# Computational Statistics-Report

## David Niederkofler, Erlend Lokna

## 2022-12-16

```
mydata<-read.table("Report2_Dataset.txt", header=FALSE)
```

# Statistical Analysis of Covariates

It is important to mention the use of notation before we proceed. We will in this section use the notation $\mathbf{x} = \{x_1, x_2, \cdots, x_n\}$ as our covariate vector.

## Ascicles

### 1.1 Model selection

Since the Ascicles - covariate has a 0-1 outcome we can assume that it is Bernoulli distributed with parameter $\theta$. A natural conjugate prior for the Bernoulli distribution is the Beta distribution. The posterior beta distribution for the parameter is given by

$$Beta(\theta|a + \sum_{i=1}^{n} x_i, b + n - \sum_{i+1}^{n} x_i)$$

### 1.2 Results

The following results where found using the posterior beta distribution with a=1 and b=1 (Uniform distributed) for the ascicles data:

```
## Posterior mean:  0.08227848
```

```
## Posterior mode:  0.07961783
```

```
## Centered 95% Confidence Interval: [ 0.05235453 , 0.1119428 ]
```

With the following HPD interval:

```
##      lower      upper
## 0.05059244 0.10956674
## attr(,"credMass")
## [1] 0.95
```

## 1. Sex

The sex of the patients is encoded in a binary variable, where 0 means *male* and 1 means *female*.

### 1.1 Model selection

We assume a Bernoulli model $Ber(\theta)$ for the sex of the patient conditional on one parameter $\theta$, the probability of the patient to be female. The density function is given by

$$f(x|\theta) = \theta^x (1-\theta)^{1-x}, \tag{1}$$

where $x \in \{0, 1\}$. As a prior distribution for $\theta$ we use the natural conjugate family of the Bernoulli distribution, namely the Beta distribution, $Beta(a, b)$, with two shape parameters $a = b = 2$ to give more weight to the middle of the interval $[0, 1]$, knowing how females and males are represented in the general population. The density is given by

$$h(\theta) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1}(1-\theta)^{b-1}, \tag{2}$$

for $\theta \in [0, 1]$.

### 1.2 Results

From the given dataset we get the sample size $n$ and the sum of the observations $s$:

```
n<-length(mydata$V6[!is.na(mydata$V6)])
s<-sum(mydata$V6)
n
```

```
## [1] 312
s
```

```
## [1] 276
```

Therefore the posterior distribution is $Beta(2 + s, 2 + n - s)$, which turns out to be $Beta(278, 38)$. From that we get

```
## Posterior mean:  0.8797468
```

```
## Posterior mode:  0.8821656
```

```
## Centered 95% Confidence Interval: [ 0.8417454 , 0.9132003 ]
```

And the HPD confidence Interval calculates to:

```
tst<-rbeta(1e5,278,38)
hdi(tst)
```

```
##     lower     upper
## 0.8427747 0.9139823
## attr(,"credMass")
## [1] 0.95
```

## 2. Spiders

The presence of spiders is encoded in a Binary variable, where 1 means spiders are present.

### 2.1 Model selection

We assume a Bernoulli model $Ber(\theta)$ for the presence of spiders in patients conditional on one parameter $\theta$, the probability of the presence of spiders in the patient. The density function is given as stated earlier. As a prior distribution for $\theta$ we use the natural conjugate family of the Bernoulli distribution, namely the Beta distribution, $Beta(a, b)$, with two shape parameters $a = b = 1$, because we have no prior information. The density is given as above.

### 2.2 Results

From the given dataset we get the sample size $n$ and the sum of the observations $s$:

```
n<-length(mydata$V9[!is.na(mydata$V9)])
s<-sum(mydata$V9)
n
```

```
## [1] 312
```

```
s
```

```
## [1] 90
```

Therefore the posterior distribution is $Beta(1+s, 1+n-s)$, which turns out to be $Beta(91, 223)$. From that we get

```
## Posterior mean:  0.2911392
```

```
## Posterior mode:  0.2898089
```

```
## Centered 95% Confidence Interval: [ 0.2410228 , 0.341131 ]
```

And the HPD confidence interval calculates to:

```
tst<-rbeta(1e5,91,223)
hdi(tst)
```

```
##     lower     upper
## 0.2401626 0.3404193
## attr(,"credMass")
## [1] 0.95
```

## 3. Hepatomegaly

The presence of hepatomegaly is encoded in a Binary variable, where 1 means hepatomegaly is present.

### 3.1 Model selection

We assume a Bernoulli model $Ber(\theta)$ for the presence of hepatomegaly in the patient, conditional on one parameter $\theta$, the probability of the presence of hepatomegaly in the patient. The density function is given as stated earlier. As a prior distribution for $\theta$ we use the natural conjugate family of the Bernoulli distribution, namely the Beta distribution, $Beta(a, b)$, with two shape parameters $a = b = 1$, because we have no prior information. The density is given as above.

### 3.2 Results

From the given dataset we get the sample size $n$ and the sum of the observations $s$:

```
n<-length(mydata$V8[!is.na(mydata$V8)])
s<-sum(mydata$V8)
n
```

```
## [1] 312
```

```
s
```

```
## [1] 160
```

Therefore the posterior distribution is $Beta(1+s, 1+n-s)$, which turns out to be $Beta(161, 153)$. From that we get

```
## Posterior mean:  0.5126582
```

```
## Posterior mode:  0.5127389
```

```
## Centered 95% Confidence Interval: [ 0.4575015 , 0.5678225 ]
```

And the HPD confidence interval calculates to:

```
tst<-rbeta(1e5,161,153)
hdi(tst)
```

```
##     lower     upper
## 0.4580133 0.5679745
## attr(,"credMass")
## [1] 0.95
```

## 4. Histologic stage

The Histologic stage of the disease is a number in $\{1, 2, 3, 4\}$, where the stage increases with severeness. We will give here the frequencies of the stages in the dataset.

```
##   1   2   3   4
##  16  67 120 109
```

We see that, most patients have been diagnosed in the last to stages of the disease.

## 5. Age

The age of the patient in days.

### 5.1 Model selection

The data seems to follow a poisson distribution $Poi(\lambda)$. Using the non informative Jeffreys prior, we can derive that the posterior for the parameter $\lambda$ is Gamma distributed.

$$\theta|x \sim Gamma(\alpha = \frac{1}{2} + \sum_{i=1}^{n} x_i, \beta = n)$$

### 5.2 Results

$$\theta|x \sim Gamma(\frac{1}{2} + s, n)$$

```
## posterior distribution: Gamma( 5700067 , 312 )
## mean: 18269.44
## variance: 58.55591
## HPD intervall:
```

```
##     lower     upper
## 18254.45 18284.53
## attr(,"credMass")
## [1] 0.95
```

## 6. Cholesterol

### 6.1 Model selection

We assume that the data is sampled from a poisson, $Poi(\lambda)$, distribution, and we use the non informative Jeffreys prior for the rate parameter in the bayesian analysis.

$$\mathbf{x} \sim Poi(\lambda)$$

$$h(\lambda) \propto \lambda^{-\frac{1}{2}}$$

**6.2 Results**

```
## posterior distribution: Gamma( 104941.5 , 312 )
## mean: 336.351
## variance: 1.078048
## HPD intervall:

##     lower    upper
## 334.3066 338.4027
## attr(,"credMass")
## [1] 0.95
```

## 7. Urine

**7.1 Model selection**

$$\mathbf{x} \sim Poi(\lambda)$$

$$h(\lambda) \propto \lambda^{-\frac{1}{2}}$$

**7.2 Results**

```
## posterior distribution: Gamma( 30271.5 , 312 )
## mean: 97.02404
## variance: 0.3109745
## HPD intervall:

##    lower    upper
## 95.95828 98.14324
## attr(,"credMass")
## [1] 0.95
```

## 8 SGOT

**8.1 Model selection**

$$\mathbf{x} \sim Poi(\lambda)$$

$$h(\lambda) \propto \lambda^{-\frac{1}{2}}$$

**8.2 Results**

```
poisson_jeffrey(mydata$V16)
```

```
## posterior distribution: Gamma( 38238.08 , 312 )
## mean: 122.5579
## variance: 0.3928139
## HPD intervall:

##    lower    upper
## 121.3332 123.7836
## attr(,"credMass")
## [1] 0.95
```

## 9. Plateles

### 9.1 Model selection

$$\mathbf{x} \sim Poi(\lambda)$$

$$h(\lambda) \propto \lambda^{-\frac{1}{2}}$$

### 9.2 Results

```
poisson_jeffrey(mydata$V18)
```

```
## posterior distribution: Gamma( 80676.5 , 312 )
## mean: 258.5785
## variance: 0.8287773
## HPD intervall:
```

```
##     lower    upper
## 256.8077 260.3794
## attr(,"credMass")
## [1] 0.95
```

## 10. Prothrombin

### 10.1 Model selection

$$\mathbf{x} \sim Poi(\lambda)$$

$$h(\lambda) \propto \lambda^{-\frac{1}{2}}$$

### 10.2 Results

```
poisson_jeffrey(mydata$V19)
```

```
## posterior distribution: Gamma( 3346.9 , 312 )
## mean: 10.72724
## variance: 0.03438219
## HPD intervall:
```

```
##     lower    upper
## 10.35639 11.08303
## attr(,"credMass")
## [1] 0.95
```

## 11. Bilirubin

The serum bilirubin of the patients is given in mg/dl.

### 11.1 Model Selection

We assume by inspecting the histogramm plot,

## Histogram of bilirubin



that the data follows a exponential distribution with parameter $\lambda$. Density is given by

$$f(x|\lambda) = \lambda e^{-\lambda x} \tag{3}$$

As a prior for $\lambda$ we use, the jeffreys non-informative prior, namely: $h(\lambda) \propto \frac{1}{\lambda}$.

### 11.2 Results

From the data we get the number of samples $n$ and the sum of the samples $s$ as

```
## [1] 312
```

```
## [1] 1015.9
```

That means the posterior distribution for $\lambda$ is $Gamma(n, s)$. Which turns out to be $Gamma(312, 1015.9)$. From that we get

```
## Posterior mean:  0.3071168
```

```
## Posterior mode:  0.3061325
```

```
## Centered 95% Confidence Interval: [ 0.2739805 , 0.3421174 ]
```

And the HPD confidence interval calculates to:

```
tst<-rgamma(1e5,312,1015.9)
hdi(tst)
```

```
##     lower     upper
## 0.2729973 0.3412094
## attr(,"credMass")
## [1] 0.95
```

## 12. Albumin

The Albumin is given in mg/dl.

### 12.1 Model selection

**Histogram of albumin**



By the histogram plot of the data we see,

that albumin could be gamma distributed with shape and rate parameters $a$ and $b$. We assume prior independence between $a$ and $b$ and use the marginal prior distributions $Gamma(0.001, 0.001)$ for both of them.

### 12.2 Results

Using OpenBUGS and MCMC methods, we get posterior information about the parameters $a$ and $b$:

```
n<-length(albumin[!is.na(albumin)])
X<-albumin
data1<-list("X","n")
params<-c("a" , "b")
inits<-list(a=1,b=1)
fit1<-bugs(data=data1,inits=list(inits),parameters.to.save=params,"model_albu.txt",n.chains=1, n.iter=2
fit1$summary
```

```
##                mean       sd   2.5%    25%    50%    75%    97.5%
## a          66.32709 4.995273  57.66  62.76  66.30  69.38  77.1405
## b          18.84450 1.424988  16.37  17.83  18.84  19.71  21.9400
## deviance  365.10028 1.880193 363.30 363.70 364.50 365.80 370.1000
```

And the HPD confidence interval for $a$ calculates to:

```
## lower upper
## 56.88 76.12
## attr(,"credMass")
## [1] 0.95
```

whereas the HPD confidence interval for $b$ is

```
## lower upper
## 16.19 21.69
## attr(,"credMass")
## [1] 0.95
```
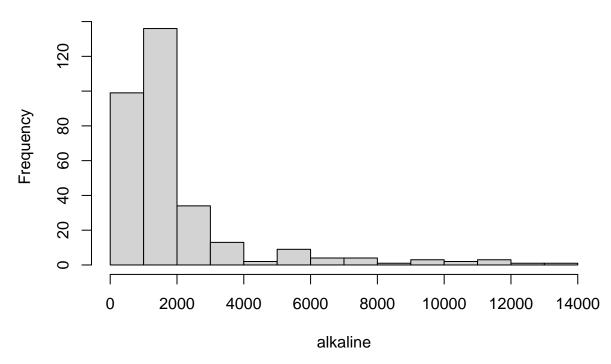
## 13. Alkaline

The data contains the units of alkaline phosphatase per liter of the patients.

### 13.1 Model selection

Since the units of alkaline per liter are integers, we assume that it is a counting process. Therefore we want to assume, that the data is poisson distributed conditional on one parameter $\lambda$. The histogram plot justifies our as-

**Histogram of alkaline**



sumptions:
The density function of a single obervation is given as

$$f(x|\lambda) = e^{-\lambda}\frac{\lambda^x}{x!} \tag{4}$$

As a prior for $\lambda$ we use the natural conjugate prior of the poisson distribution which is the gamma distribution. To not give a lot of prior information, we will use $Gamma(0.001, 0.001)$.

**13.2 Results**

From our data we get the sample size $n$ and the sum $s$ over the sample:

```
## [1] 312
```

```
## [1] 618588.6
```

The posterior distribution for $\lambda$ is given by $Gamma(s + 0.001, n + 0.001)$ which in our case results to $Gamma(618588.601, 312.001)$. This yields to:

```
## Posterior mean:  1982.649
```

```
## Posterior mode:  1982.646
```

```
## Centered 95% Confidence Interval: [ 1977.712 , 1987.593 ]
```

And the HPD confidence interval calculates to:

```r
tst<-rgamma(1e5,618588.601,312.001)
hdi(tst)
```

```
##    lower    upper
## 1977.639 1987.529
## attr(,"credMass")
## [1] 0.95
```

## 14. Triglicerides

Triglicerides of the patients in mg/dl.

**14.1 Model selection**

By the histogram plot of the data we see,

```
## Warning: NAs durch Umwandlung erzeugt
```

# Histogram of triglicerides



that triglicerides could be gamma distributed with shape and rate parameters $a$ and $b$. We assume prior independence between $a$ and $b$ and use the marginal prior distributions $Gamma(0.001, 0.001)$ for both of them.

## 14.2 Results

By OpenBUGS and MCMC methods we get posterior information about the parameters $a$ and $b$:

```
n<-length(triglicerides[!is.na(triglicerides)])
X<-triglicerides
data1<-list("X","n")
params<-c("a" , "b")
inits<-list(a=1,b=1)
fit1<-bugs(data=data1,inits=list(inits),parameters.to.save=params,"model_albu.txt",n.chains=1, n.iter=20
fit1$summary
```

```
##                   mean          sd        2.5%       25%       50%       75%
## a        4.742428e+00 0.410163943 3.988975e+00 4.457e+00 4.7310e+00 5.010e+00
## b        3.775649e-02 0.003443284 3.136975e-02 3.535e-02 3.7645e-02 4.001e-02
## deviance 2.732715e+03 1.974568784 2.731000e+03 2.731e+03 2.7320e+03 2.734e+03
##              97.5%
## a          5.58705
## b          0.04484
## deviance 2738.00000
```

And the HPD confidence interval for $a$ calculates to:

```
## lower upper
## 3.939 5.522
```

```
## attr(,"credMass")
## [1] 0.95
```

whereas the HPD confidence interval for $b$ is

```
##   lower   upper
## 0.03086 0.04419
## attr(,"credMass")
## [1] 0.95
```

# Weibull Survival Analysis

We will use a survival model, based on a hazard function, conditional on regression parameters and dependent on (now assumed) deterministic covariates. The hazard function is given by

$$\lambda(t|\alpha, \beta, \delta) = \delta \alpha t^{\alpha-1} e^{\beta^T z} \tag{5}$$

where $z$ is the covariate vector. As prior distribution for the regression parameters $\beta$ we will use normal distributions with 0 mean and $\sigma^2 = 1000$. For the parameters $\alpha$ and $\delta$ we use $Gamma(0.001, 0.001)$ prior distribution. MCMC methods and OpenBUGS help us to get inference about our parameters:

```
##                    mean           sd         2.5%          25%          50%
## alpha       2.235833e+01 6.084229e-01  2.141000e+01  2.1960e+01  2.216e+01
## beta[1]     3.239920e-02 3.100144e-02 -3.041000e-02  1.1160e-02  3.291e-02
## beta[2]    -2.185738e-02 4.880559e-04 -2.264000e-02 -2.2230e-02 -2.178e-02
## beta[3]    -6.430849e-03 3.216121e-02 -6.926625e-02 -2.9390e-02 -6.401e-03
## beta[4]     1.608548e-02 3.239408e-02 -4.832000e-02 -6.2360e-03  1.568e-02
## beta[5]     2.035891e-02 3.134577e-02 -4.080000e-02 -1.6785e-04  2.009e-02
## beta[6]     4.423746e-02 3.140013e-02 -1.719000e-02  2.3050e-02  4.433e-02
## beta[7]     2.437797e-02 3.144015e-02 -4.007000e-02  2.3915e-03  2.488e-02
## beta[8]     2.739440e-01 4.482108e-02  1.765000e-01  2.4640e-01  2.744e-01
## beta[9]    -1.269871e-02 3.206366e-03 -1.626000e-02 -1.5330e-02 -1.474e-02
## beta[10]    1.009687e-02 3.382378e-02 -5.969000e-02 -1.2275e-02  1.131e-02
## beta[11]    1.728432e-01 1.207354e-02  1.439000e-01  1.6340e-01  1.759e-01
## beta[12]    4.824459e-03 4.312141e-04  4.199000e-03  4.3420e-03  5.041e-03
## beta[13]    1.419207e-01 3.336608e-03  1.374000e-01  1.3860e-01  1.417e-01
## beta[14]   -4.616947e-02 8.579027e-03 -6.566000e-02 -4.9190e-02 -4.478e-02
## beta[15]    3.005904e-02 1.426182e-03  2.736000e-02  2.9190e-02  2.996e-02
## beta[16]    2.135307e-01 3.595282e-02  1.290000e-01  1.9320e-01  2.185e-01
## beta[17]    1.368091e-01 3.612292e-02  6.508000e-02  1.1360e-01  1.374e-01
## delta       1.181189e+04 3.510970e+03  5.973975e+03  9.3015e+03  1.146e+04
## deviance    4.632292e+04 1.173365e+03  4.446000e+04  4.5480e+04  4.608e+04
##                     75%         97.5%
## alpha       2.28300e+01  2.337000e+01
## beta[1]     5.47325e-02  9.243525e-02
## beta[2]    -2.15300e-02 -2.105000e-02
## beta[3]     1.58800e-02  5.571200e-02
## beta[4]     3.93125e-02  7.820250e-02
## beta[5]     4.10450e-02  8.198575e-02
## beta[6]     6.55200e-02  1.067025e-01
## beta[7]     4.52450e-02  8.684025e-02
## beta[8]     3.04700e-01  3.617000e-01
## beta[9]    -1.08800e-02 -6.208000e-03
## beta[10]    3.37500e-02  7.470175e-02
## beta[11]    1.83700e-01  1.867000e-01
```

```
## beta[12]   5.10700e-03   5.408000e-03
## beta[13]   1.45300e-01   1.470000e-01
## beta[14]  -3.90700e-02  -3.551000e-02
## beta[15]   3.08800e-02   3.251000e-02
## beta[16]   2.40800e-01   2.720000e-01
## beta[17]   1.61325e-01   2.056050e-01
## delta      1.38900e+04   1.967100e+04
## deviance   4.72200e+04   4.826000e+04
```

By applying Heidelberg and Welchs method to decide whether the simulated values from the markov chain come from its stationary distribution we get

```
##
##           Stationarity start     p-value
##           test           iteration
## alpha     passed     1           0.3233
## beta[1]   passed     1           0.4563
## beta[2]   passed     1           0.2782
## beta[3]   passed     1           0.4992
## beta[4]   passed     1           0.5484
## beta[5]   passed     1           0.4785
## beta[6]   passed     1           0.4974
## beta[7]   passed     1           0.2412
## beta[8]   passed     1           0.5576
## beta[9]   passed     1           0.1309
## beta[10]  passed     1           0.9451
## beta[11]  passed     1           0.1364
## beta[12]  passed     1           0.3154
## beta[13]  passed     1           0.5834
## beta[14]  passed     1           0.4611
## beta[15]  passed     1           0.1357
## beta[16]  passed     1           0.7140
## beta[17]  passed     1           0.0904
## delta     passed     1           0.6907
##
##           Halfwidth Mean      Halfwidth
##           test
## alpha     passed      2.24e+01 1.69e-02
## beta[1]   passed      3.24e-02 8.30e-04
## beta[2]   passed     -2.19e-02 1.35e-05
## beta[3]   failed     -6.43e-03 8.61e-04
## beta[4]   passed      1.61e-02 8.98e-04
## beta[5]   passed      2.04e-02 8.69e-04
## beta[6]   passed      4.42e-02 8.70e-04
## beta[7]   passed      2.44e-02 8.71e-04
## beta[8]   passed      2.74e-01 1.24e-03
## beta[9]   passed     -1.27e-02 8.89e-05
## beta[10]  passed      1.01e-02 9.19e-04
## beta[11]  passed      1.73e-01 3.35e-04
## beta[12]  passed      4.82e-03 1.20e-05
## beta[13]  passed      1.42e-01 9.25e-05
## beta[14]  passed     -4.62e-02 2.38e-04
## beta[15]  passed      3.01e-02 3.95e-05
## beta[16]  passed      2.14e-01 9.97e-04
## beta[17]  passed      1.37e-01 1.00e-03
```

```
## delta      passed      1.18e+04 9.97e+01
```

# Appendix

## Bernoulli/Beta

A natural conjugate prior for the Bernoulli distribution is the Beta distribution.

$$f(x_i|\theta) = \theta^{x_i}(1-\theta)^{1-x_i}$$

$$L(\mathbf{x}|\theta) = \theta^{\sum_{i=1}^{n} x_i}(1-\theta)^{n-\sum_{i=1}^{n} x_i}$$

$$h(\theta) = Beta(a,b)$$

We proceed by calculating the posterior distribution for $\theta$

$$h(\theta|\mathbf{x}) \propto L(\mathbf{x}|\theta)h(\theta) = \theta^{\sum_{i=1}^{n} x_i}(1-\theta)^{n-\sum_{i=1}^{n} x_i}\frac{1}{B(a,b)}\theta^{a-1}(1-\theta)^{b-1}I(0 < \theta < 1)$$

$$\propto Beta(\theta|a + \sum_{i=1}^{n} x_i, b + n - \sum_{i+1}^{n} x_i)$$

## Poisson/Gamma

If our data $X_1, \cdot, X_n$ are iid Poisson($\lambda$) distributed then a gamma($\alpha$, $\beta$) prior on $\lambda$ is a conjugate prior. The Likelyhood function is:

$$L(\lambda|\mathbf{x}) = \prod_{i=1}^{n}\frac{e^{-\lambda}\lambda^{\sum_{i=1}^{n} x_i}}{x_i!} = \frac{e^{-\lambda}\lambda^{\sum_{i=1}^{n} x_i}}{\prod_{i=1}^{n} x_i!}$$

Our gamma prior has the expression:

$$h(\lambda) = \frac{\beta^{\alpha}}{\Gamma(\alpha)}\lambda^{\alpha-1}e^{-\beta\lambda}$$

Using bayes rule we find the following posterior:

$$h(\lambda|\mathbf{x}) \propto h(\lambda)L(\mathbf{x}|\lambda) \propto \lambda^{\sum_{i=1}^{n} x_i+\alpha-1}e^{-(n+\beta)\lambda}$$

$$\propto gamma(\sum_{i=1}^{n} x_i + \alpha, n + \beta)$$

## Poisson/Jeffreys prior

The density distribution for poisson is equal to

$$f(n|\lambda) = e^{-\lambda}\frac{\lambda^n}{n!}$$

The jeffreys prior $h(\theta$ is a non informative prior distrubution for a parameter space and its proportionality is expressed as

$$h(\theta) \propto \sqrt{detI(\theta)}$$

$$I(\theta) = -E[\frac{\partial^2}{\partial\theta^2}logf(x|\theta)] = \frac{1}{\theta}$$

And the following jeffreys prior is thus

$$h(\theta) \propto \theta^{-\frac{1}{2}} I_{\theta > 0}$$

The posterior is calculated as follows

$$h(\theta|x) \propto f(\mathbf{x}|\theta)h(\theta) \propto e^{-n\theta}\theta^{-\frac{1}{2}+\sum_{i=1}^{n} x_i}$$

which is in fact a gamma distribution

$$\theta|x \sim Gamma(\alpha = \frac{1}{2} + \sum_{i=1}^{n} x_i, \beta = n)$$