

Look at one class at a time instead of 2

Discriminative: Learn $P(y|x)$ (or learn $\theta_{01} = \{0\}$)

Generative: build model of each first: Learn $P(x|y)$ & $P(y)$

$$P(y=1|x) = \frac{P(x|y=1)P(y=1)}{P(x)}$$

$$P(x) = P(x|y=1)P(y=1) + P(x|y=0)P(y=0)$$

GDA: Gaussian Discriminative Analysis

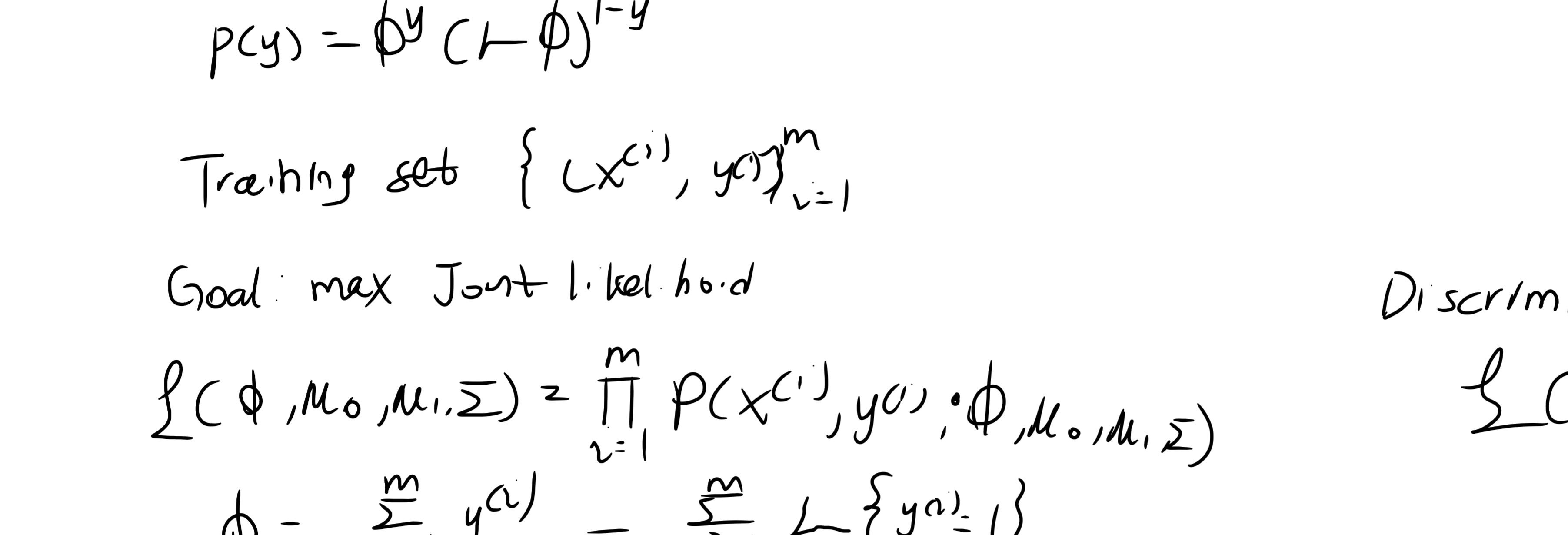
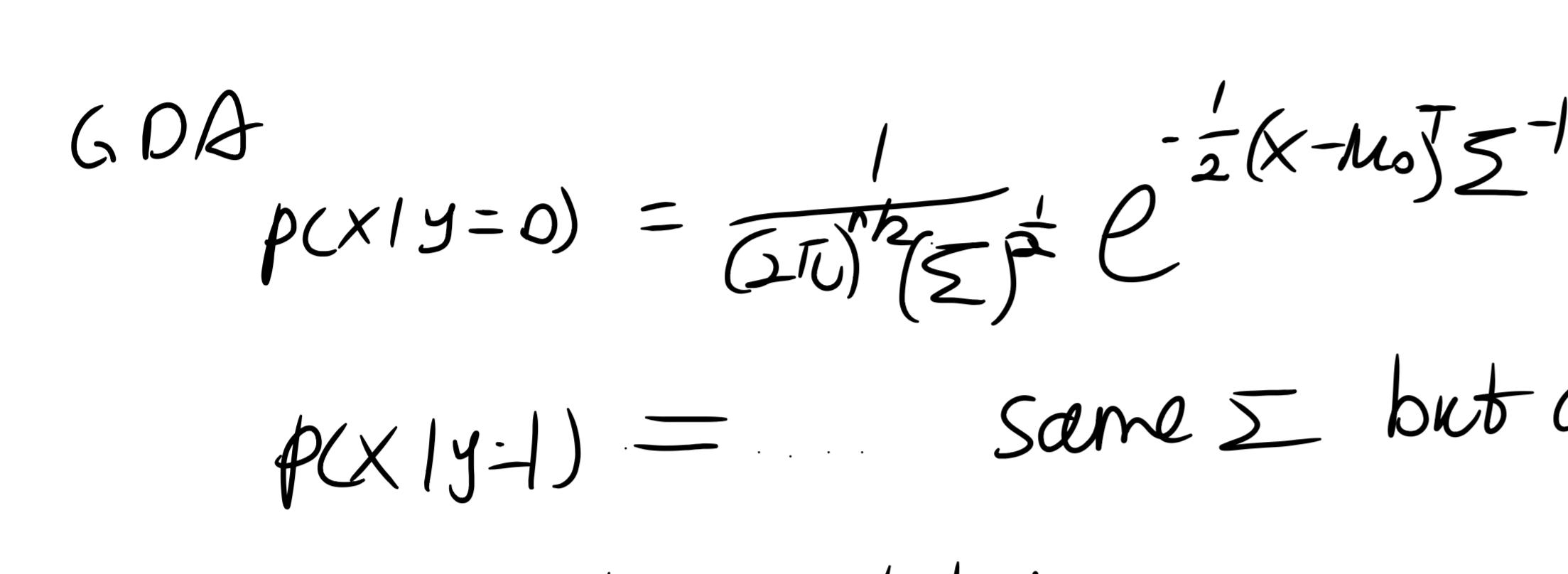
Suppose x continuous $x \in \mathbb{R}^n$ (instead of $n+1$)

Assume $P(x|y)$ is Gaussian

Multi-variate Gaussian $Z \sim N(\bar{\mu}, \Sigma)$

$$\begin{aligned} E[Z] &= \bar{\mu}, \quad \text{Cov}[Z] = E[(Z - \bar{\mu})(Z - \bar{\mu})^T] \\ &= E[ZZ^T] - (E[Z])(E[Z]^T) \end{aligned}$$

$$P(Z) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\bar{\mu})^T \Sigma^{-1} (x-\bar{\mu})}$$



$$GDA \quad P(x|y=0) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0)}$$

$$P(x|y=1) = \dots \quad \text{same } \Sigma \text{ but different } \mu$$

$$P(y) = \phi^y (1-\phi)^{1-y}$$

$$\text{Training set } \{(x^{(i)}, y^{(i)})\}_{i=1}^m$$

Goal: max joint likelihood

Discriminative

$$\mathcal{L}(\phi, \mu_0, \mu_1, \Sigma) = \prod_{i=1}^m P(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)$$

$$\mathcal{L}(\theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta)$$

$$\phi = \frac{\sum_{i=1}^m y^{(i)}}{m} = \frac{1}{m} \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\}$$

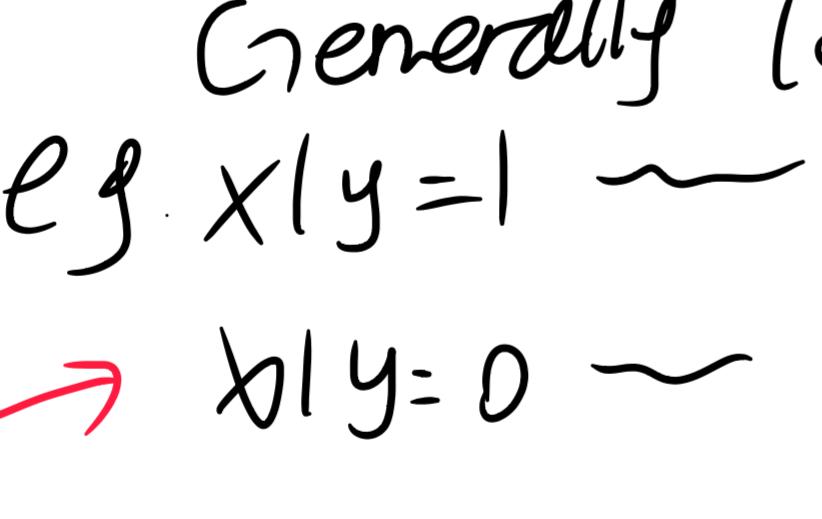
$$\mu_0 = \frac{\sum_{i=1}^m \mathbb{I}\{y^{(i)}=1 | x^{(i)}\}}{\sum_{i=1}^m \mathbb{I}\{y^{(i)}=0\}} \leftarrow \begin{array}{l} \text{sum of feature vectors} \\ \text{for } y=0 \end{array}$$

$$\mu_1 = \frac{\sum_{i=1}^m \mathbb{I}\{y^{(i)}=1 | x^{(i)}\}}{\sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\}} \leftarrow \# \text{ with } y=1$$

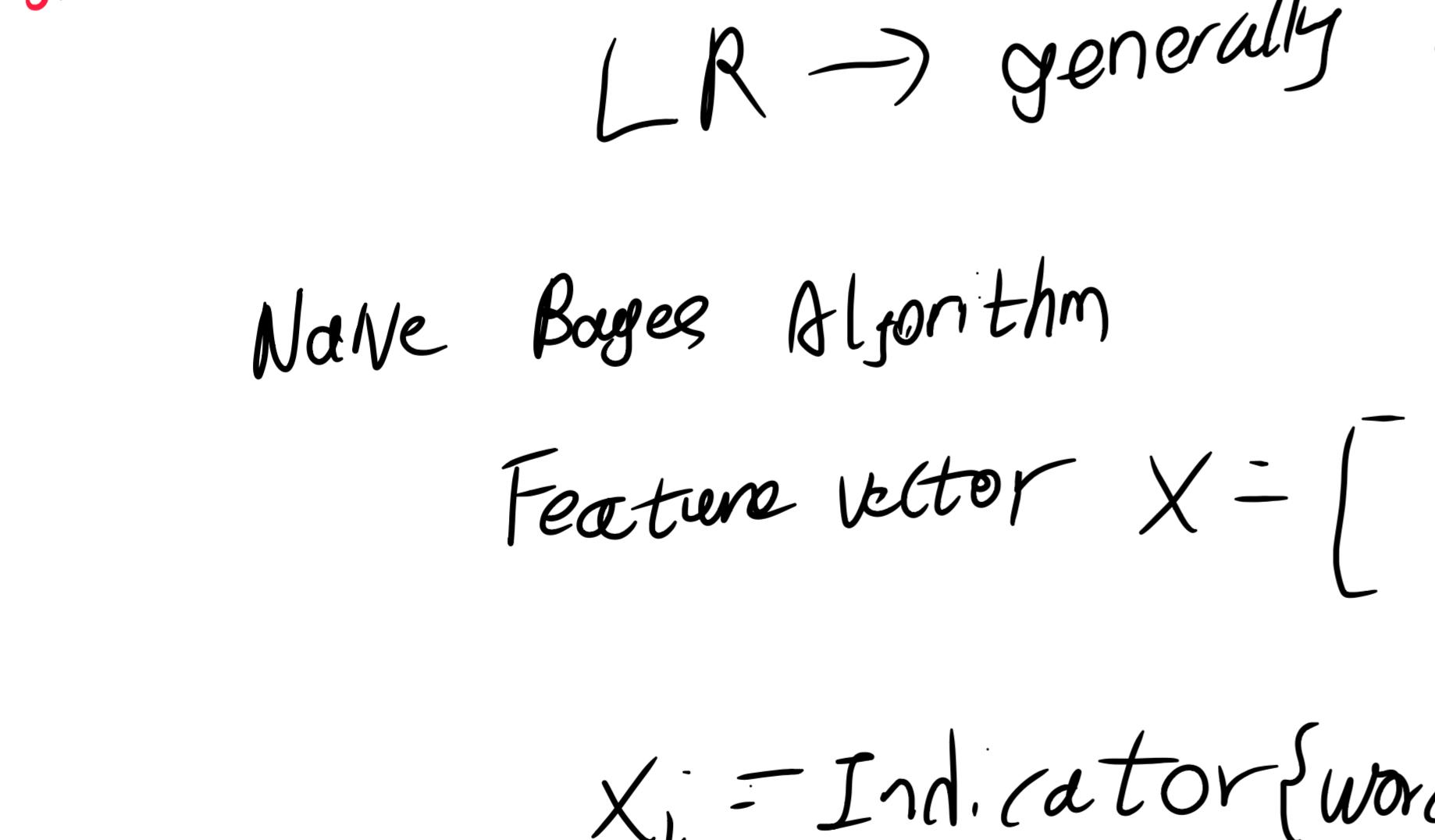
$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

$$\text{prediction: } \arg \max_y P(y|x) = \arg \max_y \frac{P(x|y)P(y)}{P(x)} = \arg \max_y P(x|y)P(y)$$

return y



Σ same? decision boundary linear? end up with quadratic terms in logistic function



Both discriminative / GDA use sigmoid function

$$\begin{aligned} \text{GDA assumes} \\ x|y=0 \sim N(\mu_0, \Sigma) \quad x|y=1 \sim N(\mu_1, \Sigma) \quad \Rightarrow \text{Logistic regression} \\ x|y=0 \sim \text{Poisson}(\lambda_0) \quad x|y=1 \sim \text{Poisson}(\lambda_1) \quad \Rightarrow P(y=1|x) \text{ is logistic} \\ y \sim \text{Bin}(\phi) \end{aligned}$$

If they both same exponentially, then logistic

GDA \rightarrow computation efficient

LR \rightarrow generally better

Naive Bayes Algorithm

$$\text{Feature vector } X = \begin{bmatrix} b \\ 0 \end{bmatrix} \quad \alpha \text{ ad a rank } \frac{1}{10000} = n$$

$$x_i = \text{Indicator}\{\text{word } i \text{ appears}\}$$

want to model $P(X|Y), P(Y)$

$$2^{10000} \times X$$

Assume X_i 's are conditionally independent

$$P(X_1 \dots X_{10000}|Y) = P(X_1|Y) P(X_2|X_1, Y) P(X_3|X_1, X_2, Y) \dots P(X_{10000}|X_1 \dots Y)$$

$$= P(X_1|Y) P(X_2|Y) \dots P(X_{10000}|Y)$$

$$= \prod_{i=1}^m P(X_i|Y)$$

Parameters

$$\phi_{j|y=1} = P(X_j=1 | Y=1)$$

$$\phi_{j|y=0} = P(X_j=1 | Y=0)$$

$$\phi_y = P(Y=1)$$

Joint likelihood

$$\mathcal{L}(\phi_y, \phi_{j|y})$$

$$\text{MLE: } \phi_y = \frac{\sum_{i=1}^m I\{Y_i=1\}}{m}$$

$$\phi_{j|y} = \frac{\sum_{i=1}^m I\{X_{ji}=1, Y_i=1\}}{\sum_{i=1}^m I\{Y_i=1\}}$$