

# STAT 153 & 248 - Time Series

## Lecture Twenty Five

Spring 2025, UC Berkeley

Aditya Guntuboyina

April 29, 2025

### 1 Nonlinear AutoRegression

In the last lecture, we started discussing nonlinear forms of autoregression for an observed time series  $y_1, \dots, y_n$ . For each  $t$ , we take  $x_t = (y_{t-1}, \dots, y_{t-p})^T$  for some integer  $p \geq 1$ .  $x_t$  can be called the covariate at time  $t$  corresponding to the response value  $y_t$ . In the context of recurrent neural network models,  $x_t$  is referred to as the input at time  $t$ .

The usual (linear) autoregression AR( $p$ ) model corresponds to:

$$\mu_t = \beta_0 + \beta^T x_t. \quad (1)$$

The loss function is  $\sum_t (y_t - \mu_t)^2$ , and the parameters  $\beta_0, \beta$  are estimated by minimizing the loss.

In nonlinear autoregression, we change the formula (1) into a nonlinear function of  $x_t$ . When  $p = 1$ , one simple nonlinear AR(1) model is:

$$\mu_t = \beta_0 + \beta_1 x_t + \beta_2 (x_t - c_1)_+ + \dots + \beta_{k+1} (x_t - c_k)_+.$$

We simplify this slightly by dropping  $x_t$  (because  $x_t = (x_t - c_0)_+ + c_0$  for all  $t$  provided  $c_0$  is smaller than all the observed values of  $x_t$ ; we will not lose anything by dropping  $x_t$ ). This leads to

$$\mu_t = \beta_0 + \beta_1 (x_t - c_1)_+ + \dots + \beta_k (x_t - c_k)_+.$$

We rewrite this equation using the following notation:

$$\begin{aligned} s_t &= (x_t - c_1, \dots, x_t - c_k)^T \\ r_t &= \sigma(s_t) \\ \mu_t &= \beta_0 + \beta^T r_t. \end{aligned} \quad (2)$$

$s_t$  is a linear function of  $x_t$  which maps the scalar  $x_t$  to the  $k \times 1$  vector  $s_t$ .  $\sigma(\cdot)$  denotes the ReLU function applied pointwise to the input. So  $r_t$  is obtained by apply the ReLU function to each coordinate of  $s_t$ . Finally  $\mu_t$  is a linear function of  $r_t$  (we shall sometimes refer to  $\mu_t$  as the output corresponding to the input  $x_t$ ).

When  $p \geq 1$ , there are multiple ways of generalizing (2). One simple way is to consider the following “additive” model (below  $x_t^{(i)} = y_{t-i}$  denotes the  $i^{\text{th}}$  coordinate of  $x_t$ )

$$\begin{aligned}
s_t^{(i)} &= (x_t^{(i)} - c_1^{(i)}, \dots, x_t^{(i)} - c_k^{(i)})^T \quad \text{for } 1 \leq i \leq p \\
s_t &= \begin{pmatrix} s_t^{(1)} \\ \cdot \\ \cdot \\ \cdot \\ s_t^{(p)} \end{pmatrix} \\
r_t &= \sigma(s_t) \\
\mu_t &= \beta_0 + \beta^T r_t
\end{aligned} \tag{3}$$

This is called an additive model because  $\mu_t$  can be written as an additive sum of separate functions of  $x_t^{(i)}$  for  $i = 1, \dots, p$ . A different (i.e., non-additive) generalization of (2) is the single-hidden layer neural network defined as follows.

$$\begin{aligned}
s_t &= Wx_t + b \\
r_t &= \sigma(s_t) \\
\mu_t &= \beta_0 + \beta^T r_t
\end{aligned} \tag{4}$$

Here  $s_t$  is again  $k \times 1$ ,  $W$  is  $k \times p$  and  $b$  is  $p \times 1$ . We shall refer to (4) as the NonLinear AR model of order  $p$ . The total number of parameters here is  $kp + k + k + 1 = kp + 2k + 1$ . When  $p$  increases by 1, the number of parameters in (4) increases by  $k$ . On the other hand, in the usual (linear), AR( $p$ ) model, the number of parameters increases only by 1 when  $p$  increases by 1. So these models have a tendency to become high-dimensional faster than the linear AR( $p$ ) models.

## 2 Recurrent Neural Network (RNN)

RNN is given by

$$\begin{aligned}
r_0 &= 0 \\
s_t &= W_r r_{t-1} + Wx_t + b \\
r_t &= \sigma_{\tanh}(s_t) \\
\mu_t &= \beta_0 + \beta^T r_t
\end{aligned} \tag{5}$$

This formula can also be written as

$$\begin{aligned}
r_0 &= 0 \\
r_t &= \sigma_{\tanh}(W_r r_{t-1} + Wx_t + b) \\
\mu_t &= \beta_0 + \beta^T r_t
\end{aligned} \tag{6}$$

Here the activation function  $\sigma_{\tanh}$  is the tanh activation function given by

$$\sigma_{\tanh}(u) := \frac{e^u - e^{-u}}{e^u + e^{-u}}.$$

The parameters now are  $W_r$  ( $k \times k$  matrix),  $W$  ( $k \times p$  matrix),  $b$  ( $k \times 1$  vector),  $\beta_0$  (scalar) and  $\beta$  ( $k \times 1$  vector).

In (6),  $r_t$  depends on all of  $x_u, u \leq t$ . To see this, just note (below  $\sigma = \sigma_{\tanh}$ )

$$\begin{aligned} r_1 &= \sigma(Wx_1 + b) && \text{because } r_0 = 0 \\ r_2 &= \sigma(W_r\sigma(Wx_1 + b) + Wx_2 + b), \\ r_3 &= \sigma(W_r\sigma(W_r\sigma(Wx_1 + b) + Wx_2 + b) + Wx_3 + b), \\ r_4 &= \sigma(W_r\sigma(W_r\sigma(W_r\sigma(Wx_1 + b) + Wx_2 + b) + Wx_3 + b) + Wx_4 + b). \end{aligned} \quad (7)$$

From the above,  $r_t$  clearly depends on all of  $x_1, \dots, x_t$ . But the strength of the dependence of  $r_t$  on  $x_s$  varies with  $s$ . To see this, observe that

$$\frac{\partial r_t}{\partial x_u} = \sigma'(s_t)W_r\sigma'(s_{t-1})W_r \dots \sigma'(s_{u+1})W_r\sigma'(s_u)W \quad \text{for } u \leq t. \quad (8)$$

Here  $\frac{\partial r_t}{\partial x_u}$  denotes the  $k \times p$  Jacobian Matrix of derivatives of  $r_t$  with respect to  $x_u$ . On the right hand side in (8),  $\sigma'(s_t)$  should be interpreted as  $k \times k$  diagonal matrices whose diagonal entries are obtained by applying the  $\sigma'(u) = \frac{d}{du}\sigma(u)$  function to each element of  $s_t$  ( $\sigma'(s_{t-1}), \dots$  are similarly defined as  $k \times k$  diagonal matrices).

As a concrete example,

$$\begin{aligned} \frac{\partial r_4}{\partial x_4} &= \sigma'(s_4)W & \frac{\partial r_4}{\partial x_3} &= \sigma'(s_4)W_r\sigma'(s_3)W & \frac{\partial r_4}{\partial x_2} &= \sigma'(s_4)W_r\sigma'(s_3)W_r\sigma'(s_2)W \\ \frac{\partial r_4}{\partial x_1} &= \sigma'(s_4)W_r\sigma'(s_3)W_r\sigma'(s_2)W_r\sigma'(s_1)W \end{aligned}$$

Note that these gradient formulae are with respect to inputs  $x_u$ , and not with respect to the parameters (which is crucial to parameter estimation). In the formula (8), it is clear that when  $u$  is much smaller than  $t$ , many more terms appear in the right hand side of (8) compared to the case when  $u$  is closer to  $t$ . Note here that  $\sigma$  is the tanh activation function:

$$\sigma(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}} \text{ so that } \sigma'(u) = 1 - \sigma^2(u) \in (0, 1].$$

Thus each  $\sigma'(\cdot)$  term will add a fractional multiplier to  $\partial r_t/\partial x_u$ . The number of these fractional multipliers will increase as  $u$  decreases in (8).

Further the matrix  $W_r$  also plays a key role in (8). For the model equation in (5) to be stable,  $W_r$  needs to have spectral radius (defined as the largest modulus of any eigenvalue) to be strictly smaller than one. In that case, each additional  $W_r$  multiplier will bring the whole term down, leading to  $\partial r_t/\partial x_u$  being small when  $u$  is much smaller than  $t$ .

This points to the following shortcoming of RNNs that more sophisticated models such as GRUs and LSTMs attempt to fix. We want  $r_t$  to represent the ideal summary of  $x_1, \dots, x_t$  that is relevant for the output  $y_t$ . However, in an RNN,  $r_t$  effectively only depends on those inputs  $x_u$  which are somewhat close to  $t$ . In this sense, the RNN can be thought of as not having a very long memory.

This ‘‘lack of long memory’’ problem with RNNs can be fixed by use of GRUs and LSTMs.

### 3 GRU (Gated Recurrent Unit)

Consider again the RNN formula (6). The basic problem with this is that  $r_t$  depends on  $r_{t-1}$  through the term  $W_r r_{t-1}$ . If  $W_r$  is a matrix with spectral radius less than 1 (which it needs

to be for stability purposes), then the multiplier  $W_r r_{t-1}$  can be thought of as “reducing”  $r_{t-1}$  by a factor of  $W_r$ . If this formula is applied repeatedly, then very soon the dependence of  $r_t$  on  $r_u$  will be very small. In order to avoid this, one needs to prevent  $r_t$  from depending on  $r_{t-1}$  only through  $W_r r_{t-1}$ .

This leads to the following idea. First construct a potential version  $\tilde{r}_t$  of  $r_t$  in the same way as (6):

$$\tilde{r}_t = \sigma(W_r r_{t-1} + W x_t + b). \quad (9)$$

This  $\tilde{r}_t$  only depends on  $r_{t-1}$  through  $W_r r_{t-1}$ . The two natural options for  $r_t$  now are:

1.  $r_t = \tilde{r}_t$ : in this case, we are back to the RNN (6).
2.  $r_t = r_{t-1}$ : in this case,  $r_t$  is exactly equal to  $r_{t-1}$ , which means that the current input  $x_t$  is ignored.

The idea behind GRU is to take a “convex-like” combination of these two options in the following way:

$$r_t = z_t r_{t-1} + (1 - z_t) \tilde{r}_t.$$

This would be exactly a convex combination if  $z_t$  were a scalar in the interval  $[0, 1]$ . But we allow  $z_t$  to be a vector interpreting the multiplication as pointwise. Thus it is better to write

$$r_t = z_t \odot r_{t-1} + (1 - z_t) \odot \tilde{r}_t. \quad (10)$$

The next thing to specify  $z_t$ . In GRU, we take

$$z_t = \sigma_{\text{sigmoid}}(W_r^z r_{t-1} + W^z x_t + b^z) \quad \text{where } \sigma_{\text{sigmoid}}(u) := \frac{1}{1 + e^{-u}}. \quad (11)$$

Because  $\sigma_{\text{sigmoid}}$  takes values between 0 and 1, the above formula ensures that the components of  $z_t$  take values in  $[0, 1]$  so that (10) represents a convex combination at the level of each individual component. Further (11) implies that  $z_t$  is also determined by  $r_{t-1}$  and  $x_t$ . The parameters  $W_r^z, W^z, b$  controlling the formula (11) are also unknown and they will be estimated along with all the other parameters of the model.

$z_t$  is sometimes referred to as a gate. It controls the closeness of  $r_t$  to  $r_{t-1}$  and  $\tilde{r}_t$ .

$r_{t-1}$  appears in two places in the formula (10): in the term  $z_t \odot r_{t-1}$  as well as in the formula (9) for  $\tilde{r}_t$ . It might be redundant to have  $r_{t-1}$  appear in both these places. To address this, GRU modifies (9) by using one more gate as follows:

$$\tilde{r}_t := \sigma(W_r(r_{t-1} \odot g_t) + W x_t + b),$$

where  $g_t$  controls the extent to which  $r_{t-1}$  is used in the formula for  $\tilde{r}_t$ . Similar to (11), the gate  $g_t$  is specified via

$$g_t = \sigma_{\text{sigmoid}}(W_r^g r_{t-1} + W^g x_t + b^g). \quad (12)$$

Putting all the formulae together, we get the following specification of the GRU model:

$$\begin{aligned} r_0 &= 0 \\ g_t &= \sigma_{\text{sigmoid}}(W_r^g r_{t-1} + W^g x_t + b^g) \\ z_t &= \sigma_{\text{sigmoid}}(W_r^z r_{t-1} + W^z x_t + b^z) \\ \tilde{r}_t &:= \sigma_{\text{tanh}}(W_r(r_{t-1} \odot g_t) + W x_t + b) \\ r_t &= z_t \odot r_{t-1} + (1 - z_t) \odot \tilde{r}_t \\ \mu_t &= \beta_0 + \beta^T r_t. \end{aligned} \quad (13)$$

$z_t$  is called the update gate while  $g_t$  is called the reset gate. The unknown parameters in this model (which need to be estimated from the data) are  $W_r^g, W^g, b^g, W_r^z, W^z, b^z, W_r, W, b, \beta_0, \beta$ .

This is a more sophisticated model compared to the RNN model (6). In fact, (6) is a special case of (13) corresponding to  $g_t = 1$  and  $z_t = 0$ . The presence of the gates  $g_t$  and  $z_t$  can alleviate the lack of long memory problem that was an issue with the RNNs.

## 4 LSTM (Long Short Term Memory)

LSTM is another modification to the basic RNN for enabling long memory. It also uses gates and has one more gate compared to the GRU. Instead of a recursion directly between  $r_{t-1}$  and  $r_t$ , the LSTM recursions are between the pairs  $(s_{t-1}, r_{t-1}) \rightarrow (s_t, r_t)$ .

We again construct a potential version  $\tilde{r}_t$  of  $r_t$  in the same way as RNN:

$$\tilde{r}_t = \sigma(W_r r_{t-1} + W x_t + b). \quad (14)$$

In GRU,  $r_t$  was defined as a convex combination of  $\tilde{r}_t$  and  $r_{t-1}$ . In LSTM,  $s_t$  is taken to be a linear combination of  $s_{t-1}$  and  $\tilde{r}_t$  with gates controlling both coefficients of the linear combination:

$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{r}_t,$$

where  $f_t$  and  $i_t$  denote gates.  $r_t$  is defined usually as  $\sigma_{\tanh}(s_t)$ . In LSTM, one also adds a gate to  $r_t$ :

$$r_t = o_t \odot \sigma_{\tanh}(s_t).$$

Putting all the terms together (and also writing the formulae for the gates), we obtain the full LSTM model:

$$\begin{aligned} r_0 &= 0 \\ f_t &= \sigma_{\text{sigmoid}}(W_r^f r_{t-1} + W^f x_t + b^f) \\ i_t &= \sigma_{\text{sigmoid}}(W_r^i r_{t-1} + W^i x_t + b^i) \\ o_t &= \sigma_{\text{sigmoid}}(W_r^o r_{t-1} + W^o x_t + b^o) \\ \tilde{r}_t &:= \sigma_{\tanh}(W_r r_{t-1} + W x_t + b) \\ s_t &= f_t \odot s_{t-1} + i_t \odot \tilde{r}_t \\ r_t &= o_t \odot \sigma_{\tanh}(s_t) \\ \mu_t &= \beta_0 + \beta^T r_t \end{aligned} \quad (15)$$

$f_t$  is called the forget gate,  $i_t$  is called the input gate and  $o_t$  is called the output gate. The unknown parameters in this model are  $W_r^f, W^f, b^f, W_r^i, W^i, b^i, W_r^o, W^o, b^o, W_r, W, b, \beta_0, \beta$ . These need to be estimated from data.

## 5 Additional Optional Reading

1. For more on GRUs, read [https://en.wikipedia.org/wiki/Gated\\_recurrent\\_unit](https://en.wikipedia.org/wiki/Gated_recurrent_unit).
2. For more on LSTMs, read [https://en.wikipedia.org/wiki/Long\\_short-term\\_memory](https://en.wikipedia.org/wiki/Long_short-term_memory).
3. A clear description of LSTM is given here: <https://www.youtube.com/watch?v=YCzL96nL7j0&t=870s>