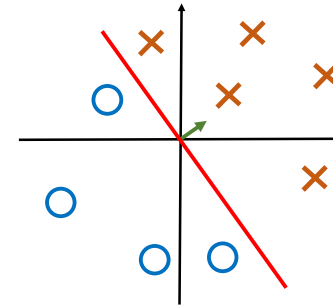


Support Vector Machines

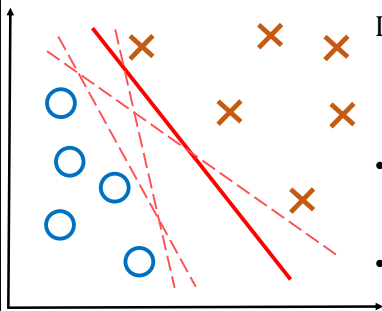
CISC 5800
Professor Daniel Leeds

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:

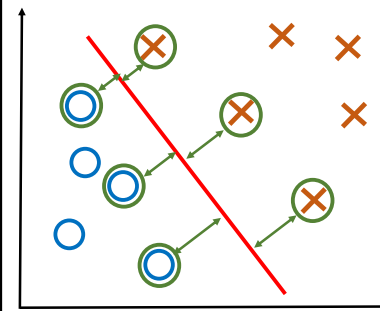
$$LL(y|x; w): \sum_i ((1 - y^i) \ln(1 - g(x^i; w)) + y^i \ln(g(x^i; w)))$$

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

Adjusted
Log Likelihood
expression

3

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**"

4

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 $w^T w$

λ derivation

Optional extra math

- $w^T x^- + b = -1$
- $w^T x^+ + b = +1$
- $x^+ = \lambda w + x^-$

- $w^T x^+ + b = +1$
- $w^T (\lambda w + x^-) + b = +1$
- $\lambda w^T w + w^T x^- + b = +1$
- $\lambda w^T w - 1 - b + b = +1$
- $\lambda = \frac{2}{w^T w}$

M derivation

Optional extra math

- $w^T x^- + b = -1$
- $w^T x^+ + b = +1$
- $x^+ = \lambda w + x^-$
- $|x^+ - x^-| = M$

- $M = |\lambda w + x^- - x^-| = |\lambda w| = \lambda |w|$
- $M = \lambda \sqrt{w^T w}$
- $M = \frac{2}{w^T w} \sqrt{w^T w} = \frac{2}{\sqrt{w^T w}}$

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

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

Support vector machine (SVM) optimization

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

Find λ that causes highest errorFind \mathbf{w} that causes lowest error given hardest λ

$$\text{Gradient ascent: } \lambda_i \leftarrow \lambda_i + \varepsilon \frac{\partial}{\partial \lambda_i} \mathcal{L}(\mathbf{x}, y; \mathbf{w}, \lambda)$$

$$\text{Gradient descent: } w_j \leftarrow w_j - \varepsilon \frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{x}, y; \mathbf{w}, \lambda)$$

11

Support vector machine (SVM) optimization

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

$$\mathcal{L}(\mathbf{x}, y; \mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{w} + \left(\sum_{i \in +1} \lambda_i (1 - (\mathbf{w}^T \mathbf{x}^i + b)) + \sum_{i \in -1} \lambda_i ((\mathbf{w}^T \mathbf{x}^i + b + 1)) \right)$$

$$\text{Gradient ascent: } \lambda_i \leftarrow \lambda_i + \varepsilon \frac{\partial}{\partial \lambda_i} \mathcal{L}(\mathbf{x}, y; \mathbf{w}, \lambda) \quad \text{Require } \lambda \geq 0$$

$$\text{Gradient descent: } w_j \leftarrow w_j - \varepsilon \frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{x}, y; \mathbf{w}, \lambda) \quad \text{If } \lambda \text{ drops below } 0, \text{ reset to } \lambda = 0$$

13

Support vector machine (SVM) optimization

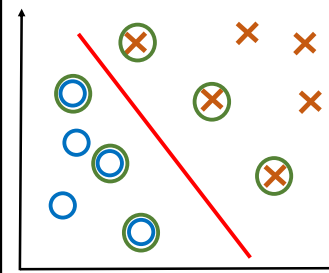
$$\operatorname{argmax}_{\lambda} \operatorname{argmin}_{\mathbf{w}} \mathbf{w}^T \mathbf{w} + \left(\sum_{i \in +1} \lambda_i (1 - (\mathbf{w}^T \mathbf{x}^i + b)) \right) + \sum_{i \in -1} \lambda_i ((\mathbf{w}^T \mathbf{x}^i + b + 1))$$

$$\text{Gradient descent: } w_j \leftarrow w_j - \varepsilon \frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{x}, y; \mathbf{w}, \lambda)$$

$$\frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{x}, y; \mathbf{w}, \lambda): 2w_j + (\sum_{i \in +1} -\lambda_i x_j^i + \sum_{i \in -1} \lambda_i x_j^i)$$

15

Alternate SVM formulation



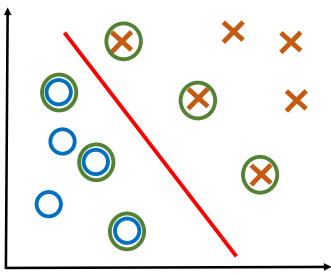
$$\mathbf{w} = \sum_i \alpha^i \mathbf{x}^i y^i$$

Support vectors \mathbf{x}_i have $\alpha_i > 0$ y_i are the data labels +1 or -1

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

17

Example



$$\mathbf{w} = \sum_i \alpha^i \mathbf{x}^i y^i \quad \begin{array}{l} \alpha^i \geq 0 \quad \forall i \\ \sum_i \alpha^i y^i = 0 \end{array}$$

$$\mathbf{x}^1 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, y^1 = +1, \alpha^1 = 0.5$$

$$\mathbf{x}^2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, y^2 = +1, \alpha^2 = 0.7$$

$$\mathbf{x}^3 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, y^3 = -1, \alpha^3 = 1$$

$$\mathbf{x}^4 = \begin{bmatrix} -0.5 \\ -3 \end{bmatrix}, y^4 = -1, \alpha^4 = 0.2$$

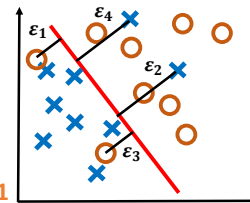
$$\mathbf{w} = 0.5 \times \begin{bmatrix} -1 \\ 1 \end{bmatrix} + 0.7 \times \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 1 \times \begin{bmatrix} -1 \\ -1 \end{bmatrix} - 0.2 \times \begin{bmatrix} -0.5 \\ -3 \end{bmatrix}$$

$$= \begin{bmatrix} -0.5 + 1 + 0.1 \\ 0.5 + 1 + 0.6 \end{bmatrix} = \begin{bmatrix} 0.6 \\ 2.1 \end{bmatrix}$$

19

Support vector machine (SVM) optimization
with slack variables

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



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

subject to

$$\mathbf{w}^T \mathbf{x} + b \geq 1 - \varepsilon^i \quad \text{for } \mathbf{x} \text{ in class 1}$$

$$\mathbf{w}^T \mathbf{x} + b \leq -1 + \varepsilon^i \quad \text{for } \mathbf{x} \text{ in class -1}$$

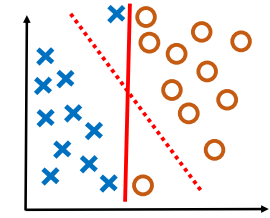
$$\varepsilon^i \geq 0 \quad \forall i$$

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

20

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

$$\mathbf{w}^T \mathbf{x} + b \geq 1 - \varepsilon^i \quad \text{for } \mathbf{x} \text{ in class 1}$$

$$\mathbf{w}^T \mathbf{x} + b \leq -1 + \varepsilon^i \quad \text{for } \mathbf{x} \text{ in class -1}$$

$$\varepsilon^i \geq 0 \quad \forall i$$

21

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

22

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

23

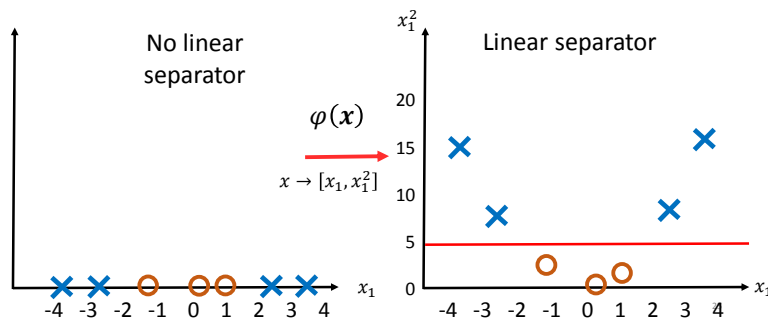
Binary $\rightarrow M$ -class classification

- 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

24

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)

$$\mathbf{w}^T \mathbf{x}^k + b = \sum_i \alpha^i y^i (\mathbf{x}^i)^T \mathbf{x}^k + b$$

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

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

N^2 values to learn for \mathbf{w} in higher-dimensional space

$$\begin{aligned} \text{Or, observe: } (\mathbf{v}^T \mathbf{x} + 1)^2 &= \mathbf{v}_1^2 x_1^2 + \dots + \mathbf{v}_N^2 x_N^2 \\ &+ \mathbf{v}_1 \mathbf{v}_2 x_1 x_2 + \dots + \mathbf{v}_{N-1} \mathbf{v}_N x_{N-1} x_N \\ &+ \mathbf{v}_1 x_1 + \dots + \mathbf{v}_N x_N \end{aligned}$$

\mathbf{v} with N elements
operating in quadratic
space

26

Quadratic mapping function *Simplified*

$$\mathbf{x} = [x_1, x_2] \rightarrow [\sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2, 1]$$

$$\mathbf{x}^i = [5, -2] \rightarrow [10, -4, 25, 4, -20, 1] \quad \mathbf{x}^k = [3, -1] \rightarrow [6, -2, 9, 1, -6, 1]$$

$$\varphi(\mathbf{x}^i)^T \varphi(\mathbf{x}^k) = 30 + 4 + 225 + 4 + 60 + 1 = 324$$

$$\text{Or, observe: } (\mathbf{x}^{iT} \mathbf{x}^k + 1)^2 = ((15 + 2) + 1)^2 = (18)^2 = 324$$

27

Mapping function(s)

• Map from low-dimensional space $\mathbf{x} = (x_1, x_2)$ to higher dimensional space $\varphi(\mathbf{x}) = (\sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2, 1)$

• N data points guaranteed to be separable in space of N-1 dimensions or more

$$\mathbf{w} = \sum_i \alpha_i \varphi(\mathbf{x}^i) y^i$$

Classifying \mathbf{x}^k :

$$\sum_i \alpha_i y^i \varphi(\mathbf{x}^i)^T \varphi(\mathbf{x}^k) + b$$

29

Kernels

Classifying \mathbf{x}^k :

$$\sum_i \alpha_i y^i \varphi(\mathbf{x}^i)^T \varphi(\mathbf{x}^k) + b$$

Kernel trick:

• Estimate high-dimensional dot product with function

$$\bullet K(\mathbf{x}^i, \mathbf{x}^k) = \varphi(\mathbf{x}^i)^T \varphi(\mathbf{x}^k)$$

Now classifying \mathbf{x}^k

$$\sum_i \alpha_i y^i K(\mathbf{x}^i, \mathbf{x}^k) + b$$

30

Radial Basis Kernel

Try projection to infinite dimensions

$$\varphi(\mathbf{x}) = [x_1, \dots, x_n, x_1^2, \dots, x_n^2, \dots, x_1^\infty, \dots, x_n^\infty]$$

$$\text{Taylor expansion: } e^x = \frac{x^0}{0!} + \frac{x^1}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^\infty}{\infty!}$$

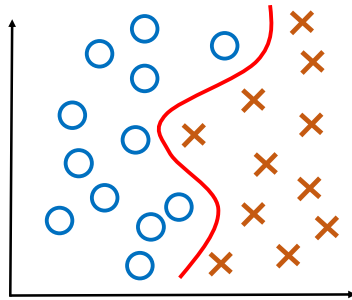
$$K(\mathbf{x}^i, \mathbf{x}^k) = \exp\left(-\frac{(\mathbf{x}^i - \mathbf{x}^k)^2}{2\sigma^2}\right)$$

$$\text{Note: } (\mathbf{x}^i - \mathbf{x}^k)^2 = (\mathbf{x}^i - \mathbf{x}^k)^T (\mathbf{x}^i - \mathbf{x}^k)$$

Draw separating plane to curve around all support vectors

31

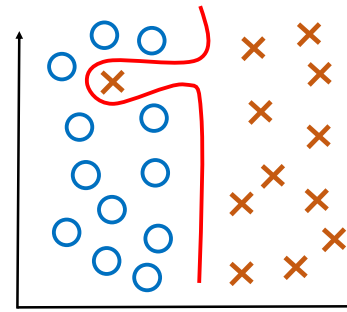
Example RBF-kernel separator



Large margin
Non-linear separation

32

Potential dangers of RBF-kernel separator



Small margin - **overfitting**
Non-linear separation

33

The power of SVM (+kernels)

Boundary defined by a few support vectors

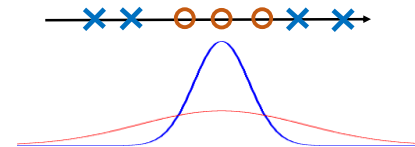
- Caused by: maximizing margin
- Causes: less overfitting
- Similar to: regularization

Kernels keep number of learned parameters in check

34

Benefits of generative methods

- $P(\mathbf{D}|\boldsymbol{\theta})$ and $P(\boldsymbol{\theta}|\mathbf{D})$ can generate non-linear boundary
- E.g.: Gaussians with multiple variances



35