

4/7/09 Math 152

- empirical cum. dist function.
explanation.

Kolmogorov Smirnov test.

R: Bees was example

- Survival functions / hazard rate.
explanation.

Guinea Pig data.

Q-Q plots.

additive and multiplicative treatment effects.

group III and V from

Bjerkedal Guinea Pig data.

Histograms, Density Curves, Stem and Leaf plots.

Measures of location.

Arith mean

median

$$\text{Var}(\log(1 - F_n(t))) \approx \frac{1}{100} \frac{1 - e^{-t}}{e^{-t}}$$

$$= \frac{1}{100} (e^t - 1)$$

$$\text{sd} \approx \frac{1}{10} \sqrt{e^t - 1}$$

e.c.d.f

x_1, \dots, x_n a batch of numbers.

$$F_n(x) = \frac{1}{n} (\#\{x_i \leq x\})$$

= proportion of observations $\leq x$.

If $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$
is the ordered batch of numbers.

$$F_n(x) = 0 \quad \text{for } x < x_{(1)}$$

$$= \frac{1}{n} \quad \text{for } x_{(1)} \leq x < x_{(2)}$$

\vdots

$$= \frac{k}{n} \quad \text{for } x_{(k)} \leq x < x_{(k+1)}$$

etc.

F_n is piecewise constant, with a jump of $\frac{1}{n}$
at each $x_{(k)}$.

(if the $x_{(k)}$ are not distinct the jump is
a multiple of $\frac{1}{n}$).

We have

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq t)$$

$$\text{where } I(x_i \leq t) = \begin{cases} 1 & x_i \leq t \\ 0 & x_i > t. \end{cases}$$

$$\text{So } \text{Var}(F_n(t)) = \frac{1}{n^2} \cdot n \cdot F(t)(1-F(t)) \\ = \frac{F(t)(1-F(t))}{n}$$

$$\log(1 - F_n(t)) \approx \log(1 - F(t))$$

$$\approx - \frac{1}{1 - F(t)} (F_n(t) - F(t)) + \dots$$

$$\text{Var}(\log(1 - F_n(t))) \approx \frac{\text{Var}(F_n(t))}{(1 - F(t))^2}$$

$$= \frac{1}{n} \frac{F(t)(1 - F(t))}{(1 - F(t))^2}$$

$$= \frac{1}{n} \left(\frac{F(t)}{1 - F(t)} \right)$$

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq t)$$

where

$$I(X_i \leq t) = \begin{cases} 1 & \text{if } X_i \leq t \\ 0 & \text{if } X_i > t \end{cases}$$

$$\text{Var}(F_n(t)) = \frac{1}{n^2} \cdot n F(t)(1 - F(t)) = \frac{F(t)(1 - F(t))}{n}$$

$$D_n = \max_{-\infty < t < \infty} |F_n(t) - F(t)|, \quad \text{dist of } D_n \text{ is independent of } F.$$

and this provides a test for

$H_0: X_1, \dots, X_n$ are i.i.d. $\sim F$.

Hazard Function

Survival function

$$S(t) = P(T > t) = 1 - F_T(t)$$

$$S_n(t) = 1 - F_n(t)$$

= proportion of data greater than t.

Hazard function

instantaneous death rate for individuals who have survived up to a given time

$$P(t \leq T \leq t + \delta \mid T \geq t) = \frac{P(t \leq T \leq t + \delta)}{P(T \geq t)}$$
$$\approx \frac{\delta f(t)}{1 - F(t)}$$

The hazard function is defined as

$$h(t) = \frac{f(t)}{1 - F(t)} = -\frac{d}{dt} \log(1 - F(t))$$
$$= -\frac{d}{dt} \log S(t)$$

→ Since $\text{Var}[\log(1 - F_n(t))] \approx \frac{1}{n} \frac{F(t)}{1 - F(t)}$

values of t s.t. $F(t) \rightarrow$ near 1
must be discarded.

→ See Calculation on next page.

$$\log(1 - F_n(t)) \approx \log(1 - F(t))$$

$$- \frac{1}{1 - F(t)} (F_n(t) - F(t)) + \dots$$

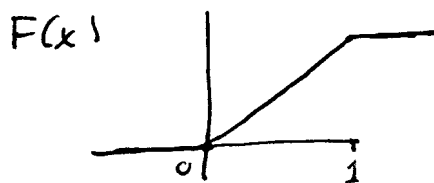
$$\text{Var}(\log(1 - F_n(t))) \approx \frac{\text{Var}(F_n(t))}{(1 - F(t))^2}$$

$$= \frac{1}{n} \frac{F(t)(1 - F(t))}{(1 - F(t))^2}$$

$$= \frac{1}{n} \left(\frac{F(t)}{1 - F(t)} \right)$$

See problems 2 and 8.

2. X_1, \dots, X_n i.i.d. $U[0, 1]$.



$$\sqrt{E((F_n(x) - x)^2)} = \sqrt{\frac{F_n(x)(1 - F_n(x))}{n}}$$

8. $\text{Var}(\log(1 - F_n(t))) \approx \frac{1}{100} \frac{1 - e^{-t}}{e^{-t}}$

$$\text{so s.d.} \approx \frac{1}{10} \sqrt{\frac{1 - e^{-t}}{e^{-t}}} = \frac{1}{100} (e^t - 1)$$

A surprising result (which we won't prove) is that the distribution of

$$D_n = \max_{-\infty < t < \infty} |F_n(t) - F(t)|$$

is independent of F .

This provides a test for

$H_0: X_1, \dots, X_n$ are i.i.d with cdf F .

R: Beeswax example.

e.g. Guinea Pig data.

as dosage increases

instantaneous mortality rates increase more ~~steadily~~ ^{quickly} and reach higher levels.

increased mortality rate sets in earlier for higher dosage groups and seems greater.

notice the large dispersion as t increases

again:

$$\text{Var} \{ \log(1 - F_n(t)) \} \approx \frac{\text{Var}(1 - F_n(t))}{(1 - F(t))^2}$$

$$= \frac{1}{n} \frac{F(t)}{(1 - F(t))}$$

so the variance becomes large as $F(t) \rightarrow 1$.

Q-Q plots

When making probability plots we plotted sample quantiles

(assigning the value $X_{(k)}$ to the $\frac{k}{n+1}$ quantile when there are n data points)

vs. the quantiles of a theoretical distribution

as a way of investigating how well the theoretical dist. fits the data.

To compare two batches of numbers,
(which may come from the same or different distributions)
we will plot
the empirical quantiles of one batch
vs. the other

$$X_{(1)} \dots X_{(n)} \quad \text{vs.} \quad Y_{(1)} \dots Y_{(n)}$$

plot $(X_{(i)}, Y_{(i)})$.

(if the batches are of different sizes, there
is an interpolation procedure).

e.g. (additive treatment effect).

(cdf's) F models a control group's data (x)
 G models a treatment group's data (y)

If the treatment ~~increases~~ changes the
expected response of each individual by
the same fixed amount h

$$\text{and } F(x_p) = p = G(y_p)$$

then we should expect $y_p = x_p + h$.

and the Q-Q plot will be a line
with slope 1 and intercept h .

We have $G(y_p) = F(y_p - h)$ so an additive
treatment effect corresponds to a shift of the cdf.

multiplicative treatment effect.

response is multiplied by a constant c .

$$y_p = c x_p$$

and the Q-Q plot will be a line with slope c and intercept 0.

$$\text{We have } G(y_p) = p = F(x_p) = F\left(\frac{y_p}{c}\right)$$

so the treatment effect corresponds to a rescaling of the c.d.f.

e.g. Q-Q plot of III vs. V from the Guinea Pig data.