

EXAM HANDOUT

Counting Outcomes:

Permutations
with
replacement

$$P_r(n, k) = \overbrace{n \cdot n \cdot \dots \cdot n}^{k \text{ terms}} = n^k$$

Permutations
without
replacement

$$P(n, k) = \overbrace{n(n-1)(n-2) \cdot \dots \cdot (n-k+1)}^{k \text{ terms}} = \frac{n!}{(n-k)!}$$

Combinations
without
replacement

$$C(n, k) = \binom{n}{k} = \frac{P(n, k)}{P(k, k)} = \frac{n!}{k!(n-k)!}$$

Combinations
with replacement

$$C_r(n, k) = \binom{k+n-1}{k} = \frac{(k+n-1)!}{k!(n-1)!}$$

Union of events:

Probability
of a union

$$P\{A \cup B\} = P\{A\} + P\{B\} - P\{A \cap B\}$$

For mutually exclusive events,
 $P\{A \cup B\} = P\{A\} + P\{B\}$

To scale this up for mutually exclusive events just add their probabilities together.
To scale this up for non-mutually exclusive events, use principle of inclusion and exclusion

Intersection of independent events

Independent
events

$$P\{E_1 \cap \dots \cap E_n\} = P\{E_1\} \cdot \dots \cdot P\{E_n\}$$

Intersection of dependent events are dealt with by conditional probability.

Intersection,
general case

$$P\{A \cap B\} = P\{B\} P\{A \mid B\}$$

Conditional Probability

DEFINITION 2.15

Conditional probability of event A given event B is the probability that A occurs when B is known to occur.

Conditional
probability

$$P\{A \mid B\} = \frac{P\{A \cap B\}}{P\{B\}}$$

For independent events, $P\{A \mid B\} = P\{A\}$

Bayes Rule

Bayes
Rule

$$P\{B \mid A\} = \frac{P\{A \mid B\} P\{B\}}{P\{A\}}$$

Total Probability

Law of Total
Probability

$$P\{A\} = \sum_{j=1}^k P\{A \mid B_j\} P\{B_j\}$$

In case of two events ($k = 2$),

$$P\{A\} = P\{A \mid B\} P\{B\} + P\{A \mid \overline{B}\} P\{\overline{B}\}$$

Bayes Rule
for two events

$$P\{B \mid A\} = \frac{P\{A \mid B\} P\{B\}}{P\{A \mid B\} P\{B\} + P\{A \mid \overline{B}\} P\{\overline{B}\}}$$

Random Variables

Distribution	Discrete	Continuous
Definition	$P(x) = P\{X = x\}$ (pmf)	$f(x) = F'(x)$ (pdf)
Computing probabilities	$P\{X \in A\} = \sum_{x \in A} P(x)$	$P\{X \in A\} = \int_A f(x)dx$
Cumulative distribution function	$F(x) = P\{X \leq x\} = \sum_{y \leq x} P(y)$	$F(x) = P\{X \leq x\} = \int_{-\infty}^x f(y)dy$
Total probability	$\sum_x P(x) = 1$	$\int_{-\infty}^{\infty} f(x)dx = 1$

If X is a discrete variable,

$$P(X \leq x) = F(x) \quad P(X \geq x) = 1 - P(X < x) = 1 - P(X \leq x-1) = 1 - F(x-1)$$

$$P(X = x) = F(x) - F(x-1) \quad P(a \leq X \leq b) = F(b) - F(a)$$

For a continuous distribution,

$$P(X = x) = 0 \quad P(X \leq x) = P(X < x) = F(x) \quad P(X \geq x) = P(X > x) = 1 - F(x)$$

$$P(a \leq X \leq b) = F(b) - F(a) \quad P(x-h \leq X \leq x+h) \approx 2h \cdot f(x) \approx P(X = x) \text{ (If } h \text{ is really small)}$$

Joint Distributions:

The following table shows how to calculate marginal distributions, check for independence and how to calculate probabilities of random vectors using joint probability distributions

Distribution	Discrete	Continuous
Marginal distributions	$P(x) = \sum_y P(x, y)$ $P(y) = \sum_x P(x, y)$	$f(x) = \int f(x, y)dy$ $f(y) = \int f(x, y)dx$
Independence	$P(x, y) = P(x)P(y)$	$f(x, y) = f(x)f(y)$
Computing probabilities	$P\{(X, Y) \in A\}$ $= \sum_{(x,y) \in A} P(x, y)$	$P\{(X, Y) \in A\}$ $= \iint_{(x,y) \in A} f(x, y) dx dy$

Properties (or Moments) of a Distribution

Discrete	Continuous
$\mathbf{E}(X) = \sum_x xP(x)$ $\text{Var}(X) = \mathbf{E}(X - \mu)^2$ $= \sum_x (x - \mu)^2 P(x)$ $= \sum_x x^2 P(x) - \mu^2$ $\text{Cov}(X, Y) = \mathbf{E}(X - \mu_X)(Y - \mu_Y)$ $= \sum_x \sum_y (x - \mu_X)(y - \mu_Y) P(x, y)$ $= \sum_x \sum_y (xy) P(x, y) - \mu_x \mu_y$	$\mathbf{E}(X) = \int x f(x) dx$ $\text{Var}(X) = \mathbf{E}(X - \mu)^2$ $= \int (x - \mu)^2 f(x) dx$ $= \int x^2 f(x) dx - \mu^2$ $\text{Cov}(X, Y) = \mathbf{E}(X - \mu_X)(Y - \mu_Y)$ $= \iint (x - \mu_X)(y - \mu_Y) f(x, y) dx dy$ $= \iint (xy) f(x, y) dx dy - \mu_x \mu_y$

Std(X) is square root of Var(X)

DEFINITION 3.9

Correlation coefficient between variables X and Y is defined as

$$\rho = \frac{\text{Cov}(X, Y)}{(\text{Std}X)(\text{Std}Y)}$$

**Properties
of
expectations**

$$\mathbf{E}(aX + bY + c) = a\mathbf{E}(X) + b\mathbf{E}(Y) + c$$

In particular,

$$\mathbf{E}(X + Y) = \mathbf{E}(X) + \mathbf{E}(Y)$$

$$\mathbf{E}(aX) = a\mathbf{E}(X)$$

$$\mathbf{E}(c) = c$$

For **independent** X and Y ,

$$\mathbf{E}(XY) = \mathbf{E}(X)\mathbf{E}(Y)$$

Properties of variances and covariances

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

$$\begin{aligned} \text{Cov}(aX + bY, cZ + dW) \\ = ac \text{Cov}(X, Z) + ad \text{Cov}(X, W) + bc \text{Cov}(Y, Z) + bd \text{Cov}(Y, W) \end{aligned}$$

$$\text{Cov}(X, Y) = \text{Cov}(Y, X)$$

$$\rho(X, Y) = \rho(Y, X)$$

In particular,

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

$$\text{Cov}(aX + b, cY + d) = ac \text{Cov}(X, Y)$$

$$\rho(aX + b, cY + d) = \rho(X, Y)$$

For independent X and Y ,

$$\text{Cov}(X, Y) = 0$$

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$

Only a large variance will allow a Variable to vary greatly from the expected value. This is shown by Chebyshev's Inequality

Chebyshev's inequality

$$P\{|X - \mu| > \varepsilon\} \leq \left(\frac{\sigma}{\varepsilon}\right)^2$$

for any distribution with expectation μ and variance σ^2 and for any positive ε .

This shows that a Variable with a high variance has bigger risk of varying from the expected amount by a large value.

Sample Statistics (Used as estimators for Population Parameters)

Mean

DEFINITION 8.3

Sample mean \bar{X} is the arithmetic average,

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}$$

Variance and Std. Deviation

DEFINITION 8.8

For a sample (X_1, X_2, \dots, X_n) , a **sample variance** is defined as

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2. \quad (8.4)$$

It measures variability among observations and estimates the population variance $\sigma^2 = \text{Var}(X)$.

Sample standard deviation is a square root of a sample variance,

$$s = \sqrt{s^2}.$$

It measures variability in the same units as X and estimates the population standard deviation $\sigma = \text{Std}(X)$.

Alt formula for Variance:

$$s^2 = \frac{\sum_{i=1}^n X_i^2 - n\bar{X}^2}{n-1}.$$

Quantiles, percentiles and quartiles

DEFINITION 8.7

A **p -quantile** of a population is such a number x that solves equations

$$\begin{cases} P\{X < x\} \leq p \\ P\{X > x\} \leq 1 - p \end{cases}$$

A **sample p -quantile** is any number that exceeds at most $100p\%$ of the sample, and is exceeded by at most $100(1-p)\%$ of the sample.

A **γ -percentile** is (0.01γ) -quantile.

First, second, and third **quartiles** are the 25th, 50th, and 75th percentiles. They split a population or a sample into four equal parts.

A **median** is at the same time a 0.5-quantile, 50th percentile, and 2nd quartile.

Shape of a distribution (comparing mean and median)

Symmetric distribution $\Rightarrow M = \mu$

Right-skewed distribution $\Rightarrow M < \mu$

Left-skewed distribution $\Rightarrow M > \mu$

IQR and outliers.

DEFINITION 8.10

An **interquartile range** is defined as the difference between the first and the third quartiles,

$$IQR = Q_3 - Q_1.$$

It measures variability of data. Not much affected by outliers, it is often used to detect them. IQR is estimated by the *sample interquartile range*

$$\widehat{IQR} = \hat{Q}_3 - \hat{Q}_1.$$

Any samples that are less than $Q_1 - 1.5(IQR)$ or more than $Q_3 + 1.5(IQR)$ can be treated as potential outliers.

Standard error of any estimator is its std deviation.

$$\left\| \begin{array}{l} \sigma(\hat{\theta}) = \text{standard error of estimator } \hat{\theta} \text{ of parameter } \theta \\ s(\hat{\theta}) = \text{estimated standard error} = \hat{\sigma}(\hat{\theta}) \end{array} \right\|$$

Graphical Statistics:

Can be used to visualize the given samples to make observations about the nature of the population

- Histograms are bar charts for columns for each bin
 - If height of bin is freq count: Frequency histogram
 - If height of bin is proportion of data: Relative Frequency histogram
 - Each sample value can have its own bin, or you can have multiple nearby values in one bin.
 - You can use histograms to guess shape of distribution.
- Stem and leaf plots
 - Choose the stem such that the values are not all limited to one stem.
 - Distribute values to each stem and sort the leaf values.
 - You can use this to calculate mean, median and guess shape of distribution
 - It can also be used to compare two distributions
- Boxplots
 - Boxplots are based on 5-point summaries
 - $< \min(X_i), \widehat{Q}_1, \widehat{M}, \widehat{Q}_3, \max(X_i) >$
 - Represent sample mean with small cross. Draw a box between sample Q1 and Q3 and draw a line for sample median. Draw whiskers to smallest sample and largest sample that is within the 1.5 IQR range. Draw dots for all samples outside the 1.5 IQR range.
 - Can be used to compare multiple distributions.



Entropy – The Mean Information

- For any random variable X with a probability mass function (pmf) p_i it is thus possible to measure the information content of each outcome

$$I(x_i) = I(p_i) = -\log_b(p_i)$$

- The expected information of an experiment (i.e. of the unknown value of X) can be computed

$$\begin{aligned} H(X) &= E[I(X)] = \sum_{x \in \mathcal{X}} p(x) I(x) \\ &= -\sum_{x \in \mathcal{X}} p(x) \log_b p(x) \end{aligned}$$

© Manfred Huber 2017

6



Relative Entropy

- Relative Entropy (Kullback-Leibler distance) is a measure of the difference of two distributions

$$\begin{aligned} D(p \parallel q) &= E_p \left[\log_b(p(x)/q(x)) \right] \\ &= \sum_x p(x) \log_b(p(x)/q(x)) \end{aligned}$$

- Measures not the difference in the amount of information but difference in the information itself
- Both distributions have to be defined over the same domain
- Is always positive and zero only if the two distributions are identical

© Manfred Huber 2017

12

Bernoulli distribution

- Used to model experiments with a binary outcome (yes/no, pass/fail, true/false)
 - Called Bernoulli Trials
- Random variable can take values 0 (fail) or 1 (pass)
- p = probability of success

$$P(x) = \begin{cases} 1 - p & \text{if } x = 0 \\ p & \text{if } x = 1 \end{cases}$$

$$E(X) = p$$

$$Var(X) = p(1 - p)$$

Binomial Distribution

- It is used to model number of success in a sequence of **independent** Bernoulli trials
- Models the probability of x successes in n trials
- p = probability of success; n = number of trials

$$P(x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

$$E(X) = np$$

$$Var(X) = np(1 - p)$$

Geometric Distribution

- It is used to model number trials needed to achieve the first success in a sequence of **independent** Bernoulli trials
- Models the probability of x^{th} successive trial resulting in a success
- p = probability of success

$$P(x) = (1 - p)^{x-1}p$$

$$E(X) = \frac{1}{p}$$

$$Var(X) = \frac{1 - p}{p^2}$$

Negative Binomial Distribution

- It is used to model number of trials needed to obtain k success in a sequence of **independent** Bernoulli trials
- Models the probability of x^{th} successive trial resulting in the k^{th} success
- p = probability of success; k = number of success

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

$$E(X) = \frac{k}{p}$$

$$Var(X) = \frac{k(1-p)}{p^2}$$

Side Note: Calculating Neg. Binomial in Practice

- If X follows Negative Binomial(k, p)
 - $P(X = x) = \text{prob of needing } x \text{ trials for } k \text{ success} = \text{prob of } k\text{th trial being success} * \text{prob of getting } k-1 \text{ success in } x-1 \text{ trials} = p * P(Y = k-1)$
 - Where Y follows Binomial($x-1, p$)
 - $P(X \geq x) = \text{prob of needing } \geq x \text{ trials for } k \text{ success} = \text{prob of } x-1 \text{ trials not having } \leq k-1 \text{ success} = P(Y \leq k-1)$
 - Where Y follows Binomial($x-1, p$)

Poisson Distribution

- It is used to model number rare events occurring within a fixed period of time
- Models the probability of x rare events occurring in a fixed period of time if we know the frequency at which the events occur on average
- λ = frequency (average number of events in a fixed time period)

$$P(x) = e^{-\lambda} \frac{\lambda^x}{x!}$$

$$E(X) = \lambda$$

$$\text{Var}(X) = \lambda$$

Side Note: Poisson Approx. of Binomial

- If the number of trials is large and the probability of success is low, then we can use Poisson Distribution to approximate the Binomial Distribution
 - Also works if probability of failure is very low
- $\lim_{\substack{n \rightarrow \infty \\ p \rightarrow 0 \\ np \rightarrow \lambda}} \binom{n}{x} p^x (1-p)^{n-x} = e^{-\lambda} \frac{\lambda^x}{x!}$
- Can use this approximation if $n \geq 30$ and $p \leq 0.05$

Uniform Distribution

- Used to model scenarios where outcome lies within a given interval (a,b) and all outcomes are equally likely.
- If interval is (0,1) it is called standard uniform distribution

$$f(x) = \frac{1}{b-a}, a < x < b$$

$$E(X) = \frac{a+b}{2}$$

$$Var(X) = \frac{(b-a)^2}{12}$$

Exponential Distribution

- Used to model time (or separation) between events occurring at frequency λ (rate at which the events occur)

$$\begin{aligned}f(x) &= \lambda e^{-\lambda x}, x > 0 \\E(X) &= \frac{1}{\lambda} \\Var(X) &= \frac{1}{\lambda^2}\end{aligned}$$

Gamma Distribution

- Used to model total time of multistage processes with α steps (shape parameter) where time of each step can be modeled as a Exponential distribution with frequency λ .

$$\begin{aligned}f(x) &= \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x} \quad \text{if } \alpha > 0 \ x > 0 \\ \Gamma(\alpha) &= \int_0^\infty y^{\alpha-1} e^{-y} dy \quad \text{if } \alpha > 0 \\ \text{also } \Gamma(\alpha) &= (\alpha - 1)! \quad \text{if } \alpha \text{ is a positive integer} \\ E(X) &= \frac{\alpha}{\lambda} \\ Var(X) &= \frac{\alpha}{\lambda^2}\end{aligned}$$

- Please note that $\text{Gamma}(1, \lambda) = \text{Exponential}(\lambda)$

Side-Note: Gamma-Poisson Formula

- Can be used to simplify calculation of probabilities of RV T with Gamma Distribution.

$$\{T > t\} = \{X < \alpha\}$$

- Where T has Gamma distribution with parameters α (number of events) and λ (frequency of each event).
- X models the number of events that occurs before time t. It has Poisson distribution with parameter λt .

So,

$$P\{T > t\} = P\{X < \alpha\}$$

$$P\{T \leq t\} = P\{X \geq \alpha\}$$

Where T has $\text{Gamma}(\alpha, \lambda)$ distribution and X has $\text{Poisson}(\lambda t)$ distribution

Normal Distribution

- Used to model a large number of scenarios
 - Sums, averages or errors: Mainly due to CLT
 - Naturally occurring phenomena
- Allows you to model a scenario on the basis of expectation μ (location parameter) and standard deviation σ (scale parameter)

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

$$E(X) = \mu$$

$$\text{Var}(X) = \sigma^2$$

Central Limit Theorem

- If you have a Random variable that is expressed as a sum of large number of independent random variables (usually ≥ 30), then you can use this theorem to model them as a Normal Distribution.

Theorem 1 (CENTRAL LIMIT THEOREM) *Let X_1, X_2, \dots be independent random variables with the same expectation $\mu = \mathbf{E}(X_i)$ and the same standard deviation $\sigma = \text{Std}(X_i)$, and let*

$$S_n = \sum_{i=1}^n X_i = X_1 + \dots + X_n.$$

As $n \rightarrow \infty$, the standardized sum

$$Z_n = \frac{S_n - \mathbf{E}(S_n)}{\text{Std}(S_n)} = \frac{S_n - n\mu}{\sigma\sqrt{n}}$$

converges in distribution to a Standard Normal random variable, that is,

$$F_{Z_n}(z) = \mathbf{P}\left\{\frac{S_n - n\mu}{\sigma\sqrt{n}} \leq z\right\} \rightarrow \Phi(z)$$

for all z .

Side-Note: Using Normal to Approx. Binomial

- A binomial distribution is a sum of n Bernoulli trials. So if n is large (≥ 30) but p is not small enough (or large enough) to use Poisson approximation ($0.05 \leq p \leq 0.95$) then we can model the binomial as a sum of Bernoulli distributions with mean p and variance $p(1 - p)$.
- So by Central Limit Theorem,

$$\text{Binomial}(n, p) \approx \text{Normal}\left(\mu = np, \sigma = \sqrt{np(1 - p)}\right)$$

- This normal distribution can be calculated by converting it to a standard normal distribution

Side-Note: Using Normal to Approx. Gamma

- Similarly if α is large enough the we can use CLT to approx. a Gamma distribution using a normal distribution

$$Gamma(\alpha, \lambda) \approx Normal\left(\frac{\alpha}{\lambda}, \sqrt{\frac{\alpha}{\lambda^2}}\right)$$