# On these review notes

1. You are responsible for the correctness of all of the formulae on this review sheet. (There are undoubtedly ytopgraphical errors :-).

2. You should know, and understand, everything in these review notes.

3. The exam format will be a series of multiple choice, short answer questions and R codes. Tedious calculations will be avoided.

4. You can bring a non-fancy (you know what I mean) scientific calculator. It must be able to take logs and raise numbers to exponents.

5. You can bring in one sheet of $8.5 \times 11$ paper filled, front and back, with formulae and notes.

# 1 Random variables

1. A **random variable** is a function from $\Omega$ to the real numbers. A random variable is a random number that is the result of an experiment governed by a probability distribution.

2. A **Bernoulli** random variable is one that takes the value 1 with probability $p$ and 0 with probability $(1 - p)$. That is, $P(X = 1) = p$ and $P(X = 0) = 1 - p$.

3. A **probability mass function** (pmf) is a function that yields the various probabilities associated with a random variable. For example, the probability mass function for a Bernoulli random variable is $f(x) = p^x(1-p)^{1-x}$ for $x = 0, 1$ as this yields $p$ when $x = 1$ and $(1 - p)$ when $x = 0$.

4. The **expected value** or (population) **mean** of a discrete random variable, $X$, with pmf $f(x)$ is
$$\mu = E[X] = \sum_x x f(x).$$
The mean of a Bernoulli variable is then $1f(1) + 0f(0) = p$.

5. The **variance** of any random variable, $X$, (discrete or continuous) is
$$\sigma^2 = E\left[(X - \mu)^2\right] = E[X^2] - E[X]^2.$$
The latter formula being the most convenient for computation. The variance of a Bernoulli random variable is $p(1 - p)$.

6. The (population) **standard deviation**, $\sigma$, is the square root of the variance.

7. A **Binomial** random variable, $X$, is obtained as the sum of $n$ Bernoulli random variables and has pmf
$$P(X = k) = \left( \begin{array}{c} n \\ k \end{array} \right) p^k(1-p)^{n-k}.$$
Binomial random variables have expected value $np$ and variance $np(1 - p)$.

8. An uniform random variable, $X$. The expected value and variance of X.

# 2 Continuous random variables

1. **Continuous** random variables take values on a continuum.

2. The probability that a continuous random variable takes on any specific value is 0.

3. Probabilities associated with continuous random variables are governed by **probability density functions** (pdfs). Areas under probability density functions correspond to probabilities. For example, if $f$ is a pdf corresponding to random variable $X$, then

$$P(a \leq X \leq b) = \int_a^b f(x)dx.$$

To be a pdf, a function must be positive and integrate to 1. That is, $\int_{-\infty}^{\infty} f(x)dx = 1$

4. If $h$ is a positive function such that $\int_{-\infty}^{\infty} h(x)dx \leq \infty$ then $f(x) = h(x)/\int_{-\infty}^{\infty} h(x)dx$ is a valid density. Therefore, if we only know a density up to a constant of proportionality, then we can figure out the exact density.

5. The expected value, or mean, of a continuous random variable, $X$, with pdf $f$, is

$$\mu = E[X] = \int_{-\infty}^{\infty} tf(t)dt.$$

6. The variance is $\sigma^2 = E[(X - \mu)^2] = E[X^2] - E[X]^2$.

7. The **distribution function**, say $F$, corresponding to a random variable $X$ with pdf, $f$, is

$$P(X \leq x) = F(x) = \int_{-\infty}^x f(t)dt.$$

(Note the common convention that $X$ is used when describing an unobserved random variable while $x$ is for specific values.)

8. The $p^{th}$ **quantile** (for $0 \leq p \leq 1$), say $X_p$, of a distribution function, say $F$, is the point so that $F(X_p) = p$. For example, the $.025^{th}$ quantile of the standard normal distribution is -1.96.

# 3 Properties of expected values and variances

The following properties hold for all expected values (discrete or continuous)

1. Expected values commute across sums: $E[X + Y] = E[X] + E[Y]$.

2. Multiplicative and additive constants can be pulled out of expected values $E[cX] = cE[X]$ and $E[c + X] = c + E[X]$.

3. For independent random variables, $X$ and $Y$, $E[XY] = E[X]E[Y]$.

4. In general, $E[h(X)] \neq h(E[X])$.

5. Variances commute across sums *for independent variables* $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$.

6. Multiplicative constants are squared when pulled out of variances $\text{Var}(cX) = c^2\text{Var}(X)$.

7. Additive constants do not change variances: $\text{Var}(c + X) = \text{Var}(X)$.

8. $E(\sum a_i Y_i) = \sum a_i E(Y_i)$, and $Var(\sum a_i Y_i)$.

# 4   The normal distribution

a. The **Bell curve** or **normal** or **Gaussian** density is the most common density. It is speci-
   fied by its mean, $\mu$, and variance, $\sigma^2$. The density is given by $f(x) = (2\pi\sigma^2)^{-1/2}\exp\{-(x-\mu)^2/2\sigma^2\}$. We write $X \sim N(\mu, \sigma^2)$ to denote that $X$ is normally distributed with mean $\mu$
   and variance $\sigma^2$.

b. The **standard normal** density, labeled $\phi$, corresponds to a normal density with mean
   $\mu = 0$ and variance $\sigma^2 = 1$.

$$\phi(z) = (2\pi)^{-1/2}\exp\{-z^2/2\}.$$

   The standard normal distribution function is usually labeled $\Phi$.

c. If $f$ is the pdf for a $N(\mu, \sigma^2)$ random variable, $X$, then note that $f(x) = \phi\{(x-\mu)/\sigma\}/\sigma$.
   Correspondingly, if $F$ is the associated distribution function for $X$, then $F(x) = \Phi\{(x-\mu)/\sigma\}$.

d. If $X$ is normally distributed with mean $\mu$ and variance $\sigma^2$ then the random variable $Z = (X-\mu)/\sigma$ is standard normally distributed. Taking a random variable subtracting its
   mean and dividing by its standard deviation is called "standardizing" a random variable.

e. If $Z$ is standard normal then $X = \mu + Z\sigma$ is normal with mean $\mu$ and variance $\sigma^2$.

f. 68%, 95% and 99% of the mass of any normal distribution lies within 1, 2 and 3
   (respectively) standard deviations from the mean.

g. $Z_\alpha$ refers to the $\alpha^{th}$ quantile of the standard normal distribution. $Z_{.90}$, $Z_{.95}$, $Z_{.975}$ and
   $Z_{.99}$ are 1.28, 1.645, 1.96 and 2.32.

h. Sums and means of normal random variables are normal (regardless of whether or not
   they are independent). You can use the rules for expectations and variances to figure
   out $\mu$ and $\sigma$.

i. The sample standard deviation of iid normal random variables, appropriated normal-
   ized, is a Chi-squared random variable (see below).

# 5   Sample means and variances

Throughout this section let $X_i$ be a collection of iid random variables with mean $\mu$ and
variance $\sigma^2$.

1. We say random variables are **iid** if they are independent and identically distributed.

2. For random variables, $X_i$, the **sample mean** is $\bar{X} = \sum_{i=1}^{n} X_i/n$.

3. $E[\bar{X}] = \mu = E[X_i]$ (does not require the independence or constant variance).

4. If the $X_i$ are iid with variance $\sigma^2$ then $\text{Var}(\bar{X}) = \text{Var}(X_i)/n = \sigma^2/n$.

5. The **sample variance** is defined to be

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

6. $\sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n X_i^2 - n\bar{X}^2$ is a shortcut formula for the numerator.

7. $\sigma/\sqrt{n}$ is called the **standard error** of $\bar{X}$. The estimated standard error of $\bar{X}$ is $S/\sqrt{n}$. Do not confuse dividing by this $\sqrt{n}$ with dividing by $n - 1$ in the calculation of $S^2$.

8. An estimator is **unbiased** if its expected value equals the parameter it is estimating.

9. $E[S^2] = \sigma^2$, which is why we divide by $n - 1$ instead of $n$. That is, $S^2$ is unbiased. However, dividing by $n - 1$ rather than $n$ does increase the variance of this estimator slightly, $\mathrm{Var}(S^2) \geq \mathrm{Var}((n - 1)S^2/n)$.

10. If the $X_i$ are normally distributed with mean $\mu$ and variance $\sigma^2$, then $\bar{X}$ is normally distributed with mean $\mu$ and variance $\sigma^2/n$.

11. The **Central Limit Theorem**. If the $X_i$ are iid with mean $\mu$ and (finite) variance $\sigma^2$ then

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

will limit to a standard normal distribution. The result is true for small sample sizes, if the $X_i$ iid normally distributed.

12. If we replace $\sigma$ with $S$; that is,

$$Z = \frac{\bar{X} - \mu}{S/\sqrt{n}},$$

then $Z$ still limits to a standard normal. If the $X_i$ are iid normally distributed, then $Z$ follows the Students $T$ distribution for small $n$.

# 6 Confidence intervals for a mean using the CLT.

1. Using the CLT, we know that

$$P\left(-Z_{1-\alpha/2} \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq Z_{1-\alpha/2}\right) = 1 - \alpha$$

for large $n$. Solving the inequalities for $\mu$, we calculated that in repeated sampling, the interval

$$\bar{X} \pm Z_{1-\alpha/2} \frac{S}{\sqrt{n}}$$

will contain $\mu$ $100(1 - \alpha)\%$ of the time.

2. The probability that $\mu$ is in an observed confidence interval is either 1 or 0. The correct interpretation is that in repeated sampling, the interval we obtain will contain $\mu$ $100(1-\alpha)\%$ of the time. (Assumes that the CLT has kicked in).

3. As $n$ increases, the interval gets narrower.

4. As $S$ increases, the interval gets wider.

5. As the **confidence level**, $(1-\alpha)$, increases, the interval gets wider.

6. Fixing the confidence level controls the **accuracy** of the interval. A 95% interval has 95% coverage regardless of the sample size. (Again, assuming that the CLT has kicked in.) Increasing $n$ will improve the precision (width) of the interval.

7. Prior to conducting a study, you can fix the **margin of error** (half width), say $\delta$, of the interval by setting $n = (Z_{1-\alpha/2}\sigma/\delta)^2$. Round up. Requires an estimate of $\sigma$.

# 7   Confidence intervals for a variance and T confidence intervals

1. If $Z$ is standard normal and $X$ is and independent Chi-squared with $df$ degrees of freedom then $\frac{Z}{\sqrt{X/df}}$ follows what is called a Student's $T$ distribution with $df$ degrees of freedom.

2. The Student's $T$ density looks like a normal density with heavier tails (so it looks more squashed down).

3. By the previous item, if the $X_i$ are iid $\mathrm{N}(\mu, \sigma^2)$ then

$$Z = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

follows a Student's $T$ distribution with $(n-1)$ degrees of freedom. Therefore if $t_{n-1,\alpha}$ is the $\alpha^{th}$ quantile of the Student's $T$ distribution then

$$\bar{X} \pm t_{n-1,1-\alpha/2}\frac{S}{\sqrt{n}}$$

is a $100(1-\alpha)\%$ confidence interval for $\mu$.

4. The Student's $T$ confidence interval assumes normality of the $X_i$. However, the $T$ distribution has quite heavy tails and so the interval is conservative and works well in many situations.

5. For large sample sizes, the Student's $T$ and CLT based intervals are nearly the same because the Student's $T$ quantiles become more and more like standard normal quantiles as $n$ increases.

# 8 The bootstrap

1. The (non-parametric) **bootstrap** can be used to calculate **percentile bootstrap confidence intervals**.

2. The **bootstrap principle** is to use the empirical distribution defined by the data to obtain an estimate of the sampling distribution of a statistic. In practice the bootstrap principle is always executed by **resampling (with replacement)** from the observed data.

3. Assume that we have $n$ data points. The bootstrap obtains a confidence interval by sampling $m$ complete data sets by drawing with replacement from the original data. The statistic of interest, say the median, is applied to all $m$ of the resampled data sets, yielding $m$ medians. The percentile confidence interval is obtained by taking the $\alpha/2$ and $1 - \alpha/2$ quantiles of the $m$ medians.

4. Make sure you do enough resamples so that your confidence interval has stabilized.

5. Bootstrap intervals are interpreted the same as frequentist intervals.

6. To guarantee coverage, the bootstrap interval requires large sample sizes.

7. There are improvements to the percentile method that are not covered in this class.

# 9 Hypothesis testing for a single mean

1. The null, or status quo, hypothesis is labeled $H_0$, the alternative $H_a$ or $H_1$ or $H_2$ ...

2. A **type I error** occurs when we falsely reject the null hypothesis. The probability of a type I error is usually labeled $\alpha$.

3. A **type II error** occurs when we falsely fail to reject the null hypothesis. A type II error is usually labeled $\beta$.

4. A **Power** is the probability that we correctly reject the null hypothesis, $1 - \beta$.

5. The $Z$ test for $H_0 : \mu = \mu_0$ versus $H_1 : \mu < \mu_0$ or $H_2 : \mu \neq \mu_0$ or $H_3 : \mu > \mu_0$ constructs a test statistic $TS = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$ and rejects the null hypothesis when

   $H_1$  $TS \leq -Z_{1-\alpha}$

   $H_2$  $|TS| \geq Z_{1-\alpha/2}$

   $H_3$  $TS \geq Z_{1-\alpha}$

   respectively.

6. The $Z$ test requires the assumptions of the CLT and for $n$ to be large enough for it to apply.

7. If $n$ is small, then a Student's $T$ test is performed exactly in the same way, with the normal quantiles replaced by the appropriate Student's $T$ quantiles and $n - 1$ df.

8. Tests define confidence intervals by considering the collection of values of $\mu_0$ for which you fail to reject a two sided test. This yields exactly the $T$ and $Z$ confidence intervals respectively.

9. Conversely, confidence intervals define tests by the rule where one rejects $H_0$ if $\mu_0$ is *not in* the confidence interval.

10. A **P-value** is the probability of getting evidence as extreme or more extreme than we actually got under the null hypothesis. For $H_3$ above, the P-value is calculated as $P(Z \geq TS_{obs}|\mu = \mu_0)$ where $TS_{obs}$ is the observed value of our test statistic. To get the P-value for $H_2$, calculate a one sided P-value and double it.

11. The P-value is equal to the **attained significance level**. That is, the smallest $\alpha$ value for which we would have rejected the null hypothesis. Therefore, rejecting the null hypothesis if a P-value is less than $\alpha$ is the same as performing the rejection region test.

12. The power of a $Z$ test for $H_3$ is given by the formula (know how this is obtained)

$$P(TS > Z_{1-\alpha}|\mu = \mu_1) = P\left(Z \geq \frac{\mu_0 - \mu_1}{\sigma/\sqrt{n}} + Z_{1-\alpha}\right).$$

Notice that power required a value for $\mu_1$, the value under the null hypothesis. Correspondingly for $H_1$ we have

$$P\left(Z \leq \frac{\mu_0 - \mu_1}{\sigma/\sqrt{n}} - Z_{1-\alpha}\right).$$

For $H_2$, the power is approximately the appropriate one sided power using $\alpha/2$.

13. Some facts about power.

    a. Power goes up as $\alpha$ goes down.

    b. Power of a one sided test is greater than the power of the associated two sided test.

    c. Power goes up as $\mu_1$ gets further away from $\mu_0$.

    d. Power goes up as $n$ goes up.

14. The prior formula can be used to calculate the sample size. For example, using the power formula for $H_1$, setting $Z_{1-\beta} = \frac{\mu_0 - \mu_1}{\sigma/\sqrt{n}} - Z_{1-\alpha}$ yields

$$n = \frac{(Z_{1-\beta} + Z_{1-\alpha})^2\sigma^2}{(\mu_0 - \mu_1)^2},$$

which gives the sample size to have power $= 1 - \beta$. This formula applies for $H_3$ also. For the two sided test, $H_2$, replace $\alpha$ by $\alpha/2$.

15. Determinants of sample size.

    a. $n$ gets larger as $\alpha$ gets smaller.

    b. $n$ gets larger as the power you want gets larger.

    c. $n$ gets lager the closer $\mu_1$ is to $\mu_0$.

16. Paired T-test

17. Use simulation to calculate type I error rate and power

# 10 Group comparisons

1. For group comparisons, make sure to differentiate whether or not the observations are paired (or matched) versus independent.

2. For paired comparisons for continuous data, one strategy is to calculate the **differences** and use the methods for testing and performing hypotheses regarding a single mean. The resulting tests and confidence intervals are called **paired Student's** $T$ tests and intervals respectively.

3. For independent groups of iid variables, say $X_i$ and $Y_i$, *with a constant variance $\sigma^2$ across* groups
$$Z = \frac{\bar{X} - \bar{Y} - (\mu_x - \mu_y)}{S_p\sqrt{\frac{1}{n_x} + \frac{1}{n_y}}}$$

   limits to a standard normal random variable as both $n_x$ and $n_y$ get large. Here
$$S_p^2 = \frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2}$$

   is the **pooled estimate** of the variance. Obviously, $\bar{X}$, $S_x$, $n_x$ are the sample mean, sample standard deviation and sample size for the $X_i$ and $\bar{Y}$, $S_y$ and $n_y$ are defined analogously.

4. If the $X_i$ and $Y_i$ happen to be normal, then $Z$ follows the Student's $T$ distribution with $n_x + n_y - 2$ degrees of freedom.

5. Therefore a $(1 - \alpha) \times 100\%$ confidence interval for $\mu_y - \mu_x$ is

$$\bar{Y} - \bar{X} \pm t_{n_x+n_y-2,1-\alpha/2}S_p\left(\frac{1}{n_x} + \frac{1}{n_y}\right)^{1/2}$$

6. The statistic

$$\frac{\bar{Y} - \bar{X} - (\mu_y - \mu_x)}{\left(\frac{\sigma_x^2 1}{n_x} + \frac{\sigma_y^2}{n_y}\right)^{1/2}}$$

approximately follows Gosset's $T$ distribution with degrees of freedom equal to

$$\frac{\left(S_x^2/n_x + S_y^2/n_y\right)^2}{\left(\frac{S_x^2}{n_x}\right)^2/(n_x - 1) + \left(\frac{S_y^2}{n_y}\right)^2/(n_y - 1)}$$

# 11   Non-parametric Tests, Permutation Test

1. Specify hypotheses for each test.

2. The assumptions of each test.

3. The power comparison of two-sample tests

# 12   Simple linear regression

1. Simple linear regression models

2. Least estimations

3. Normal equation

4. Estimates for $\beta_0$, $\beta_1$ and $\sigma_2$

5. Properties of the residuals

6. Confidence intervals for the $\beta$ estimates, $\hat{Y}$.

7. Prediction Intervals