

CS2109S Tutorial 5

Classification and Logistic Regression

(AY 24/25 Semester 2)

March 14, 2025

(Prepared by Benson)

Contents

Classification

Q1. Linear vs Non-linear Separability

Logistic Regression

Recap

Q2. Loss Function of Logistic Regression

Q4. Logistic Regression for Multi-Class Classification

Performance measures

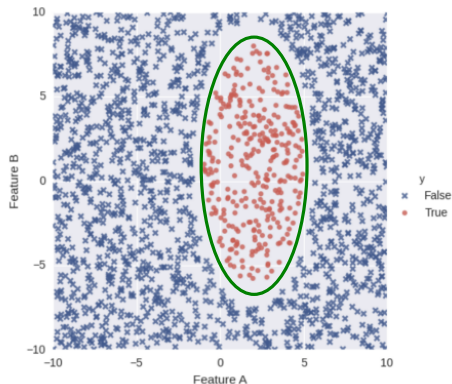
Q3. Precision, recall, F1 score and ROC curve

Q5. Evaluating Logistic Regression

Bonus. Sigmoid vs Softmax

Q1. Linear vs Non-linear Separability

Define a (minimal) set of features that will perfectly classify whether or not a bunny can be released into the wild.



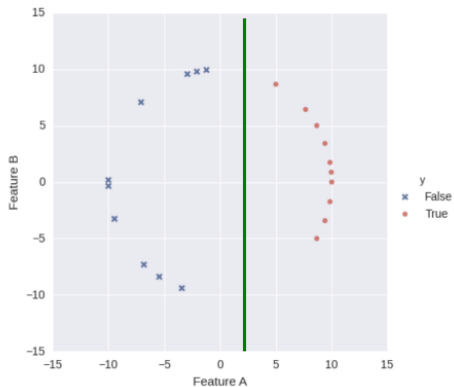
Equation of an ellipse:

$$\frac{(A - x)^2}{a^2} + \frac{(B - y)^2}{b^2} = 1$$

A minimal set of features is (A^2, B^2, A, B) .

Q1. Linear vs Non-linear Separability

Define a (minimal) set of features that will perfectly classify whether or not a bunny can be released into the wild.



A minimal set of features is (A).

Recap: Logistic Regression

Target: Solve **classification** tasks (output a **probability** between 0 and 1).

- ▶ First Try (regression): $p = \mathbf{w} \cdot \mathbf{x}$
 - ▶ FAIL: $\mathbf{w} \cdot \mathbf{x}$ can be outside $[0, 1]$.
- ▶ Second Try (odds ratio): $\frac{p}{1-p} = \mathbf{w} \cdot \mathbf{x}$
 - ▶ FAIL: $\frac{p}{1-p}$ is always non-negative, but $\mathbf{w} \cdot \mathbf{x}$ can be outside $[0, \infty)$.
- ▶ Third Try (logits): $\log\left(\frac{p}{1-p}\right) = \mathbf{w} \cdot \mathbf{x}$ (works!)
 - ▶ Rearrange terms: $p = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$

Recap: Logistic Regression

Hypothesis:

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

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Loss function:

- ▶ **Mean Squared Error** is not convex under logistic regression \Rightarrow gradient descent might not reach global minimum.
- ▶ Use **binary cross entropy loss** instead:

$$\begin{aligned} BCE(\hat{y}) &= \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases} \\ &= -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y}) \end{aligned}$$

“surprisal” in entropy!

Q2. Loss Function of Logistic Regression

- (a) Write down the probability p as a function of \mathbf{x} and calculate the derivative of $\log(p)$ with respect to each weight w_i .

$$\begin{aligned}
 p &= \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}} = \frac{1}{1 + e^{\sum_{i=1}^n -w_i x_i}} \\
 \log(p) &= \log\left(\frac{1}{1 + e^{\sum_{i=1}^n -w_i x_i}}\right) \stackrel{\log \frac{1}{x} = -\log x}{=} -\log\left(1 + e^{\sum_{i=1}^n -w_i x_i}\right) \\
 \frac{\partial \log(p)}{\partial w_i} &= -\left(\frac{1}{1 + e^{\sum_{i=1}^n -w_i x_i}} \cdot \frac{\partial}{\partial w_i} \left(1 + e^{\sum_{i=1}^n -w_i x_i}\right)\right) \quad \text{Chain Rule} \\
 &= -\left(\frac{1}{1 + e^{\sum_{i=1}^n -w_i x_i}} \cdot \left(e^{\sum_{i=1}^n -w_i x_i} \cdot (-x_i)\right)\right) \quad \text{Chain Rule} \\
 &= \frac{e^{\sum_{i=1}^n -w_i x_i}}{1 + e^{\sum_{i=1}^n -w_i x_i}} \cdot (x_i) = (1 - p)x_i
 \end{aligned}$$

Q2. Loss Function of Logistic Regression

- (b) Write down the probability $1 - p$ as a function of x and calculate the derivative of $\log(1 - p)$ with respect to each weight w_i .

$$1 - p = \frac{\frac{1}{1+e^{-w^\top x}} e^{-w^\top x}}{1 + e^{-w^\top x}} \cdot \frac{e^{w^\top x}}{e^{w^\top x}} = \frac{1}{1 + e^{w^\top x}} = \frac{1}{1 + e^{w \cdot x}} = \frac{1}{1 + e^{\sum_{i=1}^n w_i x_i}}$$

$$\log(1 - p) = \log\left(\frac{1}{1 + e^{\sum_{i=1}^n w_i x_i}}\right) = -\log\left(1 + e^{\sum_{i=1}^n w_i x_i}\right)$$

$$\begin{aligned} \frac{\partial \log(1 - p)}{\partial w_i} &= - \left(\frac{1}{1 + e^{\sum_{i=1}^n w_i x_i}} \cdot \frac{\partial}{\partial w_i} \left(1 + e^{\sum_{i=1}^n w_i x_i} \right) \right) \quad \text{Chain Rule} \\ &= - \left(\frac{1}{1 + e^{\sum_{i=1}^n w_i x_i}} \cdot (e^{\sum_{i=1}^n w_i x_i}) (x_i) \right) \quad \text{Chain Rule} \\ &= - \frac{e^{\sum_{i=1}^n w_i x_i}}{1 + e^{\sum_{i=1}^n w_i x_i}} \cdot \frac{\frac{\partial}{\partial x} (1 + e^x)}{\frac{\partial}{\partial w_i} (\sum_{i=1}^n w_i x_i)} \cdot (x_i) = - \frac{1}{1 + e^{\sum_{i=1}^n w_i x_i}} \cdot (x_i) = \boxed{-px_i} \end{aligned}$$

Q2. Loss Function of Logistic Regression

Derivative of Sigmoid Function

Let $\sigma(z) = \frac{1}{1 + e^{-z}}$. We have $\sigma'(z) = \sigma(z)(1 - \sigma(z))$.

Proof.

$$\begin{aligned}\sigma'(z) &= \frac{\frac{d}{dz}(1) \cdot (1 + e^{-z}) - (1) \cdot \frac{d}{dz}(1 + e^{-z})}{(1 + e^{-z})^2} \\ &= \frac{(0) \cdot (1 + e^{-z}) - (1) \cdot (-e^{-z})}{(1 + e^{-z})^2} \\ &= \frac{e^{-z}}{(1 + e^{-z})^2} \\ &= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}} = \boxed{\sigma(z)(1 - \sigma(z))}\end{aligned}$$

◀ quotient rule

Q2. Loss Function of Logistic Regression

Derivative of Sigmoid Function

Let $\sigma(z) = \frac{1}{1 + e^{-z}}$. We have $\sigma'(z) = \sigma(z)(1 - \sigma(z))$.

Easier solution to (a), (b):

$$\begin{aligned} \frac{\partial \log(p)}{\partial w_i} &= \underbrace{\frac{1}{p}}_{\frac{\partial \log(p)}{\partial p}} \cdot \underbrace{\frac{\partial p}{\partial w_i}}_{\substack{p = \sigma(\mathbf{w}^\top \mathbf{x}) \\ \frac{\partial p}{\partial (\mathbf{w}^\top \mathbf{x})}}} = \frac{1}{p} \cdot \cancel{p} \cdot \cancel{(1-p)} \cdot \underbrace{\frac{\partial}{\partial w_i}(\mathbf{w}^\top \mathbf{x})}_{\substack{\mathbf{w}^\top \mathbf{x} = \mathbf{w} \cdot \mathbf{x} = \sum w_i x_i \\ = x_i}} = \boxed{(1-p)x_i} \\ \frac{\partial \log(1-p)}{\partial w_i} &= \underbrace{\frac{-1}{1-p}}_{\frac{\partial \log(1-p)}{\partial p}} \cdot \frac{\partial p}{\partial w_i} = \frac{-1}{\cancel{1-p}} \cdot \cancel{p} \cdot \cancel{(1-p)} \cdot \underbrace{\frac{\partial}{\partial w_i}(\mathbf{w}^\top \mathbf{x})}_{= x_i} = \boxed{-px_i} \end{aligned}$$

Q2. Loss Function of Logistic Regression

- (c) Using results from 3(a) and 3(b), derive $\frac{\partial L}{\partial w_i}$, where L is the loss function of logistic regression model.

$$L = -y \log(p) - (1 - y) \log(1 - p)$$

$$\begin{aligned}\frac{\partial L}{\partial w_i} &= -y \frac{\partial \log(p)}{\partial w_i} - (1 - y) \frac{\partial \log(1 - p)}{\partial w_i} \\ &= -y(1 - p)x_i - (1 - y)(-px_i) \\ &= -yx_i + ypx_i + px_i - ypx_i \\ &= x_i(p - y) \\ &= x_i(h_{\mathbf{w}}(\mathbf{x}) - y)\end{aligned}$$

Q4. Logistic Regression for Multi-Class Classification

- (a) Compute the probability of an animal belonging to a certain class and classify them accordingly.

First animal: $\mathbf{x} = [1 \quad 4.2 \quad 0.4]^\top$

$$\blacktriangleright \mathbf{w}_{cat} \cdot \mathbf{x} = 1 \cdot 4.2 + 4.2 \cdot (-0.01) + 0.4 \cdot (-0.12) = 4.11$$

$$p_{cat} = \frac{1}{1 + e^{-4.11}} = 0.984$$

$$\blacktriangleright \mathbf{w}_{horse} \cdot \mathbf{x} = -6.336$$

$$p_{horse} = \frac{1}{1 + e^{6.336}} = 0.00177$$

$$\blacktriangleright \mathbf{w}_{elephant} \cdot \mathbf{x} = -1246.196$$

$$p_{elephant} = \frac{1}{1 + e^{1246.196}} \approx 0$$

$$\mathbf{w}_{cat} = [4.2 \quad -0.01 \quad -0.12]^\top$$

$$\mathbf{w}_{horse} = [-20 \quad -0.08 \quad 35]^\top$$

$$\mathbf{w}_{elephant} = [-1250 \quad 0.82 \quad 0.9]^\top$$

Weight (kg)	Length (m)
4.2	0.4
720	2.4
2350	5.5

Q4. Logistic Regression for Multi-Class Classification

- (a) Compute the probability of an animal belonging to a certain class and classify them accordingly.

Second animal: $\mathbf{x} = [1 \quad 720 \quad 2.4]^\top$

► $\mathbf{w}_{cat} \cdot \mathbf{x} = -3.288$

$$p_{cat} = \frac{1}{1 + e^{3.288}} = 0.0360$$

► $\mathbf{w}_{horse} \cdot \mathbf{x} = 6.4$

$$p_{horse} = \frac{1}{1 + e^{-6.4}} = 0.998$$

► $\mathbf{w}_{elephant} \cdot \mathbf{x} = -657.44$

$$p_{elephant} = \frac{1}{1 + e^{657.44}} \approx 0$$

$$\mathbf{w}_{cat} = [4.2 \quad -0.01 \quad -0.12]^\top$$

$$\mathbf{w}_{horse} = [-20 \quad -0.08 \quad 35]^\top$$

$$\mathbf{w}_{elephant} = [-1250 \quad 0.82 \quad 0.9]^\top$$

Weight (kg)	Length (m)
4.2	0.4
720	2.4
2350	5.5

Q4. Logistic Regression for Multi-Class Classification

- (a) Compute the probability of an animal belonging to a certain class and classify them accordingly.

Third animal: $\mathbf{x} = [1 \quad 2350 \quad 5.5]^\top$

► $\mathbf{w}_{cat} \cdot \mathbf{x} = -19.96$

$$p_{cat} = \frac{1}{1 + e^{19.96}} \approx 0$$

► $\mathbf{w}_{horse} \cdot \mathbf{x} = -15.5$

$$p_{horse} = \frac{1}{1 + e^{15.5}} \approx 0$$

► $\mathbf{w}_{elephant} \cdot \mathbf{x} = 681.95$

$$p_{elephant} = \frac{1}{1 + e^{-681.95}} \approx 1$$

$$\mathbf{w}_{cat} = [4.2 \quad -0.01 \quad -0.12]^\top$$

$$\mathbf{w}_{horse} = [-20 \quad -0.08 \quad 35]^\top$$

$$\mathbf{w}_{elephant} = [-1250 \quad 0.82 \quad 0.9]^\top$$

Weight (kg)	Length (m)
4.2	0.4
720	2.4
2350	5.5

Q4. Logistic Regression for Multi-Class Classification

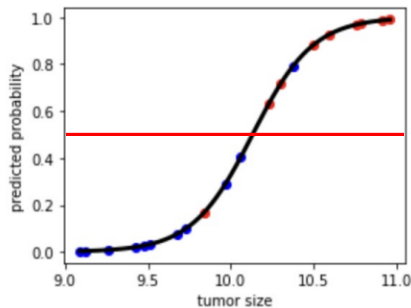
- (b) What if we want to extend the classification task to classify other animals? Can we train a new model while keeping the weights of the previous models?
- ▶ For an animal that are very distinct with the three animals, we can create a new logistic regression model without changing the previous weights.
 - ▶ For classifying a new animal that is similar with one of the classes (e.g, classifying a dog), we need to retrain the old models.

Q3. Precision, recall, F1 score and ROC curve

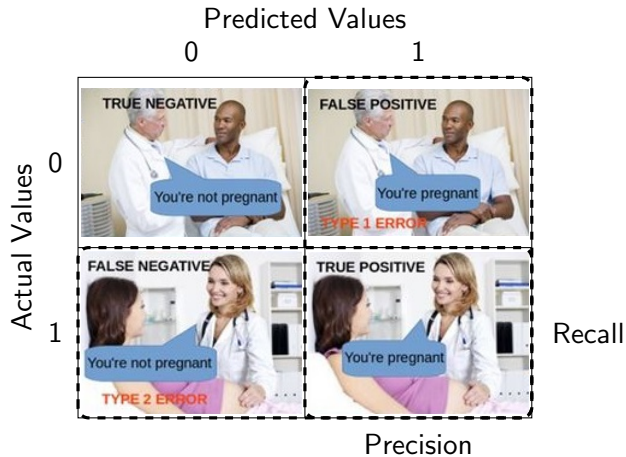
- (a) For the threshold $p = 0.5$, come up with the confusion matrix.

		Predicted Values	
		0	1
Actual Values	0	10	1
	1	1	8

Model M outputs label 1 if $M(x)$ is greater than or equal to the threshold, otherwise the model outputs 0.



Q3. Precision, recall, F1 score and ROC curve



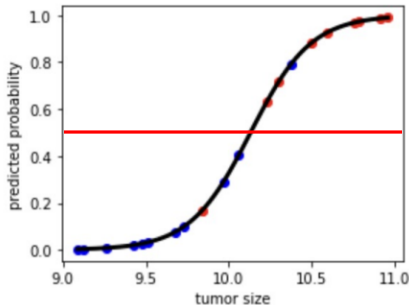
Q3. Precision, recall, F1 score and ROC curve

- (b) For the threshold $p = 0.5$, find the precision, recall and F1 score.

		Predicted Values	
		0	1
Actual Values	0	10	1
	1	1	8

$$\begin{aligned}\text{Precision} &= \frac{TP}{TP + FP} = \frac{8}{8 + 1} = \frac{8}{9} \\ \text{Recall} &= \frac{TP}{TP + FN} = \frac{8}{8 + 1} = \frac{8}{9} \\ \text{F1 Score} &= \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2}{\frac{8}{9} + \frac{8}{9}} = \frac{8}{9}\end{aligned}$$

Model M outputs label 1 if $M(x)$ is greater than or equal to the threshold, otherwise the model outputs 0.



Q3. Precision, recall, F1 score and ROC curve

- (d) In this question's case for detecting tumours, should we maximize precision or recall? Explain the reason for your choice.
- ▶ If cancer detection is being performed as a regular check up, then precision should be maximized; as we do not want to start cancer treatment on a person unless we are sure that he has cancer.
 - ▶ If cancer detection is being performed as part of cancer treatment progress monitoring, then recall should be maximized; as we do not want to stop the ongoing treatment unless we are sure that there is no residual tumour cell left in the patient.

Q5. Evaluating Logistic Regression

📌 Which of the following evaluation metrics is the least appropriate when comparing a logistic regression model's output with the target label? Explain your answer.

- (a) Accuracy
- (b) Precision, Recall
- (c) Binary Cross Entropy Loss
- (d) Mean Squared Error (MSE)

Q5. Evaluating Logistic Regression

Evaluation Metric:

- ▶ Judges the performance, doesn't care about the process.

[MRQ] Which of the following link(s) are pruned? Shade all that is/are true.

- | | | | | |
|-------------------------|-------------------------|------------------------------------|-------------------------|------------------------------------|
| <input type="radio"/> a | <input type="radio"/> b | <input type="radio"/> c | <input type="radio"/> d | <input type="radio"/> e |
| <input type="radio"/> f | <input type="radio"/> g | <input type="radio"/> h | <input type="radio"/> i | <input type="radio"/> j |
| <input type="radio"/> k | <input type="radio"/> l | <input checked="" type="radio"/> m | <input type="radio"/> n | <input type="radio"/> o |
| <input type="radio"/> p | <input type="radio"/> q | <input checked="" type="radio"/> r | <input type="radio"/> s | <input checked="" type="radio"/> t |
| <input type="radio"/> u | <input type="radio"/> v | <input type="radio"/> w | <input type="radio"/> x | <input type="radio"/> y |
| <input type="radio"/> z | | | | |
- 0 / 4

Loss Function:

- ▶ Helps with model training.
Minimized by the optimizer.

[MRQ] Which of the following link(s) are pruned? Shade all that is/are true.

- | | | | | |
|-------------------------|-------------------------|------------------------------------|-------------------------|------------------------------------|
| <input type="radio"/> a | <input type="radio"/> b | <input type="radio"/> c | <input type="radio"/> d | <input type="radio"/> e |
| <input type="radio"/> f | <input type="radio"/> g | <input type="radio"/> h | <input type="radio"/> i | <input type="radio"/> j |
| <input type="radio"/> k | <input type="radio"/> l | <input checked="" type="radio"/> m | <input type="radio"/> n | <input type="radio"/> o |
| <input type="radio"/> p | <input type="radio"/> q | <input checked="" type="radio"/> r | <input type="radio"/> s | <input checked="" type="radio"/> t |
| <input type="radio"/> u | <input type="radio"/> v | <input type="radio"/> w | <input type="radio"/> x | <input type="radio"/> y |
| <input type="radio"/> z | | | | |
- 25/26 items right 🍀 Try again!

🔍 Which options are suitable evaluation metrics?

Q5. Evaluating Logistic Regression

Evaluation metrics:

- (a) Accuracy
Gauging a model's overall performance.
- (b) Precision, Recall
Quantifies the model's ability to distinguish between positive and negative classes effectively.

Loss functions:

- (c) Binary Cross Entropy Loss
More suitable for **classification** tasks (assumes the binomial distribution).
- (d) Mean Squared Error
More suitable for **regression** tasks (assumes the normal distribution).

Answer: (b) > (a) > (c) > (d).

Bonus. Sigmoid vs Softmax

We use the sigmoid function for logistic regression in lecture:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

In multi-class logistic regression, we often use the softmax function instead:

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

- (a) Show that sigmoid function is a special case of the softmax function.
- (b) Express the derivative $\text{softmax}'(\mathbf{z})_{ik}$ in terms of $\text{softmax}(\mathbf{z})$.
- (c) Under what scenarios would you consider using softmax function instead of the sigmoid function?

Bonus. Sigmoid vs Softmax

Extra Slide

Solution.

(a) When $K = 2$,

$$\begin{aligned}\text{softmax}(\mathbf{z}) &= \begin{bmatrix} \frac{e^{z_1}}{e^{z_1} + e^{z_2}} & \frac{e^{z_2}}{e^{z_1} + e^{z_2}} \end{bmatrix}^T \\ &= \begin{bmatrix} \frac{e^{z_1 - z_2}}{e^{z_1 - z_2} + 1} & \frac{e^{z_2 - z_1}}{1 + e^{z_2 - z_1}} \end{bmatrix}^T \\ &= [\sigma(z_1 - z_2) \quad \sigma(z_2 - z_1)]^T\end{aligned}$$

which can be replaced by logistic regression where $z = z_1 - z_2$, predicting the probability of class 1.

```
1 def mysoftmax(z):  
2     softmax_class0 = torch.sigmoid(z[:, 0:1] - z[:, 1:2])  
3     return torch.hstack([softmax_class0, 1 - softmax_class0])
```


Bonus. Sigmoid vs Softmax

(b) Define δ_{ik} as 1 if $i = k$, 0 otherwise.

$$\begin{aligned}\text{softmax}'(\mathbf{z})_{ik} &= \frac{\delta_{ik} e^{z_i} \left(\sum_{j=1}^K e^{z_j} \right) - e^{z_i} e^{z_k}}{\left(\sum_{j=1}^K e^{z_j} \right)^2} \\ &= \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \cdot \left(\delta_{ik} - \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}} \right) \\ &= \text{softmax}(\mathbf{z})_i \cdot (\delta_{ik} - \text{softmax}(\mathbf{z})_k)\end{aligned}$$

Notice the similarity with the sigmoid function!

```
1 def mysoftmax_grad(z):
2     n, m = z.shape
3     z = torch.softmax(z, dim=1)
4     return z.reshape(n, m, 1)
5         * (torch.eye(m).reshape(1, m, m) - z.reshape(n, 1, m))
```

Bonus. Sigmoid vs Softmax

Extra Slide

- (c) The softmax function ensures that all output probabilities sum up to 1. It is a good idea to use the softmax function if the classes are mutually exclusive. On the other hand, use the sigmoid function if the classes are independent events.