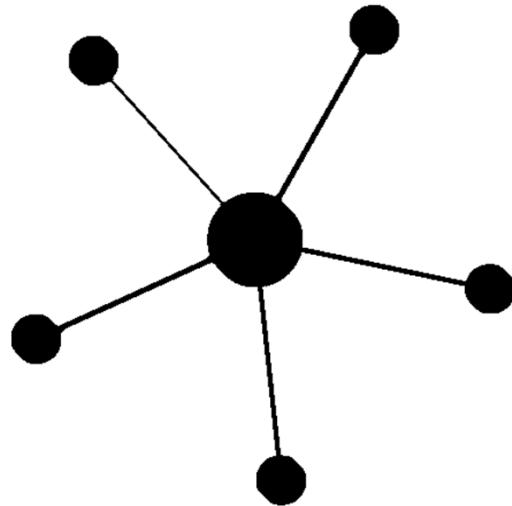


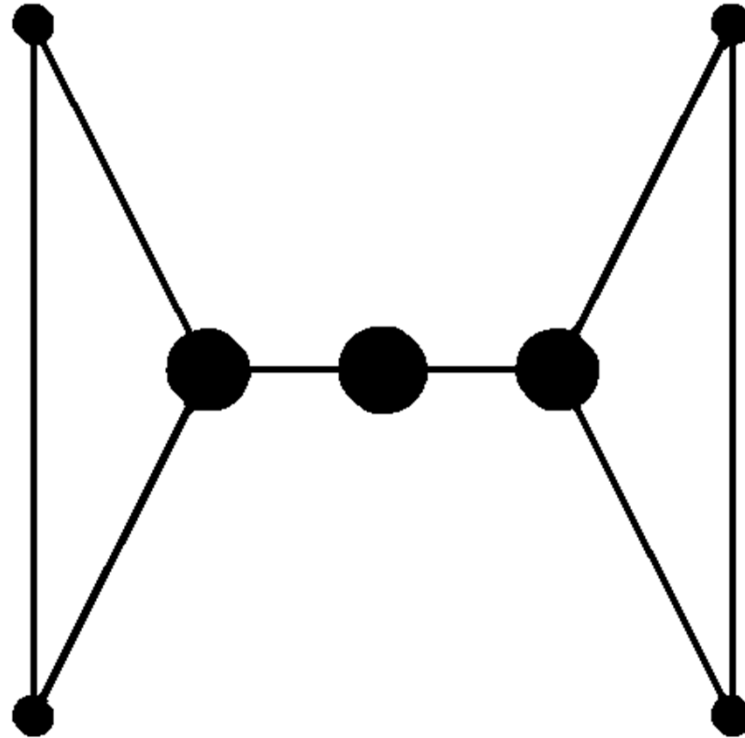
How to find “important” nodes?

- By their degree
- Hubs = important
- Degree weighted by self-consistent importance :
Google’s PageRank



How to find “important” nodes?

- By their connectivity
- Connectors = important
- Betweenness-centrality

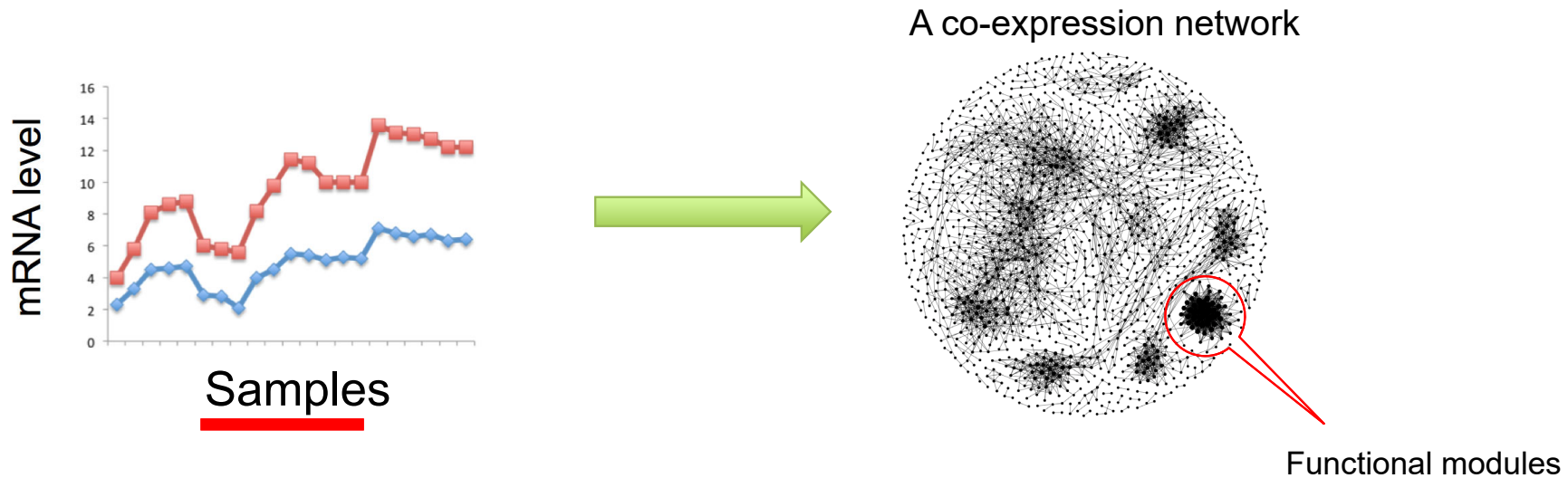


Betweenness centrality: definition

- Take a node i
- There are $(N-1)*(N-2)/2$ pairs of other nodes
- For each pair find the shortest path on the network
- If more than one shortest path, sample them equally
- Betweenness-centrality $C(i) \sim$ the number of shortest paths going through node i

To analyze
correlations in expression
for all pairs of genes:
Co-expression networks

How to construct a co-expression network?



- Start with a matrix of log2 of expression levels of N genes in K samples (conditions): for our T-cell data N=3000, K=47
- For each of $N(N-1)/2$ pairs of genes i and j calculate the correlation coefficient $\rho_{ij} = \sigma_{ij} / \sigma_i \sigma_j$ of gene levels across K samples
- Put a threshold, e.g. $\rho_{ij} > 0.85$, or otherwise select the most correlated pairs of genes (~4500 in our case). Now you have a weighted network.
- Identify densely interconnected functional modules in this network.
- Modules can be used to infer unknown functions of genes via “Guilt by Association” principle.

Co-expression network analysis exercise

- Start Gephi and open [coexpression_network_random_start.gephi](#)
- Run “Layout” → Fruchterman Reingold → Speed 10.0
- Run “Average degree”, “Network diameter”, “Modularity” in the Statistics tab in the right panel.
- Color nodes by “modularity class”:
Appearance → Nodes → Partition → Palette Icon → Modularity class
- Size nodes first by “degree”.
Appearance → Nodes → Ranking → Multiple Circles Icon → Degree
 - If the nodes are too small, select “Min size”: 10 and “Max size”:80
 - Nodes in large tightly connected clusters have large degree
- Then size nodes by “betweenness-centrality”
Appearance → Nodes → Ranking → Multiple Circles Icon → Betweenness-centrality
 - Large circles are “coordinator” genes connecting different co-expressed clusters to each other. Potentially biologically interesting

Disease-disease similarity network

- Based on the table summarizing all current medical knowledge of genes implicated in diseases:
 - Rows: 516 common human diseases
 - Columns: 25,000 human genes
 - Matrix element $D_{i\alpha} = 1$ if the gene α is known to be involved in the disease i . 0 – otherwise
- Constructed disease-disease similarity network:
 - Weight of the edge - # of shared genes between two diseases
 - Easy to construct: the adjacency matrix A of the network is simply $A = D \cdot D^+$

Disease network analysis exercise

- Start Gephi and open `disease_disease_random_start.gexi`
- Run “Layout” → Fruchterman Reingold → Speed 10.0
Observe how clusters emerge.
- Run “Average degree”, “Network diameter”, “Modularity” analysis tools in the right panel.
- Color nodes with **medical term: “disorder class”**
Appearance → Nodes → Partition → Palette Icon → Disorder class
- Then color nodes by “modularity class”. See how well it agrees with the previous color.
Appearance → Nodes → Partition → Palette Icon → Modularity class
- Size nodes first by “**degree**”.
Appearance → Nodes → Ranking → Multiple Circles Icon → Degree
 - Which disease has the largest degree?
- Size nodes by “**betweenness centrality**”
Appearance → Nodes → Ranking → Multiple Circles Icon → Degree
 - Which diseases have the largest betweenness-centrality?
These “connector” diseases linking different diseases clusters to each other. They highlight potentially interesting connections between diseases

Review for the Final Exam

Grading

Attendance 10%

Homework 20%

Midterm 1 20%

Midterm 2 20%

Final 30%

Midterm Info

- **Bring Your OWN Device (BYOD)** exam in this room on May 12, 7pm-9pm.
 - Come earlier, say 6:30pm
 - Bring **laptop** and **charger**
 - Bring **calculator** not on phone
 - Bring **UIUC ID**
- **Closed book exam**; no books, notes, phones...
- **Calculators (not on smartphones)** can be used
- The following **two printouts and lecture slides** will be provided

Name	Probability Distribution	Mean	Variance	Section in Book
Discrete				
Uniform	$\frac{1}{n}, a \leq b$	$\frac{(b + a)}{2}$	$\frac{(b - a + 1)^2 - 1}{12}$	3-5
Binomial	$\binom{n}{x} p^x (1 - p)^{n-x},$ $x = 0, 1, \dots, n, 0 \leq p \leq 1$	np	$np(1 - p)$	3-6
Geometric	$(1 - p)^{x-1} p,$ $x = 1, 2, \dots, 0 \leq p \leq 1$	$1/p$	$(1 - p)/p^2$	3-7.1
Negative binomial	$\binom{x-1}{r-1} (1 - p)^{x-r} p^r$ $x = r, r + 1, r + 2, \dots, 0 \leq p \leq 1$	r/p	$r(1 - p)/p^2$	3-7.2
Poisson	$\frac{e^{-\lambda} \lambda^x}{x!}, x = 0, 1, 2, \dots, 0 < \lambda$	λ	λ	3-9
Continuous				
Uniform	$\frac{1}{b - a}, a \leq x \leq b$	$\frac{(b + a)}{2}$	$\frac{(b - a)^2}{12}$	4-5
Normal	$\frac{1}{\sigma\sqrt{2\pi}} e^{-1/2(\frac{x-\mu}{\sigma})^2}$ $-\infty < x < \infty, -\infty < \mu < \infty, 0 < \sigma$	μ	σ^2	4-6
Exponential	$\lambda e^{-\lambda x}, 0 \leq x, 0 < \lambda$	$1/\lambda$	$1/\lambda^2$	4-8
Erlang	$\frac{\lambda^r x^{r-1} e^{-\lambda x}}{(r - 1)!}, 0 < x, r = 1, 2, \dots$	r/λ	r/λ^2	4-9.1
Gamma	$\frac{\lambda^r x^{r-1} e^{-\lambda x}}{\Gamma(r)}, 0 < x, 0 < r, 0 < \lambda$	r/λ	r/λ^2	4-9.2

What may be on the final exam?

- Probability Multiplication, Combinatorics
- Bayes Theorem
- Discrete & Continuous Random Variables
- Joint Probability Distributions, Covariation/Correlations
- Sampling distributions and parameter point estimation
- Confidence Intervals
- Hypothesis testing for one and two samples
- Other topics
- Look at Homework 1-4 for examples of problems

One-sample hypothesis testing

3. (8 points) The college bookstore tells prospective students that the average cost of its textbooks is \$52 with a standard deviation of \$4.50. A group of statistics students think that the average cost is actually higher. In order to test bookstore's claim against this alternative hypothesis, the students bought a random sample of 100 books. The mean price of this sample was \$52.80. Perform the hypothesis test at the 5% level of significance and state your decision.

What type of hypothesis should I apply?

A. Two-sided: $\mu_1 \neq \mu_0$

B. One-sided: $\mu_1 > \mu_0$

C. One-sided: $\mu_1 < \mu_0$

D. Three-sided

E. I have no idea

Get your i-clickers

3. (8 points) The college bookstore tells prospective students that the average cost of its textbooks is \$52 with a standard deviation of \$4.50. A group of statistics students think that the average cost is **actually higher**. In order to test bookstore's claim against this alternative hypothesis, the students bought a random sample of 100 books. The mean price of this sample was \$52.80. Perform the hypothesis test at the 5% level of significance and state your decision.

The standard deviation of \bar{x} in this sample is:

A. \$4.50

B. \$45

C. \$0.45

D. I have no idea

Get your i-clickers

3. (8 points) The college bookstore tells prospective students that the average cost of its textbooks is \$52 with a standard deviation of \$4.50. A group of statistics students think that the average cost is **actually higher**. In order to test bookstore's claim against this alternative hypothesis, the students bought a random sample of 100 books. The mean price of this sample was \$52.80. Perform the hypothesis test at the 5% level of significance and state your decision.

3. (8 points) The college bookstore tells prospective students that the average cost of its textbooks is \$52 with a standard deviation of \$4.50. A group of statistics students think that the average cost is **actually higher**. In order to test bookstore's claim against this alternative hypothesis, the students bought a random sample of 100 books. The mean price of this sample was \$52.80. Perform the hypothesis test at the 5% level of significance and state your decision.

Answer: Hypothesis: $\begin{cases} H_0 : \mu = 52 \\ H_1 : \mu > 52 \end{cases}$. The critical z-value can be obtained from $z^* = \frac{52.8 - 52}{4.5 / 10} = 1.78$. Since $z^* > z_\alpha = 1.65$, this test statistic lies in the rejection region for H_0 . Thus null hypothesis H_0 will be rejected and alternative hypothesis H_1 is accepted.

Two-sample hypothesis

Mating Calls. In a study of mating calls in the gray treefrogs *Hyla chrysoscelis* and *Hyla versicolor*, Gerhart (1994) reports that in a location in Louisiana the following data on the length of male advertisement calls have been collected:

	Sample size	Average duration	SD of duration	Duration range
<i>Hyla chrysoscelis</i>	43	0.65	0.18	0.36–1.27
<i>Hyla versicolor</i>	12	0.54	0.14	0.36–0.75

The two species cannot be distinguished by external morphology, but *H. chrysoscelis* are diploids while *H. versicolor* are tetraploids. The triploid crosses exhibit high mortality in larval stages, and if they attain sexual maturity, they are sterile. Females responding to the mating calls try to avoid mismatches.

Based on the data summaries provided, test whether the length of call is a discriminatory characteristic? Use $\alpha = 0.05$.

	Sample size	Average duration	SD of duration
<i>Hyla chrysoscelis</i>	43	0.65	0.18
<i>Hyla versicolor</i>	12	0.54	0.14

Based on the data summaries provided, test whether the length of call is a discriminatory characteristic? Use $\alpha = 0.05$.

	Sample size	Average duration	SD of duration
<i>Hyla chrysoscelis</i>	43	0.65	0.18
<i>Hyla versicolor</i>	12	0.54	0.14

Based on the data summaries provided, test whether the length of call is a discriminatory characteristic? Use $\alpha = 0.05$.

1. Use two-sided hypothesis
2. $z_{\{\alpha/2\}}=1.96$
3. $Z=(0.65-0.54)/\sqrt{0.18.^2/43+0.14.^2/12}=2.2516$
4. Since $Z > z_{\{\alpha/2\}}$ null hypothesis can be rejected

Confidence intervals

2. (6 points) The operations manager of a large production plant would like to estimate the mean amount of time a worker takes to assemble a new electronic component. Assume that the standard deviation of this assembly time is 3.6 minutes. After observing a sample of 100 workers assembling similar devices, the manager noticed that their average time was 16.2 minutes. Construct a 90% confidence interval for the population mean of the assembly time.

two-sided

What Z should I look up in the table?

- A. $\Phi(Z)=0.9$
- B. $\Phi(Z)=0.05$
- C. $\Phi(Z)=0.95$
- D. $\Phi(Z)=0.1$
- E. I have no idea

Get your i-clickers

2. (6 points) The operations manager of a large production plant would like to estimate the mean amount of time a worker takes to assemble a new electronic component. Assume that the standard deviation of this assembly time is 3.6 minutes. After observing a sample of 100 workers assembling similar devices, the manager noticed that their average time was 16.2 minutes. Construct a 90% confidence interval for the population mean of the assembly time.

What Z should I look up in the table?

- A. $\Phi(Z)=0.9$
- B. $\Phi(Z)=0.05$
- C. $\Phi(Z)=0.95$
- D. $\Phi(Z)=0.1$
- E. I have no idea

Get your i-clickers

2. (6 points) The operations manager of a large production plant would like to estimate the mean amount of time a worker takes to assemble a new electronic component. Assume that the standard deviation of this assembly time is 3.6 minutes. After observing a sample of 100 workers assembling similar devices, the manager noticed that their average time was 16.2 minutes. Construct a **90% confidence interval** for the population mean of the assembly time.

2. (6 points) The operations manager of a large production plant would like to estimate the mean amount of time a worker takes to assemble a new electronic component. Assume that the standard deviation of this assembly time is 3.6 minutes. After observing a sample of 100 workers assembling similar devices, the manager noticed that their average time was 16.2 minutes. Construct a **90% confidence interval** for the population mean of the assembly time.

Answer: Let μ denote the mean assembly time (in minutes). We want a 90% confidence interval for μ based on the following information: $n = 100$, $\bar{X} = 16.2$, $\alpha = 0.1$, $\sigma = 3.6$. Since σ is known, we can use normal distribution to calculate confidence interval:

$$\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 16.2 \pm (1.65) \frac{3.6}{10} = [15.61, 16.79]$$

What is X in this problem?

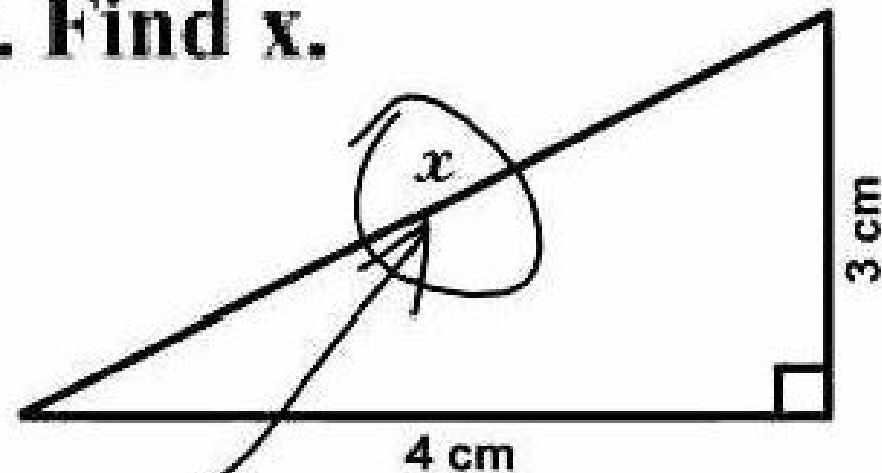
- **What is X?** Look for keywords:
 - Find the probability that....
 - What is the mean (or variance) of...
- **What are the parameters?**

Look for keywords:

- Given that...
- Assuming that...

- **Is X discrete or continuous?**

3. Find x.



Here it is

Discrete Probability Distributions

(8 points) You are doing a long series of experiments. Assume that each of your experiments has a probability of 0.02 of succeeding. Assume that your experiments are independent.

(A) (2 points) What is the probability that you first succeed on tenth experiment?

(B) (2 points) What is the probability that it requires more than five experiments for you to succeed?

(C) (2 points) What is the mean number of experiments needed to succeed once?

(D) (2 points) What is the probability that the second experiment that worked is the tenth one since you started?

2. (8 points, 2 points each) You are doing a long series of experiments. Assume that each of your experiments has a probability of 0.02 of succeeding. Assume that your experiments are independent.
- (a) What is the probability that you first succeed on tenth experiment?

$$P(X=10) = (1-0.02)^9 * 0.02 = 0.0167$$

- (b) What is the probability that it requires more than five experiments for you to succeed?

$$\begin{aligned} P(X > 5) &= 1 - P(X=1) - P(X=2) - P(X=3) - P(X=4) - P(X=5) \\ &= 1 - 0.98^0 * 0.02 - 0.98^1 * 0.02 - 0.98^2 * 0.02 - 0.98^3 * 0.02 - 0.98^4 * 0.02 = 0.9039 \\ \text{Easier solution: } P(X > 5) &= 0.98^5 = 0.9039 \end{aligned}$$

- (c) What is the mean number of experiments needed to succeed once?

$$\text{Since } X \text{ follows geometric distribution, the mean value of } X \text{ is } 1/0.02 = 50.$$

- (d) What is the probability that the second experiment that worked is the tenth one since you started
- $$\text{Probability} = 9 * 0.02 * 0.98^8 * 0.02 = 0.0031$$

Continuous Probability Distributions

(12 points) Time interval separating subsequent bus arrivals at a stop is an exponential random variable with mean 20 minutes. Steve and Andrew work at the same place and each will be late to work unless they board a bus on or before 8:40am. Steve comes to the bus stop exactly at 8am. Andrew also comes to the same bus stop but at a random time, uniformly distributed between 8am and 8:30am. Both of them take the first bus that arrives.

(a) (4 points) What is the probability that Steve will be late for work tomorrow?

(b) (4 points) What is the probability that Andrew will be late for work tomorrow?

(c) (4 points) What is the probability that Steve and Andrew will ride the same bus

(12 points) Time interval separating subsequent bus arrivals at a stop is an exponential random variable with mean 20 minutes. Steve and Andrew work at the same place and each will be late to work unless they board a bus on or before 8:40am. Steve comes to the bus stop exactly at 8am. Andrew also comes to the same bus stop but at a random time, uniformly distributed between 8am and 8:30am. Both of them take the first bus that arrives.

(a) **(4 points)** What is the probability that Steve will be late for work tomorrow?

$$\text{Answers: } P(\text{Steve late}) = 1 - P(T < 40) = 1 - \frac{1}{20} \int_0^{40} e^{-t/20} dt = e^{-2} = 0.1353$$

(b) **(4 points)** What is the probability that Andrew will be late for work tomorrow?

Answers:

$$P(\text{Andrew late}) = \int_0^{30} \frac{dx}{30} P(T \geq 40 | T > x) = \int_0^{30} \frac{dx}{30} e^{-(40-x)/20} = \frac{e^{-2}}{30} \int_0^{30} e^{x/20} dx = \frac{20e^{-2}}{30} (e^{30/20} - 1) = 0.3141$$

(c) **(4 points)** What is the probability that Steve and Andrew will ride the same bus?

Probability that Steve will not leave by the time x when Andrew comes is $\exp(-x/20)$.

It needs to be integrated over $\int_0^{30} dx/30 \exp(-x/20) =$

$$\text{Answers: } P(\text{Steve and Andrew meet}) = \int_0^{30} \frac{dx}{30} e^{-x/20} = \frac{20}{30} (1 - e^{-30/20}) = 0.5179$$

Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

WHY DO TWINS HAVE DIFFERENT FINGERPRINTS
WHY ARE AMERICANS AFRAID OF DRAGONS

WHY IS HTTPS CROSSED OUT IN RED
WHY IS THERE A LINE THROUGH HTTPS
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK
WHY IS HTTPS IMPORTANT

QUESTIONS FOUND IN GOOGLE AUTOCOMPLETE



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

WHY DO WHALES JUMP
WHY ARE WITCHES GREEN
WHY ARE THERE MIRRORS ABOVE BEDS
WHY DO I SAY UH
WHY IS SEA SALT BETTER
WHY ARE THERE TREES IN THE MIDDLE OF FIELDS
WHY IS THERE NOT A POKEMON MMO
WHY IS THERE LAUGHING IN TV SHOWS
WHY ARE THERE DOORS ON THE FREEWAY
WHY ARE THERE SO MANY SVCHOST.EXE RUNNING
WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA
WHY ARE THERE SCARY SOUNDS IN MINECRAFT
WHY IS THERE KICKING IN MY STOMACH
WHY ARE THERE TWO SLASHES AFTER HTTP
WHY ARE THERE CELEBRITIES
WHY DO SNAKES EXIST
WHY DO OYSTERS HAVE PEARLS
WHY ARE DUCKS CALLED DUCKS
WHY DO THEY CALL IT THE CLAP
WHY ARE KYLE AND CARTMAN FRIENDS
WHY IS THERE AN ARROW ON AANG'S HEAD
WHY ARE TEXT MESSAGES BLUE
WHY ARE THERE MUSTACHES ON CLOTHES
WHY ARE THERE MUSTACHES ON CARS
WHY ARE THERE MUSTACHES EVERYWHERE
WHY ARE THERE SO MANY BIRDS IN OHIO
WHY IS THERE SO MUCH RAIN IN OHIO
WHY IS OHIO WEATHER SO WEIRD

WHY DO IGUANAS DIE
WHY AREN'T THERE DINOSAUR GHOSTS

WHY AREN'T ECONOMISTS RICH
WHY DO AMERICANS CALL IT SOCCER
WHY ARE MY EARS RINGING
WHY ARE THERE SO MANY AVENGERS
WHY ARE THE AVENGERS FIGHTING THE X MEN
WHY IS WOLVERINE NOT IN THE AVENGERS

WHY ARE THERE SWARMS OF GNATS
WHY IS THERE PHLEGM
WHY ARE THERE SO MANY CROWS IN ROCHESTER, MN
WHY IS PSYCHIC WEAK TO BUG
WHY DO CHILDREN GET CANCER
WHY IS POSEIDON ANGRY WITH ODYSSEUS
WHY IS THERE ICE IN SPACE

WHY ARE THERE ANTS IN MY LAPTOP
WHY IS EARTH TILTED
WHY IS SPACE BLACK
WHY IS OUTER SPACE SO COLD
WHY ARE THERE PYRAMIDS ON THE MOON
WHY IS NASA SHUTTING DOWN



WHY IS THERE AN OWL IN MY BACKYARD
WHY IS THERE AN OWL OUTSIDE MY WINDOW
WHY IS THERE AN OWL ON THE DOLLAR BILL
WHY DO OWLS ATTACK PEOPLE
WHY ARE AK 47s SO EXPENSIVE
WHY ARE THERE HELICOPTERS CIRCLING MY HOUSE
WHY ARE THERE GODS
WHY ARE THERE TWO SPOCKS

WHY ARE THERE BRIDESMAIDS
WHY DO DYING PEOPLE REACH UP
WHY AREN'T THERE VARICOSE ARTERIES
WHY ARE OLD KUNGONS DIFFERENT



WHY ARE THERE TINY SPIDERS IN MY HOUSE
WHY DO SPIDERS COME INSIDE
WHY ARE THERE HUGE SPIDERS IN MY HOUSE
WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE
WHY ARE THERE SPIDERS IN MY ROOM
WHY ARE THERE SO MANY SPIDERS IN MY ROOM
WHY DO SPIDER BITES ITCH
WHY IS DYING SO SCARY

WHY IS MT VESUVIUS THERE
WHY DO THEY SAY T MINUS
WHY ARE THERE OBELISKS
WHY ARE WRESTLERS ALWAYS WET
WHY ARE OCEANS BECOMING MORE ACIDIC
WHY IS ARWEN DYING
WHY AREN'T MY QUAIL LAYING EGGS
WHY AREN'T MY QUAIL EGGS HATCHING
WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

WHY ARE CIGARETTES LEGAL
WHY ARE THERE DUCKS IN MY POOL
WHY IS JESUS WHITE
WHY IS THERE LIQUID IN MY EAR
WHY DO Q TIPS FEEL GOOD
WHY DO GOOD PEOPLE DIE



WHY IS PROGRAMMING SO HARD
WHY IS THERE A 0 OHM RESISTOR
WHY DO AMERICANS HATE SOCCER
WHY DO RHYMES SOUND GOOD
WHY DO TREES DIE
WHY IS THERE NO SOUND ON CNN
WHY AREN'T POKEMON REAL
WHY AREN'T BULLETS SHARP
WHY DO DREAMS SEEM SO REAL

WHY IS THERE NO GPS IN LAPTOPS
WHY DO KNEES CLICK
WHY AREN'T THERE E GRADES
WHY IS ISOLATION BAD
WHY DO BOYS LIKE ME
WHY DON'T BOYS LIKE ME
WHY IS THERE ALWAYS A JAVA UPDATE
WHY ARE THERE RED DOTS ON MY THIGHS
WHY IS LYING GOOD

WHY IS SEX SO IMPORTANT



WHY ARE DOGS AFRAID OF FIREWORKS
WHY IS THERE NO KING IN ENGLAND