

Derivation of a Class of Training Algorithms *

S. P. LUTTRELL

Royal Signals and Radar Establishment, Malvern, Worcs., WR14 3PS, U.K.

This paper presents a novel derivation of Kohonen's topographic mapping training algorithm, based upon an extension of the Linde-Buzo-Gray (LBG) algorithm for vector quantiser design. Thus a vector quantiser is designed by minimising an L_2 reconstruction distortion measure, including an additional contribution from the effect of code noise which corrupts the output of the vector quantiser. The neighbourhood updating scheme of Kohonen's topographic mapping training algorithm emerges as a special case of this code noise model. This formulation of Kohonen's algorithm is a specific instance of the "robust hidden layer principle", which stabilises the internal representations chosen by a network against anticipated noise or distortion processes.

I. INTRODUCTION

Vector quantisation theory is a generalisation of scalar quantisation theory as expressed by the Lloyd-Max equations [1, 2]. The generalisation of the Lloyd-Max equations to vector quantisation is straightforward, and the Linde-Buzo-Gray (LBG), or k means, algorithm [3] provides a means of adjusting the code vectors to locate a local minimum of a distortion measure.

In this paper, the L_2 distortion measure is extended to include the effect of corruption of the vector quantiser output code. Kohonen's topographic mapping training algorithm then emerges naturally. Some extensions of the technique are indicated.

II. THE LBG ALGORITHM

Although it is not strictly necessary, we shall restrict our attention to an L_2 distortion measure $d(\mathbf{x}, \mathbf{x}')$ defined as the following, which measures the Euclidean distance between a vector \mathbf{x} and its reconstruction after vector quantisation \mathbf{x}' :

$$d(\mathbf{x}, \mathbf{x}') \equiv \|\mathbf{x}' - \mathbf{x}\|^2 \quad (2.1)$$

We shall be concerned with the average of $d(\mathbf{x}, \mathbf{x}')$ over a training set of vectors \mathbf{x} , so we shall introduce the probability density function (PDF) $P(\mathbf{x})$ over samples \mathbf{x} selected at random from the training set.

Denote the encoding operation as $y(\mathbf{x})$, and the corresponding decoding operation as $\mathbf{x}'(y)$. The average L_2 distortion D is the average of $d(\mathbf{x}, \mathbf{x}'(y(\mathbf{x})))$ over \mathbf{x} sampled from $P(\mathbf{x})$

$$D \equiv \int d\mathbf{x} P(\mathbf{x}) \|\mathbf{x}'(y(\mathbf{x})) - \mathbf{x}\|^2 \quad (2.2)$$

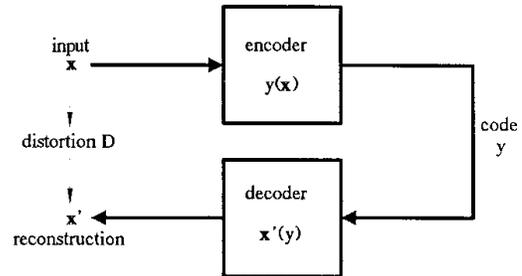


Figure 1: Representation of Equation 2.2 describing the L_2 distortion D after encoding using $y(\mathbf{x})$ and then decoding using $\mathbf{x}'(y)$. This is a single stage encoder.

Equation 2.2 is depicted schematically in Figure 1, where $\mathbf{x}'(y)$ acts as a pseudoinverse for $y(\mathbf{x})$. The optimum encoding/decoding scheme is found by appropriately varying the functions $y(\mathbf{x})$ and $\mathbf{x}'(y)$ so as to minimise D . Necessary conditions for minimisation of D are [3]:

A1) Given \mathbf{x} , the appropriate code $y = y(\mathbf{x})$ minimises the distortion $\|\mathbf{x}'(y(\mathbf{x})) - \mathbf{x}\|^2$.

A2) Given y , the appropriate reconstruction $\mathbf{x}' = \mathbf{x}'(y)$ is the centroid of those \mathbf{x} 's which satisfy $y = y(\mathbf{x})$.

Note that A1 is a nearest neighbour encoding rule. A1 and A2 imply that D is stationary (and locally minimum) with respect to variations of $y(\mathbf{x})$ and $\mathbf{x}'(y)$, respectively. Such conditions are therefore sufficient for a local minimum of D , but are only necessary for a global minimum. The LBG algorithm runs in batch training mode where the whole training set is presented before performing an update, which consists of alternately adjusting $y(\mathbf{x})$ according to A1, and then adjusting $\mathbf{x}'(y)$ according to A2.

III. TWO STAGE QUANTISATION

Now consider a scheme in which the encoding operation is broken down into two stages, $\mathbf{h}(\mathbf{x})$ followed by $y(\mathbf{h})$, where \mathbf{h} is an intermediate code. The L_2 distortion

*This paper appeared in IEEE Trans. Neural Networks, 1990, vol. 1, no. 2, pp. 229-232. Manuscript received June 2, 1989; revised November 27, 1989. An earlier version of this paper was presented at the 1989 International Conference on Neural Networks, Washington, DC, June 19-22, 1989.

is then given by

$$D \equiv \int d\mathbf{x} P(\mathbf{x}) \|\mathbf{x}'(y(\mathbf{h}(\mathbf{x}))) - \mathbf{x}\|^2 \quad (3.1)$$

Equation 3.1 is more clumsy than Equation 2.2 because there are many ways of achieving the same overall encoding operation 7 by delicately balancing the functional forms of $\mathbf{h}(\mathbf{x})$ and $y(\mathbf{h})$. In general, the functional forms of $\mathbf{h}(\mathbf{x})$ and $y(\mathbf{h})$ are not simple, although $y(\mathbf{h}(\mathbf{x}))$ itself is a vector quantiser.

The essential step is to model the *average* distortion due to the second stage of encoding $y(\mathbf{h})$ as a noise process acting on the output of the first stage $\mathbf{h}(\mathbf{x})$. Thus, consider a modified distortion defined as

$$D_1 \equiv \int d\mathbf{x} P(\mathbf{x}) \int d\mathbf{n} \pi(\mathbf{n}) \|\mathbf{x}'(\mathbf{h}(\mathbf{x}) + \mathbf{n}) - \mathbf{x}\|^2 \quad (3.2)$$

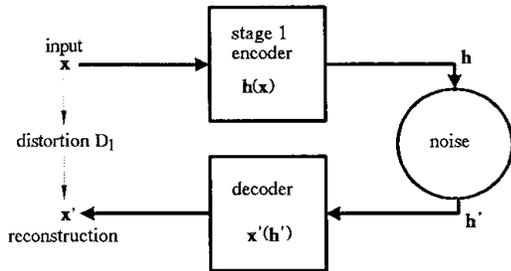


Figure 2: Representation of Equation 3.2 describing the L_2 distortion D_1 , after encoding using $\mathbf{h}(\mathbf{x})$ and modelling the distorting effect of $y(\mathbf{h})$ as a noise process that corrupts \mathbf{h} into \mathbf{h}' , and then decoding using $\mathbf{x}'(\mathbf{h}')$.

Equation 3.2 is depicted in Figure 2. This is the L_2 distortion which would arise from encoding \mathbf{x} using $\mathbf{h}(\mathbf{x})$, followed by addition of noise \mathbf{n} with PDF $\pi(\mathbf{n})$ (to produce \mathbf{h}'), followed by decoding using $\mathbf{x}'(\mathbf{h}')$. By a suitable choice of $\pi(\mathbf{n})$ we can model the average distortion due to $y(\mathbf{h})$. Note that we use a somewhat cavalier notation by using $\mathbf{x}'(\mathbf{h}')$ here and $\mathbf{x}'(y)$ elsewhere: this does not imply that the functions are the same.

The use of a modified distortion measure (D_1 , rather than D) is a specific instance of the application of the ‘‘robust hidden layer principle’’. This principle states that a network should be trained in the presence of all anticipated noise and distortion processes acting on its internal representations. The effect of this is to encourage the network to encode information in its internal representations in such a way that it is robust with respect to the damaging effects of noise and distortion. In this case we use $\pi(\mathbf{n})$ to model stochastically the effects of the distortion process.

The noise model that we have used in Equation 3.2 is additive. We could use more sophisticated noise models to take into account the different distortions introduced by $y(\mathbf{h})$ for different \mathbf{h} , in which case $\pi(\mathbf{n})$ would become

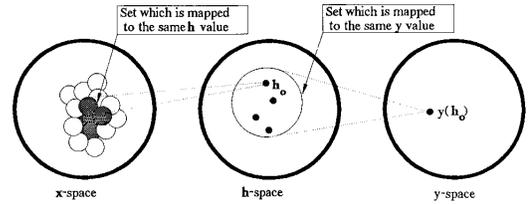


Figure 3: Set theoretic representation of Figure 2.

$\pi(\mathbf{n}|\mathbf{h})$. We do not study such an extension of the model in this paper.

In Figure 3 we show the set theoretic version of Figure 2. Note how the distortion is intuitively and consistently represented by the sizes of the subsets depicted. Neither $\mathbf{h}(\mathbf{x})$ nor $y(\mathbf{h})$ is guaranteed to have a simple form (only $y(\mathbf{h}(\mathbf{x}))$ is a complete vector quantiser), but we have nevertheless represented them as if they mapped clusters of nearby points to the same output: this is purely for diagrammatic convenience. The equivalent noise which $y(\mathbf{h})$ induces on a particular \mathbf{h} (\mathbf{h}_0 say) is then represented by the set of \mathbf{h} which satisfy $y(\mathbf{h}) = y(\mathbf{h}_0)$. In Figure 3 we show how four separate \mathbf{h} values are mapped to the same y value, and we show how each of those four \mathbf{h} values derives from a small cluster of \mathbf{x} values, and how the four clusters themselves are clustered, thus guaranteeing a small overall L_2 distortion.

Note that we have not yet said anything about the detailed form of the noise model, not even whether it is a continuous function. It is determined entirely by the particular form of the distortion introduced by $y(\mathbf{h})$.

IV. A MODIFICATION OF THE LBG ALGORITHM

Now let us minimise D_1 with respect to $\mathbf{h}(\mathbf{x})$ and $\mathbf{x}'(\mathbf{h})$. The functional derivatives of D_1 with respect to $\mathbf{h}(\mathbf{x})$ and $\mathbf{x}'(\mathbf{h})$ are given by

$$\frac{\delta D_1}{\delta \mathbf{h}(\mathbf{x})} = P(\mathbf{x}) \int d\mathbf{n} \pi(\mathbf{n}) \left. \frac{\partial \|\mathbf{x}'(\mathbf{h}) - \mathbf{x}\|^2}{\partial \mathbf{h}} \right|_{\mathbf{h}=\mathbf{h}(\mathbf{x})+\mathbf{n}} \quad (4.1)$$

$$\frac{\delta D_1}{\delta \mathbf{x}'(\mathbf{h})} = 2 \int d\mathbf{x} P(\mathbf{x}) \pi(\mathbf{h} - \mathbf{h}(\mathbf{x})) [\mathbf{x}'(\mathbf{h}) - \mathbf{x}] \quad (4.2)$$

These two derivatives imply that the two necessary conditions A1 and A2 for minimising D , must be modified to become:

B1) Given \mathbf{x} , choose $\mathbf{h}(\mathbf{x})$ to minimise $\int d\mathbf{n} \pi(\mathbf{n}) \|\mathbf{x}'(\mathbf{h}(\mathbf{x}) + \mathbf{n}) - \mathbf{x}\|^2$.

B2) Given \mathbf{h} , choose $\mathbf{x}'(\mathbf{h})$ to satisfy

$$\mathbf{x}'(\mathbf{h}) = \frac{\int d\mathbf{x} P(\mathbf{x}) \pi(\mathbf{h} - \mathbf{h}(\mathbf{x})) \mathbf{x}}{\int d\mathbf{x} P(\mathbf{x}) \pi(\mathbf{h} - \mathbf{h}(\mathbf{x}))} \quad (4.3)$$

As a consistency check, if $\pi(\mathbf{n}) = \delta(\mathbf{n})$, where $\delta(\mathbf{n})$ is a Dirac delta function (or the analogous Kronecker

delta if we consider the discrete case), then the model reduces to the encoding/decoding depicted in Figure 1 (with y replaced by \mathbf{h}), and the above conditions reduce to the standard LBG vector quantisation conditions, as expected.

It remains to determine a suitable form for the noise model $\pi(\mathbf{n})$, which depends on the distorting effect of the second stage of quantisation $y(\mathbf{h})$. We need to be more specific about $y(\mathbf{h})$, so we shall now assume that $y(\mathbf{h})$ is a minimum L_2 distortion vector quantiser for \mathbf{h} . Thus we shall minimise D_2 with respect to $y(\mathbf{h})$ and $\mathbf{h}'(y)$, where D_2 is defined as

$$D_2 \equiv \int d\mathbf{h} P_h(\mathbf{h}) \|\mathbf{h}'(y(\mathbf{h})) - \mathbf{h}\|^2 \quad (4.4)$$

$P_h(\mathbf{h})$ is the PDF of the output \mathbf{h} of the first stage of quantisation. Given the value of D_2 , and no further information about the distorting effect of $y(\mathbf{h})$, a zero mean Gaussian noise model with variance D_2 is then the appropriate maximum entropy choice for $\pi(\mathbf{n})$.

We may simplify B1 by Taylor expanding $\mathbf{x}'(\mathbf{h}(\mathbf{x}) + \mathbf{n})$

$$\mathbf{x}'(\mathbf{h}(\mathbf{x}) + \mathbf{n}) = \exp[\mathbf{n} \cdot \nabla_{\mathbf{h}}] \mathbf{x}'(\mathbf{h})|_{\mathbf{h}=\mathbf{h}(\mathbf{x})} \quad (4.5)$$

where $\nabla_{\mathbf{h}}$ is the vector differentiation operator with respect to \mathbf{h} . This yields to second order

$$\int d\mathbf{n} \pi(\mathbf{n}) \|\mathbf{x}'(\mathbf{h}(\mathbf{x}) + \mathbf{n}) - \mathbf{x}\|^2 \simeq \left[1 + \frac{1}{2} D_2 \nabla_{\mathbf{h}}^2 \right] [\mathbf{x}'(\mathbf{h}) - \mathbf{x}] \Big|_{\mathbf{h}=\mathbf{h}(\mathbf{x})} \quad (4.6)$$

where the statistics of the maximum entropy zero mean Gaussian noise model for $\pi(\mathbf{n})$ are

$$\begin{aligned} \int d\mathbf{n} \pi(\mathbf{n}) &= 1 \\ \int d\mathbf{n} \pi(\mathbf{n}) n_i &= 0 \\ \int d\mathbf{n} \pi(\mathbf{n}) n_i n_j &= D_2 \delta_{i,j} \end{aligned} \quad (4.7)$$

The first term on the right-hand side of Equation 4.6 is the conventional distortion, and the second (curvature) term arises from the output noise model $\pi(\mathbf{n})$. We shall henceforth assume that the curvature term is small (i.e., $\mathcal{O}(D_2 \nabla_{\mathbf{h}}^2) \ll 1$, effectively), so that B1 may be approximated by A1 (with y replaced by $\mathbf{h}(\mathbf{x})$). This reduces B1 to a nearest neighbour encoding rule, as before.

B2 cannot be simplified in an analogous fashion to B1. However, from Equation 4.2 we may derive a stochastic gradient descent algorithm for realising B2. Thus we draw \mathbf{x} samples at random from the training set using the $\int d\mathbf{x} P(\mathbf{x})$ factor, and gradient descent of D_1 is then achieved by updates of the following form where $\epsilon > 0$, \mathbf{x} is drawn from the probability density function $P(\mathbf{x})$, and $\mathbf{h}(\mathbf{x})$ is the nearest neighbour encoding approximation to B1:

$$\mathbf{x}'(h') \longrightarrow \mathbf{x}'(h') + \epsilon \pi(h' - \mathbf{h}(\mathbf{x})) [\mathbf{x} - \mathbf{x}'(h')] \quad (4.8)$$

This update prescription is applied to all h' for which $\pi(h' - \mathbf{h}(\mathbf{x})) > 0$. This update scheme is identical to the one that Kohonen has proposed for obtaining topographic mappings [4], and it leads to a decoding operation $\mathbf{x}'(\mathbf{h})$ that is a smoothly varying function of \mathbf{h} , although the encoding operation $\mathbf{h}(\mathbf{x})$ is not guaranteed to be a smooth function of \mathbf{x} when $\dim(\mathbf{h}) < \dim(\mathbf{x})$. Note that we have had to make an approximation (where we retain only the leading term in Equation 4.6) in order to obtain the correspondence between a minimum L_2 dis-

tortion vector quantiser (with added code noise) and Kohonen's topographic mappings.

Our noise model does not make specific predictions about the dependence of $\pi(\mathbf{n})$ on the number of updates that have elapsed. However, D_1 has multiple local minima which makes it difficult to locate the global minimum by using the gradient descent algorithm specified in Equation 4.8. The randomness of sampling from the training set makes unpredictable the order in which updates are performed according to Equation 4.8, which alleviates the local minimum problem somewhat. However, a broader $\pi(\mathbf{n})$ gives rise to a D_1 which has a softer dependence on $\mathbf{h}(\mathbf{x})$ and $\mathbf{x}'(\mathbf{h})$, so better convergence to the global minimum of D_1 is obtained by starting with a broad $\pi(\mathbf{n})$, and then progressively reducing the width of $\pi(\mathbf{n})$ as the updates proceed until the required form is obtained [4]. This is equivalent to a renormalised scheme in which $\pi(\mathbf{n})$ is fixed while the underlying space is progressively dilated as the updates proceed [5]. The renormalised scheme gives much faster convergence than Kohonen's original scheme.

A similar approach to that presented in this paper has been presented elsewhere [6]. They minimise a quadratic potential which is analogous to our distortion measure D_1 in Equation 3.2. However, our approach provides a more detailed interpretation of the neighbourhood updating aspect of the training algorithm.

V. EXTENSION TO A MULTILAYER VECTOR QUANTISER

We shall now use the robust hidden layer principle to state the form of Equation 3.2 that applies to an adjacent pair of layers in a multilayer vector quantiser. Thus

define a distortion $D(m)$ for encoding/decoding layer m of a multilayer vector quantiser

$$D(m) \equiv \int d\mathbf{x}_m P_m(\mathbf{x}_m) \int d\mathbf{n} \pi_m(\mathbf{n}) \|\mathbf{x}'_m(\mathbf{x}_{m+1}(\mathbf{x}_m) + \mathbf{n}) - \mathbf{x}_m\|^2 \quad (5.1)$$

where \mathbf{x}_m is a vector representing the state of layer m of the quantiser, $\mathbf{x}_{m+1}(\mathbf{x}_m)$ is the encoding operation used to obtain the state of layer $m+1$ from layer m , $\mathbf{x}'_m(\mathbf{x}_{m+1})$ is the decoding operation used to obtain the state of layer m from layer $m+1$, and $\pi_m(\mathbf{n})$ is a noise

process that models the distortion due to all subsequent stages of the multilayer vector quantiser. The input probability density function $P_m(\mathbf{x}_m)$ is obtained as shown where we now write \mathbf{x}_1 instead of \mathbf{x} for the original input vector:

$$P_m(\mathbf{x}_m) = \int d\mathbf{x}_1 d\mathbf{x}_2 \cdots d\mathbf{x}_{m-1} P_1(\mathbf{x}_1) \delta(\mathbf{x}_2 - \mathbf{x}_2(\mathbf{x}_1)) \cdots \delta(\mathbf{x}_m - \mathbf{x}_m(\mathbf{x}_{m-1})) \quad (5.2)$$

In practice, samples \mathbf{x}_m from $P_m(\mathbf{x}_m)$ may easily be generated by propagating samples \mathbf{x}_1 from $P_1(\mathbf{x}_1)$ through the layers of the network until layer m is reached. Thus the output of layer $m-1$ is used as the training set for determining the optimum encoding/decoding operations on layer m .

The advantage of modeling the distorting effect of all subsequent layers of the multilayer vector quantiser as a noise process is that the encoding/decoding operations associated with layer m can be optimised (almost) independently of the other layers' encoding/decoding opera-

tions. Coupling between these optimisations is achieved via the distortion models $\pi_m(\mathbf{n})$. Although this will not lead to an optimal overall encoding/decoding scheme, it is computationally much cheaper than simultaneously optimising all of the coupled layers. Finally, when the multilayer vector quantiser has been trained, and we wish to decode \mathbf{x}_m to obtain a reconstruction $\mathbf{x}_{1,\text{rec}}(\mathbf{x}_m)$ of the input vector \mathbf{x}_1 , we should use the following formula which is the centroid of the set of input vectors that map to \mathbf{x}_m :

$$\mathbf{x}_{1,\text{rec}}(\mathbf{x}_m) = \frac{[\int d\mathbf{x}_1 d\mathbf{x}_2 \cdots d\mathbf{x}_{m-1} P_1(\mathbf{x}_1) \delta(\mathbf{x}_2 - \mathbf{x}_2(\mathbf{x}_1)) \cdots \delta(\mathbf{x}_m - \mathbf{x}_m(\mathbf{x}_{m-1})) \mathbf{x}_1]}{P_m(\mathbf{x}_m)} \quad (5.3)$$

VI. FURTHER EXTENSIONS

All of the results in this paper can be reexpressed using an arbitrary (non- L_2) distortion measure. The choice of a preferred form for $d(\mathbf{x}, \mathbf{x}')$ depends on whatever structure in \mathbf{x} one deems to be important.

The LBG algorithm is the batch training version of Kohonen's vector quantisation algorithm with zero neighbourhood size. It is also possible to use Equation 4.1 and Equation 4.2 to derive a batch training version of Kohonen's topographic mapping algorithm.

We have presented elsewhere some applications of 2-stage and multistage vector quantisation to low-level image processing [5], cluster decomposition of PDF's [7], image compression [8], and time series analysis [9].

A nice example of the use of the robust hidden layer principle has been reported in [10]. The problem studied there was the design of a codebook containing 8 bit code-words that needed to be robust with respect to bit errors. $\pi(\mathbf{n})$ must therefore describe Hamming noise, which induces an 8-dimensional Hamming cube topology on the codebook. A worthwhile reduction of 3 dB in the reconstruction distortion was obtained by topographically organising the codebook in this fashion.

VII. SUMMARY AND CONCLUSIONS

We have shown how topographic mappings arise naturally in a 2-stage vector quantiser whose output is made

robust with respect to various types of distortion. The choice of neighbourhood used at the output of the first stage is determined solely by the additional distortion introduced by the second stage of the quantiser. This principle generalises directly to multistage quantisers. The advantage of our approach is that (variants of) Kohonen's self-organising neural network may be derived from first principles.

We have encapsulated this robustness argument in the robust hidden layer principle: train your network in the

presence of all anticipated distortions to its internal representations. The network will then do its best to choose internal representations that encode information in such a way that it is robust with respect to being damaged by the distortions.

A hierarchical multistage quantiser constructed along these lines is an adaptive pyramid processor in which successively cruder quantisations (or codings) of the input are extracted. Current research indicates that such neural networks show much promise [9].

-
- [1] S. P. Lloyd (1982). Least squares quantisation in PCM. *IEEE Trans. Inform. Theory*, **28**(2), 129-137.
 - [2] J. Max (1960). Quantising for minimum distortion. *IEEE Trans. Inform. Theory*, **6**(1), 7-12.
 - [3] Y. Linde, A. Buzo and R. M. Gray (1980). An algorithm for vector quantiser design. *IEEE Trans. Commun.*, **28**(1), 84-95.
 - [4] T. Kohonen (1984). *Self-Organisation and Associative Memory* (Springer-Verlag, Berlin).
 - [5] S. P. Luttrell (1988). Self organising multilayer topographic mappings. In *Proc. 2nd IEEE Conf. on Neural Networks*, (IEEE, San Diego), 93-100.
 - [6] H. Ritter and K. Schulten (1988). Kohonen's self-organising maps: Exploring their computational capabilities. In *Proc. IEEE Int. Conf. on Neural Networks*, (IEEE, San Diego), 109-116.
 - [7] S. P. Luttrell (1989). The use of Bayesian and entropic methods in neural network theory. In *Maximum Entropy and Bayesian Methods*, edited by G. Erickson, J. T. Rychert and C. R. Smith (Kluwer, Dordrecht), 363-370.
 - [8] S. P. Luttrell (1989). Image compression using a multilayer neural network. *Patt. Recogn. Lett.*, **10**(1), 1-7.
 - [9] S. P. Luttrell (1989). Hierarchical vector quantisation. *Proc. Inst. Electr. Eng. I*, **136**(6), 405-413.
 - [10] D. S. Bradburn (1989). Reducing transmission error effects using a self-organising network. In *Proc. 3rd IEEE Int. Joint Conf. Neural Networks*, (IEEE, Washington DC), 531-537.