# Finding Principal Components

1. find eigenvalues by solving: $\det(\Sigma - \lambda I) = 0$

$$\det\begin{pmatrix} 2.0 - \lambda & 0.8 \\ 0.8 & 0.6 - \lambda \end{pmatrix} = (2 - \lambda)(0.6 - \lambda) - (0.8)(0.8) = \lambda^2 - 2.6\lambda + 0.56 = 0$$
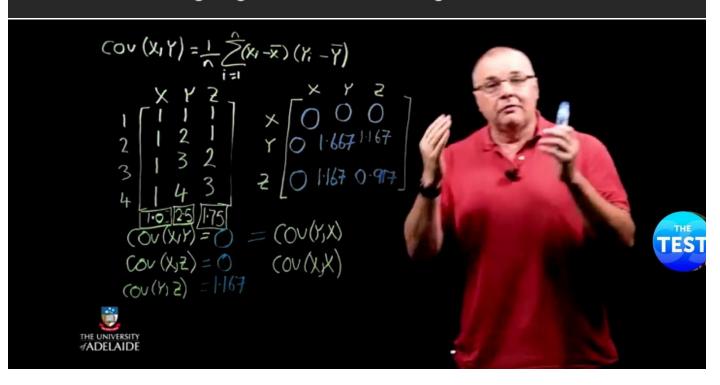
$$\{\lambda_1, \lambda_2\} = \tfrac{1}{2}\left(2.6 \pm \sqrt{2.6^2 - 4*0.56}\right) = \{2.36, 0.23\}$$

2. find $i^{th}$ eigenvector by solving: $\Sigma \, \mathbf{e}_i = \lambda_i \, \mathbf{e}_i$

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix}\begin{pmatrix} e_{1,1} \\ e_{1,2} \end{pmatrix} = 2.36\begin{pmatrix} e_{1,1} \\ e_{1,2} \end{pmatrix} \Rightarrow \begin{aligned} 2.0e_{1,1} + 0.8e_{1,2} = 2.36e_{1,1} \\ 0.8e_{1,1} + 0.6e_{1,2} = 2.36e_{1,2} \end{aligned} \Rightarrow e_{1,1} = 2.2e_{1,2}$$

$$e_1 \sim \begin{bmatrix} 2.2 \\ 1 \end{bmatrix}$$

want: $\|e_1\| = 1$

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix}\begin{pmatrix} e_{2,1} \\ e_{2,2} \end{pmatrix} = 0.23\begin{pmatrix} e_{2,1} \\ e_{2,2} \end{pmatrix} \Rightarrow e_2 = \begin{bmatrix} -0.41 \\ 0.91 \end{bmatrix}$$

$$e_1 = \begin{bmatrix} 0.91 \\ 0.41 \end{bmatrix}$$

slope: 0.454

PCA 5: finding eigenvalues and eigenvectors

$$\text{cov}(X, Y) = \frac{1}{n}\sum_{i=1}^{\hat{n}}(x_i - \bar{x})(y_i - \bar{y})$$

| | X | Y | Z |
|---|---|---|---|
| 1 | 1 | 1 | 1 |
| 2 | 1 | 2 | 1 |
| 3 | 1 | 3 | 2 |
| 4 | 1 | 4 | 3 |
| | 1.0 | 2.5 | 1.75 |

|   | X | Y | Z |
|---|---|---|---|
| X | 0 | 0 | 0 |
| Y | 0 | 1.667 | 1.167 |
| Z | 0 | 1.167 | 0.917 |

$\text{cov}(X,Y) = 0 = \text{cov}(Y,X)$

$\text{cov}(X,Z) = 0 \qquad \text{cov}(X,X)$

$\text{cov}(Y,Z) = 1.167$

THE UNIVERSITY
of ADELAIDE

THE TEST

# PCA algorithm II
# (sample covariance matrix)

- Given data $\{\mathbf{x}_1, ..., \mathbf{x}_m\}$, compute covariance matrix $\Sigma$

$$\Sigma = \frac{1}{m} \sum_{i=1}^{m} (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T \quad \text{where} \quad \bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^{m} \mathbf{x}_i$$

- **PCA** basis vectors = the eigenvectors of $\Sigma$

- Larger eigenvalue $\Rightarrow$ more important eigenvectors

26

# Reminder: Eigenvector and Eigenva

$$Ax = \lambda x$$

A: Square matrix
$\lambda$: Eigenvector or characteristic vector
$x$: Eigenvalue or characteristic value

*Show* $x = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ *is an eigenvector for* $A = \begin{bmatrix} 2 & -4 \\ 3 & -6 \end{bmatrix}$

*Solution* : $Ax = \begin{bmatrix} 2 & -4 \\ 3 & -6 \end{bmatrix}\begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

*But for* $\lambda = 0$, $\lambda x = 0\begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

*Thus, x is an eigenvector of A, and* $\lambda = 0$ *is an eigenvalue.*

27

## SVM

1) Draw decision boundary
2) Write equation for boundary $y = x \cdot b = 0$
3) Re-write equation as $w \cdot \vec{x} + b = 0$

Scale equation $\vec{w} \cdot \vec{x} + b = 0$

① Support vector
② Support vector estimate
$c(\vec{c} \cdot \vec{v})[\vec{x}] + c(-s) = -1 \Rightarrow c = \frac{1}{2}$

$w = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \end{bmatrix} \quad b = -4$

1) For all non support vectors, $\alpha = 0$
otherwise, the point lies on the margin and $\alpha > 0$

2) $\sum t_i g_i = 0 \Leftrightarrow \sum_{s.v.} \frac{1}{s.v.}$

3) $\sum \alpha_i y_i x_i = \vec{w} \Leftrightarrow \vec{w} = \sum_{s.v.}$ s.v. vectors

Calculate $\alpha \Rightarrow$ margin Width $= 2/\|w\|$

$K(\vec{u}, \vec{v}) = \varphi(u)^T \varphi(v) \Rightarrow$ No change
Polynomial $\Rightarrow K(u,v) = (u \cdot v + b)^n$
$K(u,v) = e^{-\frac{\|u-v\|^2}{2\sigma^2}} \Rightarrow$ Radial Basis Function

### Decision Tree

$H(x) = -\sum_i P(x=v) \log_2 P(x=v)$

$IG(x) = H(x) - H(Y|X)$

$H(Y|X) = -\sum_{i=1}^{n} P(x=x_i) \sum_i P(y=y_i | x=x_i) \log_2 P(y=y_i | x=x_i)$

Random Forest

SVM $\Rightarrow$ Why large margin? relative to uncertainty

Clustering: Kmeans: 1) initialize K cluster centroids randomly drawn without replacement
2) Repeat until convergence
   → Assign each data point to the cluster with the closest center
   → Assign each cluster center to be the mean of its cluster's data point
   $\mu_k \leftarrow |S_k|^{-1} \sum_{i} x_i$

* K-means always terminate.
* Good dissimilarity only useful for compactness clustering

### Boosting

1) Initialize weights $w_i = 1/N$
2) Calculate weighted error rate for each $h_i$
   $\epsilon = \sum w_i$
3) Pick best $h$ with smallest error rate
4) Calculate voting power of $h$
   $\alpha = \frac{1}{2} \ln\left(\frac{1+\epsilon}{\epsilon}\right)$
5) Finished?
   $H(x) = sign(\sum \alpha_i h_i(x))$ good enough?
6) Update weight

Boosting $\Rightarrow$ is machine learning ensemble algorithm

### Spectral Clustering / Principal Component Analysis

$W_{ij} = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$

$min \ cut = cut(A,B) = \sum W(u,v)$

k-means vs. Spectral Clustering
Applying k-means to eigenvectors

### Hierarchical Clustering

Bottom-up $\Rightarrow$ Agglomerative clustering
Top-down $\Rightarrow$ Divisive

Complete Link (max) $\rightarrow$ Tight clusters

Average Link $\rightarrow$ The most widely used measure

#### Internal
→ Intra-class similarity high / inter-class is low

#### External
Entropy, Purity

$Purity(\Omega, C) = \frac{1}{N} \sum_j max_j |w_k \cap c_j|$

$Entropy(CC_i) = \sum_{i} -P(c_j | CC_i) \log_2 P(c_j | CC_i)$

# Reminder: Eigenvector and Eigenvalue

Example 1: Find the eigenvalues of

$$A = \begin{bmatrix} 2 & -12 \\ 1 & -5 \end{bmatrix}$$

$$|\lambda I - A| = \begin{vmatrix} \lambda - 2 & 12 \\ -1 & \lambda + 5 \end{vmatrix} = (\lambda - 2)(\lambda + 5) + 12$$

$$= \lambda^2 + 3\lambda + 2 = (\lambda + 1)(\lambda + 2)$$

two eigenvalues: $-1, -2$

*Note:* The roots of the characteristic equation can be repeated. That is, $\lambda_1 = \lambda_2 = \ldots = \lambda_k$. If that happens, the eigenvalue is said to be of multiplicity k.

Example 2: Find the eigenvalues of

$$A = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

$$|\lambda I - A| = \begin{vmatrix} \lambda - 2 & -1 & 0 \\ 0 & \lambda - 2 & 0 \\ 0 & 0 & \lambda - 2 \end{vmatrix} = (\lambda - 2)^3 = 0$$
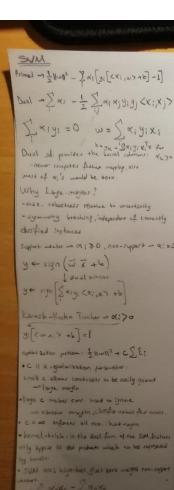
$\lambda = 2$ is an eigenvector of multiplicity 3.

# PCA algorithm II (sample covariance matrix)

**Goal:** Find $r$-dim projection that best preserves variance

1. Compute mean vector $\mu$ and covariance matrix $\Sigma$ of original points

2. Compute eigenvectors and eigenvalues of $\Sigma$

3. Select top $r$ eigenvectors

4. Project points onto subspace spanned by them:

$$y = A(x - \mu)$$

where $y$ is the new point, $x$ is the old one, and the rows of $A$ are the eigenvectors

# SVM

Primal $\rightarrow \frac{1}{2}\|w\|^2 - \sum_i \alpha_i [y_i[\langle x_i, w\rangle + b] - 1]$

Dual $\rightarrow \sum_i \alpha_i - \frac{1}{2}\sum_{ij}\alpha_i\alpha_j y_i y_j \langle x_i, x_j\rangle$

$\sum_i \alpha_i y_i = 0 \qquad w = \sum_i \alpha_i y_i x_i$

$b = y_k - \sum_i \alpha_i y_i x_i^T x_k$ for

Dual sol. provides the kernel solutions. $\alpha_k > 0$

- never computes feature mapping, also
most of $\alpha_i$'s would be zero.

## Why Large-margins?

- Max. robustness relative to uncertainty
- symmetry breaking, independent of correctly classified instances

Support vector $\rightarrow \alpha_i > 0$, non-support $\rightarrow \alpha_i = 0$

$y \leftarrow \text{sign}(\vec{w}\,\vec{x} + b)$
$\qquad \downarrow$ dual solution
$y \leftarrow \text{sign}[\sum_i \alpha_i y_i \langle x_i, x\rangle + b]$

Karush-Kuhn Tucker $\rightarrow \alpha_i > 0$

$y_i[\langle w, x_i\rangle + b] = 1$

Optimal problem: $\frac{1}{2}\|w\|^2 + C\sum_i \xi_i$

- C is a regularization parameter:

small c allows constraints to be easily ignored
$\rightarrow$ large margin

- large c makes cons. hard to ignore
$\rightarrow$ narrow margin, classifier makes few errors
- C = ∞ enforces all cons.: hard margin

- kernel-trick: in the dual form of the SVM, features only appear as dot products which can be represented by kernels.

- SVM uses hinge-loss, gives zero weight non-support vectors

$w = \sum_{i=1}^{n}\alpha_i\phi x_\phi - \sum_{n\phi}\alpha_n x_n$

Partitioning: k-means, spectral clust.
hierarchical clust. $\rightarrow$ bottom-up agglomerative

# Decision Trees

$IG(X) = H(Y) - H(Y|X)$

$H(Y) = -\sum_{i=1}^{k} P(Y=y_i)\log_2 P(Y=y_i)$

$H(Y|X) = -\sum_{i=1}^{n} P_x(x_i) H_{Y|X=x_i}$

$= \sum_{i=1}^{n} P_x(x_i)\left(-\sum_{j=1}^{m} P_{Y|X}(y_j|x_i)\log P_{Y|X}(y_j|x_i)\right)$

Decision tree will overfit:
$\quad$ in order to avoid
$\qquad \hookrightarrow$ fixed depth
$\qquad\quad$ fixed number of leaves
$\qquad\quad$ use random forest
$\qquad\quad$ pruning
$\qquad\quad$ early stopping

## High entropy | Low entropy

| High entropy | Low entropy |
|---|---|
| Y is from a uniform like distribution, flat histogram, values sampled from it are less predictable | Y is from a varied (peaks and valleys) dist. Histogram has many lows and highs, values sampled from it are more predictable |

- standard decision trees have no learning bias, training set error is always zero (noiseless)
- decision tree can represent any Boolean func.
linearly sep. $\rightarrow$ 100% accuracy $\rightarrow$ lin. classifier
$\hookrightarrow$ Boolean func. $\rightarrow$ decision tree

## Bias/Variance Tradeoff

| | LV | HV |
|---|---|---|
| LB | ⊙ | ⊙ |
| HB | ⊙ | ⊙ |

- simple learners are good
- HB, LB, low variance, don't usually overfit
- simple learn. are bad
- high bias, can't solve hard prob.

Bagging (Bootstrap Aggregating)

$\text{Var}(\text{Bagging}(L(x,D))) = \frac{\text{Var}(L(x,D))}{N}$

Random Forests: ensemble method designed for decision tree classifier

two sources of randomness: bagging and random input vectors.

$\Rightarrow$ Bagging method: each tree is grown using a bootstrap sample of training data
$\Rightarrow$ Random vector method: At each node, best split is chosen from a random sample of m attributes instead of all attributes

Bagging reduces variance by avg. bagging has little effect on bias. Can we always reduce bias? Yes: Boosting!

# Boosting

we use weak learners to create strong learner

$H(x) = \text{sign}\left(\sum_{i=1}^{T}\alpha_i \cdot h_i(\vec{x})\right)$

minimize the error:

$\epsilon_t = P_{i\sim D_t}[h_t(x_i) \neq y_i]$

$\alpha_t = \frac{1}{2}\ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$

$w = \begin{cases} \frac{1}{2}\cdot\frac{1}{\epsilon} & w_{old} \rightarrow \text{wrong} \\ \frac{1}{2}\cdot\frac{1}{1-\epsilon} & w_{old} \rightarrow \text{correct} \end{cases}$

## Boosting vs. Bagging

- resample data points, reweights data points.
- weight of each classifier is the same, dependent on class accuracy
- only variance reduction, both bias and variance reduced

## Clustering:

k-means always terminate,
is the clustering any good?
- Global dissimilarity only useful for comparing clustering.

## Spectral clustering

$w_{ij} = \begin{cases} e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}} & \rightarrow\text{centralize of neighborhood} \end{cases}$

small $\sigma$: nearby points
min. cut: $\text{cut}(A,B) = \sum_{u\in A, v\in B} W(u,v)$ with $A\cap B = \emptyset$

k-mean vs spect. clust.
Applying k-means to Laplacian eigenvectors allows us to find cluster with non-convex boundary. $\rightarrow$ applies unless $u$

hierarchical: still have to choose # of clust.

$\#\text{leaves} = (2n-3)!/[2^{(n-2)}\cdot(n-2)!]$

single link: potentially long and skinny clusters
complete link: tight clusters
avg. link: robust against noise
+ provides a hierarch. of clusters, clusters have adaptive shapes
- may have unbalanced clusters
$\qquad$ Good clustering

| Internal criterion | External criteria |
|---|---|
| intra-class similarity, the measured quality of cluster. depends obj. repr. and the sim. being measured | purity entropy of classes in clusters $\downarrow$ $\frac{1}{N}\sum_k \max_j|w_k \cap C_j|$ |

# PCA: orthogonal projection of the data onto a lower-dim. linear space that...

- maximizes variance of projected data
- minimizes mean squared distance between
- data-point, - projections

PCA #1: points in the direction of the largest variance, each subsequent principal component is orthogonal to preceding one

PCA Algorithm 1: Given centered data
$x_1, ..., x_m$ compute the principal vectors

$v_1 = \text{Argmax}_{\|w\|=1} \frac{1}{m}\sum_{i=1}^{m}\{(w^T x_i)^2\}$ 1st PCA

$x'$ PCA reconstruction: $w_1 w_1^T x_i$
PCA basis vectors = the eigenvectors of $\Sigma$.
Larger eigenvalue $\rightarrow$ more important eigenvectors

PCA Alg. II: Given data $x_1, ... x_m$ compute covariance matrix $\Sigma$

$\Sigma = \frac{1}{m}\sum_{i=1}^{m}(x_i - \bar{x})(x_i - \bar{x})^T$

eigenvalues = $|\lambda I - A|$ non solutions

PCA alg. III: (SVD of the data matrix)
Singular Value Decomp. of the centered data matrix X.

$X = [x_1, ... x_m]\in \mathbb{R}^{N\times M}$, $M$: number of ins.
$N$: dim

$X_{\text{Features}\times\text{samples}} = USV^T$

$X = U \qquad S \qquad V^T$

columns of $U$: the principal vectors
matrix $S$: diagonal, shows importance of each eigenvector
columns of $V^T$: the coefficient for reconstructing the samples