

# Breaking substitution ciphers with Markov models

Algorithms description and implementation notes

Enrico Rollando (e.rollando@outlook.com)

August 6, 2021

## Abstract

Cryptography is a useful playground for applying statistical methods. This document uses the simple substitution cipher as a toy problem and shows how to decrypt messages without knowing the key (i.e., without known plaintext). It focuses on Markov model-based approaches, reviews the underlying theory, and provides concrete implementations and usage examples.

## Introduction

Can we break a simple substitution cipher automatically? “To break” here means discovering and exploiting weaknesses to restore—at least in part—the original text.

For scope, consider plaintext and ciphertext over uppercase English letters plus space (26+1 symbols). In the base setting, space maps to itself.

A *simple substitution cipher* is a mapping function  $S_k(\text{text}) \rightarrow \text{ciphertext}$  where  $k$  is the key. It is a bijection on the alphabet (equivalently, a permutation of its symbols). Given the ciphertext only, the goal is to recover the inverse mapping  $S_k^{-1}(\text{ciphertext})$  without any known plaintext.

To guide the search, define a function  $\text{fit}(\text{text})$  that scores the current guess. Ideally this would answer “how likely is this the original plaintext?”, but a practical proxy is: “how plausible is this as English text?”. This score lets us compare different candidate decryptions and prefer the most plausible one.

A natural way to build such a score is to use language statistics: how often letters and  $n$ -grams appear in English, and how likely one character is to follow another. We then score candidate decryptions using these statistics.

There is still a problem: the number of possible keys is huge (26! for the English alphabet), far too many to try exhaustively. Other techniques are needed; this document presents a few of them.

Let us first define the cipher and briefly introduce the language model, which is the main tool used to exploit the cipher’s weakness.

## The simple substitution cipher

A simple substitution cipher [1] is a one-to-one mapping between two alphabets: the plaintext alphabet and the ciphertext alphabet. This mapping is the key of the cipher.

Spaces are sometimes removed from plaintext before encryption, as they provide information about word boundaries, which in turn helps recovery (the Patristocrat variant).

Related work also treats substitution decipherment as a statistical inference problem, where the plaintext structure and the key are latent and must be inferred jointly [2].

$P = [a_1, \dots, a_{26}]$ , example: ABCDEFGHIJKLMNOPQRSTUVWXYZ  
 $C = [c_1, \dots, c_{26}]$ , example: HQUVNCFYYPXJDIMOBWLAZTGSRKE

$$S_k : P \rightarrow C \quad (1)$$

$$S_k(\text{"THISISATEST"}) = \text{"ZYPAPAHZNAZ"} \\ S_k^{-1}(\text{"ZYPAPAHZNAZ"}) = \text{"THISISATEST"}$$

## The Language Model

In any language there are implicit rules we use when we write or interpret sentences. For example, the incomplete sentence “I like to write my letters by ...” is likely to end with “hand”. The preceding words provide context, which makes some continuations more plausible than others.

At the character level, the same idea lets us read an incomplete word like “uncomf.rt.b.y”: plausibility can be interpreted statistically as “which character is most likely here, given its context?”.

A language model is tied to a specific language and captures such regularities using frequencies at the character or word level. As an example, consider a bigram English language model  $LM_{EN}$ :

- bigram frequency:
  - how frequent is a bigram  $ab$  in English text?  $P_{LM}(ab) = \frac{\text{count(occurrences of } ab\text{)}}{\text{len(text)} - 1}$
  - for a bigram  $ab$ , how frequently does  $b$  follow  $a$ ?

$$P_{LM}(b | a) = \frac{\text{count}(ab)}{\text{count}(a)} \quad (2)$$

## The weakness

A simple substitution cipher always encrypts the same plaintext symbol in the same way. As a result, it preserves a surprising amount of structure from the plaintext: equal letters remain equal, repeated patterns remain repeated, and, more generally, the plaintext  $n$ -gram statistics are preserved **up to a relabeling of symbols** induced by the unknown key.

For example, if a ciphertext bigram is a double letter (e.g., XX), then the corresponding plaintext bigram must also be a double letter. By itself this is not enough to identify the letters, but across a whole message these constraints become informative, especially when combined with an English language model.

To exploit this, we can introduce a generative model  $G$  for English plaintext, controlled by parameters  $\Theta$ , and search for a key (or plaintext) that makes the decoded text both (i) consistent with the ciphertext under a substitution mapping and (ii) plausible under  $G$ . In practice, this is done iteratively: we refine  $\Theta$  and/or the key to increase the model score (likelihood or plausibility) of the observed ciphertext under the induced plaintext statistics.

## Results at a glance

Method	Cipher length	Iterations/time	Outcome
EM (Baum-Welch)	~120 chars	200 iters (<1 min)	Readable plaintext; some letters may need manual swaps
MH (substitution)	~40 chars	500k swaps (~1 min)	Gradual convergence to legible text; prints best candidates
MH (patristocrat)	~30 chars	5k insert/remove (~seconds)	Recovers likely spacing; ambiguity remains on very short texts
LM pickle prep	corpus-dependent	varies	N-gram log-probability tables reused across runs

# Reproducibility

## Goal & Scope

- Ciphertext alphabet: A-Z + space. Spaces are preserved for the base examples; the patristocrat variant removes them.
- Corpus: plain-text English corpus (e.g., the Brown corpus [3]), preprocessed by uppercasing, keeping only A-Z and spaces, and collapsing whitespace.
- Success criteria: recover readable plaintext and a near-bijective substitution map; report convergence traces instead of only final text.
- Deliverables: pseudocode for EM and Metropolis-Hastings, runnable helpers in `src/crypto`, and a reproducible setup (commands + package versions).

## How to reproduce

- Run from `crypto1/` so relative paths resolve (notably `resources/...`).
- Ensure the n-gram LM pickle exists (shipped as `resources/en/ngrams_space.pkl`). To regenerate it:
  - Default (Brown corpus; requires `nltk` + Brown data): `python3 -m src.crypto.demo prep-lm --out resources/en/ngrams_space.pkl`
  - Or from a local corpus file: `python3 -m src.crypto.demo prep-lm --corpus path/to/corpus.txt --out resources/en/ngrams_space.pkl`
- Run the demos:
  - EM: `python3 -m src.crypto.demo em --max-iter 200 --restarts 15 --verbose`
  - MH: `python3 -m src.crypto.demo mh --iters-sub 500000 --iters-space 5000 --verbose`
  - All: `python3 -m src.crypto.demo all --max-iter 200 --restarts 15 --iters-sub 500000 --iters-space 5000 --verbose`
- Runtime (typical laptop): EM decryption ~1-2 minutes at 200-300 iterations; MH on short texts ~seconds; pickle generation scales linearly with corpus size.

## Reproducibility setup

- Environment: Python 3.x and `numpy`.
- Project layout: run from `crypto1/` and invoke `python3 -m src.crypto.demo ...` (or use `.venv/bin/python3` if you created a virtualenv in this directory).
- Data: the demos load the language model from `resources/en/ngrams_space.pkl` by default; regenerate it with `prep-lm` as shown above if needed.
- Optional: to regenerate the pickle from the Brown corpus, install `nltk` and ensure the Brown corpus data is available locally.
- Determinism: MH uses fixed internal seeds; EM is stochastic due to random initialization, so use multiple restarts and report representative traces.

## Methods

### Breaking the cipher with Expectation-Maximization (EM / Baum-Welch)

We cast the substitution cipher as a Hidden Markov Model (HMM): hidden plaintext states emit the observed ciphertext symbols, and plaintext-to-plaintext transitions follow an English bigram language model. EM (Baum-Welch) [4] [5] then adjusts the HMM parameters to increase the likelihood of the observed ciphertext. In our setup, we typically learn the transition model ( $A$  and  $\pi$ ) from an English corpus and, for each ciphertext, re-estimate the emission model ( $B$ ) while keeping  $A, \pi$  fixed.

Implementation details: scaled forward/backward passes prevent numerical underflow. We also avoid exact zeros in  $A$  (e.g., via smoothing) so transitions are never impossible and EM updates do not get stuck.

#### Baum-Welch: the algorithm

```

1 # scaled forward/backward to avoid underflow
2 initialize A, B, PI
3 repeat until convergence:
4   alpha, scale = forward(A, B, PI, V)
5   beta = backward(A, B, V, scale)
6   gamma_state = expected_state_visits(alpha, beta)
7   gamma_trans = expected_state_transitions(alpha, beta, A, B, V)
8   if recompute_A: A = normalize_rows(gamma_trans)
9   if recompute_B: B = normalize_rows(state_emissions(gamma_state, V))
10  if recompute_PI: PI = gamma_state[0]
11  B = constrain_space(B) # enforce space->space

```

This highlights the scaled passes, the non-zero constraints on A, and the bijective tweak for the space character.

### Baum-Welch: forward-backward and EM updates

Ciphertext symbols are modeled as emissions from a hidden Markov model (HMM). Let  $H = (p_1, \dots, p_N)$  be the hidden state sequence (e.g., plaintext symbols) and  $V = (c_1, \dots, c_N)$  the observed sequence (ciphertext symbols). Parameters are  $\Theta = (\pi, A, B)$ : \*  $\pi_i = P(p_1 = i)$  initial state probability. \*  $A_{ij} = P(p_{t+1} = j | p_t = i)$  transition probability. \*  $B_i(c) = P(c_t = c | p_t = i)$  emission probability.

**Joint probability** (one particular hidden path):

$$P(V, H | \Theta) = \pi_{p_1} B_{p_1}(c_1) \prod_{t=1}^{N-1} A_{p_t p_{t+1}} B_{p_{t+1}}(c_{t+1})$$

The likelihood  $P(V | \Theta) = \sum_H P(V, H | \Theta)$  sums over all possible hidden paths, which is too expensive to do directly. Forward-backward computes it efficiently.

**Forward pass** (prefix probability). Define:

$$\alpha_t(i) = P(c_{1:t}, p_t = i | \Theta)$$

Initialize and recurse left to right:

$$\begin{aligned} \alpha_1(i) &= \pi_i B_i(c_1) \\ \alpha_t(j) &= B_j(c_t) \sum_i \alpha_{t-1}(i) A_{ij} \quad (t = 2, \dots, N) \end{aligned}$$

**Backward pass** (suffix probability). Define:

$$\beta_t(i) = P(c_{t+1:N} | p_t = i, \Theta)$$

Start at the end and recurse right-to-left:

$$\begin{aligned} \beta_N(i) &= 1 \\ \beta_t(i) &= \sum_j A_{ij} B_j(c_{t+1}) \beta_{t+1}(j) \quad (t = N-1, \dots, 1) \end{aligned}$$

**Likelihood** (total probability of the observed ciphertext):

$$P(V | \Theta) = \sum_i \alpha_N(i)$$

**Posterior marginals (E-step).** Once we have  $\alpha$  and  $\beta$ , we can ask: “how likely was the model in state  $i$  at time  $t$ ?” and “how likely was the transition  $i \rightarrow j$  at time  $t$ ?”

$$\gamma_t(i) = P(p_t = i | V, \Theta) = \frac{\alpha_t(i)\beta_t(i)}{P(V|\Theta)},$$

$$\xi_t(i, j) = P(p_t = i, p_{t+1} = j | V, \Theta) = \frac{\alpha_t(i)A_{ij}B_j(c_{t+1})\beta_{t+1}(j)}{P(V|\Theta)} \quad (t = 1, \dots, N-1)$$

**Parameter updates (M-step).** Update parameters by normalizing expected counts:

$$\pi_i \leftarrow \gamma_1(i)$$

$$A_{ij} \leftarrow \frac{\sum_{t=1}^{N-1} \xi_t(i, j)}{\sum_{t=1}^{N-1} \gamma_t(i)}$$

$$B_i(c) \leftarrow \frac{\sum_{t=1}^N \mathbf{1}[c_t = c] \gamma_t(i)}{\sum_{t=1}^N \gamma_t(i)}$$

Iterate E/M steps until convergence (or for a fixed number of iterations).

Key implementation details (see [6]): - Transition matrix  $A$  avoids exact zeros; strict positivity implies the Markov chain is irreducible and aperiodic. - Forward/backward passes are scaled each step to prevent underflow on long ciphertexts. - First-order LM keeps sparsity manageable; higher orders need smoothing and more data.

### EM implementation: a concrete example

So far we have summarized the algorithm and its Forward-Backward passes. Let’s use it to automatically decrypt some text.

### EM usage sketch (see Appendix A for full code)

1. Build bigram LM  $(A, B_{id}, \pi)$ [3]: `A, B_id, PI = build_bigram_model(corpus_text, max_iter=50, limit=200k)`.
2. Initialize emission matrix  $B$  randomly, apply `constrain_space_mapping(B)`.
3. Run `BaumWelch.process(A, B, PI, V, max_iter=200, recompute_A=False, recompute_PI=False, ext_fun=constrain_space_mapping)`.
4. Decode with `decode_with_emissions(B_hat, ciphertext)` or feed emissions to Viterbi[7] for MAP decoding.

Log-likelihoods from a representative run are in the convergence table below.

### EM convergence snapshot

Iteration	$\log P(V \Theta)$	Note
0	-374.15	random init
2	-297.81	structure starts to appear
16	-242.69	readable plaintext emerges
200	-237.72	plateau; restart to escape local optima

Scaling prevents underflow, and strictly positive  $A$  keeps the chain irreducible and aperiodic. Run multiple random restarts (different  $B$  seeds) and keep the run with the best final log-likelihood.

## EM decoded sample

- Iter 0: BZN BQWNXLN JQBZ ... (random initialization, gibberish)
- Iter 16: THE TOXUPLE WITH HAZIND AN OQEN MING ... (largely readable, a few bad mappings)
- Iter 200: THE TROUPLE WITH HAZING AN OQEN MIND ... (plateau; minor letter swaps remain)

These snapshots show how likelihood climbs and the mapping improves before stalling. Rerun with new seeds to escape plateaus.

## EM usage example: computing statistics for the English bigrams LM

Use `build_bigram_model` to estimate the bigram transition matrix  $A$  and initial state distribution  $\pi$  from a cleaned English corpus (e.g., the Brown corpus [3]). The helper trims the text to a configurable limit to keep runtimes reasonable.

## EM usage: decrypt a text

With  $A$  and  $\pi$  fixed, initialize an emission matrix  $B$ , apply `constrain_space_mapping`, and run `BaumWelch.process` with `recompute_A=False` and `recompute_PI=False`. Use `decode_with_emissions` for a quick readout or run Viterbi[7] with the learned emissions to obtain a MAP state sequence.

## Review of the Baum-Welch approach

Baum-Welch can often recover readable plaintext quickly on sufficiently long ciphertexts, although runtime and quality depend on text length and initialization. In our setup we train  $A, \pi$  from an English corpus (bigram LM), then learn an emission matrix  $B$  per ciphertext while keeping  $A, \pi$  fixed.

Notes and limitations:

- the entropy of English implies a (rough) *unicity distance* - the minimum length needed for unique recovery in an idealized setting (Shannon, 1949) [8]. A common estimate uses per-character redundancy  $D \approx 3.2$  bits:

$$U = \frac{H(K)}{D} = \frac{\log_2(26!)}{3.2} \approx 28$$

In practice, Baum-Welch typically needs substantially longer ciphertexts for stable convergence.

- Using only a first-order (bigram) language model limits context to local dependencies.
- Higher-order models (e.g., trigrams) can improve accuracy but increase computation and memory.
- A substitution cipher is a bijection, but plain Baum-Welch does not enforce a one-to-one constraint on  $B$ ; enforcing bijectivity requires constrained optimization or post-processing.
- We can bias initialization toward expected mappings (e.g., space→space, or other priors in  $B$ ). Injecting higher-level knowledge (common words, names, greetings) is possible via constraints/priors or hybrid search, but it is not built into standard Baum-Welch.

## Breaking the cipher with Metropolis-Hastings

Let's try a different angle.

Is there a way to enforce the bijective property of the cipher while solving the problem?

We briefly describe the Monte Carlo technique and the Metropolis-Hastings algorithm before using it with a new model.

## Monte Carlo

If a quantity is intractable or hard to compute directly, we can approximate it via random sampling from a related distribution, avoiding the need for closed-form mathematics.

A typical example is estimating the value of  $\pi$ . Consider a circle of radius 1 inscribed in the square  $[-1, 1] \times [-1, 1]$ . If we repeatedly generate pseudo-random points  $(x, y)$  uniformly in the square and count the fraction that fall inside the circle ( $x^2 + y^2 \leq 1$ ), that fraction estimates the area ratio:

$$\text{square area} = 2r \cdot 2r = 4$$

$$\text{circle area} = \pi \cdot r^2 = \pi$$

$$\frac{\text{circle area}}{\text{square area}} = \pi/4$$

Therefore  $\pi \approx 4 \cdot (\text{fraction of points inside the circle})$ . By the law of large numbers, this estimate converges as the number of samples grows.

```

1 limit=1_000_000 # the more the better
2 print(sum(random.uniform(-1,1)**2 + random.uniform(-1,1)**2<=1 for i in range(limit)
) *4.0/limit)
3 3.14188

```

Under “Monte Carlo” methods fall different techniques, which stochastically tweak the parameters of a model to explore unknown values and evaluate its outputs, such that the model can be improved. In our setting, this “tweaking” will take the form of a random-walk proposal over candidate cipher keys. This walk is mostly done by generating random samples and then accepting or discarding them according to a criterion, or by generating samples from some weighted distribution (importance sampling). For example, in the case of a substitution cipher, consider a random walk where at each step two letters of the cipher key are swapped and the resulting plaintext is evaluated for plausibility under a language model.

### Metropolis-Hastings, brief introduction

Metropolis-Hastings is a Markov chain Monte Carlo (MCMC) algorithm that works well for the substitution cipher case.

- consider a target distribution  $P(x)$  over solutions  $x$  (here: cipher keys) which we cannot evaluate exactly, but for which we can compute a score from the observations;
- we have a function  $\ell(x) = \log \tilde{P}(x) + C$  which is equal to  $\log P(x)$  up to an additive constant (log unnormalized score);
- we make a random walk by randomly making a change to the parameters (e.g., swapping two letters in the key) and then evaluating the result;
- the new value is accepted or rejected using an acceptance function that compares  $\tilde{P}(x')$  with  $\tilde{P}(x)$  (and, in general, also accounts for the proposal distribution);
- the proposal distribution  $q(x'|x)$  specifies how we generate the candidate  $x'$  from the current state  $x$  (for swap moves, it is typically symmetric).

In the longer term (and under some basic conditions), this walk is guaranteed to sample from the target distribution  $P(x)$  (more precisely:  $P(x)$  is the stationary distribution of the Markov chain).

### More about convergence of Markov models

A discrete-time Markov chain with transition probabilities  $P(y|x)$  has a stationary distribution  $\Pi^*$  if (see [9]):

$$\Pi^*(y) = \sum_x \Pi^*(x) P(y|x)$$

In words: if the chain starts distributed as  $\Pi^*$ , then after one step it is still distributed as  $\Pi^*$ .

If the chain is irreducible and aperiodic (and the state space is finite), then a stationary distribution exists, is unique, and the chain converges to it. In particular, if we write the distribution at time  $t$  as a row vector  $\Pi_t$ , and the transition matrix as  $A$  with entries  $A_{xy} = P(y|x)$ , then

$$\Pi_{t+1} = \Pi_t A \quad \Pi_t = \Pi_0 A^t$$

For large  $t$ ,

$$\Pi_0 A^t \xrightarrow[t \rightarrow \infty]{} \Pi^*$$

and the stationary distribution satisfies the fixed-point equation

$$\Pi^* = \Pi^* A$$

Note that convergence does not mean the chain stops moving; it means the distribution of the state stabilizes.

A strong sufficient condition is strict positivity: if  $A_{xy} > 0$  for all  $x, y$ , then the chain is automatically irreducible and aperiodic, hence it has a unique stationary distribution  $\Pi^*$ .

In Metropolis-Hastings, the goal is to construct a Markov chain whose stationary distribution is the desired target distribution over solution (e.g., cipher keys). In other words, the target distribution is the  $\Pi^*$  we want the chain to converge to.

A common way to guarantee that a chosen  $\Pi^*$  is stationary is detailed balance (time reversibility):

$$\Pi^*(x)P(y|x) = \Pi^*(y)P(x|y) \quad (3)$$

Summing both sides over  $x$  gives stationarity:

$$\sum_x \Pi^*(x)P(y|x) = \sum_x \Pi^*(y)P(x|y) = \Pi^*(y) \sum_x P(x|y) = \Pi^*(y)$$

since  $\sum_x P(x|y) = 1$ .

### Some clues behind the inner workings of Metropolis-Hastings for the substitution cipher

So far, we modeled ciphertext with an HMM and learned an emission model with Baum-Welch. However, the substitution cipher key is bijective, and Baum-Welch does not naturally enforce a one-to-one constraint on the emission matrix. Also, using higher-order language models inside an HMM quickly becomes expensive. Let's rethink the model in a way that enforces the substitution mapping directly, and makes it easy to plug in an  $n$ -gram LM.

The technique discussed here is inspired by “The Markov Chain Monte Carlo Revolution” (Diaconis, 2009) [10] which gives deeper insights into MCMC and uses Metropolis-Hastings to (among other things) decode an encrypted message.

Consider the substitution mapping  $S$  (a bijection) and its inverse  $S^{-1}$ :  $S(P) = C$  and  $S^{-1}(C) = P$ .

In Metropolis-Hastings, the state of the Markov chain is the key  $S$  itself. We move from one key to another by proposing small changes (e.g., swapping two letters in the key) and accepting/rejecting them. The chain is constructed so that its stationary distribution  $\Pi(S)$  assigns higher probability to keys that produce more plausible English plaintext.

The stationary distribution  $\Pi(S)$  does not need to be known in normalized form: it is enough to define an unnormalized score  $\tilde{\Pi}(S) \propto \Pi(S)$ , because MH only uses ratios of scores.

Using an  $n$ -gram language model (LM), we can score a candidate key  $S$  by decoding the ciphertext sequence  $(c_1, \dots, c_N)$  into candidate plaintext symbols

$$p_i = S^{-1}(c_i)$$

and then computing the LM score of the resulting plaintext under the  $n$ -gram factorization:

$$Pl(S) = \sum_{i=n}^N \log P_{LM}(p_i | p_{i-n+1:i-1})$$

In implementation we use an unnormalized weighted sum of n-gram log-weights; MH only needs score differences.

### Metropolis-Hastings (symmetric proposal): the algorithm

Metropolis et al. (1953) [11] is the special case of Metropolis-Hastings where the proposal is symmetric; it can be motivated via detailed balance.

It needs (i) a symmetric proposal distribution to generate candidate moves and (ii) a plausibility score (based on the language model) proportional to the target distribution. It works as follows:

0. Start with:

- a current state  $x$  (e.g., a key; can be random),
- a symmetric proposal  $g(y|x)$  such that  $g(y|x) = g(x|y)$ ,
- a log-score  $Pl(x)$  (plausibility), defined so that  $\exp(Pl(x)) \propto \Pi(x)$  (equivalently,  $Pl(x) = \log \tilde{\Pi}(x) + C$  for some constant  $C$ ).

1. From the current state  $x$ , draw a candidate  $y \sim g(y|x)$ .

2. Compute the acceptance probability

$$\alpha(x, y) = \min(1, \exp(Pl(y) - Pl(x)))$$

3. Accept/Reject:

- draw  $u \sim \text{Uniform}(0, 1)$ ;
- if  $u \leq \alpha(x, y)$ , accept and set  $x \leftarrow y$ ;
- otherwise, reject and keep  $x$ .

4. Repeat from step 1 (for a fixed number of iterations, or after burn-in collect samples).

The acceptance rule comes from the detailed-balance condition (3). Let  $\Pi^*$  be the desired stationary distribution and let  $T(y|x)$  be the actual transition kernel of the Metropolis chain. Condition (3) requires:

$$\Pi^*(x)T(y|x) = \Pi^*(y)T(x|y)$$

For  $y \neq x$ , the Metropolis transition can be written as:

$$T(y|x) = g(y|x)\alpha(x, y)$$

where  $g(y|x)$  is the proposal distribution and  $\alpha(x, y)$  is the acceptance probability. Plugging this into (3) gives:

$$\frac{g(y|x)\alpha(x, y)}{g(x|y)\alpha(y, x)} = \frac{\Pi^*(y)}{\Pi^*(x)}$$

If the proposal is symmetric,  $g(y|x) = g(x|y)$ , this simplifies to:

$$\frac{\alpha(x, y)}{\alpha(y, x)} = \frac{\Pi^*(y)}{\Pi^*(x)}$$

Since  $Pl(x)$  is a log-score (i.e.,  $Pl(x) = \log \tilde{\Pi}(x)$  up to an additive constant), the condition is satisfied by choosing:

$$\alpha(x, y) = \min(1, \exp(Pl(y) - Pl(x))) = \min\left(1, \frac{\Pi^*(y)}{\Pi^*(x)}\right)$$

```

1 ct = normalize(ciphertext)
2 S = random_key(alphabet)           # current bijection (state)
3 pt = decode(ct, S)                 # candidate plaintext under S
4 score = plausibility(pt, lm)       # log-score (recommended)
5
6 best_S, best_score = S, score

```

```

7
8 for _ in range(N):
9     a, b = pick_two_symbols(alphabet)
10    S2 = swap_in_key(S, a, b)           # symmetric proposal, still bijective
11    pt2 = decode(ct, S2)
12    score2 = plausibility(pt2, lm)Δ
13
14    = score2 - score                 # log acceptance ratio
15    if log(rand_uniform()) <= Δ:      # accept with prob min(1, expΔ())
16        S, pt, score = S2, pt2, score2
17
18    if score > best_score:
19        best_S, best_score = S, score

```

The proposal enforces a bijection by swapping two symbols; the plausibility function can mix multiple n-gram lengths.

Acceptance intuition: early on, many swaps are accepted; later, acceptance drops as the chain reaches higher-scoring regions. Restart from a fresh seed if trapped in low-quality states.

Plausibility should return a log score (sum of log n-gram probs) so  $\Delta$  is stable and  $\exp(\Delta)$  is well-defined.

### Metropolis-Hastings implementation: a concrete example

With this algorithm, we can exploit an LM with multiple  $n$ -grams. This is straightforward because the plausibility function only needs to rank candidates; a fully normalized probability is not required. Intuitively, a candidate plaintext should score higher when it contains  $n$ -grams that are typical for English. A simple choice is the sum of weighted  $n$ -gram frequencies, where the weight increases with  $n$  because longer contexts carry more information.

Since “frequency” is a normalized quantity, the normalization step can be omitted as long as the normalization factor (e.g., the total number of extracted  $n$ -grams) is constant across candidates; it cancels out when comparing two states in the acceptance ratio. In practice, the weighted count is treated as a log-plausibility score, so the acceptance step uses  $\exp(\Delta)$  with  $\Delta = Pl(y) - Pl(x)$ .

Furthermore, at each step an existing mapping is swapped; this maintains the bijection between the plaintext and ciphertext alphabets.

The  $n$ -gram model is sparse for larger  $n$ , so it would be inefficient to store it as a dense array. In this Python example we use a dictionary to store only the observed  $n$ -grams.

Finally, the acceptance rule allows occasional moves to lower-scoring states, which helps the chain jump between local optima. There is no single point where the algorithm is guaranteed to be “done”, so for this project we set a maximum number of iterations; another option is to stop if there is no improvement in the best score for a fixed window of iterations.

Here is my full implementation of the Metropolis algorithm and a concrete example:

### MH usage sketch (see Appendix B for full code)

- Ensure the n-gram LM pickle exists (e.g., `resources/en/ngrams_space.pkl`, generated with `python3 -m src.crypto.demo prep-lm`).
- Substitution cipher: `metropolis_substitution(ciphertext, iterations=500_000, lm_pickle=Path("resources/ngrams_space.pkl"))`.
- Patristocrat spacing: `metropolis_patristocrat(text_without_spaces, iterations=5_000, lm_pickle=Path("resources/en/ngrams_space.pkl"))`.

- CLI equivalent: `python3 -m src.crypto.demo mh --iters-sub 500000 --iters-space 5000 --verbose`.

### MH trace highlights

- In this trace, the first accepted swaps quickly lock in a few high-signal bigrams (e.g., TH, HE), and readability starts improving early.
- The acceptance rate here tracks the plausibility spread: once the score separation between candidates grows, most proposed swaps are rejected and progress comes from occasional improving moves.
- In the spacing stage of this same trace, patristocrat spacing recovers as common trigrams reappear.

### MH trace sample

- It#1: K CXJO XI HBD BKIO XS ENJHB HEN XI HBD CLSB (initial guess)
- It#15: M CFJS FI TBR BMIS FO ENJTB TEN FI TBR CLOB (bigrams emerging)
- It#219: IF YOU LOOK FOR PERFECTION YOU LL NEVER BE CONTENT (spacing recovered)

### Patristocrats

To make a substitution cipher harder to break, spaces can be removed from the ciphertext. The American Cryptogram Association calls this kind of cryptogram a “Patristocrat”. This increases difficulty, at the cost of added ambiguity when reconstructing word boundaries.

To solve a Patristocrat, it is convenient to score candidate plaintext without spaces using an LM trained on a corpus with spaces removed (alphabet A-Z only). The same algorithms as before can be used, but longer ciphertext is typically needed to provide enough statistical signal. Higher-order  $n$ -gram language models (as used in the Metropolis-Hastings scoring function) tend to work much better here than first-order/bigram models, as expected.

After the letters have been decoded, spaces can be reintroduced using an LM trained on a corpus with spaces. There are several ways to do this (e.g., dynamic programming), but with only minor changes to the previous implementation we use Metropolis-Hastings: proposals insert/remove a space at random positions and moves are accepted according to the spaced-LM plausibility score.

Practical usage: train the A-Z LM (spaces removed) and run `metropolis_patristocrat`; then train the spaced LM and run the spacing sampler.

### Review of the Metropolis approach

With enough iterations and a few random restarts, Metropolis-Hastings can often recover readable plaintext even from relatively short ciphertexts, although success depends on text length and the strength of the language-model score. The algorithm is flexible: it can incorporate higher-order  $n$ -gram language models directly in the plausibility function without changing the sampler itself. In practice, the speed of improvement depends mainly on the proposal function (how new candidates are generated) and the acceptance rule (which controls how often worse moves are tolerated). Variants of Metropolis-Hastings typically modify these two components to trade off exploration, runtime, and solution quality.

### Limitations & future work

- Short ciphertexts remain ambiguous; even with restarts, many distinct keys can yield similarly plausible plaintext. Report variability across seeds/restarts and show representative convergence traces, not only a single final output.
- EM and MH are not optimizing the same objective (HMM likelihood vs LM-based key plausibility). For a fair comparison, use synthetic ciphers with known keys and report accuracy metrics (e.g., key accuracy / character accuracy) with variance across runs and wall-clock runtime.
- Stronger language models (smoothed higher-order  $n$ -grams or word-level models) can improve recovery, but they increase data requirements and can introduce corpus bias if the plaintext domain differs.

- Enforcing bijectivity explicitly (e.g., via an assignment step to project emissions onto a permutation, or by using MH directly on the key) is likely to improve stability over unconstrained EM emissions.
- MH performance depends heavily on the proposal; better proposals (block swaps, tempered/annealed schedules, parallel tempering) can reduce mixing time and improve robustness on harder instances.

## Appendices

Full Python listings referenced in the text. Place these in `src/crypto/` or on your import path when re-running the experiments.

### Appendix A: EM helper code

```

1 """Baum-Welch helpers for bigram HMMs, including training from text and decoding utilities."""
2
3 from __future__ import annotations
4
5 import math
6 from typing import Callable, Tuple
7
8 import numpy as np
9
10 from .utils import (
11     ALPHABET,
12     SPACE_IDX,
13     char_to_index,
14     encode_text,
15     index_to_char,
16     normalize_text,
17     seed_everything,
18 )
19
20 EmissionHook = Callable[[np.ndarray], np.ndarray]
21
22 def _softmax_from_logs(logs: np.ndarray) -> np.ndarray:
23     m = float(np.max(logs))
24     shifted = logs - m
25     w = np.exp(shifted)
26     s = float(np.sum(w))
27     if s == 0.0 or not np.isfinite(s):
28         return np.full_like(w, 1.0 / len(w))
29     return w / s
30
31
32 def bigram_hmm_from_ngram_lm(
33     lm_stats: dict[int, dict[str, float]],
34     *,
35     eps: float = 1e-6,
36     alphabet: str = ALPHABET,
37 ) -> tuple[np.ndarray, np.ndarray]:
38     """
39     Convert an n-gram LM pickle (log-weights) into bigram HMM parameters.
40
41     Returns (A, PI) where PI is the initial state distribution and A is the transition matrix.
42     Missing unigrams/bigrams fall back to `log(eps)` and rows are normalized via softmax.

```

```

43     """
44     symbols = list(alphabet)
45     log_eps = math.log(eps)
46
47     unigrams = lm_stats.get(1, {})
48     log_pi = np.array([float(unigrams.get(ch, log_eps)) for ch in symbols], dtype=float)
49     PI = _softmax_from_logs(log_pi)
50
51     bigrams = lm_stats.get(2, {})
52     A = np.zeros((len(symbols), len(symbols)), dtype=float)
53     for i, a in enumerate(symbols):
54         row_logs = np.array([float(bigrams.get(a + b, log_eps)) for b in symbols],
55                             dtype=float)
56         A[i, :] = _softmax_from_logs(row_logs)
57     return A, PI
58
59 def score_bigram_logprob(text: str, bigram_log_probs: dict[str, float], *, eps: float = 1e-6)
60     -> float:
61         """Score text by summing bigram log-weights from a (possibly sparse) bigram table."""
62         log_eps = math.log(eps)
63         total = 0.0
64         for i in range(len(text) - 1):
65             total += float(bigram_log_probs.get(text[i : i + 2], log_eps))
66         return total
67
68 class BaumWelch:
69     """Minimal Baum-Welch implementation with scaling to avoid underflow."""
70
71     def __init__(self, num_hidden: int = len(ALPHABET)):
72         self.num_hidden = num_hidden
73         self.num_visible = num_hidden
74         self.T = 0
75
76     @staticmethod
77     def normalize(matrix: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
78         totals = np.sum(matrix, axis=min(1, len(matrix.shape) - 1), keepdims=True)
79         totals[totals == 0] = 1.0
80         return matrix / totals, totals
81
82     def alpha(self, A: np.ndarray, B: np.ndarray, PI: np.ndarray, V: np.ndarray):
83         alphas = np.zeros((self.T, self.num_hidden))
84         scale_facts = np.zeros(self.T)
85         alphas[0], total = self.normalize(PI * B[:, V[0]])
86         scale_facts[0] = total.item()
87         for t in range(1, self.T):
88             alphas[t], total = self.normalize(
89                 B[:, V[t]] * np.dot(alphas[t - 1], A)
90             )
91             scale_facts[t] = total.item()
92         return alphas, scale_facts
93
94     def beta(self, A: np.ndarray, B: np.ndarray, V: np.ndarray, scale_facts: np.ndarray):

```

```

95     betas = np.zeros((self.T, self.num_hidden))
96     betas[self.T - 1] = 1.0
97     for t in range(self.T - 2, -1, -1):
98         betas[t] = np.sum(betas[t + 1] * B[:, V[t + 1]] * A, axis=1) / scale_facts[t + 1]
99     return betas
100
101    def process(
102        self,
103        A: np.ndarray,
104        B: np.ndarray,
105        PI: np.ndarray,
106        V: np.ndarray,
107        max_iter: int,
108        recompute_A: bool = True,
109        recompute_B: bool = True,
110        recompute_PI: bool = True,
111        ext_fun: EmissionHook | None = None,
112        verbose_fun: Callable[[np.ndarray, np.ndarray, int, float], None] | None = None,
113        stop_on_decrease: bool = True,
114        min_delta: float = 0.0,
115        patience: int | None = None,
116    ):
117        self.T = V.size
118        log_scale_facts = None
119        stale = 0
120        for it in range(max_iter):
121            alphas, scale_facts = self.alpha(A, B, PI, V)
122            new_log_scale = float(np.sum(np.log(scale_facts)))
123            if log_scale_facts is not None:
124                delta = new_log_scale - log_scale_facts
125                if stop_on_decrease and delta < -min_delta:
126                    if verbose_fun:
127                        print("Algorithm has converged (decrease)")
128                    break
129                if abs(delta) <= min_delta:
130                    stale += 1
131                    if patience is not None and stale >= patience:
132                        if verbose_fun:
133                            print("Algorithm has converged (plateau)")
134                        break
135                    else:
136                        stale = 0
137
138            log_scale_facts = new_log_scale
139            betas = self.beta(A, B, V, scale_facts)
140            gammas = np.zeros((self.T, self.num_hidden, self.num_hidden))
141            for t in range(self.T - 1):
142                gammas[t] = (A.T * alphas[t, :]).T * (betas[t + 1] * B[:, V[t + 1]])
143                denom = float(np.sum(gammas[t]))
144                if denom != 0.0:
145                    gammas[t] /= denom
146                gammas_i = np.sum(gammas, axis=2)
147                gammas_i[self.T - 1] = alphas[self.T - 1, :]
148

```

```

149         if recompute_A:
150             gammas_t = np.sum(gammas, axis=0)
151             totals = np.sum(gammas_t, axis=1)[:, None]
152             totals[totals == 0] = 1.0
153             A = (gammas_t.T / totals).T
154
155         if recompute_B:
156             B_counts = np.zeros_like(B)
157             for v in range(self.num_visible):
158                 mask = V == v
159                 if np.any(mask):
160                     B_counts[:, v] = np.sum(gammas_i[mask], axis=0)
161                     totals = np.sum(gammas_i, axis=0)
162                     totals[totals == 0] = 1.0
163                     B = (B_counts.T / totals).T
164
165         if recompute_PI:
166             PI = gammas_i[0]
167
168         if ext_fun is not None:
169             B = ext_fun(B)
170
171         if verbose_fun:
172             verbose_fun(V, B, it, log_scale_facts)
173
174     return A, B, PI, log_scale_facts
175
176
177 def constrain_space_mapping(B: np.ndarray) -> np.ndarray:
178     """
179     Force a 1-1 mapping for space:space in the emission matrix.
180
181     Ensures space can only map to space, keeping the substitution mapping bijective for that
182     symbol.
183     """
184     B = B.copy()
185     B[SPACE_IDX, :] = 0.0
186     B[:, SPACE_IDX] = 0.0
187     B[SPACE_IDX, SPACE_IDX] = 1.0
188     row_sums = B.sum(axis=1, keepdims=True)
189     row_sums[row_sums == 0] = 1.0
190     return B / row_sums
191
192 def build_bigram_model(
193     corpus_text: str,
194     *,
195     max_iter: int = 150,
196     limit: int = 1_000_000,
197     seed: int = 13,
198 ):
199     """
200     Train a bigram language model from plain text using Baum-Welch.
201

```

```

202     Returns the transition matrix A, an identity emission matrix B, and the
203     initial state distribution PI. The text is upper-cased, cleaned, and
204     truncated to `limit` characters to keep runtimes reasonable.
205     """
206     seed_everything(seed)
207     normalized = normalize_text(corpus_text)[:limit]
208     V = np.array(encode_text(normalized))
209
210     num_states = len(ALPHABET)
211     PI = np.random.rand(num_states)
212     A = np.random.rand(num_states, num_states)
213     B = np.identity(num_states)
214
215     baum = BaumWelch(num_states)
216     A, _, _ = baum.process(A, B, PI, V, max_iter=max_iter, recompute_B=False)
217     return A, B, PI
218
219
220 def decode_with_emissions(B: np.ndarray, observations: str) -> str:
221     """
222     Simple decoder: pick the highest emission probability per symbol.
223
224     For maximum a posteriori decoding, use the Viterbi algorithm instead.
225     """
226     V = [char_to_index(ch) for ch in observations]
227     decoded_indices = [int(np.argmax(B[:, v])) for v in V]
228     return "".join(index_to_char(idx) for idx in decoded_indices)

```

## Appendix B: Metropolis-Hastings helper code

```

1 """Metropolis(-Hastings) helpers for substitution and spacing recovery."""
2
3 from __future__ import annotations
4
5 import math
6 import random
7 import string
8 from collections import Counter
9 from pathlib import Path
10 from typing import Dict
11
12 from .utils import load_pickled_language_model
13
14
15 def count_ngrams(text: str, min_dgram_len: int = 2, max_dgram_len: int = 9) -> Dict[int,
16     Dict[str, int]]:
17     """Raw n-gram counts over the requested length range."""
18     return {
19         dlen: Counter(text[idx : idx + dlen] for idx in range(len(text) - dlen + 1))
20         for dlen in range(min_dgram_len, max_dgram_len + 1)
21     }
22
23 def build_language_model(keep_spaces: bool = True, pickle_path: Path | None = None) ->
24     Dict[int, Dict[str, float]]:
25     """
26     Create the multi-gram language model.

```

```

24
25     Requires `pickle_path` pointing to an n-gram LM pickle.
26     The pickle is expected to contain a mapping `dict[int, dict[str, float]]` of n-gram
27     length to log-weights.
28     """
29     if pickle_path and pickle_path.exists():
30         return load_pickled_language_model(pickle_path)
31     raise RuntimeError("Language model pickle not found.")
32
33 def plausibility(
34     text: str,
35     lm: Dict[int, Dict[str, float]],
36     *,
37     missing: float = math.log(1e-6),
38 ) -> float:
39     """
40     Log-plausibility under the LM: sum of count * log-weight, weighted by n-gram length.
41     """
42     min_n, max_n = min(lm), max(lm)
43     cnt = count_ngrams(text, min_n, max_n)
44     total = 0.0
45     for ngram_len, counts in cnt.items():
46         lm_row = lm.get(ngram_len, {})
47         for k, freq in counts.items():
48             weight = lm_row.get(k, missing)
49             total += freq * weight * ngram_len
50     return total
51
52
53 def metropolis_substitution(
54     ciphertext: str,
55     iterations: int = int(1e6),
56     seed: int = 13,
57     lm_pickle: Path | None = None,
58     verbose: bool = False,
59 ):
60     """
61     Recover simple substitution using Metropolis on swap proposals.
62     """
63     random.seed(seed)
64     mm = build_language_model(keep_spaces=True, pickle_path=lm_pickle)
65     best = pl_f = plausibility(ciphertext, mm)
66     best_text = ciphertext
67     fixed = set(ciphertext) - set(string.ascii_uppercase)
68     used_symbols = list(set(ciphertext) - fixed)
69     all_symbols = set(string.ascii_uppercase)
70     smallprob = math.log(1e-3)
71     for i in range(iterations):
72         c1 = random.choice(used_symbols)
73         c2 = random.choice([c for c in all_symbols if c != c1])
74         candidate = ciphertext.translate(str.maketrans({c1: c2, c2: c1}))
75         pl_f_star = plausibility(candidate, mm)
76         mu = math.log(random.uniform(0, 1))

```

```

77     if mu < max(min(pl_f_star - pl_f, 0), smallprob):
78         if pl_f_star > best:
79             best = pl_f_star
80             best_text = candidate
81             if verbose:
82                 print("It#", i, candidate)
83             ciphertext = candidate
84             pl_f = pl_f_star
85             used_symbols = list(set(ciphertext) - fixed)
86     return best_text
87
88
89 def metropolis_patristocrat(
90     ciphertext: str,
91     iterations: int = 5000,
92     seed: int = 21,
93     lm_pickle: Path | None = None,
94     verbose: bool = False,
95 ):
96     """
97     Restore spaces in a patristocrat ciphertext.
98     """
99     random.seed(seed)
100    mm = build_language_model(keep_spaces=True, pickle_path=lm_pickle)
101    best = pl_f = plausibility(ciphertext, mm) / len(ciphertext)
102    ret = ciphertext
103    for i in range(iterations):
104        c = random.choice(range(1, len(ciphertext) - 2))
105        candidate = list(ciphertext)
106        if ciphertext[c] == " ":
107            candidate = candidate[:c] + candidate[c + 1 :]
108            pl_f_star = plausibility("".join(candidate), mm) / len(candidate)
109        elif ciphertext[c - 1] != " " and ciphertext[c + 1] != " ":
110            candidate[c:c] = " "
111            pl_f_star = plausibility("".join(candidate), mm) / len(candidate)
112        else:
113            pl_f_star = pl_f
114        mu = math.log(random.uniform(0, 1))
115        if mu < min(pl_f_star - pl_f, 0):
116            ciphertext = "".join(candidate)
117            pl_f = pl_f_star
118            if best < pl_f_star:
119                best = pl_f_star
120                ret = ciphertext
121                if verbose:
122                    print(ciphertext)
123    return ret, best

```

## Appendix C: Utility functions

```

1 import math
2 import pickle
3 import random
4 import re
5 from collections import Counter

```

```

6 from pathlib import Path
7 from typing import Iterable
8
9 ALPHABET = "ABCDEFGHIJKLMNOPQRSTUVWXYZ "
10 SPACE_IDX = ALPHABET.index(" ")
11
12 def load_pickled_language_model(path: Path) -> dict[int, dict[str, float]]:
13     """Load a precomputed n-gram language model from a pickle."""
14     data = pickle.loads(path.read_bytes())
15     return {int(k): v for k, v in data.items()}
16
17
18 def seed_everything(seed: int) -> None:
19     """Seed Python's random generator and numpy if available."""
20     random.seed(seed)
21     try:
22         import numpy as np
23
24         np.random.seed(seed)
25     except Exception:
26         # numpy is optional
27         pass
28
29
30 def normalize_text(text: str) -> str:
31     """
32     Keep A-Z and spaces, collapse whitespace, and uppercase.
33
34     Ensures corpus and cipher text share the same character mapping for the LM.
35     """
36     cleaned = re.sub("[^A-Z ]", " ", text.upper())
37     return re.sub("\s+", " ", cleaned).strip()
38
39
40 def char_to_index(ch: str) -> int:
41     """Map a character in ALPHABET to its ordinal index."""
42     return ALPHABET.index(ch)
43
44
45 def index_to_char(idx: int) -> str:
46     """Map an ordinal index back to a character in ALPHABET."""
47     return ALPHABET[idx]
48
49
50 def encode_text(text: Iterable[str]) -> list[int]:
51     """Convert a string of characters into indices."""
52     return [char_to_index(ch) for ch in text]
53
54
55 def decode_indices(indices: Iterable[int]) -> str:
56     """Convert indices back to a string of characters."""
57     return "".join(index_to_char(i) for i in indices)
58
59

```

```

60 def write_ngram_lm_pickle(
61     corpus_text: str,
62     out_path: str | Path,
63     *,
64     min_n: int = 1,
65     max_n: int = 7,
66     keep_spaces: bool = True,
67     limit_chars: int | None = None,
68 ) -> Path:
69     """
70     Build an n-gram language model from text and write it as a pickle.
71
72     The corpus is normalized to the same A-Z plus space alphabet as the cipher text.
73     The pickle contains `dict[int, dict[str, float]]`, mapping n-gram length to
74     n-gram log-probabilities (natural log) estimated by relative frequency (no smoothing).
75     Unseen n-grams are not stored and are handled by the scoring function via a fallback.
76     """
77     if min_n < 1 or max_n < min_n:
78         raise ValueError("Expected 1 <= min_n <= max_n")
79
80     normalized = normalize_text(corpus_text)
81     if not keep_spaces:
82         normalized = normalized.replace(" ", "")
83     if limit_chars is not None:
84         if limit_chars < 0:
85             raise ValueError("Expected limit_chars >= 0")
86         if limit_chars:
87             normalized = normalized[:limit_chars]
88
89     lm: dict[int, dict[str, float]] = {}
90     for n in range(min_n, max_n + 1):
91         counts = Counter(normalized[idx : idx + n] for idx in range(len(normalized) - n + 1))
92         total = sum(counts.values())
93         if total == 0:
94             lm[n] = {}
95             continue
96         lm[n] = {ng: math.log(cnt / total) for ng, cnt in counts.items()}
97
98     out_path = Path(out_path)
99     out_path.parent.mkdir(parents=True, exist_ok=True)
100    out_path.write_bytes(pickle.dumps(lm, protocol=pickle.HIGHEST_PROTOCOL))
101    return out_path
102
103
104 def load_brown_corpus_text(*, keep_spaces: bool = True) -> str:
105     """
106     Load the Brown corpus as a single string.
107
108     Requires `nltk` and the locally-installed Brown corpus data.
109     """
110     try:
111         from nltk.corpus import brown # type: ignore[import-not-found]
112     except ModuleNotFoundError as exc:
113         raise RuntimeError("NLTK is required to load the Brown corpus; install `nltk` or pass

```

```

114         a corpus file.") from exc
115     try:
116         sentences = brown.sents()
117     except LookupError as exc:
118         raise RuntimeError("NLTK Brown corpus data not found; install it or pass a corpus
119                         file.") from exc
120     sep = " " if keep_spaces else ""
121     return sep.join(" ".join(sent) for sent in sentences)

```

## Appendix D: Package exports

```

1 """Python package exporting the EM and MH helpers."""
2
3 from .em import (
4     BaumWelch,
5     build_bigram_model,
6     constrain_space_mapping,
7     decode_with_emissions,
8 )
9 from .mh import (
10    build_language_model,
11    count_ngrams,
12    metropolis_patristocrat,
13    metropolis_substitution,
14 )
15 from .utils import ALPHABET, SPACE_IDX, char_to_index, decode_indices, encode_text,
16                  index_to_char, normalize_text, seed_everything
17
18 __all__ = [
19     "ALPHABET",
20     "SPACE_IDX",
21     "char_to_index",
22     "index_to_char",
23     "encode_text",
24     "decode_indices",
25     "normalize_text",
26     "seed_everything",
27     "BaumWelch",
28     "build_bigram_model",
29     "constrain_space_mapping",
30     "decode_with_emissions",
31     "build_language_model",
32     "count_ngrams",
33     "metropolis_substitution",
34     "metropolis_patristocrat",
35 ]

```

## Appendix E: Demo entrypoint

```

1 """Command-line entrypoint for the Baum-Welch and Metropolis demos."""
2
3 from __future__ import annotations
4

```

```

5 import argparse
6 from pathlib import Path
7
8 import numpy as np
9
10 from .em import (
11     BaumWelch,
12     bigram_hmm_from_ngram_lm,
13     constrain_space_mapping,
14     decode_with_emissions,
15     score_bigram_logprob,
16 )
17 from .mh import metropolis_patristocrat, metropolis_substitution
18 from .utils import ALPHABET, encode_text, load_brown_corpus_text,
19     load_pickled_language_model, write_ngram_lm_pickle
20 LM_PICKLE_DEFAULT = Path("resources/en/ngrams_space.pkl")
21 EM_EXAMPLE_CIPHERTEXT = (
22     "RBO RPXTIGO VCRB BWUCJA WJ KLOJ HCJD KM SKTPQO CQ RBWR "
23     "LOKLGO VCGG CJQCQR KJ SKHCJA WGKJA WJD RPYCJA RK LTR RBCJAQ CJ CR "
24 )
25
26
27 def run_em_demo(max_iter: int = 500, verbose: bool = False, restarts: int = 15) -> None:
28     """Run the Baum-Welch experiment using a precomputed bigram LM loaded from
29         `resources/en/ngrams_space.pkl`."""
30     num_states = len(ALPHABET)
31     lm_stats = load_pickled_language_model(LM_PICKLE_DEFAULT)
32     A, PI = bigram_hmm_from_ngram_lm(lm_stats)
33
34     # Use Baum-Welch to estimate emissions for this ciphertext.
35     V = np.array(encode_text(EM_EXAMPLE_CIPHERTEXT))
36     baum = BaumWelch()
37
38     bigram = lm_stats[2]
39
40     def score_bigram(text: str) -> float:
41         return score_bigram_logprob(text, bigram)
42
43     best_overall = {"text": "", "log": -1e300, "score": -1e300}
44
45     for r in range(restarts):
46         # Start from a random emission matrix; space forced later.
47         B = np.random.rand(num_states, num_states)
48
49         best_candidate = {"text": "", "log": -1e300, "score": -1e300, "B": None}
50
51         def verbose_fun(V_, B_, it, log_norm):
52             decoded = decode_with_emissions(B_, EM_EXAMPLE_CIPHERTEXT)
53             sc = score_bigram(decoded)
54             if verbose:
55                 print(f"[r{r}] #It", it, decoded, log_norm, f"(bg_score={sc:.2f})")
56             if sc > best_candidate["score"]:
57                 best_candidate["score"] = sc

```

```

57         best_candidate["text"] = decoded
58         best_candidate["log"] = log_norm
59         best_candidate["B"] = B_.copy()
60     if log_norm > best_candidate["log"]:
61         best_candidate["log"] = log_norm
62
63     _, B_hat, _, log_norm = baum.process(
64         A,
65         B,
66         PI,
67         V,
68         max_iter=max_iter,
69         ext_fun=constrain_space_mapping,
70         verbose_fun=verbose_fun,
71         recompute_A=False,
72         recompute_PI=False,
73         stop_on_decrease=False,
74     )
75     winner_decoded = best_candidate["text"] or decode_with_emissions(B_hat,
76         EM_EXAMPLE_CIPHERTEXT)
77     winner_log = best_candidate["log"] if best_candidate["text"] else log_norm
78     winner_score = best_candidate["score"] if best_candidate["text"] else
79         score_bigram(winner_decoded)
80     # Prefer the emissions matrix at the best LM score if captured
81     if best_candidate["B"] is not None and best_candidate["text"]:
82         winner_decoded = decode_with_emissions(best_candidate["B"], EM_EXAMPLE_CIPHERTEXT)
83     if winner_score > best_overall["score"]:
84         best_overall = {"text": winner_decoded, "log": winner_log, "score": winner_score}
85
86     print(best_overall["text"])
87     print("log_norm:", best_overall["log"])
88
89 def run_mh_substitution_demo(iterations: int = 1_000_000, *, verbose: bool = False) -> None:
90     """Metropolis substitution demo."""
91     print("Metropolis substitution:")
92     best = metropolis_substitution(
93         " K CXJO XI HBD BKIO XL ENJHB HEN XI HBD CSLB ",
94         iterations=iterations,
95         lm_pickle=Path("resources/en/ngrams_space.pkl"),
96         verbose=verbose,
97     )
98     print("\nBest candidate:")
99     print(best)
100
101 def run_mh_spacing_demo(iterations: int = 5000, *, verbose: bool = False) -> None:
102     """Metropolis spacing demo."""
103     print("\nMetropolis spacing:")
104     spaced, _ = metropolis_patristocrat(
105         " IIFYOULOOKFORPERFECTIONYOUULLNEVERBECONTENT ",
106         iterations=iterations,
107         lm_pickle=Path("resources/en/ngrams_space.pkl"),
108         verbose=verbose,

```

```

109     )
110     print("\nBest spacing:")
111     print(spaced)
112
113
114 def run_prepare_lm_pickle(
115     *,
116     out_path: Path,
117     corpus_path: Path | None,
118     min_n: int,
119     max_n: int,
120     keep_spaces: bool,
121     limit_chars: int | None,
122 ) -> None:
123     """Generate an n-gram language model pickle used by the demos."""
124     if corpus_path is None:
125         corpus_text = load_brown_corpus_text(keep_spaces=keep_spaces)
126     else:
127         corpus_text = corpus_path.read_text(encoding="utf-8")
128     out = write_ngram_lm_pickle(
129         corpus_text,
130         out_path,
131         min_n=min_n,
132         max_n=max_n,
133         keep_spaces=keep_spaces,
134         limit_chars=limit_chars,
135     )
136     print(f"Wrote: {out}")
137
138
139 def main() -> int:
140     parser = argparse.ArgumentParser(description="Run the EM and MH demos from the examples.")
141     sub = parser.add_subparsers(dest="cmd")
142
143     p_em = sub.add_parser("em", help="Run Baum-Welch demo")
144     p_em.add_argument("--max-iter", type=int, default=200)
145     p_em.add_argument("--verbose", action="store_true")
146     p_em.add_argument("--restarts", type=int, default=15)
147
148     p_mh = sub.add_parser("mh", help="Run Metropolis demos")
149     p_mh.add_argument("--verbose", action="store_true")
150     p_mh.add_argument("--iters-sub", type=int, default=500_000)
151     p_mh.add_argument("--iters-space", type=int, default=5000)
152
153     p_all = sub.add_parser("all", help="Run both demos")
154     p_all.add_argument("--max-iter", type=int, default=200)
155     p_all.add_argument("--verbose", action="store_true")
156     p_all.add_argument("--iters-sub", type=int, default=500_000)
157     p_all.add_argument("--iters-space", type=int, default=5000)
158     p_all.add_argument("--restarts", type=int, default=15)
159
160     p_prep = sub.add_parser("prep-lm", help="Generate the n-gram LM pickle (default: Brown
161                             corpus)")
162     p_prep.add_argument("--out", type=Path, default=Path("resources/en/ngrams_space.pkl"))

```

```

162     p_prep.add_argument("--corpus", type=Path, default=None, help="Path to a plain-text
163         corpus file (UTF-8)")
164     p_prep.add_argument("--min-n", type=int, default=1)
165     p_prep.add_argument("--max-n", type=int, default=7)
166     p_prep.add_argument("--no-spaces", dest="keep_spaces", action="store_false", help="Remove
167         spaces from the corpus")
168     p_prep.add_argument("--limit-chars", type=int, default=0, help="Optionally truncate the
169         normalized corpus")
170     p_prep.set_defaults(keep_spaces=True)
171
172     args = parser.parse_args()
173     if args.cmd is None:
174         parser.print_help()
175         return 0
176
177     if args.cmd in {"em", "all"}:
178         run_em_demo(max_iter=args.max_iter, verbose=args.verbose, restarts=args.restarts)
179     if args.cmd in {"mh", "all"}:
180         run_mh_substitution_demo(iterations=args.iters_sub, verbose=args.verbose)
181         run_mh_spacing_demo(iterations=args.iters_space, verbose=args.verbose)
182     if args.cmd == "prep-lm":
183         limit = None if args.limit_chars == 0 else args.limit_chars
184         run_prepare_lm_pickle(
185             out_path=args.out,
186             corpus_path=args.corpus,
187             min_n=args.min_n,
188             max_n=args.max_n,
189             keep_spaces=args.keep_spaces,
190             limit_chars=limit,
191         )
192     return 0
193
194 if __name__ == "__main__":
195     raise SystemExit(main())

```

## References

- [1] Substitution cipher. (2021, July 27). In Wikipedia. [https://en.wikipedia.org/w/index.php?title=Substitution\\_cipher&oldid=1035776944](https://en.wikipedia.org/w/index.php?title=Substitution_cipher&oldid=1035776944)
- [2] Knight, Kevin & Nair, Anish & Rathod, Nishit & Yamada, Kenji. (2006). Unsupervised Analysis for Decipherment Problems. <https://dx.doi.org/10.3115/1273073.1273138>
- [3] Francis, W. Nelson & Henry Kucera. 1967. Computational Analysis of Present-Day American English. Providence, RI: Brown University Press.
- [4] Dempster, A.P.; Laird, N.M.; Rubin, D.B. (1977). “Maximum Likelihood from Incomplete Data via the EM Algorithm”. Journal of the Royal Statistical Society, Series B. 39 (1): 1-38. JSTOR 2984875. MR 0501537.
- [5] Maya R. Gupta and Yihua Chen (2011), “Theory and Use of the EM Algorithm”, Foundations and Trends® in Signal Processing: Vol. 4: No. 3, pp 223-296. <https://dx.doi.org/10.1561/2000000034>
- [6] L. R. Rabiner, “A tutorial on hidden Markov models and selected applications in speech recognition,” in Proceedings of the IEEE, vol. 77, no. 2, pp. 257-286, Feb. 1989. <https://dx.doi.org/10.1109/5.18626>

- [7] A. Viterbi, "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm," in IEEE Transactions on Information Theory, vol. 13, no. 2, pp. 260-269, April 1967, doi: 10.1109/TIT.1967.1054010.
- [8] Shannon, Claude. "Communication Theory of Secrecy Systems", Bell System Technical Journal, vol. 28(4), page 656-715, 1949. <https://doi.org/10.1002/j.1538-7305.1949.tb00928.x>
- [9] Gagniuc, Paul A. (2017). Markov Chains: From Theory to Implementation and Experimentation. USA, NJ: John Wiley & Sons. pp. 1-256. ISBN 978-1-119-38755-8
- [10] Diaconis, Persi. (2009). The Markov Chain Monte Carlo Revolution. Bulletin of the American Mathematical Society. 46. 179-205. <https://dx.doi.org/10.1090/S0273-0979-08-01238-X>
- [11] Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., & Teller, E. (1953). Equation of State Calculations by Fast Computing Machines. The Journal of Chemical Physics, 21(6), 1087–1092. <https://doi.org/10.1063/1.1699114>