

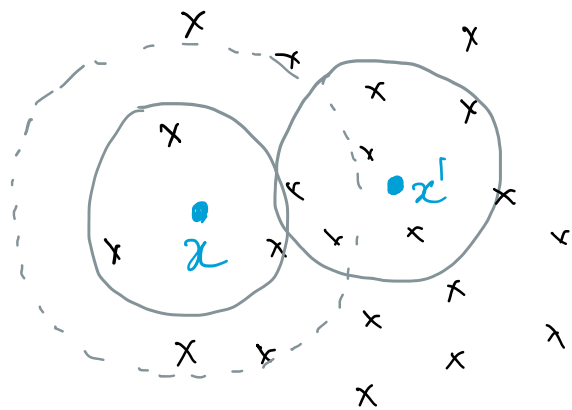
# L1: Data analysis and ML basics

- ML basics
  - o Supervised and unsupervised learning
  - o Data, model, and algorithm
  - o No-free-lunch theorem
- Training/test split and generalization error
  - o underfitting and overfitting
  - o Bias-variance trade-off
  - o Model selection and validation
- The break-down of three errors in machine learning
- Overview of kernel methods
  - o Parametric and non-parametric model

• KDE bias-variance trade-off

$f(x)$  density function in  $\mathbb{R}^d$ ,

$x_i \sim p(x) dx$ ,  $i=1, \dots, n$  i.i.d.



$$\hat{p}_\sigma(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\sigma^d} K\left(\frac{\|x_i - x\|}{\sigma}\right), \quad \sigma > 0$$

bandwidth param

$$K(z) = e^{-z^2/2}$$

other functions in  $\mathbb{R}^d$ , regular, decay

Question:  $|\hat{p}_\sigma(x) - p(x)| \leq ?$

- Bias error

$$\mathbb{E} \hat{p}_\sigma(x) = \frac{1}{\sigma^d} \int_{\mathbb{R}^d} K\left(\frac{\|y-x\|}{\sigma}\right) p(y) dy =: \bar{p}_\sigma(x)$$

If  $p$  and  $K$  are smooth,  $\sigma$  small,

$$y = x + \sigma v$$

$$\begin{aligned} \bar{p}_\sigma(x) &= \int_{\mathbb{R}^d} \frac{1}{\sigma^d} K\left(\frac{\|y-x\|}{\sigma}\right) p(y) dy \\ &= \int_{\mathbb{R}^d} K(\|v\|) \underbrace{p(x + \sigma v)}_{p(x) + \sigma \nabla p(x)^T v + O(\sigma^2)} dv \end{aligned}$$

suppose  $K(\|v\|)$  is symmetric, s.t.

$$\int_{\mathbb{R}^d} v K(\|v\|) dv = 0,$$

also  $K(\|v\|)$  is rescaled by a constant s.t.

$$\int_{\mathbb{R}^d} K(\|v\|) dv = 1$$

then

$$\begin{aligned} \bar{p}_\sigma(x) &= \int_{\mathbb{R}^d} K(\|v\|) \left( p(x) + \sigma \nabla p(x)^T v + O(\sigma^2) \right) dv \\ &= \underbrace{\left( \int_{\mathbb{R}^d} K(\|v\|) dv \right)}_1 p(x) + \sigma \nabla p(x)^T \underbrace{\left( \int_{\mathbb{R}^d} v K(\|v\|) dv \right)}_0 + O(\sigma^2) \\ &= p(x) + O(\sigma^2) \end{aligned}$$

- Variance error

$$\hat{p}_\sigma(x) - \mathbb{E} \hat{p}_\sigma(x) = ?$$

For fixed  $x \in \mathbb{R}^d$ ,

$$\hat{p}_\sigma(x) = \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{\sigma^d} K\left(\frac{\|x_i - x\|}{\sigma}\right)}_{\xi_i \text{ i.i.d. r.v.}}$$

$$\text{Var}(\xi_i) \leq \mathbb{E} \xi_i^2 = \frac{1}{\sigma^{2d}} \int_{\mathbb{R}^d} \underbrace{K\left(\frac{\|y-x\|}{\sigma}\right)}_{\|v\|}^2 p(y) dy$$

$$= \frac{1}{\sigma^{2d}} \int_{\mathbb{R}^d} K(\|v\|)^2 p(x + \sigma v) dv$$

$$\text{suppose } \int_{\mathbb{R}^d} K(\|v\|)^2 dv < \infty,$$

$$p \text{ bdd.} : |p(x)| \leq C, \forall x \in \mathbb{R}^d,$$

$$\text{then } \mathbb{E} \xi_i^2 \leq \frac{1}{\sigma^{2d}} \cdot C \cdot \int_{\mathbb{R}^d} K(\|v\|)^2 dv = O\left(\frac{1}{\sigma^{2d}}\right)$$

Then, by concentration of independent sum,

$$|\hat{p}_\sigma(x) - \mathbb{E} \hat{p}_\sigma(x)| = \tilde{O}\left(\sqrt{\frac{1}{n \sigma^{2d}}}\right)$$

- Trade-off.

$$O(\sigma^2) + O\left(n^{-1/2} \sigma^{-d/2}\right)$$

if we match

$$\sigma^2 \sim n^{-1/2} \sigma^{-d/2}$$

$$\Rightarrow \sigma \sim n^{-\frac{1}{4+d}}$$

The overall error is  $O\left(n^{-\frac{2}{2+d}}\right)$ .