

# A robust hybrid of lasso and ridge regression

AI2101: Convex Optimisation Project, Group 35

Ananthoju Pranav Sai  
R.No: AI20BTECH11004

Mulugu Vishwanath Sharma  
R.No: MA20BTECH11010

Murarisetty Adhvik  
R.No: AI20BTECH11015

Yashas Tadikamalla  
R.No: AI20BTECH11027

Nelakuditi Rahul Naga  
R.No: AI20BTECH11029

**Abstract**—Lasso and Ridge regression are the regularized versions of least squares regression using L1 and L2 penalties respectively, on the coefficient vector. To make these regressions more robust we may replace least squares with Huber's criterion which is a hybrid of squared error (for relatively small errors) and absolute error (for relatively large ones). A reversed version of Huber's criterion (Berhu criterion) can be used as a hybrid penalty function. Relatively small coefficients contribute their L1 norm to this penalty while larger ones cause it to grow quadratically. This hybrid sets some coefficients to 0 (as lasso does) while shrinking the larger coefficients the way ridge regression does. Both the Huber and reversed Huber penalty functions employ a scale parameter. We provide an objective function that is jointly convex in the regression coefficient vector and these two scale parameters.

**Index Terms**—lasso, ridge, regression, robust, huber, reversed huber, berhu

## I. INTRODUCTION

Regularization is often necessary in the context of machine learning, especially, when we are looking at regression. In this project, we focus on the regularized least squares regression problem. As the name suggests, regularized least squares (RLS) refers to the family methods which solve for least squares problem using regularization, to further constrain the solution. Regularized version of least squares can be beneficial in the following two situations: When the number of variables in the linear system exceeds the number of observations or when the model suffers from lack of generalization ability. Lasso and ridge regression are two famous versions of the RLS problem. However, they do have some disadvantages. In this project, we try to formulate a robust hybrid of lasso and ridge regression. Section II briefly discusses about lasso and ridge regression. Section III recalls the Huber function and Section IV proposes the Reversed Huber function, which will be used as loss and penalty functions respectively. Section V explains the proposed robust hybrid formulation for the regression problem. Section VI describes a technique to formulate a function that is jointly convex in both original and scale parameters. Section VII backs the theory behind it. Section VIII shows how to optimise the convex penalized regression criterion using CVXPY. Finally, Section IX demonstrates the application of the proposed formulation to diabetes problem.

## II. LASSO AND RIDGE REGRESSION

Consider the regression problem of predicting  $y \in \mathbb{R}$  based on  $z \in \mathbb{R}^d$ . Hence the training data in this case is given by the data points  $(z_i, y_i)$  for  $i = 1, 2, \dots, n$ . Suppose that each predictor vector  $z_i$  is converted into a **feature** vector  $x_i \in \mathbb{R}^p$  using a fixed function  $\phi$  as follows :

$$x_i = \phi(z_i) \quad (1)$$

Then the predictor for  $y$ , which is denoted by  $\hat{y}$  is given by :

$$\hat{y} = \mu + x^T \beta \quad (2)$$

where  $\beta \in \mathbb{R}^p$ .

In Ridge regression, we regularize the objective by adding a L2 penalty term. Thus, the loss function in Ridge regression is given by :

$$L(\theta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad (3)$$

$$L(\theta) = \sum_{i=1}^n (y_i - \mu - x_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad (4)$$

where  $\theta$  denotes all the parameters involved in the above formulation. The ridge parameter  $\lambda \in [0, \infty]$ . Defining  $\epsilon_i = y_i - \mu - x_i^T \beta$  the loss function becomes :

$$L(\theta) = \|\epsilon\|_2^2 + \lambda \|\beta\|_2^2 \quad (5)$$

We can solve the Ridge regression problem by **minimizing** the above loss function over  $\theta$ .

In Lasso regression, we replace the L2 penalty term with L1 penalty term. The loss function in Lasso regression is given by :

$$L(\theta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\beta_j| \quad (6)$$

$$L(\theta) = \sum_{i=1}^n (y_i - \mu - x_i^T \beta)^2 + \lambda \sum_{j=1}^p |\beta_j| \quad (7)$$

Defining  $\epsilon_i = y_i - \mu - x_i^T \beta$  the loss function becomes :

$$L(\theta) = \|\epsilon\|_2^2 + \lambda \|\beta\|_1 \quad (8)$$

We can solve the Lasso regression problem by **minimizing** the above loss function over  $\theta$ .

The advantage with Lasso regression is that, it can produce sparse solutions, i.e, the solutions with some or most of the  $\beta_j$ 's equal to zero. Sparsity is desirable for interpretation. Ridge regression cannot produce such sparsity, but can shrink the weights in the solution. However, such sparsity being forced into solutions from  $L1$  has some disadvantages. In case of several correlated features with large effect on the output, lasso regression tries to zero out some, perhaps all but one of them. Ridge regression, on the other hand, shares the coefficient value among the correlated features, without any selection. A belief exists that the solutions from  $L1$  penalty are less accurate than from  $L2$  penalty.

### III. HUBER FUNCTION

The least squares criterion is well suited to  $y_i$  with Gaussian distribution, but can give poor performance when  $y_i$  has a heavier tailed distribution, or what is almost the same, when there are outliers. To overcome this Huber introduced a robust estimator with a loss function, which is unaffected by very large residuals.

**Huber function** is defined as:

$$\mathcal{H}(z) = \begin{cases} z^2 & |z| \leq 1 \\ 2|z| - 1 & |z| \geq 1 \end{cases} \quad (9)$$

The function has quadratic nature for smaller values of  $z$  and is linear for larger values of  $z$ . Using this Huber function, the robust estimator can be defined as:

$$\sum_{i=1}^n \mathcal{H}_M \left( \frac{y_i - \mu - x_i^T \beta}{\sigma} \right) \quad (10)$$

Where,  $\mathcal{H}_M(z)$  is:

$$\mathcal{H}_M(z) = M^2 \mathcal{H}(z/M) = \begin{cases} z^2 & |z| \leq M \\ 2M|z| - M^2 & |z| \geq M \end{cases} \quad (11)$$

Here,  $M$  is a shape parameter and  $\sigma$  is a scale parameter.  $M$  describes the point of transition from quadratic to linear and controls the amount of robustness. Therefore, if the error is smaller than  $M\sigma$  it gets squared or else it comes under linear regime.

The Huber criteria resembles least squares as  $M$  increases, making  $\beta$  more efficient for normally distributed data but less robust. The criterion is more like  $L1$  regression for small values of  $M$ , making it more robust against outliers but less efficient for normally distributed data. Huber proposes that taking  $M = 1.35$  gives more robustness while retaining 95% statistical efficiency for normally distributed data.

For fixed values of  $M$  and  $\sigma$  and a convex penalty  $P(\cdot)$ , the following will be convex in  $(\mu, \beta) \in \mathcal{R}^{p+1}$ .

$$\sum_{i=1}^n \mathcal{H}_M \left( \frac{y_i - \mu - x_i^T \beta}{\sigma} \right) + \lambda \sum_{j=1}^p P(\beta_j) \quad (12)$$

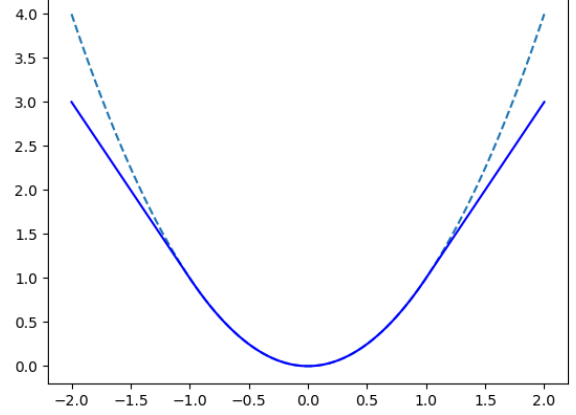


Figure 1. Huber function  $\mathcal{H}(z)$

### IV. REVERSED HUBER FUNCTION

We know that Huber function is quadratic near zero and linear for larger values. It's also a good fit for approximately normal distributed errors with heavier tails. But it is not well suited as a regularization term on regression coefficients.

It is generally preferred to use  $L_1$  regularization as it sets some  $\beta_j$  to zero making  $\beta$  sparse. This sparsity leads to savings in computation and storage. But there are some disadvantages associated with it. Hence, to overcome them, we propose a new penalty function Berhu, which uses  $L_1$  penalty for small values and  $L_2$  penalty for larger values. **Berhu function** is defined as:

$$\mathcal{B}(z) = \begin{cases} |z| & |z| \leq 1 \\ \frac{z^2+1}{2} & |z| \geq 1 \end{cases} \quad (13)$$

Then the penalty term using Berhu will be:

$$\mathcal{B}_M(z) = M \mathcal{B}_M \left( \frac{z}{M} \right) = \begin{cases} |z| & |z| \leq M \\ \frac{z^2+M^2}{2M} & |z| \geq M \end{cases} \quad (14)$$

Here  $M$  describes where the transition from linear to quadratic should occur. Note that the function  $\mathcal{B}_M(z)$  is convex in  $z$ .

### V. PROPOSED HYBRID OF LASSO AND RIDGE REGRESSION

Using Huber function for loss and Berhu function for penalty, we propose a robust hybrid of Lasso and Ridge regression. The robustness of the function is decided by value of  $M$ . The final hybrid function is:

$$\sum_{i=1}^n \mathcal{H}_M \left( \frac{y_i - \mu - x_i^T \beta}{\sigma} \right) + \lambda \sum_{j=1}^p \mathcal{B}_M \left( \frac{\beta_j}{\tau} \right) \quad (15)$$

equation (15) is jointly convex in  $\mu$  and  $\beta$  given  $M, \tau$ , and  $\sigma$ .

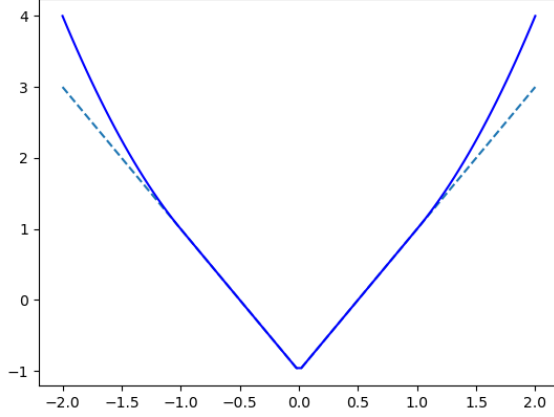


Figure 2. Berhu function  $B(z)$

## VI. CONCOMITANT SCALE ESTIMATION

The parameter  $M$  can be fixed to some value like 1.35. But there remain 3 other tuning parameters to consider:  $\lambda, \sigma$ , and  $\tau$ . Instead of doing a 3 dimensional parameter search, we try to frame a objective which is jointly convex in  $(\mu, \beta, \sigma, \tau)$ , leaving only one dimensional search over  $\lambda$ . In practice it is necessary to estimate  $\sigma$  from the data, simultaneously with  $\beta$ . Huber proposed several ways to jointly estimate  $\sigma$  and  $\beta$ . One of his ideas for robust regression is to minimize

$$n\sigma + \sum_{i=1}^n \mathcal{H}_M\left(\frac{y_i - \mu - x_i' \beta}{\sigma}\right) \sigma \quad (16)$$

over  $\beta$  and  $\sigma$ . For any fixed value of  $\sigma \in (0, \infty)$ , the minimizer  $\beta$  of (16) is the same as that of (10). The criterion (16) is however jointly convex as a function of  $(\beta, \sigma) \in \mathbb{R}^p \times (0, \infty)$  as shown in VII. Therefore convex optimization can be applied to estimate  $\beta$  and  $\sigma$  together. This removes the need for ad hoc algorithms to estimate  $\beta$  for fixed  $\sigma$  and  $\sigma$  for fixed  $\beta$ . The same idea can be used for  $\tau$ . We can reframe the penalty term as  $\tau + \mathcal{B}_M(\beta/\tau)\tau$  which is jointly convex in both  $\beta$  and  $\tau$ . Using Huber's penalty function on the regression errors and the Berhu function on the coefficients leads to a criterion of the form

$$n\sigma + \sum_{i=1}^n \mathcal{H}_M\left(\frac{y_i - \mu - x_i' \beta}{\sigma}\right) \sigma + \lambda \left[ \rho\tau + \left\{ \sum_{j=1}^p \mathcal{B}_M\left(\frac{\beta_j}{\tau}\right) \right\} \right] \quad (17)$$

where  $\lambda \in [0, \infty)$  governs the amount of regularization applied. The expression in (17) is jointly convex in  $(\mu, \beta, \sigma, \tau) \in \mathbb{R}^{p+1} \times (0, \infty)^2$ . So we can apply convex optimization to estimate  $\mu, \beta, \sigma$  and  $\tau$  together.

## VII. THEORY OF CONCOMITANT SCALE ESTIMATION

**Lemma 7.1:** let  $\rho$  be a convex and twice differential function on an interval  $\mathcal{I} \subseteq \mathbb{R}$ . Then  $\rho\left(\frac{\eta}{\sigma}\right)\sigma$  is a convex function of  $(\eta, \sigma) \in \mathcal{I} \times (0, \infty)$ .

*Proof 7.1:* Let  $\eta_0 \in \mathcal{I}$  and  $\sigma_0 \in (0, \infty)$  and parameterise  $\eta$  and  $\sigma$  linearly as  $\eta = \eta_0 + \cos(\theta) \times t$  and  $\sigma = \sigma_0 + \sin(\theta) \times t$  over  $t$  in an open interval with 0 and  $\theta \in [0, 2\pi)$ . Where  $\theta$  denotes the direction.

Now let us try to calculate the double derivative of  $\rho\left(\frac{\eta}{\sigma}\right)\sigma$  w.r.t  $t$ .

First derivative:

$$\frac{d}{dt} \rho\left(\frac{\eta}{\sigma}\right) \sigma = \rho'\left(\frac{\eta}{\sigma}\right) \frac{\cos(\theta)\sigma - \sin(\theta)\eta}{\sigma} + \rho\left(\frac{\eta}{\sigma}\right) \sin(\theta) \quad (18)$$

Second derivative:

$$\frac{d^2}{dt^2} \rho\left(\frac{\eta}{\sigma}\right) \sigma = \rho''\left(\frac{\eta}{\sigma}\right) \frac{(\cos(\theta)\sigma - \sin(\theta)\eta)^2}{\sigma^3}. \quad (19)$$

which is always  $\geq 0$ . i.e. the curvature of  $\rho\left(\frac{\eta}{\sigma}\right)\sigma$  is always non negative.

Hence proved.

**Lemma 7.2:** Let  $\rho$  be a convex function on an interval  $\mathcal{I} \subseteq \mathbb{R}$ . Then  $\rho\left(\frac{\eta}{\sigma}\right)\sigma$  is a convex function of  $(\eta, \sigma) \in \mathcal{I} \times (0, \infty)$ .

*Proof 7.2:* Take two points  $(\eta_0, \sigma_0)$  and  $(\eta_1, \sigma_1)$  both in  $\mathcal{I} \times (0, \infty)$ . Now take a point  $(\eta, \sigma)$  cutting the line segment joining the two points in  $\frac{\lambda}{1-\lambda}$  ratio internally. Where  $0 < \lambda < 1$ . Then,

$$\begin{aligned} \eta &= \lambda\eta_1 + (1-\lambda)\eta_0 \\ \sigma &= \lambda\sigma_1 + (1-\lambda)\sigma_0 \end{aligned}$$

For  $\epsilon > 0$ . Let  $\rho_\epsilon$  be a convex and twice differentiable function on  $\mathcal{I}$  that is everywhere within  $\epsilon$  of  $\rho$ . Then,

$$\begin{aligned} \rho\left(\frac{\eta}{\sigma}\right) \sigma &\geq \rho_\epsilon\left(\frac{\eta}{\sigma}\right) \sigma - \epsilon\sigma \\ &\geq \lambda\rho_\epsilon\left(\frac{\eta_1}{\sigma_1}\right) \sigma_1 + (1-\lambda)\rho_\epsilon\left(\frac{\eta_0}{\sigma_0}\right) \sigma_0 - \epsilon\sigma \\ &\geq \lambda\rho\left(\frac{\eta_1}{\sigma_1}\right) \sigma_1 + (1-\lambda)\rho\left(\frac{\eta_0}{\sigma_0}\right) \sigma_0 - \epsilon(\lambda\sigma_1 + (1-\lambda)\sigma_0 - \sigma) \end{aligned} \quad (20)$$

By taking  $\epsilon$  arbitrarily small we can say that

$$\rho\left(\frac{\eta}{\sigma}\right) \sigma \geq \lambda\rho\left(\frac{\eta_1}{\sigma_1}\right) \sigma_1 + (1-\lambda)\rho\left(\frac{\eta_0}{\sigma_0}\right) \sigma_0 \quad (21)$$

Hence proved that  $\rho\left(\frac{\eta}{\sigma}\right)\sigma$  is convex.

**Lemma 7.3:** Let  $y_i \in \mathbb{R}$  and  $x_i \in \mathbb{R}^p$  for  $i=1, 2, \dots, n$ . Let  $\rho$  be a convex function on  $\mathbb{R}$ . Then

$$n\sigma + \sum_{i=1}^n \rho\left(\frac{y_i - \mu - x_i' \beta}{\sigma}\right) \sigma \quad (22)$$

is convex in  $(\mu, \beta, \sigma) \in \mathbb{R}^{p+1} \times (0, \infty)$ .

*Proof 7.3:* The first term i.e  $n\sigma$  is linear so it is convex. Now we need to prove that the second term  $\sum_{i=1}^n \rho\left(\frac{y_i - \mu - x_i' \beta}{\sigma}\right) \sigma$  is convex.

From 7.2  $\rho\left(\frac{\eta}{\sigma}\right)\sigma$  is convex on  $(\eta, \sigma) \in \mathcal{I} \times (0, \infty)$ . But the mapping of  $\eta \rightarrow (y_i - \mu - x'_i\beta)$  is affine so  $\rho\left(\frac{y_i - \mu - x'_i\beta}{\sigma}\right)\sigma$  is convex over  $(\beta, \sigma) \in \mathbb{R}^p \times (0, \infty)$ . And sum over  $i$  preserves the convexity.

Now let us take  $\rho(z) = z^2$ .

From 7.1 we can tell that  $\rho(z) = z^2$  is convex. so by applying it to the above 7.3 we can say that

$$n\sigma + \sum_{i=1}^n \rho\left(\frac{y_i - \mu - x'_i\beta}{\sigma}\right)\sigma \quad (23)$$

is convex in  $(\mu, \beta, \sigma) \in \mathbb{R}^p \times (0, \infty)$ .

Now the equation 23 is minimized by taking  $\mu$  and  $\beta$  as their least square estimates and  $\sigma = \hat{\sigma}$ . Where,

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left(y_i - \hat{\mu} - x'_i\hat{\beta}\right)^2 \quad (24)$$

This gives rise to the usual normal distribution maximum likelihood estimates. It turns out that the equation (23) is not a simple monotone transformation of the negative log likelihood.

$$\frac{n}{2} \log(2\pi) + n \log(\sigma) + \frac{1}{2\sigma^2} \sum_{i=1}^n \left(y_i - \mu - x'_i\beta\right)^2 \quad (25)$$

The negative log-likelihood of the above equation is not convex in  $\sigma$  (for any  $\mu, \beta$ ). Huber's technique has convexified the Gaussian Log Likelihood 25 into equation (23).

Turning to the least squares case, if we substitute equation (24) in the criterion we get  $2 \left( \sum_{i=1}^n (y_i - \mu - x'_i\beta)^2 \right)^{1/2}$ . A ridge regression for both residuals and coefficients then minimizes

$$\mathcal{R}(\mu, \beta, \tau, \sigma; \lambda) = n\sigma + \sum_{i=1}^n \left( \frac{(y_i - \mu - x'_i\beta)^2}{\sigma} \right) + \lambda \left( p\tau + \sum_{j=1}^p \frac{\beta_j^2}{\tau} \right) \quad (26)$$

over  $(\mu, \beta, \sigma, \tau) \in \mathbb{R}^{p+1} \times [0, \infty)^2$  for fixed  $\lambda \in [0, \infty)$ . Minimizing over  $\sigma$  and  $\tau$  in closed form leaving

$$\min_{\sigma, \tau} \mathcal{R}(\beta, \tau, \sigma; \lambda) = 2 \left( \sum_{i=1}^n (y_i - \mu - x'_i\beta)^2 \right)^{1/2} + 2\lambda \left( \sum_{j=1}^p \beta_j^2 \right)^{1/2} \quad (27)$$

to be minimized on  $\beta$  given  $\lambda$ . In other words, after putting Huber's concomitant scale estimators into ridge regression we recover ridge regression again. The criterion is  $\|y - \mu - x'\beta\|_2 + \lambda \|\beta\|_2^2$  which gives the same trace as  $\|y - \mu - x'\beta\|_2^2 + \lambda \|\beta\|_2^2$ .

Now let us take  $\rho(z) = |z|$  we obtain a degenerate result:  $\sigma + \rho\left(\frac{\beta}{\sigma}\right)\sigma = \sigma + |\beta|$ . Even though it is convex it is minimized as  $\sigma \rightarrow 0$  without any regards to  $\beta$ . Hence Huber's device

does not yield a usable concomitant scale estimate for  $L_1$  regression.

The degeneracy for L1 loss propagates to the Huber loss HM when  $M \leq 1$ . We may write

$$\sigma + \mathcal{H}_M(z/\sigma) = \begin{cases} \sigma + z^2/\sigma & \sigma \geq |z|/M \\ \sigma + 2M|z| - M^2\sigma & \sigma \leq |z|/M \end{cases} \quad (28)$$

The minimum of 28 is attained at  $\sigma = 0$  regardless of  $z$ . If one should ever want a concomitant scale estimate when  $M \leq 1$ , then a simple fix is to use  $(1 + M^2)\sigma + \mathcal{H}_M(z/\sigma)\sigma$ .

## VIII. CVXPY IMPLEMENTATION

The Huber function  $\mathcal{H}_M(x)$  is equivalent to the quadratic program

$$\text{minimize} \quad u^2 + 2Mv \quad (29)$$

$$\text{subject to} \quad |x| \leq u + v \quad (30)$$

$$u \leq M \quad (31)$$

$$v \geq 0, \quad (32)$$

CVXPY has a built-in Huber function, but as we are using concomitant scale estimation, we represent the function  $\sigma + \mathcal{H}_M(z/\sigma)\sigma$  as quadratic program

$$\text{minimize} \quad \sigma + u^2/\sigma + 2Mv \quad (33)$$

$$\text{subject to} \quad |z| \leq u + v \quad (34)$$

$$u \leq M\sigma \quad (35)$$

$$v \geq 0, \quad (36)$$

after substituting and simplifying.

Similarly, the Berhu function  $\mathcal{B}_M(x)$  may be represented by the quadratic program

$$\text{minimize} \quad v + u^2/(2M) + u \quad (37)$$

$$\text{subject to} \quad |x| \leq v + u \quad (38)$$

$$v \leq M \quad (39)$$

$$u \geq 0. \quad (40)$$

As we are using concomitant scale estimation, we represent the function  $\tau + \mathcal{B}_M(z/\tau)\tau$  as quadratic program

$$\text{minimize} \quad \tau + v + u^2/(2M\tau) + u \quad (41)$$

$$\text{subject to} \quad |z| \leq v + u \quad (42)$$

$$v \leq M\tau \quad (43)$$

$$u \geq 0, \quad (44)$$

The code for CVXPY implementation can be found [here](#)

## IX. PRACTICAL APPLICATION - DIABETES EXAMPLE

We show and compare the results for penalized regression on the diabetes test data. Figure 3 visualises the coefficient vectors for Lasso, ridge and robust hybrid regressions, respectively. In each graph, the coefficient vector starts at  $\beta(\infty) = (0, 0, \dots, 0)$  and grows as one moves from left to right.

We note the following observations:

- For lasso regression, we can observe more sparsity.
- For ridge regression, the coefficient vectors are mostly non-sparse
- With the robust hybrid, we get a hybrid behaviour. We get some sparsity, and some diverse values for  $\beta$ 's

#### REFERENCES

- [1] Art B. Owen, "A robust hybrid of lasso and ridge regression", 2007.
- [2] Boyd, S. and Vandenberghe, L., "Convex optimization", Cambridge University Press, Cambridge, 2004.

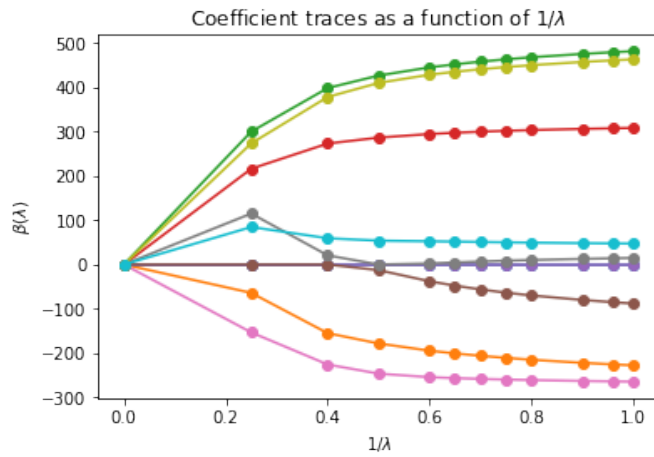
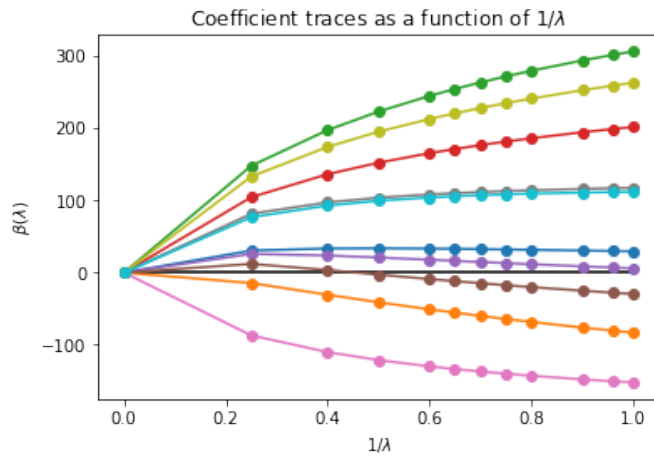
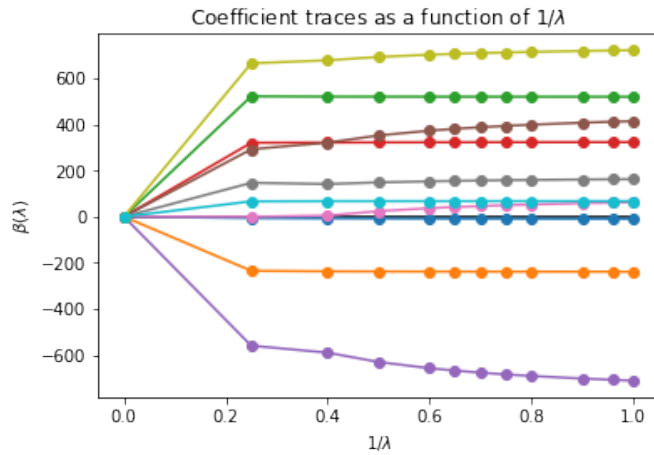


Figure 3. Plots for coefficient traces  $\beta(\lambda)$  for lasso, ridge and robust hybrid regressions, respectively