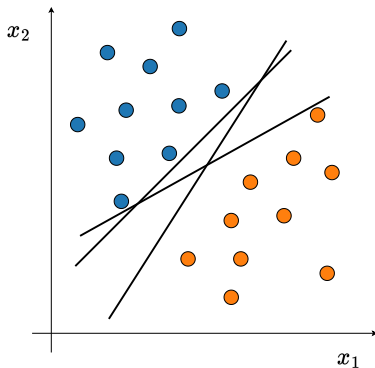


Support vector machines

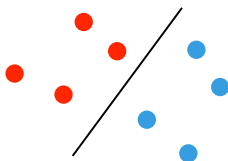
Separating hyperplanes

- **Goal:** Build hyperplanes that separate points in two classes,
- **Question:** Which is the best separating line?



Non-linear separation

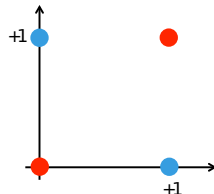
- **Question:** How to handle non-linearly separable data points?



Linearly separable



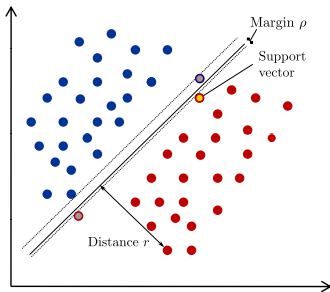
Nonlinearly separable



The xor function

Classification margin

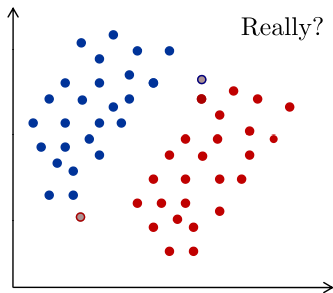
- Signed distance from an example to the separator: $r = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$
- Examples closest to the hyperplane are **support vectors**.
- **Margin** ρ is the distance between the separator and the support vector(s).



- Is this one good?
- Consider two testing samples, where are they assigned?
- This separator may have large generalization error.

Classification margin

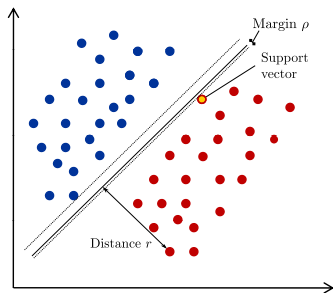
- Signed distance from an example to the separator: $r = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$
- Examples closest to the hyperplane are **support vectors**.
- **Margin** ρ is the distance between the separator and the support vector(s).



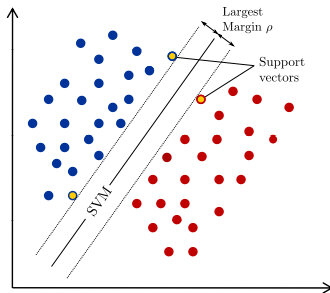
- Is this one good?
- Consider two testing samples, where are they assigned?
- This separator may have large generalization error.

Largest margin

- Maximizing the margin reduces generalization error
- Then, there are necessary support vectors on both sides of the hyperplane
- Only support vectors are important \rightarrow other examples can be discarded



Poor separation



Optimal separation

Margin formulation

In SVM, the target vectors are $y_i = +1$ or $y_i = -1$.

We are looking for parameters (\mathbf{w}, b) such that:

$$\mathbf{w}^\top \mathbf{x}_i + b \geq +1 \quad \text{if} \quad y_i = +1$$

$$\mathbf{w}^\top \mathbf{x}_i + b \leq -1 \quad \text{if} \quad y_i = -1$$

or in short: $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq +1$ and $|\mathbf{w}^\top \mathbf{x}_i + b| \geq 1$

For support vectors \mathbf{x}_s , this inequality can be forced to be an equality (parameters \mathbf{w} and b can always be renormalized to satisfy this):

$$|\mathbf{w}^\top \mathbf{x}_s + b| = 1$$

Then the margin is given by:

$$\rho = |r| = \frac{|\mathbf{w}^\top \mathbf{x}_s + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

Margin maximization problem

The margin maximization problem can be formulated as a minimization:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad \forall i = 1, \dots, N$$

This is a **quadratic** and **convex** optimization problem, subject to linear constraints \rightarrow There exists numerical solvers for this type of problem.

We can reformulate the problem using Lagrange multipliers $\alpha_i \geq 0$:

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i (y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

After cancelling derivatives of the Lagrangian $L(\mathbf{w}, b, \boldsymbol{\alpha})$ with respect to \mathbf{w} and b , we find:

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_{i=1}^N \alpha_i y_i = 0$$

Dual representation

By plugging these into the Lagrangian, we obtain the following maximization problem (called the *dual problem*):

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^{\top} \mathbf{x}_j \\ \text{subject to} \quad & \alpha_i \geq 0 \quad \text{and} \quad \sum_i \alpha_i y_i = 0, \quad \forall i = 1, \dots, N \end{aligned}$$

which can be solved with quadratic programming also.

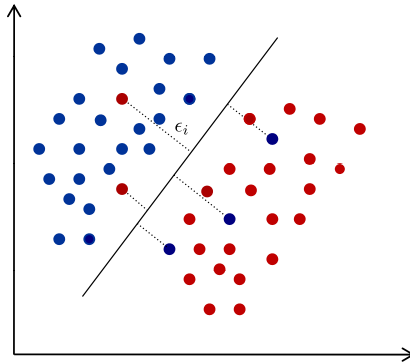
It can be shown that **non-zero α_i correspond to support vectors \mathbf{x}_i** (at test time, only support vectors are considered).

- In practice, the dual problem is more efficient to solve, and allows identifying the support vectors
- It only depends on the dot products $\mathbf{x}_i^{\top} \mathbf{x}_j$ between training points

SVM with non linearly separable data

Non linearly separable data

What if the training set is not linearly separable?



The optimization problem does not admit solutions.

Solution: allow the model to make errors ϵ_i for some x_i .

Soft margin SVM

- Relax the constraints by permitting errors to some extent:

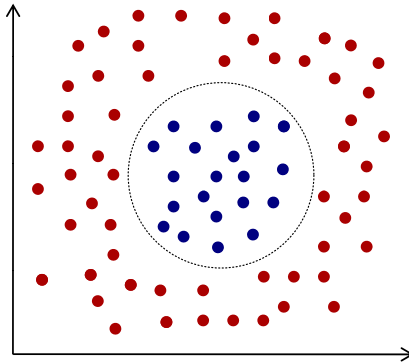
$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$

- The new term is called the **hinge loss**
- The hyperparameter $C > 0$ controls overfitting
- The dual problem becomes

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C \quad \text{and} \quad \sum_i \alpha_i y_i = 0, \quad \text{for all } i. \end{aligned}$$

- It is still a convex quadratic optimization problem.

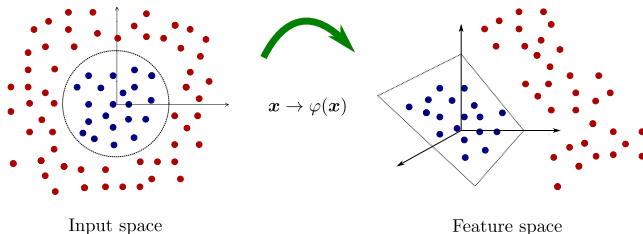
Non-linear SVM – Overview



How to obtain non-linear decision boundaries?

Higher dimension feature maps

Map features to higher dimensions



Here: $(x_1, x_2) \mapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$

The **input space** (original feature space) can always be mapped to some **higher-dimensional feature space** where the training data set is **linearly separable**, via some (non-linear) transformation $x \rightarrow \phi(x)$.

Kernel trick

- The mapping $\phi(\mathbf{x})$ may be **very difficult** to find and may require **exponentially more** parameters to learn
- However, in its dual formulation, SVM does not rely on the mapping itself, but only on the dot products between them $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$
- **The trick**: there exists¹ **kernel functions** such that

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

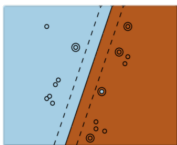
- The kernel function $K(\mathbf{x}, \mathbf{x}')$ can be used as a replacement, without knowing the actual mapping $\phi(\mathbf{x})$
- Generic tool: can also be applied to other algorithms (linear regression, logistic regression, etc.)

¹See Mercer's theorem

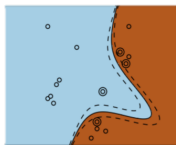
Examples of kernels

- Linear: $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$
- Polynomial: $K(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x}^\top \mathbf{x}' + \beta)^p$
- Gaussian (RBF): $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$
- Sigmoid: $K(\mathbf{x}, \mathbf{x}') = \tanh(\gamma \mathbf{x}^\top \mathbf{x}' + \beta)$

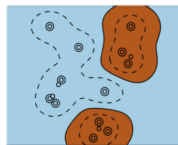
Linear Kernel



Polynomial Kernel



RBF Kernel

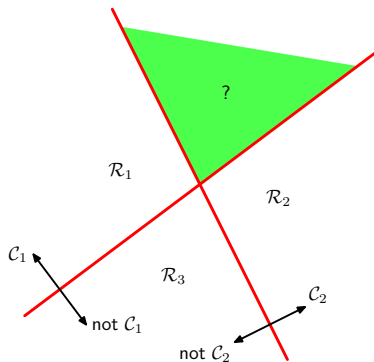


In practice, K is usually chosen with a validation set $\mathcal{D}_{\text{valid}}$

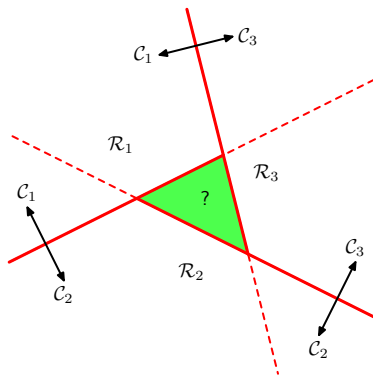
Extension to multiclass classification

Multiclass setting with K classes

One-vs-all



One-vs-one



- Indecision region: weight votes by probabilities $p(\mathcal{C}_k | \mathbf{x})$