
Support Vector Machines

Mehul Motani

Electrical & Computer Engineering

National University of Singapore

Email: motani@nus.edu.sg

AI and Machine Learning

- “AI is the new electricity. Just as 100 years ago electricity transformed industry after industry, AI will now do the same.” – Andrew Ng
- Machine learning is about learning from data.
- Supervised learning – Learning from data with labels which serve a supervisory purpose
- Unsupervised learning – Learning from data without labels allows tasks such as clustering.
- Reinforcement learning – Learning from data without labels but there is feedback from the environment.

Support Vector Machines (SVM)

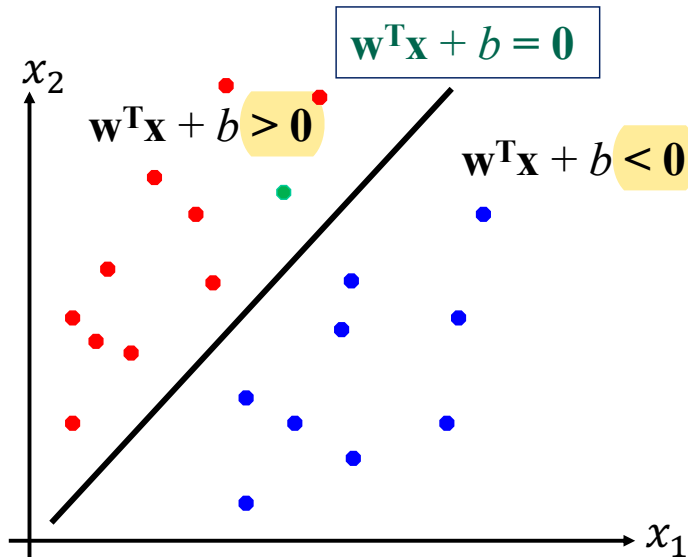
- SVM is a supervised learning algorithm
 - Useful for both classification and regression problems
- Linear SVM – Maximum-Margin Classifier
 - Formalize notion of the best linear separator
- Optimization Problem with Lagrangian Multipliers
 - Technique to solve a constrained optimization problem
- Nonlinear SVM – Extending Linear SVM with Kernels
 - Project data into higher-dimensional space to make it linearly separable.
 - create nonlinear classifiers by applying the kernel trick to maximum-margin hyperplanes.
 - Complexity: Depends only on the number of training examples, not on dimensionality of the kernel space!

SVM – A Brief History

- Pre-1980: Almost all learning methods learned linear decision surfaces.
 - Linear learning methods have nice theoretical properties
- 1980's: Decision trees and Neural Nets allowed efficient learning of non-linear decision surfaces
 - Little theoretical basis and all suffer from local minima
- 1990's: Efficient learning algorithms for non-linear functions based on computational learning theory developed
- Support Vector Machines
 - The original SVM algorithm was invented by Vapnik and Chervonenkis in 1963.
 - Nonlinear SVMs using the kernel trick were first introduced in a conference paper by Boser, Guyon and Vapnik in 1992.
 - The SVM with soft margin was proposed by Cortes and Vapnik in 1993 and published in 1995.

Supervised Learning: Linear Separators

- Binary classification can be viewed as the task of separating classes in feature space.



- Two features: x_1 and x_2
- Two classes: **red** and **blue**
- Linear separator given by line:

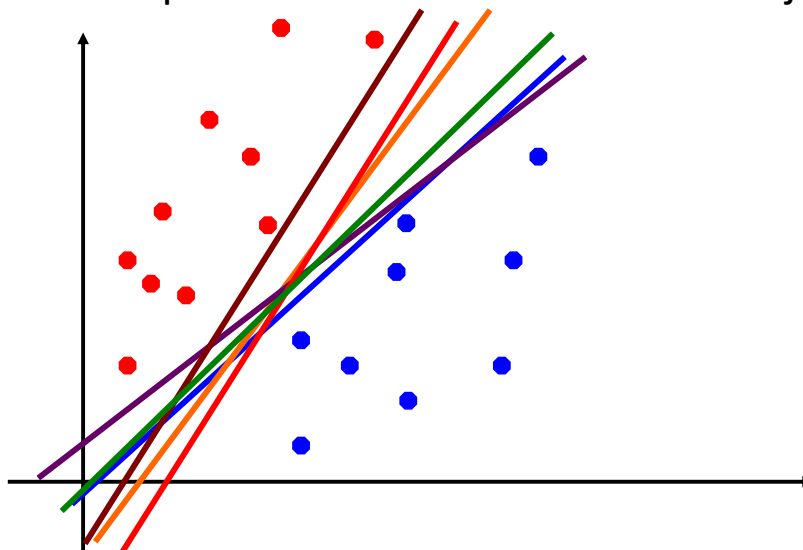
$$\mathbf{w}^T \mathbf{x} + b = 0 \quad (1)$$

$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

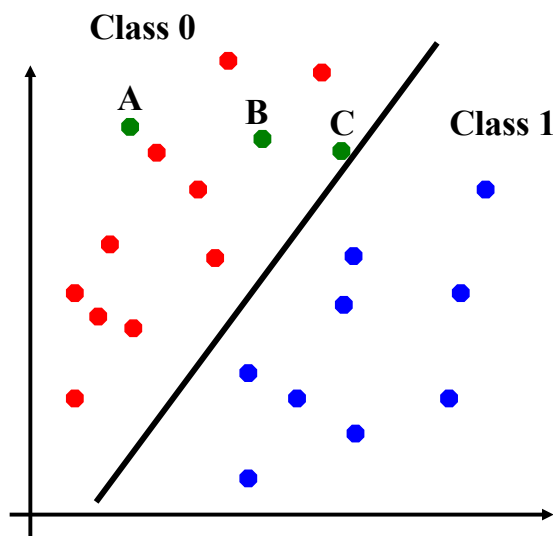
- Classification
 - $\mathbf{w}^T \mathbf{x} + b < 0 \rightarrow$ **blue (-1)**
 - $\mathbf{w}^T \mathbf{x} + b > 0 \rightarrow$ **red (+1)**
- Classifier function
 - $f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$
- New data: **Green dot** will be classified as **red class (+1)**

Linear Separators

- There are many possible linear separators!
- Which of the linear separators is optimal?
- The linear SVM solution defines an objective and finds the linear separator which maximizes that objective.



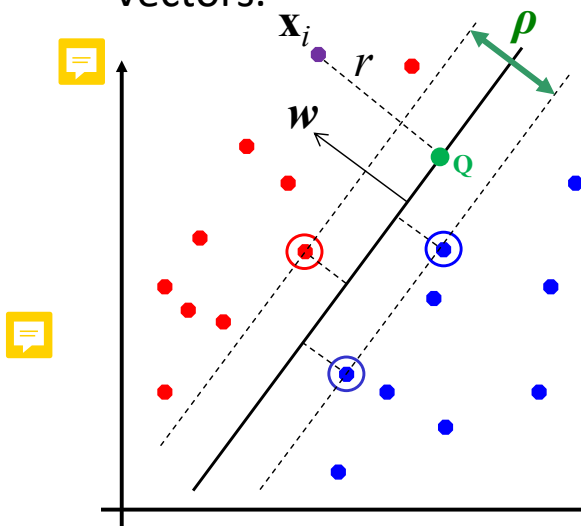
Linear Separators and Margin



- Consider three new data points: A, B, C (**green dots**), all of which are classified as class 0.
- How confident are you that point A is class 0?
- What about point C?
- What about point B?
- Intuitively, we are more confident about point A than point C.
- Intuition: if a point is far from the separating hyperplane (i.e., large margin), then we may be more confident in our prediction.

Classification Margin

- Data points closest to the hyperplane are called the **support vectors** (circled data points)
- Margin** ρ of separator is the distance between support vectors
- Note that the separator is completely defined by its support vectors.



What is the distance, r , from data point \mathbf{x}_i to the separator?

For \mathbf{x}_1 and \mathbf{x}_2 on the separating hyperplane:

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0 \Rightarrow \mathbf{w} \perp \text{hyperplane} \quad (1)$$

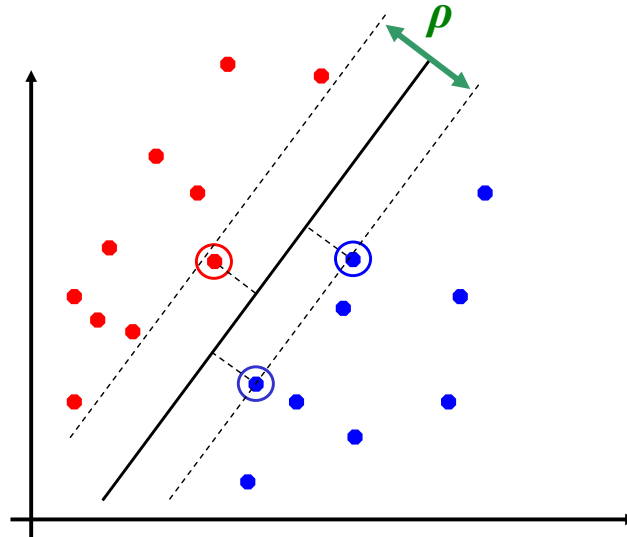
$$\mathbf{w}^T \left(\mathbf{x}_i - r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + b = 0 \Rightarrow r = \frac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|} \quad (2)$$

Point Q

Norm of \mathbf{w}

Maximum Margin Classification

- Maximizing the margin is provably good and intuitive
 - Larger margin leads to lower generalization error (Vapnik).
- Implies that only support vectors matter; other training examples can be ignored → SVM is stable and robust to outliers



Linear SVM Mathematically

- Let training set be $S = \{(\mathbf{x}_i, y_i)\}_{i=1,2,\dots,n}$ with $\mathbf{x}_i \in \mathbf{R}^d$ and $y_i \in \{-1, 1\}$.
- Suppose we have a separating hyperplane with margin ρ , weight vector \mathbf{w} and scalar b .
- Then for each training example (\mathbf{x}_i, y_i) :

$$\begin{aligned} \mathbf{w}^T \mathbf{x}_i + b &\leq -\rho/2 & \text{if } y_i = -1 \\ \mathbf{w}^T \mathbf{x}_i + b &\geq \rho/2 & \text{if } y_i = 1 \end{aligned} \Leftrightarrow y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq \rho/2 \quad (1)$$

- For every support vector \mathbf{x}_s the above inequality is an equality. After rescaling \mathbf{w} and b by $\rho/2$ in the equality, we obtain that distance between each \mathbf{x}_s and the hyperplane is

$$r = \frac{y_s(\mathbf{w}^T \mathbf{x}_s + b)}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|} \quad (2)$$

- Then the margin can be expressed through (rescaled) \mathbf{w} and b as:

$$\rho = 2r = \frac{2}{\|\mathbf{w}\|} \quad (3)$$

Linear SVMs Mathematically (cont.)

- Then we can formulate the quadratic optimization problem:

Find \mathbf{w} and b such that $\rho = \frac{2}{\|\mathbf{w}\|}$ is maximized (a)

and for all $(\mathbf{x}_i, y_i) \in S$: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$ (b)

(1)

- We can reformulate the problem in (1) as follows:

Find \mathbf{w} and b such that

$\Phi(\mathbf{w}) = \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w}$ is minimized

and for all $(\mathbf{x}_i, y_i) \in S$: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

(2)

Solving the Optimization Problem

Primal: Find \mathbf{w} and b such that

$\Phi(\mathbf{w}) = \mathbf{w}^T \mathbf{w}$ is minimized

and for all $(\mathbf{x}_i, y_i) \in S$: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

(1)

- Need to optimize a *quadratic* function subject to *linear* constraints.
- Quadratic optimization problems are a well-known class of mathematical programming problems for which several (non-trivial) algorithms exist.
- Solution involves constructing a *dual problem* where a *Lagrange multiplier* α_i is associated with every inequality constraint in the primal problem:

Dual: Find $\alpha_1 \dots \alpha_n$ such that

$Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ is maximized and

(1) $\sum \alpha_i y_i = 0$

(2) $\alpha_i \geq 0$ for all α_i

(2)

See https://en.wikipedia.org/wiki/Quadratic_programming

The Optimization Problem Solution

- Given a solution $\alpha_1 \dots \alpha_n$ to the dual problem, solution to the primal is:



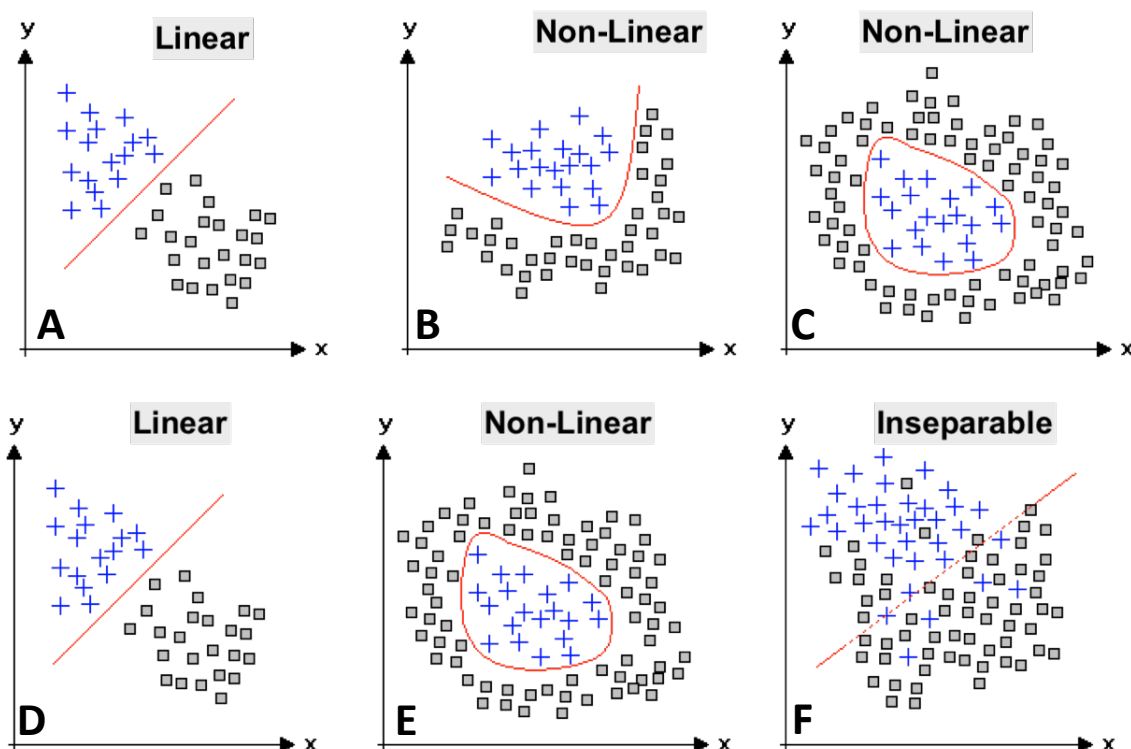
$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \quad b = y_k - \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_k \quad \text{for any } \alpha_k > 0 \quad (1)$$

- Each non-zero α_i indicates that corresponding \mathbf{x}_i is a support vector.
- Then the classifying function is (note that we don't need \mathbf{w} explicitly):

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \quad (2)$$

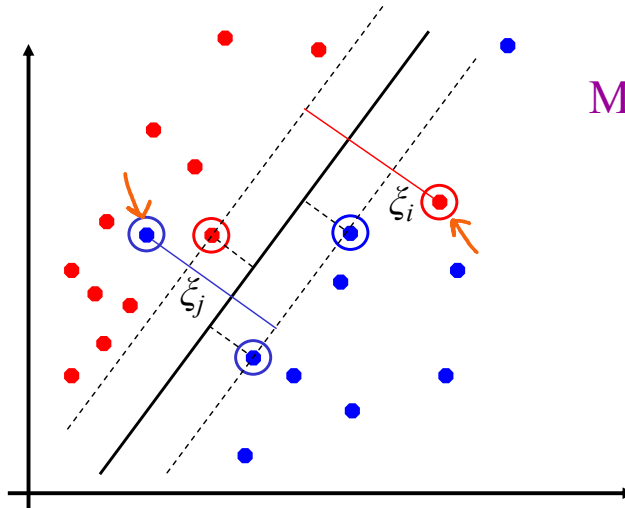
- The quantity $\mathbf{x}^T \mathbf{y}$ is called the inner product or dot product between the vector \mathbf{x} and the vector \mathbf{y} .
- Notice that the solution relies on the inner product between the test point \mathbf{x} and the support vectors \mathbf{x}_i – we will return to this later.
- Also keep in mind that solving the optimization problem involved computing the inner products $\mathbf{x}_i^T \mathbf{x}_j$ between all training points.

Linear and nonlinear data models



Soft Margin Classification

- What if the training set is not linearly separable?
- *Slack variables* ξ_i can be added to allow misclassification of difficult or noisy examples, resulting margin called *soft*.



What should our quadratic optimization criterion be?

Minimize:

$$\underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w}}_{\text{Maximize Margin}} + C \underbrace{\sum_{k=1}^R \xi_k}_{\text{Misclassification Penalty}}$$

Note: that ξ is the Greek letter Xi and is pronounced as 'zai' or 'ksi'.

Hard margin vs Soft margin

- The **hard-margin** SVM formulation:

Find \mathbf{w} and b such that
 $\Phi(\mathbf{w}) = \mathbf{w}^T \mathbf{w}$ is minimized
 and for all $(\mathbf{x}_i, y_i) \in S$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$



(1)

- Modified **soft-margin** SVM formulation with slack variables:

Find \mathbf{w} and b such that
 $\Phi(\mathbf{w}) = \mathbf{w}^T \mathbf{w} + C \sum \xi_i$ is minimized
 and for all $(\mathbf{x}_i, y_i) \in S$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$, $\xi_i \geq 0$



(2)

- Parameter C can be viewed as a way to control overfitting
 - It trades off the relative importance of maximizing the margin and fitting the training data.
 - Larger $C \rightarrow$ the more the penalty for misclassifications. This leads to smaller and smaller margins but less misclassifications. This is essentially overfitting.
 - Small $C \rightarrow$ the lower the penalty for misclassifications. This leads to larger margins but more misclassifications.

Soft Margin Classification – Solution

- Dual problem is identical to separable case:

Find $\alpha_1 \dots \alpha_N$ such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ is maximized and

(1)

(1) $\sum \alpha_i y_i = 0$

(2) $0 \leq \alpha_i \leq C$ for all α_i

- Again, \mathbf{x}_i with non-zero α_i will be support vectors.
- Solution to the soft margin SVM is:

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i$$

$$b = y_k(1 - \xi_k) - \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_k \quad \text{for any } k \text{ s.t. } \alpha_k > 0$$

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

(2)

Note: We don't need to compute \mathbf{w} explicitly for classification:

Note: If the 2-norm penalty for slack variables $C \sum \xi_i^2$ was used in primal objective, we would need additional Lagrange multipliers for slack variables...

Theoretical Justification for Maximum Margins

- VC dimension is a measure of the complexity of a classifier. The more complex the classifier, the more prone it is to overfitting.
- Vapnik proved the following:

The class of optimal linear separators has VC dimension h bounded from above as

$$h \leq \min \left\{ \left\lceil \frac{D^2}{\rho^2} \right\rceil, m_0 \right\} + 1 \quad (1)$$

where ρ is the margin, D is the diameter of the smallest sphere that can enclose all of the training examples, and m_0 is the dimensionality.

- Intuitively, this implies that regardless of dimensionality m_0 we can minimize the VC dimension by maximizing the margin ρ .
- Thus, complexity of the classifier is kept small regardless of dimensionality.

Summary of Linear SVMs

- The classifier is a *separating hyperplane*.
- Most “important” training points are support vectors; they define the hyperplane.
- Quadratic optimization algorithms can identify which training points \mathbf{x}_i are support vectors with non-zero Lagrangian multipliers α_i .
- Both in the dual formulation of the problem and in the solution, the training points appear only inside inner products:

Find $\alpha_1 \dots \alpha_N$ such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ is maximized and

(1) $\sum \alpha_i y_i = 0$

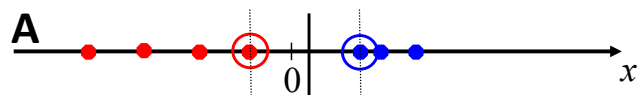
(2) $0 \leq \alpha_i \leq C$ for all α_i (1)

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \quad (2)$$

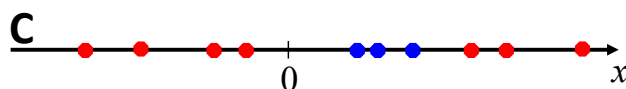
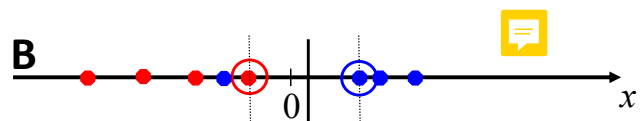
Non-linear SVMs



Consider this noisy dataset:

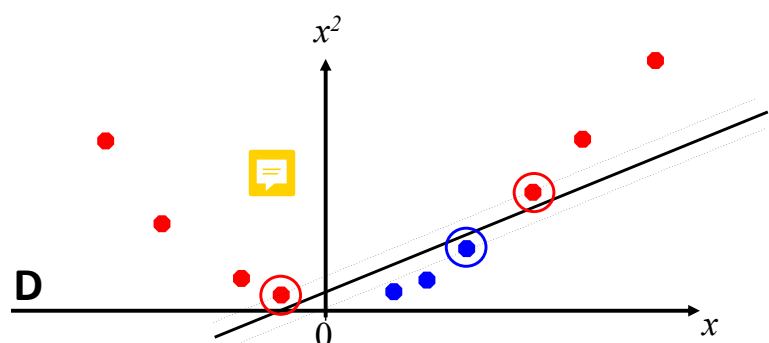


How about this dataset?



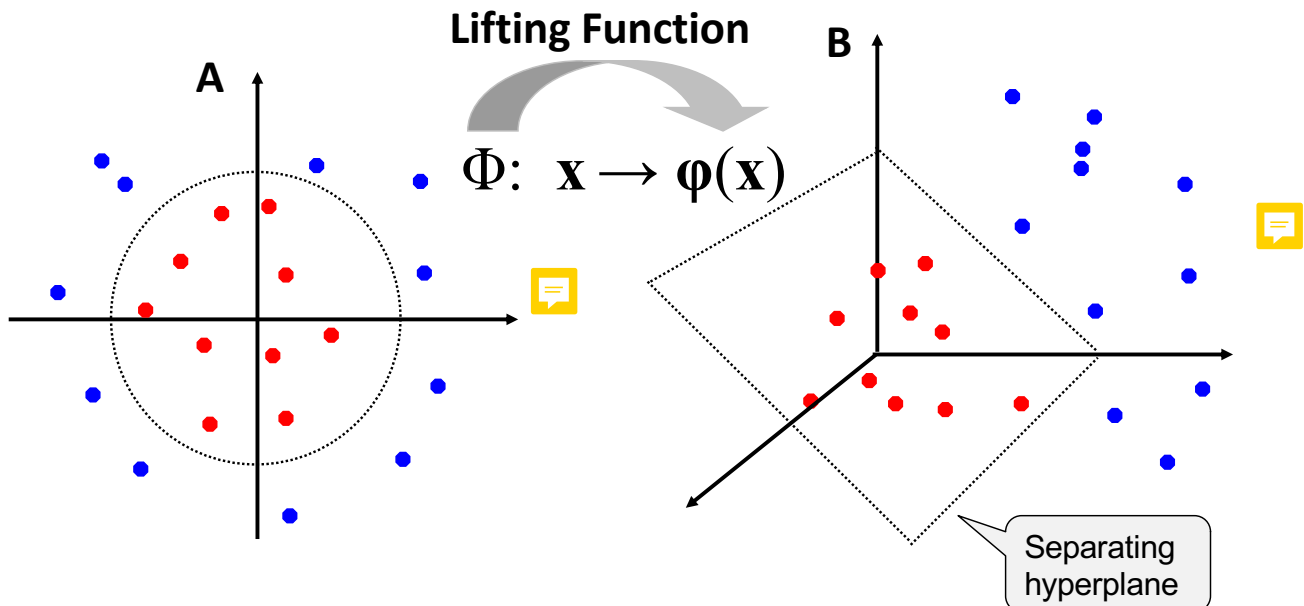
But what are we going to do if the dataset is just too hard?

How about... mapping data to a higher-dimensional space:



Non-linear SVMs: Feature spaces

- General idea: the original feature space is mapped to some higher-dimensional feature space where the training set is separable:



The “Kernel Trick”

- The linear SVM classifier relies on the inner product between vectors, for example: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
 - If every datapoint is mapped into high-dimensional space via some transformation $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$, the inner product becomes: $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$
 - A *kernel function* is a function that is equivalent to an inner product in some higher dimensional feature space.
 - Example: 2-dimensional vectors $\mathbf{x} = [x_1 \ x_2]^T$
 - Let $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$ (1)
 - Need to show that $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ for some $\phi(\mathbf{x})$
- $$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2 = 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} \quad (2)$$
- $$= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}] [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}]^T \quad (2)$$
- $$\rightarrow K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j), \text{ where } \phi(\mathbf{x}) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2]^T \quad (3)$$
- Thus, a kernel function *implicitly* maps data to a high-dimensional space (without the need to compute each $\phi(\mathbf{x})$ explicitly).

What Functions are Kernels?



- For some functions $K(\mathbf{x}_i, \mathbf{x}_j)$ checking that $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ can be cumbersome.
- Mercer's theorem: **Every semi-positive definite symmetric function is a valid kernel**
- Semi-positive definite symmetric functions correspond to a semi-positive definite symmetric Gram matrix:

$$\mathbf{K} = \begin{array}{ccccc} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & K(\mathbf{x}_1, \mathbf{x}_3) & \dots & K(\mathbf{x}_1, \mathbf{x}_n) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & K(\mathbf{x}_2, \mathbf{x}_3) & & K(\mathbf{x}_2, \mathbf{x}_n) \\ & & & & \\ \dots & \dots & \dots & \dots & \dots \\ K(\mathbf{x}_n, \mathbf{x}_1) & K(\mathbf{x}_n, \mathbf{x}_2) & K(\mathbf{x}_n, \mathbf{x}_3) & \dots & K(\mathbf{x}_n, \mathbf{x}_n) \end{array} \quad (1)$$

Check out the discussion at: <https://www.quora.com/What-is-the-kernel-trick>

Examples of Kernel Functions

- Linear: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
 - Mapping $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$, where $\phi(\mathbf{x})$ is \mathbf{x} itself
- Polynomial of power p : $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$
 - Mapping $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$, where $\phi(\mathbf{x})$ has $\binom{d+p}{p}$ dimensions, where d is the original feature space dimension.
- Gaussian (radial-basis function): $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}}$
 - Mapping $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$, where $\phi(\mathbf{x})$ is *infinite-dimensional*: every point is mapped to a *function* (a Gaussian); combination of functions for support vectors is the separator.
- Higher-dimensional space still has *intrinsic* dimensionality d (the mapping is not *onto*), but linear separators in it correspond to *non-linear* separators in original space.

Non-linear SVMs Mathematically

- Dual problem formulation:

Find $\alpha_1 \dots \alpha_n$ such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$ is maximized and

(1) $\sum \alpha_i y_i = 0$

(2) $\alpha_i \geq 0$ for all α_i

(1)

- The solution is:

$$f(\mathbf{x}) = \sum \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_j) + b$$

(2)

- Optimization techniques for finding α_i 's remain the same!

Nonlinear SVM - Summary

- In summary, linear SVM locates a separating hyperplane in the feature space and classifies points in that space
- Nonlinear SVM lifts the problem to a higher dimensional space and performs linear SVM in the higher dimensional space.
- This corresponds to a nonlinear separator in the original feature space.
- The algorithm does not need to represent the space explicitly, it does this by simply defining a kernel function, which plays the role of the inner product in the high dimensional feature space.



Properties of SVM

- Sparseness of solution when dealing with large data sets as only support vectors are used to specify the separating hyperplane
- Ability to handle large feature spaces as the complexity does not depend on the dimensionality of the feature space
- Overfitting can be controlled by soft margin approach
- Mathematically nice – a simple convex optimization problem which is guaranteed to converge to a single global solution
- Supported by theory and intuition
- SVM empirically works very well
 - Text (and hypertext) categorization, image classification,
 - Protein classification, Disease classification
 - Hand-written character recognition

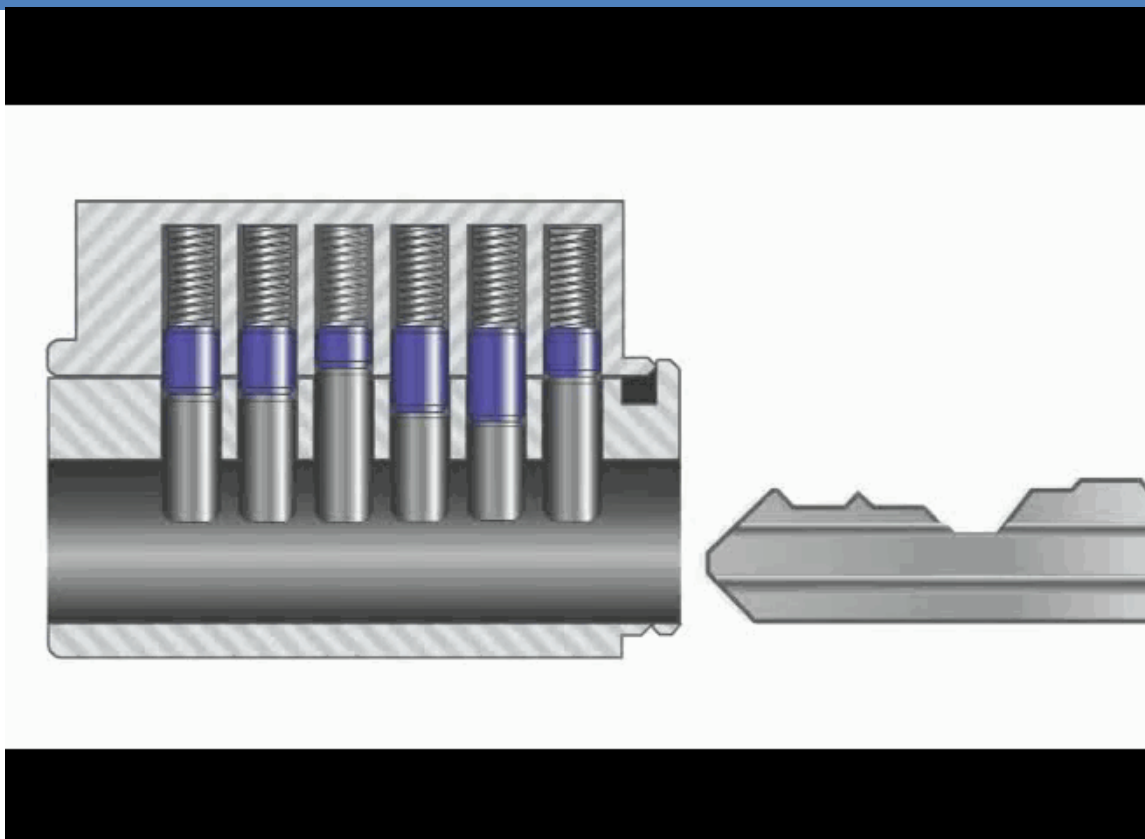
Weakness of SVM

- SVM is sensitive to noise
 - A relatively small number of mislabeled examples can dramatically decrease the performance
- Standard SVM only considers two classes
- Question: How to do multi-class classification with SVM?
- Answer: Build multiple SVMs
 1. With m classes, learn m SVM's
 - SVM 1 learns “Output = 1” vs “Output \neq 1”
 - SVM 2 learns “Output = 2” vs “Output \neq 2”
 - :
 - SVM m learns “Output = m ” vs “Output $\neq m$ ”
 2. To predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the positive region.

SVM Summary

- SVMs were originally proposed by Boser, Guyon and Vapnik in 1992 and gained increasing popularity in late 1990s.
- SVMs are currently among the best performers for a number of classification tasks ranging from text to genomic data.
- SVMs can be applied to complex data types beyond feature vectors (e.g. graphs, sequences, relational data) by designing kernel functions for such data.
- Tuning SVMs remains a black art: selecting a specific kernel and parameters is usually done in a try-and-see manner.
- Some references on VC-dimension and Support Vector Machines:
 - C.J.C. Burges. A tutorial on support vector machines for pattern recognition. Data Mining and Knowledge Discovery, 2(2):955-974, 1998.
 - The VC/SRM/SVM Bible: Statistical Learning Theory by Vladimir Vapnik, Wiley-Interscience, 1998
 - https://en.wikipedia.org/wiki/Support_vector_machine

What do data engineers and thieves have in common?



Thank you!

- Please send me your feedback and any questions you may have.
- The best way to contact me is via email:
mehul.motani@gmail.com
- Thanks for listening!