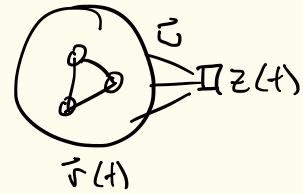


Derivation of recursive least squares

Consider an RNN with a 1D readout z :

$$z(t') = \sum_{i=1}^N w_i(t) r_i(t') \quad (1)$$



where t is the most recent timestep (so $t' \leq t$).

The goal is to choose weights $\bar{w}(t) = \hat{w}(t)$ that minimize the loss:

$$L(t) = \sum_{t'=1}^t \lambda^{t-t'} [\hat{z}(t') - z(t')]^2 + \alpha \lambda^t \sum_{i=1}^N [w_i(t)]^2 \quad (2)$$

[The second term is necessary so that the problem is well-posed (i.e. Φ^{-1} exists later) even if $t < N$.]

Plan:

- 1) Minimize $L(t)$ w.r.t. $\hat{w}(t)$ to get an equation for $\bar{w}(t)$.
- 2) Note that the resulting equation can be solved recursively, i.e.

$$\bar{w}(t) = \bar{w}(t-1) + \text{new stuff}$$

Ref: Adaptive Filter Theory, S. Haykin

②

First, minimize* (2):

$$0 = \frac{\partial L(t)}{\partial w_i}$$

$$= -2 \sum_{t'=1}^t \lambda^{t-t'} [\hat{z}(t') - z(t')] \frac{\partial z(t')}{\partial w_i} + 2\alpha \lambda^t w_i(t)$$

$$\stackrel{(1)}{=} -2 \sum_{t'=1}^t \lambda^{t-t'} [\hat{z}(t') - z(t')] r_i(t') + 2\alpha \lambda^t w_i(t)$$

$$\stackrel{(1)}{=} -2 \sum_{t'=1}^t \lambda^{t-t'} \hat{z}(t') r_i(t') + 2 \sum_{t'=1}^t \lambda^{t-t'} \sum_j w_j(t) r_j(t) + 2\alpha \lambda^t w_i(t)$$

$$\Rightarrow \sum_j [\alpha \lambda^t \delta_{ij} + \sum_{t'=1}^t \lambda^{t-t'} r_i(t') r_j(t')] \hat{w}_j(t)$$

$$= \sum_{t'=1}^t \lambda^{t-t'} \hat{z}(t') r_i(t')$$

$$\Rightarrow \sum_j \Phi_{ij}(t) \hat{w}_j(t) = u_i(t) \quad (3)$$

where we have defined

$$\Phi_{ij}(t) \equiv \alpha \lambda^t \delta_{ij} + \sum_{t'=1}^t \lambda^{t-t'} r_i(t) r_j(t) \quad (4)$$

$$u_i(t) \equiv \sum_{t'=1}^t \lambda^{t-t'} \hat{z}(t') r_i(t') \quad (5)$$

and used the Kronecker delta:

$$\delta_{ij} = \begin{cases} 1, & i=j \\ 0, & i \neq j \end{cases}$$

* Note this differs from grad. descent. Here we can set $\partial L / \partial w = 0$ since the problem is linear.

(3) Eq. (3) can be inverted to get $\hat{\tilde{w}}(t)$:

$$\hat{\tilde{w}}(t) = \Phi^{-1}(t) \tilde{u}(t) \quad (6)$$

But just inverting $\Phi(t)$ numerically at every t would be very expensive.

We can do better by noting that (4) and (5) have recursive forms:

$$\begin{aligned} \Phi_{ij}(t) &= \lambda \left[\alpha \lambda^{t-1} \delta_{ij} + \sum_{t'=1}^{t-1} \lambda^{t-t'-1} r_i(t') r_j(t') \right] \\ &\quad + r_i(t) r_j(t) \\ &= \lambda \Phi_{ij}(t-1) + r_i(t) r_j(t) \end{aligned} \quad (7)$$

$$u_i(t) = \lambda u_i(t-1) + \hat{s}_i(t) r_i(t) \quad (8)$$

What we really want, though, is a recursion relation for $\Phi^{-1}(t)$, not for $\Phi(t)$. It turns out that, since the last term in (7) has rank $1 \ll N$, this can be done easily. In general, if

$$A = B + \bar{F} \bar{F}^T \quad (9)$$

then

$$A^{-1} = B^{-1} - \frac{B^{-1} \bar{F} \bar{F}^T B^{-1}}{1 + \bar{F} B^{-1} \bar{F}} \quad (10)$$

④ Proof :

$$\begin{aligned}
 A \cdot A^{-1} &= (\mathbf{B} + \vec{r} \vec{r}^T) \cdot \left(\mathbf{B}^{-1} - \frac{\mathbf{B}^{-1} \cdot \vec{r} \vec{r}^T \cdot \mathbf{B}^{-1}}{1 + \vec{r} \cdot \mathbf{B}^{-1} \cdot \vec{r}} \right) \\
 &= \mathbb{I} + \vec{r} \vec{r}^T \cdot \mathbf{B}^{-1} \\
 &\quad - \frac{\mathbb{I} \cdot \vec{r} \vec{r}^T \cdot \mathbf{B}^{-1} - \vec{r} \vec{r}^T \cdot \mathbf{B}^{-1} \cdot \vec{r} \vec{r}^T \cdot \mathbf{B}^{-1}}{1 + \vec{r} \cdot \mathbf{B}^{-1} \cdot \vec{r}} \\
 &= \mathbb{I} + \vec{r} \vec{v}^T - \frac{\vec{r} \vec{v}^T - \vec{r} (\vec{v}^T \cdot \vec{r}) \vec{v}^T}{1 + \vec{r} \cdot \vec{v}} \\
 &= \mathbb{I} + \frac{\vec{r} \vec{v}^T - (\vec{r} \cdot \vec{v}) \vec{r} \vec{v}^T - \vec{r} \vec{v}^T + (\vec{r} \cdot \vec{v}) \vec{r} \vec{v}^T}{1 + \vec{r} \cdot \vec{v}} \\
 &= \mathbb{I}
 \end{aligned}$$

where \mathbb{I} is the identity matrix, and $\vec{v} = \vec{r}^T \cdot \mathbf{B}^{-1}$.

So, letting $A = \Phi(t)$ and $B = \lambda \Phi(t-1)$, (7) becomes

$$\Phi^{-1}(t) = \frac{1}{\lambda} \Phi^{-1}(t-1) - \frac{\Phi^{-1}(t-1) \cdot \vec{r}(t) \vec{r}^T(t) \cdot \Phi^{-1}(t-1)}{\lambda^2 + \lambda \vec{r}(t) \cdot \Phi^{-1}(t-1) \cdot \vec{r}(t)}$$

Or, letting $P(t) = \Phi^{-1}(t)$,

$$P(t) = \frac{1}{\lambda} \left[P(t-1) - \frac{P(t-1) \cdot \vec{r}(t) \vec{r}^T(t) \cdot P(t-1)}{\lambda + \vec{r}(t) \cdot P(t-1) \cdot \vec{r}(t)} \right] \quad (11)$$

⑤ For convenience, define the "gain vector":

$$\vec{k}(t) \equiv \frac{P(t-1) \cdot \vec{r}(t)}{\lambda + \vec{r}(t) \cdot P(t-1) \cdot \vec{r}(t)} \quad (12)$$

and note that

$$\begin{aligned} P(t) \cdot \vec{r}(t) &= \frac{1}{\lambda} \left[P(t-1) - \frac{P(t-1) \cdot \vec{r}(t) \vec{r}^\top(t) \cdot P(t-1)}{\lambda + \vec{r}(t) \cdot P(t-1) \cdot \vec{r}(t)} \right] \cdot \vec{r}(t) \\ &= \frac{1}{\lambda} P(t-1) \cdot \vec{r}(t) \left[1 - \frac{\vec{r}(t) \cdot P(t-1) \cdot \vec{r}(t)}{\lambda + \vec{r}(t) \cdot P(t-1) \cdot \vec{r}(t)} \right] \\ &= \frac{P(t-1) \cdot \vec{r}(t)}{\lambda + \vec{r}(t) \cdot P(t-1) \cdot \vec{r}(t)} \\ &= \vec{k}(t) \end{aligned} \quad (13)$$

With (8) and (11)-(13), the weight update equation-(6) is

$$\begin{aligned} \hat{\vec{w}}(t) &= P(t) \cdot \vec{u}(t) \\ &\stackrel{(8)}{=} P(t) \cdot \left[\lambda \vec{u}(t-1) + \hat{\chi}(t) \vec{r}(t) \right] \\ &\stackrel{(13)}{=} \lambda P(t) \cdot \vec{u}(t-1) + \hat{\chi}(t) \vec{k}(t) \\ &\stackrel{(11)}{=} \left[P(t-1) - \frac{P(t-1) \cdot \vec{r}(t) \vec{r}^\top(t) \cdot P(t-1)}{\lambda + \vec{r}(t) \cdot P(t-1) \cdot \vec{r}(t)} \right] \cdot \vec{u}(t-1) \\ &\quad + \hat{\chi}(t) \vec{k}(t) \end{aligned}$$

$$\begin{aligned}
 (6) \quad & \stackrel{(12)}{=} \vec{P}(t-1) \cdot \vec{u}(t-1) - \vec{k}(t) \vec{r}^T(t) \cdot \vec{P}(t-1) \cdot \vec{u}(t-1) \\
 & + \hat{z}(t) \vec{k}(t) \\
 & \stackrel{(6)}{=} \hat{\omega}(t-1) - \vec{k}(t) [\vec{r}(t) \cdot \hat{\omega}(t-1)] + \hat{z}(t) \vec{k}(t) \\
 & = \hat{\omega}(t-1) + \vec{k}(t) \{ \hat{z}(t) - \hat{\omega}(t-1) \cdot \vec{r}(t) \} \quad (14)
 \end{aligned}$$

Letting $\hat{\xi}(t) \equiv \hat{z}(t) - \hat{\omega}(t-1) \cdot \vec{r}(t)$, this is

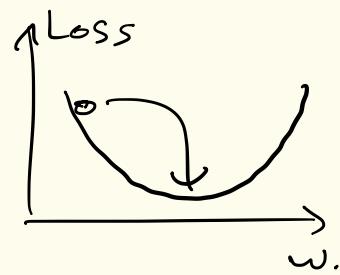
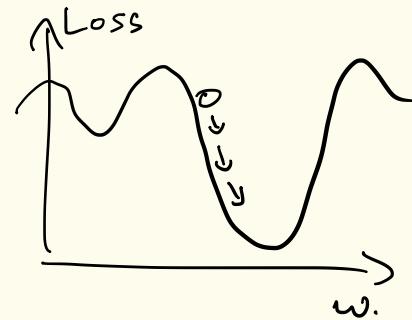
$$\boxed{
 \begin{aligned}
 \hat{\omega}(t) &= \hat{\omega}(t-1) + \xi(t) \vec{k}(t) \\
 &= \hat{\omega}(t-1) + \text{new stuff}
 \end{aligned} \quad (15)
 }$$

as promised.

⑦

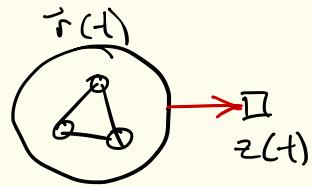
Discussion

- RLS differs from gradient-based methods (e.g. BPTT and RTRL).
 - Grad-based methods solve a nonlinear problem in many small steps.
 - RLS solves a linear problem all at once. The error is always small.
- The difficulty with RLS is how to get a basis set $\{r_i(t)\}$ for constructing the output $z(t)$ that is (i) sufficiently rich, and (ii) stable.

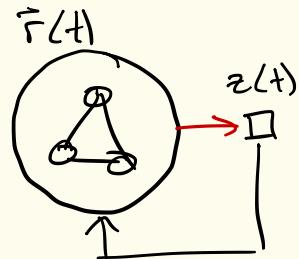


(8)

- In the echo state network (Jaeger and Haas, 2004), the RNN is driven with the target $\hat{z}(t)$, and recurrent weights are small.



- In FORCE learning (Sussillo and Abbott, 2009), recurrent weights are large, and the RNN is driven by its own output $z(t)$.
 - Using $z(t)$ rather than $\hat{z}(t)$ as feedback might make the trained network more robust since it introduces noise.
 - Large recurrent weights are good because they lead to a rich basis set $\{r_i(t)\}$. Chaos isn't a problem, though, since feedback suppresses it (Rajan et al, 2010).

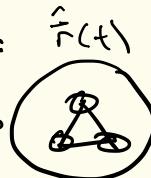


(9)

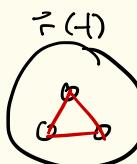
- In Full-FORCE, a teacher RNN provides targets $\hat{r}(t)$ for training recurrent weights in a student RNN (DePasquale et al., 2018).

Teacher: $\hat{r}(t)$

$$\hat{z}(t) \rightarrow$$



Student:



$$\rightarrow \square$$

Relation to Kalman filter

Rename variables:

- $\hat{w} \rightarrow \vec{x}(t)$, Target weights \rightarrow Latent var.
- $\vec{w}(t) \rightarrow \hat{\vec{x}}(t)$, Readout weights \rightarrow Est. latent var.
- $\vec{r}(t) \rightarrow \vec{c}(t)$, RNN activity \rightarrow Readout weights
- $\vec{z}(t) \rightarrow \vec{z}(t)$, RNN readout \rightarrow Readout
- $\varepsilon(t) \rightarrow \vec{v}(t)$, Readout error \rightarrow Readout noise

In this language, our problem is, given the update rules

$$\vec{x}(t+1) = \vec{x}(t)$$

$$\vec{z}(t) = \vec{c}(t) \cdot \vec{x}(t) + \vec{v}(t),$$

and assuming that $\vec{z}(t)$ and $\vec{c}(t)$ are known but $\vec{x}(t)$ is not, to find the

(10) best estimate $\hat{\vec{x}}(t)$ to reproduce $\vec{z}(t)$.

This is a special case of the Kalman filtering problem:

$$\vec{\hat{x}}(t+1) = F(t+1, t) \cdot \vec{\hat{x}}(t) + \vec{v}_1(t)$$

$$\vec{z}(t) = C(t) \cdot \vec{\hat{x}}(t) + \vec{v}_2(t)$$

where $\vec{v}_{1,2}(t)$ are white noise with known variance. The goal is then, given C, F , and $\{\vec{z}(1), \dots, \vec{z}(t)\}$, to find the best estimate $\hat{\vec{x}}(t')$ of $\vec{x}(t')$. This can have different names:

$$\begin{cases} t' = t \rightarrow \text{"filtering"} \\ t' < t \rightarrow \text{"smoothing"} \\ t' > t \rightarrow \text{"prediction"} \end{cases}$$

The solution looks very similar to the RLS algorithm defined earlier (see Haykin, ch. 10).