# Topics in ML:

what is
(1) big data!

let $\phi$ be an algorithm

$S_\phi(n)$ = storage required by $\phi$ to act on input of size $n$.

$R_\phi(n)$ = runtime —"—

$S_{max}(t)$ = max storage available at time $t$.

$\Delta t$ = doubling time

$S_{max}(t+\Delta t) = 2 S_{max}(t)$

$n_{max} = \max \{ n' : S_\phi(n') \leq S_{max} \}$

clearly $n_{max}(t+\Delta t) = 2 n_{max}(t)$

so $R_\phi(n_{max}(t+\Delta t)) = R_\phi(2 n_{max}(t))$

---

Suppose $R_\phi(n) = O(n^p)$    polynomial time!

Then $R_\phi(n_{max}(t+\Delta t)) = R_\phi(2 n_{max}(t))$

$= 2^p R_\phi(n_{max}(t))$

Now, since computer speed has doubled,
computational time increases by $2^p/2 = 2^{p-1}$.

Def$^n$  An algorithm is scalable if $p=1$

( or $R_\phi = O(n \log n)$ ).

why?

Answer: big data is about scalable algorithms

(2) **Matrix factoring**

$$X \simeq A \cdot B$$

$$(D \times N) \quad (D \times R) \cdot (R \times N)$$

→ storage cost: $X$ : $DN$

$$A, B : R(N+D)$$

→ computational cost: $X \cdot y$    $D(2N-1)$

$$A B y \quad R(2N-1) + D(2R-1)$$

**Applications**

(1) dimension reduction

(2) sparse coding

(3) collaborative filtering / recommender systems

(4) in painting

(5) compression

(6) NNs.

$D$ movies by $N$ users   $X$

$A \ B$   $D$ movies $\times$ $R$ latent features

$R$ latent features $\times$ $N$ users.

eg. recommender systems

| | alice | bob | edna | frank |
|---|---|---|---|---|
| star wars | 3 | | | |
| wall E | | 5 | | |
| transpotting | | 2 | | |
| metropolis | | | | 4 |
| citizen kane | | 0 | | 1 |

$$X =$$

# Principal Component Analysis ③

mean: $\bar{x} = \frac{1}{N}\sum_{n=1}^{N} x_n$

Covariance: $\Sigma = \frac{1}{N}\sum_{n=1}^{N}(x_n - \bar{x})(x_n - \bar{x})^T$

## projected data:

mean: $u_1^T \bar{x}$

Covariance: $\frac{1}{N}\sum_{n=1}^{N}(u_1^T x_n - u_1^T \bar{x})^2 = \frac{1}{N}\sum_{n=1}^{N}\left[u_1^T(x_n - \bar{x})\right]^2$

$$= \frac{1}{N}\sum_{n=1}^{N} u_1^T(x_n - \bar{x})(x_n - \bar{x})^T u_1$$

$$= u_1^T \Sigma u_1$$

Goal: maximize variance of projected data

ie "find most interesting direction in data"

---

$\max\limits_{u_1} \quad u_1^T \Sigma u_1 \quad$ s.t. $\quad \|u_1\|_2^2 = 1$

Lagrangian: $L = u_1^T \Sigma u_1 + \lambda_1(1 - u_1^T u_1)$

set $\partial/\partial u_1 = 0$, get $\Sigma u_1 = \lambda_1 u_1$

First principal direction: $\Sigma u_1 = \lambda_1 u_1$

$u_1$ is evector of $\Sigma$ w/ eval $\lambda_1$

$\lambda_1$ is variance of projected data

why?

# Minimum Error formulation

Define orthonormal basis $\{u_n\}$

proj$^n$ of $x_n$ onto $\underline{u}_d$ is $z_{n,d} = \underline{x}_n^T \underline{u}_d$

$$\underline{x}_n = \sum_{d=1}^{D} z_{n,d} \underline{u}_d = \sum_{d=1}^{D} (\underline{x}_n^T \underline{u}_d) \underline{u}_d$$

why?

$$\tilde{\underline{x}}_n = \sum_{d=1}^{R} a_{n,d} \underline{u}_d + \sum_{d=R+1}^{D} b_d \cdot \underline{u}_d$$

← doesn't depend on $n$.

## Approx$^n$ error:

$$J(\{a_{n,d}\}, \{b_d\}) = \frac{1}{N} \sum_{n=1}^{N} \| \underline{x}_n - \tilde{\underline{x}}_n \|_2^2$$

minimize $J$ wrt $a_{n,d}$ & $b_d$:

$$a_{n,d} = \underline{x}_n^T \underline{u}_d$$
$$b_d = \bar{\underline{x}}^T \underline{u}_d$$

---

Min wrt $u_d$:

$$\underline{x}_n - \tilde{\underline{x}}_n = \sum_{d=R+1}^{D} \left( (\underline{x}_n - \bar{\underline{x}})^T \underline{u}_d \right) \underline{u}_d$$

$\underbrace{}$ displacement vector is orthogonal to principal subspace.

Minim$^n$ wrt $\underline{y}_d$

$$J = \frac{1}{N} \sum_{n=1}^{N} \sum_{d=R+1}^{D} (\underline{x}_n^T \underline{u}_d - \bar{\underline{x}}^T \underline{u}_d)^2$$

$$= \sum_{d=R+1}^{D} \underline{u}_d^T \Sigma \underline{u}_d$$

minimize by picking evectors $\underline{u}$ w/ smallest evals.

why?

eg $D=2$, $K=1$

choose $\underline{u_2}$ to minimize $\underline{u_2}^T \Sigma \underline{u_2}$

subject to $\underline{u_2}^T \underline{u_2} = 1$

$\mathcal{L} = \underline{u_2}^T \Sigma \underline{u_2} + \lambda_2 (1 - \underline{u_2}^T \underline{u_2})$

$\hookrightarrow \Sigma \underline{u_2} = \lambda_2 \underline{u_2}$

$\rightarrow J = \lambda_2.$

Error is minimized by choosing smaller of the two evals.

---

Matrix viewpoint:

$X = [\underline{x_1}, \dots, x_N]$ $\qquad x_n = D-\text{dim}$ vec.

mean-centred data $\bar{X} = X - m$ $\qquad m = [\bar{x}, \dots, \bar{x}]$

---

Matrix factor$^n$ of $\bar{X}$.

$\bar{Z}_R = U_R^T \cdot \bar{X}$

$(R \times N) \quad (R \times D) \quad (D \times N)$

to approximate $\bar{X}$, return to original basis

$\tilde{\bar{X}} = U_R \cdot \bar{Z}_R$

for $R = D$, get $\tilde{\bar{X}} = \bar{X}$!

# Singular Value Decomposition

**Thm:** let $A \in \mathbb{R}^{m \times N}$

Then $A$ can be decomposed (not nec. uniquely) as

$$A = U \cdot D \cdot V^T$$



$m \times N$   $m \times m$   $m \times N$   $N \times N$

$$U^T U = I_m = U U^T$$

$$D \ \text{diagonal}$$

$$V^T V = I_N = V V^T$$

Elements along diagonal of $D$ are called __singular values__

Assume w.l.o.g. $m > N$.

$$\begin{cases} d_{i,i} > 0 & i \leq i \leq r \\ d_{i,i} = 0 & (r+1) \leq i \leq N \\ d_{i,j} = 0 & i \neq j \end{cases}$$

by convention order singular values from high to low.

$$d_1 \geq d_2 \geq \ \_\_\_ \ \geq d_r > d_{r+1} = \ \_\_\_ = d_N = 0.$$

## Second principal direction:

$$\max_{u_2} \quad u_2^T \Sigma u_2 \quad \text{s.t.} \quad \|u_2\|_2^2 = 1 \;\&\; \lambda \;\&\; u_2^T u_1 = 0$$

$$L = u_2^T \Sigma u_2 + \lambda(1 - u_2^T u_2) + \eta(u_2^T u_1)$$

$$\frac{\partial L}{\partial u_2} = 0 \Rightarrow 2\Sigma u_2 - 2\lambda u_2 + \eta u_1 = 0$$

orthogonality $\Rightarrow \eta = 0$ (Why?)

choose $u_2$ = evector w/ $2^{nd}$ largest eval.

More generally, decomposition $\Sigma = U \Lambda U^T$ contains all relevant info.

✓ (Why?)

---

**Spectral Theorem:** If $A$ is Hermitian then $\exists$ orthonormal basis of $V$ consisting of evecs of $A$. Evals are real.

proof: FTA applied to char poly$^n$ $\Rightarrow$ (Why?)
$\exists$ eval $\lambda_1$ & evec $\underline{e}_1$.

Then $\lambda_1 \langle \underline{e}_1, \underline{e}_1 \rangle = \langle A\underline{e}_1, \underline{e}_1 \rangle$

$\qquad = \langle \underline{e}_1, A\underline{e}_1 \rangle = \bar{\lambda}_1 \langle \underline{e}_1, \underline{e}_1 \rangle$
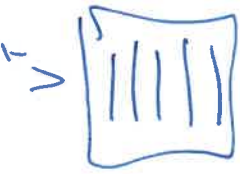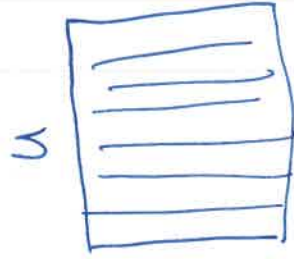
$\qquad$ so $\lambda_1 \in \mathbb{R}$.

Now consider $K = \text{span}(\underline{e}_1)^\perp$

$\qquad AK \subset K$ (Why?)

induction. ▨

## Left & right singular vectors

$$U \cdot [D] \cdot V^T$$

right singular vectors
orthogonal basis for
row space of A.

left-singular vectors
orthogonal basis for
space spanned by columns of A

## singular values & vectors

$$A v_i = d_i \, \underline{u_i}$$
$$A^T \underline{u_i} = d_i \, \underline{v_i}$$

$i = 1 \longrightarrow \min \{M, N\}$.

## Range, Null-space, Rank

* right singular vectors corresponding to vanishing singular
  values span null-space of A.

* left singular vectors corresponding to non-zero singular values
  (vectors)
  of A span range of A.

$$\Longrightarrow \; rk(A) = \# \text{ non-zero singular values}$$
$$= \# \{d_i > 0\}.$$

# Existence of SVD

Let $\sigma(\underline{u}, \underline{v}) = \underline{u}^T M \underline{v}$ where $\|\underline{u}\|_2 = 1 = \|\underline{v}\|_2$

Since $\sigma$ is continuous & $S^{M-1}$, $S^{N-1}$ are cpt, it follows that

a maximum is attained & it is $\geq 0$.

(why?!)

Denote max by $\sigma_1 \geq 0$.

Thm: $\underline{u}_1, \underline{v}_1$ are left & right-singular vectors corresponding to $\sigma_1$.

proof: $\nabla_{\underline{u},\underline{v}} \, \sigma = \nabla_{\underline{u},\underline{v}} \, \underline{u}^T M \underline{v} - \lambda_1 \nabla_{\underline{u},\underline{v}} \, \underline{u}^T \underline{u} - \lambda_2 \nabla_{\underline{u},\underline{v}} \, \underline{v}^T \underline{v} = 0.$

$\hookrightarrow \nabla_{\underline{u}} \, \underline{u}^T M \underline{v} = \lambda_1 \nabla_{\underline{u}} \, \underline{u}^T \underline{u} \quad \Rightarrow \quad M \underline{v} = 2\lambda_1 \underline{u}_1$

$\& \; \nabla_{\underline{v}} \, \underline{u}^T M \underline{v} = \lambda_2 \nabla_{\underline{v}} \, \underline{v}^T \underline{v} \quad \rightarrow \quad M^T \underline{u} = 2\lambda_2 \underline{v}$

$\overbrace{\qquad\qquad\qquad} \quad \rightarrow \quad \sigma_1 = 2\lambda_1 = 2\lambda_2.$

$M \underline{v}_1 = \sigma_1 \underline{u}_1, \; \& \; M^T \underline{u}_1 = \sigma_1 \underline{v}_1$

# Eigenvectors of $AA^T$ & $A^TA$

$A = UDV^T$

$\therefore$ $AA^T = UD^2U^T$

$A^TA = VD^2V^T$

columns of $U$ are evectors of $AA^T$

cols of $V$ ~~rows~~ $A^TA$

evals are $\sigma_i^2$ in both cases.

# Why SVD & not PCA?

suppose $D$ large.

cov. has dim $D^2$.

# Pseudo Inverse

Let $D$ be diagonal.

$$D^+_{ii} = \begin{cases} 1/D_{ii} & \text{if } D_{ii} \neq 0 \\ 0 & \text{else.} \end{cases}$$

Suppose $A = UDV^T$

Then $A^+ = VD^+U^T$ pseudo-inverse.

(1) · if $A$ square & invertible then $A^+ = A^{-1}$

(2) if $A$ is overdetermined, $A^+\underline{b}$ gives least-squares soln to $A\underline{x} \simeq \underline{b}$

(3) · if $A$ underdetermined, $A^+\underline{b}$ gives least squares soln to $A\underline{x} \simeq \underline{b}$ w/ minimal norm.