

Lecture Notes on Vector Quantization

Laurenz Wiskott
Institut für Neuroinformatik
Ruhr-Universität Bochum, Germany, EU

21 February 2013

Contents

1	General idea	2
2	Error function	2
3	Optimal assignment given the reference vectors	2
3.1	Voronoi tessellation	2
3.2	Delaunay triangulation	3
4	Optimal reference vectors given an assignment	3
5	A simple vector quantization algorithm	3
6	An online learning rule for vector quantization	3
7	Learning rates	4
7.1	Constant learning rate	4
7.2	1/ t -decaying learning rate	5
7.3	Exponentially decaying learning rate	5
8	Soft competitive learning	5
9	Growing networks +	6
10	Towards a neural implementation +	6
10.1	Euclidian distance and inner product	6
10.2	Winner takes all	6
10.3	Adaptive resonance theory	7

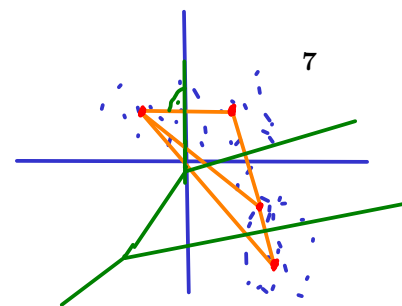
• 2007,2013 Laurenz Wiskott. This work (except for all figures from other sources, if present) is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

Homepage is <https://www.ini.rub.de/PEOPLE/wiskott/>.

More teaching material is available at <https://www.ini.rub.de/PEOPLE/wiskott/Teaching/Material/>.

11 Demo

These lecture notes are largely based on (Hertz et al., 1991; Rojas, 1996).



1 General idea

The purpose of vector quantization is to compress vectorial data.¹ The general idea is to find a good set of reference vectors and then replace each data vector simply by the index of its best reference vector. If there are much fewer reference vectors than data vectors, it is much cheaper to store or transmit just the reference vectors and the index for each data vector than the data vectors themselves. Finding the optimal vector quantization has thus two subproblems: (i) finding an optimal set of reference vectors and (ii) finding the optimal index for each data vector.

2 Error function

The performance of the vector quantization is typically measured by the quantization error, which is the mean Euclidean distance between the data vectors and their reference vectors. Let \mathbf{x}^μ indicate the data vectors, with $\mu = 1, \dots, M$, and \mathbf{r}_i indicate the reference vectors, with $i = 1, \dots, I$. The index associated with each data vector is $i^* = i^*(\mu)$, where we refer to $i^*(\mu)$ as the assignment function. The quantization error is then given by

$$\bullet E := \langle \|\mathbf{x}_\mu - \mathbf{r}_{i^*}\|^2 \rangle_\mu, \quad (1)$$

where $\langle \cdot \rangle$ indicates an averaging. Note that i^* is a function of μ .

3 Optimal assignment given the reference vectors

Assume the data vectors and reference vectors are given. What is the optimal assignment of the data vectors to the reference vectors, i.e. what is the optimal assignment function $i^*(\mu)$? It is pretty obvious and follows directly from equation (1) that each data vector should be assigned to its nearest reference vector, i.e.

$$\circ i^*(\mu) := \arg \min_i \|\mathbf{x}_\mu - \mathbf{r}_i\|^2 \quad (2)$$

is optimal.

3.1 Voronoi tessellation

Instead of deciding for each data point to which reference vector it belongs, one can also partition the space such that all points within one partition get assigned to one reference vector. The optimal reference vector then becomes a function of space and not of data vector, but otherwise it is the same, i.e.

$$i^*(\mathbf{x}) := \arg \min_i \|\mathbf{x} - \mathbf{r}_i\|^2 \quad (3)$$

is optimal. Such a partitioning is called Voronoi (or Dirichlet) tessellation (D: Voronoi-Diagramm, Dirichlet-Zerlegung) and the region closest to one reference vector is called Voronoi cell (D: Voronoi Region).

¹Important text (but not inline formulas) is set bold face; \bullet marks important formulas worth remembering; \circ marks less important formulas, which I also discuss in the lecture; + marks sections that I typically skip during my lectures.

3.2 Delaunay triangulation

If one connects all pairs of reference vectors whose Voronoi regions share a common edge, one gets a **Delaunay triangulation** (D: Delaunay-Triangulation).

A definition independent of Voronoi regions reads as follows: “For a set P of points in the (d -dimensional) Euclidean space, a Delaunay triangulation is a triangulation $DT(P)$ such that no point in P is inside the circum-hypersphere of any simplex in $DT(P)$. It is known that there exists a unique Delaunay triangulation for P , if P is a set of points in general position; that is, there exists no k -flat containing $k + 2$ points nor a k -sphere containing $k + 3$ points, for $1 \leq k \leq d - 1$ (e.g., for a set of points in \mathbb{R}^3 ; no three points are on a line, no four on a plane, no four are on a circle, and no five on a sphere).” (Wikipedia, 2013)

4 Optimal reference vectors given an assignment

In general it is a hard problem to determine **the optimal reference vectors**. However, if a **fixed assignment is given**, it is easy to determine them. They simply **lie in the center of gravity of the data vectors assigned to them** as can be seen as follows. Consider just one reference i and let \mathcal{M}_i be the set of data vectors assigned to it. The quantization error then simplifies to

$$\circ E_i := \left\langle \|\mathbf{x}_\mu - \mathbf{r}_i\|^2 \right\rangle_{\mu \in \mathcal{M}_i}. \quad (4)$$

A condition for the optimal choice of the reference \mathbf{r}_i is that the gradient of E_i with respect to \mathbf{r}_i vanishes.

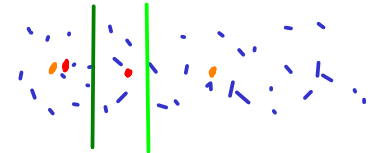
$$\circ 0 \stackrel{!}{=} \text{grad } E_i \stackrel{(4)}{=} -2 \langle \mathbf{x}_\mu - \mathbf{r}_i \rangle_{\mu \in \mathcal{M}_i} \quad (5)$$

$$\circ = 2 \langle \mathbf{r}_i - \langle \mathbf{x}_\mu \rangle_{\mu \in \mathcal{M}_i} \rangle \quad (6)$$

$$\circ \iff \mathbf{r}_i = \langle \mathbf{x}_\mu \rangle_{\mu \in \mathcal{M}_i}. \quad (7)$$

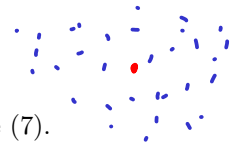
Since E_i is a quadratic error function, this is indeed the optimum.

5 A simple vector quantization algorithm



Now that we have a criterion for the optimal assignment given the reference vectors and for the optimal reference vectors given a fixed assignment, we can formulate a **simple iterative algorithm for vector quantization**. This algorithm is known as the LBG or generalized Lloyd algorithm.

1. **Initialize the reference vectors** to different data vectors.
2. **Find the assignment** of the data vectors to their nearest reference vectors.
3. **Move each reference vector to the center of gravity** of the data assigned to it, see (7).
4. **Repeat** at step 2 **until convergence**.



This algorithm is simple and converges quickly. **Unfortunately**, it is not guaranteed to converge to the optimal solution but **gets easily stuck in a local optimum**.

6 An online learning rule for vector quantization

If it is not possible to process all data simultaneously, e.g. if there is too much data, or if one prefers to process data one by one for greater biological plausibility, **one can use an online version of the algorithm above**. **Instead of moving a given reference vector directly to the center of gravity of the data vector within its Voronoi region, it gets moved by a certain fraction towards a data vector** whenever it is assigned to the reference vector. This procedure is approximately equivalent in the long run, if the fraction, determined by the so-called *learning rate* η , is small enough so that the reference vectors do not jump about too much. The algorithm now reads as follows:

1. **Initialize the reference vectors** to different data vectors.
2. **Pick a data vector \mathbf{x}_μ at random.**
3. **Determine the optimal reference vector \mathbf{r}_{i^*} ,** see (2).
4. **Move** the optimal reference vector **a little bit towards the data vector:**

$$\Delta \mathbf{r}_{i^*} = \eta(\mathbf{x}_\mu - \mathbf{r}_{i^*}) \quad (8)$$

5. **Repeat** at 2 **until** a maximum number of steps is reached or any other measure indicates **convergence.**



7 Learning rates

With a constant learning rate in the online learning rule the learning is typically too slow in the initial phase and at the end the algorithm does not converge, because the reference vectors get changed indefinitely. Thus, **it might be good to consider decreasing the learning rate over time.** However, first consider the effect of a constant learning rate.

7.1 Constant learning rate

Consider just one reference vector \mathbf{r} and let $\mathbf{x}(t)$ be a sequence of data vectors used to update \mathbf{r} . **With a constant learning rate** we have

$$\circ \quad \mathbf{r}(0) = \mathbf{r}_0, \quad (9)$$

$$\circ \quad \mathbf{r}(t) \stackrel{(8)}{=} \mathbf{r}(t-1) + \eta(\mathbf{x}(t) - \mathbf{r}(t-1)) \quad (10)$$

$$\circ \quad = (1-\eta)\mathbf{r}(t-1) + \eta\mathbf{x}(t). \quad (11)$$

By mathematical induction **we find that** for $t > 0$

$$\circ \quad \mathbf{r}(t) = (1-\eta)^t \mathbf{r}_0 + \eta \sum_{t'=1}^t (1-\eta)^{t-t'} \mathbf{x}(t'). \quad (12)$$

Proof:

$$\text{Base } (t=1): \quad \mathbf{r}(1) \stackrel{(11,9)}{=} (1-\eta)\mathbf{r}_0 + \eta\mathbf{x}(1) \quad (13)$$

$$= (1-\eta)^1 \mathbf{r}_0 + \eta \sum_{t'=1}^1 \underbrace{(1-\eta)^{1-t'}}_{=1} \mathbf{x}(t') \quad (14)$$

$$= \text{rhs of (12) for } t=1. \quad (15)$$

$$(16)$$

$$\text{Induction step } (t > 1): \quad \mathbf{r}(t) \stackrel{(11)}{=} (1-\eta)\mathbf{r}(t-1) + \eta\mathbf{x}(t) \quad (17)$$

$$\stackrel{(12)}{=} (1-\eta) \left((1-\eta)^{(t-1)} \mathbf{r}_0 + \eta \sum_{t'=1}^{t-1} (1-\eta)^{t-1-t'} \mathbf{x}(t') \right) + \eta\mathbf{x}(t) \quad (18)$$

$$= (1-\eta)^t \mathbf{r}_0 + \eta \sum_{t'=1}^{t-1} (1-\eta)^{t-t'} \mathbf{x}(t') + \eta(1-\eta)^{t-t} \mathbf{x}(t) \quad (19)$$

$$= (1-\eta)^t \mathbf{r}_0 + \eta \sum_{t'=1}^t (1-\eta)^{t-t'} \mathbf{x}(t') \quad (20)$$

$$= \text{rhs of (12) for } t > 1. \quad (21)$$

Thus, **the influence of data vectors** picked for training **decays exponentially into the past.** As a consequence **the system never converges** but it remains adaptive in case the data distribution changes. If the data distribution does not change, the reference vectors first move towards an equilibrium distribution and then wiggle around forever.

7.2 $1/t$ -decaying learning rate

To avoid the problems connected with a constant learning rate one can choose a decaying learning rate. One interesting scheme is to introduce a learning rate for each reference vector individually and let it decay like a harmonic sequence

$$\bullet \quad \eta(t) = \frac{1}{t}, \quad (22)$$

where t again indexes the training event for the one particular reference vector considered (one could also write $\eta_i(t_i)$ so that it is more explicit that each reference vector has its own learning rate and its own counter for how often it has been modified).

With such a decaying learning rate we can find by mathematical induction that

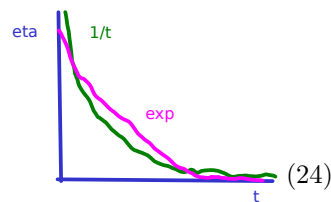
$$\bullet \quad \mathbf{r}(t) = \frac{1}{t} \sum_{t'=1}^t \mathbf{x}(t'). \quad (23)$$

Thus, with such a learning rate the reference vectors lie exactly in the center of gravity of all data vectors used to train it. However, some of the data vectors used at some time might not lie in the Voronoi region of the reference vector anymore, because the reference vectors have moved around. Also, if the data vectors are picked at random, some might have been picked multiple times while others have not been picked at all. Thus, online learning is not equivalent to the batch algorithm (LBG) considered above. Finally, this online learning does not even converge in a strict sense because the harmonic series diverges (although this is more of a theoretical issue, because the harmonic series diverges extremely slowly, see <http://www.mathematik.com/Harmonic/index.html>).

7.3 Exponentially decaying learning rate

Another proposal for a learning rate is

$$\bullet \quad \eta(t) = \eta_{\text{init}} (\eta_{\text{final}} / \eta_{\text{init}})^{t/t_{\text{final}}}, \quad (24)$$



which is exponentially decaying. η_{init} indicates the initial learning rate, η_{final} the final learning rate, and t_{final} the total number of steps used. Note that now there is just one common learning rate for all reference vectors.

Depending on the parameters the learning rate tends to be larger in the beginning and smaller at the end compared to the $1/t$ -decay. This has several advantages. Firstly, learning is more aggressive in the beginning (reference vectors jump around more), which can be interpreted as additional noise and helps the system to escape a local optimum it might have been initialized to. Secondly, data vectors used early on lose their influence rapidly so that in the end only data vectors within the final Voronoi region have significant influence on the position of the reference vector. Finally, since the learning rate is smaller at the end, the learning has better convergence. It has also been shown experimentally that the exponentially decaying learning rate usually leads to better results than $1/t$ (except if there is only one reference vector, in which case $1/t$ is optimal).

8 Soft competitive learning

So far only the nearest reference vector got updated, which is referred to as hard competitive learning. Some algorithms not only update the nearest reference vector but also the second and even more neighbors, but to a lesser degree the farther away the reference vectors are from the data vector, although it is actually not the distance that matters but the rank in the list of neighbors. This can be formalized as follows.

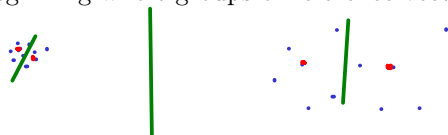
1. Initialize the reference vectors to different data vectors.
2. Pick a data vector \mathbf{x}_μ at random.

3. **Determine the nearest reference vector, the second nearest, etc.** indicated by i^*_1, i^*_2, \dots etc.
4. **Move** the reference vectors **a little bit towards the data vector depending on the rank in the nearest reference vector list**
 - o $\Delta \mathbf{r}_{i^*_k} = \eta \exp(-k/\lambda)(\mathbf{x}_\mu - \mathbf{r}_{i^*_k})$ (25)
5. **Repeat** at 2 **until** a maximum number of steps is reached or any other measure indicates **convergence**.

The factor $\exp(-k/\lambda)$ weights the adaptation according to the rank k of the reference vectors. This soft competitive learning can be combined with an exponential decay of the learning rate as well as an exponential decay of the scale λ of the neighborhood function $\exp(-k/\lambda)$.

The advantage of soft competitive learning, I guess, is that **local optima are avoided to some extent** and that **learning can be faster**, since in the beginning whole groups of reference vectors are updated.

9 Growing networks +



Bernd Fritzke has developed a growing neural gas algorithm, in which additional reference vectors are added at locations of large errors, so that the error can be reduced in a very targeted way (Fritzke, 1994).

10 Towards a neural implementation +

So far we have considered vector quantization only as an algorithmic problem. If we want to use vector quantization for computational neuroscience we have to consider also how such an algorithm could be implemented in a neural system.

10.1 Euclidian distance and inner product

In the algorithmic version the Euclidian norm is used as a measure of distance between data and reference vectors. In neural systems a more natural choice is the inner product between the input vector and the weight vector. These two measures are quite different, however, if data/input vectors as well as reference/weight vectors are normalized to length 1, Euclidian distance and inner product are closely related.

$$\begin{aligned} \overset{(\mathbf{x}-\mathbf{r})^T(\mathbf{x}-\mathbf{r})}{|\mathbf{x} - \mathbf{r}|^2} &= |\mathbf{x}|^2 - 2\mathbf{x} \cdot \mathbf{r} + |\mathbf{r}|^2 & (26) \\ &= 2(1 - \mathbf{x} \cdot \mathbf{r}). & (27) \end{aligned}$$

(since $|\mathbf{x}| = |\mathbf{r}| = 1$ by assumption)

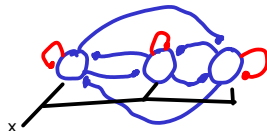
This also implies

$$\begin{aligned} \mathbf{x} \cdot \mathbf{r}_i &> \mathbf{x} \cdot \mathbf{r}_j & (28) \\ \iff |\mathbf{x} - \mathbf{r}_i| &< |\mathbf{x} - \mathbf{r}_j|. & (29) \end{aligned}$$

Thus, if input and weight vectors are normalized to one, the most excited unit is always the one with the least Euclidean distance of its weights vector to the input vector.

Even though there is a strong relationship between Euclidian distance and inner product if input and weight vectors are normalized to one, the normalization forces the data to lie on the surface of a hypersphere, which is a special topology.

10.2 Winner takes all



The vector quantization algorithm requires to determine the nearest reference vector. In a neural system this can be achieved by a winner-take-all mechanism, which can be implemented by mutual inhibition and self-excitation.

There exist also k -winner-take-all versions of this, where the first few winners get activated. It is difficult, though, to implement a decaying activation of the first winners. Soft-max might be an option here.

10.3 Adaptive resonance theory

Adaptive resonance theory (ART) is a detailed model of vector quantization, which employs new units as qualitatively new input data arrives. It is therefore similar to the growing neural gas algorithm.

11 Demo

See <http://www.demogng.de/js/demogng.html>.

References

- Fritzke, B. (1994). Growing cell structures - a self-organizing network for unsupervised and supervised learning. *Neural Networks*, 7(9):1441–1460.
- Hertz, J., Krogh, A., and Palmer, R. G. (1991). *Introduction to the Theory of Neural Computation*. Addison-Wesley, Redwood City, CA.
- Rojas, R. (1996). *Theorie der neuronalen Netze: Eine systematische Einführung*. Springer Verlag, Berlin. 4., korregierter Nachdruck.
- Wikipedia (2013). Delaunay triangulation — wikipedia, the free encyclopedia. [Online; accessed 21-February-2013].