

## 0 Essentials

### Matrix/Vector

**Vectors:** Unit vector:  $u^\top u = 1$  Orthogonal vectors:  $u^\top v = 0$  **Range, Kernel, Nullity:**  $\text{range}(\mathbf{A}) = \{\mathbf{z} | \exists \mathbf{x} : \mathbf{z} = \mathbf{A}\mathbf{x}\} = \text{span}(\text{columns of } \mathbf{A})$   
 $\text{rank}(\mathbf{A}) = \dim(\text{range}(\mathbf{A}))$   $\text{kernel}(\mathbf{A}) = \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{0}\}$  (spans nullspace)  $\text{nullity}(\mathbf{A}) = \dim(\text{kernel}(\mathbf{A}))$   
**Ranks:**  $\text{rank}(XY) \leq \text{rank}(X) \forall X \in \mathbb{R}^{m \times n}, Y \in \mathbb{R}^{n \times k}$   
eq. if  $Y \in \mathbb{R}^{n \times n}, \text{rank}(Y) = n$  **Rank-nullity Theorem:**  $\dim(\text{kernel}(\mathbf{A})) + \dim(\text{range}(\mathbf{A})) = n$

**Orthogonal mat.**  $\mathbf{A}^{-1} = \mathbf{A}^\top$ ,  $\mathbf{A}\mathbf{A}^\top = \mathbf{A}^\top \mathbf{A} = \mathbf{I}$ ,  $\det(\mathbf{A}) \in \{+1, -1\}$ ,  $\det(\mathbf{A}^\top \mathbf{A}) = 1$ , preserves inner product, norm, distance, angle, rank, matrix orthogonality **Outer Product:**  $\mathbf{u}\mathbf{v}^\top$ ,  $(\mathbf{u}\mathbf{v}^\top)_{i,j} = \mathbf{u}_i \mathbf{v}_j$   
**Inner Product:**  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n \mathbf{x}_i \mathbf{y}_i$ .  $\langle \mathbf{x} \pm \mathbf{y}, \mathbf{x} \pm \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{x} \rangle \pm 2\langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{y}, \mathbf{y} \rangle$

**( $\mathbf{u}_i^\top \mathbf{v}_j$ )  $\mathbf{v}_j = (\mathbf{v}_j \mathbf{v}_j^\top) \mathbf{u}_i$  Cross product:**  $\vec{a} \times \vec{b} = (a_2 b_3 - a_3 b_2, a_3 b_1 - a_1 b_3, a_1 b_2 - a_2 b_1)^\top$   
**Trace:**  $\text{trace}(\mathbf{XYZ}) = \text{trace}(\mathbf{ZXY})$   
**Transpose:**  $(\mathbf{A}^\top)^{-1} = (\mathbf{A}^{-1})^\top$ ,  $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$ ,  $(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top$

**Cauchy-Schwarz inequality:**  $|\langle \mathbf{u}, \mathbf{v} \rangle| \leq \|\mathbf{u}\| \|\mathbf{v}\|$   
**Jensen inequality:** for convex function  $f$ , non negative  $\lambda_i$  s.t.  $\sum_{i=1}^n \lambda_i = 1$ :  $f(\sum_{i=1}^n \lambda_i x_i) \leq \sum_{i=1}^n \lambda_i f(x_i)$  Note: for concave, inequality sign switches **Convexity:**  $f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y), \forall \theta \in [0, 1]$  **Least Squares equations:**  $\arg \min_{\beta \in \mathbb{R}^{p+1}} \|\mathbf{y} - \mathbf{X}\beta\|^2, \hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

**Einstein matrix notation:**  $(\mathbf{A} \cdot \mathbf{B})_{ij} = \sum_{k=1}^n \mathbf{A}_{ik} \cdot \mathbf{B}_{kj}$

**Kullback-Leibler:**  $KL(P||Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)}$

### Norms

- $\|\mathbf{x}\|_0 = |\{i | x_i \neq 0\}|$
- $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n \mathbf{x}_i^2} = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$
- $\|\mathbf{u} - \mathbf{v}\|_2 = \sqrt{(\mathbf{u} - \mathbf{v})^\top (\mathbf{u} - \mathbf{v})}$
- $\|\mathbf{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}; \|\mathbf{x}\|_\infty = \max_{i=1, \dots, n} |x_i|$
- $\|\mathbf{M}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n \mathbf{m}_{i,j}^2} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2} = \|\sigma(\mathbf{A})\|_2 = \sqrt{\text{trace}(\mathbf{M}^\top \mathbf{M})}$
- $\|\mathbf{M}\|_G = \sqrt{\sum_{i,j} g_{ij} x_{ij}^2}$  (weighted Frobenius)
- $\|\mathbf{M}\|_1 = \sum_{i,j} |m_{i,j}|$
- $\|\mathbf{M}\|_2 = \sigma_{\max}(\mathbf{M}) = \|\sigma(\mathbf{M})\|_\infty$  (spectral)
- $\|\mathbf{M}\|_p = \max_{\mathbf{v} \neq 0} \frac{\|\mathbf{M}\mathbf{v}\|_p}{\|\mathbf{v}\|_p}$
- $\|\mathbf{M}\|_* = \sum_{i=1}^{\min(m,n)} \sigma_i = \|\sigma(\mathbf{A})\|_1$  (nuclear)

### Derivatives

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{b}^\top \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{b}) = \mathbf{b} \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x}$$
$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = (\mathbf{A}^\top + \mathbf{A}) \mathbf{x} \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{b}^\top \mathbf{A} \mathbf{x}) = \mathbf{A}^\top \mathbf{b}$$
$$\frac{\partial}{\partial \mathbf{X}} (\mathbf{c}^\top \mathbf{X} \mathbf{b}) = \mathbf{c} \mathbf{b}^\top \quad \frac{\partial}{\partial \mathbf{X}} (\mathbf{c}^\top \mathbf{X}^\top \mathbf{b}) = \mathbf{b} \mathbf{c}^\top$$

$$\frac{\partial}{\partial \mathbf{x}} (\|\mathbf{x} - \mathbf{b}\|_2) = \frac{\mathbf{x} - \mathbf{b}}{\|\mathbf{x} - \mathbf{b}\|_2} \quad \frac{\partial}{\partial \mathbf{x}} (\|\mathbf{x}\|_2) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x}$$
$$\frac{\partial}{\partial \mathbf{X}} (\|\mathbf{X}\|_F^2) = 2\mathbf{X} \quad \frac{\partial}{\partial \mathbf{x}} \log(x) = \frac{1}{x}$$

### Eigendecomposition

$\mathbf{A} \in \mathbb{R}^{N \times N}$  then  $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$  with  $\mathbf{Q} \in \mathbb{R}^{N \times N}$ .  
if fullrank:  $\mathbf{A}^{-1} = \mathbf{Q} \mathbf{\Lambda}^{-1} \mathbf{Q}^{-1}$  and  $(\mathbf{\Lambda}^{-1})_{i,i} = \frac{1}{\lambda_i}$ .

if  $\mathbf{A}$  symmetric:  $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$  ( $\mathbf{Q}$  orthogonal). Eigenvalue  $\lambda$ : solve  $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$  Eigenvector  $\mathbf{v}$ : solve  $(\mathbf{A} - \lambda \mathbf{I}) * \mathbf{v} = \vec{0}$

### Probability / Statistics

$\bullet P(x) := Pr[X = x] = \sum_{y \in Y} P(x, y) \bullet P(x|y) := Pr[X = x | Y = y] := \frac{P(x, y)}{P(y)}$ , if  $P(y) > 0 \bullet \forall y \in Y : \sum_{x \in X} P(x|y) = 1$  (property for any fixed  $y$ )  $\bullet P(x, y) = P(x|y)P(y) \bullet \text{posterior } P(A|B) = \frac{\text{prior } P(A) \times \text{likelihood } P(B|A)}{\text{evidence } P(B)}$  (Bayes' rule)  $\bullet P(x|y) = P(x) \Leftrightarrow P(y|x) = P(y)$  (iff  $X, Y$  independent)  $\bullet P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i)$  (iff IID)  $\bullet \text{Variance } Var[X] := E[(X - \mu_x)^2] := \sum_{x \in X} (x - \mu_x)^2 P(x) = E(X^2) - E(X)^2$   $Var(aX) = a^2 Var(X) \bullet \text{expectation } \mu_x := E[X] := \sum_{x \in X} x P(x) \bullet E[X + Y] = E[X] + E[Y] \bullet \text{standard deviation } \sigma_x := \sqrt{Var[X]}$

### Lagrangian Multipliers

Minimize  $f(\mathbf{x})$  s.t.  $g_i(\mathbf{x}) \leq 0, i = 1, \dots, m$  (**inequality constr.**) and  $h_i(\mathbf{x}) = \mathbf{a}_i^\top \mathbf{x} - b_i = 0$  or  $h_i(\mathbf{x}) = \sum_w x_{w,i} - b_i = 0, i = 1, \dots, p$  (**equality constraint**)  
 $L(\mathbf{x}, \alpha, \beta) := f(\mathbf{x}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{x}) + \sum_{i=1}^p \beta_i h_i(\mathbf{x})$

### 1 Principal Component Analysis

$\mathbf{X} \in \mathbb{R}^{D \times N}$ .  $N$  observations,  $K$  rank.

- Empirical Mean:  $\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$ .
- Center Data:  $\bar{\mathbf{X}} = \mathbf{X} - [\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}}] = \mathbf{X} - \mathbf{M}$ .
- Cov.:  $\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^\top = \frac{1}{N} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top$ .
- Eigenvalue Decomposition:  $\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ .
- Select  $K < D$ , only keep  $\mathbf{U}_K, \lambda_K$ .
- Transform data onto new Basis:  $\bar{\mathbf{Z}}_K = \mathbf{U}_K^\top \bar{\mathbf{X}}$ .
- Reconstruct to original Basis:  $\tilde{\bar{\mathbf{X}}} = \mathbf{U}_K \bar{\mathbf{Z}}_K$ .

8. Reverse centering:  $\tilde{\bar{\mathbf{X}}} = \tilde{\bar{\mathbf{X}}} + \mathbf{M}$ .

For compression save  $\mathbf{U}_K, \bar{\mathbf{Z}}_K, \bar{\mathbf{x}}$ .  
 $\mathbf{U}_k \in \mathbb{R}^{D \times K}, \Sigma \in \mathbb{R}^{D \times D}, \bar{\mathbf{Z}}_K \in \mathbb{R}^{K \times N}, \bar{\mathbf{x}} \in \mathbb{R}^{D \times N}$   
**Calculation of:**  $\text{var}(X) = \frac{1}{N} \sum_{n=1}^N (X_i - \bar{X})^2$

### Iterative View

Residual  $r_i$ :  $x_i - \tilde{x}_i = \mathbf{I} - \mathbf{u} \mathbf{u}^\top x_i$   
Cov of  $r$ :  $\frac{1}{n} \sum_{i=1}^n (\mathbf{I} - \mathbf{u} \mathbf{u}^\top) x_i x_i^\top (\mathbf{I} - \mathbf{u} \mathbf{u}^\top)^\top = (\mathbf{I} - \mathbf{u} \mathbf{u}^\top) \Sigma (\mathbf{I} - \mathbf{u} \mathbf{u}^\top)^\top = \Sigma - 2 \Sigma \mathbf{u} \mathbf{u}^\top + \mathbf{u} \mathbf{u}^\top \Sigma \mathbf{u} \mathbf{u}^\top = \Sigma - \lambda \mathbf{u} \mathbf{u}^\top$ 

- Find principal eigenvector of  $(\Sigma - \lambda \mathbf{u} \mathbf{u}^\top)$
- which is the second eigenvector of  $\Sigma$
- iterating to get  $d$  principal eigenvector of  $\Sigma$

### Power Method

Power iteration:  $v_{t+1} = \frac{\mathbf{A} v_t}{\|\mathbf{A} v_t\|}, \lim_{t \rightarrow \infty} v_t = u_1$   
Assuming  $\langle u_1, v_0 \rangle \neq 0$  and  $|\lambda_1| > |\lambda_j| (\forall j \geq 2)$

### Reconstruction Proof Sketch

Given:  $\tilde{X} = U_K U_K^\top \bar{X}$  To prove: squared reconstruction error is the sum of the lowest  $D - K$  eigenvalues of  $\Sigma$ .  $\text{err} = 1/N \sum_{i=1}^N \|\tilde{x}_i - \bar{x}_i\|_2^2 = 1/N \|\tilde{X} - \bar{X}\|_F^2 = 1/N \|(U_K U_K^\top - \mathbf{I}_d) \bar{X}\|_F^2 = 1/N * \text{trace}((U_K U_K^\top - \mathbf{I}_d) \bar{X} \bar{X}^\top (U_K U_K^\top - \mathbf{I}_d)^\top) = 1/N * \text{trace}([U_K; 0] - U) \mathbf{\Lambda} ([U_K; 0] - U)^\top)$

### 2 Singular Value Decomposition

$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^\top = \sum_{k=1}^{\text{rank}(\mathbf{A})} d_{k,k} u_k (v_k)^\top$   
 $\mathbf{A} \in \mathbb{R}^{N \times P}, \mathbf{U} \in \mathbb{R}^{N \times N}, \mathbf{D} \in \mathbb{R}^{N \times P}, \mathbf{V} \in \mathbb{R}^{P \times P}$   
 $\mathbf{U}^\top \mathbf{U} = \mathbf{I} = \mathbf{V}^\top \mathbf{V}$  ( $\mathbf{U}, \mathbf{V}$  orthonormal)  
 $\mathbf{U}$  columns are eigvecs of  $\mathbf{A} \mathbf{A}^\top$ ,  $\mathbf{V}$  columns are eigvecs of  $\mathbf{A}^\top \mathbf{A}$ ,  $\mathbf{D}$  diag. elements are singular values.  
 $(\mathbf{D}^{-1})_{i,i} = \frac{1}{d_{i,i}}$  (don't forget to transpose)

- calculate  $\mathbf{A}^\top \mathbf{A}$ .
- calculate eigvals of  $\mathbf{A}^\top \mathbf{A}$ , the square root of them, in descending order, are the diagonal elements of  $\mathbf{D}$ .
- calc. eigvecs of  $\mathbf{A}^\top \mathbf{A}$  using eigvals resulting in the columns of  $\mathbf{V}$ .
- calculate the missing matrix:  $\mathbf{U} = \mathbf{A} \mathbf{V} \mathbf{D}^{-1}$ .
- normalize each column of  $\mathbf{U}$  and  $\mathbf{V}$ .

### Low-Rank approximation

Use only  $K$  largest eigvals (and corresp. eigvecs).  
 $\tilde{\mathbf{A}}_{i,j} = \sum_{k=1}^K \mathbf{U}_{i,k} \mathbf{D}_{k,k} \mathbf{V}_{j,k} = \sum_{k=1}^K \mathbf{U}_{i,k} \mathbf{D}_{k,k} (\mathbf{V}^\top)_{k,j}$ .

### Echart-Young Theorem

$\mathbf{A}_k = \arg \min_{\text{rank}(\mathbf{B})=k} \|\mathbf{A} - \mathbf{B}\|_F^2$  (not convex)

$$\min_{\text{rank}(\mathbf{B})=K} \|\mathbf{A} - \mathbf{B}\|_F^2 = \|\mathbf{A} - \mathbf{A}_K\|_F^2 = \sum_{r=K+1}^{\text{rank}(\mathbf{A})} \sigma_r^2$$
$$\min_{\text{rank}(\mathbf{B})=K} \|\mathbf{A} - \mathbf{B}\|_2 = \|\mathbf{A} - \mathbf{A}_K\|_2 = \sigma_{K+1}$$

### 3 Matrix Approximation & Reconstruction

$\min_{\text{rank}(\mathbf{B})=k} \|\sum_{(i,j) \in I} (a_{ij} - b_{ij})^2\|, I = \{(i, j) : ob.\}$

### Alternating Least Squares

$$f(U, v_i) = \sum_{(i,j) \in I} (a_{i,j} - \langle u_j, v_i \rangle)^2$$

$$f(u_i, V) = \sum_{(i,j) \in I} (a_{i,j} - \langle u_j, v_i \rangle)^2$$

Convex when fixed one.

### Convex Optimization

Def.:  $\{(x, t) | x \in \text{dom } f, f(x) \leq t\}, f : \mathbb{R}^D \rightarrow \mathbb{R}$  is convex, if  $\text{dom } f$  is a convex set, and if  $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f$ , and  $\forall \alpha \in [0, 1]: f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y})$ . Convex  $\Leftrightarrow$  Hessian p.s.d  $\Leftrightarrow$  local=global

Positive semi-definite: all principal minors (same-indexed rows and columns)  $\geq 0$

Positive definite: leading principal minors  $> 0$

### Convex Relaxation

Replace non-convex rank constraints by convex norm constraints (superset). Then project optimum back (hopefully still optimal).

$$\min_{\mathbf{B} \in P_k} \|\mathbf{A} - \mathbf{B}\|_G^2, P_k = \{\mathbf{B} : \|\mathbf{B}\|_* \leq k\} \supseteq Q_k =$$

$\{\mathbf{B} : \text{rank}(\mathbf{B}) \leq k\}$  (in fact tightest convex lower-bound  $\text{rank}(\mathbf{B}) \geq \|\mathbf{B}\|_*, \text{for } \|\mathbf{B}\|_2 \leq 1$ )

### SVD Thresholding

$\mathbf{B}^* = \text{shrink}_\tau(\mathbf{A}) = \arg \min_{\mathbf{B}} \{\|\mathbf{A} - \mathbf{B}\|_F^2 + \tau \|\mathbf{B}\|_*\}$   
Then with SVD  $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^\top, \mathbf{D} = \text{diag}(\sigma_i)$ , holds  $\mathbf{B}^* = \mathbf{U} \mathbf{D}_\tau \mathbf{V}^\top, \mathbf{D}_\tau = \text{diag}(\max\{0, \sigma_i - \tau\})$   
Iteration:  $\mathbf{B}_{t+1} = \mathbf{B}_t + \eta_t \Pi(\mathbf{A} - \text{shrink}_\tau(\mathbf{B}_t))$

### 4 Non-Negative Matrix Factorization

$\mathbf{X} \in \mathbb{Z}_{\geq 0}^{N \times M}$ , NMF:  $\mathbf{X} \approx \mathbf{U}^\top \mathbf{V}, x_{ij} = \sum_z u_{zi} v_{zj}$   
 $\langle \mathbf{u}_i, \mathbf{v}_j \rangle$  Decompose object into features: topics, face parts, etc..  $\mathbf{u}$  weights on parts,  $\mathbf{v}$  parts (bases). More interpretable (PCA: holistic repre.).

### EM for MLE for pLSA (No global opt. guarantee)

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

**Conditional independence assumption (\*):**

$$p(w|d) = \sum_z p(w, z|d) = \sum_z p(w|d, z) p(z|d) = \sum_z p(w|z) p(z|d) \text{ or } p(w|d, z) = p(w|z)$$

### Symmetric parameterization:

$$p(w, d) = \sum_z p(z) p(w|z) p(d|z)$$

$$\text{Log-Likelihood: } L(\mathbf{U}, \mathbf{V}) = \sum_{i,j} x_{i,j} \log p(w_j | d_i)$$

$$= \sum_{(i,j) \in X} \log \sum_{z=1}^K p(w_j | z) p(z | d_i)$$

$$p(w_j | z) = v_{z,j}, p(z | d_i) = u_{z,i}, \sum_j v_{z,j} = \sum_z u_{z,i} = 1$$

E-Step (optimal q: posterior of  $z$  over  $(d_i, w_j)$ ):

$$q_{zij} = \frac{p(w_j | z) p(z | d_i)}{\sum_{k=1}^K p(w_j | k) p(k | d_i)} := \frac{v_{z,j} u_{z,i}}{\sum_{k=1}^K v_{k,j} u_{k,i}}, \sum_z q_{zij} = 1$$

M-Steps:

$$p(z | d_i) = \frac{\sum_j x_{ij} q_{zij}}{\sum_j x_{ij}}, p(w_j | z) = \frac{\sum_i x_{ij} q_{zij}}{\sum_{i,l} x_{il} q_{zil}}$$

Lower Bound of  $L(\mathbf{U}, \mathbf{V})$  Jensen ineq. :

$$\sum_{i,j \in X} \sum_{z=1}^K q_{zij} (\log(v_{z,j}) + \log(u_{z,i}) - \log(q_{zij}))$$

### Latent Dirichlet Allocation

To sample a new document, we need to extend  $X$  and  $U^T$  with a new row, s.t.  $X = U^T V$ . (While pLSA fixes both dimensions)

For each  $d_i$  sample topic weights  $\mathbf{u}_i \sim \text{Dirichlet}(\alpha)$ :

$p(u_i | \alpha) = \prod_{z=1}^K u_{zi}^{\alpha_z - 1}$ , then topic  $z' \sim \text{Multi}(u_i)$ , word  $w' \sim \text{Multi}(v_{z'})$

Multinom. obsv. model on wc vec:  $p(\mathbf{x} | V, u) = \frac{1}{\prod_j x_j!} \prod_j \pi_j^{x_j}$  where  $\pi_j = \sum_z v_{z,j} u_{z,i}, l = \sum_j x_j$

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

### NMF Algorithm for quadratic cost function

$$\min_{\mathbf{U}, \mathbf{V}} J(\mathbf{U}, \mathbf{V}) = \frac{1}{2} \|\mathbf{X} - \mathbf{U}^\top \mathbf{V}\|_F^2 \text{ (non-negativity)}$$

s.t.  $\forall i, j, z : u_{zi}, v_{zj} \geq 0$

Comparison with pLSA:

- sampling model: Gaussian vs multinomial
  - objective: quadratic vs KL divergence
  - constraints: not normalized
- Alternating least squares:
- init:  $\mathbf{U}, \mathbf{V} = \text{rand}()$
  - repeat 3~4 for *maxIters*:
  - upd.  $(\mathbf{V} \mathbf{V}^\top) \mathbf{U} = \mathbf{V} \mathbf{X}^\top$ , proj.  $u_{zi} = \max\{0, u_{zi}\}$

4. update  $(\mathbf{U}\mathbf{U}^\top)\mathbf{V} = \mathbf{U}\mathbf{X}$ , proj.  $v_{zj} = \max\{0, v_{zj}\}$

## 5 Word Embeddings

**Distr. Model:**  $p_\theta(w|w') = \Pr[w \text{ in context of } w']$

**Log-likelihood:**

$$L(\theta; \mathbf{w}) = \sum_{t=1}^T \sum_{\Delta \in I} \log p_\theta(w^{(t+\Delta)} | w^{(t)})$$

**Latent Vector Model:**  $w \rightarrow (\mathbf{x}_w, b_w) \in \mathbb{R}^{D+1}$

$$p_\theta(w|w') = \frac{\exp\{\langle \mathbf{x}_w, \mathbf{x}_{w'} \rangle + b_w\}}{\sum_{v \in V} \exp\{\langle \mathbf{x}_v, \mathbf{x}_{w'} \rangle + b_v\}} \quad (\text{soft-max}).$$

**Modifications:**

$\log p_\theta(w|w') = \langle y_w, x_{w'} \rangle + b_w$ , word  $y_w$ , c'txt  $x_{w'}$   
use GloVe obj., negative sampling (logistic class.)

**GloVe (Weighted Square Loss)**

**Co-occ.**:  $\mathbf{N} = (n_{ij}) \in \mathbb{R}^{|V| \times |C|} = \#w_i \text{ in context } w_j$

**Objective:**  $H(\theta; \mathbf{N}) = \sum_{n_{ij} > 0} f(n_{ij})(\log n_{ij} - \log \exp[\langle \mathbf{x}_i, \mathbf{y}_j \rangle + b_i + d_j])^2$ ,  $f(n) = \min\{1, (\frac{n}{n_{\max}})^\alpha\}$ ,  $\alpha \in (0; 1]$  ( $= 3/4$ ) unnorm. distr.  
 $\rightarrow$  2-sided loss. cutoff  $n_{\max}$ : limit influence of high

freq.  $f(n) \xrightarrow{n \rightarrow 0} 0$ : as small counts very noisy

1. sample  $(i, j) u.a.r. s.t. n_{ij} > 0$

2.  $\mathbf{x}_i^{\text{new}} \leftarrow \mathbf{x}_i + 2\eta f(n_{ij})(\log n_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle) \mathbf{y}_j$

3.  $\mathbf{y}_j^{\text{new}} \leftarrow \mathbf{y}_j + 2\eta f(n_{ij})(\log n_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle) \mathbf{x}_i$   
embeds can model analogies and relatedness, but  
antonyms are usually not well captured.

## 6 Data Clustering & Mixture Models

**K-means Target:**  $\min_{\mathbf{U}, \mathbf{Z}} J(\mathbf{U}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{UZ}\|_F^2$   
 $= \sum_{n=1}^N \sum_{k=1}^K \|\mathbf{z}_{k,n} \mathbf{x}_n - \mathbf{u}_k\|_2^2$

1. **Initiate:** choose  $K$  centroids  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_K]$

2. **Cluster Assign:** data points to clusters.  $k^*(\mathbf{x}_n) = \arg \min_k \{\|\mathbf{x}_n - \mathbf{u}_k\|_2\}$  returns cluster  $k^*$ , whose centroid  $\mathbf{u}_{k^*}$  is closest to data point  $\mathbf{x}_n$ . Set  $\mathbf{z}_{k^*,n} = 1$ , and for  $l \neq k^*$   $\mathbf{z}_{l,n} = 0$ .

3. **Update centroids:**  $\mathbf{u}_k = \frac{\sum_{n=1}^N \mathbf{z}_{k,n} \mathbf{x}_n}{\sum_{n=1}^N \mathbf{z}_{k,n}}$ .

4. Repeat until  $\|\mathbf{Z} - \mathbf{Z}^{\text{new}}\|_0 = \|\mathbf{Z} - \mathbf{Z}^{\text{new}}\|_F^2 = 0$ .

Computational cost:  $O(k \cdot n \cdot d)$  Prior:  $p(z) = 1/K$   
**K-Means++:** 1. Choose centroid  $\mathbf{u}_1$  randomly from datapoints  $S$  2. For  $x \in S$ , calculate min. squared distance  $d_m(x)$  to existing centroids  $c_1, \dots, c_m$  3. Add new centroid  $c_{m+1}$ , chosen randomly from  $S$  with prob.  $p(x) = d_m(x) / \sum_{z \in S} d_m(z)$  4. Repeat until  $K$  centroids chosen  $\rightarrow$  proceed with K-means

**Gaussian Mixture Models (GMM)**

Gaussian  $p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$  Multivariate  
 $p(\mathbf{x}; \mu; \Sigma) = \frac{1}{|\Sigma|^{\frac{1}{2}} (2\pi)^{\frac{D}{2}}} \exp[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)]$

For GMM let  $\theta_k = (\mu_k, \Sigma_k)$ ;  $p_{\theta_k}(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)$

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

**Assignment variable (generative model):**

$z_{ij} \in \{0, 1\}$ ,  $\sum_{j=1}^k z_{ij} = 1$

Prior:  $\Pr(z_k = 1) = \pi_k \Leftrightarrow p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$

**Complete data distribution:**

$p_\theta(\mathbf{x}, \mathbf{z}) = \prod_{k=1}^K (\pi_k p_{\theta_k}(\mathbf{x}))^{z_k}$

**Posterior Probabilities:**

$$\Pr(z_k = 1 | \mathbf{x}) = \frac{\Pr(z_k=1)p(\mathbf{x}|z_k=1)}{\sum_{l=1}^K \Pr(z_l=1)p(\mathbf{x}|z_l=1)} = \frac{\pi_k p_{\theta_k}(\mathbf{x})}{\sum_{l=1}^K \pi_l p_{\theta_l}(\mathbf{x})}$$

$$\text{posterior } P(A|B) = \frac{\text{prior } P(A) \times \text{likelihood } P(B|A)}{\text{evidence } P(B)}$$

**Likelihood of observed data  $\mathbf{X}$ :**

$$p_\theta(\mathbf{X}) = \prod_{n=1}^N p_\theta(\mathbf{x}_n) = \prod_{n=1}^N (\sum_{k=1}^K \pi_k p_{\theta_k}(\mathbf{x}_n))$$

**Max. Likelihood Estimation (MLE):**

$$\arg \max_\theta \sum_{n=1}^N \log (\sum_{k=1}^K \pi_k p_{\theta_k}(\mathbf{x}_n))$$

$$\geq \sum_{n=1}^N \sum_{k=1}^K q_k [\log p_{\theta_k}(\mathbf{x}_n) + \log \pi_k - \log q_k]$$

with  $\sum_{k=1}^K q_k = 1$  by Jensen Inequality.

**Generative Model**

1. sample cluster index  $j \sim \text{Categorical}(\pi)$

2. given  $j$ , sample data  $x \sim \text{Normal}(\mu_j, \Sigma_j)$

**Expectation-Maximization (EM) for GMM**

$$\text{E-Step: } \Pr[z_{k,n}] = 1[\mathbf{x}_n] = q_{k,n} =$$

$$\frac{\pi_k^{(t-1)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t-1)}, \Sigma_k^{(t-1)})}{\sum_{j=1}^K \pi_j^{(t-1)} \mathcal{N}(\mathbf{x}_n | \mu_j^{(t-1)}, \Sigma_j^{(t-1)})}$$

$$\text{M-Step: } \mu_k^{(t)} := \frac{\sum_{n=1}^N q_{k,n} \mathbf{x}_n}{\sum_{n=1}^N q_{k,n}}, \pi_k^{(t)} := \frac{1}{N} \sum_{n=1}^N q_{k,n}$$

$$\Sigma_k^{(t)} = \frac{\sum_{n=1}^N q_{k,n} (\mathbf{x}_n - \mu_k^{(t)}) (\mathbf{x}_n - \mu_k^{(t)})^\top}{\sum_{n=1}^N q_{k,n}}$$

**K-means vs. EM** hard vs soft; spherical clusters vs covariance matrix; fast and cheap vs slow and more iteration; K-means can be used as init. for EM. K-means as a special case of GMM with covariances  $\Sigma_j = \sigma^2 I$ . limit of  $\sigma \rightarrow 0$  is K-means (hard asgmts).

**Model Order Selection (AIC / BIC for GMM)**

Trade-off between data fit (i.e. likelihood  $p(\mathbf{X}|\theta)$ ) and complexity (i.e. # of free parameters  $\kappa(\cdot)$ ). For choosing  $K$ :  $\text{AIC}(\theta|\mathbf{X}) = -\log p_\theta(\mathbf{X}) + \kappa(\theta)$

$$\text{BIC}(\theta|\mathbf{X}) = -\log p_\theta(\mathbf{X}) + \frac{1}{2} \kappa(\theta) \log N$$

# of free params, fixed covariance matrix:  $\kappa(\theta) = K \cdot D + (K - 1)$  ( $K$ : # clusters,  $D$ : dim(data) = dim( $\mu_i$ ),  $K - 1$ :  $\pi$  of # 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.

## 7 Sparse Coding

**Orthogonal Basis**

Pros: fast inverse; preserves energy. For  $\mathbf{x}$  and orthog. mat.  $\mathbf{U}$  compute  $\mathbf{z} = \mathbf{U}^\top \mathbf{x}$ . Approx  $\hat{\mathbf{x}} = \mathbf{U}\hat{\mathbf{z}}$ ,  $\hat{z}_i = z_i$  if  $|z_i| > \varepsilon$  else 0. Reconstruction Error  $\|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \sum_{d \notin \mathcal{S}} \langle \mathbf{x}, \mathbf{u}_d \rangle^2$ . Choice of base depends on signal. Fourier: global support, good for sine like waves; wavelet: local support, poor for non-vanishing signal; PCA basis optimal for given  $\Sigma$ . Stripes & check patterns: hi-freq in Fourier. Fourier:  $O(D \cdot \log D)$ , Wavelet:  $O(D)$  or  $O(D \cdot \log D)$

**Haar Wavelets (form orthogonal basis)**

scaling fnc.  $\phi(x) = [1, 1, 1, 1]$ , mother  $W(x) = [1, 1, -1, -1]$ , dilated  $W(2x) = [1, -1, 0, 0]$ , translated  $W(2x - 1) = [0, 0, 1, -1]$  Must be normalized

## 5.5 Sparse Coding

$\mathbf{U} \in \mathbb{R}^{D \times L}$  for # atoms  $L > D = \text{dim}(\text{data})$ . 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  $\mathbf{B}$  orthog. matrix •  $m([\mathbf{B}, \mathbf{u}]) \geq \frac{1}{\sqrt{D}}$  if atom  $\mathbf{u}$  is added to orthog. basis  $\mathbf{B}$  (o.n.b. = orthonormal base)

**Matching Pursuit (MP)** approximation of  $\mathbf{x}$  onto  $\mathbf{U}$ , using  $K$  entries. Objective:  $\mathbf{z}^* \in \arg \min_{\mathbf{z}} \|\mathbf{x} - \mathbf{U}\mathbf{z}\|_2$ , s.t.  $\|\mathbf{z}\|_0 \leq K$  1. init:  $z \leftarrow 0$ ,  $r \leftarrow x$  2. while  $\|\mathbf{z}\|_0 < K$  do 3. select atom *index* with smallest angle  $i^* = \arg \max_i |\langle \mathbf{u}_i, \mathbf{r} \rangle|$  4. update coefficients:  $z_{i^*} \leftarrow z_{i^*} + \langle \mathbf{u}_{i^*}, \mathbf{r} \rangle$  5. update residual:  $\mathbf{r} \leftarrow \mathbf{r} - \langle \mathbf{u}_{i^*}, \mathbf{r} \rangle \mathbf{u}_{i^*}$ .

**Exact recovery** when:  $K < 1/2(1 + 1/m(\mathbf{U}))$

**Compressive Sensing:** Compress data while gathering: •  $\mathbf{x} \in \mathbb{R}^D$ ,  $K$ -sparse in o.n.b.  $\mathbf{U}$ .  $\mathbf{y} \in \mathbb{R}^M$  with  $y_i = \langle \mathbf{w}_i, \mathbf{x} \rangle$ :  $M$  lin. combinations of signal;  $\mathbf{y} = \mathbf{W}\mathbf{x} = \mathbf{W}\mathbf{U}\mathbf{z} = \Theta\mathbf{z}$ ,  $\Theta \in \mathbb{R}^{M \times D}$  • Reconstruct  $\mathbf{x} \in \mathbb{R}^D$  from  $\mathbf{y}$ ; find  $\mathbf{z}^* \in \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0$ , s.t.  $\mathbf{y} = \Theta\mathbf{z}$  (e.g. with MP, or convex it with 1-norm: can be eq.). Given  $\mathbf{z}$ , reconstruct  $\mathbf{x} = \mathbf{U}\mathbf{z}$

Any orthogonal  $\mathbf{U}$  sufficient if: •  $\mathbf{W} =$  Gaussian random projection, i.e.  $w_{ij} \sim \mathcal{N}(0, \frac{1}{D})$  •  $M \geq cK \log(\frac{D}{K})$ , where  $c$  is some constant

## 8 Dictionary Learning

Adapt dict. to signal characteristics. Obj:  $(\mathbf{U}^*, \mathbf{Z}^*) \in \arg \min_{\mathbf{U}, \mathbf{Z}} \|\mathbf{X} - \mathbf{U} \cdot \mathbf{Z}\|_F^2$  not jointly convex but convex in either. **Matrix Fact. by Iter Greedy Min.** 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 ( $\mathbf{z}_n^{t+1} \in \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0$  s.t.  $\|\mathbf{x}_n - \mathbf{U}^t \mathbf{z}\|_2 \leq \sigma \|\mathbf{x}_n\|_2$ ) 2. Dict update step:  $\mathbf{U}^{t+1} \in \arg \min_{\mathbf{U}} \|\mathbf{X} - \mathbf{U} \mathbf{Z}^{t+1}\|_F^2$ , subj to  $\forall l \in [L]: \|\mathbf{u}_l\|_2 = 1$ . (set  $\mathbf{U} = [\mathbf{u}_1^t \dots \mathbf{u}_l^t \dots \mathbf{u}_L^t]$ ,  $\min_{\mathbf{u}_l} \|\mathbf{X} - \mathbf{U} \mathbf{Z}^{t+1}\|_F^2 = \min_{\mathbf{u}_l} \|\mathbf{R}_l^t - \mathbf{u}_l (\mathbf{z}_l^{t+1})^\top\|_F^2$  with  $\mathbf{R}_l^t = \tilde{\mathbf{U}} \Sigma \tilde{\mathbf{V}}^\top$  by  $\mathbf{u}_l^t = \tilde{\mathbf{u}}_l$ )

## 9 Neural Networks

**Activation:** scalar, non-linear  $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$ ,  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

**Neurons:**  $F_\sigma(\mathbf{x}; \mathbf{w}) = \sigma(w_0 + \sum_{i=1}^M x_i w_i) = \sigma(\mathbf{w}^\top \mathbf{x})$

**Output:** linear regression  $\mathbf{y} = \mathbf{W}^L \mathbf{x}^{L-1}$ , binary (logistic)  $y_1 = \text{P}[Y = 1 | \mathbf{x}] = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}^{L-1})}$ , multiclass

(soft-max)  $y_k = \text{P}[Y = k | \mathbf{x}] = \frac{\exp(\mathbf{w}_k^T \mathbf{x}^{L-1})}{\sum_{m=1}^K \exp(\mathbf{w}_m^T \mathbf{x}^{L-1})}$ . **Loss**

$l(y, \hat{y})$ : squared loss  $\frac{1}{2}(y - \hat{y})^2$ , cross-entropy loss  $-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$   $0 \leq \hat{y} \leq 1, y \in \{0, 1\}$

or  $y \in [0, 1]$ . **layer-wise:**  $\mathbf{x}^l = \sigma'(\mathbf{W}^{(l)} \mathbf{x}^{(l-1)})$ .

## Backpropagation

Layer-to-layer Jacobian:  $\mathbf{x} =$  prev. layer activation,  $\mathbf{x}^+ =$  next layer activation. Jacobian matrix  $\mathbf{J} = J_{ij}$  of mapping  $\mathbf{x} \rightarrow \mathbf{x}^+$ ,  $\mathbf{x}_i^+ = \sigma(\mathbf{w}_i^\top \mathbf{x})$ ,

$J_{ij} = \frac{\partial x_i^+}{\partial x_j} = w_{ij} \cdot \sigma'(\mathbf{w}_i^\top \mathbf{x})$ . Across multiple layers:

$$\frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{x}^{(l-n)}} = \mathbf{J}^{(l)} \cdot \frac{\partial \mathbf{x}^{(l-1)}}{\partial \mathbf{x}^{(l-n)}} = \mathbf{J}^{(l)} \cdot \mathbf{J}^{(l-1)} \dots \mathbf{J}^{(l-n+1)} \quad \text{and then back prop. } \nabla_{\mathbf{x}^{(l)}} \ell = \nabla_{\mathbf{y}^{(l)}} \ell \cdot \mathbf{J}^{(L)} \dots \mathbf{J}^{(l+1)}$$

$$\text{Weights: } \frac{\partial l}{\partial w_{ij}^{(l)}} = \frac{\partial l}{\partial x_i^{(l)}} \frac{\partial x_i^{(l)}}{\partial w_{ij}^{(l)}}, \quad \frac{\partial x_i^l}{\partial w_{ij}^l} =$$

$\sigma'([\mathbf{w}_i^{(l)}]^T \mathbf{x}^{(l-1)}) \cdot x_j^{(l-1)}$  (sensitivity of downstream unit · activation of up-stream unit) <

**Gradient Descent (or Deepest Descent)**

**Gradient:**  $\nabla f(\mathbf{x}) := \left( \frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_D} \right)^\top$

$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \gamma \nabla f(\mathbf{x}^{(t)})$ , usually  $\gamma \approx \frac{1}{L}$

**SGD** Assume additive obj.  $f(x) = \frac{1}{N} \sum_{n=1}^N f_n(x)$

sample  $n \in u.a.r. \{1, \dots, N\}$ , then

$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \gamma \nabla f_n(\mathbf{x}^{(t)})$ , typically  $\gamma \approx \frac{1}{L}$ .

**Neural Networks for Images (CNN)**

$$F_{n,m}(\mathbf{x}; \mathbf{w}) = \sigma(b + \sum_{k=-2}^2 \sum_{l=-2}^2 w_{k,l} x_{n+k,m+l}).$$

## 10 Deep Unsupervised Learning

**AR:** Image  $p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$  **ELBO:**

$$\mathbb{E}_{x \sim P_{\mathbf{x}}} [\mathbb{E}_{z \sim Q} \log P_g(x|z) - D_{KL}(Q(z|x) \| P(z))]$$

$Q$  enc. posterior distr.,  $P(z)$  prior distr. on latent var  $z$ ,  $P_g$  likelihood of dec. generated  $\mathbf{x}$ . Jointly trained: enc. optimize regularizer term, sample  $\mathbf{z} \sim Q$ , feed to dec., produce  $\hat{x}$  to max. reconstruction quality. Both terms diff'able, can use SGD to train end-to-end. **Repairam. trick:** use variational distr.s s.t.  $q_\phi(\mathbf{z}; \mathbf{x}) = g_\phi(\zeta; \mathbf{x})$  with eg.  $\zeta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  Example:  $\zeta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\mathbf{z} = \mu + U\zeta$  then  $z \sim \mathcal{N}(\mu, \mathbf{U}^\top \mathbf{U})$