

Histograms

①

non-uniform bin size

Important: Now area becomes important

- if uniform binsize, height of bins usually $\rightarrow N_{\text{bin}}$,
with $\sum N_{\text{bin}} = N_{\text{TOT}}$
- if non uniform: $N_{\text{bin}} = \text{area} \Rightarrow \text{height} = N_{\text{bin}} / h_{\text{bin}}$

\rightarrow results not easy to interpret

Thus \rightarrow ~~plot~~ calculate/plot pdf

$$f(x)dx = P \text{ for } x \in [x, x+dx]$$

since $P = \frac{N_j}{N_{\text{TOT}}}$

$$\Rightarrow f(x) = \frac{N_j}{N_{\text{TOT}} \cdot h_j}$$

and $\int f(x)dx \Rightarrow \sum \frac{N_j}{N_{\text{TOT}} \cdot h_j} \cdot h_j = \sum \frac{N_j}{N_{\text{TOT}}} = 1$ automatically satisfied

$$\Rightarrow E(x) = \int x f(x) dx \rightarrow \sum x_i f_i \cdot h_i = \sum x_i \frac{N_j}{N_{\text{TOT}}}$$

$$\text{Var}(x) = \int (x - E(x))^2 f(x) dx \rightarrow \sum (x_i - E(x))^2 \cdot f_i \cdot h_i$$

See: p1-new.pro: bins ordered in such a way that
 $\rightarrow h$ and σ similar $N_j \approx 10$ everywhere (if possible)

How many bins? ②

1st rule of thumb: at least 5...10 events/bin

• statistical fluctuations inside bin follow Poisson-statistics

$$\text{error} \approx \sigma = \sqrt{\text{Var}} = \sqrt{(N_j)} = \frac{1}{\sqrt{N_j}} \text{ expect. value}$$

$$\text{relative error} \approx \frac{N_j}{\sqrt{N_j}} \approx \frac{1}{\sqrt{N_j}}$$

for $N_j = 10$, relative error = 33%

Note: same argumentation for photon-noise

if specific channel measures N_j photons

$$\Rightarrow \text{relative error} = \frac{1}{\sqrt{N_j}}, \quad S/N \left(\frac{\text{inverse}}{\text{rever}} \right) \text{ quantity} \approx \sqrt{N_j}$$

(for 200 photons/channel, $S/N \approx 14$, and relative error = 3%)

2nd rule of thumb

total number of events N_{TOT}	# of bins No-bins	average N_j (if uniform)
≤ 50	5 to 7	2...10
50...100	6 to 10	5...10
100...250	7 to 12	10...20
> 250	10 to 20	12...25

③

Note: various "rules" (based on simplified assumptions)

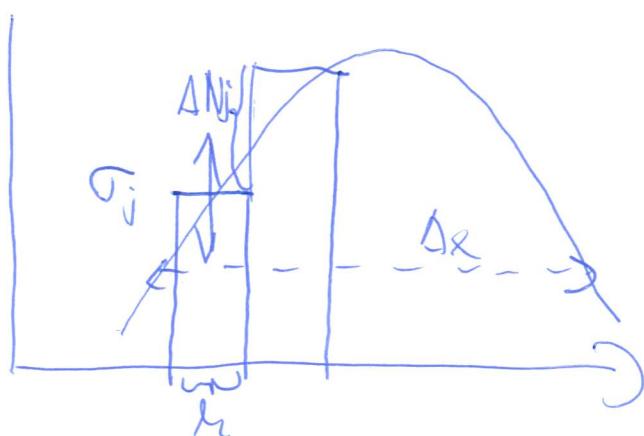
$\rightarrow h \text{ propo } \frac{1}{N_{\text{TOT}}^{\frac{1}{3}}}$ for larger samples [e.g. David Scott, Biometrika, 1909, 66, 3]

$$\text{since } N_{\text{bins}} \approx \frac{\Delta x}{h} \Rightarrow N_{\text{bins}} \approx (N_{\text{TOT}})^{\frac{1}{3}}$$

Accounting for "2nd rule" from above, here my personal suggestion (if distribution quite smooth)

N_{TOT}	$10 \dots 30$	$\Rightarrow N_{\text{bins}} \approx N_{\text{TOT}} / 5$
	≥ 30	$N_{\text{bins}} \approx 2 \cdot N_{\text{TOT}}^{\frac{1}{3}}$

Question: Why $h \approx N_{\text{TOT}}^{\frac{1}{3}}$?



h chosen in such a way that $\sigma_j \leq \Delta N_j$

$$\sigma_j \approx \sqrt{N_j} \quad (\text{Poisson}) \quad N_j \approx N_{\text{TOT}} \cdot \frac{h}{\Delta x}$$

Δx characteristic width of distribution

$$N_{j+1} \approx N_j + \frac{\Delta N_j}{\Delta x} \cdot h \Rightarrow \Delta N_j \approx N_{j+1} - N_j \approx \frac{\Delta N_j}{\Delta x} \cdot h \approx \frac{N_j}{\Delta x} \cdot h$$

From $\sigma_j = \sqrt{N_j}$:

(4)

$$\sqrt{\frac{N_{\text{tot}} \cdot h}{\Delta x}} = \frac{N_j \cdot h}{\Delta x} = N_{\text{tot}} \left(\frac{h}{\Delta x}\right)^2$$

$$\Rightarrow \left(\frac{h}{\Delta x}\right)^{3/2} = N_{\text{tot}}^{-1/2}$$

$$h = \frac{\Delta x}{N_{\text{tot}}^{1/3}} \quad \text{and} \quad N_{\text{bins}} = \frac{\Delta x}{h} \approx N_{\text{tot}}^{1/3}$$

Start of first bin

Have seen: start of first bin \rightarrow significant influence
on shape/properties
of histogram

subjective choice

more objective choice: "average shifted histogram"

(David Scott, Multivariate Density Estimation, John Wiley, 1992)

calculate histogram for various start values

(e.g., if bin-size = 10, for 0, 2, 4, 6, 8),

and average histograms

\rightarrow New histogram will have steps at Δoffset

(in the above example, the "new" binwidth will be 2)

\Rightarrow To obtain meaningful results, calculation of pdf

$(f(x) = \frac{N_j}{N_{\text{tot}} \cdot h})$ required! e.g.) 10,2 ; 10,1 ; 20,4

\rightarrow see average-shift.pro (average-shift, bins, Δoffset)