

Machine Learning from Data

Lecture 24: Spring 2021

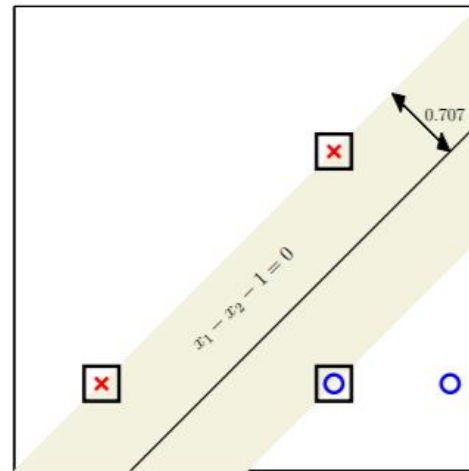
Today's Lecture

- Support Vector Machines (SVMs)
 - Why is fattest hyperplane better?
 - Non-separable Data

RECAP: The Optimal Hyperplane

The Optimal Hyperplane

The fattest hyperplane that separates the data
tolerates most measurement error



1. Can we efficiently find the fattest separator?
2. Is fatter better than thin?

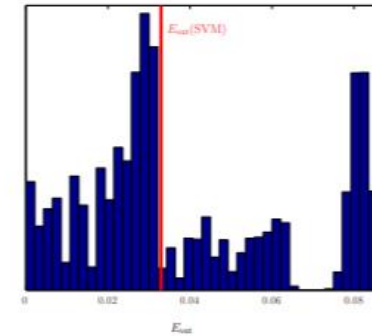
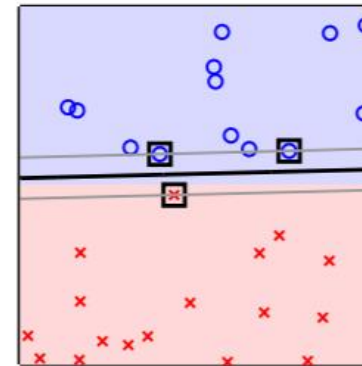
The Algorithm

Quadratic Programming:

$$\underset{b, \mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{subject to: } y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \text{for } n = 1, \dots, N.$$

Support vectors: the data points that sit on the cushion.
Using only support vectors, the classifier does not change.



PLA depends on the (random) order of data

Link to Regularization

optimal hyperplane

minimize $\frac{1}{2} \mathbf{w}^T \mathbf{w}$

subject to: $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$ for $n = 1, \dots, N$.

regularization

minimize $E_{\text{in}}(\mathbf{w})$

subject to: $\mathbf{w}^T \mathbf{w} \leq C$.

	optimal hyperplane	regularization
minimize	$\mathbf{w}^T \mathbf{w}$	E_{in}
subject to	$E_{\text{in}} = 0$	$\mathbf{w}^T \mathbf{w} \leq C$

The optimal hyperplane performs ‘automatic’ regularization.

Evidence that Larger Margin is Better

- (1) Experimental: larger margin gives lower E_{out} ; **bias** drops a little and **var** a lot.
- (2) Bound for d_{VC} can be less than $d + 1$ – fat hyperplanes generalize better.
- (3) E_{cv} bound does not explicitly depend on d .

Larger Margin is better

Generate a separable dataset ($N=20$)

- Repeat
- ①
 - ② Randomly generate 50,000 separating hyperplanes
 - ③
- | | | | |
|-----------|---------------|-----------------------------------|---|
| h_1 | \rightarrow | $r(h_1)/r_{opt} = \rho_1$ | $\left. \begin{array}{l} E_{out}(h_1) \\ E_{out}(h_2) \\ \vdots \\ E_{out}(h_{50K}) \end{array} \right\} \begin{array}{l} \text{Optimal hyperplane} \\ r_{opt} \end{array}$ |
| h_2 | \rightarrow | $r(h_2)/r_{opt} = \rho_2$ | |
| \vdots | | \vdots | |
| \vdots | | \vdots | |
| h_{50K} | \rightarrow | $r(h_{50K})/r_{opt} = \rho_{50K}$ | |

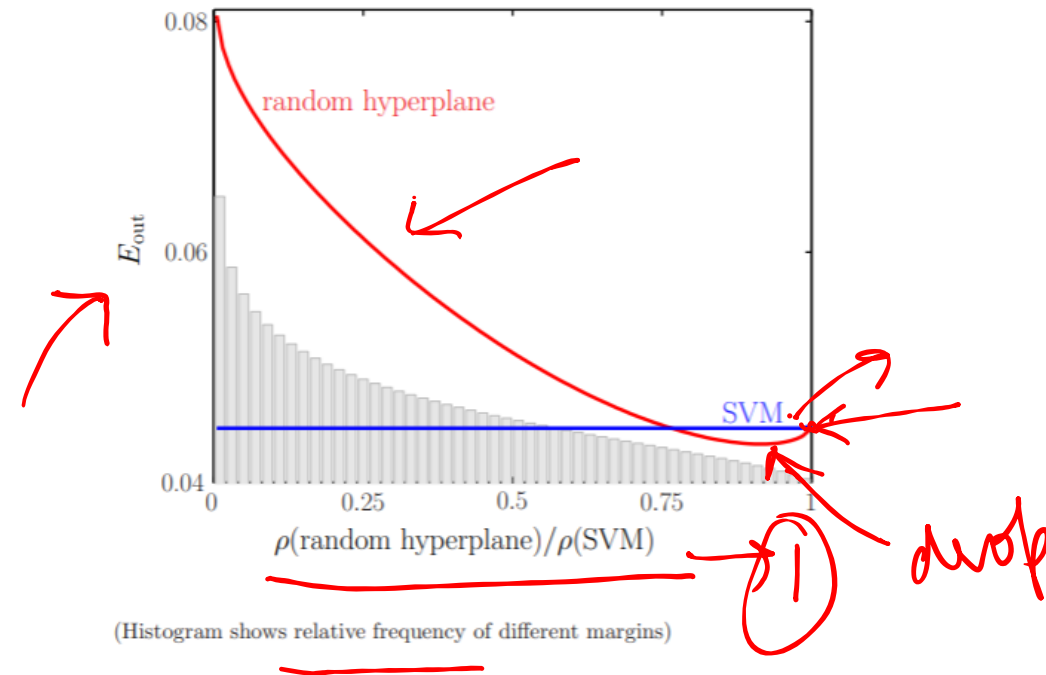
Larger Margin is Better

Generate a random separable data set ($N = 20$)

Select 50,000 random separating hyperplanes h

Compute E_{out} and $\rho(h)/\rho(SVM)$

Average over several thousands of random data sets



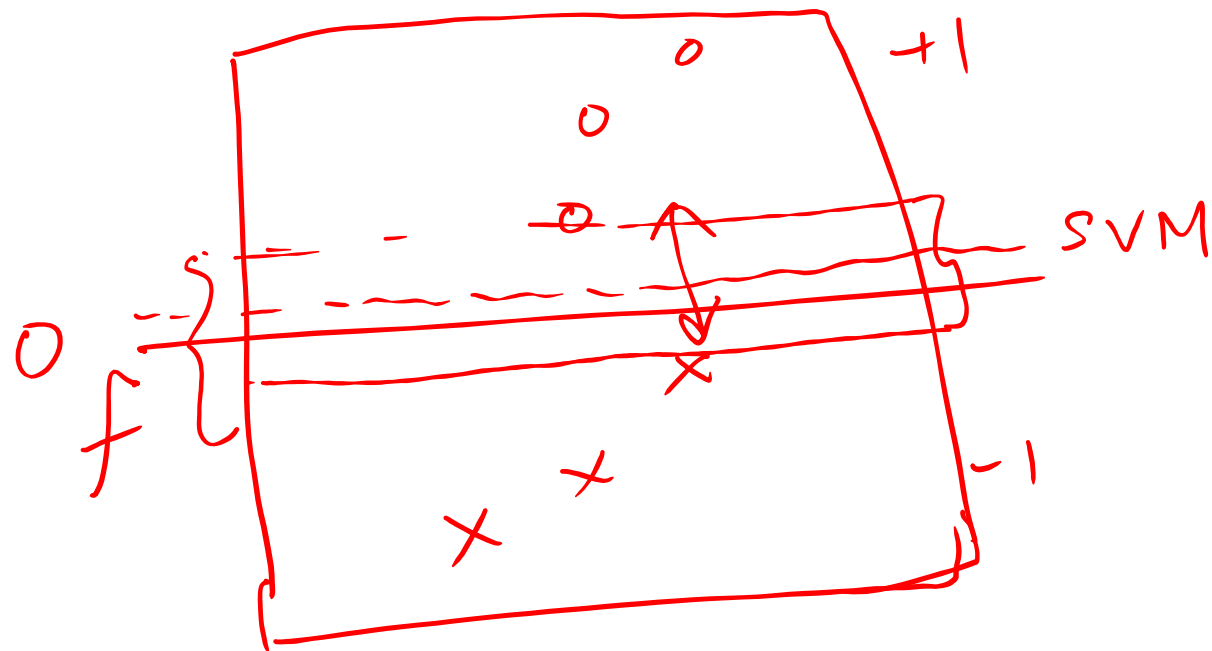
Bigger margin is generally better

Biggest is not best.

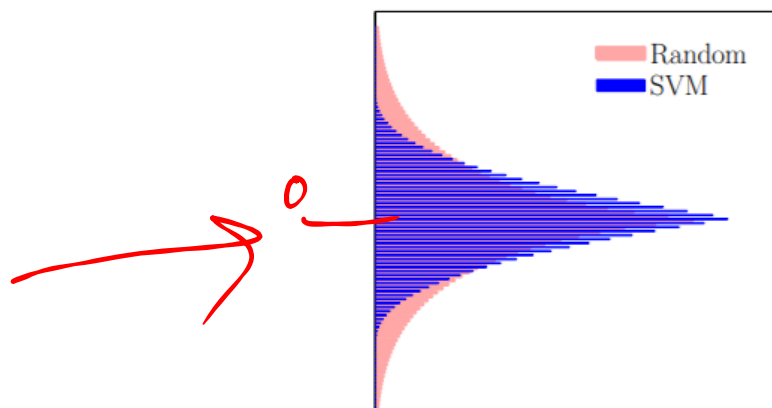
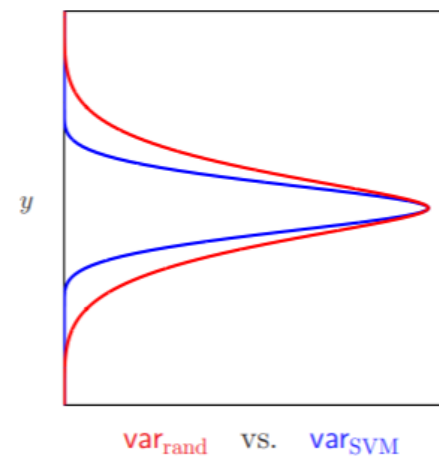
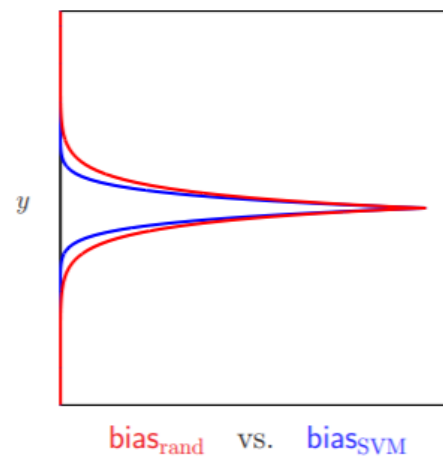
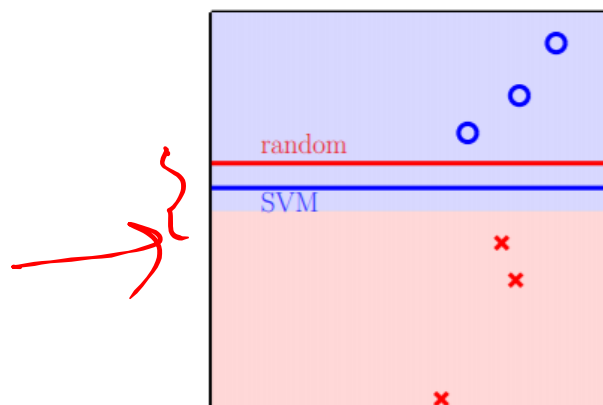
← Data other than support vectors can have role in fine-tuning

① Optimal

② Random

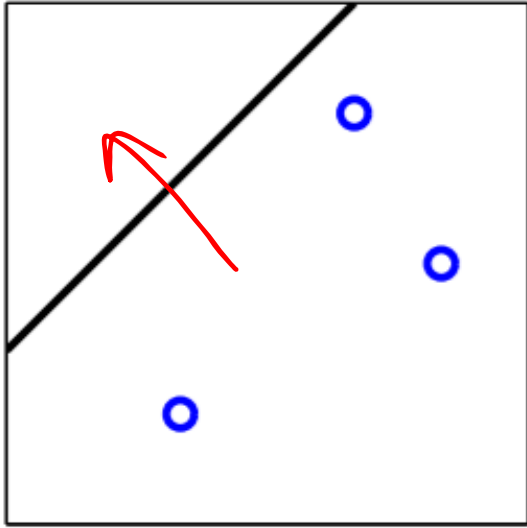


Bias and Variance

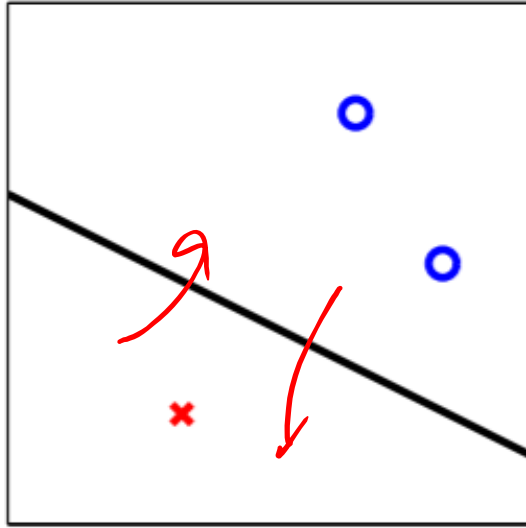


	<u>Random</u>	<u>SVM</u>	
bias	0.02	<u>0.015</u>	-0.005
var	0.059	0.038	-0.021
E_{out}	0.079	<u>0.053</u>	-0.026

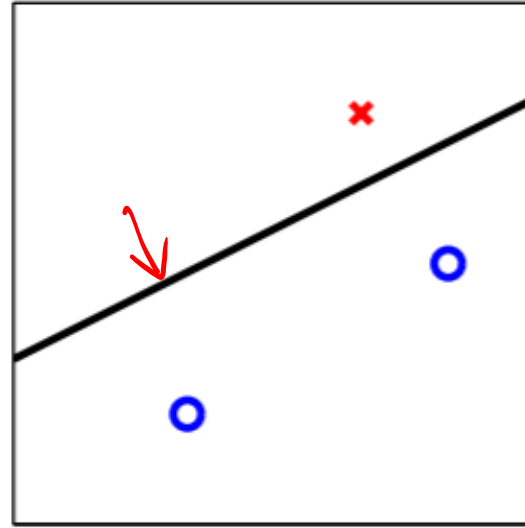
Fat Hyperplanes Shatter Fewer Points



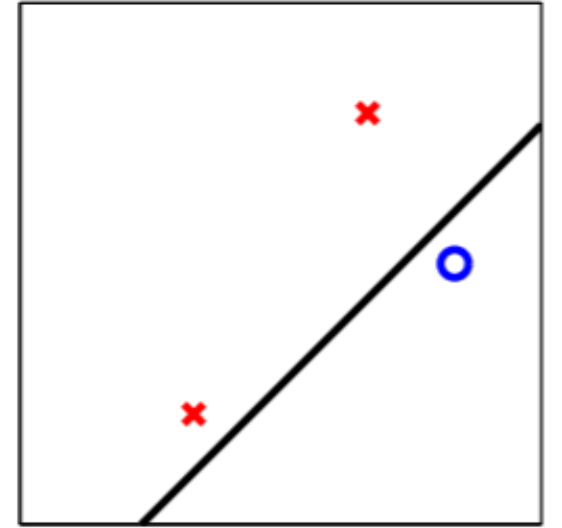
↑
doc



↑

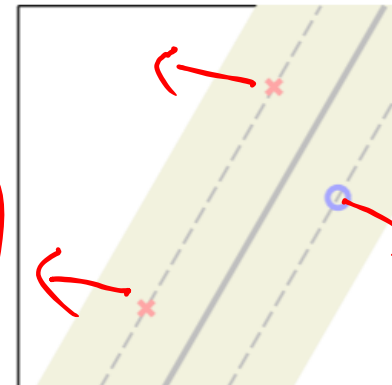
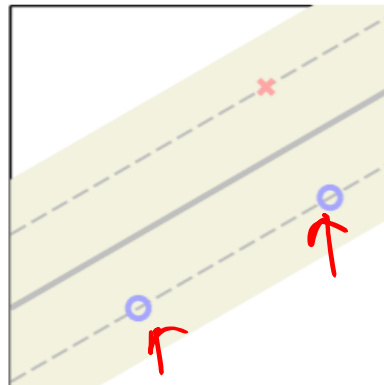
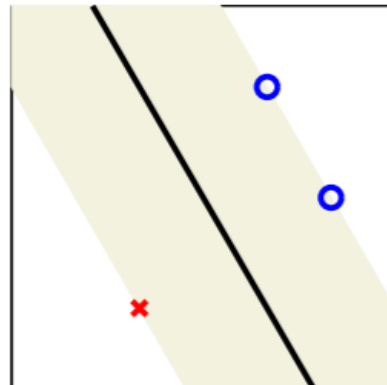
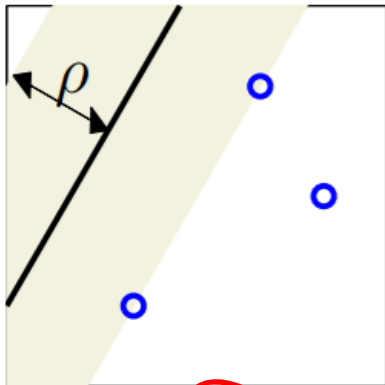
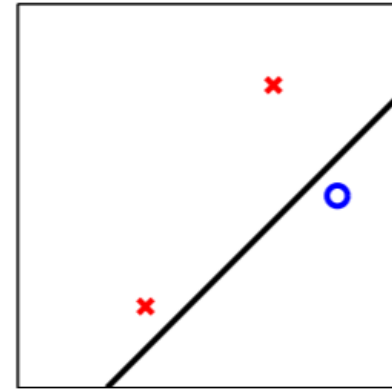
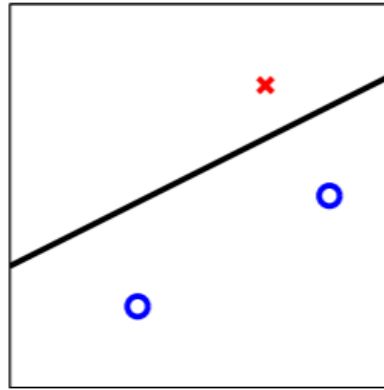
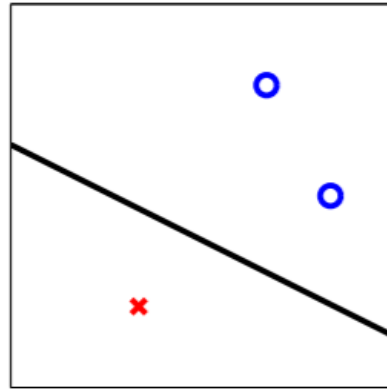
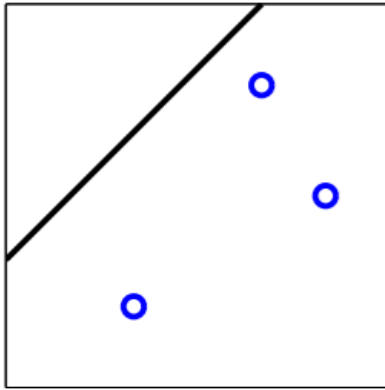


↑



↑

Fat Hyperplanes Shatter Fewer Points



*Sphere of
radius*

R

R/γ

①

②

✓

✓

Theorem

(induction)

$$d_{VC}(V) \leq$$

$$\left\lceil \frac{R^2}{n^2} \right\rceil + 1$$

↙ bigger

→ thickness of hyperplane.

↓
lower

$E_{in} \neq E_{out}$

$$, d_{VC} \leq d + 1$$

$$d_{VC} \leq \min \left(d + 1, \left\lceil \frac{R^2}{n^2} \right\rceil + 1 \right)$$

↙ independent of dimensions.

↘ n → big
 R → small

E_{cv} (error) of max-margin (fattest) hyperplane

Result to prove

E_{cv} (fattest hyperplane)
of support vectors

$E_{cv} \rightarrow$ unbiased
estimator of

$E_{out}(N-1)$

$\leq N$ are imp.

$N \leftarrow$
 $=$

$E_{cv} \rightarrow$ small $\rightarrow E_{out}$
 \downarrow control



Proof:

$$e_{L00}(x_{\text{notsv}}) = 0$$

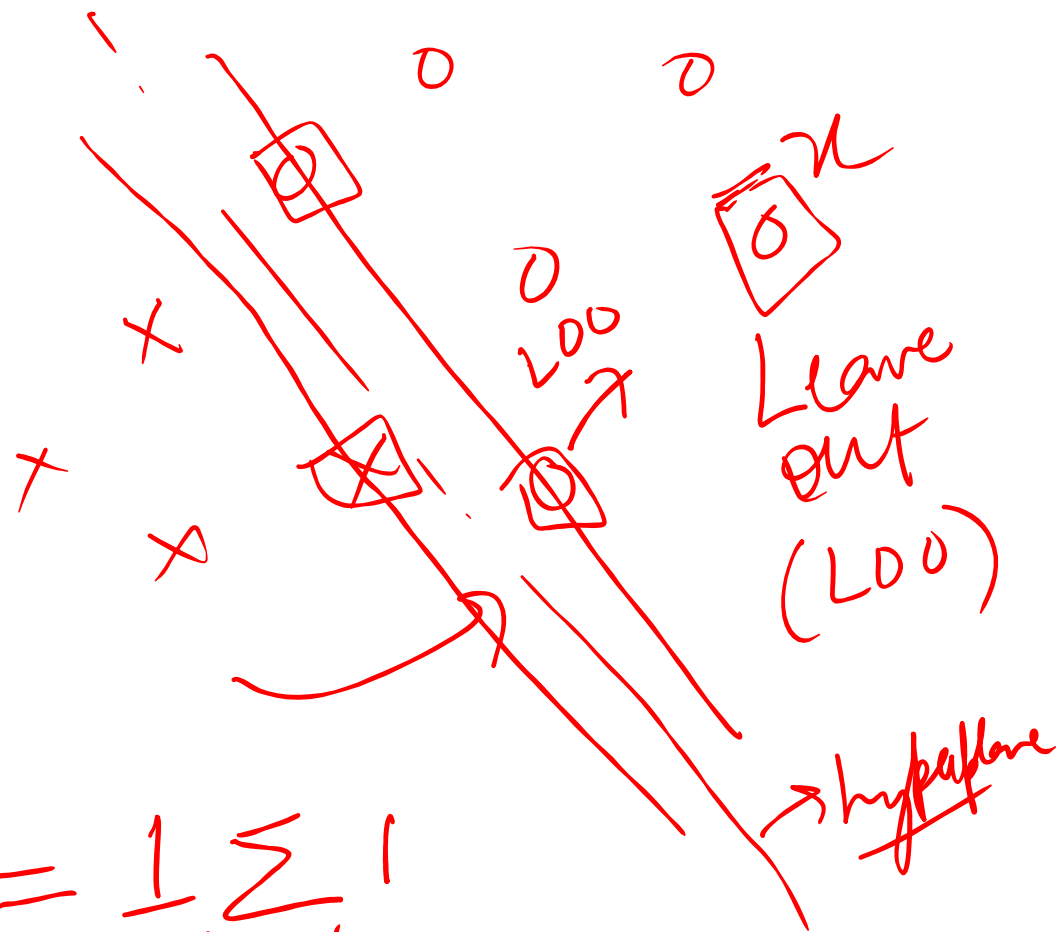
$$e_{L00}(x_{\text{sv}}) \leq 1$$

$$\therefore E_{\text{cv}} = \frac{1}{N} \sum_{n=1}^N e_{L00}(x_n)$$

$$= \frac{1}{N} \sum_{\text{SV}} e_{L00}(x_{\text{sv}}) = \frac{1}{N} \sum_{\text{SV}} 1$$

$$E_{\text{cv}} \rightarrow E_{\text{out}}$$

$$= \frac{\# \text{ of SVS}}{N}$$



Perceptron : $d_{vc} = d + 1$

↓
hyperplane

Optimal

1) Experimentally

$$2) d_{vc} \leq \min \left(d, \left\lceil \frac{RV}{n^2} \right\rceil \right) + 1$$

$$3) E_{cv} \leq \frac{\# \text{ of SVs}}{N}$$

Bias ↓ Variance ↓ ↓ (∵ $E_{out} ↓$)

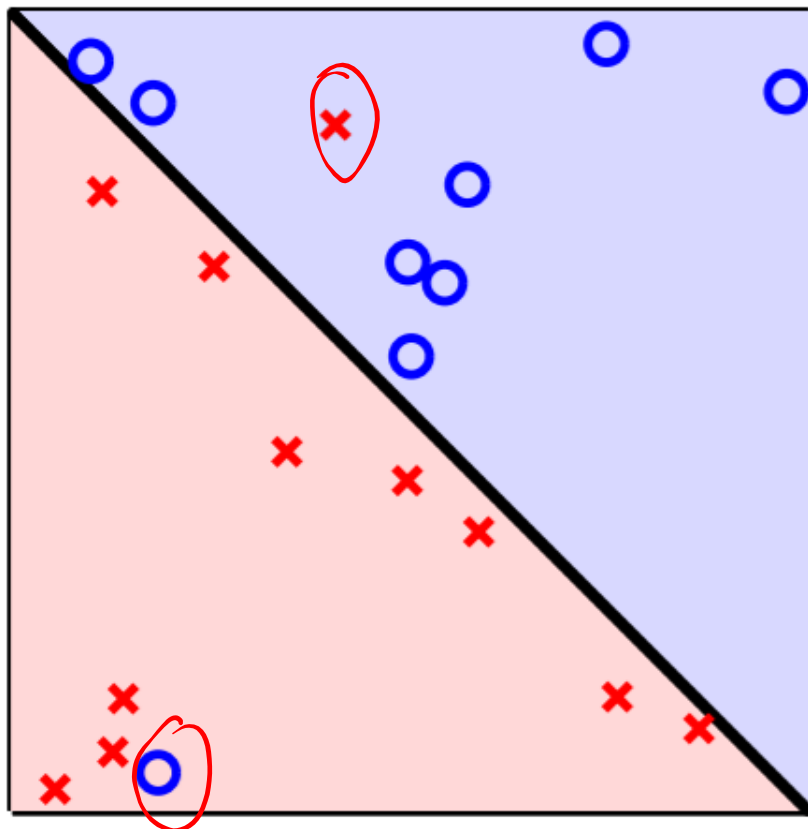
$E_{in} \rightarrow E_{out}$

(4) Algorithm: QP

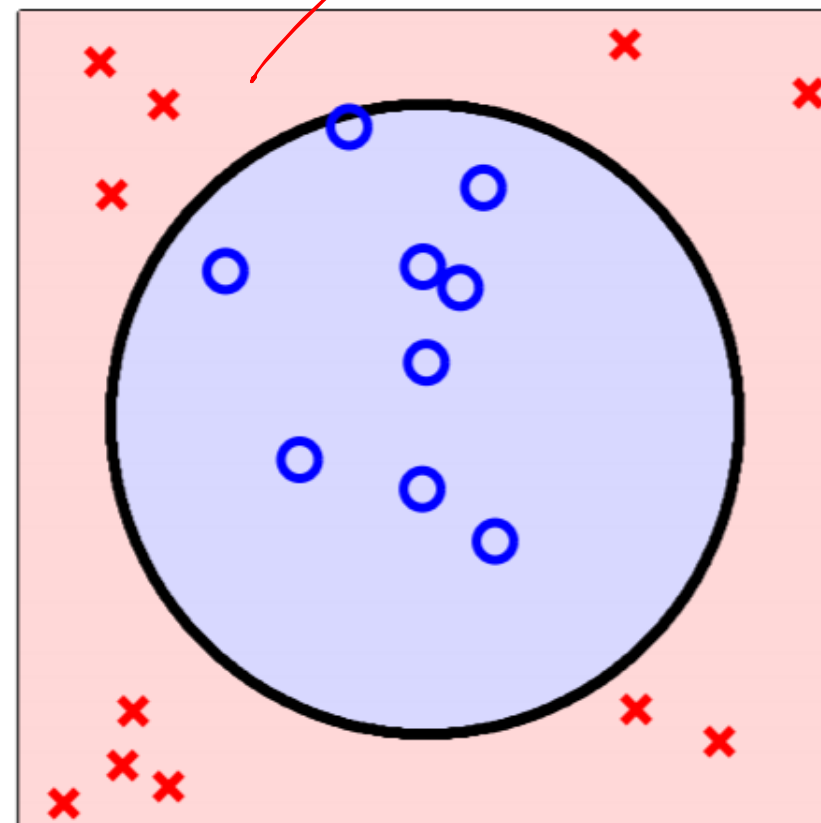
PLA

$$\bar{E}_{cv} \leq \frac{R^2}{N} \rightarrow \text{optimal}$$

Non-Separable Data

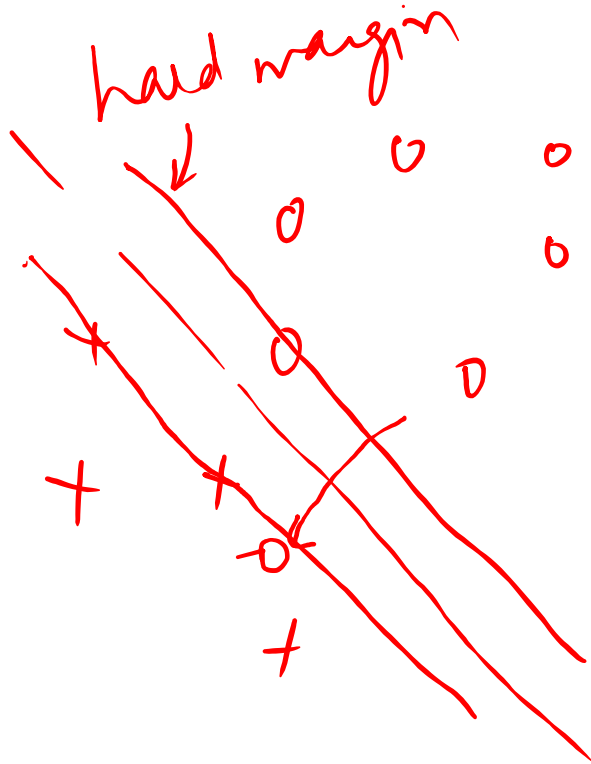


tolerate error



nonlinear transform

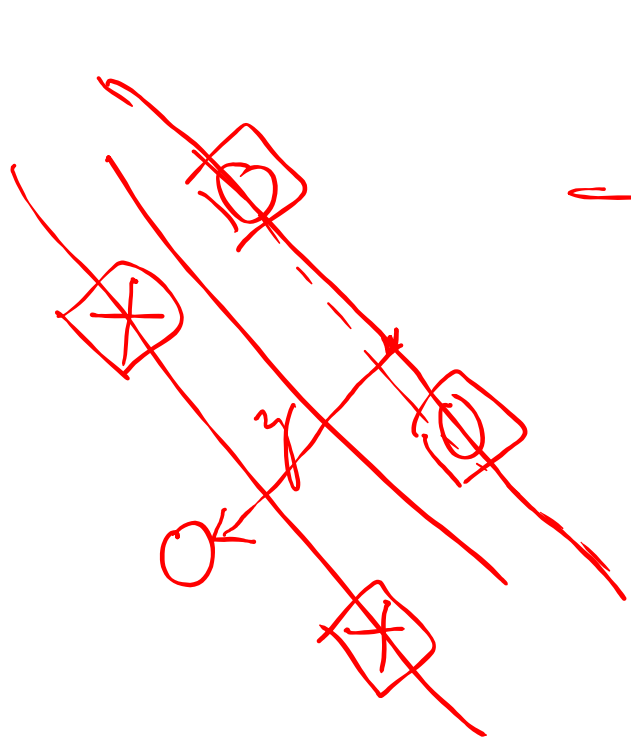
SOFT MARGIN OPTIMAL HYPERPLANE



$$\min \frac{1}{2} \omega^T \omega$$

$$\text{s.t. } y_n (\omega^T x_n + b) \geq 1 \quad \left. \vphantom{\frac{1}{2} \omega^T \omega} \right\} \underline{\text{S vs S}}$$

Original condition: $y_n (\omega^T x_n + b) = 1$
substituting



Soft Margin
optimal hyperplane

Variables: w, b, z

$$\begin{aligned} y_n(w^T x_n + b) &\geq 1 - z_n \\ z_n &\geq 0 \end{aligned} \quad \left\{ \begin{array}{l} z_n \rightarrow \infty \\ \text{penalty} \end{array} \right.$$

Optimization:

$$\begin{aligned} \min_{w, b, z} & \frac{1}{2} w^T w + C \sum z_n \\ \text{s.t.} & y_n(w^T x_n + b) \geq 1 - z_n \\ & z_n \geq 0 \end{aligned}$$

Regularization.

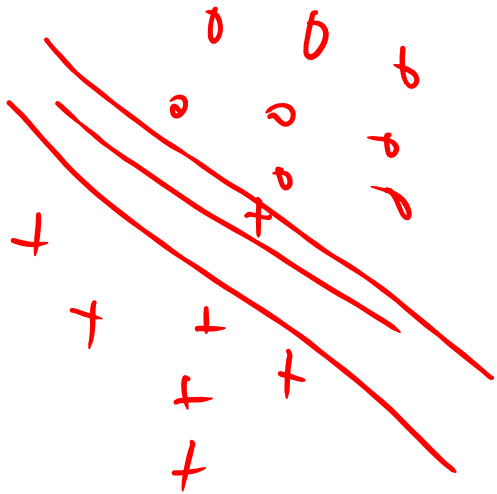
margin

margin violation

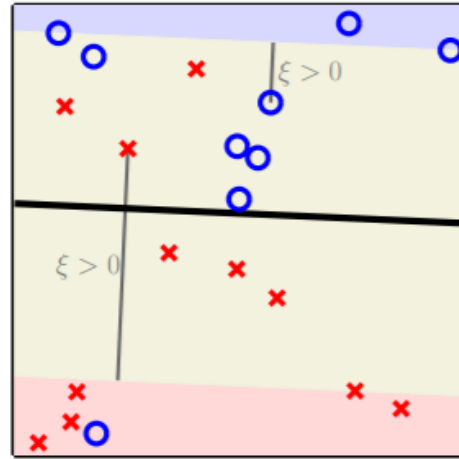
Regularization \rightarrow penalty (complexity)
Soft margin \rightarrow size of margin (robustness).

$C \rightarrow$ Regularization parameter.

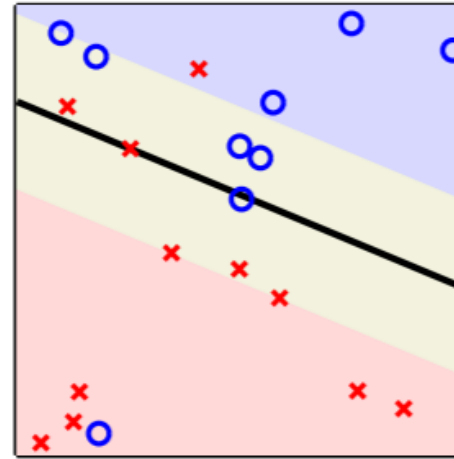
Small $C \rightarrow$ more robustness \rightarrow more regularized
large $C \rightarrow$ less " , $C \rightarrow \infty$ (Prioritizing the hard margin).



Non-Separable Data



$C = 1$

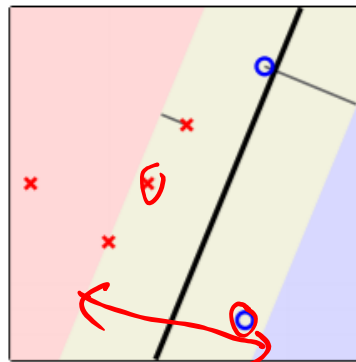


$C = 500$

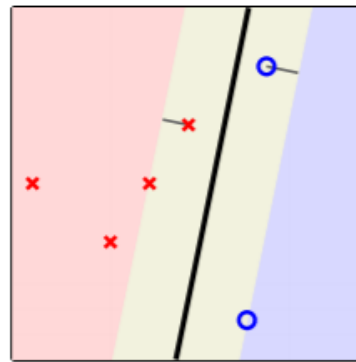
Robustness
VS
In-sample
error.

$$\begin{aligned} & \underset{b, \mathbf{w}, \xi}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ & \text{subject to:} && y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \\ & && \xi_n \geq 0 \quad \text{for } n = 1, \dots, N \end{aligned}$$

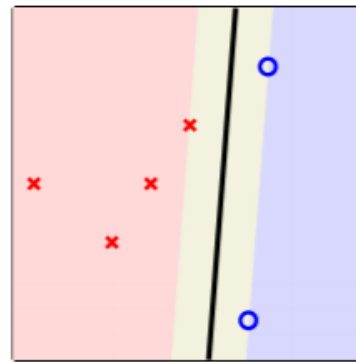
Soft Margin SVM With Separable Data



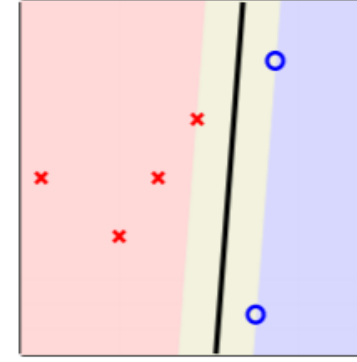
small C



medium C



large C



hard margin

$C \rightarrow \infty$ (big)

QPAL

$$\begin{aligned} &\underset{b, \mathbf{w}, \xi}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ &\text{subject to:} && y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \\ &&& \xi_n \geq 0 \end{aligned}$$

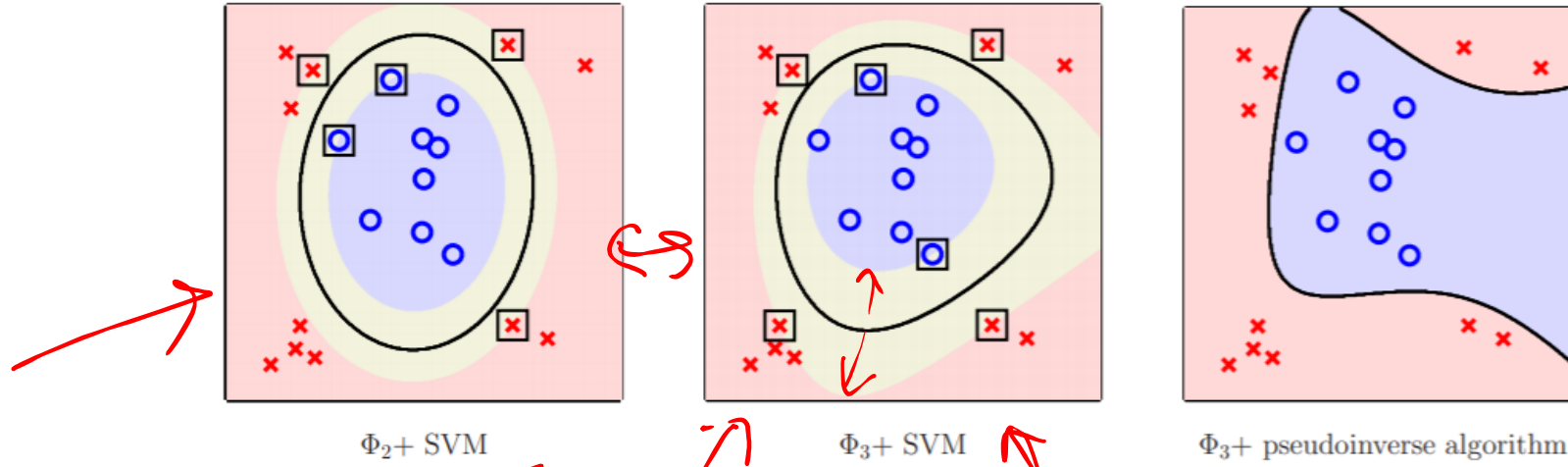
for $n = 1, \dots, N$

QP

Choice of C is
IMPORTANT

Tolerance
of in-sample
error.

Nonlinear Transform and SVM



overfit

Observations:

1. Φ_3 has almost $2\times$ the parameters of Φ_2
2. Φ_3 -SVM does not display significant overfitting compared to Φ_3 -regression
3. #support vectors did not double ✓
4. Can go to higher dimensions if #support vectors stays small or margin stays large ✓

	pseudoinverse regression		SVM	
	linear	nonlinear (ϕ)	linear	nonlinear (ϕ)
overfitting	little	lots	tiny	ok
boundary	linear	complex	linear	complex

Going to Even Higher Dimension

In higher dimension, can control overfitting with # support vectors or margin ρ

What about:

Efficiency? ✓

Infinitely many dimensions?

E_{cv}
 $\hookrightarrow E_{out}$

Thanks!