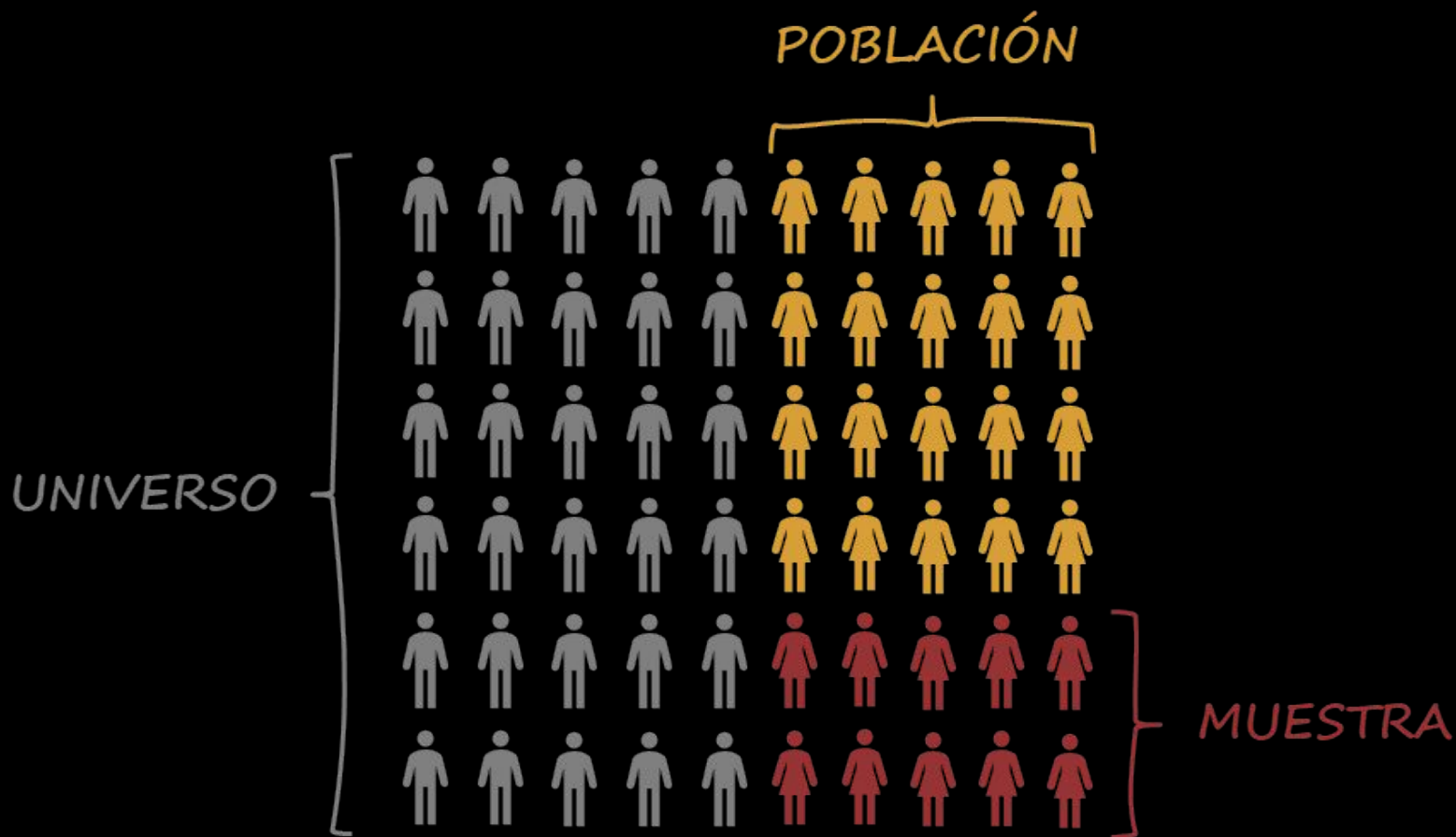


| | Outcome variable | | | | | | |
|----------------|------------------------------|-----------------------------------|-----------------------------|-----------------------------------|-----------------------------|--|--|
| Input Variable | | Nominal | Categorical (>2 Categories) | Ordinal | Quantitative Discrete | Quantitative Non-Normal | Quantitative Normal |
| | Nominal | χ^2 or Fisher's | χ^2 | χ^2 -trend or Mann - Whitney | Mann-Whitney | Mann-Whitney or log-rank ^a | Student's <i>t</i> test |
| | Categorical (2>categories) | χ^2 | χ^2 | Kruskal-Wallis ^b | Kruskal-Wallis ^b | Kruskal-Wallis ^b | Analysis of variance ^c |
| | Ordinal (Ordered categories) | χ^2 -trend or Mann - Whitney | * | Spearman rank | Spearman rank | Spearman rank | Spearman rank or linear regression ^d |
| | Quantitative Discrete | Logistic regression | * | * | Spearman rank | Spearman rank | Spearman rank or linear regression ^d |
| | Quantitative non-Normal | Logistic regression | * | * | * | Plot data and Pearson or Spearman rank | Plot data and Pearson or Spearman rank and linear regression |
| | Quantitative Normal | Logistic regression | * | * | * | Linear regression ^d | Pearson and linear regression |

| | Outcome variable | | | | | | |
|----------------|------------------------------|-----------------------------------|-----------------------------|-----------------------------------|-----------------------------|--|--|
| Input Variable | | Nominal | Categorical (>2 Categories) | Ordinal | Quantitative Discrete | Quantitative Non-Normal | Quantitative Normal |
| | Nominal | χ^2 or Fisher's | χ^2 | χ^2 -trend or Mann - Whitney | Mann-Whitney | Mann-Whitney or log-rank ^a | Student's <i>t</i> test |
| | Categorical (2>categories) | χ^2 | χ^2 | Kruskal-Wallis ^b | Kruskal-Wallis ^b | Kruskal-Wallis ^b | Analysis of variance ^c |
| | Ordinal (Ordered categories) | χ^2 -trend or Mann - Whitney | * | Spearman rank | Spearman rank | Spearman rank | Spearman rank or linear regression ^d |
| | Quantitative Discrete | Logistic regression | * | * | Spearman rank | Spearman rank | Spearman rank or linear regression ^d |
| | Quantitative non-Normal | Logistic regression | * | * | * | Plot data and Pearson or Spearman rank | Plot data and Pearson or Spearman rank and linear regression |
| | Quantitative Normal | Logistic regression | * | * | * | Linear regression ^d | Pearson and linear regression |

PRUEBAS PARAMÉTRICAS Y NO PARAMÉTRICAS

Ciencia de Datos



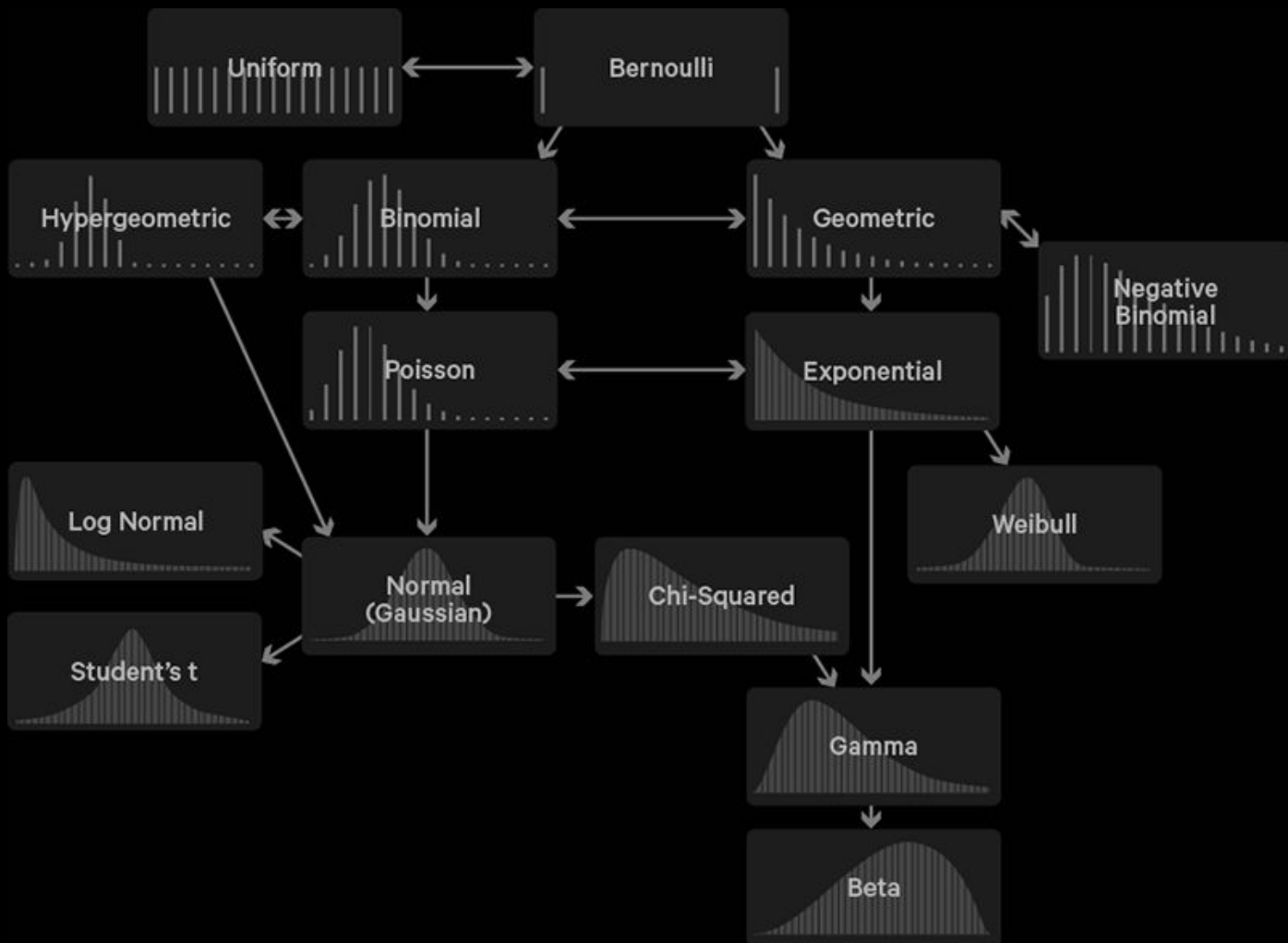
MUESTRA



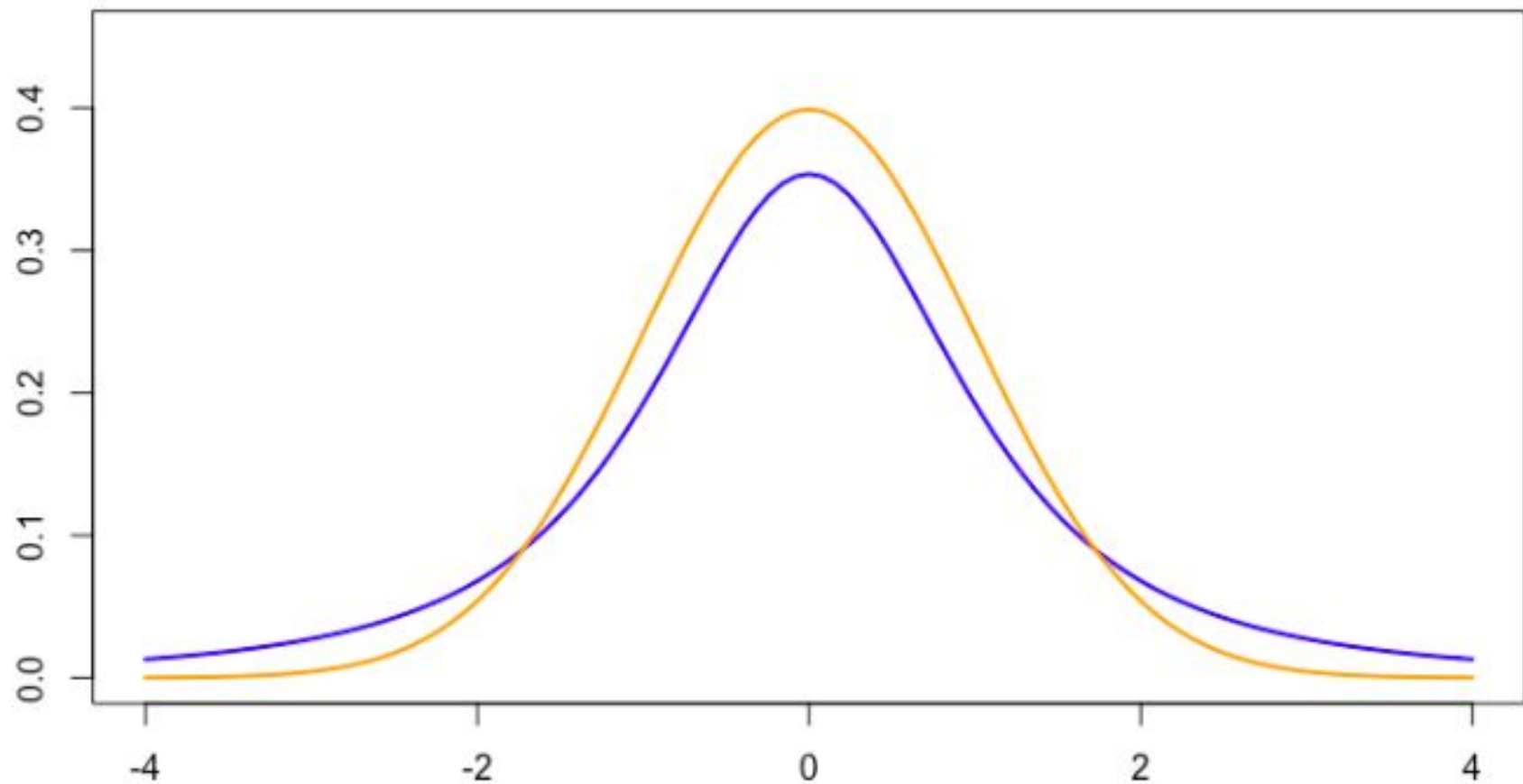
POBLACIÓN

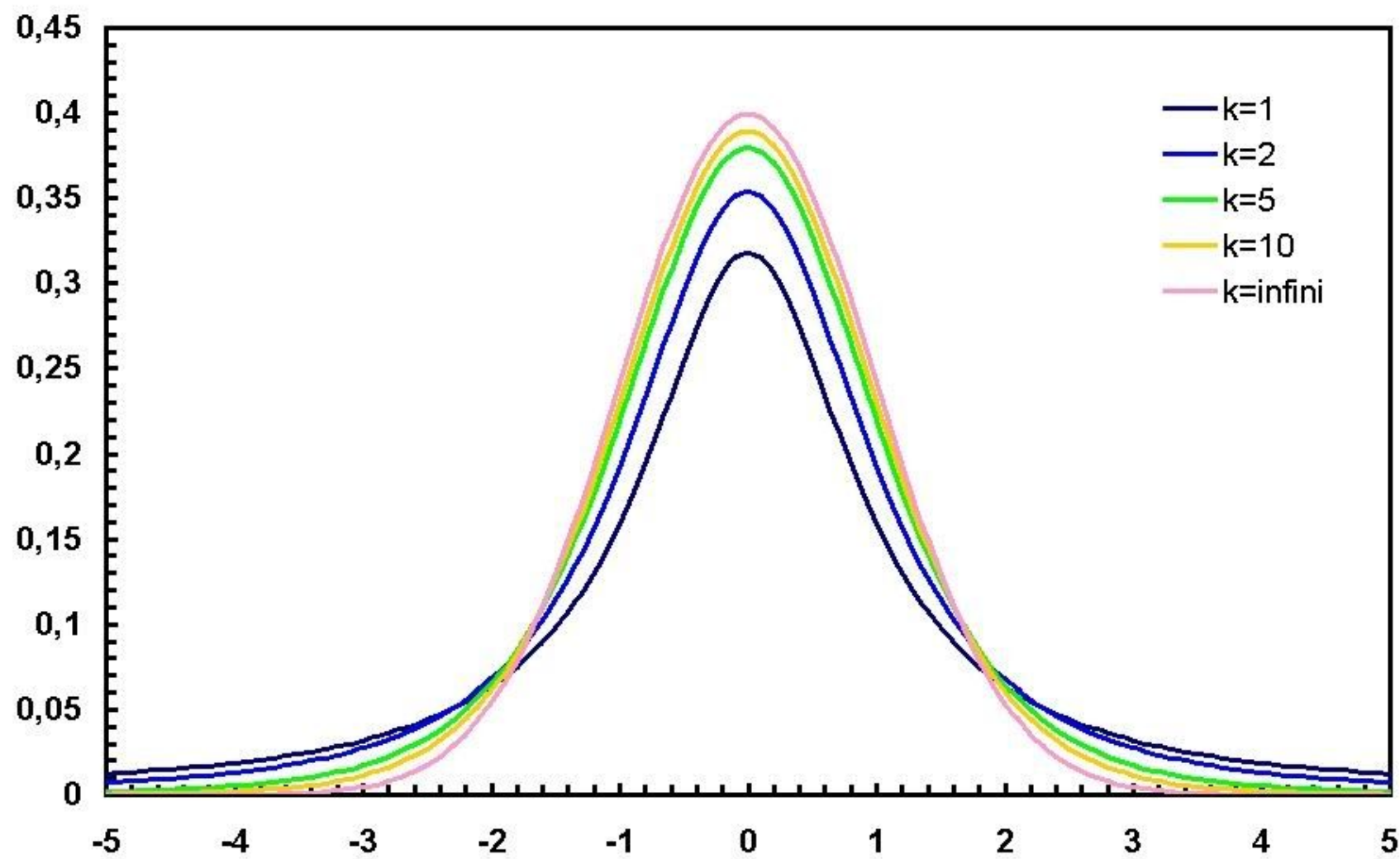


DISTRIBUCIONES DE FRECUENCIA



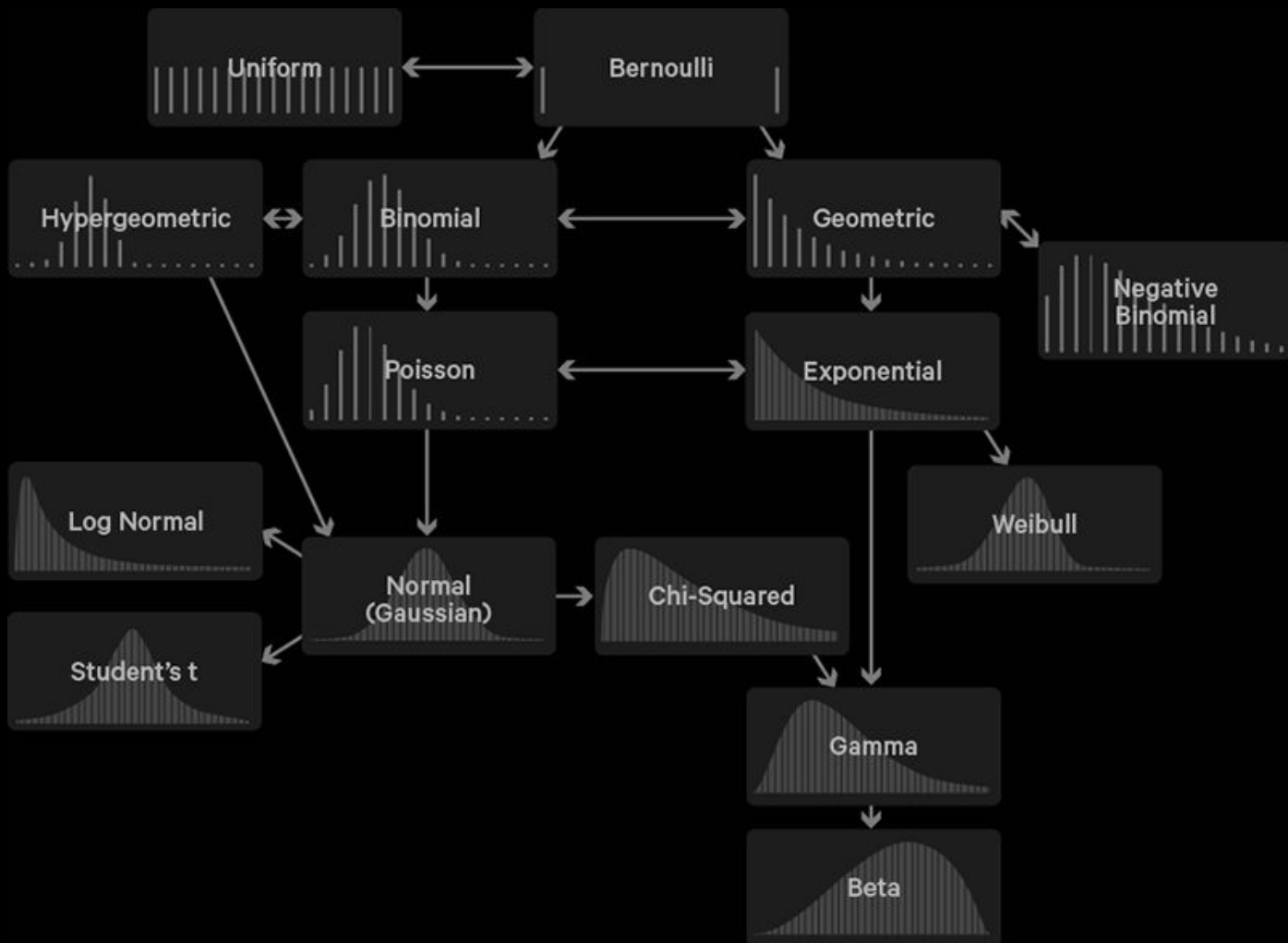
Distribución t de Student (azul) y distribución Normal estándar $N(0,1)$ (naranja)





TENEMOS UNA DIFICULTAD

DISTRIBUCIONES DE FRECUENCIA



OPCIÓN 1: HACEMOS SUPUESTOS

DISTRIBUCIÓN NORMAL

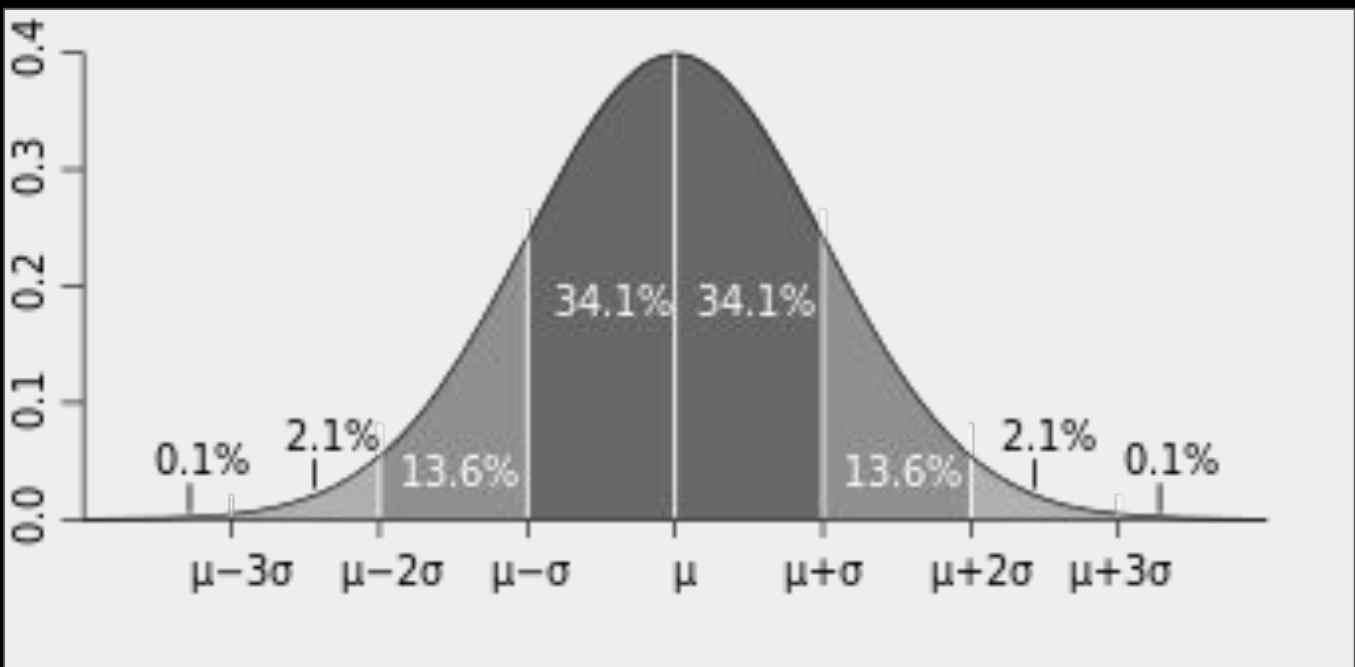


Tabla 1. Pruebas paramétricas

Prueba t de 1 muestra

Prueba t de 2 muestras

ANOVA de un solo factor

Correlación Pearson

Regresión lineal

DISTRIBUCIÓN MÁS FRECUENTES



Recomendables si

1. La muestra es grande
2. La distribución poblacional es conocida
3. Funcionan más rápidamente
4. Se usan variables de intervalo o razón

Ventajas:

Más potencia estadística

OPCIÓN 2: LIBRES DE DISTRIBUCIÓN

Tabla 2. Pruebas no paramétricas

Wilcoxon para 1 muestra

Prueba de Mann-Whitney

Kruskal-Wallis

Mediana de Wood

Prueba de Friedman

Recomendables si

1. La muestra es pequeña
2. La distribución poblacional es desconocida
3. Se usan variables categoriales u ordinales, aunque también hay para intervalo o razón

Análisis con variables numéricas:

| Análisis | Paramétrico | No paramétrico |
|---------------------------------------|---------------------------|-----------------------------|
| Describir un grupo | μ , σ^2 | Mediana, rango intercuartil |
| Comparar un grupo a un valor | T Student de una muestra | Prueba Wilcoxon |
| Comparar medias en 2 grupos | T Student de dos muestras | Mann-Whitney |
| Comparar medias en 2 grupos apareados | T Student apareada | Prueba Wilcoxon |
| Comparar medias en 3 o mas grupos | ANOVA | Kruskal-Wallis |
| Correlación entre dos variables | Pearson (lineal) | Spearman (monotónica) |

¿CÓMO DETERMINAR?

Crterios

1. Fuente de los datos
2. Información en literatura
3. Prueba de normalidad

```
1 from scipy import stats
2 rng = np.random.default_rng()
3 pts = 1000
4 a = rng.normal(0, 1, size=pts)
5 b = rng.normal(2, 1, size=pts)
6 x = np.concatenate((a, b))
7 k2, p = stats.normaltest(x)
8 alpha = 1e-3
9 print("p = {:.g}".format(p))
10 p = 8.4713e-19
11 if p < alpha: # null hypothesis: x comes from a normal distribution
12     print("The null hypothesis can be rejected")
13 else:
14     print("The null hypothesis cannot be rejected")
```

p = 1.21879e-12

The null hypothesis can be rejected

```
1 df = pd.read_csv("cars.csv")
2 odo = df[["odometer_value"]]
3 plt.hist(df[["odometer_value"]])
4 k2, p = stats.normaltest(odo)
5 alpha = 1e-3
6 print("p = ", p)
7 p = 8.4713e-19
8 if p < alpha: # null hypothesis: x comes from a normal distribution
9     print("The null hypothesis can be rejected")
10 else:
11     print("The null hypothesis cannot be rejected")
```

p = [0.]

The null hypothesis can be rejected

