

Ch02. Logistic Regression

Hwanjo Yu

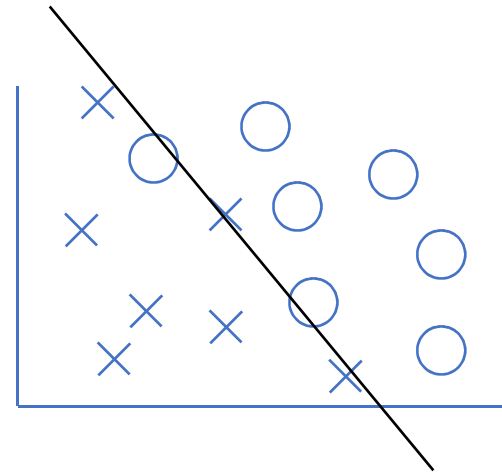
POSTECH

<http://hwanjoyu.org>

Binary Classification

Examples

- Medical diagnosis: Disease or no disease
- Cancer diagnosis: Benign or malignant
- Image classification: Cat or dog
- Image classification: Face or non-face
- Sentiment analysis: Positive or negative
- Name entity recognition: Name or not
- Speech or audio: Voice or not

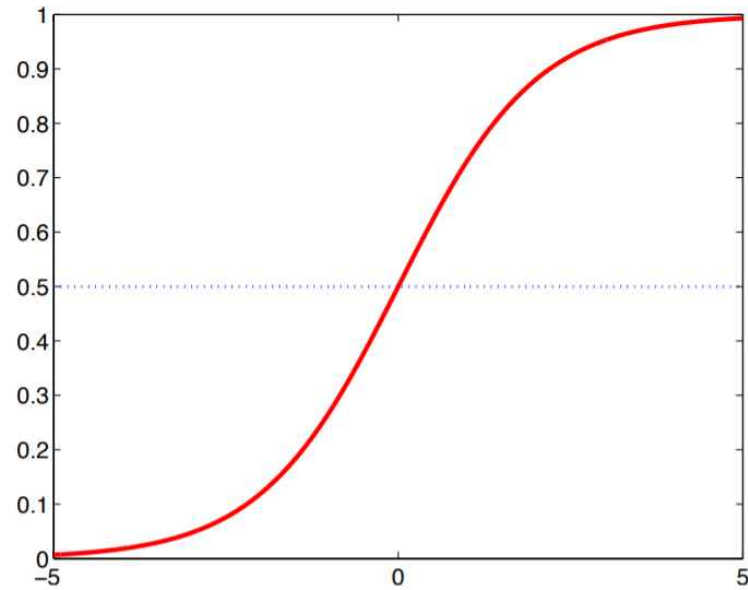


Logistic Function

Logistic function (or sigmoid function) $\sigma(\xi)$:

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

- $\sigma(\xi) \rightarrow 0$ as $\xi \rightarrow -\infty$.
- $\sigma(\xi) \rightarrow 1$ as $\xi \rightarrow \infty$.
- $\sigma(-\xi) = 1 - \sigma(\xi)$.
- $\frac{d}{d\xi} [\sigma(\xi)] = \sigma(\xi)\sigma(-\xi)$.



Maximum Likelihood Formulation

Logistic Regression

Predict a binary output $y_n \in \{0,1\}$ from an input \mathbf{x}_n .

The logistic regression models the input-output by a **conditional Bernoulli distribution**

$$\mathbb{E}[y_n|\mathbf{x}_n] = \mathbb{P}(y_n = 1|\mathbf{x}_n) = \sigma(\mathbf{w}^T \mathbf{x}_n),$$

where

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

Discriminative model: Directly model $p(y|\mathbf{x})$.

Logistic Regression: MLE

Given $\{(\mathbf{x}_n, y_n) | n = 1, \dots, N\}$, the likelihood is given by

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{w}) &= \prod_{n=1}^N p(y_n = 1|\mathbf{x}_n)^{y_n} (1 - p(y_n = 1|\mathbf{x}_n))^{1-y_n} \\ &= \prod_{n=1}^N \sigma(\mathbf{w}^T \mathbf{x}_n)^{y_n} (1 - \sigma(\mathbf{w}^T \mathbf{x}_n))^{1-y_n} \end{aligned}$$

Then, log-likelihood function is given by

$$\mathcal{L} = \sum_{n=1}^N \log p(y_n|\mathbf{x}_n) = \sum_{n=1}^N \{y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n)\},$$

where $\hat{y}_n = \sigma(\mathbf{w}^T \mathbf{x}_n)$.

- This is a nonlinear function of \mathbf{w} whose maximum cannot be computed in a closed form.
- **Iterative re-weighted least squares (IRLS)** is a popular algorithm, derived from Newton's method.

Newton's Method

Mathematical Preliminaries

- Gradient
- Hessian matrix
- Gradient descent/ascent
- Newton's method

Gradient

- Consider a real-valued function $f(\mathbf{x})$, that takes a real-valued vector $\mathbf{x} \in \mathbb{R}^D$ as an input,

$$f(\mathbf{x}): \mathbb{R}^D \rightarrow \mathbb{R}.$$

- The gradient of $f(\mathbf{x})$ is defined by

$$\nabla f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right]^T.$$

Hessian Matrix

The Hessian matrix \mathbf{H} is defined as the symmetric matrix with the (i, j) -element $\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}$,

$$\begin{aligned}\mathbf{H} = \nabla^2 f(\mathbf{x}) &= \left[\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} \right] = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_D} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_D \partial x_1} & \frac{\partial^2 f}{\partial x_D \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_D^2} \end{bmatrix} = \frac{\partial}{\partial \mathbf{x}} \left[\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right]^T \\ &= \frac{\partial}{\partial \mathbf{x}} [\nabla f(\mathbf{x})]^T\end{aligned}$$

Gradient Descent/Ascent

- The gradient descent/ascent learning is a simple (first-order) iterative method for minimization/maximization

- **Gradient descent**: iterative minimization

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \eta \left(\frac{\partial \mathcal{J}}{\partial \mathbf{w}} \right).$$

- **Gradient ascent**: iterative maximization

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \eta \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right).$$

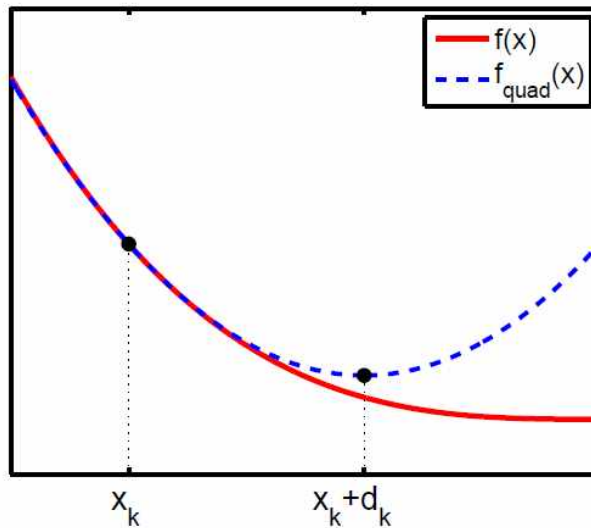
- **Learning rate**: $\eta > 0$

Taylor series expansion

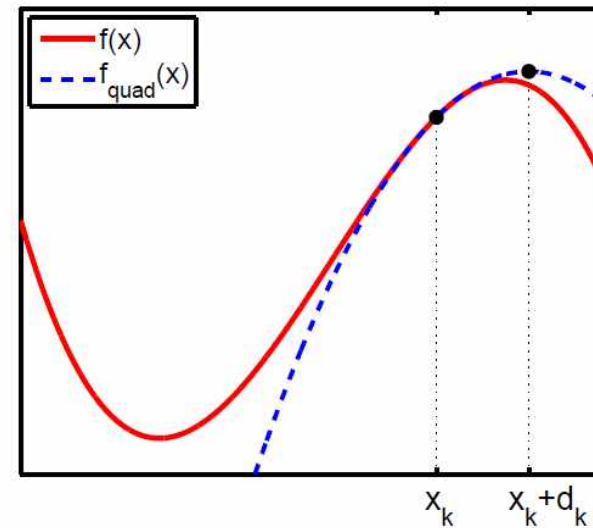
- Approximate an arbitrary function with a polynomial form.
- $f(x) = C_0 + C_1(x - a) + C_2(x - a)^2 + C_3(x - a)^3 + C_4(x - a)^4 + \dots$
 - $f(a) = C_0$
 - $f'(x) = C_1 + 2C_2(x - a) + 3C_3(x - a)^2 + 4C_4(x - a)^3 + \dots$
 - $f'(a) = C_1$
 - $f''(x) = 2 * 1C_2 + 3 * 2C_3(x - a)^1 + 4 * 3C_4(x - a)^2 + \dots$
 - $f''(a) = 2C_2$
 - $f'''(x) = 3 * 2 * 1C_3 + 4 * 3 * 2C_4(x - a)^1 + \dots$
 - $f'''(a) = 3 * 2 * 1C_3$
- $$f(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2 + \frac{f'''(a)}{3!}(x - a)^3 + \dots$$
$$= \sum_{n=1}^{\infty} \frac{1}{n!} f^{(n)}(a)(x - a)^n$$
- Must $f(x)$ be “multiply differentiable” at $x = a$ to approximate $f(x)$ using Taylor

Newton's method

The basic idea of Newton's method is to optimize the quadratic approximation of the objective function $\mathcal{J}(\mathbf{w})$ around the current point $\mathbf{w}^{(k)}$.



(a) convex



(b) non-convex

[from Murphy's book]

The second-order Taylor series expansion of $\mathcal{J}(\mathbf{w})$ at the current point $\mathbf{w}^{(k)}$ gives

$$\mathcal{J}_{\text{quad}}(\mathbf{w}) = \mathcal{J}(\mathbf{w}^{(k)}) + [\nabla \mathcal{J}(\mathbf{w}^{(k)})]^T (\mathbf{w} - \mathbf{w}^{(k)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(k)})^T \nabla^2 \mathcal{J}(\mathbf{w}^{(k)}) (\mathbf{w} - \mathbf{w}^{(k)}).$$

where $\nabla^2 \mathcal{J}(\mathbf{w}^{(k)})$ is the Hessian of $\mathcal{J}(\mathbf{w})$ evaluated at $\mathbf{w} = \mathbf{w}^{(k)}$.

Newton's method

Differentiate $\mathcal{J}_{\text{quad}}(\mathbf{w})$ w.r.t. \mathbf{w} and set it equal 0, which leads to

$$\nabla \mathcal{J}(\mathbf{w}^{(k)}) + \nabla^2 \mathcal{J}(\mathbf{w}^{(k)})\mathbf{w} - \nabla^2 \mathcal{J}(\mathbf{w}^{(k)})\mathbf{w}^{(k)} = 0.$$

Thus, we have

$$\mathbf{w} = \mathbf{w}^{(k)} - [\nabla^2 \mathcal{J}(\mathbf{w}^{(k)})]^{-1} \nabla \mathcal{J}(\mathbf{w}^{(k)}).$$

Hence the Newton's method is of the form

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - [\nabla^2 \mathcal{J}(\mathbf{w}^{(k)})]^{-1} \nabla \mathcal{J}(\mathbf{w}^{(k)}).$$

Logistic Regression: Algorithms

Logistic Regression: Gradient Ascent

The gradient ascent learning has the form

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \eta \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right),$$

where the log-likelihood

$$\mathcal{L} = \sum_{n=1}^N \log p(y_n | \mathbf{x}_n) = \sum_{n=1}^N \{y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n)\},$$

where $\hat{y}_n = \sigma(\mathbf{w}^T \mathbf{x}_n)$.

Then, the gradient

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= \sum_{n=1}^N \left[y_n \frac{\hat{y}_n'}{\hat{y}_n} \mathbf{x}_n + (1 - y_n) \frac{-\hat{y}_n'}{1 - \hat{y}_n} \mathbf{x}_n \right] = \sum_{n=1}^N \left[y_n \frac{\hat{y}_n(1 - \hat{y}_n)}{\hat{y}_n} - (1 - y_n) \frac{\hat{y}_n(1 - \hat{y}_n)}{1 - \hat{y}_n} \right] \mathbf{x}_n \\ &= \sum_{n=1}^N [y_n(1 - \hat{y}_n) - (1 - y_n)\hat{y}_n] \mathbf{x}_n = \sum_{n=1}^N [y_n - \hat{y}_n] \mathbf{x}_n. \end{aligned}$$

Hence, the gradient ascent update rule for \mathbf{w} is

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \eta \sum_{n=1}^N [y_n - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n.$$

Logistic Regression: Hessian

Calculate the Hessian:

$$\nabla^2 \mathcal{L} = \frac{\partial}{\partial \mathbf{w}} [\nabla \mathcal{L}]^T = \frac{\partial}{\partial \mathbf{w}} \left[\sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n^T \right] = \sum_{n=1}^N -\hat{y}_n (1 - \hat{y}_n) \mathbf{x}_n \mathbf{x}_n^T.$$

We set the objective function $J(\mathbf{w})$ as the negative log-likelihood (cross-entropy error):

$$J(\mathbf{w}) = -\mathcal{L}(\mathbf{w}) = -\sum_{n=1}^N \{y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n)\}.$$

Thus, the gradient and the Hessian are:

$$\begin{aligned} \nabla J(\mathbf{w}) &= -\sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n, \\ \nabla^2 J(\mathbf{w}) &= \sum_{n=1}^N \hat{y}_n (1 - \hat{y}_n) \mathbf{x}_n \mathbf{x}_n^T. \end{aligned}$$

Note that the Hessian matrix $\nabla^2 J(\mathbf{w}) = \sum_{n=1}^N \hat{y}_n (1 - \hat{y}_n) \mathbf{x}_n \mathbf{x}_n^T$ is positive definite since $\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$ is positive definite and $\hat{y}_n (1 - \hat{y}_n) \geq 0, \forall n$.

Logistic Regression: IRLS

Newton's update has the form

$$\Delta \mathbf{w} = - \underbrace{\left[\sum_{n=1}^N \hat{y}_n (1 - \hat{y}_n) \mathbf{x}_n \mathbf{x}_n^T \right]^{-1}}_{[\nabla^2 J(\mathbf{w})]^{-1}} \underbrace{\left[- \sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n^T \right]}_{\nabla J(\mathbf{w})},$$

Newton's update reduces to [iterative re-weighted least squares \(IRLS\)](#):

$$\Delta \mathbf{w} = (\mathbf{X} \mathbf{S} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{S} \mathbf{b},$$

where

$$\mathbf{S} = \begin{bmatrix} \hat{y}_1(1 - \hat{y}_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \hat{y}_N(1 - \hat{y}_N) \end{bmatrix},$$
$$\mathbf{b} = \begin{bmatrix} \frac{y_1 - \hat{y}_1}{\hat{y}_1(1 - \hat{y}_1)} \\ \vdots \\ \frac{y_N - \hat{y}_N}{\hat{y}_N(1 - \hat{y}_N)} \end{bmatrix}.$$

Logistic Regression: IRLS

Alternatively we write

$$\mathbf{w}_{k+1} = \mathbf{w}_k + (\mathbf{X}\mathbf{S}_k\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{S}_k\mathbf{b}_k,$$

Where \mathbf{S}_k and \mathbf{b}_k are computed using \mathbf{w}_k .

That is,

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k + (\mathbf{X}\mathbf{S}_k\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{S}_k\mathbf{b}_k = (\mathbf{X}\mathbf{S}_k\mathbf{X}^T)^{-1}[(\mathbf{X}\mathbf{S}_k\mathbf{X}^T)\mathbf{w}_k + \mathbf{X}\mathbf{S}_k\mathbf{b}_k] \\ &= (\mathbf{X}\mathbf{S}_k\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{S}_k[\mathbf{X}^T\mathbf{w}_k + \mathbf{b}_k].\end{aligned}$$

Algorithm: IRLS

Input: $\{(\mathbf{x}_n, y_n) | n = 1, \dots, N\}$

1. Initialize $\mathbf{w} = 0$ and $w_0 = \log(\bar{y}/(1 - \bar{y}))$
2. repeat
3. for $n = 1, \dots, N$ do
4. Compute $\eta_n = \mathbf{w}^T \mathbf{x}_n + w_0$
5. Compute $\hat{y}_n = \sigma(\eta_n)$
6. Compute $s_n = \hat{y}_n(1 - \hat{y}_n)$
7. Compute $z_n = \eta_n + (y_n - \hat{y}_n)/s_n$
8. end for
9. Construct $\mathbf{S} = \text{diag}(s_{1:N})$
10. Update $\mathbf{w} = (\mathbf{X}\mathbf{S}\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{S}\mathbf{z}$
11. until convergence
12. return \mathbf{w}

Softmax Regression for Multiclass Classification

- Model input-output by a softmax transformation of logits $\theta_k = \mathbf{w}_k^T \mathbf{x}_n$:

$$p(y_n = k | \mathbf{x}_n) = \text{softmax}(\theta_k = \mathbf{w}_k^T \mathbf{x}_n) = \frac{\exp\{\mathbf{w}_k^T \mathbf{x}_n\}}{\sum_{j=1}^K \exp\{\mathbf{w}_j^T \mathbf{x}_n\}}$$

- Given $\mathbf{Y} \in \mathbb{R}^{K \times N}$ (each column $\mathbf{y}_n \in \mathbb{R}^K$ follows the 1-of- K coding) and $\mathbf{X} \in \mathbb{R}^{D \times N}$, the likelihood is given by

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(y_n = k | \mathbf{x}_n)^{Y_{k,n}}.$$

- The log-likelihood is given by

$$\mathcal{L} = \sum_{n=1}^N \sum_{k=1}^K Y_{k,n} \log[p(y_n = k | \mathbf{x}_n)].$$

- One can apply Newton's update to derive IRLS, as in logistic regression.

Cross Entropy Error

- Cross entropy between two probability distributions p and q is $H(p; q) = -\sum_x p(x) \log q(x)$ (refer to Wikipedia)
- In the case of logistic regression, we have $p \in \{y, 1 - y\}$, $q \in \{\hat{y}, 1 - \hat{y}\}$. Then, the **cross entropy error** is calculated as

$$\mathcal{E} = \sum_{n=1}^N [-y_n \log \hat{y}_n - (1 - y_n) \log(1 - \hat{y}_n)].$$

- In the case of softmax regression, we have $p \in \{y_1, \dots, y_K\}$, $q \in \{\hat{y}_1, \dots, \hat{y}_K\}$. Then, the **cross entropy error** is calculated as

$$\mathcal{E} = \sum_{n=1}^N \sum_{k=1}^K [-y_{k,n} \log \hat{y}_{k,n}].$$

- Note, the cross entropy error \mathcal{E} is equal to the negative log-likelihood $-\mathcal{L}(\mathbf{w})$.

Cross Entropy Error

- $y = \begin{bmatrix} \text{tiger} \\ \text{lion} \\ \text{cat} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \hat{y} = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}, \mathcal{E} = -\log 0.7 = 0.36$
- $y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \hat{y} = \begin{bmatrix} 0.5 \\ 0.3 \\ 0.2 \end{bmatrix}, \mathcal{E} = -\log 0.5 = 0.69$

Multi-Label Learning

- $y = \begin{bmatrix} \text{dog} \\ \text{cat} \\ \text{sky} \\ \text{sand} \\ \text{lake} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$



Labels: dog, sand, sky

- The error function is given by

$$\begin{aligned} \mathcal{J} &= \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \underbrace{\ell(y_{k,n}, \hat{y}_{k,n})}_{\text{bss}} \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K [-y_{k,n} \log \hat{y}_{k,n} - (1 - y_{k,n}) \log(1 - \hat{y}_{k,n})] , \end{aligned}$$

- where $\hat{y}_{k,n}$ is the output of the k th logistic regression model,

$$\hat{y}_{k,n} = \sigma(\mathbf{w}_k^T \mathbf{x}_n)$$

Multi-Label Learning

- For single-label learning, $y = \begin{bmatrix} \text{tiger} \\ \text{lion} \\ \text{cat} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$, $\hat{y} = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$, $\mathcal{E} = -\log 0.7 = 0.36$
- For multi-label learning, $y = \begin{bmatrix} \text{tiger} \\ \text{lion} \\ \text{cat} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$, $\hat{y} = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$, $\mathcal{E} = -\log 0.7 - \log 0.8 - \log 0.9 = 0.69$
- When having '?' labels such as $[\mathbf{Y} \in \mathbb{R}^{K \times N}] = \begin{bmatrix} 1 & 0 & ? & \dots \\ 0 & 0 & 1 & \dots \\ 1 & ? & 0 & \dots \\ 1 & 1 & ? & \dots \\ 0 & 1 & ? & \dots \end{bmatrix}$,
- k runs only for indices associated with 0 or 1 in $\mathcal{J} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \underbrace{\ell(y_{k,n})}_{\text{bgt}} \underbrace{\hat{y}_{k,n}}_{\text{bss}}.$