# LARGE MARGIN CLASSIFIERS

David Kauchak
CS 158 – Fall 2025

1

## Admin

Assignment 5
- Experiments
- Course feedback
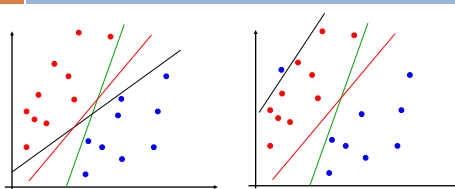
Assignment 6: due *Friday* (10/10)

Midterm: out and due by the end of the day Friday

No class or office hours next Thursday (10/9)

2

## Which hyperplane?



Two main variations in linear classifiers:
- which hyperplane they choose when the data is linearly separable
- how they handle data that is not linearly separable

3

## Linear approaches so far

Perceptron:
- separable:
- non-separable:

Gradient descent:
- separable:
- non-separable:

4

## Linear approaches so far

Perceptron:
- separable:
    - finds **some** hyperplane that separates the data
- non-separable:
    - will continue to adjust as it iterates through the examples
    - final hyperplane will depend on which examples it saw recently
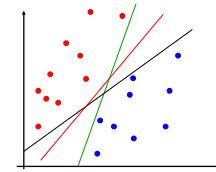
Gradient descent:
- separable and non-separable
    - finds the hyperplane that minimizes the objective function (loss + regularization)
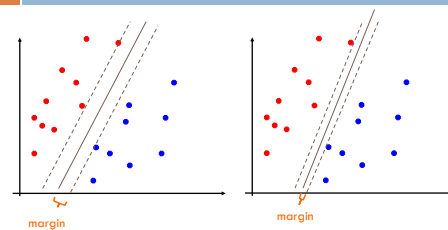
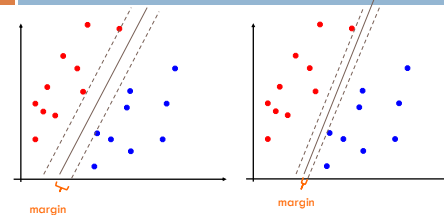Which hyperplane is this?

5

## Which hyperplane would you choose?



6

## Large margin classifiers



margin

Choose the line where the distance to the nearest point(s) is as large as possible

7

## Large margin classifiers



margin

The margin of a classifier is the distance to the closest points of either class
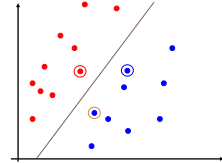
Large margin classifiers attempt to maximize this

8

## Support vectors

For any separating hyperplane, there exist some set of "closest points"
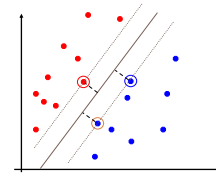
These are called the support vectors

For n dimensions, there will be *at least* n+1 support vectors

9

## Measuring the margin

The margin is the distance to the support vectors, i.e. the "closest points", on either side of the hyperplane

10

## Measuring the margin

negative examples
$$w \cdot x_i + b < 0$$

$$w \cdot x_i + b = 0$$

positive examples
$$w \cdot x_i + b > 0$$

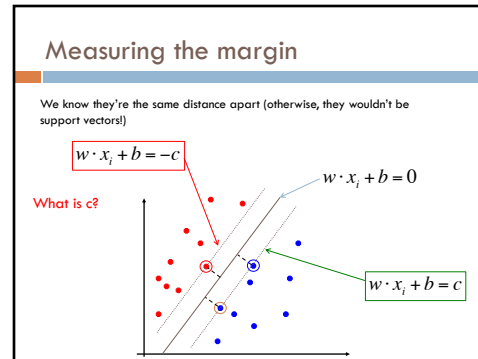11

## Measuring the margin

What are the equations for the margin lines?

negative examples
$$w \cdot x_i + b < 0$$

$$w \cdot x_i + b = 0$$

positive examples
$$w \cdot x_i + b > 0$$

12

## Measuring the margin

We know they're the same distance apart (otherwise, they wouldn't be support vectors!)

$w \cdot x_i + b = -c$

$w \cdot x_i + b = 0$

What is c?

$w \cdot x_i + b = c$

13

## Measuring the margin

Depends! If we scale w, we vary the constant without changing the separating hyperplane

$w \cdot x_i + b = -c$

$w \cdot x_i + b = 0$

$w \cdot x_i + b = c$

14

## Measuring the margin

Depends! If we scale w, we vary the constant without changing the separating hyperplane

$w \cdot x_i + b = -c$

$w \cdot x_i + b = 0$

Larger w result in larger constants

$w \cdot x_i + b = c$

15

## Measuring the margin

Depends! If we scale w, we vary the constant without changing the separating hyperplane

$w \cdot x_i + b = -c$

$w \cdot x_i + b = 0$

Smaller w result in smaller constants

$w \cdot x_i + b = c$

16

## Measuring the margin

For now, let's assume c = 1.

$$w \cdot x_i + b = -1$$

What is this distance?

$$w \cdot x_i + b = 1$$

17

## Distance from the hyperplane

How far away is this point from the hyperplane?

$f_2$

w=(1,2)

$f_1$

(-1,-2)

18

## Distance from the hyperplane

How far away is this point from the hyperplane?

$f_2$

w=(1,2)

$f_1$

$$d = \sqrt{(-1)^2 + (-2)^2} = \sqrt{5}$$

(-1,-2)

19

## Distance from the hyperplane

How far away is this point from the hyperplane?

$f_2$

w=(1,2)

(1,1)

$f_1$

20

## Distance from the hyperplane

How far away is this point from the hyperplane?

w=(1,2)

Is it?

(1,1)

$d(x) = w \cdot x + b$

21

## Distance from the hyperplane

Does that seem right?  What's the problem?

w=(1,2)

(1,1)

$d(x) = w \cdot x + b$

$= w_1 x_1 + w_2 x_2 + b$

$= 1 * 1 + 1 * 2 + 0$

$= 3?$

22

## Distance from the hyperplane

How far away is the point from the hyperplane?

w=(2,4)

(1,1)

$d(x) = w \cdot x + b$

23

## Distance from the hyperplane

How far away is the point from the hyperplane?

w=(2,4)

(1,1)

$d(x) = w \cdot x + b$

$= w_1 x_1 + w_2 x_2 + b$

$= 2 * 1 + 4 * 1$

$= 6?$

24

## Distance from the hyperplane

How far away is this point from the hyperplane?

$f_2$

w=(1,2)

(1,1)

$f_1$

$$d(x) = \frac{w \cdot x + b}{\|w\|}$$

length normalized
weight vectors

25

## Distance from the hyperplane

How far away is this point from the hyperplane?

$f_2$

w=(1,2)

(1,1)

$f_1$

$$d(x) = \frac{w \cdot x + b}{\|w\|}$$

$$= \frac{(w_1 x_1 + w_2 x_2) + b}{\sqrt{5}}$$

$$= \frac{(1*1 + 1*2) + 0}{\sqrt{5}}$$

$$= 1.34$$

26

## Distance from the hyperplane

The magnitude of the weight vector doesn't matter

$f_2$

w=(2,4)

(1,1)

$f_1$

$$d(x) = \frac{w \cdot x + b}{\|w\|}$$

length normalized
weight vectors

27

## Distance from the hyperplane

The magnitude of the weight vector doesn't matter

$f_2$

w=(0.5,1)

(1,1)

$f_1$

$$d(x) = \frac{w \cdot x + b}{\|w\|}$$

length normalized
weight vectors

28

## Measuring the margin

For now, let's just assume c = 1.

$w \cdot x_i + b = -1$

What is this distance?

$w \cdot x_i + b = 1$

29

## Measuring the margin

For now, let's just assume c = 1.

$w \cdot x_i + b = -1$

$$\frac{w \cdot x_i + b}{\|w\|} = \ ?$$

$w \cdot x_i + b = 1$

30

## Measuring the margin

For now, let's just assume c = 1.

$w \cdot x_i + b = -1$

$$\frac{w \cdot x_i + b}{\|w\|} = \frac{1}{\|w\|}$$

$w \cdot x_i + b = 1$

31

## Large margin classifier setup

Select the hyperplane with the largest margin where the points are classified correctly *and outside the margin!*

Setup as a constrained optimization problem:

$$\max_{w,b} \ \text{margin}(w,b)$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 \ \ \forall i \quad \text{what does this say?}$$

32

8

## Large margin classifier setup

Select the hyperplane with the largest margin where the points are classified correctly *and outside the margin!*

Setup as a constrained optimization problem:

$$\max_{w,b} \frac{1}{\|w\|}$$

subject to:

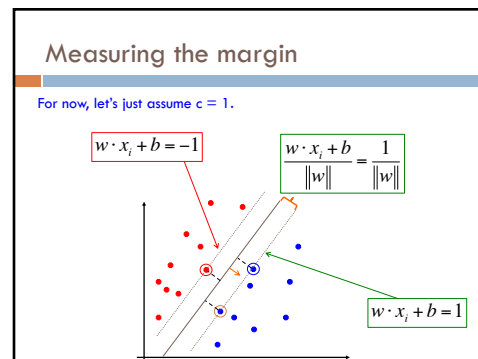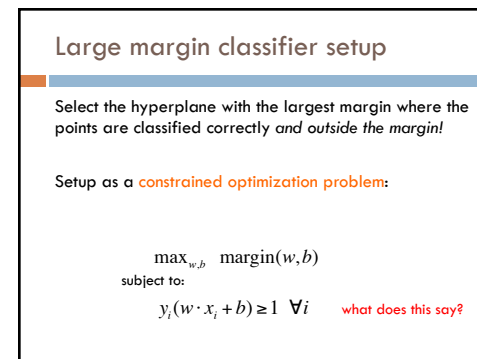$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

33

## Maximizing the margin

$$\min_{w,b} \|w\|$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

**Maximizing the margin is equivalent to minimizing $\|w\|$!**
**(subject to the separating constraints)**

34

## Maximizing the margin

The minimization criterion wants w to be as small as possible

$$\min_{w,b} \|w\|$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

The constraints:
1. make sure the data is separable
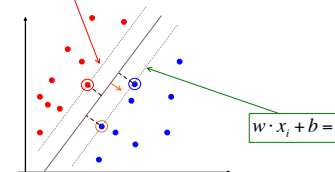2. encourages w to be larger (once the data is separable)

35

## Measuring the margin

For now, let's just assume c = 1.

$$w \cdot x_i + b = -1$$

Claim: it does not matter what c we choose for the SVM problem. Why?

$$w \cdot x_i + b = 1$$



36

9

## Measuring the margin

$w \cdot x_i + b = -c$

What is this distance?

$w \cdot x_i + b = c$

37

## Measuring the margin

$w \cdot x_i + b = -c$

$$\frac{w \cdot x_i + b}{\|w\|} = \frac{c}{\|w\|}$$

$w \cdot x_i + b = c$

38

## Maximizing the margin

$$\min_{w,b} \quad \frac{\|w\|}{c}$$

subject to:

$$y_i(w \cdot x_i + b) \geq c \quad \forall i$$

vs.

What's the difference?

$$\min_{w,b} \quad \|w\|$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

39

## Maximizing the margin

$$\min_{w,b} \quad \frac{\|w\|}{c}$$

subject to:

$$y_i(w \cdot x_i + b) \geq c \quad \forall i$$

Learn the exact same hyperplane just scaled by a constant amount

vs.

$$\min_{w,b} \quad \|w\|$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

Because of this, often see it with c = 1

40

## For those that are curious…

$$\frac{\|w\|}{c} = \frac{\sqrt{w_1^2 + w_2^2 + \ldots + w_m^2 + b^2}}{c}$$

$$= \sqrt{\left(\frac{\sqrt{w_1^2 + w_2^2 + \ldots + w_m^2}}{c}\right)^2}$$

$$= \sqrt{\frac{w_1^2 + w_2^2 + \ldots + w_m^2}{c^2}}$$

$$= \sqrt{\frac{w_1^2}{c^2} + \frac{w_2^2}{c^2} + \ldots + \frac{w_m^2}{c^2}}$$
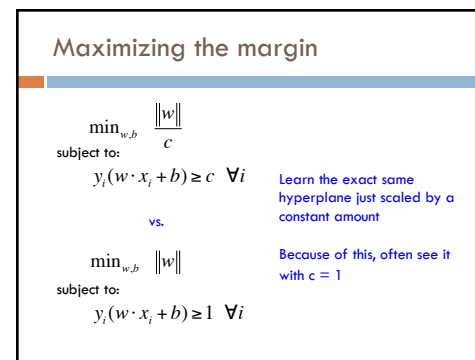
$$= \sqrt{\left(\frac{w_1}{c}\right)^2 + \left(\frac{w_2}{c}\right)^2 + \ldots + \left(\frac{w_m}{c}\right)^2} \qquad \text{scaled version of w}$$

41

## Maximizing the margin: the real problem

$$\min_{w,b} \quad \|w\|^2$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

Why the squared?

42

## Maximizing the margin: the real problem

$$\min_{w,b} \quad \|w\| = \sqrt{\sum_i w_i^2}$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

$$\min_{w,b} \quad \|w\|^2 = \sum_i w_i^2$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

Minimizing $\|w\|$ is equivalent to minimizing $\|w\|^2$

The sum of the squared weights is a convex function!

43

## Support vector machine problem

$$\min_{w,b} \quad \|w\|^2$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

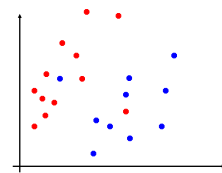This is a version of a quadratic optimization problem

Maximize/minimize a quadratic function

Subject to a set of linear constraints

Many, many variants of solving this problem (we'll see one in a bit)

44

## Soft Margin Classification

$$\min_{w,b} \quad \|w\|^2$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

What about this problem?

45

## Soft Margin Classification

$$\min_{w,b} \quad \|w\|^2$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

We'd like to learn something like this,
but our constraints won't allow it ☹

46

## Slack variables

$$\min_{w,b} \quad \|w\|^2$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

$$\min_{w,b} \quad \|w\|^2 + C\sum_i \varsigma_i$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

slack variables
(one for each example)

What effect does this have?

47

## Slack variables

$$\min_{w,b} \quad \|w\|^2 + C\sum_i \varsigma_i$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

slack penalties

48

## Slack variables

margin

trade-off between margin maximization and penalization

$$\min_{w,b} \quad \|w\|^2 + C\sum_i \varsigma_i$$

penalized by how far from "correct"

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$

$$\varsigma_i \geq 0$$

allowed to make a mistake

49

## Soft margin SVM

$$\min_{w,b} \quad \|w\|^2 + C\sum_i \varsigma_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$

$$\varsigma_i \geq 0$$

Still a quadratic optimization problem!

50

## Demo

https://cs.stanford.edu/~karpathy/svmjs/demo/

51

## Solving the SVM problem



52

## Understanding the Soft Margin SVM



$$\min_{w,b} \ \|w\|^2 + C\sum_i \varsigma_i$$

subject to:
$$y_i(w\cdot x_i + b) \ge 1 - \varsigma_i \ \ \forall i$$
$$\varsigma_i \ge 0$$

Given the optimal solution, w, b:

Can we figure out what the slack penalties are for each point?

53

## Understanding the Soft Margin SVM



What do the margin lines represent wrt w,b?

$$\min_{w,b} \ \|w\|^2 + C\sum_i \varsigma_i$$

subject to:
$$y_i(w\cdot x_i + b) \ge 1 - \varsigma_i \ \ \forall i$$
$$\varsigma_i \ge 0$$

54

## Understanding the Soft Margin SVM

$$w\cdot x_i + b = -1$$

$$w\cdot x_i + b = 1$$



$$\min_{w,b} \ \|w\|^2 + C\sum_i \varsigma_i$$

subject to:
$$y_i(w\cdot x_i + b) \ge 1 - \varsigma_i \ \ \forall i$$
$$\varsigma_i \ge 0$$

Or: $\boxed{y_i(w\cdot x_i + b) = 1}$

55

## Understanding the Soft Margin SVM

$$y_i(w\cdot x_i + b) = 1$$



$$\min_{w,b} \ \|w\|^2 + C\sum_i \varsigma_i$$

subject to:
$$y_i(w\cdot x_i + b) \ge 1 - \varsigma_i \ \ \forall i$$
$$\varsigma_i \ge 0$$

What are the slack values for points outside (or on) the margin AND correctly classified?

56

14

## Understanding the Soft Margin SVM

$$y_i(w \cdot x_i + b) = 1$$

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

0!  The slack variables have to be greater than or equal to zero
and if they're on or beyond the margin then y$_i$(wx$_i$+b) ≥ 1 already

57

## Understanding the Soft Margin SVM

$$y_i(w \cdot x_i + b) = 1$$

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

What are the slack values for points inside the margin
AND classified correctly?

58

## Understanding the Soft Margin SVM

$$y_i(w \cdot x_i + b) = 1$$

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

Difference from the point to the margin.  Which is?

$$\varsigma_i = 1 - y_i(w \cdot x_i + b)$$

59

## Understanding the Soft Margin SVM

$$y_i(w \cdot x_i + b) = 1$$

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

What are the slack values for points that are incorrectly
classified?

60

## Understanding the Soft Margin SVM

$$y_i(w \cdot x_i + b) = 1$$

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$

subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

Which is?

61

## Understanding the Soft Margin SVM

$$y_i(w \cdot x_i + b) = 1$$

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$

subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

"distance" to the hyperplane *plus* the "distance" to the margin

62

## Understanding the Soft Margin SVM

$$y_i(w \cdot x_i + b) = 1$$

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$

subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

"distance" to the hyperplane *plus* the "distance" to the margin

?

63

## Understanding the Soft Margin SVM

$$y_i(w \cdot x_i + b) = 1$$

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i$$

subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

"distance" to the hyperplane *plus* the "distance" to the margin

$$-y_i(w \cdot x_i + b)$$ Why -?

64

16

## Understanding the Soft Margin SVM

$y_i(w \cdot x_i + b) = 1$

$$\min_{w,b} \quad \|w\|^2 + C\sum_i \varsigma_i$$

subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

"distance" to the hyperplane *plus* the "distance" to the margin

$-y_i(w \cdot x_i + b)$          ?

65

## Understanding the Soft Margin SVM

$y_i(w \cdot x_i + b) = 1$

$$\min_{w,b} \quad \|w\|^2 + C\sum_i \varsigma_i$$

subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

"distance" to the hyperplane *plus* the "distance" to the margin
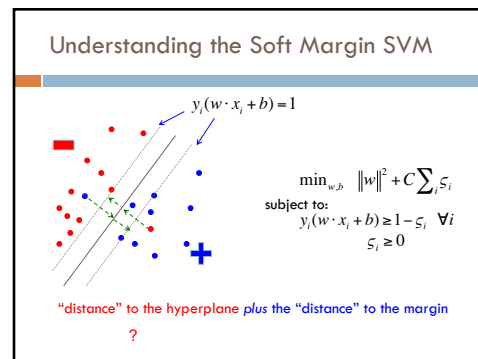
$-y_i(w \cdot x_i + b)$          1

66

## Understanding the Soft Margin SVM

$y_i(w \cdot x_i + b) = 1$

$$\min_{w,b} \quad \|w\|^2 + C\sum_i \varsigma_i$$

subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

"distance" to the hyperplane *plus* the "distance" to the margin

$\varsigma_i = 1 - y_i(w \cdot x_i + b)$

67

## Understanding the Soft Margin SVM

$$\min_{w,b} \quad \|w\|^2 + C\sum_i \varsigma_i$$

subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

$$\varsigma_i = \begin{cases} 0 & \text{if } y_i(w \cdot x_i + b) \geq 1 \\ 1 - y_i(w \cdot x_i + b) & \text{otherwise} \end{cases}$$

68

17

## Understanding the Soft Margin SVM

$$\varsigma_i = \begin{cases} 0 & if \ y_i(w \cdot x_i + b) \geq 1 \\ 1 - y_i(w \cdot x_i + b) & otherwise \end{cases}$$

$$\varsigma_i = \max(0, 1 - y_i(w \cdot x_i + b))$$
$$= \max(0, 1 - yy')$$

Does this look familiar?

69

## Hinge loss!

0/1 loss:     $l(y, y') = 1[yy' \leq 0]$

Hinge:     $l(y, y') = \max(0, 1 - yy')$

Exponential:     $l(y, y') = \exp(-yy')$

Squared loss:     $l(y, y') = (y - y')^2$

70

## Understanding the Soft Margin SVM

$$\min_{w,b} \ \|w\|^2 + C \sum_i \varsigma_i$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

$$\varsigma_i = \max(0, 1 - y_i(w \cdot x_i + b))$$

Do we need the constraints still?

71

## Understanding the Soft Margin SVM

$$\min_{w,b} \ \|w\|^2 + C \sum_i \varsigma_i$$
subject to:
$$y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i$$
$$\varsigma_i \geq 0$$

$$\varsigma_i = \max(0, 1 - y_i(w \cdot x_i + b))$$

$$\min_{w,b} \ \|w\|^2 + C \sum_i \max(0, 1 - y_i(w \cdot x_i + b))$$

Unconstrained problem!

72

18

## Understanding the Soft Margin SVM

$$\min_{w,b} \quad \|w\|^2 + C\sum_i loss_{hinge}(y_i, y_i')$$

**Does this look like something we've seen before?**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^{n} loss(yy') + \lambda \; regularizer(w,b)$$

**Gradient descent problem!**

73

## Soft margin SVM as gradient descent

$$\min_{w,b} \quad \|w\|^2 + C\sum_i loss_{hinge}(y_i, y_i')$$

**multiply through by 1/C and rearrange**
$$\min_{w,b} \quad \sum_i loss_{hinge}(y_i, y_i') + \frac{1}{C}\|w\|^2$$

**let λ=1/C**
$$\min_{w,b} \quad \sum_i loss_{hinge}(y_i, y_i') + \lambda\|w\|^2$$

**What type of gradient descent problem?**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^{n} loss(yy') + \lambda \; regularizer(w,b)$$

74

## Soft margin SVM as gradient descent

One way to solve the soft margin SVM problem is using gradient descent

$$\min_{w,b} \quad \sum_i loss_{hinge}(y_i, y_i') + \lambda\|w\|^2$$

hinge loss

L2 regularization

75

## Gradient descent SVM solver

- □ pick a starting point (w)
- □ repeat until loss doesn't decrease in all dimensions:
  - ■ pick a dimension
  - ■ move a small amount in that dimension towards decreasing loss (using the derivative)

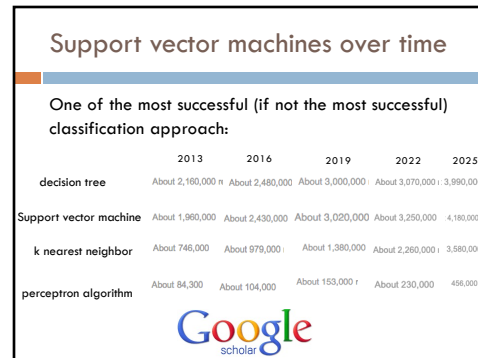$$w_i = w_i - \eta \frac{d}{dw_i}(loss(w) + regularizer(w,b))$$

$$w_j = w_j + \eta \sum_{i=1}^{n} y_i x_i \mathbb{1}[y_i(w \cdot x + b) < 1] - \eta\lambda w_j$$

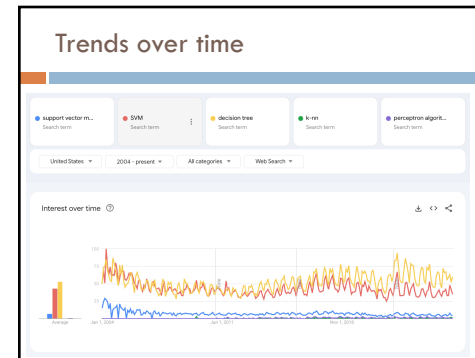hinge loss          L2 regularization

Finds the largest margin hyperplane while allowing for a soft margin

76

19

## Support vector machines over time

One of the most successful (if not the most successful) classification approach:

| | 2013 | 2016 | 2019 | 2022 | 2025 |
|---|---|---|---|---|---|
| decision tree | About 2,160,000 r | About 2,480,000 | About 3,000,000 | About 3,070,000 | 3,990,000 |
| Support vector machine | About 1,960,000 | About 2,430,000 | About 3,020,000 | About 3,250,000 | 4,180,000 |
| k nearest neighbor | About 746,000 | About 979,000 | About 1,380,000 | About 2,260,000 | 3,580,000 |
| perceptron algorithm | About 84,300 | About 104,000 | About 153,000 r | About 230,000 | 456,000 |

Google scholar

77

## Trends over time



78