



Student-ID: 11913169  
Degree Program: Master of Science in Economics (Science Track)  
Examiner: Peter Knaus, PhD  
Submission date: TBA

## Triple-Gamma-Regularization

*A Flexible Non-Convex Regularization Penalty based on the Triple-Gamma-Prior*

by

Lucas Unterweger  GitHub  
(Student-ID: 11913169)

## Abstract

Lorem ipsum...

## Acknowledgements

I like to acknowledge ...

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Theoretical Section</b>	<b>2</b>
2.1	Model Complexity . . . . .	2
2.2	Under- and Overfitting . . . . .	2
2.2.1	Regularization . . . . .	2
2.3	Short Deviation to Bayesian Shrinkage Priors . . . . .	4
<b>3</b>	<b>Literature Review</b>	<b>5</b>
3.1	Duality of the Ridge Regression . . . . .	6
3.2	Triple-Gamma-Prior by Cadonna et al. ( <a href="#">2020</a> ) . . . . .	7
<b>4</b>	<b>Model Setup and Derivation</b>	<b>8</b>
4.1	Varying the Hyperparameters . . . . .	11
4.2	Comparison to already existing Penalty Terms . . . . .	17
4.3	Restricted Strong Convexity (RSC) of the Triple-Gamma-Regularization	17
4.4	Alternative Specification (with +1 to mimic artan) . . . . .	17
4.5	Approaches to Estimation . . . . .	18
<b>5</b>	<b>Simulation Section</b>	<b>18</b>
5.1	Computational Performance . . . . .	18
5.2	Implementation as Python Package . . . . .	18
<b>6</b>	<b>Possible Extensions and Criticism</b>	<b>18</b>
<b>7</b>	<b>Conclusion</b>	<b>18</b>
<b>8</b>	<b>List of Figures</b>	<b>19</b>
<b>9</b>	<b>List of Tables</b>	<b>20</b>
<b>10</b>	<b>References</b>	<b>21</b>
	<b>References</b>	<b>iv</b>

## 1 Introduction

Willam of Ockham, born in Ockham, Surrey, probably lived between 1287 and 1348 and is nowadays recognized as a pre-eminent philosopher of the middle ages. Although his name itself is no common knowledge, a principle carrying his name is: *Ockham's Razor*. Interestingly, the main formulation of the principle (*Entia non sunt multiplicanda praeter necessitatem* [plurality should not be posited without necessity]) can not be traced back to Ockham directly, but variations of it can be found in Ockham's writings. Since then nonetheless, the principle has long been used by statisticians and other researchers as a scientific credo to capture the notion that "the simpler of two explanations is to be preferred" (Lazar, [2010](#)).

## 2 Theoretical Section

### 2.1 Model Complexity

Overfitting is a problem. Especially when dimensionality issues arise. For example, 50 observations and 20 features. Model will perform well/almost perfect on training data, but will do terrible on training data.

### 2.2 Under- and Overfitting

General problem: Fit vs. Complexity. What to choose and where is the optimal balance?

Solutions (find source, seen in LMU course): (1) More data (Not always feasible) (2) Better data (Same issue) (3) Reduce Model Complexity (Regularization) (4) Less aggressive optimization (Maybe short review; Idea: Stop optimization when generalization error degrades)

#### 2.2.1 Regularization

Focusing on option three to battle model complexity, regularization methods have been studied thoroughly since the 1990s. However, the emergence of data science - and especially machine learning - as a standalone field of study has led to a broader meaning of the term *regularization*. This phenomenon has been discussed by Kukačka et al. (2017), where the authors establish a taxonomy to distinguish between multiple different definitions. In the traditional sense, as can be seen in Hastie et al. (2009, pp. 167–170), *regularization* refers to a general class of problems of the form

$$\min_{f \in \mathcal{H}} \left\{ \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f) \right\}$$

where  $L(\cdot)$  refers to a loss function defined as some function of the true values and the predicted values and  $J(f)$  is penalty based on the chosen functional from a space of functions  $\mathcal{H}$ . In the context of penalized linear regression, this is equivalent to finding the set of risk minimizing coefficients  $\hat{\beta}$  from the set of all possible combinations of coefficients  $\beta$ . **(Missing Citation)** Thus, resulting in the general class of regularization problems of the form:

$$\min_{\beta \in \mathcal{B}} \left\{ \sum_{i=1}^N L(y_i, f_{\beta}(x_i)) + \lambda J(\beta) \right\}$$

where  $f_{\beta}(x_i)$  is a linear function of the inputs  $x_i$  parametrized by the coefficients  $\beta$ . Hence, in this setting, regularization deals with penalizing the risk function based on the value of the chosen set of coefficients.

However, this only describes a subset of *regularization* methods as stated by Kukačka et al. (2017). The authors use a more general definition of regularization:

**Definition 1. Regularization** is any supplementary technique that aims at making the model generalize better, i.e. produce better results on the test set.

Building on that, they split up the majority of *regularization* methods into (1) methods applied to the data set like transformations or modifications of the inputs, (2) methods altering the selected model family, (3) methods applied to the error/loss function  $L(y_i, f_\beta(x_i))$ , (4) methods applied to regularization/penalty term as described above and (5) alterations of the optimization procedure itself.

*Vielleicht noch genauer auf die Taxonomy eingehen?*

Unsurprisingly, this thesis is concerned with the fourth group of regularization methods, which add a penalty/regularizer term  $J(\beta)$  into the risk function, but before advancing to literature that deals with this kind of problem, it is necessary to establish a terminology which will be used throughout this thesis. **TODO**

Let  $\mathcal{D}$  be a training data set with  $n \in \mathbb{N}$  observations, where every consists of a target variable  $y_i \in \mathbb{R}$  along with a number of corresponding inputs  $x_i \in \mathbb{R}$ . Given a linear function  $f_\beta(x_i)$  of the inputs parametrized by coefficients  $\beta \in \mathcal{B}$ ,  $L(y_i, f_\beta(x_i))$  is the **Loss** function measuring the discrepancy between the actual target  $y_i$  and the output of the linear function  $f_\beta(x_i)$ . According to Vapnik (1991), the **Empirical Risk Functional** is then

$$R_{emp}(\beta) = \frac{1}{n} \sum_{i=1}^n L(y_i, f_\beta(x_i)).$$

**Regularization** in this thesis' context refers then to adding some penalty function  $J(\beta)$  dependent on the set of parameters  $\beta$ , multiplied by some weighting parameter  $\lambda$ , to the empirical risk functional  $R_{emp}(\beta)$ . Thus,  $R_{reg} = R_{emp}(\beta) + \lambda \cdot J(\beta)$ . This results in the overall optimization problem

$$\begin{aligned} & \arg \min_{\beta \in \mathcal{B}} \{R_{reg}(\beta)\} \\ &= \arg \min_{\beta \in \mathcal{B}} \{R_{emp}(\beta) + \lambda \cdot J(\beta)\} \\ &= \arg \min_{\beta \in \mathcal{B}} \left\{ \sum_{i=1}^n L(y_i, f_\beta(x_i)) + \lambda \cdot J(\beta) \right\} \end{aligned}$$

It is important to note here that as  $J(\beta)$  is only a function of the coefficients  $\beta$  and neither the targets  $y_i$  nor the inputs  $x_i$ . It only affects the generalization error of the model, not the training error given by the empirical risk functional  $R_{emp}(\beta)$ .

*Ad Solution (3) because it is relatively easy to implement: Simple approach = Start with simplest model and iteratively add one feature OR iteratively get rid of feature by feature. (Problem: Very arbitrary and hard to reasonably do) Thus, adjust risk function minimization. Important: Regularization adjusts generalization error, not training error!*

*Short break: Talk about Taxonomy by Kukačka et al 2017 on the differing use of the term "regularization"*

*What is needed:*

$$R_{\text{reg}}(\beta) = R_{\text{emp}}(\beta) + \lambda \cdot J(\beta)$$

*where  $R_{\text{emp}}(\beta) = \sum_i^n L(y_i, f(x_i, \beta))$ . Note: Regularization term does not depend on data, just on parametrization.  $\lambda$  controls the strength of regularization. Thus,  $\lambda = 0$  means simple MSE optimization and  $\lambda \rightarrow \infty$  chooses simplest model. As  $\lambda$  is set manually, this also bears some problems. However, typical solution is cross-validation.*

**Literature** (Kukačka et al., 2017)

## 2.3 Short Deviation to Bayesian Shrinkage Priors

Explain how Bayesians battle model complexity.



### 3 Literature Review

The concept of a penalized regression has been around for quite some time and been studied widely in various fields of scientific research. Arguably, this methodological approach to penalized regression started with the publication of two pieces of literature published by Arthur Hoerl and Robert Kennard in 1970 (A. E. Hoerl & Kennard, 1970a, 1970b). With these two papers the authors introduced the widely known *Ridge Regression*, which has been developed from the previously known concept of Ridge analysis. In its core, the authors were trying to tackle the problem of high variances of the regression coefficients in high-dimensional problem settings. This shrinkage estimator, which uses the squared coefficient as a penalty term, "attempt[s] to shrink the coefficients to reduce these variances, while adding some bias." (R. W. Hoerl, 2020) This closely resembles the previously discussed issue of the *Bias-Variance-Tradeoff*, which has been discussed in the *Under- and Overfitting* chapter in section 2 (Roger W. Hoerl, Arthur Hoerl's son, published a historical overview of the development of the concept of *Ridge Regression* in 2020 (R. W. Hoerl, 2020)). The closed form solution of the Ridge estimator is given by

$$\hat{\beta}_{Ridge} = (X^T X + \lambda I)^{-1} X^T y$$

, which adjusts the OLS estimator by shifting the main diagonal entries of the design matrix by  $\lambda$  ( $\lambda \geq 0$ ). It can be shown that this closed form estimator is equivalent to a Lagrangian problem of the following form (van Wieringen, 2015):

$$\hat{\beta}_{Ridge}(\lambda) = \arg \min_{\beta} \{ \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \}$$

This resembles an example of the above defined regularization framework with squared residual loss and a penalty term of the form  $\|\beta\|_2^2$ , which only depends on the parameter  $\beta$ . In case of  $\lambda$  being equal to zero, this reduces to the *maximum likelihood (ML) estimator*.

The publications of Arthur Hoerl and Robert Kennard have led to further advancements, although it took more than 25 years, in shrinkage estimation or related concepts. One concept which is almost as famous *Ridge Regression* is the *Least Absolute Shrinkage and Selection Operator*, more commonly known as *LASSO*, developed by Tibshirani (1996). He argues that the two at the time most prominent shrinkage methods - Ridge and Subset Selection - both have their drawbacks. Ridge regression on the one hand is an optimization problem which continuously shrinks coefficients towards zero, but doesn't select them in a discrete sense, which makes it hard to interpret these models. Subset Selection on the other hand chooses variables in a discrete sense - a variable either stays within the model or it doesn't - and thus creates easily interpretable models, but "[s]mall changes in the data can result in very different models being selected and this can reduce its prediction accuracy." (Tibshirani, 1996) *LASSO* is trying to combine both methods'

advantages by using  $\|\beta\|_1$  as a penalty term.

*Bit more on LASSO?*

*LASSO* and to some extend *Ridge* can be viewed as a special case of a  $l_p$ -norm regularization with corresponding values for  $p$  ( $p = 1$  for *LASSO* and  $p = 2$  for *Ridge*) (Frank & Friedman, 1993).

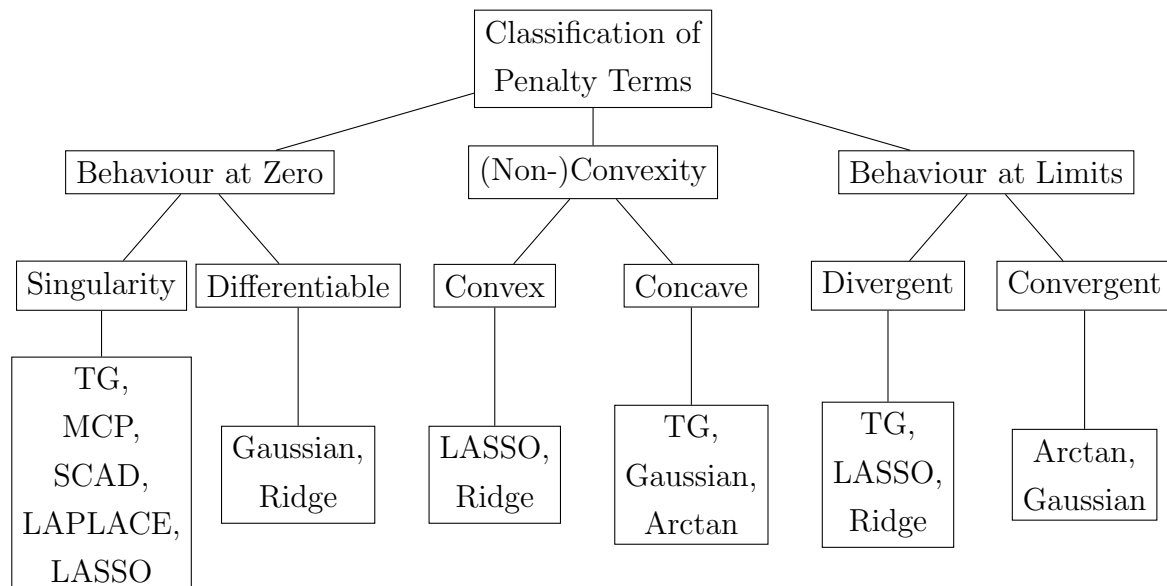
$$\|\beta\|_p = \left( \sum_{i=1}^p |\beta_i|^p \right)^{1/p}$$

An approach combining these two methods is called *Elastic Net* regularization and has been developed by

*Talk about Elastic Net*

*Talk about non-convex penalties*

*Use tree to classify regularization methods.*



Vielleicht als Tabelle mit drei Spalten.

### 3.1 Duality of the Ridge Regression

The main motivation for this thesis stems from the striking duality that exists between Bayesian shrinkage priors and approaches to regularization. Both approaches aim at

tackeling the issue of model complexity and overfitting by making it necessary for the data to be more convincing that the value of an estimate is statistically significant different from zero. Bayesian statistics use specific prior distributions with a usually a lot of mass around zero and heavy tails. Popular examples...

### 3.2 Triple-Gamma-Prior by Cadonna et al. (2020)

A new development in the area of shrinkage priors has been made by Cadonna et al. (2020)

## 4 Model Setup and Derivation

Coming to the theoretical framework of the *triple-gamma-regularization*, let's assume we have a response variable  $y$  and  $p$  predictors along with  $n$  data points. More formally, let  $y = [y_1 \ y_2 \cdots y_n]^T$  and  $x_i = [x_{i1} \ x_{i2} \cdots x_{in}]^T$  with  $\forall i \in \{1, \dots, n\} : y_i, x_i \in \mathbb{R}$ . Here,  $x_i$  is the  $i$ -th predictor, thus resulting in the design matrix  $X = [x_1 \ x_2 \cdots x_p]$ .

Starting from the Bayesian framework, the standard linear regression model is given by

$$y_i = x_i^T \cdot \beta + \varepsilon_i \quad i \in \{1, \dots, n\}$$

with the assumed distribution of  $\varepsilon_i \sim N(0, \sigma^2)$ . Thus it follows that  $y \sim N_n(X\beta, \sigma^2 I)$ . The posterior distribution of the parameter vector  $\beta$ , according to Bayes' Rule, is then proportional to the product of the likelihood of the data and the prior distribution, which can be seen in equation 1.

$$p(\beta|y, X, \sigma^2) \propto \mathcal{L}(y|\beta, \sigma^2, X) \times p(\beta) \quad (1)$$

As stated above, each data point  $y_i$  is assumed to be identically and independently drawn from a normal distribution with mean  $X\beta$  and variance  $\sigma_i^2$ , thus:

$$\begin{aligned} \mathcal{L}(\mathbf{y}|\beta, \sigma^2, X) &= \prod_i^n p(y_i|\beta, \sigma^2, X_i) \\ &= \prod_i^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma}\right)^2\right) \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^\top(\mathbf{y} - \mathbf{X}\beta)\right) \end{aligned} \quad (2)$$

The log of the likelihood function is then given by

$$\begin{aligned} \log \mathcal{L}(\mathbf{y}|\beta, \sigma^2, \mathbf{X}) &= \log\left(\frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^\top(\mathbf{y} - \mathbf{X}\beta)\right)\right) \\ &= \log\left(\frac{1}{(2\pi\sigma^2)^{n/2}}\right) + \log \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^\top(\mathbf{y} - \mathbf{X}\beta)\right) \\ &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^\top(\mathbf{y} - \mathbf{X}\beta) \\ &\propto -\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^\top(\mathbf{y} - \mathbf{X}\beta) = -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \end{aligned}$$

The marginal prior distribution for the parameter vector  $\beta$  stems from the Triple-Gamma-Prior constructed in Cadonna et al. (2020) given in Theorem 1 (a) and is given

by

$$\begin{aligned} p(\sqrt{\beta_j}|\phi^\xi, a^\xi, c^\xi) &= \frac{\Gamma(c^\xi + \frac{1}{2})}{\sqrt{2\pi\phi^\xi} \cdot B(a^\xi, c^\xi)} \cdot U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j}{2\phi^\xi}\right) \\ &\propto U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j}{2\phi^\xi}\right) \end{aligned}$$

Here,  $U(a, b, z)$  refers the confluent hyper-geometric function of the second kind which was introduced by Tricomi (1947). As this prior is specified for the parameter  $\sqrt{\beta_j}$ , we transform the prior by squaring the parameter to gain

$$p(\beta_j|\phi^\xi, a^\xi, c^\xi) \propto U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j^2}{2\phi^\xi}\right)$$

Now, assuming that the parameters are independent a priori, the prior distribution is given by

$$\begin{aligned} p(\beta) &= \prod_j^p p(\beta_j) \\ &= \prod_j^p p(\beta_j|\phi^\xi, a^\xi, c^\xi) \\ &\propto \prod_j^p U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j^2}{2\phi^\xi}\right) \\ &= \prod_j^p \frac{1}{\Gamma(c^\xi + \frac{1}{2})} \int_0^\infty e^{-(\frac{\beta_j^2}{2\phi^\xi})t} t^{c^\xi + \frac{1}{2} - 1} (1+t)^{\frac{3}{2} - a^\xi - c^\xi + \frac{1}{2} - 1} dt \\ &\propto \prod_j^p \int_0^\infty \exp\left(-\frac{\beta_j^2}{2\phi^\xi}t\right) t^{c^\xi - \frac{1}{2}} (1+t)^{1 - a^\xi - c^\xi} dt \end{aligned} \tag{3}$$

Here, in line 1 the assumption of independence between the parameters has been used to describe the distribution of the parameter vector as the product of its individual parameter distributions. In line 2, the marginal prior from Cadonna et al. (2020) has been used as the prior distribution for each individual parameter  $\beta_j$ . In line 3, scaling parameters have been removed by using the proportionality assumption. The last two lines of the derivation insert the integral representation of the confluent hyper-geometric function of the second kind,  $U(a, b, z)$ , which is valid in the case of a positive real part for the first parameter ( $\Re(a) > 0$ ) and again apply proportionality.

Taking the log of the prior distribution and using the properties of the logarithmic function yields the general result

$$\log(p(\beta)) = \log\left(\prod_j^p p(\beta_j|\phi^\xi, a^\xi, c^\xi)\right) = \sum_j^p \log(p(\beta_j|\phi^\xi, a^\xi, c^\xi))$$

A common approach to estimation in regularization settings is the *maximum a posteriori probability (MAP)* estimator (**Missing Citation**), which is defined as

$$\hat{\beta}_{MAP}(x) = \arg \max_{\beta \in \mathbb{R}^p} \{f(x|\beta)g(\beta)\}$$

where  $f(x|\beta)$  describes the the probability density function of a variable  $x$ , which is parametrized by the parameter vector  $\beta$ . The second function  $g(\beta)$  incorporates our prior information about the parameter vector  $\beta$  into the optimization problem.

Returning to our specific problem at hand, the posterior distribution of our parameter vector  $\beta$  can be retrieved by applying Bayes' theorem and the previously gained results in equations 3 and 2. Thus, the posterior distribution of the parameter vector  $\beta$  is proportional to

$$\begin{aligned} p(\beta|y, X, \sigma^2) &\propto p(y|X, \beta, \sigma) \times p(\beta) \\ &\propto \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2}(y-X\beta)^T(y-X\beta)} \times \prod_j^p U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j^2}{2\phi^\xi}\right) \end{aligned} \quad (4)$$

Making use of the monotonicity of the logarithmic function and seeing that it is easier to optimize the log-posterior, Taking the log of the posterior probability distribution, we take the log of result 4.

$$\begin{aligned} \log(\beta|X, y, \sigma^2) &= \log\left(\frac{1}{(2\pi\sigma^2)^{n/2}}\right) - \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \log\left(\prod_j^p U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j^2}{2\phi^\xi}\right)\right) \\ &= \log\left(\frac{1}{(2\pi\sigma^2)^{n/2}}\right) - \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \sum_j^p \log\left(U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j^2}{2\phi^\xi}\right)\right) \\ &\propto -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \sum_j^p \log\left(U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j^2}{2\phi^\xi}\right)\right) \end{aligned} \quad (5)$$

To align with the general specification structure of regularization problems, which can be seen from equation **missEQ**, a parameter  $\lambda$  will be multiplicatively added in front of the penalty term, which makes it possible to adjust the strength of the influence that the penalty has on the chosen parameters. By minimizing the negative log-posterior adjusted with  $\lambda$ , we can retrieve the *maximum a posteriori probability (MAP)* estimator using **Triple-Gamma-Regularization**:

$$\hat{\beta}_{MAP} = \arg \min_{\beta \in \mathbb{R}^p} \left( \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \sum_j^p -\log\left(U\left(c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j^2}{2\phi^\xi}\right)\right) \right) \quad (6)$$

### 4.1 Varying the Hyperparameters

After closer inspection of equation 6, it can easily be seen that this resembles the general penalized regression already seen in **ESL**page398<empty citation> and in section 2 as  $R(\beta) + \lambda \cdot J(\beta)$ . The first term, also called the empirical loss in machine learning literature, is the widely known residual sum of squares:

$$R(\beta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

The second part of the optimization problem can be viewed as a penalty imposed on the total risk based on the size of the estimates:

$$J_{TG}(\beta) = \sum_j^p -\log \left( U \left( c^\xi + \frac{1}{2}, \frac{3}{2} - a^\xi, \frac{\beta_j^2}{2\phi^\xi} \right) \right)$$

In contrast to the *LASSO* penalty, which uses the the absolute value of the coefficient, or the *Ridge* penalty, which uses the square of the estimate, this penalty derived from Cadonna et al. (2020) is based on the log of the confluent hyper-geometric of the second kind. Notably, this penalty term has three additional hyper-parameters:  $c^\xi$ ,  $a^\xi$  and  $\kappa_B$  as  $\phi^\xi = (2c^\xi)/(\kappa_B^2 a^\xi)$ . Here, the restrictions  $a^\xi > 0.5$  and  $0 < c^\xi < \infty$  are necessary to ensure that the penalty for a  $\beta_j$  being equal to zero remains finite and not diverges to negative infinity at zero. This results, which has already been presented and proven as part of Theorem 2 in Cadonna et al. (2020, pp. 5–6), ensures that the negative log of the hypergeometric function remains finite and thus does not produce parameter estimates which are zero for every variables. (**BESSER SCHREIBEN?**)

*Write more to fill this page*

### Variations of the Hyperparameter $a^\xi$

The first hyper-parameter which can be adjusted is  $a^\xi$ . A plot with a set of different values for  $a^\xi$  can be found in figure 1. As already mentioned earlier, the necessary restriction for this hyper-parameter is that it has to be strictly greater than 0.5 to guarantee the finiteness of the penalty. To demonstrate the effects of changes in  $a^\xi$ , the other parameters have been set to  $c^\xi = 0.1$  and  $\kappa_B = 2$ . It can easily be seen from the figure that  $a^\xi$  steers the sharpness of the penalty in small neighbourhoods around  $\beta = 0$ . As  $a^\xi$  increases, the penalty becomes smoother at  $\beta = 0$  with it eventually converging a *Gaussian Penalty* like behaviour. From a modelling perspective, this opens up the possibility of steering the degree of variable selection the penalty performs. Nonetheless, the overall structure of the penalty in the tails does not change systematically apart from a parallel shift, which can be readjusted by specifying a different weighting parameter  $\lambda$  or a different value for  $\kappa_B$  (more on effect of  $\kappa_B$  on the penalty can be found later in this chapter).

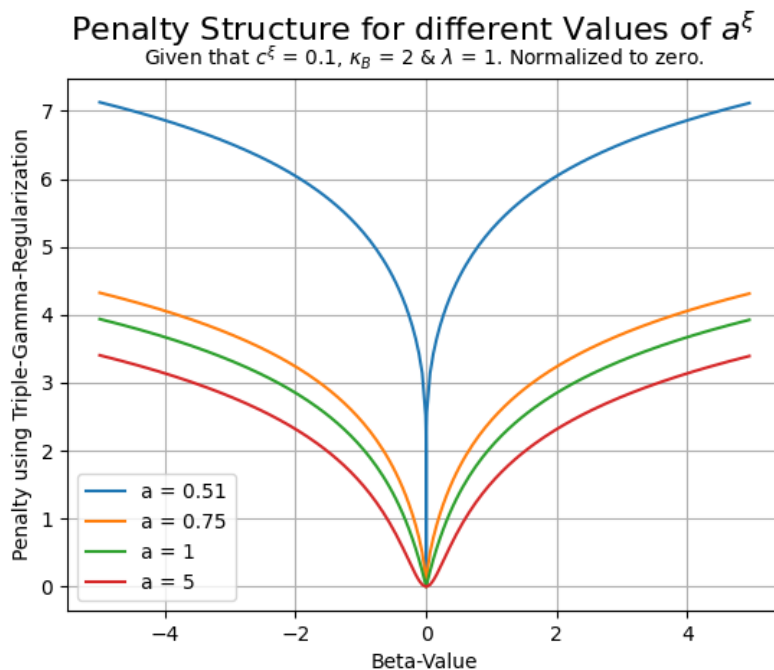


Figure 1: Triple-Gamma-Penalty using different values of  $a^\xi$

Seeing this, it is apparent that a value of  $a^\xi$  close but strictly larger than 0.5 mimics the behaviour of the *Arctan Penalty* by Wang and Zhu (2016) in small neighbourhoods of  $\beta = 0$ . Similar, large positive values for  $a^\xi$  lead to a *Gaussian Penalty* like behaviour in small neighbourhoods of  $\beta = 0$  as recently proposed by John et al. (2022).



### Variations of the Hyperparameter $c^\xi$

In contrast to the hyperparameter  $a^\xi$ , which mainly affects the behaviour at and around  $\beta = 0$ , changes in  $c^\xi$  mainly affect the behaviour in the tails. However, the effect that a change in  $c^\xi$  has on the penalty structure can be split up in two rough subsets of  $(0, \infty)$ . The effect of the first subset of values for  $c^\xi$  which are strictly greater than 0 but less or equal than 0.1 can be found in figure 2 (as already mentioned before, by definition,  $c^\xi$  has to be strictly greater than zero:  $c^\xi > 0$ ). Here, it can be seen that as the values for  $c^\xi$  become smaller, a shifting effect takes place which generally does not influence the overall structure of the penalty, but increases the amount of penalty which is added to the risk function for  $\beta$ -values which are different from zero (In a sense, this has a similar effect to changes in  $\kappa_B$ , which will be explained later). Or, to put it differently, with values of  $c^\xi$  closer to zero, the data has become even more convincing that the value is significantly different from zero.

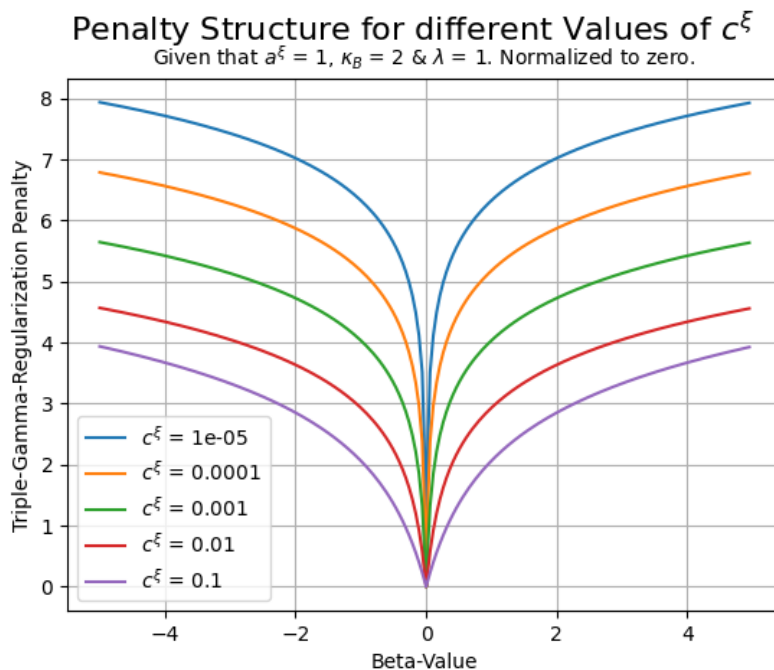


Figure 2: Triple-Gamma-Penalty using different values of  $c^\xi$  with  $0 < c^\xi \leq 0.1$

However, the more interesting effect that a change in  $c^\xi$  has on the penalty structure can be seen for values of  $c^\xi$  that are greater than 0.1. In figure 3, a plot can be found with the Triple-Gamma-Penalty for larger values of  $c^\xi$ . Again, starting from the baseline with  $c^\xi = 0.1$ , higher values for this hyper-parameter mainly change the behaviour of the penalty in the tails. A result that has already been shown by Cadonna et al. (2020) in Table 1, where multiple different hyper-parameter settings are presented, is that with an

increasing value for  $c^\xi$  and with  $a^\xi = 1$  as well as  $\kappa_B = ??$ , the Triple-Gamma-Prior converges to a *LASSO* like shrinkage behaviour. A property which carries over to the proposed regularization setting when using the proposed hyper-parameter values, thus showing that the Triple-Gamma-Penalty can be used both as a non-convex penalty as well as a *LASSO* penalty, creating increased flexibility in modelling approaches.

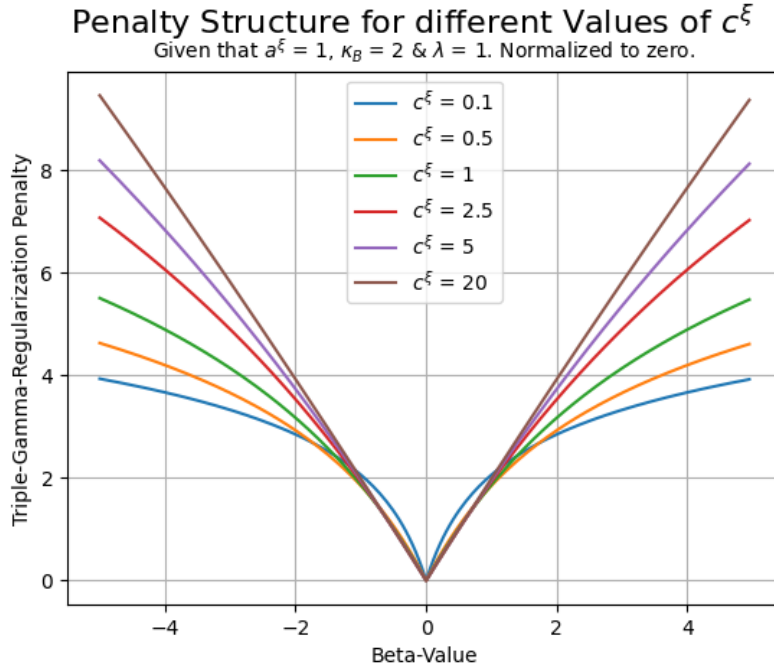
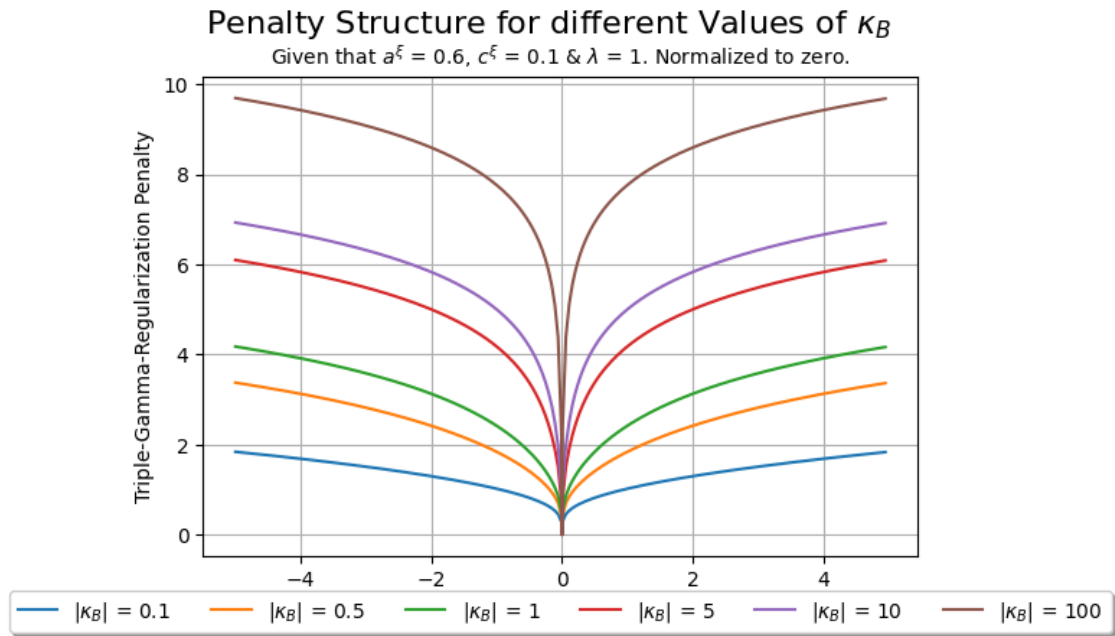


Figure 3: Triple-Gamma-Penalty using different values of  $c^\xi$  with  $c^\xi \geq 0.1$

#### *Variations of the Hyperparameter $\kappa_B$*

The third and final hyper-parameter  $\kappa_B$  enters the Triple-Gamma-Penalty  $J_{TG}(\beta)$  as part of  $\phi^\xi = \frac{2c^\xi}{\kappa_B^2 a^\xi}$  as can be seen from equation 6. Notably,  $\kappa_B$  is squared and thus only the absolute value of  $\kappa_B$ ,  $|\kappa_B|$ , influences the structure of the penalty. The overall third function value is defined as  $\frac{\beta_j^2}{2\phi^\xi}$  and by plugging in  $\phi^\xi$  we get  $\frac{\beta_j^2 \kappa_B^2 a^\xi}{4c^\xi}$ , it can be seen that a value of  $\kappa_B = 0$  leads to the entire parameter value being zero for all values of  $\beta_j$ . Hence, a change in  $\beta_j$  won't influence the penalty and furthermore won't have an influence on the overall risk minimization problem, the result being that the optimal set of parameters will only depend on the chosen loss function.

For all values of  $\kappa_B \neq 0$ , the value of the hyper-parameter will influence the penalty structure. A plot with several different values for the absolute value of  $\kappa_B$  can be found in figure 4.

Figure 4: Triple-Gamma-Penalty using different values of  $\kappa_B$

Variable	Change		Mathematical Properties	
	Positive Change	Negative Change	Defined Range	Misc.
$a^\xi$	Shifting towards <i>Gaussian</i> -like behaviour at $\beta = 0$	Shifting towards singularity at $\beta = 0$ ; Higher immediate penalty for coefficients $\beta \neq 0$	$(0.5, \infty)$	Expl.
$c^\xi$	Generally, convergence towards convexity and, given certain settings for $a^\xi$ and $\kappa_B$ , LASSO. Higher values increase the additional penalty for higher absolute values of $\beta$ .	For values smaller than 0.1, similar effect to increase in $a^\xi$	$(0, \infty)$	Expl.
$\kappa_B$	15883	5.2e-8	$(-\infty, \infty) \setminus \{0\}$	Expl.

Table 1: Summary of the effects of changes in the hyperparameters  $a^\xi$ ,  $c^\xi$  and  $\kappa_B$  on the penalty structure

## 4.2 Comparison to already existing Penalty Terms

As already mentioned in section 3, several other penalty terms have already been widely studied in the literature. Convex penalties like *Ridge* (A. E. Hoerl & Kennard, 1970b) or non-convex penalties like the *Ar(c)tan* (Wang & Zhu, 2016) and Gaussian (John et al., 2022) have managed to establish itself as prominent approaches to regularization. It is now certainly of interest to see how the *Triple-Gamma* penalty compares to the established methods. Seeing that, due to its flexibility, there is not *one Triple-Gamma* penalty, three distinct hyper-parameter setting have been chosen to represent the proposed penalty term.

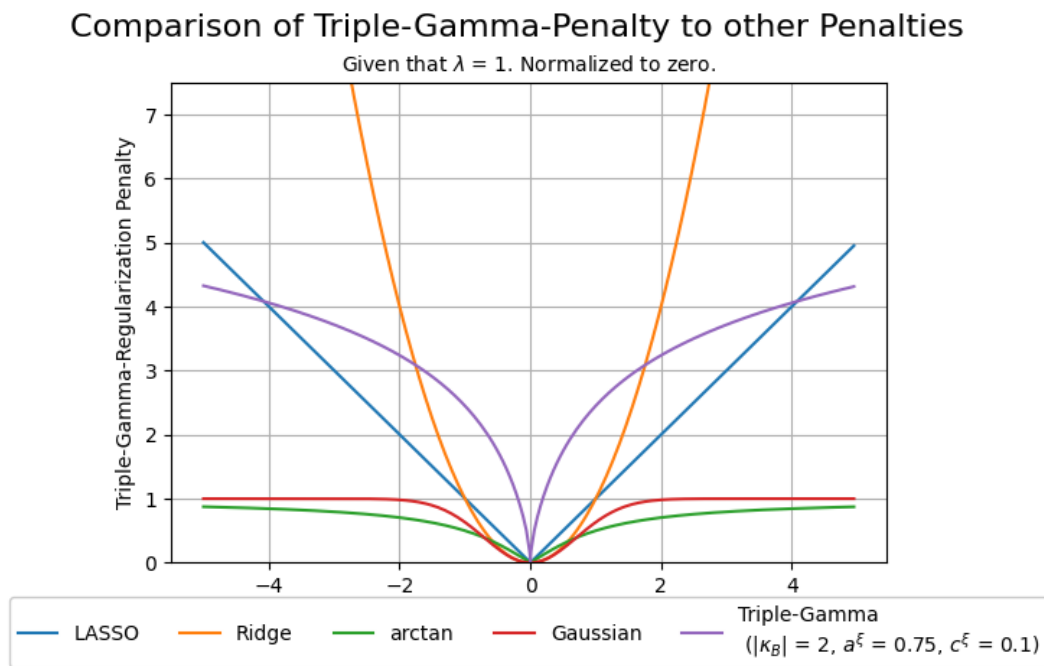


Figure 5: Basic Representation of the Triple-Gamma-Penalty using Hyperparameters  $c^\xi = 0.1, \kappa_B = 2, a^\xi = 0.75$  compared to *LASSO* and *Ridge* Penalties

## 4.3 Restricted Strong Convexity (RSC) of the Triple-Gamma-Regularization

## 4.4 Alternative Specification (with +1 to mimic artan)

As previously mentioned, the *Arctan* penalty developed by Wang and Zhu (2016) has a similar structure as the *Triple-Gamma-Penalty* in small neighbourhoods near  $\beta = 0$  when using specific hyper-parameters (see section 4.2). The distinctive structural difference between these two penalties emerge when looking at the limits when  $\lim_{\beta \rightarrow +/\infty} J_{TG}(\beta)$

## 4.5 Approaches to Estimation

## 5 Simulation Section

Use the earlier derivation, code up the functions and simulate data to check behaviour for different datasets. Compare to base OLS, Ridge and Lasso Regression?

### 5.1 Computational Performance

Talking about Gradient Descent Methods in more depth. Why Gradient Clipping. Maybe more modern estimation techniques?

See how computation times change when increasing the size of the data set or when the number of parameter changes. Use Stochastic Gradient Descent (SGD)

### 5.2 Implementation as Python Package

How to use it. Explanation of Functions. Input - Output Tables

## 6 Possible Extensions and Criticism

## 7 Conclusion

## 8 List of Figures

### List of Figures

1	Triple-Gamma-Penalty using different values of $a^\xi$ . . . . .	12
2	Triple-Gamma-Penalty using different values of $c^\xi$ with $0 < c^\xi \leq 0.1$ . . .	13
3	Triple-Gamma-Penalty using different values of $c^\xi$ with $c^\xi \geq 0.1$ . . . . .	14
4	Triple-Gamma-Penalty using different values of $\kappa_B$ . . . . .	15
5	Basic Representation of the Triple-Gamma-Penalty using Hyperparameters $c^\xi = 0.1, \kappa_B = 2, a^\xi = 0.75$ compared to <i>LASSO</i> and <i>Ridge</i> Penalties	17

**9 List of Tables****List of Tables**

1	Summary of the effects of changes in the hyperparameters $a^\xi$ , $c^\xi$ and $\kappa_B$ on the penalty structure . . . . .	16
---	---	----



**10 References**

## References

- Cadonna, A., Frühwirth-Schnatter, S., & Knaus, P. (2020). Triple the gamma—a unifying shrinkage prior for variance and variable selection in sparse state space and tvp models. *Econometrics*, 8(2), 20.
- Frank, L. E., & Friedman, J. H. (1993). A statistical view of some chemometrics regression tools. *Technometrics*, 35(2), 109–135.
- Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (Vol. 2). Springer.
- Hoerl, A. E., & Kennard, R. W. (1970a). Ridge regression: Applications to nonorthogonal problems. *Technometrics*, 12(1), 69–82.
- Hoerl, A. E., & Kennard, R. W. (1970b). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1), 55–67.
- Hoerl, R. W. (2020). Ridge regression: A historical context. *Technometrics*, 62(4), 420–425.
- John, M., Vettam, S., & Wu, Y. (2022). A novel nonconvex, smooth-at-origin penalty for statistical learning. *arXiv preprint arXiv:2204.03123*.
- Kukačka, J., Golkov, V., & Cremers, D. (2017). Regularization for deep learning: A taxonomy. *arXiv preprint arXiv:1710.10686*.
- Lazar, N. (2010). Ockham’s razor. *Wiley Interdisciplinary Reviews: Computational Statistics*, 2(2), 243–246.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1), 267–288.
- Tricomi, F. (1947). Sulle funzioni ipergeometriche confluenti [Paper for hypergeometric function of second kind.]. *Annali di Matematica Pura ed Applicata*, 26(1), 141–175. <https://doi.org/10.1007/BF02415375>
- van Wieringen, W. N. (2015). Lecture notes on ridge regression. *arXiv preprint arXiv:1509.09169*.
- Vapnik, V. (1991). Principles of risk minimization for learning theory [Definition of Empirical Risk Minimization]. *Advances in neural information processing systems*, 4.
- Wang, Y., & Zhu, L. (2016). Variable selection and parameter estimation with the atan regularization method. *Journal of Probability and Statistics*, 2016.