

Text as Data

Justin Grimmer

Associate Professor
Department of Political Science
Stanford University

November 10th, 2014

Supervised Learning Methods

1) Task

- Classify documents to pre existing categories
- Measure the proportion of documents in each category

2) Objective function

1) Penalized Regressions

- Ridge regression
- LASSO regression

2) Classification Surface \rightsquigarrow Support Vector Machines

3) Measure Proportions \rightsquigarrow Naive Bayes(ish) objective

3) Optimization

- Depends on method

4) Validation

- Obtain predicted fit for new data $f(\mathbf{X}_i, \hat{\theta})$
- Examine prediction performance \rightsquigarrow compare classification to **gold standard**

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

We represent each document i is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$.

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

We represent each document i is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$.

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2$$

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

We represent each document i is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$.

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2$$

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2 \right\}$$

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

We represent each document i is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$.

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2$$

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2 \right\} \\ &= (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} \end{aligned}$$

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

We represent each document i is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$.

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2$$

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2 \right\} \\ &= (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} \end{aligned}$$

Problem:

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

We represent each document i is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$.

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2$$

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2 \right\} \\ &= (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} \end{aligned}$$

Problem:

- J will likely be large (perhaps $J > N$)

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

We represent each document i is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$.

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2$$

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2 \right\} \\ &= (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} \end{aligned}$$

Problem:

- J will likely be large (perhaps $J > N$)
- There many correlated variables

Regression models

Suppose we have N documents, with each document i having label $y_i \in \{-1, 1\} \rightsquigarrow \{\text{liberal, conservative}\}$

We represent each document i is $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$.

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2$$

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N (y_i - \boldsymbol{\beta}' \mathbf{x}_i)^2 \right\} \\ &= (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} \end{aligned}$$

Problem:

- J will likely be large (perhaps $J > N$)
- There many correlated variables

Predictions will be **variable**

Mean Square Error

Suppose θ is some value of the true parameter

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

$$\text{Bias} = E[\hat{\theta} - \theta]$$

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

$$\text{Bias} = E[\hat{\theta} - \theta]$$

We may care about average distance from truth

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

$$\text{Bias} = E[\hat{\theta} - \theta]$$

We may care about average distance from truth

$$E[(\hat{\theta} - \theta)^2]$$

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

$$\text{Bias} = E[\hat{\theta} - \theta]$$

We may care about average distance from truth

$$E[(\hat{\theta} - \theta)^2] = E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2$$

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

$$\text{Bias} = E[\hat{\theta} - \theta]$$

We may care about average distance from truth

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2 \\ &= E[\hat{\theta}^2] - E[\hat{\theta}]^2 + E[\hat{\theta}]^2 - 2\theta E[\hat{\theta}] + \theta^2 \end{aligned}$$

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

$$\text{Bias} = E[\hat{\theta} - \theta]$$

We may care about average distance from truth

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2 \\ &= E[\hat{\theta}^2] - E[\hat{\theta}]^2 + E[\hat{\theta}]^2 - 2\theta E[\hat{\theta}] + \theta^2 \\ &= E[\hat{\theta}^2] - E[\hat{\theta}]^2 + (E[\hat{\theta}] - \theta)^2 \end{aligned}$$

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

$$\text{Bias} = E[\hat{\theta} - \theta]$$

We may care about average distance from truth

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2 \\ &= E[\hat{\theta}^2] - E[\hat{\theta}]^2 + E[\hat{\theta}]^2 - 2\theta E[\hat{\theta}] + \theta^2 \\ &= E[\hat{\theta}^2] - E[\hat{\theta}]^2 + (E[\hat{\theta}] - \theta)^2 \\ &= \text{Var}(\hat{\theta}) + \text{Bias}^2 \end{aligned}$$

Mean Square Error

Suppose θ is some value of the true parameter

Bias:

$$\text{Bias} = E[\hat{\theta} - \theta]$$

We may care about average distance from truth

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= E[\hat{\theta}^2] - 2\theta E[\hat{\theta}] + \theta^2 \\ &= E[\hat{\theta}^2] - E[\hat{\theta}]^2 + E[\hat{\theta}]^2 - 2\theta E[\hat{\theta}] + \theta^2 \\ &= E[\hat{\theta}^2] - E[\hat{\theta}]^2 + (E[\hat{\theta}] - \theta)^2 \\ &= \text{Var}(\hat{\theta}) + \text{Bias}^2 \end{aligned}$$

To reduce MSE, we are willing to induce bias to decrease variance \rightsquigarrow
methods that **shrink** coefficients toward zero

Ridge Regression

Penalty for model complexity

Ridge Regression

Penalty for model complexity

$$f(\beta, \mathbf{X}, \mathbf{Y})$$

Ridge Regression

Penalty for model complexity

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2$$

Ridge Regression

Penalty for model complexity

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^J \beta_j^2}_{\text{Penalty}}$$

Ridge Regression

Penalty for model complexity

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^J \beta_j^2}_{\text{Penalty}}$$

where:

Ridge Regression

Penalty for model complexity

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^J \beta_j^2}_{\text{Penalty}}$$

where:

- $\beta_0 \rightsquigarrow$ intercept

Ridge Regression

Penalty for model complexity

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^J \beta_j^2}_{\text{Penalty}}$$

where:

- $\beta_0 \rightsquigarrow$ intercept
- $\lambda \rightsquigarrow$ penalty parameter

Ridge Regression \rightsquigarrow Optimization

$$\beta^{\text{Ridge}} = \arg \min_{\beta} \{f(\beta, \mathbf{X}, \mathbf{Y})\}$$

Ridge Regression \rightsquigarrow Optimization

$$\begin{aligned}\beta^{\text{Ridge}} &= \arg \min_{\beta} \{f(\beta, \mathbf{X}, \mathbf{Y})\} \\ &= \arg \min_{\beta} \left\{ \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \beta_j^2 \right\}\end{aligned}$$

Ridge Regression \rightsquigarrow Optimization

$$\begin{aligned}\beta^{\text{Ridge}} &= \arg \min_{\beta} \{f(\beta, \mathbf{X}, \mathbf{Y})\} \\ &= \arg \min_{\beta} \left\{ \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \beta_j^2 \right\} \\ &= \arg \min_{\beta} \left\{ (\mathbf{Y} - \mathbf{X}'\beta)' (\mathbf{Y} - \mathbf{X}'\beta) + \lambda \beta' \beta \right\}\end{aligned}$$

Demean the data and set $\beta_0 = \bar{y} = \sum_{i=1}^N \frac{y_i}{N}$

Ridge Regression \rightsquigarrow Optimization

$$\begin{aligned}\beta^{\text{Ridge}} &= \arg \min_{\beta} \{f(\beta, \mathbf{X}, \mathbf{Y})\} \\ &= \arg \min_{\beta} \left\{ \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \beta_j^2 \right\} \\ &= \arg \min_{\beta} \left\{ (\mathbf{Y} - \mathbf{X}'\beta)' (\mathbf{Y} - \mathbf{X}'\beta) + \lambda \beta' \beta \right\} \\ &= (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_J)^{-1} \mathbf{X}'\mathbf{Y}\end{aligned}$$

Demean the data and set $\beta_0 = \bar{y} = \sum_{i=1}^N \frac{y_i}{N}$

Ridge Regression \rightsquigarrow Optimization

$$\begin{aligned}\beta^{\text{Ridge}} &= \arg \min_{\beta} \{f(\beta, \mathbf{X}, \mathbf{Y})\} \\ &= \arg \min_{\beta} \left\{ \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \beta_j^2 \right\} \\ &= \arg \min_{\beta} \left\{ (\mathbf{Y} - \mathbf{X}'\beta)' (\mathbf{Y} - \mathbf{X}'\beta) + \lambda \beta' \beta \right\} \\ &= (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_J)^{-1} \mathbf{X}'\mathbf{Y}\end{aligned}$$

Demean the data and set $\beta_0 = \bar{y} = \sum_{i=1}^N \frac{y_i}{N}$

Ridge Regression \rightsquigarrow Intuition (1)

Suppose $\mathbf{X}'\mathbf{X} = I_J$.

Ridge Regression \rightsquigarrow Intuition (1)

Suppose $\mathbf{X}'\mathbf{X} = I_J$.

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

Ridge Regression \rightsquigarrow Intuition (1)

Suppose $\mathbf{X}'\mathbf{X} = I_J$.

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \\ &= \mathbf{X}'\mathbf{Y}\end{aligned}$$

Ridge Regression \rightsquigarrow Intuition (1)

Suppose $\mathbf{X}'\mathbf{X} = I_J$.

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \\ &= \mathbf{X}'\mathbf{Y} \\ \beta^{\text{ridge}} &= (\mathbf{X}'\mathbf{X} + \lambda I_J)^{-1} \mathbf{X}'\mathbf{Y}\end{aligned}$$

Ridge Regression \rightsquigarrow Intuition (1)

Suppose $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$.

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \\ &= \mathbf{X}'\mathbf{Y} \\ \boldsymbol{\beta}^{\text{ridge}} &= (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_J)^{-1} \mathbf{X}'\mathbf{Y} \\ &= (\mathbf{I}_J + \lambda \mathbf{I}_J) \mathbf{X}'\mathbf{Y}\end{aligned}$$

Ridge Regression \rightsquigarrow Intuition (1)

Suppose $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$.

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \\ &= \mathbf{X}'\mathbf{Y} \\ \boldsymbol{\beta}^{\text{ridge}} &= (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_J)^{-1} \mathbf{X}'\mathbf{Y} \\ &= (\mathbf{I}_J + \lambda\mathbf{I}_J) \mathbf{X}'\mathbf{Y} \\ &= (\mathbf{I}_J + \lambda\mathbf{I}_J) \hat{\boldsymbol{\beta}}\end{aligned}$$

Ridge Regression \rightsquigarrow Intuition (1)

Suppose $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$.

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \\ &= \mathbf{X}'\mathbf{Y} \\ \beta^{\text{ridge}} &= (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_J)^{-1} \mathbf{X}'\mathbf{Y} \\ &= (\mathbf{I}_J + \lambda \mathbf{I}_J) \mathbf{X}'\mathbf{Y} \\ &= (\mathbf{I}_J + \lambda \mathbf{I}_J) \hat{\beta} \\ \beta_j^{\text{Ridge}} &= \frac{\hat{\beta}_j}{1 + \lambda}\end{aligned}$$

Ridge Regression \rightsquigarrow Intuition (2)

$$\beta_j \sim \text{Normal}(0, \tau^2)$$

$$y_i \sim \text{Normal}(\beta_0 + \mathbf{x}'_i \boldsymbol{\beta}, \sigma^2)$$

Ridge Regression \rightsquigarrow Intuition (2)

$$\beta_j \sim \text{Normal}(0, \tau^2)$$

$$y_i \sim \text{Normal}(\beta_0 + \mathbf{x}'_i \boldsymbol{\beta}, \sigma^2)$$

$$p(\boldsymbol{\beta} | \mathbf{X}, \mathbf{Y}) \propto \prod_{j=1}^J p(\beta_j) \prod_{i=1}^N p(y_i | \mathbf{x}_i, \boldsymbol{\beta})$$

Ridge Regression \rightsquigarrow Intuition (2)

$$\begin{aligned}\beta_j &\sim \text{Normal}(0, \tau^2) \\ y_i &\sim \text{Normal}(\beta_0 + \mathbf{x}'_i \boldsymbol{\beta}, \sigma^2)\end{aligned}$$

$$\begin{aligned}p(\boldsymbol{\beta} | \mathbf{X}, \mathbf{Y}) &\propto \prod_{j=1}^J p(\beta_j) \prod_{i=1}^N p(y_i | \mathbf{x}_i, \boldsymbol{\beta}) \\ &\propto \prod_{j=1}^J \frac{1}{\sqrt{2\pi\tau}} \exp\left(-\frac{\beta_j^2}{2\tau^2}\right) \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y_i - \beta_0 + \mathbf{x}'_i \boldsymbol{\beta})^2}{2\sigma^2}\right)\end{aligned}$$

Ridge Regression \rightsquigarrow Intuition (2)

$$\log p(\boldsymbol{\beta} | \mathbf{X}, \mathbf{Y}) = - \sum_{j=1}^J \frac{\beta_j^2}{2\tau^2} - \sum_{i=1}^N \frac{(y_i - \beta_0 + \mathbf{x}'_i \boldsymbol{\beta})^2}{2\sigma^2}$$

Ridge Regression \rightsquigarrow Intuition (2)

$$\begin{aligned}\log p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{Y}) &= -\sum_{j=1}^J \frac{\beta_j^2}{2\tau^2} - \sum_{i=1}^N \frac{(y_i - \beta_0 + \mathbf{x}'\boldsymbol{\beta})^2}{2\sigma^2} \\ -2\sigma^2 \log p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{Y}) &= \sum_{i=1}^N (y_i - \beta_0 + \mathbf{x}'\boldsymbol{\beta})^2 + \sum_{j=1}^J \frac{\sigma^2}{\tau^2} \beta_j^2\end{aligned}$$

Ridge Regression \rightsquigarrow Intuition (2)

$$\begin{aligned}\log p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{Y}) &= -\sum_{j=1}^J \frac{\beta_j^2}{2\tau^2} - \sum_{i=1}^N \frac{(y_i - \beta_0 + \mathbf{x}'\boldsymbol{\beta})^2}{2\sigma^2} \\ -2\sigma^2 \log p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{Y}) &= \sum_{i=1}^N (y_i - \beta_0 + \mathbf{x}'\boldsymbol{\beta})^2 + \sum_{j=1}^J \frac{\sigma^2}{\tau^2} \beta_j^2\end{aligned}$$

where:

Ridge Regression \rightsquigarrow Intuition (2)

$$\begin{aligned}\log p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{Y}) &= -\sum_{j=1}^J \frac{\beta_j^2}{2\tau^2} - \sum_{i=1}^N \frac{(y_i - \beta_0 + \mathbf{x}'\boldsymbol{\beta})^2}{2\sigma^2} \\ -2\sigma^2 \log p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{Y}) &= \sum_{i=1}^N (y_i - \beta_0 + \mathbf{x}'\boldsymbol{\beta})^2 + \sum_{j=1}^J \frac{\sigma^2}{\tau^2} \beta_j^2\end{aligned}$$

where:

$$- \lambda = \frac{\sigma^2}{\tau^2} \beta_j^2$$

Lasso Regression Objective Function/Optimization

Different Penalty for Model Complexity

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \underbrace{|\beta_j|}_{\text{Penalty}}$$

Lasso Regression Objective Function/Optimization

Different Penalty for Model Complexity

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \underbrace{|\beta_j|}_{\text{Penalty}}$$

- Optimization is non-linear (Absolute Value)

Lasso Regression Objective Function/Optimization

Different Penalty for Model Complexity

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \underbrace{|\beta_j|}_{\text{Penalty}}$$

- Optimization is non-linear (Absolute Value)
 - Coordinate Descent

Lasso Regression Objective Function/Optimization

Different Penalty for Model Complexity

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \underbrace{|\beta_j|}_{\text{Penalty}}$$

- Optimization is non-linear (Absolute Value)
 - Coordinate Descent
 - Start with Ridge

Lasso Regression Objective Function/Optimization

Different Penalty for Model Complexity

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \underbrace{|\beta_j|}_{\text{Penalty}}$$

- Optimization is non-linear (Absolute Value)
 - Coordinate Descent
 - Start with Ridge
 - Sub-differential, update steps

Lasso Regression Objective Function/Optimization

Different Penalty for Model Complexity

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \left(y_i - \beta_0 + \sum_{j=1}^J \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^J \underbrace{|\beta_j|}_{\text{Penalty}}$$

- Optimization is non-linear (Absolute Value)
 - Coordinate Descent
 - Start with Ridge
 - Sub-differential, update steps
- Induces **sparsity** \rightsquigarrow sets some coefficients to zero

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Suppose again $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Suppose again $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$

$$f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^J |\beta_j|$$

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Suppose again $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$

$$\begin{aligned} f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^J |\beta_j| \\ &= -2\mathbf{X}'\mathbf{Y}\boldsymbol{\beta} + \boldsymbol{\beta}'\boldsymbol{\beta} + \lambda \sum_{j=1}^J |\beta_j| \end{aligned}$$

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Suppose again $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$

$$\begin{aligned} f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^J |\beta_j| \\ &= -2\mathbf{X}'\mathbf{Y}\boldsymbol{\beta} + \boldsymbol{\beta}'\boldsymbol{\beta} + \lambda \sum_{j=1}^J |\beta_j| \end{aligned}$$

The coefficient is

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Suppose again $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$

$$\begin{aligned} f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^J |\beta_j| \\ &= -2\mathbf{X}'\mathbf{Y}\boldsymbol{\beta} + \boldsymbol{\beta}'\boldsymbol{\beta} + \lambda \sum_{j=1}^J |\beta_j| \end{aligned}$$

The coefficient is

$$\beta_j^{\text{LASSO}} = \text{sign}(\hat{\beta}_j) \left(|\hat{\beta}_j| - \lambda \right)_+$$

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Suppose again $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$

$$\begin{aligned} f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^J |\beta_j| \\ &= -2\mathbf{X}'\mathbf{Y}\boldsymbol{\beta} + \boldsymbol{\beta}'\boldsymbol{\beta} + \lambda \sum_{j=1}^J |\beta_j| \end{aligned}$$

The coefficient is

$$\beta_j^{\text{LASSO}} = \text{sign}(\hat{\beta}_j) \left(|\hat{\beta}_j| - \lambda \right)_+$$

- $\text{sign}(\cdot) \rightsquigarrow 1$ or -1

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Suppose again $\mathbf{X}'\mathbf{X} = \mathbf{I}_J$

$$\begin{aligned} f(\boldsymbol{\beta}, \mathbf{X}, \mathbf{Y}) &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^J |\beta_j| \\ &= -2\mathbf{X}'\mathbf{Y}\boldsymbol{\beta} + \boldsymbol{\beta}'\boldsymbol{\beta} + \lambda \sum_{j=1}^J |\beta_j| \end{aligned}$$

The coefficient is

$$\beta_j^{\text{LASSO}} = \text{sign}(\hat{\beta}_j) \left(|\hat{\beta}_j| - \lambda \right)_+$$

- $\text{sign}(\cdot) \rightsquigarrow 1$ or -1
- $\left(|\hat{\beta}_j| - \lambda \right)_+ = \max(|\hat{\beta}_j| - \lambda, 0)$

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Compare soft assignment

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Compare soft assignment

$$\beta_j^{\text{LASSO}} = \text{sign}(\hat{\beta}_j) (|\hat{\beta}_j| - \lambda)_+$$

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Compare soft assignment

$$\beta_j^{\text{LASSO}} = \text{sign}(\hat{\beta}_j) (|\hat{\beta}_j| - \lambda)_+$$

With hard assignment, selecting M biggest components

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Compare soft assignment

$$\beta_j^{\text{LASSO}} = \text{sign}(\hat{\beta}_j) \left(|\hat{\beta}_j| - \lambda \right)_+$$

With hard assignment, selecting M biggest components

$$\beta_j^{\text{subset}} = \hat{\beta}_j \cdot I \left(|\hat{\beta}_j| \geq |\hat{\beta}_{(M)}| \right)$$

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Compare soft assignment

$$\beta_j^{\text{LASSO}} = \text{sign}(\hat{\beta}_j) \left(|\hat{\beta}_j| - \lambda \right)_+$$

With hard assignment, selecting M biggest components

$$\beta_j^{\text{subset}} = \hat{\beta}_j \cdot I\left(|\hat{\beta}_j| \geq |\hat{\beta}_{(M)}|\right)$$

Intuition 2: Prior on coefficients \rightsquigarrow Double exponential

Lasso Regression \rightsquigarrow Intuition 1, Soft Thresholding

Compare soft assignment

$$\beta_j^{\text{LASSO}} = \text{sign}(\hat{\beta}_j) \left(|\hat{\beta}_j| - \lambda \right)_+$$

With hard assignment, selecting M biggest components

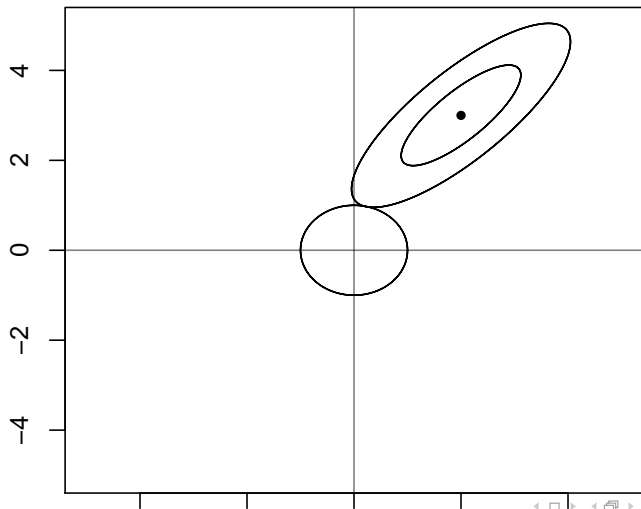
$$\beta_j^{\text{subset}} = \hat{\beta}_j \cdot I\left(|\hat{\beta}_j| \geq |\hat{\beta}_{(M)}|\right)$$

Intuition 2: Prior on coefficients \rightsquigarrow Double exponential

Why does LASSO induce sparsity?

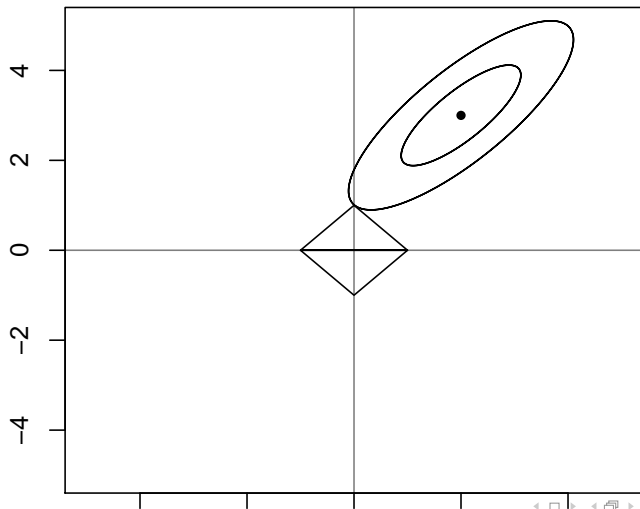
Comparing Ridge and LASSO

Ridge Regression



Comparing Ridge and LASSO

LASSO Regression



Comparing Ridge and LASSO

Contrast $\beta = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ and $\tilde{\beta} = (1, 0)$

Comparing Ridge and LASSO

Contrast $\beta = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ and $\tilde{\beta} = (1, 0)$

Under ridge:

Comparing Ridge and LASSO

Contrast $\beta = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ and $\tilde{\beta} = (1, 0)$

Under ridge:

$$\sum_{j=1}^2 \beta_j^2 = \frac{1}{2} + \frac{1}{2} = 1$$

Comparing Ridge and LASSO

Contrast $\beta = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ and $\tilde{\beta} = (1, 0)$

Under ridge:

$$\sum_{j=1}^2 \beta_j^2 = \frac{1}{2} + \frac{1}{2} = 1$$

$$\sum_{j=1}^2 \tilde{\beta}_j^2 = 1 + 0 = 1$$

Comparing Ridge and LASSO

Contrast $\beta = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ and $\tilde{\beta} = (1, 0)$

Under ridge:

$$\sum_{j=1}^2 \beta_j^2 = \frac{1}{2} + \frac{1}{2} = 1$$

$$\sum_{j=1}^2 \tilde{\beta}_j^2 = 1 + 0 = 1$$

Under LASSO

Comparing Ridge and LASSO

Contrast $\beta = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ and $\tilde{\beta} = (1, 0)$

Under ridge:

$$\sum_{j=1}^2 \beta_j^2 = \frac{1}{2} + \frac{1}{2} = 1$$

$$\sum_{j=1}^2 \tilde{\beta}_j^2 = 1 + 0 = 1$$

Under LASSO

$$\sum_{j=1}^2 |\beta_j| = \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} = \sqrt{2}$$

Comparing Ridge and LASSO

Contrast $\beta = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ and $\tilde{\beta} = (1, 0)$

Under ridge:

$$\sum_{j=1}^2 \beta_j^2 = \frac{1}{2} + \frac{1}{2} = 1$$

$$\sum_{j=1}^2 \tilde{\beta}_j^2 = 1 + 0 = 1$$

Under LASSO

$$\sum_{j=1}^2 |\beta_j| = \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} = \sqrt{2}$$

$$\sum_{j=1}^2 |\tilde{\beta}_j| = 1 + 0 = 1$$

Selecting λ

How do we determine λ ? \rightsquigarrow Cross validation (lecture on Thursday)

To the R code!

Selecting λ

How do we determine λ ? \rightsquigarrow Cross validation (lecture on Thursday)
Applying models gives score (probability) of document belong to class \rightsquigarrow
threshold to classify

To the R code!

Selecting λ

How do we determine λ ? \rightsquigarrow Cross validation (lecture on Thursday)
Applying models gives score (probability) of document belong to class \rightsquigarrow
threshold to classify
To the R code!

Assessing Models (Elements of Statistical Learning)

- **Model Selection**: tuning parameters to select final model (next week's discussion)
- **Model assessment** : after selecting model, estimating error in classification

Comparing Training and Validation Set

Text classification and model assessment

- **Replicate** classification exercise with **validation** set
- General **principle** of classification/prediction
- Compare supervised learning labels to hand labels

Confusion matrix

Comparing Training and Validation Set

Representation of Test Statistics from Dictionary week (along with some new ones)

Classification (algorithm)	Actual Label	
	Liberal	Conservative
Liberal	True Liberal	False Liberal
Conservative	False Conservative	True Conservative

Comparing Training and Validation Set

Representation of Test Statistics from Dictionary week (along with some new ones)

Classification (algorithm)	Actual Label	
	Liberal	Conservative
Liberal	True Liberal	False Liberal
Conservative	False Conservative	True Conservative

$$\text{Accuracy} = \frac{\text{TrueLib} + \text{TrueCons}}{\text{TrueLib} + \text{TrueCons} + \text{FalseLib} + \text{FalseCons}}$$

Comparing Training and Validation Set

Representation of Test Statistics from Dictionary week (along with some new ones)

Classification (algorithm)	Actual Label	
	Liberal	Conservative
Liberal	True Liberal	False Liberal
Conservative	False Conservative	True Conservative

$$\text{Accuracy} = \frac{\text{TrueLib} + \text{TrueCons}}{\text{TrueLib} + \text{TrueCons} + \text{FalseLib} + \text{FalseCons}}$$

$$\text{Precision}_{\text{Liberal}} = \frac{\text{True Liberal}}{\text{True Liberal} + \text{False Liberal}}$$

Comparing Training and Validation Set

Representation of Test Statistics from Dictionary week (along with some new ones)

Classification (algorithm)	Actual Label	
	Liberal	Conservative
Liberal	True Liberal	False Liberal
Conservative	False Conservative	True Conservative

$$\text{Accuracy} = \frac{\text{TrueLib} + \text{TrueCons}}{\text{TrueLib} + \text{TrueCons} + \text{FalseLib} + \text{FalseCons}}$$

$$\text{Precision}_{\text{Liberal}} = \frac{\text{True Liberal}}{\text{True Liberal} + \text{False Liberal}}$$

$$\text{Recall}_{\text{Liberal}} = \frac{\text{True Liberal}}{\text{True Liberal} + \text{False Conservative}}$$

Comparing Training and Validation Set

Representation of Test Statistics from Dictionary week (along with some new ones)

	Actual Label	
Classification (algorithm)	Liberal	Conservative
Liberal	True Liberal	False Liberal
Conservative	False Conservative	True Conservative

$$\text{Accuracy} = \frac{\text{TrueLib} + \text{TrueCons}}{\text{TrueLib} + \text{TrueCons} + \text{FalseLib} + \text{FalseCons}}$$

$$\text{Precision}_{\text{Liberal}} = \frac{\text{True Liberal}}{\text{True Liberal} + \text{False Liberal}}$$

$$\text{Recall}_{\text{Liberal}} = \frac{\text{True Liberal}}{\text{True Liberal} + \text{False Conservative}}$$

$$F_{\text{Liberal}} = \frac{2\text{Precision}_{\text{Liberal}}\text{Recall}_{\text{Liberal}}}{\text{Precision}_{\text{Liberal}} + \text{Recall}_{\text{Liberal}}}$$

Comparing Training and Validation Set

Representation of Test Statistics from Dictionary week (along with some new ones)

	Actual Label	
Classification (algorithm)	Liberal	Conservative
Liberal	True Liberal	False Liberal
Conservative	False Conservative	True Conservative

$$\text{Accuracy} = \frac{\text{TrueLib} + \text{TrueCons}}{\text{TrueLib} + \text{TrueCons} + \text{FalseLib} + \text{FalseCons}}$$

$$\text{Precision}_{\text{Liberal}} = \frac{\text{True Liberal}}{\text{True Liberal} + \text{False Liberal}}$$

$$\text{Recall}_{\text{Liberal}} = \frac{\text{True Liberal}}{\text{True Liberal} + \text{False Conservative}}$$

$$F_{\text{Liberal}} = \frac{2\text{Precision}_{\text{Liberal}}\text{Recall}_{\text{Liberal}}}{\text{Precision}_{\text{Liberal}} + \text{Recall}_{\text{Liberal}}}$$

ROC Curve

ROC as a measure of model performance

$$\text{Recall}_{\text{Liberal}} = \frac{\text{True Liberal}}{\text{True Liberal} + \text{False Conservative}}$$
$$\text{Recall}_{\text{Conservative}} = \frac{\text{True Conservative}}{\text{True Conservative} + \text{False Liberal}}$$

Tension:

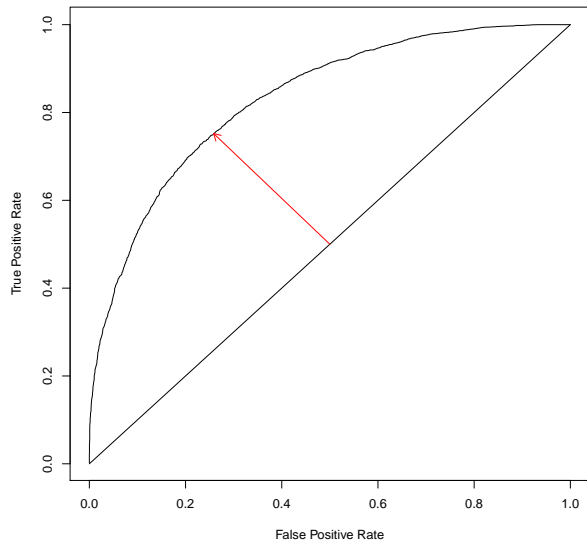
- Everything liberal: $\text{Recall}_{\text{Liberal}} = 1$; $\text{Recall}_{\text{Conservative}} = 0$
- Everything conservative: $\text{Recall}_{\text{Liberal}} = 0$; $\text{Recall}_{\text{Conservative}} = 1$

Characterize Tradeoff:

Plot True Positive Rate $\text{Recall}_{\text{Liberal}}$

False Positive Rate $(1 - \text{Recall}_{\text{Conservative}})$

Precision/Recall Tradeoff



Simple Classification Example

Analyzing house press releases

Hand Code: 1,000 press releases

- Advertising
- Credit Claiming
- Position Taking

Divide 1,000 press releases into two sets

- 500: Training set
- 500: Test set

Initial exploration: provides baseline measurement at classifier performances

Improve: through improving model fit

Example from First Model Fit

Classification (Naive Bayes)	Actual Label		
	Position Taking	Advertising	Credit Claim.
Position Taking	10	0	0
Advertising	2	40	2
Credit Claiming	80	60	306

$$\text{Accuracy} = \frac{10 + 40 + 306}{500} = 0.71$$

$$\text{Precision}_{PT} = \frac{10}{10} = 1$$

$$\text{Recall}_{PT} = \frac{10}{10 + 2 + 80} = 0.11$$

$$\text{Precision}_{AD} = \frac{40}{40 + 2 + 2} = 0.91$$

$$\text{Recall}_{AD} = \frac{40}{40 + 60} = 0.4$$

$$\text{Precision}_{Credit} = \frac{306}{306 + 80 + 60} = 0.67$$

$$\text{Recall}_{Credit} = \frac{306}{306 + 2} = 0.99$$

Fit Statistics in R

RWeka library provides **Amazing** functionality.
You can easily code them yourself

Support Vector Machines

Document i is an $J \times 1$ vector of counts

$$\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{ji})$$

Suppose we have **two** classes, C_1, C_2 .

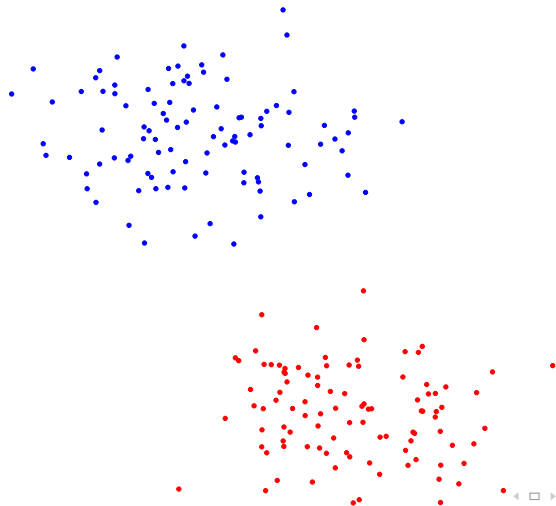
$$Y_i = 1 \text{ if } i \text{ is in class 1}$$

$$Y_i = -1 \text{ if } i \text{ is in class 2}$$

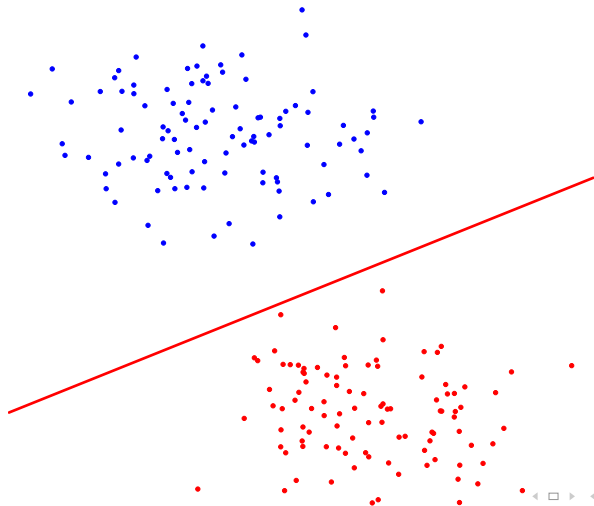
Suppose they are **separable**:

- Draw a line between groups
- Goal: identify the line **in the middle**
- **Maximum margin**

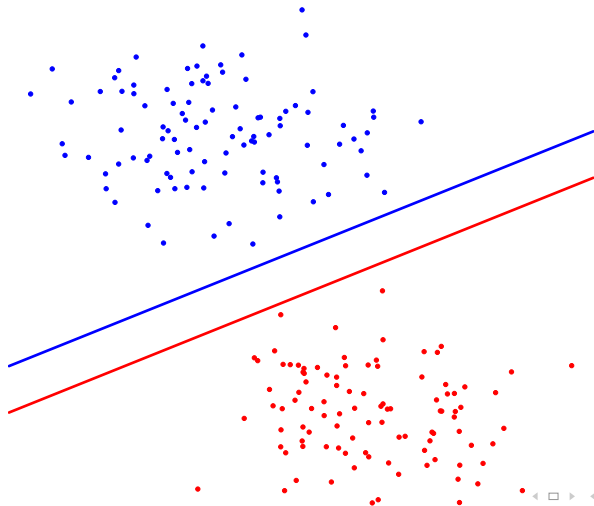
Support Vector Machines: Maximum Margin Classifier (Bishop 2006)



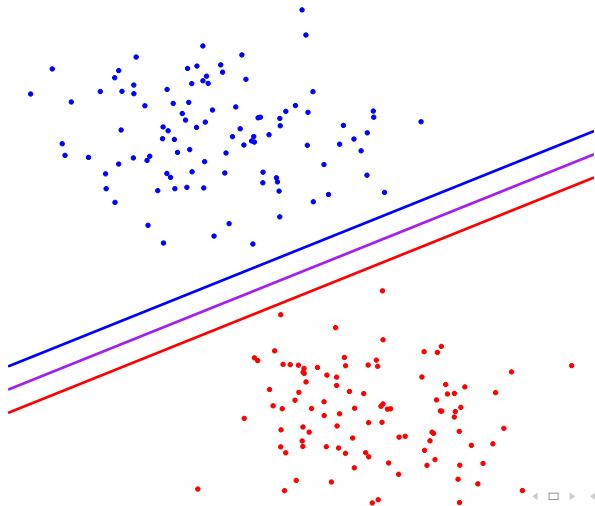
Support Vector Machines: Maximum Margin Classifier (Bishop 2006)



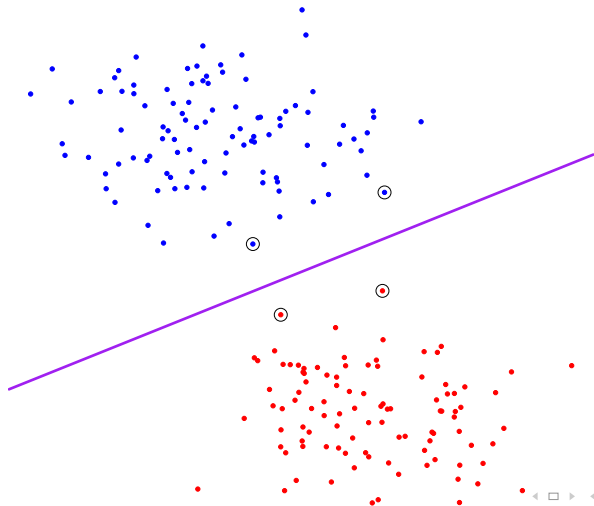
Support Vector Machines: Maximum Margin Classifier (Bishop 2006)



Support Vector Machines: Maximum Margin Classifier (Bishop 2006)



Support Vector Machines: Maximum Margin Classifier (Bishop 2006)



Support Vector Machines: Algebra (Bishop 2006)

Goal create a score to classify:

$$s(\mathbf{x}_i) = \boldsymbol{\beta}' \mathbf{x}_i + b$$

- $\boldsymbol{\beta}$ Determines orientation of surface (slope)
- b determines location (moves surface up or down)
- If $s(\mathbf{x}_i) > 0 \rightarrow$ class 1
- If $s(\mathbf{x}_i) < 0 \rightarrow$ class 2
- $\frac{|s(\mathbf{x}_i)|}{\|\boldsymbol{\beta}\|} =$ Document distance from decision surface (margin)

Support Vector Machines: Algebra (Bishop 2006)

Objective function: **maximum margin**

Support Vector Machines: Algebra (Bishop 2006)

Objective function: **maximum margin**

$\min_i [|s(\mathbf{x}_i)|]$: Point closest to decision surface

Support Vector Machines: Algebra (Bishop 2006)

Objective function: **maximum margin**

$\min_i [|s(\mathbf{x}_i)|]$: Point closest to decision surface

We want to identify β and b to maximize the margin:

Support Vector Machines: Algebra (Bishop 2006)

Objective function: **maximum margin**

$\min_i [|(s(\mathbf{x}_i)|)|]$: Point closest to decision surface

We want to identify β and b to maximize the margin:

$$\arg \max_{\beta, b} \left\{ \frac{1}{\|\beta\|} \min_i [|(s(\mathbf{x}_i)|)|] \right\}$$

Support Vector Machines: Algebra (Bishop 2006)

Objective function: **maximum margin**

$\min_i [|(s(\mathbf{x}_i)|)]$: Point closest to decision surface

We want to identify β and b to maximize the margin:

$$\arg \max_{\beta, b} \left\{ \frac{1}{\|\beta\|} \min_i [|(s(\mathbf{x}_i)|)] \right\}$$
$$\arg \max_{\beta, b} \left\{ \frac{1}{\|\beta\|} \min_i [|\beta' \mathbf{x}_i + b|] \right\}$$

Support Vector Machines: Algebra (Bishop 2006)

Objective function: **maximum margin**

$\min_i [|(s(\mathbf{x}_i)|)|]$: Point closest to decision surface

We want to identify β and b to maximize the margin:

$$\arg \max_{\beta, b} \left\{ \frac{1}{\|\beta\|} \min_i [|(s(\mathbf{x}_i)|)|] \right\}$$
$$\arg \max_{\beta, b} \left\{ \frac{1}{\|\beta\|} \min_i [|\beta' \mathbf{x}_i + b|] \right\}$$

Constrained optimization problem \rightsquigarrow Quadratic programming problem

What About Overlap? (Bishop 2006)

- Rare that classes are separable.

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

$\xi_i = 0$ if correctly classified

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

$\xi_i = 0$ if correctly classified

$\xi_i = |s(\mathbf{x}_i)|$ if incorrectly classified

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

$\xi_i = 0$ if correctly classified

$\xi_i = |s(\mathbf{x}_i)|$ if incorrectly classified

Tradeoff:

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

$$\xi_i = 0 \text{ if correctly classified}$$

$$\xi_i = |s(\mathbf{x}_i)| \text{ if incorrectly classified}$$

Tradeoff:

- Maximize margin between correctly classified groups

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

$$\xi_i = 0 \text{ if correctly classified}$$

$$\xi_i = |s(\mathbf{x}_i)| \text{ if incorrectly classified}$$

Tradeoff:

- Maximize margin between correctly classified groups
- Minimize error from misclassified documents

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

$$\xi_i = 0 \text{ if correctly classified}$$

$$\xi_i = |s(\mathbf{x}_i)| \text{ if incorrectly classified}$$

Tradeoff:

- Maximize margin between correctly classified groups
- Minimize error from misclassified documents

$$\arg \max_{\beta, b} \left\{ C \sum_{i=1}^N \xi_i + \frac{1}{\|\beta\|} \min_i [|\beta' \mathbf{x}_i + b|] \right\}$$

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

$$\xi_i = 0 \text{ if correctly classified}$$

$$\xi_i = |s(\mathbf{x}_i)| \text{ if incorrectly classified}$$

Tradeoff:

- Maximize margin between correctly classified groups
- Minimize error from misclassified documents

$$\arg \max_{\beta, b} \left\{ C \sum_{i=1}^N \xi_i + \frac{1}{\|\beta\|} \min_i [|\beta' \mathbf{x}_i + b|] \right\}$$

C captures tradeoff

What About Overlap? (Bishop 2006)

- Rare that classes are separable.
- Define:

$$\xi_i = 0 \text{ if correctly classified}$$

$$\xi_i = |s(\mathbf{x}_i)| \text{ if incorrectly classified}$$

Tradeoff:

- Maximize margin between correctly classified groups
- Minimize error from misclassified documents

$$\arg \max_{\beta, b} \left\{ C \sum_{i=1}^N \xi_i + \frac{1}{\|\beta\|} \min_i [|\beta' \mathbf{x}_i + b|] \right\}$$

C captures tradeoff

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?
 - 1) Set up K classification problems:

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?
 - 1) Set up K classification problems:
 - Compare each class to all other classes

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?
 - 1) Set up K classification problems:
 - Compare each class to all other classes
 - Problem: can lead to inconsistent results

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?
 - 1) Set up K classification problems:
 - Compare each class to all other classes
 - Problem: can lead to inconsistent results
 - Solution(?): select category with largest “score”

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?
 - 1) Set up K classification problems:
 - Compare each class to all other classes
 - Problem: can lead to inconsistent results
 - Solution(?): select category with largest “score”
 - Problem: scales are not comparable

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?
 - 1) Set up K classification problems:
 - Compare each class to all other classes
 - Problem: can lead to inconsistent results
 - Solution(?): select category with largest “score”
 - Problem: scales are not comparable
 - 2) Common solution: set up $K(K - 1)/2$ classifications

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?
 - 1) Set up K classification problems:
 - Compare each class to all other classes
 - Problem: can lead to inconsistent results
 - Solution(?): select category with largest “score”
 - Problem: scales are not comparable
 - 2) Common solution: set up $K(K - 1)/2$ classifications
 - Perform vote to select class (still suboptimal)

How to Handle Multiple Comparisons?

- Rare that we only want to classify two categories
- How to handle classification into K groups?
 - 1) Set up K classification problems:
 - Compare each class to all other classes
 - Problem: can lead to inconsistent results
 - Solution(?): select category with largest “score”
 - Problem: scales are not comparable
 - 2) Common solution: set up $K(K - 1)/2$ classifications
 - Perform vote to select class (still suboptimal)
 - 3) Simultaneous estimation possible, much slower

R Code to Run SVMs

```
library(e1071)
fit<- svm(T . , as.data.frame(tdm) , method ='C',
kernel='linear')
```

where: method = 'C' → Classification
kernel='linear' → allows for distortion of feature space. Options:

- Linear
- Polynomial
- Radial
- sigmoid

```
preds<- predict(fit, data =
as.data.frame(tdm[-c(1:no.train),]))
```

Example of SVMs in Political Science Research

Hillard, Purpura, Wilkerson: SVMs to code topic/sub topics for policy agendas project

TABLE 3. Bill Title Interannotator Agreement for Five Model Types

	SVM	MaxEnt	Boostexter	Naïve Bayes
Major topic $N = 20$	88.7% (.881)	86.5% (.859)	85.6% (.849)	81.4% (.805)
Subtopic $N = 226$	81.0% (.800)	78.3% (.771)	73.6% (.722)	71.9% (.705)

SVMs are **under utilized** in political science

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes (and next week, SVM): focused on individual document classification.

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes (and next week, SVM): focused on individual document classification.

But what if we're focused on **proportions only**?

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes (and next week, SVM): focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes (and next week, SVM): focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

Can be much more accurate than individual classifiers, requires fewer assumptions (**do not need random sample of documents**) .

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes (and next week, SVM): focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

Can be much more accurate than individual classifiers, requires fewer assumptions (**do not need random sample of documents**) .

- King and Lu (2008): derive method for characterizing causes of deaths for verbal autopsies

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes (and next week, SVM): focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

Can be much more accurate than individual classifiers, requires fewer assumptions (**do not need random sample of documents**) .

- King and Lu (2008): derive method for characterizing causes of deaths for verbal autopsies
- Hopkins and King (2010): extend the method to text documents

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes (and next week, SVM): focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes
Can be much more accurate than individual classifiers, requires fewer assumptions (**do not need random sample of documents**) .

- King and Lu (2008): derive method for characterizing causes of deaths for verbal autopsies
- Hopkins and King (2010): extend the method to text documents

Basic intuition:

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes (and next week, SVM): focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes
Can be much more accurate than individual classifiers, requires fewer assumptions (**do not need random sample of documents**) .

- King and Lu (2008): derive method for characterizing causes of deaths for verbal autopsies
- Hopkins and King (2010): extend the method to text documents

Basic intuition:

- Examine joint distribution of characteristics (without making Naive Bayes like assumption)
- Focus on distributions (only) makes this analysis possible

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of \mathbf{x}_i ?

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of \mathbf{x}_i ?

- 2^J possible vectors

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of \mathbf{x}_i ?

- 2^J possible vectors

Define:

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of \mathbf{x}_i ?

- 2^J possible vectors

Define:

$$P(\mathbf{x}) = \text{probability of observing } \mathbf{x}$$

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of \mathbf{x}_i ?

- 2^J possible vectors

Define:

$P(\mathbf{x})$ = probability of observing \mathbf{x}

$P(\mathbf{x}|C_j)$ = Probability of observing \mathbf{x} conditional on category C_j

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of \mathbf{x}_i ?

- 2^J possible vectors

Define:

$P(\mathbf{x})$ = probability of observing \mathbf{x}

$P(\mathbf{x}|C_j)$ = Probability of observing \mathbf{x} conditional on category C_j

$P(\mathbf{X}|C)$ = Matrix collecting vectors

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of \mathbf{x}_i ?

- 2^J possible vectors

Define:

$P(\mathbf{x})$ = probability of observing \mathbf{x}

$P(\mathbf{x}|C_j)$ = Probability of observing \mathbf{x} conditional on category C_j

$P(\mathbf{X}|C)$ = Matrix collecting vectors

$P(C)$ = $P(C_1, C_2, \dots, C_K)$ target quantity of interest

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [$(J \times 1)$ vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of \mathbf{x}_i ?

- 2^J possible vectors

Define:

$P(\mathbf{x})$ = probability of observing \mathbf{x}

$P(\mathbf{x}|C_j)$ = Probability of observing \mathbf{x} conditional on category C_j

$P(\mathbf{X}|C)$ = Matrix collecting vectors

$P(C)$ = $P(C_1, C_2, \dots, C_K)$ target quantity of interest

ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

$$\underbrace{P(\mathbf{x})}_{2^J \times 1} = \underbrace{P(\mathbf{x}|C)}_{2^J \times K} \underbrace{P(C)}_{K \times 1}$$

Matrix algebra problem to solve, for $P(C)$

Like Naive Bayes, requires two pieces to estimate

Complication $2^J \gg$ no. documents

Kernel Smoothing Methods (without a formal model)

- $P(\mathbf{x})$ = estimate directly from test set
- $P(\mathbf{x}|C)$ = estimate from training set
 - Key assumption: $P(\mathbf{x}|C)$ in training set is equivalent to $P(\mathbf{x}|C)$ in test set
- If true, can perform biased sampling of documents, worry less about drift...

Algorithm Summarized

- Estimate $\hat{p}(\mathbf{x})$ from test set
- Estimate $\hat{p}(\mathbf{x}|C)$ from training set
- Use $\hat{p}(\mathbf{x})$ and $\hat{p}(\mathbf{x}|C)$ to solve for $p(C)$

Assessing Model Performance

Not classifying individual documents → different standards

Mean Square Error :

$$E[(\hat{\theta} - \theta)^2] = \text{var}(\hat{\theta}) + \text{Bias}(\hat{\theta}, \theta)^2$$

Suppose we have true proportions $P(C)^{\text{true}}$. Then, we'll estimate **Root Mean Square Error**

$$\begin{aligned} \text{RMSE} &= \sqrt{\frac{\sum_{j=1}^J (P(C_j)^{\text{true}} - P(C_j))^2}{J}} \\ \text{Mean Abs. Prediction Error} &= \left| \frac{\sum_{j=1}^J (P(C_j)^{\text{true}} - P(C_j))}{J} \right| \end{aligned}$$

Visualize: plot true and estimated proportions

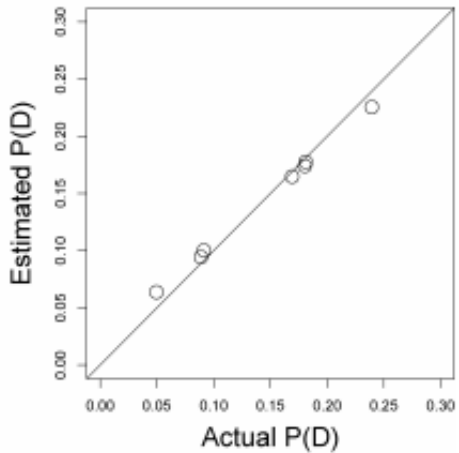


TABLE 1 Performance of Our Nonparametric Approach and Four Support Vector Machine Analyses

Percent of Blog Posts Correctly Classified				
	In-Sample Fit	In-Sample Cross-Validation	Out-of-Sample Prediction	Mean Absolute Proportion Error
Nonparametric	—	—	—	1.2
Linear	67.6	55.2	49.3	7.7
Radial	67.6	54.2	49.1	7.7
Polynomial	99.7	48.9	47.8	5.3
Sigmoid	15.6	15.6	18.2	23.2

Notes: Each row is the optimal choice over numerous individual runs given a specific kernel. Leaving aside the sigmoid kernel, individual classification performance in the first three columns does not correlate with mean absolute error in the document category proportions in the last column.

Using the House Press Release Data

Method	RMSE	APSE
ReadMe	0.036	0.056
NaiveBayes	0.096	0.14
SVM	0.052	0.084

Code to Run in R

Control file:

filename	truth	trainingset
20July2009LEWIS53.txt	4	1
26July2006LEWIS249.txt	2	0

```
tdm<- undergrad(control=control, fullfreq=F)
process<- preprocess(tdm)
output<- undergrad(process)
output$est.CSMF ## proportion in each category
output$true.CSMF ## if labeled for validation set (but not
used in training set)
```

Model Selection!