# Stat 231

October 5, 2016

Syllabus $\leq$ Friday, Oct 7.

Tutorial : 3-30 $\longrightarrow$ Surya

            STP 105.

      6-00 $\longrightarrow$ Cyntha.

         DC

# Roadmap

## MODEL SELECTION METHODS

- Graphical methods
    - Relative Frequency
    - cdf
    - Q-Q-plot
    - Run charts.

- Numerical methods: (Ch 7)
    - Observed versus Expected frequencie's

DATA.

↓
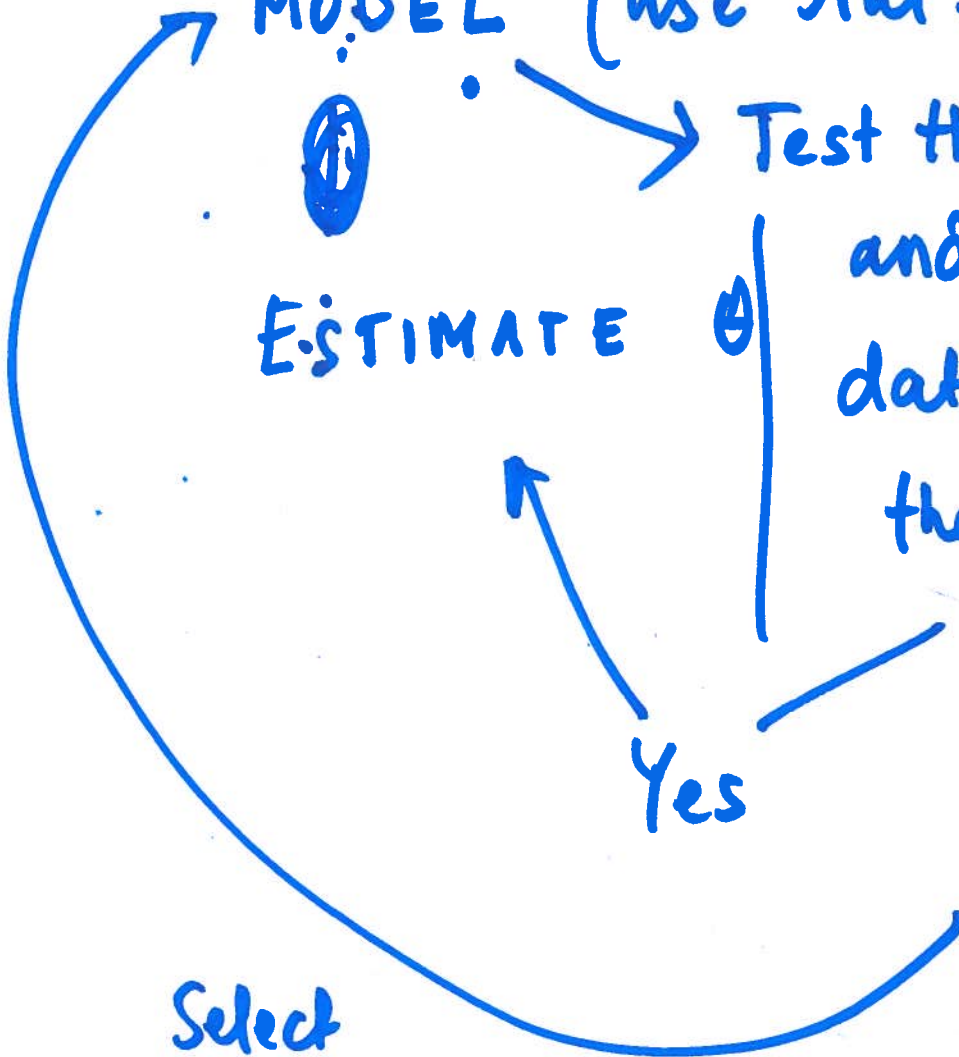
MODEL (use Stat. theory)

θ

→ Test the model
and see the
data fits
the model ?

ESTIMATE θ

Yes                    No

Select
a different model.

How to select the "right" model?

.

Graphical ways.

· Relative frequency histogram method.

Example: Suppose the model chosen is Gaussian $(r, \sigma)$.

Draw the density histogram $G(\hat{r}, \hat{\sigma})$
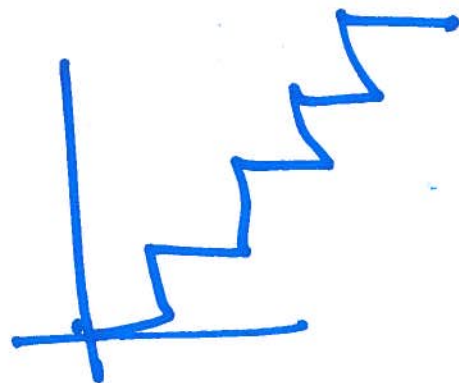$D.$



$b_{ins}$

The parameters of the superimposed distribution is chosen by the method of max. Likelihood.

- Compar· the empirical cdf with the theoretical cdf.
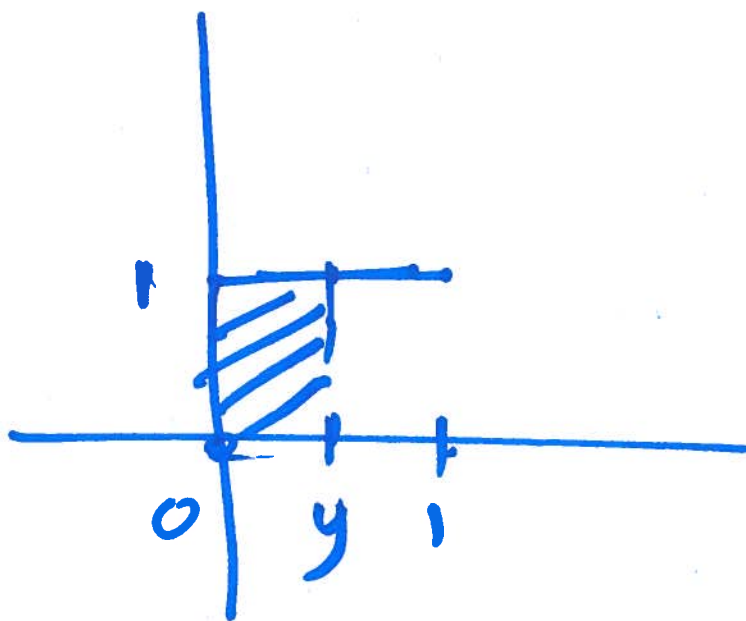
$$\{ y_1, \dots y_n \}$$

$$\hat{F}(y) = \frac{\# \text{ of obs} \leq y}{n} = ECDF$$
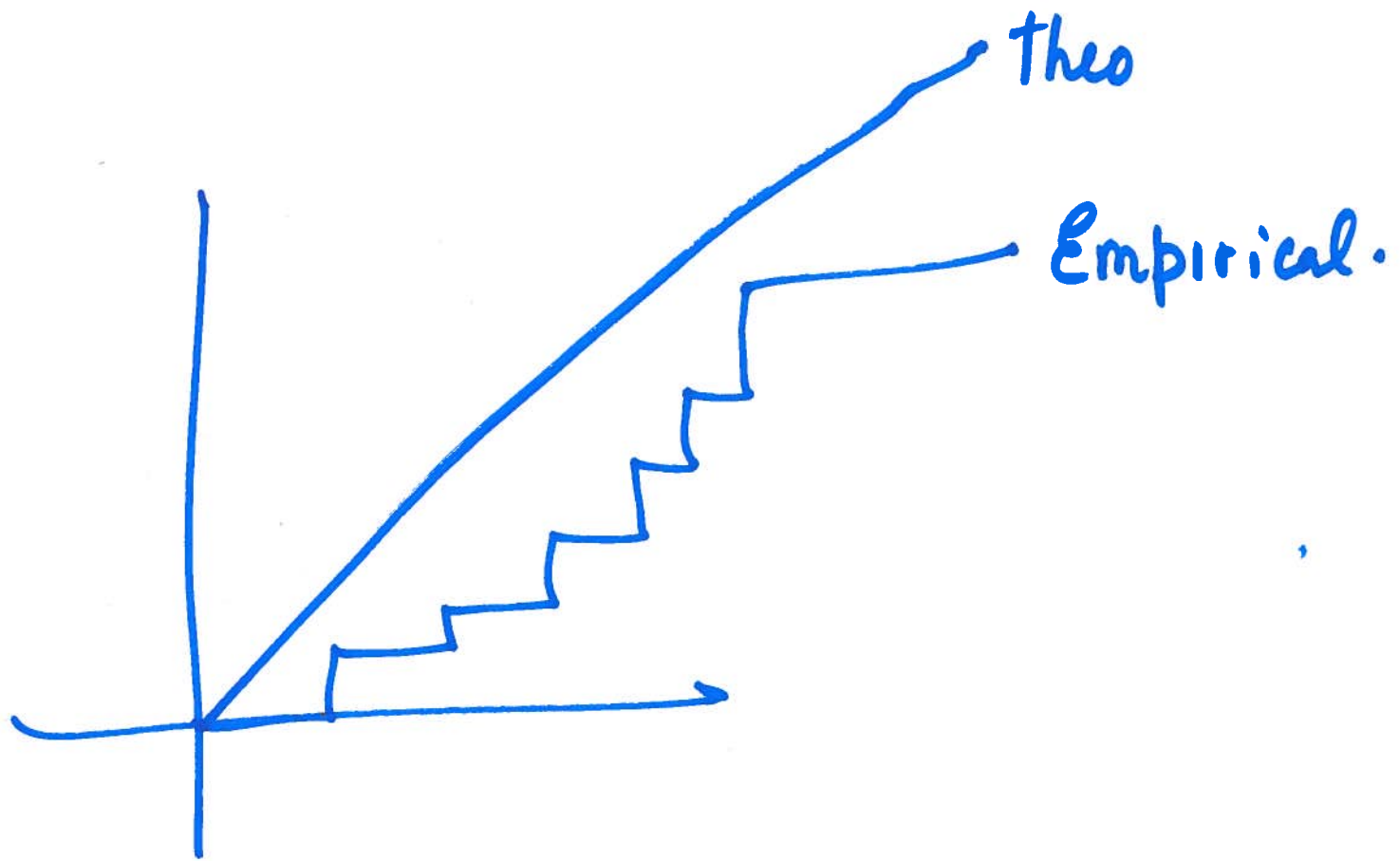
$$F(y) = P(Y \leq y)$$

Example $\{y_1, \ldots y_n\}$ DATA

Model: $Y_L \sim Uni(0,1)$

independent.

$F(y)$

$P(Y \leq y) = y$

cdf of $U(0,1) = F(y) = y$

$= 45°$ line

Superimpose the theoretical cdf on the empirical cdf and compare their shapes.

# The Q-Q-plot

Typically, the Q-Q plot is used to check for Gaussian.

Q

$$\{y_1, \ldots y_n\}$$

Arranged.

$$\downarrow$$

$$\{y_{(1)}, y_{(2)} \ldots y_{(n)}\}$$

The Q-Q plot : plots the sample quantiles against the quantiles of the standard normal distribution.
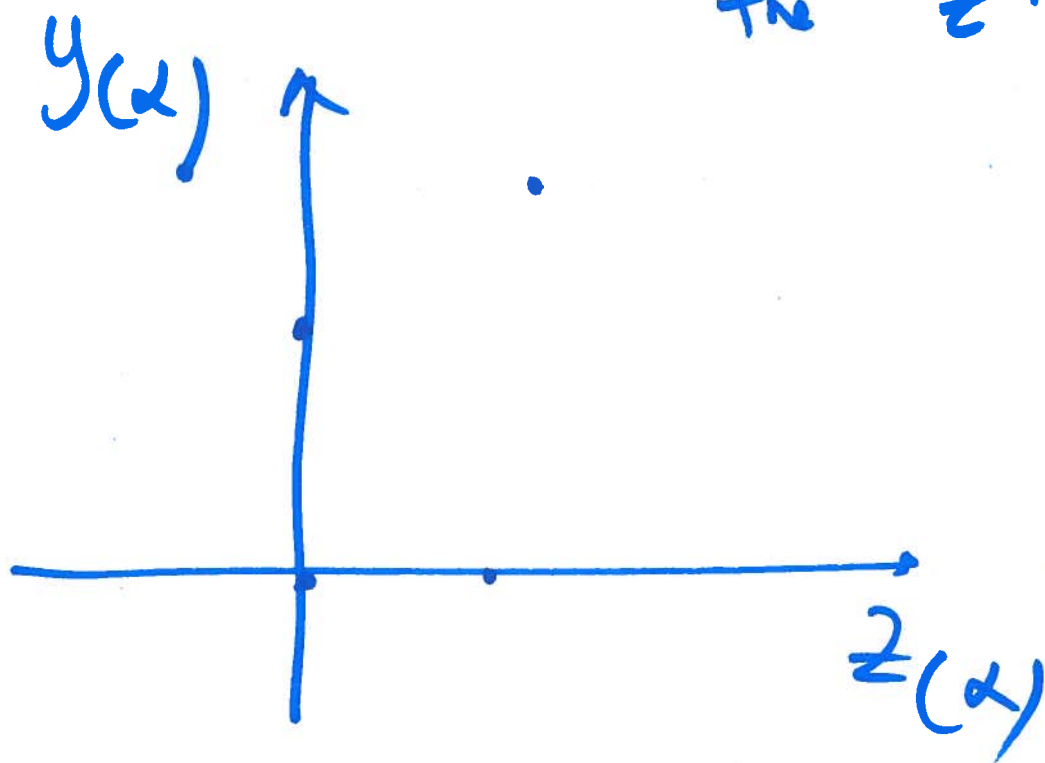
$$Z \sim \mathcal{G}(0, 1)$$

$$QQ \rightarrow \left( y_{(\alpha)}, z_{(\alpha)} \right)$$

$$y_{(\alpha)} = \alpha^{th} \text{ quantile of the data set}$$

$$z_{(\alpha)} = \alpha^{th} \text{ quantile of the } z \sim G(0,1)$$
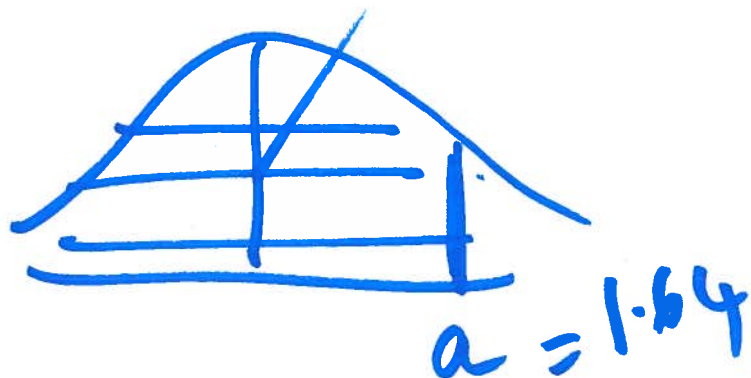
$$y_1, \ldots \ldots y_{101}$$

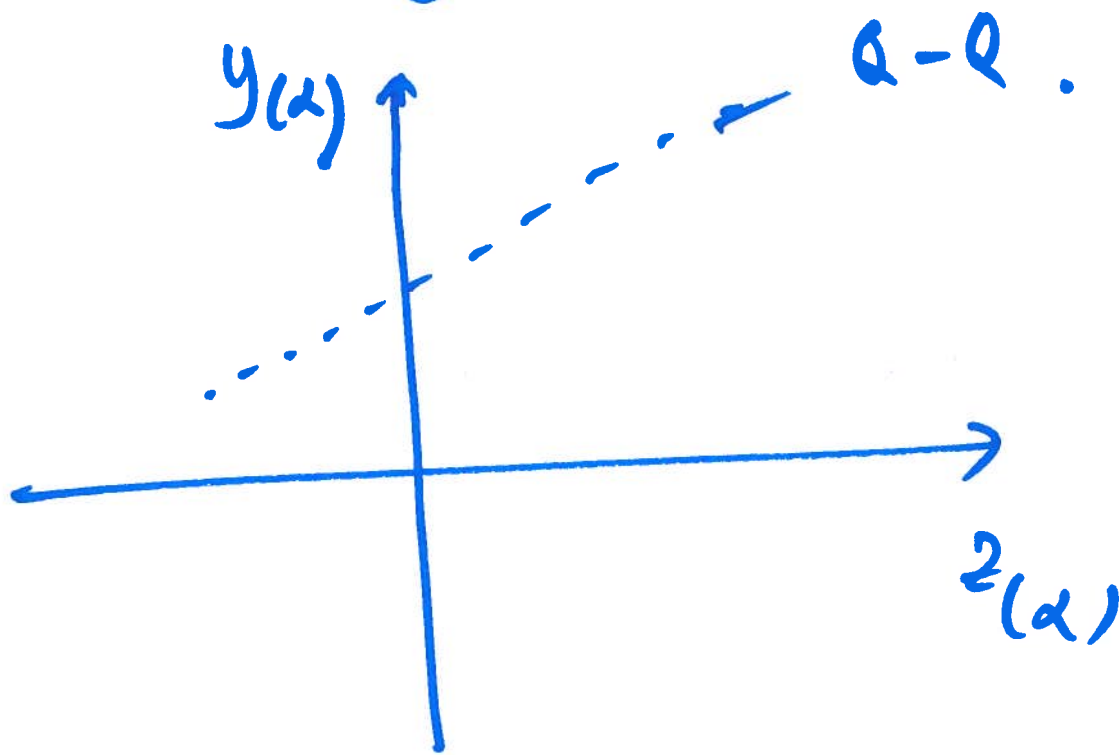Median $\longrightarrow$ y-value $\rightarrow$ median of the data

$z$-value $\longrightarrow$ 0

$95^{th}$ percentile $\rightarrow$ y-value $\rightarrow$ $95^{th}$ %ile of the data
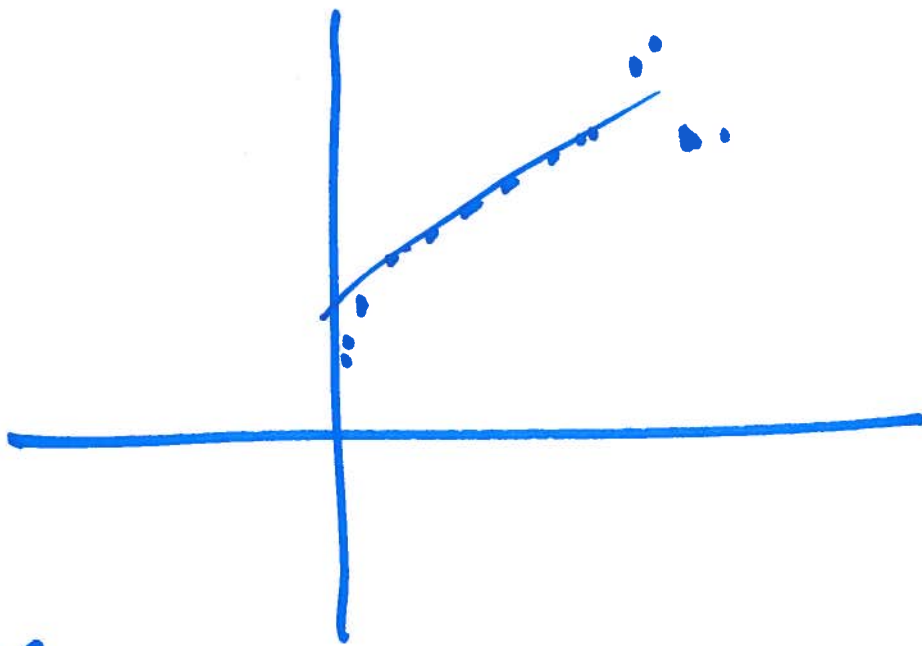
0.95

$z$-value: $95^{th}$ percentile of $z$



$a = 1.64$

If the Q-Q plot is a straight line,
that is strong evidence that the
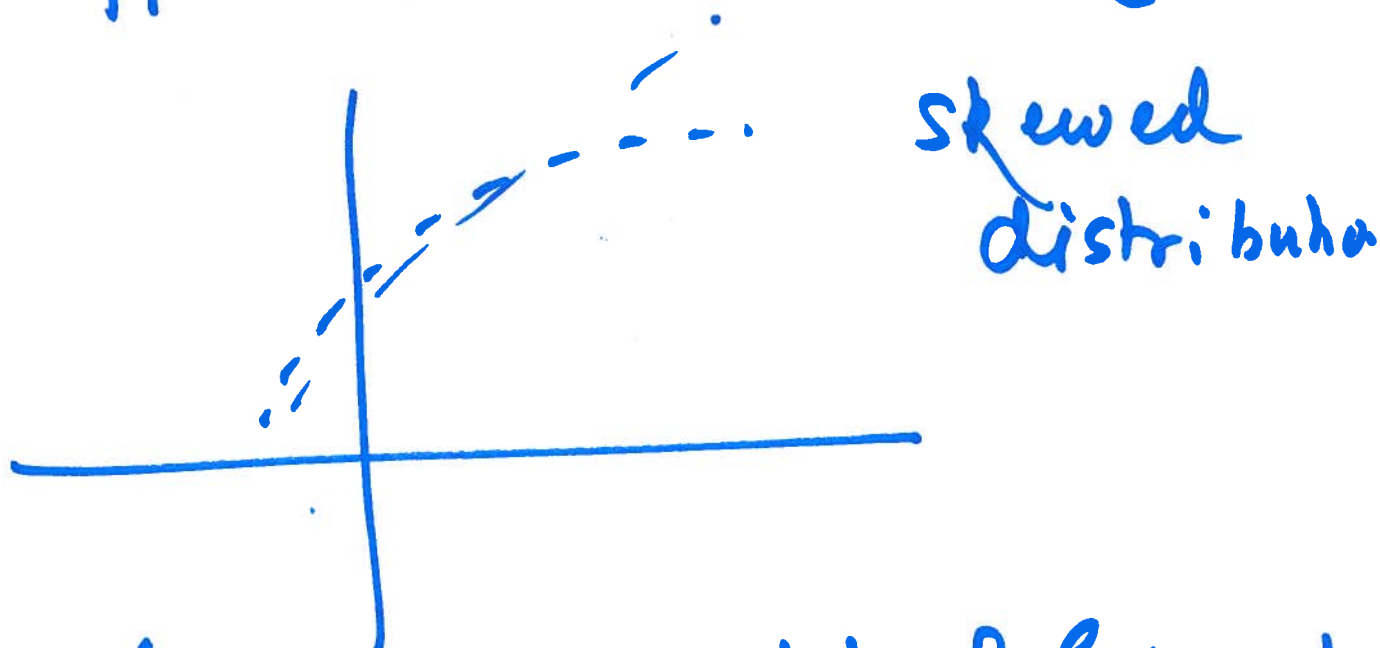data is Gaussian.

$y_{(\alpha)}$

Q-Q.

$z_{(\alpha)}$

# Notes about the Q-Q plot

- Even if the data is actually Gaussian the extreme points are typically off the line



(Why?)

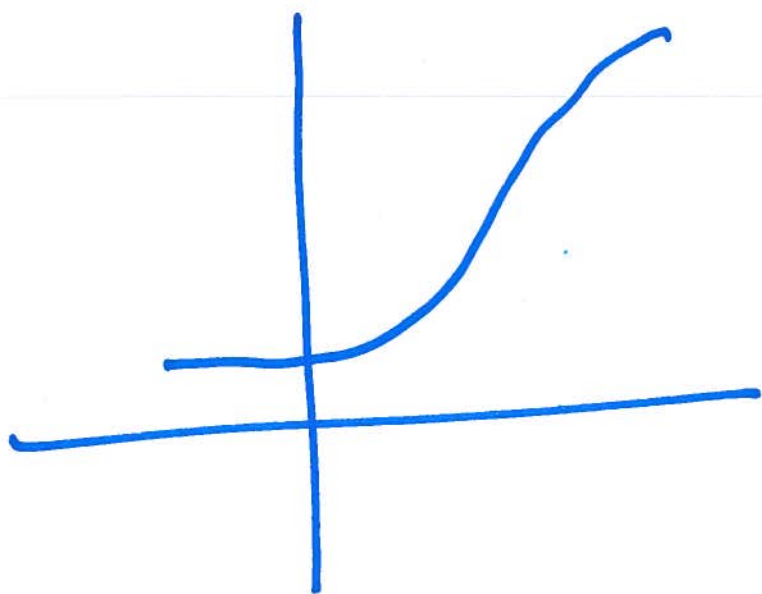• Apply the Q-Q plot to other distributions. (Exponential).

• Suppose the distribution is asymmetric.



skewed distribution

(Try drawing Q-Q plot of Exponential against Normal)

Typically, with a higher kurtosis,
the Q-Q plot looks like a $.

# Numerical methods of model checking

$Y_i = \#$ of accidents on a highway in a 1 hour period during rush-hour.

$$\{ y_1, \dots, y_n \}$$

$$\cdot \quad Y_i \sim Poi(\mu) \rightarrow \text{SUGGESTED MODEL.}$$

DATA

| # of acci | frequency | Expected frequencies |
|-----------|-----------|----------------------|
| 0 | 10 | |
| 1 | 20 | |
| 2 | 40 | |
| 3 | 20 | |
| 4 | 10 | |
| 5 | 0 | |
| 6 | 0 | |
| | 1 | |

Calculate $p$

$$P(\hat{Y} = 0) = \frac{e^{-\hat{p}} \hat{p}^0}{0!} = 0\%$$

~~$\times$ $n$~~

Expected frequency

of $0$ $\rightarrow$ $\dfrac{e^{-\hat{p}} \hat{p}^0}{0!} \times n$

$m =$ sample swe

We compare the observed with
the expected frequencies
and see whether they are "close
enough" .

---

## Continuous

Data

| Data | Freq |
| --- | --- |
| $[0, 100)$ | 20 |
| $[100, 200)$ | 30 |
| $[200, 300)$ | 20 |
| $\geq 300$ | 10 |

$Y_i \sim \text{Exp}(\mu)$

$Y_1, \ldots, Y_n$.