

Introducing non-linearity.

$$Linear = LS, Ridge, \dots = \hat{f}_{lin}(\vec{x}) = \underbrace{\vec{x}^T}_{(1 \times p)} \cdot \underbrace{\vec{\beta}}_{(p \times 1)} \quad \vec{x}^T = (x_1, \dots, x_p)$$

$$Nonlinear \hat{=} \hat{f}_{nonlin}(\vec{x}) = \underbrace{\vec{h}_m(\vec{x})}_{(1 \times p)} \cdot \vec{\beta}_{nonlin}$$

Examples: $\{h_{m,1}, \dots, h_{m,T}\}$ from $(x_1^{d_1}, x_2^{d_2}, \dots, x_p^{d_p})$

Dictionary

$$\begin{pmatrix} p+D \\ D \end{pmatrix} \quad \sum d_i < D$$

Sparse is desirable.

Other: (i) $\log(x_i)$ ---

(ii) $\sin(x_i)$ $\cos(x_i)$ $\sin(kx_i)$

Statistical vs. Linear models

$$x_1^{(i)} \dots x_p^{(i)} \mid y^{(i)} \quad i = 1, \dots, N.$$

Statistical models
defined by

$$\vec{\beta}$$

$$\left\{ \begin{array}{l} \text{Linear model: } \hat{y} = \hat{f}_{lin}(\vec{x}_{in}) \quad \vec{x}_{in} = (x_1 \dots x_p) \\ \quad \quad \quad = \vec{x}_{in}^T \cdot \vec{\beta}^* \\ \text{Nonlinear model: } \hat{y} = \hat{f}_{nonlin}(\vec{x}_{in}) \\ \quad \quad \quad = h(\vec{x}_{in})^T \cdot \vec{\beta}_h \end{array} \right.$$

Alternative is to try and describe the data "locally" (or piece-wise)

We aim first for a model \hat{f} but for a model that is different depending on \vec{x}_{in}

K -NN = K nearest neighbours

$$\{f(\vec{x}_{in})\} \rightarrow \hat{y} = \hat{f}(\vec{x}_{in})$$

For continuous variables: $\vec{x}^{(i)} \in \mathbb{R}^p$, $y \in \mathbb{R}$
 $i=1, \dots, N$

Introduce a metric in the space of inputs:

$$\|\vec{x}^{(i)} - \vec{x}^{(j)}\|$$

Distance is a
choice

Given an input \vec{x}_{in} :

① Compute all distances between \vec{x}_{in} and the samples.

$$\|\vec{x}_{in} - \vec{x}^{(i)}\| \quad \forall i=1, \dots, N$$

② Find the K -nearest neighbours which define
the ~~neighbourhood~~ neighbourhood

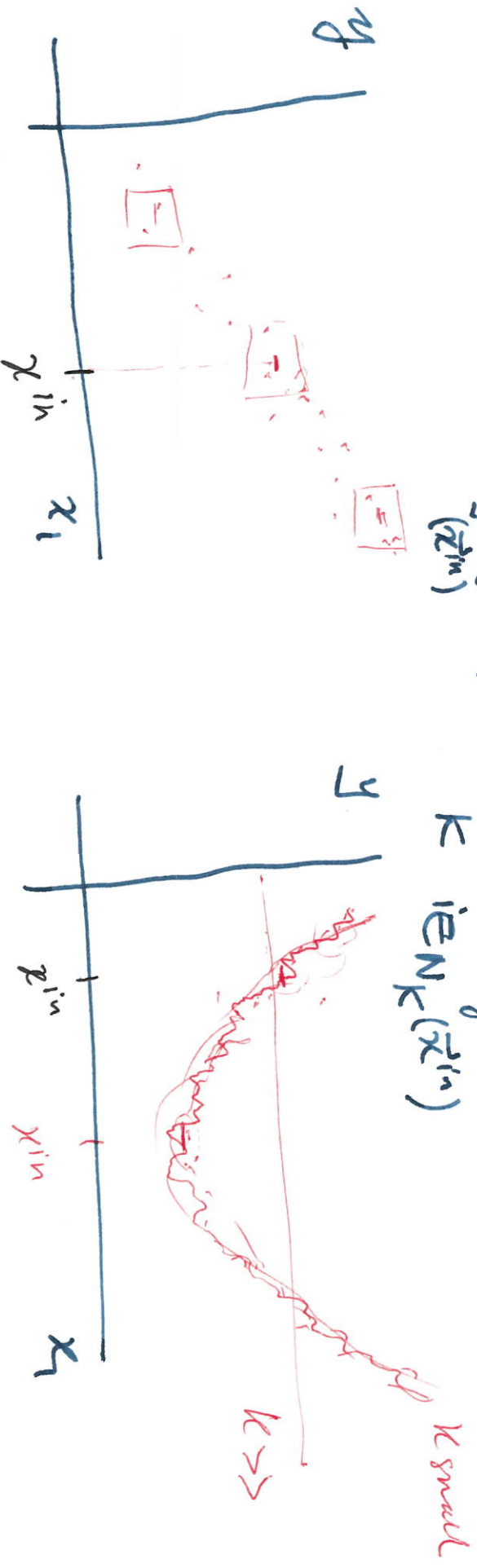
$$N_K(\vec{x}_{in})$$

K is a choice.

③

Simplest choice for predictor :

$$\hat{f}_{(\vec{x}^{in})}(\vec{x}^{in}) = \frac{1}{K} \sum_{i \in N_K(\vec{x}^{in})} y^{(i)}$$



The outcome depend on our choice :

(1) Distance and normalization matter a lot.

$$x_1 \dots x_p \quad \underbrace{10^{-10} < x_p < 10^{-5}}_{1 < x_1 < 10^6}$$

(2) It works better for small, relevant p .

(3) K matters a lot.

↳ allows to scan the big-variance tradeoff.