The note of UIUC course STAT 400 in Fall 2023, by Zory Zhang. In case of any broken math renderring in Github preview, please open this markdown file using your own markdown editor. PDF version (maybe not up-to-date) can be found here.

For robust statistics, see the page.

Reminder

• The two properties of variance can save a lot of work.

Preliminaries

The following content is based on the course STAT400@UIUC. Probability knowledge is assumed.

Textbook

Probability and Statistical Inference, Tenth Edition; by Robert V. Hogg, Elliot A. Tanis, Dale L. Zimmerman; Hoboken, NJ: Pearson; 2020.

Ch0 Matrix cookbook

- The matrix cookbook
- a for scalar, \underline{a} for vector, A for matrix.
- X for random variable, \underline{X} for random vector, \mathbf{X} for random matrix.
- $\mathbb{E}[\underline{a}^T\underline{X}] = \underline{a} \cdot \mathbb{E}\underline{X}$; $\mathbb{E}[A\underline{X}] = A\mathbb{E}\underline{X}$; $\mathbb{E}A\mathbf{X}B = A\mathbb{E}[\mathbf{X}]B$.
- $\bullet \ \ \underline{\mu}:= \mathrm{I}\!\mathrm{E}\underline{X}, \underline{v}:= \mathrm{I}\!\mathrm{E}\underline{Y}.$

Ch1 Probability

1.1 Discrete distributions

Def. (**Bernoulli distribution**, $X \sim Be(p)$) Bernoulli trial is that $A \in \mathcal{F}$, and we call the trial a success if A occurs. Bernoulli distribution is based on single Bernoulli

trial. $m_1 = m_2 = p, Var = p(1 - p)$.

Def. (**Binomial distribution**, $Y \sim Bin(n,p)$) Perform n independent Bernoulli trials with $p = \mathbb{P}(A)$, and let Bernoulli r.v.s. $X_1, X_2 \dots X_n$ be the indicator function of success of the experiments. Let $Y := \sum X_i$, then the p.m.f.

 $f_Y(k) = \binom{n}{k} p^k (1-p)^{n-k}, k = 1, 2, \dots, n$. **Binomial process** is that Y_n to be the number of successes in the first n Be(p) trials.

Def. (**Geometric distribution**, $W \sim Geom(p)$) Keep performing Bernoulli trials until the first success, and let W:= the waiting time, then $f_W(k)=(1-p)^{k-1}p, k=1,2,...,F(x)=1-(1-p)^x.$

Def. (Negative Binomial distribution, $W_r \sim NB(r,p)$) Let

 $W_r :=$ the waiting time for r successes, then

$$f_{W_r}(k) = inom{k-1}{k-r} p^{r-1} (1-p)^{k-r} p, k = 1, 2, \ldots$$

Def. (**Approximate poisson process** with param $\lambda > 0$) Let the number of **occurrences** of some event in a given **continuous interval** be counted. If

- 1. The number of occurrences in nonoverlapping subintervals are independent.
- 2. The probability of exactly one occurrence in a sufficiently short subinterval of length h is approximately λh .
- 3. The probability of more than one occurrence in a sufficiently short subinterval of length is essentially 0.

Def. (**Poisson distribution**, $X \sim poisson(\lambda)$) An approximation of number of occurances in an interval of length 1 is given. Consider partition the unit interval into n subintervals. As n becomes larger, the length of subinterval become smaller and thus more likely the last two properties hold. For a large n,

$${
m I\!P}(X=k) \sim {n\choose k} (rac{\lambda}{n})^k (1-rac{\lambda}{n})^{n-k}$$
, and it converges to $f_X(k):=rac{\lambda^k}{k!}e^{-\lambda}, k=0,1,2,\ldots$

Rmk. $M(t) = e^{\lambda(e^t - 1)}$, then $\mu = \sigma^2 = \lambda$.

E.g. Number of customers of a shop between 5-6pm.

Def. (**Hypergeometric distribution**) $f_X(k)=rac{\binom{K}{k}\binom{N-K}{n-k}}{\binom{N}{n}}$. IE $X=nN_1/N$, $VarX=nrac{N_1}{N}rac{N-K}{N}rac{N-n}{N-1}$.

Def. (**Multinomial**) Multi-outcome version of hypergeometric. $f_X(x) = \frac{n!}{\prod x_i!} \prod p_i^{x_i}$, where $\mathbb{E} X_i = np_i, Var X_i = np_i (1-p_i)$.

1.2 Continuous distributions

Def. (**Exponential distribution**, $X \sim Exp(\lambda)$) Define X as the waiting time for the first occurance in a poisson process. Then

 $F_X(x) = 1 - {
m I\!P} ext{(no occurences in [0,x])},$ therefore

$$F_X(x)=egin{cases} 0,x<0\ 1-e^{-\lambda x},x\geq 0 \end{cases}, f_X(x)=F_X'(x)=egin{cases} 0,x<0\ \lambda e^{-\lambda x},x\geq 0 \end{cases}$$

Rmk.

- 1. It's usually expressed as $\theta = \frac{1}{\lambda}, f(x) = \frac{1}{\theta}e^{-\frac{x}{\theta}}$.
- 2. $M(t) = \frac{1}{1-\theta t}$ for $t < \frac{1}{\theta}$, thus $\mu = \theta, \sigma^2 = \theta^2$, Fisher's moment coefficient of skewness is 2.
- 3. (Memoryless) $\mathbb{P}(X>s+t|X>t)=\mathbb{P}(X>s)$.
- 4. (Failture/hazard rate for positive r.v.) $r(t)=\lim_{h\downarrow 0} \tfrac{1}{h} \mathrm{I\!P}(X\leq t+h|X>t)=\tfrac{f(t)}{1-F(t)}=\lambda.$

Def. (**Gamma function**) For t > 0, $\Gamma(t) = \int_0^\infty y^{t-1} e^{-y} dy$.

Rmk. $\Gamma(t) = (-y^{t-1}e^{-y})_{y=0}^{\infty} + \int_0^{\infty} (t-1)y^{t-2}e^{-y} dy = (t-1)\Gamma(t-1)$. $\Gamma(1) = 1$. When t is an integer, $\Gamma(t) = (t-1)!$.

Def. (**Gamma distribution**, $X\sim Gamma(\alpha,\theta)$) For $x\geq 0$, $f_X(x)=\frac{1}{\Gamma(\alpha)\theta^\alpha}x^{\alpha-1}e^{-x/\theta}$. In a poisson process, the waiting time of first occurence is $Exp(\theta=1/\lambda)$, now gamma distribution describes the waiting time until α -th occurrence.

Rmk. $\mathbb{E}X = \alpha\theta, VarX = \alpha\theta^2$.

Def. (Normal/Gaussian distribution, $X \sim \mathcal{N}(\mu, \sigma^2)$) $f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$. Standard normal distribution $Y \sim \mathcal{N}(0,1)$, denote the p.d.f. as $\phi(x) = \frac{1}{\sqrt{2\pi}} \exp(\frac{-x^2}{2})$.

Rmk. To compute $I=\int_{-\infty}^{\infty}\phi(x)~\mathrm{d}x$, consider

$$egin{align} I^2 &= \int_{-\infty}^{\infty} \phi(x) \, \mathrm{d}x \int_{-\infty}^{\infty} \phi(y) \, \mathrm{d}y \ &= rac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(-rac{x^2 + y^2}{2}) \, \mathrm{d}x \, \mathrm{d}y \ &= rac{1}{2\pi} \int_{0}^{2\pi} \int_{0}^{\infty} \exp(-rac{r^2}{2}) r \, \mathrm{d}r \, \mathrm{d} heta = 1 \end{split}$$

Rmk. For normal distribution, the i-th central moment of $N(\mu, \sigma^2)$ is $\sigma_i = \begin{cases} 0, & \text{if i odd,} \\ \sigma^i(i-1)!!, & \text{o.w.} \end{cases}$

Def. (**(standard) Cauchy distribution**, $X \sim \mathcal{N}(\mu, \sigma^2)$) p.d.f. $f_X(x) = \frac{1}{\pi(1+x^2)}, F_X(x) = \arctan(x) + \frac{\pi}{2}.$

Rmk. (Heavy tail distribution without moment) Notice the $\mathbb{E} X = \int_{-\infty}^{\infty} |x| f_X(x) \; \mathrm{d}x = \infty$ does not exist.

1.3 Covariance

1.
$$\sigma_{XY} = Cov[X, Y] = \mathbb{E}[(X - \mathbb{E}X)(Y - \mathbb{E}Y)] = \mathbb{E}(XY) - \mu_X \mu_Y$$
.

$$[2.\ V[rac{X}{\sigma_X}]=1,
ho_{XY}=Cov[ilde{X}, ilde{Y}]=rac{Cov[X,Y]}{\sigma_X\sigma_Y}\in[-1,1].$$

- 3. The **covariance matrix** $V[X] = \mathbb{E}[(X \mu)(X \mu)^T] = \mathbb{E}XX^T \mu\mu^T$.
- 4. The **variance of r.v.** in the form of $\underline{a} \cdot \underline{X}$: $V[\underline{a} \cdot \underline{X}] = \underline{a}^T V[\underline{X}]\underline{a}$; $V[\underline{A}\underline{X}] = AV[\underline{X}]A^T$. When $||\underline{a}|| = 1$, it's taking the variance in the direction of \underline{a} , over the joint p.d.f. (consider variance as the width of the function picture along certain axis).

1.4 CLT

Prop. (Independent normal) $\sum_{i=1}^n c_i X_i \sim N(\sum c_i \mu_i, \sum c_i^2 \sigma_i^2)$.

Cor. When all Xi are from the same distribution, $\frac{1}{n} \sum X_i \sim N(\mu, \frac{\sigma^2}{n})$.

Thm. (CLT)
$$rac{1}{n}S_n:=rac{1}{n}\sum_{i=1}^n X_i\sim N(\mu,rac{\sigma^2}{n})$$
.

Rmk. Normal approximation work especially well on Binomial distribution when both np, n(1-p) > 5.

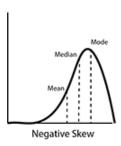
Rmk. (**Continuity correlation**) When normal approximation is applied to interger sample space, we need to adjust the event space. E.g.

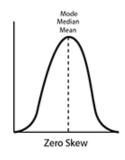
$$\mathbb{P}(3 \le X \le 5) = \mathbb{P}(2.5 \le X \le 5.5).$$

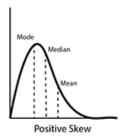
1.5 Generating function

Def. (k-th moments) $m_k := \mathbb{E} X^k = \int x^k f_X(x) \ \mathrm{d}x.$

Rmk. Third and forth moments are called skewness and kurtosis respectively.







Motivation. $G_X(s)=\mathbb{E} s^X, M_X(t)=\mathbb{E} e^{tX}, \phi_X(t)=\mathbb{E} e^{itX}$. Then $X\perp Y$, we have $G_{X+Y}=G_XG_Y$. The converse may not be true.

Rmk. $rac{\mathrm{d}^k}{\mathrm{d}t^k} M_X(t)|_{t=0} = \int x^k \, \mathrm{d}F_X(x) = m_k$.

STAT

Def. (Sample mean, variance, and moment) n samples $X_1, \ldots X_n, n$ observations $x_1, \ldots x_n$.

- 1. The sample mean $\bar{x} := \frac{1}{n} \sum_{i=1}^{n} x_i$;
- 2. The sample variance $s^2 := rac{1}{n-1} \sum_{i=1}^n (x_i ar{x})^2$;
- 3. The sample moment $\frac{1}{n} \sum_{i=1}^{n} x_i^k$.

These sample statistical quantities are contrast to population statistical quantities. Essentially, they are just random variables before observed.

2.1 Point estimation

Def. (**Point estimation**) Given a dist with some known params and a unknown param θ . The parameter space Ω is the range of all possible values of θ . Now we need to find our best guess for θ to have a good estimate of dist. Let an **estimator** of θ be $\hat{\theta} := u(X_1, X_2, \dots X_n)$. The instance $u(x_1, x_2, \dots x_n)$ is a **point estimate**.

Rmk. (Maximum likelihood estimate, MLE) For a likelihood func

$$L(\theta) := \prod_{i=1}^n f(x_i; \theta)$$
, the MLE $\hat{\theta} := \operatorname{argmax}_{\theta} L(\theta)$.

Def. (**Bias of estimator**) Estimator $\hat{\theta}$ has bias $\mathbb{E}\hat{\theta} - \theta$. E.g. S_X^2 is an **unbiased estimator** of the variances.

Rmk. (**Method of moments, MOM**) Let i-th sample moment equal to i-th theoretical moment, i from 1 to the # unknown params, then solve the equations.

Def. (Chi-squared distribution) If k iid $Z_i \sim N(0,1)$, then $\sum_{i=1}^k Z_i^2 \sim X_{(k)}^2 = Gamma(\alpha=k/2,\theta=2)$. k is called the **degree of freedom**, DF.

Prop.
$$rac{(n-1)S^2}{\sigma^2}\sim X_{(n-1)}^2.$$

Proof.

1. Assume #TODO

2.2 Confidence Interval, CI

Def. (**Student's t-distribution**) $T:=rac{Z}{\sqrt{U/r}}\sim t_r$, where r is the DF, $Z\sim N(0,1), U\sim X^2(r), Z\perp U.$ Then $f_T(t)=rac{\Gamma\left(rac{r+1}{2}
ight)}{\sqrt{\pi r}\Gamma\left(rac{r}{2}
ight)(1+t^2/r)^{(r+1)/2}}.$

Rmk. IE $T=0, VarT=\frac{v}{v-2}$. When the sample size n is relatively big, the difference between normal and t-distribution is small.

E.g.

1. If the underlying distribution is normal, then

$$Z=rac{ar{x}-\mu}{\sigma/\sqrt{n}}\sim N(0,1); T=rac{ar{x}-\mu}{s/\sqrt{n}}\sim t_{n-1}$$
. #NotCovered

- 2. The above can serve as approximation #NotCovered
- 3. "When the distribution is not normal but is unimodal (has only one mode), symmetric, and continuous, the approximation is usually quite good even for small n, such as n = 5.... In almost all cases encountered in real applications, an n of at least 30 is usually adequate."
- 4. Usually n being larger than 30 is good enough. But if the underlying distribution is badly skewed or contaminated with occasional outliers, most statisticians

would prefer to have a larger sample size—say, 50 or more—and even that might not produce good results.

Def. (**Confidence interval, CI**) A $100(1-\alpha)\%$ confidence interval is a range believed with probability α for failing to contain the parameter of interest. The answer format: we are $\alpha\%$ confident that the mean number of xxx is between xxx and xxx.

E.g.

- 1. (CI for μ when σ is known) Let $z_{\frac{\alpha}{2}}$ be the $100(1-\frac{\alpha}{2})$ -percentile of N(0,1), then an **approximate 2-sided** $100(1-\alpha)\%$ CI for μ when σ is known is $\bar{x}\pm z_{\frac{\alpha}{2}}\frac{\sigma}{\sqrt{n}}$, derived from $Z=\frac{\bar{X}-\mu}{\sigma/\sqrt{n}}\sim N(0,1)$ approximately and thus $\mathbb{P}(-z_{\frac{\alpha}{2}}\leq Z\leq z_{\frac{\alpha}{2}})\approx 1-\alpha$.
- 2. (CI for μ when σ is unknown but n is large) use s as the unbiased estimator of σ^{**} , since $\frac{\bar{x}-\mu}{s/\sqrt{n}} \sim N(0,1)$.
- 3. (**One sided CI**) Let z_{α} be the $100(1-\alpha)$ -percentile of N(0,1), then an approximate 1-sided CI for μ when σ is known can be $(\bar{x}-z_{\alpha}\frac{\sigma}{\sqrt{n}},\infty)$, from $\mathrm{IP}(Z=\frac{\bar{X}-\mu}{\sigma/\sqrt{n}}\leq z_{\alpha})\approx 1-\alpha$.
- 4. (CI for μ when σ is unknown and $\mathbf n$ is small) Let $t_{n-1,\frac{\alpha}{2}}$ be the $100(1-\frac{\alpha}{2})$ -percentile of t_{n-1} , then when σ is unknown (need estimation s), it is $\bar x \pm t_{n-1,\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$, derived from $\mathrm{IP}(-t_{\frac{\alpha}{2}} \leq T = \frac{\bar X \mu}{s/\sqrt{n}} \leq t_{\frac{\alpha}{2}}) \approx 1 \alpha$.
- 5. (Sample size for target error when σ is known or n is large) The question is to find a sample size so that we have high confident that the mean is within $\bar{X} \pm \epsilon$. In other word, the CI is within $\bar{X} \pm \epsilon$, where ϵ is called the **maximum** error of the estimate. Then $\epsilon := z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$. When we want to know what sample size n is needed for certain amount of error, we can solve n from that equation.
- 6. (Sample size for target error when σ is unknown and n is small) Note that $\mathbb{P}(\mu \in \bar{X} \pm t_{n-1,\frac{\alpha}{2}} \frac{S}{\sqrt{n}}) = 1 \alpha$, and $\mathbb{E}(t_{n-1,\frac{\alpha}{2}} \frac{S}{\sqrt{n}})^2 = t_{n-1,\frac{\alpha}{2}} \frac{\sigma^2}{\sqrt{n}}$. Now if we hope $\mathbb{E}\epsilon^2 := \mathbb{E}(t_{n-1,\frac{\alpha}{2}} \frac{S}{\sqrt{n}})^2 \le \epsilon_u^2$, then we want $\frac{n}{t_{n-1,\frac{\alpha}{2}}^2} > \frac{s^2}{\epsilon_u^2}$. We can approximate it by conducting a preliminary experiment first to have sample variance, then get a sense of the n by using (5), and then improve the gap one observation at a time until it obeys the inequality.

7. (CI for $\mu_X - \mu_Y$ when σ_X, σ_Y is known or n is large) Suppose $X_1, \dots X_{n_X}$ and $Y_1, \dots Y_{n_Y}$ independently from two normal $N(\mu_X, \sigma_X^2), N(\mu_Y, \sigma_Y^2)$.

$$W:=X-Y, \sigma_W=\sqrt{rac{\sigma_X^2}{n_X}+rac{\sigma_Y^2}{n_Y}}, Z:=rac{ar{X}-ar{Y}-(\mu_X-\mu_Y)}{\sigma_W}\sim N(0,1).$$

- 8. (CI for $\mu_X \mu_Y$ when the common σ is unknown and n is small) Suppose $X_1, \ldots X_{n_X}$ and $Y_1, \ldots Y_{n_Y}$ independently from two normal $N(\mu_X, \sigma^2), N(\mu_Y, \sigma^2)$.
 - $ullet \ Z:=rac{ar{X}-ar{Y}-(\mu_X-\mu_Y)}{\sqrt{rac{\sigma^2}{n_Y}+rac{\sigma^2}{n_Y}}}\sim N(0,1).$
 - OTAH, $U:=rac{(n_X-1)S_X^2}{\sigma^2}+rac{(n_Y-1)S_Y^2}{\sigma^2}$ is the sum of two independent chi-square r.v. thus $U\sim X^2(n_X+n_Y-2)$.
 - The independence between the sample means and sample variances implies that $Z\perp U.$ Thus $T:=rac{Z}{\sqrt{rac{U}{n_X+n_Y-2}}}\sim t_{n_X+n_Y-2}$.
 - Then let $S_p:=\sqrt{rac{(n_X-1)S_X^2+(n_Y-1)S_Y^2}{n_X+n_Y-2}}, t_0:=t_{n_X+n_Y-2,lpha/2}$, we have the CI: $ar{X}-ar{Y}\pm t_0S_p\sqrt{rac{1}{n_X}+rac{1}{n_Y}}.$
- 9. (CI for proportion p) Treat proportion as event success with probability p. Assume independent $X_1, \ldots X_n \sim Bernoulli(p)$, unbiased MLE $\hat{p} = \frac{Y}{n} := \frac{\sum X_i}{n}$, where $Y \sim Bin(n, p)$.
 - By CLT, $\frac{\hat{p}-\mathbb{E}\hat{p}}{\sqrt{Var\hat{p}}}=\frac{\hat{p}-p}{\sqrt{p(1-p)/n}}\sim N(0,1)$ approximately when n is large enough.

$$egin{aligned} & extbf{I\!P}ig(-z_{rac{lpha}{2}} \leq rac{\hat{p}-p}{\sqrt{rac{p(1-p)}{n}}} \leq z_{rac{lpha}{2}}ig) = 1-lpha \ &\iff & extbf{I\!P}\left(rac{Y}{n} - z_{rac{lpha}{2}}\sqrt{rac{p(1-p)}{n}} \leq p \leq rac{Y}{n} + z_{rac{lpha}{2}}\sqrt{rac{p(1-p)}{n}}
ight) = 1-lpha \end{aligned}$$

• Make additional approximation by replace p in end points with $\hat{p} := \frac{Y}{n}$.

10. (CI for variance)
$$(\frac{(n-1)s^2}{X_{\frac{\alpha}{2}}^2}, \frac{(n-1)s^2}{X_{1-\frac{\alpha}{2}}^2})$$
. #TODO

2.3 Hypothesis testing

Standard procedure (go back to here later):

- 1. Formulate H_0, H_1 ;
- 2. Identify a test statistic and its dist under H_0 ;
- 3. Evaluate the test statistic by calculating a p-value;
- 4. Compare p-value to α . If less than, reject H_0 ;

5. State conclusion as "there is/isn't enough evidence to show that ...".

Partition the parameter space into two nonoverlapping disjoint subsets and then use the ovserved data to decide which set the parameter belongs to. The two cases are called **null hypothesis** H_0 and **alternative hypothesis** H_1 . We don't reject H_0 by default unless we have enough evidence to support H_1 , and thus reject H_0 .

The result of hypothesis testing is either to **reject or fail to reject** H_0 . Hence, it is necessary to partition the sample spaces into two part, say C, C'. The rejection region C for H_0 is called **critical region**. Often, the partition is specified in terms of the values of a statistic called the **test statistic** (should convert into standard normal/t-distribution before testing). **Remember that critical region and test statistic should be in standard dist.**

Def. **Type 1 error**: H_0 is true but $(X_i) \in C$, i.e. get rejected. The probability of this case, denoted as $\alpha := \mathbb{P}((X_i) \in C | H_0)$, is called the **significance level**. To control this error, we may pick suitable critical region for specified significance level that we want to achieve.

Def. **Type 2 error**: H_1 is true but $(X_i) \in C'$, i.e. fail to reject H_0 . The probability of this case is denoted as $\beta := \mathbb{P}((X_i) \in C'|H_1)$. $1-\beta$ is called the **power** of a test, which is a function of some parameters when H_1 is compound. Sometimes the power function is denoted as $\beta(\theta)$, $K(\mu)$. The whole function itself depends on the critical region, which depends on α . The more power a test is, the more likely it can reject the null hypothesis when H_1 is true.

Def. **p-value**: the tail-end probability, under H_0 , of the distribution of the test statistic beyond the observed value.

- The "beyond" means at least as extreme as the observed value so that away from H_0 in the direction(s) of H_1 . When H_1 is two-sided, it is generally taken to be $2 \min(\mathbb{P}(\text{smaller than the observed}), \mathbb{P}(\text{greater than the observed}))$.
- The smaller the p-value, the less we believe in H_0 . When p-value $< \alpha$, we reject H_0 . The specific p-value is more informative than simply put whether hypothesis is rejected under certain α , since the latter cannot tell what will happen when the α is different.
- In testing, asking p-value of test statistic to be less than α is the same as asking test statistic to be in the CI based on μ . In contrast, in last session,

we ask the chance of μ being in the CI based on the test statistic.

- Test for **one mean** with known variance or large n:
 - $H_0: \mu=\mu_0$, test statistics $Z:=rac{ar{X}-\mu_0}{\sigma/\sqrt{n}}\sim N(0,1)$.
 - When $H_1: \mu \neq \mu_0$, the critical region $|z| \geq z_{\alpha/2}$ or $|\bar{x} \mu_0| \geq z_{\frac{\alpha}{2}} \sigma / \sqrt{n}$.
 - When $H_1: \mu > \mu_0$, the critical region $z \geq z_\alpha$ or $\bar{x} \mu_0 \geq z_\alpha \sigma / \sqrt{n}$.
- Test for **one mean** with unknown variance and small n:
 - Simply adopt the t-distribution CI.
- Test for **one proportion**:
 - $H_0: p=p_0$, under which test statistics $Z:=\frac{\frac{Y}{n}-p_0}{\sqrt{p_0(1-p_0)/n}}\sim N(0,1)$, whose denominator is different from the CI. That one can also be used, yet produces approximately the same result.
 - When $H_1: \mu \neq \mu_0$, the critical region $|z| \geq z_{\alpha/2}$.
- Test for **equality of two means** with both dist approximately normal:

$$ullet \ Z:=rac{ar{X}-ar{Y}-(\mu_X-\mu_Y)}{\sqrt{rac{\sigma_X^2}{n_X}+rac{\sigma_Y^2}{n_Y}}}\sim N(0,1).$$

• CI =
$$(\bar{X} - \bar{Y}) \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_X^2}{n_X} + \frac{\sigma_Y^2}{n_Y}}$$
.

- Test for **equality of two means** with unknown variance and n is small:
 - In general case, to use t distribution, we can use conservative T:

$$df = min(n_1,n_2) - 1 ext{ or Welch's T: } df = \lfloor rac{(rac{s_X^2}{n_X} + rac{s_Y^2}{n_Y})^2}{rac{1}{n_X - 1}(rac{s_X^2}{n_X})^2 + rac{1}{n_V - 1}(rac{s_Y^2}{n_V})^2}
floor.$$

- **Pooled variance**: in the case that we assume $\sigma_X = \sigma_Y$, we can use s_{pooled} as the estimator for both, where $s_{pooled}^2 := \frac{(n_X-1)s_X^2+(n_Y-1)s_Y^2}{n_X+n_Y-2}$ and $df = n_X + n_Y 2$. This assumption is reasonable **if quotient between two sample standard deviation is at most 2**.
- Matched pairs: if the data are paried, i.e. $D_i = X_i Y_i$ are approximately random sample from normal. Then df=n-1, treat as one dist directly.
- Test for **two proportion**:

$$ullet \ Z := rac{(rac{Y_1}{n_1} - rac{Y_2}{n_2}) - (p_1 - p_2)}{\sqrt{rac{p_1(1 - p_1)}{n_1} + rac{p_2(1 - p_2)}{n_2}}} \sim N(0, 1).$$

ullet To get test statistic under $H_0: p_1=p_2$, estimate p_1,p_2 with $\hat{p}:=rac{Y_1+Y_2}{n_1+n_2}$. $Z:=rac{(rac{Y_1}{n_1}-rac{Y_2}{n_2})}{\sqrt{\hat{p}(1-\hat{p})(rac{1}{n_1}+rac{1}{n_1})}}\sim N(0,1).$

- $ullet X^2=rac{(n-1)S^2}{\sigma^2}\sim X_{(n-1)}^2$
- $ullet \ H_0: \sigma = \sigma_0, H_1: \sigma < \sigma_0, C = \{X^2 < X^2_{n-1,1-lpha}\}.$

2.4 Chi Square tests

2.4.1 Goodness of Fit

- A random sample of size $n = \sum_{i=1}^k Y_i$ is classified into k catagories / cells with frequency Y_i . Ground truth cell probabilities are p_i , while the sample probabilities are $\hat{p}_i = Y_i/n$. $H_0: \forall i, p_i = p_{i0}$.
- ullet Test statistic $X^2=\sumrac{(obs-exp)^2}{exp}=\sumrac{(Y_i-n*p_{i0})^2}{n*p_{i0}}\sim X_{k-1}^2.$
- This approximate test is only appropriate when all $Y_i \geq 5$.

2.4.2 Test of Homogeneity

- A random sample of size $n = \sum_{i=1}^r n_i = \sum_{i=1}^r \sum_{j=1}^c Y_{ij}$. $H_0: r$ sub-groups of a population share the same distribution of a single categorical variable with c categories. In other word, $\forall j, p_{ij}$ are all the same among i.
- $ullet \ X^2 = \sum rac{(obs-exp)^2}{exp} = \sum_{j=1}^c \sum_{i=1}^r rac{(Y_{ij} n_i * p_{ij})}{n_i * p_{ij}} \sim X^2_{r(c-1)}.$
- Under H_0 , estimate p_{ij} by the relative frequency $\hat{p}_{ij} := \frac{\sum_i Y_{ij}}{\sum_i n_i}$, by thinking they are part of a larger experiment in which $\sum_i Y_{ij}$ is the frequency.
- ullet Test statistic $X^2 = \sum rac{(obs-exp)^2}{exp} = \sum_{j=1}^c \sum_{i=1}^r rac{(Y_{ij}-n_i*\hat{p}_{ij})^2}{n_i*\hat{p}_{ij}} \sim X^2_{(r-1)(c-1)}.$
- By selecting k and defining categories, this method can be extended to test the equality of two distributions.

2.4.3 Test of Independence

- H_0 : two categorial variables, with r, c categories respectively, are independent in the population. Y_{ij} is the frequency of $A_i \cap B_j$, the intersection of two categories.
- In other word, $H_0: p_{ij}=p_{i\cdot}p_{\cdot j}$, which means $\mathbb{P}(A_i\cap B_j)=P(A_i)P(B_j)$. Under this hypothesis, $\hat{p}_{ij}=\hat{p_{i\cdot}}\hat{p_{\cdot j}}:=\frac{Y_{i\cdot}}{n}\frac{Y_{\cdot j}}{n}$.
- ullet Test statistic $X^2 = \sum rac{(obs-exp)^2}{exp} = \sum_{j=1}^c \sum_{i=1}^r rac{(Y_{ij}-n*\hat{p}_{ij})^2}{n*\hat{p}_{ij}} \sim X^2_{(r-1)(c-1)}.$
- The same degree of freedom eventually.

Multivariate normal

- 1. Standard normal: $\underline{Z} \sim N(\underline{0}, I)$.
- 2. (Stretch) $\underline{X} = D\underline{Z}$ where $D = diag(\sigma_1, \sigma_2, \ldots)$, then $\underline{X} \sim N(\underline{0}, D^2)$.
- 3. (Rotate) $\underline{Y} = Q\underline{X}$ where $Q^TQ = I$ (orthogonal), then $\underline{Y} \sim N(\underline{0}, V)$. Given D to find V and Q, consider that by (Covariance-4),

$$V = V[Y] = V[QX] = QV[X]Q^T = QD^2Q^T$$

Eigenvalue decomposition solves this problem. Say $eigen(V) = \lambda_1, \ldots$, then $\lambda_i = \sigma_i^2$.

Kullback-Leibler divergence

a measure of how one probability distribution P is different from a second, reference probability distribution Q.