

# Étude des Fonctions de Perte par Calcul Symbolique

Gradients, Convexité et Applications

DONGMO TCHOUMENE ANITA BELVIANE - **22W2184**

DONFACK SYNTHIA CALORINE - **22U2073**

BOKOU-BOUNA-ANGE-LARISSA - **22W2188**

JIATSA ROMMEL JUNIOR - **22T2906**

Master I Data Science - Techniques d'Optimisation Avancée

2 octobre 2025

# Plan

- 1 Introduction aux Fonctions de Perte
- 2 1. Erreur Quadratique Moyenne (MSE)
- 3 2. Entropie Croisée Binaire (BCE)
- 4 3. Entropie Croisée Catégorielle (CCE)
- 5 4. Perte de Huber
- 6 Synthèse et Comparaison

# Fonctions de Perte Étudiées

Dans ce TP, nous étudions quatre fonctions de perte fondamentales :

- **Erreur Quadratique Moyenne (MSE)** - Régression
- **Entropie Croisée Binaire (BCE)** - Classification binaire
- **Entropie Croisée Catégorielle (CCE)** - Classification multi-classes
- **Perte de Huber** - Régression robuste

## Objectifs :

- 1 Calculer les gradients analytiquement
- 2 Étudier la convexité
- 3 Visualiser et appliquer sur des données réelles

# MSE : Définition et Formulation

## Formulation :

$$L_{\text{MSE}}(\mathbf{w}) = \frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

où  $\hat{y}_i = \mathbf{w}^T \mathbf{x}_i$  pour la régression linéaire.

## Forme vectorielle :

$$L_{\text{MSE}}(\mathbf{w}) = \frac{1}{2m} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

**Utilisation :** Problèmes de régression, sensible aux outliers.

# MSE : Calcul du Gradient par SymPy

## Listing – Gradient de MSE

```
1 import sympy as sp
2 import numpy as np
3
4 # Définition des symboles
5 w = sp.Matrix(sp.symbols('w1:4')) # w = [w1, w2, w3]
6 X = sp.Matrix(sp.symbols('x1:13')).reshape(4, 3) # 4x3
7 y = sp.Matrix(sp.symbols('y1:5')) # vecteur cible
8 m = sp.symbols('m', positive=True)
9
10 # Prédictions
11 y_pred = X * w
12
13 # MSE Loss
14 mse_loss = (1/(2*m)) * (y - y_pred).T * (y - y_pred)
15 mse_loss = mse_loss[0] # extraction scalaire
16
```

# MSE : Résultats du Gradient

**Gradient analytique :**

$$\nabla_{\mathbf{w}} L_{\text{MSE}} = \frac{1}{m} \mathbf{X}^T (\mathbf{X} \mathbf{w} - \mathbf{y})$$

**Matrice Hessienne :**

$$\mathbf{H}_{\text{MSE}} = \frac{1}{m} \mathbf{X}^T \mathbf{X}$$

**Propriétés :**

- Hessienne constante (indépendante de  $\mathbf{w}$ )
- $\mathbf{X}^T \mathbf{X}$  est semi-définie positive
- Si  $\text{rang}(\mathbf{X}) = p$ , alors strictement convexe

**Critère de convexité :** Une fonction est convexe si sa Hessienne est semi-définie positive.

**Pour MSE :**

$$\mathbf{H}_{\text{MSE}} = \frac{1}{m} \mathbf{X}^T \mathbf{X} \succeq 0$$

**Preuve :** Pour tout vecteur  $\mathbf{v}$  :

$$\mathbf{v}^T \mathbf{H}_{\text{MSE}} \mathbf{v} = \frac{1}{m} \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} = \frac{1}{m} \|\mathbf{X} \mathbf{v}\|_2^2 \geq 0$$

**Conclusion :** MSE est **convexe** avec un unique minimum global.

## Formulation :

$$L_{\text{BCE}}(\mathbf{w}) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)]$$

où  $\hat{p}_i = \sigma(\mathbf{w}^T \mathbf{x}_i) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}$  (fonction sigmoïde)

**Utilisation :** Classification binaire ( $y_i \in \{0, 1\}$ )

**Interprétation :** Mesure la divergence entre les distributions prédite et réelle.



# BCE : Calcul du Gradient par SymPy

## Listing – Gradient de BCE

```
1 # Fonction sigmoïde
2 def sigmoid_sympy(z):
3     return 1 / (1 + sp.exp(-z))
4
5 # Logits et probabilités prédites
6 z = X * w # logits
7 p_pred = sp.Matrix([sigmoid_sympy(zi) for zi in z])
8
9 # BCE Loss (pour un échantillon)
10 bce_loss = 0
11 for i in range(len(y)):
12     bce_loss += -(y[i] * sp.log(p_pred[i]) +
13                  (1 - y[i]) * sp.log(1 - p_pred[i]))
14 bce_loss = bce_loss / m
15
16 # Gradient
```

**Gradient analytique :**

$$\nabla_{\mathbf{w}} L_{\text{BCE}} = \frac{1}{m} \mathbf{X}^T (\hat{\mathbf{p}} - \mathbf{y})$$

où  $\hat{\mathbf{p}} = [\hat{p}_1, \hat{p}_2, \dots, \hat{p}_m]^T$

**Matrice Hessienne :**

$$\mathbf{H}_{\text{BCE}} = \frac{1}{m} \mathbf{X}^T \mathbf{D} \mathbf{X}$$

où  $\mathbf{D} = \text{diag}(\hat{p}_1(1 - \hat{p}_1), \dots, \hat{p}_m(1 - \hat{p}_m))$

**Note :**  $\hat{p}_i(1 - \hat{p}_i) > 0$  pour  $\hat{p}_i \in (0, 1)$

**Analyse de la Hessienne :**

$$\mathbf{H}_{\text{BCE}} = \frac{1}{m} \mathbf{X}^T \mathbf{D} \mathbf{X}$$

**Matrice  $\mathbf{D}$  :** - Diagonale avec  $d_{ii} = \hat{p}_i(1 - \hat{p}_i) \geq 0$  -  $\mathbf{D} \succeq 0$  (semi-définie positive)

**Preuve de convexité :**

$$\mathbf{v}^T \mathbf{H}_{\text{BCE}} \mathbf{v} = \frac{1}{m} (\mathbf{X} \mathbf{v})^T \mathbf{D} (\mathbf{X} \mathbf{v}) \geq 0$$

**Conclusion :** BCE est **convexe** avec un unique minimum global.

## Formulation :

$$L_{\text{CCE}}(\mathbf{W}) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_{i,k} \log(\hat{p}_{i,k})$$

où  $\hat{p}_{i,k} = \frac{e^{z_{i,k}}}{\sum_{j=1}^K e^{z_{i,j}}}$  (fonction softmax)

et  $z_{i,k} = \mathbf{w}_k^T \mathbf{x}_i$  (logit pour la classe  $k$ )

**Utilisation :** Classification multi-classes ( $K > 2$ )

**Encodage :**  $\mathbf{y}_i$  en one-hot encoding

# CCE : Calcul du Gradient par SymPy

## Listing – Gradient de CCE

```
1 K = 3 # nombre de classes
2 W = sp.Matrix(sp.symbols('w1:10')).reshape(3, K) # matrice des poids
3 Y = sp.Matrix(sp.symbols('y1:13')).reshape(4, K) # one-hot encoding
4
5 # Fonction softmax
6 def softmax_sympy(z_vec):
7     exp_z = sp.Matrix([sp.exp(zi) for zi in z_vec])
8     sum_exp = sum(exp_z)
9     return sp.Matrix([exp_zi / sum_exp for exp_zi in exp_z])
10
11 # Calcul des logits et probabilités
12 Z = X * W # logits pour chaque classe
13 P = sp.zeros(4, K)
14
15 for i in range(4): # pour chaque échantillon
16     z_i = Z.row(i).T
```

**Gradient analytique** : Pour la classe  $k$  :

$$\nabla_{\mathbf{w}_k} L_{\text{CCE}} = \frac{1}{m} \mathbf{X}^T (\hat{\mathbf{p}}_k - \mathbf{y}_k)$$

où  $\hat{\mathbf{p}}_k$  et  $\mathbf{y}_k$  sont les vecteurs de probabilités et labels pour la classe  $k$ .

**Forme compacte** :

$$\nabla_{\mathbf{W}} L_{\text{CCE}} = \frac{1}{m} \mathbf{X}^T (\hat{\mathbf{P}} - \mathbf{Y})$$

**Hessienne** : Plus complexe, mais reste semi-définie positive.

**Propriété fondamentale :** La fonction softmax avec entropie croisée est convexe par rapport aux logits.

**Preuve intuitive :** - La fonction  $-\log(\cdot)$  est convexe - La fonction softmax préserve la convexité - La composition preserve la convexité

**Hessienne :**

$$\mathbf{H}_{\text{CCE}} = \frac{1}{m} \mathbf{X}^T (\mathbf{D} - \hat{\mathbf{P}} \hat{\mathbf{P}}^T) \mathbf{X}$$

où  $\mathbf{D}$  est diagonale avec  $d_{ii} = \hat{p}_i(1 - \hat{p}_i)$ .

**Conclusion :** CCE est **convexe**.

# Huber : Définition et Formulation

## Formulation :

$$L_{\text{Huber}}(\mathbf{w}; \delta) = \frac{1}{m} \sum_{i=1}^m H_{\delta}(y_i - \hat{y}_i)$$

où :

$$H_{\delta}(r) = \begin{cases} \frac{1}{2}r^2 & \text{si } |r| \leq \delta \\ \delta|r| - \frac{1}{2}\delta^2 & \text{si } |r| > \delta \end{cases}$$

**Paramètre  $\delta$**  : Seuil de transition entre comportement quadratique et linéaire.

**Avantage** : Robuste aux outliers (moins sensible que MSE).



# Huber : Calcul du Gradient par SymPy

## Listing – Gradient de Huber

```
1 delta = sp.symbols('delta', positive=True)
2
3 # Fonction de Huber
4 def huber_loss_sympy(r, delta):
5     return sp.Piecewise(
6         (sp.Rational(1, 2) * r**2, sp.Abs(r) <= delta),
7         (delta * sp.Abs(r) - sp.Rational(1, 2) * delta**2, True)
8     )
9
10 # Résidus
11 residuals = y - X * w
12
13 # Huber Loss totale
14 huber_total = sum([huber_loss_sympy(residuals[i], delta)
15                    for i in range(len(residuals))]) / m
```

## Gradient/Sous-gradient :

$$\frac{\partial H_\delta(r)}{\partial r} = \begin{cases} r & \text{si } |r| < \delta \\ \delta \cdot \text{sign}(r) & \text{si } |r| > \delta \\ [-\delta, \delta] & \text{si } r = \pm\delta \end{cases}$$

## Gradient par rapport aux paramètres :

$$\nabla_{\mathbf{w}} L_{\text{Huber}} = \frac{1}{m} \mathbf{X}^T \mathbf{g}$$

où  $g_i = \frac{\partial H_\delta(r_i)}{\partial r_i}$  avec  $r_i = y_i - \mathbf{w}^T \mathbf{x}_i$

## Analyse de convexité :

**Pour**  $|r| < \delta$  :  $H_\delta(r) = \frac{1}{2}r^2$  est strictement convexe.

**Pour**  $|r| > \delta$  :  $H_\delta(r) = \delta|r| - \frac{1}{2}\delta^2$  est convexe (fonction linéaire par morceaux).

**À**  $|r| = \delta$  : Les dérivées secondes à gauche et droite sont positives.

**Hessienne généralisée** : - Définie positive dans la région quadratique - Semi-définie positive globalement

**Conclusion** : La perte de Huber est **convexe** mais non différentiable aux points de transition.

# Tableau Comparatif des Fonctions de Perte

Fonction	Domaine	Convexité	Différentiabilité	Robustesse
MSE	Régression	Convexe	$C^\infty$	Faible
BCE	Classif. binaire	Convexe	$C^\infty$	Moyenne
CCE	Classif. multi	Convexe	$C^\infty$	Moyenne
Huber	Régression	Convexe	$C^1$ non $C^2$	Élevée

## Points clés :

- Toutes les fonctions étudiées sont **convexes**
- Garantissent un **minimum global unique**
- Algorithmes de gradient convergent vers l'optimum
- Huber offre le meilleur compromis robustesse/efficacité

## Choix de la fonction de perte selon le contexte :

- **MSE** : Régression avec données propres, interprétation en termes de variance
- **BCE** : Classification binaire, probabilités calibrées
- **CCE** : Classification multi-classes, extension naturelle de BCE
- **Huber** : Régression robuste, présence d'outliers

**Optimisation** : - Gradient descent, Newton, quasi-Newton - Convergence garantie (fonctions convexes) - Choix du pas d'apprentissage critique

## Ce que nous avons accompli :

- ➊ **Gradients analytiques** calculés avec SymPy
- ➋ **Propriétés de convexité** démontrées pour chaque fonction
- ➌ **Comparaison** des caractéristiques et domaines d'application

## Prochaines étapes :

- Implémentation et visualisation sur données réelles
- Comparaison empirique des performances
- Étude de l'influence des hyperparamètres ( $\delta$  pour Huber)

**Outils utilisés :** SymPy pour le calcul symbolique, théorie de l'optimisation convexe.

Questions ?