# 2.1 Histograms

Saturday, February 5, 2022    11:18 AM

## 2.1 Histograms

Histograms
- estimate density $f$ from windows of $[x_0, x_0 + h)$
  $x \in [x_0, x_0 + h)$:
  $$f(x_0) = F'(x_0) = \lim_{h \to 0^+} \frac{F(x_0 + h) - F(x_0)}{h}$$
  $$= \lim_{h \to 0^+} \frac{P[x_0 < X < x_0 + h]}{h}$$

- Bins:
  - $B_k := [t_k, t_{k+1})$; $t_k = b_0 + hk$, $k \in \mathbb{Z}$
  - samp pts
  - bandwidth $h$
- hist defined as:
  $$\hat{f}_H(x; t_0, h) := \frac{1}{nh} \sum_{i=1}^{n} 1\{X_i \in B_k : x \in B_k\}$$
  $$= \frac{v_k}{nh} \qquad v_k := \# \text{ obs in bin } B_k$$

- analysis of $\hat{f}_H(x; t_0, h)$ as RV:
  $$v_k \sim B(n, p_k)$$
  <span style="color:red">↑ #trials  ↑prob(suc)</span>   $p_k := P[X \in B_k] = \int_{B_k} f(t) dt$

→ if $f$ cont, then by MVT:
  $$p_k = h f(\xi_{k,h}) \qquad \text{<span style=\"color:blue\">(sc } p_k = \frac{v_k}{n})</span>$$
  in $\xi_{k,h} \in (t_k, t_{k+1})$

  <span style="color:red">→ $E\{\hat{f}_H\} = nh \, E\{v_k\} = \frac{np_k}{nh}$  by mean of binomial</span>

  $$E\{\hat{f}_H(x; b_0, h)\} = \frac{np_k}{nh} = f(\xi_{k,h})$$
  $$Var\{\hat{f}_H(x; t_0, h)\} = \frac{np_k(1-p_k)}{n^2 h} = \frac{f(\xi_{k,h})(1 - hf(\xi_{k,h}))}{nh}$$

  <span style="color:red">→ $Var\{\hat{f}_H\} = \frac{1}{n^2 h^2} Var\{v_k\} = \frac{1}{n^2 h^2}(np(1-p))$</span>

Observations:
  1) if $h \to 0$, then $\xi_{k,h} \to x \Rightarrow f(\xi_{k,h}) \to f(x)$
     <span style="color:red">Note</span>
     ↳ $x \in [t_k, t_{k+1}]$, $k \downarrow$ as $h \to 0$
     & interval collapses to $x$
     ↳ <span style="color:blue">unbiased estimator of $f(x)$ asymptotically
     when $h \to 0$ (small bins)</span>
  2) if $h \to 0$, ↑var
     ↳ $nh \to \infty$ to decrease
     (more points n)
  3) $f(\xi_{k,h})(1 - hf(\xi_{k,h})) \to f(x)$ as $h \to 0$

     ↳ more variability where higher density

↳ more variability where higher density

~ to matters!

Moving Histogram "Naive Density Estimator"

- Goal: aggregate $X_1, ..., X_n$ in intervals $(x-h, x+h)$

$$f(x) = F'(x) = \lim_{h \to 0^+} \frac{F(x+h) - F(x-h)}{2h}$$

$$= \lim_{h \to 0^+} \frac{P[x-h < X < x+h]}{2h} \longrightarrow \text{symmetric derivative}$$

↳ intervals based on eval point $x$ and are centered around it
⟹ directly estimates $f(x)$ ↳ proxy $f(x_0)$

- Here, we have:

$$\hat{f}_N(x; h) := \frac{1}{2nh} \sum_{i=1}^{n} 1\{x-h < X_i < x+h\}$$

- Analysis

$$\sum 1 \sim B(n, p_{x,h})$$

$$p_{x,h} := P[x-h < X < x+h] = F(x+h) - F(x-h)$$

$$\Rightarrow \frac{1}{2nh}(np)$$

$$E[\hat{f}_N(x;h)] = \frac{F(x+h) - F(x-h)}{2h}$$

$$Var[\hat{f}_N(x;h)] = \frac{F(x+h) - F(x-h)}{4nh^2} - \frac{[F(x+h) - F(x-h)]^2}{4nh^2}$$

$$\Rightarrow Var(f) = E[(x - E(x))^2]$$
$$+ \text{mechanics}$$

Observations:
  1) if $h \to 0$, then
     1.a) $E[\hat{f}_N(x;h)] \to f(x)$

     1.b) $Var[\hat{f}_N(x;h)] \approx \frac{f(x)}{2nh} - \frac{f(x)^2}{n} \to \infty$

  2) if $h \to \infty$
     2.a) $E[\hat{f}_N(x;h)] \to 0$

     2.b) $Var[\hat{f}_N(x;h)] \to 0$
  3) $Var \to 0$   if $nh \to \infty$

Consider — should we weight points closer to $x$ more heavily?