

# **EECS 445**

## Introduction to Machine Learning

Soft Margin SVMs  
(Efficient) Feature Maps

**Prof. Kutty**

# Today's Agenda

- Recap: Hard Margin SVMs
- Section 1: Soft Margin SVMs
- Section 2: Feature maps

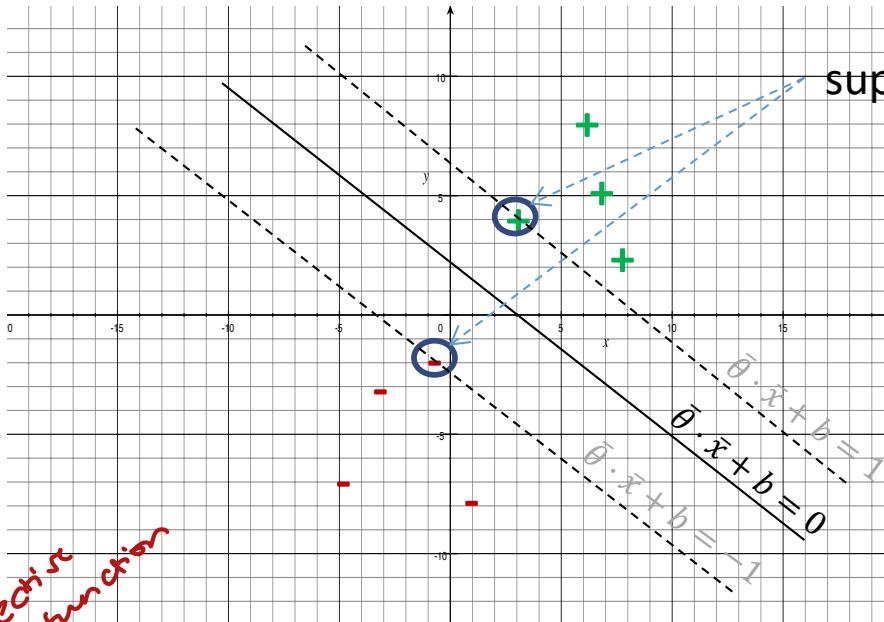
<https://forms.gle/ffiBvNbPjHF8ghi77>



## Recap: Hard Margin SVMs

# Hard Margin SVM

Assuming data are linearly separable



QP: Quadratic Program

$$\min_{\bar{\theta}, b} \frac{\|\bar{\theta}\|^2}{2} \text{ subject to } y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) \geq 1 \text{ for } i \in \{1, \dots, n\}$$

constraints

Linear classifier output by this QP:  $\text{sign}(\bar{\theta} \cdot \bar{x} + b)$

Output of this optimization problem:  $\bar{\theta}, b$

# Linear Separability

What if data are not linearly separable?

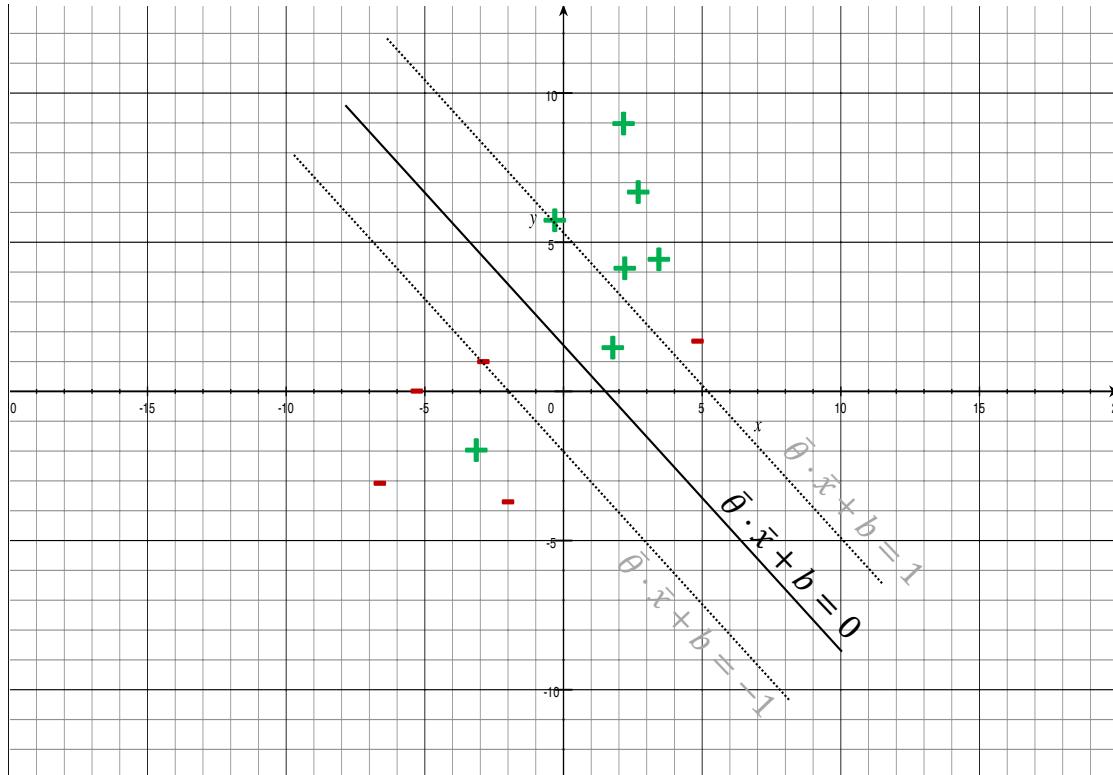
How can we handle such cases?

1. Constraints seem too restrictive
  - Fix the constraints: **Soft-Margin SVMs**
2. Map to a higher dimensional space

# Soft Margin SVMs

# Soft-Margin SVM

Suppose data are *not* linearly separable



$$\min_{\bar{\theta}, b, \xi} \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i \quad \text{slack variables}$$

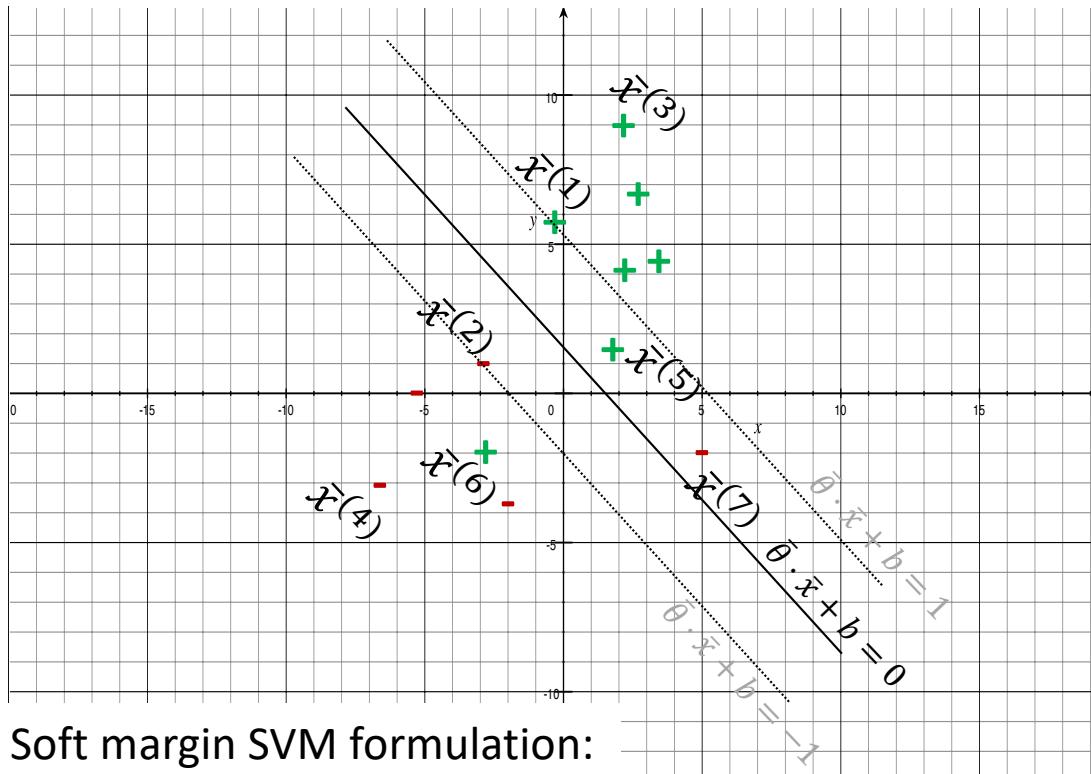
$$\bar{\theta} \in \mathbb{R}^d; b \in \mathbb{R}$$
$$\xi \in \mathbb{R}^n$$

subject to  $y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) \geq 1 - \xi_i$  soft constraints  
for  $i \in \{1, \dots, n\}$

penalty

# Soft-Margin SVM

Suppose data are *not* linearly separable



Soft margin SVM formulation:

CONSTRAINT	PENALTY ( $\xi_i$ )
$y^{(1)}(\bar{\theta} \cdot \bar{x}^{(1)} + b) = 1$	
$y^{(2)}(\bar{\theta} \cdot \bar{x}^{(2)} + b) = 1$	
$y^{(3)}(\bar{\theta} \cdot \bar{x}^{(3)} + b) > 1$	
$y^{(4)}(\bar{\theta} \cdot \bar{x}^{(4)} + b) > 1$	
$y^{(5)}(\bar{\theta} \cdot \bar{x}^{(5)} + b) < 1$	
$y^{(6)}(\bar{\theta} \cdot \bar{x}^{(6)} + b) < 1$	
$y^{(7)}(\bar{\theta} \cdot \bar{x}^{(7)} + b) < 1$	

<https://forms.gle/ffiBvNbPjHF8ghi77>

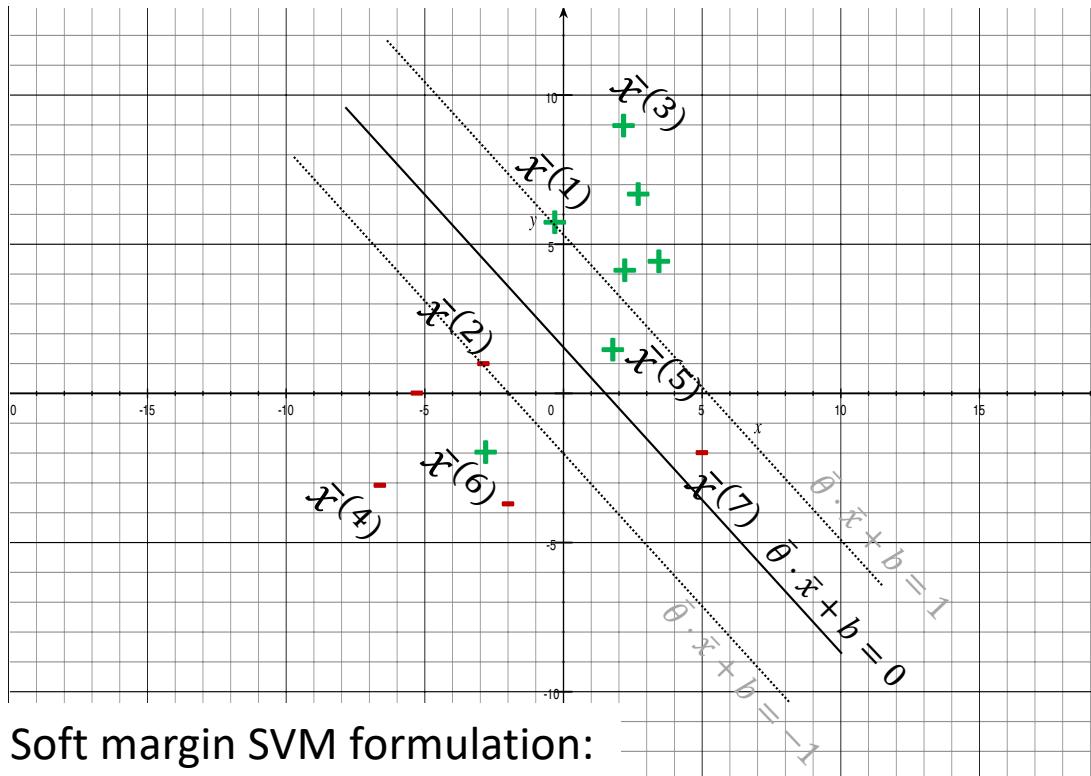
$$\min_{\bar{\theta}, b, \xi} \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i$$

subject to  $y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) \geq 1 - \xi_i$   
for  $i \in \{1, \dots, n\}$  and  $\xi_i \geq 0$



# Soft-Margin SVM

Suppose data are *not* linearly separable



CONSTRAINT	PENALTY ( $\xi_i$ )
$y^{(1)}(\bar{\theta} \cdot \bar{x}^{(1)} + b) = 1$	$\xi_1 = 0$
$y^{(2)}(\bar{\theta} \cdot \bar{x}^{(2)} + b) = 1$	$\xi_2 = 0$
$y^{(3)}(\bar{\theta} \cdot \bar{x}^{(3)} + b) > 1$	$\xi_3 = 0$
$y^{(4)}(\bar{\theta} \cdot \bar{x}^{(4)} + b) > 1$	$\xi_4 = 0$
$y^{(5)}(\bar{\theta} \cdot \bar{x}^{(5)} + b) < 1$	$0 < \xi_5 < 1$
$y^{(6)}(\bar{\theta} \cdot \bar{x}^{(6)} + b) < 1$	$\xi_6 > 1$
$y^{(7)}(\bar{\theta} \cdot \bar{x}^{(7)} + b) < 1$	$\xi_7 > 1$

Soft margin SVM formulation:

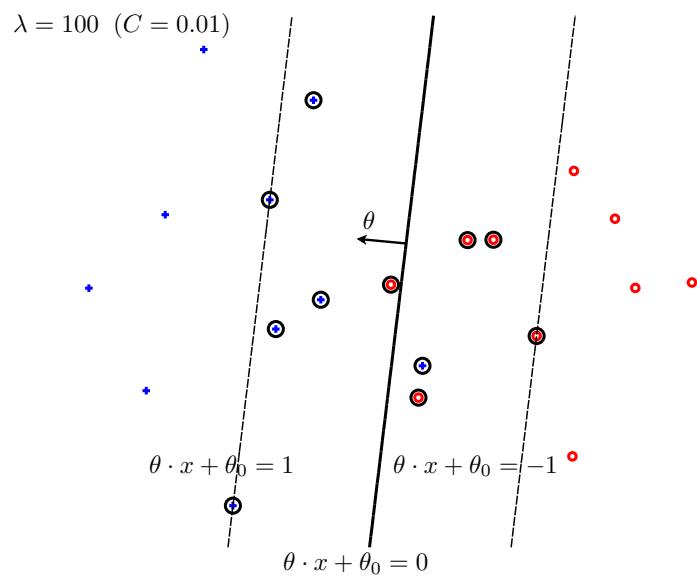
$$\min_{\bar{\theta}, b, \xi} \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i$$

subject to  $y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) \geq 1 - \xi_i$   
for  $i \in \{1, \dots, n\}$  and  $\xi_i \geq 0$

# Soft-Margin SVM

## Soft-margin SVM advantages

- can handle data that are not linearly separable
- reduce effect of outliers and less chances of overfitting



$$\min_{\bar{\theta}, b, \xi} \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i$$

subject to  $y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) \geq 1 - \xi_i$   
for  $i \in \{1, \dots, n\}$  and  $\xi_i \geq 0$

hyperparameter

Image credit: Barzilay & Jaakkola

# Soft-Margin SVM: hyperparameter

$$\min_{\bar{\theta}, b, \xi} \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i$$

$$\text{subject to } y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0 \text{ for } i \in \{1, \dots, n\}$$

Intuitively

as  $C \uparrow$ , penalty on errors/misclassifications  $\uparrow$

as  $C \rightarrow \infty$ , in the limit this is the hard margin SVM

# Soft-Margin SVM hyperparameter

for higher values of C

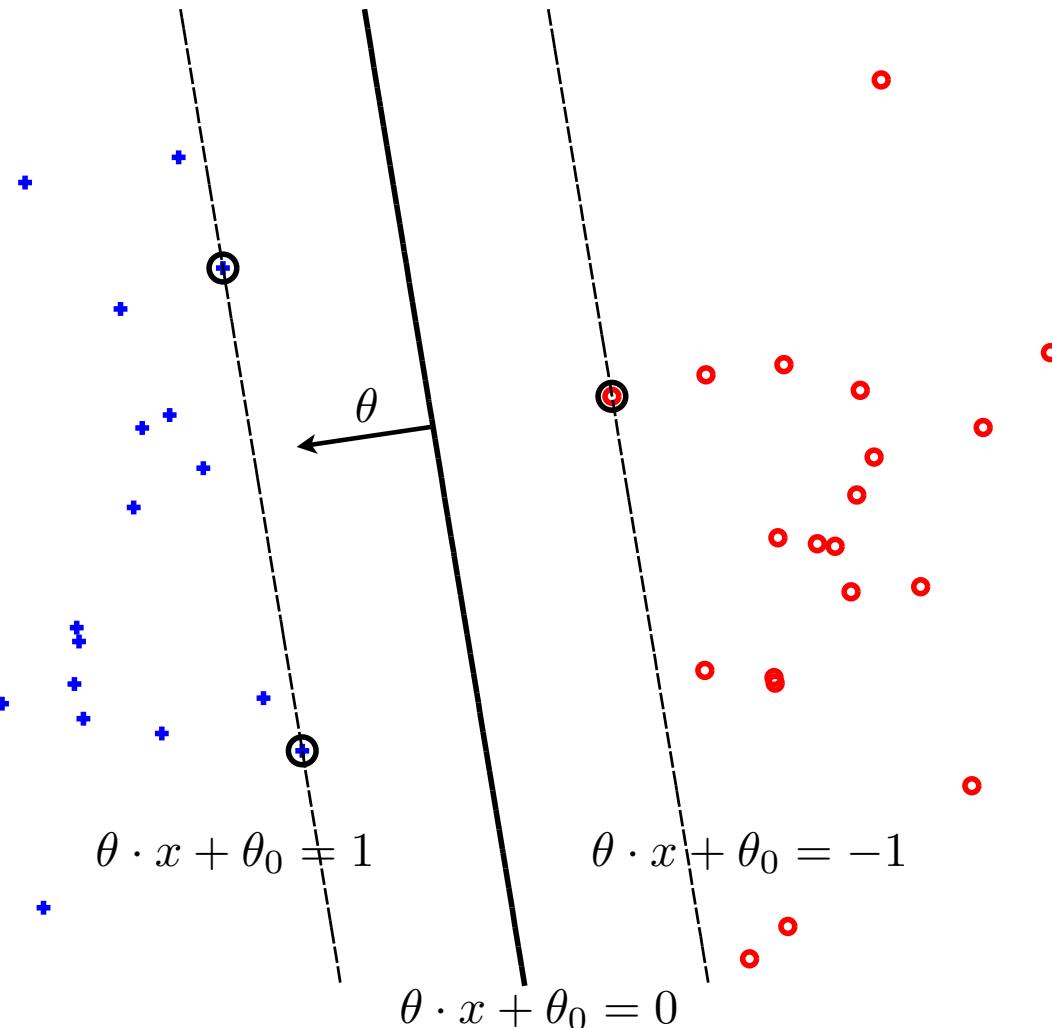


Image credit: Barzilay & Jaakkola

# Soft-Margin SVM hyperparameter

for lower values of C

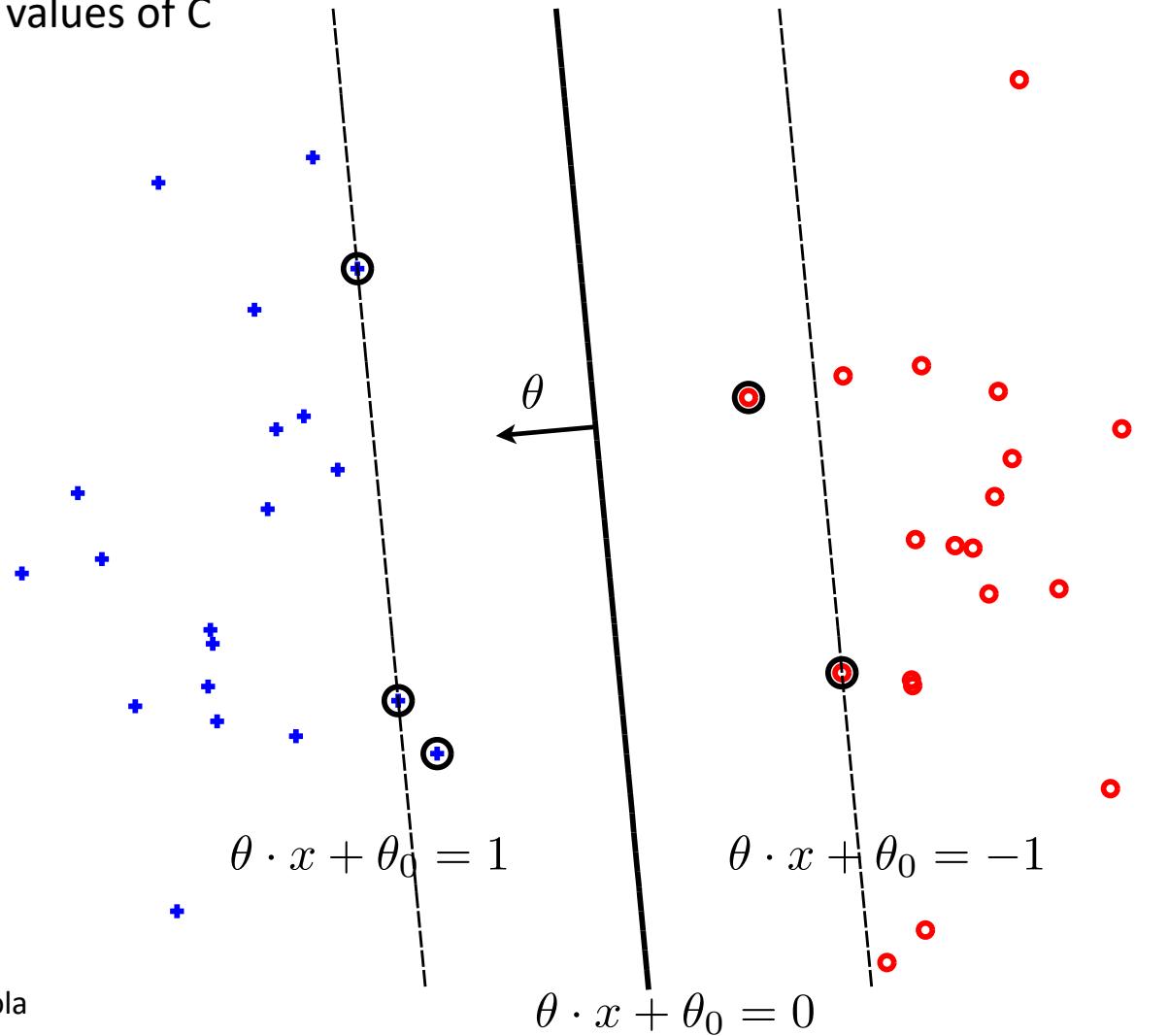


Image credit: Barzilay & Jaakkola

# Soft-Margin SVM: exercise

Claim: Soft margin SVM is an optimization problem with **hinge loss** as objective function and  $\ell_2$ -norm regularizer

$$\min_{\bar{\theta}, b, \xi} \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i \quad \text{subject to } y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

for  $i \in \{1, \dots, n\}$

Hints:

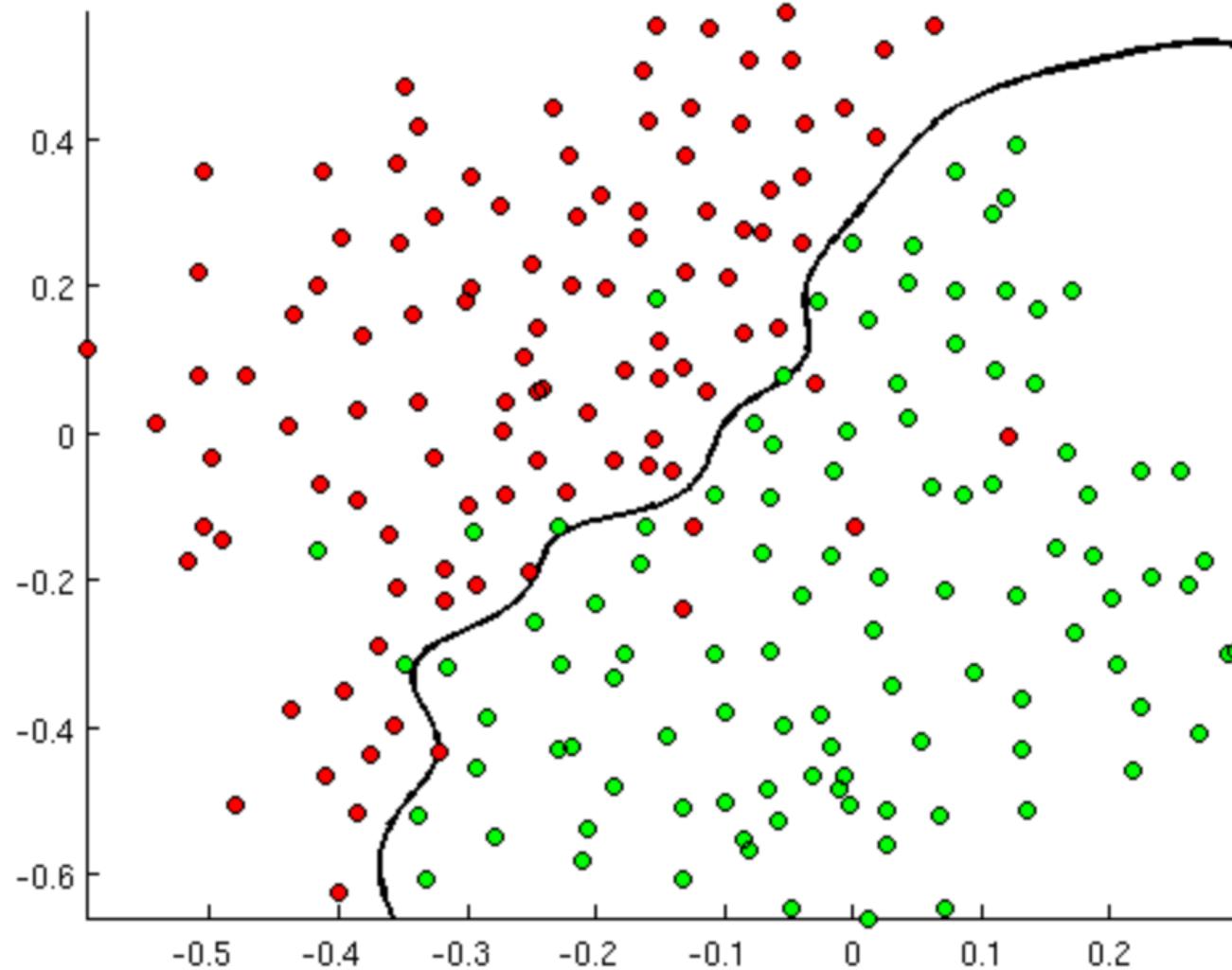
- Write  $\xi_i \geq 1 - y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b)$  and  $\xi_i \geq 0$
- Observe that the objective function includes the terms  $\min_{\xi} \sum_{i=1}^n \xi_i$

$$\min_{\bar{\theta}, b} R_n(\bar{\theta}) + \lambda Z(\bar{\theta})$$

*penalty term  
associated  
model complexity*

*empirical risk with hinge loss*

# Non-linear decision boundaries with SVMs

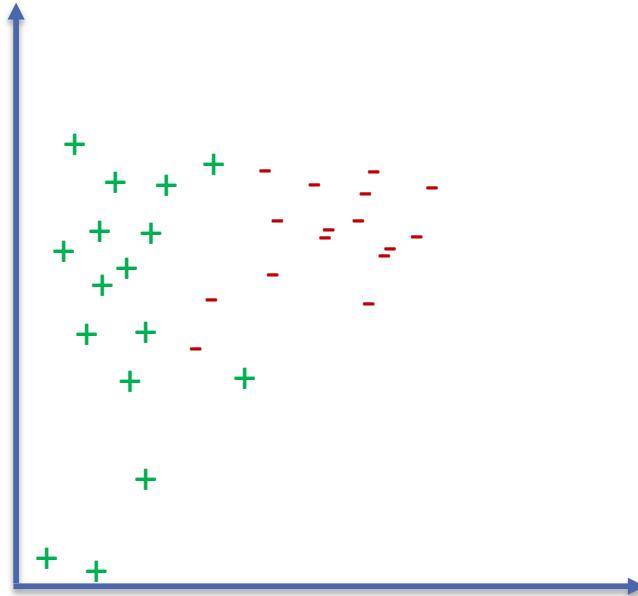


# Feature maps

<https://forms.gle/ffiBvNbPjHF8ghi77>

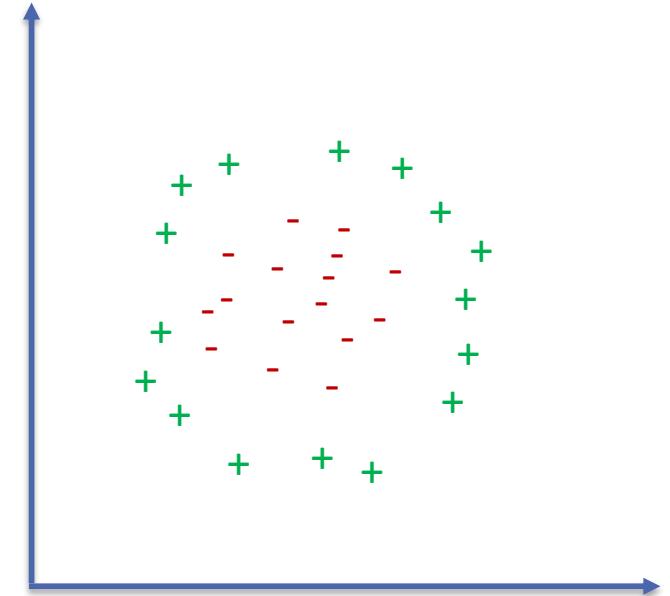


# Data that are not linearly separable



**Idea:** minimize empirical risk with hinge loss using gradient descent

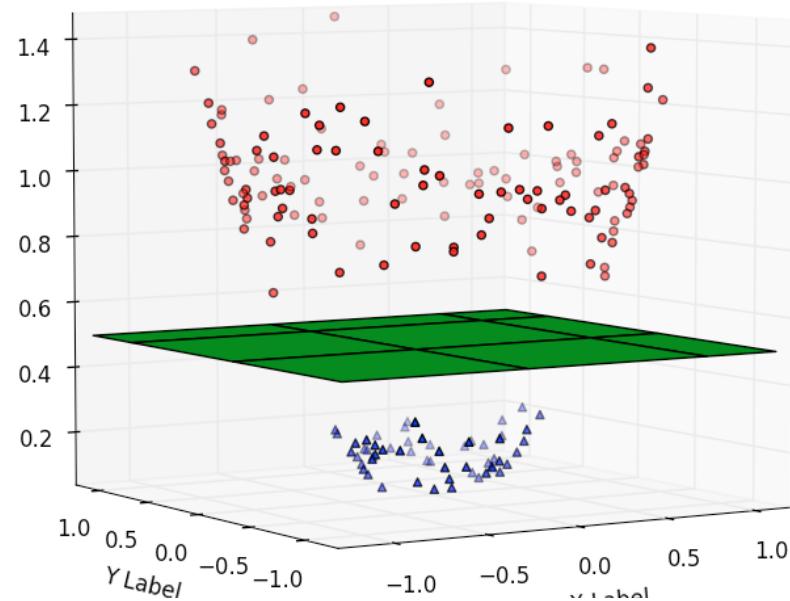
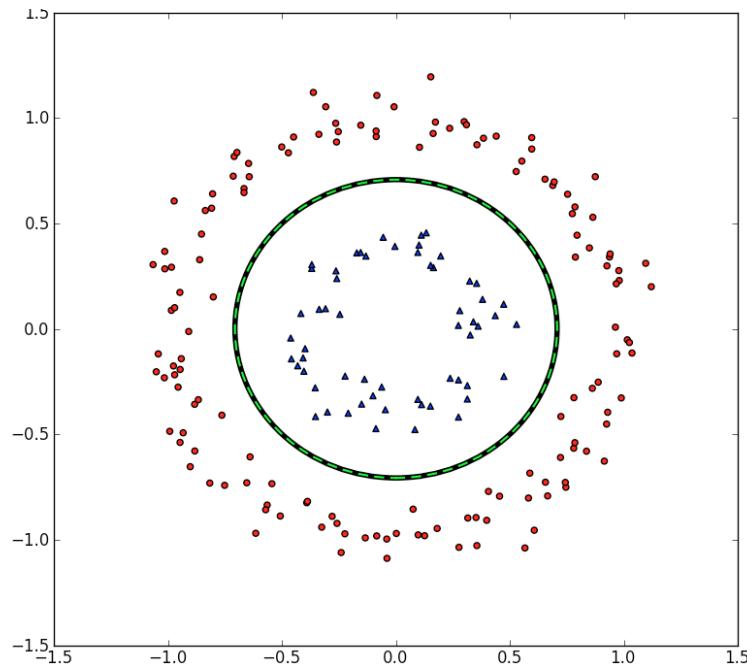
**Idea:** Soft Margin SVMs



**Idea:** feature maps

# Linear classifiers in higher dimensional spaces: idea

image source: <http://www.eric-kim.net>



map data to a higher dim. space in which there exists a separating hyperplane  
(corresponds to a non-linear decision boundary in the original space)

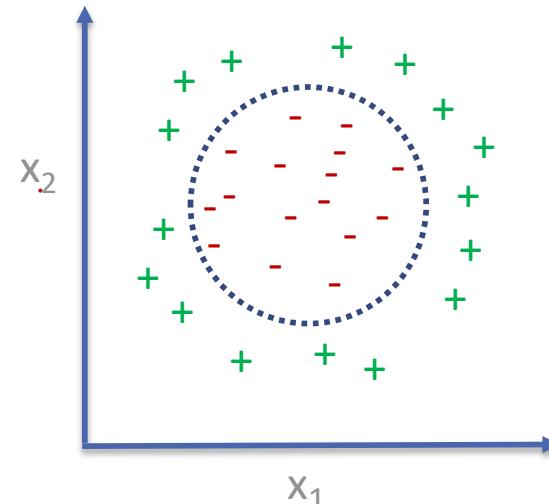
# Linear classifiers in higher dimensional spaces: exercise

$$(x_1 - 2)^2 + (x_2 - 2)^2 = 1^2$$

$$x_1^2 + x_2^2 - 4x_1 - 4x_2 + 7 = 0$$

$$\begin{bmatrix} 1 \\ 1 \\ -4 \\ -4 \end{bmatrix} \cdot \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_1 \\ x_2 \end{bmatrix} + 7 = 0$$

$$\bar{\theta} \cdot \phi(\bar{x}) + b = 0$$



$$\bar{x} \in \mathbb{R}^2$$

'ideal' decision boundary

center at  $(2,2)$

radius = 1

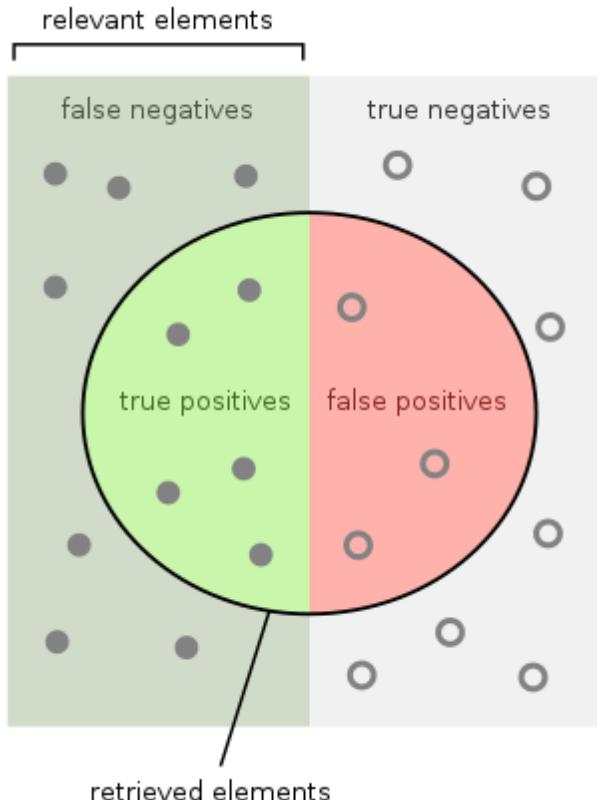
$$\phi(\bar{x}) = [x_1^2 \quad x_2^2 \quad x_1 \quad x_2 \quad 1]$$

$$\bar{\theta} = [1 \quad 1 \quad -4 \quad -4 \quad 7]$$

Linear Classifier  $\text{sign}(\bar{\theta} \cdot \phi(\bar{x}))$  separates  
positive from negative examples

# Evaluating Classifiers

image source <https://commons.wikimedia.org/w/index.php?title=File:Precisionrecall.svg&oldid=697135202>



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

## Common Performance Metrics

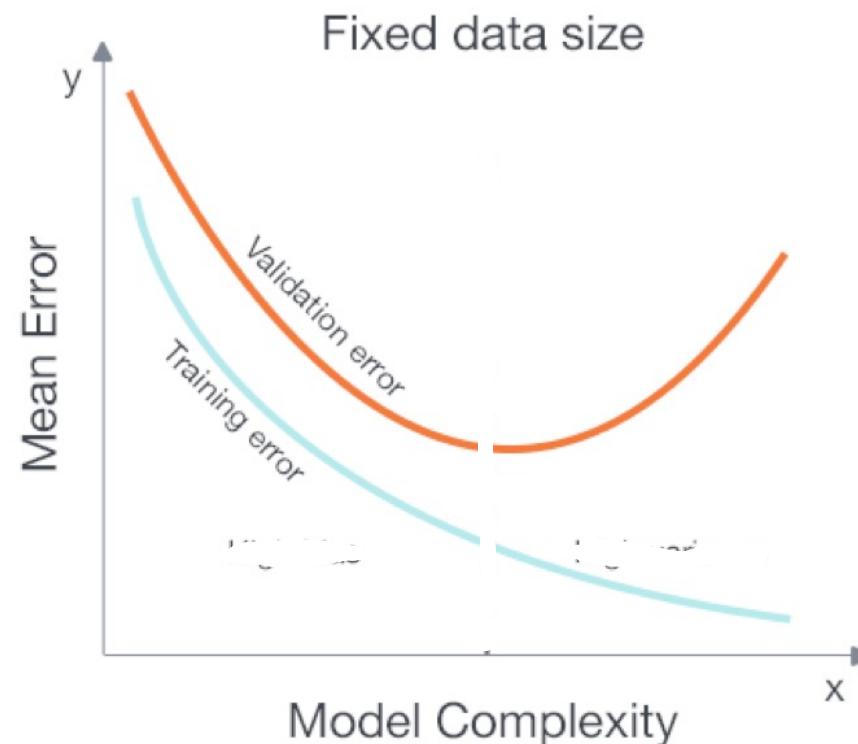
Classification Error:	$\text{FP} + \text{FN}$
Accuracy:	$(\text{TP} + \text{TN}) / \text{N}$
False Positive Rate:	$\text{FP} / (\text{TN} + \text{FP})$
True Positive Rate:	$\text{TP} / (\text{TP} + \text{FN})$
Sensitivity (Recall):	TPR
Precision:	$\text{TP} / (\text{TP} + \text{FP})$
Specificity:	$\text{TN} / (\text{TN} + \text{FP})$

F1 score:  
harmonic mean of  
precision and recall

		Predicted Class	
Actual Class	Pos.	Pos.	Neg.
		#TP	#FN
Neg.	Pos.		
	Neg.	#FP	#TN

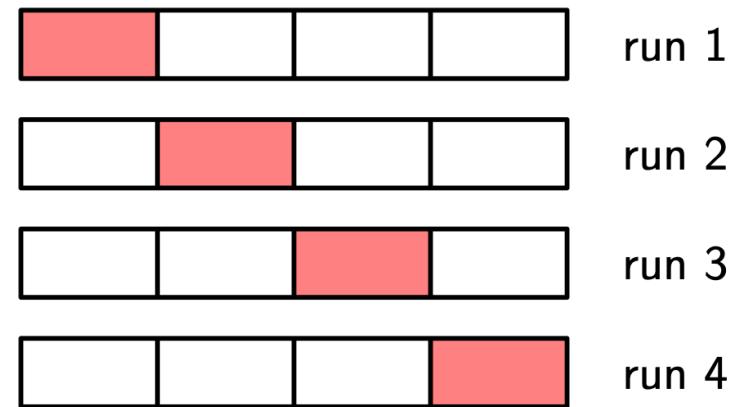
# How to measure generalization error?

- *training dataset*
- *validation dataset*
- *test dataset*



# $k$ -fold Cross Validation

The technique of  $S$ -fold cross-validation, illustrated here for the case of  $S = 4$ , involves taking the available data and partitioning it into  $S$  groups (in the simplest case these are of equal size). Then  $S - 1$  of the groups are used to train a set of models that are then evaluated on the remaining group. This procedure is then repeated for all  $S$  possible choices for the held-out group, indicated here by the red blocks, and the performance scores from the  $S$  runs are then averaged.



**Note:** we say  $k$ -fold Cross Validation

image source: Bishop 2006

# Implications for SVM

$$\min_{\bar{\theta}, b, \xi} \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i$$

subject to  $y^{(i)}(\bar{\theta} \cdot \phi(\bar{x}^{(i)}) + b) \geq 1 - \xi_i$

and  $\xi_i \geq 0$

$$S_n = \{(\bar{x}^{(i)}, y^{(i)})\}_{i=1}^n$$

↓

$$S'_n = \{(\phi(\bar{x}^{(i)}), y^{(i)})\}_{i=1}^n$$
$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^P$$

for  $i \in \{1, \dots, n\}$

**Issue:** potentially *very* inefficient

# Linear Separability

What if data are not linearly separable?

How can we handle such cases?

1. Soft-Margin SVMs
  2. Map to a higher dimensional space
  3. SVM dual and the kernel trick
- ) efficiently

# Support Vector Machines

QP formulation

$$\min_{\bar{\theta}} \frac{\|\bar{\theta}\|^2}{2} \text{ subject to } y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)}) \geq 1 \text{ for } i \in \{1, \dots, n\}$$

**Goal:** rewrite in **dual form**

$$\max_{\bar{\alpha}, \alpha_i \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \bar{x}^{(i)} \cdot \bar{x}^{(j)}$$

# Dual Formulation in General

<https://forms.gle/ffiBvNbPjHF8ghi77>



# Relating the Lagrangian to $f(\bar{w})$

original problem:  $\min_{\bar{w}} f(\bar{w})$  *primal variable*

Lagrangian:  $L(\bar{w}, \bar{\alpha}) = f(\bar{w}) + \sum_{i=1}^n \alpha_i h_i(\bar{w})$  *dual*  $\alpha_i \geq 0$

Define:  $g_p(\bar{w}) = \max_{\bar{\alpha}, \alpha_i \geq 0} L(\bar{w}, \bar{\alpha})$

Claim:  $g_p(\bar{w}) = \begin{cases} f(\bar{w}), & \text{if constraints are satisfied} \\ \infty, & \text{otherwise} \end{cases}$

Case 1: constraints are satisfied

$$\begin{aligned} g_p(\bar{w}) &= \max_{\bar{\alpha}} \left( f(\bar{w}) + \sum_{i=1}^n \alpha_i h_i(\bar{w}) \right) && \text{by definition} \\ &= \max_{\bar{\alpha}} \left( f(\bar{w}) + \alpha_1 h_1(\bar{w}) + \cdots + \alpha_j h_j(\bar{w}) + \cdots + \alpha_n h_n(\bar{w}) \right) \\ &= f(\bar{w}) \end{aligned}$$

$\alpha_i h_i(\bar{w}) = 0$

Case 2: constraints are *not* satisfied

$$\begin{aligned} g_p(\bar{w}) &= \max_{\bar{\alpha}} \left( f(\bar{w}) + \alpha_1 h_1(\bar{w}) + \cdots + \alpha_j h_j(\bar{w}) + \cdots + \alpha_n h_n(\bar{w}) \right) \\ &= \infty \end{aligned}$$

↗  
 ~~$\alpha_j h_j(\bar{w}) > 0$~~

# Primal formulation

original problem:  $\min_{\bar{w}} f(\bar{w}) \quad \text{s. t. } h_i(\bar{w}) \leq 0 \quad \text{for } i = 1, \dots, n$

Lagrangian:  $L(\bar{w}, \bar{\alpha}) = f(\bar{w}) + \sum_{i=1}^n \alpha_i h_i(\bar{w}) \quad \alpha_i \geq 0$

Define:  $g_p(\bar{w}) = \max_{\bar{\alpha}, \alpha_i \geq 0} L(\bar{w}, \bar{\alpha})$

$$g_p(\bar{w}) = \begin{cases} f(\bar{w}), & \text{if constraints are satisfied} \\ \infty, & \text{otherwise} \end{cases}$$

Note that  $\min_{\bar{w}} g_p(\bar{w}) = \min_{\bar{w}} \max_{\bar{\alpha}} L(\bar{w}, \bar{\alpha}) = \min_{\bar{w}} f(\bar{w})$   
if constraints are satisfiable

Primal formulation  $\min_{\bar{w}} \max_{\bar{\alpha}, \alpha_i \geq 0} L(\bar{w}, \bar{\alpha})$

# Primal vs Dual formulation

original problem:  $\min_{\bar{w}} f(\bar{w}) \quad \text{s. t. } h_i(\bar{w}) \leq 0 \quad \text{for } i = 1, \dots, n$

Lagrangian:  $L(\bar{w}, \bar{\alpha}) = f(\bar{w}) + \sum_{i=1}^n \alpha_i h_i(\bar{w}) \quad \alpha_i \geq 0$

Define:  $g_p(\bar{w}) = \max_{\bar{\alpha}, \alpha_i \geq 0} L(\bar{w}, \bar{\alpha})$

$$g_p(\bar{w}) = \begin{cases} f(\bar{w}), & \text{if constraints are satisfied} \\ \infty, & \text{otherwise} \end{cases}$$

Primal formulation

$$\min_{\bar{w}} \max_{\bar{\alpha}, \alpha_i \geq 0} L(\bar{w}, \bar{\alpha})$$

Dual formulation

$$\max_{\bar{\alpha}, \alpha_i \geq 0} \min_{\bar{w}} L(\bar{w}, \bar{\alpha})$$