Lec 3.  Moving Beyond Linearity: Basis Expansion.

− Transformation of $\bar{X}$ to correct non linearity.

    − Core idea:  augment / replace $X$ with additional variables.

              then  use  linear  models  in the  new space  of  derived input  features.

    − Denote  $h_m(x): R^p \rightarrow R$ ,  $m = 1, \cdots M$.

$$f(x) = \sum_{m=1}^{M} \beta_m \, h_m(X) \quad , \qquad \text{a linear expansion of } X.$$

    $\cdot$ $h_m(x) = X_m$,  $m = 1, \cdots p$.  original  linear  model.

    $\cdot$ $h_m(x) = X_j^2$  or  $h_m(x) = X_j X_k$.  polynomial  (quadratic).

    $\cdot$ $h_m(x) = \log(X_j)$, $\sqrt{X_j}$.

    $\cdot$ $h_m(x) = I(L_m \leq X_k \leq U_m)$.  indicator.  $\mathcal{E}$ (step function).

~~Assume  X  is  one  dimensional.~~  piecewise  constant.
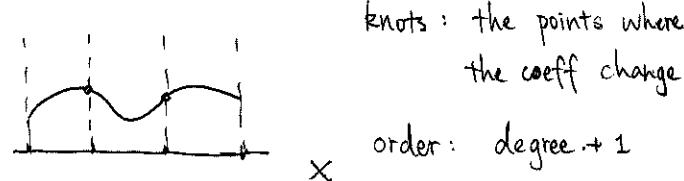
~~1. Polynomials.~~

    − Basis function:  a family of functions or transformations $\phi$ that can be applied to $X$

              e.g. Fourier series / wavelets.  regression splines.

Assume  $X$  is  1-dim.

1. Regression  splines.

 1.1 Piece-wise  polynomials.



knots: the points where
         the coeff change

order: degree + 1

    − $f(x)$.  divide $\mathcal{E}$ the domain  of  $X$  into  continuous  intervals. & represent $f$

      by a  separate polynomial in each  interval

    − These fixed knots splines $\overset{called}{\Rightarrow}$  regression  splines.

~~1.2  constraints & splines.~~

    − More  knots  → more flexible.  $K$ knots , → $K+1$  polynomials.

    − Face  discontinuous &  overfitting.

              (Show  Figure  7.3 $\mathcal{E}$  ISLR.  & 5.2.) ESL

## 1.2. Constraints.

- continuous.   [ i.e. no jump at knots ].

- First & second derivatives are continuous → smooth.
  [ i.e. a cubic splines or higher order ].   $M-2$ continuous derivatives.
  ↳ a cubic spline with $K$ knots uses $K+4$ degree of freedoms.

- boundary constraints.

  Splines can have high variance at the outer range of $x$ ~~Fig~~.  ( see Fig. 7.4  Confidence band wide).

  \* A natural spline is a regression spline with additional boundary constraints.
  the function need to be linear at the boundary.   [ may introduce bias ]

- Other types of splines :   B - splines, truncated power splines.

## 1.3  Choosing the number & location of the knots.

## 2.  Smoothing Splines.

- Put a knot at every data point, i.e.  knots = $\{x_1, x_2, \cdots, x_n\}$.   maximum set of knots

- The complexity of the fit is controlled by regularization.   i.e., among all $f(x)$. with 2 continous derivatives.  find. the one.

  ☆ $RSS(f, \lambda) = \min_{f} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \int (f''(t))^2 \, dt.$

  ($\uparrow$ fit the data.)  (smoothing)  regularization para. ($\uparrow$ penalize wiggliness)

  | | |
  |---|---|
  | $\lambda = 0$ | interpolating spline |
  | small | wiggly |
  | large | smooth. |
  | $\lambda \to \infty$ | straight line |

- It can be shown ☆ has an explicit, finite - dimensional, unique minimizer.
  ↳ a natural cubic spline with knots at the unique values of $x_i$.  $i = 1, \cdots, N$.

  $f(x) = \sum_{j=1}^{N} N_j(x) \theta_j$ .   $N_j(x)$ — an $N$-dim set of basis functions of natural spline

  ⇒ $RSS(\theta, \lambda) = (y - N\theta)^T (y - N\theta) + \lambda \theta^T \Omega_N \theta.$   $\{N\}_{ij} = N_j(x_i)$
  $\{\Omega_N\}_{jk} = \int N_j''(t) N_k''(t) \, dt.$

  ⇒ $\hat{\theta} = (N^T N + \lambda \Omega_N)^{-1} N^T y$ ,   a generalized ridge regression.

## 2.1. Degrees of Freedom & Smoother Matrices.

$$\hat{f} = N(N^T N + \lambda \Omega_N)^{-1} N^T y. \quad \Rightarrow \quad \text{A smoothing spline is a linear smoother.}$$

$$:= S_\lambda y.$$

↑ smoother matrix.

Define the effective degrees of freedom

$$df_\lambda = \text{trace}(S_\lambda).$$

Large → flexible model

small → rigid model

↓

related to $\lambda$.  how. (see. previous table)

---

\* Since $S_\lambda$. symmetric (& $^+$ semi-definite)

Rewrite $S_\lambda = N(N^T N)^{-1/2} [I + \lambda \underbrace{(N^T N)^{-\frac{1}{2}} \Omega_N (N^T N)^{-\frac{1}{2}}}_{:= A}]^{-1} (N^T N)^{-\frac{1}{2}} N^T$

$\underbrace{\phantom{N(N^T N)^{-1/2}}}_{U}$    ind of $\lambda$.  symmetric. positive semi definite.

$$= U^\circ [I + \lambda A]^{-1} U^T$$

By. spectrum decomposition.

$$A = V \, diag(v_1, \cdots v_N) V^T$$

$$VV^T = I$$

$$= (UV)^\circ \, diag\left(\frac{1}{1 + \lambda v_j}\right) (UV)^T$$

$$\Rightarrow df_\lambda = \sum_{j=1}^{N} \frac{1}{1 + \lambda v_j}$$   monotone ↓ of $\lambda$ ⇒

1−1 relationship.

no local minima. no reversals.

tuning stable. GCV. well-behaved.

---

## 2.2. Automatic selection of the smoothing para.

a. Fixing the Df.

b. Bias − Variance Tradeoff

$$Y = f(x) + \epsilon \qquad \hat{f} = S_\lambda y.$$

$$\epsilon \sim N(0, \sigma^2)$$

$$Cov(\hat{f}) = S_\lambda \, Cov(y) \, S_\lambda^T$$

$$= \sigma^2 \, S_\lambda S_\lambda^T$$

From Statistical Decisision theorey,

$$EPE(\hat{f}_\lambda) = E[Y - \hat{f}_\lambda(x)]^2$$

$$Bias(\hat{f}) = f - E(\hat{f}). \quad = f - S_\lambda f$$

$$= (I - S_\lambda) f.$$

$$= E[(Y - f) + (f - \hat{f}_\lambda(x))]^2$$

$$= \sigma^2 + E[(f - \hat{f})]^2 = E[f - E(\hat{f}(x)) + E(\hat{f}(x)) - \hat{f}(x)]^2.$$

$$= E[Bias^2(\hat{f}(x)) + Var(\hat{f}(x))] = MSE(\hat{f}_\lambda).$$

\*. Don't know the true function.   No acess to EPE. and  need an estimate.

$\Rightarrow$ K-fold CV.   Generalized CV (GCV)

(Multivariate Case).

# 3. Multidimensional Splines. & Generalized Additive Models.

Many vars.       $y = f(x_1, x_2, \cdots x_p) + \varepsilon.$       A full general $f$.
- very flexible.
- statistically impossible.

## 3.1. Multidimensional splines.

$X = (x_1, x_2).$       Tensor - product Splines

- use splines in each var.

basis   $g_{jk}(X) = h_{1j}(x_1) h_{2k}(x_2).$   $j = 1, \cdots, M_1$
- Capture interaction.

$k = 1, \cdots, M_2.$   - Explode in dim

$\Rightarrow g'(X) = \sum_{j=1}^{M_1} \sum_{k=1}^{M_2} \theta_{jk} g_{jk}(X).$

$M_1 M_2$ para.

When to use :   spatial data. ; surface ;   small $p$ , large $n$.


## 3.2. Generalized Additive Models.

$E(Y|X) = \beta_0 + f_1(x_1) + \cdots + f_p(x_p).$

$f_j$ - univariate smooth function.

Pros : · Allow non-linearity.

· Because of additive nature, can examine the effect of each $x_i$.
( interpretability )

· smoothness can be summarized by df.

$df_{total} = 1 + \sum_j df_j$

Cons :   Restricted to "Additive".

(But could add interaction term).


# 4 Moving to Dictionaries.

Assumption :   sparse. expansion.  $\rightarrow$ regularization

Fixed bases · small. predefined set of functions.

& sparsity.

Dictionary:   more flexible. by using large or overcomplete collection of functions.