

- TOPICS:
- 0. ANOVA → Non-param
 - 1. Kruskal-Wallis [H] Test → When ANOVA breaks down
 - ✓ 2. Plot: don't just see numbers
 - ✓ 3. Biased estimate → why we use n-1 for std-dev
 - ✓ { 4. Case-study: Loan granting → Feature Engg..
 - using everything we know
 - ✓ [5. Handle missing values

OPS:

- { 1 vs 2 vs 3 classes break [wed is the last class]
- = F M F M W
- ↳ optional problem solving → pick least freq-answer
- ↳ whatsapp → Tag Ajay + Assignments Team

Handwritten text in pink ink:

Two parallel diagonal lines starting from the top left.

Handwritten text in pink ink:

AnoVA

The word "AnoVA" is underlined with a pink line. A small pink pencil icon is at the end of the underline.

A one-way between groups ANOVA is used to compare the means of more than two independent groups. A one-way between groups ANOVA comparing just two groups will give you the same results at the independent t test that you conducted in [Lesson 8](#). We will use the five step hypothesis testing procedure again in this lesson.

1. Check assumptions and write hypotheses

The assumptions for a one-way between groups ANOVA are:

- 1. Samples are independent
- 2. The response variable is approximately normally distributed for each group or all group sample sizes are at least 30
- 3. The population variances are equal across responses for the group levels (if the largest sample standard deviation divided by the smallest sample standard deviation is not greater than two, then assume that the population variances are equal)

Given that you are comparing k independent groups, the null and alternative hypotheses are:

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_k$$

$$H_a : \text{Not all } \mu_i \text{ are equal}$$

In other words, the null hypothesis is that all of the groups' population means are equal. The alternative is that they are not all equal; there are at least two population means that are not equal to one another.

2. Calculate the test statistic

ANOVA uses an F test statistic. Hand calculations for ANOVAs require many steps. In this class, you will be working primarily with Minitab outputs.

Conceptually, the F statistic is a ratio: $F = \frac{\text{Between groups variability}}{\text{Within groups variability}}$. Numerically this translates to $F = \frac{MS_{\text{Between}}}{MS_{\text{Within}}}$. In other words how much do individuals in different groups vary from one another over how much do individuals within groups vary from one another.

Statistical software will compute the F ratio for

n = MXK
{ n : m samples per group
K groups

- ▶ 2: Describing Data, Part 1
- ▶ 3: Describing Data, Part 2
- ▶ 4: Confidence Intervals
- ▶ 5: Hypothesis Testing, Part 1
- ▶ 6: Hypothesis Testing, Part 2
- ▶ 7: Normal Distributions
- ▶ 8: Inference for One Sample
- ▶ 9: Inference for Two Samples
- ▼ 10: One-Way ANOVA
 - 10.1 - Introduction to the F Distribution
 - 10.2 - Hypothesis Testing
 - 10.3 - Pairwise Comparisons
 - 10.4 - Minitab: One-Way ANOVA
 - 10.5 - Example: SAT-Math Scores by Award Preference
 - 10.6 - Example: Exam Grade by Professor
 - 10.7 - Lesson 10 Summary
- ▶ 11: Chi-Square Tests
- ▶ 12: Correlation & Simple Linear Regression

Resources

-
- | | | |
|---------------------|----------------|----------|
| Datasets | Glossary | Formulas |
| Contact | Help & Support | |
| Minitab Quick Guide | | |



A one-way between groups ANOVA is used to compare the means of more than two independent groups. A one-way between groups ANOVA comparing just two groups will give you the same results at the independent t test that you conducted in [Lesson 8](#). We will use the five step hypothesis testing procedure again in this lesson.

1. Check assumptions and write hypotheses

The assumptions for a one-way between groups ANOVA are:

1. Samples are independent ✓
2. The response variable is approximately normally distributed for each group or all group sample sizes are at least 30
3. The population variances are equal across responses for the group levels (if the largest sample standard deviation divided by the smallest sample standard deviation is not greater than two, then assume that the population variances are equal)

Given that you are comparing k independent groups, the null and alternative hypotheses are:

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_k$$

$$H_a : \text{Not all } \mu_i \text{ are equal}$$

In other words, the null hypothesis is that all of the groups' population means are equal. The alternative is that they are not all equal; there are at least two population means that are not equal to one another.

2. Calculate the test statistic

ANOVA uses an F test statistic. Hand calculations for ANOVAs require many steps. In this class, you will be working primarily with Minitab outputs.

Conceptually, the F statistic is a ratio: $F = \frac{\text{Between groups variability}}{\text{Within groups variability}}$. Numerically this translates to $F = \frac{MS_{\text{Between}}}{MS_{\text{Within}}}$. In other words how much do individuals in different groups vary from one another over how much do individuals within groups vary from one another.

Statistical software will compute the F ratio for

*box-Cox
tec-times
QQ/TS/LAD*

- ▶ 2: Describing Data, Part 1
- ▶ 3: Describing Data, Part 2
- ▶ 4: Confidence Intervals
- ▶ 5: Hypothesis Testing, Part 1
- ▶ 6: Hypothesis Testing, Part 2
- ▶ 7: Normal Distributions
- ▶ 8: Inference for One Sample
- ▶ 9: Inference for Two Samples
- ▼ 10: One-Way ANOVA
 - 10.1 - Introduction to the F Distribution
 - 10.2 - Hypothesis Testing
 - 10.3 - Pairwise Comparisons
 - 10.4 - Minitab: One-Way ANOVA
 - 10.5 - Example: SAT-Math Scores by Award Preference
 - 10.6 - Example: Exam Grade by Professor
 - 10.7 - Lesson 10 Summary
- ▶ 11: Chi-Square Tests
- ▶ 12: Correlation & Simple Linear Regression

Resources

- [Datasets](#)
- [Glossary](#)
- [Formulas](#)
- [Contact](#)
- [Help & Support](#)
- [Minitab Quick Guide](#)



A one-way between groups ANOVA is used to compare the means of more than two independent groups. A one-way between groups ANOVA comparing just two groups will give you the same results at the independent t test that you conducted in [Lesson 8](#). We will use the five step hypothesis testing procedure again in this lesson.

1. Check assumptions and write hypotheses

The assumptions for a one-way between groups ANOVA are:

1. Samples are independent
2. The response variable is approximately normally distributed for each group or all group sample sizes are at least 30
3. The population variances are equal across responses for the group levels (if the largest sample standard deviation divided by the smallest sample standard deviation is not greater than two, then assume that the population variances are equal)

$$S_1 = 10 \quad S_2 = 20 \quad S_3 = 18$$

Given that you are comparing k independent groups, the null and alternative hypotheses are:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$$

H_a : Not all μ_i are equal

$$S_1 = 10 \quad S_2 = 18$$

In other words, the null hypothesis is that all of the groups' population means are equal. The alternative is that they are not all equal; there are at least two population means that are not equal to one another.

2. Calculate the test statistic

ANOVA uses an F test statistic. Hand calculations for ANOVAs require many steps. In this class, you will be working primarily with Minitab outputs.

Conceptually, the F statistic is a ratio: $F = \frac{\text{Between groups variability}}{\text{Within groups variability}}$. Numerically this translates to $F = \frac{MS_{\text{Between}}}{MS_{\text{Within}}}$. In other words how much do individuals in different groups vary from one another over how much do individuals within groups vary from one another.

- ▶ 2: Describing Data, Part 1
- ▶ 3: Describing Data, Part 2
- ▶ 4: Confidence Intervals
- ▶ 5: Hypothesis Testing, Part 1
- ▶ 6: Hypothesis Testing, Part 2
- ▶ 7: Normal Distributions
- ▶ 8: Inference for One Sample
- ▶ 9: Inference for Two Samples
- ▶ 10: One-Way ANOVA
 - 10.1 - Introduction to the F Distribution
 - 10.2 - Hypothesis Testing
 - 10.3 - Pairwise Comparisons
 - 10.4 - Minitab: One-Way ANOVA
 - 10.5 - Example: SAT-Math Scores by Award Preference
 - 10.6 - Example: Exam Grade by Professor
 - 10.7 - Lesson 10 Summary
- ▶ 11: Chi-Square Tests
- ▶ 12: Correlation & Simple Linear Regression

Resources

- [Datasets](#)
- [Glossary](#)
- [Formulas](#)
- [Contact](#)
- [Help & Support](#)
- [Minitab Quick Guide](#)



$$\left\{ \begin{array}{l} n = m \text{ per GR} \\ \times k \text{ GPS} \end{array} \right.$$

$$T_{\text{Anova}} = f = \frac{\frac{1}{k-1} \sum_{i=1}^k (\bar{x}_i - \bar{x})^2 \times \underline{\underline{m}}}{\frac{1}{n-k} \sum_{i=1}^k \sum_{j=1}^m (x_{ij} - \bar{x}_{ij})^2}$$

Statistical software will compute the F ratio for you and produce what is known as an ANOVA source

table. The ANOVA source table will give you information about the variability between groups and within groups. The table below gives you all of the formulas, but you will not be responsible for performing these calculations by hand. Minitab will do all of these calculations for you and provide you with the full ANOVA source table.

M = # obs per group

Source	df	SS	MS	F	p
Between Groups (Factor)	$k - 1$	$\sum_k n_k (\bar{x}_k - \bar{x}_{\cdot})^2$	$\frac{SS_{Between}}{df_{Between}}$	$\frac{MS_{Between}}{MS_{Within}}$	Area to the right of $F_{k-1, n-k}$
Within Groups (Error)	$n - k$	$\sum_k \sum_i (x_{ik} - \bar{x}_k)^2$	$\frac{SS_{Within}}{df_{Within}}$		
Total	$n - 1$	$\sum_k \sum_i (x_{ik} - \bar{x}_{\cdot})^2$			

Legend

k Number of groups

n Total sample size (all groups combined)

n_k Sample size of group k

\bar{x}_k Sample mean of group k

\bar{x}_{\cdot} Grand mean (i.e., mean for all groups combined)

SS Sum of squares

MS Mean square



$$\left. \begin{array}{l} gp^1 \rightarrow M_1 \\ gp^2 \rightarrow M_2 \\ \vdots \\ gp^K \rightarrow M_K \end{array} \right\}$$

$$n = M_1 + M_2 + \dots + M_K$$

→ `scipy.stats`
→ `statsmodels` →

groups. A one-way between groups ANOVA comparing just two groups will give you the same results at the independent t test that you conducted in [Lesson 8](#). We will use the five step hypothesis testing procedure again in this lesson.

1. Check assumptions and write hypotheses

The assumptions for a one-way between groups ANOVA are:

1. Samples are independent
2. The response variable is approximately normally distributed for each group or all group sample sizes are at least 30
3. The population variances are equal across responses for the group levels (if the largest sample standard deviation divided by the smallest sample standard deviation is **not** greater than two, then assume that the population variances are equal)

Given that you are comparing k independent groups, the null and alternative hypotheses are:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$$

$$H_a : \text{Not all } \mu_i \text{ are equal}$$

In other words, the null hypothesis is that all of the groups' population means are equal. The alternative is that they are not all equal; there are at least two population means that are not equal to one another.

2. Calculate the test statistic

ANOVA uses an F test statistic. Hand

$s_1 \approx s_2 \approx s_3$

- ▶ 4: Confidence Intervals
- ▶ 5: Hypothesis Testing, Part 1
- ▶ 6: Hypothesis Testing, Part 2
- ▶ 7: Normal Distributions
- ▶ 8: Inference for One Sample
- ▶ 9: Inference for Two Samples
- ▶ 10: One-Way ANOVA
 - 10.1 - Introduction to the F Distribution
 - 10.2 - Hypothesis Testing
 - 10.3 - Pairwise Comparisons
 - 10.4 - Minitab: One-Way ANOVA
 - 10.5 - Example: SAT-Math Scores by Award Preference
 - 10.6 - Example: Exam Grade by Professor
 - 10.7 - Lesson 10 Summary
- ▶ 11: Chi-Square Tests
- ▶ 12: Correlation & Simple Linear Regression

Resources

- Datasets
- Glossary
- Formulas
- Contact



10.2 - Hypothesis Testing | ST x Kruskal–Wallis one-way analysis of variance x Anscombe's quartet - Wikipedia x Bias of an estimator - Wikipedia x EDA_FE.ipynb - Colaboratory x Markov's inequality - Wikipedia x Chebyshev's inequality - Wikipedia x + en.wikipedia.org/wiki/Kruskal–Wallis_one-way_analysis_of_variance

Not logged in Talk Contributions Create account Log in



WIKIPEDIA
The Free Encyclopedia

Main page
Contents
Current events
Random article
About Wikipedia

Contact us
Donate

Contribute

Help

Learn to edit
Community portal
Recent changes
Upload file

Tools

What links here
Related changes
Special pages
Permanent link

Article Talk

Read

Edit View history

Search Wikipedia



Kruskal–Wallis one-way analysis of variance

From Wikipedia, the free encyclopedia

The **Kruskal–Wallis test** by ranks, **Kruskal–Wallis *H* test**^[1] (named after William Kruskal and W. Allen Wallis), or **one-way ANOVA on ranks**^[1] is a **non-parametric** method for testing whether samples originate from the same distribution.^{[2][3][4]} It is used for comparing two or more independent samples of equal or different sample sizes. It extends the **Mann–Whitney *U* test**, which is used for comparing only two groups. The parametric equivalent of the Kruskal–Wallis test is the **one-way analysis of variance (ANOVA)**.

A significant Kruskal–Wallis test indicates that at least one sample **stochastically dominates** one other sample. The test does not identify where this stochastic dominance occurs or for how many pairs of groups stochastic dominance obtains. For analyzing the specific sample pairs for stochastic dominance, Dunn's test,^[5] pairwise Mann–Whitney tests with Bonferroni correction,^[6] or the more powerful but less well known Conover–Iman test^[6] are sometimes used.

Since it is a nonparametric method, the Kruskal–Wallis test does not assume a **normal distribution** of the residuals, unlike the analogous one-way analysis of variance. If the researcher can make the assumptions of an identically shaped and scaled distribution for all groups, except for any difference in medians, then the null hypothesis is that the medians of all groups are equal, and the alternative hypothesis is that at least one population median of one group is different from the population median of at least one other group. Otherwise, it is impossible to say, whether the rejection of the null hypothesis comes from the shift in locations or group dispersions. This is the same issue that happens also with the Mann-Whitney test.^{[7][8][9]}

Contents [hide]

1 Method



11/11

10.2 - Hypothesis Testing | ST x Kruskal–Wallis one-way analysis of variance x Anscombe's quartet - Wikipedia x Bias of an estimator - Wikipedia x EDA_FE.ipynb - Colaboratory x Markov's inequality - Wikipedia x Chebyshev's inequality - Wikipedia x + en.wikipedia.org/wiki/Kruskal–Wallis_one-way_analysis_of_variance

Not logged in Talk Contributions Create account Log in

Article Talk Read Edit View history Search Wikipedia



WIKIPEDIA
The Free Encyclopedia

Main page
Contents
Current events
Random article
About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Related changes

Special pages

Permanent link

Kruskal–Wallis one-way analysis of variance

From Wikipedia, the free encyclopedia

The **Kruskal–Wallis test** by ranks, **Kruskal–Wallis *H* test**^[1] (named after William Kruskal and W. Allen Wallis), or **one-way ANOVA on ranks**^[1] is a **non-parametric** method for testing whether samples originate from the same distribution.^{[2][3][4]} It is used for comparing two or more independent samples of equal or different sample sizes. It extends the **Mann–Whitney *U* test**, which is used for comparing only two groups. The parametric equivalent of the Kruskal–Wallis test is the **one-way analysis of variance** (ANOVA).

A significant Kruskal–Wallis test indicates that at least one sample **stochastically dominates** one other sample. The test does not identify where this stochastic dominance occurs or for how many pairs of groups stochastic dominance obtains. For analyzing the specific sample pairs for stochastic dominance, Dunn's test,^[5] pairwise Mann–Whitney tests with Bonferroni correction,^[6] or the more powerful but less well known Conover–Iman test^[6] are sometimes used.

Since it is a nonparametric method, the Kruskal–Wallis test does not assume a **normal distribution** of the residuals, unlike the analogous one-way analysis of variance. If the researcher can make the assumptions of an identically shaped and scaled distribution for all groups, except for any difference in medians, then the null hypothesis is that the medians of all groups are equal, and the alternative hypothesis is that at least one population median of one group is different from the population median of at least one other group. Otherwise, it is impossible to say, whether the rejection of the null hypothesis comes from the shift in locations or group dispersions. This is the same issue that happens also with the Mann-Whitney test.^{[7][8][9]}

Contents [hide]

1 Method



12/12

Languages



العربية

Deutsch

Español

Français

한국어

Italiano

Português

Русский

Türkçe

文 7 more

Edit links

Method

[edit]

1. Rank all data from all groups together; i.e., rank the data from 1 to N ignoring group membership. Assign any tied values the average of the ranks they would have received had they not been tied.
2. The test statistic is given by:

$$H = (N - 1) \frac{\sum_{i=1}^g n_i (\bar{r}_{i\cdot} - \bar{r})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (r_{ij} - \bar{r})^2}, \text{ where:}$$

- \bar{N} is the total number of observations across all groups
- g is the number of groups
- n_i is the number of observations in group i
- r_{ij} is the rank (among all observations) of observation j from group i
- $\bar{r}_{i\cdot} = \frac{\sum_{j=1}^{n_i} r_{ij}}{n_i}$ is the average rank of all observations in group i
- $\bar{r} = \frac{1}{2}(N + 1)$ is the average of all the r_{ij} .

3. If the data contain no ties the denominator of the expression for H is exactly $(N - 1)N(N + 1)/12$ and $\bar{r} = \frac{N+1}{2}$. Thus

$$H = \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \left(\bar{r}_{i\cdot} - \frac{N + 1}{2} \right)^2$$

$$= \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \bar{r}_{i\cdot}^2 - 3(N + 1)$$

*N observations
g #gps*

[Download as PDF](#)[Printable version](#)

Languages

[العربية](#)[Deutsch](#)[Español](#)[Français](#)[한국어](#)[Italiano](#)[Português](#)[Русский](#)[Türkçe](#)[文 A 7 more](#)[Edit links](#)

7 External links

Method [\[edit \]](#)

1. Rank all data from all groups together; i.e., rank the data from 1 to N ignoring group membership. Assign any tied values the average of the ranks they would have received had they not been tied.

2. The test statistic is given by:

$$H = (N - 1) \frac{\sum_{i=1}^g n_i (\bar{r}_{i\cdot} - \bar{r})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (r_{ij} - \bar{r})^2}, \text{ where:}$$

- N is the total number of observations across all groups
- g is the number of groups
- n_i is the number of observations in group i
- r_{ij} is the rank (among all observations) of observation j from group i
- $\bar{r}_{i\cdot} = \frac{\sum_{j=1}^{n_i} r_{ij}}{n_i}$ is the average rank of all observations in group i
- $\bar{r} = \frac{1}{2}(N + 1)$ is the average of all the r_{ij} .

3. If the data contain no ties the denominator of the expression for H is exactly $(N - 1)N(N + 1)/12$ and $\bar{r} = \frac{N+1}{2}$. Thus

$$H = \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \left(\bar{r}_{i\cdot} - \frac{N + 1}{2} \right)^2$$

$$= \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \bar{r}_{i\cdot}^2 - 3(N + 1)$$

pearson corr. coeff x_i 's & y_i 's
vs
 $\text{SRCC} \rightarrow$ rank(x_i) & rank(y_i)
=

Languages



العربية

Deutsch

Español

Français

한국어

Italiano

Português

Русский

Türkçe

文 7 more

Edit links

Method [edit]

- Rank all data from all groups together; i.e., rank the data from 1 to N ignoring group membership. Assign any tied values the average of the ranks they would have received had they not been tied.
- The test statistic is given by:

$$H = (N - 1) \frac{\sum_{i=1}^g n_i (\bar{r}_{i\cdot} - \bar{r})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (r_{ij} - \bar{r})^2}, \text{ where:}$$

✓ \bar{r} is the total number of observations across all groups

- g is the number of groups
- n_i is the number of observations in group i
- r_{ij} is the rank (among all observations) of observation j from group i

$\bar{r}_{i\cdot} = \frac{\sum_{j=1}^{n_i} r_{ij}}{n_i}$ is the average rank of all observations in group i

$\bar{r} = \frac{1}{2}(N + 1)$ is the average of all the r_{ij} .

- If the data contain no ties the denominator of the expression for H is exactly $(N - 1)N(N + 1)/12$ and $\bar{r} = \frac{N+1}{2}$. Thus

$$H = \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \left(\bar{r}_{i\cdot} - \frac{N + 1}{2} \right)^2$$

$$= \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \bar{r}_{i\cdot}^2 - 3(N + 1)$$

 \bar{r}_1 \bar{r}_2

⋮

 \bar{r}_N

across all

groups

[Download as PDF](#)[Printable version](#)

Languages



العربية

Deutsch

Español

Français

한국어

Italiano

Português

Русский

Türkçe

文 A 7 more

Edit links

7 External links

Method [edit]

- Rank all data from all groups together; i.e., rank the data from 1 to N ignoring group membership. Assign any tied values the average of the ranks they would have received had they not been tied.

- The test statistic is given by:

$$H = (N - 1) \frac{\sum_{i=1}^g n_i (\bar{r}_{i\cdot} - \bar{r})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (r_{ij} - \bar{r})^2}, \text{ where:}$$

- N is the total number of observations across all groups
- g is the number of groups
- n_i is the number of observations in group i
- r_{ij} is the rank (among all observations) of observation j from group i

$\bar{r}_{i\cdot} = \frac{\sum_{j=1}^{n_i} r_{ij}}{n_i}$ is the average rank of all observations in group i

$\bar{r} = \frac{1}{2}(N + 1)$ is the average of all the r_{ij} .

- If the data contain no ties the denominator of the expression for H is exactly $(N - 1)N(N + 1)/12$ and $\bar{r} = \frac{N+1}{2}$. Thus

$$H = \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \left(\bar{r}_{i\cdot} - \frac{N + 1}{2} \right)^2$$

$$= \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \bar{r}_{i\cdot}^2 - 3(N + 1)$$

$$\bar{x}_i - \bar{x}$$

$$x_{ij} - \bar{x}_i$$

[Download as PDF](#)[Printable version](#)

Languages

[العربية](#)[Deutsch](#)[Español](#)[Français](#)[한국어](#)[Italiano](#)[Português](#)[Русский](#)[Türkçe](#)[文 A 7 more](#)[Edit links](#)

7 External links

Method [edit]

1. Rank all data from all groups together; i.e., rank the data from 1 to N ignoring group membership. Assign any tied values the average of the ranks they would have received had they not been tied.

2. The test statistic is given by:

$$\{ H = (N - 1) \frac{\sum_{i=1}^g n_i (\bar{r}_{i\cdot} - \bar{r})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (r_{ij} - \bar{r})^2}, \text{ where:}$$

- N is the total number of observations across all groups
- g is the number of groups
- n_i is the number of observations in group i
- r_{ij} is the rank (among all observations) of observation j from group i
- $\bar{r}_{i\cdot} = \frac{\sum_{j=1}^{n_i} r_{ij}}{n_i}$ is the average rank of all observations in group i
- $\bar{r} = \frac{1}{2}(N + 1)$ is the average of all the r_{ij} .

3. If the data contain no ties the denominator of the expression for H is exactly $(N - 1)N(N + 1)/12$ and $\bar{r} = \frac{N+1}{2}$. Thus

$$H = \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \left(\bar{r}_{i\cdot} - \frac{N + 1}{2} \right)^2$$

$$= \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \bar{r}_{i\cdot}^2 - 3(N + 1)$$

f-stat

$$\frac{1}{N(N+1)} \sum_{i=1}^k n_i r_i = \sigma^2 (1 + \frac{1}{N})$$

The last formula only contains the squares of the average ranks.

4. A correction for ties if using the short-cut formula described in the previous point can be made by dividing H by

$$1 - \frac{\sum_{i=1}^G (t_i^3 - t_i)}{N^3 - N},$$
 where G is the number of groupings of different tied ranks, and t_i is the number of tied values within group i

that are tied at a particular value. This correction usually makes little difference in the value of H unless there are a large number of ties.

5. Finally, the decision to reject or not the null hypothesis is made by comparing H to a critical value H_c obtained from a table or a software for a given significance or alpha level. If H is bigger than H_c , the null hypothesis is rejected. If possible (no ties, sample not too big) one should compare H to the critical value obtained from the exact distribution of H . Otherwise, the distribution of H can be approximated by a chi-squared distribution with $g-1$ degrees of freedom. If some n_i values are small (i.e., less than 5) the exact probability distribution of H can be quite different from this chi-squared distribution. If a table of the chi-squared probability distribution is available, the critical value of chi-squared, $\chi_{\alpha;g-1}^2$, can be found by entering the table at $g-1$ degrees of freedom and looking under the desired significance or alpha level.

6. If the statistic is not significant, then there is no evidence of stochastic dominance between the samples. However, if the test is significant then at least one sample stochastically dominates another sample. Therefore, a researcher might use sample contrasts between individual sample pairs, or *post hoc* tests using Dunn's test, which (1) properly employs the same rankings as the Kruskal-Wallis test, and (2) properly employs the pooled variance implied by the null hypothesis of the Kruskal-Wallis test in order to determine which of the sample pairs are significantly different.^[5] When performing multiple sample contrasts or tests, the Type I error rate tends to become inflated, raising concerns about multiple comparisons.

Exact probability tables [edit]

A large amount of software exists for performing the Kruskal-Wallis test. Existing software only



[Cite this page](#)[Wikidata item](#)[Print/export](#)[Download as PDF](#)[Printable version](#)

Languages

[العربية](#)[Deutsch](#)[Español](#)[Français](#)[한국어](#)[Italiano](#)[Português](#)[Русский](#)[Türkçe](#)[A 7 more](#)[Edit links](#)

3 Exact distribution of H

4 See also

5 References

6 Further reading

7 External links

Method [\[edit\]](#)

1. Rank all data from all groups together; i.e., rank the data from 1 to N ignoring group membership. Assign any tied values the average of the ranks they would have received had they not been tied.

2. The test statistic is given by:

$$H = (N - 1) \frac{\sum_{i=1}^g n_i (\bar{r}_{i\cdot} - \bar{r})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (r_{ij} - \bar{r})^2}, \text{ where:}$$

- N is the total number of observations across all groups
- g is the number of groups
- n_i is the number of observations in group i
- r_{ij} is the rank (among all observations) of observation j from group i

- $\bar{r}_{i\cdot} = \frac{\sum_{j=1}^{n_i} r_{ij}}{n_i}$ is the average rank of all observations in group i

- $\bar{r} = \frac{1}{2}(N + 1)$ is the average of all the r_{ij} .

3. If the data contain no ties the denominator of the expression for H is exactly $(N - 1)N(N + 1)/12$ and $\bar{r} = \frac{N+1}{2}$. Thus

$$H = \frac{12}{N(N+1)} \sum_{i=1}^g n_i (\bar{r}_{i\cdot} - \bar{r})^2$$

groups. A one-way between groups ANOVA comparing just two groups will give you the same results at the independent t test that you conducted in [Lesson 8](#). We will use the five step hypothesis testing procedure again in this lesson.

1. Check assumptions and write hypotheses

The assumptions for a one-way between groups ANOVA are:

1. Samples are independent
2. The response variable is approximately normally distributed for each group or all group sample sizes are at least 30
3. The population variances are equal across responses for the group levels (if the largest sample standard deviation divided by the smallest sample standard deviation is **not** greater than two, then assume that the population variances are equal)

Given that you are comparing k independent groups, the null and alternative hypotheses are:

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_k$$

$$H_a : \text{Not all } \mu_i \text{ are equal}$$

In other words, the null hypothesis is that all of the groups' population means are equal. The alternative is that they are not all equal; there are at least two population means that are not equal to one another.

2. Calculate the test statistic

ANOVA uses an F test statistic. Hand



- ▶ 4: Confidence Intervals
- ▶ 5: Hypothesis Testing, Part 1
- ▶ 6: Hypothesis Testing, Part 2
- ▶ 7: Normal Distributions
- ▶ 8: Inference for One Sample
- ▶ 9: Inference for Two Samples
- ▼ 10: One-Way ANOVA
 - 10.1 - Introduction to the F Distribution
 - 10.2 - Hypothesis Testing
 - 10.3 - Pairwise Comparisons
 - 10.4 - Minitab: One-Way ANOVA
 - 10.5 - Example: SAT-Math Scores by Award Preference
 - 10.6 - Example: Exam Grade by Professor
 - 10.7 - Lesson 10 Summary
- ▶ 11: Chi-Square Tests
- ▶ 12: Correlation & Simple Linear Regression

Resources

- Datasets Glossary
Formulas Contact



[Printable version](#)

7 External links

Languages

[العربية](#)[Deutsch](#)[Español](#)[Français](#)[한국어](#)[Italiano](#)[Português](#)[Русский](#)[Türkçe](#)[文 A 7 more](#)[Edit links](#)

Method

[\[edit \]](#)

1. Rank all data from all groups together; i.e., rank the data from 1 to N ignoring group membership. Assign any tied values the average of the ranks they would have received had they not been tied.
2. The test statistic is given by:

$$H = (N - 1) \frac{\sum_{i=1}^g n_i (\bar{r}_{i\cdot} - \bar{r})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (r_{ij} - \bar{r})^2}, \text{ where:}$$

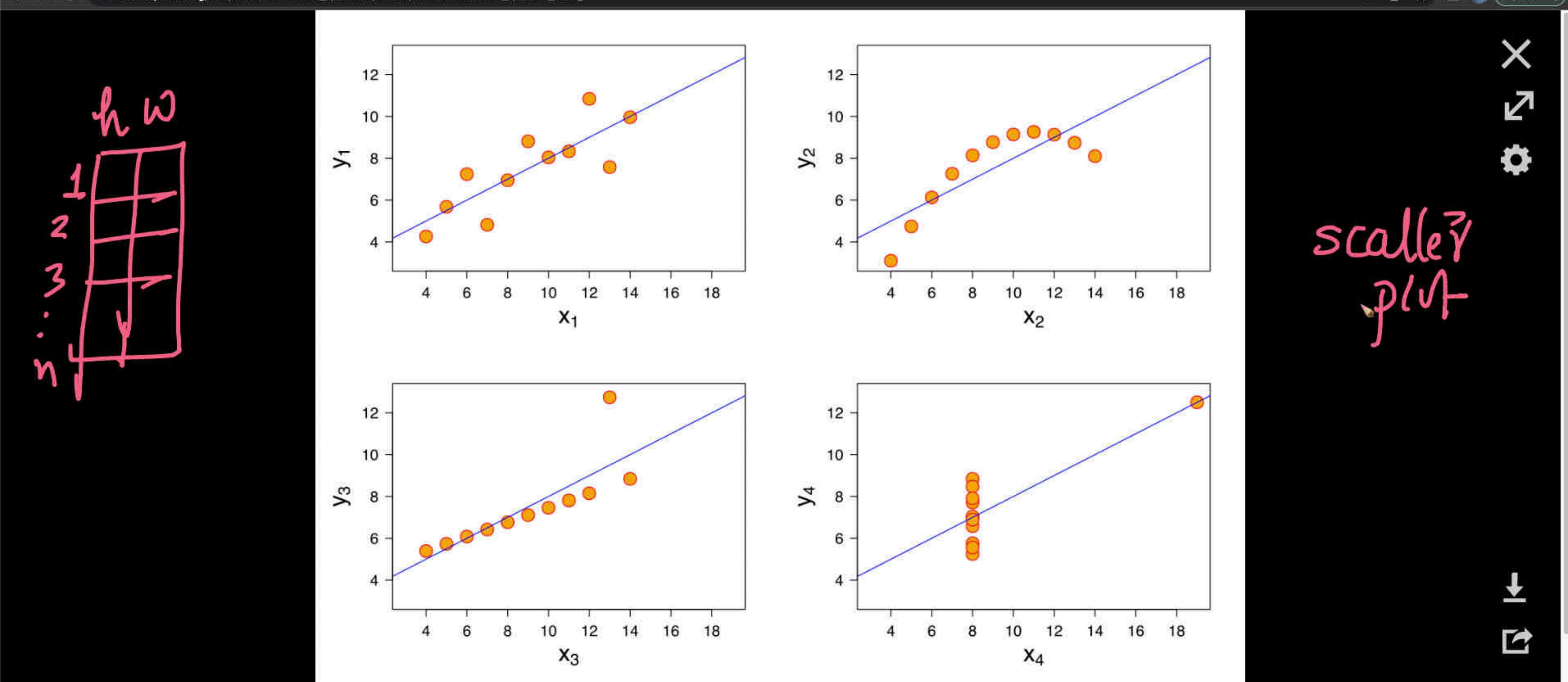
$$\chi^2 (g-1)$$

- N is the total number of observations across all groups
- g is the number of groups
- n_i is the number of observations in group i
- r_{ij} is the rank (among all observations) of observation j from group i
- $\bar{r}_{i\cdot} = \frac{\sum_{j=1}^{n_i} r_{ij}}{n_i}$ is the average rank of all observations in group i
- $\bar{r} = \frac{1}{2}(N + 1)$ is the average of all the r_{ij} .

3. If the data contain no ties the denominator of the expression for H is exactly $(N - 1)N(N + 1)/12$ and $\bar{r} = \frac{N+1}{2}$. Thus

$$H = \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \left(\bar{r}_{i\cdot} - \frac{N + 1}{2} \right)^2$$

$$= \frac{12}{N(N + 1)} \sum_{i=1}^g n_i \bar{r}_{i\cdot}^2 - 3(N + 1)$$



All four sets are identical when examined using simple summary statistics, but vary considerably when graphed

More details

[Permanent link](#)[Page information](#)[Cite this page](#)[Wikidata item](#)[Print/export](#)[Download as PDF](#)[Printable version](#)[In other projects](#)[Wikimedia Commons](#)[Languages](#)[Deutsch](#)[Español](#)[Français](#)[한국어](#)[日本語](#)[Português](#)[Русский](#)[Tiếng Việt](#)[中文](#)[11 more](#)[Edit links](#)

Data [\[edit\]](#)

For all four datasets:

Property	Value	Accuracy
Mean of x	9.8	exact
Sample variance of x : s_x^2	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of y : s_y^2	4.125	± 0.003
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression : R^2	0.67	to 2 decimal places

- The first [scatter plot](#) (top left) appears to be a simple [linear relationship](#), corresponding to two [variables](#) correlated where y could be modelled as [gaussian](#) with mean linearly dependent on x .
- The second graph (top right); while a relationship between the two variables is obvious, it is not linear, and the [Pearson correlation coefficient](#) is not relevant. A more general regression and the corresponding [coefficient of determination](#) would be more appropriate.
- In the third graph (bottom left), the modelled relationship is linear, but should have a different [regression line](#) (a [robust regression](#) would have been called for). The calculated regression is offset by the one [outlier](#) which exerts enough influence to lower the correlation coefficient from 1 to 0.816.
- Finally, the fourth graph (bottom right) shows an example when one [high-leverage point](#) is enough to produce a high correlation coefficient, even between the variables.

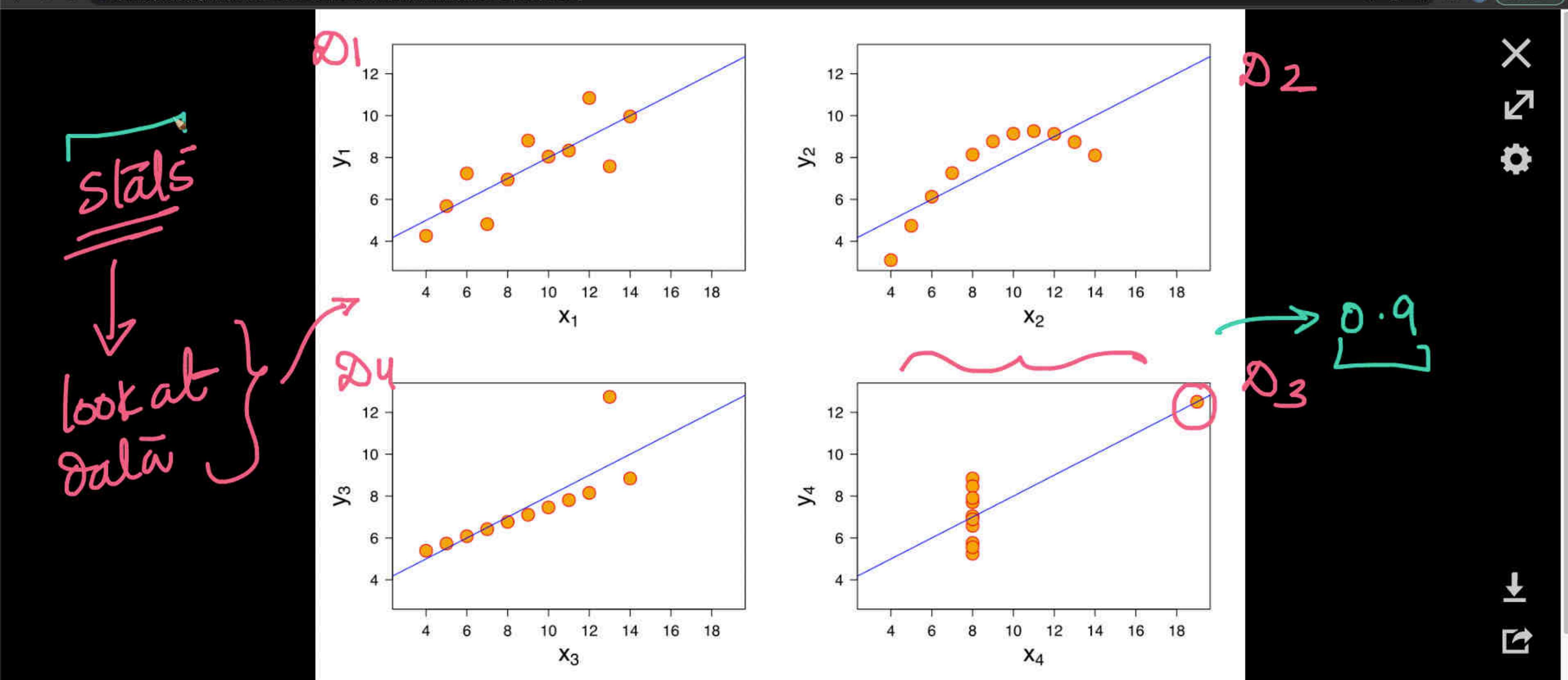
[Permanent link](#)[Page information](#)[Cite this page](#)[Wikidata item](#)[Print/export](#)[Download as PDF](#)[Printable version](#)[In other projects](#)[Wikimedia Commons](#)[Languages](#)[Deutsch](#)[Español](#)[Français](#)[한국어](#)[日本語](#)[Português](#)[Русский](#)[Tiếng Việt](#)[中文](#)[11 more](#)[Edit links](#)

Data [\[edit\]](#)

For all four datasets:

Property	Value	Accuracy
Mean of x	9 ✓	exact
Sample variance of x : s_x^2	11 ✓	exact
Mean of y	7.50 ✓	to 2 decimal places
Sample variance of y : s_y^2	4.125 ✓	± 0.003
Correlation between x and y	0.816 ✓	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression : R^2	0.67	to 2 decimal places

- The first [scatter plot](#) (top left) appears to be a simple [linear relationship](#), corresponding to two [variables](#) correlated where y could be modelled as [gaussian](#) with mean linearly dependent on x .
- The second graph (top right); while a relationship between the two variables is obvious, it is not linear, and the [Pearson correlation coefficient](#) is not relevant. A more general regression and the corresponding [coefficient of determination](#) would be more appropriate.
- In the third graph (bottom left), the modelled relationship is linear, but should have a different [regression line](#) (a [robust regression](#) would have been called for). The calculated regression is offset by the one [outlier](#) which exerts enough influence to lower the correlation coefficient from 1 to 0.816.
- Finally, the fourth graph (bottom right) shows an example when one [high-leverage point](#) is enough to produce a high correlation coefficient, even between the variables.



All four sets are identical when examined using simple summary statistics, but vary considerably when graphed

More details

Biased estimate

Data → param

$x_1 \dots x_n$

$\overbrace{\hspace{10em}}$

m

s

$\tilde{\mu}$

μ

σ

{ } (brace)

$X : x_1 \dots x_n$

estimate Θ : param

param = θ

$\hat{\theta} \approx \theta$
 $\hat{\theta} \neq \theta \rightarrow$

$\theta: \Gamma$
 $\hat{\theta}: \text{estimated using data}$

data follow some unknown distribution $P(x | \theta)$ (where θ is a fixed, unknown constant that is part of this distribution), and then we construct some estimator $\hat{\theta}$ that maps observed data to values that we hope are close to θ . The **bias** of $\hat{\theta}$ relative to θ is defined as^[1]

$$\text{Bias}(\hat{\theta}, \theta) = \text{Bias}_\theta[\hat{\theta}] = \mathbb{E}_{x|\theta}[\hat{\theta}] - \theta = \mathbb{E}_{x|\theta}[\hat{\theta} - \theta],$$

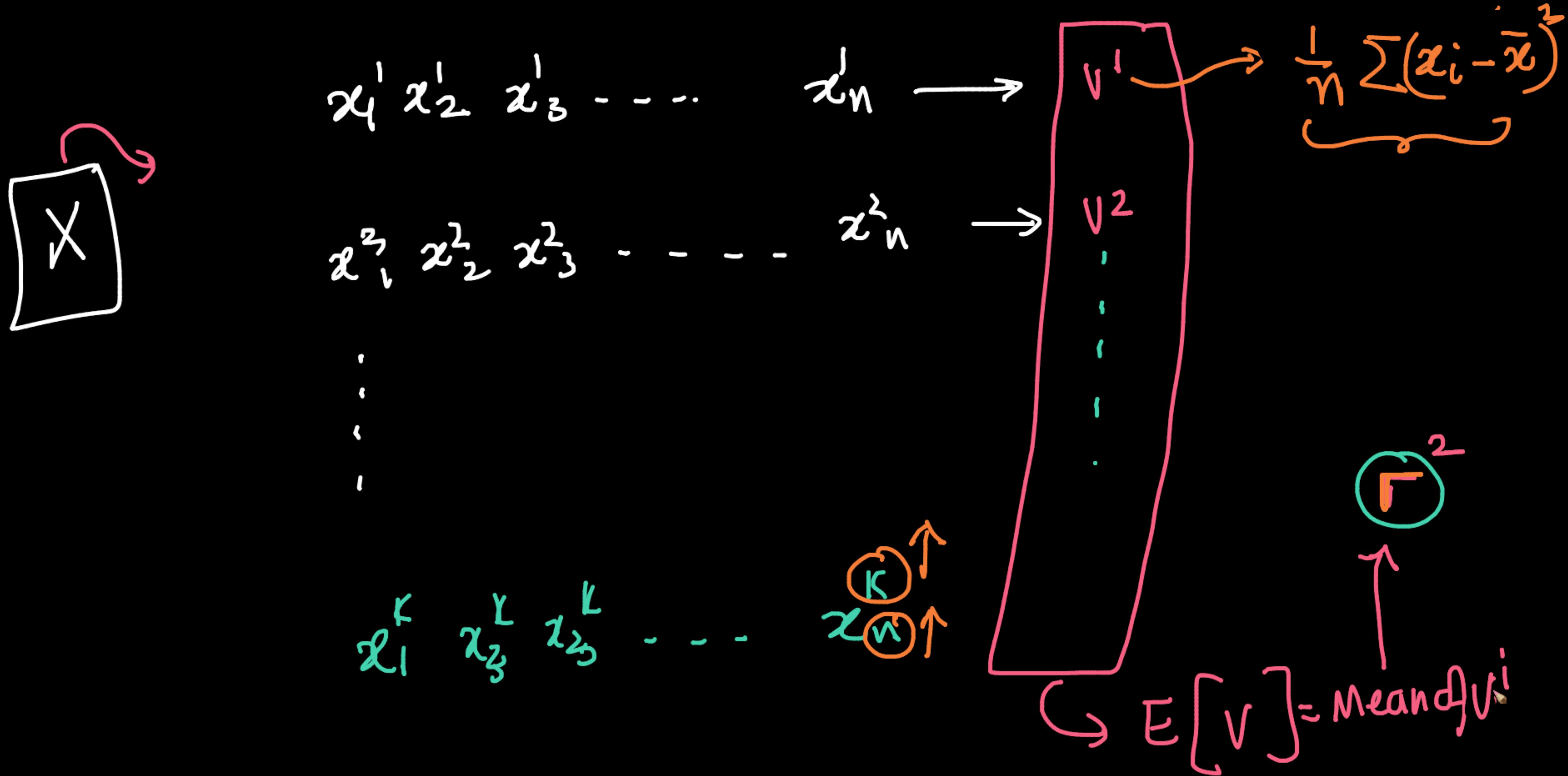
where $\mathbb{E}_{x|\theta}$ denotes **expected value** over the distribution $P(x | \theta)$ (i.e., averaging over all possible observations x). The second equation follows since θ is measurable with respect to the conditional distribution $P(x | \theta)$.

An estimator is said to be **unbiased** if its bias is equal to zero for all values of parameter θ , or equivalently, if the expected value of the estimator matches that of the parameter.^[2]

In a simulation experiment concerning the properties of an estimator, the bias of the estimator may be assessed using the **mean signed difference**.

Examples [edit]

Sample variance [edit]



want:

$$E[s] = r \quad \text{ideally}$$

$$\text{Bias}(s) = r - \bar{E}[s] \quad \text{as low as possible}$$

variance

if $\text{var} = \frac{1}{n} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2$ on a sample of n pts

ideally

$$E[\text{var}] = \sigma^2$$

But

this formula
+ Math

$$E[s^2] \neq \sigma^2$$

Dividing instead by $n - 1$ yields an unbiased estimator. Conversely, MSE can be minimized by dividing by a different number (depending on distribution), but this results in a biased estimator. This number is always larger than $n - 1$, so this is known as a **shrinkage estimator**, as it "shrinks" the unbiased estimator towards zero; for the normal distribution the optimal value is $n + 1$.

Suppose X_1, \dots, X_n are **independent and identically distributed** (i.i.d.) random variables with **expectation** μ and **variance** σ^2 .

If the **sample mean** and uncorrected **sample variance** are defined as

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

then S^2 is a biased estimator of σ^2 , because

$$E[S^2] = E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] = E\left[\frac{1}{n} \sum_{i=1}^n ((X_i - \mu) - (\bar{X} - \mu))^2\right]$$

$$= E\left[\frac{1}{n} \sum_{i=1}^n ((X_i - \mu)^2 - 2(\bar{X} - \mu)(X_i - \mu) + (\bar{X} - \mu)^2)\right]$$

$$= E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \frac{2}{n}(\bar{X} - \mu) \sum_{i=1}^n (X_i - \mu) + \frac{1}{n}(\bar{X} - \mu)^2 \sum_{i=1}^n 1\right]$$

$$= E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \frac{2}{n}(\bar{X} - \mu) \sum_{i=1}^n (X_i - \mu) + \frac{1}{n}(\bar{X} - \mu)^2 \cdot n\right]$$

Dividing instead by $n - 1$ yields an unbiased estimator. Conversely, MSE can be minimized by dividing by a different number (depending on distribution), but this results in a biased estimator. This number is always larger than $n - 1$, so this is known as a **shrinkage estimator**, as it "shrinks" the unbiased estimator towards zero; for the normal distribution the optimal value is $n + 1$.

Suppose X_1, \dots, X_n are **independent and identically distributed** (i.i.d.) random variables with **expectation** μ and **variance** σ^2 . If the **sample mean** and uncorrected **sample variance** are defined as

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

then S^2 is a biased estimator of σ^2 , because

$$\begin{aligned} E[S^2] &= E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] = E\left[\frac{1}{n} \sum_{i=1}^n ((X_i - \mu) - (\bar{X} - \mu))^2\right] \\ &= E\left[\frac{1}{n} \sum_{i=1}^n \left((X_i - \mu)^2 - 2(\bar{X} - \mu)(X_i - \mu) + (\bar{X} - \mu)^2\right)\right] \\ &= E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \frac{2}{n}(\bar{X} - \mu) \sum_{i=1}^n (X_i - \mu) + \frac{1}{n}(\bar{X} - \mu)^2 \sum_{i=1}^n 1\right] \\ &= E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \frac{2}{n}(\bar{X} - \mu) \sum_{i=1}^n (X_i - \mu) + \frac{1}{n}(\bar{X} - \mu)^2 \cdot n\right] \end{aligned}$$

then S^2 is a biased estimator of σ^2 , because

$$\begin{aligned} \mathbb{E}[S^2] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n ((X_i - \mu) - (\bar{X} - \mu))^2\right] \\ &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n ((X_i - \mu)^2 - 2(\bar{X} - \mu)(X_i - \mu) + (\bar{X} - \mu)^2)\right] \\ &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \frac{2}{n}(\bar{X} - \mu) \sum_{i=1}^n (X_i - \mu) + \frac{1}{n}(\bar{X} - \mu)^2 \sum_{i=1}^n 1\right] \\ &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \frac{2}{n}(\bar{X} - \mu) \sum_{i=1}^n (X_i - \mu) + \frac{1}{n}(\bar{X} - \mu)^2 \cdot n\right] \\ &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \frac{2}{n}(\bar{X} - \mu) \sum_{i=1}^n (X_i - \mu) + (\bar{X} - \mu)^2\right] \end{aligned}$$



To continue, we note that by subtracting μ from both sides of $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$, we get

$$\bar{X} - \mu = \frac{1}{n} \sum_{i=1}^n X_i - \mu = \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} \sum_{i=1}^n (X_i - \mu).$$

$$= E \left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - 2(\bar{X} - \mu)^2 + (\bar{X} - \mu)^2 \right]$$

$$= E \left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X} - \mu)^2 \right]$$

$$= E \left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \right] - E \left[(\bar{X} - \mu)^2 \right]$$

$$= \sigma^2 - E \left[(\bar{X} - \mu)^2 \right] = \boxed{\left(1 - \frac{1}{n} \right) \sigma^2} < \sigma^2.$$

$$\bar{E}[S^2] = \left(1 - \frac{1}{n} \right) \sigma^2$$

This can be seen by noting the following formula, which follows from the [Bienaymé formula](#), for the term in the inequality

for the expectation of the uncorrected sample variance above: $E [(\bar{X} - \mu)^2] = \frac{1}{n} \sigma^2$.

In other words, the expected value of the uncorrected sample variance does not equal the population variance σ^2 , unless multiplied by a normalization factor. The sample mean, on the other hand, is an unbiased^[3] estimator of the population mean μ .^[2]

Note that the usual definition of sample variance is $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$, and this is an unbiased estimator of the population

$$\begin{aligned} &= E \left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X} - \mu)^2 \right] \\ &= E \left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \right] - E \left[(\bar{X} - \mu)^2 \right] \\ &= \sigma^2 - E \left[(\bar{X} - \mu)^2 \right] = \left(1 - \frac{1}{n}\right) \sigma^2 < \sigma^2. \end{aligned}$$

$$\begin{array}{c} E[S^2] = \sigma^2 \\ \downarrow \\ \left(1 - \frac{1}{n}\right) \sigma^2 + \sigma^2 \end{array}$$

This can be seen by noting the following formula, which follows from the [Bienaymé formula](#), for the term in the inequality for the expectation of the uncorrected sample variance above: $E[(\bar{X} - \mu)^2] = \frac{1}{n} \sigma^2$.

In other words, the expected value of the uncorrected sample variance does not equal the population variance σ^2 , unless multiplied by a normalization factor. The sample mean, on the other hand, is an unbiased^[3] estimator of the population mean μ .^[2]

Note that the usual definition of sample variance is $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$, and this is an unbiased estimator of the population variance.

Algebraically speaking, $E[S^2]$ is unbiased because:

$$E[S^2] = E \left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \frac{1}{n-1} \sum_{i=1}^n E[(X_i - \bar{X})^2] = \frac{1}{n-1} \sum_{i=1}^n (E[X_i] - E[\bar{X}])^2 = \frac{1}{n-1} \sum_{i=1}^n (\mu - \mu)^2 = \frac{1}{n-1} \sum_{i=1}^n 0 = 0 = \sigma^2$$

$i=1$

$$\begin{aligned} &= E \left[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \right] - E \left[(\bar{X} - \mu)^2 \right] \\ &= \sigma^2 - E \left[(\bar{X} - \mu)^2 \right] = \left(1 - \frac{1}{n} \right) \sigma^2 < \sigma^2. \end{aligned}$$

σ^2

This can be seen by noting the following formula, which follows from the [Bienaymé formula](#), for the term in the inequality for the expectation of the uncorrected sample variance above: $E[(\bar{X} - \mu)^2] = \frac{1}{n}\sigma^2$.

In other words, the expected value of the uncorrected sample variance does not equal the population variance σ^2 , unless multiplied by a normalization factor. The sample mean, on the other hand, is an unbiased^[3] estimator of the population mean μ .^[2]

Note that the usual definition of sample variance is $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$, and this is an unbiased estimator of the population variance.

Algebraically speaking, $E[S^2]$ is unbiased because:

$$E[S^2] = E \left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \frac{n}{n-1} E \left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right]$$

39 / 39

for the expectation of the uncorrected sample variance above: $E[(\bar{X} - \mu)^2] = \frac{1}{n}\sigma^2$.

In other words, the expected value of the uncorrected sample variance does not equal the population variance σ^2 , unless multiplied by a normalization factor. The sample mean, on the other hand, is an unbiased^[3] estimator of the population mean μ .^[2]

Note that the usual definition of sample variance is $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$, and this is an unbiased estimator of the population variance.

Algebraically speaking, $E[S^2]$ is unbiased because:

$$\begin{aligned} E[S^2] &= E\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right] = \frac{n}{n-1} E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] \\ &= \frac{n}{n-1} \left(1 - \frac{1}{n}\right) \sigma^2 = \sigma^2, \end{aligned}$$

$E[S^2] = \sigma^2$

$n \rightarrow \infty$

where the transition to the second line uses the result derived above for the biased estimator. Thus $E[S^2] = \sigma^2$, and therefore $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is an unbiased estimator of the population variance, σ^2 . The ratio between the biased (uncorrected) and unbiased estimates of the variance is known as [Bessel's correction](#).

The reason that



on the fact that the sample mean is an ordinary

Var:

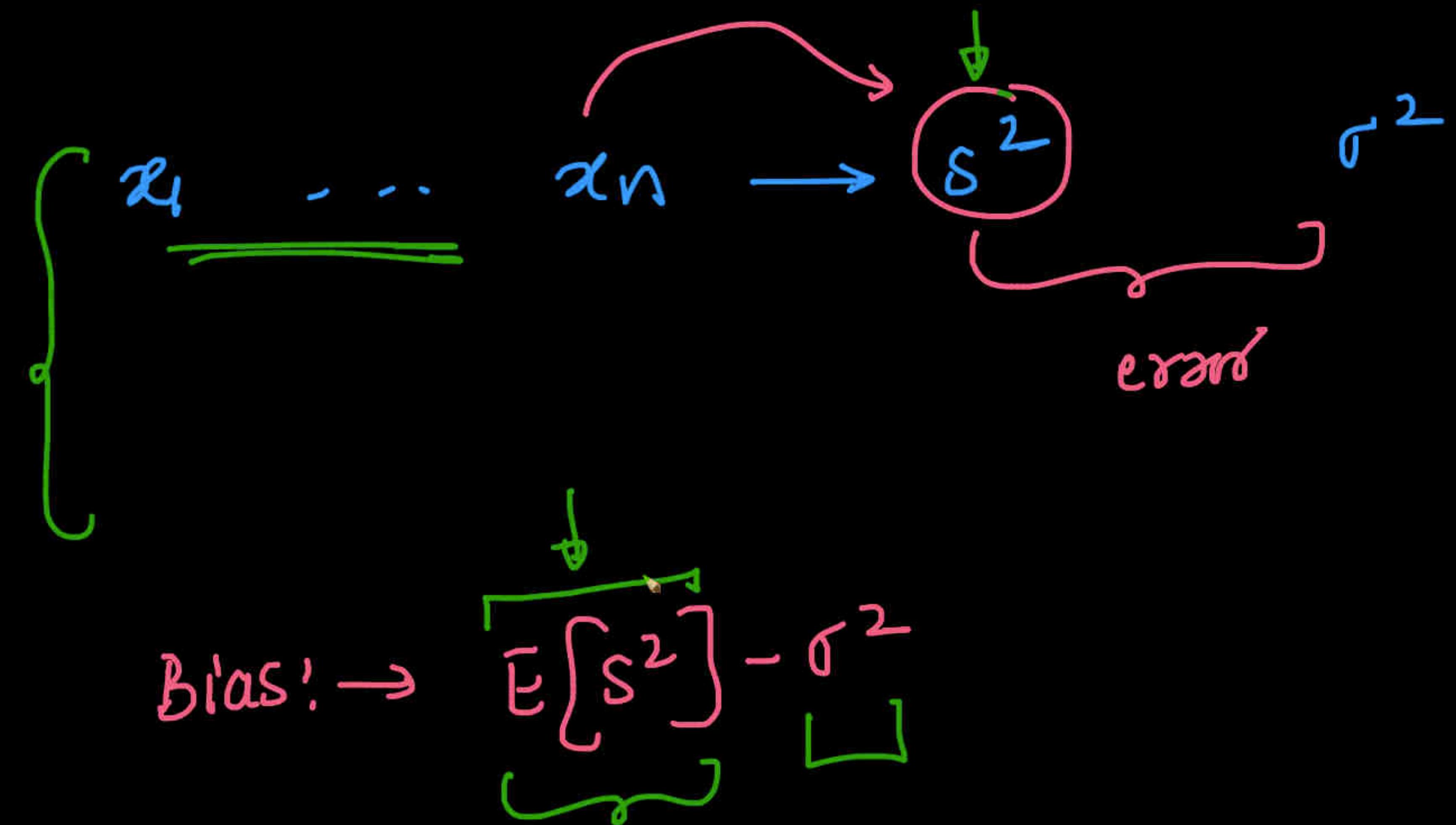
$$\left\{ \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2 \right.$$

if n is small $\Rightarrow n$ & $n-1$ will
change the value



div $\frac{n-1}{n}$ doesn't matter as much

Bias - Variance tradeoff \rightarrow ML module



+ Code + Text

Reconnect



Task: Determine the eligibility for granting Home loan.

Objective of this notebook is:

1. To understand the patterns in the data.
2. How to Handle the categorical features.
3. How to deal with missing data.
4. Feature Engineering
5. Finding the most important features while taking the decision of granting a loan application.
6. Understanding the Normalization and standardisation of the data.

▶ Load data and libraries

[] ↳ 11 cells hidden



▶ Basic Data Exploration



[] ↳ 9 cells hidden

+ Code + Text

Reconnect



Task: Determine the eligibility for granting Home loan.

Objective of this notebook is:

1. To understand the patterns in the data.
2. How to Handle the categorical features.
3. How to deal with missing data.
4. Feature Engineering
5. Finding the most important features while taking the decision of granting a loan application.
6. Understanding the Normalization and standardisation of the data.

✓ → discrete r.v

▶ Load data and libraries

[] ↳ 11 cells hidden



▶ Basic Data Exploration



[] ↳ 9 cells hidden

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=x9qeM0-oClpO Reconnect + Text

```
+ Code + Text
import seaborn as sns

import warnings
warnings.filterwarnings('ignore')

{x}
[ ] #Data: https://drive.google.com/file/d/1oJbdRpTLqPu1SIBXHkzWRaLaZbvZot7w/view?usp=sharing
# Download data
{
id = "1oJbdRpTLqPu1SIBXHkzWRaLaZbvZot7w"
path = "https://docs.google.com/uc?export=download&id=" + id
print(path)

https://docs.google.com/uc?export=download&id=1oJbdRpTLqPu1SIBXHkzWRaLaZbvZot7w
[ ] !wget "https://docs.google.com/uc?export=download&id=1oJbdRpTLqPu1SIBXHkzWRaLaZbvZot7w" -O train.csv
--2022-07-25 13:37:17-- https://docs.google.com/uc?export=download&id=1oJbdRpTLqPu1SIBXHkzWRaLaZbvZot7w
Resolving docs.google.com (docs.google.com)... 172.253.62.138, 172.253.62.102, 172.253.62.100, ...
Connecting to docs.google.com (docs.google.com)|172.253.62.138|:443... connected.
HTTP request sent, awaiting response... 303 See Other
Location: https://doc-0o-90-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbpl/d2lf91u0ibh30d
Warning: wildcards not supported in HTTP.
--2022-07-25 13:37:17-- https://doc-0o-90-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbpl
Resolving doc-0o-90-docs.googleusercontent.com (doc-0o-90-docs.googleusercontent.com)... 172.217.15.65, 2607:f8b0:40
Connecting to doc-0o-90-docs.googleusercontent.com (doc-0o-90-docs.googleusercontent.com)|172.217.15.65|:443... conn
HTTP request sent, awaiting response... 200 OK
Length: 38011 (37K) [text/csv]
```

10.2 - Hypothesis Te x | Kruskal-Wallis one-w x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab x | Markov's inequality x | Chebyshev's inequality x | scipy.stats.kruskal x | numpy.var — NumPy x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=Ej8pWNLIGL-i

+ Code + Text Reconnect

TP001302 Female No 1 Graduate Yes 7451 0 360 1 Continue

data = pd.read_csv('./train.csv')
data.shape

(614, 13)

[] data.columns

Index(['Loan_ID', 'Gender', 'Married', 'Dependents', 'Education',
'Self_Employed', 'ApplicantIncome', 'CoapplicantIncome', 'LoanAmount',
'Loan_Amount_Term', 'Credit_History', 'Property_Area', 'Loan_Status'],
dtype='object')

data.head()

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loa
0	LP001002	Male	No	0	Graduate	No	5849	0.0	NaN	
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	

47 / 47

+ Code + Text

Reconnect ▾



L

TP001202 Female No. 1 Graduate Year 7AF1 A 260 1 Commission 3

插入工具设置

7

```
data = pd.read_csv('./train.csv')  
data.shape
```

(614, 13)

```
[ ] data.columns
```

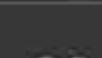
```
Index(['Loan_ID', 'Gender', 'Married', 'Dependents', 'Education',  
       'Self_Employed', 'ApplicantIncome', 'CoapplicantIncome', 'LoanAmount',  
       'Loan_Amount_Term', 'Credit_History', 'Property_Area', 'Loan_Status'],  
      dtype='object')
```

▶ `data.head()`

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loa
0	LP001002	Male	No	0	Graduate	No	5849	0.0	Nan	
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	

10.2 - Hypothesis Te x | Kruskal-Wallis one-w x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab x | Markov's inequality x | Chebyshev's inequality x | scipy.stats.kruskal x | numpy.var — NumPy x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=Ej8pWNLIGL-i

+ Code + Text Reconnect  

[] `data.dtypes`
#object => typically categorical/IDs
#Int64, Float64

{x}

Loan_ID object
Gender object
Married object
Dependents object
Education object
Self_Employed object

ApplicantIncome int64
CoapplicantIncome float64
LoanAmount float64
Loan_Amount_Term float64
Credit_History float64
Property_Area object
Loan_Status object

dtype: object

feature ↗

Categorical (discrete, r.v)

[] `data['Dependents'].value_counts()`

	value	count
0	345	
1	102	
2	101	

49 / 49

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=Ej8pWNLIGL-i

Update

+ Code + Text

Reconnect



{x}

```
[ ] data.dtypes  
{x} #object => typically categorical/IDs  
#Int64, Float64
```

Loan_ID	object
Gender	object
Married	object
Dependents	object
Education	object
Self_Employed	object
ApplicantIncome	int64
CoapplicantIncome	float64
LoanAmount	float64
Loan_Amount_Term	float64
Credit_History	float64
Property_Area	object
Loan_Status	object
dtype:	object

object

int64

float64

float64

float64

float64

object

int64

float64

object

int64

float64

object

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=Ej8pWNLIGL-i

Update

+ Code + Text

Reconnect



{x}

```
[ ] data.dtypes  
{x} #object => typically categorical/IDs  
#Int64, Float64
```

Loan_ID	object
Gender	object
Married	object
Dependents	object
Education	object
Self_Employed	object
ApplicantIncome	int64
CoapplicantIncome	float64
LoanAmount	float64
Loan_Amount_Term	float64
Credit_History	float64
Property_Area	object
Loan_Status	object
dtype:	object

```
<> [ ] data['Dependents'].value_counts()
```

0	345
1	102
2	101



+ Code + Text

Reconnect

dtype: object

```
[ ] data['Dependents'].value_counts()
```

	Dependents
0	345
1	102
2	101
3+	51

Name: Dependents, dtype: int64

```
[ ] # drop loanID column  
data = data.drop('Loan_ID', axis = 1)
```

Is App pm
y/N

C₁
C₂
C_i

Basic Data Exploration

{ 9 cells hidden

Basic Data visualization: Univariate

[] 11 cells hidden

52 / 53

+ Code + Text

Reconnect

Load data and libraries

↳ 11 cells hidden

Basic Data Exploration

[] ↳ 9 cells hidden

Basic Data visualization: Univariate

[] ↳ 11 cells hidden

Simple Feature Engineering

[] ↳ 37 cells hidden

SRCC

dep · vs income

numerous

numerous

Page-Footer 53 / 54

+ Code + Text

Reconnect ▾



✓

▼ Basic Data Exploration

100

```
[ ] data.describe()  
# only numeric features
```

ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
614.000000	614.000000	592.000000	600.000000	564.000000
5403.459283	1621.245798	146.412162	342.000000	0.842199
6109.041673	2926.248369	85.587325	65.12041	0.364878
150.000000	0.000000	9.000000	12.00000	0.000000
2877.500000	0.000000	100.000000	360.000000	1.000000
3812.500000	1188.500000	128.000000	360.000000	1.000000
5795.000000	2297.250000	168.000000	360.000000	1.000000
81000.000000	41667.000000	700.000000	480.000000	1.000000

```
[1] # categorical features
```

10.2 - Hypothesis Te x | Kruskal-Wallis one-w x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab x | Markov's inequality x | Chebyshev's inequality x | scipy.stats.kruskal x | numpy.var — NumPy x | + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=TarGSsXDG8Tp

+ Code + Text Reconnect

Basic Data Exploration

{x}

[] data.describe()
only numeric features

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	592.000000	600.000000	564.000000
mean	5403.459283	1621.245798	146.412162	342.000000	0.842199
std	6109.041673	2926.248369	85.587325	65.12041	0.364878
min	150.000000	{ 0.000000	9.000000	12.00000	0.000000
25%	2877.500000	{ 0.000000	100.000000	360.000000	1.000000
50%	3812.500000	{ 1188.500000	128.000000	360.000000	1.000000
75%	5795.000000	{ 2297.250000	168.000000	360.000000	1.000000
max	81000.000000	41667.000000	700.000000	480.000000	1.000000

< > # categorical features

55 / 56

10.2 - Hypothesis Te x | Kruskal-Wallis one-w x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab x | Markov's inequality x | Chebyshev's inequality x | scipy.stats.kruskal x | numpy.var — NumPy x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=TarGSsXDG8Tp

+ Code + Text Reconnect

Basic Data Exploration

{x}

[] data.describe()
only numeric features

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	592.000000	600.000000	564.000000
mean	5403.459283	1621.245798	146.412162	342.000000	0.842199
std	6109.041673	2926.248369	85.587325	65.12041	0.364878
min	150.000000	0.000000	9.000000	12.00000	0.000000
25%	2877.500000	0.000000	100.000000	360.000000	1.000000
50%	3812.500000	1188.500000	128.000000	360.000000	1.000000
75%	5795.000000	2297.250000	168.000000	360.000000	1.000000
max	81000.000000	41667.000000	700.000000	480.000000	1.000000

categorical features

56 / 57

10.2 - Hypothesis Te x | Kruskal-Wallis one-w x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab x | Markov's inequality x | Chebyshev's inequality x | scipy.stats.kruskal x | numpy.var — NumPy x | + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=TarGSsXDG8Tp

+ Code + Text Reconnect

Basic Data Exploration

```
[ ] data.describe()  
# only numeric features
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	592.000000	600.000000	564.000000
mean	5403.459283	1621.245798	146.412162	342.000000	0.842199
std	6109.041673	2926.248369	85.587325	65.12041	0.364878
min	150.000000	0.000000	9.000000	12.00000	0.000000
25%	2877.500000	0.000000	100.000000	360.000000	1.000000
50%	3812.500000	1188.500000	128.000000	360.000000	1.000000
75%	5795.000000	2297.250000	168.000000	360.000000	1.000000
max	81000.000000	41667.000000	700.000000	480.000000	1.000000

categorical features

Reconnect

Update

Up Down Reload Print Copy Delete More

57 / 58

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=BHXHp28hG-uh

Update

+ Code + Text

Reconnect



max 81000.000000 41667.000000 700.000000 480.000000 1.000000

↑ ↓ ⌂ ⚙ 📄 🗑 :

[] # catgeorical features
{x} data.describe(include = ['object'])

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
count	601	611	599	614	582	614	614
unique	2	2	4	2	2	3	2
top	Male	Yes	0	Graduate	No	Semiurban	Y
freq	489	398	345	480	500	233	422

[] #missing values
data.isna().sum()

Gender	13
Married	3
Dependents	15
Education	0
Self_Employed	32
ApplicantIncome	0
CoapplicantIncome	0
LoanAmount	22



10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=BHXHp28hG-uh

+ Code + Text Reconnect

max 81000.000000 41667.000000 700.000000 480.000000 1.000000

[] # catgeorical features
{x} data.describe(include = ['object'])

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
count	601	611	599	614	582	614	614
unique	2	2	4	2	2	3	2
top	Male	Yes	0	Graduate	No	Semiurban	Y
freq	489	398	345	480	500	233	422

[] #missing values
data.isna().sum()

Gender	13
Married	3
Dependents	15
Education	0
Self_Employed	32
ApplicantIncome	0
CoapplicantIncome	0
LoanAmount	22

59 / 60

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=BHXHp28hG-uh

+ Code + Text Reconnect

max 81000.000000 41667.000000 700.000000 480.000000 1.000000

[] # catgeorical features
{x} data.describe(include = ['object'])

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
count	601	611	599	614	582	614	614
unique	2	2	4	2	2	3	2
top	Male	Yes	0	Graduate	No	Semiurban	Y
freq	489	398	345	480	500	233	422

[] #missing values
data.isna().sum()

Gender	13
Married	3
Dependents	15
Education	0
Self_Employed	32
ApplicantIncome	0
CoapplicantIncome	0
LoanAmount	22

60 / 61

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=BHXHp28hG-uh Reconnect + ⚙️ 🔍

+ Code + Text

max 81000.000000 41667.000000 700.000000 480.000000 1.000000

[] # catgeorical features
{x} data.describe(include = ['object'])

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
count	601	611	599	614	582	614	614
unique	2	2	4	2	2	3	2
top	Male	Yes	0	Graduate	No	Semiurban	Y
freq	489	398	345	480	500	233	422

[] #missing values
data.isna().sum()

Gender	13
Married	3
Dependents	15
Education	0
Self_Employed	32
ApplicantIncome	0
CoapplicantIncome	0
LoanAmount	22

61 / 62

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=1YAwMv3fHa10

+ Code + Text Reconnect

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
count	601	611	599	614	582	614	614
unique	2	2	4	2	2	3	2
top	Male	Yes	0	Graduate	No	Semiurban	Y
freq	489	398	345	480	500	233	422

#missing values
data.isna().sum()

Gender	13
Married	3
Dependents	15
Education	0
Self_Employed	32
ApplicantIncome	0
CoapplicantIncome	0
LoanAmount	22
Loan_Amount_Term	14
Credit_History	50
Property_Area	0
Loan_Status	0
	dtype: int64

62 / 63

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy + colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=1YAwMv3fHa10 Reconnect + ⚙️ Update

+ Code + Text

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
count	601	611	599	614	582	614	614
unique	2	2	4	2	2	3	2
top	Male	Yes	0	Graduate	No	Semiurban	Y
freq	489	398	345	480	500	233	422

#missing values
data.isna().sum()

Gender	13
Married	3
Dependents	15
Education	0
Self_Employed	32
ApplicantIncome	0
CoapplicantIncome	0
LoanAmount	22
Loan_Amount_Term	14
Credit_History	50
Property_Area	0
Loan_Status	0
dtype: int64	

Annotations:

- A brace groups the columns "ApplicantIncome" and "CoapplicantIncome".
- The value "50" in the "Credit_History" row is circled.
- An arrow points from the circled "50" to the "Property_Area" column.
- An arrow points from the "Property_Area" column to the "Loan_Status" column.

+ Code + Text

Reconnect ▾



V

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
count	601	611	599	614	582	614	614
unique	2	2	4	2	2	3	2
top	Male	Yes	0	Graduate	No	Semiurban	Y
freq	489	398	345	480	500	233	422

```
#missing values  
data.isna().sum()
```

The diagram illustrates a sequence of functions f_1, f_2, \dots, f_n . Each function is depicted as a vertical orange line segment within a large square frame. The segments are positioned such that they appear to converge towards the right edge of the frame, suggesting a limit process.

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy + colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=1YAwMv3fHa10 Reconnect + ⚙️ Update

+ Code + Text

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
{x}	count	601	611	599	614	582	614
	unique	2	2	4	2	2	3
	top	Male	Yes	0	Graduate	No	Semiurban
	freq	489	398	345	480	500	233

#missing values
data.isna().sum()

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
{x}	count	601	611	599	614	582	614
	unique	2	2	4	2	2	3
	top	Male	Yes	0	Graduate	No	Semiurban
	freq	489	398	345	480	500	233

Gender
Married
Dependents
Education
Self_Employed
ApplicantIncome
CoapplicantIncome
LoanAmount
Loan_Amount_Term
Credit_History
Property_Area
Loan_Status
dtype: int64

0
2
3+
5+



65 / 66

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=1YAwMv3fHa10

+ Code + Text Reconnect

unique 2 2 4 2 2 3 2

	top	Male	Yes	0	Graduate	No	Semiurban	Y
{x}	freq	489	398	345	480	500	233	422

#missing values
data.isna().sum()

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status	dtype: int64
{x}	13	3	15	0	32	0	0	22	14	50	0	0	614

50 ~ 8%

[] # catgeorical and numerical columns
cat_cols = data.dtypes == 'object'
cat_cols = list(cat_cols[cat_cols].index)

66 / 67

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=1YAwMv3fHa10

+ Code + Text Reconnect

dtype: int64

[] # catgeorical and numerical columns

{ } ✓ cat_cols = (data.dtypes == 'object')

cat_cols = list(cat_cols[cat_cols].index)

{ } num_cols = data.dtypes != 'object'

num_cols = list(num_cols[num_cols].index)

cat_cols.remove('Loan_Status')

[] cat_cols

✓ { ['Gender', 'Married', 'Dependents', 'Education', 'Self_Employed', 'Property_Area'] }

<> [] num_cols

[] ['ApplicantIncome', 'CoapplicantIncome', 'LoanAmount', 'Loan Amount Term',

Reconnect

Update

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=1YAwMv3tHa10

+ Code + Text

econnect ▾

dtype: int64

```
[ ] # categorical and numerical columns
cat_cols = data.dtypes == 'object'
cat_cols = list(cat_cols[cat_cols].index)

num_cols = data.dtypes != 'object'
num_cols = list(num_cols[num_cols].index)
cat_cols.remove('Loan Status')
```

[] cat cols

```
} ['Gender',  
 'Married',  
 'Dependents',  
 'Education',  
 'Self_Employed',  
 'Property_Area']
```

[] num_cols

```
{ 'ApplicantIncome',  
  'CoapplicantIncome',  
  'LoanAmount',  
  'Loan Amount Term',
```

+ Code + Text

Reconnect



Basic Data Exploration

[] ↗ 9 cells hidden

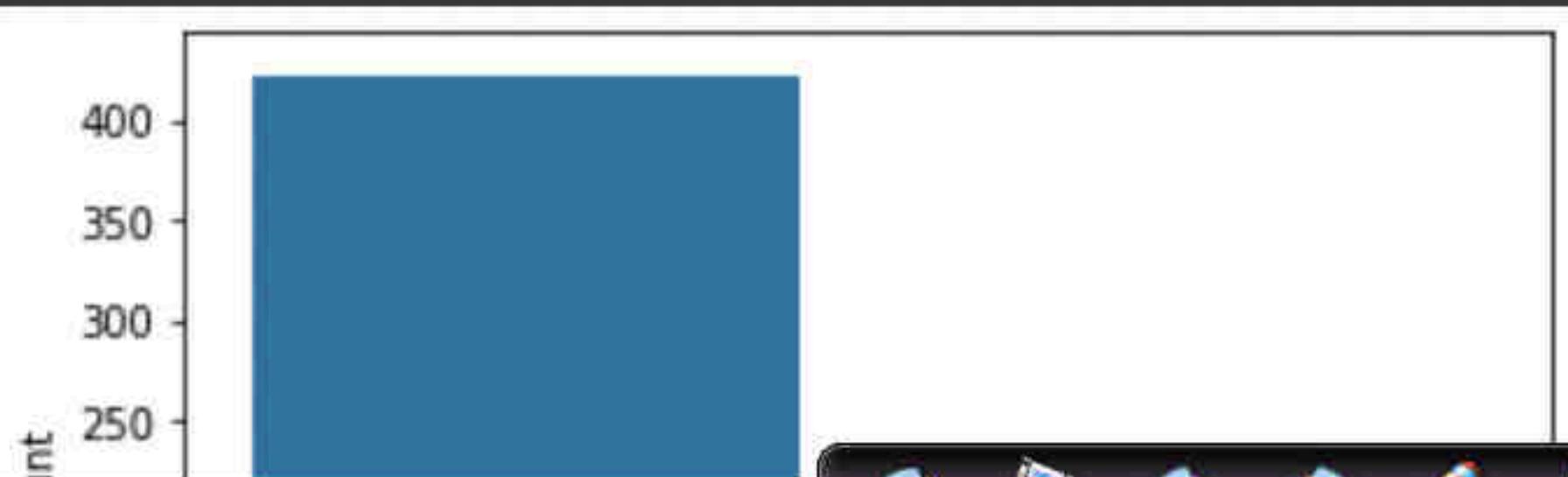
{x}

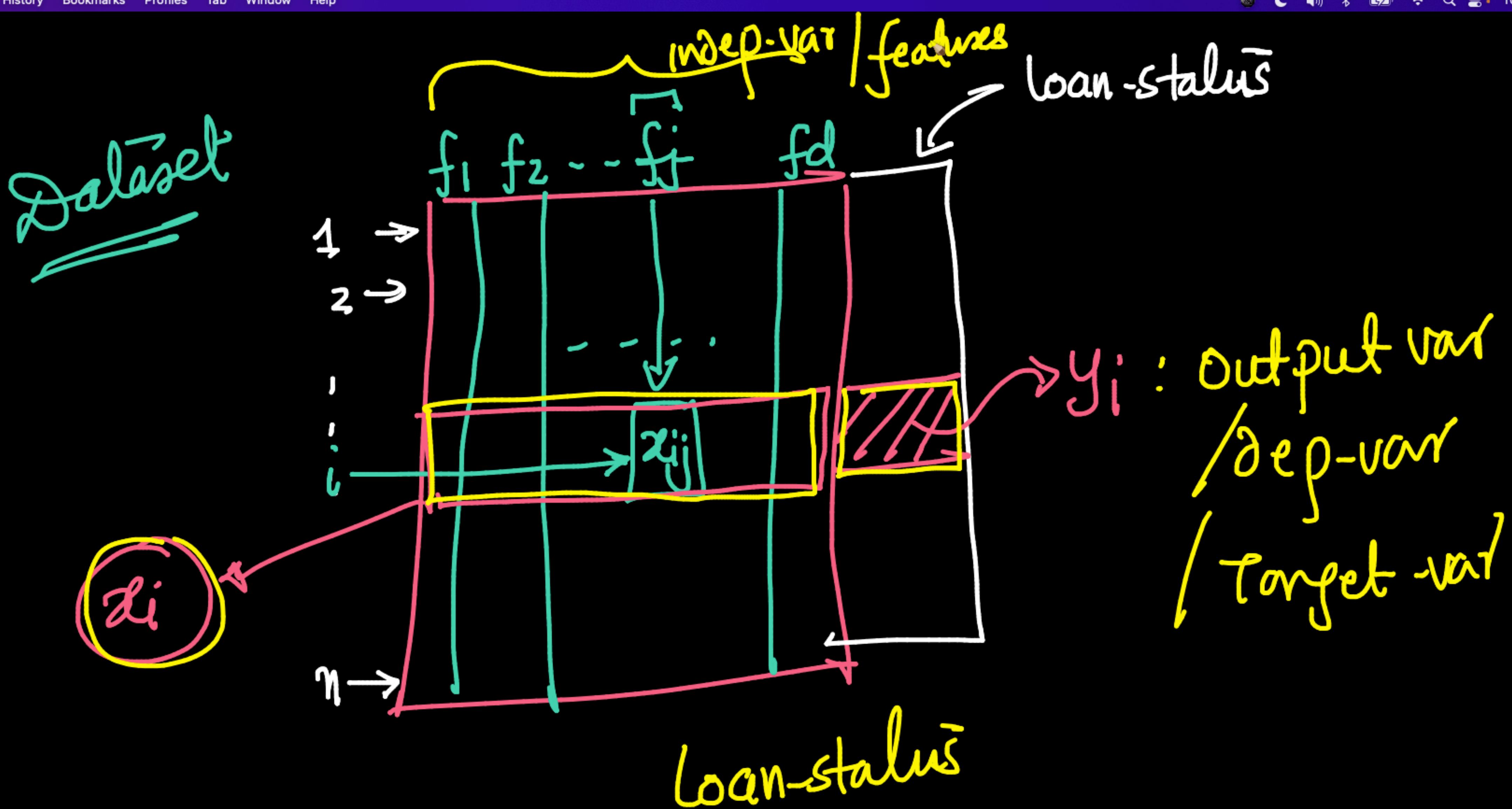
Basic Data visualization: Univariate

imbalanced dataset

```
data['Loan_Status'].value_counts()  
Y    422  
N    192  
Name: Loan_Status, dtype: int64
```

```
#Q: How many loans the company has approved in the past?  
sns.countplot(data=data, x='Loan_Status')  
plt.show()
```





binary
y.v

$y_i \in \{\text{Yes}, \text{No}\}$

↓ ↓

10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=GNmtX3um9oJR

+ Code + Text Reconnect

Basic Data Exploration

[] ↳ 9 cells hidden

{x}

Basic Data visualization: Univariate

data['Loan_Status'].value_counts()

Y 422
N 192

Name: Loan_Status, dtype: int64

imbalanced dataset

[] #Q: How many loans the company has approved in the past?
sns.countplot(data=data, x='Loan_Status')
plt.show()

72 / 72

+ Code + Text

Basic Data Exploration

[] ↗ 9 cells hidden

Basic Data visualization: Univariate

[] `data['Loan_Status'].value_counts()`

```
{ Y    422
  N    192
Name: Loan_Status, dtype: int64
```

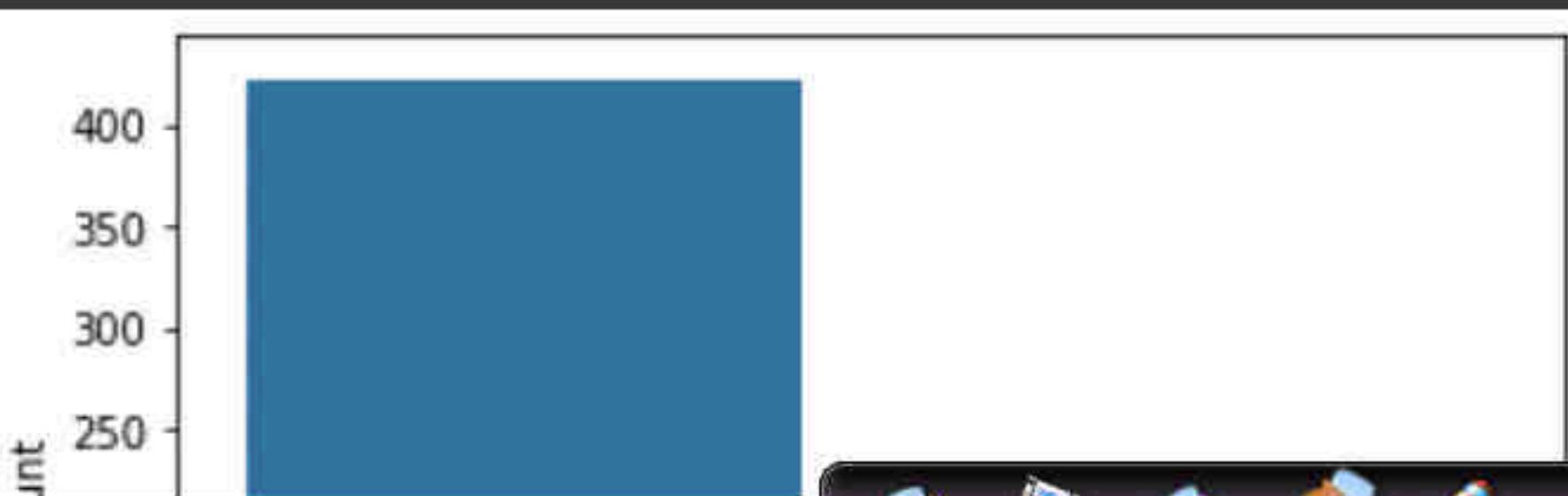
1: (000)

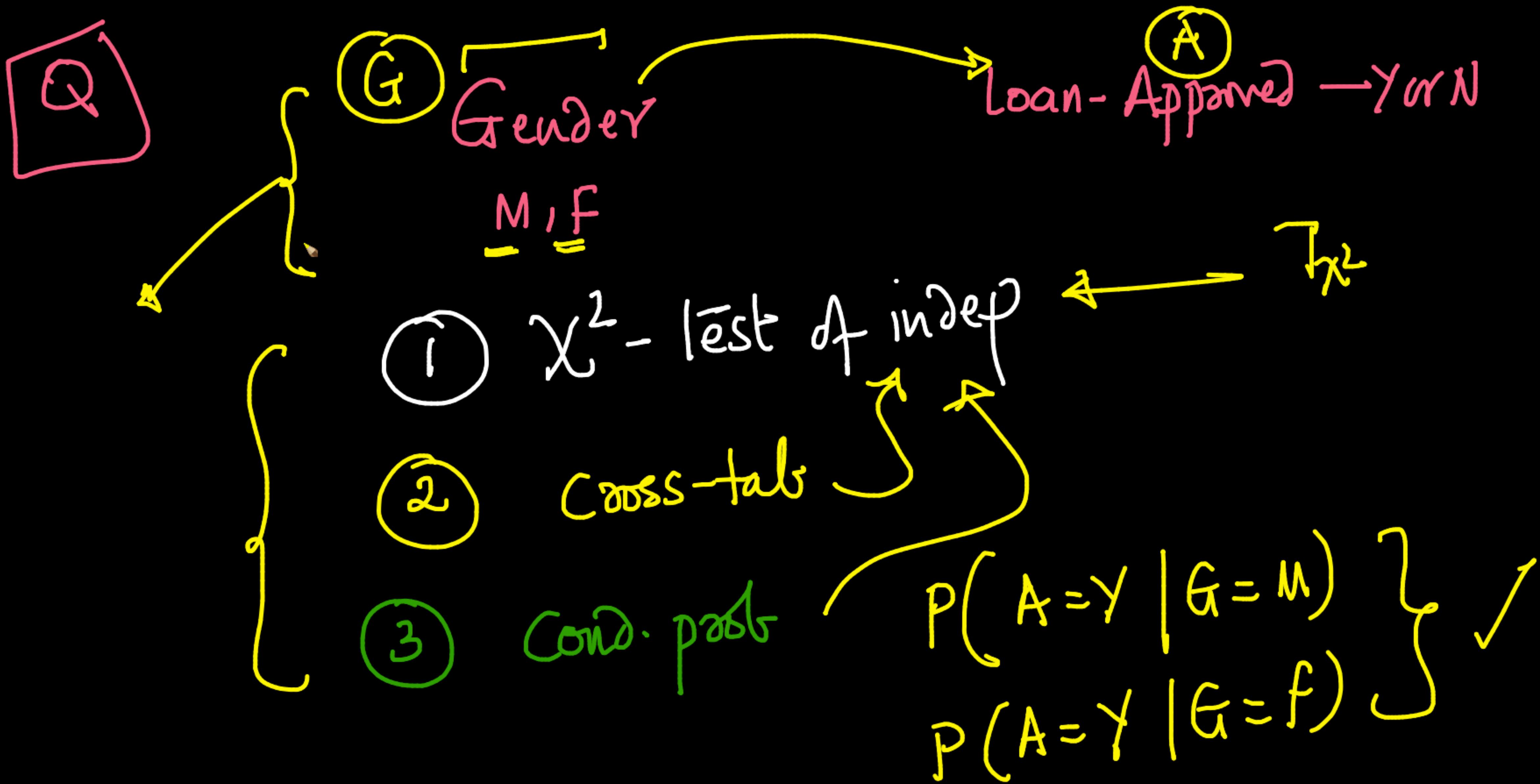
Questions

Techniques

Answers

```
[ ] #Q: How many loans the company has approved in the past?
sns.countplot(data=data, x='Loan_Status')
plt.show()
```





+ Code + Text

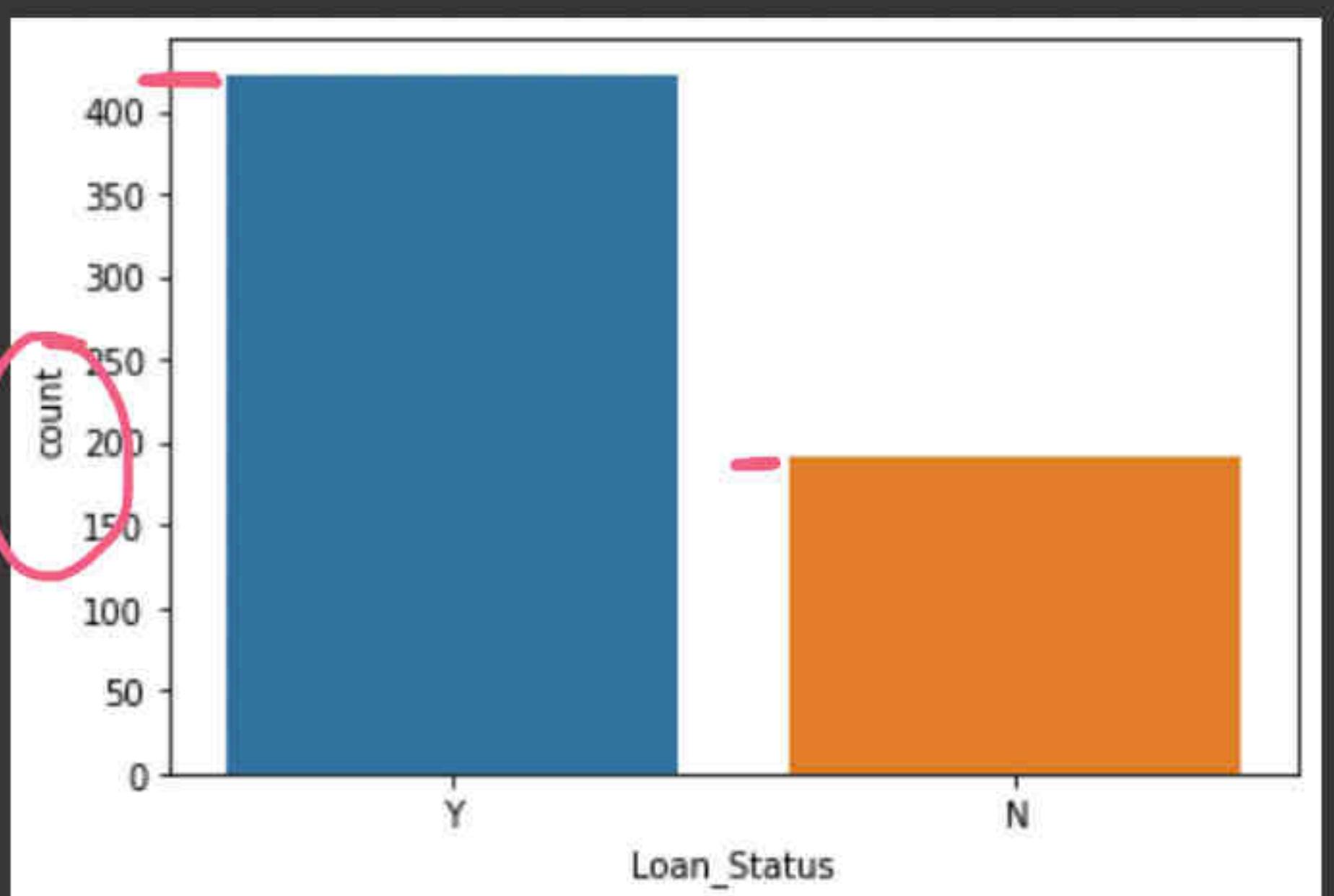
Reconnect



Y 422
N 192
Name: Loan_Status, dtype: int64



#Q: How many loans the company has approved in the past?
sns.countplot(data=data, x='Loan_Status')
plt.show()



[] target = 'Loan_Status'
data[target].value_counts()



10.2 - Hypothesis Te Kruskal-Wallis one-w Anscombe's quartet Bias of an estimator EDA_FE.ipynb - Colab Markov's inequality Chebyshev's inequality scipy.stats.kruskal numpy.var — NumPy

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=ufCVt1-DI051

+ Code + Text Reconnect

Imbalanced data

{x}

Y 422
N 192
Name: Loan_Status, dtype: int64

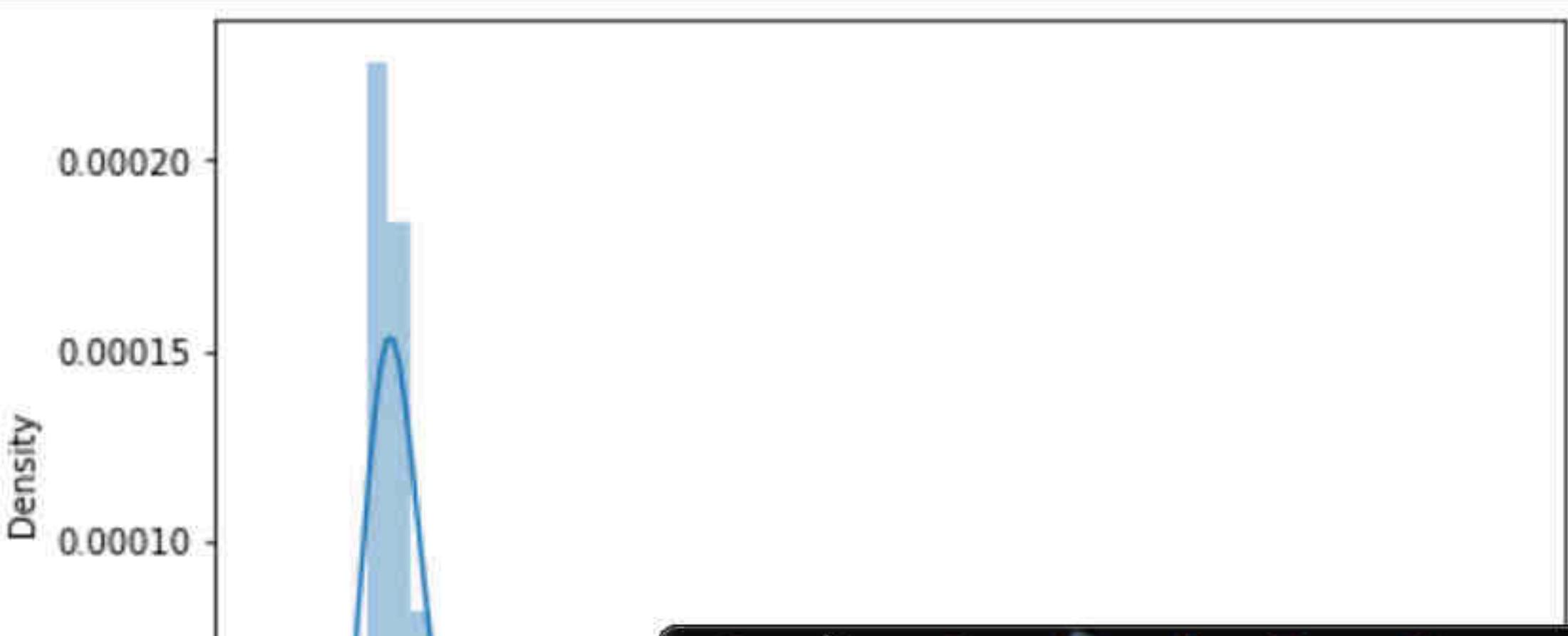
📁

Income of the applicant

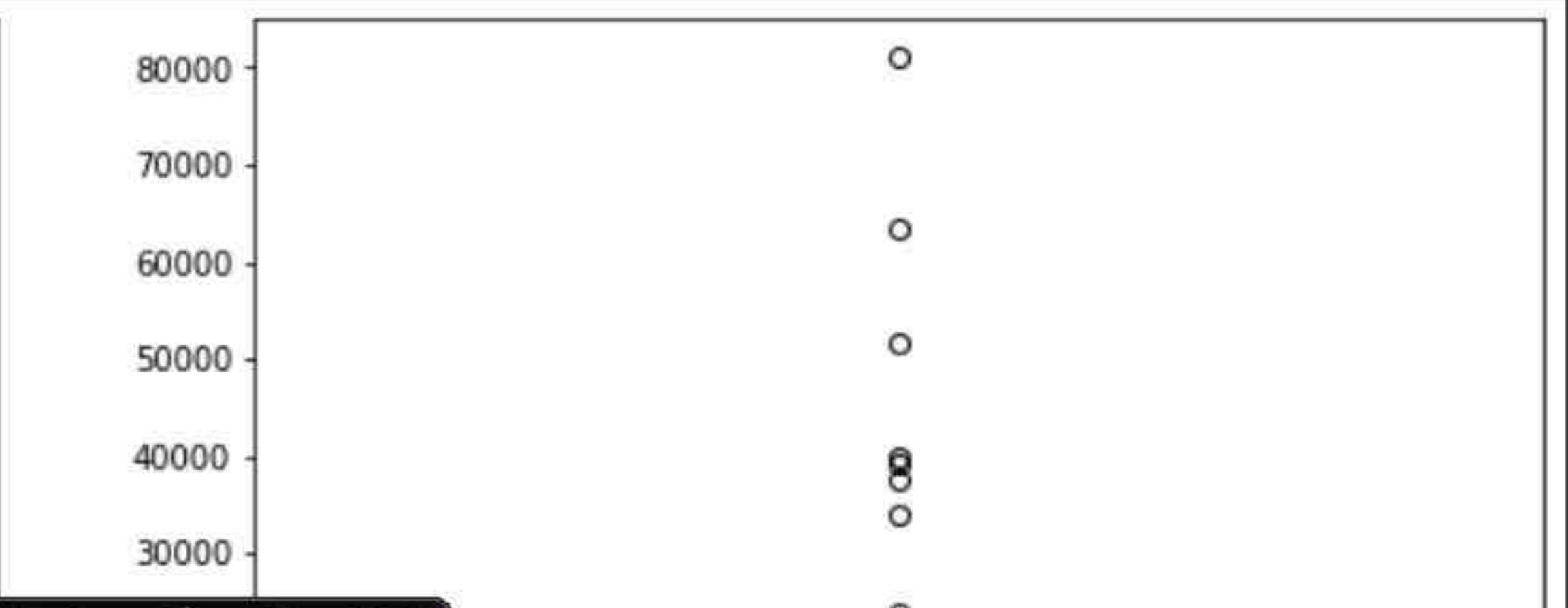
plt.subplot(121)
sns.distplot(data["ApplicantIncome"])

plt.subplot(122)
data["ApplicantIncome"].plot.box(figsize=(16,5))
plt.show()

Density



80000
70000
60000
50000
40000
30000

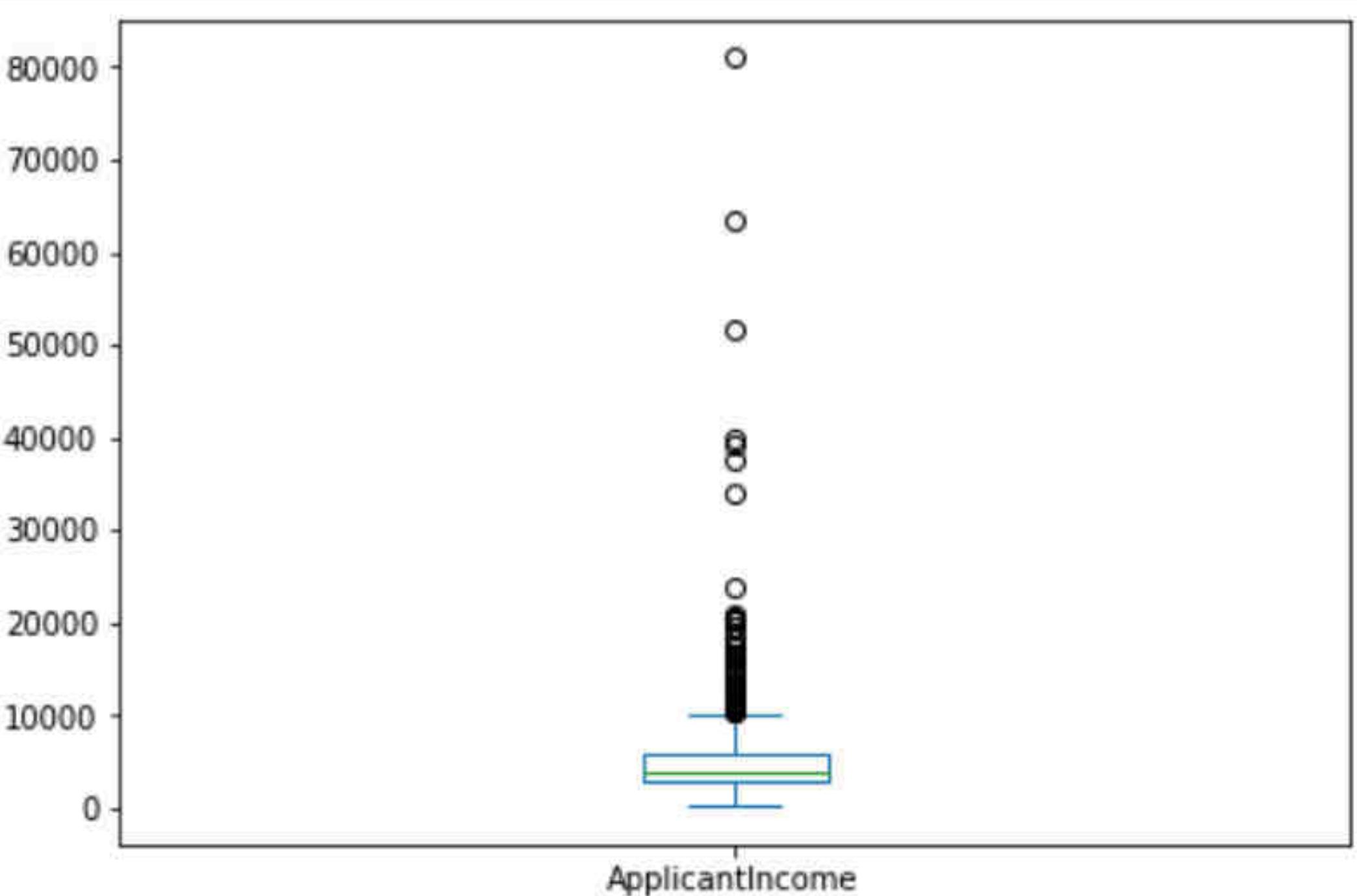
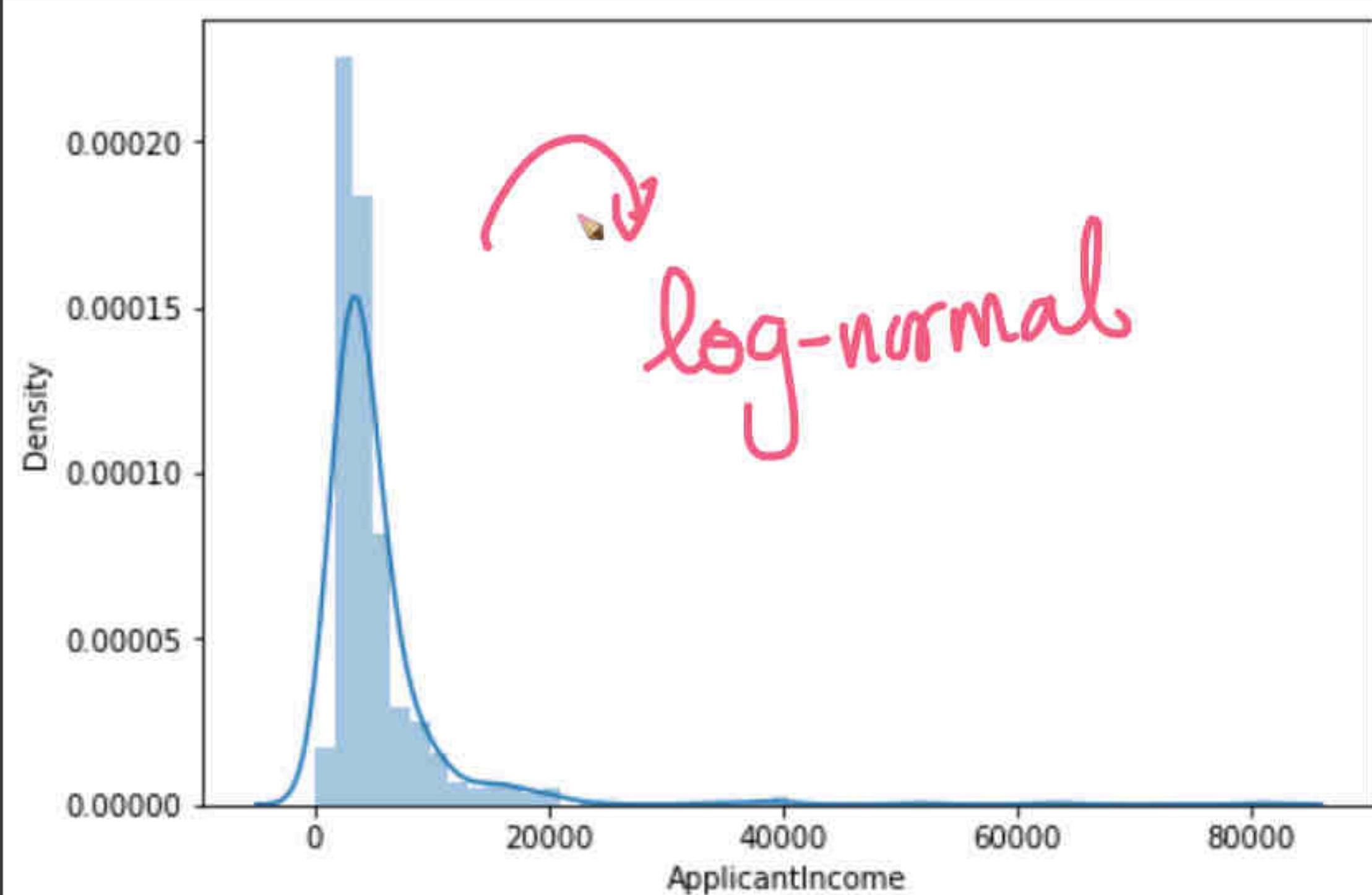


+ Code + Text

Reconnect



```
plt.subplot(122)
data[ "ApplicantIncome" ].plot.box(figsize=(16,5))
plt.show()
```



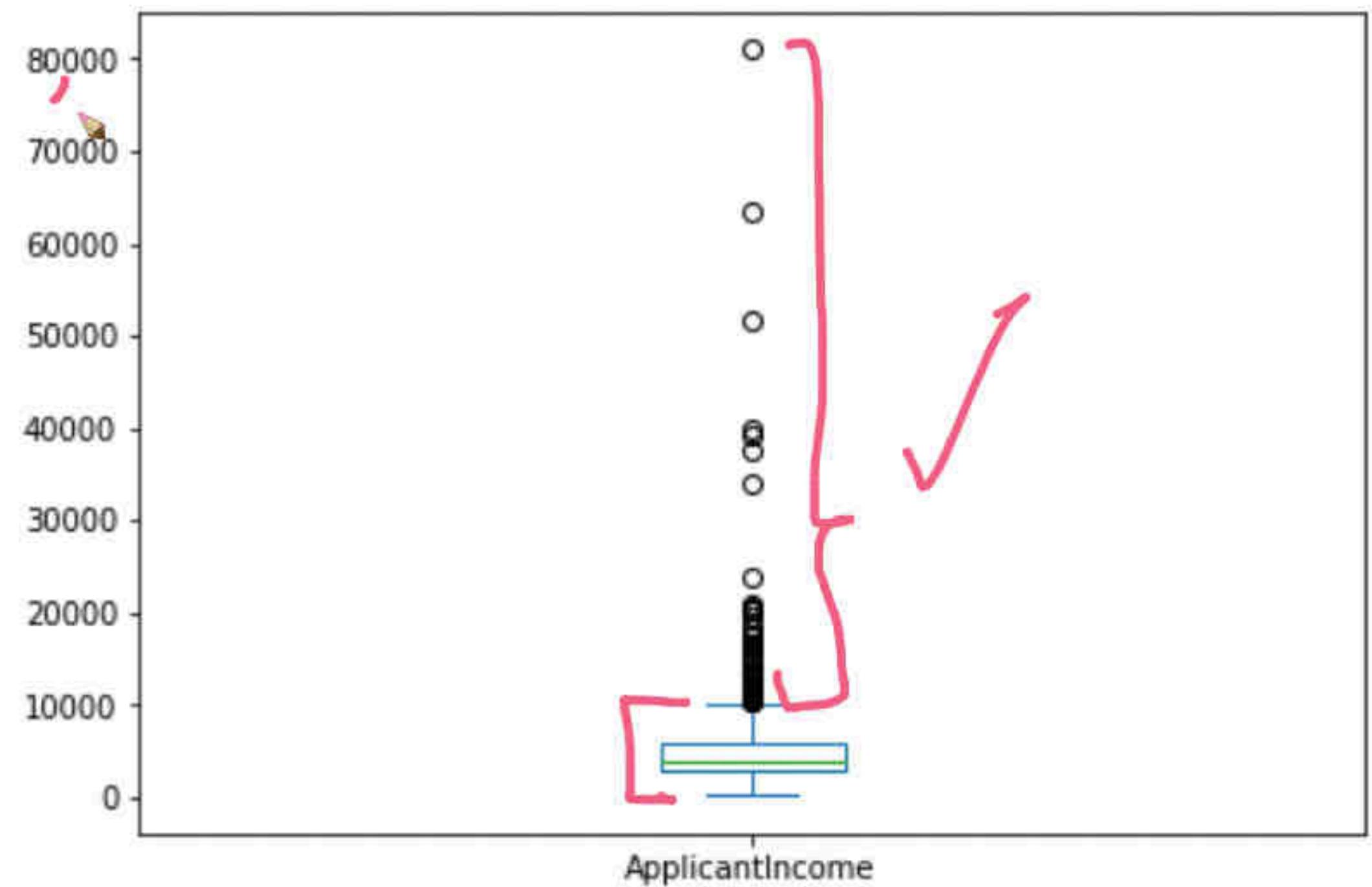
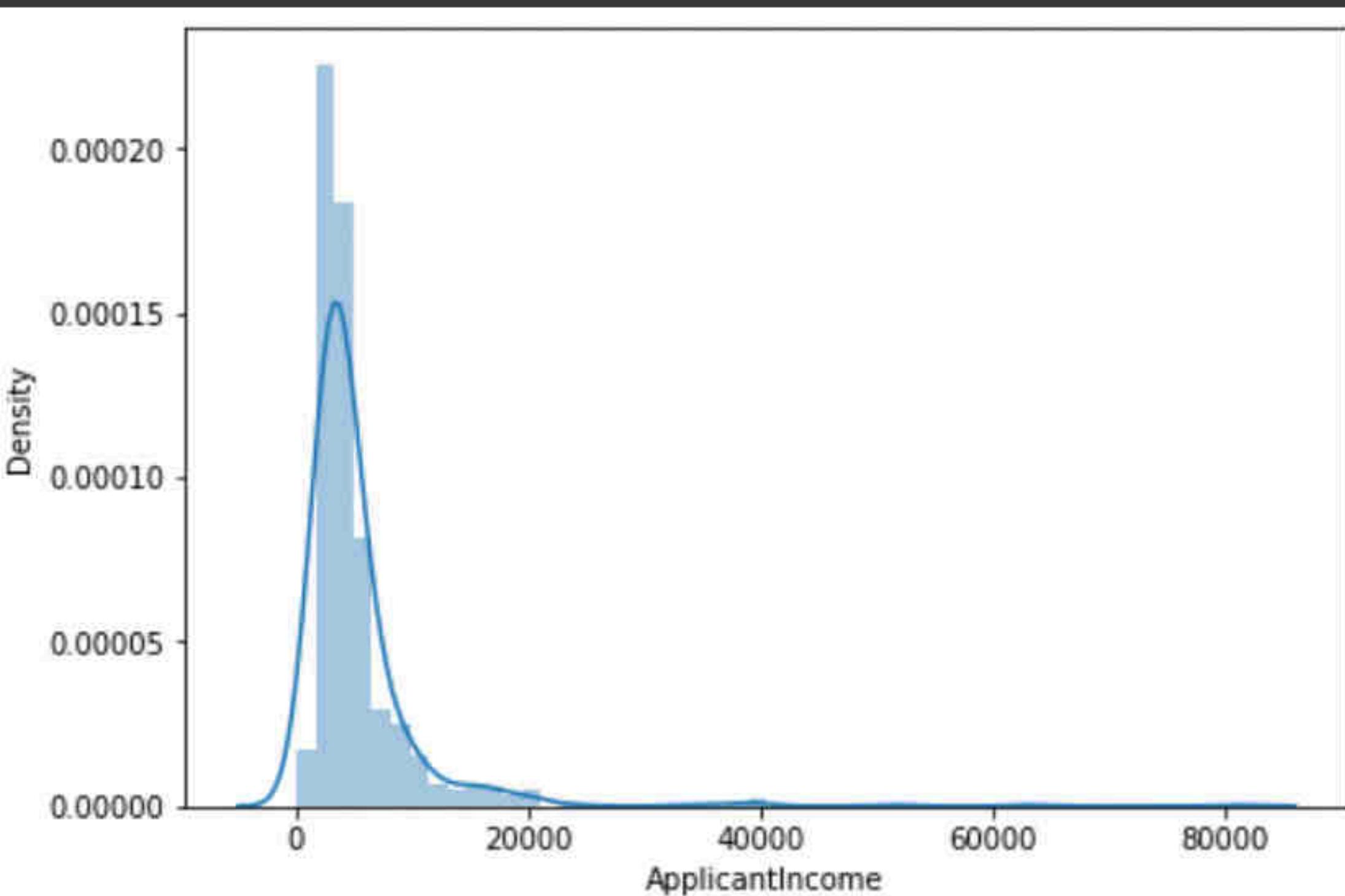
```
[ ] plt.subplot(121)
sns.distplot(np.log(data[ "ApplicantIncome" ]))
```

+ Code + Text

Reconnect



```
plt.subplot(122)
data["ApplicantIncome"].plot.box(figsize=(16,5))
plt.show()
```



```
[ ] plt.subplot(121)
sns.distplot(np.log(data["ApplicantIncome"]))
```

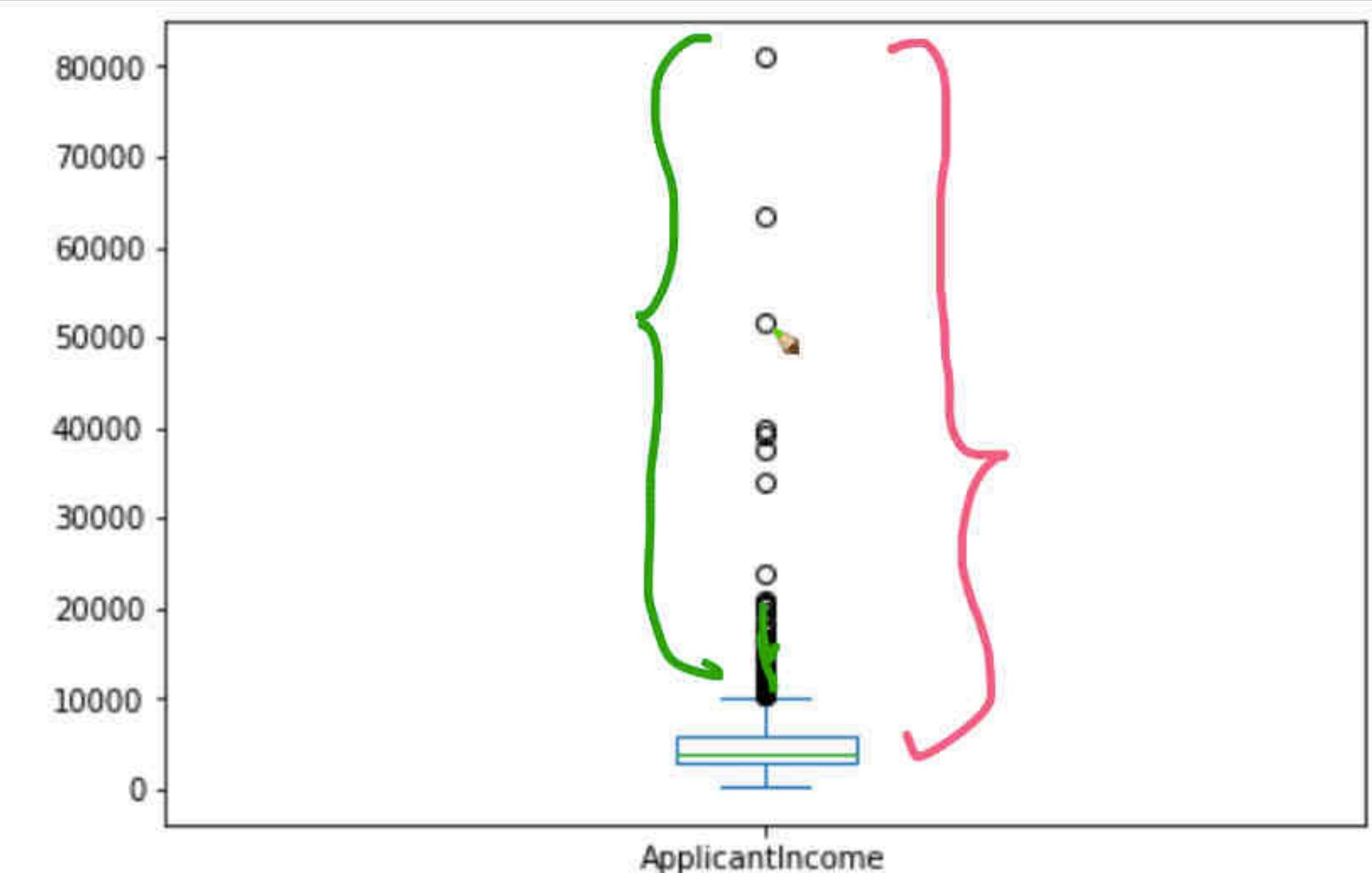
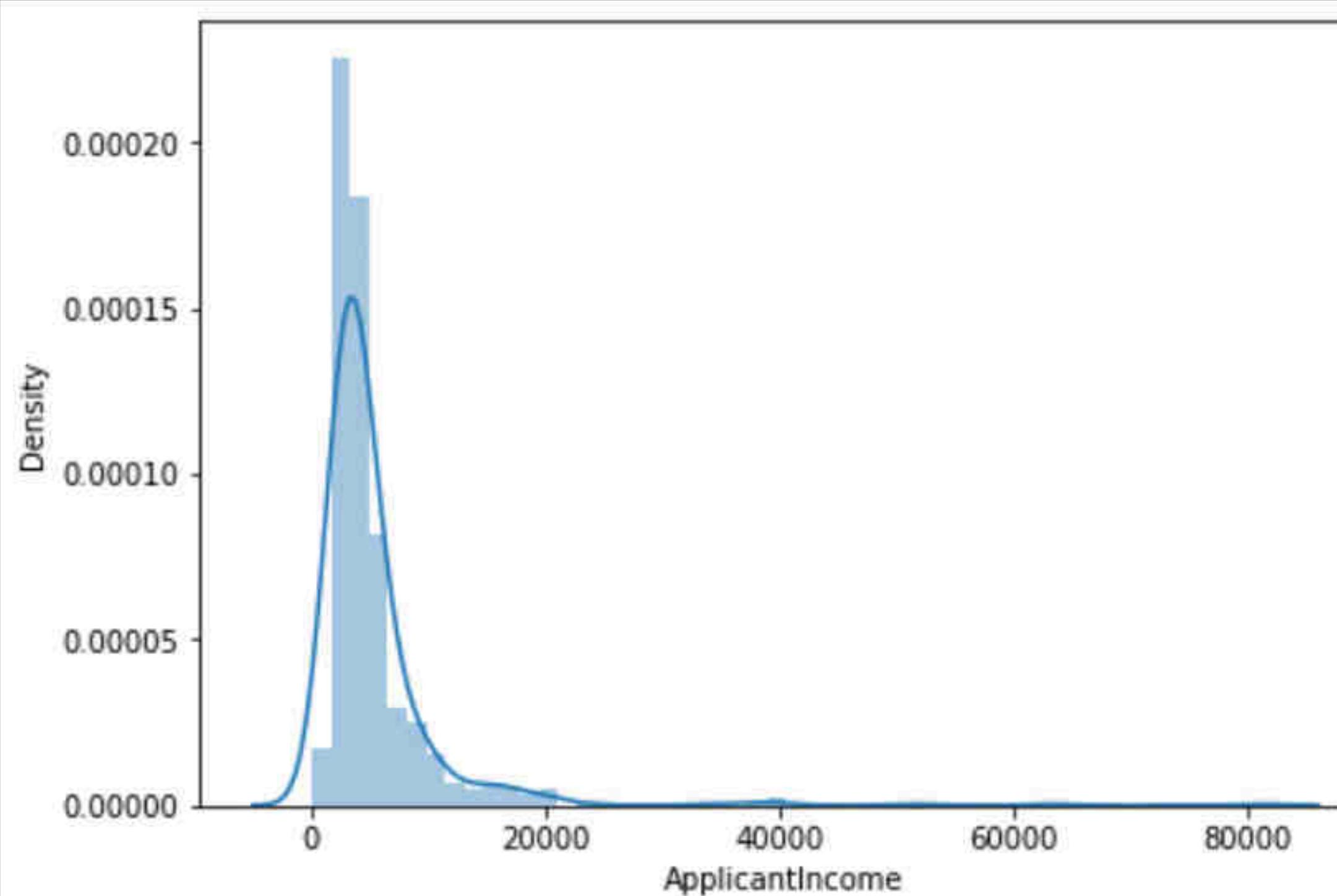
+ Code + Text

Reconnect



```
[ ] #Income of the applicant  
plt.subplot(121)  
sns.distplot(data["ApplicantIncome"] )  
  
plt.subplot(122)  
data["ApplicantIncome"].plot.box(figsize=(16,5))  
plt.show()
```

QR

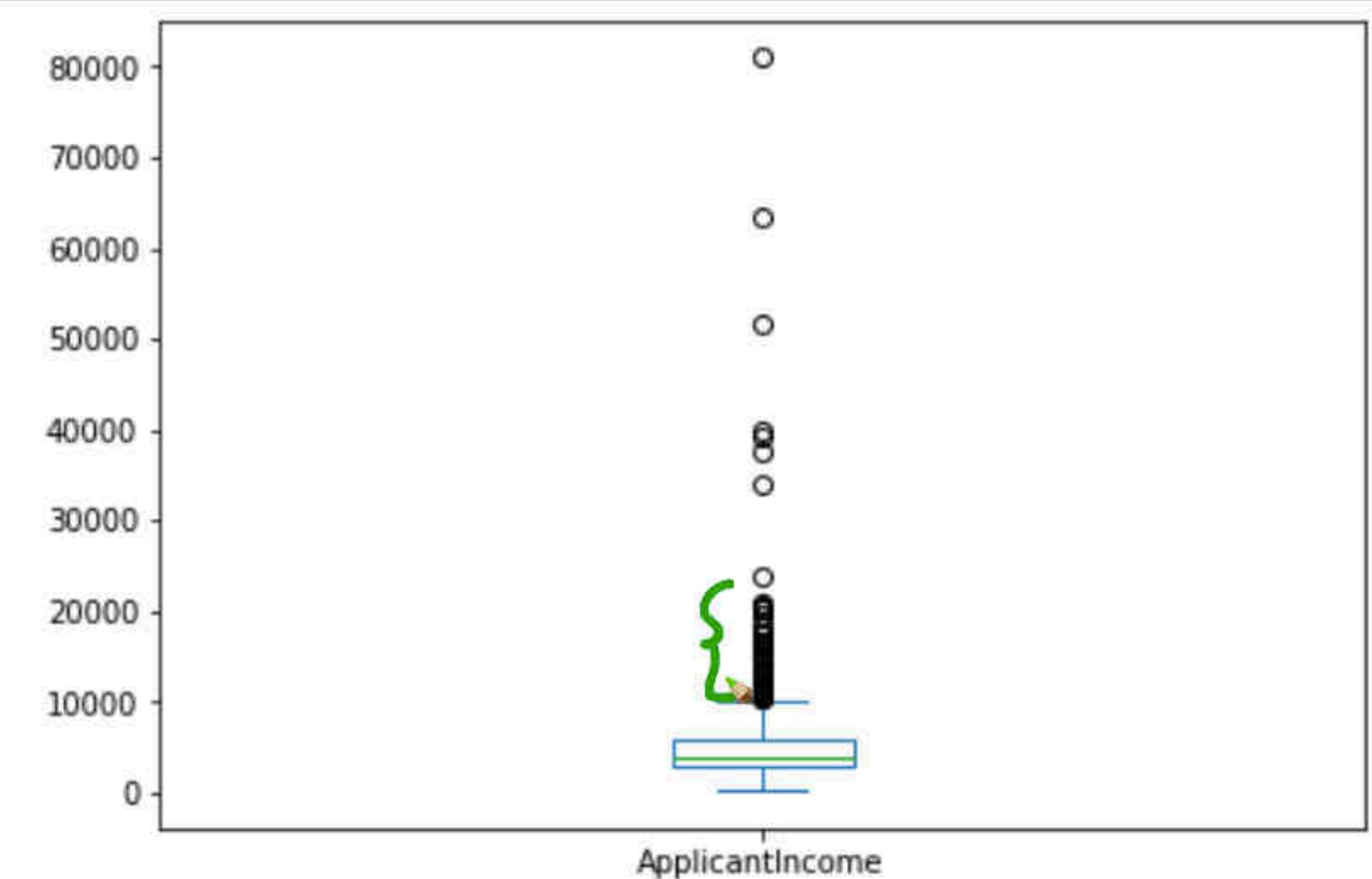
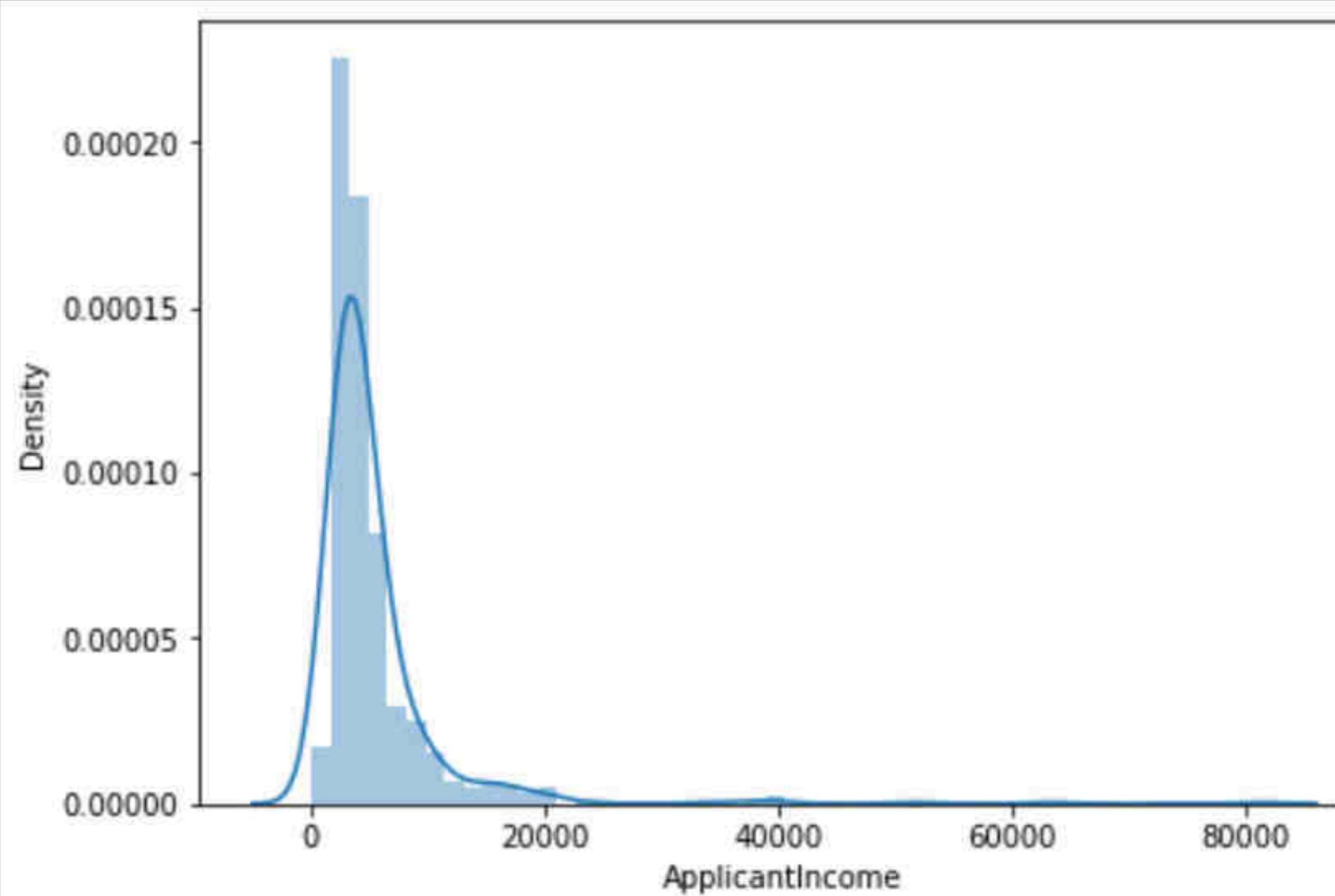


+ Code + Text

Reconnect

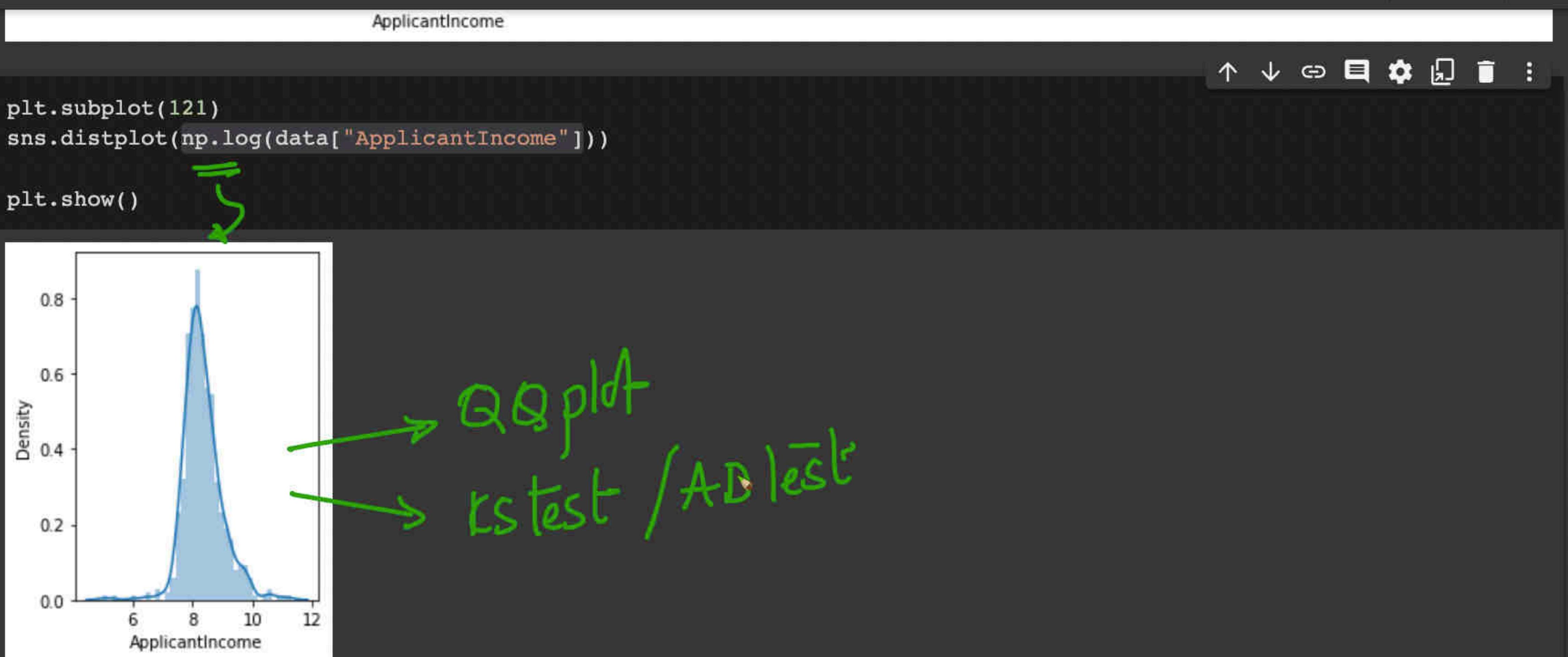


```
[ ] #Income of the applicant  
plt.subplot(121)  
sns.distplot(data["ApplicantIncome"] )  
  
{x}  
plt.subplot(122)  
data["ApplicantIncome"].plot.box(figsize=(16,5))  
plt.show()
```



+ Code + Text

econnect ▾



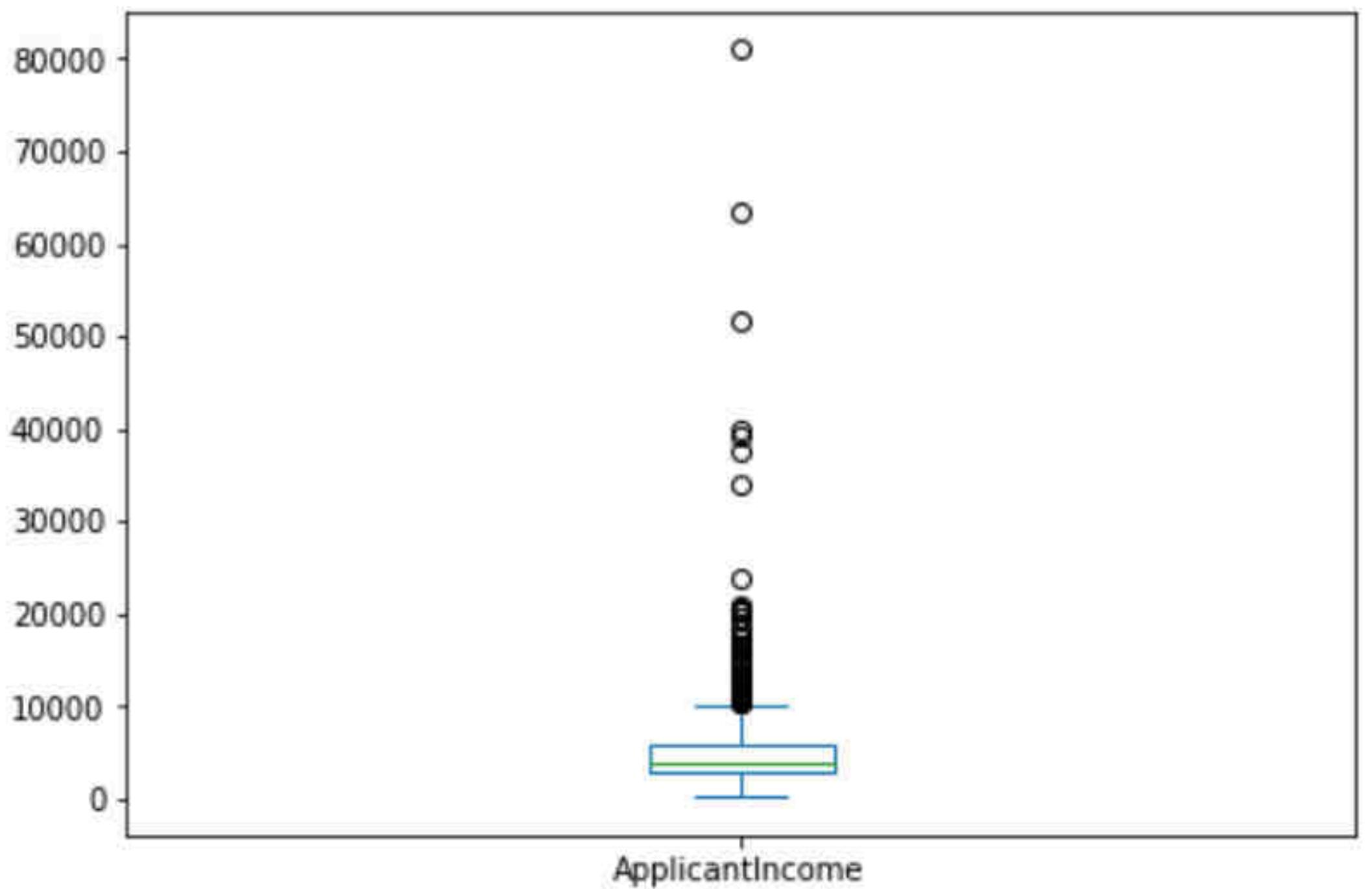
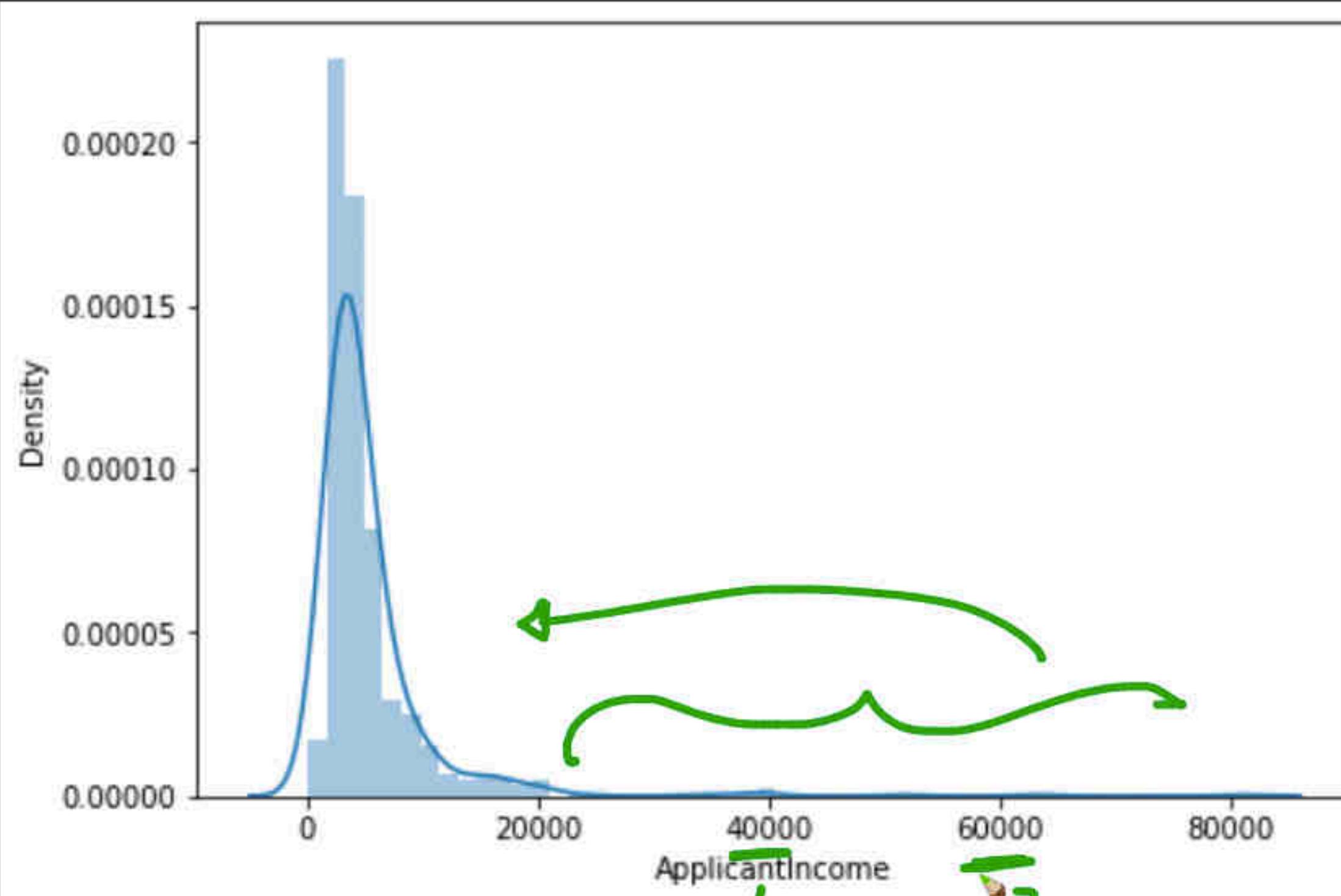
```
[ ] #Slice this data by Education
```

+ Code + Text

Reconnect



```
data[ "ApplicantIncome" ].plot.box(figsize=(16,5))  
plt.show()
```



```
plt.subplot(121)  
sns.distplot(np.log(data[ "ApplicantIncome" ]))  
  
plt.show()
```

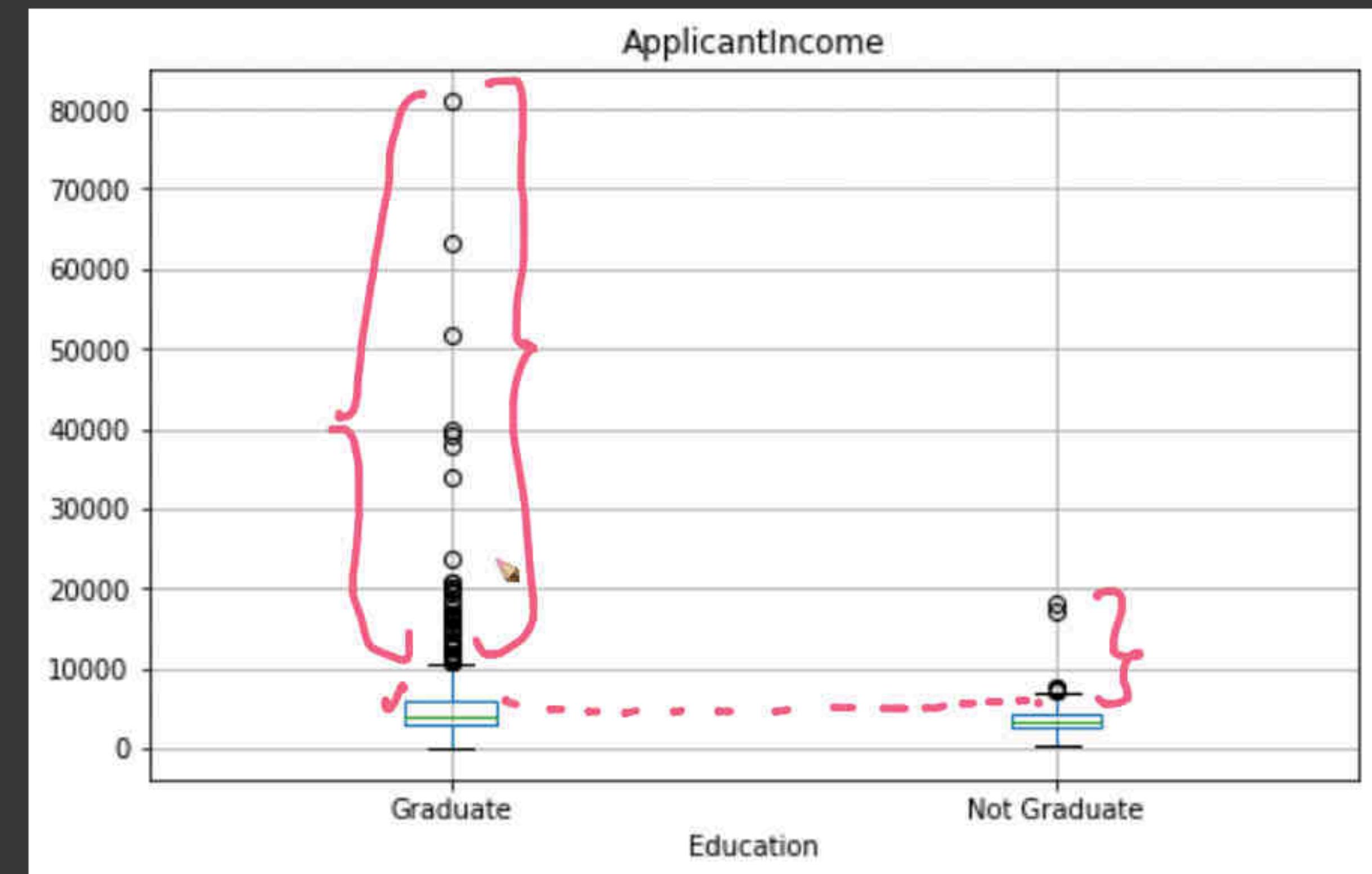


+ Code + Text

Reconnect



```
▶ data.boxplot(column='ApplicantIncome', by="Education", figsize=(8,5))
plt.suptitle("")
plt.show()
```



```
[ ] #co-applicant income
plt.subplot(121)
```

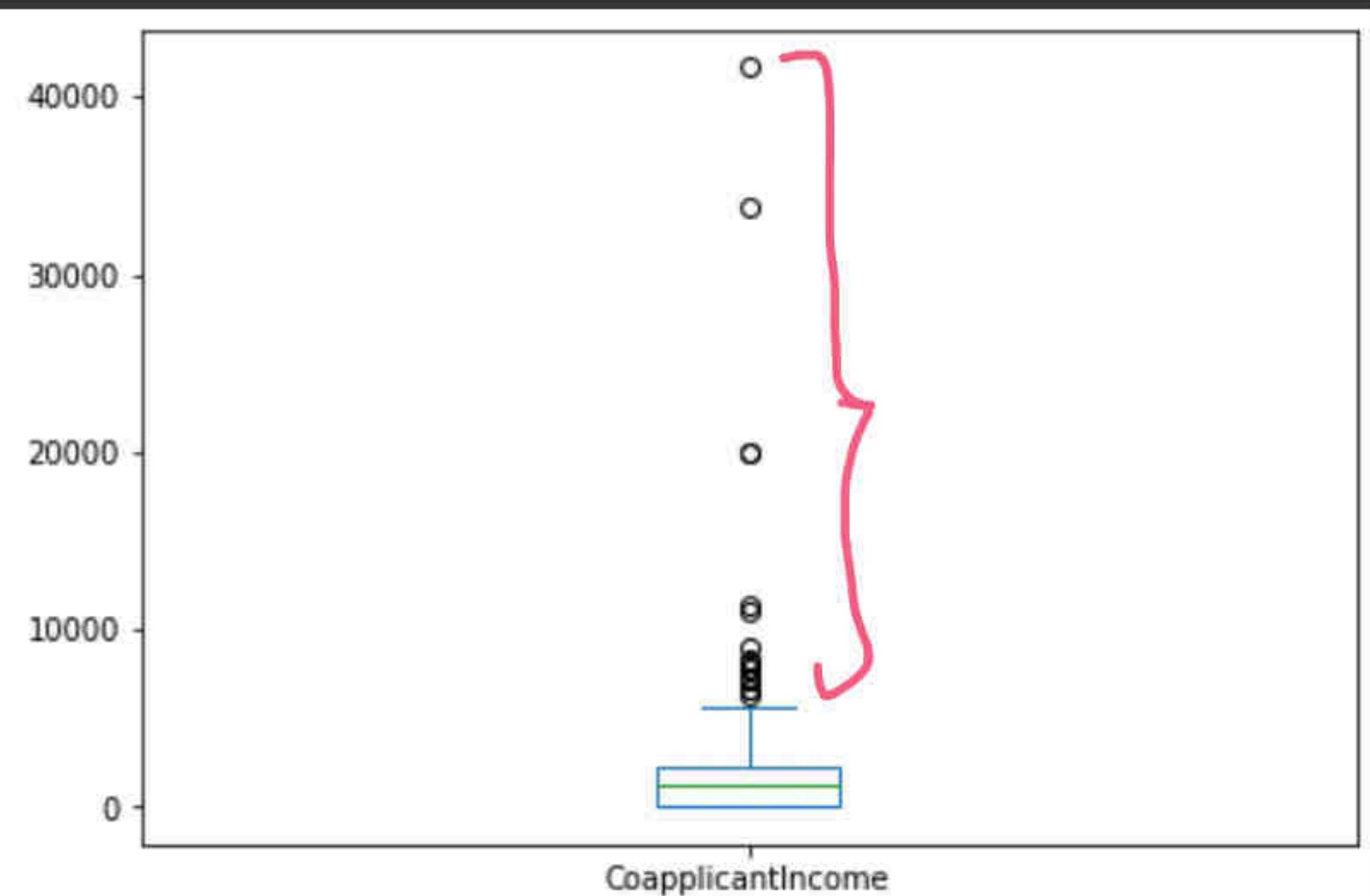
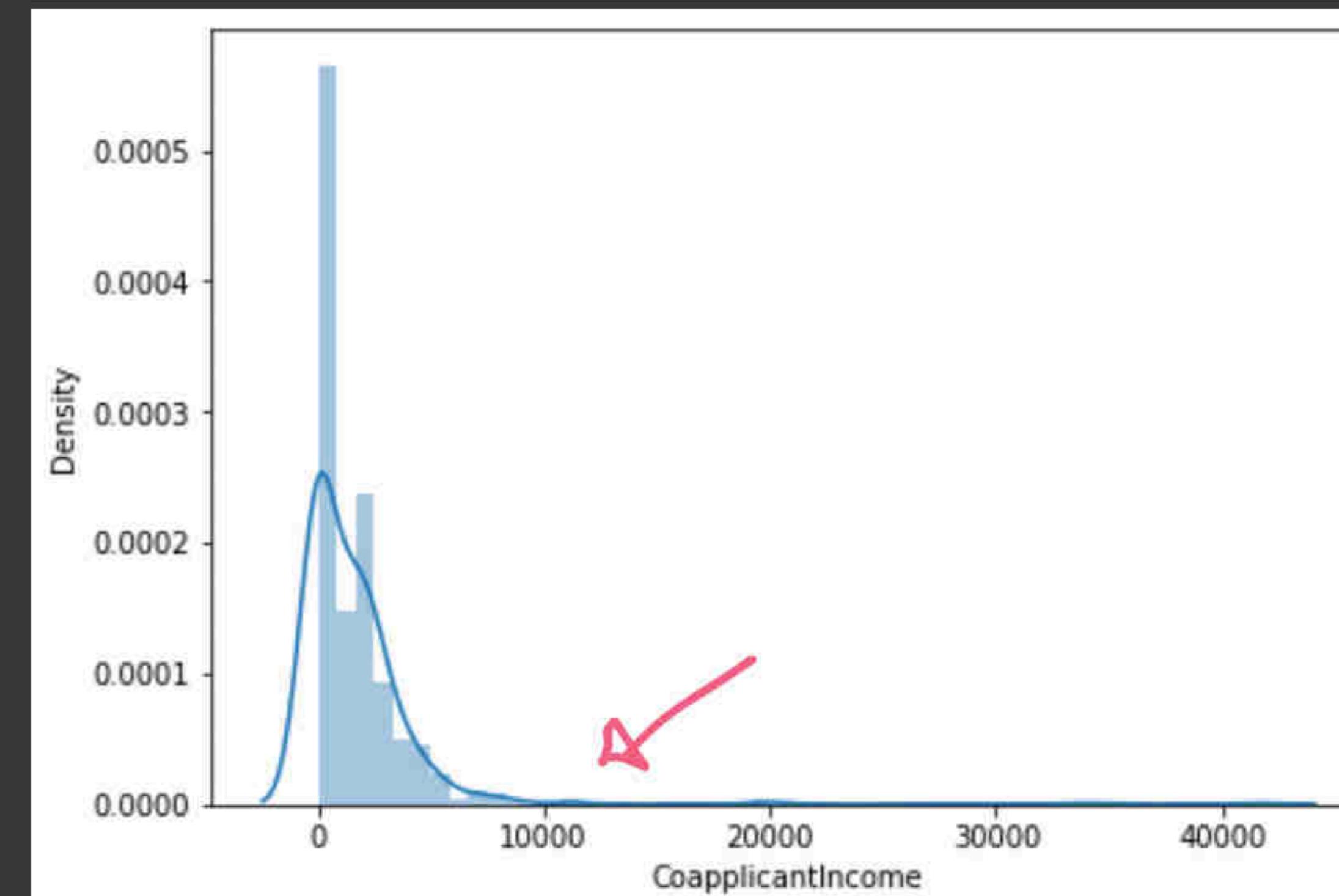


+ Code + Text

Reconnect



```
plt.subplot(122)
data["CoapplicantIncome"].plot.box(figsize=(16,5))
plt.show()
```



[] #Relation between "Loan_Status" and "Income"

10.2 - Hypothesis Testing | ST Kruskal-Wallis one-way analy... Anscombe's quartet 3 - Ansco... Bias of an estimator - Wikipedia EDA_FE.ipynb - Colaboratory scipy.stats.kruskal — SciPy v1.23.0 numpy.var — NumPy v1.23.0

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=gwludxlnKDvf

+ Code + Text Reconnect

Code Text

Search

{x}

CoapplicantIncome

0.0002
0.0001
0.0000

10000
0

CoapplicantIncome

#Relation between "Loan_Status" and "Income"

[] data.groupby("Loan_Status").mean()["ApplicantIncome"]

Loan_Status

N	5446.078125
Y	5384.068720

Name: ApplicantIncome, dtype: float64

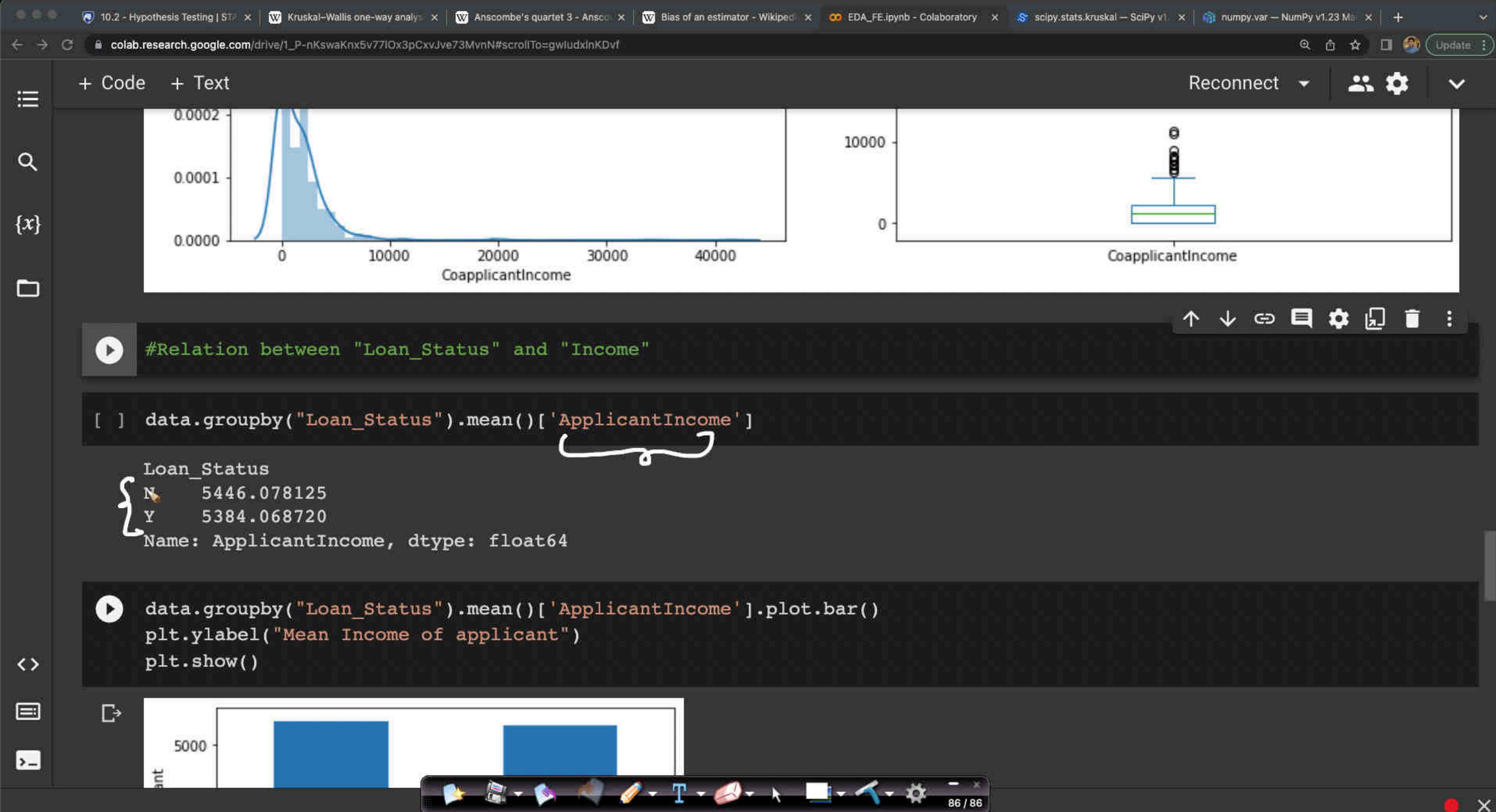
data.groupby("Loan_Status").mean()["ApplicantIncome"].plot.bar()
plt.ylabel("Mean Income of applicant")
plt.show()

5000

int

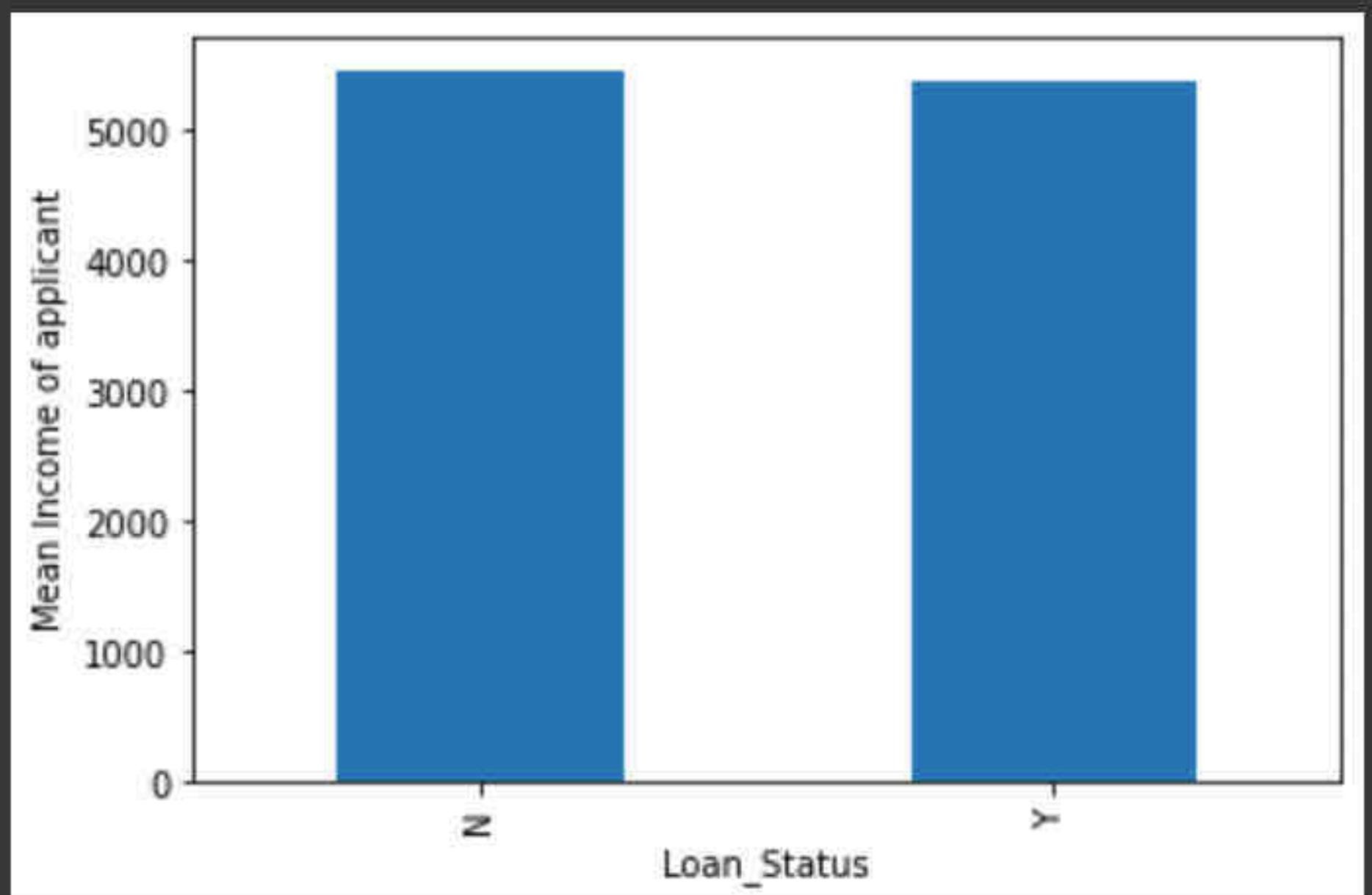
85 / 85

The screenshot shows a Jupyter Notebook interface with several plots and code cells. At the top, there are two plots: a histogram of 'CoapplicantIncome' with a blue shaded area around zero, and a box plot of 'CoapplicantIncome' with a median at approximately 10,000. Below these are two code cells. The first cell contains the code 'data.groupby("Loan_Status").mean()["ApplicantIncome"]'. The output of this cell is a table showing the mean 'ApplicantIncome' for loans with status 'N' (5446.078125) and 'Y' (5384.068720). The second cell contains code to plot the mean income by loan status using bar charts. The bottom part of the screen shows a toolbar with various icons and a status bar indicating '85 / 85'.



+ Code + Text

Reconnect ▾



Data → conclusions

- Simple ✓
- plots
- Summary ✓

Biz

► Simple Feature Engineering

[] ↳ 37 cells hidden

10.2 - Hypothesis Testing | ST x Kruskal-Wallis one-way analysis of variance by ranks x Anscombe's quartet 3 - Anscombe's quartet x Bias of an estimator - Wikipedia x EDA_FE.ipynb - Colaboratory x scipy.stats.kruskal — SciPy v1.23.0 Reference Guide x numpy.var — NumPy v1.23.0 Reference Guide x + colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=sp3SJkJmKalm

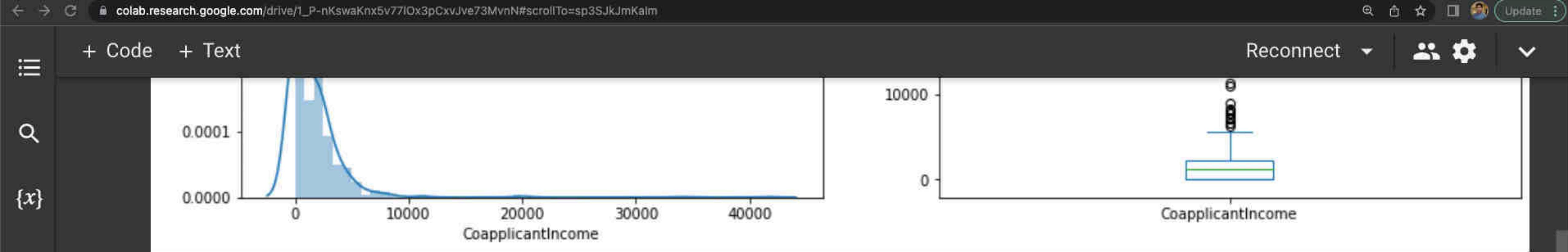
+ Code + Text Reconnect

Simple Feature Engineering

{x} [] # Feature binning: income
bins=[0,2500,4000,6000, 8000, 10000, 20000, 40000, 81000]
group=['Low', 'Average', 'medium', 'H1', 'h2', 'h3', 'h4' , 'Very high']
data['Income_bin']= pd.cut(data['ApplicantIncome'],bins,labels=group)

[] data.head()

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Disbursed
0	Male	No	0	Graduate	No	5849	0.0	Nan	141.0
1	Male	Yes	1	Graduate	No	4583	1508.0	128.0	120.0
2	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	30.0
3	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	30.0
4	Male	No	0	Graduate	No	6000	0.0	141.0	30.0

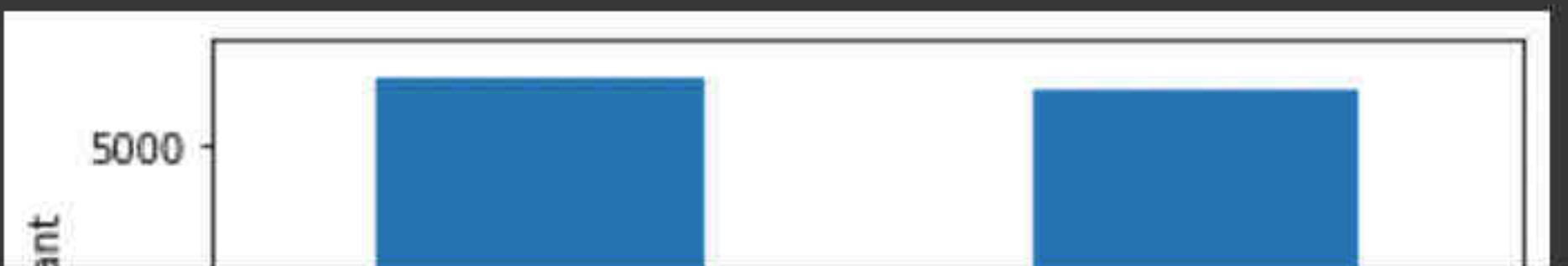


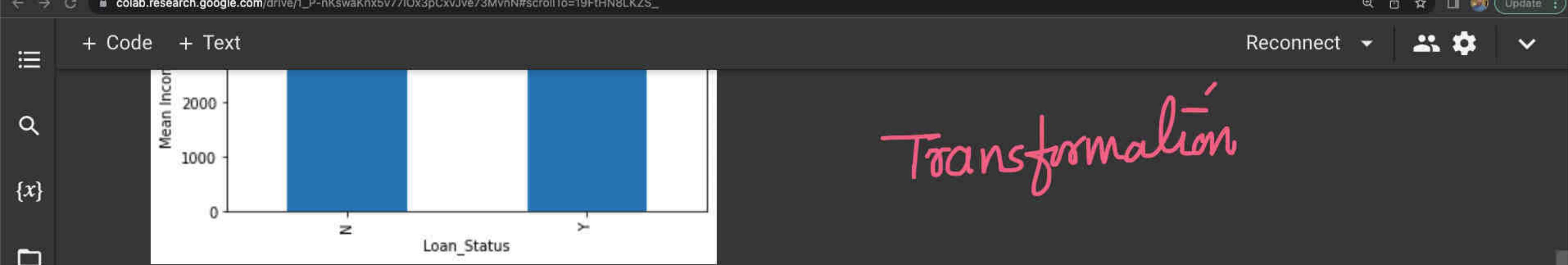
```
[ ] #Relation between "Loan_Status" and "Income"
```

```
[ ] data.groupby("Loan_Status").mean()['ApplicantIncome']
```

```
Loan_Status
N    5446.078125
Y    5384.068720
Name: ApplicantIncome, dtype: float64
```

```
[ ] data.groupby("Loan_Status").mean()['ApplicantIncome'].plot.bar()
plt.ylabel("Mean Income of applicant")
plt.show()
```



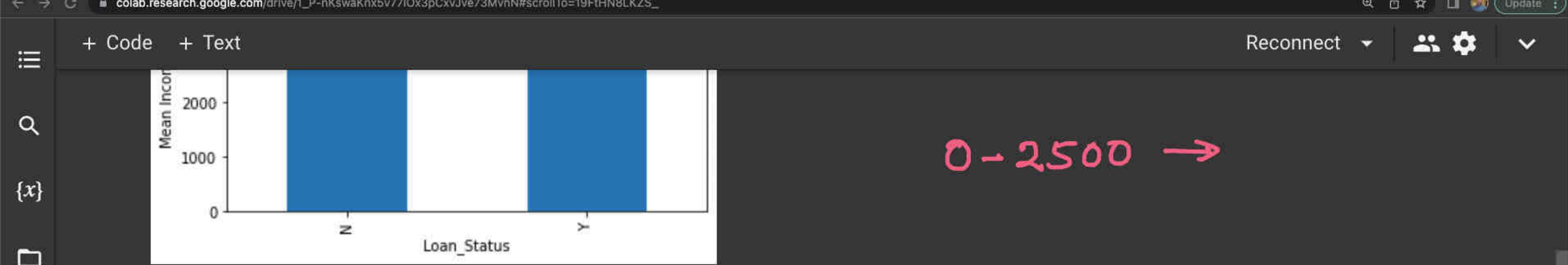


Simple Feature Engineering

```
# Feature binning: income
bins=[0,2500,4000,6000, 8000, 10000, 20000, 40000, 81000]
group=['Low', 'Average', 'medium', 'H1', 'h2', 'h3', 'h4' , 'Very high']
data[ 'Income_bin']= pd.cut(data[ 'ApplicantIncome'],bins,labels=group)
```

```
[ ] data.head()
```

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_
0	Male	No	0	Graduate	No	5849	0.0	NaN	3
1	Male	Yes	1	Graduate	No	4583	1508.0	128.0	3



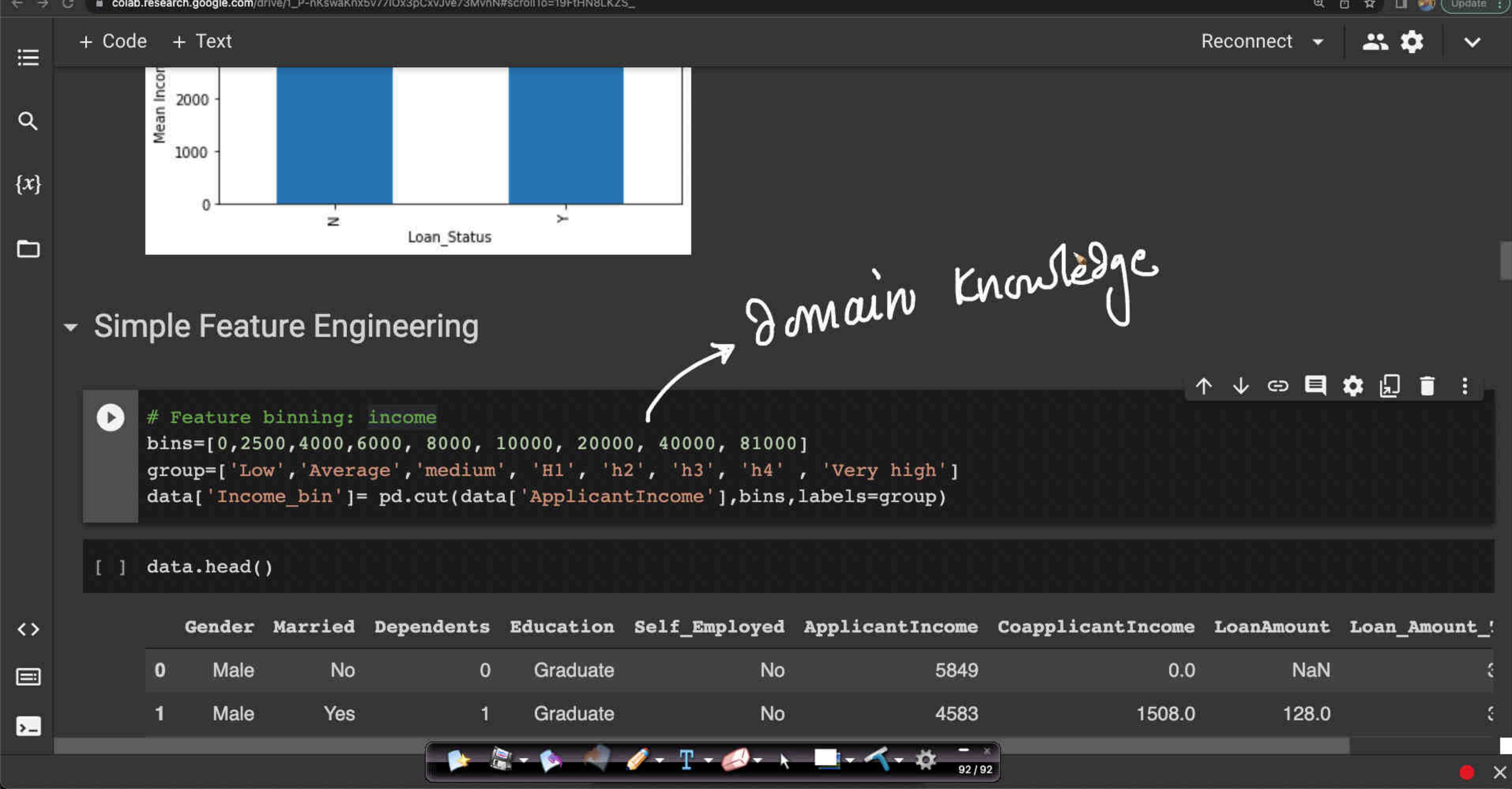
Simple Feature Engineering



```
# Feature binning: income
bins=[0,2500,4000,6000, 8000, 10000, 20000, 40000, 81000]
group=['Low', 'Average', 'medium', 'H1', 'h2', 'h3', 'h4', 'Very high']
data['Income_bin']= pd.cut(data['ApplicantIncome'],bins,labels=group)
```

```
[ ] data.head()
```

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_
0	Male	No	0	Graduate	No	5849	0.0	Nan	3
1	Male	Yes	1	Graduate	No	4583	1508.0	128.0	3





income - bins

0 → 81,000

✓ → domain-specific ↓
✓ → 5 bins: percentiles

0 - 20 21 - 40

percentiles

41 - 60 61 - 80 81 - 100

→ ML: decision trees (automatically bin)
= t(later)

10.2 - Hypothesis Testing | ST Kruskal-Wallis one-way analy... Anscombe's quartet 3 - Ansco... Bias of an estimator - Wikipedia EDA_FE.ipynb - Colaboratory scipy.stats.kruskal — SciPy v1.23.0 numpy.var — NumPy v1.23.0

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=19FtHN8LKZS_

+ Code + Text Reconnect

Loan_Status

{x}

Simple Feature Engineering

Feature binning: income
bins=[0,2500,4000,6000, 8000, 10000, 20000, 40000, 81000]
group=['Low', 'Average', 'medium', 'H1', 'h2', 'h3', 'h4' , 'Very high']
data['Income_bin']= pd.cut(data['ApplicantIncome'],bins,labels=group)

[] data.head()

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_
0	Male	No	0	Graduate	No	5849	0.0	Nan	3
1	Male	Yes	1	Graduate	No	4583	1508.0	128.0	3
2	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	3
3	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	3

94 / 94



+ Code + Text

Reconnect



▼ Incomes



{x}

```
[ ] #observed  
pd.crosstab(data["Income_bin"], data["Loan_Status"])
```



Loan_Status	N	Y
-------------	---	---

Income_bin		
------------	--	--

Low	34	74
-----	----	----

Average	67	159
---------	----	-----

medium	45	98
--------	----	----

H1	20	34
----	----	----

h2	9	22
----	---	----

h3	13	27
----	----	----

h4	3	6
----	---	---

Very high	1	2
-----------	---	---

+ Code + Text

Reconnect



Incomes

}

```
[ ] #observed  
pd.crosstab(data["Income_bin"], data["Loan_Status"])
```

Loan_Status	N	Y
Low	34	74
Average	67	159
medium	45	98
H1	20	34
h2	9	22
h3	13	27
h4	3	6
Very high	1	2



Income_bin	N	Y
Low	34	74
Average	67	159
medium	45	98
H1	20	34
h2	9	22
h3	13	27
h4	3	6

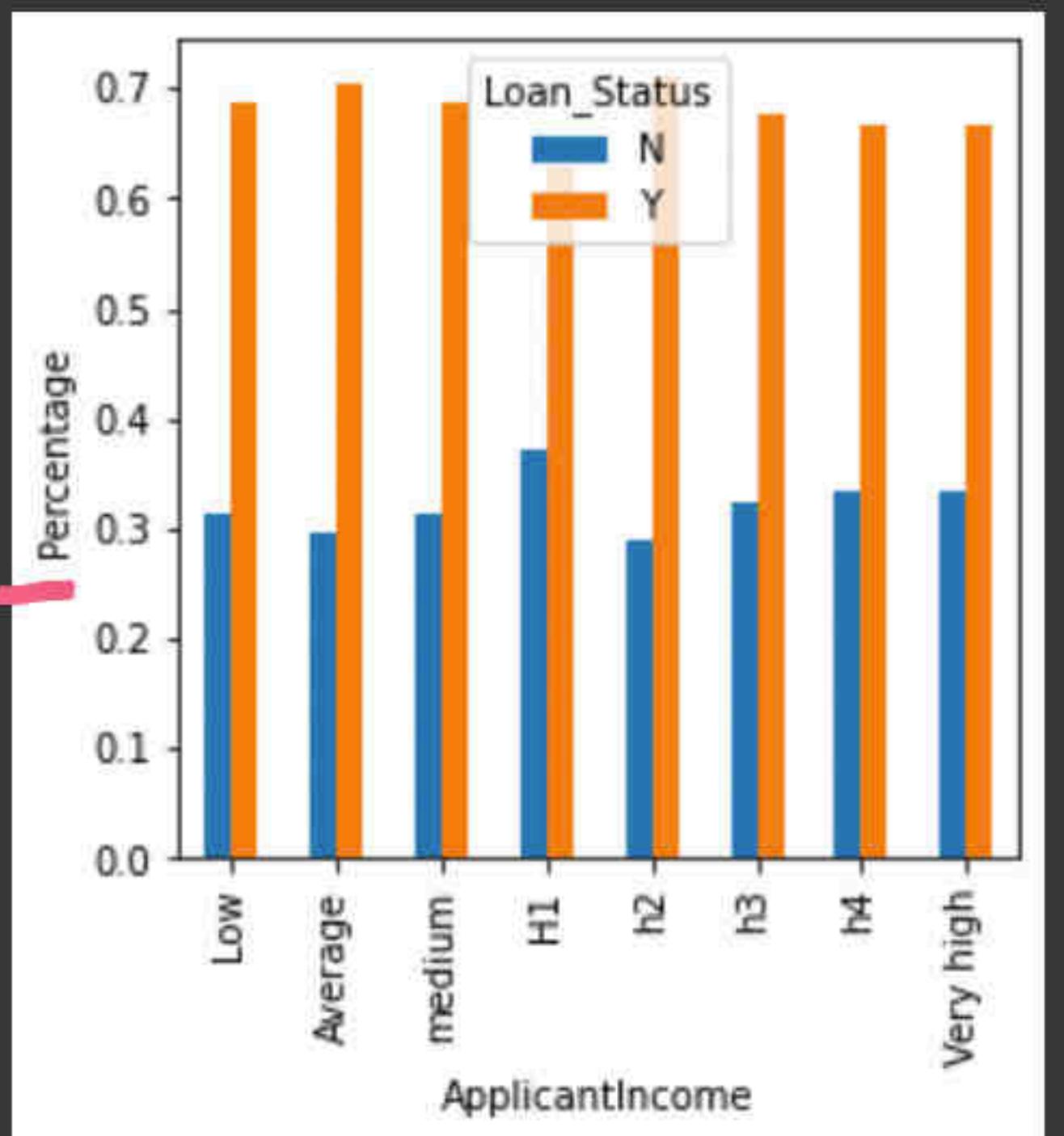
+ Code + Text

Reconnect



```
plt.xlabel("ApplicantIncome")
plt.ylabel("Percentage")
plt.show()
```

#It can be inferred that Applicant income does not affect the chances of loan approval. Which seems wrong as per our logic.



$$P(L = Y \mid \omega)$$

$$P(L = N \mid \omega)$$

```
[ ] #co-applicant income
```

```
bins=[0 1000 2000 4000 6000]
```



+ Code + Text

Reconnect



```
Income_bin.div(Income_bin.sum(axis=1),axis=0).plot(kind="bar",figsize=(4,4))  
plt.xlabel("ApplicantIncome")  
plt.ylabel("Percentage")  
plt.show()
```

{x} #It can be inferred that Applicant income does not affect the chances of loan approval. Which seems wrong as per ou



[] #co-applicant income



10.2 - Hypothesis Te x | Kruskal-Wallis one-w x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab x | scipy.stats.kruskal - x | numpy.var — NumPy x | pandas.cut — pandas x | pandas.qcut — pandas x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=19FtHN8LKZS_

+ Code + Text Reconnect

h2 9 22

h3 13 27

h4 3 6

{x}

Very high 1 2

```
[ ] Income_bin = pd.crosstab(data["Income_bin"], data["Loan_Status"])

Income_bin.div(Income_bin.sum(axis=1), axis=0).plot(kind="bar", figsize=(4,4))
plt.xlabel("ApplicantIncome")
plt.ylabel("Percentage")
plt.show()

#It can be inferred that Applicant income does not affect the chances of loan approval. Which seems wrong as per ou
```

Income Bin	N (%)	Y (%)
Very high	~32%	~68%
1	~32%	~68%
2	~32%	~68%

+ Code + Text

Reconnect

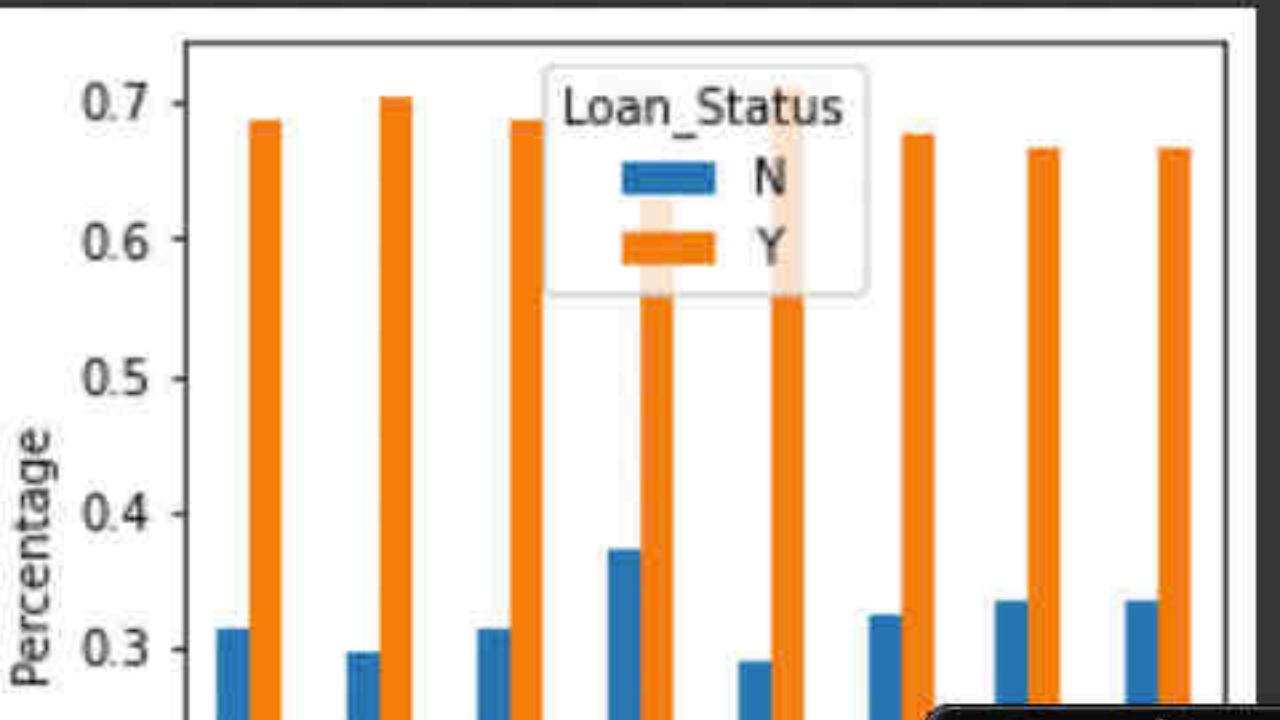


	h2	9	22
	h3	13	27
{x}	h4	3	6
	Very high	1	2

```
[ ] Income_bin = pd.crosstab(data["Income_bin"],data["Loan_Status"])

Income_bin.div(Income_bin.sum(axis=1),axis=0).plot(kind="bar",figsize=(4,4))
plt.xlabel("ApplicantIncome")
plt.ylabel("Percentage")
plt.show()
```

#It can be inferred that Applicant income does not affect the chances of loan approval. Which seems wrong as per ou



10.2 - Hypothesis x | Kruskal-Wallis one ... x | Anscombe's quart ... x | Bias of an estimate x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy x | pandas.cut — pandas x | pandas.qcut — pandas x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=9cnxWVMaK4aM

+ Code + Text Reconnect

[] CoapplicantIncome_Bin = pd.crosstab(data["CoapplicantIncome_bin"],data["Loan_Status"])
CoapplicantIncome_Bin.div(CoapplicantIncome_Bin.sum(axis = 1),axis=0).plot(kind='bar',figsize=(4,4))
plt.xlabel("CoapplicantIncome")
plt.ylabel("Percentage")
plt.show()

What's the problem here? Why co-applicant having low income is getting maximum loan approved?

CoapplicantIncome	N	Y
Low	~0.15	~0.85
Average	~0.30	~0.72
High	~0.35	~0.68

101 / 101

10.2 - Hypothesis x | Kruskal-Wallis on x | Anscombe's quart x | Bias of an estimate x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy x | pandas.cut — pandas x | pandas.qcut — pandas x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=TarGSsXDG8Tp

+ Code + Text

only numeric features

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	592.000000	600.000000	564.000000
mean	5403.459283	1621.245798	146.412162	342.000000	0.842199
std	6109.041673	2926.248369	85.587325	65.12041	0.364878
min	150.000000	{ 0.000000}	9.000000	12.000000	0.000000
25%	2877.500000	{ 0.000000}	100.000000	360.000000	1.000000
50%	3812.500000	1188.500000	128.000000	360.000000	1.000000
75%	5795.000000	2297.250000	168.000000	360.000000	1.000000
max	81000.000000	41667.000000	700.000000	480.000000	1.000000

low
0 - 1000
1000+

[] # categorical features
data.describe(include = ['object'])

	Gender	Married	Dependents	Education	Self_Employed	Property_Area	Loan_Status
count	601	611	599	614	582	614	614
unique	2	2	4	2	2	3	2

102 / 102

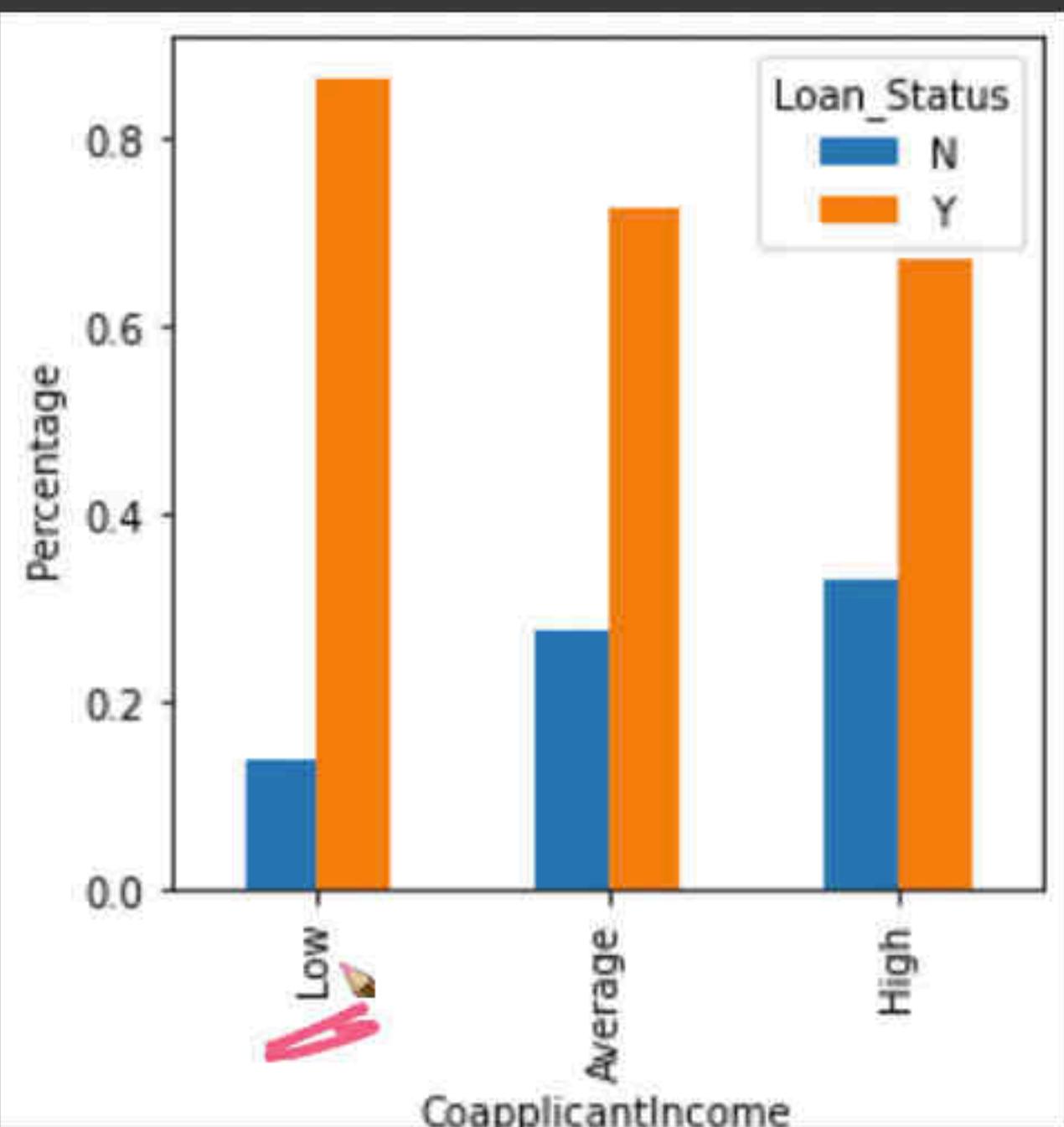
+ Code + Text

Reconnect



```
[ ] CoapplicantIncome_Bin = pd.crosstab(data[ "CoapplicantIncome_bin" ],data[ "Loan_Status" ])
CoapplicantIncome_Bin.div(CoapplicantIncome_Bin.sum(axis = 1),axis=0).plot(kind='bar',figsize=(4,4))
plt.xlabel("CoapplicantIncome")
plt.ylabel("Percentage")
plt.show()

## What's the problem here? Why co-applicant having low income is getting maximum loan approved?
```



10.2 - Hypothesis x | Kruskal-Wallis one-way... x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy x | pandas.cut — pandas x | pandas.qcut — pandas x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=TarGSsXDG8Tp

+ Code + Text Reconnect

0.0 Low Average High

CoapplicantIncome

613

{x}

```
[ ] data['CoapplicantIncome'].value_counts().head()
```

{ 0.0 273
2500.0 5
2083.0 5
1666.0 5
2250.0 3
Name: CoapplicantIncome, dtype: int64

273

```
[ ] # New feature: total household income  
data["TotalIncome"] = data["ApplicantIncome"] + data["CoapplicantIncome"]
```

```
[ ] bins = [0,3000,5000,8000,81000]  
group = ['Low', 'Average', 'High', 'Very High']  
data["TotalIncome_bin"] = pd.cut(data["TotalIncome"], bins, labels=group)
```

```
[ ] pd.crosstab(data["TotalIncome_bin"], data["Loan_Status"])
```

<https://colab.research.google.com/drive/1-PwKJwvKqxEfZ7jDx3nGxJvcZ3M4uN1#scrollTo=TDwzDQEenrhMK21>

Update

+ Code + Text

Reconnect ▾



```
{x}      Name: CoapplicantIncome, dtype: int64
```

```
# New feature: total household income  
data["TotalIncome"] = data["ApplicantIncome"] + data["CoapplicantIncome"]
```

```
[ ] bins = [0,3000,5000,8000,81000]
group = ['Low', 'Average', 'High', 'Very High']
data[ "TotalIncome_bin" ] = pd.cut(data[ "TotalIncome" ],bins,labels=group)
```

```
[ ] pd.crosstab(data[ "TotalIncome_bin" ], data[ "Loan_Status" ])
```

	N	Y
Loan_Status		
TotalIncome_bin		
Low	20	27
Average	69	154
High	61	151

10.2 - Hypothesis x | Kruskal-Wallis one-way ... x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy v1.19.0 API x | pandas.cut — pandas v1.0.3 API x | pandas.qcut — pandas v1.0.3 API x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=oDO5pryhMK21

+ Code + Text

Reconnect

Code

Text

Search

Copy

Reconnect

Code

Text

{x}

Name: CoapplicantIncome, dtype: int64

Up Down Reload Settings Copy Delete More

```
# New feature: total household income
data["TotalIncome"] = data["ApplicantIncome"] + data["CoapplicantIncome"]

[ ] bins = [0,3000,5000,8000,81000]
group = ['Low', 'Average', 'High', 'Very High']
data["TotalIncome_bin"] = pd.cut(data["TotalIncome"], bins, labels=group)
```

[] pd.crosstab(data["TotalIncome_bin"], data["Loan_Status"])

	Loan_Status	N	Y
TotalIncome_bin	Low	20	27
Average	69	154	
High	61	151	

106 / 106

10.2 - Hypothesis x | Kruskal-Wallis one ... x | Anscombe's quart ... x | Bias of an estimate x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy x | pandas.cut — pandas x | pandas.qcut — pandas x | pandas.DataFrame x | +

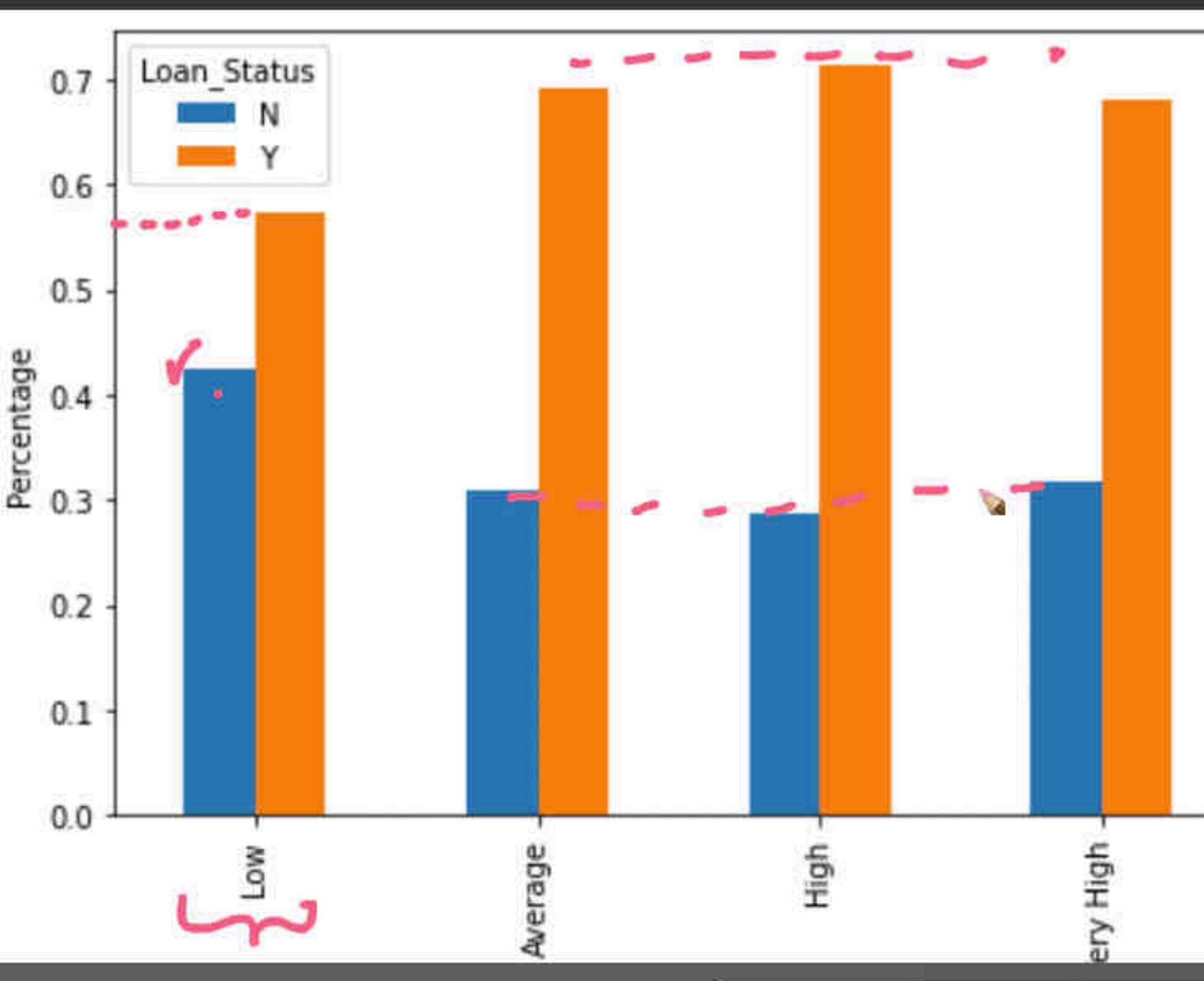
colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=oDO5pryhMK21

Reconnect  

+ Code + Text

```
plt.xlabel("TotalIncome")
plt.ylabel("Percentage")
plt.show()

{x}
# Observation: We can see that Proportion of loans getting approved for
# applicants having low Total_Income is very less as compared to that of applicants
# with Average, High and Very High Income.
```


The chart displays the percentage of loans approved ('Y') versus denied ('N') across four income categories: Low, Average, High, and Very High. The Y-axis represents the percentage from 0.0 to 0.7. The X-axis categories are Low, Average, High, and Very High. For each category, there are two bars: a blue bar for 'N' and an orange bar for 'Y'. A red dashed horizontal line is drawn at approximately 0.57. A red curly brace highlights the 'Low' category, and a red arrow points to the 'Very High' category. The chart shows that the percentage of approved loans is significantly higher for higher income groups compared to the low-income group.

TotalIncome	Loan_Status	Percentage
Low	N	~0.42
Low	Y	~0.58
Average	N	~0.31
Average	Y	~0.70
High	N	~0.29
High	Y	~0.72
Very High	N	~0.32
Very High	Y	~0.68

107 / 107

10.2 - Hypothesis x | Kruskal-Wallis one-way ... x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy v1.19.0 API x | pandas.cut — pandas v1.0.3 API x | pandas.qcut — pandas v1.0.3 API x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=hlyzygASMSB_

+ Code + Text Reconnect

TotalIncome Ver

```
[ ] data = data.drop(["Income_bin", "CoapplicantIncome_bin", "TotalIncome_bin"], axis=1)
```

Loan Amount and Loan Term

```
[ ] data['Loan_Amount_Term'].value_counts()
```

Loan_Amount_Term	Count
360.0	512
180.0	44
480.0	15
300.0	13
240.0	4
84.0	4
120.0	3
60.0	2
36.0	2
12.0	1

```
<> Name: Loan_Amount_Term, dtype: int64
```

```
[ ] data['Loan_Amount_Term'] = (data['Loan_Amount_Term']/12).astype('float')
```

① binning
② combine features

10.2 - Hypothesis x | Kruskal-Wallis on x | Anscombe's quart x | Bias of an estimate x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy x | pandas.cut — pandas x | pandas.qcut — pandas x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=hlyzygASMSB_

+ Code + Text Reconnect  

Loan Amount and Loan Term

```
{x} [ ] data['Loan_Amount_Term'].value_counts()
```

360.0 512
180.0 44
480.0 15
300.0 13
240.0 4
84.0 4
120.0 3
60.0 2
36.0 2
12.0 1
Name: Loan_Amount_Term, dtype: int64

```
[ ] data['Loan_Amount_Term'] = (data['Loan_Amount_Term']/12).astype('float')
```

```
[ ] sns.countplot(x='Loan_Amount_Term', data=data)  
plt.xlabel("Term in years")  
plt.show()  
# Observations: The count plot shows that most loans have a term less than 300 months, with the highest frequency occurring at 360 months.
```

10.2 - Hypothesis x | Kruskal-Wallis on x | Anscombe's quart x | Bias of an estimate x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy x | pandas.cut — pandas x | pandas.qcut — pandas x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=hlyzygASMSB_

+ Code + Text Reconnect ▾

Loan Amount and Loan Term

{x} [] data['Loan_Amount_Term'].value_counts()

✓ 360.0 512
180.0 44
480.0 15
300.0 13
240.0 4
84.0 4
120.0 3
60.0 2
36.0 2
12.0 1
Name: Loan_Amount_Term, dtype: int64

[] data['Loan_Amount_Term'] = (data['Loan_Amount_Term']/12).astype('float')

[] sns.countplot(x='Loan_Amount_Term', data=data)
plt.xlabel("Term in years")
plt.show()
Observations: Most loans have a term less than 300 months, with the distribution skewed towards shorter terms.

110 / 110

+ Code + Text

Reconnect



▼ Loan Amount and Loan Term

{x} [] data['Loan_Amount_Term'].value_counts()

360.0	512
180.0	44
480.0	15
300.0	13
240.0	4
84.0	4
120.0	3
60.0	2
36.0	2
12.0	1

{ } 360.0 512
180.0 44
480.0 15
300.0 13
240.0 4
84.0 4
120.0 3
60.0 2
36.0 2
12.0 1

Name: Loan_Amount_Term, dtype: int64

[] data['Loan_Amount_Term'] = (data['Loan_Amount_Term']/12).astype('float')

[] sns.countplot(x='Loan_Amount_Term', data=data)
plt.xlabel("Term in years")
plt.show()

+ Code + Text

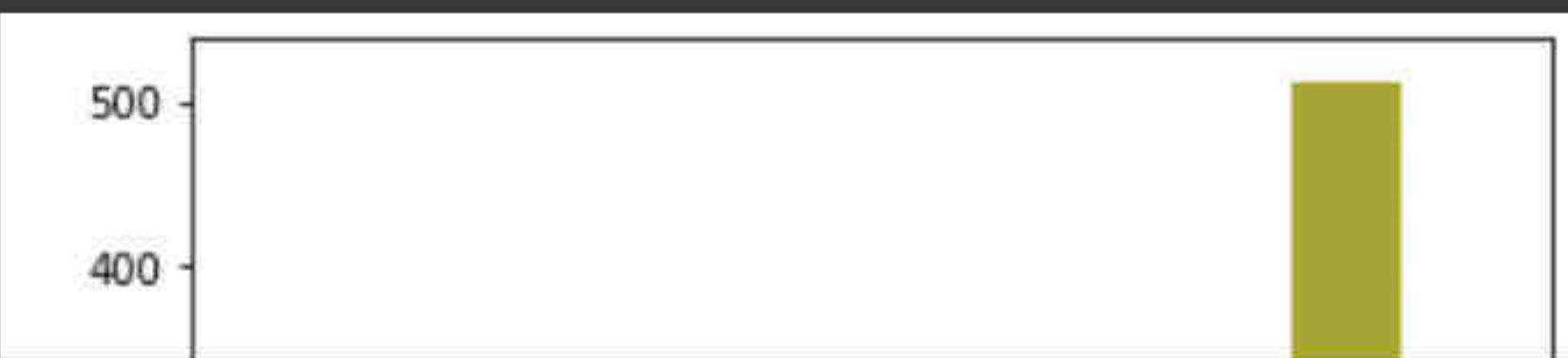
Reconnect



```
180.0  
480.0    15  
300.0    13  
240.0     4  
84.0      4  
{x} 120.0    3  
60.0      2  
36.0      2  
12.0      1  
Name: Loan_Amount_Term, dtype: int64
```

```
[ ] data['Loan_Amount_Term'] = (data['Loan_Amount_Term']/12).astype('float')
```

```
[ ] sns.countplot(x='Loan_Amount_Term', data=data)  
plt.xlabel("Term in years")  
plt.show()  
# Observation: We can clearly see that more than 90% of the loans were applied for 30 years.
```

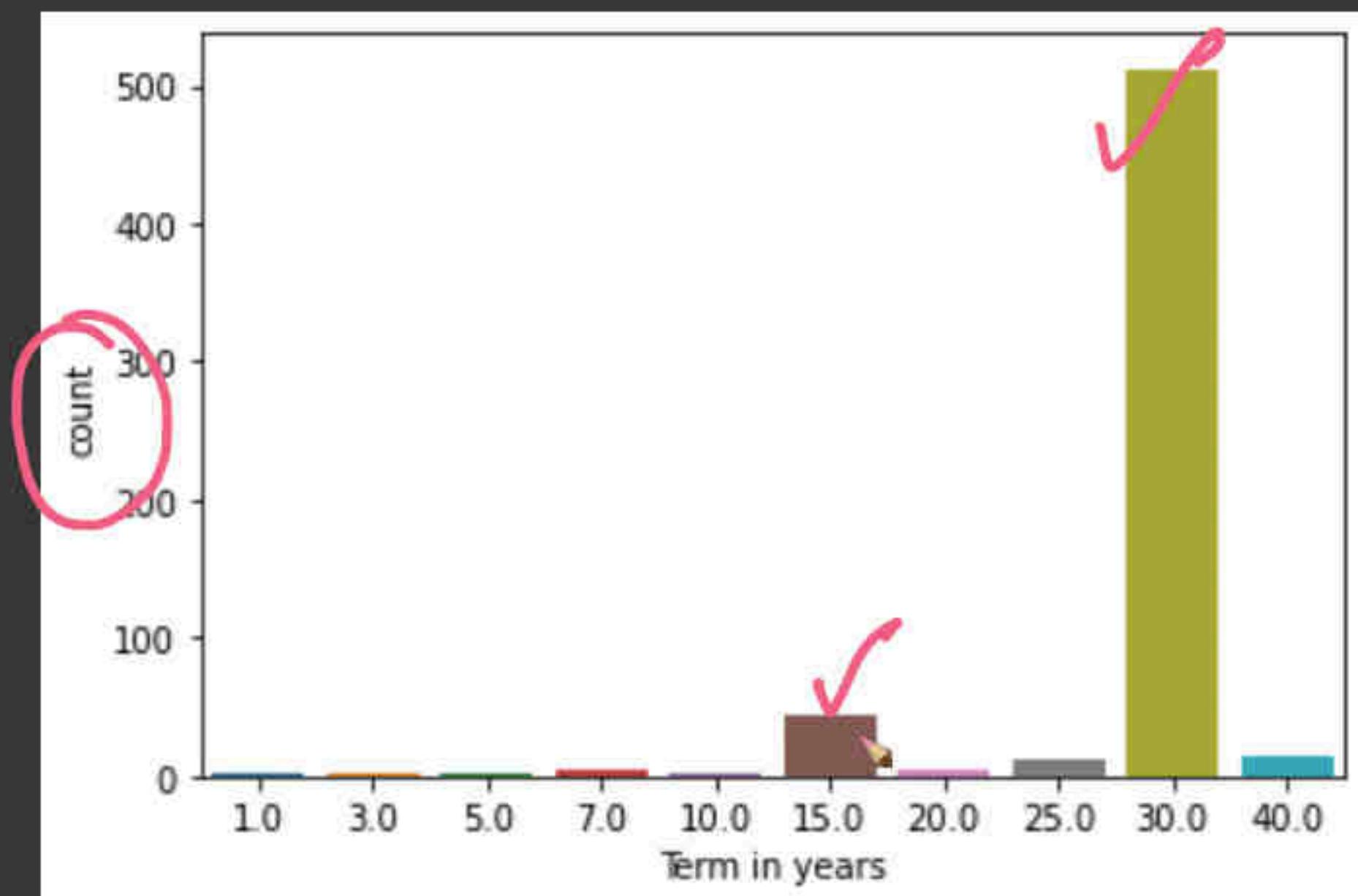


+ Code + Text

Reconnect

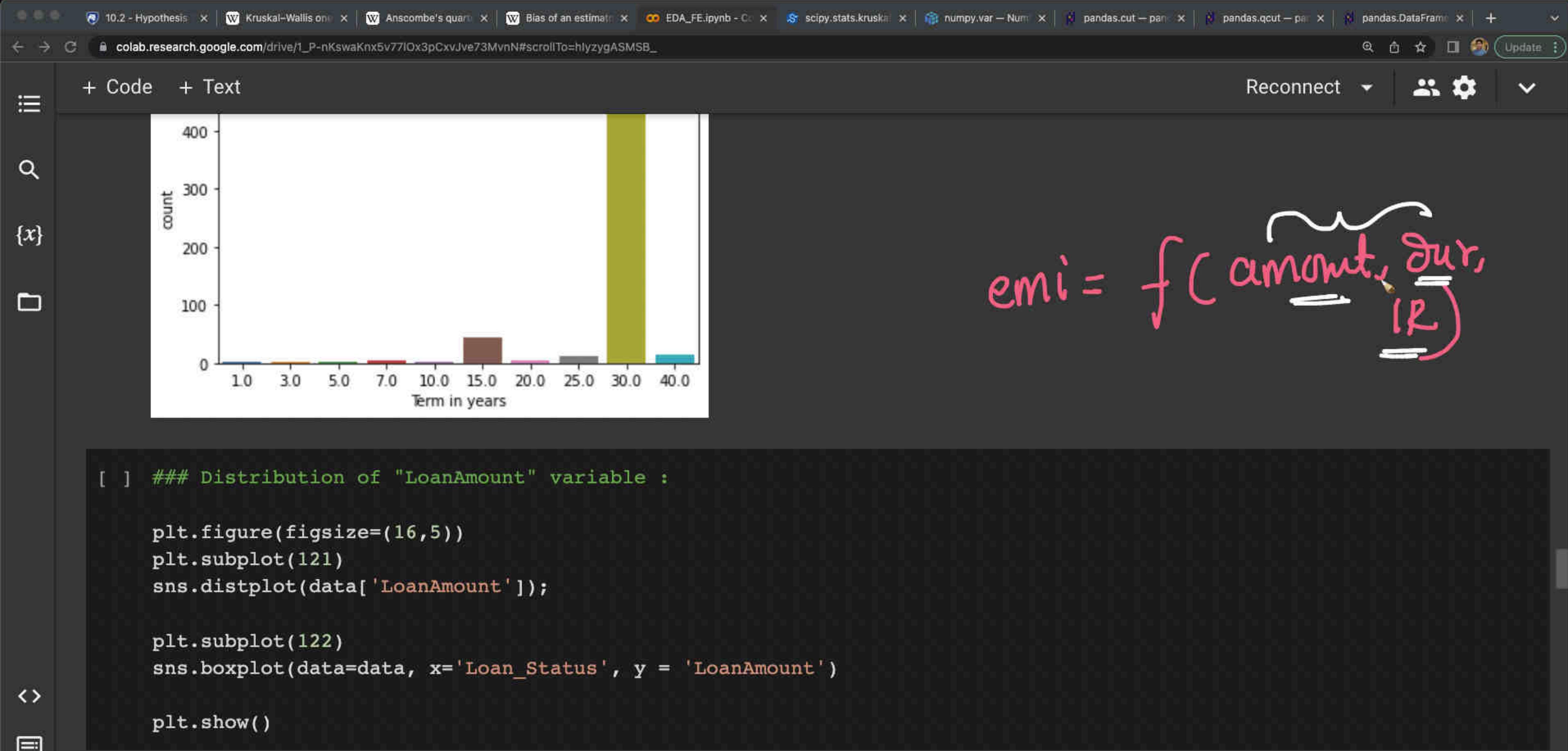


```
[ ] sns.countplot(x='Loan_Amount_Term', data=data)
plt.xlabel("Term in years")
plt.show()
# Observation: We can clearly see that more than 90% of the loans were applied for 30 years.
```



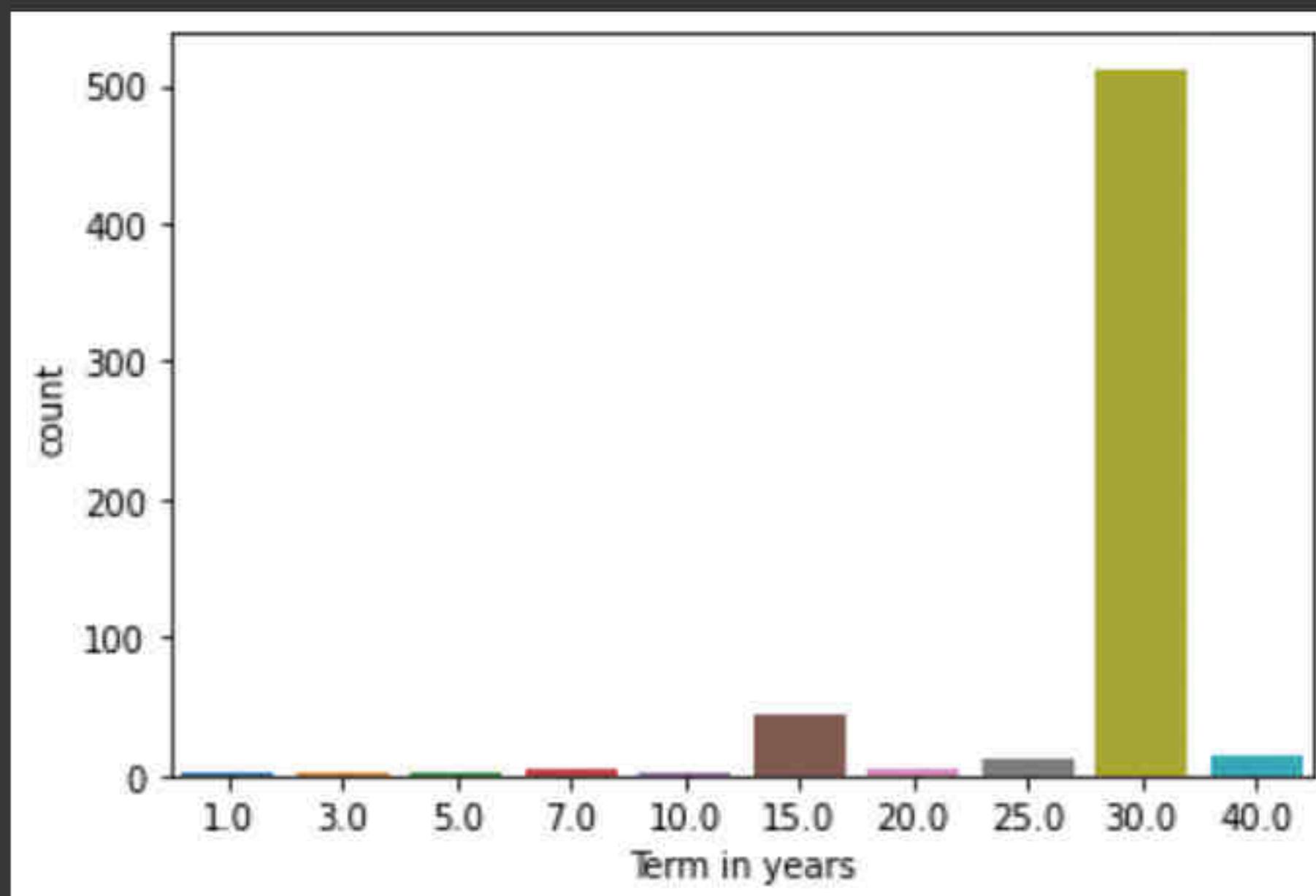
```
[ ] ### Distribution of "LoanAmount" variable :
```





$$\text{emi} = f(\text{amount}, \text{dur}, \underline{\text{IR}})$$

```
[ ] sns.countplot(x='Loan_Amount_Term', data=data)
plt.xlabel("Term in years")
plt.show()
# Observation: We can clearly see that more than 90% of the loans were applied for 30 years.
```



emi \approx $\frac{\text{amt}}{\# \text{mths}}$ ignoring IR

VS MI

emi
MI

Y/N

```
[ ] ### Distribution of "LoanAmount" variable :
```

```
plt.figure(figsize=(16, 5))
```

+ Code + Text

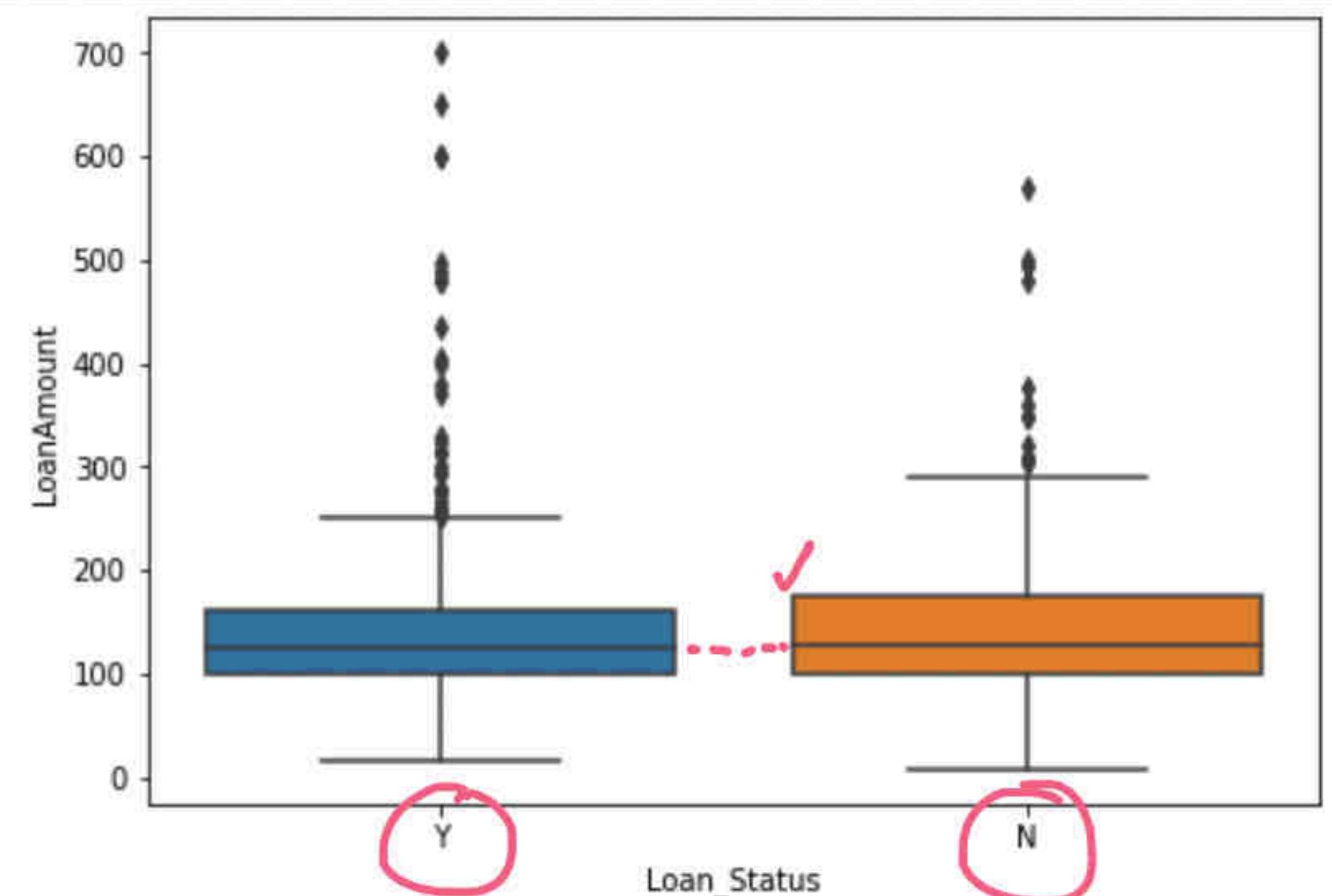
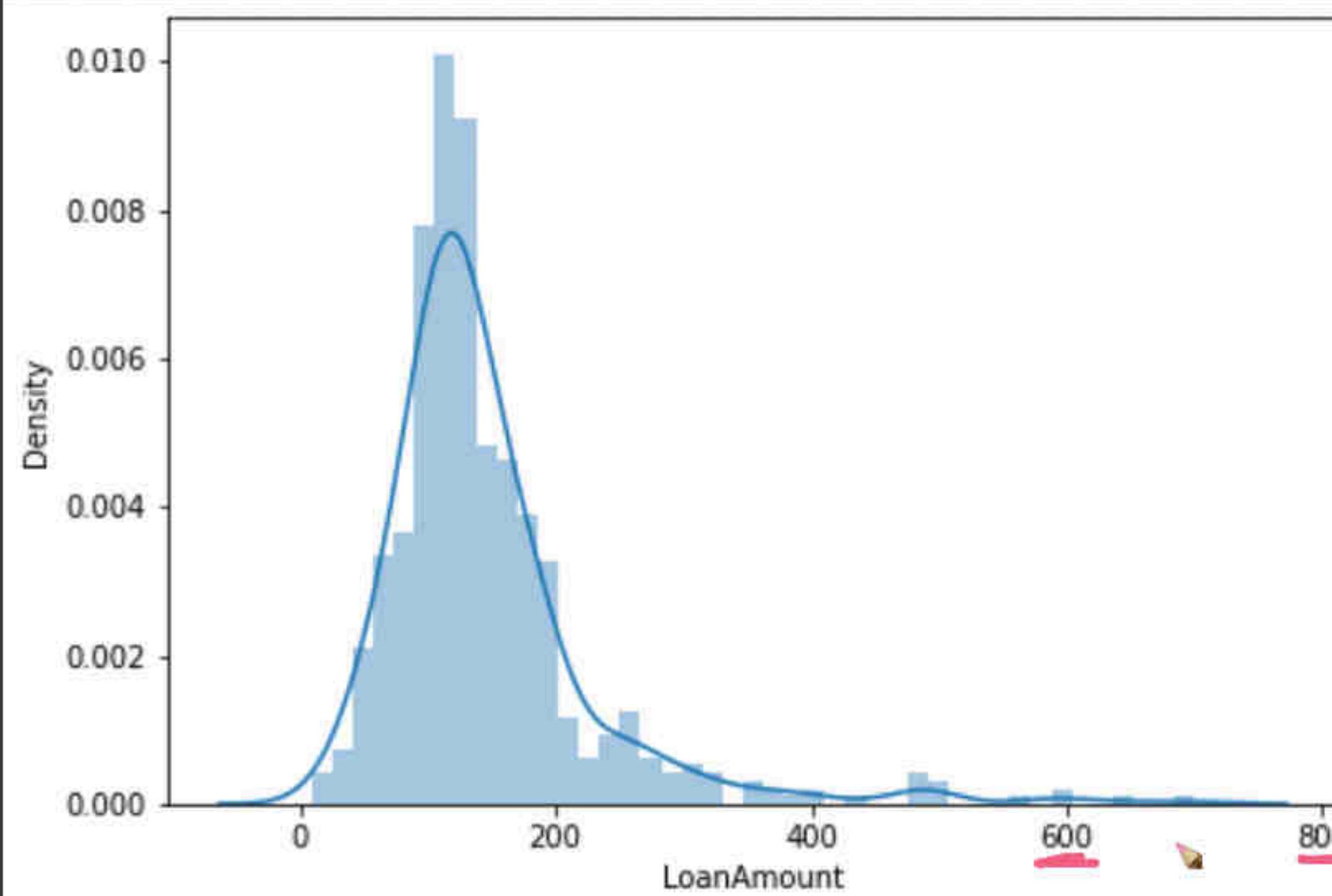
Reconnect



```
plt.subplot(121)
sns.distplot(data['LoanAmount']);

plt.subplot(122)
sns.boxplot(data=data, x='Loan_Status', y = 'LoanAmount')

plt.show()
```



10.2 - Hypothesis x | Kruskal-Wallis one-way analysis of variance by ranks x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy v1.17.0 documentation x | pandas.cut — pandas 0.24.2 documentation x | pandas.qcut — pandas 0.24.2 documentation x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=K4uMgnkQNE39

+ Code + Text Reconnect

LoanAmount

0.000

0 200 400 600 800

0

Y

N

Loan_Status

{x}

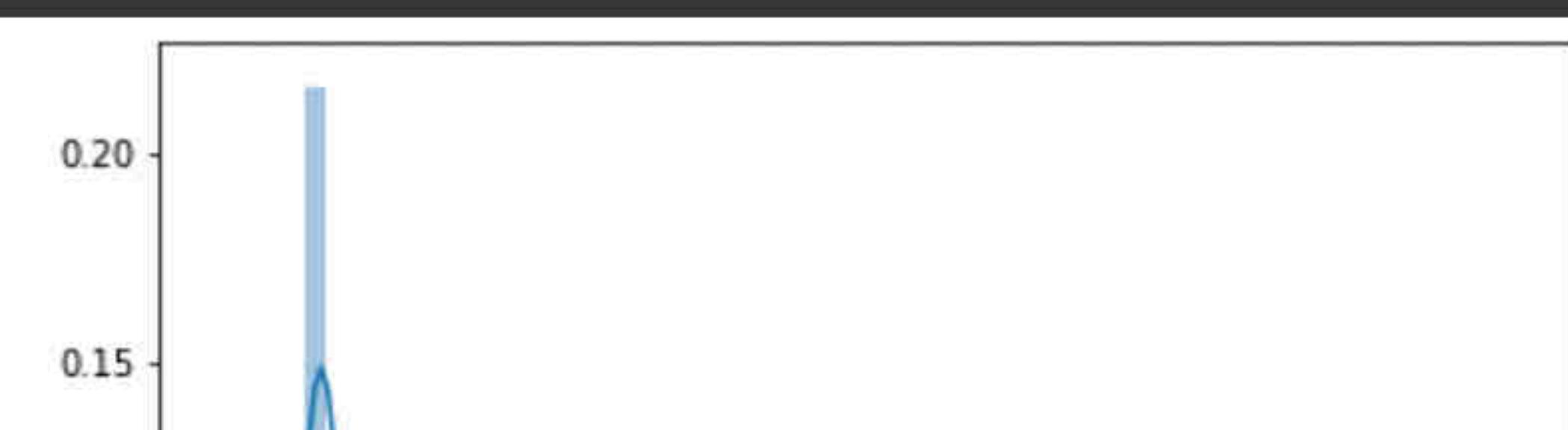
Approximate calc: ignoring interest rates as we dont know that.

data['Loan_Amount_per_year'] = data['LoanAmount']/data['Loan_Amount_Term']

[] plt.figure(figsize=(16,5))
plt.subplot(121)
sns.distplot(data['Loan_Amount_per_year']);

plt.subplot(122)
sns.boxplot(data=data, x='Loan_Status', y = 'Loan_Amount_per_year')

plt.show()



10.2 - Hypothesis x | Kruskal-Wallis one-way analysis of variance by ranks x | Anscombe's quartet x | Bias of an estimator x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy v1.18.0 documentation x | pandas.cut — pandas v1.1.3 documentation x | pandas.qcut — pandas v1.1.3 documentation x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=K4uMgnkQNE39

+ Code + Text Reconnect

```
sns.boxplot(data=data, x='Loan_Status', y = 'Loan_Amount_per_year')

plt.show()
```

{x}

□

Density

0.20
0.15
0.10
0.05
0.00

0 20 40 60 80 100 120

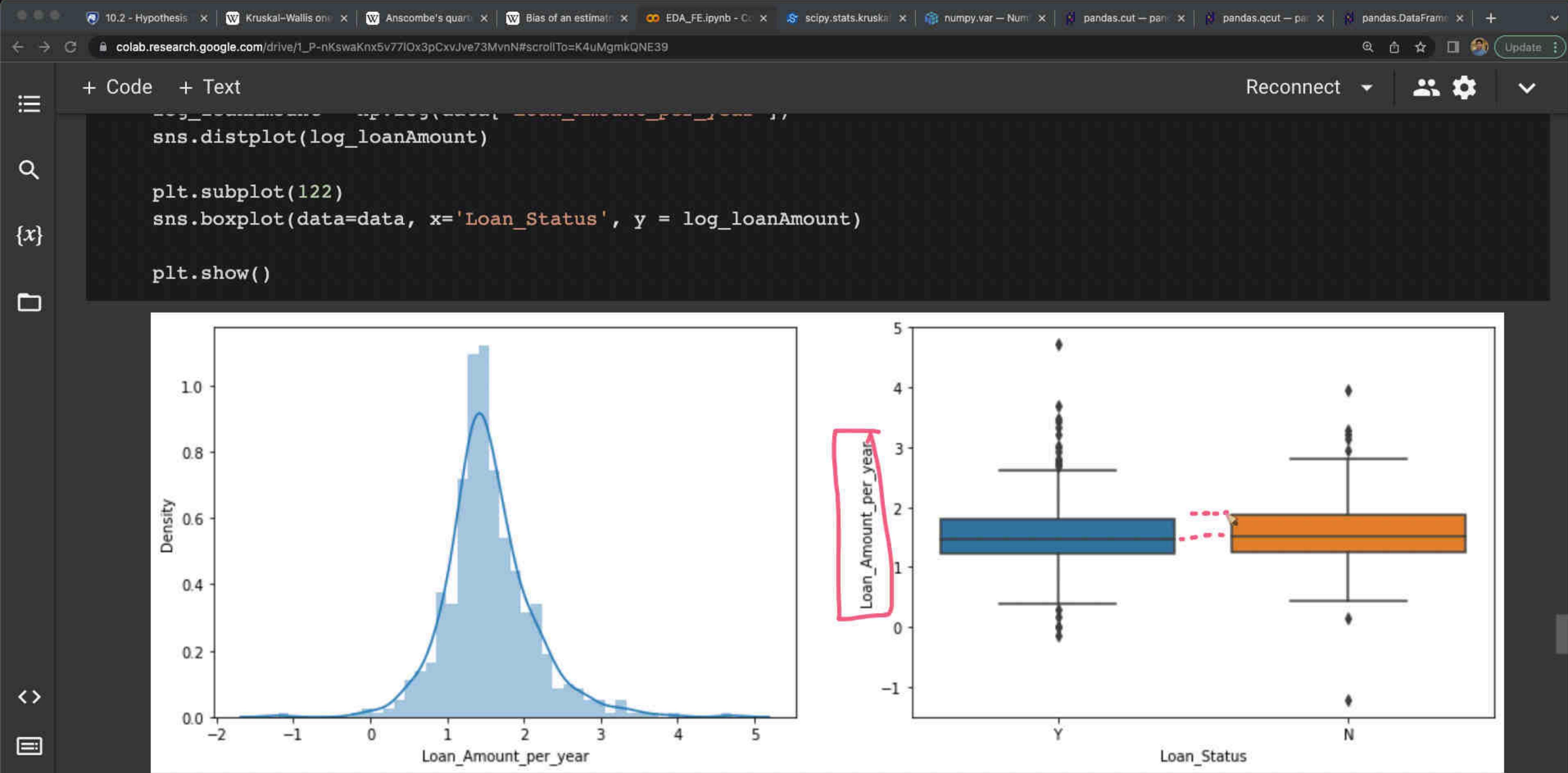
Loan_Amount_per_year

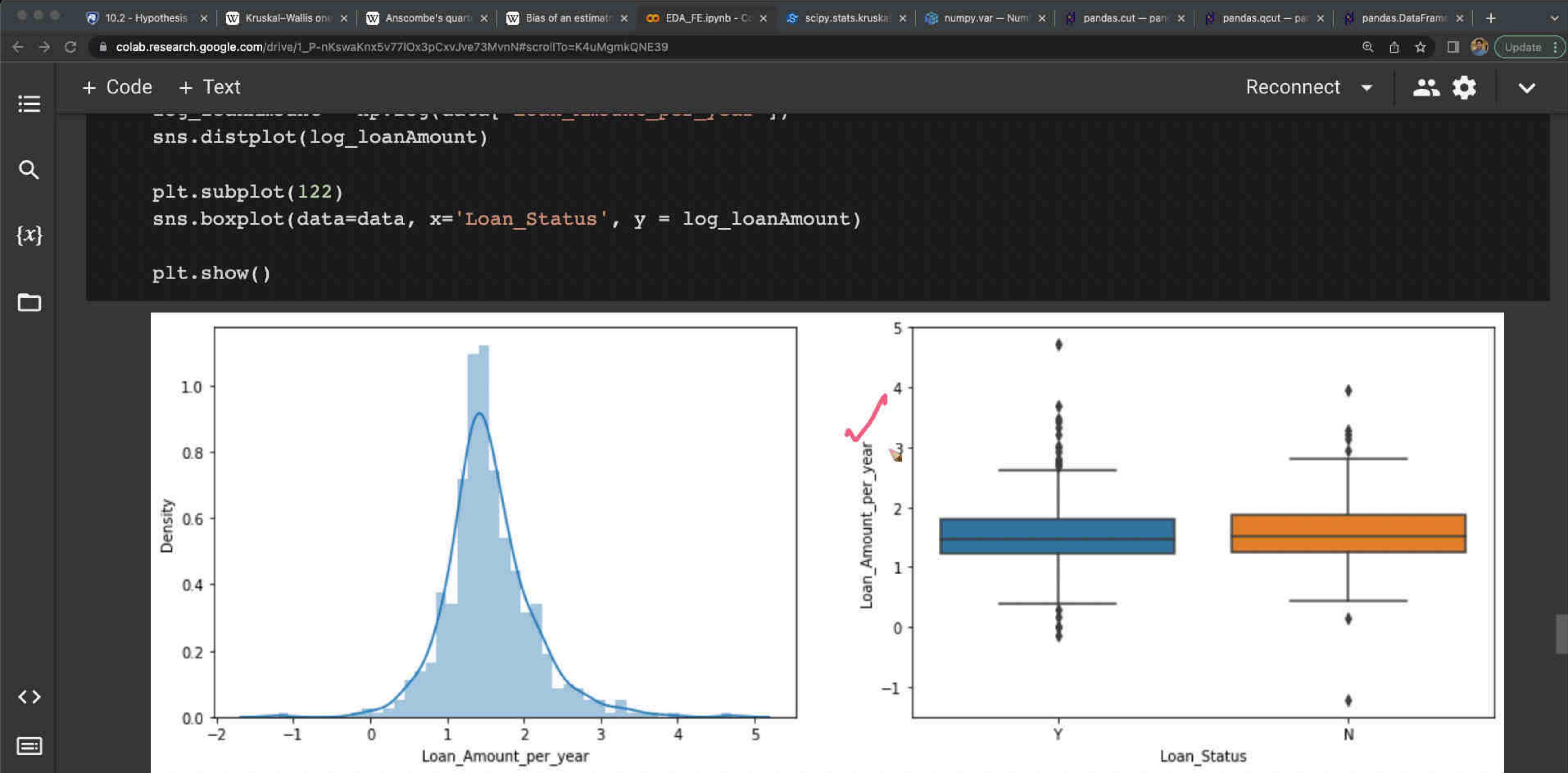
Y N

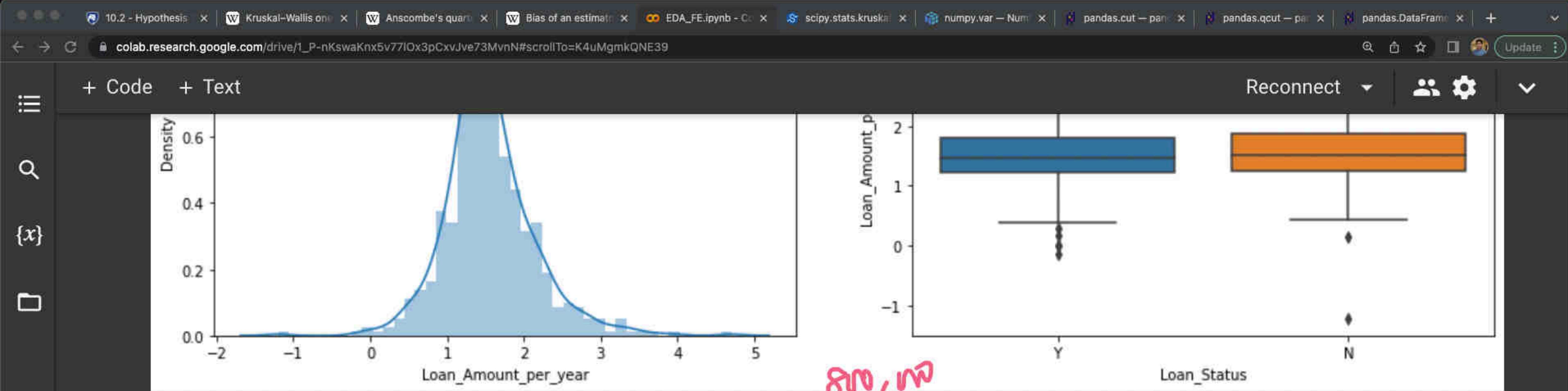
Loan_Status

118 / 118

[] # log transform
plt.figure(figsize=(16,5))
plt.subplot(121)



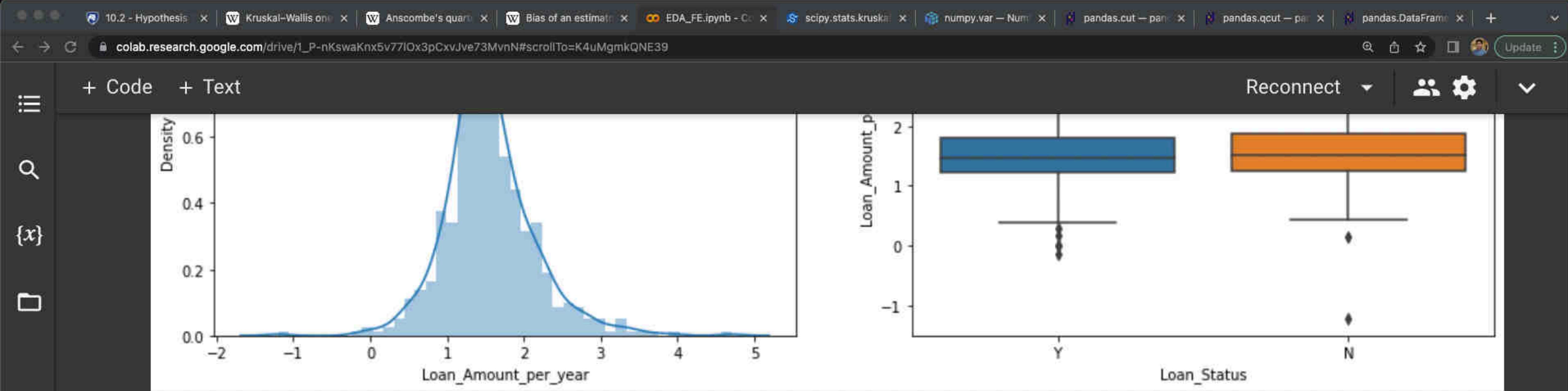




```
[ ] # Feature : Calculate the EMI based on the Loan Amount Per year.  
data['EMI'] = data['Loan_Amount_per_year']*1000/12
```

```
[ ] #Feature : Able_to_pay_EMI  
data['Able_to_pay_EMI'] = (data['TotalIncome']*0.1 > data['EMI']).astype('int')
```

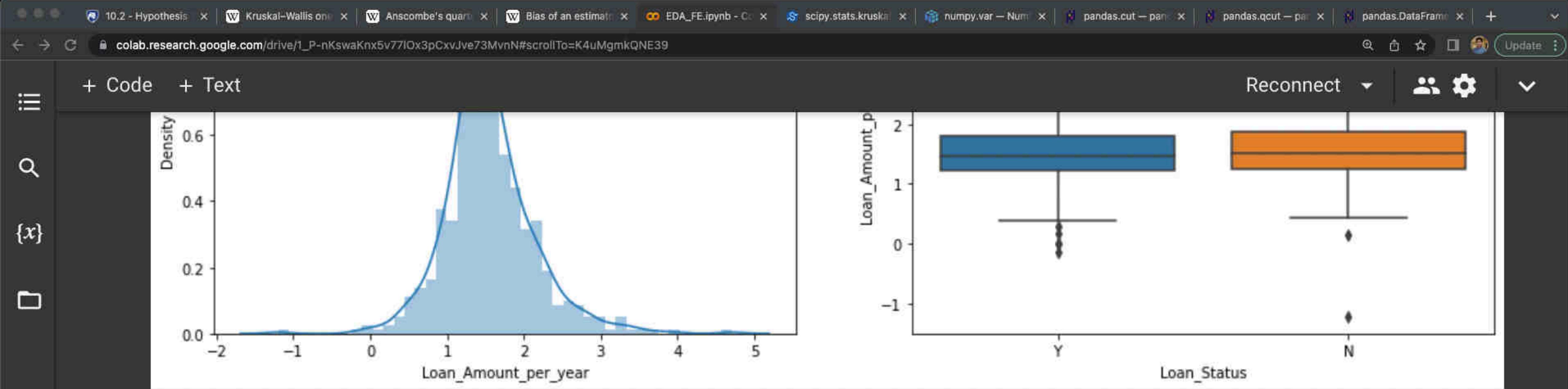
```
[ ] sns.countplot(x='Able_to_pay_EMI', data = data, hue = 'Loan_Status')  
#Observation:  
###There is 50% chance that you may get the loan approved if you cannot pay the EMI.  
###But there, is a 72% chance that you may get the loan approved if you can pay the EMI.
```



```
[ ] # Feature : Calculate the EMI based on the Loan Amount Per year.  
data['EMI'] = data['Loan_Amount_per_year']*1000/12
```

```
[ ] #Feature : Able_to_pay_EMI  
data['Able_to_pay_EMI'] = (data['TotalIncome']*0.1 > data['EMI']).astype('int')
```

```
[ ] sns.countplot(x='Able_to_pay_EMI', data = data, hue = 'Loan_Status')  
#Observation:  
###There is 50% chance that you may get the loan approved if you cannot pay the EMI.  
###But there, is a 72% chance that you may get the loan approved if you can pay the EMI.
```



```
[ ] # Feature : Calculate the EMI based on the Loan Amount Per year.
data['EMI'] = data['Loan_Amount_per_year']*1000/12
```

$$\text{emi} = \frac{\text{loan-amt} \times 1000}{\text{duration}} / 12$$

```
[ ] #Feature : Able_to_pay_EMI
data['Able_to_pay_EMI'] = (data['TotalIncome']*0.1 > data['EMI']).astype('int')
```

```
[ ] sns.countplot(x='Able_to_pay_EMI', data = data, hue = 'Loan_Status')
#Observation:
###There is 50% chance that you may get the loan approved if you cannot pay the EMI.
###But there, is a 72% chance that you may get the loan approved if you can pay the EMI.
```

10.2 - Hypothesis x | Kruskal-Wallis one ... x | Anscombe's quart ... x | Bias of an estimate x | EDA_FE.ipynb - Colab Notebooks x | scipy.stats.kruskal x | numpy.var — NumPy x | pandas.cut — pandas x | pandas.qcut — pandas x | pandas.DataFrame x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=XIWTmwhDZEcN

+ Code + Text Reconnect

Q {x} D

#Feature : Able_to_pay_EMI
data['Able_to_pay_EMI'] = (data['TotalIncome']*0.1 > data['EMI']).astype('int')

[] sns.countplot(x='Able_to_pay_EMI', data = data, hue = 'Loan_Status')
#Observation:
###There is 50% chance that you may get the loan approved if you cannot pay the EMI.
###But there, is a 72% chance that you may get the loan approved if you can pay the EMI.

<matplotlib.axes._subplots.AxesSubplot at 0x7f3e2eadc950>

Able_to_pay_EMI	Loan_Status	Count
0	Y	~50
0	N	~50
1	Y	~350
1	N	~150

124 / 124

10.2 - Hypothesis Testing | Kruskal-Wallis | Anscombe's qu... | Bias of an estimator | EDA_FE.ipynb | scipy.stats.kruskal | numpy.var — Numpy v1.20.2 API | pandas.cut — pandas 1.3.3 documentation | pandas.qcut — pandas 1.3.3 documentation | pandas.DataFrame | Pareto distribution | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=FxqCjW7laETR

+ Code + Text Reconnect  

3. Target Encoding

Appropriate encoding depends on what our task is (and) what we do next?

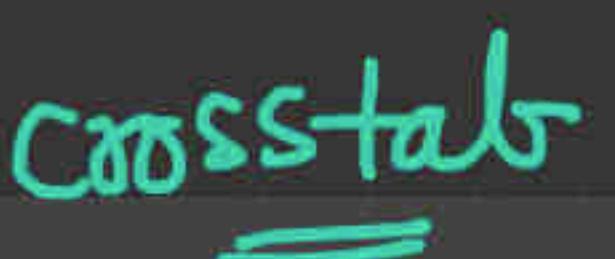
Task: Compute Correlation (PCC and SRCC) between each feature and the Loan-Status

[] ↳ 30 cells hidden

▶ Column Standardization and Normalization

- Mean centering and Variance scaling (Standard Scaling)
- MinMax Scaling

[] ↳ 1 cell hidden


$$\left[\begin{array}{c} P(\underline{\text{PI}} | M) \\ \hline \end{array} \right]$$

125 / 125



$$T_{\chi^2} = \sum_{i=1}^{\text{#cells}} \frac{(O_i - E_i)^2}{E_i} \sim \chi^2(\kappa)$$

$$(n_{rows}-1) \times (n_{cols}-1)$$



Google

Search Google or type a URL



Colaboratory



YouTube



My Drive



InterviewBit S...



Learning



GitHub



Scaler Academ...



Reduce the fil...



(4) Feed



Add shortcut



127 / 127