

# Machine learning for vision

Fall 2013

Roland Memisevic

Lecture 7, October 15, 2013



## Overcomplete codes

- ▶ The determinant is not defined for rectangular  $\mathbf{W}$ .
- ▶ But why is  $|\det \mathbf{W}|$  there anyway?
- ▶ Because it is the **normalizing constant** that allows us to express densities over  $\mathbf{x}$  using the linearly transformed  $\mathbf{W}^T \mathbf{x}$ .

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{1}{Z(\mathbf{W})} p_s(\mathbf{W}^T \mathbf{x})$$

$$Z(\mathbf{W}) = \int_{\mathbf{x}} p_s(\mathbf{W}^T \mathbf{x}) d\mathbf{x}$$

- ▶ This suggests dealing with overcomplete  $\mathbf{W}$  by trying to optimizing  $Z(\mathbf{W})$  directly.
- ▶ Unfortunately, it will be hard to compute in general.
- ▶  $Z(\mathbf{W})$  is also called *partition function*.



## Overcomplete codes

- ▶ V1 increases the dimensionality of the representation as compared to the retinal representation (by a factor of 25 or so).
- ▶ (If sparseness is a good thing, then this is probably not surprising.)
- ▶ If we want to turn the ICA model into such an overcomplete model we will need to make  $\mathbf{W}$  rectangular.
- ▶ Recall the definition of the ICA log likelihood:

$$\log L(\mathbf{w}_1, \dots, \mathbf{w}_n) = T \log |\det \mathbf{W}| + \sum_{i=1}^n \sum_{t=1}^T \log p_i(\mathbf{w}_i^T \mathbf{x}_t)$$



## Overcomplete codes

- ▶ This also allows us to define the density of images in a more general way.
- ▶ A common general formulation, retaining independence of the hidden, is:

$$p(\mathbf{x}) = \frac{1}{Z(\mathbf{W})} \exp \left( \sum_{i=1}^n G_i(\mathbf{w}_i^T \mathbf{x}) \right)$$

$$Z(\mathbf{W}) = \int_{\mathbf{x}} \prod_{i=1}^n \exp (G_i(\mathbf{w}_i^T \mathbf{x})) d\mathbf{x}$$

where  $G_i$  are any non-linear functions.



## Markov Random Fields

- ▶ We may now define the model convolutionally, by scanning small filters across a larger image.
- ▶ Define fine filters to have the shape of image patches, and write:

$$\begin{aligned} & \log p(W_1, \dots, W_n) \\ &= \sum_{x,y} \sum_{i=1}^n G\left(\sum_{\xi,\eta} w_{\xi,\eta} I(x+\xi, y+\eta)\right) - \log Z(W_1, \dots, W_n) \end{aligned}$$

- ▶ (eg. Roth & Black, 2005)

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍 🔄

## Overcomplete codes

- ▶ Good news I: In practice we often care only about computing *features*  $\mathbf{w}_i^T \mathbf{x}$  for test-data.
- ▶ Good news II: We can still *compare* the probabilities between points  $\mathbf{x}_i, \mathbf{x}_j$ , because  $Z(\mathbf{W})$  does not depend on  $\mathbf{x}$ .
- ▶ Good news III: We can do *approximate* maximum likelihood training, which often works just as well in practice.

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍 🔄

## Overcomplete codes

- ▶ Bad news I: In general, we cannot compute the log-likelihood, nor its derivative.
- ▶ Bad news II: In the absence of a computable partition function, we cannot evaluate probabilities for test-data, either.

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍 🔄

## Energy based models

- ▶ To get insights into the learning problem, it is useful to rewrite the probability of data in an even more general form:

$$p(\mathbf{x}; \mathbf{W}) = \frac{1}{Z(\mathbf{W})} q(\mathbf{x}; \mathbf{W})$$

with log-probability

$$\log p(\mathbf{x}; \mathbf{W}) = \log q(\mathbf{x}; \mathbf{W}) - \log \int_{\mathbf{x}} q(\mathbf{x}; \mathbf{W}) d\mathbf{x}$$

- ▶ We can think of  $q(\mathbf{x})$  as an unnormalized (“pre-”) probability.
- ▶ Usually,  $q(\mathbf{x}) = \prod_{i=1}^n \exp(G_i(\mathbf{w}_i^T \mathbf{x}))$

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍 🔄

## Energy based models

- ▶ For a set of IID points,  $\mathbf{x}_i$ , the log-likelihood and its derivative take the form:

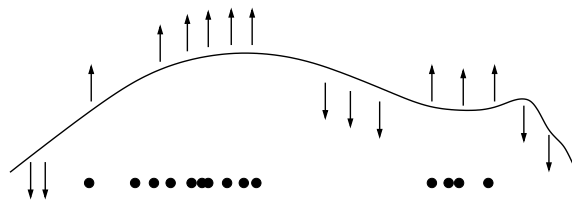
$$\begin{aligned}L(\mathbf{W}) &= \sum_i \log p(\mathbf{x}_i; \mathbf{W}) \\ &= \sum_i \log q(\mathbf{x}_i; \mathbf{W}) - N \log \int_{\mathbf{x}} q(\mathbf{x}; \mathbf{W}) d\mathbf{x}\end{aligned}$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} = \sum_i \frac{\partial \log q(\mathbf{x}_i; \mathbf{W})}{\partial \mathbf{W}} - \frac{N}{Z(\mathbf{W})} \int_{\mathbf{x}} \frac{\partial q(\mathbf{x}; \mathbf{W})}{\partial \mathbf{W}} d\mathbf{x}$$

- ▶ Intuitively, what does maximizing these two terms do with the model?



## Energy based models



- ▶ Decreasing  $q(\mathbf{x}; \mathbf{W})$ , and thereby  $p(\mathbf{x}; \mathbf{W})$ , ensures that we get a normalized probability distribution.



## Energy based models

- ▶ For a set of IID points,  $\mathbf{x}_i$ , the log-likelihood and its derivative take the form:

$$\begin{aligned}L(\mathbf{W}) &= \sum_i \log p(\mathbf{x}_i; \mathbf{W}) \\ &= \sum_i \log q(\mathbf{x}_i; \mathbf{W}) - N \log \int_{\mathbf{x}} q(\mathbf{x}; \mathbf{W}) d\mathbf{x}\end{aligned}$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} = \sum_i \frac{\partial \log q(\mathbf{x}_i; \mathbf{W})}{\partial \mathbf{W}} - \frac{N}{Z(\mathbf{W})} \int_{\mathbf{x}} \frac{\partial q(\mathbf{x}; \mathbf{W})}{\partial \mathbf{W}} d\mathbf{x}$$

- ▶ Intuitively, what does maximizing these two terms do with the model?

Maximum likelihood learning *increases*  $q(\mathbf{x}; \mathbf{W})$  at the data points, and it *decreases*  $q(\mathbf{x}; \mathbf{W})$  everywhere.



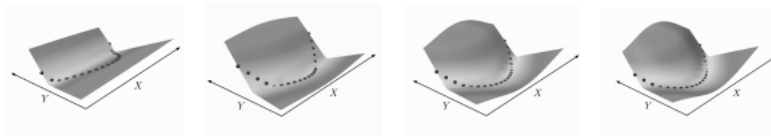
## Energy based models

- ▶ This view is even more liberating than replacing the determinant with  $Z(\mathbf{W})$ :
- ▶ In high-dimensional spaces, it can be easy to push down probabilities near the data and to ignore what's happening far away.
- ▶ This will be OK, if all we will ever see are points from high-density regions.
- ▶ Technically, the pushing-down can be done by sampling from the model distribution. This is inefficient in high-dimensional spaces. Solution: Sample only near the data.
- ▶ “Contrastive divergence” (Hinton, 2002)



## 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):



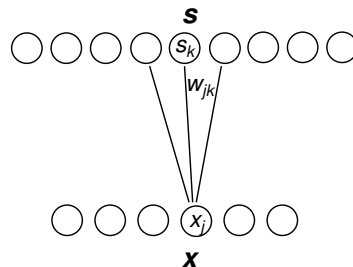
## Energy based models

- ▶ 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 are connected to pixels in a bi-partite graph.
- ▶ In other words, most feature learning models are based on a variation of the linear encoder/decoder equations.
- ▶ Pushing up the energy away from the data typically translates into a *capacity constraint* on the hidden.



## Feature learning and bi-partite networks

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



- ▶ PCA is a special case with linear dependencies and low-dimensional  $\mathbf{s}$
- ▶ ICA is a special case with linear dependencies and sparse  $\mathbf{s}$

