

Beyond the exponential family

Eric Pedersen, Gavin Simpson, David Miller
August 6th, 2016

Away from the exponential family

Most glm families (Poisson, Gamma, Gaussian, Binomial) are *exponential families*

Away from the exponential family

Most glm families (Poisson, Gamma, Gaussian, Binomial) are *exponential families*

$$f(x|\theta) \sim \exp\left(\sum_i \eta_i(\theta) T_i(x) - A(\theta)\right)$$

Away from the exponential family

Most glm families (Poisson, Gamma, Gaussian, Binomial) are *exponential families*

$$f(x|\theta) \sim \exp\left(\sum_i \eta_i(\theta) T_i(x) - A(\theta)\right)$$

- Computationally easy

Away from the exponential family

Most glm families (Poisson, Gamma, Gaussian, Binomial) are *exponential families*

$$f(x|\theta) \sim \exp\left(\sum_i \eta_i(\theta) T_i(x) - A(\theta)\right)$$

- Computationally easy
- Has sufficient statistics: easier to estimate parameter variance

Away from the exponential family

Most glm families (Poisson, Gamma, Gaussian, Binomial) are *exponential families*

$$f(x|\theta) \sim \exp\left(\sum_i \eta_i(\theta) T_i(x) - A(\theta)\right)$$

- Computationally easy
- Has sufficient statistics: easier to estimate parameter variance
- ... but it doesn't describe everything

Away from the exponential family

Most glm families (Poisson, Gamma, Gaussian, Binomial) are *exponential families*

$$f(x|\theta) \sim \exp\left(\sum_i \eta_i(\theta) T_i(x) - A(\theta)\right)$$

- Computationally easy
- Has sufficient statistics: easier to estimate parameter variance
- ... but it doesn't describe everything
- mgcv has expanded to cover many new families

Away from the exponential family

Most glm families (Poisson, Gamma, Gaussian, Binomial) are *exponential families*

$$f(x|\theta) \sim \exp\left(\sum_i \eta_i(\theta) T_i(x) - A(\theta)\right)$$

- Computationally easy
- Has sufficient statistics: easier to estimate parameter variance
- ... but it doesn't describe everything
- mgcv has expanded to cover many new families
- Lets you model a much wider range of scenarios with smooths

What we'll cover

- “Counts”: Negative binomial and Tweedie distributions
- Modelling proportions with the Beta distribution
- Robust regression with the Student's t distribution
- Ordered and unordered categorical data
- Multivariate normal data
- Modelling extra zeros with zero-inflated and adjusted families
- *NOTE*: All the distributions we're covering here have their own quirks. Read the help files carefully before using them!

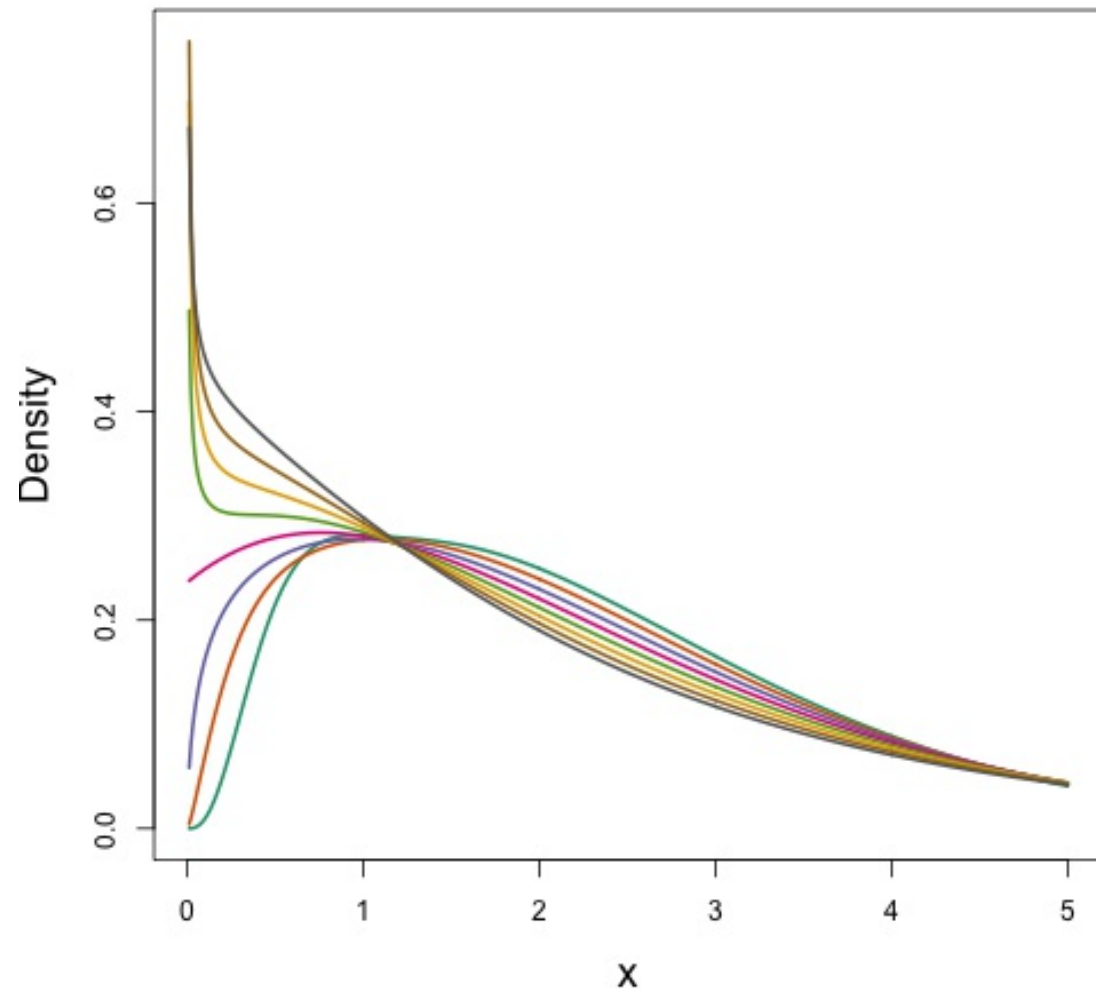
Modelling "counts"

Counts and count-like things

- Response is a count (not always integer)
- Often, it's mostly zero (that's complicated)
- Could also be catch per unit effort, biomass etc
- Flexible mean-variance relationship

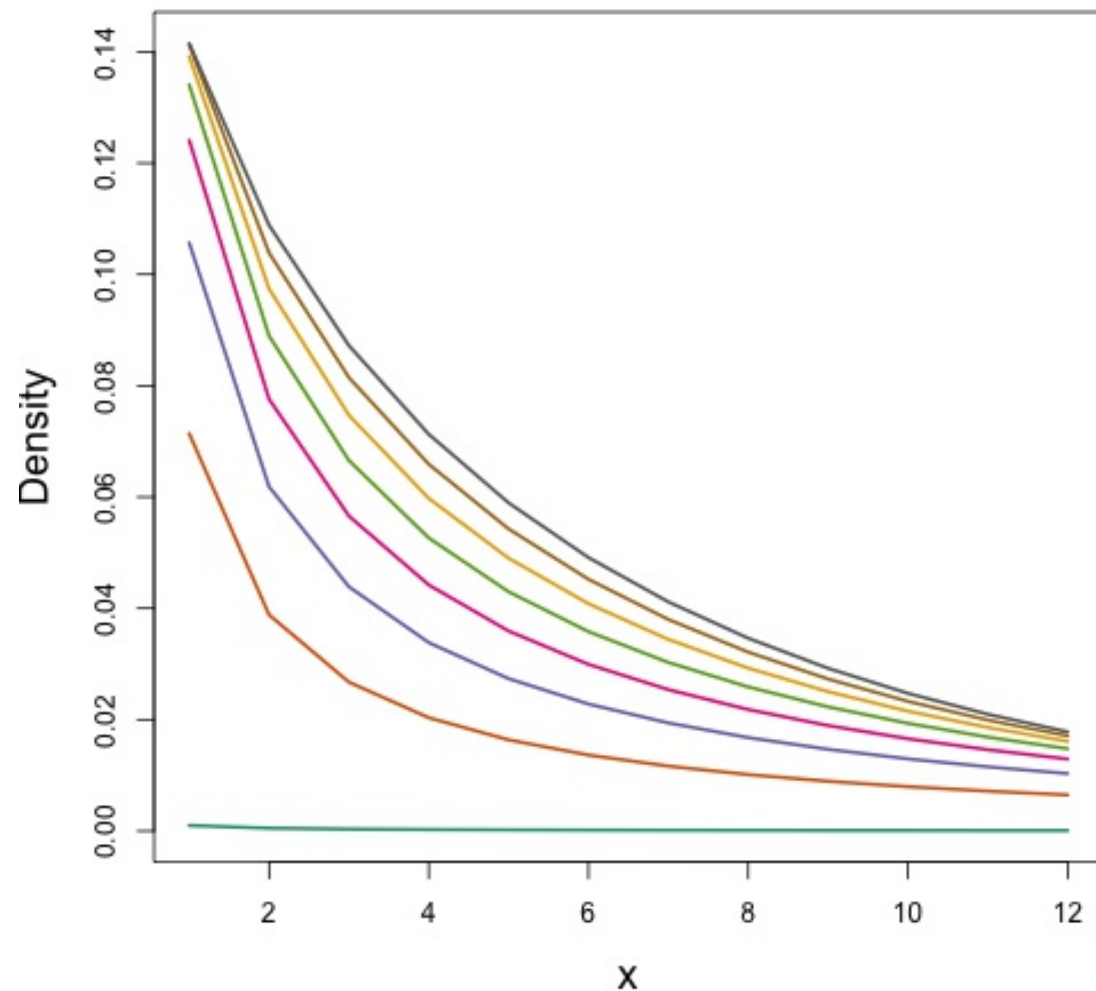


Tweedie distribution



- $\text{Var}(\text{count}) = \varphi(\text{count})^q$
- Common distributions are sub-cases:
 - $q = 1 \Rightarrow \text{Poisson}$
 - $q = 2 \Rightarrow \text{Gamma}$
 - $q = 3 \Rightarrow \text{Normal}$
- We are interested in $1 < q < 2$
- (here $q = 1.2, 1.3, \dots, 1.9$)
- `tw()`

Negative binomial



- $\text{Var}(\text{count}) = (\text{count}) + \kappa(\text{count})^2$
- Estimate κ
- Is quadratic relationship a “strong” assumption?
- Similar to Poisson:
 $\text{Var}(\text{count}) = (\text{count})$
- `nb()`

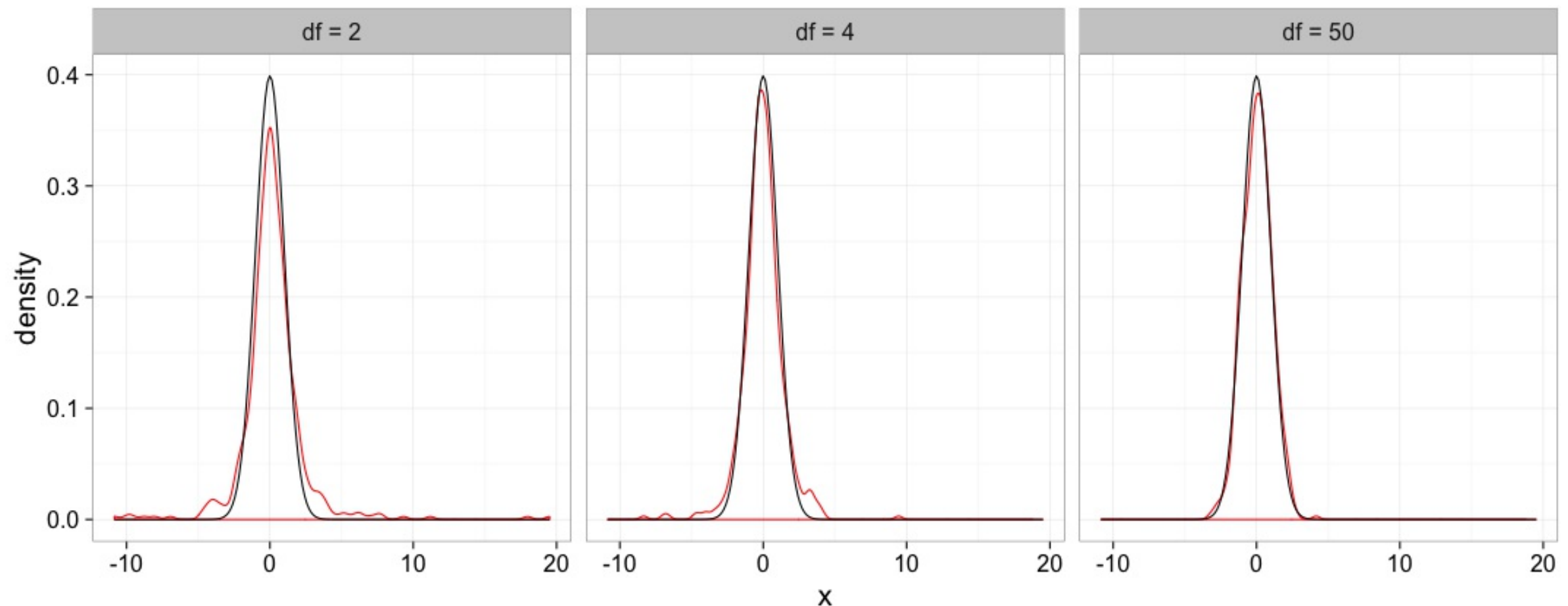
Modelling proportions

The Beta distribution

Modelling outliers

The student-t distribution

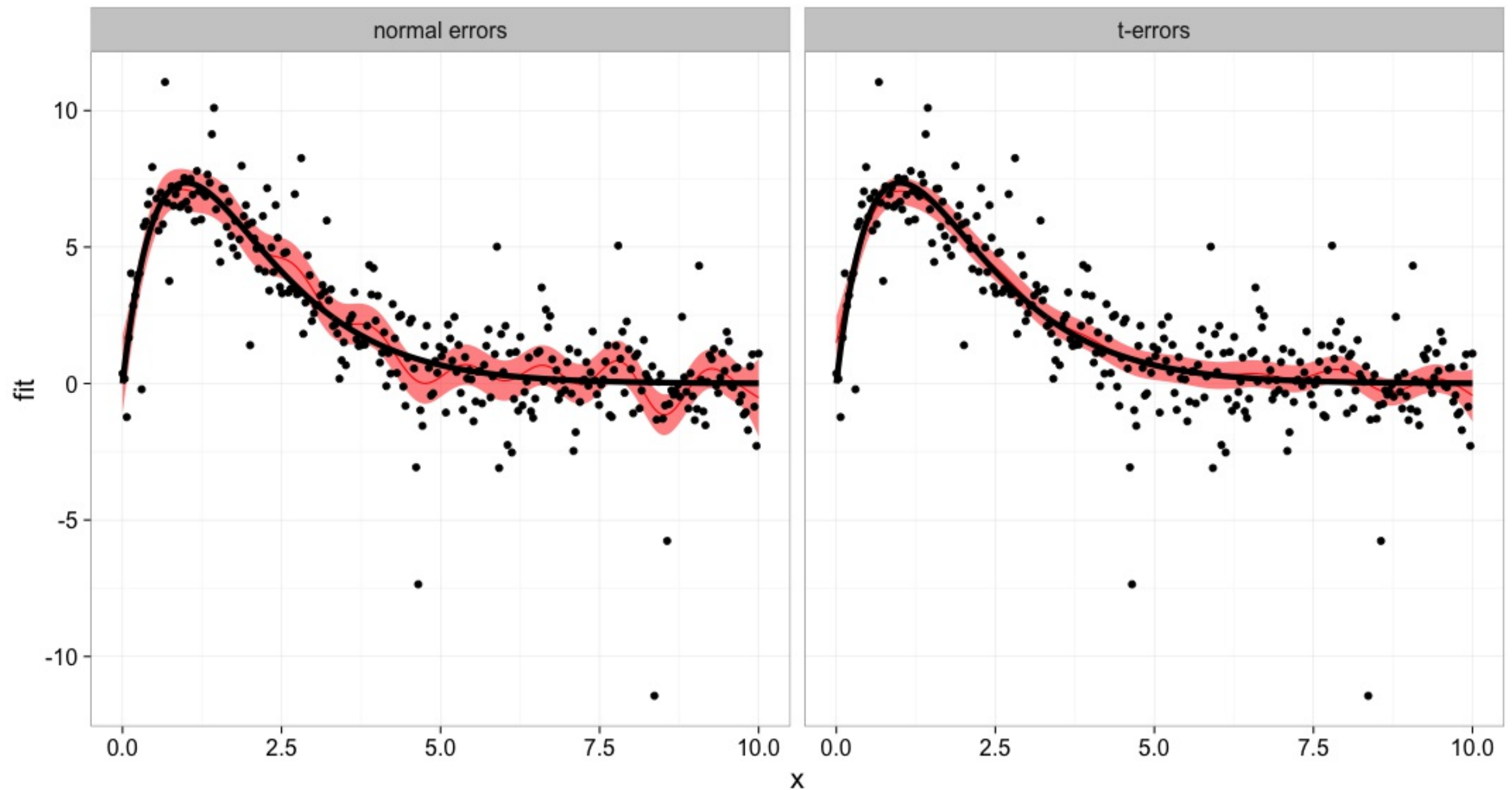
- Models continuous data w/ longer tails than normal
- Far less sensitive to outliers
- Has one extra parameter: df.
- bigger df: t dist approaches normal



The student-t distribution: Usage

```
set.seed(4)
n=300
dat = data.frame(x=seq(0,10,length=n))
dat$f = 20*exp(-dat$x)*dat$x
dat$y = 1*rt(n,df = 3) + dat$f
norm_mod = gam(y~s(x,k=20), data=dat,
family=gaussian(link="identity"))
t_mod = gam(y~s(x,k=20), data=dat, family=scat(link="identity"))
```

The student-t distribution: Usage



The student-t distribution: Usage

Family: Scaled $t(2.976, 0.968)$
Link function: identity

Formula:
 $y \sim s(x, k = 20)$

Parametric coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	2.02664	0.06853	29.57	<2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Approximate significance of smooth terms:

	edf	Ref.df	Chi.sq	p-value	
s(x)	13.27	15.71	1221	<2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

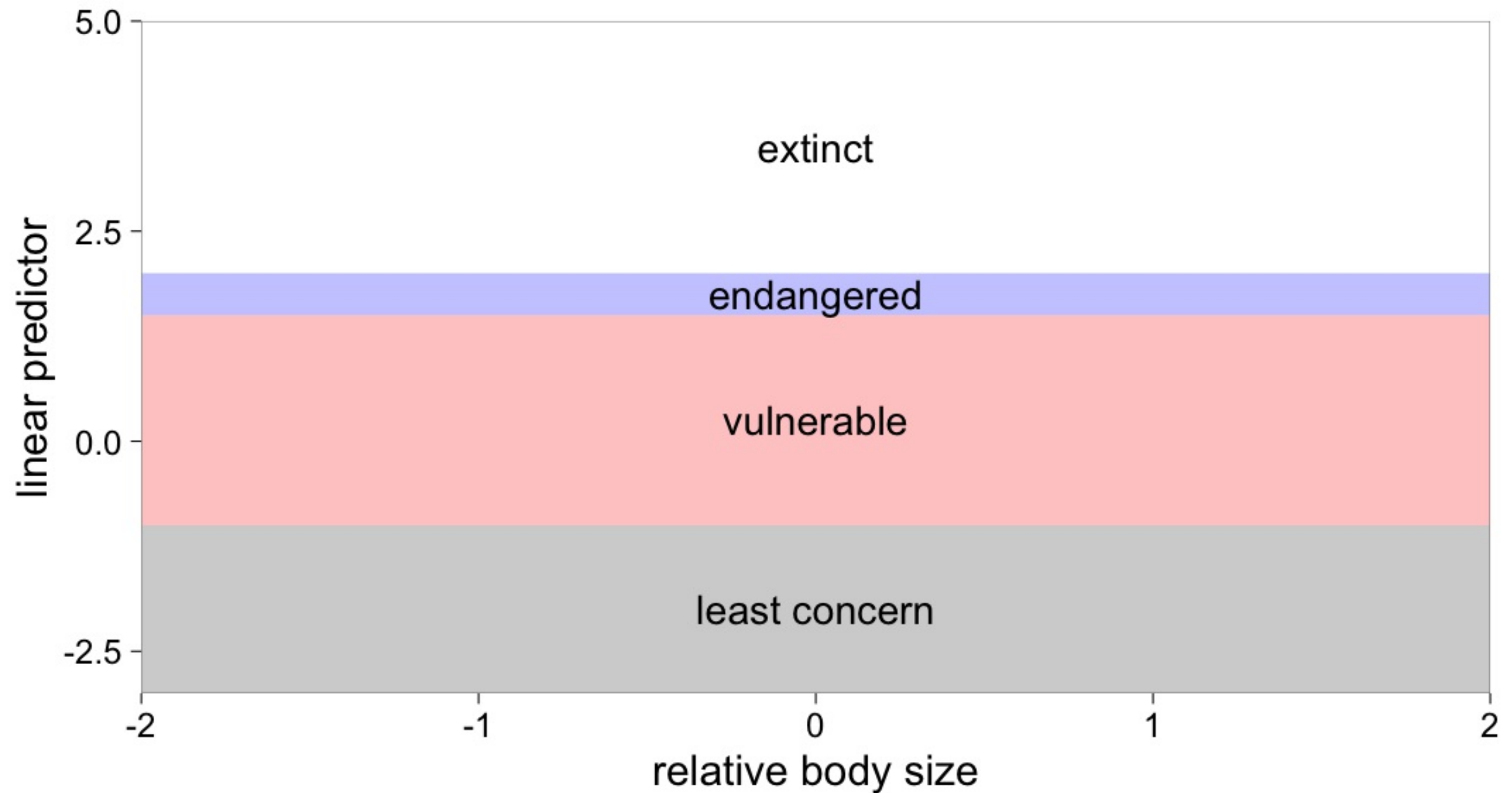
R-sq.(adj) = 0.695 Deviance explained = 63.1%
-REML = 546.75 Scale est. = 1 n = 300

Modelling multi-dimensional data

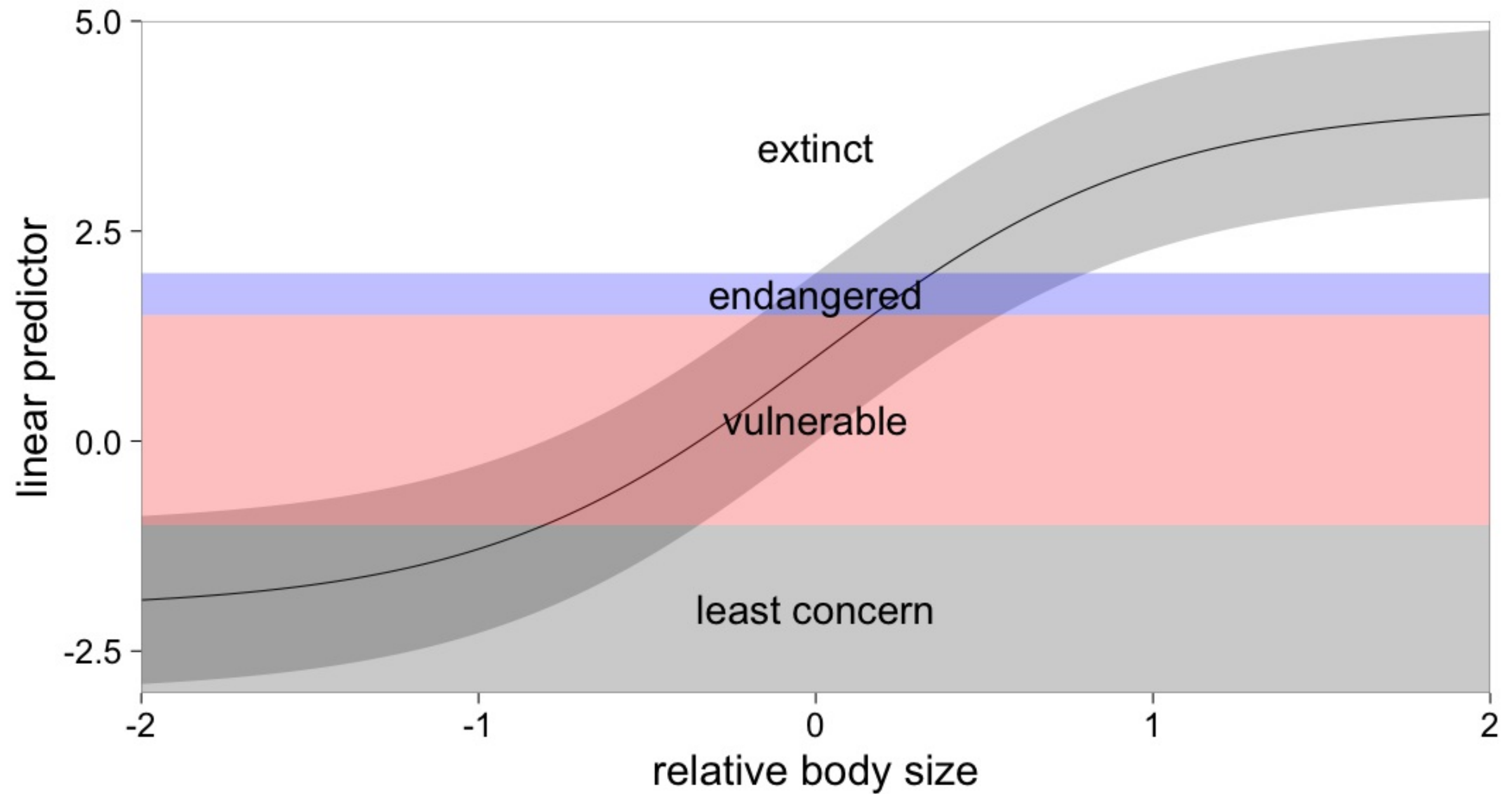
Ordered categorical data

- Assumes data are in discrete categories, and categories fall in order
- e.g.: conservation status: “least concern”, “vulnerable”, “endangered”, “extinct”
- fits a linear latent model using covariates, w/ threshold for each level
- First cut-off always occurs at -1

Ordered categorical data

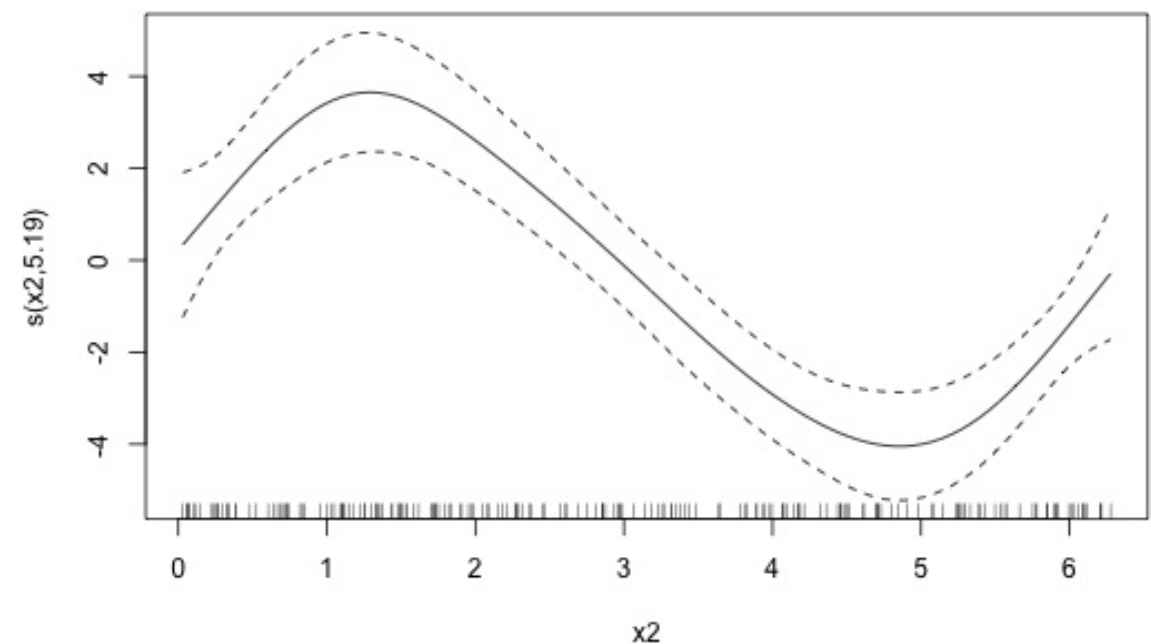
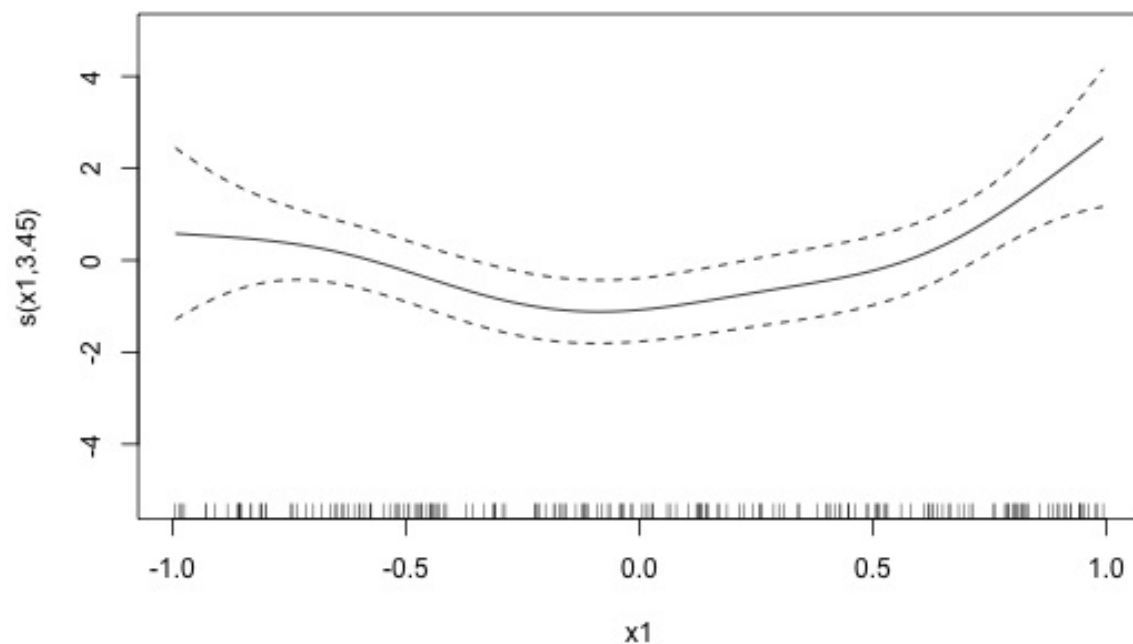


Ordered categorical data



Using ocat

```
n= 200  
dat = data.frame(x1 = runif(n,-1,1),x2=2*pi*runif(n))  
dat$f = dat$x1^2 + sin(dat$x2)  
dat$y_latent = dat$f + rnorm(n,dat$f)  
dat$y = ifelse(dat$y_latent<0,1, ifelse(dat$y_latent<0.5,2,3))  
ocat_model = gam(y~s(x1)+s(x2), family=ocat(R=3),data=dat)  
plot(ocat_model,page=1)
```



Using ocat

```
summary(ocat_model)
```

```
Family: Ordered Categorical(-1,-0.09)
Link function: identity
```

```
Formula:
y ~ s(x1) + s(x2)
```

```
Parametric coefficients:
```

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	0.5010	0.2792	1.794	0.0727

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Approximate significance of smooth terms:
```

	edf	Ref.df	Chi.sq	p-value	
s(x1)	3.452	4.282	18.67	0.00133	**
s(x2)	5.195	6.270	84.34	1.09e-15	***

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Deviance explained = 57.7%
```

```
-REML = 97.38 Scale est. = 1
```

```
n = 200
```

Using ocat

```
Error in eval(expr, envir, enclos) : invalid subscript type  
'double'
```