

Statistical Analysis of Data


Prof. Americo Cunha Jr

americocunha.org



 @AmericoCunhaJr

 @AmericoCunhaJr

 @AmericoCunhaJr

Estimators for statistical moments

X_1, \dots, X_n are independent observations of X

- Sample mean

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$$

- Sample skewness

$$\hat{\gamma}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^3}{\left(\frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu})^2 \right)^{3/2}}$$

- Sample variance

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

- Sample kurtosis

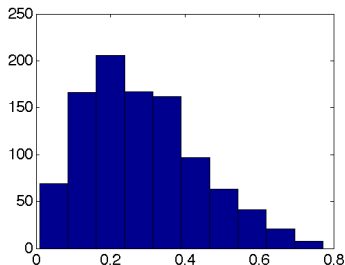
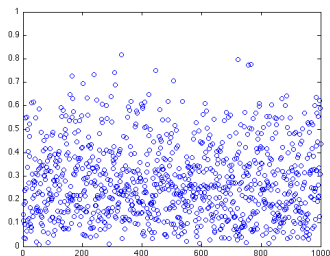
$$\hat{\beta}_2 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^4}{\left(\frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^2 \right)^2}$$

main_data_analysis1.m (1/2)

```
1  clc; clear all; close all;
2  a = 2; b=5; Ns = 1000;
3
4  X      = betarnd(a,b,Ns,1);
5  mu     = mean(X)
6  sigma2 = var(X)
7  sigma  = std(X)
8  gamma1 = skewness(X)
9  beta2  = kurtosis(X)
10
11 figure(1)
12 plot(1:Ns,X,'o')
13 ylim([0 1]);
14
15 figure(2)
16 hist(X)
17 xlim([0 1]);
```

Statistical analysis in Matlab/Octave

```
mu = 0.2849  
sigma2 = 0.0234  
sigma = 0.1528  
gamma1 = 0.5900  
beta2 = 2.8555
```



main_data_analysis2.m (2/2)

```
1  clc; clear all; close all;
2  a = 2; b=5; Ns = 1000;
3
4  rng_stream = RandStream('mt19937ar','Seed',30081984);
5  RandStream.setGlobalStream(rng_stream); % Matlab 2013
6
7  X      = betarnd(a,b,Ns,1);
8  mu     = mean(X)
9  sigma2 = var(X)
10 sigma  = std(X)
11 gamma1 = skewness(X)
12 beta2  = kurtosis(X)
13
14 figure(1)
15 plot(1:Ns,X, 'o')
16 ylim([0 1]);
17 figure(2)
18 hist(X)
19 xlim([0 1]);
```

Estimators for PDF and CDF

- Histogram

$$\hat{p}_n(x) = \sum_{m=-\infty}^{+\infty} \frac{\nu_m}{n h_m} \mathbb{1}_{\mathcal{B}_m}(x)$$

- Kernel Density Estimator

$$\hat{p}_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - X_i}{h}\right)$$

- Empirical CDF

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathcal{I}(X_i \leq x),$$

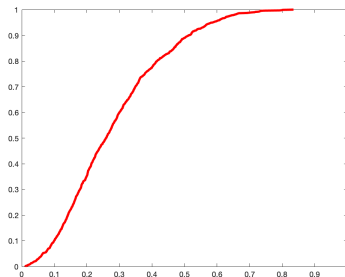
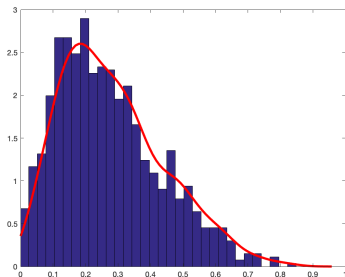
randvar_pdf.m

```
1 function [bins,freq,area] = randvar_pdf(data,numbins)
2
3     Ns = length(data);
4
5     data_max = max(data);
6     data_min = min(data);
7     binwidth = (data_max-data_min)/(numbins-1);
8
9     bins      = (data_min:binwidth:data_max);
10    freq       = histc(data,bins);
11    freq       = freq/(Ns*binwidth);
12    area       = binwidth*sum(freq);
13
14 end
```

main_histogram_ecdf.m

```
1  clc; clear; close all;
2
3  a = 2; b = 5; Ns = 1000;
4
5  X      = betarnd(a,b,Ns,1);
6  Nbins  = round(sqrt(Ns));
7
8  [X_bins,X_freq,X_area] = randvar_pdf(X,Nbins);
9  [X_ksd ,X_supp1       ] = ksdensity(X);
10 [X_cdf ,X_supp2       ] = ecdf(X);
11
12 figure(1)
13 bar(X_bins,X_freq,1.0);
14 hold on
15 plot(X_supp1,X_ksd,'r','linewidth',3)
16 xlim([0 1]);
17 hold off
18
19 figure(2)
20 plot(X_supp2,X_cdf,'r','linewidth',3)
21 xlim([0 1]); ylim([0 1]);
```


PDF and CDF estimation in Matlab



Construction of a confidence interval/envelope

- p -th quantile of distribution F_X

$$Q(p) = \inf \{x \in \mathbb{R} : p \leq F_X(x)\}, \quad 0 < p < 1$$

- envelope with probability P_c

$$\mathcal{P} \{r^- < X \leq r^+\} = P_c$$

$$r^+ = Q((1 + P_c)/2) \quad r^- = Q((1 - P_c)/2)$$

- estimation via percentiles

$X_1 < X_2 < \dots < X_n$ are independent observations of X

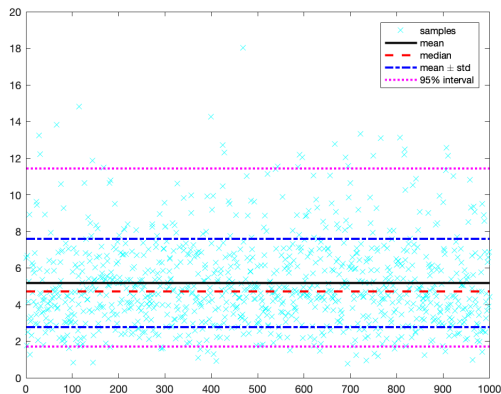
$$\hat{r}^+ = X_{n^+} \quad n^+ = \text{floor}(n(1 + P_c)/2)$$

$$\hat{r}^- = X_{n^-} \quad n^- = \text{floor}(n(1 - P_c)/2)$$

main_conf_interval.m

```
1  clc; clear; close all;
2
3  a = 5.0; b = 1.0; Ns = 1000; Pc = 95;
4
5  X = gamrnd(a,b,Ns,1);
6  mu = mean(X); sigma = std(X); mu50 = median(X);
7  r_plus = 0.5*(100 + Pc); r_minus = 0.5*(100 - Pc);
8  X_upp = prctile(X,r_plus); X_low = prctile(X,r_minus);
9
10 figure(1)
11 plot(X,'xc');
12 hold on
13 line([1 Ns],[mu mu], 'Color','k','LineStyle','-','linewidth',2);
14 line([1 Ns],[mu50 mu50], 'Color','r','LineStyle','--','linewidth',2);
15 line([1 Ns],[mu-sigma mu+sigma], 'Color','b','LineStyle','-.','linewidth',2);
16 line([1 Ns],[X_low X_low], 'Color','m','LineStyle',':','linewidth',2);
17 line([1 Ns],[mu+sigma mu+sigma], 'Color','b','LineStyle','-.','linewidth',2);
18 line([1 Ns],[X_upp X_upp], 'Color','m','LineStyle',':','linewidth',2);
19 legend('samples','mean','median','mean \pm std','95% interval')
20 hold off
```

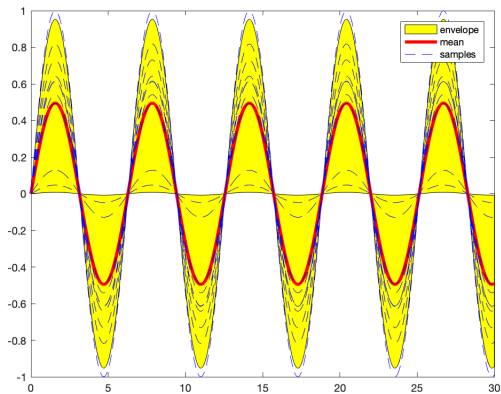
Confidence interval for a random variable



main_curve_envelope.m

```
1  clc; clear; close all;
2
3  Ns = 50; Pc = 95;
4
5  A = rand(Ns,1); x = 0:0.01:30; Y = A*sin(x);
6  r_plus = 0.5*(100 + Pc); r_minus = 0.5*(100 - Pc);
7  Y_upper = prctile(Y,r_plus); Y_low = prctile(Y,r_minus);
8
9  figure(1)
10 fh1 = plot(x,mean(Y),'r','linewidth',3);
11 hold on
12 fh2 = plot(x,Y(1:10,:),'--b','linewidth',0.5);
13 fh3 = fill([x fliplr(x)],[Y_upper fliplr(Y_low)],'y');
14 uistack(fh3,'top');
15 uistack(fh1,'top');
16 uistack(fh2,'top');
17 legend('envelope','mean','samples')
18 hold off
```

Confidence envelope for a random curve



References



V. Stodden, *Reproducing Statistical Results*, **Annual Review of Statistics and Its Application**, 2:1–19, 2015.



L. Wasserman, **All of Statistics: A Concise Course in Statistical Inference**, Springer, 2004.



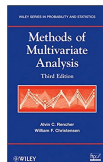
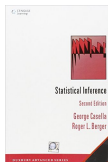
G. Casella, **Statistical Inference**, Thomson Press (India) Ltd, 2008.



J. Shao, **Mathematical Statistics**, Springer, 2nd Edition, 2007



A. C. Rencher, **Methods of Multivariate Analysis**, John Wiley & Sons, 3rd edition, 2012



How to cite this material?

A. Cunha Jr, *Statistical Analysis of Data*, 2021.



 @AmericoCunhaJr



@AmericoCunhaJr



@AmericoCunhaJr

These class notes may be shared under the terms of
Creative Commons BY-NC-ND 4.0 license,
for educational purposes only.

