# CSE 312

# Foundations of Computing II

**Lecture 22: Maximum Likelihood Estimation (MLE)**

## Agenda

- Idea: Estimation ◀
- Maximum Likelihood Estimation (example: mystery coin)
- Continuous MLE

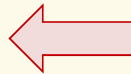# Probability vs Statistics

$\text{Ber}(p = 0.5)$ ⟹ **Probability** Given model, predict data ⟹ $P(THHTHH)$



$\text{Ber}(p = ??)$ ⟸ **Statistics** Given data, predict model ⟸ $THHTHH$

# Recall Formalizing Polls

Population size $N$, true fraction of voting in favor $p$, sample size $n$.
   **Problem:** We don't know $p$

## What type of r.v. is $X_i$?

|     |           | $\mathbb{E}[X_i]$ | $\mathrm{Var}(X_i)$ |
|-----|-----------|-------------------|---------------------|
| a.  | Bernoulli | $p$               | $p(1-p)$            |

## Polling Procedure

for $i = 1, \ldots, n$ :

1. Pick uniformly random person to call (prob: $1/N$)

2. Ask them how they will vote

$$X_i = \begin{cases} 1, & \text{voting in favor} \\ 0, & \text{otherwise} \end{cases}$$

Report our estimate of $p$:     $\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$

4

# Recall Formalizing Polls

We assume that poll answers $X_1, \ldots, X_n \sim \text{Ber}(p)$ i.i.d. for <u>unknown</u> $p$

$X_1 \qquad X_n$

**Goal:** Estimate $p$

We did this by computing $\hat{p} = \frac{1}{n} \sum_{i=1}^{n} X_i$

"$p$ hat"

# Notation – Parametric Model (discrete case)

**Definition.** A **(parametric) model** is a family of distributions indexed by a parameter $\theta$, described by a two-argument function
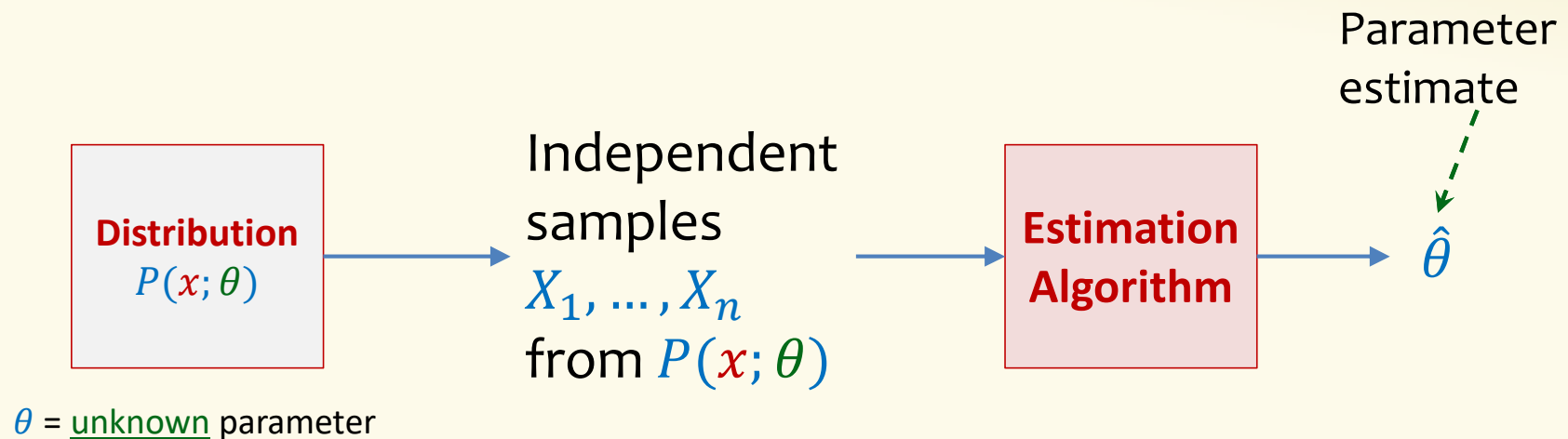
$$P(x; \theta) = \text{prob. of outcome } x \text{ when distribution has parameter } \theta$$

$$[\text{i.e., every } \theta \text{ defines a different distribution } \sum_x P(x; \theta) = 1]$$

**Examples**

- "Bernoullis": $Ber(p)$ $P(x; \theta = p) = \begin{cases} p & x = 1 \\ 1 - p & x = 0 \end{cases}$

- "Geometrics": $Geo(p)$ $P(i; \theta = p) = (1 - p)^{i-1} p \quad$ for $i \in \mathbb{N}$

# Statistics: Parameter Estimation – Workflow

Parameter estimate

$$\boxed{\begin{array}{c} \textbf{Distribution} \\ P(x; \theta) \end{array}} \longrightarrow \begin{array}{l} \text{Independent} \\ \text{samples} \\ X_1, \dots, X_n \\ \text{from } P(x; \theta) \end{array} \longrightarrow \boxed{\begin{array}{c} \textbf{Estimation} \\ \textbf{Algorithm} \end{array}} \longrightarrow \hat{\theta}$$

$\theta$ = <u>unknown</u> parameter

**Example:** coin flip distribution with unknown $\theta$ = probability of heads

Observation: $HTTHHHTHTHTTTTHTHTTTTTHT$

**Goal:** Estimate $\theta$

# Example

Suppose we have a mystery coin with some probability $p$ of coming up heads. We flip the coin 8 times, independent of other flips, and see the following sequence flips

*TTHTHTTH*

3 heads

Given this data, what would you estimate $p$ is?

8 flips

Poll: pollev.com/paulbeame028
a.  1/2
b.  5/8
c.  3/8
d.  1/4

sample mean

8

# Agenda

- Idea: Estimation
- Maximum Likelihood Estimation (example: mystery coin) ◀
- Continuous MLE

# Likelihood

Say we see outcome $HHTHH$.

You tell me your best guess about the value of the unknown parameter $\theta$ (a.k.a. $p$) is 4/5. Is there some way that you can argue "objectively" that this is the best estimate?

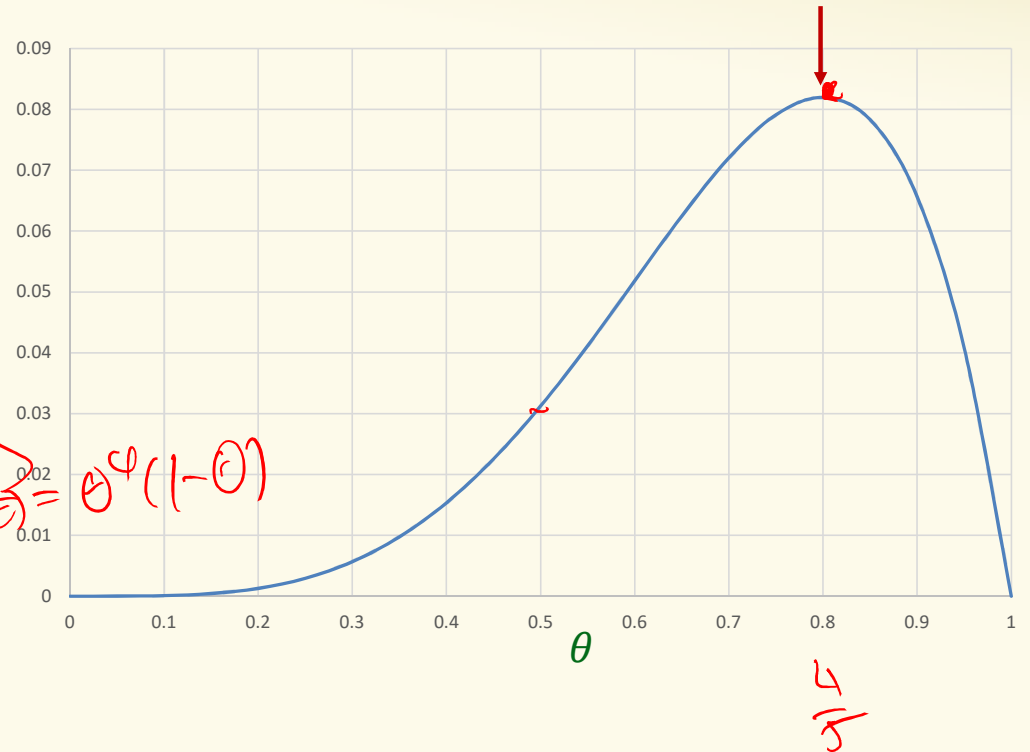# Likelihood

Say we see outcome *HHTHH*.

$$\mathcal{L}(HHTHH \mid \theta) = \theta^4(1 - \theta)$$

Probability of observing the outcome *HHTHH* if $\theta$ = prob. of heads.

For a fixed outcome *HHTHH* , this is a function of $\theta$.

$\theta \;\; \theta \;\; (1-\theta) \;\; \theta \;\; \theta$

$f(\theta) = \theta^4(1-\theta)$

$\dfrac{4}{5}$

**Max Prob of seeing HHTHH**



11

# Likelihood of Different Observations <span style="float:right">(Discrete case)</span>

**Definition.** The **likelihood** of independent observations $x_1, \ldots, x_n$ is

$$\mathcal{L}(x_1, \ldots, x_n \mid \theta) = \prod_{i=1}^{n} P(x_i; \theta)$$

*# observations* *probability* *# independent*

**Maximum Likelihood Estimation (MLE).** Given data $x_1, \ldots, x_n$, find $\hat{\theta}$ such that $\mathcal{L}(x_1, \ldots, x_n \mid \hat{\theta})$ is maximized!

$$\hat{\theta} = \underset{\theta}{\text{argmax}} \; \mathcal{L}(x_1, \ldots, x_n \mid \theta)$$

*continuous funch of $\theta$*

*log e*

Usually: Solve $\dfrac{\partial \mathcal{L}(x_1, \ldots, x_n \mid \theta)}{\partial \theta} = 0$ or $\dfrac{\partial \ln \mathcal{L}(x_1, \ldots, x_n \mid \theta)}{\partial \theta} = 0$ [+check it's a max!]

# Likelihood vs. Probability

- Fixed $\theta$: **probability** $\prod_{i=1}^{n} P(x_i; \theta)$ that dataset $x_1, \ldots, x_n$ is sampled by distribution with parameter $\theta$
  - A function of $x_1, \ldots, x_n$

- Fixed $x_1, \ldots, x_n$: **likelihood** $\mathcal{L}(x_1, \ldots, x_n | \theta)$ that parameter $\theta$ explains dataset $x_1, \ldots, x_n$.
  - A function of $\theta$

These notions are the same number if we fix _both_ $x_1, \ldots, x_n$ and $\theta$, but different role/interpretation

# Example – Coin Flips

Observe: Coin-flip outcomes $x_1, \ldots, x_n$, with $n_H$ heads, $n_T$ tails

 – i.e., $n_H + n_T = n$ 　　　　**Goal:** estimate $\theta$ = prob. heads.

$$\mathcal{L}(x_1, \ldots, x_n | \theta) = \theta^{n_H} (1 - \theta)^{n_T}$$

$$\frac{\partial}{\partial \theta} \mathcal{L}(x_1, \ldots, x_n | \theta) = ???$$

While it is possible to compute this derivative, it's not always nice since we are working with products.

# Log-Likelihood

We can save some work if we work with the **log-likelihood** instead of the likelihood directly.

**Definition.** The **log-likelihood** of independent observations $x_1, \ldots, x_n$ is

$$\ln \mathcal{L}(x_1, \ldots, x_n \mid \theta) = \ln \prod_{i=1}^{n} P(x_i; \theta) = \sum_{i=1}^{n} \ln P(x_i; \theta)$$

Useful log properties

$$\ln(ab) = \ln(a) + \ln(b)$$
$$\ln(a/b) = \ln(a) - \ln(b)$$
$$\ln(a^b) = b \cdot \ln(a)$$

## Example – Coin Flips

Observe: Coin-flip outcomes $x_1, \ldots, x_n$, with $n_H$ heads, $n_T$ tails
  – i.e., $n_H + n_T = n$

**Goal:** estimate $\theta$ = prob. heads.

$$\mathcal{L}(x_1, \ldots, x_n \mid \theta) = \theta^{n_H} (1 - \theta)^{n_T}$$

$$\ln \mathcal{L}(x_1, \ldots, x_n \mid \theta) = n_H \ln \theta + n_T \ln(1 - \theta)$$

$$\frac{\partial}{\partial \theta} \ln \mathcal{L}(x_1, \ldots, x_n \mid \theta) = n_H \cdot \frac{1}{\theta} - n_T \cdot \frac{1}{1 - \theta}$$

Want value $\hat{\theta}$ of $\theta$ s.t. $\frac{\partial}{\partial \theta} \ln \mathcal{L}(x_1, \ldots, x_n \mid \theta) = 0$

So we need $n_H \cdot \frac{1}{\hat{\theta}} - n_T \cdot \frac{1}{1 - \hat{\theta}} = 0$

Solving gives
$$\hat{\theta} = \frac{n_H}{n}$$

$(\ln x)' = \frac{1}{x}$

$\frac{n_H}{\theta} = \frac{n_T}{1-\theta}$

$(1-\theta) n_H = \theta \cdot n_T$

$n_H = \theta(n_H + n_T) = n\theta$

16

# General Recipe

1. **Input** Given $n$ i.i.d. samples $x_1, \ldots, x_n$ from parametric model with parameter $\theta$.

2. **Likelihood** Define your likelihood $\mathcal{L}(x_1, \ldots, x_n | \theta)$.
   - For discrete $\qquad \mathcal{L}(x_1, \ldots, x_n | \theta) = \prod_{i=1}^{n} P(x_i ; \theta)$

3. **Log** Compute $\ln \mathcal{L}(x_1, \ldots, x_n | \theta)$

4. **Differentiate** Compute $\frac{\partial}{\partial \theta} \ln \mathcal{L}(x_1, \ldots, x_n | \theta)$

5. **Solve for** $\hat{\theta}$ by setting derivative to $0$ and solving for max.

Generally, you need to do a second derivative test to verify it is a maximum, but we won't ask you to do that in CSE 312.

# Brain Break

## Agenda

- Idea: Estimation
- Maximum Likelihood Estimation (example: mystery coin)
- Continuous MLE ◀

# The Continuous Case

Given $n$ (independent) samples $x_1, \ldots, x_n$ from (continuous) parametric model $f(x_i; \theta)$ which is now a family of <u>densities</u>

**Definition.** The **likelihood** of independent observations $x_1, \ldots, x_n$ is

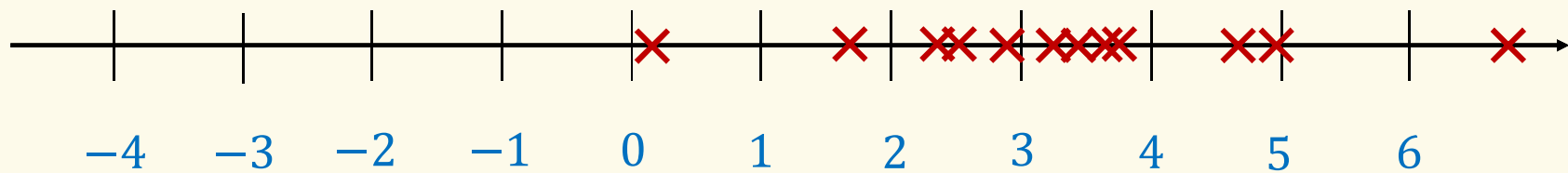$$\mathcal{L}(x_1, \ldots, x_n \mid \theta) = \prod_{i=1}^{n} f(x_i; \theta)$$

Density function! (Why?)

# Why density?
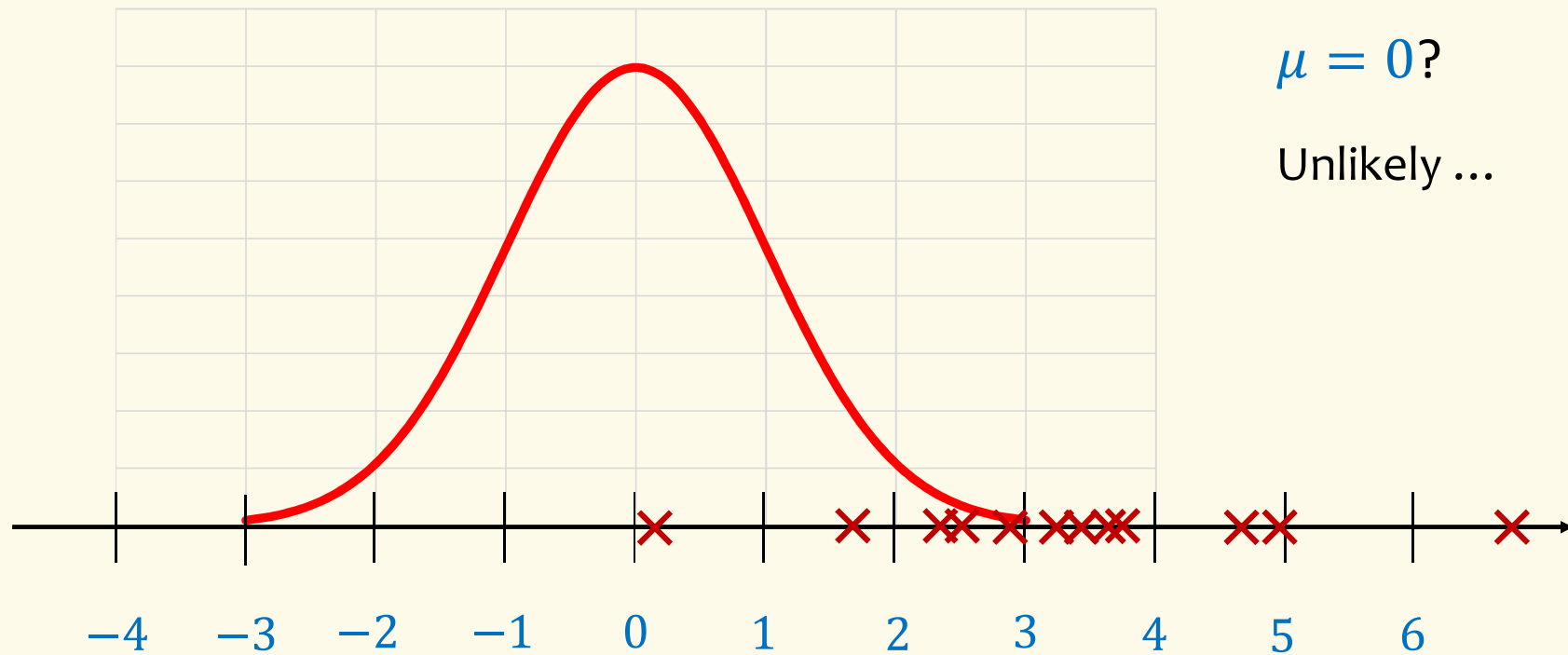
- Density ≠ probability, but:
  - For maximizing likelihood, <span style="color:red">we really only care about relative likelihoods</span>, and density captures that
  - has desired property that likelihood increases with better fit to the model

*var is fixed*

$n$ samples $x_1, \ldots, x_n \in \mathbb{R}$ from Gaussian $\mathcal{N}(\mu, 1)$. <u>Most likely</u> $\mu$?

[i.e., we are given the <u>promise</u> that the variance is 1]

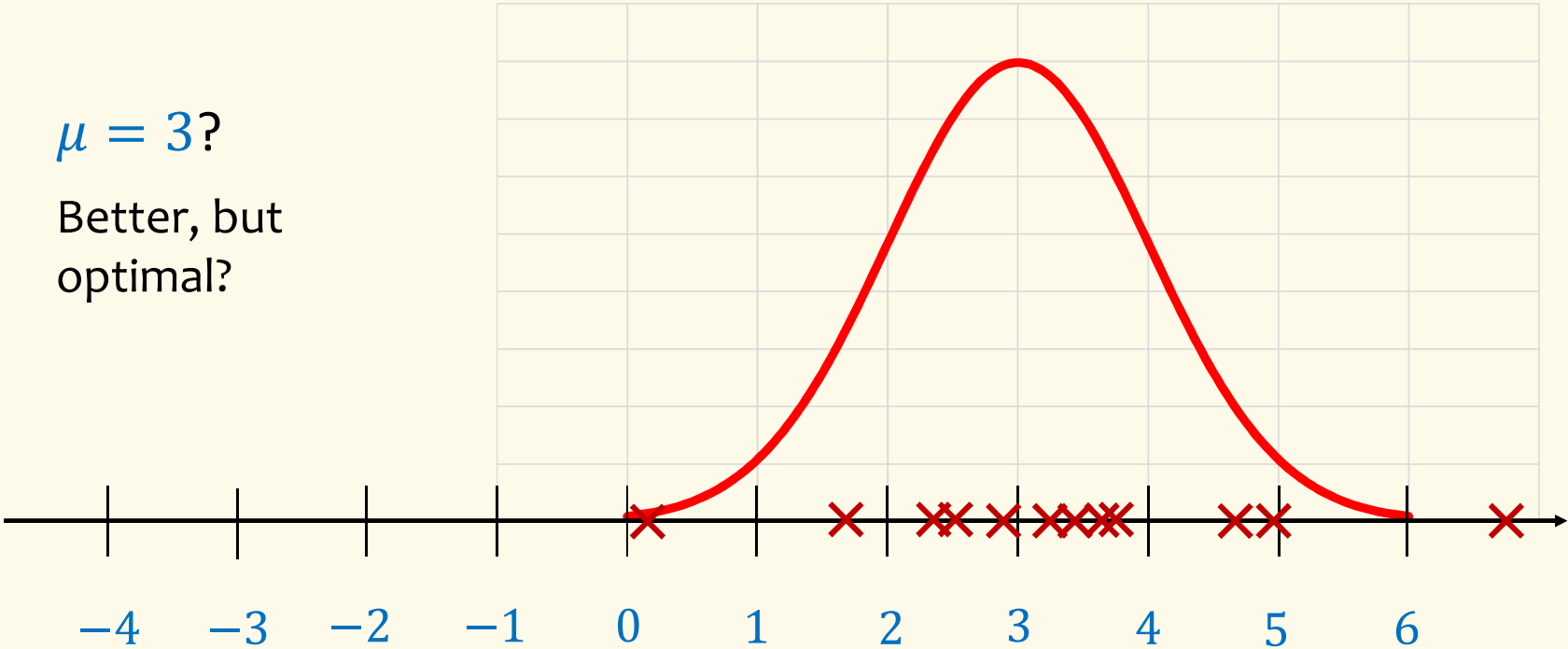$n$ samples $x_1, \ldots, x_n \in \mathbb{R}$ from Gaussian $\mathcal{N}(\mu, 1)$. Most likely $\mu$?

$\mu = 0$?

Unlikely …

$n$ samples $x_1, \ldots, x_n \in \mathbb{R}$ from Gaussian $\mathcal{N}(\mu, 1)$. <u>Most likely</u> $\mu$?

$\mu = 3$?

Better, but optimal?

# Example – Gaussian Parameters

Normal outcomes $x_1, \dots, x_n$, known variance $\sigma^2 = 1$ but *unknown* mean $\mu$

**Goal:** estimate $\theta$ = mean

**Next time:**
$$\hat{\theta} = \frac{\sum_i^n x_i}{n}$$

In other words, MLE is the *sample mean* of the data.

$P(x)(\ \ )\theta=1$

$\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}}$

Gauss $\sigma^2 = 1$

$f(x_i; \theta)$ for

$\mathcal{L}(x_1, \dots, n | \theta) = \prod_{i=1}^{n} \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i-\theta)^2}{2}} \right)$

$= \left( \frac{1}{\sqrt{2\pi}} \right)^n \prod_{i=1}^{n} e^{-\frac{(x_i-\theta)^2}{2}}$

$\ln \mathcal{L}(x_1 \dots x_n | \theta) =$

$n \ln\left(\frac{1}{\sqrt{2\pi}}\right) + \sum_{i=1}^{n} \left( -\frac{(x_i-\theta)^2}{2} \right)$

$\ln e^{g} a$

sample mean

# General Recipe

1. **Input** Given $n$ i.i.d. samples $x_1, \ldots, x_n$ from parametric model with parameter $\theta$.

2. **Likelihood** Define your likelihood $\mathcal{L}(x_1, \ldots, x_n | \theta)$.
   - For discrete $\qquad \mathcal{L}(x_1, \ldots, x_n | \theta) = \prod_{i=1}^{n} P(x_i ; \theta)$
   - For continuous $\quad \mathcal{L}(x_1, \ldots, x_n | \theta) = \prod_{i=1}^{n} f(x_i ; \theta)$

3. **Log** Compute $\ln \mathcal{L}(x_1, \ldots, x_n | \theta)$

4. **Differentiate** Compute $\frac{\partial}{\partial \theta} \ln \mathcal{L}(x_1, \ldots, x_n | \theta)$

5. **Solve for** $\hat{\theta}$ by setting derivative to $0$ and solving for max.

Generally, you need to do a second derivative test to verify it is a maximum, but we won't ask you to do that in CSE 312.