

Machine Learning I: Supervised Methods

B. Keith Jenkins

Announcements

- Slido event code: 4509123
- Midterm grading is almost completed; graded midterms will be handed out on Wednesday
- Homework 6 is due Friday
- Project assignment will be posted later this week.

Today's lecture

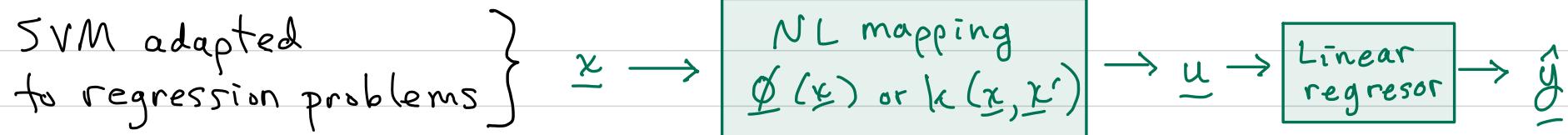
- Support vector regression (finish)
- Feature selection and dimensionality reduction (part 2)
 - Fisher's linear discriminant (FLD)
 - Criterion function and solution
 - Extensions to $d > 1$ and $C > 2$
 - Multiple classifiers
 - MDA
- Artificial neural networks (ANN) (part 1)
 - Single neuron unit
 - Activation functions

deferred

Reading

- Bishop 12.1.1 (PCA), 12.1.3 (PCA applications - optional)
- Bishop, 4.1.4 (FLD), 4.1.6 (MDA)

Support Vector Regression (SVR) [Bishop 7.1.4] (non-augmented notation)



Motivation: to have a regression model suitable to sparse data (e.g., high D' , relatively low N).

In Ridge Regression, we used a criterion fcn.: $\text{MSE} + \ell_2^2$ regularizer:

$$J_{RR}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N [\hat{y}(x_n) - y_n]^2 + \lambda \|\underline{w}\|_2^2, \lambda \geq 0$$

{ } { } { }

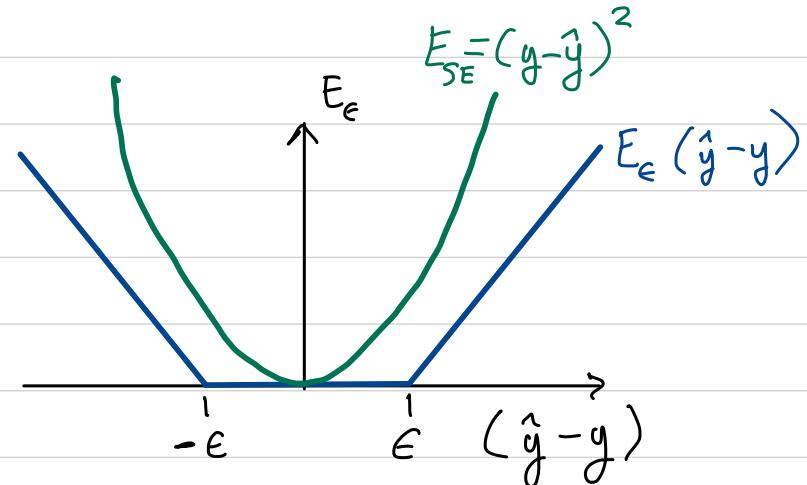
For SVR:

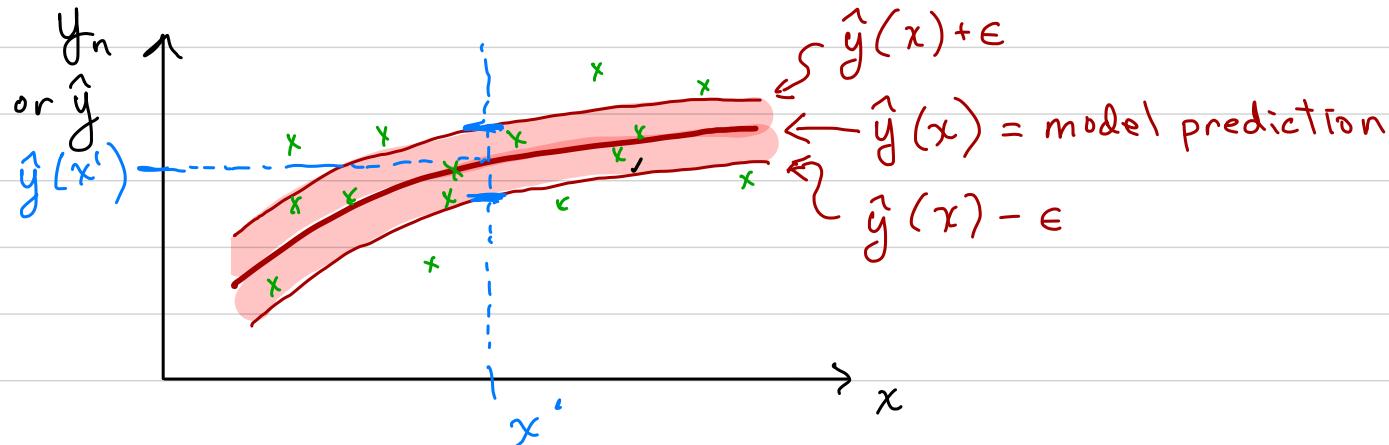
use " ϵ -insensitive" loss
instead of squared-error loss

keep ℓ_2^2 regularizer

ϵ -insensitive loss E_ϵ :

$$E_\epsilon(\hat{y} - y) = \begin{cases} 0, & \text{if } |\hat{y} - y| < \epsilon \\ |\hat{y} - y| - \epsilon & \text{if } |\hat{y} - y| \geq \epsilon \end{cases}$$





Any data points
in the shaded region
have $E_\epsilon = 0$.

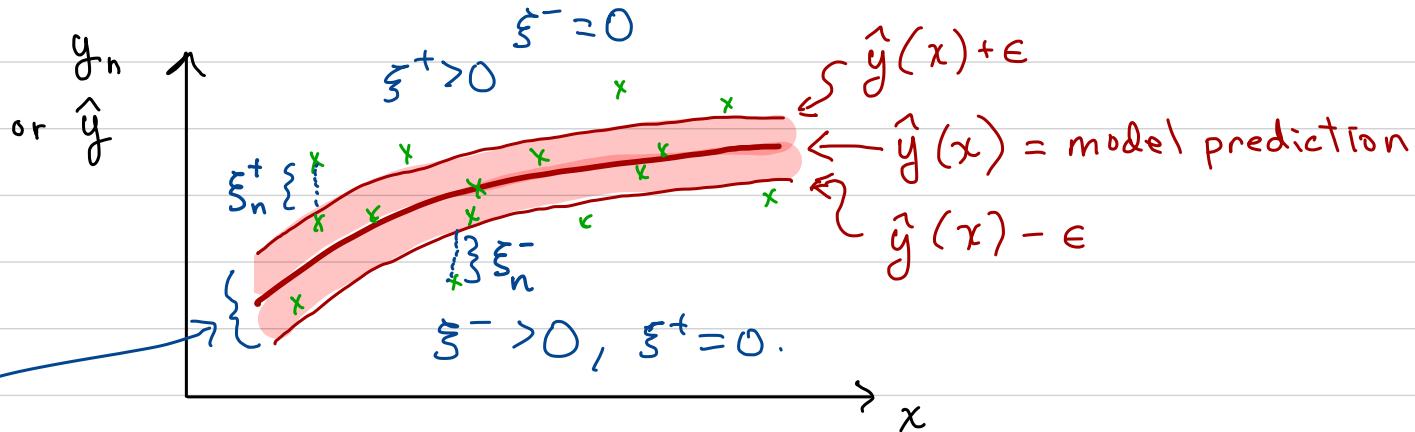
- 0 error for predictions $\hat{y}(x_n)$ within ϵ of known output y_n , means learning algorithm is less likely to overfit to noise or small variations in data.

⇒ Use error function (before adding constraints):

$$\tilde{J}(\underline{w}) = C \sum_{n=1}^N E_\epsilon [\hat{y}(x_n) - y_n] + \frac{1}{2} \|\underline{w}\|_2^2$$

\nwarrow C here by convention.

To pose as a tractable constrained optimization problem, use slack variables and constraints in place of E_ϵ :



|(1) If: $\{y_n \leq \hat{y}(x_n) + \epsilon \text{ and } y_n \geq \hat{y}(x_n) - \epsilon\}$ then: y_n is inside the tube

| and $\xi_n^+ = \xi_n^- = 0$ (where $\epsilon_\epsilon = 0$) ↔ Use as set of constraints?

$$\xi_n^+ > 0 \text{ iff } y_n > \hat{y}(x_n) + \epsilon; \quad \xi_n^- > 0 \text{ iff } y_n < \hat{y}(x_n) - \epsilon$$

Add slack variables:

$$\Rightarrow \text{we want } \{y_n \leq \hat{y}(x_n) + \epsilon + \xi_n^+ \text{ and } y_n \geq \hat{y}(x_n) - \epsilon - \xi_n^-\}$$

(5)

we can write: $\tilde{J}(\underline{w}) = C \sum_{n=1}^N (\xi_n^+ + \xi_n^-) + \frac{1}{2} \|\underline{w}\|_2^2$

subject to constraints:

(i) $\xi_n^+ \geq 0, \xi_n^- \geq 0 \quad \forall n$

ξ_n^+, ξ_n^-
prefers small

regularizer term

accommodates
points that
are outside
the tube,

(ii) $y_n \leq \hat{y}(x_n) + \epsilon + \xi_n^+ \quad \forall n$

(iii) $y_n \geq \hat{y}(x_n) - \epsilon - \xi_n^- \quad \forall n$

(2)

Yields a Lagrangian (primal form):

ϵ insensitive loss

regularizer

$$\begin{aligned}
 L(\underline{w}, w_0, \xi^+, \xi^-, \mu^+, \mu^-, \lambda^+, \lambda^-) &= C \sum_{n=1}^N (\xi_n^+ + \xi_n^-) + \frac{1}{2} \|\underline{w}\|_2^2 \\
 (3) \quad &- \sum_{n=1}^N (\mu_n^+ \xi_n^+ + \mu_n^- \xi_n^-) - \sum_{n=1}^N \lambda_n^+ (\epsilon + \xi_n^+ + \hat{y}_n - y_n) - \sum_{n=1}^N \lambda_n^- (\epsilon + \xi_n^- - \hat{y}_n + y_n) \\
 &\text{in which } \hat{y}_n = \hat{y}(\underline{x}_n) = \underline{w}^\top \phi(\underline{x}_n) + w_0
 \end{aligned}$$

(i)
(ii)
(iii)

(i) $\xi_n^+ \geq 0, \xi_n^- \geq 0 \quad \forall n$
 (ii) $y_n \leq \hat{y}(\underline{x}_n) + \epsilon + \xi_n^+ \quad \forall n$
 (iii) $y_n \geq \hat{y}(\underline{x}_n) - \epsilon - \xi_n^- \quad \forall n$
 (iv) $\mu_n^+, \mu_n^-, \lambda_n^+, \lambda_n^- \geq 0 \quad \forall n$

ϵ is user-specified (or obtained through model selection).

One can obtain the dual representation of L [see Bishop].

Support vectors are defined as the vectors that are outside the ϵ -tube or on the boundary of the ϵ -tube.

FEATURE SELECTION AND DIMENSIONALITY REDUCTION (part 2)

Why?

(S)

- Reduce overfitting
- Reduce computation
- Eliminate irrelevant features.

PCA - last lecture

From Lecture 16, p.18

How well does PCA work?

Ex: $C=2$ classes

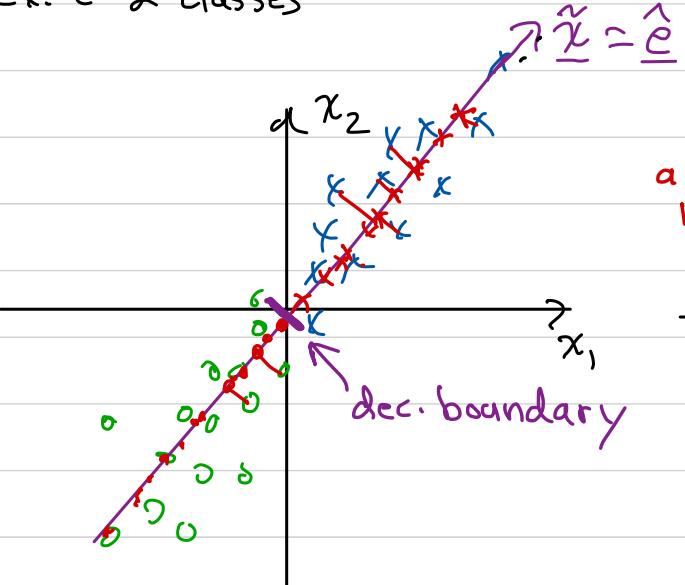


Fig. A

PCA should work well

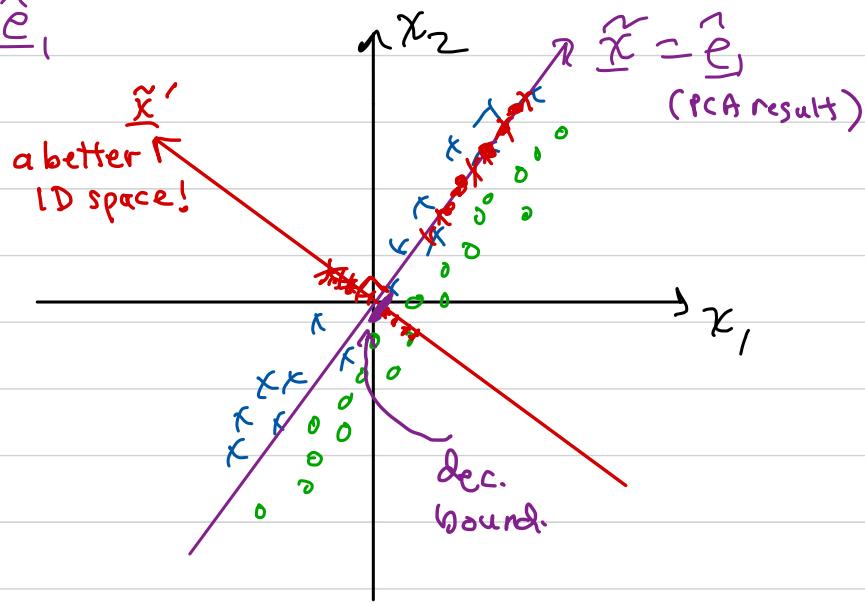


Fig. B

PCA works poorly.

but, using \tilde{x}' as 1D feature space could work well.

Why can't PCA do better on Fig. B?

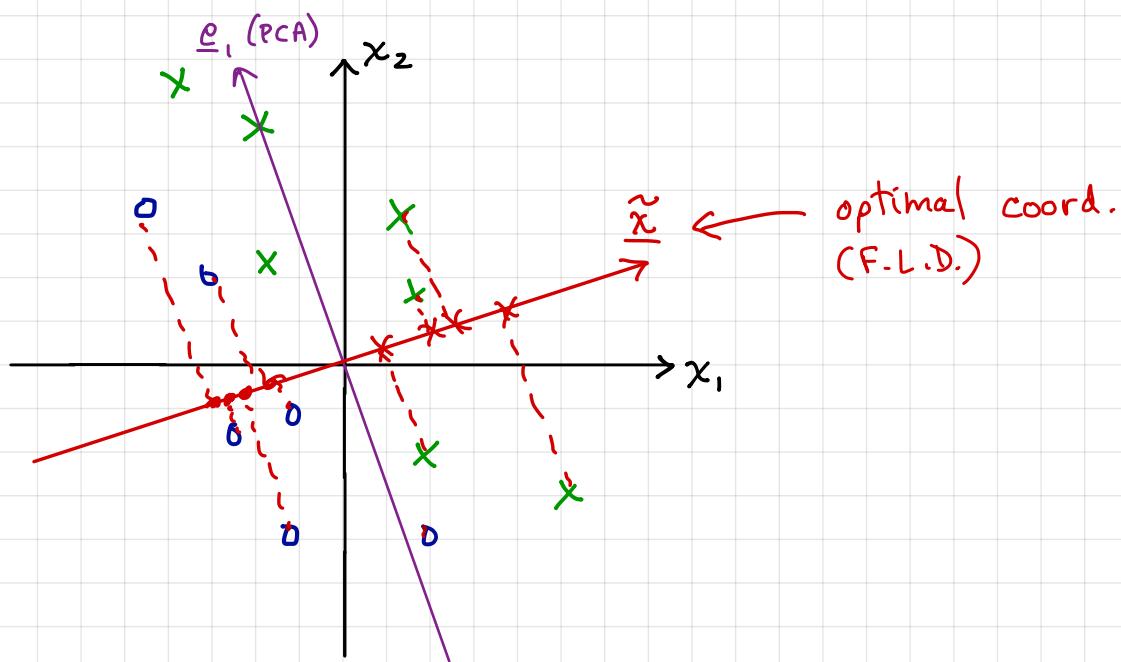
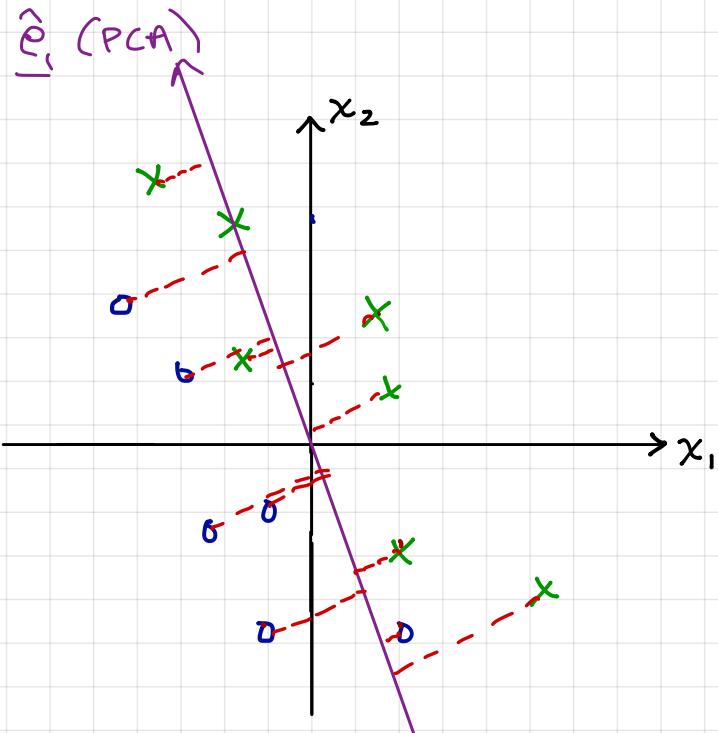
→ PCA ignores class labels

Fisher's Linear Discriminant (FLD) - Problem statement

[Bishop 4.1.4]

For classification problems.

$C=2$ classes. Dimensionality $D \rightarrow l = D'$.



Find direction of a line (1-D space) that maximizes the separability of projected data pts. (in 1-D feature space).

Criterion:

$$J = \frac{[\text{distance between projected class means}]}{[\text{some measure of variance of each projected class, combined}]}^2$$

→ Then find max of J w.r.t. direction of line (\underline{w}).

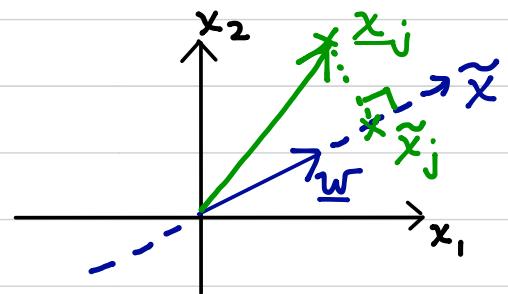
FISHER'S LINEAR DISCRIMINANT - Problem solution

\tilde{x} is new 1D space

$\underline{x}_j \rightarrow$ Linear transformation $\rightarrow \tilde{x}_j$

$$\tilde{x}_j = \underline{w}^T \underline{x}_j$$

$\|\underline{w}\| = 1$ gives pure projection



Name	Original feat. space	Reduced (1D) feat. space
Data point	\underline{x}_j	\tilde{x}_j

$$\text{Sample mean of } S_i \text{ pts.} \quad \underline{m}_i \quad \tilde{m}_i = \underline{w}^T \underline{m}_i$$

$$\text{Sample (co)variance of } S_i \text{ pts.} \quad \frac{1}{N_i} \underline{\underline{S}}_i \quad \frac{1}{N_i} \tilde{\underline{\underline{S}}}^2_i = \frac{1}{N_i} \sum_{\tilde{x}_j \in S_i} (\tilde{x}_j - \tilde{m}_i)^2$$

$\underbrace{\quad}_{= \frac{1}{N_i} \sum_{x_j \in S_i} (x_j - \underline{m}_i)(x_j - \underline{m}_i)^T}$

$$\text{Within-class scatter} \quad \underline{\underline{S}}_w = \underline{\underline{S}}_1 + \underline{\underline{S}}_2 \quad \tilde{\underline{\underline{S}}}^2_1 + \tilde{\underline{\underline{S}}}^2_2$$

$$\text{Between-class scatter} \quad \underline{\underline{S}}_B = (\underline{m}_2 - \underline{m}_1)(\underline{m}_2 - \underline{m}_1)^T \quad |\tilde{m}_2 - \tilde{m}_1|^2$$

Def Fisher's Linear Discriminant is the linear function $\underline{w}^T \underline{x}$ for which the criterion fcn.:

$$(1) \quad J(\underline{w}) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{\underline{\underline{S}}}^2_1 + \tilde{\underline{\underline{S}}}^2_2} \quad \text{is maximized.}$$

$J(\underline{w})$ can be maximized by

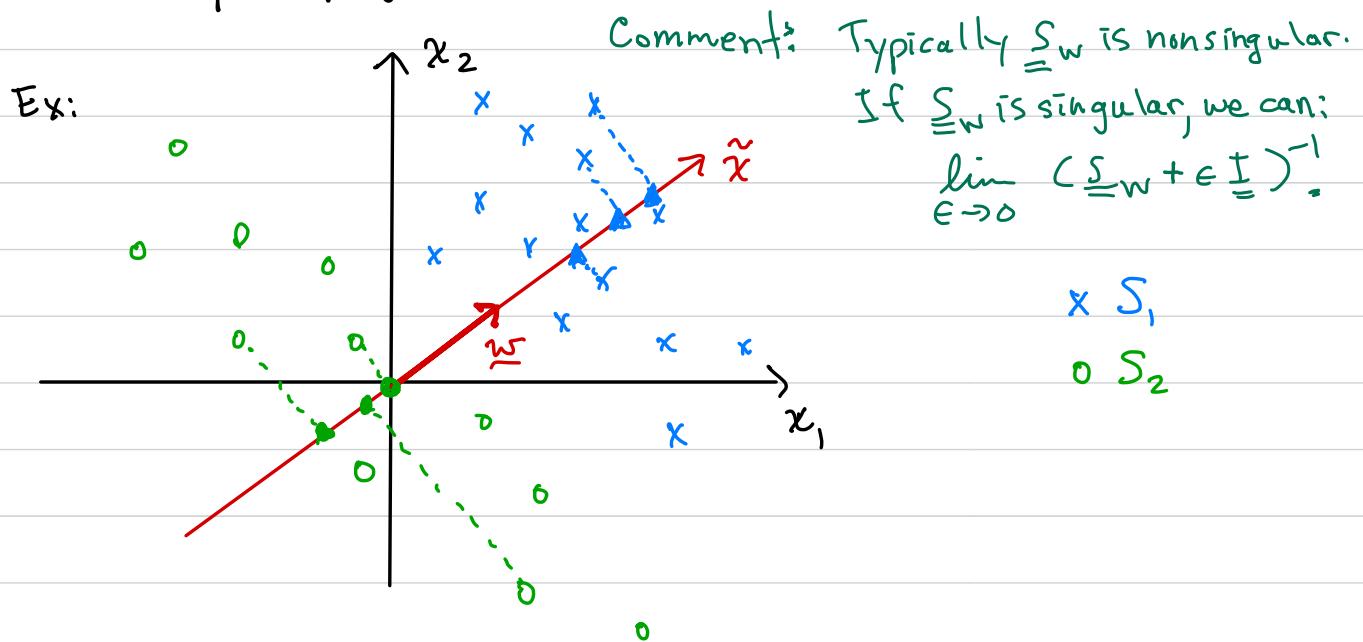
$$\nabla_{\underline{w}} J(\underline{w}) = \underline{0}.$$

Can show:

\Rightarrow if $\underline{\Sigma}_w$ is nonsingular, then

$$(2) \quad \underline{w} = \underline{\Sigma}_w^{-1} [\underline{m}_1 - \underline{m}_2] \quad \leftarrow \text{FLD solution.}$$

Note: direction of \underline{w} defines the new (1D) feature space. \therefore We can normalize \underline{w} to $\|\underline{w}\|=1$ for pure projection.



Above FLD has 2 limitations:

- (i) Always $D \rightarrow 1$ dimension.
- (ii) $C=2$ only.

\rightarrow both limitations can be removed — see below.

FLD can be extended to d final dimensions ($1 \leq d \leq D$)

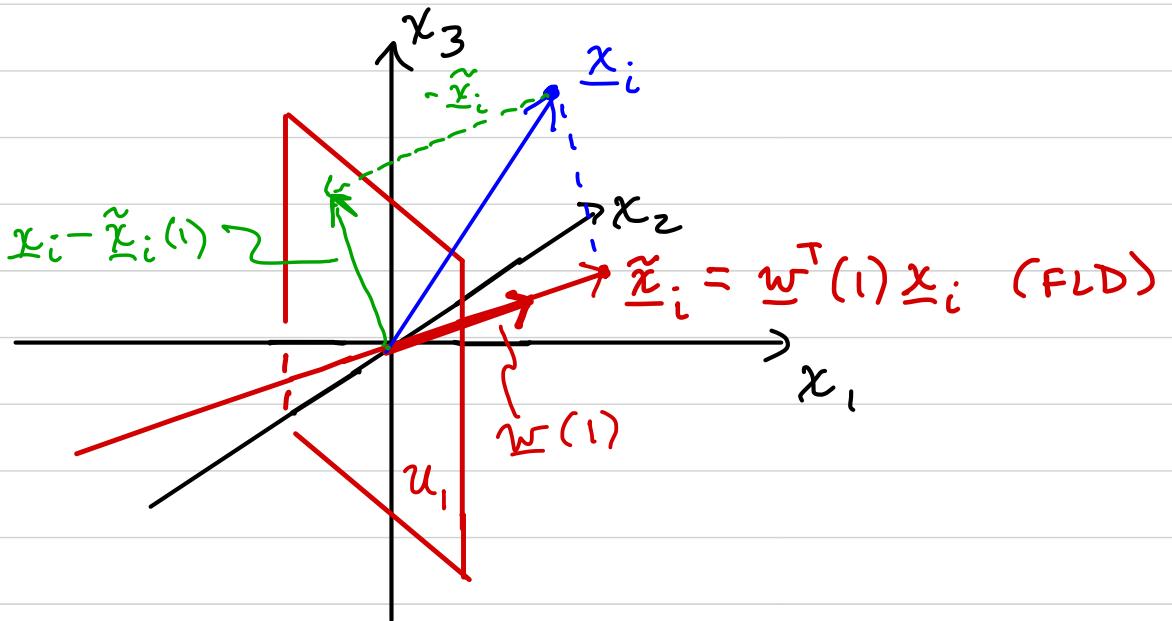
For example by:

1. Apply FLD for $D \rightarrow 1$; result is $\underline{w}(1)$

2. For $k=1$ to $(d-1)$

3. Subtract $\tilde{\underline{x}}_i \triangleq \underline{w}^T(k) \underline{x}_i$ from each data point \underline{x}_i .

(3)



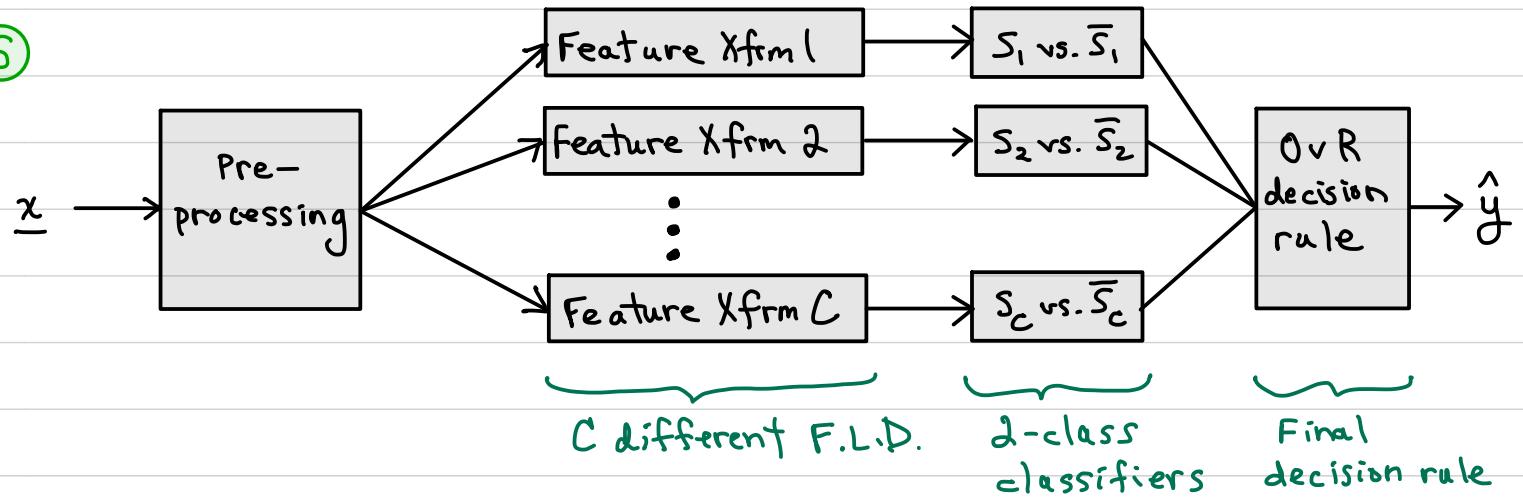
\Rightarrow work next in subspace U_k orthogonal to $\underline{w}(k)$

4. Find FLD for $(D-k) \rightarrow 1$ dimension in new subspace U_{k+1} ; result is $\underline{w}(k+1)$.

This $\underline{w}(k+1)$ becomes direction of $(k+1)^{th}$ axis $\rightarrow (k+1)^{th}$ new feature.

For $C > 2$

1. If using a set of 2-class classifiers (e.g., OvR, OvO), use FLD for each 2-class classifier, independently.
 Note that each 2-class problem could have a different feature set. For example:



OR

2. For any $C > 2$ classifier, use:

Multiple Discriminant Analysis (MDA)

Generalization of FLD to $C > 2$ classes.

for $D \geq C$

MDA will reduce dimensionality $D \rightarrow (C-1)$ (or sometimes less)

MDA optimizes a criterion J (generalization of $J(\tilde{w})$ for FLD).

Results are globally optimal (optimizes J globally)

For more information, see Bishop 4.1.6

(also Duda, Hart, & Stork 3.8.3)

Comments on PCA, FLD, MDA

1. PCA ignores class labels; FLD and MDA use class labels.
2. Algorithm (3) is a greedy algorithm.
Each feature is chosen optimally from the remaining features,
but end result isn't necessarily the best set of d features.
3. For FLD and MDA, if we add 2 assumptions:
data from each class is normally distributed;
all classes have the same covariance matrix;
then they are equivalent to linear discriminant analysis (LDA).
(many references use these terms synomously).
4. Any of these can be applied for nonlinear transformations,
using a Φ basis-set expansion or kernel substitution
(e.g., kernel PCA or kernel FLD).