# Maximum Likelihood Estimation

# Your first consulting job

- *Billionaire*: I have special coin, if I flip it, what's the probability it will be heads?
- *You*: Please flip it a few times:

flip $n=5$ times

HHTHT

- *You*: The probability is: $3/5$

$$\frac{k=3 \text{ heads}}{n=5 \text{ flips}}$$

- *Billionaire:* Why?

Independent, yani P(Xs=b)'nin P(Xt=a) üzerinde etkisi yok.

# Coin – Binomial Distribution

- **Data**: sequence *D= (HHTHT* ~~...~~ *)*, **k heads** out of **n flips**   Observed ✓

- **Hypothesis:** P(Heads) = θ,  P(Tails) = 1-θ

  - Flips are i.i.d.:  If  $X_t \in \{0,1\}$ denoting $t$-th flip

    - Independent events  $\mathbb{P}(X_t = a, X_s = b) = \mathbb{P}(X_t = a)\,\mathbb{P}(X_s = b)$

    - Identically distributed according to Binomial distribution
    
    $$\mathbb{P}(X_t = 0) = \mathbb{P}(X_s = 0) = 1 - \theta$$
    $$\mathbb{P}(X_t = 1) = \theta$$

- $P(\mathcal{D}|\theta) = P(HHTHT \,|\, \theta)$

$$= P(H|\theta)\,P(H|\theta)\,P(T|\theta)\,P(H|\theta)\,P(T|\theta)$$

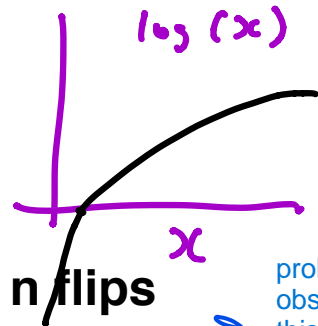$$= \theta^3 (1-\theta)^2 \quad\bigg/\quad \begin{array}{l} \text{In general} \\ P(D|\theta) = \theta^k (1-\theta)^{n-k} \end{array}$$

# Maximum Likelihood Estimation

true prob., we don't know that value.

$P(head) = \theta^*$

$\log(x)$

- **Data**: sequence *D= (HHTHT…),* **k heads** out of **n flips**
- **Hypothesis:** P(Heads) = θ,  P(Tails) = 1-θ

prob. of observing this data given that theta is true is this expression.

Likelihood
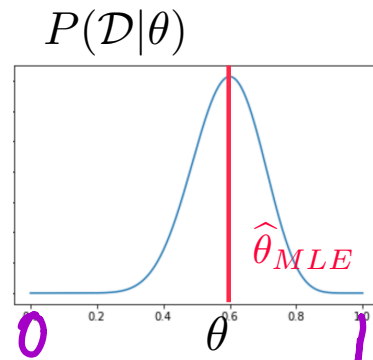
$$P(\mathcal{D}|\theta) = \theta^k (1-\theta)^{n-k}$$

- Maximum likelihood estimation (MLE): Choose θ that maximizes the probability of observed data:

$$\widehat{\theta}_{MLE} = \arg\max_{\theta} \ P(\mathcal{D}|\theta)$$

just scaling

$$= \arg\max_{\theta} \ \log P(\mathcal{D}|\theta)$$

$P(\mathcal{D}|\theta)$

$\widehat{\theta}_{MLE}$

0    0.2    0.4    $\theta$ 0.6    0.8    1.0    1

think of this as a function of theta and choose the theta that maximize this func.

BECAUSE LOG FUNCTION IS MONOTONICALLY INCREASING, MAXIMIZING THE LOG OF SOMETHING IS THE SAME AS MAXIMIZING ITSELF.

# Your first learning algorithm

$$\log(ab) = \log(a) + \log(b)$$
$$\text{for any } a, b > 0$$

türev 0

$$\widehat{\theta}_{MLE} := \arg\max_{\theta} \ \log\left(P(\mathcal{D}|\theta)\right)$$

$$= \arg\max_{\theta} \ \log\left(\theta^k (1-\theta)^{n-k}\right)$$

- Set derivative to zero:

$$\boxed{\frac{d}{d\theta} \log P(\mathcal{D}|\theta) = 0}$$

$$\rightarrow \ = \arg\max_{\theta} \ k \log(\theta) + (n-k) \log(1-\theta)$$

$$\frac{\partial}{\partial \theta}\left[ \cdot \right] = \frac{k}{\theta} + \frac{n-k}{1-\theta} \cdot (-1) = 0 \qquad \text{(multiply } \theta(1-\theta) \text{ on both sides)}$$

$$\rightarrow \ (1-\theta)k - \theta(n-k) = 0 \implies \widehat{\theta}_{MLE} = \frac{k}{n}$$

# How many flips do I need?

$$\widehat{\theta}_{MLE} = \frac{k}{n}$$

- *You*: flip the coin 5 times. *Billionaire*: I got 3 heads.

$$\widehat{\theta}_{MLE} = 3/5$$

- *You*: flip the coin 50 times. *Billionaire*: I got 20 heads.

$$\widehat{\theta}_{MLE} = 20/50 = 2/5$$

trust that answer a little bit more. Because 50 > 5

- *Billionaire:* Which one is right? Why?

# Quantifying Uncertainty

- For **n flips** and **k heads** the MLE is **unbiased** for true $\theta^*$:

$$\widehat{\theta}_{MLE} = \frac{k}{n} \qquad \mathbb{E}[\widehat{\theta}_{MLE}] = \theta^*$$

- **Expectation** describes how the estimator behaves *on average.*

$$\widehat{\theta}_{MLE} = \frac{1}{n} \sum_{t=1}^{n} \mathbb{1}\{X_t = H\}$$

$$X_t \in \{H, T\}$$

$$\mathbb{1}\{\Omega\} = \begin{cases} 1 & \text{if } \Omega = \text{true} \\ 0 & \text{o.w.} \end{cases}$$
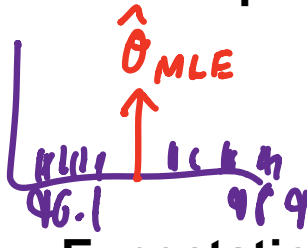
By assumption, $\exists \theta^* : \mathbb{P}(X_t = H) = \theta^*$

$$= \mathbb{E}[\mathbb{1}\{X_t = H\}]$$

= P(Xt=H) * I(Xt=H) + P(Xt=T) * I(Xt=H)

$$\mathbb{E}[\widehat{\theta}_{MLE}] = \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[\mathbb{1}\{X_t = H\}] = \frac{1}{n} \sum_{t} \theta^* = \theta^*$$

# Quantifying Uncertainty

- For **n flips** and **k heads** the MLE is **unbiased** for true $\theta^*$:

$$\widehat{\theta}_{MLE} = \frac{k}{n} \qquad \mathbb{E}[\widehat{\theta}_{MLE}] = \theta^*$$

Random · Deterministic

- **Expectation** describes how the estimator behaves *on average.*
- The **Variance** is the expected squared deviation from the mean:

$$\text{Variance}(\widehat{\theta}_{MLE}) := \mathbb{E}\left[\left(\widehat{\theta}_{MLE} - \mathbb{E}[\widehat{\theta}_{MLE}]\right)^2\right]$$

- As a rule of thumb:

$$\widehat{\theta}_{MLE} \approx \mathbb{E}[\widehat{\theta}_{MLE}] \pm \sqrt{\text{Variance}(\widehat{\theta}_{MLE})}$$

- **Exercise**: compute the $\text{Variance}(\widehat{\theta}_{MLE})$

# Expectation versus High Probability

- For **n flips** and **k heads** the MLE is **unbiased** for true θ*:

$$\widehat{\theta}_{MLE} = \frac{k}{n} \qquad \mathbb{E}[\widehat{\theta}_{MLE}] = \theta^*$$

- Expectation describes how the estimator behaves *on average.*
- For any ε>0 can we bound $\mathbb{P}(|\widehat{\theta}_{MLE} - \mathbb{E}[\widehat{\theta}_{MLE}]| \geq \epsilon)$ ?

**Markov's inequality**
For any $t > 0$ and non-negative random variable $X$

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[X]}{t}$$

- **Exercise**: Apply Markov's inequality to obtain bound.
  (Hint: set $X = |\widehat{\theta}_{MLE} - \theta^*|^2$ )

# Maximum Likelihood Estimation

Each $X_i$ is iid from $f(x;\theta)$ $\longrightarrow$ PMF

**Observe** $X_1, X_2, \ldots, X_n$ drawn IID from $f(x; \theta)$ for some "true" $\theta = \theta_*$

**Likelihood function** $L_n(\theta) = \prod_{i=1}^{n} f(X_i; \theta)$

**Log-Likelihood function** $l_n(\theta) = \log(L_n(\theta)) = \sum_{i=1}^{n} \log(f(X_i; \theta))$

**Maximum Likelihood Estimator (MLE)** $\widehat{\theta}_{MLE} = \arg\max_{\theta} L_n(\theta)$
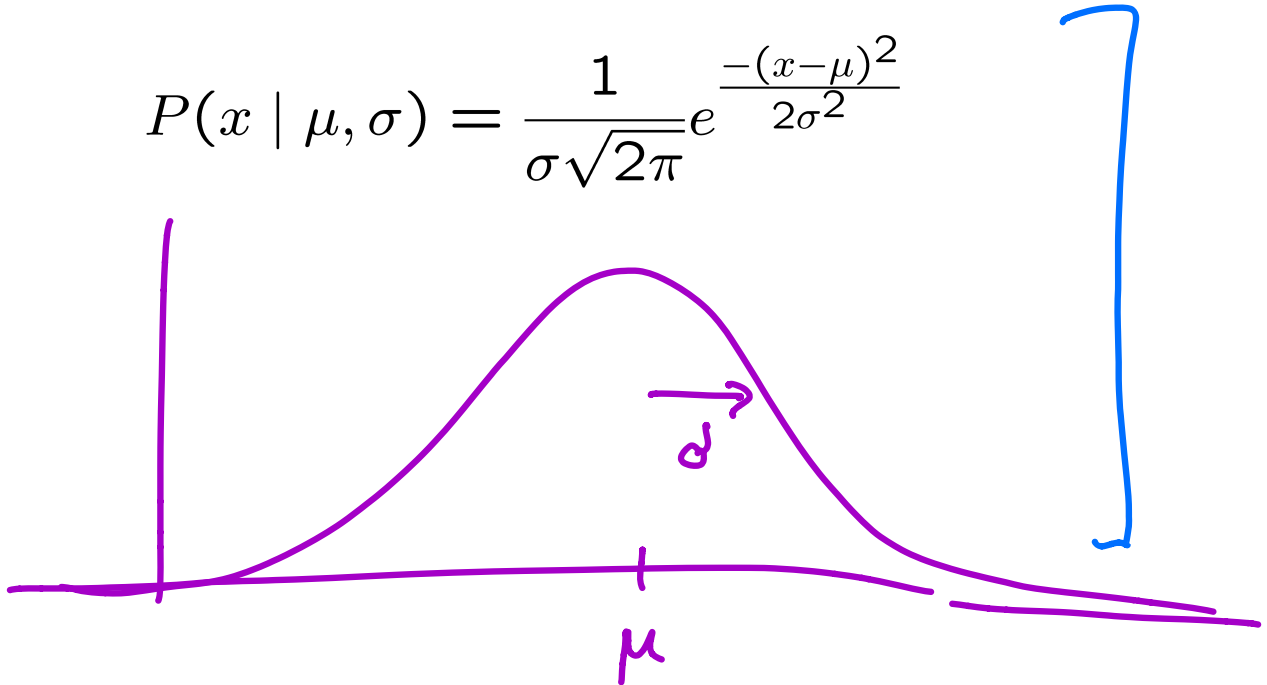
Set d/dx (log(Ln(theta))) to ZERO. $\longrightarrow$ for closed-form solutions

# What about continuous variables?

- *Billionaire*: What if I am measuring a **continuous variable**?
- *You*: **Let me tell you about Gaussians…**

$$P(x \mid \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

# Some properties of Gaussians

- affine transformation (multiplying by scalar and adding a constant)
    - $X \sim N(\mu, \sigma^2)$
    - $Y = aX + b$ ➜ $Y \sim N(a\mu + b, a^2\sigma^2)$

- Sum of Gaussians
    - $X \sim N(\mu_X, \sigma^2_X)$
    - $Y \sim N(\mu_Y, \sigma^2_Y)$
    - $Z = X + Y$ ➜ $Z \sim N(\mu_X + \mu_Y, \sigma^2_X + \sigma^2_Y)$

# MLE for Gaussian

- Prob. of i.i.d. samples $D = \{x_1, \ldots, x_n\}$ (e.g., temperature):

$$P(\mathcal{D}|\mu, \sigma) = P(x_1, \ldots, x_n|\mu, \sigma) = \prod_{i=1}^{n} P(x_i | \mu, \sigma)$$

We can draw different distributions by trying different (mu,sigma) tuples.

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

- Log-likelihood of data:

$$P(x|\mu, \sigma)$$

$$\log P(\mathcal{D}|\mu, \sigma) = -n\log(\sigma\sqrt{2\pi}) - \sum_{i=1}^{n} \frac{(x_i - \mu)^2}{2\sigma^2}$$

$$\mu = 0.7$$
$$\sigma = 2$$

- What is $\widehat{\theta}_{MLE}$ for $\theta = (\mu, \sigma^2)$? Draw a picture!

$$\mu = 0, \sigma = 1$$

$$D = \{4.5, 5.3, 4.8., \ldots.\}$$

$$\mu = 8, \sigma = 0.5$$

4.5   4.8      5.3

# Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for mean?

$$\frac{d}{d\mu} \log P(\mathcal{D}|\mu, \sigma) = \frac{d}{d\mu} \left[ -n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^{n} \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= \sum_{i=1}^{n} \frac{(x_i - \mu)}{\sigma^2} = 0$$

$$\implies \hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

# MLE for variance

- Again, set derivative to zero:

$$\frac{d}{d\sigma} \log P(\mathcal{D}|\mu, \sigma) = \frac{d}{d\sigma} \left[ -n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^{n} \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= -n \cdot \frac{1}{\sigma\sqrt{2\pi}} \cdot \sqrt{2\pi} - \sum_{i=1}^{n} \frac{(x_i - \mu)^2}{2} \cdot (-2)\sigma^{-3} = 0$$

(multiply $\sigma^3$ on both sides after setting deriv. $= 0$)

$$\longrightarrow \quad -n\sigma^2 + \sum_{i=1}^{n} (x_i - \mu)^2 = 0$$

$$\hat{\sigma}^2_{MLE} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu)^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu}_{MLE})^2$$

# Learning Gaussian parameters

- MLE:

$$\widehat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

$$\widehat{\sigma^2}_{MLE} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \widehat{\mu}_{MLE})^2$$

- MLE for the variance of a Gaussian is **biased** $\left( Exercise \right)$

$$\mathbb{E}[\widehat{\sigma^2}_{MLE}] \neq \sigma^2$$

- Unbiased variance estimator:

$$\widehat{\sigma^2}_{unbiased} = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \widehat{\mu}_{MLE})^2$$

# Maximum Likelihood Estimation

**Observe** $X_1, X_2, \ldots, X_n$ drawn IID from $f(x; \theta)$ for some "true" $\theta = \theta_*$

**Likelihood function** $\quad L_n(\theta) = \prod_{i=1}^{n} f(X_i; \theta)$

**Log-Likelihood function** $\quad l_n(\theta) = \log(L_n(\theta)) = \sum_{i=1}^{n} \log(f(X_i; \theta))$

**Maximum Likelihood Estimator (MLE)** $\quad \widehat{\theta}_{MLE} = \arg\max_{\theta} L_n(\theta)$

$$X_i \sim \exp(\lambda) \qquad P(X_i \geq t) = e^{-\lambda t}$$

$$\widehat{\theta}_{MLE} \text{ is unbiased if}$$
$$\mathbb{E}[\widehat{\theta}_{MLE}] = \theta^*$$

# Maximum Likelihood Estimation

**Observe** $X_1, X_2, \ldots, X_n$ drawn IID from $f(x; \theta)$ for some "true" $\theta = \theta_*$

**Likelihood function** $\quad L_n(\theta) = \prod_{i=1}^{n} f(X_i; \theta)$

**Log-Likelihood function** $\quad l_n(\theta) = \log(L_n(\theta)) = \sum_{i=1}^{n} \log(f(X_i; \theta))$

**Maximum Likelihood Estimator (MLE)** $\quad \widehat{\theta}_{MLE} = \arg\max_{\theta} L_n(\theta)$

$\sqrt{\text{Sample Variance}}$

Properties (under benign regularity conditions—smoothness, identifiability, etc.):

• Asymptotically consistent and normal: $\dfrac{\widehat{\theta}_{MLE} - \theta_*}{\widehat{se}} \sim \mathcal{N}(0, 1)$

• Asymptotic Optimality, minimum variance (see Cramer-Rao lower bound)

# Recap

- Learning is…
  - Collect some data
    - E.g., coin flips
  - Choose a hypothesis class or model
    - E.g., binomial
  - Choose a loss function
    - E.g., data likelihood
  - Choose an optimization procedure
    - E.g., set derivative to zero to obtain MLE
  - Justifying the accuracy of the estimate
    - E.g., Markov's inequality