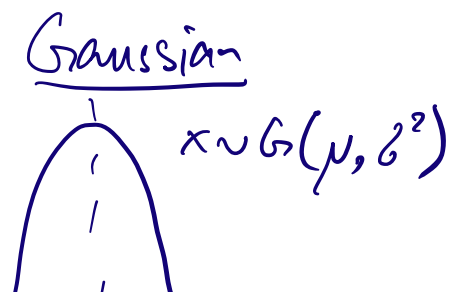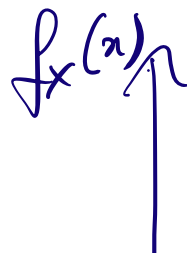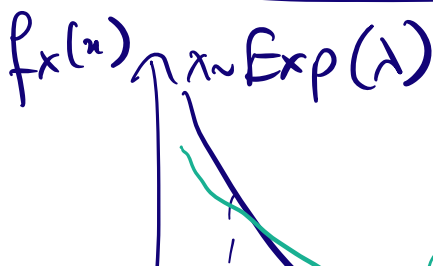# Lecture 14

- Moments
- KDE
- Statistical Inference
- Hypothesis tests
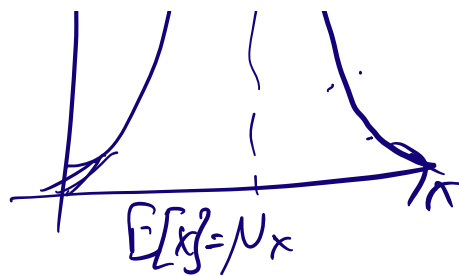
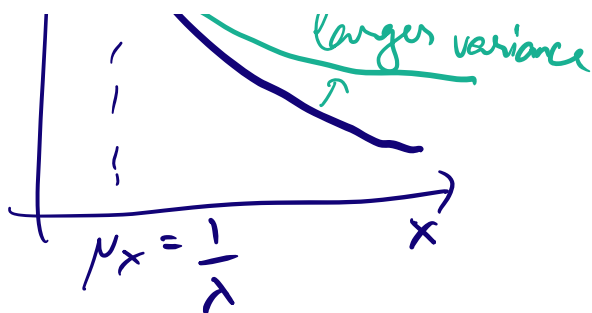$n^{th}$ central moment (discrete RV):

$$E\left[(x - \mu_x)^n\right] = \sum_x (x - \mu_x)^n \cdot p_x(x)$$
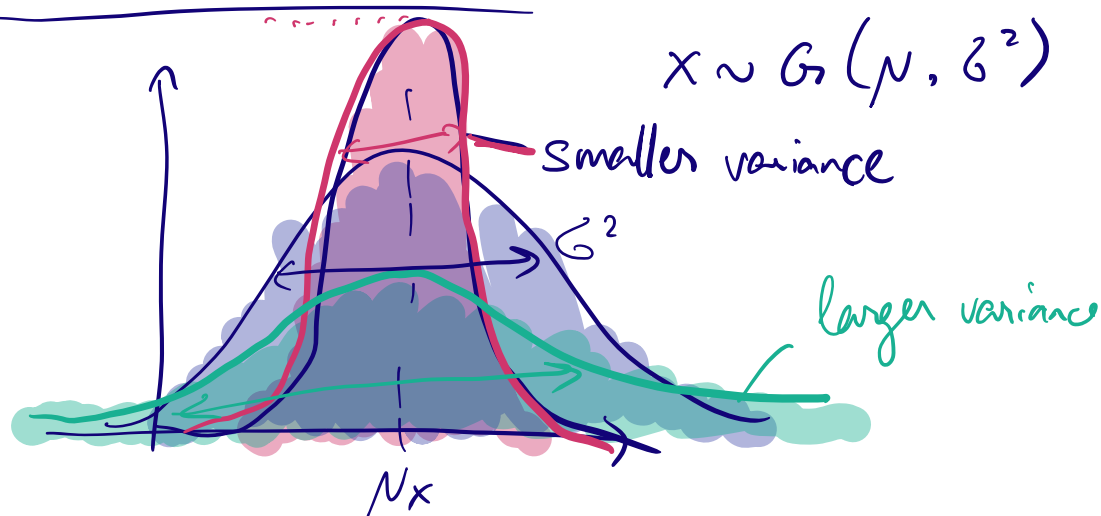
$1^{st}$ central moment $\equiv$ mean

$$\boxed{E[x] = \mu_x \neq \text{mean}}$$

$f_x(x)$    $x \sim Exp(\lambda)$

$f_x(x)$    Gaussian    $x \sim G(\mu, \sigma^2)$

$\mu_x = \dfrac{1}{\lambda}$

$E[x] = \mu_x$

---

$2^{nd}$ Central Moment $\equiv Var[x]$

---

$X \sim G(\mu, \sigma^2)$

smaller variance

$\sigma^2$

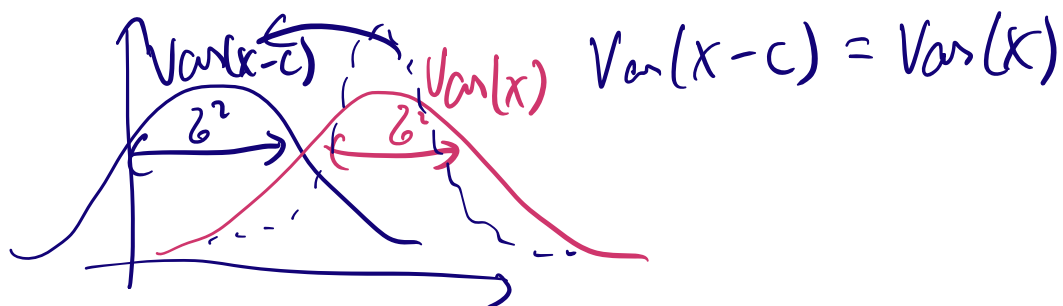larger variance

$\mu_x$

$$Var[x] = E\left[(X - \mu_x)^2\right], \text{ where } \mu_x = E[x]$$

$$= E\left[X^2 - 2\mu x + \mu_x^2\right]$$

$$= E[x^2] - 2\mu_x \underbrace{E[x]}_{= \mu_x} + \mu_x^2$$

$$= E[x^2] - \mu_x^2$$

$$= E[x^2] - (E[x])^2$$

$Var(x-c)$    $Var(x)$    $Var(x-c) = Var(x)$

$\sigma^2$    $\sigma^2$

---

## $3^{rd}$ central moments $\equiv$ skewness

$$E[(x - \mu_x)^3]$$

$G(\mu, \sigma^2)$

$\mu_x$

symmetric w.r.t $\mu_x$

skewness $= 0$

$x \sim Exp(\lambda)$

$\mu_x$

long right tail

skewness $> 0$

$\mu_x$

$x \sim Binomial(15, 0.9)$

Long left tail

$$\frac{\cdots}{\sigma} \cdots$$

Skewness $< 0$

---

$4^{th}$ central moment $\equiv$ kurtosis

$$E[(X - \mu_X)^4]$$

↗ area under the curve

measures the "volume" of the tails in the distribution



$X \sim G(0,1) \Rightarrow \text{kurtosis}(X) = 3$

kurtosis $> 3$

kurtosis $< 3$

$$E[(X - \mu_X)^4] - 3 \equiv \text{"Excess" kurtosis}$$

$$= 0 \text{ for } G(0,1)$$

---

Statistical Inference

$f_X(x) \uparrow$ ⟍  $X \sim \text{Exp}(\lambda)$

$\lambda$ is a <u>parameter</u> that we can <u>infer</u> from data !
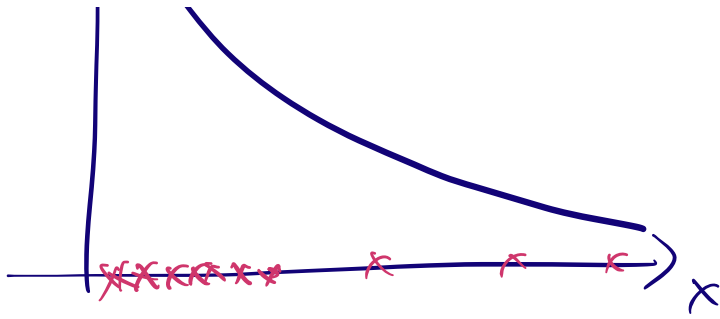
$\theta \equiv$ parameters $\leftarrow$ deterministic but unknown!

$f_x(x; \theta) \equiv$ PDF with unknown parameters $\theta$

If $x \sim Exp(\lambda)$, then $\theta = \lambda$

If $x \sim G(\mu, \sigma^2)$, then $\theta = \{\mu, \sigma\}$

$\hat{\theta} \equiv$ estimate

$\hat{\theta}$ is estimated from a data sample

$$\{x_i\}_{i=1}^{N} = \{x_1, \dots x_N\}$$

Error of that estimator:

$$\mathcal{E}_\theta(\hat\theta) = \hat\theta - \theta$$

Bias of estimator

$$b_\theta(\hat\theta) = E[\hat\theta] - \theta$$

If $E[\hat\theta] = \theta$
then $b_\theta(\hat\theta) = 0$
$\Downarrow$
unbiased
estimator!

Variance of estimator

$$Var_\theta[\hat\theta] = E[(\hat\theta - E[\hat\theta])^2]$$



$$E[\hat\mu] = \frac{1}{N}\sum_{i=1}^{N}\hat\mu_i \to \mu_x$$

Mean - Square error (MSE)

$$E[(\theta - \hat\theta)^2] = b_\theta^2(\hat\theta) + Var_\theta(\hat\theta)$$

Data samples

$$x = \{x_1, \ldots, x_N\} = \{x_i\}_{i=1}^{N}$$

sample size $= N$

Draw from the same underlying distribution **independently**:

independent & identically distributed

$$(i.i.d.)$$

Estimator for the mean:

True mean $= \mu_x$

Estimate $= \hat{\mu}_x$

$$\boxed{\hat{\mu}_x = \frac{1}{N} \sum_{i=1}^{N} x_i} \implies \hat{\mu}_x \text{ is an unbiased estimator}$$

**Proof:** of $\mu_x$

$$E[\hat{\mu}_x] = E\left[\frac{1}{N}\sum_{i=1}^{N}x_i\right] = \frac{1}{N}E\left[\sum_{i=1}^{N}x_i\right]$$

$$= \frac{1}{N}E\{x_1 + x_2 + \ldots + x_N\}$$

$$= \frac{1}{N}\sum_{i=1}^{N}\underbrace{E[x_i]}_{=\mu_x} = \frac{1}{N}\sum_{i=1}^{N}\mu_x = \mu_x$$

**Estimator for the variance:**

True variance is $\sigma_x^2$

① 
$$\boxed{S_N^2 = \frac{1}{N}\sum_{i=1}^{N}(x_i - \mu_x)^2}$$

1$^{st}$ estimator for variance

variance for R.V. x:

$$\sigma_x^2 = E[(x-\mu_x)^2] = \frac{1}{N}\sum_{i=1}^{N}(x_i - \mu_x)^2$$

$$\boxed{E[S_N^2] = \frac{N-1}{N}\cdot\sigma_x^2} \neq \sigma_x^2$$

↳ biased estimator!

Note that $N-1 \xrightarrow{N\to\infty} 1$

$$\overline{N} \longrightarrow \downarrow$$

$$\downarrow$$

Sample size N is large, then we may get unbiased estimator

② $\quad \dfrac{N}{N-1} \ S_N^2 \ = \ \boxed{S_{N-1}^2 \ = \ \dfrac{1}{N-1} \sum_{i=1}^{N} (x_i - \mu_x)^2}$

2nd estimator for variance

$$E\left[ S_{N-1}^2 \right] = \sigma_x^2 \implies \text{unbiased estimator!}$$

---

## Properties of sum of independent Gaussian R.V.s

$X = \{x_i\}_{i=1}^{N}$ $\qquad\qquad$ $Y = \{y_i\}_{i=1}^{N}$

$x \sim G(\mu_x, \sigma_x^2)$ $\qquad\qquad$ $Y \sim G(\mu_Y, \sigma_y^2)$

$$z = x + y$$

If $x$ & $y$ are independent RVs :

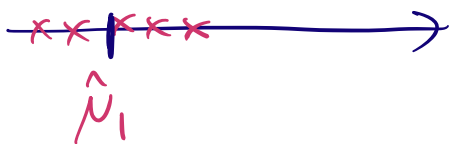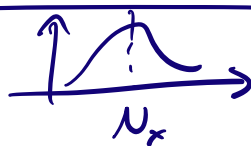$$\mu_z = \mu_x + \mu_y$$

$$\sigma_z^2 = \sigma_x^2 + \sigma_y^2$$

$$z \sim G(\mu_z, \sigma_z^2)$$

$z' = az + b$ is a Gaussian R.V.

$$\mu_{z'} = a\mu_z + b$$

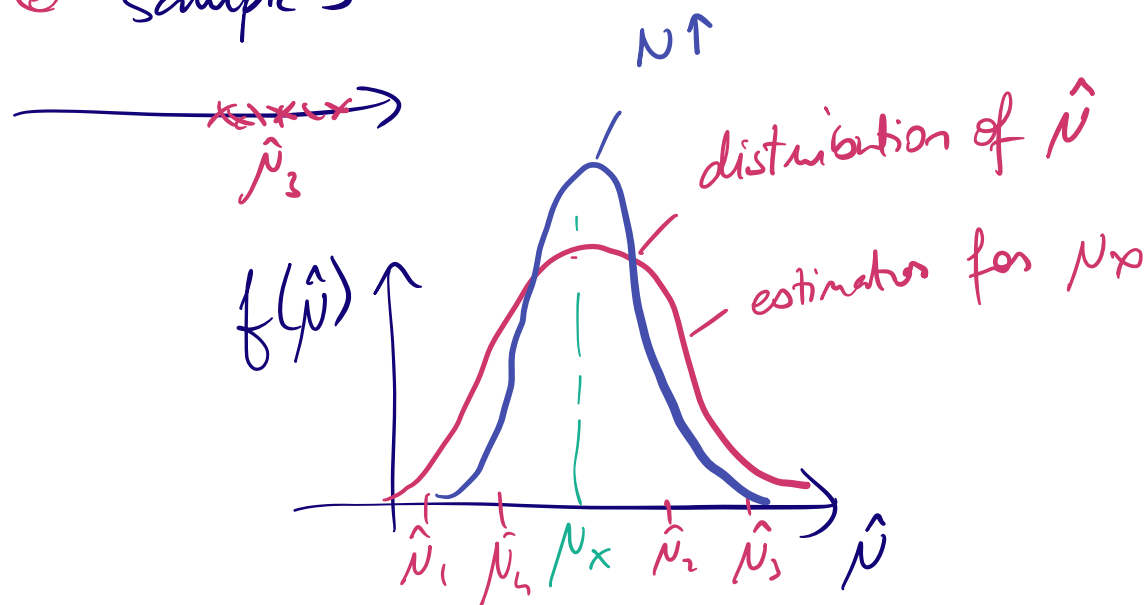$$\sigma_{z'}^2 = a^2 \sigma_z^2$$

---

① sample 1



$\mu_x$



$\hat{\mu}_1$

② sample 2



$$\hat{\mu}_i = \frac{1}{N} \sum_{j=1}^{N} x_j$$

$N_2$

② sample 3

$\hat{N}_3$

$f(\hat{\mu})$

$N\uparrow$

distribution of $\hat{N}$

estimator for $N_x$

$\hat{N}_1$ $\hat{N}_4$ $N_x$ $\hat{N}_2$ $\hat{N}_3$ $\hat{N}$

$$\hat{N}_x \sim G\left(N_x, \frac{\sigma_x^2}{N}\right)$$

$$Var[\hat{N}_x] = Var\left[\frac{1}{N}\sum_{i=1}^{N} x_i\right]$$

$$= \frac{1}{N^2} Var\left[\sum_{i=1}^{N} x_i\right]$$

$$= \frac{1}{N^2} \cdot Var[x_1 + \dots + x_N]$$

$$= \frac{1}{N^2}\sum_{i=1}^{N} \underbrace{Var[x_i]}_{\sigma_x^2}$$

$var(a\cdot x)$
$= a^2 var(x)$
for $a$ constant

$$= \frac{1}{N^2} \cdot N \cdot \sigma_x^2$$

$$= \frac{1}{N} \cdot \sigma_x^2 \xrightarrow[N \to \infty]{} 0$$

# Hypothesis Testing:

## Z-test

$$x = \{x_i\}_{i=1}^{N} \qquad Y = \{y_i\}_{i=1}^{M}$$

→ N samples for x          M samples for Y

Consider the estimator for the mean of these two sampled data:

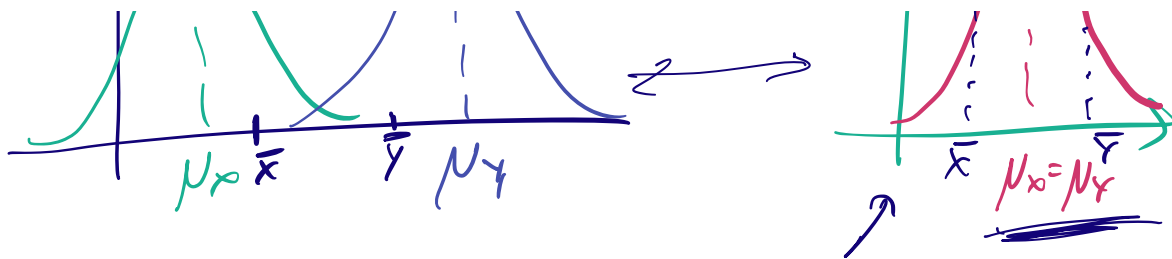$$\bar{x} = \frac{1}{N} \sum_{i=1}^{N} x_i \qquad \text{and} \qquad \bar{y} = \frac{1}{M} \sum_{j=1}^{M} y_j$$

If $\bar{x} \neq \bar{y} \implies \mu_x \neq \mu_y$

z-test checks if two sampled data are drawn from the same distribution by measuring the difference of their mean estimator

It assumes that the variance for the mean estimators are known:

$$\mathrm{Var}[x] = \mathrm{var}[Y] = \sigma^2$$

$$\boxed{T = \hat{N}_x - \hat{N}_y} \leftarrow \text{this is also Gaussian distributed}$$

$$t = \bar{x} - \bar{y}$$

under the assumption that $N_x = N_y = N$

① $\underline{E[T]} = N_T = E[\hat{N}_x - \hat{N}_y]$

$$= E[\hat{N}_x] - E[\hat{N}_y]$$

$$= N - N$$

$$= 0$$

② $\sigma_T^2 = \text{Var}[T]$

$$= \text{Var}[\hat{\mu}_x - \hat{\mu}_y]$$

$$= \text{Var}[\hat{\mu}_x + (-1) \cdot \hat{\mu}_y] \quad \Big\} \text{add}$$

$$= \text{Var}[\hat{\mu}_x] + \text{Var}[(-1) \cdot \hat{\mu}_y]$$

$$\quad\quad\quad \Big\} \begin{array}{l} \text{Var}(cx) \\ = c^2 \text{Var}(x) \end{array}$$

$$= \text{Var}[\hat{\mu}_x] + \text{Var}[\hat{\mu}_y]$$

$$= \frac{\sigma_x^2}{N} + \frac{\sigma_y^2}{M} = \sigma^2 \left( \frac{1}{N} + \frac{1}{M} \right)$$

assumption:
known variances
for x & Y
⇓
(z-test)

assumption:
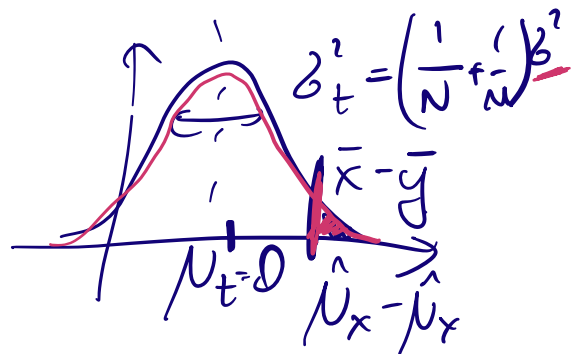equal variances
for x and Y

## Hypothesis Test :

$H_0$ : the true means are the same

$$\mu_x = \mu_y$$

$H_1$ : $\mu_x \neq \mu_y$

Steps:

① compute the statistic : $t = \hat{\mu}_x - \hat{\mu}_y$

② we know that $t$ is Gaussian distributed ; we know the mean & variance of this distribution under the null hypothesis.
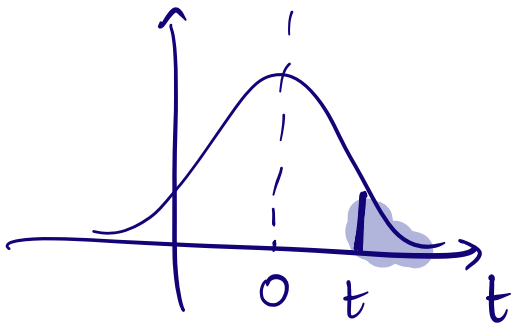
③ compute p-value

$$\sigma_t^2 = \left(\frac{1}{N} + \frac{1}{N}\right)\beta^2$$

$\bar{x} - \bar{y}$

$\mu_{t=0}$    $\hat{\mu}_x - \hat{\mu}_y$

$H_0$ :

$$T = \hat{\mu}_x - \hat{\mu}_y$$

$$T \sim G\left(0, \; \sigma^2\left(\frac{1}{N} + \frac{1}{M}\right)\right)$$

$$P(T \geqslant t \mid H_0) = P(T > t \mid H_0)$$
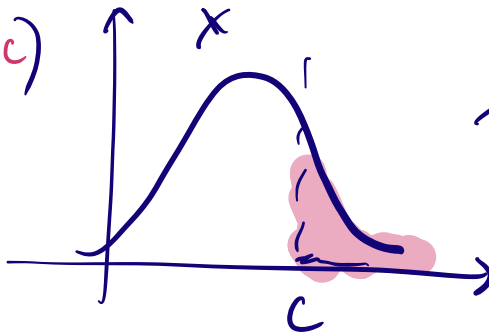
$$= Q\left(\frac{t - \mu_T}{\sigma_T}\right)$$

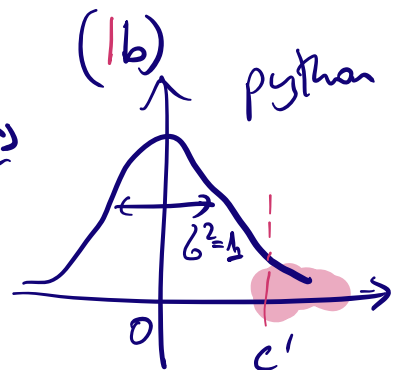$$= Q\left(\frac{t}{\sqrt{\sigma^2\left(\frac{1}{N} + \frac{1}{M}\right)}}\right)$$

(1c)

(1b) python

transform

$\sigma^2 = 1$

c

c'

(1a)  Q-function
table