

Today: 1-D Continuous

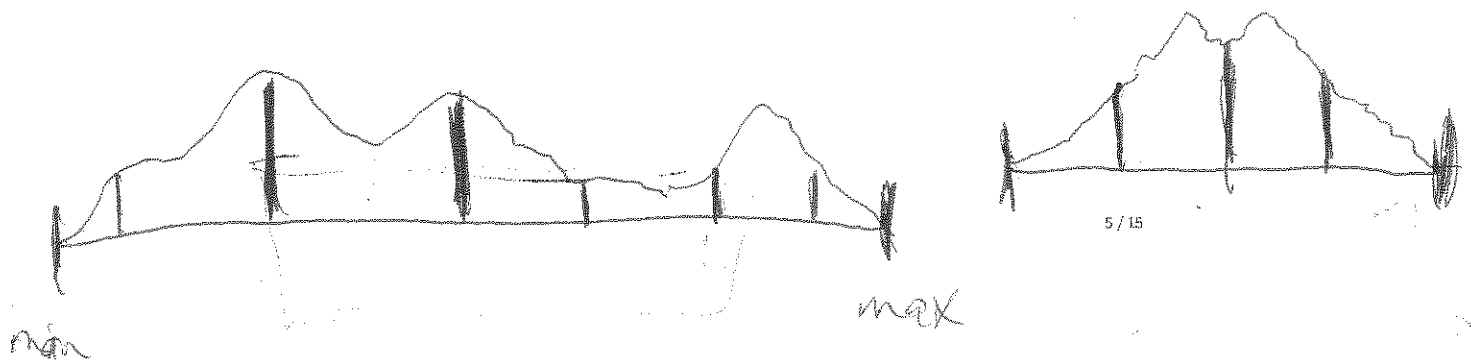
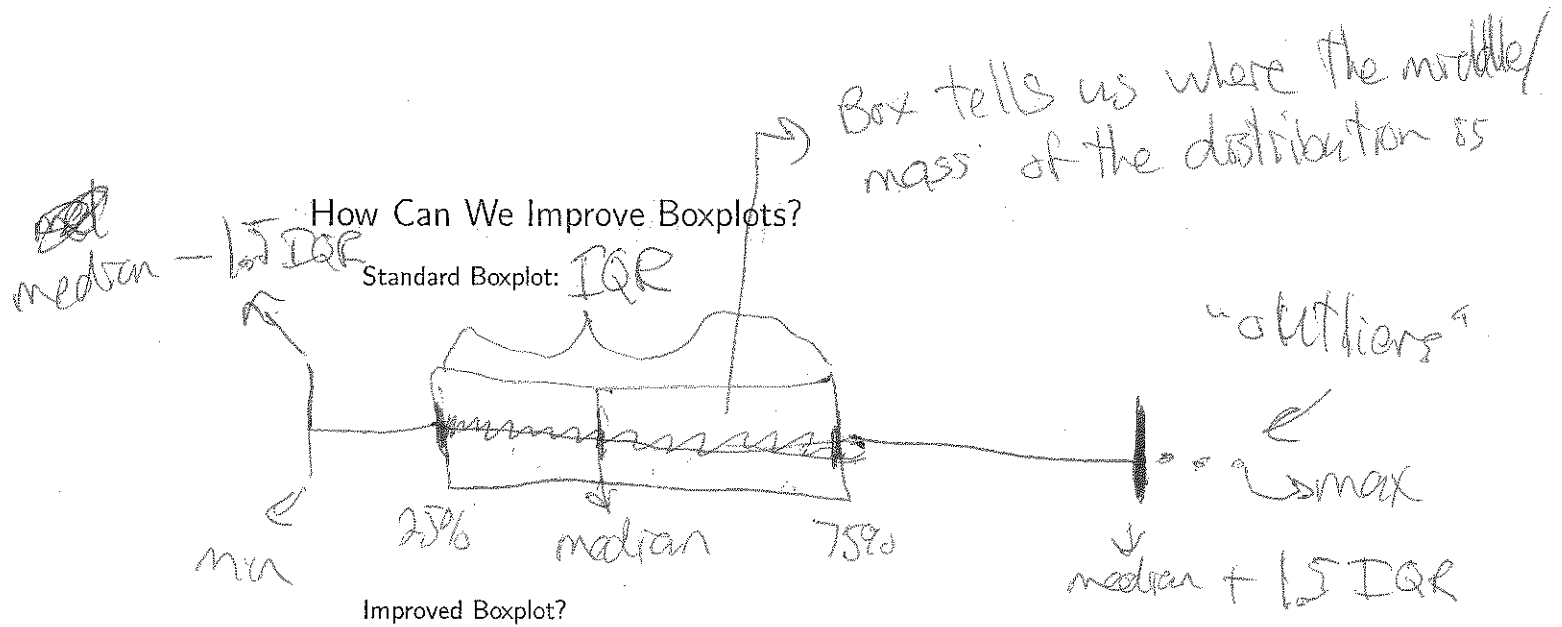
Sam Ventura
36-315

Today: Density Estimates
Box-Percentile Plots

Department of Statistics
Carnegie Mellon University

February 17, 2016

1/15



5/15

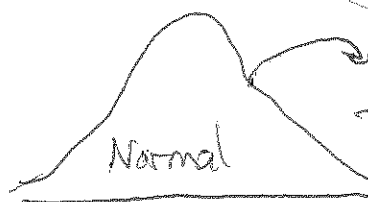
Height of line \propto density of distribution

→ "probability density function"

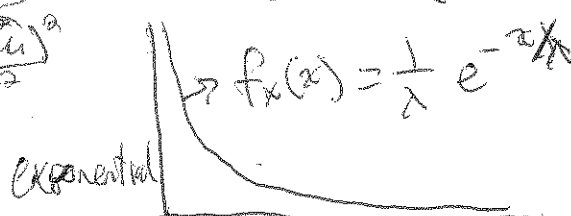
Continuous Densities → Smoothed-out function of the data's distribution

Theoretical:

X is RV, $f_X(x)$ reg $f_X(x=1)$



$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



$$f_X(x) = \frac{1}{\lambda} e^{-x/\lambda}$$

What if we don't know the underlying distribution?

→ "true distribution"

How can we estimate the empirical distribution?

↳ theoretical distn.

↳ "observed"

→ based on the data

so \hat{f} means "our estimate of f , the continuous density"

Parametric statistics: making assumption about the underlying distribution

6/15

1-D Kernel Density Estimation

→ Non-parametric → No assumptions
→ based on the data.

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$$

→ K contributes less to shape of the density estimate.

Small h : bumpy densities

large h : smoother densities

K contributes more to shape of density estimate

• n = # of observations in dataset

• h = "bandwidth" → "bin size" / "bin width"

It dictates "smoothness" vs. "rigidness" / "jaggedness" of the density estimate.

• K = "kernel function"

→ we get to choose this as well (parameter)
different functions give us different features of our resulting density estimate.

• x : the point at which we are estimating the density; should be w/in the range of the continuous variable

Ex) Boxcar kernel, Gaussian kernel, uniform kernel, etc.

7/15