

Measuring Similarity Within Species

Rajat Sirohi

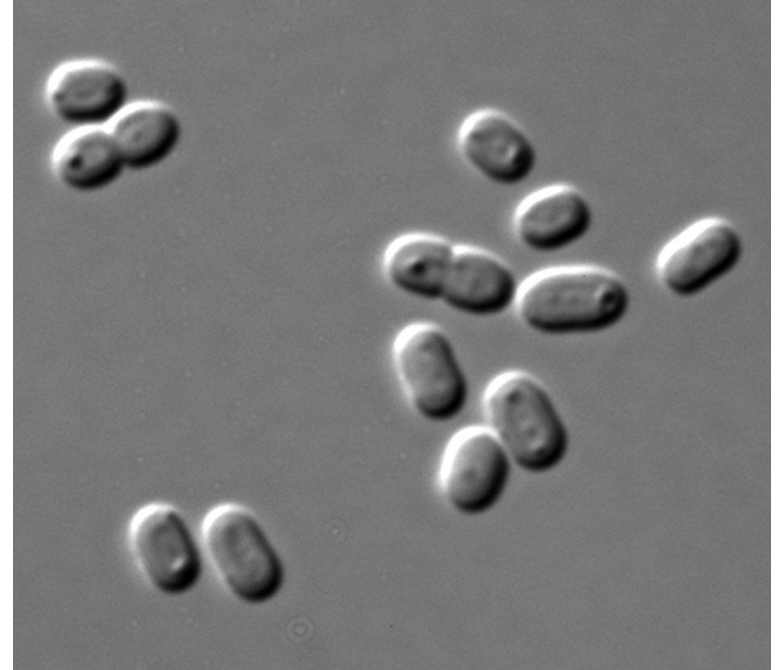
Overview and Introduction

“How similar are organisms within the same species compared to any random pair of organisms?”

- Goals
 - Devise measure of similarity
 - Calculate intraspecies similarity
- Hypotheses
 - H_0 : the two organisms are no more similar than random chance would predict
 - H_a : the two organisms belong to the same species

Data Collection

- NCBI Genome Database: GenBank
- *Synechococcus elongatus*
- 100 random chosen samples
- Conditions:
 - Reputable database ensures non-biased samples
 - Random selection process
 - 10% condition
 - Success/failure condition



Expected Distribution

N = length of genome strand

n = number of nucleotide options = 4

S = number of simulations

ψ = measure of similarity

μ = mean

σ = standard deviation

$A = a_1 a_2 a_3 \dots a_N$

$B = b_1 b_2 b_3 \dots b_N$

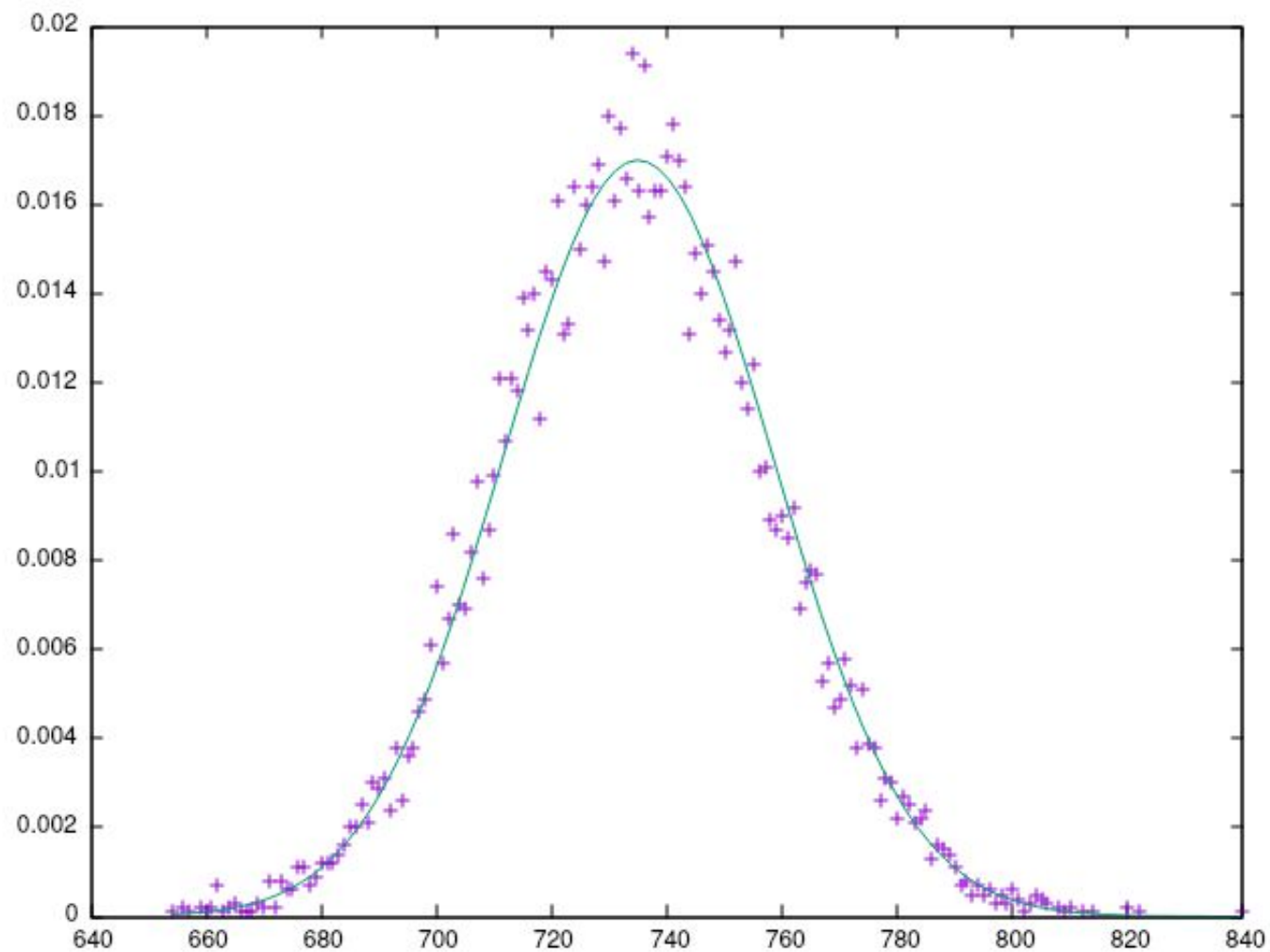
$$P(a_i == b_i) = 1/n$$

$$E(\psi) = N * P(a_i == b_i) = N/n$$

$$\text{Freq}(\psi) = (N \text{ choose } \psi) (1/n)^\psi ((n-1)/n)^{N-\psi}$$

$$\mu = Np = N/n$$

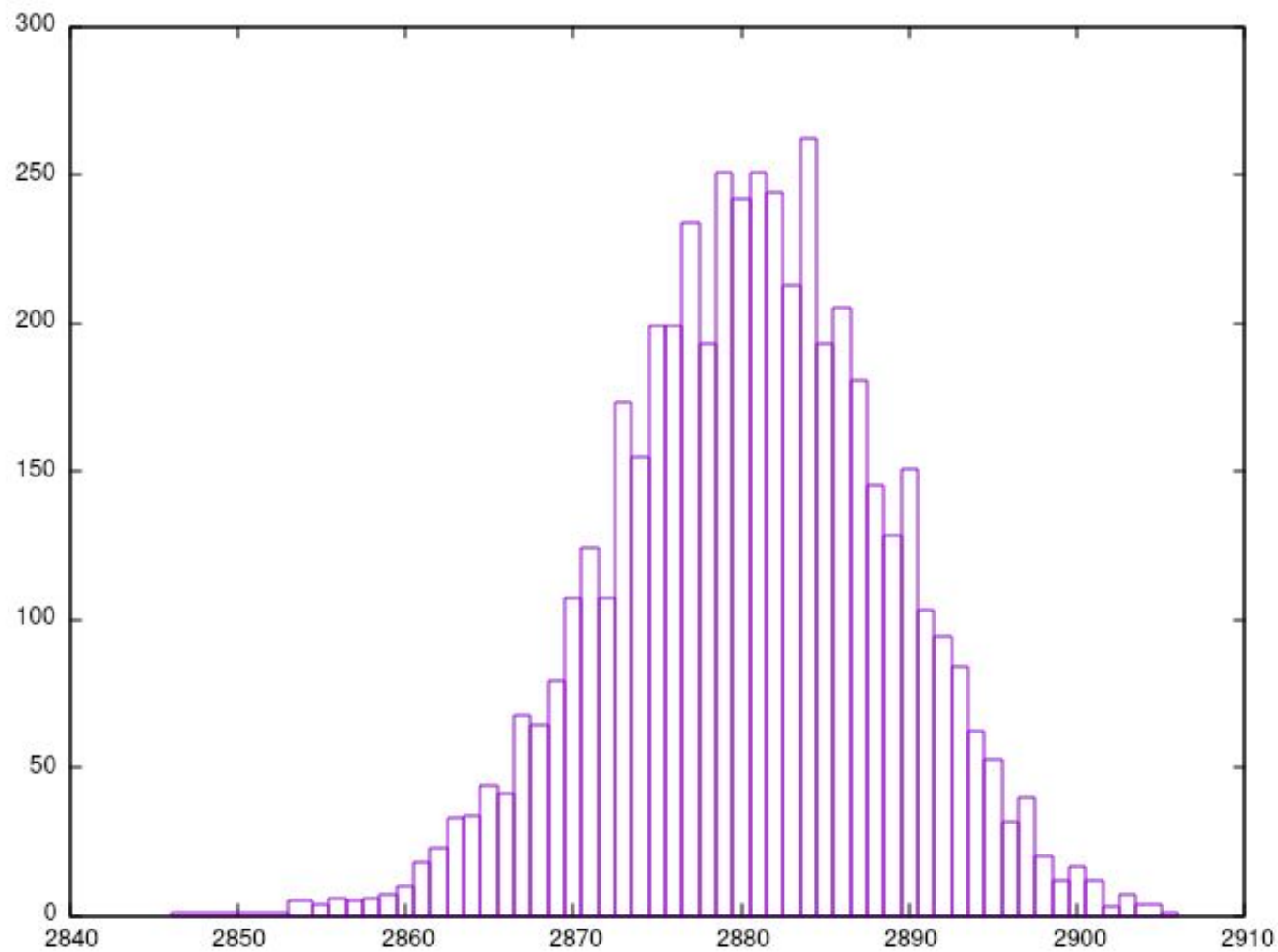
$$\sigma = \text{sqrt}(Npq) = \text{sqrt}(N * (n-1)/n/n)$$

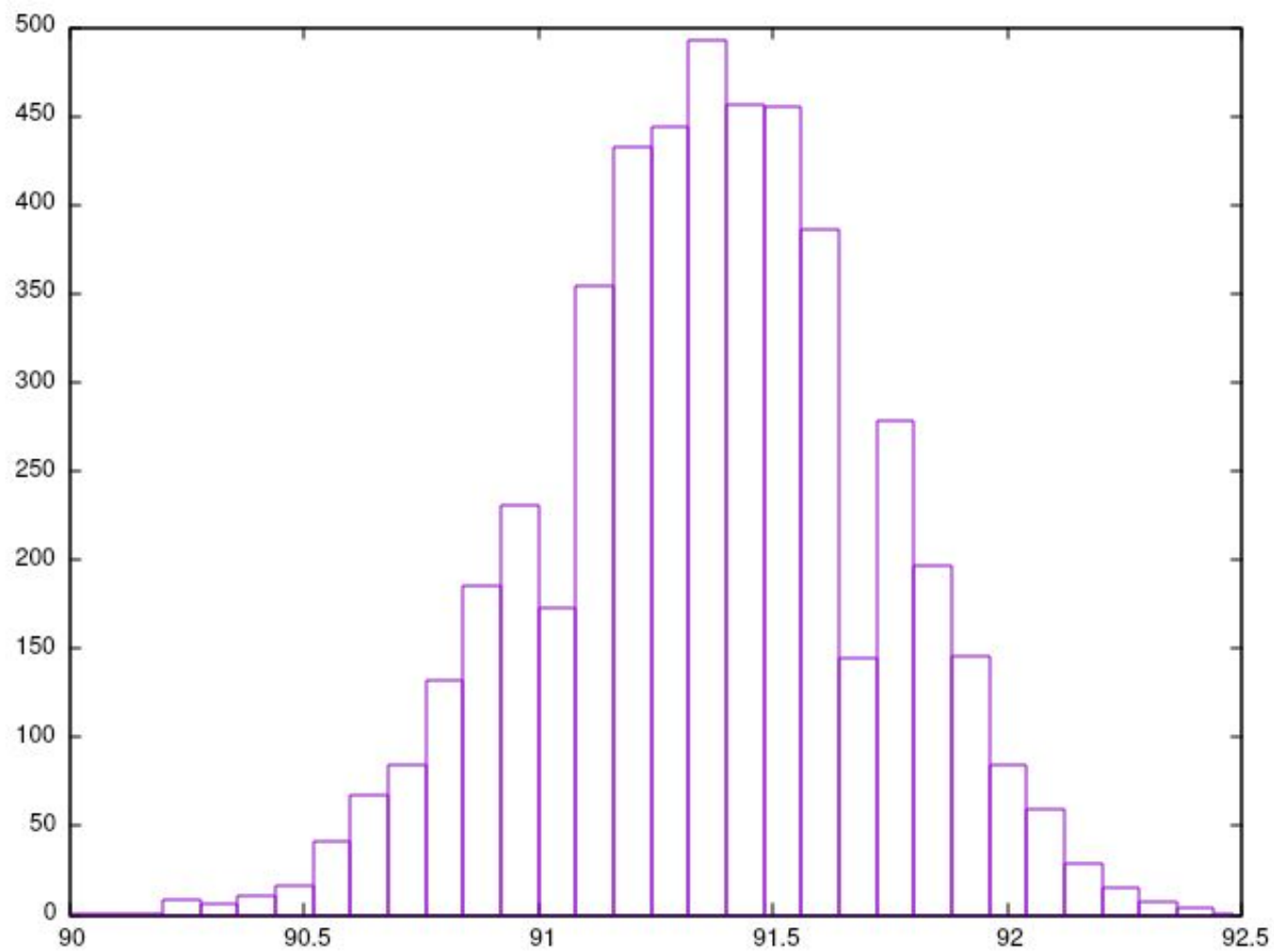


Data Analysis

- Analyzed 100 samples, resulting in $(100 \text{ choose } 2) = 4950$ comparisons
- Results
 - $\psi = 735$
 - Sigma ~ 23.479
 - $\psi_{\text{hat}} \sim 2880.503$
 - Z_score ~ 91.380
 - P_val ~ 0







Conclusions

- Given the extremely low p-value, we can reasonably **reject the null hypothesis** in favor of the alternative hypothesis with near 100% confidence
- Confidence interval: $(-1410.503, 2880.503)$



Reflection

- Extension opportunities
 - Use different measure of similarity
 - Cosine similarity
 - Compare pairs rather than singular nucleotides
 - Weighted element-wise or cosine similarity
 - Restrict random distribution to viable genomes
- Difficulties collecting data
- Unsurprising results, stepping stone for further research





Thank you