

Capturing spatial interdependence in image features: the counting grid, an epitomic representation for bags of features

Alessandro Perina, Nebojsa Jojic
Microsoft Research Redmond, USA

Abstract

In recent scene recognition research images or large image regions are often represented as disorganized “bags” of features which can then be analyzed using models originally developed to capture co-variation of word counts in text. However, image feature counts are likely to be constrained in different ways than word counts in text. For example, as a camera pans upwards from a building entrance over its first few floors and then further up into the sky Fig. 1, some feature counts in the image drop while others rise – only to drop again giving way to features found more often at higher elevations. The space of all possible feature count combinations is constrained both by the properties of the larger scene and the size and the location of the window into it. To capture such variation, in this paper we propose the use of the counting grid model. This generative model is based on a grid of feature counts, considerably larger than any of the modeled images, and considerably smaller than the real estate needed to tile the images next to each other tightly. Each modeled image is assumed to have a representative window in the grid in which the feature counts mimic the feature distribution in the image. We provide a learning procedure that jointly maps all images in the training set to the counting grid and estimates the appropriate local counts in it. Experimentally, we demonstrate that the resulting representation captures the space of feature count combinations more accurately than the traditional models, not only when the input images come from a panning camera, but even when modeling images of different scenes from the same category.

Index Terms

Bag of Features Spatial Layout Scene Analysis Bag of Features Spatial Layout Scene AnalysisB



Capturing spatial interdependence in image features: the counting grid, an epitomic representation for bags of features



Fig. 1. Feature counts change slightly as the field of view moves. For example, the abundance of the “car” features is reduced, but the counts of the features found on building facades are increased. The counting grid model accounts for such changes naturally, and it can also account for images of different scenes.

1 INTRODUCTION

A popular way to deal with diversity of imaging conditions as well as geometric variation in objects or entire scenes is to simply represent images or image regions as disordered “bags” of image features [1]–[3]. These models are particularly attractive due to the computational efficiency and simplicity achieved by ignoring spatial relationships of the image patches or object parts.

The bag of features can arise in a variety of ways. For example, after extracting local low-level features from images, these are often clustered and a discrete “codeword” is assigned to each feature descriptor. An image is then described by a histogram over the codebook entries. Ideally, these features should be highly discriminative so that most categories of images of interest are uniquely identifiable by the presence of a handful of features. In practice, however, individual features are not sufficiently discriminative, and modeling joint variation in feature counts becomes an interesting machine learning problem.

It is tempting to use here the existing discrete models, such as histograms [4], multinomial mixtures [5,6] or topic models [7,8], already extensively validated on text data, where each document is also simply represented as a count distribution over the entire vocabulary. However, the bags of features extracted from natural images have an imprint of the images’ spatial structure, which is evident when the bags from related images are considered *together*. Thus ignoring these natural constraints on the feature counts may have negative conse-

quences in classification tasks.

For an illustration, Fig. 2 provides a synthetic example starting with several images of a train station, taken as windows into the larger scene - *ii*). Just for illustrative purposes, we hand-labeled the scene with feature labels as shown in - *iii*). In a realistic application, where we may want to train a model that assigns high likelihood to images of train stations, it is likely that most available images would be taken with a narrower field of view, as simulated here. Feature extractors would presumably generalize much less effectively than our ideal features, but still enough to permit comparisons of images of *different* train stations, too. Then the question is if a learning model that captures feature count co-variation uses the training data efficiently. Assuming that a few images are taken at random from the scene, we wonder if the feature counts in these images are sufficient to predict the possible feature counts in other images of the scene. In particular, we consider images taken from the regions close to A, B, and C and ask the question if the image D would fit the so defined train station class.

The literature uses two sets of approaches to this problem. Kernel or nearest-neighbor techniques start with the comparisons of the feature counts in the test image and each of the previously studied exemplars [9]–[13]. Although this comparison can be done in many different ways, we note here that these approaches would be complicated by the fact that none of images A,B and C have the combination of all five features that are present in D (see Fig. 2-*i*). The other approach is to consider all bags of features together and generalize [1,3,5,14]–[17]. A simplest approach to this would be to simply merge the bags. In this case, there is a danger of overgeneralization. For this particular example, there is a need for interpolating between the feature count vectors for A,B,C. However, this interpolation is best performed by spatial reasoning. Across various windows into the scene we find that from the top of the window to the bottom we sometimes see roof, train, tracks, in that order, but other times we see mountain, grass, roof, train. We can infer that the grass, roof, train, tracks combination is likelier than the existence of the mountain, roof, train, tracks combination of features. Furthermore, the proportions of different features in the images carry information about the thickness of the layers of these features, which should be useful for inferring which previously unseen feature count combinations can be found elsewhere in the scene. We show in this paper that, surprisingly, not much of the spatial organization of the features in the training images needs to be retained in order to perform the spatial reasoning about which feature combinations are likely.

In Fig. 2-*iv*) we show the counting grid inferred by iterating Eqs. 12 and 13 on the label counts from 50 windows into the scene taken at random, *but avoiding all windows that contain all five of the features in D in any proportion*. Each training image was represented as a set of 2×2 feature bags (upper left, lower left, upper right, lower right). Without using the

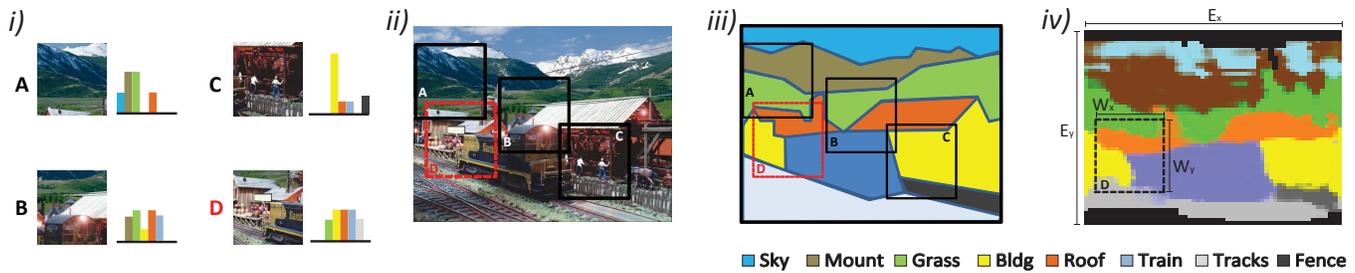


Fig. 2. Counting grid illustration. i) Images and their Bag of feature representation. ii) Images of a train station, taken as windows into the larger scene. iii) Hand labeled features. iv) Scene reconstructed (e.g., counting grid) starting from bags taken from 50 windows (see the Text for details).

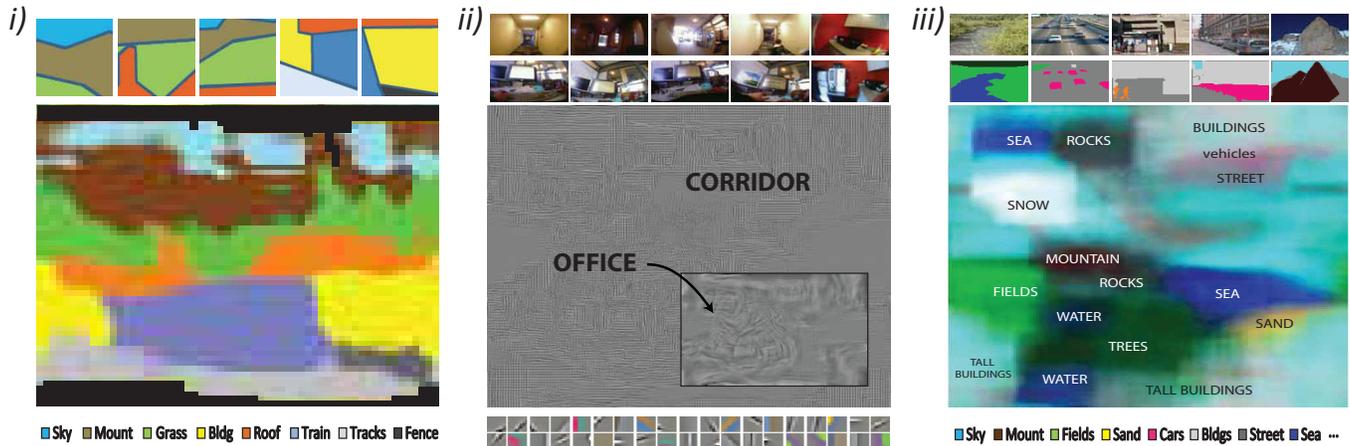


Fig. 3. Illustration of counting grids using different data at different levels of abstraction. At the top of each CG we show few examples of the training images from which we extracted the bags. The bottom row of each panel, illustrates the features. i) A counting grid learned using patches extracted from the train station toy example of Fig. 2. In this case all the bags come from windows into the same image, which is reconstructed on the grid. ii) A counting grid estimated from images taken with a wearable camera [18]. In this case we learned a dictionary of features (illustrated by the textons on the bottom) clustering image patches, as in [19]. To illustrate the office scene (e.g., the computer screens - see images on the top), we overlap these textons by as much as the patches were overlapping during feature extraction process, and then average to create a clearer visual representation. iii) A counting grid estimated starting from LabelMe annotations (see the top row) [20].

original window location information, the counting grid was computed so that for each training image, a window into the counting grid can be found so that the appropriate sections have matching histograms. The resolution of the reconstructed feature layout of the large scene goes well beyond what would be expected from a crude 2×2 tessellation of the input images (the height of each section is roughly 20% of the large scene and only the feature counts in each section were used, not their spatial layout within the section). Although none of the training examples was taken from the area close to D where all five of D's features can be seen in a single image, that part of the scene is reconstructed as well, and D's histogram can be matched well.

In this simple example, the training images are different views of a single scene. However, at the feature level, images of other train stations are likely have a similar layout, and so they could be used to learn a counting grid. In practice, we rarely have access to highly discriminative and reliable features, and so instead of the 8 fake features in our example, in our experiments we had to use hundreds of simpler automatically derived features, and infer the counting grids from related

images of different scenes. For example in Fig. 3-iii) we used as features the human-supplied labels for LabelMe [20] dataset, and in - ii) the outputs of hundreds of simple computational feature detectors applied to images from various scenes in the SenseCam dataset [18]. As opposed to the train station example, input images are not subimages of a larger single scene, but rather images of the same types of scenes. Each window into a counting grid represents a possible feature combination¹ present in the dataset (3- ii) and the model is able to reconstruct the feature layout only exploiting the spatial patterns very coarsely, but depending on feature count co-variation for most of its reasoning power. For example, despite only relying on a 2×2 tessellation of each input image, full resolution panoramas of office and corridors are visible in the CG in Fig. 3-ii).

This paper presents and extends the counting grid model [21,22]. The basic model is extended to include priors which help with overfitting issues. We also formally introduce the tessellated counting grid model and analyzed the the extreme

1. The CG model also learns a prior over the grid window usage which may prevent some combinations.

tessellations where each bag captured a single feature, collapsing the representation into the discrete epitome [23]. The paper also provides full comparisons of different algorithms on various datasets, including the effect of grid and window size variation can be found in the experimental section.

Specialization and extensions of the counting grid model already appeared in top tier conferences [18,24]–[26]. Nevertheless, in this paper we want to limit our attention to the basic variants of the counting grid model, their properties and relationships with the standard techniques for modeling bags of words in computer vision [1,3,16]. We found that our representation captured the space of possible feature count combinations for various image categories significantly better than other generalization techniques, and that our simple generative model, which can be used for unsupervised learning and clustering, too, often rivals the state of the art based on discriminative techniques that require supervision.

Related Work

Previous probabilistic approaches to scene recognition treat the spatial arrangement of image features in different ways.

In bag of words (BOW) models [1], spatial relationships among features are completely ignored in order to facilitate computational efficiency and high level of generalization. Topic models [3,7,8], for example, assign a topic to each codeword based on their co-occurrence and describe images as admixtures of topics. Another bag of word model is described in [17], where a scene model is a mixture of Gaussians model trained on the gist descriptors [30]. As we will see in this paper the basic counting grid model [22], also reduces to a (large) mixture, but with highly tied parameters, reflecting the inherent spatial structure of the data. Each bag is represented as a point in a large grid of feature counts. This latent point is a corner of a window of grid points which are uniformly combined to match the (normalized) feature counts in the image.

To capture some spatial information, it is possible to separate the bags originating in different (pre-defined or learned [31]–[33]) regions of the image. These models are sometimes referred to as spatial-BoW models [2,16]. The tessellated counting grids that we introduce here also have that flavor, although in our approach tessellation helps guide the quilting of the bags of words reconstructing the layout with sub-region accuracy. Thus tessellated counting grids capture layout-driven constraints on counts within the regions, even though this information is not directly provided during learning: The layout within a region of one image is inferred based on the feature distributions found in regions of many other images, assuming that misalignments of these images are often smaller than the size of the tessellated regions. In contrast, typical spatial-BOW models requires the modeled images to be approximately aligned. Recent approaches that relax this assumption are [16,18,34]. The former, the Reconfigurable BoW model [16], represents a scene as a collection of parts arranged in a reconfigurable pattern. Each image is divided into pre-defined regions and a latent variable specifies which “region model” (e.g., sky, grass...) is assigned to each image region. On the other hand, [18,34] represents scenes using deformable parts. In [34] a lower-resolution root filter is placed in the center of the image and a set of higher-resolution part filters arranged in a flexible spatial configuration.

It is also possible to keep the spatial arrangement of features intact, sacrificing some generalization in the basic rep-

resentation of the input, and allowing the model to capture the problems with this rigidity through various levels of uncertainty modeling. For example, the epitome-like models [23,27,28] quilts images or image patches, essentially building giant panoramas consisting of probability distributions in each location. As these are based on pixel-to-pixel comparisons they cannot generalize well in case of large geometric deformations, and so they are mostly used to model relatively small image patches, typically for synthesis, or modeling large scenes or textures that can tolerate the lack of transformation invariance beyond translation [29,35,36]. Epitomes have been employed in scene analysis, only on particular datasets where “panoramic stitching” would work, e.g., sequences taken with wearable cameras [23,28].

Various scene modeling techniques also take different approaches to representing exponential structure of natural scenes. Being ad-mixtures, rather than simple mixtures, topic models [3,7,8] are a simple example of multi *-part* or *-object* models. Other examples of componential models, are the flexible sprites model [29], which allow each image to be mapped to multiple sources and [16] in which each sector is mapped independently. On the other hand, the counting grids, epitomes [23,27,28], histogram-based approaches [1,17] are essentially mixtures because they map the entire scene to a single point: a position or a mixture component. (These models can, however, be turned into ad-mixtures.)

The main topic this paper is modeling bags of features in computer vision. In the experimental section we will mainly consider generative approaches and compare counting grids with latent Dirichlet allocation [3], mixture models [17], epitomes [23] and the reconfigurable bag of words model [16]. The nature of each generative approach just discussed is summarized in Tab. 1. It should be noted however that the models presented here can be used as components in hierarchical models, and that the basic idea of modeling intersections and laying them out on an inferred grid can be used within other machine learning techniques, including non-generative approaches.

2 IMPRINT OF SPATIAL ORGANIZATION IN DISORDERED BAGS OF WORDS

As discussed above, we would like to understand the hidden constraints that govern the often-practiced simplification of images into bags of features. This simplification has two stages. First, image features $z_{i,j}$ are extracted on a grid inside the image. These features are discrete, $z \in [1..Z]$, and they point to a codebook of features obtained by clustering the multidimensional real-valued features calculated by local image processing, e.g., SIFT [37]. Next, the feature counts are computed $c_z = \sum_i [z_i = z]$, where $[\cdot]$ is the indicator function. Only the counts c_z are then retained, and the spatial distribution z_i is typically forgotten, with the justification that establishing correspondence for individual image locations across different images of the same thing would be prohibitively expensive, and that in practice only the presence or absence of features is informative, not their spatial distribution. However, if we consider a set of such bags of words from related images we can see that the feature counts in these disordered bags of features may still indirectly follow the rules of spatial organization. For example, if the bags $\{c_z^t\}$, indexed by t are extracted from several overlapping windows from a larger image, then the spatial structure of that image is imprinted in

TABLE 1
Generative approaches to scene analysis.

Method	Componential hidden variables			Spatial structure in input		
	Single integer	Multiple integers	Continuous	BoW	Tessellated	Pixel
LDA [3,7]			✓	✓	(✓)	
Multinomial Mixtures [1,17]	✓			✓		
Spatial BoW [2]	✓				✓	
Reconfigurable BoW [16]		✓			✓	
Epitomes [23,27,28]	✓					✓
Flex. Sprites [29]		✓				✓
Counting grids	✓			✓		
Tessellated Counting grids	✓				✓	

the particular count combinations in these bags. Furthermore, the spatial layout of the features in the large image may even be recoverable from these disordered bags! If the bags $\{c_z^t\}$ are created from *all* the overlapping windows from a large image, and if the source location for each bag is known, then we can easily see that under minimal additional assumptions regarding the boundaries in the image, we can reconstruct feature indices z at each location in the large image by solving the system of linear equations that arise from the count constraints. Consider two horizontally neighboring windows: The count differences are completely determined by the feature identities of the only two columns that the two do not share. To separate the effect of the two columns, we can consider another pair of overlapping images whose count differences depend on only one of those two columns. To further break each column apart, we can consider vertically neighboring windows, etc. As long as the image has a thick enough border with only a single feature present, we can propagate these constraints until any given location's feature is uniquely determined.

In this way, we can reconstruct a large grid of features such that any of the count combinations we see in the given bags can be found in an appropriate window in this reconstruction. But this implies that the bags of features from the images of the same scene, when considered jointly, obey very strong constraints and thus taking these constraints into account will likely improve image analysis tasks that depend on the feature count representations. This insight leads to several interesting problems which we address in the next section.

- *Joint estimation of the feature layout and the matching of the bags to windows into it:* If the bags of features (feature counts) from many – but *not all* – overlapping windows from a large scene are provided, and if the original locations of these windows are *withheld*, can we still reconstruct at least some of the original spatial arrangement of the features?
- *Category modeling:* If the bags of features are not coming from the windows into a single scene, but instead from different but related images (e.g. of a particular image category or an object class), would these bags, when considered jointly, imply some spatial layout of the features, and would this layout help predict which combinations of feature counts are more likely in bags of features extracted from new images of the category in question?
- *Using more of the original structure:* Given that in practice we typically have access to the original images, can more of their spatial structure be used in learning the spatial

layout of features that would in turn constrain the bag of words representation in a useful way?

3 THE COUNTING GRID MODEL

The basic counting grid $\pi_{i,z}$ is a set of normalized counts of features indexed by z on the grid $\mathbf{i} = (i_x, i_y) \in \mathbf{E} = [1 \dots E_x] \times [1 \dots E_y]$, with $\sum_z \pi_{i,z} = 1$ everywhere on the grid [21,22]. A given bag of image features, represented by counts $\{c_z\}$ is assumed to follow a distribution found somewhere in the counting grid. In other words, the bag can be generated by firstly averaging all counts in the window $W_{\mathbf{k}}$ of size $W_x \times W_y$ placed at location \mathbf{k}

$$W_{\mathbf{k}} = [k_x, \dots, k_x + W_x - 1] \times [k_y, \dots, k_y + W_y - 1]$$

to form the histogram

$$h_{\mathbf{k},z} = \frac{1}{(W_x \cdot W_y)} \cdot \sum_{\mathbf{i} \in W_{\mathbf{k}}} \pi_{\mathbf{i},z} \quad (1)$$

and then generating the features in the bag. The sum in Eq. 1 is carried out in all the locations \mathbf{i} in the window $W_{\mathbf{k}}$. An example of counting grid geometry is illustrated in Fig. 4-ii) In other words, the position of the window in the grid is a latent variable ℓ given which the probability of the bag of features $\mathbf{c} = \{c_z\}_{z=1}^Z$ is

$$p(\mathbf{c}|\ell = \mathbf{k}) = \prod_{z=1}^Z (h_{\mathbf{k},z})^{c_z} = \alpha \cdot \prod_{z=1}^Z \left(\sum_{\mathbf{i} \in W_{\mathbf{k}}} \pi_{\mathbf{i},z} \right)^{c_z} \quad (2)$$

where the constant $\alpha = \left(\frac{1}{W_x \cdot W_y} \right)^{\sum_z c_z}$. In our notation the letter ℓ indicates the latent variable, while \mathbf{i} and \mathbf{k} a generic position in the grid.

The Bayesian network of the model is illustrated in Fig. 4-i). For a given grid p_i , it defines the following joint distribution over all bags of features $\{c_z^t\}$, indexed by t and their corresponding latent window positions ℓ^t in the counting grid

$$P(\{\mathbf{c}^t\}, \{\ell^t\}) \propto \prod_{t=1}^T \sum_{\mathbf{k} \in \mathbf{E}} \left(P(\ell^t = \mathbf{k}) \cdot \prod_{z=1}^Z \left(\sum_{\mathbf{i} \in W_{\mathbf{k}}} \pi_{\mathbf{i},z} \right)^{c_z^t} \right)$$

Where $P(\ell = \mathbf{k})$ represents the overall prior probability of a mapping location. The first sum in the RHS is performed over all the location of the counting grid, while the second over all the locations in the window placed at location $W_{\mathbf{k}}$.

To summarize the notation we will use throughout the paper, ℓ is the hidden variable that represents the mapping location in the grid; each bag (sample) is mapped to a (possibly) different location and we will use the superscript t to refer to the particular t -th bag therefore ℓ^t will represent the mapping

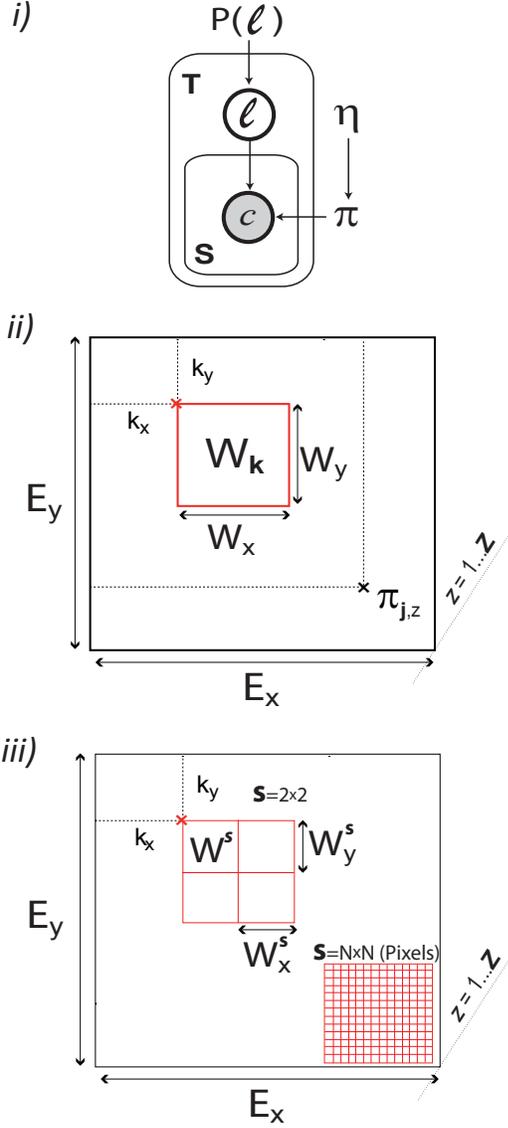


Fig. 4. i) Counting grid Generative model, ii) Counting grid Geometry, iii) Tessellated counting grid geometry.

position of the bag $c^t = \{c_z^t\}$. Since we are introducing a probabilistic model, it is interesting to estimate the prior probability $p(\ell^t = \mathbf{k})$ of mapping the t -th sample to location \mathbf{k} . In this case ℓ^t is the t -th sample's hidden variable (window location), while \mathbf{k} is a generic constant or index that represents a possible location in the grid that samples share. This distribution is a table over all values for \mathbf{k} and shared across all samples (independent of t). On the other hand, the posterior distribution $p(\ell^t = \mathbf{k} | c^t)$, or its (exact) variational counterpart $q(\ell^t = \mathbf{k})$, is a function of the counts seen in the t -th sample and capture the quality of the fit to different windows in the grid of the t -th sample in particular.

3.1 Inference and learning

To compute the log likelihood of the data, $\log P$, we need to sum over the latent variable ℓ before computing the logarithm, which, as in mixture models, or as in epitomes [27], makes it difficult to perform assignment of the latent variables while

also estimating the model parameters. Although the following is the exact EM procedure, we use the variational [38,39] notation $p(\ell^t | c^t) = q(\ell^t)$, and bound (variationally) $\log P$ (omitting the effect of additive constant that arises from α) to derive an iterative EM algorithm:

$$\begin{aligned} \log P &\geq \sum_{t=1}^T \sum_{\mathbf{k} \in \mathbf{E}} \left(q(\ell^t = \mathbf{k}) \cdot \log q(\ell^t = \mathbf{k}) \right. \\ &\quad - q(\ell^t = \mathbf{k}) \cdot \log P(\ell = \mathbf{k}) \\ &\quad \left. - q(\ell^t = \mathbf{k}) \cdot \left(\sum_z c_z^t \cdot \log h_{\mathbf{k},z} \right) \right) = B, \end{aligned} \quad (3)$$

Because of the use of fully parameterized q , optimizing the bound is equivalent to optimizing the log likelihood of the data, as long as the $q(\ell^t)$ distributions are also optimized. Keeping the model parameters fixed, optimizing these q distribution (exact E step) leads to

$$q(\ell^t = \mathbf{k}) \propto P(\ell = \mathbf{k}) \cdot \exp \left(\sum_{z=1}^Z c_z^t \cdot \log h_{\mathbf{k},z} \right), \quad (4)$$

which simply establishes that the choice of ℓ should minimize the KL divergence between the counts in the bag and the counts $h_{\mathbf{k},z}$ in the appropriate window $W_{\mathbf{k}}$ in the counting grid. For each t , the above expression is normalized over all possible window choices \mathbf{k} .

To optimize the bound B with respect to model parameters (M step) we note that the first term in Eq. 3 involves these parameters, and it requires another summation before applying the logarithm. The summation is over the grid positions \mathbf{i} within the window $W_{\mathbf{k}}$, which we can again bound using a (full) variational distribution and the Jensen's inequality:

$$\log \sum_{\mathbf{i} \in W_{\mathbf{k}}} \pi_{\mathbf{i},z} = \log \sum_{\mathbf{i} \in W_{\mathbf{k}}} r_{\mathbf{i},\mathbf{k},z}^t \frac{\pi_{\mathbf{i},z}}{r_{\mathbf{i},\mathbf{k},z}^t} \geq \sum_{\mathbf{i} \in W_{\mathbf{k}}} r_{\mathbf{i},\mathbf{k},z}^t \log \frac{\pi_{\mathbf{i},z}}{r_{\mathbf{i},\mathbf{k},z}^t}, \quad (5)$$

where $r_{\mathbf{i},\mathbf{k},z}^t$ is a distribution over locations \mathbf{i} , i.e. r is positive and $\sum_{\mathbf{i} \in W_{\mathbf{k}}} r_{\mathbf{i},\mathbf{k},z}^t = 1$. It is indexed by \mathbf{k} as the normalization is done differently in each window, it is indexed by z as it can be different for different features, and it is indexed by t as the term is inside the summation over t , so a different distribution r could be needed for different bags $\{c_z^t\}$. This distribution could be thought of as information about what proportion of the c_z features of type z was contributed by each of the different sources $\pi_{\mathbf{i},z}$ in the window $W_{\mathbf{k}}$. However, by performing constrained optimization (so that r adds up to one), we find that assuming a fixed set of parameters π , the distribution $r_{\mathbf{i},\mathbf{k},z}^t$ that maximizes the bound is the same for each bag:

$$r_{\mathbf{i},\mathbf{k},z}^t = \frac{\pi_{\mathbf{i},z}}{\sum_{\mathbf{i} \in W_{\mathbf{k}}} \pi_{\mathbf{i},z}} = \frac{\pi_{\mathbf{i},z}}{W_x \cdot W_y \cdot h_{\mathbf{k},z}}. \quad (6)$$

If we do consider distributions r as a feature mapping to the counting grid, then this result is again intuitive. If all we know is that a bag containing c_z features of type z is mapped to the grid section $W_{\mathbf{k}}$, and have no additional information about what proportions of these c_z features were contributed from different incremental counts $\pi_{\mathbf{i},z}$, then the best guess is that these proportions follow the proportions among $\pi_{\mathbf{i},z}$ inside the window.

If we assume now that r and q distributions are fixed, then combining Eq. 3 and Eq. 5 and minimizing the resulting bound

wrt parameters $\pi_{i,z}$ under the normalization constraint over features z , we obtain the update rule

$$\hat{\pi}_{i,z} \propto \sum_{t=1}^T \sum_{\mathbf{k} | \mathbf{i} \in W_{\mathbf{k}}} q(\ell^t = \mathbf{k}) \cdot c_z^t \cdot r_{i,\mathbf{k},z}^t, \quad (7)$$

which by Eq. 6 reduces to

$$\hat{\pi}_{i,z} \propto \pi_{i,z}^{old} \cdot \sum_{t=1}^T \left(c_z^t \cdot \sum_{\mathbf{k} | \mathbf{i} \in W_{\mathbf{k}}} \frac{q(\ell^t = \mathbf{k})}{h_{\mathbf{k},z}} \right), \quad (8)$$

where $\pi_{i,z}^{old}$ is the counting grid at the previous iteration.

The reader will note that in the above, we simply optimized the likelihood of the set of data for a single set of weights $\pi_{i,z}$, as Eq. 3 is the variational bound for the model with a fixed π as it was expanded in the previous section. Thus the iteration of the above equations would optimize for a set of parameters π given the observed data and ignoring the prior over π in the full network in Fig. 4-i). Of course, the Dirichlet prior with parameters η is the appropriate conjugate prior (as in LDA models) making the inclusion of its influence trivial: The parameters η_z , one for each feature, act as pseudocounts of each feature,

$$\hat{\pi}_{i,z} \propto \eta_z + \pi_{i,z}^{old} \cdot \sum_{t=1}^T \left(c_z^t \cdot \sum_{\mathbf{k} | \mathbf{i} \in W_{\mathbf{k}}} \frac{q(\ell^t = \mathbf{k})}{h_{\mathbf{k},z}} \right), \quad (9)$$

The prior elegantly precludes zero counts of any feature anywhere in π , preventing overtraining and numerical problems. The innermost sum of the equation above is carried out across all the locations \mathbf{k} whose window $W_{\mathbf{k}}$ contains the generic location \mathbf{i} indexed in the LHS. This simply reduces to summing in “shifted” windows where now \mathbf{k} represents the lower right corner. Finally, by taking derivatives with respect to prior probabilities of different locations, we can readily show that the update for the prior over locations should be updated as follows:

$$P(\ell = \mathbf{k}) \propto \sum_{t=1}^T q(\ell^t = \mathbf{k}). \quad (10)$$

This is not surprising, as mathematically, the model is a mixture of distributions h and the above is essentially an update for a mixture prior. However, if we consider the data efficiency of this update, we see that it differs dramatically from how efficiently the data is used to learn distributions h . Consider a large CG model, e.g. on a 64×64 grid, that uses relatively large windows, too, e.g. 16×16 . Then even though there are over $4k$ individual distributions π to learn, these are in fact used in aggregates of 256 at a time in each of the h distributions, which makes the parameters of the mixture’s sources highly tied. In fact, we can only tile $4 \times 4 = 16$ non-overlapping windows over the grid, and the rest of the 4096 overlapping windows are a special kind of interpolation of these 16. Thus the equivalent capacity of such a model, when compared with a simple mixture, is only 16, allowing such grids to be trained without overtraining with just an order or two more data than this capacity number. In other words, we should be able to train a model with $4k$ fractional sources π with only around $1k$ bags of words. But the equivalently efficient use of data for estimating which parts of the grid are used more than others would require a similar aggregation of fractional probabilities of individual cells, just like π distributions are aggregated into h distributions. The similar issue was resolved in epitome models by literally aggregating the updates above

within overlapping windows, to avoid overfocusing the prior probability over the $4k$ windows in our example on only those $1k$ positions where training data fell:

$$P(\ell = \mathbf{k}) \propto \sum_{t=1}^T \sum_{\mathbf{i} \in \mathbf{E}} q(\ell^t = \mathbf{i}) \cdot m_{\mathbf{k}-\mathbf{i}}. \quad (11)$$

where m is a $E_x \times E_y$ mask, with ones in the upper left corner’s $W_x \times W_y$ entries and zeros elsewhere. In our experiments the same update proved to be a valid way to overcome local minima when the prior over location is learned.² In this update, the prior $P(\ell)$ must, of course, be normalized across the locations.

The steps in Eqs. 4, 8 and 11 constitute the E and M step which can be iterated till convergence (within a desired precision τ). The learning algorithm is summarized in Alg. 1.

Algorithm 1: EM-Algorithm to learn a counting grid.

Input: Bag of features, c_z^t for each patch, counting grid size \mathbf{E} , window size \mathbf{W}

while Convergence **do**

 % E-Step ;

foreach Sample $t = 1 \dots T$ **do**

 1. Update $q(\ell^t = \mathbf{k}) \propto \exp \{ \sum_z c_z^t \log h_{\mathbf{k},z} \}$;

 % M-Step ;

 2. Update $\pi_{i,z} \propto \pi_{i,z}^{old} \cdot \sum_t c_z^t \sum_{\mathbf{k} | \mathbf{i} \in W_{\mathbf{k}}} \frac{q(\ell^t = \mathbf{k})}{h_{\mathbf{k},z}}$;

 3. Compute $h_{\mathbf{k},z} = \frac{1}{W_x \cdot W_y} \sum_{\mathbf{i} \in W_{\mathbf{k}}} \pi_{i,z}$;

 4. Update $P(\ell)$ using Eq. 10 or Eq. 11;

 5. Compute the Log-Likelihood B with Eq. 3 ;

 6. Check for convergence, e.g. $|B - B^{old}| \leq \tau$;

6. Return $\pi_{i,z}$, $P(\ell)$ and $\{q(\ell^t)\}_t$;

Starting with non-informative (but symmetry breaking) initialization, this iterative process will jointly estimate the counting grid and align all bags to it. To avoid severe local minima, it is important, however, to consider the counting grid as a torus, and consider all windowing operations accordingly, as was previously proposed for learning epitomes [23,27,28]. This prevents the problems with grid boundaries which otherwise not be crossed when more space is need to grow the layout of the features.

4 FROM COUNTING GRIDS TO FEATURE EPITOMES

We can express many other models used in vision as special cases of our framework by assuming an appropriate choice

² It is also possible to change the model in way that would allow for this update to arise naturally, in a manner equivalent to defining h distributions as arising from π distributions

TABLE 2

Relationship between counting grids and other computer vision methods

Name	E-Step	M-Step	W	S
Counting grid	Eq.4	Eq.8	$\geq 2 \times 2$	1×1
Tessellated CG	Eq.12	Eq.13	$\geq 2 \times 2$	$\geq 2 \times 2$
CG-Epitome [22]	Eq.4	Eq.15	$N_x \times N_y$	$\geq 1 \times 1$
Discr. Epit. [23]	Eq.14	Eq.15	$N_x \times N_y$	$N_x \times N_y$
Mix.Unigram [6]	Eq.4	Eq.8	1×1	1×1
Spatial BoW [2]	Eq.4	Eq.8	1×1	$\geq 2 \times 2$

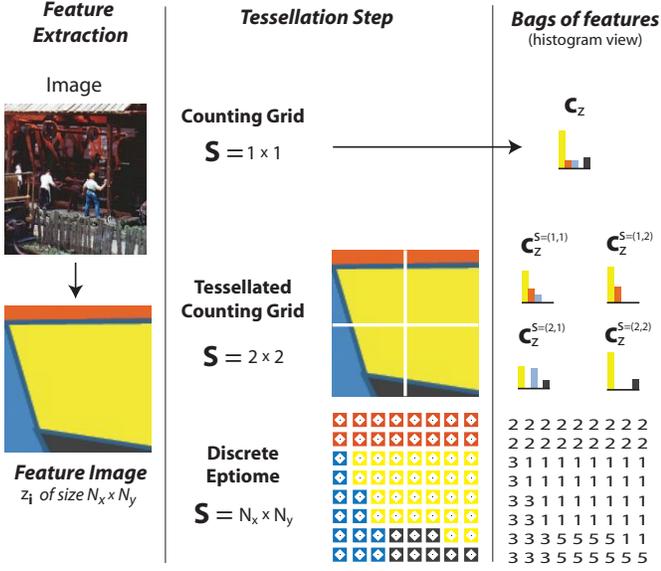


Fig. 5. Tesselation step illustration. Once the features are extracted and quantized one can decide tessellation $\mathbf{S} = S_x \times S_y$ of the image and compute the feature counts separately in each different section as illustrated by the second and third column. To the limit, when $\mathbf{S} = \mathbf{W} = N_x \times N_y$, we obtain the discrete feature epitome model.

of the tessellation \mathbf{S} of the images and the window size \mathbf{W} . A tessellation \mathbf{S} is simply a partition of the image space as illustrated in Fig. 5.

Tesselated counting grids: Algorithm 1 works remarkably well given that its task is essentially to infer an image, not from many image patches as is the case for epitome models, but only from the bag of features representation of for such patches. The task is formidable because no directionality is provided in the bag representation. Unfortunately the iterative algorithm may start to lay out the features topologically correctly, but following inconsistent directions in different parts of the counting grid, leading to local minima (This will be illustrated in the next section). However, we can modify the model and its E and M rules to deal with image representations that consist not of one, but *several* bags of words, each corresponding to a section of the image. In this case feature re-arrangement is tolerated within each region, but the regions themselves cannot move relatively to each other and the model becomes similar in spirit to [2,16]

More specifically, we define a tessellation $\mathbf{S} = S_x \times S_y$ and for each feature map z_i^t , we compute the feature counts separately in each different section $\{c_z^{t,s}\}$ being $s = s_x \times s_y$ a bi-dimensional index that runs across the sectors of \mathbf{S} . This process is illustrated in Fig. 5-ii). When inferring the mapping of the set of section bags, the window $W_{\mathbf{k}}$ is tessellated into $S_x \times S_y$ sections of size $W^{\mathbf{S}}$ indexed by $W_{\mathbf{k}}^{\mathbf{S}}$ in the same way images are tessellated. The histogram comparisons are done accordingly, in formulae:

$$q(\ell^t = \mathbf{k}) \propto P(\ell = \mathbf{k}) \cdot \exp \left(\sum_{s \in \mathbf{S}} \sum_{z=1}^Z c_z^{t,s} \log \sum_{i \in W_{\mathbf{k}}^{\mathbf{S}}} \pi_{i,z} \right), \quad (12)$$

It is important to note that all the $S_x \cdot S_y$ bags contribute to the same mapping on the grid. Therefore the tessellated counting grid model inherits the same componential nature of the counting grid while making use of more spatial information, as reported in Tab. 1.

The M step using section bags reduces to

$$\pi_{i,z} \propto \pi_{i,z}^{old} \cdot \sum_{t=1}^T \left(\sum_{s \in \mathbf{S}} c_z^{t,s} \cdot \sum_{\mathbf{k} | i \in W_{\mathbf{k}}^{\mathbf{S}}} \frac{q(\ell^t = \mathbf{k})}{h_{\mathbf{k},z}^s} \right) \quad (13)$$

The three plates in Fig. 3 show that even just considering an representation consisting of four bags of features for the 4 image sections (upper left, upper right, lower left and lower right) provides enough symmetry breaking that good counting grids can be estimated.

Discrete Epitomes: To the limit, when both tessellation and window size are equal to the images size, e.g., $\mathbf{S} = \mathbf{W} = N_x \times N_y$, we obtain the discrete feature epitome model. In this case, each bag is composed by a single feature $c_z^{t,s} = z_i^t$ and the sector index s indexes a pixel i . In the M-step, there is no re-arrangement of the features in the window and they are simply “copied” according to the mappings $q(\ell^t)$. The E-Step thus becomes:

$$q(\ell^t = \mathbf{k}) \propto P(\ell = \mathbf{k}) \cdot \exp \left(\sum_{i \in \mathbf{E}} \sum_{z=1}^Z [z_i^t = z] \cdot \log \pi_{\mathbf{k}-i,z} \right) \quad (14)$$

where $[\cdot]$ is the indicator function, equal to 1 when the equality holds, zero otherwise. Eq. 14 can be efficiently computed using FFTs [40].

The M-Step reduces to

$$\pi_{i,z} \propto \sum_{t=1}^T \sum_{\mathbf{k} \in \mathbf{E}} q(\ell^t = \mathbf{k}) \cdot [z_{i-\mathbf{k}}^t = z], \quad (15)$$

Likewise the epitome [27] and the counting grid, discrete epitomes are single-component models. However differently from the former, they are characterized by a multinomial observation model and differently from the latter, they consider the original feature layout making the model less efficient and harder to generalize.

Finally, in [23] local histograms are used as pixel descriptor. This helped to overcome the rigidity of epitome models and reach great performances on location recognition. This technique loosely correspond to $\mathbf{W} = N_x \times N_y > \mathbf{S}$.

Hybrid counting grid - epitome: Another alternative is to use the layout of features z_i^t of each image when updating the counting grid (Eq.15) while its bag of words representation to compute the mapping (Eq.4). The result is an hybrid between counting grids and epitomes and it is what has been used in the experimental section of the conference version of this paper [22]. For some dataset this strategy proved to be successful.

Relationships with other models: When $\mathbf{W} = 1 \times 1$, the model collapses into a mixture of unigrams [1] and each point in the grid $\pi_{\mathbf{k},z} = h_{\mathbf{k},z}$ is now a mixture component. If a tessellation is also enforced, the model becomes similar to the spatial BOW models introduced in [2,16].

Finally, despite the counting grid shares its focus on modeling image feature counts with LDA (and in general topic models), neither model is a generalization of another. However, by using large windows to collate many grid distributions from a large grid, the counting grid model can be thought as a

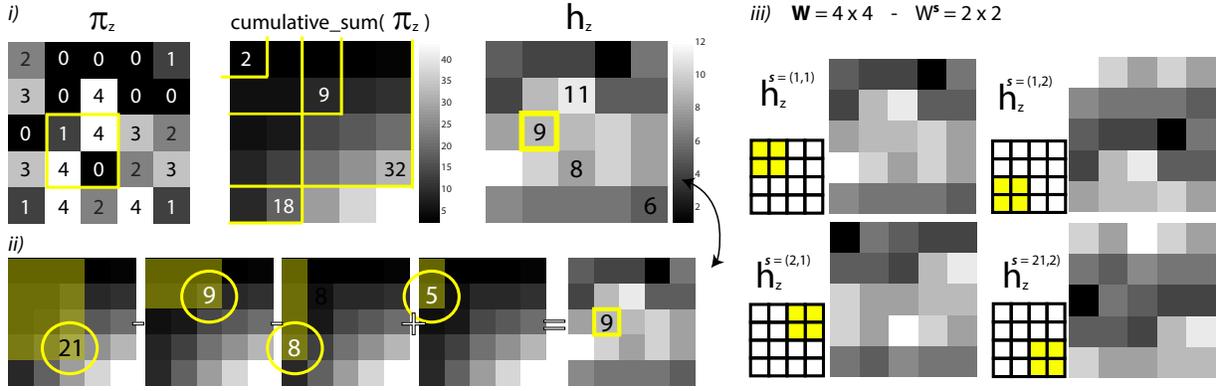


Fig. 6. Few implementation details useful for computational efficiency. Panes i) and ii) shows how h can be efficiently computed using cumulative sums of π . Panel iii) show the shifted versions of h .

very large mixture of sources without overtraining, as these sources are highly correlated: Small shifts in the grid change the window distribution only slightly. LDA model does not have this benefit, and thus has to deal with a smaller number of topics to avoid overtraining. Topic mixing cannot quite appropriately represent feature correlations due to translational camera motion.

The relationships between counting grids variations introduced in this section in terms of W , S and variational updates are summarized in Tab. 4.

5 COMPUTATIONAL COMPLEXITY AND IMPLEMENTATION

Careful examination of the steps reveals that by the efficient use of cumulative sums, all versions of the E and M steps has $\mathcal{O}(N)$ complexity in the size of the counting grid, except for the epitome version. This last version of the counting grid update utilizes the feature layout of the original images z_i^t , which requires the a convolution operation, still manageable in a $\mathcal{O}(N \log N)$ complexity.

More generally most of the updates of the E and M steps of the algorithm require computing windowed sums

$$\sum_{(i_x, i_y) \in W(k_x, k_y)} \pi_{(i_x, i_y), z} \quad (16)$$

where in the previous formula we explicated the two coordinates of the generic position indices $\mathbf{k} = (k_x, k_y)$ and $\mathbf{i} = (i_x, i_y)$. In Fig.6-i) we show a “slice” of π and we want to compute the sum in the yellow window. These sums can be done efficiently by first computing, in linear time, the cumulative sum

$$\text{cumulative_sum}(\pi_{(k_x, k_y)}) = \sum_{(i_x, i_y) \leq (k_x, k_y)} \pi_{(i_x, i_y)} \quad (17)$$

as illustrated in the second panel of Fig.6-i), and then setting

$$\begin{aligned} \sum_{(i_x, i_y) \in W(k_x, k_y)} f_{(i_x, i_y)} &= F_{(k_x+W_x+1, k_y+W_y+1)} \\ &- F_{(k_x, k_y+W_y+1)} \\ &- F_{(k_x+W_x+1, k_y)} \\ &+ F_{(k_x, k_y)} \end{aligned} \quad (18)$$

which is illustrated by Fig. 6-ii). This procedure is used to compute all window histograms h in the counting grid, as

well as in either of the M step versions Eqs. 8 and 13, which only use the counts c_z^t , and not the original feature layout $z_{i,j}^t$.

Efficiency of the computation over multiple section bags in Eqs. 12, 13 can be increased if the sections break the window uniformly along both directions. In this case, one can pre-compute the sum $\sum_{i \in W_k^s} \pi_{i,z}$ in each section and keep $h_{k_x+\tau_s, k_y+\tau_s, z}^s$'s which are shifted versions of each other as Fig. 6-iii) (remember that each sector contribute to the same mapping!).

6 LAYOUT RECONSTRUCTION

³ In scene/object classification tasks, the image features are typically clustered around hundreds of centers and image locations \mathbf{i} are associated with pointers z to these discretized features. For example, in our classification experiments below, we clustered SIFT [37] features in $Z=200$ visual words. The illustrations in Fig. 2 and Fig. 3 do not provide enough insight into how well the counting grids can be inferred when such large sets of features are considered. Visualizing the feature identities on a grid is difficult, and so, in order to simply study the properties of the counting grid estimation procedures discussed above, we have run the first set of tests on fifty 16×16 color patches taken at random from a drawing (available in Matlab: load trees) sub-sampled to the resolution of 33×40 . The drawing is illustrated in Fig. 7-i).

The patches are first transformed into feature maps z_i^t pointing to one of $Z=64$ colors obtained by approximating the color map. Then, 1×1 , 2×2 and 4×4 histograms were computed in the appropriate sections of these images to obtain the section bags of words for the algorithm defined by the appropriate equations (see Tab. 4). The algorithm is then run on each section bag representation separately, to obtain the counting grids in ii), iii), and iv). Finally, the plate v) shows the result of the combination of the counting grid E step, i.e. mapping of the windows based only on the single bag of words, Eq. 4, and the epitome M step, Eq. 15, which uses the known layout of features z_i in counting grid re-estimate under the assumption that this layout could help arrange features in the counting grid even more than a coarse tessellation.

³ In the additional material we added videos that better describe this section and the learning procedure.



Fig. 7. The source of 50 image patches taken from random locations i), and counting grids estimated by various versions of the algorithm. Most remarkably ii) is the reconstruction obtained using *only* 50 histograms of image features, and for reconstruction in iii) we used only 50 sets of 4 histograms (from 2×2 sections of the input images). iv) Result using 16 histograms, (from 4×4 sections of the input images) v) Result using cg-epitomes. In all the cases, colors were treated as unrelated 64 discrete features. SEE THE VIDEOS IN THE ADDITIONAL MATERIAL!

To visualize the different counting grids, each counting grid location \mathbf{k} was assigned the color equal to the average of the $Z=64$ colors in color map, weighted by the normalized local feature counts $\pi_{\mathbf{k},z}$. The image in ii) is therefore an attempt at reconstructing the image in i) from fifty color histograms for which we did not provide any additional information about their source: Image i) was not provided to the algorithm, nor were the locations of the images from which the fifty histograms were extracted. Note also that the algorithm is not aware of any similarities among the 64 colors, as these are treated as discrete features.

Remarkably, a lot of the spatial structure in feature distributions was reconstructed from these 50 histograms. The algorithm discovers that the dark, red and brown tones go together and that they are bordered by green. Elongated dark structures against the blue background are discovered, as is the coast/island boundary. In this sense, the counting grid provides a good model for interpolating among the original 50 histograms, as the histograms from the original image are also likely under the inferred counting grid. Using 2×2 bags as a representation of images is already sufficient to break some symmetry problems and reconstruct almost the entire scene. This improvement is also remarkable, as in this case, ostensibly very little information about the 50 image patches is used: The source image i), or locations of the 50 patches in it are again **not** available to the algorithm, and the algorithm only uses fifty sets of 4 histograms (upper left, upper right, lower left, lower right) over $Z=64$ colors found in appropriate sections to reconstruct the island and the trees. The most accurate reconstruction is obtained in v) by iterating Eqs. 4, 15), which is interesting from the epitome modeling point of view. If the counting grid is considered a feature epitome (as used at low resolutions in [23]), from which detailed feature maps z_i^t are generated, rather than simply bags of features, then the inference step that only considers the patch histograms efficiently replaces the convolutional E step of the epitome model (if it were extended to have feature distribution in each image location, rather than real-valued Gaussian models). Furthermore, in this case we also found that this combination is less prone to local minima than the epitome models or the pure counting grid inference and learning of Eqs. 12, 13. Finally we note here that in the extreme case of tessellating the patches down to individual pixels, the counting grid becomes the feature epitome model.

These results are possible, of course due to very high redundancy in images which makes, for example, the extracted 50×64 count numbers that represent the image patches used for reconstruction of ii) sufficient for this partial recovery of the $33 \times 40 \times \log 64$ parameters necessary to represent i). We next

show that these procedures can be used to analyze images that are related by the fact that they belong to the same category, rather than a large image, and that the resulting generalization over the space of possible bag of feature count distributions far surpasses the standard count models including other latent models, such as latent Dirichlet allocation [7].

7 EXPERIMENTAL SECTION

In all the experiments as visual words we used SIFT features [37] clustered into $Z = 200$ discretized features. The SIFT processing was based on 16×16 pixel patches spaced 8 pixels apart. In this way, each image was transformed into a feature map z_i and then its bag of features c_z was created.

For a fair comparison, we used our implementations of the reconfigurable part-model [16] and latent Dirichlet allocation [3] on the very same features.

In each task, unless specified, we employed the dataset author’s training/testing/validation protocol. To classify a test image we learned a model per class and we assigned the test samples to the class that gives the lowest free energy.

We considered counting grids of various complexities with grid size $\mathbf{E} = [2$ (e.g., 2×2), $\mathbf{3}$ (e.g., 3×3), \dots , $\mathbf{10}$, $\mathbf{15}$, $\mathbf{20}$, \dots , $\mathbf{40}]^4$ and window size $\mathbf{W} = [2, 4, 6, \dots]$, limiting the tests only to the combinations with overall capacity $\kappa = \frac{E_x \cdot E_y}{W_x \cdot W_y}$ between 1.5 and $T/2$, where T is the number of training samples. We considered $\mathbf{S} = [1 \times 1, 2 \times 2, 4 \times 4, N_x \times N_y]$ and we updated $P(\ell)$ with Eq. 11.

The capacity κ is roughly equivalent to the number of LDA topics as it represents the number of independent windows that they can be fit in the grid; we compared the results using this parallelism [18,22].

7.1 Scene Classification

Scene classification task is useful to shows that counting grids can generalize well even when the most basic spatial interpolation assumption is not perfectly met. In particular we will empirically demonstrate that each individual image can be thought a “window” in a larger visual word represented by the counting grid. This has been previously illustrated in Fig. 2 where sampling windows gave rise to the features combinations present in the dataset⁵.

4. We only considered squared counting grids; where $\mathbf{E} = \mathbf{N}$ stands for $\mathbf{E} = N \times N$. The same holds for the window.

5. With the prior $P(\ell)$ possibly preventing to pick some combination

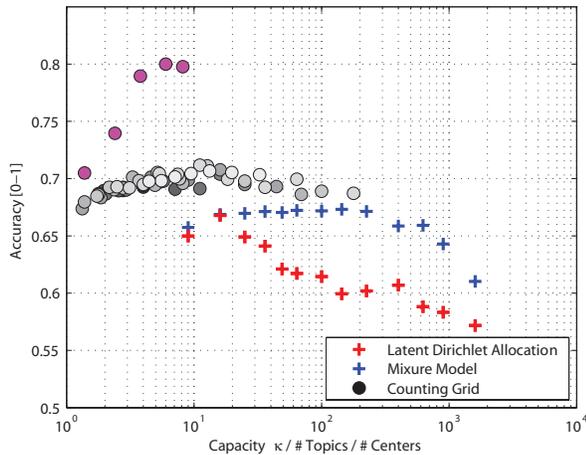


Fig. 8. 15-Scenes classification results. Using pink circles, we also reported the results of [22] which are computed using the hybrid CG-Epitome (see Tab.4). Due to the presence of many training images, the method generalizes very well

TABLE 3
15-Scenes dataset results

Method	Citation	Tessellation	Accuracy
Mixture Model		1×1	59,88%
LDA	[3]	1×1	65,12%
Rec. Part Model	[16]	4×4	74,98%
Spatial BoW		4×4	73,26%
Counting Grid		1×1	72,21%
Tess. Counting Grid		4×4	74,48%
Hybrid CG-Epitome		$N_x \times N_y$	82,79%
Spatial Pyramid Kernel	[10]	4×4	79,93%

As datasets we considered the 15-Scenes [10] and the 67-Indoor Scenes [41]. Classification accuracy on the former are reported in Fig. 8. On the x-axis we reported the different model complexities, in term of capacity κ , whereas on the y-axis we reported the accuracy. As the same κ can be obtained with different choices of \mathbf{E} and \mathbf{W} , we specified the counting grid size \mathbf{E} by using gray levels, the lighter the marker color the bigger the grid.

As Fig. 8 shows, counting grids performed better than latent Dirichlet allocation. The accuracy regularly increased with κ , independently from the Grid size \mathbf{E} . It also worth noticing that $P(\ell)$ helped to prevent overtraining for big capacities κ . In the same figure, we also reported the results of hybrid CG-Epitome approach which comprises the basic CG’s E-step and the epitome M-step. For efficiency reasons, we only considered $\kappa = 1.5, 2.5, 4, 6, 8$. In this version $\mathbf{W} = N_x \times N_y$ and the grid size is unequivocally determined by κ , therefore we used pink markers to show the results. The hybrid CG generalized very well, probably because of the abundance of training data. In Tab. 3 we reported a numerical comparison with other models and some discriminative baseline. For counting grids as well as [3] and [16], we used 3-Fold crossvalidation on the training set to pick a model complexity.

As second dataset, we considered the 67-indoor scene [41] (we did not use the annotations). Results are reported in Tab. 4, where the tessellated counting grid outperformed all the other generative approaches.

TABLE 4
MIT 67 Indoor Scenes dataset results

Method	Citation	Tessellation	Accuracy
Mixture Model		1×1	14,31%
LDA	[3]	1×1	24,53%
Rec. Part Model	[16]	4×4	25,32%
Spatial BoW		4×4	20,94%
Counting Grid		1×1	25,42%
Tess. Counting Grid		4×4	28,32%
Hybrid CG-Epitome		$N_x \times N_y$	16,21%
Spatial Pyramid Kernel	[10]	4×4	32,12%

TABLE 5
SenseCam dataset results

Method	Citation	Tessellation	Accuracy
Mixture Model		1×1	41,19%
LDA	[3]	1×1	57,05%
Rec. Part Model	[16]	4×4	58,17%
Spatial BoW		4×4	49,10%
Counting Grid		1×1	55,32%
Tess. Counting Grid		4×4	59,83%
Hybrid CG-Epitome		$N_x \times N_y$	39,40%
Spatial Pyramid Kernel	[10]	4×4	52,76%

7.2 Place Classification

Recently in [18] a 32-classes dataset have been introduced. This dataset is a subset of the whole visual input of a subject who wore a wearable camera for few weeks. Images in the dataset exhibit dramatic viewing angle, scale, illumination variations and a lot of foreground objects, and clutter. Each category presents images taken in a particular *place* such as house rooms or office environments, or outdoors locations. Some images for each class are shown in Fig. 9.

The task here is place classification. As validation protocol, we used 10-folds cross evaluation. Results are summarized in Fig. 10. In the bag-of-words scenario, e.g., $\mathbf{S} = 1 \times 1$, latent Dirichlet allocation [3] performed better than regular counting grids and mixture models. This can be explained with local minima issues as some classes have a very limited number of training samples and the counting grid simply cannot well recover the panoramic structure (although this is not perfectly evident or recoverable) of half of the classes. Once we provide some directionality information (coarse tessellations $\mathbf{S} = 2 \times 2$) counting grids can better exploit the panorama and they outperformed significantly LDA [3] and its naive tessellated extension which learns a model in each sector, summing the \mathbf{S} likelihoods. Finally in the last panel (Fig. 10-iii) we compared $\mathbf{S} = 4 \times 4$ tessellated counting grid, again the tessellated latent Dirichlet allocation and the Reconfigurable part model [16] which uses the same spatial information. Finer tessellations didn’t help recognition but neither hurt up to $\mathbf{S} = 6 \times 6$. To the limit, when $\mathbf{S} = N_x \times N_y$, accuracy does not exceed 30%.

As final experiment on this dataset, we repeated the experiment only using 13 training images per class as previously done in [18]. Here we want to test the robustness of the models in overtrain regimes. We reported the final accuracy in Tab. 5.

Summarizing counting grids map images onto a bigger real estate, where they lay out the features into a 2D window and stitch overlapping windows trying to recover the panoramic nature of the scene. This fits the qualities of the data acquired by a wearable camera and indeed our model largely outperform [3,16].



Fig. 9. Images from the SenseCam dataset.

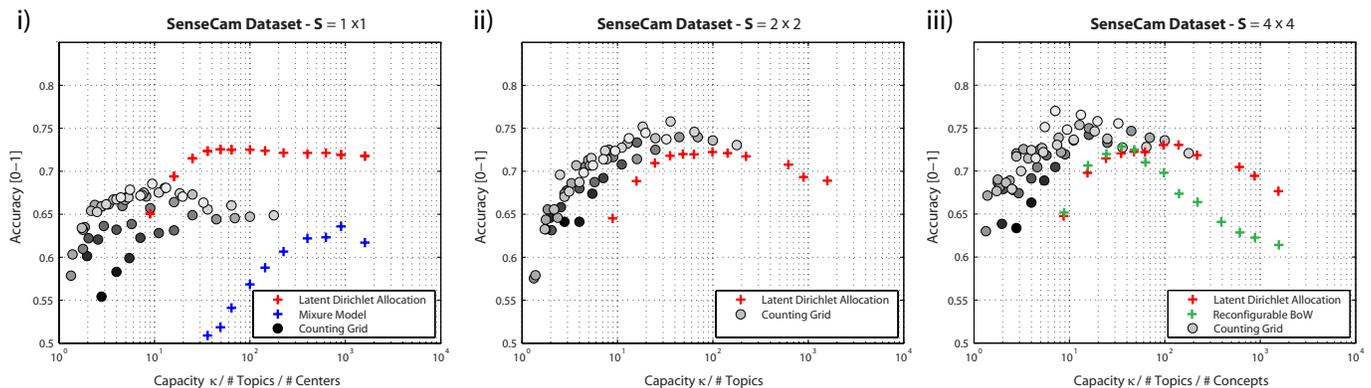


Fig. 10. Results for SenseCam dataset. i) $S = 1 \times 1$. ii) Tessellated version $S = 2 \times 2$. iii) Tessellated version $S = 4 \times 4$ and comparison with the reconfigurable bag of words model [16] and with latent Dirichlet allocation using the same tessellation.

7.3 Wearable Camera Sequences

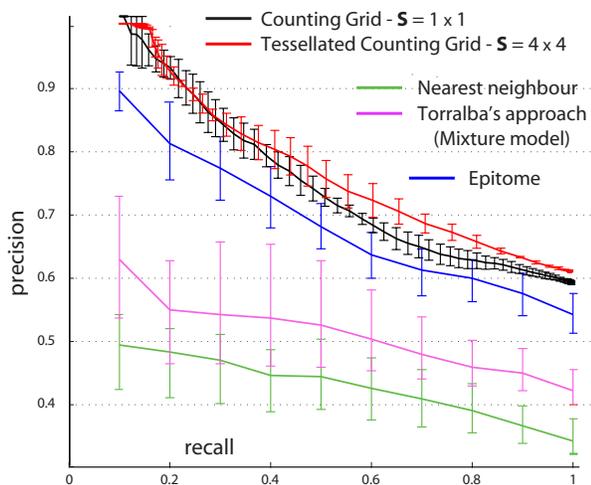


Fig. 11. Results on Torralba Dataset. We reported the results of [17] (Torralba's approach) and [23] (Epitome) from the original papers. We followed the evaluation procedure of [17] and the error bars indicate variability in accuracy across different image sequences.

We considered the sequences of [17]. This data represents the perfect fit for our model as the true panoramic structure

of each scene or place, can actually be recovered. The dataset is composed by 7 video sequences acquired with a wearable camera.

The original paper [17] is based on learning a Gaussian mixture model for each class, using Gist [30] as image descriptor. In addition to [17], we also compared with Epitomes [23] which was, among applications of epitome, one of the most successful. The method of [23] uses a low resolution epitome with each low res image location represented by a histogram of features. This method combines several cues: RGB (local) histograms, disparity features and Gist. For what concern counting grids, we only used quantized SIFT and we set the complexity of the model using cross-validation considering only models with capacity $2 \leq \kappa \leq 10$.

After training a model for each scene $l = 1 \dots C$ our goal is to compute the place posterior probabilities for every frame t of the test sequence, given all the previous images $P(l^t = k | \mathbf{c}^{1:t})$. This can be easily achieved using the forward-backwards procedure [42]

$$P(l^t = k | \mathbf{c}^{1:t}) \propto p(\mathbf{c}^t | l^t = k) \cdot \sum_j P(l^t = k | l^{t-1} = j) \cdot P(l^{t-1} | \mathbf{c}^{1:t-1}) \quad (19)$$

We fixed the observation log likelihood to the negative free energy given by our model (Eq. 3) while we used EM estimate the transition matrix and the place posteriors. When using HMM, the observation likelihood may be dominated by the

TABLE 6
Where was I?

Method	Citation	Tessellation	Accuracy
LDA	[3]	1×1	76.80%
Rec. Part Model	[16]	4×4	74.63%
Mixture Model		1×1	70.37%
Counting Grid		1×1	76.21%
Tess. Counting Grid		4×4	81.45%

transition prior. To balance the contribution we re-scaled the likelihood terms using a constant γ , chosen via cross-evaluation [43].

Results are presented in Fig. 11; the improvement wrt [17,23] is significant. The tessellation marginally helped because *i*) training data is abundant and *ii*) the metaphor upon which CGs are based, the “moving camera”, perfectly fits here. Indeed the spatial layout can be at least piece-wise recovered also from a single bag.

Tessellation finer than $S = 2 \times 2$ did not hurt. Latent Dirichlet allocation [7] and Rec-Bow [16] performances were slightly inferior of [23] and we did not report it in the graph for the sake of clarity.

We have also investigated what happens if we equally scale \mathbf{E} and \mathbf{W} . We considered counting grids of size $\mathbf{E} = \sigma \cdot [8, 10, 12, 15, 18, 24]$, $\mathbf{W} = \sigma \cdot \mathbf{6}$ and three scales $\sigma = 1, 2, 3$ and we run the same experiment on Torralba’s sequences. Results are shown in Fig. 12, where each row represents a different scale.

Results are easily interpretable, counting grids are not very sensitive to the choice of \mathbf{E} and \mathbf{W} and what really matters is their ratio κ . This can also be evinced by Fig. 8 and Fig. 10 where complexities characterized by similar κ performed equally well. Higher variances for large κ , indicate local minima issues.

In general, once the window is “sufficiently big” for spatial interpolation, scaled models learn “scaled” versions of the scene, which are quantitatively (and qualitatively) very similar. The real estate is too big and the model learn multiple copies of the same scene.

We have finally considered a day worth of images from (1800 images ca.) from the SenseCam collection [44] and repeated the same test, combining counting grids an hidden Markov models. During this day, the camera bearer visited 20 of the 32 labeled locations of the full dataset [18], nevertheless we trained models with all the 32 classes as a-priori we cannot know the locations visited during a day. As for Torralba sequences, our goal is to compute the place posterior probabilities at the instant t , given all the previous images, Eq. 19 We used at most 30 images per class to learn the models. Results are reported in Tab. 6. We run [17] using a mixture of dirichlet model over quantized sift histograms (our very same features). For sake of completeness we also implemented the method of [17] extracting the original descriptors from whole images and within the four sectors. In both cases, the performance was below 50%.

7.4 Image clustering on SenseCam

As final test, we analyzed the same subset of SenseCam, divided in 10 categories used in [28]. The images of this subset are suitable for epitomes as they can actually be stitched together using pixels, therefore a comparison with [23,27,28] is

TABLE 7
Unsupervised place clustering

Method	Citation	Tessellation	Accuracy
Epitome	[27]	$N_x \times N_y$	69,42%
Stel Epitome	[28]	$N_x \times N_y$	73,06%
LDA	[3]	1×1	74,32 %
Counting Grid		1×1	82.34%
Tess. Counting Grid		4×4	83.94%
Hybrid CG-Epitome		$N_x \times N_y$	86,6%
Feature Epitome		$N_x \times N_y$	69,93%

fair.

As the spirit of the data collection is to provide summary of the subject’s life, we have trained the counting grids in an unsupervised way (combining images of all categories together) and then investigated if the images are separated in the counting grid in accordance to the human labeling. We compared with other “visual summarization” approaches that lay out the visual input on a larger grid, the epitomic approaches [23,27,28] which clusters pixel measurements within an epitome. While in the standard epitomes images are mapped into the epitome by means of pixel wise comparisons, here we are placing bags in a 2 dimensional space, i.e. an image is mapped in a particular spot if its bag-of-word representation agrees with the images mapped in the neighborhood. To make the comparison fair, we fixed the complexity of the counting grids to the one used for epitomes in [28] (e.g., $\kappa = 14$). Upon learning, each test image is labeled by the label of the closest mapped training image. The results are reported in Table 7.

The counting grid model is so far the best performing model on this task.

8 CONCLUSIONS

We introduce the counting grid model of images which captures natural constraints on image feature histograms by assuming that these can be represented by averaging of feature distributions from a window into the grid. In this way, the flexibility of the bag of words representation is indirectly enriched by the spatial constraints of epitome-like models. By observing the actual observation model, we see that the counting grid model is not attempting to capture the spatial constraints explicitly as has been often done in the past. In fact, we can view the counting grid as producing a large mixture of histograms whose parameters are constrained in a way that is a natural consequence of the fact that images from which the features are collected live in an ordered 2D space. Despite their simplicity, both conceptual and algorithmic (the matlab code for counting grid estimation fits half a page), and that the ultimate parametrization used for likelihood computation is simply a set of histograms, this generative model significantly outperforms other histogram-based representations in a variety of tasks and is often approaching the discriminative state of the art (and the features extracted from the generative model can often be used within discriminative models to further improve them [45]). Computationally, the algorithm is efficient and the computational steps also lend themselves to further improvement of the model to add more scale/rotation reasoning. Experiments show that, despite the apparent need of setting \mathbf{E} and \mathbf{W} , the algorithm is only sensitive to their ratio. For what concern performances, counting grids, especially in their tessellated version, outperformed standard bag of words approaches in computer vision [3,16,17,23] across most of the datasets considered. Finally we observe that a variety of

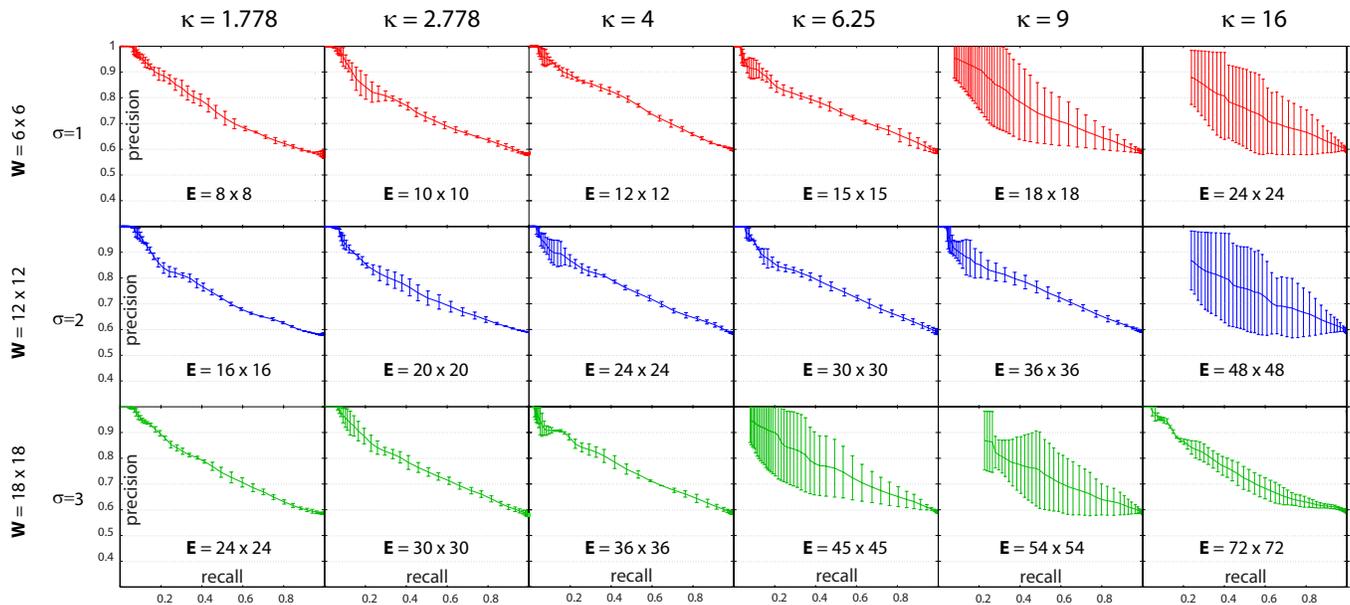


Fig. 12. Scaled counting grids (1×1 -case).

methods are based on latent dirichlet allocation and we would like the community considered our method as “basis” to solve complex problems or perform complex analysis.

REFERENCES

- [1] G. Csurka, C. R. Dance, L. Fan, J. Willamowski, and C. Bray, “Visual categorization with bags of keypoints,” in *In Workshop on Statistical Learning in Computer Vision, ECCV, 2004*, pp. 1–22.
- [2] J. Yang, Y. G. Jiang, A. G. Hauptmann, and C. W. Ngo, “Evaluating bag-of-visual-words representations in scene classification,” in *MIR '07: Proceedings of the international workshop on Workshop on multimedia information retrieval*. New York, NY, USA: ACM, 2007, pp. 197–206. [Online]. Available: <http://dx.doi.org/10.1145/1290082.1290111>
- [3] L. Fei-Fei and P. Perona, “A bayesian hierarchical model for learning natural scene categories,” in *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 2005*, pp. 524–531.
- [4] P. Langley, W. Iba, and K. Thompson, “An analysis of bayesian classifiers,” in *Annual Conference on Artificial Intelligence*. MIT Press, 1992, pp. 223–228.
- [5] N. Bouguila, “Count data modeling and classification using finite mixtures of distributions,” *Neural Networks, IEEE Transactions on*, vol. 22, no. 2, pp. 186–198, 2011.
- [6] K. Nigam, J. Lafferty, and A. McCallum, “Using maximum entropy for text classification,” in *IJCAI - Workshop on Machine Learning for Information Filtering*, 1999.
- [7] D. Blei, A. Ng, and M. Jordan, “Latent dirichlet allocation,” *Journal of machine Learning Research*, vol. 3, pp. 993–1022, 2003.
- [8] T. Hofmann, “Probabilistic latent semantic indexing,” in *Proceedings of the annual international ACM conference on Research and development in information retrieval (SIGIR)*, 1999, pp. 50–57.
- [9] A. Bosch, A. Zisserman, and X. Munoz, “Image classification using random forests and ferns,” in *Proceedings of International Conference on Computer Vision (ICCV), 2007*, pp. 1–8.
- [10] S. Lazebnik, C. Schmid, and J. Ponce, “Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories,” in *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 2006*, pp. 2169–2178.
- [11] O. Boiman, E. Shechtman, and M. Irani, “In defense of nearest-neighbor based image classification,” in *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 2008*, pp. 1–8.
- [12] J. Vogel and B. Schiele, “Semantic modeling of natural scenes for content-based image retrieval,” *International Journal of Computer Vision*, vol. 72, pp. 133–157, 2007.
- [13] A. Oliva and A. Torralba, “Modeling the shape of the scene: A holistic representation of the spatial envelope,” *Int. J. Comput. Vision*, vol. 42, no. 3, pp. 145–175, 2001.
- [14] A. Bosch, A. Zisserman, and X. Munoz, “Scene classification via pls,” in *Proceedings of European Conference on Computer Vision (ECCV), 2006*, pp. 517–530.
- [15] M. Boutell, J. Luo, and C. Brown, “Scene parsing using region-based generative models,” *Multimedia, IEEE Transactions on*, vol. 9, no. 1, pp. 136–146, 2007.
- [16] S. Parizi, J. Oberlin, and P. Felzenszwalb, “Reconfigurable models for scene recognition,” in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, 2012, pp. 2775–2782.
- [17] A. Torralba, K. P. Murphy, W. T. Freeman, and M. A. Rubin, “Context-based vision system for place and object recognition,” in *ICCV*, 2003, pp. 273–280.
- [18] A. Perina and N. Jojic, “Spring lattice counting grids: Scene recognition using deformable positional constraints,” in *ECCV (6)*, ser. Lecture Notes in Computer Science, A. W. Fitzgibbon, S. Lazebnik, P. Perona, Y. Sato, and C. Schmid, Eds., vol. 7577. Springer, 2012, pp. 837–851.
- [19] A. Coates, A. Y. Ng, and H. Lee, “An analysis of single-layer networks in unsupervised feature learning,” in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2011, Fort Lauderdale, USA, April 11-13, 2011*, 2011, pp. 215–223. [Online]. Available: <http://www.jmlr.org/proceedings/papers/v15/coates11a/coates11a.pdf>
- [20] B. C. Russell, A. B. Torralba, K. P. Murphy, and W. T. Freeman, “Labelme: A database and web-based tool for image annotation,” *International Journal of Computer Vision*, vol. 77, no. 1-3, pp. 157–173, 2008.
- [21] N. Jojic and A. Perina, “Multidimensional counting grids: Inferring word order from disordered bags of words,” in *UAI 2011, Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence, Barcelona, Spain, July 14-17, 2011*, 2011, pp. 547–556.
- [22] A. Perina and N. Jojic, “Image analysis by counting on a grid,” in *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 2011*, pp. 1985–1992.
- [23] K. Ni, A. Kannan, A. Criminisi, and J. Winn, “Epitomic location recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 12, pp. 2158–2167, 2009.
- [24] A. Perina and N. Jojic, “Capturing layers in image collections

- with componential models: from the layered epitome to the componential counting grid," in *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [25] M. R. Amer and S. Todorovic, "Sum-product networks for modeling activities with stochastic structure," in *CVPR*, 2012, pp. 1314–1321.
- [26] P. Lovato, A. Perina, N. Sebe, O. Zandonà, A. Montagnini, M. Bicego, and M. Cristani, "Tell me what you like and i'll tell you what you are: Discriminating visual preferences on flickr data," in *ACCV*, 2012, pp. 45–56.
- [27] N. Jovic, B. J. Frey, and A. Kannan, "Epitomic analysis of appearance and shape," in *Proceedings of International Conference on Computer Vision (ICCV)*, 2003, pp. 34–41.
- [28] N. Jovic, A. Perina, and V. Murino, "Structural epitome: a way to summarize ones visual experience," in *Advances in Neural Information Processing Systems*, 2010, pp. 1027–1035.
- [29] N. Jovic and B. Frey, "Learning flexible sprites in video layers," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2001, pp. 199–206.
- [30] A. Oliva and A. Torralba, "Building the gist of a scene: The role of global image features in recognition," *Progress in Brain Research: Visual perception*, vol. 155, pp. 23–36, 2006.
- [31] Y. Jiang, J. Yuan, and G. Yu, "Randomized spatial partition for scene recognition," in *Proceedings of the 12th European conference on Computer Vision - Volume Part II*, ser. ECCV'12. Berlin, Heidelberg: Springer-Verlag, 2012, pp. 730–743.
- [32] A. Perina, M. Cristani, and V. Murino, "Learning natural scene categories by selective multi-scale feature extraction," *Image Vision Comput.*, vol. 28, no. 6, pp. 927–939, 2010.
- [33] A. Perina, N. Jovic, U. Castellani, M. Cristani, and V. Murino, "Object recognition with hierarchical stel models," in *Proceedings of European Conference on Computer Vision (ECCV)*, 2010, pp. 15–28.
- [34] M. Pandey and S. Lazebnik, "Scene recognition and weakly supervised object localization with deformable part-based models," in *ICCV*, 2011, pp. 1307–1314.
- [35] X. Chu, S. Yan, L. Li, K.-L. Chan, and T. Huang, "Spatialized epitome and its applications," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, 2010, pp. 311–318.
- [36] L. Bazzani, M. Cristani, A. Perina, and V. Murino, "Multiple-shot person re-identification by chromatic and epitomic analyses," *Pattern Recognition Letters*, vol. 33, no. 7, pp. 898–903, 2012.
- [37] D. Lowe, "Object recognition from local scale-invariant features," in *Proceedings of International Conference on Computer Vision (ICCV)*, 1999, pp. 1150–1157.
- [38] M. Jordan, Z. Ghahramani, T. Jaakkola, and L. Saul, "An introduction to variational methods for graphical models," *Machine Learning*, vol. 37, no. 2, pp. 183–233, 1999.
- [39] R. M. Neal and G. E. Hinton, "A view of the em algorithm that justifies incremental, sparse, and other variants," *Learning in graphical models*, pp. 355–368, 1999.
- [40] B. Frey and N. Jovic, "Transformation-invariant clustering using the EM algorithm," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 1, pp. 1 – 17, 2003.
- [41] A. Quattoni and A. Torralba, "Recognizing indoor scenes," in *CVPR*, 2009, pp. 413–420.
- [42] L. Rabiner, "A tutorial on Hidden Markov Models and selected applications in speech recognition," *Proc. of IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [43] L. Deng and D. O'Shaughnessy, *Speech Processing: A Dynamic and Optimization-oriented Approach*, ser. Signal Processing and Communications Series. Marcel Dekker Incorporated, 2003. [Online]. Available: http://books.google.ca/books?id=136wRmFT_t8C
- [44] A. Perina and N. Jovic, "In the sight of my wearable camera: Classifying my visual experience," In *2nd IEEE Workshop on Egocentric Vision*, in conjunction with CVPR, Tech. Rep., 2012. [Online]. Available: <http://arxiv.org/abs/1304.7236>
- [45] A. Perina, M. Cristani, U. Castellani, V. Murino, and N. Jovic, "Free energy score spaces: Using generative information in discriminative classifiers," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 7, pp. 1249–1262, 2012.