

Probability And Statistics

Michael Harmon

November 6, 2018

Contents

1	Probability & Random Variables	2
1.1	Probability	2
1.2	Random Variables	2
2	Probability Distributions	4
2.1	Discrete Distributions	4
2.2	Continuous Distributions	5
2.3	Approximation Theorems	6
3	Statistics & Estimators	7
3.1	Statistics	7
3.2	Convergence Theorems	8
3.3	Maximum Likelihood Esimators	9
3.4	Bayesian Estimators	10
3.5	Evaluating Estimators	11
4	Confidence Intervals	11
5	Statistical Testing	11
5.1	Hypothesis Testing	11
5.2	χ^2 “Goodness Of Fit” Tests	11
5.3	AB Testing	11

1 Probability & Random Variables

1.1 Probability

Sample Space (\mathcal{S}): The set of all possible outcomes.

Event: Any collection of possible outcomes, $E \in \mathcal{S}$ from the sample space.

Sigma Algebra: Σ is a collection of subsets of \mathcal{S} such that,

1. $\emptyset \in \Sigma$
2. If $A \in \Sigma$, then $A^c \in \Sigma$.
3. If $A_1, A_2, \dots \in \Sigma$ then $\cup_{i=1}^{\infty} A_i \in \Sigma$

Probability Function: Given a sample space \mathcal{S} and a Sigma algebra Σ then a probability function with domain Σ satisfies

1. $P(A) \geq 0, \quad \forall A \in \Sigma$
2. $P(\mathcal{S}) = 1$
3. If $A_1, A_2, \dots \in \Sigma$ are pairwise disjoint then $P(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$

Probabilities can be thought of as a frequency of occurrence.

Conditional Probabilities

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad P(B) > 0$$

Let A_1, A_2, \dots, A_n be a partition of \mathcal{S} then $P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$.

Theorem 1 *Bayes Theorem:*

Let A_1, A_2, \dots , be a partition of the space \mathcal{S} then,

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{P(B)} \tag{1.1}$$

$$= \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^{\infty} P(B|A_j)P(A_j)} \tag{1.2}$$

Two events A and B are **statistically independent** if

$$P(A \cap B) = P(A)P(B) \quad \text{or} \quad P(A|B) = P(A) \tag{1.3}$$

1.2 Random Variables

Random Variable: A function $X : \mathcal{S} \rightarrow \mathbb{R}$

Example: The sum of a roll of two die.

Probabilities can be induced by a random variable.

$$P_X(X = X_i), \quad \text{Discrete}$$

or

$$p(x), \quad \text{Continuous}$$

Cumulative Distribution Function: $F_X(x) = P_X(X < x)$

Two random variables X, Y are **identically distributed** if

$$\forall A \in \Sigma: \quad P(X \in A) = P(Y \in A)$$

Expectation Of Random Variables:

$$E(X) = \sum_i X_i P(X_i) \quad (1.4)$$

$$E(X) = \int x p(x) dx \quad (1.5)$$

Variance Of Random Variables: (w/ mean μ_X)

$$\text{Var}(X) = E[(X - \mu_X)^2] = E[X^2] - (E[X])^2 \quad (1.6)$$

Covariance: (of X and Y)

$$\text{cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

If X and Y are statistically independent then,

$$\text{cov}(X, Y) = E_X[(X - \mu_X)]E_Y[(Y - \mu_Y)] \quad (1.7)$$

$$= 0 \quad (1.8)$$

However, the converse is not true.

Correlation: (of X and Y)

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X^2 \sigma_Y^2}$$

Note that, $|\text{corr}(X, Y)| \leq 1$. We also remark that,

$$\text{Var}(aX \pm bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) \pm 2ab \text{cov}(X, Y)$$

so that if X and Y are statistically independent,

$$\text{Var}(aX \pm bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y)$$

Conditional Expectation Of Random Variables:

$$E(X|Y = y) = \sum_i X_i P(X_i|Y = y) \quad (1.9)$$

Moment Generating Function

$$M_X(t) = E_X[e^{tX}] \quad (1.10)$$

Note: $M_X^{(n)}(0) = E[X^n]$. The MGF has issues with existence.

Characteristic Function

$$\phi_X(t) = E_X[e^{itX}] \quad (1.11)$$

Note: $(-i)^n \phi_X^{(n)}(0) = E[X^n]$ and for X_i independent,

$$\phi_{\sum_{i=1}^n X_i}(t) = \prod_{i=1}^n \phi_{X_i}(t)$$

This characteristic function always exists!

Theorem 2 Chebyshev Inequality

Let X be a random variable and $g(X)$ be a non-decreasing function of X then $\forall r > 0$,

$$P(g(X) > r) \leq \frac{E[g(X)]}{r} \quad (1.12)$$

Convex Function: A function f is convex if $\forall x_1, x_2 \in X, \forall t \in [0, 1]$ then

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2) \quad (1.13)$$

Theorem 3 Jensen's Inequality

If X is a random variable and ϕ is a convex function then,

$$\phi(E[X]) \leq E[\phi(X)] \quad (1.14)$$

2 Probability Distributions

Probability distributions defined for both continuous and random variables. Random variables with discrete values have **discrete distributions**, while random variables with continuous values have **continuous distributions**.

2.1 Discrete Distributions

Bernoulli Distribution

A Bernoulli random variable binary outcome,

$$x = \begin{cases} 1, & \text{prob. } p \\ 0, & \text{prob. } 1-p \end{cases}$$

Distribution:

$$P(x|p) = p^x(1-p)^{(1-x)} \quad (2.1)$$

Mean: $E[x] = p$

Variance: $\text{Var}(x) = p(1-p)$

Binomial Distribution

A binomial random variable y is defined as the sum of n independent Bernoulli random variables, x_i all with prob. p :

$$y = \sum_{i=1}^n x_i$$

Distribution: The distribution is given as a function of $y = k$ (success), where $k \leq n$ where (n is the number of trials).

$$P(y = k | n, p) = \frac{n!}{(n-k)!k!} p^k (1-p)^{(n-k)} \quad (2.2)$$

Note: The product comes from the independence of the trials.

Mean: $E[x] = np$

Variance: $\text{Var}(x) = np(1-p)$

Note: These results can come from the definition of i.i.d property of the n Bernoulli trials and linearity of Var.

Geometric Distribution

A geometric random variable x is defined as the number $x = k$ of i.i.d Bernoulli trials *until* a success.

Distribution:

$$P(x = k) = p(1-p)^{k-1} \quad (2.3)$$

Mean: $E[x] = \frac{1}{p}$

Variance: $\text{Var}(x) = \frac{(1-p)}{p^2}$

Note: A geometric random variable is **memoryless**, i.e. if $s > t$ then $(P(x > s | x > t) = P(x > s - t))$.

Poisson Distribution

A Poisson random variable used to x is defined as the number of occurrences within a fixed time interval, given that the “average” number of occurrences is λ .

Distribution:

$$P(x | \lambda) = \frac{\lambda^x e^{-\lambda}}{x!}, \quad x = 1, 2, \dots \quad (2.4)$$

Mean: $E[x] = \lambda$

Variance: $\text{Var}(x) = \lambda$

2.2 Continuous Distributions

Beta Distribution

Distribution:

$$P(x | \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^\alpha (1-x)^{\beta-1} \quad (2.5)$$

For $0 < x < 1$, $\alpha > 0$ and $\beta > 0$.

Mean: $E[x] = \frac{\alpha}{\alpha+\beta}$

Variance: $\text{Var}(x) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$

Where,

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx \quad (2.6)$$

$$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \quad (2.7)$$

Note: Then Beta distribution is useful for deriving other distributions.

Normal Distribution

Distribution:

$$P(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/\sigma^2} \quad (2.8)$$

Mean: $E[x] = \mu$

Variance: $\text{Var}(x) = \sigma^2$

Note: Then,

$$P(|x - \mu| \leq \sigma) \simeq 0.67 \quad (2.9)$$

$$P(|x - \mu| \leq 2\sigma) \simeq 0.95 \quad (2.10)$$

$$P(|x - \mu| \leq 3\sigma) \simeq 0.99 \quad (2.11)$$

$$(2.12)$$

Student t Distribution

Arises from estimating mean of $N(\mu, \sigma^2)$ population, but where sample size is small and σ^2 is unknown. The **degrees of freedom** ($df > 2$) is

$$df = n - 1$$

Distribution:

Mean: $E[t] = 0$

Variance: $\text{Var}(t) = \frac{df}{df-2}$

Note: The t -distribution has fatter tails than normal distribution and is used for statistical significance between sample means and confidence intervals:

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

χ^2 Distribution

Let Z_1, \dots, Z_k be indep. normally distributed random variables then,

$$Q = \sum_{i=1}^k Z_i^2 \quad (2.13)$$

is χ^2 distributed with k degrees of freedom.

Distribution:

$$P(x, k) = \frac{1}{2^{k/2}\Gamma(k/2)} x^{k/2-1} e^{-x/2}, \quad x \in [0, \infty) \quad (2.14)$$

Mean: $E[x] = k$

Variance: $\text{Var}(x) = 2k$

Note: Is used for χ^2 “Goodness of fit” test.

2.3 Approximation Theorems

Theorem 4 *Normal Approx. To Binomial Random Variables*

When the number of trials, n , is sufficiently large then the probability of x success in n trials each having probability p can be approximated with,

$$\text{Binomial}(n, p)(x) \simeq \frac{1}{2\pi np(1-p)} e^{-(x-np)^2/np(1-p)} \quad (2.15)$$

Or $N_{\mu, \sigma^2}(x)$, where, $\mu = np$ and $\sigma^2 = np(1-p)$.

Theorem 5 *Normal Approx. To Poisson Random Variables*

When the average number of occurrences, λ , is sufficiently large then the probability of k occurrences can be approximated with,

$$\text{Poisson}(x, \lambda) \simeq \frac{1}{2\pi\lambda} e^{-(x-\lambda)^2/\lambda} \quad (2.16)$$

Or $N_{\mu, \sigma^2}(x)$, where, $\mu = \lambda$ and $\sigma^2 = \lambda$.

3 Statistics & Estimators

3.1 Statistics

A **statistic** is any function of a sample. An **Estimator** is any function of a sample *that is used to estimate a population parameter*.

Sample mean:

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

Sample variance

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

Theorem 6 Let X_1, \dots, X_n be a random sample from $N(\mu, \sigma^2)$ and $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$. Then,

1. \bar{X} and S^2 are independent random variables.
2. $\bar{X} \sim N(\mu, \sigma^2/n)$
3. $(n-1) \frac{S^2}{\sigma^2} \sim \chi_{n-1}^2$

Let X_1, \dots, X_n be a random sample from $N(\mu, \sigma^2)$ then

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} = \frac{(\bar{X} - \mu)(\sigma/\sqrt{n})}{\sqrt{S^2/\sigma^2}} = \frac{N(0, 1)}{\sqrt{\chi_{n-1}^2/(n-1)}} \quad (\text{Distributionally!})$$

That is,

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim \text{Student t-distribution with } n-1 \text{ deg. of freedom}$$

Unbiased Estimator: An estimator $T(X)$ for a parameter (Θ) is unbiased if,

$$E[T(X)] = \Theta$$

Theorem 7 Let X_1, \dots, X_n be a random sample with mean μ and variance S^2 , then

1. $E[\bar{X}] = \mu$
2. $E[S^2] = \sigma^2$
3. $Var(\bar{X}) = \frac{\sigma^2}{n}$

Sufficient Statistic: A statistic $T(\mathbf{x})$ is **sufficient statistic** for θ if the conditional distribution of the sample \mathbf{x} given $T(\mathbf{x})$ does not depend on θ .

Theorem 8 Sufficiency Principle If $T(\mathbf{x})$ is a sufficient statistic for θ , then any inference about θ should depend on the sample \mathbf{x} only through $T(\mathbf{x})$.

Theorem 9 Factorization Theorem Let $f(\mathbf{x}|\theta)$ denote the joint pdf or pm of sample \mathbf{X} . A statistic $T(\mathbf{x})$ is a sufficient statistic for θ iff $\exists g(\mathbf{t}|\theta)$ and $h(\mathbf{x})$ s.t. $\forall \mathbf{x}$ and θ ,

$$f(\mathbf{x}|\theta) = g(T(\mathbf{x})|\theta) h(\mathbf{x})$$

3.2 Convergence Theorems

Convergence In Probability Let X_1, \dots, X_n be a sequence of random variables then $X_i \rightarrow^P X$. If $\forall \epsilon > 0, \lim_{n \rightarrow \infty} P(|X_i - X| > \epsilon) = 0$.

Absolute In Convergence Let X_1, \dots, X_n be a sequence of random variables then $X_i \rightarrow^{A.S.} X$. If $\forall \epsilon > 0, P(\lim_{n \rightarrow \infty} |X_i - X| > \epsilon) = 0$.

Absolute convergence is equivalent to pointwise convergence.

Convergence In Distribution Let X_1, \dots, X_n be a sequence of random variables then $X_i \rightarrow^D X$. If $\lim_{n \rightarrow \infty} F_{X_n}(X) = F(X)$.

Convergence Relationships

$$X_n \rightarrow^{A.S.} X \Rightarrow X_n \rightarrow^P X$$

$$X_n \rightarrow^P X \Rightarrow X_n \rightarrow^D X$$

Consistent Estimator: An estimator $T(\mathbf{x})$ for θ is **consistent** if it converges in probability.

Consistency is the minimum requirement for an estimator!

Theorem 10 Let X_1, \dots, X_n be a sequence of random variables such that $X_n \rightarrow^P X$ and $h(X)$ is continuous. Then $h(X_n) \rightarrow^P h(X)$.

Theorem 11 S^2 is a consistent estimator (if $S_n \rightarrow 0$)

Proof:

$$\lim_{n \rightarrow \infty} P(|S_n^2 - \sigma^2| > \epsilon) \leq \lim_{n \rightarrow \infty} \frac{E(|S_n^2 - \sigma^2|^2)}{\epsilon} \quad (3.1)$$

$$= \lim_{n \rightarrow \infty} \frac{Var(S_n^2)}{\epsilon} \quad (3.2)$$

$$= 0 \quad (3.3)$$

by Chebyshev's theorem.

Theorem 12 *Weak Law Of Large Numbers*

Let X_1, \dots, X_n be i.i.d of random variables with $E[x_i] = \mu$ and $\text{Var}(x_i) = \sigma^2 < \infty$ then $x_n \rightarrow^P \mu$ **Proof:**

$$\lim_{n \rightarrow \infty} P(|\bar{x} - \mu| > \epsilon) \leq \lim_{n \rightarrow \infty} \frac{E(|\bar{x} - \mu|^2)}{\epsilon^2} \quad (3.4)$$

$$= \lim_{n \rightarrow \infty} \frac{\text{Var}(\bar{x})}{\epsilon^2} \quad (3.5)$$

$$= \lim_{n \rightarrow \infty} \frac{\sigma_n^2}{n \epsilon} \quad (3.6)$$

$$= 0 \quad (3.7)$$

Theorem 13 *Central Limit Theorem*

Let X_1, X_2, \dots , be i.i.d of random variables with $E[x_i] = \mu$ and $0 < \text{Var}(x_i) = \sigma^2 < \infty$ then,

$$\frac{X_n - \mu}{\sigma/\sqrt{n}} \rightarrow^D N(0, 1)$$

3.3 Maximum Likelihood Esimators

The **maximum likelihood estimator** is the value of a population distribution θ that maximizes the probability of observing the sample. We can find the MLE from a random sample x_1, x_2, \dots, x_n from $f(x|\theta)$ then the **likelihood** function is defined as,

$$L(\theta | x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i | \theta)$$

Then we can find the MLE $\hat{\theta}$ such that,

$$\frac{\partial L(\theta | x_1, x_2, \dots, x_n)}{\partial \theta} = 0 \quad \text{or} \quad \frac{\partial \log(L(\theta | x_1, \dots))}{\partial \theta} = 0$$

Theorem 14 *Invariance Property Of MLE*

If $\hat{\theta}$ is the MLE of θ , then for any function $\tau(\theta)$, the MLE of $\tau(\theta)$ is $\tau(\hat{\theta})$.

Example: X_1, X_2, \dots, X_n i.i.d. Bernoulli(ϕ). Find the MLE of ϕ .

The likelihood function is,

$$L(\phi | x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i | \phi) \quad (3.8)$$

$$= \prod_{i=1}^n \phi^{x_i} (1 - \phi)^{1-x_i} \quad (3.9)$$

$$= \phi^{\sum_{i=1}^n x_i} (1 - \phi)^{\sum_{i=1}^n (1-x_i)} \quad (3.10)$$

$$= \phi^{n\bar{x}} (1 - \phi)^{n(1-\bar{x})} \quad (3.11)$$

Or,

$$\log(L(\phi | \dots)) = n\bar{x} \log(\phi) + n(1-\bar{x}) \log(1-\phi) \quad (3.12)$$

So,

$$\frac{\partial \log(L)}{\partial \theta} = \frac{n \bar{x}}{\phi} - \frac{n(1 - \bar{x})}{(1 - \phi)} = 0 \quad (3.13)$$

$$\Rightarrow n \bar{x} (1 - \hat{\phi}) - n \hat{\phi} (1 - \bar{x}) = 0 \quad (3.14)$$

$$\Rightarrow \hat{\phi} = \bar{x} \quad (3.15)$$

This is the average of overall positive outcomes. We can test that this is the maximum, by taking the second derivative:

$$\frac{\partial^2 \log(L(\hat{\phi}))}{\partial \theta^2} = \frac{-n \bar{x}}{\hat{\phi}^2} - \frac{n(1 - \bar{x})}{(1 - \hat{\phi})^2} \Big|_{\hat{\phi}=\bar{x}} \quad (3.16)$$

$$= \frac{-n \bar{x}}{\bar{x}^2} - \frac{n(1 - \bar{x})}{(1 - \bar{x})^2} \quad (3.17)$$

$$= \frac{-n(1 - \bar{x}) - n \bar{x}}{\bar{x}(1 - \bar{x})} \quad (3.18)$$

The MLE has issues with existence.

Theorem 15 *The MLE is a consistent estimator
It also has optimal variance, but can be difficult to compute.*

3.4 Bayesian Estimators

In the **classical approach or frequentist**, θ is known, but fixed. x_1, x_2, \dots, x_n are drawn from a population index by *theta* and knowledge about the value of θ is obtained. In a **Bayesian approach**, θ is a quantity whose variation can be described by a probability distribution called a prior, $P(\theta)$. A sample is taken from a population and used to update the prior distribution, now called the posterior distribution, $P(\theta | \mathbf{x})$.

Let $f(\theta | \mathbf{x})$ be the sampling distribution then,

$$P(\theta | \mathbf{x}) = \frac{P(\mathbf{x} | \theta) P(\theta)}{m(\mathbf{x})}, \quad \text{where} \quad m(\mathbf{x}) = \int P(\mathbf{x} | \theta) P(\theta) d\theta$$

The **Bayesian estimator** could then be taken to be the expected value:

$$\hat{\theta} = E(\theta | \mathbf{x})$$

This requires us to calculate the full posterior distribution. Instead, one could use a **Bayesian estimator** that is called the **Maximum A-Posteriori (MAP)**:

$$\hat{\theta} = \max_{\theta} P(\theta | \mathbf{x})$$

Note: Bayesian estimators are ALWAYS biased due to their choice of prior, however, they can reduce the variance in our estimators.

Example: Let $y \sim \text{bin}(n, p)$ and $p \sim \text{beta}(\alpha, \beta)$, then

$$P(p | y) = \frac{P(y | p) P(p)}{m(y)} \quad (3.19)$$

$$= \binom{n}{k} p^y (1 - p)^{n-y} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\beta)} p^{\alpha-1} (1 - p)^{\beta-1} \quad (3.20)$$

$$= \text{beta}(y + \alpha, n - y + \beta) \quad (3.21)$$

This means the Bayesian estimate of p is,

$$\hat{p}_{\text{Bayes}} = \frac{y + \alpha}{y + \alpha + (n - y + \beta)} \quad (3.22)$$

$$= \frac{y + \alpha}{\alpha + \beta + n} \quad (3.23)$$

$$= \left(\frac{n}{\alpha + \beta + n} \right) \frac{y}{n} + \left(\frac{\alpha + \beta}{\alpha + \beta + n} \right) \left(\frac{\alpha}{\alpha + \beta} \right) \quad (3.24)$$

Note that this is a linear combination of the sample mean and the prior mean. However, as the sample size grows the contribution of the prior mean grows smaller and we get more confident in the sample mean.

In the limit as the sample size $n \rightarrow \infty$ Bayesian and classical estimators should converge to the same estimator.

Indeed, the maximum likelihood estimator is the same thing as a maximum a-posteriori estimator with uniform prior!

3.5 Evaluating Estimators

For continuous random variables

$$\text{MSE} = E_{\theta}(\hat{\theta} - \theta)^2 \quad (3.25)$$

$$= \quad (3.26)$$

4 Confidence Intervals

5 Statistical Testing

5.1 Hypothesis Testing

5.2 χ^2 “Goodness Of Fit” Tests

5.3 AB Testing