# Distribution Plots

Distribution plots display a single continuous feature and help visualize properties such as deviation and average values.

3 main distribution plot types:
- ○  Rug Plot
- ○  Histogram
- ○  KDE Plot

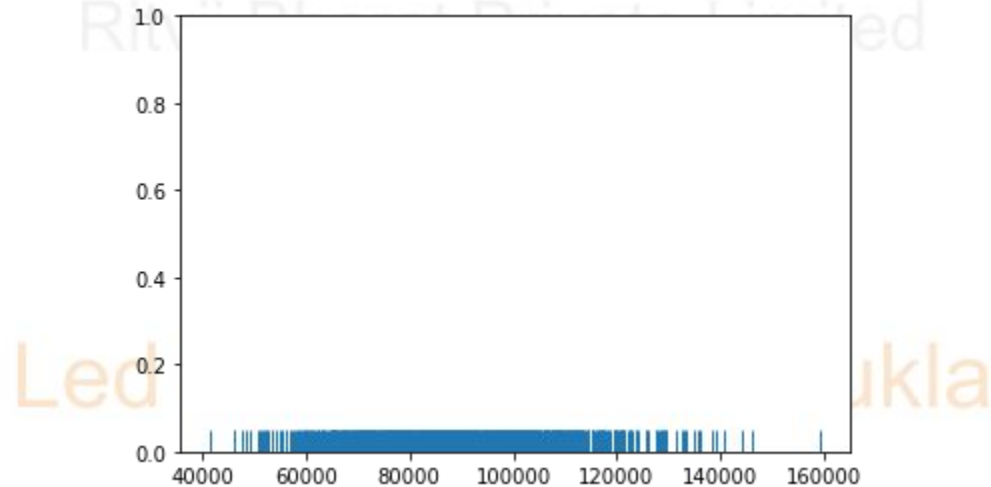Let's explore the distribution of employee salaries.

- One way is a rug plot which is the simplest distribution plot and merely adds a dash or tick line for every single value.
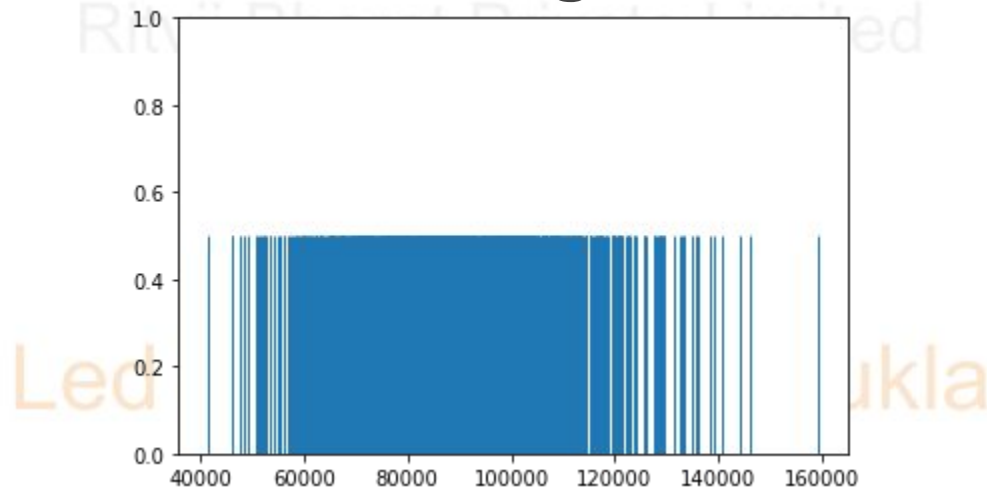- The y-axis does not really have a meaning.

1. Adds a tick for every salary value
2. Optionally adjust height of ticks
3. Y-axis not interpretable



22

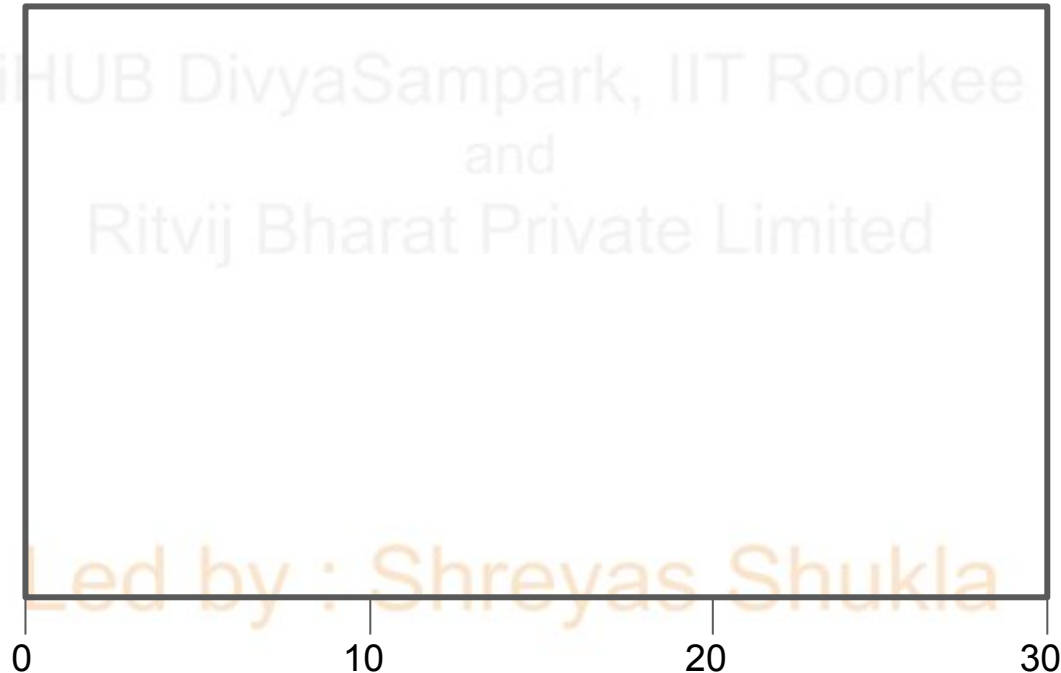1. Highest salary near $160,000
2. Many salaries between $60k - $120k

Many ticks could be right on top of eachother, we can't tell!If we count how many ticks there are per various x-ranges, we can create a histogram.

Let's explore a simple example

0          10          20          30

We place the rug plot ticks

# Choose a number of "bins", we'll pick 3

# Count ticks per bin

# Create a bar as high as count

# Create a bar as high as count

# Create a bar as high as count

# Histogram is complete

# Y-axis can also be normalized as percent

Changing number of bins shows more detail instead of general trends.

Changing number of bins shows more detail instead of general trends.



34

Changing number of bins shows more detail instead of general trends.

Seaborn also allows us to add on a KDE plot curve on top of a histogram.

Let's explore what a KDE plot is and how it is constructed.

Kernel Density Estimation (KDE) is a method of **estimating** a probability density function of a random variable.

In simpler terms, it is a way of estimating a continuous probability curve for a finite data sample.

KDE plots are best understood by visualizing their "construction".

Let's start with a rug plot….

# Mastering Machine Learning with Python
## (27th Aug 2024 - 18th Oct 2024)



Led by : Shreyas Shukla

You can change the kernel and bandwidth used which can make your KDE show more or less of the variance contained in the data.
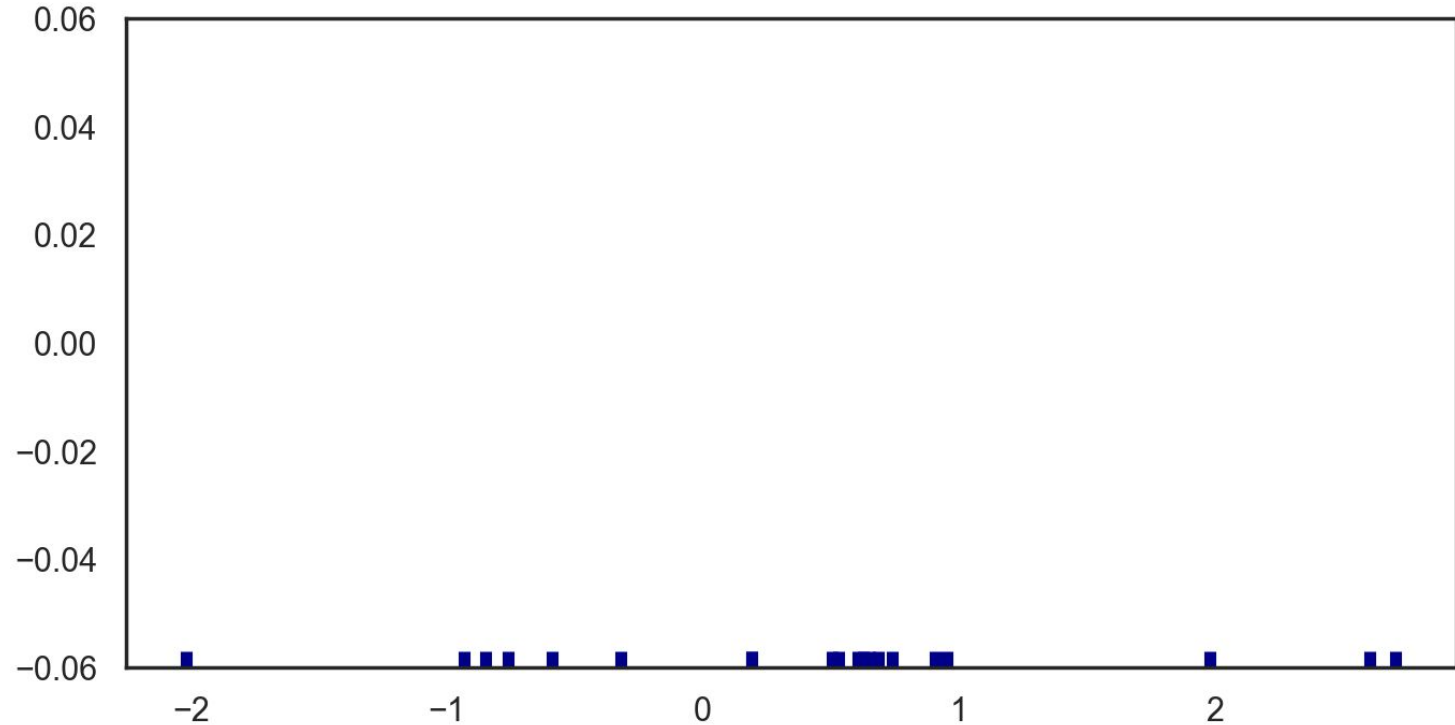
iHUB DivyaSampark, IIT Roorkee
and
Ritvij Bharat Private Limited

Led by : Shreyas Shukla

# Mastering Machine Learning with Python
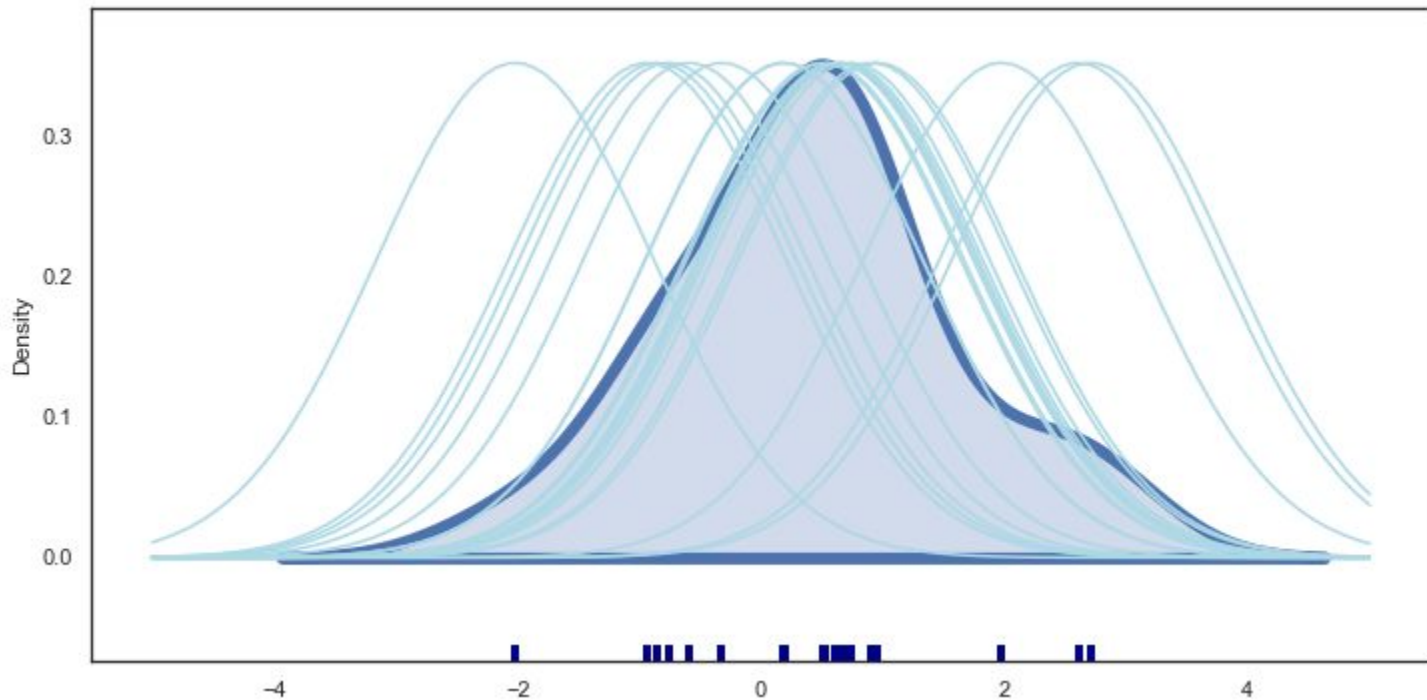## (27th Aug 2024 - 18th Oct 2024)

iHUB DivyaSampark, IIT Roorkee
and
Ritvij Bharat Private Limited

Let's Code !!

Led by : Shreyas Shukla

# Categorical Plots

Statistical Estimation within Categories
Part One: Understanding the Plots

The categorical plots display a statistical metrics **per** a category.

- For example mean value per category or a count of the number of rows **per** category.
- It is the visualization equivalent of a groupby() call.

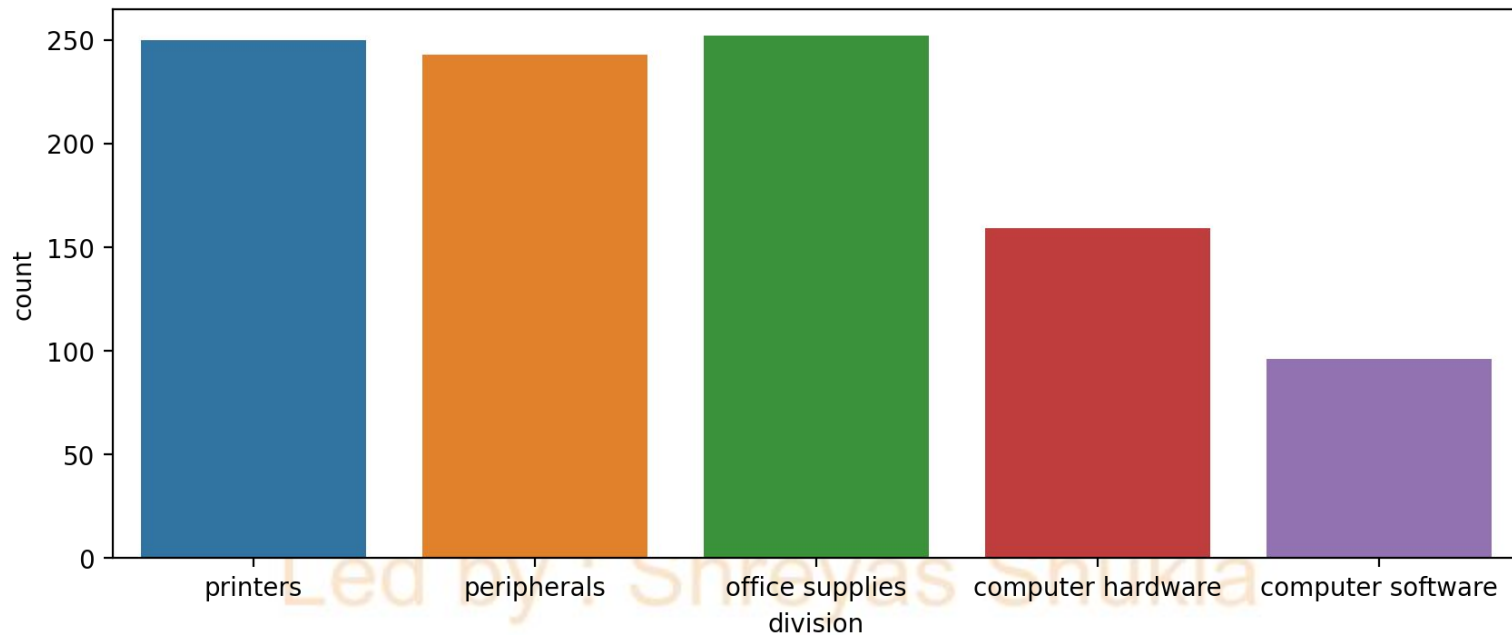Led by : Shreyas Shukla

Two main types of plots for this are:
- countplot()
  - Counts number of rows per category.
- barplot()
  - General form of displaying any chosen metric per category.

Led by : Shreyas Shukla

# Countplot for corporate divisions

# Countplot for education level

# Countplot with additional hue separation

The barplot is the general form that allows you to choose any measure or estimator for the y axis.

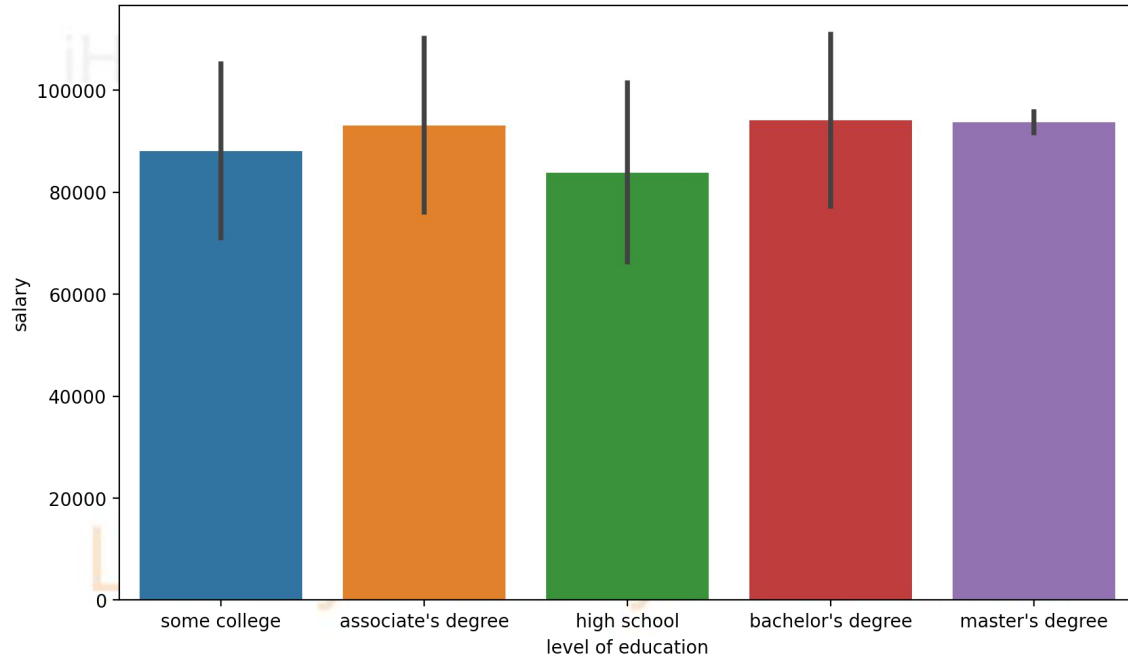We could plot the mean value and standard deviation per category instead.

# Caution!

- ○ Be very careful with these plots, since the bar is filled and continuous, a viewer may interpret continuity along the y-axis which may be incorrect!
- ○ Always make sure to add additional labeling and explanation for these plots!

Led by : Shreyas Shukla

# Barplot showing mean and SD bar

A simple table is probably better.

| level of education | mean | std |
|---|---|---|
| associate's degree | 93156.41 | 17066.06 |
| bachelor's degree | 94133.76 | 17007.09 |
| high school | 83887.35 | 17674.44 |
| master's degree | 93718.00 | 2497.63 |
| some college | 88115.84 | 17076.28 |

Let's explore coding out these plots with seaborn!

# Mastering Machine Learning with Python
## (27th Aug 2024 - 18th Oct 2024)

iHUB DivyaSampark, IIT Roorkee

and

Ritvij Bharat Private Limited

Let's Code !!

Led by : Shreyas Shukla

We've explored distribution plots for a single feature, but what if we want to compare distributions across categories?

For example, instead of the distribution of everyone's salary, we can compare the distributions of salaries **per** level of education.

We will first separate out each category, then create the distribution visualization.

Let's explore what plot types we have available….

Distribution within Categories
- ○ Boxplot
- ○ Violinplot
- ○ Swarmplot

Let's explore understanding these plots on the previous salary dataset.

The Boxplot displays the distribution of a continuous variable. It does this through the use of quartiles.

Quartiles separate out the data into 4 equal number of data points :
- 25% of data points are in bottom quartile.
- 50th percentile (Q2) is the median.

# Boxplot on single feature:



salary

# Median is 50th percentile. Median splits data in half

- IQR defines the box width
- 50% of all data points are inside the box

# Q1 is the 25th percentile below which are 25% of data points

# Q3 is the 75th percentile. 25% of all points are above Q3

The "whiskers" are defined by 1.5 × IQR

# Outside of the whiskers are outliers

Boxplot gives statistical distribution information in a visual format:

# Boxplot can be oriented vertically or horizontally.

We can create a box plot **per** category!

The violin plot plays a similar role as the box plot. It displays the probability density across the data using a KDE.
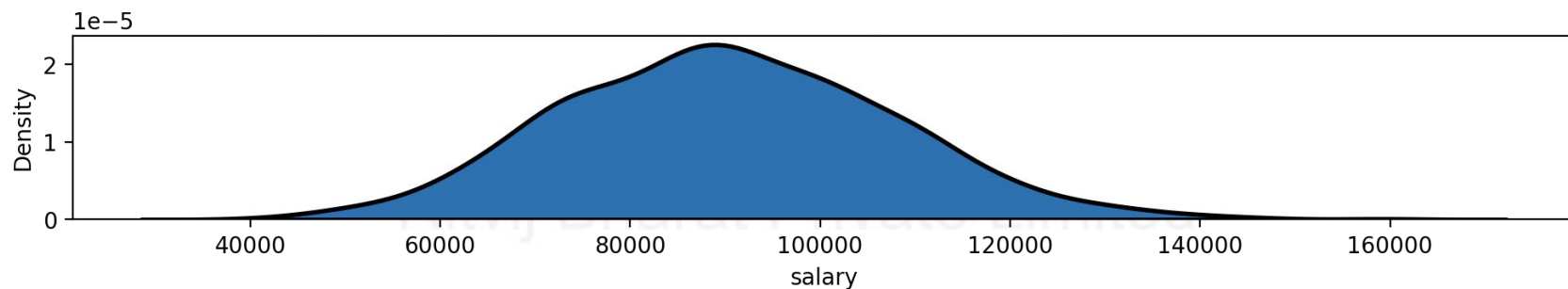
We can imagine it as a mirrored KDE plot.

The violin plot plays a similar role as the box plot. It displays the probability density across the data using a KDE.

We can imagine it as a mirrored KDE plot.

We take the KDE of a single feature:

# We could then "mirror" it:

# Then combine it to get the violin plot:

Then combine it to get the violin plot:



salary

The violin plots can then be created **per** category:

A few more less common categorical distribution plot, that is,  swarmplot .

Let's quickly explore this plot type…
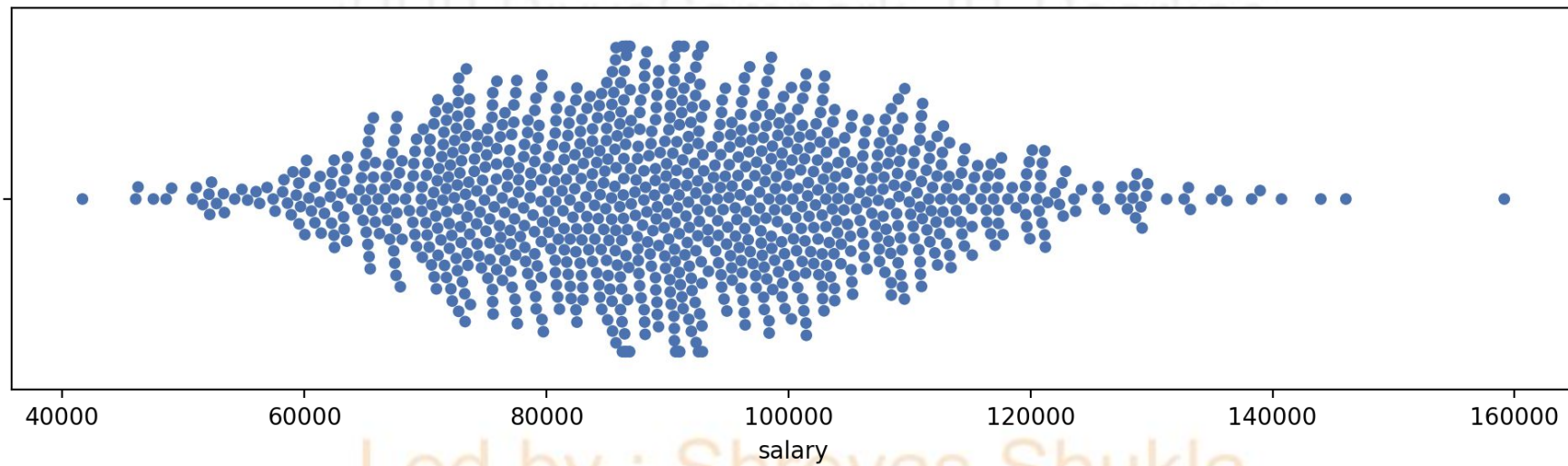
Led by : Shreyas Shukla

The swarmplot simply shows all the data points in the distribution.

(For very large data sets, it won't show all the points, but will display the general distribution of them.)

**That is it!!**

**Let's code !!** DivyaSampark, IIT Roorkee

and

Ritvij Bharat Private Limited

Led by : Shreyas Shukla