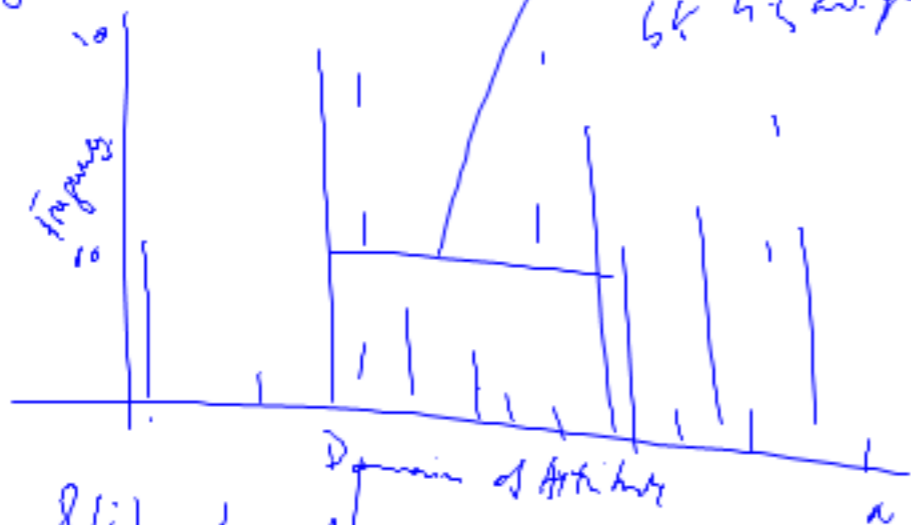


Lecture 23

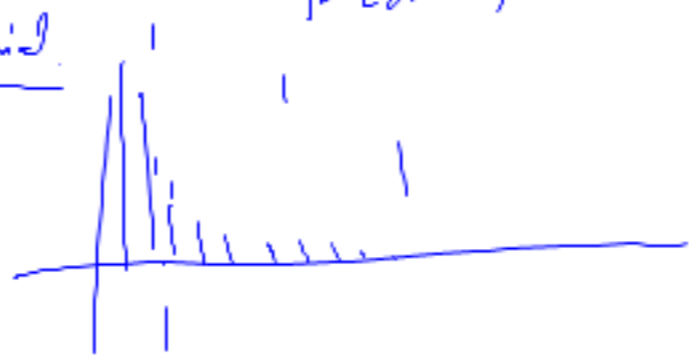
Histogram



$f(i)$ - true freq.

$$\hat{f}(i) = \sum_{j=0,1}^i f(j)$$

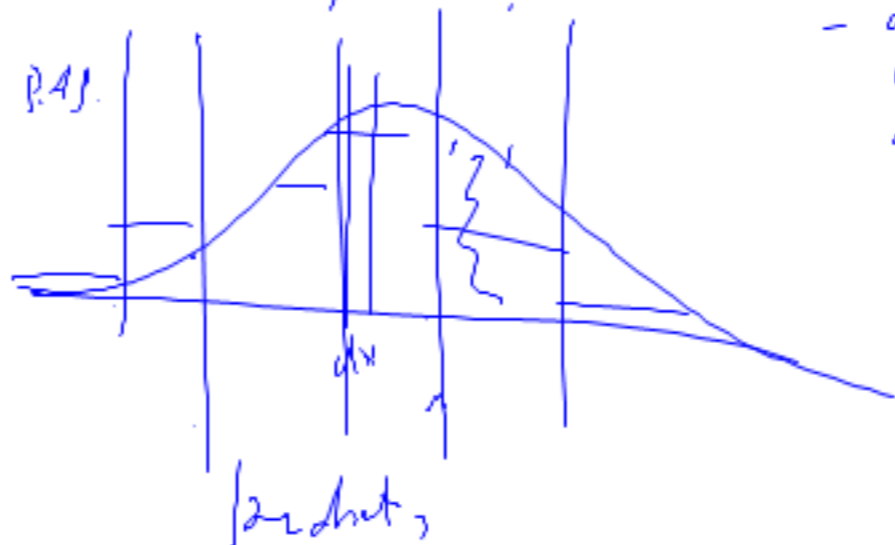
Binomial



History

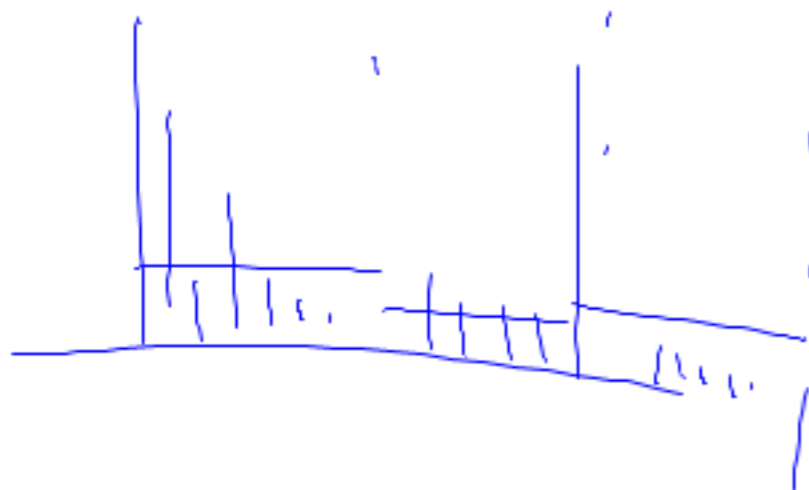
- Statistics
 - as nonparametric approximators of densities
- $N(\mu, \sigma)$

When CDF exact
is required
- develop theory to say
how good histogram
approx is.



Why deal with histograms?

- Space



Assumptions made by histograms
• there is an order on the attribute values

• frequencies, in the vertical axis, change smoothly.

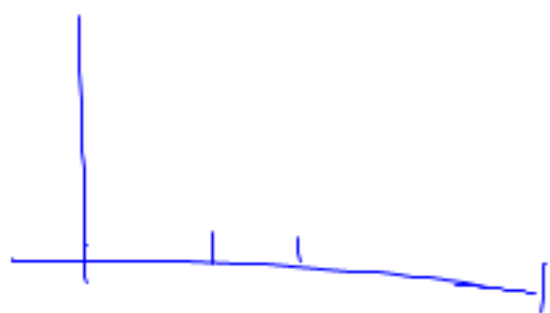
↑ usually true in statistics



Histograms in statistics are used as function approximators for p.d.f.

Types of histograms

- Equi-width

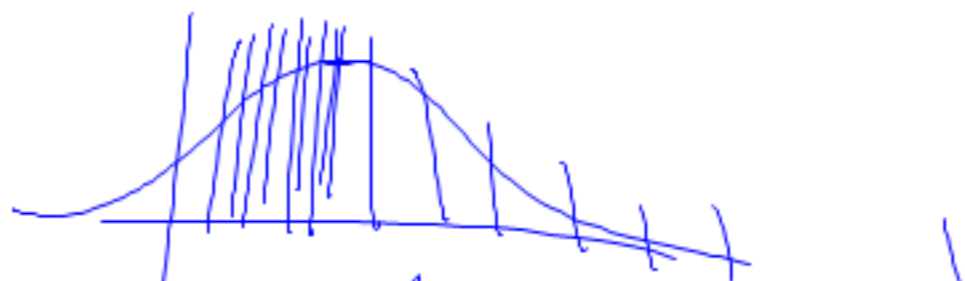


10 buckets of equal size.



- Equidepth histograms

- ensure the sum of freq in a bucket is about the same



- good choice for selectivity estimation
C.d.f.

Size of Join Computation

$R \bowtie S$

$$|R \bowtie S| = \sum_i \rho_i \delta_i$$

ρ_i
 δ_i) approx

$|R \bowtie S| \approx$

$$= \sum_{b \in B} \hat{\rho}_b \cdot \hat{\delta}_b \cdot |B|$$

