## Abstract

With the resurgence of interest in neural networks, representation learning has re-emerged as a central focus in artificial intelligence. Representation learning refers to the discovery of useful encodings of data that make domain-relevant information explicit. *Factorial representations* identify underlying independent causal factors of variation in data. A factorial representation is compact and faithful, makes the causal factors explicit, and facilitates human interpretation of data. Factorial representations support a variety of applications, including the generation of novel examples, indexing and search, novelty detection, and transfer learning.

This article surveys various constraints that encourage a learning algorithm to discover factorial representations. I dichotomize the constraints in terms of unsupervised and supervised inductive bias. Unsupervised inductive biases exploit assumptions about the environment, such as the statistical distribution of factor coefficients, assumptions about the perturbations a factor should be invariant to (e.g. a representation of an object can be invariant to rotation, translation or scaling), and assumptions about how factors are combined to synthesize an observation. Supervised inductive biases are constraints on the representations based on additional information connected to observations. Supervisory labels come in variety of types, which vary in how strongly they constrain the representation, how many factors are labeled, how many observations are labeled, and whether or not we know the associations between the constraints and the factors they are related to.

This survey brings together a wide variety of models that all touch on the problem of learning factorial representations and lays out a framework for comparing these models based on the strengths of the underlying supervised and unsupervised inductive biases.

# A Survey of Inductive Biases for Factorial Representation-Learning

Karl Ridgeway

November 2016

## Contents

# 1   Introduction

With the advent of modern deep neural nets, representation learning has re-emerged as a primary focus in machine learning. Representation learning refers to the discovery of useful representations of data. Data samples are drawn from some environment — for example, photographs of faces or audio recordings of music. A data representation is a means for transforming these raw examples into some new space. Typically, we have some particular task in mind, and the data representation allows us to easily recover attributes from the data relevant to the task at hand. Suppose the task is to create a searchable index of the dataset. The best representation for this task might summarize the face photographs as a set of attributes — ethnicity, hair color, or eye color. We might similarly represent music by its tempo, genre, artist, or key signature. A good representation makes explicit relevant information about the task. While ethnicity and hair color are, in some way, represented in the raw pixels of the image, as is tempo in the audio recording, the more explicit representation makes acting on the information much easier.

## 1.1   Desiderata for Representations

Different tasks and environments place different demands on representations, but there are a few characteristics commonly associated with good representations which are useful across a wide variety of tasks. Good representations are both *compact* and *faithful* to information represented in the input. They also *explicitly* represent the attributes required for the task at hand. Finally, they are *interpretable* by humans.

### 1.1.1   Compact and Faithful

A faithful representation preserves the information in the observation [49], with as little distortion as possible. A perfectly faithful representation can of course be achieved by not doing any transformation on the input — this is why good representations also need to encode the information compactly. This is equivalent to the representation being an efficient compression of the input.

    The faithfulness of a representation can only be evaluated in terms of the task it's intended to be used with. Different tasks have different demands on how much of the input need be represented. A classification task might just need to know object category. An image compression program should represent all of the non-noise features in the image. A representation of medical test results should probably be fully invertible, meaning we can perfectly reconstruct the input from the representation.

### 1.1.2 Explicit

The representation should explicitly encode the attributes necessary for whatever task is at hand, without interference from other attributes or sources of noise. Attributes can be far removed from the low-level observation; for example, facial expression is observed as a complex conjunction of muscle movements around the mouth, the eyes, and jaw. Attributes are often combined in complex ways with other attributes. Consider how the glasses in Figure 1.1 interact with identity, pose, and lighting conditions — they fit onto and cast a shadow on the face. Explicit representations of attributes should be *invariant* to these kinds of variations — identity should not change if the image is slightly corrupted by noise or if the light source is moved, even though many pixels on the face change. Explicit representations are also invariant to small local changes / noise in the input.



Figure 1.1: Random draws from the dataset of subjects wearing various styles of glasses.

A representation of an attribute that makes some information explicit should have a *structure* that reflects the information represented: A categorical, discrete, or continuous attribute needs to be modeled as a variable of the appropriate type. Additionally, attributes can be single- or multi-dimensional. For example, color can be represented in a single dimension (frequency) or as a three-dimensional vector (red/green/blue channels). The choice of variable type depends on the environment and task.

### 1.1.3 Interpretable

Many tasks rely on data representations being *interpretable* by humans. Representations are often better off if they can be easily interpreted by humans, regardless of whether humans will use them directly or not. For example, the task of predicting face gender will be much easier if the presence of facial hair and makeup is directly represented in the code.

Representations that are interpretable need to be unique. A representation that can be exchanged for another that is equally valid is more difficult to attach meaning to. They should also be identifiable, meaning that they are tied to some named concept or attribute.

I have described three desirable properties of representations — compactness/faithfulness, expressivity, and intepretability. Now, I will describe a type of representation, called a factorial representation, that can help achieve these goals.

## 1.2 Factorial Representations

In a *factorial representation*, the attributes are statistically independent. A *factor* is one of these independent attributes. Imagine a set of images of landscapes containing trees and sheep. In the generative story of these landscapes, the location and number of sheep is largely independent of the location and number of trees. Another common example is that of blind source separation — in a room full of people speaking, we can assume that the probability of any one person speaking is independent of others speaking simultaneously. This independence assumption places restrictions on the environment — that the factors are truly independent — and on the task, which should depend on these independent causes. For the right environment and task, factorial representations provide numerous benefits: they provide a useful bias for learning [62] and provide more compact/faithful, explicit, and interpretable representations. Horace Barlow argued for the biological and computational utility of factorial representations [6, 7].

### 1.2.1 Statistical Independence

When variables are statistically independent, the joint probability distribution is factorial. Consider the example set of colored shapes in Figure 1.2. Each observation can be encoded as a combination of a shape and color. The occurrence probability of any shape or color is one-fourth, and the occurrence probability of any combination of the two, one-sixteenth. This two-factor representation is *factorial* because the joint distribution $P(shape, color)$ can be factorized into the product of the distributions of the components $P(shape)P(color)$ — consequently, the component factors are statistically independent of one another, and have no redundancy.
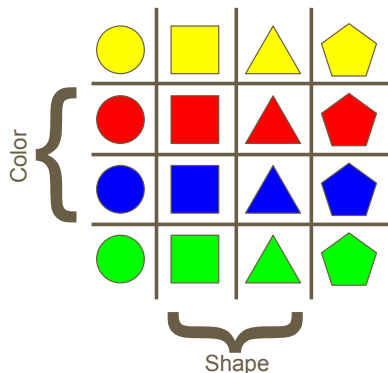


Figure 1.2: An example dataset with two independent factors of variation: shape and color.

### 1.2.2 Useful Properties of Factorial Representations

A factorial representation that is faithful to the factors modeled minimizes redundancy between them. Therefore, these representations are also maximally compact. Factorial representations are also more interpretable: if the representation explicitly captures the true independent generative factors, these factors should be unique and identifiable. The shape and color example illustrates this nicely — there is no other way to describe these images that is as compact, faithful, explicit, and interpretable as "shape and color". Independent factors need to be invariant with respect to one another — no change in shape should alter our impression of an image's color.

These properties make factorial representations very useful for certain types of tasks, such as:

- **Novel example generation**. An invertible factorial representation can be used to generate novel examples not found in the original dataset. In fact, we can manipulate the representation to generate examples of the full Cartesian product of factor combinations — we could generate men with the women's styles of glasses in Figure 1.1.

- **Novelty Detection**. Since factors are independent, any statistical dependence found between factors in a new dataset indicates a new situation [5].

- **Search**. A factorial representation also supports efficient searching — in the example shown in Figure 1.2, specifying a shape narrows a search down to one-quarter of the dataset.

- **Prediction**. More generally, they are useful as features used as input to supervised learning methods, as observed by [8] for linear supervised tasks.

- **Compression**. Factorial representations are compact: a representation that minimizes redundancy also minimizes the cost associated with storing the representation, according to the minimum description length principle [62].

### 1.2.3 Factorial Representations in Non-Factorial Datasets and Environments

The independence assumption of a factorial representation is easily violated. It can be violated by the method used to sample from the environment, or by the environment itself. Consider the case of missing data: If the dataset in Figure 1.2 were missing just one of the 16 shape/color combinations, then $P(shape, color) \neq P(shape)P(color)$. More generally, datasets can have a *sampling bias*, in which certain combinations of factors are systematically missing, or sampled more than other combinations. This imbalance in factors can of course also be caused by noise. The bias

imposed by a factorial representation can help to learn better representations in the case of datasets with a sampling bias or noise.

Often, the environment itself is non-factorial. Datasets have natural dependencies — for example, facial hair and gender are naturally dependent. A factorial representation that teases apart facial hair and gender will over-generate: it will assume that facial hair and gender are independent, so a woman with a beard will be just as likely as a man with a beard, according to the representation.

For some tasks, representing these dependencies is critical, so a factorial representation is inappropriate. For example, a person's identity might be naturally correlated with the style of glasses worn by that person. If our task is to generate example "typical" images of people (what we might expect them to look like on any given day), we want to consider a person's glasses style as part of their identity. For other tasks, we might want to ignore this dependency. If our task is high-security face verification, we want to be sure that our program will not falsely accept an impostor who wears the same glasses as the user. In this setting, imposing the bias of a factorial representation is helpful.

# 2 Inductive Biases

In the last section, we noted that imposing the bias of a factorial representation can result in improved model explicitness as well as better interpretability. This paper explores different *inductive biases* that can be used to learn factorial representations for a variety of environments and tasks.

In this section, I will highlight several different approaches for biasing representations to be more factorial. Each of these biases makes a different assumption about the generative environment. First, I will look at biases that assume different *distributions* of coefficients for a factor. Next, I will look at biases that allow a factor to be *invariant* to changes in the output that are not related to the factor. Another class of bias look at the ways that factors can *combine* to generate the observed data. Finally, I will look at how to apply *supervision* as a bias to improve factorial representations. Since supervision is mostly orthogonal to, and can be arbitrarily combined with the other types of bias, I have given it its own dimension in the model chart in Figure 2.1. Unsupervised biases leverage assumptions about the generative environment to learn factorial representations, while supervised biases are specific to a type of prior domain knowledge.

Here, I focus on inductive biases that are highly generic, meaning that they could apply to almost any observation domain, such as images, audio, or medical test results. The biases considered here allow us to learn representations that are more explicit, and more interpretable — these biases allow us to gain a better understanding of the underlying causal structure, regardless of the observation domain. I will therefore avoid describing inductive biases that are too specific to apply to a range of domains.

For consistency, I will introduce each unsupervised inductive bias the same way. First, I will describe the **generative assumptions** of the model. Then, I will describe an **inductive bias** that is consistent with this environment, and the corresponding learned **representation structure**. Finally, I will show some **examples** of the inductive bias in action, and mention other models that use a similar inductive bias. When possible, I will also show examples that incorporate multiple inductive biases simultaneously.

By matching the inductive biases of a model with the characteristics of the environment and the requirements of the task, it is possible to discover factorial representations in many settings.

## 2.1 Distributional Bias

In this section, I will explore some inductive biases for discovering factorial representations that are motivated by the statistical distribution of the factors in the environment. To highlight the differences between these distribution biases, I will mostly focus on methods that also use a *linear* bias — an assumption that factors are linearly related to, and are linearly combined in, the observation. Keeping this bias constant makes it easier to highlight the differences in distributional bias. Where applicable, I will also highlight non-linear models that incorporate these distributional biases.

I will describe two different assumptions about the distribution of coefficients for factors that each lead to a type of inductive bias. In some environments, factor coefficients follow a *Gaussian* distribution. I will describe an inductive bias that finds factorial representations when the environment has Gaussian-distributed factors. I will show how naturalistic environments are often non-Gaussian, in that they are sparse. Likewise, I will describe a bias that corresponds to sparse environments. In sparse environments, I show that the sparsity bias produces representations that are better aligned to the true causal structure of the environment than representations with Gaussian-distributed factors.
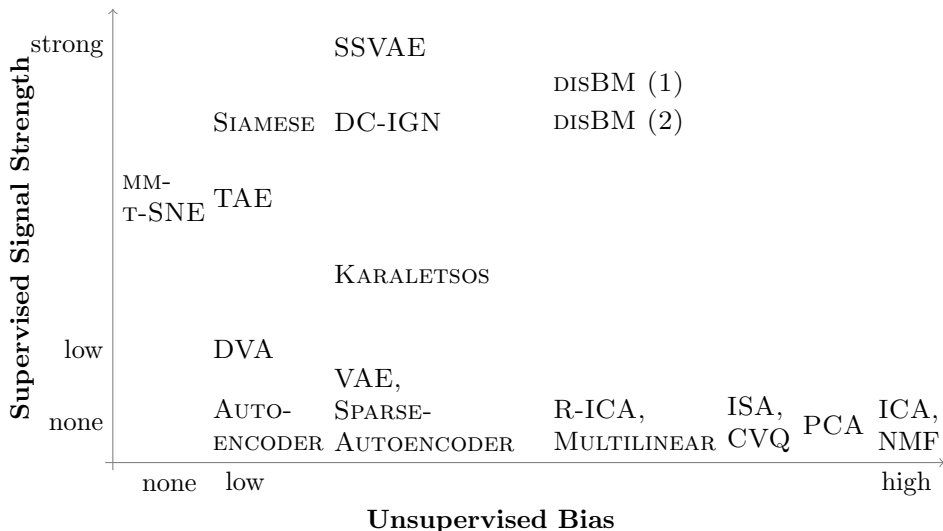
Figure 2.1: Model organization chart

### 2.1.1 Gaussian Distribution of Factors

**Generative Assumptions.** A factorial representation should ideally have non-redundant, statistically independent factors. Since two factors that are jointly Gaussian-distributed are independent when they are uncorrelated, one way to find independent factors in a Gaussian environment is to reduce or minimize the correlation between all pairs of factors. Principle Components Analysis, or PCA, is a model that learns representations with zero correlation between the latent factors.

#### 2.1.1.1 PCA

**Inductive Bias.** PCA's mapping function is linear: a $(1 \times D)$-dimensional input vector x is transformed into its $(1 \times N)$-dimensional latent representation z by multiplying $x$ with a $(D \times N)$ matrix W:

$$\mathbf{z} = \mathbf{x}\mathbf{W} \tag{1}$$

We can also compute the *reconstruction*, or the prediction of x given z.

$$\tilde{\mathbf{x}} = \mathbf{z}\mathbf{W}^{-1} \tag{2}$$

PCA minimizes *reconstruction cost*, or the euclidean distance between the input and the reconstruction:

$$\mathcal{L}_{\mathrm{PCA}} = \|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 \tag{3}$$

PCA also restricts the the dimensions of $\mathbf{z}$ to be uncorrelated — the columns of $\mathbf{W}$ are all orthogonal. Two variables that are jointly Gaussian are independent if they are uncorrelated. Therefore, decorrelation finds independent components in the case of Gaussian-distributed causes. PCA is further constrained to consider only components that are linearly related to the input.[1] PCA produces a list of components, ordered by the amount of variance in the input that they explain. By considering only the top $M$ components, one can trade off between compactness and faithfulness of the representation.

**Representation Structure.** PCA's representation is very simple — each observation is mapped to a vector of scalar factor coefficients $\mathbf{z}$.

Unfortunately, representations learned by PCA are not unique, in that representations are unique only to rotations in the latent space [10]. This means that the representation learned is often just one sample of many possible representations that are equally valid according to PCA. This presents some obvious problems when we try to interpret the results of PCA, since one solution can be exchanged for other, equally valid solutions.

---

[1]A family of related models share similar characteristics — probabilistic PCA adds isotropic Gaussian observation noise to equation 1. Factor analysis extends this further by adding diagonal instead of isotropic noise. Linear autoencoders are essentially a variant of PCA trained with stochastic gradient descent.

### 2.1.1.2 Non-PCA

There are other approaches for reducing or removing correlations in representations. For example, [11] added a "cross-covariance penalty" to the loss function of an autoencoder. This extra penalty term seeks to directly minimize the pairwise correlations between hidden units. In [15], the author used direct inhibitory ("anti-Hebbian") feedback connections between the representation factors, which has the effect of of decreasing correlations between factors.

A variational auto-encoder (VAE) [28] adds a the bias of Gaussian-distributed factors to an autoencoder by imposing a Gaussian prior on the representation, in which the covariance matrix is (typically) diagonal. This corresponds to assuming that the factors are Gaussian and uncorrelated/independent. Autoencoders are nonlinear models and thus do not share the linear bias of PCA, so the VAE appears to the left of PCA in the model chart in Figure 2.1.

### 2.1.2 Non-Gaussian Distribution of Factors

**Generative Assumptions.** The factors in many environments are non-Gaussian. In this section, I will show methods for finding factorial representations in non-Gaussian settings that also solve the uniqueness problem of PCA. First I will describe a theoretical motivation for why decorrelation is insufficient in non-Gaussian settings, which leads to a generic method for finding factorial representations called Independent Components Analysis or ICA. Then, I will describe a subset of non-Gaussian distributions, called sparse distributions, that are appropriate for most naturalistic environments. These distributions lead to a specific type of ICA called SPARSE ICA. I will then describe some Non-ICA methods for discovering non-Gaussian independent factors.

To find independent factors that are not Gaussian-distributed, it is essential to impose a stronger bias than decorrelation. Correlation between two random variables $\rho(x, y)$ depends only on the first moment of the joint PDF $p(x, y)$. However, independence is a generalized measure of non-correlation: two variables are independent if and only if $\rho(x^p, y^q) = 0$ for all positive integer values $p$ and $q$. For jointly Gaussian-distributed variables, there is no higher-order correlation structure, and so non-correlation is a sufficient condition for independence. However, the factors of variation in many naturalistic environments are non-Gaussian. Finding the true independent causes in these environments requires a stronger inductive bias.

### 2.1.2.1 ICA

**Inductive Bias.** An approach to finding truly independent factors is to find the most non-Gaussian factors that are also faithful to the input data. Intuitively, if a learned component is actually the sum of more than one of the true independent, non-Gaussian factors, that sum will look more Gaussian, according to the central limit theorem. Only the true non-Gaussian source factors will exhibit maximal non-Gaussianity. Maximization of non-Gaussianity has a connection to information theory as well: it is equivalent to maximizing the mutual information between the observed data and the learned representation [9, 21], and forms the basis of a group of models called Independent Components Analysis or ICA. There are many variants of ICA, including the original neural networks approach [19], information-theoretic approaches [9, 21], and statistical approaches [1]. ICA models, like PCA, generally assume a linear relationship between observation and representation, but there has been some research into developing non-linear ICA-like methods [50].

### 2.1.2.2 Sparse ICA

**Generative Assumptions.** Many factors found in naturalistic environments are sparsely distributed. For example, consider the set of photographs of objects. There are many possible objects, but most photographs will only display a small number of them, and any given object will only occur in a relatively small number of photographs. Sparsity can be interpreted in two ways — only a limited number of factors are present for any given observation, and for any factor, it is present in only relatively few observations. Of course, not all components in naturalistic images are sparse: For example, the mean luminance of images tends to follow a roughly uniform distribution, which is sub-Gaussian. However, sparsity is well-matched to most features of naturalistic images relevant to tasks such as object classification.

**Inductive Bias.** Environments with sparse factors offer a nice opportunity to employ a stronger independence bias. Sparse distributions are a type of non-Gaussian distribution called a *super-Gaussian* distribution. These distributions have higher kurtosis (fourth moment) than a Gaussian. Examples of super-Gaussian distributions include the Laplacian (used in L1/lasso regularization) and Cauchy distributions. Sparsity is also used in the field of sparse coding to find unique solutions to over-complete problems. We can transform PCA into ICA by adding a

sparsity penalty $s(\mathbf{z})$ to the loss function from Equation 3:

$$\mathcal{L}_{\text{ICA}} = \|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 + \lambda \; s(\mathbf{z}) \tag{4}$$

$\lambda$ is a parameter governing the loss trade-off between reconstruction and sparsity.

The sparsity penalty corresponds to minimizing the negative log of the probability of $\mathbf{z}$ under some super-Gaussian prior in a bayesian setting. For example, let's assume $s(\mathbf{z})$ corresponds to an L1 penalty, $\|\mathbf{z}\|_1$. The L1 penalty is equivalent to a super-Gaussian Laplacian prior. First, we assume independence of the factors, so the prior $p(\mathbf{z})$ can be factorized as:

$$p(\mathbf{z}) = \prod_i p(\mathbf{z_i}) \tag{5}$$

The Laplacian prior probability density is typically defined as[2]:

$$p_{\text{Laplace}}(\mathbf{z}) \propto \prod_i \exp -\lambda |\mathbf{z}_i| \tag{6}$$

This laplace prior is centered/peaked at zero (if it were not, then the resulting distribution would not be sparse), and $\lambda$ represents a shape parameter and is equivalent to the parameter governing the loss tradeoff. For the purpose of optimization, we often minimize the negative logarithm of Equation 6, which leads us to the formula for $s(\mathbf{z})$:

$$\lambda \; s(\mathbf{z}) = \sum_i \lambda |\mathbf{z}_i| = \lambda \; \|\mathbf{z}\|_1 \tag{7}$$

The formulation of ICA in Equation 4 is called Reconstruction ICA [34]. The reconstruction loss in Equation 4 can also be replaced with a restriction that the representation be fully invertible, and helps to prevent degenerate solutions. Of course not all ICA methods use sparse factors — some methods explicitly seek sub-Gaussian factors, or some other non-Gaussian characteristic. Likewise, not all methods that employ sparsity can be reduced to ICA; for example, sparse autoencoders are not ICA because they allow for non-linear relationships between observation and representation. In the next section, I will describe some of these non-ICA methods for non-Gaussian factors.
**Representation Structure.** The representation of ICA is exactly like that of PCA — each observation is mapped to a vector of factor coefficients.

### 2.1.2.3 Non-ICA

There are a few other methods for finding non-Gaussian independent factors that are not ICA, in that they do not explicitly maximize non-Gaussianity of factors. Like SPARSE ICA, these algorithms all make some assumption about the environment; for example, they assume that the factors are binary, or make some assumption about the generative process of the environment.

We can incorporate the bias of a sparse representation into an autoencoder architecture by adding a sparsity penalty [37] to get a sparse autoencoder.

In [46], binary factorial representations are learned. In this paper, the reconstruction error is minimized alongside an additional penalty term: the predictability of a factor given all the other factors. This penalty term is made tractable by the assumption that all representations have binary coefficients. This penalty works in non-Gaussian settings because it looks at dependencies beyond correlations.

A cooperative vector quantizer [63, 16] is composed of several vector quantizer units, each of which make a weighted contribution to the output. The vector quantizer can be viewed as implementing sparsity — only one unit within the VQ can be active at any time. The combination of outputs from several vector quantizers models an environment in which several groups of sparse factors contribute to each observation.

In [45], the author found that factorial representations can be discovered in non-sparse settings by designing the mapping function to discourage over-cooperation between factors, and to closely reflect the generative structure of the environment.

Austerweil and Griffiths [3] investigate how humans might select features to use for category learning. They use an Indian Buffet Process prior, which corresponds to an assumption that humans will learn independent features — features that covary will be grouped together. They find that their model predicts human categorization behavior in experiments.

In the next section, I will compare the bias of Gaussian-distributed factors with the bias of sparsely distributed factors on a few example sparse datasets.

---

[2]For simplicity, I have omitted the scaling factor

### 2.1.3   Comparison of Gaussian and Sparsity Bias in a Sparse Environment

The relationship of PCA and ICA can be seen in the model organization chart in Figure 2.1. Both PCA and ICA have fairly high inductive bias due to their linear mapping assumption. The difference in inductive bias between the two models is obvious when comparing equations 3 and 4 — ICA has an extra term (sparsity) in the objective.

A problem with PCA/decorrelation bias is that is not well-matched to environments where the independent components are sparse. An illustration of the problem can be found in [45], which I have recreated here in Figure 2.2 on the left. The synthetic dataset consists of $(11 \times 11)$-dimensional vectors of numbers, where each number is either on (1, white) or off (0, black). One factor controls the position of the black rectangle on the left, and a separate, independent, factor controls the position of the black rectangle on the right. The factors in this synthetic environment are sparse — only one rectangle on each side is active at a time. Also, the dataset is fully factorial, showing all combinations of positions once each. I trained both a PCA and ICA model on this dataset. At the top right in Figure 2.2, we see the mean vector (left) and components learned by PCA — the rows of the $\mathbf{W^{-1}}$ matrix. In this visualization, the pixels are scaled so that grey represents a neutral (zero) weight, black represents a negative coefficient, and white represents a positive coefficient. The basis learned by PCA is compact and faithful — only four factors (plus the mean vector) can explain all the variation in the dataset. However, the representation is not interpretable — the generative process that synthesized the dataset is not at all apparent in the representation. In other words, the bias of decorrelation is insufficient in explaining the true generative process of the environment. At the bottom right in Figure 2.2, we see an ICA decomposition of the same dataset[3]. The components shown are again the rows of the $\mathbf{W^{-1}}$ matrix, which in the parlance of ICA is called the "mixing matrix". The four ICA components reflect the causal structure of this environment better — each component only describes activity on one side of the image. The ICA representation is just as compact and faithful as the PCA representation since all original examples can be reconstructed from the four components.
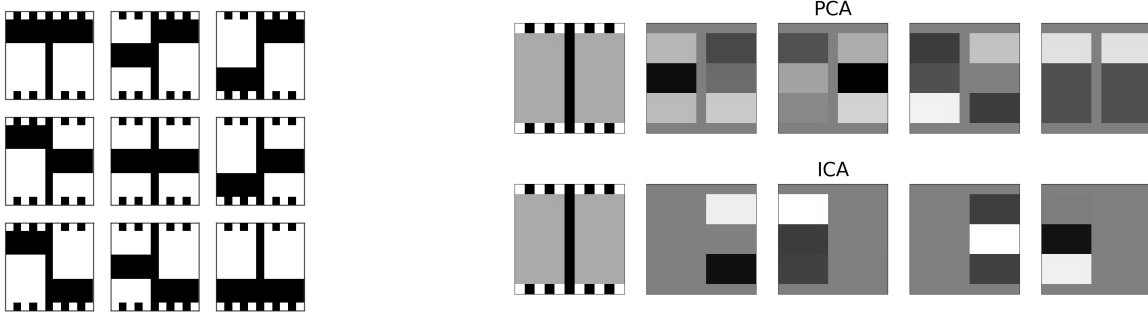


Figure 2.2: (left) A synthetic dataset like the one from [45]. Nine 121-dimensional test data samples, where the positions of the black rectangles on the left and right sides is controlled by two independent processes. (top right) The mean vector and four components learned by a PCA model. (bottom right) The mean vector and four components learned by an ICA model.

Another example where ICA improves interpretability over PCA is shown in Figure 2.3, using the naturalistic face dataset from the introduction. Once again, I trained both a PCA and ICA model on the dataset. On the left are the 24 top components (again, rows of $\mathbf{W^{-1}}$, scaled in the same way as Figure 2.2) learned by PCA, and on the right 24 components learned by ICA. In this dataset, there are several independent sources of variation including identity, lighting conditions, and glasses style. Subjectively, we can see that each PCA component seems to highlight a mixture of person identity, glasses style, and lighting conditions. The ICA components tend to focus on one cause — if glasses are emphasized then identity and lighting conditions are de-emphasized. It is also notable that unlike PCA, ICA does not produce an ordered list of components in terms of how much variance they explain. Note that the ICA components do not isolate the causal factors perfectly — ICA is still limited by the inductive bias of the linear mapping between observation and representation, so it cannot disentangle glasses from the rest of the face.

A less constrained environment is shown in Figure 2.4, which shows samples from a dataset of nature images from [23]. Each image is split into $32 \times 32$ patches, and the patches were used to train PCA and ICA models. In this environment, each image patch can be interpreted as having been generated by a number of localized factors that operate independently at different scales — the crest of a wave is generated independently from the overall lighting of the scene. In the center is the PCA basis and on the right is the ICA basis, with components scaled the same way as the other figures. These images are all from [23]. The PCA components all code for global features at different

---

[3]Using the FastICA[21] algorithm

Figure 2.3: (left) The top 24 components for the glasses data set. (right) 24 ICA components for the glasses dataset.

scales. The ICA components also seem to operate at different scales, but are much more localized, better matching the generative process of the environment.
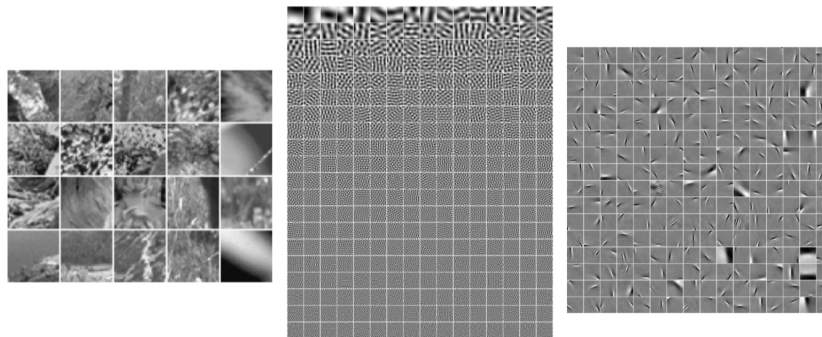


Figure 2.4: (left) Samples from a data set of nature images. Results are from [23]. (center) A visualization of the PCA basis of this dataset. (right) A visualization of an ICA basis for this dataset.

I have shown how the bias of factorial representations can help learn better representations given a factorial environment. For sparse environments observed in naturalistic images, decorrelation is an insufficient bias to learn factorial representations. The stronger bias of ICA with sparsity is shown to improve representations. ICA solutions are also unique, which make them more interpretable than non-unique PCA solutions. Next I will describe a bias that assumes factor representations should be insensitive to variations in their expression that are unrelated to the factor.

## 2.2 Invariance Bias

Factors can vary in how they are expressed. Objects can appear at various locations in an image, or be rotated and scaled. Linear methods like PCA and ICA are unable to express this type of variation, and tend to learn redundant representations of factors to describe all the ways in which a factor can be realized. In these environments, we would prefer a representation of a factor to be *invariant* to random fluctuations in how that factor is expressed.

**Generative Assumptions.** Many factors can be expressed in the input in more than one way. In a nature scene, animals might appear in different parts of the image. It would be nice to be able to represent the presence of the animal as one coherent factor, rather than have a code for each animal in each location. Likewise, a texture pattern in an image might show up in various phases, or under rotations. In the introductory example of shapes and colors, the "shape" factor is invariant to color, and vice versa.

I use the term *view* to refer to one of the possible realizations of a factor — e.g. one location, rotation angle, or one translation. Views of a factor can *compete* or *cooperate* to generate an observation. For example, there is competition between locations for an animal to appear in the nature scene: the locations are mutually exclusive. However, views could exist at only the extremes of locations, and observations would interpolate between these views.

In this case, the views cooperate, and are not mutually exclusive.

I will first show an inductive bias that that is consistent with an environment where factors compete, and views within a factor cooperate. Then I will show an inductive bias consistent with an environment where factors cooperate and views within a factor compete.

### 2.2.1 ISA

**Inductive Bias.** A method for learning invariant factors, where views cooperate is to *pool* similar views together. In the case of ISA, a view corresponds to a linear component, of the same dimensionality as the input, just like in PCA and ICA. The magnitude of all the coefficients in the pool is the coefficient on the factor. Sparsity is then applied at the pooled layer, and has the effect of making the pools compete with each other. The subspace formed by the pool represents all the views of the factor. Independent Subspace Analysis (ISA)[22] is an extension of ICA that implements this pooling method. In the generative model of ISA, we first choose some small number of factors to be present and then generate the observation by linearly combining elements from within the pools representing the active factors.

In this model, there are two layers — the first is a linear layer, and the second pools the linear components together to form the final representation. The first layer of N coefficients $\mathbf{z}$ is, like ICA, linearly related to the input $\mathbf{x}$:

$$\mathbf{z} = \mathbf{x}\mathbf{W} \tag{8}$$

A grouping function $g(\mathbf{z})$ pools together some subset of $\mathbf{z}$, in this case by their Euclidean norm:

$$g(\mathbf{z}_{...}) = \sqrt{\sum_i \mathbf{z}_i^2} \tag{9}$$

The second layer, the set of pooled components, is called $\hat{\mathbf{z}}$. In this case, $\hat{\mathbf{z}}$ is composed of $N/D$ pools of size $D$, but the size and number of pools is arbitrary:

$$\hat{\mathbf{z}} = [g(\mathbf{z}_{1..D}), g(\mathbf{z}_{D+1..2D}), \ ... \ , g(\mathbf{z}_{(N-1)D..ND})] \tag{10}$$

The loss function of ISA is identical to that of ICA, except that the sparsity loss $s(\cdot)$ is applied to $\hat{\mathbf{z}}$ instead of $\mathbf{z}$:

$$\mathcal{L}_{\text{ISA}} = \|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 + \lambda \ s(\hat{\mathbf{z}}) \tag{11}$$

Again, $\lambda$ governs the trade-off between reconstruction loss and the sparsity penalty. The grouping structure of ISA is illustrated in Figure 2.5.
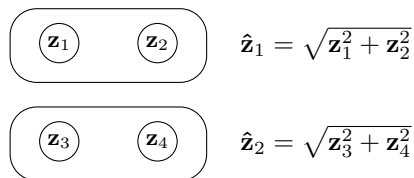


Figure 2.5: Illustration of the grouping structure of ISA. The circles represent the coefficients z.

Sparsity is applied at the pooled layer $\hat{\mathbf{z}}$, so only the pools, not the individual components within a pool, are encouraged to be sparse. This has the side-effect of encouraging the components of each pool to be *more dependent* on each other, and increase cooperation. Intuitively, if the components of a pool are independent, then their sum in Equation 9 should be more Gaussian, according to the central limit theorem. The sparsity objective eliminates this possibility, and so encourages the components in the pool to be more dependent.

**Representation Structure.** ISA has a more complex structure than ICA, in that it has two layers. Each factor in the second layer $\hat{\mathbf{z}}$ is accompanied by a vector of coefficients in its corresponding *subspace* in $\mathbf{z}$. The representation is therefore composed of both a scalar coefficient in $\hat{\mathbf{z}}$ as well as the coordinate in its corresponding z subspace.

**Examples.** A visualization of some ISA pools, from [23], is shown in Figure 2.6. The authors applied ISA to the same dataset of nature images that they used for ICA in Figure 2.4. A visualization of the learned components inside the $\mathbf{W}^{-1}$ matrix is shown in Figure 2.6. In this case, the number of components per pool was set to four, and the pools are ordered so that the sparsest ones, according to $s(\cdot)$, are at the top. The analysis in [23] indicates that the pools are less sensitive to phase than features learned by ICA, meaning that in this case, the factors in ISA have
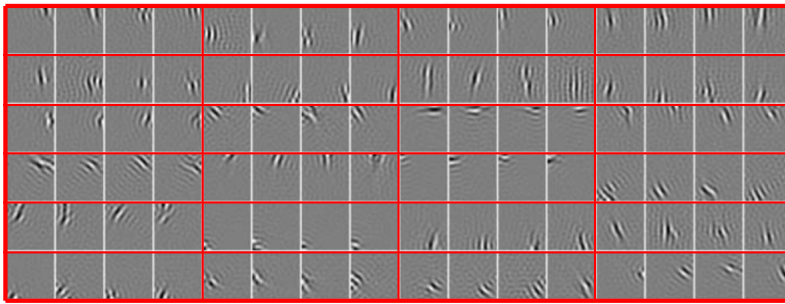
Figure 2.6: Components learned by independent subspace analysis for the dataset of $32\times32$ patches of naturalistic images from [23]. Every four consecutive components in a row corresponds to one pooled component. Positive weights are shown as light regions, and negative weights are shown as dark, while weights close to zero are gray.

more phase-invariance. In an ICA solution, the components are highly localized in space and operate at different scales (e.g. Figure 2.4). In ISA, each pool's components are at the same scale, localized within the same region of the image, and seem to be oriented in the same way.

There are, of course, many other types of inductive bias for invariant features, especially for image processing, where there are known classes of variation. Models often have hard-wired invariance for these types of variations. These biases are helpful for learning in their respective domains, and are used in many of the models in later sections. However, they are domain-specific adaptations: they are more concerned with specifics of the *environment* than with the *representation*. Therefore, their purpose is largely orthogonal to the purpose of this work, which is to describe how to learn factorial representations in any environment.

The inductive bias of ISA is weaker than ICA because it introduces a non-linearity into the mapping function, and reduces the number of components that sparsity is applied to. ISA can be reduced to ICA by setting the number of groups equal to the number of first-layer components. Therefore, it appears to the left of ICA on the model chart in Figure 2.1.

Next, I will demonstrate an inductive bias that is appropriate for an environment in which factors cooperate with each other and components within a factor compete.

**Other Methods With Cooperative Pools.** Factor models with cooperative pools are used in a variety of contexts. Pooling can be used in an overlapping spatial grid to create a topographic map representation, as in Kohonen's adaptive-self organizing maps [30]. Topographic maps lay out the first level of the latent representation as a grid. Pooling happens in overlapping sub-regions of this spatial grid. This encourages each sub-region to be more similar, so the filters in each grid cell have spatial coherency. A topographic map model incorporating sparsity can be found in [27].

### 2.2.2 Cooperative Vector Quantizer

**Inductive Bias.** In some environments, views compete and factors cooperate. In [63], each factor is represented by a Vector Quantizer (VQ), which corresponds to a pool of views that compete to explain the observation. An array of VQs cooperate to generate the output. Each VQ is composed of a pool of linear components, just like in ISA. This architecture is called a Cooperative Vector Quantizer (CVQ). Unlike ISA, in the CVQ, the components in the pool are mutually exclusive: only one of them may be active at a time. Additionally, there is no competition imposed between the factors. They are all assumed to be present all the time.

The inductive bias of CVQ is the same as ISA. The difference is only in whether sparsity/competition is applied within each pool or on the pools themselves.

**Representation Structure.** The representation structure of CVQ is identical to that of ISA: each factor is represented by a pool of components, each of which is linearly related to the input.

**Examples.** In [16], a CVQ is trained on a dataset of overlapping shapes, reproduced here in Figure 2.7. In the dataset, each $6 \times 6$ black-and-white image is generated by combining three white shapes: a cross, a diagonal line, and an empty square, in a number of different locations. Most examples have all three shapes, although in some the shapes are completely overlapping. Therefore, this example has three cooperating factors, and each factor is realized in only one component. The CVQ was trained with three pools/VQs. The components of both pools are shown on the right in Figure 2.7. The first pool clearly represents the diagonal line, in various locations in the image. The second represents the empty square, and the last one represents the cross. In this example, the model has learned to represent shape in a way that is invariant to location.
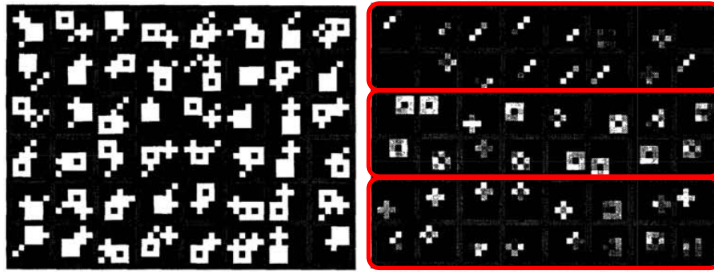
Figure 2.7: (left) Samples from the shapes dataset used for training in [16]. (right) A representation learned by CVQ, also from [16]. Each red box corresponds to one pool of a VQ.

### 2.2.3 Other Inductive Biases for Invariance

In [4], a transformed Indian Buffet Process prior is applied to

Some models encourage representations to be invariant to minor perturbations in the input. A contractive regularization penalty can be used to encourage factors to be insensitive to local directions of variation in the observation space [42]. Denoising autoencoders [57] also encourage representations to invariant to noise added to the observation.

Next, I will describe a setting in which factors can vary in how they are expressed, but it is the interaction between factors that is responsible for generating observations.

## 2.3 Combination Bias

I will highlight three different modeling biases that determine how factors can interact in generating observations. The simplest and strongest inductive bias for factor combination is the *linear* bias. In a linear model, observations are generated through a weighted summation of factor components. The coefficients can be either negative or positive, so a factor can cancel the contributions of other factors. PCA and ICA are both examples of this linear bias. First, I will discuss a bias that allows for *multiplicative interactions* to occur between a small set of factors. Environments where factors are mutually exclusive in explaining portions of the observation are called *functional parts*; I will describe some biases that work well in these environments. Finally, some environments are composed of a *hierarchy of layers* of factors, which generate observations at increasing levels of detail and decreasing levels of abstraction.

### 2.3.1 Multilinear

**Generative Assumptions.** In some environments, a small number of factors are known to *always* be present, and their interactions are responsible for generating observations. Each factor has a number of views, and view interacts with the views of the other factors. For example, the face dataset from the introductory section can be expressed as the interaction of person identity and glasses style — the interactions between these two factors can account of a large portion of the variation in images. A person's identity modulates where the glasses sit on the face, and the glasses style modulates the shadow cast on the person's face.

**Inductive Bias.** These environments are the motivation behind the general class of MULTILINEAR models, in which the multiplicative interactions between factors account for the variability of the observations. First, I will describe the simplest case of a MULTILINEAR model, called a *bilinear* model, in which there are only two factors. I will describe how MULTILINEAR models can be combined with an ICA-like sparsity penalty. Finally, I will show an example with three interacting factors.

#### 2.3.1.1 Bilinear Models

In a bilinear model [53], two factors interact with each other to generate each observation. Like ICA and CVQ, each factor in a bilinear model is represented by a pool of coefficients. Each coefficient of a factor interacts multiplicatively with the coefficients from the other factor. Each combination of coefficients is associated with a component vector, of the same dimensionality as the input. The weight on a component is determined by the product of the two coefficients associated with it, and the observation is generated by a summation of these weighted components. In [53], these factors are named "content" and "style", but they can be any two factors. For example, consider the environment of upright characters and italicized characters. The appearance of a character is completely determined by the interaction of two factors: the identity of the character, and its style.

Each factor in a bilinear model is represented by a "pooled" structure like in ISA. Suppose we are modeling two factors, a and b, with representations $\mathbf{z^a}, \mathbf{z^b}$. There are $M$ components in the pool $\mathbf{z^a}$, and $N$ components in pool $\mathbf{z^b}$. The weight matrix $\mathbf{W}$ is three-dimensional ($N \times M \times D$), and each component of W is weighted by the interaction between $\mathbf{z^a}$ and $\mathbf{z^b}$:

$$\tilde{\mathbf{x}} = \sum_i^M \sum_j^N \mathbf{W}_{ij} \mathbf{z}_i^{\mathbf{a}} \mathbf{z}_j^{\mathbf{b}} \tag{12}$$

As in PCA, bilinear models minimize reconstruction loss. We can interpret this model as a PCA model where the representation $\mathbf{Z}$ is a matrix that is factored into the interactions between the two vectors $\mathbf{z^a}$ and $\mathbf{z^b}$: $\mathbf{Z} = \mathbf{z^a}(\mathbf{z^b})^{\mathrm{T}}$. This model is essentially conjunctive: for a given coefficient $\mathbf{Z}_{ij}$ to be non-zero, both $\mathbf{z}_i^{\mathbf{a}}$ and $\mathbf{z}_j^{\mathbf{b}}$ have to be non-zero.

The bilinear model encourages similarity in the rows and columns of the component matrix $\mathbf{W}$. To help understand this effect, consider the diagram in Figure 2.8. The coefficient $\mathbf{z}_1^{\mathbf{a}}$, acts as a "gate" on the components $\mathbf{W}_{11}$ and $\mathbf{W}_{21}$. If the $\mathbf{z}_1^{\mathbf{a}}$ coefficient is zero, then the entire column is turned off. The gating property of $\mathbf{z}_1^{\mathbf{a}}$ encourages $\mathbf{W}_{11}$ and $\mathbf{W}_{21}$ to both be active for the same inputs, and thus more correlated and more similar to each other. In this way, the bilinear model encourages similarity in both the rows and columns of Figure 2.8.

This model is has some characteristics that are similar to CVQ: each factor is represented by a pool of linear components, and each factor is assumed to be present.

**Representation Structure.** MULTILINEAR representations use pools to represent factors, like in ICA and CVQ. However, unlike CVQ and ISA, the components associated with a pool are encouraged be consistent from pool to pool.
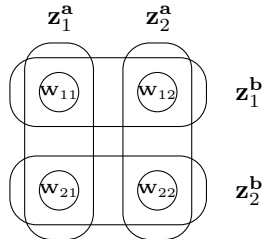


Figure 2.8: Illustration of the grouping structure of a bilinear representation with two factors a and b, each composed of two sub-factors. The circles represent the components $\mathbf{W}$. The weight on a component is determined by the product of the $\mathbf{z^a}$ and $\mathbf{z^b}$ coefficients of the groups it is a part of.

**Examples of Bilinear Models.** A simple example of a bilinear model can be found in [17], which also includes a sparsity constraint on $\mathbf{a}$ and $\mathbf{b}$. Sparsity here means that, for each factor pool, we sample only a small number of components to be present for any observation. The sparse bilinear loss function looks like:

$$\mathcal{L}_{\mathrm{SPARSE-BILINEAR}} = \|\mathbf{x} - \tilde{\mathbf{x}}\|^2 + \lambda_1 s(\mathbf{z^a}) + \lambda_2 s(\mathbf{z^b}) \tag{13}$$

The $\lambda_1$ and $\lambda_2$ settings govern the tradeoff between reconstruction loss and the sparsity penalties on a and b.

In [17], the authors compared ICA with a sparse bilinear model. The dataset used for training is generated by randomly selecting $10 \times 10$ image patch locations, drawn from $512 \times 512$ naturalistic images. For each image patch location, a set of image patches is generated by shifting the patch location $\pm 2$ pixels horizontally. The authors construct a sparse bilinear model, where the image patch feature is represented by a pool which they name $\mathbf{x}$, and the translation is modeled by another pool named $\mathbf{y}$. The model was trained using a custom training procedure that encouraged $\mathbf{x}$ to represent image patch features, and $\mathbf{y}$ to represent translations. This training procedure also used some *supervised* bias, utilizing knowledge of the relationship between an image feature and its translations. Supervised bias will be covered thoroughly in a later section. The authors also trained an ICA model on this dataset.

On the right of Figure 2.9, in the box labeled "Example bilinear basis", we can see some sample components (from the $\mathbf{W}$ matrix) corresponding to these two factors. The rows correspond to the $\mathbf{x}$ factor (image patch feature), and the columns correspond to the $\mathbf{y}$ factor (horizontal translations). The components in each row are consistent: they represent the same image patch feature, at different translations. The components in each column are consistent: they represent the image patch feature at a given translation.

On the left of Figure 2.9, in the box labeled "Example linear basis", two components of the ICA model are shown that are most similar to the $\mathbf{x}$ components of the examples chosen for the bilinear model. ICA learns to represent similar components, but would redundantly represent different translations of an image patch feature rather than grouping them together.

Figure 2.10 shows the response of $\mathbf{x}$ and $\mathbf{y}$ for a sample image patch, labeled "Canonical patch". Note that both $\mathbf{x}$ and $\mathbf{y}$ are sparse — very few coefficients are significantly nonzero. According to the model, the canonical patch is a mixture of $\mathbf{x}_3$ and $\mathbf{x}_9$ and $\mathbf{y}_1$. A translated version of the canonical patch is also shown, labeled "Translated patch". The $\mathbf{y}$ vector of this patch is different from the canonical patch — instead of activating $\mathbf{y}_1$ it activates $\mathbf{y}_{1,2,7,8}$, representing a translation.
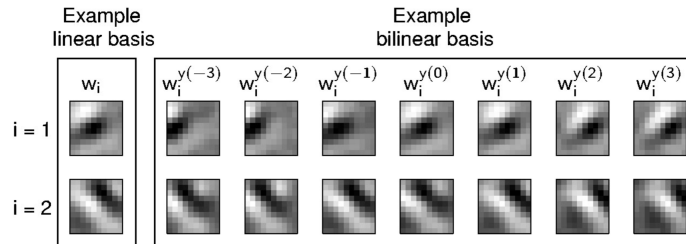


Figure 2.9: Visualization of an example linear model, on the left, and a corresponding bilinear model, on the right, from [17]. Negative component pixels are represented as dark regions, and positive pixels as light regions.
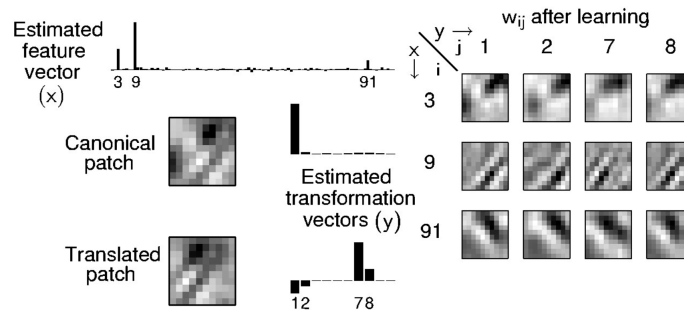


Figure 2.10: A sample component learned by a bilinear model from [17]. See text for a detailed explanation.

**Other Bilinear Models.** Not all bilinear models assume a sparse environment, or learn to represent translations. For example, Tensor Analyzers [51] are bilinear models that assume Gaussian-distributed factors, and were applied to model lighting conditions and identity on a face dataset.

### 2.3.1.2 Trilinear Models

An example of a *trilinear* model, with 3-way multiplicative interactions, is shown in Figure 2.11. This example comes from [13]. The model is trained on the Toronto Face Dataset, composed of grayscale $48 \times 48$ images. The two major factors of variation in this dataset are identity and emotion. The authors also model a third factor, which I will call "group". The third factor allows the model to represent local consistency without requiring global consistency. Identity and expression are consistent within a group, but are not necessarily consistent outside of the group. This makes the modeling easier: we do not have to enforce global consistency of factors. The groups work on $5 \times 5$ blocks of bilinear components. In this case, the model is trained in an unsupervised fashion.

The factors themselves are what is known as "spike-and-slab" variables: each variable has a binary "spike" sub-variable, indicating whether the variable is on, and a continuous "slab" sub-variable, indicating the weight of the factor. A the coefficient for a component is only nonzero when all three "spike" variables are 1; in this way, the value of a factor acts as a "gate" on the other factors. In each block, the rows represent identity and the columns represent emotion. Additionally, the blocks provide local consistency: the identity in the first row of the left-hand block does not correspond to the identity of the first row of the right-hand block. A similar MULTILINEAR model is used in a supervised setting [40] to separate face identity, emotion, and pose factors.

MULTILINEAR models are a good choice for environments where a few known factors are known to exist, and are known to have strong interactions. In terms of inductive bias, MULTILINEAR models are weaker than ISA. They allow for multiplicative interactions between coefficients, whereas ISA only allows linear combinations. This is reflected in the model chart in Figure 2.1. A sparse bilinear model can be reduced to ISA by setting the all weights of one factor pool to 1. It can be reduced to ICA by reducing the model to one factor.
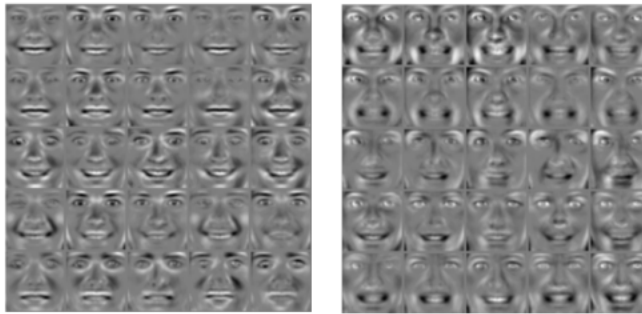
Figure 2.11: Two sample blocks of components of a tri-linear model from [13].

Next I will describe a type of environment where factors represent non-overlapping parts of the observation.

### 2.3.2 Functional Parts

**Generative Assumptions.** In some environments, each factor is responsible for generating only a part of a given observation. For example, a face can be broken down into a number of independent parts: eyes, eyebrows, mouth, nose, hair, etc. There is no overlap between any two parts, so there are no interactions between the factors. Models suited to these environments are often called "parts-based" or "functional parts" [52] models. In a linear environment, factor components are combined through summation to construct the output. In a functional parts environment, factor components are *concatenated* to construct the output.

I will describe two methods for learning functional parts models: first, by restricting coefficients on factors be only be positive, and second, by imposing competition between factors to explain parts of the input observation.

#### 2.3.2.1 Non-Negative Matrix Factorization (NMF)

**Inductive Bias.** One model capable of learning parts-based representations is non-negative matrix factorization (NMF) [35]. Like PCA, NMF reconstructs an observation by multiplying a weight matrix with a representation vector. However, both the weight matrix and representation vector are constrained to be non-negative, which corresponds to the following loss function:

$$\mathcal{L}_{\text{NMF}} = \|\mathbf{x} - \mathbf{z}\mathbf{W}^{-1}\|_2^2 \quad s.t. \ \mathbf{W}^{-1} \geq 0, \ \mathbf{z} \geq 0 \tag{14}$$

In this case, only positive coefficients, and positive component vectors are allowed. In PCA, we can generate a zero output for one of the observed dimensions by summing a positive component with a negative component, which cancel each other out. In NMF, this cooperation is not allowed: the only way to get a zero output is for both components to be zero. NMF eliminates a type of cooperation between factors, which encourages them to represent independent portions of the input.

This corresponds to the generative model — a functional-part factor cannot remove anything from another functional part; it can only contribute positively to the observation. This bias leads to highly localized components. **Representation Structure.** The structure of representations learned by NMF is the same as ICA — each factor's coefficient is one scalar value.
**Examples.** Figure 2.12 shows two example parts-based representations of faces from [35]. The representations were learned on a database of $19 \times 19$ face images. A visualization of the components (again, the rows of $\mathbf{W}^{-1}$) for each model are shown. Also shown are the z weights and reconstruction for an example face, marked "Original". While the reconstructions are about equally good for both models, the NMF representation is clearly parts-based — some factors correspond to eyes, noses, or a mouth. In PCA, as before, the factors are highly global.

The inductive bias of NMF is fairly high. It is still essentially linear, with the exception of the non-negativity constraint. It obviously has higher bias than PCA, since it is a subset of PCA, but it is unclear whether its bias is higher than ICA. It therefore resides at the same coordinate as ICA in the model chart in Figure 2.1.

#### 2.3.2.2 Other Models with Functional-Parts Bias

NMF is not the only means for discovering parts-based representations. In [54], a factorial parts-based model is learned by restricting each element of the observation to correspond to only one explanatory factor of a vector quantizer. This can be viewed as another application of the notion of sparsity: Instead of sparse factor coefficients
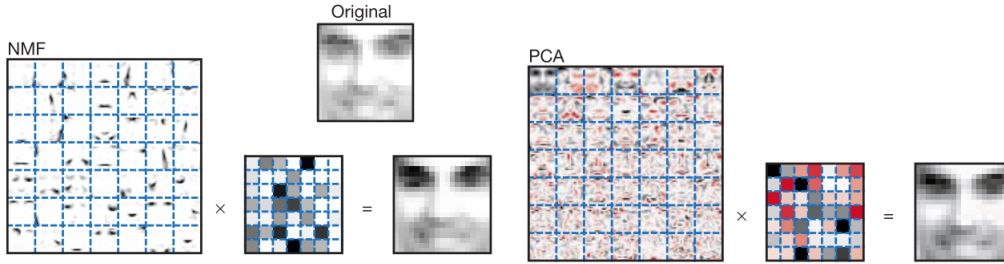
17

Figure 2.12

Figure 2.13: Visualization of a face representation learned by non-negative matrix factorization (NMF), compared to PCA. Positive component pixels are shown as dark regions, zero pixels are shown as white, and negative pixels shown as red. Note that the NMF representation has no negative components.

as we saw in ICA, sparsity is applied to the incoming weights from the factors to the observed feature, so that each feature in an observation is generated by only one factor. This basic approach is also used in [43].

This "functional parts" bias can easily lead to non-factorial representations if used in improper environments. In the face example, the face parts are, of course, not independent. The morphology of the face is highly influenced by gender and ethnicity, and of course correlations due to symmetry. However, there are tasks where we want to ignore these types of correlations and treat these parts as independent anyway.

Next, I will explore a type of bias that assumes that observations are generated by a sequence of decreasingly abstract layers.

### 2.3.3  Hierarchical Layers

**Generative Assumptions.**   Some environments have multiple levels of abstraction. For example, a nature scene might be composed of a foreground and background. First, we choose the locations of the foreground and background. The foreground is composed of a group of trees and rocks, which are sampled independently. We then choose a shape of each rock, which is independent of the shape of the other rocks. In this way, each layer represents a set of choices that are conditionally independent of each other given the higher level layer.

#### 2.3.3.1  Recursive ICA (R-ICA)

**Inductive Bias.**   Recursive ICA (or R-ICA) [48] is a multi-layered version of ICA, where the sparse ICA algorithm is applied iteratively until the desired number of layers is reached. The first representation layer of R-ICA, $\mathbf{z^1}$, is just the output of ICA. The first layer is then prepared for the next layer of ICA by transforming the outputs $\mathbf{z^1}$ by a function $t(\cdot)$. The purpose of $t(\cdot)$ is to reformat the highly-non-Gaussian coefficients of $\mathbf{z^1}$ to be Gaussian. ICA is then run on the Gaussianized first layer to form the second layer $\mathbf{z^2}$:

$$
\begin{aligned}
\mathbf{z^1} &= \mathbf{ICA}[\ \mathbf{x}\ ] \\
\mathbf{z^2} &= \mathbf{ICA}[\ t(\mathbf{z^1})\ ]
\end{aligned}
\tag{15}
$$

This procedure is repeated until the desired number of layers is reached. The details of the Gaussian transformation procedure can be found in [48], and involves discarding the sign of the representation, transforming it to a uniform distribution by applying the CDF of the prior, and then transforming that result back to a Gaussian by using the inverse CDF of the Gaussian distribution. The signs are discarded because they do not carry any redundancy in the signal, and removing them helps to reveal nonlinear dependencies.

**Representation Structure.**   The structure of R-ICA is multi-layered, and can be extended to an arbitrary number of layers.

**Examples.**   An example of R-ICA is shown in Figure 2.14. The algorithm was trained on a dataset of naturalistic images. A sample of Gabor-like ICA components is shown, along with responses from the second layer of ICA, arranged by both the position of the Gabor-like filter or the frequency and orientation. Some second-layer factors are strongly position-oriented, whereas others seem to be oriented towards patterns of frequency and orientation.
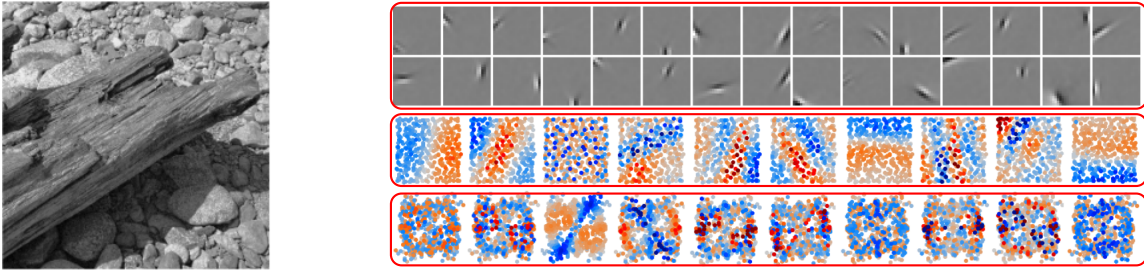
Figure 2.14: (left) A sample from the dataset of naturalistic images taken from [48]. (top right) Sample components from the first layer of an ICA solution, trained on $20 \times 20$ image patches. Each component is a gabor-like filter, with dark regions representing negative pixels and light regions positive. (center right) Weights from the first layer of ICA to a selection of the second-layer factors. Each square represent the weights from all first-layer factors to one second-layer factor. Each dot in a square represents the weight of one of the components of the first layer, where warmer colors represent a positive weight and cooler colors negative. The dots are arranged based on the location of the center of the filter. (bottom right) Factor coefficients from the second layer of ICA. The dots in this panel arranged by the frequency and orientation of the fitted Gabor filter, in polar coordinates.

### 2.3.3.2 Other Models with Hierarchical Layers Bias

A similar approach for learning higher-order statistical regularities using a hierarchical ICA framework is explored in [26]. ISA can also be made into a hierarchy [33], to form hierarchical layers of invariant factors. An R-ICA representation is less restrictive than PCA and ICA; it enables the discovery of more abstract factors. ISA is also a multi-layered method, and in fact we can construct ISA-like representations with an R-ICA model. Since R-ICA is a superset that includes ISA, it has lower inductive bias than ISA. However, R-ICA does not model any multiplicative interactions, we cannot say whether it has lower or higher inductive bias than a Multilinear model, and therefore appears in the same location on the chart in Figure 2.1. A layered or "deep" mixture of experts model [14] was found empirically to represent the location of an object in the first layer of experts, and the identity of an object in the second layer.

These models are all trained sequentially; the first layer is estimated, then its output is used as input into the next layer. A neural network is a variant of a hierarchical model in which the entire model is trained simultaneously.

I have so far described three different unsupervised inductive biases for factorial representation learning: distributional bias, invariance bias, and combination bias. I will now describe the another type of inductive bias, the use of a supervisory signal, which can be combined arbitrarily with the last three.

## 2.4 Supervision Bias

Supervisory data are another source of bias, in which we use additional information connected to observations to constrain the representation of a factor. For an identity classification task on a face dataset, the face identity would be a very useful factor to represent explicitly, and in a way that is not confounded with any other non-identity factor.

A supervised signal can make representation learning completely unnecessary: the representation is already available. However, the supervisory signal might not always be available, or it might be weak. We may need a model that works for new data that are not yet labeled. The supervisory signal may specify only one of many factors in the representation, and the others still need to be learned. Supervisory signal might be available for only a portion of the training data. Supervisory signals can also come in weaker forms that provide constraints on the representation, but do not specify it directly.

Supervised factors help make representations more interpretable, to the extent that the factors described by the supervisory data are interpretable. Supervision is also useful for learning factorial representations in the case of a non-factorial dataset (sampling bias) or environment.

Supervised constraints come in several different types. Which type to use depends on both the availability of data as well as on the type of representation sought.

### 2.4.1 Constraint Types

Various supervisory constraints are listed in Table 1. They are sorted by how strongly they constrain the solution, with the strongest constraints at the top. Specifically, I consider how big the constrained space of possible learned

solutions is, given the value for one of the observations. In this table, $n$ represents the dimensionality of the representation, and $z(a), z(b), z(c), z(d)$ stand in for representations of different observations $x(a), x(b), x(c), x(d)$.

Table 1: Types of constraints.

| Type | Example Constraint | Size of constrained space, given a |
|---|---|---|
| direct | $z(a) = \phi$ | - |
| equality | $\|z(a) - z(b)\| = 0$ | 1 |
| distance | $\|z(a) - z(b)\| = \delta$ | $S(z(a), \delta)$ |
| inequality | $\|z(a) - z(c)\| < \|z(a) - z(b)\|$ | $\mathbb{R}^n + B(z(a), \mathbb{R}^+)$ |
| analogy | $z(a) - z(b) = z(c) - z(d)$ | $\mathbb{R}^{2n}$ |

For each type of constraint, I will go into some detail into the types of tasks and representations that the constraint is well-matched to. I will also discuss the strength of the constraint, and the reasoning for its position in Table 1. Then I will present examples from the literature that illustrate the incorporation of these constraints into models.

### 2.4.1.1   Direct

Direct constraints are most commonly used for attributes of data that we can easily assign a value to — a binary attribute, categorical label, or some scalar quantity. For example, the class of an image of a handwritten digit is a categorical variable of a closed set, zero through nine. Direct constraints are the strongest constraints considered here, and explicitly state the setting for a factor.

### 2.4.1.2   Equality

The most basic non-direct constraint is an *equality*, which simply groups examples together that are equal with regard to a particular factor. Equality constraints can be simpler to generate than direct constraints — for example, we can ask humans if two face images are the same person or different identities, rather than asking them to label identity for each image in the dataset.

Equality constraints are useful for learning open-class representations of a categorical variable. For example, a face database might have an identity label associated with each image. If the task is to learn a representation useful for face verification of subjects that are not in the training database, it is not so useful to simply learn a categorical representation indicating which of the N people are represented in an image. Instead, we'd prefer a representation that could accommodate new identities. Instead of directly constraining the representation, we can constrain the representation to satisfy the equalities in the database.

An equality constraint is simply the specification that two examples are equal. This is a strong constraint, since given one observation, the other is constrained to occupy exactly the same point.

### 2.4.1.3   Distance

Another common constraint type is a *distance*, or, equivalently, a *similarity* constraint. Distance constraints are commonly used in one of three ways: to quantify some scalar difference between two objective factor values (e.g. difference in age between people), to specify an objective similarity between two objects (e.g. $L_2$ vector distance), or to specify a human's subjective similarity rating between two objects.

Objective similarities can come from differences between very high dimensional feature vectors. These are notoriously difficult to visualize, and so a representation that respects these distances is useful. For example, images tend to be very high dimensional, and it is difficult to visualize the relationships between many images simultaneously. Another common source of distance information are object co-occurrences — for example, we can rate word similarity based on how often they occur together in a sentence.

Subjective similarities come from humans. For example, we might want to create a visualization of a dataset of birds, designed to help us understand human judgments of bird similarity. This visualization could be a Cartesian space where each point represents a bird, and distance between two birds is reflective of human judgments of similarity. This visualization could help us understand what features humans attend to when making similarity judgments. In this case, a direct constraint would be a specific coordinate value for each bird. However, human judgment of image similarity is known to be complex, and is based on the simultaneous comparison of many features. Therefore, this labeling task is too difficult for humans. In this setting, a direct or equality constraint is not appropriate.

We could collect human judgments of similarity of pairs of bird pictures, and use these judgments to constrain our representation. This is the motivation for solutions like multidimensional scaling (MDS) [31]. MDS creates a

Cartesian space in which distance between the pair is reflective of similarity ratings from humans. Each pair of objects with a distance constrains the solution, without specifying exactly the coordinates for the pictures.

Many models dealing with distances are designed to learn one space that satisfies all constraints. However, this one-space assumption is not valid in the case of multiple factors which might contribute to a comparison. For example, consider the case of *intransitive similarities*. For example, a red bird is similar to a red fish with respect to color. The red bird is also similar to a blue bird with respect to animal type. However, the red fish is not similar to a blue bird. If we correctly factorize the space, such that distance constraints are assigned to color and animal type, it is easy to satisfy these constraints.

A *distance* constraint is almost as strong as an equality constraint. Given two observations x(a), x(b) and a distance between them $\delta$, the constraint on the representation is $\|z(a) - z(b)\| = \delta$. Given a representation of the first observation z(a), the other point z(b) is constrained to lie on the sphere around z(a), with radius equal to the distance $\delta$, or $S(z(a), \delta)$. This is reflected in Table 1.

#### 2.4.1.4 Inequality

Human judgments of distances are known to be unreliable — different users will employ different scales to judge objects, and sequential effects can cause drift in judgments based on the order of presentation of stimuli.

A solution to these issues is to collect an *inequality* constraint, instead of a distance. Inequality constraints are commonly used to collect human judgments of similarity, while avoiding the issues with direct similarity judgments. To acquire an inequality constraint, we present the human labeler with an anchor object x(a), and ask him which of two alternatives, x(b) or x(c), is most similar to the anchor object. If the user chooses x(c), then the resulting inequality is: $\text{sim}(x(a), x(c)) > \text{sim}(x(a), x(b))$. We can use this as a constraint on the representation z to encourage distance to represent (dis-)similarity: $\|z(a) - z(c)\| < \|z(a) - z(b)\|$. The constraint is insensitive to the magnitude of the differences, and is therefore insensitive to different scales or sequential effects.

Given one point, an inequality has two degrees of freedom. Given the point z(a), z(b) is free to lie anywhere in the space $\mathbb{R}^n$, and z(c) is constrained to lie inside the open ball around a with a positive radius in $\mathbb{R}^+$, $B(z(a), \mathbb{R}^+)$. Since each choice of z(b) maps to one open ball constraining z(c) rather than many, we add the sizes of the sets $\mathbb{R}^n$ and $B(z(a), \mathbb{R}^+)$ together. Since $\mathbb{R}^n > S(z(a), \delta)$, inequalities are a weaker constraint than distances, as reflected in its position in Table 1.

#### 2.4.1.5 Analogy

Another class of constraints are called *analogies*. Humans are adept at analogical reasoning, and are able to extrapolate answers to hypothetical questions such as "What if this face were viewed from a different angle?" or "What if this piece of music were played at a higher tempo?".

Analogy constraints take the form x(a) : x(b) :: x(c) : x(d), or, the relationship between observations x(a) and x(b) is the same as the relationship between x(c) and x(d). They are a means of capturing a relationship between observations.

Analogy constraints can be convenient in situations where direct constraints are not possible, but we have access to pairs of observations where we expect their relationship to be the same. For example, we might have images of two people before and after completing high school, or before and after completing college. We can interpret two of these pairs as an analogy constraint, and use them to learn age as a factor in a representation.

Analogies are typically interpreted as vector-difference constraints on the representation. For example, "x(a):x(b) :: x(c):x(d)" can be interpreted as $z(a) - z(b) = z(c) - z(d)$. Given the representation z(a), z(b) is free to lie anywhere in the space $\mathbb{R}^n$. Given z(a) and z(b), z(c) is still free to lie anywhere else in the space $\mathbb{R}^n$. However, given z(a), z(b), and z(c), we can determine the exact location of z(d), so it does not constrain the space further. So, the size of the constrained set is $\mathbb{R}^{2n}$. Since $\mathbb{R}^{2n} > \mathbb{R}^n$, we can see that analogies are a weaker constraint than inequalities, and are the lowest level of supervision in Table 1.

#### 2.4.1.6 Summary of Constraint Types

Direct constraints are useful for learning representations where factors have an unambiguous category or quantity. Equality constraints are almost as strong, are easier to obtain, and can be useful for learning open-class representations of a categorical variable. Objective measures of similarity are best modeled as distances, and human judgments of similarity are best modeled as inequalities. Analogies are the weakest constraint, and are applicable in situations where only the relationships among observations are available.

The best representation for a task might not correspond to the type of constraint available. For example, we might have direct constraints for face identity in the form of categorical labels. However, we might not want to

represent faces with a categorical variable. Instead, we might choose a representation where distance is related to the similarity of faces, so that we can understand which faces are most similar. In this case, the direct constraints can be transformed into equality constraints.

I will present models showing examples across the entire spectrum of constraint types. Each model shown will present some evidence for factorial representations. Before presenting these examples, I need to discuss how I will evaluate the factorial representations in the papers.

### 2.4.2 Evaluation of Representations

In a supervised context, we have the opportunity to directly measure some properties of factorial representations. There are three broad categories of evidence that demonstrate factorial representations:

- **Distillation.** Information relating to a supervised factor should be distilled into the factor representation. This distillation should make the information more easily accessible, as evidenced through results such as improved classification performance, or example reconstructions that traverse the factor manifold showing variation in that factor.

- **Disentangling.** Information not relating to a supervised factor should no longer be present in the factor's representation. We should be able to see a decrease in performance on a *non-target* task using the representation of a factor. We should also find evidence of *invariance* to other factors of variation.

- **Cross-Over Effects.** In the case of partially-labeled factors, disentangling should have the effect of reducing nuisance variation on even the non-supervised factors. We should therefore expect to see better factor distillation in even non-supervised factors, as evidenced by improved classification performance, or example reconstructions.

Next, I will describe example models where constraints are used to learn factorial representations. For each model, I will summarize the constraints, the method used for incorporating the constraint, and show examples that demonstrate evidence of factorial representations, where available. I also describe some modes that combine several different types of constraints.

### 2.4.3 Examples of Models

#### 2.4.3.1 Semi-Supervised VAE (Direct)

Direct constraints are most commonly applied by directly encouraging the representation to mirror the labels. For example, in [29], a variational auto-encoder (VAE) is modified to also include a categorical label for each data point, to form a Semi-Supervised VAE (SSVAE).

In the standard VAE [28] observations are generated by first choosing a representation z, from a zero-mean unit Gaussian distribution with diagonal covariance — corresponding to a bias of independent Gaussian factors. The observation is then generated through a non-linear function of the representation (a neural net), with parameters $\theta$.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|0, \mathbf{I}) \qquad p_\theta(\mathbf{x}|\mathbf{z}) = f(\mathbf{x}; \mathbf{z}, \theta) \tag{16}$$

In the SSVAE [29], we first generate a categorical label from the multinomial distribution $\mathrm{Cat}(y|\pi)$, with parameters $\pi$. We independently sample a $\mathbf{z}$. The observed data $\mathbf{x}$ are then generated conditionally on both $y$ and $\mathbf{z}$:

$$p(y) = \mathrm{Cat}(y|\pi) \qquad p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|0, \mathbf{I}) \qquad p_\theta(\mathbf{x}|y, \mathbf{z}) = f(\mathbf{x}; y, \mathbf{z}, \theta) \tag{17}$$

The SSVAE can account for *partially-labeled* datasets by marginalizing over $y$ for datapoints where no label is available. The representation is then composed of two factors, $y$, the class label, and z, the remainder of the representation. The optimization then follows a procedure similar to that of the VAE.

The authors trained an SSVAE on the standard MNIST dataset, which is composed of small black-and-white images of hand-written digits. The authors show evidence of factor **distillation**: SSVAE is more accurate than VAE in classifying the digit identity given an image. The authors also demonstrate **cross-over effects**: by representing digit identity using $y$, z is encouraged to represent *other* sources of variation. Figure 2.15 shows an example of this effect. The model is trained using a two-dimensional z vector. Each image block are samples from the model with $y$ clamped to represent the digit 2, 3, and 4. Within each block, the dimensions of z are varied smoothly. Nearby samples correspond to similar *styles* of handwriting: the left region corresponds to non-slanted styles, and the right region represents slanted styles.

Figure 2.15: Samples of handwriting styles from an SSVAE trained on the MNIST digit dataset.

**Similar Models.** Many other semi-supervised models use direct constraints in this fashion. For example, in [39], the authors add class labels to part of the representation of a ladder network, a form of deep denoising autoencoder. They report state-of-the-art performance on semi-supervised MNIST and CIFAR-10 tasks. There are a few examples of using direct constraints alongside with higher encoder inductive biases. For example, there are supervised versions of ICA [44], and NMF [58].

Direct constraints are the most straightforward constraints, and are appropriate for types of factors that are easily quantified, such as binary, categorical, or numerical variables. Next I will describe a weaker kind of constraint, where the exact value of a factor is not known.

### 2.4.3.2   Inverse Graphics Network (Equality)

The Deep Convolutional Inverse Graphics Network (DC-IGN) [32] uses equality constraints on a set of synthetically-generated face images. An autoencoder is trained to explicitly represent a number of factors, including azimuth angle of the face, elevation angle, and azimuth angle of the light source. The faces are synthetically generated, so the authors have a unique ability to put constraints on *all* of the factors of variation in the images.

The DC-IGN representation is split up to represent three factors explicitly, and to represent the remaining factors using a single pool of neurons. The first factor neuron represents the face azimuth, the second the face elevation, and the third represents the light source azimuth. The remaining pool of neurons represents other nuisance factors, including face identity, morphology, and expression.

The authors developed a custom auto-encoder training procedure that varies one factor at a time while holding all other factors of variation constant. The representations of the non-varying factors are encouraged by the learning procedure to remain constant. In the training, the pool representing nuisance factors is treated as one factor.

Figure 2.16 shows some sample reconstructions from the network, which shows evidence for factor **distillation** and **disentangling**. On the left, we see some sample input images and reconstructions. The remainder of the grid is formed by manipulating the representation of elevation, changing it smoothly over a range of values, and leaving the other factors unchanged. We can see in the reconstructions that face elevation is changing, but the face azimuth, light source, and identity remain constant. The authors demonstrate variation in the same way on the other factors as well. The authors did not explore cross-over effects.

On the right, the ground truth elevation of a sample of synthetic images is plotted against the learned representation of elevation. We can see that the two variables are highly correlated.
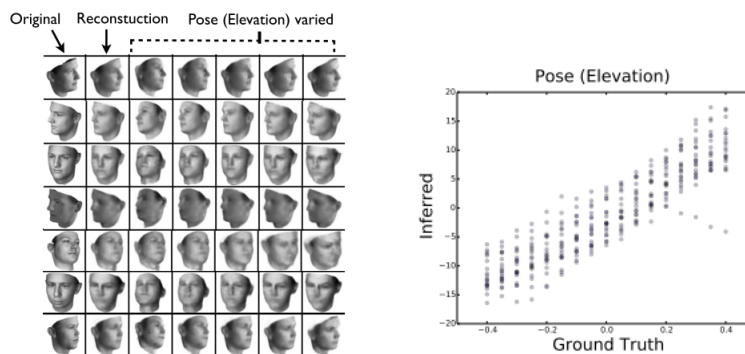


Figure 2.16: (left) Varying the elevation of random test samples. (right) A scatter plot showing the inferred elevation factor representation vs. the ground truth setting used to render the face image.

| model | task | phone embed. phn | phone embed. spk | speaker embed. phn | speaker embed. spk |
|---|---|---|---|---|---|
| MFSC7 | - | 24.5 | 32.9 | 24.5 | 32.9 |
| Sia7 | single | 10.9 | 46.0 | 46.4 | 23.9 |
| | double | 10.5 | 45.9 | 45.4 | 9.3 |
| Sia15 | single | 9.7 | 47.1 | 45.8 | 12.4 |
| | double | 10.2 | 46.6 | 45.3 | 8.7 |
| Tri7 | single | 10.0 | 46.0 | 45.0 | 10.0 |
| | double | 11.5 | 45.0 | 45.7 | 9.4 |
| Tri15 | single | 9.8 | 46.9 | 44.6 | 9.4 |
| | double | 10.7 | 46.2 | 44.7 | 8.1 |

Figure 2.17: Phonetic and Speaker ABX discrimination results from [61].

### 2.4.3.3   Siamese Network (Equality)

As noted, equality constraints are useful for learning open-set representations of categorical variables. This procedure is also called a *learned similarity metric*. The method is shown off in the Siamese network of [18], which learns a low-dimensional embedding of MNIST digits where similarity corresponds to likelihood of being in the same digit class. [12] used a Siamese network to build an embedding representation of face identity. Instead of the clamping procedure, Siamese networks are trained with a loss which minimizes distance for same-category examples and maximizes distance for out-of-category examples. Given two examples $\mathbf{x}(a), \mathbf{x}(b)$, whose factor values of factor f are f(a) and f(b), and their corresponding representations $\mathbf{z}^{f}(a), \mathbf{z}^{f}(b)$, the loss function in [18] is:

$$\mathcal{L}_{\text{SIAMESE}} = \begin{cases} \|\mathbf{z}^{f}(a) - \mathbf{z}^{f}(b)\|_2^2 & \text{if } f(a) = f(b) \\ \left[max(0, m - \|\mathbf{z}^{f}(a) - \mathbf{z}^{f}(b)\|_2^2\right] & \text{otherwise} \end{cases} \tag{18}$$

In [61], this technique is extended to represent multiple factors simultaneously as separate embedding spaces. The loss functions for the embeddings are simply summed together to form the final loss function.

The authors of [61] performed experiments using a speech dataset of consisting of different speakers speaking the same set of words. One embedding was trained to represent speaker identity, and anther was trained to represent the target phonemes. Instead of the euclidean distance differences in Equation 18, they used cosine similarity. The networks were trained on a subset of the Librispeech dataset, which consisted of 360 hours of read speech from 920 speakers. To test the model, the authors performed an ABX task to estimate classification accuracy. The error rates reported are the percentage of correct answers to the ABX tests. In the tests, the examples are triphones (a phoneme in the context of two other phonemes). To test phonetic discrimination, the A and B examples were from the same speaker, while the A and X examples matched on phonetic content but were pronounced by different speakers. This allows the testing to directly measure **disentanglement** as well as **distillation**.

Figure 2.17 shows a table of results from [61]. In the table, "Sia" stands for "Siamese" and "Tri" stands for "Triamese", or a network which was trained using a loss function that took both a positive and negative sample into account simultaneously. The number after the model name indicates how many stacked frames were included in the input. "MFSC7" is a baseline model that simply concatenates 7 MFCC frames together and uses that as the embedding. The "task" refers to whether the speaker and phonetic embeddings were trained separately ("single") or together ("double"). Performance is reported for both embeddings on both tasks, with lower error rates indicating better performance.

We can see strong evidence of **distillation**, as the trained embeddings learn to perform better at predicting their target factor. We also see evidence of **disentanglement**, as the trained embeddings perform worse on the non-target factor. We also see some indirect evidence for **cross-over effects**, since the "double" models tend to perform better than the "single" models, especially on speaker identification.

### 2.4.3.4   Multiple-Maps t-SNE (Distance)

A model that accounts for distances with multiple factors is *multiple-maps t-SNE* (MM T-SNE) [55], which is an extension of t-distributed Stochastic Neighbor Embedding (t-SNE) [36]. In t-SNE, distances between pairs of observations $i$ and $j$ are converted into probabilities; they form a joint distribution P over all pairs of objects where $p_{ij}$ is proportional to the similarity between objects $i$ and $j$.

The dataset used in [55] is composed of human-labeled word associations. Humans were given a prompt word and were asked to name associated words. 5,019 words were used as prompts, and 10,617 total words were labeled in associations. Each word and its associations are counted as one word co-occurrence. Co-occurrences were made symmetric and normalized to probabilities. In this case, the factors roughly correspond to topics, which can lead to intransitive similarities. For example, the word "tie" might co-occur frequently with "tuxuedo" when the topic is

24

"formal wear". The word "tie" might also co-occur frequently with "rope" when the topic is related to ropes and knots. However, "rope" will likely not co-occur with "tuxuedo". A factorial representation should separate these topics and allow the model to represent intransitive similarities.

The goal is to learn a word-representation $\mathbf{z}$ in which distances are reflective of the probabilities in P. To facilitate this goal, the learning procedure defines a probability distribution Q over distances in $\mathbf{z}$, and minimizes the KL-divergence $KL(P\|Q)$.

Q is chosen so that distances are proportional to a Student-t distribution with one degree of freedom[4]. The probability $q$ between words $\mathbf{x}(i)$ and $\mathbf{x}(j)$ with representations $\mathbf{z}(i)$, $\mathbf{z}(j)$ is defined as:

$$q_{ij} = \frac{(1 + \|\mathbf{z}(i) - \mathbf{z}(j)\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{z}(k) - \mathbf{z}(l)\|^2)^{-1}} \tag{19}$$

The top term means that the probability is inversely proportional to the distance between $\mathbf{z(i)}$ and $\mathbf{z(j)}$, and the bottom normalization term ensures that $q_{ij}$ behaves like a probability distribution.

This model is extended to include multiple factors. Each factor defines its own space in which all the words are mapped. For each word $i$ and factor $f$, the model learns an *importance weight* $\pi_f^i$. Importance weights are between zero and one, and sum to 1 over all the factors for each word. For a given combination of words $\mathbf{x}_i, \mathbf{x}_j$, the distance calculation is multiplied by the product of the two factor-specific importance weights. This has the effect of only penalizing the model for not respecting distances when both words have high importance weight for a factor. Q is updated to sum over all the factors:

$$q_{ij} = \frac{\sum_f \pi_f^i \pi_f^j (1 + \|\mathbf{z}_f(i) - \mathbf{z}_f(j)\|^2)^{-1}}{\sum_k \sum_{l \neq k} \sum_{f'} \pi_{(f')}^k \pi_{(f')}^l (1 + \|\mathbf{z}_{f'}(k) - \mathbf{z}_{f'}(l)\|^2)^{-1}} \tag{20}$$

Optimization finds the best settings for $z$ as well as $\pi$ such that the original co-occurrence probabilities in P are respected. In this way, the model automatically assigns words to factors — it does not depend on any prior knowledge of factors or associations between factors and similarity ratings, this makes the constraints in this model *factor-agnostic*. The number of factors must be known in advance.

Figure 2.18 shows some example maps from this model. The model was trained with 40 maps/factors, which are all two-dimensional. The dataset contains many topics, so there are typically multiple topics per factor. While every map contains every word, words with an importance weight below 0.1 are not shown. The top map corresponds to "sports" and "fashion", whereas the bottom map corresponds to a few topics, including criminal justice, kitchen equipment, and ropes/knots. We see that "tie" is placed near "tuxuedo" in the top map, and near "rope" in the bottom map. Some of the maps have some topic coherency, which is evidence for **distillation**. Not every map has every topic in it, which is evidence for a **disentangled** representation.

In the case of MM T-SNE, the fact that the constraints are factor-agnostic severely weakens the constraint space by multiplying the size of the space by the number of possible factors that a constraint could be assigned to, or $S(\mathbf{z}(a), \delta)$ times the number of factors/maps.

### 2.4.3.5 Karaletsos et al. 2015 (Inequality)

Triplet inequalities are used to learn a factorial representation for a variational autoencoder in [25]. An oracle (which could be human or machine) answers questions in the form of "is x(i) more similar to x(j) or x(l) in terms of factor $f$?". The answer to this question is a triplet $t_{i,j,l}^f$, interpreted as a factor-specific inequality: $\text{sim}_f(\text{x}(i), \text{x}(j)) > \text{sim}_f(\text{x}(i), \text{x}(l))$. The factors for each inequality are specified in the training signal. The number of factors and the triplets associated with each factor are specified in advance.

The probability of the inequality is modeled as a softmax function, with a factor-specific dis-similarity function D whose range only includes values greater than or equal to zero:

$$p(t_{i,j,l}^f) = Ber(t_{i,j,l}^f) = \frac{e^{-D_{i,j}^f}}{e^{-D_{i,j}^f} + e^{-D_{i,l}^f}} \tag{21}$$

The model discovers an $N$-dimensional embedding representation, which serves as the basis for the the $D_{a,b}^f$ function. In the paper, $D_{a,b}^f$ is defined in terms of the Jenson-Shannon divergences between each of the N dimensions of a and b, although in practice, an approximation to the true JS divergence is used.

---

[4]The choice of Q is important, because it determines which distances are attended to most during the learning. The t-distribution is chosen instead of a Gaussian because it allows points that are only slightly similar in high dimensional space to become far apart in the low dimensional representation, a solution for the crowding problem of high-dimensional spaces.
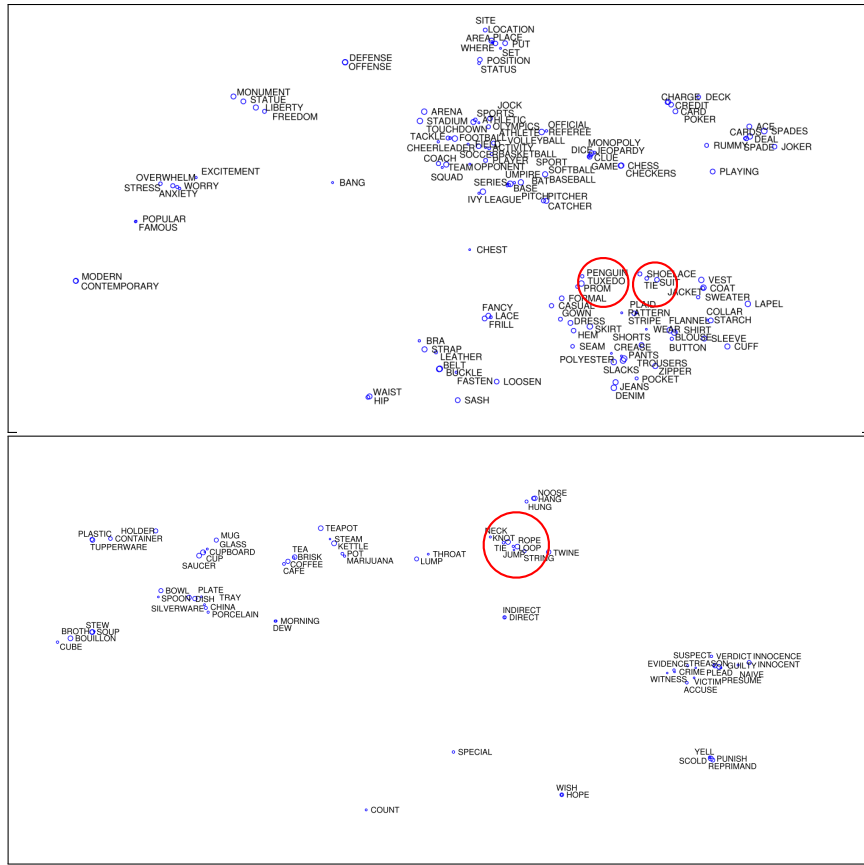
Figure 2.18: Visualization of two two-dimensional maps learned via MM T-SNE, trained on human-generated word associatation data.

The model learns which of the $N$ dimensions should contribute to explaining each inequality. A factor-specific *mask* selectively attends to dimensions of the representation that are relevant to the factor and hides irrelevant dimensions. Each factor-specific mask variable $\mathbf{m^f}$ is $N$-dimensional, and constrained to lie in the interval $[0,1]^N$. The distance function $D$ takes the mask into account when comparing the dimensions of the representation of two examples:

$$\overset{\mathbf{m^f}}{D_{a,b}} = \sum_{n=1}^{N} \mathbf{m_n^f} D_{a,b}^n \tag{22}$$

The optimization involves maximizing $p(t_{i,j,l})$ for all the factor-specific triplets, alongside the reconstruction error and representation regularization terms in the VAE. The masks are learned during the optimization. In this way, the model learns to switch on an appropriate subspace that represents factor $\mathbf{f}$ to explain the triplets drawn with respect to that factor.

The authors test their model on the Yale Faces dataset, composed of 2,414 images of 38 individuals under varying lighting conditions. The azimuth and elevation of the light sources are varied in the images. Each triplet is composed of three randomly selected images. The three factors of variation — identity, light source azimuth, and light source elevation — were used to create triplet inequalities out of the three randomly chosen images. For example, a light source azimuth triplet would provide information about which two images had more similar azimuths.

The authors demonstrate factor **distillation**: this masked model performs better on identity classification, azimuth prediction, and elevation prediction, compared to a model trained with no constraints. Figure 2.19 shows a visualization of the three factor-specific masks learned on the face dataset. The masks learn to factorize the space, such that only a subset of the dimensions of z contribute to decisions comparing identity or light source position. A dimension-reduced visualization of the dimensions in these masked subspaces is shown on the right in Figure 2.19. We can see that the subspaces group similar examples together, as would be expected.

**Similar Models.** A very similar approach is taken in [56], which demonstrated the model's ability to learn factorial representations on other datasets, such as stylized fonts and images of shoes. Factor-agnostic triplets are used to
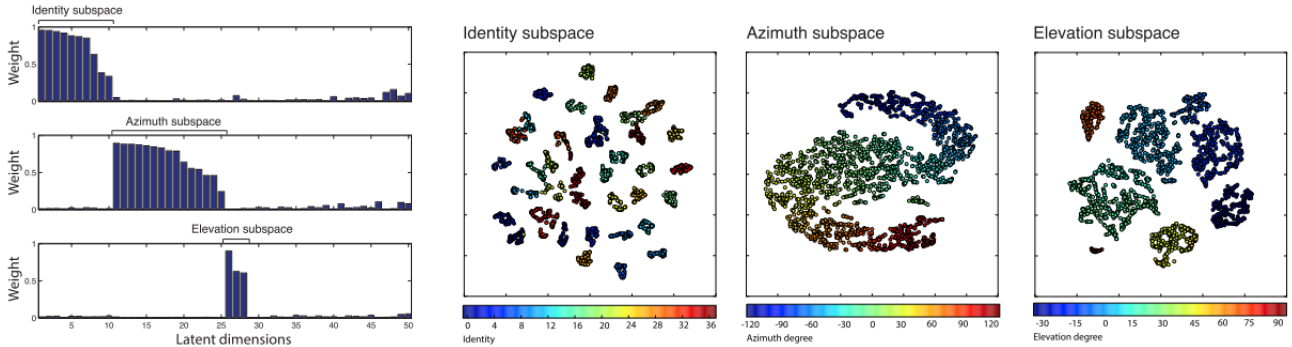
Figure 2.19: (left) A visualization of three factor-specific masks learned on the Yale Faces dataset, from [25]. (right) A t-SNE visualization of the identity, azimuth, and elevation subspaces as derived from the masked dimensions on the right. Only dimensions with a weight greater than 0.2 were used as input for this visualization. Colors represent the actual values for the factor of interest.

generate multi-factor map representations in [2]. Triplet constraints have also been used to learn highly accurate single-factor embedding models [47, 59].

### 2.4.3.6  Deep Visual Analogy-Making (Analogy and Equality)

In Deep Visual Analogy-Making (DVA) [41], two different types of constraints are used to learn a factorial representation in an autoencoder architecture. *Equality* constraints are used to enforce factorization of the representation into different attributes and *analogy* constraints are used to shape the subspace representing each factor. Each of these constraints is incorporated into an autoencoder architecture via its own constraint-specific loss function.

Assume we have data observations $x(a), x(b), x(c), x(d)$ in the analogy $x(a) : x(b) :: x(c) : x(d)$. The encoder function $e(\cdot)$ maps an observation to its representation. The generator function $g(\cdot)$ maps a representation back to the observational space. The loss is defined in terms of the square difference (in the observation space) of example $d$ with the analogical prediction for $d$ given $a, b, c$:

$$\mathcal{L}_{\text{ANALOGY\_ADD}} = \|x(d) - g(e[x(b)] - e[x(a)] + e[x(c)])\|_2^2 \tag{23}$$

This corresponds to a vector-addition interpretation of the analogy. However, this interpretation does not make sense for operations such as rotations. Rotations are circular, so a sufficient rotation should cause the point to end up back where it started. To address this issue, the authors include two additional interpretations of the analogy. One uses multiplicative interactions between $e(b) - e(a)$ and $e(c)$ to determine the vector that is added to $e(c)$. The other uses a neural network whose inputs are $e(b) - e(a)$ and $e(c)$. The neural net perform the best, because it is the most flexible of the three models.

The authors also use an additional loss, like Equation 23, but in the representation-space rather than the output-space, to enforce the relationships between the representations.

To learn a factorial representation, the authors use an additional penalty based on equality constraints. Suppose we have two non-overlapping sets of factors called f1 and f2, where $f1 \cap f2 = \emptyset$. Additionally, all of the factors f are included: $f1 \cup f2 = f$. We choose three example observations a, b, and c, where a and c share the f1 factors and b and c share the f2 factors. This corresponds to two factor-specific equalities: $a^{f1} = c^{f1}$ and $b^{f2} = c^{f2}$.

A binary mask $s$, of the same dimension as the representation, serves to select representation dimensions corresponding to f1 vs dimensions corresponding to f2. The mask $s$ is not learned; it is specified in advance. In this case, we expect to be able to reconstruct c using the dimensions corresponding to f1 from a and the dimensions corresponding to f2 from b. This corresponds to the following loss function:

$$\mathcal{L}_{\text{EQL}} = \|x(c) - e(s \cdot e[x(a)] + (1 - s) \cdot e[x(b)]\|_2^2 \tag{24}$$

To demonstrate the model, the authors train on a dataset of images of video game characters, called "sprites", and rotated 3D car renderings. An example sequence of reconstructions of analogies is shown in Figure 2.20 (a), in which the sequence of pose changes from the character on top is transferred to the character on the bottom. Each character has an identity and a pose. Three models were built: the first, trained with only analogy constraints, the second trained to separate identity from pose using equality constraints, and the third trained to classify identity with direct constraints and to predict pose with equality constraints. According to the constraint hierarchy in Table 1, the

model with the strongest constraints is the third model, followed by the second, then first. The prediction results, shown in Figure 2.20, from the paper follow this pattern. The authors demonstrate factor **distillation** by improved reconstruction performance on test analogies.
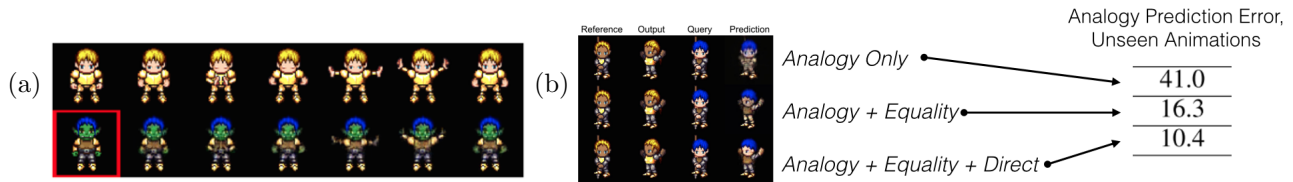


Figure 2.20: Example predictions from an autoencoder trained using analogies, from [41]. (a) "Animation transfer" results from the sprites dataset. (b) Analogy prediction results for models trained on three different combinations of constraints.

**Similar Models.** Visual analogies are also used to learn representations in [24]. Neural net word representations are known to obey analogical reasoning [38], but these models are not explicitly trained using analogy constraints.

### 2.4.3.7 Transforming Auto-Encoders (Distance and Equality)

An approach that combines equality and distance constraints is to teach a model to *transform* an observation with one setting for a factor, into another observation with a different known setting for that factor. The discrepancy between the factor values between the input and output is fed as input into the model. This approach also uses equality constraints, because it requires that all *other* factors remain equal while varying the target factor.

For example, the Transforming Auto-Encoder (TAE) [20] is based on a series auto-encoder capsules. Each capsule is presented an image as input, and is trained to reconstruct a translated version of that input image. The amount of horizontal and vertical translation is fed as input into the network during training. An illustration of three of this capsules is reproduced here on the left in Figure 2.21. The units in red represent *recognition* units, and the units in green represent *generation* units. The variables $x$ and $y$ represent the translation of the input image. $\Delta x$ and $\Delta y$ are fed as inputs to the network, and represent the $x$ and $y$ translation of the output image. The $p$ variable represents the probability that the object represented by the capsule is present in the image, and is multiplied with the output of the generation units to produce the output of the capsule.

The authors demonstrate factor **distillation** on translated MNIST digits and on stereo image pairs of rotated cars. In each case, they find that the TAE learns a representation of the translation or rotation that strongly correlates to the true value. This distillation is demonstrated in the scatterplot of x outputs for a capsule in Figure 2.21 on the right.

**Similar Models.** In [64], a representation of faces is learned that separates view from identity, and is trained to transform images in a similar way, using direct constraints on pose and equality constraints on identity. A similar approach is taken by, [60], which generates a sequence of transformed views of an input image using a recurrent network.
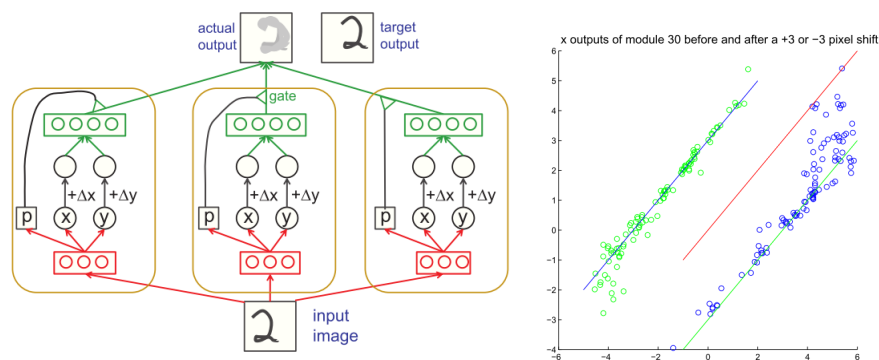


Figure 2.21: (left) An illustration of three Transforming Auto-Encoder capsules, from [20]. (right) A scatterplot showing, on the vertical axis, the x output of a capsule for an input image, and on the horizontal axis, the x output for the same capsule if the image is translated +3 (blue) or -3 (green) pixels in the x direction.
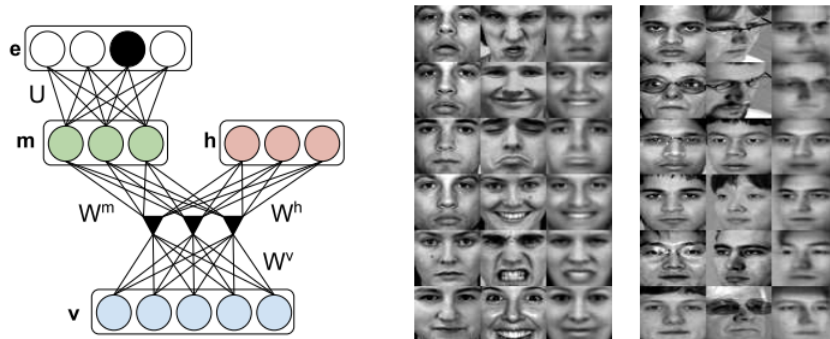
Figure 2.22: (left) A diagram of the DISBM model architecture from [40]. (right) Example reconstructions where factors are mixed between observations. The left block of images represents transfer of expression, and the right block represents transfer of pose. In each row, the left hand image represents the identity. This identity is combined with the expression or pose of the middle image to form the image on the right.

#### 2.4.3.8 Disentangling Boltzmann Machines (Direct and Equality)

The Disentangling Boltzmann Machine (DISBM) [40] adds supervisory constraints to a Restricted Boltzmann Machine with a multilinear bias. In the paper, the authors explore various combinations of constraints, and explore the effect of constraining only some of the factors in the representation.

Figure 2.22 shows a diagram of the model architecture. The hidden representation is divided into two sets of units, to model two factors. In the diagram, $\mathbf{v}$ are the visible units, $\mathbf{h}$ are a set of hidden units associated with one factor, and $\mathbf{m}$ are associated with another. *Direct constraints* are applied by adding an extra layer on top of $\mathbf{m}$ whose job is to predict the class label of that factor. This encourages the units in $\mathbf{m}$ to represent information related to that class. The version of the model with direct constraints is labeled DISBM (1) in Figure 2.1.

*Equality constraints* are applied in two different ways. First, the authors use a "clamping" procedure in which the units corresponding to a factor are encouraged to be constant for the pair. For example, if $\mathbf{h}$ represents face identity, we can clamp the $\mathbf{h}$ units for two images of the same individual. This is similar in spirit to the clamping procedure we saw used in the IGN [32].

The second method for incorporating equality constraints involves adding a penalty term to encourage the $\mathbf{h}$ representations to be more similar for the same individual. This is the same approach to the loss used by the SIAMESE network.

The version of the model with direct and equality constraints is labeled DISBM (1) in Figure 2.1, and the version with only equality constraints is DISBM (2).

The authors tested their models on two face datasets, the Toroto Face Dataset and Multi-PIE. Both of these datasets have individuals with differing identities. For both datasets, $\mathbf{h}$ units correspond to identity. For TFD, the $\mathbf{m}$ units correspond to expression, and for Multi-PIE they correspond to pose. On the right hand side of Figure 2.22, we can see an example of pose and expression transfer between examples. The identity units and pose or expression units are combined from two different images to form a new image that demonstrates the disentangling of the factors.

The authors performed experiments with different levels of supervision, and investigated disentangling: the effect of supervising one factor on disentangling of the other. They report performance on both prediction of the target factor value, as well as prediction of the value of the *non-target factor*. If the representation is more factorial, then performance on the target factor should improve and performance on the non-target factor should degrade. The tables in Figure 2.23 show predictive performance on each factor subspace, for the Multi-PIE dataset on the top and the TFD dataset on the bottom.

In each table, the authors explore five different models. For example, for Multi-PIE, "Naive" refers to a mutlilinear model without any supervision. "Labels (Pose)" is a model where the *pose* variable is supervised with direct constraints. "Clamp (ID)" uses the clamping procedure to learn identity. "Labels (Pose) + Clamp (ID)" is a combination of those two, and "Manifold (Both)" uses the penalty term to learn both factors. The same naming convention applies to TFD, but with "Pose" replaced with "Expression".

For each model and pool of units, the authors report test accuracy on a prediction task. Test accuracy on pose and emotion prediction is a percentage correct, and accuracy on verification is reported as a rate, from zero to one.

In the Multi-PIE table, we can see evidence of **distillation**: the supervised models generally have higher performance than the "Naive" model. We also see that constraining a factor improves its performance — the pose units of "Labels (Pose)" predict pose better than the "Naive" model. We see a similar effect with "Clamp (ID)". Also, as

| MODEL | POSE UNITS FOR POSE EST. | POSE UNITS FOR VERIFICATION | ID UNITS FOR POSE EST. | ID UNITS FOR VERIFICATION |
|---|---|---|---|---|
| NAIVE | $96.60 \pm 0.23$ | $0.583 \pm 0.004$ | $95.79 \pm 0.37$ | $0.640 \pm 0.005$ |
| LABELS (POSE) | $\mathbf{98.07} \pm 0.12$ | $0.485 \pm 0.005$ | $86.55 \pm 0.23$ | $0.656 \pm 0.004$ |
| CLAMP (ID) | $97.18 \pm 0.15$ | $0.509 \pm 0.005$ | $57.37 \pm 0.45$ | $0.922 \pm 0.003$ |
| LABELS (POSE) + CLAMP (ID) | $97.68 \pm 0.17$ | $0.504 \pm 0.006$ | $49.08 \pm 0.50$ | $0.934 \pm 0.002$ |
| MANIFOLD (BOTH) | $\mathbf{98.20} \pm 0.12$ | $0.469 \pm 0.005$ | $8.68 \pm 0.38$ | $\mathbf{0.975} \pm 0.002$ |

| MODEL | EXPR. UNITS FOR EMOTION REC. | EXPR. UNITS FOR VERIFICATION | ID UNITS FOR EMOTION REC. | ID UNITS FOR VERIFICATION |
|---|---|---|---|---|
| NAIVE | $79.50 \pm 2.17$ | $0.835 \pm 0.018$ | $79.81 \pm 1.94$ | $0.878 \pm 0.012$ |
| LABELS (EXPR) | $83.55 \pm 1.63$ | $0.829 \pm 0.021$ | $78.26 \pm 2.58$ | $0.917 \pm 0.006$ |
| CLAMP (ID) | $81.30 \pm 1.47$ | $0.803 \pm 0.013$ | $59.47 \pm 2.17$ | $\mathbf{0.978} \pm 0.025$ |
| LABELS (EXPR) + CLAMP (ID) | $82.97 \pm 1.85$ | $0.799 \pm 0.013$ | $59.55 \pm 3.04$ | $\mathbf{0.978} \pm 0.024$ |
| MANIFOLD (BOTH) | $\mathbf{85.43} \pm 2.54$ | $0.513 \pm 0.011$ | $43.27 \pm 7.45$ | $0.951 \pm 0.025$ |

Figure 2.23: (top) Model results for the Multi-PIE dataset. (bottom) Model results for the Toronto Face dataset.

predicted, test accuracy on the prediction of the non-target factor goes down, indicating **disentangling**.

We can also see a **cross-over** effect: the pose units of "Clamp (ID)" perform better at pose estimation than they do in the "Naive" model. This effect is similar to what we observed in the SSVAE — adding supervised constraints to one factor tends to remove related variability from other factors.

The same pattern of results generally holds for TFD. We see supervised models performing better than "Naive", and we see decreasing performance on non-target tasks. We can also see cross-over effects: the ID units of "Labels (Expr)" perform better at verification than the ID Units of "Naive". In three out of four cases, best test accuracy seems to come from manifold training.

### 2.4.4 Desiderata for Supervised Models

I have discussed various constraints, which have been incorporated into example models to help discover factorial representations. However, the models vary along four other dimensions as well. In order to completely characterize the models, I summarize these dimensions of variation, since they are also relevant for constructing factorial representations with supervisory signals. Table 2 shows each model laid out on these dimensions.

Table 2: Listing of supervised models.

| Example | Constraint | Factor Agnostic | Encoder | Prop. Dataset Labeled | Prop. Factors Labeled | Generator |
|---|---|---|---|---|---|---|
| SSVAE [29] | direct | | VAE | ◐ | ◐ | ✓ |
| DISBM (1) [40] | direct + equality | | multilinear | ◐ | ◐ | ✓ |
| DC-IGN [32] | equality | | VAE | ● | ● | ✓ |
| SIAMESE [61] | equality | | neural net | ● | ● | |
| DISBM (2) | equality | | multilinear | ◐ | ◐ | ✓ |
| TAE [20] | distance + equality | | neural net | ● | ◐ | ✓ |
| MM T-SNE [55] | distance | ✓ | lookup table | ● | ● | |
| KARALETSOS [25] | inequality | | VAE | ● | ◐ | ✓ |
| DVA [41] | analogy + equality | | neural net | ● | ● | ✓ |

#### 2.4.4.1 Factor-Agnostic vs Factor-Specific

Constraints can be *factor-specific*, meaning each label is directly tied to some known factor that is part of the representation. Constraints can also be *factor-agnostic*: the constraint has to be satisfied by some factor, but the particular factor is not indicated by the label. For example, people may attend to different features when rating bird image similarity: one person may attend primarily to head color and another person might attend to tail color. In each case, the label collected from the human is simply a similarity judgment, without any indication of which features the person was attending to. In this case, the factors being compared are hidden, and the model has to infer which factor a similarity constraint is associated with. The only model fully described here that was factor-agnostic was MM T-SNE. However, there are other models that use factor-agnostic constraints, such as [2].

Note that a factor-agnostic version of any constraint multiplies the size of the constrained space by the number of factors in the representation.

### 2.4.4.2   Encoder

Each model has some kind of *encoder* — a means for mapping an observation to its representation. Some models, such as MM T-SNE, simply have a fixed lookup table. This avoids additional inductive bias but also disallows looking up a representation from an un-labeled observation. This model does not make use of direct features of the low-level data observations — it uses only the distance constraints. Therefore, without distance labels for an observation, there is no representation for that observation.

   The other models I discussed all have some kind of functional encoder. The various encoders in these models all have different inductive biases: The SSVAE, DC-IGN, and KARALETSOS models are all based on the variational auto-encoder, which in turn uses the bias of Gaussian- distributed factors. The DISBM model is based on the multilinear architecture, one of the combination biases described earlier. The SIAMESE network, DVA and TAE are based on neural net architectures, and have lower unsupervised inductive bias for the encoder.

### 2.4.4.3   Proportion of Dataset Labeled

The strength of constraint is influenced by how much supervisory data we have. If we have labels for each observation, we say the task is *fully-supervised*, and is marked with a ● in the table. A task where the dataset has labels for only a partial subset of the data is called *semi-supervised*, and is marked with a ◖.

### 2.4.4.4   Proportion of Factors Labeled

Stronger constraints also have labels for a larger portion of factors in the representation. In most naturalistic environments, we cannot hope to label *all* factors of variation — for example, consider the factors in the face and glasses dataset. It is easy to come up with labels for face identity and glasses style. It is more difficult to label the position of the light source and the face's pose, and there are likely more factors of variation beyond those which we might have trouble even identifying. The more represented factors that are labeled, the stronger the constraints on the representation. Models that make use of a full set of factor labels are marked with a ●, and models that use a partial set of factor labels are marked with a ◖.

### 2.4.4.5   Generator

Some tasks, such as image generation, require a mapping from a representation back to the observation. Models suited to these tasks include a functional *generator*; a function that maps the representation back to the observation. All of the models except the SIAMESE network and MM T-SNE have a functional generator component.

   A functional generator puts a heavy burden on the representation; it needs to be able to reconstruct an input from its representation. In the case when only a subset of all factors are labeled, this means that the other factors are unsupervised. In the case of SSVAE and DISBM, we observed that cross-over effects of constraining one factor can help in learning the unsupervised factors.

   A functional generator also makes producing a fully-labeled dataset very difficult. It means that we need to know the full set of factors that contributed to generating the observations. This explains why models that have both a generator and a fully labeled dataset (DC-IGN and DVA) use synthetically-created datasets.

### 2.4.5   Hierarchy of Constraints

With so many different ways that supervision can vary, how can we rank strength of supervision across different models? In this paper, I use the type of constraint (direct, equality, etc) as the primary indicator of supervision strength. Concerns related to the proportion of examples labeled in the dataset, or the proportion of factors labeled, are fairly dataset-specific, but I am interested in the differences between constraints, not the particular test environments of the papers. Whether the constraint is factor-agnostic or not is important. However, as noted earlier, this concern will not move a constraint type up or down in the overall ranking, but will make any given constraint type slightly weaker. A model can also include supervisory information using multiple constraint types. When constraint types are mixed, the ranking of a model in the chart is determined by the weakest type of constraint used.

# 3    Discussion

## 3.1    Summary of Biases

A wide variety of biases can help discover factorial representations, and each bias comes with an associated set of generative assumptions. I demonstrated how sparsity bias in ICA can help to learn factorial representations when the generative environment is sparse. Likewise, an invariance bias helps bring coherency in factors that can be expressed in a variety of ways. Different combination biases are also shown to lead to factorial representations when the environment is matched to the generative assumptions.

Supervisory bias is used to shape representations to match prior knowledge about the generative environment. Different types of supervised constraints are associated with different types of prior knowledge, and vary in how strongly the constrain the representation.

The biases described here can be combined together in various ways. For example, sparse distribution bias is combined with invariance bias in the ISA model, with multilinear bias in a sparse bilinear model, and with hierarchical layers in R-ICA. Likewise, supervised constraints can be combined with an encoder with distributional, invariance, or combination biases. For example, DC-IGN combines the bias of Gaussian-distributed, independent, factors with supervised equality constraints. DISBM combines multilinear factor combination bias with equality and direct constraints. In general, models with supervised bias make use of low unsupervised bias — with enough supervised constraints, super-strong unsupervised bias is no longer necessary.

The key to discovering factorial representations is to match unsupervised inductive bias to the environment, and, if the unsupervised sources of bias are not sufficient, use additional supervisory signals to learn.

## 3.2    Effect of Bias on Factorization

Since there are numerous examples of different combinations of biases, the natural question is: how do they compare with each other? Which bias leads to representations where factors are most distilled and disentangled from one another? Can we lay out a space where different levels of supervised constraints correspond with degrees of factorization?

Unfortunately, the models cannot be compared to one another directly given the evidence from the papers. They make different assumptions about the environment, work in different data domains (e.g. image, audio), and use different data sets.

However, there is some indication that stronger constraints lead to more cross-over effects, which indicates better factorization. Each model demonstrates some evidence of factorization in the representation, and most models show evidence for distillation and disentangling. Evidence of a cross-over benefit is found for the SSVAE, DISBM, and, to a lesser extent, SIAMESE models. Their positions in the model chart in Figure 2.1, indicate that these models all have fairly strong supervised constraints, and DISBM has a strong combination bias. While it is possible that the other papers could have shown cross-over effects but did not report them, it seems plausible that higher levels of bias should cause stronger factorization in the representation.

Another hint that stronger constraints lead to better factorization is in the datasets used for training in the papers. The only papers that used partially labeled datasets were SSVAE and DISBM, which also used the strongest biases.

The model with the least evidence for factorization was MM T-SNE, despite the use of the fairly strong distance constraints. This can be attributed to the use of a dataset with a very high number of factors (in this case, topics), factor-agnostic constraints, and a very large space of inputs (thousands of words).

## 3.3    Directions for Future Research

### 3.3.1    Factorial Representations in Non-Factorial Environments

None of the papers considered the effect of data set sampling bias or non-factorial environments. Is a supervisory signal always necessary to learn factorial representations, or are some unsupervised inductive biases sufficient? An interesting future project could investigate the effect of the level of factorization of a dataset on a model's ability discover a factorial representation. For example, suppose we could vary the amount of correlation between identity and presence of glasses in a face dataset, and use a variable number of supervised constraint examples to learn a factorial representation. How many supervised examples are required to learn in the case when there is no correlation, in the case of some correlation, or in the case where identity and glasses are completely correlated? Does this relationship change with constraint type? Does this relationship change with our choice of unsupervised inductive bias?

### 3.3.2  New Combinations of Biases

The model chart in Figure 2.1 hints at some interesting unexplored combinations of biases.

Several of the unsupervised biases could be combined with an autoencoder to further reduce factor combination bias and help with invariance. For example, we could combine the bias of multilinear combination with an autoencoder architecture. Multilinear models are pretty wasteful in terms of model size — they need to learn one component vector for every combination of latent factors. An autoencoder could be both more efficient in representing these combinations, and have a lower bias in terms of how the combinations are combined. Sparse autoencoders combine the sparsity bias with an autoencoder. We could similarly extend an autoencoder to use pooled/block sparsity of ISA to find multi-dimensional sparse factors. We could even experiment with merging different combination biases, e.g. use the non-negativity bias of NMF in a MULTILINEAR model to create a functional parts model where the parts can interact multiplicatively.

We could similarly create new combinations of supervised biases. For example, we could apply the "learning to transform" model of TAE to direct constraints, and learn a model that represents the presence of glasses by transforming an image of a person with glasses to an image of the same person without.

Relatively few papers use factor-agnostic constraints, despite their potential usefulness. It is difficult to identify all the dimensions humans might use to compare examples. None of the papers that use factor-agnostic constraints learn a representation with a functional encoder. It would therefore be interesting to extend a model like KARALETSOS to handle factor-agnostic triplets.

# References

[1] Shun-ichi Amari and Andrzej Cichocki. Adaptive blind signal processing-neural network approaches. *Proceedings of the IEEE*, 86(10):2026–2048, 1998.

[2] Ehsan Amid and Antti Ukkonen. Multiview Triplet Embedding : Learning Attributes in Multiple Maps. *Proceedings of The 32nd International Conference on Machine Learning*, 37:1472–1480, 2015.

[3] Jl Austerweil and Tl Griffiths. The effect of distributional information on feature learning. *... -First Annual Conference of the ...*, pages 2765–2770, 2009.

[4] Joseph L Austerweil and Thomas L Griffiths. Learning invariant features using the transformed indian buffet process. In *Advances in neural information processing systems*, pages 82–90, 2010.

[5] H.B. Barlow, T.P. Kaushal, and G.J. Mitchison. Finding Minimum Entropy Codes. *Neural Computation*, 1(3):412–423, 1989.

[6] Horace B Barlow. Possible principles underlying the transformations of sensory messages. 1961.

[7] Horrace Barlow. Unsupervised Learning. *Neural Computation*, 1:295–311, 1989.

[8] Suzanna Becker. Unsupervised learning procedures for neural networks. *International Journal of Neural Systems*, 2(01n02):17–33, 1991.

[9] Anthony J. Bell and Terrence J. Sejnowski. An Information-Maximization Approach to Blind Separation and Blind Deconvolution. *Neural Computation*, 7(6):1129–1159, 1995.

[10] Christopher M Bishop. Pattern recognition. *Machine Learning*, 128, 2006.

[11] Brian Cheung, Jesse a. Livezey, Arjun K. Bansal, and Bruno a. Olshausen. Discovering Hidden Factors of Variation in Deep Networks. *arXiv preprint arXiv:1412.6583*, page 12, 2014.

[12] S Chopra, R Hadsell, and LeCun Y. Learning a similiarty metric discriminatively, with application to face verification. *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 349–356, 2005.

[13] Guillaume Desjardins, Aaron Courville, and Yoshua Bengio. Disentangling Factors of Variation via Generative Entangling. *arXiv:1210.5474 [cs, stat]*, (2), 2012.

[14] David Eigen, Marc Aurelio Ranzato, and Ilya Sutskever. Learning Factored Representations in a Deep Mixture of Experts. *ICLRws-2014*, pages 1–8, 2014.

[15] P. Foldiak. Forming sparse representations by local anti-Hebbian learning. *Biological Cybernetics*, 64(2):165–170, 1990.

[16] Zoubin Ghahramani, G Tesauro, D S Touretzky, and T K Leen. Factorial Learning and the EM Algorithm. *Advances in Neural Information Processing Systems*, (Mdl):617–624, 1995.

[17] David B Grimes. Bilinear Sparse Coding for Invariant Vision. 73:47–73, 2005.

[18] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2:1735–1742, 2006.

[19] Jeanny Herault and Christian Jutten. Space or time adaptive signal processing by neural network models. In *Neural networks for computing*, volume 151, pages 206–211. AIP Publishing, 1986.

[20] Geoffrey E. Hinton, Alex Krizhevsky, and Sida D. Wang. Transforming auto-encoders. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 6791 LNCS(PART 1):44–51, 2011.

[21] Aapo Hyvarinen. Fast and robust fixed-point algorithms for independent component analysis. *IEEE transactions on Neural Networks*, 10(3):626–634, 1999.

[22] Aapo Hyvärinen and Patrik Hoyer. Emergence of phase-and shift-invariant features by decomposition of natural images into independent feature subspaces. *Neural computation*, 12(7):1705–1720, 2000.

[23] Aapo Hyvärinen, Jarmo Hurri, and Patrick O Hoyer. *Natural Image Statistics: A Probabilistic Approach to Early Computational Vision.*, volume 39. Springer Science & Business Media, 2009.

[24] Sung Ju Hwang, Kristen Grauman, and Fei Sha. Analogy-preserving Semantic Embedding for Visual Object Categorization. *Proceedings of The 30th International Conference on Machine Learning*, 28:639–647, 2013.

[25] Theofanis Karaletsos, Serge Belongie, and Gunnar Rätsch. Bayesian representation learning with oracle constraints. *Iclr*, pages 1–9, 2015.

[26] Yan Karklin and Michael S Lewicki. Learning higher-order structures in natural images. *Network: Computation in Neural Systems*, 14(3):483–499, 2003.

[27] Koray Kavukcuoglu, Marc'Aurelio Ranzato, Rob Fergus, and Yann LeCun. Learning invariant features through topographic filter maps. *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2009*, pages 1605–1612, 2009.

[28] Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. (Ml):1–14, 2013.

[29] Dp Kingma, Dj Rezende, and Max Welling. Semi-supervised Learning with Deep Generative Models. *arXiv preprint arXiv: . . .*, pages 1–9, 2014.

[30] Teuvo Kohonen. Emergence of invariant-feature detectors in the adaptive-subspace self-organizing map. *Biological cybernetics*, 75(4):281–291, 1996.

[31] Joseph B Kruskal. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, 29(1):1–27, 1964.

[32] TD Kulkarni and W Whitney. Deep Convolutional Inverse Graphics Network. *Advances in Neural Information Processing Systems*, pages 2539–2547, 2015.

[33] Quoc V. Le, Will Y. Zou, Serena Y. Yeung, and Andrew Y. Ng. Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 3361–3368, 2011.

[34] Qv Le, Alexandre Karpenko, Jiquan Ngiam, and Ay Ng. ICA with reconstruction cost for efficient overcomplete feature learning. *Advances in Neural . . .*, pages 1–9, 2011.

[35] D D Lee and H S Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–91, 1999.

[36] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.

[37] Alireza Makhzani and Brendan J Frey. Winner-Take-All Autoencoders. *Nips*, pages 2773–2781, 2015.

[38] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.

[39] Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semi-supervised learning with ladder networks. In *Advances in Neural Information Processing Systems*, pages 3546–3554, 2015.

[40] Scott Reed, Kihyuk Sohn, Yuting Zhang, and Honglak Lee. Learning to Disentangle Factors of Variation with Manifold Interaction. *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, (May):1431–1439, 2014.

[41] Scott E. Reed, Yi Zhang, Yuting Zhang, and Honglak Lee. Deep Visual Analogy-Making. *Advances in Neural Information Processing Systems*, pages 1252–1260, 2015.

[42] Salah Rifai, Yoshua Bengio, Aaron Courville, Pascal Vincent, and Mehdi Mirza. Disentangling factors of variation for facial expression recognition. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 7577 LNCS(PART 6):808–822, 2012.

[43] David A Ross and Richard S Zemel. Multiple Cause Vector Quantization. *Learning*, pages 2–9, 2003.

[44] Yoshinori Sakaguchi, Seiichi Ozawa, and Manabu Kotani. Feature extraction using supervised independent component analysis by maximizing class distance. In *Neural Information Processing, 2002. ICONIP'02. Proceedings of the 9th International Conference on*, volume 5, pages 2502–2506. IEEE, 2002.

[45] Eric Saund. A Multiple Cause Mixture Model for Unsupervised Learning. *Neural Computation*, 7(1):51–71, 1995.

[46] Jürgen Schmidhuber. Learning Factorial Codes by Predictability Minimization. *Neural Computation*, 4(6):863–879, 1992.

[47] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 815–823, 2015.

[48] Honghao Shan, Lingyun Zhang, and Garrison W Cottrell. Recursive ICA. *Advances in Neural Information Processing Systems 19*, pages 1273–1280, 2006.

[49] Paul Smolensky. Tensor product variable binding and the representation of symbolic structures in connectionist systems. *Artificial intelligence*, 46(1-2):159–216, 1990.

[50] Anisse Taleb and Christian Jutten. Source separation in post-nonlinear mixtures. *IEEE Transactions on signal Processing*, 47(10):2807–2820, 1999.

[51] Yichuan Tang and Geoffrey Hinton. Tensor Analyzers. *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 28, 2013.

[52] J Tenenbaum. Functional parts. In *Proceedings of the 16th Annula Conference of the Congnitive Science Society*, 1994.

[53] J B Tenenbaum and W T Freeman. Separating style and content with bilinear models. *Neural computation*, 12(6):1247–1283, 2000.

[54] Joshua B Tenenbaum and Emanuel Todorov. Factorial Learning by Clustering Features. *Advances in Neural Information Processing Systems 7*, pages 561–568, 1995.

[55] Laurens Van Der Maaten and Geoffrey Hinton. Visualizing non-metric similarities in multiple maps. *Machine Learning*, 87(1):33–55, 2012.

[56] Andreas Veit, Serge Belongie, and Theofanis Karaletsos. Disentangling Nonlinear Perceptual Embeddings With Multi-Query Triplet Networks. 2016.

[57] Pascal Vincent and Hugo Larochelle. Extracting and Composing Robust Features with Denoising Autoencoders. 2008.

[58] Changhu Wang, Shuicheng Yan, Lei Zhang, and Hongjiang Zhang. Non-negative semi-supervised learning. In *AISTATS*, pages 575–582, 2009.

[59] Jiang Wang, Yang Song, Thomas Leung, Chuck Rosenberg, Jingbin Wang, James Philbin, Bo Chen, and Ying Wu. Learning fine-grained image similarity with deep ranking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1386–1393, 2014.

[60] Jimei Yang, Ming-Hsuan Yang, Scott E. Reed, and Honglak Lee. Weakly-supervised Disentangling with Recurrent Transformations for 3D View Synthesis. *Nips*, pages 1–9, 2015.

[61] Neil Zeghidour, Gabriel Synnaeve, Nicolas Usunier, and Emmanuel Dupoux. Joint Learning of Speaker and Phonetic Similarities with Siamese Networks. *Interspeech 2016*, pages 1295–1299, 2016.

[62] Richard S Zemel. *A minimum description length framework for unsupervised learning*. PhD thesis, University of Toronto, 1993.

[63] Richard Stanley Zemel. A Minimum Description Length Framework for Unsupervised Learning. 1993.

[64] Zhenyao Zhu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Multi-View Perceptron : a Deep Model for Learning Face Identity and View Representations. *Advances in Neural Information Processing Systems*, pages 1–9, 2014.