

Bells and whistles in neural net training

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Project Add WeChat edu_assist_pro

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Tricks in training neural networks

<https://eduassistpro.github.io/>

Assignment Project Exam Help

There are various tricks that people use when training neural networks:

- ▶ Regularization
- ▶ Dropout: Adjusting the probability of dropping out
- ▶ Optimization methods: Adjusting the learning rate
- ▶ Initialization: Using particular forms of initialization

Regularization

<https://eduassistpro.github.io/>

Neural networks can be regularized in a similar way as linear models. Neural networks can also with **Frobenius norm**, which is a trivial extension to L_2 norm for matrices. In fact, in machine learning it is just referred to as L_2 regularization.

<https://eduassistpro.github.io/>

$$\mathcal{L} = \sum_{i=1}^N \ell^{(i)} + \lambda_{z \rightarrow y} \|\Theta^{(z \rightarrow y)}\|_F^2$$

[Add WeChat edu_assist_pro](#)

where $\|\Theta\|_F^2 = \sum_{i,j} \theta_{i,j}^2$ is the squared **Frobenius norm**, which generalizes the L_2 norm to matrices. The bias parameters b are not regularized, as they do not contribute to the classifier to the inputs.

L2 regularization

- Compute the gradient $\frac{\partial \mathcal{L}}{\partial \theta} = \sum_{n=1}^N \frac{\partial \ell^{(n)}}{\partial \theta}$

Assignment Project Exam Help

Add WeChat edu_assist_pro

- Update the weights $\theta = \theta - \eta \left(\frac{\partial \mathcal{L}}{\partial \theta} + \lambda \theta \right)$

Add WeChat edu_assist_pro

- “Weigh decay factor”: λ is a tunable hyper parameter that pulls a weight back when it has become too big
- Question: Does it matter which layer θ is from when computing the regularization term?

L1 regularization

- L1 regularization <https://eduassistpro.github.io/>

$$\mathcal{L} = \sum_{i=1}^N \ell^{(i)} + \lambda_{z \rightarrow y} \|\Theta^{(z \rightarrow y)}\|_1$$

- Compute the

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- update the weights

$$\theta = \theta - \eta \left(\sum_{i=1}^N \frac{\partial \ell^{(i)}}{\partial \theta} + \lambda \operatorname{sign}(\theta) \right)$$

Comparison of L1 and L2

<https://eduassistpro.github.io/>

- ▶ In L1 regularization, the weights shrink by a constant amount toward 0. In L2 regularization, the weights shrink by an amount which is proportional to w .
- ▶ When a particular weight has a large absolute value, L1 regularization shrinks it more than L2 regularization. In L1 regularization, the weights are driven toward zero, while in L2 regularization, the weights are driven toward a small value.
- ▶ The net result is that L1 regularization tends to keep the weight of the network in a relatively small number of high-importance connections, while the other weights are driven toward zero. So L1 regularization effectively does *feature selection*.

Dropout

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ Randomly drops a certain percentage of the **feature** over-reliance **co-adaptation** working together is to avoid overfitting. **goal**

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Dropout

- Dropout can be ac

Assignment Project Exam Help

$$\mathbf{z}^{(1)} = g(\Theta^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

Assignment Project Help edu_assist_pro

$$\mathbf{m}^1 \sim \text{Bernoulli}(r)$$
$$\tilde{\mathbf{z}}^{(1)} = \mathbf{m}^1 \odot \mathbf{z}^{(1)}$$

https://eduassistpro.github.io/

$$\mathbf{m}^2 \sim \text{Bernoulli}(r)$$

Add WeChat edu_assist_pro

$$\tilde{\mathbf{z}}^{(2)} = \mathbf{m}^2 \odot \mathbf{z}^{(2)}$$

$$\mathbf{y} = \Theta^{(3)}\tilde{\mathbf{z}}^{(2)}$$

where \mathbf{m}^1 and \mathbf{m}^2 are mask vectors. The values of the elements in these vectors are either 1 or 0, drawn from a Bernoulli distribution with parameter r (usually $r = 0.5$)

Optimization methods

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ SGD with Momentum
- ▶ AdaGrad
- ▶ Root Mean Sq
- ▶ Adam

Add WeChat edu_assist_pro

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

SGD with Momentum

<https://eduassistpro.github.io/>

- At each timestep t , compute $\nabla_{\theta} \mathcal{L}$, and then compute the momentum as follows:

[Add WeChat edu_assist_pro](#)

<https://eduassistpro.github.io/>

[Add WeChat edu_assist_pro](#)

- The momentum term increases for dimensions whose gradient point in the same directions and reduces updates for dimensions whose gradient change directions.

AdaGrad

- <https://eduassistpro.github.io/>
- Keep a running sum of the squared gradient $V_{\nabla_{\theta}}$. When updating the weight of this theta, divide the gradient by the square root of this term

[Add WeChat edu_assist_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

e.g., $\epsilon = 10^{-8}$

- The net effect is to slow down the update for weights with large gradient and accelerate the update for weights with small gradient

Root Mean Square Prop (RMSProp)

<https://eduassistpro.github.io/>

- A minor adjustment of AdaGrad. Instead of letting the sum of squared gradient continuously grow, we let it

[Add WeChat edu_assist_pro](#)

$$V = 0$$

<https://eduassistpro.github.io/>

[Add WeChat edu_assist_pro](#)

e.g. $\beta \approx 0.9, \eta = 0.001, \epsilon = 10^{-8}$

Adaptive Moment Estimation (Adam)

- Weight update at t

Assignment Project Exam Help

$V_0 = 0, S_0 = 0,$
 $V_t = \beta_1 V_{t-1} + (1 - \beta_1) m$

$S_t = \beta_2 S_{t-1} + (1 - \beta_2) S^{Prop}$

$V_t = \overline{\beta_1}$
 $S_t^{corrected} = \frac{S_t}{\beta_2^t}$

$$\theta_j = \theta_j - \eta \frac{V_t^{corrected}}{\sqrt{S_t^{corrected} + \epsilon}}$$

- Adam combines Momentum and RMSProp

Initialization

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Xavier Initialization:

Add WeChat edu_assist_pro

$$\Theta = \left[\begin{array}{c} \text{https://eduassistpro.github.io/} \\ \text{Add WeChat edu_assist_pro} \end{array} \right]_{(l+1)}$$

where $n^{(l)}$ is the number of input units to Θ (and $n^{(l+1)}$ is the number of output units from Θ)

Neural net in PyTorch

```
from torch import nn
class Net(nn.Module):
    """subclass from nn.Module is
    important for inspecting the parameters"""

    def __init__(self, in_dim=25, out_dim=10, ch_size=1):
        super(Net, self).__init__()
        self.in_dim = in_dim
        self.out_dim = out_dim
        self.softmax = nn.Softmax()

    def forward(self, input_matrix):
        logit = self.linear(input_matrix)
        #return raw score, not normalized score
        return logit

    def xentropy_loss(self, input_matrix, target_label_vec):
        loss = nn.CrossEntropyLoss()
        logits = self.forward(input_matrix)
        return loss(logits, target_label_vec)
```

Use optimizers in Pytorch

```
import torch.optim as optim
net = Net(input_dim, output_dim)
optimizer = optim.Adam(net.parameters(), lr=lr_rate)
for epoch in range(epochs):
    total_nll = 0
    for batch in batch_loader(train_data_loader):
        optimizer.zero_grad_() #zero the grad buffer
        vectors = torch.cat([feature_vectors, label_vectors], dim=1)
        feat_vec = torch.nn.functional.linear(vectors, weight)
        label_vec = torch.nn.functional.linear(vectors, bias)
        feat_list = list(feat_vec.data)
        label_list = list(label_vec.data)
        x = torch.Tensor(feat_list)
        y = torch.LongTensor(label_list)
        loss = net.xentropy_loss(x, y)
        total_nll += loss
        loss.backward()
        optimizer.step()
torch.save(net.state_dict(), net_path)
```