### **Information-Theoretic Quantization and Its Connection to Classification**



Brian M. Kurkoski Japan Advanced Institute of Science and Technology

**12 February 2019** 12th International ITG Conference on Systems, Communication and Coding - SCC2019 **Rostock, Germany** 



## Just Enough Information Theory



### What is the best code we can design?

 $R < C = \max_{p_{X}(x)} I(X; Y)$ Code rate < Channel Capacity Claude Shannon: mutual information is the highest achievable rate



# **Highest Achievable Rate for a Communications over a Quantized Channel**



Given a channel, find the quantizer Q which maximizes the achievable rate:  $C = \max_{Q} I(\mathsf{X};\mathsf{Z})$ 

We will fix the input distribution  $p_{X}(x)$ . Jointly optimizing Q and  $p_{X}(x)$  is a much more difficult problem.













Given a continuous-output channel, we want to create a discrete version

• For example, digital circuits deal with discrete values.





Given a continuous-output channel, we want to create a discrete version

- For example, digital circuits deal with discrete values.

X. Ma, X. Zhang, H. Yu, and A. Kavcic, "Optimal quantization for soft-decision decoding revisited," in Proc. Int. Symp. Inform. Theory Appl., Xian, China, Oct. 2002

• A quantizer Q maps real values Y to discrete values  $Z \in \{1, \ldots, M\}$ 

• How to choose the "quantization boundaries" to  $\max I(X; Z)$ ?





### "Quantization" of Discrete Memoryless Channel





Given a discrete memoryless channel and input distribution  $p_{XY}(x, y)$ , find the quantizer Q which maximizes mutual information:

$$Q^* = \arg\max_Q I(X; Z)$$

with  $|\mathsf{Z}| < |\mathsf{Y}|$ .

 $|Z| \ge |Y|$  is trivial.

"Channel downgrading" in polar code design.



**Factor Graph** 

Encoder-side function *f* 



**Decoder-side lookup-table** 



 $L_i$  is a noisy version of  $X_i$ 

Z is a noisy version of  $X_3$ 

J. Lewandowsky, M. Stark, and G. Bauch, "Information bottleneck graphs for receiver design," in Proceedings of IEEE International Symposium on Information Theory. IEEE, July 2016, pp. 2888–2892.







# Handwriting Recognition is Typical Example of Classification

noisy observation Y

input X

0  $\rightarrow$ 0000 0 1  $\rightarrow$ 2 222222  $\rightarrow$ 3 333333  $\rightarrow$ 4 844444  $\rightarrow$ 5 ந 5 **5** 5 5 5  $\rightarrow$ 6 6666  $\rightarrow$ 6 6 6 7  $\rightarrow$ М 8 8 8 8 C  $\rightarrow$ 8 9999999 g

Deep neural networks solve this problem well, but are a bit complicated  $_{\rm Wikimedia/Josef \, Steppan \, using \, MNIST \, data \, set}$ 

estimate Z

0	$\diamond$	$\mathcal{O}$	0	0	0	0	0	$\rightarrow$	0
/	1	/	1	1	١	/	1	$\rightarrow$	1
ዲ	2	2	2	2	2	2_	ン	$\rightarrow$	2
3	3	3	3	3	3	3	З	$\rightarrow$	3
4	4	4	4	4	4	4	4	$\rightarrow$	4
5	5	5	5	5	5	5	5	$\rightarrow$	5
6	6	Q	6	6	6	6	b	$\rightarrow$	6
7	7	7	7	7	7	7	7	$\rightarrow$	7
8	8	8	8	8	8	8	8	$\rightarrow$	8
٩	η	٩	9	9	9	9	9	$\rightarrow$	9



## Simple Example of Classification







# Simple Example of Classification





1 '

## **Goal and Outline**

Discuss tools for classification which can be applied to channel quantization Can find effective or even optimal methods for channel quantization Outline

- K-Means algorithm with Euclidean metrics
- K-Means with generalized metrics
  - "convexity theorem" by Burshtein et al
- Backwards channel perspective
- Three particular cases:
  - Quantization of arbitrary DMC using K-Means • Optimal quantization of binary input DMC using dynamic programming • Quantization of binary-input, continuous-output channels



K-means algorithm partitions n observations into kclusters — each observation belongs to the cluster with the nearest mean. Can also be seen as vector quantization.

Attempts to minimize mean-squared of quantization

$$\min_{Q} E\left[(\mathsf{X} - Q(\mathsf{X}))^2\right]$$

Not optimal, but works well. Widely used in machine learning.

# **K-Means Algorithm**





Contents lists available at ScienceDirect

### Pattern Recognition Letters

journal homepage: www.elsevier.com/locate/patrec

### Data clustering: 50 years beyond K-means $\stackrel{\text{\tiny{}^{\diamond}}}{}$

### Anil K. Jain \*

Department of Computer Science and Engineering, Michigan State University, East Lansing, Michigan 48824, USA Department of Brain and Cognitive Engineering, Korea University, Anam-dong, Seoul, 136-713, Korea

ARTICLE INFO

ABSTRACT





K-means algorithm partitions n observations into kclusters — each observation belongs to the cluster with the nearest mean. Can also be seen as vector quantization.

Attempts to minimize mean-squared of quantization

$$\min_{Q} E\left[(\mathsf{X} - Q(\mathsf{X}))^2\right]$$

Not optimal, but works well. Widely used in machine learning.

# **K-Means Algorithm**





### Data clustering: 50 years beyond K-means $\stackrel{\text{\tiny{}^{\diamond}}}{}$

### Anil K. Jain \*

Department of Computer Science and Engineering, Michigan State University, East Lansing, Michigan 48824, USA Department of Brain and Cognitive Engineering, Korea University, Anam-dong, Seoul, 136-713, Korea

ARTICLE INFO

ABSTRACT





### **Two Observations**

- Given reconstruction points  $\mathbf{m}_1, \mathbf{m}_2, \ldots, \mathbf{m}_K$ , the reconstruction region should minimize the average error. Or, any point y should be **assigned** to the closest point  $\mathbf{m}_i$ .
- Within some region  $\mathcal{R}_i$ , its reconstruction point  $\mathbf{m}_i$  should be **up**dated to minimize the average error for that region.







## **K-Means Algorithm with Euclidean Distance** Metric

given *n*-dimensional data set, randomly choose K means (centroids)



- Assignment step K clusters consists of data points closest to its mean in Euclidean distance
- **Update step** move the mean to the center of 3. the cluster

Nearest in Euclidean-distance sense. What about other metrics?







## **K-Means with Generalized Metrics**

Generalize objective function:

 $\frac{\text{Euclidean distance}}{\min_{Q} E\left[(\mathsf{X} - Q(\mathsf{X}))^2\right]}$ 

Change the Assignment/Update "local" function in the to match the objective function.





KL divergence, Gini index, Hamming error, Itakura-Saito, logistic loss, etc.





Two independent streams: Tishby et al, Banerjee et al did not cite the Chou stream.



# **Optimality of Assignment Step**

Consider a broad class of convex objective functions. There exists an optimal assignment such that each cluster is a convex set in the "backward channel"

The Annals of Statistics 1992, Vol. 20, No. 3, 1637–1646

### MINIMUM IMPURITY PARTITIONS

By David Burshtein, Vincent Della Pietra, Dimitri Kanevsky and Arthur Nádas

### 2. Results.

THEOREM 1. For any  $C: \mathscr{X} \to \mathscr{C}$  there exists a  $\tilde{C}: \mathscr{U} \to \mathscr{C}$  such that  $\Psi(\tilde{C}(Y)) \leq \Psi(C)$  and such that  $\tilde{C}^{-1}(c)$  is convex for all  $c \in \mathscr{C}$ .

Significance: for a wide variety of objective functions, the search can be restricted to convex clusters.





### **Original Problem: Maximization of Mutual Information**





Given a discrete memoryless channel and input distribution  $p_{XY}(x, y)$ , find the quantizer Q which maximizes mutual information:

$$Q^* = \arg\max_Q I(X; Z)$$

with  $|\mathsf{Z}| < |\mathsf{Y}|$ .

 $|\mathsf{Z}| \ge |\mathsf{Y}|$  is trivial.



20

## **Problem Setup and Backwards Channel**



 $X \in \{1, 2, 3\}$ 

 $\mathbf{u}_y = |\operatorname{Pr}(\boldsymbol{\lambda})|$ 

Justify this later.

Assume  $p_{XY}(x, y)$  is known; X is discrete. Examples show Y is discrete, but results can be extend to continuous case. Running example:

- $Y \in \{red, lime, yellow, orange, green, brown\}$
- Work with the backward channel:

$$\mathsf{X} = 1 | \mathsf{Y} = y), \dots, \Pr(\mathsf{X} = J | \mathsf{Y} = y)$$





$$\mathbf{u}_y = \Big[ \Pr(\mathsf{X} = 1 | \mathsf{Y} = y), \Pr(\mathsf{X} = 1 | \mathsf{Y} = y) \Big]$$











$$\mathbf{u}_{y} = \left[ \Pr(\mathsf{X} = 1 | \mathsf{Y} = y), \Pr(\mathsf{Y} = 1 | \mathsf{Y} = y), \Pr(\mathsf{Y} = 1 | \mathsf$$

$$\mathbf{x}$$

$$\mathbf{u}_{red} = [1, 0, 0]$$

$$\mathbf{u}_{lime} = [0.8, 0.2, 0]$$

$$\mathbf{u}_{yellow} = [0.3, 0.2, 0.5]$$

$$\mathbf{u}_{orange} = [0.33, 0.33, 0.34]$$

$$\mathbf{u}_{green} = [0.3, 0.3, 0.4]$$

$$\mathbf{u}_{brown} = [0.8, 0.1, 0.1]$$







$$\mathbf{u}_y = \Big[ \Pr(\mathsf{X} = 1 | \mathsf{Y} = y), \Pr(\mathsf{X} = 1 | \mathsf{Y} = y) \Big]$$









**Backward Channel** Pr(X | Y) as a Vector  $\mathbf{u}_y = \left\lceil \Pr(\mathsf{X} = 1 | \mathsf{Y} = y), \Pr(\mathsf{X} = 2 | \mathsf{Y} = y), \Pr(\mathsf{X} = 3 | \mathsf{Y} = y) \right\rceil$ 







### Visualizing the "Means"









### Visualizing the "Means"

 $p_{X|Z}(x|z)$  as well as  $p_{X|Y}(x|y)$ 





### From Mutual Information to KL Divergence

Recall: 
$$\mathbf{u}_y = | \Pr(\mathsf{X} = 1)$$

- Random vector versions:  $\mathbf{U} = \begin{bmatrix} \Pr(X = 1) \\ \mathbf{V} = \begin{bmatrix} \Pr(X = 1) \end{bmatrix}$
- Then, the following holds:

$$I(X; Y) - I(X; Z) =$$

 $D(\cdot || \cdot)$  is the Kullback-Leiber divergence and E is expectation.

$$Q^* = \arg\max_Q I(\mathsf{X};\mathsf{Z}) =$$

Thus, maximization of mutual information is minimization of KL divergence.

 $1|Y = y), Pr(X = 2|Y = y), \dots, Pr(X = J|Y = y)$ 

$$Y), \Pr(X = 2|Y), \dots, \Pr(X = J|Y)]$$
$$Z), \Pr(X = 2|Y), \dots, \Pr(X = J|Z)]$$

### $= E(D(\mathbf{U}||\mathbf{V}))$

 $\arg\min_{Q} E(D(\mathbf{U}||\mathbf{V}))$ 





### Three Particular Cases

- 1. General discrete quantization:
  - 1. K-Means with KL Divergence metrics "KL-Means"
  - 2. Information bottleneck method
- 3. Binary-input, continuous-output quantization with arbitrary noise

2. Binary-input, discrete output quantization: Dynamic programming quantization



### K-Means With KL Divergence Metric "KL-Means algorithm" replace Euclidean distance with KL distance

Min KL divergence = max. mutual information

 $Q^* = \arg\max_{Q} I(\mathsf{X};\mathsf{Z}) = \arg\min_{Q} E(D(\mathsf{U}||\mathsf{V}))$ 

A. Zhang and B. Kurkoski, "Low-Complexity Quantization of Non-Binary Input DMCs" ISITA 2016.







### K-Means With KL Divergence Metric "KL-Means algorithm" replace Euclidean distance with KL distance

Min KL divergence = max. mutual information

 $Q^* = \arg\max_{Q} I(\mathsf{X};\mathsf{Z}) = \arg\min_{Q} E(D(\mathsf{U}||\mathsf{V}))$ 

TEXT MINING WITH INFORMATION-THEORETIC CLUSTERING

Motivated by the success of hybrid information-retrieval algorithms, the authors report on the development of their hybrid clustering scheme. Scheme experiments on data in a reduced vector space model indicate a higher performance level over several existing clustering algorithms.

Α.

"n text mining and information retrieval, as partitions for the second step of the procedure.







# **K-Means With KL Divergence Metric**

### "KL-Means algorithm" replace Euclidean distance with KL distance

Min KL divergence = max. mutual information

$$Q^* = \arg\max_Q I(X; Z) = \arg\min_Q E(Q)$$

Numerical results show tradeoff:

- increasing number of quantizer outputs
- decreases the loss of mutual information

### A. Zhang and B. Kurkoski, "Low-Complexity Quantization of Non-Binary Input DMCs" ISITA 2016.







## Information Bottleneck Method

For Markov chain with joint distribution:



Tishby, Pereira, Bialek (1999) gave the information bottleneck method which finds:

$$\min_{p_{\mathsf{Z}|\mathsf{Y}}(z|y)} I(\mathsf{Y};\mathsf{Z}) - \beta I(\mathsf{X};\mathsf{Z})$$

When:

- $\beta$  finite: probabilistic clustering  $p_{Z|Y}(z|y)$
- $\beta \to \infty$ : hard clustring,  $p_{\mathsf{Z}|\mathsf{Y}}(z|y) = 0$  or 1

Q is the probability distribution  $p_{Z|Y}(z|y)$ 

Image credit: Boris Epshtein & Lena Gorelick

### Bottleneck



Cover and Thomas, *Elements of* Information Theory, 2006. Problem 7.52 on page 234.









## Information Bottleneck Method

 $\beta \geq 0$  sweeps a kind of rate-distortion curve.

When  $\beta \to \infty$  we get hard clustering. In fact, the information bottleneck method becomes Kmeans algorithm with KL divergence metric. "KL means algorithm"

B. M. Kurkoski, "On the relationship between the KL means algorithm and the information bottleneck method," in 11th International ITG Conference on Systems, Communications and Coding (SCC2017), (Hamburg, Germany), pp. 1-5, February 2017.





# **Quantizer Design for Binary Input DMC**



- Consider a binary-input channel
- (Burshtein et al)

• Convex in the space:

 $u_y = [\mathbf{F}]$ 

- all boundaries



B. Kurkoski and H. Yagi, "Quantization of Binary-Input Discrete Memoryless Channels," IEEE Trans on Information Theory, May 2014.

• For an optimal quantizer, there must exist an convex cluster

$$\Pr[\mathsf{X} = 0 | \mathsf{Y} = y], \Pr[\mathsf{X} = 1 | \mathsf{Y} = y]]$$

• **Dynamic programming** approach is suggested: search over

• Algorithm is optimal, complexity is  $|\mathcal{Y}|^3$ 



## Structure of Optimal Solution



variance = 0.07

### Structure of Optimal Solution

variance = 0.3





Consider binary-input, continuous-output  $y \in \mathbb{R}$  channel:

- For BI-AWGN solution is obvious, i.e. threshold at 0
- Arbitrary, data-dependent noise  $p_{Y|X}(y|+1)$  different from  $p_{Y|X}(y|-1)$
- Quantization to **one bit**

• Using Burshtein et al Theorem, will not restrict ourselves to a single threshold





### Quantizer:

$$z = \begin{cases} 0 & y < a & & 0.6 \\ 1 & y \ge a & & 0.4 \\ & & 0.2 \\ & & & 0.2 \\ & & & 0.2 \\ & & & 0.2 \\ & & & & 0.2 \\ & & & & 0.2 \\ & & & & & 0.2 \\ & & & & & 0.2 \\ & & & & & & 0.2 \\ & & & & & & & 0.2 \\ & & & & & & & 0.2 \\ & & & & & & & 0.2 \\ & & & & & & & 0.2 \\ & & & & & & & & 0.2 \\ & & & & & & & & & 0.2 \\ & & & & & & & & & 0.2 \\ & & & & & & & & & 0.2 \\ & & & & & & & & & 0.2 \\ & & & & & & & & & 0.2 \\ & & & & & & & & & 0.2 \\ & & & & & & & & & & 0.2 \\ & & & & & & & & & 0.2 \\ & & & & & & & & & & 0.2 \\ \end{array}$$

Example: data-dependent Gaussian noise mixtures

$$\max I(X;Z) = 0.493 \text{ at}$$

$$a^* = -0.153$$

$$0.4930$$

$$0.4$$

$$0.4$$

$$0.4$$

$$0.3$$

0.8

0

-6

### **Quantization By Threshold Search**





## **Backward Channel** Pr(X = x | Y = y)





### **Convex Backward Channel Quantizer**

Threshold  $\tilde{a}$ . Backward channel quantizer Q:





### **Optimal Quantizer in Backward Channel**



Happens to be concave. But not in general





### Even Backwards Channel is Hard





Mutual information is not convex in  $\tilde{a}$ 



# **Optimal Convex Forward Quantizer**

Lemma 2 If the channel log-likelihood ratio satisfies:

$$\log \frac{\Pr(y|X=1)}{\Pr(y|X=0)} \le \log \frac{\Pr(y'|X=1)}{\Pr(y'|X=0)}$$

for all y < y', then there exists an optimal forward channel quantizer  $Q^*$ which is a convex quantizer.

Consequences:

- convex.
- The BI-AWGN channel satisfies this condition

B. M. Kurkoski and H. Yagi, "Single-bit quantization of binary-input, continuous-output channels," in Proceedings of IEEE International Symposium on Information Theory, (Aachen, Germany), pp. 2088-2092, June 2017

• For many well-behaved channels, the optimal forward channel quantizer is



### **Channel Quantization and Classification**

### Information Theory



- to quantize channels
- For binary-input, discrete-output channels, optimality is possible
- For binary-input, continuous output channels, still work to be done

### Machine Learning

Before 2002, little interest in channel quantization to maximize mutual information.

• But in the 1990s the pattern recognition community had developed necessary tools

• K-means with KL divergence/information bottleneck is effective but suboptimal way

