

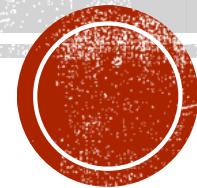
LINEAR MODELS: PART 2

CS 412 Introduction to Machine Learning

Prof. Zheleva

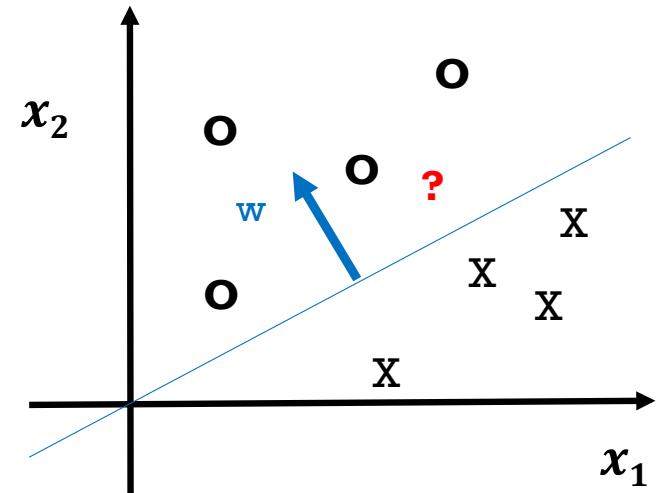
February 20, 2017

Reading Assignment: CML: 7



BINARY CLASSIFICATION VIA HYPERPLANES

- A classifier is a hyperplane that separates positive from negative examples (\mathbf{w}, b)
- At test time, we check on what side of the hyperplane examples fall
$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$
- This is a **linear classifier model**
 - Because the prediction is a linear combination of feature values \mathbf{x}
 - Perceptron is one example



LEARNING A LINEAR CLASSIFIER AS AN OPTIMIZATION PROBLEM

- Structural risk minimization framework
 - Tradeoff between low training error and a solution that is simple

$$\min_{w,b} L(w, b) = \min_{w,b} \left(\sum_{n=1}^N l(y_n, w \cdot x_n + b) + \lambda R(w, b) \right)$$

Variables over which we are optimizing the objective function

Objective function

Loss function
Measures how well classifier fits training data

Regularizing function
Helps with avoiding overfitting
Penalizes complex functions
 λ is a hyperparameter



APPROXIMATING THE 0/1 LOSS WITH CONVEX SURROGATE LOSS FUNCTIONS

- Examples (with $b=0$)

Zero/one: $\ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0]$

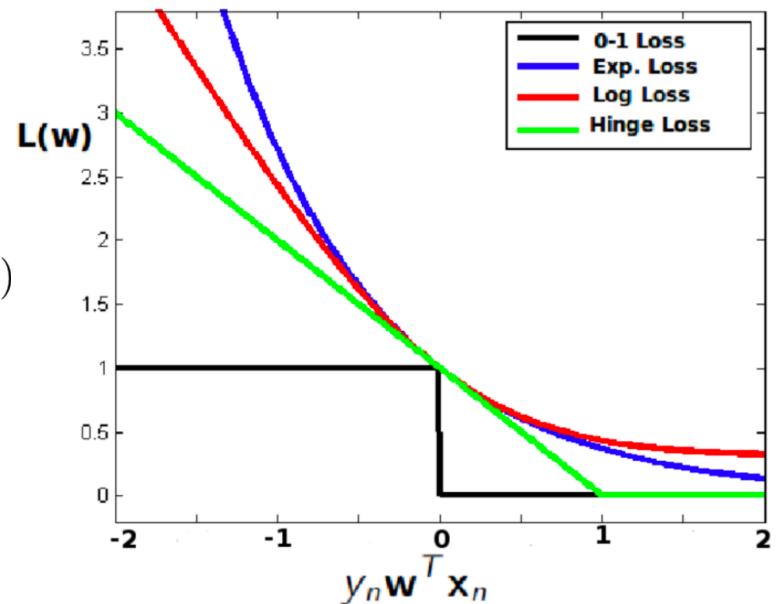
Hinge: $\ell^{(\text{hng})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$

Logistic: $\ell^{(\log)}(y, \hat{y}) = \frac{1}{\log 2} \log (1 + \exp[-y\hat{y}])$

Exponential: $\ell^{(\exp)}(y, \hat{y}) = \exp[-y\hat{y}]$

Squared: $\ell^{(\text{sqr})}(y, \hat{y}) = (y - \hat{y})^2$

- All are convex upper bounds on 0/1 loss
- Some smooth, some more sensitive to outliers



GRADIENT DESCENT ALGORITHM

\mathcal{F} : Objective function to minimize
K: Number of steps
 η_1 : first step size

Algorithm 21 GRADIENTDESCENT($\mathcal{F}, K, \eta_1, \dots$)

```
1:  $z^{(0)} \leftarrow \langle o, o, \dots, o \rangle$                                 // initialize variable we are optimizing
2: for  $k = 1 \dots K$  do
3:    $g^{(k)} \leftarrow \nabla_z \mathcal{F}|_{z^{(k-1)}}$                             // compute gradient at current location
4:    $z^{(k)} \leftarrow z^{(k-1)} - \eta^{(k)} g^{(k)}$                          // take a step down the gradient
5: end for
6: return  $z^{(K)}$ 
```



NOW LET'S CALCULATE GRADIENTS FOR MULTIVARIATE OBJECTIVES

- Consider the following learning objective
 - Loss function: Logistic loss
 - Regularizer: Squared l_2 norm

$$\min_{w,b} L(w,b) = \min_{w,b} \left(\sum_{n=1}^N \exp(-y_n(w \cdot x_n + b)) + \frac{\lambda}{2} \|w\|_2^2 \right)$$

- What do we need to do to run gradient descent?
 - Find gradient for w and b



TODAY: LINEAR MODELS CONTD.

- Subgradients
- Closed form solutions to optimization problems
- (Linear) Support Vector Machines



SUBGRADIENTS

- Problem: some objective functions are not differentiable everywhere
 - Hinge loss (at $y\hat{y}=1$)
 - l_1 norm (when $w_d=0$)
- Solution: subgradient optimization
 - Let's ignore the problem, and just try to apply gradient descent anyway!!
 - We will just differentiate by parts...
 - Where function is differentiable, take the derivative
 - Where function is not differentiable, the subderivative is the set of all possible derivative-like lines



SUBGRADIENT OF HINGE LOSS

$$\partial_w \max\{0, 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b)\} \quad (7.24)$$

$$= \partial_w \begin{cases} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases} \quad (7.25)$$

$$= \begin{cases} \partial_w 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ \partial_w 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases} \quad (7.26)$$

$$= \begin{cases} \mathbf{0} & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ -y_n \mathbf{x}_n & \text{otherwise} \end{cases} \quad (7.27)$$



SUBGRADIENT DESCENT OF HINGE LOSS

Algorithm 22 HINGE_{REGULARIZED}GD($\mathbf{D}, \lambda, MaxIter$)

```
1:  $w \leftarrow \langle o, o, \dots o \rangle$  ,  $b \leftarrow o$                                 // initialize weights and bias
2: for  $iter = 1 \dots MaxIter$  do
3:    $g \leftarrow \langle o, o, \dots o \rangle$  ,  $g \leftarrow o$                                 // initialize gradient of weights and bias
4:   for all  $(x, y) \in \mathbf{D}$  do
5:     if  $y(w \cdot x + b) \leq 1$  then
6:        $g \leftarrow g + y \ x$                                          // update weight gradient
7:        $g \leftarrow g + y$                                          // update bias derivative
8:     end if
9:   end for
10:   $g \leftarrow g - \lambda w$                                          // add in regularization term
11:   $w \leftarrow w + \eta g$                                          // update weights
12:   $b \leftarrow b + \eta g$                                          // update bias
13: end for
14: return  $w, b$ 
```



WHAT IS THE PERCEPTRON OPTIMIZING?

Algorithm 5 PERCEPTRONTRAIN($D, MaxIter$)

```
1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$  // initialize weights
2:  $b \leftarrow 0$  // initialize bias
3: for  $iter = 1 \dots MaxIter$  do
4:   for all  $(x,y) \in D$  do
5:      $a \leftarrow \sum_{d=1}^D w_d x_d + b$  // compute activation for this example
6:     if  $ya \leq 0$  then
7:        $w_d \leftarrow w_d + yx_d$ , for all  $d = 1 \dots D$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:     end if
10:   end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 
```

Loss function is a variant of the hinge loss
 $\max(0, y_n(\mathbf{w}^T \mathbf{x}_n + b))$



CLOSED FORM OPTIMIZATION

- Consider the following learning objective
 - Loss function: Squared error loss
 - Regularizer: Squared l_2 norm

$$\min_{w,b} L(w, b) = \min_{w,b} \left(\frac{1}{2} \sum_{n=1}^N (y_n - (w \cdot x_n + b))^2 + \frac{\lambda}{2} \|w\|_2^2 \right)$$

- Gradient descent would work
 - Can we do better?



OPTIMIZATION AS MATRIX OPERATIONS

- Training data is matrix \mathbf{X} of size $N \times D$
 - $X_{n,d}$: value of d th feature on n th example
- Labels represented as a column vector \mathbf{Y} of dimension N
- Weights represented as a column vector \mathbf{w} of size D
- Matrix-vector product $\mathbf{a} = \mathbf{X}\mathbf{w}$ has dimension N (*bias omitted*)
- Predictions

$$a_n = [\mathbf{X}\mathbf{w}]_n = \sum_d X_{n,d} w_d$$



EXPANDED MATRIX MULTIPLICATION

$$\hat{Y} = [Xw]_n = \sum_d X_{n,d} w_d$$

$$\underbrace{\begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,D} \\ x_{2,1} & x_{2,2} & \dots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,D} \end{bmatrix}}_X \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}}_w = \underbrace{\begin{bmatrix} \sum_d x_{1,d} w_d \\ \sum_d x_{2,d} w_d \\ \vdots \\ \sum_d x_{N,d} w_d \end{bmatrix}}_{\hat{Y}} \approx \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\hat{Y}}$$



USE CALCULUS TO MINIMIZE

$$\min_w L(w) = \min_w \left(\frac{1}{2} \|Xw - Y\|^2 + \frac{\lambda}{2} \|w\|_2^2 \right)$$

- Take gradient
- Set derivative to 0

$$\begin{aligned}\nabla_w \mathcal{L}(w) &= X^\top (Xw - Y) + \lambda w \\ &= X^\top Xw - X^\top Y + \lambda w \\ &= (X^\top X + \lambda I)w - X^\top Y\end{aligned}$$

$$\begin{aligned}(X^\top X + \lambda I)w - X^\top Y &= 0 \\ \iff (X^\top X + \lambda I_D)w &= X^\top Y \\ \iff w &= (X^\top X + \lambda I_D)^{-1} X^\top Y\end{aligned}$$

- Optimal weights computed by a few matrix multiplications and inverses!



MATH REVIEW | MATRIX MULTIPLICATION AND INVERSION

If \mathbf{A} and \mathbf{B} are matrices, and \mathbf{A} is $N \times K$ and \mathbf{B} is $K \times M$ (the inner dimensions must match), then the matrix product \mathbf{AB} is a matrix \mathbf{C} that is $N \times M$, with $\mathbf{C}_{n,m} = \sum_k \mathbf{A}_{n,k} \mathbf{B}_{k,m}$. If \mathbf{v} is a vector in \mathbb{R}^D , we will treat it as a *column vector*, or a matrix of size $D \times 1$. Thus, \mathbf{Av} is well defined if \mathbf{A} is $D \times M$, and the resulting product is a vector \mathbf{u} with $u_m = \sum_d \mathbf{A}_{d,m} v_d$.

Aside from matrix product, a fundamental matrix operation is inversion. We will often encounter a form like $\mathbf{Ax} = \mathbf{y}$, where \mathbf{A} and \mathbf{y} are known and we want to solve for \mathbf{x} . If \mathbf{A} is square of size $N \times N$, then the inverse of \mathbf{A} , denoted \mathbf{A}^{-1} , is also a square matrix of size $N \times N$, such that $\mathbf{AA}^{-1} = \mathbf{I}_N = \mathbf{A}^{-1}\mathbf{A}$. I.e., multiplying a matrix by its inverse (on either side) gives back the identity matrix. Using this, we can solve $\mathbf{Ax} = \mathbf{y}$ by multiplying both sides by \mathbf{A}^{-1} on the left (recall that order matters in matrix multiplication), yielding $\mathbf{A}^{-1}\mathbf{Ax} = \mathbf{A}^{-1}\mathbf{y}$ from which we can conclude $\mathbf{x} = \mathbf{A}^{-1}\mathbf{y}$. Note that not all square matrices are invertible. For instance, the all zeros matrix does not have an inverse (in the same way that $1/0$ is not defined for scalars). However, there are other matrices that do not have inverses; such matrices are called **singular**.



CLOSED FORM VS. GRADIENT DESCENT

- Matrix solution gives exact solution (no iterations!)
- Gradient descent gives progressively better solutions
 - Converges to the optimum at a rate $1/k$
- Is the exact solution always more efficient?
 - One step of gradient descent takes $O(ND)$, K steps take $O(KND)$
 - Matrix solution takes $O(D^3 + D^2N)$ (when $N > D$, then $O(D^2N)$)
 - reconstructing $X^T X$ takes $O(D^2N)$,
 - inversion takes $O(D^3)$
 - final multiplication takes $O(ND)$
 - If $K > D$, then matrix solution is more efficient

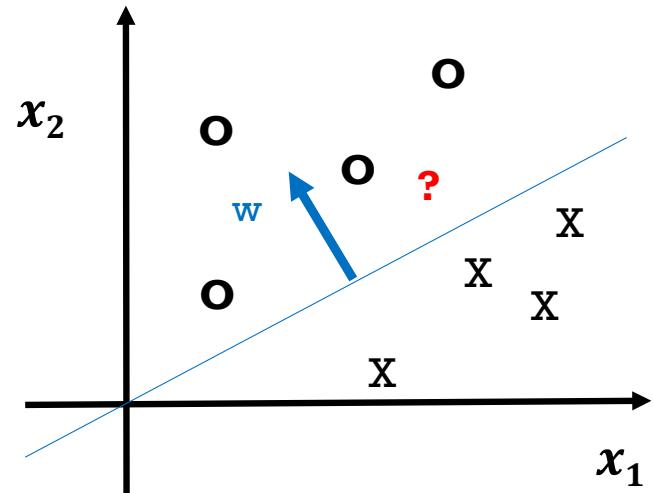


RECAP: CLASSIFICATION VIA HYPERPLANES

- A classifier is a hyperplane that separates positive from negative examples (\mathbf{w}, b)
- At test time, we check on what side of the hyperplane examples fall

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

- This is a **linear classifier model**
 - Because the prediction is a linear combination of feature values \mathbf{x}
 - Perceptron is one example



REVISITING DATASET MARGINS

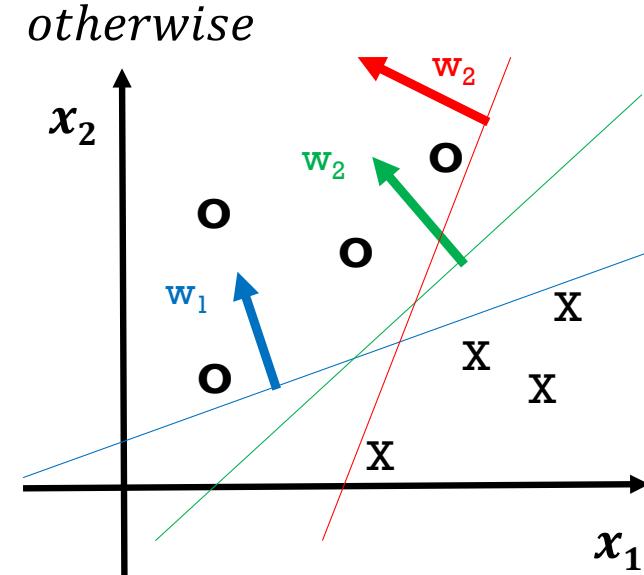
- Hyperplane margin on D : Distance between the hyperplane (w, b) and the nearest point in D

$$\text{margin}(D, w, b) = \begin{cases} \min_{(x,y) \in D} y(w \cdot x + b) & \text{if } w \text{ separates } D \\ -\infty & \text{otherwise} \end{cases}$$

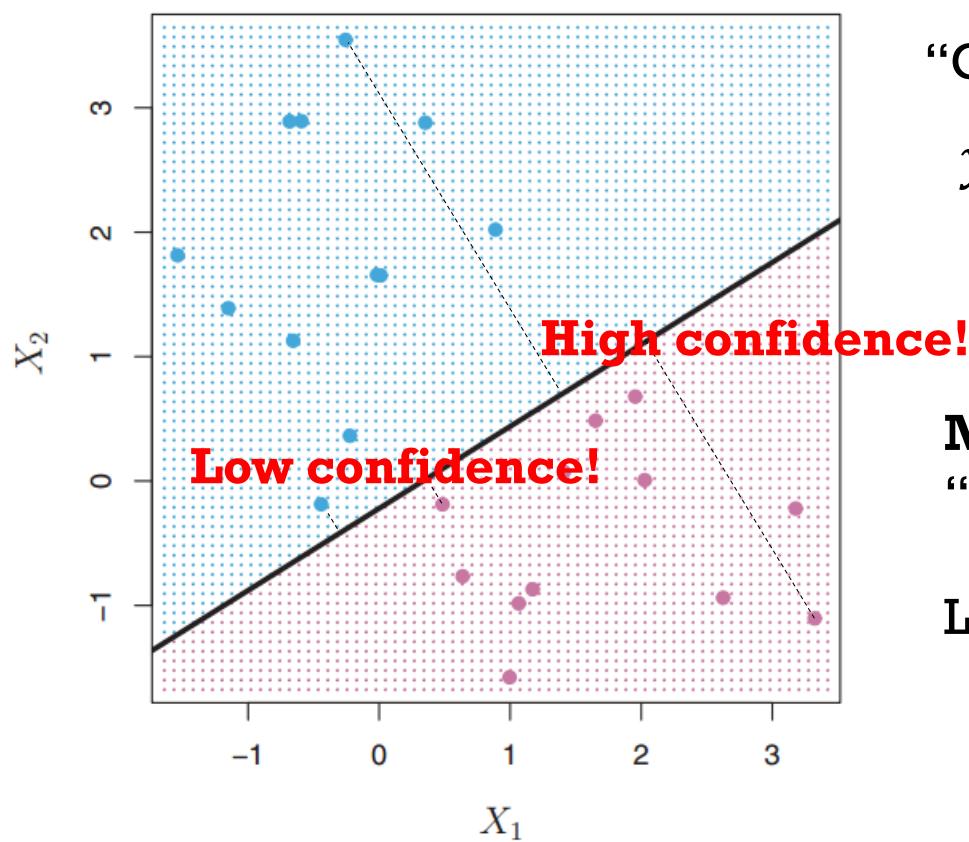
- Dataset margin: Largest attainable margin on D

$$\text{margin}(D) = \sup_{w,b} \text{margin}(D, w, b)$$

- Which hyperplane is best?



MARGIN AS “CONFIDENCE” OF PREDICTIONS



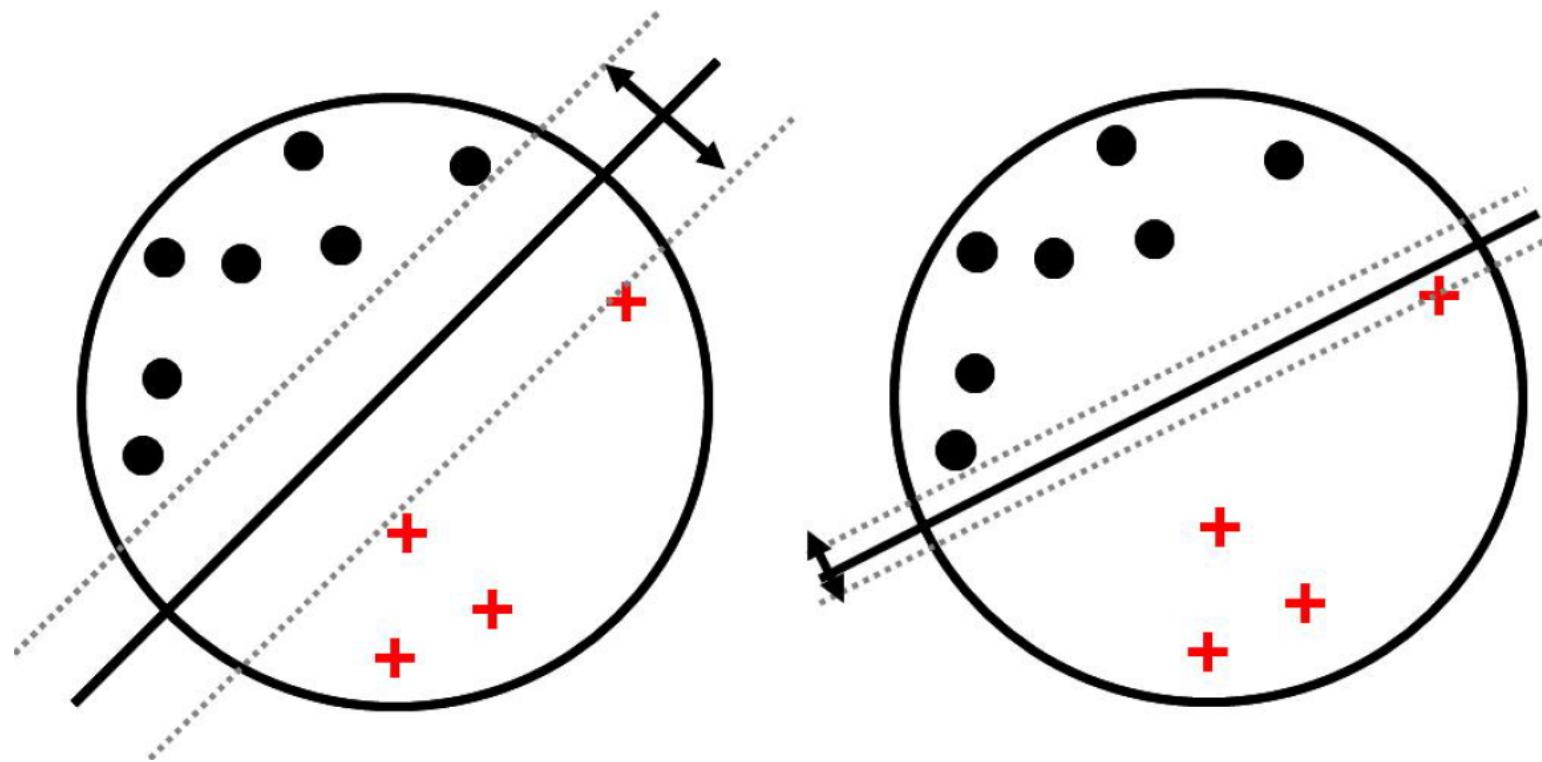
“Confidence” =

$$y_n(w_1x_{n1} + \dots w_dx_{nd} + b)$$

Margin is the minimum
“confidence”

Let's maximize this!



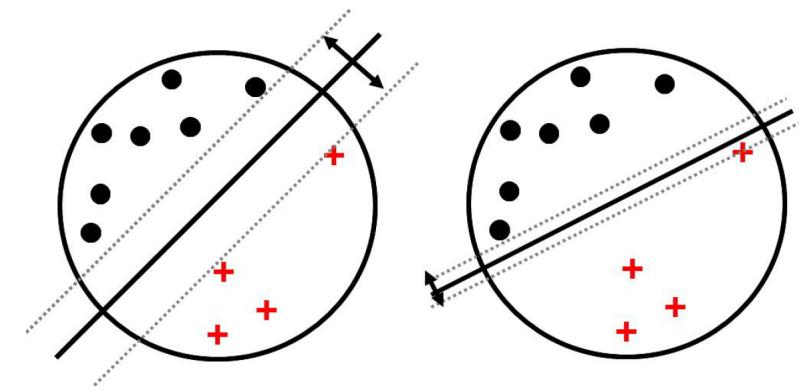


SUPPORT VECTOR MACHINES

- Goal: find a hyperplane with largest possible margin
- Constrained optimization problem:

$$\min_{w,b} \frac{1}{\gamma(w, b)}$$

Subject to $y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 \ (\forall n)$



- What's the problem with this formulation?



HARD-MARGIN VS SOFT-MARGIN SVM

- **Hard-margin SVM:** if *data not linearly separable*, feasible region is empty
 - No parameters can satisfy the constraints
- **Soft-margin SVM:** use slack parameters

$$\min_{w,b,\xi} \underbrace{\frac{1}{\gamma(w,b)}}_{\text{Large margin}} + C \sum_n \xi_n$$

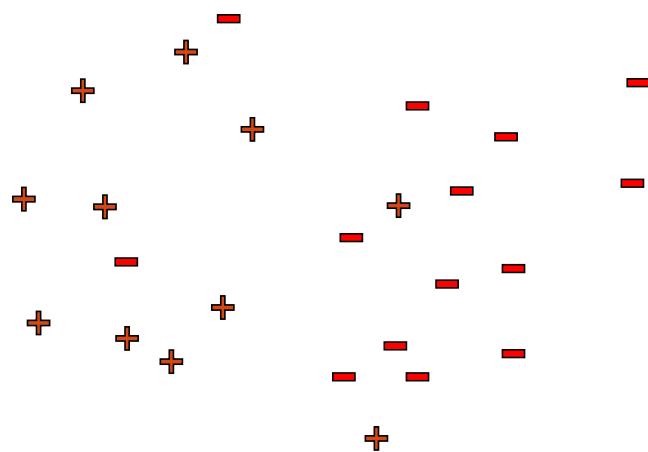
Small slack:
Amount you need to pay, so a point is not misclassified

$$\begin{aligned} \text{Subject to } y_n(w \cdot x_n + b) &\geq 1 - \xi_n \quad (\forall n) \\ \xi_n &\geq 0 \end{aligned}$$



“SUPPORT VECTORS” OF SVM

- Which examples influence the margin and decision boundaries?

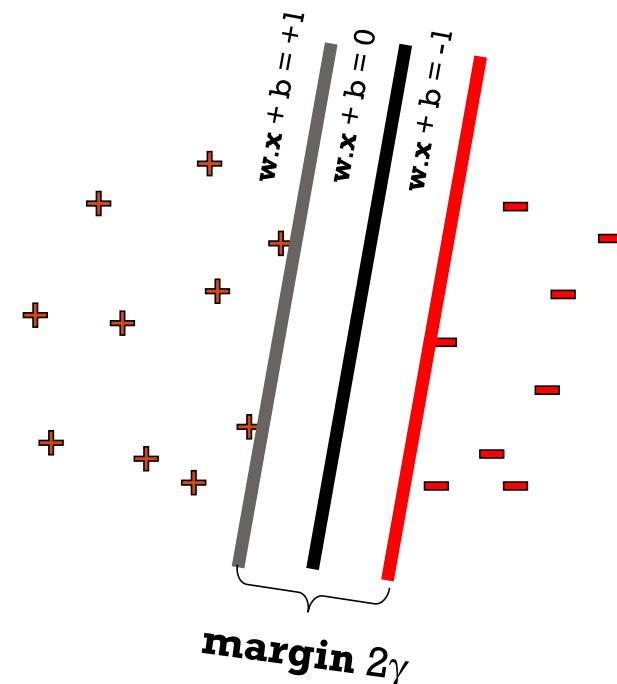


SVM OPTIMIZATION

- Need to be able to measure margin size
 - Hard-margin SVM: optimal hyperplane halfway between closest positive and negative points

$$d^+ = \frac{1}{\|w\|} w \cdot x^+ + b - 1$$

$$d^- = -\frac{1}{\|w\|} w \cdot x^- - b + 1$$



COMPUTING THE MARGIN

$$\begin{aligned}d^+ &= \frac{1}{\|w\|} w \cdot x^+ + b - 1 & \gamma &= \frac{1}{2} [d^+ + d^-] \\d^- &= -\frac{1}{\|w\|} w \cdot x^- - b + 1 & &= \frac{1}{2} \left[\frac{1}{\|w\|} w \cdot x^+ + b - 1 - \frac{1}{\|w\|} w \cdot x^- - b + 1 \right] \\&&&= \frac{1}{2} \left[\frac{1}{\|w\|} w \cdot x^+ - \frac{1}{\|w\|} w \cdot x^- \right] \\&&&= \frac{1}{2} \left[\frac{1}{\|w\|} (+1) - \frac{1}{\|w\|} (-1) \right] \\&&&= \frac{1}{\|w\|}\end{aligned}$$

Maximizing the margin is equivalent to minimizing the l_2 norm!



SVM OPTIMIZATION

- Maximizing the margin is equivalent to minimizing the l_2 norm

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

Large margin Small slack:

Does this remind
you of something
you have seen?

Subject to $y_n(w \cdot x_n + b) \geq 1 - \xi_n \quad (\forall n)$

$$\xi_n \geq 0 \quad (\forall n)$$

- If you were given w and b , can you recover the slacks?



SVM OPTIMIZATION

- Recovering slacks when given \mathbf{w} and b

$$\xi_n = \begin{cases} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 \\ 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases}$$

Does this remind
you of something
you have seen?

- SVM as unconstrained optimization problem

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n l^{(hinge)}(y_n(\mathbf{w} \cdot \mathbf{x}_n + b))$$



Large margin Small slack:



SUMMARY

- Subgradient descent
 - When functions are not differentiable
- Closed form solutions
- Support vector machines



ANNOUNCEMENTS

- Homework
 - Academic honesty expectation



ACKNOWLEDGEMENTS

- These slides use materials by Marine Carpuat and Brian Ziebart

