

Section 9 Solutions

1 The Most Important Features (17 points)

Let's explore saliency, a measure of how important a feature is for classification. We define the saliency of the i th input feature for a given example (\mathbf{x}, y) to be the absolute value of the partial derivative of the log likelihood, with respect to that input feature $|\frac{\partial LL}{\partial x_i}|$. Below, we show both input images and the corresponding saliency of their features (in this case, pixels) for an image classification model:



Consider a trained logistic regression classifier with weights θ that predicts binary class labels, y . In this question, we allow the values of \mathbf{x} to be real numbers, which doesn't change the algorithm at all (neither training nor testing).

- a. (4 points) What is the Log Likelihood of a single training example (\mathbf{x}, y) for this logistic regression classifier?

The log-likelihood here is the same as it was when we covered logistic regression in class:

$$LL(\theta) = y \cdot \log \sigma(\theta^T \cdot \mathbf{x}) + (1 - y) \log [1 - \sigma(\theta^T \cdot \mathbf{x})]$$

- b. (8 points) Calculate the saliency of a single feature (x_i) for one training example (\mathbf{x}, y) .

We can calculate the saliency for a single feature as follows. Note that we are taking the derivative with respect to x_i , not θ_i , which is different from the usual derivative we find for MLE optimization.

$$\begin{aligned}
 LL(\theta) &= y \log z + (1 - y) \log (1 - z) && \text{where } z = \sigma(\theta^T \cdot \mathbf{x}) \\
 \frac{\partial LL}{\partial x_i} &= \frac{\partial LL}{\partial z} \cdot \frac{\partial z}{\partial x_i} && \text{chain rule} \\
 &= \left(\frac{y}{z} - \frac{1-y}{1-z} \right) \cdot (z(1-z)\theta_i) && \text{partial derivatives} \\
 \text{saliency} &= \left| \left(\frac{y}{z} - \frac{1-y}{1-z} \right) z(1-z)\theta_i \right|
 \end{aligned}$$

- c. (5 points) Show that the ratio of saliency for features i and j is the ratio of the absolute value of their weights $\frac{|\theta_i|}{|\theta_j|}$.

We can take the ratio as follows using our expression above.

$$\begin{aligned}
 \text{saliency for feature } i, S_i &= \left| \left(\frac{y}{z} - \frac{1-y}{1-z} \right) z(1-z)\theta_i \right|, \text{ and same for } S_j \\
 \frac{S_i}{S_j} &= \frac{\left| \left(\frac{y}{z} - \frac{1-y}{1-z} \right) z(1-z)\theta_i \right|}{\left| \left(\frac{y}{z} - \frac{1-y}{1-z} \right) z(1-z)\theta_j \right|} = \frac{S_i}{S_j} = \frac{|\theta_i|}{|\theta_j|} \text{ by elimination}
 \end{aligned}$$

2 Timing Attack (23 points)

In this problem we will see how to crack a password in linear time by measuring how long the password check takes to execute (see code below).

```
# An insecure string comparison
def does_password_match(guess, password):
    n_guess = len(guess)
    n_password = len(password)
    if n_guess != n_password:
        return False # 4 lines executed to get here
    for i in range(n_guess):
        if guess[i] != password[i]:
            return False # 6 + 2i lines executed to get here
    return True # 5 + 2n lines executed to get here
```

Assume that our server takes T ms to execute any line in the code where $T \sim N(\mu = 5, \sigma^2 = 0.5)$ milliseconds. The amount of time taken to execute a line is always independent of other lines.

On our site, all passwords only use lower case letters and are between 5 and 10 letters long, inclusive. A hacker is trying to crack the root password which is “gobayes” by carefully measuring how long the code takes to tell her that her guesses are incorrect.

- a. (5 points) What is the distribution for the time it takes to execute k lines of code?

Let T_k be the amount of time to execute k lines. $T_k = \sum_{i=1}^k X_i$ where X_i is the amount of time to execute line i . $X_i \sim N(\mu = 5, \sigma^2 = 0.5)$. Since T_k is the sum of k independent normals:

$$T_k \sim N(\mu = 5k, \sigma^2 = 0.5k)$$

- b. (7 points) First, the hacker needs to find the length of the password. What is the probability that the time taken to check a guess of correct length (when the server executes 6 lines) is longer than the time taken to check a guess of an incorrect length (when the server only executes 4 lines)? Assume the first letter of the guess does not match the password’s first letter. Hint: $P(A > B) = P(A - B > 0)$.

Time to run 6 lines of code: $T_6 \sim N(\mu = 30, \sigma^2 = 3)$

Time to run 4 lines of code: $T_4 \sim N(\mu = 20, \sigma^2 = 2)$ Then we apply a linear transform to T_4 so we can subtract it from T_6 by “adding” two Normal RVs.

$$-T_4 \sim N(\mu = -20, \sigma^2 = 2)$$

$$T_6 - T_4 \sim N(\mu = 10, \sigma^2 = 5)$$

$$\begin{aligned} P(T_6 > T_4) &= P(T_6 - T_4 > 0) \\ &= 1 - F_{T_6 - T_4}(0) \\ &= 1 - \Phi\left(\frac{0 - 10}{\sqrt{5}}\right) \approx 1.0 \end{aligned}$$

- c. (8 points) Now that our hacker knows the length of the password, to get the actual string, she will try to determine one letter at a time, starting with the first letter. To start, the hacker tries the string “aaaaaaa” and sees that it takes 27ms. Based on this timing, how much more probable is it that first character did not match (server executes 6 lines) than the first character did match (server executes 8 lines)? Assume that all letters in the alphabet are equally likely to be the first letter.

Let M be the event that the first letter matched. The problem’s wording indicates that we are looking for a ratio between the probability of M^C to the probability of M , conditioned on observing a code execution time of 27ms. Let $T_?$ be the time it took to check the password, with the number of lines executed being unknown. To calculate this ratio, we can apply Bayes’ Theorem:

$$\begin{aligned} \frac{P(M^C|T_? = 27)}{P(M|T_? = 27)} &= \frac{f(T_? = 27|M^C)P(M^C)}{f(T_? = 27|M)P(M)} \\ &= \frac{f(T_? = 27|M^C)\frac{25}{26}}{f(T_? = 27|M)\frac{1}{26}} \\ &= 25 \cdot \frac{f(T_? = 27|M^C)}{f(T_? = 27|M)} \end{aligned}$$

From here, we can reason about how $T_?$ is distributed if we know M or M^C . If M happened, then $T_? = T_8$; otherwise, $T_? = T_6$. So we can re-write the above as:

$$\begin{aligned} &= 25 \cdot \frac{f(T_6 = 27)}{f(T_8 = 27)} \\ &= 25 \cdot \frac{\frac{1}{\sqrt{6\pi}}e^{-\frac{(27-30)^2}{6}}}{\frac{1}{\sqrt{8\pi}}e^{-\frac{(27-40)^2}{8}}} \\ &= 25 \cdot \frac{\sqrt{8}}{\sqrt{6}} \cdot \frac{e^{-\frac{9}{6}}}{e^{-\frac{169}{8}}} \\ &\approx 9.6 \text{ million} \end{aligned}$$

- d. (3 points) If it takes the hacker 6 guesses to find the length of the password, and 26 guesses per letter to crack the password string, how many attempts does she need to crack our password, “gobayes”? Yikes!

Since the password is length 7: $6 + 7 \cdot 26 = 188$.

3 Bayesian Carbon Dating (34 points)

We are able to know the age of ancient artifacts using a process called carbon dating. This process involves a lot of uncertainty! Living things have a constant proportion of a molecule called C14 in them. When living things die those molecules start to decay. The time to decay in years, T , of a C14 molecule is distributed as an exponential. $T \sim \text{Exp}(\lambda = 1/8267)$.

- a. (5 points) Consider a single C14 molecule. What is the probability that it decays within 500 years?

Using the CDF of the Exponential, $P(T \leq 500) = 1 - e^{-\frac{1}{8267} * 500} \approx 0.0587$

- b. (8 points) C14 molecules decay independently. A particular sample started with 100 molecules. What is the probability that exactly 95 are left after 500 years? Let p be your answer to part a.

Because we have a fixed number of molecules, and each molecule could independently decay with a consistent probability, the number of molecules that decay can be modeled using a Binomial.

$$X \sim \text{Bin}(n = 100, p = 0.0587)$$

$$P(X = 5) = \binom{100}{5} (p)^5 (1 - p)^{95}$$

- c. (6 points) Write pseudocode for a function `pr_measure_given_age(m, age)` which returns $P(M = m | A = \text{age})$, the probability that exactly m molecules are left out of the original 100 after exactly age number of years.

Consider each of the 100 original molecules as an independent trials of a C14 molecule. We can model `pr_measure_given_age(m, age)` as a binomial with parameters 100 and probability of a single element left after some age.

```
def pr_measure_given_age(m, age):
    # Find the exponential probability
    p = 1 - np.exp(-(1/8267) * age)

    # Find the amount of successes after 100 trials.
    pr_m_given_age = scipy.stats.binom.pmf(100 - m, 100, p)

    return pr_m_given_age
```

- d. (15 points) You observe 95 C14 molecules in a sample. You assume that the sample originally had 100 C14 molecules when it died. Write pseudocode for a function `age_belief()` that returns the full probability distribution $P(A = i | M = 95)$, where $A = i$ is the event that the sample died i years ago (here age is a discrete random variable). Use the function `pr_measure_given_age(m, age)` from part c. Your prior belief is that the sample's age was between $A = 500$ and $A = 600$ inclusive and that every year in that range is equally likely.

First, apply Bayes' Theorem:

$$P(A = i | M = 95) = \frac{P(M = 95 | A = i)P(A = i)}{P(M = 95)}$$

Our goal is to write code that calculates this equation for all i . What we know:

- $P(M = 95 | A = i)$ can be found with `pr_measure_given_age(95, i)`.
- $P(A = i) = \frac{1}{101}$; the denominator is 101 because our prior is over 500 to 600 **inclusive**.
- $P(M = 95)$ can be disregarded since it's a normalization term (we'll normalize the returned array anyway).
- To return a full probability distribution, we'll need to loop over all possible values for A

With that, here's the code:

```
def age_belief(m = 95):
    probs = dict()

    # loop over all possible ages
    for age in range(500, 600 + 1):
        p_m_given_a = pr_m_given_age(m, age)
        probs[age] = p_m_given_a * (1 / 101)

    # normalize
    probs_sum = sum(probs.values())
    for age in range(500, 600 + 1):
        probs[age] = probs[age] / probs_sum

    return probs
```