

## Linear Algebra

$\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top \mathbf{v} = \sum_{i=1}^d u_i v_i$ ,  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$   
 $\langle \mathbf{ax}, \mathbf{y} \rangle = \mathbf{a} \langle \mathbf{x}, \mathbf{y} \rangle$ ,  $\langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle$   
 $\langle \mathbf{x}, \mathbf{x} \rangle = \|\mathbf{x}\|_2^2$   
 $\mathbf{AB} = \mathbf{C} \Rightarrow \sum_{k=1}^n a_{ik} b_{kj}$

## Orthogonality

$\mathbf{u}, \mathbf{v}$  are orthogonal iff  $\langle \mathbf{u}, \mathbf{v} \rangle = 0$   
 $\mathbf{A}$  orthog.  $\Rightarrow \mathbf{AA}^\top = \mathbf{A}^\top \mathbf{A} = \mathbf{I}$ ,  $\mathbf{A}^\top = \mathbf{A}^{-1}$   
 $\mathbf{A}$  has orthonormal rows, i.e.  $\|\mathbf{a}_i\|_2 = 1$   
 $\|\mathbf{Ax}\|_2 = \|\mathbf{x}\|_2$ , dist. & energy preserved  
**basis change**  $\mathbf{x} = \sum_{i=1}^n \langle \mathbf{x}, \mathbf{u}_i \rangle \mathbf{u}_i = \mathbf{U}^\top \mathbf{x}$   
and  $\|\langle \mathbf{x}, \mathbf{u}_i \rangle \mathbf{u}_i\|_2^2 = |\langle \mathbf{x}, \mathbf{u}_i \rangle|^2$   
 $\det(\mathbf{A}) \in \{-1, 1\}$ ,  $\det(\mathbf{A}^\top \mathbf{A}) = 1$

## Symmetry

$\mathbf{A}^\top = \mathbf{A}$ . All eigenvalues of  $\mathbf{A} > 0$ .  
 $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$  iff  $\mathbf{A}$  symmetric where  $\mathbf{U}$  is orthogonal.

## Positive Semi-Definit

$\mathbf{v}^\top \mathbf{A} \mathbf{v} \geq 0 \forall \mathbf{v}$ ,  $\mathbf{A} = \mathbf{B}^\top \mathbf{B} \Rightarrow \mathbf{A}$  is p.s.t.  
 $\mathbf{A}$  pst iff all its eigenvalues  $\geq 0$

## Rank

dimension of vector space generated by  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .  
 $\text{rank}(\mathbf{A}) = \# \sigma > 0$ ,  $\text{rank}(\mathbf{A}) \leq \min(m, n)$   
 $\text{rank}(\mathbf{AB}) = \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\}$

## Norms

$\|\mathbf{x}\|_0 = |\{i | x_i \neq 0\}|$ ,  $\|\mathbf{x}\|_p = (\sum_{i=1}^N |x_i|^p)^{\frac{1}{p}}$

nuclear  $\|\mathbf{X}\|_* = \sum_{i=1}^{\min(m, n)} \sigma_i$

eucl.  $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^N \mathbf{x}_i^2} = \sqrt{\mathbf{x}^\top \mathbf{x}} = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sigma_1$

frob.  $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^M \sum_{j=1}^N |\mathbf{A}_{i,j}|^2} =$

$\sqrt{\text{trace}(\mathbf{A}^\top \mathbf{A})} = \sum_{i=1}^{\min(m, n)} \sigma_i^2$ ,  $\mathbf{A} \in \mathbb{R}^{M \times N}$   
w.  $\text{trace}(\mathbf{ABC}) = \text{trace}(\mathbf{CAB})$  (cyclic)

## Determinant

$\det(\mathbf{A}) = a_{11} a_{22} - a_{12} a_{21} = \prod_i \sigma_i$   
 $\det(\mathbf{A}^\top) = \det(\mathbf{A})$ ,  $\det(\mathbf{A}^{-1}) = (\det(\mathbf{A}))^{-1}$

## Eigenvectors & Eigenvalues

$\sigma \in \mathbb{R}$  s.t.  $\mathbf{A} \mathbf{u} = \sigma \mathbf{u}$  Find eigenvals: solve  $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$

Find eigenvcs  $\mathbf{u}_i$ : solve  $(\mathbf{A} - \sigma_i \mathbf{I}) \mathbf{u}_i = \mathbf{0}$

## SVD - Singular Value Decomposition

$\mathbf{A} = \mathbf{U} \cdot \mathbf{D} \cdot \mathbf{V}^\top \forall \mathbf{A} \in \mathbb{R}^{m \times n}$   
 $\mathbf{U} \in \mathbb{R}^{m \times m}$  orth. matrix ( $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_m$ )  
 $\mathbf{D} \in \mathbb{R}^{m \times n}$  diagonal matrix  
 $\mathbf{V} \in \mathbb{R}^{n \times n}$  orth. matrix ( $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_n$ )  
 $\mathbf{U}$  has eigenvcs for  $\mathbf{AA}^\top$ ,  $\mathbf{V}$  for  $\mathbf{A}^\top \mathbf{A}$ .  
 $\mathbf{D}$  has eigenvalues for  $\mathbf{AA}^\top$  and  $\mathbf{A}^\top \mathbf{A}$ ,  
are sqrt of entries of  $\mathbf{\Lambda}$ :  $\mathbf{AA}^\top = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$

## Eckart-Young Theorem

$\arg \min_{\hat{\mathbf{X}}: \text{rank}(\hat{\mathbf{X}})=k} \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \mathbf{U} \Sigma_k \mathbf{V}^\top$   
 $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_s)$ ,  $\sigma_1 \geq \dots \geq \sigma_s \geq 0$   
 $\min_{\mathbf{r}(\mathbf{B})} \|\mathbf{A} - \mathbf{B}\|_F^2 = \|\mathbf{A} - \mathbf{A}_k\|_F^2 = \sum_{r=k+1}^{\mathbf{r}(\mathbf{A})} \sigma_r^2$   
 $\min_{\mathbf{r}(\mathbf{B})} \|\mathbf{A} - \mathbf{B}\|_2 = \|\mathbf{A} - \mathbf{A}_k\|_2 = \sigma_{k+1}$

Those are not convex optimisation problems (combination of matrices of rank  $k$  is usually not rank  $k$ )

## Probabilities

### Expectation

$\mathbb{E}[X] = \int_{\Omega} x f(x) dx = \int_{\omega} x P[X=x] dx$   
 $\mathbb{E}_{Y|X}[Y] = \mathbb{E}_Y[Y|X]$   
 $\mathbb{E}_X[f(Y)] = f(Y)$  if  $P(Y|X) = P(Y)$

### Variance & Covariance

$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$   
 $\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] \quad XY \text{ iid}$   
 $\text{Var}[\alpha X] = \alpha^2 \text{Var}[X]$   
 $\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$

### Conditional Probabilities

$P[X|Y] = \frac{P[X, Y]}{P[Y]}$ ,  $P[\bar{X}|Y] = 1 - P[X|Y]$

**Bayes:**  $P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$

### Distributions

$\mathcal{N}(x|\mu, \sigma^2) = (\sqrt{2\pi\sigma^2})^{-1} \exp^{-(x-\mu)^2/(2\sigma^2)}$   
 $\mathcal{N}(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^\top \Sigma^{-1}(\mathbf{x}-\mu)}$

$\text{Exp}(x|\lambda) = \lambda e^{-\lambda x}$ ,  $\text{Ber}(x|\theta) = \theta^x (1-\theta)^{(1-x)}$

### Kullback-Leibler Divergence

divergence of distr  $P$  &  $Q$   $D_{KL}(P||Q) = -\sum_{x \in \mathcal{X}} P(x) \log \frac{Q(x)}{P(x)} = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}$

### Matrix Derivatives

$\frac{\partial(\mathbf{b}^\top \mathbf{x})}{\partial \mathbf{x}} = \frac{\partial(\mathbf{x}^\top \mathbf{b})}{\partial \mathbf{x}} = \mathbf{b}$ ,  $\frac{\partial(\mathbf{x}^\top \mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{x}$ ,  
 $\frac{\partial(\mathbf{b}^\top \mathbf{Ax})}{\partial \mathbf{x}} = \mathbf{A}^\top \mathbf{b}$ ,  $\frac{\partial(\mathbf{x}^\top \mathbf{Ax})}{\partial \mathbf{x}} = (\mathbf{A}^\top + \mathbf{A})\mathbf{x}$ ,  
 $\frac{\partial(\mathbf{c}^\top \mathbf{Xb})}{\partial \mathbf{X}} = \mathbf{cb}^\top$ ,  $\frac{\partial(\mathbf{c}^\top \mathbf{X}^\top \mathbf{b})}{\partial \mathbf{X}} = \mathbf{bc}^\top$ ,  $\frac{\partial \log |\mathbf{A}|}{\partial \mathbf{A}} = \mathbf{A}^{-\top}$ ,  
 $\frac{\partial(\|\mathbf{x} - \mathbf{b}\|_2)}{\partial \mathbf{x}} = \frac{\mathbf{x} - \mathbf{b}}{\|\mathbf{x} - \mathbf{b}\|_2}$ ,  $\frac{\partial(\|\mathbf{x}\|_2^2)}{\partial \mathbf{x}} = \frac{\partial(\mathbf{x}^\top \mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{x}$ ,  
 $\frac{\partial(\|\mathbf{X}\|_F^2)}{\partial \mathbf{X}} = 2\mathbf{X}$ ,  $\frac{\partial(\|\mathbf{Ax} - \mathbf{b}\|_2^2)}{\partial \mathbf{x}} = 2(\mathbf{A}^\top \mathbf{Ax} - \mathbf{A}^\top \mathbf{b})$

$\mathbf{X}^\top \mathbf{X}$ : invertible if eigvals  $\neq 0$  but instable if big ratio last/first eigval.

### Derivatives

**chain rule**  $\frac{d}{dx} f(g(x)) = f'(g(x))g'(x)$   
 $\frac{d}{dx} \cos(x) = -\sin(x)$ ,  $\frac{d}{dx} \sin(x) = \cos(x)$   
 $\frac{d}{dx} \exp(x) = \exp$ ,  $\frac{d}{dx} \log(x) = \frac{1}{x}$

### Optimization

#### Gradient Descent

$\theta^{\text{new}} \leftarrow \theta^{\text{old}} - \eta \nabla_{\theta} \mathcal{L}(\theta)$   
Convergence isn't guaranteed.  
Less zigzag by adding momentum:  
 $\theta^{(l+1)} \leftarrow \theta^{(l)} - \eta \nabla_{\theta} \mathcal{L} + \mu(\theta^{(l)} - \theta^{(l-1)})$

## Stochastic Gradient Descent

Assume  $\mathcal{L}(\theta) = \sum_{n=1}^N \mathcal{L}_n(\theta)$   
 $\theta^{\text{new}} \leftarrow \theta^{\text{old}} - \eta \nabla_{\theta} \mathcal{L}_n(\theta)$   $n \sim \text{uniform}$   
this requires the function to be  $L$ -smooth to avoid oscillations

## Jensen's Inequality

$\sigma(\frac{\sum a_i x_i}{\sum a_i}) \geq \frac{\sum a_i \sigma(x_i)}{\sum a_i}$  if  $\sigma$  concave (log)  
for distribution  $a$ :  $\sum_i a_i = 1$

## Convexity

$f$  convex if  $\forall x_1, x_2 \in \mathcal{X}, \forall t \in [0, 1]$   
 $f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$   
if  $f$  is twice diff., it is convex iff it's double Hessian is positive semi-definite  
 $\mathbf{A}$  convex function has a unique minimum. sum of conv. is conv.

## Linear Autoencoder

Encoder  $\mathbf{C} \in \mathbb{R}^{k \times m}$ , Decoder  $\mathbf{D} \in \mathbb{R}^{m \times k}$   
lin. map:  $\mathbf{F}: \mathbb{R}^m \rightarrow \mathbb{R}^m$  w. limited rank.  
performs low rank approximation  
 $\ell(\mathbf{x}; \theta) = \frac{1}{2} \|\mathbf{x} - \hat{\mathbf{x}}(\theta)\|^2$   
 $\hat{\mathbf{x}}(\theta) = \mathbf{DCx}$ ,  $\theta = (\mathbf{C}, \mathbf{D})$

### Reconstruction error

$J(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) = \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - \hat{\mathbf{x}}_i(\theta)\|^2 = \frac{1}{2n} \|\mathbf{X} - \hat{\mathbf{X}}(\theta)\|_F^2$

### Optimal Solution

$\mathbf{C}^* = \mathbf{U}_k^\top$ ,  $\mathbf{D}^* = \mathbf{U}_k$ , s.t.  $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^\top$   
 $\hat{\mathbf{X}} = \mathbf{D}^* \mathbf{C}^* \mathbf{X} = \mathbf{U}_k \mathbf{U}_k^\top (\mathbf{U} \Sigma \mathbf{V}^\top) = \dots = \mathbf{U} \Sigma_k \mathbf{V}^\top$

opt. by EY. not the only optimal solution, i.e. limited interpretability

## Weight Sharing

$\mathbf{D} = \mathbf{C}^\top$  reduces ambiguity but not modeling power. Mapping unique.

## PCA - Principle Component Analysis

project data  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  to basis of orthogonal components.  
Centralise data by subtracting the mean  $\bar{\mathbf{X}} = \mathbf{X} - [\bar{x} \dots \bar{x}]$ ,  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$ .

### Variance-Covariance Matrix

$\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^\top = \frac{1}{N} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top$ .  
symmetric  $\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ ,  $\mathbf{U}$  orth.  
New orthogonal basis  $\mathbf{U}_K$ ,  $K \ll D$

$\bar{\mathbf{Z}}_K = \mathbf{U}_K^\top \bar{\mathbf{X}}$  and  $\tilde{\tilde{\mathbf{X}}} = \mathbf{U}_K \bar{\mathbf{Z}}_K$  opt. rec.

### Iterative View

Residual  $\mathbf{r}_i$ :  $\mathbf{x}_i - \tilde{\mathbf{x}}_i = \mathbf{I} - \mathbf{u} \mathbf{u}^\top \mathbf{x}_i$   
Covariance  $\frac{1}{N} \sum_{i=1}^N \mathbf{r}_i \mathbf{r}_i^\top = \Sigma - \lambda \mathbf{u} \mathbf{u}^\top$   
1<sup>st</sup> eigvec of  $\Sigma - \lambda \mathbf{u} \mathbf{u}^\top = 2^{\text{nd}} eigvec of  $\Sigma$   
get  $d$  princ. eigvecs of  $\Sigma$  by iteration$

### Power Method

$\mathbf{v}_{t+1} = \frac{\mathbf{A} \mathbf{v}_t}{\|\mathbf{A} \mathbf{v}_t\|}$ ,  $\lim_{t \rightarrow \infty} \mathbf{v}_t = \mathbf{u}_1$   
assume  $\langle \mathbf{u}_1, \mathbf{v}_0 \rangle \neq 0$  and  $|\lambda_1| > |\lambda_j|$

## 1-D PCA to a line

**Line**  $\mu + \mathbb{R} \mathbf{u} \equiv \{\mathbf{v} \in \mathbb{R}^m : \exists z \text{ s.t. } \mathbf{v} = \mu + z \mathbf{u}\}$   
 $\hat{\mathbf{x}} = \mu + \langle \mathbf{x} - \mu, \mathbf{u} \rangle \mathbf{u} = \arg \min_{\hat{\mathbf{x}} \in \mu + \mathbb{R} \mathbf{u}} \|\mathbf{x} - \hat{\mathbf{x}}\|^2$   
**Center the data** to find unique  $\mathbf{u}$   
 $\mathbf{u} \leftarrow \arg \min [\frac{1}{n} \sum_{i=1}^n \|\langle \mathbf{u}, \mathbf{x}_i \rangle \mathbf{u} - \mathbf{x}_i\|^2]$   
yields  $\mathbf{u} \leftarrow \arg \max [\frac{1}{n} \sum_{i=1}^n \langle \mathbf{u}, \mathbf{x}_i \rangle^2] = \arg \max [\mathbf{u}^\top (\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top) \mathbf{u}] = \arg \max [\mathbf{u}^\top \frac{1}{n} \mathbf{X} \mathbf{X}^\top \mathbf{u}] = \arg \max [\mathbf{u} \Sigma \mathbf{u}^\top]$   
**Lagrangian Optimization**  
 $\mathcal{L}(\mathbf{u}, \lambda) = \mathbf{u}^\top \Sigma \mathbf{u} + \lambda (\mathbf{u}^\top \mathbf{u} - 1)$

$\nabla_{\mathbf{u}} \mathcal{L}(\mathbf{u}, \lambda) \stackrel{!}{=} 0 \Leftrightarrow \Sigma \mathbf{u} = \lambda \mathbf{u}$ , ie  $\mathbf{u}$  princ. vec

## Matrix Approximation & Reconst

**Exact Rec.**  $\min_{\mathbf{B}} \|\mathbf{B}\|_*$  s.t.  $\|\mathbf{A} - \mathbf{B}\|_G = 0$

**Approx. Rec.**  $\min_{\mathbf{B}} \|\mathbf{A} - \mathbf{B}\|_G^2$  s.t.  $\|\mathbf{B}\|_* \leq k$

## SVD (e.g recommendation system)

$\mathbf{U}$  users to factors associations  $\mathbf{V}$  items to factors associations  $\Sigma$  level of strength of each factor. **Limitations** matrix incomplete. can't run SVD. Lower rank might not be the best solution. Ratings should not be outside the expected range

## Beyond SVD

optimize  $\min_{\text{rk}(\mathbf{B})=k} [\sum_{(i,j) \in \mathcal{I}} (a_{ij} - b_{ij})^2]$   
 $\mathcal{I} = \{(i, j)\}$  observed values

$\|\mathbf{X}\|_G = \sqrt{\sum_{i,j} g_{ij} x_{ij}^2}$ ,  $g_{ij} \in \{0, 1\}$

$\min_{\text{rk}(\mathbf{B})} \|\mathbf{A} - \mathbf{B}\|_G^2$  NP hard

## Alternating Least Square

parametrize  $\mathbf{B}_{rk(k)} = \mathbf{U} \mathbf{V}$ ,  $\mathbf{U} \in \mathbb{R}^{m \times k}$ ,  $\mathbf{V} \in \mathbb{R}^{k \times n}$

$f(\mathbf{U}, \mathbf{V}) = \frac{1}{|\mathcal{I}|} \sum_{(i,j) \in \mathcal{I}} (a_{ij} - \langle \mathbf{u}_i, \mathbf{v}_j \rangle)^2 = \sum_i [\sum_{j: (i,j) \in \mathcal{I}} (a_{ij} - \langle \mathbf{u}_i, \mathbf{v}_j \rangle)^2] = \sum_i f(\mathbf{U}, \mathbf{v}_i)$ .  
 $f$  is convex for fixed  $\mathbf{U}$  in  $\mathbf{V}$  and vice-versa but not jointly

**Regularise**  $\Omega(\mathbf{U}, \mathbf{V}) = \|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2$

**Minimise**  $f(\mathbf{U}, \mathbf{V}) + \mu \Omega(\mathbf{U}, \mathbf{V})$  alternate between  $\mathbf{U}$  and  $\mathbf{V}$ , ie:

$f(\mathbf{U}, \mathbf{v}_i) = \sum_{(i,j) \in \mathcal{I}} (a_{ij} - \langle \mathbf{u}_i, \mathbf{v}_i \rangle)^2$

## Convex Relaxation

$\text{rank}(\mathbf{B}) \geq \|\mathbf{B}\|_*$  for  $\|\mathbf{B}\|_2 \leq 1$   
 $\min_{\mathbf{B} \in \mathcal{P}_k} \|\mathbf{A} - \mathbf{B}\|_G^2$ ,  $\mathcal{P}_k = \{\mathbf{B} : \|\mathbf{B}\|_* \leq k\}$   
 $\mathcal{P}_k \cap \mathcal{Q}_k = \{\mathbf{B} : \text{rank} \mathbf{B} = k\}$   
 $\mathbf{B}^* = \text{shrink}_{\tau}(\mathbf{A}) = \arg \min_{\mathbf{B}} [\frac{1}{2} \|\mathbf{A} - \mathbf{B}\|_F^2 + \tau \|\mathbf{B}\|_*]$  thus  
 $\mathbf{B}^* = \mathbf{U} \mathbf{D}_{\tau} \mathbf{V}^\top$ ,  $\mathbf{D}_{\tau} = \text{diag}(\max(0, \sigma_i - \tau))$   
**Shrinkage Iterations** for  $\eta \geq 0$   
 $\mathbf{B}_{t+1} = \mathbf{B}_t + \eta \mathbf{I} (\mathbf{A} - \text{shrink}_{\tau}(\mathbf{B}_t))$   
 $\Pi(\mathbf{X}) = x_{ij}$  if  $(i, j) \in \mathcal{I}$ , 0 otw.  
generalisation guarantees  $\Pi(\mathbf{A}) = \mathbf{A}^*$

ie  $\mathbf{B}^* = \arg \min_{\mathbf{B}} \{\|\mathbf{B}\|_*\}$  s.t.  $\Pi(\mathbf{A} - \mathbf{B}) = 0$

## Topic Model

find low-dim representation of document from a corpus **Preprocessing** vocabulary extraction/tokenisation, filtering (of too frequent or rare words), normalisation (stemming) e.g. argue(d) reduce to arg **bag of words** counts of cooc. ignores order. sparse!  
**pLSA**

$p(w) = \sum_{z=1}^K p(w|z)p(z|d)$ , for word  $w$  in doc  $d$  and given topics  $z \in \{1 \dots K\}$   
**assumption**  $p(w|d, z) = p(w|z)$   
**Log-Likel.**  $x_{ij} = \# \text{ of } w_j \text{ in } d_i$ ,  $\mathbf{X} = x_{ij}$   
 $\ell(\mathbf{U}, \mathbf{V}) = \sum_{ij} x_{ij} \log p(w_j | d_i) = \sum_{(i,j) \in \mathcal{X}} x_{ij} \log \sum_{z=1}^K p(w_j | z) p(z | d_i) = \sum_{(i,j) \in \mathcal{X}} x_{ij} \log \sum_{z=1}^K v_{zj} u_{zi}$ , with  $u_{iz} \geq 0$ ,  $\sum_z u_{iz} = 1$  and  $v_{jz} \geq 0$ ,  $\sum_j v_{jz} = 1$

**Lower Bound** w. Jensen's inequ.  
 $q_{zij} = P(w_j \text{ in } d_i \text{ is from } z)$ ,  $\sum_z q_{ijz} = 1$

$\sum_{ij} x_{ij} \log \sum_{z=1}^K q_{ijz} \frac{u_{zi} v_{zj}}{q_{ijz}} \geq g(\mathbf{X} | \mathbf{U}, \mathbf{V}) =$

$\sum_{ij} x_{ij} \sum_{z=1}^K q_{ijz} [\log u_{zi} + \log v_{zj} - \log q_{zij}]$

**Lagrangian**  $\mathcal{L}_{\mathbf{U}, \mathbf{V}}(\alpha, \beta) = -g(\mathbf{X} | \mathbf{U}, \mathbf{V}) + \sum_j \alpha_j (\sum_z u_{zi} - 1) + \sum_z \beta_z (\sum_j v_{zj} - 1)$

## Optimal solution (EM)

E:  $q_{zij} = \frac{u_{zi} v_{zj}}{\sum_{k=1}^K u_{zi} v_{kj}} = \frac{p(w_j | z) p(z | d_i)}{\sum_{k=1}^K p(w_j | k) p(k | d_i)}$

M:  $u_{zi} = \frac{\sum_j x_{ij} q_{zij}}{\sum_j x_{ij}}$ ,  $v_{zj} = \frac{\sum_i x_{ij} q_{zij}}{\sum_{i,j} x_{ij} q_{zij}}$

conv. guaranteed but not global opt. fixed docs and words. Add a topic?

## Latent Dirichlet Allocation

$p(\mathbf{u}_i | \alpha) \propto \prod_{z=1}^K u_{zi}^{\alpha_z - 1}$  generate topic weights. for doc. w. length  $l = \sum_j x_j$

$p(\mathbf{x}, \mathbf{V}, \mathbf{u}) = \frac{l!}{\prod_j x_j!} \prod_j (\sum_z v_{zj} u_z)^{x_j}$

$p(\mathbf{x} | \mathbf{V}, \alpha) = \int p(\mathbf{x} | \mathbf{V}, \mathbf{u}) p(\mathbf{u} | \alpha) d\mathbf{u}$

## Non-Negative Matrix Factorization

factorize count matrix  $\in \mathbb{Z}_{\geq 0}^{N \times M}$   
 $\mathbf{X} = \mathbf{U}^\top \mathbf{V}$ .  $\mathbf{U}, \mathbf{V}$  non-neg. entries and  $L_1$  column normalized.

Useful to model non-negative data like images (ink) and leads to part-based representation. pLSA is a kind of NMF.

## Word Embeddings

**latent vect model**  $w \rightarrow (\mathbf{x}_w, b_w) \in \mathbb{R}^{d+1}$

## Context Models

semantic from by co-occurrences, e.g. skip-gram  $p(w|w')$   $w$  in context of  $w'$

**Log-likelihood**  
 $\mathcal{L}(\theta|\mathbf{w}) = \sum_{t=1}^T \sum_{\delta \in \mathcal{I}} \log p_{\theta}(w^{(t+\delta)}|w^{(t)})$   
 $\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|\mathbf{w})$  large cardinality

**Log-bilinear model**  
 $\log p(w|w') = \langle \mathbf{x}_w, \mathbf{x}_{w'} \rangle + b_w + c_{w'} \uparrow$   
 $b_w \uparrow \Rightarrow p_{\theta}(w|w') \uparrow$   
 $\mathcal{L}(\mathbf{x}_w, \mathbf{x}_{w'}) \downarrow \Rightarrow p_{\theta}(w|w') \uparrow$

**Softmax**  $p_{\theta}(w|w') = \frac{\exp[\langle \mathbf{x}_w, \mathbf{x}_{w'} \rangle + b_w]}{Z_{\theta}(w')}$

where  $Z_{\theta}(w') = \sum_{v \in \mathcal{V}} \exp[\langle \mathbf{x}_v, \mathbf{x}_{w'} \rangle + b_v]$

**Context Vectors** input & output context vectors to get rid of the bilinearity.  $\log p_{\theta}(w|w') = \langle \mathbf{x}_w, \mathbf{y}_{w'} \rangle + b_w$

**Negative Sampling**  
 $\Delta^+$  observed pair,  $\Delta^- \sim p_n$  where  $(w_i, w_j) \sim p_n(i, j)$  rand. context words  $w_j \propto P(w_j)^{\alpha}$ ,  $\alpha = \frac{3}{4}$  oversample by  $k \leq 20$

**Maximize logistic Regression**  
 $\mathcal{L}(\theta) = \sum_{(i,j) \in \Delta^+} \log \sigma(\langle \mathbf{x}_i, \mathbf{y}_j \rangle) + \sum_{(i,j) \in \Delta^-} \log \sigma(-\langle \mathbf{x}_i, \mathbf{y}_j \rangle)$ ,  $\sigma(z) = \frac{1}{1 + \exp(-z)}$

**Bayes optimal discriminant for  $\mathcal{L}$**   
 $\langle \mathbf{x}_i, \mathbf{y}_j \rangle = \log \frac{p(w_i, w_j)}{p_n(w_i, w_j)} + \log \frac{\kappa}{1 - \kappa}$ ,  $\kappa = \frac{1}{k+1}$

**pointwise mutual information**  
for  $k=1$  and  $p_n(w_i, w_j) = p(w_i)p(w_j)$   
 $\langle \mathbf{x}_i, \mathbf{y}_j \rangle \approx \text{PMI}(w_i, w_j)$

**GloVe**  
co-occurrence matrix  $\mathbf{N} = \{n_{ij} | i, j \in \mathcal{V}\}$   
 $n_{ij} = \# \text{ occ. of } w_i \in \mathcal{V} \text{ in ctx of } w_j \in \mathcal{C}$   
 $\mathbf{N}$  is sparse & computed in one pass.  
**Objective** least square  $\mathcal{H}(\theta, \mathbf{N}) = \sum_{i,j} f(n_{ij})(\log n_{ij} - \log \tilde{p}_{\theta}(w_i|w_j))^2$   
**unnormalised** distribution (model):  $\tilde{p}_{\theta}(w_i|w_j) = \exp[\langle \mathbf{x}_i, \mathbf{y}_j \rangle + b_i + c_j]$ ,  
 $n_{ij}$  is the target and  $f$  weight func  
 $f(n) = \min[1, (\frac{n}{n_{\max}})^{\alpha}]$ ,  $\alpha \in [0, 1]$ ,  
limits the influence of large and small noisy counts

**Solves  $\min_{\mathbf{X}, \mathbf{Y}} \|\mathbf{N} - \mathbf{X}^T \mathbf{Y}\|_F^2$**  if  $f=1$  But non-convex, hard to find optimum, full gradient descent too expensive to compute.  $\Rightarrow$  use SGD

**Mixtures**

**K-Means**  
**Objective**  $J(\mathbf{U}, \mathbf{Z}) = \sum_i^N \sum_j^K z_{ij} \|\mathbf{x}_i - \mathbf{u}_j\|^2$   
 $= \|\mathbf{X} - \mathbf{U}^T \mathbf{Z}\|^2$ ,  $\sum_{k=1}^K z_{ki} = 1$   
 $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  data matrix  
 $\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_K] \in \mathbb{R}^{D \times K}$  centroids  
**Recursive optimal assignments**  
 $E: z_{ij}^* = \begin{cases} 1 & \text{if } j = \arg \min_k \|\mathbf{x}_i - \mathbf{u}_k\|^2 \end{cases}$

**M:**  $\nabla_{\mathbf{u}} J(\mathbf{U}, \mathbf{Z}) \stackrel{!}{=} 0 \Rightarrow \mathbf{u}_j^* = \frac{\sum_{i=1}^N z_{ij} \mathbf{x}_i}{\sum_{i=1}^N z_{ij}}$

guaranteed convergence but non-convex objective

K-means solves  $\arg \min_{\mathbf{Z}} \|\mathbf{X} - \mathbf{U} \mathbf{Z}\|_F^2$

**K-Means++**  
Init. with incremental  $D^2$  sampling  
sample first centroid rand.  $U_1 = \{\mathbf{x}_I\}$   
 $D_i = \min_{\mathbf{u} \in U_k} \|\mathbf{x}_i - \mathbf{u}\|$   $U_{k+1} = U_k \cup \{\mathbf{x}_I\}$   
for  $I \sim \text{Cat}(\mathbf{p})$ ,  $p_i = D_i^2 / \sum_{j=1}^N D_j^2$

More expensive but better results.

**K-Means Core Sets**  
sample core set of  $m$  centroids  
 $I \sim \text{Cat}(\mathbf{p})$ ,  $p_i = 1/2N + D_i^2/2 \sum_{j=1}^N D_j^2$

where  $D_i^2 = \|\mathbf{x}_i - \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i\|^2$   
give each sample weight  $1/mp_i$  and weighted K-Mean on this core-set.

**Finite Mixture Model**  
Probabilistic assignment to clusters.  
 $p(\mathbf{x}, \theta) = \sum_{j=1}^K \pi_j p(\mathbf{x}; \theta_j)$   
where  $\theta = (\pi, \theta_1, \dots, \theta_K)$   
and  $\pi \geq 0$ ,  $\sum_{i=1}^K \pi_i = 1$

**e.g Gaussian Mixture Model**  
 $p(\mathbf{x}; \theta_j) = p(\mathbf{x}; \mu_j, \Sigma_j)$  (normal dist.)  
Complete data distribution:  
 $p(\mathbf{x}, \mathbf{z}; \theta) = \prod_{j=1}^K (\pi_j p(\mathbf{x}; \theta_j))^{z_j}$  where  $z_j$  latent and  $P(z_j = 1) = \pi_j$   
Posterior Probabilities:  $p(\mathbf{z}_k = 1 | \mathbf{x}) = \frac{p(\mathbf{z}_k = 1) p(\mathbf{x} | \mathbf{z}_k = 1)}{\sum_{l=1}^K p(\mathbf{z}_l = 1) p(\mathbf{x} | \mathbf{z}_l = 1)} = \frac{\pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{l=1}^K \pi_l \mathcal{N}(\mathbf{x} | \mu_l, \Sigma_l)}$

**Maximum Likelihood for MM**  
 $\arg \max_{\theta} \sum_{i=1}^N \log[\sum_{j=1}^K \pi_j p(\mathbf{x}_i; \theta_j)]$   
has no closed form solution  
 $\log[\sum_{j=1}^K \pi_j p(\mathbf{x}_i; \theta_j)] = \log[\sum_{j=1}^K q_j \frac{\pi_j p(\mathbf{x}_i; \theta_j)}{q_j}] \geq \sum_{j=1}^K q_j [\log p(\mathbf{x}; \theta_j) + \log \pi_j - \log q_j]$

**Lagrangian**  
 $\mathcal{L} = \max_{\mathbf{q}} \{ \sum_{j=1}^K q_j [\log p(\mathbf{x}; \theta_j) + \log \pi_j - \log q_j] + \lambda (\sum_{j=1}^K q_j - 1) \}$

**E-Step** compute assignments  
 $\nabla q_j \stackrel{!}{=} 0 \Rightarrow q_j^* = \frac{\pi_j p(\mathbf{x}; \theta_j)}{\sum_{l=1}^K \pi_l p(\mathbf{x}; \theta_l)} = p(z_j = 1 | \mathbf{x})$

**M-Step** optimize clusters  
 $\mu_j^* = \frac{\sum_{i=1}^N q_{ij} \mathbf{x}_i}{\sum_{i=1}^N q_{ij}}$ ,  $\pi_j := \frac{1}{N} \sum_{i=1}^N q_{ij}$ , and  
 $\Sigma_j = \frac{\sum_{i=1}^N q_{ij} (\mathbf{x}_i - \mu_j)(\mathbf{x}_i - \mu_j)^T}{\sum_{i=1}^N q_{ij}}$  EM requires more steps and computation to reach convergence than K-Mean

**Singularities of GMM** happen when a cluster shrinks to fit exactly one data point.  $\log \text{likelihood} \rightarrow \infty$ . Bad convergence

**Neural Networks**  
 $F^{\sigma}(\mathbf{x}; \mathbf{W}) = \sigma(\mathbf{W} \mathbf{x}) \Rightarrow F_j^{\sigma}(\mathbf{x}; \mathbf{W}) = \sigma(\mathbf{w}_j^T \mathbf{x})$ ,  
w.  $\mathbf{W} = (\mathbf{w}_1 \cdots \mathbf{w}_m)^T$  mapping  $\mathbb{R}^n \rightarrow \mathbb{R}^m$  between 2 layers  
 $\mathbf{x}^{(l)} = \sigma(\mathbf{W}^{(l)} \mathbf{x}^{(l-1)})$ ,  $1 \leq l \leq L$

**Activation Functions**  
**Sig.**  $\sigma(x) = \frac{1}{1 + \exp^{-x}}$ ,  $\nabla_x \sigma(x) = \sigma(x)(1 - \sigma(x))$ ,  
 $\sigma^{-1}(x) = \log(1/(1-x))$ ,  $1 - \sigma(x) = \sigma(-x)$   
**ReLU**  $R(x) = \max\{0, x\}$  has simple derivative on  $\mathbb{R} - \mathbf{0}$  and reduces vanishing gradient problem

**Output Layer**  
**Linear Regression**  $\mathbf{y} = \mathbf{W}^{(L)} \mathbf{x}^{(L-1)}$   
**Logistic** binary classif. (one output)  
 $y_1 = P(Y = 1 | \mathbf{x}) = \frac{1}{1 + \exp[-\mathbf{w}^T \mathbf{x}]}$

**Soft-Max** K-Multiclass  
 $y_k = P(Y = k | \mathbf{x}) = \frac{e[\mathbf{w}_k^T \mathbf{x}]}{\sum_{j=1}^K e[\mathbf{w}_j^T \mathbf{x}]}$

**Loss-Functions**  
**Squared**  $l(y^*; y) = \frac{1}{2} (y^* - y)^2$   
**Cross-Entropy** for classification  
 $l(y^*; y) = -y^* \log y - (1 - y^*) \log(1 - y)$   
**Empirical Risk**  
 $\mathcal{L}(\theta; \mathcal{X}) = \frac{1}{T} \sum_{t=1}^T l(y_t; y(\mathbf{x}_t; \theta))$  for weights  $\theta = (\mathbf{W}^{(1)} \dots \mathbf{W}^{(L)})$  and training data  $\mathcal{X} = \{(\mathbf{x}_t, y_t), 1 \leq t \leq T\}$

**Regularization**  
favors smaller weights  
**L<sub>2</sub>:**  $\mathcal{L}_{\lambda}(\theta; \mathcal{X}) = \mathcal{L}(\theta; \mathcal{X}) + \frac{\lambda}{2} \|\theta\|_2^2$   
**Dropout** training with noise  
**Backpropagation**  
costs  $\mathcal{O}(n)$  for NN with  $n$  nodes  
 $\frac{\partial x_i^{(l)}}{\partial x_k^{(l-n)}} = \sum_j \frac{\partial x_j^{(l)}}{\partial x_j^{(l-1)}} \frac{\partial x_j^{(l-1)}}{\partial x_k^{(l-n)}} = \sum_j \mathbf{J}_{ij}^{(l)} \frac{\partial x_j^{(l-1)}}{\partial x_k^{(l-n)}}$   
 $\frac{\partial x_i^{(l)}}{\partial \mathbf{x}^{(l-n)}} = \sum_j \mathbf{J}^{(l)} \frac{\partial \mathbf{x}^{(l-1)}}{\partial \mathbf{x}^{(l-n)}} = \mathbf{J}^{(l)} \dots \mathbf{J}^{(l-n+1)}$   
Backprop:  $\nabla_{\mathbf{x}^{(l)}}^T l = \nabla_{\mathbf{T}^l}^T l \cdot \mathbf{J}^{(l)} \dots \mathbf{J}^{(l+1)}$   
 $\frac{\partial l}{\partial w_{ij}^{(l)}} = \frac{\partial l}{\partial x_i^{(l)}} \frac{\partial x_i^{(l)}}{\partial w_{ij}^{(l)}} = \sigma'(\mathbf{w}_i^{(l)T} \mathbf{x}^{(l-1)}) x_j^{(l-1)}$

**CNN**  
**Convolutional Layers**  $F_{n,m}(\mathbf{x}; \mathbf{w}) = \sigma(b + \sum_{k=-i}^i \sum_{l=-i}^i w_{kl} x_{n+k, m+l})$   
Weight sharing and shift-invariant filtering thus less parameters and computational power required  
**Pooling** Take avg or max over wind. Reduce size or extract features.

**Generative Models**  
**VAE - Variational Autoencoder**  
find  $\mathbf{x} \in \mathbb{R}^n$  by sampling  $\mathbf{z} \in \mathbb{R}^m$  from simple distribution and set  $\mathbf{x} = F_{\theta}(\mathbf{z})$  where  $F_{\theta} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ , deterministic DNN, since  $\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] = \mathbb{E}_{\mathbf{z}}[f(F_{\theta}(\mathbf{z}))]$   
Requires  $F_{\theta}^{-1}$  often impossible to get.  
**ELBO - evidence lower bound**  
learn parameters of distribution  $p_{\theta}$  instead of deterministic  $F_{\theta}$   
 $\log p_{\theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x})] = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log \frac{p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z}) q_{\theta}(\mathbf{z} | \mathbf{x})}{p_{\theta}(\mathbf{z} | \mathbf{x}) q_{\theta}(\mathbf{z} | \mathbf{x})}] = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - KL(q_{\phi}(\mathbf{z} | \mathbf{x}) | p_{\theta}(\mathbf{z} | \mathbf{x})) \geq \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - KL(q_{\phi}(\mathbf{z} | \mathbf{x}) | p_{\theta}(\mathbf{z}))] = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [-\log q_{\phi}(\mathbf{z} | \mathbf{x}) + \log p_{\theta}(\mathbf{x}, \mathbf{z})] = ELBO(\phi, \theta) = \mathcal{L}(\mathbf{x}; \theta)$   
maxim. wrt  $\phi$  for inference model  
maxim. wrt  $\theta$  for generative model  
**stochastic approximation**  
Update for generative model:  
 $\nabla_{\theta} \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] = \mathbb{E}_{q_{\phi}} [\nabla_{\theta} \log p(\mathbf{x} | \mathbf{z})] \approx \frac{1}{L} \sum_{r=1}^L \nabla_{\theta} \log p(\mathbf{x} | \mathbf{z}^{(r)})$ ,  $\mathbf{z}^{(r)} \sim q_{\phi}(\cdot | \mathbf{x})$  (by Monte Carlo approximation)  
Update for inference model:  
 $\nabla_{\phi} \mathbb{E}_{q_{\phi}} [\mathcal{L}(\mathbf{x}, \mathbf{z})] = \int \mathcal{L}(\mathbf{x}, \mathbf{z}) \nabla_{\phi} q_{\phi}(\mathbf{z} | \mathbf{x}) d\mathbf{z} = \mathbb{E}_{q_{\phi}} [\mathcal{L}(\mathbf{x}, \mathbf{z}) \nabla_{\phi} \log q_{\phi}(\mathbf{z} | \mathbf{x})]$  hard!  
variance in gradient usually high  
**Reparametrization Trick**  
 $\mathbf{z} = g_{\phi}(\zeta | \mathbf{x})$ ,  $\zeta \sim$  simple distribution  
 $\nabla_{\phi} \mathbb{E}_{q_{\phi}} [\mathcal{L}(\mathbf{x}, \mathbf{z})] \approx \frac{1}{L} \sum_{r=1}^L \nabla_{\phi} \mathcal{L}(\mathbf{x}, g_{\phi}(\zeta^{(r)}))$   
easier to sample and differentiable by  $\phi$  since the  $\mathbf{z}$  are now deterministic.  
backprop. can be used

**GAN-Generative Adversarial Net.**  
generator  $G$  and discriminator  $D$   
learn to fool each others.  
 $\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$   
Hard to train, might not converge/learn to generate a few samples generates sharper images since it can predict high frequency details. But can only generate/sample

**Sparse Coding**  
**Basis Transformation**  
signals often allow sparse representation. Vanishing coeffs due to regularity. **Find orthonormal dictionary**  $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_L\}$  and change of basis.  
 $\mathbf{x} = \mathbf{U} \mathbf{z} \rightarrow \mathbf{z} = \mathbf{U}^T \mathbf{x} \rightarrow \hat{\mathbf{x}} = \mathbf{U} \hat{\mathbf{z}}$

$\hat{\mathbf{x}} = \sum_{d \in \sigma} \hat{z}_d(\mathbf{x}) \mathbf{u}_d$ ,  $\hat{z}_d(\mathbf{x}) = \langle \mathbf{x}, \mathbf{u}_d \rangle$   
Error  $\|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \sum_{d \notin \sigma} \langle \mathbf{x}, \mathbf{u}_d \rangle^2$  e.g. fourier transform denoises well but not for localised signals **O**( $D \log D$ )

**Haar Wavelets**  
 $\psi_{n,k}(t) = 2^{n/2} \psi(2^n t - k)$ ,  $0 \leq k \leq 2^n$   
good for localized signal but poor for denoising smooth signal **O**( $D \log D$ )

**Overcomplete Dictionaries**  
dictionaries generalisation  
**overcompleteness**  
 $\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_L]$ ,  $\mathbf{U} \in \mathbb{R}^{D \times L}$ ,  $L > D$   
 $\mathbf{z} \in \mathbb{R}^L$  contains coeffs. of signal  $\mathbf{x} \in \mathbb{R}^D$  in base  $\mathbf{U}$ .  $\mathbf{z}$  is not unique.  
Search  $\mathbf{z}^* = \arg \min_{\mathbf{z} \in \mathbb{R}^L} \|\mathbf{z}\|_0$  s.t.  $\mathbf{U} \mathbf{z} = \mathbf{x}$   
Problem is NP-hard/ill-posed.

**Coherence**  
 $L/D \uparrow \Rightarrow \text{sparsity} \uparrow$ ,  $\mathbf{u}_i$  dependency  $\uparrow$   
 $m(\mathbf{U}) = \max_{i,j: i \neq j} |\mathbf{u}_i^T \mathbf{u}_j|$  coherence  
 $m(\mathbf{B}) = 0$  for orthog.  $\mathbf{B}$   
 $m([\mathbf{B} \mathbf{u}] \geq \frac{1}{\sqrt{D}}$  if  $\mathbf{u}$  added to  $\mathbf{B}$

**Matching Pursuit** greedy algorithm at each step chose dimension with max projection onto residual.  
Init  $\mathbf{r}_0 = \mathbf{x}$ ,  $\hat{\mathbf{x}}_0 = \mathbf{0}$ . Then repeat:  
find  $j^* = \arg \max_j |\langle \mathbf{r}_i, \mathbf{u}_j \rangle|$   
 $\hat{\mathbf{x}}_{i+1} \leftarrow \hat{\mathbf{x}}_i + \langle \mathbf{r}_i, \mathbf{u}_{j^*} \rangle \mathbf{u}_{j^*}$   
 $\mathbf{r}_{i+1} \leftarrow \mathbf{r}_i - \langle \mathbf{r}_i, \mathbf{u}_{j^*} \rangle \mathbf{u}_{j^*}$   
**Convergence** greedily reduces residual energy at each step.  
 $\|\mathbf{r}_i\|_2^2 = \|\mathbf{r}_{i+1}\|_2^2 + |\langle \mathbf{r}_i, \mathbf{u}_{j^*} \rangle|^2$  (by energy conservation and linearity)

**Convex Optimization** with  $l_1$ -norm  
 $\mathbf{z}^* = \arg \min_{\mathbf{z} \in \mathbb{R}^L} \|\mathbf{z}\|_1$  s.t.  $\mathbf{U} \mathbf{z} = \mathbf{x}$   
can approximate  $l_0$  even same results

**Dictionary Learning**  
learn one dictionary for  $\mathbf{x}_1 \cdots \mathbf{x}_N$  obj.  
 $(\mathbf{U}^*, \mathbf{Z}^*) \in \arg \min_{\mathbf{U}, \mathbf{Z}} \|\mathbf{x} - \mathbf{U} \mathbf{Z}\|_F^2$  not jointly convex in  $(\mathbf{U}, \mathbf{Z})$

**Greedy Convex Minimization**  
**1.Coding**  $\mathbf{Z}^{t+1} \in \arg \min_{\mathbf{Z}} \|\mathbf{X} - \mathbf{U}^T \mathbf{Z}\|_F^2$  where  $\mathbf{Z}$  sparse and  $\mathbf{U}$  fixed  
column separable,  $\forall n = 1 \cdots N$   
 $\mathbf{z}_n^{t+1} \in \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0 \text{ s.t. } \|\mathbf{x}_n - \mathbf{U}^T \mathbf{z}\|^2 \leq \sigma \|\mathbf{x}_n\|_2$   
**2.Update**  $\mathbf{U}^{t+1} \in \arg \min_{\mathbf{U}} \|\mathbf{X} - \mathbf{U} \mathbf{Z}^{t+1}\|_F^2$  where  $\|\mathbf{u}_l\|_2 = 1, \forall l$  and  $\mathbf{Z}$  fixed.  
not separable  $\mathbf{R}_l^t$  residual of atom  $\mathbf{u}_l$   
 $\|\mathbf{X} - [\mathbf{u}_1^t \cdots \mathbf{u}_l^t \cdots \mathbf{u}_L^t] \mathbf{Z}^{t+1}\|_F^2 = \|\mathbf{X} - (\sum_{\ell \neq l} \mathbf{u}_{\ell}^t (\mathbf{z}_{\ell}^{t+1})^T + \mathbf{u}_l^t \mathbf{z}_{\ell}^{t+1 T})\|_F^2 = \|\mathbf{R}_{\ell}^t - \mathbf{u}_{\ell}^t (\mathbf{z}_{\ell}^{t+1 T})\|$  approx. by SVD.  
 $\mathbf{R}_{\ell}^t = \tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{V}}^T$ ,  $\mathbf{u}_{\ell}^t = \tilde{\mathbf{u}}_1$  first. sing vec.