

Week 2: Describing Probability Distributions and Probability Distributions with Multiple Variables

① Expected Value: It can mean many things: like $E[X]$

Measures of Center: → Mean (Average) = $\frac{\text{summation weighted values}}{\text{Number of weighted values}}$

- Mean
- Mode
- Median

Discrete: $E[X] = \sum_x x p_X(x)$ (PMF)

Continuous: $E[X] = \int_{-\infty}^{\infty} x \cdot f(x) dx$ (PDF)

Mean doesn't work when data contains outliers.

→ Median = middle value of ordered dataset.

If the number of values is even, then median = average of the two data in the middle

→ Mode = most frequent data in the dataset

Sometimes, we can have a multi-modal distribution; where many data have the same frequency that appears to be the highest.

→ Expected Value of a Function rather than random variables:

$$\bullet E[X] = x_1 p(x_1) + x_2 p(x_2) + \dots + x_n p(x_n)$$

$$\bullet E[g(X)] = g(x_1) p(x_1) + g(x_2) p(x_2) + \dots + g(x_n) p(x_n)$$

$$\bullet E[aX+b] = a E[X] + b$$

$$\bullet E[X_1 + X_2] = E[X_1] + E[X_2]$$

Measures of Spread:

Variance

Standard deviation

Since deviation ($x - E[X]$) can give wrong indications because of negative and positive signs. To measure variance, we get the average of squared deviations of all values $(x - E[X])^2$.

average value or expected value

$$\text{Variance} = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

same ↪ $\text{Var}(X) = E[(X - \bar{X})^2]$

$$\text{Var}(aX + b) = a^2 \text{Var}(X)$$

④ Standard Deviation:

To cover a drawback of the variance, which is having the units squared (not practical), we use standard deviation as the square root of the variance.

$$\text{std}(X) = \sqrt{\text{Var}(X)} = \sigma(X)$$

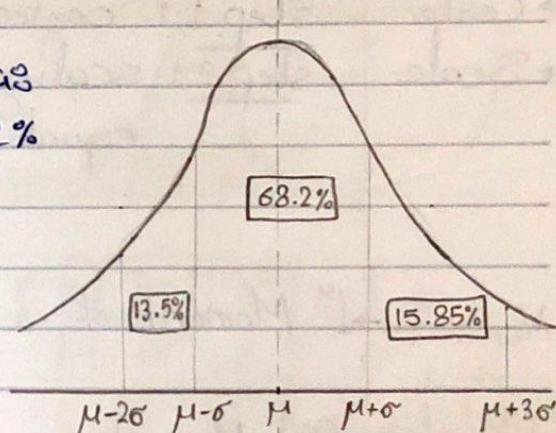
Gaussian Distribution for 2 Variables looks like a bell.

For normally distributed data

→ Between $[\mu - \sigma, \mu + \sigma]$ lies 68.2% of the data.

→ Between $[\mu - 2\sigma, \mu + 2\sigma]$ lies 95 % of the data.

→ Between $[\mu - 3\sigma, \mu + 3\sigma]$ lies 99.7 % of the data.



Normal or Gaussian Distribution

Total Response Time of a Computer System

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

$$R = T + L$$

Processing Time

Network (connection) Latency

we suppose that

$$T \sim N(10, 2^2)$$

independent

$$L \sim N(5, 1^2)$$

$$\mu_R = E[R] = E[T] + E[L] = \mu_T + \mu_L = 10 + 5 = 15$$

Therefore, the summation of two gaussians is still a gaussian.

$$\sigma_R^2 = \text{Var}(R) = \text{Var}(T+L) = \text{Var}(T) + \text{Var}(L)$$

$$\sigma_R^2 = \sigma_T^2 + \sigma_L^2 = 4 + 1$$

$$R \sim N(15, 5)$$

\rightarrow In general: $W = aX + bY$

independent $\begin{cases} X \sim N(\mu_X, \sigma_X^2) \\ Y \sim N(\mu_Y, \sigma_Y^2) \end{cases}$

then $W \sim N(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2)$

④ Standardize Distribution:

- Center step 1: centering the dataset about zero. ($X \rightarrow X - \mu$)
- Scale step 2: scaling the dataset to standard deviation equals to 1. ($X - \mu \rightarrow \frac{X - \mu}{\sigma}$)

Moments → 1st Moment of a random variable distribution is:

of a
Distribution → 2nd Moment:

→ Kth Moment:

$$E[X] = p_1 x_1 + p_2 x_2 + \dots + p_n x_n$$

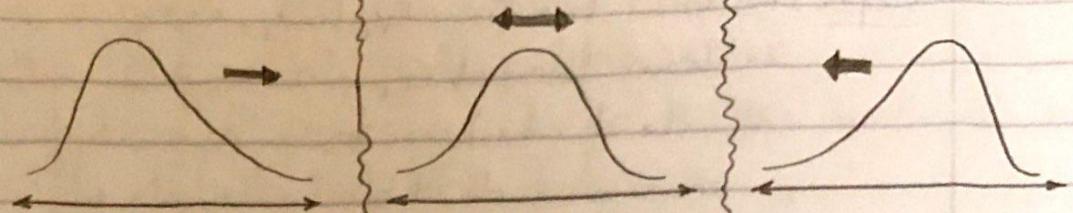
$$E[X^2] = p_1 x_1^2 + p_2 x_2^2 + \dots + p_n x_n^2$$

$$E[X^k] = p_1 x_1^k + p_2 x_2^k + \dots + p_n x_n^k$$

Skewness: is detected using the third moment of standardized event or distribution

$$E \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] = \text{Skewness}$$

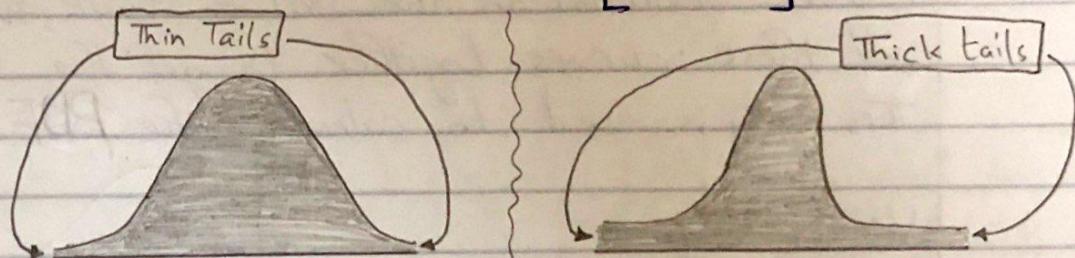
- Positively Skewed
- Right Skewed
- Not Skewed
- Zero Skewed
- Negatively Skewed
- Left Skewed



$$E\left[\left(\frac{X-\mu}{\sigma}\right)^3\right] > 0 \quad E\left[\left(\frac{X-\mu}{\sigma}\right)^3\right] = 0 \quad E\left[\left(\frac{X-\mu}{\sigma}\right)^3\right] < 0$$

Sometimes, there are different curves, which have similar variance, standard deviation, expected value and even skewness. Therefore, we use Kurtosis; it's the fourth moment of expected value of standardized distribution:

$$\text{Kurtosis} = E\left[\left(\frac{X-\mu}{\sigma}\right)^4\right]$$



$$E\left[\left(\frac{X-\mu}{\sigma}\right)^4\right] = \text{small}$$

$$E\left[\left(\frac{X-\mu}{\sigma}\right)^4\right] = \text{Large}$$

The k^{th} quantile ($q_{(k/100)}$) is the value that leaves $k\%$ of the values to the left and $(100-k)\%$ to the right.

• 25% Quantile \rightarrow 1st quartile

Q_1

• 50% Quantile \rightarrow 2nd quartile (median)

Q_2

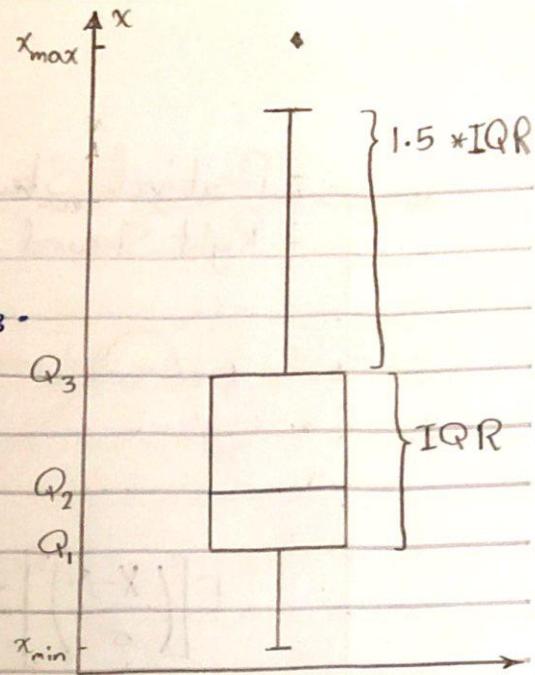
• 75% Quantile \rightarrow 3rd quartile

Q_3

$$P(X \leq q_{(k/100)}) = \frac{k}{100}$$

Box Plots:

- steps:
 - 1- draw box from Q_1 to Q_3 .
 - 2- draw line at Q_2 .
 - 3- Extend whiskers to maximum of 1.5 times the interquartile range ($IQR \times 1.5$).
 - 4- Represent the outliers with rhombus dots.



- insights:
 - 1- Skewness
 - 2- Presence of outliers
 - 3- Spread (Dispersion) Analysis

Kernel Density Estimation:

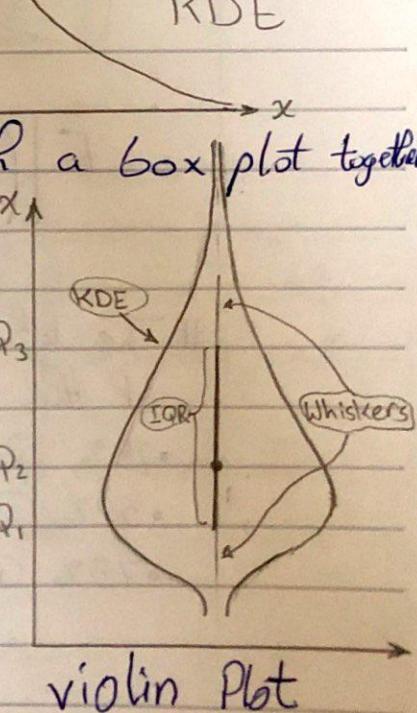
It is done by representing values with gaussian curves centered on the value itself, then add all these curves together to get the KDE graph. This way is used to calculate the PDF of datasets.

Violin Plots:

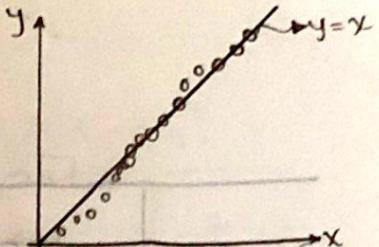
It is simply putting KDE with a box plot together

- # Many model assumes that data follows gaussian (normal) distribution:
 - Linear Regression
 - Logistic Regression
 - Gaussian Naive Bayes
 - Others

, and so does some data science tests.



Quantile-Quantile (QQ) Plots



It is a graph of Theoretical - Sample quantiles as scatter plot with a line ($y=x$). The more the distribution of points is closer to the line the more normally distributed the data is.

Joint Distribution (Discrete):

It is easy to plot a feature, but when we have two features to plot on the same graph it gets harder.

Example:

Get the joint distribution graph of Features X and Y for event E; where :

E → 2 6-sided dices were rolled

X → Number on dice 1

Y → Sum of two numbers appeared

Solution:

	Y					
	12	11	10	9	8	7
12				6,6		
11				5,6	6,5	
10			4,6	5,5	6,4	
9		3,6	4,5	5,4	6,3	
8	2,6	3,5	4,4	5,3	6,2	
7	1,6	2,5	3,4	4,3	5,2	6,1
6	1,5	2,4	3,3	4,2	5,1	
5	1,4	2,3	3,2	4,1		
4	1,3	2,2	3,1			
3	1,2	2,1				
2	1,1					
1						

• each box has probability of $\frac{1}{36}$

• If we replace each box with the value of its probability, we will get the joint distribution of X & Y.

For a continuous dataset:

- we use a scatter plot
- Expected Value ($E[X], E[Y]$) → coordinate
- We can calculate the variances of each Feature on its own: $\text{Var}(X) = E[X^2] - E[X]^2$

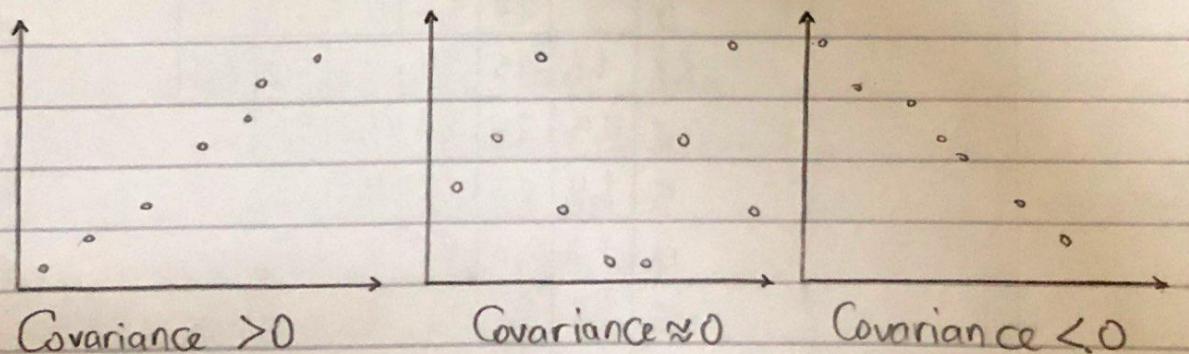
Discrete Joint Distribution	Continuous Joint Distribution
<ul style="list-style-type: none">- Assign Probability to individual outcomes.- Countable possible outcomes.	<ul style="list-style-type: none">- Assign Probability to ranges or intervals.- Infinite number of possible outcomes.

Marginal distribution is getting the probability of values of one feature (variable) from multi-variable graph or joint distribution.

Marginal distribution is useful in calculating conditional probability.

* Covariance:

A measure of relation between two variables:



in standardized
Graph →

$$(\sum xy) > 0$$

$$(\sum xy) \approx 0$$

$$(\sum xy) < 0$$

For equal Probabilities

$$\text{Covariance}(X, Y) = \frac{\sum_{i=1}^n (x_i - \mu_x) \cdot (y_i - \mu_y)}{n}$$

$Y = y$

Missused!!

General

$$\text{Cov}(X, Y) = \sum_{i=1}^n P_{XY}(x_i, y_i) \cdot (x_i - \mu_x) \cdot (y_i - \mu_y)$$

$$= E[XY] - E[X]E[Y]$$

→ In case, we have more than 2 variables: A, B, C, D and E for example. We use something called "Covariance Matrix"; where:

	A	B	C	D	E
A	$\text{Var}(A)$	$\text{Cov}(B, A)$	$\text{Cov}(C, A)$	$\text{Cov}(D, A)$	$\text{Cov}(E, A)$
B	$\text{Cov}(A, B)$	$\text{Var}(B)$	$\text{Cov}(C, B)$	$\text{Cov}(D, B)$	$\text{Cov}(E, B)$
C	$\text{Cov}(A, C)$	$\text{Cov}(B, C)$	$\text{Var}(C)$	$\text{Cov}(D, C)$	$\text{Cov}(E, C)$
D	$\text{Cov}(A, D)$	$\text{Cov}(B, D)$	$\text{Cov}(C, D)$	$\text{Var}(D)$	$\text{Cov}(E, D)$
E	$\text{Cov}(A, E)$	$\text{Cov}(B, E)$	$\text{Cov}(C, E)$	$\text{Cov}(D, E)$	$\text{Var}(E)$

Diagonal of variances and rest of covariances, hence, $\text{Cov}(X, Y) \equiv \text{Cov}(Y, X)$.

We use the sign of Covariance only.

* Correlation Coefficients

Standardized Covariance

$$\text{Correlation Coefficient} = \frac{\text{Cov}(X, Y)}{\sigma_X \cdot \sigma_Y}$$

range: $(-1, 1)$

we can use the sign and the value of the correlation coefficient.

Multivariable Gaussian Distribution

For two variables X & Y that are independent and normally distributed, the graph would look like a bell shape. The equation would be:

$$f_{XY}(x, y) = f_X(x) f_Y(y)$$

~~$\Sigma \rightarrow \text{covariance matrix}$~~

$$f_{XY}(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} e^{-0.5 \left(\frac{(x-\mu_x)^2}{\sigma_x^2} + \frac{(y-\mu_y)^2}{\sigma_y^2} \right)}$$

~~$\text{Cov}(X, Y) = 0$~~
~~If X & Y are independent~~

$$f_{XY}(x, y) = \frac{1}{2\pi \det \Sigma^{1/2}} \exp \left(-\frac{1}{2} ([x, y] - \mu)^T \Sigma^{-1} ([x, y] - \mu) \right)$$

~~$\sigma' = \sqrt{\text{Variance}}$~~

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{(x-\mu)^2}{\sigma^2} \right)}$$