

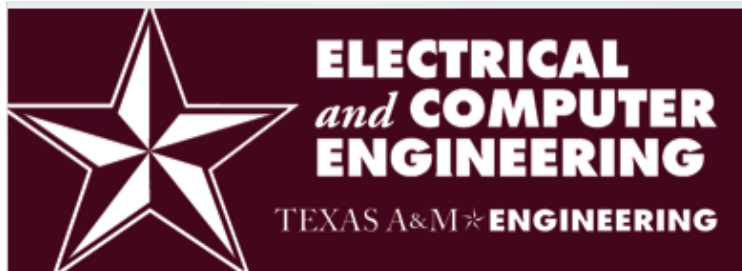
Materials Informatics

Lecture 4: Review of Probability and Statistics

Ulisses Braga Neto

Department of Electrical and Computer Engineering

Texas A&M University



Random Variables

- A *random variable* X is a (measurable) function $X : S \rightarrow \mathbb{R}$, that is, it assigns to each outcome of the experiment a real number.
- The *probability distribution function* (PDF) of a r.v. X is the function $F_X : \mathbb{R} \rightarrow [0, 1]$ given by

$$F_X(a) = P(\{X \leq a\}), \quad a \in \mathbb{R}.$$

- Properties of a PDF:
 1. F_X is non-decreasing: $a \leq b \Rightarrow F_X(a) \leq F_X(b)$.
 2. $\lim_{a \rightarrow -\infty} F_X(a) = 0$ and $\lim_{a \rightarrow +\infty} F_X(a) = 1$
 3. F_X is right-continuous: $\lim_{b \rightarrow a+} F_X(b) = F_X(a)$.

Probability Density Functions

- The notion of a *probability density function* (pdf) is fundamental in probability theory. However, it is a secondary notion to that of a PDF. In fact, all r.v.'s must have a PDF, but not all r.v.'s have a pdf.
- If F_X is everywhere continuous and differentiable, then X is said to be a *continuous* r.v. and the pdf f_X of X is given by:

$$f_X(a) = \frac{dF_X}{dx}(a), \quad a \in \mathbb{R}.$$

- Probability statements about X can then be made in terms of integration of f_X . For example,

$$F_X(a) = \int_{-\infty}^a f_X(x)dx, \quad a \in \mathbb{R}$$

$$P(a \leq X \leq b) = \int_a^b f_X(x)dx, \quad a, b \in \mathbb{R}$$

Useful Continuous R.V.'s

- Uniform(a, b)

$$f_X(x) = \frac{1}{b-a}, \quad a < x < b.$$

- The univariate Gaussian($\mu, \sigma > 0$):

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad x \in \mathbb{R}.$$

- Exponential($\lambda > 0$):

$$f_X(x) = \lambda e^{-\lambda x}, \quad x \geq 0.$$

- Gamma($\lambda > 0, t > 0$):

$$f_X(x) = \frac{\lambda e^{-\lambda x} (\lambda x)^{t-1}}{\Gamma(t)}, \quad x \geq 0.$$

- Beta(a, b):

$$f_X(x) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1}, \quad 0 < x < 1.$$

Probability Mass Function

- If F_X is not everywhere differentiable, then X does not have a pdf.
- A useful particular case is that where X assumes countably many values and F_X is a *staircase function*. In this case, X is said to be a *discrete* r.v.
- One defines the *probability mass function* (PMF) of a discrete r.v. X to be:

$$p_X(a) = P(\{X = a\}) = F_X(a) - F_X(a^-), \quad a \in \mathbb{R}.$$

(Note that for any continuous r.v. X , $P(\{X = a\}) = 0$ so there is no PMF to speak of.)

Useful Discrete R.V.'s

● Bernoulli(p):

$$p_X(0) = P(\{X = 0\}) = 1 - p$$

$$p_X(1) = P(\{X = 1\}) = p$$

● Binomial(n, p):

$$p_X(k) = P(\{X = k\}) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k = 0, 1, \dots, n$$

● Poisson (λ):

$$p_X(k) = P(\{X = k\}) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k = 0, 1, \dots$$

● Geometric(p):

$$p_X(k) = P(\{X = k\}) = (1 - p)^{k-1} p, \quad k = 1, 2, \dots$$

Expectation

- Expectation is a fundamental concept in probability theory, which has to do with the intuitive concept of "averages."
- The mean value of a r.v. X is an average of its values weighted by their probabilities. If X is a continuous r.v., this corresponds to:

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx$$

If X is discrete, then this can be written as:

$$E[X] = \sum_{i: p_X(x_i) > 0} x_i p_X(x_i)$$

Expectation of Functions of R.V.'s

- Given a r.v. X , and a (measurable) function $g : \mathbb{R} \rightarrow \mathbb{R}$, then $g(X)$ is also a r.v. If there is a pdf, it can be shown that

$$E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx$$

If X is discrete, then this becomes:

$$E[g(X)] = \sum_{i: p_X(x_i) > 0} g(x_i) p_X(x_i)$$

- Immediate corollary: $E[aX + c] = aE[X] + c$.
- (Joint R.V.'s)** If $g : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is measurable, then (continuous case):

$$E[g(X, Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{XY}(x, y) dx dy$$

Other Properties of Expectation

- Linearity: $E[a_1X_1 + \dots a_nX_n] = a_1E[X_1] + \dots a_nE[X_n]$.
- if X and Y are *uncorrelated*, then $E[XY] = E[X]E[Y]$ (independence always implies uncorrelatedness, but the converse is true only in special cases; e.g. jointly Gaussian or multinomial random variables).
- If $X \geq Y$ then $E[X] \geq E[Y]$.
- If X is non-negative,

$$E[X] = \int_0^{\infty} P(\{X > x\}) dx$$

- Markov's Inequality: If X is non-negative, then for $a > 0$,

$$P(\{X \geq a\}) \leq \frac{E[X]}{a}$$

- Cauchy-Schwarz Inequality: $E[XY] \leq \sqrt{E[X^2]E[Y^2]}$

Conditional Expectation

- If X and Y are jointly continuous r.v.'s and $f_Y(y) > 0$, we define:

$$E[X|Y = y] = \int_{-\infty}^{\infty} x f_{X|Y}(x, y) dx = \int_{-\infty}^{\infty} x \frac{f_{XY}(x, y)}{f_Y(y)} dx$$

If X, Y are jointly discrete and $p_Y(y_j) > 0$, this can be written as:

$$E[X|Y = y_j] = \sum_{i: p_X(x_i) > 0} x_i p_{X|Y}(x_i, y_j) = \sum_{i: p_X(x_i) > 0} x_i \frac{p_{XY}(x_i, y_j)}{p_Y(y_j)}$$

- Conditional expectations have all the properties of usual expectations. For example,

$$E \left[\sum_{i=1}^n X_i | Y = y \right] = \sum_{i=1}^n E[X_i | Y = y]$$

$E[Y|X]$ is a Random Variable

- Given a r.v. X , the mean $E[X]$ is a deterministic parameter.
- Now, $E[X|Y]$ is not random w.r.t. to X , but it is a function of the r.v. Y , so it is also a r.v. One can show that its mean is precisely $E[X]$:

$$E[E[X|Y]] = E[X]$$

What this says in the continuous case is:

$$E[X] = \int_{-\infty}^{\infty} E[X|Y = y] f_Y(y) dy$$

and, in the discrete case,

$$E[X] = \sum_{i: p_Y(y_i) > 0} E[X|Y = y_i] P(\{Y = y_i\})$$

- Computing $E[X|Y = y]$ first is often easier than finding $E[X]$ directly.

Variance

- The mean $E[X]$ is a good “guess” of the value of a r.v. X , but by itself it can be misleading. The *variance* $\text{Var}(X)$ says how the values of X are spread around the mean $E[X]$:

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

- Property: $\text{Var}(aX + c) = a^2\text{Var}(X)$
- Chebyshev's Inequality: For any $\epsilon > 0$,

$$P(\{|X - E[X]| \geq \epsilon\}) \leq \frac{\text{Var}(X)}{\epsilon^2}$$

Conditional Variance

- If X and Y are jointly-distributed r.v.'s , we define:

$$\text{Var}(X|Y) = E[(X - E[X|Y])^2|Y] = E[X^2|Y] - (E[X|Y])^2$$

- Conditional Variance Formula:

$$\text{Var}(X) = E[\text{Var}(X|Y)] + \text{Var}(E[X|Y])$$

This breaks down the *total variance* in a "within-rows" component and a "across-rows" component.

Covariance

- Is $\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2)$?
- The covariance of two r.v.'s X and Y is given by:

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$$

- Two variables X and Y are uncorrelated if and only if $\text{Cov}(X, Y) = 0$. Jointly Gaussian X and Y are independent if and only if they are uncorrelated (in general, independence implies uncorrelatedness but not vice-versa).
- It can be shown that

$$\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) + 2\text{Cov}(X_1, X_2)$$

- So the variance is distributive over sums if all variables are *pair-wise uncorrelated*.

Correlation Coefficient

- The correlation coefficient ρ between two r.v.'s X and Y is given by:

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

- Properties:

1. $-1 \leq \rho(X, Y) \leq 1$
2. X and Y are uncorrelated iff $\rho(X, Y) = 0$.
3. (Perfect linear correlation):

$$\rho(X, Y) = \pm 1 \Leftrightarrow Y = a \pm bX, \text{ where } b = \frac{\sigma_y}{\sigma_x}$$

Vector Random Variables

- A vector r.v. or *random vector* $\mathbf{X} = (X_1, \dots, X_d)$ takes values in \mathbb{R}^d . Its distribution is the joint distribution of the component r.v.'s.
- The mean of \mathbf{X} is the vector $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)$, where $\mu_i = E[X_i]$.
- The *covariance matrix* Σ is a $d \times d$ matrix given by:

$$\Sigma = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T]$$

where $\Sigma_{ii} = \text{Var}(X_i)$ and $\Sigma_{ij} = \text{Cov}(X_i, X_j)$.

Properties of the Covariance Matrix

- Matrix Σ is real symmetric and thus diagonalizable:

$$\Sigma = UDU^T$$

where U is the matrix of eigenvectors and D is the diagonal matrix of eigenvalues.

- All eigenvalues are nonnegative (Σ is *positive semi-definite*). In fact, except for “degenerate” cases, all eigenvalues are positive, and so Σ is invertible (Σ is said to be *positive definite* in this case).
- (*Whitening or Mahalanobis transformation*) It is easy to check that the random vector

$$\mathbf{Y} = \Sigma^{-\frac{1}{2}}(\mathbf{X} - \boldsymbol{\mu}) = D^{-\frac{1}{2}}U^T(\mathbf{X} - \boldsymbol{\mu})$$

has zero mean and covariance matrix \mathbf{I}_d (so that all components of \mathbf{Y} are zero-mean, unit-variance, and uncorrelated).

The Multivariate Gaussian

- The multivariate Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix Σ (assuming Σ invertible, so that also $\det(\Sigma) > 0$) corresponds to the multivariate pdf

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^p \det(\Sigma)}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

We write $X \sim \mathcal{N}_p(\boldsymbol{\mu}, \Sigma)$.

- The multivariate gaussian has elliptical contours of the form

$$(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = c^2, \quad c > 0$$

The axes of the ellipsoids are given by the eigenvectors of Σ and the length of the axes are proportional to its eigenvalues.

Properties of The Multivariate Gaussian

- The density of each component X_i is univariate gaussian $\mathcal{N}(\mu_i, \Sigma_{ii})$.
- The components of $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$ are independent *if and only if* they are uncorrelated, i.e., Σ is a diagonal matrix.
- The whitening transformation $\mathbf{Y} = \Sigma^{-\frac{1}{2}}(\mathbf{X} - \mu)$ produces another multivariate gaussian $\mathbf{Y} \sim \mathcal{N}(0, \mathbf{I}_p)$.
- In general, if \mathbf{A} is a nonsingular $p \times p$ matrix and \mathbf{c} is a p -vector, then $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{c} \sim \mathcal{N}_p(\mathbf{A}\mu + \mathbf{c}, \mathbf{A}\Sigma\mathbf{A}^T)$.
- The r.v.'s $\mathbf{A}\mathbf{X}$ and $\mathbf{B}\mathbf{X}$ are independent if and only if $\mathbf{A}\Sigma\mathbf{B}^T = 0$.
- If \mathbf{Y} and \mathbf{X} are jointly Gaussian, then the distribution of \mathbf{Y} given \mathbf{X} is again Gaussian.
- The regression $E[\mathbf{Y}|\mathbf{X}]$ is a linear function of \mathbf{X} .

Variability/Reproducibility

- Results can be expected to have **reproducibility** only if **variability** is taken into account in **experimental design**.
- Even if the subject and assay are clearly specified, there are two sources of variability
 - **Technical Variability** : it comprises random changes from an experiment (microarray) to the next.
 - **Populational Variability** : it comprises random changes from an individual or sample to the next.

Statistical Significance

- Suppose there are two conditions A and B (e.g., ordinary vs. superconductivity) under study, and one measurement Y (e.g. a QSPR).
- Suppose further that there are n replicate specimens, $n/2$ under condition A and $n/2$ under condition B (since A and B are represented by the same number of samples, this is called a **balanced design**).

Statistical Significance - II

- The question of interest is:
"Based on the set of n replicates, can we conclude that Y is significantly different between A and B ?"
- To examine this question, let us assume that there is indeed a difference between A and B . Suppose for example that in truth we have

$Y = 100$, under A

$Y = 200$, under B .

Statistical Significance - III

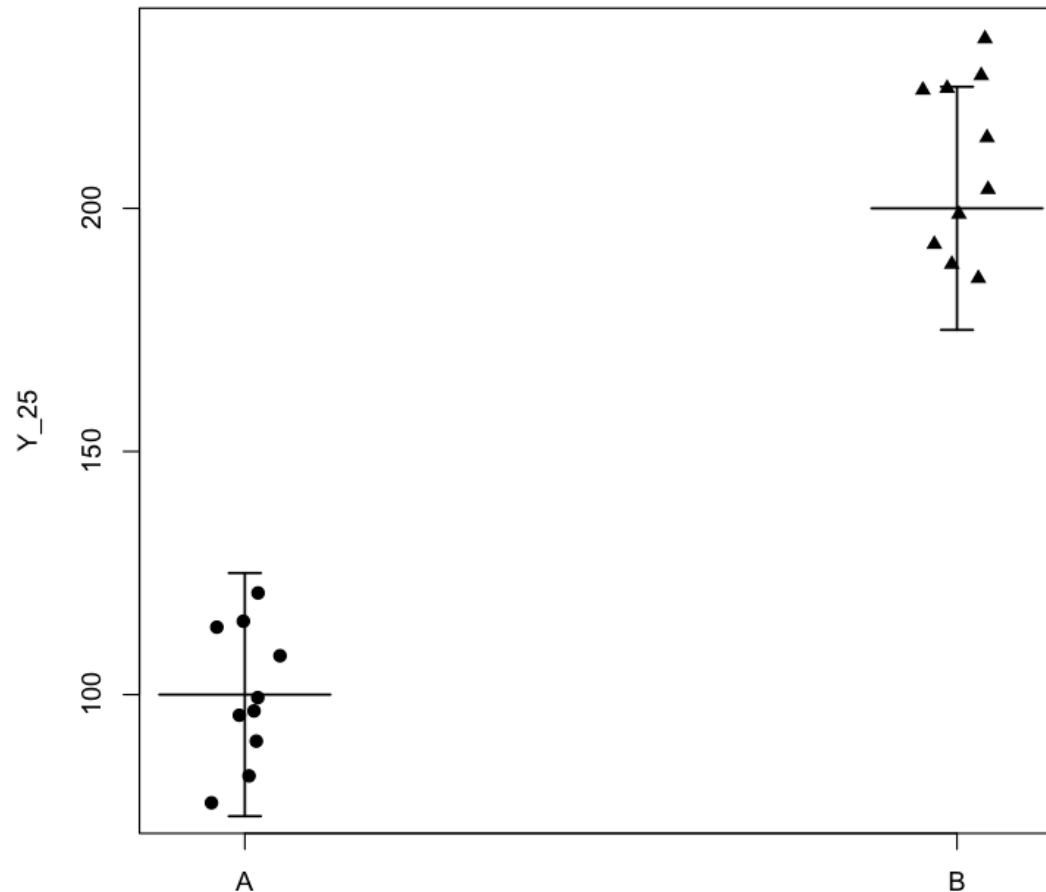
- The **fold-change** is $200/100 = 2$.
- If there were no variability, then with $n=2$ (1 replicate for each condition) we would be able to conclude there was a difference.
- But clearly there will be some variability, both technical and populational. Let us model this by using **Gaussian** distributions

$$Y \sim N(100, \sigma^2) \quad \text{under A}$$

$$Y = N(200, \sigma^2) \quad \text{under B}$$

where σ^2 is the variance (assumed equal, for the moment) in A and B.

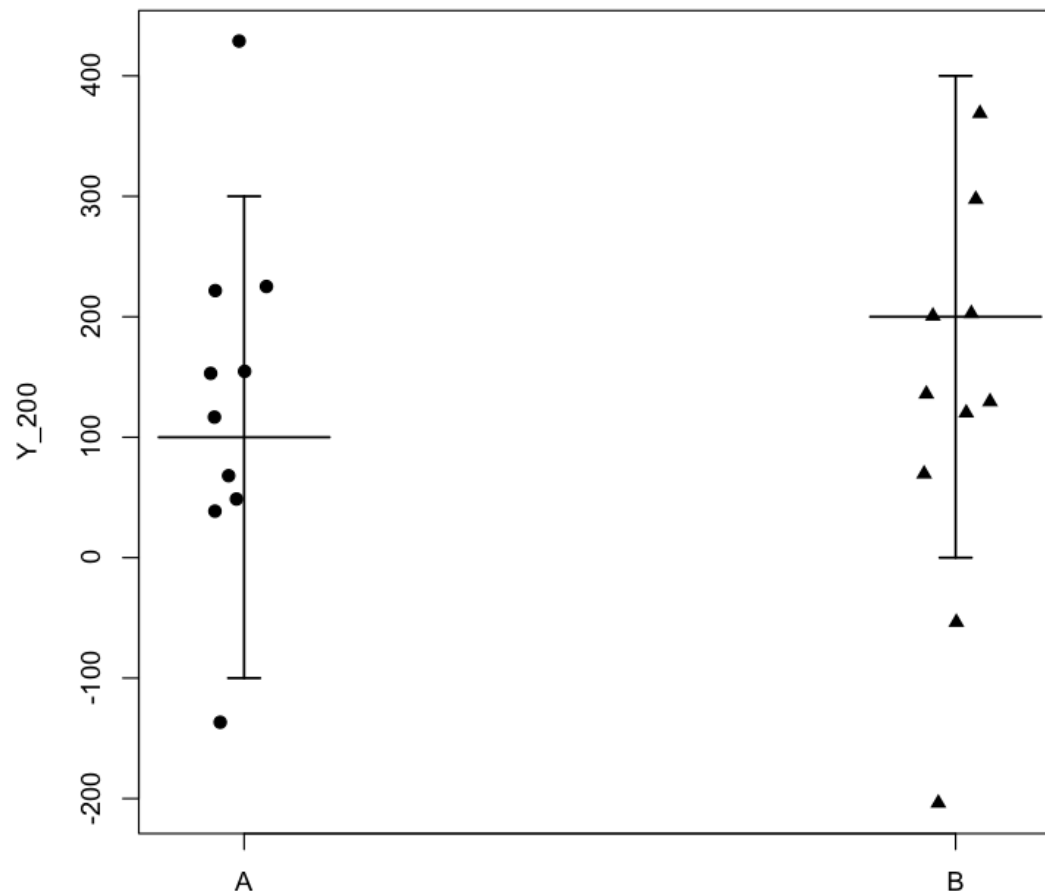
Statistical Significance - IV



- It seems clear here that there is indeed differential expression in Y between A and B.

Statistical Significance - V

- But what if variability were larger, e.g. $\sigma = 200$?



There does not seem to be a difference now.

Weakness of Fold Change

- The fold-change is the same in the two cases. How can one explain this discrepancy?
- The answer is that fold-change by itself **is not a reproducible criterion for discrimination**, because it does not take into account the inherent variability in the data.
- In practice, the fold-change would be **estimated** as the ratio between the **sample means** based on the data at hand, which would not change any of these conclusions.

Hypothesis Tests

- The matter can be formalized with the notion of **hypothesis tests**.
- Let μ_A and μ_B be the expression level of Y under conditions A and B, respectively. We would like to test the **null hypothesis**

$$H_0 : \mu_A = \mu_B$$

against the **alternative hypothesis**

$$H_1 : \mu_A \neq \mu_B$$

at a given **significance level** $100 \times (1 - \alpha)\%$.

Hypothesis Tests - II

- To do this, one computes a **test statistic** T , which is simply some function of the data.
- T is called a "statistic" because it is a function of the data. Because the data is random (it has variability), T is also a **random variable**.
- The null hypothesis is said to be **rejected**, and the alternative hypothesis accepted, if the observed T has an atypical value under the null hypothesis, i.e. T falls in the **rejection region** R of the test, which is defined such that

$$P_{H_0}(T \in R) = \alpha$$

The difference is called **statistically significant**.

The p-value

- The p-value is simply the probability, under the null hypothesis, of observing a more atypical (e.g. larger in magnitude) value of T than the one actually observed, T_0 :

$$p = P_{H_0}(|T| > |T_0|)$$

- Clearly, in this example, the larger $|T_0|$ is, the smaller p is.
- If $p < 1-\alpha$ then H_0 is rejected and the test is statistically significant.

Misconceptions about p-value

- The p-value is **NOT** the probability that the null hypothesis is true.
- Similarly, subtracting the p-value from 1 does **NOT** give the probability that the alternative hypothesis is true.
- The threshold $p < 0.05$ is nothing other than a pure convention.
- The observed value of T itself can be useful (it measures the so-called **effect size**).

Error Rates

- An error rate is a probability of making a mistake in the hypothesis test.
- Error **Type-I** (α): Probability of rejecting H_0 when it is true; **false positive**.
- Error **Type II** (β): Probability of not rejecting H_0 when it is false; **false negative**.

H_0 true H_0 false

Reject H_0

Not Reject H_0

α	$1 - \beta$
$1 - \alpha$	β

Sensitivity/Specificity

- The **sensitivity** of a test, also called **statistical power**, is $1 - \beta$. The more sensitive a test is, the smaller the differences it can **detect**.
- Sensitivity is a function of sample size. The larger it is, the more sensitive the test will be.
- The **specificity** of a test is simply the significance level α . The more specific a test is, the fewest false positives it will produce.
- Sensitivity and specificity are conflicting requirements. One can make one larger by decreasing the other. A superior test will be more sensitive at the same specificity.

Two-Sample t-test

- The most common test used in differential expression studies, by far, is the two-sample t-test. It assumes Gaussian data, but is quite robust against failure of this assumption.
- The version known as Welch's applies to conditions that have different variances.
- The test statistic in this case is

$$T = \frac{\bar{x}_A - \bar{x}_B}{\sqrt{\hat{\sigma}_A^2/n_A + \hat{\sigma}_B^2/n_B}}$$

Two-Sample t-test - II

- In the previous formula:

$\bar{x}_.$ = sample mean

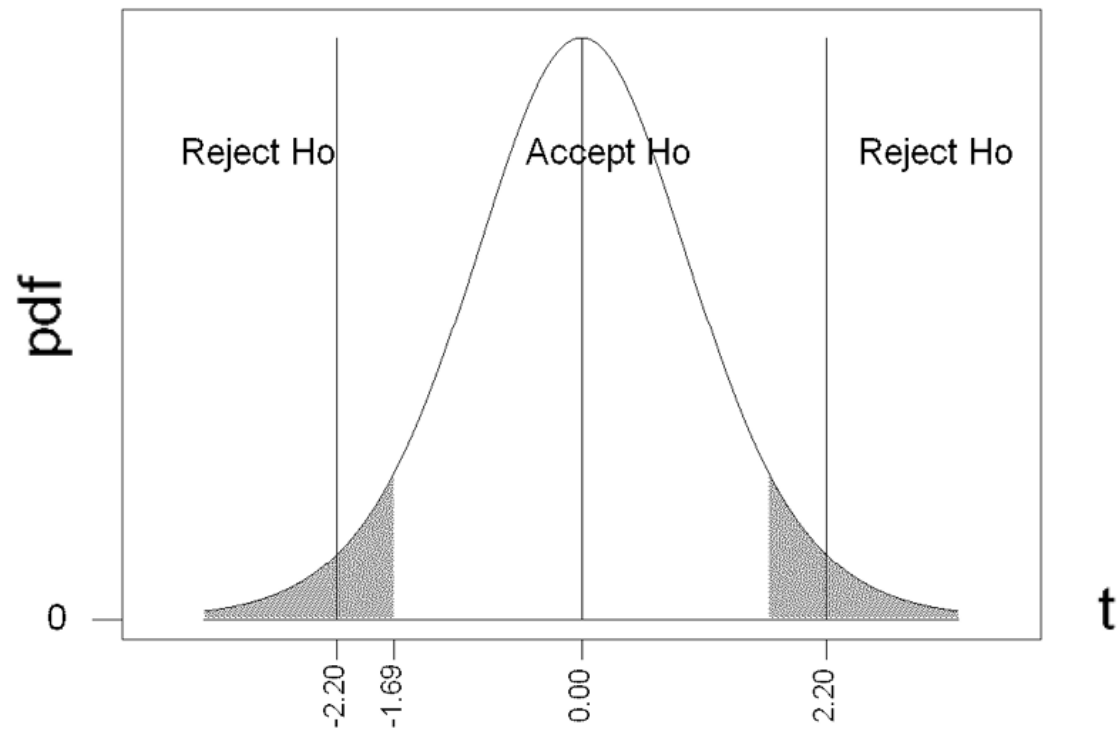
$n_.$ = number of samples

$\hat{\sigma}_.^2$ = sample variance

- Note that T takes into account the variance in the data, which the fold-change does not.
- The smaller the variances, and/or the larger the number of samples are, the larger the magnitude of T is, the smaller the p-value is, and the larger the evidence for rejecting H_0 .

Two-Sample t-test - III

- For Welch's test, under the null hypothesis, T has an approximate **Student's t distribution**, with a noninteger number of degrees of freedom.



(From <http://junior.apk.net/~pmathews/tot/abstract/MTBExecMacros.html>)

Numerical Example

- Same simulated data as before, $\sigma = 25$

```
> t.test(Y_25~cond)
```

```
Welch Two Sample t-test
```

```
data: Y_25 by cond
```

```
t = -13.0557, df = 17.992, p-value =  
1.292e-10
```

```
alternative hypothesis: true difference in  
means is not equal to 0
```

```
95 percent confidence interval:
```

```
-120.53782 -87.12028
```

```
sample estimates:
```

```
mean in group A mean in group B
```

```
102.0391
```

```
205.8682
```

Numerical Example - II

- Same simulated data as before, $\sigma = 200$

```
> t.test(Y~cond)
```

```
Welch Two Sample t-test
```

```
data: Y_200 by cond
```

```
t = 0.0734, df = 17.806, p-value = 0.9423
```

```
alternative hypothesis: true difference in  
means is not equal to 0
```

```
95 percent confidence interval:
```

```
-142.2260 152.5159
```

```
sample estimates:
```

```
mean in group A mean in group B
```

```
131.9043
```

```
126.7593
```

Nonparametric Test

- The t-test is founded on a parametric Gaussian assumption.
- A non parametric alternative can be obtained by considering not the numeric values of the measurement, but only their ranking.
- If most points in A are among the top and those in B are among the bottom of the rankings, intuitively there is discrimination.
- The resulting test is called the **Wilcoxon rank-sum test**.

Numerical Example

- Same simulated data as before, $\sigma = 25$

```
> wilcox.test(Y_25~cond)
```

```
Wilcoxon rank sum test
```

```
data: Y_25 by cond
```

```
W = 0, p-value = 1.083e-05
```

```
alternative hypothesis: true location shift  
is not equal to 0
```

- Note that the p-value is larger than for the t-test; if the Gaussianity assumption is satisfied (it is for this data), then Welch's t-test is more powerful than Wilcoxon's test (and the equal-variance t-test is **uniformly most powerful**).

Numerical Example - II

- Same simulated data as before, $\sigma = 200$

```
> wilcox.test(Y_200~cond)
```

```
Wilcoxon rank sum test
```

```
data: Y by cond
```

```
W = 48, p-value = 0.9118
```

```
alternative hypothesis: true location shift  
is not equal to 0
```


Multiple Testing Issue

- Suppose there are 1000 variables to be tested and none of them are significantly different between the conditions.
- If we apply the standard significance level of 0.05, and assuming the tests are independent, we expect to get $1000 \times 0.05 = 50$ false positives. (why?)
- This is clearly an acceptable situation. If I have 20 significant differences at the 95% confidence level, there is a good chance that all of them are false positives.
- This happens because of the **multiple tests**.

Bonferroni Correction

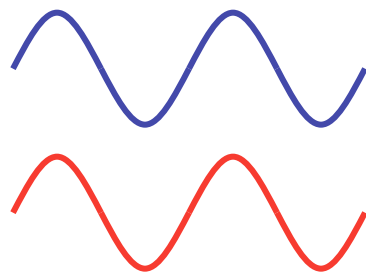
- One possibility to correct this situation is to enlarge the p-values by multiplying them by the number of tests. This is known as the **Bonferroni correction**.
- If one starts with 1000 variables, to get a corrected $p < 0.05$ one has to get $p < 0.00005$.
- However, it can be shown that the Bonferroni correction is **conservative** (it reduces the false positives too much, creating an excess of false negatives).

Ranking by Effect Size

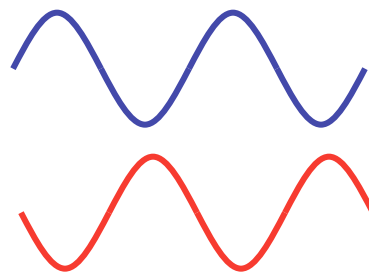
- There are other less-conservative alternatives to the Bonferroni correction, such as the "false discovery rate". But the situation is a bit chaotic. No one knows for sure which correction method is the best.
- One alternative in practice is to ignore the p-values and simply rank the variables by effect sizes.

Correlation

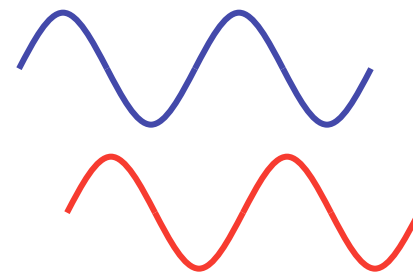
- Correlation is a measure of the coincidence (or lack there of) of directionality of change across samples between two variables.



high positive
correlation



high negative
correlation



small
correlation (in
magnitude)

Correlation is **NOT** causality.

Sample Correlation Coefficient

- Correlation can be measured by **Pearson's correlation coefficient**

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

- We have $-1 \leq \rho \leq 1$ and

$\rho \rightarrow 1 \Rightarrow$ large positive correlation

$\rho \rightarrow -1 \Rightarrow$ large negative correlation

$\rho \approx 0 \Rightarrow$ small correlation

Sample Correlation Coefficient

- Pearson's correlation coefficient measures **linear** association.
- It is quite sensitive to outliers, therefore plotting the data is always recommended.
- The **coefficient of determination**

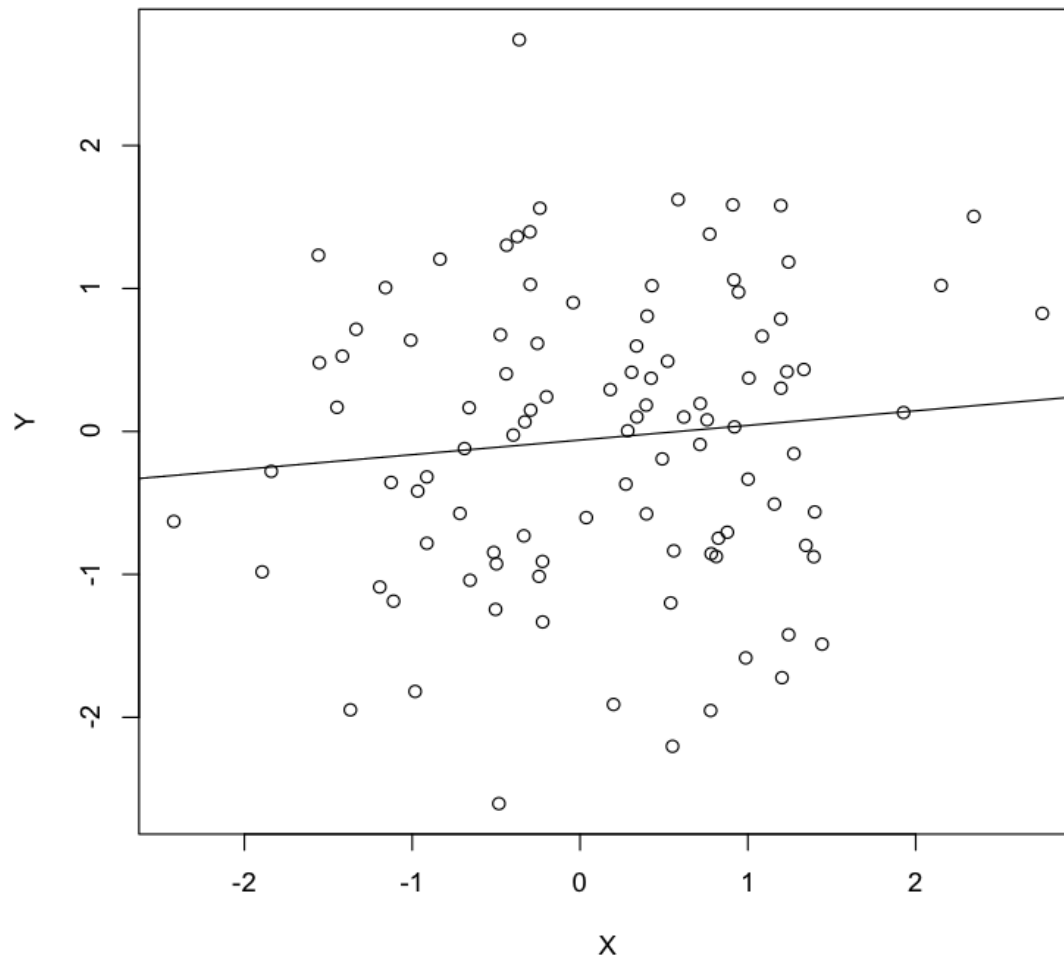
$$R^2 = \rho^2$$

is a measure between 0 and 1 that measures the **magnitude** of association.

- There are other nonparametric correlation coefficients (as in the case of t-tests).

Numerical Example - II

- With $\text{cor} = 0.1$



Numerical Example - II

- With $\text{cor} = 0.8$

