

# CEE 616: Probabilistic Machine Learning

## M2 Linear Methods: L2b Logistic Regression

**Jimi Oke**

UMassAmherst

---

College of Engineering

Thu, Sep 25, 2025

# Outline

- 1 Introduction
- 2 Logistic regression model
- 3 MLE
- 4 Optimization
- 5 Outlook

# Logistic regression

# Logistic regression

The logistic regression model has the form:

# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

- **Binary logistic regression:**  $C = 2$ :

# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

- **Binary logistic regression:**  $C = 2$ :

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x})) \quad (2)$$

# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

- **Binary logistic regression:**  $C = 2$ :

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x})) \quad (2)$$

where  $\sigma$  is the sigmoid function and  $\boldsymbol{\theta} = \mathbf{w} = (b, w_1, w_2, \dots, w_D)$



# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

- **Binary logistic regression:**  $C = 2$ :

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x})) \quad (2)$$

where  $\sigma$  is the sigmoid function and  $\boldsymbol{\theta} = \mathbf{w} = (b, w_1, w_2, \dots, w_D)$

- **Multinomial logistic regression:**

# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

- **Binary logistic regression:**  $C = 2$ :

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x})) \quad (2)$$

where  $\sigma$  is the sigmoid function and  $\boldsymbol{\theta} = \mathbf{w} = (b, w_1, w_2, \dots, w_D)$

- **Multinomial logistic regression:**  $C > 2$

# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

- **Binary logistic regression:**  $C = 2$ :

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x})) \quad (2)$$

where  $\sigma$  is the sigmoid function and  $\boldsymbol{\theta} = \mathbf{w} = (b, w_1, w_2, \dots, w_D)$

- **Multinomial logistic regression:**  $C > 2$

$$p(y|\mathbf{x}; \boldsymbol{\theta}) =$$

# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

- **Binary logistic regression:**  $C = 2$ :

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x})) \quad (2)$$

where  $\sigma$  is the sigmoid function and  $\boldsymbol{\theta} = \mathbf{w} = (b, w_1, w_2, \dots, w_D)$

- **Multinomial logistic regression:**  $C > 2$

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Cat}(y|\mathcal{S}(\mathbf{W}\mathbf{x})) \quad (3)$$

# Logistic regression

The logistic regression model has the form:

$$p(y|\mathbf{x}; \boldsymbol{\theta}), \quad \mathbf{x} \in \mathbb{R}^D, y \in \{1, \dots, C\} \quad (1)$$

- **Binary logistic regression:**  $C = 2$ :

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x})) \quad (2)$$

where  $\sigma$  is the sigmoid function and  $\boldsymbol{\theta} = \mathbf{w} = (b, w_1, w_2, \dots, w_D)$

- **Multinomial logistic regression:**  $C > 2$

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Cat}(y|\mathcal{S}(\mathbf{W}\mathbf{x})) \quad (3)$$

where  $\mathcal{S}$  is the softmax function and  $\boldsymbol{\theta} = \mathbf{W}$ .

# Definitions

# Definitions

- Input vector:  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$

# Definitions

- Input vector:  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$
- Weights (or weight vector):  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$



# Definitions

- Input vector:  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$
- Weights (or weight vector):  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$
- Bias:  $b$

# Definitions

- Input vector:  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$
- Weights (or weight vector):  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$
- Bias:  $b$  (absorbed into weight vector)

# Definitions

- Input vector:  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$
- Weights (or weight vector):  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$
- Bias:  $b$  (absorbed into weight vector)
- Sigmoid function:

# Definitions

- Input vector:  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$
- Weights (or weight vector):  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$
- Bias:  $b$  (absorbed into weight vector)
- Sigmoid function:

$$\sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} \quad (4)$$

- Log-odds or logit:  $\mathbf{w}^\top \mathbf{x}$  (binary case);  $\mathbf{W}\mathbf{x}$  (multinomial)

# Definitions

- Input vector:  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$
- Weights (or weight vector):  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$
- Bias:  $b$  (absorbed into weight vector)
- Sigmoid function:

$$\sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} \quad (4)$$

- Log-odds or logit:  $\mathbf{w}^\top \mathbf{x}$  (binary case);  $\mathbf{W}\mathbf{x}$  (multinomial)
- Softmax function:

$$\mathcal{S}(\mathbf{W}\mathbf{x}) = \left[ \frac{e^{\mathbf{w}_1^\top \mathbf{x}}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^\top \mathbf{x}}}, \dots, \frac{e^{\mathbf{w}_C^\top \mathbf{x}}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^\top \mathbf{x}}} \right] \quad (5)$$

# Definitions

- Input vector:  $\mathbf{x} = (1, x_1, x_2, \dots, x_D)$
- Weights (or weight vector):  $\mathbf{w} = (b, w_1, w_2, \dots, w_D)$
- Bias:  $b$  (absorbed into weight vector)
- Sigmoid function:

$$\sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} \quad (4)$$

- Log-odds or logit:  $\mathbf{w}^\top \mathbf{x}$  (binary case);  $\mathbf{W}\mathbf{x}$  (multinomial)
- Softmax function:

$$\mathcal{S}(\mathbf{W}\mathbf{x}) = \left[ \frac{e^{\mathbf{w}_1^\top \mathbf{x}}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^\top \mathbf{x}}}, \dots, \frac{e^{\mathbf{w}_C^\top \mathbf{x}}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'}^\top \mathbf{x}}} \right] \quad (5)$$

where  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C]$  is a  $C \times D$  weight matrix

# Logistic (sigmoid) function

For  $w_1 > 0$ , the logistic function increases w.r.t.  $x$ .

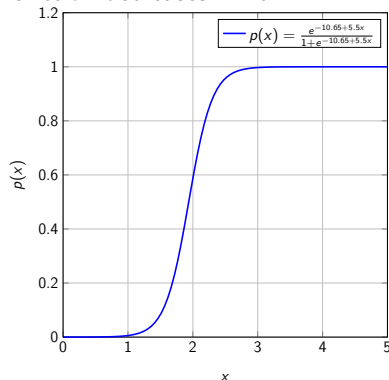
# Logistic (sigmoid) function

For  $w_1 > 0$ , the logistic function increases w.r.t.  $x$ . For  $w_1 < 0$ , the logistic function decreases w.r.t.  $x$ .



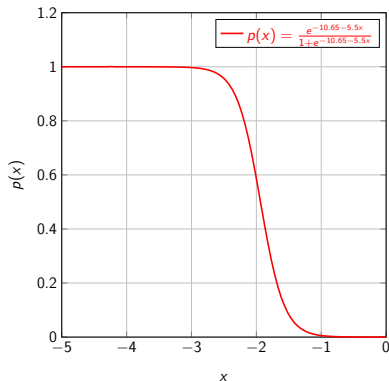
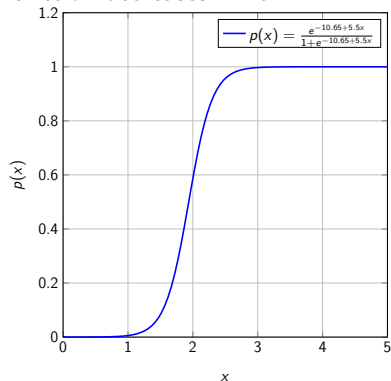
# Logistic (sigmoid) function

For  $w_1 > 0$ , the logistic function increases w.r.t.  $x$ . For  $w_1 < 0$ , the logistic function decreases w.r.t.  $x$ .



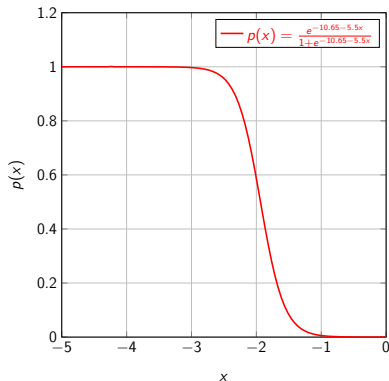
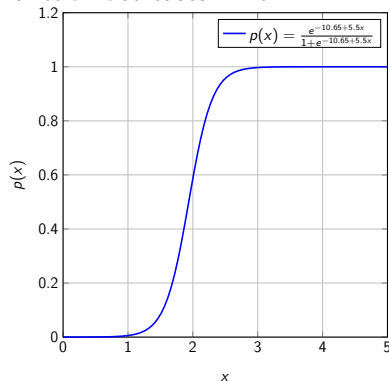
# Logistic (sigmoid) function

For  $w_1 > 0$ , the logistic function increases w.r.t.  $x$ . For  $w_1 < 0$ , the logistic function decreases w.r.t.  $x$ .



# Logistic (sigmoid) function

For  $w_1 > 0$ , the logistic function increases w.r.t.  $x$ . For  $w_1 < 0$ , the logistic function decreases w.r.t.  $x$ .



What happens when  $b$  is increased or decreased?

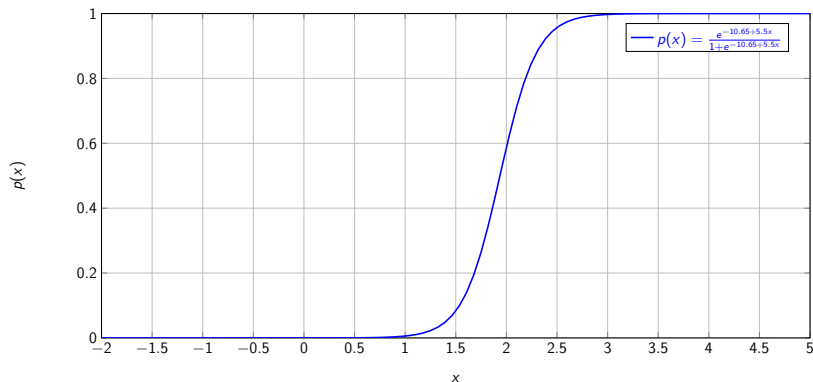
# Logistic function (cont.)

# Logistic function (cont.)

$b$  shifts the curve left or right (adjusts average fitted probabilities).

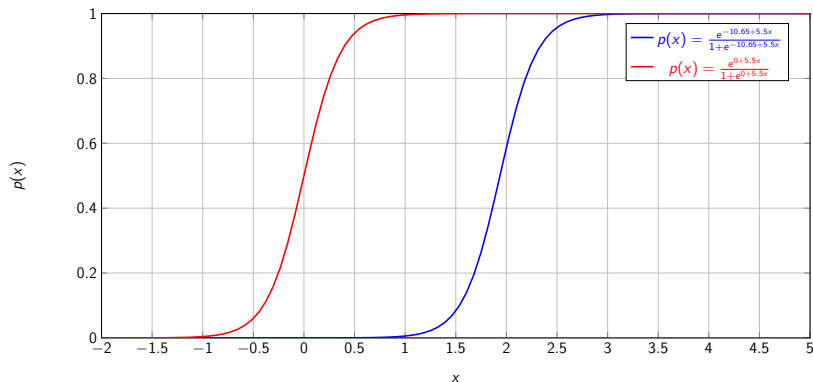
# Logistic function (cont.)

$b$  shifts the curve left or right (adjusts average fitted probabilities).



# Logistic function (cont.)

$b$  shifts the curve left or right (adjusts average fitted probabilities).



# Logistic function (cont.)

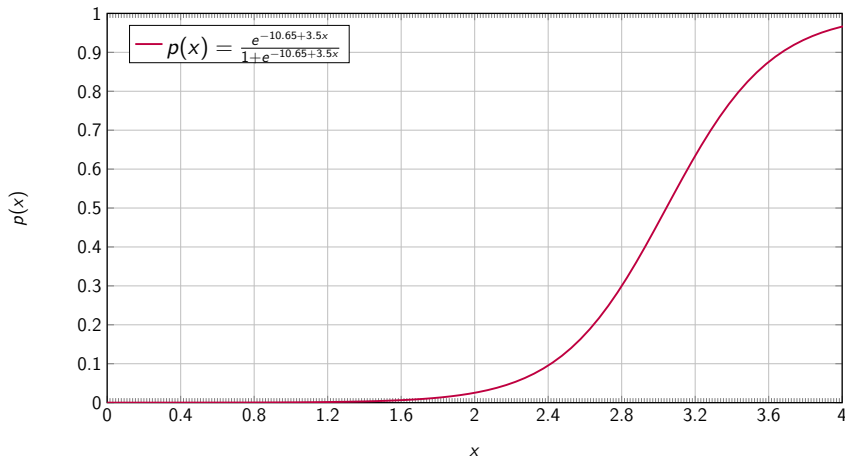


# Logistic function (cont.)

$w_1$  adjusts the steepness of the curve: as  $\beta_1$  increases, the curve becomes steeper

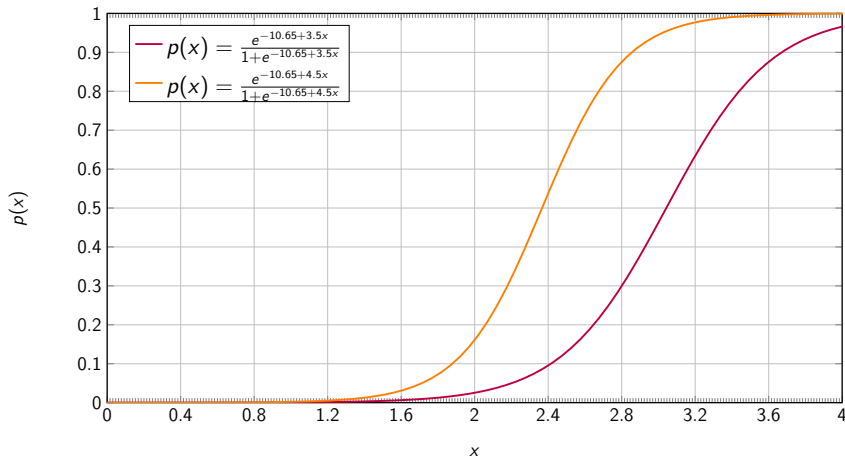
# Logistic function (cont.)

$w_1$  adjusts the steepness of the curve: as  $\beta_1$  increases, the curve becomes steeper



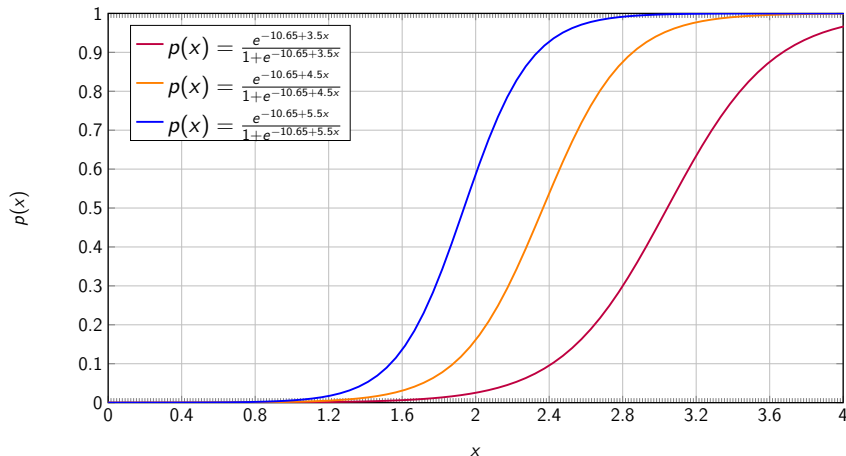
# Logistic function (cont.)

$w_1$  adjusts the steepness of the curve: as  $\beta_1$  increases, the curve becomes steeper



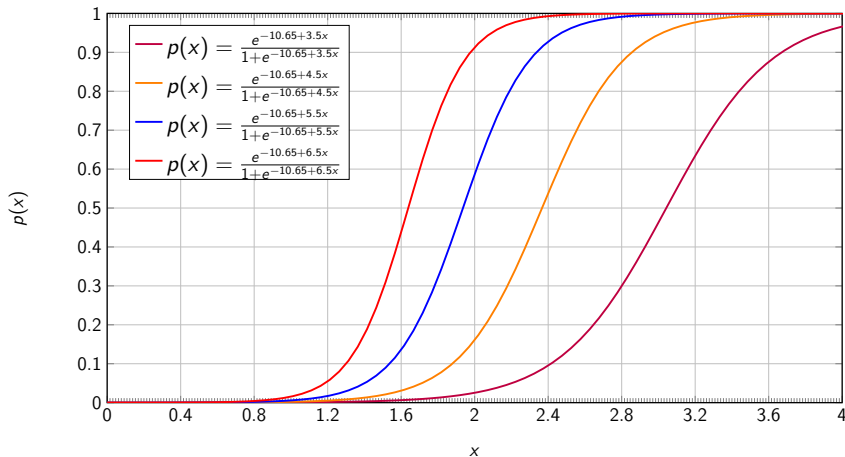
# Logistic function (cont.)

$w_1$  adjusts the steepness of the curve: as  $\beta_1$  increases, the curve becomes steeper



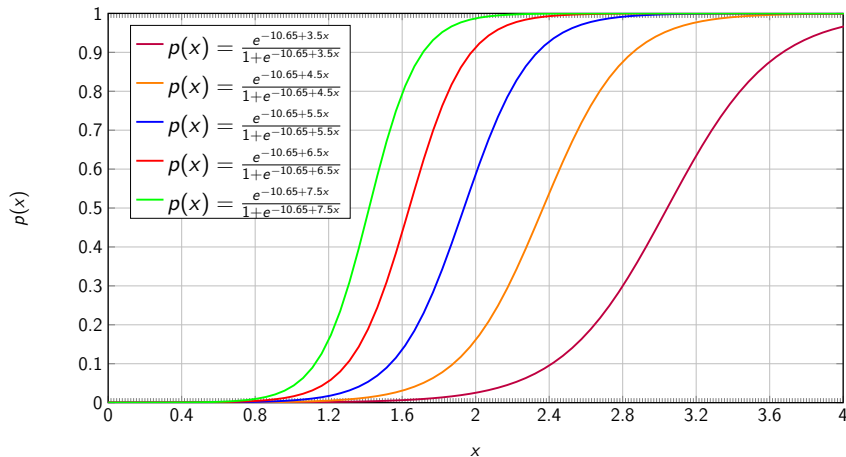
# Logistic function (cont.)

$w_1$  adjusts the steepness of the curve: as  $\beta_1$  increases, the curve becomes steeper



# Logistic function (cont.)

$w_1$  adjusts the steepness of the curve: as  $\beta_1$  increases, the curve becomes steeper



# Odds ratio and the logit function

# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:



# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:

$$OR = \frac{p(x)}{1 - p(x)} = e^{b + w_1 x} \quad (6)$$

which is considered as the relative likelihood of success ( $p(x)$ ).

# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:

$$OR = \frac{p(x)}{1 - p(x)} = e^{b + w_1 x} \quad (6)$$

which is considered as the relative likelihood of success ( $p(x)$ ).

Taking the log of the odds ratio yields the log-odds or **logit** function:

# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:

$$OR = \frac{p(x)}{1 - p(x)} = e^{b + w_1 x} \quad (6)$$

which is considered as the relative likelihood of success ( $p(x)$ ).

Taking the log of the odds ratio yields the log-odds or **logit** function:

$$\text{logit}(p(x)) =$$

# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:

$$OR = \frac{p(\mathbf{x})}{1 - p(\mathbf{x})} = e^{b + w_1 x} \quad (6)$$

which is considered as the relative likelihood of success ( $p(x)$ ).

Taking the log of the odds ratio yields the log-odds or **logit** function:

$$\text{logit}(p(x)) = \log\left(\frac{p(x)}{1 - p(x)}\right) = b + w_1 x \quad (7)$$

# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:

$$OR = \frac{p(x)}{1 - p(x)} = e^{b + w_1 x} \tag{6}$$

which is considered as the relative likelihood of success ( $p(x)$ ).

Taking the log of the odds ratio yields the log-odds or **logit** function:

$$\text{logit}(p(x)) = \log\left(\frac{p(x)}{1 - p(x)}\right) = b + w_1 x \tag{7}$$

## Notes

# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:

$$OR = \frac{p(x)}{1 - p(x)} = e^{b + w_1 x} \quad (6)$$

which is considered as the relative likelihood of success ( $p(x)$ ).

Taking the log of the odds ratio yields the log-odds or **logit** function:

$$\text{logit}(p(x)) = \log\left(\frac{p(x)}{1 - p(x)}\right) = b + w_1 x \quad (7)$$

## Notes

- The logit function is linear in  $x$

# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:

$$OR = \frac{p(x)}{1 - p(x)} = e^{b + w_1 x} \quad (6)$$

which is considered as the relative likelihood of success ( $p(x)$ ).

Taking the log of the odds ratio yields the log-odds or **logit** function:

$$\text{logit}(p(x)) = \log\left(\frac{p(x)}{1 - p(x)}\right) = b + w_1 x \quad (7)$$

## Notes

- The logit function is linear in  $x$
- The inverse of the logit function yields the logistic function

# Odds ratio and the logit function

From the logistic function, we can obtain the **odds ratio** ( $OR$ ) as:

$$OR = \frac{p(\mathbf{x})}{1 - p(\mathbf{x})} = e^{b + w_1 x} \quad (6)$$

which is considered as the relative likelihood of success ( $p(\mathbf{x})$ ).

Taking the log of the odds ratio yields the log-odds or **logit** function:

$$\text{logit}(p(\mathbf{x})) = \log\left(\frac{p(\mathbf{x})}{1 - p(\mathbf{x})}\right) = b + w_1 x \quad (7)$$

## Notes

- The logit function is linear in  $x$
- The inverse of the logit function yields the logistic function
- In the generalized linear framework, logit is the *link function* between the predictors and the mean response



# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

## Logistic function

This is the model equation for simple logistic regression:

# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

## Logistic function

This is the model equation for simple logistic regression:

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{b+w_1x}}{1 + e^{b+w_1x}} \quad (8)$$

# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

## Logistic function

This is the model equation for simple logistic regression:

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{b+w_1x}}{1 + e^{b+w_1x}} \quad (8)$$

where  $p(y = 1|\mathbf{x}, \boldsymbol{\theta}) \in [0, 1]$

# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

## Logistic function

This is the model equation for simple logistic regression:

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{b+w_1x}}{1 + e^{b+w_1x}} \quad (8)$$

where  $p(y = 1|\mathbf{x}, \boldsymbol{\theta}) \in [0, 1]$

The logistic function is a member of the class of **sigmoid** functions (S-shaped curves) and can also be written:

# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

## Logistic function

This is the model equation for simple logistic regression:

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{b+w_1x}}{1 + e^{b+w_1x}} \quad (8)$$

where  $p(y = 1|\mathbf{x}, \boldsymbol{\theta}) \in [0, 1]$

The logistic function is a member of the class of **sigmoid** functions (S-shaped curves) and can also be written:

$$p(y = 1|x, \boldsymbol{\theta}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}$$

# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

## Logistic function

This is the model equation for simple logistic regression:

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{b+w_1x}}{1 + e^{b+w_1x}} \quad (8)$$

where  $p(y = 1|\mathbf{x}, \boldsymbol{\theta}) \in [0, 1]$

The logistic function is a member of the class of **sigmoid** functions (S-shaped curves) and can also be written:

$$p(y = 1|x, \boldsymbol{\theta}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} = (1 + e^{-\mathbf{w}^\top \mathbf{x}})^{-1}$$



# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

## Logistic function

This is the model equation for simple logistic regression:

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{b+w_1x}}{1 + e^{b+w_1x}} \quad (8)$$

where  $p(y = 1|\mathbf{x}, \boldsymbol{\theta}) \in [0, 1]$

The logistic function is a member of the class of **sigmoid** functions (S-shaped curves) and can also be written:

$$p(y = 1|x, \boldsymbol{\theta}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} = (1 + e^{-\mathbf{w}^\top \mathbf{x}})^{-1} = \sigma(\mathbf{w}^\top \mathbf{x}) \quad (9)$$

# Logistic regression

Logistic regression is an approach for modeling the *probability* of a *multinomial* response.

In the simple case, we consider a binomial (or binary) response.

## Logistic function

This is the model equation for simple logistic regression:

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{b+w_1x}}{1 + e^{b+w_1x}} \quad (8)$$

where  $p(y = 1|\mathbf{x}, \boldsymbol{\theta}) \in [0, 1]$

The logistic function is a member of the class of **sigmoid** functions (S-shaped curves) and can also be written:

$$p(y = 1|x, \boldsymbol{\theta}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} = (1 + e^{-\mathbf{w}^\top \mathbf{x}})^{-1} = \sigma(\mathbf{w}^\top \mathbf{x}) \quad (9)$$

where  $\mathbf{w}^\top = (b, w_1)$

# Example 1: Binomial logistic regression with single predictor

# Example 1: Binomial logistic regression with single predictor

## Credit card defaults

# Example 1: Binomial logistic regression with single predictor

## Credit card defaults

- Data: A simulated data set containing information on ten thousand customers.

# Example 1: Binomial logistic regression with single predictor

## Credit card defaults

- Data: A simulated data set containing information on ten thousand customers.
- Question: Can we predict which customers will default on their credit card debt (based on income, etc)?

Four variables:

# Example 1: Binomial logistic regression with single predictor

## Credit card defaults

- Data: A simulated data set containing information on ten thousand customers.
- Question: Can we predict which customers will default on their credit card debt (based on income, etc)?

Four variables:

- **default**: A factor with levels *No* and *Yes* indicating whether the customer defaulted on their debt

# Example 1: Binomial logistic regression with single predictor

## Credit card defaults

- Data: A simulated data set containing information on ten thousand customers.
- Question: Can we predict which customers will default on their credit card debt (based on income, etc)?

Four variables:

- **default**: A factor with levels *No* and *Yes* indicating whether the customer defaulted on their debt
- **student**: A factor with levels *No* and *Yes* indicating whether the customer is a student



# Example 1: Binomial logistic regression with single predictor

## Credit card defaults

- Data: A simulated data set containing information on ten thousand customers.
- Question: Can we predict which customers will default on their credit card debt (based on income, etc)?

Four variables:

- **default**: A factor with levels *No* and *Yes* indicating whether the customer defaulted on their debt
- **student**: A factor with levels *No* and *Yes* indicating whether the customer is a student
- **balance**: The average balance that the customer has remaining on their credit card after making their monthly payment

# Example 1: Binomial logistic regression with single predictor

## Credit card defaults

- Data: A simulated data set containing information on ten thousand customers.
- Question: Can we predict which customers will default on their credit card debt (based on income, etc)?

Four variables:

- **default**: A factor with levels *No* and *Yes* indicating whether the customer defaulted on their debt
- **student**: A factor with levels *No* and *Yes* indicating whether the customer is a student
- **balance**: The average balance that the customer has remaining on their credit card after making their monthly payment
- **income**: Income of customer

# Example 1: Binomial logistic regression (cont.)

# Example 1: Binomial logistic regression (cont.)

We want to model the probability of **default** based on the **balance** predictor.

# Example 1: Binomial logistic regression (cont.)

We want to model the probability of **default** based on the **balance** predictor.  
The estimated coefficients from a computer program are:

# Example 1: Binomial logistic regression (cont.)

We want to model the probability of `default` based on the `balance` predictor.  
The estimated coefficients from a computer program are:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-1.065e+01	3.612e-01	-29.49	<2e-16	***
balance	5.499e-03	2.204e-04	24.95	<2e-16	***

# Example 1: Binomial logistic regression (cont.)

We want to model the probability of `default` based on the `balance` predictor.  
The estimated coefficients from a computer program are:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-1.065e+01	3.612e-01	-29.49	<2e-16 ***
balance	5.499e-03	2.204e-04	24.95	<2e-16 ***

Note that the null hypothesis for the tests is:  $H_0 : \mathbf{w}_i = 0$  (i.e. no dependence on the corresponding predictor)

## Example 1: Binomial logistic regression (cont.)



# Example 1: Binomial logistic regression (cont.)

The estimated model is:

# Example 1: Binomial logistic regression (cont.)

The estimated model is:

$$\hat{p}(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{(-10.65+0.0055x)}}{1 + e^{(-10.65+0.0055x)}} =$$

# Example 1: Binomial logistic regression (cont.)

The estimated model is:

$$\hat{p}(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{(-10.65+0.0055x)}}{1 + e^{(-10.65+0.0055x)}} = \frac{1}{1 + e^{(10.65-0.0055x)}} =$$

# Example 1: Binomial logistic regression (cont.)

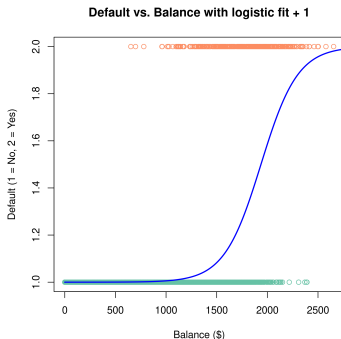
The estimated model is:

$$\hat{p}(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{e^{(-10.65+0.0055x)}}{1 + e^{(-10.65+0.0055x)}} = \frac{1}{1 + e^{(10.65-0.0055x)}} = \sigma(1+e^{(10.65-0.0055x)}) \quad (10)$$

# Example 1: Binomial logistic regression (cont.)

The estimated model is:

$$\hat{p}(y = 1 | \mathbf{x}; \boldsymbol{\theta}) = \frac{e^{(-10.65 + 0.0055x)}}{1 + e^{(-10.65 + 0.0055x)}} = \frac{1}{1 + e^{(10.65 - 0.0055x)}} = \sigma(1 + e^{(10.65 - 0.0055x)}) \quad (10)$$



# Example 1: Binomial logistic regression (cont.)

# Example 1: Binomial logistic regression (cont.)

Recall the model:

# Example 1: Binomial logistic regression (cont.)

Recall the model:

$$\hat{p}(y = 1|x) = \frac{e^{-10.65+0.0055x}}{1 + e^{-10.65+0.0055x}}$$



# Example 1: Binomial logistic regression (cont.)

Recall the model:

$$\hat{p}(y = 1|x) = \frac{e^{-10.65+0.0055x}}{1 + e^{-10.65+0.0055x}}$$

- ① How would  $\hat{p}$  (the predicted probability) change if  $x$  were to increase by \$100?

# Example 1: Binomial logistic regression (cont.)

Recall the model:

$$\hat{p}(y = 1|x) = \frac{e^{-10.65+0.0055x}}{1 + e^{-10.65+0.0055x}}$$

- ① How would  $\hat{p}$  (the predicted probability) change if  $x$  were to increase by \$100?
- ② What about if  $x$  were to decrease by \$100

# Example 1: Binomial logistic regression (cont.)

Recall the model:

$$\hat{p}(y = 1|x) = \frac{e^{-10.65+0.0055x}}{1 + e^{-10.65+0.0055x}}$$

- ① How would  $\hat{p}$  (the predicted probability) change if  $x$  were to increase by \$100?
- ② What about if  $x$  were to decrease by \$100?
- ③ There are 333 defaults out of 10000 observations. What is the impact of  $b$ ?

# Multiple logistic regression

# Multiple logistic regression

In multiple logistic regression, we predict a binary response using *multiple predictors*:

# Multiple logistic regression

In multiple logistic regression, we predict a binary response using *multiple predictors*:

$$\log \left( \frac{p(x)}{1 - p(x)} \right) =$$

# Multiple logistic regression

In multiple logistic regression, we predict a binary response using *multiple predictors*:

$$\log \left( \frac{p(x)}{1 - p(x)} \right) = b + w_1 x_1 + \cdots + \mathbf{w}_D x_D =$$

# Multiple logistic regression

In multiple logistic regression, we predict a binary response using *multiple predictors*:

$$\log \left( \frac{p(x)}{1 - p(x)} \right) = b + w_1 x_1 + \cdots + w_D x_D = \mathbf{w}^\top \mathbf{X} \quad (11)$$



# Multiple logistic regression

In multiple logistic regression, we predict a binary response using *multiple predictors*:

$$\log \left( \frac{p(\mathbf{x})}{1 - p(\mathbf{x})} \right) = b + w_1 x_1 + \cdots + w_D x_D = \mathbf{w}^\top \mathbf{X} \quad (11)$$

where  $\mathbf{x} = (x_1, \dots, x_D)$ .

# Multiple logistic regression

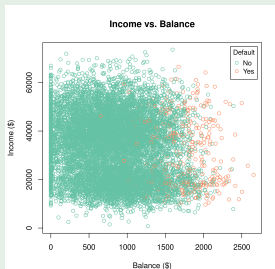
In multiple logistic regression, we predict a binary response using *multiple predictors*:

$$\log \left( \frac{p(\mathbf{x})}{1 - p(\mathbf{x})} \right) = b + w_1 x_1 + \cdots + w_D x_D = \mathbf{w}^\top \mathbf{X} \quad (11)$$

where  $\mathbf{x} = (x_1, \dots, x_D)$ .

## Activity: Binomial logistic regression with multiple predictors

Using the Default dataset, predict the probability of **default** based on **balance** and **student**. Comment on your results and interpret the coefficient estimates.



# Decision boundary

# Decision boundary

This is the line that defines the probability threshold  $\tau$  for class assignment:

# Decision boundary

This is the line that defines the probability threshold  $\tau$  for class assignment:  
In 1-D:

# Decision boundary

This is the line that defines the probability threshold  $\tau$  for class assignment:  
In 1-D:

$$x^* : p(y = 1|x = x^*, \theta) = \tau \quad (12)$$

# Decision boundary

This is the line that defines the probability threshold  $\tau$  for class assignment:  
In 1-D:

$$x^* : p(y = 1 | x = x^*, \theta) = \tau \quad (12)$$

Typically,  $\tau = 0.5$

# Estimation of logistic regression coefficients



# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

$$\mathcal{L}(\theta) = \prod_{i=1}^N p_{c_i}(x_i; \theta) \quad (13)$$

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

$$\mathcal{L}(\theta) = \prod_{i=1}^N p_{c_i}(x_i; \theta) \quad (13)$$

- where  $p_{c_i}(x_i; \theta) = \Pr(G = c_i | X = x_i; \theta)$

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

$$\mathcal{L}(\theta) = \prod_{i=1}^N p_{c_i}(x_i; \theta) \quad (13)$$

- where  $p_{c_i}(x_i; \theta) = \Pr(G = c_i | X = x_i; \theta)$
- In the two-class case:  $\theta = \mathbf{w} = \{b, w_1\}$

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

$$\mathcal{L}(\theta) = \prod_{i=1}^N p_{c_i}(x_i; \theta) \quad (13)$$

- where  $p_{c_i}(x_i; \theta) = \Pr(G = c_i | X = x_i; \theta)$
- In the two-class case:  $\theta = \mathbf{w} = \{b, w_1\}$
- Thus, we can write the conditional probabilities as:

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

$$\mathcal{L}(\theta) = \prod_{i=1}^N p_{c_i}(x_i; \theta) \quad (13)$$

- where  $p_{c_i}(x_i; \theta) = \Pr(G = c_i | X = x_i; \theta)$
- In the two-class case:  $\theta = \mathbf{w} = \{b, w_1\}$
- Thus, we can write the conditional probabilities as:



# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

$$\mathcal{L}(\theta) = \prod_{i=1}^N p_{c_i}(x_i; \theta) \quad (13)$$

- where  $p_{c_i}(x_i; \theta) = \Pr(G = c_i | X = x_i; \theta)$
- In the two-class case:  $\theta = \mathbf{w} = \{b, w_1\}$
- Thus, we can write the conditional probabilities as:

$$p_1(x; \theta) = p(x; \theta) \quad (14)$$

$$p_2(x; \theta) = 1 - p(x; \theta) \quad (15)$$

- It is also convenient to encode  $c_i$  using a 0/1 response  $y_i$ :

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

$$\mathcal{L}(\theta) = \prod_{i=1}^N p_{c_i}(x_i; \theta) \quad (13)$$

- where  $p_{c_i}(x_i; \theta) = \Pr(G = c_i | X = x_i; \theta)$
- In the two-class case:  $\theta = \mathbf{w} = \{b, w_1\}$
- Thus, we can write the conditional probabilities as:

$$p_1(x; \theta) = p(x; \theta) \quad (14)$$

$$p_2(x; \theta) = 1 - p(x; \theta) \quad (15)$$

- It is also convenient to encode  $c_i$  using a 0/1 response  $y_i$ :

# Estimation of logistic regression coefficients

The method of **maximum likelihood** is used to estimate logistic regression coefficients  $\mathbf{w}$

- The likelihood function  $\mathcal{L}(\theta)$  represents the support provided by a sample for a given parameter  $\theta$ :

$$\mathcal{L}(\theta) = \prod_{i=1}^N p_{c_i}(x_i; \theta) \quad (13)$$

- where  $p_{c_i}(x_i; \theta) = \Pr(G = c_i | X = x_i; \theta)$
- In the two-class case:  $\theta = \mathbf{w} = \{b, w_1\}$
- Thus, we can write the conditional probabilities as:

$$p_1(x; \theta) = p(x; \theta) \quad (14)$$

$$p_2(x; \theta) = 1 - p(x; \theta) \quad (15)$$

- It is also convenient to encode  $c_i$  using a 0/1 response  $y_i$ :

$$y_i = \begin{cases} 1, & \text{when } c_i = \text{Class 1} \\ 0, & \text{when } c_i = \text{Class 2} \end{cases} \quad (16)$$

# Log-likelihood function for logistic regression

# Log-likelihood function for logistic regression

The principle of maximum likelihood dictates that the best parameter estimates are those that maximize the likelihood function.

# Log-likelihood function for logistic regression

The principle of maximum likelihood dictates that the best parameter estimates are those that maximize the likelihood function.

- Equivalently, we *minimize* the **negative log-likelihood** function  $NLL(\theta)$ :

# Log-likelihood function for logistic regression

The principle of maximum likelihood dictates that the best parameter estimates are those that maximize the likelihood function.

- Equivalently, we *minimize* the **negative log-likelihood** function  $NLL(\theta)$ :

$$NLL(\theta) = - \sum_i \log p_{c_i}(x_i; \theta) \quad (17)$$

# Log-likelihood function for logistic regression

The principle of maximum likelihood dictates that the best parameter estimates are those that maximize the likelihood function.

- Equivalently, we *minimize* the **negative log-likelihood** function  $NLL(\theta)$ :

$$NLL(\theta) = - \sum_i \log p_{c_i}(x_i; \theta) \quad (17)$$

- In the binomial case, this simplifies to:



# Log-likelihood function for logistic regression

The principle of maximum likelihood dictates that the best parameter estimates are those that maximize the likelihood function.

- Equivalently, we *minimize* the **negative log-likelihood** function  $NLL(\theta)$ :

$$NLL(\theta) = - \sum_i \log p_{c_i}(x_i; \theta) \quad (17)$$

- In the binomial case, this simplifies to:

$$NLL(\mathbf{w}) = - \sum_i [y_i \log p(x_i; \mathbf{w}) + (1 - y_i) \log(1 - p(x_i; \mathbf{w}))] \quad (18)$$

# Log-likelihood function for logistic regression

The principle of maximum likelihood dictates that the best parameter estimates are those that maximize the likelihood function.

- Equivalently, we *minimize* the **negative log-likelihood** function  $NLL(\theta)$ :

$$NLL(\theta) = - \sum_i \log p_{c_i}(x_i; \theta) \quad (17)$$

- In the binomial case, this simplifies to:

$$NLL(\mathbf{w}) = - \sum_i [y_i \log p(x_i; \mathbf{w}) + (1 - y_i) \log(1 - p(x_i; \mathbf{w}))] \quad (18)$$

- Recall that we model  $p(x_i; \mathbf{w})$  as:

# Log-likelihood function for logistic regression

The principle of maximum likelihood dictates that the best parameter estimates are those that maximize the likelihood function.

- Equivalently, we *minimize* the **negative log-likelihood** function  $NLL(\theta)$ :

$$NLL(\theta) = - \sum_i \log p_{\theta_i}(x_i; \theta) \quad (17)$$

- In the binomial case, this simplifies to:

$$NLL(\mathbf{w}) = - \sum_i [y_i \log p(x_i; \mathbf{w}) + (1 - y_i) \log(1 - p(x_i; \mathbf{w}))] \quad (18)$$

- Recall that we model  $p(x_i; \mathbf{w})$  as:

$$p(x_i) = \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \quad (19)$$

# Log-likelihood function for logistic regression

# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

$$\text{NLL}(\mathbf{w}) = - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( 1 - \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) \right]$$

# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

$$\begin{aligned}\text{NLL}(\mathbf{w}) &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( 1 - \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) \right] \\ &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( \frac{1}{1 + e^{b+w_1x_i}} \right) \right]\end{aligned}$$

# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

$$\begin{aligned}\text{NLL}(\mathbf{w}) &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( 1 - \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) \right] \\ &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( \frac{1}{1 + e^{b+w_1x_i}} \right) \right] \\ &= - \sum_i \left[ y_i \log (e^{b+w_1x_i}) - y_i \log (1 + e^{b+w_1x_i}) \right]\end{aligned}$$



# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

$$\begin{aligned}\text{NLL}(\mathbf{w}) &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( 1 - \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) \right] \\ &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( \frac{1}{1 + e^{b+w_1x_i}} \right) \right] \\ &= - \sum_i \left[ y_i \log (e^{b+w_1x_i}) - y_i \log (1 + e^{b+w_1x_i}) \right. \\ &\quad \left. + \cancel{(1 - y_i) \log(1)}^0 - (1 - y_i) \log (1 + e^{b+w_1x_i}) \right]\end{aligned}$$

# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

$$\begin{aligned}\text{NLL}(\mathbf{w}) &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( 1 - \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) \right] \\&= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( \frac{1}{1 + e^{b+w_1x_i}} \right) \right] \\&= - \sum_i \left[ y_i \log (e^{b+w_1x_i}) - y_i \log (1 + e^{b+w_1x_i}) \right. \\&\quad \left. + \cancel{(1 - y_i) \log(1)}^0 - (1 - y_i) \log (1 + e^{b+w_1x_i}) \right] \\&= - \sum_i \left[ y_i (b + w_1x_i) - y_i \log (1 + e^{b+w_1x_i}) - \log (1 + e^{b+w_1x_i}) \right]\end{aligned}$$

# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

$$\begin{aligned}
 \text{NLL}(\mathbf{w}) &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( 1 - \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) \right] \\
 &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( \frac{1}{1 + e^{b+w_1x_i}} \right) \right] \\
 &= - \sum_i \left[ y_i \log (e^{b+w_1x_i}) - y_i \log (1 + e^{b+w_1x_i}) \right. \\
 &\quad \left. + \cancel{(1 - y_i) \log(1)}^0 - (1 - y_i) \log (1 + e^{b+w_1x_i}) \right] \\
 &= - \sum_i \left[ y_i (b + w_1x_i) - y_i \log (1 + e^{b+w_1x_i}) - \log (1 + e^{b+w_1x_i}) \right. \\
 &\quad \left. + y_i (1 + e^{b+w_1x_i}) \right]
 \end{aligned}$$

# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

$$\begin{aligned}
 \text{NLL}(\mathbf{w}) &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( 1 - \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) \right] \\
 &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( \frac{1}{1 + e^{b+w_1x_i}} \right) \right] \\
 &= - \sum_i \left[ y_i \log (e^{b+w_1x_i}) - y_i \log (1 + e^{b+w_1x_i}) \right. \\
 &\quad \left. + \cancel{(1 - y_i) \log(1)}^0 - (1 - y_i) \log (1 + e^{b+w_1x_i}) \right] \tag{20} \\
 &= - \sum_i \left[ y_i (b + w_1x_i) - \cancel{y_i \log (1 + e^{b+w_1x_i})} - \log (1 + e^{b+w_1x_i}) \right. \\
 &\quad \left. + y_i (1 + e^{b+w_1x_i}) \right] \\
 \text{NLL}(\mathbf{w}) &= - \sum_i [y_i (b + w_1x_i) - \log (1 + e^{b+w_1x_i})]
 \end{aligned}$$

# Log-likelihood function for logistic regression

Substituting (19) into (18), we obtain:

$$\begin{aligned}
 \text{NLL}(\mathbf{w}) &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( 1 - \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) \right] \\
 &= - \sum_i \left[ y_i \log \left( \frac{e^{b+w_1x_i}}{1 + e^{b+w_1x_i}} \right) + (1 - y_i) \log \left( \frac{1}{1 + e^{b+w_1x_i}} \right) \right] \\
 &= - \sum_i \left[ y_i \log (e^{b+w_1x_i}) - y_i \log (1 + e^{b+w_1x_i}) \right. \\
 &\quad \left. + \cancel{(1 - y_i) \log(1)}^0 - (1 - y_i) \log (1 + e^{b+w_1x_i}) \right] \tag{20} \\
 &= - \sum_i \left[ y_i (b + w_1x_i) - \cancel{y_i \log (1 + e^{b+w_1x_i})} - \log (1 + e^{b+w_1x_i}) \right. \\
 &\quad \left. + \cancel{y_i (1 + e^{b+w_1x_i})} \right] \\
 \text{NLL}(\mathbf{w}) &= - \sum_i [y_i (b + w_1x_i) - \log (1 + e^{b+w_1x_i})]
 \end{aligned}$$

# Maximizing the log-likelihood

# Maximizing the log-likelihood

To find  $\hat{\mathbf{w}}$ , we find the derivative of  $\text{NLL}(\mathbf{w})$ , set it to zero and solve the resulting **score equations**:

# Maximizing the log-likelihood

To find  $\hat{\mathbf{w}}$ , we find the derivative of  $\text{NLL}(\mathbf{w})$ , set it to zero and solve the resulting **score equations**:

$$\frac{\partial \text{NLL}}{\partial \mathbf{w}} = -\frac{\partial}{\partial \mathbf{w}} \sum_i [y_i (b + w_1 x_i) - \log(1 + e^{b + w_1 x_i})]$$



# Maximizing the log-likelihood

To find  $\hat{\mathbf{w}}$ , we find the derivative of  $\text{NLL}(\mathbf{w})$ , set it to zero and solve the resulting **score equations**:

$$\frac{\partial \text{NLL}}{\partial \mathbf{w}} = - \frac{\partial}{\partial \mathbf{w}} \sum_i [y_i (b + w_1 x_i) - \log(1 + e^{b + w_1 x_i})]$$
$$\begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} = - \begin{pmatrix} \sum_i \left[ y_i - \frac{e^{b + w_1 x_i}}{1 + e^{b + w_1 x_i}} \right] \\ \sum_i \left[ x_i y_i - \frac{x_i (e^{b + w_1 x_i})}{1 + e^{b + w_1 x_i}} \right] \end{pmatrix}$$

# Maximizing the log-likelihood

To find  $\hat{\mathbf{w}}$ , we find the derivative of  $\text{NLL}(\mathbf{w})$ , set it to zero and solve the resulting **score equations**:

$$\frac{\partial \text{NLL}}{\partial \mathbf{w}} = - \frac{\partial}{\partial \mathbf{w}} \sum_i [y_i (b + w_1 x_i) - \log(1 + e^{b + w_1 x_i})]$$

$$\begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} = - \begin{pmatrix} \sum_i \left[ y_i - \frac{e^{b + w_1 x_i}}{1 + e^{b + w_1 x_i}} \right] \\ \sum_i \left[ x_i y_i - \frac{x_i (e^{b + w_1 x_i})}{1 + e^{b + w_1 x_i}} \right] \end{pmatrix}$$

$$\begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} = - \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix}$$

# Maximizing the log-likelihood

To find  $\hat{\mathbf{w}}$ , we find the derivative of  $\text{NLL}(\mathbf{w})$ , set it to zero and solve the resulting **score equations**:

$$\begin{aligned}\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= -\frac{\partial}{\partial \mathbf{w}} \sum_i [y_i (b + w_1 x_i) - \log(1 + e^{b + w_1 x_i})] \\ \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} &= - \begin{pmatrix} \sum_i \left[ y_i - \frac{e^{b + w_1 x_i}}{1 + e^{b + w_1 x_i}} \right] \\ \sum_i \left[ x_i y_i - \frac{x_i (e^{b + w_1 x_i})}{1 + e^{b + w_1 x_i}} \right] \end{pmatrix} \\ \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} &= - \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}\end{aligned}\tag{21}$$

# Maximizing the log-likelihood

To find  $\hat{\mathbf{w}}$ , we find the derivative of  $\text{NLL}(\mathbf{w})$ , set it to zero and solve the resulting **score equations**:

$$\begin{aligned}\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= -\frac{\partial}{\partial \mathbf{w}} \sum_i [y_i (b + w_1 x_i) - \log(1 + e^{b + w_1 x_i})] \\ \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} &= - \begin{pmatrix} \sum_i \left[ y_i - \frac{e^{b + w_1 x_i}}{1 + e^{b + w_1 x_i}} \right] \\ \sum_i \left[ x_i y_i - \frac{x_i (e^{b + w_1 x_i})}{1 + e^{b + w_1 x_i}} \right] \end{pmatrix} \\ \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} &= - \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}\end{aligned}\tag{21}$$

This is a system of two **nonlinear** equations in  $\mathbf{w}$

# Maximizing the log-likelihood

To find  $\hat{\mathbf{w}}$ , we find the derivative of  $\text{NLL}(\mathbf{w})$ , set it to zero and solve the resulting **score equations**:

$$\begin{aligned}\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= -\frac{\partial}{\partial \mathbf{w}} \sum_i [y_i (b + w_1 x_i) - \log(1 + e^{b + w_1 x_i})] \\ \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} &= - \begin{pmatrix} \sum_i \left[ y_i - \frac{e^{b + w_1 x_i}}{1 + e^{b + w_1 x_i}} \right] \\ \sum_i \left[ x_i y_i - \frac{x_i (e^{b + w_1 x_i})}{1 + e^{b + w_1 x_i}} \right] \end{pmatrix} \\ \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} &= - \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}\end{aligned}\tag{21}$$

This is a system of two **nonlinear** equations in  $\mathbf{w}$  which can be solved via the **Newton-Raphson** method.

# Maximizing the log-likelihood

To find  $\hat{\mathbf{w}}$ , we find the derivative of  $\text{NLL}(\mathbf{w})$ , set it to zero and solve the resulting **score equations**:

$$\begin{aligned}\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= -\frac{\partial}{\partial \mathbf{w}} \sum_i [y_i (b + w_1 x_i) - \log(1 + e^{b + w_1 x_i})] \\ \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} &= - \begin{pmatrix} \sum_i \left[ y_i - \frac{e^{b + w_1 x_i}}{1 + e^{b + w_1 x_i}} \right] \\ \sum_i \left[ x_i y_i - \frac{x_i (e^{b + w_1 x_i})}{1 + e^{b + w_1 x_i}} \right] \end{pmatrix} \\ \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} &= - \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}\end{aligned}\tag{21}$$

This is a system of two **nonlinear** equations in  $\mathbf{w}$  which can be solved via the **Newton-Raphson** method.

Alternatively, we can use the **gradient descent** approach to directly minimize NLL.

# Maximum likelihood estimation (MLE) in logistic regression

# Maximum likelihood estimation (MLE) in logistic regression

Recall the negative log-likelihood function for the binomial logistic regression case:



# Maximum likelihood estimation (MLE) in logistic regression

Recall the negative log-likelihood function for the binomial logistic regression case:

$$\text{NLL}(\mathbf{w}) = - \sum_i [y_i (b + w_1 x_i) - \log (1 + e^{b + w_1 x_i})] \quad (22)$$

# Maximum likelihood estimation (MLE) in logistic regression

Recall the negative log-likelihood function for the binomial logistic regression case:

$$\text{NLL}(\mathbf{w}) = - \sum_i [y_i (b + w_1 x_i) - \log (1 + e^{b + w_1 x_i})] \quad (22)$$

The optimal  $\hat{\mathbf{w}}$  which minimizes  $\text{NLL}(\mathbf{w})$  is the **maximum likelihood estimate**.

# Maximum likelihood estimation (MLE) in logistic regression

Recall the negative log-likelihood function for the binomial logistic regression case:

$$\text{NLL}(\mathbf{w}) = - \sum_i [y_i (b + w_1 x_i) - \log (1 + e^{b + w_1 x_i})] \quad (22)$$

The optimal  $\hat{\mathbf{w}}$  which minimizes  $\text{NLL}(\mathbf{w})$  is the **maximum likelihood estimate**.

Also recall the derivative of NLL:

# Maximum likelihood estimation (MLE) in logistic regression

Recall the negative log-likelihood function for the binomial logistic regression case:

$$\text{NLL}(\mathbf{w}) = - \sum_i [y_i (b + w_1 x_i) - \log (1 + e^{b + w_1 x_i})] \quad (22)$$

The optimal  $\hat{\mathbf{w}}$  which minimizes  $\text{NLL}(\mathbf{w})$  is the **maximum likelihood estimate**.

Also recall the derivative of NLL:

$$\nabla_{\mathbf{w}} \text{NLL} =$$

# Maximum likelihood estimation (MLE) in logistic regression

Recall the negative log-likelihood function for the binomial logistic regression case:

$$\text{NLL}(\mathbf{w}) = - \sum_i [y_i (b + w_1 x_i) - \log (1 + e^{b + w_1 x_i})] \quad (22)$$

The optimal  $\hat{\mathbf{w}}$  which minimizes  $\text{NLL}(\mathbf{w})$  is the **maximum likelihood estimate**.

Also recall the derivative of NLL:

$$\nabla_{\mathbf{w}} \text{NLL} = \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} =$$

# Maximum likelihood estimation (MLE) in logistic regression

Recall the negative log-likelihood function for the binomial logistic regression case:

$$\text{NLL}(\mathbf{w}) = - \sum_i [y_i (b + w_1 x_i) - \log (1 + e^{b + w_1 x_i})] \quad (22)$$

The optimal  $\hat{\mathbf{w}}$  which minimizes  $\text{NLL}(\mathbf{w})$  is the **maximum likelihood estimate**.

Also recall the derivative of NLL:

$$\nabla_{\mathbf{w}} \text{NLL} = \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} = \begin{pmatrix} - \sum_i [y_i - p(x_i)] \\ - \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \quad (23)$$

# Maximum likelihood estimation (MLE) in logistic regression

Recall the negative log-likelihood function for the binomial logistic regression case:

$$\text{NLL}(\mathbf{w}) = - \sum_i [y_i (b + w_1 x_i) - \log (1 + e^{b + w_1 x_i})] \quad (22)$$

The optimal  $\hat{\mathbf{w}}$  which minimizes  $\text{NLL}(\mathbf{w})$  is the **maximum likelihood estimate**.

Also recall the derivative of NLL:

$$\nabla_{\mathbf{w}} \text{NLL} = \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} = \begin{pmatrix} - \sum_i [y_i - p(x_i)] \\ - \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \quad (23)$$

We can use either Newton-Raphson or gradient *descent* to *minimize* NLL.

# NLL and entropy



# NLL and entropy

We can show that the NLL is equal to the sum of the **binary cross entropy** of  $y_i$  and  $p(y = 1|\mathbf{x}_i)$  over  $N$ :

# NLL and entropy

We can show that the NLL is equal to the sum of the **binary cross entropy** of  $y_i$  and  $p(y = 1|\mathbf{x}_i)$  over  $N$ :

$$\mathbb{H}_i(y_i, p_i) = -[y_i \log p_i + (1 - y_i) \log(1 - p_i)] \quad (24)$$

# NLL and entropy

We can show that the NLL is equal to the sum of the **binary cross entropy** of  $y_i$  and  $p(y = 1|\mathbf{x}_i)$  over  $N$ :

$$\mathbb{H}_i(y_i, p_i) = -[y_i \log p_i + (1 - y_i) \log(1 - p_i)] \quad (24)$$

Note that  $p_i = \sigma(\mathbf{w}^\top \mathbf{x}_i)$ .

- Binary cross-entropy quantifies how far your predicted probabilities are from the actual binary labels.

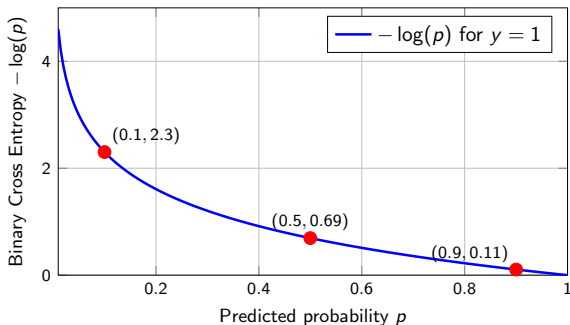
# Binary Cross Entropy vs. $p$ (when $y=1$ )

# Binary Cross Entropy vs. $p$ (when $y=1$ )

For a true label  $y = 1$ , the binary cross entropy is  $\mathbb{H}(1, p) = -\log(p)$ .

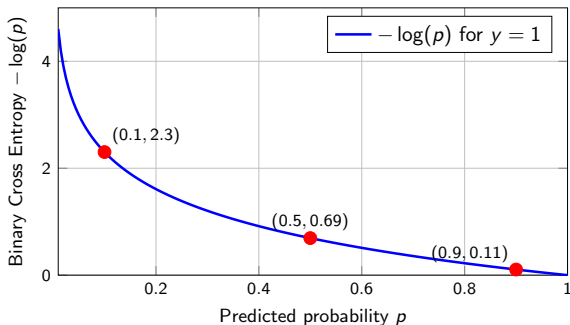
# Binary Cross Entropy vs. $p$ (when $y=1$ )

For a true label  $y = 1$ , the binary cross entropy is  $\mathbb{H}(1, p) = -\log(p)$ . This function penalizes predictions that are far from the true label:



# Binary Cross Entropy vs. $p$ (when $y=1$ )

For a true label  $y = 1$ , the binary cross entropy is  $\mathbb{H}(1, p) = -\log(p)$ . This function penalizes predictions that are far from the true label:

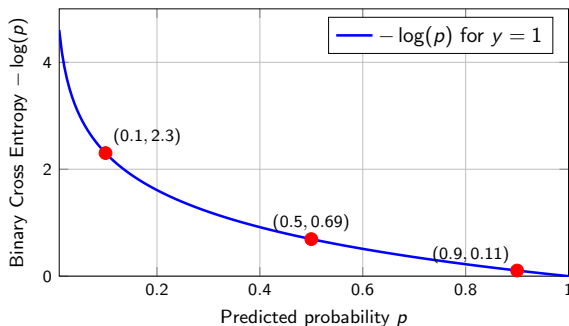


## Key observations:

- As  $p \rightarrow 1$ : cross entropy  $\rightarrow 0$  (low penalty for correct prediction)

# Binary Cross Entropy vs. $p$ (when $y=1$ )

For a true label  $y = 1$ , the binary cross entropy is  $\mathbb{H}(1, p) = -\log(p)$ . This function penalizes predictions that are far from the true label:



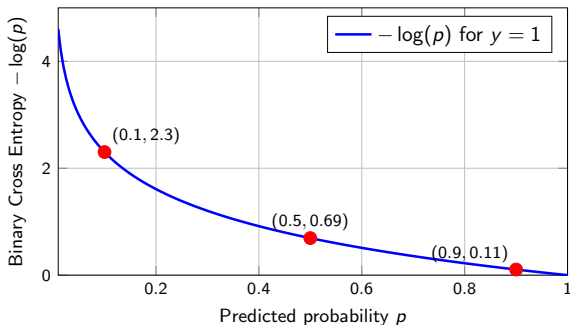
## Key observations:

- As  $p \rightarrow 1$ : cross entropy  $\rightarrow 0$  (low penalty for correct prediction)
- As  $p \rightarrow 0$ : cross entropy  $\rightarrow \infty$  (high penalty for incorrect prediction)



# Binary Cross Entropy vs. $p$ (when $y=1$ )

For a true label  $y = 1$ , the binary cross entropy is  $\mathbb{H}(1, p) = -\log(p)$ . This function penalizes predictions that are far from the true label:



## Key observations:

- As  $p \rightarrow 1$ : cross entropy  $\rightarrow 0$  (low penalty for correct prediction)
- As  $p \rightarrow 0$ : cross entropy  $\rightarrow \infty$  (high penalty for incorrect prediction)
- The function is convex, ensuring unique minimum in optimization

# Example: predicting spam

# Example: predicting spam

- a If the email is spam ( $y=1$ ) and you predict 90% probability of spam ( $p=0.9$ ), find the binary cross entropy (BCE):

# Example: predicting spam

- a If the email is spam ( $y=1$ ) and you predict 90% probability of spam ( $p=0.9$ ), find the binary cross entropy (BCE):

$$\text{BCE} = -[1 \times \log(0.9) + 0 \times \log(0.1)] = -\log(0.9) \approx 0.105 \quad (\text{low loss - good!})$$

# Example: predicting spam

- a If the email is spam ( $y=1$ ) and you predict 90% probability of spam ( $p=0.9$ ), find the binary cross entropy (BCE):

$$\text{BCE} = -[1 \times \log(0.9) + 0 \times \log(0.1)] = -\log(0.9) \approx 0.105 \quad (\text{low loss - good!})$$

- b If the email is spam ( $y=1$ ) but you predict only 10% probability of spam ( $p=0.1$ ), find the BCE.

# Example: predicting spam

- a If the email is spam ( $y=1$ ) and you predict 90% probability of spam ( $p=0.9$ ), find the binary cross entropy (BCE):

$$\text{BCE} = -[1 \times \log(0.9) + 0 \times \log(0.1)] = -\log(0.9) \approx 0.105 \quad (\text{low loss - good!})$$

- b If the email is spam ( $y=1$ ) but you predict only 10% probability of spam ( $p=0.1$ ), find the BCE.

$$\text{BCE} = -[1 \times \log(0.1) + 0 \times \log(0.9)] = -\log(0.1) \approx 2.303 \quad (\text{high loss - bad!})$$

# Gradient descent for MLE in logistic regression

# Gradient descent for MLE in logistic regression

This approach only requires the first derivative:



# Gradient descent for MLE in logistic regression

This approach only requires the first derivative:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \rho \nabla \text{NLL}(\mathbf{w}_k) \quad (25)$$

# Gradient descent for MLE in logistic regression

This approach only requires the first derivative:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \rho \nabla \text{NLL}(\mathbf{w}_k) \quad (25)$$

Thus, to find  $\hat{\mathbf{w}}$  we iterate using:

# Gradient descent for MLE in logistic regression

This approach only requires the first derivative:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \rho \nabla \text{NLL}(\mathbf{w}_k) \quad (25)$$

Thus, to find  $\hat{\mathbf{w}}$  we iterate using:

$$\begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} =$$

# Gradient descent for MLE in logistic regression

This approach only requires the first derivative:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \rho \nabla \text{NLL}(\mathbf{w}_k) \quad (25)$$

Thus, to find  $\hat{\mathbf{w}}$  we iterate using:

$$\begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} + \rho \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \quad (26)$$

# Gradient descent for MLE in logistic regression

This approach only requires the first derivative:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \rho \nabla \text{NLL}(\mathbf{w}_k) \quad (25)$$

Thus, to find  $\hat{\mathbf{w}}$  we iterate using:

$$\begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} + \rho \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \quad (26)$$

Because the negative log-likelihood is *convex*, and thus a *minimization* problem, we *descend* the function and thus *subtract* the scaled derivative.

# Gradient descent for MLE in logistic regression

This approach only requires the first derivative:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \rho \nabla \text{NLL}(\mathbf{w}_k) \quad (25)$$

Thus, to find  $\hat{\mathbf{w}}$  we iterate using:

$$\begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} + \rho \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \quad (26)$$

Because the negative log-likelihood is *convex*, and thus a *minimization* problem, we *descend* the function and thus *subtract* the scaled derivative.

## Note

- The gradient descent method does not require a second derivative

# Gradient descent for MLE in logistic regression

This approach only requires the first derivative:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \rho \nabla \text{NLL}(\mathbf{w}_k) \quad (25)$$

Thus, to find  $\hat{\mathbf{w}}$  we iterate using:

$$\begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} + \rho \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \quad (26)$$

Because the negative log-likelihood is *convex*, and thus a *minimization* problem, we *descend* the function and thus *subtract* the scaled derivative.

## Note

- The gradient descent method does not require a second derivative
- However, it may require more iterations to converge than Newton-Raphson

# Newton-Raphson approach for MLE in logistic regression



# Newton-Raphson approach for MLE in logistic regression

The optimal point  $\hat{\mathbf{w}}$  is given by the root of the equation  $\nabla_{\mathbf{w}} \text{NLL} = 0$ .

# Newton-Raphson approach for MLE in logistic regression

The optimal point  $\hat{\mathbf{w}}$  is given by the root of the equation  $\nabla_{\mathbf{w}} \text{NLL} = 0$ .

Applying Newton-Raphson, the update step is:

# Newton-Raphson approach for MLE in logistic regression

The optimal point  $\hat{\mathbf{w}}$  is given by the root of the equation  $\nabla_{\mathbf{w}} \text{NLL} = 0$ .

Applying Newton-Raphson, the update step is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}_{\mathbf{w}_k}^{-1}(\text{NLL}) \nabla_{\mathbf{w}_k} \text{NLL}(\mathbf{w}_k) \quad (27)$$

# Newton-Raphson approach for MLE in logistic regression

The optimal point  $\hat{\mathbf{w}}$  is given by the root of the equation  $\nabla_{\mathbf{w}} \text{NLL} = 0$ .

Applying Newton-Raphson, the update step is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}_{\mathbf{w}_k}^{-1}(\text{NLL}) \nabla_{\mathbf{w}_k} \text{NLL}(\mathbf{w}_k) \quad (27)$$

The operator  $\mathbf{H}^{-1}$  represents the inverse **Hessian** (second derivative) matrix of NLL with respect to  $\mathbf{w}$ :

# Newton-Raphson approach for MLE in logistic regression

The optimal point  $\hat{\mathbf{w}}$  is given by the root of the equation  $\nabla_{\mathbf{w}} \text{NLL} = 0$ .

Applying Newton-Raphson, the update step is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}_{\mathbf{w}_k}^{-1}(\text{NLL}) \nabla_{\mathbf{w}_k} \text{NLL}(\mathbf{w}_k) \quad (27)$$

The operator  $\mathbf{H}^{-1}$  represents the inverse **Hessian** (second derivative) matrix of NLL with respect to  $\mathbf{w}$ :

$$\mathbf{H}_{\mathbf{w}_k}(\text{NLL}) = \nabla_{\mathbf{w}_k}^2 \text{NLL} = \begin{pmatrix} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} \\ \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1 \partial b} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} \end{pmatrix} \quad (28)$$

# Newton-Raphson approach for MLE in logistic regression

The optimal point  $\hat{\mathbf{w}}$  is given by the root of the equation  $\nabla_{\mathbf{w}} \text{NLL} = 0$ .

Applying Newton-Raphson, the update step is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}_{\mathbf{w}_k}^{-1}(\text{NLL}) \nabla_{\mathbf{w}_k} \text{NLL}(\mathbf{w}_k) \quad (27)$$

The operator  $\mathbf{H}^{-1}$  represents the inverse **Hessian** (second derivative) matrix of NLL with respect to  $\mathbf{w}$ :

$$\mathbf{H}_{\mathbf{w}_k}(\text{NLL}) = \nabla_{\mathbf{w}_k}^2 \text{NLL} = \begin{pmatrix} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} \\ \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1 \partial b} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} \end{pmatrix} \quad (28)$$

Note that (27) is just the matrix representation of the 1-D case:

# Newton-Raphson approach for MLE in logistic regression

The optimal point  $\hat{\mathbf{w}}$  is given by the root of the equation  $\nabla_{\mathbf{w}} \text{NLL} = 0$ .

Applying Newton-Raphson, the update step is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}_{\mathbf{w}_k}^{-1}(\text{NLL}) \nabla_{\mathbf{w}_k} \text{NLL}(\mathbf{w}_k) \quad (27)$$

The operator  $\mathbf{H}^{-1}$  represents the inverse **Hessian** (second derivative) matrix of NLL with respect to  $\mathbf{w}$ :

$$\mathbf{H}_{\mathbf{w}_k}(\text{NLL}) = \nabla_{\mathbf{w}_k}^2 \text{NLL} = \begin{pmatrix} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} \\ \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1 \partial b} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} \end{pmatrix} \quad (28)$$

Note that (27) is just the matrix representation of the 1-D case:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\text{NLL}'(\mathbf{w}_k)}{\text{NLL}''(\mathbf{w}_k)} \quad (29)$$

# Newton-Raphson approach for MLE (cont.)



# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b w_1} = \sum_i x_i p(x_i)(1 - p(x_i)) \quad (31)$$

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} = \sum_i x_i p(x_i)(1 - p(x_i)) \quad (31)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} = \sum_i x_i^2 p(x_i)(1 - p(x_i)) \quad (32)$$

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} = \sum_i x_i p(x_i)(1 - p(x_i)) \quad (31)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} = \sum_i x_i^2 p(x_i)(1 - p(x_i)) \quad (32)$$

The complete update can then be shown as:

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} = \sum_i x_i p(x_i)(1 - p(x_i)) \quad (31)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} = \sum_i x_i^2 p(x_i)(1 - p(x_i)) \quad (32)$$

The complete update can then be shown as:

$$\begin{pmatrix} \mathbf{w}_{0,k+1} \\ \mathbf{w}_{1,k+1} \end{pmatrix} =$$

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} = \sum_i x_i p(x_i)(1 - p(x_i)) \quad (31)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} = \sum_i x_i^2 p(x_i)(1 - p(x_i)) \quad (32)$$

The complete update can then be shown as:

$$\begin{pmatrix} \mathbf{w}_{0,k+1} \\ \mathbf{w}_{1,k+1} \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} -$$

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} = \sum_i x_i p(x_i)(1 - p(x_i)) \quad (31)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} = \sum_i x_i^2 p(x_i)(1 - p(x_i)) \quad (32)$$

The complete update can then be shown as:

$$\begin{pmatrix} \mathbf{w}_{0,k+1} \\ \mathbf{w}_{1,k+1} \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} - \left[ \begin{pmatrix} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} \\ \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1 \partial b} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} \end{pmatrix}^{-1} \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} \right]_{\mathbf{w}_k} \quad (33)$$



# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} = \sum_i x_i p(x_i)(1 - p(x_i)) \quad (31)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} = \sum_i x_i^2 p(x_i)(1 - p(x_i)) \quad (32)$$

The complete update can then be shown as:

$$\begin{pmatrix} \mathbf{w}_{0,k+1} \\ \mathbf{w}_{1,k+1} \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} - \left[ \begin{pmatrix} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} \\ \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1 \partial b} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} \end{pmatrix}^{-1} \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} \right]_{\mathbf{w}_k} \quad (33)$$

Alternatively:

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} = \sum_i p(x_i)(1 - p(x_i)) \quad (30)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} = \sum_i x_i p(x_i)(1 - p(x_i)) \quad (31)$$

$$\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} = \sum_i x_i^2 p(x_i)(1 - p(x_i)) \quad (32)$$

The complete update can then be shown as:

$$\begin{pmatrix} \mathbf{w}_{0,k+1} \\ \mathbf{w}_{1,k+1} \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{0,k} \\ \mathbf{w}_{1,k} \end{pmatrix} - \left[ \begin{pmatrix} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b^2} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial b \partial w_1} \\ \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1 \partial b} & \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial w_1^2} \end{pmatrix}^{-1} \begin{pmatrix} \frac{\partial \text{NLL}}{\partial b} \\ \frac{\partial \text{NLL}}{\partial w_1} \end{pmatrix} \right]_{\mathbf{w}_k} \quad (33)$$

Alternatively:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \left[ \left( \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right)^{-1} \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} \right]_{\mathbf{w}_k} \quad (34)$$

# Compact matrix representation of derivatives

# Compact matrix representation of derivatives

Recall in 1D:

# Compact matrix representation of derivatives

Recall in 1D:  $\mathbf{w} = \begin{pmatrix} b \\ w_1 \end{pmatrix}$ .

# Compact matrix representation of derivatives

Recall in 1D:  $\mathbf{w} = \begin{pmatrix} b \\ w_1 \end{pmatrix}$ . If we denote:  $\mathbf{x}_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}$

# Compact matrix representation of derivatives

Recall in 1D:  $\mathbf{w} = \begin{pmatrix} b \\ w_1 \end{pmatrix}$ . If we denote:  $\mathbf{x}_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}$  Then we can write:

# Compact matrix representation of derivatives

Recall in 1D:  $\mathbf{w} = \begin{pmatrix} b \\ w_1 \end{pmatrix}$ . If we denote:  $\mathbf{x}_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}$  Then we can write:

$$\begin{aligned} \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \begin{pmatrix} -\sum_i [y_i - p(x_i)] \\ -\sum_i [x_i (y_i - p(x_i))] \end{pmatrix} = - \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}^T \begin{pmatrix} y_1 - p(x_1) \\ y_2 - p(x_2) \\ \vdots \\ y_n - p(x_n) \end{pmatrix} \\ &= -\mathbf{X}^T(\mathbf{y} - \mathbf{p}) \end{aligned} \quad (35)$$



# Compact matrix representation of derivatives

Recall in 1D:  $\mathbf{w} = \begin{pmatrix} b \\ w_1 \end{pmatrix}$ . If we denote:  $\mathbf{x}_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}$  Then we can write:

$$\begin{aligned} \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \begin{pmatrix} -\sum_i [y_i - p(x_i)] \\ -\sum_i [x_i (y_i - p(x_i))] \end{pmatrix} & \begin{pmatrix} y_1 - p(x_1) \\ y_2 - p(x_2) \\ \vdots \\ y_n - p(x_n) \end{pmatrix} & (35) \\ &= -\mathbf{X}^T (\mathbf{y} - \mathbf{p}) \end{aligned}$$

# Compact matrix representation of derivatives

Recall in 1D:  $\mathbf{w} = \begin{pmatrix} b \\ w_1 \end{pmatrix}$ . If we denote:  $\mathbf{x}_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}$  Then we can write:

$$\begin{aligned} \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \begin{pmatrix} -\sum_i [y_i - p(x_i)] \\ -\sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \\ &= -\mathbf{X}^T (\mathbf{y} - \mathbf{p}) \end{aligned} \tag{35}$$

# Compact matrix representation of derivatives

Recall in 1D:  $\mathbf{w} = \begin{pmatrix} b \\ w_1 \end{pmatrix}$ . If we denote:  $\mathbf{x}_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}$  Then we can write:

$$\begin{aligned} \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \begin{pmatrix} -\sum_i [y_i - p(x_i)] \\ -\sum_i [x_i (y_i - p(x_i))] \end{pmatrix} = - \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}^T \begin{pmatrix} y_1 - p(x_1) \\ y_2 - p(x_2) \\ \vdots \\ y_n - p(x_n) \end{pmatrix} \\ &= -\mathbf{X}^T (\mathbf{y} - \mathbf{p}) \end{aligned} \quad (35)$$

# Compact matrix representation of derivatives (cont.)

# Compact matrix representation of derivatives (cont.)

We can also decompose the Hessian as:

# Compact matrix representation of derivatives (cont.)

We can also decompose the Hessian as:

$$\begin{aligned}\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} &= \begin{pmatrix} \sum_i p(x_i)(1 - p(x_i)) & \sum_i x_i p(x_i)(1 - p(x_i)) \\ \sum_i x_i p(x_i)(1 - p(x_i)) & \sum_i x_i^2 p(x_i)(1 - p(x_i)) \end{pmatrix} \\ &= \mathbf{X}^\top \mathbf{S} \mathbf{X}\end{aligned}\tag{36}$$

# Compact matrix representation of derivatives (cont.)

We can also decompose the Hessian as:

$$\begin{aligned}\frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} &= \begin{pmatrix} \sum_i p(x_i)(1 - p(x_i)) & \sum_i x_i p(x_i)(1 - p(x_i)) \\ \sum_i x_i p(x_i)(1 - p(x_i)) & \sum_i x_i^2 p(x_i)(1 - p(x_i)) \end{pmatrix} \\ &= \mathbf{X}^\top \mathbf{S} \mathbf{X}\end{aligned}\tag{36}$$

where  $\mathbf{S}$  is a diagonal  $N \times N$  matrix:

# Compact matrix representation of derivatives (cont.)

We can also decompose the Hessian as:

$$\begin{aligned} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} &= \begin{pmatrix} \sum_i p(x_i)(1 - p(x_i)) & \sum_i x_i p(x_i)(1 - p(x_i)) \\ \sum_i x_i p(x_i)(1 - p(x_i)) & \sum_i x_i^2 p(x_i)(1 - p(x_i)) \end{pmatrix} \\ &= \mathbf{X}^\top \mathbf{S} \mathbf{X} \end{aligned} \quad (36)$$

where  $\mathbf{S}$  is a diagonal  $N \times N$  matrix:

$$\mathbf{S} = \begin{pmatrix} p(x_1)(1 - p(x_1)) & 0 & \dots & 0 \\ 0 & p(x_2)(1 - p(x_2)) & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & 0 & p(x_n)(1 - p(x_n)) \end{pmatrix} \quad (37)$$



# Compact matrix representation of derivatives (cont.)

We can also decompose the Hessian as:

$$\begin{aligned} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} &= \begin{pmatrix} \sum_i p(x_i)(1 - p(x_i)) & \sum_i x_i p(x_i)(1 - p(x_i)) \\ \sum_i x_i p(x_i)(1 - p(x_i)) & \sum_i x_i^2 p(x_i)(1 - p(x_i)) \end{pmatrix} \\ &= \mathbf{X}^\top \mathbf{S} \mathbf{X} \end{aligned} \quad (36)$$

where  $\mathbf{S}$  is a diagonal  $N \times N$  matrix:

$$\mathbf{S} = \begin{pmatrix} p(x_1)(1 - p(x_1)) & 0 & \dots & 0 \\ 0 & p(x_2)(1 - p(x_2)) & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & 0 & p(x_n)(1 - p(x_n)) \end{pmatrix} \quad (37)$$

and  $\mathbf{X}$  is defined as before:

# Compact matrix representation of derivatives (cont.)

We can also decompose the Hessian as:

$$\begin{aligned} \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} &= \begin{pmatrix} \sum_i p(x_i)(1-p(x_i)) & \sum_i x_i p(x_i)(1-p(x_i)) \\ \sum_i x_i p(x_i)(1-p(x_i)) & \sum_i x_i^2 p(x_i)(1-p(x_i)) \end{pmatrix} \\ &= \mathbf{X}^\top \mathbf{S} \mathbf{X} \end{aligned} \quad (36)$$

where  $\mathbf{S}$  is a diagonal  $N \times N$  matrix:

$$\mathbf{S} = \begin{pmatrix} p(x_1)(1-p(x_1)) & 0 & \cdots & 0 \\ 0 & p(x_2)(1-p(x_2)) & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & 0 & p(x_n)(1-p(x_n)) \end{pmatrix} \quad (37)$$

and  $\mathbf{X}$  is defined as before:

$$\mathbf{X}^\top = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix} \quad (38)$$

# Compact matrix representation (cont.)

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k \\ &= \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{X}) \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}\end{aligned}\tag{39}$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})\tag{40}$$

In this form, the estimation is also called iteratively reweighted least squares (IRLS).

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{X}) \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}\end{aligned}\tag{39}$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})\tag{40}$$

In this form, the estimation is also called iteratively reweighted least squares (IRLS).

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}\end{aligned}\tag{39}$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})\tag{40}$$

In this form, the estimation is also called iteratively reweighted least squares (IRLS).

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}\end{aligned}\tag{39}$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})\tag{40}$$

In this form, the estimation is also called iteratively reweighted least squares (IRLS).



# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k$$
(39)

$$= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})$$
(40)

In this form, the estimation is also called iteratively reweighted least squares (IRLS).

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k \\ &= \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{X}) \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}\end{aligned}\tag{39}$$

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\begin{aligned} \mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k \\ &= \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{X}) \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \end{aligned} \tag{39}$$

where the adjusted response  $\mathbf{z}$  is given as:

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\begin{aligned} \mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k \\ &= \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{X}) \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \end{aligned} \tag{39}$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p}) \tag{40}$$

# Compact matrix representation (cont.)

Putting the previous results together, we can express the update step as:

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{H}^{-1} \mathbf{g}_k \\ &= \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{X}) \mathbf{w}_k + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}\end{aligned}\tag{39}$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})\tag{40}$$

In this form, the estimation is also called iteratively reweighted least squares (IRLS).

# OLS, WLS and IRLS

# OLS, WLS and IRLS

- Recall the OLS estimate:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (41)$$

# OLS, WLS and IRLS

- Recall the OLS estimate:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (41)$$

if  $\mathbf{y}$  is the response.



# OLS, WLS and IRLS

- Recall the OLS estimate:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (41)$$

if  $\mathbf{y}$  is the response.

- Also recall that the weighted least squares (WLS) is given by:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y} \quad (42)$$

# OLS, WLS and IRLS

- Recall the OLS estimate:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (41)$$

if  $\mathbf{y}$  is the response.

- Also recall that the weighted least squares (WLS) is given by:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y} \quad (42)$$

- In logistic regression, the coefficients can be found via the Newton-Raphson update, which can be specified as:

# OLS, WLS and IRLS

- Recall the OLS estimate:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (41)$$

if  $\mathbf{y}$  is the response.

- Also recall that the weighted least squares (WLS) is given by:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y} \quad (42)$$

- In logistic regression, the coefficients can be found via the Newton-Raphson update, which can be specified as:

$$\mathbf{w}_{k+1} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \quad (43)$$

# OLS, WLS and IRLS

- Recall the OLS estimate:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (41)$$

if  $\mathbf{y}$  is the response.

- Also recall that the weighted least squares (WLS) is given by:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y} \quad (42)$$

- In logistic regression, the coefficients can be found via the Newton-Raphson update, which can be specified as:

$$\mathbf{w}_{k+1} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \quad (43)$$

where  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p}) \quad (44)$$

# OLS, WLS and IRLS

- Recall the OLS estimate:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (41)$$

if  $\mathbf{y}$  is the response.

- Also recall that the weighted least squares (WLS) is given by:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y} \quad (42)$$

- In logistic regression, the coefficients can be found via the Newton-Raphson update, which can be specified as:

$$\mathbf{w}_{k+1} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \quad (43)$$

where  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p}) \quad (44)$$

- Note that the update step is identical in form to the WLS estimator
- However,  $\mathbf{W}$  and  $\mathbf{z}$  change in each iteration, hence the name iteratively reweighted least squares (IRLS)

# Summary

- The binary logistic regression model is given by:

$$p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} \quad (45)$$

- The negative log-likelihood of a sample of  $N$  observations in the binomial response case is:

$$\text{NLL}(\mathbf{w}) = \sum_i [y_i (b + w_1 x_i) - \log(1 + e^{b + w_1 x_i})] \quad (46)$$

- Based on the principle of maximum likelihood, the estimate  $\hat{\mathbf{w}}$  is given by the minimizing NLL.
- This can be solved via gradient descent or Newton-Raphson.

# Summary (cont.)

## Gradient descent update for logistic regression

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \lambda \nabla \text{NLL}(\mathbf{w}_k) \quad (47)$$

## Newton-Raphson update for logistic regression

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}_{\mathbf{w}_k}^{-1}(\text{NLL}) \nabla_{\mathbf{w}_k} \text{NLL}(\mathbf{w}_k) \quad (48)$$

This can be rewritten as:

$$\mathbf{w}_{k+1} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \quad (49)$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p}) \quad (50)$$

# Summary (cont.)

## Gradient descent update for logistic regression

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \lambda \nabla \text{NLL}(\mathbf{w}_k) \quad (47)$$

## Newton-Raphson update for logistic regression

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}_{\mathbf{w}_k}^{-1}(\text{NLL}) \nabla_{\mathbf{w}_k} \text{NLL}(\mathbf{w}_k) \quad (48)$$

This can be rewritten as:

$$\mathbf{w}_{k+1} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \quad (49)$$

where the adjusted response  $\mathbf{z}$  is given as:

$$\mathbf{z} = \mathbf{X} \mathbf{w}_k + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p}) \quad (50)$$

In this form, the estimation is also called iteratively reweighted least squares (IRLS).



# Other considerations

# Other considerations

- MAP estimation: weight decay/regularization to make NLL convex (have unique solution). We define the **penalized negative log-likelihood** PNLL as:

$$\text{PNLL}(\mathbf{w}) = \text{NLL}(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w} \quad (51)$$

# Other considerations

- MAP estimation: weight decay/regularization to make NLL convex (have unique solution). We define the **penalized negative log-likelihood** PNLL as:

$$\text{PNLL}(\mathbf{w}) = \text{NLL}(\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w} \quad (51)$$

where  $\lambda$  is the decay parameter.

- Thus:  $\nabla_{\mathbf{w}} \text{PNLL}(\mathbf{w}) = \mathbf{g}(\mathbf{w}) + 2\lambda \mathbf{w}$

# Other considerations

- MAP estimation: weight decay/regularization to make NLL convex (have unique solution). We define the **penalized negative log-likelihood** PNLL as:

$$\text{PNLL}(\mathbf{w}) = \text{NLL}(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w} \quad (51)$$

where  $\lambda$  is the decay parameter.

- Thus:  $\nabla_{\mathbf{w}} \text{PNLL}(\mathbf{w}) = \mathbf{g}(\mathbf{w}) + 2\lambda \mathbf{w}$
- And:  $\nabla_{\mathbf{w}}^2 \text{PNLL}(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + 2\lambda \mathbf{I}$

# Other considerations

- MAP estimation: weight decay/regularization to make NLL convex (have unique solution). We define the **penalized negative log-likelihood** PNLL as:

$$\text{PNLL}(\mathbf{w}) = \text{NLL}(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w} \quad (51)$$

where  $\lambda$  is the decay parameter.

- Thus:  $\nabla_{\mathbf{w}} \text{PNLL}(\mathbf{w}) = \mathbf{g}(\mathbf{w}) + 2\lambda \mathbf{w}$
- And:  $\nabla_{\mathbf{w}}^2 \text{PNLL}(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + 2\lambda \mathbf{I}$

# Reading assignments

- **PMLCE 9.2**
- **PMLI 10.1-3**
- **ESL 4.4**