

MLPR Tutorial¹ Sheet 2

Reminders: Attempt the tutorial questions *before* your tutorial. We strongly recommend you discuss these questions — and the course in general — with your peers. You could work at ML-Base, or arrange to meet up with people from your tutorial group (e.g., through Learn) or other people you have met in lectures. Detailed answers will be made available after the week's tutorials, and can be discussed further on Hypothesis.

This tutorial is largely just maths. While we ask for a small numerical experiment in the middle, there's no real data, and no machine learning. However, throughout the course we will derive models and algorithms that use multivariate Gaussian distributions. And other machine learning methods share some of the same maths. We've put this material on the tutorial, because it's useful stuff, and you need to work through it at your own pace. In the mean time, you're also working on assignment 1, which involves some data!

1. Warm-up exercise:

If \mathbf{a} and \mathbf{b} are $D \times 1$ column vectors and M is a $D \times D$ symmetric matrix, show that

$$\mathbf{a}^\top M \mathbf{b} = \mathbf{b}^\top M \mathbf{a}.$$

You wouldn't need to show this result in an exam unless you were explicitly asked to. In some working (like for the next question) you could just state that it's true for symmetric M .

2. Identifying a Gaussian:

As part of a derivation, we may need to identify the probability density function of a vector \mathbf{x} , up to a constant with respect to \mathbf{x} . For example:

$$p(\mathbf{x}) \propto \exp\left(-\mathbf{x}^\top A \mathbf{x} - \mathbf{x}^\top \mathbf{c}\right),$$

where A is a symmetric invertible matrix. As this distribution is proportional to the exponential of a quadratic in \mathbf{x} , it is a Gaussian: $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$.

Identify which Gaussian \mathbf{x} comes from by identifying the mean $\boldsymbol{\mu}$ and covariance Σ in terms of A and \mathbf{c} . The easiest method is to compare $p(\mathbf{x})$ to the standard form for the multivariate Gaussian PDF (given in class).

The answer you should be able to show is:

$$\Sigma = \frac{1}{2} A^{-1}, \quad \boldsymbol{\mu} = -\frac{1}{2} A^{-1} \mathbf{c}.$$

3. Creating a 2D multivariate Gaussian, and a simple experiment:

The first element of a vector has $p(x_1) = \mathcal{N}(x_1; m, \sigma^2)$.

A second element is generated according to the following process:

$$x_2 = \alpha x_1 + v, \quad v \sim \mathcal{N}(0, n^2).$$

Here x_2 depends on x_1 , but the noise term v is independent of x_1 .

Recall that a linear combination of Gaussian values is Gaussian distributed.

- a) The joint distribution of the vector $\mathbf{x} = [x_1 \ x_2]^\top$ is Gaussian, and so takes the form $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$. Identify $\boldsymbol{\mu}$ and Σ .

1. Parts of this tutorial sheet are based on previous versions by Amos Storkey, Charles Sutton, and Chris Williams

- b) Turning to a computer: pick a setting for each of the parameters m, σ, α , and n , and simulate samples from the above process. Estimate the mean and covariance from the samples. Do your estimates of the mean and covariance agree with their theoretical values?

Putting a standard error on your estimates of the means should be straightforward. You may have to use some creativity to put error bars on your estimates of the covariances.

4. Sampling Gaussians, and matrix decompositions:

[This question has quite a lot of detail on computation and linear algebra, which you can skim over on first reading. You don't actually need to understand the decompositions in detail, or know how to call them in Matlab or Python, to be able to answer some of the questions.]

In lectures we saw that we can sample from a multivariate Gaussian $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ by drawing a vector of standard normals, $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$, and setting $\mathbf{x} = A\mathbf{v}$, for a matrix A where $AA^T = \Sigma$.

The lower-triangular Cholesky decomposition will decompose a symmetric positive-definite² covariance into $\Sigma = LL^T$.

Matlab/Octave: `L = chol(Sigma, 'lower');`

NumPy: `L = np.linalg.cholesky(Sigma)`

This decomposition can be used to draw samples, with $A=L$ above.

A triangular decomposition makes computing most things we might want to know about a covariance quick and easy. Cholesky decompositions are widely used. We can quickly find the determinant: $|L| = \prod_d L_{dd}$ where $|\Sigma| = |L|^2$, or more frequently $\log|L| = \sum_d \log L_{dd}$. We can also solve linear systems³: $L^{-1}\mathbf{b}$ takes similar time to a matrix-vector multiply $L\mathbf{b}$. In Matlab/Octave replace `inv(L)*b` with `L\b` and `inv(Sigma)*b` with `L'\(L\b)`. In Python use `scipy.linalg.solve_triangular`, and `scipy.linalg.cho_solve`.

- a) Sometimes instead of decomposing the covariance matrix, we have the Cholesky decomposition of the precision matrix, $\Sigma^{-1} = CC^T$, where C is lower-triangular. How would we use C to sample from $\mathcal{N}(\mathbf{0}, \Sigma)$?
- b) Real symmetric matrices, like covariance matrices, also have a decomposition of the following form⁴:

$$\Sigma = Q\Lambda Q^T,$$

where Λ is a diagonal matrix of eigenvalues, and the columns of Q are the eigenvectors of Σ .

- i) Describe how to sample from $\mathcal{N}(\mathbf{0}, \Sigma)$ using this decomposition.
- ii) Q is an orthogonal matrix, corresponding to a rigid rotation (and possibly a reflection). Describe geometrically (perhaps in 2D) how your sampling process transforms a cloud of points drawn from a standard normal.
- c) Yet another possible decomposition is the principal square root⁵: $\Sigma = \Sigma^{1/2}\Sigma^{1/2}$,

2. Unfortunately, it doesn't work for semi-definite covariances with zero determinant.

3. There are faster and more numerically stable ways to solve $A^{-1}\mathbf{c}$ than inverting A and multiplying \mathbf{c} by A^{-1} . Matlab's `A\b` and NumPy's `np.linalg.solve(A, c)` do something sensible and should be preferred to `inv(A)*c`. However, specialized methods that use properties of A , or a cached decomposition may be better in some situations.

4. https://en.wikipedia.org/wiki/Eigendecomposition_of_a_matrix#Real_symmetric_matrices

We can use this decomposition even when the determinant is zero.

5. Non-examinable: $\Sigma^{1/2} = Q\Lambda^{1/2}Q^T$, using Q and Λ from the eigendecomposition, and $\Lambda^{1/2}$ simply replaces each eigenvalue on the diagonal with its square root. However, it would be better to compute it with `sqrtm`.

where $\Sigma^{1/2}$ is symmetric. None of the decompositions discussed so far are the same. In this part we'll try to understand how they're related.

- i) Consider two different decompositions $\Sigma = AA^\top = BB^\top$. We'll assume the matrices are full rank so that we can write $B = AU$. Show that $UU^\top = \mathbb{I}$, the identity matrix, which means that U is an orthogonal matrix.
- ii) Explain geometrically why if computing Av from $v \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$ is a way to sample from $\mathcal{N}(\mathbf{0}, \Sigma)$, computing $Bv = AUv$ will be as well.