# MLE/MAP

# +

# Naïve Bayes

Matt Gormley
Lecture 17
Mar. 20, 2019

# Reminders

- **Homework 5: Neural Networks**
  - **Out: Fri, Mar 1**
  - **Due: Fri, Mar 22 at 11:59pm**
- **Homework 6: Learning Theory / Generative Models**
  - **Out: Fri, Mar 22**
  - **Due: Fri, Mar 29 at 11:59pm (1 week)**
      **TIP: Do the readings!**


- **Today's In-Class Poll**
  - **http://p17.mlcourse.org**

# MLE AND MAP

# Likelihood Function

- Suppose we have N **samples** $D = \{x^{(1)}, x^{(2)}, \ldots, x^{(N)}\}$ from a **random variable** X

- The **likelihood** function:
  - <u>Case 1</u>: X is **discrete** with *pmf* $p(x|\theta)$
    $$L(\theta) = p(x^{(1)}|\theta)\, p(x^{(2)}|\theta) \ldots p(x^{(N)}|\theta)$$
  - <u>Case 2</u>: X is **continuous** with *pdf* $f(x|\theta)$
    $$L(\theta) = f(x^{(1)}|\theta)\, f(x^{(2)}|\theta) \ldots f(x^{(N)}|\theta)$$

In both cases (discrete / continuous), the **likelihood** tells us how likely one sample is relative to another

- The **log**-likelihood function:
  - <u>Case 1</u>: X is **discrete** with *pmf* $p(x|\theta)$
    $$\ell(\theta) = \log p(x^{(1)}|\theta) + \ldots + \log p(x^{(N)}|\theta)$$
  - <u>Case 2</u>: X is **continuous** with *pdf* $f(x|\theta)$
    $$\ell(\theta) = \log f(x^{(1)}|\theta) + \ldots + \log f(x^{(N)}|\theta)$$

# Likelihood Function

- Suppose we have N **samples** $D = \{(x^{(1)}, y^{(1)}), \ldots, (x^{(N)}, y^{(N)})\}$ from a pair of **random variables** X, Y

- The **conditional likelihood** function:
  - Case 1: Y is **discrete** with *pmf* $p(y \mid x, \theta)$
    $L(\theta) = p(y^{(1)} \mid x^{(1)}, \theta) \ldots p(y^{(N)} \mid x^{(N)}, \theta)$
  - Case 2: Y is **continuous** with *pdf* $f(y \mid x, \theta)$
    $L(\theta) = f(y^{(1)} \mid x^{(1)}, \theta) \ldots f(y^{(N)} \mid x^{(N)}, \theta)$

- The **joint likelihood** function:
  - Case 1: X and Y are **discrete** with *pmf* $p(x, y \mid \theta)$
    $L(\theta) = p(x^{(1)}, y^{(1)} \mid \theta) \ldots p(x^{(N)}, y^{(N)} \mid \theta)$
  - Case 2: X and Y are **continuous** with *pdf* $f(x, y \mid \theta)$
    $L(\theta) = f(x^{(1)}, y^{(1)} \mid \theta) \ldots f(x^{(N)}, y^{(N)} \mid \theta)$

# Likelihood Function

- Suppose we have N **samples** $D = \{(x^{(1)}, y^{(1)}), \ldots, (x^{(N)}, y^{(N)})\}$ from a pair of **random variables** X, Y

- The **joint likelihood** function:
  - Case 1: X and Y are **discrete** with *pmf* $p(x,y|\theta)$
    $$L(\theta) = p(x^{(1)}, y^{(1)}|\theta) \ldots p(x^{(N)}, y^{(N)}|\theta)$$

  - Case 2: X and Y are **continuous** with *pdf* $f(x,y|\theta)$
    $$L(\theta) = f(x^{(1)}, y^{(1)}|\theta) \ldots f(x^{(N)}, y^{(N)}|\theta)$$

  - Case 3: Y is **discrete** with *pmf* $p(y|\beta)$ and
    X is **continuous** with *pdf* $f(x|y,\alpha)$
    $$L(\alpha, \beta) = f(x^{(1)}| y^{(1)}, \alpha) p(y^{(1)}|\beta) \ldots f(x^{(N)}| y^{(N)}, \alpha) p(y^{(N)}|\beta)$$

  - Case 4: Y is **continuous** with *pdf* $f(y|\beta)$ and
    X is **discrete** with *pmf* $p(x|y,\alpha)$
    $$L(\alpha, \beta) = p(x^{(1)}| y^{(1)}, \alpha) f(y^{(1)}|\beta) \ldots p(x^{(N)}| y^{(N)}, \alpha) f(y^{(N)}|\beta)$$
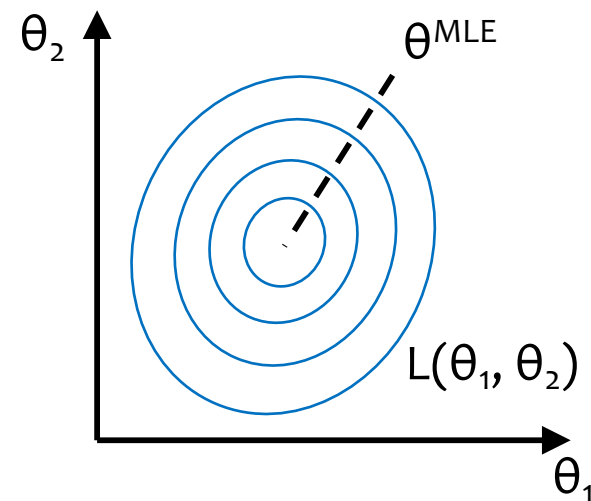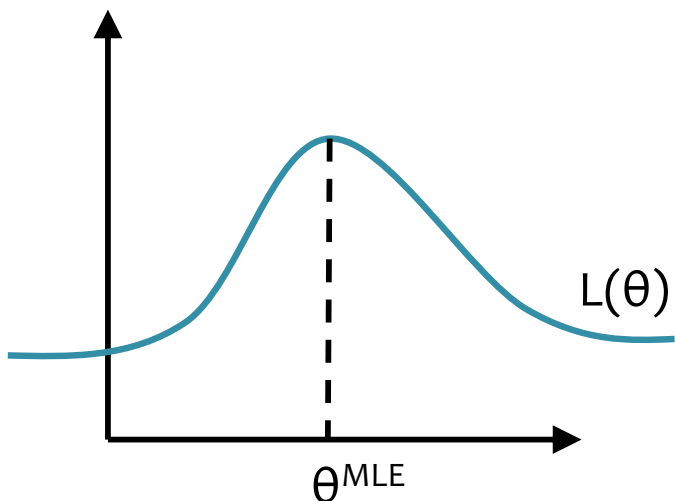
**Mixed discrete/continuous!**

19

# MLE

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likelihood Estimation:**

Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\text{MLE}} = \underset{\boldsymbol{\theta}}{\text{argmax}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)



$L(\theta)$

$\theta^{\text{MLE}}$



$\theta_2$

$\theta^{\text{MLE}}$

$L(\theta_1, \theta_2)$

$\theta_1$

# MLE

What does maximizing likelihood accomplish?

- There is only a finite amount of probability mass (i.e. sum-to-one constraint)

- MLE tries to allocate **as much** probability mass **as possible** to the things we have observed…

  …**at the expense** of the things we have **not** observed

# Recipe for Closed-form MLE

1. Assume data was generated i.i.d. from some model (i.e. write the generative story)
$$x^{(i)} \sim p(x|\theta)$$

2. Write log-likelihood
$$\ell(\theta) = \log p(x^{(1)}|\theta) + \dots + \log p(x^{(N)}|\theta)$$

3. Compute partial derivatives (i.e. gradient)
$$\partial\ell(\theta)/\partial\theta_1 = \dots$$
$$\partial\ell(\theta)/\partial\theta_2 = \dots$$
$$\dots$$
$$\partial\ell(\theta)/\partial\theta_M = \dots$$

4. Set derivatives to zero and solve for $\theta$
$$\partial\ell(\theta)/\partial\theta_m = 0 \text{ for all } m \in \{1, \dots, M\}$$
$$\theta^{MLE} = \text{solution to system of } M \text{ equations and } M \text{ variables}$$

5. Compute the second derivative and check that $\ell(\theta)$ is concave down at $\theta^{MLE}$

# MLE

Example: MLE of Exponential Distribution

Goal:

- pdf of Exponential$(\lambda)$: $f(x) = \lambda e^{-\lambda x}$
- Suppose $X_i \sim$ Exponential$(\lambda)$ for $1 \leq i \leq N$.
- Find MLE for data $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$

Steps:

- First write down log-likelihood of sample.
- Compute first derivative, set to zero, solve for $\lambda$.
- Compute second derivative and check that it is concave down at $\lambda^{\text{MLE}}$.

# MLE

- pdf of Exponential($\lambda$): $f(x) = \lambda e^{-\lambda x}$
- Suppose $X_i \sim$ Exponential($\lambda$) for $1 \leq i \leq N$.
- Find MLE for data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

## Example: MLE of Exponential Distribution

- First write down log-likelihood of sample.

$$\ell(\lambda) = \sum_{i=1}^{N} \log f(x^{(i)}) \tag{1}$$

$$= \sum_{i=1}^{N} \log(\lambda \exp(-\lambda x^{(i)})) \tag{2}$$

$$= \sum_{i=1}^{N} \log(\lambda) + -\lambda x^{(i)} \tag{3}$$

$$= N \log(\lambda) - \lambda \sum_{i=1}^{N} x^{(i)} \tag{4}$$

# MLE

- pdf of Exponential($\lambda$): $f(x) = \lambda e^{-\lambda x}$
- Suppose $X_i \sim$ Exponential($\lambda$) for $1 \le i \le N$.
- Find MLE for data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

## Example: MLE of Exponential Distribution

- Compute first derivative, set to zero, solve for $\lambda$.

$$\frac{d\ell(\lambda)}{d\lambda} = \frac{d}{d\lambda} N \log(\lambda) - \lambda \sum_{i=1}^{N} x^{(i)} \quad (1)$$

$$= \frac{N}{\lambda} - \sum_{i=1}^{N} x^{(i)} = 0 \quad (2)$$

$$\Rightarrow \lambda^{\text{MLE}} = \frac{N}{\sum_{i=1}^{N} x^{(i)}} \quad (3)$$

# MLE

**In-Class Exercise**

Show that the MLE of parameter $\phi$ for N samples drawn from Bernoulli($\phi$) is:

$$\phi_{MLE} = \frac{\text{Number of } x_i = 1}{N}$$

**Steps to answer:**

1. Write log-likelihood of sample

2. Compute derivative w.r.t. $\phi$

3. Set derivative to zero and solve for $\phi$

# MLE

**Question:**

Assume we have N samples $x^{(1)}$, $x^{(2)}$, ..., $x^{(N)}$ drawn from a Bernoulli($\phi$).

What is the **log-likelihood** of the data $\ell(\phi)$?

Assume $N_1$ = # of ($x^{(i)}$ = 1)
$N_0$ = # of ($x^{(i)}$ = 0)

**Answer:**

A. $l(\phi) = N_1 \log(\phi) + N_0 (1 - \log(\phi))$

B. $l(\phi) = N_1 \log(\phi) + N_0 \log(1-\phi)$

C. $l(\phi) = \log(\phi)^{N_1} + (1 - \log(\phi))^{N_0}$

D. $l(\phi) = \log(\phi)^{N_1} + \log(1-\phi)^{N_0}$

E. $l(\phi) = N_0 \log(\phi) + N_1 (1 - \log(\phi))$

F. $l(\phi) = N_0 \log(\phi) + N_1 \log(1-\phi)$

G. $l(\phi) = \log(\phi)^{N_0} + (1 - \log(\phi))^{N_1}$

H. $l(\phi) = \log(\phi)^{N_0} + \log(1-\phi)^{N_1}$

I. $l(\phi)$ = the most likely answer

# MLE

**Question:**

Assume we have N samples $x^{(1)}$, $x^{(2)}, \ldots, x^{(N)}$ drawn from a Bernoulli($\phi$).

What is the **derivative** of the log-likelihood $\partial \ell(\mathbf{\theta})/\partial \theta$?

Assume $N_1$ = # of $(x^{(i)} = 1)$

$N_0$ = # of $(x^{(i)} = 0)$

**Answer:**

A. $\partial \ell(\mathbf{\theta})/\partial \theta = \phi^{N_1} + (1 - \phi)^{N_0}$

B. $\partial \ell(\mathbf{\theta})/\partial \theta = \phi / N_1 + (1 - \phi) / N_0$

C. $\partial \ell(\mathbf{\theta})/\partial \theta = N_1 / \phi + N_0 / (1 - \phi)$

D. $\partial \ell(\mathbf{\theta})/\partial \theta = \log(\phi) / N_1 + \log(1 - \phi) / N_0$

E. $\partial \ell(\mathbf{\theta})/\partial \theta = N_1 / \log(\phi) + N_0 / \log(1 - \phi)$

# Learning from Data (Frequentist)

*Whiteboard*

- Optimization for MLE

- Examples: 1D and 2D optimization

- Example: MLE of Bernoulli

- Example: MLE of Categorical

- Aside: Method of Langrange Multipliers

# MLE vs. MAP

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likelihood Estimation:**

Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\mathrm{MLE}} = \underset{\boldsymbol{\theta}}{\mathrm{argmax}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)

**Principle of Maximum *a posteriori* (MAP) Estimation:**

Choose the parameters that maximize the posterior of the parameters given the data.

$$\boldsymbol{\theta}^{\mathrm{MAP}} = \underset{\boldsymbol{\theta}}{\mathrm{argmax}} \prod_{i=1}^{N} p(\boldsymbol{\theta}|\mathbf{x}^{(i)})$$

Maximum *a posteriori* (MAP) estimate

# MLE vs. MAP

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likelihood Estimation:**

Choose the parameters that maximize the likelihood of the data.

$$\boldsymbol{\theta}^{\text{MLE}} = \underset{\boldsymbol{\theta}}{\arg\max} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)

**Principle of Maximum *a posteriori* (MAP) Estimation:**

Choose the parameters that maximize the posterior of the parameters given the data.

Prior

$$\boldsymbol{\theta}^{\text{MAP}} = \underset{\boldsymbol{\theta}}{\arg\max} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

Maximum *a posteriori* (MAP) estimate

# MLE vs. MAP

Suppose we have data $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$

**Principle of Maximum Likeli**
Choose the parameters that
of the data.

$$\boldsymbol{\theta}^{\text{MLE}} = \text{argm}$$

Maximum Likelihood Estimate (MLE)

> **Important!**
> Usually the parameters are **continuous,** so the prior is a probability **density** function

**Principle of Maximum *a posteriori* (MAP) Estimation:**
Choose the parameters that maximize the posterior
of the parameters given the data.

Prior

$$\boldsymbol{\theta}^{\text{MAP}} = \underset{\boldsymbol{\theta}}{\text{argmax}} \prod_{i=1}^{N} p(\mathbf{x}^{(i)}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

Maximum *a posteriori* (MAP) estimate

# Learning from Data (Bayesian)

*Whiteboard*

&ndash; *maximum a posteriori* (MAP) estimation

&ndash; Optimization for MAP

&ndash; Example: MAP of Bernoulli&mdash;Beta

# Takeaways

- One view of what ML is trying to accomplish is **function approximation**
- The principle of **maximum likelihood estimation** provides an alternate view of learning

- **Synthetic data** can help **debug** ML algorithms
- Probability distributions can be used to **model** real data that occurs in the world
  (don't worry we'll make our distributions more interesting soon!)

# Learning Objectives

**MLE / MAP**

*You should be able to…*

1.  Recall probability basics, including but not limited to: discrete and continuous random variables, probability mass functions, probability density functions, events vs. random variables, expectation and variance, joint probability distributions, marginal probabilities, conditional probabilities, independence, conditional independence

2.  Describe common probability distributions such as the Beta, Dirichlet, Multinomial, Categorical, Gaussian, Exponential, etc.

3.  State the principle of maximum likelihood estimation and explain what it tries to accomplish

4.  State the principle of maximum a posteriori estimation and explain why we use it

5.  Derive the MLE or MAP parameters of a simple model in closed form

# NAÏVE BAYES

# Naïve Bayes Outline

- **Real-world Dataset**
  - Economist vs. Onion articles
  - Document → bag-of-words → binary feature vector
- **Naive Bayes: Model**
  - Generating synthetic "labeled documents"
  - Definition of model
  - Naive Bayes assumption
  - Counting # of parameters with / without NB assumption
- **Naïve Bayes: Learning from Data**
  - Data likelihood
  - MLE for Naive Bayes
  - MAP for Naive Bayes
- **Visualizing Gaussian Naive Bayes**

# Naïve Bayes

- Why are we talking about Naïve Bayes?
  - It's **just another decision function** that fits into our "big picture" recipe from last time
  - But it's our first **example of a Bayesian Network** and provides a *clearer* picture of **probabilistic learning**
  - Just like the other Bayes Nets we'll see, it **admits a closed form solution** for MLE and MAP
  - So learning is **extremely efficient** (just counting)

# Fake News Detector

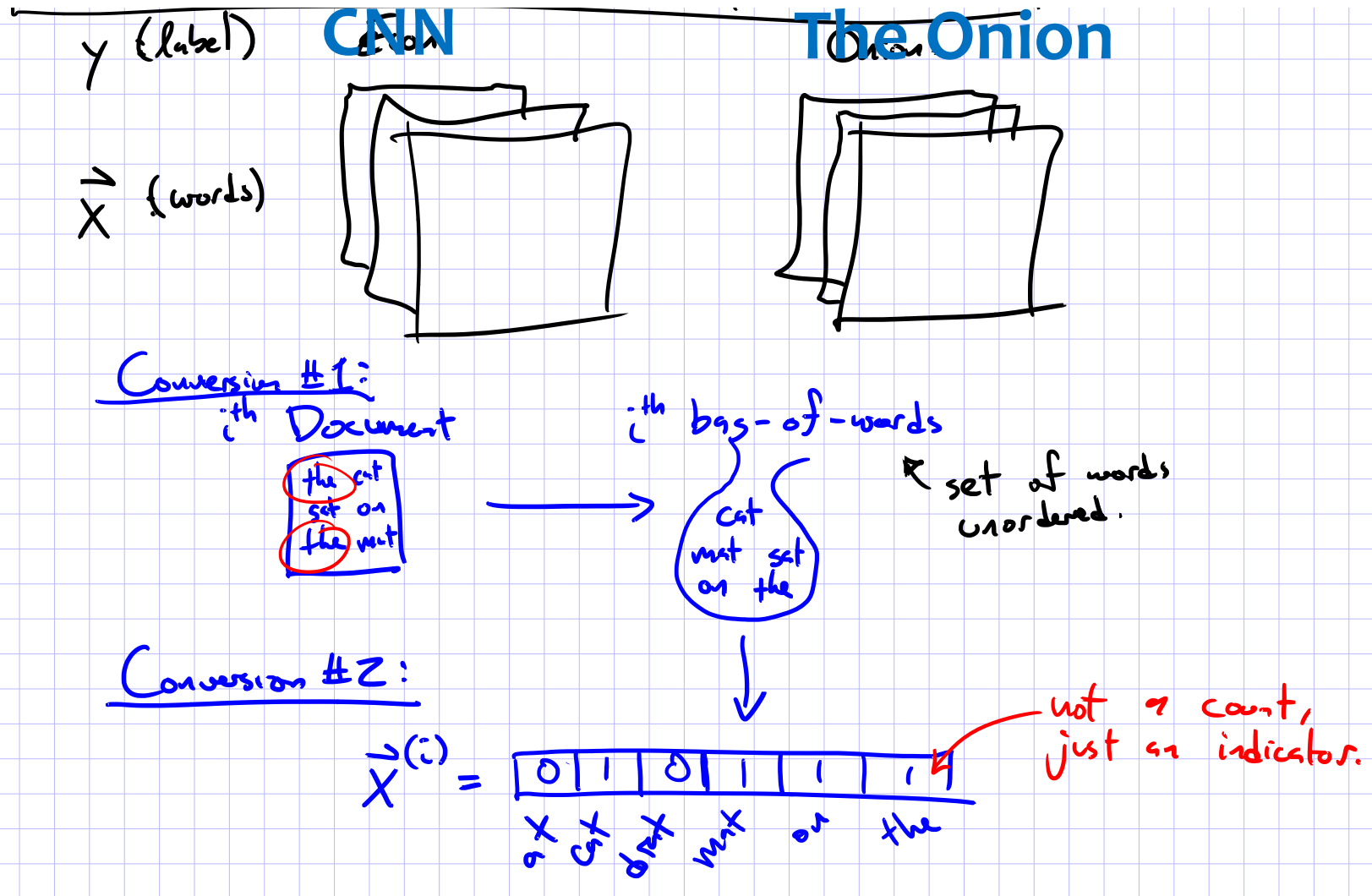**Today's Goal:** To define a generative model of emails of two different classes (e.g. real vs. fake news)

## CNN



## The Onion

# Fake News Detector



We can pretend the natural process generating these vectors is stochastic…
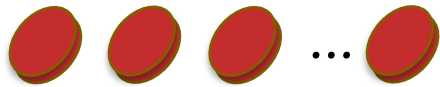
# Naive Bayes: Model

*Whiteboard*

- Document → bag-of-words → binary feature vector

- Generating synthetic "labeled documents"

- Definition of model

- Naive Bayes assumption

- Counting # of parameters with / without NB assumption
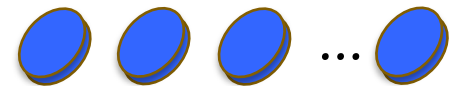
# Model 1: Bernoulli Naïve Bayes

Flip weighted coin

If HEADS, flip
each red coin

If TAILS, flip
each blue coin

| $y$ | $x_1$ | $x_2$ | $x_3$ | ... | $x_M$ |
|-----|-------|-------|-------|-----|-------|
| 0 | 1 | 0 | 1 | ... | 1 |
| 1 | 0 | 1 | 0 | ... | 1 |
| 1 | 1 | 1 | 1 | ... | 1 |
| 0 | 0 | 0 | 1 | ... | 1 |
| 0 | 1 | 0 | 1 | ... | 0 |
| 1 | 1 | 0 | 1 | ... | 0 |

We can **generate** data in this fashion. Though in practice we never would since our data is **given**.

Instead, this provides an explanation of **how** the data was generated (albeit a terrible one).

Each red coin corresponds to an $x_m$