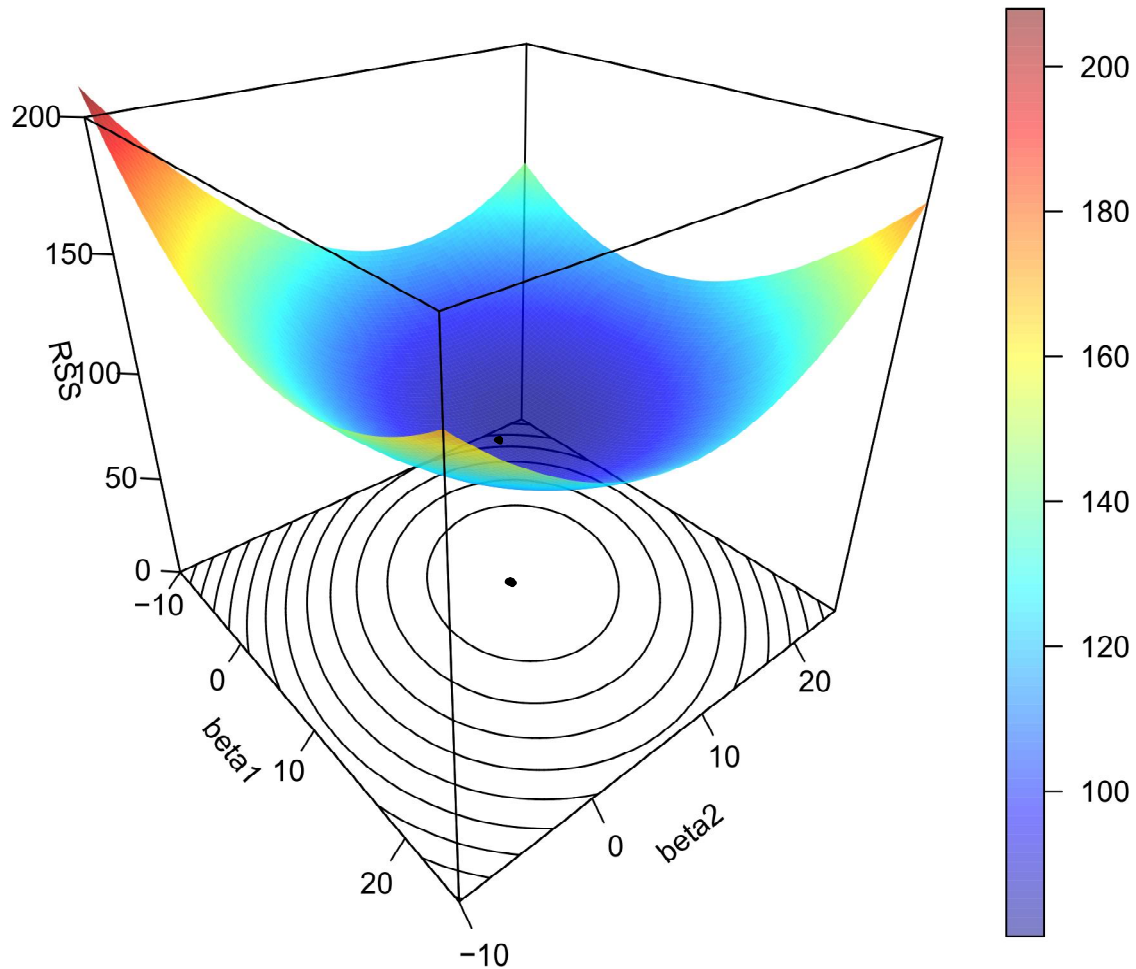


LINEAR METHODS FOR REGRESSION: REGULARIZATION

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

LEAST SQUARES



$$\hat{\beta}_{LS} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|Y - X\beta\|_2^2$$

REGULARIZATION

Another way to control bias and variance is through **regularization** or **shrinkage**.

The idea is to make your estimates of β 'smaller', rather than set them to zero

(which is what all subsets does)

One way to do this is called **ridge regression**¹:

$$\hat{\beta}_{\text{ridge},t} = \underset{\|\beta\|_2^2 \leq t}{\operatorname{argmin}} \|Y - \mathbb{X}\beta\|_2^2$$

for any $t \geq 0$.

Compare this to **least squares**

$$\hat{\beta}_{LS} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|Y - \mathbb{X}\beta\|_2^2$$

$$\text{min } \|Y - \mathbb{X}\beta\|_2^2$$

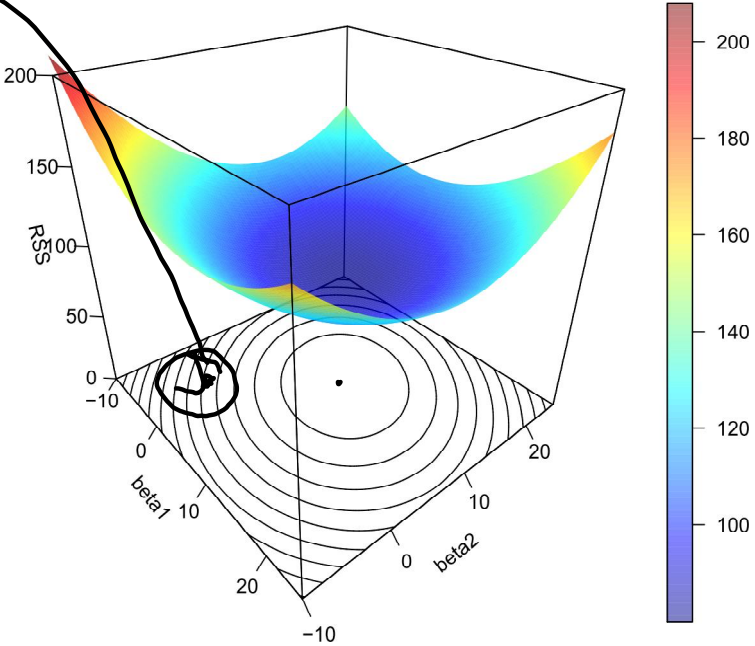
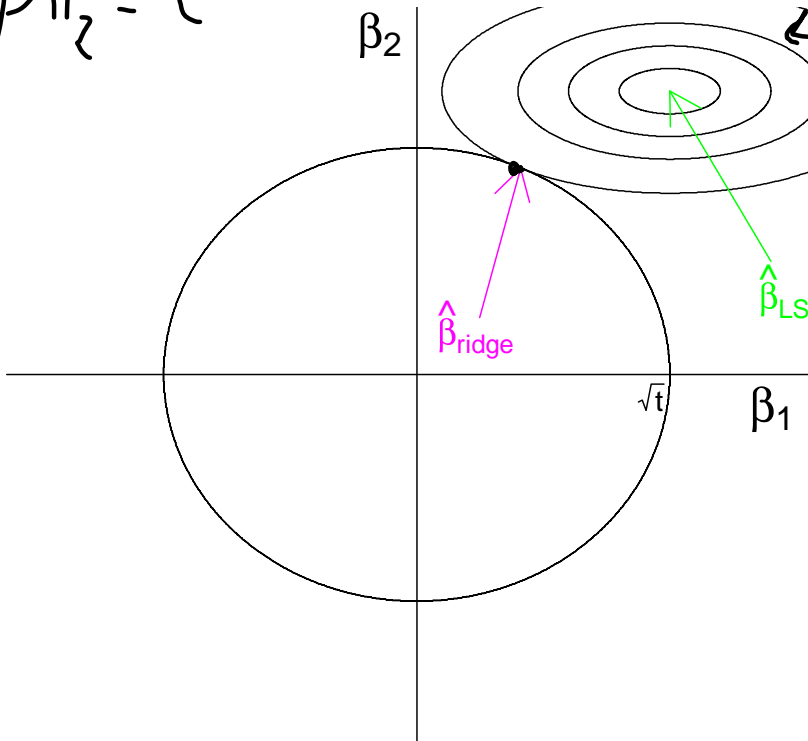
SUBJECT TO

$$\|\beta\|_2^2 \leq t$$

¹Hoerl, Kennard (1970)

GEOMETRY OF RIDGE REGRESSION IN \mathbb{R}^2

$$\|\beta\|_2^2 \leq t$$



RIDGE REGRESSION

$$t \text{ SMALL} \Rightarrow \frac{1}{t} \text{ LARGE} = \lambda$$

$$\Rightarrow \lambda \|\beta\|_2^2 \gg \|Y - X\beta\|_2^2 \Rightarrow \|\beta\|_2^2 \text{ SMALL}$$

An equivalent way to write

$$t \text{ LARGE} \Rightarrow \frac{1}{t} \text{ SMALL} = \lambda$$

$$\hat{\beta}_{\text{ridge},t} = \underset{\|\beta\|_2^2 \leq t}{\operatorname{argmin}} \|Y - X\beta\|_2^2 \quad \lambda \|\beta\|_2^2 \ll \|Y - X\beta\|_2^2 \quad (1)$$

$$\Rightarrow \|\beta\|_2^2 \text{ LARGE}(\lambda)$$

is in the **Lagrangian** form

$$\hat{\beta}_{\text{ridge},\lambda} = \underset{\beta}{\operatorname{argmin}} \left(\|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right) \quad (2)$$

For every λ' there is a unique t' (and vice versa) that makes

$$t = \|\hat{\beta}_{\lambda}\|_2^2$$

$$\hat{\beta}_{\text{ridge},\lambda'} = \hat{\beta}_{\text{ridge},t'}$$

$$\frac{1}{t} \sim \lambda$$

REGULARIZATION AND STANDARDIZATION

The coefficient vector isn't invariant to rescaling.

If an intercept is included, do not penalize it:

$$\min_{\beta_0, \beta} \sum_{i=1}^n (Y_i - \beta_0 + \beta^\top X_i)^2 + \lambda \|\beta\|_2^2$$

The usual way of addressing this in regression is:

- Standardize all covariates for which scale is meaningful:

$$x_j \leftarrow \frac{(x_j - \text{mean}(x_j))}{\text{sd}(x_j)}$$

$$\|\beta\|_2^2 \leq t$$

(So, don't standardize indicator variables, for instance)

- Standardize the response $Y \leftarrow Y - \text{mean}(Y)$
- Don't include an intercept

(It would have been equal to $\text{mean}(Y)$)

RIDGE REGRESSION

Observe:

- $\lambda = 0$ (or $t = \infty$) makes $\hat{\beta}_{\text{ridge},\lambda=0} = \hat{\beta}_{LS}$
- Any $\lambda > 0$ (or $t < \infty$) penalizes larger values of β , effectively shrinking them.

Note: λ and t are known as **tuning parameters**

(Alternatively, hyper-parameters)

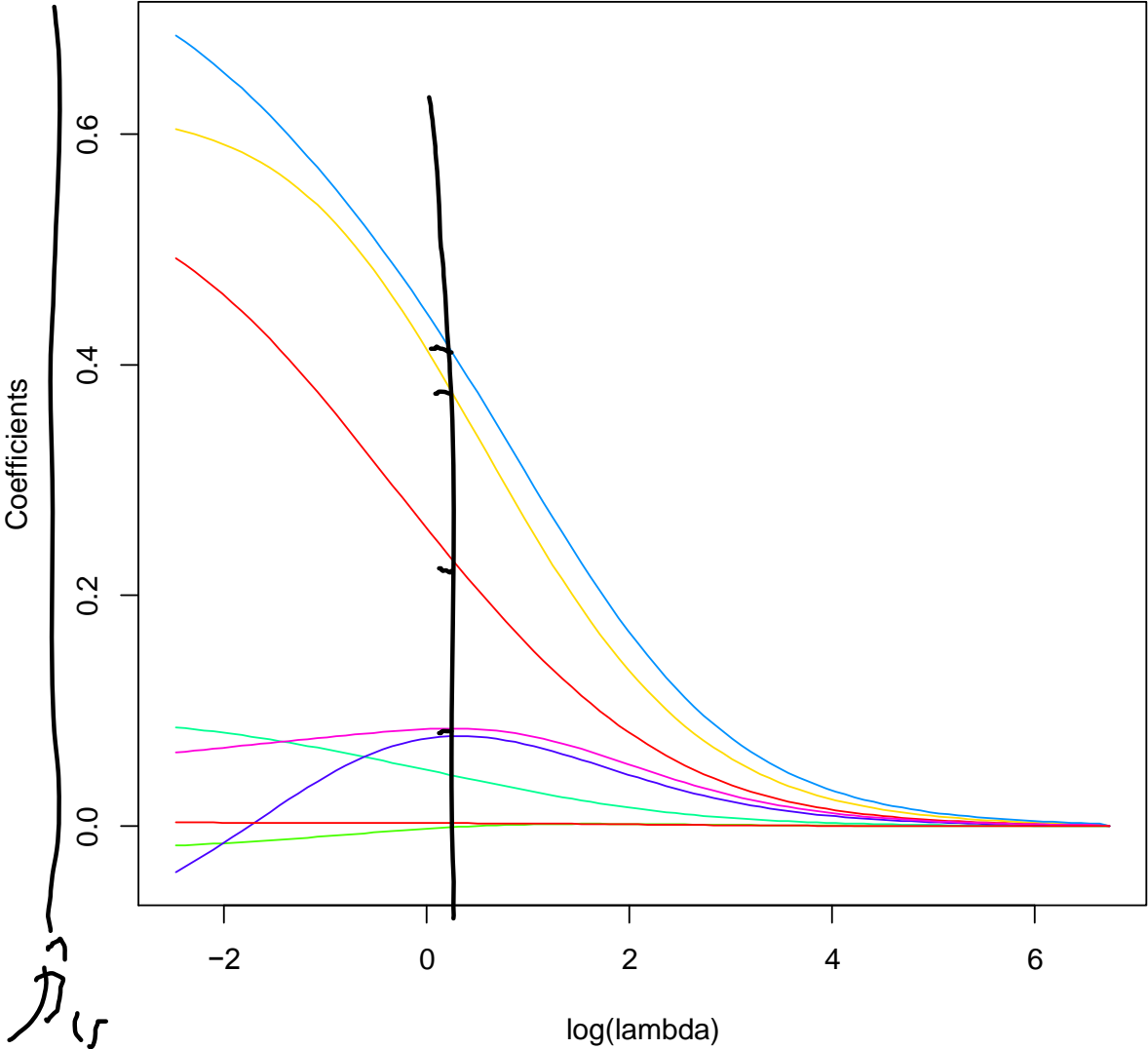
However we think about it, we have produced a **suite** of solutions

$$\{\hat{\beta}_{\text{ridge},\lambda} : \lambda \in [0, \infty)\}$$

What do these solutions look like?

RIDGE REGRESSION PATH

"COEFFICIENT PATH"



RIDGE REGRESSION

REMINDER: The least squares solution can be written:

$$\hat{\beta}_{\text{LS}} = (\mathbb{X}^T \mathbb{X})^\dagger \mathbb{X}^T Y. \quad \|Y - \mathbb{X}\beta\|_2^2$$

However, if $\text{rank}(\mathbb{X}) < p$, then $\hat{\beta}_{\text{LS}}$ is not unique. In fact,

$$\forall b \in \{b : \mathbb{X}b = 0\}$$

$\hat{\beta}_{\text{LS}} + b$ is a valid least squares solution.

It turns out through differential calculus, we can write out the ridge regression solution as well:

$$\hat{\beta}_{\text{ridge}, \lambda} = (\mathbb{X}^T \mathbb{X} + \underbrace{(\lambda I)}_{\lambda > 0})^{-1} \mathbb{X}^T Y$$

Quite similar. However, the λ can make all the difference..

REGULARIZATION - RIDGE REGRESSION

Using the **SVD** ($\mathbb{X} = UDV^T$), we can look even deeper.

LEFT SVECS \uparrow RIGHT SVECS

$$\hat{\beta}_{LS} = VD^{-1}U^T Y$$

$$= \sum_{j=1}^p \mathbf{v}_j \left(\frac{1}{d_j} \right) \boxed{\mathbf{u}_j^T Y}$$

\swarrow j^{th} column of V

$$\hat{\beta}_{\text{ridge}, \lambda} = V(D^2 + \lambda I)^{-1} D U^T Y$$

$$= \sum_{j=1}^p \mathbf{v}_j \left(\frac{d_j}{d_j^2 + \lambda} \right) \mathbf{u}_j^T Y.$$

\hookrightarrow KEEPS AWAY FROM 0

Similarly $= U \underbrace{D V^T V D^{-1}} U^T Y$

$$\mathbb{X} \hat{\beta}_{LS} = \underbrace{U U^T}_{\text{PROJECTION}} Y$$

$$= \sum_{j=1}^p \mathbf{u}_j \left(\frac{1}{d_j} \right) \mathbf{u}_j^T Y$$

$$\mathbb{X} \hat{\beta}_{\text{ridge}, \lambda} = U D (D^2 + \lambda I)^{-1} D U^T Y$$

$$= \sum_{j=1}^p \mathbf{u}_j \left(\frac{d_j^2}{d_j^2 + \lambda} \right) \mathbf{u}_j^T Y.$$

\Rightarrow Ridge shrinks the data by an additional factor of λ .

RIDGE REGRESSION: A BAYESIAN APPROACH

Suppose we specify the likelihood as

$$Y_i \sim N(X_i^\top \beta, \sigma^2)$$

and put a prior distribution of $\beta \sim N(0, \tau^2 I)$.

Then we have the following posterior (making some conditional independence assumptions)

$$p(\beta | Y, X, \sigma^2, \tau^2) \propto p(Y | X, \beta, \sigma^2) p(\beta | \tau^2).$$

After kernel matching, we find that the posterior mode/mean is

$$\hat{\beta}_{\text{ridge}, \lambda = \sigma^2 / \tau^2}$$

RIDGE REGRESSION IN A NEW SPACE

WOODBURY IDENTITY

Note the matrix identity

$$(A - BC^{-1}E)^{-1}BC^{-1} = A^{-1}B(C - EA^{-1}B)^{-1}$$

(Henderson, Searle (1980), equation (13))

Then,

$$\hat{\beta}_{\text{ridge},\lambda} = \overbrace{(\mathbb{X}^T \mathbb{X} + \lambda I)^{-1} \mathbb{X}^T}^H Y = \overbrace{\mathbb{X}^T (\mathbb{X} \mathbb{X}^T + \lambda I)^{-1}}^{H'} Y$$

$p \times p$ $n \times n$

Now, the inversion is in n -space instead of p , which could be a substantial savings

This seemingly innocuous change can have large implications

RIDGE IN A NEW SPACE: COMPUTATIONS

The ridge solution solves either the **normal equations**

$$(\mathbb{X}^\top \mathbb{X} + \lambda I) \hat{\beta} = \mathbb{X}^\top Y$$

or the **adjoint problem**

$$\mathbb{X}^\top (\mathbb{X} \mathbb{X}^\top + \lambda I)^{-1} Y$$

The 'heavy lifting' in each case is done with the inversion

- $\mathbb{X}^\top \mathbb{X} \in \mathbb{R}^{p \times p} \implies$ takes p^3 computations, p^2 space
- $\mathbb{X} \mathbb{X}^\top \in \mathbb{R}^{n \times n} \implies$ takes n^3 computations, n^2 space

CONCLUSION: Depending on the relative size of n and p , this could be substantial savings

However, a much deeper realization is possible..

(KERNEL) RIDGE REGRESSION

Suppose we want to predict at X , then

$$\hat{f}(X) = X^\top \hat{\beta}_{\text{ridge}, \lambda} = X^\top \mathbb{X}^\top (\mathbb{X}\mathbb{X}^\top + \lambda I)^{-1} Y$$

Also,

$$\mathbb{X}\mathbb{X}^\top = \begin{bmatrix} \langle X_1, X_1 \rangle & \langle X_1, X_2 \rangle & \cdots & \langle X_1, X_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle X_n, X_1 \rangle & \langle X_n, X_2 \rangle & \cdots & \langle X_n, X_n \rangle \end{bmatrix}$$

$X_1^\top X_2 = \langle X_1, X_2 \rangle \Rightarrow X_1^\top X_1 = \langle X_1, X_1 \rangle = \|X_1\|_2^2$

and

$$X^\top \mathbb{X}^\top = [\langle X, X_1 \rangle, \langle X, X_2 \rangle, \dots, \langle X, X_n \rangle]$$

where $\langle X, X' \rangle = X^\top X'$ is the Euclidean inner product.

If we transform $X_i \mapsto \Phi(X_i)$, and the range of Φ is equipped with an inner product, we can use $\langle \Phi(X_i), \Phi(X_{i'}) \rangle$

Inserting Φ is known as **kernelization** or a **kernel trick**

kernelization def.
 $\langle X, X' \rangle = \langle \Phi(X), \Phi(X') \rangle$

(KERNEL) RIDGE REGRESSION

EXAMPLE: Suppose $X = (\text{income}, \text{height})^\top$

Then we could specify the map

$$\Phi(X)^\top = (\text{income}, \text{height}, \text{income} * \text{height}, \text{income}^2, \text{height}^2)$$

The induced **feature matrix** is then

$$X = \begin{bmatrix} \Phi(X_1) \\ \vdots \\ \Phi(X_n) \end{bmatrix} \in \mathbb{R}^{n \times 5}$$

(KERNEL) RIDGE REGRESSION

Ordinarily, this would mean we need to solve the **normal equation** inversion for $p = 5$

- $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{p \times p} \implies$ takes p^3 computations, p^2 space

However, using the **kernel trick** we can solve instead

- $\mathbf{X}\mathbf{X}^T \in \mathbb{R}^{n \times n} \implies$ takes n^3 computations, n^2 space

which is fixed in \mathfrak{p}

IMPLICATION: We can add essentially arbitrary nonlinearity without paying higher computational, storage cost!

We will return to this again with **support vector machines (SVM)**

Ridge in practice

RIDGE REGRESSION: THE TUNING PARAMETER

We can use a degrees of freedom based risk estimator to choose λ

The degrees of freedom of $\hat{\beta}_{\text{ridge},\lambda}$ can be seen to be

$$\text{df} = \text{trace} \left[\mathbb{X}(\mathbb{X}^\top \mathbb{X} + \lambda I)^{-1} \mathbb{X}^\top \right] = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

(As $\lambda \rightarrow 0$, we get the number of parameters)

A common, classic choice is **generalized cross-validation** (GCV), which has the form:

$$\text{GCV}(\hat{\beta}) = \frac{\hat{\mathbb{P}}\ell_{\hat{\beta}}}{(1 - \text{df}(\hat{\beta})/n)^2}$$

(Golub, Heath, Wahba (1979))

Note that this looks a lot like AIC with unknown variance, but with $\log(1 - \text{df}/n)$ as penalty

RIDGE REGRESSION: THE TUNING PARAMETER

To see this last claim, observe

$$\log \left(\text{GCV}(\hat{\beta}) \right) \propto \log(\hat{R}_{\text{train}}) - 2 \log(1 - \text{df}(\hat{\beta})/n)$$

VERSUS

$$\text{AIC}(\hat{\beta}) \propto \log(\hat{R}_{\text{train}}) + 2n^{-1} \text{df}(\hat{\beta})$$

RIDGE REGRESSION: THE TUNING PARAMETER

Nowadays, using K -fold cross-validation is common

Think of CV_K as a function of λ , and pick its **minimum**:

$$\hat{\lambda} = \underset{\lambda \geq 0}{\operatorname{argmin}} CV_K(\lambda)$$

Now, we report $\hat{\beta}_{\text{ridge}, \hat{\lambda}}$ as our estimator

RIDGE REGRESSION: COMPUTATION

There are several ways to compute ridge regression

We can follow any conventional least squares solving technique (i.e.: QR factorization, Cholesky Decomposition, SVD,...):

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\beta = \mathbf{X}^T \mathbf{Y}$$

Alternatively, we can actually solve it using **lm** in **R** if we make the following augmentation

$$\tilde{\mathbf{Y}} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{n+p} \text{ and } \tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I} \end{bmatrix}$$

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I} \end{bmatrix}^T \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I} \end{bmatrix} = \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$$

RIDGE REGRESSION IN R

We will concentrate on a slightly more complicated way, as it will make things easier later.

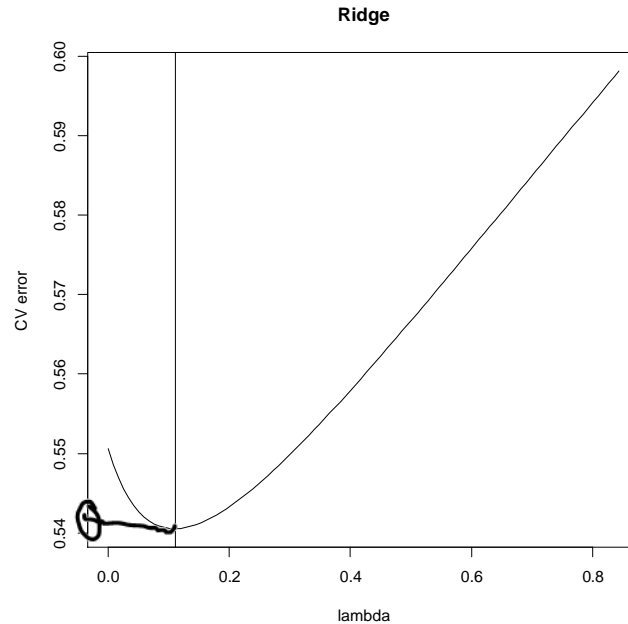
```
install.packages('glmnet')  
library(glmnet)  
ridge.out = cv.glmnet(x=X,y=Y,alpha=0)
```

RIDGE
REGRESSION

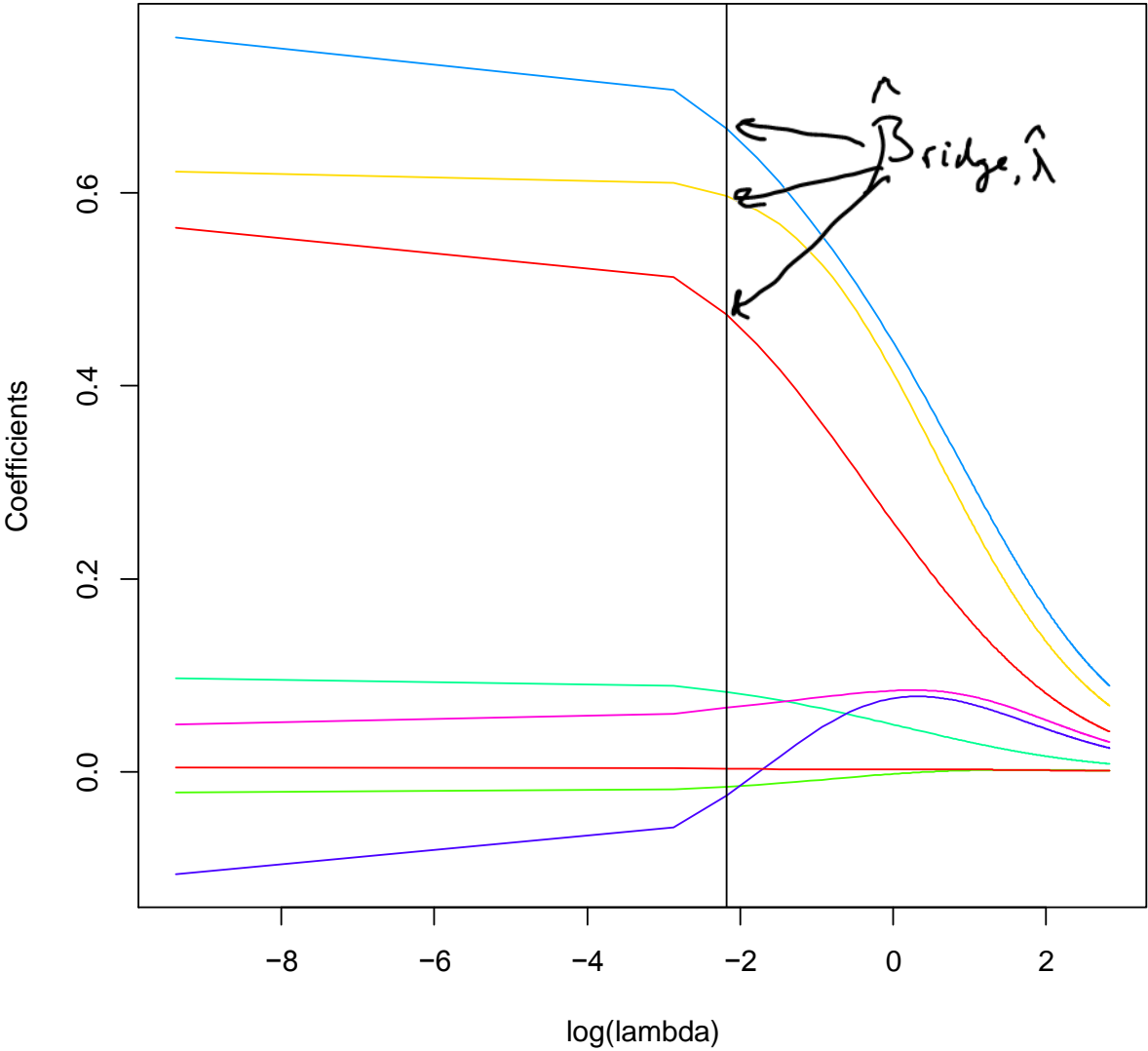
RIDGE REGRESSION: CV PLOT

```
X = as.matrix(X)
ridge.out = cv.glmnet(x=X,y=Y,alpha=0)
plot(ridge.out$lambda,ridge.out$cvm,
      xlab='lambda',ylab='CV error',main='Ridge',type='l')
abline(v=ridge.out$lambda[which.min(ridge.out$cvm)])
```

$\Lambda = [\lambda_1, \dots, \lambda_p]$ $CV_K(\Lambda)$



RIDGE REGRESSION PATH



CAN WE GET THE BEST OF BOTH WORLDS?

To recap:

- Forward, backward, and all subsets regression offer good tools for model selection.

(but the optimization problem is nonconvex)

- Ridge regression provides regularization, which trades off bias and variance and also stabilizes multicollinearity.

(problem is convex, but doesn't do model selection)

RIDGE REGRESSION $\min ||\mathbb{Y} - \mathbb{X}\beta||_2^2$ subject to $||\beta||_2 \leq t$

BEST LINEAR $\min ||\mathbb{Y} - \mathbb{X}\beta||_2^2$ subject to $||\beta||_0 \leq t$

REGRESSION MODEL

($||\beta||_0 =$ the number of nonzero elements in β)

AN INTUITIVE IDEA

RIDGE REGRESSION $\min ||\mathbb{Y} - \mathbb{X}\beta||_2^2$ subject to $||\beta||_2^2 \leq t$

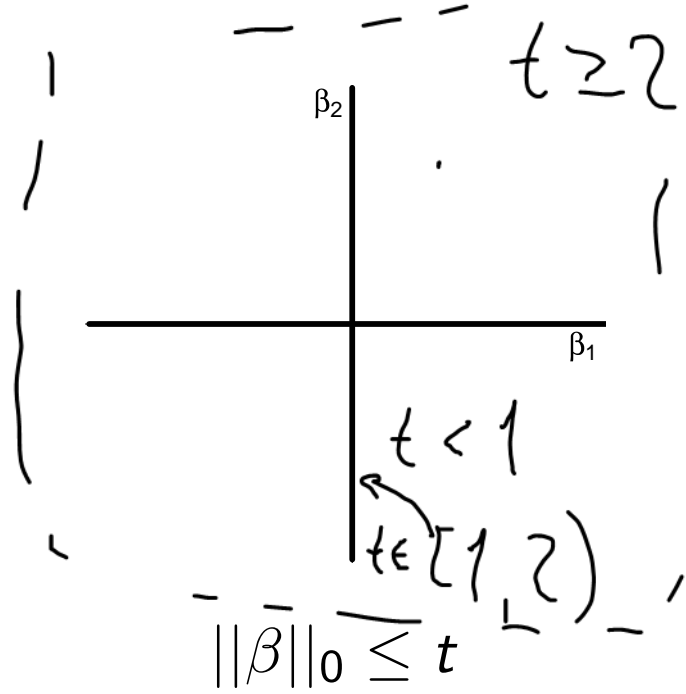
BEST LINEAR REGRESSION MODEL $\min ||\mathbb{Y} - \mathbb{X}\beta||_2^2$ subject to $||\beta||_0 \leq t$

($||\beta||_0 =$ the number of nonzero elements in β)

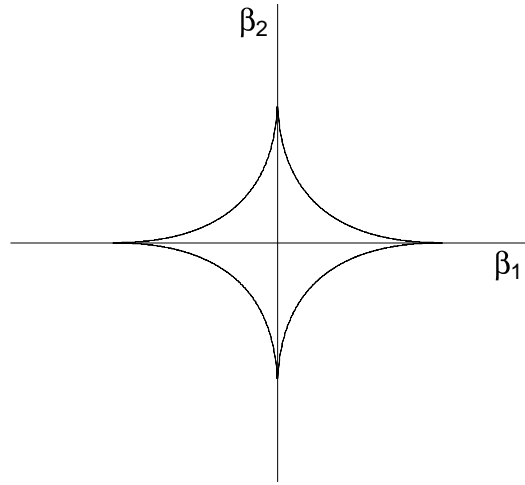
| | BEST LINEAR REGRESSION MODEL | RIDGE REGRESSION |
|----------------------------------|---|-----------------------------|
| Computationally Feasible? | No | Yes |
| Does Model Selection? | Yes | No |

Can we ‘interpolate’ $||\beta||_2$ and $||\beta||_0$ to find a method that does both?

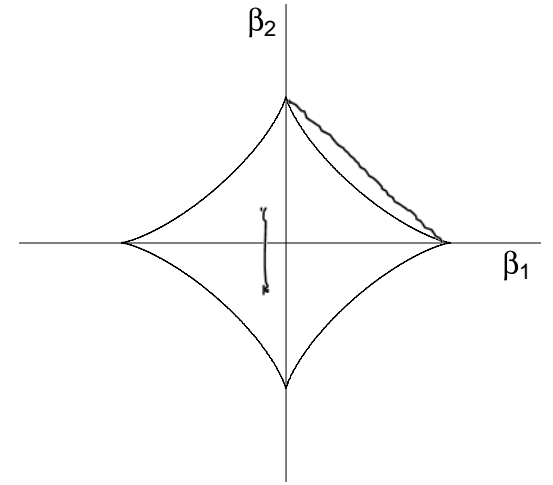
GEOMETRY OF REGULARIZATION IN \mathbb{R}^2 : CONVEXITY



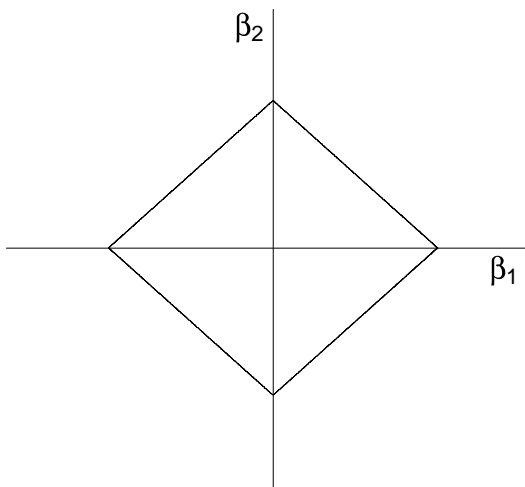
$$\|\beta\|_0 \leq t$$



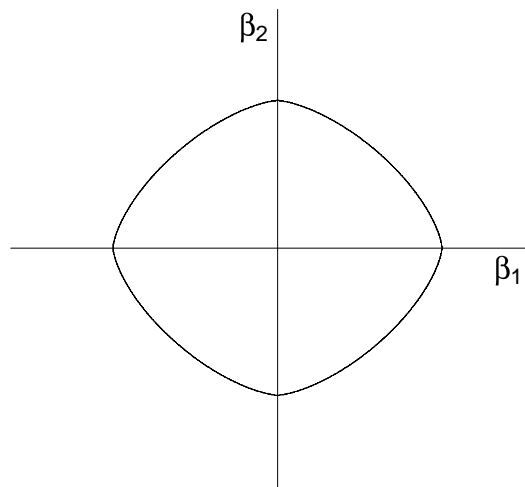
$$\|\beta\|_{\frac{1}{2}} \leq t$$



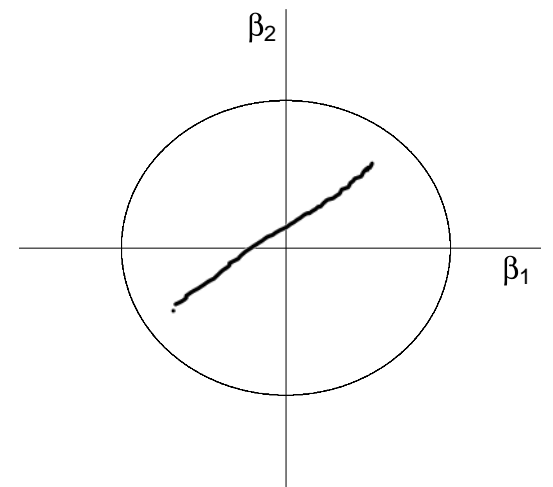
$$\|\beta\|_{\frac{3}{4}} \leq t$$



$$\|\beta\|_1 \leq t$$

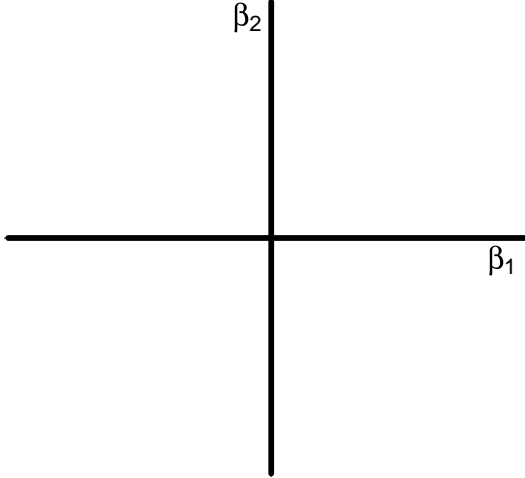


$$\|\beta\|_{\frac{3}{2}} \leq t$$

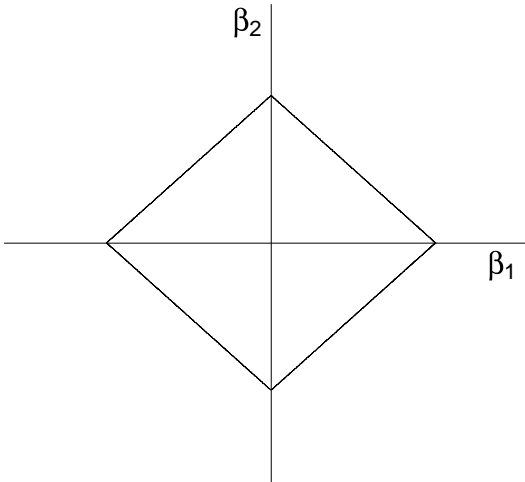


$$\|\beta\|_2 \leq t$$

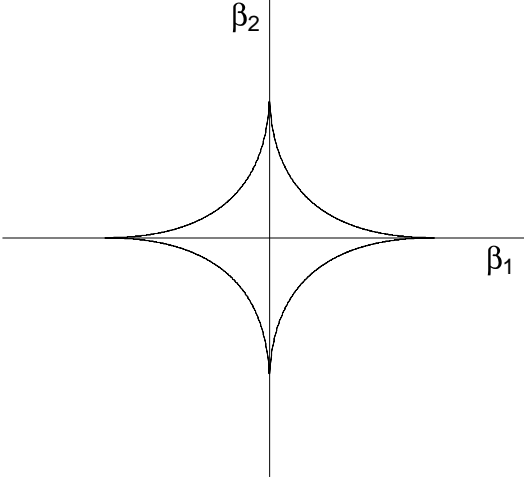
GEOMETRY OF REGULARIZATION IN \mathbb{R}^2 : CONVEXITY



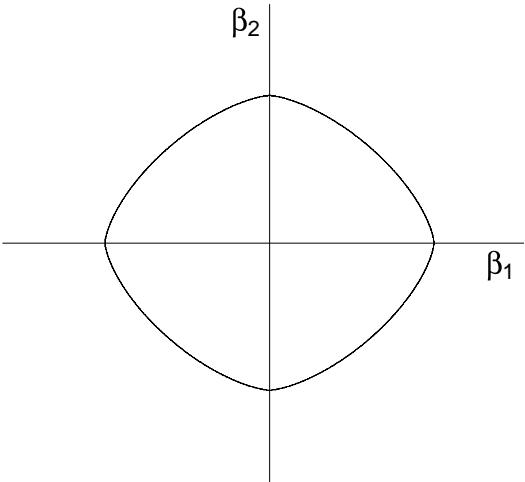
$$\|\beta\|_0 \leq t$$



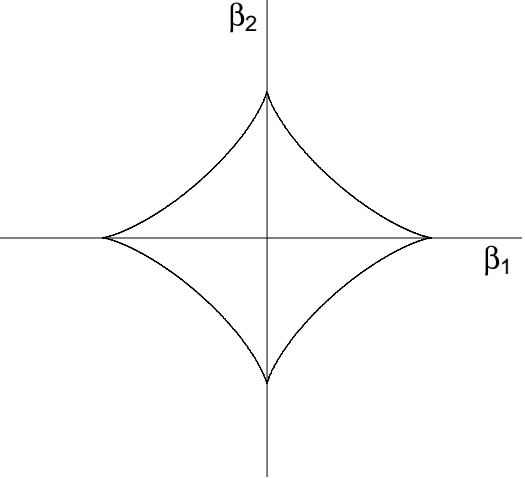
$$\|\beta\|_1 \leq t$$



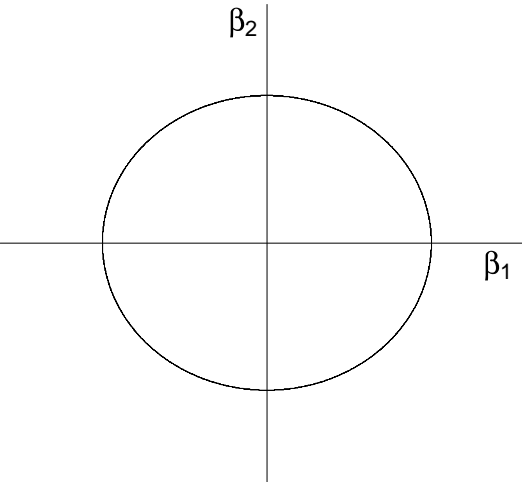
$$\|\beta\|_{\frac{1}{2}} \leq t$$



$$\|\beta\|_{\frac{3}{2}} \leq t$$

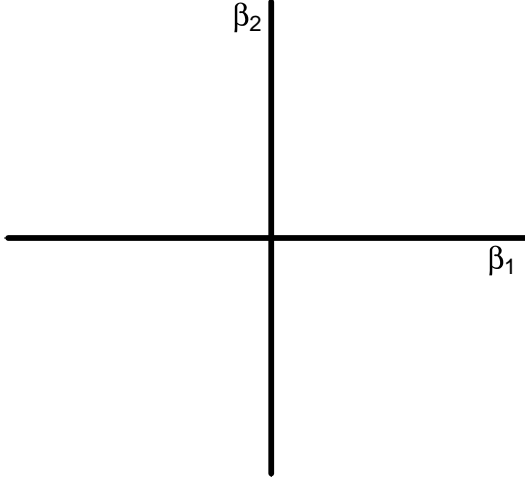


$$\|\beta\|_{\frac{3}{4}} \leq t$$

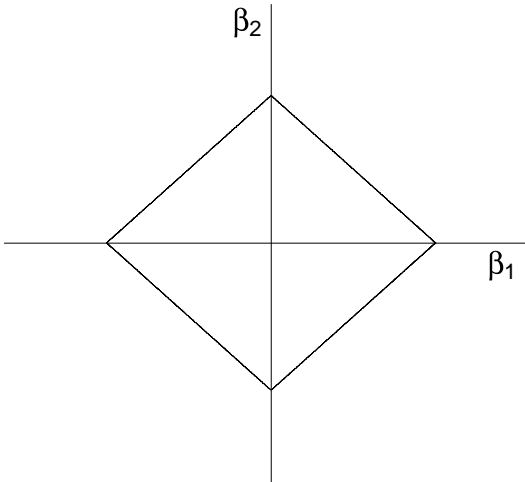


$$\|\beta\|_2 \leq t$$

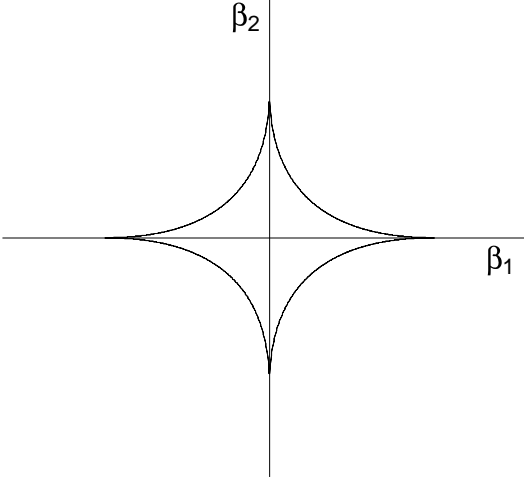
GEOMETRY OF REGULARIZATION IN \mathbb{R}^2 : CONVEXITY



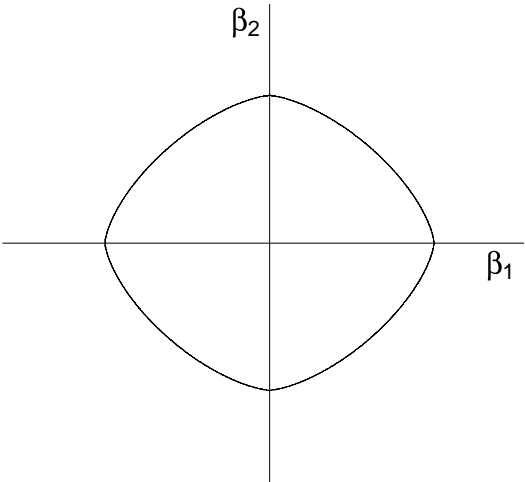
$$\|\beta\|_0 \leq t$$



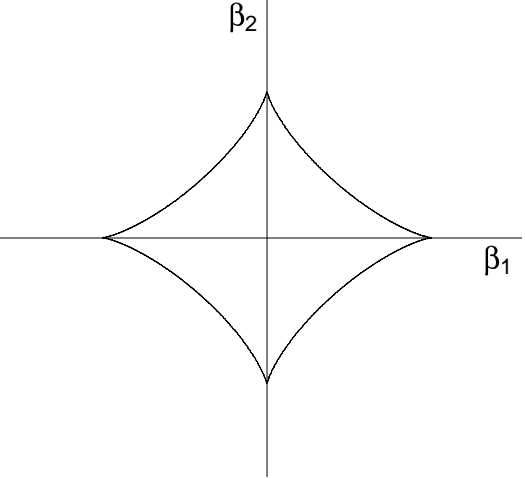
$$\|\beta\|_1 \leq t$$



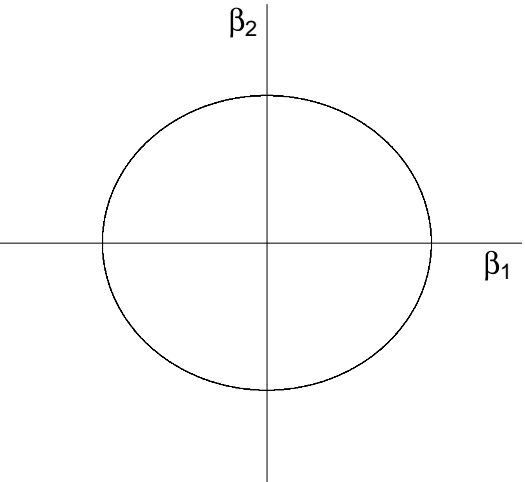
$$\|\beta\|_{\frac{1}{2}} \leq t$$



$$\|\beta\|_{\frac{3}{2}} \leq t$$

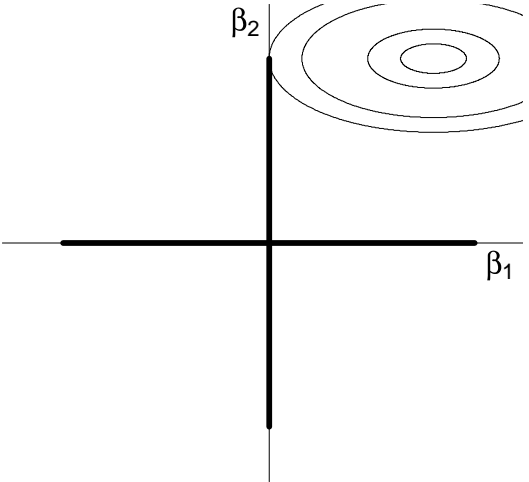


$$\|\beta\|_{\frac{3}{4}} \leq t$$

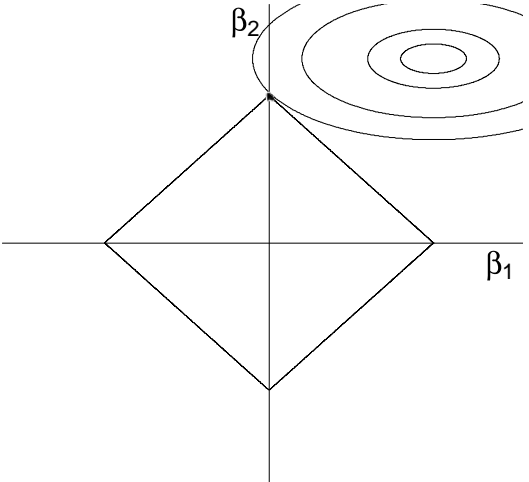


$$\|\beta\|_2 \leq t$$

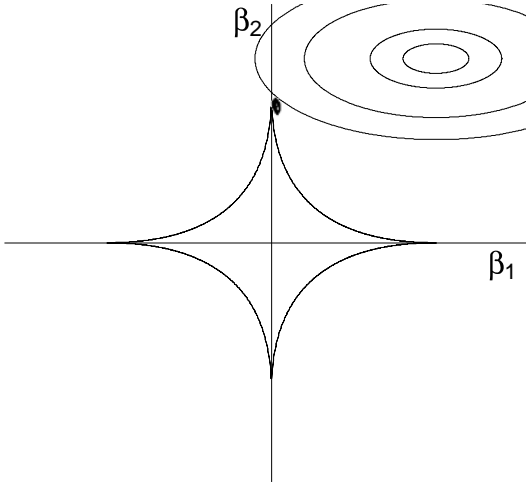
GEOMETRY OF REGULARIZATION IN \mathbb{R}^2 : MODEL SELECTION



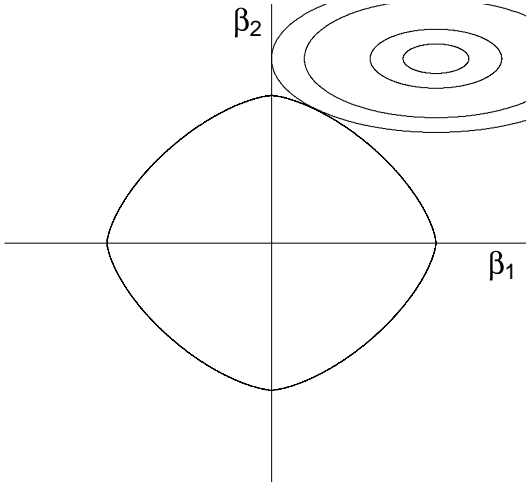
$$\|\beta\|_0 \leq t$$



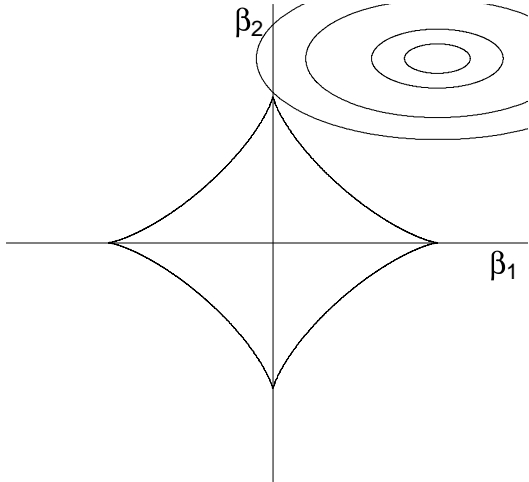
$$\|\beta\|_1 \leq t$$



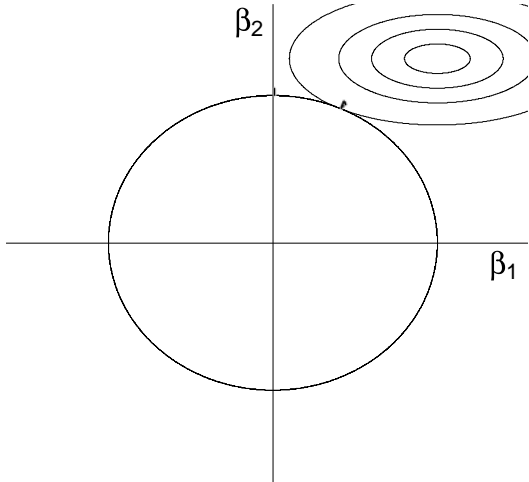
$$\|\beta\|_{\frac{1}{2}} \leq t$$



$$\|\beta\|_{\frac{3}{2}} \leq t$$

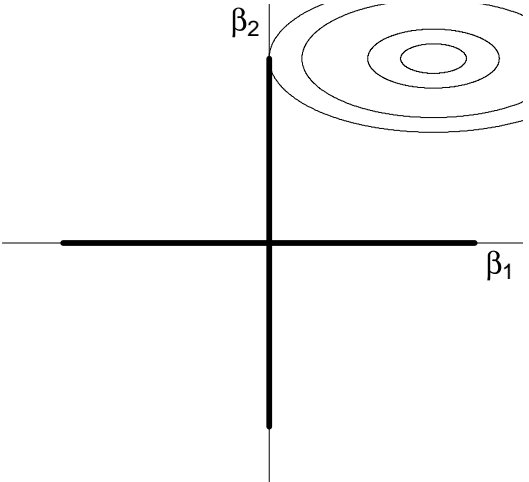


$$\|\beta\|_{\frac{3}{4}} \leq t$$

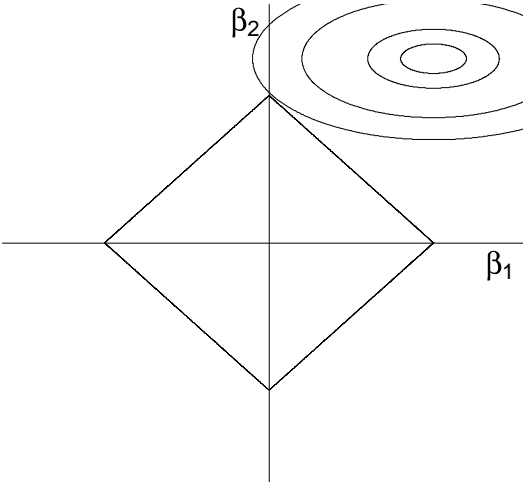


$$\|\beta\|_2 \leq t$$

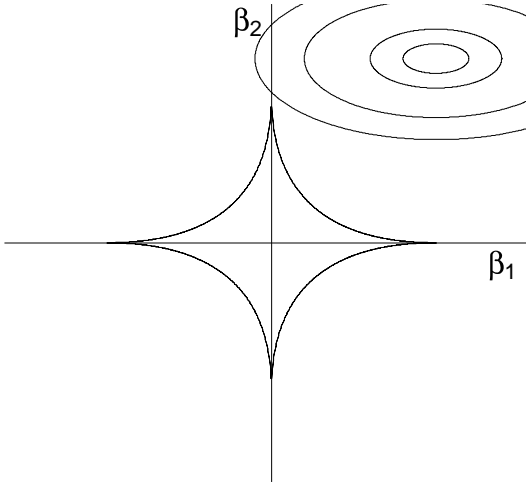
GEOMETRY OF REGULARIZATION IN \mathbb{R}^2 : MODEL SELECTION



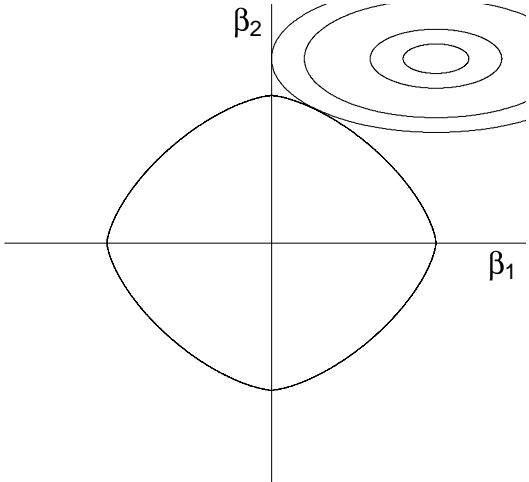
$\|\beta\|_0 \leq t$



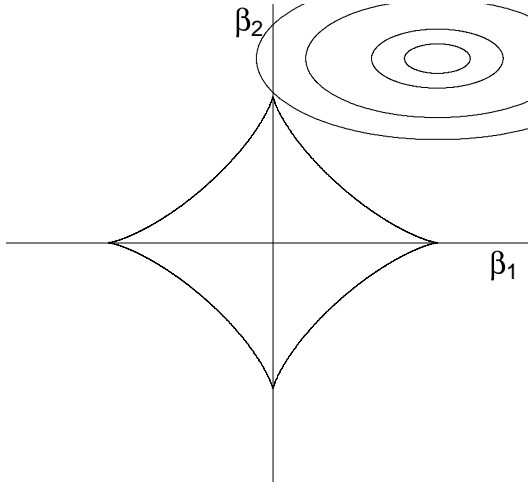
$\|\beta\|_1 \leq t$



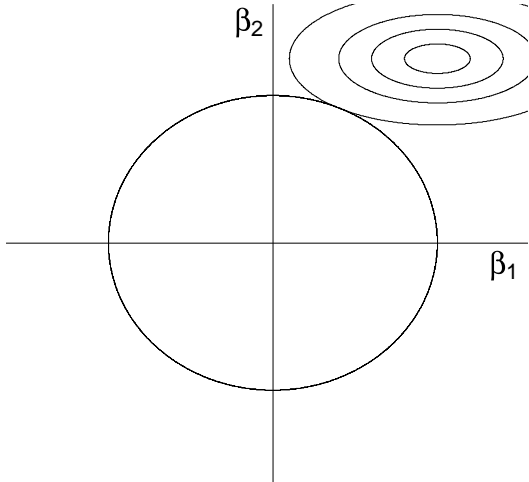
$\|\beta\|_{\frac{1}{2}} \leq t$



$\|\beta\|_{\frac{3}{2}} \leq t$

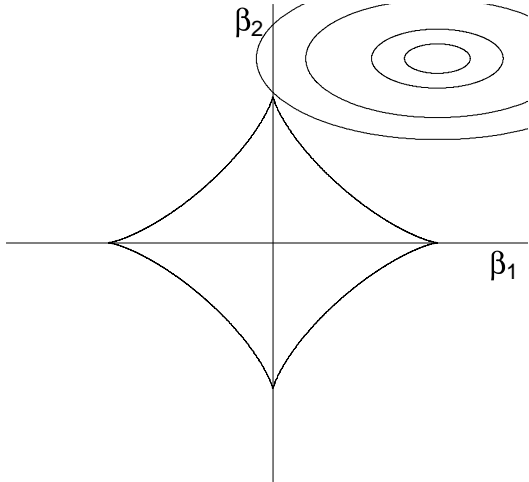
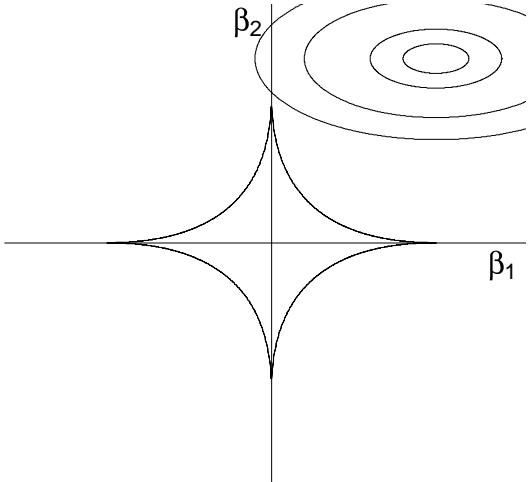
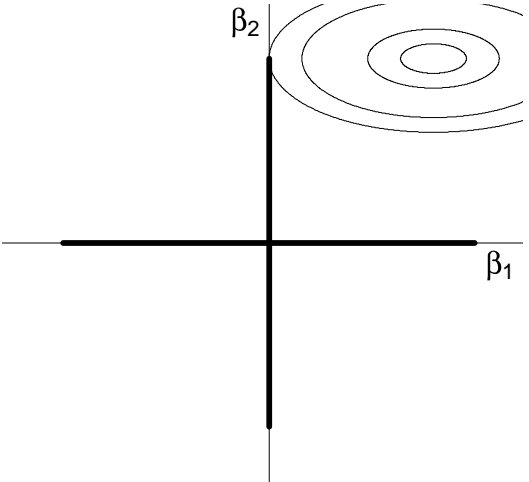


$\|\beta\|_{\frac{3}{4}} \leq t$

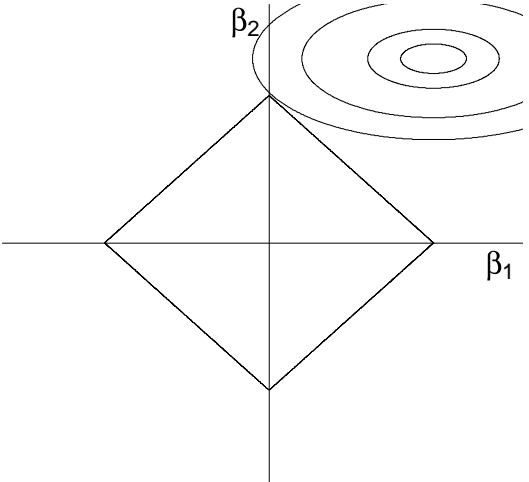


$\|\beta\|_2 \leq t$

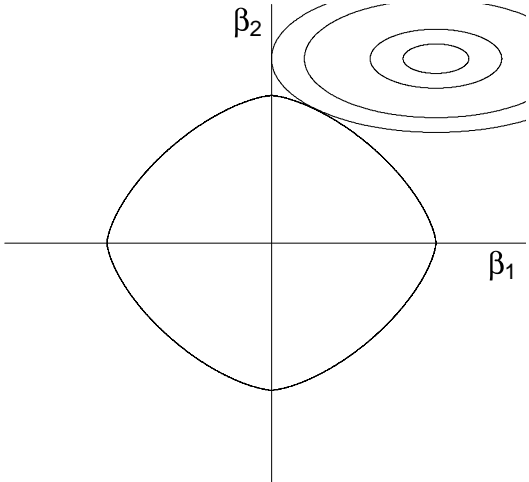
GEOMETRY OF REGULARIZATION IN \mathbb{R}^2 : MODEL SELECTION



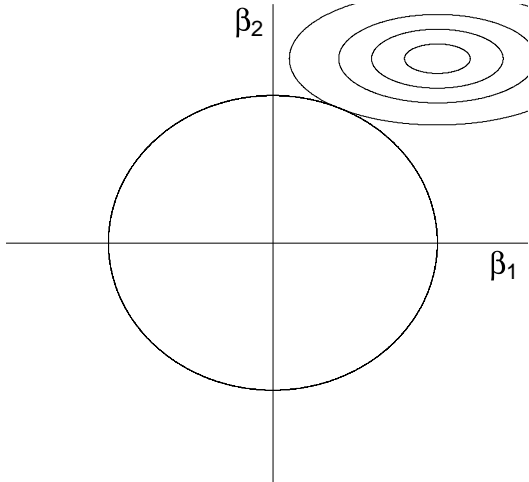
$$\|\beta\|_0 \leq t$$



$$\|\beta\|_{\frac{1}{2}} \leq t$$



$$\|\beta\|_{\frac{3}{4}} \leq t$$

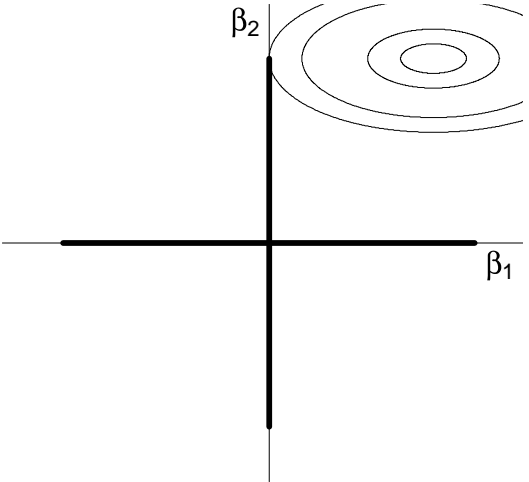


$$\|\beta\|_1 \leq t$$

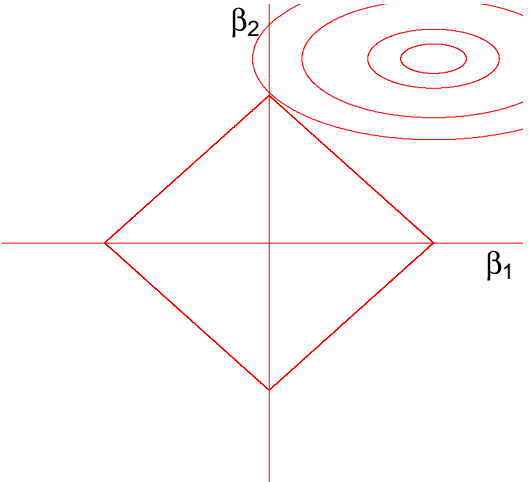
$$\|\beta\|_{\frac{3}{2}} \leq t$$

$$\|\beta\|_2 \leq t$$

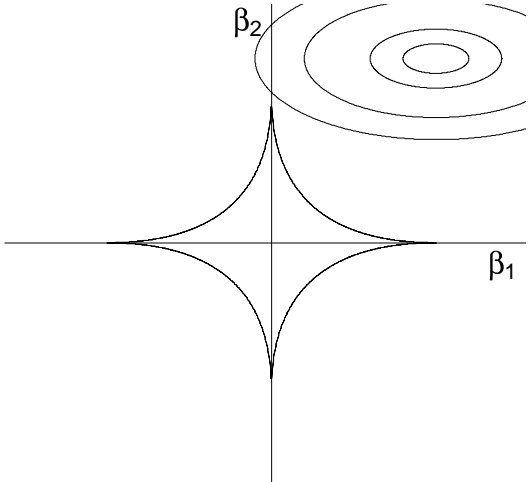
GEOMETRY OF REGULARIZATION IN \mathbb{R}^2 : BOTH



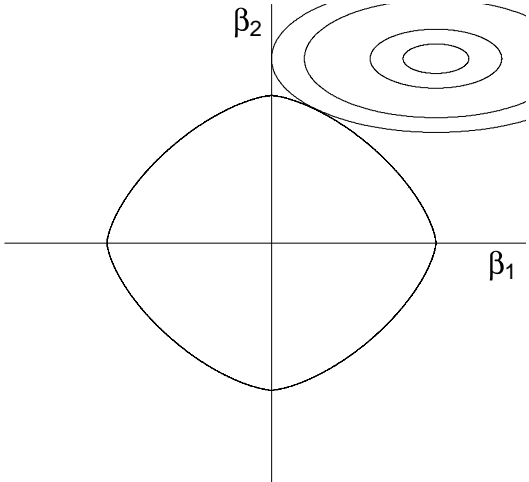
$$\|\beta\|_0 \leq t$$



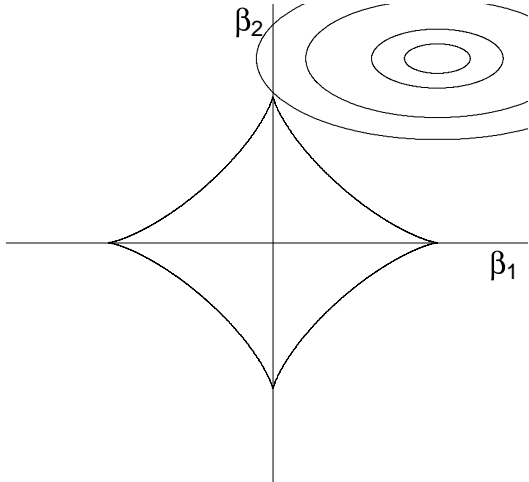
$$\|\beta\|_1 \leq t$$



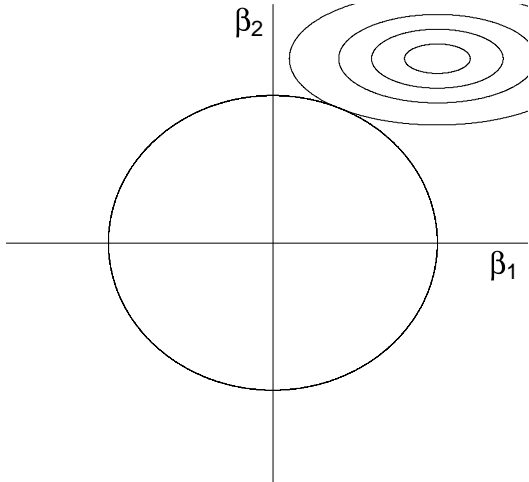
$$\|\beta\|_{\frac{1}{2}} \leq t$$



$$\|\beta\|_{\frac{3}{2}} \leq t$$



$$\|\beta\|_{\frac{3}{4}} \leq t$$

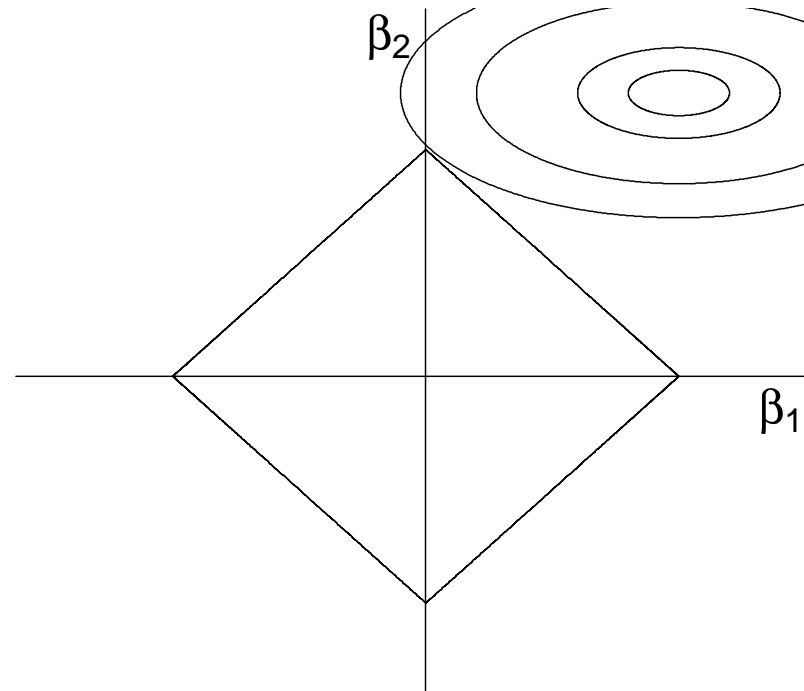


$$\|\beta\|_2 \leq t$$

SUMMARY

| | CONVEX? | CORNERS? | |
|---------------------------|---------|----------|---|
| $\ \beta\ _0$ | No | Yes | |
| $\ \beta\ _{\frac{1}{2}}$ | No | Yes | |
| $\ \beta\ _{\frac{3}{4}}$ | No | Yes | |
| $\ \beta\ _1$ | Yes | Yes | ✓ |
| $\ \beta\ _{\frac{3}{2}}$ | Yes | No | |
| $\ \beta\ _2$ | Yes | No | |

THE BEST OF BOTH WORLDS: $\|\beta\|_1$



This regularization set...

- ... is convex (computationally efficient)
- ... has corners (performs model selection)

Lasso in practice

ℓ_1 -REGULARIZED REGRESSION

Related methods are known as

- **LASSO**: The covariates are recorded
- **BASIS PURSUIT**: The covariates are frames comprised of various bases
- **COMPRESSED SENSING**: The covariates are random draws from some distribution

The estimator satisfies

$$\hat{\beta}_{lasso}(t) = \underset{\|\beta\|_1 \leq t}{\operatorname{argmin}} \|\mathbb{Y} - \mathbb{X}\beta\|_2^2$$

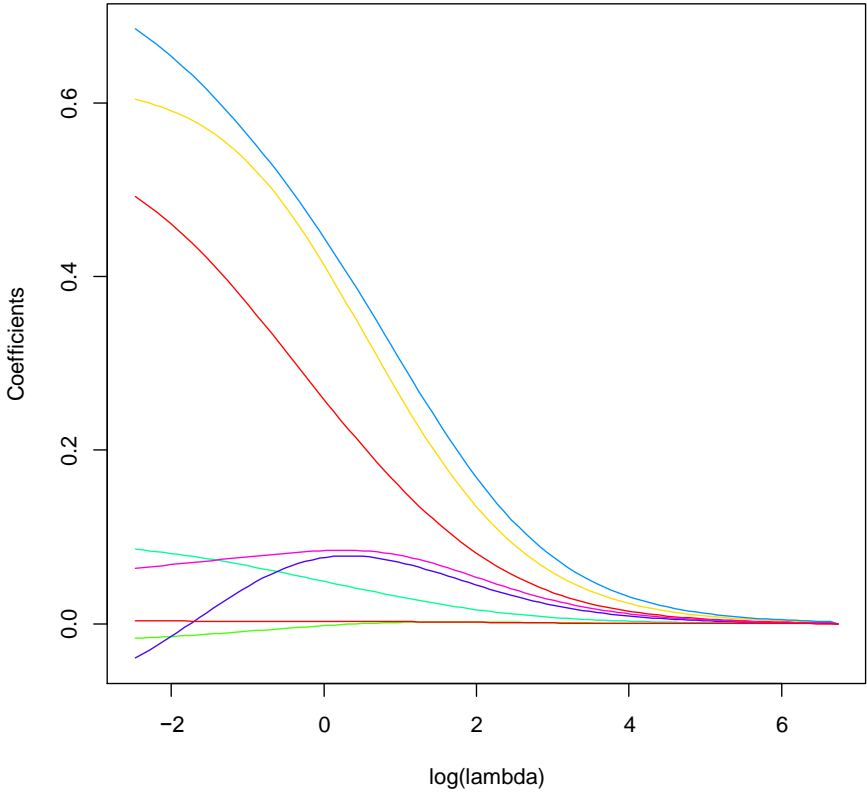
In its corresponding Lagrangian dual form:

$$\hat{\beta}_{lasso}(\lambda) = \underset{\beta}{\operatorname{argmin}} \left(\|\mathbb{Y} - \mathbb{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right)$$

(Note that if $\operatorname{rank}(X) < p$, then the objective function is not strictly convex. There are now an infinite number of possible lasso solutions. (all must have the same fitted value and $\|\cdot\|_1$))

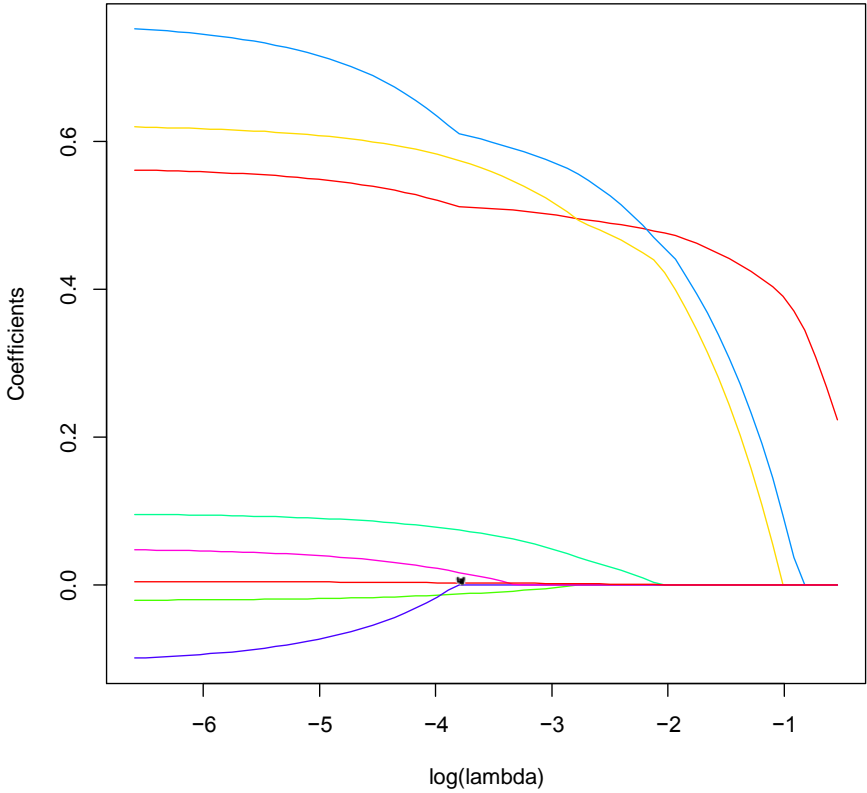
LASSO REGRESSION PATH

$\|\beta\|_2 \rightarrow 0$



Ridge

$\|\beta\|_1 \rightarrow 0$



Lasso

THE LASSO IN R: GLMNET

Luckily, we already know how to lasso.

Just change the 'alpha =0' to 'alpha =1', and you're lassoing.

```
lasso.out = glmnet(x=as.matrix(X), y=Y, alpha=1)
```

#Note: glmnet automatically scales X

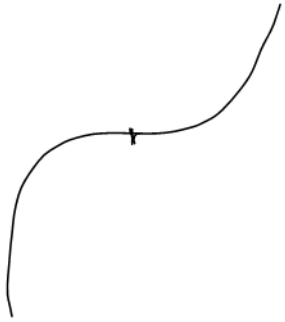
glmnet uses gradient descent to quickly fit the lasso solution

It can...

- handle other likelihoods than Gaussian
- supports/exploits sparse matrices (e.g. for text processing)
- use warm restarts for the grid of λ to produce more stable fits/faster computations

(See (Friedman et al. (2007) for details))

OPTIMALITY CONDITIONS: REVIEW



$$\text{minimize } F(x) \quad (3)$$

$$\text{subject to } x \in \mathbb{R}^p \quad (4)$$

Search for x_* such that $\nabla F|_{x_*} = 0$

- Turns a geometric problem into an algebraic problem: solve for the point where the gradient vanishes
- Is necessary for optimality of x_* . Is sufficient if F is convex and smooth.

GRADIENT DESCENT: INTUITION

A summary:

1. Start with some initial x^0
2. Propose x to reduce $F(x)$
3. Alternate between 1. and 2. until the objective function doesn't change (much).

Algorithmically, the implementations tend to look like

$$x[k + 1] \leftarrow x[k] + \alpha_k v[k],$$

where

- $x[k]$ is the current value of the minimizing parameter
- $v[k]$ is a direction that (hopefully) reduces F
- α_k is a relaxation term.

GRADIENT DESCENT

Assume $\exists x_* \in D$ such that $\nabla F|_{x_*} = 0$

Define the map

$$\psi(x) = x - \alpha \nabla F|_x$$

(Recall the general form $x[k+1] \leftarrow x[k] + \alpha_k v[k]$)

If ψ is **contractive**, ie

$$\|\psi(x) - \psi(x')\| \leq c \|x - x'\|$$

where $c \in [0, 1)$, then...

Gradient descent is guaranteed to converge

GRADIENT DESCENT: CONVERGENCE PROOF

$$\|x[k+1] - x_*\| = \|x[k] - \alpha \nabla F|_{x[k]} - x_*\| \quad (5)$$

$$= \|\psi(x[k]) - \psi(x_*)\| \quad (6)$$

$$\leq c \|x[k] - x_*\| \quad (7)$$

$$\vdots \quad (8)$$

$$\leq c^{k+1} \|x[0] - x_*\| \quad (9)$$

$$(10)$$

(This means we get exponential convergence²)

Important fact: If F is $2\times$ differentiable, contractivity means F is convex on D

²Optimization people call this linear convergence due to equation (7)

GRADIENT DESCENT EXAMPLE

If we look at multiple regression via least squares we get:

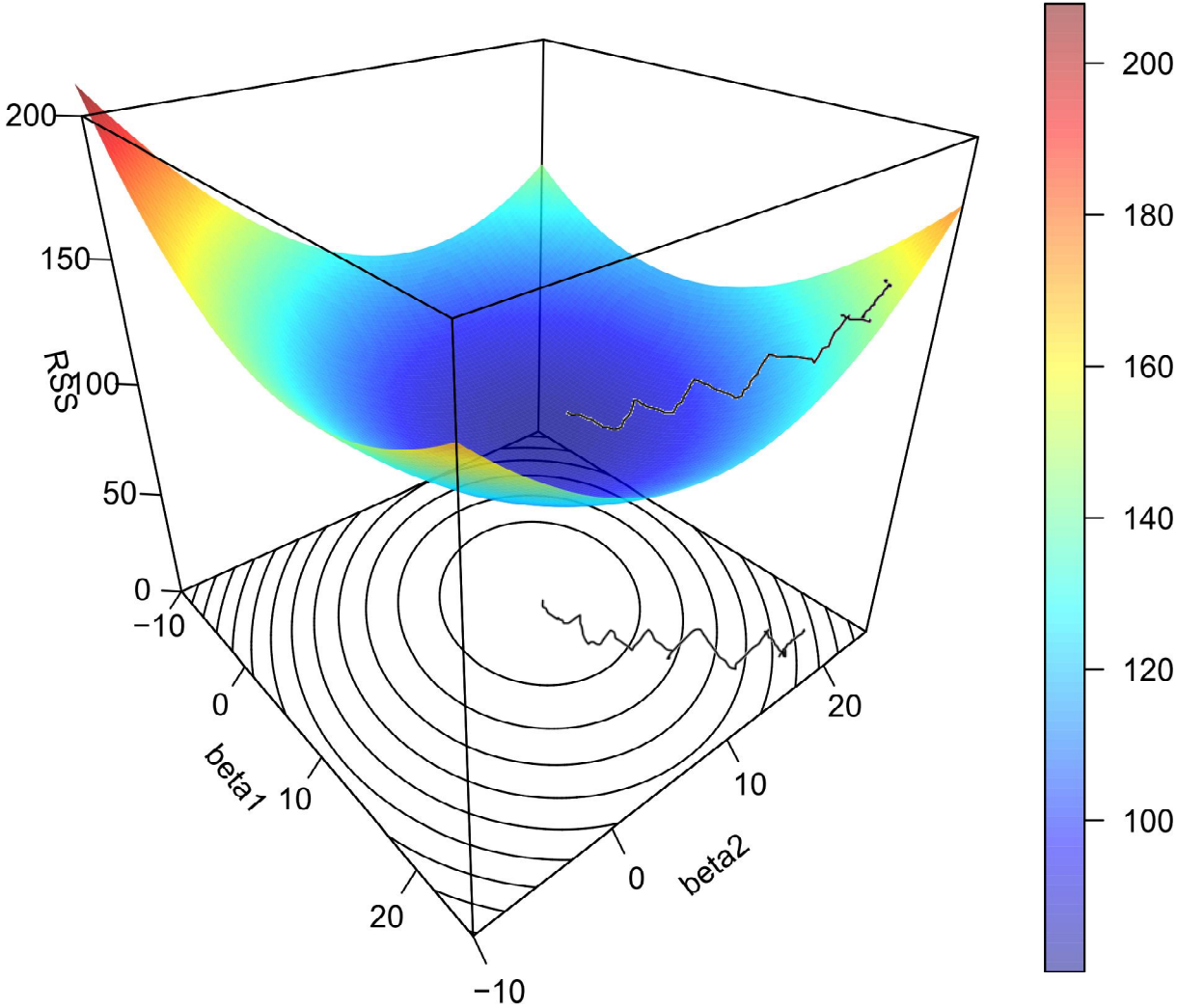
$$\begin{aligned} F(x) &\rightarrow \|Y - X\beta\|_2^2 \\ \min_{\beta} \|Y - X\beta\|_2^2 &\Rightarrow \frac{\partial}{\partial \beta_j} \|Y - X\beta\|_2^2 \\ &= \frac{\partial}{\partial \beta_j} \sum_{i=1}^n (Y_i - X_i^\top \beta)^2 \\ &= 2 \sum_{i=1}^n (Y_i - X_i^\top \beta) X_{ij} \end{aligned}$$

Hence, we will cycle over j and make the update $k = 1, \dots, K$ iterations:

$$\hat{\beta}_j^{k+1} = \hat{\beta}_j^k - \alpha \sum_{i=1}^n (Y_i - X_i^\top \hat{\beta}^k) X_{ij}$$

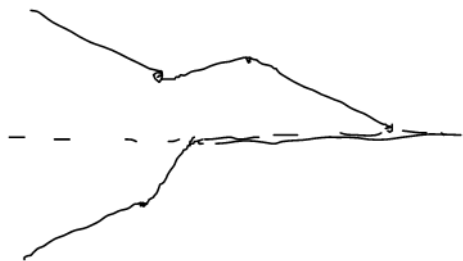
BATCH G.D.
STOCHASTIC G.D.

GRADIENT DESCENT EXAMPLE



With $RSS = \|Y - X\beta\|_2^2$ for $p = 2$

THE LASSO IN R: LARS



Alternatively, the `lars` package exploits the fact that the coefficient profiles are piecewise linear, which leads to an algorithm with the same computational cost as the full least-squares fit on the data (See Osborne et al. (2000) for details on the convex optimization, Efron et al. (2004) for the LARS algorithm)

CHOOSING THE TUNING PARAMETER FOR LASSO

Of course, just like in Ridge, we need a way of choosing this tuning parameter.

We can just use **cross-validation** again, though this is still an area of active research:

Homrighausen, D. and McDonald, D.J. *Leave-one-out cross-validation is risk consistent for lasso*, Machine Learning

Homrighausen, D. and McDonald, D.J. *Risk consistency of cross-validation for lasso-type procedures*, Journal of Machine Learning Research

Homrighausen, D. and McDonald, D.J. *The lasso, persistence, and cross-validation*, (2013) International Conference on Machine Learning, JMLR 28(3), 1031–1039.

CHOOSING THE TUNING PARAMETER FOR LASSO

For cross-validation, the heavy lifting has been done for us

```
cv.glmnet(x=as.matrix(X),y=Y,alpha=1)
```

```
cv.lars(x=as.matrix(X),y=Y,type='lasso')
```

$$\|a\|_{\infty} = \max_i |a_i|$$

Note that for the grid λ , we need only look over the interval

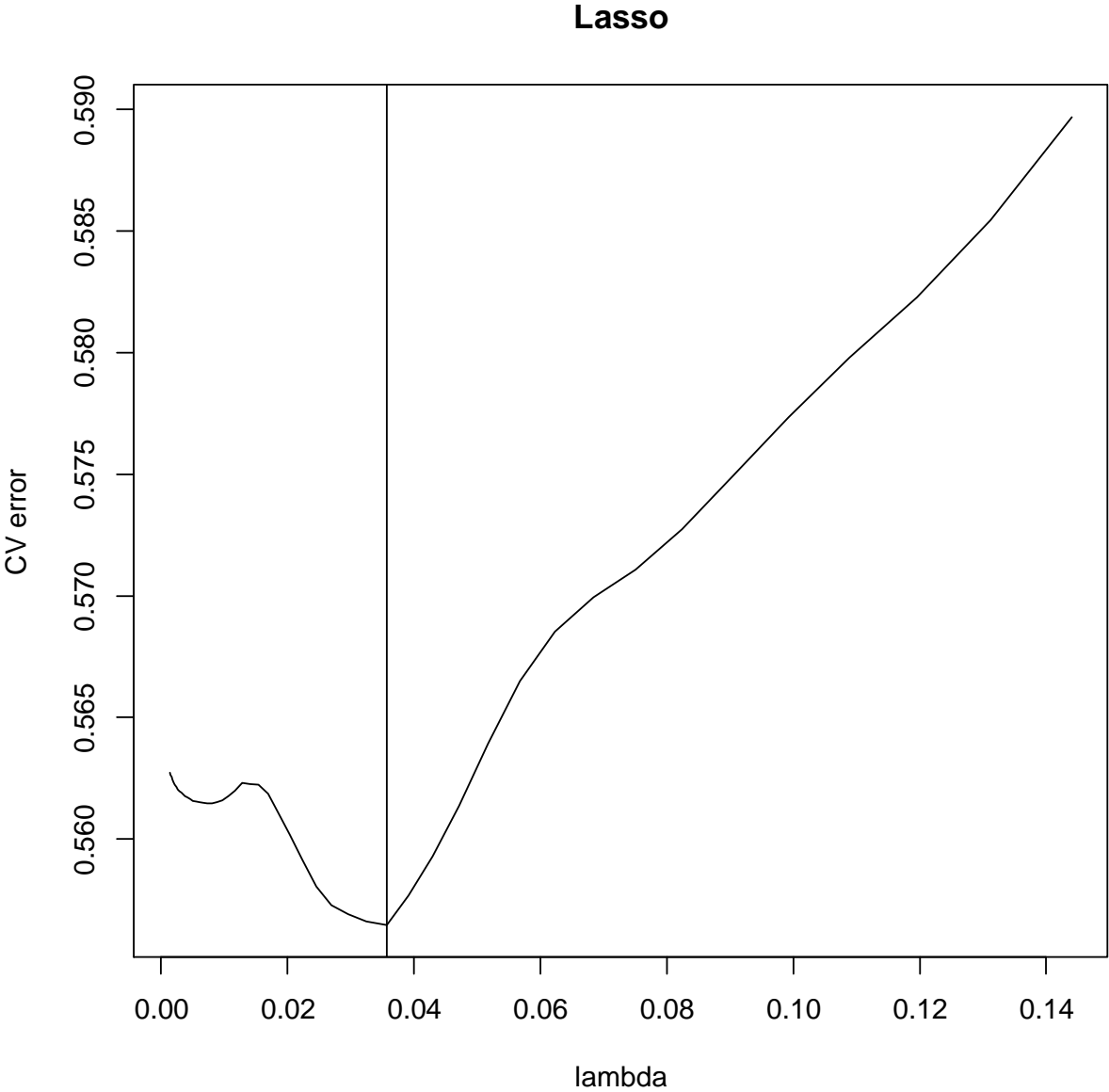
$$\left[0, \|\mathbb{X}^T Y\|_{\infty}\right) \quad \hat{\beta}_{\text{Lasso}, \lambda} = 0 \quad \lambda \geq \|\mathbb{X}^T Y\|_{\infty}$$

A grid of t has a similar restriction $[0, t_0)$, where

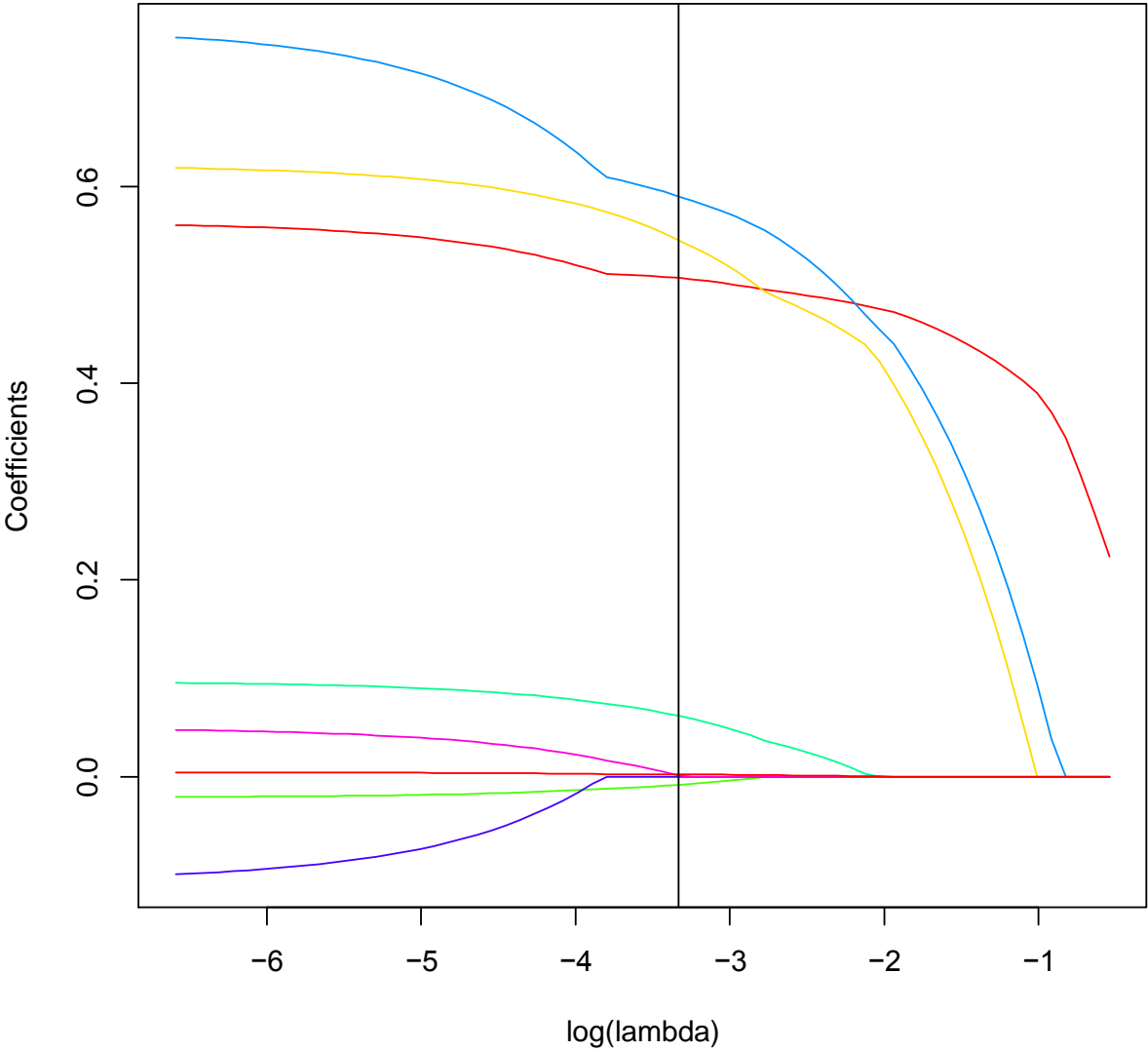
$$t_0 = \min_{\{b: \mathbb{X}b=0\}} \|\hat{\beta}_{\text{LS}} + b\|_1$$



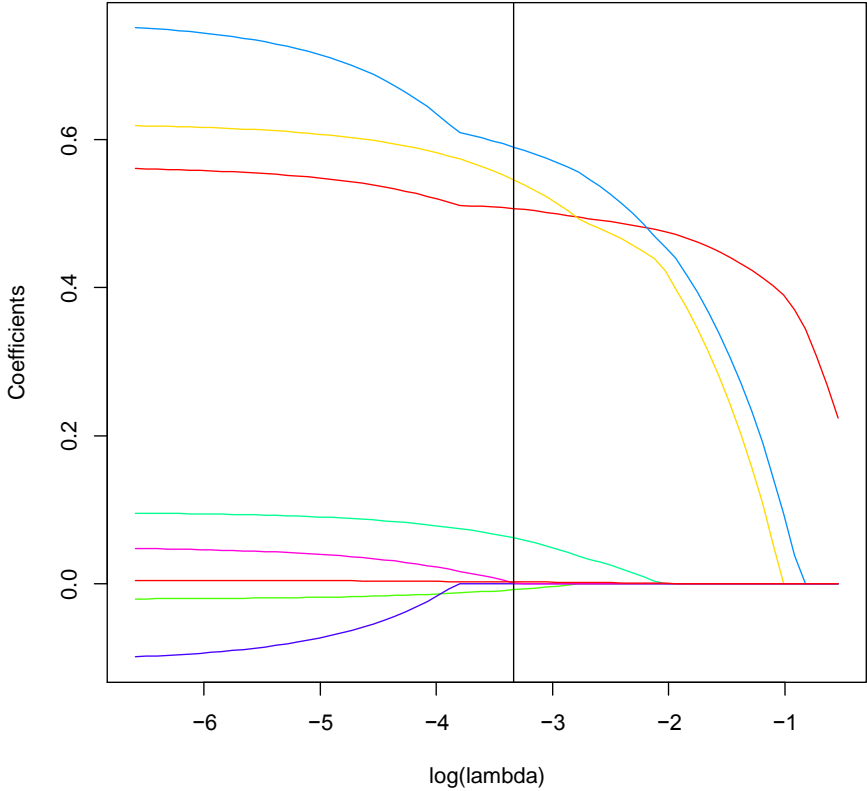
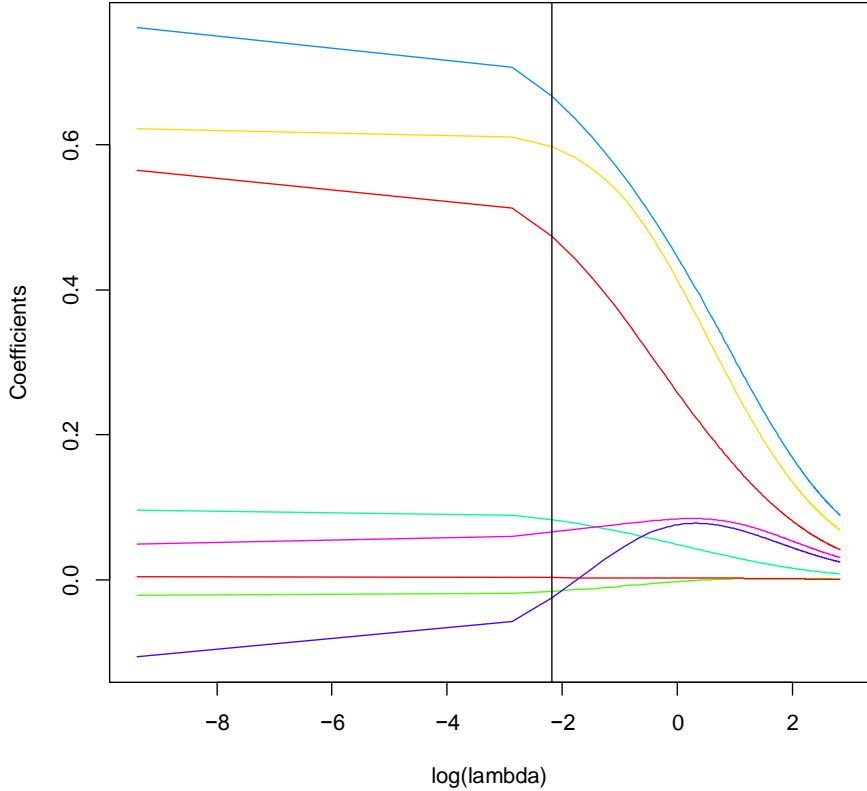
THE LASSO IN R



LASSO REGRESSION PATH



COMPARISON: REGRESSION PATH



Vertical line at minimum CV tuning parameter

COMPARISON OF LARS AND GLMNET

There are two main problems with **glmnet**

- In practice, the λ interval looks like $[\epsilon \|\mathbb{X}^\top Y\|_\infty, \|\mathbb{X}^\top Y\|_\infty)$ for a small ϵ . Sometimes, this results in finding a boundary solution.
- The iterative nature sometimes results in bad coefficient vectors (such as having more than $\min\{n, p\}$ nonzero coefficients, which is impossible³)

$$\|\hat{\beta}^{(k)} - \hat{\beta}^{(k+1)}\| < \epsilon$$

There are two main problems with **lars**

- It is slow(er)
- It doesn't support other likelihoods

³This is not quite true (Tibshirani (2013), Lemma 13). However, see Lemma 15 in same paper: For any \mathbb{X} , λ and *almost* all Y , the column space of $\mathbb{X}_{\mathcal{S}}$ is the same for every \mathcal{S} , where $\mathcal{S} = \{j : |\hat{\beta}_{\lambda,j}| > 0\}$

FLAVORS OF LASSO

- Grouped lasso (Yuan and Lin (2007), Meier et al. (2008)), where variables are included or excluded in groups.
- Refitted lasso (e.g. Lederer 2013). Takes the estimated model from lasso and fits the full least squares solution on selected covariates (less bias, more variance).
- Dantzig selector (Candes, Tao (2007)), a slightly modified version of the lasso
- The elastic net (Zou, Hastie (2005)), generally used for correlated variables that combines a ridge/lasso penalty. Included in **glmnet**. Fixes non-uniqueness problem of lasso (although, see Tibshirani (2013)).
- SCAD (Fan and Li (2005)), a non-convex version of lasso that adds a more severe variable selection penalty $\|y - \tilde{X}\beta\|_2 \rightarrow \|y - \tilde{X}\beta\|_2$
- $\sqrt{\text{lasso}}$ (Belloni et al. (2011)), claims to be tuning parameter free (but isn't). Uses $\|\cdot\|_2$ instead of $\|\cdot\|_2^2$ for the loss.
- Generalized lasso (Tibshirani, Taylor (2011)). Adds various additional penalty matrices to the penalty term (ie: $\|D\beta\|_1$).