

Chapter 15

Linear Factor Models and Auto-Encoders

Linear factor models are generative unsupervised learning models in which we imagine that some unobserved factors \mathbf{h} explain the observed variables \mathbf{x} through a linear transformation. Auto-encoders are unsupervised learning methods that learn a representation of the data, typically obtained by a non-linear parametric transformation of the data, i.e., from \mathbf{x} to \mathbf{h} , typically a feedforward neural network, but not necessarily. They also learn a transformation going backwards from the representation to the data, from \mathbf{h} to \mathbf{x} , like the linear factor models. Linear factor models therefore only specify a parametric decoder, whereas auto-encoder also specify a parametric encoder. Some linear factor models, like PCA, actually correspond to an auto-encoder (a linear one), but for others the encoder is implicitly defined via an inference mechanism that searches for an \mathbf{h} that could have generated the observed \mathbf{x} .

The idea of auto-encoders has been part of the historical landscape of neural networks for decades (LeCun, 1987; Bourlard and Kamp, 1988; Hinton and Zemel, 1994) but has really picked up speed in recent years. They remained somewhat marginal for many years, in part due to what was an incomplete understanding of the mathematical interpretation and geometrical underpinnings of auto-encoders, which are developed further in Chapters 17 and 20.11.

An auto-encoder is simply a neural network that tries to copy its input to its output. The architecture of an auto-encoder is typically decomposed into the following parts, illustrated in Figure 15.1: KEY
IDEA

- an input, \mathbf{x}
- an encoder function f
- a “code” or internal representation $\mathbf{h} = f(\mathbf{x})$

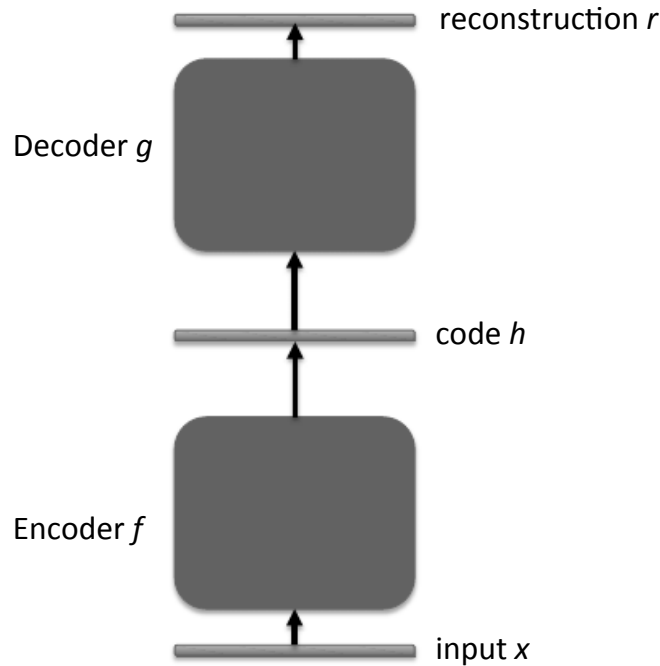


Figure 15.1: General schema of an auto-encoder, mapping an input \mathbf{x} to an output (called reconstruction) \mathbf{r} through an internal representation or code \mathbf{h} . The auto-encoder has two components: the encoder f (mapping \mathbf{x} to \mathbf{h}) and the decoder g (mapping \mathbf{h} to \mathbf{r}).

- a decoder function g
- an output, also called “reconstruction” $\mathbf{r} = g(\mathbf{h}) = g(f(\mathbf{x}))$
- a loss function L computing a scalar $L(\mathbf{r}, \mathbf{x})$ measuring how good of a reconstruction \mathbf{r} is of the given input \mathbf{x} . The objective is to minimize the expected value of L over the training set of examples $\{\mathbf{x}\}$.

15.1 Regularized Auto-Encoders

Predicting the input may sound useless: what could prevent the auto-encoder from simply copying its input into its output? In the 20th century, this was achieved by constraining the architecture of the auto-encoder to avoid this, by forcing the dimension of the code \mathbf{h} to be smaller than the dimension of the input \mathbf{x} .

Figure 15.2 illustrates the two typical cases of auto-encoders: undercomplete vs overcomplete, i.e., with the dimension of the representation \mathbf{h} respectively smaller vs larger than the input \mathbf{x} . Whereas early work with auto-encoders, just like PCA, uses the undercompleteness – i.e. a bottleneck in the sequence of layers – to avoid learning the identity function, more recent work allows overcomplete

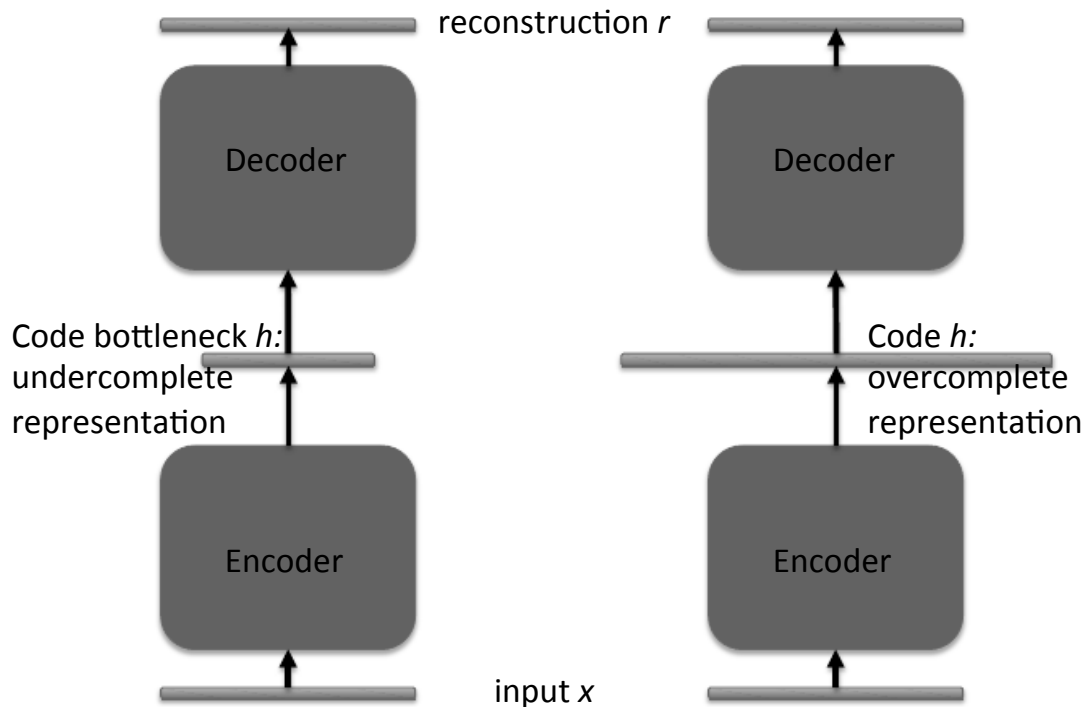


Figure 15.2: Left: undercomplete representation (dimension of code h is less than dimension of input x). Right: overcomplete representation. Overcomplete auto-encoders require some other form of regularization (instead of the constraint on the dimension of h) to avoid the trivial solution where $r = x$ for all x .

representations. What we have learned in recent years is that it is possible to make the auto-encoder meaningfully capture the structure of the input distribution even if the representation is overcomplete, with other forms of constraint or regularization. In fact, once you realize that auto-encoders can capture the input distribution (indirectly, not as an explicit probability function), you also realize that it should need more capacity as one increases the complexity of the distribution to be captured (and the amount of data available): it should not be limited by the input dimension. This is a problem in particular with the shallow auto-encoders, which have a single hidden layer (for the code). Indeed, that hidden layer size controls both the dimensionality reduction constraint (the code size at the bottleneck) and the capacity (which allows to learn a more complex distribution).

Besides the **bottleneck** constraint, alternative constraints or regularization methods have been explored and can guarantee that the auto-encoder does something useful and not just learn some trivial identity-like function:

- **Sparsity of the representation or of its derivative:** even if the intermediate representation has a very high dimensionality, the effective local dimensionality (number of degrees of freedom that capture a coordinate sys-

tem among the probable x 's) could be much smaller if most of the elements of h are zero (or any other constant, such that $\|\frac{\partial h_i}{\partial \mathbf{x}}\|$ is close to zero). When $\|\frac{\partial h_i}{\partial \mathbf{x}}\|$ is close to zero, h_i does not participate in encoding local changes in \mathbf{x} . There is a geometrical interpretation of this situation in terms of *manifold learning* that is discussed in more depth in Chapter 17. The discussion in Chapter 16 also explains how an auto-encoder naturally tends towards learning a coordinate system for the actual factors of variation in the data. At least four types of “auto-encoders” clearly fall in this category of sparse representation:

- **Sparse coding** (Olshausen and Field, 1996) has been heavily studied as an unsupervised feature learning and feature inference mechanism. It is a linear factor model rather than an auto-encoder, because it has no explicit parametric encoder, and instead uses an iterative optimization procedure to compute the maximally likely code. Sparse coding looks for representations that are both sparse and explain the input through the decoder. Instead of the code being a parametric function of the input, it is considered like free variable that is obtained through an optimization, i.e., a particular form of inference:

$$\mathbf{h}^* = f(\mathbf{x}) = \arg \min_{\mathbf{h}} L(g(\mathbf{h}), \mathbf{x}) + \lambda \Omega(\mathbf{h}) \quad (15.1)$$

where L is the reconstruction loss, f the (non-parametric) encoder, g the (parametric) decoder, $\Omega(\mathbf{h})$ is a sparsity regularizer, and in practice the minimization can be approximate. Sparse coding has a manifold or geometric interpretation that is discussed in Section 15.8. It also has an interpretation as a directed graphical model, described in more details in Section 19.3. To achieve sparsity, the objective function to optimize includes a term that is minimized when the representation has many zero or near-zero values, such as the L1 penalty $|\mathbf{h}|_1 = \sum_i |h_i|$.

- An interesting variation of sparse coding combines the freedom to choose the representation through optimization and a parametric encoder. It is called **predictive sparse decomposition** (PSD) (Kavukcuoglu *et al.*, 2008a) and is briefly described in Section 15.8.2.
- At the other end of the spectrum are simply **sparse auto-encoders**, which combine with the standard auto-encoder schema a sparsity penalty which encourages the output of the encoder to be sparse. These are described in Section 15.8.1. Besides the L1 penalty, other sparsity penalties that have been explored include the Student-t penalty (Olshausen and Field, 1996; Bergstra, 2011), TODO: should the t be in

math mode, perhaps?

$$\sum_i \log(1 + \alpha^2 h_i^2)$$

(i.e. where αh_i has a Student-t prior density) and the KL-divergence penalty (Lee *et al.*, 2008; Goodfellow *et al.*, 2009; Larochelle and Bengio, 2008a)

$$- \sum_i (t \log h_i + (1 - t) \log(1 - h_i)),$$

with a target sparsity level t , for $h_i \in (0, 1)$, e.g. through a sigmoid non-linearity.

- **Contractive autoencoders** (Rifai *et al.*, 2011b), covered in Section 15.10, explicitly penalize $\|\frac{\partial \mathbf{h}}{\partial \mathbf{x}}\|_F^2$, i.e., the sum of the squared norm of the vectors $\frac{\partial h_i(\mathbf{x})}{\partial \mathbf{x}}$ (each indicating how much each hidden unit h_i responds to changes in \mathbf{x} and what direction of change in \mathbf{x} that unit is most sensitive to, around a particular \mathbf{x}). With such a regularization penalty, the auto-encoder is called **contractive**¹ because the mapping from input \mathbf{x} to representation \mathbf{h} is encouraged to be contractive, i.e., to have small derivatives in all directions. Note that a sparsity regularization indirectly leads to a contractive mapping as well, when the non-linearity used happens to have a zero derivative at $h_i = 0$ (which is the case for the sigmoid non-linearity).
- **Robustness to injected noise or missing information:** if noise is injected in inputs or hidden units, or if some inputs are missing, while the neural network is asked to *reconstruct the clean and complete input*, then it cannot simply learn the identity function. It has to capture the structure of the data distribution in order to optimally perform this reconstruction. Such auto-encoders are called *denoising auto-encoders* and are discussed in more detail in Section 15.9.

15.2 Denoising Auto-encoders

There is a tight connection between the denoising auto-encoders and the contractive auto-encoders: it can be shown (Alain and Bengio, 2013) that in the limit of small Gaussian injected input noise, the denoising reconstruction error is equivalent to a contractive penalty on the reconstruction function that maps \mathbf{x} to $\mathbf{r} = g(f(\mathbf{x}))$. In other words, since both \mathbf{x} and

¹ A function $f(\mathbf{x})$ is contractive if $\|f(\mathbf{x}) - f(\mathbf{y})\| < \|\mathbf{x} - \mathbf{y}\|$ for nearby \mathbf{x} and \mathbf{y} , or equivalently if its derivative $\|f'(\mathbf{x})\| < 1$.

$\mathbf{x} + \epsilon$ (where ϵ is some small noise vector) must yield the same target output \mathbf{x} , the reconstruction function is encouraged to be insensitive to changes in all directions ϵ . The only thing that prevents reconstruction \mathbf{r} from simply being a constant (completely insensitive to the input \mathbf{x}), is that one also has to reconstruct correctly for different training examples \mathbf{x} . However, the auto-encoder can learn to be approximately constant around training examples \mathbf{x} while producing a different answer for different training examples. As discussed in Section 17.4, if the examples are near a low-dimensional manifold, this encourages the representation to vary only on the manifold and be locally constant in directions orthogonal to the manifold, i.e., the representation locally captures a (not necessarily Euclidean, not necessarily orthogonal) coordinate system for the manifold. In addition to the denoising auto-encoder, the *variational auto-encoder* (Section 20.9.3) and the *generative stochastic networks* (Section 20.11) also involve the injection of noise, but typically in the representation-space itself, thus introducing the notion of \mathbf{h} as a *latent variable*.

- **Pressure of a Prior on the Representation:** an interesting way to generalize the notion of regularization applied to the representation is to introduce in the cost function for the auto-encoder a log-prior term

$$-\log P(\mathbf{h})$$

which captures the assumption that we would like to find a representation that has a simple distribution (if $P(\mathbf{h})$ has a simple form, such as a factorized distribution²), or at least one that is simpler than the original data distribution. Among all the encoding functions f , we would like to pick one that

1. can be inverted (easily), and this is achieved by minimizing some reconstruction loss, and
2. yields representations \mathbf{h} whose distribution is “simpler”, i.e., can be captured with less capacity than the original training distribution itself.

The sparse variants described above clearly fall in that framework. The variational auto-encoder (Section 20.9.3) provides a clean mathematical framework for justifying the above pressure of a top-level prior when the objective is to model the data generating distribution.

From the point of view of regularization (Chapter 7), adding the $-\log P(\mathbf{h})$ term to the objective function (e.g. for encouraging sparsity) or adding a contractive penalty do not fit the traditional view of a prior on the parameters. Instead,

²all the sparse priors we have described correspond to a factorized distribution

the prior on the latent variables acts like a *data-dependent prior*, in the sense that it depends on the particular values \mathbf{h} that are going to be sampled (usually from a posterior or an encoder), based on the input example \mathbf{x} . Of course, indirectly, this is also a regularization on the parameters, but one that depends on the particular data distribution.

15.3 Representational Power, Layer Size and Depth

Nothing in the above description of auto-encoders restricts the encoder or decoder to be shallow, but in the literature on the subject, most trained auto-encoders have had a single hidden layer which is also the representation layer or code³

For one, we know by the usual universal approximator abilities of single hidden-layer neural networks that a sufficiently large hidden layer can represent any function with a given accuracy. This observation justifies overcomplete auto-encoders: in order to represent a rich enough distribution, one probably needs many hidden units in the intermediate representation layer. We also know that Principal Components Analysis (PCA) corresponds to an undercomplete auto-encoder with no intermediate non-linearity, and that PCA can only capture a set of directions of variation that are the same everywhere in space. This notion is discussed in more details in Chapter 17 in the context of manifold learning.

For two, it has also been reported many times that training a deep neural network, and in particular a deep auto-encoder (i.e. with a deep encoder and a deep decoder) is more difficult than training a shallow one. This was actually a motivation for the initial work on the *greedy layerwise unsupervised pre-training procedure*, described below in Section 16.1, by which we only need to train a series of shallow auto-encoders in order to initialize a deep auto-encoder. It was shown early on (Hinton and Salakhutdinov, 2006) that, if trained properly, such deep auto-encoders could yield much better compression than corresponding shallow or linear auto-encoders (which are basically doing the same as PCA, see Section 15.6 below). As discussed in Section 16.7, deeper architectures can be in some cases exponentially more efficient (both in terms of computation and statistically) than shallow ones. However, because we can usefully pre-train a deep net by training and stacking shallow ones, it makes it interesting to consider single-layer (or at least shallow and easy to train) auto-encoders, as has been done in most of the literature discussed in this chapter.

³as argued in this book, this is probably not a good choice, and we would like to independently control the constraints on the representation, e.g. dimension and sparsity of the code, and the capacity of the encoder.

15.4 Reconstruction Distribution

The above “parts” (encoder function f , decoder function g , reconstruction loss L) make sense when the loss L is simply the squared reconstruction error, but there are many cases where this is not appropriate, e.g., when \mathbf{x} is a vector of discrete variables or when $P(\mathbf{x} | \mathbf{h})$ is not well approximated by a Gaussian distribution⁴. Just like in the case of other types of neural networks (starting with the feedforward neural networks, Section 6.3.2), it is convenient to define the loss L as a negative log-likelihood over some target random variables. This probabilistic interpretation is particularly important for the discussion in Sections 20.9.3, 20.10 and 20.11 about generative extensions of auto-encoders and stochastic recurrent networks, where the output of the auto-encoder is interpreted as a probability distribution $P(\mathbf{x} | \mathbf{h})$, for reconstructing \mathbf{x} , given hidden units \mathbf{h} . This distribution captures not just the expected reconstruction but also the *uncertainty* about the original \mathbf{x} (which gave rise to \mathbf{h} , either deterministically or stochastically, given \mathbf{h}). In the simplest and most ordinary cases, this distribution factorizes, i.e., $P(\mathbf{x} | \mathbf{h}) = \prod_i P(x_i | \mathbf{h})$. This covers the usual cases of $x_i | \mathbf{h}$ being Gaussian (for unbounded real values) and $x_i | \mathbf{h}$ having a Bernoulli distribution (for binary values x_i), but one can readily generalize this to other distributions, such as mixtures (see Sections 3.10.6 and 6.3.2).

Thus we can generalize the notion of *decoding function* $g(\mathbf{h})$ to *decoding distribution* $P(\mathbf{x} | \mathbf{h})$. Similarly, we can generalize the notion of *encoding function* $f(\mathbf{x})$ to *encoding distribution* $Q(\mathbf{h} | \mathbf{x})$, as illustrated in Figure 15.3. We use this to capture the fact that noise is injected at the level of the representation \mathbf{h} , now considered like a latent variable. This generalization is crucial in the development of the variational auto-encoder (Section 20.9.3) and the generalized stochastic networks (Section 20.11).

We also find a stochastic encoder and a stochastic decoder in the RBM, described in Section 20.2. In that case, the encoding distribution $Q(\mathbf{h} | \mathbf{x})$ and $P(\mathbf{x} | \mathbf{h})$ “match”, in the sense that $Q(\mathbf{h} | \mathbf{x}) = P(\mathbf{h} | \mathbf{x})$, i.e., there is a unique joint distribution which has both $Q(\mathbf{h} | \mathbf{x})$ and $P(\mathbf{x} | \mathbf{h})$ as conditionals. This is not true in general for two independently parametrized conditionals like $Q(\mathbf{h} | \mathbf{x})$ and $P(\mathbf{x} | \mathbf{h})$, although the work on generative stochastic networks (Alain *et al.*, 2015) shows that learning will tend to make them compatible asymptotically (with enough capacity and examples).

⁴See the link between squared error and normal density in Sections 5.6 and 6.3.2

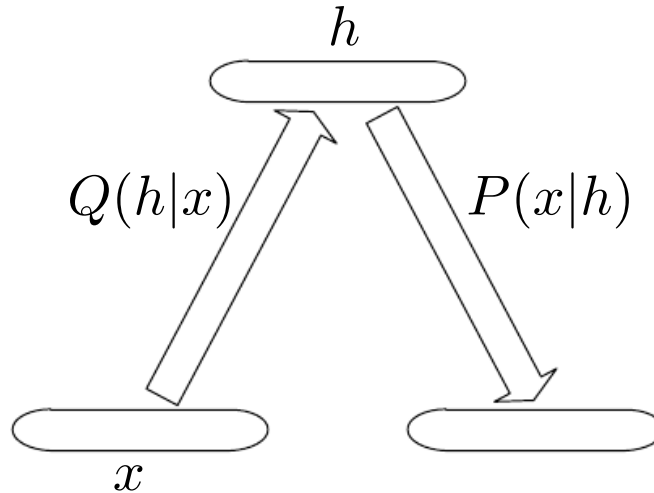


Figure 15.3: Basic scheme of a stochastic auto-encoder, in which both the encoder and the decoder are not simple functions but instead involve some noise injection, meaning that their output can be seen as sampled from a distribution, $Q(\mathbf{h} | \mathbf{x})$ for the encoder and $P(\mathbf{x} | \mathbf{h})$ for the decoder. RBMs are a special case where $P = Q$ (in the sense of a unique joint corresponding to both conditionals) but in general these two distributions are not necessarily conditional distributions compatible with a unique joint distribution $P(\mathbf{x}, \mathbf{h})$.

15.5 Linear Factor Models

Now that we have introduced the notion of a probabilistic decoder, let us focus on a very special case where the latent variable \mathbf{h} generates \mathbf{x} via a linear transformation plus noise, i.e., classical linear factor models, which do not necessarily have a corresponding parametric encoder.

The idea of discovering explanatory factors that have a simple joint distribution among themselves is old, e.g., see Factor Analysis (see below), and has been explored first in the context where the relationship between factors and data is linear, i.e., we assume that the data was generated as follows. First, sample the real-valued factors,

$$\mathbf{h} \sim P(\mathbf{h}), \quad (15.2)$$

and then sample the real-valued observable variables given the factors:

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \text{noise} \quad (15.3)$$

where the noise is typically Gaussian and diagonal (independent across dimensions). This is illustrated in Figure 15.4.

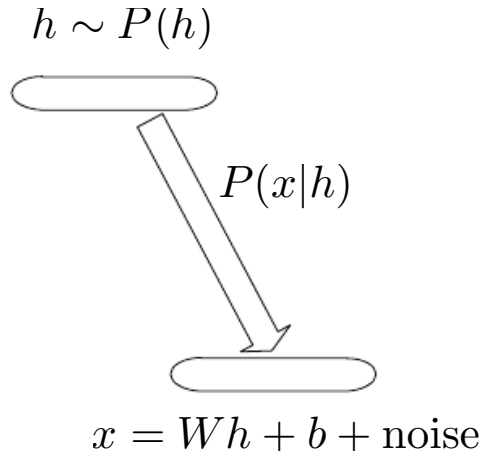


Figure 15.4: Basic scheme of a linear factors model, in which we assume that an observed data vector \mathbf{x} is obtained by a linear combination of latent factors \mathbf{h} , plus some noise. Different models, such as probabilistic PCA, factor analysis or ICA, make different choices about the form of the noise and of the prior $P(\mathbf{h})$.

15.6 Probabilistic PCA and Factor Analysis

Probabilistic PCA (Principal Components Analysis), factor analysis and other linear factor models are special cases of the above equations (15.2 and 15.3) and only differ in the choices made for the prior (over latent, not parameters) and noise distributions.

In factor analysis (Bartholomew, 1987; Basilevsky, 1994), the latent variable prior is just the unit variance Gaussian

$$\mathbf{h} \sim \mathcal{N}(0, \mathbf{I})$$

while the observed variables x_i are assumed to be *conditionally independent*, given \mathbf{h} , i.e., the noise is assumed to be coming from a diagonal covariance Gaussian distribution, with covariance matrix $\boldsymbol{\psi} = \text{diag}(\boldsymbol{\sigma}^2)$, with $\boldsymbol{\sigma}^2 = (\sigma_1^2, \sigma_2^2, \dots)$ a vector of per-variable variances.

The role of the latent variables is thus to *capture the dependencies* between the different observed variables x_i . Indeed, it can easily be shown that \mathbf{x} is just a Gaussian-distribution (multivariate normal) random variable, with

$$\mathbf{x} \sim \mathcal{N}(\mathbf{b}, \mathbf{W}\mathbf{W}^\top + \boldsymbol{\psi})$$

where we see that the weights \mathbf{W} induce a dependency between two variables x_i and x_j through a kind of auto-encoder path, whereby x_i influences $\hat{\mathbf{h}}_k = \mathbf{W}_k \mathbf{x}$ via w_{ki} (for every k) and $\hat{\mathbf{h}}_k$ influences x_j via w_{kj} .

In order to cast PCA in a probabilistic framework, we can make a slight modification to the factor analysis model, making the conditional variances σ_i

equal to each other. In that case the covariance of \mathbf{x} is just $\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}$, where σ^2 is now a scalar, i.e.,

$$\mathbf{x} \sim \mathcal{N}(\mathbf{b}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})$$

or equivalently

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \sigma\mathbf{z}$$

where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$ is white noise. Tipping and Bishop (1999) then show an iterative EM algorithm for estimating the parameters \mathbf{W} and σ^2 .

What the probabilistic PCA model is basically saying is that the covariance is mostly captured by the latent variables \mathbf{h} , up to some small residual *reconstruction error* σ^2 . As shown by Tipping and Bishop (1999), probabilistic PCA becomes PCA as $\sigma \rightarrow 0$. In that case, the conditional expected value of \mathbf{h} given \mathbf{x} becomes an orthogonal projection onto the space spanned by the d columns of \mathbf{W} , like in PCA. See Section 17.1 for a discussion of the “inference” mechanism associated with PCA (probabilistic or not), i.e., recovering the expected value of the latent factors h_i given the observed input \mathbf{x} . That section also explains the very insightful *geometric and manifold interpretation* of PCA.

However, as $\sigma \rightarrow 0$, the density model becomes very sharp around these d dimensions spanned the columns of \mathbf{W} , as discussed in Section 17.1, which would not make it a very faithful model of the data, in general (not just because the data may live on a higher-dimensional manifold, but more importantly because the real data manifold may not be a flat hyperplane - see Chapter 17 for more).

15.6.1 ICA

TODO: do we really want to put every linear factor model in the auto-encoder chapter? if latent variable models are auto-encoders, what deep probabilistic model would not be an auto-encoder? Independent Component Analysis (ICA) is among the oldest representation learning algorithms (Herault and Ans, 1984; Jutten and Herault, 1991; Comon, 1994; Hyvärinen, 1999; Hyvärinen *et al.*, 2001). It is an approach to modeling linear factors that seeks non-Gaussian projections of the data. Like probabilistic PCA and factor analysis, it also fits the linear factor model of Eqs. 15.2 and 15.3. What is particular about ICA is that unlike PCA and factor analysis it *does not assume that the latent variable prior is Gaussian*. It only assumes that it is *factorized*, i.e.,

$$P(\mathbf{h}) = \prod_i P(h_i). \quad (15.4)$$

Since there is no parametric assumption behind the prior, we are really in front of a so-called *semi-parametric model*, with parts of the model being parametric

($P(\mathbf{x} | \mathbf{h})$) and parts being non-specified or non-parametric ($P(\mathbf{h})$). In fact, this

typically yields to *non-Gaussian* priors: if the priors were Gaussian, then one could not distinguish between the factors \mathbf{h} and a rotation of \mathbf{h} . Indeed, note that if

$$\mathbf{h} = \mathbf{U}\mathbf{z}$$

with \mathbf{U} an orthonormal (rotation) square matrix, i.e.,

$$\mathbf{z} = \mathbf{U}^\top \mathbf{h},$$

then, although \mathbf{h} might have a $\text{Normal}(0, \mathbf{I})$ distribution, the \mathbf{z} *also have a unit covariance*, i.e., they are uncorrelated:

$$\text{Var}[\mathbf{z}] = \mathbb{E}[\mathbf{z}\mathbf{z}^\top] = \mathbb{E}[\mathbf{U}^\top \mathbf{h}\mathbf{h}^\top \mathbf{U}] = \mathbf{U}^\top \text{Var}[\mathbf{h}] \mathbf{U} = \mathbf{U}^\top \mathbf{U} = \mathbf{I}.$$

In other words, imposing independence among Gaussian factors does not allow one to disentangle them, and we could as well recover any linear rotation of these factors. It means that, given the observed \mathbf{x} , even though we might assume the right generative model, PCA cannot recover the original generative factors. However, if we assume that the latent variables are *non-Gaussian*, then we can recover them, and this is what ICA is trying to achieve. In fact, under these generative model assumptions, the true underlying factors can be recovered (Comon, 1994). In fact, many ICA algorithms are looking for projections of the data $\mathbf{s} = \mathbf{V}\mathbf{x}$ such that they are *maximally non-Gaussian*. An intuitive explanation for these approaches is that although the true latent variables \mathbf{h} may be non-Gaussian, almost any linear combination of them will look more Gaussian, because of the central limit theorem. Since linear combinations of the x_i 's are also linear combinations of the h_j 's, to recover the h_j 's we just need to find the linear combinations that are maximally non-Gaussian (while keeping these different projections orthogonal to each other).

There is an interesting connection between ICA and sparsity, since the dominant form of non-Gaussianity in real data is due to sparsity, i.e., concentration of probability at or near 0. Non-Gaussian distributions typically have more mass around zero, although you can also get non-Gaussianity by increasing skewness, asymmetry, or kurtosis.

Like PCA can be generalized to non-linear auto-encoders described later in this chapter, ICA can be generalized to a non-linear generative model, e.g., $\mathbf{x} = f(\mathbf{h}) + \text{noise}$. See Hyvärinen and Pajunen (1999) for the initial work on non-linear ICA and its successful use with ensemble learning by Roberts and Everson (2001); Lappalainen *et al.* (2000).

15.6.2 Sparse Coding as a Generative Model

One particularly interesting form of non-Gaussianity arises with distributions that are sparse. These typically have not just a peak at 0 but also a fat tail⁵. Like the other linear factor models (Eq. 15.3), sparse coding corresponds to a linear factor model, but one with a “sparse” latent variable \mathbf{h} , i.e., $P(\mathbf{h})$ puts high probability at or around 0. Unlike with ICA (previous section), the latent variable prior is parametric. For example the factorized Laplace density prior is

$$P(\mathbf{h}) = \prod_i P(h_i) = \prod_i \frac{\lambda}{2} e^{-\lambda|h_i|} \quad (15.5)$$

and the factorized Student-t prior is

$$P(\mathbf{h}) = \prod_i P(h_i) \propto \prod_i \frac{1}{1 + \frac{h_i^2}{\nu}}^{\frac{\nu+1}{2}}. \quad (15.6)$$

Both of these densities have a strong preference for near-zero values but, unlike the Gaussian, accomodate large values. In the standard sparse coding models, the reconstruction noise is assumed to be Gaussian, so that the corresponding reconstruction error is the squared error.

Regarding sparsity, note that the actual value $h_i = 0$ has zero measure under both densities, meaning that the posterior distribution $P(\mathbf{h} \mid \mathbf{x})$ will not generate values $\mathbf{h} = 0$. However, sparse coding is normally considered under a maximum a posteriori (MAP) inference framework, in which the inferred values of \mathbf{h} are those that maximize the posterior, and these tend to often be zero if the prior is sufficiently concentrated around 0. The inferred values are those defined in Eq. 15.1, reproduced here,

$$\mathbf{h} = f(\mathbf{x}) = \arg \min_{\mathbf{h}} L(g(\mathbf{h}), \mathbf{x}) + \lambda \Omega(\mathbf{h})$$

where $L(g(\mathbf{h}), \mathbf{x})$ is interpreted as $-\log P(\mathbf{x} \mid g(\mathbf{h}))$ and $\Omega(\mathbf{h})$ as $-\log P(\mathbf{h})$. This MAP inference view of sparse coding and an interesting probabilistic interpretation of sparse coding are further discussed in Section 19.3.

To relate the generative model of sparse coding to ICA, note how the prior imposes not just sparsity but also independence of the latent variables h_i under $P(\mathbf{h})$, which may help to separate different explanatory factors, unlike PCA, factor analysis or probabilistic PCA, because these rely on a Gaussian prior, which yields a factorized prior under any rotation of the factors, multiplication by an orthonormal matrix, as demonstrated in Section 15.6.1.

⁵with probability going to 0 as the values increase in magnitude at a rate that is slower than the Gaussian, i.e., less than quadratic in the log-domain.

See Section 17.2 about the manifold interpretation of sparse coding.

TODO: relate to and point to Spike-and-slab sparse coding (Goodfellow *et al.*, 2012) (section?)

15.7 Reconstruction Error as Log-Likelihood

Although traditional auto-encoders (like traditional neural networks) were introduced with an associated training loss, just like for neural networks, that training loss can generally be given a probabilistic interpretation as a conditional log-likelihood of the original input \mathbf{x} , given the representation \mathbf{h} .

We have already covered negative log-likelihood as a loss function in general for feedforward neural networks in Section 6.3.2. Like prediction error for regular feedforward neural networks, reconstruction error for auto-encoders does not have to be squared error. When we view the loss as negative log-likelihood, we interpret the reconstruction error as

$$L = -\log P(\mathbf{x} \mid \mathbf{h})$$

where \mathbf{h} is the representation, which may generally be obtained through an encoder taking \mathbf{x} as input.

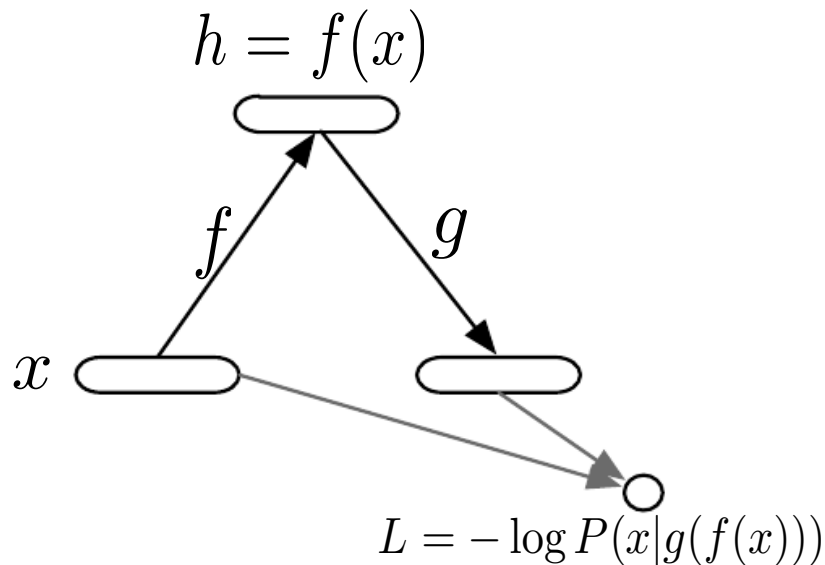


Figure 15.5: The computational graph of an auto-encoder, which is trained to maximize the probability assigned by the decoder g to the data point \mathbf{x} , given the output of the encoder $\mathbf{h} = f(\mathbf{x})$. The training objective is thus $L = -\log P(\mathbf{x} \mid g(f(\mathbf{x})))$, which ends up being squared reconstruction error if we choose a Gaussian reconstruction distribution with mean $g(f(\mathbf{x}))$, and cross-entropy if we choose a factorized Bernoulli reconstruction distribution with means $g(f(\mathbf{x}))$.

An advantage of this view is that it immediately tells us what kind of loss function one should use depending on the nature of the input. If the input is real-valued and unbounded, then squared error is a reasonable choice of reconstruction error, and corresponds to $P(\mathbf{x} \mid \mathbf{h})$ being Normal. If the input is a vector of bits, then cross-entropy is a more reasonable choice, and corresponds to $P(\mathbf{x} \mid \mathbf{h}) = \prod_i P(x_i \mid \mathbf{h})$ with $x_i \mid \mathbf{h}$ being Bernoulli-distributed. We then view the decoder $g(\mathbf{h})$ as computing the *parameters* of the reconstruction distribution, i.e., $P(\mathbf{x} \mid \mathbf{h}) = P(\mathbf{x} \mid g(\mathbf{h}))$.

Another advantage of this view is that we can think about the training of the decoder as estimating the conditional distribution $P(\mathbf{x} \mid \mathbf{h})$, which comes handy in the probabilistic interpretation of denoising auto-encoders, allowing us to talk about the distribution $P(\mathbf{x})$ explicitly or implicitly represented by the auto-encoder (see Sections 15.9, 20.9.3 and 20.10 for more details). In the same spirit, we can rethink the notion of encoder from a simple function to a conditional distribution $Q(\mathbf{h} \mid \mathbf{x})$, with a special case being when $Q(\mathbf{h} \mid \mathbf{x})$ is a Dirac at some particular value. Equivalently, thinking about the encoder as a distribution corresponds to *injecting noise* inside the auto-encoder. This view is developed further in Sections 20.9.3 and 20.11.

15.8 Sparse Representations

Sparse auto-encoders are auto-encoders which learn a sparse representation, i.e., one whose elements are often either zero or close to zero. Sparse coding was introduced in Section 15.6.2 as a linear factor model in which the prior $P(\mathbf{h})$ on the representation $\mathbf{h} = f(\mathbf{x})$ encourages values at or near 0. In Section 15.8.1, we see how ordinary auto-encoders can be prevented from learning a useless identity transformation by using a sparsity penalty rather than a bottleneck. The main difference between a sparse auto-encoder and sparse coding is that sparse coding has no explicit parametric encoder, whereas sparse auto-encoders have one. The “encoder” of sparse coding is the algorithm that performs the approximate inference, i.e., looks for

$$\mathbf{h}^*(\mathbf{x}) = \arg \max_{\mathbf{h}} \log P(\mathbf{h} \mid \mathbf{x}) = \arg \min_{\mathbf{h}} \frac{\|\mathbf{x} - (\mathbf{b} + \mathbf{W}\mathbf{h})\|^2}{\sigma^2} - \log P(\mathbf{h}) \quad (15.7)$$

where σ^2 is a reconstruction variance parameter (which should equal the average squared reconstruction error⁶), and $P(\mathbf{h})$ is a “sparse” prior that puts more prob-

⁶but can be lumped into the regularizer λ which controls the strength of the sparsity prior, defined in Eq. 15.8, for example.

ability mass around $\mathbf{h} = 0$, such as the Laplacian prior, with factorized marginals

$$P(h_i) = \frac{\lambda}{2} e^{\lambda|h_i|} \quad (15.8)$$

or the Student-t prior, with factorized marginals

$$P(h_i) \propto \frac{1}{\left(1 + \frac{h_i^2}{\nu}\right)^{\frac{\nu+1}{2}}}. \quad (15.9)$$

The advantages of such a non-parametric encoder and the sparse coding approach over sparse auto-encoders are that

1. it can in principle minimize the combination of reconstruction error and log-prior better than any parametric encoder,
2. it performs what is called *explaining away* (see Figure 13.8), i.e., it allows to “choose” some “explanations” (hidden factors) and inhibits the others.

The disadvantages are that

1. computing time for encoding the given input \mathbf{x} , i.e., performing inference (computing the representation \mathbf{h} that goes with the given \mathbf{x}) can be substantially larger than with a parametric encoder (because an optimization must be performed *for each example* \mathbf{x}), and
2. the resulting encoder function could be non-smooth and possibly too non-linear (with two nearby \mathbf{x} ’s being associated with very different \mathbf{h} ’s), potentially making it more difficult for the downstream layers to properly generalize.

In Section 15.8.2, we describe PSD (Predictive Sparse Decomposition), which combines a non-parametric encoder (as in sparse coding, with the representation obtained via an optimization) and a parametric encoder (like in the sparse auto-encoder). Section 15.9 introduces the Denoising Auto-Encoder (DAE), which puts pressure on the representation by requiring it to extract information about the underlying distribution and where it concentrates, so as to be able to denoise a corrupted input. Section 15.10 describes the Contractive Auto-Encoder (CAE), which optimizes an explicit regularization penalty that aims at making the representation as insensitive as possible to the input, while keeping the information sufficient to reconstruct the training examples.

15.8.1 Sparse Auto-Encoders

A sparse auto-encoder is simply an auto-encoder whose training criterion involves a sparsity penalty $\Omega(\mathbf{h})$ in addition to the reconstruction error:

$$L = -\log P(\mathbf{x} \mid g(\mathbf{h})) + \Omega(\mathbf{h}) \quad (15.10)$$

where $g(\mathbf{h})$ is the decoder output and typically we have $\mathbf{h} = f(\mathbf{x})$, the encoder output.

We can think of that penalty $\Omega(\mathbf{h})$ simply as a regularizer or as a log-prior on the representations \mathbf{h} . For example, the sparsity penalty corresponding to the Laplace prior ($\frac{\lambda}{2} e^{-\lambda|h_i|}$) is the absolute value sparsity penalty (see also Eq. 15.8 above):

$$\begin{aligned} \Omega(\mathbf{h}) &= \lambda \sum_i |h_i| \\ -\log P(\mathbf{h}) &= \sum_i \log \frac{\lambda}{2} + \lambda|h_i| = \text{const} + \Omega(\mathbf{h}) \end{aligned} \quad (15.11)$$

where the constant term depends only of λ and not \mathbf{h} (which we typically ignore in the training criterion because we consider λ as a hyperparameter rather than a parameter). Similarly (as per Eq. 15.9), the sparsity penalty corresponding to the Student-t prior (Olshausen and Field, 1997) is

$$\Omega(\mathbf{h}) = \sum_i \frac{\nu + 1}{2} \log(1 + \frac{h_i^2}{\nu}) \quad (15.12)$$

where ν is considered to be a hyperparameter.

The early work on sparse auto-encoders (Ranzato *et al.*, 2007a, 2008) considered various forms of sparsity and proposed a connection between sparsity regularization and the partition function gradient in energy-based models (see Section TODO). The idea is that a regularizer such as sparsity makes it difficult for an auto-encoder to achieve zero reconstruction error everywhere. If we consider reconstruction error as a proxy for energy (unnormalized log-probability of the data), then minimizing the training set reconstruction error forces the energy to be low on training examples, while the regularizer prevents it from being low everywhere. The same role is played by the gradient of the partition function in energy-based models such as the RBM (Section TODO).

However, the sparsity penalty of sparse auto-encoders does not need to have a probabilistic interpretation. For example, Goodfellow *et al.* (2009) successfully used the following sparsity penalty, which does not try to bring h_i all the way down to 0, but only towards some low target value such as $\rho = 0.05$.

$$\Omega(\mathbf{h}) = \sum_i \rho \log h_i + (1 - \rho) \log(1 - h_i) \quad (15.13)$$

where $0 < h_i < 1$, usually with $h_i = \text{sigmoid}(a_i)$. This is just the cross-entropy between the Bernoulli distributions with probability $p = h_i$ and the target Bernoulli distribution with probability $p = \rho$.

One way to achieve *actual zeros* in \mathbf{h} for sparse (and denoising) auto-encoders was introduced in Glorot *et al.* (2011c). The idea is to use a half-rectifier (a.k.a. simply as “rectifier”) or ReLU (Rectified Linear Unit, introduced in Glorot *et al.* (2011b) for deep supervised networks and earlier in Nair and Hinton (2010a) in the context of RBMs) as the output non-linearity of the encoder. With a prior that actually pushes the representations to zero (like the absolute value penalty), one can thus indirectly control the average number of zeros in the representation. ReLUs were first successfully used for *deep feedforward networks* in Glorot *et al.* (2011a), achieving for the first time the ability to *train fairly deep supervised networks without the need for unsupervised pre-training*, and this turned out to be an important component in the 2012 object recognition breakthrough with deep convolutional networks (Krizhevsky *et al.*, 2012b).

Interestingly, the “regularizer” used in sparse auto-encoders does not conform to the classical interpretation of regularizers as priors on the parameters. That classical interpretation of the regularizer comes from the MAP (Maximum A Posteriori) point estimation (see Section 5.5.1) of parameters associated with the Bayesian view of parameters as random variables and considering the joint distribution of data \mathbf{x} and parameters $\boldsymbol{\theta}$ (see Section 5.7):

$$\arg \max_{\boldsymbol{\theta}} P(\boldsymbol{\theta} \mid \mathbf{x}) = \arg \max_{\boldsymbol{\theta}} (\log P(\mathbf{x} \mid \boldsymbol{\theta}) + \log P(\boldsymbol{\theta}))$$

where the first term on the right is the usual data log-likelihood term and the second term, the log-prior over parameters, incorporates the preference over particular values of $\boldsymbol{\theta}$.

With regularized auto-encoders such as sparse auto-encoders and contractive auto-encoders, instead, the regularizer corresponds to a *log-prior over the representation, or over latent variables*. In the case of sparse auto-encoders, predictive sparse decomposition and contractive auto-encoders, the regularizer specifies a *preference over functions* of the data, rather than over parameters. This makes such a regularizer *data-dependent*, unlike the classical parameter log-prior. Specifically, in the case of the sparse auto-encoder, it says that we prefer an encoder whose output produces values closer to 0. Indirectly (when we marginalize over the training distribution), this is also indicating a preference over parameters, of course.

15.8.2 Predictive Sparse Decomposition

TODO: we have too many forward refs to this section. There are 150 lines about PSD in this section and at least 20 lines of forward references to this section

in this chapter, some of which are just 100 lines away. Predictive sparse decomposition (PSD) is a variant that combines sparse coding and a parametric encoder (Kavukcuoglu *et al.*, 2008b), i.e., it has both a parametric encoder and iterative inference. It has been applied to unsupervised feature learning for object recognition in images and video (Kavukcuoglu *et al.*, 2009, 2010b; Jarrett *et al.*, 2009a; Farabet *et al.*, 2011), as well as for audio (Henaff *et al.*, 2011). The representation is considered to be a free variable (possibly a latent variable if we choose a probabilistic interpretation) and the training criterion combines a sparse coding criterion with a term that encourages the optimized sparse representation \mathbf{h} (after inference) to be close to the output of the encoder $f(\mathbf{x})$:

$$L = \arg \min_{\mathbf{h}} (||\mathbf{x} - g(\mathbf{h})||^2 + \lambda ||\mathbf{h}||_1 + \gamma ||\mathbf{h} - f(\mathbf{x})||^2) \quad (15.14)$$

where f is the encoder and g is the decoder. Like in sparse coding, for each example \mathbf{x} an iterative optimization is performed in order to obtain a representation \mathbf{h} . However, because the iterations can be initialized from the output of the encoder, i.e., with $\mathbf{h} = f(\mathbf{x})$, only a few steps (e.g. 10) are necessary to obtain good results. Simple gradient descent on \mathbf{h} has been used by the authors. After \mathbf{h} is settled, both g and f are updated towards minimizing the above criterion. The first two terms are the same as in L1 sparse coding while the third one encourages f to predict the outcome of the sparse coding optimization, making it a better choice for the initialization of the iterative optimization. Hence f can be used as a parametric approximation to the non-parametric encoder implicitly defined by sparse coding. It is one of the first instances of *learned approximate inference* (see also Sec. 19.6). Note that this is different from separately doing sparse coding (i.e., training g) and then training an approximate inference mechanism f , since both the encoder and decoder are trained together to be “compatible” with each other. Hence the decoder will be learned in such a way that inference will tend to find solutions that can be well approximated by the approximate inference. TODO: this is probably too much forward reference, when we bring these things in we can remind people that they resemble PSD, but it doesn’t really help the reader to say that the thing we are describing now is similar to things they haven’t seen yet A similar example is the variational auto-encoder, in which the encoder acts as approximate inference for the decoder, and both are trained jointly (Section 20.9.3). See also Section 20.9.4 for a probabilistic interpretation of PSD in terms of a variational lower bound on the log-likelihood.

In practical applications of PSD, the iterative optimization is only used during training, and f is used to compute the learned features. It makes computation fast at recognition time and also makes it easy to use the trained features f as initialization (unsupervised pre-training) for the lower layers of a deep net. Like other unsupervised feature learning schemes, PSD can be stacked greedily, e.g.,

training a second PSD on top of the features extracted by the first one, etc.

15.9 Denoising Auto-Encoders

The Denoising Auto-Encoder (DAE) was first proposed (Vincent *et al.*, 2008, 2010) as a means of forcing an auto-encoder to learn to capture the data distribution without an explicit constraint on either the dimension or the sparsity of the learned representation. It was motivated by the idea that in order to fully capture a complex distribution, an auto-encoder needs to have at least as many hidden units as needed by the complexity of that distribution. Hence its dimensionality should not be restricted to the input dimension.

The principle of the denoising auto-encoder is deceptively simple and illustrated in Figure 15.6: the encoder sees as input a corrupted version of the input, but the decoder tries to reconstruct the clean uncorrupted input.

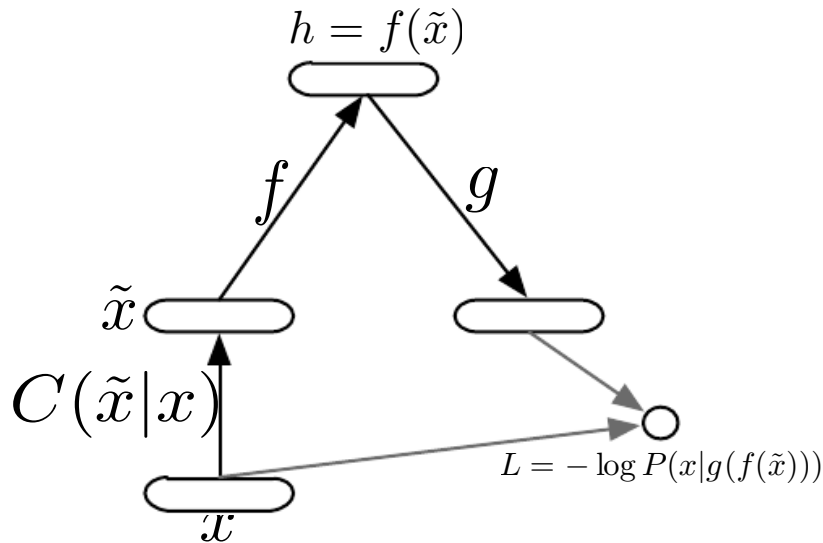


Figure 15.6: The computational graph of a denoising auto-encoder, which is trained to reconstruct the clean data point \mathbf{x} from its corrupted version $\tilde{\mathbf{x}}$, i.e., to minimize the loss $L = -\log P(\mathbf{x} | g(f(\tilde{\mathbf{x}})))$, where $\tilde{\mathbf{x}}$ is a corrupted version of the data example \mathbf{x} , obtained through a given corruption process $C(\tilde{\mathbf{x}} | \mathbf{x})$.

Mathematically, and following the notations used in this chapter, this can be formalized as follows. We introduce a corruption process $C(\tilde{\mathbf{x}} | \mathbf{x})$ which represents a conditional distribution over corrupted samples $\tilde{\mathbf{x}}$, given a data sample \mathbf{x} . The auto-encoder then learns a *reconstruction distribution* $P(\mathbf{x} | \tilde{\mathbf{x}})$ estimated from training pairs $(\mathbf{x}, \tilde{\mathbf{x}})$, as follows:

1. Sample a training example $\mathbf{x} = \mathbf{x}$ from the data generating distribution (the

training set).

2. Sample a corrupted version $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}$ from the conditional distribution $C(\tilde{\mathbf{x}} \mid \mathbf{x} = \mathbf{x})$.
3. Use $(\mathbf{x}, \tilde{\mathbf{x}})$ as a training example for estimating the auto-encoder reconstruction distribution $P(\mathbf{x} \mid \tilde{\mathbf{x}}) = P(\mathbf{x} \mid g(\mathbf{h}))$ with \mathbf{h} the output of encoder $f(\tilde{\mathbf{x}})$ and $g(\mathbf{h})$ the output of the decoder.

Typically we can simply perform gradient-based approximate minimization (such as minibatch gradient descent) on the negative log-likelihood $-\log P(\mathbf{x} \mid \mathbf{h})$, i.e., the denoising reconstruction error, using back-propagation to compute gradients, just like for regular feedforward neural networks (the only difference being the corruption of the input and the choice of target output).

We can view this training objective as performing stochastic gradient descent on the denoising reconstruction error, but where the “noise” now has two sources:

1. the choice of training sample \mathbf{x} from the data set, and
2. the random corruption applied to \mathbf{x} to obtain $\tilde{\mathbf{x}}$.

We can therefore consider that the DAE is performing stochastic gradient descent on the following expectation:

$$-E_{\mathbf{x} \sim Q(\mathbf{x})} E_{\tilde{\mathbf{x}} \sim C(\tilde{\mathbf{x}} \mid \mathbf{x})} \log P(\mathbf{x} \mid g(f(\tilde{\mathbf{x}})))$$

where $Q(\mathbf{x})$ is the training distribution.

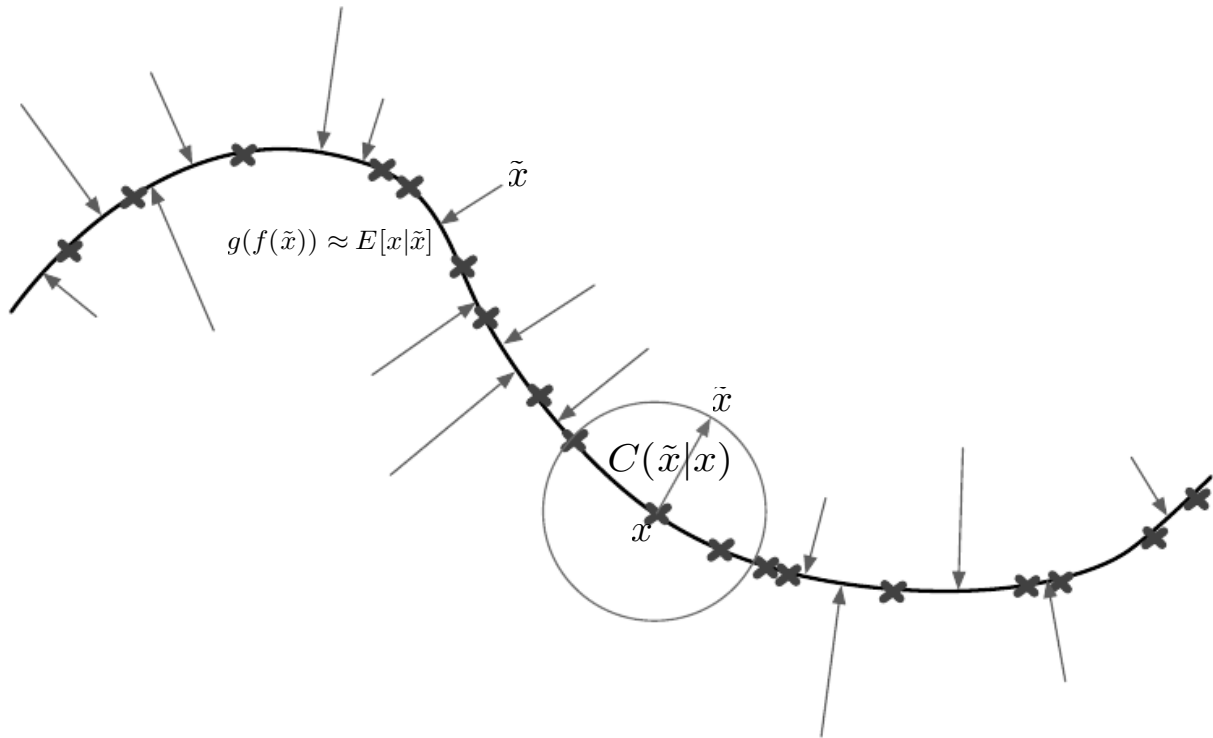


Figure 15.7: A denoising auto-encoder is trained to reconstruct the clean data point \mathbf{x} from its corrupted version $\tilde{\mathbf{x}}$. In the figure, we illustrate the corruption process $C(\tilde{\mathbf{x}} | \mathbf{x})$ by a grey circle of equiprobable corruptions, and grey arrows for the corruption process acting on examples \mathbf{x} (red crosses) lying near a low-dimensional manifold near which probability concentrates. When the denoising auto-encoder is trained to minimize the average of squared errors $\|g(f(\tilde{\mathbf{x}})) - \mathbf{x}\|^2$, the reconstruction $g(f(\tilde{\mathbf{x}}))$ estimates $\mathbb{E}[\mathbf{x} | \tilde{\mathbf{x}}]$, which approximately points orthogonally towards the manifold, since it estimates the center of mass of the clean points \mathbf{x} which could have given rise to $\tilde{\mathbf{x}}$. The auto-encoder thus learns a vector field $g(f(\tilde{\mathbf{x}})) - \mathbf{x}$ (the green arrows) and it turns out that this vector field estimates the gradient field $\frac{\partial \log Q(\mathbf{x})}{\partial \mathbf{x}}$ (up to a multiplicative factor that is the average root mean square reconstruction error), where Q is the unknown data generating distribution.

15.9.1 Learning a Vector Field that Estimates a Gradient Field

As illustrated in Figure 15.7, a very important property of DAEs is that their training criterion makes the auto-encoder learn a vector field $(g(f(\tilde{\mathbf{x}})) - \mathbf{x})$ that estimates the gradient field (or *score*) $\frac{\partial \log Q(\mathbf{x})}{\partial \mathbf{x}}$, as per Eq. 15.15. A first result in this direction was proven by Vincent (2011a), showing that minimizing squared reconstruction error in a denoising auto-encoder with Gaussian noise was related to *score matching* (Hyvärinen, 2005a), making the denoising criterion a regularized form of score matching called *denoising score matching* (Kingma and LeCun, 2010a). Score matching is an alternative to maximum likelihood and provides a consistent estimator. It is discussed further in Section 18.4. The denoising version

is discussed in Section 18.5.

The connection between denoising auto-encoders and score matching was first made (Vincent, 2011a) in the case where the denoising auto-encoder has a particular parametrization (one hidden layer, sigmoid activation functions on hidden units, linear reconstruction), in which case the denoising criterion actually corresponds to a regularized form of score matching on a Gaussian RBM (with binomial hidden units and Gaussian visible units). The connection between ordinary auto-encoders and Gaussian RBMs had previously been made by Bengio and Delalleau (2009), which showed that contrastive divergence training of RBMs was related to an associated auto-encoder gradient, and later by Swersky (2010), which showed that non-denoising reconstruction error corresponded to score matching plus a regularizer.

The fact that the denoising criterion yields an estimator of the score for general encoder/decoder parametrizations has been proven (Alain and Bengio, 2012, 2013) in the case where the corruption and the reconstruction distributions are Gaussian (and of course \mathbf{x} is continuous-valued), i.e., with the squared error denoising error

$$||g(f(\tilde{\mathbf{x}})) - \mathbf{x}||^2$$

and corruption

$$C(\tilde{\mathbf{x}} = \tilde{\mathbf{x}}|\mathbf{x}) = \mathcal{N}(\tilde{\mathbf{x}}; \mu = \mathbf{x}, \Sigma = \sigma^2 I)$$

with noise variance σ^2 .

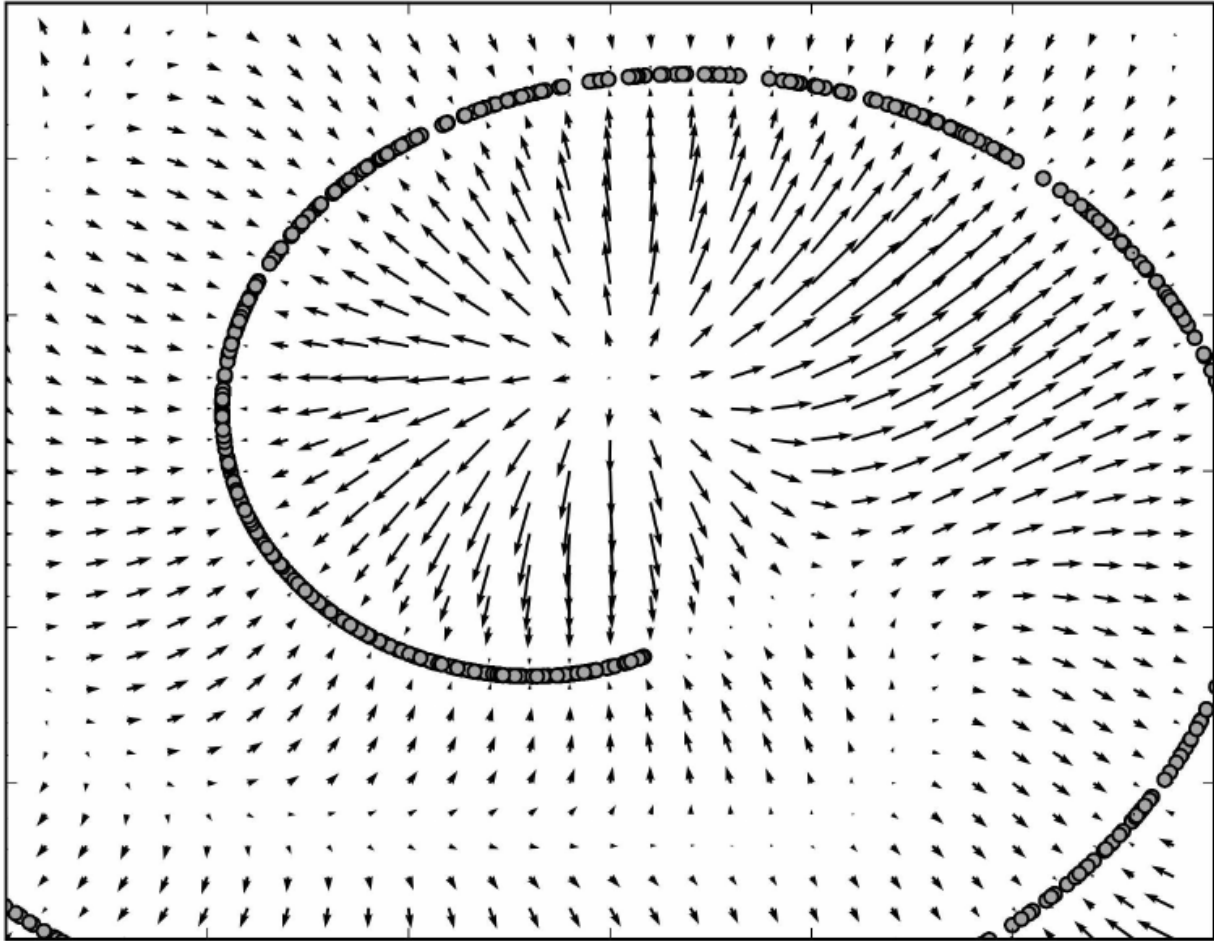


Figure 15.8: Vector field learned by a denoising auto-encoder around a 1-D curved manifold near which the data (orange circles) concentrates in a 2-D space. Each arrow is proportional to the reconstruction minus input vector of the auto-encoder and points towards higher probability according to the implicitly estimated probability distribution. Note that the vector field has zeros at both peaks of the estimated density function (on the data manifolds) and at troughs (local minima) of that density function, e.g., on the curve that separates different arms of the spiral or in the middle of it.

More precisely, the main theorem states that $\frac{g(f(\mathbf{x})) - \mathbf{x}}{\sigma^2}$ is a consistent estimator of $\frac{\partial \log Q(\mathbf{x})}{\partial \mathbf{x}}$, where $Q(\mathbf{x})$ is the data generating distribution,

$$\frac{g(f(\mathbf{x})) - \mathbf{x}}{\sigma^2} \rightarrow \frac{\partial \log Q(\mathbf{x})}{\partial \mathbf{x}}, \quad (15.15)$$

so long as f and g have sufficient capacity to represent the true score (and assuming that the expected training criterion can be minimized, as usual when proving consistency associated with a training objective).

Note that in general, there is no guarantee that the reconstruction $g(f(\mathbf{x}))$ minus the input \mathbf{x} corresponds to the gradient of something (the estimated score should be the gradient of the estimated log-density with respect to the input

\mathbf{x}). That is why the early results (Vincent, 2011a) are specialized to particular parametrizations where $g(f(\mathbf{x})) - \mathbf{x}$ is the derivative of something. See a more general treatment by Kamyshanska and Memisevic (2015).

Although it was intuitively appealing that in order to denoise correctly one must capture the training distribution, the above consistency result makes it mathematically very clear in what sense the DAE is capturing the input distribution: it is estimating the gradient of its energy function (i.e., of its log-density), i.e., learning to point towards more probable (lower energy) configurations. Figure 15.8 (see details of experiment in Alain and Bengio (2013)) illustrates this. Note how the norm of reconstruction error (i.e. the norm of the vectors shown in the figure) is related to but *different* from the energy (unnormalized log-density) associated with the estimated model. The energy should be low only where the probability is high. The reconstruction error (norm of the estimated score vector) is low where probability is near a peak of probability (or a trough of energy), but it can also be low at *maxima* of energy (minima of probability).

Section 20.10 continues the discussion of the relationship between denoising auto-encoders and probabilistic modeling by showing how one can *generate* from the distribution implicitly estimated by a denoising auto-encoder. Whereas (Alain and Bengio, 2013) generalized the score estimation result of Vincent (2011a) to arbitrary parametrizations, the result from Bengio *et al.* (2013b), discussed in Section 20.10, provides a probabilistic – and in fact generative – interpretation to every denoising auto-encoder.

15.10 Contractive Auto-Encoders

The Contractive Auto-Encoder or CAE (Rifai *et al.*, 2011a,c) introduces an explicit regularizer on the code $\mathbf{h} = f(\mathbf{x})$, encouraging the derivatives of f to be as small as possible:

$$\Omega(\mathbf{h}) = \left\| \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right\|_F^2 \quad (15.16)$$

which is the squared Frobenius norm (sum of squared elements) of the Jacobian matrix of partial derivatives associated with the encoder function. Whereas the denoising auto-encoder learns to contract the reconstruction function (the composition of the encoder and decoder), the CAE learns to specifically contract the encoder. See Figure 17.13 for a view of how contraction near the data points makes the auto-encoder capture the manifold structure.

If it weren't for the opposing force of reconstruction error, which attempts to make the code \mathbf{h} keep all the information necessary to *reconstruct training examples*, the CAE penalty would yield a code \mathbf{h} that is constant and does not

depend on the input \mathbf{x} . The compromise between these two forces yields an auto-encoder whose derivatives $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$ are tiny in most directions, except those that are needed to reconstruct training examples, i.e., the directions that are tangent to the manifold near which data concentrate. Indeed, in order to distinguish (and thus, reconstruct correctly) two nearby examples on the manifold, one must assign them a different code, i.e., $f(\mathbf{x})$ must vary as \mathbf{x} moves from one to the other, i.e., in the direction of a tangent to the manifold.

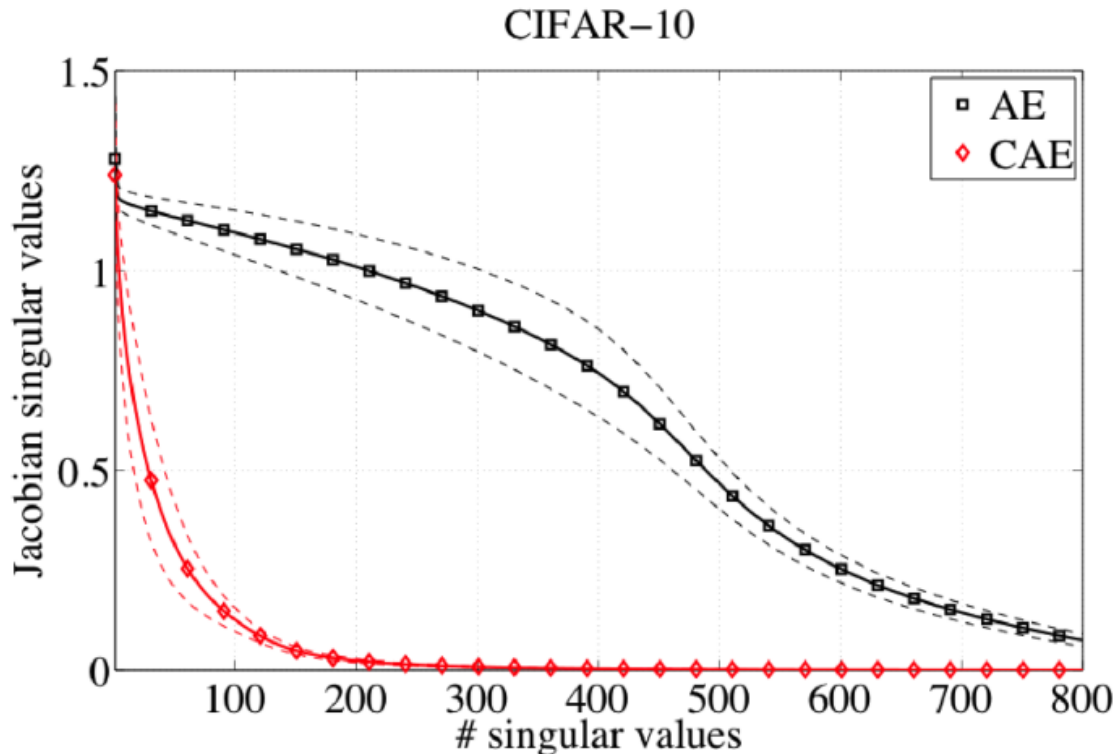


Figure 15.9: Average (over test examples) of the singular value spectrum of the Jacobian matrix $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$ for the encoder f learned by a regular auto-encoder (AE) versus a contractive auto-encoder (CAE). This illustrates how the contractive regularizer yields a smaller set of directions in input space (those corresponding to large singular value of the Jacobian) which provoke a response in the representation \mathbf{h} while the representation remains almost insensitive for most directions of change in the input.

What is interesting is that this penalty forces more strongly the representation to be invariant in directions orthogonal to the manifold. This can be seen clearly by comparing the singular value spectrum of the Jacobian $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$ for different auto-encoders, as shown in Figure 15.9. We see that the CAE manages to concentrate the sensitivity of the representation in fewer dimensions than a regular (or sparse) auto-encoder. Figure 17.3 illustrates tangent vectors obtained by a CAE on the MNIST digits dataset, showing that the leading tangent vectors correspond to small deformations such as translation. More impressively, Figure 15.10 shows

tangent vectors learned on 32×32 color (RGB) CIFAR-10 images by a CAE,

compared to the tangent vectors by a non-distributed representation learner (a mixture of local PCAs).

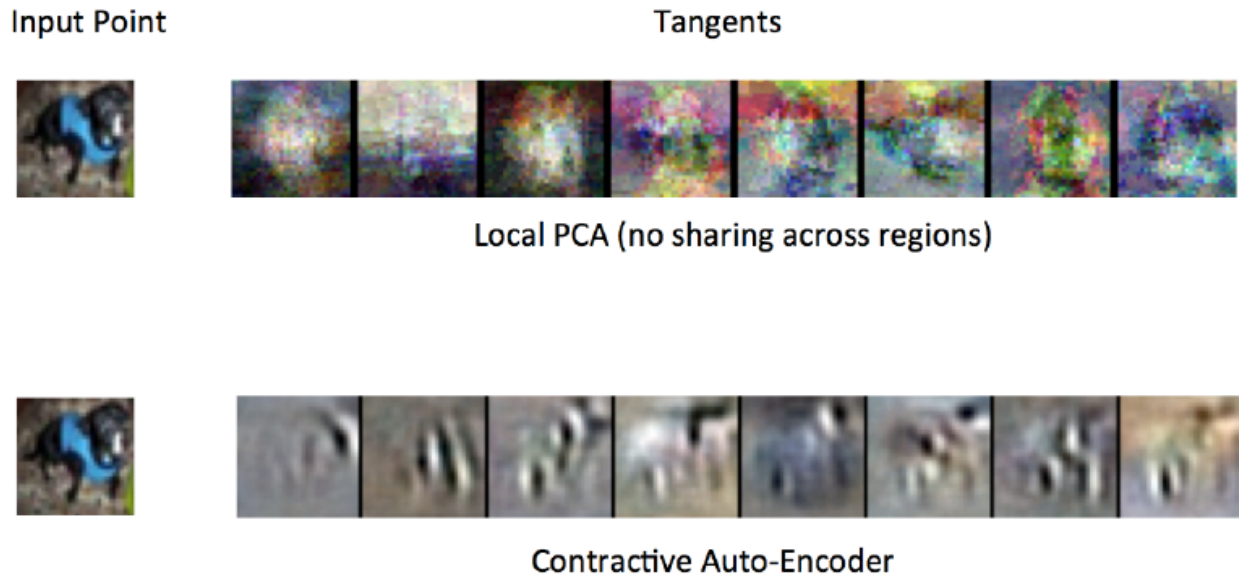


Figure 15.10: Illustration of tangent vectors (bottom) of the manifold estimated by a contractive auto-encoder (CAE), at some input point (left, CIFAR-10 image of a dog). See also Fig. 17.3. Each image on the right corresponds to a tangent vector, either estimated by a local PCA (equivalent to a Gaussian mixture), top, or by a CAE (bottom). The tangent vectors are estimated by the leading singular vectors of the Jacobian matrix $\frac{\partial h}{\partial x}$ of the input-to-code mapping. Although both local PCA and CAE can capture local tangents that are different in different points, the local PCA does not have enough training data to meaningful capture good tangent directions, whereas the CAE does (because it exploits parameter sharing across different locations that share a subset of active hidden units). The CAE tangent directions typically correspond to moving or changing parts of the object (such as the head or legs), which corresponds to plausible changes in the input image.

One practical issue with the CAE regularization criterion is that although it is cheap to compute in the case of a single hidden layer auto-encoder, it becomes much more expensive in the case of deeper auto-encoders. The strategy followed by Rifai *et al.* (2011a) is to separately pre-train each single-layer auto-encoder stacked to form a deeper auto-encoder. However, a deeper encoder could be advantageous in spite of the computational overhead, as argued by Schulz and Behnke (2012).

Another practical issue is that the contraction penalty on the encoder f could yield useless results if the decoder g would exactly compensate (e.g. by being scaled up by exactly the same amount as f is scaled down). In Rifai *et al.* (2011a), this is compensated by tying the weights of f and g , both being of the form of an affine transformation followed by a non-linearity (e.g. sigmoid), i.e., the weights of g and the transpose of the weights of f .