

Essential Probability Distributions for Data Science

Contents

1	Bernoulli Distribution	2
2	Binomial Distribution	4
3	Poisson Distribution	4
4	Normal (Gaussian) Distribution	6
5	Uniform Distribution	6
6	Exponential Distribution	8
7	Gamma Distribution	8
8	Beta Distribution	9
9	Student's t-Distribution	10
10	Relationships and Quick Reference Table	11
11	Bernoulli and Binomial for Proportion Tests	13
11.1	Bernoulli/Binomial Overview	13
11.2	Example: Binomial Test in Python	13
12	Normal Distribution for z-Tests	14
12.1	Normal Overview	14
12.2	Example: z-Test for a Mean (Approximation)	14
13	t-Distribution for t-Tests	15
13.1	Student's t Overview	15
13.2	Example: One-Sample t-Test in Python	15
14	Chi-square Distribution for Goodness-of-Fit or Independence	16
14.1	Chi-square Overview	16
14.2	Example: Chi-square Goodness-of-Fit	16

15 F-Distribution for ANOVA	17
15.1 F-Distribution Overview	17
15.2 Example: One-Way ANOVA	17
16 Poisson Distribution for Rate Tests	18
16.1 Poisson Overview	18
16.2 Example: Test if Observed Count Matches a Poisson Rate	18
17 Summary of Distributions in Hypothesis Testing	19

Introduction

This document summarizes the most important probability distributions used in data science and statistics, with:

- **Formulas** (PMF or PDF)
- **Mean and Variance**
- **Key Relationships**
- **Example Python Code with Plots** using NumPy, Matplotlib/Seaborn

All code snippets assume you have the following libraries:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 sns.set_style("whitegrid")

```

1 Bernoulli Distribution

Use Case:

- Models a single trial with two outcomes (e.g., success/failure).
- Data science examples: binary classification outcomes, coin toss.

PMF

$$P(X = x) = p^x (1 - p)^{(1-x)}, \quad x \in \{0, 1\}, \quad 0 \leq p \leq 1$$

Mean: $\mu = p$

Variance: $\sigma^2 = p(1 - p)$

Python Code with Plot

```
1 p = 0.3
2 samples_bernoulli = np.random.binomial(n=1, p=p, size=1000)
3
4 sns.histplot(samples_bernoulli, discrete=True, stat='probability')
5 plt.title(f"Bernoulli (p={p}) ")
6 plt.xlabel("Value")
7 plt.ylabel("Probability")
8 plt.show()
```

2 Binomial Distribution

Use Case:

- Number of successes in n independent Bernoulli trials, each with success probability p .
- Common in A/B testing, success/failure counting.

PMF

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k = 0, 1, \dots, n$$

Mean: $\mu = np$

Variance: $\sigma^2 = np(1 - p)$

Relationships

- A **Bernoulli** distribution is a special case of **Binomial** when $n = 1$.
- For large n and small p (with $np = \lambda$ fixed), Binomial approximates Poisson(λ).
- For large n , Binomial approximates the **Normal** distribution with mean np and variance $np(1 - p)$.

Python Code with Plot

```
1 n, p = 10, 0.3
2 samples_binomial = np.random.binomial(n=n, p=p, size=1000)
3
4 sns.histplot(samples_binomial, discrete=True, stat='probability')
5 plt.title(f"Binomial (n={n}, p={p}) ")
6 plt.xlabel("Number_of_Successes")
7 plt.ylabel("Probability")
8 plt.show()
```

3 Poisson Distribution

Use Case:

- Models the number of events in a given interval with constant average rate λ .
- Example: number of website hits per minute.

PMF

$$P(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad k = 0, 1, 2, \dots$$

Mean: $\mu = \lambda$

Variance: $\sigma^2 = \lambda$

Relationships

- For large λ , $\text{Poisson}(\lambda)$ approximates $\text{Normal}(\lambda, \lambda)$.
- $\text{Binomial}(n, p)$ with large n and small p (where $np = \lambda$) approximates $\text{Poisson}(\lambda)$.

Python Code with Plot

```
1 lam = 5
2 samples_poisson = np.random.poisson(lam=lam, size=1000)
3
4 sns.histplot(samples_poisson, discrete=True, stat='probability')
5 plt.title(f"Poisson(  ={lam}) ")
6 plt.xlabel("Number_of_Events")
7 plt.ylabel("Probability")
8 plt.show()
```

4 Normal (Gaussian) Distribution

Use Case:

- Continuous distribution, the classic “bell curve.”
- By the Central Limit Theorem, sums/averages of large samples of i.i.d. variables tend to be Normal.

PDF

$$f(x) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right), \quad -\infty < x < \infty$$

Mean: μ

Variance: σ^2

Relationships

- Binomial(n, p) approximates Normal($np, np(1 - p)$) for large n .
- Poisson(λ) approximates Normal(λ, λ) for large λ .
- The t-distribution converges to Normal as $\nu \rightarrow \infty$.

Python Code with Plot

```
1 mu, sigma = 0, 1
2 samples_normal = np.random.normal(loc=mu, scale=sigma, size=1000)
3
4 sns.histplot(samples_normal, stat='density', kde=True, color='blue')
5 plt.title(f"Normal ( μ={mu}, σ={sigma} )")
6 plt.xlabel("Value")
7 plt.ylabel("Density")
8 plt.show()
```

5 Uniform Distribution

Use Case:

- Constant probability across an interval $[a, b]$.
- Often used for simulation or uninformative priors in Bayesian settings.

PDF

$$f(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$$

Mean: $\mu = \frac{a+b}{2}$

Variance: $\sigma^2 = \frac{(b-a)^2}{12}$

Python Code with Plot

```
1 a, b = 0, 1
2 samples_uniform = np.random.uniform(low=a, high=b, size=1000)
3
4 sns.histplot(samples_uniform, stat='density', kde=True, color='green')
5 plt.title(f"Uniform(a={a}, b={b}) ")
6 plt.xlabel("Value")
7 plt.ylabel("Density")
8 plt.show()
```

6 Exponential Distribution

Use Case:

- Models the time between events in a Poisson process (rate λ).
- Has the memoryless property.

PDF

$$f(x) = \lambda e^{-\lambda x}, \quad x \geq 0$$

Mean: $\mu = \frac{1}{\lambda}$

Variance: $\sigma^2 = \frac{1}{\lambda^2}$

Relationships

- The Γ distribution with shape parameter $k = 1$ is exactly Exponential(λ).
- The sum of k i.i.d. Exponential(λ) variables is Gamma(k, λ).

Python Code with Plot

```
1 lam_exp = 2 # Rate
2 samples_exponential = np.random.exponential(scale=1/lam_exp, size=1000)
3
4 sns.histplot(samples_exponential, stat='density', kde=True)
5 plt.title(f"Exponential( = {lam_exp} )")
6 plt.xlabel("Value")
7 plt.ylabel("Density")
8 plt.show()
```

7 Gamma Distribution

Use Case:

- Generalization of Exponential.
- Used for waiting times, Bayesian priors (e.g., for Poisson rate parameters).

PDF (Shape-Rate Parameterization)

$$f(x) = \frac{\lambda^k}{\Gamma(k)} x^{k-1} e^{-\lambda x}, \quad x \geq 0, \quad k > 0, \lambda > 0$$

Mean: $\mu = \frac{k}{\lambda}$

Variance: $\sigma^2 = \frac{k}{\lambda^2}$

Relationships

- $\text{Gamma}(k = 1, \lambda) = \text{Exponential}(\lambda)$.
- $\text{Chi-square}(\nu) = \text{Gamma}(\frac{\nu}{2}, \frac{1}{2})$.

Python Code with Plot

```
1 k, lam_gamma = 2.0, 1.0
2 samples_gamma = np.random.gamma(shape=k, scale=1/lam_gamma, size=1000)
3
4 sns.histplot(samples_gamma, stat='density', kde=True)
5 plt.title(f"Gamma (k={k}, λ = {lam_gamma}) ")
6 plt.xlabel("Value")
7 plt.ylabel("Density")
8 plt.show()
```

8 Beta Distribution

Use Case:

- Models values strictly in $[0, 1]$.
- Often used as a prior over probabilities (e.g., for Bernoulli/Binomial).

PDF

$$f(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}, \quad 0 \leq x \leq 1, \quad \alpha > 0, \beta > 0$$

where $B(\alpha, \beta)$ is the Beta function.

Mean: $\mu = \frac{\alpha}{\alpha + \beta}$

Variance:

$$\sigma^2 = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

Python Code with Plot

```
1 alpha, beta_ = 2, 5
2 samples_beta = np.random.beta(alpha, beta_, size=1000)
3
4 sns.histplot(samples_beta, stat='density', kde=True)
5 plt.title(f"Beta (α={alpha}, β={beta_}) ")
6 plt.xlabel("Value")
7 plt.ylabel("Density")
8 plt.show()
```

9 Student's t-Distribution

Use Case:

- Used when population standard deviation is unknown, especially with small sample sizes.
- Common in hypothesis testing (t-tests) and confidence intervals.

PDF (for ν degrees of freedom)

$$f(t) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi} \Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

Mean: 0 (for $\nu > 1$)

Variance: $\frac{\nu}{\nu-2}$ (for $\nu > 2$)

Relationship

- As $\nu \rightarrow \infty$, the t-distribution approaches Normal(0, 1).

Python Code with Plot

```
1 df = 10 # degrees of freedom
2 samples_t = np.random.standard_t(df, size=1000)
3
4 sns.histplot(samples_t, stat='density', kde=True)
5 plt.title(f"t-Distribution(df={df}) ")
6 plt.xlabel("Value")
7 plt.ylabel("Density")
8 plt.show()
```

10 Relationships and Quick Reference Table

Key Relationships

1. **Bernoulli** is a special case of **Binomial** ($n = 1$).
2. **Binomial**(n, p) \approx **Poisson**($\lambda = np$) if n is large and p is small.
3. **Poisson**(λ) \approx **Normal**($\mu = \lambda, \sigma^2 = \lambda$) if λ is large.
4. **Binomial**(n, p) \approx **Normal**($np, np(1 - p)$) for large n .
5. **Gamma**($k = 1, \lambda$) = **Exponential**(λ).
6. **Chi-square**(ν) = **Gamma**($\frac{\nu}{2}, \frac{1}{2}$).
7. **t-Distribution** \rightarrow Normal as $\nu \rightarrow \infty$.

Quick Reference Table

Distribution	Discrete / Continuous	Support	Parameters	Mean	Variance
Bernoulli	Discrete	$\{0, 1\}$	p	p	$p(1 - p)$
Binomial	Discrete	$\{0, 1, \dots, n\}$	n, p	np	$np(1 - p)$
Poisson	Discrete	$\{0, 1, 2, \dots\}$	λ	λ	λ
Uniform	Continuous	$[a, b]$	a, b	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
Normal	Continuous	$(-\infty, \infty)$	μ, σ^2	μ	σ^2
Exponential	Continuous	$[0, \infty)$	λ	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
Gamma	Continuous	$[0, \infty)$	k, λ	$\frac{k}{\lambda}$	$\frac{k}{\lambda^2}$
Beta	Continuous	$[0, 1]$	α, β	$\frac{\alpha}{\alpha + \beta}$	$\frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$
t-Distribution	Continuous	$(-\infty, \infty)$	ν (dof)	0 (for $\nu > 1$)	$\frac{\nu}{\nu - 2}$ (for $\nu > 2$)

Final Remarks

- These distributions are **central** to data science and statistics.
- Understanding their **properties** and **relationships** helps in choosing appropriate models.
- Python's `numpy` and `seaborn` libraries make it convenient to generate random samples and visualize them.

Introduction

In hypothesis testing, we use **probability distributions** to:

- Derive test statistics (e.g., z-statistic, t-statistic, chi-square statistic, etc.).
- Compute p -values, confidence intervals, and critical regions.

This document highlights several common distributions and shows how they underpin different hypothesis tests, with **Python code examples** (using `numpy` and `scipy.stats`) to illustrate practical usage.

All code snippets assume:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from scipy import stats
5
6 sns.set_style("whitegrid")
```

11 Bernoulli and Binomial for Proportion Tests

11.1 Bernoulli/Binomial Overview

- **Bernoulli:** Single trial with two outcomes (0 or 1).
- **Binomial:** Number of successes in n Bernoulli trials with success probability p .

Hypothesis Testing Context:

- **One-Proportion Test (Binomial Test):**

$$H_0 : p = p_0 \quad \text{vs.} \quad H_a : p \neq p_0 \quad (\text{or one-sided})$$

- We often use a binomial distribution or normal approximation to test whether the true probability p differs from a hypothesized p_0 .

11.2 Example: Binomial Test in Python

```
1 # Suppose we observe 40 successes out of 100 trials.
2 # We want to test if p = 0.5.
3
4 observed_successes = 40
5 n = 100
6 p0 = 0.5
7
8 # Using scipy.stats.binom_test (deprecated in newer SciPy versions);
9 # for newer versions, use proportions_ztest from statsmodels.
10 p_value = stats.binom_test(observed_successes, n=n, p=p0, alternative='two-
    sided')
11 print("Binomial_test_p-value:", p_value)
12
13 # If p_value < alpha (e.g., 0.05), we reject H0 that p = 0.5.
```

Note: In newer versions of SciPy, `binom_test` is deprecated. You can use `statsmodels.stats.proportions` for a z-based approximation.

12 Normal Distribution for z-Tests

12.1 Normal Overview

- **Normal (Gaussian)** distribution is used when sample sizes are large (by the Central Limit Theorem) or when population standard deviation is known.

Hypothesis Testing Context:

- **z-Test for Means:**

$$H_0 : \mu = \mu_0 \quad \text{vs.} \quad H_a : \mu \neq \mu_0$$

- Typically used when population variance (σ^2) is known or n is sufficiently large.

12.2 Example: z-Test for a Mean (Approximation)

In Python, pure z-tests are less common than t-tests, but we can approximate with:

```
1 # Synthetic data: sample from a Normal(  =0,  =1)
2 np.random.seed(42)
3 data = np.random.normal(loc=0, scale=1, size=100)
4
5 # Suppose we want to test if mean = 0 (H0:  =0).
6 # We'll approximate with a z-test if we assume  =1 known.
7
8 sample_mean = np.mean(data)
9 n = len(data)
10 sigma = 1.0 # known population std dev
11 z_stat = (sample_mean - 0) / (sigma / np.sqrt(n))
12
13 # Two-sided p-value using Normal
14 p_value = 2 * (1 - stats.norm.cdf(abs(z_stat)))
15 print("Z-statistic:", z_stat)
16 print("p-value:", p_value)
```

If $p_value \leq \alpha$ (common choice 0.05), we reject H_0 .

13 t-Distribution for t-Tests

13.1 Student's t Overview

- Used when the population standard deviation is unknown and estimated by the sample.
- Common in small-sample scenarios or standard parametric tests in practice.

Hypothesis Testing Context:

- One-sample t-Test:

$$H_0 : \mu = \mu_0$$

- Two-sample t-Test (independent or paired).

13.2 Example: One-Sample t-Test in Python

```
1 # Suppose we have data from an unknown distribution,
2 # and we want to test if the mean is 5.
3
4 data = np.array([4.9, 5.1, 5.3, 4.8, 5.2, 5.4, 5.0])
5 t_stat, p_val = stats.ttest_1samp(data, popmean=5)
6 print("One-sample_t-test_statistic:", t_stat)
7 print("p-value:", p_val)
8
9 # If p_val < alpha, reject H0: mean is 5.
```

14 Chi-square Distribution for Goodness-of-Fit or Independence

14.1 Chi-square Overview

- χ^2 distribution arises from summing squares of Normal(0,1) variables.
- Often used in **goodness-of-fit, independence tests** (contingency tables).

Hypothesis Testing Context:

- **Chi-square Goodness-of-Fit:**

H_0 : The data follow a specified distribution.

- **Chi-square Test of Independence** in contingency tables.

14.2 Example: Chi-square Goodness-of-Fit

```
1 # Suppose we observe frequencies in 4 categories:
2 observed = np.array([18, 22, 15, 25])
3
4 # Suppose expected probabilities for these categories
5 # are [0.25, 0.25, 0.25, 0.25].
6 n = np.sum(observed)
7 expected = n * np.array([0.25, 0.25, 0.25, 0.25])
8
9 chi_stat, p_value = stats.chisquare(f_obs=observed, f_exp=expected)
10 print("Chi-square_statistic:", chi_stat)
11 print("p-value:", p_value)
12
13 # If p_value < alpha, reject H0 that distribution matches [0.25, 0.25, 0.25,
    0.25].
```


15 F-Distribution for ANOVA

15.1 F-Distribution Overview

- Ratio of two scaled chi-square distributions is an F-distribution.
- Used in **ANOVA** (analysis of variance) to compare means of multiple groups.

Hypothesis Testing Context (One-Way ANOVA):

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_k \quad \text{vs.} \quad H_a : \text{at least one mean is different.}$$

15.2 Example: One-Way ANOVA

```
1 # Suppose we have 3 groups of data:
2 group1 = np.random.normal(loc=5, scale=1, size=30)
3 group2 = np.random.normal(loc=5.5, scale=1, size=30)
4 group3 = np.random.normal(loc=6, scale=1, size=30)
5
6 f_stat, p_val = stats.f_oneway(group1, group2, group3)
7 print("F-statistic:", f_stat)
8 print("p-value:", p_val)
9
10 # If p_val < alpha, reject H0 that all group means are equal.
```

16 Poisson Distribution for Rate Tests

16.1 Poisson Overview

- Models the number of events in a time/space interval with rate λ .
- Hypothesis tests often compare observed vs. expected counts at a certain rate.

Hypothesis Testing Context:

- Test whether the observed count matches a $\text{Poisson}(\lambda_0)$ with a known or hypothesized rate λ_0 .

16.2 Example: Test if Observed Count Matches a Poisson Rate

While not as standard as z/t-tests, we can build a likelihood-based test or use a simpler approximation:

```
1 # Observed count in a fixed interval
2 observed_count = 12
3 # Hypothesized rate
4 lambda_0 = 10
5
6 # Under H0, X ~ Poisson(lambda_0).
7 # Probability of observing something "as extreme or more extreme"
8 # can be computed as p-value.
9 # For a two-sided test, we might do:
10
11 p_lower = stats.poisson.cdf(observed_count, mu=lambda_0)
12 p_upper = 1 - stats.poisson.cdf(observed_count - 1, mu=lambda_0)
13 p_value_two_sided = 2 * min(p_lower, p_upper)
14
15 print("Poisson_two-sided_p-value:", p_value_two_sided)
```

If $\text{p_value_two_sided} < \alpha$, we reject $H_0 : \lambda = \lambda_0$.

17 Summary of Distributions in Hypothesis Testing

Which Distribution for Which Test?

- **Binomial** (*or Normal approx.*): Testing proportions (e.g., yes/no outcomes).
- **Normal (z-Test)**: Large-sample mean tests or known population variance.
- **t-Distribution (t-Test)**: Small-sample mean tests, unknown variance.
- **Chi-square**: Goodness-of-fit, independence tests in contingency tables.
- **F-Distribution (ANOVA)**: Comparing means of 2 or more groups.
- **Poisson**: Rate-based tests (counts over time/space).

General Steps in Hypothesis Testing

1. **Formulate** H_0 and H_a (null and alternative hypotheses).
2. **Choose** appropriate test statistic and distribution (z, t, χ^2 , F, etc.).
3. **Compute** the test statistic from sample data.
4. **Obtain** the p -value or compare to a critical value.
5. **Decide** whether to reject or fail to reject H_0 based on α level.

Key Insight: Each test statistic is tied to a **reference distribution** (Normal, t, Chi-square, F, etc.) that tells us how extreme our sample results are if H_0 were true.

Conclusion

These distributions and tests form the backbone of classical hypothesis testing. Mastering their assumptions and usage is crucial for drawing valid inferences from data.