

Section #7 Solutions

Based on the work of many CS109 staffs

The Honor Code

We have decided that automated tools should be used to identify if two submissions are suspiciously similar. (N.B. these tools are never used as a basis for community standards cases — those always require professional human opinion.) However, automated tools are never perfect. As active CS109 students we would like to explore the chance of a false positive in automated tools. For this task, we will consider Breakout, a CS106A assignment where students implement Breakout.

1. Beta Sum

What is the distribution of the sum of 100 IID Betas? Let X be the sum

$$X = \sum_{i=1}^{100} X_i \quad \text{where each } X_i \sim \text{Beta}(a = 3, b = 4)$$

Note the variance of a Beta:

$$\text{Var}(X_i) = \frac{ab}{(a+b)^2(a+b+1)} \quad \text{where } X_i \sim \text{Beta}(a, b)$$

By the Central Limit Theorem, the sum of equally weighted IID random variables will be Normally distributed. We calculate the expectation and variance of X_i using the beta formulas:

$$\begin{aligned} E(X_i) &= \frac{a}{a+b} && \text{Expectation of a Beta} \\ &= \frac{3}{7} \approx 0.43 \end{aligned}$$

$$\begin{aligned} \text{Var}(X_i) &= \frac{ab}{(a+b)^2(a+b+1)} && \text{Variance of a Beta} \\ &= \frac{3 \cdot 4}{(3+4)^2(3+4+1)} \\ &= \frac{12}{49 \cdot 8} \approx 0.03 \end{aligned}$$

$$\begin{aligned} X &\sim N(\mu = n \cdot E[X_i], \sigma^2 = n \cdot \text{Var}(X_i)) \\ &\sim N(\mu = 43, \sigma^2 = 3) \end{aligned}$$

2. Single Match

Say there are 1000 decision points when writing Breakout. Assume: Each decision point is binary. Each student makes all 1000 decisions. For each decision there is a probability p that a student takes the more popular choice. What is the probability distribution for the number of matching decisions (we are going to refer to this as the “score”)? If possible, could you approximate this probability?

Let A_i be the event that decision point i is matched. We note that a match occurs when both students make the more popular choice or when both students make the less popular choice. $P(A_i) = P(\text{Both more popular}) + P(\text{Both less popular}) = p^2 + (1 - p)^2$.

Let M be a random variable for the number of matches. It is easy to see that each of the 1000 decisions is an independent Bernoulli experiment with probability of success $p' = p^2 + (1 - p)^2$. Therefore $M \sim \text{Bin}(1000, p')$.

We can use a Normal distribution to approximate a binomial. We approximate $M \sim \text{Bin}(1000, p')$ with Normal random variable $Y \sim N(1000p', 1000p'(1 - p'))$.

3. Maximum Match

When we look at two assignments, the probability of a false match is exceedingly small. What would the max similarity score look like when we compare one student to 5000 historical breakout submissions? Let X_i be the similarity score between a student who worked on their own and student i . Let Y be the highest match score between the student and all other submissions:

$$Y = \max_i X_i$$

The Central Limit Theorem tells us about the distribution of the sum of IID random variables. A more obscure theorem, the Fisher-Tippett-Gnedenko theorem, tells us about the *max* of IID random variables. It says that the max of IID exponential or normal random variables will be a “Gumbel” random variable.

$Y \sim \text{Gumbel}(\mu, \beta)$	The max of IID vars
$f(Y = k) = \frac{1}{\beta} e^{-(z+e^{-z})}$ where $z = \frac{k - \mu}{\beta}$	The Gumbel PDF

You are given data of 1300 students’ max scores from three quarters (we believe they all worked independently): $y^{(1)} \dots y^{(1300)}$. Set up (but do not solve) simultaneous equations we could solve to find the values of μ and β .

For this problem, we use Maximum Likelihood Estimator (MLE) to estimate the parameters $\theta = (\mu, \beta)$.

$$\begin{aligned}
 L(\theta) &= \prod_{i=1}^n f(Y^{(i)} = y^{(i)} \mid \theta) \\
 LL(\theta) &= \log \prod_{i=1}^n f(Y^{(i)} = y^{(i)} \mid \theta) \\
 &= \sum_{i=1}^n \log f(Y^{(i)} = y^{(i)} \mid \theta) \\
 &= \sum_{i=1}^n \log \frac{1}{\beta} e^{-(z_i + e^{-z_i})} && \text{where } z_i = \frac{y^{(i)} - \mu}{\beta} \\
 &= \sum_{i=1}^n \log \frac{1}{\beta} + \sum_{i=1}^n -(z_i + e^{-z_i}) \\
 &= -n \log(\beta) + \sum_{i=1}^n -(z_i + e^{-z_i})
 \end{aligned}$$

Now we must choose the values of $\theta = (\mu, \beta)$ that maximize our log-likelihood function. First, we need to find the first derivative of the log-likelihood function with respect to our parameters.

$$\begin{aligned}
 \frac{\partial LL(\theta)}{\partial \mu} &= \frac{\partial}{\partial \mu} \left[-n \log(\beta) + \sum_{i=1}^n -(z_i + e^{-z_i}) \right] \\
 &= \sum_{i=1}^n \frac{\partial}{\partial \mu} \left[-(z_i + e^{-z_i}) \right] \\
 &= \sum_{i=1}^n \frac{\partial}{\partial z_i} \left[-(z_i + e^{-z_i}) \right] \frac{\partial z_i}{\partial \mu} && \text{By the Chain Rule} \\
 &= \sum_{i=1}^n \left[-1 + e^{-z_i} \right] \left[-\frac{1}{\beta} \right] \\
 &= \frac{1}{\beta} \sum_{i=1}^n 1 - e^{-z_i}
 \end{aligned}$$

$$\begin{aligned}
\frac{\partial LL(\theta)}{\partial \beta} &= \frac{\partial}{\partial \beta} \left[-n \log(\beta) + \sum_{i=1}^n -(z_i + e^{-z_i}) \right] \\
&= -\frac{n}{\beta} + \sum_{i=1}^n \frac{\partial}{\partial \beta} \left[-(z_i + e^{-z_i}) \right] \\
&= -\frac{n}{\beta} + \sum_{i=1}^n \frac{\partial}{\partial z_i} \left[-(z_i + e^{-z_i}) \right] \frac{\partial z_i}{\partial \beta} && \text{By the Chain Rule} \\
&= -\frac{n}{\beta} + \sum_{i=1}^n \left[-1 + e^{-z_i} \right] \left[\frac{\mu - y^{(i)}}{\beta^2} \right] && \text{Where the last term equals } \frac{\partial z_i}{\partial \beta}
\end{aligned}$$

We want to find a simultaneous solution for both, but this is algebraically not possible. We will instead use an approximate method (gradient ascent) to solve for these, which will be taught next week.