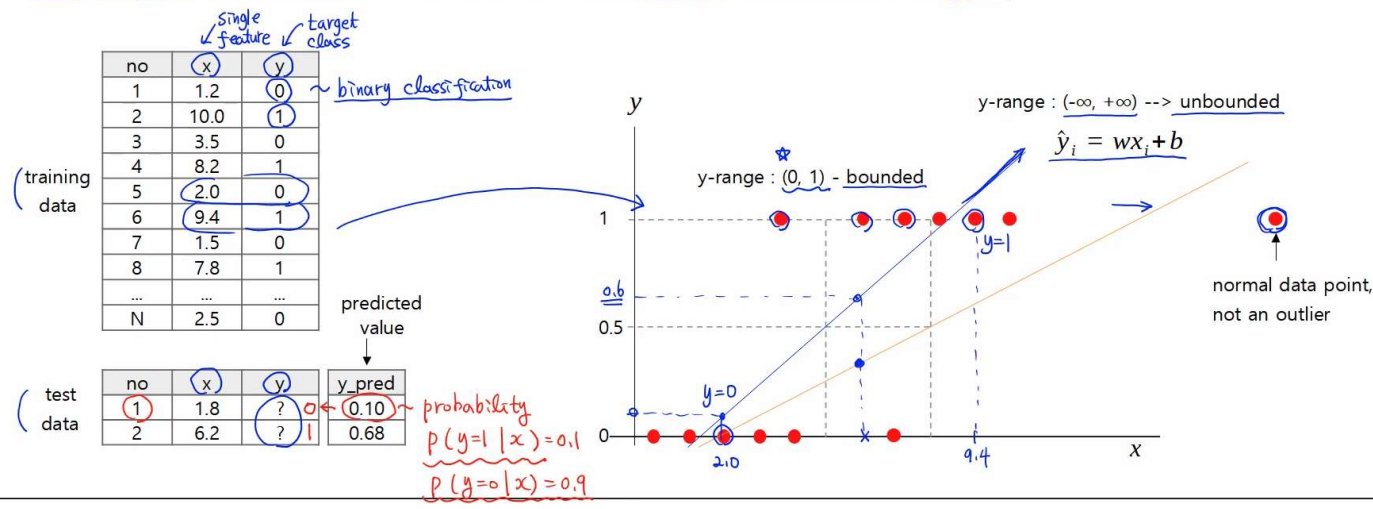
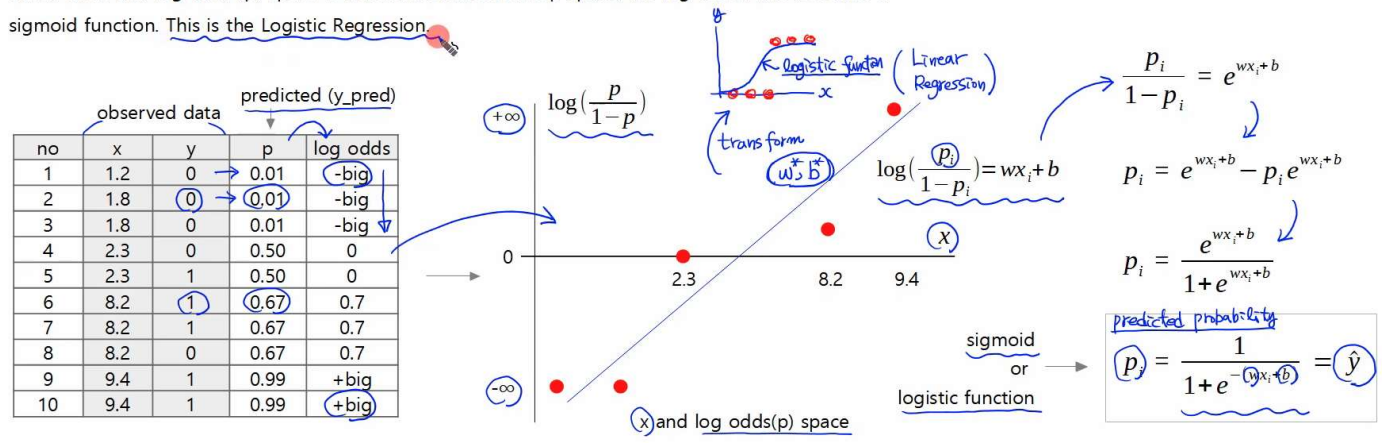
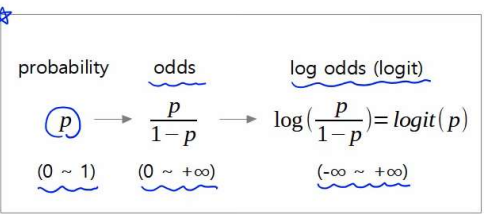


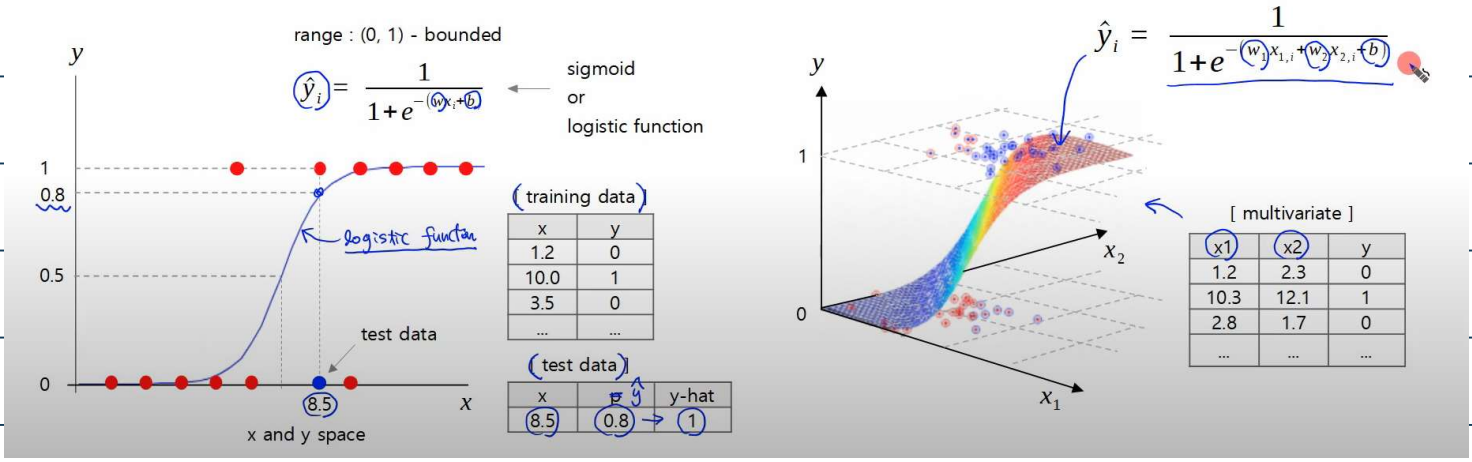
Can we use a linear regression model to predict the y values for following data? The answer is no



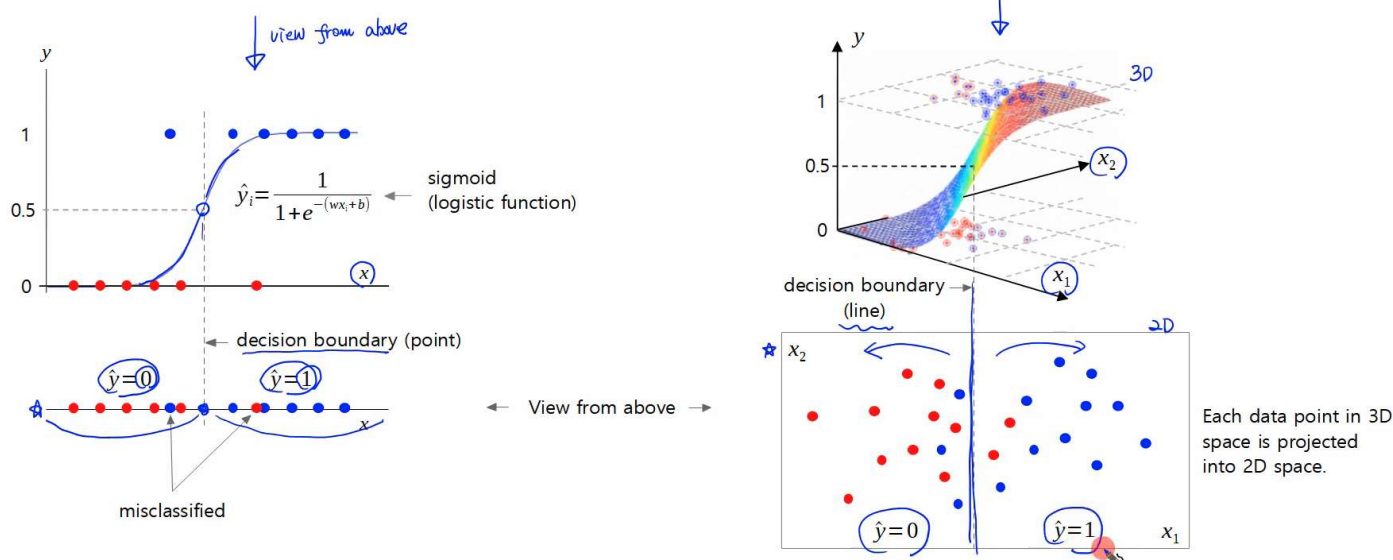
- The y value is bounded to be 0 or 1, and the probability p is bounded to be between 0 and 1.
- Linear regression can be applied to the values unbounded from $-\infty$ to $+\infty$.
- Therefore, it is not appropriate to apply linear regression to the data below.
- However, linear regression can be applied by converting the range of the predicted value to $-\infty \sim +\infty$.
- If the predicted value is expressed in log odds (logit), the range can be converted from $-\infty$ to $+\infty$.
- Linear regression can be applied in x and log odds (p) space instead of x and p space.
- When the x and log odds (p) space is transformed into x and p space, the regression line becomes a sigmoid function. This is the Logistic Regression.



- Binary classification can be performed by fitting the logistic (or sigmoid) function to the training data in x and y space. w and b can be estimated.
- Use the estimated logistic function to predict the class y of the test data.



- If the sigmoid value is greater than 0.5, the y class is predicted to be 1, otherwise it is predicted to be 0.
- The point where the sigmoid value is 0.5 is called the decision boundary.



- In linear regression, maximum likelihood estimation (MLE) was used to generate the objective function. In the same way, Logistic Regression can also use MLE to create an objective function that minimizes binary cross entropy.
- You can add a regularization term to the objective function in the same way as linear regression.

	y	\hat{y}	$P(y x) = \hat{y}^y (1-\hat{y})^{1-y}$
True	1	0.8	$\rightarrow 0.8^1 \times 0.2^0 = 0.8$
False	1	0.2	$\rightarrow 0.2^1 \times 0.8^0 = 0.2$
True	0	0.2	$\rightarrow 0.2^0 \times 0.8^1 = 0.8$

The larger this value, the better the prediction.

Objective function (binary cross entropy)

$$\max_{w,b} \sum_i \log(L(w,b))$$

$$= \min_{w,b} \sum_i [y_i \cdot \log \hat{y}_i + (1-y_i) \cdot \log (1-\hat{y}_i)]$$

$$= \min_{w,b} \sum_i BCE$$

Regularization

Lasso

$$\min_{w,b} \sum_i BCE + \lambda \sum_{j=0}^k |w_j|$$

Ridge

$$\min_{w,b} \sum_i BCE + \lambda \sum_{j=0}^k w_j^2$$

The large the $P(y|x)$, the better the prediction performance.

$$L(w,b) = \prod_{i=1}^N P(y_i | x_i; w, b)$$

$$\log(L(w,b)) = \sum_{i=1}^N \log(\hat{y}_i^{y_i} (1-\hat{y}_i)^{1-y_i})$$

$$\log(L(w,b)) = \sum_{i=1}^N [y_i \cdot \log \hat{y}_i + (1-y_i) \cdot \log (1-\hat{y}_i)]$$

view from above

sigmoid (logistic function)

decision boundary (point)

misclassified

View from above

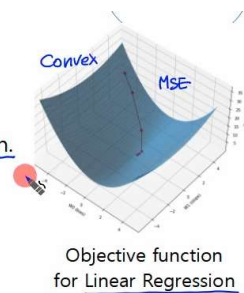
3D

2D

Each data point in 3D space is projected into 2D space.

Convexity of the objective function

- Binary cross entropy, the objective function of logistic regression, is a convex function and has a unique solution.
- Since the second derivative of binary cross entropy is always greater than 0, we can see that it is a convex function.



$$J(w,b) = \sum_{i=1}^N [-y_i \cdot \log \hat{y}_i - (1-y_i) \cdot \log (1-\hat{y}_i)]$$

$$\hat{y}_i(x|w,b) = \frac{1}{1+e^{-(wx+b)}}$$

$$\hat{y}_i = \frac{1}{1+e^{-t}}, \quad (t = wx_i + b)$$

$$J_0 = -y \cdot \log \hat{y}_t - (1-y) \cdot \log (1-\hat{y}_t)$$

→ binary cross entropy

$$\frac{d}{dt} \hat{y}_t = \frac{e^{-t}}{(1+e^{-t})^2} = \frac{1+e^{-t}}{(1+e^{-t})^2} - \frac{1}{(1+e^{-t})^2}$$

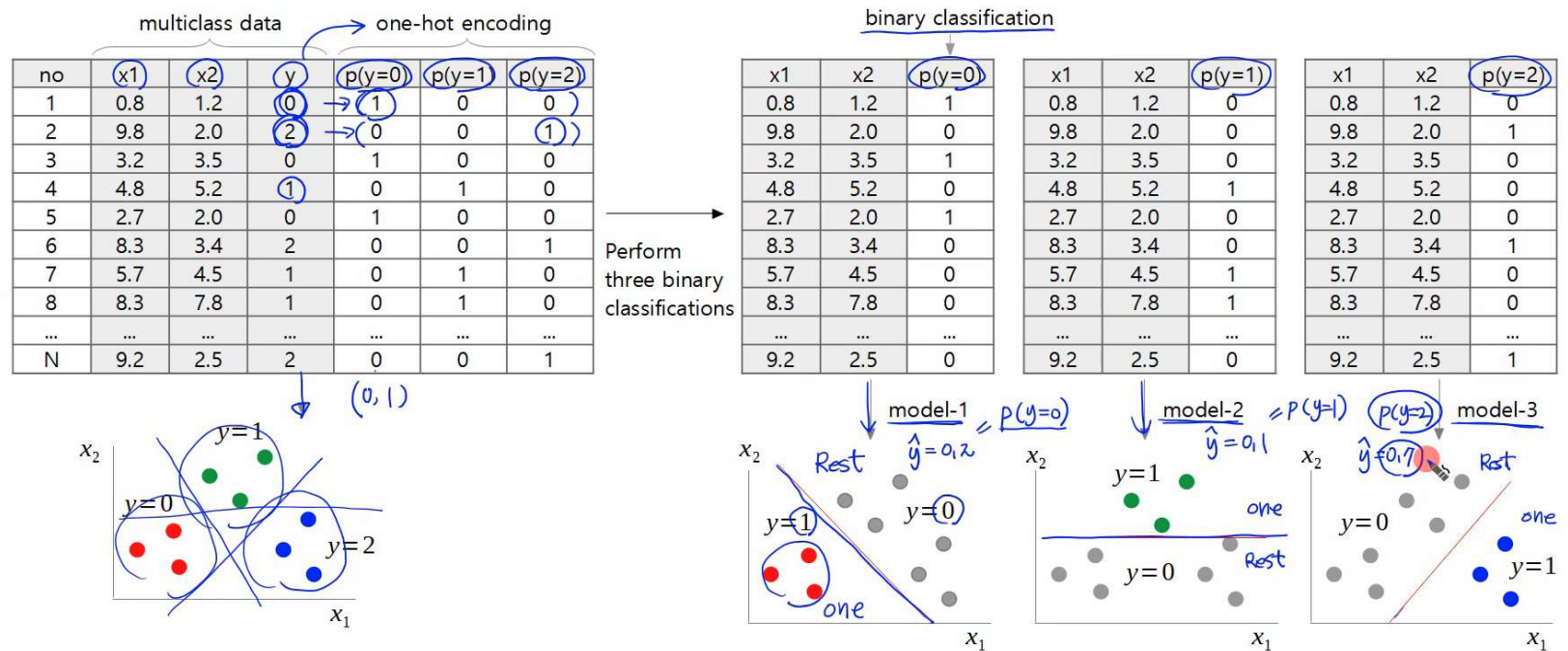
$$= \frac{1}{1+e^{-t}} \left(1 - \frac{1}{1+e^{-t}} \right) = \hat{y}_t (1-\hat{y}_t)$$

$$\frac{d}{dt} J_0 = -\frac{y}{\hat{y}_t} \hat{y}_t (1-\hat{y}_t) + (1-y) \frac{\hat{y}_t (1-\hat{y}_t)}{1-\hat{y}_t} = -y + \hat{y}_t$$

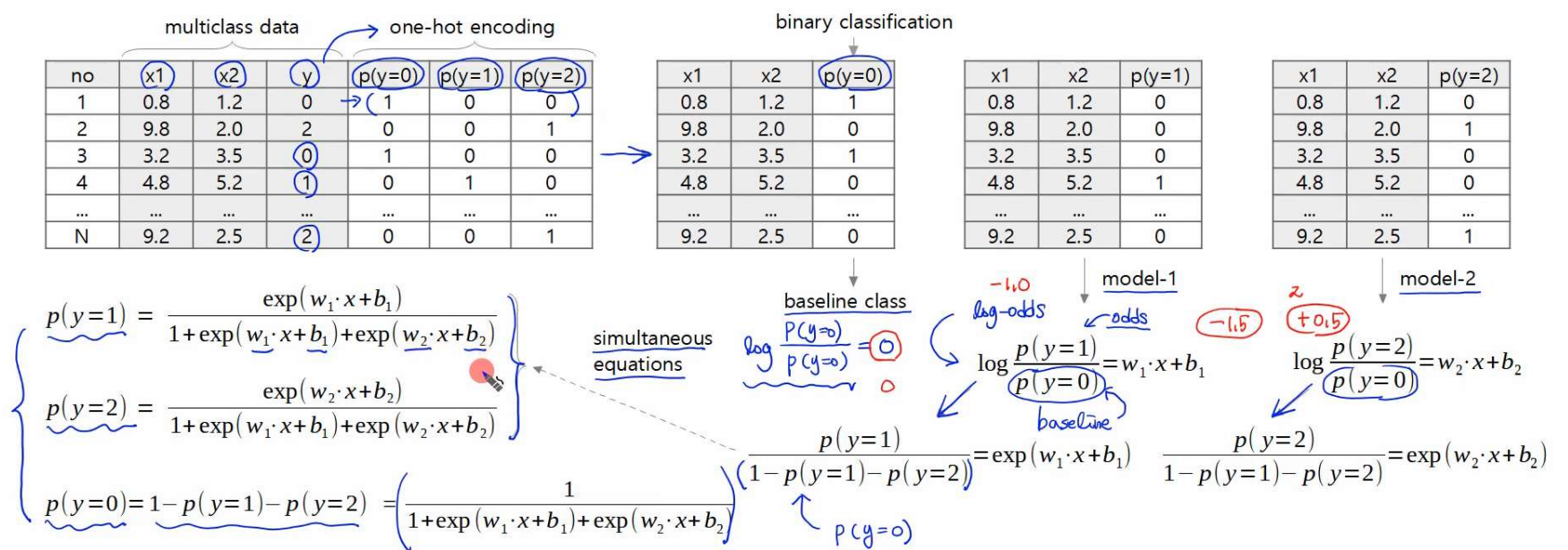
$$\frac{d^2}{dt^2} J_0 = \hat{y}_t (1-\hat{y}_t) = \frac{1}{1+e^{-t}} - \frac{1}{(1+e^{-t})^2} = \frac{e^{-t}}{(1+e^{-t})^2} > 0$$

■ Multiclass classification – OvR (One-vs-Rest) = OVA (one-vs-All) (Multinomial or Softmax Logistic Regression)

- Data with three classes ($y = [0, 1, 2]$) is one-hot encoded and then trained through three binary classification models (model-1, 2, 3).
- After training is complete, the test data is input into the three models to predict the probability that y is 0, 1, and 2, respectively, and the class with the highest probability is selected as the final class.



- Let's assume that the data with three classes ($y=0,1,2$) is trained with two binary classification models (model-1, 2) as shown below.
- Set a baseline class and perform binary classification on the remaining classes. Any class can be a baseline.
- While OvR actually requires training multiple binary models, this method estimates all parameters simultaneously with single model.



■ Multinomial Logistic Regression (or Softmax Regression)

- Softmax function

Binary classification

$$\begin{cases} p(y=1|x) = \frac{1}{1 + \exp(wx+b)} : \text{Logistic function} \\ p(y=0|x) = 1 - p(y=1|x) \end{cases}$$

$w_0=0, b_0=0$

$$p(y=1) = \frac{\exp(w_1 \cdot x + b_1)}{1 + \exp(w_1 \cdot x + b_1) + \exp(w_2 \cdot x + b_2)} \rightarrow \frac{\exp(w_1 \cdot x + b_1)}{\exp(w_0 \cdot x + b_0) + \exp(w_1 \cdot x + b_1) + \exp(w_2 \cdot x + b_2)} = \frac{\exp(w_1 \cdot x + b_1)}{\sum_{k=0}^{C-1} \exp(w_k \cdot x + b_k)}$$

softmax

(C : the number of classes)

$$p(y=2) = \frac{\exp(w_2 \cdot x + b_2)}{1 + \exp(w_1 \cdot x + b_1) + \exp(w_2 \cdot x + b_2)} = \frac{\exp(w_2 \cdot x + b_2)}{\sum_{k=0}^{C-1} \exp(w_k \cdot x + b_k)}$$

$$p(y=0) = \frac{1}{1 + \exp(w_1 \cdot x + b_1) + \exp(w_2 \cdot x + b_2)} = \frac{\exp(w_0 \cdot x + b_0)}{\sum_{k=0}^{C-1} \exp(w_k \cdot x + b_k)}$$

■ Loss function for Softmax Regression

Binary classification → Logistic function → MLE → binary cross entropy

Loss function for Softmax Regression

- Just like we did for binary classification, we can use MLE to obtain the loss function for softmax regression. The loss function is cross entropy, which is a general expression for binary cross entropy.
- Regularization can also be applied to the cross entropy loss function.

$y \rightarrow \hat{y}$ probability $P(y|x) = \prod_{k=0}^{C-1} \hat{y}_{i,k}^{1(y_i=k)}$

	k=0	1	2	
True 1	0.1	0.7	0.2	$\rightarrow 0.1^0 \times 0.7^1 \times 0.2^0 = 0.7$
False 1	0.8	0.1	0.1	$\rightarrow 0.8^0 \times 0.1^1 \times 0.1^0 = 0.1$
True 0	0.8	0.1	0.1	$\rightarrow 0.8^1 \times 0.1^0 \times 0.1^0 = 0.8$

The larger this value, the better the prediction.

(Softmax)

$$\hat{y}_{i,k} = \frac{\exp(w_k \cdot x_i + b_k)}{\sum_{k=0}^{C-1} \exp(w_k \cdot x_i + b_k)}$$

i : data index, k : class index
C: the number of classes

$$L_i(w, b) = \prod_{k=0}^{C-1} \hat{y}_{i,k}^{1(y_i=k)}$$

$1(y_i=k)$: indicator function

$$\log(L_i(w, b)) = \sum_{k=0}^{C-1} \log(\hat{y}_{i,k}^{1(y_i=k)})$$

$$= \sum_{k=0}^{C-1} 1(y_i=k) \cdot \log(\hat{y}_{i,k})$$

$1(y_i=k) \rightarrow y_{i,k}$

$$J(w, b) = \sum_{i=0}^{N-1} \left[\sum_{k=0}^{C-1} y_{i,k} \cdot \log(\hat{y}_{i,k}) \right]$$

N : the number of data points
A general expression for binary cross entropy.
C = 2 → binary cross entropy.

• Objective function

$$\begin{aligned} \max_{w,b} J(w, b) \\ &= \min_{w,b} \sum_{i=0}^{N-1} \sum_{k=0}^{C-1} y_{i,k} \cdot \log(\hat{y}_{i,k}) \\ &= \min_{w,b} \sum_i CE_i \end{aligned}$$

CE: cross entropy
softmax

• Regularization

$$\min_{w,b} \sum_i CE_i + \lambda \sum_{k=0}^{C-1} |w_k| \quad (\text{Lasso})$$

$$\min_{w,b} \sum_i CE_i + \lambda \sum_{k=0}^{C-1} w_k^2 \quad (\text{Ridge})$$

• comparison: binary classification

$$\hat{y}_i = \frac{1}{1 + \exp(-(w x_i + b))}$$

$$P(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$$

indicator function

Let $\begin{cases} \hat{y}_{k=0} = 1 - \hat{y} \\ \hat{y}_{k=1} = \hat{y} \end{cases}$