

## Topic 7: Trending Data Models

INSTRUCTOR: DANIEL L. PIMENTEL-ALARCÓN

© COPYRIGHT 2017

## 7.1 Subspaces

In many applications we observe data  $\mathbf{x} \in \mathbb{R}^D$  according to:

$$\mathbf{x} = \mathbf{U}\mathbf{c}^* + \boldsymbol{\epsilon}, \quad (7.1)$$

where  $\mathbf{U} \in \mathbb{R}^{D \times R}$  is a basis of a known subspace  $\mathcal{U}$ ,  $\boldsymbol{\epsilon}$  represents noise, and  $\mathbf{c}^* \in \mathbb{R}^R$  represents the *noiseless* coefficient of  $\mathbf{x}$  with respect to (w.r.t.)  $\mathbf{U}$ .

**Example 7.1** (Iris flowers dataset). In the Iris Flowers dataset (also known as Fisher's or Anderson's) we can find the width and length (in centimeters) of the sepals and petals of 50 samples (flowers). Let  $\mathbf{x}_i \in \mathbb{R}^4$  be the vector containing the 4 values corresponding to the  $i^{\text{th}}$  flower. It turns out that the columns  $\mathbf{x}_1, \dots, \mathbf{x}_{50}$  lie mostly in a 2-dimensional subspace  $\mathcal{U}$  (see Figure 7.1 to build some intuition). This means that we can model this dataset as

$$\mathbf{x}_i = \mathbf{U}\mathbf{c}_i^* + \boldsymbol{\epsilon}_i$$

where  $\mathbf{U} \in \mathbb{R}^{4 \times 2}$  is a basis of  $\mathcal{U}$ ,  $\mathbf{c}_i^* \in \mathbb{R}^2$  is the coefficient of  $\mathbf{x}_i$  in this basis, and  $\boldsymbol{\epsilon}_i \in \mathbb{R}^4$  represents small noise.

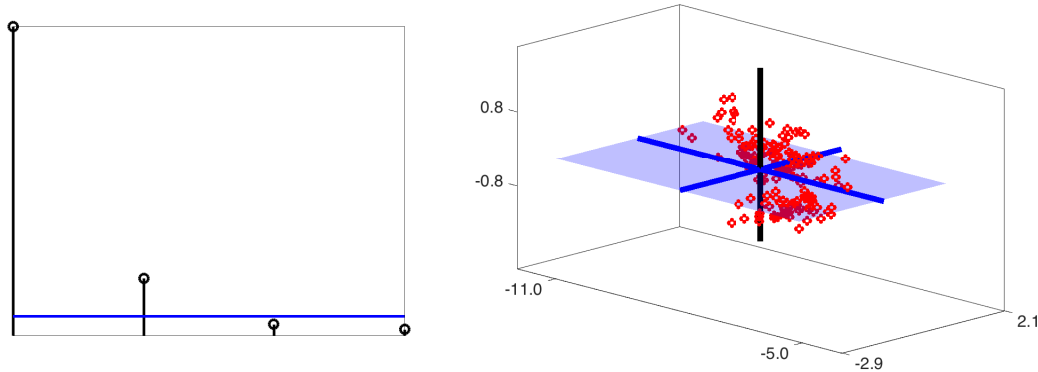


Figure 7.1: **Left:** Singular value decay of Iris Flowers dataset; blue line indicates 5% of the variance in the dataset. **Right:** Iris dataset lies mostly in 2-dimensional subspace. See Example 7.1.

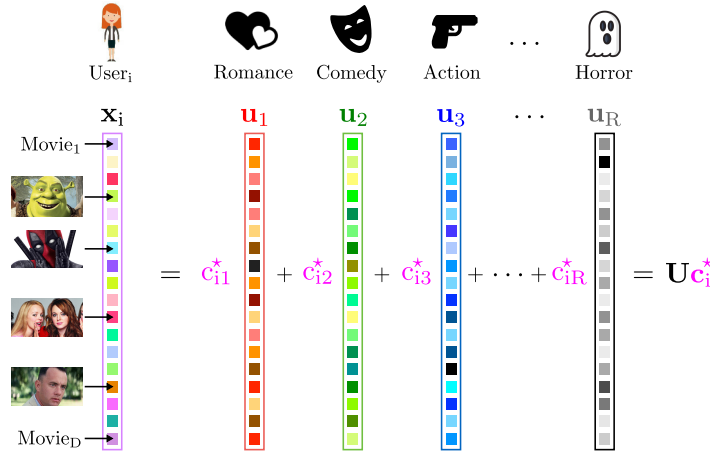


Figure 7.2: See Example 7.2

**Example 7.2** (Netflix and recommender systems). Netflix wants to know which movies you *would* like, so it can make you good recommendations, and you enjoy more your experience.

Consider the collection of all  $D$  movies in Netflix. Each movie has a component of romance, comedy, action, horror, etc. Let  $R \ll D$  be the number of components, and let  $\mathbf{u}_r \in \mathbb{R}^D$  be the vector containing the  $r^{\text{th}}$  component of all movies. For example, the  $j^{\text{th}}$  entry in  $\mathbf{u}_1$  contains the amount of romance in the  $j^{\text{th}}$  movie.

On the other hand, we can think of each Netflix user as a combination of a romance-lover with a comedy-lover with an action-lover, etc. Let  $\mathbf{c}_i^* \in \mathbb{R}^R$  be the vector containing the *coefficients* of the  $i^{\text{th}}$  user. For example, the first entry in  $\mathbf{c}_i^*$  indicates how much user  $i$  likes romance. See Figure 7.2 to build some intuition.

Hence we can model the vector of *preferences*  $\mathbf{x}_i \in \mathbb{R}^D$  indicating how much user  $i$  likes each movie as

$$\mathbf{x}_i = \mathbf{U}\mathbf{c}_i^* + \boldsymbol{\epsilon}_i,$$

where  $\mathbf{U} \in \mathbb{R}^{D \times R}$  is the matrix formed with  $\mathbf{u}_1, \dots, \mathbf{u}_R$  as columns, and  $\boldsymbol{\epsilon}$  represents a small error term. This way, the vectors  $\mathbf{x}_1, \dots, \mathbf{x}_N$  of preferences of all users lie in the  $R$ -dimensional subspace  $\mathcal{U} = \text{span}[\mathbf{U}]$ . If Netflix can find the coefficients  $\mathbf{c}_i^*$ , it can make good predictions about what user  $i$  will like.

There is nothing special about movies. You can change movies for clothes, toys, food, etc., and the *features* {romance, comedy, etc.} for {color, size, etc.}, and you obtain a recommender system for any item. If you make good recommendations, users are more likely to buy. You can see why Amazon, Apple, Yelp, and all these big companies have great interest in this problem and are offering LOTS of money to people who work on this.

One of the first goals in these problems is to find  $\mathbf{c}^*$  based on the *noisy* observation  $\mathbf{x}$ . This can be posed as an estimation problem.

**Example 7.3.** Suppose we observe  $\mathbf{x} \in \mathbb{R}^D$  according to (7.1), where  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ . Then

$$\mathbf{x} \sim \mathcal{N}(\mathbf{U}\mathbf{c}^*, \sigma^2 \mathbf{I}).$$

Our goal is to find  $\mathbf{c}^*$ . Let's use our estimation techniques from before to find the maximum likelihood estimator (MLE):

$$\begin{aligned}
 \hat{\mathbf{c}}_{ML} &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} p(\mathbf{x}|\mathbf{c}) = \arg \max_{\mathbf{c} \in \mathbb{R}^R} \underbrace{\frac{1}{(\sqrt{2\pi}\sigma)^D}}_{\text{constant}} e^{-\frac{1}{2\sigma^2}(\mathbf{x}-\mathbf{Uc})^\top(\mathbf{x}-\mathbf{Uc})} \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} \log e^{-\frac{1}{2\sigma^2}(\mathbf{x}-\mathbf{Uc})^\top(\mathbf{x}-\mathbf{Uc})} = \arg \max_{\mathbf{c} \in \mathbb{R}^R} - \underbrace{\frac{1}{2\sigma^2}}_{\text{constant}} (\mathbf{x}-\mathbf{Uc})^\top(\mathbf{x}-\mathbf{Uc}) \\
 &= \arg \min_{\mathbf{c} \in \mathbb{R}^R} (\mathbf{x}-\mathbf{Uc})^\top(\mathbf{x}-\mathbf{Uc}) = \arg \min_{\mathbf{c} \in \mathbb{R}^R} \mathbf{x}^\top \mathbf{x} - 2\mathbf{c}^\top \mathbf{U}^\top \mathbf{x} + \mathbf{c}^\top \mathbf{U}^\top \mathbf{Uc}. \tag{7.2}
 \end{aligned}$$

We now proceed using elemental techniques from optimization: take derivative, set to zero and solve for the desired variable. To learn more about how to take derivatives w.r.t. vectors and matrices see *Old and new matrix algebra useful for statistics* by Thomas P. Minka. Taking derivative w.r.t.  $\mathbf{c}$  and setting to zero, we obtain the following:

$$-2\mathbf{U}^\top \mathbf{x} + 2\mathbf{U}^\top \mathbf{Uc} = \mathbf{0},$$

and solving for  $\mathbf{c}$ :

$$\hat{\mathbf{c}}_{ML} = (\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{x}.$$

In Example 7.3 we found the MLE of  $\mathbf{c}$  when the noise  $\epsilon$  was gaussian. In general, finding the MLE requires several lucky breaks, like:

- Knowing the distribution of the noise  $\epsilon$ .
- Being able to write down the likelihood explicitly. What if  $\epsilon = \sum_{k=1}^K \epsilon_k$ , where  $\epsilon_k \sim \text{Gamma}(\alpha_k, \beta_k)$  with different  $\beta_k$ 's? Then we know that  $\mathbf{x}$  has the distribution of the sum of Gammas, but this distribution has no known closed form.
- Being able to manipulate the likelihood algebraically so we can solve for its maximum in closed form or at least its derivatives, to use numerical optimization. Plus having a convex likelihood, so it has a unique maximum. For example, if  $\epsilon$  is a mixture of gaussians, the likelihood will generally be non-convex.

However, not everything is lost. Luckily, the model in (7.1) provides a very natural and general way to estimate  $\mathbf{c}^*$  as the coefficients of the projection of  $\mathbf{x}$  onto the subspace  $\mathcal{U} = \text{span}[\mathbf{U}]$ . See Figure 7.3 to build some intuition. Recall that the projection  $\hat{\mathbf{x}}_{\mathcal{U}} \in \mathcal{U}$  of  $\mathbf{x}$  is defined as

$$\hat{\mathbf{x}}_{\mathcal{U}} := \mathbf{x} \in \mathcal{U} \text{ such that } \|\mathbf{x} - \hat{\mathbf{x}}\| \leq \|\mathbf{x} - \tilde{\mathbf{x}}\| \text{ for every } \tilde{\mathbf{x}} \in \mathcal{U},$$

and the coefficients  $\hat{\mathbf{c}}_{\mathcal{U}}$  of  $\hat{\mathbf{x}}_{\mathcal{U}}$  w.r.t.  $\mathbf{U}$  are simply

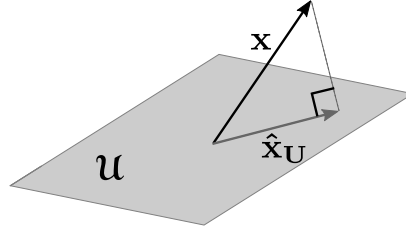
$$\hat{\mathbf{c}}_{\mathcal{U}} := \mathbf{c} \in \mathbb{R}^R \text{ such that } \hat{\mathbf{x}}_{\mathcal{U}} = \mathbf{Uc}.$$

Luckily, we can find  $\hat{\mathbf{c}}_{\mathcal{U}}$  directly as:

$$\hat{\mathbf{c}}_{\mathcal{U}} = \arg \min_{\mathbf{c} \in \mathbb{R}^R} \|\mathbf{x} - \mathbf{Uc}\|_2^2 = \arg \min_{\mathbf{c} \in \mathbb{R}^R} (\mathbf{x} - \mathbf{Uc})^\top (\mathbf{x} - \mathbf{Uc}).$$

Notice that this is the same as (7.2), and hence

$$\hat{\mathbf{c}}_{\mathcal{U}} = (\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{x}.$$

Figure 7.3: Projection  $\hat{\mathbf{x}}_U$  of vector  $\mathbf{x}$  onto subspace  $U$ .

Notice that  $\hat{\mathbf{c}}_U$  also provides an estimate of  $\mathbf{x}^* = U\mathbf{c}^*$ , the *noiseless* counterpart of  $\mathbf{x}$ , via

$$\hat{\mathbf{x}}_U = U\hat{\mathbf{c}}_U = U(U^T U)^{-1} U^T \mathbf{x},$$

where we recognize  $\mathbf{P}_U = U(U^T U)^{-1} U^T \in \mathbb{R}^{D \times D}$  to be the *projector operator* onto  $U$ .

**Example 7.4** (Linear regression). In linear regression we want to infer how a variable  $y$  (e.g., glucose or cholesterol level) depends on an other collection of variables  $\mathbf{x} \in \mathbb{R}^M$  (e.g., diet, weight, income, education, sex, age) based on a sequence of *training* examples  $y_1, \dots, y_N$ . The main idea is to model  $y_i$  as a linear combination of  $\mathbf{x}_i$ , i.e.,

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta}^* + \epsilon_i,$$

where  $\epsilon_i$  models noise and the goal is to find  $\boldsymbol{\beta}^* \in \mathbb{R}^M$ , which determines the dependency of  $y$  on  $\mathbf{x}$ . Equivalently, in vector form we can write

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\epsilon},$$

where  $\mathbf{y} = [y_1 \dots y_N]^T$ ,  $\mathbf{X}$  is the  $N \times M$  matrix formed with  $\{\mathbf{x}_i^T\}_{i=1}^N$  as rows and  $\boldsymbol{\epsilon} = [\epsilon_1 \dots \epsilon_N]^T$ . Notice that this is equivalent to (7.1), only with slightly different notation:  $\mathbf{y} \equiv \mathbf{x}$ ,  $\mathbf{X} \equiv U$ ,  $\boldsymbol{\beta}^* \equiv \mathbf{c}^*$ . It follows from our previous discussion that

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y},$$

which is the well-known result from linear regression.

**Remark 7.1.** We point out that unlike  $\hat{\mathbf{c}}_{ML}$ , the form of  $\hat{\mathbf{c}}_U$  does not depend on the distribution of  $\boldsymbol{\epsilon}$ . In the case where  $\boldsymbol{\epsilon}$  is gaussian, these two estimators  $\hat{\mathbf{c}}_{ML}$  and  $\hat{\mathbf{c}}_U$  happen to be the same. However, in general this is not the case.

## 7.2 Sparsity

In the previous section we learned how to find the  $\mathbf{c}$  that minimizes the error  $\|\mathbf{x} - U\mathbf{c}\|_2^2$ . In many modern applications we have prior information that  $\mathbf{c}$  should be *sparse*.



Figure 7.4: See Example 7.5. **Left:** Many astronomy images, like the Magellanic Clouds depicted here, are sparse. **Right:** The Wavelet transform of an MRI image is very sparse.

**Example 7.5.** Let  $\mathbf{C} \in \mathbb{R}^{\sqrt{D} \times \sqrt{D}}$  be a *sparse* image. For example an astronomy image, or the Wavelet transform of an MRI image (see Figure 7.4). In many applications, rather than observing  $\mathbf{C}$ , we observe a *transformation* of  $\mathbf{C}$ . For example, telescopes will produce the Fourier transform of  $\mathbf{C}$ . This is equivalent to observing  $\mathbf{x}$  as in (7.1), where  $\mathbf{c} \in \mathbb{R}^D$  is the *vectorized* version of  $\mathbf{C}$  (i.e., the vector formed by stacking the columns in  $\mathbf{C}$ ), and  $\mathbf{U}$  determines the transformation. In cases like these, we know *a priori* that the desired solution  $\mathbf{c}$  should be sparse.

We saw before that bayesian approaches *bias* our solutions towards our *prior* information. In this case, our prior information is that the solution should be sparse, so we will use a prior distribution on  $\mathbf{c}$  that reflects this. One such popular distribution is the Laplacian distribution:

$$p(\mathbf{c}) = \prod_{r=1}^R \frac{\lambda}{2} e^{-\lambda |c_r|} = \left(\frac{\lambda}{2}\right)^R e^{-\lambda \mathbf{c}^T \text{sign}(\mathbf{c})}.$$

Because of the shape of this density, a vector  $\mathbf{c}$  with i.i.d. Laplacian entries is likely to have a lot of small values and a few large ones. This way, a Laplacian prior *favors* sparse solutions.

Notice that the Laplacian density has much heavier tails than the gaussian (see Figure 7.5; the Laplacian's exponential decay is  $e^{-|c_r|}$ , while the gaussian's is  $e^{-c_r^2}$ ). This allows the Laplacian to produce *some* deviated values, i.e., mostly small values, and a few large ones. In contrast, since the gaussian density has lighter tails, is very unlikely to produce deviated values. A vector  $\mathbf{c}$  with i.i.d.  $\mathcal{N}(0, \sigma^2)$  entries will either have most entries be zero (if  $\sigma$  is small) or most entries be large (if  $\sigma$  is large). Hence a gaussian is not a good *prior* to model sparsity. However, gaussian *likelihoods* are good to model noise. Luckily, as we will see in the next example, we can mix and match: we can use a Laplacian *prior* together with a gaussian *likelihood* in a bayesian setting in order to produce a sparse solution.

**Example 7.6** (Laplacian prior). Suppose we observe  $\mathbf{x}$  as in (7.1), with  $R = D$ ,  $\mathbf{U}$  orthonormal,  $\mathbf{c}^*$  sparse, and  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . Then

$$\mathbf{x} \sim p(\mathbf{x}|\mathbf{c}) = \mathcal{N}(\mathbf{U}\mathbf{c}^*, \sigma^2 \mathbf{I}),$$

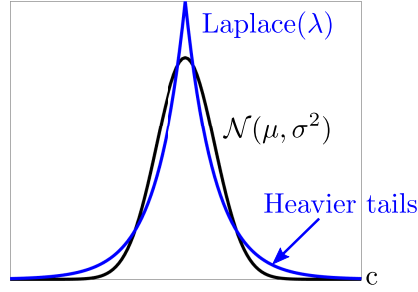


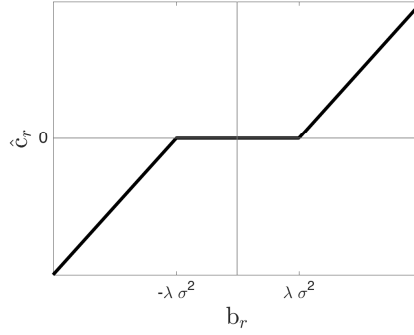
Figure 7.5: The Laplacian density is a good model for sparsity priors, as it would produce mostly small values, yet it would also produce a few large ones. This is because there is enough probability in its *heavy* tails. In contrast, the gaussian density is very unlikely to produce large values (deviated from the mean).

and we want to estimate  $\mathbf{c}^*$ . Since we know  $\mathbf{c}^*$  is sparse, we can induce a Laplacian prior. Then our maximum a posteriori (MAP) estimator is

$$\begin{aligned}
 \hat{\mathbf{c}} &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} p(\mathbf{c}|\mathbf{x}) = \arg \max_{\mathbf{c} \in \mathbb{R}^R} \log(p(\mathbf{x}|\mathbf{c})p(\mathbf{c})) \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} \log \left( \frac{1}{\sqrt{(2\pi)^D}} e^{-\frac{1}{2\sigma^2}(\mathbf{x}-\mathbf{Uc})^\top(\mathbf{x}-\mathbf{Uc})} \cdot \left(\frac{\lambda}{2}\right)^R e^{-\lambda \mathbf{c}^\top \text{sign}(\mathbf{c})} \right) \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} \log e^{-\frac{1}{2\sigma^2}(\mathbf{x}-\mathbf{Uc})^\top(\mathbf{x}-\mathbf{Uc}) - \lambda \mathbf{c}^\top \text{sign}(\mathbf{c})} + \underbrace{\log \frac{\lambda^R}{2^R \sqrt{(2\pi)^D}}}_{\text{constant } \kappa} \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} -\frac{1}{2\sigma^2}(\mathbf{x}-\mathbf{Uc})^\top(\mathbf{x}-\mathbf{Uc}) - \lambda \mathbf{c}^\top \text{sign}(\mathbf{c}) + \kappa \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} -\frac{1}{2\sigma^2}(\mathbf{x}^\top \mathbf{x} - 2\mathbf{c}^\top \mathbf{U}^\top \mathbf{x} + \mathbf{c}^\top \underbrace{\mathbf{U}^\top \mathbf{U}}_{\mathbf{I}} \mathbf{c}) - \lambda \mathbf{c}^\top \text{sign}(\mathbf{c}) + \kappa \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} -\frac{1}{2\sigma^2}(\mathbf{x}^\top \mathbf{x} - 2\mathbf{c}^\top \mathbf{U}^\top \mathbf{x} + \mathbf{c}^\top \mathbf{c}) - \lambda \mathbf{c}^\top \text{sign}(\mathbf{c}) + \kappa \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} -\frac{1}{2\sigma^2}(\mathbf{x}^\top \underbrace{\mathbf{U}\mathbf{U}^\top}_{\mathbf{I}} \mathbf{x} - 2\mathbf{c}^\top \underbrace{\mathbf{U}^\top \mathbf{x}}_{=: \mathbf{b}} + \mathbf{c}^\top \mathbf{c}) - \lambda \mathbf{c}^\top \text{sign}(\mathbf{c}) + \kappa \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} -\frac{1}{2\sigma^2}(\mathbf{b}^\top \mathbf{b} - 2\mathbf{c}^\top \mathbf{b} + \mathbf{c}^\top \mathbf{c}) - \lambda \mathbf{c}^\top \text{sign}(\mathbf{c}) + \kappa \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} -\frac{1}{2\sigma^2}(\mathbf{b} - \mathbf{c})^\top (\mathbf{b} - \mathbf{c}) - \lambda \mathbf{c}^\top \text{sign}(\mathbf{c}) + \kappa \\
 &= \arg \max_{\mathbf{c} \in \mathbb{R}^R} \sum_{r=1}^R \left( -\frac{1}{2\sigma^2}(\mathbf{b}_r - \mathbf{c}_r)^2 - \lambda |\mathbf{c}_r| \right) + \kappa. \tag{7.3}
 \end{aligned}$$

Next we need to study two cases:

- (i) If  $\mathbf{c}_r \neq 0$ , the derivative of (7.3) w.r.t.  $\mathbf{c}_r$  is  $\frac{1}{\sigma^2}(\mathbf{b}_r - \mathbf{c}_r) - \lambda \text{sign}(\mathbf{c}_r)$ . Setting to zero and solving for  $\mathbf{c}_r$  we obtain  $\mathbf{c}_r = \mathbf{b}_r - \lambda \sigma^2 \text{sign}(\mathbf{c}_r)$ . Notice that if  $\text{sign}(\mathbf{c}_r) \neq \text{sign}(\mathbf{b}_r)$ , then  $(\mathbf{b}_r - \mathbf{c}_r)$  gets bigger, which is exactly the opposite of what we want. Thus we can replace  $\text{sign}(\mathbf{c}_r)$  with  $\text{sign}(\mathbf{b}_r)$  to

Figure 7.6: Soft-threshold function  $\hat{c}_r = \text{sign}(b_r) \max(|b_r| - \lambda\sigma^2, 0)$ . See Example 7.6.

obtain:

$$c_r = b_r - \lambda\sigma^2 \text{sign}(b_r).$$

Hence, if  $c_r \neq 0$ , the *cost* associated to the  $r^{\text{th}}$  term is

$$\begin{aligned} \left( \frac{1}{2\sigma^2} (b_r - c_r)^2 + \lambda |c_r| \right) \Big|_{c_r = b_r - \lambda\sigma^2 \text{sign}(b_r)} &= \frac{1}{2\sigma^2} \underbrace{(b_r - b_r + \lambda\sigma^2 \text{sign}(b_r))^2}_{\lambda^2 \sigma^4} + \lambda |b_r - \lambda\sigma^2 \text{sign}(b_r)| \\ &= \frac{\lambda^2 \sigma^2}{2} + \lambda |b_r - \lambda\sigma^2 \text{sign}(b_r)|. \end{aligned} \quad (7.4)$$

(ii) On the other hand, if  $c_r = 0$ , then the *cost* associated to the  $r^{\text{th}}$  term is

$$\left( \frac{1}{2\sigma^2} (b_r - c_r)^2 + \lambda |c_r| \right) \Big|_{c_r=0} = \frac{b_r^2}{2\sigma^2}. \quad (7.5)$$

Observing that

$$\begin{aligned} (7.4) &\geq (7.5) && \text{if } |b_r| \leq \lambda\sigma^2, \\ (7.4) &< (7.5) && \text{if } |b_r| > \lambda\sigma^2 \end{aligned}$$

it follows that the minimizer is

$$\hat{c}_r = \begin{cases} 0 & \text{if } |b_r| \leq \lambda\sigma^2 \\ b_r - \lambda\sigma^2 \text{sign}(b_r) & \text{if } |b_r| > \lambda\sigma^2 \end{cases} = \text{sign}(b_r) \max(|b_r| - \lambda\sigma^2, 0),$$

which is often known as the *soft-threshold* function (see Figure 7.6).

### 7.3 Mixtures and the EM algorithm

Previously we talked about *mixtures*, where

$$x_i \stackrel{iid}{\sim} p(x|\theta) = \sum_{k=1}^K \rho_k p_k(x|\theta_k^*),$$

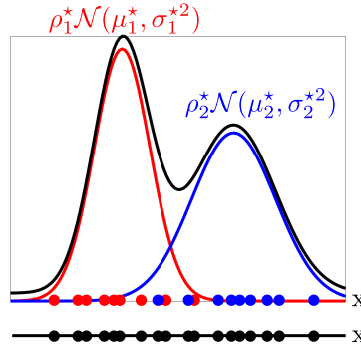


Figure 7.7: In a *mixture* each sample  $x_i$  is drawn according to one out of  $K$  distributions  $p_k(x|\theta_k^*)$ , and we want to estimate  $\theta^* = [\theta_1^* \dots \theta_K^*]^\top$ . The challenge is that we don't know which distribution generated each  $x_i$ , and so we don't know which  $x_i$ 's should be used to estimate  $\theta_k^*$ . In this illustration we only observe the *black* points, which are drawn according to  $p_k(x|\theta_k^*) = \mathcal{N}(\mu_k^*, \sigma_k^{*2})$ . If we knew which point is red and which one is blue, we could use only the red to estimate  $\{\mu_1^*, \sigma_1^{*2}\}$  and only the blue to estimate  $\{\mu_2^*, \sigma_2^{*2}\}$ . The challenge is that we don't know the colors, and so we also have to estimate which distribution each point belongs to. Here  $\rho_k^* \in [0, 1]$  models the fraction of points that correspond to the  $k^{\text{th}}$  distribution, which also has to be estimated.

that is,  $x_i \stackrel{iid}{\sim} p_k(x|\theta_k^*)$  with probability  $\rho_k^*$ ,  $k = 1, \dots, K$ . See Figure 7.7 to build some intuition. The goal is to estimate  $\theta^* = [\theta_1^* \dots \theta_K^*]^\top$ . However, the MLE is not so easy to compute, and taking log does not simplify things because the log of a sum does not factor nicely. Furthermore,  $p(x|\theta)$  is not even convex. This means that there may be multiple solutions (local maximums). Estimating mixtures is an active and challenging field of research.

There are several computational (iterative) methods to try to solve this problem, like the well-known gradient descent, or Newton's method. An other option is to use the *expectation-maximization* (EM) algorithm, which we will focus on. The intuition behind EM is that often things would greatly simplify if we had *side* information. For instance, consider the additional variable

$$z_i^* = k \quad \text{if } x_i \sim p_k(x|\theta_k^*).$$

In words,  $z_i^*$  indicates which of the  $K$  distributions the  $i^{\text{th}}$  sample corresponds to. We don't know  $\mathbf{z}^* = [z_1^* \dots z_N^*]^\top$ , but if we did, things would greatly simplify, because we could *split* our data into

$$\mathbf{X}_k^* = \{x_i : z_i^* = k\}, \quad k = 1, \dots, K,$$

and then compute the MLE of each  $\theta_k^*$  using only  $\mathbf{X}_k^*$ .

**Example 7.7.** Consider a mixture of gaussians, i.e.,  $x_1, \dots, x_N$  drawn i.i.d. according to

$$p(x|\theta) = \sum_{k=1}^K \rho_k^* \frac{1}{\sqrt{2\pi\sigma_k^*}} e^{-\frac{1}{2} \left( \frac{x - \mu_k^*}{\sigma_k^*} \right)^2}.$$

This time our usual tricks from elementary optimization (taking derivatives and setting to zero) won't work, because the log term won't factor nicely. However, if we knew  $\mathbf{z}$ , we would be able to compute



the MLE of each  $\theta_k^*$  using only  $\mathbf{X}_k^*$  to obtain

$$\hat{\rho}_k = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{z_i^*=k\}}, \quad \hat{\mu}_k = \frac{1}{|\mathbf{X}_k^*|} \sum_{x_i \in \mathbf{X}_k^*} x_i, \quad \hat{\sigma}_k^2 = \frac{1}{|\mathbf{X}_k^*|} \sum_{x_i \in \mathbf{X}_k^*} (x_i - \hat{\mu}_k)^2.$$

With this in mind, the EM algorithm incorporates the *missing data*  $\mathbf{z}$  in the *complete-data* log-likelihood:

$$\log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}).$$

Recall that we don't know  $\mathbf{z}^*$ , so the idea is to average  $\log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$  over  $\mathbf{z}$  to obtain the *expected* complete-data log-likelihood

$$q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_t) := E_{\mathbf{z} | \mathbf{x}, \hat{\boldsymbol{\theta}}_t} [\log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})],$$

where the expectation is taken given  $\mathbf{x}$  and the current estimate  $\hat{\boldsymbol{\theta}}_t$ . Intuitively,  $q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_t)$  is the *simplified* likelihood that is easy to maximize. The EM algorithm, summarized in Algorithm 1, essentially consists on iteratively maximizing  $q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_t)$ .

---

**Algorithm 1:** Expectation-Maximization (EM)

---

**Input:** Data  $x_1, \dots, x_N$ , number of gaussians  $K$ .

**Initialize**  $\hat{\boldsymbol{\theta}}_t$ , for example randomly.

**Repeat until convergence**

- **E-step.** Compute the expected complete log-likelihood:

$$q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_t) = E_{\mathbf{z} | \mathbf{x}, \hat{\boldsymbol{\theta}}_t} [\log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})].$$

- **M-step.** Maximize the expected complete log-likelihood:

$$\hat{\boldsymbol{\theta}}_{t+1} := \arg \max_{\boldsymbol{\theta}} q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_t).$$

**Output:**  $\hat{\boldsymbol{\theta}}_{t+1}$ .

---

The EM algorithm is attractive when  $q$  is easily computed and maximized.

**Example 7.8.** Consider a mixture of gaussians, as in Example 7.7. We can write  $p(x_i, z_i | \boldsymbol{\theta})$  as

$$p(x_i, z_i | \boldsymbol{\theta}) = p(x_i | z_i, \boldsymbol{\theta}) p(z_i | \boldsymbol{\theta}) = \prod_{k=1}^K \left( \rho_k \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2} \left( \frac{x_i - \mu_k}{\sigma_k} \right)^2} \right)^{\mathbb{1}_{\{z_i=k\}}}.$$

It follows that

$$\begin{aligned}
 \log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) &= \log \prod_{i=1}^N \prod_{k=1}^K \left( \rho_k \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2} \left( \frac{x_i - \mu_k}{\sigma_k} \right)^2} \right)^{\mathbb{1}_{\{z_i=k\}}} \\
 &= \sum_{i=1}^N \sum_{k=1}^K \log \left( \rho_k \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2} \left( \frac{x_i - \mu_k}{\sigma_k} \right)^2} \right)^{\mathbb{1}_{\{z_i=k\}}} \\
 &= \sum_{i=1}^N \sum_{k=1}^K \mathbb{1}_{\{z_i=k\}} \log \rho_k \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2} \left( \frac{x_i - \mu_k}{\sigma_k} \right)^2}.
 \end{aligned}$$

Taking expectation we have:

$$q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_t) = E_{z_i | x_i, \hat{\boldsymbol{\theta}}_t} [p(x_i, z_i | \boldsymbol{\theta})] = \sum_{i=1}^N \sum_{k=1}^K E_{z_i | x_i, \hat{\boldsymbol{\theta}}_t} [\mathbb{1}_{\{z_i=k\}}] \log \rho_k \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2} \left( \frac{x_i - \mu_k}{\sigma_k} \right)^2},$$

where

$$E_{z_i | x_i, \hat{\boldsymbol{\theta}}_t} [\mathbb{1}_{\{z_i=k\}}] = p(z_i = k | x_i, \hat{\boldsymbol{\theta}}_t) = \frac{p(x_i | z_i = k, \hat{\boldsymbol{\theta}}_t) p(z_i = k | \hat{\boldsymbol{\theta}}_t)}{p(x_i | \hat{\boldsymbol{\theta}}_t)} = \frac{\hat{\rho}_{k_t} \frac{1}{\sqrt{2\pi}\hat{\sigma}_{k_t}} e^{-\frac{1}{2} \left( \frac{x_i - \hat{\mu}_{k_t}}{\hat{\sigma}_{k_t}} \right)^2}}{\sum_{k=1}^K \hat{\rho}_{k_t} \frac{1}{\sqrt{2\pi}\hat{\sigma}_{k_t}} e^{-\frac{1}{2} \left( \frac{x_i - \hat{\mu}_{k_t}}{\hat{\sigma}_{k_t}} \right)^2}}$$

can be easily computed. At this point we have already completed the E-step. To obtain  $\hat{\boldsymbol{\theta}}_{t+1} = \arg \max_{\boldsymbol{\theta}} q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_t)$  in the M-step, we can use our usual tricks: take derivative, set to zero, and solve for the maximizer, which yields:

$$\begin{aligned}
 \hat{\rho}_{k_{t+1}} &= \frac{1}{N} \sum_{i=1}^N p(z_i = k | x_i, \hat{\boldsymbol{\theta}}_t), \\
 \hat{\mu}_{k_{t+1}} &= \frac{1}{N \hat{\rho}_{k_{t+1}}} \sum_{i=1}^N p(z_i = k | x_i, \hat{\boldsymbol{\theta}}_t) x_i, \\
 \hat{\sigma}_{k_{t+1}}^2 &= \frac{1}{N \hat{\rho}_{k_{t+1}}} \sum_{i=1}^N p(z_i = k | x_i, \hat{\boldsymbol{\theta}}_t) (x_i - \hat{\mu}_{k_{t+1}})^2.
 \end{aligned}$$

**Exercise 7.1.** Derive the EM algorithm for a mixture of *multivariate* gaussians.

**Remark 7.2.** The EM algorithm is useful only if  $q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}_t)$  is easy to compute and maximize. Otherwise we would need to maximize an ugly function, which is what we are trying to avoid. Moreover, we would need to do so at each step! Also, determining which side information to use to simplify the problem involves a bit of art, and needs to be tailored case by case. Luckily, there are many cases where this can be done.

The EM algorithm is widely used to solve non-convex problems with missing data (in our case, the *missing*

data is  $\mathbf{z}$ ). One of the reasons of EM's great success is that its estimates are guaranteed to get better with each iteration.

**Proposition 7.1** (EM-convergence). Let  $\hat{\boldsymbol{\theta}}_t$  be the estimate produced by the EM algorithm at time  $t$ . Then

$$p(\mathbf{x}|\hat{\boldsymbol{\theta}}_{t+1}) \geq p(\mathbf{x}|\hat{\boldsymbol{\theta}}_t).$$

*Proof.* We want to show that  $p(\mathbf{x}|\hat{\boldsymbol{\theta}}_{t+1}) \geq p(\mathbf{x}|\hat{\boldsymbol{\theta}}_t)$ , or equivalently, that

$$\log p(\mathbf{x}|\hat{\boldsymbol{\theta}}_{t+1}) - \log p(\mathbf{x}|\hat{\boldsymbol{\theta}}_t) \geq 0.$$

We know from elementary probability that  $p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) = p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x}|\boldsymbol{\theta})$ , so

$$\begin{aligned} \log p(\mathbf{x}|\hat{\boldsymbol{\theta}}_{t+1}) - \log p(\mathbf{x}|\hat{\boldsymbol{\theta}}_t) &= \log p(\mathbf{x}, \mathbf{z}|\hat{\boldsymbol{\theta}}_{t+1}) - \log p(\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_{t+1}) - \log p(\mathbf{x}, \mathbf{z}|\hat{\boldsymbol{\theta}}_t) + \log p(\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_t) \\ &= \log p(\mathbf{x}, \mathbf{z}|\hat{\boldsymbol{\theta}}_{t+1}) - \log p(\mathbf{x}, \mathbf{z}|\hat{\boldsymbol{\theta}}_t) + \log \frac{p(\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_t)}{p(\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_{t+1})}. \end{aligned}$$

Next notice that  $\log p(\mathbf{x}|\boldsymbol{\theta})$  is not a function of  $\mathbf{z}$ , so taking expectations on both sides,

$$\begin{aligned} \log p(\mathbf{x}|\hat{\boldsymbol{\theta}}_{t+1}) - \log p(\mathbf{x}|\hat{\boldsymbol{\theta}}_t) &= \underbrace{\mathbb{E}_{\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_t} [\log p(\mathbf{x}, \mathbf{z}|\hat{\boldsymbol{\theta}}_{t+1})]}_{q(\hat{\boldsymbol{\theta}}_{t+1}|\hat{\boldsymbol{\theta}}_t)} - \underbrace{\mathbb{E}_{\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_t} [\log p(\mathbf{x}, \mathbf{z}|\hat{\boldsymbol{\theta}}_t)]}_{q(\hat{\boldsymbol{\theta}}_t|\hat{\boldsymbol{\theta}}_t)} + \underbrace{\mathbb{E}_{\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_t} \left[ \log \frac{p(\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_t)}{p(\mathbf{z}|\mathbf{x}, \hat{\boldsymbol{\theta}}_{t+1})} \right]}_{\text{KL divergence} \geq 0} \\ &\geq q(\hat{\boldsymbol{\theta}}_{t+1}|\hat{\boldsymbol{\theta}}_t) - q(\hat{\boldsymbol{\theta}}_t|\hat{\boldsymbol{\theta}}_t) \geq 0, \end{aligned}$$

where the last inequality follows because  $\hat{\boldsymbol{\theta}}_{t+1} = \arg \max_{\boldsymbol{\theta}} q(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}_t)$  by definition, so

$$q(\hat{\boldsymbol{\theta}}_{t+1}|\hat{\boldsymbol{\theta}}_t) = \max_{\boldsymbol{\theta}} q(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}_t) \geq q(\hat{\boldsymbol{\theta}}_t|\hat{\boldsymbol{\theta}}_t).$$

□

## 7.4 Subspace Clustering

In subspace clustering we observe column vectors  $\mathbf{x}_1, \dots, \mathbf{x}_N$  lying in the union of  $K$   $R$ -dimensional subspaces of  $\mathbb{R}^D$ ,  $\mathbf{u}_1^*, \dots, \mathbf{u}_K^*$ . The goal is to estimate the underlying subspaces from  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , and to cluster the columns according to the subspaces. See Figure ?? to build some intuition.

This is a problem with active research that has attracted a lot of attention in recent years due to its wide range of applications, including computer vision, network inference and recommenders systems, among others []. We will now analyze it in light of the topics we discussed before.

### 7.4.1 Subspace Clustering as a Mixture

First, subspace clustering can be posed as a mixture problem

$$\mathbf{x}_i \stackrel{iid}{\sim} \sum_{k=1}^K \rho_k \mathcal{N}(\mathbf{0}, \mathbf{U}_k \mathbf{U}_k^T),$$

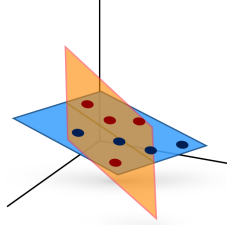


Figure 7.8: Some vectors (blue) lie in one subspace, other vectors (red) lie in an other subspace.

where  $\mathbf{U}_k \in \mathbb{R}^{D \times R}$  is a basis of  $\mathcal{U}_k$ .

**Exercise 7.2.** Derive an EM algorithm for subspace clustering.

### 7.4.2 Subspace Clustering using Sparsity

We will now demonstrate the power of sparsity for clustering applications. In particular we will discuss the main principle behind the state-of-the-art *sparse subspace clustering* algorithm (SSC) [1]. The key idea is to try to write each  $\mathbf{x}_i$  as a linear combination of *other* columns from the dataset, i.e.,

$$\mathbf{x}_i = \mathbf{X}_i \mathbf{c}_i, \quad (7.6)$$

where  $\mathbf{X}_i$  is the  $D \times (N - 1)$  matrix containing all the columns except  $\mathbf{x}_i$ , and  $\mathbf{c}_i \in \mathbb{R}^{N-1}$  is the vector of *coefficients* of  $\mathbf{x}_i$  with respect to  $\mathbf{X}_i$ . In general, there may be infinitely many vectors  $\mathbf{c}_i$  satisfying (7.6). However, if  $\mathbf{x}_i \in \mathcal{U}_k$ , and  $\mathbf{X}_i$  contains a basis of  $\mathcal{U}_k$ , then we can write  $\mathbf{x}_i$  as a *sparse* linear combination of the columns in  $\mathbf{X}_i$ . For example, if  $i = N$  and the first  $R$  columns of  $\mathbf{X}_i$  span  $\mathcal{U}_k$ , then we can write

$$\mathbf{x}_i = \mathbf{X}_i \mathbf{c}_i^* = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_R & \mathbf{x}_{R+1} & \cdots & \mathbf{x}_{N-1} \end{bmatrix} \begin{bmatrix} c_{i1}^* \\ \vdots \\ c_{iR}^* \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (7.7)$$

for some  $[c_{i1}^* \cdots c_{iR}^*]^T$ , containing the coefficients of  $\mathbf{x}_i$  with respect to the basis  $[\mathbf{x}_1 \dots \mathbf{x}_R]$ .

Notice that if  $\mathbf{X}_i$  contains multiple basis of  $\mathcal{U}_k$  (which is generally the case in subspace clustering applications), there may be multiple *sparse representations* of  $\mathbf{x}_i$ . Continuing with our example, if  $[\mathbf{x}_2 \dots \mathbf{x}_{R+1}]$  are *also* a basis of  $\mathcal{U}_k$ , then we could *also* write

$$\mathbf{x}_i = \mathbf{X}_i \mathbf{c}_i^* = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_{R+1} & \mathbf{x}_{R+2} & \cdots & \mathbf{x}_{N-1} \end{bmatrix} \begin{bmatrix} 0 \\ c_{i2}^* \\ \vdots \\ c_{iR+1}^* \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

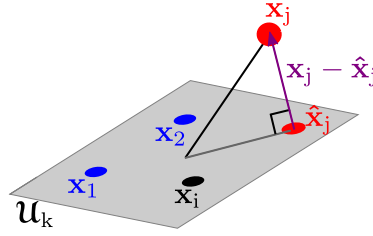


Figure 7.9: Suppose we find that  $\mathbf{x}_i$  can be written as a linear combination of  $\mathbf{x}_1, \mathbf{x}_2$ , as in (7.7). This means that  $\mathbf{x}_i, \mathbf{x}_1, \mathbf{x}_2$  lie in the same subspace  $\mathcal{U}_k$ . But more importantly, this means that  $\mathbf{x}_1, \mathbf{x}_2$  are a basis of  $\mathcal{U}_k$ . In other words, we have already identified  $\mathcal{U}_k$ . We can then find *all* the columns that lie in  $\mathcal{U}_k$  by simply projecting onto  $\mathcal{U}_k$  and checking residuals. If a column lies in  $\mathcal{U}_k$ , like  $\mathbf{x}_i$ , then its residual  $\mathbf{x}_i - \hat{\mathbf{x}}_i$  will be zero. If a column does not lie in  $\mathcal{U}_k$ , like  $\mathbf{x}_j$ , its residual  $\mathbf{x}_j - \hat{\mathbf{x}}_j$  will be nonzero.

for some *other*  $[c_{i_2}^* \cdots c_{i_{R+1}}^*]^\top$ , containing the coefficients of  $\mathbf{x}_i$  with respect to the basis  $[\mathbf{x}_2 \cdots \mathbf{x}_{R+1}]$ .

Luckily, regardless of the particular *sparse*  $\mathbf{c}_i^*$ , the nonzero coefficients of  $\mathbf{c}_i^*$  determine  $R$  columns in  $\mathbf{X}_i$  that correspond to the same subspace as  $\mathbf{x}_i$ . But more importantly, these  $R$  columns form a basis of  $\mathcal{U}_k$ . In other words, we have already identified  $\mathcal{U}_k$ , and hence we can identify *all* the columns in  $\mathbf{X}$  that lie in  $\mathcal{U}_k$  (by simply projecting onto  $\mathcal{U}_k$ , and checking residuals; see Figure 7.9 to build some intuition). Hence the goal is to find a *sparse*  $\mathbf{c}_i^*$ . After that, we can *prune*  $\mathbf{X}$  (i.e., remove all the columns corresponding to  $\mathcal{U}_k$ ) and repeat this procedure to find the columns corresponding to each of the subspaces (which is what we ultimately want).

So the goal is to find a *sparse*  $\mathbf{c}_i^*$  such that  $\mathbf{x}_i = \mathbf{X}_i \mathbf{c}_i^*$ . Ideally, this could be done by solving:

$$\hat{\mathbf{c}}_i = \arg \min_{\mathbf{c}_i: \|\mathbf{c}_i\|_0=R} \|\mathbf{x}_i - \mathbf{X}_i \mathbf{c}_i\|,$$

or equivalently,

$$\hat{\mathbf{c}}_i = \arg \min_{\mathbf{c}_i} \|\mathbf{x}_i - \mathbf{X}_i \mathbf{c}_i\| + \|\mathbf{c}_i\|_0, \quad (7.8)$$

where  $\|\cdot\|_0$  is the so-called  $\ell_0$ -norm (which is actually *not* a norm! Can you show this?), defined as the number of nonzero entries in a vector. Sadly, the  $\ell_0$ -norm is non-convex, and minimizing it is NP-hard.

Whenever one wants to minimize an ugly, non-convex function  $f$  (as in (7.8)), it is often practical to minimize *the next best convex thing* instead, i.e., minimize an other function  $g$  that is both convex and similar to  $f$ . This is usually called *convex relaxation*. In our case,  $g = \ell_1$ -norm (defined as the sum of the absolute values of the entries in a vector) is both convex and similar to  $f = \ell_0$ -norm. Hence, instead of (7.8), we can solve

$$\hat{\mathbf{c}}_i = \arg \min_{\mathbf{c}_i} \|\mathbf{x}_i - \mathbf{X}_i \mathbf{c}_i\| + \lambda \|\mathbf{c}_i\|_1, \quad (7.9)$$

where  $\lambda$  is a parameter that needs to be tuned to achieve the desired level of sparsity. Sadly, (7.9) has no closed-form solution, but luckily it is convex, and can be solved efficiently using numerical methods, such as `[1]`. In fact, devising better methods to solve  $\ell_1$ -type problems like (7.9) is a prolific active field of research.