

nonlinear change of variables (using equation 3.47) to determine $p(\mathbf{x})$. Learning the model then proceeds as usual, using maximum likelihood.

The motivation for this approach is that by choosing $p(\mathbf{h})$ to be independent, we can recover underlying factors that are as close as possible to independent. This is commonly used, not to capture high-level abstract causal factors, but to recover low-level signals that have been mixed together. In this setting, each training example is one moment in time, each x_i is one sensor's observation of the mixed signals, and each h_i is one estimate of one of the original signals. For example, we might have n people speaking simultaneously. If we have n different microphones placed in different locations, ICA can detect the changes in the volume between each speaker as heard by each microphone, and separate the signals so that each h_i contains only one person speaking clearly. This is commonly used in neuroscience for electroencephalography, a technology for recording electrical signals originating in the brain. Many electrode sensors placed on the subject's head are used to measure many electrical signals coming from the body. The experimenter is typically only interested in signals from the brain, but signals from the subject's heart and eyes are strong enough to confound measurements taken at the subject's scalp. The signals arrive at the electrodes mixed together, so ICA is necessary to separate the electrical signature of the heart from the signals originating in the brain, and to separate signals in different brain regions from each other.

As mentioned before, many variants of ICA are possible. Some add some noise in the generation of \mathbf{x} rather than using a deterministic decoder. Most do not use the maximum likelihood criterion, but instead aim to make the elements of $\mathbf{h} = \mathbf{W}^{-1}\mathbf{x}$ independent from each other. Many criteria that accomplish this goal are possible. Equation 3.47 requires taking the determinant of \mathbf{W} , which can be an expensive and numerically unstable operation. Some variants of ICA avoid this problematic operation by constraining \mathbf{W} to be orthogonal.

All variants of ICA require that $p(\mathbf{h})$ be non-Gaussian. This is because if $p(\mathbf{h})$ is an independent prior with Gaussian components, then \mathbf{W} is not identifiable. We can obtain the same distribution over $p(\mathbf{x})$ for many values of \mathbf{W} . This is very different from other linear factor models like probabilistic PCA and factor analysis, that often require $p(\mathbf{h})$ to be Gaussian in order to make many operations on the model have closed form solutions. In the maximum likelihood approach where the user explicitly specifies the distribution, a typical choice is to use $p(h_i) = \frac{d}{dh_i}\sigma(h_i)$. Typical choices of these non-Gaussian distributions have larger peaks near 0 than does the Gaussian distribution, so we can also see most implementations of ICA as learning sparse features.