

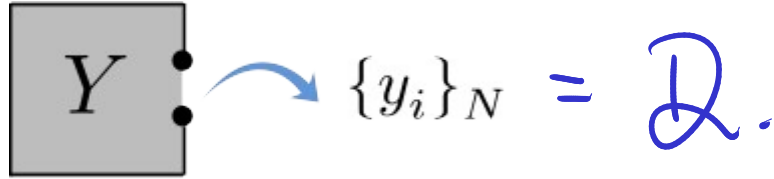


# Statistics and Data Science for Engineers E178 / ME276DS

Statistical inference:  
Point estimation

# Statistical inference

No inputs



Inference: Statement based on data.

Assumption: Sampling is iid.

↳  $Y$  does not change b/w samples  
↳ Samples are independent.

Given  $D$ .

## Three types of inferences

✓ 1) Point estimation  $\rightarrow$  : "My best guess for some parameter  $\theta$  of  $P_Y$  is  $\hat{\theta}_N$ "

✓ 2) Confidence intervals : "Parameter  $\theta$  lies in the interval  $I$  with confidence  $\gamma$ "

✓ 3) Hypothesis tests :  $H_0$  : \_\_\_\_\_  
 $H_1$  : \_\_\_\_\_  
null hypothesis

" $H_0$  is rejected in favor of  $H_1$ "

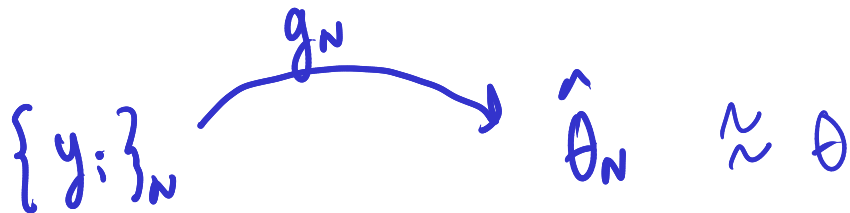
" $H_0$  is not rejected in favor of  $H_1$ "



# Point estimation

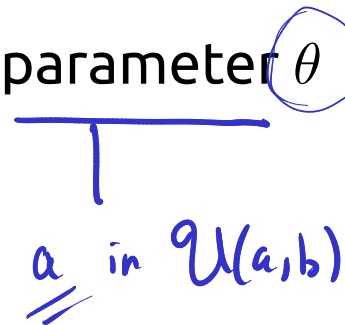
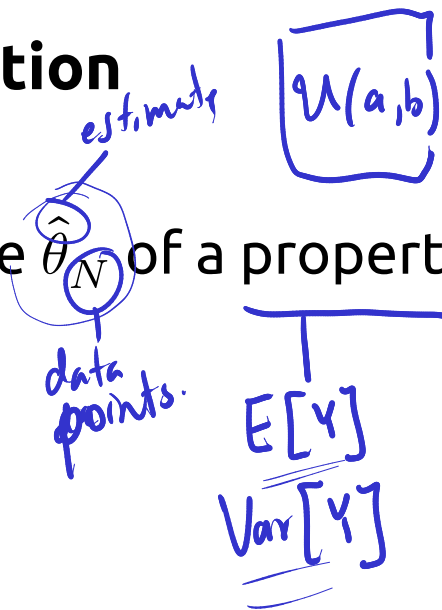
Given  $\mathcal{D} = \{\underline{y_i}\}_N \stackrel{\text{iid}}{\sim} Y$ , find a best estimate  $\hat{\theta}_N$  of a property or parameter  $\theta$

Estimator  $\hat{\theta}_N = g_N(y_1, \dots, y_N)$



Is  $g_N$  a good estimator?

Means that "expect  $\hat{\theta}_N \approx \theta$ "



Make an assumption about  $Y$

$$\{Y_i\}_n \sim \text{iid } Y$$

"Estimator"

$g_n$

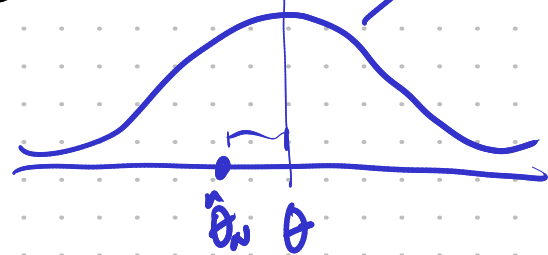
$\{Y_i\}_n$

$\hat{\Theta}_n$

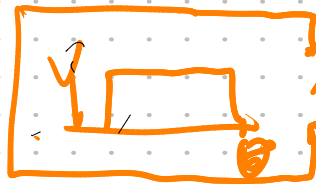


$E[\hat{\Theta}_n]$

$\hat{\Theta}_n$



$\hat{\Theta}_n$



$$\{y_1, \dots, y_{100}\} = D$$

$g_n^1$

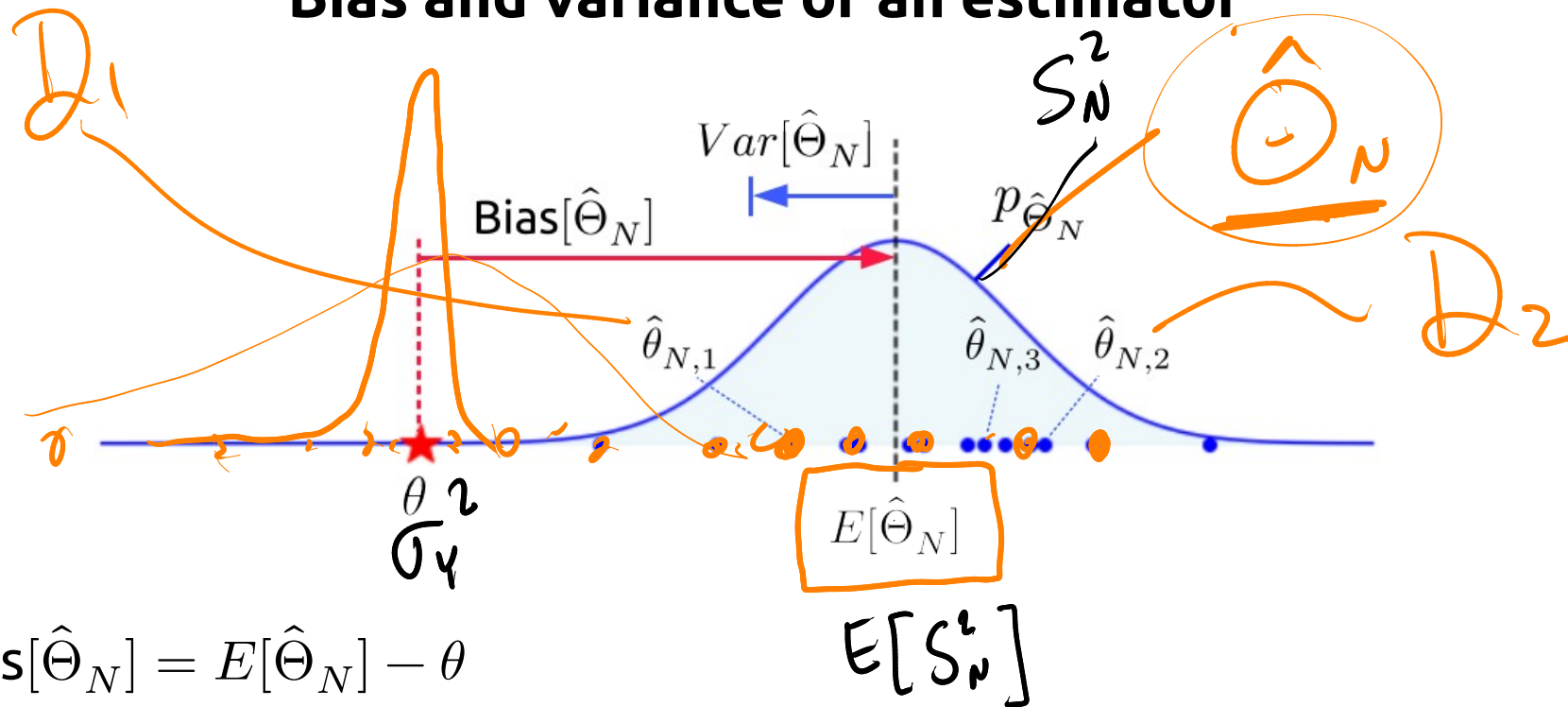
$$\max(y_1, \dots, y_{100}) = \hat{\Theta}_n$$

$g_n^2$

$$\max(y_1, \dots, y_{100}) + 10 \text{ cm.}$$

$\theta \dots$  size of largest apple in the orchard

# Bias and variance of an estimator



- $\text{Bias}[\hat{\Theta}_N] = E[\hat{\Theta}_N] - \theta$

- $\text{Var}[\hat{\Theta}_N] = E \left[ (\hat{\Theta}_N - E[\hat{\Theta}_N])^2 \right]$

## Estimating the mean

The sample mean:

$$\hat{\mu}_N = g_N(y_1, \dots, y_N) = \frac{1}{N} \sum_{i=1}^N y_i$$
$$\bar{Y}_N = g_N(Y_1, \dots, Y_N) = \frac{1}{N} \sum_{i=1}^N Y_i$$

$$\begin{aligned} E[\bar{Y}_N] &= E\left[\frac{1}{N} \sum_{i=1}^N Y_i\right] = \frac{1}{N} \sum E[Y_i] \\ &= \frac{1}{N} \sum E[Y] = \frac{1}{N} N E[Y] = E[Y] = \mu_Y \end{aligned}$$

## Estimating the mean

✓ •  $\text{Bias} [\bar{Y}_N] = 0$

Proof

✓ •  $\text{Var} [\bar{Y}_N] = \frac{\sigma_Y^2}{N}$

Prove:  $\text{Var} [\bar{Y}_N] = \text{Var} \left[ \frac{1}{N} \sum Y_i \right]$

$$\sim \left( \frac{1}{N} \right)^2 \sum \underbrace{\text{Var} [Y_i]}_{\sigma_Y^2} \sim \frac{N \sigma_Y^2}{N^2} = \frac{\sigma_Y^2}{N}$$



$$N=1 \quad (\bullet)$$

$$\bar{Y}_N = Y$$

$$\text{Var}[\bar{Y}_N] = \sigma_Y^2$$

$$N=10 \quad (\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet)$$

$$\text{Var}[\bar{Y}_N] = \frac{\sigma_Y^2}{10}$$

Var :  $\frac{\sigma_Y^2}{N}$

Std :  $\frac{\sigma_Y}{\sqrt{N}}$

# Estimating the variance

$$\sigma_Y^2$$

Unbiased sample variance:

$$\hat{\sigma}_N^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \hat{\mu}_N)^2$$

$$S_N^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y}_N)^2$$

$$\text{Var}[Y] \approx \frac{1}{N} \sum (y - \hat{\mu}_N)^2$$

## Estimating the variance

✓ •  $\text{Bias}[S_N^2] = 0$

Proof in the reader.

✗ •  $\text{Var}[S_N^2] =$  complicated.

$Y$  is Gaussian.

⇓

$S_N^2 \sim \chi^2$  distribution

Biased sample variance:

$$\tilde{S}_N^2 = \frac{1}{N} \sum_{i=1}^N (Y_i - \bar{Y}_N)^2$$

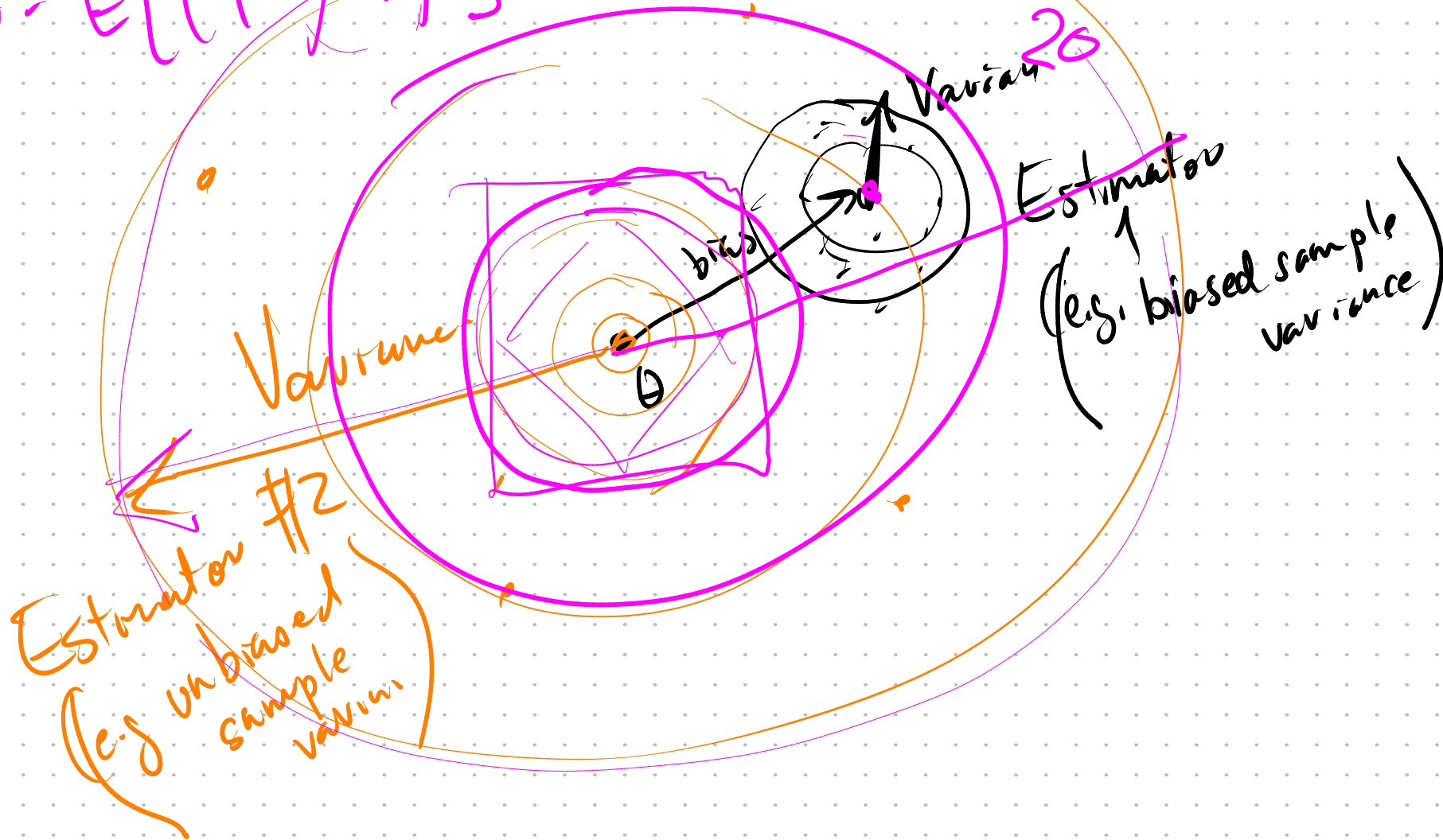
$$\bullet \text{ Bias}[\tilde{S}_N^2] = E[\tilde{S}_N^2] - \sigma_Y^2 = \frac{N-1}{N} \sigma_Y^2 - \sigma_Y^2 = -\frac{1}{N} \sigma_Y^2$$

$$\tilde{S}_N^2 = \frac{N-1}{N} S_N^2 \rightarrow E[\tilde{S}_N^2] = \frac{N-1}{N} E[S_N^2] = \frac{N-1}{N} \cdot \sigma_Y^2$$

$$\bullet \text{ Var}[\tilde{S}_N^2] = \text{complicated.}$$

smaller than  $\text{Var}[S_N^2]$ .

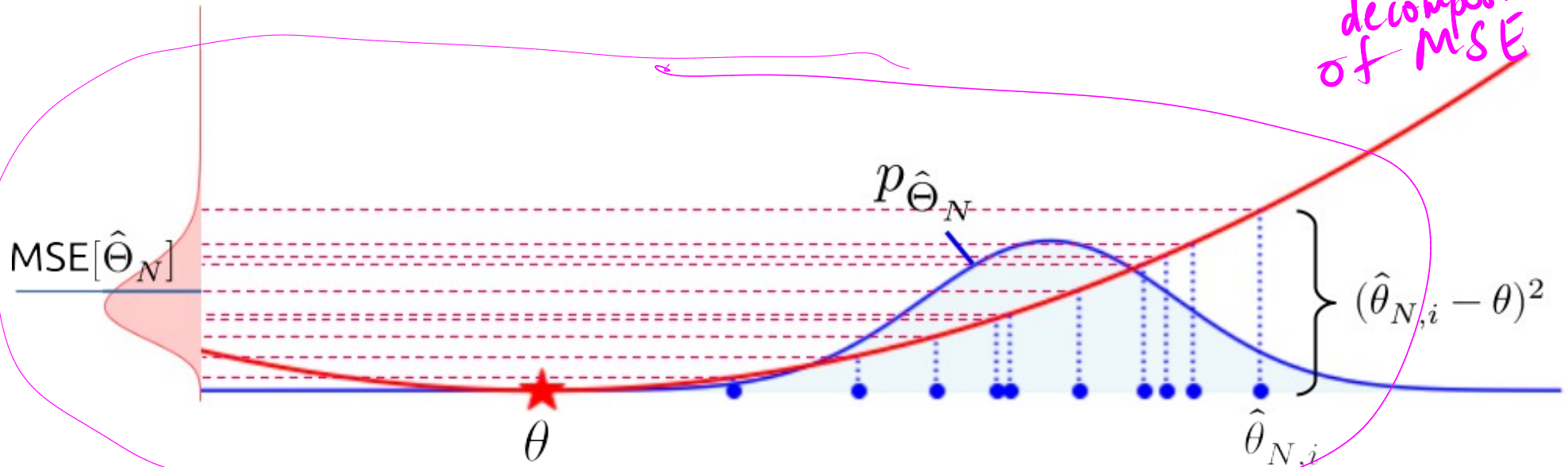
$$\text{Var} = E[(Y - \mu_Y)^2]$$



$Var[\hat{\Theta}_N] = E[(\hat{\Theta}_N - E[\hat{\Theta}_N])^2]$  **Mean squared error (MSE)**

$$\begin{aligned} \text{MSE}[\hat{\Theta}_N] &= E[(\hat{\Theta}_N - \theta)^2] \\ &= Var[\hat{\Theta}_N] + (\text{Bias}[\hat{\Theta}_N])^2 \end{aligned}$$

prove this  
Bias / Variance  
decomposition  
of MSE



sample

$$\text{MSE}[\bar{Y}_N] = \text{Var}[\bar{Y}_N] + (\text{Bias}[\bar{Y}_N])^2$$

$$\frac{\sigma_Y^2}{N} + 0 = \frac{\sigma_N^2}{N}.$$

$$\text{MSE}[S_N^2]$$



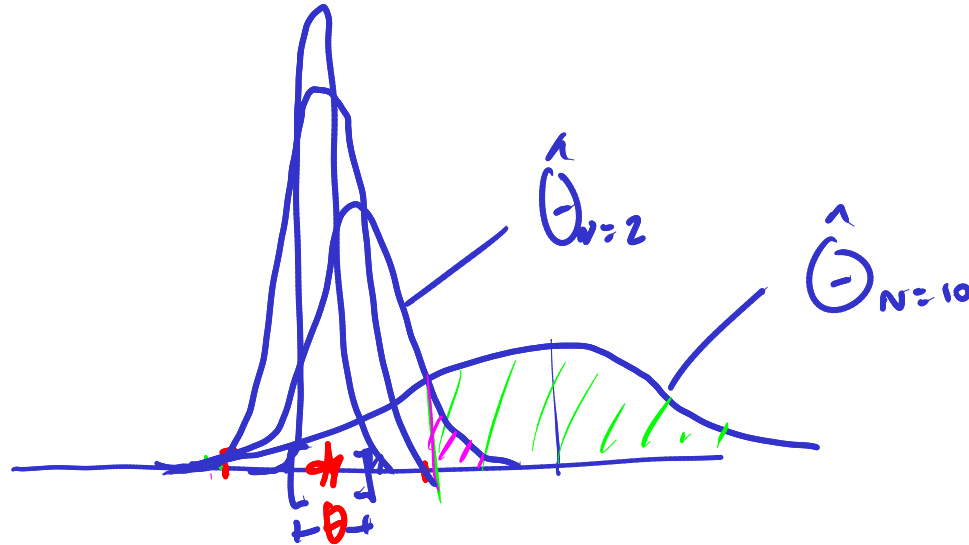
$$\text{MSE}[\tilde{S}_N^2]$$

# Asymptotic properties

- Asymptotic unbiasedness:  $\lim_{N \rightarrow \infty} \text{Bias}[\hat{\Theta}_N] = 0$
- Consistency:  $\lim_{N \rightarrow \infty} P(|\hat{\Theta}_N - \theta| \geq \epsilon) = 0$

$$\epsilon = 0.1$$

$$\epsilon = 0.01$$



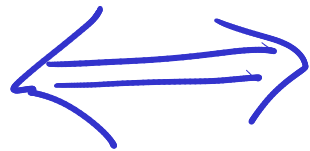


$$MSE = Var + Bias^2.$$

As  $N \rightarrow \infty$

$$Var \rightarrow 0$$

$$Bias \rightarrow 0$$



$$MSE \rightarrow 0.$$

if  $MSE \rightarrow 0$   
as  $N \rightarrow \infty$



consistent.

MSE is "stronger"  
than consistent

consistent



$$MSE \rightarrow 0$$

$$\text{as } N \rightarrow \infty.$$

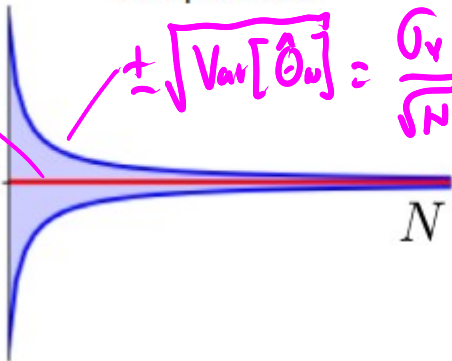
...  
"usually" it does.

$$E[\bar{Y}_n]$$

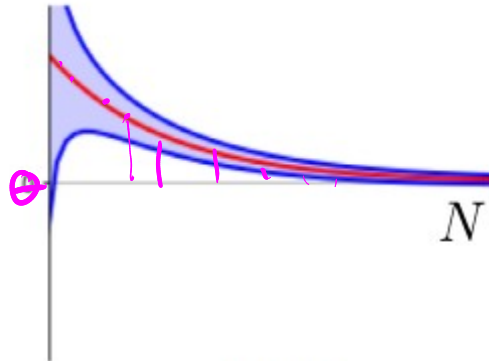
Sample mean

$$\pm \sqrt{\text{Var}[\bar{\theta}_n]} = \frac{\sigma_v}{\sqrt{n}}$$

$\mu_v$



Unbiased,  
Consistent



Biased ,  
Asymptotically unbiased,  
Consistent



Unbiased,  
Inconsistent

# Maximum Likelihood Estimation (MLE)

$$\hat{\theta}_N = g_N(\mathcal{D}) \quad \dots \text{general point estimation}$$

- MLE:

1) Pick a parametrization.

2) Solve:  $\hat{\underline{\theta}}_{\text{MLE}} = \underset{\underline{\theta}}{\operatorname{argmax}} \mathcal{L}(\underline{\theta}; \mathcal{D})$

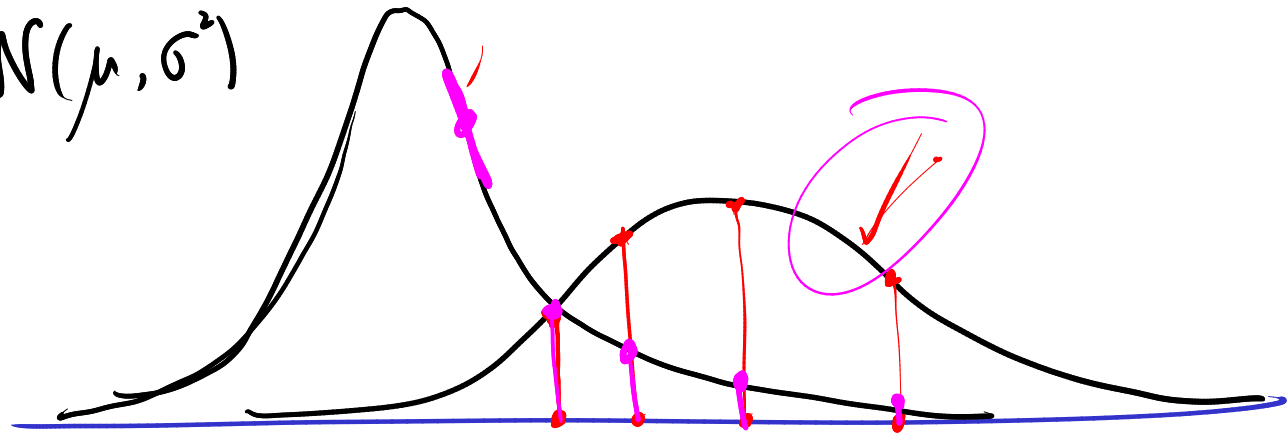
- Likelihood:  $\mathcal{L}(\underline{\theta}; \mathcal{D}) = \prod_{i=1}^N p_Y(y_i; \underline{\theta})$

$Y \sim \mathcal{N}(\mu, \sigma^2), Y \sim \mathcal{B}(p), \dots$

$\underline{\theta} = (\mu, \sigma^2), \quad \underline{\theta} = p.$

Likelihood: Probability of obtaining the dataset  $\mathcal{D}$  if parameters are  $\underline{\theta}$

$$Y \sim N(\mu, \sigma^2)$$



## Example

- 4 marbles in a bag, all either black or white
- pick 5 times with replacement  $\mathcal{D} = \{1, 0, 1, 0, 0\}$

black white.



Estimate the number of black marbles in the bag.

$$\theta \dots \{0, 1, 2, 3, 4\}.$$

$$p = \theta/4 = \{0, 1/4, 1/2, 3/4, 1\}.$$

1. Pick a family:  $Y \sim \mathcal{B}(p)$

2. Solve:  $\underset{\theta}{\text{maximize}} \prod_{i=1}^5 p_Y(\theta; y_i)$

$$p_Y(\theta; y_i) = p^{y_i} (1-p)^{1-y_i} = \begin{cases} p = \theta/4 & y_i = 1 \\ 1-p = 1 - \frac{\theta}{4} & y_i = 0. \end{cases}$$

$$\underset{\theta}{\text{maximize}} \prod_{i=1}^5 p_Y(\theta; y_i) = \underset{\theta}{\text{maximize}} \overbrace{\left(\frac{\theta}{4}\right)^2 \left(1 - \frac{\theta}{4}\right)^3}^{J(\theta)}.$$

$\theta$	0	1	2	3	4
$J(\theta)$	0	$\frac{2^2}{4^5}$	$\frac{3^2}{4^5}$	$\frac{9}{4^5}$	0

$\frac{1}{4^2} \cdot \frac{3^3}{4^3}$

$\frac{2^2 \cdot 2^3}{4^5}$

$\frac{3^2 \cdot 1^3}{4^5}$

# Log-likelihood

$$\begin{aligned}\hat{\underline{\theta}}_{\text{MLE}} &= \operatorname{argmax}_{\underline{\theta}} \mathcal{L}(\underline{\theta}; \mathcal{D}) \\ &= \operatorname{argmax}_{\underline{\theta}} \ln \mathcal{L}(\underline{\theta}; \mathcal{D}) \\ &= \operatorname{argmax}_{\underline{\theta}} \ln \left( \prod_{i=1}^N p_Y(y_i; \underline{\theta}) \right) \\ &= \operatorname{argmax}_{\underline{\theta}} \sum_{i=1}^N \ln p_Y(y_i; \underline{\theta})\end{aligned}$$

## Example: Gaussian data

Assume:  $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$  ,  $\underline{\theta} = (\underline{\mu}_Y, \underline{\sigma}_Y^2)$

$$\begin{aligned}(\hat{\mu}_{\text{MLE}}, \hat{\sigma}_{\text{MLE}}^2) &= \underset{\mu, \sigma^2}{\operatorname{argmax}} \sum_{i=1}^N \ln p_Y(y_i; \mu, \sigma^2) \\&= \underset{\mu, \sigma^2}{\operatorname{argmax}} \sum_{i=1}^N \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{1}{2} \frac{(y_i - \mu)^2}{\sigma^2} \right) \right) \\&= \underset{\mu, \sigma^2}{\operatorname{argmax}} \left( \overbrace{-\frac{N}{2} \ln(2\pi\sigma^2)} - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right) \\&= \underset{\mu, \sigma^2}{\operatorname{argmin}} \left( \frac{N}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right)\end{aligned}$$



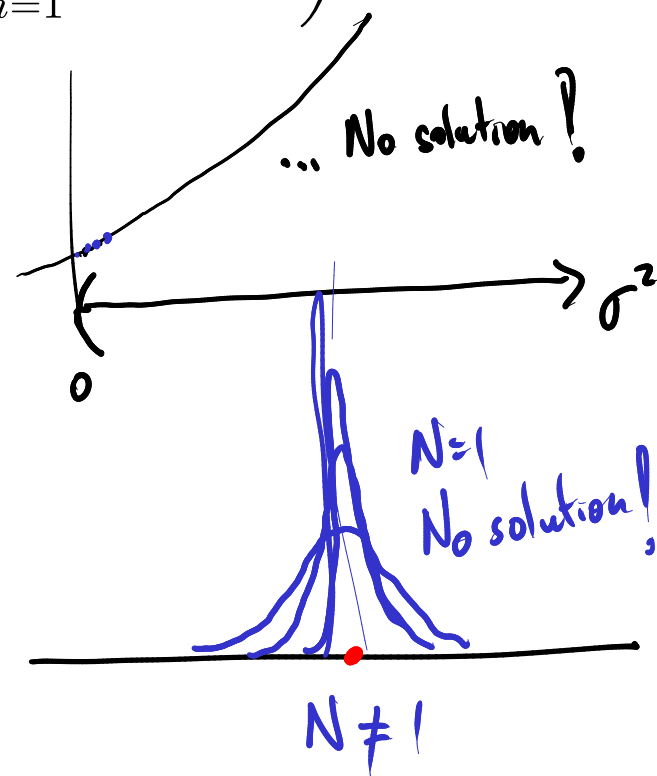
## Example: Gaussian data

$$(\hat{\mu}_{\text{MLE}}, \hat{\sigma}_{\text{MLE}}^2) = \underset{\mu, \sigma^2}{\operatorname{argmin}} \left( \frac{N}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right) \dots \text{Convex} \checkmark$$

subject to:  $\sigma^2 > 0$ .

$$J(\mu, \sigma^2) = \frac{N}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2$$

$$\nabla J(\mu, \sigma^2) = \left( \frac{\partial J}{\partial \mu}, \frac{\partial J}{\partial \sigma^2} \right) = 0$$



## Example: Gaussian data

$$\begin{aligned}\frac{\partial J}{\partial \mu} &= \frac{\partial}{\partial \mu} \left( \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right) \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^N \frac{\partial}{\partial \mu} (y_i - \mu)^2 \\ &= -\frac{1}{\sigma^2} \sum_{i=1}^N (y_i - \mu) \\ &= \frac{N\mu}{\sigma^2} - \frac{1}{\sigma^2} \sum_{i=1}^N y_i = 0\end{aligned}$$

$$N\mu = \sum_{i=1}^N y_i$$

$$\therefore \hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N y_i = \hat{\mu}_w!$$

## Example: Gaussian data

$$\frac{\partial J}{\partial \sigma^2} = \frac{\partial}{\partial \sigma^2} \left( \frac{N}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right)$$

$$= \frac{N}{2\sigma^2} - \frac{1}{2\sigma^4} \sum_{i=1}^N (y_i - \mu)^2$$

$$= \frac{1}{2\sigma^2} \left( N - \frac{1}{\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right) = 0$$

$$\sigma_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\mu}_N)^2 \dots \text{biased sample variance!}$$

# Properties of MLE

*N is finite.*

- MLE has no finite-sample properties.
  - not necessarily unbiased
  - not necessarily minimum MSE.
- MLE has good asymptotic properties.
  - consistent
  - usually asymptotically unbiased ..