбуждения.

Задача определения параметров речевого тракта в классе моделей авторегрессии рассматривалась в работах [3-7]. Модель авторегрессии имеет вид [8]:

$$A(z^{-1}) w_n = \alpha_n, \quad (4)$$

$$w_n = V^d a_n, \quad (5)$$

$$A(z^{-1}) = 1 - \sum_{i=1}^p a_i z^{-i} \quad (6)$$

- стационарный оператор авторегрессии p -го порядка ($AR(p)$); $V = 1 - z^{-1}$ - разностный оператор со сдвигом назад; z^{-1} - оператор сдвига назад такой, что $z^{-1} w_n = w_{n-1}$, $i = 1, 2, \dots, p$; $\{a_i\}$ - временной ряд; $\{\alpha_n\}$ - белый гауссовский шум с дисперсией σ_α^2 ; n - номер отсчета. При аппроксимации, как правило, используют $d = 0$ или $d = 1$, предполагая, что с учетом (1) и (2) уравнение речеобразования имеет вид

$$a_n = H(z^{-1}) e_n = H_{ГМ}(z^{-1}) H_{РГ}(z^{-1}) H_{МЗД}(z^{-1}) e_n, \quad (7)$$

а оценку $H(z^{-1})$, учитывая (1) и (4)-(7), ищут в виде:

$$H_{ГМ}(z^{-1}) H_{МЗД}(z^{-1}) V^{-d} A^{-1}(z^{-1}) = V^{-d} \Lambda^{-1}(z^{-1}). \quad (8)$$

Отсюда следует, что если $H_{ГМ}(z^{-1}) H_{МЗД}(z^{-1}) \neq V^{-d}$, то $p > P_T$. В самом деле, $H_{ГМ}(z^{-1})$ имеет нули, каждый из которых можно вы-

проксимировать бесконечным числом полюсов. Это следует из того, что процесс скользящего среднего первого порядка можно представить в виде процесса авторегрессии бесконечного порядка [8]. Поэтому, сравнивая долю и прямую части (8), получаем $p > p_T$. Следовательно, $A(z^{-1}) = E(z^{-1})E_0(z^{-1})$, где $E_0(z^{-1})$ - стационарный оператор $\Delta P(p - p_T)$. На практике полагают, что $p - p_T \leq 5$. При этом по $A(z^{-1})$ нельзя найти $E(z^{-1})$, поскольку оценить можно только $A(z^{-1})$. Тем не менее, вычислив корни $z_{a1} = \alpha_{a1} \pm i\beta_{a1}$ ($1 = 1, 2, \dots, p/2$) полинома $A(z^{-1})$, их частоты и полюсы $\langle F_{a1}, \Delta F_{a1} \rangle$ в соответствии с (3), используя априорную информацию о частотах и полосах $\langle F_j, \Delta F_j \rangle$ ($j = 1, 2, \dots, p_T/2$) формант для конкретных звуков, можно отождествить некоторые из вычисленных частот F_{a1} с формантными частотами F_j . Но строгий алгоритм отождествления построить достаточно трудно. К тому же при совпадении нулей $H_{ГЛ}(z^{-1})$ с полюсами $H_{PT}(z^{-1})$ некоторые из формант могут быть утрачены.

Применение оператора γ к речевому сигналу приводит к подъему спектра его верхних частот и удалению из сигнала постоянной составляющей, что соответствует совместному учету наклона спектра $H_{ГЛ}(z^{-1})$ и $H_{НАК}(z^{-1})$, который предполагает равным -6 дБ/октаву. В действительности этот наклон изменяется во времени, поэтому используют адаптивный оператор $V_a = 1 - az^{-1}$, где $a = V_1/V_0$, а V_0, V_1 - автоковариации ряда $\{x_n\}$ для задержек 0 и 1. Однако, как показывают проведенные эксперименты, этого недостаточно для исключения влияния источника возбуждения и излучения.

Для устранения этого недостатка оценивание параметров проводят на интервалах полного смыкания голосовых связок [3,4]. Поиск таких интервалов является совсем непростой задачей. Если к тому же учесть, что у ряда дикторов эти интервалы отсутствуют, то желаемый результат не будет достигнут.

В данной работе предлагается иной подход для устранения влияния источника возбуждения и излучения без поиска интервалов полного смыкания. А именно параметры речевого тракта предлагается оценивать в классе моделей авторегрессии со стационарной сезонной разностью первого порядка и стационарной разностью первого порядка

$$A(z^{-1}) \nabla \nabla_{\tau} x_n = \alpha_n, \quad (9)$$

где ∇_k - оператор сезонной разности первого порядка такой, что $\nabla_k a_n = a_n - a_{n-k} = (1-z^{-k}) a_n$, т.е. оценку $H(z^{-1})$ предлагается искать в виде:

$$H_{ГМ}(z^{-1})H_{МЗД}(z^{-1})E^{-1}(z^{-1}) = \nabla^{-1}\nabla_k^{-1} \Lambda^{-1}(z^{-1}), \quad (10)$$

где вместо ∇ возможно использование ∇_n . Из (10) следует, что при $\Lambda(z^{-1}) = E(z^{-1})$

$$H_{ГМ}(z^{-1})H_{МЗД}(z^{-1}) = \nabla^{-1}\nabla_k^{-1}. \quad (11)$$

Наоборот, если (11) верно, то $\Lambda(z^{-1}) = E(z^{-1})$, и отождествлять корни $\Lambda(z^{-1})$ и $E(z^{-1})$ не нужно. Проведенные эксперименты доказали, что (10) и (11) достаточно хорошо выполняются.

Нетрудно показать, что максимальное устранение влияния источника возбуждения достигается при $k = N_{от} = T_0 F_n + 1$ ($N_{от}$ - число отсчетов в тактовом интервале). Поэтому оценивание параметра k эквивалентно определению интервала T_0 , который можно найти, например, методом обратной фильтрации [6].

Обозначим $w_n = \nabla \nabla_k a_n$, тогда (9) можно переписать в виде: $\Lambda(z^{-1})w_n = a_n$. Оценивание параметров оператора $\Lambda(z^{-1})$ сводится к решению системы нормальных уравнений:

$$\sum_{i=1}^p a_i R_{j-i} = R_j, \quad 1 \leq j \leq p,$$

где

$$R_j = \sum_{n=1}^{N_{1-j}} w'_n w'_{n-j}, \quad N_1 = N-k-1, \quad 0 \leq j \leq p,$$

$$w'_n = w_n \cdot v_n, \quad v_n = 0,54 + 0,46 \cos \left[\pi \frac{2(n-1)-N_1}{N_1} \right].$$

Для нахождения корней уравнения $\Lambda(z^{-1}) = 0$ первоначально использовались: процедура Шиллера, предложенная в [9] и переведенная на ФОРТРАН для ЭВМ "Минск-32", а также стандартная процедура Берстоу. Однако при степени полинома выше семи эти программы давали существенную ошибку. Поэтому была написана программа, реализующая метод Берстоу [10] с некоторыми модификациями, позволяющая значительно уменьшить ошибку.

§ 3. Описание эксперимента и результаты

Оценивание формантных параметров проводилось при помощи шести моделей, которые можно условно обозначить: $\Lambda(z^{-1})a_n$, $\Lambda(z^{-1})\nabla a_n$, $\Lambda(z^{-1})\nabla_a a_n$, $\Lambda(z^{-1})\nabla_k a_n$, $\Lambda(z^{-1})\nabla\nabla_k a_n$, $\Lambda(z^{-1})\nabla_a \nabla_k a_n$ (название моделей ясно из обозначения). Первые три, как отмечалось, уже применялись для нахождения параметров речевого тракта, остальные — предлагаются в данной работе. Эксперименты проводились на синтезированном и реальном сигналах.

Исследовались два вида возбуждения: 1) синусоидальное, 2) возбуждение δ -импульсами. В первом случае импульсы голосовых связок моделировались в виде:

$$y_i = \begin{cases} c \sin \frac{\pi(1-i)}{2\alpha(N_{OT}-1)}, & i=1, 2, \dots,]\alpha N_{OT}[, \\ \frac{c}{2} \left\{ 1 + \sin \frac{\pi}{2\beta} \left[\beta - \alpha + \frac{i-1}{N_{OT}-1} \right] \right\}, & i=] \alpha N_{OT}[+ 1, \dots,](\alpha+2\beta)N_{OT}[, \\ 0, & i=](\alpha+2\beta)N_{OT}[+ 1, \dots, N_{OT} . \end{cases}$$

Параметры α и β позволяют варьировать длительность интервалов открытия и закрытия голосовой щели и наклон спектра источника, который, например, при $\alpha = 0,23$ и $\beta = 0,25$ равен примерно -12 дБ/октаву. При δ -возбуждении $y_i = 1$ для $i=1$ и $y_i = 0$ для $i = 2, \dots, N_{OT}$. Длина синтезируемого ряда $N =]\lambda N_{OT}[$ (λ — обязательно целое). К импульсам возбуждения добавлялся гауссовский белый шум: $\tilde{y}_n = y_n + \eta$, $n = 1, 2, \dots, N$, где $\eta \sim N(0, \sigma_\eta)$. В экспериментах $\sigma_\eta = 0,01$ при $c = 1$. Наконец, сигнал центрировался:

$$y_n^* = \tilde{y}_n - \frac{1}{N} \sum_{n=1}^N \tilde{y}_n .$$

Для моделирования речевого тракта использовалось пять формант. Каждое из пяти формантных колебаний задавалось моделью АР (2), которой соответствует фильтр с передаточной функцией

$$H_i(z^{-1}) = 1 / (1 - b_{1i} z^{-1} - b_{2i} z^{-2}) .$$

Параметры b_{1i} и b_{2i} находятся по частоте F_i и полосе ΔF_i i -й форманты из соотношений: $b_{1i} = 2\rho_i \cos \omega_i T$, $b_{2i} = -2\rho_i^2$, $\omega_i = 2\pi F_i$, $\rho_i = e^{-\Delta F_i \pi T}$. Поэтому

$$H_{PT}(z^{-1}) = \prod_{i=1}^5 H_{i1}(z^{-1}) = 1 / (1 - \sum_{i=1}^{10} e_i z^{-i}),$$

а параметры e_i нетрудно найти из произведения пяти полиномов этого порядка. Таким образом, модель речевого тракта является сложной десятого порядка.

Значения частот и полос формант, используемых при моделировании шести русских гласных, были взяты из [2]. Синтез сигнала осуществлялся по уравнению

$$a_n = \sum_{i=1}^{10} e_i a_{n-i} + y_n^*, \quad 1 \leq n \leq N. \quad (12)$$

Для исследования надежности метода в условиях аналого-цифрового преобразования предусматривалась возможность добавления к a_n равномерного на интервале $(-\xi, \xi)$ шума u_n , т.е. $a_n^* = a_n + u_n$. Семизрядный преобразователь перекрывает динамический диапазон примерно 42 дБ, а его ошибка квантования $\xi = 2^{-8}$. Поэтому $\sigma_u = \xi / \sqrt{2} = 0,002255$. При моделировании использовалось значение $\sigma_u = 0,003$.

Пусть F_{ij} - заданное при синтезе значение частоты i -й форманты для j -й гласной, \hat{F}_{ij} - ее оценка, тогда $\delta_{ij} = |F_{ij} - \hat{F}_{ij}| / F_{ij}$ - относительная ошибка оценивания i -й форманты для j -й гласной,

$\delta_{.j} = \sum_{i=1}^5 \delta_{ij} / 5$ - средняя ошибка оценивания частоты форманты для каждой из шести гласных, $\delta_{i.} = \sum_{j=1}^6 \delta_{ij} / 6$ - средняя ошибка оценивания каждой из пяти формантных частот, $\delta = \sum_{i=1}^5 \sum_{j=1}^6 \delta_{ij} / 30$ - ошибка оценивания частоты форманты, усредненная по шести гласным.

Если $\{y_n\}$ - последовательность δ -импульсов, то $\{a_n\}$ - последовательность импульсных откликов модели речевого тракта, полученная из (12). Поэтому можно определить ошибку оценивания формантных частот по импульсному отклику в классе моделей $AR(10)$. Очевидно, эта ошибка будет совпадать с ошибкой оценивания максимумов спектра сигнала $\{a_n\}$, содержащего только полюса. В табл. I и 2 приведены значения $\delta_{i.}$, $\delta_{.j}$ и δ для различных σ_u и σ_η при δ -возбуждении и порядке аппроксимирующей модели $p = 10$. Из таблиц видно, что при одном и том же σ_η и различных σ_u ошибки оценивания совпадают. Следовательно, ошибки квантования практически

Т а б л и ц а 1

Средняя ошибка (в %) оценивания частоты форманты
для шести гласных в классе моделей $\Lambda(z^{-1})a_n$
при δ -возбуждении; $p = 10$

| № п/п | σ_u | σ_η | $ a _{\delta_1}$ | $ o _{\delta_2}$ | $ y _{\delta_3}$ | $ э _{\delta_4}$ | $ ы _{\delta_5}$ | $ и _{\delta_6}$ | Средняя ошибка |
|-------|------------|---------------|------------------|------------------|------------------|------------------|------------------|------------------|-------------------|
| 1 | 0 | 0 | 0,84 | 1,62 | 2,74 | 0,71 | 2,49 | 1,33 | 1,62 |
| 2 | 0 | 0,01 | 0,9 | 1,66 | 3,01 | 0,49 | 2,47 | 1,43 | 1,66 |
| 3 | 0,003 | 0 | 0,84 | 1,62 | 2,74 | 0,71 | 2,49 | 1,33 | 1,62 |
| 4 | 0,003 | 0,01 | 0,9 | 1,66 | 3,01 | 0,49 | 2,47 | 1,43 | 1,66 |

Т а б л и ц а 2

Средняя ошибка (в %) оценивания каждой из пяти
формантных частот в классе моделей $\Lambda(z^{-1})a_n$
при δ -возбуждении; $p = 10$

| № п/п | σ_u | σ_η | F_1/δ_1 | F_2/δ_2 | F_3/δ_3 | F_4/δ_4 | F_5/δ_5 | Средняя ошибка |
|-------|------------|---------------|----------------|----------------|----------------|----------------|----------------|-------------------|
| 1 | 0 | 0 | 6,5 | 0,69 | 0,56 | 0,29 | 0,08 | 1,62 |
| 2 | 0 | 0,01 | 6,63 | 0,95 | 0,46 | 0,18 | 0,09 | 1,66 |
| 3 | 0,003 | 0 | 6,5 | 0,69 | 0,56 | 0,29 | 0,08 | 1,62 |
| 4 | 0,003 | 0,01 | 6,63 | 0,95 | 0,46 | 0,18 | 0,09 | 1,66 |

не влияют на ошибки оценивания. Средняя ошибка δ при $\delta_\eta = 0,01$ несколько больше, чем при $\delta_\eta = 0$. Это объясняется тем, что при $\delta_\eta = 0,01$ последовательность $\{a_n\}$ уже не является последовательностью импульсных откликов, поскольку $\{y_n\}$ - "замуленные" δ -импульсы. В результате вместо формантных частот оцениваются максимумы спектра сигнала, что и приводит к увеличению ошибки, которая будет тем больше, чем больше сигнал возбуждения отличается от δ -импульса.

В табл. 3 и 4 приведены ошибки оценивания для шести моделей ($1 - \Lambda(z^{-1})a_n$; $2 - \Lambda(z^{-1})\nabla a_n$; $3 - \Lambda(z^{-1})\nabla_a a_n$; $4 - \Lambda(z^{-1})\nabla_k a_n$; $5 - \Lambda(z^{-1})\nabla \nabla_k a_n$; $6 - \Lambda(z^{-1})\nabla_a \nabla_k a_n$) при синусоидальном возбуждении и различных σ_u и σ_η ($p = 10$). Из таблиц видно, что ошибки квантования практически не влияют на ошибки оценивания. Значения ошибок δ примерно одинаковы для одного и того же класса моделей при различных операторах ∇ и ∇_a . Это объясняется тем, что для гласных звуков $\nabla \approx \nabla_a$, так как $a \approx 1$.

Т а б л и ц а 3

Средняя ошибка (в %) оценивания частоты форманты для шести русских гласных в классе моделей: $\Lambda(z^{-1})_{a_n}(1)$;

$$\Lambda(z^{-1})_{\nabla a_n}(2); \quad \Lambda(z^{-1})_{\nabla a_n}(3); \quad \Lambda(z^{-1})_{\nabla k a_n}(4);$$

$\Lambda(z^{-1})_{\nabla \nabla k a_n}(5); \quad \Lambda(z^{-1})_{\nabla a \nabla k a_n}(6)$ при синусоидальном возбуждении; $p = 10$

| σ_{η} | σ_{κ} | Но дель | $ a $ | $ b $ | $ y $ | $ z $ | $ m $ | $ n $ | Сред- нее δ |
|-----------------|-------------------|------------|---------------|---------------|---------------|---------------|---------------|---------------|-----------------------|
| | | | $\delta_{.1}$ | $\delta_{.2}$ | $\delta_{.3}$ | $\delta_{.4}$ | $\delta_{.5}$ | $\delta_{.6}$ | |
| 0 | 0 | 1 | 26,59 | 26,46 | 33,02 | 29,30 | 33,46 | 26,12 | 29,16 |
| | | 2 | 24,48 | 26,43 | 8,11 | 27,30 | 9,50 | 26,70 | 20,42 |
| | | 3 | 24,33 | 27,96 | 8,08 | 27,3 | 9,51 | 26,58 | 20,63 |
| | | 4 | 25,66 | 25,49 | 4,91 | 24,26 | 25,17 | 23,08 | 21,43 |
| | | 5 | 0,46 | 1,31 | 2,70 | 0,61 | 1,49 | 1,71 | 1,37 |
| | | 6 | 0,97 | 1,23 | 2,57 | 0,70 | 1,48 | 1,82 | 1,46 |
| 0 | 0,003 | 1 | 26,59 | 26,48 | 33,15 | 29,30 | 33,30 | 26,12 | 29,16 |
| | | 2 | 24,52 | 26,43 | 8,09 | 27,28 | 9,91 | 26,53 | 20,46 |
| | | 3 | 24,38 | 26,80 | 8,08 | 27,31 | 9,93 | 26,53 | 20,51 |
| | | 4 | 25,67 | 25,51 | 4,24 | 24,26 | 25,23 | 23,07 | 21,33 |
| | | 5 | 0,45 | 1,31 | 2,68 | 0,63 | 1,47 | 1,70 | 1,37 |
| | | 6 | 0,97 | 1,23 | 2,57 | 0,70 | 1,47 | 1,83 | 1,46 |
| 0,01 | 0 | 1 | 26,35 | 27,20 | 33,65 | 28,15 | 30,53 | 25,34 | 28,54 |
| | | 2 | 5,88 | 6,87 | 8,91 | 27,46 | 7,88 | 26,73 | 13,95 |
| | | 3 | 6,72 | 7,27 | 8,88 | 27,58 | 7,89 | 26,74 | 14,18 |
| | | 4 | 8,46 | 5,88 | 5,82 | 23,90 | 11,35 | 22,95 | 13,06 |
| | | 5 | 1,43 | 1,31 | 3,17 | 1,13 | 0,85 | 1,82 | 1,62 |
| | | 6 | 1,07 | 1,23 | 3,05 | 1,19 | 0,85 | 2,09 | 1,58 |
| 0,01 | 0,003 | 1 | 26,35 | 27,13 | 33,72 | 28,15 | 30,58 | 25,34 | 28,55 |
| | | 2 | 5,89 | 6,91 | 8,55 | 27,46 | 7,87 | 26,33 | 13,95 |
| | | 3 | 6,74 | 7,25 | 8,84 | 27,59 | 7,88 | 26,73 | 14,17 |
| | | 4 | 8,45 | 5,90 | 5,82 | 23,90 | 11,37 | 22,94 | 13,06 |
| | | 5 | 1,44 | 1,31 | 3,20 | 1,14 | 0,79 | 1,82 | 1,62 |
| | | 6 | 1,07 | 1,23 | 3,08 | 1,20 | 0,79 | 2,09 | 1,58 |

Средняя ошибка оценивания частоты форманты δ при различных σ_{κ} и σ_{η} составляет величину: в классе моделей AP(10) примерно 30%, в классе моделей AP(10) со стационарной разностью перво-

Т а б л и ц а 4

Ошибка (в %) оценивания пяти формантных частот,
усредненная по шести русским гласным, в классе
моделей: $A(z^{-1})v_{a_n}(1)$; $A(z^{-1})v_{a_n}(2)$; $A(z^{-1})v_{a_n}(3)$;

$A(z^{-1})v_{a_n}(4)$; $A(z^{-1})v_{a_n}(5)$; $A(z^{-1})v_{a_n}(6)$

при синусоидальном возбуждении; $p = 10$

| σ_{η} | σ_{κ} | Мо- дель | F_1 $\delta_{1.}$ | F_2 $\delta_{2.}$ | F_3 $\delta_{3.}$ | F_4 $\delta_{4.}$ | F_5 $\delta_{5.}$ | Сред- нее δ |
|-----------------|-------------------|-------------|------------------------|------------------------|------------------------|------------------------|------------------------|-----------------------|
| 0 | 0 | 1 | 18,75 | 15,72 | 6,34 | 100,00 | 4,97 | 29,16 |
| | | 2 | 13,79 | 11,60 | 5,24 | 67,33 | 4,17 | 20,42 |
| | | 3 | 15,93 | 10,64 | 5,2 | 67,32 | 4,16 | 20,63 |
| | | 4 | 4,88 | 21,55 | 6,98 | 69,3 | 4,44 | 21,43 |
| | | 5 | 2,18 | 1,20 | 1,12 | 0,87 | 1,46 | 1,37 |
| | | 6 | 2,60 | 0,94 | 1,23 | 0,84 | 1,68 | 1,46 |
| 0 | 0,003 | 1 | 18,81 | 15,64 | 7,04 | 100,00 | 4,97 | 29,16 |
| | | 2 | 13,79 | 11,47 | 5,27 | 67,33 | 4,44 | 20,46 |
| | | 3 | 13,80 | 11,61 | 5,33 | 67,34 | 4,45 | 20,51 |
| | | 4 | 4,88 | 21,54 | 7,10 | 68,74 | 4,40 | 21,33 |
| | | 5 | 2,07 | 1,05 | 1,34 | 1,04 | 1,41 | 1,37 |
| | | 6 | 2,62 | 0,86 | 1,38 | 1,07 | 1,35 | 1,46 |
| 0,01 | 0 | 1 | 17,64 | 13,92 | 7,60 | 100,00 | 3,53 | 28,54 |
| | | 2 | 15,85 | 11,31 | 4,37 | 35,23 | 3,07 | 13,95 |
| | | 3 | 16,44 | 11,62 | 4,65 | 35,12 | 3,08 | 14,18 |
| | | 4 | 7,11 | 10,28 | 4,39 | 40,07 | 3,42 | 13,06 |
| | | 5 | 2,19 | 1,27 | 1,17 | 2,46 | 1,02 | 1,62 |
| | | 6 | 2,21 | 1,15 | 1,12 | 2,47 | 0,95 | 1,58 |
| 0,01 | 0,003 | 1 | 17,61 | 13,9 | 7,63 | 100,00 | 3,59 | 28,55 |
| | | 2 | 15,85 | 11,31 | 4,33 | 35,20 | 3,08 | 13,95 |
| | | 3 | 16,44 | 11,62 | 4,61 | 35,10 | 3,09 | 14,17 |
| | | 4 | 7,26 | 10,5 | 4,34 | 40,04 | 3,83 | 13,06 |
| | | 5 | 2,21 | 1,27 | 1,20 | 2,39 | 1,01 | 1,62 |
| | | 6 | 2,18 | 1,15 | 1,10 | 2,50 | 0,95 | 1,58 |

го порядка 14-21%, в классе моделей $AR(10)$ со стационарной сезон-
ной разностью первого порядка 13-21%, а в классе моделей $AR(10)$
со стационарной сезонной разностью первого порядка и стационар -

ной разностью первого порядка 1,4-1,7%. Взятие сезонной разности позволяет уменьшить ошибку в 1,4-2,3 раза. Примерно такой же результат дает взятие первой разности. Совместное же их использование понижает ошибку оценивания в 17,6-21,4 раза.

Из табл.4 видно, что ошибка оценивания частоты каждой форманты для всех гласных лежит в пределах 0,84-2,6% при использовании предлагаемых моделей. В то же время она может достигать 100% для других моделей, что говорит о потере форманты (в экспериментах при потере форманты ошибка полагалась равной 100%). Поэтому выполнение неравенства $\delta_{j,1} \geq 201$ для табл.3 означает утерю i форманты для j -й гласной. Аналогично выполнение неравенства $\delta_{j,1} \geq 50j/3$ для табл.4 означает утерю i -й форманты для j гласных.

Табл. 5 и рис. 1 и 2 дают наглядное представление о том, чего позволяет достигнуть все шесть моделей при оценивании параметров авторегрессии $\{a_1\}$ автокорреляций импульсного отклика ($r_1 = R_1/R_0$) и логарифмической амплитудно-частотной характеристики модели речевого тракта для гласной /а/. В первой строке табл.5 приведены заданные при синтезе параметры, во второй - их оценки, полученные по импульсному отклику (при δ -возбуждении). В строках 3-8 приведены оценки параметров для различных моделей ($p = 10$) при синусоидальном возбуждении. Наилучшие оценки в последних двух строках.

Рис. 2 иллюстрирует влияние источника на оценивание частотной характеристики. Хорошо видно, что F4 и F5 различимы только для предлагаемых моделей, а сходство оценки частотной характеристики для этих моделей с частотной характеристикой, полученной по импульсному отклику (рис.1), очевидно.

В табл.6-9 приведены ошибки оценивания для различных моделей при $p = 12, 14$. Во всех случаях ошибка оценивания в классе предлагаемых моделей меньше, чем в классе применявшихся ранее. Кроме того, ошибка оценивания в классе предлагаемых моделей при $p = 10$ меньше, чем ошибка оценивания в классе ранее использовавшихся при $p = 12, 14$. Как уже упоминалось, при $p > p_0 = 10$ возникает задача отождествления. Эта задача решалась вручну.

Описанный метод проверялся на реальном речевом сигнале. В ЭВМ "Минск-32" через семизрядный аналого-цифровой преобразователь вводилась фраза "Вода в луже медленно убывает" при $F_s = 10$ кГц. Длительность интервала анализа T_a составила 25,6 мсек, а сдвиг

от одного интервала к другому был равен 16 мсек. Для примера на рис.3 приведена картинка видимой речи, полученная по тестовой фразе.

Параметр k находился при помощи метода обратной фильтрации. Значение частот полюсов аппроксимирующих моделей при $p = 10, 12, 14$ выводилось в виде графика для получения траекторий формантных частот. На рис.4 приведены траектории пяти формантных частот, полученных по параметрам модели $A(z^{-1})\nabla\nabla_k a_n$ при $p = 10$. Их значения хорошо согласуются с данными, имеющимися в литературе не только для гласных, но и для звонких согласных. Для модели $A(z^{-1})\nabla_a \nabla_k$ результаты аналогичны.

Для $p > 10$ формантные частоты по полюсам моделей определяются неоднозначно (так как при $F_n = 10$ кГц имеем $T_T = 10$ [?]). Неоднозначность устраняется при $p = 10$, однако для моделей $A(z^{-1})a_n$, $A(z^{-1})\nabla a_n$, $A(z^{-1})\nabla_a a_n$, $A(z^{-1})\nabla_k a_n$, в отличие от предлагаемых, происходит утеря формант под влиянием источника возбуждения. Наибольшее число потерь формант происходит при использовании модели $A(z^{-1})a_n$. Модели $A(z^{-1})\nabla_a a_n$, $A(z^{-1})\nabla_a a_n$ и $A(z^{-1})\nabla_k a_n$ дают примерно одинаковые, но лучшие, чем модель $A(z^{-1})a_n$, результаты.

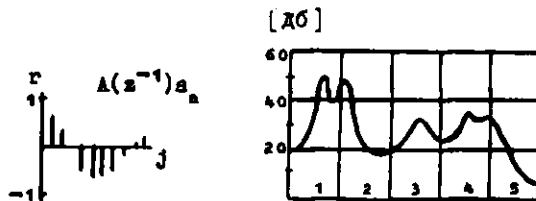


Рис.1. Автокорреляционная функция импульсного отклика и оценка логарифмической амплитудно-частотной характеристики модели речевого тракта, полученная при помощи модели $A(z^{-1})a_n$. Гласная /а/; $p = 10$; $\sigma_k = 0$; $\sigma_n = 0$; $\sigma_\eta = 0$; δ -возбуждение.

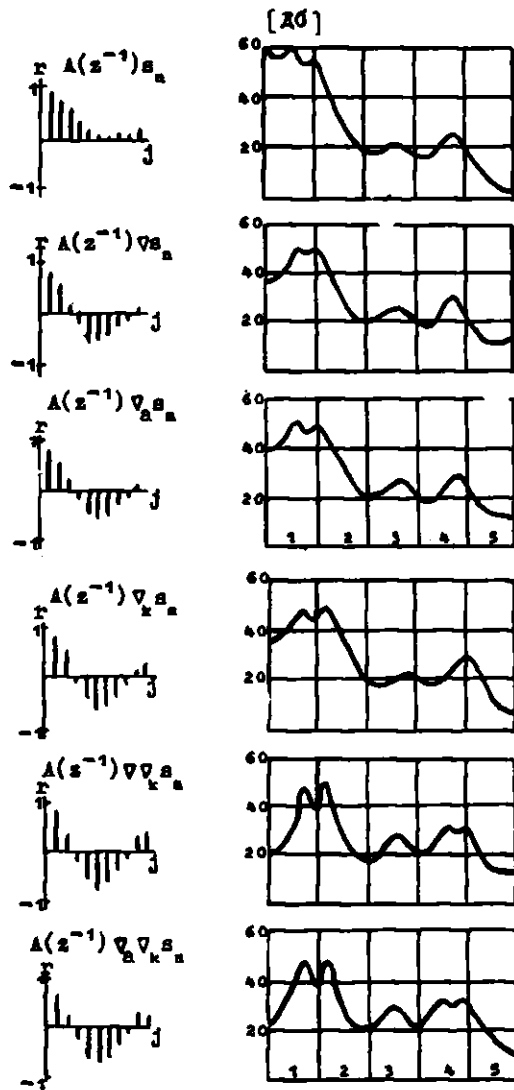


Рис.2. Автокорреляционные функции оценок импульсного отклика и оценки амплитудно-частотной характеристики речевого тракта, полученные при помощи шести моделей для гласной /а/; $p = 10$, $\sigma_{\eta} = 0$, $\sigma_{\eta} = 0$; синусоидальное возбуждение.

Т а б л и ц а 5

Параметры авторегрессии, заданные при синтезе, и их оценки, полученные при помощи различных моделей при опускательном и 6-возбужденных; $p = 10$, $\sigma_n = 0$, $\sigma_k = 0$

| Модель | a_1 | a_2 | a_3 | a_4 | a_5 | a_6 | a_7 | a_8 | a_9 | a_{10} | Вед. возбуждения |
|----------------------------------|--------|---------|--------|---------|--------|---------|--------|---------|--------|----------|------------------|
| Можная, $\mathcal{K}(z^{-1})$ | 0,7388 | -0,4892 | 0,9084 | -1,3633 | 0,6494 | -1,1881 | 0,8081 | -0,3273 | 0,3706 | -0,6714 | |
| $\Lambda(z^{-1})e_n$ | 0,7350 | -0,5062 | 0,9017 | -1,3876 | 0,6164 | -1,1791 | 0,8190 | -0,3434 | 0,3325 | -0,6830 | 6 |
| $\Lambda(z^{-1})e_n$ | 2,7276 | -3,0524 | 2,4581 | -2,8246 | 3,1146 | -2,5205 | 2,0774 | -1,5590 | 0,6114 | -0,0478 | |
| $\Lambda(z^{-1})\nabla e_n$ | 1,5897 | -0,8673 | 0,5382 | -1,1036 | 0,7763 | -0,4303 | 0,4777 | -0,0848 | 0,3763 | 0,1669 | |
| $\Lambda(z^{-1})\nabla e_n$ | 1,6233 | -0,8785 | 0,5186 | -1,0896 | 0,7791 | -0,3967 | 0,4624 | -0,0745 | 0,4070 | 0,2029 | вн |
| $\Lambda(z^{-1})\nabla e_n$ | 1,7648 | -1,2704 | 0,9956 | -1,5926 | 1,3451 | -0,8925 | 0,8545 | -0,4621 | -0,108 | 0,0907 | |
| $\Lambda(z^{-1})\nabla e_n$ | 0,7264 | -0,3745 | 0,7070 | -1,2140 | 0,4846 | -0,9642 | 0,5642 | -0,2229 | 0,2613 | -0,6168 | |
| $\Lambda(z^{-1})\nabla e_n$ | 0,7755 | -1,3761 | 0,7036 | -1,2028 | 0,5073 | -0,9172 | 0,6584 | -0,2056 | 0,2409 | -0,5690 | |

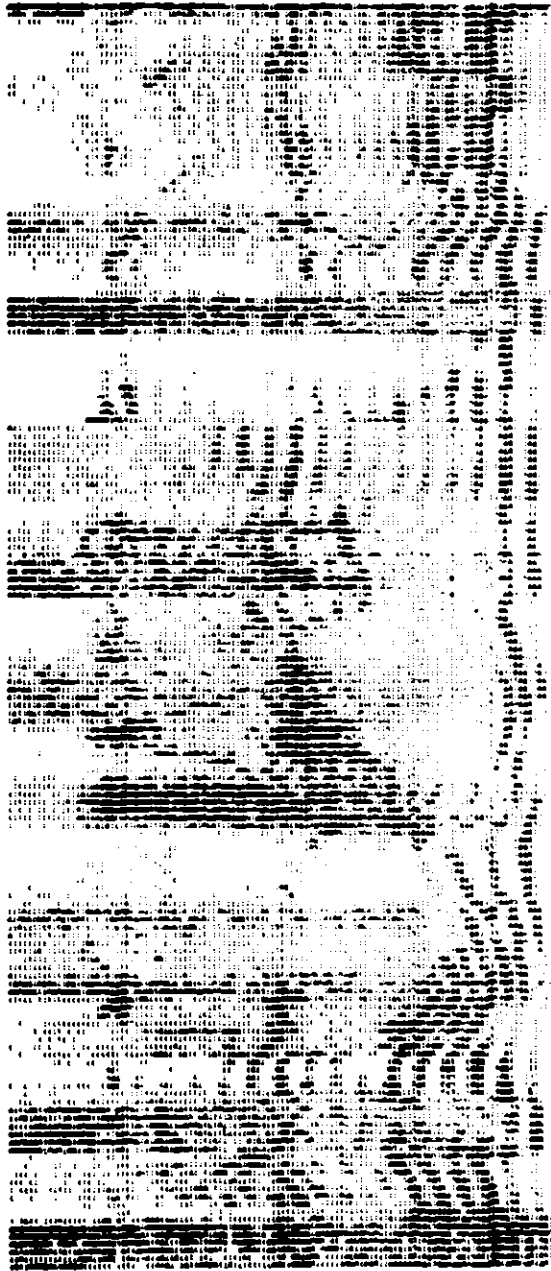
Т а б л и ц а 6

Средняя ошибка (в %) оценивания частоты форманты для шести русских гласных в классе моделей: $\Lambda(z^{-1})s_n(1)$, $\Lambda(z^{-1})\nabla s_n(2)$,

$\Lambda(z^{-1})\nabla_a s_n(3)$, $\Lambda(z^{-1})\nabla_{\kappa} s_n(4)$, $\Lambda(z^{-1})\nabla\nabla_{\kappa} s_n(5)$,

$\Lambda(z^{-1})\nabla_a\nabla_{\kappa} s_n(6)$ при синусоидальном возбуждении и различных значениях σ_{η} и σ_{κ} ; $p = 12$

| σ_{η} | σ_{κ} | Ио- дель | /а/ $\delta_{.1}$ | /о/ $\delta_{.2}$ | /у/ $\delta_{.3}$ | /э/ $\delta_{.4}$ | /ы/ $\delta_{.5}$ | /и/ $\delta_{.6}$ | Сред- нее δ |
|-----------------|-------------------|-------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|--------------------------|
| 0 | 0 | 1 | 0,81 | 3,61 | 12,26 | 1,27 | 6,04 | 10,62 | 5,77 |
| | | 2 | 1,05 | 1,64 | 4,66 | 1,66 | 2,09 | 3,88 | 2,50 |
| | | 3 | 1,18 | 1,78 | 4,56 | 1,71 | 2,09 | 3,88 | 2,53 |
| | | 4 | 0,41 | 1,90 | 4,80 | 0,67 | 1,28 | 2,49 | 1,91 |
| | | 5 | 0,92 | 1,41 | 2,94 | 1,12 | 0,8 | 1,44 | 1,44 |
| | | 6 | 0,97 | 1,28 | 2,90 | 1,10 | 0,8 | 1,41 | 1,41 |
| 0 | 0,003 | 1 | 0,81 | 3,36 | 12,42 | 1,27 | 6,06 | 10,63 | 5,76 |
| | | 2 | 1,05 | 1,63 | 4,51 | 1,66 | 2,08 | 3,88 | 2,49 |
| | | 3 | 1,18 | 1,77 | 4,45 | 1,70 | 2,08 | 3,88 | 2,52 |
| | | 4 | 0,41 | 1,87 | 4,74 | 0,66 | 1,24 | 2,49 | 1,91 |
| | | 5 | 0,92 | 1,40 | 2,95 | 1,12 | 0,79 | 1,44 | 1,44 |
| | | 6 | 0,96 | 1,29 | 2,93 | 1,10 | 0,79 | 1,41 | 1,41 |
| 0,01 | 0 | 1 | 1,54 | 4,23 | 13,45 | 1,69 | 5,48 | 10,20 | 6,10 |
| | | 2 | 1,00 | 1,43 | 4,9 | 1,87 | 1,67 | 4,17 | 2,51 |
| | | 3 | 1,03 | 1,45 | 4,78 | 1,88 | 1,66 | 4,17 | 2,50 |
| | | 4 | 1,28 | 2,21 | 3,04 | 2,02 | 0,80 | 2,77 | 2,02 |
| | | 5 | 1,74 | 2,11 | 2,43 | 0,99 | 1,11 | 1,66 | 1,67 |
| | | 6 | 1,48 | 2,04 | 2,38 | 1,06 | 1,12 | 1,88 | 1,66 |
| 0,01 | 0,003 | 1 | 1,54 | 4,52 | 13,46 | 1,69 | 5,48 | 10,20 | 6,15 |
| | | 2 | 1,00 | 1,43 | 4,90 | 1,87 | 1,66 | 4,17 | 2,51 |
| | | 3 | 1,03 | 1,44 | 4,76 | 1,87 | 1,66 | 4,17 | 2,49 |
| | | 4 | 1,28 | 2,22 | 3,10 | 2,03 | 0,81 | 2,76 | 2,03 |
| | | 5 | 1,75 | 2,10 | 2,48 | 0,98 | 1,11 | 1,66 | 1,68 |
| | | 6 | 1,50 | 2,03 | 2,45 | 1,06 | 1,11 | 1,88 | 1,67 |



В | О | Д | А | В | Л | У | Ж | Е | М | Е | Д | Л | Е | Н | Н | Ю | У | Б | Ы | В | А | Л | А

Рис. 3. Визуальная речь для фразы "Буде в дуже медженно убываля".

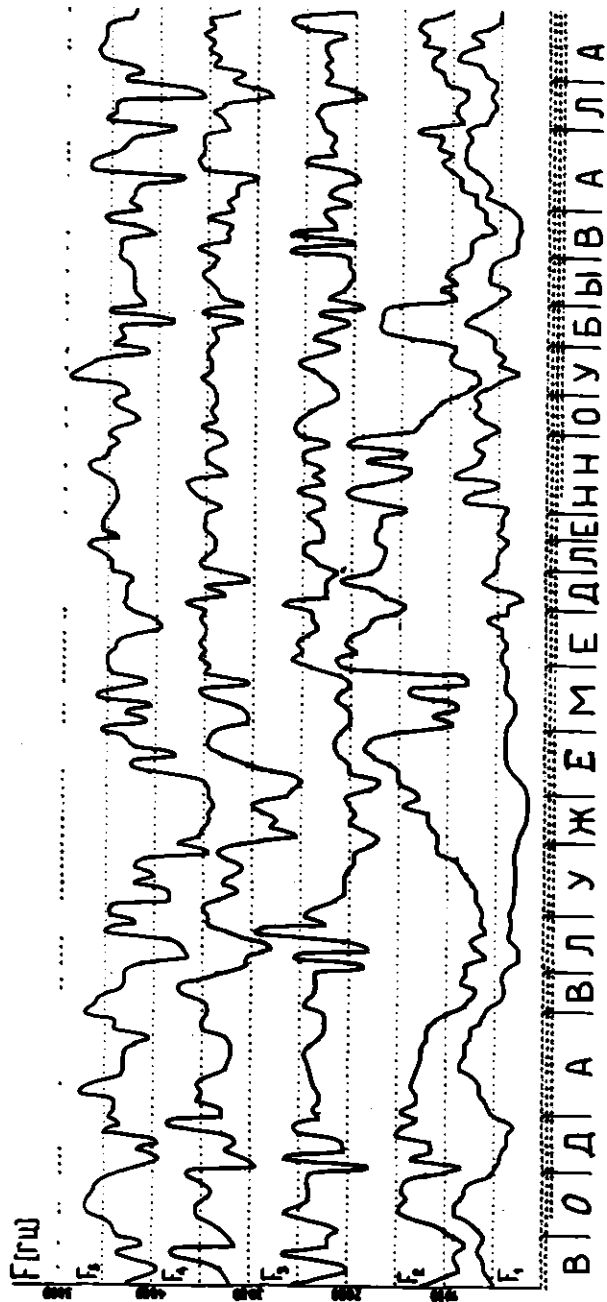


Рис. 4. Трассировки формантных частот, полученные по параметрам модели $\lambda(s^{-1}) \nu \nu_2, \nu$, при $p = 10$.

Т а б л и ц а 7

Ошибка (в %) оценивания формантных частот,
 усредненная по шести русским гласным, в классе
 моделей: $\Lambda(z^{-1})a_n(1)$; $\Lambda(z^{-1})\nabla a_n(2)$; $\Lambda(z^{-1})\nabla a_n(3)$;
 $\Lambda(z^{-1})\nabla a_n(4)$; $\Lambda(z^{-1})\nabla\nabla a_n(5)$; $\Lambda(z^{-1})\nabla a_n(6)$
 при синусоидальном возбуждении и различных значениях
 σ_η и σ_k ; при $p = 12$

| σ_η | σ_k | Мо- дель | F_1 | F_2 | F_3 | F_4 | F_5 | Сред- нее δ |
|---------------|------------|-------------|---------------|---------------|---------------|---------------|---------------|--------------------------|
| | | | $\delta_{1.}$ | $\delta_{2.}$ | $\delta_{3.}$ | $\delta_{4.}$ | $\delta_{5.}$ | |
| 0 | 0 | 1 | 18,76 | 5,19 | 1,40 | 1,27 | 2,27 | 5,77 |
| | | 2 | 5,94 | 3,50 | 1,20 | 0,99 | 0,86 | 2,50 |
| | | 3 | 6,02 | 3,48 | 1,18 | 1,05 | 0,93 | 2,53 |
| | | 4 | 1,85 | 2,60 | 2,67 | 1,46 | 1,06 | 1,91 |
| | | 5 | 2,71 | 0,96 | 1,68 | 1,38 | 0,44 | 1,44 |
| | | 6 | 2,63 | 0,93 | 1,63 | 1,41 | 0,45 | 1,41 |
| 0 | 0,003 | 1 | 18,76 | 5,16 | 1,38 | 1,25 | 2,24 | 5,76 |
| | | 2 | 5,83 | 3,47 | 1,18 | 0,98 | 0,88 | 2,49 |
| | | 3 | 5,96 | 3,45 | 1,14 | 1,05 | 0,94 | 2,52 |
| | | 4 | 1,80 | 2,58 | 2,71 | 1,41 | 1,03 | 1,91 |
| | | 5 | 2,71 | 0,96 | 1,71 | 1,37 | 0,43 | 1,44 |
| | | 6 | 2,63 | 0,93 | 1,66 | 1,41 | 0,44 | 1,41 |
| 0,01 | 0 | 1 | 17,21 | 6,80 | 1,59 | 2,28 | 2,62 | 6,10 |
| | | 2 | 6,15 | 3,77 | 1,48 | 0,50 | 0,63 | 2,51 |
| | | 3 | 6,14 | 3,73 | 1,47 | 0,50 | 0,63 | 2,50 |
| | | 4 | 2,07 | 1,95 | 3,22 | 2,27 | 0,59 | 2,02 |
| | | 5 | 2,00 | 1,67 | 2,32 | 1,85 | 0,54 | 1,67 |
| | | 6 | 1,98 | 1,71 | 2,22 | 1,76 | 0,62 | 1,66 |
| 0,01 | 0,003 | 1 | 17,21 | 6,95 | 1,62 | 2,30 | 2,67 | 6,17 |
| | | 2 | 6,15 | 3,77 | 1,47 | 0,51 | 0,62 | 2,51 |
| | | 3 | 6,14 | 3,70 | 1,46 | 0,50 | 0,64 | 2,49 |
| | | 4 | 2,07 | 1,95 | 3,25 | 2,27 | 0,63 | 2,03 |
| | | 5 | 2,00 | 1,65 | 2,35 | 1,85 | 0,56 | 1,68 |
| | | 6 | 1,98 | 1,73 | 2,25 | 1,76 | 0,63 | 1,67 |

Т а б л и ц а 8

Средняя ошибка (в %) оценивания частоты форманты для шести русских гласных в классе моделей: $\Lambda(z^{-1})_{a_n}(1)$, $\Lambda(z^{-1})_{\nabla a_n}(2)$,
 $\Lambda(z^{-1})_{\nabla a_n}(3)$, $\Lambda(z^{-1})_{\nabla a_n}(4)$, $\Lambda(z^{-1})_{\nabla a_n}(5)$,
 $\Lambda(z^{-1})_{\nabla a_n}(6)$ при синусоидальном возбуждении и различных значениях σ_η и σ_κ ; $p = 14$.

| σ_η | σ_κ | Мо- дель | /а/ $\delta_{.1}$ | /о/ $\delta_{.2}$ | /у/ $\delta_{.3}$ | /э/ $\delta_{.4}$ | /ы/ $\delta_{.5}$ | /и/ $\delta_{.6}$ | Сред- нее δ |
|---------------|-----------------|-------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|-----------------------|
| 0 | 0 | 1 | 0,92 | 1,44 | 12,91 | 2,08 | 6,09 | 5,11 | 4,76 |
| | | 2 | 2,02 | 1,33 | 3,96 | 2,98 | 0,83 | 4,19 | 2,55 |
| | | 3 | 1,20 | 1,37 | 4,00 | 2,98 | 0,88 | 4,18 | 2,44 |
| | | 4 | 1,70 | 1,30 | 1,83 | 1,23 | 0,86 | 0,87 | 1,30 |
| | | 5 | 1,14 | 1,05 | 1,63 | 1,15 | 1,24 | 1,11 | 1,22 |
| | | 6 | 1,21 | 1,04 | 1,65 | 1,09 | 1,24 | 1,18 | 1,23 |
| 0 | 0,003 | 1 | 0,93 | 1,56 | 12,94 | 2,08 | 6,09 | 5,10 | 4,78 |
| | | 2 | 2,02 | 1,32 | 3,93 | 2,98 | 0,84 | 4,18 | 2,55 |
| | | 3 | 1,20 | 1,39 | 3,93 | 2,98 | 0,88 | 4,18 | 2,43 |
| | | 4 | 1,70 | 1,31 | 1,86 | 1,24 | 0,86 | 0,87 | 1,31 |
| | | 5 | 1,13 | 1,04 | 1,62 | 1,16 | 1,24 | 1,11 | 1,22 |
| | | 6 | 1,21 | 1,02 | 1,61 | 1,11 | 1,24 | 1,18 | 1,23 |
| 0,01 | 0 | 1 | 2,42 | 17,62 | 11,93 | 9,54 | 10,77 | 4,98 | 9,54 |
| | | 2 | 1,43 | 1,49 | 5,36 | 2,99 | 0,98 | 3,99 | 2,71 |
| | | 3 | 1,51 | 1,53 | 5,36 | 2,99 | 0,97 | 4,08 | 2,74 |
| | | 4 | 2,18 | 1,73 | 1,32 | 1,54 | 1,19 | 1,88 | 1,64 |
| | | 5 | 1,85 | 1,86 | 1,42 | 1,43 | 0,99 | 1,84 | 1,57 |
| | | 6 | 2,06 | 1,86 | 1,42 | 1,44 | 0,99 | 2,01 | 1,63 |
| 0,01 | 0,003 | 1 | 1,33 | 17,62 | 11,95 | 9,33 | 10,77 | 4,98 | 9,33 |
| | | 2 | 1,43 | 1,50 | 5,36 | 2,98 | 0,96 | 4,00 | 2,71 |
| | | 3 | 1,51 | 1,53 | 5,36 | 2,99 | 0,97 | 4,08 | 2,74 |
| | | 4 | 2,17 | 1,73 | 1,32 | 1,53 | 1,19 | 1,88 | 1,64 |
| | | 5 | 1,85 | 1,86 | 1,42 | 1,41 | 0,99 | 1,84 | 1,56 |
| | | 6 | 2,06 | 1,86 | 1,42 | 1,44 | 0,99 | 2,01 | 1,63 |

Ошибка (в %) оценивания пяти формантных частот,
 усредненная по шести русским гласным, в классе
 моделей: $\Lambda(z^{-1})_{a_n}(1)$; $\Lambda(z^{-1})_{\nabla a_n}(2)$; $\Lambda(z^{-1})_{\nabla a_n}(3)$;
 $\Lambda(z^{-1})_{\nabla_k a_n}(4)$; $\Lambda(z^{-1})_{\nabla \nabla_k a_n}(5)$; $\Lambda(z^{-1})_{\nabla a_n \nabla_k a_n}(6)$
 при синусоидальном возбуждении и различных значениях
 σ_k и σ_n ; при $p = 14$.

| σ_n | σ_k | Мо- дель | F_1 $\delta_{1.}$ | F_2 $\delta_{2.}$ | F_3 $\delta_{3.}$ | F_4 $\delta_{4.}$ | F_5 $\delta_{5.}$ | Сред- нее δ |
|------------|------------|-------------|------------------------|------------------------|------------------------|------------------------|------------------------|-----------------------|
| 0 | 0 | 1 | 14,81 | 4,68 | 1,90 | 1,60 | 0,81 | 4,76 |
| | | 2 | 6,39 | 4,01 | 0,46 | 1,00 | 0,90 | 2,55 |
| | | 3 | 6,10 | 3,68 | 0,49 | 1,03 | 0,86 | 2,44 |
| | | 4 | 1,89 | 1,96 | 0,82 | 0,91 | 0,91 | 1,3 |
| | | 5 | 2,04 | 1,60 | 0,52 | 1,01 | 0,96 | 1,22 |
| | | 6 | 2,12 | 1,60 | 0,51 | 0,97 | 0,96 | 1,23 |
| 0 | 0,003 | 1 | 14,81 | 4,67 | 1,96 | 1,66 | 0,83 | 4,78 |
| | | 2 | 6,39 | 3,98 | 0,47 | 0,98 | 0,93 | 2,55 |
| | | 3 | 6,04 | 3,68 | 0,50 | 1,00 | 0,87 | 2,43 |
| | | 4 | 1,96 | 1,93 | 0,83 | 0,91 | 0,90 | 1,31 |
| | | 5 | 2,04 | 1,60 | 0,50 | 0,97 | 0,97 | 1,22 |
| | | 6 | 2,12 | 1,59 | 0,50 | 0,95 | 0,98 | 1,23 |
| 0,01 | 0 | 1 | 16,23 | 19,32 | 6,84 | 4,33 | 0,99 | 9,54 |
| | | 2 | 6,61 | 4,13 | 1,16 | 0,81 | 0,83 | 2,71 |
| | | 3 | 6,71 | 4,15 | 1,16 | 0,82 | 0,85 | 2,74 |
| | | 4 | 1,33 | 2,34 | 1,54 | 1,82 | 1,18 | 1,64 |
| | | 5 | 1,75 | 1,78 | 1,71 | 1,41 | 1,18 | 1,57 |
| | | 6 | 2,04 | 1,87 | 1,76 | 1,34 | 1,15 | 1,63 |
| 0,01 | 0,003 | 1 | 15,17 | 19,28 | 6,81 | 4,39 | 0,99 | 9,33 |
| | | 2 | 6,61 | 4,12 | 1,15 | 0,82 | 0,83 | 2,71 |
| | | 3 | 6,71 | 4,15 | 1,16 | 0,83 | 0,85 | 2,74 |
| | | 4 | 1,33 | 2,34 | 1,54 | 1,81 | 1,18 | 1,64 |
| | | 5 | 1,75 | 1,78 | 1,70 | 1,40 | 1,18 | 1,56 |
| | | 6 | 2,04 | 1,87 | 1,76 | 1,32 | 1,16 | 1,63 |

Таким образом,

1. Оценивание формантных частот в классе моделей авторегрессии со стационарной сезонной разностью первого порядка и стационарной разностью первого порядка позволяет по сравнению с оценкой:

а) в классе моделей авторегрессии при одинаковом порядке моделей понизить ошибку оценивания в среднем в 18-21 раз при порядке, равном 10; в 4 раза, при порядке, равном 12; и в 4-6 раз при порядке, равном 14;

б) в классе моделей авторегрессии со стационарной разностью первого порядка при одинаковом порядке моделей понизить ошибку оценивания в среднем в 9-15 раз при порядке, равном 10, и в 2 раза при порядке, равном 12 и 14;

в) в классах моделей авторегрессии и авторегрессии со стационарной разностью первого порядка при одинаковой ошибке получить описание, экономичнее более чем в 1,4 раза, и избежать решения задачи отжестотвления полюсов модели с полюсами передаточной функции речевого тракта.

2. Ошибка квантования семипразрядного аналого-цифрового преобразователя практически не влияет на ошибку оценивания формантных частот.

3. Для выделения параметров речевого тракта рекомендуется использовать модель десятого порядка при частоте квантования сигнала 10 кГц.

4. Предложенный метод позволяет определять резонансные частоты речевого тракта с ошибкой, равной примерно 1,4-1,7%.

Л и т е р а т у р а

1. КИЛЬМАНОВ А.В. Алгоритм классификации тон/кум, основанный на критерии адекватности модели авторегрессии. - В кн.: Методы обработки информации. (Вычислительные системы, вып. 74.) Новосибирск, 1978, с. 123-148.

2. ФАНТ Г. Акустическая теория речеобразования. М., "Наука", 1964.

3. ЮВОВСКИЙ В.С. Аппроксимация отклика системы в z -пространстве и формантный анализ речи. - В кн.: Вычислительные системы, вып. 37. Новосибирск, 1969, с. 22-37.

4. СОБАКИН А.Н. Об определении формантных параметров голосового тракта по речевому сигналу с помощью ЦМ. - "Акустический журнал", 1972, т. XVII, вып. I, с. 106-114.

5. ATAL B.S., HANAUER S.L. Speech analysis and synthesis by linear prediction of the speech wave. - "J. Acoust. Soc. Amer.", 1971, v. 50, N 2, pt 2, p. 637-655.

6. MARKEL J.D. Application of digital inverse filter for automatic formant and F_0 analysis. - "IEEE Trans. Audio Electroacoust.", 1973, v.AU-21, N 3, p.154-160.

7. WAKITA H. Direct estimation of the vocal tract shape by inverse filtering of acoustic speech waveforms. - "IEEE Trans. Audio Electroacoust.", 1973, v.AU-21, N 5, p.417-427.

8. БОКС Др., ДИВЕНКИНС Г. Анализ временных рядов, прогноза и управление. Том I. М., "Мир", 1974.

9. ЛОЗОВСКИЙ В.С. Процедура решения полиномиальных уравнений методом Муллера. - В кн.: Вычислительные системы. Вып.44. Новосибирск, 1971, с. 155-161.

10. КУНЦ К.С. Численный анализ. Киев, "Техника", 1964.

Поступила в ред.-над.отд.

14 апреля 1978 года