

Essentials

Norms

$$\|M\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n m_{i,j}^2} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2} = \sqrt{\text{trace}(A^T A)}$$

$$\|M\|_1 = \sum_{i,j} |m_{i,j}| \quad \|M\|_2 = \sigma_{\max}(M)$$

$$\|M\|_p = \max_{\mathbf{v} \neq 0} \frac{\|M\mathbf{v}\|_p}{\|\mathbf{v}\|_p} \quad \|M\|_* = \sum_{i=1}^{\min(m,n)} \sigma_i$$

(Nuclear Norm) $\|M\|_* = \sum_i \sigma_i$, σ_i : singular value of M .

rank

$$\text{rank}(XY) \leq \text{rank}(X)$$

$$\text{rank}(XY) = \text{rank}(X), \forall Y \in \mathbb{R}^{n \times n}, \text{rank}(Y) = n$$

$$\text{rank}(A) = \text{rank}(UDV^T) = \text{rank}(D)$$

Derivatives

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{b}^T \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{b}) = \mathbf{b} \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{x}) = 2\mathbf{x}$$

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A}^T + \mathbf{A}) \mathbf{x} \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{b}^T \mathbf{A} \mathbf{x}) = \mathbf{A}^T \mathbf{b}$$

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{c}^T \mathbf{X} \mathbf{b}) = \mathbf{c} \mathbf{b}^T \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{c}^T \mathbf{X}^T \mathbf{b}) = \mathbf{b} \mathbf{c}^T$$

$$\frac{\partial}{\partial \mathbf{x}} (\|\mathbf{x} - \mathbf{b}\|_2) = \frac{\mathbf{x} - \mathbf{b}}{\|\mathbf{x} - \mathbf{b}\|_2} \quad \frac{\partial}{\partial \mathbf{x}} (\|\mathbf{x}\|_2^2) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{x}) = 2\mathbf{x}$$

$$\frac{\partial}{\partial \mathbf{x}} (\|\mathbf{X}\|_F^2) = 2\mathbf{X}$$

$$\frac{\partial}{\partial \mathbf{x}} (\|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2^2) = 2(\mathbf{A}^T \mathbf{A} \mathbf{x} - \mathbf{A}^T \mathbf{b})$$

Eigenvalue / -vectors

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

- solve $\det(\mathbf{A} - \lambda \mathbf{I}) \stackrel{!}{=} 0$ resulting in $\{\lambda_i\}_i$
- $\forall \lambda_i$: solve $(\mathbf{A} - \lambda_i \mathbf{I}) \mathbf{x}_i = \mathbf{0}$, for \mathbf{x}_i .
- $\mathbf{A} \in \mathbb{R}^{N \times N}$ 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: $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T$ (\mathbf{Q} orthogonal).

Convexity

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y), \forall \lambda \in [0, 1]$$

Distribution

Dirichlet: $p(u_i|\alpha) = \prod_{z=1}^K u_{zi}^{\alpha_k - 1}$

Multinomial: $p(x|\pi) = \frac{n!}{\prod_j x_j!} \prod_j \pi_j^{x_j}$

Gaussian: $f(x) = \frac{1}{2\pi\sqrt{|\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$

KL-Divergence

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

$$x \sim \mathcal{N}(\mu, \sigma^2) \Rightarrow \int_{-\infty}^{\infty} p(x) \log p(x) dx = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2}$$

Lagrangian Multipliers

Minimize $f(\mathbf{x})$ s.t. $g_i(\mathbf{x}) \leq 0, i = 1, \dots, m$ (**inequality constr.**) and $h_i(\mathbf{x}) = \mathbf{a}_i^T \mathbf{x} - b_i = 0, i = 1, \dots, p$ (**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. $\lambda \geq 0$. Note: $\max_{\lambda, \nu} D(\lambda, \nu) \leq \min_{\mathbf{x}} f(\mathbf{x})$, equality

if dom f and f convex

2 Principle Component Analysis

Algorithm Implementation

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

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

Iterative View

Idea: optimal direction = principle eigenvector of the sample Residual r_i : $x_i - \tilde{x}_i = I - uu^T x_i$

Cov of r : $\frac{1}{n} \sum_{i=1}^n (I - uu^T) x_i x_i^T (I - uu^T)^T = (I - uu^T) \Sigma (I - uu^T)^T = \Sigma - \lambda uu^T$

- Find principal eigenvector of $(\Sigma - \lambda uu^T)$, which is the second eigenvector of Σ
- iterating to get d principal eigenvector of Σ

Power Method

Power iteration: $v_{t+1} = \frac{Av_t}{\|Av_t\|}, \lim_{t \rightarrow \infty} v_t = u_1$

Assuming $\langle u_1, v_0 \rangle \neq 0$ and $|\lambda_1| > |\lambda_j| (\forall j \geq 2)$

Reconstruction Error

$$\min_{\text{rank}(B)=K} \frac{1}{N} \|\mathbf{A} - \mathbf{B}\|_F^2 = \frac{1}{N} \|\mathbf{A} - \mathbf{A}_k\|_F^2 = \sum_{r=k+1}^{\text{rank}(\mathbf{A})} \lambda_r, \lambda_r \text{ is the eigenvalue of } \Sigma = \frac{1}{N} \mathbf{A} \mathbf{A}^T.$$

In SVD, there is no $\frac{1}{N}$

Interpret eigenvalues as the variance in the dimension specified by the corresponding eigenvector.

Others

- If $\mathbf{A} = \mathbf{B} \mathbf{B}^T$ then \mathbf{A} is semi-positive definite. Proof: $v^T \mathbf{A} v = v^T \mathbf{B} \mathbf{B}^T v = \|\mathbf{B}^T v\|_2^2 \geq 0$
- Spectral Theorem: Matrix \mathbf{A} is diagonalizable by an orthogonal matrix \Leftrightarrow it is symmetric
- Compared to general linear autoencoder, PCA is unique and interpretable.
- Compare power methods and SVD.

Power: good for small k . Robust and conceptually easy.

SVD: good for mid-sized problem. Leverage wealth of numerical techniques, e.g. QR decomposition.

3 Support Vector Discriminant

$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T = \sum_{k=1}^{\text{rank}(\mathbf{A})} d_{k,k} u_k (v_k)^T$

$\mathbf{U}^T \mathbf{U} = \mathbf{I} = \mathbf{V}^T \mathbf{V}$ (\mathbf{U}, \mathbf{V} orthonormal)

$\mathbf{U}_{:,i}$: eigenvectors of $\mathbf{A} \mathbf{A}^T$, $\mathbf{V}_{:,i}$: eigenvectors of $\mathbf{A}^T \mathbf{A}$, \mathbf{D}_{ii} : singular values.

Algorithm Implementation

- eigenvalues of $\mathbf{A}^T \mathbf{A}$, in descending order, are $D_{ii} = \sqrt{A_{ii}}$
- eigenvectors of $\mathbf{A}^T \mathbf{A} \rightarrow \mathbf{V}$.

$\mathbf{B} = \mathbf{A} \mathbf{V} \mathbf{D}^{-1}$. 4. normalize each column of \mathbf{U} and \mathbf{V} .

Low-Rank approximation

$$\mathbf{A} = \sum_{i=1}^k d_i u_i v_i^T = \mathbf{U}_k \mathbf{D}_k \mathbf{V}_k^T$$

Echart-Young Theorem

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

$$\min_{\text{rank}(B)=K} \|\mathbf{A} - \mathbf{B}\|_2 = \|\mathbf{A} - \mathbf{A}_k\|_2 = \sigma_{k+1}$$

4 Matrix Reconstruction

Two formalization

- (Low-rank) $\min_{X: \text{rank}(X) \leq k} \|\mathbf{A} - \mathbf{X}\|_{\mathcal{I}}^2$
Regularize $\Rightarrow L(\mathbf{U}, \mathbf{V}) = \|\mathbf{A} - \mathbf{U}^T \mathbf{V}\|_{\mathcal{I}}^2 + \lambda \|\mathbf{U}\|_F^2 + \lambda \|\mathbf{V}\|_F^2 = \sum_{(i,j) \in \mathcal{I}} (a_{ij} - u_i^T v_j)^2 + \lambda \sum_{i=1}^m \|u_i\|^2 + \lambda \sum_{j=1}^n \|v_j\|^2$
- non-convex w.r.t (\mathbf{U}, \mathbf{V}) but convex w.r.t \mathbf{U} and \mathbf{V} .
- (Exact matrix) $\min_X \text{rank}(X)$, s.t. $\|\mathbf{A} - \mathbf{X}\|_{\mathcal{I}} = 0$
- the rank function (objective) is not convex, it is not smooth
- the constraint is very stringent.

ALS

For the first formalization: low-rank.

$$u_i = (\sum_{j:(i,j) \in \mathcal{I}} v_j v_j^T + \lambda I_k)^{-1} \sum_{j:(i,j) \in \mathcal{I}} a_{ij} v_j$$

Interpret: given low-dimensional representations of the items, compute independently the best representation of each user

Computational complexity for u_i is $O(n_i k^2 + k^3)$, k features, n_i number of items evaluated by user i .

$$v_j = (\sum_{i:(i,j) \in \mathcal{I}} u_i u_i^T + \lambda I_k)^{-1} \sum_{i:(i,j) \in \mathcal{I}} a_{ij} u_i$$

Interpret: given low-dimensional representations of the users, compute independently the best representation of each item.

Nuclear Norm

(Exact form) $\min_B \|\mathbf{B}\|_*, \text{s.t. } \|\mathbf{A} - \mathbf{B}\|_G = 0$

(Approximate) $\min_B \|\mathbf{A} - \mathbf{B}\|_G^2, \text{s.t. } \|\mathbf{B}\|_* \leq r$

Property1: $\text{rank}(B) \geq \|\mathbf{B}\|_*, \forall \|\mathbf{B}\|_2 \leq 1$

Property2: convexity. Proof:
goal: $\|\lambda \mathbf{A} + (1 - \lambda) \mathbf{B}\|_* \leq \lambda \|\mathbf{A}\|_* + (1 - \lambda) \|\mathbf{B}\|_*$
key1: write $\lambda \mathbf{A} + (1 - \lambda) \mathbf{B} = \mathbf{U}_\lambda \mathbf{D}_\lambda \mathbf{V}_\lambda^T$
key2: prove $\text{trace}(\mathbf{U}_\lambda^T \mathbf{A} \mathbf{V}_\lambda) \leq \sum_j \sigma_j(\mathbf{A}) = \|\mathbf{A}\|_*$, using Cauchy-Schwartz inequality.

SVD shrinkage

Idea: SVD thresholding + projection.

Objective: $\mathbf{B}^* = \arg_B \min \{\|\mathbf{B}\|_* + \frac{1}{2\tau} \|\mathbf{B}\|_F^2\}, \text{s.t. } \Pi_{\mathcal{I}}(\mathbf{A} - \mathbf{B}) = x_{ij} \forall (i, j) \in \mathcal{I}$

Algorithm: $\mathbf{B}_{t+1} = \mathbf{B}_t + \eta \Pi_{\mathcal{I}}(\mathbf{A} - \text{shrink}_{\tau}(\mathbf{B}_t))$
 shrink_{τ} : all singular values are reduced by at least τ .

5 Non-Negative Matrix Factorization

pLSA

Conditional independence: $p(w|d, z) = p(w|z)$

Objective: $L(\mathbf{U}, \mathbf{V}) = \sum_{i,j} x_{ij} \log p(w_j|d_i) = \sum_{i,j} x_{ij} \log \sum_z p(w_j|z) p(z|d_i) = \sum_{i,j} x_{ij} \log \sum_z v_{zj} u_{zi}$
 $v_{jz} := p(w_j|z), u_{zi} := p(z|d_i), \sum_z v_{jz} = \sum_k u_{ki} = 1$

EM

Latent variable: q_{zij} probability of w_j in d_i generated by topic $z, \sum_z q_{zij} = 1$

Lower bound: $\log \sum_{z=1}^K q_{zij} \frac{u_{zi} v_{zj}}{q_{zij}} \geq \sum_{z=1}^K q_{zij} [\log u_{zi} + \log v_{zj} - \log q_{zij}]$

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

M-Step: $u_{zi} = p(z|d_i) = \frac{\sum_j x_{ij} q_{zij}}{\sum_z \sum_j x_{ij} q_{zij} = \sum_j x_{ij}}, v_{zj} = p(w_j|z) = \frac{\sum_i x_{ij} q_{zij}}{\sum_{i,l} x_{il} q_{zil}}$

Latent Dirichlet Allocation

Objective:

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

$$p(x|\mathbf{V}, \mathbf{u}) = \text{Multi}(x|\pi), \pi_j := \sum_z v_{zj} u_z$$

Generative model:

- for d_i , sample $u_i \sim \text{Dirichlet}(\alpha)$
- for word slot t : (1) sample topic $z^t \sim \text{Multi}(u_i)$; (2) sample word $w^t \sim \text{Multi}(v_{z^t})$

NMF Algorithm for quadratic cost function

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

Constraints: non-negativity, normalization.

Projected ALS 1. randomly init \mathbf{U}, \mathbf{V} 2. repeat for enough times:

- update $\mathbf{U} : (\mathbf{V} \mathbf{V}^T) \mathbf{U} = \mathbf{V} \mathbf{X}^T, \text{proj. } u_{zi} = \max\{0, u_{zi}\}$
- update $(\mathbf{U} \mathbf{U}^T) \mathbf{V} = \mathbf{U} \mathbf{X}, \text{proj. } v_{zj} = \max\{0, v_{zj}\}$

6 Word Embeddings

(Skip-gram model) Log-likelihood:

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

Latent vector model

log-bilinear: $\log p(w|w') = \langle \mathbf{x}_w, \mathbf{z}_{w'} \rangle + b_w + \text{const}$

Negative Sampling We change the objective function as $L(\theta) = \sum_{(i,j) \in \Delta^+} \log \sigma(\langle \mathbf{x}_i, \mathbf{z}_j \rangle) + \sum_{(i,j) \in \Delta^-} \log \sigma(-\langle \mathbf{x}_i, \mathbf{z}_j \rangle)$

distribution: ratio to the appearance frequency.
Number: given

GloVe (Weighted Square Loss)

Objective: $\min H(\theta; \mathbf{N})$

$$= \sum_{(i,j)} f(n_{ij}) (\log n_{ij} - \log \exp[\langle \mathbf{x}_i, \mathbf{y}_j \rangle + b_i + d_j])^2$$

weighting: $f(n) = \min\{1, (\frac{n}{n_{\max}})^{\alpha}\}, \alpha \in (0; 1]$.

SGD solution

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

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

SGD
GD: $\mathbf{w}_{k+1} = \mathbf{w}_k - \gamma_k \frac{1}{n} \sum_i \nabla f_i(\mathbf{w}_k)$
SGD: $\mathbf{w}_{k+1} = \mathbf{w}_k - \gamma_k \nabla f_{r_k}(\mathbf{w}_k)$, r_k random.

7 K-means
Objective $\min_{\mathbf{U}, \mathbf{Z}} J(\mathbf{U}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{UZ}^T\|_F^2$
Constraints $\sum_j z_{ij} = 1, z_{ij} \in \{0, 1\}$

Strategy alternating (EM) minimization
- $z_{ij} = 1, j = \arg \min_k \|x_i - u_k\|^2$ (closed centroid)
- $\mathbf{u}_j = \frac{\sum_{i=1}^N z_{ij} \mathbf{x}_i}{\sum_{i=1}^N z_{ij}}$

Computational cost of each iteration: $O(k \cdot n \cdot d)$
Quadratic convergence rate.
K-means++ 1. $\mathcal{U}_1 = \{x_I\}, I \sim \text{Uniform}[1 : N]$
2. $\mathcal{U}_{k+1} = \mathcal{U}_{k+1} \cup \{x_I\}, I \sim \text{Categorical}(\mathbf{p}), p_i = \frac{D_i^2}{\sum_{j=1}^N D_j^2}$, D_i^2 the closest distance of node i to all the centroids.
more expensive but theoretical $O(\log K)$ guarantee

Core set $p_i = \frac{1}{2N} + \frac{D_i^2}{2 \sum_{j=1}^N D_j^2}$, $D_i^2 = \|x_i - \frac{1}{N} \sum_j x_j\|^2$
selected sample weighed by $\frac{1}{mp_i}$
 ϵ -approximation guarantee: $m \propto \frac{dk \log k + \log 1/\delta}{\epsilon^2}$

8 Gaussian Mixture Models (GMM)
Mixture model
Complete data distribution:

$p_\theta(\mathbf{x}, \mathbf{z}) = \prod_{k=1}^K (\pi_k p_{\theta_k}(\mathbf{x}))^{z_k}$
- $p_\theta(\mathbf{x}) = \sum_{k=1}^K \pi_k p_{\theta_k}(\mathbf{x})$
- $\Pr(z_k = 1) = \pi_k \Leftrightarrow p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$
Generation: given π , generate x
1. sample cluster index $j \sim \text{Categorical}(\pi)$
2. sample x from the j -th component.

Inference: given x , infer z .
 $\Pr(z_k = 1 | \mathbf{x}) = \frac{\pi_k p_{\theta_k}(\mathbf{x})}{\sum_{l=1}^K \pi_l p_{\theta_l}(\mathbf{x})}$

Objective: find $\pi, \{\theta_k\}$
MLE estimate: $\hat{\pi}, \hat{\theta} = \arg \max_{\pi, \theta} \log p_\theta(\mathbf{X}) = \log \prod_{n=1}^N p_\theta(\mathbf{x}_n) = \sum_{n=1}^N \log (\sum_{k=1}^K \pi_k p_{\theta_k}(\mathbf{x}_n))$
 $\geq \sum_{n=1}^N \sum_{k=1}^K q_k [\log p_{\theta_k}(\mathbf{x}_n) + \log \pi_k - \log q_k]$
with $\sum_{k=1}^K q_k = 1$ by Jensen Inequality.
EM method Objective Lagrangian

Expectation (Two-side equal): $q_j^* = \frac{\pi_j p(\mathbf{x}; \theta_j)}{\sum_{l=1}^K \pi_l p(\mathbf{x}; \theta_l)}$
 $\pi_j^* := \frac{1}{N} \sum_{i=1}^N q_{ij}^*$
Maximization: Gaussian model: $\theta_j = (\mu_j, \Sigma_j)$
- $\mu_j^* := \frac{\sum_{i=1}^N q_{ij} \mathbf{x}_i}{\sum_{i=1}^N q_{ij}}$, $\Sigma_j^* = \frac{\sum_{i=1}^N q_{ij} (\mathbf{x}_i - \mu_j)(\mathbf{x}_i - \mu_j)^T}{\sum_{i=1}^N q_{ij}}$

To K-means
With covariances $\Sigma_j = \sigma^2 * I$, $\sigma \rightarrow 0$

Model Selection
AIC($\theta | \mathbf{X}$) = $-\log p_\theta(\mathbf{X}) + \kappa(\theta)$
BIC($\theta | \mathbf{X}$) = $-\log p_\theta(\mathbf{X}) + \frac{1}{2} \kappa(\theta) \log N$
BIC penalizes complexity more.
Example for GMM
- fixed covariance matrix: $\kappa(\theta) = K \cdot D + (K - 1)$
- full covariance matrices: $\kappa(\theta) = K \cdot (D + \frac{D(D+1)}{2}) + (K - 1)$

9 Convolutional Neural Networks
Neurons:
 $\mathbf{x} = \sigma(w_0 + \sum_{i=1}^M x_i w_i)$. $\mathbf{x}^l = \sigma^l(\mathbf{W}^{(l)} \mathbf{x}^{(l-1)})$,
 $\mathbf{y} = \sigma^{(L)}(\mathbf{W}^{(L)} \sigma^{(L-1)}(\dots(\sigma^{(1)}(\mathbf{W}^{(1)} \mathbf{x}) \dots))$
Activation: ReLu $\max\{x, 0\}$, Sigmoid: $s(x) = \frac{1}{1+e^{-x}}$, $s'(x) = s(x)(1 - s(x))$

Output:
Regression: Loss: MSE, $\hat{\mathbf{y}} = \mathbf{W}^L \mathbf{x}^{L-1}$
Classification: Loss: cross entropy:
 $-\sum_{i=1}^N \sum_{j=1}^K y_{ij} \log(\hat{p}_{ij})$
Logistic: $\hat{y}_1 = \Pr[Y = 1 | \mathbf{x}] = \frac{1}{1 + \exp[-(\mathbf{w}_1^T \mathbf{x}^{L-1})]}$

Soft-max: $\hat{y}_k = \Pr[Y = k | \mathbf{x}] = \frac{\exp[(\mathbf{w}_k^L \mathbf{x}^{L-1})]}{\sum_{m=1}^K \exp[(\mathbf{w}_m^L \mathbf{x}^{L-1})]}$
Objective $\mathcal{L}_\lambda(\theta; \mathcal{X}) = \frac{1}{N} \sum_{i=1}^N l(y_i, \hat{y}_i) + \frac{\lambda}{2} \|\theta\|_2^2$

Backpropagation
SGD $\theta \leftarrow (1 - \eta \lambda) \theta - \eta \nabla_\theta l(y_i^*; y(x_i; \theta))$
Jacobian $\mathbf{x}_i^+ = \sigma(\mathbf{w}_i^T \mathbf{x})$, $J_{ij} = \frac{\partial \mathbf{x}_i^+}{\partial \mathbf{x}_j} = w_{ij} \cdot \sigma'(\mathbf{w}_i^T \mathbf{x})$.
 $\frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{x}^{(l-n)}} = \mathbf{J}^{(l)} \cdot \mathbf{J}^{(l-1)} \dots \mathbf{J}^{(l-n+1)}$
 $\nabla_{\mathbf{x}^{(l)}}^T \ell = \nabla_{\mathbf{y}}^T \ell \cdot \mathbf{J}^{(L)} \dots \mathbf{J}^{(l+1)}$

$\frac{\partial l}{\partial w_{ij}^{(l)}} = \frac{\partial l}{\partial x_i^{(l)}} \frac{\partial x_i^{(l)}}{\partial w_{ij}^{(l)}} = \frac{\partial l}{\partial x_i^{(l)}} \sigma'([\mathbf{w}_i^{(l)}]^T \mathbf{x}^{(l-1)}) x_j^{(l-1)}$
CNN

Convolution layer $F_{n,m}(\mathbf{x}; \mathbf{w}) = \sigma(b + \sum_{k=-2}^2 \sum_{l=-2}^2 w_{k,l} x_{n+k, m+l})$. To reduce dimension of convolution, use {max, avg}-pooling

10 Generative models
VAE
Objective
 $\log p_\theta(x) = \mathbb{E}_{h \sim q_\phi} [\log p_\theta(x)] = \mathbb{E}_h [\log p_\theta(x|h)] - D[q_\phi(h|x) \| p_\theta(h)] + D[q_\phi(h|x) \| p_\theta(h|x)]$
= ELBO + $D[q_\phi(h|x) \| p_\theta(h|x)]$

ELBO $\log p(x; \theta) = \log \int p(x, h) dh \geq \int q(h|x) [\log p(x|h) + \log p(h) - \log q(h|x)] dh = \int q(h|x) \log p(x|h) dh - D[q(h|x) \| p(h)]$
Stochastic approximation: $\nabla_\theta \mathbb{E}_{q_\phi} [\mathcal{L}(\mathbf{x}, \mathbf{h})] \approx$

$\frac{1}{L} \sum_{r=1}^L \nabla_\theta \log p_\theta(\mathbf{x} | \mathbf{h}^{(r)})$, $\mathbf{h}^{(r)} \sim q_\phi(\cdot | \mathbf{x})$, i.i.d
In the log-likelihood context: $\mathbb{E}_{q_\phi} [\mathcal{L}(\mathbf{x}, \mathbf{h})] = B(q_\phi, \mathbf{x})$
Re-parametrization trick $q_\phi(h; x) = g_\phi(\zeta; x)$, $\zeta \sim \text{simple distribution}$
GAN
Objective: $\min_G \max_D \mathbb{E}_{x \sim P_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim P_z(z)} [\log(1 - D(G(z)))]$ or
 $l(\theta, \phi) := \mathbb{E}_{\tilde{p}_\theta} [y \ln q_\phi(\mathbf{x}) + (1 - y) \ln(1 - q_\phi(\mathbf{x}))]$
- θ : Generator, $\tilde{p}_\theta(\mathbf{x}, y = 1) = p(y = 1) \cdot p(\mathbf{x})$, $\tilde{p}_\theta(\mathbf{x}, y = 0) = p(y = 0) \cdot p_\theta(\mathbf{x})$
- ϕ : Discriminator. $q_\phi: \mathbf{x} \mapsto [0; 1]$

Saddle-point problem SGD as a heuristic solution (may diverge!)
 $\theta^{t+1} = \theta^t - \eta \nabla_\theta l(\theta^t, \phi^t)$, minimize: minus
 $\phi^{t+1} = \phi^t + \eta \nabla_\phi l(\theta^{t+1}, \phi^t)$, maximize: plus

Autoregressive model
 $p(x_1, x_2, \dots, x_m) = \prod_{t=1}^m p(x_t | x_{1:t-1})$

11 Sparse Coding
Orthogonal Basis
For \mathbf{x} and o.n.b. \mathbf{U} compute $\mathbf{z} = \mathbf{U}^T \mathbf{x}$. Approx $\hat{\mathbf{x}} = \mathbf{U} \hat{\mathbf{z}}$, $\hat{z}_i = z_i$ if $|z_i| > \epsilon$ else 0. Reconstruction Error $\|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \sum_{d \notin \sigma} \langle \mathbf{x}, \mathbf{u}_d \rangle^2$.

Wavelet Basis
Mother wavelet: $\psi(t) = 1, 0 \leq t \leq \frac{1}{2}; = 0, \frac{1}{2} \leq t \leq 1$.
 $\psi_{n,k}(t) = 2^{\frac{n}{2}} \psi(2^n t - k), 0 \leq k < 2^n$
Coherence • $m(\mathbf{U}) = \max_{i,j: i \neq j} |\mathbf{u}_i^T \mathbf{u}_j|$
• $m([\mathbf{B}, \mathbf{u}]) \geq \frac{1}{\sqrt{D}}$ if atom \mathbf{u} is added to orthogonal basis \mathbf{B}

Matching Pursuit (MP)
1. init: $\hat{\mathbf{x}}_0 \leftarrow \mathbf{x}$, $\mathbf{r}_0 \leftarrow \mathbf{x}$ 2. select the basis $j^* = \arg \max_j |\langle \mathbf{u}_j, \mathbf{r}_i \rangle|$ 3. update coefficients: $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \langle \mathbf{u}_{j^*}, \mathbf{r}_i \rangle \mathbf{u}_{j^*}$ 4. update residual: $\mathbf{r}_{i+1} \leftarrow \mathbf{r}_i - \langle \mathbf{u}_{j^*}, \mathbf{r}_i \rangle \mathbf{u}_{j^*}$.

Convergence $\frac{\|\mathbf{r}_{i+1}\|_2^2}{\|\mathbf{r}_i\|_2^2} = 1 - |\langle \frac{\mathbf{r}_i}{\|\mathbf{r}_i\|_2}, \mathbf{u}_{j^*} \rangle|^2$
 $\exists \mu_{min} \in (0, 1]$ (for overcomplete dictionary)
 $\|\mathbf{r}_i\|_2^2 \leq (1 - \mu_{min})^i \|\mathbf{r}_0\|_2^2$

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). 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

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

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 ($\mathbf{z}_n^{t+1} \in \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0$ s.t. $\|\mathbf{x}_n - \mathbf{U}^t \mathbf{z}\|_2 \leq \sigma \|\mathbf{x}_n\|_2$)
2. Dict update step: $\mathbf{U}^{t+1} \in \arg \min_{\mathbf{U}} \|\mathbf{X} - \mathbf{U} \mathbf{Z}^{t+1}\|_F^2$. $\forall l \in 1, 2, \dots, L$:
set $\mathbf{U} = [\mathbf{u}_1^t \dots \mathbf{u}_l^t \dots \mathbf{u}_L^t]$
 $\min_{\mathbf{u}_l} \|\mathbf{X} - \mathbf{U} \mathbf{Z}^{t+1}\|_F^2 = \|(\mathbf{X} - \sum_{e \neq l} \mathbf{u}_e^t (\mathbf{z}_e^{t+1})^T) - \mathbf{u}_l \mathbf{z}_l^{t+1}\|_F^2 = \|\mathbf{R}_l^t - \mathbf{u}_l (\mathbf{z}_l^{t+1})^T\|_F^2$
Doing 1-SVD on $\mathbf{R}_l^t = \sum_i \sigma_i \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T$
update $\mathbf{u}_l^{t+1} = \tilde{\mathbf{u}}_1$, $\mathbf{z}_l^{t+1} = \sigma_1 \tilde{\mathbf{v}}_1$