

# Homework 6

SHREYA CHETAN PAWASKAR

[pawaskas@uci.edu](mailto:pawaskas@uci.edu)

12041645

## Question 1:

a)

- Let  $\pi_{k\ell} = p(z_{t+1} = \ell \mid z_t = k)$ .
- Using the Markov structure of the model:

$$p(z_1, z_2, z_3) = p(z_1) p(z_2 \mid z_1) p(z_3 \mid z_2)$$

- Substitute the values as 0.8 and 0.9,

$$p(z_1 = z_2 = z_3) = 0.5\pi_{11}\pi_{11} + 0.5\pi_{22}\pi_{22} = 0.725$$

- Answer = 0.725**
- 

b)

- Let  $\phi_{ka} = p(x_t = a \mid z_t = k)$ .
- Given the variables  $z_1, x_1$  and  $x_2$  are conditionally independent, let's begin by assessing the impact of the observation

$$p(z_1 \mid x_1 = 1) \propto p(z_1) p(x_1 = 1 \mid z_1) \propto [\phi_{11}, \phi_{21}] \propto [0.25, 0.7]$$

- We apply the Markov state dynamics.

$$p(z_2 \mid x_1 = 1) = \sum_{z_1} p(z_2 \mid z_1) p(z_1 \mid x_1 = 1) \propto [\pi_{11}\phi_{11} + \pi_{21}\phi_{21}, \pi_{12}\phi_{11} + \pi_{22}\phi_{21}] \propto [0.27, 0.68]$$

$$p(z_2 \mid x_1 = 1) = \text{Cat}(z_2 \mid [0.284, 0.716])$$

- Then we apply the observation model:

$$p(x_2 = 1 \mid x_1 = 1) = \sum_{z_2} p(x_2 = 1 \mid z_2) p(z_2 \mid x_1 = 1) = 0.284\phi_{11} + 0.716\phi_{21} \approx 0.57$$

- When  $x_1$  is unobserved then the probability that  $x_2 = 1$  is lower:

$$p(z_2) = \sum_{z_1} p(z_2 \mid z_1) p(z_1) = \text{Cat}(z_2 \mid [0.45, 0.55])$$

$$p(x_2 = 1) = \sum_{z_2} p(x_2 = 1 \mid z_2) p(z_2) = 0.45\phi_{11} + 0.55\phi_{21} \approx 0.50$$

- Answer = 0.5**
-

c)

- From the structure of the HMM model and using Bayes' rule, we have:

$$\begin{aligned}
 p(z_1 | x_1, x_2) &= \sum_{z_2} p(z_1, z_2 | x_1, x_2) \propto \sum_{z_2} p(z_1, z_2, x_1, x_2) \\
 &\propto p(z_1) p(x_1) \sum_{z_2} p(z_2 | z_1) p(x_2 | z_2)
 \end{aligned}$$

- Evaluating the second term gives

$$p(x_2 = 4 | z_1) \propto [\pi_{11}\phi_{14} + \pi_{12}\phi_{24}, \pi_{21}\phi_{14} + \pi_{22}\phi_{24}] \propto [0.220, 0.115]$$

- Next we have:

$$p(z_1 | x_1 = 2, x_2 = 4) \propto p(x_1 = 2 | z_1) p(x_2 = 4 | z_1) \propto [0.055, 0.0115]$$

- Finally, normalizing these terms gives us:

$$p(z_1 = 1 | x_1 = 2, x_2 = 4) \approx 0.83.$$

Since neither roll equals 1, it is likely that the fair die is in use.

---

d)

- The most likely state sequence is  $z_t = 2$  for all  $t$ , as  $p(z_1)$  is uniform. Also  $\pi_{22} \geq \pi_{k\ell}$  for all  $k, \ell$ .
  - $\phi_{21} \geq \phi_{ka}$  for all  $k, a$ , so the corresponding most likely observation sequence is  $x_t = 1$  for all  $t$
  - The most likely configuration of the HMM is that the biased die is used at all times, and all rolls = 1.
-

## Question 2:

Code:

```
import numpy as np
import random
import matplotlib.pyplot as plt
from text_utils import *
from hmmlearn import hmm

load_from_saves = False
model_to_save = False

# Load the data
train_text = read_text_file("aliceTrain.txt")
test_text = read_text_file("aliceTest.txt")
train_data = text_to_num(train_text)
test_data = text_to_num(test_text)

num_valid_chars = 27
max_num_em_steps = 500
em_change_tolerance = 1e-6

m_values = [1, 5, 10, 15, 20, 30, 40, 50, 60]
models = []
training_log_likelihoods = []

if model_to_save:
    for M in m_values:
        model = hmm.CategoricalHMM(n_components=M, n_iter=max_num_em_steps,
tol=em_change_tolerance)
        model.fit(np.reshape(train_data, (-1, 1)))
        models.append(model)
        training_log_likelihoods = list(model.monitor_.history)
        model_save = {
            "n_features": model.n_features,
            "emissionprob_": model.emissionprob_,
            "transmat_": model.transmat_,
            "startprob_": model.startprob_,
            "tll": np.asarray(list(model.monitor_.history))
        }
        np.save(f"model_m_{M}.npy", model_save)
else:
```

```

for M in m_values:
    data = np.load(f"model_m_{M}.npy", allow_pickle=True).item()
    model = hmm.CategoricalHMM(n_components=M, verbose=True)
    model.n_features = data["n_features"]
    model.emissionprob_ = data["emissionprob_"]
    model.transmat_ = data["transmat_"]
    model.startprob_ = data["startprob_"]
    t11 = data["t11"]
    models.append(model)
    training_log_likelihoods.append(t11.tolist())

AIC_train_values = np.zeros(len(m_values))
BIC_train_values = np.zeros(len(m_values))
test_log_likelihoods = np.zeros(len(m_values))
max_train_log_likelihoods = np.zeros(len(m_values))

for index_of_M in range(len(m_values)):
    M = m_values[index_of_M]
    num_param = (M * (M - 1)) + (M * (num_valid_chars - 1)) + (M - 1)
    AIC_train_values[index_of_M] =
np.max(np.asarray(training_log_likelihoods[index_of_M])) - num_param
    BIC_train_values[index_of_M] =
np.max(np.asarray(training_log_likelihoods[index_of_M])) - 0.5 * num_param *
np.log(train_data.shape[0])
    test_log_likelihoods[index_of_M] = models[index_of_M].score(np.reshape(test_data,
(-1, 1)))
    max_train_log_likelihoods[index_of_M] =
np.max(np.asarray(training_log_likelihoods[index_of_M]))

# Plot log-likelihood, AIC, BIC, and test log-likelihood
plt.plot(m_values, max_train_log_likelihoods, label="Train Log Likelihood")
plt.plot(m_values, AIC_train_values, label="AIC")
plt.plot(m_values, BIC_train_values, label="BIC")
plt.plot(m_values, test_log_likelihoods, label="Test Log Likelihood")
plt.xlabel("Number of States")
plt.ylabel("Log Likelihood")
plt.legend()
plt.show()

# Identify best models based on AIC and BIC
best_model_AIC = np.argmax(AIC_train_values)

```

```

best_model_BIC = np.argmax(BIC_train_values)

print(f"\nBest Model according to AIC: M = {m_values[best_model_AIC]}")
print(f"Best Model according to BIC: M = {m_values[best_model_BIC]}\n")

test_indices = [0, best_model_AIC, best_model_BIC, len(m_values) - 1]
test_indices.sort()
sample_size = 500

for i in test_indices:
    sample, _ = models[idx].sample(sample_size)
    sample_text = num_to_text(sample.flatten())
    print(f"Sample from HMM with M = {m_values[i]}:\n{sample_text}\n")

erase_rate = 0.2
erase_ind = np.random.uniform(size=(len(test_text),)) < erase_rate
num_erased = np.sum(erase_ind)
noisy_text = np.copy(test_text)
noisy_text[erase_ind] = "*"
noisy_data = text_to_num(noisy_text)

denoised_text_list = []
num_corrected = np.zeros(len(test_indices))

for _char, i in enumerate(test_indices):
    M = m_values[i]
    model = models[i]
    erase_prob = np.eye(num_valid_chars, num_valid_chars + 1)
    erase_prob[:, :-1] *= (1.0 - erase_rate)
    erase_prob[:, -1] = erase_rate
    obs_mat = np.matmul(model.emissionprob_, erase_prob)
    noise_fix_model = hmm.CategoricalHMM(n_components=M, verbose=True)
    noise_fix_model.n_features = num_valid_chars + 1
    noise_fix_model.emissionprob_ = obs_mat
    noise_fix_model.transmat_ = model.transmat_
    noise_fix_model.startprob_ = model.startprob_
    noisy_prob = noise_fix_model.predict_proba(np.reshape(noisy_data, (-1, 1)))

    denoised_data = np.copy(noisy_data)
    # Find the most probable missing letters
    for i in range(noisy_data.shape[0]):

```

```

    if noisy_text[i] == "*":
        letter_probability = np.matmul(model.emissionprob_.T, noisy_prob[i])
        letter = np.argmax(letter_probability)
        denoised_data[i] = letter

num_corrected[_char] = num_erased - np.sum(denoised_data != test_data)
denoised_text
denoised_text = num_to_text(denoised_data.flatten())
denoised_text_list.append(denoised_text)

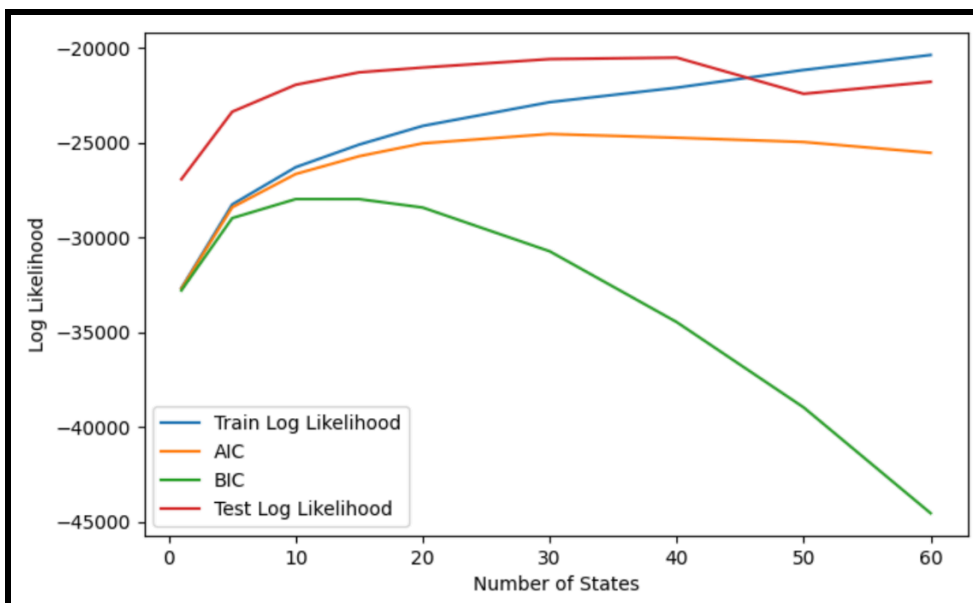
sample_size = 500

print("Denoised Outputs and Correction Metrics")
print("Original Noisy Text:")
print("".join(noisy_text[:sample_size].tolist()))
print()

for i in range(len(test_indices)):
    percent_correct = (float(num_corrected[i]) / float(num_erased)) * 100
    print(f"M = {m_values[test_indices[i]]} | Letters Corrected: {int(num_corrected[i])} | Percentage Corrected: {percent_correct:.2f}%")
    print(denoised_text_list[i][:sample_size])

```

A) Plot:



B)

- The number of free parameters in a multinomial hidden Markov model (HMM) consists of the initial state distribution, the state transition matrix, and the emission distribution matrix.
  - Including all of them, the number of free parameters is  

$$d = (M - 1) + M(M - 1) + M(W - 1).$$
  - As plotted in the above graph, the Bayesian Information Criterion (BIC) imposes a larger penalty on the number of parameters than the Akaike Information Criterion (AIC).
  - Consequently, the BIC favors simpler models with fewer parameters compared to the AIC, which is more lenient towards complex models with a larger number of parameters.
- 

C)

- AIC is largest for M=30 and BIC is largest for M=10.
  - Test log-likelihood is largest for =30s predicted by AIC
  - But it's very low for M=10
  - In this case, AIC performed better.
  - Theoretical results suggest that AIC is often better for moderate datasets, but it can be inconsistent, and thus inferior to BIC for very large datasets.
- 

d):

**HMM with M = 1**

seuhhs\_aeo\_\_sls\_ended\_oshiuka\_na\_ms\_\_eea\_wel\_a\_ernl\_is\_\_eani\_wbe\_lmte\_\_oeauhi\_r  
uinstg\_h\_oh\_ieo\_ywttie\_abytos\_hao\_oui\_ietdeierdysoehitashi\_r\_ntp\_edes\_hcresestm  
awweesia\_oss\_onco\_htgnhheaefverr\_\_d\_ahewia\_tahnsa\_ait\_\_\_h\_elt\_n\_ennehhmsrdi\_e\_k  
nnwnaug\_l\_uglienarlii\_st\_o\_m\_\_gttryotat\_enisget\_ednnak\_py\_er\_lees\_igiposiiu\_  
ylvpcio\_er\_e\_dotovweauxmaeirehelrregleo\_prasedts\_mwtd\_haueleegr\_\_\_ie\_so\_rd\_rut  
mi\_\_idlut\_eteteiciian diba\_te\_\_daes\_oe\_m\_whrfwmonlste\_out\_irl\_icesu\_eca\_\_oEU\_se\_  
\_et\_\_h\_ee\_saddormnno\_j\_ywe

**HMM with M = 20**

yonedfeave\_dailg\_ce\_at\_anrefire\_hut\_ande\_howde\_we\_cor\_sousonlyg\_eesbe\_tid\_yer\_mh  
ass\_aalem\_oy\_ant\_doy\_br\_tre\_linomet\_ir\_tawf\_tal\_anl\_any\_tovilf\_bplee\_aluch\_wind  
\_a\_thayestind\_bis\_he\_attrith\_bet\_mve\_he\_host\_mro\_the\_thy\_limy\_dhetherthowkesdik  
e\_i\_msk\_toot\_yautef\_be\_bre\_tau\_alig\_hit\_the\_maketath\_tes\_han\_ang\_it\_dhe\_let\_tol  
g\_ad\_hacy\_hat\_ikt\_swe\_hepd\_an\_sure\_asuaalilt\_bu\_toumh\_wort\_shed\_he\_tory\_pes\_you  
\_le\_ald\_aplene\_ta\_thyatowd\_it\_gher\_no\_liowte\_laa\_nhime\_dont\_sounysrent\_tomo\_he\_  
shirn\_an\_shar\_wint\_pbe\_ggl

**HMM with M = 30**

planned\_thint\_i\_den\_tangt\_i\_thymwas\_le\_thastine\_alld\_bwer\_the\_tee\_widis\_dene\_ni  
\_hare\_thed\_to\_plot\_uouoy\_dor\_os\_a\_mastingn\_vosed\_ot\_hasicg\_mars\_ttdce\_outh\_alit  
e\_mavle\_hutisk\_maplinh\_i\_pat\_letreadxot\_i\_ore\_sar\_coemsar\_ud\_the\_shicg\_en\_it\_ou  
led\_nersawe\_surlayhats\_maid\_a\_doene\_putweriersads\_miicg\_the\_whe\_feavine\_neay\_id  
\_thkd\_udhad\_wug\_o\_satse\_qhe\_ereptor\_a\_hand\_muld\_sas\_ot\_sadint\_therdnw\_a\_dis\_an\_  
iss\_yoy\_trerie\_verius\_ot\_isper\_thklece\_che\_cup\_hor\_imyd\_y\_wouoed\_lin\_i\_i\_loagai  
d\_dem\_inwond\_sveand\_itfe\_d

**HMM with M = 60**

tauebed\_i\_cany\_the\_yortitert\_astring\_at\_lepse\_getghked\_mapt\_a\_rile\_sormoulded\_a  
 s\_to\_alice\_sughts\_that\_om\_is\_ther\_fou\_sint\_ok\_the\_hith\_thrust\_the\_nou\_yoct\_on\_  
 untes\_co\_thes\_vavy\_be\_she\_she\_vamh\_noulk\_on\_it\_a\_ksovene\_il\_moull\_alike\_ware\_no  
 r\_cou\_yen\_you\_lubp\_evothe\_come\_it\_to\_thint\_hetned\_the\_soup\_emeng\_twer\_thed\_tork  
 e\_hety\_hath\_shoy\_mo\_mou\_reep\_shing\_on\_up\_med\_at\_dcace\_on\_alace\_tave\_trked\_beeke  
 r\_dirnie\_the\_mis\_i\_said\_raad\_usfoich\_tciupon\_i\_ware\_and\_alice\_wory\_youryord\_bus  
 ses\_piut\_twe\_toce\_faly\_rem

- The unigram model captures only the basic frequencies of letters.
- More complicated models learn states which emit a number of common short words, but are still very far from real English text.
- A much larger training corpus would be needed to learn more realistic models.

**e) Code:**

```

erase_rate = 0.2
erase_ind = np.random.uniform(size=(len(test_text),)) < erase_rate
num_erased = np.sum(erase_ind)
noisy_text = np.copy(test_text)
noisy_text[erase_ind] = "*"
noisy_data = text_to_num(noisy_text)

denoised_text_list = []
num_corrected = np.zeros(len(test_indices))

for _char, i in enumerate(test_indices):
    M = m_values[i]
    model = models[i]
    erase_prob = np.eye(num_valid_chars, num_valid_chars + 1)
    erase_prob[:, :-1] *= (1.0 - erase_rate)
    erase_prob[:, -1] = erase_rate
    obs_mat = np.matmul(model.emissionprob_, erase_prob)
    noise_fix_model = hmm.CategoricalHMM(n_components=M, verbose=True)
    noise_fix_model.n_features = num_valid_chars + 1
    noise_fix_model.emissionprob_ = obs_mat
    noise_fix_model.transmat_ = model.transmat_
    noise_fix_model.startprob_ = model.startprob_
    noisy_prob = noise_fix_model.predict_proba(np.reshape(noisy_data, (-1, 1)))

    denoised_data = np.copy(noisy_data)
    # Find the most probable missing letters

```



```

for i in range(noisy_data.shape[0]):
    if noisy_text[i] == "*":
        letter_probability = np.matmul(model.emissionprob_.T, noisy_prob[i])
        letter = np.argmax(letter_probability)
        denoised_data[i] = letter

num_corrected[_char] = num_erased - np.sum(denoised_data != test_data)
denoised_text
denoised_text = num_to_text(denoised_data.flatten())
denoised_text_list.append(denoised_text)

```

f)

- The given equations represent the marginal and conditional probability distributions, where the last equality indicates that the letters are conditionally independent given the Markov states.

$$p(x_t | y) = \sum_{z_t} p(x_t, z_t | y) = \sum_{z_t} p(x_t | z_t, y) p(z_t | y) = \sum_{z_t} p(x_t | z_t, y_t) p(z_t | y)$$

where the last equality means letters  $x_t$  are conditionally independent given the Markov states  $z_t$ .

- The result then follows because  $p(x_t | z_t, y_t = *) = p(x_t | z_t)$
- This follows from the uninformative nature of  $p(y_t = * | x_t)$ .  
The marginal probability distribution of the letters depends solely on the transition probabilities between the Markov states and the emission probabilities of the letters given those states

g)

- The minimum-error decision rule chooses  $\hat{x}_t = \arg \max_{x_t} p(x_t | y)$  to maximize the marginal posterior distribution from question C.
- It selects the option that has the highest marginal posterior probability, which is the probability of that option being true.

h)

- The models with M=1, 10, 30, and 60 states correct 19.50 %, 45.29 %, 48.62 % and 51.43% of the erased characters, respectively.
- Choosing one of the 27 characters at random would correct 3.7%.
- The unigram model improves on this by always picking the most common character, space, but the HMM models do much better:

Original, Noisy Text:

\*he\_king\_an\*\_queen\_of\*he\*r\*\*\*w\*re\_eat\*d\*on\_thei\*\_thron\*\_wh\*n\_th\*y\*arr\*ve\*\_with  
\_a\_g\*eat\_crowd\_a\*sembl\*d\*a\*out\_th\*\*\_al\*\_sorts\*of\*\*i\*\*l\*\_b\*\*ds\*\_nd\_beast\*\_a\*\_el  
l\_as\_the\_whole\_pack\*of\_card\*\_th\*\_knave\*\*\*\_\*tanding\_bef\*r\*\_them\*i\*\_chains\*\_it\*\_  
a\_soldie\*\*o\*\_each\*\_ide\*\_o\_g\*ar\*\*\*im\*\_nd\*\_ear\*\_he\*k\*\*g\*\_a\*\_the\*white\*r\*bbi\*\_with  
\_a\_trumpe\*\_in\_one\*ha\*d\_and\_a\*scrol\*\_of\_parchment\*in\_the\_othe\*\_i\*\_th\*\_very\*middl  
e\_of\*the\*\_ourt\*wa\*\_a\*\_a\*l\*\_w\*\*h\*\_l\*r\*e\*d\*sh\_o\*\_ta\*\*s\_upon\_it\_th\*y\*l\*oked\_s\*\_g\*\_  
od\_that\_it\*\_ade\_ali\*e\_qui\*

M = 1

Letters Corrected: 375

Percentage Corrected: 19.50%

\_he\_king\_an\_\_queen\_of\_he\_r\_\_w\_re\_\_eat\_d\_on\_thei\_\_thron\_\_wh\_n\_th\_y\_arr\_ve\_\_with  
\_a\_g\_eat\_crowd\_a\_sembl\_d\_a\_out\_th\_\_al\_\_sorts\_of\_\_i\_\_l\_\_b\_\_ds\_\_nd\_beast\_\_a\_\_el  
l\_as\_the\_whole\_pack\_of\_card\_\_th\_\_knave\_\_tanding\_bef\_r\_\_them\_i\_\_chains\_\_it\_\_  
a\_soldie\_\_o\_\_each\_\_ide\_\_o\_g\_ar\_\_im\_\_nd\_\_ear\_\_he\_k\_\_g\_\_a\_\_the\_white\_r\_bbi\_\_with  
\_a\_trumpe\_\_in\_one\_ha\_d\_and\_a\_scrol\_\_of\_parchment\_in\_the\_othe\_\_i\_\_th\_\_very\_middl  
e\_of\_the\_\_ourt\_wa\_\_a\_\_a\_l\_\_w\_\_h\_\_l\_r\_e\_d\_sh\_o\_\_ta\_\_s\_upon\_it\_th\_y\_l\_oked\_s\_\_g\_\_  
od\_that\_it\_\_ade\_ali\_e\_qui\_\_

M = 20

Letters Corrected: 871

Percentage Corrected: 45.29%

the\_king\_and\_queen\_of\_he\_re\_\_wore\_heat\_dhon\_theit\_thrond\_when\_they\_arrevet\_with  
\_a\_gheat\_crowd\_atsembled\_alout\_thet\_ale\_sorts\_of\_tinile\_breds\_and\_beaste\_at\_tel  
l\_as\_the\_whole\_pack\_of\_carde\_the\_knave\_\_he\_atanding\_bef\_re\_them\_it\_chains\_tite\_  
a\_soldie\_\_ou\_each\_tide\_to\_ghar\_thim\_and\_tear\_thetking\_tau\_the\_white\_rebbie\_with  
\_a\_trumpet\_in\_one\_hand\_and\_a\_scrol\_d\_of\_parchment\_in\_the\_othe\_t\_it\_the\_very\_middl  
e\_of\_the\_tourt\_wau\_a\_tauld\_wich\_a\_lorte\_dish\_ot\_tares\_upon\_it\_they\_lookd\_so\_gh  
od\_that\_it\_tade\_alice\_quin

M = 30

Letters Corrected: 935

Percentage Corrected: 48.62%

the\_king\_and\_queen\_of\_he\_re\_\_were\_deated\_on\_theid\_throne\_when\_they\_arr\_ver\_with  
\_a\_gteat\_crowd\_a\_sembl\_d\_alout\_thed\_all\_sorts\_of\_tinkle\_berds\_and\_beaste\_at\_del  
l\_as\_the\_whole\_pack\_of\_carde\_the\_knave\_t\_e\_tanding\_befere\_them\_it\_chains\_tite\_  
a\_soldie\_\_ot\_each\_aide\_to\_geare\_\_im\_and\_dear\_the\_king\_sad\_the\_white\_rebbie\_with  
\_a\_trumper\_in\_one\_hand\_and\_a\_scrolg\_of\_parchment\_in\_the\_other\_it\_the\_very\_middl  
e\_of\_the\_uourt\_war\_a\_sakle\_winh\_a\_larke\_dish\_ot\_tares\_upon\_it\_they\_leoked\_si\_gt  
od\_that\_it\_sade\_aline\_quin

M = 60

Letters Corrected: 989

Percentage Corrected: 51.43%

the\_king\_and\_queen\_of\_hearle\_ware\_beated\_on\_theid\_throng\_whin\_they\_arrever\_with  
\_a\_gbeat\_crowd\_a\_sembledea\_out\_thet\_aly\_sortsiof\_singly\_beads\_and\_beaste\_at\_wel  
l\_as\_the\_whole\_pack\_of\_cardh\_the\_knave\_tae\_otanding\_befere\_them\_it\_chains\_wite\_  
a\_soldie\_\_ou\_each\_wide\_yo\_grare\_\_im\_and\_bear\_theeking\_had\_the\_whiterrubbie\_with

\_a\_trumped\_in\_one\_haad\_and\_atscroly\_of\_parchment\_in\_the\_othed\_it\_the\_very\_middl  
e\_of\_the\_dourt\_war\_a\_ealle\_warh\_a\_larte\_d\_sh\_on\_tates\_upon\_it\_they\_lookd\_so\_go  
od\_that\_it\_wade\_alice\_quin