



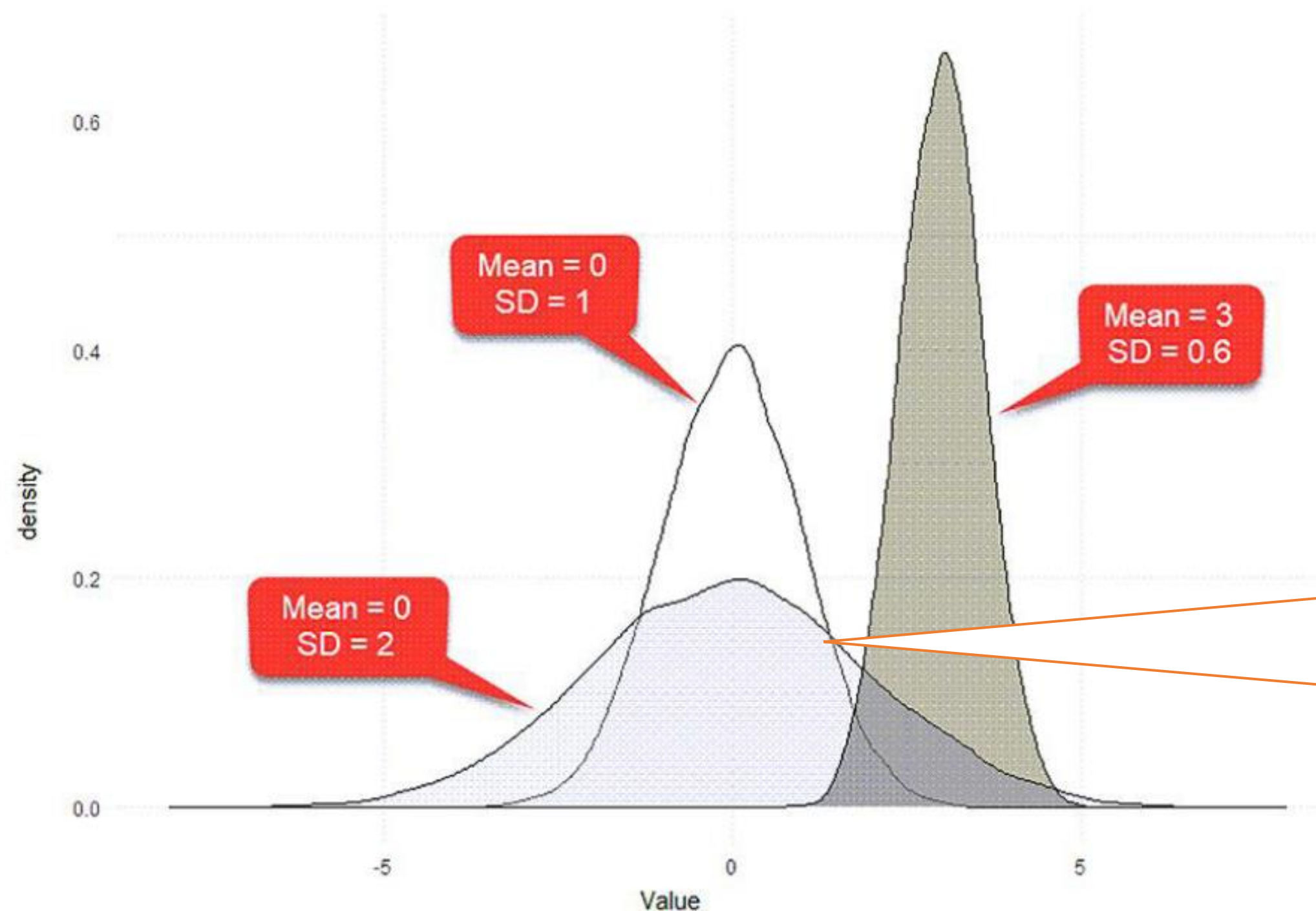
Dealing with outliers

Outliers handled 3 ways

- Remove from dataset
- Make more “palatable” with appropriate scaling and/or transformation
- Analyze with specialized techniques that focus on detecting & predicting outliers

Standard deviation method

- Works best for Gaussian distributions $\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$
- A univariate Gaussian



Why does height decrease as the distribution becomes wide?

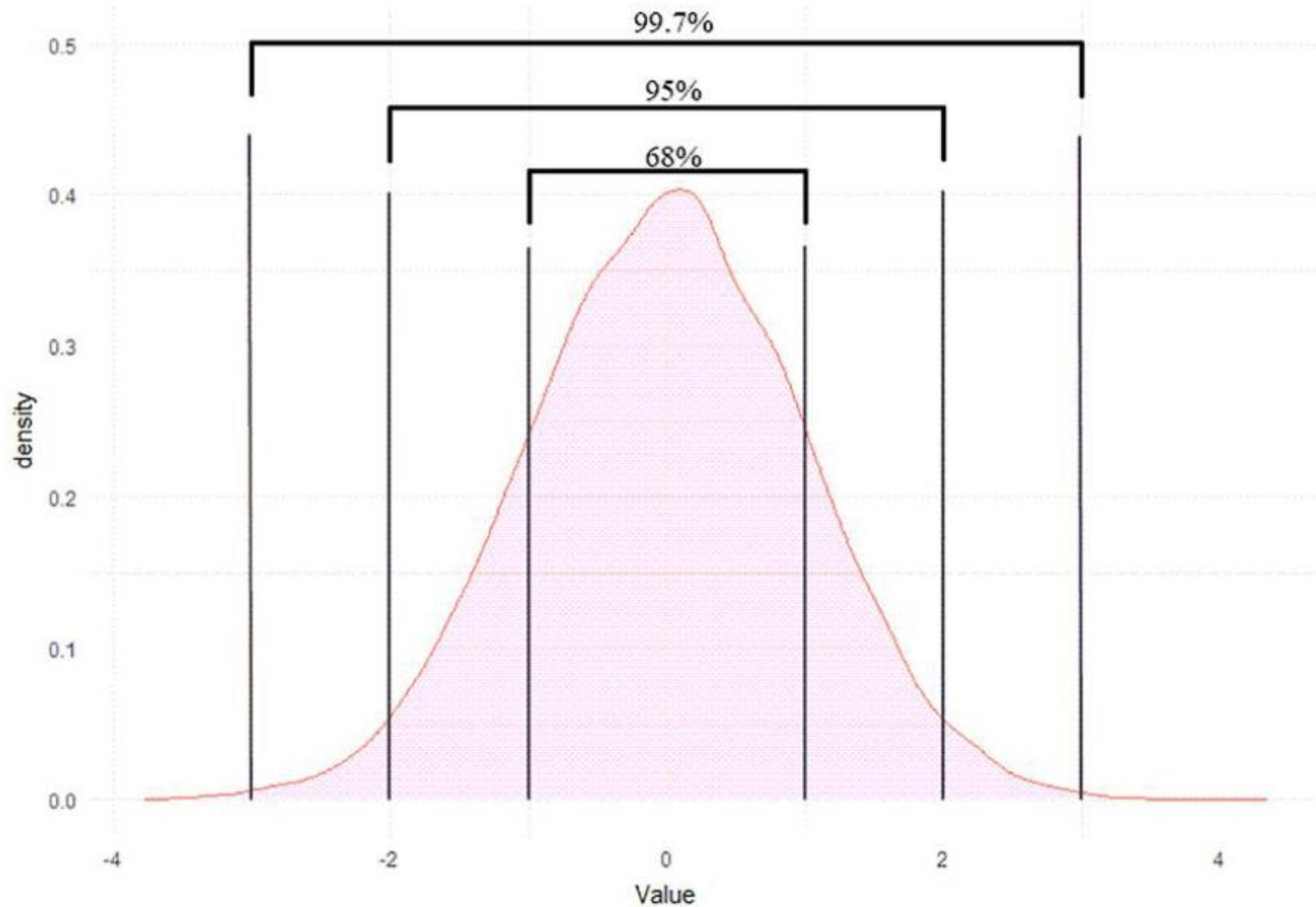
Standard deviation method (contd.)

- Standard deviation is the typical deviation of feature value from mean

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2} = \frac{\|x - \mu \mathbf{1}\|}{\sqrt{n}}$$

- <https://www.geogebra.org/calculator/ve2earn>

Empirical Formula for Gaussian Distribution



Outlier v/s anomaly

- Used interchangeably
- Difference between outlier and anomaly
 - Outlier is determined by technical considerations
 - Anomaly is additionally determined by business considerations
- Outlier is generally
 - 2.5 Standard deviations away
 - 3 standard deviations away
 - Or some SD per our choice suiting the problem

Removing outliers during data preprocessing

- Data farther than 3 standard deviation from mean
 - Impacts mean – pulling it in the direction of outlier
 - Skews prediction
- How to deal?
 - Delete first
 - Apply Standard Scaler for the rest

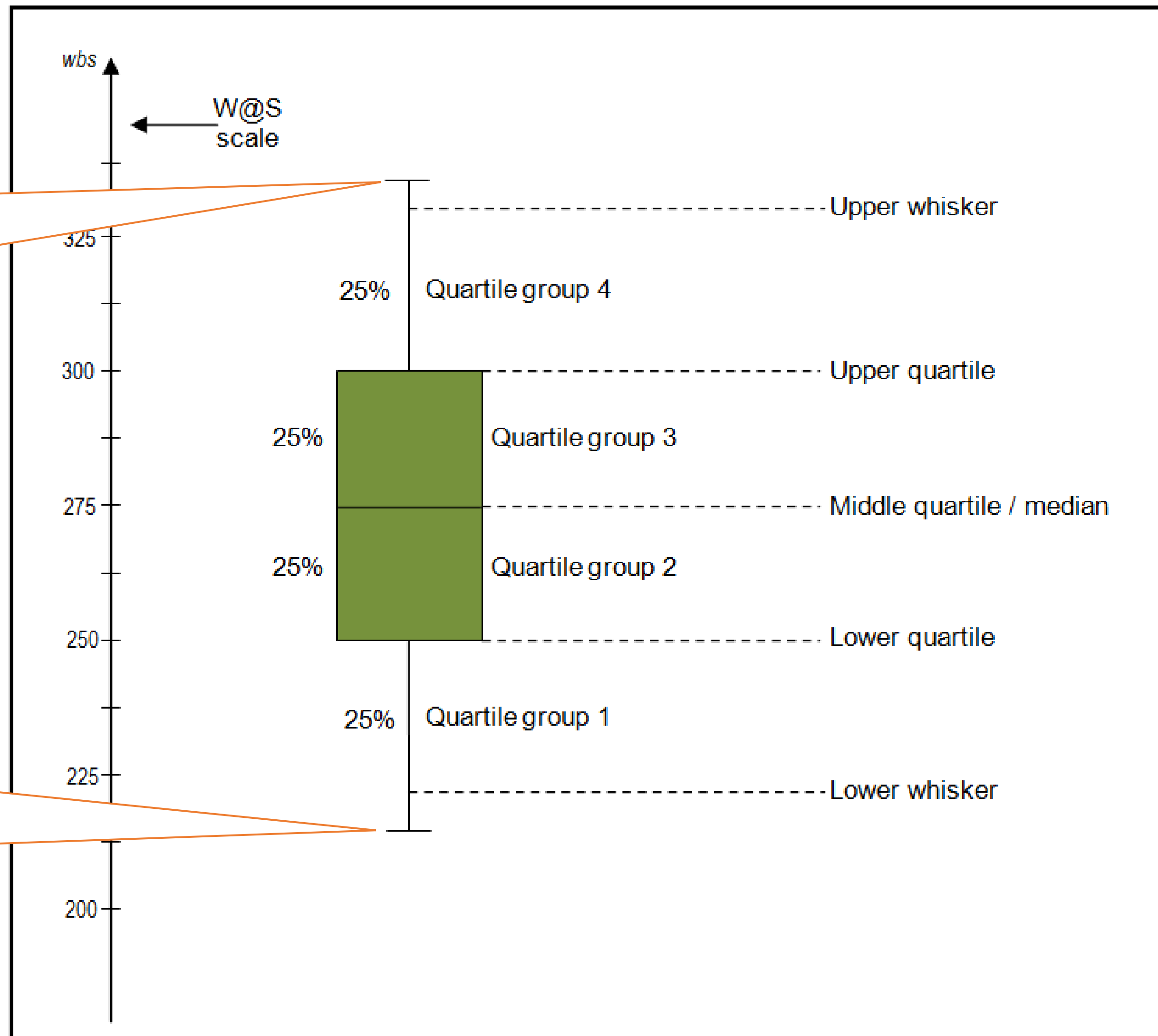
**Convention to
treat any
feature
transformation
as function**

$$\phi(x) = z = \frac{x - \mu}{\sigma}$$

Retaining outliers

- Box plot analysis
- Very extreme values can be removed
- Then Robust Scaler

**Do we care
about max in
the context of
outliers?**

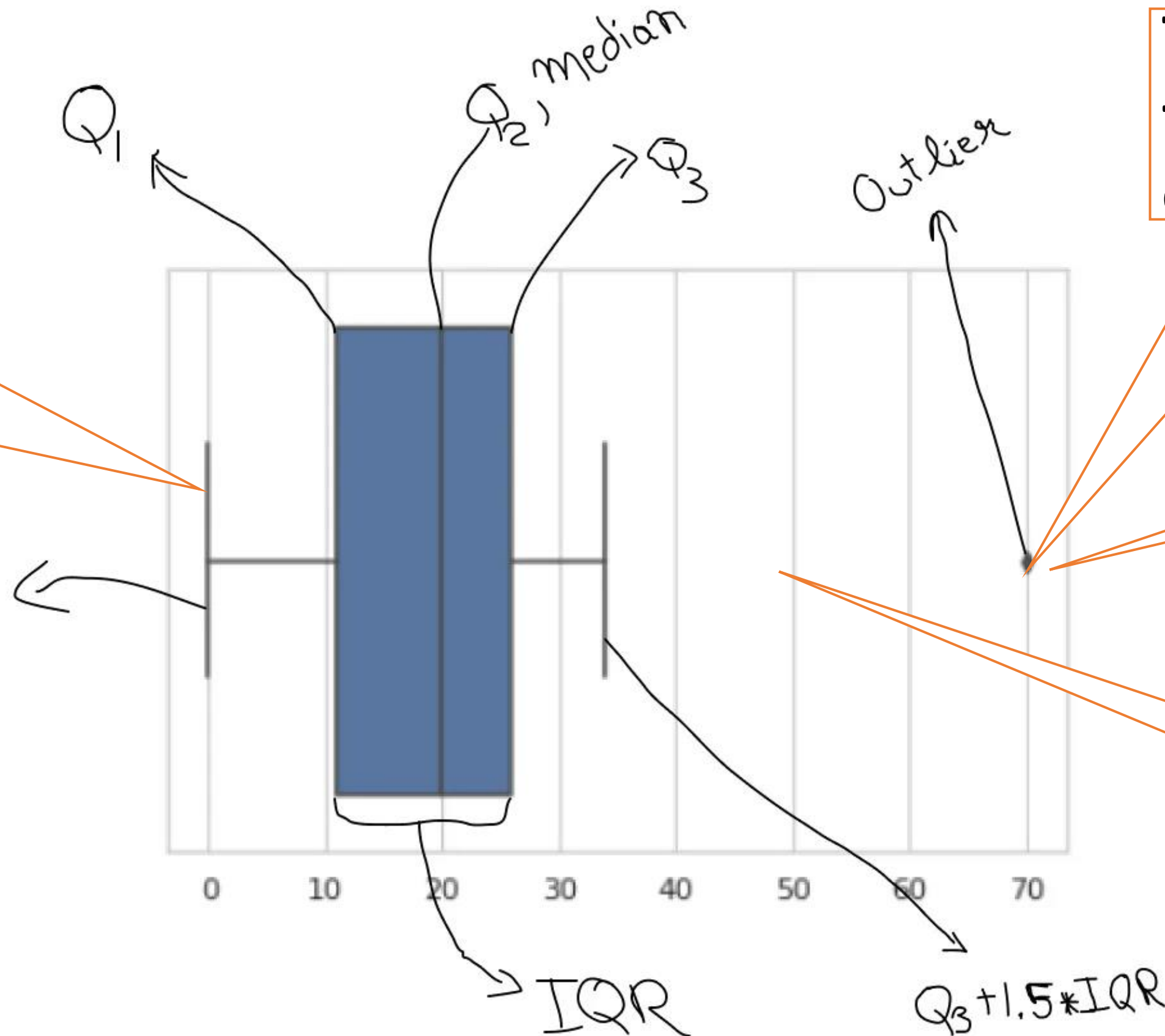


**Do we care
about min in
the context of
outliers?**

Retaining outliers – Box Plots

Internal fences are created like this. Anything outside this is outlier

$$Q_1 - 1.5 * IQR$$



This is max, but this is also outlier

We MAY decide to delete this

We may retain some points from this region

Retaining outliers – Why 1.5 IQR?

1.0

Lower Bound:

$$\begin{aligned} &= Q1 - 1 * IQR \\ &= Q1 - 1 * (Q3 - Q1) \\ &= -0.675\sigma - 1 * (0.675 - [-0.675])\sigma \\ &= -0.675\sigma - 1 * 1.35\sigma \\ &= -2.025\sigma \end{aligned}$$

Upper Bound:

$$\begin{aligned} &= Q3 + 1 * IQR \\ &= Q3 + 1 * (Q3 - Q1) \\ &= 0.675\sigma + 1 * (0.675 - [-0.675])\sigma \\ &= 0.675\sigma + 1 * 1.35\sigma \\ &= 2.025\sigma \end{aligned}$$

1.5

Lower Bound:

$$\begin{aligned} &= Q1 - 1.5 * IQR \\ &= Q1 - 1.5 * (Q3 - Q1) \\ &= -0.675\sigma - 1.5 * (0.675 - [-0.675])\sigma \\ &= -0.675\sigma - 1.5 * 1.35\sigma \\ &= -2.7\sigma \end{aligned}$$

Upper Bound:

$$\begin{aligned} &= Q3 + 1.5 * IQR \\ &= Q3 + 1.5 * (Q3 - Q1) \\ &= 0.675\sigma + 1.5 * (0.675 - [-0.675])\sigma \\ &= 0.675\sigma + 1.5 * 1.35\sigma \\ &= 2.7\sigma \end{aligned}$$

- 1.7 IQR = 3 Standard deviation
- 1.5 IQR = 2.7 SD captures 99.65% data in Gaussian

Retaining outliers during data pre-processing

- Apply Robust Scaler and retain data

$$\phi(x) = \frac{x - Q_2}{Q_3 - Q_1}$$

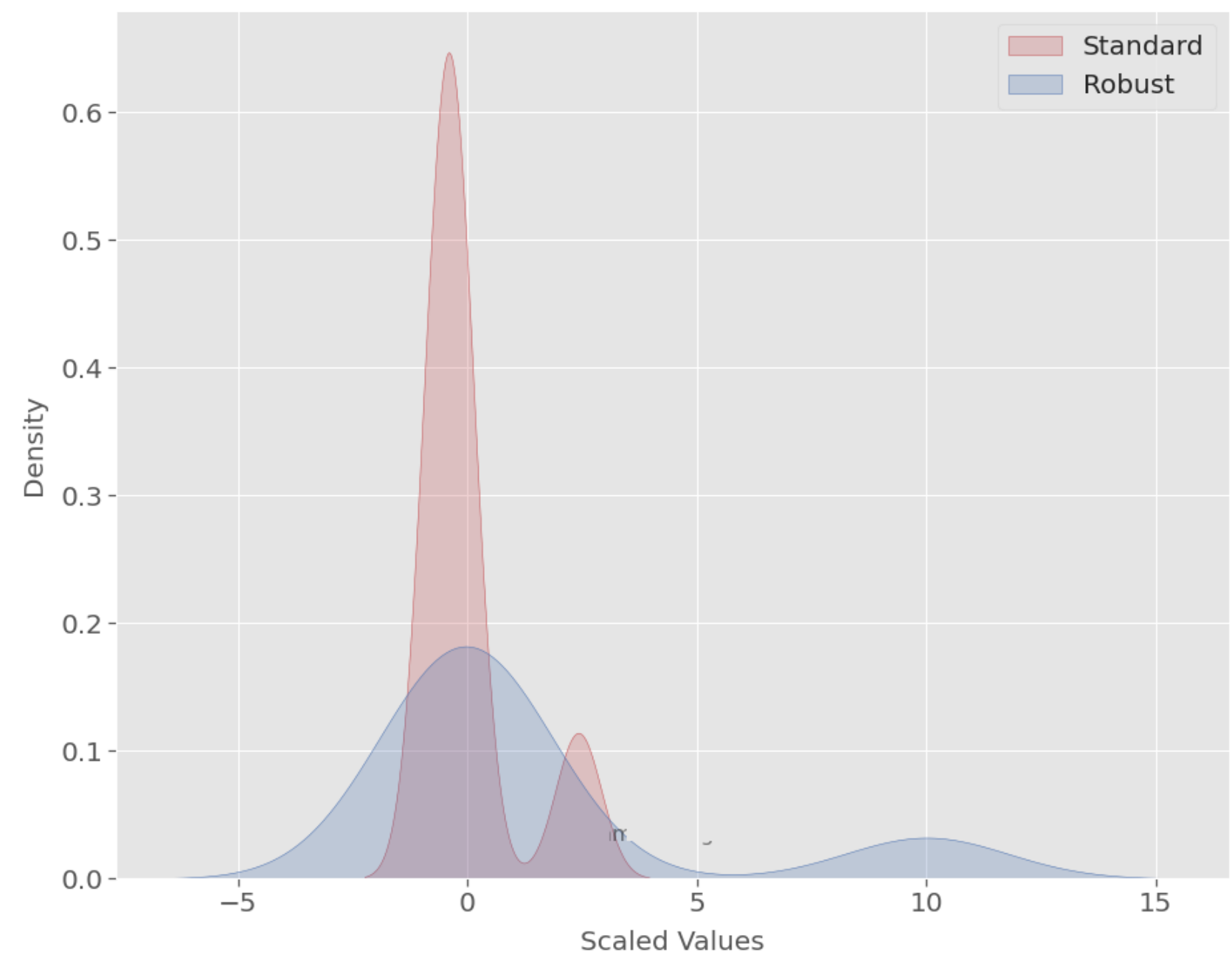
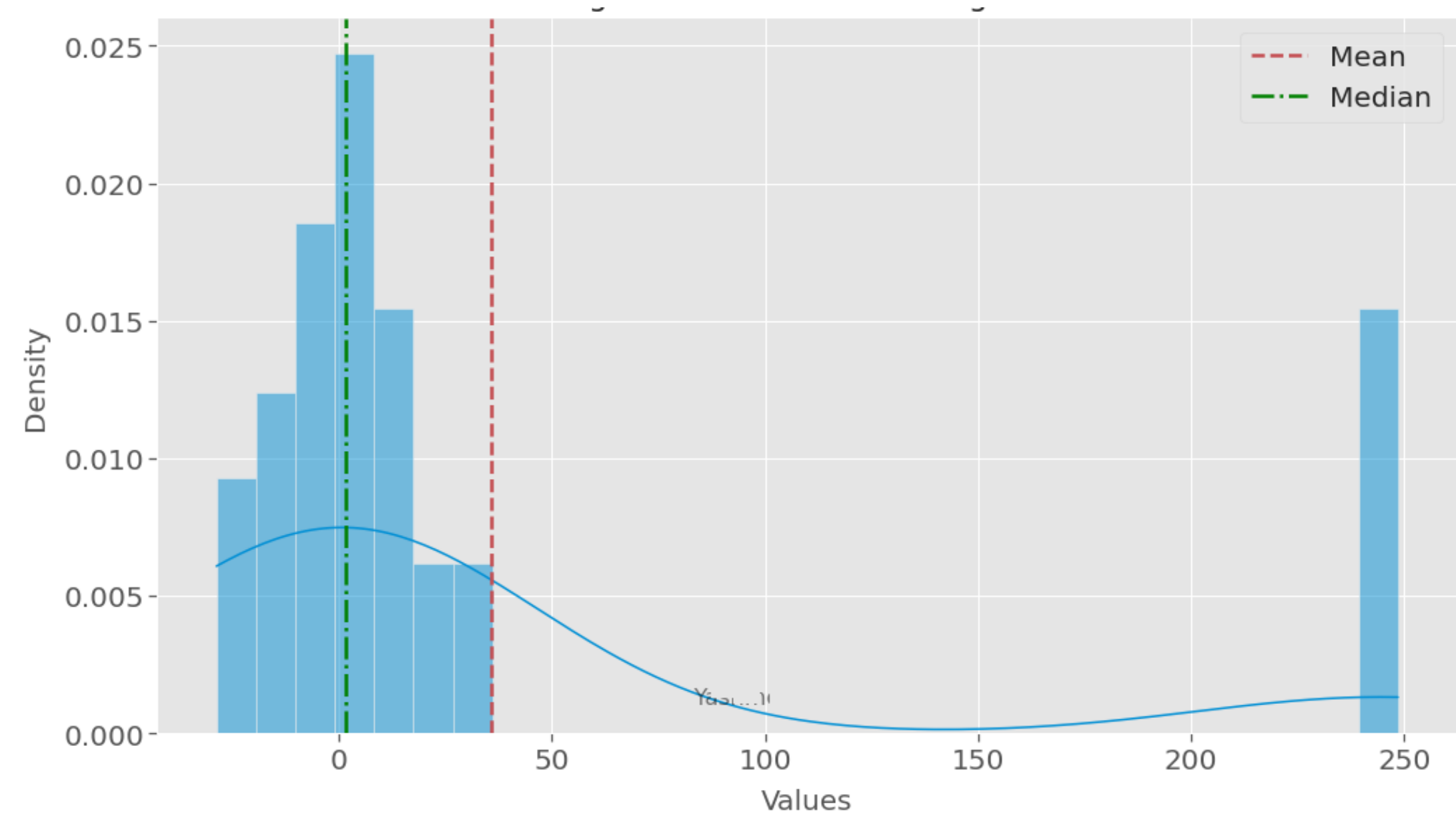
- Q_2 = Median, $Q_3 - Q_1$ = IQR
- Median is not severely affected by the outlier

Why Robust Scaler?

	No outliers	With outliers
Min	-28.72	-28.72
Max	32.45	248.51
Range	61.17	277.23
Mean	0.92	35.71
Median	-0.38	1.50
Standard Deviation	15.47	87.64
IQR	17.23	24.30

	Standard	Robust
Min	-0.75	-1.24
Max	2.46	10.17
Range	3.21	11.41
Mean	-0.00	1.41
Median	-0.40	0.00
Standard Deviation	1.01	3.61
IQR	0.28	1.00

- Robust Scaler provides wide range for feature
- Standard Scaler shrinks the range due to outliers
- More variance in scaled data is good for prediction



Sample sessional problem

- A dataset has 3 features
- Outliers were deleted on first two features
- Standard scaler was applied on the first feature
- MinMax Scaler on second feature
- Third feature ranges between 1500 to 10,000
 - Robust Scaler was applied

Are 3 features ML ready now? *yes*

- Clue:
 - What is the range of Standard, Minmax & Robust Scaler?

MinMaxScaler

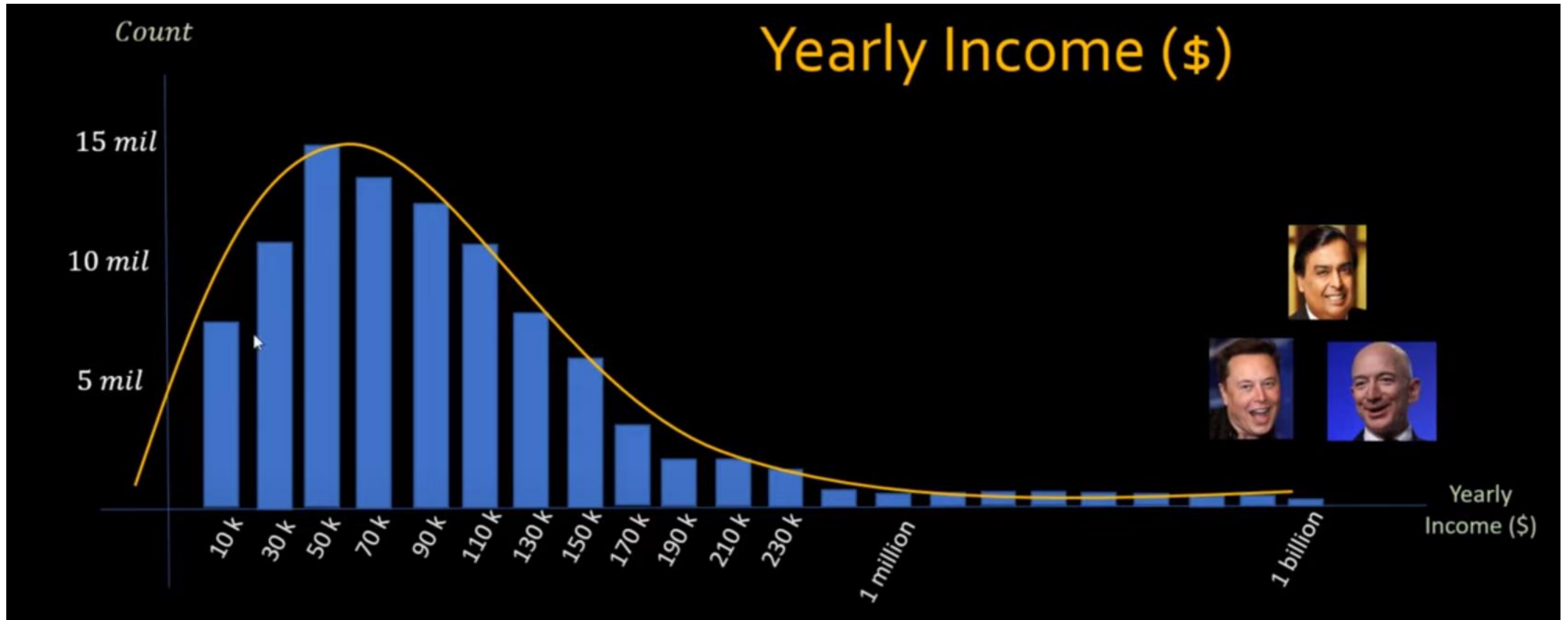
$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

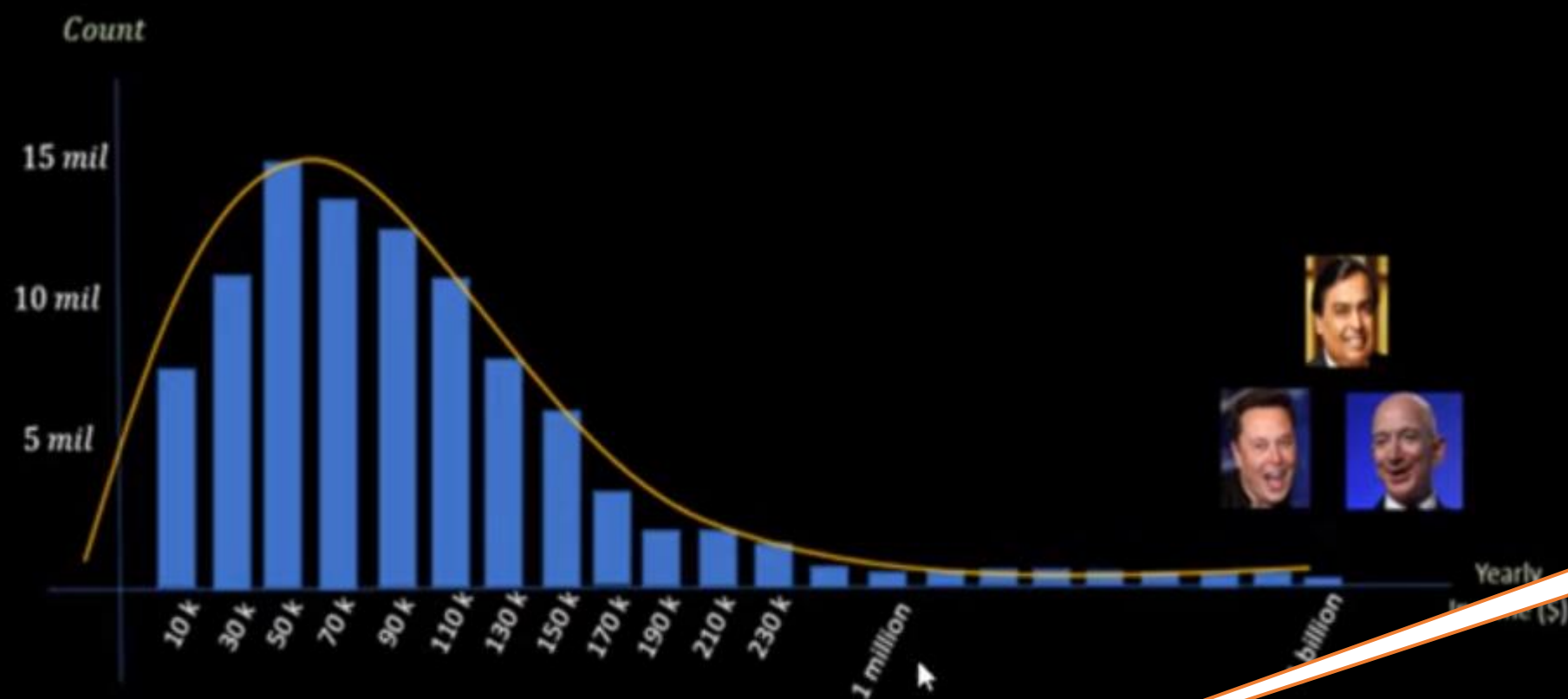
- Compress data in 0-1 range
- Applied to images
- Compress inliers in 0-0.5 when outliers are present
- https://scikit-learn.org/stable/auto_examples/preprocessing/plot_all_scaling.html



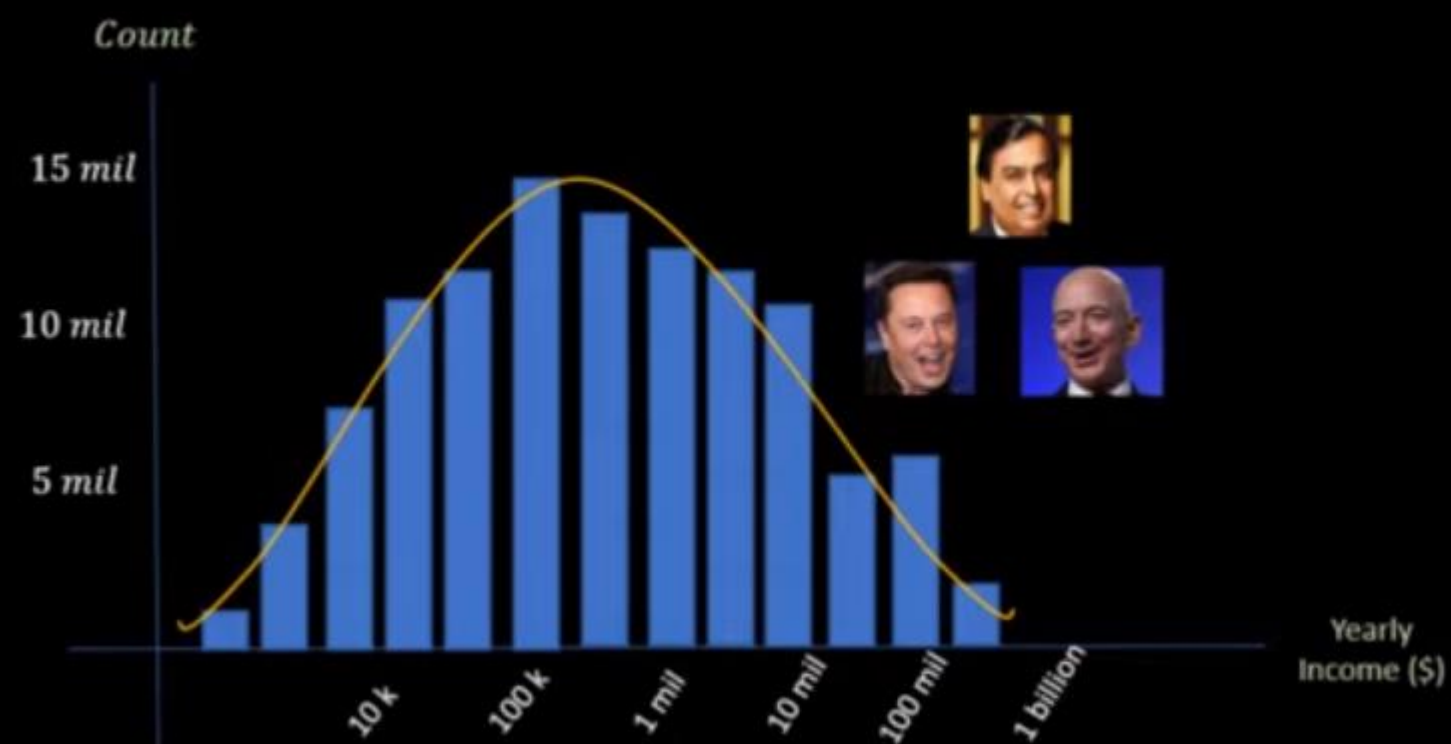
Distribution Transformations

Log Normal Distribution



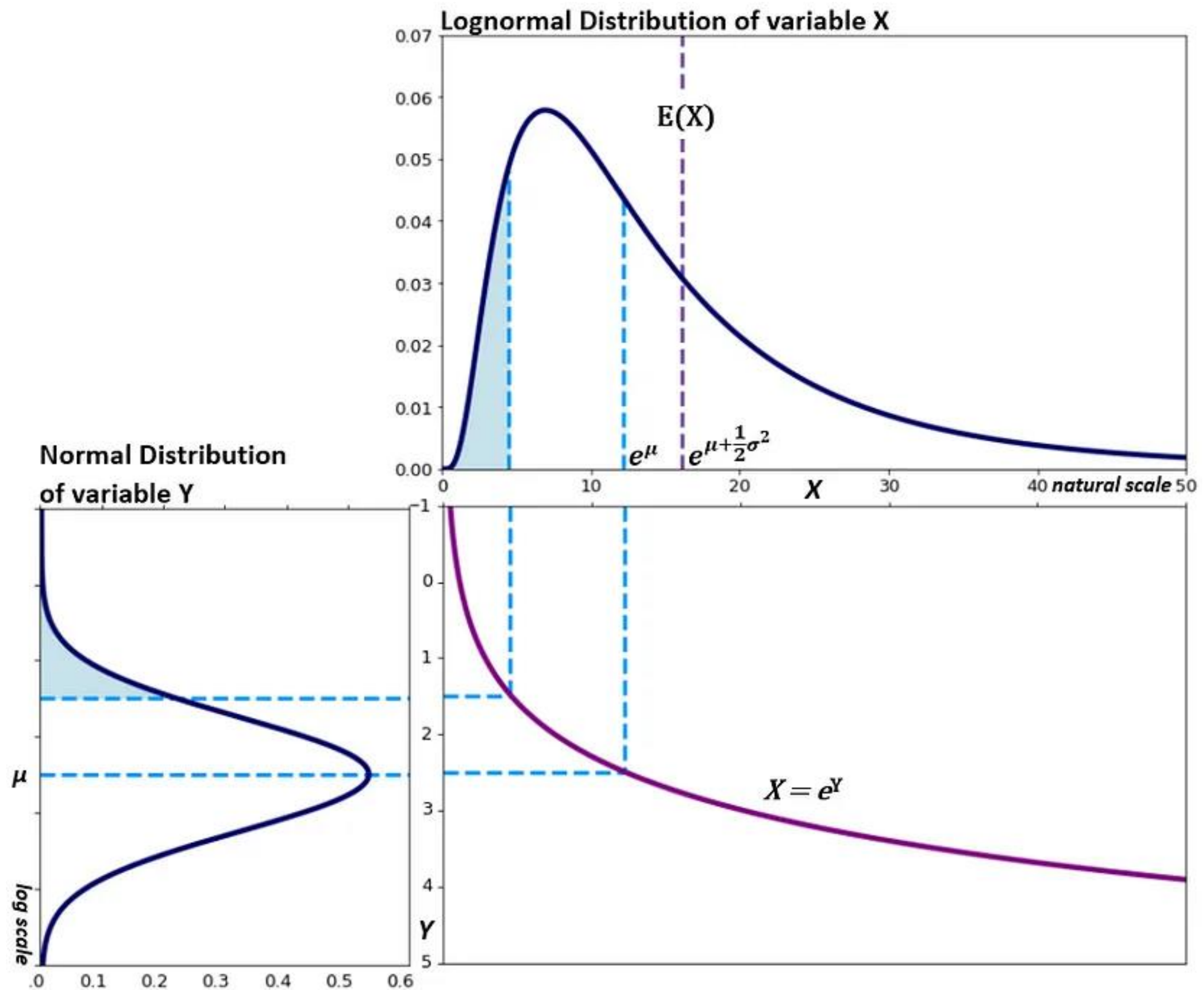


$\log(\text{income})$



$\log(1+\text{income})$

If you get a normal distribution by applying a log function to a dataset then dataset is log normally distributed



Generic Goal of Distribution Transforms

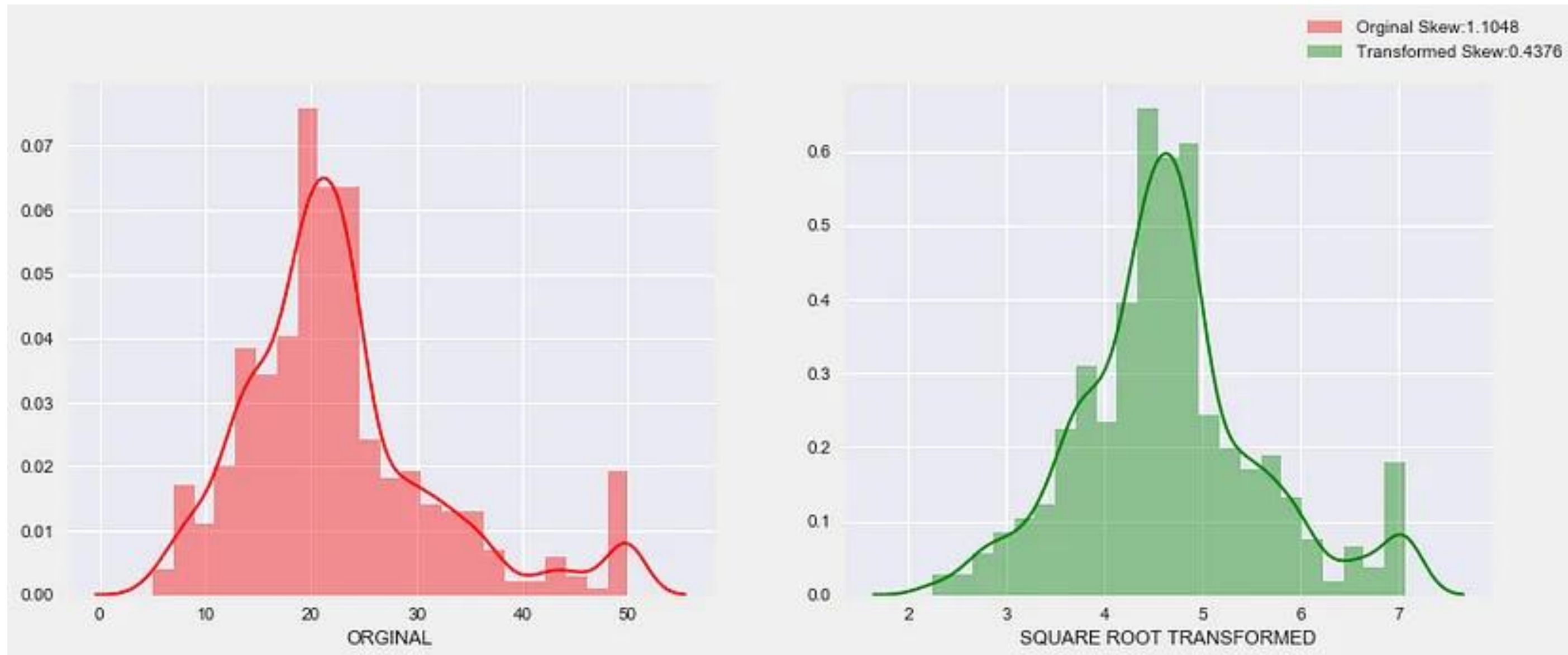
- Change the feature/target distribution to Gaussian
- Important for ML algorithms that have an underlying assumption that **target** variable is distributed normally
 - Linear Regression
 - Logistic Regression
- Theoretically **features** need not be Gaussian
- But numerically stable, faster convergence when features are Gaussian in a standardized range

Checking Gaussian-ness

- Visually with Seaborn `distplot()`
- Deviation from symmetric Gaussian
 - `pandas.skew()`: 0, <0 or >0
- Visually with QQ plot – Very reliable (`scipy.stats`)
- Normality tests
- What to apply when?
 - Heavily right skewed -> Log Transformer
 - Left skewed – Square Transformer
- Function Transformer in Sklearn

Gaussian Transformation

- Square Root Transformation
 - For removing slight right skewedness
 - Weaker than log transformation



Power Transforms

- Log Transform is a member of family: Power Transform
- Variance stabilizing transformations
- Remove skewness, Creates Gaussian
- Power Transforms $\phi(x, \lambda)$
 - Box Cox Transforms ($x > 0$)
 - Yeo-Johnson Transforms (for any x)
- Lambda is hyper param. Tune for feature transformation

Power Transforms

Box Cox

$$\phi(x, \lambda) = \begin{cases} \frac{x^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0 \\ \log x, & \text{if } \lambda = 0 \end{cases} \quad \lambda \in [-5, 5] \quad x > 0$$

Remember this

Yeo Johnson

$$\phi(x, \lambda) = \begin{cases} \frac{(x+1)^\lambda - 1}{\lambda}, & \text{if } x \geq 0 \text{ and } \lambda \neq 0 \\ \log(x + 1), & \text{if } x \geq 0 \text{ and } \lambda = 0 \\ -\frac{(-y+1)^{2-\lambda} - 1}{2-\lambda} & \text{if } x < 0 \text{ and } \lambda \neq 2 \\ -\log(-y + 1) & \text{if } x < 0 \text{ and } \lambda = 2 \end{cases}$$

Don't even try to remember this



QUESTIONS



Thank You!