

Introduction to Machine Learning

Feature Selection

Andres Mendez-Vazquez

January 26, 2023

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
- Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Outline

1

Introduction

● Feature Engineering

- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Why Feature Engineering?

As always we love simple linear models

- Easy to analyze
- Unique solution

Why Feature Engineering?

As always we love simple linear models

- Easy to analyze
- Unique solution

Definition

- Feature engineering (or feature extraction) is the process of using domain knowledge to extract features (characteristics, properties, attributes) from raw data.

Therefore

We have several attempts for this

- Feature Selection
 - ▶ Shrinkage Methods
- Feature Generation
 - ▶ Fisher Linear Discriminant
 - ▶ Principal Component Analysis
 - ▶ Singular Value Decomposition
 - ▶ Autoencoders

Basically Feature Engineering

Feature Selection

- Selection of a compact set to avoid the peaking phenomena

Basically Feature Engineering

Feature Selection

- Selection of a compact set to avoid the peaking phenomena

Feature Generation

- Generate Richer Features
 - ▶ And select the best ones

Outline

1

Introduction

- Feature Engineering
- **What is Feature Selection?**
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
- Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

What is this?

Main Question

“Given a number of features, how can one select the most important of them so as to reduce their number and at the same time retain as much as possible of their class discriminatory information? “

What is this?

Main Question

“Given a number of features, how can one select the most important of them so as to reduce their number and at the same time retain as much as possible of their class discriminatory information? “

Why is important?

- 1 If we selected features with little discrimination power, the subsequent design of a classifier would lead to poor performance.

What is this?

Main Question

“Given a number of features, how can one select the most important of them so as to reduce their number and at the same time retain as much as possible of their class discriminatory information? “

Why is important?

- 1 If we selected features with little discrimination power, the subsequent design of a classifier would lead to poor performance.
- 2 if information-rich features are selected, the design of the classifier can be greatly simplified.

What is this?

Main Question

“Given a number of features, how can one select the most important of them so as to reduce their number and at the same time retain as much as possible of their class discriminatory information? “

Why is important?

- 1 If we selected features with little discrimination power, the subsequent design of a classifier would lead to poor performance.
- 2 if information-rich features are selected, the design of the classifier can be greatly simplified.

Therefore

We want features that lead to

What is this?

Main Question

“Given a number of features, how can one select the most important of them so as to reduce their number and at the same time retain as much as possible of their class discriminatory information? “

Why is important?

- 1 If we selected features with little discrimination power, the subsequent design of a classifier would lead to poor performance.
- 2 if information-rich features are selected, the design of the classifier can be greatly simplified.

Therefore

We want features that lead to

- 1 Large between-class distance.

What is this?

Main Question

“Given a number of features, how can one select the most important of them so as to reduce their number and at the same time retain as much as possible of their class discriminatory information? “

Why is important?

- 1 If we selected features with little discrimination power, the subsequent design of a classifier would lead to poor performance.
- 2 if information-rich features are selected, the design of the classifier can be greatly simplified.

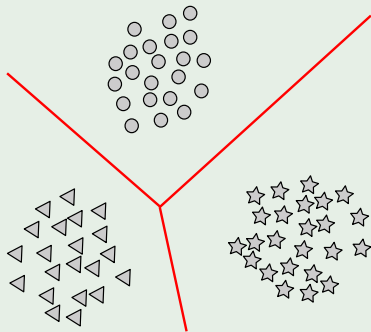
Therefore

We want features that lead to

- 1 Large between-class distance.
- 2 Small within-class variance.

Then

Basically, we want nice separated and dense clusters!!!



Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?

● Preprocessing

- Outlier Removal
- Finding Multivariate Outliers
- Data Normalization
- Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

However, Before That...

It is necessary to do the following

- 1 Outlier removal.

However, Before That...

It is necessary to do the following

- ① Outlier removal.
- ② Data normalization.

However, Before That...

It is necessary to do the following

- ① Outlier removal.
- ② Data normalization.
- ③ Deal with missing data.

However, Before That...

It is necessary to do the following

- ➊ Outlier removal.
- ➋ Data normalization.
- ➌ Deal with missing data.

Actually

PREPROCESSING!!!

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?

● Preprocessing

- Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Outliers

Definition

An outlier is defined as a point that lies very far from the mean of the corresponding random variable.

Outliers

Definition

An outlier is defined as a point that lies very far from the mean of the corresponding random variable.

Note: We use the standard deviation

Outliers

Definition

An outlier is defined as a point that lies very far from the mean of the corresponding random variable.

Note: We use the standard deviation

Example

For a normally distributed random

Outliers

Definition

An outlier is defined as a point that lies very far from the mean of the corresponding random variable.

Note: We use the standard deviation

Example

For a normally distributed random

- 1 A distance of two times the standard deviation covers 95% of the points.

Outliers

Definition

An outlier is defined as a point that lies very far from the mean of the corresponding random variable.

Note: We use the standard deviation

Example

For a normally distributed random

- 1 A distance of two times the standard deviation covers 95% of the points.
- 2 A distance of three times the standard deviation covers 99% of the points.

Outliers

Definition

An outlier is defined as a point that lies very far from the mean of the corresponding random variable.

Note: We use the standard deviation

Example

For a normally distributed random

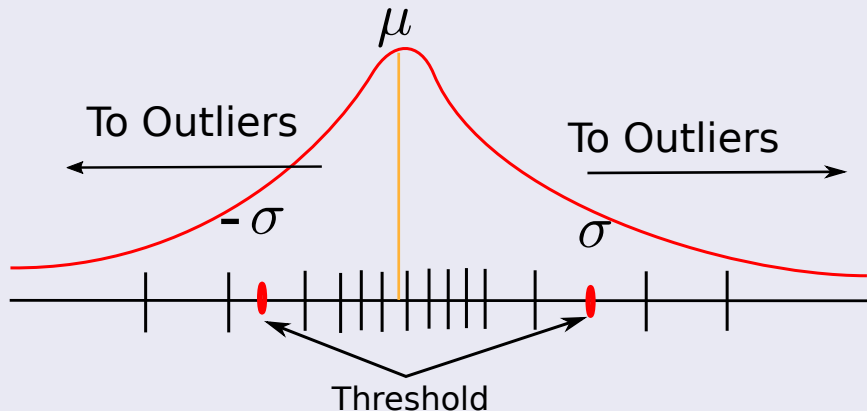
- ➊ A distance of two times the standard deviation covers 95% of the points.
- ➋ A distance of three times the standard deviation covers 99% of the points.

Note

Points with values very different from the mean value produce large errors during training and may have disastrous effects. These effects are even worse when the outliers, and they are the result of noisy measurements.

For example, we can use the standard deviation

For a set of samples $x_1, x_2, x_3, \dots \in \mathbb{R}$



Outlier Removal

Important

Then removing outliers is the biggest importance.

Outlier Removal

Important

Then removing outliers is the biggest importance.

Therefore

You can do the following

Outlier Removal

Important

Then removing outliers is the biggest importance.

Therefore

You can do the following

- 1 If you have a small number \Rightarrow discard them!!!

Outlier Removal

Important

Then removing outliers is the biggest importance.

Therefore

You can do the following

- ① If you have a small number \Rightarrow discard them!!!
- ② Adopt cost functions that are not sensitive to outliers:

Outlier Removal

Important

Then removing outliers is the biggest importance.

Therefore

You can do the following

- 1 If you have a small number \Rightarrow discard them!!!
- 2 Adopt cost functions that are not sensitive to outliers:
- 3 For more techniques
 - 1 Huber, P.J. "Robust Statistics," JohnWiley and Sons, 2nd Ed 2009.

An improvement over SD

We can do the following

- Estimate the “middle” of the data, the sample median

An improvement over SD

We can do the following

- Estimate the “middle” of the data, the sample median

Using order statistics for the samples $x_{(1)} \leq x_{(2)} \leq x_{(3)} \leq \dots \leq x_{(n)}$

$$\text{Med}(\mathbf{x}) = \begin{cases} x_{(m)} & n \text{ is odd for } n = 2m - 1 \\ \frac{x_{(m)} + x_{(m+1)}}{2} & n \text{ is even} \end{cases}$$

Then, it is possible

To define a MAD estimator

$$MADN = \frac{Med \{ |x - Med(x)| \}}{0.6745}$$

Then, it is possible

To define a MAD estimator

$$MADN = \frac{Med\{|x - Med(x)|\}}{0.6745}$$

This can be compared with the standard deviation

- The scale constant (approximately 0.6745) is the inverse of the standard normal distribution function evaluated at 3/4.

Something quite interesting

If we use the thresholding to eliminate using SD vs *MADN*

- We have that MADN for example, when deleting large outliers, you go from $0.53 \rightarrow 0.50$

Something quite interesting

If we use the thresholding to eliminate using SD vs *MADN*

- We have that MADN for example, when deleting large outliers, you go from 0.53 \rightarrow 0.50

Instead with SD you have

- 5.30 to 0.69... bad!!!
- MADN is not influenced very much by the presence of a large outlier

Something quite interesting

If we use the thresholding to eliminate using SD vs *MADN*

- We have that MADN for example, when deleting large outliers, you go from 0.53 \rightarrow 0.50

Instead with SD you have

- 5.30 to 0.69... bad!!!
- MADN is not influenced very much by the presence of a large outlier

Why not always use the median and MAD?

- These estimates have statistical poorer performance when outliers do not exist.

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?

● Preprocessing

- Outlier Removal
- **Finding Multivariate Outliers**
- Data Normalization
- Methods
- Missing Data
- Using EM
- Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

We can do the following

Algorithm

Input: An $N \times d$ data set $Data$

Output: Candidate Outliers

We can do the following

Algorithm

Input: An $N \times d$ data set $Data$

Output: Candidate Outliers

- 1 Calculate the sample mean μ and sample covariance matrix Σ .

We can do the following

Algorithm

Input: An $N \times d$ data set $Data$

Output: Candidate Outliers

- 1 Calculate the sample mean μ and sample covariance matrix Σ .
- 2 Let M be $N \times 1$ vector consisting of square of the Mahalanobis distance to μ .

We can do the following

Algorithm

Input: An $N \times d$ data set $Data$

Output: Candidate Outliers

- 1 Calculate the sample mean μ and sample covariance matrix Σ .
- 2 Let M be $N \times 1$ vector consisting of square of the Mahalanobis distance to μ .
- 3 Find points O in M whose values are greater than

We can do the following

Algorithm

Input: An $N \times d$ data set $Data$

Output: Candidate Outliers

- 1 Calculate the sample mean μ and sample covariance matrix Σ .
- 2 Let M be $N \times 1$ vector consisting of square of the Mahalanobis distance to μ .
- 3 Find points O in M whose values are greater than

$$\chi_d^2(0.05)$$

We can do the following

Algorithm

Input: An $N \times d$ data set $Data$

Output: Candidate Outliers

- 1 Calculate the sample mean μ and sample covariance matrix Σ .
- 2 Let M be $N \times 1$ vector consisting of square of the Mahalanobis distance to μ .
- 3 Find points O in M whose values are greater than

$$\chi_d^2(0.05)$$

- 4 Return O .

How?

Get the Sample Mean per feature k

$$m_i = \frac{1}{N} \sum_{k=1}^N x_{ki}$$

How?

Get the Sample Mean per feature k

$$\mathbf{m}_i = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_{ki}$$

Get the Sample Variance per feature k

$$v_i = \frac{1}{N-1} \sum_{k=1}^N (\mathbf{x}_{ki} - \mathbf{m}_i) (\mathbf{x}_{ki} - \mathbf{m}_i)^T$$

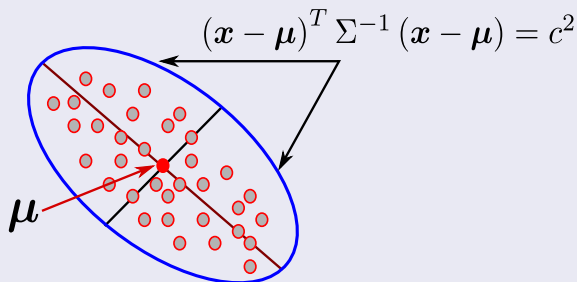
Mahalanobis Distance

We have

$$M(x) = \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

Thus

Setting $M(x)$ to a constant c defines a multidimensional ellipsoid with centroid at μ



As Johnson and Wichern (2007, p. 155, Eq. 4-8) state

The solid ellipsoid of \mathbf{x} vectors satisfying

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq \chi_d^2(\alpha)$$

has a probability $1 - \alpha$.

How?

We know that

χ_d^2 is defined as the distribution of the sum $\sum_{i=1}^d Z_i^2$ where Z_i 's are independent $N(0, 1)$ random variables.

How?

We know that

χ_d^2 is defined as the distribution of the sum $\sum_{i=1}^d Z_i^2$ where Z_i 's are independent $N(0, 1)$ random variables.

Additionally, if we assume that Σ is positive definite and $\Sigma \in \mathbb{R}^{d \times d}$

$$\Sigma = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$

- 1 \mathbf{u}_i are the orthonormal eigenvectors of Σ
- 2 λ_i are the corresponding real eigenvalues

Then

Something Notable

$$\Sigma^{-1} = \sum_{i=1}^d \frac{1}{\lambda} \mathbf{u}_i \mathbf{u}_i^T$$

Then

Something Notable

$$\Sigma^{-1} = \sum_{i=1}^d \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$$

Now, if our data matrix element $X \sim N_d(\boldsymbol{\mu}, \Sigma)$

We have

$$\Sigma^{-1} \mathbf{u}_i = \frac{1}{\lambda_i} \mathbf{u}_i$$

Therefore

We have that

$$(X - \mu)^T \Sigma^{-1} (X - \mu) = \sum_{i=1}^d \frac{1}{\lambda_i} (X - \mu)^T \mathbf{u}_i \mathbf{u}_i^T (X - \mu)$$

Therefore

We have that

$$(X - \mu)^T \Sigma^{-1} (X - \mu) = \sum_{i=1}^d \frac{1}{\lambda_i} (X - \mu)^T \mathbf{u}_i \mathbf{u}_i^T (X - \mu)$$

Then

$$(X - \mu)^T \Sigma^{-1} (X - \mu) = \sum_{i=1}^d \left[\frac{1}{\sqrt{\lambda_i}} \mathbf{u}_i^T (X - \mu) \right]^2 = \sum_{i=1}^d Z_i^2$$

Therefore

If we define

$$\mathbf{Z} = \begin{pmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_d \end{pmatrix}, A_{d \times d} = \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} \mathbf{u}_1^T \\ \frac{1}{\sqrt{\lambda_2}} \mathbf{u}_2^T \\ \vdots \\ \frac{1}{\sqrt{\lambda_d}} \mathbf{u}_d^T \end{pmatrix}$$

Therefore

If we define

$$\mathbf{Z} = \begin{pmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_d \end{pmatrix}, A_{d \times d} = \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} \mathbf{u}_1^T \\ \frac{1}{\sqrt{\lambda_2}} \mathbf{u}_2^T \\ \vdots \\ \frac{1}{\sqrt{\lambda_d}} \mathbf{u}_d^T \end{pmatrix}$$

We know that $(X - \boldsymbol{\mu}) \sim N_d(0, \Sigma)$

- Then, we have $\mathbf{Z} = A(X - \boldsymbol{\mu}) \sim N_d(0, A\Sigma A^T)$

Therefore

Something Notable

$$A \Sigma A^T = \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} \mathbf{u}_1^T \\ \frac{1}{\sqrt{\lambda_2}} \mathbf{u}_2^T \\ \vdots \\ \frac{1}{\sqrt{\lambda_d}} \mathbf{u}_d^T \end{pmatrix} \left[\sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^T \right] \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} \mathbf{u}_1 & \frac{1}{\sqrt{\lambda_2}} \mathbf{u}_2 & \cdots & \frac{1}{\sqrt{\lambda_d}} \mathbf{u}_d \end{pmatrix}$$

Therefore

Something Notable

$$A \Sigma A^T = \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} \mathbf{u}_1^T \\ \frac{1}{\sqrt{\lambda_2}} \mathbf{u}_2^T \\ \vdots \\ \frac{1}{\sqrt{\lambda_d}} \mathbf{u}_d^T \end{pmatrix} \left[\sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^T \right] \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} \mathbf{u}_1 & \frac{1}{\sqrt{\lambda_2}} \mathbf{u}_2 & \cdots & \frac{1}{\sqrt{\lambda_d}} \mathbf{u}_d \end{pmatrix}$$

Therefore

$$A \Sigma A^T = \begin{pmatrix} \sqrt{\lambda_1} \mathbf{u}_1^T \\ \sqrt{\lambda_2} \mathbf{u}_2^T \\ \vdots \\ \sqrt{\lambda_d} \mathbf{u}_d^T \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} \mathbf{u}_1 & \frac{1}{\sqrt{\lambda_2}} \mathbf{u}_2 & \cdots & \frac{1}{\sqrt{\lambda_d}} \mathbf{u}_d \end{pmatrix} = I$$

Therefore

We have that Z_1, Z_2, \dots, Z_d are independent standard normal variables

- $(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$ has a χ_d^2 -distribution.

Therefore

We have that Z_1, Z_2, \dots, Z_d are independent standard normal variables

- $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$ has a χ_d^2 -distribution.

Finally, the $P \left((\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq c^2 \right)$

- It is the probability assigned to the ellipsoid $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq c^2$ by the density $N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Therefore

We have $P\left((\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq \chi_d^2(\alpha)\right) = 1 - \alpha$

Basically $\chi_d^2(\alpha)$ is the the critical chi-square value that makes possible the probability $1 - \alpha$

Therefore

We have $P\left((\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq \chi_d^2(\alpha)\right) = 1 - \alpha$

Basically $\chi_d^2(\alpha)$ is the the critical chi-square value that makes possible the probability $1 - \alpha$

Basically

- We assume that if $1 - \alpha = .95$ is the data with probability of not being an outlier!!!

Algorithm

The Partial Code

```
def OutlierRemoval(self, Data):  
    SampleMean = Data.mean(1)  
    SampleCov = Data - SampleMean  
    SampleCov = np.cov(SampleCov.T)  
    Mahalonobis = (Data - SampleMean)*  
                  np.inv(SampleCov)*  
                  ((Data - SampleMean).T)  
  
    # Something else here  
    # Here you can use chi2.isf(\alpha,dim)
```

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?

● Preprocessing

- Outlier Removal
- Finding Multivariate Outliers
- **Data Normalization**
- Methods
- Missing Data
- Using EM
- Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Data Normalization

In the real world

- In many practical situations a designer is confronted with features whose values lie within different dynamic ranges.

Data Normalization

In the real world

- In many practical situations a designer is confronted with features whose values lie within different dynamic ranges.

For Example

- We can have two features with the following ranges

$$x_i \in [0, 100,000]$$

$$x_j \in [0, 0.5]$$

Data Normalization

In the real world

- In many practical situations a designer is confronted with features whose values lie within different dynamic ranges.

For Example

- We can have two features with the following ranges

$$x_i \in [0, 100,000]$$

$$x_j \in [0, 0.5]$$

Thus

- Many classification machines will be swamped by the first feature!!!

Data Normalization

We have the following situation

- Features with large values may have a larger influence in the cost function than features with small values.

Data Normalization

We have the following situation

- Features with large values may have a larger influence in the cost function than features with small values.

Data Normalization

We have the following situation

- Features with large values may have a larger influence in the cost function than features with small values.

Thus!!!

- This does not necessarily reflect their respective significance in the design of the classifier.

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?

● Preprocessing

- Outlier Removal
- Finding Multivariate Outliers
- Data Normalization

● Methods

- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Min-Max Method

Be Naive

- For each feature $i = 1, \dots, d$ obtain the \max_i and the \min_i such that

$$\hat{x}_{ik} = \frac{x_{ik} - \min_i}{\max_i - \min_i} \quad (1)$$

Min-Max Method

Be Naive

- For each feature $i = 1, \dots, d$ obtain the \max_i and the \min_i such that

$$\hat{x}_{ik} = \frac{x_{ik} - \min_i}{\max_i - \min_i} \quad (1)$$

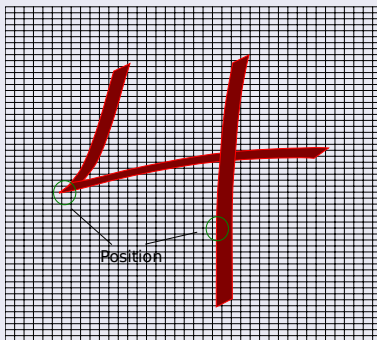
Problem

- This simple normalization will send everything to a unitary sphere!!!
 - ▶ However, it works for certain type of data in Deep Learning

However

Even though this can happen there have been reports that it can work...

- When data does not depend on single values as:



Gaussian Method

Use the idea of

Everything is Gaussian...

Gaussian Method

Use the idea of

Everything is Gaussian...

Thus

- For each feature set...

Gaussian Method

Use the idea of

Everything is Gaussian...

Thus

- For each feature set...

$$\textcircled{1} \quad \bar{x}_k = \frac{1}{N} \sum_{i=1}^N x_{ik}, \quad k = 1, 2, \dots, d$$

Gaussian Method

Use the idea of

Everything is Gaussian...

Thus

- For each feature set...

① $\bar{x}_k = \frac{1}{N} \sum_{i=1}^N x_{ik}, \quad k = 1, 2, \dots, d$

② $\sigma_k^2 = \frac{1}{N-1} \sum_{i=1}^N (x_{ik} - \bar{x}_k)^2, \quad k = 1, 2, \dots, d$

Gaussian Method

Use the idea of

Everything is Gaussian...

Thus

- For each feature set...

$$\textcircled{1} \quad \bar{x}_k = \frac{1}{N} \sum_{i=1}^N x_{ik}, \quad k = 1, 2, \dots, d$$

$$\textcircled{2} \quad \sigma_k^2 = \frac{1}{N-1} \sum_{i=1}^N (x_{ik} - \bar{x}_k)^2, \quad k = 1, 2, \dots, d$$

Thus

$$\hat{x}_{ik} = \frac{x_{ik} - \bar{x}_k}{\sigma} \quad (2)$$

Gaussian Method

Thus

- All new features have zero mean and unit variance.

Gaussian Method

Thus

- All new features have zero mean and unit variance.

Further

- Other linear techniques limit the feature values in the range of $[0, 1]$ or $[-1, 1]$ by proper scaling.

Gaussian Method

Thus

- All new features have zero mean and unit variance.

Further

- Other linear techniques limit the feature values in the range of $[0, 1]$ or $[-1, 1]$ by proper scaling.

However

- We can non-linear mapping. For example the softmax scaling.

Soft Max Scaling

Softmax Scaling

- It consists of two steps

Soft Max Scaling

Softmax Scaling

- It consists of two steps

First one

$$y_{ik} = \frac{x_{ik} - \bar{x}_k}{\sigma} \quad (3)$$

Soft Max Scaling

Softmax Scaling

- It consists of two steps

First one

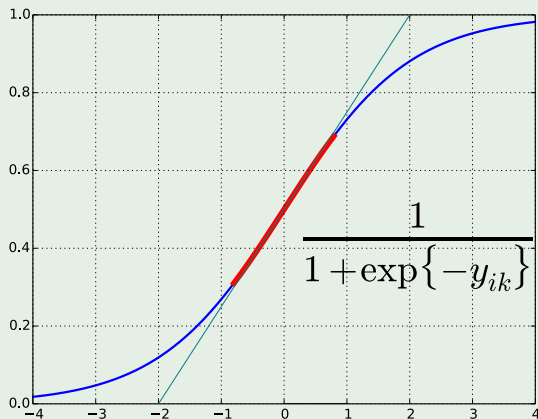
$$y_{ik} = \frac{x_{ik} - \bar{x}_k}{\sigma} \quad (3)$$

Second one

$$\hat{x}_{ik} = \frac{1}{1 + \exp \{-y_{ik}\}} \quad (4)$$

Explanation

Notice the red area is almost flat!!!



Actually

Thus, we have that

- The red region represents values of y inside of the region defined by the mean and variance (small values of y).
- Then, if we have those values x behaves as a linear function.

Actually

Thus, we have that

- The red region represents values of y inside of the region defined by the mean and variance (small values of y).
- Then, if we have those values x behaves as a linear function.

And values too away from the mean

- They are squashed by the exponential part of the function.

If you want a more complex analysis

A more complex analysis

- You can use a Taylor's expansion

$$x = f(y) = f(a) + f'(y)(y - a) + \frac{f''(y)(y - a)^2}{2} + \dots \quad (5)$$

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- **Missing Data**
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Missing Data

This can happen

In practice, certain features may be missing from some feature vectors.

Missing Data

This can happen

In practice, certain features may be missing from some feature vectors.

Examples where this happens

- 1 Social sciences - incomplete surveys.

Missing Data

This can happen

In practice, certain features may be missing from some feature vectors.

Examples where this happens

- 1 Social sciences - incomplete surveys.
- 2 Remote sensing - sensors go off-line.

Missing Data

This can happen

In practice, certain features may be missing from some feature vectors.

Examples where this happens

- 1 Social sciences - incomplete surveys.
- 2 Remote sensing - sensors go off-line.
- 3 etc.

Missing Data

This can happen

In practice, certain features may be missing from some feature vectors.

Examples where this happens

- 1 Social sciences - incomplete surveys.
- 2 Remote sensing - sensors go off-line.
- 3 etc.

Note

Completing the missing values in a set of data is also known as imputation.

Some traditional techniques to solve this problem

Use zeros and risked it!!!

The idea is not to add anything to the features

Some traditional techniques to solve this problem

Use zeros and risked it!!!

The idea is not to add anything to the features

The sample mean/unconditional mean

Does not matter what distribution you have use the sample mean

$$\bar{x}_i = \frac{1}{N} \sum_{k=1}^N x_{ik} \quad (6)$$

Some traditional techniques to solve this problem

Use zeros and risked it!!!

The idea is not to add anything to the features

The sample mean/unconditional mean

Does not matter what distribution you have use the sample mean

$$\bar{x}_i = \frac{1}{N} \sum_{k=1}^N x_{ik} \quad (6)$$

Find the distribution of your data

Use the mean from that distribution. For example, if you have a beta distribution

$$\bar{x}_i = \frac{\alpha}{\alpha + \beta} \quad (7)$$

The MOST traditional

Drop it

- Remove that data
 - ▶ Still you need to have a lot of data to have this luxury

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- **Missing Data**
 - **Using EM**
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Something more advanced

Split data samples in two set of variables

$$\mathbf{x}_{complete} = \begin{pmatrix} \mathbf{x}_{observed} \\ \mathbf{x}_{missed} \end{pmatrix} \quad (8)$$

Something more advanced

Split data samples in two set of variables

$$\mathbf{x}_{complete} = \begin{pmatrix} \mathbf{x}_{observed} \\ \mathbf{x}_{missed} \end{pmatrix} \quad (8)$$

Generate the following probability distribution

$$P(\mathbf{x}_{missed}|\mathbf{x}_{observed}, \Theta) = \frac{P(\mathbf{x}_{missed}, \mathbf{x}_{observed}|\Theta)}{P(\mathbf{x}_{observed}|\Theta)} \quad (9)$$

Something more advanced

Split data samples in two set of variables

$$\mathbf{x}_{complete} = \begin{pmatrix} \mathbf{x}_{observed} \\ \mathbf{x}_{missed} \end{pmatrix} \quad (8)$$

Generate the following probability distribution

$$P(\mathbf{x}_{missed} | \mathbf{x}_{observed}, \Theta) = \frac{P(\mathbf{x}_{missed}, \mathbf{x}_{observed} | \Theta)}{P(\mathbf{x}_{observed} | \Theta)} \quad (9)$$

where

$$p(\mathbf{x}_{observed} | \Theta) = \int_{\mathcal{X}} p(\mathbf{x}_{complete} | \Theta) d\mathbf{x}_{missed} \quad (10)$$

We can use a Roulette based algorithm

Basically, we use the data to obtain a multivariate version of the data

- Then, we use the α_i in a roulette based algorithm to select a sample
 - ▶ Then, we generate $x_{missed} \sim p_j(x|\theta) + Var(x)$

We can use a Roulette based algorithm

Basically, we use the data to obtain a multivariate version of the data

- Then, we use the α_i in a roulette based algorithm to select a sample
 - ▶ Then, we generate $x_{missed} \sim p_j(x|\theta) + Var(x)$

This is the most simple

- What about something more complex?

For this, we can do

We have the following joint probability

$$f(x_{\text{missed}}, x_{\text{observed}} | \theta)$$

For this, we can do

We have the following joint probability

$$f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta)$$

Thus, the complete log likelihood

$$\ell(\theta) = \log f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta)$$

For this, we can do

We have the following joint probability

$$f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta)$$

Thus, the complete log likelihood

$$\ell(\theta) = \log f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta)$$

Therefore, we have when integrating the missing (Yes! Marginalization)

$$l_{\mathbf{x}_{\text{missed}}}(\theta) = \log \int f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta) d\mathbf{x}_{\text{missed}}$$

Here, it is quite interesting to observe

We have a ratio like this

$$\log \frac{f(\mathbf{x}_{missed}, \mathbf{x}_{observed} | \theta)}{f(\mathbf{x}_{missed}, \mathbf{x}_{observed} | \theta_t)}$$

Here, it is quite interesting to observe

We have a ratio like this

$$\log \frac{f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta)}{f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta_t)}$$

Basically, we can use this function on the EM

$$\begin{aligned} Q(\theta | \theta_t) &= E_{\theta_t} \left[\log \frac{f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta)}{f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta_t)} \right] \\ &= \int \log \frac{f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta)}{f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}} | \theta_t)} f(\mathbf{x}_{\text{observed}} | \mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}} \end{aligned}$$

In this case

Why this ratio?

- Actually, because we want the missing data to be estimated by the observed one

In this case

Why this ratio?

- Actually, because we want the missing data to be estimated by the observed one

Actually... There is something quite interesting

- Kullback–Leibler Divergence!!! Yes this ratio is similar

Actually the Kullback–Leibler Divergence

Definition

- For probability distributions P and Q defined on the same probability space, \mathcal{X} , the Kullback–Leibler divergence is defined as

$$KL(P \parallel Q) = \int p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

Actually the Kullback–Leibler Divergence

Definition

- For probability distributions P and Q defined on the same probability space, \mathcal{X} , the Kullback–Leibler divergence is defined as

$$KL(P \parallel Q) = \int p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

Thus, we have that Q is actually a KL version!!!

$$\begin{aligned} Q(\theta|\theta_t) &= \int \log \frac{f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}}|\theta)}{f(\mathbf{x}_{\text{missed}}, \mathbf{x}_{\text{observed}}|\theta_t)} f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}} \\ &= \int \log \left(\underbrace{\frac{f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta) f(\mathbf{x}_{\text{missed}}|\theta)}{f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) f(\mathbf{x}_{\text{missed}}|\theta_t)}}_{\frac{p(x)}{q(x)}} \right) \underbrace{f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}}}_{p(x)} \end{aligned}$$

A Small Problem and Fixing it

We have from EM

- We have that $\mathcal{L}(\Theta) \geq \mathcal{L}(\Theta_n)$

A Small Problem and Fixing it

We have from EM

- We have that $\mathcal{L}(\Theta) \geq \mathcal{L}(\Theta_n)$

So the new Q does not have this difference only the KL

- Basically, the Q lacks a way to enforce this regularization

A Small Problem and Fixing it

We have from EM

- We have that $\mathcal{L}(\Theta) \geq \mathcal{L}(\Theta_n)$

So the new Q does not have this difference only the KL

- Basically, the Q lacks a way to enforce this regularization

A simple solution

- As in Ridge Regression

Basically, we can integrate this

Generate a new Q , $l_y(\Theta) - l_y(\Theta_n)$ (EM) and KL Divergence for a Q

$$\begin{aligned} Q(\theta|\theta_t) &= \log f(\mathbf{x}_{\text{missed}}|\theta) \int f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}} - \dots \\ &\quad \log f(\mathbf{x}_{\text{missed}}|\theta_t) \int f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}} + \dots \\ &\quad \int_{\theta_t} \log \frac{f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta)}{f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t)} f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}} \end{aligned}$$

Basically, we can integrate this

Generate a new Q , $l_y(\Theta) - l_y(\Theta_n)$ (EM) and KL Divergence for a Q

$$\begin{aligned} Q(\theta|\theta_t) &= \log f(\mathbf{x}_{\text{missed}}|\theta) \int f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}} - \dots \\ &\quad \log f(\mathbf{x}_{\text{missed}}|\theta_t) \int f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}} + \dots \\ &\quad \int_{\theta_t} \log \frac{f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta)}{f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t)} f(\mathbf{x}_{\text{observed}}|\mathbf{x}_{\text{missed}}, \theta_t) d\mathbf{x}_{\text{observed}} \end{aligned}$$

Using a little bit of notation

$$Q(\theta|\theta_t) = l_y(\theta) - l_y(\theta_t) - KL\left(f_{\theta_t}^{\mathbf{x}_{\text{missed}}} \parallel f_{\theta}^{\mathbf{x}_{\text{missed}}}\right)$$

KL-divergence is minimized for $\theta = \theta_t$, actually zero!!!

Then when differentiating the Q divergence

$$\left. \frac{\partial Q(\theta|\theta_t)}{\partial \theta} \right|_{\theta=\theta_y} = \left. \frac{\partial l_{x_{missed}}(\theta)}{\partial \theta} \right|_{\theta=\theta_y}$$

KL-divergence is minimized for $\theta = \theta_t$, actually zero!!!

Then when differentiating the Q divergence

$$\left. \frac{\partial Q(\theta|\theta_t)}{\partial \theta} \right|_{\theta=\theta_y} = \left. \frac{\partial l_{x_{missed}}(\theta)}{\partial \theta} \right|_{\theta=\theta_y}$$

Thus define the iteration as

$$\theta_{t+1} = \arg \max_{\theta} Q(\theta|\theta_t)$$

It is possible to see that

Something Notable

$$Q(\theta_{t+1}|\theta_t) + l_y(\theta_t) + KL(f_{\theta_t}^{x_{missed}} \parallel f_{\theta_t}^{x_{missed}}) = l_y(\theta_{t+1})$$

It is possible to see that

Something Notable

$$Q(\theta_{t+1}|\theta_t) + l_y(\theta_t) + KL\left(f_{\theta_t}^{x_{missed}} \parallel f_{\theta_t}^{x_{missed}}\right) = l_y(\theta_{t+1})$$

Then

$$l_y(\theta_{t+1}) \geq l_y(\theta_t) + 0 + 0$$

It is possible to see that

Something Notable

$$Q(\theta_{t+1}|\theta_t) + l_y(\theta_t) + KL\left(f_{\theta_t}^{x_{miss}} \parallel f_{\theta_t}^{x_{miss}}\right) = l_y(\theta_{t+1})$$

Then

$$l_y(\theta_{t+1}) \geq l_y(\theta_t) + 0 + 0$$

Thus

- The log-likelihood never decreases after a combined *E-step* and *M-step*.

Here, everything looks great but...

We need to know to which distribution could come the result

- Thus, we have that we assume that the missing data can come from two distributions!!!

Here, everything looks great but...

We need to know to which distribution could come the result

- Thus, we have that we assume that the missing data can come from two distributions!!!

Start from the simple

- We assume a two possible sources of the information for the missing data.

Thus, we can device the following Likelihood

We can consider a sample $Y = \{Y_1, \dots, Y_n\}$ from individual densities

$$f(y|\alpha, \mu) = \alpha \phi(y - \mu) + (1 - \alpha) \phi(y)$$

Thus, we can device the following Likelihood

We can consider a sample $Y = \{Y_1, \dots, Y_n\}$ from individual densities

$$f(y|\alpha, \mu) = \alpha \phi(y - \mu) + (1 - \alpha) \phi(y)$$

Where, we can impose the following distribution

$$\phi(y) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{y^2}{2} \right\}$$

- With both α and μ are both unknown, but $0 < \alpha < 1$.

Incomplete observation

The likelihood function becomes

$$L_{x_{\text{missed}}}(\alpha, \mu) = \prod_{i=1}^N \alpha \phi(y_i - \mu) + (1 - \alpha) \phi(y_i)$$

Incomplete observation

The likelihood function becomes

$$L_{x_{\text{missed}}}(\alpha, \mu) = \prod_{i=1}^N \alpha \phi(y_i - \mu) + (1 - \alpha) \phi(y_i)$$

This is a quite unpleasant function

- But suppose we knew which observations came from which population?

What?

Let $X = \{x_1, \dots, x_n\}$ be i.i.d. with $P(x_i = 1) = \alpha$

- Then, we play the hierarchical idea

What?

Let $X = \{x_1, \dots, x_n\}$ be i.i.d. with $P(x_i = 1) = \alpha$

- Then, we play the hierarchical idea

Hierarchy

$$y_i \sim N(\mu, 1) \text{ if } x_i = 1$$

$$y_i \sim N(0, 1) \text{ if } x_i = 0$$

What?

Let $X = \{x_1, \dots, x_n\}$ be i.i.d. with $P(x_i = 1) = \alpha$

- Then, we play the hierarchical idea

Hierachy

$$y_i \sim N(\mu, 1) \text{ if } x_i = 1$$

$$y_i \sim N(0, 1) \text{ if } x_i = 0$$

i.e x_i allows to indicate to which distribution y_i belongs

- Then we need the marginal distribution of Y .

Thus

The Complete Data Likelihood is

$$L_{x,y}(\alpha, \mu) = \prod_{i=1}^N \alpha^{x_i} \phi(y_i - \mu)^{x_i} (1 - \alpha)^{1-x_i} \phi(y_i)^{1-x_i}$$

Thus

The Complete Data Likelihood is

$$L_{x,y}(\alpha, \mu) = \prod_{i=1}^N \alpha^{x_i} \phi(y_i - \mu)^{x_i} (1 - \alpha)^{1-x_i} \phi(y_i)^{1-x_i}$$

Or given that $\phi(y_i)$ does not contain any parameter

$$L_{x,y}(\alpha, \mu) \propto \alpha^{\sum x_i} (1 - \alpha)^{n - \sum x_i} \prod_{i=1}^N \phi(y_i - \mu)^{x_i}$$

Then taking logarithms

We have that

$$l_{x,y}(\alpha, \mu) = \sum x_i \log \alpha + \left(n - \sum x_i\right) \log(1 - \alpha) - \sum \frac{x_i (y_i - \mu)^2}{2}$$

Then taking logarithms

We have that

$$l_{x,y}(\alpha, \mu) = \sum x_i \log \alpha + \left(n - \sum x_i\right) \log(1 - \alpha) - \sum \frac{x_i (y_i - \mu)^2}{2}$$

Therefore, if we differentiate

$$\hat{\alpha} = \frac{1}{\sum x_i}, \hat{\mu} = \frac{\sum x_i y_i}{\sum x_i}$$

Then taking logarithms

We have that

$$l_{x,y}(\alpha, \mu) = \sum x_i \log \alpha + \left(n - \sum x_i\right) \log(1 - \alpha) - \sum \frac{x_i (y_i - \mu)^2}{2}$$

Therefore, if we differentiate

$$\hat{\alpha} = \frac{1}{\sum x_i} \sum x_i, \hat{\mu} = \frac{\sum x_i y_i}{\sum x_i}$$

We have seen this formulations

- The EM algorithm for the Mixture of Gaussian's

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- **Missing Data**
 - Using EM
- **Matrix Completion**
 - The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Example

We have two matrices

- Data Matrix X
- Missing Data M

$$M_{ij} = \begin{cases} 0 & X_{ij} \text{ is missing} \\ 1 & X_{ij} \text{ is not missing} \end{cases}$$

Example

We have two matrices

- Data Matrix X
- Missing Data M

$$M_{ij} = \begin{cases} 0 & X_{ij} \text{ is missing} \\ 1 & X_{ij} \text{ is not missing} \end{cases}$$

Therefore, we have

- $X = (X_{obs}, X_{mis})$

Example

We have two matrices

- Data Matrix X
- Missing Data M

$$M_{ij} = \begin{cases} 0 & X_{ij} \text{ is missing} \\ 1 & X_{ij} \text{ is not missing} \end{cases}$$

Therefore, we have

- $X = (X_{obs}, X_{mis})$

This comes from

- “Bayes and multiple imputation” by RJA Little, DB Rubin (2002)

We can use the following optimization

We can do the following

$$\min_{M_{ij}=1} \|X - AB\|_F$$

We can use the following optimization

We can do the following

$$\min_{M_{ij}=1} \|X - AB\|_F$$

Clearly an initial matrix decomposition, where

$$M_{ij}x_{ij} \approx \sum_{k=1}^K a_{ik}b_{kj}$$

We can use the following optimization

We can do the following

$$\min_{M_{ij}=1} \|X - AB\|_F$$

Clearly an initial matrix decomposition, where

$$M_{ij}x_{ij} \approx \sum_{k=1}^K a_{ik}b_{kj}$$

So the total error to be minimized is

$$\min_{M_{ij}=1} \|X - AB\|_F = \sqrt{\sum_{i=1}^N \sum_{j=1}^M \left[M_{ij}x_{ij} - \sum_{k=1}^K a_{ik}b_{kj} \right]^2}$$

- $K \ll N, M$

This can be regularized

Using the following ideas

$$\min_{M_{ij}=1} \|X - AB\|_F + \lambda \left[\|A\|^2 + \|B\|^2 \right]$$

This can be regularized

Using the following ideas

$$\min_{M_{ij}=1} \|X - AB\|_F + \lambda \left[\|A\|^2 + \|B\|^2 \right]$$

Therefore, once the minimization is achieved

- We finish with two dense matrices A, B that can be used to obtain the elements with entries $M_{ij} = 0$

There are many other methods for this

For example

- Moritz Hardt. Understanding Alternating Minimization for Matrix Completion. FOCS, pages 651–660, 2014.
- Moritz Hardt, Mary Wootters. Fast matrix completion without the condition number. COLT, pages 638–678, 20
- Raghunandan H Keshavan, Andrea Montanari, and Sewoong Oh, Matrix completion from noisy entries, The Journal of Machine Learning Research 99 (2010), 2057–2078.
- Stephen J Wright, Robert D Nowak, and M´ario AT Figueiredo, Sparse reconstruction by separable approximation, Signal Processing, IEEE Transactions on 57 (2009), no. 7, 2479–2493.

Outline

- 1 Introduction
 - Feature Engineering
 - What is Feature Selection?
 - Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
 - Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

- 2 Feature Selection
 - Feature Selection
 - Feature selection based on statistical hypothesis testing
 - Example
 - Application of the t -Test in Feature Selection
 - Example
 - Considering Feature Sets
 - Scatter Matrices
 - What to do with it?
 - Sequential Backward Selection

- 3 Shrinkage Methods
 - Introduction
 - Intuition from Overfitting
 - The Idea of Regularization
 - Ridge Regression
 - Standardization of Data
 - The LASSO
 - The Lagrangian Version of the LASSO

THE PEAKING PHENOMENON

Remember

Normally, to design a classifier with good generalization performance, we want the number of sample N to be larger than the number of features d .

THE PEAKING PHENOMENON

Remeber

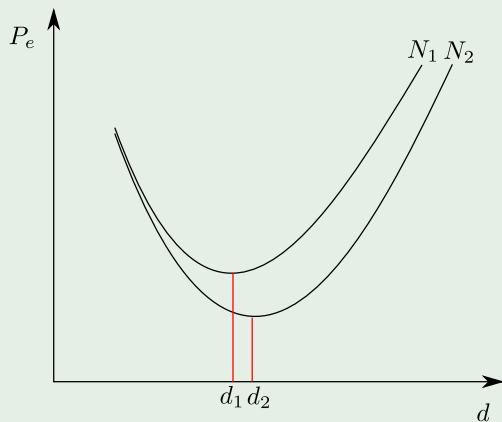
Normally, to design a classifier with good generalization performance, we want the number of sample N to be larger than the number of features d .

What?

The intuition, the larger the number of samples vs the number of features, the smaller the error P_e

Graphically

For $N_2 \gg N_1$



Let us explain

Something Notable

Let's look at the following example from the paper:

- “A Problem of Dimensionality: A Simple Example” by G.A. Trunk

THE PEAKING PHENOMENON

Assume the following problem

We have two classes ω_1, ω_2 such that

$$P(\omega_1) = P(\omega_2) = \frac{1}{2} \quad (11)$$

THE PEAKING PHENOMENON

Assume the following problem

We have two classes ω_1, ω_2 such that

$$P(\omega_1) = P(\omega_2) = \frac{1}{2} \quad (11)$$

Both Classes have the following Gaussian distribution

- ① $\omega_1 \Rightarrow \mu$ and $\Sigma = I$
- ② $\omega_2 \Rightarrow -\mu$ and $\Sigma = I$

THE PEAKING PHENOMENON

Assume the following problem

We have two classes ω_1, ω_2 such that

$$P(\omega_1) = P(\omega_2) = \frac{1}{2} \quad (11)$$

Both Classes have the following Gaussian distribution

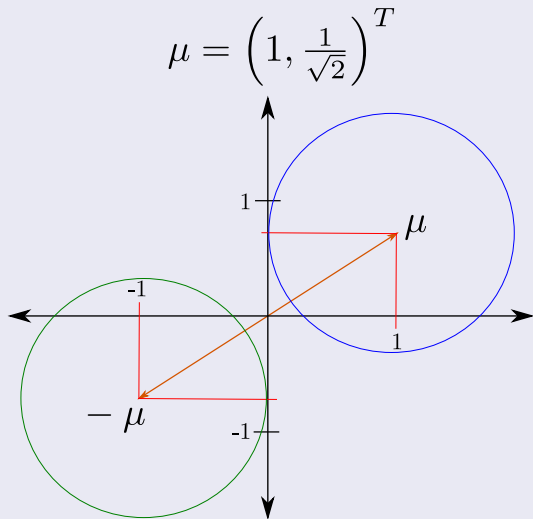
- ① $\omega_1 \Rightarrow \mu$ and $\Sigma = I$
- ② $\omega_2 \Rightarrow -\mu$ and $\Sigma = I$

Where

$$\mu = \left[1, \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{3}}, \dots, \frac{1}{\sqrt{d}} \right]$$

Example

The μ for \mathbb{R}^2



THE PEAKING PHENOMENON

Properties of the features

Since the features are jointly Gaussian and $\Sigma = I$, the involved features are statistically independent.

THE PEAKING PHENOMENON

Properties of the features

Since the features are jointly Gaussian and $\Sigma = I$, the involved features are statistically independent.

We use the following rule to classify

if for any vector x , we have that

THE PEAKING PHENOMENON

Properties of the features

Since the features are jointly Gaussian and $\Sigma = I$, the involved features are statistically independent.

We use the following rule to classify

if for any vector x , we have that

$$\textcircled{1} \quad \|x - \mu\|^2 < \|x + \mu\|^2 \text{ or } z \equiv x^T \mu > 0 \text{ then } x \in \omega_1.$$

THE PEAKING PHENOMENON

Properties of the features

Since the features are jointly Gaussian and $\Sigma = I$, the involved features are statistically independent.

We use the following rule to classify

if for any vector x , we have that

- 1 $\|x - \mu\|^2 < \|x + \mu\|^2$ or $z \equiv x^T \mu > 0$ then $x \in \omega_1$.
- 2 $z \equiv x^T \mu < 0$ then $x \in \omega_2$.

A little bit of algebra

For the first case

$$\begin{aligned}(x - \mu)^T (x - \mu) &< (x + \mu)^T (x + \mu) \\ x^t x - 2x^T \mu + \mu^T \mu &< x^t x + 2x^T \mu + \mu^T \mu \\ 0 &< x^T \mu \equiv z\end{aligned}$$

A little bit of algebra

For the first case

$$\begin{aligned}x^t x - 2x^T \mu + \mu^T \mu &< x^t x + 2x^T \mu + \mu^T \mu \\ 0 &< x^T \mu \equiv z\end{aligned}$$

A little bit of algebra

For the first case

$$0 < \mathbf{x}^T \boldsymbol{\mu} \equiv z$$

We have then two cases

- 1 Known mean value $\boldsymbol{\mu}$.
- 2 Unknown mean value $\boldsymbol{\mu}$.

A little bit of algebra

For the first case

$$\begin{aligned}\|x - \mu\|^2 &< \|x + \mu\|^2 \\ (x - \mu)^T (x - \mu) &< (x + \mu)^T (x + \mu) \\ x^t x - 2x^T \mu + \mu^T \mu &< x^t x + 2x^T \mu + \mu^T \mu \\ 0 &< x^T \mu \equiv z\end{aligned}$$

We have then two cases

A little bit of algebra

For the first case

$$\begin{aligned}\|x - \mu\|^2 &< \|x + \mu\|^2 \\ (x - \mu)^T (x - \mu) &< (x + \mu)^T (x + \mu) \\ x^t x - 2x^T \mu + \mu^T \mu &< x^t x + 2x^T \mu + \mu^T \mu \\ 0 &< x^T \mu \equiv z\end{aligned}$$

We have then two cases

- 1 Known mean value μ .

A little bit of algebra

For the first case

$$\begin{aligned}\|x - \mu\|^2 &< \|x + \mu\|^2 \\ (x - \mu)^T (x - \mu) &< (x + \mu)^T (x + \mu) \\ x^T x - 2x^T \mu + \mu^T \mu &< x^T x + 2x^T \mu + \mu^T \mu \\ 0 &< x^T \mu \equiv z\end{aligned}$$

We have then two cases

- 1 Known mean value μ .
- 2 Unknown mean value μ .

Known mean value μ

Given that z is a linear combination of independent Gaussian Variables

- 1 It is a Gaussian variable.

Known mean value μ

Given that z is a linear combination of independent Gaussian Variables

- 1 It is a Gaussian variable.
- 2 $E[z] = \sum_{i=1}^d \mu_i E(x_i) = \sum_{i=1}^d \frac{1}{\sqrt{i}} \frac{1}{\sqrt{i}} = \sum_{i=1}^d \frac{1}{i} = \|\mu\|^2.$

Known mean value μ

Given that z is a linear combination of independent Gaussian Variables

- 1 It is a Gaussian variable.
- 2 $E[z] = \sum_{i=1}^d \mu_i E(x_i) = \sum_{i=1}^d \frac{1}{\sqrt{i}} \frac{1}{\sqrt{i}} = \sum_{i=1}^d \frac{1}{i} = \|\mu\|^2.$
- 3 $\sigma_z^2 = \|\mu\|^2.$

Known mean value μ

Given that z is a linear combination of independent Gaussian Variables

- 1 It is a Gaussian variable.
- 2 $E[z] = \sum_{i=1}^d \mu_i E(x_i) = \sum_{i=1}^d \frac{1}{\sqrt{i}} \frac{1}{\sqrt{i}} = \sum_{i=1}^d \frac{1}{i} = \|\mu\|^2.$
- 3 $\sigma_z^2 = \|\mu\|^2.$

Why the first statement?

Given that each feature of x

It can be seen as random variable with mean $\frac{1}{\sqrt{i}}$ and variance 1 with no correlation between each other.

Why the first statement?

Given that each feature of x

It can be seen as random variable with mean $\frac{1}{\sqrt{i}}$ and variance 1 with no correlation between each other.

What about the variance of z ?

$$=E \left[z^2 - 2z \|\mu\|^2 + \|\mu\|^4 \right]$$

$$=E \left[z^2 \right] - \|\mu\|^4$$

$$=E \left[\left(\sum_{i=1}^d \mu_i x_i \right) \left(\sum_{i=1}^d \mu_i x_i \right) \right] - \left(\sum_{i=1}^d \frac{1}{i^2} + \sum_{j=1}^d \sum_{\substack{h=1 \\ j \neq h}}^d \frac{1}{i} \times \frac{1}{j} \right)$$

Why the first statement?

Given that each feature of x

It can be seen as random variable with mean $\frac{1}{\sqrt{i}}$ and variance 1 with no correlation between each other.

What about the variance of z ?

$$\begin{aligned} &= E[z^2] - \|\mu\|^4 \\ &= E\left[\left(\sum_{i=1}^d \mu_i x_i\right) \left(\sum_{i=1}^d \mu_i x_i\right)\right] - \left(\sum_{i=1}^d \frac{1}{i^2} + \sum_{\substack{j=1 \\ j \neq h}}^d \sum_{h=1}^d \frac{1}{i} \times \frac{1}{j}\right) \end{aligned}$$

Why the first statement?

Given that each feature of x

It can be seen as random variable with mean $\frac{1}{\sqrt{i}}$ and variance 1 with no correlation between each other.

What about the variance of z ?

$$= E \left[\left(\sum_{i=1}^d \mu_i x_i \right) \left(\sum_{i=1}^d \mu_i x_i \right) \right] - \left(\sum_{i=1}^d \frac{1}{i^2} + \sum_{j=1}^d \sum_{\substack{h=1 \\ j \neq h}}^d \frac{1}{i} \times \frac{1}{j} \right)$$

Why the first statement?

Given that each feature of x

It can be seen as random variable with mean $\frac{1}{\sqrt{i}}$ and variance 1 with no correlation between each other.

What about the variance of z ?

$$\begin{aligned} \text{Var}(z) &= E \left[\left(z - \|\mu\|^2 \right)^2 \right] \\ &= E \left[z^2 - 2z \|\mu\|^2 + \|\mu\|^4 \right] \\ &= E \left[z^2 \right] - \|\mu\|^4 \\ &= E \left[\left(\sum_{i=1}^d \mu_i x_i \right) \left(\sum_{i=1}^d \mu_i x_i \right) \right] - \left(\sum_{i=1}^d \frac{1}{i^2} + \sum_{j=1}^d \sum_{\substack{h=1 \\ j \neq h}}^d \frac{1}{i} \times \frac{1}{j} \right) \end{aligned}$$

Thus

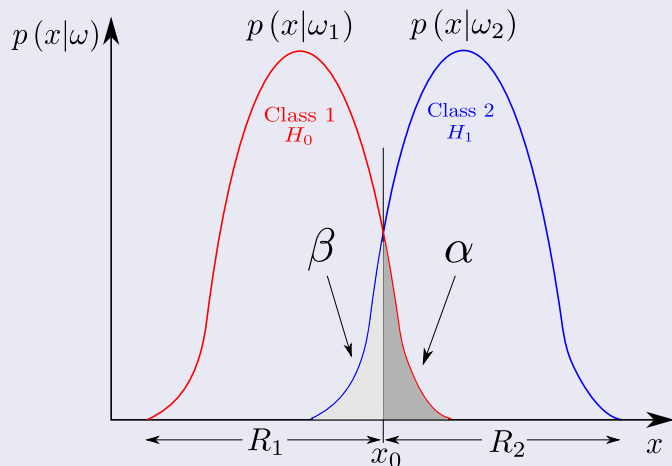
But, given that $x_i^2 \sim \chi_1^2 \left(\frac{1}{i} \right)$, with mean

$$E \left[x_i^2 \right] = 1 + \frac{1}{i} \quad (12)$$

Remark: The rest is for you to solve so $\sigma_z^2 = \|\boldsymbol{\mu}\|^2$.

Remember the P_e

We have then...



We get the probability of error

We know that the error is coming from the following equation

$$P_e = \frac{1}{2} \int_{-\infty}^{x_0} p(z|\omega_2) d\mathbf{x} + \frac{1}{2} \int_{x_0}^{\infty} p(z|\omega_1) d\mathbf{x} \quad (13)$$

We get the probability of error

We know that the error is coming from the following equation

$$P_e = \frac{1}{2} \int_{-\infty}^{x_0} p(z|\omega_2) d\mathbf{x} + \frac{1}{2} \int_{x_0}^{\infty} p(z|\omega_1) d\mathbf{x} \quad (13)$$

But, we have equiprobable classes

$$= \int_{x_0}^{\infty} p(z|\omega_1) d\mathbf{x}$$

We get the probability of error

We know that the error is coming from the following equation

$$P_e = \frac{1}{2} \int_{-\infty}^{x_0} p(z|\omega_2) dx + \frac{1}{2} \int_{x_0}^{\infty} p(z|\omega_1) dx \quad (13)$$

But, we have equiprobable classes

$$\begin{aligned} P_e &= \frac{1}{2} \int_{-\infty}^{x_0} p(z|\omega_2) dx + \frac{1}{2} \int_{x_0}^{\infty} p(z|\omega_1) dx \\ &= \int_{x_0}^{\infty} p(z|\omega_1) dx \end{aligned}$$

Thus, we have that

Now, given that z is a sum of Gaussian

$$\text{exp term} = -\frac{1}{2\|\boldsymbol{\mu}\|^2} \left[\left(z - \|\boldsymbol{\mu}\|^2 \right)^2 \right] \quad (14)$$

Thus, we have that

Now, given that z is a sum of Gaussian

$$\text{exp term} = -\frac{1}{2\|\boldsymbol{\mu}\|^2} \left[\left(z - \|\boldsymbol{\mu}\|^2 \right)^2 \right] \quad (14)$$

Because we have the rule

We can do a change of variable to a normalized z

$$P_e = \int_{b_d}^{\infty} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{z^2}{2} \right\} dz \quad (15)$$

Known mean value μ

The probability of error is given by

$$P_e = \int_{b_d}^{\infty} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{z^2}{2} \right\} dz \quad (16)$$

Known mean value μ

The probability of error is given by

$$P_e = \int_{b_d}^{\infty} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{z^2}{2} \right\} dz \quad (16)$$

Where

$$b_d = \sqrt{\sum_{i=1}^d \frac{1}{i}} \quad (17)$$

How?

Known mean value μ

Thus

When the series b_d tends to infinity as $d \rightarrow \infty$, the probability of error tends to **zero** as the number of features increases.

Unknown mean value μ

For This, we use the maximum likelihood

$$\hat{\mu} = \frac{1}{N} \sum_{k=1}^N s_k \mathbf{x}_k \quad (18)$$

where

Unknown mean value μ

For This, we use the maximum likelihood

$$\hat{\mu} = \frac{1}{N} \sum_{k=1}^N s_k \mathbf{x}_k \quad (18)$$

where

① $s_k = 1$ if $\mathbf{x}_k \in \omega_1$

Unknown mean value μ

For This, we use the maximum likelihood

$$\hat{\mu} = \frac{1}{N} \sum_{k=1}^N s_k \mathbf{x}_k \quad (18)$$

where

- ① $s_k = 1$ if $\mathbf{x}_k \in \omega_1$
- ② $s_k = -1$ if $\mathbf{x}_k \in \omega_2$

Unknown mean value μ

For This, we use the maximum likelihood

$$\hat{\mu} = \frac{1}{N} \sum_{k=1}^N s_k \mathbf{x}_k \quad (18)$$

where

- ① $s_k = 1$ if $\mathbf{x}_k \in \omega_1$
- ② $s_k = -1$ if $\mathbf{x}_k \in \omega_2$

Now, we have a problem z is no more a Gaussian variable

Still, if we select d large enough and knowing that $z = \sum x_i \hat{\mu}_i$, then for the central limit theorem, we can consider z to be Gaussian.

Unknown mean value μ

For This, we use the maximum likelihood

$$\hat{\mu} = \frac{1}{N} \sum_{k=1}^N s_k \mathbf{x}_k \quad (18)$$

where

- ① $s_k = 1$ if $\mathbf{x}_k \in \omega_1$
- ② $s_k = -1$ if $\mathbf{x}_k \in \omega_2$

Now, we have a problem z is no more a Gaussian variable

Still, if we select d large enough and knowing that $z = \sum x_i \hat{\mu}_i$, then for the central limit theorem, we can consider z to be Gaussian.

With mean and variance

- ① $E[z] = \sum_{i=1}^d \frac{1}{i}.$

Unknown mean value μ

For This, we use the maximum likelihood

$$\hat{\mu} = \frac{1}{N} \sum_{k=1}^N s_k \mathbf{x}_k \quad (18)$$

where

- ① $s_k = 1$ if $\mathbf{x}_k \in \omega_1$
- ② $s_k = -1$ if $\mathbf{x}_k \in \omega_2$

Now, we have a problem z is no more a Gaussian variable

Still, if we select d large enough and knowing that $z = \sum x_i \hat{\mu}_i$, then for the central limit theorem, we can consider z to be Gaussian.

With mean and variance

- ① $E[z] = \sum_{i=1}^d \frac{1}{i}$.
- ② $\sigma_z^2 = \left(1 + \frac{1}{N}\right) \sum_{i=1}^d \frac{1}{i} + \frac{d}{N}$.

Unknown mean value μ

Thus

$$b_d = \frac{E[z]}{\sigma_z} \quad (19)$$

Unknown mean value μ

Thus

$$b_d = \frac{E[z]}{\sigma_z} \quad (19)$$

Thus, using P_e

- It can now be shown that $b_d \rightarrow 0$ as $d \rightarrow \infty$ and the probability of error tends to $\frac{1}{2}$ for any finite number N .

Finally

Case I

- If for any d the corresponding PDF is known, then we can perfectly discriminate the two classes by arbitrarily increasing the number of features.

Finally

Case I

- If for any d the corresponding PDF is known, then we can perfectly discriminate the two classes by arbitrarily increasing the number of features.

Case II

- If the PDF's are not known, then the arbitrary increase of the number of features leads to the maximum possible value of the error rate, that is, $\frac{1}{2}$.

Finally

Case I

- If for any d the corresponding PDF is known, then we can perfectly discriminate the two classes by arbitrarily increasing the number of features.

Case II

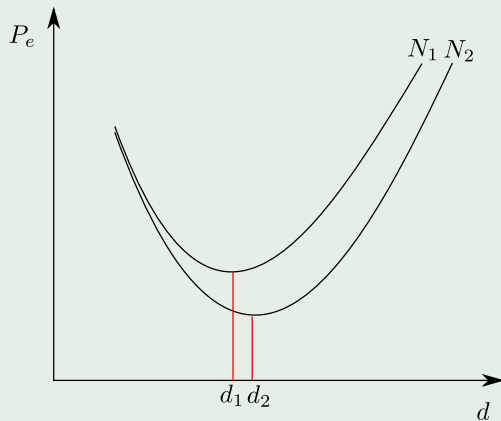
- If the PDF's are not known, then the arbitrary increase of the number of features leads to the maximum possible value of the error rate, that is, $\frac{1}{2}$.

Thus

- Under a limited number of training data we must try to keep the number of features to a relatively low number.

Graphically

For $N_2 \gg N_1$, minimum at $d = \frac{N}{\alpha}$ with $\alpha \in [2, 10]$



Back to Feature Selection

The Goal

- 1 Select the “optimum” number d of features.

Back to Feature Selection

The Goal

- 1 Select the “optimum” number d of features.
- 2 Select the “best” d features.

Back to Feature Selection

The Goal

- 1 Select the “optimum” number d of features.
- 2 Select the “best” d features.

Why? Large d has a three-fold disadvantage:

- High computational demands.

Back to Feature Selection

The Goal

- 1 Select the “optimum” number d of features.
- 2 Select the “best” d features.

Why? Large d has a three-fold disadvantage:

- High computational demands.
- Low generalization performance.

Back to Feature Selection

The Goal

- 1 Select the “optimum” number d of features.
- 2 Select the “best” d features.

Why? Large d has a three-fold disadvantage:

- High computational demands.
- Low generalization performance.
- Poor error estimates

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
- Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- **Feature Selection**
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Back to Feature Selection

Given N

d must be large enough to learn what makes classes different and what makes patterns in the same class similar

Back to Feature Selection

Given N

d must be large enough to learn what makes classes different and what makes patterns in the same class similar

In addition

d must be small enough not to learn what makes patterns of the same class different

Back to Feature Selection

Given N

d must be large enough to learn what makes classes different and what makes patterns in the same class similar

In addition

d must be small enough not to learn what makes patterns of the same class different

In practice

In practice, $d < N/3$ has been reported to be a sensible choice for a number of cases

Thus

Oh!!!

Once d has been decided, choose the d most informative features:

Thus

Oh!!!

Once d has been decided, choose the d most informative features:

Best: Large between class distance, Small within class variance.

Thus

Oh!!!

Once d has been decided, choose the d most informative features:

Best: Large between class distance, Small within class variance.

The basic philosophy

- 1 Discard individual features with poor information content.

Thus

Oh!!!

Once d has been decided, choose the d most informative features:

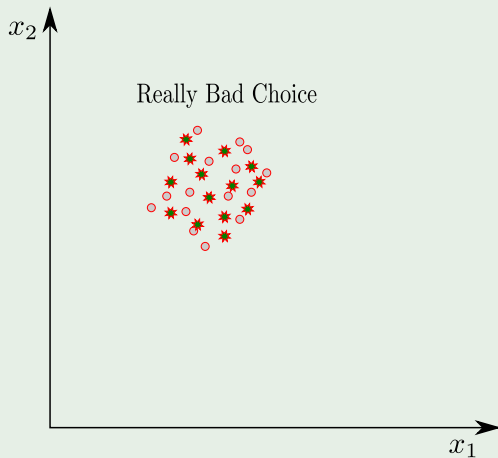
Best: Large between class distance, Small within class variance.

The basic philosophy

- 1 Discard individual features with poor information content.
- 2 The remaining information rich features are examined jointly as vectors

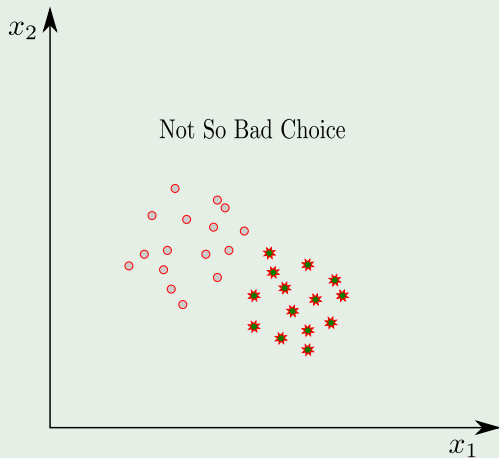
Example

Thus, we want to avoid choices



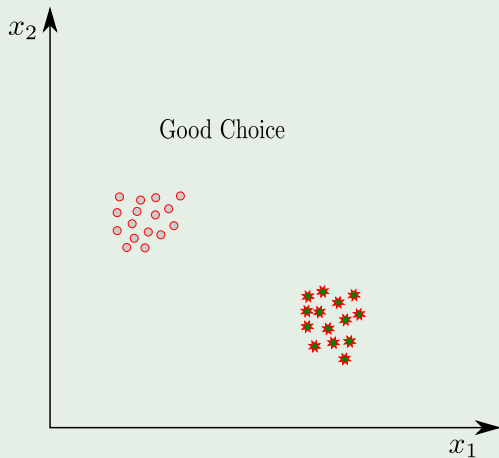
Example

Better Choice



Example

What We Want to Have



Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
- Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- **Feature selection based on statistical hypothesis testing**
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Using Statistics

Simplicity First Principles - Marcus Aurelius

- A first step in feature selection is to look at each of the generated features independently.

Using Statistics

Simplicity First Principles - Marcus Aurelius

- A first step in feature selection is to look at each of the generated features independently.
- Then, test their discriminatory capability for the problem at hand.

Using Statistics

Simplicity First Principles - Marcus Aurelius

- A first step in feature selection is to look at each of the generated features independently.
- Then, test their discriminatory capability for the problem at hand.

For this, we can use the following hypothesis testing

Assume the samples for two classes ω_1 , ω_2 are vectors of random variables.

Using Statistics

Simplicity First Principles - Marcus Aurelius

- A first step in feature selection is to look at each of the generated features independently.
- Then, test their discriminatory capability for the problem at hand.

For this, we can use the following hypothesis testing

Assume the samples for two classes ω_1, ω_2 are vectors of random variables.

- 1 H_1 : The values of the feature differ significantly

Using Statistics

Simplicity First Principles - Marcus Aurelius

- A first step in feature selection is to look at each of the generated features independently.
- Then, test their discriminatory capability for the problem at hand.

For this, we can use the following hypothesis testing

Assume the samples for two classes ω_1, ω_2 are vectors of random variables.

- 1 H_1 : The values of the feature differ significantly
- 2 H_0 : The values of the feature do not differ significantly

Using Statistics

Simplicity First Principles - Marcus Aurelius

- A first step in feature selection is to look at each of the generated features independently.
- Then, test their discriminatory capability for the problem at hand.

For this, we can use the following hypothesis testing

Assume the samples for two classes ω_1, ω_2 are vectors of random variables.

- 1 H_1 : The values of the feature differ significantly
- 2 H_0 : The values of the feature do not differ significantly

Meaning

H_0 is known as the null hypothesis and H_1 as the alternative hypothesis.

Hypothesis Testing Basics

We need to represent these ideas in a more mathematical way

For this, given an unknown parameter θ :

Hypothesis Testing Basics

We need to represent these ideas in a more mathematical way

For this, given an unknown parameter θ :

$$H_1 : \theta \neq \theta_0$$

$$H_0 : \theta = \theta_0$$

Hypothesis Testing Basics

We need to represent these ideas in a more mathematical way

For this, given an unknown parameter θ :

$$H_1 : \theta \neq \theta_0$$

$$H_0 : \theta = \theta_0$$

We want to generate a q

That measures the quality of our answer under our knowledge of the sample features x_1, x_2, \dots, x_N .

Hypothesis Testing Basics

We need to represent these ideas in a more mathematical way

For this, given an unknown parameter θ :

$$H_1 : \theta \neq \theta_0$$

$$H_0 : \theta = \theta_0$$

We want to generate a q

That measures the quality of our answer under our knowledge of the sample features x_1, x_2, \dots, x_N .

We ask for

- 1 Where a D (Acceptance Interval) is an interval where q lies with high probability under hypothesis H_0 .

Hypothesis Testing Basics

We need to represent these ideas in a more mathematical way

For this, given an unknown parameter θ :

$$H_1 : \theta \neq \theta_0$$

$$H_0 : \theta = \theta_0$$

We want to generate a q

That measures the quality of our answer under our knowledge of the sample features x_1, x_2, \dots, x_N .

We ask for

- 1 Where a D (Acceptance Interval) is an interval where q lies with high probability under hypothesis H_0 .
- 2 Where \overline{D} , the complement or critical region, is the region where we reject H_0 .

Hypothesis Testing Basics

We need to represent these ideas in a more mathematical way

For this, given an unknown parameter θ :

$$H_1 : \theta \neq \theta_0$$

$$H_0 : \theta = \theta_0$$

We want to generate a q

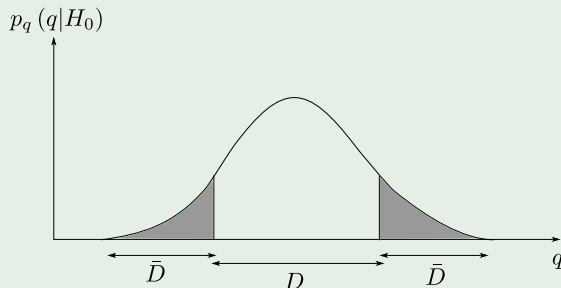
That measures the quality of our answer under our knowledge of the sample features x_1, x_2, \dots, x_N .

We ask for

- 1 Where a D (Acceptance Interval) is an interval where q lies with high probability under hypothesis H_0 .
- 2 Where \overline{D} , the complement or critical region, is the region where we reject H_0 .

Example

Acceptance and critical regions for hypothesis testing. The area of the shaded region is the probability of an erroneous decision.



Known Variance Case

Assume

Be x a random variable and x_i the resulting experimental samples.

Known Variance Case

Assume

Be x a random variable and x_i the resulting experimental samples.

Let

① $E[x] = \mu$

② $E[(x - \mu)^2] = \sigma^2$

Known Variance Case

Assume

Be x a random variable and x_i the resulting experimental samples.

Let

$$① \quad E[x] = \mu$$

$$② \quad E[(x - \mu)^2] = \sigma^2$$

We can estimate μ using

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad (20)$$

Known Variance Case

It can be proved that the

\bar{x} is an unbiased estimate of the mean of x .

Known Variance Case

It can be proved that the

\bar{x} is an unbiased estimate of the mean of x .

In a similar way

The variance of $\sigma_{\bar{x}}^2$ of \bar{x} is

$$E[(\bar{x} - \mu)^2] = E\left[\left(\frac{1}{N} \sum_{i=1}^N x_i - \mu\right)^2\right] = E\left[\left(\frac{1}{N} \sum_{i=1}^N (x_i - \mu)\right)^2\right] \quad (21)$$

Known Variance Case

It can be proved that the

\bar{x} is an unbiased estimate of the mean of x .

In a similar way

The variance of $\sigma_{\bar{x}}^2$ of \bar{x} is

$$E[(\bar{x} - \mu)^2] = E\left[\left(\frac{1}{N} \sum_{i=1}^N x_i - \mu\right)^2\right] = E\left[\left(\frac{1}{N} \sum_{i=1}^N (x_i - \mu)\right)^2\right] \quad (21)$$

Which is the following

$$E[(\bar{x} - \mu)^2] = \frac{1}{N^2} \sum_{i=1}^N E[(x_i - \mu)^2] + \frac{1}{N^2} \sum_i \sum_{j \neq i} E[(x_i - \mu)(x_j - \mu)] \quad (22)$$

Known Variance Case

Because independence

$$E[(x_i - \mu)(x_j - \mu)] = E[x_i - \mu] E[x_j - \mu] = 0 \quad (23)$$

Known Variance Case

Because independence

$$E[(x_i - \mu)(x_j - \mu)] = E[x_i - \mu] E[x_j - \mu] = 0 \quad (23)$$

Thus

$$\sigma_{\bar{x}}^2 = \frac{1}{N} \sigma^2 \quad (24)$$

Note: the larger the number of measurement samples, the smaller the variance of \bar{x} around the true mean.

What to do with it

Now, you are given a $\hat{\mu}$ the estimated parameter (In our case the mean sample)

Thus:

$$H_1 : E[x] \neq \hat{\mu}$$

$$H_0 : E[x] = \hat{\mu}$$

What to do with it

Now, you are given a $\hat{\mu}$ the estimated parameter (In our case the mean sample)

Thus:

$$H_1 : E[x] \neq \hat{\mu}$$

$$H_0 : E[x] = \hat{\mu}$$

We define q

$$q = \frac{\bar{x} - \hat{\mu}}{\frac{\sigma}{N}} \quad (25)$$

What to do with it

Now, you are given a $\hat{\mu}$ the estimated parameter (In our case the mean sample)

Thus:

$$H_1 : E[x] \neq \hat{\mu}$$

$$H_0 : E[x] = \hat{\mu}$$

We define q

$$q = \frac{\bar{x} - \hat{\mu}}{\frac{\sigma}{N}} \quad (25)$$

Recalling the central limit theorem

The probability density function of \bar{x} under H_0 is approx Gaussian $N(\hat{\mu}, \frac{\sigma}{N})$

Thus

Thus

q under H_0 is approx $N(0, 1)$

Thus

Thus

q under H_0 is approx $N(0, 1)$

Then

We can choose an acceptance level ρ with interval $D = [-x_\rho, x_\rho]$ such that q lies on it with probability $1 - \rho$.

Final Process

First Step

- Given the N experimental samples of x , compute \bar{x} and then q .

Final Process

First Step

- Given the N experimental samples of x , compute \bar{x} and then q .

Second One

- Choose the significance level ρ .

Final Process

First Step

- Given the N experimental samples of x , compute \bar{x} and then q .

Second One

- Choose the significance level ρ .

Third One

- Compute from the corresponding tables for $N(0,1)$ the acceptance interval $D = [-x_\rho, x_\rho]$ with probability $1 - \rho$.

Final Process

Final Step

If $q \in D$ decide H_0 , if not decide H_1 .

Final Process

Final Step

If $q \in D$ decide H_0 , if not decide H_1 .

Second one

- Basically, all we say is that we expect the resulting value q to lie in the high-percentage $1 - \rho$ interval.

Final Process

Final Step

If $q \in D$ decide H_0 , if not decide H_1 .

Second one

- Basically, all we say is that we expect the resulting value q to lie in the high-percentage $1 - \rho$ interval.
- If it does not, then we decide that this is because the assumed mean value is not “correct.”

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
- Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- **Feature selection based on statistical hypothesis testing**
 - **Example**
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Example

Let us consider an experiment with a random variable x of $\sigma = 0.23$

- Assume N to be equal to 16 and $\bar{x} = 1.35$
- Adopt $\rho = 0.05$

Example

Let us consider an experiment with a random variable x of $\sigma = 0.23$

- Assume N to be equal to 16 and $\bar{x} = 1.35$
- Adopt $\rho = 0.05$

We will test if the hypothesis $\hat{\mu} = 1.4$ is true

$$P \left\{ -1.97 < \frac{\bar{x} - \hat{\mu}}{0.23/4} < 1.97 \right\} = 0.95$$

Example

Let us consider an experiment with a random variable x of $\sigma = 0.23$

- Assume N to be equal to 16 and $\bar{x} = 1.35$
- Adopt $\rho = 0.05$

We will test if the hypothesis $\hat{\mu} = 1.4$ is true

$$P \left\{ -1.97 < \frac{\bar{x} - \hat{\mu}}{0.23/4} < 1.97 \right\} = 0.95$$

Therefore, we accept the hypothesis

- We have $1.237 < \hat{\mu} < 1.463$

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
- Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Application of the t -Test in Feature Selection

Very Simple

Use the difference $\mu_1 - \mu_2$ for the testing.

Application of the t -Test in Feature Selection

Very Simple

Use the difference $\mu_1 - \mu_2$ for the testing.

Note Each μ correspond to a class ω_1, ω_2

Application of the t -Test in Feature Selection

Very Simple

Use the difference $\mu_1 - \mu_2$ for the testing.

Note Each μ correspond to a class ω_1, ω_2

Thus, What is the logic?

Basically, if we have two classes... we must see different μ' s.

Application of the t -Test in Feature Selection

Very Simple

Use the difference $\mu_1 - \mu_2$ for the testing.

Note Each μ correspond to a class ω_1, ω_2

Thus, What is the logic?

Basically, if we have two classes... we must see different μ' s.

Assume that the variance of the feature values is the same in both

$$\sigma_1^2 = \sigma_2^2 = \sigma^2 \quad (26)$$

What is the Hypothesis?

A very simple one

$$H_1 : \Delta\mu = \mu_1 - \mu_2 \neq 0$$

$$H_0 : \Delta\mu = \mu_1 - \mu_2 = 0$$

What is the Hypothesis?

A very simple one

$$H_1 : \Delta\mu = \mu_1 - \mu_2 \neq 0$$

$$H_0 : \Delta\mu = \mu_1 - \mu_2 = 0$$

The new random variable is

$$z = x - y \tag{27}$$

where x , y denote the random variables corresponding to the values of the feature in the two classes.

What is the Hypothesis?

A very simple one

$$H_1 : \Delta\mu = \mu_1 - \mu_2 \neq 0$$

$$H_0 : \Delta\mu = \mu_1 - \mu_2 = 0$$

The new random variable is

$$z = x - y \quad (27)$$

where x , y denote the random variables corresponding to the values of the feature in the two classes.

Properties

- $E[z] = \mu_1 - \mu_2$
- $\sigma_z^2 = 2\sigma^2$

Then

It is possible to prove that z follows the distribution

$$N\left(\mu_1 - \mu_2, \frac{2\sigma^2}{N}\right) \quad (28)$$

Then

It is possible to prove that z follows the distribution

$$N\left(\mu_1 - \mu_2, \frac{2\sigma^2}{N}\right) \quad (28)$$

So

We can use the following

$$q = \frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{s_z \sqrt{\frac{2}{N}}} \quad (29)$$

Then

It is possible to prove that z follows the distribution

$$N\left(\mu_1 - \mu_2, \frac{2\sigma^2}{N}\right) \quad (28)$$

So

We can use the following

$$q = \frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{s_z \sqrt{\frac{2}{N}}} \quad (29)$$

where

$$s_z^2 = \frac{1}{2N - 2} \left(\sum_{i=1}^N (x_i - \bar{x})^2 + \sum_{i=1}^N (y_i - \bar{y})^2 \right) \quad (30)$$

Now

It can be shown that $\frac{s_z^2(2N-2)}{\sigma^2}$ follows

- A Chi-Square distribution with $2N - 2$ degrees of freedom.

Now

It can be shown that $\frac{s_z^2(2N-2)}{\sigma^2}$ follows

- A Chi-Square distribution with $2N - 2$ degrees of freedom.

Testing

- q turns out to follow a Chi-Square distribution with $2N - 2$ degrees of freedom

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
- Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

We have two classes

The sample measurements of a feature in two classes are

class ω_1	3.5	3.7	3.9	4.1	3.4	3.5	4.1	3.8	3.6	3.7
class ω_2	3.2	3.6	3.1	3.4	3.0	3.4	2.8	3.1	3.3	3.6

We have two classes

The sample measurements of a feature in two classes are

class ω_1	3.5	3.7	3.9	4.1	3.4	3.5	4.1	3.8	3.6	3.7
class ω_2	3.2	3.6	3.1	3.4	3.0	3.4	2.8	3.1	3.3	3.6

Now, we want to know if the feature is informative enough

$$H_1 : \Delta\mu = \mu_1 - \mu_2 \neq 0$$

$$H_0 : \Delta\mu = \mu_1 - \mu_2 = 0$$

We have two classes

The sample measurements of a feature in two classes are

class ω_1	3.5	3.7	3.9	4.1	3.4	3.5	4.1	3.8	3.6	3.7
class ω_2	3.2	3.6	3.1	3.4	3.0	3.4	2.8	3.1	3.3	3.6

Now, we want to know if the feature is informative enough

$$H_1 : \Delta\mu = \mu_1 - \mu_2 \neq 0$$

$$H_0 : \Delta\mu = \mu_1 - \mu_2 = 0$$

Again, we choose $\rho = 0.05$

$$\omega_1 : \bar{x} = 3.73, \hat{\sigma}_1^2 = 0.0601$$

$$\omega_2 : \bar{y} = 3.25, \hat{\sigma}_2^2 = 0.0672$$

Then

For $N = 10$

- $s_z^2 = \frac{1}{2} (\hat{\sigma}_1^2 + \hat{\sigma}_2^2)$
- $q = \frac{(\bar{x} - \bar{y} - 0)}{s_z \sqrt{\frac{2}{N}}}$

Then

For $N = 10$

- $s_z^2 = \frac{1}{2} (\hat{\sigma}_1^2 + \hat{\sigma}_2^2)$
- $q = \frac{(\bar{x} - \bar{y} - 0)}{s_z \sqrt{\frac{2}{N}}}$

We have $q = 4.25$

- We have $20 - 2 = 18$ degrees of freedom and significance level 0.05

Then

For $N = 10$

- $s_z^2 = \frac{1}{2} (\hat{\sigma}_1^2 + \hat{\sigma}_2^2)$
- $q = \frac{(\bar{x} - \bar{y} - 0)}{s_z \sqrt{\frac{2}{N}}}$

We have $q = 4.25$

- We have $20 - 2 = 18$ degrees of freedom and significance level 0.05

Then, $D = [-2.10, 2.10]$

- $q = 4.25$ is outside of D , we decide $H_1 : \Delta\mu = \mu_1 - \mu_2 \neq 0$

Finally

The means μ_1 and μ_2 are significantly different with $\alpha = 0.05$

- The Feature is selected

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- **Considering Feature Sets**
 - Scatter Matrices
 - What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Considering Feature Sets

Something Notable

- The emphasis so far was on individually considered features.

Considering Feature Sets

Something Notable

- The emphasis so far was on individually considered features.

But

- That is, two features may be rich in information, but if they are highly correlated we need not consider both of them.

Considering Feature Sets

Something Notable

- The emphasis so far was on individually considered features.

But

- That is, two features may be rich in information, but if they are highly correlated we need not consider both of them.

Then

- Combine features to search for the “best” combination after features have been discarded.

What to do?

Possible

- Use different feature combinations to form the feature vector.

What to do?

Possible

- Use different feature combinations to form the feature vector.
- Train the classifier, and choose the combination resulting in the best classifier performance.

What to do?

Possible

- Use different feature combinations to form the feature vector.
- Train the classifier, and choose the combination resulting in the best classifier performance.

However

- A major disadvantage of this approach is the high complexity.

What to do?

Possible

- Use different feature combinations to form the feature vector.
- Train the classifier, and choose the combination resulting in the best classifier performance.

However

- A major disadvantage of this approach is the high complexity.
- Also, local minimum may give misleading results.

What to do?

Possible

- Use different feature combinations to form the feature vector.
- Train the classifier, and choose the combination resulting in the best classifier performance.

However

- A major disadvantage of this approach is the high complexity.
- Also, local minimum may give misleading results.

Better

- Adopt a class separability measure and choose the best feature combination against this cost.

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- **Scatter Matrices**
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Scatter Matrices

Definition

- These are used as a measure of the way data are scattered in the respective feature space.

Scatter Matrices

Definition

- These are used as a measure of the way data are scattered in the respective feature space.

Within-class Scatter Matrix

$$S_w = \sum_{i=1}^C P_i S_i \quad (31)$$

- where C is the number of classes.

Scatter Matrices

Definition

- These are used as a measure of the way data are scattered in the respective feature space.

Within-class Scatter Matrix

$$S_w = \sum_{i=1}^C P_i S_i \quad (31)$$

- where C is the number of classes.

where

$$① \quad S_i = E \left[(\mathbf{x} - \boldsymbol{\mu}_i) (\mathbf{x} - \boldsymbol{\mu}_i)^T \right]$$

Scatter Matrices

Definition

- These are used as a measure of the way data are scattered in the respective feature space.

Within-class Scatter Matrix

$$S_w = \sum_{i=1}^C P_i S_i \quad (31)$$

- where C is the number of classes.

where

- $S_i = E \left[(\mathbf{x} - \boldsymbol{\mu}_i) (\mathbf{x} - \boldsymbol{\mu}_i)^T \right]$
- P_i the a priori probability of class ω_i defined as $P_i \cong n_i/N$.

Scatter Matrices

Definition

- These are used as a measure of the way data are scattered in the respective feature space.

Within-class Scatter Matrix

$$S_w = \sum_{i=1}^C P_i S_i \quad (31)$$

- where C is the number of classes.

where

- $S_i = E \left[(\mathbf{x} - \boldsymbol{\mu}_i) (\mathbf{x} - \boldsymbol{\mu}_i)^T \right]$
- P_i the a priori probability of class ω_i defined as $P_i \cong n_i/N$.
 - n_i is the number of samples in class ω_i .

Scatter Matrices

Between-class scatter matrix

$$S_b = \sum_{i=1}^C P_i (\mathbf{x} - \boldsymbol{\mu}_0) (\mathbf{x} - \boldsymbol{\mu}_0)^T \quad (32)$$

Scatter Matrices

Between-class scatter matrix

$$S_b = \sum_{i=1}^C P_i (\mathbf{x} - \boldsymbol{\mu}_0) (\mathbf{x} - \boldsymbol{\mu}_0)^T \quad (32)$$

Where

$$\boldsymbol{\mu}_0 = \sum_{i=1}^C P_i \boldsymbol{\mu}_i \quad (33)$$

The global mean.

Scatter Matrices

Between-class scatter matrix

$$S_b = \sum_{i=1}^C P_i (\mathbf{x} - \boldsymbol{\mu}_0) (\mathbf{x} - \boldsymbol{\mu}_0)^T \quad (32)$$

Where

$$\boldsymbol{\mu}_0 = \sum_{i=1}^C P_i \boldsymbol{\mu}_i \quad (33)$$

The global mean.

Mixture scatter matrix

$$S_m = E \left[(\mathbf{x} - \boldsymbol{\mu}_0) (\mathbf{x} - \boldsymbol{\mu}_0)^T \right] \quad (34)$$

Note: it can be proved that $S_m = S_w + S_b$

Criterion's

First One

$$J_1 = \frac{\text{trace}\{S_m\}}{\text{trace}\{S_w\}} \quad (35)$$

- It takes large values when samples in the d -dimensional space are well clustered around their mean, within each class, and the clusters of the different classes are well separated.

Criterion's

First One

$$J_1 = \frac{\text{trace}\{S_m\}}{\text{trace}\{S_w\}} \quad (35)$$

- It takes large values when samples in the d -dimensional space are well clustered around their mean, within each class, and the clusters of the different classes are well separated.

Other Criteria are

① $J_2 = \frac{|S_m|}{|S_w|}$

Criterion's

First One

$$J_1 = \frac{\text{trace}\{S_m\}}{\text{trace}\{S_w\}} \quad (35)$$

- It takes large values when samples in the d -dimensional space are well clustered around their mean, within each class, and the clusters of the different classes are well separated.

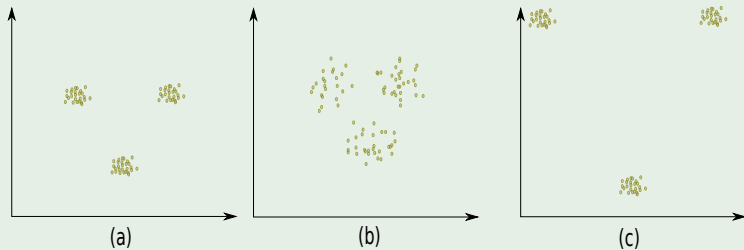
Other Criteria are

- 1 $J_2 = \frac{|S_m|}{|S_w|}$
- 2 $J_3 = \text{trace}\{S_w^{-1}S_m\}$

Example

We have

- Classes with
 - ▶ (a) small within-class variance and small between-class distances,
 - ▶ (b) large within- class variance and small between-class distances,
 - ▶ (c) small within-class variance and large between-class distances.



Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- **What to do with it?**
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

What to do with it

We want to avoid

High Complexities

What to do with it

We want to avoid

High Complexities

As for example

- 1 Select a class separability

What to do with it

We want to avoid

High Complexities

As for example

- 1 Select a class separability
- 2 Then, get all possible combinations of features

$$\binom{m}{l}$$

with $l = 1, 2, \dots, m$

What to do with it

We want to avoid

High Complexities

As for example

- 1 Select a class separability
- 2 Then, get all possible combinations of features

$$\binom{m}{l}$$

with $l = 1, 2, \dots, m$

We can do better

- 1 Sequential Backward Selection

What to do with it

We want to avoid

High Complexities

As for example

- 1 Select a class separability
- 2 Then, get all possible combinations of features

$$\binom{m}{l}$$

with $l = 1, 2, \dots, m$

We can do better

- 1 Sequential Backward Selection
- 2 Sequential Forward Selection

What to do with it

We want to avoid

High Complexities

As for example

- 1 Select a class separability
- 2 Then, get all possible combinations of features

$$\binom{m}{l}$$

with $l = 1, 2, \dots, m$

We can do better

- 1 Sequential Backward Selection
- 2 Sequential Forward Selection
- 3 Floating Search Methods

What to do with it

We want to avoid

High Complexities

As for example

- 1 Select a class separability
- 2 Then, get all possible combinations of features

$$\binom{m}{l}$$

with $l = 1, 2, \dots, m$

We can do better

- 1 Sequential Backward Selection
- 2 Sequential Forward Selection
- 3 Floating Search Methods

However these are sub-optimal methods

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

For example: Sequential Backward Selection

We have the following example

Given x_1, x_2, x_3, x_4 and we wish to select two of them

For example: Sequential Backward Selection

We have the following example

Given x_1, x_2, x_3, x_4 and we wish to select two of them

Step 1

Adopt a class separability criterion, C , and compute its value for the feature vector $[x_1, x_2, x_3, x_4]^T$.

For example: Sequential Backward Selection

We have the following example

Given x_1, x_2, x_3, x_4 and we wish to select two of them

Step 1

Adopt a class separability criterion, C , and compute its value for the feature vector $[x_1, x_2, x_3, x_4]^T$.

Step 2

Eliminate one feature, you get

$$[x_1, x_2, x_3]^T, [x_1, x_2, x_4]^T, [x_1, x_3, x_4]^T, [x_2, x_3, x_4]^T,$$

For example: Sequential Backward Selection

You use your criterion C

Thus the winner is $[x_1, x_2, x_3]^T$

For example: Sequential Backward Selection

You use your criterion C

Thus the winner is $[x_1, x_2, x_3]^T$

Step 3

Now, eliminate a feature and generate $[x_1, x_2]^T, [x_1, x_3]^T, [x_2, x_3]^T,$

For example: Sequential Backward Selection

You use your criterion C

Thus the winner is $[x_1, x_2, x_3]^T$

Step 3

Now, eliminate a feature and generate $[x_1, x_2]^T, [x_1, x_3]^T, [x_2, x_3]^T,$

Use criterion C

To select the best one

Complexity of the Method

Complexity

Thus, starting from m , at each step we drop out one feature from the “best” combination until we obtain a vector of l features.

Complexity of the Method

Complexity

Thus, starting from m , at each step we drop out one feature from the “best” combination until we obtain a vector of l features.

Thus, we need

$1 + 1/2((m + 1)m - l(l + 1))$ combinations

Complexity of the Method

Complexity

Thus, starting from m , at each step we drop out one feature from the “best” combination until we obtain a vector of l features.

Thus, we need

$1 + 1/2((m + 1)m - l(l + 1))$ combinations

However

- The method is sub-optimal

Complexity of the Method

Complexity

Thus, starting from m , at each step we drop out one feature from the “best” combination until we obtain a vector of l features.

Thus, we need

$1 + 1/2((m + 1)m - l(l + 1))$ combinations

However

- The method is sub-optimal
- It suffers of the so called nesting-effect

Complexity of the Method

Complexity

Thus, starting from m , at each step we drop out one feature from the “best” combination until we obtain a vector of l features.

Thus, we need

$1 + 1/2((m + 1)m - l(l + 1))$ combinations

However

- The method is sub-optimal
- It suffers of the so called nesting-effect
 - ▶ Once a feature is discarded, there is no way to reconsider that feature again.

Similar Problem

For

- Sequential Forward Selection

Similar Problem

For

- Sequential Forward Selection

We can overcome this by using

- Floating Search Methods

Similar Problem

For

- Sequential Forward Selection

We can overcome this by using

- Floating Search Methods

A more elegant methods are the ones based on

- Dynamic Programming
- Branch and Bound

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

● Introduction

- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,

Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,
- It possibly produces lower prediction error than the full model.

Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,
- It possibly produces lower prediction error than the full model.

However given process

- it often exhibits high variance,

Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,
- It possibly produces lower prediction error than the full model.

However given process

- it often exhibits high variance,
- It does not reduce the prediction error of the full model.

Shrinkage Methods

By retaining a subset of the predictors and discarding the rest

- Subset Selection produces a model that is interpretable,
- It possibly produces lower prediction error than the full model.

However given process

- it often exhibits high variance,
- It does not reduce the prediction error of the full model.

Therefore

- Shrinkage methods are more continuous avoiding high variability.

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- **Intuition from Overfitting**
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

The house example

Imagine the following data set



Now assume that we use LSE

For the fitting

$$\frac{1}{2} \sum_{i=1}^N (h_{\mathbf{w}}(x_i) - y_i)^2$$

Now assume that we use LSE

For the fitting

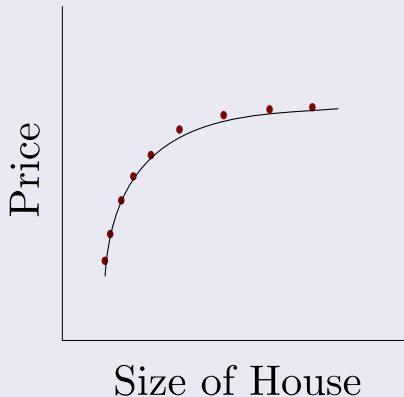
$$\frac{1}{2} \sum_{i=1}^N (h_{\mathbf{w}}(x_i) - y_i)^2$$

We can then run one of our machine to see what minimize better the previous equation

Question: Did you notice that I did not impose any structure to $h_{\mathbf{w}}(x)$?

Then, First fitting

What about using $h_1(x) = w_0 + w_1x + w_2x^2$?



Second fitting

What about using $h_2(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4 + w_5x^5$?



Therefore, we have a problem

We get weird overfitting effects!!!

What do we do? What about minimizing the influence of w_3, w_4, w_5 ?

Therefore, we have a problem

We get weird overfitting effects!!!

What do we do? What about minimizing the influence of w_3, w_4, w_5 ?

How do we do that?

$$\min_w \frac{1}{2} \sum_{i=1}^N (h_w(x_i) - y_i)^2$$

What about integrating those values to the cost function? Ideas

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- **The Idea of Regularization**
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

We have

Regularization intuition is as follow

Small values for parameters $w_0, w_1, w_2, \dots, w_n$

We have

Regularization intuition is as follow

Small values for parameters $w_0, w_1, w_2, \dots, w_n$

It implies

- 1 "Simpler" function
- 2 Less prone to overfitting

We can do the previous idea for the other parameters

We can do the same for the other parameters

$$\min_w \frac{1}{2} \sum_{i=1}^N (h_w(x_i) - y_i)^2 + \sum_{i=1}^d \lambda_i w_i^2 \quad (36)$$

We can do the previous idea for the other parameters

We can do the same for the other parameters

$$\min_w \frac{1}{2} \sum_{i=1}^N (h_w(x_i) - y_i)^2 + \sum_{i=1}^d \lambda_i w_i^2 \quad (36)$$

However handling such many parameters can be so difficult

Combinatorial problem in reality!!!

Better, we can

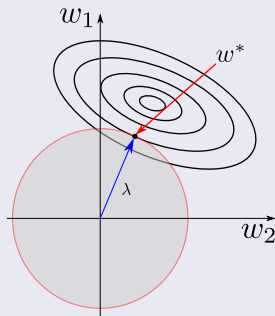
We better use the following

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (h_{\mathbf{w}}(x_i) - y_i)^2 + \lambda \sum_{i=1}^d w_i^2 \quad (37)$$

Graphically

Geometrically Equivalent to

$$\sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \lambda \sum_{i=1}^{d+1} w_i^2$$



Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- **Ridge Regression**
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

Ridge Regression

Equation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + \lambda \sum_{j=1}^d w_j^2 \right\}$$

Ridge Regression

Equation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + \lambda \sum_{j=1}^d w_j^2 \right\}$$

Here

- $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage

Therefore

The Larger $\lambda \geq 0$

- The coefficients are shrunk toward zero (and each other).

Therefore

The Larger $\lambda \geq 0$

- The coefficients are shrunk toward zero (and each other).

This is also used in Neural Networks

- where it is known as weight decay

This is also can be written

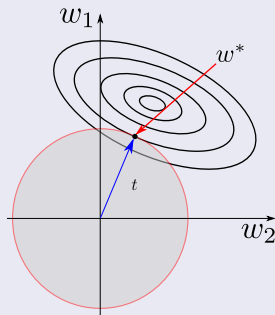
Optimization Solution

$$\begin{aligned} & \arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 \\ & \text{subject to } \sum_{j=1}^d w_j^2 < t \end{aligned}$$

Graphically

Geometrically Equivalent to

$$\begin{aligned} \arg \min \quad & \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 \\ \text{subject to} \quad & \sum_{i=1}^{d+1} w_i^2 < t \end{aligned}$$



Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- **Standardization of Data**
- The LASSO
 - The Lagrangian Version of the LASSO

Remarks

We have the following

- The Ridge solutions are not equivariant under scaling of the inputs.

Remarks

We have the following

- The Ridge solutions are not equivariant under scaling of the inputs.

Thus, the need to standardize the input data

- Before Solving the optimization:

$$\arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{ij} w_j \right)^2$$

subject to $\sum_{j=1}^d w_j^2 < t$

Here

Notice that w_0 is not being penalized

- Penalizing w_0 would make the procedure depend on the origin chosen for y_i .

Therefore

We can center the Data

- Thus, each x_{ij} gets replaced by $x_{ij} - \bar{x}_j$.

Therefore

We can center the Data

- Thus, each x_{ij} gets replaced by $x_{ij} - \bar{x}_j$.

Then, we estimate w_0

$$w_0 = \frac{1}{N} \sum_{i=1}^N y_i$$

Thus after centering the Data

Now, Given a data matrix \mathbf{X} with d dimensions

$$Loss_{RSS}(\lambda) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$$

Thus after centering the Data

Now, Given a data matrix \mathbf{X} with d dimensions

$$Loss_{RSS}(\lambda) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$$

The Ridge Regression solution is equivalent to

$$\hat{\mathbf{w}}^{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}.$$

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- **The LASSO**
 - The Lagrangian Version of the LASSO

Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\mathbf{w}}^{garrote} = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + N\lambda \sum_{j=1}^d w_j$$

s.t. $w_j > 0 \ \forall j$

Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\mathbf{w}}^{garrote} = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + N\lambda \sum_{j=1}^d w_j$$

s.t. $w_j > 0 \ \forall j$

This is quite derivable

However, Tibshirani realized that you could get a more flexible model by using the absolute value at the constraint!!!

Least Absolute Shrinkage and Selection Operator (LASSO)

It was introduced by Robert Tibshirani in 1996 based on Leo Breiman's nonnegative garrote

$$\hat{\mathbf{w}}^{garrote} = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 + N \lambda \sum_{j=1}^d w_j$$

s.t. $w_j > 0 \ \forall j$

This is quite derivable

However, Tibshirani realized that you could get a more flexible model by using the absolute value at the constraint!!!

Robert Tibshirani proposed the use of the L_1 norm

$$\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|$$

The Final Optimization Problem

LASSO

$$\begin{aligned}\hat{\mathbf{w}}^{LASSO} &= \arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2 \\ \text{s.t. } &\sum_{i=1}^d |w_i| \leq t\end{aligned}$$

The Final Optimization Problem

LASSO

$$\hat{\mathbf{w}}^{LASSO} = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^d x_{ij} w_j \right)^2$$
$$\text{s.t. } \sum_{i=1}^d |w_i| \leq t$$

This is not derivable

More advanced methods are necessary to solve this problem!!!

Outline

1

Introduction

- Feature Engineering
- What is Feature Selection?
- Preprocessing
 - Outlier Removal
 - Finding Multivariate Outliers
 - Data Normalization
 - Methods
- Missing Data
 - Using EM
 - Matrix Completion
- The Peaking Phenomena

2

Feature Selection

- Feature Selection
- Feature selection based on statistical hypothesis testing
 - Example
- Application of the t -Test in Feature Selection
 - Example
- Considering Feature Sets
- Scatter Matrices
- What to do with it?
 - Sequential Backward Selection

3

Shrinkage Methods

- Introduction
- Intuition from Overfitting
- The Idea of Regularization
- Ridge Regression
- Standardization of Data
- The LASSO
 - The Lagrangian Version of the LASSO

The Lagrangian Version

The Lagrangian

$$\hat{\mathbf{w}}^{LASSO} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left(y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i| \right\}$$

The Lagrangian Version

The Lagrangian

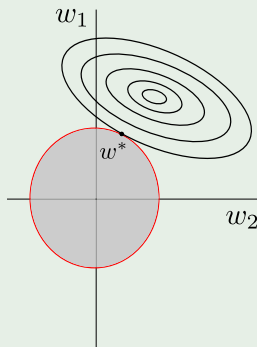
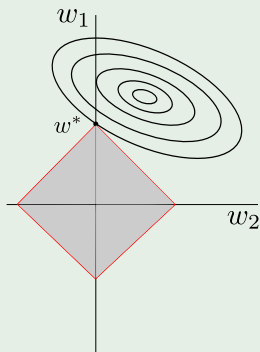
$$\hat{\mathbf{w}}^{LASSO} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left(y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i| \right\}$$

However

You have other regularizations as $\|\mathbf{w}\|_2 = \sqrt{\sum_{i=1}^d |w_i|^2}$

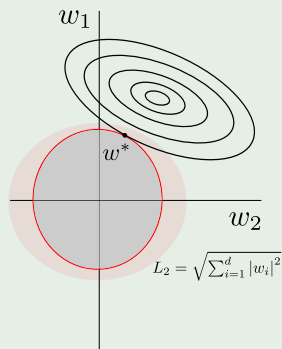
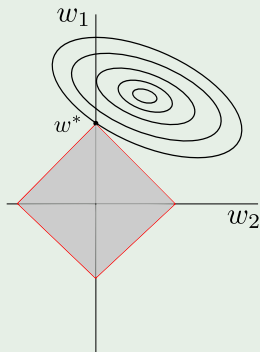
Graphically

The first area correspond to the L_1 regularization and the second one?



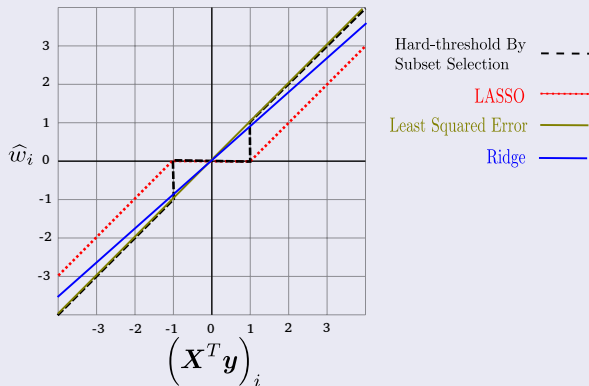
Graphically

Yes the circle defined as $\|w\|_2 = \sqrt{\sum_{i=1}^d |w_i|^2}$



For Example

In the Case of X is a Orthogonal Matrix



The seminal paper by Robert Tibshirani

An initial study of this regularization can be seen in

“Regression Shrinkage and Selection via the LASSO” by Robert Tibshirani
- 1996

This out the scope of this class

However, it is worth noticing that the most efficient method for solving LASSO problems is

“Pathwise Coordinate Optimization” By Jerome Friedman, Trevor Hastie, Holger Ho and Robert Tibshirani

This out the scope of this class

However, it is worth noticing that the most efficient method for solving LASSO problems is

“Pathwise Coordinate Optimization” By Jerome Friedman, Trevor Hastie, Holger Ho and Robert Tibshirani

Nevertheless

It will be a great seminar paper!!!

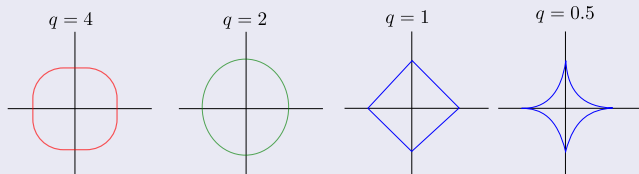
Furthermore

We can generalize ridge regression and the lasso, and view them as Bayes estimates

$$\hat{\mathbf{w}}^{LASSO} = \arg \min_{\mathbf{w}} \left\{ \sum_{i=1}^N \left(y_i - \mathbf{x}^T \mathbf{w} \right)^2 + \lambda \sum_{i=1}^d |w_i|^q \right\} \text{ with } q \geq 0$$

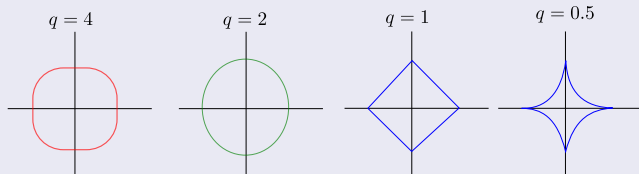
For Example

We have when $d = 2$



For Example

We have when $d = 2$



Here, when $q > 1$

- You are having a derivable Lagrangian, but you lose the LASSO properties

Therefore

Zou and Hastie (2005) introduced the elastic- net penalty

$$\lambda \sum_{i=1}^d \left\{ \alpha w_i^2 + (1 - \alpha) |w_i| \right\}$$

Therefore

Zou and Hastie (2005) introduced the elastic- net penalty

$$\lambda \sum_{i=1}^d \left\{ \alpha w_i^2 + (1 - \alpha) |w_i| \right\}$$

This is Basically

- A Compromise Between the Ridge and LASSO.