



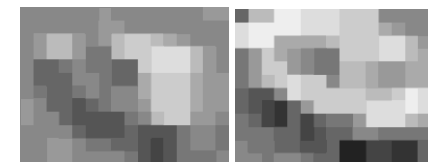
Unsupervised Object Counting without Object Recognition

Takayuki Katsuki, Tetsuro Morimura (IBM Research - Tokyo), and
Tsuyoshi Idé (IBM T. J. Watson Research Center)

How many vehicles are there in this road?



Identification of independent objects is often impossible when we handle



Very low resolution

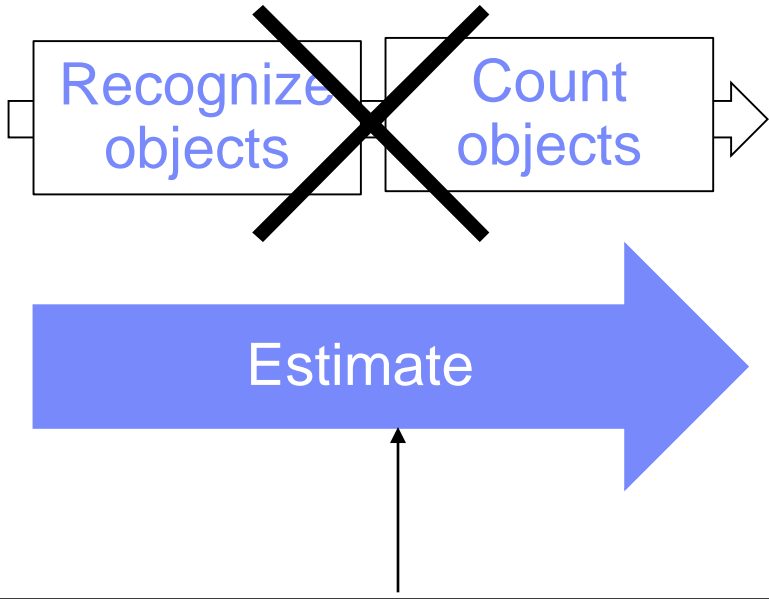


Overlapped



Count objects without recognizing any objects or using any labeled training data

Observation \mathcal{X}



of objects

$$h \in \{0, 1\}^\infty, \sum_{d=0}^\infty h_d = 1$$

We use the training data $\mathbf{X} \equiv \{x_1, x_2, \dots, x_N\}$ **without any label information as to the count.**

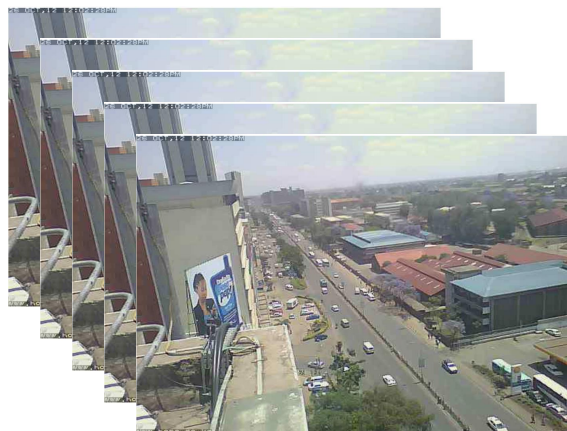




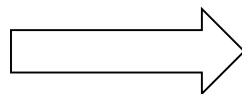
Formalize the problem as a density estimation using a particular type of mixture model

Learning phase

Training data without labels

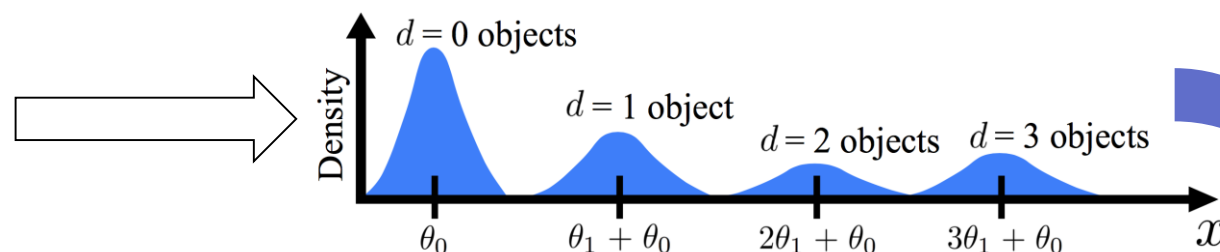


Feature
extraction



\mathcal{X}_1
 \mathcal{X}_2
 \mathcal{X}_3
 \mathcal{X}_4
 \mathcal{X}_5
:
:

Learn a Gaussian mixture model (GMM) whose mixture index is equated with the count of the objects



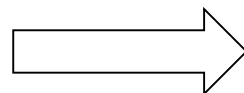
Our idea is to use the **stick-breaking process** as a constraint to make it possible to interpret the mixture indexes as the count.

Runtime phase

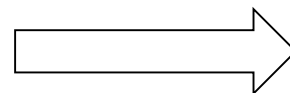
New observation



Feature
extraction

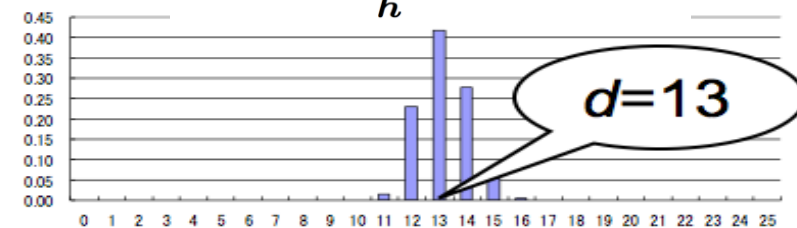


\mathcal{X}



To find the count for a new observation, we pick the cluster of the highest posterior for the count given \mathcal{X} :

$$h^* \equiv \operatorname{argmax}_h p(h|x, \mathbf{X})$$





Experimental Results: Vehicle-Counting from Web Camera Images

We confirmed that the proposed framework can apply to the task of counting vehicles in the web camera images.



(a) Nationkimathi



(b) Westistg



(c) Ukulima

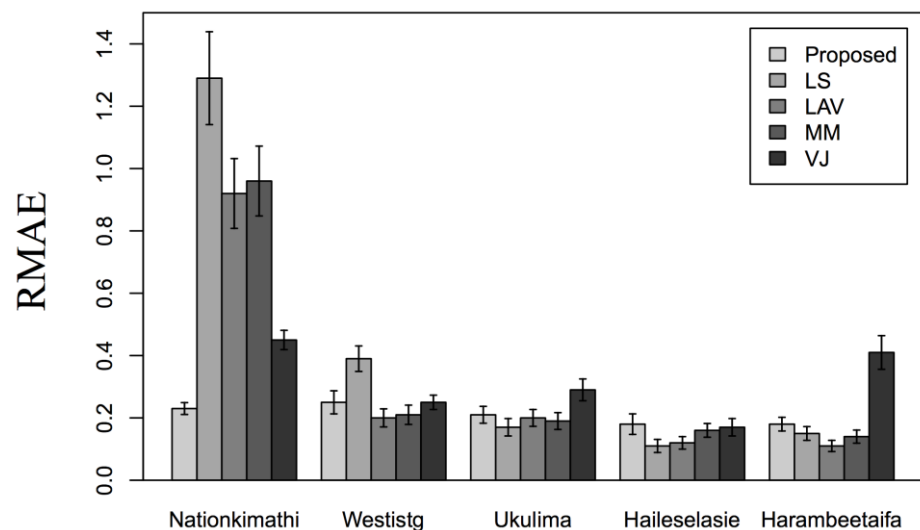


(d) Haileselasie



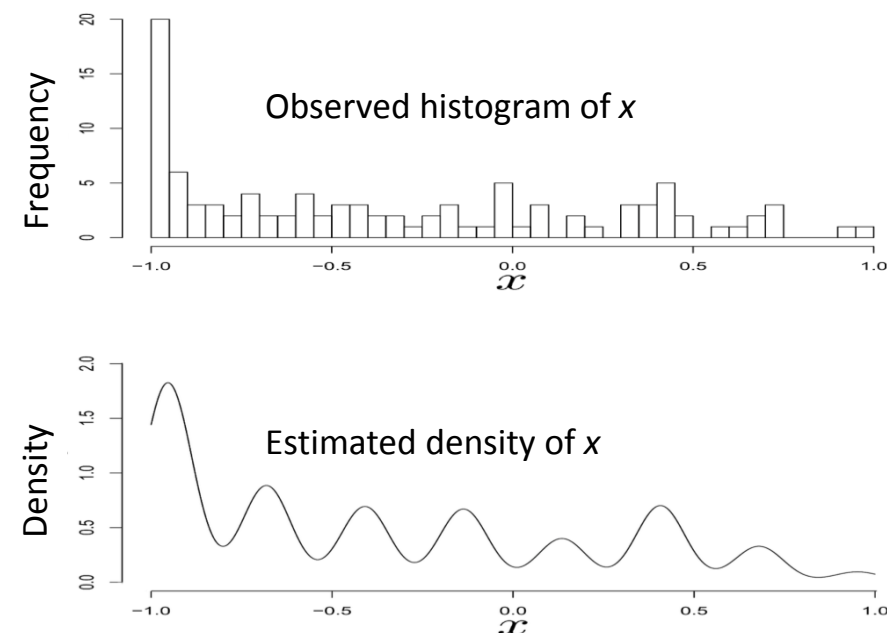
(e) Harambeetaifa

We compared our unsupervised approach with several *supervised* alternatives in RMAEs for all of the camera locations (smaller is better). Error bars represent the standard error.



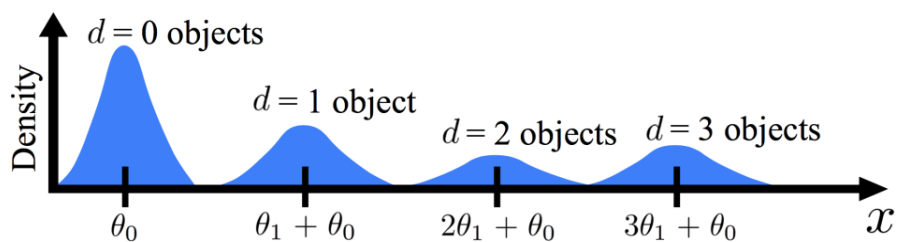
We can see that the overall performance of our method is comparable to or even better than those of the supervised alternatives. This is rather surprising, because our method does *not* use any labeled training data.

Comparison of the estimated $p(x)$ distribution with the true one created from the data.



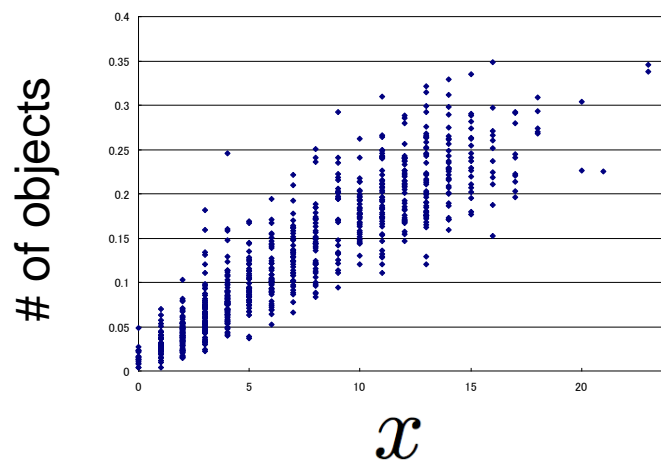
Gaussian Mixture Model whose d -th component is responsible for x having d number of objects

The proposed GMM has **a restriction on its mean parameter**:



$$p(x|\mathbf{h}, \boldsymbol{\theta}, \beta) \equiv \prod_{d=0}^{\infty} \mathcal{N}(x|\theta_1 d + \theta_0, \beta^{-1})^{h_d} \\ = \frac{\exp\left(-\frac{\beta}{2} \sum_{d=0}^{\infty} h_d (x - \theta_1 d - \theta_0)^2\right)}{(2\pi\beta^{-1})^{\frac{1}{2}}}$$

We loosely assume that feature x is a good enough feature in the sense that it is proportional to the count.

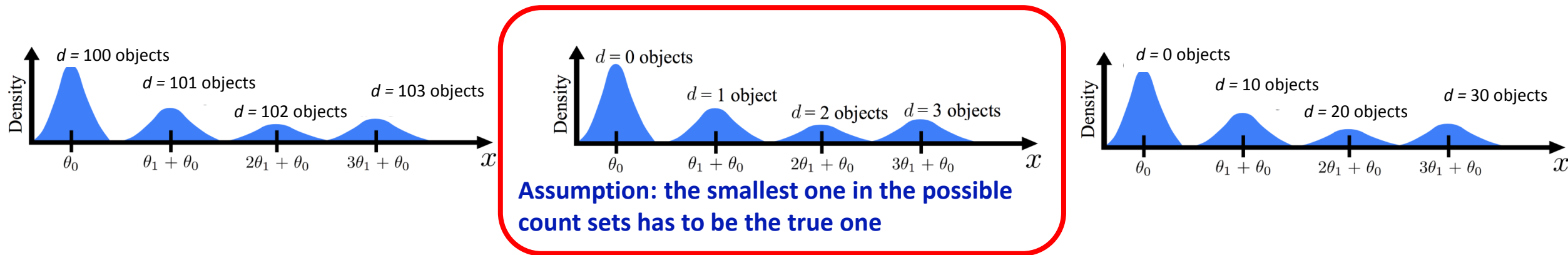


*However, the counting results of the proposed GMM without any additional constraint will become **linearly proportional to the true count.***



Regularize density estimation results of the GMM by stick-breaking process prior

The likelihood of the count h in the proposed GMM is invariant with respect to the simultaneous translation of x and ϑ_0 , as well as the simultaneous scaling between count d and ϑ_1 , such as,



We introduce the stick-breaking process as the prior for the count, which can represent the desired property of the count set: the assigned count values for the observations are consecutive natural numbers from zero.

Stick-breaking Process
$$p(\mathbf{h}|\mathbf{v}) \equiv \prod_{d=0}^{\infty} \left(v_d \prod_{k=0}^{d-1} (1 - v_k) \right)^{h_d}$$

Interestingly, in traditional Bayesian nonparametric literature, this nature, which is each component is always associated with the same index, is known as a drawback; that is, it can cause the solution to get stuck at a local minimum in practical use.



Estimate posterior through efficient variational Inference algorithm

Using the proposed model and the conjugate priors, we can rewrite the posterior as

$$p(\mathbf{h}|x, \mathbf{X}) \propto \int p(x, \mathbf{h}, \mathbf{X}, \mathbf{H}, \boldsymbol{\theta}, \beta, \mathbf{v}) d\mathbf{H} d\boldsymbol{\theta} d\beta d\mathbf{v},$$

where $p(x, \mathbf{h}, \mathbf{X}, \mathbf{H}, \boldsymbol{\theta}, \beta, \mathbf{v}) \equiv p(x|\mathbf{h}, \boldsymbol{\theta}, \beta)$

$$\times p(\mathbf{h}|\mathbf{v}) \left[\prod_{n=1}^N p(x_n|\mathbf{h}_n, \boldsymbol{\theta}, \beta) p(\mathbf{h}_n|\mathbf{v}) \right] p(\boldsymbol{\theta}) p(\beta) p(\mathbf{v})$$

We approximate the posterior in a factorized form:

$$q(\mathbf{h}, \mathbf{H}, \boldsymbol{\theta}, \beta, \mathbf{v}) \equiv q(\mathbf{h}, \mathbf{H}) q(\boldsymbol{\theta}) q(\beta, \mathbf{v})$$

We then identify the optimal q that minimizes

$$D_{\text{KL}}(q||p) \equiv \int q(\ln q - \ln p) d\mathbf{H} d\boldsymbol{\theta} d\beta d\mathbf{v}$$

Finally, thanks to the conjugate modeling, we can get the iterative updating equations as

$$q(\mathbf{h}) q(\mathbf{H}) = \text{Categorical}(\mathbf{h}|\boldsymbol{\mu}_{\mathbf{h}}) \times \prod_{n=1}^N \text{Categorical}(\mathbf{h}_n|\boldsymbol{\mu}_{\mathbf{h}_n}),$$

$$q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}}), \text{ and}$$

$$q(\beta, \mathbf{v}) = \text{Gamma}(\beta|a_{\beta}, b_{\beta}) \text{Beta}(\mathbf{v}_d|a_{\mathbf{v}_d}, b_{\mathbf{v}_d}),$$

We stop the VB iterations when this condition is satisfied

$$\frac{(D_{\text{KL}}(q||p) - D_{\text{KL}}(q'||p))^2}{D_{\text{KL}}(q'||p)^2} < 10^{-10}$$

We can estimate the number of objects in the new observation as

$$\mathbf{h}^* \simeq \underset{\mathbf{h}}{\text{argmax}} q(\mathbf{h})$$



Future work

- Including many other features and introducing a non-linear relationship in the proposed GMM would be an important research area
- Applying the proposed approach to other applications, such as crowd counting and cell counting, would be another promising area of study

