

СИНТЕТИЧЕСКАЯ ТЕЛЕФОНΙΑ

УДК 534.78: 621.39: 681.513.6

ИССЛЕДОВАНИЕ АЛГОРИТМА КАЛМАНОВСКОЙ ФИЛЬТРАЦИИ РЕЧЕВОГО СИГНАЛА, НАБЛЮДАЕМОГО В ШУМЕ

В.Г. Санников, профессор МТУСИ, к.т.н.; tes_mtuci@mail.ru

Ключевые слова: речевой сигнал, шум, идентификация, показатель затухания, оптимальная линейная фильтрация, отношение сигнал/шум.

Введение. В системе речевой связи, простейшая схема которой приведена на рис. 1, речевой сигнал (РС) $x(t)$ формируется под действием сигнала $u(t)$ источника голосового возбуждения (ИГВ) на выходе голосового тракта (ГТ) модели речеобразования [1–3]. При этом наблюдению подлежит не сам РС, а его зашумленная копия $r(t)=x(t)+n(t)$, где $n(t)$ – реализация шума передачи.

В случае, когда уровень шума больше или равен уровню РС, а отношение сигнал/шум (ОСШ_{вх}) меньше или равно 0 дБ (мобильная телефония), качество восприятия РС по показателю разборчивости резко ухудшается. Для улучшения отношения сигнал/шум (ОСШ_{вых}) на выходе оптимального линейного фильтра (ОЛФ) по сравнению с ОСШ_{вх} и, соответственно, для повышения качества восприятия речи актуальной становится задача оптимальной фильтрации наблюдаемой смеси $r(t)$.

Известны примеры решения задачи оптимальной линейной фильтрации применительно к обработке зашумленной речи [1–3]. Однако большинство из них основаны на представлении РС $x(t)$ моделью авторегрессии, ограниченной полюсной моделью ГТ, а рекуррентная оценка параметров авторегрессии осуществляется на основе метода наименьших квадратов (МНК), обладающего бесконечной памятью. Такая оценка справедлива для случайных сигналов, относящихся к классу стационарных.

Речевые сигналы входят в класс нестационарных случайных сигналов. Поэтому для точной оценки динамики изменения во времени статистики более адекватен МНК с взвешиванием (с конечной памятью) [4], характеризуемый «множителем забывания» v . Кроме того, в [1] указывается, что более адекватной моделью голосового тракта является его нуль-полюсная модель или модель РС, построенная на основе авторегрессии и скользящего среднего (АРСС модель РС). Наиболее заметно это при анализе и синтезе назализованных звуков речи.

Поскольку в литературе вопросам фильтрации зашумленной речи, учитывающим отмеченные выше особенности модели речеобразования, уделено незначительное внимание, актуально решение следующих задач:

- теоретическое и экспериментальное исследование ОЛФ зашумленной речи с учетом нуль-полюсной или АРСС модели РС;
- оценка параметров ОЛФ на основе МНК с конечной памятью;

- оптимизация «множителя забывания» v_{opt} , обеспечивающего экстремум показателю качества оптимальной фильтрации зашумленной речи.

При реализации алгоритмов фильтрации сигналов часто применяют цифровые методы, а сигналы рассматриваются в дискретном времени. Оптимальная фильтрация предполагает наличие априорных сведений о структуре ГТ модели речеобразования. Поэтому сначала определим АРСС модель РС.

Идентификация АРСС модели РС. Нуль-полюсная или АРСС модель РС определяется разностным уравнением

$$x_t = \sum_{i=1}^m a_{i,t} x_{t-i} + \sum_{j=0}^{m-1} b_{j,t} u_{t-j}, \tag{1}$$

показывающим РС как отклик $x_t, t = 0, 1, 2, \dots$, устойчивой нестационарной динамической линейной системы, характеризующей ГТ модели речеобразования (рис. 1). На ее вход действует случайная последовательность $u_t, t = 0, 1, 2, \dots$, вырабатываемая ИГВ.

Идентификация (в узком смысле) представляет собой процесс определения вектора параметров модели РС (1) по результатам измерения значений $\{x_t\}, t = 0, 1, 2, \dots$. Пусть параметры $\{a_{i,t}\}, i = 1, m$, и $\{b_{j,t}\}, j = 0, m-1$, модели (1) неизвестны.

Введем обобщенный вектор параметров

$$\begin{aligned} \mathbf{a}_t &= [a_{1,t}, \dots, a_{m,t}, b_{0,t}, \dots, b_{m-1,t}]^T = \\ &= [a_{1,t}, \dots, a_{m,t}, a_{m+1,t}, a_{m+2,t}, \dots, a_{2m,t}]^T \end{aligned} \tag{2}$$

и расширенный вектор результатов измерений

$$\mathbf{s}_t = [x_{t-1}, \dots, x_{t-m}, u_t, \dots, u_{t-m+1}]^T, \tag{3}$$

где T – знак векторного транспонирования.

Для оценки искомого вектора параметров (2) модели (1) воспользуемся рекуррентным алгоритмом МНК с взвешиванием, при котором минимизации подвергается взвешенная сумма квадратов «невязки» [5]:

$$\overline{\sigma}_t^2(v) = \sum_{n=0}^t v^{t-k} (x_n - \mathbf{s}_n^T \mathbf{a}_n)^2, \tag{4}$$

где $v \leq 1$ – «множитель забывания», характеризующий конечную память алгоритма МНК.

В результате последовательная (рекуррентная) оценка $\hat{\mathbf{a}}_n, n = 0, 1, 2, \dots$, вектора параметров (2) АРСС модели РС (1) осуществляется следующим образом [4]:

$$\hat{\mathbf{a}}_n = \hat{\mathbf{a}}_{n-1} + \mathbf{k}_{a,n} (x_n - \hat{\mathbf{s}}_{n-1}^T \hat{\mathbf{a}}_{n-1}), \quad n = 1, 2, \dots; \tag{5}$$

$$\mathbf{k}_{a,n} = \mathbf{P}_{n-1} \hat{\mathbf{s}}_{n-1} / (v + c_n), \quad c_n = \hat{\mathbf{s}}_{n-1}^T \mathbf{P}_{n-1} \hat{\mathbf{s}}_{n-1}; \tag{6}$$

$$\mathbf{P}_n = (\mathbf{E}_{2m-1} - \mathbf{k}_{a,n} \hat{\mathbf{s}}_{n-1}^T) \mathbf{P}_{n-1} / v \tag{7}$$

с начальными условиями: $b_{0,0} = a_{m+1,0} = 1, \mathbf{a}_0 = [0, \dots, 0, 1, 0, \dots, 0]^T, \mathbf{P}_0 = p \mathbf{E}_{2m-1}$. Здесь \mathbf{E}_{2m-1} – единичная $(2m-1) \times (2m-1)$ матрица; p – достаточно большое положительное число (выбрано

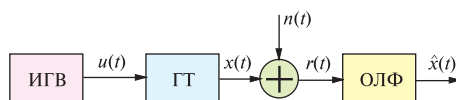


Рис. 1

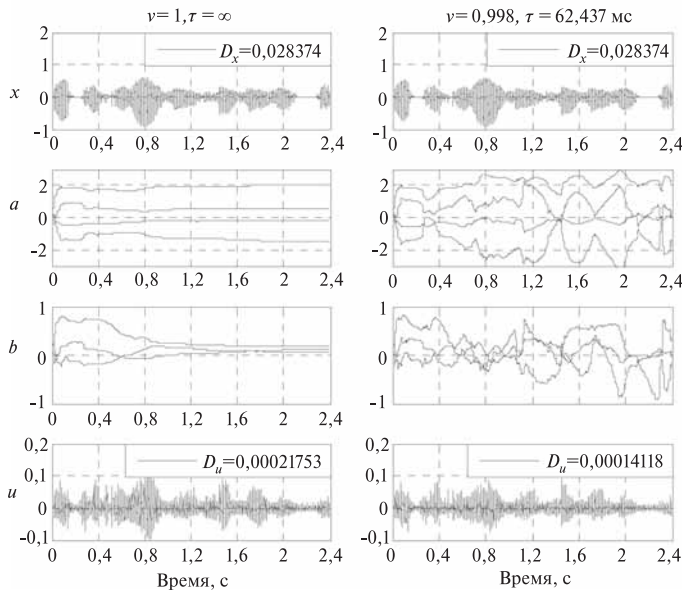


Рис. 2

$p=200$); $\mathbf{k}_{a,n}$ – вектор коррекции сигнала голосового возбуждения; $\mathbf{P}_n = \mathbf{R}_{s,n}^{-1}$ – матрица размера $(2m-1) \times (2m-1)$, обратная корреляционной матрице $\mathbf{R}_{s,n} = M\{\hat{\mathbf{s}}_n \hat{\mathbf{s}}_n^T\}$ случайного вектора $\hat{\mathbf{s}}_n$; M – знак математического ожидания.

Заметим, что вектор \mathbf{s}_t в (3) содержит ненаблюдаемые величины $u_{t-1}, \dots, u_{t-m+1}$. Их оценки определяют по оценкам параметров АРСС:

$$\hat{u}_t = (x_t - \hat{\mathbf{s}}_{t-1}^T \hat{\mathbf{a}}_t), \tag{8}$$

где $\hat{\mathbf{s}}_t = [x_{t-1}, \dots, x_{t-m}, \hat{u}_{t-1}, \dots, \hat{u}_{t-m+1}]$.

Результаты экспериментальной проверки алгоритма идентификации АРСС модели реального РС приведены на рис. 2. В качестве фонетически сбалансированной обрабатывалась стандартная фраза [6]: «Эти жирные сазаны ушли под палубу» (верхние графики) с дисперсией $D_x=0,28374$. Интервал дискретизации выбирался равным $\Delta t = 125$ мкс. Графики АР $\{a_{i,t}\}$, $i = 1,4$, и СС $\{b_{j,t}\}$, $j = 1,3$, параметров, а также оценки \hat{u}_t сигнала ИГВ изображены в левой и правой колонках с различными «множителями забывания»: 1) $\nu=1$ с интервалом памяти МНК $\tau \rightarrow \infty$ (левая колонка); 2) $\nu=0,998$ с $\tau=62,437$ мкс (правая колонка). Здесь $\nu = \exp(-\Delta t / \tau)$, $\tau = \Delta t / \ln(1 / \nu)$.

Анализ графиков рис. 2 позволяет сделать следующие выводы. При $\nu=1$ параметры АРСС модели РС существенно изменяются только на незначительном временном интервале, после которого они практически постоянны и не отслеживают изменение статистики РС. При $\nu=0,998$ динамика параметров существенно изменяется: они более детально описывают изменение статистики РС, приводя к уменьшению дисперсии сигнала ИГВ от $D_u=0,00021753$ при $\nu=1$ до $D_u=0,00014118$ при $\nu=0,998$.

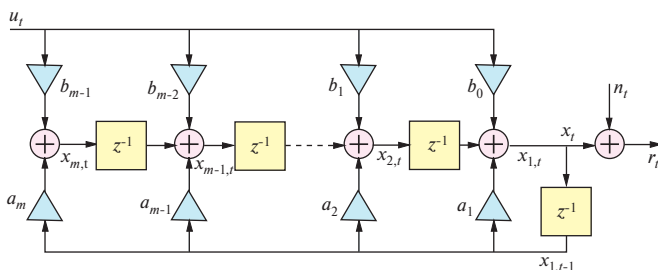


Рис. 3

Для конструктивного использования оптимальной калмановской фильтрации зашумленной речи, определяемой векторно-матричными уравнениями состояния и наблюдения, необходимо от скалярной модели ГТ (1) перейти к ее векторно-матричному представлению.

Уравнения состояния и наблюдения ГТ. Учитывая, что параметры АРСС модели РС можно считать постоянными на интервале длительностью 20 мс [1], построим векторную модель РС в пространстве состояний. Для этого введем переменные состояния ГТ $\{x_{k,t}\}$, $k = 1, m$, следующим образом:

$$\begin{cases} x_{m,t} = (a_m x_{t-1} + b_{m-1} u_t); \\ x_{m-1,t} = (a_{m-1} x_{t-1} + b_{m-2} u_t) + a_m x_{t-2} + b_{m-1} u_{t-1} = \\ = (a_{m-1} x_{t-1} + b_{m-2} u_t) + x_{m,t-1}; \\ x_{m-2,t} = (a_{m-2} x_{t-1} + b_{m-3} u_t) + a_{m-1} x_{t-2} + b_{m-2} u_{t-1} + a_m x_{t-3} + \\ + b_{m-1} u_{t-2} = (a_{m-2} x_{t-1} + b_{m-3} u_t) + x_{m-1,t-1}; \\ \dots \\ x_{1,t} = (a_1 x_{t-1} + b_0 u_t) + a_2 x_{t-2} + b_1 u_{t-1} + \dots + a_m x_{t-m} + \\ + b_{m-1} u_{t-m+1} = (a_1 x_{t-1} + b_0 u_t) + x_{2,t-1}, \quad b_0 = 1. \end{cases} \tag{9}$$

Сравнивая последнюю строку в (9) и соотношение (1), заключаем, что $x_t \equiv x_{1,t}$, $x_{t-1} \equiv x_{1,t-1}$. Теперь, с учетом шумов наблюдения, модель ГТ представляется:

- уравнением состояния

$$\mathbf{x}_t = \mathbf{A} \mathbf{x}_{t-1} + \mathbf{b} u_t, \tag{10}$$

- уравнением наблюдения

$$r_t = \mathbf{h}^T \mathbf{x}_t + n_t = x_t + n_t, \tag{11}$$

где $\mathbf{x}_t = [x_{1,t}, x_{2,t}, \dots, x_{m,t}]^T$ – вектор состояния ГТ в момент t ;

$$\mathbf{A} = \begin{bmatrix} a_1 & 1 & 0 & \dots & 0 \\ a_2 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ a_{m-1} & 0 & 0 & \dots & 1 \\ a_m & 0 & 0 & \dots & 0 \end{bmatrix}$$

– переходная матрица ГТ размера $m \times m$, составленная из параметров авторегрессии (АР) модели (1); $\mathbf{b} = [1, b_1, \dots, b_{m-1}]^T$ – вектор параметров скользящего среднего (СС) модели (1); $\mathbf{h} = [1, 0, 0, \dots, 0]^T$ – вектор наблюдения РС; n_t – гауссовская шумовая последовательность с корреляционной функцией $R_{n,t} = \sigma_{n,t}^2 \delta_{t,\tau}$; $\delta_{t,\tau}$ – символ Кронекера, n_t и u_t некоррелированы.

Структурная схема ГТ, соответствующая уравнениям состояния (10) и наблюдения (11), показана на рис. 3.

Оптимальная линейная фильтрация РС. Задача фильтрации состоит в том, чтобы по совокупности последовательных измерений $\{r_1, r_2, \dots, r_t\} = \{r_k\}$, $k = \overline{1, t}$, найти оптимальную оценку $\hat{\mathbf{x}}_t$ вектора состояния \mathbf{x}_t , удовлетворяющую критерию минимума текущей среднеквадратической погрешности фильтрации $\sigma_t^2 = M\{\mathbf{e}_t^T \mathbf{e}_t\}$, где $\mathbf{e}_t = (\mathbf{x}_t - \hat{\mathbf{x}}_t)$.

Для синтеза ОЛФ, в соответствии с полученной моделью ГТ, можно воспользоваться хорошо разработанной теорией ОЛФ в пространстве состояний, приводящей к фильтру Калмана, работа которого определяется соотношениями [4]:

$$\begin{aligned} \hat{\mathbf{x}}_t &= \mathbf{A}_{t-1} \hat{\mathbf{x}}_{t-1} + \mathbf{k}_{x,t} [r_t - \mathbf{h}^T \mathbf{A}_{t-1} \hat{\mathbf{x}}_{t-1}] = \\ &= \mathbf{A}_{t-1} \hat{\mathbf{x}}_{t-1} + \mathbf{k}_{x,t} [r_t - (a_{1,t-1} \hat{x}_{1,t-1} + \hat{x}_{2,t-1})], \quad t = 1, 2, \dots; \end{aligned} \tag{12}$$

$$\begin{aligned} \mathbf{k}_{x,t} &= \mathbf{P}_{t|t-1} \mathbf{h} (R_{n,t} + \mathbf{h}^T \mathbf{P}_{t|t-1} \mathbf{h})^{-1} = \\ &= \mathbf{V}_t \mathbf{h} R_{n,t}^{-1} = [V_{11}, V_{21}, \dots, V_{m1}]^T / \sigma_{n,t}^2; \end{aligned} \tag{13}$$

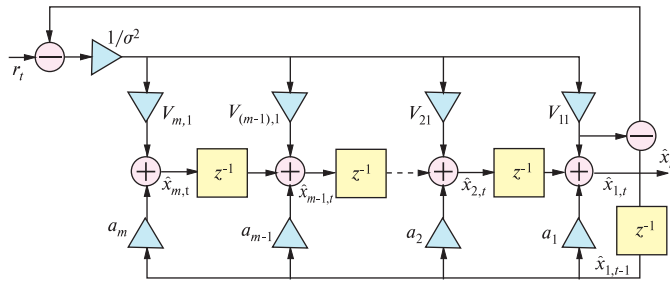


Рис. 4

$$P_{t|t-1} = A_{t-1} V_{t-1} A_{t-1}^T + b_{t-1} R_{u,t-1} b_{t-1}^T; \tag{14}$$

$$V_t = (E_m - k_{x,t} h^T) P_{t|t-1} \tag{15}$$

с начальными условиями: $x_0 = M\{x_0\} = \bar{x}_0$, $V_0 = M\{(x_0 - \bar{x}_0)(x_0 - \bar{x}_0)^T\}$. Здесь $k_{x,t}$ – векторный коэффициент усиления ОЛФ; $P_{t|t-1}$ – корреляционная матрица погрешности экстраполяции на один шаг (априорная корреляционная матрица погрешности фильтрации); $V_t = [V_{11}, V_{21}, \dots, V_{m1}]^T$ – апостериорная корреляционная матрица погрешности фильтрации; $R_{u,t}$ – корреляционная матрица сигнала возбуждения ИГВ; E_m – единичная $m \times m$ матрица; $\sigma_{n,t}^2$ – текущая дисперсия шума наблюдения.

Структурная схема ОЛФ для оценки РС, наблюдаемого на фоне шума, равномерно распределенного в полосе частот сигнала, показана на рис. 4. Здесь полюса системной функции цифрового ОЛФ (нижняя часть схемы) определяются АР параметрами $\{a_j\}$, $j = 1, m$, модели ГТ, в то время как ее нули характеризуются не СС параметрами $\{b_j\}$, $j = 0, m - 1$, модели ГТ (рис. 3), а параметрами $\{V_{j1}\}$, $j = 1, m$, первого столбца матрицы погрешности фильтрации (верхняя часть схемы).

Результаты экспериментальной проверки алгоритма оптимальной линейной (калмановской) фильтрации зашумленного РС представлены рис. 5 и таблицей, в которую сведены величины v_{opt} , $\tau_{opt} = \Delta t / \ln(1 / v_{opt})$ и максимальные

m	2	4	6	8	10	12	14
v_{opt}	0,9990	0,9950	0,9930	0,9916	0,9915	0,9915	0,9915
τ_{opt} , мс	124,94	24,937	17,795	14,818	14,643	14,643	14,643
ОСШ _{вых} , дБ	5,6768	6,5224	7,1988	7,2769	7,6234	7,4988	7,2981

значения ОСШ_{вых} в зависимости от m . Результаты получены для ОСШ_{вх} = 0 дБ.

Как и ранее обработке подвергалась зашумленная фраза: «Эти жирные сазаны ушли под палубу» с различными ОСШ_{вх} = $10 \lg(\sigma_x^2 / \sigma_n^2)$, дБ. Экспериментальные данные получены для различного числа АР и СС параметров модели ГТ, а также для разных значений «множителя забывания» v МНК с взвешиванием. Видно, что при фиксированной памяти моделей ГТ и ОЛФ, определяемой параметром m , значение ОСШ_{вых} = $10 \lg[\sigma_x^2 / (x_t - x_t)^2]$ изменяется так, что при некоторой величине v_{opt} принимает максимальное значение.

Наличие v_{opt} объясняется следующим. При $v \rightarrow 1$ память алгоритма МНК бесконечна ($\tau \rightarrow \infty$), а корреляционная матрица речи в (7) и параметры предсказания в (5) значительно усреднены и не успевают отслеживать изменение статистики РС. При уменьшении v и τ возрастает вес краткосрочной корреляции (P_n в (7) увеличивается). Следовательно, корреляционная матрица речи в (7) и параметры предсказания в (5) начинают лучше отслеживать изменение статистики РС. Это приводит к уменьшению погрешности предсказания и фильтрации, а также возрастанию ОСШ_{вых}. Однако при дальнейшем уменьшении v и τ дисперсии параметров предсказания возрастают из-за увеличения в них вредной доли сигнала голосового возбуждения, что приводит к ухудшению предсказания и в результате к уменьшению ОСШ оптимальной фильтрации РС.

Графическая иллюстрация результатов эксперимента приведена на рис. 5. Здесь в правой колонке (сверху вниз) показаны фрагменты: исходного РС $\{x_t\}$, зашумленного РС

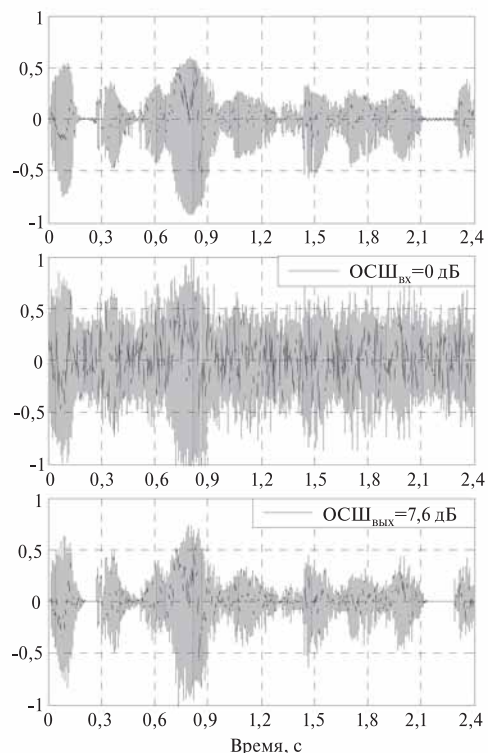
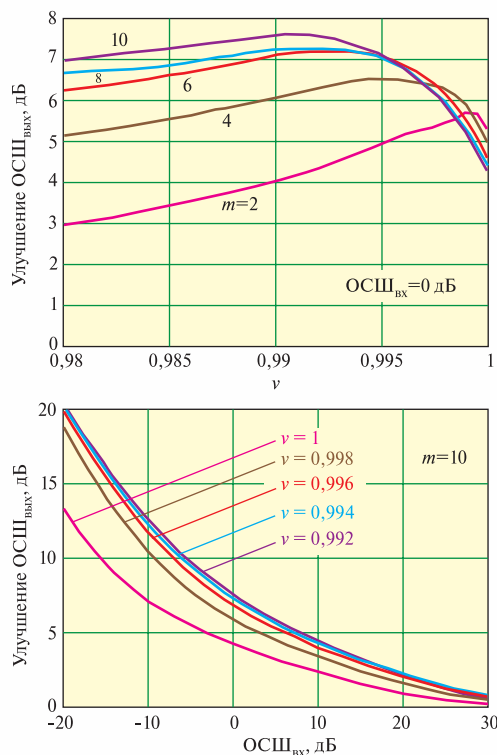


Рис. 5

$\{r_t\}$ и оценки РС $\{x_t\}$, наблюдаемой на выходе ОЛФ. Нетрудно заметить, что фильтр значительно подавляет шум, особенно в паузах речи.

В левой колонке (верхние графики) показаны зависимости ОСШ_{вых} от «множителя забывания» ν при различных m и ОСШ_{вх} = 0 дБ. Явно виден экстремальный характер этих зависимостей, оптимальные значения которых приведены в таблице.

На нижнем графике показаны кривые улучшения ОСШ_{вых} (ОСШ_{вых} – ОСШ_{вх}), обеспечиваемого ОЛФ в зависимости от ОСШ_{вх} при различных ν и $m=10$. Из графиков следует, что при фиксированном ОСШ_{вх} МНК при $\nu=1$ уступает по эффективности фильтрации МНК с $\nu<1$. Так, при ОСШ_{вх} = 0 дБ и $\nu=1$ (интервал памяти $\tau \rightarrow \infty$) ОСШ_{вых} = 4,4 дБ, в то время как при $\nu=0,992$ (интервал памяти $\tau=14,85$ мс) ОСШ_{вых} = 7,62 дБ, т.е. на 3,22 дБ больше.

Заключение. Дальнейшее совершенствование методов оптимальной фильтрации зашумленной речи с повышенной помехозащищенностью возможно только при использовании априорных сведений, более адекватных моделям формирования и слухового приема реальных речевых сигналов.

Это направление развивается в лаборатории речевой информатики и связи МТУСИ и представляет значительный интерес для идентификации личности говорящего по голосу и речи, являющейся одной из приоритетных задач, решаемых на базе комбинированной обработки речевых сигналов, повышения разборчивости речи и голосовой биометрии.

ЛИТЕРАТУРА

1. **Маркел Дж. Д., Грэй А.Х.** Линейное предсказание речи: Пер. с англ./ Под ред. Ю.Н. Прохорова и В.С. Звездина. – М.: Связь, 1980. – 308 с.
2. **Прохоров Ю.Н.** Статистические модели и рекуррентное предсказание речевых сигналов. – М.: Радио и связь, 1984. – 240 с.
3. **Шелухин О.И., Лукьянцев Н.Ф.** Цифровая обработка и передача речи/ Под ред. О.И. Шелухина. – М.: Радио и связь, 2000. – 456 с.
4. **Сейдж Э., Мелс Дж.** Теория оценивания и её применение в связи и управлении: Пер. с англ./ Под ред. Б.Р. Левина. – М.: Связь, 1976. – 496 с.
5. **Санников В.Г.** Статистический анализ методов формирования речевых сигналов. – М.: МТУСИ, 2005. – 140 с.
6. ГОСТ Р 50840-95. Передача речи по трактам связи. Методы оценки качества, разборчивости и узнаваемости.

Получено после доработки 11.03.11