

## Probability Distributions (Review)

### Discrete Distributions

We want to model the counts of two *independent* processes described by the random variables  $X$  and  $Y$ , with the following distributions:

$$X \sim \text{Binomial}(n, p)$$

$$Y \sim \text{Poisson}(\lambda)$$

1. What are the following expected values  $\mathbb{E}[X]$ ,  $\mathbb{E}[Y]$ ,  $\mathbb{E}[X + Y]$ ?
2. What is the probability of  $X < Y$  when  $n = 10, p = 0.5, \lambda = 5$ ?
3. What is the *conditional* probability distribution of  $\mathbb{P}(X + Y|Y)$ ?

### Continuous Distributions

We have a density function with the following form for  $x \in (-\infty, \infty)$ :

$$f(x) \propto e^{-x^2}$$

1. Estimate the proportionality constant for the density function to be a valid probability density. You are allowed to use symbolic algebra tools (e.g. `sympy`) for this task if needed (Hint: a probability density should integrate to 1). Numerical answers are also acceptable.
2. What is the expected value,  $\mathbb{E}[X]$ , of the random variable  $X$  with the above density function?
3. What is the variance of  $X$ ? (Hint: use the relationship  $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ )

## Likelihoods and P-values: Linear Regression

Linear regression is a central model for relating predictor variables ( $\mathbf{X}$ ) to response variables ( $\mathbf{y}$ ). The  $N$  ( $x_i, y_i$ ) identically and independently distributed (i.i.d) observations are related through the following simplified probability model:

$$y_i \sim \mathcal{N}(x_i\beta, \sigma)$$

Using this definition:

1. Write the expression for the log-likelihood of  $N$  datapoints in terms of  $\beta, \sigma$  using the normal density function.
2. Using the data provided in `data/pset1/linregA.tsv`, calculate the log-likelihood of the data with  $\beta = 0, \sigma = 1$ . Calculate the log-likelihood at the value of  $\hat{\beta}$  which maximizes the log-likelihood.
3. Calculate the log-likelihood when  $\beta = 0$ , and evaluate the log-likelihood ratio ( $\Lambda$ ). Comment on if the
4. For computational estimation, many approaches use “least-squares”, which tries to find the  $\beta$  which minimizes  $\sum_i (y_i - x_i\beta)^2$ . Calculate the least-squares estimate of  $\beta$  based on the current data. Comment on how this is similar (both mathematically and empirically) to the maximum-likelihood estimator.

## Multiple Testing Correction and False Discovery Rates

There are two relevant error rates in statistical genetics, the family-wise error rate (FWER) and the false-discovery rate (FDR):

$$FWER = P(R > 0 | H_0)$$

$$FDR = P(R/n \leq r | H_0)$$

, where  $R$  is a random variable representing the number of rejections of the null hypothesis and  $n$  is the total number of tests.

1. Assuming that you are conducting  $n = 10^6$  independent hypothesis tests, at a marginal testing threshold of  $\alpha = 0.05$  what is the FWER and FDR respectively? (Hint: consider the distribution of null p-values as Uniform and properties of the Bernoulli/Binomial distribution to show this analytically).
2. At what marginal threshold of  $\alpha$  does the FWER drop below 0.05 for the first time? What about the FDR? How is this related to the Bonferroni-threshold of  $5 \times 10^{-8}$ ?
3. Plot the FWER and FDR as a function of the marginal p-value threshold  $\alpha$  and indicate the Bonferroni threshold on the plot as a vertical line.

## Mixture Models

We want to model how genetic effects  $\beta$  are distributed in an experiment, and from some early looks it appears that some portion of them appear to be much higher or lower than you might expect under the standard  $\mathcal{N}(0, 1)$ . You choose to model each  $\beta$  as coming from a mixture of normal distributions:

$$\vec{\beta} \stackrel{\text{iid}}{\sim} (1 - \pi_0)\mathcal{N}(0, 1) + \pi_0\mathcal{N}(0, \sigma^2)$$

,

where the fraction  $\pi_0$  can be treated as an estimate of the proportion of effects that are not “null”.

1. Write down the log-likelihood function for the observed effects (e.g.  $\log(\mathbb{P}(\vec{\beta}|\pi_0, \sigma^2))$ )
2. Using the data in `data/pset1/sumstats_testA.tsv`, plot the log-likelihood as a function of  $\pi_0$  when  $\sigma^2 = 2$ . Provide an estimate of the value of  $\pi_0$  that maximizes the log-likelihood. (Hint: use `logsumexp` to avoid numerical underflow.)
3. Plot the log-likelihood as a function of  $\pi_0$  with  $\sigma^2 = 2$  for the data in `data/pset1/sumstats_testB.tsv`. How does the curve appear relative to the first dataset? Would you estimate there are greater or fewer *non-null* signals in this second dataset?
4. ♦ Fixing the value of  $\pi_0$  at 0.25, what is the value of  $\sigma^2$  which maximizes the log-likelihood? Is this different than the assumed value of 2? Plot the two-dimensional log-likelihood surface for all appropriate values of  $\pi_0, \sigma^2$  (Note: limit  $\sigma^2 < 10$ ). Comment on the likelihood-surface.
5. ♦♦ Up to now we have assumed that the non-null variance parameter ( $\sigma^2$ ) is known, but in practice it is not. If we knew the latent component assignment parameters ( $Z_i \in \{0, 1\}$ ) per  $\beta_i$ , we could estimate the following two parameters straightforwardly:

$$\begin{aligned}\pi_0 &= \frac{1}{n} \sum \delta(Z_i = 1) \\ \sigma^2 &= \frac{1}{N_k} \sum (\beta_i^2) \delta(Z_i = 1)\end{aligned}$$

, where  $N_k$  is the number of  $Z_i = 1$  out of the whole dataset. One decent guess of  $Z_i$ , might be  $P(Z_i = 1|\beta_i, \pi_0, \sigma^2)$ ! Show the form of this probability distribution for a single outcome  $\beta_i$  (Hint: use Bayes Rule ...)

6. ♦♦♦ Using the estimates of  $\gamma_i = P(Z_i = 1|\beta_i, \pi_0, \sigma^2)$ , we can actually revise the parameter estimates:

$$\begin{aligned}\pi_0 &= \frac{1}{n} \sum_i \gamma_i \\ \sigma^2 &= \frac{1}{N_k} \sum (\beta_i^2) \gamma_i\end{aligned}$$

, where  $N_k = \sum_i \gamma_i$ . Starting from initial parameter guesses of  $\pi_0 = 0.5, \sigma^2 = 5$ , calculate what the updated parameters ( $\pi_0', \sigma^{2'}$ ) would be (using the initial guesses to calculate  $\gamma_i$  for `data/pset1/sumstats_testA.tsv`). Calculate the log-likelihood using the initial parameters and the updated parameters, which one is higher? What if you repeat the process one more step (using  $\pi_0', \sigma^{2'}$  to re-calculate  $\gamma_i$  and update the parameters again), does the likelihood increase?<sup>1</sup>

---

<sup>1</sup>If you've gotten to this step, you've implemented your first EM-algorithm!