

Before the layered perceptron...

Fisher's Linear Discriminant Function

• Preliminaries: Gradients

- Let $f(\vec{w})$ be a scalar function of the vector $\vec{w} = [w_1, w_2, \dots, w_d]^T$.
- Then $\nabla_{\vec{w}} f(\vec{w})$ is a vector with components $\nabla_{\vec{w}} f(\vec{w}) = \left[\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \dots, \frac{\partial f}{\partial w_d} \right]^T$
- Example: $f(\vec{w}) = \vec{w}^T \underline{A} \vec{w} = \sum_{i=1}^d w_i \sum_{j=1}^d a_{ij} w_j$

Then

$$\begin{aligned}
 \left(\nabla_{\vec{w}} \vec{w}^T \underline{A} \vec{w} \right)_n &= \frac{\partial}{\partial w_n} \vec{w}^T \underline{A} \vec{w} \\
 &= \frac{\partial}{\partial w_n} \sum_{i=1}^d \sum_{j=1}^d w_i a_{ij} w_j \\
 &= \sum_{j=1, j \neq n}^d a_{nj} w_j + \sum_{i=1, i \neq n}^d w_i a_{in} + 2 w_n a_{nn} \\
 &= \sum_{j=1}^d a_{nj} w_j + \sum_{i=1}^d w_i a_{in} \\
 &= \underbrace{\vec{w}^T \underline{a}_{nr}}_{\substack{\text{n}^{\text{th}} \text{ row} \\ \text{of } \underline{A}}} + \underbrace{\underline{w}^T \underline{a}_{nc}}_{\substack{\text{n}^{\text{th}} \text{ column} \\ \text{of } \underline{A}}}
 \end{aligned}$$

- Special case: $\underline{A} = \underline{A}^T$ (symmetric matrix)

$$\Rightarrow \underline{a}_{nr} = \underline{a}_{nc}$$

$$\left(\nabla_{\vec{w}} \vec{w}^T \underline{A} \vec{w} \right)_n = 2 \underline{a}_{nr}^T \vec{w}$$

and:

$$\nabla_{\vec{w}} \vec{w}^T \underline{A} \vec{w} = 2 \underline{A} \vec{w} \quad (\text{Der})$$

● Problem Statement: Data in d -dimensions

Q: How do we project this data into one dimension so that there is maximal separation?

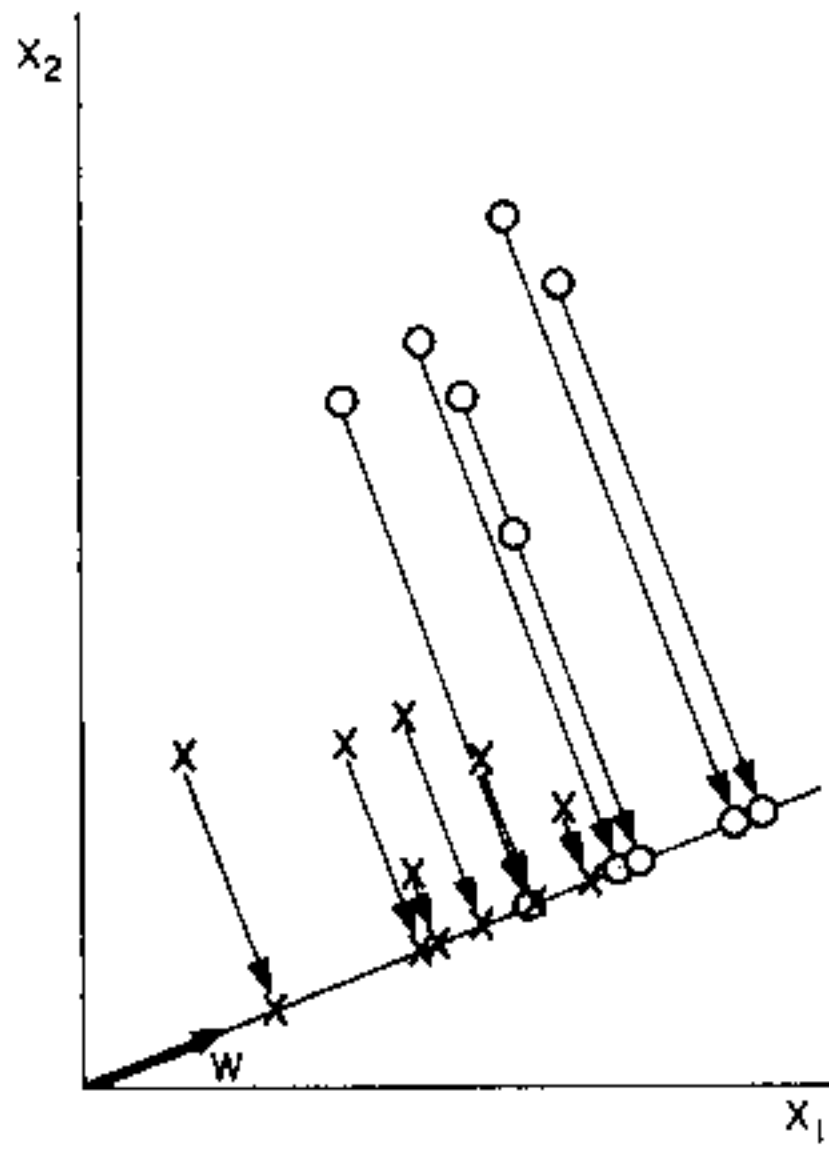
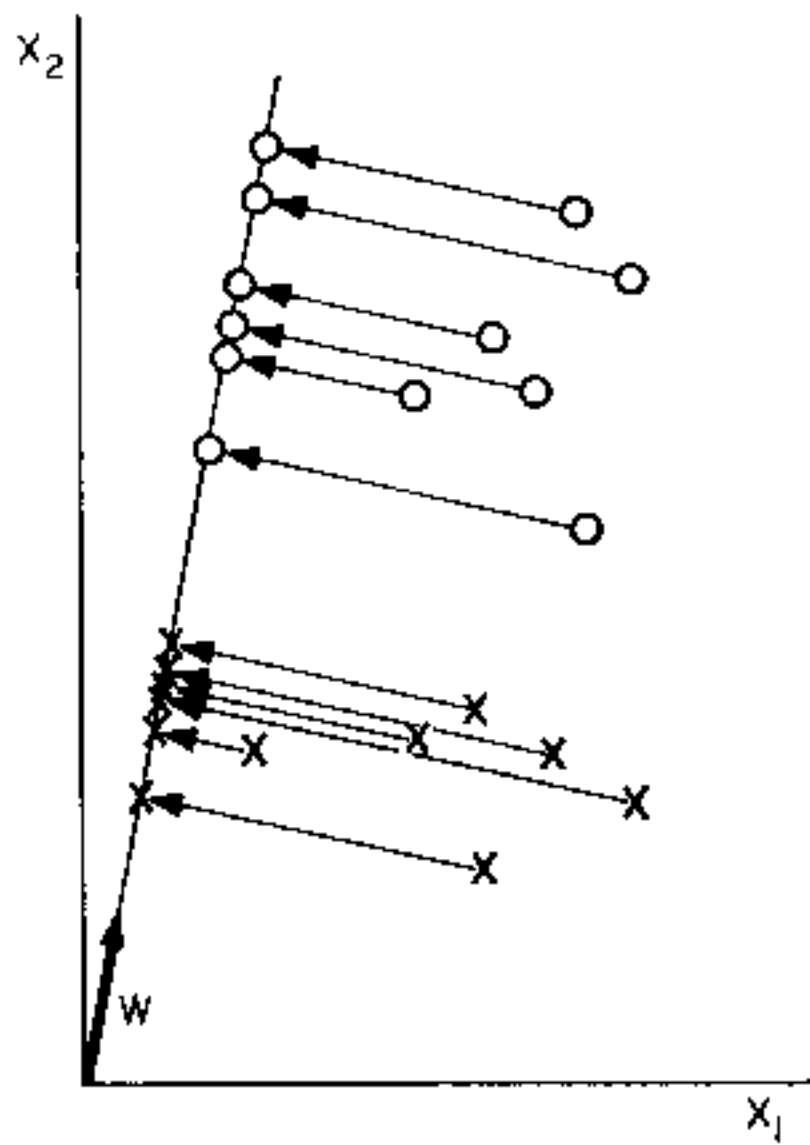
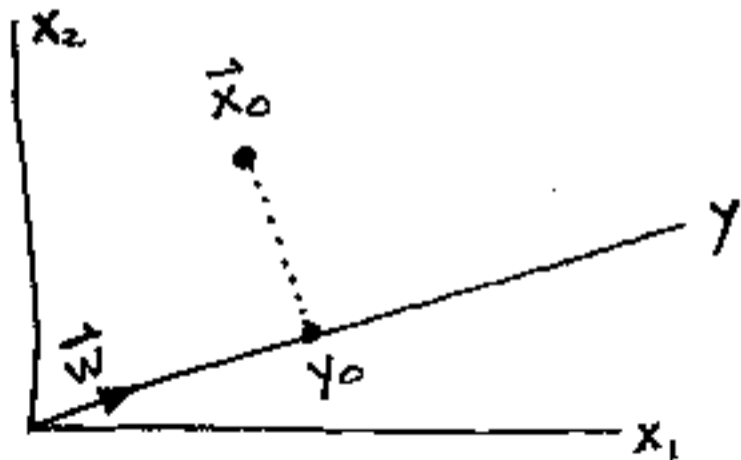


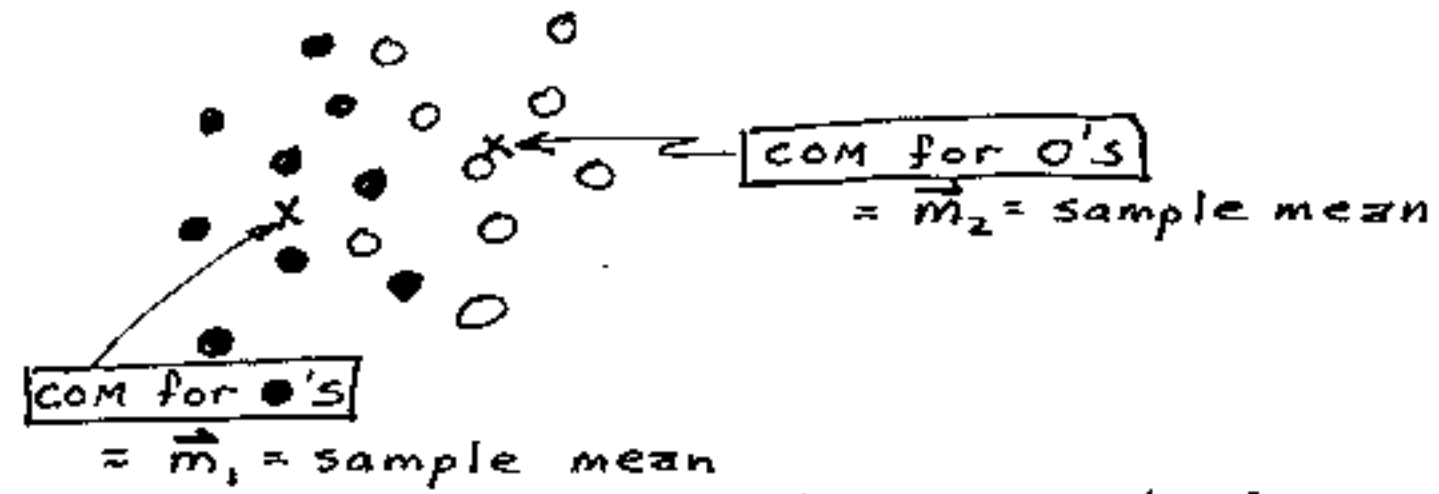
FIGURE 4.6. Projection of samples onto a line.

To project a vector onto \vec{w} , result is

$$y_0 = \vec{w}^T \vec{x}_0 \quad ; \quad \|\vec{w}\| = 1$$



Center of mass of data



n_1 = # data points in class 1 (\mathcal{X}_1)
 n_2 = " " " " " 2 (\mathcal{X}_2)

$$\vec{m}_1 = \frac{1}{n_1} \sum_{\vec{x} \in \mathcal{X}_1} \vec{x} \quad ; \quad \vec{m}_2 = \frac{1}{n_2} \sum_{\vec{x} \in \mathcal{X}_2} \vec{x}$$

The projection onto some given \vec{w} is

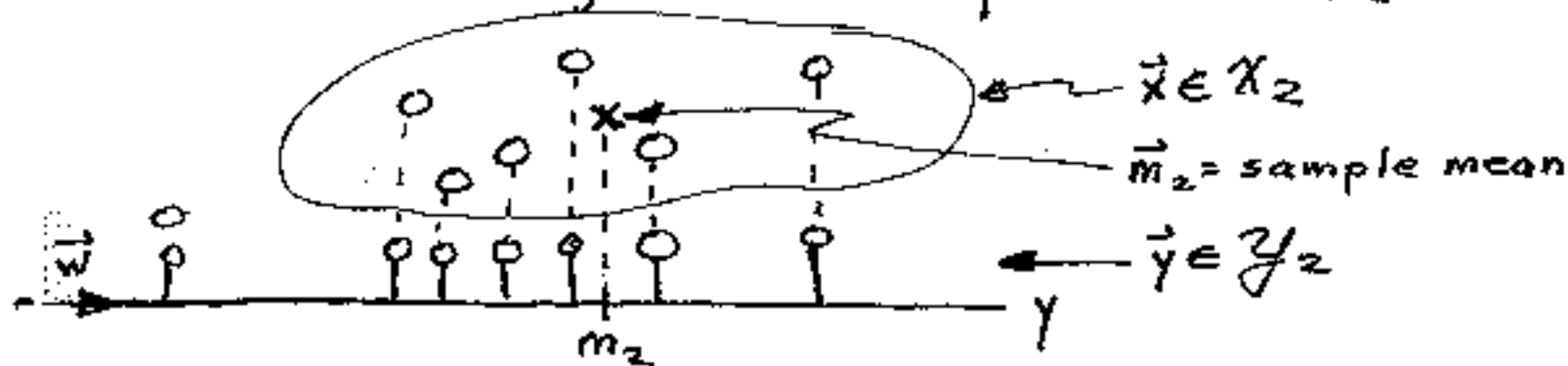
$$\tilde{m}_1 = \vec{w}^T \vec{m}_1 \quad ; \quad \tilde{m}_2 = \vec{w}^T \vec{m}_2$$

Project all the points for this \vec{w}

$$y = \vec{w}^T \vec{x}$$

If $\vec{x} \in \mathcal{X}_1 \Rightarrow y \in \mathcal{Y}_1$. If $\vec{x} \in \mathcal{X}_2 \Rightarrow y \in \mathcal{Y}_2$

Illustration: Projection of points in \mathcal{X}_2



This 1-D data has a ^{sample} mean and ^{sample} variance

$$\begin{aligned} \tilde{m}_2 &= \frac{1}{n_2} \sum_{y \in \mathcal{Y}_2} y = \frac{1}{n_2} \sum_{\vec{x} \in \mathcal{X}_2} \vec{w}^T \vec{x} \\ &= \vec{w}^T \left[\frac{1}{n_2} \sum_{\vec{x} \in \mathcal{X}_2} \vec{x} \right] \\ &= \vec{w}^T \vec{m}_2 \end{aligned}$$

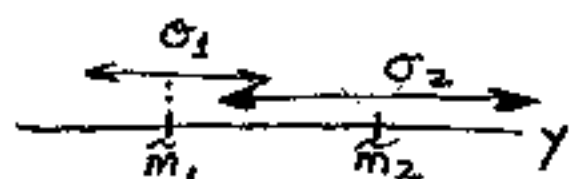
$$\sigma_2^2 = \frac{1}{n_2} \sum_{y \in \mathcal{Y}_2} (y - m_2^2); \quad (\text{some use } \frac{1}{n_2 - 1})$$

σ_2 measures the data spread



$\sigma_2 =$ _{sample} standard deviation

For a given \vec{w} , projected data looks like



These four parameters are functions of \vec{w} .

Q: How do we define "maximally separated data"?

A: Fisher's answer: Find the maximum value of

$$\frac{|m_1 - m_2|^2}{\sigma_1^2 + \sigma_2^2}$$

A small alteration. Define the "scatter" in each class as

$$\tilde{S}_1^2 = \sum_{y \in \mathcal{Y}_1} (y - m_1)^2 = n_1 \sigma_1^2$$

$$\tilde{S}_2^2 = \sum_{y \in \mathcal{Y}_2} (y - m_2)^2 = n_2 \sigma_2^2$$

The "Fisher distance" between the data sets is

$$J(\vec{w}) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{S_1^2 + S_2^2} \quad (\text{FD})$$

● Problem: Find the \vec{w} that maximizes $J(\vec{w})$

● Answer: (To be derived)

$$\vec{w} = \underline{S}_w^{-1} (\vec{m}_1 - \vec{m}_2) \quad (\text{FV})$$

where

$$\underline{S}_w = \underline{S}_1 + \underline{S}_2$$

and

$$\underline{S}_1 = \sum_{\vec{x} \in \mathcal{X}_1} (\vec{x} - \vec{m}_1)(\vec{x} - \vec{m}_1)^T; \quad \underline{S}_2 = \sum_{\vec{x} \in \mathcal{X}_2} (\vec{x} - \vec{m}_2)(\vec{x} - \vec{m}_2)^T$$

Proof:

$$\begin{aligned}
 \tilde{S}_1^2 &= \sum_{\vec{x} \in \mathcal{X}_1} (y - \tilde{m}_1)^2 \\
 &= \sum_{\vec{x} \in \mathcal{X}_1} (\vec{w}^T \vec{x} - \vec{w}^T \tilde{m}_1)^2 \\
 &= \sum_{\vec{x} \in \mathcal{X}_1} [\vec{w}^T (\vec{x} - \tilde{m}_1)]^2 \\
 &= \sum_{\vec{x} \in \mathcal{X}_1} \vec{w}^T (\vec{x} - \tilde{m}_1) (\vec{x} - \tilde{m}_1)^T \vec{w}
 \end{aligned}$$

$$\begin{aligned}
 &= \vec{w}^T \left[\sum_{\vec{x} \in \mathcal{X}_1} (\vec{x} - \tilde{m}_1) (\vec{x} - \tilde{m}_1)^T \right] \vec{w} \\
 &= \vec{w}^T \underline{S}_1 \vec{w} \quad ; \quad \underline{S}_1 = \sum_{\vec{x} \in \mathcal{X}_1} (\vec{x} - \tilde{m}_1) (\vec{x} - \tilde{m}_1)^T
 \end{aligned}$$

Scatter matrix

Similarly

$$\tilde{S}_2^2 = \vec{w}^T \underline{S}_2 \vec{w} \quad ; \quad \underline{S}_2 = \sum_{\vec{x} \in \mathcal{X}_2} (\vec{x} - \tilde{m}_2) (\vec{x} - \tilde{m}_2)^T$$

$$\begin{aligned}
 \text{Thus: } \tilde{S}_1^2 + \tilde{S}_2^2 &= \vec{w}^T \underline{S}_1 \vec{w} + \vec{w}^T \underline{S}_2 \vec{w} \\
 &= \vec{w}^T \underline{S}_w \vec{w} \quad ; \quad \text{(Den)}
 \end{aligned}$$

$$\underline{S}_w = \underline{S}_1 + \underline{S}_2 = \text{"within class scatter matrix"} \quad \text{(WI)}$$

Also:

$$\begin{aligned}
 |\tilde{m}_1 - \tilde{m}_2|^2 &= |\vec{w}^T (\tilde{m}_1 - \tilde{m}_2)|^2 \\
 &= \vec{w}^T (\tilde{m}_1 - \tilde{m}_2) (\tilde{m}_1 - \tilde{m}_2)^T \vec{w} \\
 &= \vec{w}^T \underline{S}_B \vec{w} \quad \text{(Num)}
 \end{aligned}$$

where

$$\underline{S}_B = (\tilde{m}_1 - \tilde{m}_2) (\tilde{m}_1 - \tilde{m}_2)^T = \text{between-class scatter matrix (BE)}$$

The Fisher distance (from (FD) on p.5) is

$$J(\vec{w}) = \frac{|\vec{m}_1 - \vec{m}_2|^2}{S_1^2 + S_2^2}$$

Substitute (Den) and (Num) from p.6):

$$J(\vec{w}) = \frac{\vec{w}^T \underline{S}_B \vec{w}}{\vec{w}^T \underline{S}_W \vec{w}} \quad (FD)$$

Note: Maximizing $J(\vec{w})$ minimizes within class scatter ($\vec{w}^T \underline{S}_W \vec{w}$) and maximized between class scatter ($\vec{w}^T \underline{S}_B \vec{w}$).

Note: Both numerator & den in (FD) above are scalars

Strategy to maximize $J(\vec{w})$: Take first derivative & set to zero.

$$\nabla_{\vec{w}} J(\vec{w}) = \frac{[\nabla_{\vec{w}} \vec{w}^T \underline{S}_B \vec{w}] \vec{w}^T \underline{S}_W \vec{w} - [\nabla_{\vec{w}} \vec{w}^T \underline{S}_W \vec{w}] \vec{w}^T \underline{S}_B \vec{w}}{(\vec{w}^T \underline{S}_W \vec{w})^2}$$

Using (Der) on p.1 \Rightarrow
$$= \frac{2 \underline{S}_B \vec{w} \vec{w}^T \underline{S}_W \vec{w} - 2 \underline{S}_W \vec{w} \vec{w}^T \underline{S}_B \vec{w}}{(\vec{w}^T \underline{S}_W \vec{w})^2}$$

Set $\nabla_{\vec{w}} J(\vec{w}) = \vec{0} =$ zero vector. This gives

$$\underline{S}_B \vec{w} \vec{w}^T \underline{S}_W \vec{w} = \underline{S}_W \vec{w} \vec{w}^T \underline{S}_B \vec{w}$$

(Note: Since \underline{S}_W is symmetric, $\underline{S}_W^{-1} = (\underline{S}_W^{-1})^T$).

This is an equality when (FV) on p.5 is true. Then

$$\vec{w} \vec{w}^T = \underline{S}_W^{-1} (\vec{m}_1 - \vec{m}_2)(\vec{m}_1 - \vec{m}_2) \underline{S}_W^{-1} = \underline{S}_W^{-1} \underline{S}_B \underline{S}_W^{-1}$$

Substituting:

$$\underline{S}_B (\underline{S}_W^{-1} \underline{S}_B \underline{S}_W^{-1}) \underline{S}_W \vec{w} = \underline{S}_W (\underline{S}_W^{-1} \underline{S}_B \underline{S}_W^{-1}) \underline{S}_B \vec{w}$$
$$\Rightarrow \underline{S}_B \underline{S}_W^{-1} \underline{S}_B \vec{w} = \underline{S}_B \underline{S}_W^{-1} \underline{S}_B \vec{w}$$

and the proof is done!

Summary: Fisher discriminant

- Find Fisher vector, \vec{w} , to maximize

$$J(\vec{w}) = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{S_1^2 + S_2^2}$$

- Solution is Fisher vector

$$\Rightarrow \vec{w} = \underline{S}_w^{-1} (\vec{m}_1 - \vec{m}_2)$$

where

- $\underline{S}_w = \underline{S}_1 + \underline{S}_2$ \Leftarrow within class scatter matrix

$$\underline{S}_1 = \sum_{\vec{x} \in \mathcal{X}_1} (\vec{x} - \vec{m}_1)(\vec{x} - \vec{m}_1)^T; \quad \underline{S}_2 = \sum_{\vec{x} \in \mathcal{X}_2} (\vec{x} - \vec{m}_2)(\vec{x} - \vec{m}_2)^T$$

$$\vec{m}_1 = \frac{1}{n_1} \sum_{\vec{x} \in \mathcal{X}_1} \vec{x}; \quad \vec{m}_2 = \frac{1}{n_2} \sum_{\vec{x} \in \mathcal{X}_2} \vec{x}$$

Assumptions:

1. \underline{S}_w is not singular.

● Necessary condition: $n_1 + n_2 \geq d$ (elaborate)

More generally, \underline{S}_w must be "well conditioned"

The condition of a matrix, \underline{A} , is

$$C(\underline{A}) = \frac{|\lambda_{\max}|}{|\lambda_{\min}|}$$

For singular matrices, $C(\underline{A}) = \infty$.

● Necessary condition: All data can not lie on a plane with dimension $< d$.

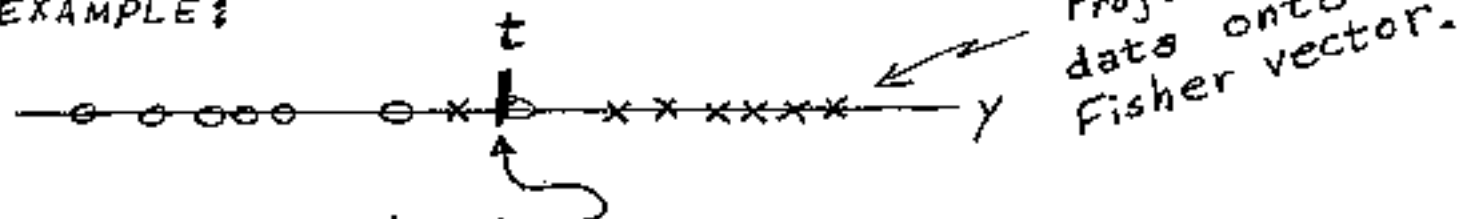
(Equivalent to non-singular requirement)

2. The data classes are linear separable.

USING THE FISHER VECTOR FOR CLASSIFICATION (9)

Q: After projection onto the Fisher vector, how do we classify?

EXAMPLE:



A: Set a threshold, t .

Above t , we have an "x", below an "o."

Q: Where do we place t ?

A: Depends on your desired performance.

Example: x = thermonuclear attack

o = no attack

Couched in classical hypothesis testing:

H_0 : $\square = o$ = no attack

H_1 : $\square = x$ = attack

where \square is the observation.

$\Rightarrow \alpha$ = FALSE ALARM PROBABILITY

$$= \Pr[H_1 \mid \square = o]$$

= Prob [Announcing H_1 GIVEN OBSERVATION IS "NO ATTACK"]

$\Rightarrow \beta$ = DETECTION PROBABILITY

$$= \Pr[H_1 \mid \square = x]$$

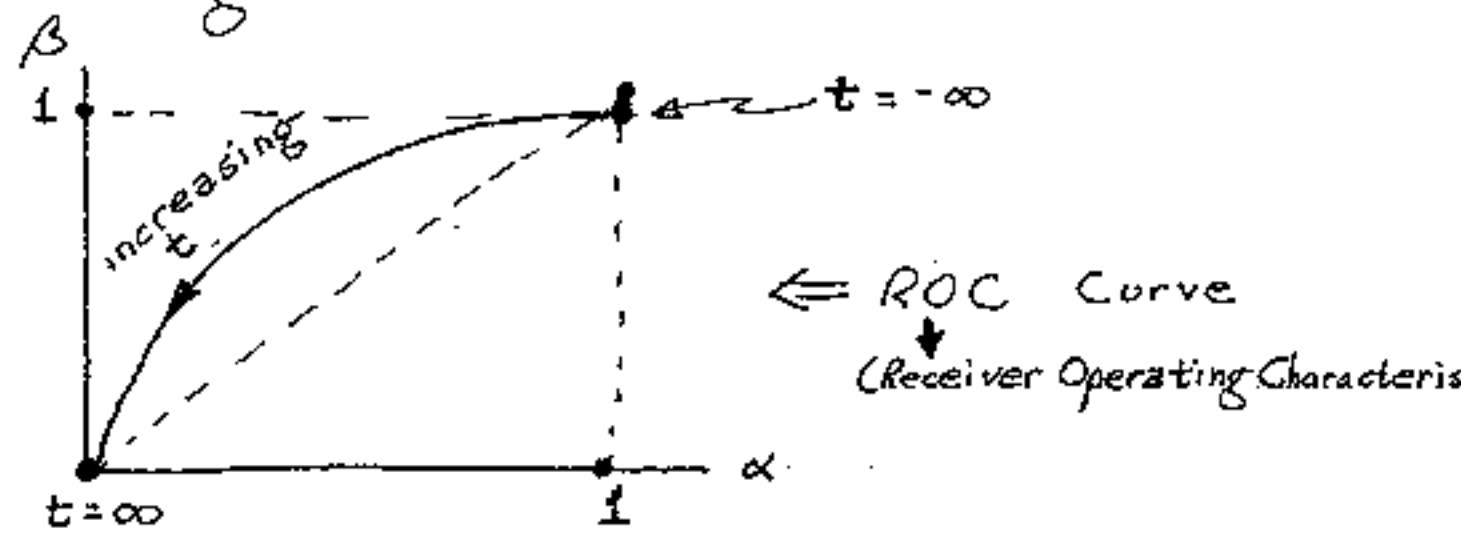
= Prob [Announcing H_1 GIVEN OBSERVATION IS "ATTACK"]

To get $\beta = 1$, set $t = -\infty$. But $\alpha = 1$

To get $\alpha = 0$, set $t = \infty$. But $\beta = 0$

Always "no attack"

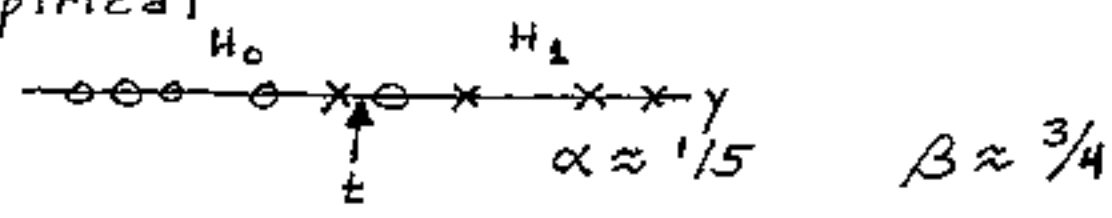
Increasing β increases α
Decreasing α decreases β



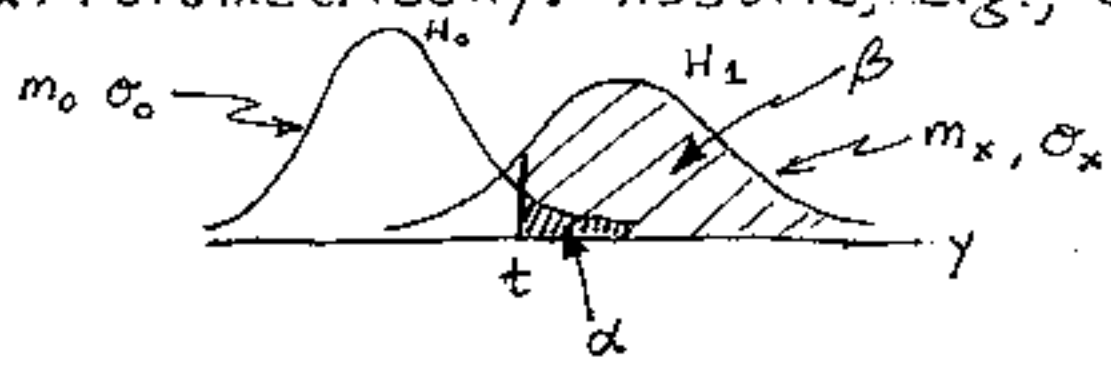
Choosing t specifies your point on the ROC curve.

Q: How is α & β estimated?

A1: Empirical



A2: Parametrically. Assume, e.g., Gaussian



Mahalanobis Distance Classifier

(11)

Multivariate Normal Density

$$f_{\vec{x}}(\vec{x}) = \frac{1}{(2\pi)^{d/2} |\underline{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\vec{x} - \vec{m})^T \underline{\Sigma}^{-1} (\vec{x} - \vec{m})}$$

$$\vec{m} = E[\vec{X}] \approx \frac{1}{n} \sum_m \vec{X}_m \leftarrow \text{AVERAGE}$$

$$\sigma_{ij} = E[(X_i - m_i)(X_j - m_j)] = (\underline{\Sigma})_{ij}$$

$$\underline{\Sigma} = E[(\vec{X} - \vec{m})(\vec{X} - \vec{m})^T]$$

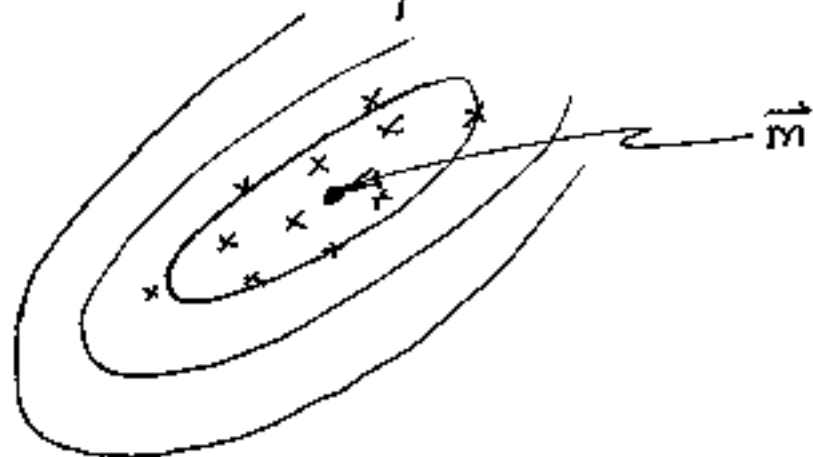
$$\approx \frac{1}{n} \sum_m (\vec{X}_m - \vec{m})(\vec{X}_m - \vec{m})^T$$

Q: What do contours of $f_{\vec{x}}(\vec{x})$ look like?

A: Same as

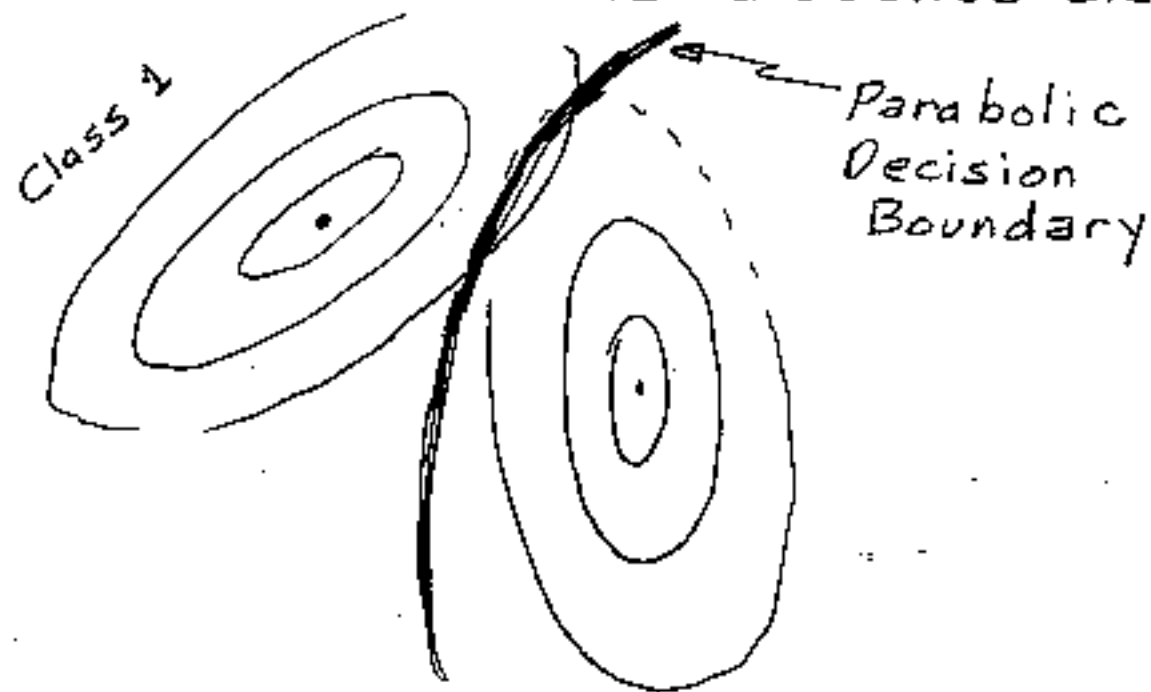
$$r^2 = (\vec{x} - \vec{m})^T \underline{\Sigma}^{-1} (\vec{x} - \vec{m})$$

In 2-D, it's ellipse centered at \vec{m} .



r = Mahalanobis distance

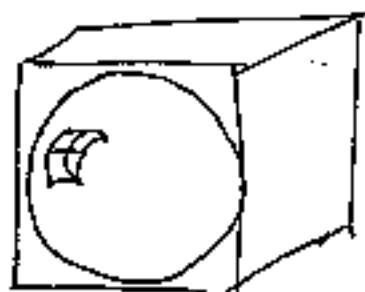
Minimum Mahalanobis distance classifier



Can extend to 3 or more classes.



← To distinguish circle, we need, say, 100 data points in 2-d.



← In 3-d, need 1000 pts

In d dimensions $\Rightarrow 10^d \Leftarrow$ wow!

Thus: If there are a lot of features, one might one to choose the "best" subset.

But, if there are n features, and we want to choose the best k , an exhaustive search would require

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \quad \text{evaluations}$$

Alternate idea: Choose the best k features

But: The k best features are not the best k features!

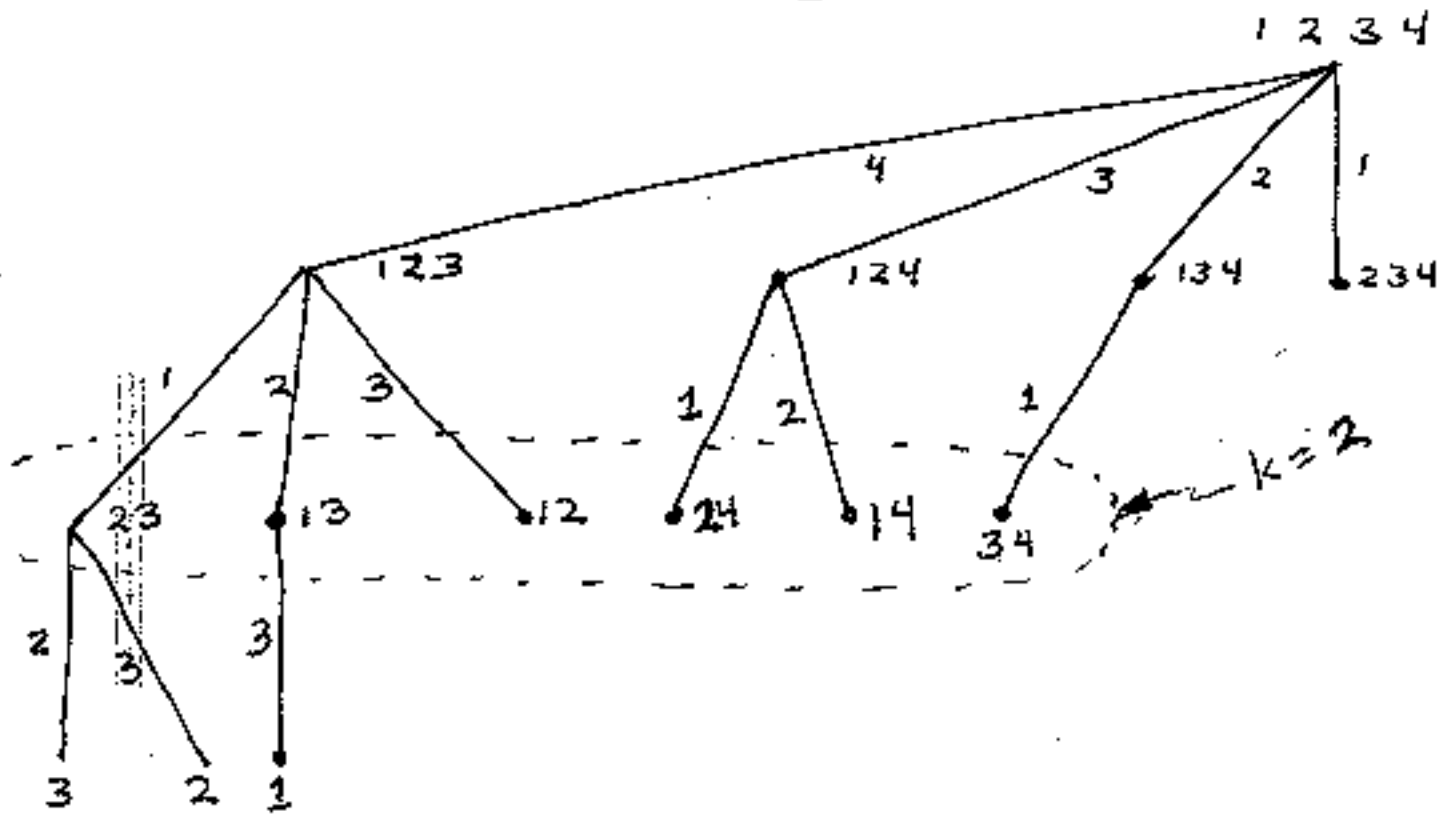
(Elaborate)

Suboptimal Choice: Sequential Selection

1. Choose best feature, a
 2. " feature that goes best with a (call it b).
 3. Choose feature c that goes best with a and b .
- etc.

Note: Increasing by a feature increases Fisher distance.

Dynamic Programming for Optimality (15)



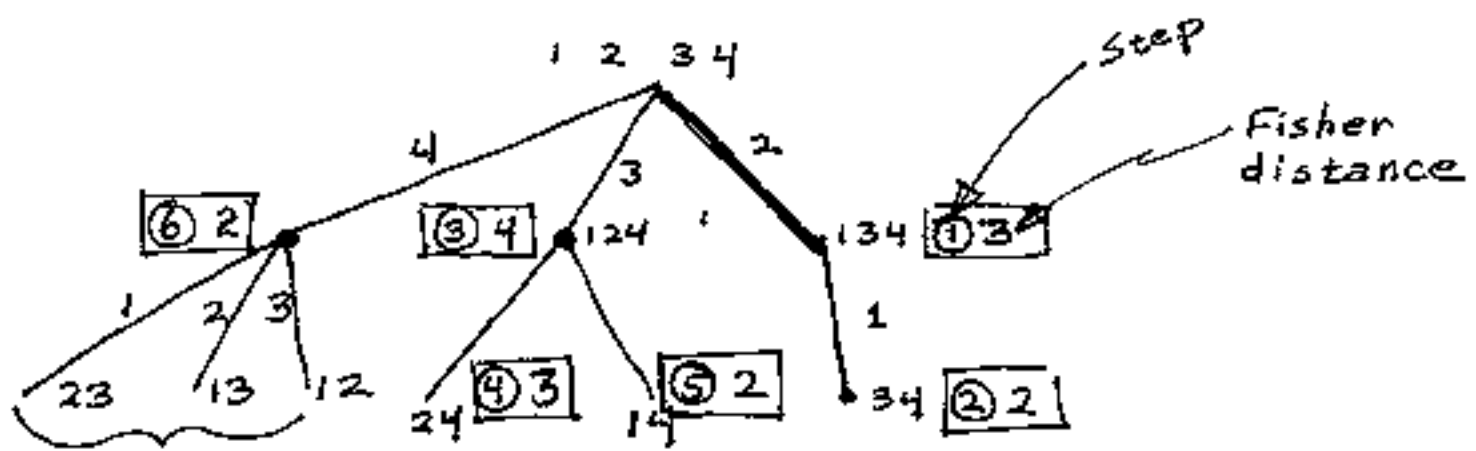
$$\binom{4}{0}$$

$$\binom{4}{1}$$

$$\binom{4}{2} = 6$$

$$\binom{4}{3} = 4$$

Choose two best features

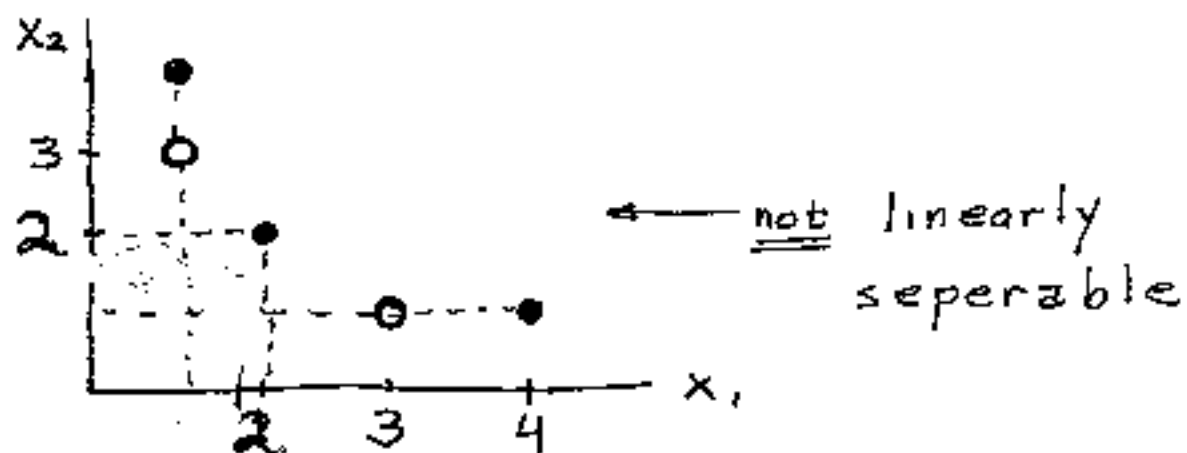


Don't need to search these

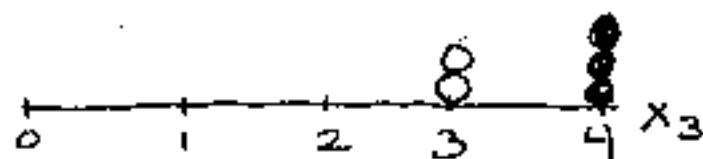
Functional Feature Alteration

(16)

Ex



Let $X_3 = X_1 \cdot X_2$



X_1, X_2 are sides of rectangle
 X_3 is area.

PRINCIPAL COMPONENTS / EIGEN DECOMPOSITION / K-L Analysis

$$\vec{y} = \sum \vec{x} \quad \leftarrow d \text{ dimensions}$$

The eigenvalues, λ_n , and eigenvectors, are defined by

$$\lambda_n \vec{\psi}_n = \underline{\Sigma} \vec{\psi}_n \quad ; \quad 1 \leq n \leq d$$

They have the following properties:

(a) Orthogonality : $\vec{\psi}_n^T \vec{\psi}_m = \delta[n-m] = \begin{cases} 1 & ; n=m \\ 0 & ; n \neq m \end{cases}$

(b) $\underline{\Sigma} = \sum_{n=1}^d \lambda_n \vec{\psi}_n \vec{\psi}_n^T$

Any vector can be expressed as

$$\vec{x} = \sum_{n=1}^d \alpha_n \vec{\psi}_n$$

where

$$\alpha_n = \vec{x}^T \vec{\psi}_n$$

Thus:
$$\begin{aligned} \vec{y} = \sum \vec{x} &= \sum \sum_{n=1}^d \alpha_n \vec{\psi}_n \\ &= \sum_{n=1}^d \alpha_n \sum \vec{\psi}_n \\ &= \sum_{n=1}^d \alpha_n \lambda_n \vec{\psi}_n \quad \Leftarrow \end{aligned}$$

PRINCIPAL COMPONENTS

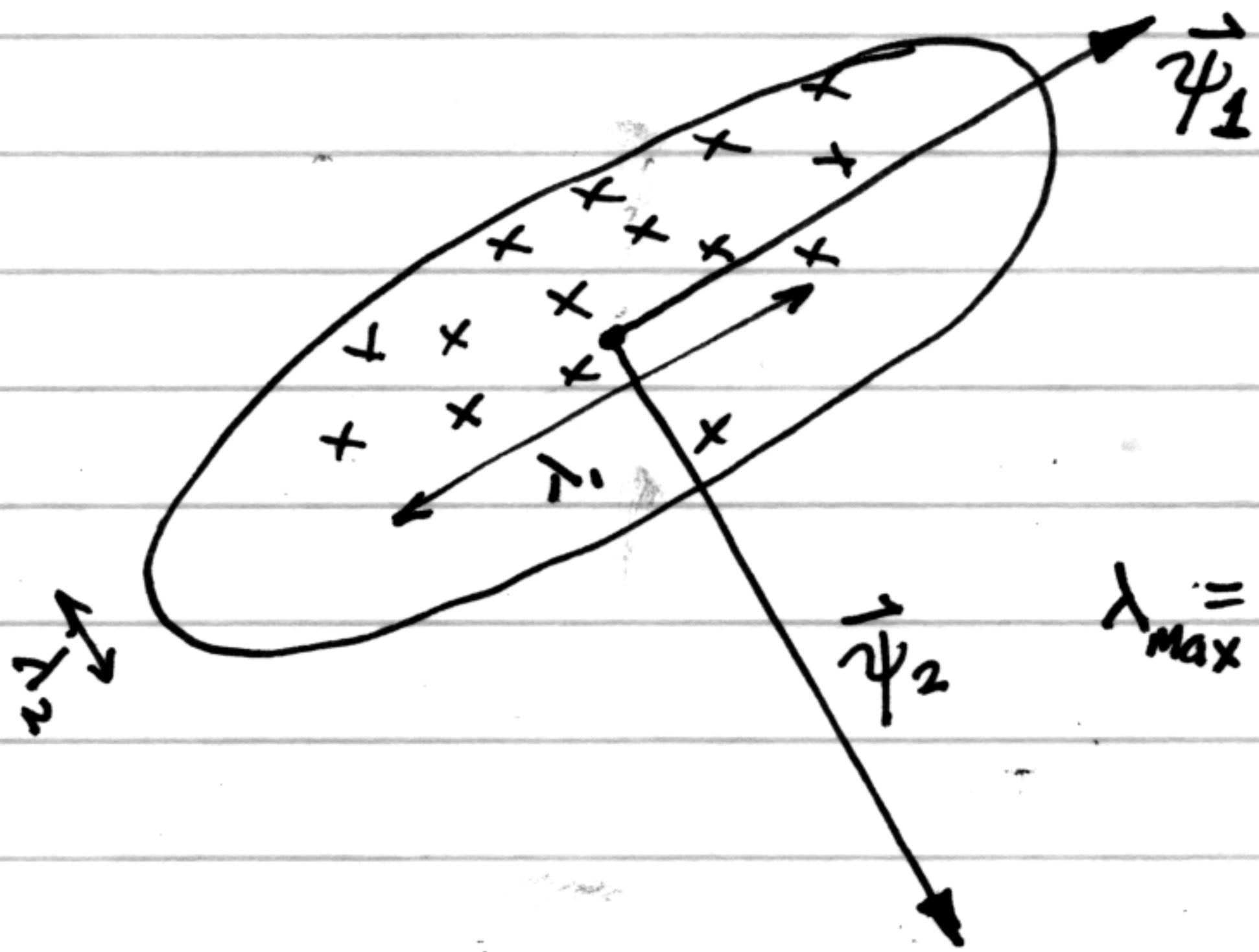
Idea: Many of λ_n 's are such that $|\lambda_n| \ll |\lambda_{\max}|$

Interpretation: If $\underline{\Sigma}$ is estimated by data:

$$\underline{\Sigma} \approx \frac{1}{d-1} \sum_{n=1}^d (\vec{X}_n - \vec{m})(\vec{X}_n - \vec{m})^T$$

$$\vec{m} \approx \frac{1}{d} \sum_{n=1}^d \vec{X}_n$$

Relation of eigenvalues to ellipses



Note: $\underline{\Sigma}$ is positive semi-definite, so

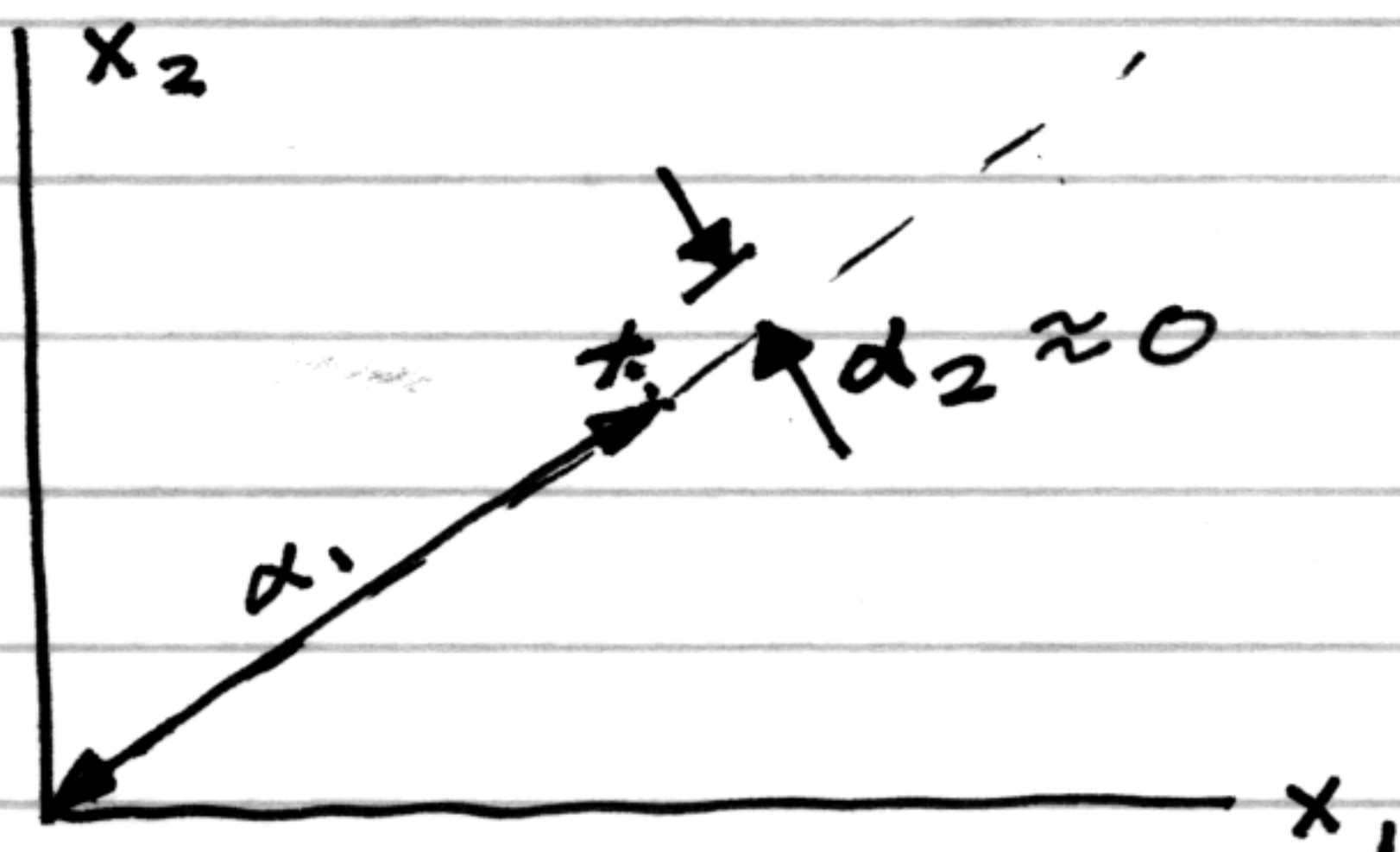
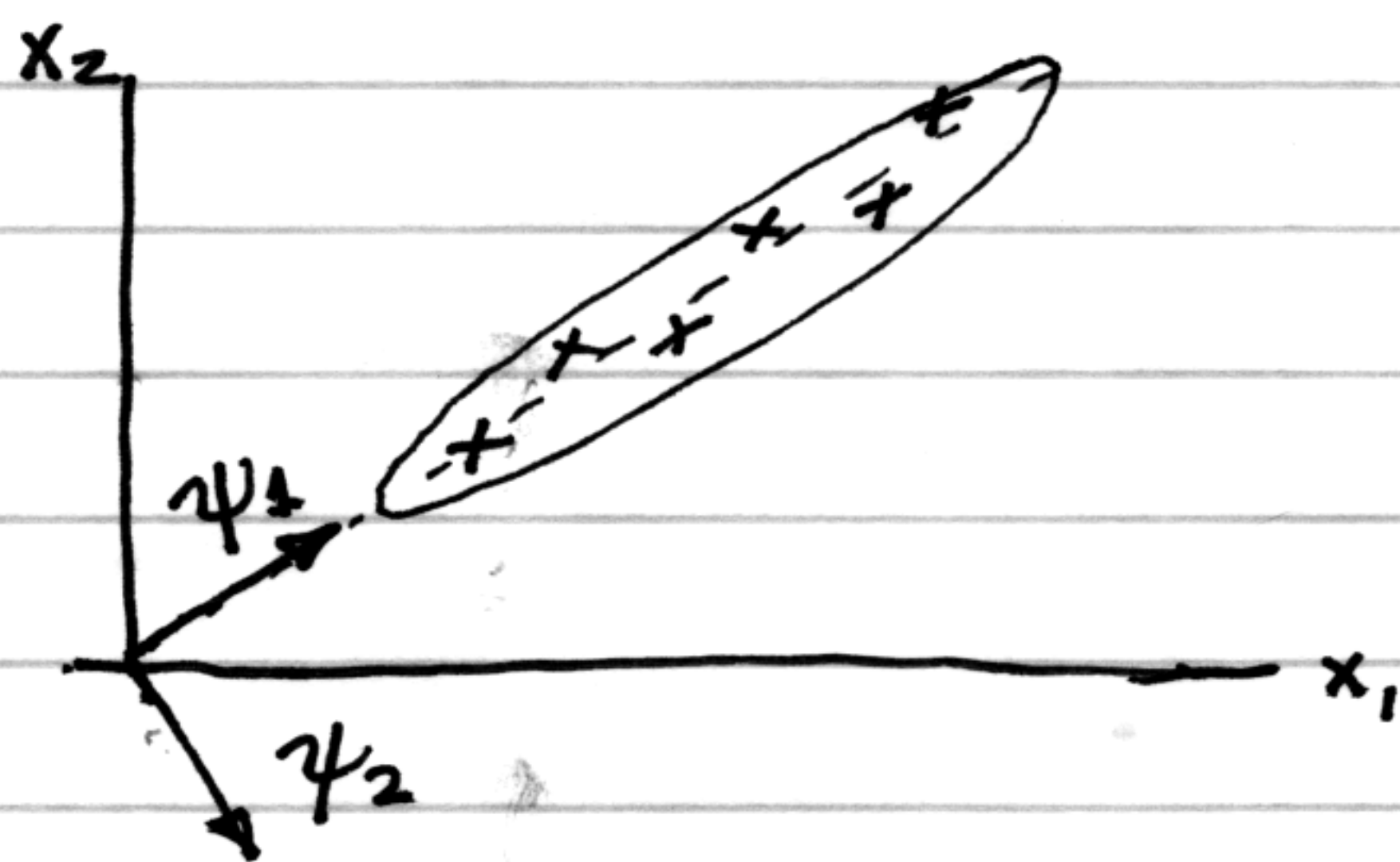
$$\lambda_{\max} = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$$

If ellipsoid is "flat" in certain dimensions, the corresponding λ_n 's are "small".
Then, for all of the \vec{x} 's in this flat ellipse,

$$\alpha_n \approx 0$$

Example:

(19)



Thus, if λ_n 's are "small" for $n > \hat{d}$,

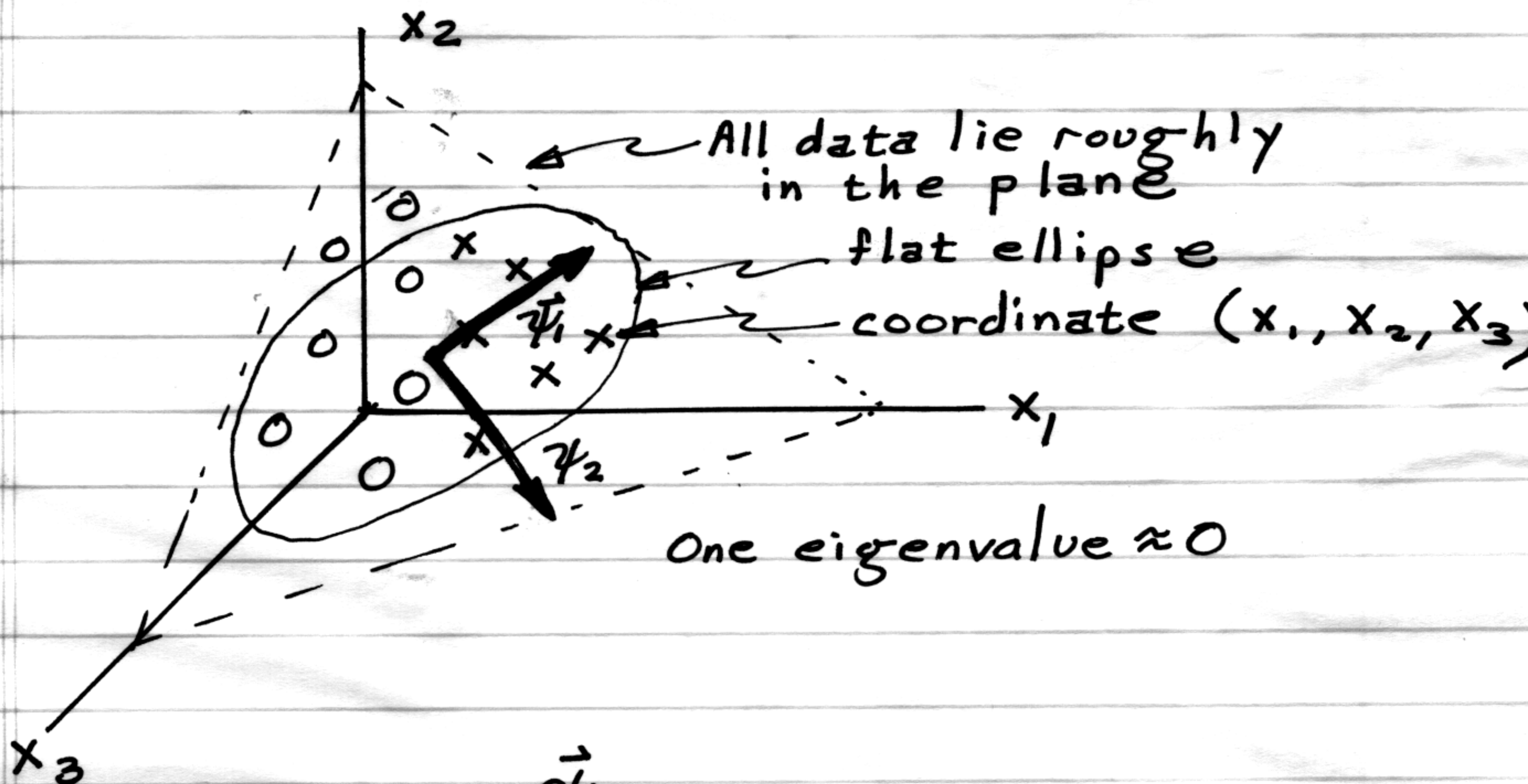
$$\vec{x} = \sum_{n=1}^d \alpha_n \vec{\psi}_n \approx \sum_{n=1}^{\hat{d}} \alpha_n \vec{\psi}_n$$

The \hat{d} numbers, $\{\alpha_n \mid 1 \leq n \leq \hat{d}\}$ have about the same information as $\{x_m \mid 1 \leq m \leq d\}$.

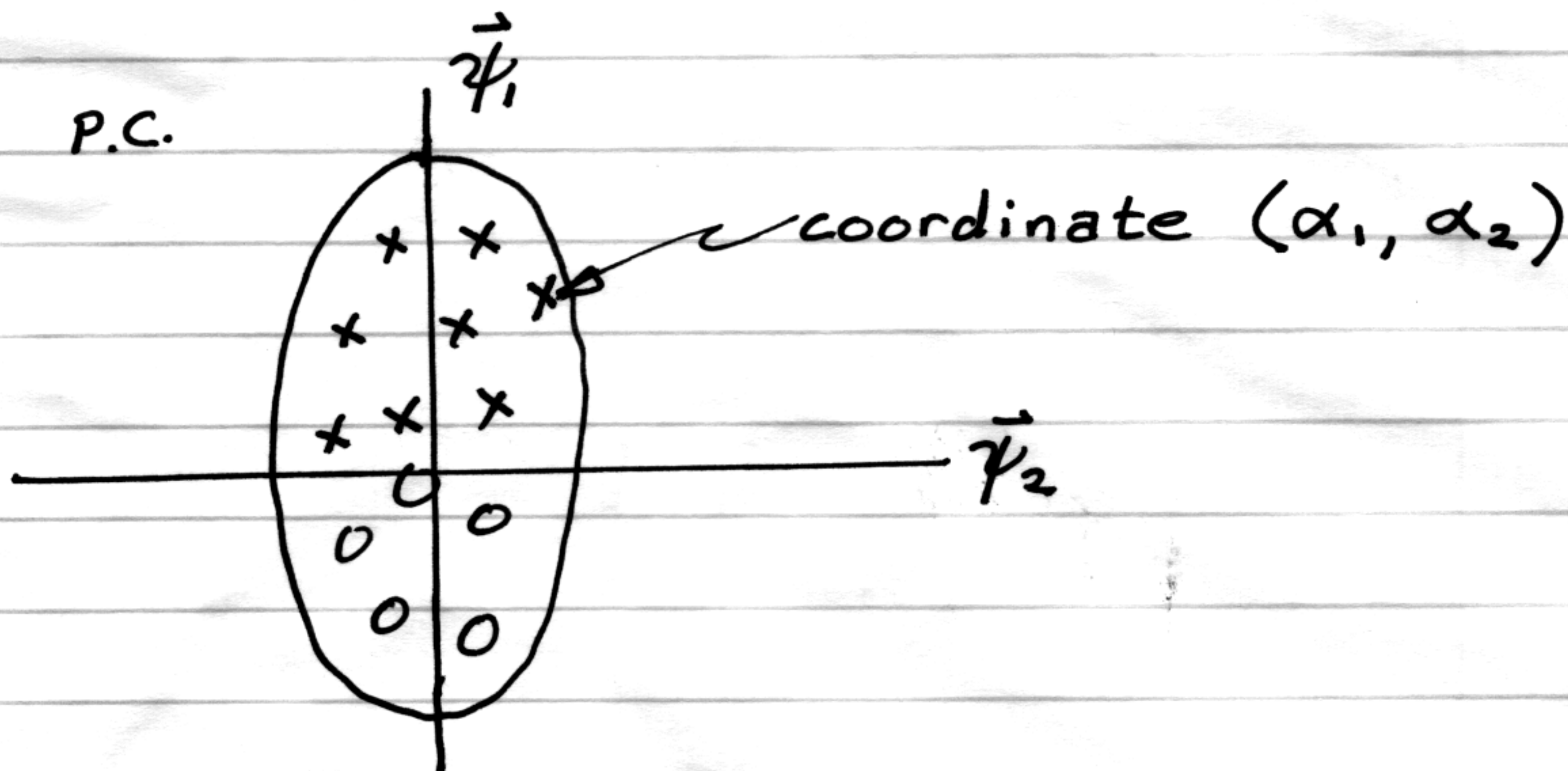
We have reduced the dimension by $d - \hat{d} \Rightarrow$ Good news against the curse.

Use $\vec{\alpha}$ vectors in place of \vec{x} vectors when training.

Geometrical Interpretation



After P.C.



Note:

- Principal Components (a.k.a. Karhunen-Loeve) Analysis can be used for
FEATURE EXTRACTION

- Choosing the k best features is
FEATURE SELECTION

→ Both reduce dimensions, but are different.

1. All features are used in PC

2. The classification label is not used in P.C.

(22)

Use in data visualization.

Problem: We have d -dimensional data.

We desire to view it projected onto a two-dimensional plane so that the data is "spread out" the most.

Solution: Use the two principal components. Plot the data in the resulting two dimensional plane.