$\underline{\underline{BH}}$ STRATOVAN

■ POINT SETS - DISTRIBUTIONS -

— DISTANCE, METRIC - PCA etc.

[ WHY? Histograms/distributions viewed as points in high-dim. space ]

• "Affine Invariant Norm" for point sets:

→ Given: $n$ points $P_i = (x_1^i, x_2^i, \ldots x_K^i)^T$, $i = 1 \ldots n$, in $K$-dim. space

→ Norm '$\|\cdot\|$' defines square of length of vector $V = (x_1, \ldots, x_K)^T$ :
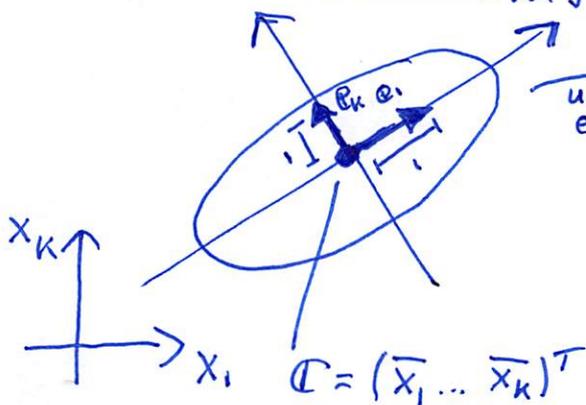
$$\|v\|^2 = V^T Q V , \qquad \text{where}$$

$$Q = n\,(D^T D)^{-1} , \qquad \text{where}$$

$$D = \begin{pmatrix} x_1^1 - \overline{x_1} & \cdots & x_K^1 - \overline{x_K} \\ \vdots & & \vdots \\ x_1^n - \overline{x_1} & \cdots & x_K^n - \overline{x_K} \end{pmatrix} , \quad \text{where}$$

$$\overline{x_j} = \frac{1}{n} \sum_{i=1}^{n} x_j^i , \quad j = 1 \ldots K$$

→ $Q$ has REAL eigenvalues $\lambda_1, \ldots, \lambda_K$

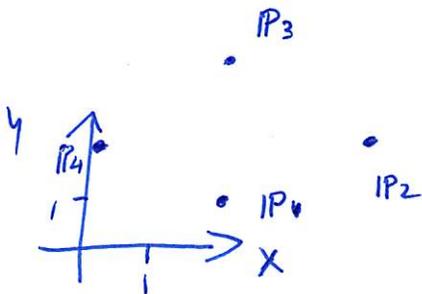( and [normalized] corresponding eigenvectors $e_1, \ldots, e_K$ )



unit ellipse

$x_K$

$x_1$ $\quad \mathbb{C} = (\overline{x_1} \ldots \overline{x_K})^T$

→ Condition for points $p$ on UNIT ellipsoid:

$$\boxed{p^T Q p = 1}$$

$$\boxed{\ell_1^2 = 1/\lambda_1 , \ldots, \ell_K^2 = 1/\lambda_K}$$

are the squared lengths of ellipsoid's principal axes.

$\qquad$ <u>BH</u> $\qquad$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ STRATOVAN
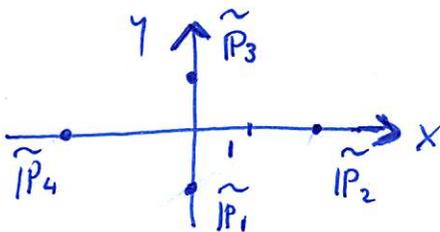
- <u>Simple 2D example</u>

$$n = 4:$$

$$P_1 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \quad P_2 = \begin{pmatrix} 4 \\ 2 \end{pmatrix}, \quad P_3 = \begin{pmatrix} 2 \\ 3 \end{pmatrix}, \quad P_4 = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$$

(i) $\quad C = \dfrac{1}{4} \displaystyle\sum_{i=1}^{4} P_i = \dfrac{1}{4}\begin{pmatrix} 8 \\ 8 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$

subtract mean $C$

$$\boxed{\tilde{P_i} = P_i - C}$$

(ii) Perform <u>mean-subtraction</u> and <u>compute</u> $\underline{D}$ ($P_i \mapsto \tilde{P_i}$):

$$\underline{D} = \begin{pmatrix} [ & \tilde{P_1} & ] \\ [ & \tilde{P_2} & ] \\ [ & \tilde{P_3} & ] \\ [ & \tilde{P_4} & ] \end{pmatrix} = \begin{pmatrix} 0 & -1 \\ 2 & 0 \\ 0 & 1 \\ -2 & 0 \end{pmatrix}$$

(iii) Compute $Q$:

$$D^T D = \begin{pmatrix} 0 & 2 & 0 & -2 \\ -1 & 0 & 1 & 0 \end{pmatrix}\begin{pmatrix} 0 & -1 \\ 2 & 0 \\ 0 & 1 \\ -2 & 0 \end{pmatrix} = \begin{pmatrix} 8 & 0 \\ 0 & 2 \end{pmatrix}$$

$$\Rightarrow (D^T D)^{-1} = \begin{pmatrix} 1/8 & 0 \\ 0 & 1/2 \end{pmatrix}$$

$$\Rightarrow Q = n \cdot (D^T D)^{-1} = 4 \cdot \begin{pmatrix} 1/8 & 0 \\ 0 & 1/2 \end{pmatrix} = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix}$$
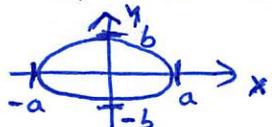
(iv) <u>Eigenvalues & eigenvectors of $Q$:</u>

$$\begin{vmatrix} \frac{1}{2}-\lambda & 0 \\ 0 & 2-\lambda \end{vmatrix} = 0 \Rightarrow \lambda_1 = \frac{1}{2}, \quad \lambda_2 = 2$$

$$\begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix}\begin{pmatrix} e_x \\ e_y \end{pmatrix} = \frac{1}{2}\begin{pmatrix} e_x \\ e_y \end{pmatrix} \Rightarrow e_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}; \quad \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix}\begin{pmatrix} e_x \\ e_y \end{pmatrix} = 2\begin{pmatrix} e_x \\ e_y \end{pmatrix} \Rightarrow e_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

(v) <u>Length of ellipse's principal axes:</u>

$$\ell_1^2 = 1/\lambda_1 = 2 \Rightarrow \ell_1 = \sqrt{2}; \quad \ell_2^2 = 1/\lambda_2 = \frac{1}{2} \Rightarrow \ell_2 = \frac{\sqrt{2}}{2}$$

(vi) Ellipse: $\qquad \boxed{\left(\dfrac{x}{a}\right)^2 + \left(\dfrac{y}{b}\right)^2 = 1}$ $\quad$ Here: $\begin{array}{l} a = \sqrt{2} \\ b = \sqrt{2}/2 \end{array}$ :
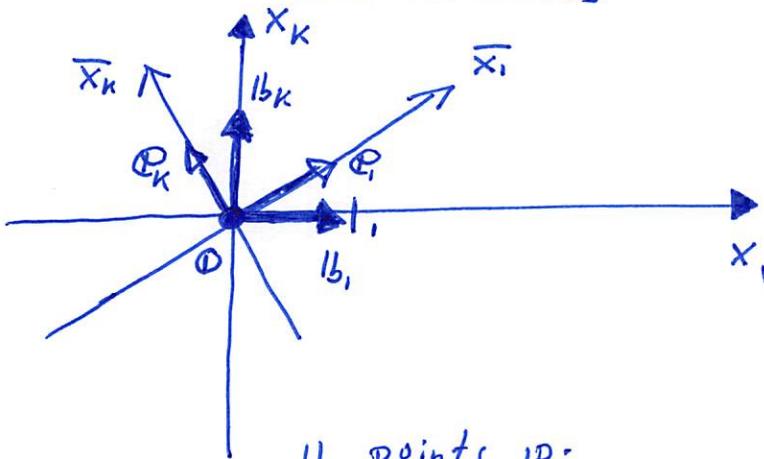
# ■ USING THE PCA-based EIGENDIRECTIONS

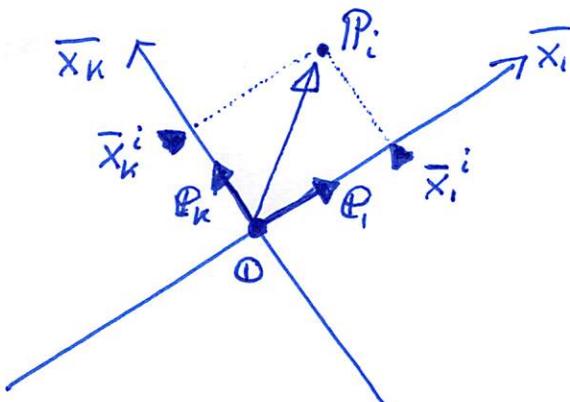## OF A POINT SET TO DEFINE "TIGHT" OBJECT/

## POINT SET-ALIGNED BOUNDING BOXES (PSABBs)

• ILLUSTRATION

[ MEAN SUBTRACTION
 HAS ALREADY BEEN DONE..]



points $P_i$
must be
expressed w.r.d.
eigensystem;
then the MINIMAL
PSABB is computed.



1) Original system:
$$\{0,\, b_1,\, ..,\, b_K\}$$
$\rightarrow$ coordinates $x_1, .., x_K$

2) "Eigen system":
$$\{0,\, e_1,\, ..,\, e_K\}$$
$\rightarrow$ coordinates $\overline{x}_1, .., \overline{x}_K$

3) Basis vectors $b_i$ & $e_i$
are all normalized.

4) Coordinates of $P_i$
w.r.t. $e_1, .., e_K$:
$$\overline{x}_1^{\,i} = P_i \bullet e_1$$
$$\vdots$$
$$\overline{x}_K^{\,i} = P_i \bullet e_K$$

$$\Rightarrow \overline{x}_j^{\,i} = P_i \bullet e_j\,, \quad !$$
$$j = 1 .. K$$

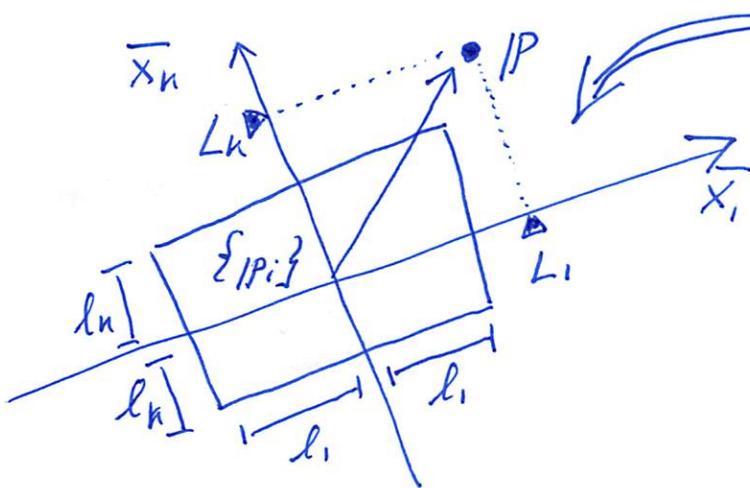5) Establish minimal **PSABB** for set of $\{P_i\}$:

→ Determine the points $P_1^{max}, \ldots,$ $P_K^{max}$ with __maximal value for coordinate $\overline{x_j}$__: (absolute value!)

→ __minimal PSABB's__ edge lengths can be computed; compute values for $l_1, \ldots, l_K$.

6) __Utilization__: Can quickly determine whether a __"new"__ point $P$ (expressed relative to $e_1, \ldots, e_K$ — after mean subtraction) "is inside the point set $\{P_i\}$ or not:
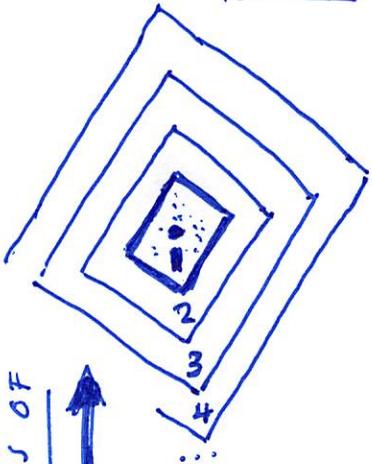
• here: $|L_1| > l_1$

$\Rightarrow$ "$P \notin \{P_i\}$"

• Can quickly __"classify"__ point $P$ as being of type $\{P_i\}$ — or not.

BH

■ CAN USE POINT DENSITY & PCA-based METRIC OF DEFINED BY THE DISTINCT CLUSTERS TO ALSO DETERMINE TO WHAT POINT CLUSTER A NEW POINT $P$ IS CLOSER TO:
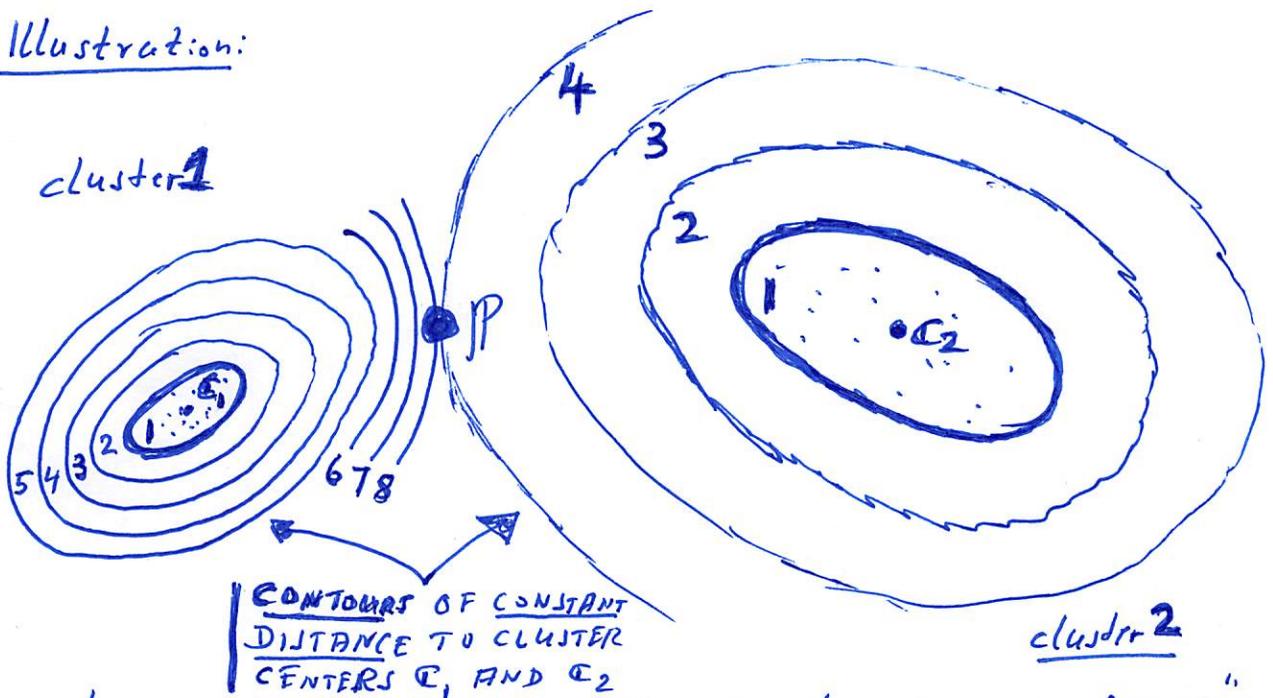
- <u>Idea</u>: The <u>eigenvalues</u> & <u>eigenvectors</u> associated with a cluster <u>define cluster-specific off-sets</u> (or contours of 'distance from cluster center') — that can be used to define a point $P$'s distance to various clusters / cluster centers:
[<u>SEE REF</u>: Louis Feng, Ingrid Hotz, Bernd Hamann, Ken Joy — "Anisotropic Noise Samples"]

- <u>Illustration</u>:

cluster **1**

cluster **2**

<u>CONTOURS</u> OF <u>CONSTANT</u> <u>DISTANCE</u> TO CLUSTER CENTERS $C_1$ AND $C_2$

<u>here</u>: "$P$ is <u>8</u> units away from $C_1$, 4 units away from $C_2$;
⇒ $P$ is <u>closer</u> to cluster **2** !!!"

EFFICIENT USE OFF-SETS OF PSABBs :: APPROXIMATION:

BH