

Machine learning for vision

Winter 2015

Roland Memisevic

Lecture 6, February 17, 2015



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 optimize $Z(\mathbf{W})$ directly.
- ▶ Unfortunately, it will be hard to compute in general.
- ▶ $Z(\mathbf{W})$ is also called *partition function*.



Overcomplete codes

- ▶ In area V1, the dimensionality of the retinal representation gets increased (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 general formulation, retaining independence of the hiddens, is

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

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

where G is a non-linear function.



Markov Random Fields

- ▶ A separate advantage is that we may now define the model convolutionally, by scanning small filters across a larger image:

$$\begin{aligned} & \log p(\mathbf{x}) \\ &= \sum_{x,y} \sum_i G\left(\sum_{\xi,\eta} \mathbf{w}_{\xi,\eta}^i I(x+\xi, y+\eta)\right) - \log Z \end{aligned}$$

- ▶ This is known as field-of-experts (Roth & Black, 2005).

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍

Roland Memisevic

Machine learning for vision

Energy based models

- ▶ To gain intuitive insights into the learning problem, it can be useful to rewrite the probability of data in an even more general form:

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

with the log-likelihood

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

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

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍

Roland Memisevic

Machine learning for vision

Overcomplete codes

- ▶ Bad news: For most models, we cannot compute the (log-)likelihood, nor its derivative.
- ▶ Good news (I): We can still *compare* the probabilities between observations $\mathbf{x}_i, \mathbf{x}_j$, because $Z(\mathbf{W})$ does not depend on \mathbf{x} .
- ▶ Good news (II): We can do *approximate* maximum likelihood training, which often seems to work just as well in practice.
- ▶ Good news (III): In practice, we usually need the *features* $\mathbf{w}_i^T \mathbf{x}$ for test-data not probabilities.

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍

Roland Memisevic

Machine learning for vision

Energy based models

- ▶ For a set of IID points, \mathbf{x}_i , the log-likelihood and its derivative can be written

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

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

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

◀ ▶ ⏪ ⏩ ⏴ ⏵ ⏶ ⏷ ⏸ ⏹ ⏺ ⏻ ⏼ ⏽ ⏾ ⏿ 🔍

Roland Memisevic

Machine learning for vision

Energy based models

- ▶ For a set of IID points, \mathbf{x}_i , the log-likelihood and its derivative can be written

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

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

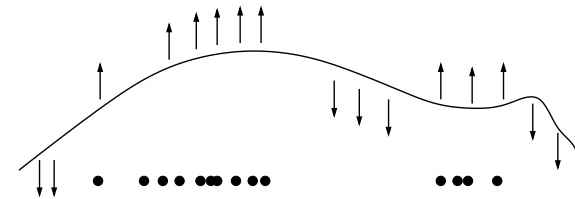
- ▶ Intuitively, what does maximizing these two terms do with the model?

Maximum likelihood learning will *increase* $q(\mathbf{x}; \mathbf{W})$ at the data points, and *decrease* $q(\mathbf{x}; \mathbf{W})$ everywhere.

Energy based models, contrastive divergence

- ▶ This view is even more liberating than replacing the determinant with $Z(\mathbf{W})$:
- ▶ In high-dimensional spaces, pushing down probabilities *near* the data may be good enough if pushing down everywhere is too expensive.
- ▶ This will be OK, if all we will ever see are test-cases from high-density regions.
- ▶ Technically, the pushing-down terms can also be derived by sampling from the model distribution, and restricting the sampling region to the vicinity of training examples is an approximation known as “contrastive divergence” (Hinton, 2002)

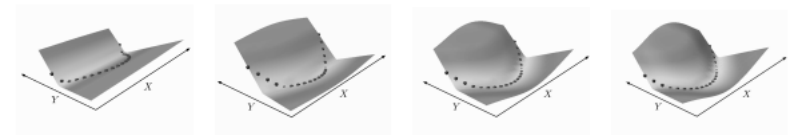
Energy based models



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

Energy based models, non-probabilistic

- ▶ 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)



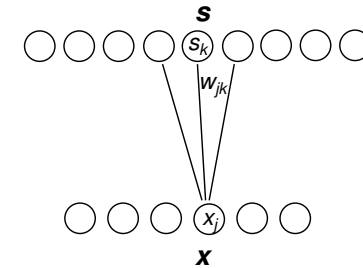
Energy based models, comments

- ▶ In practice, it is common to define energies as $-q(\mathbf{x}; W)$, in which case we want to *minimize* energy near the data.
- ▶ For feature learning energy-based models practically always contain hidden variables that 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.
- ▶ The effect of pushing down (or up) the energy away from the data usually corresponds to (and is the result of) a *capacity constraint* on the hiddens.



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}
- ▶ Other special cases: restricted Boltzmann machines, autoencoder networks, sparse coding, k-means clustering

