# Newton Methods for Neural Networks: Gauss Newton Matrix-vector Product

Chih-Jen Lin

National Taiwan University

Last updated: June 1, 2020

# Outline

# Outline

# Outline

1. **Backward setting**
   - Jacobian evaluation
   - Gauss-Newton Matrix-vector products

2. Forward + backward settings
   - R operator
   - Gauss-Newton matrix-vector product

# Jacobian Evaluation: Convolutional Layer I

- For an instance $i$ the Jacobian can be partitioned into $L$ blocks according to layers

$$J^i = \begin{bmatrix} J^{1,i} & J^{2,i} & \dots & J^{L,i} \end{bmatrix}, \quad m = 1, \dots, L, \qquad (1)$$

where

$$J^{m,i} = \begin{bmatrix} \dfrac{\partial \mathbf{z}^{L+1,i}}{\partial \text{vec}(W^m)^T} & \dfrac{\partial \mathbf{z}^{L+1,i}}{\partial (\mathbf{b}^m)^T} \end{bmatrix}.$$

- The calculation seems to be very similar to that for the gradient.

# Jacobian Evaluation: Convolutional Layer II

- For the convolutional layers, recall for gradient we have

$$\frac{\partial f}{\partial W^m} = \frac{1}{C}W^m + \frac{1}{l}\sum_{i=1}^{l}\frac{\partial \xi_i}{\partial W^m}$$

and

$$\frac{\partial \xi_i}{\partial \text{vec}(W^m)^T} = \text{vec}\left(\frac{\partial \xi_i}{\partial S^{m,i}}\phi(\text{pad}(Z^{m,i}))^T\right)^T$$

# Jacobian Evaluation: Convolutional Layer III

- Now we have

$$\frac{\partial \boldsymbol{z}^{L+1,i}}{\partial \mathrm{vec}(W^m)^T} = \begin{bmatrix} \frac{\partial z_1^{L+1,i}}{\partial \mathrm{vec}(W^m)^T} \\ \vdots \\ \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial \mathrm{vec}(W^m)^T} \end{bmatrix}$$

$$= \begin{bmatrix} \mathrm{vec}(\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}} \phi(\mathrm{pad}(Z^{m,i}))^T)^T \\ \vdots \\ \mathrm{vec}(\frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}} \phi(\mathrm{pad}(Z^{m,i}))^T)^T \end{bmatrix}$$

# Jacobian Evaluation: Convolutional Layer IV

- If $\boldsymbol{b}^m$ is considered, the result is

$$
\left[ \frac{\partial \boldsymbol{z}^{L+1,i}}{\partial \text{vec}(W^m)^T} \quad \frac{\partial \boldsymbol{z}^{L+1,i}}{\partial (\boldsymbol{b}^m)^T} \right]
$$

$$
= \begin{bmatrix} \text{vec} \left( \frac{\partial z_1^{L+1,i}}{\partial S^{m,i}} \left[ \phi(\text{pad}(Z^{m,i}))^T \ \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m} \right] \right)^T \\ \vdots \\ \text{vec} \left( \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}} \left[ \phi(\text{pad}(Z^{m,i}))^T \ \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m} \right] \right)^T \end{bmatrix}.
$$

# Jacobian Evaluation: Convolutional Layer V

- We can see that it's more complicated than gradient.
- Gradient is a vector but Jacobian is a matrix

# Jacobian Evaluation: Backward Process I

- For gradient, earlier we need a backward process to calculate

$$\frac{\partial \xi_i}{\partial S^{m,i}}$$

- Now what we need are

$$\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}, \ldots, \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}$$

- The process is similar

# Jacobian Evaluation: Backward Process II

- If with RELU activation function and max pooling, for gradient we had

$$\frac{\partial \xi_i}{\partial \mathsf{vec}(S^{m,i})^T}$$

$$= \left( \frac{\partial \xi_i}{\partial \mathsf{vec}(Z^{m+1,i})^T} \odot \mathsf{vec}(I[Z^{m+1,i}])^T \right) P_{\mathsf{pool}}^{m,i}.$$

# Jacobian Evaluation: Backward Process III

- Assume that
$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m+1,i})}$$

are available.

$$\frac{\partial z_j^{L+1,i}}{\partial \text{vec}(S^{m,i})^T}$$

$$= \left( \frac{\partial z_j^{L+1,i}}{\partial \text{vec}(Z^{m+1,i})^T} \odot \text{vec}(I[Z^{m+1,i}])^T \right) P_{\text{pool}}^{m,i},$$

$$j = 1, \ldots, n_{L+1}.$$

# Jacobian Evaluation: Backward Process IV

- These row vectors can be written together as a matrix

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(S^{m,i})^T}$$

$$= \left( \frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m+1,i})^T} \odot \left( \mathbb{1}_{n_{L+1}} \text{vec}(I[Z^{m+1,i}])^T \right) \right) P_{\text{pool}}^{m,i}.$$

# Jacobian Evaluation: Backward Process V

- For gradient, we use

$$\frac{\partial \xi_i}{\partial S^{m,i}}$$

to have

$$\frac{\partial \xi_i}{\partial \text{vec}(Z^{m,i})^T} = \text{vec} \left( (W^m)^T \frac{\partial \xi_i}{\partial S^{m,i}} \right)^T P^m_\phi P^m_{\text{pad}}$$

and pass it to the previous layer

# Jacobian Evaluation: Backward Process VI

- Now we need to generate

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m,i})^T}$$

and pass it to the previous layer.

- Now we have

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m,i})^T} = \begin{bmatrix} \text{vec}\left((W^m)^T \frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}\right)^T P_\phi^m P_{\text{pad}}^m \\ \vdots \\ \text{vec}\left((W^m)^T \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}\right)^T P_\phi^m P_{\text{pad}}^m \end{bmatrix}.$$

# Jacobian Evaluation: Fully-connected Layer I

- We do not discuss details, but list all results below

$$\frac{\partial \boldsymbol{z}^{L+1,i}}{\partial \mathrm{vec}(W^m)^T} =$$

$$\left[ \mathrm{vec}\left( \frac{\partial z_1^{L+1,i}}{\partial \boldsymbol{s}^{m,i}} (\boldsymbol{z}^{m,i})^T \right) \quad \ldots \quad \mathrm{vec}\left( \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial \boldsymbol{s}^{m,i}} (\boldsymbol{z}^{m,i})^T \right) \right]^T$$

# Jacobian Evaluation: Fully-connected Layer II

$$\frac{\partial z^{L+1,i}}{\partial (b^m)^T} = \frac{\partial z^{L+1,i}}{\partial (s^{m,i})^T},$$

$$\frac{\partial z^{L+1,i}}{\partial (s^{m,i})^T} = \frac{\partial z^{L+1,i}}{\partial (z^{m+1,i})^T} \odot \left( \mathbb{1}_{n_{L+1}} I[z^{m+1,i}]^T \right)$$

$$\frac{\partial z^{L+1,i}}{\partial (z^{m,i})^T} = \frac{\partial z^{L+1,i}}{\partial (s^{m,i})^T} W^m$$

# Jacobian Evaluation: Fully-connected Layer III

- For layer $L + 1$, if using the squared loss and the linear activation function, we have

$$\frac{\partial \boldsymbol{z}^{L+1,i}}{\partial (\boldsymbol{s}^{L,i})^T} = \mathcal{I}_{n_{L+1}}.$$

# Gradient versus Jacobian I

- Operations for gradient

$$\frac{\partial \xi_i}{\partial \mathrm{vec}(S^{m,i})^T}$$

$$= \left( \frac{\partial \xi_i}{\partial \mathrm{vec}(Z^{m+1,i})^T} \odot \mathrm{vec}(I[Z^{m+1,i}])^T \right) P_{\mathrm{pool}}^{m,i}.$$

$$\frac{\partial \xi_i}{\partial W^m} = \frac{\partial \xi_i}{\partial S^{m,i}} \phi(\mathrm{pad}(Z^{m,i}))^T$$

$$\frac{\partial \xi_i}{\partial \mathrm{vec}(Z^{m,i})^T} = \mathrm{vec}\left( (W^m)^T \frac{\partial \xi_i}{\partial S^{m,i}} \right)^T P_\phi^m P_{\mathrm{pad}}^m,$$

# Gradient versus Jacobian II

- For Jacobian we have

$$
\frac{\partial z^{L+1,i}}{\partial \text{vec}(S^{m,i})^T}
$$

$$
= \left( \frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m+1,i})^T} \odot \left( \mathbb{1}_{n_{L+1}} \text{vec}(I[Z^{m+1,i}])^T \right) \right) P_{\text{pool}}^{m,i}.
$$

$$
\frac{\partial z^{L+1,i}}{\partial \text{vec}(W^m)^T} = \begin{bmatrix} \text{vec}(\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}} \phi(\text{pad}(Z^{m,i}))^T)^T \\ \vdots \\ \text{vec}(\frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}} \phi(\text{pad}(Z^{m,i}))^T)^T \end{bmatrix}
$$

# Gradient versus Jacobian III

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m,i})^T}$$

$$= \begin{bmatrix} \text{vec}\left((W^m)^T \frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}\right)^T P_\phi^m P_{\text{pad}}^m \\ \vdots \\ \text{vec}\left((W^m)^T \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}\right)^T P_\phi^m P_{\text{pad}}^m \end{bmatrix}.$$

# Implementation I

- For gradient we did

$$\Delta \leftarrow \mathrm{mat}(\mathrm{vec}(\Delta)^T P_{\mathrm{pool}}^{m,i})$$

$$\frac{\partial \xi_i}{\partial W^m} = \Delta \cdot \phi(\mathrm{pad}(Z^{m,i}))^T$$

$$\Delta \leftarrow \mathrm{vec}\left((W^m)^T \Delta\right)^T P_\phi^m P_{\mathrm{pad}}^m$$

$$\Delta \leftarrow \Delta \odot I[Z^{m,i}]$$

- Now for Jacobian we have similar settings but there are some differences

# Implementation II

- We don't really store the Jacobian:

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(W^m)^T} = \begin{bmatrix} \text{vec}(\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}} \phi(\text{pad}(Z^{m,i}))^T)^T \\ \vdots \\ \text{vec}(\frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}} \phi(\text{pad}(Z^{m,i}))^T)^T \end{bmatrix}$$

- Recall Jacobian is used for matrix-vector products

$$G^S v = \frac{1}{C} v + \frac{1}{|S|} \sum_{i \in S} \left( (J^i)^T \left( B^i (J^i v) \right) \right) \qquad (2)$$

# Implementation III

- The form

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(W^m)^T} = \begin{bmatrix} \text{vec}(\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}\phi(\text{pad}(Z^{m,i}))^T)^T \\ \vdots \\ \text{vec}(\frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}\phi(\text{pad}(Z^{m,i}))^T)^T \end{bmatrix}$$

is like the product of two things

# Implementation IV

- If we have

$$\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}, \ldots, \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}, \text{ and } \phi(\text{pad}(Z^{m,i}))$$

  probably we can do the matrix-vector product without multiplying these two things out

- We will talk about this again later

- Thus our Jacobian evaluation is solely on obtaining

$$\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}, \ldots, \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}$$

# Implementation V

- Further we need to take all data (or data in the selected subset) into account
- In the end what we have is the following procedure
- In the beginning

$$\Delta \in R^{d^{m+1}a^{m+1}b^{m+1} \times n_{L+1} \times l}$$

This corresponds to

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m+1,i})^T} \odot \left( \mathbb{1}_{n_{L+1}} \text{vec}(I[Z^{m+1,i}])^T \right), \forall i = 1, \dots, l$$

# Implementation VI

- We then calculate

$$\Delta \leftarrow \mathsf{mat} \left( \begin{bmatrix} (P_{\mathsf{pool}}^{m,1})^T \mathsf{vec}(\Delta_{:,:,1}) \\ \vdots \\ (P_{\mathsf{pool}}^{m,l})^T \mathsf{vec}(\Delta_{:,:,l}) \end{bmatrix} \right)_{d^{m+1} \times a_{\mathsf{conv}}^m b_{\mathsf{conv}}^m n_{L+1} l}$$

- Recall that the pooling matrices are different across instances

# Implementation VII

- The above operation corresponds to

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(S^{m,i})^T}$$

$$= \left( \frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m+1,i})^T} \odot \left( \mathbb{1}_{n_{L+1}} \text{vec}(I[Z^{m+1,i}])^T \right) \right) P_{\text{pool}}^{m,i}.$$

- Now we get

$$\left[ \frac{\partial z_1^{L+1,1}}{\partial S^{m,1}} \quad \cdots \quad \frac{\partial z_{n_{L+1}}^{L+1,1}}{\partial S^{m,1}} \quad \cdots \quad \frac{\partial z_{n_{L+1}}^{L+1,l}}{\partial S^{m,l}} \right]$$

$$\in R^{d^{m+1} \times a_{\text{conv}}^m b_{\text{conv}}^m n_{L+1} l}$$

# Implementation VIII

- Next

$$V \leftarrow \text{vec}((W^m)^T \Delta) \in R^{hhd^m a^m_{\text{conv}} b^m_{\text{conv}} n_{L+1} l \times 1}$$

- This is same as

$$\text{vec}\left((W^m)^T \left[ \frac{\partial z_1^{L+1,1}}{\partial S^{m,1}} \quad \cdots \quad \frac{\partial z_{n_{L+1}}^{L+1,1}}{\partial S^{m,1}} \quad \cdots \quad \frac{\partial z_{n_{L+1}}^{L+1,l}}{\partial S^{m,l}} \right]\right).$$

# Implementation IX

- Now $V$ is a big vector like

$$\begin{bmatrix} \boldsymbol{v}_1^1 \\ \vdots \\ \boldsymbol{v}_{n_{L+1}}^1 \\ \vdots \\ \boldsymbol{v}_{n_{L+1}}^l \end{bmatrix}$$

Note that "$\boldsymbol{v}$" here is not the vector in matrix-vector products. We happen to use the same symbol

# Implementation X

- We then calculate

$$
\Delta \leftarrow \mathrm{mat} \left( \begin{bmatrix} \begin{bmatrix} (\mathbf{v}_1^1)^T P_\phi^m P_{\mathrm{pad}}^m \\ \vdots \\ (\mathbf{v}_{n_{L+1}}^1)^T P_\phi^m P_{\mathrm{pad}}^m \\ \vdots \\ (\mathbf{v}_{n_{L+1}}^l)^T P_\phi^m P_{\mathrm{pad}}^m \end{bmatrix} \end{bmatrix} \right)_{d^m a^m b^m \times n_{L+1} \times l}
$$

This corresponds to

$$
\frac{\partial \mathbf{z}^{L+1,i}}{\partial \mathrm{vec}(Z^{m,i})^T}, i = 1, \ldots, l
$$

# Implementation XI

- Finally,

$$\Delta \leftarrow \Delta \odot$$

$$\left[ \underbrace{I[Z^{m,1}] \cdots I[Z^{m,1}]}_{n_{L+1}} \cdots \underbrace{I[Z^{m,l}] \cdots I[Z^{m,l}]}_{n_{L+1}} \right] \quad (3)$$

This means

$$\frac{\partial z^{L+1,i}}{\partial \text{vec}(Z^{m,i})^T} \odot \left( \mathbb{1}_{n_{L+1}} \text{vec}(I[Z^{m,i}])^T \right), \forall i = 1, \dots, l$$

- Let's check the code

# Implementation XII

```
dzdS{m} = vTP(model, net, m, num_data,
          dzdS{m}, 'pool_Jacobian');

dzdS{m} = reshape(dzdS{m},
          model.ch_input(m+1), []);

V = model.weight{m}' * dzdS{m};
dzdS{m-1} = vTP(model, net, m, num_data,
            V, 'phi_Jacobian');

% vTP_pad
```

# Implementation XIII

```
dzdS{m-1} = reshape(dzdS{m-1},
  model.ch_input(m), model.ht_pad(m),
  model.wd_pad(m), []);
p = model.wd_pad_added(m);
dzdS{m-1} = dzdS{m-1}(:, p+1:p+model.ht_inpu
  p+1:p+model.wd_input(m), :);

dzdS{m-1} =
  reshape(dzdS{m-1}, [], nL, num_data)
  .* reshape(net.Z{m} > 0, [], 1, num_data);
```

# Implementation XIV

- In the last line for doing (3), we don't need to repeat each $I[Z^{m,i}]$ $n_{L+1}$ times. For .*, MATLAB does the expansion automatically

# Discussion I

- For doing several CG steps, we should store

$$\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}, \ldots, \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}$$

The memory cost is

$$l \times n_{L+1} \times \left( \sum_{m=1}^{L^c} d^{m+1} a_{\text{conv}}^m b_{\text{conv}}^m + \sum_{m=L^c+1}^{L} n_{m+1} \right) \quad (4)$$

- It is proportional to
  - Number of classes

# Discussion II

- Number of data for the subsampled Hessian
- The reason is that it's used for all CG steps (Jacobian matrix remains the same)
- Recalculating them at each CG is too expensive
- We will show some complexity analysis later
- Thus subsequently we will consider a different approach to reduce the memory consumption

# Outline

# Gauss-Newton Matrix-Vector Products I

- We check

$$Gv$$

  though the situation of using $G^S$ (i.e., a subset of data) is the same

- The Gauss-Newton matrix

$$G = \frac{1}{C}\mathcal{I} + \frac{1}{l}\sum_{i=1}^{l}\begin{bmatrix}(J^{1,i})^T\\\vdots\\(J^{L,i})^T\end{bmatrix}B^i\begin{bmatrix}J^{1,i} & \dots & J^{L,i}\end{bmatrix}$$

# Gauss-Newton Matrix-Vector Products II

- The Gauss-Newton matrix vector product

$$G\mathbf{v}$$

$$=\frac{1}{C}\mathbf{v} + \frac{1}{l}\sum_{i=1}^{l}\begin{bmatrix}(J^{1,i})^T\\\vdots\\(J^{L,i})^T\end{bmatrix} B^i \begin{bmatrix}J^{1,i} & \dots & J^{L,i}\end{bmatrix}\begin{bmatrix}\mathbf{v}^1\\\vdots\\\mathbf{v}^L\end{bmatrix}$$

$$=\frac{1}{C}\mathbf{v} + \frac{1}{l}\sum_{i=1}^{l}\begin{bmatrix}(J^{1,i})^T\\\vdots\\(J^{L,i})^T\end{bmatrix}\left(B^i \sum_{m=1}^{L} J^{m,i}\mathbf{v}^m\right),$$

$$(5)$$

# Gauss-Newton Matrix-Vector Products III

where
$$\boldsymbol{v} = \begin{bmatrix} \boldsymbol{v}^1 \\ \vdots \\ \boldsymbol{v}^L \end{bmatrix}$$

- Each $\boldsymbol{v}^m, m = 1, \ldots, L$ has the same length as the number of variables (including bias) at the $m$th layer.

# Gauss-Newton Matrix-Vector Products IV

- For the convolutional layers,

$$
\begin{aligned}
& J^{m,i} v^m \\
& = \begin{bmatrix}
\mathrm{vec}\left( \frac{\partial z_1^{L+1,i}}{\partial S^{m,i}} \left[ \phi(\mathrm{pad}(Z^{m,i}))^T \ \mathbb{1}_{a_{\mathrm{conv}}^m b_{\mathrm{conv}}^m} \right] \right)^T v^m \\
\vdots \\
\mathrm{vec}\left( \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}} \left[ \phi(\mathrm{pad}(Z^{m,i}))^T \ \mathbb{1}_{a_{\mathrm{conv}}^m b_{\mathrm{conv}}^m} \right] \right)^T v^m
\end{bmatrix} \\
& \in R^{n_{L+1} \times 1}
\end{aligned}
$$

- This formulation is fine, but we need

# Gauss-Newton Matrix-Vector Products V

- - a for loop to generate $n_{L+1}$ vectors
  - the product between a matrix and a vector $v^m$
- Is there a way to avoid a for loop?
- For a language like MATLAB/Octave, we hope to avoid for loops
- Also we hope the code can be simpler and shorter
- We use the following property

$$\text{vec}(AB)^T \text{vec}(C) = \text{vec}(A)^T \text{vec}(CB^T)$$

# Gauss-Newton Matrix-Vector Products VI

- The first element is

$$
\text{vec} \left( \underbrace{\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}}_{A} \underbrace{\left[ \phi(\text{pad}(Z^{m,i}))^T \ \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m} \right]}_{B} \right)^T \underbrace{\boldsymbol{v}^m}_{\text{vec}(C)}
$$

$$
= \frac{\partial z_1^{L+1,i}}{\partial \text{vec}(S^{m,i})^T} \times
$$

$$
\text{vec} \left( \text{mat}(\boldsymbol{v}^m)_{d^{m+1} \times (h^m h^m d^m + 1)} \begin{bmatrix} \phi(\text{pad}(Z^{m,i})) \\ \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}^T \end{bmatrix} \right).
$$

# Gauss-Newton Matrix-Vector Products VII

- If all elements are considered together

$$J^{m,i} v^m$$

$$= \frac{\partial z^{L+1,i}}{\partial \text{vec}(S^{m,i})^T} \times$$

$$\text{vec} \left( \text{mat}(v^m)_{d^{m+1} \times (h^m h^m d^m + 1)} \begin{bmatrix} \phi(\text{pad}(Z^{m,i})) \\ \mathbb{1}^T_{a^m_{\text{conv}} b^m_{\text{conv}}} \end{bmatrix} \right).$$

$$(6)$$

This involves

- One matrix-matrix product

# Gauss-Newton Matrix-Vector Products VIII

- One matrix-vector product
- After deriving (6), from (5), we sum results of all layers

$$\sum_{m=1}^{L} J^{m,i} \boldsymbol{v}^m$$

- Next we calculate

$$\boldsymbol{q}^i = B^i \left( \sum_{m=1}^{L} J^{m,i} \boldsymbol{v}^m \right). \tag{7}$$

# Gauss-Newton Matrix-Vector Products IX

- This is usually easy
- We mentioned earlier that if the squared loss is used

$$B^i = \begin{bmatrix} 2 & & \\ & \vdots & \\ & & 2 \end{bmatrix}$$

is a diagonal matrix

# Gauss-Newton Matrix-Vector Products X

- Finally, we calculate

$$
\begin{aligned}
& (J^{m,i})^T \boldsymbol{q}^i \\
&= \Bigg[ \mathrm{vec} \left( \frac{\partial z_1^{L+1,i}}{\partial S^{m,i}} \left[ \phi(\mathrm{pad}(Z^{m,i}))^T \; \mathbb{1}_{a_{\mathrm{conv}}^m b_{\mathrm{conv}}^m} \right] \right) \cdots \\
& \quad \mathrm{vec} \left( \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}} \left[ \phi(\mathrm{pad}(Z^{m,i}))^T \; \mathbb{1}_{a_{\mathrm{conv}}^m b_{\mathrm{conv}}^m} \right] \right) \Bigg] \boldsymbol{q}^i
\end{aligned}
$$

# Gauss-Newton Matrix-Vector Products XI

$$= \sum_{j=1}^{n_{L+1}} q_j^i \text{vec}\left(\frac{\partial z_j^{L+1,i}}{\partial S^{m,i}} \left[\phi(\text{pad}(Z^{m,i}))^T \; \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}\right]\right)$$

$$= \text{vec}\left(\sum_{j=1}^{n_{L+1}} q_j^i \left(\frac{\partial z_j^{L+1,i}}{\partial S^{m,i}} \left[\phi(\text{pad}(Z^{m,i}))^T \; \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}\right]\right)\right)$$

$$= \text{vec}\left(\left(\sum_{j=1}^{n_{L+1}} q_j^i \frac{\partial z_j^{L+1,i}}{\partial S^{m,i}}\right) \left[\phi(\text{pad}(Z^{m,i}))^T \; \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}\right]\right)$$

# Gauss-Newton Matrix-Vector Products XII

$$= \text{vec}\Bigg( \text{mat}\Bigg( \bigg( \frac{\partial \mathbf{z}^{L+1,i}}{\partial \text{vec}(S^{m,i})^T} \bigg)^T \mathbf{q}^i \Bigg)_{d^{m+1} \times a_{\text{conv}}^m b_{\text{conv}}^m} \times$$

$$\Big[ \phi(\text{pad}(Z^{m,i}))^T \ \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m} \Big] \Bigg). \tag{8}$$

A matrix-vector product and then a matrix-matrix product

# Gauss-Newton Matrix-Vector Products XIII

- Similar to the results of the convolutional layers, for the fully-connected layers we have

$$J^{m,i} \boldsymbol{v}^m = \frac{\partial \boldsymbol{z}^{L+1,i}}{\partial (\boldsymbol{s}^{m,i})^T} \text{mat}(\boldsymbol{v}^m)_{n_{m+1} \times (n_m+1)} \begin{bmatrix} \boldsymbol{z}^{m,i} \\ \mathbb{1}_1 \end{bmatrix}.$$

$$(J^{m,i})^T \boldsymbol{q}^i = \text{vec} \left( \left( \frac{\partial \boldsymbol{z}^{L+1,i}}{\partial (\boldsymbol{s}^{m,i})^T} \right)^T \boldsymbol{q}^i \left[ (\boldsymbol{z}^{m,i})^T \ \mathbb{1}_1 \right] \right).$$

# Implementation I

- As before, we must handle all instances together
- We discuss only

$$
\begin{bmatrix}
\sum_{m=1}^{L} J^{m,1} \mathbf{v}^m \\
\vdots \\
\sum_{m=1}^{L} J^{m,l} \mathbf{v}^m
\end{bmatrix}
\in R^{n_{L+1} l \times 1}
$$

- Following earlier derivation

# Implementation II

$$\begin{bmatrix} J^{m,1} \boldsymbol{v}^m \\ \vdots \\ J^{m,l} \boldsymbol{v}^m \end{bmatrix} = \begin{bmatrix} \frac{\partial \boldsymbol{z}^{L+1,1}}{\partial \text{vec}(S^{m,1})^T} \text{vec} \left( \text{mat}(\boldsymbol{v}^m) \begin{bmatrix} \phi(\text{pad}(Z^{m,1})) \\ \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}^T \end{bmatrix} \right) \\ \vdots \\ \frac{\partial \boldsymbol{z}^{L+1,l}}{\partial \text{vec}(S^{m,l})^T} \text{vec} \left( \text{mat}(\boldsymbol{v}^m) \begin{bmatrix} \phi(\text{pad}(Z^{m,l})) \\ \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}^T \end{bmatrix} \right) \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial \boldsymbol{z}^{L+1,1}}{\partial \text{vec}(S^{m,1})^T} \boldsymbol{p}^{m,1} \\ \vdots \\ \frac{\partial \boldsymbol{z}^{L+1,l}}{\partial \text{vec}(S^{m,l})^T} \boldsymbol{p}^{m,l} \end{bmatrix},$$

# Implementation III

- We have

$$\mathrm{mat}(\boldsymbol{v}^m) \in R^{d^{m+1} \times (h^m h^m d^m + 1)}$$

  and

$$\boldsymbol{p}^{m,i} = \mathrm{vec}\left(\mathrm{mat}(\boldsymbol{v}^m) \begin{bmatrix} \phi(\mathrm{pad}(Z^{m,i})) \\ \mathbb{1}^T_{a^m_{\mathrm{conv}} b^m_{\mathrm{conv}}} \end{bmatrix}\right). \qquad (9)$$

# Implementation IV

- All
$$\boldsymbol{p}^{m,i}, i = 1, \ldots, l$$
can be calculated by a matrix-matrix product

$$\mathrm{mat}(\boldsymbol{v}^m) \begin{bmatrix} \phi(\mathrm{pad}(Z^{m,1})) & \cdots & \phi(\mathrm{pad}(Z^{m,l})) \\ \mathbb{1}^T_{a^m_{\mathrm{conv}} b^m_{\mathrm{conv}}} & \cdots & \mathbb{1}^T_{a^m_{\mathrm{conv}} b^m_{\mathrm{conv}}} \end{bmatrix}$$
$$\in R^{d^{m+1} \times a^m_{\mathrm{conv}} b^m_{\mathrm{conv}} l};$$

# Implementation V

- To get

$$
\begin{bmatrix}
\frac{\partial z^{L+1,1}}{\partial \mathrm{vec}(S^{m,1})^T} p^{m,1} \\
\vdots \\
\frac{\partial z^{L+1,l}}{\partial \mathrm{vec}(S^{m,l})^T} p^{m,l}
\end{bmatrix},
$$

  we need $l$ matrix-vector products

- There is no good way to transform it to matrix-matrix operations

# Implementation VI

- At this moment we calculate

$$J^{m,i} \boldsymbol{v}^m = \frac{\partial \boldsymbol{z}^{L+1,i}}{\partial \mathrm{vec}(S^{m,i})^T} \boldsymbol{p}^{m,i}, \ i = 1, \ldots, l. \quad (10)$$

by summing up all rows of the following matrix

$$\left[ \frac{\partial z_1^{L+1,i}}{\partial \mathrm{vec}(S^{m,i})} \cdots \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial \mathrm{vec}(S^{m,i})} \right]_{d^{m+1} a_{\mathrm{conv}}^m b_{\mathrm{conv}}^m \times n_{L+1}} \odot$$

$$\left[ \boldsymbol{p}^{m,i} \cdots \boldsymbol{p}^{m,i} \right]_{d^{m+1} a_{\mathrm{conv}}^m b_{\mathrm{conv}}^m \times n_{L+1}}.$$

and extend this to cover all instances together

# Implementation VII

- The code (convolutional layers) is like

```
for m = LC : -1 : 1
    var_range = var_ptr(m) : var_ptr(m+1) - 1;
    ab = model.ht_conv(m)*model.wd_conv(m);
    d = model.ch_input(m+1);

    p = reshape(v(var_range), d, []) *
        [net.phiZ{m}; ones(1, ab*num_data)];
    p = sum(reshape(net.dzdS{m}, d*ab, nL,
            []) .*
            reshape(p, d*ab, 1, [])),1);
```

# Implementation VIII

```
    Jv = Jv + p(:);
end
```

# Outline

# Outline

# Reverse versus Forward Autodiff I

- We mentioned before that two types of autodiff are forward and reverse modes

- For the Jacobian evaluation, at layer $m$,

$$J^{m,i} = \left[ \frac{\partial \boldsymbol{z}^{L+1,i}}{\partial \mathrm{vec}(W^m)^T} \quad \frac{\partial \boldsymbol{z}^{L+1,i}}{\partial (\boldsymbol{b}^m)^T} \right],$$

  naturally we follow the gradient calculation to use the reverse mode

- But this may not be a good decision

- We will show a solution of using the forward mode

# R Operator I

- Consider $g(\boldsymbol{\theta}) \in R^{k \times 1}$. Following Pearlmutter (1994), we define

$$\mathcal{R}_{\boldsymbol{v}}\{g(\boldsymbol{\theta})\} \equiv \frac{\partial g(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T} \boldsymbol{v} = \begin{bmatrix} \nabla g_1(\boldsymbol{\theta})^T \boldsymbol{v} \\ \vdots \\ \nabla g_k(\boldsymbol{\theta})^T \boldsymbol{v} \end{bmatrix}. \qquad (11)$$

- Note that

$$\begin{bmatrix} \nabla g_1(\boldsymbol{\theta})^T \\ \vdots \\ \nabla g_k(\boldsymbol{\theta})^T \end{bmatrix}$$

is the Jacobian of $g(\boldsymbol{\theta})$

# R Operator II

- This definition can be extended to a matrix $M(\boldsymbol{\theta}) \in R^{k \times t}$ by

$$\mathcal{R}_{\boldsymbol{v}}\{M(\boldsymbol{\theta})\} \equiv \mathsf{mat}\left(\mathcal{R}_{\boldsymbol{v}}\{\mathsf{vec}(M(\boldsymbol{\theta}))\}\right)_{k \times t}$$

$$=\mathsf{mat}\left(\frac{\partial \mathsf{vec}(M(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}^T}\boldsymbol{v}\right)_{k \times t} = \begin{bmatrix} \nabla M_{11}^T \boldsymbol{v} & \cdots & \nabla M_{1t}^T \boldsymbol{v} \\ \vdots & \ddots & \vdots \\ \nabla M_{k1}^T \boldsymbol{v} & \cdots & \nabla M_{kt}^T \boldsymbol{v} \end{bmatrix}$$

- Clearly,

$$\mathcal{R}_{\boldsymbol{v}}\{M(\boldsymbol{\theta})\} = \left(\mathcal{R}_{\boldsymbol{v}}\{M(\boldsymbol{\theta})^T\}\right)^T. \qquad (12)$$

# R Operator III

- If $h(\cdot)$ is a scalar function, we let

$$h(M(\boldsymbol{\theta})) = \begin{bmatrix} h(M_{11}) & \cdots & h(M_{1t}) \\ \vdots & \ddots & \vdots \\ h(M_{k1}) & \cdots & h(M_{kt}) \end{bmatrix}$$

and

$$h'(M(\boldsymbol{\theta})) = \begin{bmatrix} h'(M_{11}) & \cdots & h'(M_{1t}) \\ \vdots & \ddots & \vdots \\ h'(M_{k1}) & \cdots & h'(M_{kt}) \end{bmatrix}.$$

# R Operator IV

- Because

$$\nabla(h(M_{ij}(\boldsymbol{\theta})))^T \boldsymbol{v} = h'(M_{ij})\nabla(M_{ij})^T \boldsymbol{v},$$

we have

$$\mathcal{R}_{\boldsymbol{v}}\{h(M(\boldsymbol{\theta}))\} = h'(M(\boldsymbol{\theta})) \odot \mathcal{R}_{\boldsymbol{v}}\{M(\boldsymbol{\theta})\}, \quad (13)$$

where $\odot$ stands for the Hadamard product.

- If $M(\boldsymbol{\theta})$ and $T(\boldsymbol{\theta})$ have the same size,

$$\mathcal{R}_{\boldsymbol{v}}\{M(\boldsymbol{\theta}) + T(\boldsymbol{\theta})\} = \mathcal{R}_{\boldsymbol{v}}\{M(\boldsymbol{\theta})\} + \mathcal{R}_{\boldsymbol{v}}\{T(\boldsymbol{\theta})\}. \quad (14)$$

# R Operator V

- Lastly, we have

$$\mathcal{R}_{\boldsymbol{v}}\{U(\boldsymbol{\theta})M(\boldsymbol{\theta})\} = \mathcal{R}_{\boldsymbol{v}}\{U(\boldsymbol{\theta})\}M(\boldsymbol{\theta}) + U(\boldsymbol{\theta})\mathcal{R}_{\boldsymbol{v}}\{M(\boldsymbol{\theta})\}$$
$$(15)$$

Proof: Note that

$$\left(\mathcal{R}\{U(\boldsymbol{\theta})M(\boldsymbol{\theta})\}\right)_{ij} = \nabla\left((U(\boldsymbol{\theta})M(\boldsymbol{\theta}))_{ij}\right)^T \boldsymbol{v}. \quad (16)$$

With

$$(U(\boldsymbol{\theta})M(\boldsymbol{\theta}))_{ij} = \sum_{p=1}^{m} U_{ip}M_{pj}, \quad (17)$$

# R Operator VI

we have both $U_{ip} \in R^1$ and $M_{pj} \in R^1$. Then,

$$\nabla \left( U_{ip} M_{pj} \right)^T \boldsymbol{v} = \left( (\nabla U_{ip})^T \boldsymbol{v} \right) M_{pj} + U_{ip} \left( (\nabla M_{pj})^T \boldsymbol{v} \right).$$

- For simplicity, subsequently we use $\mathcal{R}\{g(\boldsymbol{\theta})\}$ to be $\mathcal{R}_{\boldsymbol{v}}\{g(\boldsymbol{\theta})\}$

# R Operator for $J^i v$ I

- We have

$$J^i v = \mathcal{R}\{z^{L+1,i}\}.$$

- Now assume

$$\mathcal{R}\{Z^{m,i}\}$$

is available from the previous layer
- We consider the following forward operations
- From (15), we have

$$\mathcal{R}\{\phi(\text{pad}(Z^{m,i}))\}$$
$$= \text{mat}\left(P_\phi^{m,i} P_{\text{pad}}^{m,i} \mathcal{R}\{\text{vec}\left(Z^{m,i}\right)\}\right)_{h^m h^m d^m \times a_{\text{conv}}^m b_{\text{conv}}^m}$$

# R Operator for $J^i v$ II

- From (14), (15), we have

$$\mathcal{R}\{S^{m,i}\}$$
$$=\mathcal{R}\{W^m \phi(\text{pad}(Z^{m,i})) + \boldsymbol{b}^m \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}^T\}$$
$$=\mathcal{R}\{W^m \phi(\text{pad}(Z^{m,i}))\} + \mathcal{R}\{\boldsymbol{b}^m \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}^T\}$$
$$=\mathcal{R}\{W^m\}\phi(\text{pad}(Z^{m,i})) + W^m \mathcal{R}\{\phi(\text{pad}(Z^{m,i}))\} +$$
$$\quad \mathcal{R}\{\boldsymbol{b}^m\} \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}^T$$
$$=V_W^m \phi(\text{pad}(Z^{m,i})) + W^m \mathcal{R}\{\phi(\text{pad}(Z^{m,i}))\} +$$
$$\quad \boldsymbol{v}_b^m \mathbb{1}_{a_{\text{conv}}^m b_{\text{conv}}^m}^T,$$

# R Operator for $J^i v$ III

where we use

$$\mathcal{R}\{W^m\} = V_W^m,$$
$$\mathcal{R}\{\boldsymbol{b}^m\} = \boldsymbol{v}_b^m.$$

- Note that

$$\boldsymbol{v} = \begin{bmatrix} \boldsymbol{v}^1 \\ \vdots \\ \boldsymbol{v}^L \end{bmatrix},$$

and each $\boldsymbol{v}^m, m = 1, \ldots, L$ has the same length as the number of variables (including bias) at the $m$th layer.

# R Operator for $J^i v$ IV

- We further split $v^m$ to $V_W^m$ (a matrix form) and $v_b^m$
- From (13), we have

$$\mathcal{R}\{\sigma(S^{m,i})\} = \sigma'(S^{m,i}) \odot \mathcal{R}\{S^{m,i}\}. \qquad (18)$$

- From (15), we have

$$\begin{aligned}
&\mathcal{R}\{Z^{m+1,i}\} \\
=&\mathcal{R}\{P_{\text{pool}}^{m,i}\sigma(S^{m,i})\} \\
=&\text{mat}\left(P_{\text{pool}}^{m,i}\mathcal{R}\{\text{vec}\left(\sigma(S^{m,i})\right)\}\right)_{d^{m+1}\times a^{m+1}b^{m+1}}.
\end{aligned}$$

# R Operator for $J^i v$ V

- We can continue this process until we get

$$J^i v = \mathcal{R}\{z^{L+1,i}\}.$$

- Clearly, we do not need to store

$$\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}, \ldots, \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}$$

  as before

# Outline

# Gauss-Newton Matrix-vector Product I

- From the above discussion, we have known how to calculate

$$J^i \boldsymbol{v}$$

- Calculate

$$B^i(J^i \boldsymbol{v})$$

  is known to be easy

# Gauss-Newton Matrix-vector Product II

- Now for
$$(J^i)^T(B^i J^i \mathbf{v}),$$
if we define
$$\mathbf{u} = B^i J^i \mathbf{v},$$
then
$$(J^i)^T \mathbf{u} = \left(\frac{\partial \mathbf{z}^{L+1,i}}{\partial \boldsymbol{\theta}^T}\right)^T \mathbf{u}.$$

- But earlier the gradient calculation is

$$(J^i)^T \nabla_{\mathbf{z}^{L+1,i}} \xi(\mathbf{z}^{L+1,i}; \mathbf{y}^i, Z^{1,i}) = \left(\frac{\partial \mathbf{z}^{L+1,i}}{\partial \boldsymbol{\theta}^T}\right)^T \frac{\partial \xi_i}{\partial \mathbf{z}^{L+1,i}}$$

# Gauss-Newton Matrix-vector Product III

- Thus the same backward procedure can be used
- All we need is to replace

$$\frac{\partial \xi_i}{\partial \boldsymbol{z}^{L+1,i}}$$

with

$$\boldsymbol{u}$$

# Complexity Analysis I

- We have known from past slides that matrix-matrix products are the bottleneck (though in our cases some slow MATLAB functions are also bottlenecks in practice)

- For simplicity, in our analysis we just count the number of matrix-matrix products

- Approaches solely by backward settings: if

$$\frac{\partial z_1^{L+1,i}}{\partial S^{m,i}}, \cdots, \frac{\partial z_{n_{L+1}}^{L+1,i}}{\partial S^{m,i}}$$

# Complexity Analysis II

stored, then

$$n_{L+1} \times 3 + \#CG \times 2$$

If not, then

$$\#CG \times (n_{L+1} \times 3 + 2)$$

- Note that "3" comes from one product in the forward process and two in the backward process (the same as the situation in Gradient calculation)

# Complexity Analysis III

- If using R operators, then

$$\#\text{CG} \times (3 + 2),$$

where "3" are from the forward process

$$W^m \phi(\text{pad}(Z^{m,i})),$$

and

$$V_W^m \phi(\text{pad}(Z^{m,i})), W^m \mathcal{R}\{\phi(\text{pad}(Z^{m,i}))\},$$

and "2" are from the backward process

# Discussion I

- At this moment in the Python code we are not using the forward mode for $Jv$

- It was not available before

- However, since version 2.10 released in January 2020, this functionality is provided:

  `https://www.tensorflow.org/api_docs/python/tf/autodiff/ForwardAccumulator`

- It will be interesting to do the implementation and make a comparison