# Non-parametric density estimation.
$\quad\quad\quad \hookrightarrow$ histogram

Consider iid r.v.s. $\{X_i\}_{i=1}^{n} \sim p(x)$ ← unknown

For simplicity, we consider $X_i \in [0,1]$.
$p(x)$ is non-zero only within $[0,1]$.
$|p'(x)| \leq L$ bounded first derivative

histogram partition $[0,1]$ into $M$ equispaced bins

$$B_1 = \left[0, \tfrac{1}{M}\right), \quad B_2 = \left[\tfrac{1}{M}, \tfrac{2}{M}\right) \dots$$

$$B_{M-1} = \left[\tfrac{M-2}{M}, \tfrac{M-1}{M}\right), \quad B_M = \left[\tfrac{M-1}{M}, 1\right)$$

For any value $x \in B_\ell$, the density estimate given by the histogram is

$$\hat{p}_n(x) = \frac{\#\text{ of obs. with } B_\ell}{n \text{ binwidth}}$$

$$= M \sum_{i=1}^{n} I(X_i \in B_\ell) \quad\quad I \to \text{indicator function}$$

$$E\left[\hat{p}_n(x)\right] = M\, E\left(\sum_{i=1}^{n} \frac{I(X_i \in B_\ell)}{n}\right)$$

$$= M\, P(X_i \in B_\ell) \quad \xrightarrow{\text{using true density}}$$

$$= M \int_{(\ell-1)/M}^{\ell/M} p(u)\, du = M\left[F_X\left(\tfrac{\ell}{M}\right) - F_X\left(\tfrac{\ell-1}{M}\right)\right]$$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \hookrightarrow$ true CDF

$$E\left(\hat{p}_n(x)\right) = \frac{F_X\left(\ell/m\right) - F_{X^-}\left(\frac{\ell-1}{m}\right)}{\frac{\ell}{m} - \frac{\ell-1}{m}} = p(x^*)$$

Taylor series exp
$$f(x) = f(x_0) + (x - x_0) f'(x^*)$$
where $x^* \in [x_0, x]$

$$F_X\left(\ell/m\right) = F_X\left(\frac{\ell-1}{m}\right) + \left(\frac{\ell}{m} - \frac{\ell-1}{m}\right) p(x^*)$$

$p \to$ derivative of $F_X$

$$x^* \in \left[\frac{\ell-1}{m}, \frac{\ell}{m}\right)$$

$$E\left(\hat{p}_n(x)\right) = p(x^*)$$

$$\text{bias}\left(\hat{p}_n(x)\right) = E\left(\hat{p}_n(x)\right) - p(x)$$

$$= p(x^*) - p(x)$$

$$\left(\begin{array}{l} p(x^*) = p(x) + (x^* - x) p'(x^{**}) \\ \qquad\qquad\qquad\qquad x^{**} \in [x^*, x] \end{array}\right)$$

$$\text{bias}\left(\hat{p}_n(x)\right) = p'(x^{**})(x^* - x)$$

$$\leq |p'(x^{**})| \underline{|x^* - x|}$$

$$\leq L \times 1/m = L/m$$

As $L \uparrow$, bias $\uparrow$   As $M \uparrow$, bias drops.

$$Var\left(\hat{p}_n(x)\right) = Var\left[\frac{M}{n}\sum_{i=1}^{r} I(X_i \in B_e)\right]$$

$$= \frac{M^2}{n^2}\sum_{i=1}^{n} Var\left[I(X_i \in B_e)\right] \quad \text{due to indep} \atop \text{of } \{X_i\}_{i=1}^{n}$$

$$\left(\begin{array}{c} I(X_i \in B_e) \longrightarrow \text{Bernoulli R.V. with} \\ \text{parameter} \quad \underline{P(X_i \in B_e)} \end{array}\right)$$

$$= \frac{M^2}{n^2} \times \not{n} \times \underline{P(X_i \in B_e)(1-P(X_i \in B_e))}$$

$$P(X_i \in B_e) = p(x^*)/M$$

$$Var\left(\hat{p}_n(x)\right) = \frac{M^2}{n}\frac{p(x^*)}{M}\left(1 - \frac{p(x^*)}{M}\right)$$

$$\leq \frac{M}{n}p(x^*) + \frac{\not{M^2}}{n\not{M^2}}p^2(x^*)$$

$$= \frac{M}{n}p(x^*) + \frac{p^2(x^*)}{n}$$

Var. increases with $M$ and drops with $n$.

$$MSE \atop (\hat{p}_n(x)) \leq \text{bias}^2 + Var$$

$$= \boxed{\frac{L^2}{M^2}} + \boxed{\frac{M}{n}}p(x^*) + \frac{p^2(x^*)}{n}$$

$$\partial MSE/\partial M = \frac{L^2(-2)}{M^3} + \frac{p(x^*)}{n} = 0$$

$$\frac{p(x^*)}{n} = \frac{2L^2}{m^3}$$

$$M = \left(\frac{2L^2}{p(x^*)} n\right)^{1/3} \longrightarrow O(n^{1/3})$$

$$\text{binwidth} = \frac{1}{M} = O(n^{-1/3})$$

Exact value of M is not computable because $p(x^*)$ is not known!

Plug in the optimal M into upper bound for MSE.

$$MSE_{opt} \leq O(n^{-2/3}) + \frac{n^{-1/3}}{n} p(x^*)$$
$$+ p^2(x^*)/n$$

$$= O(n^{-2/3})$$

A hist. approaches the true density at the error rate $O(n^{-2/3})$ in terms of MSE as long as # of bins i.e. M is $O(n^{1/3})$.