# CS281 Section 4: Factor Analysis and PCA

Scott Linderman

At this point we have seen a variety of machine learning models, with a particular emphasis on models for supervised learning. In particular, we have discussed linear regression for fitting continuous outputs given their corresponding features, and classification methods which learn a linear decision boundary, like logistic regression. Finally, we've extended these methods to include nonlinear functions of the input using generalized linear models.

In many settings, however, we do not have a set of outputs and features, we just have a set of data. This is the domain of **unsupervised machine learning**. Our goals in this setting are less clear. Rarely is our data entirely random. We often believe there is some latent, low-dimensional structure in the data that is corrupted by noise. Unsupervised machine learning is about finding this latent structure, and today we will discuss some of the most widely used methods for doing so.

1. Principal Component Analysis (PCA)

   Imagine we are presented with a bunch of data $\{x_n\}$, where each $x_n$ lives in $\mathbb{R}^D$. For example, in Figure 1 we have a cloud of points in $\mathbb{R}^2$. In many cases we believe the data is actually lower dimensional, 1-dimensional in this case. That is, the 2-D data points are well approximated by a line, and each data point is described by a 1-D value indicating the location on the line. How could we find the line that best explains these data? This is very similar to linear regression, except in this case we do not know the latent position along the line.

   (a) Heuristic motivation

   One way is to assume the data is Gaussian distributed and find the MLE covariance. This corresponds to the sample covariance. Then, intuitively, the best line will be parallel to the long axis of the ellipse corresponding to the covariance matrix. That is, the best line will be parallel to the eigenvector of the covariance matrix with the largest eigenvalue. Each point will be approximated by its projection onto this eigenvector.

   This is the intuition behind PCA: We find the eigenvectors of the sample covariance matrix, and then summarize the observed data as projections onto the top eigenvectors, or **principal components**. The size of this projection for a given datapoint $x_n$ is called the **score**.

   (b) Minimizing reconstruction error

   Can we formalize the process by which we arrived at this algorithm? Under what objective function are the eigenvectors with largest eigenvalues the "best" principal components? It turns out that given a data matrix $X \in \mathbb{R}^{D \times N}$ (columns correspond to $x_n$), PCA finds principal components $W \in \mathbb{R}^{D \times L}$ and scores $Z^{L \times N}$ that minimize the
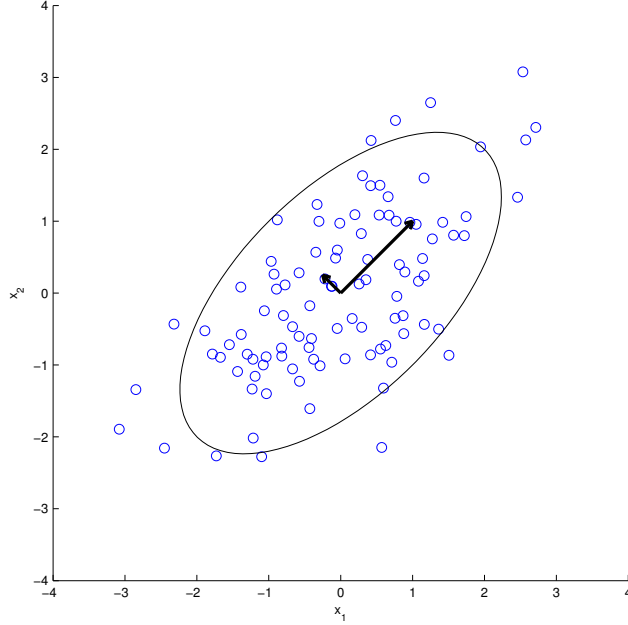
Figure 1: A cloud of data in $\mathbb{R}^2$ that we would like to approximate with a lower dimensional representation.

reconstruction error

$$J(\boldsymbol{W}, \boldsymbol{Z}) = ||\boldsymbol{X} - \boldsymbol{W}\boldsymbol{Z}||_F^2,$$

where $\boldsymbol{W}$ is constrained to be orthonormal and the error is measured with respect to the Frobenius norm.

Let us solve for the optimal first principal component $\boldsymbol{w}_1$ and "scores" $\boldsymbol{z}_1$ under the objective given in the first characterization

$$\begin{aligned}
J(\boldsymbol{w}_1, \boldsymbol{z}_1) &= \frac{1}{N} \sum_n ||\boldsymbol{x}_n - z_{1,n} \boldsymbol{w}_1||^2 \\
&= \frac{1}{N} \sum_n \left[ \boldsymbol{x}_n^T \boldsymbol{x}_n - 2 z_{1,n} \boldsymbol{w}_1^T \boldsymbol{x}_n + z_{1,n}^2 \boldsymbol{w}_1^T \boldsymbol{w}_1 \right] \\
&= \frac{1}{N} \sum_n \left[ \boldsymbol{x}_n^T \boldsymbol{x}_n - 2 z_{1,n} \boldsymbol{w}_1^T \boldsymbol{x}_n + z_{1,n}^2 \right].
\end{aligned}$$

Taking partial derivatives with respect to $z_{1,n}$ yields

$$\begin{aligned}
\frac{\partial J(\boldsymbol{w}_1, \boldsymbol{z}_1)}{\partial z_{1,n}} &= \frac{1}{N} \left[ -2 \boldsymbol{w}_1^T \boldsymbol{x}_n + 2 z_{1,n} \right] = 0 \\
&\implies z_{1,n} = \boldsymbol{w}_1^T \boldsymbol{x}_n.
\end{aligned}$$

Hence the optimal score $z_{1,n}$ is the projection of the data onto the first principal component.

2

Plug this back in to get

$$J(w_1) = \frac{1}{N} \sum_n \left[ x_n^T x_n - 2z_{1,n} w_1^T x_n + z_{1,n}^2 \right]$$

$$= \frac{1}{N} \sum_n \left[ x_n^T x_n - 2z_{1,n}^2 + z_{1,n}^2 \right]$$

$$= \text{const.} - \frac{1}{N} \sum_n z_{1,n}^2$$

$$= \text{const.} - \text{var}[z_1],$$

where the last line follows because

$$\mathbb{E}[z_1] = \mathbb{E}[X^T w]$$

$$= \mathbb{E}[X^T] w$$

$$= 0.$$

---

**Note:** It is tempting to make the following substitution

$$J(w_1) = \frac{1}{N} \sum_n \left[ x_n^T x_n - 2z_{1,n} w_1^T x_n + z_{1,n}^2 \right]$$

$$= \frac{1}{N} \sum_n \left[ x_n^T x_n - 2(w_1^T x_n)(w_1^T x_n) + z_{1,n}^2 \right]$$

$$= \frac{1}{N} \sum_n \left[ x_n^T x_n - 2x_n^T (w_1 w_1^T) x_n) + z_{1,n}^2 \right],$$

and then claim that $w_1 w_1^T = 1$ since $w$ has unit norm, *but this is false!*. Instead, $w_1 w_1^T = \tilde{W}$ is a rank-one matrix with entries $\tilde{W}_{i,j} = w_i w_j$.

---

Returning to our objective function, we see that minimizing $J(w_1)$ yields a first principal component $w_1$ in the direction that maximizes the variance of the projected data. What is this variance? Plugging back in for $z_n$ yields

$$\text{var}[z_1] = \frac{1}{N} \sum_n w_1^T x_n x_n^T w_1 = w_1^T \Sigma w_1,$$

where $\Sigma$ is the empirical covariance matrix. If the magnitude of $w_1$ is unconstrained then the variance is trivially maximized when $||w_1|| \to \infty$, but by the orthonormality of $W$ we have the constraint $||w_1|| = 1$. This can be added to our objective function with a Lagrange multiplier

$$J(w_1) = w_1^T \Sigma w_1 + \lambda_1 (w_1^T w_1 - 1).$$

Taking partial derivatives and setting to zero yields

$$\Sigma w_1 = \lambda_1 w_1.$$

We see that the optimal $w_1$ is an eigenvector of the covariance matrix. Left multiplying by $w_1^T$ we see that the variance of the projected data is the eigenvalue corresponding to $w_1$, so to maximize the variance we choose the eigenvector with the largest $\lambda_1$.

(c) Generative model approach

It would be nice if we could formulate PCA as a generative model and thereby glean some intuition for why the eigenvectors of the empirical covariance matrix are good principal components. If we were to invert the process described above, we would first sample projections $z_n$ and then transform them with $W$ in order to get the corresponding observed data $x_n$. There was no noise in the reconstruction error formulation; the error only stemmed from the fact that we had fewer principal components than dimensions ($L < D$). Hence, in the generative model, $x_n$ would be deterministic given $z_n$ and $W$.

We can relax this assumption a bit and add isotropic Gaussian noise such that

$$x_n \sim \mathcal{N}(Wz_n, \sigma^2 I),$$

where $W$ is still constrained to be orthogonal and where we are ignoring mean shifts for simplicity. If we assume

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

then we have a model known as **probabilistic PCA**. It turns out that in the limit of $\sigma^2 \to 0$ the MLE estimate of $W$ and $z_n$ recovers the classical PCA solution.

2. Factor Analysis

Factor analysis (FA) is another dimensionality reduction technique with a long history in statistics, psychology, and other fields. It turns out that both PCA and FA can be viewed as special cases of the generative model described above. In factor analysis, however, we have the following model:

$$z_n \sim \mathcal{N}(0, I),$$
$$x_n \mid z_n \sim \mathcal{N}(Wz_n, \Psi),$$

where $\Psi$ is restricted to be diagonal. This is just like probabilistic PCA except that here we allow different dimensions of our observations to have different variances. Instead of principal components, we will call the $z_n$'s latent **factors**. The model is best motivated with an example.

Suppose we have a set of of ratings vectors $\{x_n\}_{n=1}^N$ that $N$ users gave to each of $D$ jokes. Our goal is to predict a user's rating on a joke she hasn't heard, based on her ratings for other jokes and other users' ratings for the new joke. This is the **collaborative filtering** problem. We will assume that each user has a latent vector of preferences $z_n$ along $L$ dimensions, e.g. "length of the joke," "knock-knock-iness," or "years of math required to understand the joke." Each joke will have a weighting of these dimensions, $w_d \in \mathbb{R}^L$. Users' ratings are weighted sums, $w_d^T z_n$, of how well the joke aligns to their latent preferences, plus noise. If we knew the preference dimensions, weights, and user preferences, we could easily describe

the distribution over joke ratings in terms of additive noise, but otherwise we have to model the potentially complex distribution over ratings directly.

Take the simplest case where

$$x_n = W z_n + \mu + \text{noise},$$

i.e. $x$ is a linear function of $z$, with $W \in \mathbb{R}^{D \times L}$. Furthermore, assume that both the $z$'s and the noise follow multivariate Gaussian distributions.

$$z_n \sim \mathcal{N}(z_n \mid \mu_0, \Sigma_0)$$
$$x_n \mid z_n \sim \mathcal{N}(W z_n + \mu, \Psi).$$

Since both are Gaussian, their joint distribution will be Gaussian and the marginal distribution over $x$ will be Gaussian.

Our goal is to explain the data in terms of latent variables $z$. In the Gaussian case this means explaining covariance in $x$ in terms of covariance in $z$, rather than in terms of noise. Hence we restrict $\Psi$ to be diagonal.

$$\Psi = \text{diag}\left( \begin{bmatrix} \sigma_{11}^2, & \cdots, & \sigma_{DD}^2 \end{bmatrix} \right).$$

Under this restriction, the model is specified by

$$\theta = \{\mu_0, \Sigma_0, W, \mu, \Psi\}$$
$$= L + \frac{L(L+1)}{2} + LD + D + D \text{ parameters.}$$

This is actually still overparameterized because we can take $\mu_0 = 0$ and $\Sigma_0 = I$ without loss of generality by pushing them into $\mu$ and $W$, respectively. This leaves us with:

$$\theta = \{W, \mu, \Psi\}$$
$$= LD + D + D = O(LD) \text{ parameters.}$$

With a standard trick of augmenting our factors with a constant we can roll the $\mu$ into $W$. By these arguments we have arrived at the generative model for factor analysis posited above.

(a) Learning
In order to learn the parameters $W$ and $\Psi$ and infer the latent factors $z_n$ of the model, we make use of the EM algorithm. In the E step we update the latent factors given the current weights and noise matrices, and in the M step we set the weights and noise matrices to their MAP estimates under the current factors.
One benefit of the EM algorithm is that it is easy to handle missing data. We simply estimate it during the E step.

(b) Unidentifiability
Just as in mixture models, the learned weight matrix is not uniquely defined. We can rotate by any orthogonormal matrix $R$ and the likelihood will not change. Hence one must take care when interpreting the factors.