

Filtros Adaptativos do tipo LMS

O uso do critério supervisionado MMSE em conjunto com uma estrutura linear caracteriza o paradigma de filtragem de Wiener. É importante ter em mente, porém, que o filtro de Wiener baseia-se em duas hipóteses: (i) os sinais envolvidos são WSS e (ii) as médias estatísticas R_x e p_{xd} são conhecidos. Em muitos casos práticos, dois aspectos violam estas suposições:

- a necessidade de operar em tempo real: requer métodos capazes de conjuntamente efetuar a aquisição e a otimização.

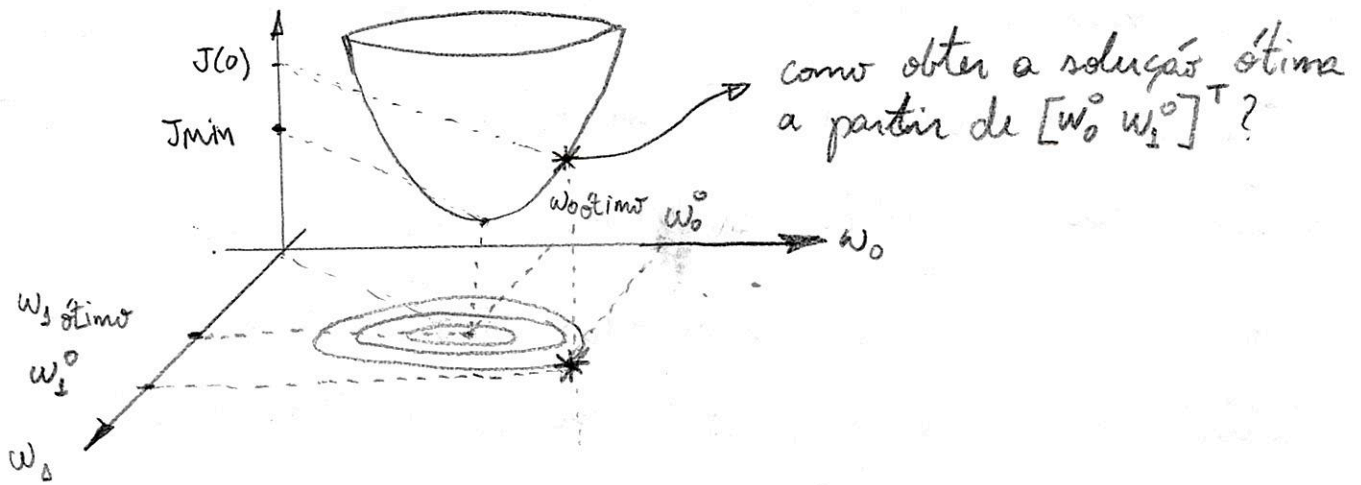
- a presença de sinais não-estacionários: inibe o uso de uma solução fechada, uma vez que não há mais valores fixos de correlação estatística.

Este novo cenário nos leva à fronteira entre filtragem ótima e adaptativa. Neste último caso, as soluções para o problema de filtragem (linear) serão determinadas de maneira iterativa/recursiva enquanto os dados são obtidos.

* Algoritmo do gradiente Determinístico.

- conhecido como steepest-descent algorithm

$$J(\underline{w}) = E\{e^2(n)\}$$



- em vez de obter diretamente $\underline{w}^{\text{ótimo}} = \underline{R}_x^{-1} \underline{p}_x$, partiremos de uma condição inicial $\underline{w}(0)$ e iterativamente atualizaremos o vetor de parâmetros na forma $\underline{w}(n+1) = f\{\underline{w}(n)\}$ até um certo critério de parada ser atingido.

- o método do gradiente expressa a função $J(\underline{w})$ apenas em termos de 1ª derivada (vetor gradiente)

$$J(\underline{w}) \Big|_{\underline{w}_{i+1}} = J(\underline{w}) \Big|_{\underline{w}_i} + \frac{\partial J(\underline{w})}{\partial \underline{w}^T} \Big|_{\underline{w}_i} \cdot \underline{\Delta w} + \frac{1}{2} \underline{\Delta w}^T \frac{\partial^2 J(\underline{w})}{\partial \underline{w}_i \partial \underline{w}_i^T} \Big|_{\underline{w}_i} \cdot \underline{\Delta w} + \dots \quad (\text{Taylor})$$

$\underline{\Delta w} = \underline{w}_{i+1} - \underline{w}_i$

- para obtermos o maior decréscimo no valor de $J(\underline{w})$, a variação $\underline{\Delta w}$ deve ser colinear ao vetor gradiente, mas com sentido oposto (se forem ortogonais, o produto escalar é nulo); (havendo um ângulo entre $\underline{\Delta w}$ e o vetor gradiente, a redução acontece cada vez mais conforme ele tende a 0°).

$$\Rightarrow \underline{\Delta w} = - \text{constante} \cdot \frac{\partial J(\underline{w})}{\partial \underline{w}^T}$$

Ora, $\nabla J(\underline{w}) = -2 \underline{p} \times d + 2 R \times \underline{w}$, como deduzido durante a discussão sobre filtragem ótima. Com isto,

$$\begin{aligned} \underline{w}_{i+1} &= \underline{w}_i - \text{constante} \frac{\partial J(\underline{w})}{\partial \underline{w}^T} \\ &= \underline{w}_i - 2 \text{constante} [R \times \underline{w}_i - \underline{p} \times d] \end{aligned}$$

Substituindo a constante por $\mu/2$, chegamos a:

$$\boxed{\underline{w}_{i+1} = \underline{w}_i - \mu (R \times \underline{w}_i - \underline{p} \times d)}$$

μ = tamanho do passo a ser dado na direção definida pelo gradiente da função custo $J(\underline{w})$.

Análise de Convergência:

• os pontos de equilíbrio de um sistema dinâmico são invariantes com relação ao processo iterativo. Em outras palavras, eles são as soluções da equação $\underline{w}_{i+1} = \underline{w}_i$. Logo,

$$\mu [R \times \underline{w} - \underline{p} \times d] = 0 \Rightarrow \underline{w} = \underline{R_x^{-1} p \times d}$$

- o sistema possui um único ponto de equilíbrio e que corresponde à própria solução de Wiener.

O sistema em questão é estável?

• para verificarmos a estabilidade do sistema dinâmico, primeiro reescrevemos a equação de atualização do vetor de coeficientes da seguinte forma:

$$\underline{w}_{i+1} = [\mathbf{I} - \mu \mathbf{R}_x] \underline{w}_i + \mu \mathbf{p}_{\times d}$$

Para que o sistema seja estável, é preciso que os autovalores de $[\mathbf{I} - \mu \mathbf{R}_x]$ estejam dentro da circunferência de raio unitário.

Seja $\mathbf{B} = [\mathbf{I} - \mu \mathbf{R}_x]$, então, $\lambda_{\mathbf{B}} = 1 - \mu \lambda_{\mathbf{R}_x}$

Para que $|\lambda_{\mathbf{B}}| < 1$, $|1 - \mu \lambda_{\mathbf{R}_x}| < 1$.

O caso mais restritivo está associado ao autovalor de \mathbf{R}_x de maior módulo.

Assim, $|1 - \mu \lambda_{\mathbf{R}_x}^{\max}| < 1 \Rightarrow -1 < 1 - \mu \lambda_{\mathbf{R}_x}^{\max} < 1$

Como $\lambda_{\mathbf{R}_x}^{\max}$ e μ são números não-negativos, chegamos

a
$$\mu < \frac{2}{\lambda_{\mathbf{R}_x}^{\max}}$$
.

\Rightarrow O algoritmo do gradiente determinístico não possui um sentido prático para aplicação em filtragem, pois pressupõe conhecimento total das informações estatísticas em \mathbf{R}_x e $\mathbf{p}_{\times d}$, além de estacionariedade.

Sua derivação e análise, porém, serve de base para os algoritmos estocásticos que veremos.

* Método de Newton

- baseado na aproximação de 2ª ordem de $J(\underline{w})$
- obtemos um $\underline{\Delta w}$ que nos leva diretamente ao ponto de mínimo, ou seja, no ponto no qual $\frac{\partial J}{\partial \underline{w}^T}(\underline{w}_{i+1}) = 0$.

Derivando a expansão de Taylor de $J(\underline{w})$:

$$0 = \underbrace{\frac{\partial J(\underline{w})}{\partial \underline{w}^T} \Big|_{\underline{w}_i}}_{\text{vetor gradiente em } \underline{w}_i} + \underbrace{\frac{\partial^2 J(\underline{w})}{\partial \underline{w}^T \partial \underline{w}} \Big|_{\underline{w}_i}}_{\text{matriz hessiana em } \underline{w}_i : H(J(\underline{w}))} \cdot \underline{\Delta w}, \quad \underline{\Delta w} = \underline{w}_{i+1} - \underline{w}_i$$

$$\begin{aligned} \Rightarrow \underline{w}_{i+1} &= \underline{w}_i - H(J(\underline{w}))^{-1} \nabla J(\underline{w}) \\ &= \underline{w}_i - H(J(\underline{w}))^{-1} [-2 \underline{p}_{xd} + 2 R_x \underline{w}_i] \end{aligned}$$

Quem é $H(J(\underline{w}))$?

$\frac{\partial E\{e^2(n)\}}{\partial \underline{w}^T \partial \underline{w}}$ - cada elemento (i, j) desta matriz corresponde

$$a \quad H_{ij}(J(\underline{w})) = \frac{\partial E\{e^2(n)\}}{\partial w_i \partial w_j}$$

$$\frac{\partial E\{e^2(n)\}}{\partial w_i} = E \left\{ z e(n) \frac{\partial e(n)}{\partial w_i} \right\} = 2 E \{ e(n) x(n-i) \}$$

$$\text{Logo, } \frac{\partial E\{e^2(n)\}}{\partial w_i w_j} = \frac{\partial [2 E\{e(n)x(n-i)\}]}{\partial w_j} = 2 E\{x(n-i)x(n-j)\} = 2r_x(i-j)$$

Portanto, $\boxed{H(J(\underline{w})) = 2 \cdot R_x}$

Retomando à equação de atualização do vetor de coeficientes:

$$\underline{w}_{i+1} = \underline{w}_i - \frac{1}{2} R_x^{-1} (-2 p_{xd} + 2 R_x \underline{w}_i)$$

Note que podemos escrever:

$$\begin{aligned} \underline{w}_{i+1} &= I \underline{w}_i + R_x^{-1} p_{xd} - I \underline{w}_i \\ &= \underline{R_x^{-1} p_{xd}} \end{aligned}$$

Ou seja, para qualquer ponto de partida inicial, chegamos ao ponto ótimo em um único passo.

* Algoritmo do gradiente estocástico (LMS)

• 1959/60: Widrow e Hoff propuseram um algoritmo iterativo de busca de solução ótima MSE (Wiener) combinado com a ideia de aproximação estocástica (Robbins e Monro, 1951).

Características:

- * aquisições dos sinais e otimizações do filtro realizadas concomitentemente.
- * possibilidade de se trabalhar em contextos não-estecionários.
- * simplicidade computacional
- * boas propriedades de convergência.

Least-Mean-Square (LMS): realiza a adaptação do vetor de coeficientes com base no gradiente estocástico.

$\hat{R}_x = \underline{x}(n) \underline{x}^T(n)$ } os valores esperados (médios) são substituídos
 $\hat{p}_{xd} = \underline{x}(n) d(n)$ } pelos valores instantâneos a partir de
estimativas não-polarizados.

Substituindo as estimativas na equação do gradiente determinístico:

$$\begin{aligned}\underline{w}_{i+1} &= \underline{w}_i - \mu (\hat{R}_x \underline{w}_i - \hat{p}_{xd}) \\ &= \underline{w}_i - \mu \underline{x}(i) \underline{x}^T(i) \underline{w}_i + \mu \underline{x}(i) d(i)\end{aligned}$$

$$\text{Mas, } \underline{x}^T(i) \underline{w}(i) = \underline{w}^T(i) \underline{x}(i) = \hat{d}(i)$$

$$\text{Então, } \underline{w}_{i+1} = \underline{w}_i + \mu \underline{x}(i) \underbrace{(d(i) - \hat{d}(i))}_{e(n)}$$

Finalmente, mudando a notação para o índice n , obtemos:

$$\underline{w}(n+1) = \underline{w}(n) + \mu \underline{x}(n) e(n)$$

- Este algoritmo pode ser igualmente obtido minimizando-se diretamente, através do método do gradiente, a função custo $J(n) = e^2(n)$.

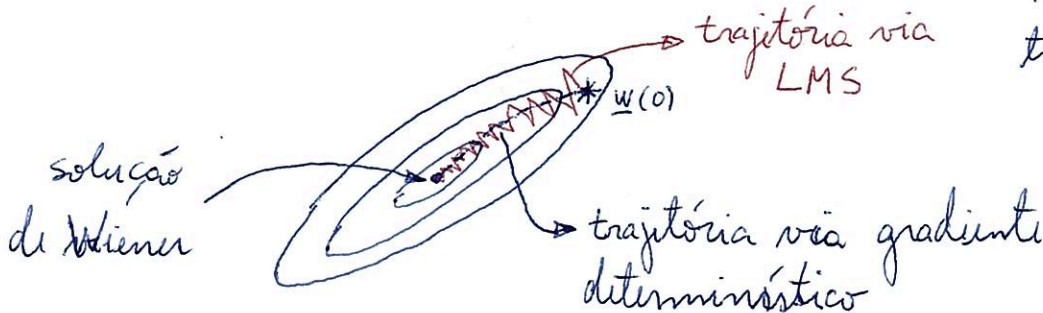
Início: $\underline{w}(0) = [0 \ 0 \ \dots \ 0]^T$

$$\underline{x}(0) = [0 \ 0 \ \dots \ 0]^T$$

$x(n)$ entre \rightarrow $x(n-M+1)$ sai

A cada observação $x(n)$, obtemos $e(n)$ (conhecendo a saída desejada $d(n)$) e atualizo o vetor de coeficientes.

$$\underline{w}(n+1) = f(x(n), x(n-1), x(n-2), \dots, x(0)) \rightarrow \text{ponho uma memória de todos os amostras recebidos.}$$



Propriedades do LMS:

- o gradiente estocástico pode ser visto como uma estimativa não-polarizada do gradiente ideal, dado um vetor fixo de coeficientes \underline{w} .

• convergência na média:

Seja $\underline{\Delta w}(n) = \underline{w}(n) - \underline{w}_{\text{ótimo}}$. Então, $E\{\underline{\Delta w}(n)\} \rightarrow 0$ para $n \rightarrow \infty$

Isso quer dizer que a tendência de se observar - quando o tamanho do passo é adequadamente escolhido - é um comportamento oscilatório em torno de $\underline{w}_{\text{ótimo}}$.

contudo, podemos ter situações distintas quanto à amplitude das oscilações:



Por isso, também é interessante observar o tamanho do desvio em relação ao coeficiente (ou, semelhantemente, ao erro) ótimo.

\Rightarrow Teoria da independência: hipótese de que os vetores $\underline{x}(n)$, para todo $n = 0, \dots, k$, são estatisticamente independentes

\rightarrow usada na análise de convergência do vetor de coeficientes

\rightarrow não é rigorosamente válida quando $\underline{x}(n)$ consiste de elementos de uma linha de atraso, como é o caso de um filtro FIR.

→ possibilita a obtenção de limitantes para o valor do passo de adaptação que evitem a divergência do algoritmo:

- $0 < \mu < 2/\lambda_{\max}$, em que λ_{\max} é o maior autovalor de R_x

- $0 < \mu < 2/\text{tr}(R_x)$

$\text{tr}(R_x) = M r_x(0)$, em que $r_x(0) = E\{x(n)x(n)\} = E\{x^2(n)\} = \sigma_x^2$ (média nula)

- erro em excesso e desajuste

- valor esperado de $E\{e^2(n)\}$ à medida que $n \rightarrow \infty$ (situações de convergência)

$$e(n) = d(n) - \underline{w}^T \underline{x}(n)$$

$$= d(n) - \{ \underline{w}_0 + \underline{\Delta w}(n) \}^T \underline{x}(n) = d(n) - \underline{w}_0^T \underline{x}(n) - \underline{\Delta w}(n)^T \underline{x}(n)$$

$$= e_0(n) - \underline{\Delta w}(n)^T \underline{x}(n)$$

Então; $e^2(n) = e_0^2(n) - 2e_0(n) \underline{\Delta w}(n)^T \underline{x}(n) + \underline{\Delta w}(n)^T \underline{x}(n) \underline{x}^T(n) \underline{\Delta w}(n)$

$\underbrace{\hspace{10em}}_{\text{ortogonais } E\{e_0(n)\underline{x}(n)\} = 0}$

Aplicando o operador de esperança, chega-se a:

$$E\{e^2(n)\} = J_{\text{MIN}} + E\{ \underline{\Delta w}(n)^T R_x \underline{\Delta w}(n) \}$$

↳ erro residual devido à natureza estocástica do algoritmo.

A partir de $J_{\text{EXC}} = E\{ \underline{\Delta w}(n)^T R_x \underline{\Delta w}(n) \}$, é mais comum avaliar o desempenho do algoritmo por meio de uma grandeza chamada desajuste (misadjustment):

$$M \triangleq \frac{J_{\text{EXC}}}{J_{\text{MIN}}} \approx \frac{1}{2} \mu M \sigma_x^2$$

• Velocidade de convergência (constante de tempo)

$$\tau \approx \frac{1}{2\mu \lambda_{AVG}}, \text{ em que } \lambda_{AVG} = \frac{1}{M} \sum_{i=1}^M \lambda_i = \text{m\u00e9dia dos autovalores de } R_x$$

Observe que:

* ao aumentarmos o tamanho do passo (μ), \u00e9 poss\u00edvel reduzir a constante de tempo τ , o que significa que o tempo de acomodac\u00e3o do algoritmo LMS foi reduzido. Ou seja, o aumento do passo proporciona um aumento na velocidade de converg\u00eancia do algoritmo.

* por outro lado, quanto maior o valor de μ , maior o desajuste, ou, equivalentemente, maior o erro em excesso, o que significa que perdemos qualidade ou precis\u00e3o na aproxima\u00e7\u00e3o dos coeficientes \u00f3timos.

Constata\u00e7\u00e3o: um passo vari\u00e1vel $\mu(n)$ seria ideal

\u2192 valor de μ elevado para erros altos - enquanto estamos longe da solu\u00e7\u00e3o \u00f3tima, permitimos valores mais elevados do passo para chegar mais r\u00e1pido ao alvo.

\u2192 valor de μ pequeno para erros baixos - quando j\u00e1 aconteceu a "converg\u00eancia" para a regi\u00e3o pr\u00f3xima \u00e0 solu\u00e7\u00e3o \u00f3tima, reduzimos μ para diminuir o desajuste.

→ no caso estacionário, início o algoritmo com um valor de passo e gradativamente o reduz, proporcionando maior rapidez no início e precisão no final.