# 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:

- *You*: The probability is:

- *Billionaire:* Why?

# Coin – Binomial Distribution

- **Data**: sequence *D= (HHTHT…),* **k heads** out of **n flips**
- **Hypothesis:** P(Heads) = θ,  P(Tails) = 1-θ
  - Flips are i.i.d.:
    - Independent events
    - Identically distributed according to Binomial distribution

- $P(\mathcal{D}|\theta) =$
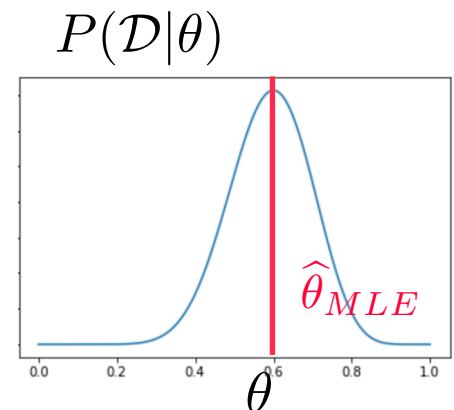
# Maximum Likelihood Estimation

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

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

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

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

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

# Your first learning algorithm

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

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

• Set derivative to zero:

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

# How many flips do I need?

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

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

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

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

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

- *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.*
- 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:

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

- **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

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)** $\widehat{\theta}_{MLE} = \arg\max_{\theta} L_n(\theta)$

# 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)
  - $X \sim N(\mu, \sigma^2)$
  - $Y = aX + b \quad \Rightarrow \quad Y \sim N(a\mu+b, a^2\sigma^2)$

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

# MLE for Gaussian

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

$$P(\mathcal{D}|\mu,\sigma) = P(x_1,\ldots,x_n|\mu,\sigma)$$

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

- Log-likelihood of data:

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

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

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

- What's MLE for mean?

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

# MLE for variance

- Again, set derivative to zero:

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

# 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**

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

- 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) $\quad \widehat{\theta}_{MLE} = \arg\max_{\theta} L_n(\theta)$

# Maximum Likelihood Estimation

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

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

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

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

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