

## 1 Essentials

### 1.1 Matrix/Vector Derivatives

•  $\nabla_x(u \cdot v) = u \cdot \nabla_x(v) + v \cdot \nabla_x(u)$  •  $\nabla_x(g(u)) = \nabla_u(g(u)) \cdot \nabla_x(u)$  •  $\nabla_x(a^T x) = \nabla_x(x^T a) = a$  •  $\nabla_x(b^T A x) = A^T b$  •  $\nabla_x(x^T A x) = (A + A^T)x = \text{sym. } 2Ax$  •  $\nabla_x(a^T x x^T b) = (ab^T + ba^T)x$  •  $\nabla_x \|x\|^2 = 2x$  •  $\nabla_x \|x\| = \frac{x}{\|x\|}$

### 1.2 Norms

**$l_0$ :**  $\|x\|_0 := |\{i | x_i \neq 0\}|$  **Nuclear:**  $\|X\|_* = \sum_{i=1}^{\min(m,n)} \sigma_i$   
 **$p$ -norm:**  $\|x\|_p := \left(\sum_{i=1}^N |x_i|^p\right)^{\frac{1}{p}}$  **Frobenius:**  $\|A\|_F := \sqrt{\sum_{i,j} |A_{i,j}|^2} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}$  ( $\sigma_i$  is the  $i$ -th singularvalue)

### 1.3 Eigenvalue / -vectors

Eigenvalue Problem:  $Ax = \lambda x$

1. solve  $\det(A - \lambda I) \stackrel{!}{=} 0$  resulting in  $\{\lambda_i\}_i$
2.  $\forall \lambda_i$ : solve  $(A - \lambda_i I)x_i = 0$ ,  $x_i$  is the  $i$ -th eigenvector.
3. (opt.) normalize eigenvector  $q_i$ :  $q_i^{\text{norm}} = \frac{1}{\|q_i\|_2} q_i$ .

### 1.4 Eigendecomposition

$A \in \mathbb{R}^{N \times N}$ ,  $A = U \Lambda U^{-1}$ ,  $Q \in \mathbb{R}^{N \times N}$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ ,  $U^T = U^{-1}$ ,  $A$  symmetric then  $A^{-1} = U \Lambda^{-1} U^{-1}$

### 1.5 Probability / Statistics

•  $P(x|y) := \frac{P(x,y)}{P(y)}$ , if  $P(y) > 0$  •  $\sum_{x \in X} P(x|y) = 1$  •  $P(x,y) = P(x|y)P(y)$  •  $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$  (Bayes' rule) •  $P(x|y) = P(x) \Leftrightarrow P(y|x) = P(y)$  (iff  $X, Y$  independent) •  $P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i)$  (iff IID)

### 2 Dimensionality Reduction / PCA

$X \in \mathbb{R}^{D \times N}$ .  $N$  observations,  $K$  properties. Target:  $\tilde{X} \in \mathbb{R}^{K \times N}$ .

1. Empirical Mean:  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$
  2. Center Data:  $\bar{X} = X - [\bar{x}, \dots, \bar{x}] = X - M$
  3. Cov. Matrix:  $\Sigma = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T = \frac{1}{N} \bar{X} \bar{X}^T$
  4. Eig.Dec.:  $\Sigma = U \Lambda U^T$
  5. Select  $K < D$ , keep first  $K$  ew. and ev.  $\Rightarrow U_K, \lambda_K$
  6. Transform data onto new Basis:  $\bar{Z}_K = U_K^T \bar{X}$
  7. Reconstruct to original Basis:  $\tilde{\bar{X}} = U_K \bar{Z}_K$
  8. Reverse centering:  $\tilde{X} = \tilde{\bar{X}} + M$
- For compression save  $U_K, \bar{Z}_K, \bar{x}$ .
  - $U_k \in \mathbb{R}^{D \times K}$ ,  $\Sigma \in \mathbb{R}^{D \times D}$ ,  $\bar{Z}_K \in \mathbb{R}^{K \times N}$ ,  $\bar{X} \in \mathbb{R}^{D \times N}$

### 3 SVD

•  $A = U D V^T = \sum_{k=1}^{\text{rank}(A)} d_{k,k} u_k(v_k)^T$  •  $A \in \mathbb{R}^{N \times P}$ ,  $U \in \mathbb{R}^{N \times N}$ ,  $D \in \mathbb{R}^{N \times P}$ ,  $V \in \mathbb{R}^{P \times P}$  •  $U^T U = V^T V = I$  (cols. orthonormal) • cols. of  $U$  are ev. of  $AA^T$  (row sim. result of

PCA of  $A$ ),  $V$  of  $A^T A$  (col. sim.),  $D = \text{diag}(\sigma_i)$ ,  $\sigma_i^2 = \lambda_i$  for  $u_i, v_i$  • cols. of  $V$  where  $\sigma_i = 0$  span  $\text{null}(A)$ ,  $U$  where  $\sigma_i > 0$  span  $\text{range}(A)$  •  $A$  sym. then  $A = U D U^T$  and  $u_i$  ev. of  $A$  •  $A_k = \sum_{i=1}^k u_i \sigma_i v_i^T$  •  $\min_{\text{rank}(B)=k} \|A - B\|_2 = \|A - A_k\|_2 = \sigma_{k+1}$  •  $\min_{\text{rank}(B)=k} \|A - B\|_F^2 = \|A - A_k\|_F^2 = \sigma_{k+1}^2 + \dots + \sigma_r^2$   
1. calculate  $A^T A$ .  
2. calculate eigenvalues of  $A^T A$ , the square root of them, in descending order, are the diagonal elements of  $D$ .  
3. calculate eigenvectors of  $A^T A$  using the eigenvalues resulting in the columns of  $V$ .  
4. calculate the missing matrix:  $U = A V D^{-1}$ . Can be checked by calculating the eigenvectors of  $A A^T$ .  
5. normalize each column of  $U$  and  $V$ .

### 4 K-means Algorithm

**Target:**  $\min_{U,Z} J(U,Z) = \sum_{n=1}^N \sum_{k=1}^K z_{k,n} \|x_n - u_k\|_2^2 = \|X - UZ\|_F^2$  1.  $U = [u_1^{(0)}, \dots, u_k^{(0)}]$  2.  $k^*(x_n) = \arg \min_k \{\|x_n - u_k^{(t-1)}\|_2\}$  Set  $z_{j,n}^{(t)} = 1$  if  $j = k^*$  else 0. 3.  $u_k^{(t)} = \frac{\sum_{n=1}^N z_{k,n}^{(t)} x_n}{\sum_{n=1}^N z_{k,n}^{(t)}}$ .  
4. stops if  $\|u_k^{(t-1)} - u_k^{(t)}\| < \varepsilon \forall k$ .

### 4.1 Clustering Stability

multiple runs return similar clusters • dist. between clust. (same data):  $d(C, C') := \min_{\Pi} \frac{1}{2} \|Z - \Pi(Z')\|_F^2$ ,  $\Pi(Z') =$  row perm. of  $Z'$  • arbitrary sets  $X, X'$  of size  $N, N'$ :  $r := \frac{1}{N'} \min_{\Pi} \{\sum_{n=1}^{N'} \mathbb{I}_{\{\Pi(\phi(x'_n)) \neq z'_n\}}\}$  ( $\phi$ : multi-class classifier trained on  $(X, Z)$ ) • for  $K$  clusters: stability :=  $1 - \frac{r}{r_{\text{rand}}}$  (1 good, 0 bad), rand. clust. of equal size:  $r_{\text{rand}} = \frac{K-1}{K}$ .

### 5 Gaussian Mixture Models (GMM)

For GMM let  $\theta_k = (\mu_k, \Sigma_k)$ ;  $p_{\theta_k}(x) = \mathcal{N}(x | \mu_k, \Sigma_k) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_k|}} \exp(-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k))$

**Mixture Models:**  $p_{\theta}(x) = \sum_{k=1}^K \pi_k p_{\theta_k}(x)$

**Assignment variable (generative model):**

$z_k \in \{0, 1\}$ ,  $\sum_{k=1}^K z_k = 1$ ,  $\Pr(z_k = 1) = \pi_k \Leftrightarrow p(z) = \prod_{k=1}^K \pi_k^{z_k}$

**Complete data distribution:**  $p_{\theta}(x, z) = \prod_{k=1}^K (\pi_k p_{\theta_k}(x))^{z_k}$

**Posterior Probabilities:**

$\Pr(z_k = 1 | x) = \frac{\Pr(z_k=1)p(x|z_k=1)}{\sum_{i=1}^K \Pr(z_i=1)p(x|z_i=1)} = \frac{\pi_k p_{\theta_k}(x)}{\sum_{i=1}^K \pi_i p_{\theta_i}(x)}$

**Likelihood of observed data  $X$ :**  $p_{\theta}(X) = \prod_{n=1}^N p_{\theta}(x_n) = \prod_{n=1}^N (\sum_{k=1}^K \pi_k p_{\theta_k}(x_n))$

**MLE:**  $\arg \max_{\theta} \sum_{n=1}^N \log(\sum_{k=1}^K \pi_k p_{\theta_k}(x_n))$

$\log\left(\sum_{k=1}^K \frac{q_k \pi_k p_{\theta_k}(x_n)}{q_k}\right) \geq \sum_{k=1}^K q_k [\log p_{\theta_k}(x_n) + \log \pi_k - \log q_k]$

with  $\sum_{k=1}^K q_k = 1$  by Jensens inequality

### 5.1 Expectation-Maximization (EM) for GMM

1. Initialize  $\pi_k^{(0)}, \mu_k^{(0)}, \Sigma_k^{(0)}$  for  $k = 1, \dots, K$  and  $t = 1$ .
2. E-Step:  $q_{k,n}^* = \Pr[z_{k,n} = 1 | x_n]$
3. M-Step:  $\mu_k^* := \frac{\sum_{n=1}^N q_{k,n} x_n}{\sum_{n=1}^N q_{k,n}}$  &  $\pi_k^* := \frac{1}{N} \sum_{n=1}^N q_{k,n}$   
&  $\Sigma_k^* = \frac{\sum_{n=1}^N q_{k,n} (x_n - \mu_k^*)(x_n - \mu_k^*)^T}{\sum_{n=1}^N q_{k,n}}$
4. stop if  $\|\log p_{\theta_{(t-1)}} - \log p_{\theta_{(t)}}\| < \varepsilon$

### 5.2 Model Order Selection (AIC / BIC for GMM)

Trade-off between data fit (i.e. likelihood  $p(X|\theta)$ ) and complexity (i.e. # of free parameters  $\kappa(\cdot)$ ). For choosing  $K$ : • **Akaike Information Criterion:**  $\text{AIC}(\theta|X) = -\log p_{\theta}(X) + \kappa(\theta)$  • **Bayesian Information Criterion:**  $\text{BIC}(\theta|X) = -\log p_{\theta}(X) + \frac{1}{2} \kappa(\theta) \log N$  • # of free params: fixed covariance matrix:  $\kappa(\theta) = K \cdot D + (K - 1)$  ( $K$ : # clusters,  $D$ :  $\dim(\text{data}) = \dim(\mu_i)$ ,  $K - 1$ : # free clusters), full covariance matrix:  $\kappa(\theta) = K(D + \frac{D(D+1)}{2}) + (K - 1)$ . • Compare AIC/BIC for different  $K$  – the smaller the better. BIC penalizes complexity more.

### 6 Word Embeddings

**Distributional Model:**  $p_{\theta}(w|w') = \Pr[w \text{ occurs close to } w']$

**Log-likelihood:**  $L(\theta; w) = \sum_{t=1}^T \sum_{\Delta \in I} \log p_{\theta}(w^{(t+\Delta)} | w^{(t)})$

**Latent Vector Model:**  $w \mapsto (x_w, b_w) \in \mathbb{R}^{D+1}$

$p_{\theta}(w|w') = \frac{\exp[\langle x_w, x_{w'} \rangle + b_w]}{\sum_{v \in V} \exp[\langle x_v, x_{w'} \rangle + b_v]}$ . Modifications: • split vocab in main vocab  $V$ , context vocab  $C$ :  $p_{\theta}(w|w') = \langle x_{w'}, y_w \rangle + b_w$ , word embed.  $x_w$ , context embed.  $y_w$  • use GloVe objective

### 6.1 GloVe (Weighted Square Loss)

**Co-occurrence Matrix:**  $N = (n_{ij}) \in \mathbb{R}^{|V| \times |C|} \leftrightarrow \#w_i \text{ in c'txt } w_j$   
**Objective:**  $H(\theta; N) = \sum_{n_{ij} > 0} f(n_{ij}) (\log n_{ij} - \log \exp[\langle x_i, y_j \rangle + b_i + d_j])^2$  with  $f(n) = \min\{1, (\frac{n}{n_{\max}})^{\alpha}\}$ ,  $\alpha \in (0; 1]$ .

unnormalized distribution  $\rightarrow$  two-sided loss function

**SGD:** 1.  $x_i^{\text{new}} \leftarrow x_i + 2\eta f(n_{ij})(\log n_{ij} - \langle x_i, y_j \rangle) y_j$   
2.  $y_j^{\text{new}} \leftarrow y_j + 2\eta f(n_{ij})(\log n_{ij} - \langle x_i, y_j \rangle) x_i$

### 7 Non-Negative Matrix Factorization (NMF) / pLSA

**Context Model:**  $p(w|d) = \sum_{z=1}^K p(w|z)p(z|d)$

**Conditional independence assumption (\*):**  $p(w|d) = \sum_z p(w, z|d) = \sum_z p(w|z)p(z|d) \stackrel{*}{=} \sum_z p(w|z)p(z|d)$

**Symmetric parameterization:**  $p(w, d) = \sum_z p(z)p(w|z)p(d|z)$

### 7.1 EM for pLSA:

1.  $X = x_{i,j} = \# \text{occ. of } w_j \text{ in doc. } d_i$
2. Log-Likelihood:  $L(U, V) = \sum_{i,j} x_{i,j} \log p(w_j | d_i) = \sum_{(i,j) \in X} \log \sum_{z=1}^K p(w_j | z) p(z | d_i) = \sum_{(i,j) \in X} \log \sum_{z=1}^K v_{zj} u_{zi}$

3. E-Step (optimal q):  $q_{zij} = \frac{v_{zj}u_{zi}}{\sum_{k=1}^K v_{kj}u_{ki}}$

4. M-Steps:  $p(z|d_i) = \frac{\sum_j x_{ij}q_{zij}}{\sum_j x_{ij}} \quad \& \quad p(w_j|z) = \frac{\sum_i x_{ij}q_{zij}}{\sum_{i,l} x_{il}q_{zlj}}$

7.2 NMF Algorithm for quadratic cost function

•  $\mathbf{X} \in \mathbb{Z}_{\geq 0}^{N \times M}$  • NMF:  $\mathbf{X} \approx \mathbf{U}^\top \mathbf{V}$ ,  $x_{ij} = \sum_z u_{zi}v_{zj} = \langle \mathbf{u}_i, \mathbf{v}_j \rangle$

$\min_{\mathbf{U}, \mathbf{V}} J(\mathbf{U}, \mathbf{V}) = \frac{1}{2} \|\mathbf{X} - \mathbf{U}^\top \mathbf{V}\|_F^2$  s.t.  $\forall i, j, z \ u_{zi}, v_{zj} \geq 0$

1. init:  $\mathbf{U}, \mathbf{V} = \text{rand}()$  2. repeat for  $\text{maxIters}$ : 3. update  $\mathbf{U}$ :  $(\mathbf{V}\mathbf{V}^\top)\mathbf{U} = \mathbf{V}\mathbf{X}^\top$  4. project  $u_{zi} = \max\{0, u_{zi}\}$  5. update  $\mathbf{V}$ :  $(\mathbf{U}\mathbf{U}^\top)\mathbf{V} = \mathbf{U}\mathbf{X}$  6. project  $v_{zj} = \max\{0, v_{zj}\}$

8 Convolutional Neural Networks

**Neurons:**  $F_\sigma(\mathbf{x}; \mathbf{w}) = \sigma(w_0 + \sum_{i=1}^M x_i w_i)$ . **Output:** linear regression;  $\mathbf{y} = \mathbf{W}^L \mathbf{x}^{L-1}$ , binary classification;  $y_1 = \text{P}[Y = 1 | \mathbf{x}] = \frac{1}{1 + \exp[-\langle \mathbf{w}_1^L, \mathbf{x}^{L-1} \rangle]}$ , multiclass;  $y_k = \text{P}[Y = k | \mathbf{x}] = \frac{\exp[\langle \mathbf{w}_k^L, \mathbf{x}^{L-1} \rangle]}{\sum_{m=1}^K \exp[\langle \mathbf{w}_m^L, \mathbf{x}^{L-1} \rangle]}$ .

**Loss function**  $l(y, \hat{y})$ : squared loss;  $\frac{1}{2}(y - \hat{y})^2$ , cross-entropy loss;  $-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$ .

8.1 Neural Networks for Images

Translation invariance of images  $\rightarrow$  neurons compute same fct, shift invariant filters; weights defined as filter masks, e.g. convolution:  $F_{n,m}(\mathbf{x}; \mathbf{w}) = \sigma(b + \sum_{k=-2}^2 \sum_{l=-2}^2 w_{k,l} x_{n+k, m+l})$ . To reduce dimension of convolution, use  $\{\text{max}, \text{avg}\}$ -pooling

9 Optimization

9.1 Coordinate Descent (update the  $d$ -th coord. per step)

1. init:  $\mathbf{x}^{(0)} \in \mathbb{R}^D$  2. for  $t = 0$  to  $\text{maxIter}$ : 3. sample u.a.r.  $d \sim \{1, \dots, D\}$  4.  $\mathbf{u}^* = \arg \min_{\mathbf{u} \in \mathbb{R}} f(x_1^{(t)}, \dots, x_{d-1}^{(t)}, u, x_{d+1}^{(t)}, \dots, x_D^{(t)})$  5.  $\mathbf{x}_d^{(t+1)} = \mathbf{u}^*$  and  $\mathbf{x}_i^{(t+1)} = \mathbf{x}_i^{(t)}$  for  $i \neq d$

9.2 Gradient Descent (or Deepest Descent)

**Gradient:**  $\nabla f(\mathbf{x}) := \left( \frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_D} \right)^\top$  1. init:  $\mathbf{x}^{(0)} \in \mathbb{R}^D$  2. for  $t = 0$  to  $\text{maxIter}$ :  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \gamma \nabla f(\mathbf{x}^{(t)})$ , usually  $\gamma \approx \frac{1}{t}$

9.3 Stochastic Gradient Descent (SGD)

Assume **Additive Objective**;  $f(x) = \frac{1}{N} \sum_{n=1}^N f_n(x)$  1. init:  $\mathbf{x}^{(0)} \in \mathbb{R}^D$  2. for  $t = 0$  to  $\text{maxIter}$ : 3. sample u.a.r.  $n \sim \{1, \dots, N\}$  4.  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \gamma \nabla f_n(\mathbf{x}^{(t)})$ , usually stepsize  $\gamma \approx \frac{1}{t}$ .

9.4 Projected Gradient Descent (Constrained Opt.)

minimize  $f(x)$ ,  $x \in Q$  (constraint). **Project**  $x$  onto  $Q$ :  $P_Q(\mathbf{x}) = \arg \min_{\mathbf{y} \in Q} \|\mathbf{y} - \mathbf{x}\|$ , **Projected Gradient Update**:  $\mathbf{x}^{(t+1)} = P_Q[\mathbf{x}^{(t)} - \gamma \nabla f(\mathbf{x}^{(t)})]$ ,  $\mathbf{x}^{(t+1)}$  is unique if  $Q$  convex.

9.5 Lagrangian Multipliers

Minimize  $f(\mathbf{x})$  s.t.  $g_i(\mathbf{x}) \leq 0$ ,  $i = 1, \dots, m$  (**inequality constr.**) and  $h_i(\mathbf{x}) = \mathbf{a}_i^\top \mathbf{x} - b_i = 0$ ,  $i = 1, \dots, p$  (**equality constraint**)

**Lagrangian:**  $L(\mathbf{x}, \lambda, \nu) := f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{i=1}^p \nu_i h_i(\mathbf{x})$

**Dual function:**  $D(\lambda, \nu) := \inf_{\mathbf{x}} L(\mathbf{x}, \lambda, \nu) \in \mathbb{R}$

**Dual Problem:**  $\max_{\lambda, \nu} D(\lambda, \nu)$  s.t.  $\lambda \geq \mathbf{0}$ . Note:  $\max_{\lambda, \nu} D(\lambda, \nu) \leq \min_{\mathbf{x}} f(\mathbf{x})$ , equality if  $\text{dom } f$  and  $f$  convex

9.6 Convex Optimization

$Q$  convex:  $\forall \mathbf{x}, \mathbf{y} \in Q : \forall 0 \leq \alpha \leq 1 : \alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in Q$   
 $f : \mathbb{R}^D \rightarrow \mathbb{R}$  is convex, if  $\text{dom } f$  is a convex set, and if  $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f : \forall 0 \leq \alpha \leq 1 : f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y})$ . local=global min, **Convergence:**  $f(\mathbf{x}^{(t)}) - f(\mathbf{x}^*) \leq \frac{c}{t}$ .  
**Subgradient**  $g \in \mathbb{R}^D$  of  $f$  at  $\mathbf{x}$ :  $f(\mathbf{y}) \geq f(\mathbf{x}) + g^\top (\mathbf{y} - \mathbf{x}) \ \forall \mathbf{y}$   
Convergence:  $f(\mathbf{x}^{(t)}) - f(\mathbf{x}^*) \leq \frac{c}{\sqrt{t}}$ .

10 Sparse Coding

10.1 Orthogonal Basis

For  $\mathbf{x}$  and orthogonal  $\mathbf{U}$  compute  $\mathbf{z} = \mathbf{U}^\top \mathbf{x}$ . Approx  $\hat{\mathbf{x}} = \mathbf{U} \hat{\mathbf{z}}$ ,  $\hat{z}_i = z_i$  if  $|z_i| > \epsilon$  else 0. Energy preserving  $\|\mathbf{U} \mathbf{z}\| = \|\mathbf{z}\|$ . Reconstruction Error  $\|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \sum_{d \notin \sigma} \langle \mathbf{x}, \mathbf{u}_d \rangle^2$ .

10.2 Overcomplete Basis

$\mathbf{U} \in \mathbb{R}^{D \times L}$  for  $\text{dim}(\text{data}) = L > D = \# \text{ atoms}$ . Decoding involved  $\rightarrow$  add constraint  $\mathbf{z}^* \in \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0$  s.t.  $\mathbf{x} = \mathbf{U} \mathbf{z}$ . NP-hard  $\rightarrow$  approximate with 1-norm (convex) or with MP.

**Coherence** •  $m(\mathbf{U}) = \max_{i,j: i \neq j} |\mathbf{u}_i^\top \mathbf{u}_j|$  •  $m(\mathbf{B}) = 0$  if  $\mathbf{B}$  orthogonal matrix •  $m([\mathbf{B}, \mathbf{u}]) \geq \frac{1}{\sqrt{D}}$  if atom  $\mathbf{u}$  is added to orthogonal basis  $\mathbf{B}$  (o.n.b. = orthonormal base)

10.3 Dictionary Learning

Adapt the dictionary to signal characteristics. Objective:  $(\mathbf{U}^*, \mathbf{Z}^*) \in \arg \min_{\mathbf{U}, \mathbf{Z}} \|\mathbf{X} - \mathbf{U} \mathbf{Z}\|_F^2$  not jointly convex but convex in 1 argument.

**Matrix Factorization by Iter Greedy Minimization 1.** Coding step:  $\mathbf{Z}^{t+1} \in \arg \min_{\mathbf{Z}} \|\mathbf{X} - \mathbf{U}^t \mathbf{Z}\|_F^2$  subject to  $\mathbf{Z}$  being sparse 2. Dictionary update step:  $\mathbf{U}^{t+1} \in \arg \min_{\mathbf{U}} \|\mathbf{X} - \mathbf{U} \mathbf{Z}^{t+1}\|_F^2$ , subject to  $\forall l \in [L] : \|\mathbf{u}_l\|_2 = 1$

11 Robust PCA

• Idea: Approximate  $\mathbf{X}$  with  $\mathbf{L} + \mathbf{S}$ ,  $\mathbf{L}$  is low-rank,  $\mathbf{S}$  is sparse.  
•  $\min_{\mathbf{L}, \mathbf{S}} \text{rank}(\mathbf{L}) + \mu \|\mathbf{S}\|_0$ , s. t.  $\mathbf{L} + \mathbf{S} = \mathbf{X}$ . As non-convex, change to  $\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \mu \|\mathbf{S}\|_1$  (not the same in general)  
• Perfect recovery ( $\mathbf{X} = \mathbf{S}_0 + \mathbf{L}_0 = \mathbf{S}^* + \mathbf{L}^*$ ) is not possible if  $\mathbf{S}$  is low-rank,  $\mathbf{L}$  is sparse, or  $\mathbf{X}$  is low-rank and sparse at the same time. Principal components of  $\mathbf{L}$  must not be sparse.

11.1 Alternating Direction Method of Multipliers (ADMM)

$\min_{\mathbf{x}_1, \mathbf{x}_2} f_1(\mathbf{x}_1) + f_2(\mathbf{x}_2)$  s. t.  $\mathbf{A}_1 \mathbf{x}_1 + \mathbf{A}_2 \mathbf{x}_2 = \mathbf{b}$ ,  $f_1, f_2$  convex • Augmented Lagrangian:  $L_\rho(\mathbf{x}_1, \mathbf{x}_2, \lambda) = f_1(\mathbf{x}_1) + f_2(\mathbf{x}_2) + \lambda^\top (\mathbf{A}_1 \mathbf{x}_1 + \mathbf{A}_2 \mathbf{x}_2 - \mathbf{b}) + \frac{\rho}{2} \|\mathbf{A}_1 \mathbf{x}_1 + \mathbf{A}_2 \mathbf{x}_2 - \mathbf{b}\|_2^2$

• ADMM:  $\mathbf{x}_1^{(t+1)} := \arg \min_{\mathbf{x}_1} L_\rho(\mathbf{x}_1, \mathbf{x}_2^{(t)}, \lambda^{(t)})$ ,  $\mathbf{x}_2^{(t+1)} := \arg \min_{\mathbf{x}_2} L_\rho(\mathbf{x}_1^{(t+1)}, \mathbf{x}_2, \lambda^{(t)})$ ,  $\lambda^{(t+1)} := \lambda^{(t)} + \rho(\mathbf{A}_1 \mathbf{x}_1^{(t+1)} + \mathbf{A}_2 \mathbf{x}_2^{(t+1)} - \mathbf{b})$  • ADDM for RPCA:  $L_\rho(\mathbf{L}, \mathbf{S}, \lambda) = \|\mathbf{L}\|_* + \mu \|\mathbf{S}\|_1 + \langle \lambda, \text{vec}(\mathbf{L} + \mathbf{S} - \mathbf{X}) \rangle + \frac{\rho}{2} \|\mathbf{L} + \mathbf{S} - \mathbf{X}\|_F^2$   
•  $\mathbf{L}^{t+1} = \arg \min_{\mathbf{L}} (L_\rho(\mathbf{L}, \mathbf{S}^t, \lambda^t))$  •  $\mathbf{S}^{t+1} = \arg \min_{\mathbf{S}} (L_\rho(\mathbf{L}^{t+1}, \mathbf{S}, \lambda^t))$  •  $\lambda^{t+1} = \lambda^t + \rho \text{vec}(\mathbf{L}^{t+1} + \mathbf{S}^{t+1} - \mathbf{X})$