

1 Essentials

Derivatives

$$\begin{aligned}\frac{\partial}{\partial \mathbf{x}}(\mathbf{a}^\top \mathbf{x}) &= \mathbf{a} & \frac{\partial}{\partial \mathbf{x}}(\mathbf{A}\mathbf{x}) &= \mathbf{A}^\top & \frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^\top \mathbf{A}) &= \mathbf{A} \\ \frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^\top \mathbf{x}) &= 2\mathbf{x} & \frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^\top \mathbf{A}\mathbf{x}) &= (\mathbf{A} + \mathbf{A}^\top)\mathbf{x} \\ \frac{\partial}{\partial \mathbf{X}}(\mathbf{c}^\top \mathbf{X}\mathbf{b}) &= \mathbf{c}\mathbf{b}^\top & \frac{\partial}{\partial \mathbf{X}}(\mathbf{c}^\top \mathbf{X}^\top \mathbf{b}) &= \mathbf{b}\mathbf{c}^\top \\ \frac{\partial}{\partial \mathbf{X}}(\|\mathbf{X}\|_F^2) &= 2\mathbf{X} & \frac{\partial}{\partial \mathbf{X}}(\log(\det(\mathbf{X}))) &= (\mathbf{X}^\top)^{-1} \\ \frac{\partial}{\partial \mathbf{a}}((\mathbf{x} - \mathbf{a})^\top \mathbf{W}(\mathbf{x} - \mathbf{a})) &= -2\mathbf{W}(\mathbf{x} - \mathbf{a}) \\ \frac{\partial}{\partial \mathbf{X}}(\mathbf{a}^\top \mathbf{X}^{-1}\mathbf{b}) &= -(\mathbf{X}^\top)^{-1}\mathbf{a}\mathbf{b}^\top (\mathbf{X}^\top)^{-1} \\ \frac{\partial}{\partial \mathbf{X}}(\text{Tr}(\mathbf{A}\mathbf{X})) &= \frac{\partial}{\partial \mathbf{X}}(\text{Tr}(\mathbf{X}\mathbf{A})) = \mathbf{A}^\top \\ \frac{\partial}{\partial \mathbf{X}}(\text{Tr}(\mathbf{A}\mathbf{X}^\top)) &= \frac{\partial}{\partial \mathbf{X}}(\text{Tr}(\mathbf{X}^\top \mathbf{A})) = \mathbf{A} \\ \frac{\partial}{\partial \mathbf{X}}(\text{Tr}(\mathbf{X}^\top \mathbf{A}\mathbf{X})) &= (\mathbf{A} + \mathbf{A}^\top)\mathbf{X}\end{aligned}$$

Jacobian

$$J_F = \begin{bmatrix} \nabla_{\mathbf{x}}^\top F_1 \\ \vdots \\ \nabla_{\mathbf{x}}^\top F_m \end{bmatrix} \text{ where } F: \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ for } \mathbf{x} \in \mathbb{R}^n$$

LinAlg

- A is **psd** if $\mathbf{v}^\top \mathbf{A} \mathbf{v} \geq 0$, note: if $\mathbf{A} = \mathbf{B}^\top \mathbf{B}$ then A is psd, $\text{psd} \Leftrightarrow \lambda \geq 0$
- **orthogonal matrix**: $\mathbf{U}^\top \mathbf{U} = \mathbf{U}\mathbf{U}^\top = \mathbf{I}$
- 1. unit length columns/rows 2. orthogonal columns/rows 3. square matrix

properties: preserve norm $\|\mathbf{U}\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$, $\det(\mathbf{U}) = \pm 1$, full rank because invertible.

- $\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x}\|_2 \cdot \|\mathbf{y}\|_2 \cdot \cos(\theta)$
- **spectral theorem**: Matrix A is diagonalizable by orthogonal matrix iff A is symmetric.
- A is degenerate/non-invertible $\Leftrightarrow \det(\mathbf{A}) = 0$.
- **rank-nullity theorem**: $\dim(\text{kernel}(\mathbf{A})) + \dim(\text{range}(\mathbf{A})) = n$ where n is input dimension.
- **determinant**: $\det(\text{diag}(a_1, \dots, a_n)) = a_1 * \dots * a_n$

$$\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = ad - bc$$

- **cauchy-schwarz**: $\mathbf{x}^\top \mathbf{y} \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$
- **trace** of a matrix is equal to the sum of its eigenvalues. Trace operator is linear and cyclic.
- eigenvectors corresponding to distinct eigenvalues of a symmetric matrix are orthogonal to each other.
- a function is convex iff its Hessian ∇^2 is psd.
- f convex $\Leftrightarrow \forall \lambda \in [0, 1], x, y \in X. f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$
- **jensens inequality**: $E[f(X)] \geq f(E[X])$ for f convex (\leq for f concave).

Eigendecomposition

$\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ where $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_m)$ with eigenvalues $\lambda_1 \geq \dots \geq \lambda_m$, $\mathbf{U} = (\mathbf{u}_1 | \dots | \mathbf{u}_m)$ eigenvectors orthogonal. Only exists if Σ is symmetric.

2 Eigenvalue / -vectors

Eigenvalue Problem: $\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$

1. solve $\det(\mathbf{A} - \lambda \mathbf{I}) \stackrel{!}{=} 0$ resulting in $\{\lambda_i\}_i$
2. $\forall \lambda_i$: solve $(\mathbf{A} - \lambda_i \mathbf{I})\mathbf{x}_i = \mathbf{0}$, for \mathbf{x}_i .

Norms

$$\begin{aligned}\|\mathbf{A}\|_F &= \sqrt{\sum_{i=1}^m \sum_{j=1}^n \mathbf{a}_{ij}^2} = \sqrt{\text{trace}(\mathbf{A}^\top \mathbf{A})} = \\ &= \sqrt{\text{trace}(\mathbf{A}\mathbf{A}^\top)} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2} \\ \|\mathbf{A}\|_2 &= \sup\{\|\mathbf{A}\mathbf{x}\| : \|\mathbf{x}\| = 1\} = \sigma_1 \\ \|\mathbf{M}\|_* &= \sum_{i=1}^{\min(m,n)} \sigma_i\end{aligned}$$

Probability / Statistics

- $\text{Var}(X) = E[(X - \mu)^2]$ • sample variance: $\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$

2 Principle Component Analysis

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

1. Empirical Mean: $\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$.
2. Center Data: $\bar{\mathbf{X}} = \mathbf{X} - [\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}}] = \mathbf{X} - \mathbf{M}$.
3. Cov.: $\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^\top = \frac{1}{N} \bar{\mathbf{X}}\bar{\mathbf{X}}^\top$.
4. Eigenvalue Decomposition: $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$.
5. Select $K < D$, only keep $\mathbf{U}_K = [\mathbf{u}_1, \dots, \mathbf{u}_K]$.
6. Transform data onto new Basis: $\bar{\mathbf{Z}}_K = \mathbf{U}_K^\top \bar{\mathbf{X}}$.
7. Reconstruct to original Basis: $\tilde{\mathbf{X}} = \mathbf{U}_K \bar{\mathbf{Z}}_K$.
8. Reverse centering: $\tilde{\mathbf{X}} = \tilde{\mathbf{X}} + \mathbf{M}$.

For compression save $\mathbf{U}_K, \bar{\mathbf{Z}}_K, \bar{\mathbf{x}}$.

$\mathbf{U}_K \in \mathbb{R}^{D \times K}, \Sigma \in \mathbb{R}^{D \times D}, \bar{\mathbf{Z}}_K \in \mathbb{R}^{K \times N}, \bar{\mathbf{x}} \in \mathbb{R}^{D \times N}$

- Reconstruction error = sum of discarded eigenvalues.

- The transformed dataset $\bar{\mathbf{Z}}_K$ has diagonal covariance matrix.

Iterative View

Residuals $r_i = x_i - \tilde{x}_i = (\mathbf{I} - \mathbf{U}\mathbf{U}^\top)\mathbf{x}_i$

Cov. matrix of residuals: $\frac{1}{n} \sum_{i=1}^n r_i r_i^\top = \frac{1}{n} \sum_{i=1}^n (\mathbf{I} - \mathbf{U}\mathbf{U}^\top) \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{I} - \mathbf{U}\mathbf{U}^\top)^\top = \dots = \Sigma - \lambda \mathbf{U}\mathbf{U}^\top$

1. Find principal eigenvector of $(\Sigma - \lambda \mathbf{U}\mathbf{U}^\top)$
2. which is the 2nd principal eigenvector of Σ
3. iterating to get d principal eigenvector of Σ
- the principal eigenvector of Σ points in the direction with largest variance of the data ($\arg\max_{\|\mathbf{u}\|=1} \mathbf{u}^\top \Sigma \mathbf{u}$).

Power iteration

find principal eigenvector of A: $\mathbf{v}_{t+1} = \frac{\mathbf{A}\mathbf{v}_t}{\|\mathbf{A}\mathbf{v}_t\|}$, $\lim_{t \rightarrow \infty} \mathbf{v}_t = \mathbf{u}_1$

3 SVD

$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^\top = \sum_{i=1}^{s=\min\{m,n\}} \sigma_i \mathbf{u}_i (\mathbf{v}_i)^\top$
 $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{U} \in \mathbb{R}^{m \times m}, \mathbf{D} \in \mathbb{R}^{m \times n}, \mathbf{V} \in \mathbb{R}^{n \times n}$
 \mathbf{U}, \mathbf{V} orthogonal, $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_s), \sigma_1 \geq \dots \geq \sigma_s \geq 0$. Rank(A) = #non-zero singular values.

Eckart-Young Theorem

$\min_{\text{rank}(\mathbf{B})=K} \|\mathbf{A} - \mathbf{B}\|_F^2 = \|\mathbf{A} - \mathbf{A}_K\|_F^2 = \sum_{r=K+1}^{\text{rank}(\mathbf{A})} \sigma_r^2$
note: use $\text{Tr}(\mathbf{A}^\top \mathbf{A})$ def of F-norm to prove.
 $\min_{\text{rank}(\mathbf{B})=K} \|\mathbf{A} - \mathbf{B}\|_2 = \|\mathbf{A} - \mathbf{A}_K\|_2 = \sigma_{K+1}$

4 Matrix Approximation & Reconstruction

problem of SVD is unobserved entries. Want to only consider I = observed entries.

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

Convex Relaxation

• $\text{rank}(\mathbf{B}) \geq \|\mathbf{B}\|_*$ for $\|\mathbf{B}\|_2 \leq 1$. proof: $\text{rank}(\mathbf{B}) = \#\{\sigma_i > 0\} = \sum_{i:\sigma_i > 0} 1 \geq \sum_{i:\sigma_i > 0} \sigma_i = \|\mathbf{B}\|_*$ since $\sigma_1 = \|\mathbf{B}\|_2 \leq 1$. Thus $Q_k = \{\mathbf{B} : \text{rank}(\mathbf{B}) \leq k\} \subseteq P_k = \{\mathbf{B} : \|\mathbf{B}\|_* \leq k\}$ a convex relaxation (in fact convex hull). • Singular Value Thresholding: if $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ then $\text{shrink}_\tau(\mathbf{A}) = \mathbf{U}\mathbf{D}_\tau \mathbf{V}^\top$ where $\mathbf{D}_\tau = \text{diag}(\max\{0, \sigma_i - \tau\})$. $\mathbf{B}_{t+1} = \mathbf{B}_t + \eta_t \pi(\mathbf{A} - \text{shrink}_\tau(\mathbf{B}_t))$ where π zeros out unobserved entries.

Alternating Least Squares

reparametrize $\mathbf{B} = \mathbf{U}\mathbf{V}$ where $\mathbf{U} \in \mathbb{R}^{m \times k}, \mathbf{V} \in \mathbb{R}^{k \times n}$. If we have product of two matrices, rank cannot be bigger than smallest rank occurring in it. $\min [\sum_{(i,j) \in I} (a_{ij} - u_i^\top v_j)^2]$. ALS = optimize over u_i 's while keeping v_j 's fixed and vice versa.

5 Non-Negative Matrix Factorization

pLSA

- co-occurrence matrix $X = x_{ij}$ #occurrences of word w_j in doc d_i . • $p(w|d) = \sum_z p(w, z|d) = \sum_z p(w|z)p(z|d) = \sum_z p(w|z)p(z|d)$.
- log-likelihood: $\sum_{i,j} x_{ij} \log p(w_j|d_i)$

E-Step (optimal q):

$$q_{zij} = p(z|w_j, d_i) = \frac{p(w_j|z)p(z|d_i)}{\sum_{k=1}^K p(w_j|k)p(k|d_i)}$$

M-Steps:

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

NMF Algorithm for quadratic cost function

$\mathbf{X} \in \mathbb{Z}_{\geq 0}^{N \times M}$, NMF: $\mathbf{X} \approx \mathbf{U}^\top \mathbf{V}, x_{ij}$

$$\min_{\mathbf{U}, \mathbf{V}} J(\mathbf{U}, \mathbf{V}) = \frac{1}{2} \|\mathbf{X} - \mathbf{U}^\top \mathbf{V}\|_F^2 = \frac{1}{2} \sum (x_{ij} - u_i^\top v_j)^2$$

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

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

• vector form:

$$(\sum_i u_i u_i^\top) v_j = \sum_i x_{ij} u_i, (\sum_j v_j v_j^\top) u_i = \sum_j x_{ij} v_j$$

note for derivation: If want matrix form, use trace def. If want vector only form, use sum representation of objective.

6 Word Embeddings

Skip-gram model:

$p_\theta(w|w') = \Pr[w \text{ occurs in context of } w']$

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

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

- **add context embeddings**: more flexibility $\log p_\theta(w|w') = \langle \mathbf{x}_{w'}, \mathbf{y}_w \rangle + b_w$, word embeddings \mathbf{y}_w , context embeddings $\mathbf{x}_{w'}$.

- **negative sampling (logistic classification)**: avoids having to compute normalization Z
 $\sum_{(i,j) \in \Delta^+} \log \sigma(\mathbf{x}_i^\top \mathbf{y}_j) + \sum_{(i,j) \in \Delta^-} \log \sigma(-\mathbf{x}_i^\top \mathbf{y}_j)$

GloVe

Co-occurrence Matrix:

$\mathbf{N} = (n_{ij}) \in \mathbb{R}^{|V| \times |C|} = \#$ of word w_i in context w_j

Objective: (Weighted Square Loss) $H(\theta; \mathbf{N})$

$$= \sum_{n_{ij} > 0} f(n_{ij}) (\log n_{ij} - \log \bar{p}_\theta(w_i | w_j))^2$$

with $\bar{p}_\theta(w_i | w_j) = \exp[\langle \mathbf{x}_i, \mathbf{y}_j \rangle + b_i + c_j]$ unnormalized! and $f(n) = \min\{1, (\frac{n}{n_{\max}})^\alpha\}$, $\alpha \in (0; 1]$.

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

7 K-means Algorithm

Hard assignments $Z \in \{0, 1\}^{N \times K}$, Centroids $\mathbf{U} = (\mathbf{u}_1 | \dots | \mathbf{u}_K) \in \mathbb{R}^{D \times K}$, data $\mathbf{X} \in \mathbb{R}^{D \times N}$

Target: $\min_{\mathbf{U}, \mathbf{Z}} J(\mathbf{U}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{U}\mathbf{Z}^\top\|_F^2$

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

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

3. Repeat until Z doesn't change anymore. Computational cost: $O(K \cdot N \cdot D)$

8 Gaussian Mixture Models (GMM)

let $\theta_k = (\boldsymbol{\mu}_k, \Sigma_k)$, $p(\mathbf{x}; \theta_k) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k)$

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

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

2. given j , sample data $\mathbf{x} \sim \text{Normal}(\boldsymbol{\mu}_j, \Sigma_j)$

Latent variables: data point \mathbf{x}_i belongs to cluster z_i . $p(z_i = j) = \pi_j$

Max. Likelihood Estimation (MLE):

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

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

Expectation-Maximization (EM) for GMM

E-Step: (posterior over latent variables)

$$q_{k,n}^* = \Pr(z_n = k | \mathbf{x}_n) = \frac{1}{Z} p(z_n = k) p(\mathbf{x}_n | z_n = k) =$$

$$\frac{\pi_k p(x_{n_i}; \theta_k)}{\sum_{i=1}^N \pi_i p(x_{n_i}; \theta_i)}$$

M-Step: $\mu_k^* := \frac{\sum_{n=1}^N q_{k,n} x_n}{\sum_{n=1}^N q_{k,n}}$, $\pi_k^* := \frac{1}{N} \sum_{n=1}^N q_{k,n}$

$$\Sigma_k^* = \frac{\sum_{n=1}^N q_{k,n} (x_n - \mu_k^*)(x_n - \mu_k^*)^T}{\sum_{n=1}^N q_{k,n}}$$

Gaussian distribution

Standard deviation σ , Mean μ

$$p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

Covariance matrix Σ , Mean μ

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

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 :

Akaike Information Criterion: $AIC(\theta|\mathbf{X}) = -\log p_\theta(\mathbf{X}) + \kappa(\theta)$

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

of free params for GMM: $\kappa(\theta) = KD + K \frac{D(D+1)}{2} + (K-1)$.

9 Sparse Coding

Orthogonal Basis

For x and orthogonal $U = (u_1 | \dots | u_D) \in \mathbb{R}^{D \times D}$. Encoding $\mathbf{z} = \mathbf{U}^T \mathbf{x}$. Decoding $\mathbf{x} = \mathbf{U}\mathbf{z} = \mathbf{U}\mathbf{U}^T \mathbf{x} = \mathbf{x}$. Thresholding $\hat{\mathbf{x}} = \mathbf{U}\hat{\mathbf{z}}$, $\hat{z}_i = z_i$ if $|z_i| > \epsilon$ else 0. Reconstr. Error $\|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \|\sum_{d \notin \sigma} \langle \mathbf{x}, \mathbf{u}_d \rangle \mathbf{u}_d\|^2 = \sum_{d \notin \sigma} \langle \mathbf{x}, \mathbf{u}_d \rangle^2$. Proof: $\|\sum x_i\|^2 = \sum \|x_i\|^2$ if x_i orthogonal. Advantages: efficient inverse, energy/length preservation.

Haar wavelets:

scaling: $[1, 1, 1, 1]$, mother wavelet: $[1, 1, -1, -1]$, dilated: $[1, -1, 0, 0]$, translated: $[0, 0, 1, -1]$.

Fourier vs Wavelets:

• Fourier: good for periodic signals, global support, no time info/only frequencies themselves.

• Wavelets: good for localized signals like abrupt changes/irregularities, represents a signal in time and frequency domain.

Overcomplete Basis

$\mathbf{U} \in \mathbb{R}^{D \times L}$ for # atoms $= L > D = \dim(\text{data})$. $\mathbf{x} \in \mathbb{R}^D$, $\mathbf{z} \in \mathbb{R}^L$ because $\mathbf{U}\mathbf{z} = \mathbf{x}$. Encoding ill-posed problem \rightarrow add constraint $\mathbf{z}^* \in \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0$ s.t. $\mathbf{x} = \mathbf{U}\mathbf{z}$. NP-hard and non-convex \rightarrow approx with Basis pursuit i.e relax to 1-norm (convex) or with Matching Pursuit.

Coherence

$$m(\mathbf{U}) = \max_{i,j:i \neq j} |\mathbf{u}_i^T \mathbf{u}_j|$$

Matching Pursuit (MP) a greedy approximating algorithm. Objective: $\mathbf{z}^* \in \arg \min_{\mathbf{z}} \|\mathbf{x} - \mathbf{U}\mathbf{z}\|_2$,

s.t. $\|\mathbf{z}\|_0 \leq K$. **init:** $z \leftarrow 0, r \leftarrow x$. **2.** while $\|\mathbf{z}\|_0 < K$ do **3.** select atom with max correlation $d^* = \arg \max_d |\langle \mathbf{u}_d, \mathbf{r} \rangle|$ **4.** update coefficients: $z_{d^*} \leftarrow z_{d^*} + \langle \mathbf{u}_{d^*}, \mathbf{r} \rangle$ **5.** update residual: $\mathbf{r} \leftarrow \mathbf{r} - \langle \mathbf{u}_{d^*}, \mathbf{r} \rangle \mathbf{u}_{d^*}$.

Compressive Sensing: Compress data while gathering: • $\mathbf{x} \in \mathbb{R}^D$, K -sparse in o.n.b. **U.** $\mathbf{y} \in \mathbb{R}^M$ corresponds to M lin. combinations/measurements of signal; $\mathbf{y} = \mathbf{W}\mathbf{x} = \mathbf{W}\mathbf{U}\mathbf{z}$ $\mathbf{z}, \theta \in \mathbb{R}^{M \times D}$ • 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). Given \mathbf{z} , reconstruct \mathbf{x} via $\mathbf{x} = \mathbf{U}\mathbf{z}$. Sufficient conditions: • \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

Dictionary Learning

Adapt the dictionary to signal characteristics.

Objective: $(\mathbf{U}^*, \mathbf{Z}^*) \in \arg \min_{\mathbf{U}, \mathbf{Z}} \|\mathbf{X} - \mathbf{U} \cdot \mathbf{Z}\|_F^2$

K-SVD (Iter Greedy Minimization): **1.** Coding step: $\mathbf{Z}^{t+1} \in \arg \min_{\mathbf{Z}} \|\mathbf{X} - \mathbf{U}^t \mathbf{Z}\|_F^2 = \sum_{i=1}^N \|x_i - \mathbf{U}^t z_i\|^2$ s.t. $\|z_i\|_0 \leq K$. Use any pursuit algorithm. **2.** Dict update step: $\mathbf{U}^{t+1} \in \arg \min_{\mathbf{U}} \|\mathbf{X} - \mathbf{U}\mathbf{Z}^{t+1}\|_F^2$, s.t. $\forall l \in [L] : \|\mathbf{u}_l\|_2 = 1$. idea: update one atom u_l at a time. $\min_{u_l} \|\mathbf{X} - \mathbf{U}\mathbf{Z}\|_F^2 = \min_{u_l} \|\mathbf{X} - (\sum_{e \neq l} u_e z_e^T + u_l z_l^T)\|_F^2 = \min_{u_l} \|\mathbf{R}_l - \mathbf{u}_l (\mathbf{z}_l)^T\|_F^2$ where \mathbf{z}_l is the l -th row of matrix \mathbf{Z} and \mathbf{R}_l is the residual due to atom u_l . Use SVD $\mathbf{R}_l = \sum_i \sigma_i \tilde{u}_i \tilde{v}_i^T$ then $\mathbf{u}_l^* = \tilde{\mathbf{u}}_1$ and $\mathbf{z}_l^* = \tilde{\mathbf{v}}_1$ (use power iteration for efficiency).

10 Neural Networks

$F^\sigma : \mathbb{R}^n \rightarrow \mathbb{R}^m, F_j^\sigma(x) = \sigma(w_j^T x)$ for $j = 1, \dots, m$

Activation fuctions: logistic function $\sigma(x) = \frac{1}{1+e^{-x}}$, $\sigma'(x) = \sigma(x)(1 - \sigma(x))$, ReLu $\phi(z) = \max(0, z)$

Output layer: linear regression: $\hat{\mathbf{y}} = \mathbf{W}^L \mathbf{x}^{L-1}$ binary classification (logistic):

$$\hat{y}_1 = P[Y = 1|\mathbf{x}] = \frac{1}{1 + \exp[-\langle \mathbf{w}_1^L, \mathbf{x}^{L-1} \rangle]}$$

multiclass (soft-max):

$$\hat{y}_k = P[Y = k|\mathbf{x}] = \frac{\exp[\langle \mathbf{w}_k^L, \mathbf{x}^{L-1} \rangle]}{\sum_{m=1}^K \exp[\langle \mathbf{w}_m^L, \mathbf{x}^{L-1} \rangle]}.$$

Loss function squared loss: $\frac{1}{2} (y - \hat{y})^2$

cross-entropy loss: $-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$.

Regularization: add l_2 -regularizer to objective or add drop-out layers.

Units and Layers: layer-to-layer fwd. prop. notation: $\mathbf{x}^{(l)} = \sigma^{(l)}(\mathbf{W}^{(l)} \mathbf{x}^{(l-1)})$ where $y = x^{(L)}$ is the output activation vector.

Backpropagation

• Use SGD to optimize over weights: $\theta \leftarrow \theta - \eta \nabla_{\theta} l(y_i; y(x_i; \theta))$ for $t = \{1, \dots, T\}$ • We want to

know $\partial l / \partial w_{ij}^{(l)}$ i.e how does changing weights affect the loss. • three steps: **1.** how does output y affect loss **2.** how do activities of units affect each other resp. y . **3.** how do weights affect activities of units. • **1.** $\nabla_y l = \frac{\partial}{\partial y} l(y^*, y) = \dots$

2. $\frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{x}^{(l-n)}} = \mathbf{J}^{(l)} \cdot \frac{\partial \mathbf{x}^{(l-1)}}{\partial \mathbf{x}^{(l-n)}} = \mathbf{J}^{(l)} \cdot \mathbf{J}^{(l-1)} \dots \mathbf{J}^{(l-n+1)}$ where \mathbf{x} = prev. layer activation, \mathbf{x}^+ = next layer activation. Jacobian matrix $\mathbf{J} = J_{ij}$ of mapping

$\mathbf{x} \rightarrow \mathbf{x}^+$, $\mathbf{x}_i^+ = \sigma(\mathbf{w}_i^T \mathbf{x})$, $J_{ij} = \frac{\partial x_i^+}{\partial x_j} = w_{ij} \cdot \sigma'(\mathbf{w}_i^T \mathbf{x})$.

3. $\frac{\partial x_i^+}{\partial w_{ij}} = \sigma'(\mathbf{w}_i^T \mathbf{x}) x_j$

• Perform forward pass to compute activities for all units. Compute gradient of objective wrt output layer activities. Propagate gradient info back from output to inputs. Compute local gradients of activities wrt weights.

Convolutional Neural Networks

• **Convolution step:** primary purpose is to extract features from the input image. Parameter/weight sharing = a kernel is used on multiple locations of the image with the same weights. Sparse interactions = by making kernel smaller than input. Discrete convolution operator $s[i, j] = (I * K)[i, j] = \sum_m \sum_n I[m, n] K[i - m, j - n]$ where I is the image and K the kernel. Note that arguments are commutative.

• **Pooling step:** reduce dim of each feature map e.g by max, sum, average over a predefined spatial neighborhood. Why? Scale invariant representation of image, less params/less overfitting.

11 Deep Unsupervised Learning

Autoencoders

learn low-dim representation $\mathbf{z} \in \mathbb{R}^d$ for given data. Linear autoencoder with weights $\mathbf{C} \in \mathbb{R}^{d \times m}$ (encoder) and $\mathbf{D} \in \mathbb{R}^{m \times d}$ (decoder). objective $\min \frac{1}{2n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\mathbf{C}\mathbf{x}_i\|^2$. Frobenius norm optimal approx (in this case) via SVD $\mathbf{X}^T \mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T$ then $\mathbf{C}^* = \mathbf{U}_d^T$ and $\mathbf{D}^* = \mathbf{U}_d$ (first d columns of \mathbf{U}).

Variational Autoencoders (VAEs)

Put a gaussian prior on distribution of continuous latent vector \mathbf{z} : $p(\mathbf{z}_l) \sim \mathcal{N}(\mu_l, \Sigma_l)$ for $l \in \{1, \dots, L\}$. This allows to easily generate new data points.

log-likelihood: $\log p_\theta(x) = \log \int p_\theta(x|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} = E_q[\log p_\theta(x|\mathbf{z})] + E_q[\log \frac{p(\mathbf{z})}{q(\mathbf{z})}]$ where $-D_{KL}(q||p) = E_q[\log \frac{p(\mathbf{z})}{q(\mathbf{z})}]$. KL-divergence tells how much two distributions diverge. Optimal $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x})$.

NN approach: **1.** recognition/inference model:

learns variational distribution $q(\mathbf{z})$ i.e given \mathbf{x} it returns params of normal distribution (μ_l, Σ_l) for $l = 1, \dots, L$ from which we than can sample the \mathbf{z}_l 's. **2.** generative model: implements $p_\theta(x|\mathbf{z})$ and deterministically maps \mathbf{z} to \mathbf{x} (reconstruction).

Autoregressive Models

generate output one variable at a time based on chain rule $p(x_{1:m}) = \prod_{i=1}^m p(x_i|x_{1:t-1})$.

PixelCNN: $n \times n$ image with pixels x_1, \dots, x_{n^2} . Generate pixel x_i by conditioning on previously generated pixels x_1, \dots, x_{i-1} . Use a masking filter for implem.

RNN: observed sequence x_1, \dots, x_T and corresponding labels y_1, \dots, y_T . Use feedbackloop $h_t = f(h_{t-1}, x_t)$ with hidden state h_t . LSTM units to avoid vanishing gradient problem.

PixelRNN: use RNN for mapping x_1, \dots, x_{i-1} to x_i . Row LSTM: convolve along each row from top to bottom (triangular receptive field i.e loss of context). Diagonal BiLSTM: convolve along the diagonal (receptive field includes all previously generated pixels).

12 Notes

• Proof that eigenvectors of symmetric matrix are orthogonal: $\lambda \langle x, y \rangle = \langle \lambda x, y \rangle = \langle Ax, y \rangle = \langle x, A^T y \rangle = \langle x, Ay \rangle = \langle x, \mu y \rangle = \mu \langle x, y \rangle \Rightarrow (\lambda - \mu) \langle x, y \rangle = 0$ where $\lambda \neq \mu$.

• K-means vs GMM: k-means has hard assignments, cheaper to train (less params). GMM has soft assignments, more expressive because cluster is described by a MVN i.e shape is defined by an arbitrary covariance matrix and not restricted to spherical clusters. GMM is generative model i.e we can do outlier detection, generate data points, uncertainty estimation.

• SVD and PCA: if $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$. Columns of \mathbf{U} are eigenvectors of $\mathbf{A}\mathbf{A}^T$. Columns of \mathbf{V} are eigenvectors of $\mathbf{A}^T \mathbf{A}$. Eigenvalues of $\mathbf{A}^T \mathbf{A}$ and $\mathbf{A}\mathbf{A}^T$ are singular values² of \mathbf{A} .

• orthogonal Haar Basis for 4-dim signals:

$$\mathbf{U} = \frac{1}{2} \begin{bmatrix} 1 & 1 & \sqrt{2} & 0 \\ 1 & 1 & -\sqrt{2} & 0 \\ 1 & -1 & 0 & \sqrt{2} \\ 1 & -1 & 0 & -\sqrt{2} \end{bmatrix}$$