

NEURAL NETWORKS AND DEEP LEARNING

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

OVERVIEW

Neural networks are models for supervised learning

Linear combinations of **features** are passed through a non-linear transformation in successive layers

At the top layer, the resulting **latent factors** are fed into an algorithm for predictions

(Most commonly via least squares or logistic regression)

(Chapter 11 in ESL is a good introductory reference for neural networks)

Background

BACKGROUND

Neural networks have come about in 3 “waves” of research

The first was an attempt to model the mechanics of the human brain

Through psychological and anatomical experimentation, it appeared the brain worked by

- taking atomic units known as **neurons**, which can either be “on” or “off”
- putting them in **networks** with each other, where the **signal** is given by which neurons are “on” at a given time

Crucially, a neuron itself interprets the status of other neurons

BACKGROUND

After the development of parallel, distributed computation in the 1980s, the artificial intelligence view was diminished

Neural networks became popular machine learning approaches

Among the growing popularity of SVMs and boosting/bagging in the late 1990s, neural networks again fell out of favor

This was due to many of the problems we'll discuss (non convexity being the main one)

BACKGROUND

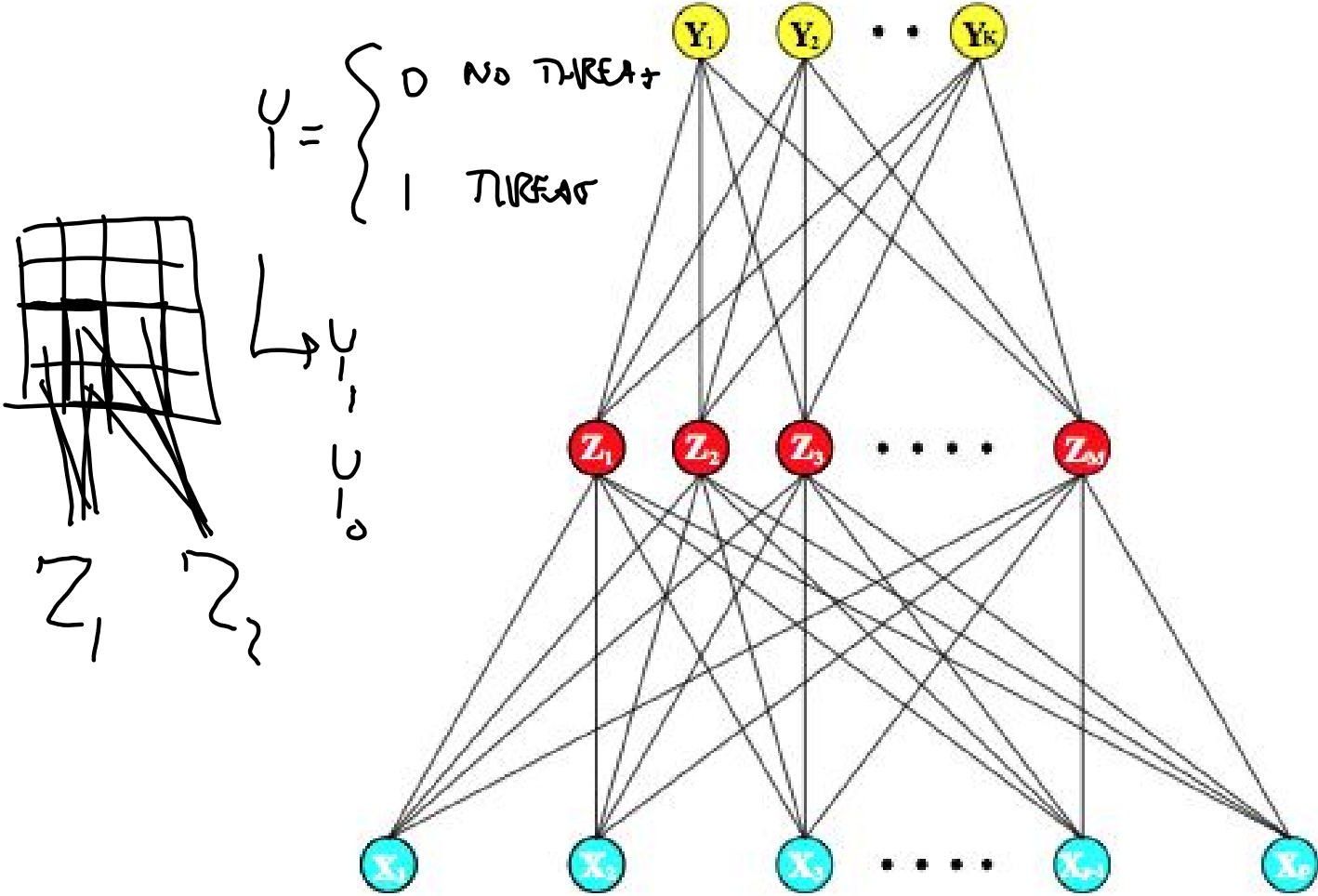
In the mid 2000's, new approaches for **initializing** neural networks became available

These approaches are collectively known as **deep learning**

Together, some state-of-the-art performance on various classification tasks have been accomplished via neural networks

High level overview

HIGH LEVEL OVERVIEW



$$Y | X \sim \mathcal{N}(X^T \beta, \sigma^2)$$

$$\beta \sim \mathcal{N}(0, \Sigma)$$

$$\Sigma \sim \text{IP}$$

FIGURE: Single hidden layer neural network. Note the similarity to latent factor models

NONPARAMETRIC REGRESSION

Suppose $Y \in \mathbb{R}$ and we are trying to nonparametrically fit the regression function

$$\mathbb{E}Y|X = f_*(X) \stackrel{?}{=} \beta^\top X$$

A common approach (particularly when p is small) is to specify

- A **fixed basis**, $(\phi_k)_{k=1}^\infty$
- A tuning parameter K

NONPARAMETRIC REGRESSION

We follow this prescription:

1. Write¹

$$f_*(X) = \sum_{k=1}^{\infty} \beta_k \phi_k(X) \quad \text{"="} \quad \sum \beta_k \chi_k$$

where $\beta_k = \langle f_*, \phi_k \rangle$

f_* HAS ω DERIVATIVES
2. Truncate this expansion² at K

$$|\beta_k| \sim k^{-\omega}$$

$$f_*^K(X) = \sum_{k=1}^K \beta_k \phi_k(X)$$

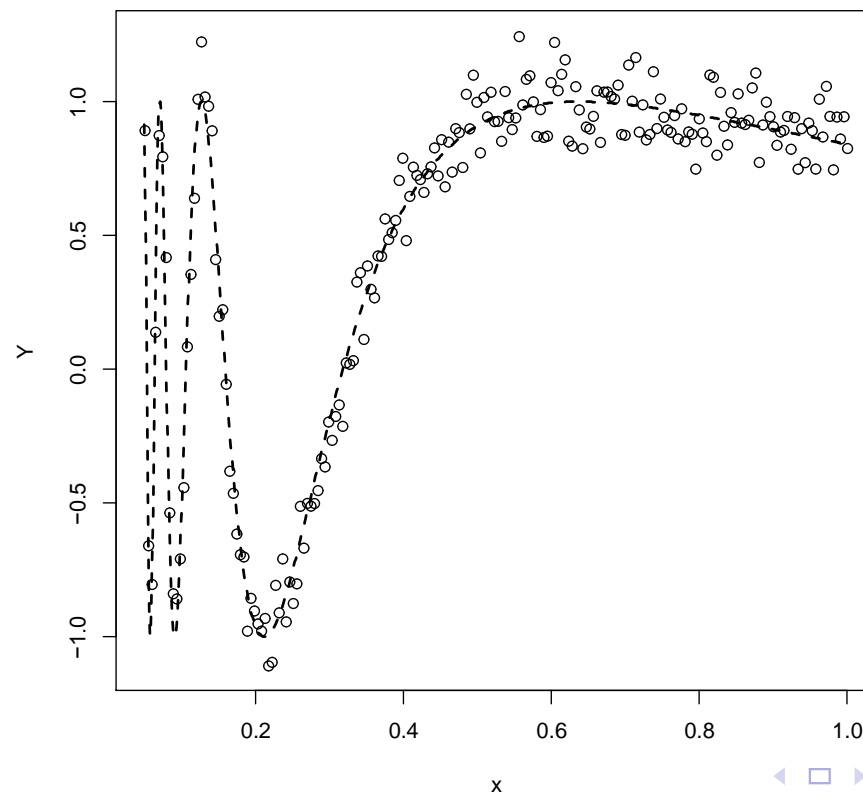
3. Estimate β_k with least squares

¹Technically, f_* might not be in the span of the basis, in which case we have incurred an irreducible approximation error. Here, I'll just write f_* as the projection of f_* onto that span

²Often higher k are more **rough** \Rightarrow this is a smoothness assumption \equiv

NONPARAMETRIC REGRESSION: EXAMPLE

```
x = seq(.05,1,length=200)
Y = sin(1/x) + rnorm(100,0,.1)
plot(x,Y)
xTest = seq(.05,1,length=1000)
lines(xTest,sin(1/xTest),col='black',lwd=2,lty=2)
```



"DOPPLER"

NONPARAMETRIC REGRESSION: EXAMPLE

require(splines)

X = bs(x,df=20)

plot(x,Y)

lines(xTest,sin(1/xTest),col='black',lwd=2,lty=2)

matlines(x=x,X,lty=2,type='l',col='blue')

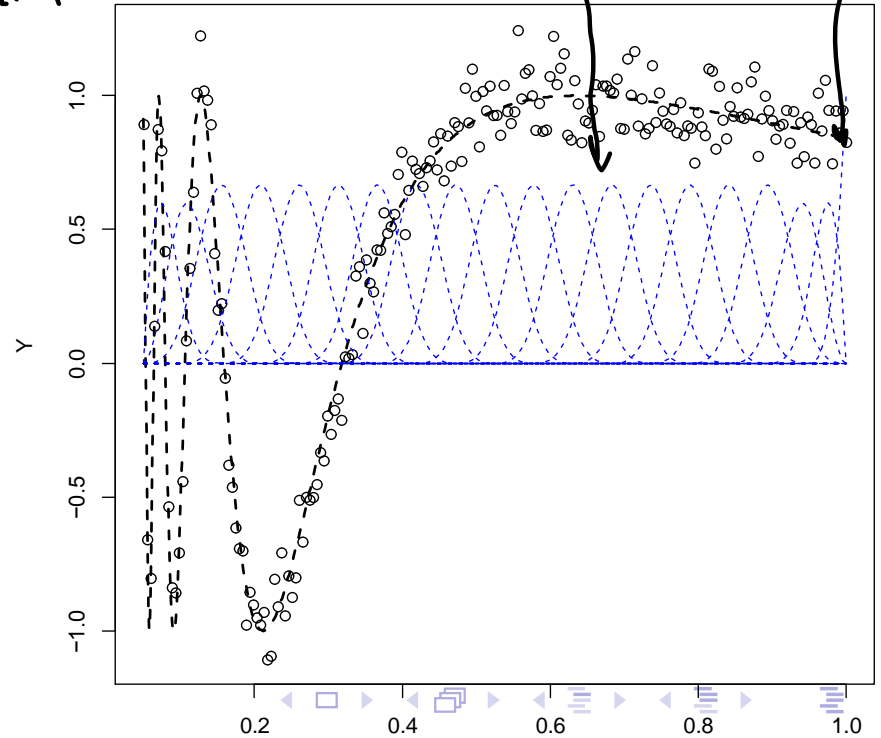
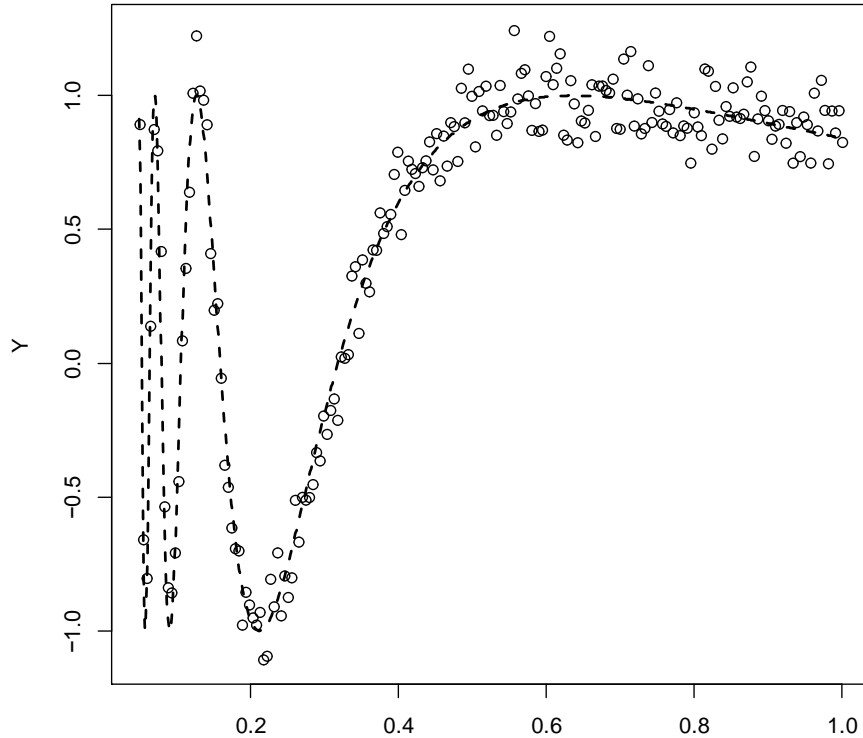
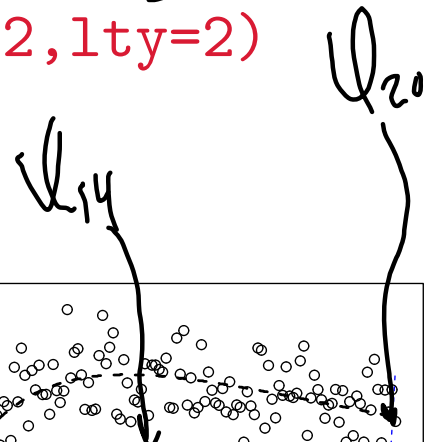
SPLINES

$$\Phi = \begin{bmatrix} \psi_1(x_1) & \dots & \psi_{20}(x_1) \\ \vdots & & \vdots \\ \psi_1(x_n) & \dots & \psi_{20}(x_n) \end{bmatrix}$$

$$(\Phi^T \Phi)^{-1} \Phi^T Y \quad \text{smoothing}$$

$$(\Phi^T \Phi + \lambda I)^{-1} \Phi^T Y$$

$n \times n$



NONPARAMETRIC REGRESSION: EXAMPLE

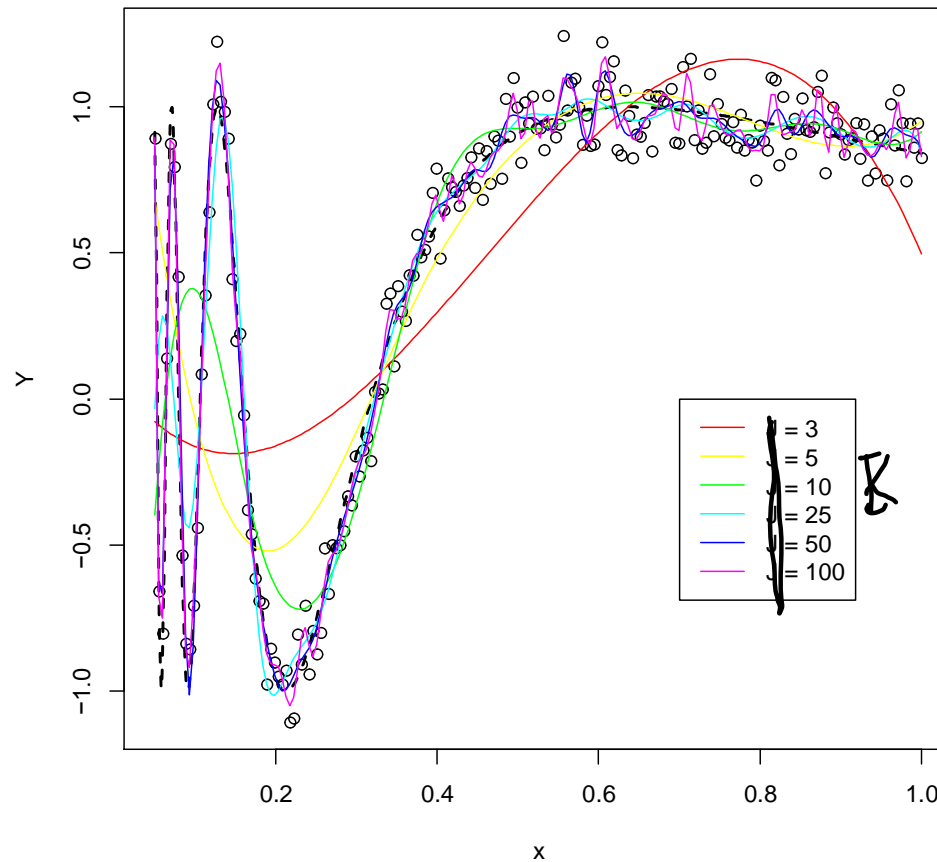
```
require(splines)
```

```
X = bs(x,df=K)
```

```
Yhat = predict(lm(Y~.,data=X))
```

BIAS

$$\sum_{k=1}^{\infty} |\beta_k|^2$$



NONPARAMETRIC REGRESSION

The weaknesses of this approach are:

- The basis is fixed and independent of the data
- If p is large, then nonparametrics doesn't work well at all
(See previous discussion on curse of dimensionality)
- If the basis doesn't 'agree' with f_* , then K will have to be large to capture the structure
- What if parts of f_* have substantially different structure?



$$\sum_{k=1}^K \langle f_*, \psi_k \rangle \psi_k$$

An alternative would be to have the data **tell** us what kind of basis to use

HIGH LEVEL OVERVIEW

Let $\mu(X) = \mathbb{E}Y|X$

Write L as the **link function**

A basic³ neural network can be phrased

$$L(\mu(X)) = \beta_0 + \sum_{k=1}^K \beta_k \sigma(\alpha_{k0} + \alpha_k^\top X)$$

COMPARE: A nonparametric GLM would have the form

$$L(\mu(X)) = \beta_0 + \sum_{k=1}^K \beta_k \phi_k(X)$$

³Here basic indicates that there are much more complex versions, not that neural networks are simple in any way

NEURAL NETWORKS: DEFINITIONS



$$L(\mu(X)) = \beta_0 + \sum_{k=1}^K \beta_k \sigma(\alpha_{k0} + \alpha_k^T X)$$

The main components are $\alpha_1 = (\alpha_{11}, \alpha_{12}, \alpha_{13}, 0, \dots, 0)^T$, $z_1 = \sigma(\alpha_{10} + \alpha_1^T X)$

$\hookrightarrow \sum_{j=1}^p \alpha_{1j} x_j = \alpha_{11} x_1 + \alpha_{12} x_2 + \alpha_{13} x_3$

- The derived features $Z_k = \sigma(\alpha_{k0} + \alpha_k^T X)$ and are called the **hidden units**
 - ▶ The function σ is called the **activation function** and is very often $\sigma(u) = (1 + e^{-u})^{-1}$
(This particular $\sigma(u)$ is known as the **sigmoid function**)
 - ▶ The parameters $\beta_0, \beta_k, \alpha_{k0}, \alpha_k$ are estimated from the data.
- The number of hidden units K is a tuning parameter

HIGH LEVEL OVERVIEW

EXAMPLE: If $L(\mu) = \mu$, then we are doing regression:

$$\mu(X) = \beta_0 + \sum_{k=1}^K \beta_k \sigma \left(\alpha_{k0} + \sum_{j=1}^p \alpha_{kj} X_j \right)$$

but in a **transformed space**

TWO OBSERVATIONS:

- The σ function generates a **feature map**
- If $\sigma(u) = u$, then neural networks reduce to classic least squares

Let's discuss each of these..

OBSERVATION 1: FEATURE MAP

We start with p **covariates**

We generate K **features**

EXAMPLE: GLMs with a **feature** transformation

$$\begin{aligned}\Phi(X) &= (1, x_1, x_2, \dots, x_p, x_1^2, x_2^2, \dots, x_p^2, x_1x_2, \dots, x_{p-1}x_p) \in \mathbb{R}^K \\ &= (\phi_1(X), \dots, \phi_K(X))\end{aligned}$$

Before feature map:

$$L(\mu(X)) = \beta_0 + \sum_{j=1}^p \beta_j x_j$$

After feature map:

$$L(\mu(X)) = \beta^\top \Phi(X) = \sum_{k=1}^K \beta_k \phi_k(X)$$

$\nearrow z_k$

OBSERVATION 1: FEATURE MAP

For **neural networks** write:

$$Z_k = \sigma \left(\alpha_{k0} + \sum_{j=1}^p \alpha_{kj} x_j \right) = \sigma \left(\alpha_{k0} + \alpha_k^\top X \right)$$

Then we have

$$\Phi(X) = (1, Z_1, \dots, Z_K)^\top \in \mathbb{R}^{K+1}$$

and

$$\mu(X) = \beta^\top \downarrow \Phi(X) = \beta_0 + \sum_{k=1}^K \beta_k \sigma \left(\alpha_{k0} + \sum_{j=1}^p \alpha_{kj} x_j \right)$$

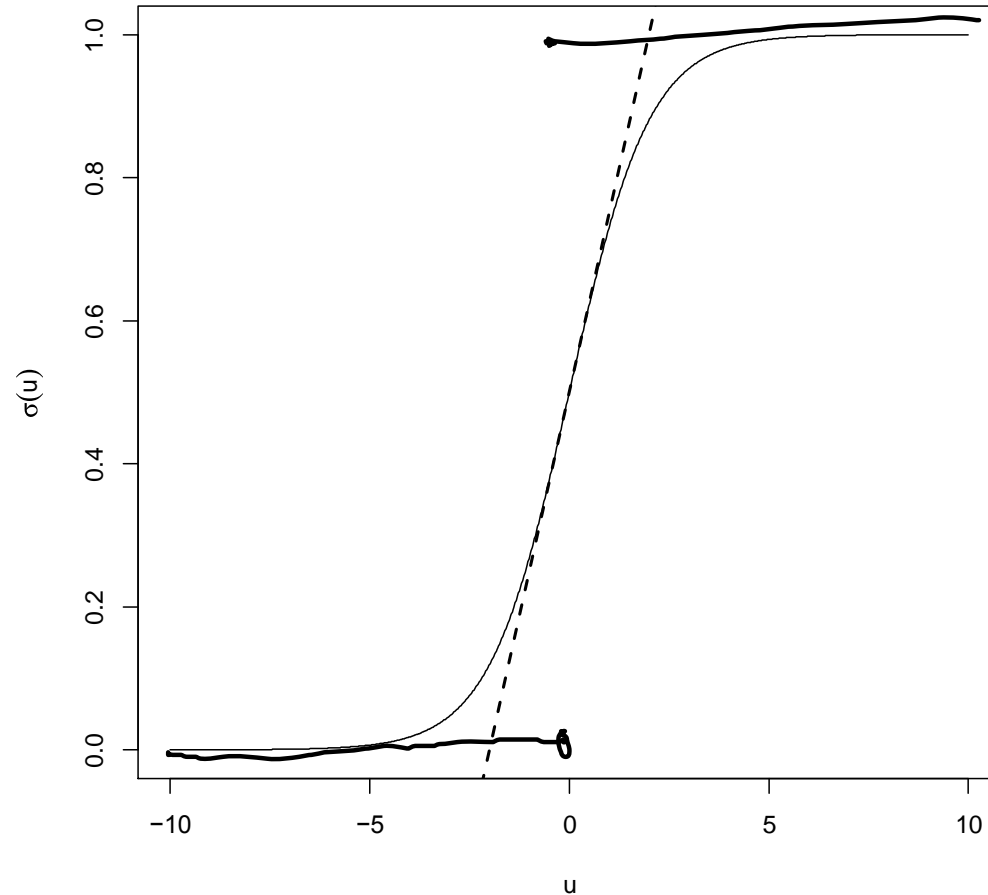
OBSERVATION 2: ACTIVATION FUNCTION

If $\sigma(u) = u$ is linear, then we recover classical methods

$$\begin{aligned}L(\mu(X)) &= \beta_0 + \sum_{k=1}^K \beta_k \sigma(\alpha_{k0} + \alpha_k^\top X) \\&= \beta_0 + \sum_{k=1}^K \beta_k (\alpha_{k0} + \alpha_k^\top X) \\&= \underbrace{\beta_0 + \sum_{k=1}^K \beta_k \alpha_{k0}}_{\text{red}} + \underbrace{\sum_{k=1}^K \beta_k \alpha_k^\top X}_{\text{blue}} \\&= \gamma_0 + \gamma^\top X \\&= \gamma_0 + \sum_{j=1}^p \gamma_j^\top x_j\end{aligned}$$

OBSERVATION 2: ACTIVATION FUNCTION

Plot of **sigmoid** activation function



If we look at a plot of the sigmoid function, it is quite linear near 0, but has nonlinear behavior further from the origin

HIERARCHICAL MODEL

A neural network can be phrased as a hierarchical model

$$Z_k = \sigma(\alpha_{k0} + \alpha_k^\top X) \quad (k=1, \dots, K)$$

$$W_g = \beta_{g0} + \beta_g^\top Z \quad (g=1, \dots, G)$$

$$\mu_g(X) = L^{-1}(W_g)$$

The output depends on the application, where we map W_g to the appropriate space:

- **REGRESSION:** The link function is $L(u) = u$
(here, $G = 1$)
- **CLASSIFICATION:** With G classes, we are modeling $\pi_g = \mathbb{P}(Y = g|X)$ and $L = \text{logit}$:

$$\hat{\pi}_g(X) = \frac{e^{W_g}}{\sum_{g'=1}^G e^{W_{g'}}} \quad \text{and} \quad \hat{Y} = \hat{g}(X) = \arg \max_g \hat{\pi}_g(X)$$

(This is called the **softmax** function for historical reasons)

TRAINING NEURAL NETWORKS

Neural networks have **many** (**MANY**) unknown parameters

(They are usually called **weights** in this context)

These are

- α_{k0}, α_k for $k = 1, \dots, K$ (total of $K(p + 1)$ parameters)
- β_{g0}, β_g for $g = 1, \dots, G$ (total of $G(K + 1)$ parameters)

TOTAL PARAMETERS: $\asymp Kp + GK = K(p + G)$

TRAINING NEURAL NETWORKS

The most common loss functions are

- REGRESSION:

$$\hat{R} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

- CLASSIFICATION: Cross-entropy

$$\hat{R} = - \sum_{i=1}^n \sum_{g=1}^G Y_{ig} \log(\hat{\pi}_g(X_i))$$

- ▶ Here, Y_{ig} is an indicator variable for the g^{th} class. In other words $Y_i \in \mathbb{R}^G$
(In fact, this means that Neural networks very seamlessly incorporate the idea of having multivariate response variables, even in regression)
- ▶ With the softmax + cross-entropy, neural networks is a linear multinomial logistic regression model in the hidden units

TRAINING NEURAL NETWORKS

The usual approach to minimizing \hat{R} is via **gradient descent**

This is known as **back propagation**

Due to the hierarchical form, derivatives can be formed using the chain rule and then computed via a forward and backward sweep

TRAINING NEURAL NETWORKS

We'll need some derivatives to implement the **gradient descent**

$$\mu(X) = \beta_0 + \sum_{k=1}^K \beta_k \sigma \left(\alpha_{k0} + \sum_{j=1}^p \alpha_{kj} x_j \right)$$

Derivatives:

$$\frac{\partial \mu}{\partial \beta_k} = \sigma(\alpha_{k0} + \alpha_k^\top X) = Z_k$$

$$\frac{\partial \mu}{\partial \alpha_{kj}} = \beta_k \sigma'(\alpha_{k0} + \alpha_k^\top X) x_j$$

NEURAL NETWORKS: BACK-PROPAGATION

For squared error, let $\hat{R}_i = (Y_i - \hat{Y}_i)^2$

Then

$$\frac{\partial \hat{R}_i}{\partial \beta_k} = -2(Y_i - \hat{Y}_i)Z_{ik}$$

$$\frac{\partial \hat{R}_i}{\partial \alpha_{kj}} = -2(Y_i - \hat{Y}_i)\beta_k \sigma'(\alpha_{k0} + \alpha_k^\top X_i)X_{ij}$$

Given these derivatives, a gradient descent update can be found

$$\hat{\beta}_k^{t+1} = \hat{\beta}_k^t - \gamma_t \sum_{i=1}^n \frac{\partial \hat{R}_i}{\partial \beta_k} \Bigg|_{\hat{\beta}_k^t}$$

$$\hat{\alpha}_{kj}^{t+1} = \hat{\alpha}_{kj}^t - \gamma_t \sum_{i=1}^n \frac{\partial \hat{R}_i}{\partial \alpha_{kj}} \Bigg|_{\hat{\alpha}_{kj}^t}$$

(γ_t is called the **learning rate**, this needs to be set)

NEURAL NETWORKS: BACK-PROPAGATION

Returning to

$$\frac{\partial \hat{R}_i}{\partial \beta_k} = -2(Y_i - \hat{Y}_i)Z_{ik} = a_i Z_{ik}$$

$$\frac{\partial \hat{R}_i}{\partial \alpha_{kj}} = -2(Y_i - \hat{Y}_i)\beta_k \sigma'(\alpha_{k0} + \alpha_k^\top X_i) X_{ij} = b_{ki} X_{ij}$$

Direct substitution of a_i into b_{ki} gives

$$b_{ki} = a_i \beta_k \sigma'(\alpha_{k0} + \alpha_k^\top X_i)$$

These are the **back-propagation equations**

NEURAL NETWORKS: BACK-PROPAGATION

ADVANTAGES:

- It's updates only depend on **local** information in the sense that if objects in the hierarchical model are unrelated to each other, the updates aren't affected
(This helps in many ways, most notably in parallel architectures)
- It doesn't require second-derivative information
- As the updates are only in terms of \hat{R}_i , the algorithm can be run in either **batch** or **online** mode

DOWN SIDES:

- It can be very slow
- Need to choose the **learning rate** γ_t

NEURAL NETWORKS: OTHER ALGORITHMS

There are a few alternative variations on the fitting algorithm

Many are using more general versions of non-Hessian dependent optimization algorithms

(For example: conjugate gradient)

The most popular are

- RESILIENT BACK-PROPAGATION
- MODIFIED GLOBALLY CONVERGENT VERSION

REGULARIZING NEURAL NETWORKS

As usual, we don't actually want the global minimizer of the training error (particularly since there are so many parameters)

Instead, some regularization is included, with some combination of:

- a complexity penalization term
- early stopping on the back propagation algorithm used for fitting

REGULARIZING NEURAL NETWORKS

Explicit regularization comes in a couple of flavors

- **WEIGHT DECAY:** This is like ridge regression in that we penalize the squared Euclidean norm of the weights

$$\rho(\alpha, \beta) = \sum \beta^2 + \sum \alpha^2$$

- **WEIGHT ELIMINATION:** This encourages more shrinking of small weights

$$\rho(\alpha, \beta) = \sum \frac{\beta^2}{1 + \beta^2} + \sum \frac{\alpha^2}{1 + \alpha^2}$$

NOTE: In either case, we now solve:

$$\min \hat{R} + \lambda \rho(\alpha, \beta)$$

This can be done efficiently by augmenting the gradient descent derivatives

COMMON PITFALLS

There are three areas to watch out for

- **NONCONVEXITY:** The neural network optimization problem is non convex. This makes any numerical solution highly dependant on the initial values. These must be
 - ▶ chosen carefully
 - ▶ regenerated several times to check sensitivity
- **SCALING:** Be sure to standardize the covariates before training
- **NUMBER OF HIDDEN UNITS (K):** It is generally better to have too many hidden units than too few (regularization can eliminate some).

(Later, we will see this can include adding multiple hidden layers)

STARTING VALUES

The quality of the neural network predictions is very dependent on the starting values

As noted, the sigmoid function is nearly linear near the origin.

Hence, starting values for the weights are generally randomly chosen near 0. Care must be chosen as:

- Weights equal to 0 will encode a symmetry that keeps the back propagation algorithm from changing solutions
- Weights that are large tend to produce bad solutions (overfitting)

This is like putting a prior on linearity and demanding the data add any nonlinearity

STARTING VALUES

Once several starting values + back-propagation pairs are run, we must sift through the output

Some common choices are:

- Choose the solution that minimizes **training error**
- Choose the solution that minimizes the **penalized** training error
- Average the solutions across runs

(This is the recommended approach as it brings a model averaging/Bayesian flair)