## Top row

$U = A$   $L = I$   $V(2) = L*U(1)$   $P = I$   vectful; many RHSs inexpensive $O(n^2)$
$\downarrow$   $\downarrow$   even tho $LU = O(n^3)$, nonsingular
$Lc = Pb = [\,]$   not! sparse bc doubles # per row addt. op.
$c = [\,] = Ux$   $x = [\,]$   $\%_i \times \%^n = \%_f$

General iterative useful; sparse, large, approx sol ok
Project onto each eqn/line perpendicularly & iteratively get closer to intersection
$$x_{k+1} = x_k + w \frac{b_i - a_i \cdot x_k}{|a_i|^2} a_i^T$$
$\measuredangle = \cos^{-1} \frac{a_1 \cdot a_2}{|a_1||a_2|}$  closer to 90°, converge faster
plot eqns: #x, #y = #   ⊙ $y=0$   ⊙ $y=0$
                                 $x=0$      $x=0$

## Second row

osc if opp signs,   diff initial guess,   $V_o = [v_o \ w_o \dots]$
$e_i = A^k v_o$   diff ans   preferably sparse   $U_{kn} = A \hat{U}_k$
$e_{min} = A^{-k} \hat{v}_o$   if big for efficiency   $\hat{U}_{kn} = GS(V_{k+1})$
$(Q - \mu I) e_i = (\lambda_i - \mu I) v_o$   CAN converge but test $p \to p+1$
$\uparrow$ inverse

Hotelling if normal: $B = A - \lambda_i e_1 e_1^T$
Wielandt: $B = A - \frac{e_1 a_p}{e_{1p}}$   $e_2 = e_B + \frac{a_p e_B}{\lambda_2 - \lambda_1} \frac{e_1}{e_{1p}}$
$P = $ largest elem in $e_1$

$\hat{\hat{Q}} = d\left(Q + \frac{1}{n}\begin{bmatrix} \cdot \\ \cdot \end{bmatrix}\begin{bmatrix}\cdot & \cdot\end{bmatrix}\right) + (1-d)\frac{1}{n}\begin{bmatrix}\cdot\\\cdot\end{bmatrix}\begin{bmatrix}\cdot & \cdot\end{bmatrix}^T$
$e_1$ of $\hat{\hat{Q}}^T$ | $\underbrace{\hat{Q}}$  ↳ Perr. Prob.  $\mathbb{R}^\oplus$, $sq \to 1$ largest $\lambda_1$
$\uparrow d \to$ smaller $\lambda_2$, converge faster, but can affect ranking

## Third row

$e^{At} = S\begin{bmatrix} e^{\lambda_1 t} \\ & e^{\lambda_2 t} \\ & & \dots\end{bmatrix} S^{-1}$   $\frac{dx}{dt} = Ax$   $x(t) = \Sigma c_i e^{\lambda_i t}$
$x(t) = e^{At} x_o$   $x_o = \Sigma c_i e_i(i)$
$e^{A(t)} = I + A + \frac{At^2}{2!} + \dots$   $x_o = (e^{At})^{-1} x(t)$   $C = S^{-1} x_o$
let $x_i = e_i e^{\lambda_i t}$
$\frac{dx_i}{dt} = \frac{d}{dt} e_i e^{\lambda_i t}$
$= \lambda_i e_i e^{\lambda_i t} = A e_i e^{\lambda_i t}$
$x(t) = (S^{-1} x_o)\begin{bmatrix} e^{\lambda t}\end{bmatrix} S$

phase space! check $\frac{dx}{dt}$ @ $(\pm 1, \pm 1)$
$\lambda_I \neq 0$, osc
$\lambda_R > 0, e^{\lambda t} \to \infty$
$\lambda_I = f$
$\lambda_R < 0, e^{\lambda t} \to 0$
$\lambda_f = $ growth speed
$\uparrow$ real roots

not normal ($e_i$ not $\perp$)? non-normality: $\frac{|A^TA - AA^T|}{|A|^2}$
$B = e^{At}$   $B^T B x_o = \lambda_B x_o$   $\lambda_R << \lambda_i < 0$
Amp factor $= \frac{|x(t)|^2}{|x_o|^2} = \max \lambda_B$

Jordan! nxn mat don't have n lin ind. e_i's, no $S^{-1}$ $\begin{bmatrix} \lambda_{\dots}\end{bmatrix}$
① A of A   for $S^{-1}As =$
② $(A - \lambda I) v_1 e_o$
③ $v_o(A - \lambda I) e_o$
④ $v_3 = (A - \lambda I) e_o$   $J = M^{-1} A M$

## Fourth row

$3\times2$ more eqns than unkn. = overdet = prob no sol bc contradictory = min $|Ax-b|^2$  $x = (A^TA)^{-1} A^T b = (A^TA)^{-1} A^T \cancel{b} A^T b = A^+ b$
$2\times3$ less eqns than unkn. = underdet = ∞ sols = min $|x|$   $x = A^+ b$
ill cond = multiple min $|Ax-b|$ so also min $|x|$   $x = (A^TA)^+ A^T b = A^+ b$

① QR decomp   $A = QR$
↓ same adv.   $A^TA x = A^Tb = R^TR = R^T Q^T b$
as $LU$   $\nabla Rx = Q^Tb$

$A = [a_1 \quad a_2]$
$\hat{q}_1 = \frac{a_1}{|a_1|}$
$q_2 = a_2 - (a_2 \cdot \hat{q}_1) \hat{q}_1$
$\hat{q}_2 = \frac{q_2}{|q_2|}$

## Fifth row

$A = U\Sigma V^T$   smaller dim $A^TA$ ⊕ $AA^T$
$A^+ = V\Sigma^+ U^T$   $\lambda$; # $e_i$ of smaller
$\sigma_i = \sqrt{\lambda_i} \to$ into $\Sigma$
if $A^TA$, $[e_i \dots] = V$
if $AA^T$, $[e_i \dots] = U$
$U = \frac{Av_i}{\sigma_i}$   $V = \frac{A^Tu_i}{\sigma_i}$
remaining cols: random guess & GS to find ($\perp$ to previous cols)

$A = u_1\sigma_1 v_1^T + \dots + u_k \sigma_k v_k^T$  keep only k terms
img compression quality $= \sum_{ij}(A'_{ij} - A_{ij})^2$
for mxn, storage $= \frac{km + kn + k}{mxn}$

$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & u_2 \\ 0 & 0 & 1\end{bmatrix}^T$ $\begin{bmatrix} U V^T & 0 \\ 0 & 0 \\ 0 & 1\end{bmatrix}$ $\begin{bmatrix} V\Sigma V^T & 0 \\ 0 & 1\end{bmatrix}$ $x = Mx$
$\begin{bmatrix} \cos\theta & -\sin\theta & 0 \\ \sin\theta & \cos\theta & 0 \\ 0 & 0 & 1\end{bmatrix}$  $RS = A$

$C = $ sensitivity of $\frac{|\delta b|}{|b|}$ at max given some $\frac{|\delta b|}{|b|}$ rel error
$= \frac{\sigma_1}{\sigma_{min}} = |A||A^{-1}|$   $\frac{|\delta x|}{|x|} \leq C \frac{|\delta b|}{|b|}$
$b_{max\,er.} = v_1$   $\delta b = \frac{b}{|\delta b|/|b|}$   $\frac{\delta_k}{\sigma_1} > \frac{|\delta b|}{|b|}$  $k = $ eff rank
$b_{min\,er} = v_4$
$\frac{\max}{\text{stretch}} \frac{|Ax|}{|x|} \leq |A| = \sigma_1$ @ $v_1$   $|A^{-1}| = \frac{1}{\sigma_{min}}$

## PCA section

⓪ remove row means from each elem
PCA / multivariate $F = \begin{bmatrix} X \\ Y\end{bmatrix}$
$U$ vecs = PC's ($e_i$'s of C)   $V$ vecs = nondim time series coeffs   $\sigma_i = $ how well PC explains
$C = \frac{1}{N} FF^T$ nondiag $C_{ij} = $ cov between rows i & j   $\frac{\lambda_i}{\Sigma \lambda_i} = \frac{\sigma_i}{\Sigma \sigma_i} = \%$ var explained
diag $C_{ii} = $ row i var   $(Ce_i = \lambda_i e_i, \text{ trace}(C))$

MCA
$U$ vecs = struct of 1st dataset ($X$ in $\frac{1}{N} XY^T$) most correlated w the $V$ vec $Y$ struct.
$\frac{\sigma_i}{\Sigma \sigma_i} = \%$ of total cov by that $U_i$ vec & $V_j$ vec ($X$ struct & $Y$ struct.)

$T = U^T F = \Sigma V^T = $ amplitudes of PC's in data @ $t=n$
$\binom{\text{proj of}}{f_n \text{ onto } PC_i} f_n \cdot u_i = t_{in}$
$F_{reconst} = UT = U\Sigma V^T$ linear comb of PC's w weights $t_{in}$
$(:,k)(k,:)$

→ If cov. PC (for inter-datasets) is smaller/not as strong as intra-dataset var. explained by other PC's. Bc after 1st PC, all rest have to be orthogonal, then if cov. isn't orthogonal, multivariate will miss it. BUT MCA ignores/loses intra-dataset variability to focus on inter-dataset cov.

Total var of dataset: sum of sqs of its elems divided by M = sum of each datavec's var Represents variability in [amount of students taking courses]. If vanishes, datavecs are constant over time.

$N = \frac{1}{N} FF^T = C$   $\frac{1}{N}[\quad][\quad] = \frac{1}{N}[\quad] = [\quad]$

Total cov of 2 datasets: sum of sqs of the elems of C Represents overall covariability of [dataset X] to [dataset Y]. If vanishes, X & Y are ind.

Multivariate:
"U1 (PC1) indicates as students taking A61, A62, C72 increase, C71 decreases. PC2 shows anticorr between A61 and rest, but sigma small so negligible."
"V1 row, given signs in PC1, above or below avg at time. + in PC1, + in V1 row elem n mean above avg at time n."
"ind cov bc PC1 only 1 dataset and PC2 only other dataset"

MCA:
"U1 shows opening hours of both pools vary together with salary (V1,1) but not covid (V2,1)=0"

## Distances section

$L_r = \left[\sum_{i=1}^{dim}(x_i - y_i)^r\right]^{\frac{1}{r}}$
$L_2 = \sqrt{\Sigma(x_i - y_i)^2} = |x - y|$
$L_1 = \sum_i^{dim} |x_i - y_i|$
$L_\infty = \max |x_i - y_i|$
cos dist $= \Theta = \cos^{-1}\left(\frac{x \cdot y}{|x||y|}\right) = \cos^{-1}\left(\frac{\Sigma x_i y_i}{\sqrt{\Sigma x_i^2}\sqrt{\Sigma y_i^2}}\right) = \cos^{-1}(\text{corr. of } x \& y)$
Jaccard - sets $= 1 - \frac{|x \cap y|}{|x \cup y|}$
Edit - ordered groups $= \leq d_x^{add} d_y^{del} = 5$
Hamming - same size ordered groups = # same elems

Curse of dim!
$L_2^2 = (x_i - y_i)^2 \to $ avg dist across all/many dims
$\Theta = \cos^{-1}(\text{corr.}) \to \cos^{-1}(0)$ for random long seqs $\to 90°$
problem bc clustering needs diff dists

## Bottom middle

$\lambda_i$ poly: $\lambda^2 - (a+d)\lambda + (ad - bc) = 0$
$e_i$ ref: $[A - \lambda_i I \mid 0]$   $\text{diag}^{-1} = \begin{bmatrix} \frac{1}{\#} & \\ & \ddots\end{bmatrix}$
$|a|^2 = a^T a = \Sigma a_i^2$   $\det(A) = |A|$
normal: $A^TA = AA^T$, $e_i \perp$
symm: normal, $A = A^T$, $e_i$ are $\mathbb{R}$
orthonormal: symm, $A^TA = AA^T = I$
nonsingular: has $A^{-1}$, sq & $\det(A) \neq 0$

## Bottom right

min $x^TAx$  st $x'x = 1$  & $x^T\begin{bmatrix}\cdot\\\cdot\end{bmatrix} = 0$
$x^TAx + \lambda(1 - x^Tx) = 0$
$Ax = \lambda x$
min $\lambda$
max $(Ax)^T(Ax)$  st $x^Tx = 1$
$x^TA^TAx + \lambda(1 - x^Tx) = 0$
$A^TAx = \lambda x$
max $\lambda$
min $(Ax - b)^T(Ax - b)$  st $x^Tx = 1$
$x^TA^TAx - b^TAx - x^TA^Tb + b^Tb + \lambda(1 - x^Tx) = 0$
$A^TAx - b^TA = \lambda x$
$x = (A^TA - \lambda I)^{-1} b^TA$

$$x_c = \frac{1}{n}\sum_j^n x_j \quad \text{→ strings don't have}$$

$$var_c = \sum_j^n \sum_i^{dim}(x_{ij} - x_{ci})^2$$

radius = max $|x_j - x_c|$

diameter = max $|x_j - x_k|$

density = $\dfrac{n}{radius}$

Ward $\Delta var = \dfrac{N_A N_B}{N_A + N_B}|x_B - x_A|^2$

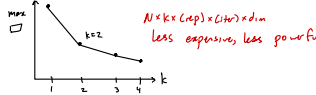$\qquad = var_{AB} - (var_A + var_B)$

single = min $|x_A - x_B|$

<span style="color:red">dist between all pts = $n(n-1) = O(n^2)$
from each new centroid = $n-1, n-2$
& find which pair to merge = $O(n^2\log n)$</span>

1) plot pts
2) manually pick & calc
3) calc & compare unsure ones

☐ of just merged cluster

| | | | | | |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 0 |
| (1 2) | | (3 4) | | 5 | 0.7 |
| (1 2) | | (3 4 5) | | | 1.4 |
| (1 2) | | 3 4 5 | | | 3.3 |

3.3 ☐
1.4
0.7

k pts → replicates
1st random
m+1 as far as possible
m+2 smallest dist as far as possible
n-k points
add pt to clustroid
adjust clustroid (or) adjust @ end
new clustroids as k pts
↳ iter. assign til converge

<span style="color:red">$N \times k \times (rep) \times (iter) \times dim$
less expensive, less powerful</span>

max ☐
k=2
k=3
1 2 3 4  k

SOM for each x datapt
→ nearest $m_j = m_j + \eta\binom{ker}{nearest}(x_i - m_j)$
→ others $m_e = m_e + \eta\binom{ker}{others}(x_i - m_e)$

1) calc $dist_{ij}$
2) $W_{ij} = e^{\left(\frac{-dist_{ij}}{dist_{avg}}\right)}$
3) D = sum over W cols (diag)
4) normalise $L = D - W$
$\quad L' = D^{-1}L = 1 - D^{-1}W \to$ better spectral gap
5) $\lambda$ value

0   k=3
$\lambda_{min} \to \lambda_i$

6) k-means in $[e_2, e_3 ... e_k]$

$e_3$ clums
$x' = \begin{bmatrix} e_2^T \\ e_3^T \end{bmatrix}$  – lose orig. dim
+ smaller dim $e_i$?
∴ + exclude noise $e_i$
– calc k $\lambda$'s tho (slower method) bc L sym
$e_2$ clums

turns out $x^T L x = \frac{1}{2}\sum\sum W_{ij}(x_i - x_j)^2$
so min $x^T L x$   $Lx = \lambda x$
wrt $x$ ← [..]  $\lambda \geq 0$ so
$2^{nd}$ smallest $\lambda$

---

$w' = \begin{bmatrix} w \\ -b \end{bmatrix}$   $z' = \begin{bmatrix} x \\ 1 \end{bmatrix}$

pt 1
$w^T z = 0.6$   y=1  ok
pt 2
$w^T z = 6$   y=-1  no
$w = w + \eta(-1)z = [\quad]$
iter. thru pts
plot: $w_1 x_1 + w_2 x_2 + w_3^{(-b)} = 0$
$\quad x_2 = \frac{-w_1}{w_2}x_1 - \frac{w_3}{w_2}$

fails XOR:

$w' = \begin{bmatrix} w \\ -b \end{bmatrix}$   $x' = \begin{bmatrix} x \\ 1 \end{bmatrix} = a'$
$0 \times N$
$z^2 = w^2 a^1$   $a^2 = \sigma(z^2)$
$z^3 = w^3 a^2$   $a^3 = \sigma(z^3)$

| | | $\sigma'(z)$ |
|---|---|---|
| Relu | $\begin{cases}0, z<0 \\ z, z\geq 0\end{cases}$ | $\begin{cases}0, z<0 \\ z, z\geq 0\end{cases}$ |
| tanh | $\frac{e^z - e^{-z}}{e^z + e^{-z}}$ | $1 - \left(\frac{e^z - e^{-z}}{e^z + e^{-z}}\right)^2$ |
| sigmoid pick | $(1+e^{-z})^{-1}$ | $e^{-z}(1+e^{-z})^{-2}$ |
| softmax pick | $\max \frac{e^{z_i}}{\sum_i e^{z_i}}$ | |

if lots of data:

| parameters | training | testing |
|---|---|---|
| | | |

Overfit? 1) more data 2) smaller network bc small can't perfectly fit training data 3) add dummy data 4) terminate optimization early, 5) add regularization term bc w's and b's go up when overfit

$$C = C_0 + \lambda\left\{\sum_{k,i,j}(w_{ij}^\ell)^2 + b_i^2\right\}$$
w's & b's ↑ when overfit

① get δ's   backprop gives grad
$$\delta^L = \frac{dC}{dz^L} = \frac{dC}{da^L}\cdot\frac{da^L}{dz^L}$$
$$\frac{dC}{da^L} = \frac{d}{da^L}\left[\tfrac{1}{2}(y-a^L)^2\right] = a^L - y = \nabla_{a^L}C$$
$$\frac{da^L}{dz^L} = \frac{d}{dz^L}(1+e^{-z^L})^{-1} = e^{-z^L}(1+e^{-z^L})^{-2} = \sigma'(z^L)$$
$$\delta^L = \nabla_{a^L}C \cdot \sigma'(z^L)$$
$$\delta^\ell = (w^{\ell+1})^T \delta^{\ell+1} \cdot \sigma'(z^\ell)$$

② adj w   steepest descent
$$\nabla_{w^\ell}C = \delta^\ell(a^{\ell-1})^T \longrightarrow \text{calc avg gradient for minibatch,}$$
$\qquad\qquad\uparrow$ from training X     go through all minibatches
$$w^\ell = w^\ell - \eta\nabla_{w^\ell}C$$

③ each time thru all pts = epoch

perturb each wt
$C_1 = C(w_{ij}^k)$
$C_2 = C(w_{ij}^k + \delta w_{ij}^k)$
⇓ 2x feedforward
$\frac{dC}{dw_{ij}^k} \approx \frac{C_2 - C_1}{\delta w_{ij}^k}$
BUT need to repeat feedforward for each $w_{ij}^k$ inefficient

---

$$\frac{d^2x}{dt^2} = a\frac{dx}{dt} + bx$$
define $y = \frac{dx}{dt}$
$$\frac{dy}{dt} = ay + bx$$
$$\frac{dx}{dt} = y$$
$$\frac{d}{dt}\begin{bmatrix}x\\y\end{bmatrix} = A\begin{bmatrix}x\\y\end{bmatrix}$$
$$A = \begin{bmatrix}0 & 1\\a & b\end{bmatrix}$$

$\frac{dx}{dt} = ax + by$   (1)
$\frac{dy}{dt} = cx + dy$   (2)
solve for y (1)
$y = \frac{dx}{dt}\frac{1}{b} - \frac{ax}{b}$   (3)
sub (2)
$\frac{dy}{dt} = cz + \frac{dx}{dt}\frac{d}{b} - \frac{ad}{b}x$
$\frac{d^2x}{dt^2} = a\left(\frac{dx}{dt}\right) + b\left(\frac{dy}{dt}\right)$
$\quad = a^2x + aby + bcx + bdy$
plug in (3) for y
$\frac{d^2x}{dt^2} = (a+d)\frac{dx}{dt} + (bc-ad)x$

# Decision trees

Most informative near top branches    Overfit? limit growth, pruning

takes into acc both yes's & no's

$$I_G(\text{interesting}) = \frac{N_{yes}}{N} \sum \text{impurity for yes}$$
$$+ \frac{N_{no}}{N} \sum \text{impurity for no}$$

$$P(\checkmark | yes) = \frac{9 \checkmark yes}{12 \, yes}$$

max @ 50%, min @ 0 or 100%   "impurity" $= P(1-P)$

$$I_G = \frac{N_{yes}}{N}\left[ P(\checkmark|yes)(1 - P(\checkmark|yes)) + P(\times|yes)(1 - P(\times|yes)) \right]$$

$$+ \frac{N_{no}}{N}\left[ P(\checkmark|no)(1 - P(\checkmark|no)) + P(\times|no)(1 - P(\times|no)) \right]$$

## multi labels/endings

$$I_G(A) = \frac{N_{yes}}{N} \sum_{i=1}^{L} P(\ell|yes)(1 - P(\ell|yes)) + \frac{N_{no}}{N} \sum_{i=1}^{L} P(\ell|no)(1 - P(\ell|no))$$

---

### discrete

$N_{yes} = 32$    $P(\checkmark|yes) = \frac{16}{32}$    $P(\times|yes) = \frac{16}{32}$

$N_{no} = 26$    $P(\checkmark|no) = \frac{11}{26}$    $P(\times|no) = \frac{15}{26}$

$N = \quad + \quad =$

$I_G = -( \quad + \quad ) + -( \quad + \quad )$  lowest $I_G$ top

_____ ignore top attrib's yes rows

_____ ignore top attrib's no rows

determine endings based on majority

### cont & discrete labels

$I_G$ for each threshold split

$I_G$'s for subsequent subset thresholds

### cont & cont labels

Min var on each leaf
around each leaf's mean

"impurity" $= \frac{N_c}{N} Var_< + \frac{N_>}{N} Var_>$
("$I_G$")
$$Var(y<) = \frac{1}{N_c} \sum_{i}^{N_c} (y_i - y_{avg<})^2$$

$$\Rightarrow \overset{actual}{split} = \frac{x_{A last} + x_{B first}}{2}$$

With 2 continuous attrib., calc var.'s for all possible splits based on attrib. 1 and then for all possible splits based on attrib. 2 and then pick lowest split & attrib.