

# Introduction to Machine Learning

## Logistic Regression

Andres Mendez-Vazquez

January 26, 2023

# Outline

## 1 Logistic Regression

- Introduction
- Constraints
- The Initial Model
- The Two Case Class
- Graphic Interpretation
- Fitting The Model
  - The Two Class Case
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

# Outline

## 1 Logistic Regression

### ● Introduction

- Constraints
- The Initial Model
- The Two Case Class
- Graphic Interpretation
- Fitting The Model
  - The Two Class Case
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

Assume the following

Let  $Y_1, Y_2, \dots, Y_N$  independent random variables

Taking values in the set  $\{0, 1\}$

Now, you have a set of fixed vectors

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$$

Mapped to a series of numbers by a weight vector  $\mathbf{w}$

$$\mathbf{w}^T \mathbf{x}_1, \mathbf{w}^T \mathbf{x}_2, \dots, \mathbf{w}^T \mathbf{x}_N$$

Assume the following

Let  $Y_1, Y_2, \dots, Y_N$  independent random variables

Taking values in the set  $\{0, 1\}$

Now, you have a set of fixed vectors

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$$

Mapped to a series of numbers by a weight vector  $w$

$$w^T \mathbf{x}_1, w^T \mathbf{x}_2, \dots, w^T \mathbf{x}_N$$

Assume the following

Let  $Y_1, Y_2, \dots, Y_N$  independent random variables

Taking values in the set  $\{0, 1\}$

Now, you have a set of fixed vectors

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$$

Mapped to a series of numbers by a weight vector  $\mathbf{w}$

$$\mathbf{w}^T \mathbf{x}_1, \mathbf{w}^T \mathbf{x}_2, \dots, \mathbf{w}^T \mathbf{x}_N$$

## In our simplest form [1, 2]

There is a suspected relation

Between  $\theta_i = P(Y_i = 1)$  and  $w^T x_i$

- Here  $Y$  is the random variable and  $y$  is the value that the random variable can take.

Thus we have

$$y = \begin{cases} 1 & w^T x + e > 0 \\ 0 & \text{else} \end{cases}$$

Note: Where  $e$  is an error with a certain distribution!!!

## In our simplest form [1, 2]

There is a suspected relation

Between  $\theta_i = P(Y_i = 1)$  and  $\mathbf{w}^T \mathbf{x}_i$

- Here  $Y$  is the random variable and  $y$  is the value that the random variable can take.

Thus we have

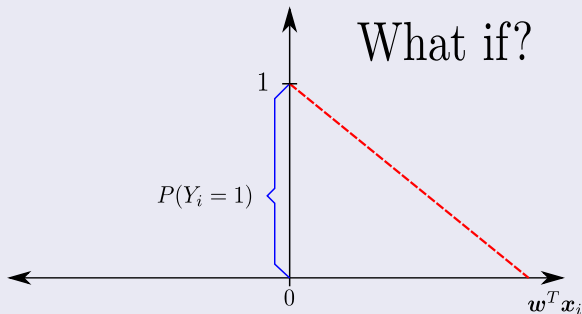
$$y = \begin{cases} 1 & \mathbf{w}^T \mathbf{x} + e > 0 \\ 0 & \text{else} \end{cases}$$

**Note:** Where  $e$  is an error with a certain distribution!!!



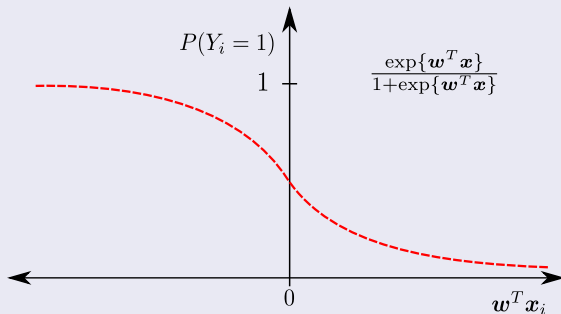
## For Example, Graphically

We have



It is better to use a logit version

We have



# Logit Distribution

PDF with support  $z \in (-\infty, \infty)$ ,  $\mu$  location and  $s$  scale

$$p(x|\mu, s) = \frac{\exp\left\{-\frac{z-\mu}{s}\right\}}{s\left(1 + \exp\left\{-\frac{z-\mu}{s}\right\}\right)^2}$$

With a CDF

$$P(Y < z) = \int_{-\infty}^z p(y|\mu, s) dy = \frac{1}{1 + \exp\left\{-\frac{z-\mu}{s}\right\}}$$

# Logit Distribution

PDF with support  $z \in (-\infty, \infty)$ ,  $\mu$  location and  $s$  scale

$$p(x|\mu, s) = \frac{\exp\left\{-\frac{z-\mu}{s}\right\}}{s\left(1 + \exp\left\{-\frac{z-\mu}{s}\right\}\right)^2}$$

With a CDF

$$P(Y < z) = \int_{-\infty}^z p(y|\mu, s) dy = \frac{1}{1 + \exp\left\{-\frac{z-\mu}{s}\right\}}$$

# Outline

## 1 Logistic Regression

- Introduction
- **Constraints**
- The Initial Model
- The Two Case Class
- Graphic Interpretation
- Fitting The Model
  - The Two Class Case
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

# In Bayesian Classification

## Assignment of a pattern

It is performed by using the posterior probabilities,  $P(\omega_i|\mathbf{x})$

And given  $K$  classes, we want:

$$\sum_{i=1}^K P(\omega_i|\mathbf{x}) = 1$$

Such that each

$$0 \leq P(\omega_i|\mathbf{x}) \leq 1$$

# In Bayesian Classification

## Assignment of a pattern

It is performed by using the posterior probabilities,  $P(\omega_i|\mathbf{x})$

And given  $K$  classes, we want

$$\sum_{i=1}^K P(\omega_i|\mathbf{x}) = 1$$

Such that each

$$0 \leq P(\omega_i|\mathbf{x}) \leq 1$$

# In Bayesian Classification

## Assignment of a pattern

It is performed by using the posterior probabilities,  $P(\omega_i|\mathbf{x})$

And given  $K$  classes, we want

$$\sum_{i=1}^K P(\omega_i|\mathbf{x}) = 1$$

Such that each

$$0 \leq P(\omega_i|\mathbf{x}) \leq 1$$



# Observation

This is a typical example of the discriminative approach

- Where the distribution of data is of no interest.
  - ▶ In the Logistic Regression the Distribution is imposed over the output!!!

# Outline

## 1 Logistic Regression

- Introduction
- Constraints
- **The Initial Model**
- The Two Case Class
- Graphic Interpretation
- Fitting The Model
  - The Two Class Case
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

# The Model

We have the following under the extended features

$$\log \frac{P(\omega_1|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_1^T \mathbf{x}$$

$$\log \frac{P(\omega_2|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_2^T \mathbf{x}$$

$$\vdots$$

$$\log \frac{P(\omega_{K-1}|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_{K-1}^T \mathbf{x}$$

# The Model

We have the following under the extended features

$$\log \frac{P(\omega_1|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_1^T \mathbf{x}$$

$$\log \frac{P(\omega_2|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_2^T \mathbf{x}$$

$$\vdots$$

$$\log \frac{P(\omega_{K-1}|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_{K-1}^T \mathbf{x}$$

# The Model

We have the following under the extended features

$$\log \frac{P(\omega_1|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_1^T \mathbf{x}$$

$$\log \frac{P(\omega_2|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_2^T \mathbf{x}$$

$$\vdots$$

$$\log \frac{P(\omega_{K-1}|\mathbf{x})}{P(\omega_K|\mathbf{x})} = \mathbf{w}_{K-1}^T \mathbf{x}$$

## Further

We have

The model is specified in terms of  $K - 1$  log-odds or logit transformations.

And

Although the model uses the last class as the denominator in the odds-ratios.

The choice of denominator is arbitrary.

- However, because the estimates are equivariant under this choice.
  - ▶ The action taken in a decision problem should not depend on transformation on the measurement used

## Further

We have

The model is specified in terms of  $K - 1$  log-odds or logit transformations.

And

Although the model uses the last class as the denominator in the odds-ratios.

The choice of denominator is arbitrary.

- However, because the estimates are equivariant under this choice.
  - ▶ The action taken in a decision problem should not depend on transformation on the measurement used

## Further

We have

The model is specified in terms of  $K - 1$  log-odds or logit transformations.

And

Although the model uses the last class as the denominator in the odds-ratios.

The choice of denominator is arbitrary

- However, because the estimates are equivariant under this choice.
  - ▶ The action taken in a decision problem should not depend on transformation on the measurement used



Now

How do we find the terms?

$$P(\omega_1|\mathbf{x}), P(\omega_2|\mathbf{x}), \dots, P(\omega_K|\mathbf{x})$$

It is possible to show that

We have that, for  $l = 1, 2, \dots, K - 1$

$$\frac{P(\omega_l | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \exp \left\{ \mathbf{w}_l^T \mathbf{x} \right\}$$

Therefore

$$\sum_{l=1}^{K-1} \frac{P(\omega_l | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \sum_{l=1}^{K-1} \exp \left\{ \mathbf{w}_l^T \mathbf{x} \right\}$$

Thus

$$\frac{1 - P(\omega_K | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \sum_{l=1}^{K-1} \exp \left\{ \mathbf{w}_l^T \mathbf{x} \right\}$$

It is possible to show that

We have that, for  $l = 1, 2, \dots, K - 1$

$$\frac{P(\omega_l | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \exp \{ \mathbf{w}_l^T \mathbf{x} \}$$

Therefore

$$\sum_{l=1}^{K-1} \frac{P(\omega_l | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \sum_{l=1}^{K-1} \exp \{ \mathbf{w}_l^T \mathbf{x} \}$$

Thus

$$\frac{1 - P(\omega_K | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \sum_{l=1}^{K-1} \exp \{ \mathbf{w}_l^T \mathbf{x} \}$$

It is possible to show that

We have that, for  $l = 1, 2, \dots, K - 1$

$$\frac{P(\omega_l | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \exp \{ \mathbf{w}_l^T \mathbf{x} \}$$

Therefore

$$\sum_{l=1}^{K-1} \frac{P(\omega_l | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \sum_{l=1}^{K-1} \exp \{ \mathbf{w}_l^T \mathbf{x} \}$$

Thus

$$\frac{1 - P(\omega_K | \mathbf{x})}{P(\omega_K | \mathbf{x})} = \sum_{l=1}^{K-1} \exp \{ \mathbf{w}_l^T \mathbf{x} \}$$

# Basically

We have, (Take a look at the board)

$$P(\omega_K | \mathbf{x}) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp \{ \mathbf{w}_l^T \mathbf{x} \}}$$

Then

$$\frac{P(\omega_i | \mathbf{x})}{\frac{1}{1 + \sum_{l=1}^{K-1} \exp \{ \mathbf{w}_l^T \mathbf{x} \}}} = \exp \{ \mathbf{w}_i^T \mathbf{x} \}$$

For  $i = 1, 2, \dots, K-1$

$$P(\omega_i | \mathbf{x}) = \frac{\exp \{ \mathbf{w}_i^T \mathbf{x} \}}{1 + \sum_{l=1}^{K-1} \exp \{ \mathbf{w}_l^T \mathbf{x} \}}$$

# Basically

We have, (Take a look a the board)

$$P(\omega_K|\mathbf{x}) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp\{\mathbf{w}_l^T \mathbf{x}\}}$$

Then

$$\frac{P(\omega_i|\mathbf{x})}{\frac{1}{1 + \sum_{l=1}^{K-1} \exp\{\mathbf{w}_l^T \mathbf{x}\}}} = \exp\{\mathbf{w}_i^T \mathbf{x}\}$$

For  $i=0, 1, \dots, K-1$

$$P(\omega_i|\mathbf{x}) = \frac{\exp\{\mathbf{w}_i^T \mathbf{x}\}}{1 + \sum_{l=1}^{K-1} \exp\{\mathbf{w}_l^T \mathbf{x}\}}$$

# Basically

We have, (Take a look a the board)

$$P(\omega_K|\mathbf{x}) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp\{\mathbf{w}_l^T \mathbf{x}\}}$$

Then

$$\frac{P(\omega_i|\mathbf{x})}{\frac{1}{1 + \sum_{l=1}^{K-1} \exp\{\mathbf{w}_l^T \mathbf{x}\}}} = \exp\{\mathbf{w}_i^T \mathbf{x}\}$$

For  $i = 1, 2, \dots, k - 1$

$$P(\omega_i|\mathbf{x}) = \frac{\exp\{\mathbf{w}_i^T \mathbf{x}\}}{1 + \sum_{l=1}^{K-1} \exp\{\mathbf{w}_l^T \mathbf{x}\}}$$

Additionally

For  $K$

$$P(\omega_K | \mathbf{x}) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp \{ \mathbf{w}_l^T \mathbf{x} \}}$$

Easy to see

They sum to one.



Additionally

For  $K$

$$P(\omega_K|\mathbf{x}) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp\{\mathbf{w}_l^T \mathbf{x}\}}$$

Easy to see

They sum to one.

## A Note in Notation

Given all these parameters, we summarized them

$$\Theta = \{w_1, w_2, \dots, w_{K-1}\}$$

Therefore

$$P(w_l | X = x) = p_l(x | \Theta)$$

## A Note in Notation

Given all these parameters, we summarized them

$$\Theta = \{\boldsymbol{w}_1, \boldsymbol{w}_2, \dots, \boldsymbol{w}_{K-1}\}$$

Therefore

$$P(\omega_l | \boldsymbol{X} = \boldsymbol{x}) = p_l(\boldsymbol{x} | \Theta)$$

# Outline

## 1 Logistic Regression

- Introduction
- Constraints
- The Initial Model
- **The Two Case Class**
- Graphic Interpretation
- Fitting The Model
  - The Two Class Case
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

## In the two class case

We have

$$P_1(\omega_1|\mathbf{x}) = \frac{\exp\{\mathbf{w}^T \mathbf{x}\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}}$$
$$P_2(\omega_2|\mathbf{x}) = \frac{1}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}}$$

A similar model

$$P_1(\omega_1|\mathbf{x}) = \frac{\exp\{-\mathbf{w}^T \mathbf{x}\}}{1 + \exp\{-\mathbf{w}^T \mathbf{x}\}}$$
$$P_2(\omega_2|\mathbf{x}) = \frac{1}{1 + \exp\{-\mathbf{w}^T \mathbf{x}\}}$$

## In the two class case

We have

$$P_1(\omega_1|\mathbf{x}) = \frac{\exp\{\mathbf{w}^T \mathbf{x}\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}}$$
$$P_2(\omega_2|\mathbf{x}) = \frac{1}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}}$$

A similar model

$$P_1(\omega_1|\mathbf{x}) = \frac{\exp\{-\mathbf{w}^T \mathbf{x}\}}{1 + \exp\{-\mathbf{w}^T \mathbf{x}\}}$$
$$P_2(\omega_2|\mathbf{x}) = \frac{1}{1 + \exp\{-\mathbf{w}^T \mathbf{x}\}}$$

# Outline

## 1 Logistic Regression

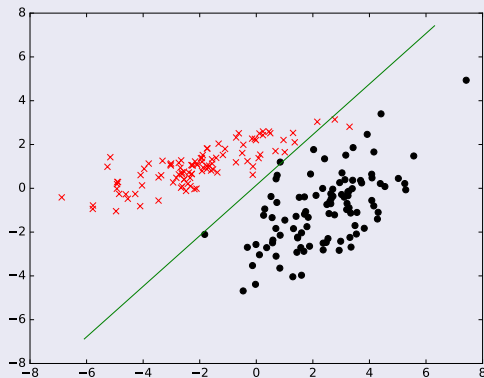
- Introduction
- Constraints
- The Initial Model
- The Two Case Class
- **Graphic Interpretation**
- Fitting The Model
  - The Two Class Case
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

We have the following split

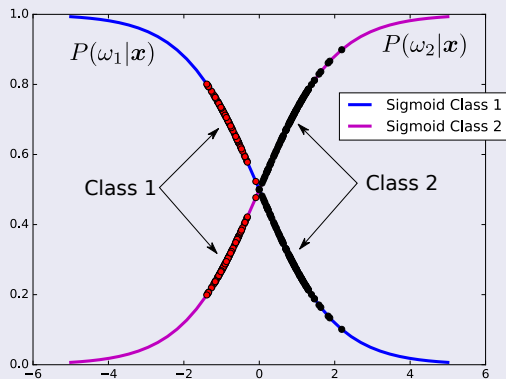
Using  $f(x) = w^T x$





Then, we have the mapping to

We have



# Outline

## 1 Logistic Regression

- Introduction
- Constraints
- The Initial Model
- The Two Case Class
- Graphic Interpretation
- **Fitting The Model**
  - The Two Class Case
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

# A Classic application of Maximum Likelihood

Given a sequence of samples iid

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$$

We have the following pdf

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \Theta) = \prod_{i=1}^N p(\mathbf{x}_i | \Theta)$$

# A Classic application of Maximum Likelihood

Given a sequence of samples iid

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$$

We have the following pdf

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \Theta) = \prod_{i=1}^N p(\mathbf{x}_i | \Theta)$$

## $P(g = v|X)$ Distribution

Where  $P(g = v|X)$  completely specify the conditional distribution

We have a multinomial distribution which under the log-likelihood of  $N$  observations:

$$\mathcal{L}(\Theta) = \log p(x_1, x_2, \dots, x_N | \Theta) = \log \prod_{i=1}^N p_{g_i}(x_i | \theta) = \sum_{i=1}^N \log p_{g_i}(x_i | \theta)$$

Where

$g_i$  represent the class that  $x_i$  belongs.

$$g_i = \begin{cases} 1 & \text{if } x_i \in \text{Class 1} \\ 2 & \text{if } x_i \in \text{Class 2} \end{cases}$$

## $P(g = v|X)$ Distribution

Where  $P(g = v|X)$  completely specify the conditional distribution

We have a multinomial distribution which under the log-likelihood of  $N$  observations:

$$\mathcal{L}(\Theta) = \log p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \Theta) = \log \prod_{i=1}^N p_{g_i}(\mathbf{x}_i | \theta) = \sum_{i=1}^N \log p_{g_i}(\mathbf{x}_i | \theta)$$

Where

$g_i$  represent the class that  $\mathbf{x}_i$  belongs.

$$g_i = \begin{cases} 1 & \text{if } \mathbf{x}_i \in \text{Class 1} \\ 2 & \text{if } \mathbf{x}_i \in \text{Class 2} \end{cases}$$

## $P(g = v|X)$ Distribution

Where  $P(g = v|X)$  completely specify the conditional distribution

We have a multinomial distribution which under the log-likelihood of  $N$  observations:

$$\mathcal{L}(\Theta) = \log p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \Theta) = \log \prod_{i=1}^N p_{g_i}(\mathbf{x}_i | \theta) = \sum_{i=1}^N \log p_{g_i}(\mathbf{x}_i | \theta)$$

Where

$g_i$  represent the class that  $\mathbf{x}_i$  belongs.

$$g_i = \begin{cases} 1 & \text{if } \mathbf{x}_i \in \text{Class 1} \\ 2 & \text{if } \mathbf{x}_i \in \text{Class 2} \end{cases}$$

## $P(g = v|X)$ Distribution

Where  $P(g = v|X)$  completely specify the conditional distribution

We have a multinomial distribution which under the log-likelihood of  $N$  observations:

$$\mathcal{L}(\Theta) = \log p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \Theta) = \log \prod_{i=1}^N p_{g_i}(\mathbf{x}_i | \theta) = \sum_{i=1}^N \log p_{g_i}(\mathbf{x}_i | \theta)$$

Where

$g_i$  represent the class that  $\mathbf{x}_i$  belongs.

$$g_i = \begin{cases} 1 & \text{if } \mathbf{x}_i \in \text{Class 1} \\ 2 & \text{if } \mathbf{x}_i \in \text{Class 2} \end{cases}$$



## $P(g = v|X)$ Distribution

Where  $P(g = v|X)$  completely specify the conditional distribution

We have a multinomial distribution which under the log-likelihood of  $N$  observations:

$$\mathcal{L}(\Theta) = \log p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \Theta) = \log \prod_{i=1}^N p_{g_i}(\mathbf{x}_i | \theta) = \sum_{i=1}^N \log p_{g_i}(\mathbf{x}_i | \theta)$$

Where

$g_i$  represent the class that  $\mathbf{x}_i$  belongs.

$$g_i = \begin{cases} 1 & \text{if } \mathbf{x}_i \in \text{Class 1} \\ 2 & \text{if } \mathbf{x}_i \in \text{Class 2} \end{cases}$$

# How do we integrate this into a Cost Function?

Clearly, we have two distributions

We need to represent the distributions into the functions  $p_{g_i}(\mathbf{x}_i|\theta)$ .

Why not to have all the distributions into this function

$$p_{g_i}(\mathbf{x}_i|\theta) = \prod_{l=1}^{K-1} p(\mathbf{x}_i|\omega_l)^{I\{\mathbf{x}_i \in \omega_l\}}$$

It is easy with the two classes

Given that we have a binary situation!!!

# How do we integrate this into a Cost Function?

Clearly, we have two distributions

We need to represent the distributions into the functions  $p_{g_i}(\mathbf{x}_i|\theta)$ .

Why not to have all the distributions into this function

$$p_{g_i}(\mathbf{x}_i|\theta) = \prod_{l=1}^{K-1} p(\mathbf{x}_i|\mathbf{w}_l)^{I\{\mathbf{x}_i \in \omega_l\}}$$

It is easy with the two classes

Given that we have a binary situation!!!

# How do we integrate this into a Cost Function?

Clearly, we have two distributions

We need to represent the distributions into the functions  $p_{g_i}(\mathbf{x}_i|\theta)$ .

Why not to have all the distributions into this function

$$p_{g_i}(\mathbf{x}_i|\theta) = \prod_{l=1}^{K-1} p(\mathbf{x}_i|\mathbf{w}_l)^{I\{\mathbf{x}_i \in \omega_l\}}$$

It is easy with the two classes

Given that we have a binary situation!!!

# Outline

## 1 Logistic Regression

- Introduction
- Constraints
- The Initial Model
- The Two Case Class
- Graphic Interpretation
- **Fitting The Model**
  - **The Two Class Case**
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

Given the two case

We have then a Bernoulli distribution

$$p_1(\mathbf{x}_i|\mathbf{w}) = \left[ \frac{\exp\{\mathbf{w}^T \mathbf{x}\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}} \right]^{y_i}$$
$$p_2(\mathbf{x}_i|\mathbf{w}) = \left[ \frac{1}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}} \right]^{1-y_i}$$

With

$y_i = 1$  if  $\mathbf{x}_i \in \text{Class 1}$

$y_i = 0$  if  $\mathbf{x}_i \in \text{Class 2}$

Given the two case

We have then a Bernoulli distribution

$$p_1(\mathbf{x}_i|\mathbf{w}) = \left[ \frac{\exp\{\mathbf{w}^T \mathbf{x}\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}} \right]^{y_i}$$
$$p_2(\mathbf{x}_i|\mathbf{w}) = \left[ \frac{1}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}} \right]^{1-y_i}$$

With

$y_i = 1$  if  $\mathbf{x}_i \in \text{Class 1}$

$y_i = 0$  if  $\mathbf{x}_i \in \text{Class 2}$

# Outline

## 1 Logistic Regression

- Introduction
- Constraints
- The Initial Model
- The Two Case Class
- Graphic Interpretation
- **Fitting The Model**
  - The Two Class Case
  - **The Final Log-Likelihood**
  - The Newton-Raphson Algorithm
  - Matrix Notation

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion



We have the following

## Cost Function

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N \{y_i \log p_1(\mathbf{x}_i|\mathbf{w}) + (1 - y_i) \log (1 - p_1(\mathbf{x}_i|\mathbf{w}))\}$$

After some reductions

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N \left\{ y_i \mathbf{w}^T \mathbf{x}_i - \log \left( 1 + \exp \left\{ \mathbf{w}^T \mathbf{x}_i \right\} \right) \right\}$$

We have the following

### Cost Function

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N \{y_i \log p_1(\mathbf{x}_i|\mathbf{w}) + (1 - y_i) \log (1 - p_1(\mathbf{x}_i|\mathbf{w}))\}$$

### After some reductions

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N \left\{ y_i \mathbf{w}^T \mathbf{x}_i - \log \left( 1 + \exp \left\{ \mathbf{w}^T \mathbf{x}_i \right\} \right) \right\}$$

Now, we derive and set it to zero

We have

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^N \mathbf{x}_i \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) = 0$$

Which are  $d+1$  equations nonlinear

$$\sum_{i=1}^N \mathbf{x}_i (y_i - p(\mathbf{x}_i | \mathbf{w})) = \begin{pmatrix} \sum_{i=1}^N 1 \times \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) \\ \sum_{i=1}^N x_1^i \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) \\ \sum_{i=1}^N x_2^i \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) \\ \vdots \\ \sum_{i=1}^N x_d^i \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) \end{pmatrix} = 0$$

It is known as a scoring function.

Now, we derive and set it to zero

We have

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^N \mathbf{x}_i \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) = 0$$

Which are  $d + 1$  equations nonlinear

$$\sum_{i=1}^N \mathbf{x}_i (y_i - p(\mathbf{x}_i | \mathbf{w})) = \begin{pmatrix} \sum_{i=1}^N 1 \times \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) \\ \sum_{i=1}^N x_1^i \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) \\ \sum_{i=1}^N x_2^i \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) \\ \vdots \\ \sum_{i=1}^N x_d^i \left( y_i - \frac{\exp\{\mathbf{w}^T \mathbf{x}_i\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}_i\}} \right) \end{pmatrix} = 0$$

It is known as a scoring function.

# Finally

In other words you

①  $d + 1$  nonlinear equations in  $w$ .

② For example, from the first equation:

$$\sum_{i=1}^N y_i = \sum_{i=1}^N p(x_i | w)$$

- ▶ The expected number of class ones matches the observed number.
- ▶ And hence also class two.

# Finally

In other words you

- 1  $d + 1$  nonlinear equations in  $w$ .
- 2 For example, from the first equation:

$$\sum_{i=1}^N y_i = \sum_{i=1}^N p(x_i | w)$$

- ▶ The expected number of class ones matches the observed number.
- ▶ And hence also class twos.

# Finally

## In other words you

- 1  $d + 1$  nonlinear equations in  $\mathbf{w}$ .
- 2 For example, from the first equation:

$$\sum_{i=1}^N y_i = \sum_{i=1}^N p(\mathbf{x}_i | \mathbf{w})$$

- The expected number of class ones matches the observed number.

► And hence also class twos.

# Finally

## In other words you

- 1  $d + 1$  nonlinear equations in  $\mathbf{w}$ .
- 2 For example, from the first equation:

$$\sum_{i=1}^N y_i = \sum_{i=1}^N p(\mathbf{x}_i | \mathbf{w})$$

- ▶ The expected number of class ones matches the observed number.
- ▶ And hence also class twos.



# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - **Fitting The Model**
    - The Two Class Case
    - The Final Log-Likelihood
    - **The Newton-Raphson Algorithm**
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

## To solve the previous equations [3]

We use the Newton-Raphson Method to find the roots or zeros

It comes from the first Taylor Approximation

$$f(x+h) \approx f(x) + hf'(x)$$

Thus we have for a root  $r$  of function  $f$

We have

- Assume a good estimate of  $r$ ,  $x_0$
- Thus we have  $r = x_0 + h$
- Or  $h = r - x_0$

## To solve the previous equations [3]

We use the Newton-Raphson Method to find the roots or zeros

It comes from the first Taylor Approximation

$$f(x + h) \approx f(x) + hf'(x)$$

Thus we have for a root  $r$  of function  $f$

We have

- 1 Assume a good estimate of  $r$ ,  $x_0$

2 Thus we have  $r = x_0 + h$

3 Or  $h = r - x_0$

## To solve the previous equations [3]

We use the Newton-Raphson Method to find the roots or zeros

It comes from the first Taylor Approximation

$$f(x + h) \approx f(x) + hf'(x)$$

Thus we have for a root  $r$  of function  $f$

We have

- 1 Assume a good estimate of  $r$ ,  $x_0$

2 Thus we have  $r = x_0 + h$

3 Or  $h = r - x_0$

## To solve the previous equations [3]

We use the Newton-Raphson Method to find the roots or zeros

It comes from the first Taylor Approximation

$$f(x+h) \approx f(x) + hf'(x)$$

Thus we have for a root  $r$  of function  $f$

We have

- 1 Assume a good estimate of  $r$ ,  $x_0$
- 2 Thus we have  $r = x_0 + h$

Or  $h = r - x_0$

## To solve the previous equations [3]

We use the Newton-Raphson Method to find the roots or zeros

It comes from the first Taylor Approximation

$$f(x + h) \approx f(x) + hf'(x)$$

Thus we have for a root  $r$  of function  $f$

We have

- ① Assume a good estimate of  $r$ ,  $x_0$
- ② Thus we have  $r = x_0 + h$
- ③ Or  $h = r - x_0$

We have then

From Taylor

$$0 = f(r) = f(x_0 + h) \approx f(x_0) + hf'(x_0)$$

Thus, as long as  $f'(x_0)$  is not close to 0

$$h \approx -\frac{f(x_0)}{f'(x_0)}$$

Thus

$$r = x_0 + h \approx x_0 - \frac{f(x_0)}{f'(x_0)}$$

We have then

From Taylor

$$0 = f(r) = f(x_0 + h) \approx f(x_0) + hf'(x_0)$$

Thus, as long  $f'(x_0)$  is not close to 0

$$h \approx -\frac{f(x_0)}{f'(x_0)}$$

Thus

$$r = x_0 + h \approx x_0 - \frac{f(x_0)}{f'(x_0)}$$



We have then

From Taylor

$$0 = f(r) = f(x_0 + h) \approx f(x_0) + hf'(x_0)$$

Thus, as long  $f'(x_0)$  is not close to 0

$$h \approx -\frac{f(x_0)}{f'(x_0)}$$

Thus

$$r = x_0 + h \approx x_0 - \frac{f(x_0)}{f'(x_0)}$$

We have our final improving

We have

$$x_1 \approx x_0 - \frac{f(x_0)}{f'(x_0)}$$

Then, on the scoring function

For this, we need the Hessian of the function

$$\frac{\partial^2 \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = - \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \left[ \frac{\exp \{ \mathbf{w}^T \mathbf{x}_i \}}{1 + \exp \{ \mathbf{w}^T \mathbf{x}_i \}} \right] \left[ 1 - \frac{\exp \{ \mathbf{w}^T \mathbf{x}_i \}}{1 + \exp \{ \mathbf{w}^T \mathbf{x}_i \}} \right]$$

Thus, we have our starting point  $\mathbf{w}^{old}$

$$\mathbf{w}^{new} = \mathbf{w}^{old} - \left( \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} \right)^{-1} \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}}$$

Then, on the scoring function

For this, we need the Hessian of the function

$$\frac{\partial^2 \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = - \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \left[ \frac{\exp \{ \mathbf{w}^T \mathbf{x}_i \}}{1 + \exp \{ \mathbf{w}^T \mathbf{x}_i \}} \right] \left[ 1 - \frac{\exp \{ \mathbf{w}^T \mathbf{x}_i \}}{1 + \exp \{ \mathbf{w}^T \mathbf{x}_i \}} \right]$$

Thus, we have at a starting point  $\mathbf{w}^{old}$

$$\mathbf{w}^{new} = \mathbf{w}^{old} - \left( \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} \right)^{-1} \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}}$$

# Outline

## 1 Logistic Regression

- Introduction
- Constraints
- The Initial Model
- The Two Case Class
- Graphic Interpretation
- **Fitting The Model**
  - The Two Class Case
  - The Final Log-Likelihood
  - The Newton-Raphson Algorithm
- **Matrix Notation**

## 2 More on Optimization Methods

- Can we do better?
- Using Cholesky Decomposition
  - Cholesky Decomposition
  - The Proposed Method
- Quasi-Newton Method
  - The Second Order Approximation
  - The BFGS Algorithm
- A Neat Trick: Coordinate Ascent
  - Coordinate Ascent Algorithm
- Conclusion

# We can rewrite all as matrix notations

## Assume

- 1 Let  $\mathbf{y}$  denotes the vector of  $y_i$
- 2  $X$  is the data matrix  $N \times (d + 1)$
- 3  $\mathbf{p}$  the vector of fitted probabilities with the  $i^{th}$  element  $p(x_i|w^{old})$
- 4  $W$  a  $N \times N$  diagonal matrix of weights with the  $i^{th}$  diagonal element

$$p(x_i|w^{old}) [1 - p(x_i|w^{old})]$$

# We can rewrite all as matrix notations

## Assume

- 1 Let  $\mathbf{y}$  denotes the vector of  $y_i$
- 2  $X$  is the data matrix  $N \times (d + 1)$
- 3  $\mathbf{p}$  the vector of fitted probabilities with the  $i^{th}$  element  $p(x_i|w^{old})$
- 4  $W$  a  $N \times N$  diagonal matrix of weights with the  $i^{th}$  diagonal element

$$p(x_i|w^{old}) [1 - p(x_i|w^{old})]$$

# We can rewrite all as matrix notations

## Assume

- 1 Let  $\mathbf{y}$  denotes the vector of  $y_i$
- 2  $X$  is the data matrix  $N \times (d + 1)$
- 3  $\mathbf{p}$  the vector of fitted probabilities with the  $i^{th}$  element  $p(\mathbf{x}_i | \mathbf{w}^{old})$
- 4  $W$  a  $N \times N$  diagonal matrix of weights with the  $i^{th}$  diagonal element

$$p(\mathbf{x}_i | \mathbf{w}^{old}) [1 - p(\mathbf{x}_i | \mathbf{w}^{old})]$$



# We can rewrite all as matrix notations

## Assume

- ① Let  $\mathbf{y}$  denotes the vector of  $y_i$
- ②  $X$  is the data matrix  $N \times (d + 1)$
- ③  $\mathbf{p}$  the vector of fitted probabilities with the  $i^{th}$  element  $p(\mathbf{x}_i | \mathbf{w}^{old})$
- ④  $W$  a  $N \times N$  diagonal matrix of weights with the  $i^{th}$  diagonal element

$$p(\mathbf{x}_i | \mathbf{w}^{old}) [1 - p(\mathbf{x}_i | \mathbf{w}^{old})]$$

Then, we have

For each updating term

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = X^T (\mathbf{y} - \mathbf{p})$$
$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = -X^T W X$$

Then, we have

Then, the Newton Step

$$\begin{aligned} \mathbf{w}^{new} &= \mathbf{w}^{old} + \left( X^T W X \right)^{-1} X^T (\mathbf{y} - \mathbf{p}) \\ &= I \mathbf{w}^{old} + \left( X^T W X \right)^{-1} X^T I (\mathbf{y} - \mathbf{p}) \\ &= \left( X^T W X \right)^{-1} X^T W X \mathbf{w}^{old} + \left( X^T W X \right)^{-1} X^T W W^{-1} (\mathbf{y} - \mathbf{p}) \\ &= \left( X^T W X \right)^{-1} X^T W \left[ X \mathbf{w}^{old} + W^{-1} (\mathbf{y} - \mathbf{p}) \right] \\ &= \left( X^T W X \right)^{-1} X^T W \mathbf{z} \end{aligned}$$

Then, we have

### Then, the Newton Step

$$\begin{aligned}\mathbf{w}^{new} &= \mathbf{w}^{old} + \left(X^T W X\right)^{-1} X^T (\mathbf{y} - \mathbf{p}) \\ &= I \mathbf{w}^{old} + \left(X^T W X\right)^{-1} X^T I (\mathbf{y} - \mathbf{p}) \\ &= \left(X^T W X\right)^{-1} X^T W X \mathbf{w}^{old} + \left(X^T W X\right)^{-1} X^T W W^{-1} (\mathbf{y} - \mathbf{p}) \\ &= \left(X^T W X\right)^{-1} X^T W \left[X \mathbf{w}^{old} + W^{-1} (\mathbf{y} - \mathbf{p})\right] \\ &= \left(X^T W X\right)^{-1} X^T W \mathbf{z}\end{aligned}$$

Then, we have

### Then, the Newton Step

$$\begin{aligned}\mathbf{w}^{new} &= \mathbf{w}^{old} + \left(X^T W X\right)^{-1} X^T (\mathbf{y} - \mathbf{p}) \\ &= I \mathbf{w}^{old} + \left(X^T W X\right)^{-1} X^T I (\mathbf{y} - \mathbf{p}) \\ &= \left(X^T W X\right)^{-1} X^T W X \mathbf{w}^{old} + \left(X^T W X\right)^{-1} X^T W W^{-1} (\mathbf{y} - \mathbf{p}) \\ &= \left(X^T W X\right)^{-1} X^T W \left[X \mathbf{w}^{old} + W^{-1} (\mathbf{y} - \mathbf{p})\right] \\ &= \left(X^T W X\right)^{-1} X^T W \mathbf{z}\end{aligned}$$

Then, we have

### Then, the Newton Step

$$\begin{aligned}\mathbf{w}^{new} &= \mathbf{w}^{old} + (X^T W X)^{-1} X^T (\mathbf{y} - \mathbf{p}) \\ &= I \mathbf{w}^{old} + (X^T W X)^{-1} X^T I (\mathbf{y} - \mathbf{p}) \\ &= (X^T W X)^{-1} X^T W X \mathbf{w}^{old} + (X^T W X)^{-1} X^T W W^{-1} (\mathbf{y} - \mathbf{p}) \\ &= (X^T W X)^{-1} X^T W [X \mathbf{w}^{old} + W^{-1} (\mathbf{y} - \mathbf{p})] \\ &= (X^T W X)^{-1} X^T W \mathbf{z}\end{aligned}$$

Then, we have

### Then, the Newton Step

$$\begin{aligned}\mathbf{w}^{new} &= \mathbf{w}^{old} + (X^T W X)^{-1} X^T (\mathbf{y} - \mathbf{p}) \\ &= I \mathbf{w}^{old} + (X^T W X)^{-1} X^T I (\mathbf{y} - \mathbf{p}) \\ &= (X^T W X)^{-1} X^T W X \mathbf{w}^{old} + (X^T W X)^{-1} X^T W W^{-1} (\mathbf{y} - \mathbf{p}) \\ &= (X^T W X)^{-1} X^T W [X \mathbf{w}^{old} + W^{-1} (\mathbf{y} - \mathbf{p})] \\ &= (X^T W X)^{-1} X^T W \mathbf{z}\end{aligned}$$

Then

We have

Re-expressed the Newton step as a weighted least squares step.

With a the adjusted response as

$$z = Xw^{old} + W^{-1}(y - p)$$



Then

We have

Re-expressed the Newton step as a weighted least squares step.

With a the adjusted response as

$$z = X\mathbf{w}^{old} + W^{-1}(\mathbf{y} - \mathbf{p})$$

# This New Algorithm

It is known as

Iteratively Re-weighted Least Squares or IRLS

After all at each iteration, it solves

A weighted Least Square Problem

$$w^{new} \leftarrow \arg \min_w (z - Xw)^T W (z - Xw)$$

# This New Algorithm

It is known as

Iteratively Re-weighted Least Squares or IRLS

After all at each iteration, it solves

A weighted Least Square Problem

$$\mathbf{w}^{new} \leftarrow \arg \min_{\mathbf{w}} (\mathbf{z} - \mathbf{X}\mathbf{w})^T \mathbf{W} (\mathbf{z} - \mathbf{X}\mathbf{w})$$

# Observations

Good Starting Point  $w = 0$

However, convergence is never guaranteed!!!

However:

- Typically the algorithm does converge, since the log-likelihood is concave.
- But overshooting can occur.

# Observations

## Good Starting Point $w = 0$

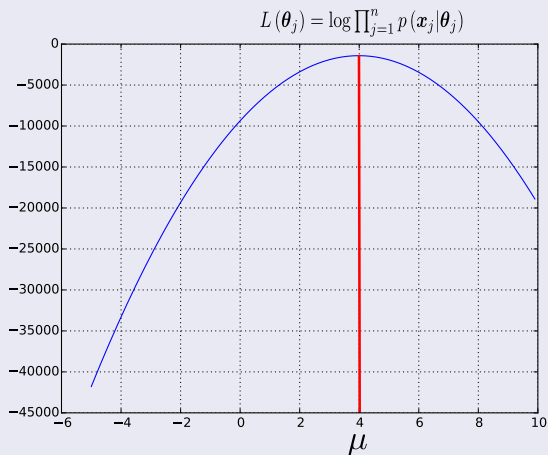
However, convergence is never guaranteed!!!

## However

- Typically the algorithm does converge, since the log-likelihood is concave.
- But overshooting can occur.

# Observations

$$L(\theta_j) = \log \prod_{j=1}^n p(\mathbf{x}_j | \theta_j)$$

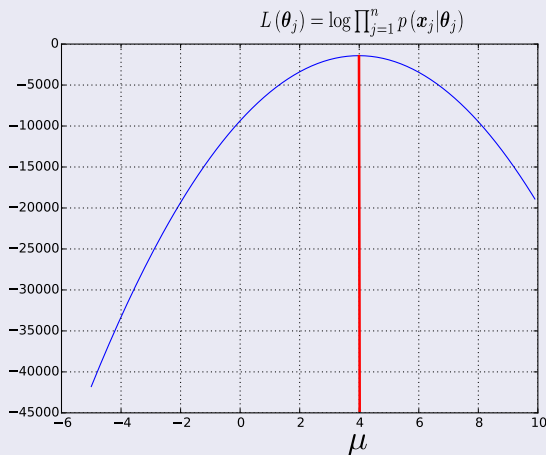


Having Solve the Problem

Perfect!!!

# Observations

$$L(\theta_j) = \log \prod_{j=1}^n p(\mathbf{x}_j | \theta_j)$$



Halving Solve the Problem

Perfect!!!

# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion



# The final question

After all, we always want to have a better solution

- We know that  $\left(\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T}\right)^{-1}$  takes  $O(d^3)$ .... and we want something better!!!

We have the following methods

- Colesky Decomposition
- Quasi-Newton Method
- Coordinate Ascent

# The final question

After all, we always want to have a better solution

- We know that  $\left(\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T}\right)^{-1}$  takes  $O(d^3)$ .... and we want something better!!!

We have the following methods

- Colesky Decomposition
- Quasi-Newton Method
- Coordinate Ascent

# The final question

After all, we always want to have a better solution

- We know that  $\left(\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T}\right)^{-1}$  takes  $O(d^3)$ .... and we want something better!!!

We have the following methods

- Colesky Decomposition
- Quasi-Newton Method
- Coordinate Ascent

# The final question

After all, we always want to have a better solution

- We know that  $\left(\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T}\right)^{-1}$  takes  $O(d^3)$ .... and we want something better!!!

We have the following methods

- Colesky Decomposition
- Quasi-Newton Method
- Coordinate Ascent

# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - **Using Cholesky Decomposition**
    - Cholesky Decomposition
    - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

We can decompose the matrix

Given  $A = X^T W X$  and  $Y = X^T W z$ , you have

$$Ax = Y$$

We want to obtain

$$x = A^{-1}Y$$

We can decompose the matrix

Given  $A = X^T W X$  and  $Y = X^T W z$ , you have

$$Ax = Y$$

We want to obtain

$$x = A^{-1}Y$$

# This can be seen as a system of linear equations

As you can see

- We start with a set of linear equations with  $d + 1$  unknowns:

$$x_1, x_2, \dots, x_{d+1} \begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1d+1}x_{d+1} & = y_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2d+1}x_{d+1} & = y_2 \\ \vdots & \vdots \\ a_{d+11}x_1 + a_{d+12}x_2 + \dots + a_{d+1d+1}x_{d+1} & = y_{d+1} \end{cases}$$

Thus

- A set of values for  $x_1, x_2, \dots, x_n$  that satisfy all of the equations simultaneously is said to be a solution to these equations.



# This can be seen as a system of linear equations

As you can see

- We start with a set of linear equations with  $d + 1$  unknowns:

$$x_1, x_2, \dots, x_{d+1} \begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1d+1}x_{d+1} & = y_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2d+1}x_{d+1} & = y_2 \\ \vdots & \vdots \\ a_{d+11}x_1 + a_{d+12}x_2 + \dots + a_{d+1d+1}x_{d+1} & = y_{d+1} \end{cases}$$

Thus

- A set of values for  $x_1, x_2, \dots, x_n$  that satisfy all of the equations simultaneously is said to be a solution to these equations.

# This can be seen as a system of linear equations

As you can see

- We start with a set of linear equations with  $d + 1$  unknowns:

$$x_1, x_2, \dots, x_{d+1} \begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1d+1}x_{d+1} & = y_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2d+1}x_{d+1} & = y_2 \\ \vdots & \vdots \\ a_{d+11}x_1 + a_{d+12}x_2 + \dots + a_{d+1d+1}x_{d+1} & = y_{d+1} \end{cases}$$

Thus

- A set of values for  $x_1, x_2, \dots, x_n$  that satisfy all of the equations simultaneously is said to be a solution to these equations.

# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
      - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

# What is the Cholesky Decomposition? [4]

It is a method that factorize a matrix

- $A \in \mathbb{R}^{d+1 \times d+1}$  is a positive definite Hermitian matrix

Positive definite matrix

$$x^T A x > 0 \text{ for all } x \in \mathbb{R}^{d+1 \times d+1}$$

Hermitian matrix in the Real Domain (Symmetric Matrix)

$$A = A^T$$

# What is the Cholesky Decomposition? [4]

It is a method that factorize a matrix

- $A \in \mathbb{R}^{d+1 \times d+1}$  is a positive definite Hermitian matrix

Positive definite matrix

$$\mathbf{x}^T A \mathbf{x} > 0 \text{ for all } \mathbf{x} \in \mathbb{R}^{d+1 \times d+1}$$

Hermitian matrix in the Real Domain (Symmetric Matrix)

$$A = A^T$$

# What is the Cholesky Decomposition? [4]

It is a method that factorize a matrix

- $A \in \mathbb{R}^{d+1 \times d+1}$  is a positive definite Hermitian matrix

Positive definite matrix

$$\mathbf{x}^T A \mathbf{x} > 0 \text{ for all } \mathbf{x} \in \mathbb{R}^{d+1 \times d+1}$$

Hermitian matrix in the Real Domain (Symmetric Matrix)

$$A = A^T$$

Therefore

Cholesky decomposes  $A$  into lower or upper triangular matrix and their conjugate transpose

$$A = LL^T$$

$$A = R^T R$$

Thus, we can use the Cholesky decomposition

- The Cholesky decomposition is of order  $O(d^3)$  and requires  $\frac{1}{6}d^3$  FLOP operations.

Therefore

Cholesky decomposes  $A$  into lower or upper triangular matrix and their conjugate transpose

$$A = LL^T$$

$$A = R^T R$$

Thus, we can use the Cholesky decomposition

- The Cholesky decomposition is of order  $O(d^3)$  and requires  $\frac{1}{6}d^3$  FLOP operations.



# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

We have

The matrices  $A \in \mathbb{R}^{d+1 \times d+1}$  and  $X = A^{-1}$

$$AX = I$$

From Cholesky, the decomposition of  $A$

$$R^T R X = I$$

If we define  $RX = B$

$$R^T B = I$$

We have

The matrices  $A \in \mathbb{R}^{d+1 \times d+1}$  and  $X = A^{-1}$

$$AX = I$$

From Cholensky, the decomposition of  $A$

$$R^T R X = I$$

If we define  $BX = B$

$$R^T B = I$$

We have

The matrices  $A \in \mathbb{R}^{d+1 \times d+1}$  and  $X = A^{-1}$

$$AX = I$$

From Cholensky, the decomposition of  $A$

$$R^T R X = I$$

If we define  $RX = B$

$$R^T B = I$$

## Now

If  $B = (R^T)^{-1} = L^{-1}$  for  $L = R^T$

- 1 We note that the inverse of the lower triangular matrix  $L$  is lower triangular.
- 2 The diagonal entries of  $L^{-1}$  are the reciprocal of diagonal entries of  $L$

$$\begin{pmatrix} a_{1,1} & 0 & \cdots & 0 \\ a_{2,1} & a_{2,2} & 0 & \vdots \\ \vdots & \vdots & \ddots & 0 \\ a_{d+1,1} & a_{d+1,2} & \cdots & a_{d+1,d+1} \end{pmatrix} \begin{pmatrix} b_{1,1} & 0 & \cdots & 0 \\ b_{2,1} & b_{2,2} & 0 & \vdots \\ \vdots & \vdots & \ddots & 0 \\ b_{d+1,1} & b_{d+1,2} & \cdots & b_{d+1,d+1} \end{pmatrix} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{pmatrix}$$

Now

Now, we construct the following matrix  $S$  with entries

$$s_{i,j} = \begin{cases} \frac{1}{l_{i,i}} & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

Now, we have

The matrix  $S$  is the correct solution to upper diagonal element of the matrix  $B$

i.e.  $s_{ij} = b_{ij}$  for  $i \leq j \leq d+1$

Then, we use backward substitution to solve  $AX = S$  at equation  $d+1 = n$

Assuming:

$$X = [x_1, x_2, \dots, x_{d+1}]$$

$$S = [s_1, s_2, \dots, s_{d+1}]$$

Now, we have

The matrix  $S$  is the correct solution to upper diagonal element of the matrix  $B$

i.e.  $s_{ij} = b_{ij}$  for  $i \leq j \leq d+1$

Then, we use backward substitution to solve  $x_{i,j}$  at equation  $R\mathbf{x}_i = \mathbf{s}_i$

Assuming:

$$X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{d+1}]$$

$$S = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_{d+1}]$$



# Back Substitution

## Back substitution

Since  $R$  is upper-triangular, we can rewrite the system  $R\mathbf{x}_i = \mathbf{s}_i$  as

$$r_{1,1}x_{1,i} + r_{1,2}x_{2,i} + \dots + r_{1,d-1}x_{d-1,i} + r_{1,d}x_{d,i} + r_{1,d+1}x_{d+1,i} = s_{1,i}$$

$$r_{2,2}x_{2,i} + \dots + r_{2,d-1}x_{d-1,i} + r_{2,d}x_{d,i} + r_{2,d+1}x_{d+1,i} = s_{2,i}$$

$$\vdots$$

$$r_{d-1,d-1}x_{d-1,i} + r_{d-1,d}x_{d,i} + r_{d-1,d+1}x_{d+1,i} = s_{d-1,i}$$

$$r_{d,d}x_{d,i} + r_{d,d+1}x_{d+1,i} = s_{d,i}$$

$$r_{d+1,d+1}x_{d+1,i} = s_{d+1,i}$$

Then

We solve only for  $x_{ij}$  such that

- We have  $i < j \leq N$  (Upper triangle elements).

• In our case the same value given that we live on the reals.

Then

We solve only for  $x_{ij}$  such that

- We have  $i < j \leq N$  (Upper triangle elements).

$$x_{ji} = \overline{x_{ij}}$$

- In our case the same value given that we live on the reals.

# Complexity

## Equation solving requires

- $\frac{1}{3}(d+1)^3$  multiply operations.

The total number of multiply operations for matrix inverse

- Including Cholesky decomposition is  $\frac{1}{2}(d+1)^3$

Therefore

- We have complexity  $O(d^3)$ !!! Per iteration!!! But actually  $\frac{1}{2}(d+1)^3$  multiply operations

# Complexity

## Equation solving requires

- $\frac{1}{3}(d+1)^3$  multiply operations.

## The total number of multiply operations for matrix inverse

- Including Cholesky decomposition is  $\frac{1}{2}(d+1)^3$

## Therefore

- We have complexity  $O(d^3)$ !!! Per iteration!!! But actually  $\frac{1}{2}(d+1)^3$  multiply operations

# Complexity

## Equation solving requires

- $\frac{1}{3}(d+1)^3$  multiply operations.

## The total number of multiply operations for matrix inverse

- Including Cholesky decomposition is  $\frac{1}{2}(d+1)^3$

## Therefore

- **We have complexity  $O(d^3)$ !!! Per iteration!!! But actually  $\frac{1}{2}(d+1)^3$  multiply operations**

# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - **Quasi-Newton Method**
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

However

As Good Computer Scientists

We want to obtain a better complexity as  $O(d^2)$ !!!

We can obtain such improvements

Using Quasi Newton Methods

Let's us to develop the solution

- For the most popular one
  - Broyden-Fletcher-Goldfarb-Shanno (BFGS) method



However

As Good Computer Scientists

We want to obtain a better complexity as  $O(d^2)$ !!!

We can obtain such improvements

Using Quasi Newton Methods

Let's now develop the solution

- For the most popular one
  - ▶ Broyden-Fletcher-Goldfarb-Shanno (BFGS) method

However

As Good Computer Scientists

We want to obtain a better complexity as  $O(d^2)$ !!!

We can obtain such improvements

Using Quasi Newton Methods

Let's us to develop the solution

- For the most popular one
  - ▶ Broyden-Fletcher-Goldfarb-Shanno (BFGS) method

# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

# The Second Order Approximation

We have

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_k)^T \mathbf{H} f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k)$$

We develop a new equation based on the previous idea by using  $\mathbf{x} = \mathbf{x}_k + \mathbf{p}$

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k) \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}$$

Here

- $\mathbf{H}_k$  is an  $d+1 \times d+1$  symmetric positive definite matrix that will be updated through the entire process

# The Second Order Approximation

We have

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_k)^T \mathbf{H} f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k)$$

We develop a new equation based in the previous idea by using  $\mathbf{x} = \mathbf{x}_k + \mathbf{p}$

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)\mathbf{p} + \frac{1}{2}\mathbf{p}^T \mathbf{H}_k \mathbf{p}$$

Here

- $\mathbf{H}_k$  is an  $d+1 \times d+1$  symmetric positive definite matrix that will be updated through the entire process

# The Second Order Approximation

We have

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_k)^T \mathbf{H} f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k)$$

We develop a new equation based in the previous idea by using  $\mathbf{x} = \mathbf{x}_k + \mathbf{p}$

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)\mathbf{p} + \frac{1}{2}\mathbf{p}^T \mathbf{H}_k \mathbf{p}$$

Here

- $\mathbf{H}_k$  is an  $d + 1 \times d + 1$  symmetric positive definite matrix that will be updated through the entire process

## For the BFGS

Then, the inverse update of it  $H_k = B_k^{-1}$

In BFGS we go directly for the inverse by setting up:

$$\begin{aligned} \min_H & \|H - H_k\| \\ \text{s.t. } & H = H^T \\ & Hy_k = s_k \end{aligned}$$

with

$$\begin{aligned} s_k &= \mathbf{x}_{k+1} - \mathbf{x}_k \\ y_k &= \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k) \end{aligned}$$

Then, we have

A unique solution will be

$$H_{k+1} = \left(I - \rho_k s_k y_k^T\right) H_k \left(I - \rho_k y_k s_k^T\right) + \rho s_k s_k^T \quad (1)$$

where  $\rho_k = \frac{1}{y_k^T s_k}$



Then, we have

A unique solution will be

$$H_{k+1} = \left( I - \rho_k s_k y_k^T \right) H_k \left( I - \rho_k y_k s_k^T \right) + \rho s_k s_k^T \quad (1)$$

where  $\rho_k = \frac{1}{y_k^T s_k}$

# Complexity of Generating $H_{k+1}$

We notice that the complexity of calculating

$$s_k s_k^T, s_k s_k^T, s_k y_k^T$$

- It is  $O(d^2)$

Why? For Example

$$\begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_d \end{pmatrix} \begin{pmatrix} s_1 & s_2 & \cdots & s_d \end{pmatrix} = \begin{pmatrix} s_1^2 & s_1 s_2 & \cdots & s_1 s_d \\ s_2 s_1 & s_2^2 & \cdots & s_2 s_d \\ \vdots & \vdots & \ddots & \vdots \\ s_d s_1 & s_d s_2 & \cdots & s_d^2 \end{pmatrix} \quad \text{-Equal to } d^2$$

# Complexity of Generating $H_{k+1}$

We notice that the complexity of calculating

$$s_k s_k^T, s_k s_k^T, s_k y_k^T$$

- It is  $O(d^2)$

Why? For Example

$$\begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_d \end{pmatrix} \begin{pmatrix} s_1 & s_2 & \cdots & s_d \end{pmatrix} = \begin{pmatrix} s_1^2 & s_1 s_2 & \cdots & s_1 s_d \\ s_2 s_1 & s_2^2 & \cdots & s_2 s_d \\ \vdots & \vdots & \ddots & \vdots \\ s_d s_1 & s_d s_2 & \cdots & s_d^2 \end{pmatrix} \text{ -Equal to } d^2 \text{ m}$$

Then, we have that

Computation  $\rho_k = \frac{1}{y_k^T s_k}$  as a constant

Complexity  $O(d^2)$

Now for  $(J - \rho_k H_k)$

Complexity  $O(d^2)$

Then, we have that

Computation  $\rho_k = \frac{1}{y_k^T s_k}$  as a constant

Complexity  $O(d^2)$

Now for  $(J - \rho_k H_k)$

Complexity  $O(d^2)$

Then, we have that

Computation  $\rho_k = \frac{1}{y_k^T s_k}$  as a constant

Complexity  $O(d^2)$

Now for  $(I - \rho_k y_k s_k^T)$

Complexity  $O(d^2)$

Finally, we have

If we expand the equation  $(I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T)$

$$(I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) = (H_k - \rho_k s_k [y_k^T H_k]) - (\rho_k [H_k y_k] s_k^T - \rho_k^2 s_k [y_k^T H_k y_k] s_k^T)$$

Where

- $[y_k^T H_k]$  and  $[H_k y_k]$  has complexity  $O(d^2)$
- Thus  $s_k [y_k^T H_k]$  and  $[y_k^T H_k y_k] s_k^T$  has complexity  $O(d^2)$

Thus

- $H_{k+1}$  has a complexity of  $O(d^2)$  by avoiding inverting a Hessian Matrix

Finally, we have

If we expand the equation  $(I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T)$

$$(I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) = (H_k - \rho_k s_k [y_k^T H_k]) - (\rho_k [H_k y_k] s_k^T - \rho_k^2 s_k [y_k^T H_k y_k] s_k^T)$$

Where

- 1  $[y_k^T H_k]$  and  $[H_k y_k]$  has complexity  $O(d^2)$
- 2 Thus  $s_k [y_k^T H_k]$  and  $[y_k^T H_k y_k]$  has complexity  $O(d^2)$

Thus

- $H_{k+1}$  has a complexity of  $O(d^2)$  by avoiding inverting a Hessian Matrix



Finally, we have

If we expand the equation  $(I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T)$

$$(I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) = (H_k - \rho_k s_k [y_k^T H_k]) - (\rho_k [H_k y_k] s_k^T - \rho_k^2 s_k [y_k^T H_k y_k] s_k^T)$$

Where

- ①  $[y_k^T H_k]$  and  $[H_k y_k]$  has complexity  $O(d^2)$
- ② Thus  $s_k [y_k^T H_k]$  and  $[y_k^T H_k y_k] s_k^T$  has complexity  $O(d^2)$

Thus

- $H_{k+1}$  has a complexity of  $O(d^2)$  by avoiding inverting a Hessian Matrix

# Problem

There is no magic formula to find an initial  $H_0$

We can use specific information about the problem:

- For instance by setting it to the inverse of an approximate Hessian calculated by finite differences at  $x_0$
- In our case, we have  $\frac{\partial \mathcal{L}(w)}{\partial w \partial w^T}$ , or in matrix format  $X^T W X$ , we could get initial setup

# Problem

There is no magic formula to find an initial  $H_0$

We can use specific information about the problem:

- For instance by setting it to the inverse of an approximate Hessian calculated by finite differences at  $x_0$

• In our case, we have  $\frac{\partial \mathcal{L}(w)}{\partial w \partial w^T}$ , or in matrix format  $X^T W X$ , we could get initial setup

# Problem

There is no magic formula to find an initial  $H_0$

We can use specific information about the problem:

- For instance by setting it to the inverse of an approximate Hessian calculated by finite differences at  $x_0$
- In our case, we have  $\frac{\partial \mathcal{L}(w)}{\partial w \partial w^T}$ , or in matrix format  $X^T W X$ , we could get initial setup

# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - **Quasi-Newton Method**
    - The Second Order Approximation
    - **The BFGS Algorithm**
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

# Algorithm (BFGS Method)

## Quasi-Newton Algorithm

- Starting point  $x_0$ , Convergence tolerance  $e$ , Inverse Hessian approximation  $H_0$

1  $k \leftarrow 0$

2 while  $\|\nabla f(x_{k+1})\| > e$

3     Compute search direction  $p_k = -H_k \nabla f(x_{k+1})$

4     Set  $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is obtained from a linear search (Under Wolfe conditions).

5     Define  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$

6     Compute  $H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho_k s_k s_k^T$

7  $k \leftarrow k + 1$

# Algorithm (BFGS Method)

## Quasi-Newton Algorithm

- Starting point  $x_0$ , Convergence tolerance  $e$ , Inverse Hessian approximation  $H_0$
- 1  $k \leftarrow 0$
  - 2 while  $\|\nabla f(x_{k+1})\| > e$ 
    - 3 Compute search direction  $p_k = -H_k \nabla f(x_{k+1})$
    - 4 Set  $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is obtained from a linear search (Under Wolfe conditions).
    - 5 Define  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
    - 6 Compute  $H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho s_k s_k^T$
    - 7  $k \leftarrow k + 1$

# Algorithm (BFGS Method)

## Quasi-Newton Algorithm

- Starting point  $x_0$ , Convergence tolerance  $e$ , Inverse Hessian approximation  $H_0$
- 1  $k \leftarrow 0$
- 2 while  $\|\nabla f(x_{k+1})\| > e$ 
  - 3 Compute search direction  $p_k = -H_k \nabla f(x_{k+1})$
  - 4 Set  $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is obtained from a linear search (Under Wolfe conditions).
  - 5 Define  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
  - 6 Compute  $H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho_k s_k s_k^T$
  - 7  $k \leftarrow k + 1$



# Algorithm (BFGS Method)

## Quasi-Newton Algorithm

- Starting point  $x_0$ , Convergence tolerance  $e$ , Inverse Hessian approximation  $H_0$
- 1  $k \leftarrow 0$
- 2 while  $\|\nabla f(x_{k+1})\| > e$
- 3     Compute search direction  $p_k = -H_k \nabla f(x_{k+1})$
- 4     Set  $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is obtained from a linear search (Under Wolfe conditions).
- 5     Define  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
- 6     Compute  $H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho s_k s_k^T$
- 7      $k \leftarrow k + 1$

# Algorithm (BFGS Method)

## Quasi-Newton Algorithm

- Starting point  $x_0$ , Convergence tolerance  $e$ , Inverse Hessian approximation  $H_0$
- 1  $k \leftarrow 0$
- 2 while  $\|\nabla f(x_{k+1})\| > e$
- 3     Compute search direction  $p_k = -H_k \nabla f(x_{k+1})$
- 4     Set  $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is obtained from a linear search (Under Wolfe conditions).
- 5     Define  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
- 6     Compute  $H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho s_k s_k^T$
- 7      $k \leftarrow k + 1$

# Algorithm (BFGS Method)

## Quasi-Newton Algorithm

- Starting point  $x_0$ , Convergence tolerance  $e$ , Inverse Hessian approximation  $H_0$
- 1  $k \leftarrow 0$
- 2 while  $\|\nabla f(x_{k+1})\| > e$
- 3     Compute search direction  $p_k = -H_k \nabla f(x_{k+1})$
- 4     Set  $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is obtained from a linear search (Under Wolfe conditions).
- 5     Define  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
- 6     Compute  $H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho_k s_k s_k^T$
- 7      $k \leftarrow k + 1$

# Algorithm (BFGS Method)

## Quasi-Newton Algorithm

- Starting point  $x_0$ , Convergence tolerance  $e$ , Inverse Hessian approximation  $H_0$
- 1  $k \leftarrow 0$
- 2 while  $\|\nabla f(x_{k+1})\| > e$
- 3     Compute search direction  $p_k = -H_k \nabla f(x_{k+1})$
- 4     Set  $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is obtained from a linear search (Under Wolfe conditions).
- 5     Define  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
- 6     Compute  $H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho s_k s_k^T$

$k \leftarrow k + 1$

# Algorithm (BFGS Method)

## Quasi-Newton Algorithm

- Starting point  $x_0$ , Convergence tolerance  $e$ , Inverse Hessian approximation  $H_0$
- 1  $k \leftarrow 0$
- 2 while  $\|\nabla f(x_{k+1})\| > e$
- 3     Compute search direction  $p_k = -H_k \nabla f(x_{k+1})$
- 4     Set  $x_{k+1} = x_k + \alpha_k p_k$  where  $\alpha_k$  is obtained from a linear search (Under Wolfe conditions).
- 5     Define  $s_k = x_{k+1} - x_k$  and  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
- 6     Compute  $H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho s_k s_k^T$
- 7      $k \leftarrow k + 1$

# Complexity

The

- Cost of update or inverse update is  $O(d^2)$  operations per iteration.

For More and better versions (With a Hessian Approximation)

- Nocedal, Jorge & Wright, Stephen J. (1999). Numerical Optimization. Springer-Verlag.

# Complexity

## The

- Cost of update or inverse update is  $O(d^2)$  operations per iteration.

## For More and better versions (With a Hessian Approximation)

- Nocedal, Jorge & Wright, Stephen J. (1999). Numerical Optimization. Springer-Verlag.

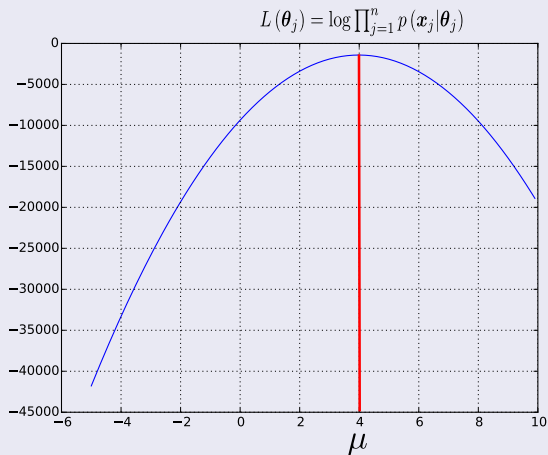
# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - **A Neat Trick: Coordinate Ascent**
    - Coordinate Ascent Algorithm
  - Conclusion



Given the following [5]

Because the likelihood is concave



## Caution

Here, we change the labeling to  $y_i = \pm 1$  with

$$p(y_i = \pm 1 | \mathbf{x}, \mathbf{w}) = \sigma(y \mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp\{-y \mathbf{w}^T \mathbf{x}\}}$$

Thus, we have the following log-likelihood under regularization

$$\mathcal{L}(\mathbf{w}) = - \sum_{i=1}^N \log \{1 + \exp\{-y_i \mathbf{w}^T \mathbf{x}_i\}\} - \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

It is possible to get a Gradient Descent

$$\nabla_{\mathbf{w}} l(\mathbf{w}) = \sum_{i=1}^N \left\{ 1 - \frac{1}{1 + \exp\{-y_i \mathbf{w}^T \mathbf{x}_i\}} \right\} y_i \mathbf{x}_i - \lambda \mathbf{w}$$

## Caution

Here, we change the labeling to  $y_i = \pm 1$  with

$$p(y_i = \pm 1 | \mathbf{x}, \mathbf{w}) = \sigma(y \mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp\{-y \mathbf{w}^T \mathbf{x}\}}$$

Thus, we have the following log likelihood under regularization  $\lambda > 0$

$$\mathcal{L}(\mathbf{w}) = - \sum_{i=1}^N \log \left\{ 1 + \exp \left\{ -y_i \mathbf{w}^T \mathbf{x}_i \right\} \right\} - \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

It is possible to get a Gradient Descent

$$\nabla_{\mathbf{w}} l(\mathbf{w}) = \sum_{i=1}^N \left\{ 1 - \frac{1}{1 + \exp\{-y_i \mathbf{w}^T \mathbf{x}_i\}} \right\} y_i \mathbf{x}_i - \lambda \mathbf{w}$$

## Caution

Here, we change the labeling to  $y_i = \pm 1$  with

$$p(y_i = \pm 1 | \mathbf{x}, \mathbf{w}) = \sigma(y \mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp\{-y \mathbf{w}^T \mathbf{x}\}}$$

Thus, we have the following log likelihood under regularization  $\lambda > 0$

$$\mathcal{L}(\mathbf{w}) = - \sum_{i=1}^N \log \left\{ 1 + \exp \left\{ -y_i \mathbf{w}^T \mathbf{x}_i \right\} \right\} - \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

It is possible to get a Gradient Descent

$$\nabla_{\mathbf{w}} l(\mathbf{w}) = \sum_{i=1}^N \left\{ 1 - \frac{1}{1 + \exp\{-y_i \mathbf{w}^T \mathbf{x}_i\}} \right\} y_i \mathbf{x}_i - \lambda \mathbf{w}$$

# Danger Will Robinson!!!

Gradient descent using resembles the Perceptron learning algorithm Problem!!! It will always converge for a suitable step size, regardless of whether the classes are separable!!!

# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

Here

We will simplify our work

By stating the algorithm for coordinate ascent

Then

A more precise version will be given

Here

We will simplify our work

By stating the algorithm for coordinate ascent

Then

A more precise version will be given



# Coordinate Ascent

## Algorithm

- Input  $Max$ , an initial  $w_0$
- ①  $counter \leftarrow 0$
- ② while  $counter < Max$
- ③     for  $i \leftarrow 1, \dots, d$
- ④         Randomly pick  $i$
- ⑤         Compute a step size  $\delta^*$  by approximately  
          maximize  $\arg \min_{\delta} f(x + \delta e_i)$
- ⑥          $x_i \leftarrow x_i + \delta^*$

where

$$e_i = (0 \ \dots \ 0 \ 1 \leftarrow i \ 0 \ \dots \ 0)^T$$

# Coordinate Ascent

## Algorithm

- Input  $Max$ , an initial  $w_0$
- ①  $counter \leftarrow 0$
- ② while  $counter < Max$
- ③     for  $i \leftarrow 1, \dots, d$
- ④         Randomly pick  $i$
- ⑤         Compute a step size  $\delta^*$  by approximately  
          maximize  $\arg \min_{\delta} f(x + \delta e_i)$
- ⑥          $x_i \leftarrow x_i + \delta^*$

## Where

$$e_i = \begin{pmatrix} 0 & \cdots & 0 & 1 \leftarrow i & 0 & \cdots & 0 \end{pmatrix}^T$$

## In the case of Logistic Regression

Thus, we can optimize each  $w_k$  alternatively by a coordinate-wise Newton update

$$w_k^{new} = w_k^{old} + \frac{-\lambda w_k^{old} + \sum_{i=1}^N \left\{ 1 - \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right\} y_i x_{ik}}{\lambda + \sum_{i=1}^N x_{ik}^2 \left( \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right) \left( 1 - \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right)}$$

Complexity of this update

Item	Complexity
$\sum_{i=1}^N \left\{ 1 - \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right\} y_i x_{ik}$	$O(N)$
$\sum_{i=1}^N x_{ik}^2 \left( \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right) \left( 1 - \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right)$	$O(N)$
Total Complexity	$O(N)$
For all the dimensions	$O(Nd)$

## In the case of Logistic Regression

Thus, we can optimize each  $w_k$  alternatively by a coordinate-wise Newton update

$$w_k^{new} = w_k^{old} + \frac{-\lambda w_k^{old} + \sum_{i=1}^N \left\{ 1 - \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right\} y_i x_{ik}}{\lambda + \sum_{i=1}^N x_{ik}^2 \left( \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right) \left( 1 - \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right)}$$

Complexity of this update

Item	Complexity
$\sum_{i=1}^N \left\{ 1 - \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right\} y_i x_{ik}$	$O(N)$
$\sum_{i=1}^N x_{ik}^2 \left( \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right) \left( 1 - \frac{1}{1 + \exp\{-y_i w^T x_i\}} \right)$	$O(N)$
Total Complexity	$O(N)$
For all the dimensions	$O(Nd)$

# Outline

- 1 Logistic Regression
  - Introduction
  - Constraints
  - The Initial Model
  - The Two Case Class
  - Graphic Interpretation
  - Fitting The Model
    - The Two Class Case
    - The Final Log-Likelihood
    - The Newton-Raphson Algorithm
    - Matrix Notation
- 2 More on Optimization Methods
  - Can we do better?
  - Using Cholesky Decomposition
    - Cholesky Decomposition
    - The Proposed Method
  - Quasi-Newton Method
    - The Second Order Approximation
    - The BFGS Algorithm
  - A Neat Trick: Coordinate Ascent
    - Coordinate Ascent Algorithm
  - Conclusion

We have the following Complexities per iteration

## Complexities

Method	Per Iteration	Convergence Rate
Cholesky Decomposition	$\frac{d^3}{2} = O(d^3)$	Quadratic
Quasi-Newton BFGS	$O(d^2)$	Super-linearly
Coordinate Ascent	$O(Nd)$	Not established



T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. Springer Series in Statistics, Springer New York, 2009.



F. E. Harrell Jr, *Regression modeling strategies: with applications to linear models, logistic and ordinal regression, and survival analysis*. Springer, 2015.



J. Nocedal and S. Wright, *Numerical optimization*. Springer Science & Business Media, 2006.



A. Krishnamoorthy and D. Menon, “Matrix inversion using cholesky decomposition,” in *2013 signal processing: Algorithms, architectures, arrangements, and applications (SPA)*, pp. 70–72, IEEE, 2013.



J. Dhanani and K. Rana, “Logistic regression with stochastic gradient ascent to estimate click through rate,” in *Information and Communication Technology for Sustainable Development*, pp. 319–326, Springer, 2018.

