

Statistical Learning

Master Spécialisé Intelligence Artificielle de Confiance (IAC)
@ Centrale Supélec en partenariat avec l'IRT SystemX
2024/2025.

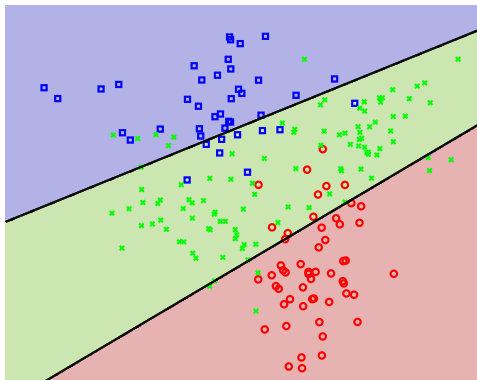
FAÏCEL CHAMROUKHI



 chamroukhi.com

- 1 Logistic Regression
 - Iteratively Reweighted Least Squares (IRLS)
- 2 Multi-class logistic regression
 - IRLS for Multi-class logistic regression

Multi-class Logistic Regression



- The data are represented by a random pair $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ where \mathbf{X} is a vector of descriptors for some variable of interest Y
 - The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(\mathbf{x})$ is a good approximation of the true output y
 - In a **classification** problem : typically $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(X) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|X, h)$ can be computed in terms of $P_\theta(Y - h(X))$.

- Data : a random sample $(\mathbf{X}_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to “minimize” the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(X)$.
 - ↪ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$
 - ↪ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}
- Data-Scientist's “**Toolbox**” : {Data, loss, hypothesis, algorithm}

- The data are represented by a random pair $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ where \mathbf{X} is a vector of descriptors for some variable of interest Y
 - The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(\mathbf{x})$ is a good approximation of the true output y
 - In a **classification** problem : typically $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)
- ⇨ We will mainly focus on parametric probabilistic models of the form

$$Y = h(X) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|X, h)$ can be computed in terms of $P_\theta(Y - h(X))$.

- Data : a random sample $(\mathbf{X}_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to “minimize” the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(X)$.
 - ⇨ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$
 - ⇨ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}
- Data-Scientist's “**Toolbox**” : {Data, loss, hypothesis, algorithm}

- The data are represented by a random pair $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ where \mathbf{X} is a vector of descriptors for some variable of interest Y
- The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(\mathbf{x})$ is a good approximation of the true output y
- In a **classification** problem : typically $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)

⇨ We will mainly focus on parametric probabilistic models of the form

$$Y = h(X) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|X, h)$ can be computed in terms of $P_\theta(Y - h(X))$.

- Data : a random sample $(\mathbf{X}_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to “minimize” the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(X)$.
 - ⇨ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$
 - ⇨ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}
- Data-Scientist's “**Toolbox**” : {Data, loss, hypothesis, algorithm}

- The data are represented by a random pair $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ where \mathbf{X} is a vector of descriptors for some variable of interest Y
 - The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(\mathbf{x})$ is a good approximation of the true output y
 - In a **classification** problem : typically $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|\mathbf{X}, h)$ can be computed in terms of $P_\theta(Y - h(\mathbf{X}))$.

- Data : a random sample $(\mathbf{X}_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- Data-Scientist's role : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to "minimize" the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(\mathbf{X})$.
 - ↪ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$
 - ↪ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}
- Data-Scientist's "**Toolbox**" : {Data, loss, hypothesis, algorithm}

- The data are represented by a random pair $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ where \mathbf{X} is a vector of descriptors for some variable of interest Y
 - The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(\mathbf{x})$ is a good approximation of the true output y
 - In a **classification** problem : typically $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)
- ↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|\mathbf{X}, h)$ can be computed in terms of $P_\theta(Y - h(\mathbf{X}))$.

- **Data** : a random sample $(\mathbf{X}_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to "minimize" the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(\mathbf{X})$.
 - ↪ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$
 - ↪ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}
- Data-Scientist's "**Toolbox**" : {Data, loss, hypothesis, algorithm}

- The data are represented by a random pair $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ where X is a vector of descriptors for some variable of interest Y
 - The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(x)$ is a good approximation of the true output y
 - In a **classification** problem : typically $X \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)
- ↔ We will mainly focus on parametric probabilistic models of the form

$$Y = h(X) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|X, h)$ can be computed in terms of $P_\theta(Y - h(X))$.

- Data : a random sample $(X_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (x_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to “minimize” the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(X)$.
 - ↔ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$
 - ↔ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}
- Data-Scientist's “**Toolbox**” : {Data, loss, hypothesis, algorithm}

- The data are represented by a random pair $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ where \mathbf{X} is a vector of descriptors for some variable of interest Y
- The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(\mathbf{x})$ is a good approximation of the true output y
- In a **classification** problem : typically $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)

↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|\mathbf{X}, h)$ can be computed in terms of $P_\theta(Y - h(\mathbf{X}))$.

- Data : a random sample $(\mathbf{X}_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to “minimize” the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(\mathbf{X})$.

↪ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$

↪ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}

- Data-Scientist's “**Toolbox**” : {Data, loss, hypothesis, algorithm}

- The data are represented by a random pair $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ where \mathbf{X} is a vector of descriptors for some variable of interest Y
- The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(\mathbf{x})$ is a good approximation of the true output y
- In a **classification** problem : typically $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)

↪ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|\mathbf{X}, h)$ can be computed in terms of $P_\theta(Y - h(\mathbf{X}))$.

- Data : a random sample $(\mathbf{X}_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to “minimize” the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(\mathbf{X})$.
 - ↪ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$
 - ↪ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}
- Data-Scientist's “**Toolbox**” : {Data, loss, hypothesis, algorithm}

- The data are represented by a random pair $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ where \mathbf{X} is a vector of descriptors for some variable of interest Y
 - The objective is **Prediction**, i.e. to seek for a prediction function $h : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\hat{y} = h(\mathbf{x})$ is a good approximation of the true output y
 - In a **classification** problem : typically $\mathbf{X} \in \mathcal{X} \subset \mathbb{R}^p$ and $Y \in \mathcal{Y} = \{0, 1\}, \{-1, +1\}$ (binary classification) or $\{1, \dots, K\}$ (multiclass classification)
- ⇨ We will mainly focus on parametric probabilistic models of the form

$$Y = h(\mathbf{X}) + \epsilon, \epsilon \sim p_\theta$$

with the conditional distr. $P(Y|\mathbf{X}, h)$ can be computed in terms of $P_\theta(Y - h(\mathbf{X}))$.

- Data : a random sample $(\mathbf{X}_i, Y_i)_{i=1}^n$ with observed values $\mathcal{D}_n = (\mathbf{x}_i, y_i)_{i=1}^n$
- **Data-Scientist's role** : given the **data**, choose a **prediction function** h from a class \mathcal{H} that attempts to “minimize” the prediction error for of all possible data (**risk**) $R(h)$, under a **loss** function ℓ measuring the error of predicting Y by $h(\mathbf{X})$.
 - ⇨ minimize the **empirical risk** (data- \mathcal{D}_n -driven) $R_n(h)$
 - ⇨ Minimizing $R_n(h)$ may require an optimization **algorithm** \mathcal{A}
- **Data-Scientist's “Toolbox”** : {Data, loss, hypothesis, algorithm}

Def. Classifier or classification rule

$$h: \mathcal{X} \rightarrow \mathcal{Y}$$
$$x \mapsto h(x)$$

is a decision/prediction function, parametric or not, linear or not, ...

Example : Linear predictors

$$h: \mathbb{R}^p \rightarrow \mathbb{R}$$
$$x \mapsto \langle x, \theta \rangle = \theta^T x$$

The **predicted** values of Y_i 's for new covariates $X_i = x_i$ s correspond to

$$\hat{y}_i = h(x_i)$$

Example : Linear predictors (cont.) : $\hat{y}_i = \langle x_i, \theta \rangle = \theta^T x_i$

Q : How good we are in prediction on a particular pair (x, y) ?

Def. Classifier or classification rule

$$h: \mathcal{X} \rightarrow \mathcal{Y}$$
$$x \mapsto h(x)$$

is a decision/prediction function, parametric or not, linear or not, ...

Example : Linear predictors

$$h: \mathbb{R}^p \rightarrow \mathbb{R}$$
$$x \mapsto \langle x, \theta \rangle = \theta^T x$$

The **predicted** values of Y_i 's for new covariates $X_i = x_i$ s correspond to

$$\hat{y}_i = h(x_i)$$

Example : Linear predictors (cont.) : $\hat{y}_i = \langle x_i, \theta \rangle = \theta^T x_i$

Q : How good we are in prediction on a particular pair (x, y) ?

Def. Classifier or classification rule

$$h: \mathcal{X} \rightarrow \mathcal{Y}$$
$$x \mapsto h(x)$$

is a decision/prediction function, parametric or not, linear or not, ...

Example : Linear predictors

$$h: \mathbb{R}^p \rightarrow \mathbb{R}$$
$$x \mapsto \langle x, \theta \rangle = \theta^T x$$

The **predicted** values of Y_i 's for new covariates $X_i = x_i$ s correspond to

$$\hat{y}_i = h(x_i)$$

Example : Linear predictors (cont.) : $\hat{y}_i = \langle x_i, \theta \rangle = \theta^T x_i$

Q : How good we are in prediction on a particular pair (x, y) ?

Def. Classifier or classification rule

$$h: \mathcal{X} \rightarrow \mathcal{Y}$$
$$x \mapsto h(x)$$

is a decision/prediction function, parametric or not, linear or not, ...

Example : Linear predictors

$$h: \mathbb{R}^p \rightarrow \mathbb{R}$$
$$x \mapsto \langle x, \theta \rangle = \theta^T x$$

The **predicted** values of Y_i 's for new covariates $X_i = x_i$ s correspond to

$$\hat{y}_i = h(x_i)$$

Example : Linear predictors (cont.) : $\hat{y}_i = \langle x_i, \theta \rangle = \theta^T x_i$

Q : How good we are in prediction on a particular pair (x, y) ?

Def. Loss function

$$\begin{aligned} \ell: \mathcal{Y} \times \mathcal{Y} &\rightarrow \mathbb{R} \\ (y, h(x)) &\mapsto \ell(y, h(x)) \end{aligned}$$

It measures how good we are on a particular (x, y) pair.

Examples of loss functions in classification

- “0-1” loss : $\ell(y, h(x)) = \mathbb{1}_{h(x) \neq y}$
- logarithmic loss : $\ell(y, h_\theta(x)) = -\log(p_\theta(x, y))$
Denoting $\ell(y, h(x)) = \phi(yh(x))$
- Hinge loss $\phi_{\text{hinge}}(u) = (1 - u)_+$
- Logistic loss $\phi_{\text{logistic}}(u) = \log(1 + \exp(-u))$
- Exponential loss $\phi_{\text{exp}}(u) = \exp(-u)$

Def. Loss function

$$\begin{aligned}\ell: \mathcal{Y} \times \mathcal{Y} &\rightarrow \mathbb{R} \\ (y, h(x)) &\mapsto \ell(y, h(x))\end{aligned}$$

It measures how good we are on a particular (x, y) pair.

Examples of loss functions in classification

- “0-1” loss : $\ell(y, h(x)) = \mathbb{1}_{h(x) \neq y}$
- logarithmic loss : $\ell(y, h_\theta(x)) = -\log(p_\theta(x, y))$
Denoting $\ell(y, h(x)) = \phi(yh(x))$
- Hinge loss $\phi_{\text{hinge}}(u) = (1 - u)_+$
- Logistic loss $\phi_{\text{logistic}}(u) = \log(1 + \exp(-u))$
- Exponential loss $\phi_{\text{exp}}(u) = \exp(-u)$

Def. Loss function

$$\begin{aligned} \ell: \mathcal{Y} \times \mathcal{Y} &\rightarrow \mathbb{R} \\ (y, h(x)) &\mapsto \ell(y, h(x)) \end{aligned}$$

It measures how good we are on a particular (x, y) pair.

Examples of loss functions in classification

- “0-1” loss : $\ell(y, h(x)) = \mathbb{1}_{h(x) \neq y}$
- logarithmic loss : $\ell(y, h_\theta(x)) = -\log(p_\theta(x, y))$
Denoting $\ell(y, h(x)) = \phi(yh(x))$
- Hinge loss $\phi_{\text{hinge}}(u) = (1 - u)_+$
- Logistic loss $\phi_{\text{logistic}}(u) = \log(1 + \exp(-u))$
- Exponential loss $\phi_{\text{exp}}(u) = \exp(-u)$

Def. Loss function

$$\begin{aligned} \ell: \mathcal{Y} \times \mathcal{Y} &\rightarrow \mathbb{R} \\ (y, h(x)) &\mapsto \ell(y, h(x)) \end{aligned}$$

It measures how good we are on a particular (x, y) pair.

Examples of loss functions in classification

- “0-1” loss : $\ell(y, h(x)) = \mathbb{1}_{h(x) \neq y}$
- logarithmic loss : $\ell(y, h_\theta(x)) = -\log(p_\theta(x, y))$
Denoting $\ell(y, h(x)) = \phi(yh(x))$
- Hinge loss $\phi_{\text{hinge}}(u) = (1 - u)_+$
- Logistic loss $\phi_{\text{logistic}}(u) = \log(1 + \exp(-u))$
- Exponential loss $\phi_{\text{exp}}(u) = \exp(-u)$

Examples of loss functions in classification

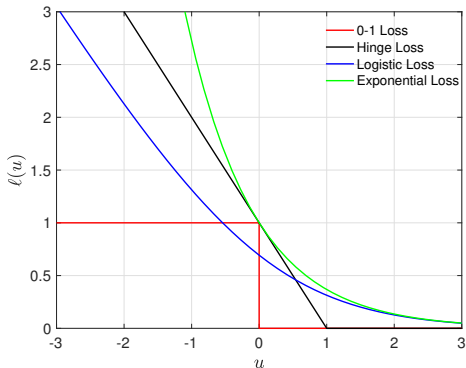


FIGURE – Some loss functions in classification : curves of $\ell(u)$ for $u = yh(x)$; $y \in \{-1, 1\}$.
[plot_losses_classification.m]

For $y \in \{-1, 1\}$, with $u = yh(x)$:

- “0-1” loss : $\ell(u) = \mathbb{1}_{\text{sign}(u) \neq 1}$
- Hinge loss $\ell_{\text{hinge}}(u) = (1 - u)_+$
- Logistic loss $\ell_{\text{logistic}}(u) = \log(1 + \exp(-u))$
- Exponential loss $\ell_{\text{exp}}(u) = \exp(-u)$

- **Risk** : the *Expected loss* :

$$R(h) = \mathbb{E}_P[\ell(Y, h(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(x)) dP(x, y)$$

- ↪ the error of approximating Y by model/hypothesis $h(X)$ as measured by a chosen loss function $\ell(Y, h(X))$ given the pair (X, Y) with (unknown) joint distribution P ,
- ↪ prediction error : measures the generalization performance of the function h .
- **"0-1" Risk** : Under the "0-1"-loss $\ell(y, h(x)) = \mathbb{1}_{h(x) \neq y}$:

$$R(h) = \mathbb{E}_P[\mathbb{1}_{h(x) \neq y}] = \mathbb{P}(h(X) \neq Y). = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{1}_{h(x) \neq y} dP(x, y)$$

↪ This is the most used risk in classification

Q : what is the best function h ? or equivalently, when the risk $R(h)$ is optimal?

- Risk : the *Expected loss* :

$$R(h) = \mathbb{E}_P[\ell(Y, h(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(x)) dP(x, y)$$

- ↔ the error of approximating Y by model/hypothesis $h(X)$ as measured by a chosen loss function $\ell(Y, h(X))$ given the pair (X, Y) with (unknown) joint distribution P ,
- ↔ prediction error : measures the generalization performance of the function h .

- **“0-1” Risk** : Under the **“0-1”-loss** $\ell(y, h(x)) = \mathbb{1}_{h(x) \neq y}$:

$$R(h) = \mathbb{E}_P[\mathbb{1}_{h(X) \neq Y}] = \mathbb{P}(h(X) \neq Y). = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{1}_{h(x) \neq y} dP(x, y)$$

↔ This is the most used risk in classification

Q : what is the best function h ? or equivalently, when the risk $R(h)$ is optimal ?

- Risk : the *Expected loss* :

$$R(h) = \mathbb{E}_P[\ell(Y, h(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(x)) dP(x, y)$$

- ↔ the error of approximating Y by model/hypothesis $h(X)$ as measured by a chosen loss function $\ell(Y, h(X))$ given the pair (X, Y) with (unknown) joint distribution P ,
- ↔ prediction error : measures the generalization performance of the function h .

- **"0-1" Risk** : Under the "0-1"-loss $\ell(y, h(x)) = \mathbb{1}_{h(x) \neq y}$:

$$R(h) = \mathbb{E}_P[\mathbb{1}_{h(x) \neq y}] = \mathbb{P}(h(X) \neq Y). = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{1}_{h(x) \neq y} dP(x, y)$$

- ↔ This is the most used risk in classification

Q : what is the best function h ? or equivalently, when the risk $R(h)$ is optimal ?

Theorem (The Bayes classifier)

Under the (0-1)-loss, $\ell(Y, h(X)) = \mathbb{1}_{h(X) \neq Y}$, the classification function $h^*(x)$ minimizing the risk (the Bayes classifier)

$$R(h) = \mathbb{P}(Y \neq h(X)) = \int_{\mathcal{X}} \mathbb{P}(Y \neq h(X) | X = x) dP_X(x)$$

is given by

$$\forall x \in \mathcal{X}, \quad h^*(x) = \arg \max_{k \in \mathcal{Y}} \mathbb{P}(Y = k | X = x).$$

Def. Decision boundaries

The decision boundary between each pair of classes k and ℓ , $(k, \ell) \in \mathcal{Y} \times \mathcal{Y}$ is defined by

$$\eta_{k,\ell}(x) = \{x : \mathbb{P}(Y = k | X = x) = \mathbb{P}(Y = \ell | X = x)\}$$

Proof. Optimal classifier.

Given $X = x$, the conditional risk under the 0-1 loss is

$$\begin{aligned}r(h|X = x) &= \mathbb{E}_{Y|X=x}[\ell(Y, h(X))|X = x] = \mathbb{E}_{Y|X=x}[\mathbb{1}_{Y \neq h(X)}|X = x] \\ &= \mathbb{P}[Y \neq h(X)|X = x] \\ &= 1 - \mathbb{P}[Y = h(X)|X = x].\end{aligned}$$

By noting that

$$\begin{aligned}\min_{k \in \mathcal{Y}} r(h|X = x) &= -1 + \max_{k \in \mathcal{Y}} \mathbb{P}(Y = k|X = x) \\ \arg \min_{k \in \mathcal{Y}} r(h|X = x) &= \arg \max_{k \in \mathcal{Y}} \mathbb{P}(Y = k|X = x)\end{aligned}$$

we see that $h^*(x) = \arg \max_{k \in \mathcal{Y}} \mathbb{P}(Y = k|X = x)$ achieves the minimized risk $r(h|X = x)$.

Then the risk $R(h^*) = \mathbb{E}_X[-1 + \max_{k \in \mathcal{Y}} \mathbb{P}(Y = k|X = x)]$ is Bayes. \square

Goal : estimate h^* , knowing only the data sample $D_n = (X_i, Y_i)_{i=1}^n$ and loss ℓ .

- Then *Expected loss* $R(h)$ depends on the joint distribution P of the pair (X, Y) .
In real situations P is unknown, as we only have a sample $D_n = (X_i, Y_i)_{1 \leq i \leq n}$,

↪ We attempt to minimize the **Empirical Risk**

$$R_n(h) = \frac{1}{n} \sum_{i=1}^n l(Y_i, h(X_i))$$

to estimate h^* (within a family \mathcal{H})

↪ ERM : $\hat{h}_n = \arg \min_{h \in \mathcal{H}} R_n(h)$ is the **ERM** of h

- 0-1 Risk : Under the 0-1 loss (standard in classification) : $\ell_{0-1}(y, h(x)) = \mathbb{1}_{y \neq h(x)}$,
the empirical 0-1 risk is

$$R_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \neq h(X_i)}$$

ERM and MLE : Conditional maximum likelihood risks :

- MLE (density estimation framework) : We seek for an estimator of the parameters θ of the joint distribution $p_\theta(x, y)$.

- In discriminative learning (eg. logistic regression), we are interested in estimating the conditional distribution $P(Y|X)$, rather than the joint distribution $P(X, Y)$.
- Consider the log-loss : $\ell(y, h_\theta(x)) = -\log(p_\theta(y|x))$. We therefore have the conditional log-likelihood risks

$$R(\theta) = -\mathbb{E}[\log p_\theta(Y|X)] \text{ and } R_n(\theta) = -\frac{1}{n} \sum_{i=1}^n \log p_\theta(y_i|x_i).$$

- For an i.i.d sample $\{(x_i, y_i)_{i=1}^n\}$, the conditional log-likelihood function of θ is :

$$\log L(\theta) = \sum_{i=1}^n \log p_\theta(y_i|x_i)$$

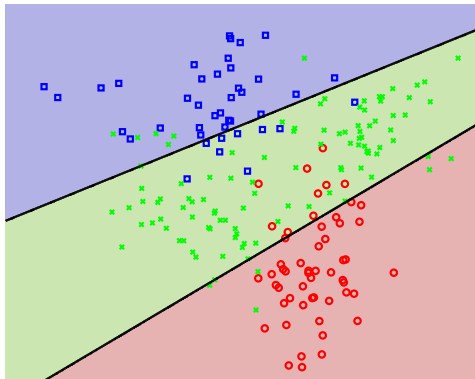
Then

$$R_n(\theta) = -\frac{1}{n} \sum_{i=1}^n \log p_\theta(y_i|x_i) = -\frac{1}{n} \log L(\theta)$$

↪ With this log-loss, ERM coincides with conditional MLE.

- Liner classifier : Consider $\mathcal{H} = \{h_\theta(x) = \alpha + \beta^T x\}$, the set of linear functions in x

Multi-class Logistic Regression



- We model the random pair (\mathbf{X}, Y) where $\mathbf{X}_i \in \mathcal{X} \subset \mathbb{R}^d$ is the predictor and the response $Y \in \mathcal{Y} = \{0, 1\}$ is the class label of \mathbf{X}
- Logistic Regression : Probabilistic Discriminative approach to model $\mathbb{P}(Y|\mathbf{X})$ as

$$\mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}) = \text{Logistic}(\mathbf{x}^T \boldsymbol{\theta}) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})} .$$

- $Y|\mathbf{X} = \mathbf{x}$ is Bernoulli with probability of success $\pi_{\boldsymbol{\theta}}(\mathbf{x})$, i.e.

$$\forall y \in \{0, 1\}, \mathbb{P}_{\boldsymbol{\theta}}(Y = y|\mathbf{X} = \mathbf{x}) = \pi_{\boldsymbol{\theta}}(\mathbf{x})^y (1 - \pi_{\boldsymbol{\theta}}(\mathbf{x}))^{1-y}$$

where $\pi(\mathbf{x}; \boldsymbol{\theta}) = \mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}; \boldsymbol{\theta})$ is the sigmoid function.

- Classification rule : We have $h(x)$ is defined as

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}) = \text{Logistic}(\mathbf{x}^T \boldsymbol{\theta}) > \frac{1}{2}, \\ 0, & \text{otherwise.} \end{cases} \quad \text{Eq. : } h_{\boldsymbol{\theta}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \boldsymbol{\theta}^T \mathbf{x} > 0, \\ 0, & \text{otherwise.} \end{cases}$$

- The latter comes from the linear boundary $\{x : \log \frac{\mathbb{P}(Y=1|\mathbf{X}=\mathbf{x})}{\mathbb{P}(Y=0|\mathbf{X}=\mathbf{x})} = \beta_0 + \boldsymbol{\beta}^T \mathbf{x} = 0\}$
- The parameter vector of the model $\boldsymbol{\theta} = (\beta_0, \boldsymbol{\beta}^T)^T \in \mathbb{R}^{d+1}$
- **Q** : Fit $\boldsymbol{\theta}$ from the training data.

- We model the random pair (\mathbf{X}, Y) where $\mathbf{X}_i \in \mathcal{X} \subset \mathbb{R}^d$ is the predictor and the response $Y \in \mathcal{Y} = \{0, 1\}$ is the class label of \mathbf{X}
- Logistic Regression : Probabilistic Discriminative approach to model $\mathbb{P}(Y|\mathbf{X})$ as

$$\mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}) = \text{Logistic}(x^T \theta) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})} .$$

- $Y|X = x$ is Bernoulli with probability of success $\pi_\theta(x)$, i.e.

$$\forall y \in \{0, 1\}, \mathbb{P}_\theta(Y = y|X = x) = \pi_\theta(x)^y (1 - \pi_\theta(x))^{1-y}$$

where $\pi(x; \theta) = \mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}; \theta)$ is the sigmoid function.

- Classification rule : We have $h(x)$ is defined as

$$h_\theta(x) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X = x) = \text{Logistic}(x^T \theta) > \frac{1}{2}, \\ 0, & \text{otherwise.} \end{cases} \quad \text{Eq. : } h_\theta(x) = \begin{cases} 1 & \text{if } \theta^T x > 0, \\ 0, & \text{otherwise.} \end{cases}$$

- The latter comes from the linear boundary $\{x : \log \frac{\mathbb{P}(Y=1|\mathbf{X}=\mathbf{x})}{\mathbb{P}(Y=0|\mathbf{X}=\mathbf{x})} = \beta_0 + \boldsymbol{\beta}^T \mathbf{x} = 0\}$
- The parameter vector of the model $\theta = (\beta_0, \boldsymbol{\beta}^T)^T \in \mathbb{R}^{d+1}$
- Q : Fit θ from the training data.

- We model the random pair (\mathbf{X}, Y) where $\mathbf{X}_i \in \mathcal{X} \subset \mathbb{R}^d$ is the predictor and the response $Y \in \mathcal{Y} = \{0, 1\}$ is the class label of \mathbf{X}
- Logistic Regression : Probabilistic Discriminative approach to model $\mathbb{P}(Y|\mathbf{X})$ as

$$\mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}) = \text{Logistic}(x^T \theta) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})} .$$

- $Y|X = x$ is Bernoulli with probability of success $\pi_{\theta}(\mathbf{x})$, i.e.

$$\forall y \in \{0, 1\}, \mathbb{P}_{\theta}(Y = y|X = x) = \pi_{\theta}(x)^y (1 - \pi_{\theta}(x))^{1-y}$$

where $\pi(\mathbf{x}; \theta) = \mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}; \theta)$ is the sigmoid function.

- Classification rule : We have $h(x)$ is defined as

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X = x) = \text{Logistic}(x^T \theta) > \frac{1}{2}, \\ 0, & \text{otherwise.} \end{cases} \quad \text{Eq. : } h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x > 0, \\ 0, & \text{otherwise.} \end{cases}$$

- The latter comes from the linear boundary $\{x : \log \frac{\mathbb{P}(Y=1|\mathbf{X}=\mathbf{x})}{\mathbb{P}(Y=0|\mathbf{X}=\mathbf{x})} = \beta_0 + \boldsymbol{\beta}^T \mathbf{x} = 0\}$
- The parameter vector of the model $\theta = (\beta_0, \boldsymbol{\beta}^T)^T \in \mathbb{R}^{d+1}$
- Q : Fit θ from the training data.

- We model the random pair (\mathbf{X}, Y) where $\mathbf{X}_i \in \mathcal{X} \subset \mathbb{R}^d$ is the predictor and the response $Y \in \mathcal{Y} = \{0, 1\}$ is the class label of \mathbf{X}
- Logistic Regression : Probabilistic Discriminative approach to model $\mathbb{P}(Y|\mathbf{X})$ as

$$\mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}) = \text{Logistic}(x^T \theta) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})} .$$

- $Y|X = x$ is Bernoulli with probability of success $\pi_{\theta}(\mathbf{x})$, i.e.

$$\forall y \in \{0, 1\}, \mathbb{P}_{\theta}(Y = y|X = x) = \pi_{\theta}(x)^y (1 - \pi_{\theta}(x))^{1-y}$$

where $\pi(\mathbf{x}; \theta) = \mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}; \theta)$ is the sigmoid function.

- Classification rule : We have $h(x)$ is defined as

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X = x) = \text{Logistic}(x^T \theta) > \frac{1}{2}, \\ 0, & \text{otherwise.} \end{cases} \quad \text{Eq. : } h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x > 0, \\ 0, & \text{otherwise.} \end{cases}$$

- The latter comes from the linear boundary $\{x : \log \frac{\mathbb{P}(Y=1|\mathbf{X}=\mathbf{x})}{\mathbb{P}(Y=0|\mathbf{X}=\mathbf{x})} = \beta_0 + \boldsymbol{\beta}^T \mathbf{x} = 0\}$
- The parameter vector of the model $\theta = (\beta_0, \boldsymbol{\beta}^T)^T \in \mathbb{R}^{d+1}$
- Q : Fit θ from the training data.

- We model the random pair (\mathbf{X}, Y) where $\mathbf{X}_i \in \mathcal{X} \subset \mathbb{R}^d$ is the predictor and the response $Y \in \mathcal{Y} = \{0, 1\}$ is the class label of \mathbf{X}
- Logistic Regression : Probabilistic Discriminative approach to model $\mathbb{P}(Y|\mathbf{X})$ as

$$\mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}) = \text{Logistic}(\mathbf{x}^T \boldsymbol{\theta}) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})} .$$

- $Y|X = x$ is Bernoulli with probability of success $\pi_{\boldsymbol{\theta}}(\mathbf{x})$, i.e.

$$\forall y \in \{0, 1\}, \mathbb{P}_{\boldsymbol{\theta}}(Y = y|X = x) = \pi_{\boldsymbol{\theta}}(\mathbf{x})^y (1 - \pi_{\boldsymbol{\theta}}(\mathbf{x}))^{1-y}$$

where $\pi(\mathbf{x}; \boldsymbol{\theta}) = \mathbb{P}(Y = 1|\mathbf{X} = \mathbf{x}; \boldsymbol{\theta})$ is the sigmoid function.

- Classification rule : We have $h(x)$ is defined as

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X = x) = \text{Logistic}(\mathbf{x}^T \boldsymbol{\theta}) > \frac{1}{2}, \\ 0, & \text{otherwise.} \end{cases} \quad \text{Eq. : } h_{\boldsymbol{\theta}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \boldsymbol{\theta}^T \mathbf{x} > 0, \\ 0, & \text{otherwise.} \end{cases}$$

- The latter comes from the linear boundary $\{x : \log \frac{\mathbb{P}(Y=1|\mathbf{X}=\mathbf{x})}{\mathbb{P}(Y=0|\mathbf{X}=\mathbf{x})} = \beta_0 + \boldsymbol{\beta}^T \mathbf{x} = 0\}$
- The parameter vector of the model $\boldsymbol{\theta} = (\beta_0, \boldsymbol{\beta}^T)^T \in \mathbb{R}^{d+1}$
- **Q** : Fit $\boldsymbol{\theta}$ from the training data.

Linear decision boundary :

$$\begin{aligned}\eta_{1,0}(\mathbf{x}) &= \{\mathbf{x} : h_1(\mathbf{x}) = h_0(\mathbf{x})\} \\ &= \{\mathbf{x} : \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = \mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})\} \\ &= \{\mathbf{x} : \log \frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(\mathbf{X} = \mathbf{x})} = 0\} \\ &= \{\mathbf{x} : \log \frac{\frac{\exp(\beta_0 + \boldsymbol{\beta}^\top \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^\top \mathbf{x})}}{\frac{1}{1 + \exp(\beta_0 + \boldsymbol{\beta}^\top \mathbf{x})}} = 0\} \\ &= \{\mathbf{x} : \beta_0 + \boldsymbol{\beta}^\top \mathbf{x} = 0\}\end{aligned}$$

↪ Maximum conditional likelihood.

- The conditional log-likelihood function :

$$\begin{aligned}L(\boldsymbol{\theta}) &= \log \mathbb{P}(Y_1 = y_1, \dots, Y_n = y_n | \mathbf{X}_1 = \mathbf{x}_1, \dots, \mathbf{X}_n = \mathbf{x}_n; \boldsymbol{\theta}) \\&= \log \prod_{i=1}^n \mathbb{P}(Y_i = y_i | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta}) \\&= \log \prod_{i=1}^n \mathbb{P}(Y_i = 1 | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta})^{y_i} \mathbb{P}(Y_i = 0 | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta})^{1-y_i} \\&= \sum_{i=1}^n y_i \log \pi(\mathbf{x}_i; \boldsymbol{\theta}) + (1 - y_i) \log (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta})) \\&= \sum_{i=1}^n y_i (\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_i) - \log(1 + \exp(\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_i)) \\&= \sum_{i=1}^n y_i (1, \mathbf{x}_i)^\top \boldsymbol{\theta} - \log\{1 + \exp((1, \mathbf{x}_i)^\top \boldsymbol{\theta})\}. \\&= \sum_{i=1}^n y_i \tilde{\mathbf{x}}_i^\top \boldsymbol{\theta} - \log\{1 + \exp(\tilde{\mathbf{x}}_i^\top \boldsymbol{\theta})\}.\end{aligned}$$

- A concave function in $\theta \leftrightarrow$ Global maximization
- However, it does not admit a closed-form solution
 \leftrightarrow Numerical optimization : Iterative Reweighted Least Squares (IRLS) Algorithm.

- Conditional log-likelihood

$$\log L(\boldsymbol{\theta}) = \sum_{i=1}^n y_i (1, \mathbf{x}_i)^\top \boldsymbol{\theta} - \log \{1 + \exp((1, \mathbf{x}_i)^\top \boldsymbol{\theta})\}.$$

- Conditional ERM : Consider the log-loss :

$$\ell(y, h_\theta(x)) = -\log(p_\theta(y|x))$$

and the hypothesis

$$h_Y(\mathbf{X}; \boldsymbol{\theta}) = \mathbb{P}_\theta(Y|\mathbf{X}) = \pi_\theta(\mathbf{X})^Y (1 - \pi_\theta(\mathbf{X}))^{1-Y}$$

- The corresponding conditional empirical risk is by definition

$$\begin{aligned} R_n(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \ell(y_i, h_\theta(x_i)) \\ &= -\frac{1}{n} \sum_{i=1}^n \log p_\theta(y_i|x_i) \\ &= -\frac{1}{n} \log L(\boldsymbol{\theta}) \end{aligned}$$

↪ With the log-loss, the conditional ERM coincides with conditional MLE.

Newton-Raphson iteration : $\theta^{(t+1)} = \theta^{(t)} - [\nabla^2 L(\theta^{(t)})]^{-1} \nabla L(\theta^{(t)})$

- Let $\tilde{\mathbf{x}}_i = (1, \mathbf{x}_i^\top)^\top$, then : $L(\theta) = \sum_{i=1}^n y_i \tilde{\mathbf{x}}_i^\top \theta - \log\{1 + \exp(\tilde{\mathbf{x}}_i^\top \theta)\}$.
- Gradient vector :

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \theta} &= \sum_{i=1}^n \left[\frac{\partial}{\partial \theta} y_i \tilde{\mathbf{x}}_i^\top \theta - \frac{\partial}{\partial \theta} \log(1 + \exp(\tilde{\mathbf{x}}_i^\top \theta)) \right] = \sum_{i=1}^n y_i \tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_i \pi(\mathbf{x}_i; \theta) \\ &= \sum_{i=1}^n \tilde{\mathbf{x}}_i (y_i - \pi(\mathbf{x}_i; \theta)) \end{aligned} \quad (1)$$

- Hessian matrix :

$$\begin{aligned} \frac{\partial^2 L(\theta)}{\partial \theta \partial \theta^\top} &= - \sum_{i=1}^n \tilde{\mathbf{x}}_i \frac{\partial}{\partial \theta^\top} \left\{ \frac{\exp(\tilde{\mathbf{x}}_i^\top \theta)}{1 + \exp(\tilde{\mathbf{x}}_i^\top \theta)} \right\} - \sum_{i=1}^n \tilde{\mathbf{x}}_i \frac{\tilde{\mathbf{x}}_i^\top \exp(\tilde{\mathbf{x}}_i^\top \theta)}{(1 + \exp(\tilde{\mathbf{x}}_i^\top \theta))^2} \\ &= - \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \pi(\mathbf{x}_i; \theta) (1 - \pi(\mathbf{x}_i; \theta)) \end{aligned} \quad (2)$$

- The Newton-Raphson iterative update of θ has therefore the following expression :

$$\theta^{(t+1)} = \theta^{(t)} + \left[\sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \pi(\mathbf{x}_i; \theta^{(t)}) (1 - \pi(\mathbf{x}_i; \theta^{(t)})) \right]^{-1} \sum_{i=1}^n \tilde{\mathbf{x}}_i (y_i - \pi(\mathbf{x}_i; \theta^{(t)}))$$

Newton-Raphson iteration : $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - [\nabla^2 L(\boldsymbol{\theta}^{(t)})]^{-1} \nabla L(\boldsymbol{\theta}^{(t)})$

■ Let $\tilde{\boldsymbol{x}}_i = (1, \boldsymbol{x}_i^\top)^\top$, then : $L(\boldsymbol{\theta}) = \sum_{i=1}^n y_i \tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta} - \log\{1 + \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})\}$.

■ Gradient vector :

$$\begin{aligned} \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= \sum_{i=1}^n \left[\frac{\partial}{\partial \boldsymbol{\theta}} y_i \tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta} - \frac{\partial}{\partial \boldsymbol{\theta}} \log(1 + \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})) \right] = \sum_{i=1}^n y_i \tilde{\boldsymbol{x}}_i - \tilde{\boldsymbol{x}}_i \pi(\boldsymbol{x}_i; \boldsymbol{\theta}) \\ &= \sum_{i=1}^n \tilde{\boldsymbol{x}}_i (y_i - \pi(\boldsymbol{x}_i; \boldsymbol{\theta})) . \end{aligned} \quad (1)$$

■ Hessian matrix :

$$\begin{aligned} \frac{\partial^2 L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} &= - \sum_{i=1}^n \tilde{\boldsymbol{x}}_i \frac{\partial}{\partial \boldsymbol{\theta}^\top} \left\{ \frac{\exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})}{1 + \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})} \right\} - \sum_{i=1}^n \tilde{\boldsymbol{x}}_i \frac{\tilde{\boldsymbol{x}}_i^\top \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})}{(1 + \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta}))^2} \\ &= - \sum_{i=1}^n \tilde{\boldsymbol{x}}_i \tilde{\boldsymbol{x}}_i^\top \pi(\boldsymbol{x}_i; \boldsymbol{\theta}) (1 - \pi(\boldsymbol{x}_i; \boldsymbol{\theta})) \end{aligned} \quad (2)$$

■ The Newton-Raphson iterative update of $\boldsymbol{\theta}$ has therefore the following expression :

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \left[\sum_{i=1}^n \tilde{\boldsymbol{x}}_i \tilde{\boldsymbol{x}}_i^\top \pi(\boldsymbol{x}_i; \boldsymbol{\theta}^{(t)}) (1 - \pi(\boldsymbol{x}_i; \boldsymbol{\theta}^{(t)})) \right]^{-1} \sum_{i=1}^n \tilde{\boldsymbol{x}}_i (y_i - \pi(\boldsymbol{x}_i; \boldsymbol{\theta}^{(t)}))$$

Newton-Raphson iteration : $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - [\nabla^2 L(\boldsymbol{\theta}^{(t)})]^{-1} \nabla L(\boldsymbol{\theta}^{(t)})$

■ Let $\tilde{\boldsymbol{x}}_i = (1, \boldsymbol{x}_i^\top)^\top$, then : $L(\boldsymbol{\theta}) = \sum_{i=1}^n y_i \tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta} - \log\{1 + \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})\}$.

■ Gradient vector :

$$\begin{aligned} \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= \sum_{i=1}^n \left[\frac{\partial}{\partial \boldsymbol{\theta}} y_i \tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta} - \frac{\partial}{\partial \boldsymbol{\theta}} \log(1 + \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})) \right] = \sum_{i=1}^n y_i \tilde{\boldsymbol{x}}_i - \tilde{\boldsymbol{x}}_i \pi(\boldsymbol{x}_i; \boldsymbol{\theta}) \\ &= \sum_{i=1}^n \tilde{\boldsymbol{x}}_i (y_i - \pi(\boldsymbol{x}_i; \boldsymbol{\theta})) . \end{aligned} \quad (1)$$

■ Hessian matrix :

$$\begin{aligned} \frac{\partial^2 L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} &= - \sum_{i=1}^n \tilde{\boldsymbol{x}}_i \frac{\partial}{\partial \boldsymbol{\theta}^\top} \left\{ \frac{\exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})}{1 + \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})} \right\} - \sum_{i=1}^n \tilde{\boldsymbol{x}}_i \frac{\tilde{\boldsymbol{x}}_i^\top \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta})}{(1 + \exp(\tilde{\boldsymbol{x}}_i^\top \boldsymbol{\theta}))^2} \\ &= - \sum_{i=1}^n \tilde{\boldsymbol{x}}_i \tilde{\boldsymbol{x}}_i^\top \pi(\boldsymbol{x}_i; \boldsymbol{\theta}) (1 - \pi(\boldsymbol{x}_i; \boldsymbol{\theta})) \end{aligned} \quad (2)$$

■ The Newton-Raphson iterative update of $\boldsymbol{\theta}$ has therefore the following expression :

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \left[\sum_{i=1}^n \tilde{\boldsymbol{x}}_i \tilde{\boldsymbol{x}}_i^\top \pi(\boldsymbol{x}_i; \boldsymbol{\theta}^{(t)}) (1 - \pi(\boldsymbol{x}_i; \boldsymbol{\theta}^{(t)})) \right]^{-1} \sum_{i=1}^n \tilde{\boldsymbol{x}}_i (y_i - \pi(\boldsymbol{x}_i; \boldsymbol{\theta}^{(t)}))$$

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \left[\sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}) (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)})) \right]^{-1} \sum_{i=1}^n \tilde{\mathbf{x}}_i (y_i - \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}))$$

Matrix form the NR iteration update :

Let

- $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)^\top$ matrix whose rows are the augmented input vectors $(1, \mathbf{x}_i^\top)$
- $\mathbf{y} = (y_1, \dots, y_n)^\top$ the vector on binary labels y_i
- $\mathbf{p} = (\pi(\mathbf{x}_1; \boldsymbol{\theta}), \dots, \pi(\mathbf{x}_n; \boldsymbol{\theta}))^\top$ the vector of logistic probabilities
- $\mathbf{W} = \text{diag}(\mathbf{p} \odot (\mathbf{1}_n - \mathbf{p}))$ diagonal matrix with $(\mathbf{W})_{ii} = \pi(\mathbf{x}_i; \boldsymbol{\theta}) (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta}))$
- $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)} + (\mathbf{W}^{(t)})^{-1}(\mathbf{y} - \mathbf{p}^{(t)})$ the current approximate response

Then

$$\Leftrightarrow \text{Vectorial form of the Gradient : } \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \tilde{\mathbf{X}}^\top (\mathbf{y} - \mathbf{p}^{(t)})$$

$$\Leftrightarrow \text{Vectorial form of the Hessian matrix : } \frac{\partial^2 L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = -\tilde{\mathbf{X}}^\top \mathbf{W} \tilde{\mathbf{X}}$$

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \left[\sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}) (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)})) \right]^{-1} \sum_{i=1}^n \tilde{\mathbf{x}}_i (y_i - \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}))$$

Matrix form the NR iteration update :

Let

- $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)^\top$ matrix whose rows are the augmented input vectors $(1, \mathbf{x}_i^\top)$
- $\mathbf{y} = (y_1, \dots, y_n)^\top$ the vector on binary labels y_i
- $\mathbf{p} = (\pi(\mathbf{x}_1; \boldsymbol{\theta}), \dots, \pi(\mathbf{x}_n; \boldsymbol{\theta}))^\top$ the vector of logistic probabilities
- $\mathbf{W} = \text{diag}(\mathbf{p} \odot (\mathbf{1}_n - \mathbf{p}))$ diagonal matrix with $(\mathbf{W})_{ii} = \pi(\mathbf{x}_i; \boldsymbol{\theta}) (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta}))$
- $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)} + (\mathbf{W}^{(t)})^{-1}(\mathbf{y} - \mathbf{p}^{(t)})$ the current approximate response

Then

$$\Leftrightarrow \text{Vectorial form of the Gradient : } \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \tilde{\mathbf{X}}^\top (\mathbf{y} - \mathbf{p}^{(t)})$$

$$\Leftrightarrow \text{Vectorial form of the Hessian matrix : } \frac{\partial^2 L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = -\tilde{\mathbf{X}}^\top \mathbf{W} \tilde{\mathbf{X}}$$

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \left[\sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}) (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)})) \right]^{-1} \sum_{i=1}^n \tilde{\mathbf{x}}_i (y_i - \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}))$$

Matrix form the NR iteration update :

Let

- $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)^\top$ matrix whose rows are the augmented input vectors $(1, \mathbf{x}_i^\top)$
- $\mathbf{y} = (y_1, \dots, y_n)^\top$ the vector on binary labels y_i
- $\mathbf{p} = (\pi(\mathbf{x}_1; \boldsymbol{\theta}), \dots, \pi(\mathbf{x}_n; \boldsymbol{\theta}))^\top$ the vector of logistic probabilities
- $\mathbf{W} = \text{diag}(\mathbf{p} \odot (\mathbf{1}_n - \mathbf{p}))$ diagonal matrix with $(\mathbf{W})_{ii} = \pi(\mathbf{x}_i; \boldsymbol{\theta}) (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta}))$
- $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)} + (\mathbf{W}^{(t)})^{-1}(\mathbf{y} - \mathbf{p}^{(t)})$ the current approximate response

Then

↪ Vectorial form of the Gradient : $\frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \tilde{\mathbf{X}}^\top (\mathbf{y} - \mathbf{p}^{(t)})$

↪ Vectorial form of the Hessian matrix : $\frac{\partial^2 L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = -\tilde{\mathbf{X}}^\top \mathbf{W} \tilde{\mathbf{X}}$

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \left[\sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}) (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)})) \right]^{-1} \sum_{i=1}^n \tilde{\mathbf{x}}_i (y_i - \pi(\mathbf{x}_i; \boldsymbol{\theta}^{(t)}))$$

Matrix form the NR iteration update :

Let

- $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)^\top$ matrix whose rows are the augmented input vectors $(1, \mathbf{x}_i^\top)$
- $\mathbf{y} = (y_1, \dots, y_n)^\top$ the vector on binary labels y_i
- $\mathbf{p} = (\pi(\mathbf{x}_1; \boldsymbol{\theta}), \dots, \pi(\mathbf{x}_n; \boldsymbol{\theta}))^\top$ the vector of logistic probabilities
- $\mathbf{W} = \text{diag}(\mathbf{p} \odot (\mathbf{1}_n - \mathbf{p}))$ diagonal matrix with $(\mathbf{W})_{ii} = \pi(\mathbf{x}_i; \boldsymbol{\theta}) (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta}))$
- $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)} + (\mathbf{W}^{(t)})^{-1}(\mathbf{y} - \mathbf{p}^{(t)})$ the current approximate response

Then

$$\hookrightarrow \text{Vectorial form of the Gradient : } \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \tilde{\mathbf{X}}^\top (\mathbf{y} - \mathbf{p}^{(t)})$$

$$\hookrightarrow \text{Vectorial form of the Hessian matrix : } \frac{\partial^2 L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = -\tilde{\mathbf{X}}^\top \mathbf{W} \tilde{\mathbf{X}}$$

Then we get the Matrix form :

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \left[\nabla^2 L(\boldsymbol{\theta}^{(t)}) \right]^{-1} \nabla L(\boldsymbol{\theta}^{(t)}) \quad (3)$$

$$= \boldsymbol{\theta}^{(t)} + \left(\tilde{\mathbf{X}}^\top \mathbf{W} \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^\top (\mathbf{y} - \mathbf{p}^{(t)})$$

$$= \boldsymbol{\theta}^{(t)} + \left(\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T (\mathbf{y} - \mathbf{p}^{(t)}) \quad (4)$$

$$= \left(\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}} \right)^{-1} \left[\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}} \boldsymbol{\theta}^{(t)} + \tilde{\mathbf{X}}^T (\mathbf{y} - \mathbf{p}^{(t)}) \right] \quad (5)$$

$$= \left(\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \left[\mathbf{W}^{(t)} \tilde{\mathbf{X}} \boldsymbol{\theta}^{(t)} + (\mathbf{y} - \mathbf{p}^{(t)}) \right] \quad (6)$$

$$= \left(\tilde{\mathbf{X}}^\top \mathbf{W}^{(t)} \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^\top \mathbf{W}^{(t)} \tilde{\mathbf{y}} \quad (7)$$

Algorithm 1 Pseudo Code for Training Logistic Regression IRLS.

Inputs : n sample $(\mathbf{x}_i, y_i)_{i=1}^n$ arranged as $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$ and $\mathbf{y} = (y_1, \dots, y_n)^\top$

Construct $\tilde{\mathbf{X}} = [\mathbf{1}_n, \mathbf{X}]$

Initialization : $\boldsymbol{\theta}^{(0)}$; set $t \leftarrow 0$ (IRLS iteration)

while increment in log-likelihood $> \epsilon$ (eg. 1e-6) **do**

$$\mathbf{p}^{(t)} = (\pi(\mathbf{x}_1; \boldsymbol{\theta}^{(t)}), \dots, \pi(\mathbf{x}_n; \boldsymbol{\theta}^{(t)}))^\top = \exp(\tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)}) \odot (\mathbf{1}_n + \exp(\tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)}))$$

$$\mathbf{W}^{(t)} = \text{diag}(\mathbf{p}^{(t)} \odot (\mathbf{1}_n - \mathbf{p}^{(t)}))$$

$$\tilde{\mathbf{z}} = \tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)} + (\mathbf{W}^{(t)})^{-1}(\mathbf{y} - \mathbf{p}^{(t)})$$

$$\boldsymbol{\theta}^{(t+1)} = (\tilde{\mathbf{X}}^\top \mathbf{W}^{(t)} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \mathbf{W}^{(t)} \tilde{\mathbf{z}}$$

% Convergence test

$$\text{log-lik} = \sum \{\mathbf{y} \odot (\tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)}) - \log(\mathbf{1}_n + \exp(\tilde{\mathbf{X}}\boldsymbol{\theta}^{(t)}))\}. \text{ \% log-likelihood.}$$

end

Result: $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}^{(t)}$ the MLE of $\boldsymbol{\theta}$

Algorithm 2 Pseudo Code for Predicting with Logistic Regression.

Inputs : Test sample $(\mathbf{x}_i)_{i=1}^n$ arranged as $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, and parameter vector $\boldsymbol{\theta}$

Construct $\tilde{\mathbf{X}} = [\mathbf{1}_n, \mathbf{X}]$

$\text{probs} = \exp(\tilde{\mathbf{X}}\boldsymbol{\theta}) \odot (\mathbf{1}_n + \exp(\tilde{\mathbf{X}}\boldsymbol{\theta}))$ % Conditional probabilities

$\hat{\mathbf{y}} = \mathbb{1}_{\text{probs} \geq 1/2}$ % Predicted labels using Bayes rule (arg max)

Result: $\hat{\mathbf{y}}$ the predicted class labels

Def. Decision boundaries

The decision boundary between each pair of classes k and ℓ , $(k, \ell) \in \mathcal{Y} \times \mathcal{Y}$ is defined by

$$\eta_{k,\ell}(\mathbf{x}) = \{\mathbf{x} : \mathbb{P}(Y = k | \mathbf{X} = \mathbf{x}) = \mathbb{P}(Y = \ell | \mathbf{X} = \mathbf{x})\}$$

- Plugin classifier : Prediction by the Bayes' decision rule

$$\hat{h}(\mathbf{x}) = \arg \max_{k \in \mathcal{Y}} \mathbb{P}(Y = k | \mathbf{X} = \mathbf{x}; \hat{\theta}) \quad (8)$$

- Plugin Decision boundaries : The decision boundary between each pair of classes k and ℓ is defined by

$$\eta_{k,\ell}(\mathbf{x}; \hat{\theta}) = \{\mathbf{x} : \mathbb{P}(Y = k | \mathbf{X} = \mathbf{x}; \hat{\theta}) = \mathbb{P}(Y = \ell | \mathbf{X} = \mathbf{x}; \hat{\theta})\}$$

ERM vs MLE : Logistic Regression : $y \in \{0, 1\}$ with $p_{\theta}(y|\mathbf{x}) = \pi_{\theta}(\mathbf{x})^y(1 - \pi_{\theta}(\mathbf{x}))^{1-y}$, and $\pi_{\theta}(\mathbf{x}) = \sigma(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}$ is the logistic function.

$$R_n(\theta) = -\frac{1}{n} \sum_{i=1}^n \log p_{\theta}(y_i|x_i) = -\frac{1}{n} \underbrace{\sum_{i=1}^n y_i(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i) - \log(1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}_i))}_{\text{Conditional log-likelihood}}$$

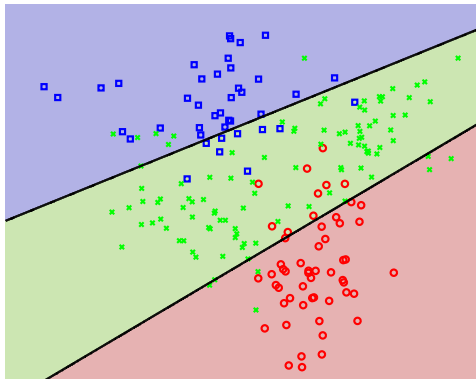
- Conditional log-likelihood ($y_i \in \{0, 1\}$)
 $\log L(\boldsymbol{\theta}) = \sum_{i=1}^n y_i (1, \mathbf{x}_i)^\top \boldsymbol{\theta} - \log\{1 + \exp((1, \mathbf{x}_i)^\top \boldsymbol{\theta})\}$.
- Conditional ERM : Consider the logistic loss : $\ell(y, h_\theta(x)) = \log(1 + \exp(-y_i h_\theta(x)))$, $y_i \in \{-1, +1\}$ and the hypothesis $h_\theta(\mathbf{X}) = \beta_0 + \beta^T \mathbf{X}$
- The corresponding conditional empirical risk is by definition

$$\begin{aligned}
 R_n(h) &= \frac{1}{n} \sum_{i=1}^n \ell(y_i, h_\theta(x_i)) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i h_\theta(x_i)}) \\
 &= \frac{1}{n} \sum_{i=1}^n \log \left(\frac{1 + e^{y_i h_\theta(x_i)}}{e^{y_i h_\theta(x_i)}} \right) = -\frac{1}{n} \sum_{i=1}^n \log \left(\frac{e^{y_i h_\theta(x_i)}}{1 + e^{y_i h_\theta(x_i)}} \right) \\
 &= -\frac{1}{n} \sum_{i=1}^n \left\{ y_i h_\theta(x_i) - \log \left(1 + e^{y_i h_\theta(x_i)} \right) \right\}; y_i \in \{-1, +1\} \\
 &= -\frac{1}{n} \begin{cases} \sum_{i=1}^n \left\{ y_i h_\theta(x_i) - \log \left(1 + e^{h_\theta(x_i)} \right) \right\} & ; y_i = 1 \\ \sum_{i=1}^n \left\{ -h_\theta(x_i) - \log \left(1 + e^{-h_\theta(x_i)} \right) \right\} & ; y_i = -1 \end{cases} \\
 &= -\frac{1}{n} \begin{cases} \sum_{i=1}^n \left\{ y_i h_\theta(x_i) - \log \left(1 + e^{h_\theta(x_i)} \right) \right\} & ; y_i = 1 \\ \sum_{i=1}^n \left\{ -\log \left(1 + e^{h_\theta(x_i)} \right) \right\} & ; y_i = -1 \end{cases} \\
 &= -\frac{1}{n} \log L(\boldsymbol{\theta})
 \end{aligned}$$

↔ With the logistic loss, the conditional ERM coincides with conditional MLE.

Multi-class logistic regression

Multi-class Logistic Regression



- $\mathbf{X} \in \mathcal{X} = \mathbb{R}^d$ and $Y_i \in \mathcal{Y} = \{1, \dots, K\}$
- Conditional (Discriminative) model : for $k = 1, \dots, K - 1$

$$\mathbb{P}(Y = k|\mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\alpha_k + \boldsymbol{\beta}_k^T \mathbf{x})}{1 + \sum_{\ell=1}^{K-1} \exp(\alpha_\ell + \boldsymbol{\beta}_\ell^T \mathbf{x})} = \pi_k(\mathbf{x}_i; \boldsymbol{\theta})$$

- for $k = K$, $\mathbb{P}(Y = K|\mathbf{x}; \boldsymbol{\theta}) = 1 - \sum_{k=1}^{K-1} \mathbb{P}(Y = k|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(\alpha_\ell + \boldsymbol{\beta}_\ell^T \mathbf{x})}$.
This is equivalent to setting $(\alpha_K, \boldsymbol{\beta}_K^T)^T = \mathbf{0}$.

- Link function : for $k = 1, \dots, K$

$$\log \frac{\mathbb{P}(Y = k|\mathbf{x}; \boldsymbol{\theta})}{\mathbb{P}(Y = K|\mathbf{x}; \boldsymbol{\theta})} = \alpha_k + \boldsymbol{\beta}_k^T \mathbf{x}$$

- The model parameter : $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K)$ with $\boldsymbol{\theta}_k = (\alpha_k, \boldsymbol{\beta}_k^T)^T$ ($k = 1, \dots, K - 1$)
- Maximum conditional likelihood estimation : The conditional log-likelihood of $\boldsymbol{\theta}$

$$\begin{aligned} L(\boldsymbol{\theta}) &= \log \prod_{i=1}^n \mathbb{P}(Y_i | \mathbf{x}_i; \boldsymbol{\theta}) = \log \prod_{i=1}^n \prod_{k=1}^K \mathbb{P}(Y_i = k | \mathbf{x}_i; \boldsymbol{\theta})^{y_{ik}} \\ &= \sum_{i=1}^n \sum_{k=1}^K y_{ik} \log \pi_k(\mathbf{x}_i; \boldsymbol{\theta}) \end{aligned}$$

where we have used the notation $y_{ik} = \mathbb{1}_{y_i \neq k}$, i.e. $y_{ik} = 1$ iff $y_i = k$

- This log-likelihood is convex but can not be maximized in a closed form.
- The Newton-Raphson (NR) algorithm :

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \left[\frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right]_{\boldsymbol{\theta} = \boldsymbol{\theta}^{(t)}}^{-1} \frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta} = \boldsymbol{\theta}^{(t)}}$$

- The Newton-Raphson algorithm is an iterative numerical optimization algorithm
- starts from an initial arbitrary solution $\theta^{(0)}$, and updates the estimation of θ
- A single NR update is given by :

$$\theta^{(t+1)} = \theta^{(t)} - \left[\frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta \partial \theta^T} \right]^{-1} \frac{\partial \mathcal{L}(\theta)}{\partial \theta} \quad (9)$$

where the Hessian and the gradient of $\mathcal{L}(\theta)$ (which are respectively the second and first derivative of $\mathcal{L}(\theta)$) are evaluated at $\theta = \theta^{(t)}$.

- NR can be stopped when the relative variation of $\mathcal{L}(\theta)$ is below a prefixed threshold.

- **Gradient vector** : $\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = \left(\left(\frac{\partial \mathcal{L}(\theta)}{\partial \theta_1} \right)^T, \dots, \left(\frac{\partial \mathcal{L}(\theta)}{\partial \theta_{K-1}} \right)^T \right)^T$ where $\forall k \in [K-1]$:

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta_k} = \sum_{i=1}^n (y_{ik} - \pi_k(\mathbf{x}_i; \theta)) \mathbf{x}_i = \mathbf{X}^T (\mathbf{y}_k - \mathbf{p}_k)$$

- i) $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$: $n \times (d+1)$ matrix whose rows are the inputs \mathbf{x}_i ,
 - ii) $\mathbf{y}_k = (y_{1k}, \dots, y_{nk})^T$: $n \times 1$ vector of indicator variables y_{ik}
 - iii) $\mathbf{p}_k = (\pi_k(\mathbf{x}_1; \theta), \dots, \pi_k(\mathbf{x}_n; \theta))^T$: $n \times 1$ vector of logistic probabilities
- Vectorized form of the gradient of $\mathcal{L}(\theta)$ for all the logistic components :

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = \begin{pmatrix} \mathbf{X}^T & 0 & \dots & 0 \\ 0 & \mathbf{X}^T & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{X}^T \end{pmatrix} \begin{pmatrix} \mathbf{y}_1 - \mathbf{p}_1 \\ \mathbf{y}_2 - \mathbf{p}_2 \\ \vdots \\ \mathbf{y}_{K-1} - \mathbf{p}_{K-1} \end{pmatrix} = \tilde{\mathbf{X}}^T (\mathbf{Y} - \mathbf{P}) \quad (10)$$

- i) $\mathbf{Y} = (\mathbf{y}_1^T, \dots, \mathbf{y}_{K-1}^T)^T$: $n \times (K-1)$ column vector
- ii) $\mathbf{P} = (\mathbf{p}_1^T, \dots, \mathbf{p}_{K-1}^T)^T$: $n \times (K-1)$ column vector
- iii) $\tilde{\mathbf{X}} = (\mathbf{X}^T, \dots, \mathbf{X}^T)^T$: $(n \times (K-1))$ by $(d+1)$ matrix of $K-1$ copies of \mathbf{X} .

- **Hessian matrix** : composed of $(K - 1) \times (K - 1)$ block matrices $\left\{ \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_\ell \partial \boldsymbol{\theta}_k^T} \right\}_{k, \ell=1}^{K-1}$

$$\frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} = \begin{pmatrix} \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1 \partial \boldsymbol{\theta}_1^T} & \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1 \partial \boldsymbol{\theta}_2^T} & \cdots & \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1 \partial \boldsymbol{\theta}_{K-1}^T} \\ \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_2 \partial \boldsymbol{\theta}_1^T} & & \cdots & \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1 \partial \boldsymbol{\theta}_{K-1}^T} \\ \vdots & \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_\ell \partial \boldsymbol{\theta}_k^T} & & \vdots \\ \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_{K-1} \partial \boldsymbol{\theta}_1^T} & & \cdots & \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_{K-1} \partial \boldsymbol{\theta}_{K-1}^T} \end{pmatrix}$$

where each block matrix is of dimension $(d + 1) \times (d + 1)$ and is given by :

$$\begin{aligned} \frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_\ell \partial \boldsymbol{\theta}_k^T} &= - \sum_{i=1}^n \pi_k(\mathbf{x}_i; \boldsymbol{\theta}) (\delta_{k\ell} - \pi_\ell(\mathbf{x}_i; \boldsymbol{\theta})) \mathbf{x}_i \mathbf{x}_i^T \\ &= -\mathbf{X}^T \mathbf{W}_{k\ell} \mathbf{X} \end{aligned}$$

- $\mathbf{W}_{k\ell}$: $n \times n$ diagonal matrix whose diagonal elements are $\pi_k(\mathbf{x}_i; \boldsymbol{\theta}) (\delta_{k\ell} - \pi_\ell(\mathbf{x}_i; \boldsymbol{\theta}))$ for $i = 1, \dots, n$.

- For all the logistic components ($k, \ell = 1, \dots, K - 1$), the Hessian takes the form :

$$\frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} = -\tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}} \quad (11)$$

→ $\mathbf{W} : (n \times (K - 1))$ by $(n \times (K - 1))$ matrix composed of $(K - 1) \times (K - 1)$ block matrices, each block is $\boldsymbol{\theta}_{k\ell}$ ($k, \ell = 1, \dots, K - 1$).

- It can be shown that the Hessian matrix for the multi-class logistic regression model is positive semi definite and therefore the log-likelihood is concave.

$$\begin{aligned} \mathbf{H} &= - \begin{pmatrix} \mathbf{X}^T \mathbf{W}_{1,1} \mathbf{X} & \dots & \mathbf{X}^T \mathbf{W}_{1,K-1} \mathbf{X} \\ \vdots & \ddots & \vdots \\ \mathbf{X}^T \mathbf{W}_{K-1,1} \mathbf{X} & \dots & \mathbf{X}^T \mathbf{W}_{K-1,K-1} \mathbf{X} \end{pmatrix} \\ &= - \begin{pmatrix} \mathbf{X}^T & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \mathbf{X}^T \end{pmatrix} \begin{pmatrix} \mathbf{W}_{1,1} & \dots & \mathbf{W}_{1,K-1} \\ \vdots & \ddots & \vdots \\ \mathbf{W}_{K-1,1} & \dots & \mathbf{W}_{K-1,K-1} \end{pmatrix} \begin{pmatrix} \mathbf{X} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \mathbf{X} \end{pmatrix} \\ &= -\tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}} \end{aligned}$$

The NR algorithm in this case can therefore be reformulated as the IRLS

$$\begin{aligned}\boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} - \left[\frac{\partial^2 \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right]_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}^{-1} \frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \\ &= \boldsymbol{\theta}^{(t)} + (\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T (\mathbf{Y} - \mathbf{P}^{(t)}) \\ &= (\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}})^{-1} \left[\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}} \boldsymbol{\theta}^{(t)} + \tilde{\mathbf{X}}^T (\mathbf{Y} - \mathbf{P}^{(t)}) \right] \\ &= (\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \left[\mathbf{W}^{(t)} \tilde{\mathbf{X}} \boldsymbol{\theta}^{(t)} + (\mathbf{Y} - \mathbf{P}^{(t)}) \right] \\ &= (\tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{W}^{(t)} \tilde{\mathbf{Y}}\end{aligned}$$

where $\tilde{\mathbf{Y}} = \tilde{\mathbf{X}} \boldsymbol{\theta}^{(t)} + (\mathbf{W}^{(t)})^{-1} (\mathbf{Y} - \mathbf{P}^{(t)})$ which yields in the Iteratively Reweighted Least Squares (IRLS) algorithm.

Tasks :

- Implement (from the scratch) each of the following functions and apply them to the given data :
 - ▶ `train_reglog` and `predict_reglog`
 - ▶ `irls` should be in a separate function

Datasets :

- ▶ Training data `Xtrain.txt` and `ytrain.txt`
 - ▶ Testing data : `Xtest.txt`
 - Plot the results by highlighting the classification and the generative model for each class
 - compare your results to those you could obtain by using standard packages
- ```
from sklearn.linear_model import LogisticRegression
or GLM from statsmodels
```