

# Visual feature learning

Winter 2013

Roland Memisevic

Lecture 9, February 19, 2013



## Energy based models

- ▶ We may eliminate the partition function altogether and define the model as an “energy landscape” that we form through learning.
- ▶ This gives us even more freedom in devising schemes that push or pull on the energy landscape.
- ▶ (LeCun, et al. 2006)

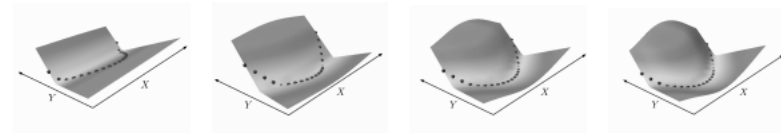


## Last time

- ▶ Limitations of ICA
- ▶ Overcomplete codes
- ▶ Energy based models



## Energy based models

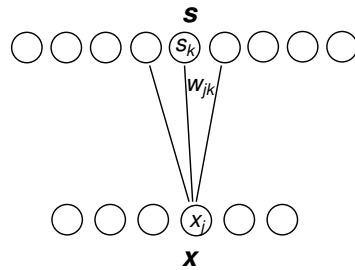


- ▶ It is common to define energies as  $-q(\mathbf{x}; W)$ , in which case we want to *minimize* energy near the data.
- ▶ Energy based models can be used in a variety of tasks, but for feature learning, they practically always involve hidden variables which can be visualized conveniently as a bi-partite graph.
- ▶ Pushing up the energy away from the data typically translates into a *capacity constraint* on the hidden variables.



## Feature learning and bi-partite networks

$$\mathbf{s} = W^T \mathbf{x}$$

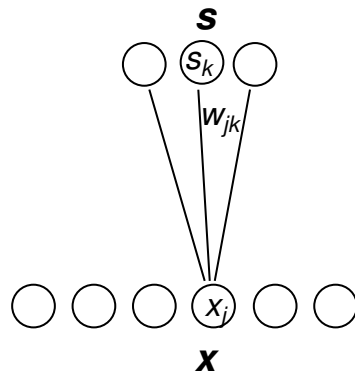


- ▶ Most feature learning models are based on a variation of the encoder/decoder equations.
- ▶ PCA is a special case with linear dependencies and low-dimensional  $\mathbf{s}$

## Feature learning models

$$\mathbf{s} = W^T \mathbf{x}$$

$$\mathbf{x} = W \mathbf{s}$$

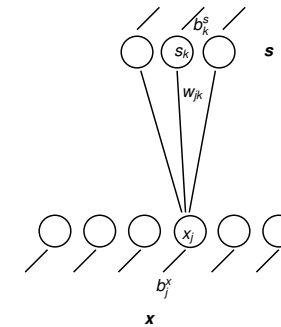


### PCA

- ▶ Training: minimize  $\|\mathbf{x} - W\mathbf{s}\|^2 = \|\mathbf{x} - W^T W \mathbf{x}\|$  s.t. orthogonality constraint.

## Feature learning

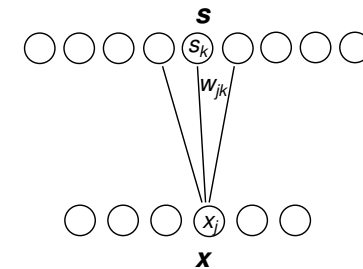
$$\mathbf{s} = W^T \mathbf{x} + \mathbf{b}^s$$



- ▶ In practice, one typically adds *bias terms* to obtain affine not linear dependencies.
- ▶ But one often drops these in derivations.

## Feature learning models

$$\mathbf{x} = W \mathbf{s}$$



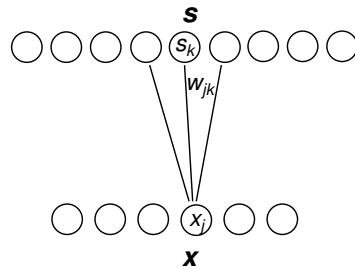
### Sparse coding (Olshausen/Fields)

- ▶ Energy:  $\sum_i \|\mathbf{x}_i - W\mathbf{s}_i\|^2 + \sum_i |\mathbf{s}_i|$  (optimize wrt.  $\mathbf{s}_i$  and  $W$ )
- ▶ Inference: *Search* for  $\mathbf{s}$  that minimizes the same cost.

## Feature learning models

$$\mathbf{s} = \text{sigmoid}(W^T \mathbf{x})$$

$$\mathbf{x} = \text{sigmoid}(W \mathbf{s})$$



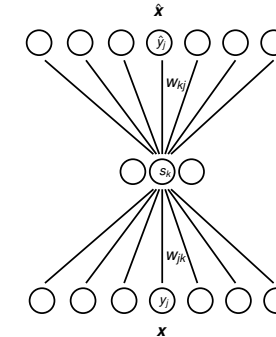
### (Binary) Restricted Boltzmann Machine (RBM)

- ▶ Density:  $p(\mathbf{x}, \mathbf{s}) = \frac{1}{Z} \exp(\sum_{jk} w_{jk} x_j s_k)$
- ▶ Contrastive divergence learning (Hebbian/anti-hebbian learning rule)
- ▶ Probability model with  $G_i(\mathbf{x}) = (1 + \exp(\mathbf{w}_i^T \mathbf{x}))$

## Feature learning models

$$\mathbf{s} = \text{sigmoid}(W^T \mathbf{x})$$

$$\mathbf{x} = W \mathbf{s}$$



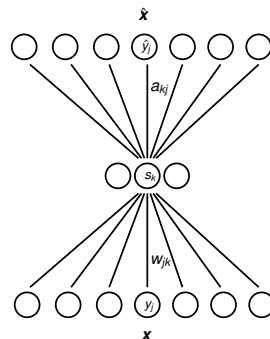
### Autoencoder

- ▶ Energy:  $\sum_i \|\mathbf{x}_i - W \text{sigmoid}(W^T \mathbf{x}_i)\|^2$
- ▶ Can add a sparsity penalty for  $\mathbf{s}$ , or *corrupt inputs during training* (Vincent et al., 2008)

## Feature learning models

$$\mathbf{s} = \text{sigmoid}(A^T \mathbf{x})$$

$$\mathbf{x} = W \mathbf{s}$$



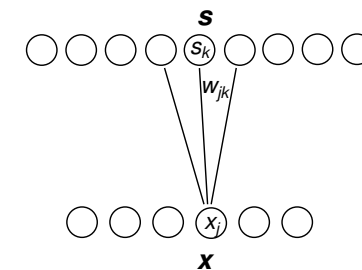
### Autoencoder with "untied" weights

- ▶ Energy:  $\sum_i \|\mathbf{x}_i - A \text{sigmoid}(W^T \mathbf{x}_i)\|^2$
- ▶ Can add a sparsity penalty for  $\mathbf{s}$ , or *corrupt inputs during training* (Vincent et al., 2008)

## Feature learning models

$$\mathbf{s} = \text{wta}(W^T \mathbf{x})$$

$$\mathbf{x} = W \mathbf{s}$$



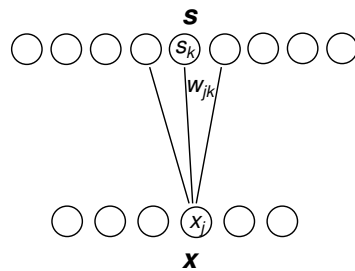
### k-means clustering

- ▶ Energy:  $\sum_i \|\mathbf{x}_i - W \mathbf{s}\|^2$  with  $\mathbf{s}$  inferred as follows:
- ▶ Inference: Set  $s_i = 1$  for the  $i$  with maximal  $\mathbf{w}_i^T \mathbf{x}$ , 0 otherwise ("winner-takes-all")

## Feature learning models

$$\mathbf{s} = \text{wta}(\mathbf{W}^T \mathbf{x})$$

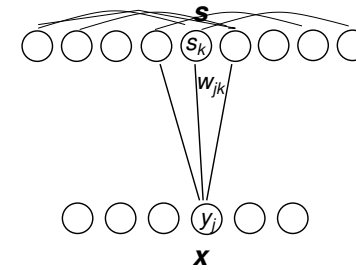
$$\mathbf{x} = \mathbf{W}\mathbf{s}$$



### Mixtures of Gaussians

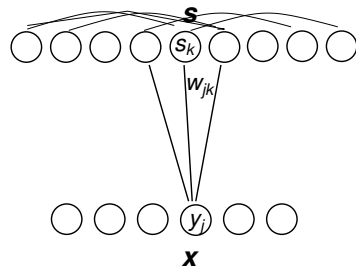
- ▶ Probabilistic version of kmeans.
- ▶ Density:
 
$$p(\mathbf{x}) = \sum_k p(s_k) p(\mathbf{x} | s_k) = \sum_k p(s_k) \mathcal{N}(\mathbf{x} | \mathbf{W}_{\cdot k}, \Sigma_k)$$
- ▶ Inference: Bayes' rule.

## Lateral interactions



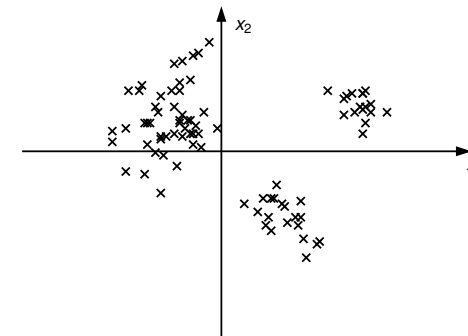
- ▶ In some of these models, inference requires **lateral interactions** or **feedback**.
- ▶ This means that hidden units need to talk to each other to figure out their joint activity pattern.
- ▶ In most cases, the interactions are *inhibitory* and lead to *competition*.
- ▶ Examples include sparse coding, k-means, MoG.

## Lateral interactions



- ▶ *End-stopping* is one known effect that can be explained with competition among simple cells.
- ▶ A process related to lateral interaction is top-down feedback from higher processing layers.

## K-means clustering



- ▶ K-means is traditionally a clustering algorithm.
- ▶ Learning: Fit  $K$  prototypes  $\mu_k$  (these will be the rows of  $\mathbf{W}$ ) to training data-points  $\mathbf{x}_n$ .
- ▶ Inference: Given a point, find the nearest prototype.

## K-means clustering

- ▶ Think of  $\mathbf{s}$  as a one-hot encoding of the discrete variable representing the index of the cluster center.
- ▶ There are  $K$  prototypes  $\mu_1, \dots, \mu_K$  that represent the  $K$  clusters. The dimensionality of the prototypes is the same as that of the data  $\mathbf{x}$ .
- ▶ Assume we *knew* the cluster assignments  $\mathbf{s}_n$  for each point  $\mathbf{x}_n$ .
- ▶ The  $K$ -means objective function measures the *average distance between points  $\mathbf{x}$  and their representatives*:

$$J = \sum_{n=1}^N \sum_{k=1}^K s_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍

Roland Memisevic

Visual feature learning

## K-means clustering

### Finding the optimal $\mathbf{s}_n$

- ▶ Note that given the  $\mu_k$ , we can optimize all the  $\mathbf{s}_n$  independently, because the objective is just the sum over  $n$ .
- ▶ But the squared error will be smallest if we set  $s_{nk} = 1$  for whichever  $\mu_k$  is *closest*.
- ▶ Formally, to optimize all  $\mathbf{s}_n$ , given the set of  $\mu_k$ , set:

$$s_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍

Roland Memisevic

Visual feature learning

## K-means clustering

- ▶ *Learning* amounts to finding *both* the prototypes  $\mu_k$  and the assignments  $\mathbf{s}_n$  for each point, so as to minimize  $J$ .
- ▶ This seems like a tricky optimization problem, because the  $\mathbf{s}_n$  are discrete and the  $\mu_k$  are continuous.
- ▶ But learning gets easy if we decouple learning the  $\mathbf{s}_n$  from learning the  $\mu_k$ .
- ▶ This gives rise to a *block coordinate-descent method*, which is a special case of the *EM-algorithm* for training mixtures of Gaussians.

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍

Roland Memisevic

Visual feature learning

## K-means clustering

### Finding the optimal $\mu_k$

- ▶ Given  $S$ ,  $J$  is a quadratic function of  $\mu_k$  which we can minimize by setting the derivative to zero:

$$2 \sum_{n=1}^N s_{nk} (\mathbf{x}_n - \mu_k) = 0$$

- ▶ Solving for  $\mu_k$  yields:

$$\mu_k = \frac{\sum_n s_{nk} \mathbf{x}_n}{\sum_n s_{nk}}$$

- ▶ This solution has a simple interpretation: Set each  $\mu_k$  to the mean of all points currently assigned to cluster  $k$  !

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍

Roland Memisevic

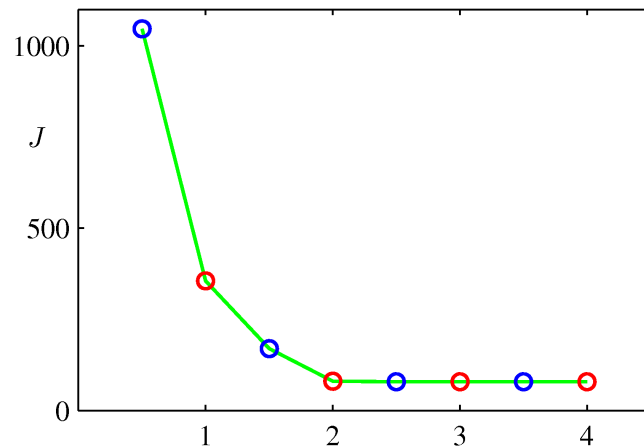
Visual feature learning

## K-means clustering

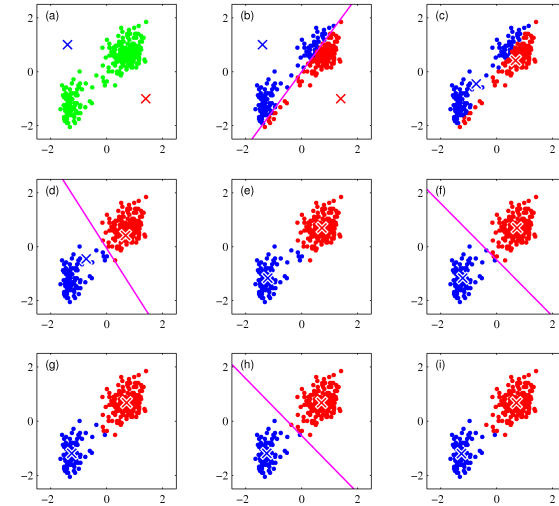
- ▶ Learning amounts to iterating inference of the  $\mathbf{s}_n$ , and adapting the parameters  $\mu_k$  until there are no more changes.
- ▶ This training procedure always converges:  $J$  is positive, and every step either decreases it or leaves it unchanged.
- ▶ But there can be local minima.
- ▶ One way to deal with this is to try multiple runs with different initializations for the parameters  $\mu_k$  and to pick the solution with the lowest final cost.
- ▶ (Bishop, 2006):



## The value of $J$ as learning progresses



## K-means example



- ▶ K-means with  $K = 2$ .



## K-means inference

- ▶ Given the trained model, we can infer the cluster-centers for new test-data points  $\mathbf{x}$  not seen during training: Pick the nearest  $\mu_k$  like we did during training:

$$s_k = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

- ▶ So inference involves some *non-linear* interactions between the hidden.
- ▶ The set of all  $K$  prototypes  $\mu_k$  is sometimes called *codebook*.
- ▶ Clustering and  $K$ -means are also known as *vector quantization*.



## K-means via online learning

- ▶ The reconstruction error for training point  $\mathbf{x}$  may be written

$$E(W) = \frac{1}{2}(\mathbf{x} - \mathbf{w}_{s(\mathbf{x})})^2$$

- ▶ Its gradient is

$$\frac{\partial E(W)}{\partial \mathbf{w}_i} = -(\mathbf{x} - \mathbf{w}_{s(\mathbf{x})})\delta_{s(\mathbf{x}),i}$$

- ▶ So we can use the online learning rule:

$$\mathbf{w}_{s(\mathbf{x})} \leftarrow \mathbf{w}_{s(\mathbf{x})} + \eta(\mathbf{x} - \mathbf{w}_{s(\mathbf{x})})$$



## Online K-means and Hebbian learning

- ▶ We can interpret the online k-means updates as:

### Hebb-rule + competition + unlearning

- ▶ To this end write the update as

$$\Delta \mathbf{w}_k = \eta \delta_{ks(\mathbf{x})}(\mathbf{x} - \mathbf{w}_k)$$

where

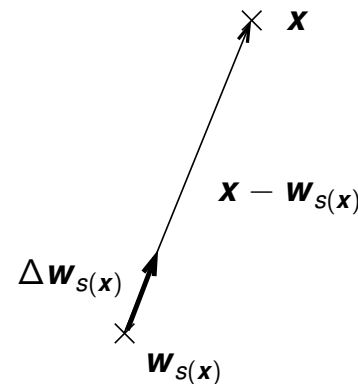
$$\delta_{ks(\mathbf{x})} = \begin{cases} 1 & \text{if } s(\mathbf{x}) = k \\ 0 & \text{else} \end{cases}$$

is the “post-synaptic activity” determined by competition (“winner takes all” rule)

- ▶ There are two learning terms:



## Geometry of online K-means



- ▶  $\Delta \mathbf{w}_{s(\mathbf{x})} = \eta(\mathbf{x} - \mathbf{w}_{s(\mathbf{x})})$  moves the winning weight vector towards the observation.



## K-means and Hebbian learning

1. A Hebbian term:

$$\delta_{ks(\mathbf{x})}\mathbf{x}$$

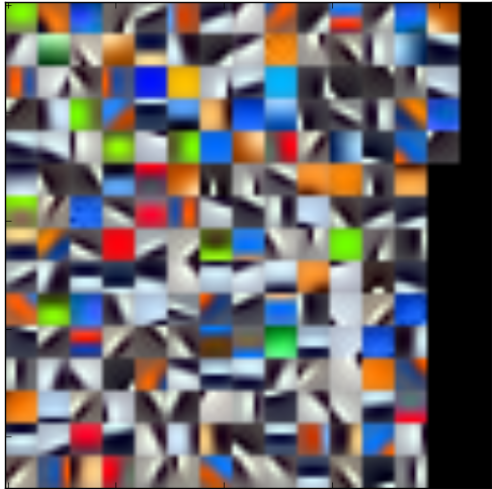
2. An “unlearning” term:

$$-\delta_{ks(\mathbf{x})}\mathbf{w}_k$$

- ▶ The positive term decreases the energy near the data.
- ▶ The unlearning term increases the energy everywhere.
- ▶ “Hebb-rule + competition + unlearning” are present (not surprisingly) in a wide variety of learning algorithms, including contrastive divergence learning for RBMs.



# K-means features from natural image patches



Navigation icons: back, forward, search, and other controls.