

Flexible Regression and Smoothing

Continuous Distributions

Mikis Stasinopoulos, Rob Rigby, Gillian Heller,
and Fernanda De Bastiani

BCAM – Basque Center for Applied Mathematics, Bilbao, Basque Country - Spain, October 2019

- 1 Types of continuous distributions
- 2 Distributions on \mathcal{R}
 - Location and scale family
 - Two parameter on \mathcal{R}
 - Three parameter on \mathcal{R}
 - Four parameter on \mathcal{R}
- 3 Choosing an appropriate distribution
- 4 Fitting a distribution
- 5 Example: DAX returns data
- 6 End

Types of continuous distributions

$(-\infty, \infty)$: real line \mathbb{R} [Chapter 4 of Rigby *et al.* (2019)]

$(0, \infty)$: positive real line \mathbb{R}^+ [Chapter 5]

$(0, 1)$ real line on interval $\mathbb{R}_{(0,1)}$ (not containing zero or one)
[Chapter 6]

Location and scale family

For these distributions, if

$$Y \sim \mathcal{D}(\mu, \sigma, \nu, \tau)$$

then

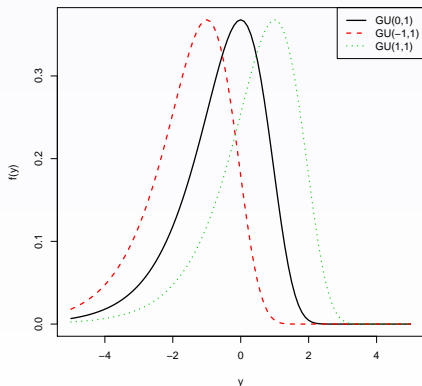
$$\varepsilon = (Y - \mu)/\sigma \sim \mathcal{D}(0, 1, \nu, \tau),$$

i.e. $Y = \mu + \sigma\varepsilon$, so Y is a scaled and shifted version of the random variable ε .

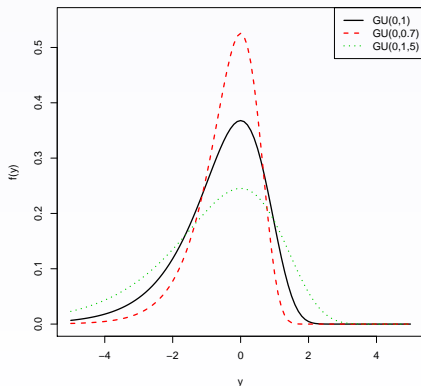
All $(-\infty, \infty)$ distributions in GAMLSS except $\text{exGAUS}(\mu, \sigma, \tau)$ are *location-scale* families

Location-scale family

(a) changing the location parameter



(b) changing the scale parameter



Two parameter on \mathcal{R}

Two parameter distributions on $(-\infty, \infty)$

Gumbel $GU(\mu, \sigma)$ left skew

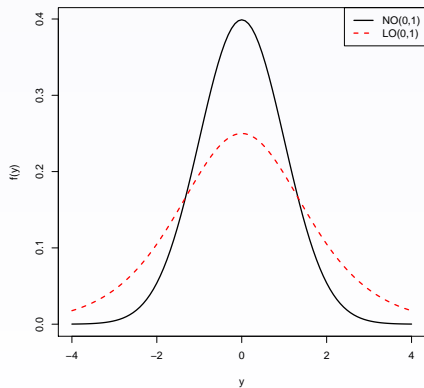
Logistic $LO(\mu, \sigma)$ leptokurtic

Normal $NO(\mu, \sigma)$ and $NO2(\mu, \sigma)$

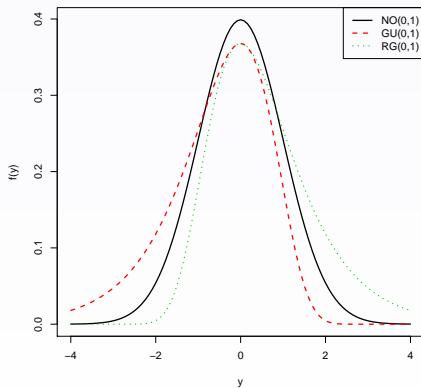
Reverse Gumbel $RG(\mu, \sigma)$ right skew

Two parameter on \mathcal{R}

(a) NO and LO distributions



(b) NO, GU, and RG distributions



Three parameter on \mathcal{R}

exponential Gaussian : $\text{exGAUS}(\mu, \sigma, \nu)$ for modelling right skew data,

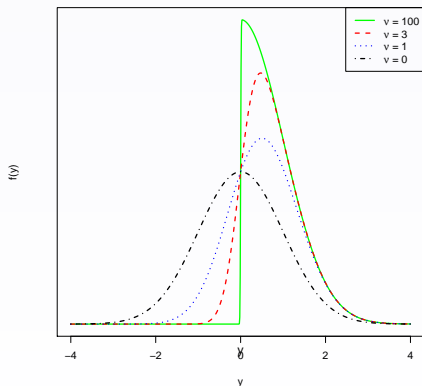
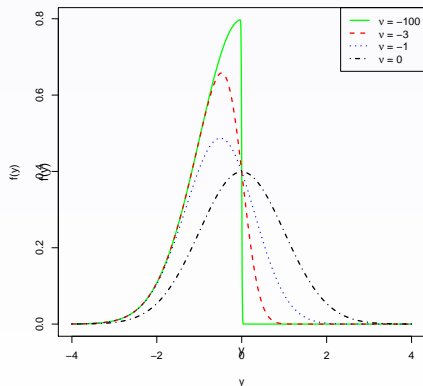
normal family: $\text{NOF}(\mu, \sigma, \nu)$ for modelling mean and variance relationships following the power law;

power exponential: $\text{PE}(\mu, \sigma, \nu)$ and $\text{PE2}(\mu, \sigma, \nu)$ for modelling leptokurtotic and platykurtotic data;

***t* family**: $\text{TF}(\mu, \sigma, \nu)$ and $\text{TF2}(\mu, \sigma, \nu)$ for modelling leptokurtotic data;

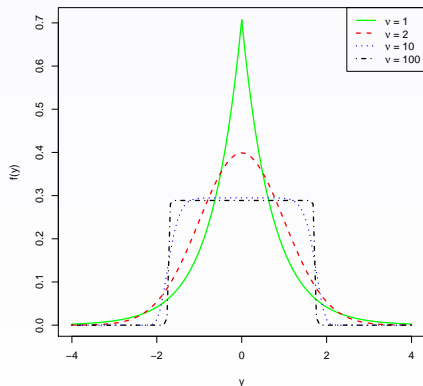
skew normal: $\text{SN1}(\mu, \sigma, \nu)$ and $\text{SK2}(\mu, \sigma, \nu)$ for modelling skewness in data.

Three parameter on \mathcal{R} : skew normal type 1

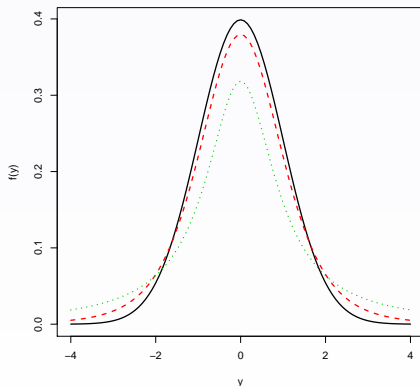


Three parameter on \mathfrak{R} : PE and TF distributions

(a) power exponential



(b) the t-family



Four parameter on \mathfrak{R}

exponential generalised beta type 2: $\text{EGB2}(\mu, \sigma, \nu, \tau)$ skewness and leptokurtosis;

generalised t : $\text{GT}(\mu, \sigma, \nu, \tau)$ kurtosis;

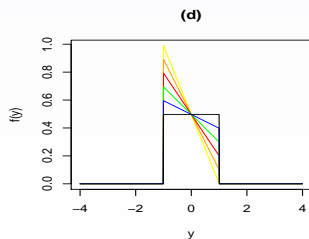
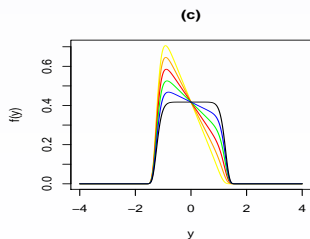
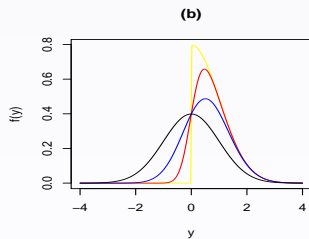
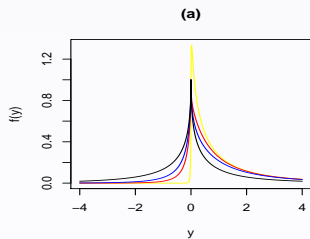
Johnson's SU: $\text{JSU}(\mu, \sigma, \nu, \tau)$ and $\text{JSUo}(\mu, \sigma, \nu, \tau)$ skewness and leptokurtosis;

normal-exponential- t : $\text{NET}(\mu, \sigma, \nu, \tau)$, robustly location and scale;

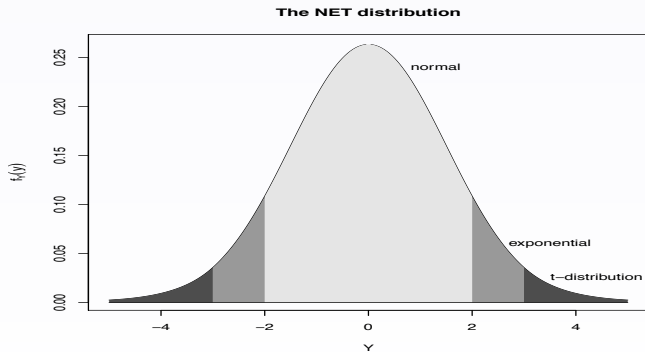
skew exponential power: $\text{SEP1}(\mu, \sigma, \nu, \tau)$, $\text{SEP2}(\mu, \sigma, \nu, \tau)$, $\text{SEP3}(\mu, \sigma, \nu, \tau)$ and $\text{SEP4}(\mu, \sigma, \nu, \tau)$ skewness and leptoplity;

sinh-arcsinh: $\text{SHASH}(\mu, \sigma, \nu, \tau)$, $\text{SHASHo}(\mu, \sigma, \nu, \tau)$ and $\text{SHASHo2}(\mu, \sigma, \nu, \tau)$ skewness and leptoplity;

skew t : $\text{ST1}(\mu, \sigma, \nu, \tau)$, $\text{ST2}(\mu, \sigma, \nu, \tau)$, $\text{ST3}(\mu, \sigma, \nu, \tau)$, $\text{ST4}(\mu, \sigma, \nu, \tau)$, $\text{ST5}(\mu, \sigma, \nu, \tau)$ and $\text{SST}(\mu, \sigma, \nu, \tau)$ skewness and leptokurtosis.

Example of four parameter on \mathfrak{R} : SEP1

The NET distribution



Summary of GAMLSS family distributions defined on $(-\infty, +\infty)$

Distributions	family	no par.	skewness	kurtosis
Exp. Gaussian	exGAUS	3	positive	-
Exp.G. beta 2	EGB2	4	both	lepto
Gen. t	GT	4	(symmetric)	lepto
Gumbel	GU	2	(negative)	-
Johnson's SU	JSU, JSU _o	4	both	lepto
Logistic	LO	2	(symmetric)	(lepto)
Normal-Expon.- t	NET	2,(2)	(symmetric)	lepto
Normal	NO-NO2	2	(symmetric)	
Normal Family	NOF	3	(symmetric)	(meso)

Summary of GAMLSS family distributions defined on $(-\infty, +\infty)$

Power Expon.	PE-PE2	3	(symmetric)	both
Reverse Gumbel	RG	2	positive	
Sinh Arcsinh	SHASH,	4	both	both
Skew Exp. Power	SEP1-SEP4	4	both	both
Skew t	ST1-ST5, SST	4	both	lepto
t Family	TF	3	(symmetric)	lepto

Choosing an appropriate distribution

global deviance

$$\text{GDEV} = -2\hat{\ell}(\hat{\boldsymbol{\theta}}) = -2 \sum_{i=1}^n \log f(y_i | \hat{\boldsymbol{\theta}}) \quad (1)$$

generalized Akaike information criterion

$$\text{GAIC}(k) = \text{GDEV} + k \cdot df \quad (2)$$

special cases

$$\text{AIC} = \text{GDEV} + 2 df$$

$$\text{SBC} = \text{GDEV} + \log n \cdot df .$$

Choosing an appropriate distribution

prediction global deviance

$$\text{VDEV} = \text{TDEV} = -2\tilde{\ell}(\tilde{\boldsymbol{\theta}}) = -2\sum_{i=1}^{n_p} \log f(\tilde{y}_i|\tilde{\boldsymbol{\theta}})$$

$\tilde{\mathbf{y}}$ indicates the response variable values at the validation (or test) sample,

$$\tilde{\boldsymbol{\theta}} = (\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\sigma}}, \tilde{\boldsymbol{\nu}}, \tilde{\boldsymbol{\mu}})$$

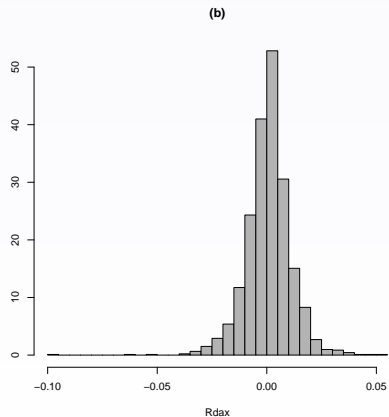
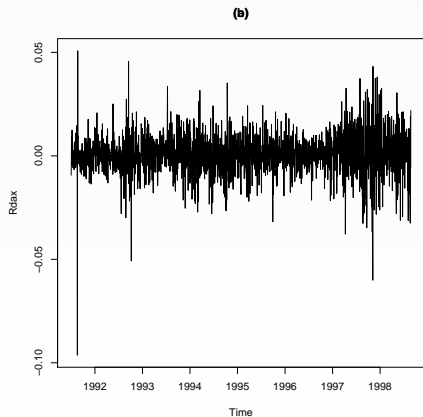
are the fitted parameters for the validation (or test) sample (which in **gamlss** can be obtained using the functions `predict()`),

How to fit distributions in R

- `optim()` or `mle()` R functions requiring initial parameter values
- `gamlssML()` fits a distribution (using MLE) on the response¹
 - `gamlss()`: fits a distribution on the response using RS or CG or mixed algorithms,
 - `histDist()` fits a distribution using `gamlssML()`, and plots a histogram of the response with the fitted distribution¹
 - `fitDist()` fits a set of distributions to the response and chooses the one with the smallest GAIC¹
- `chooseDist()` fits a set of distributions on a fitted model and chooses the one with the smallest GAIC

¹with no explanatory variables

DAX returns data



DAX returns data: `fitDist()`

```
f1 <- fitDist(Rdax, type="realline")
```

```
f1$fits
```

GT	SEP2	PE2	PE	JSU	JSUo	SEP1	SEP4
-11967.6	-11965.5	-11962.4	-11962.4	-11961.4	-11961.4	-11961.3	-11961.2
SEP3	TF2	TF	ST4	ST1	EGB2	ST5	ST2
-11960.8	-11960.6	-11960.6	-11960.1	-11959.8	-11959.7	-11959.4	-11959.2
ST3	SST	SHASH	SHASHo2	SHASHo	LO		
-11958.8	-11958.8	-11957.5	-11956.2	-11956.2	-11932.1		
NET	SN1	SN2	NO	exGAUS	GU	RG	
-11919.1	-11759.8	-11739.1	-11733.2	-11733.1	-11262.1	-10249.8	

```
f1$failed
```

```
list()
```

DAX returns data: chooseDist()

```
t1 <- chooseDist(m1, type = "realline")
```

```
minimum GAIC(k= 2 ) family: GT
```

```
minimum GAIC(k= 3.84 ) family: GT
```

```
minimum GAIC(k= 7.53 ) family: PE2
```

```
t1
```

```
          2          3.84          7.53
```

```
NO -11733.21 -11729.53 -11722.15
```

```
GU -11262.12 -11258.44 -11251.06
```

```
RG -10249.88 -10246.20 -10238.82
```

```
....
```

```
ST5 -11959.48 -11952.12 -11937.36
```

```
SST -11958.87 -11951.51 -11936.75
```

```
GT -11967.68 -11960.32 -11945.56
```

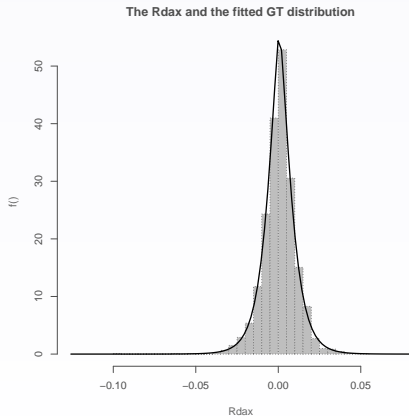
DAX returns data: chooseDist()

```

getOrder(t1,1)[1:6]
      GT      SEP2      PE2      PE      JSU      JSUo
-11967.68 -11965.53 -11962.46 -11962.46 -11961.48 -11961.48
mf <- update(m1, family="GT")
summary(mf)
fitted(f1,"mu")[1]
[1] 0.0006838042
fitted(f1,"sigma")[1]
[1] 0.009610674
fitted(f1,"nu")[1]
[1] 5.96016
fitted(f1, "tau")[1]
[1] 1.417194
fh<-histDist(Rdax, family=GT, nbins=30, line.col="black")

```

DAX returns data fitted model



END

for more information see

www.gamlss.org