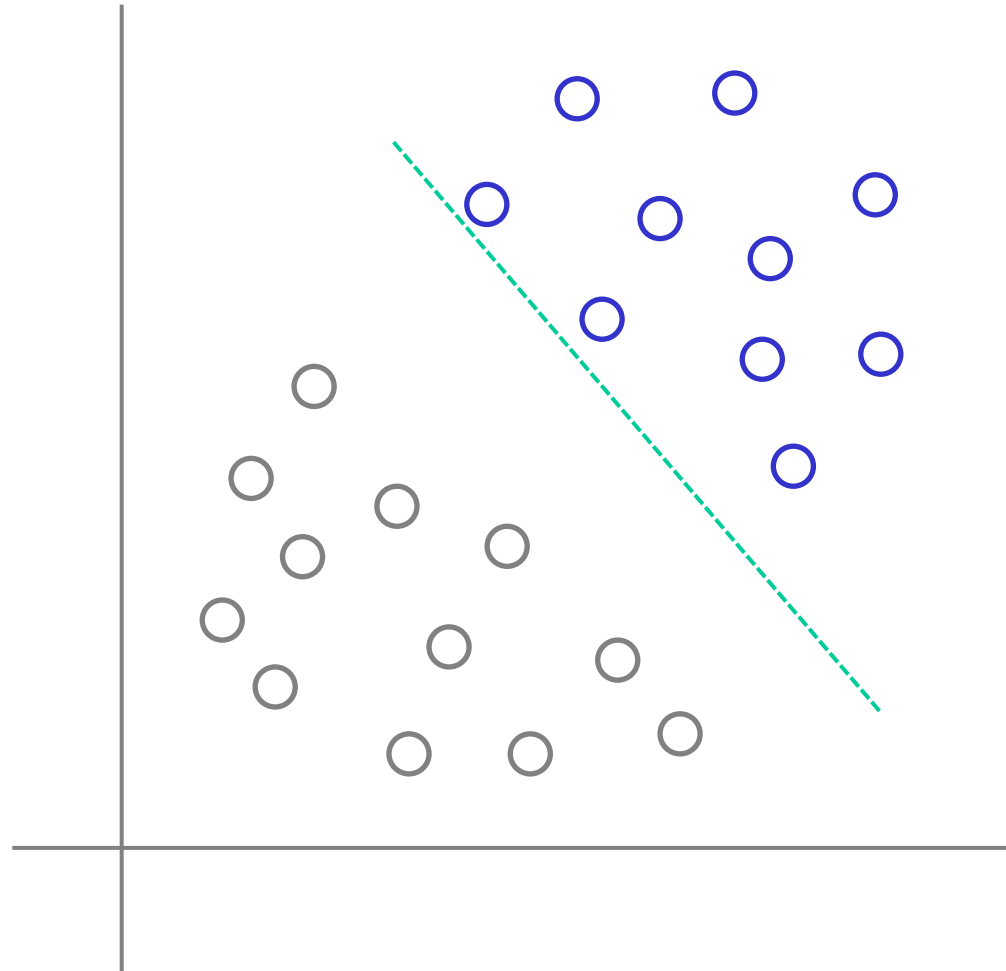


CSE 5526: Introduction to Neural Networks

Support Vector Machines (SVM)

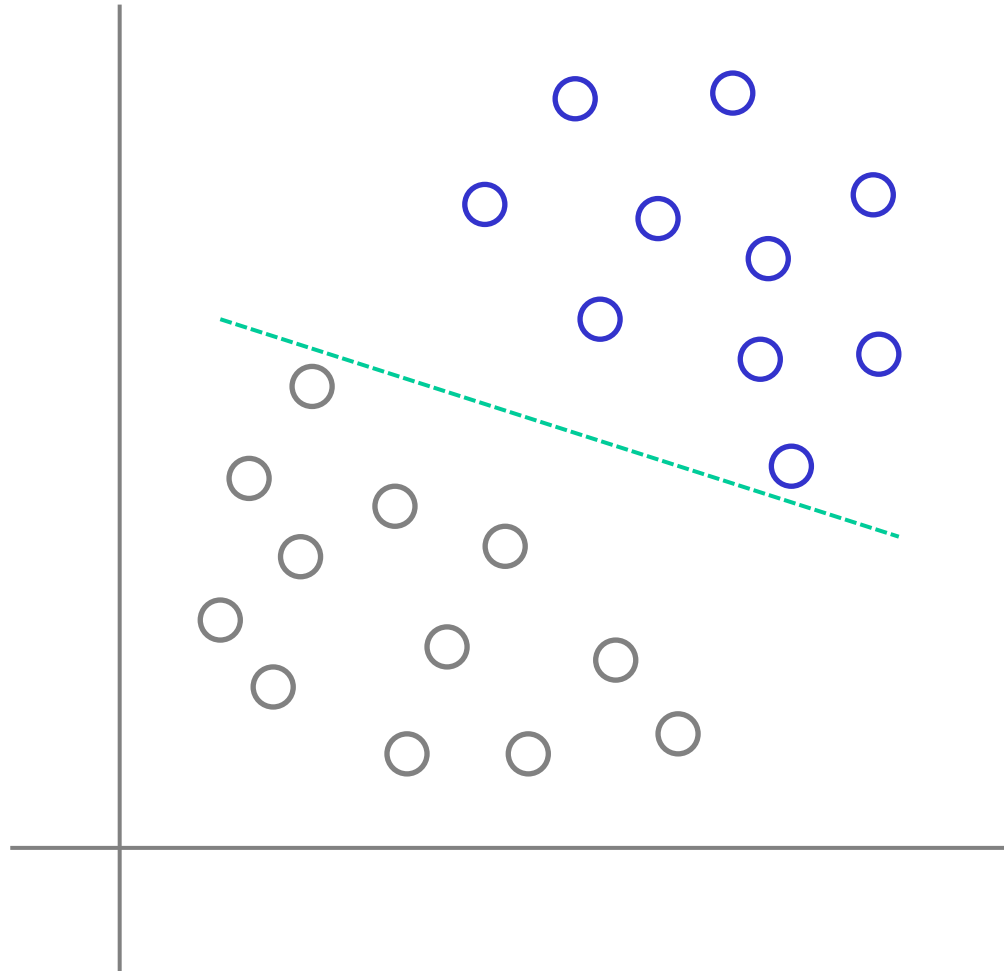
Perceptrons find any separating hyperplane

Depends on initialization and ordering of training points



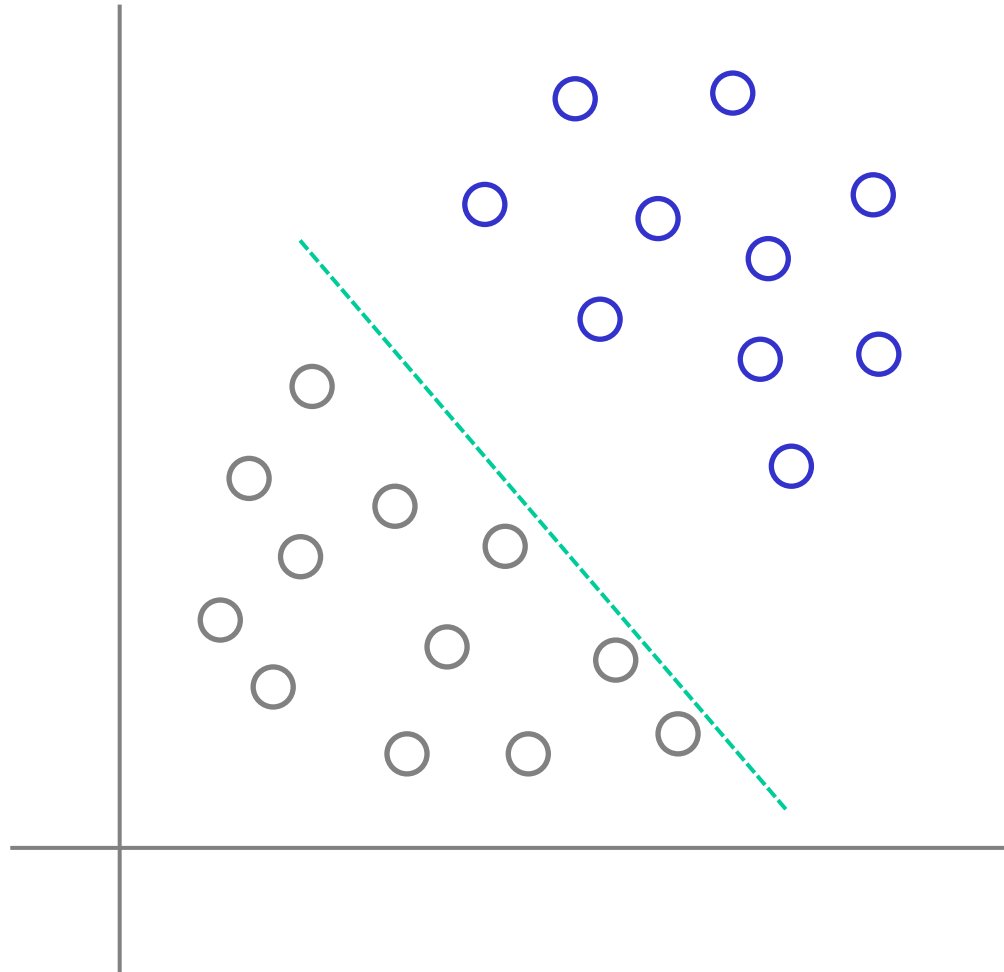
Perceptrons find any separating hyperplane

Depends on initialization and ordering of training points



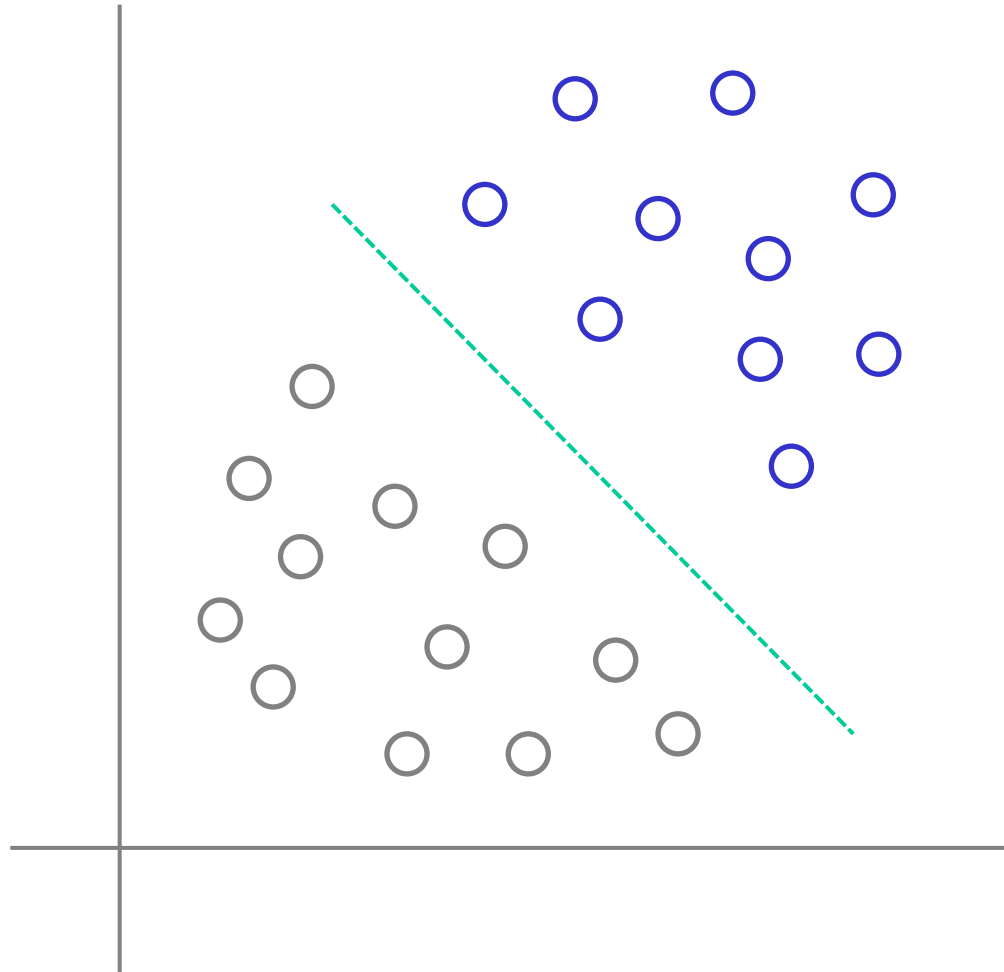
Perceptrons find any separating hyperplane

Depends on initialization and ordering of training points



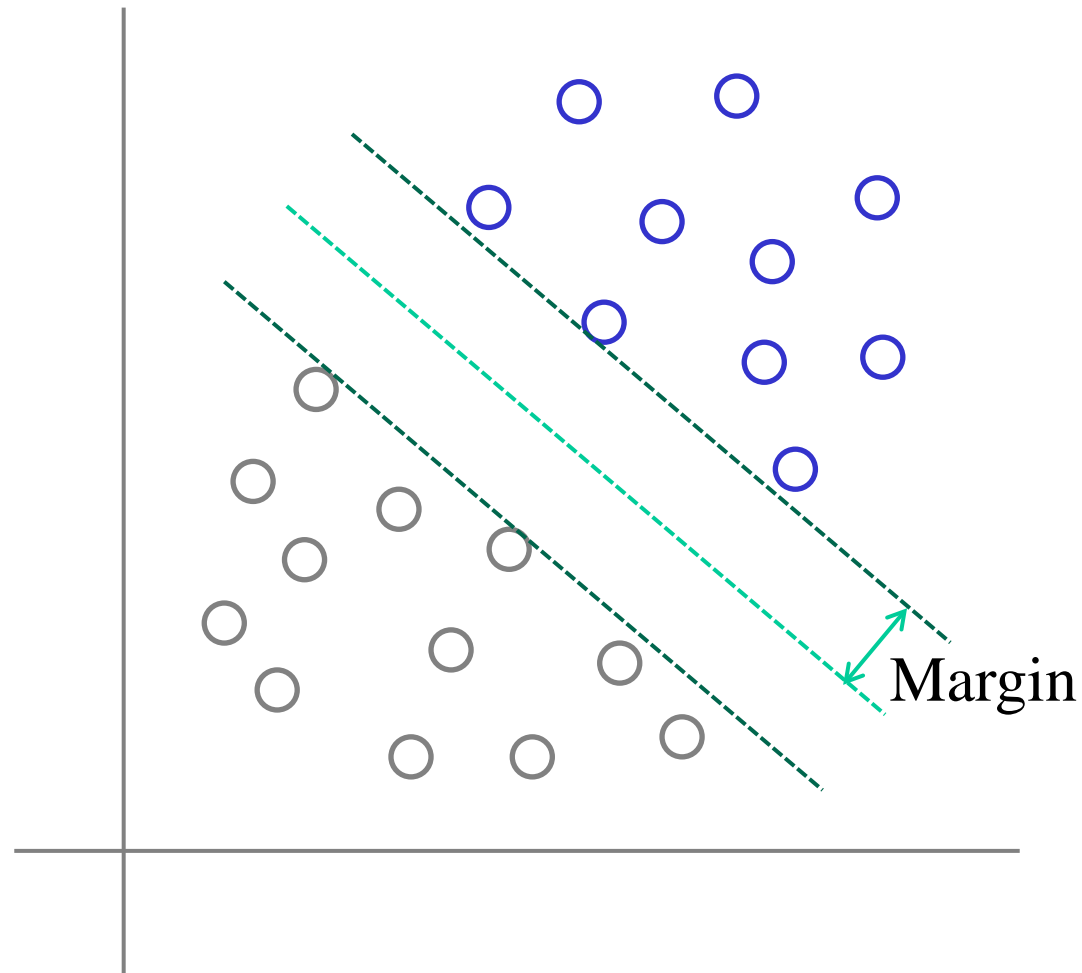
Perceptrons find any separating hyperplane

Depends on initialization and ordering of training points



But the maximum margin hyperplane generalizes the best to new data

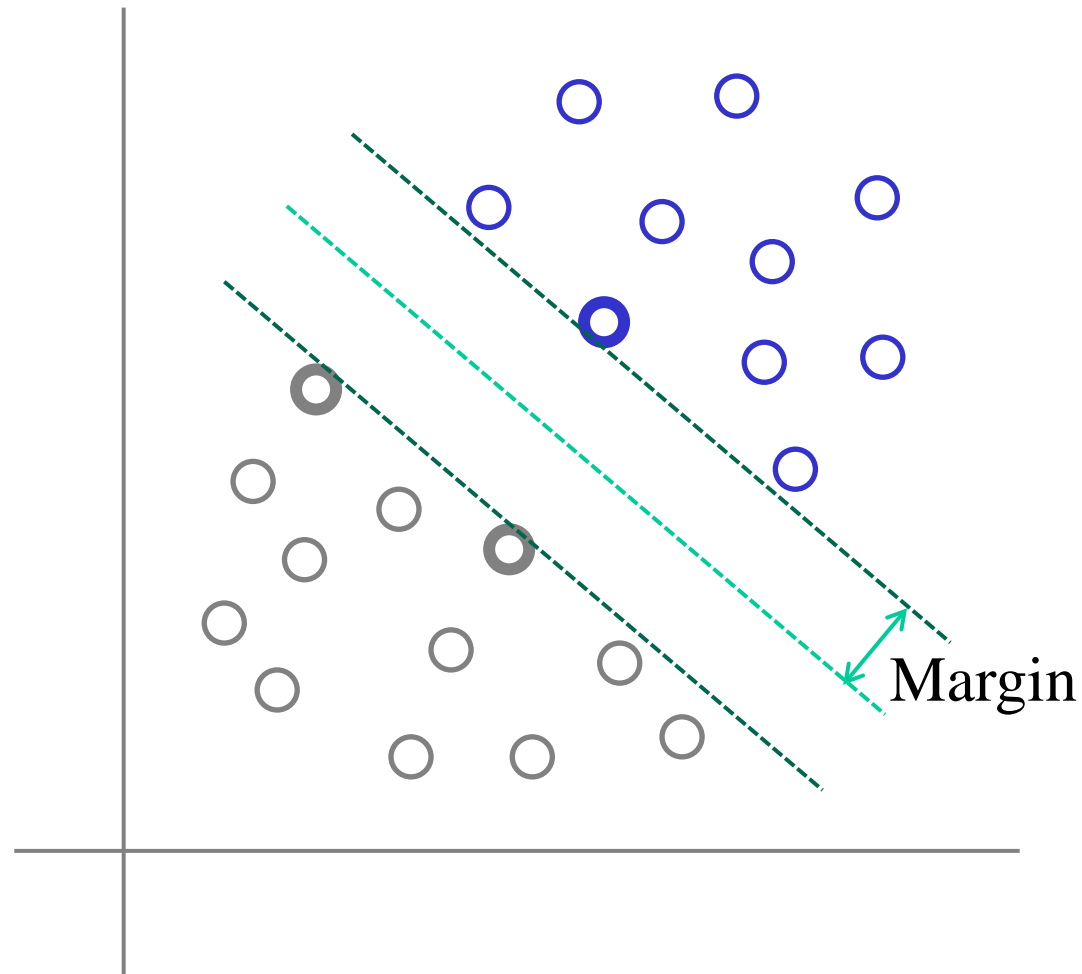
According to computational/statistical learning theory



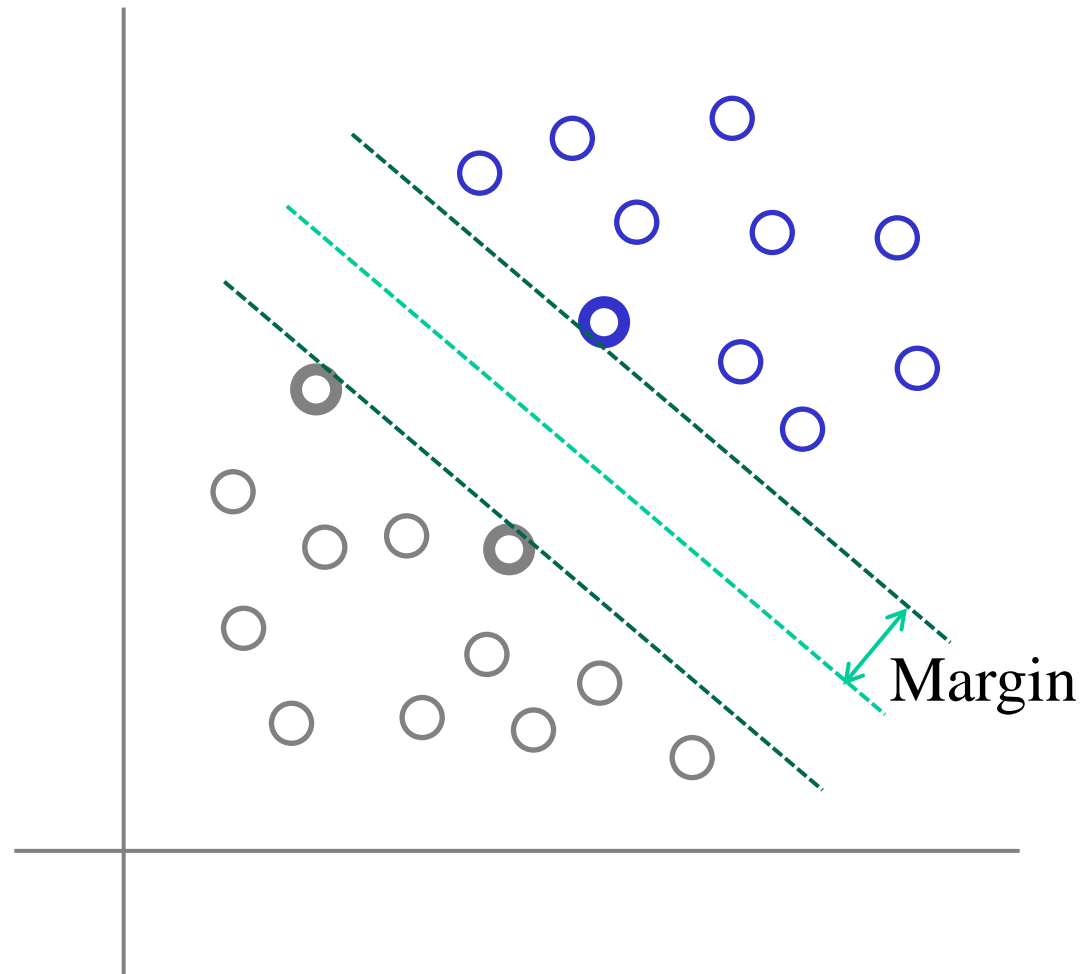
But the maximum margin hyperplane generalizes the best to new data

- According to computational learning theory
- Also known as statistical learning theory
- We won't get into the details of that
- Recall from the perceptron convergence proof
 - We assumed the existence of a best hyperplane \mathbf{w}_0
 - Which provided the maximum margin α
 - Such that $d_p \mathbf{w}_0^T \mathbf{x}_p \geq \alpha$ for all training points p
- The SVM actually finds this hyperplane

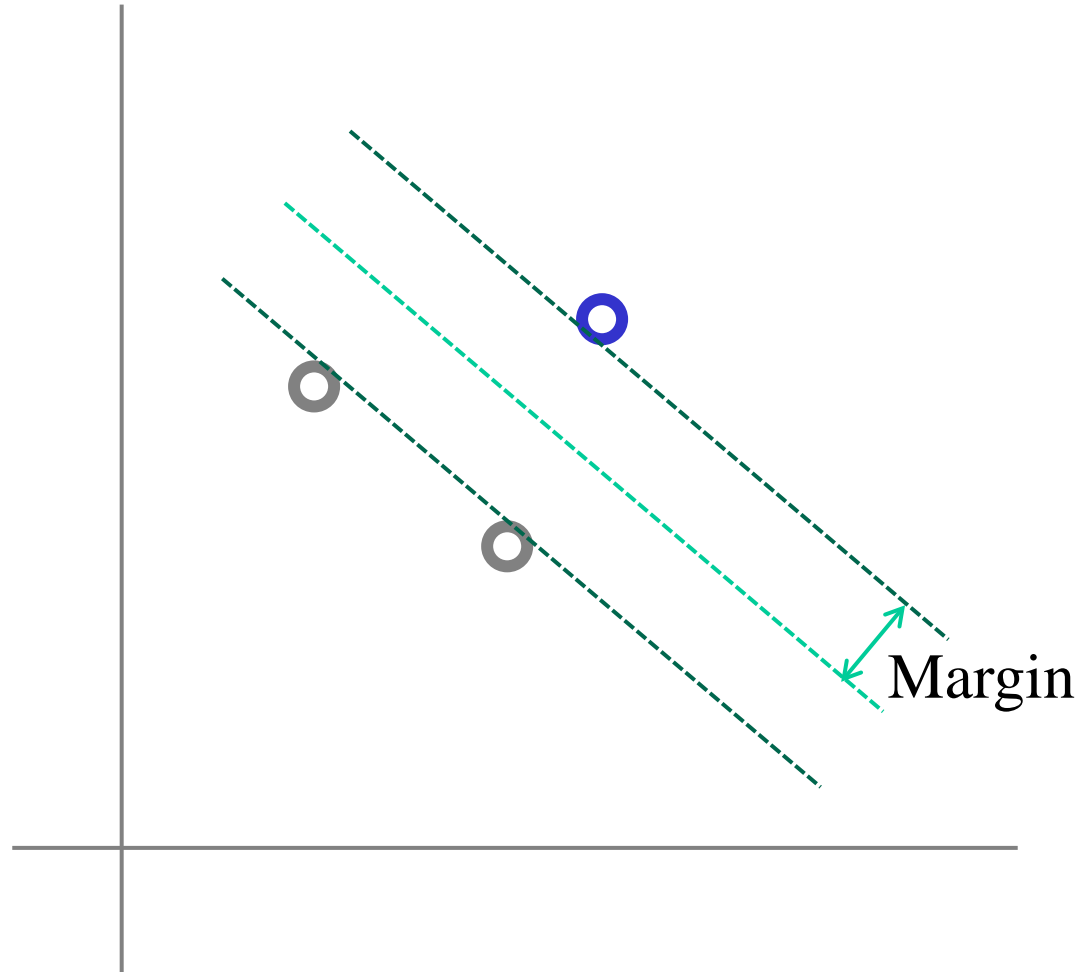
The maximum margin only depends on certain points, the support vectors



The maximum margin only depends on certain points, the support vectors



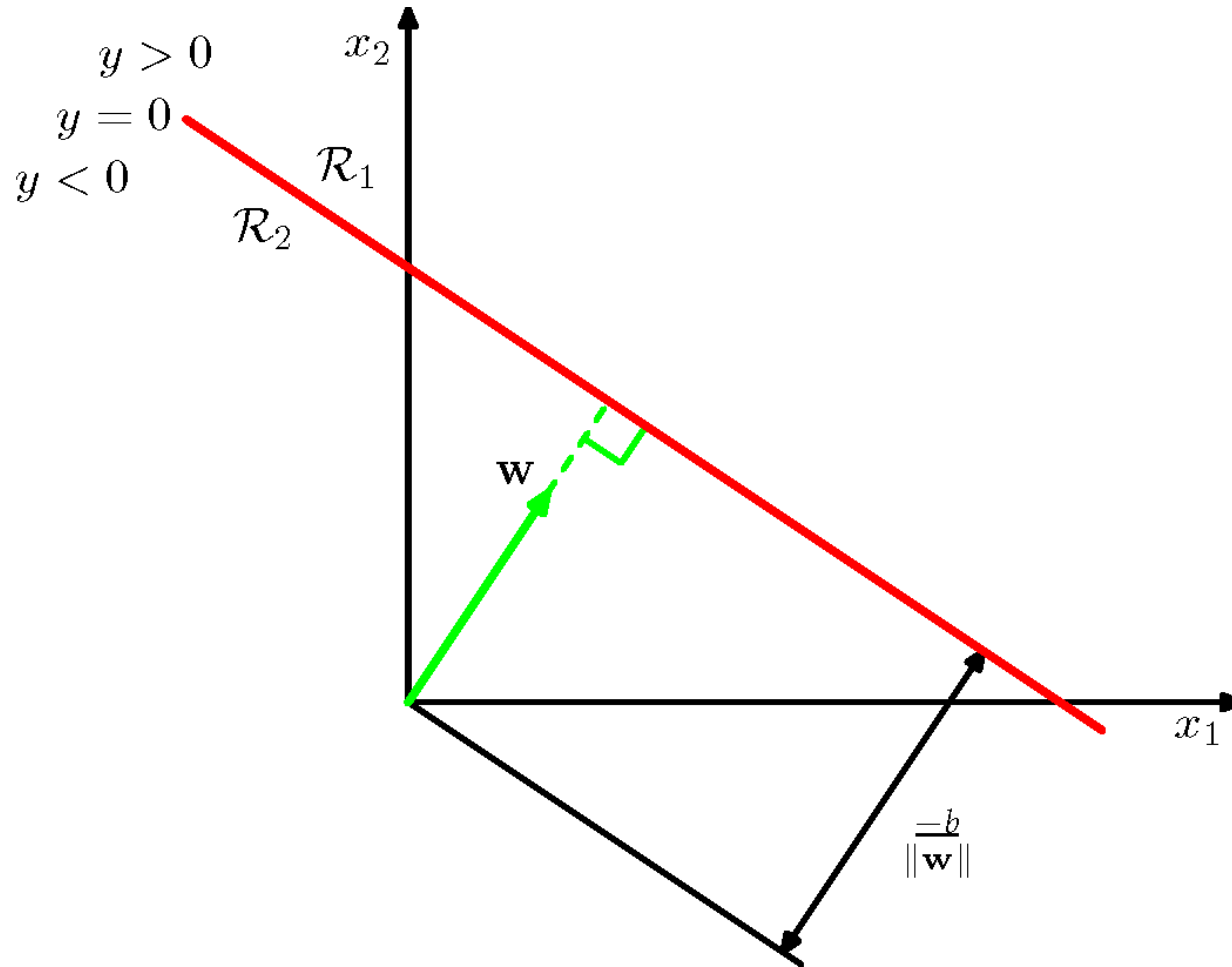
The maximum margin only depends on certain points, the support vectors



Maximum margin problem

- Given a set of data from two classes $\{\mathbf{x}_p, d_p\}$
 - $\mathbf{x}_p \in \mathbb{R}^D$ and $d_p \in \{-1, 1\}$
 - Assume the classes are linearly separable for now
- Find the hyperplane that separates them
 - with maximum margin
- Equation of general linear discriminant function
$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$
- Find \mathbf{w} and b that give maximum margin
- How can we quantify margin?

\mathbf{w} is perpendicular to the hyperplane,
 b defines its distance from the origin



\mathbf{w} is perpendicular to the hyperplane

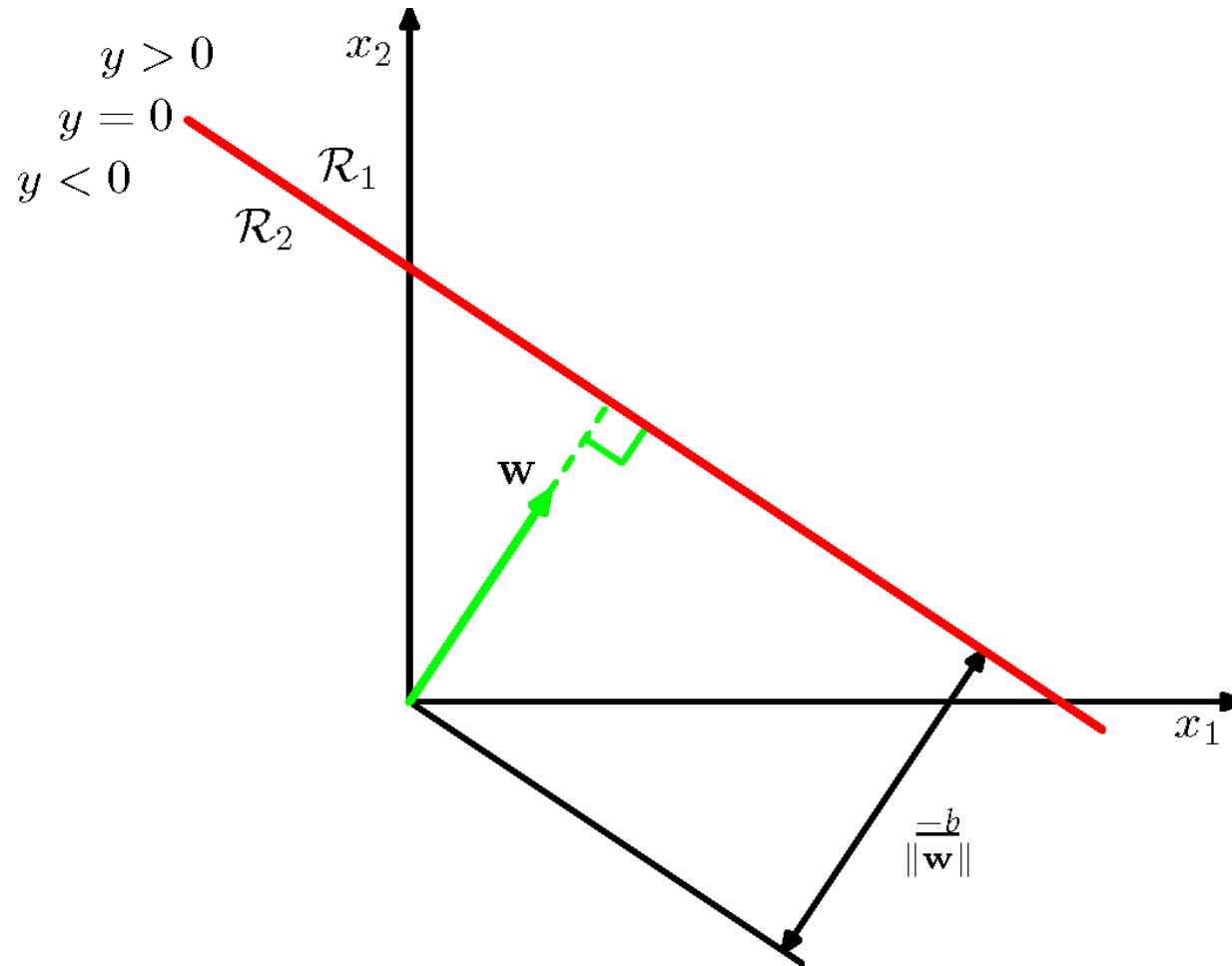
- Consider two points on the hyperplane \mathbf{x}_A and \mathbf{x}_B
- Then $y(\mathbf{x}_A) = y(\mathbf{x}_B) = 0$ by definition
- So $0 = y(\mathbf{x}_A) - y(\mathbf{x}_B) = \mathbf{w}^T (\mathbf{x}_A - \mathbf{x}_B)$
- $\mathbf{x}_A - \mathbf{x}_B$ is a vector pointing along the hyperplane
- So \mathbf{w} is perpendicular to the hyperplane

b defines the hyperplane's distance from the origin

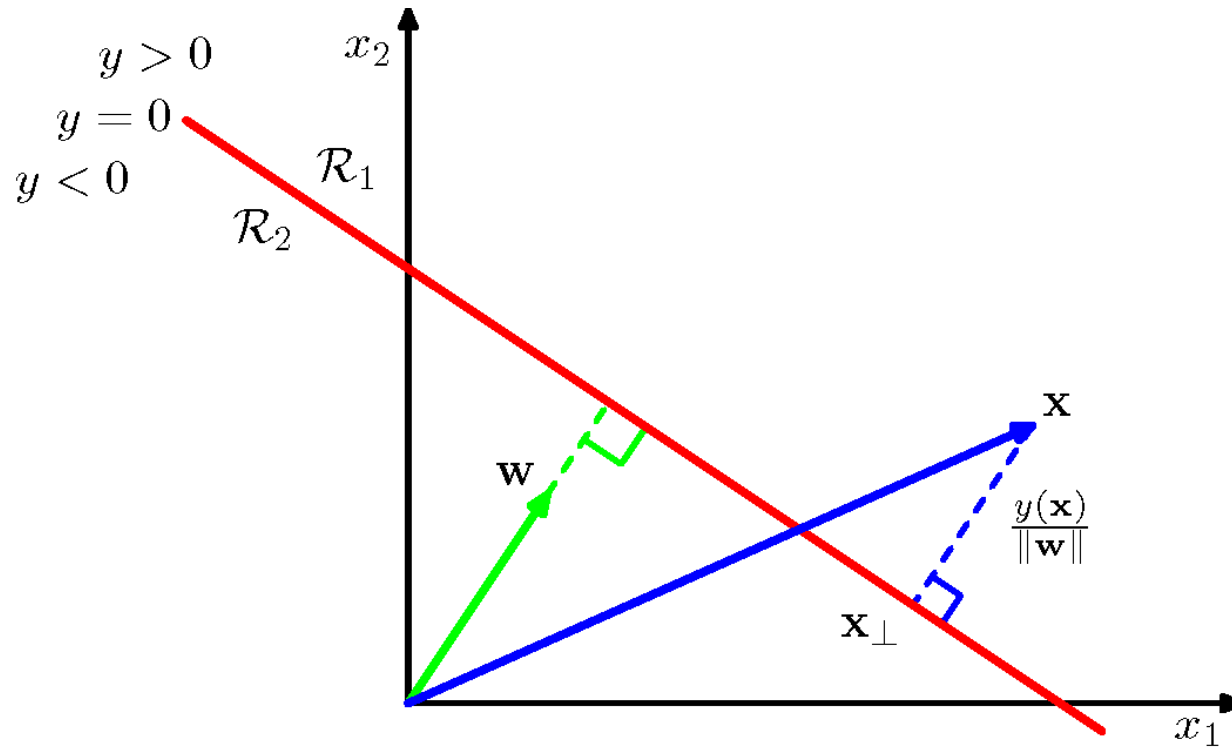
- Consider a general point \mathbf{x}
- Its distance to the origin is $D = \frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|}$
- If \mathbf{x} is on the hyperplane, then $y(\mathbf{x}) = 0$
 - So $\mathbf{w}^T \mathbf{x} = -b$
- So the distance from the hyperplane to the origin is

$$D = -\frac{b}{\|\mathbf{w}\|}$$

\mathbf{w} is perpendicular to the hyperplane,
 b defines its distance from the origin



The distance from point \mathbf{x}
to the hyperplane is $y(\mathbf{x})/\|\mathbf{w}\|$



The distance from point \mathbf{x}
to the hyperplane is $y(\mathbf{x})/\|\mathbf{w}\|$

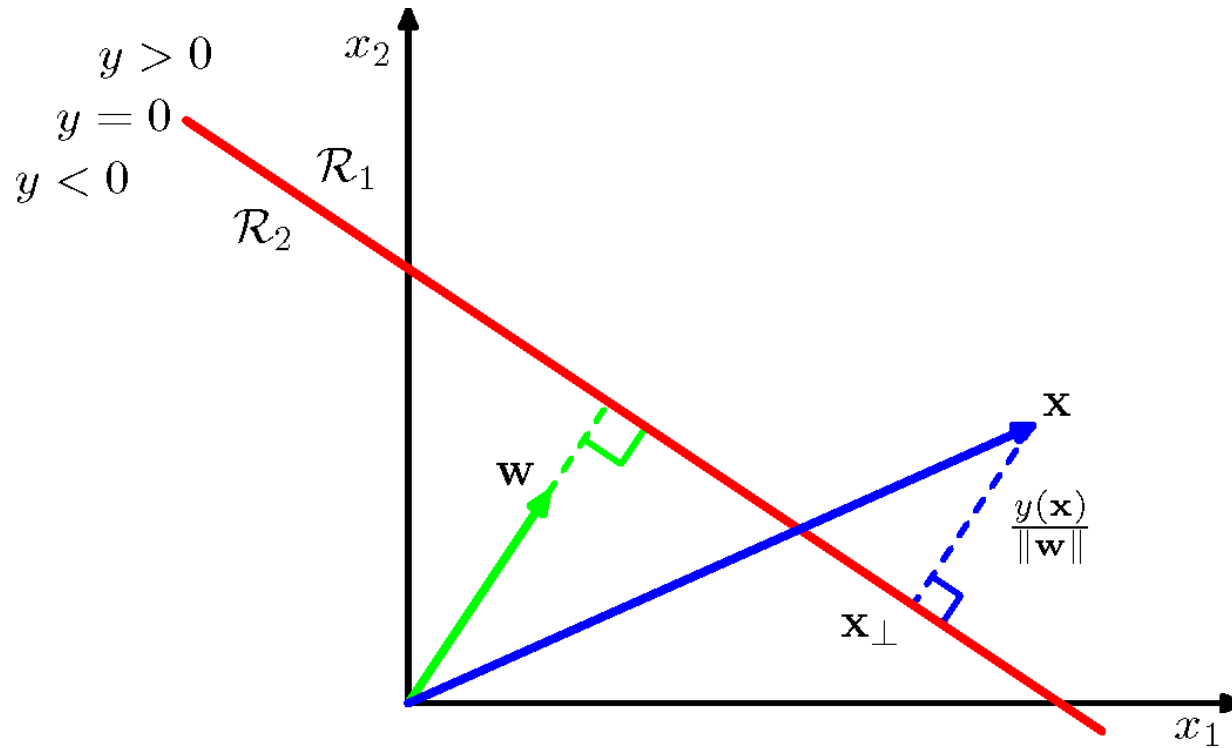
- Consider a point \mathbf{x} and its projection onto the hyperplane \mathbf{x}_\perp so that $\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$
- We want to find r , the distance to the hyperplane
- Multiply both sides by \mathbf{w}^T and add b

$$\mathbf{w}^T \mathbf{x} + b = \mathbf{w}^T \mathbf{x}_\perp + b + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|}$$

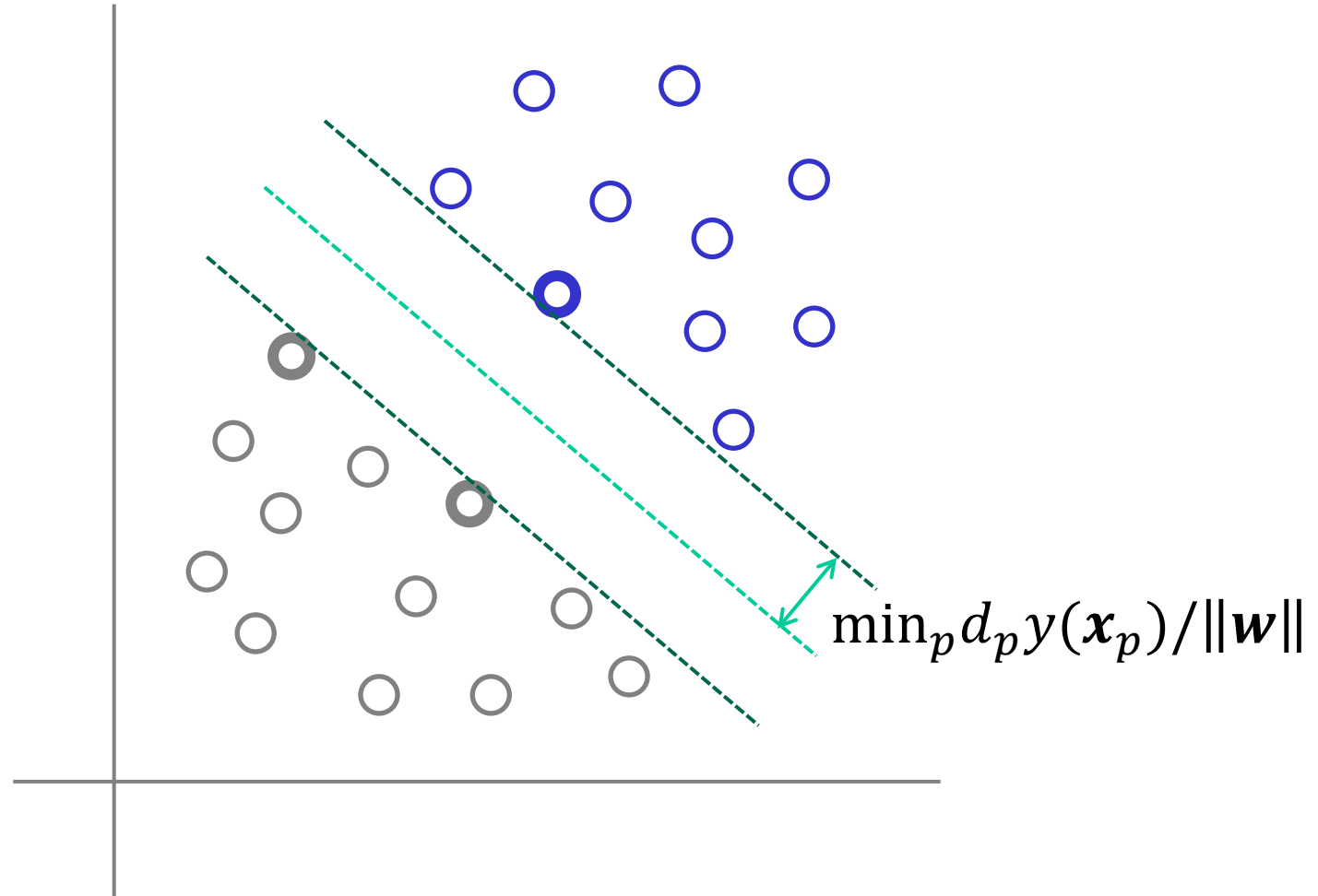
$$y(\mathbf{x}) = y(\mathbf{x}_\perp) + r \|\mathbf{w}\|$$

$$r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

The distance from point \mathbf{x}
to the hyperplane is $y(\mathbf{x})/\|\mathbf{w}\|$



The maximum margin hyperplane is farthest from all of the data points



The maximum margin hyperplane is farthest from all of the data points

- The margin is defined as

$$\begin{aligned}\alpha &= \min_p d_p \frac{y(\mathbf{x}_p)}{\|\mathbf{w}\|} \\ &= \frac{1}{\|\mathbf{w}\|} \min_p d_p y(\mathbf{x}_p)\end{aligned}$$

- We want to find \mathbf{w} and b that maximize the margin

$$\operatorname{argmax}_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \min_p d_p y(\mathbf{x}_p)$$

- Solving this problem is hard as it is written

We are free to choose a rescaling of \mathbf{w}

- If we replace \mathbf{w} by $a\mathbf{w}$ and b with ab
- Then the margin is unchanged

$$\min_p d_p \frac{a\mathbf{w}^T \mathbf{x}_p + ab}{a\|\mathbf{w}\|} = \min_p d_p \frac{\mathbf{w}^T \mathbf{x}_p + b}{\|\mathbf{w}\|}$$

- So choose a such that $\min_p d_p (a\mathbf{w}^T \mathbf{x}_p + ab) = 1$
- Which means that for all points

$$d_p (\mathbf{w}^T \mathbf{x}_p + b) \geq 1$$

Maximum margin constrained optimization problem

- Then the maximum margin optimization becomes

$$\begin{aligned} \operatorname{argmax}_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \min_p d_p y(x_p) \\ = \operatorname{argmax}_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \end{aligned}$$

- With the constraints $d_p(\mathbf{w}^T \mathbf{x}_p + b) \geq 1$

Maximum margin constrained optimization problem

- Which is equivalent to

$$\operatorname{argmin}_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to } d_p(\mathbf{w}^T \mathbf{x}_p + b) \geq 1$$

- This is a well studied type of problem
 - A quadratic program with linear inequality constraints

Detour: Lagrange multipliers

solve constrained optimization problems

- Want to maximize a function $f(x_1, x_2)$
- Subject to the equality constraint $g(x_1, x_2) = 0$
- Could solve $g(x_1, x_2) = 0$ for x_1 in terms of x_2
 - But that is hard to do in general (i.e., on computers)
- Or could use Lagrange multipliers
 - Which are easier to use in general (i.e., on computers)

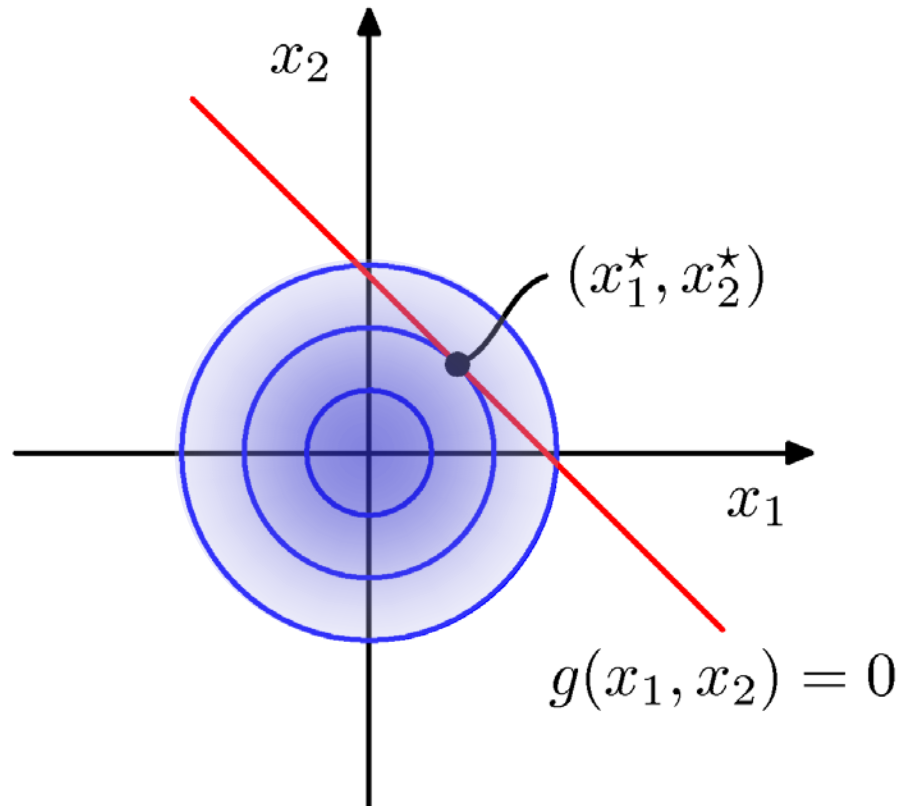
Lagrange multipliers with general \mathbf{x}

- In general, we can write

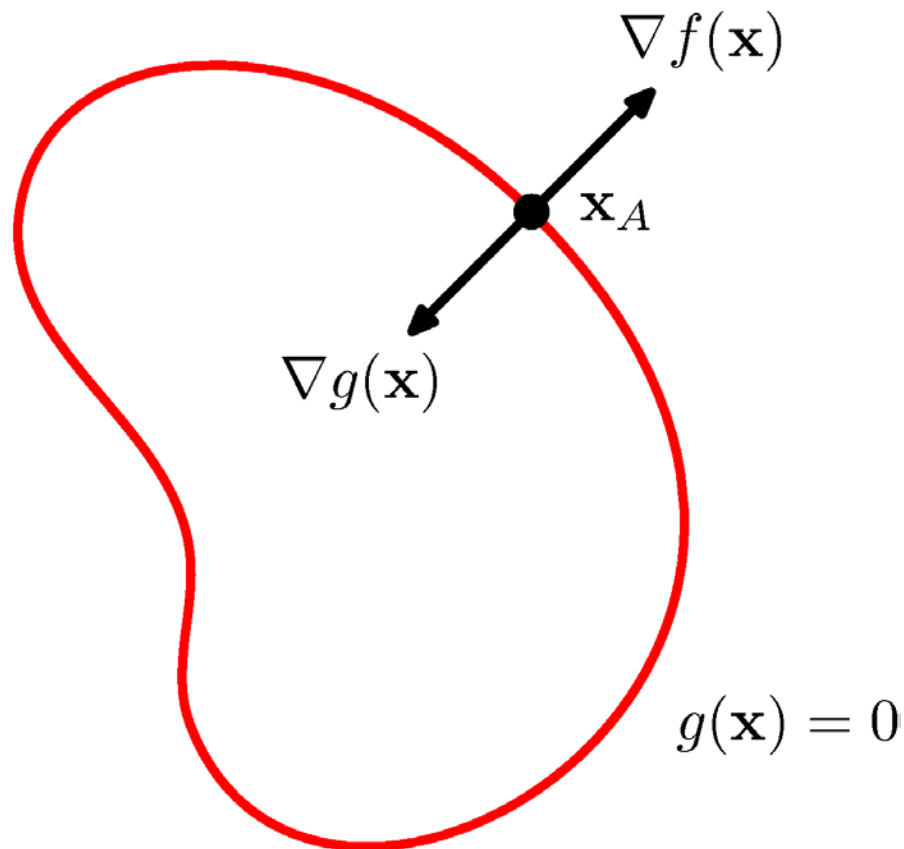
$$\max_{\mathbf{x}} f(\mathbf{x}) \text{ subject to } g(\mathbf{x}) = 0$$

- Constraint $g(\mathbf{x}) = 0$ defines a $D - 1$ dimensional surface for D dimensional \mathbf{x}

Example: Maximize $f(\mathbf{x}) = 1 - x_1^2 - x_2^2$
subject to $g(\mathbf{x}) = x_1 + x_2 - 1 = 0$



Gradients of g and f
are orthogonal to surface at solution point



Gradients of g and f are orthogonal to surface at maximum of f

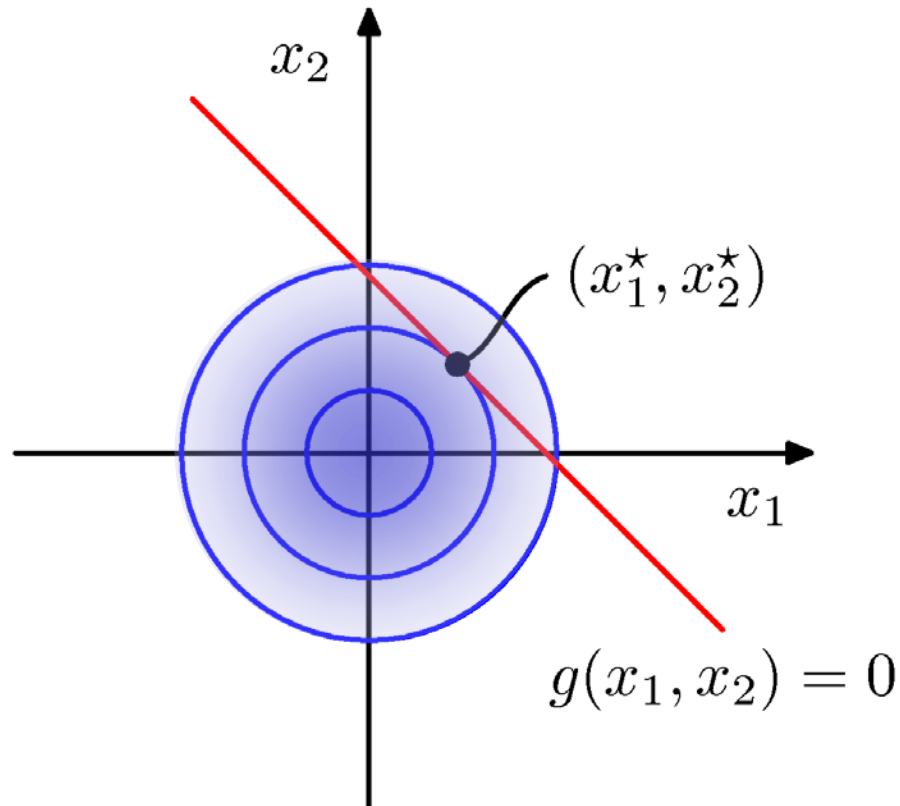
- For g because on all points on the surface $g(\mathbf{x}) = 0$
- For f because if it wasn't, you could move along the surface in the direction of the gradient to find a better maximum
- Thus ∇f and ∇g are (anti-)parallel
- And there must exist a scalar λ such that
$$\nabla f + \lambda \nabla g = 0$$

The Lagrangian function captures the constraints on \mathbf{x} and on the gradients

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

- Setting gradient of L with respect to \mathbf{x} to 0 gives
$$\nabla f + \lambda \nabla g = 0$$
- Setting partial of L with respect to λ to 0 gives
$$g(\mathbf{x}) = 0$$
- Thus stationary points of L solve the constrained optimization problem

Example: Maximize $f(\mathbf{x}) = 1 - x_1^2 - x_2^2$
subject to $g(\mathbf{x}) = x_1 + x_2 - 1 = 0$



Example: Maximize $f(\mathbf{x}) = 1 - x_1^2 - x_2^2$
subject to $g(\mathbf{x}) = x_1 + x_2 - 1 = 0$

- So the Lagrangian function is

$$\begin{aligned} L(\mathbf{x}, \lambda) &= f(\mathbf{x}) + \lambda g(\mathbf{x}) \\ &= 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1) \end{aligned}$$

- The conditions for L to be stationary are

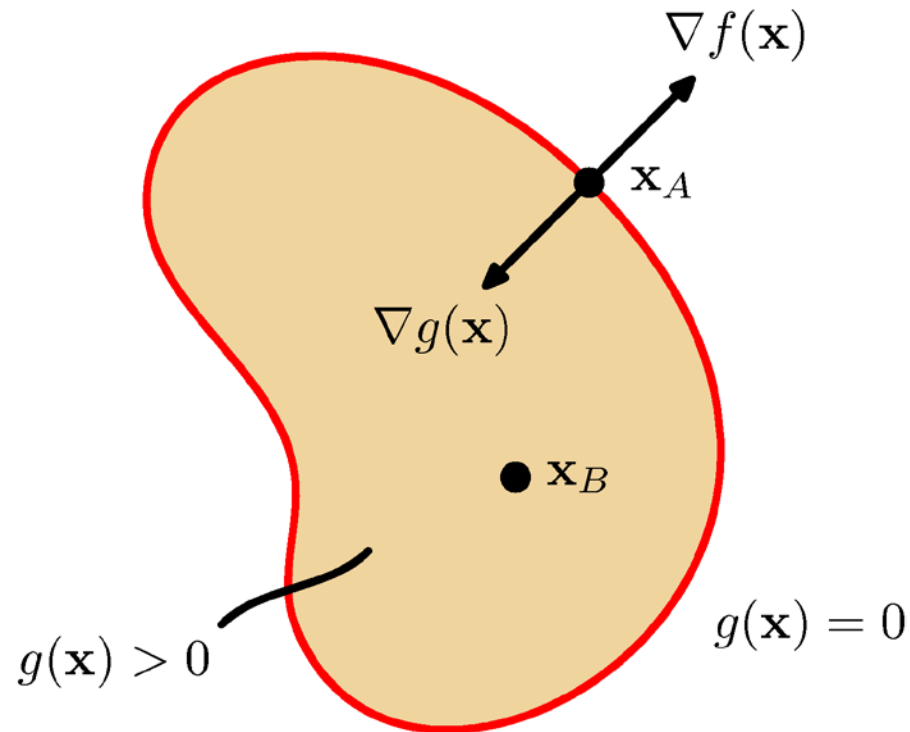
$$\partial L / \partial x_1 = -2x_1 + \lambda = 0$$

$$\partial L / \partial x_2 = -2x_2 + \lambda = 0$$

$$\partial L / \partial \lambda = x_1 + x_2 - 1 = 0$$

- Can solve to find $\lambda = 1, x_1 = x_2 = \frac{1}{2}$

Lagrange multipliers can also be used with inequality constraints $g(\mathbf{x}) \geq 0$



Lagrange multipliers can also be used with inequality constraints $g(\mathbf{x}) \geq 0$

- Now two kinds of solutions:
- If $g(\mathbf{x}) > 0$, then the solution only depends on $f(\mathbf{x})$
 - Inside constraint surface with $\nabla f = 0$
 - Stationary point of $L(\mathbf{x}, \lambda)$ with $\lambda = 0$
 - Constraint $g(\mathbf{x})$ is said to be inactive
- If $g(\mathbf{x}) = 0$, then same as before (with equality constraint)
 - On boundary of constraint surface with ∇f pointing out
 - Stationary point of $L(\mathbf{x}, \lambda)$ with $\lambda > 0$
 - Constraint $g(\mathbf{x})$ is said to be active

Lagrange multipliers can also be used with inequality constraints $g(\mathbf{x}) \geq 0$

- In either case, $\lambda g(\mathbf{x}) = 0$
- Thus maximizing $f(\mathbf{x})$ subject to $g(\mathbf{x}) \geq 0$ is obtained by optimizing $L(\mathbf{x}, \lambda)$ WRT \mathbf{x} and λ subject to

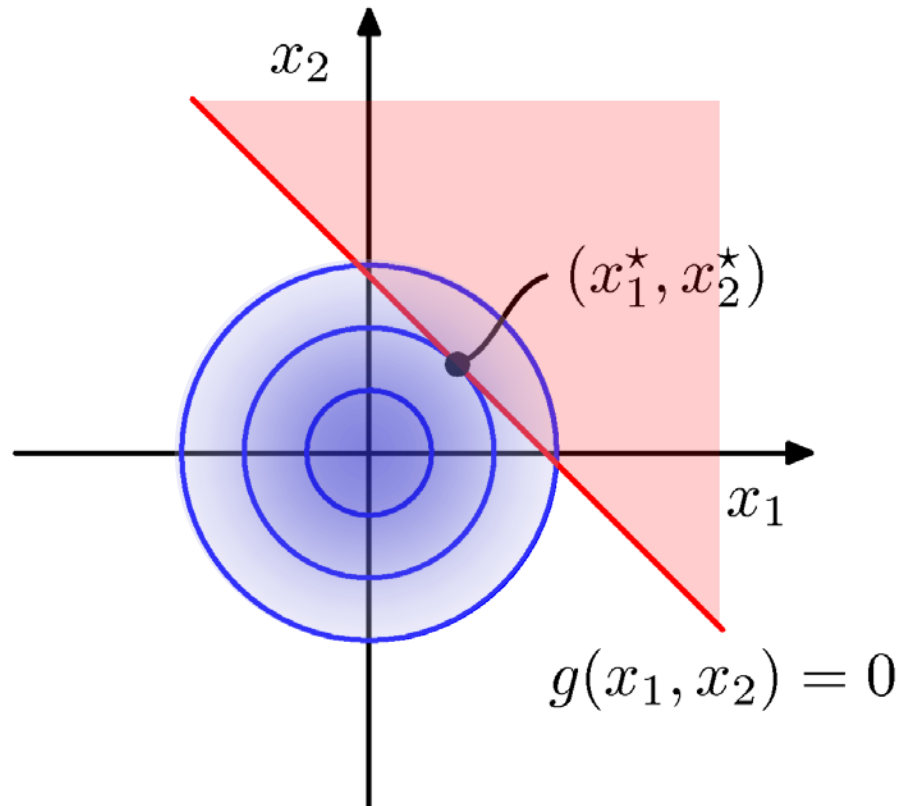
$$g(\mathbf{x}) \geq 0$$

$$\lambda \geq 0$$

$$\lambda g(\mathbf{x}) = 0$$

- These are known as the Karush-Kuhn-Tucker (KKT) conditions

Example: Maximize $f(\mathbf{x}) = 1 - x_1^2 - x_2^2$
subject to $g(\mathbf{x}) = x_1 + x_2 - 1 \geq 0$



Example: Maximize $f(\mathbf{x}) = 1 - x_1^2 - x_2^2$
subject to $g(\mathbf{x}) = x_1 + x_2 - 1 \geq 0$

- So the Lagrangian function is

$$L(\mathbf{x}, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1)$$

- The conditions for L to be stationary are

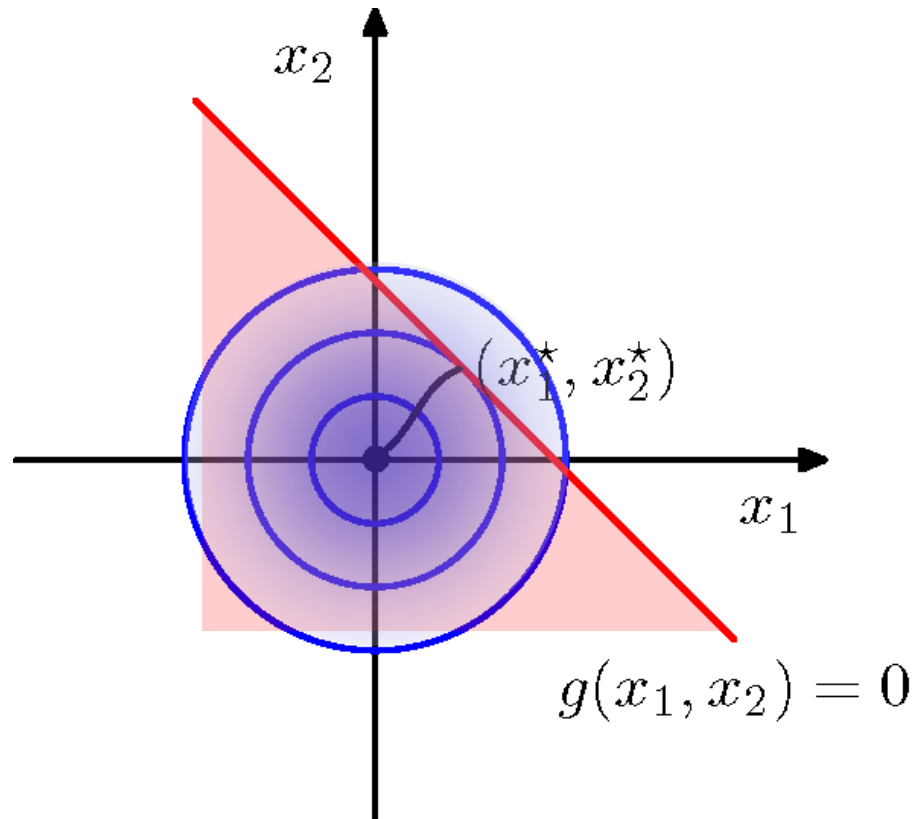
$$\partial L / \partial x_1 = -2x_1 + \lambda = 0$$

$$\partial L / \partial x_2 = -2x_2 + \lambda = 0$$

$$\partial L / \partial \lambda = x_1 + x_2 - 1 = 0$$

- Can solve to find $\lambda = 1, x_1 = x_2 = \frac{1}{2}$
 - Which still satisfies KKT conditions

Example: Maximize $f(\mathbf{x}) = 1 - x_1^2 - x_2^2$
subject to $g(\mathbf{x}) = -x_1 - x_2 + 1 \geq 0$



Example: Maximize $f(\mathbf{x}) = 1 - x_1^2 - x_2^2$
subject to $g(\mathbf{x}) = -x_1 - x_2 + 1 \geq 0$

- So the Lagrangian function is

$$L(\mathbf{x}, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(-x_1 - x_2 + 1)$$

- The conditions for L to be stationary are

$$\partial L / \partial x_1 = -2x_1 - \lambda = 0$$

$$\partial L / \partial x_2 = -2x_2 - \lambda = 0$$

$$\partial L / \partial \lambda = x_1 + x_2 - 1 = 0$$

- Can solve to find $\lambda = -1$
 - which does not satisfy KKT condition $\lambda \geq 0$
- Instead use unconstrained solution $x_1 = x_2 = 0$
 - which does satisfy KKT conditions

Multiple constraints

each get their own Lagrange multiplier

- Maximize $f(\mathbf{x})$ subject to $g_i(\mathbf{x}) = 0$ and $h_j(\mathbf{x}) \geq 0$
- Leads to the Lagrangian function

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_i \lambda_i g_i(\mathbf{x}) + \sum_j \mu_j h_j(\mathbf{x})$$

- Still solve for $\nabla L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = 0$
- Trickier in general to figure out which $h_j(\mathbf{x})$ constraints should be active

Minimizing $f(\mathbf{x})$ with an inequality constraint requires a slightly different Lagrangian

- Minimize WRT \mathbf{x}

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

- Still subject to

$$g(\mathbf{x}) \geq 0$$

Summary of Lagrange multipliers with multiple inequality constraints

- Goal: maximize $f(\mathbf{x})$ subject to $g_i(\mathbf{x}) \geq 0$
- Write down Lagrangian function

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \sum_i \lambda_i g_i(\mathbf{x})$$

- Find points where $\nabla L(\mathbf{x}, \lambda) = 0$
- Keep points that satisfy constraints $g_i(\mathbf{x}) \geq 0$
- Figure out which KKT conditions should be active
 - Don't need to try all 2^I combinations for SVMs
 - Because $f(\mathbf{x})$ and $g(\mathbf{x})$ form a “quadratic program”

Back to SVMs: Maximum margin solution is a fixed point of the Lagrangian function

- Recall, the maximum margin hyperplane is

$$\operatorname{argmin}_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to } d_p(\mathbf{w}^T \mathbf{x}_p + b) \geq 1$$

- Minimization of a quadratic function subject to multiple linear inequality constraints
- Will use Lagrange multipliers, a_p , to write Lagrangian function

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_p a_p (d_p(\mathbf{w}^T \mathbf{x}_p + b) - 1)$$

- Note that \mathbf{x}_p and d_p are fixed for the optimization