#### Xuezhi Wang

Computer Science Department Carnegie Mellon University

10701-recitation, Apr 22

・ロト ・聞 と ・ ヨ と ・ ヨ と 。

ъ

### Outline



- Linear Regression
- Regularization
- Probabilistic Intepretation

#### 2 Model Selection

- Variable Selection
- Model selection

ъ

Ridge/Lasso Regression

Model Selection

Linear Regression Regularization Probabilistic Intepretation

# Outline



- Linear Regression
- Regularization
- Probabilistic Intepretation

2 Model Selection

- Variable Selection
- Model selection

(日)

э

э

Linear Regression Regularization Probabilistic Intepretation

### Linear Regression

Data X:  $N \times P$  matrix, Target y:  $N \times 1$  vector

• N samples, each sample has P features

Want to find  $\theta$  so that y and  $X\theta$  are as close as possible

• Pick  $\theta$  that minimizes the cost function

$$L = \frac{1}{2} \sum_{i} (y_i - X_i \theta)^2 = \frac{1}{2} ||y - X\theta||^2$$

use gradient descent

$$heta_j^{t+1} = heta_j^t - step * rac{\partial L}{\partial heta_j} = heta_j^t - step * \sum_i (y_i - X_i heta) (-X_{ij})$$

< < >> < </>

Linear Regression Regularization Probabilistic Intepretation

# Linear Regression

Matrix form:

$$L = \frac{1}{2} \sum_{i} (y_i - X_i \theta)^2 = \frac{1}{2} ||y - X\theta||^2$$
$$= \frac{1}{2} (y - X\theta)^\top (y - X\theta)$$
$$= \frac{1}{2} (y^\top y - y^\top X\theta - \theta^\top X^\top y + \theta^\top X^\top X\theta)$$

Take derivative w.r.t.  $\theta$ 

$$\frac{\partial L}{\partial \theta} = \frac{1}{2} (-2X^{\top}y + 2X^{\top}X\theta) = 0$$

Hence we get

$$\theta = (X^{\top}X)^{-1}X^{\top}y$$

ヘロト 人間 とくほとく ほとう

æ

Linear Regression Regularization Probabilistic Intepretation

# Linear Regression

Comparison of iterative methods and matrix methods:

- matrix methods achieve solution in a single step, but can be infeasible for real-time data, or large amount of data.
- iterative methods can be used in large practical problems, but need to decide learning rate

Any problems?

- Data X is an  $N \times P$  matrix
- Usually N > P, i.e., number of data points larger than feature dimensions. And usually X is of full column rank.
- Under this case  $X^{\top}X$  have rank *P*, i.e., invertible
- What if X has less than full column rank?

イロト イポト イヨト イヨト

Linear Regression Regularization Probabilistic Intepretation

# Outline



- Linear Regression
- Regularization
- Probabilistic Intepretation

#### 2 Model Selection

- Variable Selection
- Model selection

ヘロト ヘアト ヘヨト ヘ

프 🕨 🗉 프

Linear Regression Regularization Probabilistic Intepretation

# Regularization: $\ell_2$ norm

Ridge Regression:

$$\min_{\theta} \frac{1}{2} \sum_{i} (y_i - X_i \theta)^2 + \lambda ||\theta||_2^2$$

Solution is given by:

$$\theta = (\boldsymbol{X}^{\top}\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}^{\top}\boldsymbol{y}$$

- Results in a solution with small  $\theta$
- Solves the problem that  $X^{\top}X$  is not invertible

イロト イヨト イヨト イ

Linear Regression Regularization Probabilistic Intepretation

# Regularization: $\ell_1$ norm

Lasso Regression:

$$\min_{\theta} \frac{1}{2} \sum_{i} (y_i - X_i \theta)^2 + \lambda ||\theta||_1$$

Solution is given by taking subgradient:

$$\sum_{i} (\mathbf{y}_i - \mathbf{X}_i \theta) (-\mathbf{X}_{ij}) + \lambda t_j$$

where  $t_i$  is the subgradient of  $\ell_1$  norm,

$$t_j = sign(\theta_j)$$
 if  $\theta_j \neq 0, t_j \in [-1, 1]$  otherwise

- Sparse solution, i.e., θ will be a vector with more zero coordinates.
- Good for high-dimensional problems

★ E → ★ E →

Linear Regression Regularization Probabilistic Intepretation

### Solving Lasso regression

Efron et al. proposed LARS (least angle regression) which computes the LASSO path efficiently

#### Forward stagewise algorithm

- Assume X is standardized and y is centered
- choose small  $\epsilon$ 
  - Start with initial residual r = y, and  $\theta_1 = ... = \theta_P = 0$
  - Find the predictor Z<sub>j</sub> (*j*th column of X) most correlated with r
  - Update  $\theta_j \leftarrow \theta_j + \delta_j$ , where  $\delta_j = \epsilon \cdot \operatorname{sign}(Z_i^{\top} r)$
  - Set  $r \leftarrow r \delta_j Z_j$ , repeat steps 2 and 3

ヘロト 人間 とくほとく ほとう

Regularization

Comparison of Ridge and Lasso regression:

# Two-dimensional case: contour plots for d = 2 $\sum_{i=1}^{\infty} (y_i - f(x_i, \mathbf{w}))^2$ $w_2$ $\mathbf{w}^{\star}$ w. $\lambda \|\mathbf{w}\|^2$ $\lambda \sum |w_j|$ lasso

ridge regression

 $\tilde{w}_1$ 

Linear Regression Regularization Probabilistic Intepretation

# Comparison of Ridge and Lasso regression:

Higher dimensional case:



イロト 不得 とくほと くほとう

ъ

Ridge/Lasso Regression

Model Selection

Linear Regression Regularization Probabilistic Intepretation

# Choosing $\lambda$



Standard practice now is to use cross-validation

イロト 不得 とくほ とくほとう

æ

Linear Regression Regularization Probabilistic Intepretation

### Outline



- Linear Regression
- Regularization
- Probabilistic Intepretation

2 Model Selection

- Variable Selection
- Model selection

イロト イポト イヨト イヨト

э

### Probabilistic Intepretation of Linear regression

Assume  $y_i = X_i \theta + \epsilon_i$ , where  $\epsilon$  is the random noise. Assume  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ 

$$p(y_i|X_i;\theta) = rac{1}{\sqrt{2\pi}\sigma} \exp\{-rac{(y_i - X_i\theta)^2}{2\sigma^2}\}$$

Since data points are i.i.d, we have the data likelihood

$$L(\theta) = \prod_{i=1}^{N} p(y_i | X_i; \theta) \propto \exp\{-\frac{\sum_{i=1}^{N} (y_i - X_i \theta)^2}{2\sigma^2}\}$$

The log likelihood is:

$$\ell( heta) = -rac{\sum_{i=1}^{N}(y_i - X_i heta)^2}{2\sigma^2} + ext{const}$$

Maximizing the log-likelihood is equivalent to minimize  $\sum_{i=1}^{N} (y_i - X_i \theta)^2$ , i.e., the loss function in LR!

Ridge/Lasso Regression Model Selection Probabilistic Intepretation

#### Probabilistic Intepretation of Ridge regression

Assume a Gaussian prior on  $\theta \sim \mathcal{N}(\mathbf{0}, \tau^2 I)$ , i.e.,

$$p(\theta) \propto \exp\{-\theta^{\top}\theta/2\tau^2\}$$

Now get the MAP estimate of  $\theta$ 

$$p(\theta|X, y) \propto p(y|X; \theta) p(\theta) = \exp\{-\frac{\sum_{i=1}^{N} (y_i - X_i \theta)^2}{2\sigma^2}\} \exp\{-\theta^{\top} \theta / 2\tau^2\}$$

The log likelihood is:

$$\ell(\theta|X, y) = -\frac{\sum_{i=1}^{N} (y_i - X_i \theta)^2}{2\sigma^2} - \theta^{\top} \theta / 2\tau^2 + \text{const}$$

which matches  $\min_{\theta} \frac{1}{2} \sum_{i} (y_i - X_i \theta)^2 + \lambda ||\theta||_2^2$ , where  $\lambda$  is a constant associated with  $\sigma^2, \tau^2$ .

Ridge/Lasso Regression Model Selection Probabilistic Intepretation

### Probabilistic Intepretation of Lasso regression

Assume a Laplace prior on  $\theta_i \stackrel{iid}{\sim} Laplace(0, t)$ , i.e.,

 $p(\theta_i) \propto \exp\{-|\theta_i|/t\}$ 

Now get the MAP estimate of  $\theta$ 

$$p(\theta|X, y) \propto p(y|X; \theta) p(\theta) = \exp\{-\frac{\sum_{i=1}^{N} (y_i - X_i \theta)^2}{2\sigma^2}\} \exp\{-\sum_i |\theta_i|/t\}$$

The log likelihood is:

$$\ell(\theta|X, y) = -\frac{\sum_{i=1}^{N} (y_i - X_i \theta)^2}{2\sigma^2} - \sum_i |\theta_i|/t + \text{const}$$

which matches  $\min_{\theta} \frac{1}{2} \sum_{i} (y_i - X_i \theta)^2 + \lambda ||\theta||_1$ , where  $\lambda$  is a constant associated with  $\sigma^2$ , *t*.

Variable Selection Model selection

# Outline

#### Ridge/Lasso Regression

- Linear Regression
- Regularization
- Probabilistic Intepretation

#### 2 Model Selection

- Variable Selection
- Model selection

<ロト < 回 > < 回 > .

프 🕨 🗆 프

### Variable Selection

- Consider "best" subsets, order O(2<sup>P</sup>) (combinatorial explosion)
- Stepwise selection
  - A new variable may be added into the model even with a small improvement in LMS
  - When applying stepwise to a perturbation of the data, probably have different set of variables enter into the model at each stage
- LASSO produces sparse solutions, which takes care of model selection
  - we can even see when variables jump into the model by looking at the LASSO path

ヘロト ヘ戸ト ヘヨト ヘヨト

Variable Selection Model selection

### Outline

#### Ridge/Lasso Regression

- Linear Regression
- Regularization
- Probabilistic Intepretation

#### 2 Model Selection

- Variable Selection
- Model selection

イロト イヨト イヨト イ

프 🕨 🗉 프



Suppose you have data  $Y_1, ..., Y_n$  and you want to model the distribution of *Y*. Some popular models are:

- the Exponential distribution:  $f(y; \theta) = \theta e^{-\theta y}$
- the Gaussian distribution:  $f(y; u, \sigma^2) \sim \mathcal{N}(u, \sigma^2)$

• ...

How do you know which model is better?

イロト イポト イヨト イヨト

э.

Suppose we have models  $M_1, ..., M_k$  where each model is a set of densities:

$$M_j = \{ p(y; heta_j) : heta_j \in \Theta_j \}$$

We have data  $Y_1, ..., Y_n$  drawn from some density *f* (not necessarily drawn from these models). Define

$$\mathsf{AIC}(j) = \ell_j(\hat{ heta}_j) - 2d_j$$

where  $\ell_j(\theta_j)$  is the log-likelihood, and  $\hat{\theta}_j$  is the parameter that maximizes the log-likelihood.  $d_j$  is the dimension of  $\Theta_j$ .

イロト イヨト イヨト イ

#### Bayesian Information Criterion We choose *j* to maximize

$$\mathsf{BIC}_j = \ell_j(\hat{\theta}_j) - \frac{d_j}{2}\log n$$

which is similar to AIC but the penalty is harsher, hence BIC tends to choose simpler models.

イロン 不得 とくほ とくほとう

ъ

Variable Selection Model selection

### Simple example

Let

$$Y_1, ..., Y_n \sim \mathcal{N}(\mu, 1)$$

we want to compare two model:

 $M_0 : \mathcal{N}(0, 1) \text{ and } M_1 : \mathcal{N}(u, 1)$ 

<ロト <回 > < 注 > < 注 > 、

3

Variable Selection Model selection

### Simple example: AIC

The log-likelihood is

$$\ell = \log \prod_{i} e^{-(Y_{i}-u)^{2}/2} = -\sum_{i} (Y_{i}-u)^{2}/2$$
$$AIC_{0} = -\sum_{i} Y_{i}^{2}/2 - 0$$
$$AIC_{1} = -\sum_{i} (Y_{i}-\bar{Y})^{2}/2 - 2 = -\sum_{i} Y_{i}^{2}/2 + \frac{n}{2}\bar{Y}^{2} - 2$$

we choose model 1 if  $AIC_1 > AIC_0$  i.e.,

$$-\sum_{i}Y_{i}^{2}/2+rac{n}{2}ar{Y}^{2}-2>-\sum_{i}Y_{i}^{2}/2$$

or  $\bar{Y} > \sqrt{\frac{4}{n}}$ .

ヘロト 人間 とくほとくほとう

ъ

Variable Selection Model selection

#### Simple example: BIC

$$BIC_{0} = -\sum_{i} \frac{Y_{i}^{2}}{2} - \frac{0}{2} \log n = -\sum_{i} \frac{Y_{i}^{2}}{2}$$
$$AIC_{1} = -\sum_{i} \frac{(Y_{i} - \bar{Y})^{2}}{2} - \frac{1}{2} \log n = -\sum_{i} \frac{Y_{i}^{2}}{2} + \frac{n}{2} \bar{Y}^{2} - \frac{1}{2} \log n$$

we choose model 1 if  $BIC_1 > BIC_0$  i.e.,

$$-\sum_{i} Y_{i}^{2}/2 + n/2\bar{Y}^{2} - \frac{1}{2}\log n > -\sum_{i} Y_{i}^{2}/2$$

or  $\bar{Y} > \sqrt{\frac{\log n}{n}}$ .

# Comparison

Generally speaking,

- AIC/CV finds the most predictive model
- BIC find the true model with high probability, i.e., BIC assumes that one of the models is true and that you are trying to find the model most likely to be true in the Bayesian sense.

イロト イポト イヨト イヨト