

STA 302 / 1001

Note Title

10/20/2011

Test has been interpretation / response to SAS
output + there is / conceptual / deviation
function

Page 1 is part 1.

Adjusted assignment 1 made on Blackboard.

Normal quantile plots

- assess whether data (assumed to be i.i.d. realizations of a r.v.) come from a normal distribution

Order statistics of r_1, \dots, r_n

$$r_{(1)} = \min \{ r_1, \dots, r_n \}$$

$$r_{(2)} = \text{second smallest}$$

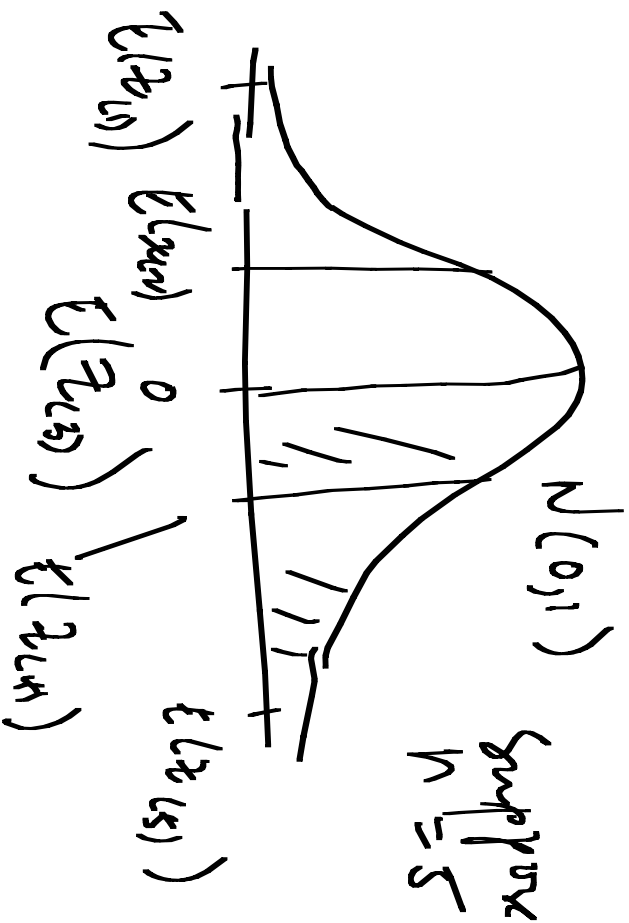
⋮

$$r_{(n)} = \max \{ r_1, \dots, r_n \}$$

Normal scores for samples of size n

$$\overline{E}(Z_{(1)}), \dots, E(Z_{(n)})$$

where Z_1, \dots, Z_n are a sample of size n from $N(0,1)$ distribution



Equal probability
between each
normal score

If r_1, \dots, r_n are random sample from $N(\mu, \sigma^2)$

$$Z_i = \frac{r_i - \mu}{\sigma} \Rightarrow \text{Sample from } N(0, 1)$$

$$Z_{(i)} = \frac{r_{(i)} - \mu}{\sigma}$$

So if r_i from $N(\mu, \sigma^2)$ distribution

$$\frac{r_{(i)} - \mu}{\sigma} \approx E(Z_{(i)})$$

$$\text{or } r_{(i)} = \sigma E(Z_{(i)}) + \mu$$

So plot $r_{(i)}$ vs $E(z_{(i)})$ should be approx. straight.

How to get $E(z_{(i)})$

$$\text{Then } E(z_{(i)}) = \Phi^{-1}\left(\frac{i}{n}\right)$$

Invest
stock
normal cdf

↪ value s.t prob $\leq i/n$

Flow problem when $i = n$, $\Phi^{-1} \left(\frac{n}{n} \right) = \infty$

Most common advice:

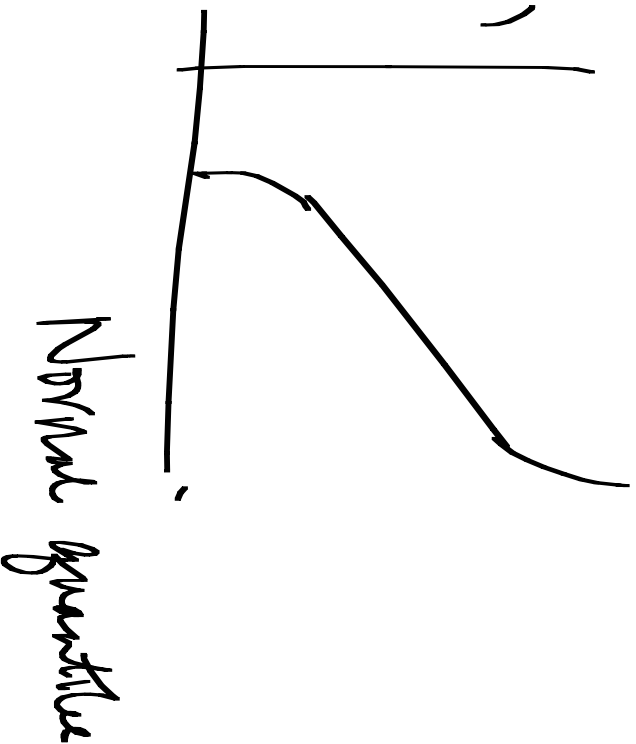
$$\underline{P_{\text{hom}}} \quad E(Z_{\text{vis}}) \doteq \Phi^{-1} \left(\frac{i - 3/8}{n + 1/4} \right)$$

- has been shown to work well in simulations

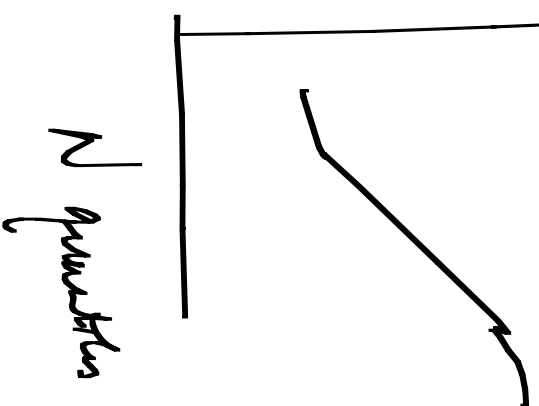
If assumption of normality is relaxed slightly to allow for heavier tails but still symmetric,

F - and t-tests and CIs for parameters are still approximately correct (p-value close, CI coverage close) why? Estimators of parameters are linear combinations of v.i.'s and CLT

Heavy tails: $r(i)$

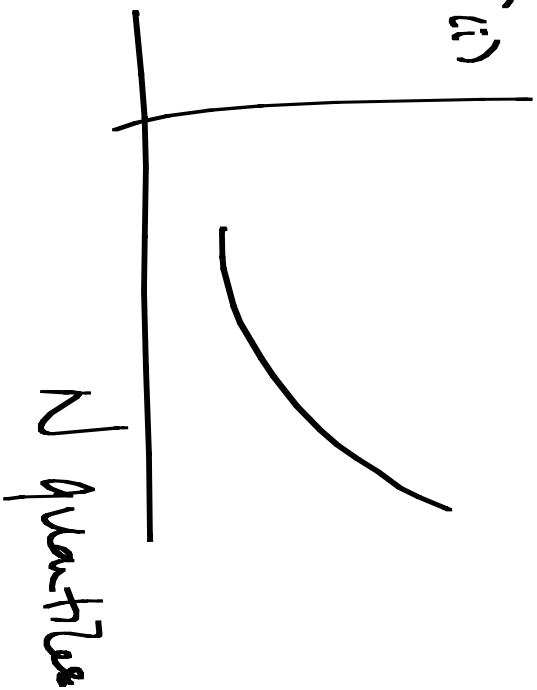


$r(i)$ Light tails



Skewed

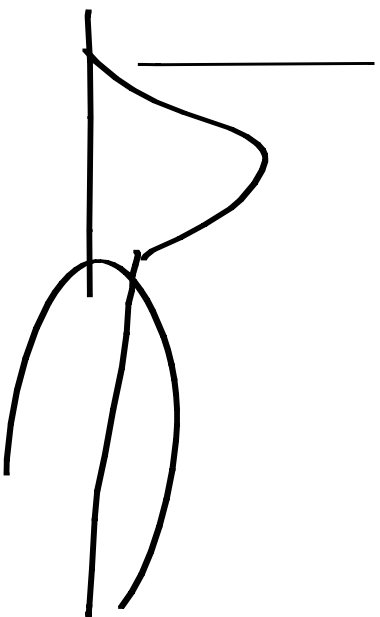
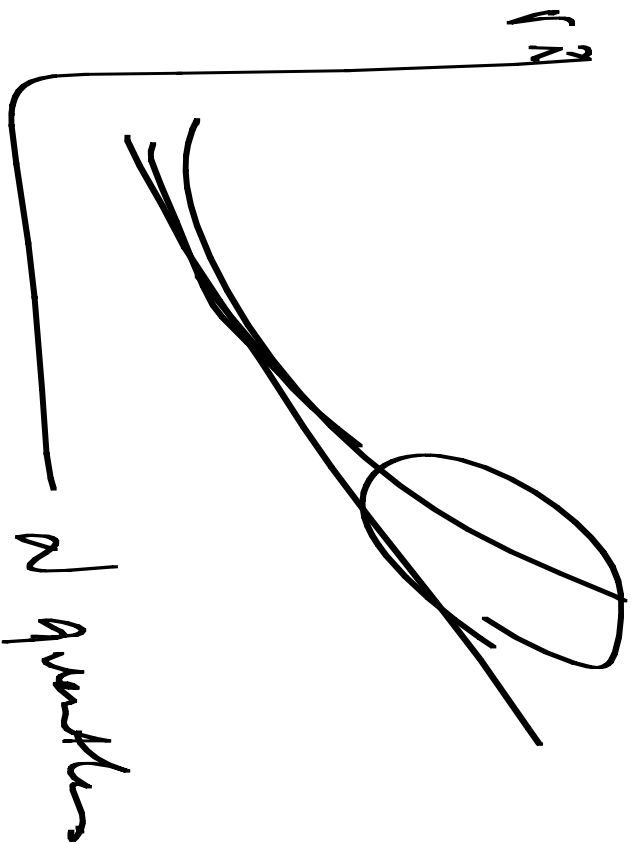
$r(i)$



Mean Example
for
sites

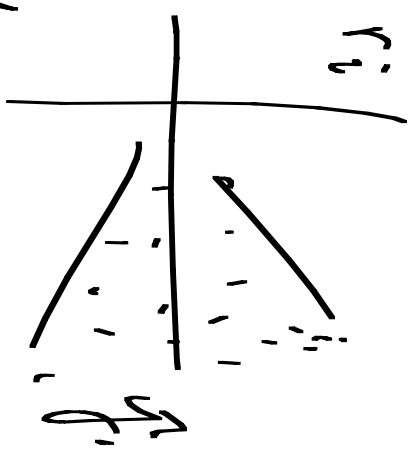
- residual plots (vs t_{true} or \hat{y})
show curvature
- this might indicate that a polynomial model would be better

- but in this case, that doesn't fix the problem
- The real problem: data collected over time, violation of 5-M condition of uncorrelated errors.
(to fix this: STA 457)



TAKES FOREMANNS TO STABILIZE THE VARIANCE

(Theoretical result: "Delta Method")



If γ has a distribution with mean μ and variance σ_γ^2

Then $Z = f(\gamma)$ has mean and variance approximately equal to

$$E(Z) \approx f(\mu)$$

("1st order linear

$$\text{Var}(Z) \approx \sigma_T^2 [f'(\mu)]^2 \quad \text{"approximation"}$$

"Proof"

$$Z = f(T) \quad \text{Taylor series expansion about } \mu$$

$$= f(\mu) + (T-\mu) f'(\mu) + \text{remainder}$$

$$E(Z) = f(\mu) + E(T-\mu) f'(\mu) + E(\text{remainder})$$

$$\approx f(\mu)$$

$$\text{Var}(Z) \approx E[(Z - f(\mu))^2] + \text{different remainder term}$$

$$= E[(T-\mu) f'(\mu)]^2 + \dots$$

$$= [f'(\mu)]^2 E[(Y-\mu)^2]$$

$$= \sigma^2 [f'(\mu)]^2$$

This result gets used to derive variances
 stabilizing transformations

In SLR, $E(Y_i) = \beta_0 + \beta_1 X_i = \mu_i$
 and $\text{Var}(Y_i) = \sigma^2$

doesn't depend on i ,
 can't depend on μ_i

Suppose Y_i has mean μ_i and variance
proportional to a function of μ_i
i.e. $\text{Var}(Y_i) \propto V(\mu_i)$, say

Want to find a transformation $Z = f(Y)$ so
 $\text{Var}(Z) \approx \text{const}$

Want f such that $\text{Var}(Z) \propto [f'(\mu)]^2 V(\mu)$

$$\text{Want } [f'(\mu)]^2 = \frac{c}{V(\mu)}, \quad c_{\text{constant}}$$

$$f'(\mu) \propto \frac{1}{\sqrt{V(\mu)}}$$

$$f(\mu) \propto \int \frac{1}{\sqrt{V(\mu)}} d\mu$$

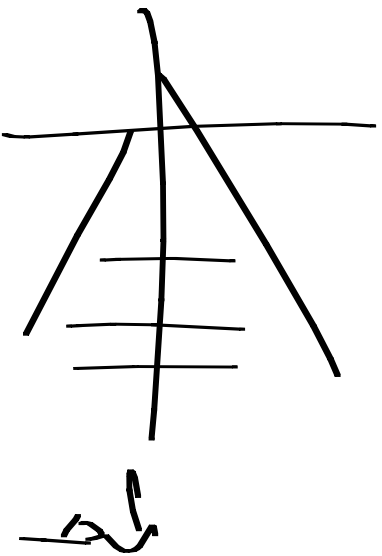
For example
(Poisson)
constants
for

$$\text{Sample } E(t_i) = \mu_i \\ \text{and } \text{Var}(t_i) = \mu_i$$

$$V(\mu) = \mu \\ f'(\mu) \propto \frac{1}{\sqrt{\mu}}$$

$$f(\mu) \propto \sqrt{\mu}$$

If $\text{Var}(X)$ is linearly proportional to $E(X)$, then $Z = \sqrt{X}$ has variance that's approximately constant.



Another example:
Exponential
Gaussian

$$E(Y) = \mu$$

$$\text{Var}(Y) \propto \mu^2$$

$$V(\mu) = \mu^2$$

$$f'(\mu) \propto \frac{1}{\mu}$$

$$f(\mu) \propto \log(\mu)$$

NB!

$$n \log'' = n''$$

(natural logarithm)

Always for me

(+ SAS)

If $\text{Var}(Y) \uparrow$ faster than linearly in $E(Y)$,
log transformation of Y stabilizes the variance

In regression, these are the most useful transformations

- log stronger than $\sqrt{\quad}$

- reciprocal: $f(Y) = 1/Y$ is even stronger

Exercise: How is $\text{Var}(Y)$ related to $E(Y)$
If reciprocal transformation is
appropriate.

What if data includes 0's or negative numbers?
and log transformation is appropriate
Use $\log(Y + K)$, K a constant
of your choice