

Introduction to deep learning in R

Toby Dylan Hocking

Tenure-track Assistant Professor, 2018–present,

Director of Machine Learning Research Lab,

School of Informatics, Computing, and Cyber Systems,

Northern Arizona University

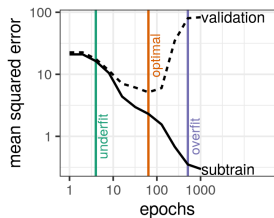
toby.hocking@nau.edu

toby.hocking@r-project.org

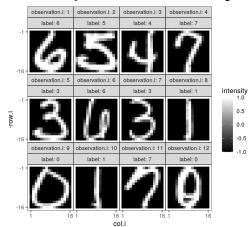
April 18, 2023

Introduction and overview

Example 1: avoiding overfitting in regression, overview of concepts



Example 2: classifying images of digits, coding demos





Summary and quiz questions

Machine learning intro: image classification example

ML is all about learning predictive functions $f(x) \approx y$, where

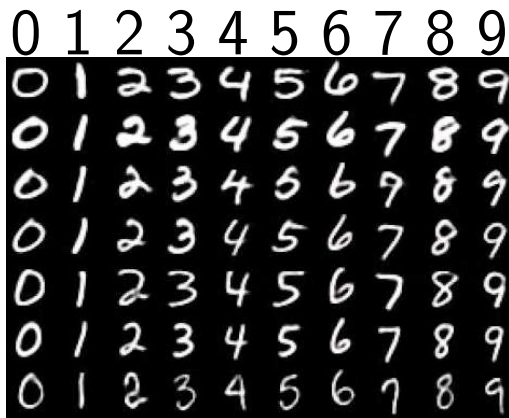
- ▶ Inputs/features x can be easily computed using traditional algorithms. For example, matrix of pixel intensities in an image.
- ▶ Outputs/labels y are what we want to predict, typically more difficult/costly to measure than inputs. For example, to get an image class label, you may have to ask a human.
- ▶ Input x = image of digit, output $y \in \{0, 1, \dots, 9\}$,
 - this is a classification problem with 10 classes.

 $f(\text{image of } 0) = 0$,  $f(\text{image of } 1) = 1$

- ▶ Traditional/unsupervised algorithm: I give you a pixel intensity matrix $x \in \mathbb{R}^{28 \times 28}$, you code a function f that returns one of the 10 possible digits. Q: how to do that?

Supervised machine learning algorithms

I give you a training data set with paired inputs/outputs, e.g.



Your job is to code an algorithm, `LEARN`, that infers a function f from the training data. (you don't code f)

Source: github.com/cazala/mnist

Advantages of supervised machine learning

Learning Algorithm	Train data	Learned function	Predictions on test data
--------------------	------------	------------------	--------------------------

Learn(



) → g

$g(\text{0}) = 0$

$g(\text{1}) = 1$

$g(\text{1}) = 1$

Learn(



) → h

$h(\text{0}) = 0$

$h(\text{0}) = 0$

$h(\text{1}) = 1$

- ▶ Input $x \in \mathbb{R}^{28 \times 28}$, output $y \in \{0, 1, \dots, 9\}$ types the same!
- ▶ Can use same learning algorithm regardless of pattern.
- ▶ Pattern encoded in the labels (not the algorithm).
- ▶ Useful if there are many un-labeled data, but few labeled data (or getting labels is long/costly).
- ▶ State-of-the-art accuracy (if there is enough training data).

Sources: github.com/cazala/mnist, github.com/zalando-research/fashion-mnist

Overview of tutorial

In this tutorial we will discuss two kinds of problems, which differ by the type of the output/label/ y variable we want to predict.

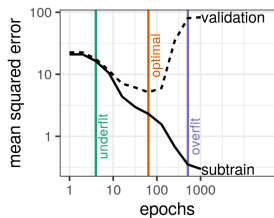
- ▶ Regression, y is a real number.
- ▶ Classification, y is an integer representing a category.

The rest of the tutorial will focus on three examples:

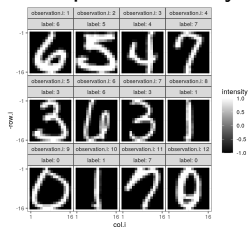
1. Regression with a single input, to demonstrate how to avoid overfitting.
2. Classification of digit images, to demonstrate how to compare machine learning algorithms in terms of test/prediction accuracy.

Introduction and overview

Example 1: avoiding overfitting in regression, overview of concepts



Example 2: classifying images of digits, coding demos

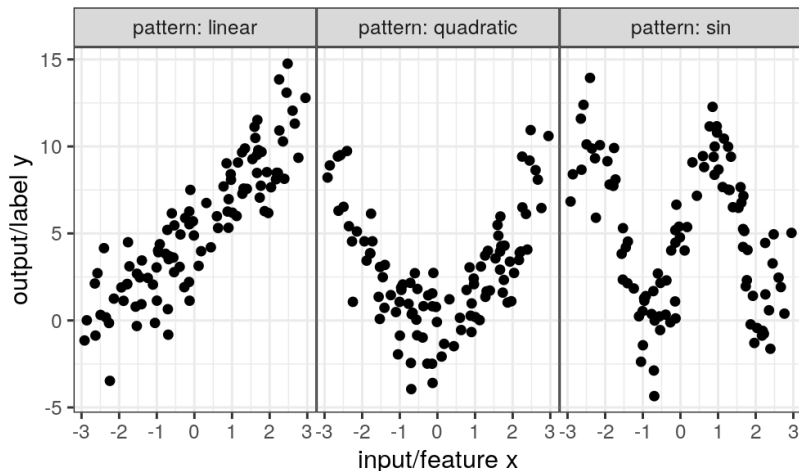


Summary and quiz questions

Goal of this section: demonstrate how to avoid overfitting

- ▶ The goal of supervised machine learning is to get accurate predictions on new/unseen/held-out test data.
- ▶ The data used during learning are called the train set.
- ▶ Any machine learning algorithm could overfit, if not used properly. Overfitting means better predictions on the train set than on a held-out validation/test set. (BAD)
- ▶ To learn a model which does NOT overfit (GOOD), you need to divide your train set into subtrain/validation sets (subtrain used as input to gradient descent algorithm, validation set used to control model complexity, for example by selecting the best number of iterations of gradient descent).
- ▶ Code for figures in this section:
<https://github.com/tdhock/2023-res-baz-az/blob/main/figure-overfitting.R>

Three different data sets/patterns



- ▶ We illustrate this using a single input/feature $x \in \mathbb{R}$.
- ▶ We use a regression problem with outputs $y \in \mathbb{R}$.
- ▶ Goal is to learn a function $f(x) \in \mathbb{R}$.

K-fold cross-validation for splitting data

- ▶ One way to split is via K-fold cross-validation.
- ▶ Each row is assigned a fold ID number from 1 to K.
- ▶ For each for ID, those data are held out, and other data are kept.
- ▶ Popular relative to other splitting methods because of simplicity and fairness (each row is held out one time).

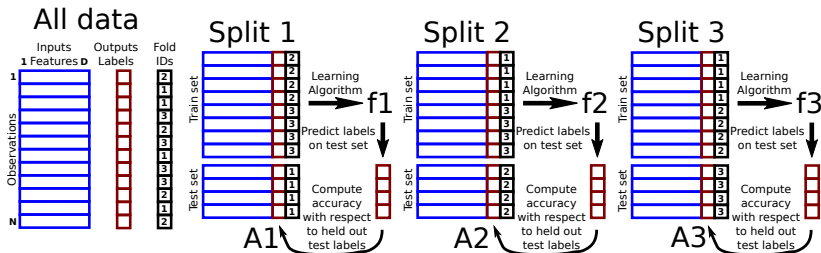
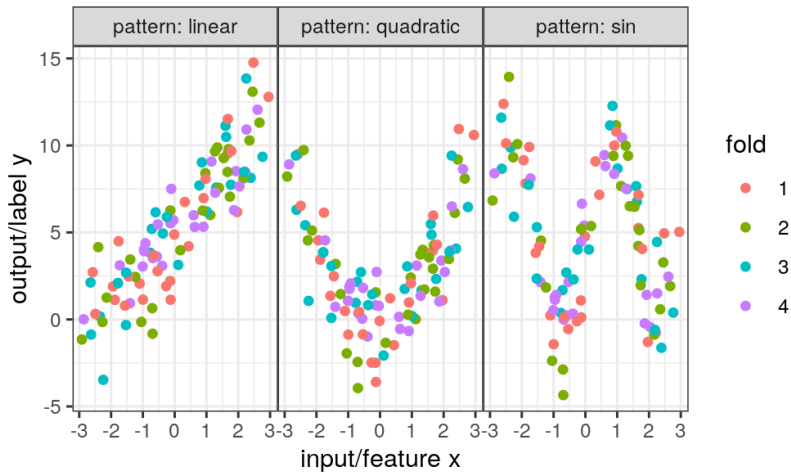


Illustration of 4-fold cross-validation

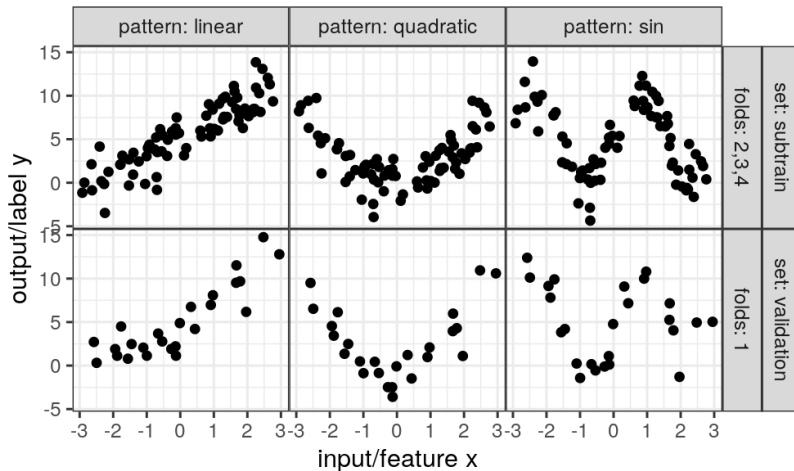


Randomly assign each observation a fold ID from 1 to 4.

Neural network learning algorithm

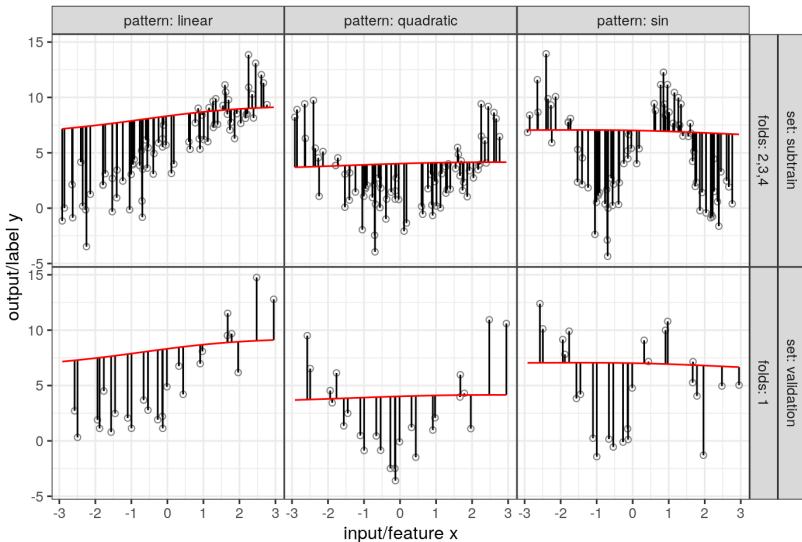
- ▶ We will fit a neural network to these data.
- ▶ The neural network learns how to predict the outputs from the inputs.
- ▶ The learning algorithm is gradient descent, which iteratively minimizes the loss of the predictions with respect to the labels in the subtrain set.
- ▶ We also compute the loss on the validation set, so we can select the number of gradient descent iterations that gives the best predictions on new data (avoiding overfitting).

Illustration of subtrain/validation split



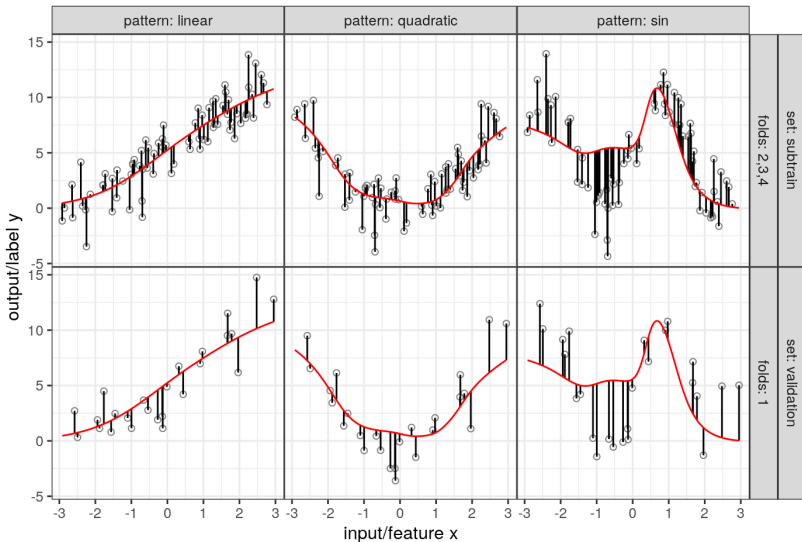
- ▶ For validation fold 1, all observations with that fold ID are considered the validation set.
- ▶ All other observations are considered the subtrain set.

Neural network, 20 hidden units, 1 gradient descent iterations



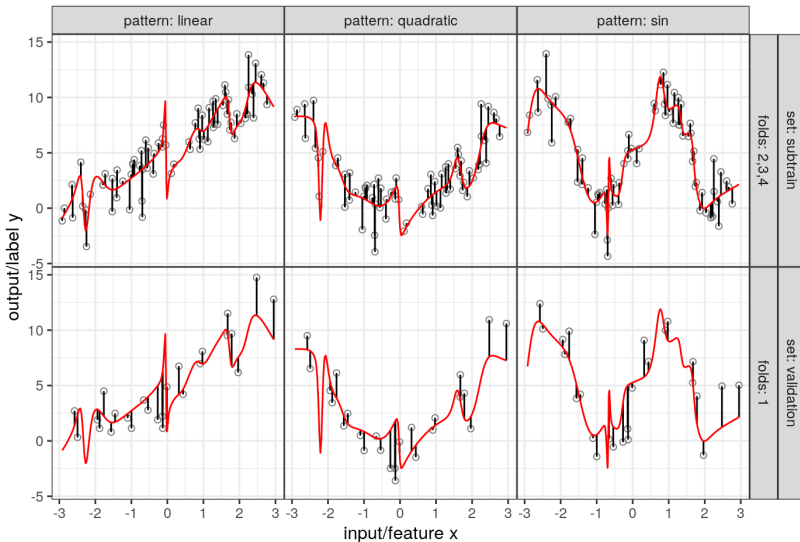
Data=grey dots, predictions=red curve, loss=black line segments.

Neural network, 20 hidden units, 10 gradient descent iterations



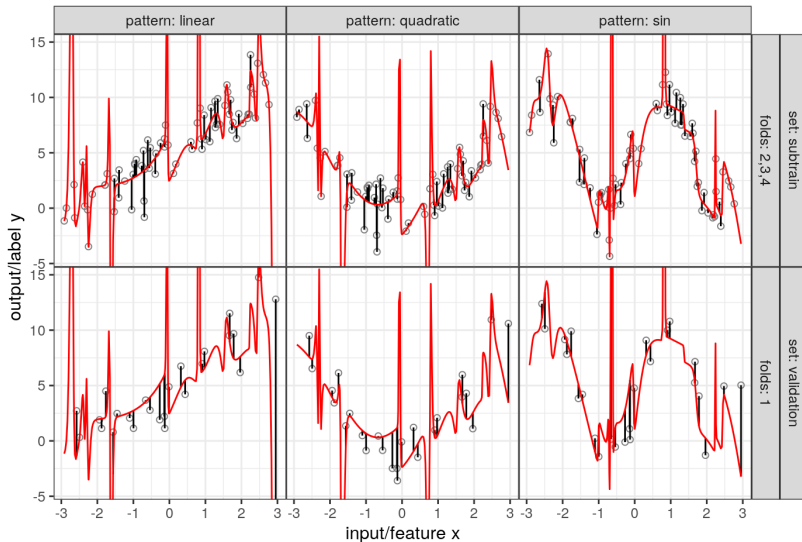
Data=grey dots, predictions=red curve, loss=black line segments.

Neural network, 20 hidden units, 100 gradient descent iterations



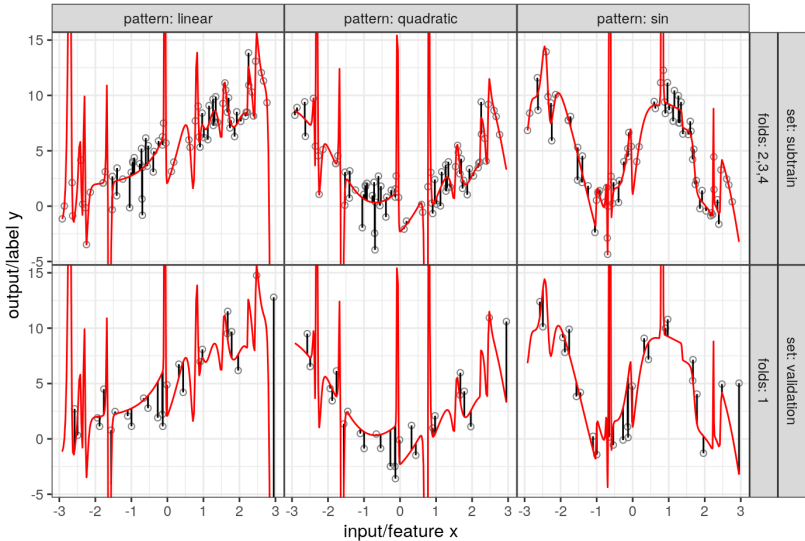
Data=grey dots, predictions=red curve, loss=black line segments.

Neural network, 20 hidden units, 1000 gradient descent iterations



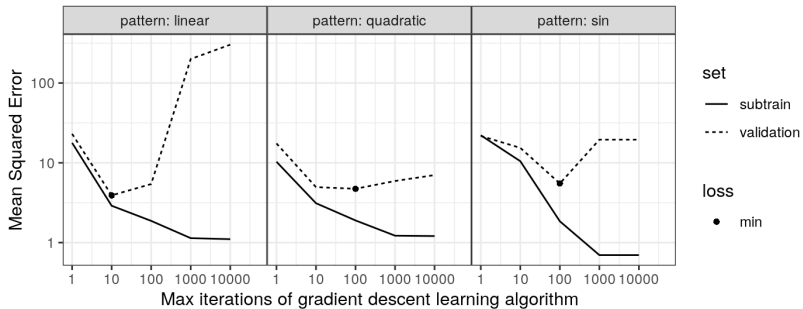
Data=grey dots, predictions=red curve, loss=black line segments.

Neural network, 20 hidden units, 10000 gradient descent iterations



Data=grey dots, predictions=red curve, loss=black line segments.

Neural network, 20 hidden units



Different number of iterations best for different data.

Neural network prediction function

For an input feature vector \mathbf{x} , the prediction function for a neural network with L functions to learn:

$$f(\mathbf{x}) = f_L[\cdots f_1[\mathbf{x}]]. \quad (1)$$

For each layer $l \in \{1, \dots, L\}$, we have the following function for predicting the units/features in the next layer:

$$f_l(t) = A_l(\mathbf{W}_l^\top t), \quad (2)$$

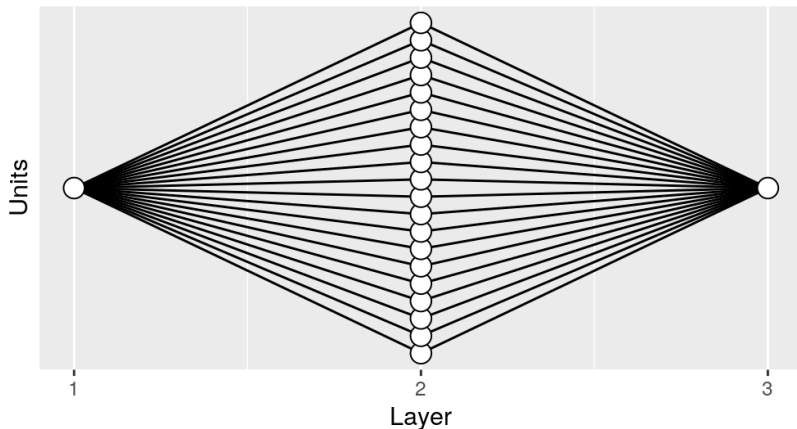
The hyper-parameters which must be fixed prior to learning:

- ▶ Number of layers or functions to learn L .
- ▶ Activation functions A_l (classically sigmoid, typically ReLU).
- ▶ Number of units per layer $(u_0, u_1, \dots, u_{L-1}, u_L)$. u_0 is the number of input features, u_L is the number of outputs, and u_1, \dots, u_{L-1} are the numbers of hidden units in the intermediate layers (hyper-parameters, larger for more powerful learner).
- ▶ Sparsity pattern in the weight matrices $\mathbf{W}_l \in \mathbb{R}^{u_l \times u_{l-1}}$.

Network for 1 input, 1 output, 1 hidden layer

Neural network diagrams show how each unit (node) is computed by applying the weights (edges) to the values of the units at the previous layer.

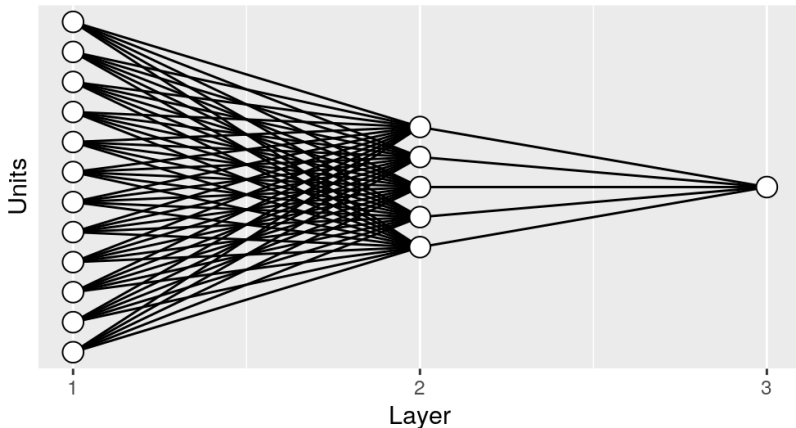
Number of units: 1,20,1



Network for 12 inputs, 1 output, 1 hidden layer

Neural network diagrams show how each unit (node) is computed by applying the weights (edges) to the values of the units at the previous layer.

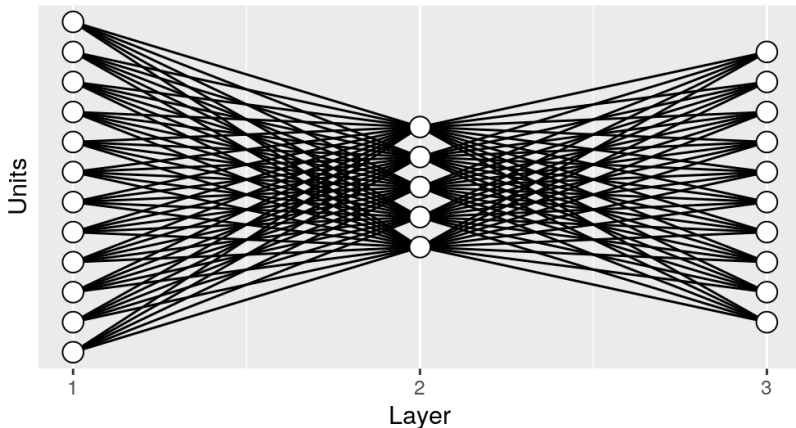
Number of units: 12,5,1



Network for 12 inputs, 10 outputs, 1 hidden layer

Neural network diagrams show how each unit (node) is computed by applying the weights (edges) to the values of the units at the previous layer.

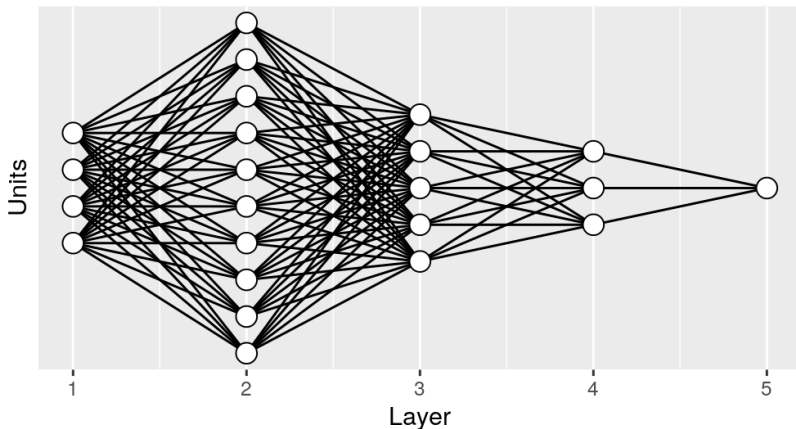
Number of units: 12,5,10



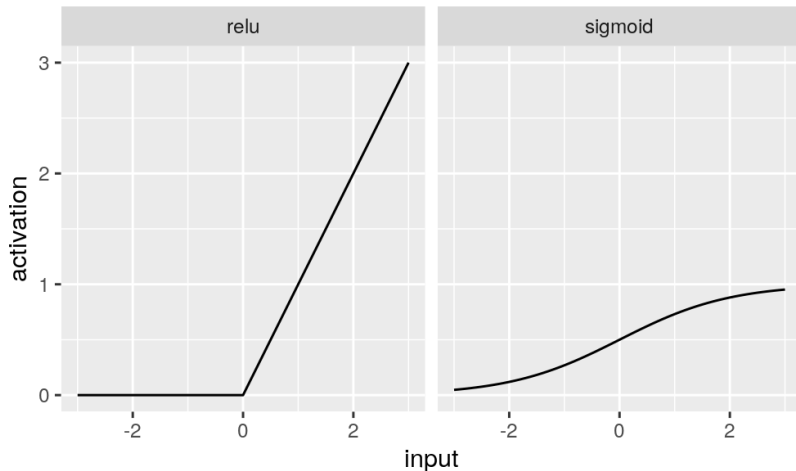
Network for 4 inputs, 1 output, 3 hidden layers

Neural network diagrams show how each unit (node) is computed by applying the weights (edges) to the values of the units at the previous layer.

Number of units: 4,10,5,3,1



Non-linear activation functions A_l



Each layer except the last should have a activation function A_l which is not linear (last layer activation should be identity/linear).

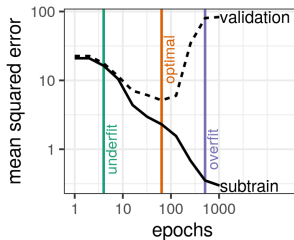
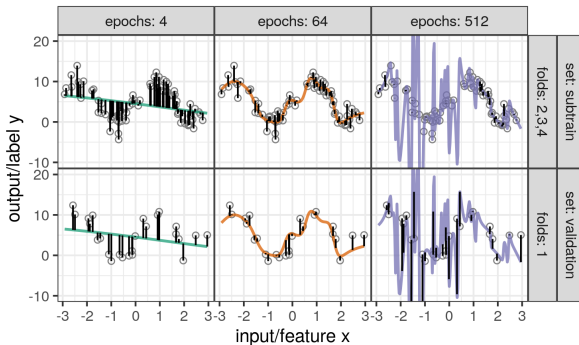
Gradient Descent Learning

The neural network prediction function $f(\mathbf{x}) = f_L[\cdots f_1[\mathbf{x}]]$ has $l \in \{1, \dots, L\}$ component functions to learn:

$$f_l(t) = A_l(\mathbf{W}_l^\top t), \quad (3)$$

The weight matrices $\mathbf{W}_l \in \mathbb{R}^{u_l \times u_{l-1}}$ are learned using gradient descent.

- ▶ A loss function $\mathcal{L}[f(\mathbf{x}), y]$ computes how bad are predictions with respect to labels y (ex: mean squared error for regression, cross entropy loss for classification).
- ▶ In each **iteration** of gradient descent, the weights are updated in order to get better predictions on subtrain data.
- ▶ An **epoch** computes gradients on all subtrain data; there can be from 1 to $N(\text{subtrain size})$ iterations per epoch.



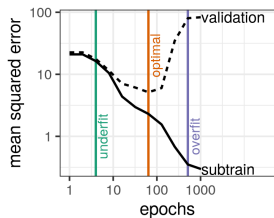
Citation: Hocking TD. *Introduction to machine learning and neural networks* in book *Land Carbon Cycle Modeling: Matrix Approach, Data Assimilation, and Ecological Forecasting*, edited by Luo Y, published by Taylor and Francis, 2022.

Summary of how to avoid overfitting

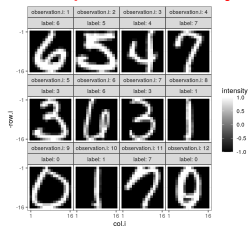
- ▶ Happens when subtrain error/loss decreases but validation error increases (as a function of some hyper-parameter)
- ▶ Here the hyper-parameter is the number of iterations of gradient descent, and overfitting starts after a certain number of iterations.
- ▶ To maximize prediction accuracy you need to choose a hyper-parameter with minimal validation error/loss.
- ▶ This optimal hyper-parameter will depend on the data set.
- ▶ To get optimal prediction accuracy in any machine learning analysis, you always need to do this, because you never know the best hyper-parameters in advance.

Introduction and overview

Example 1: avoiding overfitting in regression, overview of concepts



Example 2: classifying images of digits, coding demos



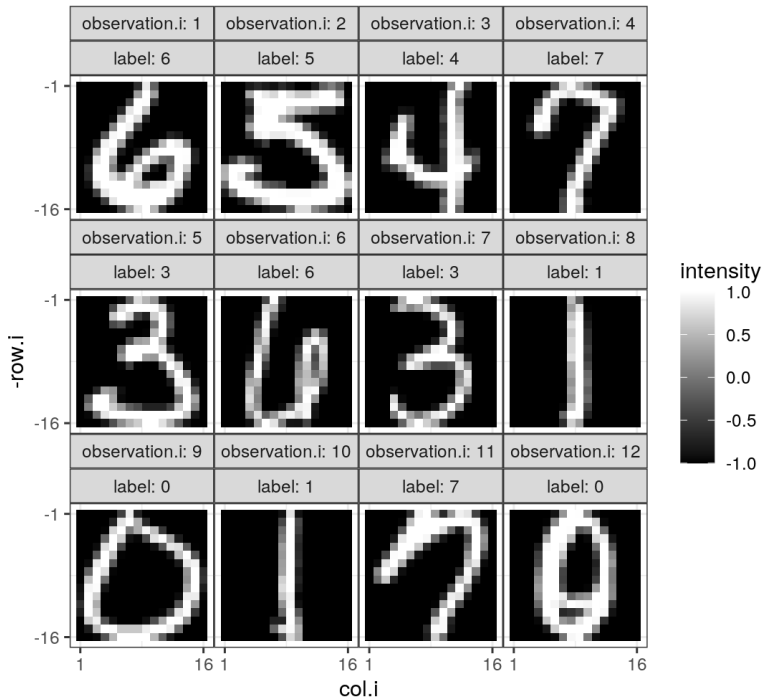
Summary and quiz questions

Image classification

- ▶ A problem in computer vision, one of the most popular/successful application domains of machine learning.
- ▶ Input: image file $x \in \mathbb{R}^{h \times w \times c}$ where h is the height in pixels, w is the width, c is the number of channels, e.g. RGB image $c = 3$ channels.
- ▶ In this tutorial we use images with $h = w = 16$ pixels and $c = 1$ channel (grayscale, smaller values are darker).
- ▶ Output: class/category y (from a finite set).
- ▶ In this tutorial there are ten image classes $y \in \{0, 1, \dots, 9\}$, one for each digit.



- ▶ Want to learn f such that $f(\text{image of 0}) = 0$, $f(\text{image of 1}) = 1$, etc.
- ▶ Code for figures in this section:
<https://github.com/tdhock/2023-res-baz-az/blob/main/figure-validation-loss.R>



Representation of digits in CSV

- ▶ Each image/observation is one row.
- ▶ First column is output/label/class to predict.
- ▶ Other 256 columns are inputs/features (pixel intensity values).

Data from

<https://web.stanford.edu/~hastie/ElemStatLearn/datasets/zip.train.gz>

```
1:  6 -1 -1  ... -1.000 -1.000  -1
2:  5 -1 -1  ... -0.671 -0.828  -1
3:  4 -1 -1  ... -1.000 -1.000  -1
4:  7 -1 -1  ... -1.000 -1.000  -1
5:  3 -1 -1  ... -0.883 -1.000  -1
6:  6 -1 -1  ... -1.000 -1.000  -1
...
```

Demo: reading CSV, plotting digits,

<https://github.com/tdhock/2023-res-baz-az/blob/main/2023-04-19-deep-learning.Rmd>

Converting R data to torch tensors

Use array function with all columns except first as data.

```
zip.dt <- data.table::fread("zip.train.gz")
zip.X.array <- array(
  data = unlist(zip.dt[,-1]),
  dim = c(nrow(zip.dt), 1, 16, 16))
zip.X.tensor <- torch::torch_tensor(zip.X.array)
zip.y.tensor <- torch::torch_tensor(
  zip.dt$V1+1L, torch::torch_long())
```

Need to specify dimensions of input/X array:

- ▶ Observations: same as the number of rows in the CSV table.
- ▶ Channels: 1 (greyscale image, would be 3 for RGB image).
- ▶ Pixels wide: 16.
- ▶ Pixels high: 16.

For output/y need to add 1 in R, and specify long int type.

Linear model R code

```
n.features <- 16*16
n.classes <- 10
linear.model <- torch::nn_sequential(
  torch::nn_flatten(),
  torch::nn_linear(n.features, n.classes))
pred.tensor <- linear.model(zip.X.tensor)
```

- ▶ First layer must specify shape of inputs (here 16x16x1).
- ▶ `nn_flatten` converts any shape to a single dimension of units (here, convert each image from 1x16x16-array to 256-vector).
- ▶ `nn_linear` uses all units/features in the previous layer (256) to predict each unit in the next layer (10).
- ▶ There are ten possible classes for an output.

Computing loss, gradient, parameter updates

```
loss.fun <- torch::nn_cross_entropy_loss()
loss.tensor <- loss.fun(pred.tensor, zip.y.tensor)
step.size <- 0.1#also known as learning rate.
optimizer <- torch::optim_sgd(
  linear.model$parameters, lr=step.size)
optimizer$zero_grad()
loss.tensor$backward()
optimizer$step()
```

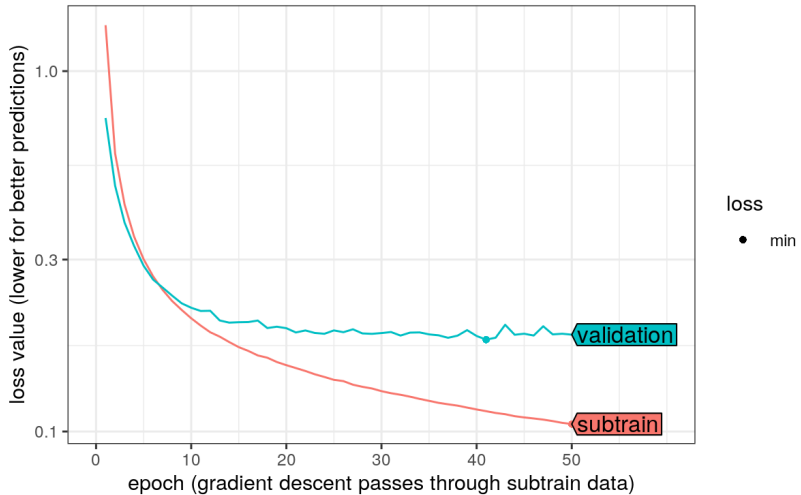
- ▶ `loss.fun` is the cross-entropy loss for multi-class classification, which is directly optimized/minimized in each iteration of the gradient descent learning algorithm.
- ▶ `optimizer` is the version of the gradient descent learning algorithm to use.
- ▶ `backward` method computes gradients.
- ▶ `step` method updates model parameters based on gradients.

Gradient Descent learning algorithm

```
gradient_descent <-  
  function(index.list, model, n_epochs, gradient.set){  
    loss.dt.list <- list()  
    for(epoch in seq(1, n_epochs)){  
      take_steps(index.list[[gradient.set]], model)  
      epoch.loss.dt <- loss_each_set(index.list, model)  
      loss.dt.list[[paste(epoch)]] <-  
        data.table(epoch, epoch.loss.dt)  
    }  
    rbindlist(loss.dt.list)  
  }
```

- ▶ `take_steps` sub-routine updates model parameters.
- ▶ `loss_each_set` computes loss and error rate on gradient set and held-out set.

Linear model



Demo: splitting data, gradient descent loop.

Dense (fully connected) neural network R code

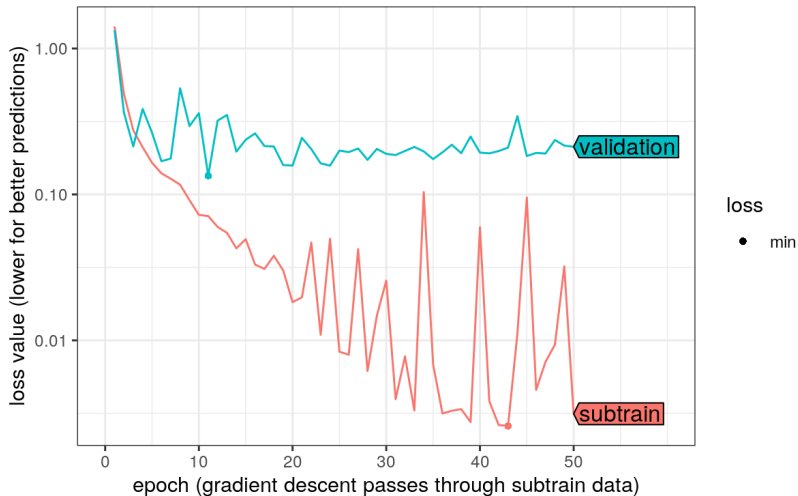
```
one.hidden.layer <- torch::nn_sequential(  
  torch::nn_flatten(),  
  torch::nn_linear(n.features, n.hidden.units),  
  torch::nn_relu(),  
  torch::nn_linear(n.hidden.units, n.classes))  
two.hidden.layers <- torch::nn_sequential(  
  torch::nn_flatten(),  
  torch::nn_linear(n.features, n.hidden.1),  
  torch::nn_relu(),  
  torch::nn_linear(n.hidden.1, n.hidden.2),  
  torch::nn_relu(),  
  torch::nn_linear(n.hidden.2, n.classes))
```

Use for loop to implement dense network

```
new_fully_connected_units <- function(units.per.layer){  
  seq.args <- list(torch::nn_flatten())  
  for(output.i in seq(2, length(units.per.layer))){  
    input.i <- output.i-1  
    seq.args[[length(seq.args)+1]] <- torch::nn_linear(  
      units.per.layer[[input.i]],  
      units.per.layer[[output.i]])  
    if(output.i<length(units.per.layer)){  
      seq.args[[length(seq.args)+1]] <- torch::nn_relu()  
    }  
  }  
  do.call(torch::nn_sequential, seq.args)  
}
```

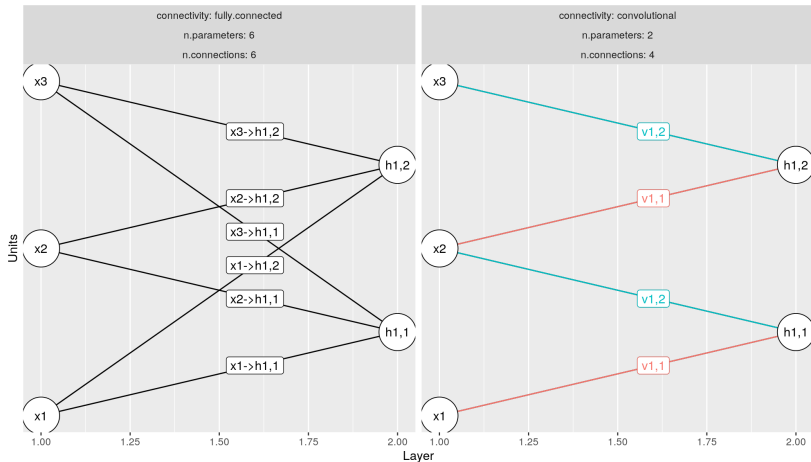
- ▶ input a vector of units per layer, for example `c(256,1000,100,10)`.
- ▶ Begin with flatten.
- ▶ Linear followed by relu in each layer except last.

Dense (fully connected) neural network with 8 hidden layers



Fully connected/dense vs convolutional/sparse network

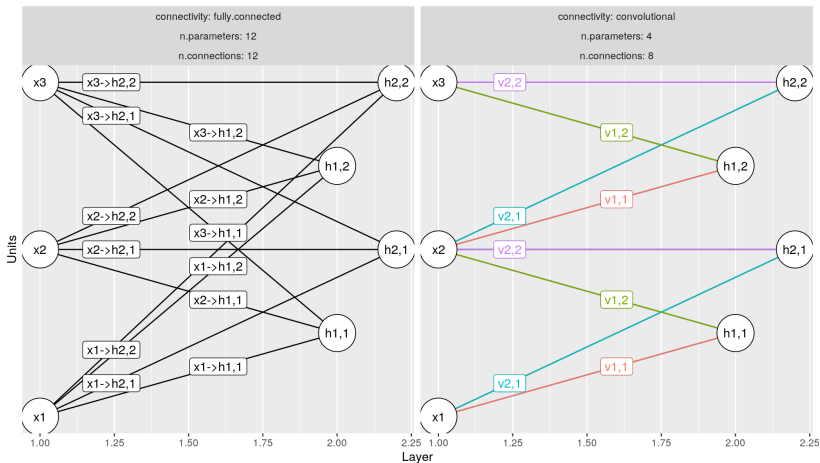
Number of units: 3, 2 filters: 1 kernel.size: 2



```
torch.nn.conv1d(in_channels=1,  
out_channels=1, kernel_size=2)
```

Convolutional with two filters/output channels

Number of units: 3,4 filters: 2 kernel.size: 2



```
torch.nn.conv1d(in_channels=1,  
out_channels=2, kernel_size=2)
```

Fully connected vs convolutional, two filters/channels

$$\begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ w_{2,1} & w_{2,2} & w_{2,3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \text{ (fully connected)}$$

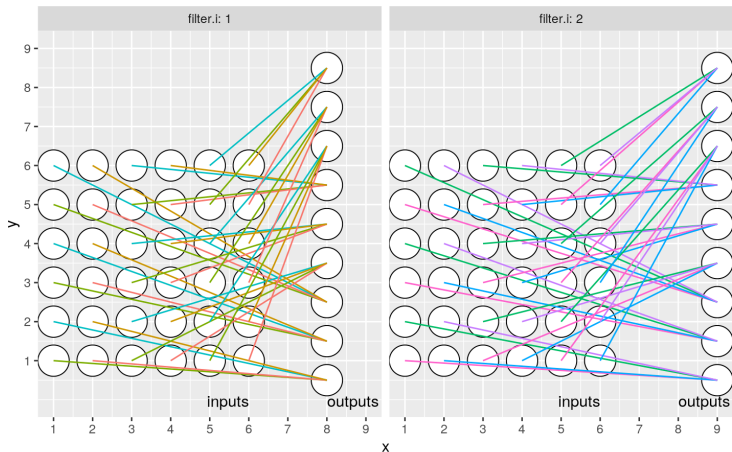
$$\begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} v_1 & v_2 & 0 \\ 0 & v_1 & v_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \text{ (convolutional)}$$

$$\begin{bmatrix} h_{1,1} \\ h_{1,2} \\ h_{2,1} \\ h_{2,2} \end{bmatrix} = \begin{bmatrix} v_{1,1} & v_{1,2} & 0 \\ 0 & v_{1,1} & v_{1,2} \\ v_{2,1} & v_{2,2} & 0 \\ 0 & v_{2,1} & v_{2,2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \text{ (conv, two filters)}$$

- ▶ Weight sharing: same weights used to compute different output units.
- ▶ Sparsity: zeros in weight matrix.

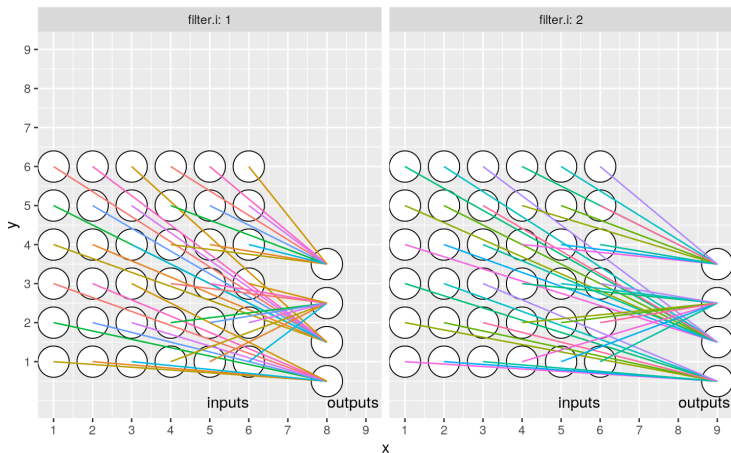
2D convolutional kernel for 6x6 pixel image, kernel size=2

6x6 pixel image input, convolutional kernel size = stride = 2,
n.filters/channels = 2, weights per channel = 4, outputs per channel = 9



2D convolutional kernel for 6x6 pixel image, kernel size=3

6x6 pixel image input, convolutional kernel size = stride = 3,
n.filters/channels = 2, weights per channel = 9, outputs per channel = 4

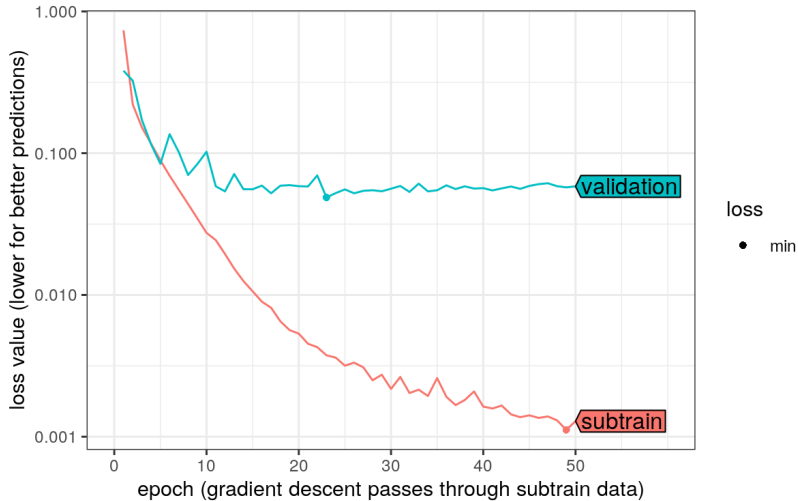


Sparse (convolutional) model R code

```
seq2flat <- torch::nn_sequential(  
  torch::nn_conv2d(  
    in_channels = 1, out_channels = 10, kernel_size = 4),  
  torch::nn_relu(),  
  torch::nn_flatten(),  
  torch::nn_linear(conv.hidden.units, last.hidden.units),  
  torch::nn_relu(),  
  torch::nn_linear(last.hidden.units, n.classes))
```

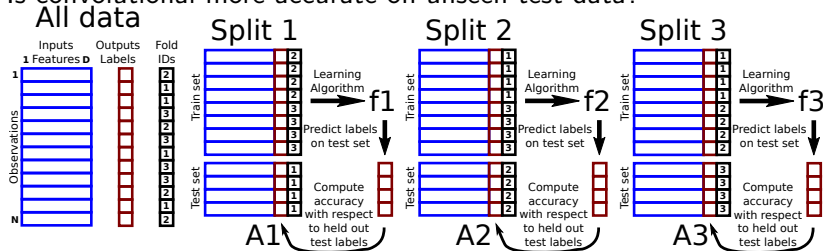
- ▶ Two hidden layers: one convolutional, one linear.
- ▶ Sparse: few inputs are used to predict each unit in `nn_conv2d`.
- ▶ Exploits structure of image data to make learning easier/faster.

Convolutional neural network



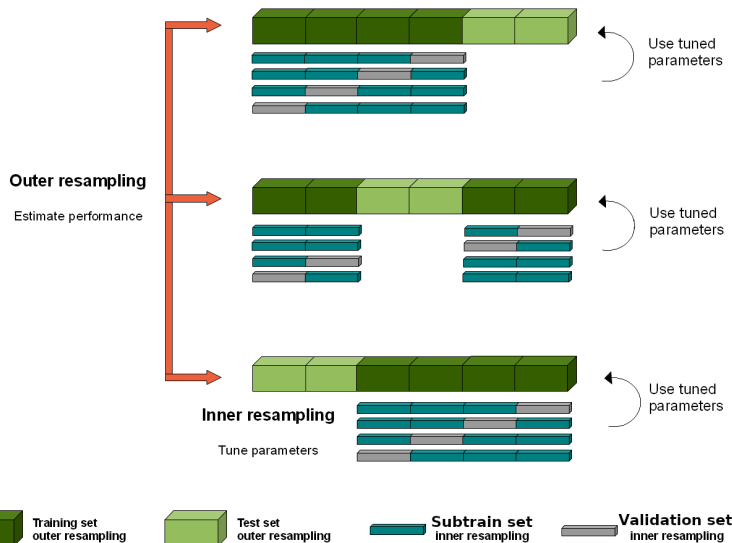
K-fold cross-validation for model evaluation

Is convolutional more accurate on unseen test data?



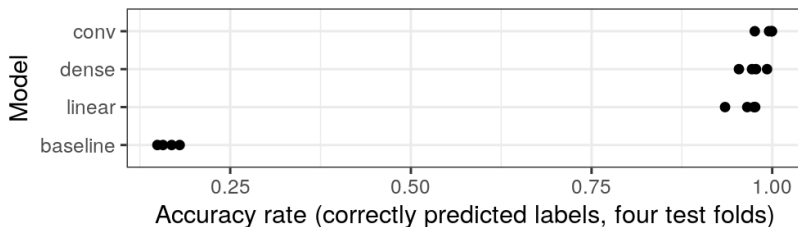
- ▶ Randomly assign a fold ID from 1 to K to each observation.
- ▶ Hold out the observations with the Split ID as test set.
- ▶ Use the other observations as the train set.
- ▶ Run learning algorithm on train set (including hyper-parameter selection), outputs learned function (f1-f3).
- ▶ Finally compute and plot the prediction accuracy (A1-A3) with respect to the held-out test set.

Two kinds of cross-validation must be used



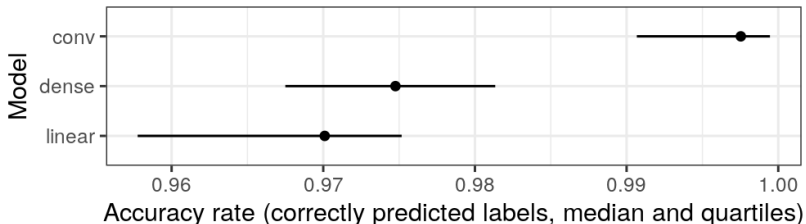
Source: https://mlr.mlr-org.com/articles/tutorial/nested_resampling.html

Accuracy rates for each test fold



- ▶ Always a good idea to compare with the trivial/featureless baseline model which always predicts the most frequent class in the train set. (ignoring all inputs/features)
- ▶ Here we see that the featureless baseline is much less accurate than the three learned models, which are clearly learning something non-trivial.
- ▶ Code for test accuracy figures:
<https://github.com/tdhock/2023-res-baz-az/blob/main/figure-test-accuracy.R>

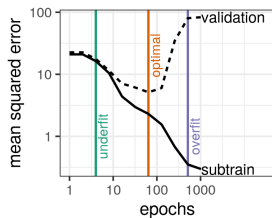
Zoom to learned models



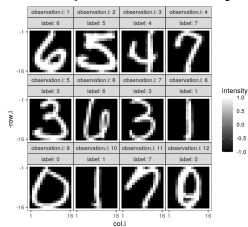
- ▶ Dense neural network slightly more accurate than linear model, convolutional significantly more accurate than others.
- ▶ Conclusion: convolutional neural network should be preferred for most accurate predictions in these data.
- ▶ Maybe not the same conclusion in other data sets, with the same models. (always need to do cross-validation experiments to see which model is best in any given data set)
- ▶ Maybe other models/algorithms would be even more accurate in these data. (more/less layers, more/less units, completely different algorithm such as random forests, boosting, etc)

Introduction and overview

Example 1: avoiding overfitting in regression, overview of concepts



Example 2: classifying images of digits, coding demos



Summary and quiz questions

Summary

We have studied

- ▶ Two kinds of machine learning problems, regression y =real number, classification y =integer category.
- ▶ Splitting a data set into train/test/subtrain/validation sets for learning hyper-parameters and evaluating prediction accuracy.
- ▶ Overfitting and how to avoid it using early stopping regularization (choosing the best number of iterations/epochs using a held-out validation set). Also, the validation set can and should be used to learn other hyper-parameters (number of hidden units/layers, learning rate/step size, batch size).
- ▶ How to use a held-out test set to compare prediction accuracy of learning algorithms with each other and to a featureless baseline.

Quiz questions

- ▶ When using a design matrix to represent machine learning inputs, what does each row and column represent?
- ▶ When splitting data into train/test sets, what is the purpose of each set? When splitting a train set into subtrain/validation sets, what is the purpose of each set?
- ▶ In order to determine if any non-trivial predictive relationship between inputs and output has been learned, a comparison with a featureless baseline that ignores the inputs must be used. How do you compute the featureless baseline predictions?
- ▶ How can you tell if machine learning model predictions are underfitting or overfitting?
- ▶ Many learning algorithms require input of the number of iterations or epochs. How should this parameter be chosen?

Thanks for participating!

- ▶ Contact: toby.hocking@nau.edu, toby.hocking@r-project.org
- ▶ Take my Deep Learning class if you want to know more about these topics,
<https://github.com/tdhock/cs499-599-fall-2022>
- ▶ The NAU SICCS Machine Learning Research lab is recruiting!
<http://ml.nau.edu/>

