

# 7750: Mathematical Foundations of Machine Learning

Instructor: Ashwin Pananjady (ashwinpm@gatech.edu)

Core problem in supervised learning

$$\begin{array}{ccc} & & \swarrow \text{"noise"} \\ Y = & f_0(X) + \epsilon & \\ \uparrow & \uparrow & \\ \text{response} & \text{covariates/predictors/features} & \\ \in \mathbb{R} & \in \mathbb{R}^d & \end{array}$$

E.g.  
 $X$ : pixels of image  
 $Y$ : label.

Take samples/data points  $\{x_i, y_i\}_{i=1}^n$ .

$y_i = f_0(x_i) + \epsilon_i$ , modeled as i.i.d draws of training data.

Goal: ① Fit a function  $f$  to training data  
② Predict on new sample  $\tilde{x}$  with  $f(\tilde{x})$ .

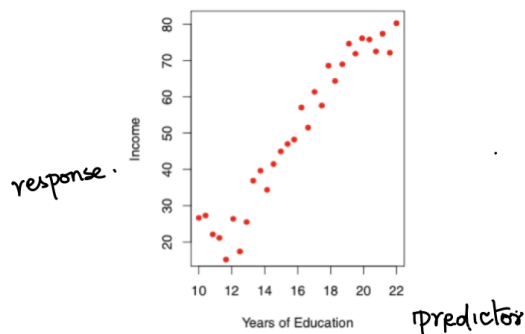
## Modeling / representation

What are reasonable assumptions on  $f_0$ ? Posit  
model class  $\mathcal{F}$ .

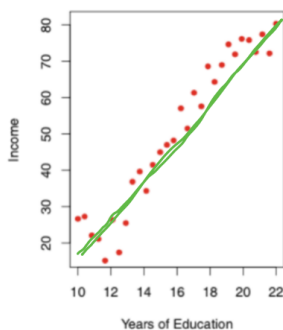
↳ Physical constraints

↳ Domain expertise

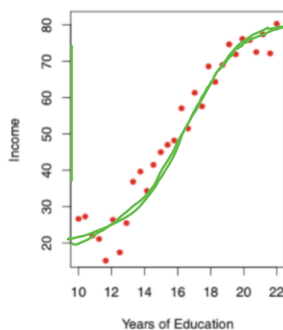
ISL.



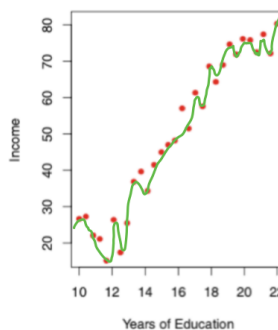
$\mathcal{F}$ : linear



$\mathcal{F}$ : cubic



$\mathcal{F}$ : polynomial high degree.



$\mathcal{F}$ .

simple

complex.

Today's focus: Properties of linear representations

$$f(x) = \sum_k \alpha_k h_k(x).$$

Example: Linear regression

$$\mathcal{F} = \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R} \mid f(x) = w^T x + b, \right. \\ \left. w \in \mathbb{R}^d, b \in \mathbb{R} \right\}.$$

$w_j$  measures "rate of change" in response when feature  $j$  is varied.

### Example: Polynomial regression

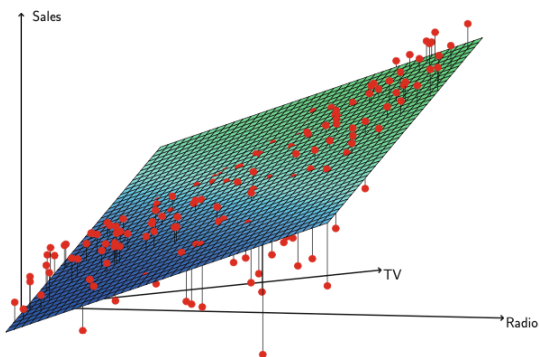
Let  $P_1(x), \dots, P_M(x)$  denote all monomials of degree  $\leq l$ , e.g.,  $1, x_1, \dots, x_d, x_1^2, x_2^2, \dots, x_1 x_2, \dots$

$$\mathcal{F} = \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R} \mid f(x) = \sum_{j=1}^M \alpha_j \cdot P_j(x), \alpha_j \in \mathbb{R} \ \forall j=1, \dots, M \right\}.$$

Ex! How big is  $M$  as a function of  $(d, l)$ ?  $\left[ \binom{d+l-1}{l}, \text{ will prove in HW!} \right]$

Note: These are both linear representations, in terms of basis functions.

Useful to think geometrically whenever these are introduced. (More on this in HW).



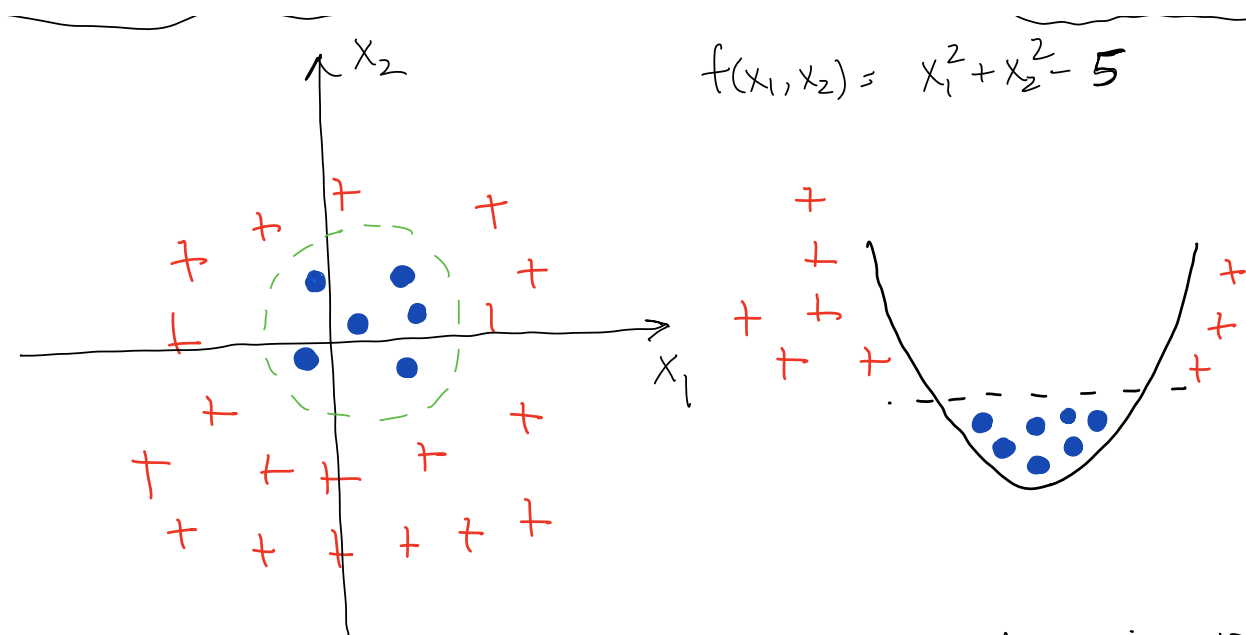
Advertising dataset, ISL.

Linear regression fits a hyperplane to data.

E.g.  $y = \text{"Sales"}$

$x_1 = \text{TV}$

$x_2 = \text{Radio.}$



More examples of polynomial approximation in 1D

Example 1 : Taylor Series with monomial basis fns.

Analytic functions on an open set  $\mathcal{D}$  can be described by Taylor series  $f(x) = \sum_{j=0}^{\infty} \alpha_j (x-x_0)^j$  for all  $x \in \mathcal{D}$ .

$$e^x = \sum_{j=0}^{\infty} \frac{1}{j!} x^j.$$

$$\log(1+x) = \sum_{j=0}^{\infty} \frac{(-1)^{j+1}}{j} x^j.$$

$$\sin x = \sum_{j=0}^{\infty} \frac{(-1)^j}{(2j+1)!} x^{2j+1}$$

Taylor's theorem :

Let  $f: \mathbb{R} \rightarrow \mathbb{R}$  be  $k$ -times differentiable at  $a \in \mathbb{R}$

Then there exists  $h_k: \mathbb{R} \rightarrow \mathbb{R}$  s.t.

$$\lim_{x \rightarrow a} h_k(x) = 0 \quad \text{and.}$$

$$f(x) = \sum_{j=0}^k \frac{f^{(j)}(a)}{j!} (x-a)^j + h_k(x) (x-a)^k.$$

Example 2 (Aside): Interpolating polynomials

Given  $\{x_i, y_i\}_{i=1}^n$ , is there a unique  
( $n-1$ )-degree polynomial that interpolates?

$$f(x) = \sum_{k=1}^n y_k \cdot \prod_{\substack{j \neq k \\ 1 \leq j \leq n}} \frac{(x - x_j)}{(x_k - x_j)}$$

Lagrange's theorem:  $f$  is unique interpolation by degree,  
( $n-1$ ) polynomial if  $\{x_i\}_{i=1}^n$  distinct

HW1 will walk you through proof.

### Example 3: Polynomial Splines

Key idea: Interpolate with different polynomials between data points, maintaining smoothness / stability

$$f(x_k) = y_k, \quad k = 1, \dots, n.$$

$f(x)$  is piecewise  $l$ -th order polynomial with

"kinks" at  $x_1, \dots, x_n$ .

$f(x)$  has  $l-1$  continuous derivatives at  $\{x_i\}_{i=1}^n$ .

See figures in Prof. Romberg's notes

Ex: We have  $n$  distinct points  $(x_i, y_i)_{i=1}^n$  and want to interpolate with cubic spline. How many unknowns and how many constraints?

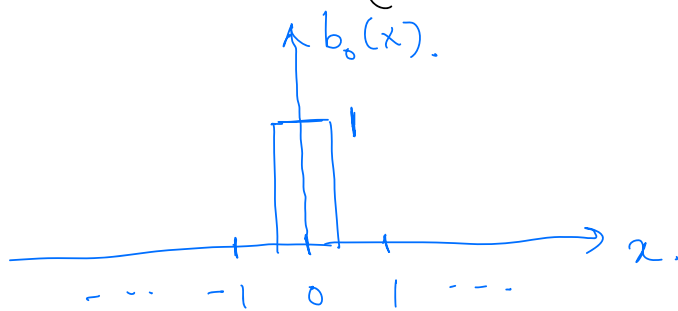
Can we think of splines through basis fn. expansions? Yes!

Basis functions are so-called "B-splines".

Q: Suppose we want to find basis functions for polynomial splines with  $l=0$ , where  $x_i$  are integers.

A: Use the basis functions  $\{b_0(x-k)\}_{k \in \mathbb{Z}}$ ,

where

$$b_0(x) = \begin{cases} 1 & , -1/2 \leq x < 1/2. \\ 0 & \text{otherwise} \end{cases}$$


Then such a spline interpolation is given

by 
$$f(x) = \sum_{k \in \mathbb{Z}} y_k \cdot b_0(x - x_k).$$

Two key observations:

→ The basis function is centered at your training points, this extends to different sampling patterns.

→ The training data does still influence the basis functions you choose, even if  $l=0$ .

Q: How to extend to  $l$ -th order spline?

$$\text{For } l=1, \quad b_1(x) = \begin{cases} x+1, & -1 \leq t \leq 0 \\ 1-x, & 0 \leq t \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Working with this basis will allow piecewise linear interpolation between equi-space data on integers.

$$\text{In general, } b_l(x) = \overbrace{b_0 * \dots * b_0}^{l \text{ times}}(x)$$

is basis for  $l$ -th order splines.

Example 4: Fourier series: We will not cover this in detail, but it is foundational material to know if you haven't been exposed to it. See Prof. Romberg's notes.



## Key takeaways

- Once we pick a basis (i.e. collection of basis functions), continuous-valued functions can be represented as a sequence of real values. This unlocks the linear algebra toolbox.

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \approx \begin{bmatrix} h_1(x_1) & \dots & h_M(x_1) \\ h_1(x_2) & \dots & h_M(x_2) \\ \vdots & & \vdots \\ h_1(x_n) & \dots & h_M(x_n) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_M \end{bmatrix}$$

- Picking a basis is nontrivial, depends on application
- Assessing in what fashion we want to approximate  $f$  is also important.
- Always a tradeoff: Large  $M \Rightarrow$  more complex