

ML2R Theory Nuggets

Centering Data- and Kernel Matrices

Christian Bauckhage*
Machine Learning Rhine-Ruhr
Fraunhofer IAIS
St. Augustin, Germany

Pascal Welke†
Machine Learning Rhine-Ruhr
University of Bonn
Bonn, Germany

ABSTRACT

We discuss the notion of centered data matrices and show how to compute them using centering matrices. As centering matrices have many applications in data science and machine learning, we have a look at one such application and discuss how they allow for centering kernel matrices.

1 CENTERING DATA MATRICES

In data mining, machine learning, or pattern recognition, we often have to analyze or process a sample of n data points $\mathbf{x}_i \in \mathbb{R}^m$. This, in turn, frequently requires us to consider *zero mean* or *centered* data.

A simple canonical example of such a setting is the computation of the $m \times m$ (biased) sample covariance matrix

$$C = \frac{1}{n} \sum_{i=1}^n [\mathbf{x}_i - \bar{\mathbf{x}}][\mathbf{x}_i - \bar{\mathbf{x}}]^\top \quad (1)$$

where the vector

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (2)$$

denotes the sample mean.

While this is a specific example, we note its following general aspect: The operation of subtracting the sample mean from each of the given data points

$$\mathbf{z}_i = \mathbf{x}_i - \bar{\mathbf{x}} \quad (3)$$

causes the resulting transformed data \mathbf{z}_i to have the following mean

$$\begin{aligned} \bar{\mathbf{z}} &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i = \frac{1}{n} \sum_{i=1}^n [\mathbf{x}_i - \bar{\mathbf{x}}] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{x}} = \bar{\mathbf{x}} - \frac{1}{n} \cdot n \cdot \bar{\mathbf{x}} = \mathbf{0} \end{aligned} \quad (4)$$

In other words, the transformation in (3) causes the transformed data to be “of zero mean” or to be “centered at $\mathbf{0}$ ”.

With respect to the sample covariance matrix in our introductory example, this means that we could also compute it as

$$C = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^\top \quad (5)$$

What does any of this have to do with data matrices? Well, we may gather the n given data points $\mathbf{x}_i \in \mathbb{R}^m$ in an $m \times n$ data matrix

$$X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \quad (6)$$

and apply linear algebra to facilitate procedures such as the one we just went through. If we did this, we should hope that there is a linear algebraic mechanism to produce a centered data matrix

$$Z = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n] \quad (7)$$

Indeed, linear algebra provides such a mechanism and we next discuss how it works.

In what follows, we let $\mathbf{1} \in \mathbb{R}^n$ denote the n -dimensional vector of all ones.

If n data points have been gathered in a data matrix X , it is an easy exercise to verify that their sample mean in (2) can also be computed like this

$$\bar{\mathbf{x}} = \frac{1}{n} X \mathbf{1} \quad (8)$$

When working with the vector of all ones, we can moreover consider the outer product $\bar{\mathbf{x}} \mathbf{1}^\top$ to obtain an $m \times n$ matrix

$$\bar{\mathbf{x}} \mathbf{1}^\top = [\bar{\mathbf{x}}, \bar{\mathbf{x}}, \dots, \bar{\mathbf{x}}] \quad (9)$$

all of whose columns are copies of $\bar{\mathbf{x}}$.

This matrix now allows us to perform data centering, i.e. to subtract the sample mean from each of the given data points, by means of just a single matrix subtraction, namely

$$Z = X - \bar{\mathbf{x}} \mathbf{1}^\top \quad (10)$$

Given what we saw in (4), it is clear that the column mean of the resulting matrix Z will be $\mathbf{0}$. In other words, Z is a centered version of our original data matrix X . Regarding our introductory example of the sample covariance matrix, this means that we can now compute it as

$$C = \frac{1}{n} Z Z^\top \quad (11)$$

But let us keep going and develop even deeper insights into data centering! If we plug (8) into (10), we obtain the following crucial result

$$Z = X - \frac{1}{n} X \mathbf{1} \mathbf{1}^\top = X \left[I - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right] \equiv X J \quad (12)$$

Here, I denotes the $n \times n$ identity matrix and the matrix

$$J = I - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \quad (13)$$

is called a **centering matrix** because it centers the columns of matrix X at $\mathbf{0}$.

Centering matrices have a couple of important or convenient characteristics. For instance, it is easy to show that

* 0000-0001-6615-2128

† 0000-0002-2123-3781

- (1) J is square, i.e. $J \in \mathbb{R}^{n \times n}$
- (2) J is symmetric, i.e. $J = J^\top$
- (3) J is idempotent, i.e. $J^k = J$ for all $1 \leq k \in \mathbb{N}$
- (4) J is positive semi-definite, i.e. $\mathbf{y}^\top J \mathbf{y} \geq 0$ for all $\mathbf{y} \in \mathbb{R}^n$
- (5) J is singular, i.e. *not* invertible, i.e. J^{-1} does not exist
- (6) because of (4), all eigenvalues of J are non-negative
- (7) because of (3), all eigenvalues of J are either 0 or 1
- (8) because of (2) and (3), J is an orthogonal projection matrix
- (9) the trace of J amounts to

$$\text{tr}[J] = n \cdot \left(1 - \frac{1}{n}\right) = n - 1$$

Regarding our introductory example of computing the sample covariance matrix, we can make use of some of the above properties to obtain

$$C = \frac{1}{n} Z Z^\top \quad (14)$$

$$= \frac{1}{n} [X J] [X J]^\top \quad (15)$$

$$= \frac{1}{n} X J J^\top X^\top \quad (16)$$

$$= \frac{1}{n} X J J X^\top \quad (17)$$

$$= \frac{1}{n} X J X^\top \quad (18)$$

Since this result may look curious at first sight, we will briefly verify it from a different point of view.

If we expand the outer products in the summation in (1), we obtain the following fairly well known decomposition of the sample covariance matrix

$$C = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \frac{2}{n} \sum_{i=1}^n \mathbf{x}_i \bar{\mathbf{x}}^\top + \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \quad (19)$$

$$= \frac{1}{n} X X^\top - 2 \bar{\mathbf{x}} \bar{\mathbf{x}}^\top + \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \quad (20)$$

$$= \frac{1}{n} X X^\top - \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \quad (21)$$

If we next plug (8) into this expression, we realize that the sample covariance matrix can also be written as

$$C = \frac{1}{n} X X^\top - \frac{1}{n^2} [X \mathbf{1}] [X \mathbf{1}]^\top \quad (22)$$

$$= \frac{1}{n} X X^\top - \frac{1}{n^2} X \mathbf{1} \mathbf{1}^\top X^\top \quad (23)$$

$$= \frac{1}{n} X \left[X^\top - \frac{1}{n} \mathbf{1} \mathbf{1}^\top X^\top \right] \quad (24)$$

$$= \frac{1}{n} X \left[I - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right] X^\top \quad (25)$$

$$= \frac{1}{n} X J X^\top \quad (26)$$

which does indeed corroborate the result in (18).

2 CENTERING KERNEL MATRICES

Many algorithms for data mining, machine learning, and pattern recognition work with Gram matrices

$$G = X^\top X \quad (27)$$

which are matrices whose elements are given by $[G]_{ij} = \mathbf{x}_i^\top \mathbf{x}_j$ and thus only depend on inner products between given data points. This allows for invoking the **kernel trick** where we replace inner products by evaluations of a Mercer kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \boldsymbol{\varphi}_i \mid \boldsymbol{\varphi}_j \rangle \quad (28)$$

Here, $\boldsymbol{\varphi}_i \equiv \boldsymbol{\varphi}(\mathbf{x}_i)$ where $\boldsymbol{\varphi} : \mathbb{R}^m \rightarrow \mathbb{H}$ is some generally unknown transformation from the data space into a potentially infinite dimensional feature space. The function $\langle \cdot \mid \cdot \rangle : \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{R}$ denotes the inner product in that feature space.

To stress the whole point of our following discussion, we remark that it sounds formidable to work with unknown transformations into potentially infinite dimensional spaces. Recall, however, that we only do this *implicitly* when using *explicitly* given kernel functions. In other words, while it may be impossible to practically compute the inner product on the right of equation (28), it is usually possible to evaluate the kernel function on its left hand side.

With respect to Gram matrices, this is to say that invoking the kernel trick is to replace the $n \times n$ matrix G by an $n \times n$ kernel matrix K where

$$[K]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) \quad (29)$$

Conceptually, or from the point of view of “pen-and-paper math”, this is the same as introducing a feature space data matrix

$$\Phi = \left[\boldsymbol{\varphi}_1, \boldsymbol{\varphi}_2, \dots, \boldsymbol{\varphi}_n \right] \quad (30)$$

and then replacing G by an $n \times n$ matrix $\Phi^\top \Phi$ where

$$[\Phi^\top \Phi]_{ij} = \langle \boldsymbol{\varphi}_i \mid \boldsymbol{\varphi}_j \rangle \quad (31)$$

Now, if we kernelize an algorithm that requires us to work with centered data $\mathbf{z}_i = \mathbf{x}_i - \bar{\mathbf{x}} \in \mathbb{R}^m$, we must make sure that their feature space representations $\boldsymbol{\psi}_i = \boldsymbol{\varphi}_i - \bar{\boldsymbol{\varphi}} \in \mathbb{H}$ are centered, too.

Again, this sounds formidable, because we just emphasized that it may be practically impossible to compute the underlying transformations. How are we then supposed to compute the feature space mean required for feature space centering? But let us do some “pen-and-paper math” to see what we are actually dealing with.

On a conceptual or symbolic level, we can easily work with the feature space data matrix Φ in (30). Our insights from the previous section then tell us that

$$\Psi = \Phi J = \Phi - \frac{1}{n} \Phi \mathbf{1} \mathbf{1}^\top \quad (32)$$

will be a centered matrix. However, if we cannot even write computer code to compute Φ , we certainly cannot write computer code to center it. Yet, what we are actually interested in when kernelizing an algorithm is to compute the Gramian $\Psi^\top \Psi$ and we will see next that this is practically feasible.

In fact, this is easy to see: Using the definition in (32) together with general properties of the centering matrix J , we have

$$\Psi^\top \Psi = [\Phi J]^\top [\Phi J] \quad (33)$$

$$= J^\top \Phi^\top \Phi J \quad (34)$$

$$= J \Phi^\top \Phi J \quad (35)$$

At this point, we note that it is actually possible to compute the entries of the feature space Gramian $\Phi^\top \Phi$, because

$$[\Phi^\top \Phi]_{ij} = \langle \varphi_i | \varphi_j \rangle = K(\mathbf{x}_i, \mathbf{x}_j) = [K]_{ij} \quad (36)$$

In other words, we may not be able to practically compute matrix Φ but we are able to compute $\Phi^\top \Phi$ because its entries can be obtained by evaluating a computable kernel function. All in all, we therefore have the crucial result that

$$\Psi^\top \Psi = J K J \quad (37)$$

At this point, we have successfully shown that it is feasible to practically compute (and therefore work with) matrix $\Psi^\top \Psi$. **But let us once again keep going and develop deeper insights into what it means to center a kernel matrix.**

To begin with, we emphasize, that we have an equation that relates the kernel matrix K to the feature space data matrix Φ , namely

$$K = \Phi^\top \Phi \quad (38)$$

Given that we can compute $\Psi^\top \Psi$ for the centered feature space data matrix Ψ , we now define

$$K_c = \Psi^\top \Psi \quad (39)$$

and henceforth call K_c the *centered* kernel matrix.

Concerning the centered kernel matrix, we observe the following

$$K_c = J K J \quad (40)$$

$$= \left[I - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right] K \left[I - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right] \quad (41)$$

$$= \left[K - \frac{1}{n} \mathbf{1} \mathbf{1}^\top K \right] \left[I - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right] \quad (42)$$

$$= K - \frac{1}{n} \mathbf{1} \mathbf{1}^\top K - \frac{1}{n} K \mathbf{1} \mathbf{1}^\top + \frac{1}{n^2} \mathbf{1} \mathbf{1}^\top K \mathbf{1} \mathbf{1}^\top \quad (43)$$

While this looks interesting, we may again wonder if there is an intuition behind this expression?

To develop a better understanding of the nature of matrix K_c , we next study its individual entries $[K_c]_{ij}$.

Just as we can conceptually work with the feature space data matrix Φ , we can conceptually work with the feature space sample mean

$$\bar{\varphi} = \frac{1}{n} \Phi \mathbf{1} = \frac{1}{n} \sum_{i=1}^n \varphi_i \quad (44)$$

Using this feature space sample mean, we can write the individual entries of K_c as

$$[K_c]_{ij} = [\Psi^\top \Psi]_{ij} = \langle \psi_i | \psi_j \rangle = \langle \varphi_i - \bar{\varphi} | \varphi_j - \bar{\varphi} \rangle \quad (45)$$

And, if we expand these inner products between two centered feature space points, we obtain

$$\begin{aligned} \langle \varphi_i - \bar{\varphi} | \varphi_j - \bar{\varphi} \rangle &= \langle \varphi_i | \varphi_j \rangle - \langle \varphi_i | \bar{\varphi} \rangle \\ &\quad - \langle \bar{\varphi} | \varphi_j \rangle + \langle \bar{\varphi} | \bar{\varphi} \rangle \end{aligned} \quad (46)$$

Using the definition of the feature space sample mean in (44) and the bilinearity of the inner product, this can also be written as

$$\begin{aligned} \langle \varphi_i - \bar{\varphi} | \varphi_j - \bar{\varphi} \rangle &= \langle \varphi_i | \varphi_j \rangle \\ &\quad - \frac{1}{n} \sum_{j=1}^n \langle \varphi_i | \varphi_j \rangle \\ &\quad - \frac{1}{n} \sum_{i=1}^n \langle \varphi_i | \varphi_j \rangle \\ &\quad + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \langle \varphi_i | \varphi_j \rangle \end{aligned} \quad (47)$$

Given this expansion, we can now resort to (28) and write all the inner products on the right hand side of (47) in terms of kernel functions

$$\begin{aligned} \langle \varphi_i - \bar{\varphi} | \varphi_j - \bar{\varphi} \rangle &= K(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad - \frac{1}{n} \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad - \frac{1}{n} \sum_{i=1}^n K(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) \end{aligned} \quad (48)$$

At this point, we observe that all terms on the right of (48) can be expressed in terms of our original kernel matrix K . In particular, we have

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{e}_i^\top K \mathbf{e}_j \quad (49)$$

$$\sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{e}_i^\top K \mathbf{1} \quad (50)$$

$$\sum_{i=1}^n K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{1}^\top K \mathbf{e}_j \quad (51)$$

$$\sum_{i=1}^n \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{1}^\top K \mathbf{1} \quad (52)$$

where $\mathbf{e}_i, \mathbf{e}_j \in \mathbb{R}^n$ denote the i -th and j -th standard basis vector. But this is to say that we can express the (i, j) entry of the centered kernel matrix as

$$[K_c]_{ij} = \mathbf{e}_i^\top K \mathbf{e}_j - \frac{1}{n} \mathbf{e}_i^\top K \mathbf{1} - \frac{1}{n} \mathbf{1}^\top K \mathbf{e}_j + \frac{1}{n^2} \mathbf{1}^\top K \mathbf{1} \quad (53)$$

Since this is beginning to resemble the expression in (43), we further note the following: The expression in (49) is but the (i, j) entry of K . The expression (50) is the i -th component of a column vector, namely $K \mathbf{1}$. The expression in (51) is the j -th component of a row vector, namely $\mathbf{1}^\top K$. And $\mathbf{1}^\top K \mathbf{1}$ is a scalar. Turning these latter objects into $n \times n$ matrices can be accomplished like this: $K \mathbf{1} \mathbf{1}^\top$, $\mathbf{1} \mathbf{1}^\top K$, and $\mathbf{1}^\top K \mathbf{1} \mathbf{1}^\top = \mathbf{1} \mathbf{1}^\top K \mathbf{1} \mathbf{1}^\top$. Using these matrices, we therefore have

$$K_c = K - \frac{1}{n} K \mathbf{1} \mathbf{1}^\top - \frac{1}{n} \mathbf{1} \mathbf{1}^\top K + \frac{1}{n^2} \mathbf{1} \mathbf{1}^\top K \mathbf{1} \mathbf{1}^\top \quad (54)$$

which corresponds exactly to the result we found in equation (43).

3 SUMMARY OF MAIN RESULTS

In data mining, pattern recognition, and machine learning, we often work with data matrices

$$X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$$

whose n columns vectors \mathbf{x}_i have a mean of $\bar{\mathbf{x}}$. For many tasks, however, we need to turn such data matrices into centered data matrices

$$Z = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n]$$

whose n columns vectors $\mathbf{z}_i = \mathbf{x}_i - \bar{\mathbf{x}}$ have a mean of $\mathbf{0}$. Considering the centering matrix

$$J = I - \frac{1}{n} \mathbf{1}\mathbf{1}^\top$$

we can compute the centered version Z of data matrix X as follows

$$Z = XJ$$

The “outer product” of a centered matrix $Z = XJ$ with itself can be computed in terms of the original matrix X , namely

$$\begin{aligned} ZZ^\top &= [XJ][XJ]^\top \\ &= XJJ^\top X^\top \\ &= XJJX^\top \\ &= XJX^\top \end{aligned}$$

This uses the fact that the centering matrix J is symmetric and idempotent.

The “inner product” of a centered matrix $Z = XJ$ with itself can also be computed in terms of the original matrix X , namely

$$\begin{aligned} Z^\top Z &= [XJ]^\top [XJ] \\ &= J^\top X^\top XJ \\ &= JX^\top XJ \end{aligned}$$

This again uses the symmetry of the centering matrix J .

When working with algorithms that involve computations with Gram matrices $X^\top X$, we can invoke the kernel trick and replace the Gramian by a kernel matrix K . The corresponding centered kernel matrix is given by

$$K_c = JKJ$$

4 A SHORT NOTE ON IMPLEMENTATION

As it is customary for this series to provide simple *NumPy* implementations of the proposed methods, we will do just that. A centering matrix for n data points can be obtained with

```
import numpy as np
def get_centering_matrix(n):
    return np.eye(n) - 1. / n * np.ones([n,1])
```

While this seems like a straightforward implementation of (13) at first glance, we note one slight difference: *NumPy* allows for so called *broadcasting* of the vector `np.ones([1, n])` implementing 1 to `np.eye(n)` implementing the $n \times n$ matrix I if the dimensions are compatible. This allows to implement (13) without explicitly creating an $n \times n$ matrix filled with ones, which becomes very useful for larger problem instances.

5 OUTLOOK TO APPLICATION EXAMPLES

As we saw in our introductory example, centered data matrices play a crucial role in computing sample covariance matrices. They are thus important for methods such as PCA [7] or CCA [10] which build on top of covariance matrices. But centered data matrices occur in other contexts, too. For instance, they may form the basis of simple yet efficient clustering algorithms [1, 2].

Centered Gram- or kernel matrices play crucial roles in multidimensional scaling [8] or kernel PCA [9]. Especially the latter allows for interesting approaches to text analysis [3–6] which we discuss in our ML2R Coding Nuggets series.

ACKNOWLEDGMENTS

This material was produced within the Competence Center for Machine Learning Rhine-Ruhr (**ML2R**) which is funded by the Federal Ministry of Education and Research of Germany (grant no. 01IS18038C). The authors gratefully acknowledge this support.

REFERENCES

- [1] C. Bauckhage, E. Brito, K. Cvejosi, C. Ojeda, R. Sifa, and S. Wrobel. 2017. Ising Models for Binary Clustering via Adiabatic Quantum Computing. In *Proc. EMM-CVPR*. Springer.
- [2] C. Bauckhage, C. Ojeda, R. Sifa, and S. Wrobel. 2018. Adiabatic Quantum Computing for Kernel k=2 Means Clustering. In *Proc. KDML-LWDA*.
- [3] M. Beekma, M. van Gompel, F. Kunneman, L. Onrust, B. Regnerus, D. Vinke, E. Brito, C. Bauckhage, and R. Sifa. 2018. Detecting and Correcting Spelling Errors in High-quality Dutch Wikipedia Text. *Computational Linguistics in the Netherlands J.* 8 (2018), 122–137.
- [4] E. Brito, B. Georgiev, D. Domingo-Fernandez, C.T. Hoyt, and C. Bauckhage. 2019. RatVec: A General Approach for Low-dimensional Distributed Vector Representations via Rational Kernels. In *Proc. KDML-LWDA*.
- [5] E. Brito, R. Sifa, and C. Bauckhage. 2017. KPCA Embeddings: An Unsupervised Approach to Learn Vector Representations of Finite Domain Sequences. In *Proc. KDML-LWDA*.
- [6] V. Gupta, S. Giesselbach, S. Rüping, and C. Bauckhage. 2019. Improving Word Embeddings Using Kernel PCA. In *Proc. Workshop on Representation Learning for NLP @ ACL*.
- [7] I.T. Jolliffe. 1986. *Principal Component Analysis*. Springer.
- [8] J.B. Kruskal and M. Wish. 1978. *Multidimensional Scaling*. SAGE Publications.
- [9] B. Schölkopf, A. Smola, and K.-R. Müller. 1998. Nonlinear Component Analysis as a Kernel Eigenvalue Problem. *Neural Computation* 10 (1998).
- [10] B. Thompson. 1985. *Canonical Correlation Analysis*. SAGE Publications.