

Exercise sheet 2: Training and Regularization

In this homework, we will train neural networks on the Boston housing dataset. For this, we will use the modular framework developed in Homework 1. A first part of the homework will analyze the parameters of the network before and after training. A second part of the homework will test some regularization penalties and their effect on the generalization error.

Boston Housing Dataset

The following code extracts the Boston housing dataset in a way that is already partitioned into training and test data. The data is normalized such that each dimension has mean 0 and variance 1.

```
In [1]: import utils

Xtrain,Ttrain,Xtest,Ttest = utils.getBostonHousingData()
```

Neural Network Regressor

In this homework, we will consider a very simple architecture consisting of one linear layer, a ReLU layer applying a nonlinear function element-wise, and a pooling layer computing a fixed weighted sum of the activations obtained in the previous layer. The architecture is shown below:

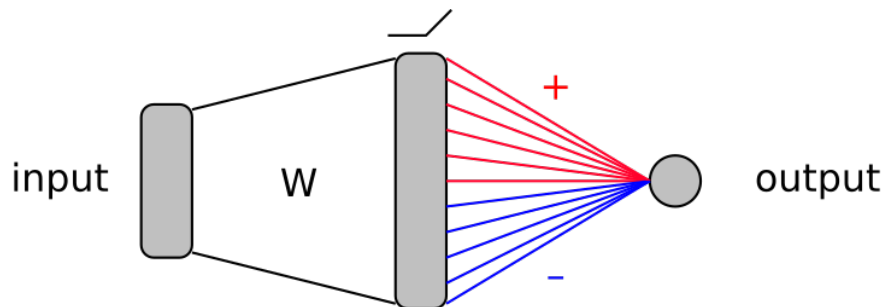


Figure 1: Diagram of the Neural Network Regressor used in this homework

The function `getarch` shown below generates an architecture specific to the Boston housing dataset, with 13 input features, h intermediate neurons, and one output. It takes as a parameter the first layer type, which is usually `layers.Linear`. Later in this homework, we will replace it by variants of `layers.Linear` that incorporate weight norm penalties.

```
In [2]: import layers

def getarch(FirstLayerType):
    h = 100
    return layers.Sequential([
        FirstLayerType(13,h,seed=0),
        layers.ReLU(),
        layers.Pooling(h)
    ])
```

The class `NeuralNetworkRegressor` shown below takes a neural network architecture, trains it on the data using gradient descent with momentum, and predict new data points. Because the dataset is rather small, we do not need to use stochastic gradient descent. We choose instead the batch variant, where we follow the gradient of the true mean square error.

```
In [3]: import sklearn,sklearn.metrics

class NeuralNetworkRegressor:
```

```

def __init__(self, arch): self.arch = arch

def fit(self, X, T, nbit=10000):

    for i in range(nbit):
        self.arch.backward(2*(self.arch.forward(X)-T))
        for l in self.arch.layers: l.compute_grad(); l.pgradstep(0.01)

def predict(self, X):
    return self.arch.forward(X)

def score(self, X, T):
    return sklearn.metrics.r2_score(T, self.predict(X))

```

Neural Network Performance vs. Baselines

We compare the performance of the neural network on the Boston housing data to two other regressors: a random forest and a support vector regression model with RBF kernel. We use the scikit-learn implementation of these models, with their default parameters.

```
In [4]: import sklearn.ensemble, sklearn.svm
```

```

rfr = sklearn.ensemble.RandomForestRegressor()
rfr.fit(Xtrain, Ttrain)

svr = sklearn.svm.SVR()
svr.fit(Xtrain, Ttrain)

nnr = NeuralNetworkRegressor(getarch(layers.Linear))
nnr.fit(Xtrain, Ttrain)

```

```
In [5]: def pretty(name, model):
        return '> %10s | R2train: %6.3f | R2test: %6.3f'%(name, model.score(Xtrain, Ttrain), model.score(Xtest, Ttest))

print(pretty('RForest', rfr))
print(pretty('SVR', svr))
print(pretty('NNreg', nnr))
```

```

> RForest | R2train: 0.967 | R2test: 0.831
> SVR | R2train: 0.913 | R2test: 0.758
> NNreg | R2train: 0.996 | R2test: 0.845

```

The neural networks performs on par with the other regression models although it has in principle the added ability to learn its own internal features. We would therefore expect a well-trained neural network to perform better.

Gradient, and Parameter Norms (20 P)

As a first step towards improving the neural network model, we will measure proxy quantities, that will then be used to regularize the model. We consider the following three quantities:

- $\|W\|_{\text{Frob}} = \sqrt{\sum_{i=1}^d \sum_{j=1}^h w_{ij}^2}$
- $\|W\|_{\text{Mix}} = h^{-0.5} \sqrt{\sum_{i=1}^d \left(\sum_{j=1}^h |w_{ij}| \right)^2}$
- $\text{Grad} = \frac{1}{N} \sum_{n=1}^N \|\nabla_{\mathbf{x}} f(\mathbf{x}_n)\|$

where d is the number of input features, h is the number of neurons in the hidden layer, and W is the matrix of weights in the first layer. In order for the model to generalize well, the last quantity (Grad) should be prevented from becoming too large. Because the latter depends on the data distribution, we rely instead on the inequality $\text{Grad} \leq \|W\|_{\text{Mix}} \leq \|W\|_{\text{Frob}}$, that we can prove for this model, and will try to control the weight norms instead.

Tasks:

- Implement the function `WMix(nn)` that receives the neural network as input and returns $\|W\|_{\text{Mix}}$.
- Implement the function `Grad(nn,X)` that receives the neural network and some dataset as input, and computes the averaged gradient norm (`Grad`).

```
In [6]: def WFrob(nn):
        W = nn.arch.layers[0].W
        return (W**2).sum()**.5

        def WMix(nn):
            # -----
            # TODO: Replace by your code
            # -----
            import solution; return solution.MixNorm(nn)
            # -----

        def Grad(nn,X):
            # -----
            # TODO: Replace by your code
            # -----
            import solution; return solution.AvgGrad(nn,X)
            # -----
```

The following experiment computes these three quantities before and after training the model.

```
In [7]: def fullpretty(name,nn):
        return pretty(name,nn) + ' | WFrob: %7.3f | WMix: %7.3f | Grad: %7.3f'%(WFrob(nn),WMix(nn),Grad(n

        nnr = NeuralNetworkRegressor(getarch(layers.Linear))
        print(fullpretty('Before',nnr))
        nnr.fit(Xtrain,Ttrain)
        print(fullpretty('After',nnr))

> Before | R2train: -0.043 | R2test: -0.006 | WFrob: 9.834 | WMix: 7.845 | Grad: 0.617
> After | R2train: 0.996 | R2test: 0.845 | WFrob: 18.752 | WMix: 15.018 | Grad: 1.201
```

We observe that the inequality $\text{Grad} \leq \|W\|_{\text{Mix}} < \|W\|_{\text{Frob}}$ also holds empirically. We also observe that these quantities tend to increase as training proceeds. This is a typical behavior, as the network starts rather simple and becomes complex as more and more variations in the training data are being captured.

Frobenius-Norm Penalty (10 P)

The first penalty we experiment with is the squared Frobenius norm. We consider the new objective $J^{\text{Frob}}(\theta) = \text{MSE}(\theta) + \lambda \cdot \|W\|_{\text{Frob}}^2$, where the first term is the original mean square error objective and where the second term is the added penalty. We hardcode the penalty coefficient to $\lambda = 10^{-3}$. In principle, for maximum performance and fair comparison between the methods, several of them should be tried (also for other models), and selected based on some validation set. Here, for simplicity, we omit this step.

Tasks:

- Create a new layer, that replaces the original `layers.Linear` layer and incorporates the proposed penalty. (*Hint: This can be achieved by extending `layers.Linear`, and rewriting its method `compute_grad`.*)

```
In [8]: # -----
        # TODO: Replace by your code
        # -----
        import solution
        LinearFrob = solution.LinearFrob
        # -----
```

The code below trains a neural network with the new first layer, and compares the performance with the previous models.

```
In [9]: nnfrob = NeuralNetworkRegressor(getarch(LinearFrob))
        nnfrob.fit(Xtrain,Ttrain)

In [10]: print(pretty('RForest',rfr))
         print(pretty('SVR',svr))
         print(fullpretty('NN',nnr))
         print(fullpretty('NN+Frob',nnfrob))

> RForest | R2train:  0.967 | R2test:  0.831
>      SVR | R2train:  0.913 | R2test:  0.758
>      NN | R2train:  0.996 | R2test:  0.845 | WFrob:  18.752 | WMix:  15.018 | Grad:   1.201
> NN+Frob | R2train:  0.967 | R2test:  0.833 | WFrob:   5.119 | WMix:   4.139 | Grad:   0.806
```

As we can see, the penalty has the effect of strongly reducing the Frobenius norm of the weight matrix. However, the test accuracy is not improved.

Mixed-Norm Penalty (10 P)

A downside of the Frobenius norm is that it is not a very tight upper bound of the gradient, that is, penalizing it does not penalize specifically high gradient. Instead, other useful properties of the model could be negatively affected by it. In the following, we experiment with the mixed-norm regularizer, which is a tighter bound of the gradient.

The objective function is now given by $J^{\text{Mix}}(\theta) = \text{MSE}(\theta) + \lambda \cdot \|W\|_{\text{Mix}}^2$, where the first and second terms are again the original objective and the penalty term. As for the previous exercise, we hardcode the latter coefficient to 10^{-3} .

Tasks:

- Using the same technique as before, create a new layer that incorporates the Mixed norm penalty.

```
In [11]: # -----
         # TODO: Replace by your code
         # -----
         import solution
         LinearMix = solution.LinearMix
         # -----
```

The code below trains a neural network with the new penalty, and compares the performance with the previous models.

```
In [12]: nnmix = NeuralNetworkRegressor(getarch(LinearMix))
        nnmix.fit(Xtrain,Ttrain)

In [13]: print(pretty('RForest',rfr))
         print(pretty('SVR',svr))
         print(fullpretty('NN',nnr))
         print(fullpretty('NN+Frob',nnfrob))
         print(fullpretty('NN+Mix',nnmix))

> RForest | R2train:  0.967 | R2test:  0.831
>      SVR | R2train:  0.913 | R2test:  0.758
>      NN | R2train:  0.996 | R2test:  0.845 | WFrob:  18.752 | WMix:  15.018 | Grad:   1.201
> NN+Frob | R2train:  0.967 | R2test:  0.833 | WFrob:   5.119 | WMix:   4.139 | Grad:   0.806
> NN+Mix  | R2train:  0.978 | R2test:  0.856 | WFrob:   8.188 | WMix:   4.214 | Grad:   0.902
```

It is interesting to observe that this mixed norm penalty more selectively reduced the mixed norm and the gradient, and has let the Frobenius norm take higher values. The mixed-norm model is also the one that produces the highest test set accuracy.