

Machine Learning I: Supervised Methods

B. Keith Jenkins

Announcements

- Slido event code: 4509075
- Homework 6 will be posted later this week.
- Midterm exams are being graded
 - Estimate completion by Monday
 - If you thought the exam was too long, remember that class will be graded primarily on the curve.

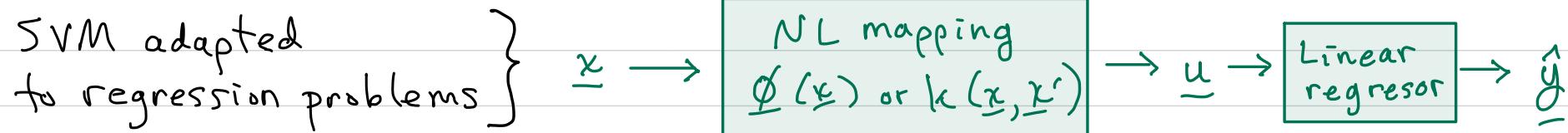
Reading

- Bishop 5.0-5.3
 - Feedforward artificial neural networks (ANN)

Today's lecture

- Support vector regression (part 2) (Jenkins)
- Validation and test sets (Sagar)
 - Model selection
 - Cross validation
 - Cross validation for model selection
- Introduction to VC dimension (Sagar)
 - Classifier complexity, d.o.f., constraints
- Feature selection and dimensionality reduction (part 1) (if time) (Sagar)
 - PCA

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:

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

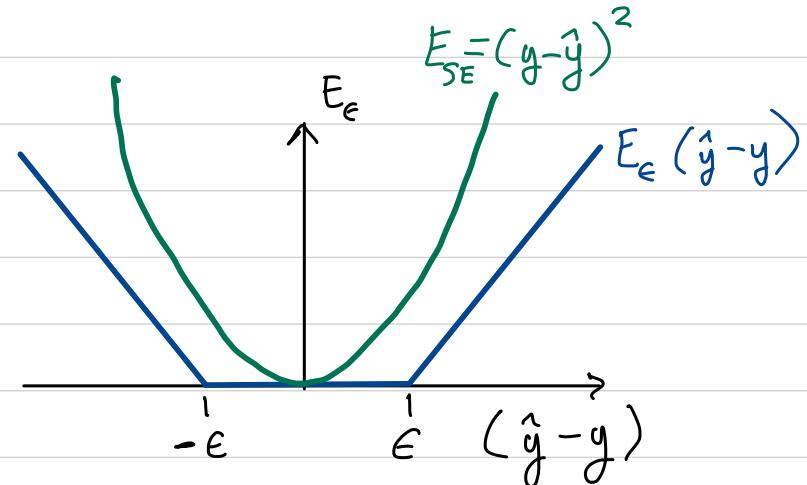
For SVR:

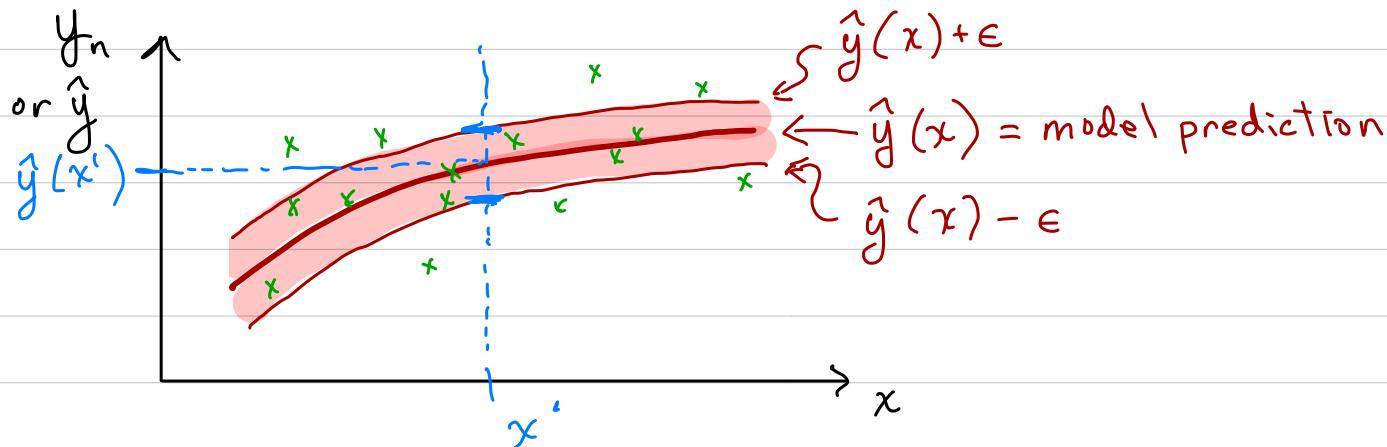
use " ϵ -insensitive" loss
instead of squared-error loss

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):

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

↑ C here by convention.

* Pages 4-5 were deferred to Lecture 17 due to technical difficulties with Webex audio.

VALIDATION AND ERROR ESTIMATION, INCLUDING MODEL SELECTION

How can we:

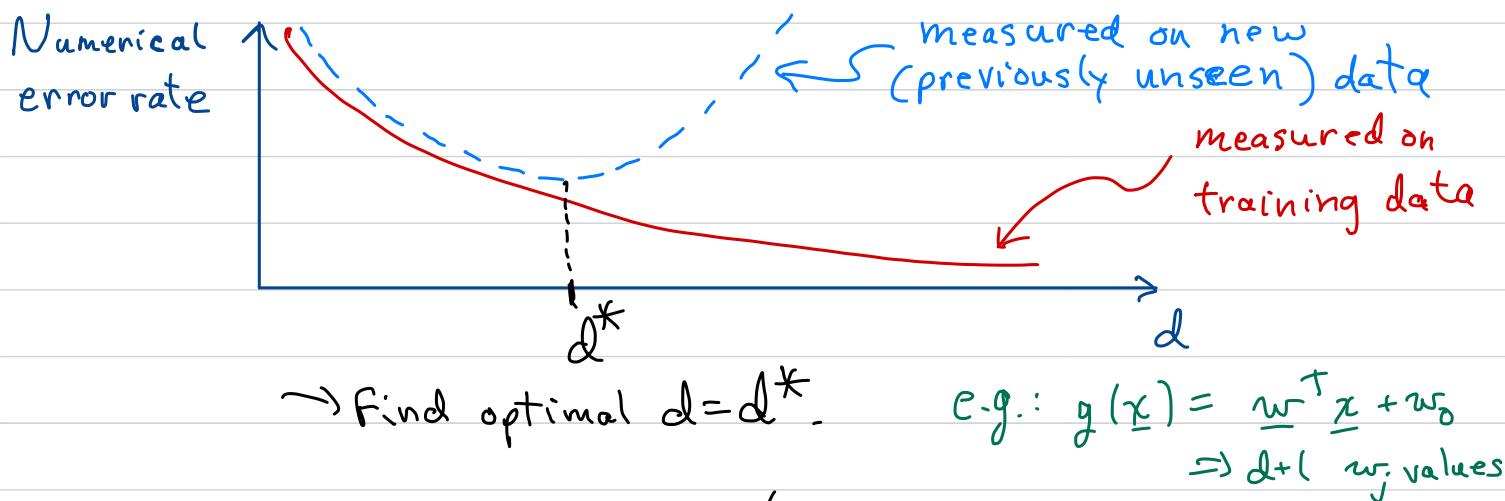
- (i) - Estimate error rate on unseen data ?
- (ii) - Choose parameter values and dimensionality for a high-accuracy result without overfitting ?

(ii) Validation (data) set - for model selection or evaluation ("validation")

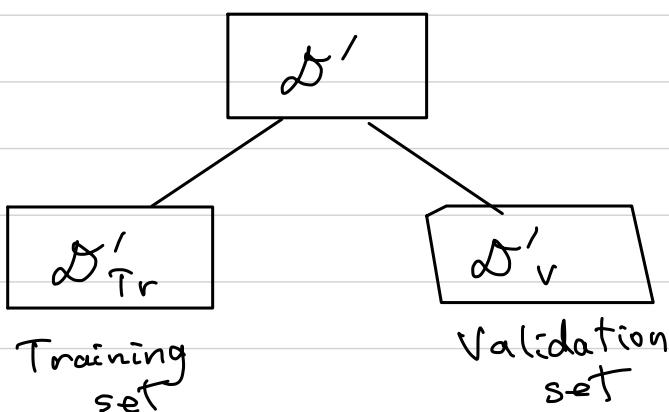
Ex: Linear classifier

Parameter: $d = \# \text{dimensions (variable)}$

[d here is #dimensions, not necessarily polynomial degree.]



Given a labelled dataset = \mathcal{D}'



Note:

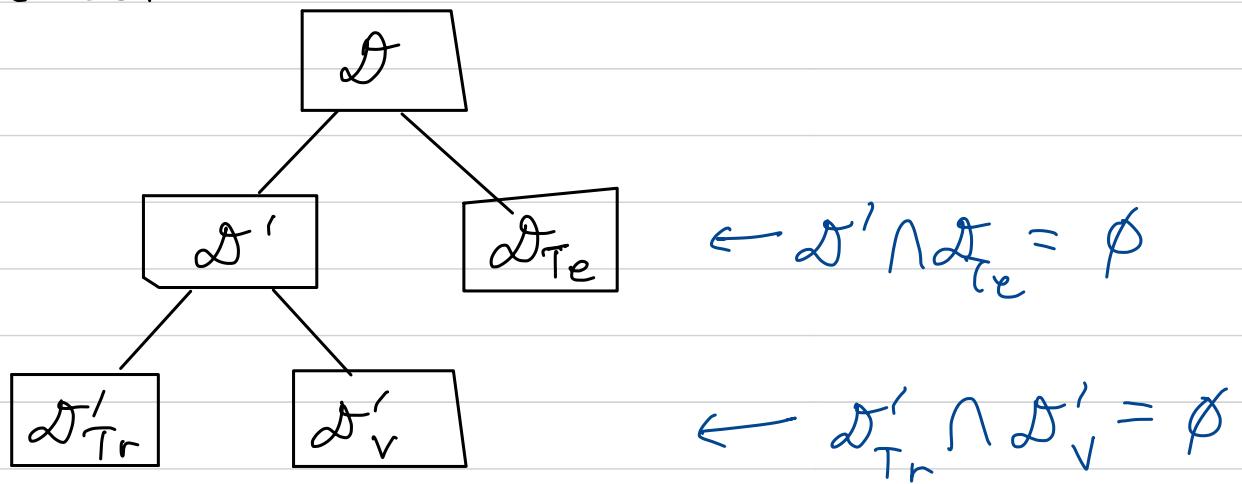
$$\mathcal{D}'_{Tr} \cap \mathcal{D}'_V = \emptyset.$$

(i) Test (data) set - for estimating error rate of final classifier.

→ To get a good estimate of error rate on unknowns ("out of sample error"), use a test set, \mathcal{D}_{Te} , that has not been seen or used by:

- training algorithm
- model selection / validation set
- person designing or optimizing the M.L. system (until the system is fully designed)

Procedure:



VALUATION FOR MODEL SELECTION

Model selection:

Ex: Suppose use SVM with RBF kernel:

$$k(x, x_n) = \exp\left\{-\gamma \|x - x_n\|^2\right\}, \quad \gamma > 0.$$

We want to find best choice of γ .

[Comment: consider γ here to be a dummy variable for a parameter that we want to optimize by using model selection.]

Suppose we use a search over values of γ .

→ Use a validation set

For $k = 0, 1, 2, \dots, K$

Use model γ : $\gamma = \gamma_{\min} + k \Delta \gamma$

Initialize training procedure

Train classifier using γ on $D_{Tr}' \rightarrow \hat{w}(\gamma)$

Evaluate classifier on D_V'

Store error ϵ_γ

Next k

Pick value $\gamma^* = \arg \min_p \epsilon_p$

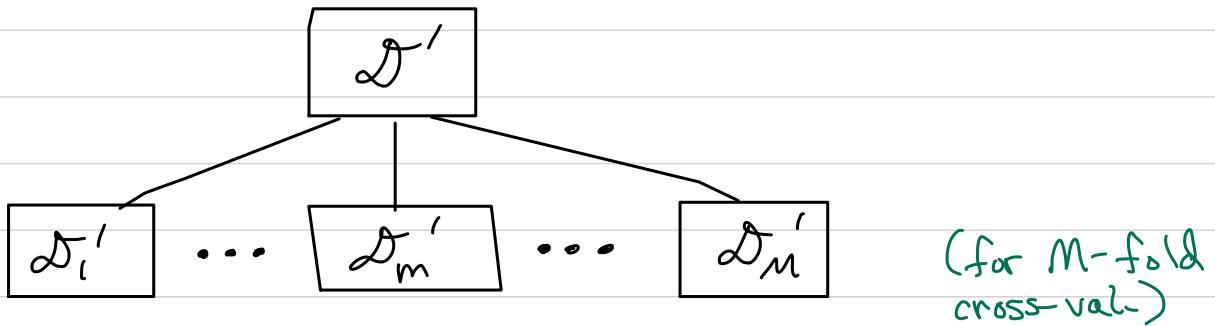
If we want to use data efficiently (keeping D_{Tr}' as large as possible), we can use cross validation instead.

CROSS VALIDATION FOR MODEL SELECTION

→ Re-use (carefully) data points in \mathcal{D}_{Tr}' , \mathcal{D}_V' .

$\mathcal{D}' \rightarrow M \sim$ equal-sized subsets.

(Chosen randomly (without replacement),
but typically subject to constraint that
% of data pts. in each class \sim same in
each subset as in \mathcal{D}' .) ("stratified")



→ all subsets \mathcal{D}_i' are disjoint.

For $i = 1, \dots, M$

Let $\mathcal{D}_V' = \mathcal{D}_i'$, $\mathcal{D}_{Tr}' = \bigcup_{j \neq i} \mathcal{D}_j'$

(i = "fold"
of cross-validation)

For $k = 0, 1, 2, \dots, K$

Use model p : $p = p_{min} + k \Delta p$

- crucial step { Similar to validation above } →
- Initialize training procedure
 - Train classifier using p on \mathcal{D}_{Tr}'
 - Evaluate classifier on \mathcal{D}_V'
 - Store error E_{pi}
 - Next k
 - Optionally store $p_i^* = \arg \min_p E_{pi}$
 - Next i

If run cross-val. only once, could:

$$\hat{E}_p = \frac{1}{M} \sum_{i=1}^M E_{p_i}$$

$$p^* = \operatorname{argmin}_p \hat{E}_p$$

Typically run a number of times (new partitioning of $\mathcal{D}' \rightarrow \mathcal{D}'_i$ each time) to get better estimates.

Run T times. ($t = 1, 2, \dots, T$)

Define:

$$\bar{E}_p = \frac{1}{MT} \sum_{t=1}^T \sum_{i=1}^M \hat{E}_{p_i}^{(t)}$$

= average error for parameter value p.

Then: pick

$$\overrightarrow{p^*} = \operatorname{argmin}_p \bar{E}_p$$

Then, how find best $\underline{w} = \underline{w}^*$ for $p = p^*$?

→ Train again using $p = p^*$, and using \mathcal{D}' (call training data).

THEN, FOR ERROR ESTIMATION (ON UNSEEN DATA)

- (1) Use \mathcal{D}_{Te} to test classifier (or M.L. system) that uses p^*, \underline{w}^* (the chosen model). [one pass method]
- or
- (2) Add an outer loop of cross validation for $\mathcal{D} \xrightarrow{\mathcal{D}'} \mathcal{D}_{Te}$

DEGREES OF FREEDOM, CLASSIFIER COMPLEXITY, AND VC DIMENSION

Earlier: d.o.f. = #indep. variables or parameters (during learning)

#constraints \approx #training data points (e.g., in MSE problem like pseudoinverse)

Risk of overfitting if #constr. < $(3-10) \times$ d.o.f.

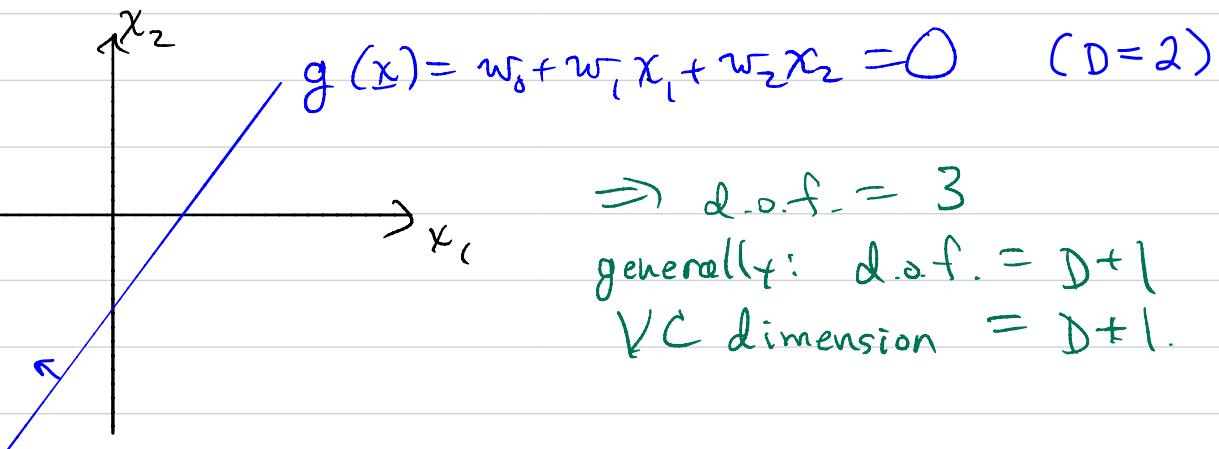
Instead of d.o.f. (for classifier complexity), a more rigorous metric is

VC dimension [Vapnik-Chervonenkis],

which quantifies the "flexibility" or "capacity" of a decision boundary (or dec. regions).

Examples:

1. Linear boundary in $(D+1)$ space [D features] (2-class problem)



2. Polynomial boundary $[g(\underline{x}) \stackrel{S_1}{\geq} 0, \text{ with } g(\underline{x}) = \text{polyn. fcn. of } \underline{x}]$ ($C=2$
classes)

If g has $D'+1$ weights, then: $\text{VC dim} = D'+1$.

Note that: higher VCdim \Rightarrow higher capacity
of dec. boundary.

$\therefore \text{VCdim} \leftrightarrow N$ is important
(overfit, underfit etc.)

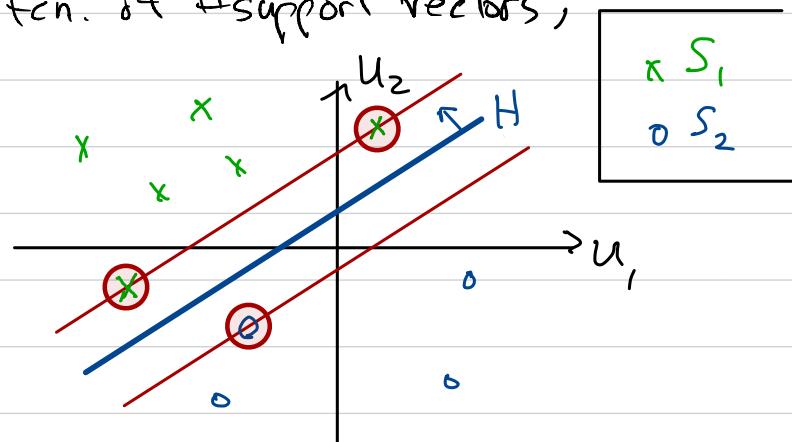
\Rightarrow can tell us how good we expect
the generalization
(performance on unknowns) to be.

What about SVM?

- For SVM, because of its additional constraints included in the criterion function, its VCdim can be $< D'+1$, and e.g. can be a fcn. of #support vectors, instead of #weights.

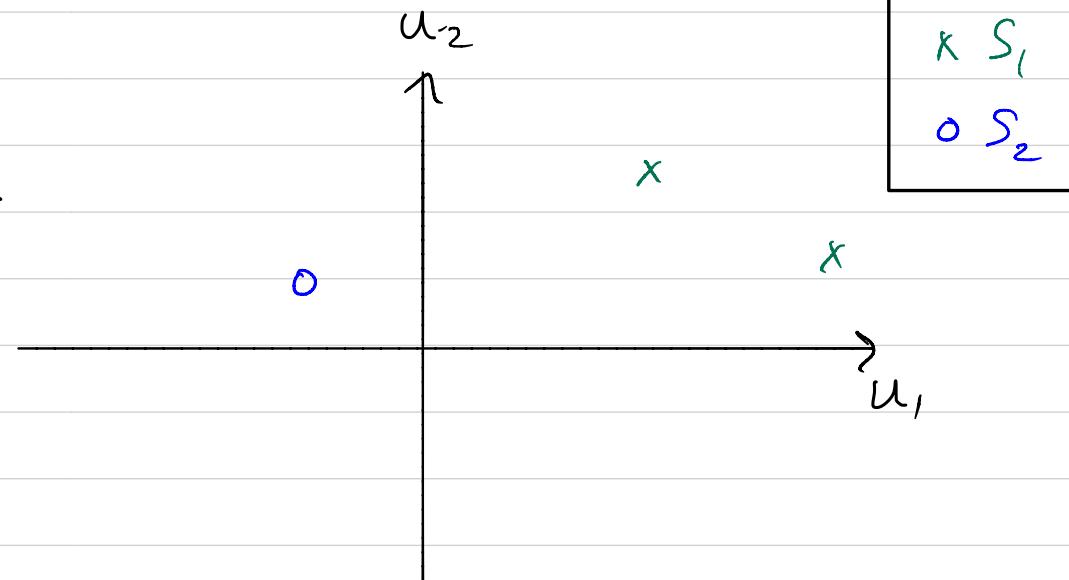
- The support vectors are the data points that lie on the margin boundary (after training).
(linearly separable case)

○ Support vector

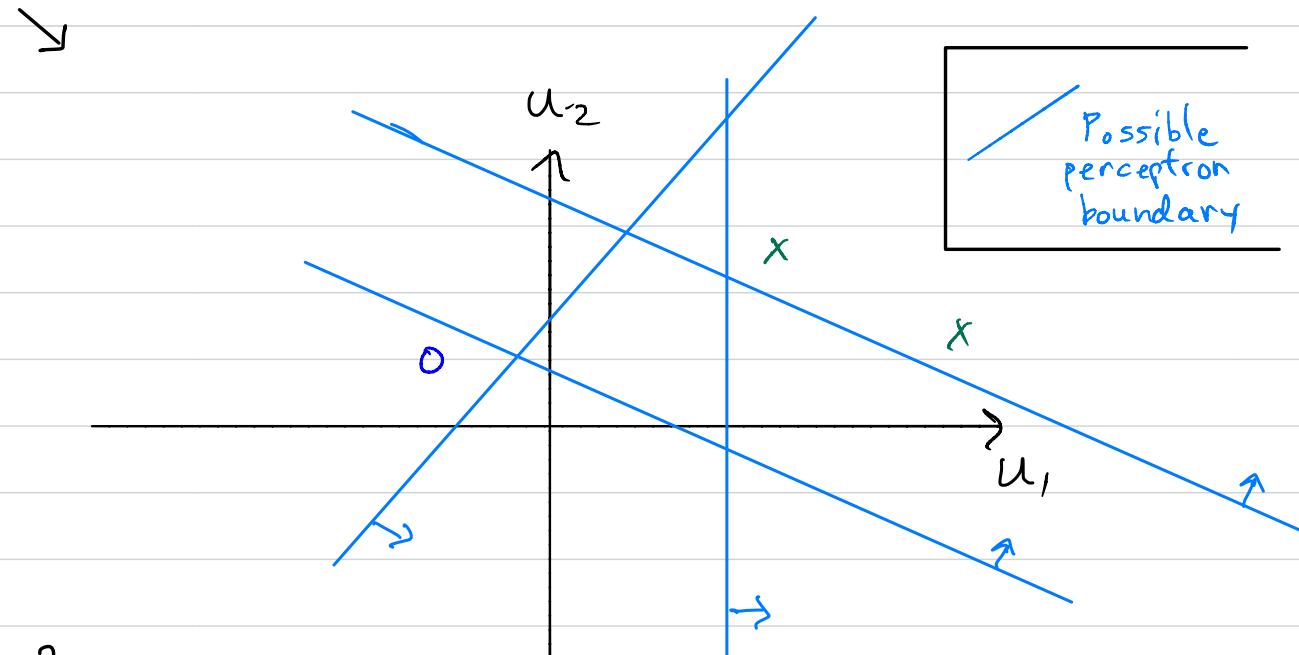


Intuitive example:

- Given 3 data points in $D=2$ unaugmented feature space \longrightarrow

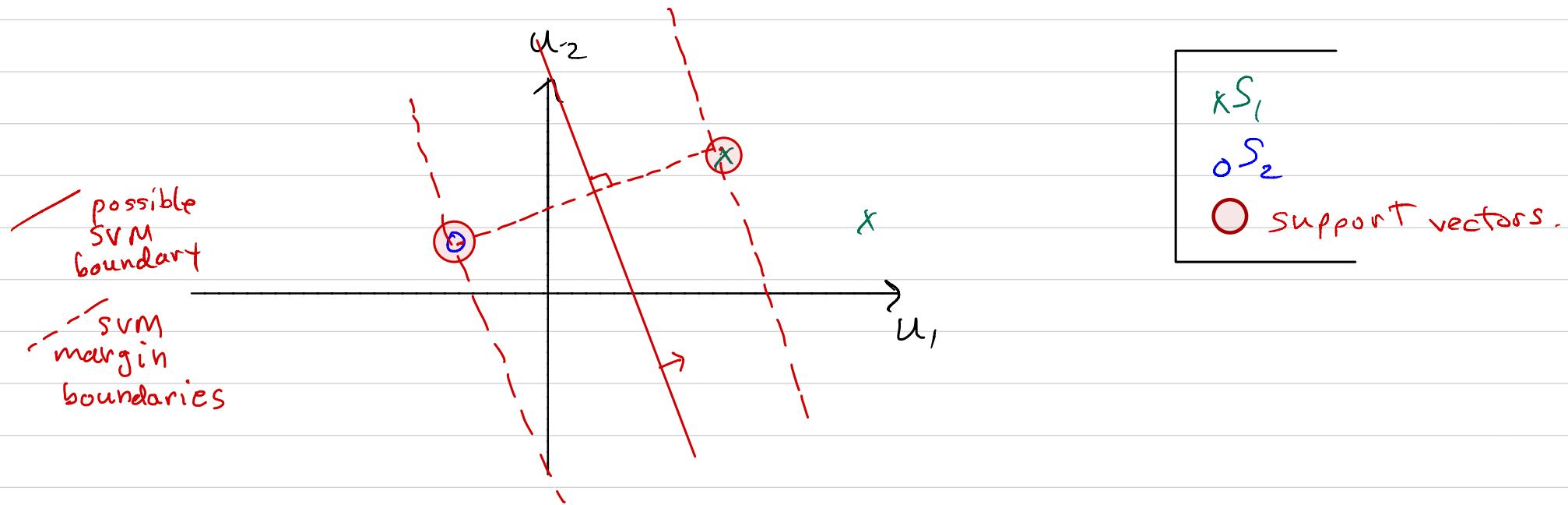


- What decision boundaries could perceptron converge to?



- What decision boundaries could SVM converge to?





- Thus, SVM can generalize much better with sparsely populated feature space ($(\# \text{weights}) \times (3-10) \leq N$), than many other classifiers (e.g., pseudoinv-or perceptron).
- This allows $\not\perp$ mapping to higher dimensional space in cases that would yield too many weights using other classifiers (that add no constraints).

FEATURE SELECTION AND DIMENSIONALITY REDUCTION

Why?

1. Reduce computation requirements
2. Reduce likelihood of overfitting
- by balancing d.o.f. and N.
3. Reduce correlation among features
4. Eliminate irrelevant features

Effect: reduce feature space dimension $D \rightarrow D' < D$.

Two overall scenarios

(1) Feature selection.

Choose some features to eliminate (using some criterion)
 $\Rightarrow D' < D$, features of D' are a subset of features of D .

(2) Feature-space transformation

Transform to a new, lower dimensional, feature space.

- PCA: keep D' principal components. ($D' \leq D$).

- FLD: find an optimal transformation

using FLD criterion

(binary classification problems).

- MDA: extends FLD to $C > 2$ classes.

$\Rightarrow D' < D$ typically. Features of D' are linear combinations of features of D .

PRINCIPAL COMPONENTS ANALYSIS (PCA) [Bishop 12.1]

Also called: Karhunen-Loeve (KL) Transform

Idea: Project data (of all classes) onto a lower-dimensional space.

$D \rightarrow D'$, $D' < D$.

using a linear transform.



PCA's criterion:

- Minimize MSE between projected points \tilde{x}_n and original points \underline{x}_n .

PCA can be computed by orthonormal transformation

Need $\hat{\Sigma}$ = covariance matrix. (will be covered next week)

Use an estimate based on training dataset -

(1) sample covariance matrix: Let \hat{m} = sample mean = $\frac{1}{N} \sum_{i=1}^N \underline{x}_i$

$$\hat{\Sigma} \triangleq \frac{1}{N} \sum_{i=1}^N (\underline{x}_i - \hat{m})(\underline{x}_i - \hat{m})^T \quad (\hat{\Sigma} \text{ is real and symmetric})$$

(2) Find eigenvectors \underline{e}_d and eigenvalues λ_d of $\hat{\Sigma}$ [1]

(3) Keep " corresponding to D' largest eigenvalues.

(4) Transform all data pts. \underline{x}_i to new feature space:

$$\tilde{\underline{x}}_i = \underline{E}^T \underline{x}_i, \quad \underline{E} \triangleq [\underline{e}_1 \ \underline{e}_2 \ \dots \ \underline{e}_{D'}] \text{ such that}$$

$\{\underline{e}_1, \underline{e}_2, \dots, \underline{e}_{D'}\}$ is orthonormal. (each \underline{e}_j is normalized to unit length)

[1] $\hat{\Sigma}$ is assumed positive definite (more later), so its eigenvalues are all positive.

Suppose we have 2-class classification problem with $D=2$.

We use PCA to reduce feature space to $D=1$.

Then we use a linear classifier in the 1D feature space.

How well will the classifier perform on the following 2 datasets?



How well does PCA work?

Ex: $C=2$ classes

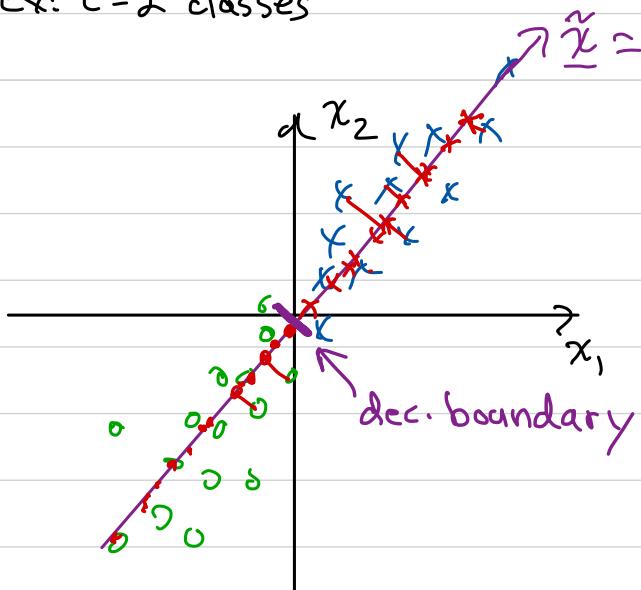


Fig. A

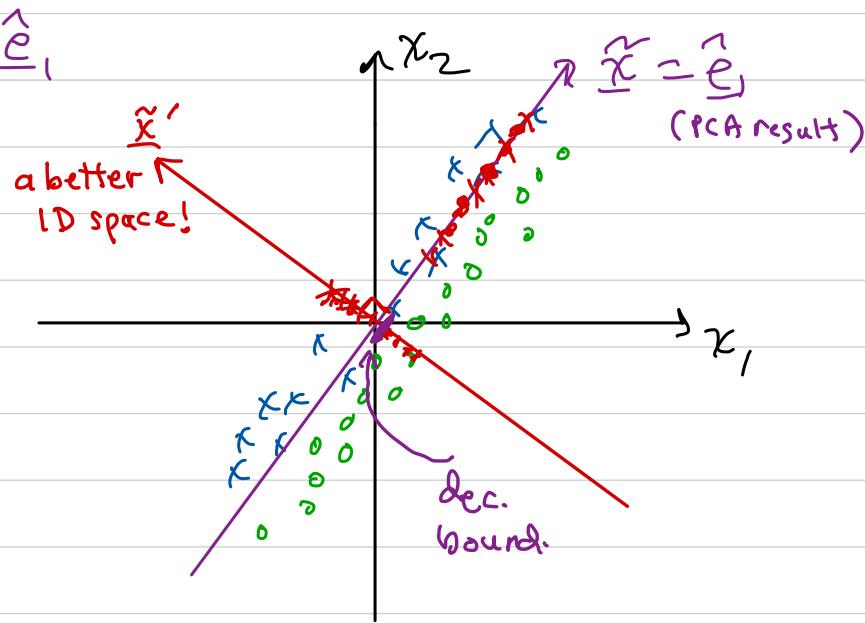


Fig. B

o | PCA should work well

PCA works poorly.

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