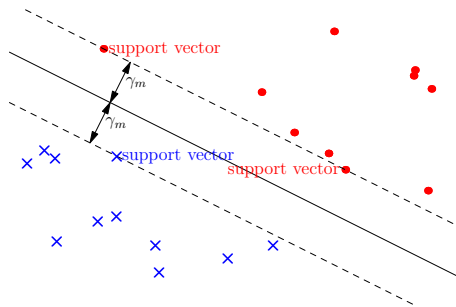


Constrained Optimisation, Lagrangians, Dual Algorithm

Linear SVM

- SVMs find hyperplane which separates two sets of linearly separable data with maximum margin



Extended Feature Spaces

- If we map into an extended feature space

$$\mathbf{x} = (x_1, x_2, \dots, x_p) \rightarrow \vec{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_r(\mathbf{x}))$$

$$r \gg p$$

- The optimisation problem becomes

$$\min_{\vec{w}, b} \frac{\|\vec{w}\|^2}{2} \quad \text{subject to } y_k (\vec{w}^T \vec{\phi}(\mathbf{x}_k) - b) \geq 1 \text{ for all } k = 1, 2, \dots, m$$

- \vec{w} is an r dimensional vector (lying in the extended feature space)

Constrained Optimisation

- Recall that when we try to solve an optimisation problem with constraints we can add Lagrange multipliers

- The optimisation problem becomes

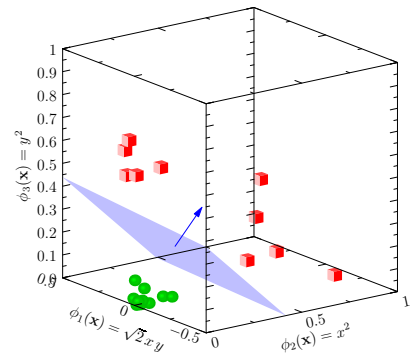
$$\min_{\vec{w}, b} \frac{\|\vec{w}\|^2}{2} \quad \text{subject to } y_k (\vec{w}^T \vec{\phi}(\mathbf{x}_k) - b) \geq 1 \text{ for all } k = 1, 2, \dots, m$$

- Is equivalent to finding the extremal point of the Lagrangian

$$\mathcal{L}(\vec{w}, b, \alpha) = \frac{\|\vec{w}\|^2}{2} + \sum_{k=1}^m \alpha_k (y_k (\vec{w}^T \vec{\phi}(\mathbf{x}_k) - b) - 1)$$

- subject to $\alpha_k \geq 0$

- Recap
- Constrained Optimisation
- Duality



Finding Maximum Margin Hyperplane

- We showed to find the maximum margin hyper-plane we can solve

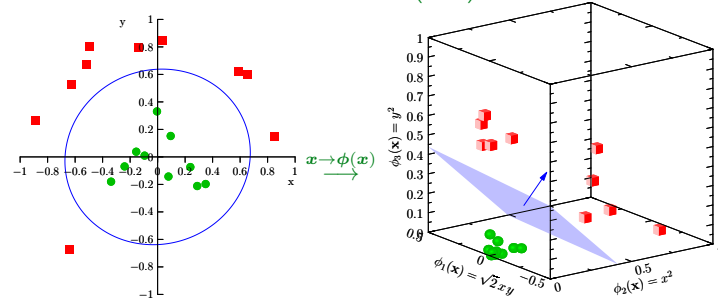
$$\min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2} \quad \text{subject to } y_k (\mathbf{w}^T \mathbf{x}_k - b) \geq 1 \text{ for all } k = 1, 2, \dots, m$$

(I've got rid of the hats on \mathbf{w} and b because life is too short)

- This is a quadratic programme (a quadratic function with linear constraints)
- Quadratic programmes have a unique solution
- This is generally true of convex function optimisation

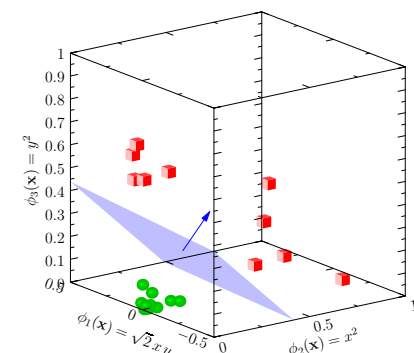
Non-linearly Separation of Data

$$\mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix} \rightarrow \vec{\phi}(\mathbf{x}) = \begin{pmatrix} 2xy \\ x^2 \\ y^2 \end{pmatrix}$$



Outline

- Recap
- Constrained Optimisation
- Duality



- Suppose we have a problem

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to } g(\mathbf{x}) = 0$$

- A standard procedure is to define the Lagrangian

$$\mathcal{L}(\mathbf{x}, \alpha) = f(\mathbf{x}) - \alpha g(\mathbf{x})$$

where α is known as a Lagrange multiplier

- In the extended space (\mathbf{x}, α) we have to solve

$$\max_{\alpha} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \alpha)$$

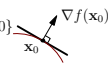
Note on Gradients

- Note that for any function $f(\mathbf{x})$ we can Taylor expand around \mathbf{x}_0

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \nabla_{\mathbf{x}} f(\mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^T \mathbf{H} (\mathbf{x} - \mathbf{x}_0) + \dots$$

where \mathbf{H} is a matrix of second derivative known as the Hessian

- If we consider points perpendicular to $\nabla_{\mathbf{x}} f(\mathbf{x}_0)$ which go through \mathbf{x}_0 these will have values

$$f(\mathbf{x}) = f(\mathbf{x}_0) + O(\|\mathbf{x} - \mathbf{x}_0\|^2)$$


thus $\nabla_{\mathbf{x}} f(\mathbf{x})$ is always orthogonal to the contour lines

Example

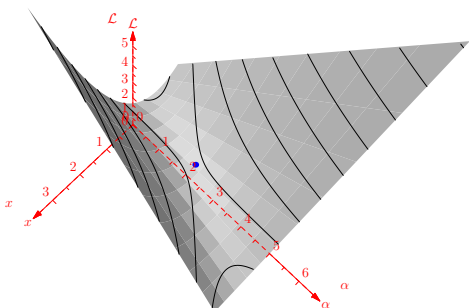
- Minimise $f(\mathbf{x}) = x^2 + 2y^2 - xy$
- Subject to $g(\mathbf{x}) = x - 2y - 3 = 0$
- Writing $\mathcal{L} = f(\mathbf{x}) - \alpha g(\mathbf{x})$
- Condition for minima is $\nabla_{\mathbf{x}} \mathcal{L} = 0$

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} 2x - y \\ -x + 4y \end{pmatrix} = \alpha \nabla_{\mathbf{x}} g(\mathbf{x}) = \alpha \begin{pmatrix} 1 \\ -2 \end{pmatrix}$$

$$\text{and } \frac{\partial \mathcal{L}}{\partial \alpha} = g(\mathbf{x}) = x - 2y - 3 = 0$$

- Solving simultaneous equations gives minima at $(x, y) = (\frac{3}{4}, -\frac{9}{8})$ with $\alpha = \frac{21}{8}$

Saddle-Point $y = -9/8$



- The optimisation problem is

$$\max_{\alpha} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \alpha) \quad \text{where } \mathcal{L}(\mathbf{x}, \alpha) = f(\mathbf{x}) - \alpha g(\mathbf{x})$$

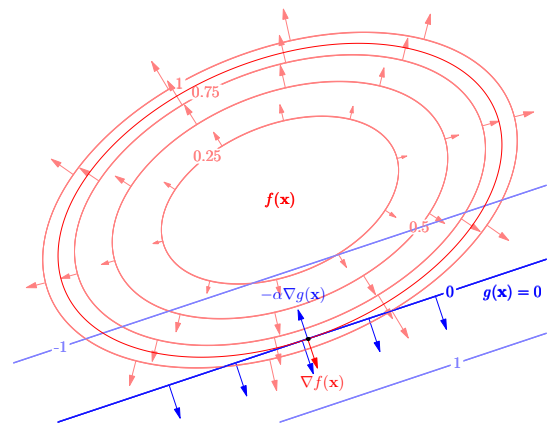
- Assuming differentiability

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \alpha) = \nabla_{\mathbf{x}} f(\mathbf{x}) - \alpha \nabla_{\mathbf{x}} g(\mathbf{x}) = 0$$

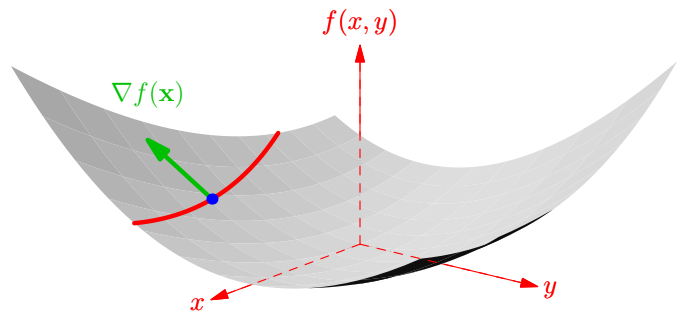
$$\frac{\partial \mathcal{L}}{\partial \alpha} = g(\mathbf{x}) = 0$$

- The second condition is just the constraint

Constrained Optima



Surface



Multiple Constraints

- Given an optimisation problem with multiple constraints

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to } g_k(\mathbf{x}) = 0 \text{ for } k = 1, 2, \dots, m$$

- We introduce multiple Lagrange multipliers

$$\mathcal{L}(\mathbf{x}, \alpha) = f(\mathbf{x}) - \sum_{k=1}^m \alpha_k g_k(\mathbf{x})$$

- The condition for an optima is $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \alpha) = 0$ which implies

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \sum_{k=1}^m \alpha_k \nabla_{\mathbf{x}} g_k(\mathbf{x})$$

plus the original constraints $\frac{\partial \mathcal{L}(\mathbf{x}, \alpha)}{\partial \alpha_k} = g_k(\mathbf{x}) = 0$

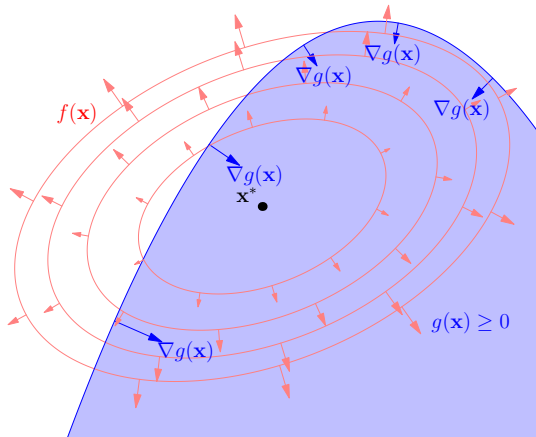
Example

- Minimise $f(\mathbf{x}) = x^2 + 2y^2 + 5z^2 - xy - xz$ subject to $g_1(\mathbf{x}) = x - 2y - z - 3 = 0$ and $g_2(\mathbf{x}) = 2x + 3y + z - 2 = 0$
- Writing $\mathcal{L}(\mathbf{x}, \alpha) = f(\mathbf{x}) - \alpha_1 g_1(\mathbf{x}) - \alpha_2 g_2(\mathbf{x})$
- Condition for minima is $\nabla_{\mathbf{x}} \mathcal{L} = 0$ or $\nabla_{\mathbf{x}} f(\mathbf{x}) = \sum_{k=1}^2 \alpha_k \nabla_{\mathbf{x}} g_k(\mathbf{x})$

$$\begin{pmatrix} 2x - y - z \\ -x + 4y \\ 10z - x \end{pmatrix} = \alpha_1 \begin{pmatrix} 1 \\ -2 \\ -1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix}$$

and $\frac{\partial \mathcal{L}}{\partial \alpha_i} = g_i(\mathbf{x}) = 0$
- Solving simultaneous equations gives minima at $(\frac{37}{20}, -\frac{11}{20}, -\frac{1}{20})$ with $\alpha_1 = 3$ and $\alpha_2 = \frac{13}{20}$

Inside Region



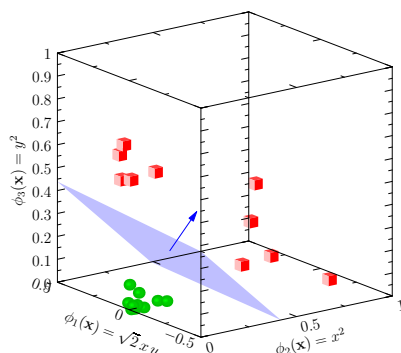
KKT Conditions

- To minimise $f(\mathbf{x})$ subject to $g(\mathbf{x}) \geq 0$

$$\mathcal{L}(\mathbf{x}, \alpha) = f(\mathbf{x}) - \alpha g(\mathbf{x})$$
- Then $\nabla_{\mathbf{x}} \mathcal{L} = 0$ or
$$\nabla_{\mathbf{x}} \mathcal{L} = \nabla_{\mathbf{x}} f(\mathbf{x}) - \alpha \nabla_{\mathbf{x}} g(\mathbf{x}) = 0$$
- where either
 - $\alpha = 0$ and the solutions in the interior or
 - $\alpha > 0$ and $g(\mathbf{x}) = 0$, i.e. the solution is on the boundary
- These conditions are known as the Karush-Kuhn-Tucker conditions

Outline

- Recap
- Constrained Optimisation
- Duality



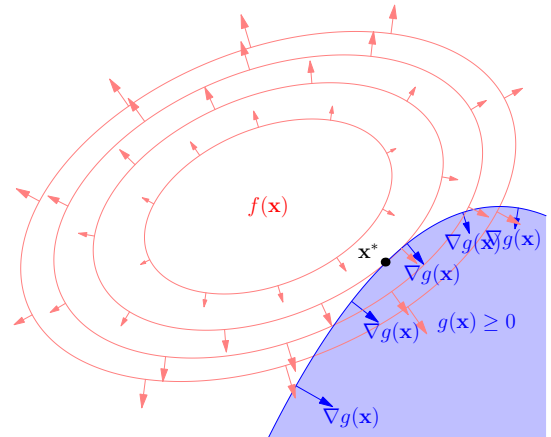
Inequality Constraints

- Suppose we have the problem

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to } g(\mathbf{x}) \geq 0$$

- Looks much more complicated, but
- Only two things can happen
 - Either the minimum of $f(\mathbf{x})$, \mathbf{x}^* , satisfies $g(\mathbf{x}^*) > 0$
 - We then have an unconstrained optimisation problem
 - Otherwise, it satisfies $g(\mathbf{x}^*) = 0$
 - We have a constrained optimisation problem

On the Boundary



Many Inequalities

- Given the problem
$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to } g_k(\mathbf{x}) \geq 0 \text{ for } k = 1, 2, \dots, m$$

- We introduce multiple Lagrange multipliers

$$\mathcal{L}(\mathbf{x}, \alpha) = f(\mathbf{x}) - \sum_{k=1}^m \alpha_k g_k(\mathbf{x})$$

- The condition for an optima is

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \sum_{k=1}^m \alpha_k \nabla_{\mathbf{x}} g_k(\mathbf{x})$$

- Plus the constraints that either $\alpha_k = 0$ or $\alpha_k > 0$ and $g_k(\mathbf{x}) = 0$

Back to the SVM

- We showed that the quadratic programming problem can be written as

$$\max_{\alpha} \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \alpha) \quad \text{subject to } \alpha_k \geq 0$$

- Where the Lagrangian is

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{k=1}^m \alpha_k (y_k (\mathbf{w}^T \mathbf{x}_k - b) - 1)$$

- $\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_{k=1}^m \alpha_k y_k \mathbf{x}_k = 0$ implies that $\mathbf{w}^* = \sum_{k=1}^m \alpha_k y_k \mathbf{x}_k$
- $\frac{\partial \mathcal{L}}{\partial b} = \sum_{k=1}^m \alpha_k y_k = 0$ implies $\sum_{k=1}^m \alpha_k y_k = 0$

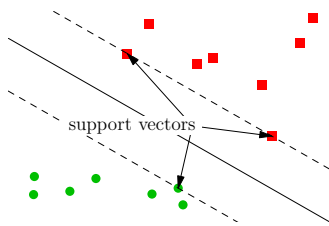
- Substituting in \mathbf{w}^* the optimisation condition becomes

$$\max_{\alpha} \sum_{k=1}^m \alpha_k - \frac{1}{2} \sum_{k,l=1}^m \alpha_k \alpha_l y_k y_l \mathbf{x}_k^{\top} \mathbf{x}_l$$

- Subject to $\sum_{k=1}^m \alpha_k y_k = 0$ and $\alpha_k \geq 0$ for all k
- This is also a quadratic programming problem!
- It is known as the *dual* form and depends on the number of training examples (it is solved by a standard package)
- It doesn't involve \mathbf{w} and involves only the dot products $\mathbf{x}_k^{\top} \mathbf{x}_l$

Support Vectors

- Where $\alpha_k > 0$ the constraints are exactly met so $y_k(\mathbf{x}_k^{\top} \mathbf{w} - b) = 1$
- These data points are known as *support vectors*



- In high dimensions there can be many support vectors

Dealing with Slack Variables

- Our new Lagrangian is

$$\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{k=1}^n \xi_k - \sum_{k=1}^m \alpha_k (y_k (\mathbf{w}^{\top} \mathbf{x}_k - b) - 1 + \xi_k) - \sum_{k=1}^m \beta_k \xi_k$$

- Minimise with respect to ξ_k

$$\frac{\partial \mathcal{L}}{\partial \xi_k} = C - \alpha_k - \beta_k = 0$$

- But $\beta_k \geq 0 \Rightarrow \alpha_k \leq C$
- Thus $0 \leq \alpha_k \leq C$

Classifying New Data

- Having trained the SVM we now have to use it
- Given a new input \mathbf{x} we decide on the class

$$\text{sgn}(\bar{\mathbf{w}}^{\top} \vec{\phi}(\mathbf{x}) - b) \quad \text{but} \quad \bar{\mathbf{w}} = \sum_{k=1}^m \alpha_k y_k \vec{\phi}(\mathbf{x}_k)$$

- In the dual representation this becomes

$$\text{sgn} \left(\sum_{k=1}^m \alpha_k y_k K(\mathbf{x}_k, \mathbf{x}) - b \right)$$

where we only need to sum over the non-zero α_k (i.e. the support vectors SVs)

- We can either work in the n -dimensional weight space

$$\min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2} \quad \text{subject to } y_k (\mathbf{w}^{\top} \mathbf{x}_k - b) \geq 1 \text{ for all } k = 1, 2, \dots, m$$

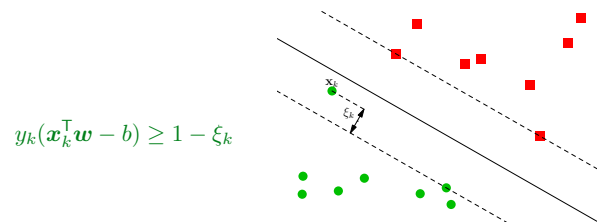
- Or the m -dimensional Lagrange multiplier space

$$\max_{\alpha} \sum_{k=1}^m \alpha_k - \frac{1}{2} \sum_{k,l=1}^m \alpha_k \alpha_l y_k y_l \mathbf{x}_k^{\top} \mathbf{x}_l \quad \text{with } \sum_{k=1}^m \alpha_k y_k = 0 \text{ \& } \alpha_k \geq 0$$

- Both require sophisticated quadratic programming algorithms to solve
- Which is easiest depends on whether m is greater or less than p

Soft Margins

- Recall we can relax constraints by introducing *slack variables*, $\xi_k \geq 0$



- Minimise $\frac{\|\mathbf{w}\|^2}{2} + C \sum_{k=1}^m \xi_k$
- subject to $\xi_k \geq 0$

Extended Feature Space

- In the extended feature space

$$\min_{\vec{\mathbf{w}}, b} \frac{\|\vec{\mathbf{w}}\|^2}{2} \quad \text{subject to } y_k (\vec{\mathbf{w}}^{\top} \vec{\phi}(\mathbf{x}_k) - b) \geq 1 \text{ for all } k = 1, 2, \dots, m$$

- Giving the dual problems

$$\max_{\alpha} \sum_{k=1}^m \alpha_k - \frac{1}{2} \sum_{k,l=1}^m \alpha_k \alpha_l y_k y_l \vec{\phi}(\mathbf{x}_k)^{\top} \vec{\phi}(\mathbf{x}_l) \quad \text{with } \sum_{k=1}^m \alpha_k y_k = 0 \text{ \& } \alpha_k \geq 0$$

- When $\vec{\phi}(\mathbf{x}_k)^{\top} \vec{\phi}(\mathbf{x}_l) = K(\mathbf{x}_k, \mathbf{x}_l)$ then

$$\max_{\alpha} \sum_{k=1}^m \alpha_k - \frac{1}{2} \sum_{k,l=1}^m \alpha_k \alpha_l y_k y_l K(\mathbf{x}_k, \mathbf{x}_l) \quad \text{with } \sum_{k=1}^m \alpha_k y_k = 0 \text{ \& } \alpha_k \geq 0$$

Conclusion

- We can solve for the maximum-margin hyper-plane either in the primal form (space of weights and bias) or the dual form (space of Lagrange multipliers)
- In the dual form the solution only depends on the dot product $\mathbf{x}_i^{\top} \mathbf{x}_j$ or $\vec{\phi}(\mathbf{x}_i)^{\top} \vec{\phi}(\mathbf{x}_j)$
- If $K(\mathbf{x}, \mathbf{y}) = \vec{\phi}(\mathbf{x})^{\top} \vec{\phi}(\mathbf{y})$ we never have to explicitly compute $\vec{\phi}(\mathbf{x})$ —the kernel trick
- This allows us to work in a very high dimensional space (giving low bias) while finding a maximum-margin hyper-plane (simple machine with low variance) magic!