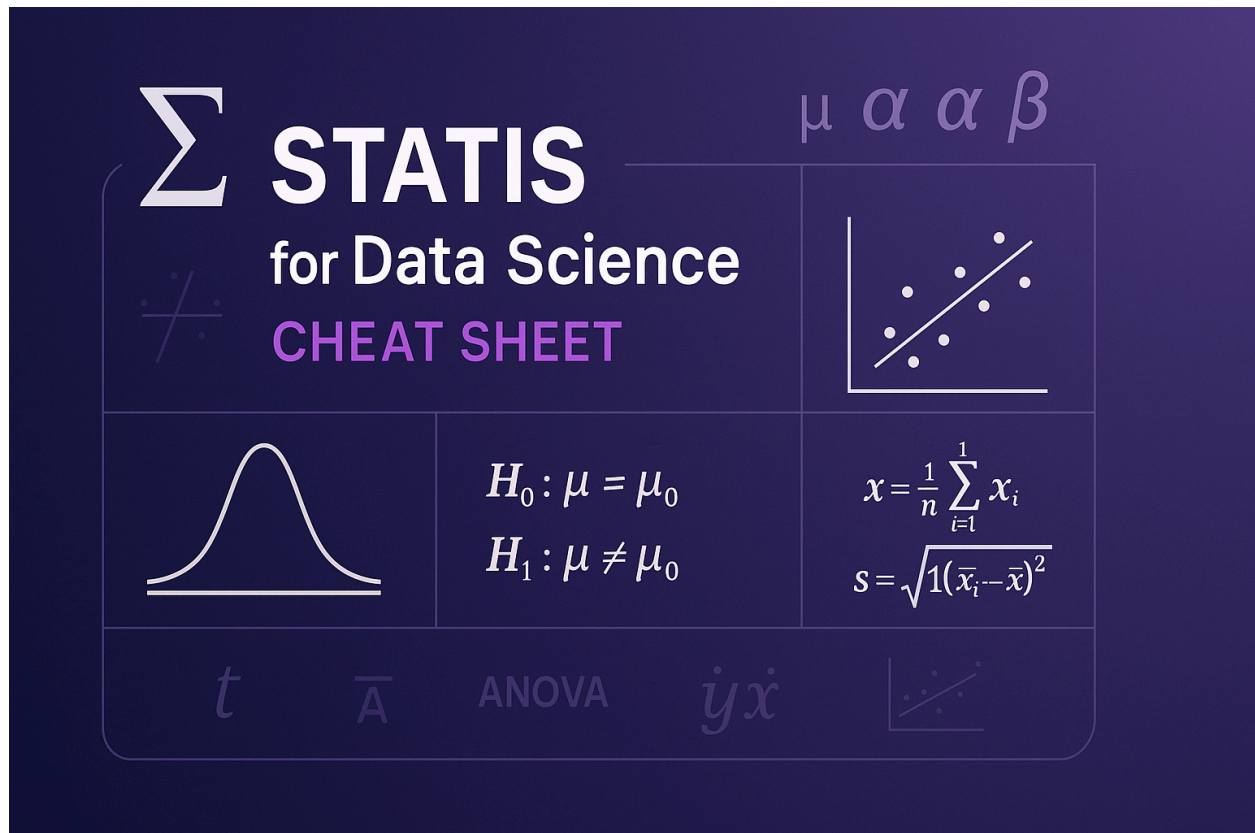


Statistics for Data Science - Comprehensive Cheat Sheet



Descriptive Statistics

Measures of Central Tendency

Mean (Average)

$$\mu = \Sigma x / n$$

Population mean

$$\bar{x} = \Sigma x / n$$

Sample mean

Median

- Middle value when data is ordered
- Less sensitive to outliers than mean

- Better for skewed distributions

Mode

- Most frequently occurring value
- Can have multiple modes (bimodal, multimodal)
- Useful for categorical data

Relationship

- Normal distribution: Mean = Median = Mode
- Right skew: Mean > Median > Mode
- Left skew: Mode > Median > Mean

Measures of Variability

Range

Range = Max - Min

Variance

$\sigma^2 = \sum (x - \mu)^2 / N$ # Population variance

$s^2 = \sum (x - \bar{x})^2 / (n-1)$ # Sample variance (Bessel's correction)

Standard Deviation

$\sigma = \sqrt{\sigma^2}$ # Population std dev

$s = \sqrt{s^2}$ # Sample std dev

Coefficient of Variation

$CV = (\sigma / \mu) \times 100\%$ # Relative variability

Interquartile Range (IQR)

$IQR = Q3 - Q1$ # Middle 50% of data

Outliers: < $Q1 - 1.5 \times IQR$ or > $Q3 + 1.5 \times IQR$

Measures of Shape

Skewness

- Measures asymmetry of distribution
- Skewness = 0: Symmetric
- Skewness > 0: Right-tailed (positive skew)
- Skewness < 0: Left-tailed (negative skew)

Kurtosis

- Measures tail heaviness
- Kurtosis = 3: Normal distribution (mesokurtic)
- Kurtosis > 3: Heavy tails (leptokurtic)
- Kurtosis < 3: Light tails (platykurtic)

Percentiles & Quartiles

Percentiles

P_k = Value below which k% of data falls

Quartiles

Q1 = 25th percentile (First quartile)

Q2 = 50th percentile (Median)

Q3 = 75th percentile (Third quartile)

Five-number summary

Min, Q1, Median, Q3, Max

Probability Distributions

Discrete Distributions

Bernoulli Distribution

Single trial with two outcomes (success/failure)

$P(X = 1) = p$ # Probability of success

$P(X = 0) = 1 - p$ # Probability of failure

Mean = p

Variance = $p(1-p)$

Binomial Distribution

n independent Bernoulli trials

$P(X = k) = C(n, k) \times p^k \times (1-p)^{(n-k)}$

Mean = np

Variance = $np(1-p)$

Standard Deviation = $\sqrt{np(1-p)}$

Use when:

- Fixed number of trials
- Each trial is independent
- Constant probability of success
- Two possible outcomes

Poisson Distribution

Number of events in fixed interval

$P(X = k) = (\lambda^k \times e^{(-\lambda)}) / k!$

Mean = λ

Variance = λ

Standard Deviation = $\sqrt{\lambda}$

Use when:

- Events occur independently
- Average rate is constant
- Rare events over time/space

Continuous Distributions

Normal Distribution

Bell-shaped, symmetric distribution

$$X \sim N(\mu, \sigma^2)$$

Standard Normal Distribution

$$Z = (X - \mu) / \sigma \quad \# \text{ Z-score transformation}$$

$$Z \sim N(0, 1)$$

Properties:

- 68% within 1 standard deviation
- 95% within 2 standard deviations
- 99.7% within 3 standard deviations

Central Limit Theorem

Sample means approach normal distribution as n increases

Student's t-Distribution

Similar to normal but heavier tails

Used when:

- Small sample sizes ($n < 30$)
- Population standard deviation unknown
- Degrees of freedom = $n - 1$

As df increases, approaches normal distribution

Chi-Square Distribution

Right-skewed distribution

Used for:

- Goodness of fit tests
- Test of independence
- Variance testing

$\chi^2 = \sum ((\text{Observed} - \text{Expected})^2 / \text{Expected})$
 $df = (\text{rows} - 1) \times (\text{columns} - 1)$ # For contingency tables

F-Distribution

Ratio of two chi-square distributions
Used for:
- ANOVA (Analysis of Variance)
- Comparing variances
- Regression analysis

$F = (s_1^2 / \sigma_1^2) / (s_2^2 / \sigma_2^2)$
 df_1 = numerator degrees of freedom
 df_2 = denominator degrees of freedom

Sampling and Estimation

Sampling Methods

Simple Random Sampling
- Every element has equal probability
- Use random number generation

Stratified Sampling
- Divide population into strata
- Sample from each stratum
- Ensures representation

Cluster Sampling
- Divide into clusters
- Randomly select clusters
- Sample all elements in selected clusters

Systematic Sampling

- Select every kth element
- $k = N/n$ (population size / sample size)

Central Limit Theorem

For sample means:

$$\mu_{\bar{x}} = \mu \quad \# \text{ Mean of sample means}$$

$$\sigma_{\bar{x}} = \sigma / \sqrt{n} \quad \# \text{ Standard error of mean}$$

Conditions:

- Sample size $n \geq 30$ (or population is normal)
- Samples are independent
- 10% condition: $n < 0.1N$ (for sampling without replacement)

Confidence Intervals

For population mean (σ known)

$$CI = \bar{x} \pm z_{(\alpha/2)} \times (\sigma/\sqrt{n})$$

For population mean (σ unknown)

$$CI = \bar{x} \pm t_{(\alpha/2, df)} \times (s/\sqrt{n})$$

For population proportion

$$CI = \hat{p} \pm z_{(\alpha/2)} \times \sqrt{(\hat{p}(1-\hat{p}))/n}$$

Margin of Error

$$ME = \text{Critical Value} \times \text{Standard Error}$$

Common Confidence Levels:

$$90\%: z = 1.645$$

$$95\%: z = 1.96$$

$$99\%: z = 2.576$$

Hypothesis Testing

Hypothesis Testing Framework

Step 1: State hypotheses

H_0 : Null hypothesis (status quo)

H_1 : Alternative hypothesis (what we want to prove)

Step 2: Choose significance level

$\alpha = 0.05$ (common choice)

Step 3: Calculate test statistic

Step 4: Find p-value or critical value

Step 5: Make decision

Step 6: State conclusion in context

Types of Errors

Type I Error (α)

- Reject true null hypothesis
- False positive
- $P(\text{Type I Error}) = \alpha$

Type II Error (β)

- Fail to reject false null hypothesis
- False negative
- $P(\text{Type II Error}) = \beta$

Power of Test

Power = $1 - \beta$

- Probability of correctly rejecting false H_0

One-Sample Tests

One-Sample t-Test

Test population mean when σ unknown

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0 \text{ (two-tailed)}$$

$$H_1: \mu > \mu_0 \text{ (right-tailed)}$$

$$H_1: \mu < \mu_0 \text{ (left-tailed)}$$

$$t = (\bar{x} - \mu_0) / (s/\sqrt{n})$$

$$df = n - 1$$

Assumptions:

- Random sample
- Normal distribution or $n \geq 30$
- Independent observations

One-Sample Proportion Test

Test population proportion

$$H_0: p = p_0$$

$$H_1: p \neq p_0$$

$$z = (\hat{p} - p_0) / \sqrt{(p_0(1-p_0)/n)}$$

Assumptions:

- Random sample
- $np_0 \geq 10$ and $n(1-p_0) \geq 10$
- Independent observations

Two-Sample Tests

Two-Sample t-Test

Independent samples (equal variances)

$$t = (\bar{x}_1 - \bar{x}_2) / (s_p \times \sqrt{1/n_1 + 1/n_2})$$

```
s_p = sqrt(((n1-1)s1^2 + (n2-1)s2^2) / (n1+n2-2)) # Pooled std dev  
df = n1 + n2 - 2
```

```
# Independent samples (unequal variances - Welch's t-test)  
t = (x̄1 - x̄2) / sqrt(s1^2/n1 + s2^2/n2)
```

```
# Paired samples  
t = (d̄ - μ_d) / (s_d/sqrt(n))  
df = n - 1
```

Two-Sample Proportion Test

```
# Compare two proportions  
z = (p̂1 - p̂2) / sqrt(p̂(1-p̂)(1/n1 + 1/n2))  
  
p̂ = (x1 + x2) / (n1 + n2) # Pooled proportion  
  
# Assumptions:  
- Independent samples  
- Large sample sizes
```

ANOVA (Analysis of Variance)

```
# Compare means of 3+ groups  
H0: μ1 = μ2 = ... = μk  
H1: At least one mean is different
```

```
F = MSB / MSW
```

```
MSB = SSB / (k-1)      # Mean Square Between  
MSW = SSW / (N-k)      # Mean Square Within
```

```
# Where:  
k = number of groups  
N = total sample size
```

Post-hoc tests (if F significant):

- Tukey's HSD
- Bonferroni correction
- Scheffé test

Chi-Square Tests

Goodness of Fit Test

Test if data follows expected distribution

H_0 : Data follows specified distribution

H_1 : Data does not follow specified distribution

$$\chi^2 = \sum ((O_i - E_i)^2 / E_i)$$

df = categories - 1

Assumptions:

- Expected frequency ≥ 5 for each category
- Independent observations

Test of Independence

Test relationship between two categorical variables

H_0 : Variables are independent

H_1 : Variables are dependent

$$\chi^2 = \sum ((O_{ij} - E_{ij})^2 / E_{ij})$$

df = (rows - 1) \times (columns - 1)

$E_{ij} = (\text{row total} \times \text{column total}) / \text{grand total}$

Effect size: Cramér's V

$$V = \sqrt{\chi^2 / (n \times \min(r-1, c-1))}$$

Correlation and Regression

Correlation Analysis

Pearson Correlation Coefficient

$$r = \frac{\sum((x_i - \bar{x})(y_i - \bar{y}))}{\sqrt{(\sum(x_i - \bar{x})^2 \times \sum(y_i - \bar{y})^2)}}$$

Interpretation:

$r = 1$: Perfect positive correlation

$r = 0$: No linear correlation

$r = -1$: Perfect negative correlation

$|r| \geq 0.7$: Strong correlation

$0.3 \leq |r| < 0.7$: Moderate correlation

$|r| < 0.3$: Weak correlation

Spearman Rank Correlation

- Non-parametric alternative

- Based on ranks, not raw values

- Detects monotonic relationships

Simple Linear Regression

Model: $y = \beta_0 + \beta_1 x + \varepsilon$

Least Squares Estimates:

$$\beta_1 = \frac{\sum((x_i - \bar{x})(y_i - \bar{y}))}{\sum(x_i - \bar{x})^2} \quad \# \text{ Slope}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x} \quad \# \text{ Intercept}$$

Coefficient of Determination

$$R^2 = SSR / SST = 1 - SSE / SST$$

$$SST = \sum(y_i - \bar{y})^2 \quad \# \text{ Total Sum of Squares}$$

$$SSR = \sum(\hat{y}_i - \bar{y})^2 \quad \# \text{ Regression Sum of Squares}$$

$$SSE = \sum(y_i - \hat{y}_i)^2 \quad \# \text{ Error Sum of Squares}$$

Standard Error of Estimate

$$s_e = \sqrt{(SSE / (n-2))}$$

Assumptions:

- Linearity
- Independence
- Homoscedasticity (constant variance)
- Normality of residuals

Multiple Linear Regression

Model: $y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_kx_k + \varepsilon$

Adjusted R^2

$$R^2_{adj} = 1 - ((1-R^2)(n-1)) / (n-k-1)$$

F-test for overall significance

$$F = (R^2/k) / ((1-R^2)/(n-k-1))$$

t-test for individual coefficients

$$t = \beta_j / SE(\beta_j)$$

Multicollinearity detection:

- Variance Inflation Factor (VIF)
- $VIF > 10$ indicates multicollinearity

Non-Parametric Tests

When to Use Non-Parametric Tests

- Data is not normally distributed
- Ordinal data
- Small sample sizes

- Presence of outliers
- Violated assumptions of parametric tests

Common Non-Parametric Tests

Mann-Whitney U Test (Wilcoxon Rank-Sum)

- Alternative to two-sample t-test
- Compares medians of two independent groups

Wilcoxon Signed-Rank Test

- Alternative to paired t-test
- Compares medians of paired samples

Kruskal-Wallis Test

- Alternative to one-way ANOVA
- Compares medians of 3+ independent groups

Friedman Test

- Alternative to repeated measures ANOVA
- Compares medians of 3+ related groups

Sign Test

- Tests median of single population
- Uses only direction of differences

Effect Size and Power Analysis

Effect Size Measures

Cohen's d (standardized mean difference)

$$d = (\mu_1 - \mu_2) / \sigma_{\text{pooled}}$$

Interpretation:

d = 0.2: Small effect

d = 0.5: Medium effect

$d = 0.8$: Large effect

Eta-squared (η^2) for ANOVA

$$\eta^2 = SSB / SST$$

Cramér's V for Chi-square

$$V = \sqrt{(\chi^2 / (n \times \min(r-1, c-1)))}$$

R^2 for regression

- Proportion of variance explained

Power Analysis

Power = $P(\text{Reject } H_0 \mid H_0 \text{ is false})$

Power = $1 - \beta$

Factors affecting power:

- Effect size (larger = more power)
- Sample size (larger = more power)
- Significance level (α) (larger = more power)
- Population variance (smaller = more power)

Sample size calculation:

$$n = (z_{\alpha/2} + z_{\beta})^2 \times \sigma^2 / (\mu_1 - \mu_0)^2$$

Bayesian Statistics Basics

Bayes' Theorem

$$P(A|B) = P(B|A) \times P(A) / P(B)$$

Where:

$P(A|B)$ = Posterior probability

$P(B|A)$ = Likelihood

$P(A)$ = Prior probability
 $P(B)$ = Marginal probability

Bayesian vs Frequentist:

Frequentist: Parameters are fixed, data is random

Bayesian: Parameters are random, data is observed

Bayesian Inference

Prior beliefs + Data \rightarrow Posterior beliefs

Posterior \propto Likelihood \times Prior

Credible Intervals

- Bayesian equivalent of confidence intervals
- Probability that parameter lies in interval

Bayesian Hypothesis Testing

- Bayes Factor
- Posterior probability of hypotheses

Time Series Analysis

Components of Time Series

Trend: Long-term movement

Seasonality: Regular periodic patterns

Cyclical: Long-term fluctuations (non-regular)

Irregular/Random: Unpredictable fluctuations

Decomposition Models:

Additive: $Y(t) = \text{Trend} + \text{Seasonal} + \text{Irregular}$

Multiplicative: $Y(t) = \text{Trend} \times \text{Seasonal} \times \text{Irregular}$

Time Series Tests

Stationarity Tests:

- Augmented Dickey-Fuller (ADF) Test
- Phillips-Perron Test
- KPSS Test

Autocorrelation Function (ACF)

- Measures correlation between observations at different lags

Partial Autocorrelation Function (PACF)

- Correlation between observations k periods apart, controlling for intermediate observations

Experimental Design

Principles of Experimental Design

Randomization

- Random assignment to treatments
- Controls for confounding variables

Replication

- Multiple observations per treatment
- Increases precision and power

Blocking

- Group similar experimental units
- Controls for known sources of variation

Factorial Design

- Multiple factors studied simultaneously
- Can detect interactions between factors

A/B Testing

Steps:

1. Define hypothesis and metrics
2. Determine sample size
3. Randomize users to treatments
4. Collect data
5. Analyze results
6. Draw conclusions

Key Considerations:

- Statistical significance vs practical significance
- Multiple testing corrections
- Minimum detectable effect
- Statistical power

Common Statistical Mistakes

Interpretation Errors

Correlation \neq Causation

- Correlation does not imply causation
- Consider confounding variables
- Use experimental design for causal inference

p-hacking

- Multiple testing without correction
- Cherry-picking significant results
- Use Bonferroni or FDR corrections

Survivorship Bias

- Only analyzing successful cases
- Ignoring failures or dropouts

Simpson's Paradox

- Trend reverses when data is aggregated
- Consider lurking variables

Assumption Violations

Check assumptions before analysis:

- Normality (Q-Q plots, Shapiro-Wilk test)
- Independence (residual plots)
- Homoscedasticity (Levene's test)
- Linearity (scatterplots)

Solutions:

- Data transformations
- Non-parametric alternatives
- Robust statistical methods

Statistical Software Commands

Python (scipy.stats)

```
import scipy.stats as stats
import numpy as np

# Descriptive statistics
np.mean(data)
np.median(data)
np.std(data, ddof=1) # Sample std dev

# Hypothesis tests
stats.ttest_1samp(data, popmean)
stats.ttest_ind(group1, group2)
stats.chi2_contingency(contingency_table)
stats.f_oneway(group1, group2, group3)
```

```
# Distributions
stats.norm.pdf(x, loc=mu, scale=sigma)
stats.norm.cdf(x, loc=mu, scale=sigma)
stats.norm.ppf(q, loc=mu, scale=sigma)
```

R Commands

```
# Descriptive statistics
mean(data)
median(data)
sd(data)
summary(data)

# Hypothesis tests
t.test(data, mu=mu0)
t.test(group1, group2)
chisq.test(contingency_table)
aov(response ~ factor)

# Distributions
dnorm(x, mean, sd) # Density
pnorm(x, mean, sd) # CDF
qnorm(p, mean, sd) # Quantile
```