

Lecture 13. Introduction to Support Vector Machines

COMP90051 Statistical Machine Learning

Semester 2, 2015

Lecturer: Andrey Kan

Content is based on slides
provided by Ben Rubinstein



THE UNIVERSITY OF
MELBOURNE

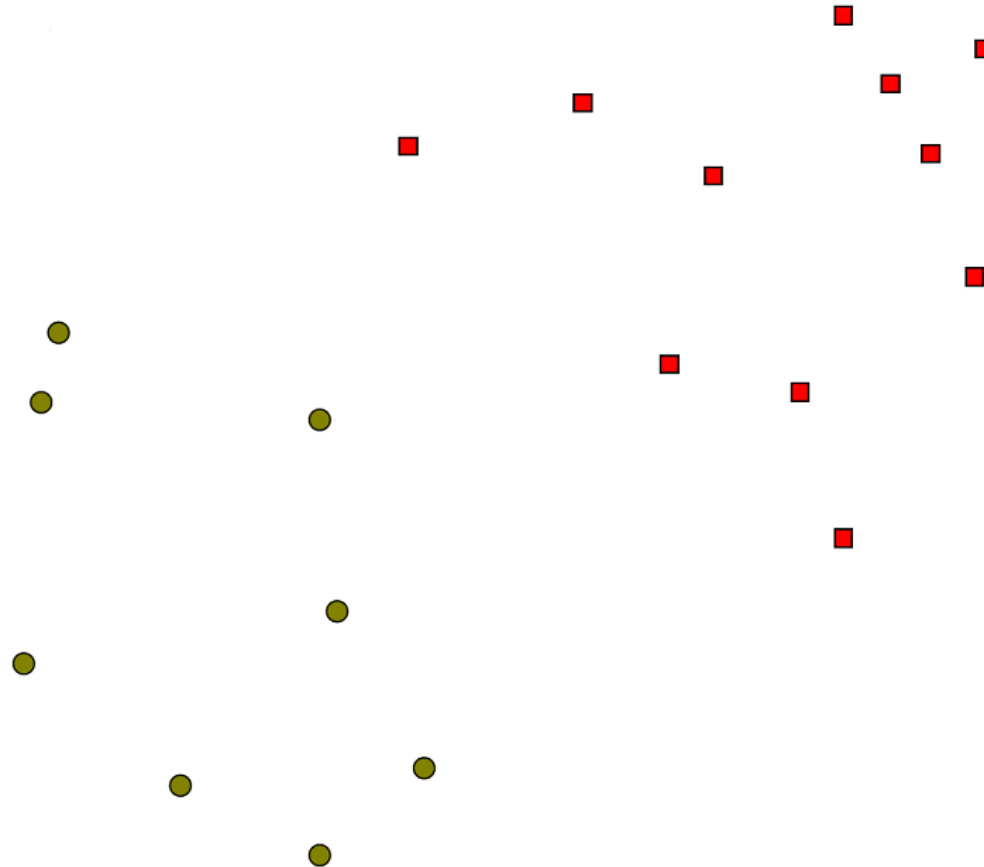
Copyright
University of
Melbourne

Linear Support Vector Machines (SVM)

Yet another linear model that can be regularised

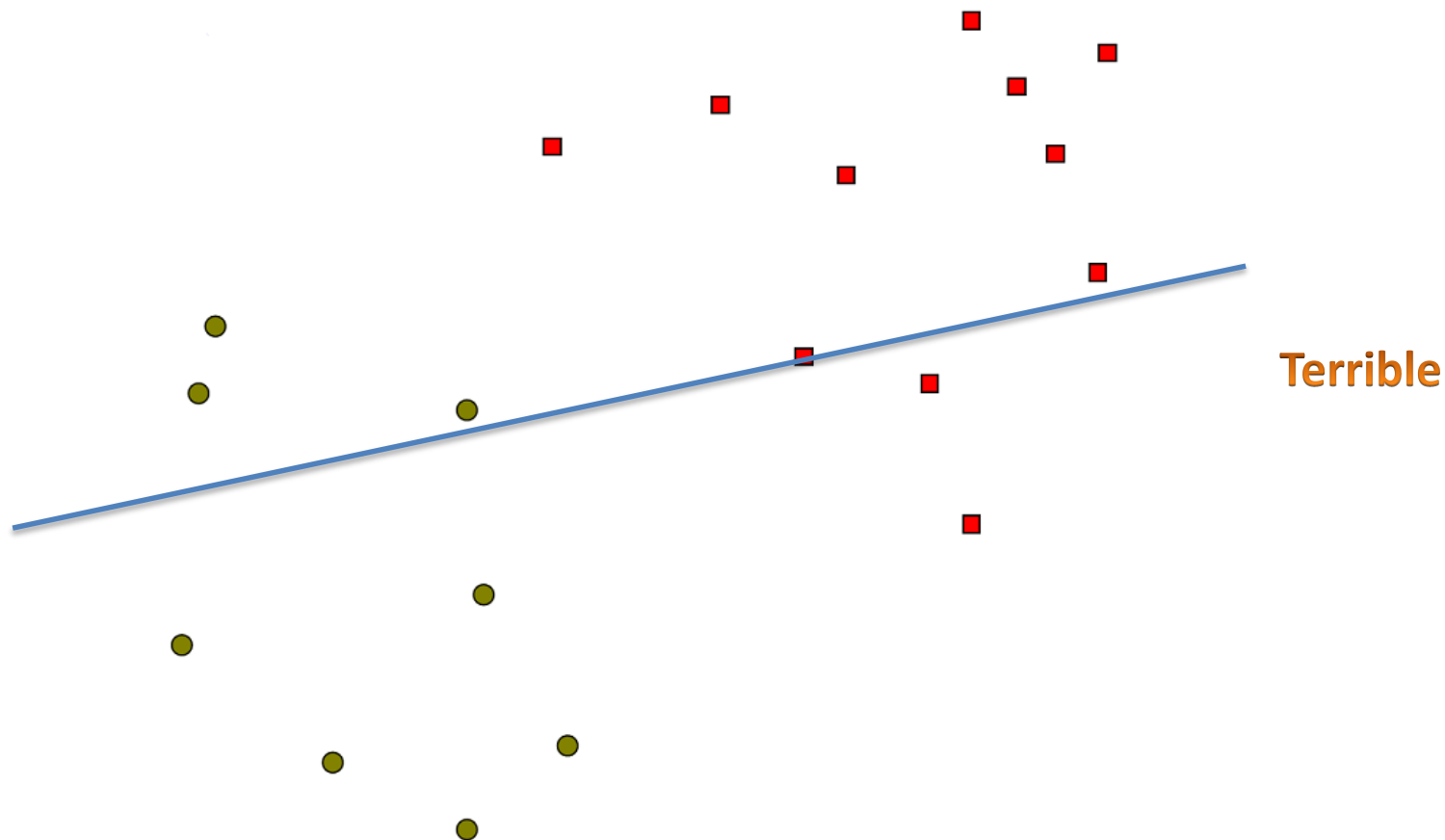
Support Vector Machines: Intuition

- What linear classifier to choose?



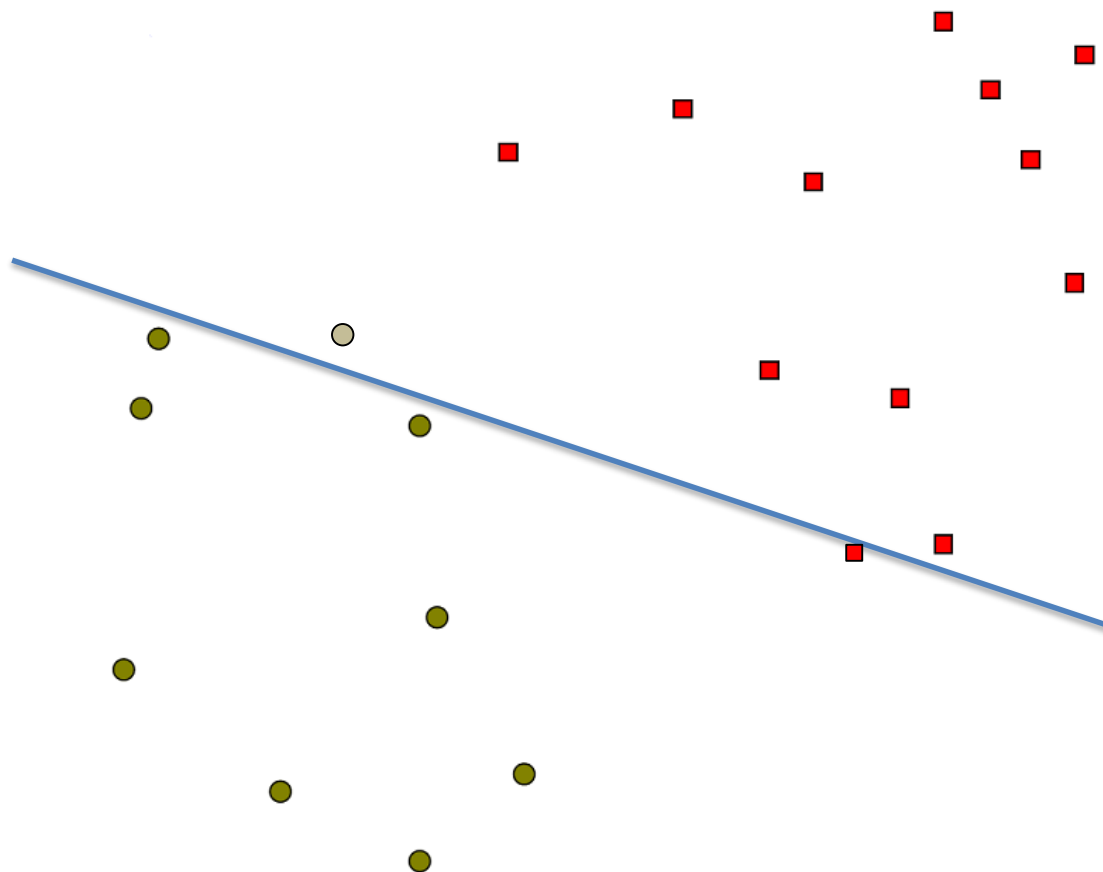
Support Vector Machines: Intuition

- What linear classifier to choose?



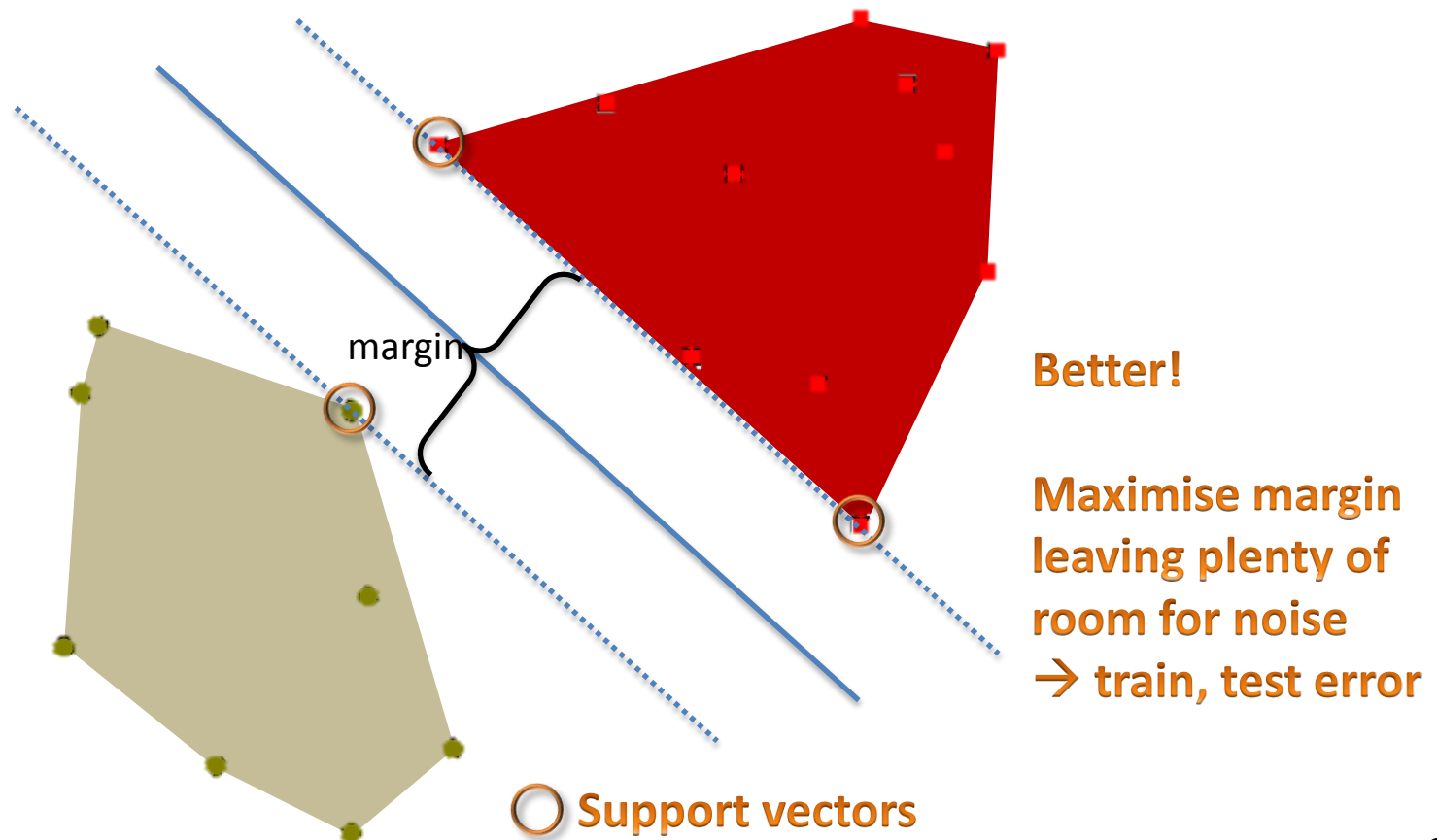
Support Vector Machines: Intuition

- What linear classifier to choose?



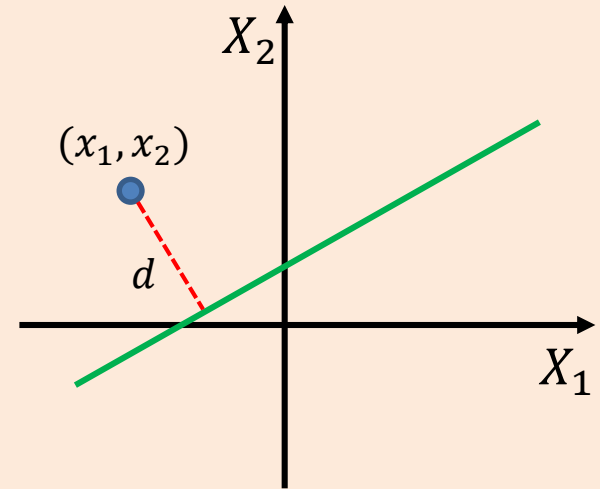
Support Vector Machines: Intuition

- What linear classifier to choose?



(very brief) Revision of geometry

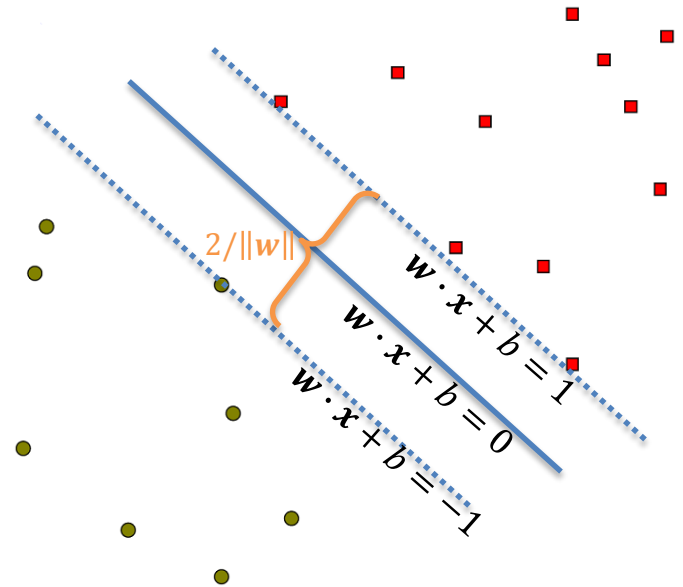
- Line = set of points (x_1, x_2) that satisfy $w_1x_1 + w_2x_2 + b = 0$
- Plane $\rightarrow w_1x_1 + w_2x_2 + w_3x_3 + b = 0$
- Hyperplane $\rightarrow \sum_{i=1}^n w_ix_i + b = 0$
- Note: we can multiply w_i, b by a constant without changing the (hyper-)plane



- Distance from point \mathbf{x} to (hyper-)plane is $d = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|}$
- We can rescale so that $\mathbf{w}^T \mathbf{x} + b = 1$, and get $d = \frac{1}{\|\mathbf{w}\|}$

Linear SVM parameters

- Goal: max margin hyperplane
- Hyperplane parameters
 - * Normal vector \mathbf{w}
 - * Offset b
- Scale parameters so that
 - * **Negative** data: $\mathbf{w} \cdot \mathbf{x} + b \leq -1$
 - * **Positive** data: $\mathbf{w} \cdot \mathbf{x} + b \geq +1$
- Then margin is $2/\|\mathbf{w}\|$



Linear SVM algorithm

- If linearly separable data: the optimal classifier is

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \forall i$

Hard-margin SVM

- In general, relax to allow training errors

- * Intuitively minimise: $\frac{1}{2} \|\mathbf{w}\|^2 + C \cdot \text{average}(\text{loss}(y_i, \hat{y}_i))$
- * Formally, introduce “slacks” $\xi_i \geq 0$'s

$$\min_{\mathbf{w}, b, \xi} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{n} \sum_{i=1}^n l(\xi_i) \right)$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$

Soft-margin SVM

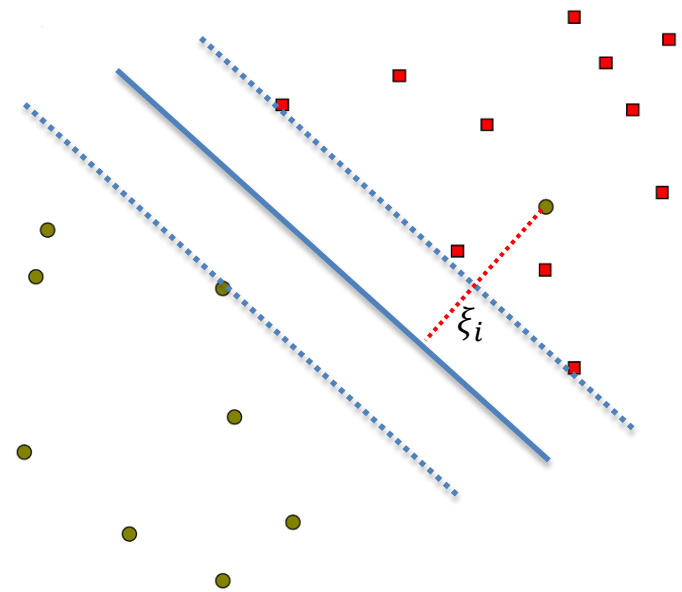
Soft-margin SVM

“Slack” variables: it is OK for some points to be on the wrong side

Still, closer they are to the margin the better

$$\min_{\mathbf{w}, b, \xi} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{n} \sum_{i=1}^n l(\xi_i) \right)$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$



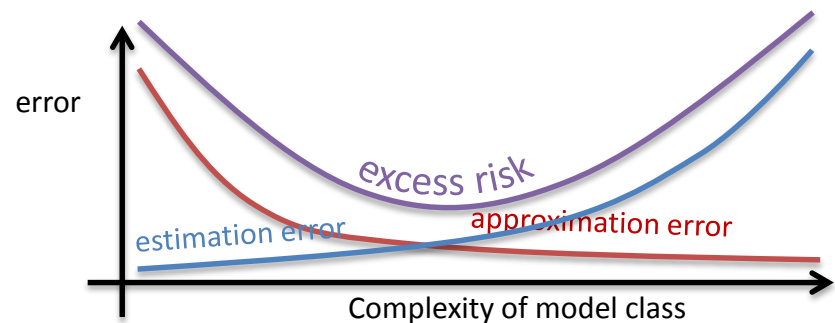
Déjà vu?

Soft-margin SVM:
$$\min_{\mathbf{w}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{model class complexity}} + C \cdot \underbrace{\text{average}(\text{loss}(y_i, \hat{y}_i))}_{\text{training error}} \right)$$

Ridge regression:
$$\min_{\mathbf{w}} \left(\underbrace{\lambda \|\mathbf{w}\|^2}_{\text{model class complexity}} + \underbrace{\sum_i (y_i - \mathbf{X}_i \cdot \mathbf{w})^2}_{\text{training error}} \right)$$

Soft-margin parameter C

- Varies effective complexity
- Theory: should grow like $O(\sqrt{n})$
- Practice: set with cross-validation



Training the SVM

- Loss function
 - * Must be convex to make minimisation tractable
 - * Typically: linear (hinge), squared
- Solving the minimisation
 - * Convex \rightarrow unique global optimum guaranteed
 - * Could use gradient descent, but there are better ways
 - * Constrained Quadratic Program \rightarrow standard
 - * Many tricks to improve speed or memory (chunking, SMO)
- Non-linear SVM coming up (“kernelisation”)

SVM summary (so far)

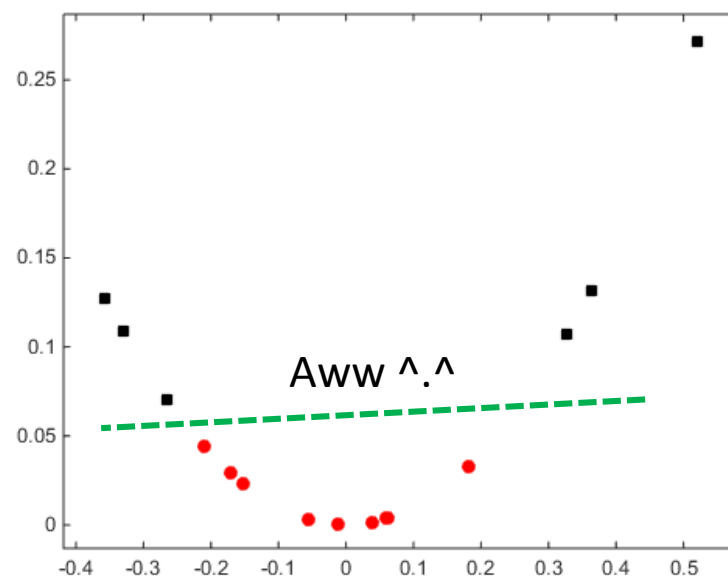
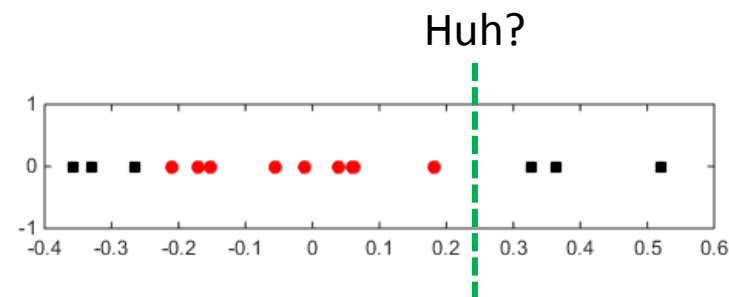
Non-parametric	There are essentially parameters per support vector (next time)
Discriminative	SVM models $P(Y X)$ – the decision boundary – nothing about $P(X)$
Regularised, frequentist	Extent of regularisation controlled by varying the soft-margin parameter; Justification is via theoretical bounds on risk
Computational efficiency	Mediocre cubic polynomial in #data, or #features, for training
Accuracy (aka statistical efficiency)	Typically very good empirically; Exists nice theory on error bounds; Much better if non-linear in some cases.
Packages	http://svmlight.joachims.org http://www.csie.ntu.edu.tw/~cjlin/libsvm Pyhton, R packages, MatLab, Weka (Java), etc.

The Kernel Trick

Very powerful method that enables effective use of linear methods for non-linear data

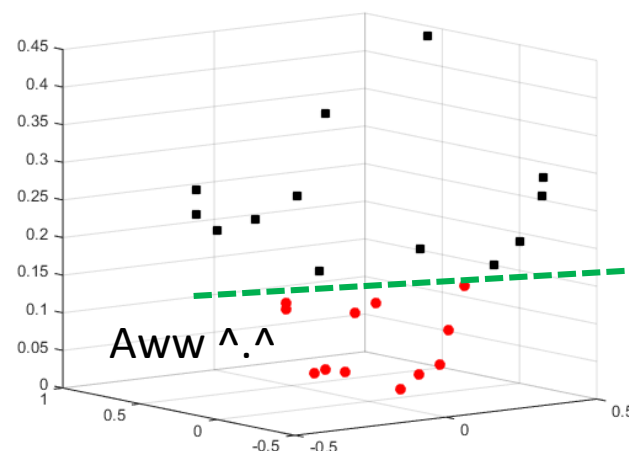
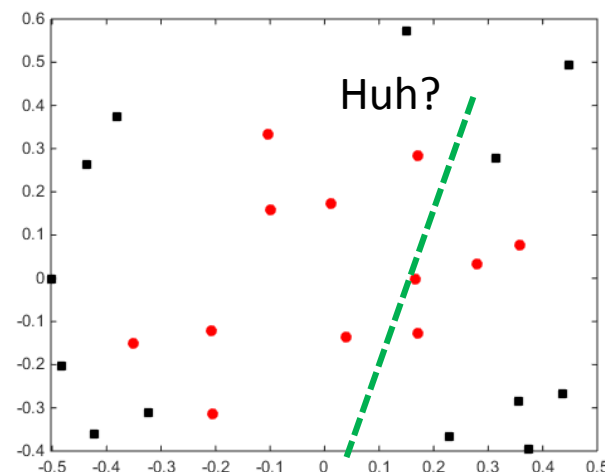
Feature space mapping: example 1

- Consider a binary classification problem
- Each example has one feature $[x_1]$
- Not linearly separable
- Now 'add' a feature $x_2 = x_1^2$
- Each point is now $[x_1, x_1^2]$
- Linearly separable!



Feature space mapping: example 2

- Consider a binary classification problem
- Each example has features $[x_1, x_2]$
- Not linearly separable
- Now 'add' a feature $x_3 = x_1^2 + x_2^2$
- Each point is now $[x_1, x_2, x_1^2 + x_2^2]$
- Linearly separable!



Observation 1

- A non-linear boundary in original feature space can become a linear boundary in a transformed space
- Denote this transformation as $\varphi(\mathbf{x})$
- For example, $\varphi([x_1, x_2]) = [x_1, x_2, x_1^2 + x_2^2]$
- Note $\varphi(\mathbf{x})$ has more dimensions than \mathbf{x}
 - * Commonly used $\varphi(\mathbf{x})$ often have very large dimensionality
 - * Sometimes $\varphi(\mathbf{x})$ have infinite dimensionality!

Naïve workflow

- Choose/design a linear model
- Choose/design a transformation $\varphi(\mathbf{x})$
- For each training example, and for each new instance compute $\varphi(\mathbf{x})$
- Train classifier/Do predictions
- Problem: impractical/impossible to compute $\varphi(\mathbf{x})$ for high/infinite-dimensional $\varphi(\mathbf{x})$

Observation 2

- For a large class of linear models, both parameter estimation and computing predictions depend on data only in a form of dot product
- Dot product $\mathbf{x}\mathbf{y} = \sum_{i=1}^d x_i y_i$
- In transformed space $\varphi(\mathbf{x})\varphi(\mathbf{y}) = \sum_{i=1}^{d'} \varphi(\mathbf{x})_i \varphi(\mathbf{y})_i$
- Kernel is a function that can be expressed as a dot product in some feature space $K(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})\varphi(\mathbf{y})$

Kernel trick: example

- For many transformations there is a shortcut computation of the dot product!
- For example, consider two points in original space $\mathbf{x} = [x_1]$ and $\mathbf{y} = [y_1]$
- Consider transformation $\varphi(\mathbf{x}) = [x_1^2, \sqrt{2c}x_1, c]$
- Then $\varphi(\mathbf{x})\varphi(\mathbf{y}) = (x_1^2y_1^2 + 2cx_1y_1 + c^2)$
- This simplifies to $\varphi(\mathbf{x})\varphi(\mathbf{y}) = (x_1y_1 + c)^2$
- Here $K(\mathbf{x}, \mathbf{y}) = (x_1y_1 + c)^2$ is called kernel

Kernel trick

- For a pair of data points in original space kernel $K(\mathbf{x}, \mathbf{y})$ can be used to compute dot product in the transformed space
- Computing the transformation $\varphi(\mathbf{x})$ as such is not necessary!
- Kernel can be computed in $O(d)$, whereas computing $\varphi(\mathbf{x})$ requires $O(d')$, where $d' \gg d$ or $d' = \infty$

Summary

- Linear SVM aims to maximise a separating margin
- Soft-margin SVM allows some instances to be misclassified
- Soft-margin SVM includes regularisation parameter
- When data is not linearly separable, transforming the feature space can make it linearly separable
- Many linear methods only require data in form of dot products
- Kernel trick – fast computation of dot products in very high dimensional feature spaces