

# Machine learning for vision

Fall 2015

Roland Memisevic

Lecture 8, February 24, 2015

Navigation icons

Roland Memisevic

Machine learning for vision

## Classic $K$ -means clustering

- ▶ Define  $\mathbf{s}_n$  as the one-hot encoding of the discrete variable representing the index of the nearest cluster center for  $\mathbf{x}_n$ .
- ▶ It is also useful to think of a matrix  $S$  with entries  $s_{nk}$ , holding the one-hot vectors in its rows.
- ▶ 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*:

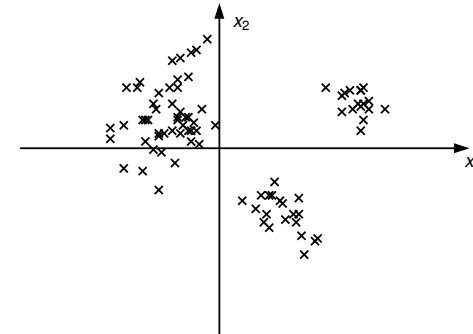
$$J = \sum_{n=1}^N \sum_{k=1}^K s_{nk} \|\mathbf{x}_n - \mathbf{w}_k\|^2$$

Navigation icons

Roland Memisevic

Machine learning for vision

## Review of classic (GOF) $K$ -means clustering



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

Navigation icons

Roland Memisevic

Machine learning for vision

## Classic $K$ -means clustering

- ▶ *Learning* amounts to finding *both* the prototypes  $\mathbf{w}_k$  *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  $\mathbf{w}_k$  are continuous.
- ▶ But learning gets easy (*on paper*) if we decouple learning the  $\mathbf{s}_n$  from learning the  $\mathbf{w}_k$ .
- ▶ In practice, it may be much better not to do this but to train online instead, as we will see later.

Navigation icons

Roland Memisevic

Machine learning for vision

## Classic $K$ -means clustering

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

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

- ▶ Given the  $\mathbf{w}_k$ , 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  $\mathbf{w}_k$  is *closest*.
- ▶ Formally, to optimize all  $\mathbf{s}_n$ , given the set of  $\mathbf{w}_k$ , set:

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

Roland Memisevic

Machine learning for vision

## Classic $K$ -means clustering

- ▶ Learning amounts to iterating inference for the  $\mathbf{s}_n$ , and adapting the parameters  $\mathbf{w}_k$  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  $\mathbf{w}_k$  and to pick the solution with the lowest final cost.

Roland Memisevic

Machine learning for vision

## Classic $K$ -means clustering

### Finding the optimal $\mathbf{w}_k$

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

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

- ▶ Solving for  $\mathbf{w}_k$  yields:

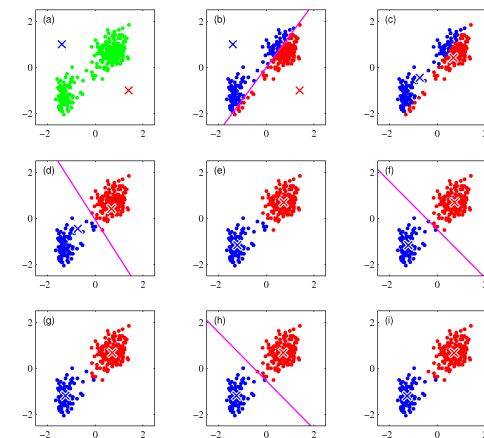
$$\mathbf{w}_k = \frac{\sum_n s_{nk} \mathbf{x}_n}{\sum_n s_{nk}}$$

- ▶ This solution has a simple interpretation: Set  $\mathbf{w}_k$  to the mean of all points currently assigned to cluster  $k$ .

Roland Memisevic

Machine learning for vision

## Classic $K$ -means example ( $K = 2$ )

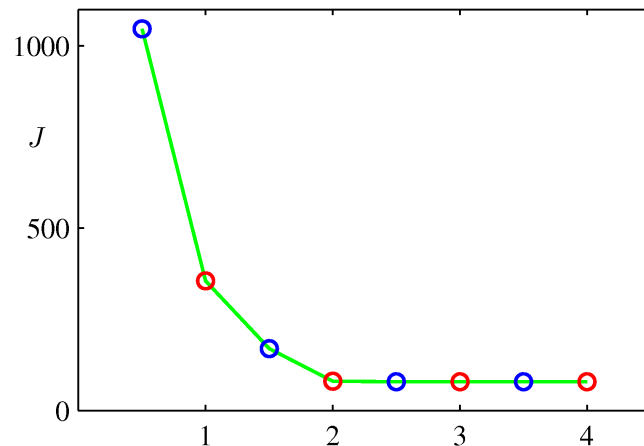


this and most of the following images from: (Bishop, 2006)

Roland Memisevic

Machine learning for vision

## The value of $J$ as learning progresses



## $K$ -means inference

- ▶ Given the trained model, we can infer the cluster-center for a new test-data point  $\mathbf{x}$  not seen during training, by finding the nearest  $\mathbf{w}_k$  like during training:

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

- ▶ The set of all  $K$  prototypes  $\mathbf{w}_k$  is called *codebook*.
- ▶ Clustering and  $K$ -means are also known as *vector quantization*.

## Classic $K$ -means clustering

end of review

## Why is $K$ -means useful?

- ▶ Since  $K$ -means “tiles” space into locally constant regions, an arbitrary non-linear function (up to the tiling resolution,  $K$ ) can be represented with a subsequent linear layer.
- ▶  $K$ -means distributes cluster-centers in space, such that their density is roughly proportional to the data density.
- ▶ So it will resolve high-density regions well, at the cost of low-density regions.

## $K$ -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:

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

- ▶ (Here, it is easier to think of  $s(x)$  as index rather than one-hot vector.)

## 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

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

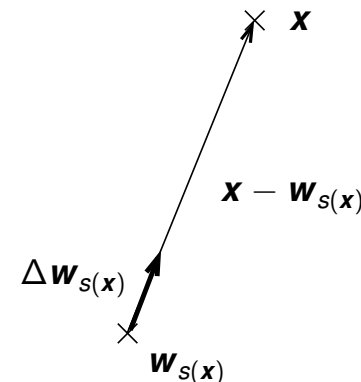
where

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

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.

## Hebbian K-means in 9 lines of code

```
import numpy
def kmeans(W, X, numepochs, learningrate=0.01, batchsize=100):
    X2 = (X**2).sum(1)[:, None]
    for epoch in range(numepochs):
        for i in range(0, X.shape[0], batchsize):
            D = -2*numpy.dot(W, X[i:i+batchsize,:].T) + (W**2).sum(1)[:, None] + X2[i:i+batchsize].T
            S = (D==D.min(0)[None,:]).astype("float").T
            W += learningrate * (numpy.dot(S.T, X[i:i+batchsize,:]) - S.sum(0)[:, None] * W)
    return W
```

## Self-organizing maps

- ▶ We obtain the *self-organizing map* (SOM) aka *Kohonen network* by changing the k-means update from

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

into

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

where  $h_{kj}$  is some smooth neighborhood function, that will let hiddenes near the winning hidden learn, too.

- ▶ This requires hiddenes to be arranged in space in some way (commonly 2-D).

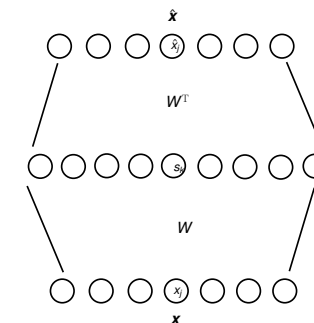
## K-means features learned from natural image patches



## K-means as autoencoder

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

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



- ▶ The cost is  $\sum_i \|\mathbf{x}_i - W_{\text{wta}}(W^T \mathbf{x}_i)\|^2$
- ▶ wta = “winner takes all”
- ▶ Weights are “tied”: Recognition weights are the transpose of generative weights.

## The $K$ -means energy function

- ▶ We can think of  $K$ -means as energy based learning, if we define the energy function as

$$E(\mathbf{x}) = \|\mathbf{x} - W_{\text{wta}}(W^T \mathbf{x})\|^2 = \|\mathbf{x} - \mathbf{w}_{s(\mathbf{x})}\|^2$$

- ▶ Since far from the cluster-centers the energy goes to  $\infty$ ,  $K$ -means has low energy everywhere else.
- ▶ In other words,  $K$ -means doesn't have the *capacity* to produce an arbitrary energy surface with low energy far away from data.

## Winner-takes-all and lateral interactions

- ▶ The winner-takes-all function may be defined as

$$\text{wta}(\mathbf{x}) = \text{onehot}\left(\arg \min_k \|\mathbf{x} - \mathbf{w}_k\|^2\right)$$

where  $\mathbf{w}_k$  is a column of  $W$ .

- ▶ The squared distance can be written as

$$\|\mathbf{x} - \mathbf{w}_k\|^2 = \mathbf{x}^T \mathbf{x} + \mathbf{w}_k^T \mathbf{w}_k - 2\mathbf{w}_k^T \mathbf{x}$$

- ▶ If all  $\mathbf{w}_k$  have the same norm, inference amounts to finding the hidden unit which maximizes

$$\mathbf{w}_k^T \mathbf{x}$$

which is the usual “simple cell” response plus competition.

## The $K$ -means energy function

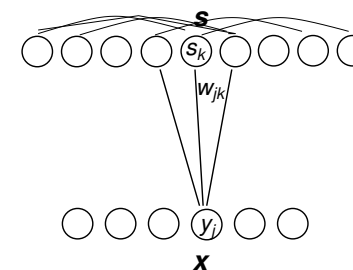
- If we define the un-normalized probability for a point as

$$q(\mathbf{x}_n) = \exp(-E(\mathbf{x}_n))$$

we obtain a density model (which is just the superposition of  $K$  bumps).

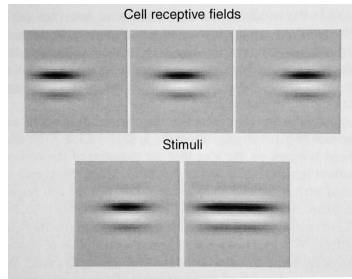
- ▶ So the probability goes to zero far from any cluster center.
- ▶ Unlike RBMs and many probabilistic models, there is no need to lower density away from the data in this model: “Negative updates are *built in*.”
- ▶ (LeCun, 2006)

## Lateral interactions



- ▶ Since  $w_{ta}$  is a function of all the hidden, inference requires the hidden to talk to each other.
- ▶ This is commonly referred to as *lateral interactions*.
- ▶ To induce competition, the interactions need to be *inhibitory*.

## End-stopping



from: Natural Image Statistics (Hyvarinen, Hurri, Hoyer; 2009)

- ▶ For the simple cell in the top middle, a linear model would predict the bottom right stimulus to give at least as large a response as the bottom left stimulus.
- ▶ But for actual neural responses it may give a weaker response.
- ▶ This effect is known as *end-stopping*, and it may be counted as evidence for lateral inhibition.