# Statistical Learning with Sparsity Generalizations of loss functions and lasso penalties

Boen Jiang

Applied Statistics @ Fudan University

September 16, 2025

## Lasso: recap

## Least Absolute Shrinkage and Selection Operator (lasso)

Recall that the lasso estimate is defined by

$$\hat{\beta}^{\mathsf{lasso}}(\lambda) = \operatorname*{argmin}_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2N} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \|_2^2 + \lambda \| \boldsymbol{\beta} \|_1 \right\},$$

where  $\lambda \geq 0$  is a tuning parameter that controls the amount of shrinkage.

For the red part of the objective function, we can generalize the square loss to other loss functions:

- Negative log-likelihood for GLMs
- Negative log-partial-likelihood for Cox models
- Hinge loss for SVM

For the blue part of the objective function, we can generalize the lasso penalty to other penalties:

- Elastic net
- Group lasso
- Fused lasso

## Part I

- Negative log-likelihood for GLMs
- Negative log-partial-likelihood for Cox models
- Hinge loss for SVM

# Binary response variable

- For simplicity, we can still use the linear model for a binary outcome
- Linear probability model  $y_i = x_i^{\top} \beta + \varepsilon_i$ ,  $E(\varepsilon_i \mid x_i) = 0$

$$P(y_i = 1 \mid x_i) = E(y_i \mid x_i) = x_i^{\top} \beta.$$

Easy interpretation

$$\frac{\partial P(y_i = 1 \mid x_i)}{\partial x_{ij}} = \beta_j$$

However, there are two defects:

- Heteroskedasticity  $\operatorname{Var}(y_i \mid x_i) = x_i^{\top} \beta (1 x_i^{\top} \beta).$
- Not natural for binary outcome because probability is bounded between zero and one.

## **GLMs**

A generalized linear model is made up of a linear predictor

$$\eta_i = \beta_0 + \beta_1 x_{1i} + \ldots + \beta_p x_{pi}$$

and two functions:

• link function that describes how the mean,  $\mathbb{E}(Y_i) = \mu_i$ , depends on the linear predictor

$$g(\mu_i) = \eta_i$$

• variance function that describes how the variance,  $Var(Y_i)$ , depends on the mean

$$\mathsf{Var}(Y_i) = \phi V(\mu)$$

where the dispersion parameter  $\phi$  is a constant

#### Link functions

We can use a monotone transformation to force the linear predictor to lie within the interval [0,1]:

$$P(y_i = 1|x_i) = g(x_i^{\top}\beta).$$

Here, the inverse of g is called the **link function**.

There are some canonical choices of g:

- ullet Logit link:  $g(t)=rac{e^t}{1+e^t}=rac{1}{1+e^{-t}}$ , c.f. standard logistic distribution
- Probit link:  $g(t) = \Phi(t)$ , c.f. standard normal distribution
- Complementary log-log link:  $g(t) = 1 e^{-e^t}$ , c.f. standard log-Weibull distribution
- Cauchit link:  $g(t) = \frac{1}{\pi}\arctan(t) + \frac{1}{2}$ , c.f. standard Cauchy distribution

## Modelling Binomial Data

Suppose

$$Y_i \sim \text{Binomial}(n_i, p_i)$$

and we wish to model the proportions  $Y_i/n_i$ . Then

$$\mathbb{E}(Y_i/n_i) = p_i \quad \text{Var}(Y_i/n_i) = \frac{1}{n_i} p_i (1 - p_i)$$

So our variance function is

$$V(\mu_i) = \mu_i (1 - \mu_i)$$

Our link function must map from  $(0,1) \to (-\infty,\infty)$ . A common choice is

$$g(\mu_i) = \operatorname{logit}(\mu_i) = \operatorname{log}\left(\frac{\mu_i}{1 - \mu_i}\right)$$

## Modelling Poisson Data

Suppose

$$Y_i \sim \text{Poisson } (\lambda_i)$$

Then

$$\mathbb{E}(Y_i) = \lambda_i \quad Var(Y_i) = \lambda_i$$

So our variance function is

$$V(\mu_i) = \mu_i$$

Our link function must map from  $(0,\infty) \to (-\infty,\infty)$ . A natural choice is

$$g(\mu_i) = \log(\mu_i)$$

# One-parameter Canonical Exponential Family

• Canonical exponential family for  $k = 1, y \in \mathbb{R}$ 

$$f_{ heta}(y) = \exp\left(rac{y heta - b( heta)}{\phi} + c(y,\phi)
ight)$$

for some known functions  $b(\cdot)$  and  $c(\cdot, \cdot)$ .

- If  $\phi$  is known, this is a one-parameter exponential family with  $\theta$  being the canonical parameter.
- If  $\phi$  is unknown, this may/may not be a two-parameter exponential family.  $\phi$  is called dispersion parameter.
- ullet We assume that  $\phi$  is known.
- The function g that links the mean  $\mu$  to the canonical parameter  $\theta$  is called Canonical Link:

$$g(\mu) = \theta$$

## Expectation

Note that

$$\ell(\theta) = \frac{Y\theta - b(\theta)}{\phi} + c(Y; \phi),$$

Therefore

$$\frac{\partial \ell}{\partial \theta} = \frac{Y - b'(\theta)}{\phi}$$

It yields

$$0 = \mathbb{E}\left(\frac{\partial \ell}{\partial heta}\right) = \frac{\mathbb{E}(Y) - b'( heta)}{\phi}$$

which leads to

$$\mathbb{E}(Y) = \mu = b'(\theta)$$

#### Variance

On the other hand we have

$$\frac{\partial^2 \ell}{\partial \theta^2} + \left(\frac{\partial \ell}{\partial \theta}\right)^2 = -\frac{b''(\theta)}{\phi} + \left(\frac{Y - b'(\theta)}{\phi}\right)^2$$

and from the previous result,

$$\frac{Y-b'( heta)}{\phi}=\frac{Y-\mathbb{E}(Y)}{\phi}$$

Together, with the second identity, this yields

$$0 = -\frac{b''(\theta)}{\phi} + \frac{\operatorname{var}(Y)}{\phi^2}$$

which leads to

$$var(Y) = V(Y) = b''(\theta)\phi$$

## Negative log-likelihood

Now we consider minimize negative log-likelihood with a penalty:

$$\underset{\beta_{0},\beta}{\mathsf{minimize}} \left\{ -\frac{1}{\textit{N}} \mathcal{L} \left(\beta_{0},\beta; \textbf{\textit{y}}, \textbf{\textit{X}}\right) + \lambda \|\beta\| \right\}$$

where the type of norm is specified in the problem. We consider the linear model as an example of GLM. Assuming  $Y|X=x\sim\mathcal{N}(\mu(x),\sigma^2)$ . Then:

$$\mathcal{L}(\beta_0, \beta; \mathbf{y}, \mathbf{X}) = -\sum_{i=1}^{N} \frac{(y_i - \beta_0 - \beta x_i)^2}{2\sigma^2} + c = -\frac{\|\mathbf{y} - \beta_0 - \beta \mathbf{X}\|_2^2}{2\sigma^2 N} + c$$

where c is a constant that does not depend on  $\beta_0$  and  $\beta$ .

Hence, negative log-likelihood is equivalent to the square error loss in this case.

#### Remarks on negative log-likelihood

Why? Under regular conditions, the Fisher information matrix is positive definite, so the negative log-likelihood is a convex function of the parameter.

12 / 12

## An example for classification

Suppose we take  $Y_i \in \{+1, -1\}$ , namely,  $P(Y_i = 1) = \pi_i, P(Y_i = -1) = 1 - \pi_i$ , and  $\operatorname{logit}(\pi_i) = \beta_0 + \beta_1^T x_i, i = 1, ..., n$ .

Then the likelihood function is

$$\mathcal{L}(\pi|y_1,...,y_n) = \prod_{i=1}^n \frac{1}{1 + \left(\frac{1-\pi_i}{\pi_i}\right)^{y_i}}.$$

After some trivial algebras, we can get the following negative log-likelihood function with a penalty:

$$\frac{1}{N} \sum_{i=1}^{N} \log \left( 1 + e^{-y_i \left( \beta_0 + \beta^\top x_i \right)} \right) + \lambda \|\beta\|_1.$$

Here,  $y_i \left(\beta_0 + \beta^\top x_i\right)$  can be interpreted as the **margin** of the *i*-th observation, for which positive values indicate correct classification and negative values indicate incorrect classification.

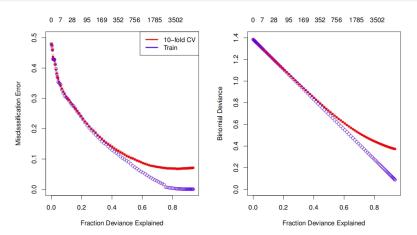
13 / 126

- The 20 Newsgroups data set is a collection of approximately 20,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups. It was originally collected by Ken Lang, for his Newsweeder: Learning to filter netnews paper, though he does not explicitly mention this collection. The 20 newsgroups collection has become a popular data set for experiments in text applications of machine learning techniques, such as text classification and text clustering.
- Koh, Kim, Boyd (2007) analysis the data set by logistic regression interior point method.
- Freidman, Hastie, Tibshirani (2010) use this data to illustrate their glmnet via coordinate descent.
- A team at Renmin U and Tsinghua U (2022) developed a Bayesian method for classification and summerization, their work published on Journal of Machine Learning Research.

Table 1: Class Groups, Classes and the Numbers of Documents in Each Class

Class	Class	Class Name	No.	No.
Group	Number		Training	Test
Computer Science	1	comp.graphics	584	389
	2	comp.os.ms-windows.misc	591	394
	3	comp.sys.ibm.pc.hardware	590	392
	4	comp.sys.mac.hardware	578	385
	5	comp.windows.x	593	395
For Sale	6	misc.forsale	585	390
Auto & Sports	7	rec.autos	594	396
	8	rec.motorcycles	598	398
	9	rec.sport.baseball	597	397
	10	rec.sport.hockey	600	399
Science	11	sci.crypt	595	396
	12	sci.electronics	591	393
	13	sci.med	594	396
	14	sci.space	593	394
Politics	15	talk.politics.guns	546	364
	16	talk.politics.mideast	564	376
	17	talk.politics.misc	465	310
Religion	18	alt.atheism	480	319
	19	soc.religion.christian	599	398
	20	talk.religion.misc	377	251

- We have N=11314 documents that we want to classify into two different groups  $(Y \in \{-1, +1\})$ .
- The features are defined as the set of **trigrams** (with some restrictions). In NLP, trigrams mean a sequence of three adjacent elements from a string of tokens. We have p=777811 features in total.
- Each document contains an average of 425 nonzero features. So this is a sparse problem.
- We want to perform  $\ell_1$  regularized logistic regression.



You can see that overfitting occurs when  $\lambda$  is too small, or equivalently, fraction deviance explained is too large, namely, the model is too "saturated".

The **fraction deviance explained**  $(D_{\lambda}^2)$  is then defined by:

$$\begin{split} D_{\lambda}^2 &= \frac{\mathrm{Dev_{null}} - \mathrm{Dev_{\lambda}}}{\mathrm{Dev_{null}}} \\ R^2 &= \frac{\mathrm{SS_{tot}} - \mathrm{SS_{res}}}{\mathit{SS_{tot}}} \end{split}$$

Deviance:  $(\text{Dev}_{\lambda})$  is defined as minus twice the difference in the log-likelihood for a model fit with parameter  $\lambda$  and the "saturated model" (having  $\hat{y} = y_i$ ).

#### Remarks on GOF statistics

It also reminds me conditional MSE in Guorong Dai's setup:

$$\frac{\mathbb{E}\{(Y-g(U))^2|S=s\}}{\mathbb{E}\{(Y-d(X))^2|S=s\}}.$$

So when we comparing two models, a natural choice is to find a proper ratio of the "errors" of two models.

# Computational techniques

ullet With current estimate  $\left( ilde{eta}_0, ilde{eta} 
ight)$ , we form the quadratic function:

$$Q(\beta_0, \beta) = \frac{1}{2N} \sum_{i=1}^{N} w_i \left( z_i - \beta_0 - \beta^T x_i \right)^2 + C\left( \tilde{\beta}_0, \tilde{\beta} \right),$$

• *C* denotes a constant independent of  $(\beta_0, \beta)$ ,  $z_i$  and  $w_i$  are defined as:

$$z_i = \tilde{\beta}_0 + \tilde{\beta}^T x_i + \frac{y_i - \tilde{p}(x_i)}{\tilde{p}(x_i)(1 - \tilde{p}(x_i))}, \quad \text{and} \quad w_i = \tilde{p}(x_i)(1 - \tilde{p}(x_i))$$

where  $\tilde{p}(x_i)$  is the current estimate for  $\Pr(Y = 1 \mid X = x_i)$ 

• Each outer loop then amounts to a weighted lasso regression

## Multiple outcomes

#### Setting

 $Y \in \{1,...,K\}$  for K > 2 classes. There are two natural ways for reduction to binary classification in general:

- OvO (One versus One): all  $\binom{K}{2}$  pairs of classes samples are used to fit  $\binom{K}{2}$  binary classifiers, then the predicted class is the one which is predicted the most.
- OvA (One versus All): treat all other classes as a single negative class.

#### Drawbacks

- OvO: computationally exhaustive and cases where same amount of votes for more classes.
- OvA: imbalance amounts positive and negative observations.

## Multinomial distribution

- Suppose we have nomial variable: a categorical variable that does not have intrinsic ordering or ranking, e.g., gender, colors, marital status, race, blood types
- Multinomial distribution

$$y \sim \text{Multinomial}(\pi_1, \dots, \pi_K), \quad \sum_{k=1}^K \pi_k = 1$$

y takes values in  $\{1,...,K\}$  with probability  $P(y=k)=\pi_k$ .

• We want to model  $y_i$  based on covariates  $x_i$ .

$$y_i \mid x_i \sim \text{ Multinomial } \left[1; \left\{\pi_1\left(x_i\right), \dots, \pi_K\left(x_i\right)\right\}\right],$$
 
$$\sum_{i=1}^K \pi_k\left(x_i\right) = 1 \text{ for all } x_i.$$

$$\pi_{y_i}(x_i) = \prod_{k=1}^{K} \{\pi_k(x_i) \text{ if } y_i = k\} = \prod_{k=1}^{K} \{\pi_k(x_i)\}^{\mathbb{I}(y_i = k)}.$$

# Multinomial logistic/softmax regression

• View the category K as the reference level, we model the ratios of the probabilities of category k and K

$$\log \frac{\pi_k(x_i)}{\pi_K(x_i)} = \beta_{0k} + x_i^{\top} \beta_k \quad (k = 1, ..., K - 1)$$
$$\pi_k(x_i) = \pi_K(x_i) e^{\beta_{0k} + x_i^{\top} \beta_k}$$

• Denote  $\beta_K = 0 \rightsquigarrow \text{reference level}$ 

$$\sum_{k=1}^{K} \pi_{k}(x_{i}) = 1 \Longrightarrow \pi_{K}(x_{i}) \sum_{k=1}^{K} e^{\beta_{0k} + x_{i}^{\top} \beta_{k}} = 1$$

$$\Longrightarrow \pi_{K}(x_{i}) = 1 / \sum_{k=1}^{K} e^{\beta_{0k} + x_{i}^{\top} \beta_{k}}$$

$$\Longrightarrow \pi_{k}(x_{i}) = \frac{e^{\beta_{0k} + x_{i}^{\top} \beta_{k}}}{\sum_{l=1}^{K} e^{\beta_{0l} + x_{i}^{\top} \beta_{l}}}$$

# Multinomial logistic with penalty

• In traditional multinomial regression, one category must be chosen as a reference group (i.e., have  $\beta_k$  set to  $\mathbf{0}$  ) or else the problem is not identifiable  $\leadsto \{\beta_{kj}+c_j\}_{k=1}^K$  and  $\{\beta_{kj}\}_{k=1}^K$  produce the same likelihood

Instead of traditional multinomial logistic regression, we can consider the following over specified version:

$$P(Y = k \mid X = x) = \frac{e^{\beta_{0k} + \beta_k^T x}}{\sum_{\ell=1}^{K} e^{\beta_{0\ell} + \beta_{\ell}^T x}}.$$

- we regularize the coefficients, and the regularized solutions are not equivariant under base/reference changes,
- the regularization automatically eliminates the redundancy
- The penalty term is  $\lambda \sum_{k=1}^{K} \|\beta_k\|_1$

# Multinomial logistic with vectorization

Log-likelihood form is

$$\log P\left(Y = y_i \mid x_i\right) = \beta_{0,y_i} + \beta_{y_i}^{\top} x_i - \log \left(\sum_{k=1}^K e^{\beta_{0k} + \beta_k^{\top} x_i}\right).$$

- Let  $r_{ik} := \mathbb{I}(y_i = k), \mathbf{R}_{N \times K} := (r_{ij})$
- We have  $\beta_{0,y_i} + \beta_{y_i}^\top x_i = \sum_{k=1}^K r_{ik} \left( \beta_{0k} + \beta_k^\top x_i \right)$
- Hence

$$\log P(Y = y_i \mid x_i) = \sum_{k=1}^K r_{ik} \left(\beta_{0k} + \beta_k^\top x_i\right) - \log \left(\sum_{k=1}^K e^{\beta_{0k} + \beta_k^\top x_i}\right)$$

# Multinomial logistic with penalty

Then we can write the log-likelihood in the more explicit form

$$\frac{1}{N} \sum_{i=1}^{N} w_i \left[ \sum_{k=1}^{K} r_{ik} \left( \beta_{0k} + \beta_k^T x_i \right) - \log \left\{ \sum_{k=1}^{K} e^{\beta_{0k} + \beta_k^T x_i} \right\} \right]$$

- The weights  $w_i$  are used to adjust the contribution of each observation to the likelihood,  $w_i = 1$  by default.
- Then for any candidate solution  $\{\tilde{\beta}_{kj}\}_{k=1}^K$ , the criterion to resolve the choice of  $c_i$  is the penalty term.

$$c_j = \operatorname*{arg\,min}_{c \in \mathbb{R}} \left\{ \sum_{k=1}^K \left| ilde{eta}_{kj} - c 
ight| 
ight\} = \operatorname{median} \left\{ ilde{eta}_{1j}, \ldots, ilde{eta}_{Kj} 
ight\},$$

for 
$$j = 1, ..., p$$
.

# Case study: Handwritten digits

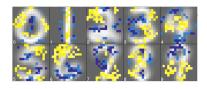
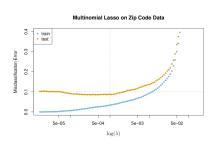


Figure: Source: Trevor Hastie, Robert Tibshirani, and Martin Wainwright. Statistical learning with sparsity: the Lasso and generalizations. CRC Press, 2015, page 38.

- Handwritten Digit Recognition with a Back-Propagation Network is published in 1989.
- LeNet-5: The first part includes two convolutional layers and two pooling layers which are placed alternatively. The second part consists of three fully connected layers.

# Case study: Handwritten digits

- We have N=7921 gray-scale images of p=256 pixels representing handwritten digits from 0 yo 9, namely,  $Y \in \{0,...,9\}$ .
- Each one of the p features represents the intensity in a [0,1]-scale of the corresponding pixel (0 black, 1 white).
- We can fit a 10-classes lasso multinomial model.
- Deep networks indeed have better performance.



# Computational techniques

Linear predictor:

$$Z_{ik} = \beta_{0k} + \beta_k^{\top} x_i, \quad p_{ik} = \Pr(Y = k \mid x_i) = \frac{e^{Z_{ik}}}{\sum_{m=1}^{K} e^{Z_{im}}}.$$

Log-likelihood (without penalty):

$$\mathcal{L}(\beta) = \sum_{i=1}^{N} \sum_{k=1}^{K} r_{ik} Z_{ik} - \sum_{i=1}^{N} \log \left( \sum_{m=1}^{K} e^{Z_{im}} \right).$$

Lasso problem:

$$\min_{\beta} -\frac{1}{N}\mathcal{L}(\beta) + \lambda \sum_{k=1}^{K} \|\beta_k\|_{1}.$$

• We use coordinate descent (BCD) algorithm to solve it (profiling)

# Computational techniques

- Fix other parameters  $\left\{ \tilde{\beta}_{0k'}, \tilde{\beta}_{k'} \right\}_{k' \neq k}$ , update  $(\beta_{0k}, \beta_k)$
- Let  $\ell_{ik} = y_{ik}Z_{ik} \log\left(\sum_{j=1}^{K} e^{Z_{ij}}\right)$  and  $p_k\left(x_i\right) = \frac{e^{Z_{ik}}}{\sum_{i=1}^{K} e^{Z_{ij}}}$ 
  - $\bullet \ \frac{\partial \ell_{ik}}{\partial Z_{ik}} = y_{ik} p_k(x_i)$
  - $\bullet \frac{\partial^2 \ell_{ik}}{\partial Z_i^2} = -p_k(x_i)(1-p_k(x_i)) =: -w_{ik}$
- Heuristically, by Taylor expansion, we can derive a quadratic objective function as

$$Q_{k}\left(\beta_{0k},\beta_{k}\right) = -\frac{1}{2N}\sum_{i=1}^{N}w_{ik}\left(h_{ik} - \beta_{0k} - \beta_{k}^{T}x_{i}\right)^{2} + C\left(\left\{\tilde{\beta}_{0k},\tilde{\beta}_{k}\right\}_{k=1}^{K}\right),$$

where C denotes a constant independent of  $(\beta_{0k}, \beta_k)$ , and

$$h_{ik} = \tilde{eta}_{0k} + \tilde{eta}_k^T x_i + rac{y_{ik} - ilde{p}_k\left(x_i
ight)}{ ilde{p}_k\left(x_i
ight)\left(1 - ilde{p}_k\left(x_i
ight)\right)} ext{ and } w_{ik} = ilde{p}_k\left(x_i
ight)\left(1 - ilde{p}_k\left(x_i
ight)\right)$$

#### Count outcomes

• A random variable Y is  $Poisson(\lambda)$  if its probability mass function is

$$P(Y = k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad (k = 0, 1, 2, ...)$$

- If  $Y \sim \mathsf{Poisson}(\lambda)$ , then  $\mathbb{E}(Y) = \mathsf{Var}(Y) = \lambda$ .
- If  $Y_1, \ldots, Y_K$  are mutually independent with  $Y_k \sim \text{Poisson}(\lambda_k)$

$$Y_1+\cdots+Y_K\sim \mathsf{Poisson}(\lambda),$$
 
$$(Y_1,\ldots,Y_K)\mid Y_1+\cdots+Y_K=n\sim \mathsf{Multinomial}\left(n,\left(\frac{\lambda_1}{\lambda},\ldots,\frac{\lambda_K}{\lambda}\right)\right),$$
 where  $\lambda=\lambda_1+\cdots+\lambda_K$ .

#### Some extensions of Poisson distribution

- The Poisson distribution restricts that the mean must be the same as the variance. It cannot capture the feature of over-dispersed data with variance larger than the mean.
- Negative binomial distribution: a scale-mixture of Poisson.

$$P(Y=k) = \frac{\Gamma(k+\theta)}{\Gamma(k+1)\Gamma(\theta)} \left(\frac{\theta}{\mu+\theta}\right)^{\theta} \left(\frac{\mu}{\mu+\theta}\right)^{k}, \quad (k=0,1,2,\ldots)$$

with  $\mathbb{E}(Y) = \mu$  and  $Var(Y) = \mu + \mu^2/\theta$ .

Zero-inflated Poisson: a mixture of Poisson and point mass at zero.

$$\begin{split} P(Y=k) &= \begin{cases} p + (1-p)e^{-\lambda}, & \text{if } k = 0\\ (1-p)e^{-\lambda}\frac{\lambda^k}{k!}, & \text{if } k = 1, 2, \dots \end{cases} \\ \mathbb{E}(Y) &= (1-p)\lambda \text{ and } \text{Var}(Y) = (1-p)\lambda(1+p\lambda) \end{split}$$

# Poisson regression model / log-linear model

•

$$\begin{cases} Y_i \mid X_i & \sim \mathsf{Poisson}\left(\lambda_i\right), \\ \lambda_i & = \lambda\left(X_i, \beta\right) = e^{\beta_0 + X_i^\top \beta}. \end{cases}$$

$$\mathbb{E}\left(Y_i \mid X_i\right) = \mathsf{var}\left(Y_i \mid X_i\right) = e^{\beta_0 + X_i^\top \beta}.$$

This model is also called log-linear model

$$\log \mathbb{E}(Y_i \mid X_i) = \beta_0 + X_i^{\top} \beta$$

Interpretation: conditional log mean ratio

$$\log \frac{\mathbb{E}(Y_i \mid \dots, X_{ij} + 1, \dots)}{\mathbb{E}(Y_i \mid \dots, X_{ij}, \dots)} = \beta_j$$

# Poisson regression with penalty

Likelihood:

$$L(\beta) = \prod_{i=1}^{n} e^{-\lambda_i} \frac{\lambda_i^{y_i}}{y_i!} \propto \prod_{i=1}^{n} e^{-\lambda_i} \lambda_i^{y_i}$$
$$\log L(\beta) = \sum_{i=1}^{n} \left( -\lambda_i + y_i \log \lambda_i \right) = \sum_{i=1}^{n} \left( -e^{\beta_0 + x_i^{\top} \beta} + y_i (\beta_0 + x_i^{\top} \beta) \right).$$

• The penalty term is  $\lambda \|\beta\|_1$ .

$$-\frac{1}{N}\sum_{i=1}^{N}\left\{y_{i}\left(\beta_{0}+\boldsymbol{\beta}^{\top}\boldsymbol{x}_{i}\right)-e^{\beta_{0}+\boldsymbol{\beta}^{\top}\boldsymbol{x}_{i}}\right\}+\lambda\|\boldsymbol{\beta}\|_{1}$$

• Take derivatives on  $\beta_0$  and set it to zero  $\leadsto \frac{1}{N} \sum_{i=1}^N e^{\hat{\beta}_0 + \hat{\beta}^\top x_i} = \bar{y}$ 

# Poisson regression for rates modeling

• If observation windows have different lengths  $T_i$ , then

$$\mathbb{E}\left[y_{i}\mid X_{i}=x_{i}\right]=T_{i}\mu\left(x_{i}\right)$$

where  $\mu(x_i)$  rate per unit time interval.

• 6 months  $\sim$  yearly visit to doctor has T=6.

•

$$\log \mathbb{E}[Y \mid X = x, T] = \underbrace{\log T}_{\text{"offset"}} + \beta_0 + \beta^{\top} x$$

• The terms  $\log T_i$  for each observation require no fitting, and are called an offset.

# Case study: distribution smoothing

- *N* count variables  $\{y_k\}_{k=1}^N$  coming from a *N*-cell multinomial distribution.
- $\mathbf{r} = \{r_k\}_{k=1}^N = \{y_k / \sum_{k=1}^N y_k\}_{k=1}^N$  vector of proportions.
- Issue: r could be sparse. Want to regularize it toward a more stable distribution  $u = \{u_k\}_{k=1}^N$ .

•

 We want a distribution q which is approximately equal to our observed proportions but at the same time as close as possible to a nominal distribution u.

# Case study: distribution smoothing

Why this problems falls in Poisson model framework?

• The previous minimization problem minimize  $\sum_{q\in\mathbb{R}^N,q_k\geq 0}^N q_k \log\left(\frac{q_k}{u_k}\right)$  such that  $\|q-r\|_\infty \leq \delta, \sum_{k=1}^N q_k = 1$  has Lagrange dual

$$\underset{\beta_0,\alpha}{\operatorname{maximize}} \left\{ \sum_{k=1}^{N} r_k \left[ \log u_k + \beta_0 + \alpha_k - u_k e^{\beta_0 + \alpha_k} \right] - \delta \|\alpha\|_1 \right\}$$

• This is equivalent to fitting a Poisson model with offset  $\log u_k$ , individual parameter  $\alpha_k$  and design matrix  $X = \mathbb{I}_{N \times N}$ .

## Case study: distribution smoothing

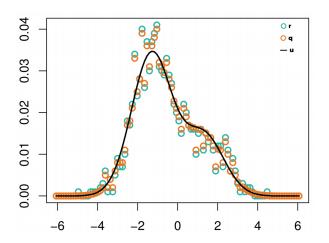


Figure: Source: Trevor Hastie, Robert Tibshirani, and Martin Wainwright. Statistical learning with sparsity

#### Time-to-event data

- Survival analysis in biostatistics
  - Outcome denotes the Survival time or the time to the recurrent of the disease
- Duration analysis in econometrics
  - Outcome denotes the weeks unemployed or days until the next arrest after being released from incarceration
- Time-to-event-data
  - Non-negative
  - May be censored, resulting in inadequate tail information

## Survival function

- Medical studies interested in time to death T of sick patients, usually characterized by the survival function  $S(t) := \mathbb{P}(T > t)$ , the probability of surviving beyond a certain time t.
- Some patients drop out the study or die because of unrelated causes: we call this situation a censoring time *C*.
- $Y := \min(C, T)$  is the observed outcome variable, together with an indicator  $\delta := \mathbb{I}\{Y = T\}$  of whether the patient died correctly (because of the studied illness).

## Hazard function and Cox model

• **Hazard function**: Instantaneous probability of death at time *t*, given survival up till *t*.

$$h(t) = \lim_{\delta \to 0} \frac{\mathbb{P}(Y \in \{t, t + \delta\} \mid Y \ge t)}{\delta} = \frac{f(t)}{S(t)}$$

where f(t) density of T.

• Cox's model treats special cases of hazard functions:

$$h(t|x) = h_0(t)e^{\beta^{\top}x}$$

where x represents e.g. gene expressions and  $h_0(t)$  is **baseline** hazard: hazard for one individual with  $x = 0 \rightsquigarrow$  semiparametric

ullet The coefficient eta represents the multiplicative effect of the covariates on the hazard

$$rac{h(t|\cdots,x_j+1,\cdots)}{h(t|\cdots,x_j,\cdots)}=e^{eta_j} \leadsto \mathsf{hazard} \; \mathsf{ratio} \; (\mathsf{HR})$$

# The hazard of hazard ratio (Hernán, 2010)

- Treatment  $Z_i$  and time-to-event outcome  $Y_i$ .
- Hazard: the event rate at time t conditional on survival until time t or later

$$\lim_{\Delta t \to 0} \frac{P\big(t \leq Y < t + \Delta t | Y \geq t\big)}{\Delta t}.$$

- Survival analysis assumes models for hazard, e.g., Cox models, additive hazard models, etc.
- Hazard ratio between the treatment and control compares

$$\lim_{\Delta t \to 0} \frac{P(t \le Y(1) < t + \Delta t | Y(1) \ge t)}{\Delta t}, \quad \lim_{\Delta t \to 0} \frac{P(t \le Y(0) < t + \Delta t | Y(0) \ge t)}{\Delta t}$$

which compares the populations  $\{i: Y_i(1) \ge t\}$  and  $\{i: Y_i(0) \ge t\}$ .

Hazard ratio has a built-in selection bias.

## Risk sets and partial likelihood

• Go back to the Cox model.

$$h(t|x) = h_0(t)e^{\beta^{\top}x}$$

- Denote by  $R_i := \{j \mid y_j \ge y_i\}$  the risk set of subject i (individuals which are still in the study when subject i dies).
- $\bullet$  Cox (1972; 1975) proposed the partial likelihood. He argued that only a part of the likelihood is useful when we are interested in the parameter  $\beta$  only
- The partial likelihood of subject *i* is given by

$$\frac{h(y_i|x_i)}{\sum_{j\in R_i}h(y_j|x_j)} = \frac{e^{\beta^\top x_i}}{\sum_{j\in R_i}e^{\beta^\top x_j}}$$

Note that baseline hazard  $h_0$  has no effect here, does not depend on actual death times

## Interpretation of Partial Likelihood

- Let the "event" = "die" for convenience
- The probability that a subject with covariate value  $X_{(j)}$  dies at time  $t_j$ , given that one and only one subject in the risk set  $R_j$  dies at time  $t_j$ , is

$$P(\text{subject with } X_{(j)} \text{ dies at time } t_j | \text{dies in } R_j \text{ but don't know who})$$

$$= \frac{P(\text{subject with } X_{(j)} \text{ dies at time } t_j | \text{ survives to time } t_j)}{P(\text{dies at time } t_j \text{ and know in the risk set } R_j | \text{ in the set } R_j \text{ survive to time } t_j)}$$

$$= \frac{h(t_j | X_{(j)})}{\sum_{i \in R_j} h(t_j | X_i)}$$

$$= \frac{e^{\beta^\top X_{(j)}}}{\sum_{i \in R_j} e^{\beta^\top X_i}}.$$

• We sequentially consider in time order  $t_1, t_2, ..., t_k$  and then we obtain the likelihood  $L_p(\beta)$ 

## Cox model with penalty

The log-partial-Likelihood is

$$\ell_{m{p}}(m{eta};m{x},m{\delta}) = \sum_{\substack{\delta_i = 1 \ ext{died "correctly"}}} \log \left[ rac{e^{m{eta}^ op \mathbf{x}_i}}{\sum_{j \in R_i} e^{m{eta}^ op \mathbf{x}_j}} 
ight]$$

with corresponding  $\ell_1$ -penalized CPH problem:

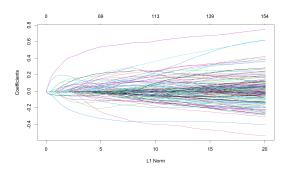
$$\underset{\beta}{\mathsf{minimize}} \left\{ -\sum_{\delta_i = 1} \log \left[ \frac{e^{\beta^\top x_i}}{\sum_{j \in R_i} e^{\beta^\top x_j}} \right] + \lambda \|\beta\|_1 \right\}$$

## Case study: lymphoma data

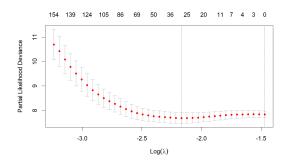
• We want to estimate the hazard function S for N=240 Lymphoma patients with p=7399 variables measuring gene expressions. 102 of these samples are right censored, i.e.

$$Y = \min(T, C) = C$$

• They use the  $\ell_1$ -penalized CPH problem to find  $\hat{\beta}(\lambda_{\min})$ .



## Case study: lymphoma data



- Simon, Friedman, Hastie and Tibshirani (2011) give details for an algorithm based on coordinate-descent
- Implemented in the R package glmnet

## Kaplan-Meier estimator

We use the Kaplan-Meier estimator of survivor function S(t): let  $\hat{\eta}(x) := \hat{\beta} \left( \lambda_{\min} \right)^{\top} x$ , then

$$\widehat{S}(t) = \prod_{i: y_i \leq t} \left( 1 - rac{e^{\widehat{\eta}(x_i)}}{\sum_{j \in R_i} e^{\widehat{\eta}(x_j)}} 
ight)$$

is an estimate of S(t). We use these in the following plot.

In computing the score  $\hat{\eta}(x_i)^{(k)}$  for the observations in fold k, we use the coefficient vector  $\hat{\beta}^{(-k)}$  computed with those observations omitted. Doing this for all K folds, we obtain the "pre-validated" dataset

$$\left\{\left(\hat{\eta}\left(x_{i}\right)^{(k)},y_{i},\delta_{i}\right)\right\}_{i=1}^{N}.$$

### Pre-validation

- The Cox model is a semi-parametric model, and the proportional hazard assumption is crucial.
- The proportional hazard assumption is not always satisfied, and the Cox model may not be the best choice.
- Pre-validation is a method to check the proportional hazard assumption.
- The idea is to split the data into two parts,
  - one part  $\rightsquigarrow x_i$
  - the other part  $\rightsquigarrow \hat{\beta}^{(k)}$
  - repeat for k = 1, ..., K folds
  - obtain the pre-validated dataset  $\{(\hat{\eta}(x_i)^{(k_i)}, y_i, \delta_i)\}_{i=1}^N$
- If the proportional hazard assumption holds, the estimated coefficients should be similar.

## Case study: lymphoma data

 Log-rank test is used to formally test whether the survival curves are statistically different

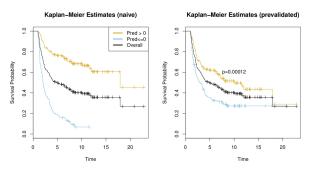
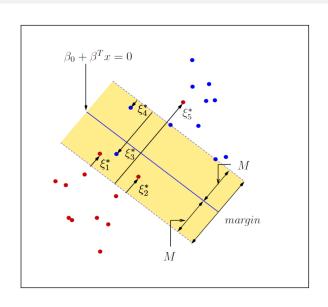


Figure 3.7 The black curves are the Kaplan-Meier estimates of S(t) for the Lymphoma data. In the left plot, we segment the data based on the predictions from the Cox proportional hazards lasso model, selected by cross-validation. Although the tuning parameter is chosen by cross-validation, the predictions are based on the full training set, and are overly optimistic. The right panel uses prevalidation to build a prediction on the entire dataset, with this training-set bias removed. Although the separation is not as strong, it is still significant. The spikes indicate censoring times. The p-value in the right panel comes from the log-rank test.

# Support vector machine



# SVM: A toy model

- Suppose we have a classification rule :  $\{x: f(x) = \beta_0 + x^T \beta\}$ , here  $x = (x_1, x_2)^T \in \mathbb{R}^2$ .
- Consider the following lines:

$$I_1: \beta_0 + \beta_1 x_1 + \beta_2 x_2 = 1,$$
  

$$I_{-1}: \beta_0 + \beta_1 x_1 + \beta_2 x_2 = -1$$
  

$$I_0: \beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0.$$

Then, by geometry,

$$d(I_0,I_1)=d(I_0,I_{-1})=rac{1}{\sqrt{eta_1^2+eta_2^2}}=rac{1}{\|eta\|_2}.$$

• So if the dataset is linear separable, by setting the closest point  $x_i$  to the classification boundary  $I_0$  as  $f(x_i) = 0$ , we have  $\exists \beta_0, \beta$  such that

$$f(x_i) = \begin{cases} > 0, & \text{if } y_i = +1, \\ < 0, & \text{if } y_i = -1. \end{cases}$$

## SVM geometric interpretation

Consider the Boundary  $B = \{x \in \mathbb{R}^p \mid f(x) = 0\}$ , where

$$f(x) = \beta_0 + \beta^\top x$$

Then the distance between the boundary and the point  $x_0$  is

$$\operatorname{dist}(x_0, B) = \inf_{z \in B} \|z - x_0\|_2 = \frac{|f(x_0)|}{\|\beta\|_2}$$

So we find that the optimal separating plane  $f^*(x) = 0$  has margin

$$M_2^* = \max_{\beta_0, \beta} \left\{ \min_{i \in \{1, \dots, n\}} \frac{y_i f(x_i, \beta_0, \beta)}{\|\beta\|_2} \right\}$$

# SVM with slack variable $\xi_i$

- We allow some points to be misclassified, and introduce a slack variable  $\xi = (\xi_1, ..., \xi_n), \xi_i \geq 0$  for each observation.
- $y_i(\beta_0 + x_i^T \beta)$  represents the "distance" of  $x_i$  to the classification boundary  $\leadsto$  a linear classifier is scale invariant, the solution coefficients can be rescaled, WLOG, we can set  $\|\beta\|_2 = 1$
- Foundamentally, we want  $y_i(\beta_0 + x_i^T \beta) \geq M$ ,
- But we allow some points to be misclassified, so we relax the constraint to  $y_i(\beta_0 + x_i^T \beta) \ge M(1 \xi_i)$ .
- $\sum_{i=1}^{N} \xi_i \leq C$  will introduce a bias-variance trade-off:
  - C = 0: hard margin SVM, no misclassification allowed,
  - *C* large: wide margin, introduce large bias and low variance in classification,
  - C small: narrow margin, introduce small bias and high variance in classification.

## Optimization problem of SVM

- Denote *M* as half width of the yellow part in the illustrator as the **margin** of the classifier.
- Objective:

$$\max_{\beta_0,\beta,\{\xi_i\}_{i=1}^N} M$$

• Constraints:  $y_i \underbrace{\left(\beta_0 + \beta^\top x_i\right)}_{f(x_i, \beta_0, \beta)} \ge M(1 - \xi_i), \forall i$ , where  $\xi_i \ge 0, \forall i, \sum_{i=1}^N \xi_i \le C, \|\beta\|_2 = 1$ 

Note that by the constraints, we have

$$\xi_i \geq 1 - y_i f(x_i, \beta_0, \beta), \quad \xi_i \geq 0.$$

so 
$$\sum_{i=1}^{N} \xi_i \ge \sum_{i=1}^{N} [1 - y_i(\beta_0 + \beta^T x_i)]_+$$
.

## Optimization problem of SVM

 By writing the Lagrangian equivalent of the original minimization problem (SVM), we get:

minimize 
$$\left\{ \frac{1}{N} \sum_{i=1}^{N} [1 - y_i f(x; \beta_0, \beta)]_+ + \lambda \|\beta\|_2^2 \right\}$$

Decreasing  $\lambda$  corresponds to decreasing C and  $f(x_i; \beta_0, \beta) = \beta_0 + \beta^T x_i$ .

- Hinge loss  $[1 y_i f(x; \beta_0, \beta)]_+$  is zero when if  $x_i$  lies on the correct side of the margin. For data on the wrong side of the margin, the function's value is proportional to the distance from the margin.
- $\ell_1$  penalized version:

$$\underset{\beta_0,\beta}{\operatorname{minimize}} \left\{ \frac{1}{N} \sum_{i=1}^{N} \left[ 1 - y_i f\left(x_i; \beta_0, \beta\right) \right]_+ + \lambda \|\beta\|_1 \right\}.$$

### RKHS and kernel trick

#### Kernel

Let  $\mathcal{X}$  be a non-empty set. A function  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  is called a kernel if there exists an  $\mathbb{R}$ -Hilbert space and a map  $\phi: \mathcal{X} \to \mathcal{H}$  such that  $\forall x, x' \in \mathcal{X}$ .

$$k(x, x') := \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}.$$

#### **RKHS**

Let  $\mathcal H$  be a Hilbert space of  $\mathbb R$ -valued functions defined on a non-empty set  $\mathcal X$ . A function  $k:\mathcal X\times\mathcal X\to\mathbb R$  is called a reproducing kernel of  $\mathcal H$ , and  $\mathcal H$  is a reproducing kernel Hilbert space, if k satisfies

- $\forall x \in \mathcal{X}, \quad k(\cdot, x) \in \mathcal{H}$ ,
- $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}, \langle f, k(\cdot, x) \rangle_{\mathcal{H}} = f(x)$  (the reproducing property).

## Examples of kernels

- If  $K(x,y) = (\langle x,y \rangle + c)^2 = (x_1y_1 + x_2y_2 + c)^2 = \langle \Phi(x), \Phi(y) \rangle$ , then  $\Phi(x) = (x_1, x_2, \sqrt{2c}x_1, \sqrt{2c}x_2, c)^T$ .
- There are many other kernels,
  - Polynomial kernel:  $K(x,y) = (\langle x,y \rangle + c)^d$ ,
  - Gaussian kernel:  $K(x, y) = \exp(-\|x y\|^2/2\sigma^2)$ ,
  - Laplacian kernel:  $K(x, y) = \exp(-\|x y\|/\sigma)$ ,
  - Sigmoid kernel:  $K(x, y) = \tanh(\kappa \langle x, y \rangle + \theta)$ .
- The linear SVM can be generalized using a kernel to create nonlinear boundaries.
- $\|\beta\|_2 \rightsquigarrow \|\beta\|_{\mathcal{H}_K}$
- We'll revisit this variation in Group Lasso

## SVM vs logistic regression

Penalized logistic:

$$\underset{\beta_0,\beta}{\text{minimize}} \left\{ \frac{1}{N} \sum_{i=1}^{N} \underbrace{\log \left( 1 + e^{-y_i f(x_i, \beta_0, \beta)} \right)}_{\text{logistic loss}} + \lambda \|\beta\| \right\}$$

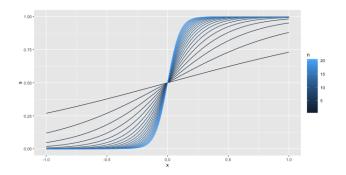
Penalized svm:

$$\underset{\beta_{0},\beta}{\operatorname{minimize}} \left\{ \frac{1}{N} \sum_{i=1}^{N} \underbrace{\left[1 - y_{i} f\left(x; \beta_{0}, \beta\right)\right]_{+}}_{\operatorname{hinge loss}} + \lambda \|\beta\| \right\}$$

 Data is separable: there exists a hyperplane that separates the two cases. In this cases logistic regression has a problem:

$$P(Y = 1 \mid X = x) = \frac{e^{\beta_0 + \beta^\top x}}{1 + e^{\beta_0 + \beta^\top x}}$$

## Problem of logistic regression



Problem: When  $p \gg N$ , the points are almost always separable.

## Relationship between SVM and logistic regression

Consider the problem

$$\underset{\beta_0,\beta}{\operatorname{minimize}} \left\{ \frac{1}{N} \sum_{i=1}^{N} \log \left( 1 + e^{-y_i f(x_i,\beta_0,\beta)} \right) + \lambda \|\beta\|_2^2 \right\}$$

Let  $\left(\tilde{\beta}_0(\lambda), \tilde{\beta}(\lambda)\right)$  be the solution, then Rosset et al. (2004) showed that

$$M_2^* = \lim_{\lambda \to 0} \left\{ \min_{i \in \{1, \dots, N\}} \frac{y_i f\left(x_i, \tilde{\beta}_0(\lambda), \tilde{\beta}(\lambda)\right)}{\|\tilde{\beta}(\lambda)\|_2} \right\}$$

So for  $\lambda \to 0$  we have that the  $\ell_2$ -regularized logistic regression corresponds to the SVM solution.

## Relationship between SVM and logistic regression

In particular, if  $\left( \breve{\beta}_0, \breve{\beta} \right)$  solve the SVM problem for C=0, then we have that:

$$\lim_{\lambda \to 0} \frac{\tilde{\beta}(\lambda)}{\|\tilde{\beta}(\lambda)\|_2} = \tilde{\beta}$$

Note that the division by the  $\ell_2$  norm of  $\tilde{\beta}(\lambda)$  makes sure that the solution on the SVM problem does not blow up.

- As  $\lambda \to 0$ , logistic regression and SVM solutions coincide
- SVM leads to a more stable numerical method for computing the solution in the solution is most dense
- Logistic regression is more useful in the sparser part of the solution path

## Part II

- Elastic net
- Group lasso/overlap group lasso
- Fused lasso





# Case study: comparison of lasso and elastic net on highly correlated variables

In microarray studies, **groups of genes in the same biological pathway** tend to be expressed (or not) together, and hence measures of their expression tend to be strongly correlated.

#### Simulation setup

- 1 2 sets of 3 variables, pairwise correlations around 0.97 in each group
- 2 sample size: N = 100,
- data are simulated as follows:

$$Z_1, Z_2 \sim N(0,1)$$
 independent  $Y = 3Z_1 - 1.5Z_2 + 2\epsilon$ ,  $\epsilon \sim N(0,1)$   $X_j(j=1,2,3) = Z_1 + \xi/5$ ,  $\xi_j \sim N(0,1)$   $X_i(j=4,5,6) = Z_2 + \xi/5$ ,  $\xi_i \sim N(0,1)$ 

# Case study: comparison of lasso and elastic net on highly correlated variables

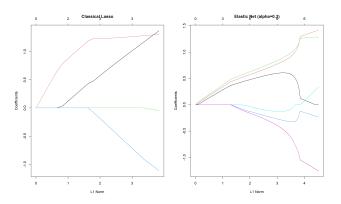
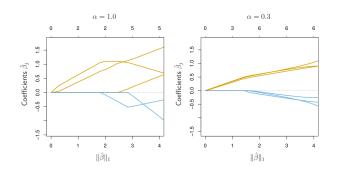


Figure: The lasso estimates (  $\alpha=1$  ), as shown in the left panel, exhibit somewhat erratic behavior as the regularization parameter  $\lambda$  is varied. In the right panel, the elastic net with ( $\alpha=0.3$ ) includes all the variables, and the correlated groups are pulled together.

# Case study: comparison of lasso and elastic net on highly correlated variables



- lasso estimates exhibit erratic behavior as  $\lambda$  varies: one variable is excluded and the correlations among variables are not clear
- elastic net includes all variables and correlated groups are pulled together, sharing values approximately equally.
- the difference between my plot and the book's plot is due to the random seed.

65 / 126

### Elastic net

Recap, the elastic net problem is defined as

$$\underset{(\beta_0,\beta)\in\mathbb{R}\times\mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^{N} \left( y_i - \beta_0 - x_i^T \beta \right)^2 + \lambda \left[ \frac{1}{2} (1-\alpha) \|\beta\|_2^2 + \alpha \|\beta\|_1 \right] \right\}$$

• Denote R as the objective function of elastic net, then for  $\beta_j > 0$ ,

$$\frac{\partial R}{\partial \beta_{j}} = \frac{1}{N} \sum_{i=1}^{N} \left( y_{i} - \beta_{0} - \boldsymbol{x}_{j}^{\top} \boldsymbol{\beta} \right) (-x_{ij}) + \lambda \left[ (1 - \alpha)\beta_{j} + \alpha \operatorname{sgn}(\beta_{j}) \right]$$

$$= \frac{1}{N} \sum_{i=1}^{N} (y_{i} - \beta_{0} - \sum_{k \neq j} x_{ik} \beta_{k} - x_{ij} \beta_{j}) (-x_{ij}) + \dots$$

$$= -\frac{1}{N} \sum_{i=1}^{N} r_{ij} x_{ij} + \frac{1}{N} \sum_{i=1}^{N} x_{ij}^{2} \beta_{j} + \lambda \left[ (1 - \alpha)\beta_{j} + \alpha \operatorname{sgn}(\beta_{j}) \right].$$

## Elastic net

Setting the derivatives to zero yields

$$\left[\frac{1}{N}\sum_{i=1}^{N}x_{ij}^{2} + \lambda(1-\alpha)\right]\beta_{j} + \alpha\lambda\operatorname{sgn}(\beta) = \frac{1}{N}\sum_{i=1}^{N}r_{ij}x_{ij}$$

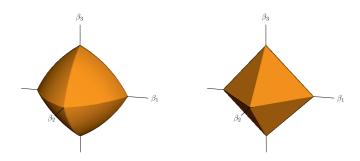
As we did in lasso case, here coordinate descent have a close form update for each  $\beta_i$ .

$$\widehat{\beta}_{j} = \frac{\mathcal{S}_{\lambda\alpha} \left( \sum_{i=1}^{N} r_{ij} x_{ij} \right)}{\sum_{i=1}^{N} x_{ij}^{2} + \lambda (1 - \alpha)},$$

where  $r_{ij} = y_i - \tilde{\beta}_0 - \sum_{k \neq j} x_{ik} \hat{\beta}_k$  and  $S_{\mu}(z) := \text{sign}(z)(z - \mu)_+$ .

- In practice, group structure may not be as evident as the previous 'ideal' model, this example does capture the main idea of elastic net.
- By adding ridge penalty to lasso penalty, elastic net automatically controls for strong within-group correlations.

## Elastic net



**Figure 4.2** The elastic-net ball with  $\alpha=0.7$  (left panel) in  $\mathbb{R}^3$ , compared to the  $\ell_1$  ball (right panel). The curved contours encourage strongly correlated variables to share coefficients (see Exercise 4.2 for details) the following pages for details

# Why does elastic net promote grouping?

## Grouping effect (Zou and Hastie, 2005)

Given data  $(\boldsymbol{y},\boldsymbol{X})$  and penalty parameters  $(\lambda_1,\lambda_2)$ , the response  $\boldsymbol{y}$  is centered and the predictor  $\boldsymbol{X}$  are standardized. Let  $\hat{\beta}(\lambda_1,\lambda_2)$  be the naive elastic net estimate. Suppose that  $\hat{\beta}_i(\lambda_1,\lambda_2)\hat{\beta}_j(\lambda_1,\lambda_2)>0$ . Then

$$D_{\lambda_1,\lambda_2}(i,j) = \frac{1}{\|\mathbf{y}\|_2} \left| \hat{\beta}_i\left(\lambda_1,\lambda_2\right) - \hat{\beta}_j\left(\lambda_1,\lambda_2\right) \right| \leq \frac{1}{\frac{\lambda_2}{2}} \sqrt{2(1-\rho)}$$

where  $\rho = \mathbf{x}_i^T \mathbf{x}_i$ , the sample correlation.

The unitless quantity  $D_{\lambda_1,\lambda_2}(i,j)$  describes the difference between the coefficient paths of predictors i and j. If  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are highly correlated, i.e.  $\rho \approx 1$  (if  $\rho \approx -1$  then consider  $-\mathbf{x}_j$ ), grouping effect says that the difference between the coefficient paths of predictor i and predictor j is almost 0

## Grouping effect

Consider the overall minimization

$$\min_{\beta} \underbrace{\| \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \|^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|^2}_{L(\boldsymbol{\beta})}.$$

Then the optimality at  $\beta_i$  and  $\beta_j$  is

$$\begin{cases} \frac{\partial L}{\partial \beta_{i}} = -2\boldsymbol{x}_{i}^{T} \left(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\right) + 2\lambda_{2}\beta_{i} + \lambda_{1}\operatorname{sgn}\left(\beta_{i}\right) = 0\\ \frac{\partial L}{\partial \beta_{j}} = -2\boldsymbol{x}_{j}^{T} \left(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\right) + 2\lambda_{2}\beta_{j} + \lambda_{1}\operatorname{sgn}\left(\beta_{j}\right) = 0 \end{cases}$$

Substracting the two equations, we have

$$2\left(\boldsymbol{x_{j}}-\boldsymbol{x_{i}}\right)^{\top}\left(\boldsymbol{y}-\boldsymbol{X}\boldsymbol{\beta}\right)+2\lambda_{2}\left(\beta_{i}-\beta_{j}\right)+\lambda_{1}\left(\underbrace{\operatorname{sgn}\left(\beta_{i}\right)-\operatorname{sgn}\left(\beta_{j}\right)}_{=0,\text{ by assumption}}\right)=0$$

# Grouping effect (cont'd)

By the previous equation, we have

$$(\beta_i - \beta_j) = \frac{1}{\lambda_2} (\mathbf{x}_i - \mathbf{x}_j)^{\top} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

Then, we have

$$(\beta_i - \beta_j)^2 \le \frac{1}{\lambda_2^2} \|\mathbf{x_i} - \mathbf{x_j}\|^2 \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$$
 (by Cauchy)

Then we can bound the inequality by parts. For  $\|\mathbf{x_i} - \mathbf{x_j}\|^2$ , by centered dataset, we have  $\|\mathbf{x_i}\|^2 = 1, i = 1, ..., p$  and  $\mathbf{x_i}^T \mathbf{x_j} = \rho$ , so

$$\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2} = \|\mathbf{x}_{i}\|^{2} + \|\mathbf{x}_{j}\|^{2} - 2\mathbf{x}_{i}^{\mathsf{T}}\mathbf{x}_{j} = 2(1 - \rho).$$

## Grouping effect (cont'd)

For  $\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2$ , by optimization,

$$\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 + \lambda_1 \|\hat{\boldsymbol{\beta}}\|_1 + \lambda_2 \|\hat{\boldsymbol{\beta}}\|_2^2 = L(\hat{\boldsymbol{\beta}}) \le L(0) = \|\mathbf{y}\|_2^2.$$

then,

$$\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 \le \|\mathbf{y}\|_2^2 - \lambda_1 \|\hat{\boldsymbol{\beta}}\|_1 - \lambda_2 \|\hat{\boldsymbol{\beta}}\|_2^2 \le \|\mathbf{y}\|_2^2.$$

Combining all the upper bounds, we have

$$|\beta_i - \beta_j| \leqslant \frac{1}{\lambda_2} \sqrt{2(1-p)} ||y||_2,$$

which completes the proof.

### Quadratic Form

Assume the design matrix has two centered and standardized columns, with inner products given by

$$\mathbf{x}_1^{\top} \mathbf{x}_1 = \mathbf{x}_2^{\top} \mathbf{x}_2 = 1, \quad \mathbf{x}_1^{\top} \mathbf{x}_2 = \rho, \quad -1 \le \rho \le 1.$$

For the coefficient vector  $\beta = (\beta_1, \beta_2)^{\mathsf{T}}$ , the squared error can be written as a quadratic form

$$L(\beta) = \frac{1}{2} \|y - X\beta\|_2^2 = \frac{1}{2} \left(\beta - \hat{\beta}_{\text{OLS}}\right)^{\top} G\left(\beta - \hat{\beta}_{\text{OLS}}\right) + \text{ const },$$

where the Gram matrix is

$$G = X^{\top}X = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$$

and the contours of the loss function, centered at  $\hat{\beta}_{OLS},$  are given by the equation

$$\left(\beta - \hat{\beta}_{\rm OLS}\right)^{\top} G\left(\beta - \hat{\beta}_{\rm OLS}\right) = \text{const}.$$

## Eigenvalues

Perform eigen-decomposition of G. Solve for eigenvalues  $\mu$ :

$$\det(G-\mu I) = \det\left( egin{array}{cc} 1-\mu & 
ho \ 
ho & 1-\mu \end{array} 
ight) = (1-\mu)^2 - 
ho^2 = 0.$$

The solutions are

$$\mu_1 = 1 + \rho, \quad \mu_2 = 1 - \rho.$$

The corresponding (unnormalized) eigenvectors can be chosen as

$$v^{(1)}=egin{pmatrix}1\\1\end{pmatrix}$$
 (in the direction  $y=x$ ) ,  $v^{(2)}=egin{pmatrix}1\\-1\end{pmatrix}$   $(y=-x)$  .

## Ellipse

In the eigenvector basis, let  $u=Q^{\top}\left(\beta-\hat{\beta}_{\mathrm{OLS}}\right)$ , where Q is composed of the unit eigenvectors as columns. The quadratic form becomes

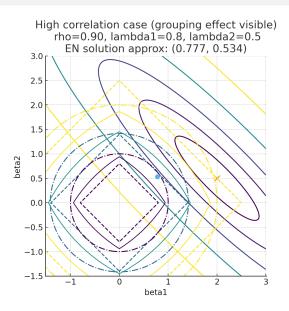
$$(\beta - \hat{\beta})^{\top} G(\beta - \hat{\beta}) = \mu_1 u_1^2 + \mu_2 u_2^2.$$

Fixing the level = c (i.e., the contour), we obtain the ellipse equation

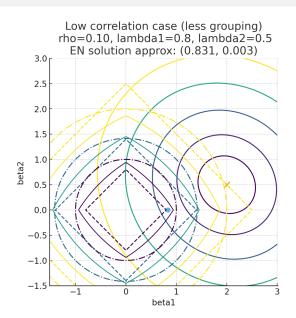
$$\frac{u_1^2}{c/\mu_1} + \frac{u_2^2}{c/\mu_2} = 1.$$

From this, the semi-axis length in the *i*-th eigenvector direction is  $\sqrt{c/\mu_i}$ . Thus, the smaller the eigenvalue, the longer the semi-axis in the corresponding direction. Noting that  $\mu_2=1-\rho\approx 0$ , the semi-axis in the y=-x direction is very long, and the loss function changes slowly in this direction. In other words, if the two variables are highly correlated, their coefficients can vary over a wide range with almost no change in the loss function value.

## Grouping effect: an illustration



## Grouping effect: an illustration



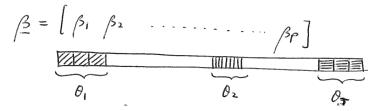
#### Group lasso

#### An example for group lasso

Genes and proteins often lie in known pathways, an investigator may be more interested in which pathway are related to an outcome than whether particular individual genes are.

- Groups of covariates should be selected into or out of a model together
- Desirable to have all coefficients within a group become nonzero (or zero) simultaneously

We use group lasso penalty (Yuan and Lin, 2006) for such situations.



78 / 126

## Group lasso

- Consider linear regression model involving J groups of covariates, where j=1,...,J
- Vector  $\mathbf{Z_j} = (X_{j1}, ..., X_{jp_j})^T \in \mathbb{R}^{p_j}$  represents the covariates in group j•  $p_j$  does not need to be the identical for all j
- Goal: predict real-valued response  $Y \in \mathbb{R}$  based on collection of covariates  $(Z_1,...,Z_J)$
- Linear model  $\mathbb{E}(Y|Z)$  takes the form  $\theta_0 + \sum_{j=1}^J \mathbf{Z}_j^T \theta_j$ , where  $\theta_j \in \mathbb{R}^{p_j}$ .

## Group lasso

Given a collection of N samples  $\{(y_i, z_{i1}, z_{i2}, \dots, z_{iJ})\}_{i=1}^{N}$  the group lasso solves the convex problem:

$$\underset{\theta_0 \in \mathbb{R}, \theta_j \in \mathbb{R}^{\rho_j}}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^{N} \left( y_i - \theta_0 - \sum_{j=1}^{J} z_{ij}^T \theta_j \right)^2 + \lambda \sum_{j=1}^{J} \left\| \theta_j \right\|_2 \right\}$$

Where  $\|\theta_i\|_2$  is the Euclidean norm.

This is group generalization of the lasso with properties:

- Depending on  $\lambda \geq 0$  either the entire vector  $\hat{\theta}_j$  will be zero, or all its elements will be nonzero.
- When  $p_j=1, j=1,...,J$ , then we have  $\|\theta_j\|_2=|\theta_j|$ , so reduces to the ordinary lasso.

### Lasso vs group lasso

If  $oldsymbol{eta}$  is individually sparse, then very likely all  $\|oldsymbol{ heta}_j\|_2$  will have value. But if  $oldsymbol{eta}$  is group sparse, then only a few groups  $\|oldsymbol{ heta}_j\|_2$  is activated.

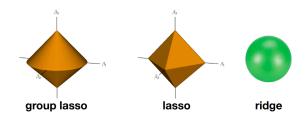
#### Lasso penalty:

- $\bullet |\beta|_1 = \sum_{j=1}^p |\beta_j|,$
- $\hat{\beta}$  is sparse.

#### Group lasso penalty:

- $\bullet \ \beta = (\theta_1, ..., \theta_J)$
- By notation abuse, denote  $\|\beta\|_{2,1} = \sum_{j=1}^{p} \|\theta\|_{2}$ , denote  $\theta = (\|\theta_1\|_1, ..., \|\theta_J\|_2)$ .
- $\bullet$   $\hat{\theta}$  is sparse.

## An unit ball example for group lasso penalty



- ullet Left unit ball can be characterized as  $\sqrt{eta_1^2+eta_2^2}+|eta_3|\leq 1$
- Middle unit ball can be characterized as  $|\beta_1| + |\beta_2| + |\beta_3| \leq 1$
- $\bullet$  Right unit ball can be characterized as  $\sqrt{\beta_1^2+\beta_2^2+\beta_3^2} \leq 1$

## Multilevel factors in regression

- A predictor can be a categorical factor with multiple levels.
- Example: one continuous predictor X and a 3-level factor  $G \in \{g_1, g_2, g_3\}$ .
- Model for the mean:

$$\mathbb{E}(Y \mid X, G) = X\beta + \sum_{k=1}^{3} \theta_{k} \mathbb{I}_{k}[G]$$

 Interpretation: regression in X with different intercepts depending on G.

## Dummy variable representation

- Introduce dummy vector  $Z = (Z_1, Z_2, Z_3)$ , where  $Z_k = \mathbb{I}_k[G]$ .
- Model becomes:

$$\mathbb{E}(Y \mid X, Z) = X\beta + Z^T \theta, \quad \theta = (\theta_1, \theta_2, \theta_3).$$

- If G has no predictive power  $\Rightarrow \theta = 0$ .
- ullet Otherwise, coefficients in  $\theta$  are typically nonzero.

## General form with multiple factors

• With multiple group variables  $G_1, \ldots, G_J$ :

$$\mathbb{E}(Y \mid X, G_1, \dots, G_J) = \beta_0 + X^T \beta + \sum_{j=1}^J Z_j^T \theta_j$$

- Variable selection often done at group level, not individual coefficients.
- **Group Lasso** is designed for this purpose.

# Aliasing and coding issues

• In unpenalized regression, dummy variables in a set sum to one

$$\mathbb{E}(Y \mid X, G_1, \dots, G_J) = \beta_0 + X^T \beta + \sum_{j=1}^J Z_j^T \theta_j$$

- This causes aliasing with the intercept (unidentifiable)
- Usual fix: use contrasts (e.g., sum-to-zero coding)
- With group lasso:
  - No aliasing concern due to  $\ell_2$  penalties
  - Penalty enforces coefficients in a group to sum to zero

## Group-lasso and factors

Problem:

$$\min_{\theta_0, \{\theta_j\}} \ \frac{1}{2} \sum_{i=1}^{N} \left( y_i - \theta_0 - \sum_{j} z_{ij}^T \theta_j \right)^2 + \lambda \sum_{j} \|\theta_j\|_2.$$

• If group j is a factor coded by dummies with  $\sum_k z_{ij,k} = 1$  for each i, then

at optimum 
$$\mathbf{1}^{\top}\theta_j = \mathbf{0}$$

• Proof idea: for any scalar c replace  $\theta_0 \mapsto \theta_0 + c$ ,  $\theta_j \mapsto \theta_j - c\mathbf{1}$ . This leaves residuals unchanged but changes the penalty. Choosing  $c = \frac{\mathbf{1}^{\top}\theta_j}{p_j}$  minimizes  $\|\theta_j - c\mathbf{1}\|_2$  and forces the group coefficients to sum to zero. So optimal  $\theta_j$  must have zero sum.

## Multitask learning with group lasso

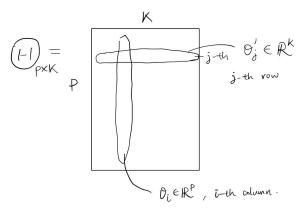
- Multivariate response  $\mathbf{Y} \in \mathbb{R}^{N \times K}$ , predictor matrix  $\mathbf{X} \in \mathbb{R}^{N \times p}$ .
- Matrix of coefficients  $\Theta \in \mathbb{R}^{p \times K}$  with, matrix of errors  $\boldsymbol{E} \in \mathbb{R}^{N \times K}$ .
- $Y \in \mathbb{R}^K$  may be correlated, for example, taking Y as K movies ratings of N users, then the ratings of different movies are correlated.
- ullet Goal: estimate ullet by minimizing the following objective function:

$$\underset{\boldsymbol{\Theta} \in \mathbb{R}^{p \times K}}{\operatorname{minimize}} \left\{ \frac{1}{2} \| \mathbf{Y} - \mathbf{X} \boldsymbol{\Theta} \|_F^2 + \lambda \left( \sum_{j=1}^p \left\| \boldsymbol{\theta}_j' \right\|_2 \right) \right\}$$

where  $||A||_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$ .  $\theta_j'$  is the *j*th row of  $\Theta$ , which means the coefficients of the *j*th task.

## Multitask learning with group lasso

- If we only use  $\sum_{j=1}^{p}\sum_{k=1}^{K}|\Theta_{jk}|$ , namely,  $\|\Theta\|_1$ , then we have no way of controlling group sparsity.
- $\sum_{j=1}^{p} \|\theta_j'\|_2$  can do group sparsity because once the  $j^{\text{th}}$  row is actuated, all elements on the row are activated.



## Computation for group lasso

We ignore the intercept and rewrite the optimization problem as follows:

$$\min_{\left(\theta_{1}, \dots, \theta_{j}\right)} \left\{ \frac{1}{2} \left\| \mathbf{y} - \sum_{j=1}^{J} \mathbf{Z}_{j} \boldsymbol{\theta}_{j} \right\|_{2}^{2} + \lambda \sum_{j=1}^{J} \left\| \boldsymbol{\theta}_{j} \right\|_{2} \right\}.$$

By taking sub-derivative on the objective function and letting the derivative to be zero, we have the following estimating equation:

$$-oldsymbol{Z_j}^T \left( oldsymbol{y} - \sum_{j=1}^J oldsymbol{Z_j} \hat{ heta}_j 
ight) + \lambda \hat{oldsymbol{s}}_j = 0,$$

$$\text{where } \hat{s}_j \in \partial \|\hat{\theta}_j\|_2 = \begin{cases} \frac{\hat{\theta}_j}{\|\theta_j\|_2}, & \text{if } \hat{\theta}_j \neq 0, \\ \text{any } \boldsymbol{v} \text{ s.t. } \|\boldsymbol{v}\|_2 \leq 1, & \text{if } \hat{\theta}_j = 0. \end{cases}$$

## Computation for group lasso

Denote  $r_j = y - \sum_{k \neq j} Z_k \hat{\theta}_k$ , then we have the estimating equation

$$-\mathbf{Z}_{\mathbf{j}}^{\mathsf{T}}\left(\mathbf{r}_{\mathbf{j}}-\mathbf{Z}_{\mathbf{j}}\hat{\theta}_{j}\right)+\lambda\hat{s}_{j}=0.$$

By coordinate descent lemma, we have

$$\hat{\theta}_j = \begin{cases} \left( \mathbf{Z}_j^T \mathbf{Z}_j + \frac{\lambda}{\|\widehat{\theta}_j\|_2} \mathbf{I} \right)^{-1} \mathbf{Z}_j^T \mathbf{r}_j, & \text{ if } \|\mathbf{Z}_j^T \mathbf{r}_j\|_2 \ge \lambda, \\ \mathbf{0}, & \text{ otherwise.} \end{cases}$$

This is not a closed form solution, one can use iterative methods to solve the equation, or add some assumptions on  $Z_j$ , for example, if  $Z_j$  is orthogonal, then the solution is closed form.

## Computation for group lasso

We have

$$\hat{\theta}_{j} = \left( \mathbf{Z}_{j}^{\mathsf{T}} \mathbf{Z}_{j} + \frac{\lambda}{\left\| \hat{\theta}_{j} \right\|_{2}} \mathbf{I} \right)^{-1} \mathbf{Z}_{j}^{\mathsf{T}} \mathbf{r}_{j} \cdot \mathbf{I} \left( \| \mathbf{Z}_{j}^{\mathsf{T}} \mathbf{r}_{j} \|_{2} \geq \lambda \right)$$

ullet Further assume that  $oldsymbol{Z_j}$  is orthonormal, then  $oldsymbol{Z_j}^Toldsymbol{Z_j}=oldsymbol{I}$ , we have

$$\widehat{ heta}_j = \left(1 - rac{\lambda}{\left\|oldsymbol{Z}_j^{oldsymbol{ au}} oldsymbol{r}_j
ight\|_2}
ight)_+ oldsymbol{Z}_j^{oldsymbol{ au}} oldsymbol{r}_j$$

# Computation for group lasso with dummy variables

- For a factor with  $p_j$  levels, the dummy matrix  $Z_j$  has columns that are level indicators.
- Then

$$Z_j^T Z_j = \operatorname{diag}(n_{j,1}, \ldots, n_{j,p_j}),$$

where  $n_{j,k}$  = number of observations in level k.

- Without normalization, groups with many observations produce larger values of  $\|Z_i^T r\|_2$  and are more easily selected.
- Normalization  $(Z_j^T Z_j = I) \rightsquigarrow$  standardize dummy columns (divide by  $\sqrt{n_{j,k}}$ )

## Sparse group lasso

- When a group is included in a group-lasso fit, all the coefficients in that group are nonzero.

#### An example for sparse group lasso

Although a biological pathway may be implicated in the progression of a particular type of cancer, not all gene in the pathway need be active.

## Sparse group lasso

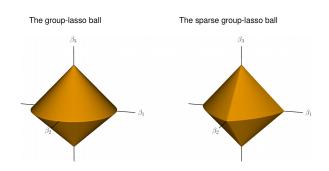
In order to achieve within-group sparsity, augment with additional  $\ell_1$ -penalty, leading to the convex program:

$$\underset{\left\{\theta_{j} \in \mathbb{R}^{p_{j}}\right\}_{j=1}^{J}}{\operatorname{minimize}} \left\{ \frac{1}{2} \left\| \mathbf{y} - \sum_{j=1}^{J} \mathbf{Z}_{j} \theta_{j} \right\|_{2}^{2} + \lambda \sum_{j=1}^{J} \left[ (1 - \alpha) \left\| \theta_{j} \right\|_{2} + \alpha \left\| \theta_{j} \right\|_{1} \right] \right\},$$

where  $\alpha \in [0, 1]$ .

- $\alpha = 0$ , reduces to the group lasso.
- $\alpha = 1$ , reduces to the lasso.

## Sparse group lasso constraint region



- Left unit ball can be characterized as  $\sqrt{\beta_1^2+\beta_2^2}+|\beta_3|\leq 1.$
- Right unit ball can be characterized as  $(1-\alpha)\sqrt{\beta_1^2+\beta_2^2}+\alpha(|\beta_1|+|\beta_2|)+(1-\alpha)|\beta_3|+\alpha|\beta_3|\leq 1.$

## Overlap group lasso

- Sometimes variables can belong to more than one group.
- Genes can belong to more than one biological pathway
- For example, we divide 5 variables into 2 groups,

$$Z_1 = (X_1, X_2, X_3), \quad Z_2 = (X_3, X_4, X_5).$$

- If we simply replicate variable  $X_3$ , then use group lasso  $\rightsquigarrow X_3$  will be selected to the model with higher probability
- If we simply replicate parameter then use sparse group lasso  $\rightsquigarrow X_3$  will be selected to the model only when both two groups are selected
- So replicate variable is preferred.

### Overlap group lasso

- $\nu_j \in \mathbb{R}^p$  is a vector which is zero everywhere except in those positions corresponding to member of the group j.
- $V_i \subseteq \mathbb{R}^p$  subspace of possible vectors.
- ullet For  $X=(X_1,\ldots,X_p)$  the coefficient vector is given by  $eta=\sum_{j=1}^J 
  u_j$
- The overlap group lasso solves the problem:

$$\operatorname{minimize}_{\nu_{j} \in \mathcal{V}_{j}, j=1, \dots, J} \left\{ \frac{1}{2} \left\| \mathbf{y} - \mathbf{X} \left( \sum_{j=1}^{J} \nu_{j} \right) \right\|_{2}^{2} + \lambda \sum_{j=1}^{J} \left\| \nu_{j} \right\|_{2} \right\}$$

• Jacob et al., (2009) showed that, the equivalent optimization problem can be put in the form:

$$\underset{\beta \in \mathbb{R}^p}{\mathsf{minimize}} \left\{ \frac{1}{2} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \|_2^2 + \lambda \Omega_{\mathcal{V}}(\boldsymbol{\beta}) \right\}, \quad \Omega_{\mathcal{V}}(\boldsymbol{\beta}) := \inf_{\substack{\nu_j \in \mathcal{V}_j \\ \boldsymbol{\beta} = \sum_{j=1}^J \nu_j}} \sum_{j=1}^J \left\| \nu_j \right\|_2$$

#### Basis functions

- Suppose for the moment that we have just a single feature X and we are interested in estimating  $\mathbb{E}(Y \mid X) = f(X)$
- A common approach for extending the linear model  $f(X) = X\beta$  is to augment X with additional, known functions of X:

$$f(X) = \sum_{m=1}^{M} \beta_m h_m(X),$$

where the  $\{h_m\}$  are called basis functions

• Because the basis functions  $\{h_m\}$  are prespecified and the model is linear in the new variables, ordinary least squares approaches can be used (at least in low-dimensional settings)

#### Basis functions

- Polynomial regression
  - Not only does this introduce bias, but it also results in extremely high variance near the edges of the range of x
  - Runge's phenomenon
- local basis functions, which ensure that a given observation affects only the nearby fit, not the fit of the entire line → splines
  - Cubic splines
  - Natural cubic splines

#### Additive models

#### Additive models

Additive models are based on approximating the regression function by sums of the form:

$$f(x) = f(x_1, \ldots, x_J) \approx \sum_{j=1}^J f_j(x_j), \quad f_j \in \mathcal{F}_j, \quad j = 1, \ldots, J$$

- $\mathcal{F}_i$  are fixes set of univariate function classes
- Each  $\mathcal{F}_j$  assumed to be a subset of  $L^2(\mathbb{P}_j)$
- $\mathbb{P}_j$  is the distribution of covariate  $X_j$  equipped with squared  $L^2(\mathbb{P}_j)$  norm

$$\|f_j\|_2^2 := \mathbb{E}\left[f_j^2(X_j)\right]$$

• Some theoretical results need  $\mathcal{F}$  to be the Sobelev class of functions on [a, b]. (Buhlmann and van de Geer, 2010)

#### Additive models

Best additive approximation to regression function  $\mathbb{E}(Y \mid X = x)$  solves problem:

$$\underset{f_{j}\in\mathcal{F}_{j}j=1,...,J}{\mathsf{minimize}} \mathbb{E}\left[\left(Y-\sum_{j=1}^{J}f_{j}\left(X_{j}\right)\right)^{2}\right],\mathcal{F}_{j}\subseteq L^{2}\left(\mathbb{P}_{j}\right),j=1,\ldots,J$$

The optimal solution  $(\tilde{f}_1, \dots, \tilde{f}_J)$  is characterized by the **backfitting** equations:

$$ilde{f_j}\left(x_j
ight) = \mathbb{E}\left[Y - \sum_{k 
eq j} ilde{f_k}\left(X_k
ight) \mid X_j = x_j
ight], ext{ for } j = 1, \dots, J$$

or 
$$ilde{f_j} = \mathcal{P}_j(Y - \sum_{k \neq j} ilde{f_k}(X_k))$$

# Partial Residuals and Backfitting for Linear Models

The general form of a linear regression model is

$$\mathbb{E}[Y \mid \vec{X} = \vec{x}] = \beta_0 + \vec{\beta} \cdot \vec{x} = \sum_{i=0}^{p} \beta_i x_i$$

Suppose we don't condition on all of  $\vec{X}$  but just one component of it, say  $X_k$ . What is the conditional expectation of Y?

$$\mathbb{E}[Y \mid X_k = x_k] = \mathbb{E}[\mathbb{E}[Y \mid X_1, X_2, \dots X_k, \dots X_p] \mid X_k = x_k]$$

$$= \mathbb{E}\left[\sum_{j=0}^p \beta_j X_j \mid X_k = x_k\right]$$

$$= \beta_k x_k + \mathbb{E}\left[\sum_{j \neq k} \beta_j X_j \mid X_k = x_k\right]$$

# Partial Residuals and Backfitting for LM/GAMs

Rearranging gives

$$\beta_k x_k = \mathbb{E} \left[ Y \mid X_k = x_k \right] - \mathbb{E} \left[ \sum_{j \neq k} \beta_j X_j \mid X_k = x_k \right]$$
$$= \mathbb{E} \left[ Y - \left( \sum_{j \neq k} \beta_j X_j \right) \mid X_k = x_k \right]$$

The expression in the expectation is the  $k^{th}$  partial residual. Let's introduce a symbol for this, say  $Y^{(k)}$ .

$$\beta_k x_k = \mathbb{E}\left[Y^{(k)} \mid X_k = x_k\right]$$

→ Gauss-Seidel type of algorithm for fitting linear models.

"a popular version of coordinate descent is known as backfitting and is used to fit generalized additive models"

# Sparse additive models (SPAM)

• For  $\lambda \ge 0$  type of k best sparse approximation:

$$\underset{f_{j} \in \mathcal{F}_{j}=1,...,J}{\operatorname{minimize}} \mathbb{E} \left[ \left( Y - \sum_{j \in \mathcal{S}} f_{j} \left( X_{j} \right) \right)^{2} \right]$$

where  $S \subset \{1,\dots,J\} \leadsto 0$ -norm constraint on the number of nonzero components

 SPAM combines ideas from sparse linear modeling and additive nonparametric regression

$$\underset{f_{j} \in \mathcal{F}_{j}=1,...,J}{\operatorname{minimize}} \left\{ \mathbb{E}\left[\left(Y - \sum_{j=1}^{J} f_{j}\left(X_{j}\right)\right)^{2}\right] + \lambda \sum_{j=1}^{J} \left\|f_{j}\right\|_{2} \right\}, \quad \left\|f\right\|_{2} = \sqrt{\mathbb{E}\left[f^{2}\left(X_{j}\right)\right]}$$

This idea was originally proposed by Ravikumar et al. (2009)

#### **COSSO**

• COSSO method uses combination of the  $\ell_1$ -norm with the Hilbert norm:

$$\|f\|_{\mathcal{H},1} := \sum_{j=1}^J \|f_j\|_{\mathcal{H}}$$

COSSO's objective function is given by

$$\min_{f_{j} \in \mathcal{H}_{j} j = 1, ..., J} \left\{ \frac{1}{N} \sum_{i=1}^{N} \left( y_{i} - \sum_{j=1}^{J} f_{j} \left( x_{ij} \right) \right)^{2} + \lambda_{\mathcal{H}} \sum_{j=1}^{J} \left\| f_{j} \right\|_{\mathcal{H}_{j}} \right\}$$

- By Pythagorean theorem, the  $j^{\text{th}}$  coordinate function  $\widehat{f_j}$  in any optimal COSSO solution can be written in the form  $\widehat{f_j}(\cdot) = \sum_{i=1}^N \widehat{\theta}_{ij} \mathcal{R}_j (\cdot, x_{ij})$ , for a suitably chosen weight vector  $\widehat{\theta}_j \in \mathbb{R}^N \rightsquigarrow \text{dimension reduction}$
- ullet Gram matrix  $oldsymbol{R}_{j} \in \mathbb{R}^{N imes N}$  with entries  $\left(oldsymbol{R}_{j}
  ight)_{ii'} = \mathcal{R}_{j}\left(x_{ij}, x_{i'j}
  ight)$

#### COSSO with RKHS

• Let  $\mathcal{H}_{j} = \operatorname{span}\{\mathcal{R}_{j}\left(\cdot, x_{ij}\right), i = 1, \dots, N\}$ , we have

$$\begin{split} \|\widehat{f}_{j}\|_{\mathcal{H}_{j}}^{2} &= \left\langle \sum_{i=1}^{N} \widehat{\theta}_{ij} \mathcal{R}_{j} \left( \cdot, x_{ij} \right), \sum_{i'=1}^{N} \widehat{\theta}_{i'j} \mathcal{R}_{j} \left( \cdot, x_{i'j} \right) \right\rangle \\ &= \sum_{i=1}^{N} \sum_{i'=1}^{N} \widehat{\theta}_{ij} \widehat{\theta}_{i'j} \mathcal{R}_{j} \left( x_{ij}, x_{i'j} \right) = \widehat{\boldsymbol{\theta}}_{j}^{T} \boldsymbol{R}_{j} \widehat{\boldsymbol{\theta}}_{j} \end{split}$$

The COSSO optimization problem can be written as

$$\mathsf{minimize}_{\theta_j \in \mathbb{R}^N, j = 1, \dots, J} \left\{ \frac{1}{N} \left\| \mathbf{y} - \sum_{j=1}^J \mathbf{R}_j \theta_j \right\|_2^2 + \tau \sum_{j=1}^J \sqrt{\theta_j^\mathsf{T} \mathbf{R}_j \theta_j} \right\}$$

## Computational techniques for COSSO

• Introducing  $\gamma \in \mathbb{R}^J$ , an equivalent formulation of COSSO is

$$\min_{\substack{f_j \in \mathcal{H}_j, j = 1, \dots J \\ \gamma \ge 0}} \left\{ \frac{1}{N} \sum_{i=1}^{N} \left( y_i - \sum_{j=1}^{J} f_j(x_{ij}) \right)^2 + \sum_{j=1}^{J} \frac{1}{\gamma_j} \|f_j\|_{\mathcal{H}_j}^2 + \lambda \sum_{j=1}^{J} \gamma_j \right\}$$

• if we set  $\lambda = \tau^2/4 \leadsto$  equivalent to original COSSO formulation

#### alternates between two steps

- For  $\gamma_j$  fixed, the problem results in an additive-spline fit
- ② With the fitted additive spline fixed, updating the vector of coefficients  $\gamma = (\gamma_1, \dots, \gamma_J)$  amounts to a nonnegative lasso problem.  $\mathbf{g}_j := \mathbf{R}_j \mathbf{\theta}_j / \gamma_j \in \mathbb{R}^N$ , where  $\mathbf{f}_j = \mathbf{R}_j \mathbf{\theta}_j$

$$\min_{\gamma \geq 0} \left\{ \frac{1}{N} \| \boldsymbol{y} - \boldsymbol{G} \gamma \|_2^2 + \lambda \| \gamma \|_1 \right\}$$

#### Multiple Penalization

- Multiple ways of enforcing sparsity for a nonparametric problem. (SPAM backfitting, COSSO).
- SPAM backfitting base on a combination of  $\ell_1$ -norm:  $||f||_{N,1} := \sum_{i=1}^J ||f_i||_N$  with  $||f_i||_N^2 := \frac{1}{N} \sum_{i=1}^J f_i^2(x_{ij})$
- COSSO method uses combination of the  $\ell_1$ -norm with the Hilbert norm:

$$\|f\|_{\mathcal{H},1} := \sum_{j=1}^J \|f_j\|_{\mathcal{H}}$$

Instead of focusing on only one regularizer, one might consider the more general family of estimators

$$\min_{\substack{f_{j} \in \mathcal{H}_{j} \\ j=1,...,J}} \left\{ \frac{1}{N} \sum_{i=1}^{N} \left( y_{i} - \sum_{j=1}^{J} f_{j}\left(x_{ij}\right) \right)^{2} + \lambda_{\mathcal{H}} \sum_{j=1}^{J} \left\| f_{j} \right\|_{\mathcal{H}_{j}} + \lambda_{N} \sum_{j=1}^{J} \left\| f_{j} \right\|_{N} \right\}$$

## Why Multiple Penalization?

#### Two Penalties in Sparse Additive Models

$$L(f) = \frac{1}{2n} \sum_{i=1}^{n} (y_i - \bar{y} - f(x_i))^2 + \lambda_n ||f||_{n,1} + \rho_n ||f||_{H,1}$$

- $\|f\|_{n,1} = \sum_{j=1}^d \|f_j\|_{L^2(n)}$ 
  - Encourages **sparsity**: many  $f_i$  are set to zero.
  - Controls **model selection** when  $d \gg n$ .
- $\bullet \|f\|_{H,1} = \sum_{j=1}^{d} \|f_j\|_{H_j}$ 
  - Encourages **smoothness**: prevents overfitting in each  $f_i$ .
  - Controls function complexity via RKHS norms.

## Impact of Each Penalty

#### Impact of Each Penalty

- **Only sparsity penalty**: selects variables but risks wiggly, overfit functions.
- Only smoothness penalty: yields smooth fits but many irrelevant  $f_j$  remain nonzero.
- **Both penalties**: balance variable selection & smooth estimation, yielding minimax-optimal rates.

#### Fused lasso

The fused LASSO (signal approximator) solves the problem

$$\underset{\theta \in \mathbb{R}^n}{\mathsf{minimize}} \left\{ \sum_{i=1}^n \left( Y_i - \theta_i \right)^2 + \lambda_1 \sum_{i=1}^n |\theta_i| + \lambda_2 \sum_{i=2}^n |\theta_i - \theta_{i-1}| \right\}.$$

More generally one can use the penalty

$$\lambda_2 \sum_{i \sim j} |\theta_i - \theta_j|,$$

where  $\sim$  is a relation depending on the problem at hand.

- Fused term :  $\sum_{i=2}^{n} |\theta_i \theta_{i-1}| \rightsquigarrow$  encourages neighboring coefficients  $\theta_i$  to be similar
- Lasso term :  $\sum_{i=1}^{n} |\theta_i| \rightsquigarrow$  encourages sparsity in the coefficients  $\theta_i$

#### Variation of fused lasso

Consider the Fused LASSO with a constant term  $\theta_0$ , which means:

$$\min_{\theta_0,\theta} L(\theta_0,\theta) \triangleq \min_{\theta_0,\theta} \frac{1}{2} \sum_{i=1}^n (y_i - \theta_0 - \theta_i)^2 + \lambda_1 \sum_{i=1}^n |\theta_i| + \lambda_2 \sum_{i=2}^n |\theta_i - \theta_{i-1}|,$$

Then we can get  $\hat{\theta}_0$  by directly differentiate L by  $\theta_0$ , get:

$$\hat{\theta}_0 = \frac{1}{N} \sum_{i=1}^{N} y_i - \frac{1}{N} \sum_{i=1}^{N} \hat{\theta}_i.$$

### Variation of fused lasso (cont'd)

We conclude that

Median 
$$(\theta_i, 1 \le i \le N) = 0$$
.

Since if Median  $(\theta_i, 1 \le i \le N) = m \ne 0$ , we can change all of the  $\theta_i$  to  $\theta_i - m$ ,  $1 \le i \le N$  and change  $\theta_0$  to  $\theta_0 + m$ , making the median to 0 and get:

$$L(\theta_{0} + m, \theta - m \cdot 1) = \sum_{i=1}^{N} (y_{i} - (\theta_{0} + m) - (\theta_{i} - m))^{2} + \lambda_{1} \sum_{i=1}^{N} |\theta_{i} - m| + \lambda_{2} \sum_{i=2}^{N} |\theta_{i} - \theta_{i-1}|,$$

which makes the loss function smaller.

#### Variation of fused lasso (cont'd)

Consider the Fused LASSO with a more general type:

$$\min_{\beta_0,\beta} \left\{ \frac{1}{2} \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij} \beta_j \right)^2 + \lambda_1 \sum_{j=1}^{p} |\beta_j| + \lambda_2 \sum_{j=2}^{p} |\beta_j - \beta_{j-1}| \right\}.$$

We have  $\hat{\beta}_0$  satisfies:

$$\sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{p} x_{ij} \beta_j \right) = 0 \Rightarrow \hat{\beta}_0 = \frac{1}{N} \sum_{i=1}^{N} y_i - \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{p} x_{ij} \beta_j,$$

if  $x_{ii}$  and  $y_i$  are centered then

$$\hat{\beta}_0 = \bar{y} - \sum_{i=1}^{P} \bar{x}_{i,j} \beta_j = 0 - 0 = 0$$

we can actually omit  $\beta_0$ , independent with the choice of  $\beta_i$ ,  $1 \le i \le p$ .

## Computational techniques for fused lasso

#### Lemmas for fused lasso

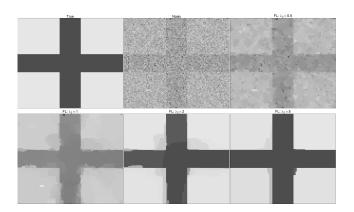
- $\bullet \ \widehat{\theta}_i\left(\lambda_1,\lambda_2\right) = \mathcal{S}_{\lambda_1}\left(\widehat{\theta}_i\left(0,\lambda_2\right)\right) \text{ for each } i=1,\ldots,N.$
- ② Suppose that for some value of  $\lambda$  and some index  $i \in \{1, \ldots, N-1\}$ , the optimal solution satisfies  $\widehat{\theta}_i(\lambda) = \widehat{\theta}_{i+1}(\lambda)$ . Then for all  $\lambda' > \lambda$ , we also have  $\widehat{\theta}_i(\lambda') = \widehat{\theta}_{i+1}(\lambda')$ .

Hence we can focus on the optimization problem:

$$\underset{\theta \in \mathbb{R}^N}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \theta_i)^2 + \lambda \sum_{i=2}^N |\theta_i - \theta_{i-1}| \right\}.$$

- reparametrize
- start from  $\lambda = 0$  and increase  $\lambda$  until all  $\theta_i$  are fused
- use lagrangian duality to solve the problem

## Case study: total variation denoising



• The idea here is that there exists a "true" image, but we only see a noisy image, from which we would like to recover the true image.

## Trend Filtering

The fused LASSO is a special case of trend filtering, which fits piecewise polynomials to data. The idea is to minimize a criterion of the form

$$\hat{\beta} = \arg\min_{\beta} \frac{1}{2} \sum_{i=1}^{n} (y_i - \beta_i)^2 + \lambda \left\| D^{(k+1)} \beta \right\|_1$$

where  $D^{(k+1)}$  is the (k+1) th order difference operator. Specifically, the second difference operator is given by

$$D^{(2)} = \left[ \begin{array}{cccccc} -1 & 2 & -1 & 0 & \cdots & 0 \\ 0 & -1 & 2 & -1 & \cdots & 0 \\ \vdots & & & \ddots & \vdots \\ 0 & 0 & \cdots & -1 & 2 & -1 \end{array} \right]$$

And the loss function is

$$\hat{\beta} = \arg\min_{\beta} \frac{1}{2} \sum_{i=1}^{n} (y_i - \beta_i)^2 + \lambda \sum_{i=1}^{n-2} |\beta_i - 2\beta_{i+1} + \beta_{i+2}|$$

## Change Point Detection

The penalty is assumed that the observations occur at evenly spaced points. For arbitrary input points  $x_1 < x_2 < \cdots < x_n$ , we can also use the following penalty <sup>1</sup>:

$$\frac{1}{2} \sum_{i=1}^{N} (y_i - \theta_i)^2 + \lambda \sum_{i=1}^{n-2} \left| \frac{\theta_i - \theta_{i+1}}{x_i - x_{i+1}} - \frac{\theta_{i+1} - \theta_{i+2}}{x_{i+1} - x_{i+2}} \right|$$

It encourages the slopes of the adjacent linear segments is the same, leading to piecewise linear fits.

• Here  $x_i$  stands for the time stamp/input feature of observation  $y_i$ .

<sup>&</sup>lt;sup>1</sup>Tibshirani, R.J. (2014). Adaptive piecewise polynomial estimation via trend filtering. Annals of Statistics, 42(1), 285-323.

#### Isotonic regression

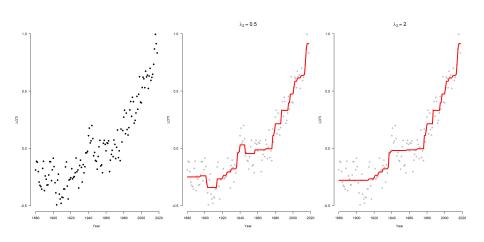
• Classical isotonic regression problem is to fit a nondecreasing sequence to a given sequence of observations  $(y_1, \ldots, y_N)$  by solving the optimization problem

$$\underset{\theta \in \mathbb{R}^N}{\text{minimize}} \left\{ \sum_{i=1}^N (y_i - \theta_i)^2 \right\} \text{ subject to } \theta_1 \leq \theta_2 \leq \ldots \leq \theta_N$$

• Nearly isotonic regression is a natural relaxation, in which we introduce a regularization parameter  $\lambda \geq 0$ , and instead solve the penalized problem

$$\underset{\boldsymbol{\theta} \in \mathbb{R}^{N}}{\text{minimize}} \left\{ \frac{1}{2} \sum_{i=1}^{N} (y_{i} - \theta_{i})^{2} + \lambda \sum_{i=1}^{N-1} (\theta_{i} - \theta_{i+1})_{+} \right\}$$

## Case study: global warming



#### **SCAD**

 A variety of nonconvex penalties have been proposed: one of the earliest and most influential was the smoothly clipped absolute deviations (SCAD) penalty:

$$P_{\lambda,\gamma}(\theta) = \begin{cases} \lambda |\theta| & \text{if } |\theta| \le \lambda \\ \frac{2\gamma\lambda|\theta| - \theta^2 - \lambda^2}{2(\gamma - 1)} & \text{if } \lambda < |\theta| < \gamma\lambda \\ \frac{\lambda^2(\gamma + 1)}{2} & \text{if } |\theta| \ge \gamma\lambda \end{cases}$$

for  $\gamma > 2$ 

• Note that SCAD coincides with the lasso until  $|\theta|=\lambda$ , then smoothly transitions to a quadratic function until  $|\theta|=\gamma\lambda$ , after which it remains constant for all  $|\theta|>\gamma\lambda$ 

#### MCP

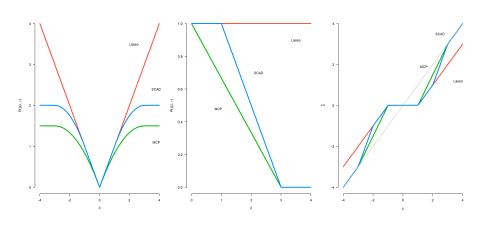
• The MC+ penalty on each coordinate is defined by

$$P_{\lambda,\gamma}(\theta) := \int_0^{|\theta|} \left(1 - \frac{x}{\lambda \gamma}\right)_+ dx = \begin{cases} |\theta| - \frac{|\theta|^2}{2\lambda \gamma}, & |\theta| < \lambda \gamma, \\ \frac{\lambda \gamma}{2}, & |\theta| \ge \lambda \gamma. \end{cases}$$

• For squared-error loss we pose the (nonconvex) optimization problem

$$\underset{\boldsymbol{\beta} \in \mathbb{R}^{p}}{\operatorname{minimize}} \left\{ \frac{1}{2} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \|_{2}^{2} + \sum_{j=1}^{p} P_{\lambda,\gamma} \left( \beta_{j} \right) \right\},$$

## SCAD, MCP and lasso in 1 dimension



# Questions or comments?

#### References

- Statistical Learning with Sparsity
- Linear models and extensions, Peng Ding
- The Elements of Statistical Learning
- Course slides for sparse learning in ETH Zurich
- Purdue ECE695Notes
- U Iowa High-Dimensional Data Analysis (BIOS 7240) course notes
- Survival analysis lecture notes, SYSU