

1 FP of VNN: Column Version

Have L hidden layers
for k = 0

$$\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$$

for k from 1 to L

a) Hidden layer pre-activation

$$\mathbf{a}^{(k)}(\mathbf{x}) = W^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x}) + \mathbf{b}^{(k)}$$

b) Hidden layer activation

$$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{f}(\mathbf{a}^{(k)}(\mathbf{x}))$$

for k = L+1

Output layer activation

$$\mathbf{a}^{(L+1)}(\mathbf{x}) = W^{(L+1)}\mathbf{h}^{(L)}(\mathbf{x}) + \mathbf{b}^{(L+1)}$$

$$\hat{\mathbf{y}} = \mathbf{h}^{(L+1)}(\mathbf{x}) = \text{softmax}(\mathbf{a}^{(L+1)}(\mathbf{x}))$$

2 BP of VNN: Column Version

$$\nabla_{\mathbf{a}^{(L+1)}(\mathbf{x})} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y}) = \hat{\mathbf{y}} - \mathbf{y}$$

For k from L+1 to 1

a) compute gradient of hidden layer parameter

$$\nabla_{W^{(k)}} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y}) = \nabla_{\mathbf{a}^{(k)}(\mathbf{x})} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y}) * \mathbf{h}^{(k-1)}(\mathbf{x})^T$$

$$\nabla_{\mathbf{b}^{(k)}} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y}) = \nabla_{\mathbf{a}^{(k)}(\mathbf{x})} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y})$$

b) compute gradient of hidden layer below(stop when it comes to $\mathbf{x} = \mathbf{h}^{(0)}(\mathbf{x})$)

$$\nabla_{\mathbf{h}^{(k-1)}(\mathbf{x})} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y}) = (W^{(k)})^T * (\nabla_{\mathbf{a}^{(k)}(\mathbf{x})} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y}))$$

c) compute gradient of hidden layer below(before activation)

$$\nabla_{\mathbf{a}^{(k-1)}(\mathbf{x})} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y}) = \nabla_{\mathbf{h}^{(k-1)}(\mathbf{x})} \text{Loss}(\hat{\mathbf{y}}, \mathbf{y}) \odot [\dots, f'(a^{(k-1)}(\mathbf{x})_j), \dots]$$

3 FP of VNN: Row Version

Have L hidden layers
for k = 0

$$\mathbf{d}^{(0)}(\mathbf{v}) = \mathbf{v}$$

for k from 1 to L

a) Hidden layer pre-activation

$$\mathbf{z}^{(k)}(\mathbf{v}) = \mathbf{d}^{(k-1)}(\mathbf{v})U^{(k)} + \mathbf{c}^{(k)}$$

b) Hidden layer activation

$$\mathbf{d}^{(k)}(\mathbf{v}) = \mathbf{g}(\mathbf{z}^{(k)}(\mathbf{v}))$$

for k = L+1

Output layer activation

$$\mathbf{z}^{(L+1)}(\mathbf{v}) = \mathbf{d}^{(L)}(\mathbf{v})U^{(L+1)} + \mathbf{c}^{(L+1)}$$

$$\hat{\mathbf{o}} = \mathbf{d}^{(L+1)}(\mathbf{x}) = \text{softmax}(\mathbf{z}^{(L+1)}(\mathbf{v}))$$

4 BP of VNN: Row Version

$$\nabla_{\mathbf{z}^{(L+1)}(\mathbf{v})} CE(\hat{\mathbf{o}}, \mathbf{o}) = \hat{\mathbf{o}} - \mathbf{o}$$

For k from L+1 to 1

a) compute gradient of hidden layer parameter

$$\nabla_{U^{(k)}} CE(\hat{\mathbf{o}}, \mathbf{o}) = \mathbf{d}^{(k-1)}(\mathbf{v})^T * \nabla_{\mathbf{z}^{(k)}(\mathbf{v})} CE(\hat{\mathbf{o}}, \mathbf{o})$$

$$\nabla_{\mathbf{c}^{(k)}} CE(\hat{\mathbf{o}}, \mathbf{o}) = \nabla_{\mathbf{z}^{(k)}(\mathbf{v})} CE(\hat{\mathbf{o}}, \mathbf{o})$$

b) compute gradient of hidden layer below(stop when it comes to $\mathbf{v} = \mathbf{d}^{(0)}(\mathbf{v})$)

$$\nabla_{\mathbf{d}^{(k-1)}(\mathbf{v})} CE(\hat{\mathbf{o}}, \mathbf{o}) = (\nabla_{\mathbf{z}^{(k)}(\mathbf{x})} CE(\hat{\mathbf{o}}, \mathbf{o})) * (U^{(k)})^T$$

c) compute gradient of hidden layer below(before activation)

$$\nabla_{\mathbf{z}^{(k-1)}(\mathbf{v})} CE(\hat{\mathbf{o}}, \mathbf{o}) = \nabla_{\mathbf{d}^{(k-1)}(\mathbf{x})} CE(\hat{\mathbf{o}}, \mathbf{o}) \odot [\dots, g'(z^{(k-1)}(\mathbf{v})_j), \dots]$$

5 VNN: Example for L=1

$$\nabla_{\mathbf{z}_2} CE = \hat{\mathbf{o}} - \mathbf{o}$$

$$\nabla_{U_2} CE = \mathbf{d}_1^T * \nabla_{\mathbf{z}_2} CE$$

$$\nabla_{\mathbf{c}_2} CE = \nabla_{\mathbf{z}_2} CE$$

$$\nabla_{\mathbf{d}_1} CE = \nabla_{\mathbf{z}_2} CE * U_2^T$$

$$\nabla_{\mathbf{z}_1} CE = \nabla_{\mathbf{d}_1} CE \odot \sigma'(\mathbf{z}_1)$$

$$\nabla_{U_1} CE = \mathbf{v}^T * \nabla_{\mathbf{z}_1} CE$$

$$\nabla_{\mathbf{c}_1} CE = \nabla_{\mathbf{z}_1} CE$$

Notice: $\sigma'(\mathbf{z}_1) = \text{sigmoid}_{grad}(\mathbf{d}_1)$