

# Deriving the E-Step and M-Step for Hidden Markov Models

From First Principles to Update Equations

*A step-by-step mathematical derivation with intuition,  
connecting each formula to the project's implementation*

## Contents

<b>1 The Problem: Why Do We Need EM?</b>	<b>2</b>
1.1 What We Want . . . . .	2
1.2 Why Direct Maximisation Fails . . . . .	2
<b>2 The EM Framework: The <math>\mathcal{Q}</math>-Function</b>	<b>3</b>
2.1 The Key Idea . . . . .	3
2.2 Why Maximising $\mathcal{Q}$ Works . . . . .	3
<b>3 Expanding the Complete-Data Log-Likelihood</b>	<b>4</b>
3.1 The Joint Probability Factorisation . . . . .	4
3.2 The $\mathcal{Q}$ -Function in Terms of $\gamma$ and $\xi$ . . . . .	4
3.2.1 Introduce indicator variables . . . . .	4
3.2.2 Take the expectation . . . . .	4
<b>4 Deriving the E-Step: Why <math>\gamma</math> and <math>\xi</math> Take Those Forms</b>	<b>6</b>
4.1 Deriving $\gamma_t(i)$ . . . . .	6
4.2 Deriving $\xi_t(i, j)$ . . . . .	6
4.3 Consistency Check: $\gamma$ from $\xi$ . . . . .	7
<b>5 Deriving the M-Step: The Update Equations</b>	<b>8</b>
5.1 Deriving $\hat{A}_{ij}$ : The Transition Matrix Update . . . . .	8
5.2 Deriving $\hat{\pi}$ : The Initial Distribution Update . . . . .	9
5.3 Deriving $\hat{\mu}_k$ : The Emission Mean Update . . . . .	10
5.4 Deriving $\hat{\Sigma}_k$ : The Emission Covariance Update . . . . .	12
<b>6 Connection to the Implementation</b>	<b>13</b>
<b>7 Summary: The Complete Derivation Map</b>	<b>14</b>
7.1 All Update Equations at a Glance . . . . .	14

# 1 The Problem: Why Do We Need EM?

## 1.1 What We Want

We have observed data  $\mathbf{O} = (O_1, \dots, O_T)$  (daily sector returns) and we want to find the HMM parameters  $\lambda = (\mathbf{A}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \boldsymbol{\pi})$  that **maximise the marginal likelihood**:

$$\hat{\lambda} = \arg \max_{\lambda} P(\mathbf{O} | \lambda) \quad (1)$$

## 1.2 Why Direct Maximisation Fails

The marginal likelihood requires summing over all hidden state sequences:

$$P(\mathbf{O} | \lambda) = \sum_{\mathbf{Q}} P(\mathbf{O}, \mathbf{Q} | \lambda) \quad (2)$$

where the sum is over all  $K^T$  possible state sequences  $\mathbf{Q} = (q_1, \dots, q_T)$ .

If we try to take the derivative and set it to zero:

$$\frac{\partial}{\partial \lambda} \log \underbrace{\sum_{\mathbf{Q}} P(\mathbf{O}, \mathbf{Q} | \lambda)}_{\text{sum of exponentials}} = 0$$

The log of a sum has no closed-form solution. The parameters are entangled inside the summation in a way that prevents direct optimisation.

### The Core Difficulty

If we *knew* the hidden states  $\mathbf{Q}$ , maximisation would be trivial (just count transitions and compute sample means). But  $\mathbf{Q}$  is unknown. If we *knew* the parameters  $\lambda$ , we could infer  $\mathbf{Q}$ . But  $\lambda$  is also unknown. This chicken-and-egg problem is exactly what EM resolves.

## 2 The EM Framework: The $\mathcal{Q}$ -Function

### 2.1 The Key Idea

Instead of maximising  $\log P(\mathbf{O} | \lambda)$  directly, EM works with the **complete-data log-likelihood**  $\log P(\mathbf{O}, \mathbf{Q} | \lambda)$ , which is much easier to handle because the log can be pushed inside the product.

Since  $\mathbf{Q}$  is unknown, we take its **expectation** under the posterior  $P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}})$  computed using our current best guess  $\lambda^{\text{old}}$ .

#### The $\mathcal{Q}$ -Function (Auxiliary Function)

$$\mathcal{Q}(\lambda, \lambda^{\text{old}}) = \mathbb{E}_{\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}}} [\log P(\mathbf{O}, \mathbf{Q} | \lambda)] = \sum_{\mathbf{Q}} P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}}) \log P(\mathbf{O}, \mathbf{Q} | \lambda) \quad (3)$$

- $\lambda^{\text{old}}$ : current parameters (used to compute the expectation)
- $\lambda$ : new parameters (what we optimise over)

### 2.2 Why Maximising $\mathcal{Q}$ Works

**Theorem 2.1** (EM Monotonicity). *If  $\lambda^{\text{new}} = \arg \max_{\lambda} \mathcal{Q}(\lambda, \lambda^{\text{old}})$ , then:*

$$\log P(\mathbf{O} | \lambda^{\text{new}}) \geq \log P(\mathbf{O} | \lambda^{\text{old}})$$

*Each EM iteration is guaranteed to increase (or maintain) the marginal log-likelihood.*

#### Proof Sketch (Jensen's Inequality)

Write:

$$\begin{aligned} \log P(\mathbf{O} | \lambda) &= \log \sum_{\mathbf{Q}} P(\mathbf{O}, \mathbf{Q} | \lambda) \\ &= \log \sum_{\mathbf{Q}} P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}}) \frac{P(\mathbf{O}, \mathbf{Q} | \lambda)}{P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}})} \\ &\geq \sum_{\mathbf{Q}} P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}}) \log \frac{P(\mathbf{O}, \mathbf{Q} | \lambda)}{P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}})} \quad (\text{Jensen's inequality: log is concave}) \\ &= \underbrace{\sum_{\mathbf{Q}} P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}}) \log P(\mathbf{O}, \mathbf{Q} | \lambda)}_{\mathcal{Q}(\lambda, \lambda^{\text{old}})} - \underbrace{\sum_{\mathbf{Q}} P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}}) \log P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}})}_{\text{entropy; constant w.r.t. } \lambda} \end{aligned}$$

The second term doesn't depend on  $\lambda$ , so maximising the lower bound over  $\lambda$  reduces to maximising  $\mathcal{Q}(\lambda, \lambda^{\text{old}})$ . Since the bound is tight at  $\lambda = \lambda^{\text{old}}$ , any increase in  $\mathcal{Q}$  translates to an increase in  $\log P(\mathbf{O} | \lambda)$ .

### 3 Expanding the Complete-Data Log-Likelihood

#### 3.1 The Joint Probability Factorisation

From the HMM's conditional independence assumptions:

$$P(\mathbf{O}, \mathbf{Q} | \lambda) = \pi_{q_1} \cdot b_{q_1}(O_1) \cdot \prod_{t=2}^T A_{q_{t-1}, q_t} \cdot b_{q_t}(O_t) \quad (4)$$

Taking the logarithm:

#### Complete-Data Log-Likelihood

$$\log P(\mathbf{O}, \mathbf{Q} | \lambda) = \underbrace{\log \pi_{q_1}}_{(I)} + \underbrace{\sum_{t=2}^T \log A_{q_{t-1}, q_t}}_{(II)} + \underbrace{\sum_{t=1}^T \log b_{q_t}(O_t)}_{(III)} \quad (5)$$

#### Why This is Easier

The log of a product became a **sum** of logs. Each term depends on the hidden state at one or two time steps only. This separability is what makes the M-step tractable.  
Compare with  $\log P(\mathbf{O} | \lambda) = \log \sum_{\mathbf{Q}} P(\mathbf{O}, \mathbf{Q} | \lambda)$ , where the log is *outside* the sum—completely intractable.

#### 3.2 The $Q$ -Function in Terms of $\gamma$ and $\xi$

Now we compute  $\mathcal{Q}(\lambda, \lambda^{\text{old}}) = \mathbb{E}_{\mathbf{Q}|\mathbf{O}, \lambda^{\text{old}}} [\log P(\mathbf{O}, \mathbf{Q} | \lambda)]$  by taking the expectation of each term in Eq. (5).

##### 3.2.1 Introduce indicator variables

To handle the expectation formally, define indicator variables:

$$\mathbb{1}[q_t = i] = \begin{cases} 1 & \text{if state at time } t \text{ is } i \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

Then:

$$\log \pi_{q_1} = \sum_{i=1}^K \mathbb{1}[q_1 = i] \log \pi_i \quad (7)$$

$$\sum_{t=2}^T \log A_{q_{t-1}, q_t} = \sum_{t=2}^T \sum_{i=1}^K \sum_{j=1}^K \mathbb{1}[q_{t-1} = i, q_t = j] \log A_{ij} \quad (8)$$

$$\sum_{t=1}^T \log b_{q_t}(O_t) = \sum_{t=1}^T \sum_{k=1}^K \mathbb{1}[q_t = k] \log b_k(O_t) \quad (9)$$

##### 3.2.2 Take the expectation

Under  $P(\mathbf{Q} | \mathbf{O}, \lambda^{\text{old}})$ , the expectations of the indicators are exactly  $\gamma$  and  $\xi$ :

### This is Where $\gamma$ and $\xi$ Come From

$$\mathbb{E}[\mathbb{1}[q_t = i]] = P(q_t = i \mid \mathbf{O}, \lambda^{\text{old}}) = \gamma_t(i) \quad (10)$$

$$\mathbb{E}[\mathbb{1}[q_{t-1} = i, q_t = j]] = P(q_{t-1} = i, q_t = j \mid \mathbf{O}, \lambda^{\text{old}}) = \xi_t(i, j) \quad (11)$$

**This is why the E-step computes  $\gamma$  and  $\xi$ :** they are exactly the expected sufficient statistics needed to evaluate  $\mathcal{Q}$ .

Substituting:

### $\mathcal{Q}$ -Function Expanded

$$\mathcal{Q}(\lambda, \lambda^{\text{old}}) = \underbrace{\sum_{i=1}^K \gamma_1(i) \log \pi_i}_{\text{(I): initial state}} + \underbrace{\sum_{t=2}^T \sum_{i=1}^K \sum_{j=1}^K \xi_t(i, j) \log A_{ij}}_{\text{(II): transitions}} + \underbrace{\sum_{t=1}^T \sum_{k=1}^K \gamma_t(k) \log b_k(O_t)}_{\text{(III): emissions}} \quad (12)$$

### The Crucial Observation

The three terms in  $\mathcal{Q}$  involve **different parameters**: Term (I) involves only  $\boldsymbol{\pi}$ , Term (II) involves only  $\mathbf{A}$ , and Term (III) involves only  $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ . We can **maximise each term independently**.

## 4 Deriving the E-Step: Why $\gamma$ and $\xi$ Take Those Forms

We've just shown that the E-step *needs*  $\gamma_t(i)$  and  $\xi_t(i, j)$ . Now let us derive their formulas.

### 4.1 Deriving $\gamma_t(i)$

We want:

$$\gamma_t(i) = P(q_t = i \mid \mathbf{O}, \lambda)$$

#### Derivation of $\gamma_t(i)$

By Bayes' theorem:

$$\gamma_t(i) = P(q_t = i \mid \mathbf{O}, \lambda) = \frac{P(q_t = i, \mathbf{O} \mid \lambda)}{P(\mathbf{O} \mid \lambda)} \quad (13)$$

Now split  $\mathbf{O}$  into past observations  $(O_1, \dots, O_t)$  and future observations  $(O_{t+1}, \dots, O_T)$ :

$$P(q_t = i, \mathbf{O} \mid \lambda) = P(O_1, \dots, O_t, q_t = i, O_{t+1}, \dots, O_T \mid \lambda) \quad (14)$$

$$= \underbrace{P(O_1, \dots, O_t, q_t = i \mid \lambda)}_{\alpha_t(i)} \cdot \underbrace{P(O_{t+1}, \dots, O_T \mid q_t = i, \lambda)}_{\beta_t(i)} \quad (15)$$

The second equality uses the fact that, given  $q_t = i$ , the future observations are conditionally independent of the past (by the Markov property and output independence). This is the product  $\alpha_t(i) \cdot \beta_t(i)$ .

The denominator is just the total probability of the data:

$$P(\mathbf{O} \mid \lambda) = \sum_{j=1}^K P(q_t = j, \mathbf{O} \mid \lambda) = \sum_{j=1}^K \alpha_t(j) \beta_t(j) \quad (16)$$

Substituting into Eq. (13):

$$\boxed{\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^K \alpha_t(j) \beta_t(j)}} \quad (17)$$

#### Intuition

$\alpha_t(i)$  summarises all evidence from the **past** that we're in state  $i$  at time  $t$ .  $\beta_t(i)$  summarises all evidence from the **future**. Multiplying them combines both directions of evidence. Dividing by the total normalises to a proper probability.

### 4.2 Deriving $\xi_t(i, j)$

We want:

$$\xi_t(i, j) = P(q_t = i, q_{t+1} = j \mid \mathbf{O}, \lambda)$$

#### Derivation of $\xi_t(i, j)$

Again, by Bayes' theorem:

$$\xi_t(i, j) = \frac{P(q_t = i, q_{t+1} = j, \mathbf{O} \mid \lambda)}{P(\mathbf{O} \mid \lambda)} \quad (18)$$

Expand the numerator by splitting  $\mathbf{O}$  at time  $t$ :

$$\begin{aligned} P(q_t = i, q_{t+1} = j, \mathbf{O} \mid \lambda) \\ = P(O_1, \dots, O_t, q_t = i \mid \lambda) \cdot P(q_{t+1} = j \mid q_t = i, \lambda) \cdot P(O_{t+1} \mid q_{t+1} = j, \lambda) \\ \cdot P(O_{t+2}, \dots, O_T \mid q_{t+1} = j, \lambda) \end{aligned} \quad (19)$$

Identifying each factor:

$$= \underbrace{\alpha_t(i)}_{\text{past evidence}} \cdot \underbrace{A_{ij}}_{\text{transition}} \cdot \underbrace{b_j(O_{t+1})}_{\text{emission at } t+1} \cdot \underbrace{\beta_{t+1}(j)}_{\text{future evidence}} \quad (20)$$

The denominator is  $P(\mathbf{O} \mid \lambda) = \sum_{i'} \sum_{j'} \alpha_t(i') A_{i'j'} b_{j'}(O_{t+1}) \beta_{t+1}(j')$ .

Therefore:

$$\boxed{\xi_t(i, j) = \frac{\alpha_t(i) A_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i'=1}^K \sum_{j'=1}^K \alpha_t(i') A_{i'j'} b_{j'}(O_{t+1}) \beta_{t+1}(j')}} \quad (21)$$

### Intuition for $\xi$

To compute the probability of the specific transition  $i \rightarrow j$  at time  $t$ , we need four things:

1.  $\alpha_t(i)$ : evidence from days  $1, \dots, t$  that we're in state  $i$
2.  $A_{ij}$ : the probability of transitioning from  $i$  to  $j$
3.  $b_j(O_{t+1})$ : how well state  $j$  explains tomorrow's observation
4.  $\beta_{t+1}(j)$ : evidence from days  $t + 2, \dots, T$  that we were in state  $j$  at  $t + 1$

Multiply all four, then normalise.

### 4.3 Consistency Check: $\gamma$ from $\xi$

Notice that  $\gamma_t(i)$  is the **marginal** of  $\xi_t(i, j)$ :

$$\gamma_t(i) = \sum_{j=1}^K \xi_t(i, j) \quad (22)$$

This makes sense: the probability of being in state  $i$  at time  $t$  equals the sum of probabilities of being in state  $i$  and transitioning to *any* state  $j$ .

## 5 Deriving the M-Step: The Update Equations

We now maximise  $\mathcal{Q}(\lambda, \lambda^{\text{old}})$  from Eq. (12) with respect to each group of parameters. Since the three terms are separable, we handle each independently.

### 5.1 Deriving $\hat{A}_{ij}$ : The Transition Matrix Update

We maximise Term (II) of  $\mathcal{Q}$ :

$$\max_{\mathbf{A}} \sum_{t=2}^T \sum_{i=1}^K \sum_{j=1}^K \xi_t(i, j) \log A_{ij} \quad \text{subject to} \quad \sum_{j=1}^K A_{ij} = 1 \quad \forall i$$

#### Derivation via Lagrange Multipliers

The constraint is that each row of  $\mathbf{A}$  sums to 1. We introduce a Lagrange multiplier  $\lambda_i$  for each row  $i$ :

$$\mathcal{L} = \sum_{t=2}^T \sum_{i=1}^K \sum_{j=1}^K \xi_t(i, j) \log A_{ij} - \sum_{i=1}^K \lambda_i \left( \sum_{j=1}^K A_{ij} - 1 \right) \quad (23)$$

Take the derivative with respect to  $A_{ij}$  and set to zero:

$$\frac{\partial \mathcal{L}}{\partial A_{ij}} = \frac{\sum_{t=2}^T \xi_t(i, j)}{A_{ij}} - \lambda_i = 0 \quad (24)$$

Solving for  $A_{ij}$ :

$$A_{ij} = \frac{\sum_{t=2}^T \xi_t(i, j)}{\lambda_i} \quad (25)$$

Now determine  $\lambda_i$  using the constraint  $\sum_j A_{ij} = 1$ :

$$\sum_{j=1}^K A_{ij} = \frac{\sum_{j=1}^K \sum_{t=2}^T \xi_t(i, j)}{\lambda_i} = 1 \quad \Rightarrow \quad \lambda_i = \sum_{j=1}^K \sum_{t=2}^T \xi_t(i, j) = \sum_{t=2}^T \gamma_t(i) \quad (26)$$

where we used  $\sum_j \xi_t(i, j) = \gamma_t(i)$ . Substituting back into Eq. (25):

$$\hat{A}_{ij} = \frac{\sum_{t=2}^T \xi_t(i, j)}{\sum_{t=2}^T \gamma_t(i)} = \frac{\text{expected number of transitions } i \rightarrow j}{\text{expected number of times in state } i} \quad (27)$$

#### Intuition: Relative Frequency

If the soft counts say: “state Bull appeared 100 times (expected), and 90 of those transitioned to Bull, 7 to Sideways, 3 to Bear,” then:

$$\hat{A}_{\text{Bull}, \text{Bull}} = \frac{90}{100} = 0.90, \quad \hat{A}_{\text{Bull}, \text{Side}} = \frac{7}{100} = 0.07, \quad \hat{A}_{\text{Bull}, \text{Bear}} = \frac{3}{100} = 0.03$$

This is the same formula as the maximum likelihood estimate for a multinomial, but with **soft** counts from  $\gamma$  and  $\xi$  instead of hard counts.

## 5.2 Deriving $\hat{\pi}$ : The Initial Distribution Update

The derivation is nearly identical. Maximise Term (I):

$$\max_{\boldsymbol{\pi}} \sum_{i=1}^K \gamma_1(i) \log \pi_i \quad \text{subject to} \quad \sum_{i=1}^K \pi_i = 1$$

### Derivation

Lagrangian:  $\mathcal{L} = \sum_i \gamma_1(i) \log \pi_i - \mu (\sum_i \pi_i - 1)$

Taking derivative and setting to zero:

$$\frac{\partial \mathcal{L}}{\partial \pi_i} = \frac{\gamma_1(i)}{\pi_i} - \mu = 0 \quad \Rightarrow \quad \pi_i = \frac{\gamma_1(i)}{\mu}$$

Using  $\sum_i \pi_i = 1$ :  $\mu = \sum_i \gamma_1(i) = 1$ , so:

$$\hat{\pi}_i = \gamma_1(i) \quad (28)$$

This is intuitive: the probability of starting in state  $i$  is simply the posterior probability of being in state  $i$  at  $t = 1$ .

### 5.3 Deriving $\hat{\mu}_k$ : The Emission Mean Update

Now we maximise Term (III). For Gaussian emissions,  $b_k(O_t) = \mathcal{N}(\mathbf{r}_t; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ , so:

$$\log b_k(\mathbf{r}_t) = -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{r}_t - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{r}_t - \boldsymbol{\mu}_k) \quad (29)$$

Term (III) of  $\mathcal{Q}$ , collecting only terms involving  $\boldsymbol{\mu}_k$  for a specific regime  $k$ :

$$\mathcal{Q}_k^{(\mu)} = -\frac{1}{2} \sum_{t=1}^T \gamma_t(k) (\mathbf{r}_t - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{r}_t - \boldsymbol{\mu}_k) + \text{const} \quad (30)$$

#### Derivation of $\hat{\mu}_k$

Take the derivative with respect to  $\boldsymbol{\mu}_k$  and set to zero.

Using the matrix calculus identity  $\frac{\partial}{\partial \boldsymbol{\mu}} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{M} (\mathbf{x} - \boldsymbol{\mu}) = -2\mathbf{M}(\mathbf{x} - \boldsymbol{\mu})$ :

$$\frac{\partial \mathcal{Q}_k^{(\mu)}}{\partial \boldsymbol{\mu}_k} = -\frac{1}{2} \sum_{t=1}^T \gamma_t(k) \cdot (-2) \boldsymbol{\Sigma}_k^{-1} (\mathbf{r}_t - \boldsymbol{\mu}_k) = \mathbf{0} \quad (31)$$

$$= \sum_{t=1}^T \gamma_t(k) \boldsymbol{\Sigma}_k^{-1} (\mathbf{r}_t - \boldsymbol{\mu}_k) = \mathbf{0} \quad (32)$$

Since  $\boldsymbol{\Sigma}_k^{-1}$  is non-singular, we can multiply both sides by  $\boldsymbol{\Sigma}_k$ :

$$\sum_{t=1}^T \gamma_t(k) (\mathbf{r}_t - \boldsymbol{\mu}_k) = \mathbf{0} \quad (33)$$

$$\sum_{t=1}^T \gamma_t(k) \mathbf{r}_t = \boldsymbol{\mu}_k \sum_{t=1}^T \gamma_t(k) \quad (34)$$

Solving for  $\boldsymbol{\mu}_k$ :

$$\boxed{\hat{\boldsymbol{\mu}}_k = \frac{\sum_{t=1}^T \gamma_t(k) \mathbf{r}_t}{\sum_{t=1}^T \gamma_t(k)}} \quad (35)$$

#### Intuition: Weighted Average

$\hat{\boldsymbol{\mu}}_k$  is the weighted average of all observations, where the weight of observation  $\mathbf{r}_t$  is  $\gamma_t(k)$ —the probability that  $\mathbf{r}_t$  was generated by regime  $k$ .

**Concrete example:** Suppose we have 3 days of returns and  $K = 2$  regimes:

$t$	$\mathbf{r}_t$	$\gamma_t(\text{Bull})$	$\gamma_t(\text{Bear})$
1	+0.5%	0.9	0.1
2	+0.3%	0.8	0.2
3	-0.4%	0.1	0.9

$$\hat{\mu}_{\text{Bull}} = \frac{0.9 \times 0.5 + 0.8 \times 0.3 + 0.1 \times (-0.4)}{0.9 + 0.8 + 0.1} = \frac{0.45 + 0.24 - 0.04}{1.8} = \frac{0.65}{1.8} = +0.361\%$$

$$\hat{\mu}_{\text{Bear}} = \frac{0.1 \times 0.5 + 0.2 \times 0.3 + 0.9 \times (-0.4)}{0.1 + 0.2 + 0.9} = \frac{0.05 + 0.06 - 0.36}{1.2} = \frac{-0.25}{1.2} = -0.208\%$$

Day 3 ( $-0.4\%$ ) barely affects  $\hat{\mu}_{\text{Bull}}$  (weight 0.1) but dominates  $\hat{\mu}_{\text{Bear}}$  (weight 0.9). The soft assignments naturally route each observation to the regime that likely generated it.

## 5.4 Deriving $\hat{\Sigma}_k$ : The Emission Covariance Update

Collecting terms from  $\mathcal{Q}$  that involve  $\Sigma_k$ :

$$\mathcal{Q}_k^{(\Sigma)} = -\frac{1}{2} \sum_{t=1}^T \gamma_t(k) [\log |\Sigma_k| + (\mathbf{r}_t - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{r}_t - \boldsymbol{\mu}_k)] + \text{const} \quad (36)$$

### Derivation of $\hat{\Sigma}_k$

We use two matrix calculus identities. Let  $\mathbf{S} = \Sigma_k^{-1}$  (the precision matrix):

1.  $\frac{\partial}{\partial \mathbf{S}} \log |\mathbf{S}^{-1}| = -\mathbf{S}^{-1}$ , equivalently  $\frac{\partial}{\partial \Sigma_k} \log |\Sigma_k| = \Sigma_k^{-1}$
2.  $\frac{\partial}{\partial \Sigma_k^{-1}} \mathbf{x}^T \Sigma_k^{-1} \mathbf{x} = \mathbf{x} \mathbf{x}^T$

It is cleaner to differentiate with respect to  $\mathbf{S} = \Sigma_k^{-1}$ . Rewriting  $\mathcal{Q}_k^{(\Sigma)}$ :

$$\mathcal{Q}_k^{(\Sigma)} = -\frac{1}{2} \sum_{t=1}^T \gamma_t(k) \left[ -\log |\mathbf{S}| + (\mathbf{r}_t - \boldsymbol{\mu}_k)^T \mathbf{S} (\mathbf{r}_t - \boldsymbol{\mu}_k) \right] \quad (37)$$

Taking the derivative with respect to  $\mathbf{S}$  and setting to zero:

$$\frac{\partial \mathcal{Q}_k^{(\Sigma)}}{\partial \mathbf{S}} = -\frac{1}{2} \sum_{t=1}^T \gamma_t(k) \left[ -\mathbf{S}^{-1} + (\mathbf{r}_t - \boldsymbol{\mu}_k)(\mathbf{r}_t - \boldsymbol{\mu}_k)^T \right] = \mathbf{0} \quad (38)$$

$$\sum_{t=1}^T \gamma_t(k) \mathbf{S}^{-1} = \sum_{t=1}^T \gamma_t(k) (\mathbf{r}_t - \boldsymbol{\mu}_k)(\mathbf{r}_t - \boldsymbol{\mu}_k)^T \quad (39)$$

$$\Sigma_k \cdot \sum_{t=1}^T \gamma_t(k) = \sum_{t=1}^T \gamma_t(k) (\mathbf{r}_t - \boldsymbol{\mu}_k)(\mathbf{r}_t - \boldsymbol{\mu}_k)^T \quad (40)$$

Solving:

$$\hat{\Sigma}_k = \frac{\sum_{t=1}^T \gamma_t(k) (\mathbf{r}_t - \hat{\boldsymbol{\mu}}_k)(\mathbf{r}_t - \hat{\boldsymbol{\mu}}_k)^T}{\sum_{t=1}^T \gamma_t(k)} \quad (41)$$

### Intuition: Weighted Sample Covariance

This is the standard sample covariance formula, but with soft weights  $\gamma_t(k)$ . Only observations that “belong to” regime  $k$  (high  $\gamma_t(k)$ ) contribute significantly to that regime’s covariance estimate.

## 6 Connection to the Implementation

All of the above is executed internally by `hmmlearn` when `model.fit(X)` is called (`models.py`, line 174). The project wraps this with multiple initialisations:

```

1 # models.py, lines 161-183
2 for i in range(n_init):
3     model = hmm.GaussianHMM(
4         n_components=self.n_regimes,          # K = 3
5         covariance_type=self.covariance_type, # 'full',
6         n_iter=self.n_iter,                  # max 100 EM iterations
7         tol=self.tol,                      # convergence threshold 1e-4
8     )
9     model.fit(X)                         # <-- EM runs here!
10    # Internally:
11    #     for iter in range(n_iter):
12    #         alpha, beta = forward_backward(X) --> E-step
13    #         gamma, xi = compute_posteriors(alpha, beta)
14    #         update A, mu, Sigma using gamma, xi --> M-step
15    #         if improvement < tol: break
16
17    self._regularize_covariance(model)    # Sigma_k += lambda * I
18    score = model.score(X)              # final log P(0|lambda)

```

After EM finishes, the learned parameters are stored as:

```

1 # models.py, lines 321-340
2 model.transmat_   # A_hat: the transition matrix
3 model.means_      # mu_hat: emission means (K x D array)
4 model.covars_     # Sigma_hat: emission covariances (K x D x D)

```

Then the covariance regularisation (`models.py`, lines 104-120) applies the stabilising correction:

```

1 # After EM, prevent singular covariance:
2 for k in range(self.n_regimes):
3     n_features = model.covars_[k].shape[0]
4     model.covars_[k] += self.regularization * np.eye(n_features)
5     # This is: Sigma_k <-- Sigma_k + lambda * I

```

## 7 Summary: The Complete Derivation Map

### The Logical Flow of the Entire Derivation

1. We want  $\hat{\lambda} = \arg \max_{\lambda} \log P(\mathbf{O} | \lambda)$  but can't optimise directly
2. Jensen's inequality gives us a tractable lower bound: the  $\mathcal{Q}$ -function
3. Writing out  $\mathcal{Q}$  requires the complete-data log-likelihood  $\log P(\mathbf{O}, \mathbf{Q} | \lambda)$
4. Taking expectation over  $\mathbf{Q}$  introduces two sufficient statistics:
  - $\gamma_t(i) = \mathbb{E}[\mathbb{1}(q_t = i)]$  — derived via Bayes' rule from  $\alpha_t(i)$  and  $\beta_t(i)$
  - $\xi_t(i, j) = \mathbb{E}[\mathbb{1}(q_t = i, q_{t+1} = j)]$  — derived similarly
5. This is the E-step: compute  $\gamma$  and  $\xi$  using forward-backward with current  $\lambda^{\text{old}}$
6.  $\mathcal{Q}$  separates into three independent terms (initial, transitions, emissions)
7. Maximising each via Lagrange multipliers / calculus gives:
  - $\hat{A}_{ij}$ : ratio of expected transitions  $i \rightarrow j$  to expected time in state  $i$
  - $\hat{\mu}_k$ :  $\gamma$ -weighted average of observations
  - $\hat{\Sigma}_k$ :  $\gamma$ -weighted sample covariance
8. This is the M-step: update parameters using the formulas above
9. Iterate E-step  $\rightarrow$  M-step until convergence

### 7.1 All Update Equations at a Glance

#### Complete Set of EM Update Equations for Gaussian HMM

$$\hat{\pi}_i = \gamma_1(i) \quad (\text{Initial})$$

$$\hat{A}_{ij} = \frac{\sum_{t=2}^T \xi_t(i, j)}{\sum_{t=2}^T \gamma_t(i)} \quad (\text{Transitions})$$

$$\hat{\mu}_k = \frac{\sum_{t=1}^T \gamma_t(k) \mathbf{r}_t}{\sum_{t=1}^T \gamma_t(k)} \quad (\text{Means})$$

$$\hat{\Sigma}_k = \frac{\sum_{t=1}^T \gamma_t(k) (\mathbf{r}_t - \hat{\mu}_k)(\mathbf{r}_t - \hat{\mu}_k)^T}{\sum_{t=1}^T \gamma_t(k)} \quad (\text{Covariances})$$

where:

$$\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{\sum_j \alpha_t(j) \beta_t(j)} \quad (\text{E-step 1})$$

$$\xi_t(i, j) = \frac{\alpha_t(i) A_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i'} \sum_{j'} \alpha_t(i') A_{i'j'} b_{j'}(O_{t+1}) \beta_{t+1}(j')} \quad (\text{E-step 2})$$