

IV. Modeling and Model Selection

One of the major themes in this course has been systematic methods for inference under the assumption that the unknowns obeyed some kind of structure. For instance, we have seen many ways to estimate a function from its samples (using least-squares, for example) when that function is assumed to live in a known subspace. We have also seen examples of estimating the probability density function of a random variable/vector when that pdf is assumed to have some parametric form (e.g. the MLE). Finally, we have used pdfs to model the unknowns themselves (Gaussian estimation, MAP estimation, etc).

In this final chapter of the course, we will learn a little bit about how to fit models to data. Problems of this type are typically referred to as “unsupervised learning” problems in the ML literature, but we will shy away from this terminology. We will also encounter, for purposes of illustration, new models for data.

In the very first set of notes in this course, we made the distinction between *geometric* models, and *probabilistic* models. The former is basically the specification of some set (e.g. subspace) that the signal/feature vector/parameters live in or near. The latter is a density function that indicated which values of the vector are expected to be more likely than others.

We will start with the most fundamental geometric model selection problem: how to fit a subspace to a bunch of points.

Principal Components Analysis (PCA)

On many occasions over the last three months, we have made the assumption that some unknown vector that we were interested in belonged to a subspace. PCA is a way in which we can *learn* this subspace from examples.

We are given data points $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$. We would like to find the $K < D$ dimensional subspace that comes closest to containing these points — maybe this is a good fit, maybe it isn't, but the point here is that we can set up an optimization program and solve it (and in the process determine how good the fit is).

PCA gives an easy two-step procedure for solving this problem:

1. Form the sample covariance matrix

$$\hat{\mathbf{R}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x} \mathbf{x}^T \quad (1)$$

2. Factor $\hat{\mathbf{R}}$ using an eigenvalue decomposition

$$\hat{\mathbf{R}} = \hat{\mathbf{V}} \mathbf{\Lambda} \hat{\mathbf{V}}^T.$$

Then take the first K eigenvectors (the eigenvectors corresponding to the K largest eigenvalues) — $\text{span}\{\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_K\}$ is the optimal subspace.

Let's make this precise.

Derivation of Principal Components Analysis

We want to find a subspace $\mathcal{T} \subset \mathbb{R}^D$ with $\dim(\mathcal{T}) = K$ that comes closest to containing $\mathbf{x}_1, \dots, \mathbf{x}_N$. To make the notion of “closest” precise, we will score a subspace \mathcal{T} by looking at the sum of the squared distances between the \mathbf{x}_n and \mathcal{T} :

$$\text{Score}(\mathcal{T}) = \sum_{n=1}^N \text{dist}(\mathbf{x}_n, \mathcal{T})^2.$$

We know that given an orthobasis $\mathbf{u}_1, \dots, \mathbf{u}_K$ for \mathcal{T} , the closest point in \mathcal{T} is

$$\hat{\mathbf{x}}_n = \sum_{k=1}^K \langle \mathbf{x}_n, \mathbf{u}_k \rangle \mathbf{u}_k = \mathbf{U} \mathbf{U}^T \mathbf{x}_n,$$

with distance

$$\|\mathbf{x}_n - \hat{\mathbf{x}}_n\|_2^2 = \|(\mathbf{I} - \mathbf{U} \mathbf{U}^T) \mathbf{x}_n\|_2^2,$$

where $\mathbf{U} = [\mathbf{u}_1 \ \cdots \ \mathbf{u}_K]$. Thus we can replace our search for a subspace with a search for an *orthobasis*; PCA solves

$$\underset{\mathbf{U} \in \mathbb{R}^{D \times K}}{\text{minimize}} \quad \sum_{n=1}^N \|(\mathbf{I} - \mathbf{U} \mathbf{U}^T) \mathbf{x}_n\|_2^2 \quad \text{subject to} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}. \quad (2)$$

Note that the solution to the program above is certainly not unique — any orthobasis for the same subspace will lead to exactly the same functional value.

Notice that because $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, we have

$$(\mathbf{I} - \mathbf{U} \mathbf{U}^T)(\mathbf{I} - \mathbf{U} \mathbf{U}^T) = (\mathbf{I} - \mathbf{U} \mathbf{U}^T),$$

and thus

$$\|(\mathbf{I} - \mathbf{U}\mathbf{U}^T)\mathbf{x}_n\|_2^2 = \mathbf{x}_n^T \mathbf{x}_n - \mathbf{x}_n^T \mathbf{U}\mathbf{U}^T \mathbf{x}_n = \|\mathbf{x}_n\|_2^2 - \|\mathbf{U}^T \mathbf{x}_n\|_2^2,$$

and since $\|\mathbf{x}_n\|_2^2$ does not depend on \mathbf{U} , we can re-write the PCA program as

$$\underset{\mathbf{U} \in \mathbb{R}^{D \times K}}{\text{maximize}} \quad \sum_{n=1}^N \|\mathbf{U}^T \mathbf{x}_n\|_2^2 \quad \text{subject to} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}.$$

Since¹ $\|\mathbf{U}^T \mathbf{x}_n\|_2^2 = \text{trace}(\mathbf{U}^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{U})$, we have

$$\begin{aligned} \sum_{n=1}^N \|\mathbf{U}^T \mathbf{x}_n\|_2^2 &= \sum_{n=1}^N \text{trace}(\mathbf{U}^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{U}) \\ &= \text{trace} \left(\mathbf{U}^T \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{U} \right) \\ &= N \text{trace}(\mathbf{U}^T \hat{\mathbf{R}} \mathbf{U}), \end{aligned}$$

where $\hat{\mathbf{R}}$ is the sample covariance defined in (1). We now know that PCA is equivalent to solving

$$\underset{\mathbf{U} \in \mathbb{R}^{D \times K}}{\text{maximize}} \quad \text{trace}(\mathbf{U}^T \hat{\mathbf{R}} \mathbf{U}) \quad \text{subject to} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}.$$

By construction, $\hat{\mathbf{R}}$ is symmetric positive semi-definite, so it has an eigenvalue decomposition

$$\mathbf{R} = \hat{\mathbf{V}} \mathbf{\Lambda} \hat{\mathbf{V}}^T.$$

¹If you are uncomfortable with this equality, you should take the time right now to sit down and derive it.

We have

$$\text{trace}(\mathbf{U}^T \hat{\mathbf{R}} \mathbf{U}) = \text{trace}(\mathbf{U}^T \hat{\mathbf{V}} \mathbf{\Lambda} \hat{\mathbf{V}}^T \mathbf{U}) = \text{trace}(\mathbf{W}^T \mathbf{\Lambda} \mathbf{W}),$$

where $\mathbf{W} = \hat{\mathbf{V}}^T \mathbf{U}$, and note² that $\mathbf{W}^T \mathbf{W} = \mathbf{I}$. Thus we can solve (2) by solving

$$\underset{\mathbf{W} \in \mathbb{R}^{D \times K}}{\text{maximize}} \quad \text{trace}(\mathbf{W}^T \mathbf{\Lambda} \mathbf{W}) \quad \text{subject to} \quad \mathbf{W}^T \mathbf{W} = \mathbf{I},$$

and then taking $\hat{\mathbf{U}} = \hat{\mathbf{V}} \hat{\mathbf{W}}$ to get the minimizer in (2).

We can show that the last maximization program above is equivalent to a simple linear program that we can solve by inspection. Let $\mathbf{w}_1, \dots, \mathbf{w}_K$ be the columns of \mathbf{W} . Then

$$\begin{aligned} \text{trace}(\mathbf{W}^T \mathbf{\Lambda} \mathbf{W}) &= \sum_{k=1}^K \mathbf{w}_k^T \mathbf{\Lambda} \mathbf{w}_k \\ &= \sum_{k=1}^K \sum_{d=1}^D (w_k[d])^2 \lambda_d \\ &= \sum_{d=1}^D h_d \lambda_d, \quad h_d = \sum_{k=1}^K (w_k[d])^2 = \sum_{k=1}^K (W[d, k])^2. \end{aligned}$$

The h_d , $d = 1, \dots, D$ above are the sums of the squares of the *rows* of \mathbf{W} . It is clear that $h_d \geq 0$. It is also true that

$$\sum_{d=1}^D h_d = K,$$

² $\hat{\mathbf{V}}$ is a $D \times D$ orthonormal matrix, so $\hat{\mathbf{V}} \hat{\mathbf{V}}^T = \mathbf{I}$.

as the fact that the norm of each columns of \mathbf{W} is 1 means that

$$\sum_{d=1}^D \sum_{k=1}^K W[d, k]^2 = \sum_{k=1}^K \left(\sum_{d=1}^D W[d, k]^2 \right) = \sum_{k=1}^K 1 = K.$$

Finally, it is also true that $h_d \leq 1$. Here is why: since the columns of \mathbf{W} are orthonormal, they can be considered as part of an orthonormal basis for all of \mathbb{R}^D . That is, there is a (and actually there are many) $D \times (D - K)$ matrix \mathbf{W}_0 such that the columns of

$$\mathbf{W}' = [\mathbf{W} \quad \mathbf{W}_0]$$

form an orthonormal basis for \mathbb{R}^D . Since \mathbf{W}' is square, $\mathbf{W}'\mathbf{W}'^T = \mathbf{I}$, meaning the sum of the squares of each row are equal to 1. Thus

$$h_d = \sum_{k=1}^K W[d, k]^2 = \sum_{k=1}^K W'[d, k]^2 \leq \sum_{k=1}^D W'[d, k]^2 = 1.$$

With these constraints on the h_d , let's see how large we can make the quantity of interest:

$$\underset{\mathbf{h} \in \mathbb{R}^D}{\text{maximize}} \sum_{d=1}^D h_d \lambda_d \quad \text{subject to} \quad \sum_{d=1}^D h_d = K, \quad 0 \leq h_d \leq 1.$$

This is a linear program, but we can intuit the answer. Since all of the λ_d are positive, we want to have their weights (i.e. the h_d) as large as possible for the largest entries. Since the weights are constrained to be less than 1, and their sum is K , this simply means we assign a weight of 1 to the K largest terms, and 0 to the others:

$$\hat{h}_d = \begin{cases} 1, & j = 1, \dots, K, \\ 0, & \text{otherwise} \end{cases}$$

This means that the sum of the squares of the entries in the rows of the corresponding $\hat{\mathbf{W}}$ are 1 for the first K , and zero below — there are many matrices with orthonormal columns which fit the bill, but a specific one which does is

$$\hat{\mathbf{W}} = \begin{bmatrix} \mathbf{I}_{K \times K} \\ \mathbf{0}_{(D-K) \times K} \end{bmatrix}. \quad (3)$$

Taking $\hat{\mathbf{U}} = \hat{\mathbf{V}}\hat{\mathbf{W}}$, this results in

$$\hat{\mathbf{U}} = [\hat{\mathbf{v}}_1 \quad \hat{\mathbf{v}}_2 \quad \cdots \quad \hat{\mathbf{v}}_K].$$

Including an affine offset

PCA is easily modified to include an affine offset (so we fit a subspace that has been shifted away from the origin) that we learn as well. Now, we simply solve

$$\underset{\mathbf{U} \in \mathbb{R}^{D \times K}, \boldsymbol{\mu} \in \mathbb{R}^D}{\text{minimize}} \quad \sum_{n=1}^N \|(\mathbf{I} - \mathbf{U}\mathbf{U}^T)(\mathbf{x}_n - \boldsymbol{\mu})\|_2^2$$

in place of (2) above. With \mathbf{U} fixed, we can take the gradient of the functional above with respect to $\boldsymbol{\mu}$ to get the optimality condition for $\hat{\boldsymbol{\mu}}$:

$$\mathbf{0} = -2 \sum_{n=1}^N (\mathbf{I} - \mathbf{U}\mathbf{U}^T)(\mathbf{x}_n - \hat{\boldsymbol{\mu}}) = -2(\mathbf{I} - \mathbf{U}\mathbf{U}^T) \left(\sum_{n=1}^N \mathbf{x}_n - N\hat{\boldsymbol{\mu}} \right),$$

which means that $\hat{\boldsymbol{\mu}}$ is the sample mean³:

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

³We are using the terms “sample mean” and “sample covariance” just because of the forms of these terms — we do not need to interpret the \mathbf{x}_n as outcomes of random variables for this to all make sense.

This answer is independent of \mathbf{U} , so we can just subtract off $\boldsymbol{\mu}$ and proceed as before.

It is also easy to see that PCA can be interpreted as searching over an orthobasis, an expansion $\boldsymbol{\alpha}_n \in \mathbb{R}^K$ in that orthobasis, and an affine offset $\boldsymbol{\mu}$. We just of course know, because of all the hard work we did in September, that with \mathbf{U} fixed, $\hat{\boldsymbol{\alpha}}_n = \mathbf{U}^T \mathbf{x}_n$.

We now summarize the results in this section:

PCA Theorem

$$\underset{\boldsymbol{\mu}, \mathbf{U}, \{\boldsymbol{\alpha}_n\}}{\text{minimize}} \sum_{n=1}^N \|\mathbf{x}_n - \boldsymbol{\mu} - \mathbf{U} \boldsymbol{\alpha}_n\|_2^2, \quad \text{subject to } \mathbf{U}^T \mathbf{U} = \mathbf{I},$$

has solution

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_n, \quad \hat{\mathbf{U}} = [\hat{\mathbf{v}}_1 \ \cdots \ \hat{\mathbf{v}}_K], \quad \hat{\boldsymbol{\alpha}}_n = \hat{\mathbf{U}}^T (\mathbf{x}_n - \hat{\boldsymbol{\mu}}),$$

where $\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_K$ are the eigenvectors corresponding to the K largest eigenvalues of

$$\hat{\mathbf{R}} = \sum_{i=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}})(\mathbf{x}_n - \hat{\boldsymbol{\mu}})^T.$$

Structured Matrix Factorization

PCA is the simplest example of a class of problems that are known as **structured matrix factorization**. The idea is as follows. Given a $D \times N$ matrix \mathbf{X} , we want to find a $D \times R$ matrix \mathbf{B} and a $R \times N$ matrix \mathbf{C} such that

$$\mathbf{X} \approx \mathbf{BC},$$

under the constraints that \mathbf{B} and \mathbf{C} have some kind of *structure*.

To interpret PCA in this framework, we take $R = K$,

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_N],$$

$\mathbf{B} = \mathbf{U}$ and $\mathbf{C} = [\boldsymbol{\alpha}_1 \ \boldsymbol{\alpha}_2 \ \cdots \ \boldsymbol{\alpha}_N]$. The structure imposed is that the columns of \mathbf{U} are orthonormal.

The following is a more general statement of the PCA theorem (without the affine offset), and is proved in essentially the same way. It is one of the fundamental results in modern linear algebra.

Eckart-Young Theorem: Let \mathbf{X} be a $D \times N$ matrix with SVD $\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$. Then the best rank- R approximation to \mathbf{X} ,

$$\underset{\mathbf{Y} \in \mathbb{R}^{D \times N}}{\text{minimize}} \ \|\mathbf{X} - \mathbf{Y}\|_F^2 \quad \text{subject to} \quad \text{rank}(\mathbf{Y}) = R,$$

is

$$\hat{\mathbf{Y}} = \mathbf{U}_R \boldsymbol{\Sigma}_R \mathbf{V}_R^T,$$

where $\boldsymbol{\Sigma}_R$ is the $R \times R$ diagonal matrix containing the R largest singular values of \mathbf{X} , \mathbf{U}_R is a $D \times R$ matrix whose columns are the corresponding left singular vectors, and \mathbf{V}_R contains the corresponding right singular vectors.

We can also replace the Frobenius norm $\|\cdot\|_F^2$ (sum of the squares of the singular values) in the optimization above with the spectral norm $\|\cdot\|$ (largest singular value). We could also re-write the program as

$$\underset{\substack{\mathbf{B} \in \mathbb{R}^{D \times R} \\ \mathbf{C} \in \mathbb{R}^{R \times N}}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{BC}\|_F^2.$$

The rank constraint is now implicit in the size of the matrices \mathbf{B} and \mathbf{C} (as when had for PCA). A solution to the above is then

$$\hat{\mathbf{B}} = \mathbf{U}_R \boldsymbol{\Sigma}_R^{1/2}, \quad \hat{\mathbf{C}} = \boldsymbol{\Sigma}_R^{1/2} \mathbf{V}_R^T.$$

Of course, this solution is not unique, as we may divvy up the $\boldsymbol{\Sigma}_R$ in an arbitrary way, i.e.

$$\hat{\mathbf{B}} = \mathbf{U}_R \boldsymbol{\Sigma}_R, \quad \hat{\mathbf{C}} = \mathbf{V}_R^T,$$

or even apply a general invertible $R \times R$ matrix \mathbf{A} to the left and right:

$$\hat{\mathbf{B}} = \mathbf{U}_R \boldsymbol{\Sigma}_R^{1/2} \mathbf{A}, \quad \hat{\mathbf{C}} = \mathbf{A}^{-1} \boldsymbol{\Sigma}_R^{1/2} \mathbf{V}_R^T,$$

etc.

As we saw in the previous section, it is easy (and most natural) to make the factors $\hat{\mathbf{B}}$ or $\hat{\mathbf{C}}$ have *orthogonal* columns or rows, respectively.

What are we doing when we perform an approximation of this kind? Suppose that the columns of \mathbf{X} are data points in \mathbb{R}^D :

$$\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_N]$$

and suppose that there exists a small $R \ll \min(D, N)$, such that

$$\begin{aligned} \mathbf{X} &\approx \mathbf{BC} \\ &= \begin{bmatrix} \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{c}_1 & \mathbf{c}_2 & \cdots & \mathbf{c}_N \end{bmatrix}, \end{aligned}$$

i.e. $\mathbf{x}_n \approx \mathbf{B}\mathbf{c}_n$. This means that each \mathbf{x}_n is a different (linear) combination (specified by the entries in \mathbf{c}_n) of a small number of latent features (specified by the columns of \mathbf{B}). So in solving this optimization algorithm, we are learning a basis for the \mathbf{x}_n .

Example: Eigenfaces

The data points \mathbf{x}_n are D -pixel images⁴ of faces (D is on the order of thousands):



There are many such images (N on the order of 1000s), but \mathbf{X} can be approximated very well with rank $R = 200$. Here are the columns of $\hat{\mathbf{B}}$:

⁴Example from <http://kixor.net/school/2008spring/comp776/assn3/>.



Example: Latent Semantic Indexing

LSI is used to automatically group documents together by their semantic content. Here, the data we are given is a *document term matrix*. We parse N documents (encyclopedia/wikipedia entries, newspaper article, etc.) and count the number of times D different words appear; the tally for word d in document n is the entry $X[d, n]$:

$$X[d, n] = \# \text{ times word } d \text{ appeared in document } n.$$

The rows of \mathbf{X} tend to be correlated, so we can approximate $\mathbf{X} \approx \mathbf{BC}$ for $r \ll \min(D, N)$. This means that word frequency counts in a document are basically a linear combination of a small number of “feature counts” given by the columns of \mathbf{B} .

But in both these cases, it is hard to interpret what the latent factors really tell us about the data, other than being able to assist us in some simple dimensionality reduction. This is because we have not enforced the proper structure on the factors \mathbf{B}, \mathbf{C} .

Non-negative matrix factorization

If we constrain the factors to have positive entries, we get different (and perhaps more interesting) results:

$$\underset{\substack{\mathbf{B} \in \mathbb{R}^{D \times R} \\ \mathbf{C} \in \mathbb{R}^{R \times N}}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{BC}\|_F^2, \quad \mathbf{B} \geq \mathbf{0}, \mathbf{C} \geq \mathbf{0}. \quad (4)$$

The idea behind this is that in many cases (including the face and LSI examples above), what we are really looking for is a mixture of *positive* factors. This problem is called *non-negative matrix factorization* (NNMF).

Solving (4) is much more involved than in the unconstrained case. Not only is there no closed-form solution, but the optimization is non-convex, and comes with all the complications there-of.

But notice that if we fix one of the factors, say $\mathbf{B} = \mathbf{B}^{(0)}$, then the program

$$\underset{\mathbf{C} \in \mathbb{R}^{R \times N}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{B}^{(0)}\mathbf{C}\|_F^2, \quad \mathbf{C} \geq \mathbf{0}.$$

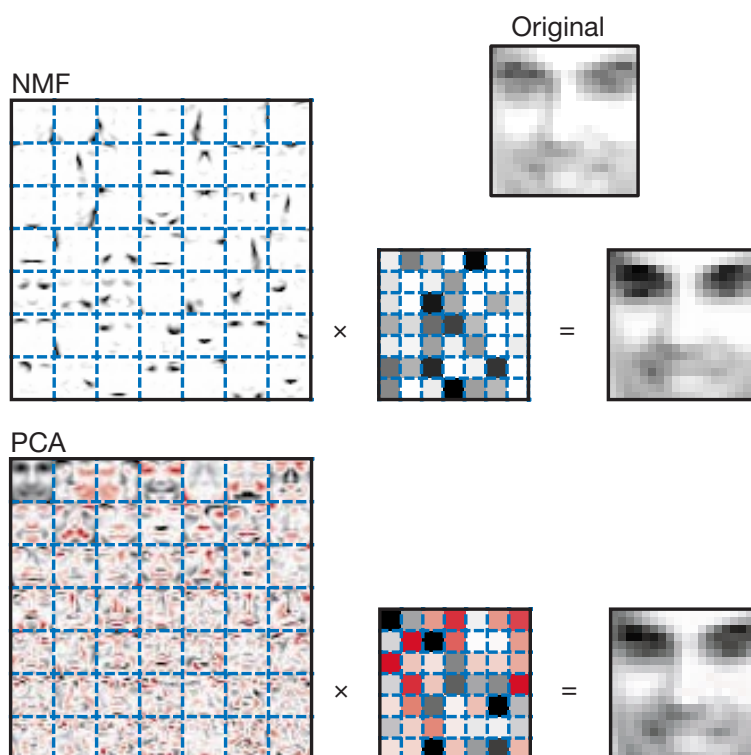
is separable by column; we can solve

$$\underset{\mathbf{c}_n \in \mathbb{R}^R}{\text{minimize}} \quad \|\mathbf{x}_n - \mathbf{B}^{(0)}\mathbf{c}_n\|_2^2 \quad \text{subject to} \quad \mathbf{c}_n \geq \mathbf{0},$$

for $n = 1, \dots, N$ independently. Each of these subproblems is a quadratics program, and so it is solvable with off-the-shelf software. The same thing is true if we pin down $\mathbf{C} = \mathbf{C}^{(0)}$ and optimize over \mathbf{B} . So NNMF algorithms typically alternate back and forth between solving for \mathbf{B} and solving for \mathbf{C} .

Here are some interesting results from the original paper on NNMF⁵

Faces:



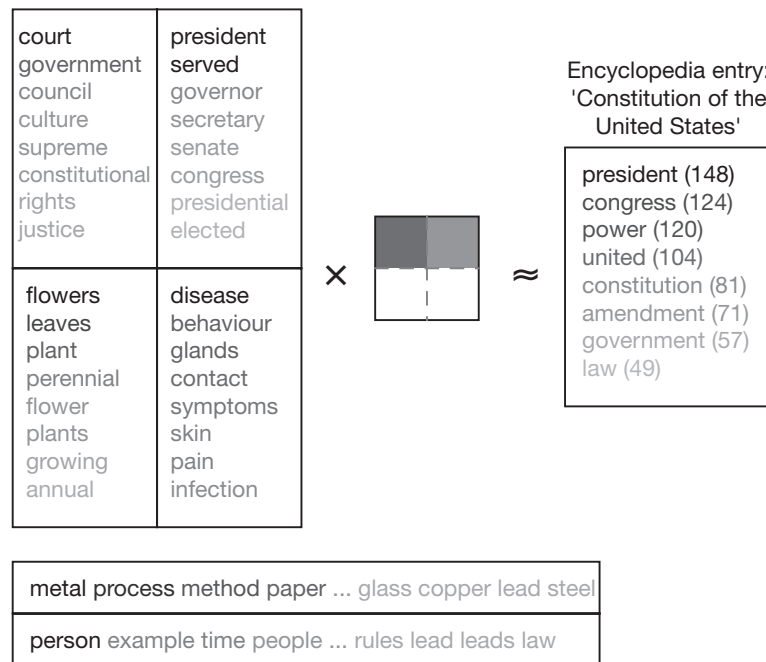
The top is essentially the eigenfaces experiment, but where we restrict the factors (\mathbf{B}) and mixture coefficients (\mathbf{C}) to be positive. Notice that, as compared to the standard PCA experiment below it, that the factors actually look like facial features.

⁵Lee and Seung, “Learning the parts of objects by non-negative matrix factorization,” *Nature*, October 1999.

Latent semantic indexing:

Here is an example from the same article for LSI.

$N = 30,991$ articles from an encyclopedia were parsed,
 $D = 15,276$ words counted



Now the columns of $\hat{\mathbf{B}}$ can be interpreted as *groupings* of words that appear together frequently.

On the top right, we see example weights for four different factors that decomposed the (word count of) the US Constitution ... the two factors (columns of \mathbf{B}) on top, learned while training on a separate data set, were given significant weights (entries in \mathbf{c}), while the two factors on bottom had very small weights.

The example on bottom shows the words associated with the two factors whose entry in the row corresponding to the word 'lead' was the greatest.