

PART-I → Data & Programming Ideas.

Q1) What is the difference b/w long and wide data? What pandas function are used to transform one data type to the other?

Sol:- Long data is "long" - typically there will be 3 columns, id, variable, and value, which would look something like

ID	Variable	Value
Barry Bonds	HR	762
Pete Rose	H	4256
Barry Bonds	RBI	1996

Wide format data will have a row for each ID, with variables as columns.

Pandas use pivot to go from long to wide, and melt to go from wide to long.

Q2) What is the difference between a pandas dataframe and a series?

Sol: Series is a type of list in pandas, which can take integer values, string values, double values and more.

→ Series can only contain single list with index whereas dataframe can be made of more than one series or we can say that a dataframe is a collection of series that can be used to analyse the data.

Q3) What is the difference

Q3) What are the types of joins and how do they differ from one another?

Sol:- There are mainly 4 types of joins -

- i) Inner Join
- ii) Left Join
- iii) Outer Join
- iv) Right Join

Inner Join - Inner Join gives Returns a dataframe with only those rows that have common characteristics gives us rows with the intersection for index values that are contained in both datasets.

Two dataframes A and B,

A:	a	b	B:	c	d
index			index		
0			0		
1			2		
2			3		
3			5		
4					

Inner Join \Rightarrow	index	a	b	c	d
	0	-	-	-	-
	2	-	-	-	-
	3	-	-	-	-

Left Join - Returns a dataframe containing all the rows of the left data frame. Contains values that are index on the left

A.join(B, how='left')

index	a	b	c d
0	-	-	- -
1	-	-	NA NA
2	-	-	- -
3	-	-	- -
4	-	-	NA NA

Right Join - Similar to left join. The only difference is that all the rows of the right dataframe are taken as it is and only those of the left dataframe that are common to both

index	a	b	c	d
0	-	-	-	-
2	-	-	-	-
3	-	-	-	-
5	NA	NA	-	-

Outer Join - Returns all those records which either have a match in the left or right dataframe.

Index	a	b	c	d
0	-	-	-	-
1	-	-	NA	NA
2	-	-	-	-
3	-	-	-	-
4	-	-	-	-
5	-	-	NA	NA
NA	NA	-	-	-

Q4.) What is groupby and what does it do?

Sol: groupby's find all the pairs of matching values in a set of columns, and then perform aggregations on them.

In other words, it is a function used to split the data into groups based on some criteria

e.g:-

	Animal	Man Speed
0	Falcon	380
1	Falcon	370
2	Peregrine	24
3	Peregrine	26

df.groupby(['Animal']).mean()

Animal	Man Speed
Falcon	375
Peregrine	25

Q5.) When I do arithmetic operations on two series with partial overlapping indices, what is the index of new series? What values are there for indices in both the series, and ones that are only in one?

Sol: The index of the new series is a union of indices of both the series. The value for the indices in both the series is the resulting operation whereas for the ones that are only in one is NaN.

area = pd.Series ({'Alaska': 1723337, 'Texas': 695662, 'California': 423967}, name = 'area')

Population = pd.Series ({'California': 38332521, 'Texas': 26448193, 'New York': 19651127, name = 'pop'})

area | population

Alaska	NaN
California	0.011060
New York	NaN
Texas	0.026303

PART-II → PROBABILITY.

Q1.) Define a sample space. What are the 3 axioms of probability?

Sol: The sample space of an experiment or a random trial is a set of all possible outcomes or results of that experiment.

Probability is a function that maps events to \mathbb{R} and follows 3 axioms -

$$i) P(E) \in [0,1]$$

$$ii) P(S) = 1$$

iii) If we have sequence of disjoint sets E_1, E_2, \dots, E_n
 $E_i \cap E_j = \emptyset$

$$P(E_1 \cup E_2 \cup \dots \cup E_n) = P(E_1) + P(E_2) + \dots + P(E_n)$$

Q2.) What does it mean when 2 events are independent?

Sol: Two events are independent if the occurrence of one does not affect the probability of the other.

If E_1 and E_2 are 2 independent events then,

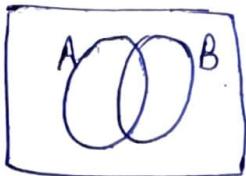
$$P(E_1 \cap E_2) = P(E_1) P(E_2)$$

x, y are independent iff

$$P(X \in A, Y \in B) = P(X \in A) P(Y \in B)$$

$$p(x, y) = p_x(x) p_y(y), f(x, y) = f_x(x) f_y(y)$$

Q3.) What is the probability of the union of 2 events (not necessarily disjoint)?
 Let A and B be 2 events that are not disjoint.



$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Q4.) What is Conditional Probability? What is Baye's Theorem?
 The conditional probability of an event B is the probability that the event will occur given the knowledge that an event A has already occurred.

$$P(B|A) = \frac{P(A \text{ and } B)}{P(A)} = \frac{P(A \cap B)}{P(A)}$$

Bayes Theorem is a way of finding a probability when we know certain other probabilities.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad P(B|A) = \frac{P(A \cap B)}{P(A)}$$

$$\Rightarrow P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Q5.) Suppose I roll 2 six-sided dice. X_1 is the result of the first die roll, X_2 is the result of the second die roll and $X = X_1 + X_2$. Calculate the foll.:

$$P(X=5)$$

$$P(X_1=2 | X=6)$$

$$P(X=6 \cup X_2=3)$$

$$\text{Sol: } P(X_1 + X_2 = 5) = P\left\{(1,4) \cup (2,3) \cup (3,2) \cup (4,1)\right\} \\ = 1 \cdot \frac{1}{36} + \frac{1}{36} + \frac{1}{36} + \frac{1}{36} = \frac{1}{9}$$

$$P(X=5) = \frac{1}{9}$$

$$P(x_1=2 | x=6) = \frac{P(x_1 \cap x)}{P(x)} = \frac{P(x_1=2 \cap x=6)}{P(x=6)}$$

$$P(x_1=2 \cap x=6) = P(\{(2,4)\}) = \frac{1}{36}$$

$$\begin{aligned} P(x=6) &= P(\{(1,5), (2,4), (3,3), (4,2), (5,1)\}) \\ &= \frac{5}{36} \end{aligned}$$

$$P(x_1=2 | x=6) = \frac{\frac{1}{36}}{\frac{5}{36}} = \frac{1}{5}$$

$$\begin{aligned} P(x=6 \cup x_2=3) &= P(x=6) + P(x_2=3) - P(x=6 \cap x_2=3) \\ &= \frac{5}{36} + \frac{6}{36} - \frac{1}{36} = \frac{10}{36} = \frac{5}{18}. \end{aligned}$$

Q6) What is a random variable? What is the difference b/w discrete and continuous random variable?

Sol:- A random variable, usually written X is a variable whose possible values are numerical outcomes of a random phenomenon.

There are 2 types of random variables -

Discrete Random Variable - which takes finite values or countably infinite values.

Continuous Random Variables - takes every value in at least one interval,

Q7.) Suppose ~~the~~ a random variable X has the following density -

$$f(x) = \begin{cases} cx^2 & -1 \leq x \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

Calculate c so that this is a proper density and calculate the expectation and variance of X .

$$\text{Sol: } \int_{-1}^1 n^2 dn = \left[\frac{n^3}{3} \right]_{-1}^1 = \frac{2}{3}c = 1$$

$$\Rightarrow c = \frac{3}{2}$$

$$\therefore f(n) = \begin{cases} \frac{3}{2}n^2, & -1 \leq n \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

$$E[X] = \int_{-\infty}^{\infty} n f(n) dn = \int_{-1}^1 n \cdot \frac{3}{2} n^2 dn = \frac{3}{2} \left[\frac{n^4}{4} \right]_{-1}^1$$

$$E[X] = \frac{3}{2} \left(\frac{1}{4} - \frac{1}{4} \right) = 0$$

$$\text{Var}[X] = E[X^2] - (E[X])^2 = E[X^2] = \int_{-\infty}^{\infty} n^2 f(n) dn$$

$$\Rightarrow \text{Var}[X] = \int_{-1}^1 n^2 \cdot \frac{3}{2} n^2 dn = \frac{3}{2} \left[\frac{n^5}{5} \right]_{-1}^1 = \cancel{\frac{3}{2}} \frac{3}{5}$$

Q8) Suppose a discrete random variable X has the following p.m.f.

$$P(n) = \begin{cases} 0.25 & n=0 \\ 0.5 & n=1 \\ 0.25 & n=2 \end{cases}$$

Calculate the expectation and variance of X.

$$\text{Sol: } E[X] = \sum_{n=0}^2 n \cdot f(n) = 0 \cdot (0.25) + 1 \cdot (0.5) + 2 \cdot (0.25) = 1$$

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

$$E[X^2] = 0^2(0.25) + 1^2(0.5) + 2^2(0.25) = 0 + 0.5 + 1 = 1.5$$

$$\Rightarrow \text{Var}[X] = 1.5 - 1 = 0.5$$

89) If X and Y have the joint density

$$f(x, y) = \begin{cases} x+y & x, y \in [0, 1] \\ 0 & \text{otherwise} \end{cases}$$

Calculate the covariance and correlation of X and Y .

Sol: $f_x(x) = \int_0^1 (x+y) dy = x + \frac{1}{2}$

$$f_y(y) = \int_0^1 (x+y) dx = y + \frac{1}{2}$$

$$E[X] = \int_0^1 x \left(x + \frac{1}{2} \right) dx = \left[\frac{x^3}{3} + \frac{x^2}{2} \right]_0^1 = \frac{7}{12}$$

$$E[Y] = \int_0^1 y \left(y + \frac{1}{2} \right) dy = \left[\frac{y^3}{3} + \frac{y^2}{4} \right]_0^1 = \frac{7}{12}$$

$$E[X^2] = \int_0^1 x^2 \left(x + \frac{1}{2} \right) dx = \left[\frac{x^4}{4} + \frac{x^3}{6} \right]_0^1 = \frac{5}{12}$$

$$\text{Var}[X] = \frac{5}{12} - \left(\frac{7}{12} \right)^2 = \frac{11}{144}$$

$$\Rightarrow \text{SD}[X] = \sqrt{\frac{11}{144}} \quad \text{Var}[Y] = \frac{11}{144} \quad \text{SD}[Y] = \sqrt{\frac{11}{144}}$$

$$\begin{aligned} E[XY] &= \int_0^1 \int_0^1 xy * (x+y) dx dy = \int_0^1 \int_0^1 (x^2 y + xy^2) dx dy \\ &= \int_0^1 \int_0^1 x^2 y dx dy + \int_0^1 \int_0^1 xy^2 dx dy \\ &= \int_0^1 \left[\frac{x^3}{3} \right]_0^1 y dy + \int_0^1 \left[\frac{x^2}{2} \right]_0^1 y^2 dy \\ &= \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{3} = \frac{2}{6} = \frac{1}{3} \end{aligned}$$

$$\begin{aligned}\therefore \text{Cov}(x, y) &= E[xy] - E[x]E[y] \\ &= \frac{1}{3} - \frac{7}{12} \cdot \frac{7}{12} \\ &= -\frac{1}{144}\end{aligned}$$

$$\text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\text{SD}(x)\text{SD}(y)} = \frac{-\frac{1}{144}}{\frac{11}{144}} = -\frac{1}{11}$$

$$\text{Corr}(x, y) = -\frac{1}{11}$$

PART-III : STATISTICS

Q1) Suppose we have a sample x_1, \dots, x_n of random variable with the following density:

$$f(n) = \begin{cases} an * \exp\left(-\frac{ax^2}{2}\right), & n > 0 \\ 0, & \text{otherwise} \end{cases}$$

Sol: $L(a | x_1, x_2, \dots, x_n) = \prod_{i=1}^n ax_i * \exp\left(-\frac{ax_i^2}{2}\right)$

$$l(a | x_1, x_2, \dots, x_n) = an \log(a) + \sum_{i=1}^n -\frac{a}{2} x_i^2 + \sum_{i=1}^n \log(x_i)$$

$$\Rightarrow \frac{dl}{da} = 0 \Rightarrow \frac{n}{a} - \frac{1}{2} \sum_{i=1}^n x_i^2 + 0 = 0$$

$$\Rightarrow \boxed{\hat{a} = \frac{2n}{\sum_{i=1}^n x_i^2}}$$

Q2.) What is a confidence interval? what is the difference between a one sided lower/upper confidence interval and a two-sided confidence interval.

Sol:- Confidence Interval is a statistic you calculate from a sample that tries to quantify a reasonable range of values for parameter of interest

$(1-\alpha)100\%$ confidence interval means that for $(1-\alpha)100\%$ of the samples, the true value will be contained in the confidence interval.

Two sided confidence bounds -

When we use a 2 sided confidence interval, we are looking at a closed interval where a certain percentage of the population is likely to lie i.e. we determine the values, b/w which a specified percentage of the population lies.

One sided confidence bounds

One sided confidence bounds are essentially an open ended version of 2 sided bounds. A one-sided bound defines the point where a certain percentage of the population is either higher or lower than the defined point.

An upper one-sided confidence bound defines that a certain percentage of population is less than.

A lower one-sided bound defines a point that a certain percentage of population is greater than.

$$\text{Lower : } \left[\bar{x} - z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}, \infty \right)$$

$$\text{Upper : } \left(-\infty, \bar{x} + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \right]$$

Q3.) What is a hypothesis test? Understand the following terms

- Critical region
- level (of a test)
- p-value
- power.

Sol:- Hypothesis test is a test where you try to see whether you have sufficient evidence to reject a hypothesis about ~~about~~ a parameter you are interested in, in your data.

Critical region is the region of values that corresponds to the rejection of the null hypothesis at some chosen probability level α .

Level of Test - α - probability of rejecting the null hypothesis when it is true.

p-value - probability of observing sample data ^{at least} as extreme as the actually obtained test statistic.

power - Power is the probability of not making a Type-II error.

Q4.) Suppose I have a normally distributed data x_1, x_2, \dots, x_{10} with $\bar{x} = 1$ with unknown variance σ^2 , but an estimate of $S^2 = 0.8$ from the data. How do I get a $(1-\alpha)100\%$ confidence interval (do not use bootstrap) with $\alpha = 0.05$? What if I wanted to test

$$\begin{aligned} H_0: \mu &= 0 \\ H_A: \mu &\neq 0 \end{aligned} \quad \text{at } \alpha = 0.05$$

Sol: $P \left(\bar{x} \pm t_{\alpha/2, n-1} \cdot \frac{S}{\sqrt{n}} \right) = 1 \pm t_{0.025, 9} \cdot \frac{\sqrt{0.8}}{\sqrt{10}}$
 $= 1 \pm 2.262 \cdot \frac{\sqrt{0.8}}{\sqrt{10}} \quad (0.36, 1.64)$

We can reject H_0 , as its not in the confidence interval.

Q5.) Describe the bootstrap procedure (write out an algorithm in pseudo code)

Sol:-

Step 1: Resample data with replacement x_1, x_2, \dots, x_n .

Step 2: Create m bootstrap samples for x_1, x_2, \dots, x_n

Step 3: Calculate MLE $\hat{\theta}$ for each bootstrap sample

Step 4: Calculate the lower and upper bound by sampling percentile to calculate the confidence interval

$$\text{lower-bound} = \text{np.percentile}(\hat{\theta}, \alpha/2)$$

$$\text{upper-bound} = \text{np.percentile}(\hat{\theta}, 100 - \alpha/2)$$

PART-IV : Machine Learning and scikit learn.

Q1) How do ridge regression and LASSO differ from the traditional linear regression? What is the difference b/w LASSO & Ridge?

Sol:- Linear Regression - The linear regression gives an estimate which minimises the sum of square error.

$$\underset{i=1}{\text{Min}} \sum^n (y_i - \beta x_i - \alpha)^2$$

$$\underset{\beta}{\text{Min}} \sum_{i=1}^N [y_i - \beta_1 x_{i,1} - \beta_2 x_{i,2} - \dots - \beta_d x_{i,d}]^2$$

$$= \underset{\beta}{\text{Min}} \sum_{i=1}^N (Y - X\beta)^T (Y - X\beta)$$

Ridge Regression - The ridge regression gives an estimate which minimises the sum of squared error, as well as satisfy the constraint $\sum_{j=1}^n \beta_j^2 \leq S$.

$$\underset{\beta}{\text{Min}} \sum_{i=1}^N (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta$$

Lasso Regression - Lasso regression gives an estimate which minimises the sum of the squared error as well as satisfy the constraint $\sum_{j=1}^n |\beta_j| \leq S$

$$\underset{\beta}{\text{Min}} (Y - X\beta)^T (Y - X\beta) + \lambda \|\beta\|_1$$

Q2) what kinds of problems do we use logistic regression for?
 What is the relationship b/w the covariates, regression coefficients and the predicted probabilities in logistic regression

Sol:- We use logistic regression for binary classification problems.

$$X = \begin{bmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_n \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad y_i = 0 \text{ or } 1$$

$$\ln\left(\frac{p_i}{1-p_i}\right) = \vec{x}_i^T \beta = \sum_{j=1}^d x_{ij} \beta_j$$

$$p_i = P(y_i = 1 \mid \vec{x}_i)$$

↑
observed
for
person i

predictors
for unit i

$$\frac{p_i}{1-p_i} = \exp(\vec{x}_i^T \beta)$$

$$p_i = (1-p_i) \exp(\vec{x}_i^T \beta)$$

$$\Rightarrow p_i = \frac{\exp(\vec{x}_i^T \beta)}{1 + \exp(\vec{x}_i^T \beta)}$$

Q3) Describe the procedure of Cross Validation?

Sol:- Cross validation is a resampling procedure used to evaluate machine learning models on a limited data sample.

Randomly partition my training set into k parts of approx same size i.e. k -folds.

Fit the model for each i in the grid, k -times.

Algorithm -

1.) shuffle the dataset randomly

2.) split dataset into k groups.

3.) for each unique group -

a) Take the group as a held out test data set.

b) fit a model on training set and evaluate it on the test set

c) fit a model on the training and evaluate on the test set.

d) Retain the evaluation score and discard the model.

4.) Summarize the skill of model using the sample of model evaluation scores.

Q4.) Describe at high level, how a binary tree works?
i.e. given a datapoint, how do you end up with a prediction?
What are the parts of a tree? Nodes, leaves etc

Q5.) What are boosting and bagging?

Boosting and Bagging come under ensemble methods.

Sol:- The general principle of ENSEMBLE METHODS is to combine the predictions of several models. These are built with a given learning algorithm in order to remove the robustness over a single model.

2 types -

- Parallel Ensemble Methods - In these methods, the base learners are generated in parallel simultaneously.
e.g:- when deciding the movie you want to watch, you may ask multiple friends for suggestions and probably watch the movie which got the highest votes.
- Sequential Ensemble Methods - In this technique, different learners learn sequentially with each learner fitting simple model to the data.

Then the data is analyzed for errors. The goal is to solve for net error from the prior model. The overall performance can be boosted by weighing previously mislabeled samples with higher weight.

Bagging - (Boot Strap Aggregating).

Bagging, a parallel ensemble methods, is a way to decrease the variance of the prediction model by generating additional data in the training stage. This is produced by random sampling with replacement from the original set.

These multisets of data are used to train multiple models. As a result we end with ensemble of diff. methods. The average of predictions from all the different models is used.

Boosting-

Boosting is a sequential ensemble method that in general decreases the bias error and builds strong, predictive models. The term 'Boosting' refers to family of algorithms which convert a weak learner to a strong learner.

In each iteration, data points that are mispredicted are identified and their weights are increased so that the next learner pays extra attention to get them right.

During training, the algorithm allocates weight to each resulting model. A learner with good prediction results on the training stage will be assigned a higher weight than a poor one.

Q6.) What is a random forest and how does it differ from a binary decision tree?