

# Machine Learning II

Prof. Dr. Tim Downie

Lecture 8 – 24th November 2023

Projection Pursuit Regression

# Contents

- ▶ Adaptation from a GAM to the projection pursuit regression model
- ▶ With one ridge function
- ▶ Multiple ridge functions
- ▶ Similarity between PPR and ... ?

## Text books:

- ▶ Hastie, Tibshirani and Friedman (2009), *The Elements of Statistical Learning (2nd edition)*. Springer-Verlag.  
<https://hastie.su.domains/ElemStatLearn><sup>1</sup>
- ▶ W.N. Venables and B.D. Ripley (2002). *Modern Applied Statistics with S* (fourth edition, Springer). *Modern Data Science with R*. Chapman and Hall.  
<https://www.stats.ox.ac.uk/pub/MASS4/>

---

<sup>1</sup>Online version updated in 2017

## Generalised additive Models

In supervised learning, given an input vector  $X = (X_1, X_2, \dots, X_p)^\top$  and outcome variable  $Y$  (Target), the aim is to find a predictor function  $f(X)$  which predicts  $Y$  well.

A multiple linear regression model has the form:

$$Y \approx f(X) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p.$$

A non-linear equivalent is the generalised additive model (GAM: ML2 Lec 3).

$$Y \approx f(X) = \beta_0 + g_1(x_1) + g_2(x_2) + \dots + g_p(x_p).$$

In a GAM, each  $g_m$  is allowed to be a different type of regression function on one variable. E.g.  $g_1$  could be a quadratic function,  $g_2$  a smoothing spline function and  $g_p$  a loess function.

## Projection pursuit regression model

The GAM Model restricts each function  $g_m$  to “respond to” just one variable, corresponding to one dimension of the  $X$  vector. To get more flexibility we can apply  $g_m$  to a linear combination of  $X_1, X_2, \dots, X_p$ .

Let  $V_1$  be any linear combination of the elements of  $X$ . We can define  $V_1$  using the scalar product with the *direction vector*  $\omega_1$ :

$$V_1 = \omega_1^\top X \tag{1}$$

A direction vector is normalised to be a unit vector, i.e.  $\|\omega_1\| = 1$ , which ensures the expression in Equation 1 is unique up to the sign.

The data  $X_1, X_2, \dots, X_p$  are projected on to a line in  $p$ -dimensional space along the direction vector  $\omega_1$ . The value of  $V_1$  is then the signed distance of the projected point from the origin.

Graphical example of projection and  $\omega$  direction vectors in two dimensions  
 $p = 2$ . Diagram on blackboard

We can fit the model

$$Y \approx f(X) = \beta_0 + \beta_1 g_1(V_1), \quad \text{where } V_1 = \omega_1^\top X.$$

The function  $g_1$  is called a ridge function, it “responds to” the vector  $X$  in the direction of  $\omega_1$ , which graphically looks like a ridge.

For any observed element or prediction,  $V_1$  is a scalar value, so  $g_1$  can be any univariate line/curve fitting method. Popular choices are loess and spline smoothing methods, as with GAMs.

The projection pursuit regression (PPR) model, allows for multiple projections and multiple smoothing functions, specified by  $M$

$$Y \approx f(X) = \beta_0 + \beta_1 g_1(V_1) + \cdots + \beta_M g_M(V_M),$$

where  $V_m = \omega_m^\top X$ ,  $m = 1, \dots, M$ .

Unlike the principal components approach, there is **no requirement** that  $V_1, \dots, V_M$  are all mutually orthogonal.

Appropriate combinations of direction vectors and ridge functions allow for a large range of predictor functions to be modelled including interactions and higher order polynomials without needing to specify the required form in advance, see Hastie et al. for the mathematical details.



The algorithm used to fit a PPR models closely resembles the back fitting algorithm for GAMs.

- 1) Calculate the constant parameter  $\hat{\beta}_0 = \bar{y}$ .
- 2) Initialise  $M = 1$  and  $\omega_1$ . A common initial direction vector is  $\omega_1^\top = \frac{1}{\sqrt{p}}(1, 1, \dots, 1)$ .
- 3) Iteratively repeat steps a) to d) using  $m = 1, \dots, M$ .
  - a) Obtain the  $m$ -th partial residuals  $\eta_{i,m}$ .
  - b) Calculate  $V_{i,m} = \omega_m^\top \eta_{i,m}$ , which are the  $m$ -th partial residuals projected on to direction  $\omega_m$ .
  - c) Obtain the  $m$ -th ridge function  $g_m(V_{\cdot,m})$ .
  - d) Update  $\beta_m$  and  $\omega_m$  to minimise the sum of squared errors

$$SSE = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{m=1}^M \beta_m g_m(\omega_m^\top \mathbf{x}_i) \right)^2$$

- 4) Increase  $M$  by one, to add a new ridge function, and initialise  $\omega_M$ .  
Repeat steps 3) a)-d) above

## Comments on the algorithm

- ▶ Reminder the  $m^{\text{th}}$ -partial residuals are defined as:

$$\eta_{im} = y_i - (\beta_0 + \sum_{k \neq m} \beta_k g_k(\omega_k^\top \mathbf{x}_i)).$$

The fitted values of all the other components are subtracted from the observed outcome values.

- ▶ Step 3 c) depends on the chosen regression/smoothing method covered in various ML1/ML2 lectures.
- ▶ The precise ordering of: i) improve  $g_m$ , ii) improve  $\beta_m, \omega_m$  and iii) increase  $M$  is not critical to the algorithm results and so we will not consider this in detail.
- ▶ The minimisation in Step 3 d) is done using Gauß-Newton optimisation. An important algorithm for function minimisation, which we will investigate in more detail in the coming weeks.

## Example

In ML1 Lecture 8, we fitted a regression tree model to a data set called `cpus` from the MASS library. A relative performance measure and characteristics of 209 CPUs from 1987.

The data include the variables:

`name`: manufacturer and model.

`syct`: cycle time in nanoseconds.

`mmin`: minimum main memory in kilobytes.

`mmax`: maximum main memory in kilobytes.

`cach`: cache size in kilobytes.

`chmin`: minimum number of channels.

`chmax`: maximum number of channels.

`perf`: published performance on a benchmark mix relative to an IBM-370/158-3.

As `perf` is a heavily skewed variable we will use  $y = \log_{10}(\text{perf})$  as our outcome variable.

A PPR model is fitted on training data and the MSE compared to a linear model and a RPART tree model using test data.

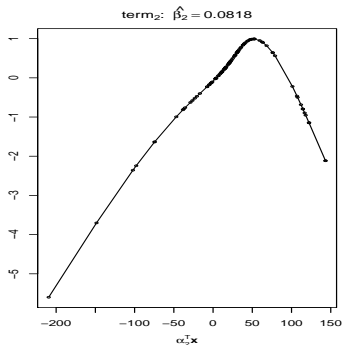
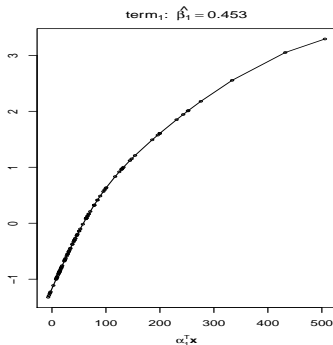
```
cpus.ppr <- ppr(log10(perf) ~ syct+mmin+mmax+cach+chmin+chmax,  
               data=cpus[cpus.train,], nterms=2, sm="spline")
```

- ▶ The function called `ppr` is part of the MASS library.
- ▶ The model formula specifies the variables to include, which uses similar notation to `lm` and `rpart`.
- ▶ The model will be fitted with two components  $M = 2$ : `nterms=2`.
- ▶ The type of non-linear regression to be used is spline smoothing, specified by `sm="spline"`

The direction vectors are

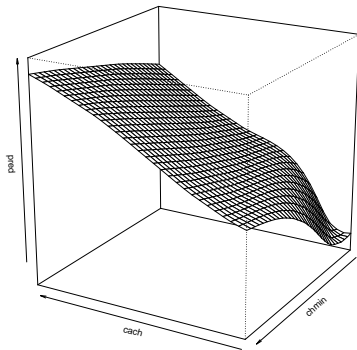
	term 1	term 2
syct	-0.013314130	0.128598978
mmin	0.006151632	-0.004512696
mmax	0.002680398	0.002356409
cach	0.939999246	-0.894044023
chmin	0.322811306	0.210153382
chmax	0.109416569	0.374108577

The two fitted functions are



The two largest components in  $\omega_1$  correspond to the variables `cach` and `chmin`.

A plot of the first function  $g_1$  against these two variables demonstrates the ridge effect.



## Comparison of models

Using the test data, the following MSE values were obtained

Multiple linear regression	0.0460
Rpart tree	0.0565
PPR	0.0290

The PPR model gives in this one example substantially better prediction on the test data set.

The disadvantage is that it is not easy to interpret the results, unlike with a linear model or tree model.

## Similarity between PPR and artificial neural networks

A simple artificial neural network (NN) with one hidden layer is often portrayed graphically in the following form.

Diagram on blackboard



The input nodes correspond to the  $p$  elements of one data vector  $\mathbf{X}$ .

The  $m$ -th node in the hidden layer takes a linear combination of  $\mathbf{X}$ :

$$Z_m = \omega_m^\top \mathbf{X}_i.$$

The output of each node is linear combination  $Z_m$  after transformation of an activation function:  $g_m(Z_m)$

The output layer takes the activations applies another linear combination and outputs the results.

$$\hat{y} = \phi_0 + \phi_1 g_1(Z_1) + \cdots + \phi_M g_M(Z_M).$$

The PPR model is a specific type of NN model.

The main difference is that with a NN the activation functions are specified in advance. With PPR functions are fitted to the combined data.