

Natural Image Statistics and Low-complexity Feature Selection

Manuela Vasconcelos Nuno Vasconcelos
Statistical Visual Computing Laboratory,
University of California, San Diego
La Jolla, CA 92039

Abstract

Low-complexity feature selection is analyzed in the context of visual recognition. It is hypothesized that high-order dependences of bandpass features contain little information for discrimination of natural images. This hypothesis is characterized formally, by introduction of the concepts of conjunctive interference and decomposability order of a feature set. Necessary and sufficient conditions, for the feasibility of low-complexity feature selection, are then derived in terms of these concepts. It is shown that the intrinsic complexity of feature selection is determined by the decomposability order of the feature set, not its dimension. Feature selection algorithms are then derived for all levels of complexity, and shown to be approximated by existing information theoretic methods, which they consistently outperform. The new algorithms are also used to objectively test the hypothesis of low decomposability order, through comparison of classification performance. It is shown that, for image classification, the gain of modeling feature dependencies has strongly diminishing returns: best results are obtained under the assumption of decomposability order 1. This suggests a generic law for bandpass features extracted from natural images: that the effect, on the dependence of any two features, of observing any other feature is constant across image classes.

Index Terms

Feature Extraction and Construction, Low-complexity, Natural image statistics, Information Theory, Feature Discrimination vs Dependence, Image databases, Object recognition, Texture, Perceptual reasoning

I. INTRODUCTION

Natural image statistics have been a subject of substantial recent research in computer and biological vision [1]–[14]. For computer vision, good models of image statistics enable algorithms tuned to the scenes that matter the most. Tuning to natural statistics can be accomplished through priors that favor solutions consistent with them [15]–[18], or optimal solutions derived from probability models which enforce this consistency [19]–[23]. The idea of optimal tuning to natural statistics also has a long history in biological vision [5], [24]–[26], where this tuning is frequently used to justify neural computations. In fact, various recent advances in computational modeling of biological vision follow from connections between neural function and properties of natural stimulus statistics [8], such as sparseness [12], [27], independence [13], [14], compliance with certain probability models [28], or optimal statistical estimation [29], [30].

Although natural images are quite diverse, their convolution with banks of band-pass functions gives rise to frequency *coefficients* with remarkably stable statistical properties [1]–[4], [6]–[8], [10]. This is illustrated by Fig. 1 a), which presents three images, the histograms of one coefficient of their wavelet decomposition, and the histogram of that coefficient conditioned on its parent. The different visual appearance of the images affects the scale (variance) of the marginal distribution, *but not its shape or that of the conditional distribution*, which is a bow-tie for all classes. This canonical pattern is simply rescaled to match the marginal statistics of each class. These type of properties have been exploited in various image processing domains, including compression [1], [2], [6], [19], de-noising [15], [16], [18], [22], retrieval [21], saliency [31], extraction of intrinsic images [20], separation of reflections [32] and inpainting [17], [18]. In fact, the study of image statistics has a complementary relationship with the development of vision algorithms. Typically, an hypothesis is advanced for the statistics, an algorithm derived under that hypothesis, and applied to natural images. If the algorithm performs well, the hypothesis is *validated*.

This *indirect validation paradigm* is useful in two ways. First, it avoids the estimation of complex statistical quantities. For example, hypotheses on high-order statistics are difficult to verify experimentally, due to the well known difficulties of estimating such statistics [33]. Instead, it is usually easier to 1) derive an algorithm that is *optimal* if the hypothesis holds, and 2) apply it to a specific vision problem, such as object recognition [34], where performance can be *easily*

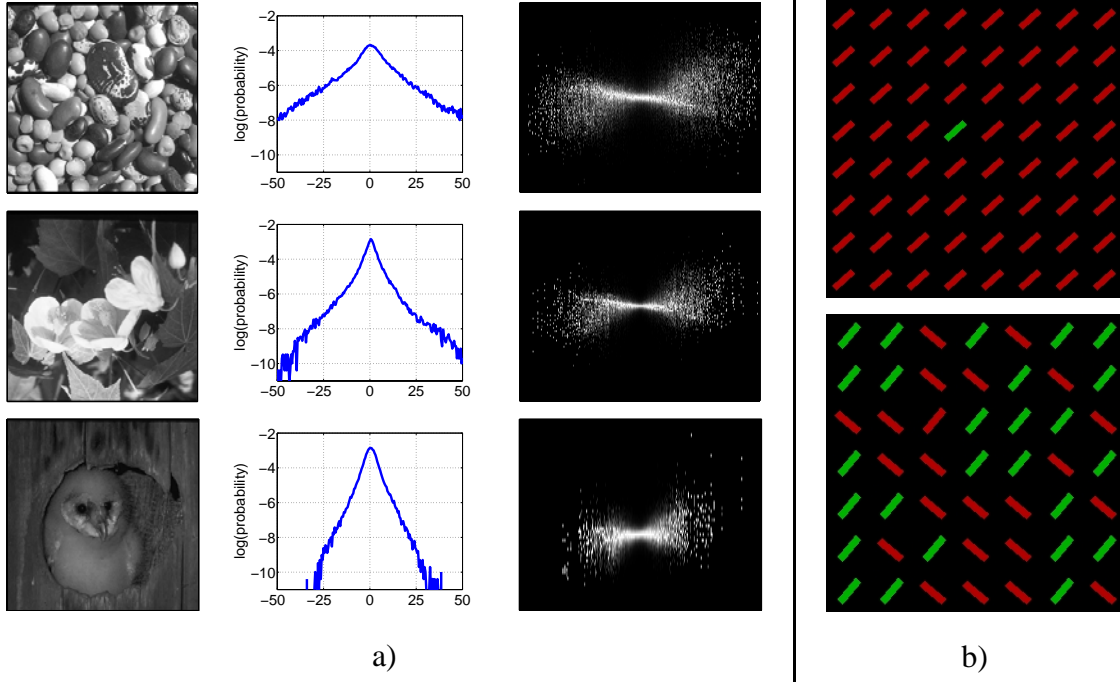


Fig. 1. a) Constancy of natural image statistics. Left: three images. Center: each plot presents the histogram of the same coefficient from a wavelet decomposition of the image on the left. Right: conditional histogram of the coefficient conditioned on the value of the co-located coefficient of immediately coarser scale (its parent). b) Biological vision frequently disregards feature dependences. Top: a stimulus that differs from its surrounds by a single feature (color) is salient. Bottom: differences in feature conjunctions (color and orientation) are not.

quantified. If the algorithm performs poorly there is reason to question the hypothesis, otherwise there is concrete evidence in its support. The second advantage of indirect validation is that it produces new vision algorithms which, *under the hypothesis*, are *optimally* tuned to the image statistics. If the hypothesis holds, these algorithms can outperform the state-of-the-art.

In this work, we adopt the indirect validation paradigm to study the *discriminant power of the statistical dependencies of frequency coefficients extracted from natural images*. While simple inspection of the histograms of Fig. 1 a) shows that these dependences exist, their constancy across image classes suggests the hypothesis that *high-order dependences contain little information for image discrimination*. This hypothesis is supported by what is known about biological vision, where it has long been argued that the early visual system dismisses feature dependences in the solution of discriminant tasks, such as visual search [35], [36]. This is illustrated by Fig. 1 b), which presents a classical example of the inability of pre-attentive vision to process feature conjunctions. When, as on the top, an object (colored bar) differs from

a background of distractors (other colored bars) in terms of a single feature (color), it can be easily discriminated (it *pops out*). However if, as on the bottom, the object differs from the distractors by a conjunction of two features (color and orientation, the bar on the 3rd row and 3rd column), there is no percept of *pop-out*. Current explanations attribute this phenomena to independent feature processing [35]–[40].

For computer vision, where models of feature dependences require estimation of high-dimensional densities, such dependences are a dominant source of complexity. A formal characterization of their role in image discrimination is, therefore, a pre-requisite for *optimal image classification with reduced complexity*. Since optimal classification requires discriminant features, we study dependences in the context of feature selection. In the spirit of indirect validation, we 1) develop optimal feature selection algorithms under the hypothesis that high-order dependences are uninformative for discrimination, and 2) evaluate their image classification performance.

The contributions of this effort are in three areas. The first is a rigorous characterization of the role of image statistics in optimal feature selection with low-complexity. We equate complexity with the dimensionality of the probability densities to be estimated, and adopt an information theoretic definition of optimality widely used in the literature [41]–[65]. We then derive, for each level of complexity, *the necessary and sufficient condition (on the statistics) for optimal feature selection with that complexity*. This condition depends exclusively on a quantity denoted as the *conjunctive interference* within the set of features \mathbf{X} , which roughly measures how, on average, the dependence between two disjoint feature subsets $\mathbf{A}, \mathbf{B} \subset \mathbf{X}$ is affected by the remaining features in \mathbf{X} . It is shown (see Theorem 1) that if this conjunctive interference is constant across classes, the complexity of the optimal solution is determined by the dimension of the subsets \mathbf{A}, \mathbf{B} , rather than that of \mathbf{X} . Hence, *the smaller the set size for which conjunctive interference is non-discriminant, the smaller the intrinsic complexity of feature selection*.

The second contribution, which follows from the theoretical analysis, is a new family of feature selection algorithms. These algorithms optimize simplified costs at all levels of complexity, and are (locally) optimal when conjunctive interference is non-discriminant at their complexity level. This family generalizes a number of low-complexity information theoretic methods [41]–[64] previously shown to outperform many state-of-the-art feature selection techniques [48], [58]. The impressive empirical performance of the previous methods is explained by the fact that they *approximate* the algorithms now derived. Nevertheless, there is a gain in replacing the

approximations with the optimal algorithms: experiments on various datasets show that the latter *consistently outperform the previous methods*, sometimes by a significant margin.

The final contribution, in the spirit of indirect validation, is the use of the feature selection algorithms to indirectly characterize the image statistics. Given that the different algorithms are optimal only when conjunctive interference is non-discriminant at their complexity level, a comparison of feature selection performance identifies the complexity at which conjunctive interference ceases to affect image discrimination. Algorithms with less than this complexity are sub-optimal, and performance levels off once it is reached. We present evidence for the hypothesis that this “leveling off” effect occurs at very low complexity levels. While simply modeling marginal densities is, in general, not enough to guarantee optimal feature selection, *there appears to be little gain in estimating more than the densities of pairs of coefficients*.

The paper is organized as follows. Section II reviews information theoretic feature selection. Section III introduces a basic decomposition of the information theoretic cost, and shows that independent feature selection can be optimal even for highly dependent feature sets. The decomposition is refined in Section IV, which formally defines conjunctive interference, and introduces a measure of the intrinsic complexity of a feature set (*decomposability order*). Section V introduces the new family of (locally) optimal algorithms, and discusses connections to prior methods. Finally, the experimental protocol for indirect validation of the decomposability hypothesis is introduced in Section VI, and experimental results discussed in Section VII. A very preliminary version of the work, focusing mostly on the theoretical connections between conjunctive interference and low complexity feature selection, has appeared in [64].

II. INFOMAX FEATURE SELECTION

We start by introducing the information theoretic optimality criterion adopted in this work, and reviewing its previous uses in the feature selection literature.

A. Definitions

A classifier $g : \mathcal{X} \rightarrow \mathcal{L} = \{1, \dots, M\}$ maps a feature vector $\mathbf{x} = (x_1, \dots, x_N)^T \in \mathcal{X} \subset \mathbb{R}^N$ into a class label $i \in \mathcal{L}$. Feature vectors result from a transformation $T : \mathcal{Z} \rightarrow \mathcal{X}$ of observation vectors $\mathbf{z} = (z_1, \dots, z_D)$ in measurement space $\mathcal{Z} \subset \mathbb{R}^D$. Observations are samples from random process \mathbf{Z} , of probability distribution $P_{\mathbf{Z}}(\mathbf{z})$ on \mathcal{Z} , feature vectors samples

from process \mathbf{X} , of distribution $P_{\mathbf{X}}(\mathbf{x})$ on \mathcal{X} , and labels samples from random variable Y , of distribution $P_Y(i)$ in \mathcal{L} . Given class i , observations have class-conditional density $P_{\mathbf{Z}|Y}(\mathbf{z}|i)$, and class-posterior probabilities determined by Bayes rule, $P_{Y|\mathbf{Z}}(i|\mathbf{z}) = P_{\mathbf{Z}|Y}(\mathbf{z}|i)P_Y(i)/P_{\mathbf{Z}}(\mathbf{z})$. The classification problem is uniquely defined by $\mathcal{C} = \{\mathcal{Z}, P_{\mathbf{Z}|Y}(\mathbf{z}|i), P_Y(i), i \in \mathcal{L}\}$. T induces class-conditional densities, $P_{\mathbf{X}|Y}(\mathbf{x}|i)$, in \mathcal{X} and defines a new classification problem $\mathcal{C}_{\mathcal{X}} = \{\mathcal{X}, P_{\mathbf{X}|Y}(\mathbf{x}|i), P_Y(i), i \in \mathcal{L}\}$. We define as optimal the spaces of maximum mutual information between features and class label.

Definition 1: Given a classification problem \mathcal{C} and a set \mathcal{S} of range spaces for the feature transforms under consideration, the infomax space is

$$\mathcal{X}^* = \arg \max_{\mathcal{X} \in \mathcal{S}} I(Y; \mathbf{X}) \quad (1)$$

where

$$I(\mathbf{X}; Y) = \sum_i \int_{\mathcal{X}} p_{\mathbf{X}, Y}(\mathbf{x}, i) \log \frac{p_{\mathbf{X}, Y}(\mathbf{x}, i)}{p_{\mathbf{X}}(\mathbf{x})p_Y(i)} d\mathbf{x}. \quad (2)$$

is the mutual information (MI) between \mathbf{X} and Y .

Infomax is closely related to the minimization of Bayes classification error, and has a number of relevant properties for low-complexity feature selection, some which are reviewed in Appendix I. In what follows, \mathbf{z} is a vector of image pixels, and \mathbf{x} the result of a bandpass transformation (e.g. a wavelet, Gabor, or windowed Fourier transform), followed by selection of N coefficients.

B. Previous Infomax approaches to feature selection

Information theoretic feature selection has been used for text categorization [41]–[44], creation of semantic ontologies [45], analysis of genomic microarrays [46], [47], classification of electroencephalograms (EEG) [49], [50] and sonar pulses [53], [54], medical diagnosis [51], audio-visual speech recognition [56] and visualization [57]. In computer vision, it been used for face detection [58], object recognition [59], [61], and image retrieval [62]–[64]. These approaches can be grouped into four classes. Algorithms in first class approximate (2) with

$$M(\mathbf{X}; Y) = \sum_{k=1}^D I(X_k; Y), \quad (3)$$

where $I(X_k; Y)$ is the MI between feature X_k and class label Y . $M(\mathbf{X}; Y)$ is a measure of the discriminant information conveyed by individual features. It is denoted as *marginal*

mutual information (MMI), and its maximization as *marginal infomax*. It is popular in text categorization [41]–[43] mostly due to its computational simplicity. It has, nevertheless, been shown to sometimes outperform methods which account for feature dependences [45], [51], [56].

Algorithms in the second class combine an heuristic extension of marginal infomax, originally proposed in [53], and the classical greedy strategy of sequential forward feature selection [66], where one feature is selected at a time. Denoting by $\mathbf{X}^* = \{X_1^*, \dots, X_k^*\}$ the set of previously selected features, and X a candidate feature, the selected feature is

$$X_{k+1}^* = \arg \max_X \{I(X; Y) - f(X, \mathbf{X}^*)\}, \quad (4)$$

where $f(\cdot)$ is a dependence measure, ranging from a hard rejection of dependent features [53] to continuous penalties. The most popular is [47], [48], [51], [52], [54], [55]

$$f(X, \mathbf{X}^*) = \xi \sum_{i=1}^k I(X; X_i^*), \quad (5)$$

where ξ controls the strength of the dependence penalty. Various information theoretic costs are either special cases of this [47], [48], or extensions that automatically determine ξ [55].

Algorithms in the third class optimize costs closer to (1), once again through sequential forward search. One proposal is to select the feature X which maximizes $I(X, X_i^*; Y)$, $i \in \{1, \dots, k\}$ [57]. This is a low complexity approximation to $I(\mathbf{X}; Y)$, which only considers pairs of features. Because it does not rely on a modular decomposition of the MI, it is somewhat inefficient. An alternative, proposed in [58], [60], addresses this problem by relying on

$$X_{k+1}^* = \arg \max_X \min_i I(Y; X|X_i^*) = \arg \max_X \min_i [I(X, X_i^*; Y) - I(X_i^*; Y)], \quad (6)$$

where we have used (31). This is equivalent (see (34)) to

$$X_{k+1}^* = \arg \max_X \{I(X; Y) + \min_i [I(X; X_i^*|Y) - I(X; X_i^*)]\}. \quad (7)$$

We will show that (4) and (7) are simplifications of (1) which disregard important components for image discrimination. Nevertheless, extensive empirical studies have shown that they can beat state-of-the-art methods [48], [58], such as boosting [67], [68] or decision trees [69].

The final class has a single member, the algorithm of [65]. Unlike the other classes, it sequentially eliminates features from \mathbf{X} . This elimination is based on the concept of a Markov blanket [70]: if there is a set of features \mathbf{M} (called a Markov blanket), such that X is conditionally independent of $(\mathbf{X} \cup Y) - \mathbf{M} - \{X\}$ given \mathbf{M} , the feature X can be removed from \mathbf{X} without any

loss of information about Y . While theoretically sound, this method has a number of practical shortcomings which are acknowledged by its authors: the Markov blanket condition is much stronger than what is really needed (conditional independence of X from Y given \mathbf{M}), there may not be a full Markov blanket for a feature, and when there is one it can be difficult to find. To overcome these problems, [65] uses various heuristics which only involve feature pairs. The assumptions, with respect to the feature statistics, underlying these heuristics are not clear.

III. OPTIMALITY OF MARGINAL INFOMAX

To gain some intuition on the feasibility of low-complexity feature selection, we start by investigating the conditions under which marginal infomax is identical to (1).

A. Features vs conjunctions

For this, we note that the MI can be decoupled into contributions from individual features and feature conjunctions.

Lemma 1: Let $\mathbf{X} = (X_1, \dots, X_D)$ be any feature set, and $\mathbf{X}_{1,k} = (X_1, \dots, X_k)$. Then

$$I(\mathbf{X}; Y) = M(\mathbf{X}; Y) + C(\mathbf{X}; Y), \quad (8)$$

where $M(\mathbf{X}; Y)$ is the MMI of (3) and

$$C(\mathbf{X}; Y) = \sum_{k=2}^D [I(X_k; \mathbf{X}_{1,k-1}|Y) - I(X_k; \mathbf{X}_{1,k-1})]. \quad (9)$$

Proof: See Appendix II. ■

The terms $I(X_k; \mathbf{X}_{1,k-1}|Y) - I(X_k; \mathbf{X}_{1,k-1})$ measure how the MI between features is affected by knowledge of the class label. They quantify the discriminant information due to feature dependences. $C(\mathbf{X}; Y)$ is referred to as the *conjunctive component* of the MI (CCMI). A consequence of Lemma 1 is that, if $C(\mathbf{X}, Y) = 0, \forall \mathcal{X} \in \mathcal{S}$, then (1) reduces to marginal infomax

$$\mathcal{X}^* = \arg \max_{\mathcal{X} \in \mathcal{S}} \sum_k I(X_k; Y). \quad (10)$$

Due to the non-negativity of MI, (10) has a simple solution: order the X_k by decreasing $I(X_k; Y)$ and select the largest N . While (1) involves combinatorial search and high-dimensional density estimation, (10) only requires a linear search based on marginal density estimates. Hence, a null CCMI is a sufficient condition for low-complexity feature selection.

B. The role of natural image statistics

To obtain some intuition on how the CCMI is affected by the dependency structure of \mathbf{X} , we consider the classification of two Gaussian features $\mathbf{X} = (X_1, X_2)$ with

$$P_{\mathbf{X}|Y}(\mathbf{x}|i) = \frac{1}{\sqrt{4\pi^2|\boldsymbol{\Sigma}_i|}} e^{-\frac{1}{2}\mathbf{x}^T\boldsymbol{\Sigma}_i^{-1}\mathbf{x}}, \quad i \in \{1, 2\},$$

$$\boldsymbol{\Sigma}_i = \begin{bmatrix} \epsilon_i & \gamma_i \\ \gamma_i & \eta_i \end{bmatrix}, \quad \boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2.$$

Gaussianity reduces all class-conditional dependences to two parameters, the correlation coefficients $\rho_i = \gamma_i/\sqrt{\epsilon_i\eta_i}$. It is relatively straightforward to measure the relative strength

$$R(\mathbf{X}; Y) = \frac{C(\mathbf{X}; Y)}{M(\mathbf{X}; Y)} \quad (11)$$

of the MI components as a function of these parameters. If the variances ϵ_i and η_i are held constant, fixing the marginal distributions, then $R(\mathbf{X}; Y)$ is proportional to $C(\mathbf{X}; Y)$, allowing the study of how the latter depends on the ρ_i . By repeating the experiment with different ϵ_i and η_i it is also possible to infer how this dependence is affected by the MMI, $M(\mathbf{X}; Y)$. The graph of $R(\mathbf{X}; Y)$ vs. ρ_i , for fixed MMI, is the *CCMI surface* associated with the latter. While natural images statistics are not Gaussian, this procedure provides intuition on how the MI is affected by feature dependences. We consider two common scenarios for pairs of bandpass coefficients.

- **S1:** two features that are active/inactive for the same images (e.g. a wavelet coefficient and its parent). X_1 and X_2 have equal variance ($\epsilon_i = \eta_i = \nu_i$) and are inactive for one class ($\nu_2 = 1$) but active for the other ($\nu_1 > 1$). The CCMI surface is measured for various activity levels (by controlling ν_1).
- **S2:** each feature active for one class but not the other, e.g., X_1 (X_2) horizontally (vertically) tuned and class 1 (2) predominantly composed of horizontal (vertical) lines. The variances are $\epsilon_1 = \eta_2 = \nu$ and $\epsilon_2 = \eta_1 = 1$. The CCMI surface is measured for various ν .

Fig. 2 presents the corresponding CCMI surfaces, suggesting three main conclusions. First, *the CCMI can be close to zero even when the features are very strongly dependent*. Note that all surfaces are approximately zero along the line $\rho_1 = \rho_2 = \rho$, independently of either ρ (dependence strength) or the MMI. Second, *the importance of the CCMI in (8) increases with the diversity of the dependence across classes*, i.e. with $|\rho_1 - \rho_2|$. Third, *this increase is inversely proportional to the MMI*. While, for small MMI, a significant difference between the ρ_i makes

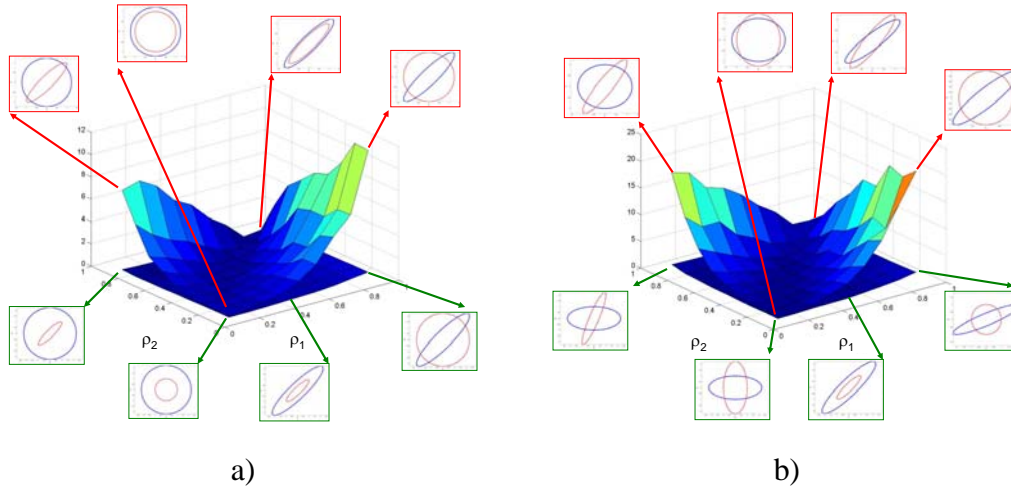


Fig. 2. $R(\mathbf{X}, Y)$ as a function of the class-conditional correlations ρ_i , for a binary Gaussian problem. The inserts show one standard deviation contours of the two Gaussian classes for various values of (ρ_1, ρ_2) . The plots report to scenarios **S1** (a) and **S2** (b). In both cases, different surfaces report to different values of ν , the variable that controls the marginal discrimination. All MIs were evaluated by replacing expectations with sample means, obtained from a sample of 10,000 points per class.

$R(\mathbf{X}, Y)$ large, this is not the case for large MMI. Overall, (8) (and Fig. 2) shows that 1) the relevance of feature dependences to the solution of (1) increases with their inter-class variability, but 2) this variability only boosts the importance of features that are not discriminant per se.

In summary, $C(\mathbf{X}, Y) = 0$ is a sufficient condition for optimal feature selection with low-complexity. It *does not require feature independence, but simply that the discriminant power of feature dependences is small*. As seen in Fig. 1 a), this hypothesis is not unreasonable for natural images. We will evaluate it in Section VII. For now, we consider a series of extensions that bridge the gap between (1) and (10).

IV. DECOMPOSITIONS OF THE CONJUNCTIVE COMPONENT

If feature conjunctions are discriminant, it is unlikely that this will hold for *all* conjunctions. For example, wavelet coefficients are dependent on their immediate neighbors (in space, scale, or orientation), but the dependence decays quickly [71]. Hence, $C(\mathbf{X}, Y)$ *should not require modeling dependences between all coefficients*. We next derive conditions for the optimality of infomax costs that only account for dependences within low-dimensional feature subsets.

A. Decompositions of the mutual information

We start by considering the decomposition of $I(\mathbf{X}, Y)$ for a given feature set \mathbf{X} . We group the D features into a collection of disjoint subsets of cardinality l

$$\mathcal{C}_l = \{\mathbf{C}_1, \dots, \mathbf{C}_{\lceil D/l \rceil}\}, \quad (12)$$

where¹

$$\mathbf{C}_i = \begin{cases} \{X_{(i-1)l+1}, \dots, X_{il}\}, & \text{if } i < \lceil D/l \rceil, \\ \{X_{(i-1)l+1}, \dots, X_D\}, & \text{if } i = \lceil D/l \rceil. \end{cases} \quad (13)$$

and $\lceil x \rceil = \inf\{m \in \mathbb{Z} | x \leq m\}$, and derive the conditions under which the CCMI is totally determined by the dependencies within each \mathbf{C}_i . This is based on the following decomposition.

Lemma 2: Consider the decomposition of \mathbf{X} into a subset collection \mathcal{C}_l , as in (12). Then

$$C(\mathbf{X}, Y) = \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}, Y) - I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}) \right]. \quad (14)$$

where \mathbf{C}_i are as in (13), $\tilde{\mathbf{C}}_{i,k}$ is the subset of features in \mathbf{C}_i whose index is smaller than k , and $\mathbf{C}_1^{i-1} = (\mathbf{C}_1, \dots, \mathbf{C}_{i-1})$.

Proof: See Appendix III. ■

This decomposition offers an explanation for why, in the absence of statistical regularities, low complexity feature selection is impossible [72]. Note that, although \mathbf{C}_1^{i-1} shares no elements with $\{X_k\}$ or $\tilde{\mathbf{C}}_{i,k}$, the state of the features of the former affects the dependences between those in the latter. Hence, *the discriminant information due to the dependences between X_k and $\tilde{\mathbf{C}}_{i,k}$ depends on the state of \mathbf{C}_1^{i-1} , and is impossible to compute with low complexity.* We refer to these indirect dependence relationships, i.e. that the state of a subset of features interferes with the dependence between two other non-overlapping subsets, as *2nd-order components of dependence*. This is opposed to direct dependences between subsets, which are referred to as *1st-order components*, or dependences within subsets, that we denote of *0th order*. The *conjunctive interference* within a feature set is the overall difference between the 1st and 2nd order dependences of its subsets.

Definition 2: Consider the decomposition of \mathbf{X} into a subset collection \mathcal{C}_l , as in Lemma 2. The conjunctive interference within \mathbf{X} , with respect to \mathcal{C}_l , is

$$CI(\mathbf{X}; \mathcal{C}_l) = \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}) - I(X_k; \tilde{\mathbf{C}}_{i,k}) \right]. \quad (15)$$

¹What follows could be extended to subsets \mathbf{C}_i of different cardinality, but this would complicate the notation and is omitted.

Conjunctive interference is a *differential measure of dependence*. It measures how, across the feature set, the dependence between two sets of features (e.g. $(X_k, \tilde{\mathbf{C}}_{i,k})$) *changes* with the observation of a third, non-overlapping, set (\mathbf{C}_1^{i-1}). Since if (\mathbf{A}, \mathbf{B}) is independent of \mathbf{C} then $I(\mathbf{A}; \mathbf{B} | \mathbf{C}) = I(\mathbf{A}; \mathbf{B})$, it follows that conjunctive interference within \mathbf{X} (with respect to decomposition \mathcal{C}_l) is null when $(X_k, \tilde{\mathbf{C}}_{i,k})$ is independent of \mathbf{C}_1^{i-1} for all valid i and k . We next show that this is not a necessary condition for low-complexity evaluation of the MI. It suffices that the conjunctive interference does not depend on the class.

Theorem 1: Consider the decomposition of \mathbf{X} into \mathcal{C}_l , as in (12). Then

$$I(\mathbf{X}; Y) = M(\mathbf{X}; Y) + C_{\mathcal{C}_l}(\mathbf{X}; Y) \quad (16)$$

with $M(\mathbf{X}, Y)$ as in (3), and

$$C_{\mathcal{C}_l}(\mathbf{X}; Y) = \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} [I(X_k; \tilde{\mathbf{C}}_{i,k} | Y) - I(X_k; \tilde{\mathbf{C}}_{i,k})], \quad (17)$$

if and only if

$$CI(\mathbf{X}; \mathcal{C}_l) = \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}, Y) - I(X_k; \tilde{\mathbf{C}}_{i,k} | Y) \right]. \quad (18)$$

Proof: See Appendix IV. ■

When (18) holds, (16) is equivalent to (8), with $C_{\mathcal{C}_l}(\mathbf{X}; Y)$ playing the role of $C(\mathbf{X}; Y)$. In particular, (16) replaces each of the terms

$$I(X_k; \mathbf{X}_1^{k-1} | Y) - I(X_k; \mathbf{X}_1^{k-1}) \quad (19)$$

of (9) by a sum, over i , of terms of the form

$$I(X_k; \tilde{\mathbf{C}}_{i,k} | Y) - I(X_k; \tilde{\mathbf{C}}_{i,k}). \quad (20)$$

While (19) quantifies the discriminant information due to dependences between X_k and the *entire* set of $X_j, j < k$, (20) restricts this measure to dependences between X_k and subset $\tilde{\mathbf{C}}_{i,k}$. Hence, (20) requires density estimates of dimension at most $l + 1$. Since density estimation has exponential complexity on feature space dimension, the complexity difference between (16) and (8) can be very significant if $l \ll D$. To illustrate this, we analyze a simple example.

Example 1: Let $D = 6, l = 2$. Then, $\mathbf{C}_1 = \{X_1, X_2\}$, $\mathbf{C}_2 = \{X_3, X_4\}$, and $\mathbf{C}_3 = \{X_5, X_6\}$, and $C_{\mathcal{C}_l}(\mathbf{X}; Y)$ is the sum of the terms in the third column of Table I. These terms measure

TABLE I
TERMS OF (17) AND (15) WHEN $D = 6$, AND $l = 2$.

k	i	$I(X_k; \tilde{\mathbf{C}}_{i,k} Y) - I(X_k; \tilde{\mathbf{C}}_{i,k})$	$I(X_k; \tilde{\mathbf{C}}_{i,k} \mathbf{C}_1, \dots, \mathbf{C}_{i-1}) - I(X_k; \tilde{\mathbf{C}}_{i,k})$
2	1	$I(X_2; X_1 Y) - I(X_2; X_1)$	$I(X_2; X_1) - I(X_2; X_1) = 0$
3	1	$I(X_3; \mathbf{C}_1 Y) - I(X_3; \mathbf{C}_1)$	$I(X_3; \mathbf{C}_1) - I(X_3; \mathbf{C}_1) = 0$
4	1	$I(X_4; \mathbf{C}_1 Y) - I(X_4; \mathbf{C}_1)$	$I(X_4; \mathbf{C}_1) - I(X_4; \mathbf{C}_1) = 0$
4	2	$I(X_4; X_3 Y) - I(X_4; X_3)$	$I(X_4; X_3 \mathbf{C}_1) - I(X_4; X_3)$
5	1	$I(X_5; \mathbf{C}_1 Y) - I(X_5; \mathbf{C}_1)$	$I(X_5; \mathbf{C}_1) - I(X_5; \mathbf{C}_1) = 0$
5	2	$I(X_5; \mathbf{C}_2 Y) - I(X_5; \mathbf{C}_2)$	$I(X_5; \mathbf{C}_2 \mathbf{C}_1) - I(X_5; \mathbf{C}_2)$
6	1	$I(X_6; \mathbf{C}_1 Y) - I(X_6; \mathbf{C}_1)$	$I(X_6; \mathbf{C}_1) - I(X_6; \mathbf{C}_1) = 0$
6	2	$I(X_6; \mathbf{C}_2 Y) - I(X_6; \mathbf{C}_2)$	$I(X_6; \mathbf{C}_2 \mathbf{C}_1) - I(X_6; \mathbf{C}_2)$
6	3	$I(X_6; X_5 Y) - I(X_6; X_5)$	$I(X_6; X_5 \mathbf{C}_1, \mathbf{C}_2) - I(X_6; X_5)$

discriminant information due to dependences within \mathbf{C}_1 , \mathbf{C}_2 , and \mathbf{C}_3 , (0^{th} order components), and between X_3 and \mathbf{C}_1 , X_4 and \mathbf{C}_1 , X_5 and \mathbf{C}_1 , X_5 and \mathbf{C}_2 , X_6 and \mathbf{C}_1 , and X_6 and \mathbf{C}_2 (1^{st} order). Hence, (16) requires joint density estimates of up to three features. On the other hand, (8) requires densities of up to six features and is *three orders of magnitude* more complex.

B. Decompositions for low-complexity feature selection

Theorem 1 only holds for the decomposition of \mathbf{X} according to (12)-(13). This is not sufficient for feature selection algorithms, which usually evaluate the MI of various subsets of \mathbf{X} . For this, the theorem must be expanded to *all* possible feature subsets of \mathbf{X} . The extension of the necessary and sufficient condition of (18) to all such subsets is denoted as l -decomposability.

Definition 3: A feature set \mathbf{X} is l -decomposable, or decomposable at order l , if and only if

$$CI(\mathbf{W}; \mathcal{C}_l) = \sum_{k=2}^{|\mathbf{W}|} \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(W_k; \tilde{\mathbf{C}}_{i,k}|\mathbf{C}_1^{i-1}, Y) - I(W_k; \tilde{\mathbf{C}}_{i,k}|Y) \right], \quad \forall \mathbf{W} \in \mathcal{S}(\mathbf{X}) \quad (21)$$

where \mathcal{C}_l and $\tilde{\mathbf{C}}_{i,k}$ are built from \mathbf{W} , as in (12)-(13), and $\mathcal{S}(\mathbf{X})$ is the set of all subsets of \mathbf{X} . Since (18) holds for any feature subset \mathbf{W} of a l -decomposable set \mathbf{X} , simple application of Theorem 1 shows that the same is true for (16).

Corollary 1: Let \mathbf{X} be an l -decomposable feature set, \mathbf{W} a subset of \mathbf{X} , and \mathcal{C}_l a collection of disjoint subsets \mathbf{C}_i of cardinality l built from \mathbf{W} , as in (12)-(13). Then

$$I(\mathbf{W}; Y) = M(\mathbf{W}; Y) + C_{\mathcal{C}_l}(\mathbf{W}; Y) \quad (22)$$

with

$$M(\mathbf{W}; Y) = \sum_{k=1}^{|\mathbf{W}|} I(W_k; Y) \quad (23)$$

$$C_{\mathcal{C}_i}(\mathbf{W}; Y) = \sum_{k=2}^{|\mathbf{W}|} \sum_{i=1}^{\lceil k-1/l \rceil} [I(W_k; \tilde{\mathbf{C}}_{i,k}|Y) - I(W_k; \tilde{\mathbf{C}}_{i,k})], \quad (24)$$

where $\tilde{\mathbf{C}}_{i,k}$ the subset of features in \mathbf{C}_i whose index is smaller than k , and $\mathbf{C}_1^{i-1} = (\mathbf{C}_1, \dots, \mathbf{C}_{i-1})$.

Hence, for an l -decomposable set, it is equivalent to use (2) or (16) as feature selection cost.

Corollary 2: If \mathbf{X} is l -decomposable, then the solution of (1) is the identical to that of

$$\mathcal{X}^* = \arg \max_{\mathcal{X} \in \mathcal{S}} \left\{ \sum_k I(X_k; Y) + \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} [I(X_k; \tilde{\mathbf{C}}_{i,k}|Y) - I(X_k; \tilde{\mathbf{C}}_{i,k})] \right\}. \quad (25)$$

In summary, *the infomax subset of an l -decomposable \mathbf{X} can be computed with density estimates of dimension $l+1$.* When $l = D$ there is only one possibility for \mathcal{C}_l , namely $\mathcal{C}_l = \{\mathbf{X}\}$, and (16) is equal to (8). Hence, all feature sets are at least D -decomposable and, in the worst case, feature selection has exponential complexity in the cardinality of \mathbf{X} . However, depending on the decomposability order of \mathbf{X} , this bound may be very loose. *The intrinsic complexity of feature selection is determined by the decomposability order l of the feature set, not its cardinality.*

V. LOW COMPLEXITY INFOMAX FEATURE SELECTION ALGORITHMS

In this section, we derive a family of infomax feature selection algorithms, based on the theoretical characterization above.

A. A new family of algorithms

When \mathbf{X} is l -decomposable, the infomax space is given by (25). When l -decomposability does not hold, (25) provides a *low-complexity approximation* to the optimal solution. In this case, l is denoted as the *order of the approximation*, and we refer to the true decomposability order as l^* . Since all feature sets are (at least) D -decomposable, the optimal solution can always be attained if (25) is solved for all values of l . This suggests 1) developing a family of algorithms parameterized by l , 2) solving the feature selection problem for all l , and 3) retaining the best solution. Note that, given l , (25) can be solved by existing feature selection strategies. In our implementation, we use the popular (greedy) strategy of sequential forward feature selection [66],

Algorithm 1 (approximate infomax of order l)

Input: feature set $\mathbf{X} = \{X_1, \dots, X_D\}$, order l , and target number of features N .

set $\mathbf{X}^* = \mathbf{C}_1 = \{X_1^*\}$ where $X_1^* = \arg \max_{X_k \in \mathbf{X}} I(X_k; Y)$, $k = 2$, and $i = 1$.

repeat

for $X_r \notin \mathbf{X}^*$ **do**

$$\delta_r = I(X_r; Y) + \sum_{p=1}^{\lceil k-1/l \rceil} \left[I(X_r; \tilde{\mathbf{C}}_{p,k}|Y) - I(X_r; \tilde{\mathbf{C}}_{p,k}) \right], \quad (26)$$

end for

 let $r^* = \arg \max_r \delta_r$.

if $k - 1$ is not a multiple of l **then**

 let $\mathbf{C}_i = \mathbf{C}_i \cup \mathbf{X}_{r^*}$,

else

 set $i = i + 1$, $\mathbf{C}_i = \mathbf{X}_{r^*}$.

end if

 set $\mathbf{X}^* = \cup_i \mathbf{C}_i$, $k = k + 1$,

until $k = N$

Output: \mathbf{X}^* .

which leads² to Algorithm 1. The MIs of (26) are computed with histograms. When b histogram bins are used per feature, the algorithm can be implemented in $O[D(b^l/l)N^2]$ time. Since N is usually small, the complexity is dominated by b and l , increasing exponentially with the latter.

B. Comparison to other infomax methods

The main novelty of Algorithm 1 is the use of (26) as sequential feature selection rule. In addition to the theoretical motivation above, this rule is interesting in two ways. First, it has an intuitive interpretation: it favors features of 1) large MI with the class label, 2) low MI with previously selected features, and 3) large MI with those features given image class. This enforces three principles that are always at play in feature selection:

- 1) *discrimination*: each selected feature must be as discriminant as possible,

²It is worth stressing that the algorithm does not guarantee the best approximation for any l , since the greedy selection of a feature limits the feature groupings of subsequent steps. This is a known limitation of sequential forward selection, e.g., shared by all algorithms of Section II. It can sometimes be circumvented with heuristics such as floating search [66], [73], [74].

- 2) *diversity*: the selected features must not be redundant,
- 3) *reinforcement* : unless this redundancy is, itself, discriminant.

Second, it unifies many algorithms previously proposed for information theoretic feature selection.

In fact, *the first three classes of Section II are special cases of the family now proposed*. Methods in the first class, marginal infomax, only use the first term of (26). Slightly abusing notation, we refer to this as the approximate infomax algorithm of order 0. It enforces the principle of discrimination, but not diversity or reinforcement, and does not guarantee a compact representation: exactly identical features are selected in consecutive steps, wasting some of the available dimensions. The second and third classes are approximations to (26) with $l = 1$, in which case (26) can be written as

$$I(X; Y) + \sum_{i=1}^{k-1} [I(X; X_i^* | Y) - I(X; X_i^*)]. \quad (27)$$

Algorithms in the second class, based on (4), simply discard the terms which account for the discriminant power of feature dependencies ($I(X; X_i^* | Y)$), failing to enforce the principle of reinforcement. This can be overkill, since discriminant dependences can be crucial for fine discrimination between otherwise similar classes. On the other hand, by relying on (7), the algorithms in the third class approximate the summation of (27) by its smallest term.

The excellent empirical performance [48], [58] of algorithms in the second and third classes suggests two hypotheses. The first is that the infomax approximation of first order ($l = 1$) is sufficient for many problems of practical interest. The second is that, even for this approximation, many terms of (27) are neglectable. It is, nevertheless, puzzling that excellent results have been achieved with two very different approximations: the average MI between features (max-relevance min-redundancy (mRMR) method [48]), and the minimum of the differential MI terms [58]. It is also unclear why these would be the only sensible simplifications. Given that both the minimum differential term and the average of the negative terms perform well, why not consider the smallest of the negative terms, their sum (as proposed in [47], [48], [51], [52], [54], [55]), or the median of the differential terms? Table II presents a number of such alternatives to (27), as well as their empirical performance on a set of experiments to be discussed in Section VII.

TABLE II

POSSIBLE ALTERNATIVES TO THE COST OF (27), THEIR RELATION TO THE LITERATURE, AND PERFORMANCE (AVERAGE AND STANDARD DEVIATION OF PRECISION-RECALL AREA) ON EXPERIMENTS OF SECTION VII.

	Cost	feature selection method	PRA
$\Delta(l=0)$	$I(X; Y)$	Marginal infomax	49.7 ± 15.1
$\Delta(l=1)$	$I(X; Y) + \sum_{i=1}^{k-1} [I(X; X_i^* Y) - I(X; X_i^*)]$	approximate infomax order 1	56.3 ± 17.7
$\Delta(l=2)$	$I(X; Y) + \sum_{i=1}^{\lceil (k-1)/2 \rceil} [I(X; \tilde{\mathbf{C}}_{i,k} Y) - I(X; \tilde{\mathbf{C}}_{i,k})]$	approximate infomax order 2	55.2 ± 17.0
δ_{\min}	$I(X; Y) + \min_i [I(X; X_i^* Y) - I(X; X_i^*)]$	method of [58]	54.1 ± 17.6
δ_{med}	$I(X; Y) + \text{median}_i [I(X; X_i^* Y) - I(X; X_i^*)]$		52.9 ± 16.9
δ_{\max}	$I(X; Y) + \max_i [I(X; X_i^* Y) - I(X; X_i^*)]$		52.6 ± 18.3
α_{\min}	$I(X; Y) + \min_i I(X; X_i^* Y)$		49.0 ± 13.5
β_{\min}	$I(X; Y) - \max_i I(X; X_i^*)$		53.5 ± 17.1
β_{avg}	$I(X; Y) - \frac{1}{k-1} \sum_{i=1}^{k-1} I(X; X_i^*)$	mRMR method of [48]	53.4 ± 15.6
α	$I(X; Y) + \sum_{i=1}^{k-1} I(X; X_i^* Y)$		50.2 ± 16.9
β	$I(X; Y) - \sum_{i=1}^{k-1} I(X; X_i^*)$	method of [54] with $\xi = 1$	53.4 ± 15.7

VI. IMAGE STATISTICS AND LOW DECOMPOSABILITY ORDER

In this section, we develop an indirect procedure to validate the hypothesis that band-pass features extracted from natural images have low decomposability order.

A. l -decomposability and image statistics

From Definition 3, \mathbf{X} is l -decomposable if the conjunctive interference (with respect to subsets of cardinality l) within any of its subsets $\mathbf{W} \subset \mathbf{X}$ is non-discriminant. This can be illustrated by returning to Example 1, for which the terms of (15) are the entries in 4th column of Table I. Note that the non-trivially zero entries (identified by boldface k and i) measure how the dependences in \mathbf{C}_2 are affected by \mathbf{C}_1 ($k=4, i=2$); how the dependences in $X_5 \cup \mathbf{C}_2$ are affected by \mathbf{C}_1 ($k=5, i=2$); how the dependences in $X_6 \cup \mathbf{C}_2$ are affected by \mathbf{C}_1 ($k=6, i=2$); and how the dependences in \mathbf{C}_3 are affected by $\mathbf{C}_1 \cup \mathbf{C}_2$ ($k=6, i=3$). $CI(\mathbf{X}; \mathcal{C}_l)$ is the sum of these measures and, for l -decomposability to hold, must not be affected by knowledge of the class Y .

In addition to this, l -decomposability requires (18) to hold for any subset $\mathbf{W} \subset \mathbf{X}$. For example, $\mathbf{W} = (X_1, X_3, X_5, X_6)$ produces a table similar to Table I, with a single non-trivially zero entry, $I(X_6; X_5|X_1, X_3) - I(X_6; X_5)$. l -decomposability requires that the interference of

(X_1, X_3) on the dependence between X_5 and X_6 be non-discriminant. Other subsets of four features give rise to similar constraints on the interference between feature pairs. Hence, in this example, l -decomposability requires all pairwise interferences to be non-discriminant.

In general, l -decomposability holds if and only if the conjunctive interference (with respect to subsets of cardinality l) within any subset \mathbf{W} of \mathbf{X} is not affected by knowledge of the class label Y . As in Fig. 2, *this does not mean that conjunctive interference is non-existent, but simply that it does not change across classes*. Overall, the sufficient condition for l -decomposability is similar to the sufficient condition for the optimality of marginal infomax. While, in that case, image statistics must satisfy $C(\mathbf{X}; Y) = 0$, i.e. that *no* dependences in \mathbf{X} are discriminant, in this case the constraints only affect 2^{nd} order *subset* dependences: *l -decomposability does not impose constraints on subset dependencies of 0^{th} or 1^{st} order, and neither does it impose that there are no 2^{nd} order subset dependencies. It only requires these dependencies to be such that the conjunctive interference $CI(\mathbf{X}; \mathcal{C}_l)$ is non-discriminant*. This is much less restrictive than what is required for the optimality of marginal infomax. As in that case, the consistency of the statistics of Fig. 1 a) suggests that, for natural images, the hypothesis that l -decomposability holds for small l is not unreasonable. We next turn to the problem of determining this value.

B. Indirect validation of the low-order decomposability hypothesis

If \mathbf{X} is l^* -decomposable, the infomax feature set can be found with (25), using $l = l^*$. For approximation orders $l \neq l^*$, the problems of (25) make looser assumptions about feature dependences as l increases. $l = 0$ assumes that no feature dependences are discriminant, $l = 1$ that only dependences within feature pairs are important, and so forth, up to $l = D$ where all dependences are accounted for. The decomposability order of the feature set can be determined with recourse to the indirect validation paradigm: the error of classifiers designed on the spaces produced by (25) is expected to decrease with l , leveling off at $l = l^*$. If this produces a consistent estimate of l^* across a number of classification problems, there is strong empirical evidence that \mathbf{X} is l^* -decomposable. If this is repeatedly observed for transformations in a certain class, e.g. wavelets, there is strong evidence that all feature sets in the class are l^* -decomposable.

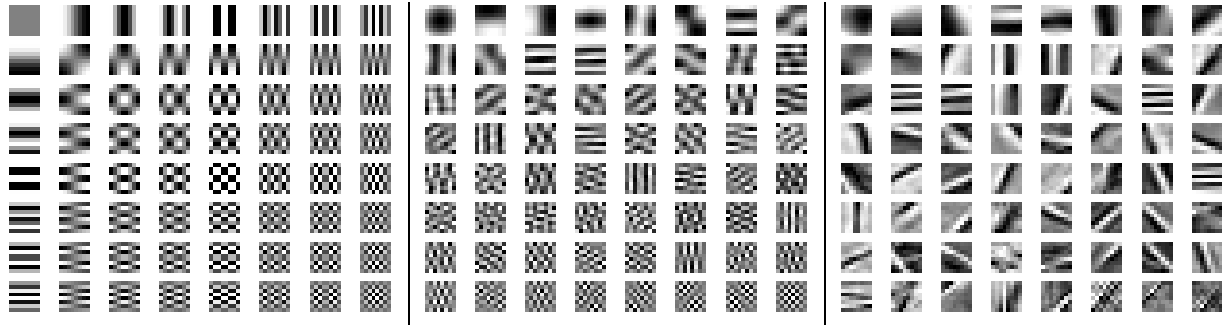


Fig. 3. Basis functions for DCT (left), PCA (center) and ICA (right).

VII. EXPERIMENTS

In this work, we hypothesize that transformations into sets of bandpass frequency coefficients have low decomposability order. We rely on indirect validation to test this hypothesis.

A. Experimental protocol

All experiments were performed with the Brodatz and Corel image databases. Brodatz is a standard benchmark for texture classification under controlled imaging conditions, and with no distractors. Corel is a standard evaluation set for recognition from unconstrained scenes (e.g. no control over lighting or object pose, cluttered backgrounds). Brodatz contains sets of 9 patches from 112 gray-scale textures, in a total of 1008 images. One patch of each texture was used for testing and the remaining 8 for training. From Corel, we selected 15 image classes³ each containing 100 color images. Train and test sets were then created by assigning each image to the test set with probability 0.2. Evaluation was based on precision and recall (PR), using the test images as queries to a database containing the training set. The PR curve was summarized by its integral, the PR area (PRA). In all experiments, feature vectors were extracted from localized image neighborhoods and classification based on (30) with Gauss mixture class-conditional densities. A Gauss mixture with a fixed number of components was learned for each image (results were qualitatively similar for various values, we report on 8 components), examples were assumed independent in (30), and class priors uniform. Four transformations

³“Arabian horses”, “Auto racing”, “Owls”, “Roses”, “Ski scenes”, “religious stained glass”, “sunsets and sunrises”, “coasts”, “Divers and diving”, “Land of the pyramids” (pictures of Egypt), “English country gardens”, “fireworks”, “Glaciers and mountains”, “Mayan and Aztec ruins”, and “Oil Paintings”.

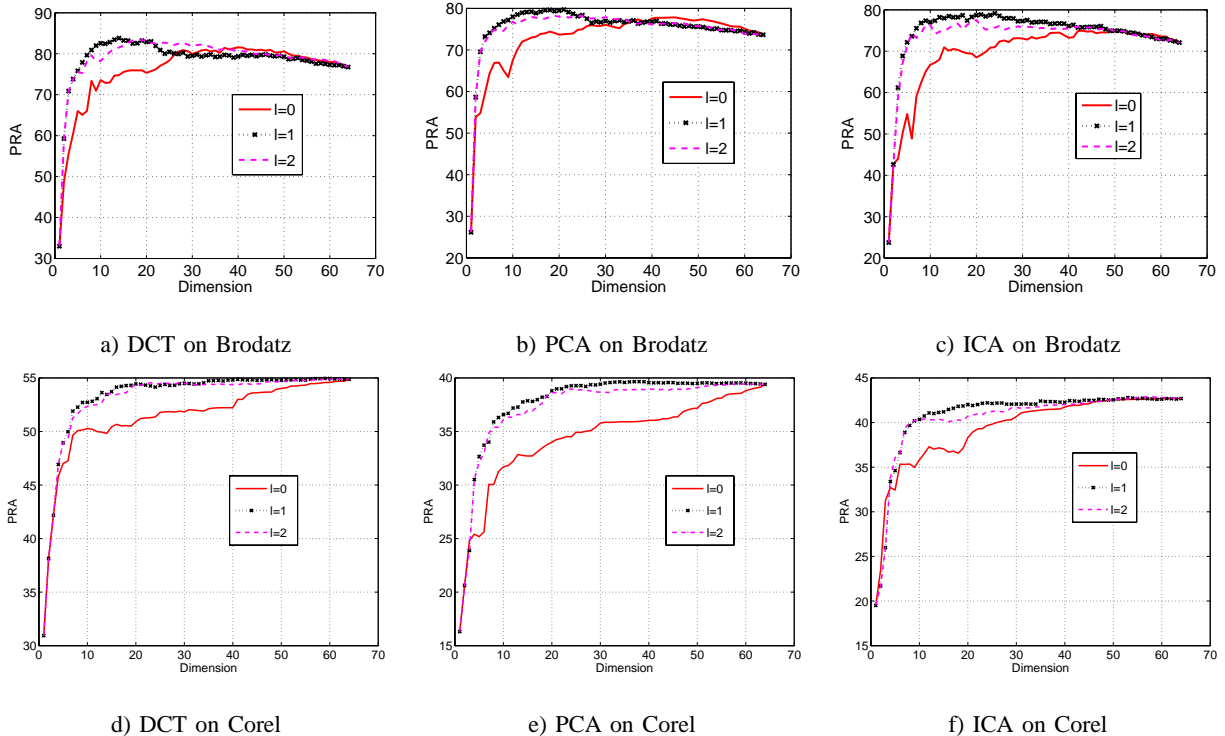


Fig. 4. PRA as a function of the number of features selected by approximate infomax $0 \leq l \leq 2$, for the DCT, PCA, and ICA feature sets, on Brodatz and Corel.

were considered: the discrete cosine transform (DCT), principal component analysis (PCA), independent component analysis (ICA), and a wavelet representation (WAV). The feature space had $D = 64$ per color channel (three layers of wavelet decomposition and 8×8 image blocks) and the observations were extracted with a sliding window. PCA and ICA were learned from 100,000 random training examples. Fig. 3 compares the basis functions learned on Brodatz with those of the DCT.

B. Decomposability order

The decomposability order of all datasets was studied with the indirect validation paradigm. Because the computational cost is exponential on the approximation order l , it is (at this point in time) only feasible to consider small values of this parameter. We have limited all experiments to the range $0 \leq l \leq 2$. Fig. 4 presents the PRA curves obtained with different l for DCT,

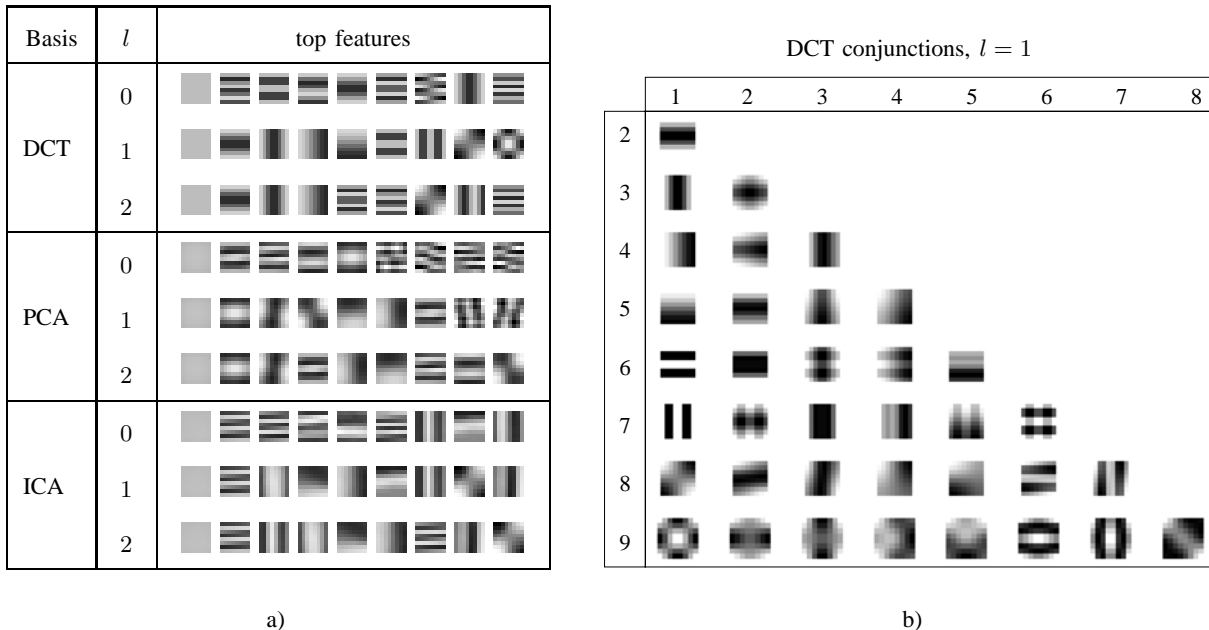


Fig. 5. a) Top nine features (in decreasing order, from left to right) selected on Brodatz for the three representations and $0 \leq l \leq 2$. b) Conjunctions of features that contribute to (16) for the optimal feature set on Brodatz with DCT features and $l = 1$. The basis function at row i and column j of the table was produced by averaging features i and j of the optimal set of a).

PCA, and ICA⁴. The most striking observation is that, for *all* databases and transformations, $l = 1$ is superior to $l = 0$, but there is no advantage of $l = 2$ over $l = 1$. The constancy of this result suggests that *all feature sets are 1-decomposable*. To understand this constancy, we analyzed the feature rankings in detail. Fig. 5 a) presents the top 9 features selected, on Brodatz, for each transformation and value of l . For $l = 0$, the top features are nearly identical: all have very high frequency and do not appear to capture perceptually interesting image structure. This indicates that marginal statistics are not enough for these problems. The solution obtained with $l = 1$ is superior: not only the features appear to be detectors of perceptually relevant image attributes but the same holds for their pair-wise conjunctions. This is shown in Fig. 5 b), which presents the optimal pairwise DCT conjunctions. While individual features detect smooth regions, horizontal and vertical bars, horizontal and vertical edges, horizontal and vertical parallel segments, corners, and rounded spots, the set of conjunctions includes detectors of crosses, T-

⁴Qualitatively identical results were obtained with the wavelet and are omitted for brevity.

and L-junctions, grids, oriented lines, etc⁵. The fact that, for $l = 1$, *features are selected not only by individual discriminant power, but also by the discriminant power of pairwise conjunctions, makes a significant difference for both classification accuracy (Fig. 4) and perceptual relevance of the visual attributes they detect (Fig. 5)*. Finally, there is no benefit in considering $l = 2$: both classification performance and perceptual relevance decrease slightly. Because the union of individual features and pairwise conjunctions is very discriminant, the gain of triplets is small. On the other hand, all dimensionality problems (complexity of density estimation, exponential increase in training data requirements) are compounded, and the overall result is a loss.

C. Comparison to previous methods

The classification performance of (25) with $0 \leq l \leq 2$ - costs $\Delta(l)$ - was compared to that of each of the other costs in Table II. For each transformation and image database, classification performance was summarized by the average PRA of the first N features. N was chosen to guarantee that the number of features needed for optimal performance was available, but not a lot more (all methods perform similarly when N is close to the total number of features). Using Fig. 4 for guidance, we chose $N = 15$ for Brodatz, and $N = 20$ for Corel. Table II presents the average and standard deviation of the PRA achieved (across datasets and transforms) by each cost. To facilitate the comparison we divided (for each dataset and transform) the average PRA of each cost by that achieved with $\Delta(l = 1)$. The average of this measure (across datasets and transforms) is denoted as the normalized average PRA (NAPRA) score of the cost. Fig. 6 presents a boxplot of the NAPRA score of each cost, across databases and transformations. A number of interesting observations can be made. The first is that $\Delta(l = 1)$ produced the best feature set in *all* cases. The second overall best performer was $\Delta(l = 2)$, followed by the three costs previously proposed in the literature: δ_{\min} [58], β_{avg} [48], and β [54]. On average, there was no substantial difference between these three costs, although δ_{\min} performed best. The fact that these are the best approximations to $\Delta(l = 1)$ (among those we evaluated) is a possible explanation for their impressive performance in previous experimental comparisons [48], [58]. Of the remaining costs, δ_{median} and δ_{\max} performed somewhat worse, but clearly above marginal infomax ($\Delta(l = 0)$), while α and α_{\min} did not consistently beat the latter.

⁵In fact, the set of conjunctions is much larger than that shown. While the table only includes pair-wise feature *averages*, the set includes all functions of the same feature pairs.

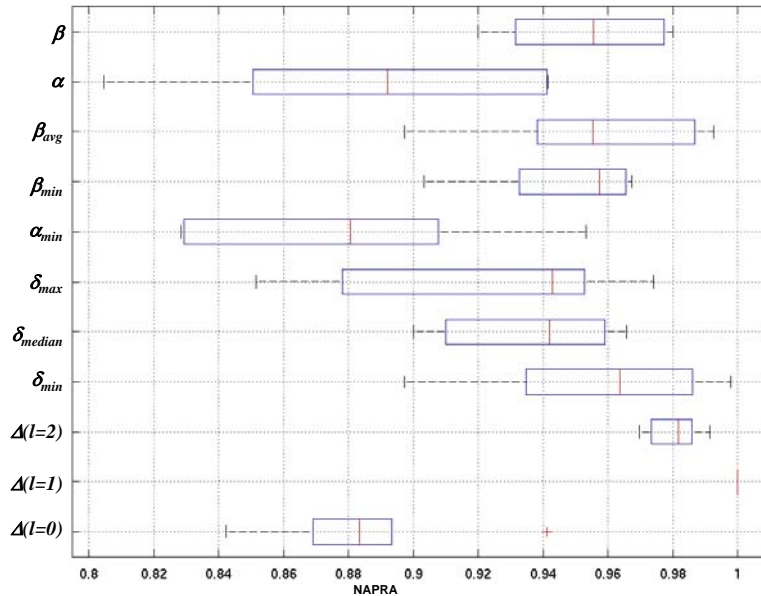


Fig. 6. NAPRA scores for the different costs of Table II, across feature transforms and databases. Box lines indicate lower, median, and upper quartile values. Dashed lines show the extent of the rest of the data.

Returning to the indirect validation paradigm, these results provide information about the importance, for discrimination, of various aspects of the feature statistics. The first interesting observation is the *average performance of marginal infomax: close to 90% of the best*. This suggests that, for natural images, *most discriminant information is contained in marginal feature statistics*. Given that marginal infomax is the *only* method which does not require joint density estimates, *it may be the best solution for recognition problems with strict constraints on time or computation*. It is also interesting to investigate which terms of $\Delta(l=1)$ are responsible for its performance gain over $\Delta(l=0)$. One observation from Fig. 6 is that this gain is very non-linear on the differential terms $\delta_i = I(X; X_i^*|Y) - I(X; X_i^*)$. In particular, the inclusion of a single term, be it the largest (δ_{\max}), median (δ_{median}) or most negative (δ_{\min}), is sufficient to achieve at least half of the total gain, with δ_{\min} achieving 2/3. Hence, while it is important to include one differential term, *the exact choice may not be very important*. This flexibility could be significant when there are complexity constraints. While computing an arbitrary δ_i has linear complexity on the number of features, the search for the best term has quadratic complexity. It follows that *the inclusion of an arbitrary differential term may be a good intermediate solution* (complexity quadratic in histogram bins but linear on features) *between marginal infomax and approximate*

infomax of order 1. On the other hand, finding the best δ_i [58] requires more computation than evaluating $\Delta(l = 1)$ (due to the search after all terms are computed) and has no advantage.

As an alternative to the differential terms δ_i , Fig. 6 shows that gains can also be obtained by adding terms of each mutual information type - $\alpha_i = I(X; X_i^*|Y)$ and $\beta_i = I(X; X_i^*)$ - to $\Delta(l = 0)$. Here, it appears that the β_i are much more important than the α_i : by themselves, the α_i do not even produce a consistent improvement over marginal infomax. On the other hand, the inclusion of the best β_i (cost β_{\min}), does not perform nearly as well as the inclusion of the best δ_i (cost δ_{\min}). In fact, the latter performed better than all the α -only, or β -only approaches considered. Yet, the gains of the β -only costs could, once again, be interesting if there are complexity constraints. Note that, unlike the α terms, they do not depend explicitly on the class Y . They could, thus, be learned from a generic collection of natural images, *independently of the particular recognition problem under consideration*. In this case, the *complexity of the β costs would be equivalent to that of marginal infomax!* While it is currently unclear if the performance would remain as in Fig. 6 (where all β_i were estimated from the training sets used for classifier design), this is an interesting topic for further research.

As discussed in Section II, there is a large literature on β costs, mostly focusing on the role of the parameter ξ of (5) [47], [48], [51]–[55]. Fig. 6 suggests that, for natural images, this discussion is inconsequential: similar performance was obtained with only one β_i (cost β_{\min}), their average (cost β_{avg}), or sum (cost β). Different ξ only affected the variance of the NAPRA score, which was smallest for $\xi = 1$. The increased variance of the other weights might explain various, sometimes conflicting, claims for their success [47], [48], [51], [52], [54], [55].

In summary, the infomax approximation of order 1 ($\Delta(l = 1)$) outperforms the previous low complexity methods. It is worth emphasizing that the discussion above is based on the *average* performance of the different costs, across datasets and transformations. One important point is that *all previous methods exhibited “breakdown” modes*, i.e. combinations of transformation/database on which their performance was well below average. This can be seen from the limit intervals (dashed lines) of Fig. 6. In almost all cases, the lower bound is close to the average performance of marginal infomax. The only salient exceptions are $\Delta(l = 1)$, which always performed best, and $\Delta(l = 2)$, which has small variance. These observations suggest that *the main role of the summation in (26) is to assure robustness. While simplifications of this rule can perform well for certain datasets, they compromise generalization.*

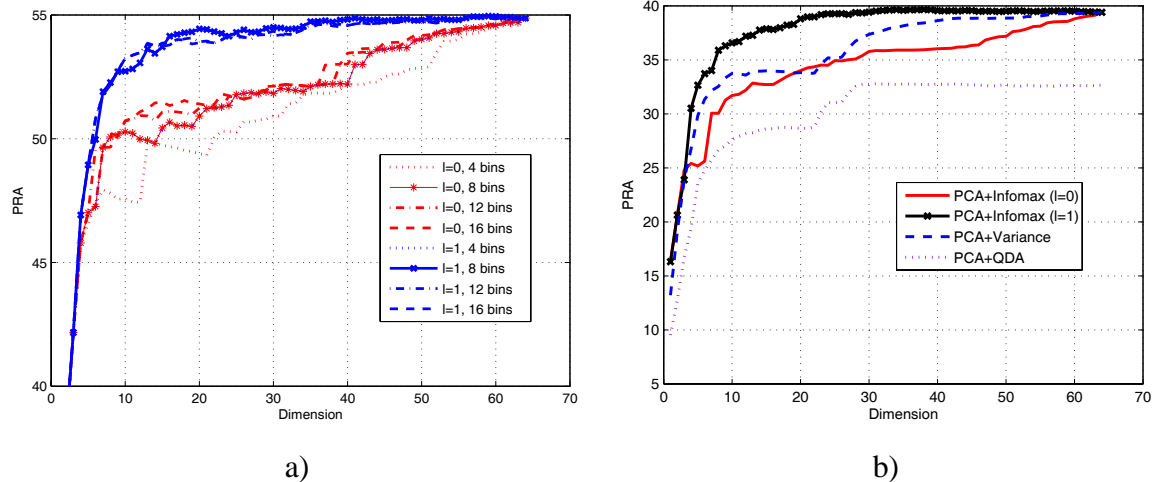


Fig. 7. a): PRA curves for the DCT on Corel using $\Delta(l = i), i \in \{0, 1\}$, and various numbers b of histogram bins. b) comparison of the PRA curves obtained with infomax and popular methods of equivalent complexity, for PCA.

D. Robustness

Assuming that bandpass transforms are indeed 1-decomposable, we performed some experiments to determine the robustness to the parameter that determines the complexity: the number of histogram bins/axis b (recall that the complexity of approximate infomax of order l is $O(b^l)$). In particular, we repeated the experiment with $l = 0$ and $l = 1$ for values of b in $[4, 16]$. Fig. 7 a) presents PRA curves from Corel, with DCT features⁶, showing that recognition accuracy is quite insensitive to this parameter. For both values of l , 8 bins are sufficient to achieve accuracy very close to the highest. A loss only occurs for $b = 4$ and, as expected, is more significant for $l = 1$, where the density estimates are two-dimensional.

E. Comparison with scalable feature selection methods

To place the results above in a larger feature selection context, we compared the infomax algorithms with two widely popular methods of *similar complexity*: PCA, and its combination with quadratic discriminant analysis [75] (PCA+QDA). Because these methods project all examples onto the PCA subspace, we restricted the comparison to the infomax subset of PCA features. Although PCA is frequently combined with the Euclidean or Mahalanobis distances and a nearest neighbor classifier, the popular ‘‘Eigenfaces’’ technique [76], preliminary experiments

⁶Similar results were obtained on Brodatz and are omitted.

showed better performance for a Gauss mixture classifier on the PCA subspace. This is identical to the classifier adopted for the infomax features, but relies on feature ranking by variance, rather than MI. PCA+QDA is an extension of the popular ‘‘Fisherfaces’’ method [77], and equivalent to (30) when the PCA coefficients are Gaussian. It was implemented by fitting a multivariate Gaussian to each training image, and using (30) to classify all test images.

Fig. 7 b) compares, on Corel, the PRA curves of PCA+Variance and PCA+QDA, with those previously presented for infomax ($l \in \{0, 1\}$) on PCA space⁷. PCA+QDA performs significantly worse than all other approaches. This is not surprising, given the strong non-Gaussianity of the distributions of Fig. 1 a). With the Gaussian mixture classifier, maximum variance and marginal infomax have similar performance⁸, but infomax with $l = 1$ is substantially better. For example, in the PCA case, energy compaction requires about 30 features to reach the accuracy that infomax ($l = 1$) achieves with only 10! For the DCT, the ratio is even larger, closer to 4/1. Visual inspection of recognition results shows significant improvement for queries from classes that share visual attributes with other classes in the database (see [78] for examples).

VIII. DISCUSSION

We have studied the hypothesis that high-order dependences of bandpass features contain little information for image discrimination. The hypothesis was characterized formally, by introduction of the concepts of conjunctive interference and decomposability order, and derivation of necessary and sufficient conditions for the feasibility of low-complexity feature selection, in terms of these concepts. It was shown that the intrinsic complexity of feature selection is determined by the decomposability order of the feature set: the infomax subset of an l -decomposable set can be computed with density estimates of dimension $l+1$. A family of (locally) optimal feature selection algorithms was then derived for all levels of complexity, and its performance characterized in two ways. Theoretically, it was shown that various previous information theoretic feature selection algorithms are approximations to the ones now derived. Experimentally, the latter were shown to consistently outperform the former, for a diverse set of images and feature transformations.

⁷Once again, similar results were obtained on Brodatz and are omitted.

⁸While maximum variance is somewhat superior to marginal infomax in Fig. 7 b), we have seen no consistent differences between the two criteria across all feature spaces.

Following the indirect validation paradigm, the new feature selection algorithms were used to objectively test the hypothesis of low decomposability order for natural image features. This has shown that, while there is a non-negligible classification gain in modeling feature dependencies (in all cases $l = 1$ outperformed $l = 0$), this gain has diminishing returns. For certain, the benefits of modeling dependencies between triplets ($l = 2$) over pairs ($l = 1$) are, at most, marginal. While it is possible that there may be some $l > 2$ with substantially better performance than $l = 1$, the consistent lack of improvement from $l = 1$ to $l = 2$, across the imagery and features considered in our experiments, suggests that this is unlikely. Unfortunately, limitations in computation and database size currently prevent us from experimenting with $l > 2$.

A detailed investigation of the $l = 1$ case has shown that, when pairwise dependences are modeled, the gains are very nonlinear on the rigor of this modeling. In particular, simple modeling of marginal statistics performs fairly well (within 90% of the top performance) and the inclusion of a single pairwise differential term, as proposed in [58], can capture as much as 2/3 of what remains. On the other hand, the simple inclusion of so-called β terms, as proposed in [47], [48], [51], [52], [54], [55], can also work well. Since β terms do not depend on the particular classification problem under analysis, they could conceivably be learned from a generic image database. In this case, it should be possible to account for dependences with feature selection algorithms that only require the estimation of marginal densities. This remains an interesting topic for further research. The main benefit of accounting for all terms of order 1 seems to be a significant increase in robustness. While the previously proposed approximations can perform very well in some cases, and reasonably well on average, they have all exhibited “breakdown” modes (combinations of features and image databases where performance was similar to that of marginal statistics). The large variance of their classification performance could explain previous conflicting claims for the superiority of different approximations [47], [48], [51], [52], [54], [55]. On the other hand, the algorithms now proposed performed very robustly, consistently achieving the best results on all datasets. It would, therefore, be speculative to 1) propose that some of the terms of the simplified MI of order 1 are more important than others, or 2) make generic statements about the details of the dependence structure encoded by these terms.

What can, thus, be said about the structure of the dependencies of bandpass features extracted from natural images? 1-decomposability means that, for natural images, the conjunctive interference between individual features is not discriminant. Or in other words, that *the effect, on the*

dependence between two features, of observing any other feature is constant across image classes. This is a significantly more precise statement than the hypothesis, that feature dependences are constant across classes, with which we started. Although our analysis is limited to features extracted from natural images, this conclusion also appears sensible for modalities such as audio, speech, or language. For example, in the language context, it would imply that the effect of observing a word, on the dependence between two other words, is constant across document classes. This simply suggests that second order dependences between words are determined by language, not the specific document classes. It is a fairly mild constraint on the structure of text, e.g. much milder than the common bag-of-words model. It is interesting to note that some experimental observations similar to the ones that we report for images have been made for text categorization. These include reports of successful application of marginal infomax [41], [42], and reports of improved performance by the criterion of (6) [44].

APPENDIX I

OPTIMALITY CRITERIA FOR FEATURE SELECTION

In the most general sense, the optimal feature space for a classification problem $\mathcal{C}_{\mathcal{X}}$ is

$$\mathcal{X}^* = \arg \min_{\mathcal{X} \in \mathcal{S}} J(\mathcal{C}_{\mathcal{X}}). \quad (28)$$

where $J(\cdot)$ is a cost, and \mathcal{S} the set of range spaces for the transforms under consideration.

A. Minimum Bayes error features

One measure of goodness of $\mathcal{C}_{\mathcal{X}}$ is the lowest possible probability of error achievable in \mathcal{X} , usually referred to as the *Bayes error* [33]

$$L_{\mathcal{X}}^* = 1 - E_{\mathbf{x}}[\max_i P_{Y|\mathbf{X}}(i|\mathbf{x})], \quad (29)$$

where $E_{\mathbf{x}}$ is the expectation with respect to $P_{\mathbf{X}}(\mathbf{x})$. It depends only on \mathcal{X} , not the classifier itself, and there is at least one classifier that achieves this bound, the *Bayes* decision rule

$$g^*(\mathbf{x}) = \arg \max_i P_{Y|\mathbf{X}}(i|\mathbf{x}). \quad (30)$$

While it is natural to define \mathcal{X}^* as the space of *minimum Bayes error*, it has long been known that the resulting optimization can be difficult. For example, sequential feature selection is not easy in this setting since the $\max(\cdot)$ non-linearity of (29) makes it impossible to decompose the

new cost ($E_{\mathbf{X}_n}[\max_i P_{Y|\mathbf{X}_n}(i|\mathbf{x}_n)]$) as a function of the previous best ($E_{\mathbf{X}_c}[\max_i P_{Y|\mathbf{X}_c}(i|\mathbf{x}_c)]$), and a function of the candidate set \mathbf{X}_a , where \mathbf{X}_c is the best current subset and $\mathbf{X}_n = (\mathbf{X}_a, \mathbf{X}_c)$.

B. Infomax features

The infomax formulation has a number of appealing properties.

Lemma 3: Let $\langle f(i) \rangle_Y = \sum_i P_Y(i) f(i)$, and $KL[p||q] = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$ be the relative entropy between p and q , with integrals replaced by summations for discrete random variables.

The following properties hold for the MI, as defined in (2).

- 1) for any two random vectors \mathbf{X} and \mathbf{Z} , $I(\mathbf{X}; \mathbf{Z}) \geq 0$, with equality if and only if \mathbf{X} and \mathbf{Z} are statistically independent. Furthermore $I(\mathbf{X}; \mathbf{Z}) = KL[P_{\mathbf{X}, \mathbf{Z}}(\mathbf{x}, \mathbf{z}) || P_{\mathbf{X}}(\mathbf{x}) P_{\mathbf{Z}}(\mathbf{z})]$.
- 2) $I(\mathbf{X}; Y) = \langle KL [P_{\mathbf{X}|Y}(\mathbf{x}|i) || P_{\mathbf{X}}(\mathbf{x})] \rangle_Y$
- 3) $I(\mathbf{X}; Y) = H(Y) - H(Y|\mathbf{X})$, where $H(Y) = -\langle \log P_Y(i) \rangle_Y$ is the entropy of Y , and $H(Y|\mathbf{X}) = -E_{\mathbf{X}} [\langle \log P_{Y|\mathbf{X}}(i|\mathbf{x}) \rangle_{Y|\mathbf{X}}]$ the posterior entropy of Y given \mathbf{X} .
- 4) If $\mathbf{X}_{1,k} = \{X_1, \dots, X_k\}$, then

$$I(\mathbf{X}_{1,k}; Y) - I(\mathbf{X}_{1,k-1}; Y) = I(X_k; Y|\mathbf{X}_{1,k-1}), \quad (31)$$

where

$$I(X; Y|\mathbf{Z}) = \sum_i \int P_{X,Y;\mathbf{Z}}(x, i, \mathbf{z}) \log \frac{P_{X,Y|\mathbf{Z}}(x, i|\mathbf{z})}{P_{X|\mathbf{Z}}(x|\mathbf{z}) P_{Y|\mathbf{Z}}(i|\mathbf{z})} dx d\mathbf{z}$$

Proof: All proofs are either available in [79], or straightforward consequences of (2). ■

Property 1), and the fact that the $KL[p||q]$ is a measure of similarity between distributions p and q , show that $I(\mathbf{X}; \mathbf{Z})$ is a measure of dependence between \mathbf{X} and \mathbf{Z} . For this reason, we frequently refer to $I(\mathbf{X}; \mathbf{Z})$ as the *dependence* between \mathbf{X} and \mathbf{Z} . Property 2) implies that

$$\mathcal{X}^* = \arg \max_{\mathcal{X} \in \mathcal{S}} \langle KL [P_{\mathbf{X}|Y}(\mathbf{x}|i) || P_{\mathbf{X}}(\mathbf{x})] \rangle_Y \quad (32)$$

i.e. that infomax feature selection is inherently discriminant: it rewards spaces where the class densities are on average well separated from the mean density. This is a sensible way to quantify the intuition that optimal discriminant transforms are those that best separate the different classes.

From Property 3), it follows that $\mathcal{X}^* = \arg \min_{\mathcal{X} \in \mathcal{S}} H(Y|\mathbf{X})$. Since entropy is a measure of uncertainty, this implies that the infomax space minimizes the uncertainty about which class is responsible for the observed features. It also establishes a formal connection to the minimization

of Bayes error since, in both cases, the optimal space is

$$\mathcal{X}^* = \arg \max_{\mathcal{X} \in \mathcal{S}} E_{\mathbf{X}} [\phi (P_{Y|\mathbf{X}}(1|\mathbf{X}), \dots, P_{Y|\mathbf{X}}(M|\mathbf{X}))],$$

where $\phi(p_1, \dots, p_M)$ is one of two functions, $\max(p_i)$ and $\langle \log p_i \rangle$, which are both convex and have a number of similar properties (including co-located maxima and minima in the unconstrained probability simplex, and interesting relationships between gradients). In fact, there are a number of problems for which the two optimal solutions are identical [62], [80]. Property 4) probably has the greatest practical significance, and justifies the adoption of infomax over the minimization of Bayes error. It enables *modular* decompositions of the MI, which are central to the efficient implementation of sequential search, and intuitive. In particular, if \mathbf{X}^* is the current set of selected features, it shows that the feature to be selected at the next step should be

$$X^* = \arg \max_{k|X_k \notin \mathbf{X}^*} I(X_k; Y | \mathbf{X}^*), \quad (33)$$

i.e. the one that most reduces the uncertainty about Y , given \mathbf{X}^* . This implies that X^* should 1) be discriminant and 2) have small redundancy with previously selected features.

APPENDIX II

PROOF OF LEMMA 1

From the chain rule of MI [79], $I(\mathbf{X}, Y) = \sum_{k=1}^D I(X_k; Y | \mathbf{X}_{1,k-1})$. Using the equality

$$\begin{aligned} I(X; Y | \mathbf{Z}) &= E_{X,Y,\mathbf{Z}} \left[\log \frac{P_{X,Y|\mathbf{Z}}(x, y | \mathbf{z})}{P_{X|\mathbf{Z}}(x | \mathbf{z}) P_{Y|\mathbf{Z}}(y | \mathbf{z})} \right] \\ &= E_{X,Y,\mathbf{Z}} \left[\log \frac{P_{X,Y}(x, y)}{P_X(x) P_Y(y)} + \log \frac{P_{X,Y|\mathbf{Z}}(x, y | \mathbf{z}) P_Y(y)}{P_{X,Y}(x, y) P_{Y|\mathbf{Z}}(y | \mathbf{z})} + \log \frac{P_X(x)}{P_{X|\mathbf{Z}}(x | \mathbf{z})} \right] \\ &= I(X; Y) + E_{X,Y,\mathbf{Z}} \left[\log \frac{P_{X|Y,\mathbf{Z}}(x | y, \mathbf{z})}{P_{X|Y}(x | y)} \right] - I(X; \mathbf{Z}) \\ &= I(X; Y) + E_{X,Y,\mathbf{Z}} \left[\log \frac{P_{X,\mathbf{Z}|Y}(x, \mathbf{z} | y)}{P_{X|Y}(x | y) P_{\mathbf{Z}|Y}(\mathbf{z} | y)} \right] - I(X; \mathbf{Z}) \\ &= I(X; Y) + I(X; \mathbf{Z} | Y) - I(X; \mathbf{Z}) \end{aligned} \quad (34)$$

with $X = X_k$ and $\mathbf{Z} = \mathbf{X}_{1,k-1}$, leads to

$$I(\mathbf{X}, Y) = \sum_{k=1}^D I(X_k; Y) - \sum_{k=2}^D [I(X_k; \mathbf{X}_{1,k-1}) - I(X_k; \mathbf{X}_{1,k-1} | Y)]$$

and the lemma follows.

APPENDIX III
PROOF OF LEMMA 2

By recursive application of the chain rule of mutual information

$$\begin{aligned}
I(X_k; \mathbf{X}_{1,k-1}|Y) &= I(X_k; \mathbf{C}_1, \dots, \tilde{\mathbf{C}}_{\lceil k-1/l \rceil, k}|Y) \\
&= I(X_k; \tilde{\mathbf{C}}_{\lceil k-1/l \rceil, k} | \mathbf{C}_1, \dots, \mathbf{C}_{\lceil k-1/l \rceil-1}, Y) + I(X_k; \mathbf{C}_1, \dots, \mathbf{C}_{\lceil k-1/l \rceil-1} | Y) \\
&= \sum_{i=1}^{\lceil k-1/l \rceil} I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1, \dots, \mathbf{C}_{i-1}, Y) = \sum_{i=1}^{\lceil k-1/l \rceil} I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}, Y)
\end{aligned}$$

where $\mathbf{C}_1^k = \{\mathbf{C}_1, \dots, \mathbf{C}_k\}$. Similarly,

$$I(X_k; \mathbf{X}_{1,k-1}) = \sum_{i=1}^{\lceil k-1/l \rceil} I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}).$$

The lemma follows from (9).

APPENDIX IV
PROOF OF THEOREM 1

Combining Lemmas 1 and 2,

$$\begin{aligned}
I(\mathbf{X}; Y) &= \sum_{k=1}^D I(X_k; Y) + \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}, Y) - I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}) \right] \\
&= \sum_{k=1}^D I(X_k; Y) - \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k}) - I(X_k; \tilde{\mathbf{C}}_{i,k} | Y) \right] \\
&\quad + \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}, Y) - I(X_k; \tilde{\mathbf{C}}_{i,k} | Y) \right] \\
&\quad - \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}) - I(X_k; \tilde{\mathbf{C}}_{i,k}) \right].
\end{aligned}$$

It follows that

$$I(\mathbf{X}; Y) = \sum_{k=1}^D I(X_k; Y) + \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | Y) - I(X_k; \tilde{\mathbf{C}}_{i,k}) \right]$$

if and only if

$$\sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}) - I(X_k; \tilde{\mathbf{C}}_{i,k}) \right] = \sum_{k=2}^D \sum_{i=1}^{\lceil k-1/l \rceil} \left[I(X_k; \tilde{\mathbf{C}}_{i,k} | \mathbf{C}_1^{i-1}, Y) - I(X_k; \tilde{\mathbf{C}}_{i,k} | Y) \right]$$

and the theorem follows from the definition of $CI(\mathbf{X}; \mathcal{C}_l)$ in (15).

ACKNOWLEDGMENTS

The authors thank the comments of three anonymous reviewers, which substantially improved clarity of the presentation. This work was funded by NSF awards IIS-0448609 and-0534985.

REFERENCES

- [1] R. Clarke, *Transform Coding of Images*. Academic Press, 1985.
- [2] S. Mallat, "A Theory for Multiresolution Signal Decomposition: the Wavelet Representation," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. Vol. 11, pp. 674–693, July 1989.
- [3] D. Ruderman, "The statistics of natural images," *Network: Computation in Neural Systems*, vol. 5, no. 4, pp. 517–548, 1994.
- [4] D. Field, "Relations between the statistics of natural images and the response properties of cortical cells," *J. Opt. Soc. Am. A*, vol. 4, no. 12, pp. 2379–2394, 1987.
- [5] —, "What is the goal of sensory coding?" *Neural Computation*, vol. 6, no. 4, pp. 559–601, January 1989.
- [6] R. Buccigrossi and E. Simoncelli, "Image compression via joint statistical characterization in the wavelet domain," *IEEE Transactions on Image Processing*, vol. 8, pp. 1688–1701, 1999.
- [7] J. Huang and D. Mumford, "Statistics of Natural Images and Models," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Fort Collins, Colorado*, 1999.
- [8] E. Simoncelli and B. Olshausen, "Natural image statistics and neural representation," *Annual Review of Neuroscience*, vol. 24, pp. 1193–1216, 2001.
- [9] A. Torralba and A. Oliva, "Depth estimation from image structure," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 9, pp. 1226–1238, 2002.
- [10] A. Srivastava, A. Lee, E. Simoncelli, and S. Zhu, "On advances in statistical modeling of natural images," *Journal of Mathematical Imaging and Vision*, vol. 18, pp. 17–33, 2003.
- [11] F. Long and D. Purves, "Natural scene statistics as the universal basis of color context effects." *Proceedings of the National Academy of Sciences of the United States*, vol. 100, no. 25, pp. 15 190–15 193, 2003.
- [12] B. Olshausen and D. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *Nature*, vol. 381, pp. 607–609, 1996.
- [13] A. Bell and T. Sejnowski, "The independent components of natural scenes are edge filters," *Vision Research*, vol. 37, no. 23, pp. 3327–3328, December 1997.
- [14] J. H. van Hateren and D. L. Ruderman, "Independent component analysis of natural image sequences yields spatiotemporal filters similar to simple cells in primary visual cortex," in *Proc. Royal Society ser. B*, vol. 265, 1998, pp. 2315–2320.
- [15] J. Portilla, V. Strela, M. Wainwright, and E. Simoncelli, "Image Denoising using Scale Mixtures of Gaussians in the Wavelet Domain," *IEEE Trans. on Image Processing*, vol. 12, no. 11, pp. 1338–1351, November 2003.
- [16] P. Moulin and L. Juan, "Analysis of multiresolution image denoising schemes using Generalized Gaussian and complexity priors," *IEEE Trans. on Information Theory*, vol. Vol. 45, pp. 909–919, April 1999.
- [17] A. Levin, A. Zomet, and Y. Weiss, "Learning how to inpaint from global image statistics," *International Conference on Computer Vision*, pp. 305–312, 2003.
- [18] S. Roth and M. Black, "Fields of experts: A framework for learning image priors," *IEEE Conf. on Computer Vision and Pattern Recognition*, vol. 2, pp. 860–867, 2005.

- [19] N. Farvardin and J. Modestino, "Optimum Quantizer Performance for a Class of Non-Gaussian Memoryless Sources," *IEEE Trans. on Information Theory*, May 1984.
- [20] Y. Weiss, "Deriving intrinsic images from image sequences," *International Conference on Computer Vision*, vol. 2, pp. 68–75, 2001.
- [21] M. Do and M. Vetterli, "Wavelet-based Texture Retrieval Using Generalized Gaussian Density and Kullback-Leibler Distance," *IEEE Trans. on Image Processing*, vol. Vol. 11, no. 2, pp. 146–158, February 2002.
- [22] S. Chang, B. Yu, and M. Vetterli, "Adaptive Wavelet Thresholding for Image Denoising and Compression," *IEEE Trans. on Image Processing*, vol. 9, no. 9, pp. 1532–1546, September 2000.
- [23] M. Heiler and C. Schnorr, "Natural image statistics for natural image segmentation," *International journal of computer vision*, vol. 63, no. 1, pp. 5–19, 2005.
- [24] F. Attneave, "Informational Aspects of Visual Perception," *Psychological Review*, vol. 61, pp. 183–193, 1954.
- [25] H. Barlow, "The Coding of Sensory Messages," in *Current Problems in Animal Behaviour*, W. Thorpe and O. Zangwill, Eds. Cambridge University Press, 1961, pp. 331–360.
- [26] ———, "Redundancy Reduction Revisited," *Network: Computation in Neural Systems*, vol. 12, pp. 241–253, 2001.
- [27] B. Olshausen and D. Field, "Sparse coding with an overcomplete basis set: A strategy employed by V1?" *Vision Research*, vol. 37, pp. 3311–3325, 1997.
- [28] O. Schwartz and E. Simoncelli, "Natural signal statistics and sensory gain control," *Nature Neuroscience*, vol. 4, pp. 819–825, 2001.
- [29] S. Deneve, P. Latham, and A. Pouget, "Reading population codes: a neural implementation of ideal observers," *Nature Neuroscience*, vol. 2, pp. 740–745, 1999.
- [30] A. Pouget, P. Dayan, and R. Zemel, "Information processing with population codes," *Nat Rev Neurosci*, vol. 1, no. 2, pp. 125–32, 2000.
- [31] D. Gao and N. Vasconcelos, "Discriminant saliency for visual recognition from cluttered scenes," in *Advances in Neural Information Processing Systems 17*, L. K. Saul, Y. Weiss, and L. Bottou, Eds. Cambridge, MA: MIT Press, 2005, pp. 481–488.
- [32] A. Levin, A. Zomet, and Y. Weiss, "Separating reflections from a single image using local features," *IEEE Conf. on Computer Vision and Pattern Recognition*, vol. 1, 2004.
- [33] R. Duda, P. Hart, and D. Stork, *Pattern Classification*. John Wiley and Sons, 2001.
- [34] N. Vasconcelos and G. Carneiro, "What is the Role of Independence for Visual Recognition?" in *Proc. European Conference on Computer Vision, Copenhagen, Denmark, 2002*.
- [35] A. Treisman and G. Galade, "A Feature Integration Theory of Attention," *Cognitive Psychology*, vol. 12, no. 1, pp. 97–136, 1980.
- [36] A. Treisman and S. Gormican, "Feature analysis in early vision: Evidence from search asymmetries," *Psychological Review*, vol. 95, no. 1, pp. 15–48, 1988.
- [37] A. Treisman and S. Sato, "Conjunction search revisited," *Journal of Experimental Perception and Performance*, vol. 16, pp. 459–478, 1990.
- [38] K. Cave and J. Wolfe, "Modeling the role of parallel processing in visual search," *Cognitive Psychology*, vol. 22, no. 5, pp. 225–271, 1990.
- [39] J. Wolfe, "Guided search 2.0: A revised model of visual search," *Psychonomic Bulletin & Review*, vol. 1, no. 2, pp. 202–238, 1994.

- [40] J. Wolfe and T. Horowitz, "What attributes guide the deployment of visual attention and how do they do it?" *Nature Reviews Neuroscience*, vol. 5, pp. 495–501, 2004.
- [41] D. Lewis, "Feature selection and feature extraction for text categorization," *Proceedings of the workshop on Speech and Natural Language*, pp. 212–217, 1992.
- [42] S. Dumais and H. Chen, "Hierarchical classification of Web content," *Proc. international ACM SIGIR conference on Research and development in information retrieval*, pp. 256–263, 2000.
- [43] S. Dumais, "Using SVMs for text categorization," *IEEE Intelligent Systems*, vol. 13, no. 4, pp. 21–23, 1998.
- [44] G. W. F. Lochovsky and Q. Yang, "Feature selection with conditional mutual information maximin in text categorization," in *Proceedings of the Thirteenth ACM conference on Information and knowledge management*. ACM Press New York, NY, USA, 2004, pp. 342–349.
- [45] Y. Seo, A. Ankolekar, and K. Sycara, "Feature Selection for Extracting Semantically Rich Words," Carnegie Mellon University, the Robotics Institute, Tech. Rep. CMU-RI-TR-04-18, March 2004.
- [46] E. Xing, M. Jordan, and R. Karp, "Feature selection for high-dimensional genomic microarray data," *Proceedings of the Eighteenth International Conference on Machine Learning*, pp. 601–608, 2001.
- [47] C. Ding and H. Peng, "Minimum redundancy feature selection from microarray gene expression data," *Proceedings of the IEEE Bioinformatics Conference, 2003.*, pp. 523–528, 2003.
- [48] H. Peng, F. Long, and C. Ding, "Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, and Min-Redundancy," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 27, no. 8, pp. 1226–1238, 2005.
- [49] P. Zarjam, M. Mesbah, and M. Boashash, "An optimal feature set for seizure detection systems for newborn EEG signals," *Proceedings of the International Symposium on Circuits and Systems*, vol. 5, 2003.
- [50] E. Grall-Maes and P. Beuseroy, "Mutual information-based feature extraction on the time-frequencyplane," *IEEE Transactions on Signal Processing*, vol. 50, no. 4, pp. 779–790, 2002.
- [51] G. Tourassi, E. Frederick, M. Markey, and C. F. Jr, "Application of the mutual information criterion for feature selection in computer-aided diagnosis," *Medical Physics*, vol. 28, p. 2394, 2001.
- [52] T. Chow and D. Huang, "Estimating optimal feature subsets using efficient estimation of high-dimensional mutual information," *IEEE Transactions on Neural Networks*, vol. 16, no. 1, pp. 213–224, 2005.
- [53] G. Barrows and J. Sciortino, "A Mutual Information Measure for Feature Selection with Application to Pulse Classification," in *IEEE Intern. Symposium on Time-Frequency and Time-Scale Analysis*, 1996.
- [54] R. Battiti, "Using Mutual Information for Selecting Features in Supervised Neural Net Learning," *IEEE Trans. on Neural Networks*, vol. 5, no. 4, pp. 537–550, July 1994.
- [55] M. Kwak and C. Choi, "Input feature selection for classification problems," *IEEE Transactions on Neural Networks*, vol. 13, no. 1, pp. 143–159, 2002.
- [56] P. Scanlon, G. Potamianos, V. Libal, and S. Chu, "Mutual Information Based Visual Feature Selection for Lipreading," *Proc. Int. Conf. Spoken Language Processing*, pp. 857–860, 2004.
- [57] H. Yang and J. Moody, "Data Visualization and Feature Selection: New Algorithms for Nongaussian Data," in *Proc. Neural Information Proc. Systems*, Denver, USA, 2000.
- [58] F. Fleuret, "Fast Binary Feature Selection with Conditional Mutual Information," *The Journal of Machine Learning Research*, vol. 5, pp. 1531–1555, 2004.

- [59] S. Ullman, M. Vidal-Naquet, and E. Sali, “Visual features of intermediate complexity and their use in classification,” *Nature Neuroscience*, vol. 5, no. 7, pp. 1–6, 2002.
- [60] M. Vidal-Naquet and S. Ullman, “Object Recognition with Informative Features and Linear Classification,” in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Madison, Wisconsin, 2003.
- [61] F. Jurie and B. Triggs, “Creating Efficient Codebooks for Visual Recognition,” in *International Conference on Computer Vision*, vol. 1, 2005.
- [62] N. Vasconcelos, “Feature Selection by Maximum Marginal Diversity,” in *Neural Information Processing Systems, Vancouver, Canada*, 2002.
- [63] —, “Feature Selection by Maximum Marginal Diversity: Optimality and Implications for Visual Recognition,” in *Proc. IEEE Computer Vision and Pattern Recognition Conf.*, Madison, Wisconsin, 2003.
- [64] N. Vasconcelos and M. Vasconcelos, “Scalable Discriminant Feature Selection for Image Retrieval and Recognition,” in *Proc. IEEE Computer Vision and Pattern Recognition Conf.*, Washington DC, 2004.
- [65] D. Koller and M. Sahami, “Toward Optimal Feature Selection,” in *Proc. International Conference on Machine Learning*, Bari, Italy, 1996.
- [66] A. Jain and D. Zongker, “Feature Selection: Evaluation, Application, and Small Sample Performance,” *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 19, no. 2, pp. 153–158, February 1997.
- [67] Y. Freund and R. Schapire, “A decision-theoretic generalization of on-line learning and an application to boosting,” *Journal of Computer and System Sciences*, vol. 55, no. 1, pp. 119–139, 1997.
- [68] R. Schapire, Y. Freund, P. Bartlett, and W. Lee, “Boosting the Margin: A New Explanation for the Effectiveness of Voting Methods,” *The Annals of Statistics*, vol. 26, no. 5, pp. 1651–1686, 1998.
- [69] C. Ratanamahatana and D. Gunopulos, “Feature Selection for the Naive Bayesian Classifier Using Decision Trees,” *Applied Artificial Intelligence*, vol. 17, no. 5, pp. 475–487, 2003.
- [70] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, 1988.
- [71] H. Schneiderman and T. Kanade, “Object Detection Using the Statistics of Parts,” *International Journal of Computer Vision*, vol. 56, no. 3, pp. 151–177, 2004.
- [72] T. Cover and J. V. Campenhout, “On the Possible Orderings in the Measurement Selection Problem,” *IEEE Trans. on Systems, Man, and Cybernetics*, vol. 7, no. 9, September 1977.
- [73] P. Pudil, J. Novovičová, and J. Kittler, “Floating search methods in feature selection,” *Pattern Recognition Letters*, vol. 15, no. 11, pp. 1119–1125, 1994.
- [74] S. Li, L. Zhu, Z. Zhang, A. Blake, H. Zhang, and H. Shum, “Statistical Learning of Multi-View Face Detection,” *Proceedings of the 7th European Conference on Computer Vision*, 2002.
- [75] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning : data mining, inference, and prediction*. Springer, NY, 2001.
- [76] M. Turk and A. Pentland, “Eigenfaces for Recognition,” *Journal of Cognitive Neuroscience*, vol. 3, 1991.
- [77] P. Belhumeur, J. Hespanha, and D. Kriegman, “Eigenfaces vs. Fisherfaces: Recognition Using Class Specific Linear Projection,” *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 19, no. 7, pp. 711–720, July 1997.
- [78] *Comparison of Infomax and Maximum Variance Features*, <http://www.svcl.ucsd.edu/projects/infomax/examples.htm>.
- [79] T. Cover and J. Thomas, *Elements of Information Theory*. John Wiley, 1991.
- [80] M. Vasconcelos and N. Vasconcelos, “Some relationships between minimum bayes error and information theoretical feature extraction,” *Proceedings of SPIE*, vol. 5807, p. 284, 2005.



Manuela Vasconcelos received the licenciatura in electrical engineering from the Universidade do Porto, Portugal, in 1986, a MS from the Massachusetts Institute of Technology in 1993, and a PhD from Harvard University in 2003. Since 2003 she has been a visiting researcher at the Statistical Visual Computing Laboratory, in the University of California San Diego. Her research interests are in the area of computer vision.



Nuno Vasconcelos received the licenciatura in electrical engineering and computer science from the Universidade do Porto, Portugal, in 1988, and the MS and PhD degrees from the Massachusetts Institute of Technology in 1993 and 2000, respectively. From 2000 to 2002, he was a member of the research staff at the Compaq Cambridge Research Laboratory, which in 2002 became the HP Cambridge Research Laboratory. In 2003, he joined the Electrical and Computer Engineering Department at the University of California, San Diego, where he heads the Statistical Visual Computing Laboratory. He is the recipient of a US National Science Foundation CAREER award, a Hellman Fellowship, and has authored more than 50 peer-reviewed publications. His work spans various areas, including computer vision, machine learning, signal processing and compression, and multimedia systems. He is a member of the IEEE.