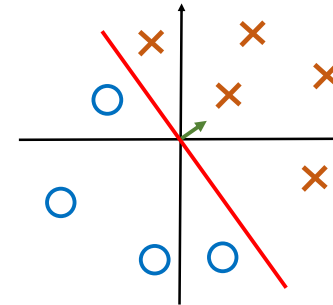


Support Vector Machines

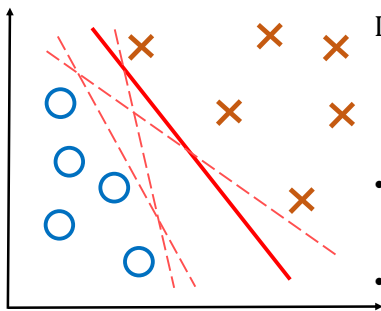
CISC 5800
Professor Daniel Leeds

Separating boundary, defined by \mathbf{w}



- Separating **hyperplane** splits **class 0** and **class 1**
- Plane is defined by line \mathbf{w} perpendicular to plane
- Is data point \mathbf{x} in class 0 or class 1? $\mathbf{w}^T \mathbf{x} + b > 0$ class **1**
 $\mathbf{w}^T \mathbf{x} + b < 0$ class **0**

But, where do we place the boundary?

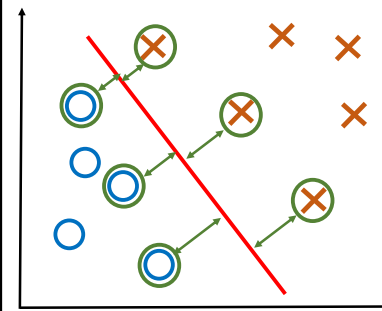


Logistic classifier:

$$LL(y|x; \mathbf{w}): \sum_i (y^i - 1) \mathbf{w}^T \mathbf{x}^i - \log(1 + e^{-\mathbf{w}^T \mathbf{x}^i})$$

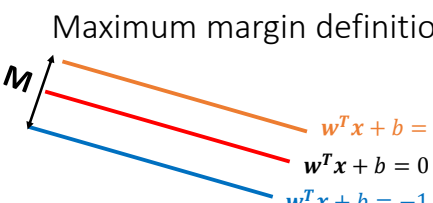
- Each data point \mathbf{x}^i considered for boundary \mathbf{w}
- Outlier data pulls boundary towards it

Max margin classifiers



- Focus on boundary points
- Find largest margin between boundary points on both sides
- Works well in practice
- We can call the boundary points **"support vectors"**

Maximum margin definitions



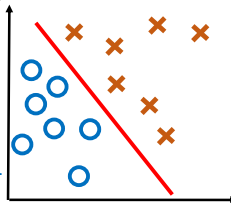
Classify as +1 if $w^T x + b \geq 1$
 Classify as -1 if $w^T x + b \leq -1$
 Undefined if $-1 < w^T x + b < 1$

- M is the margin width
- x^+ is a +1 point closest to boundary, x^- is a -1 point closest to boundary
- $x^+ = \lambda w + x^-$
- $|x^+ - x^-| = M$

$$M = \frac{2}{\sqrt{w^T w}}$$

maximize M minimize $w^T w$

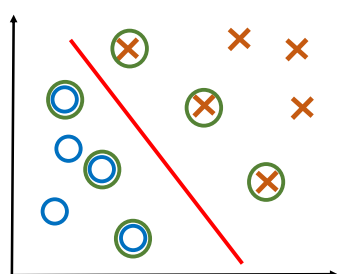
Support vector machine (SVM) optimization



$\operatorname{argmin}_w w^T w$
 subject to
 $w^T x + b \geq 1$ for x in class 1
 $w^T x + b \leq -1$ for x in class -1

$$\operatorname{argmin}_w w^T w + \left(\sum_{i \in +1} \lambda_i (1 - (w^T x^i + b)) \right) +$$

Alternate SVM formulation

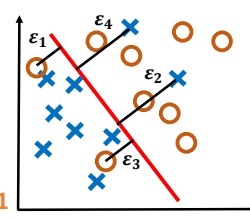


$$w = \sum_i \alpha^i x^i y^i$$

Support vectors x_i have $\alpha_i > 0$
 y_i are the data labels +1 or -1

$$\alpha^i \geq 0 \quad \forall i \quad \sum_i \alpha^i y^i = 0$$

Support vector machine (SVM) optimization with slack variables



What if data not **completely** linearly separable?

$\operatorname{argmin}_{w,b} w^T w + C \sum_i \varepsilon^i$
 subject to
 $w^T x + b \geq 1 - \varepsilon^i$ for x in class 1
 $w^T x + b \leq -1 + \varepsilon^i$ for x in class -1
 $\varepsilon^i \geq 0 \quad \forall i$

Each error ε^i is penalized based on distance from separator

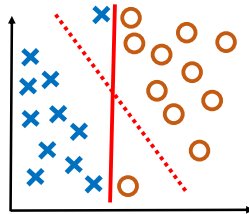
Support vector machine (SVM) optimization with slack variables

Example: Linearly separable but with narrow margins

$$\operatorname{argmin}_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i$$

subject to

$$\begin{aligned} \mathbf{w}^T \mathbf{x} + b &\geq 1 - \varepsilon_i && \text{for } \mathbf{x} \text{ in class 1} \\ \mathbf{w}^T \mathbf{x} + b &\leq -1 + \varepsilon_i && \text{for } \mathbf{x} \text{ in class -1} \\ \varepsilon_i &\geq 0 \quad \forall i \end{aligned}$$



14

Hyper-parameters for learning

$$\operatorname{argmin}_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i$$

Optimization constraints: C influences tolerance for label errors versus narrow margins

$$w_j \leftarrow w_j + \varepsilon x_j^i (y^i - g(\mathbf{w}^T \mathbf{x}^i)) - \frac{w_j}{\lambda}$$

Gradient ascent:

- ε influences effect of individual data points in learning
- T number of training examples, L number of loops through data – balance learning and over-fitting

Regularization: λ influences the strength of your prior belief

15

Parameter counts

Each data point \mathbf{x}^i has N features (presuming classify with $\mathbf{w}^T \mathbf{x}^i + b$)

Separator: \mathbf{w} and b

- N elements of \mathbf{w} , 1 value for b : $N+1$ parameters **OR**
- t support vectors $\rightarrow t$ non-zero α^i , 1 value for b : $t+1$ parameters

16

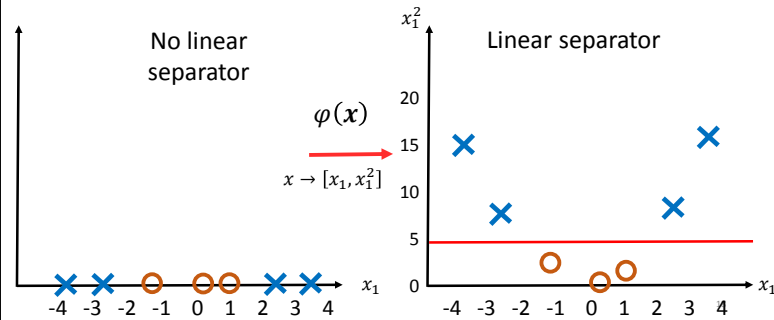
Binary $\rightarrow M$ -class classification

- Learn boundary for class m vs all other classes
 - Only need $M-1$ separators for M classes – M^{th} class is for data outside of classes 1, 2, 3, ..., $M-1$
- Find boundary that gives highest margin for data points \mathbf{x}^i

17

Classifying with additional dimensions

Note: More dimensions makes it easier to separate T
training points: training error minimized, may risk over-fit



Quadratic mapping function (math) $w^T x^k + b = \sum_i \alpha^i y^i (x^i)^T x^k + b$

$x_1, x_2, x_3, x_4 \rightarrow x_1, x_2, x_3, x_4, x_1^2, x_2^2, \dots, x_1 x_2, x_1 x_3, \dots, x_2 x_4, x_3 x_4$

N features $\rightarrow N + N + \frac{N \times (N-1)}{2} \approx N^2$ features

N^2 values to learn for w in higher-dimensional space

Or, observe: $(v^T x + 1)^2 = v_1^2 x_1^2 + \dots + v_N^2 x_N^2 + v_1 v_2 x_1 x_2 + \dots + v_{N-1} v_N x_{N-1} x_N + v_1 x_1 + \dots + v_N x_N$

v with N elements
operating in quadratic
space

19