

COMP 551 - Applied Machine Learning

Lecture 10 – Support vector machines

William L. Hamilton

(with slides and content from Joelle Pineau)

* Unless otherwise noted, all material posted for this course are copyright of the instructor, and cannot be reused or reposted without the instructor's written permission.

MiniProject 1

- Average ~90%
- Common errors:
 - Incorrect gradient descent implementation (15 points).
 - Incorrect closed-form implementation (15 points).
 - Typos and clarity issues (2-5 points).

High-level views of classification

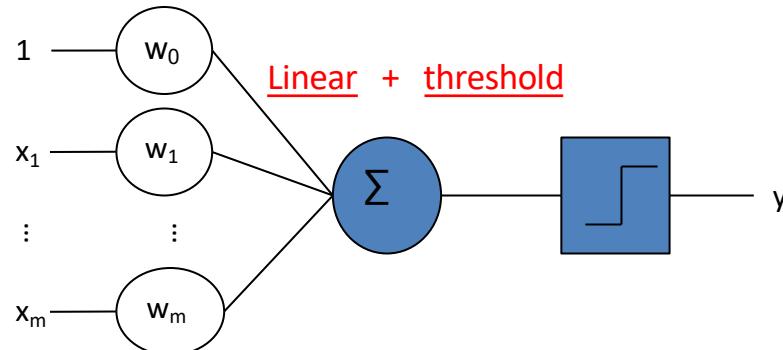
- Probabilistic
 - Goal: Estimate $P(y | x)$, i.e. the conditional probability of the target variable given the feature data.
- Decision boundaries
 - Goal: Partition the feature space into different regions, and classify points based on the region where they lie.

Outline

- Perceptrons
 - Definition
 - Perceptron learning rule
 - Convergence
- Margin & max margin classifiers
- Linear Support Vector Machines
 - Formulation as optimization problem
 - Generalized Lagrangian and dual
- Non-linear Support Vector Machines (next class)

A simple linear classifier

- Given a binary classification task: $\{\mathbf{x}_i, y_i\}_{i=1:n}$, $y_i = \{-1, +1\}$.
- The perceptron (Rosenblatt, 1957) is a classifier of the form:
$$h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x}) = \{+1 \text{ if } \mathbf{w}^T \mathbf{x} \geq 0; -1 \text{ otherwise}\}$$
 - The decision boundary is $\mathbf{w}^T \mathbf{x} = 0$.
 - An example $\langle \mathbf{x}_i, y_i \rangle$ is classified correctly if and only if: $y_i(\mathbf{w}^T \mathbf{x}_i) > 0$.



Perceptron learning rule (Rosenblatt, 1957)

- Consider the following procedure:

Initialize w_j , $j=0:m$ randomly,

While any training examples remain incorrectly classified

 Loop through all misclassified examples x_i

 Perform the update: $w \leftarrow w + \alpha y_i x_i$

 where α is the learning rate (or step size).

- Intuition: For misclassified positive examples, increase $w^T x$, and reduce it for negative examples.

Gradient-descent learning

- The perceptron learning rule can be interpreted as a **gradient descent procedure**, with optimization criterion:

$$\text{Err}(\mathbf{w}) = \sum_{i=1:n} \{ 0 \text{ if } y_i \mathbf{w}^T \mathbf{x}_i \geq 0; -y_i \mathbf{w}^T \mathbf{x} \text{ otherwise } \}$$

Gradient-descent learning

- The perceptron learning rule can be interpreted as a **gradient descent procedure**, with optimization criterion:

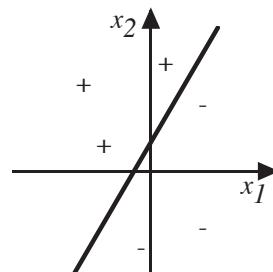
$$\text{Err}(\mathbf{w}) = \sum_{i=1:n} \{ 0 \text{ if } y_i \mathbf{w}^T \mathbf{x}_i \geq 0; -y_i \mathbf{w}^T \mathbf{x} \text{ otherwise } \}$$

- For **correctly classified examples**, the error is zero.
- For **incorrectly classified examples**, the error tells by how much $\mathbf{w}^T \mathbf{x}$ is on the wrong side of the decision boundary.
- The error is zero when all examples are classified correctly.

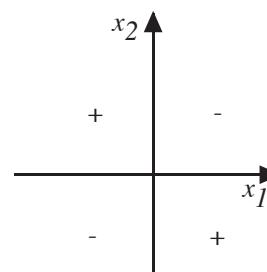
Linear separability

The data is **linearly separable** if and only if there exists a \mathbf{w} such that:

- For all examples, $y_i \mathbf{w}^T \mathbf{x}_i > 0$
- Or equivalently, the loss is zero for some set of parameters (\mathbf{w}).



Linearly separable



Not linearly separable

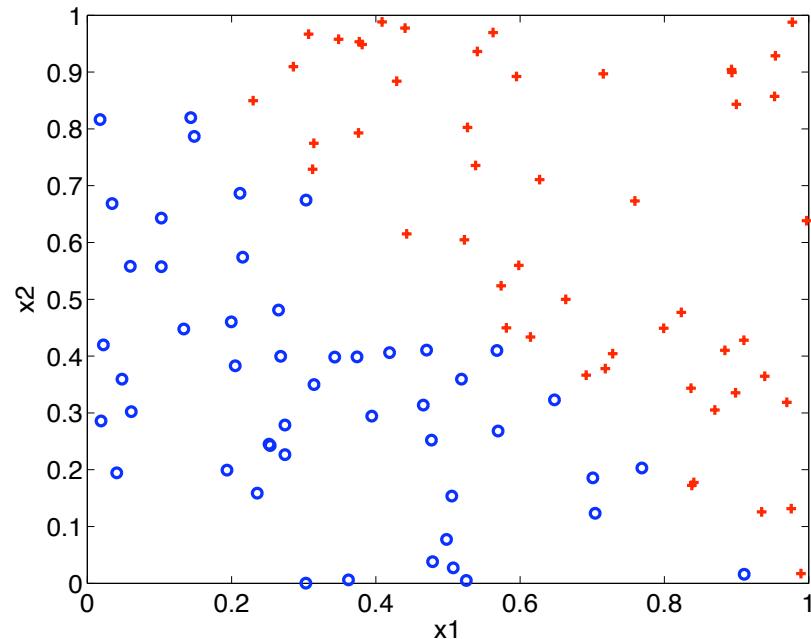
Perceptron convergence theorem

- The basic theorem:
 - If the perceptron learning rule is applied to a **linearly separable dataset**, a solution will be found after some finite number of updates.

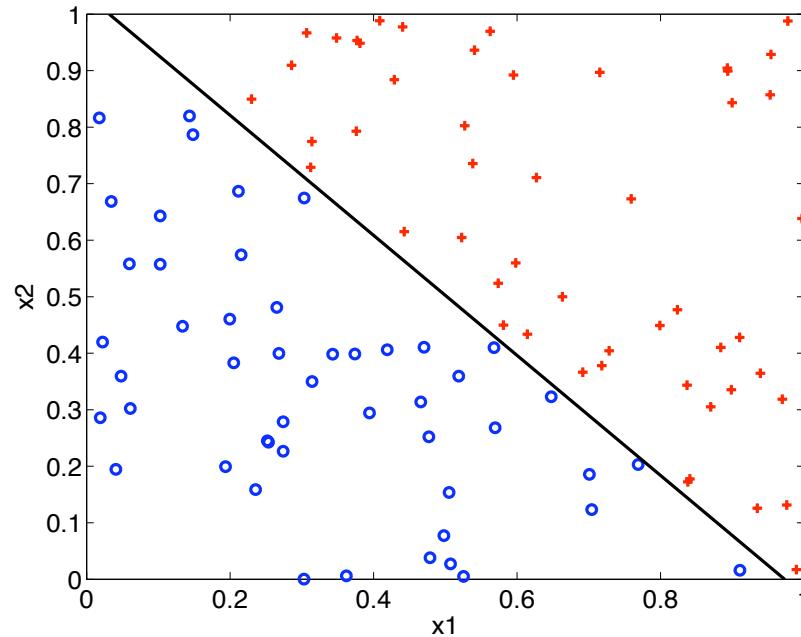
Perceptron convergence theorem

- The basic theorem:
 - If the perceptron learning rule is applied to a **linearly separable dataset**, a solution will be found after some **finite number of updates**.
- Additional comments:
 - The number of updates depends on the dataset, on the learning rate, and on the initial weights.
 - If the data is not linearly separable, there will be oscillation (which can be detected automatically).
 - Decreasing the learning rate to 0 can cause the oscillation to settle on some particular solution.

Perceptron learning example



Perceptron learning example



Weight as a combination of input vectors

- Recall perceptron learning rule:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha y_i \mathbf{x}_i$$

- If initial weights are zero, then at any step, the weights are a linear combination of feature vectors of the examples:

$$\mathbf{w} = \sum_{i=1:n} \alpha_i y_i \mathbf{x}_i$$

where α_i is the sum of step sizes used for all updates applied example i .

Weight as a combination of input vectors

- Recall perceptron learning rule:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha y_i \mathbf{x}_i$$

- If initial weights are zero, then at any step, the weights are a linear combination of feature vectors of the examples:

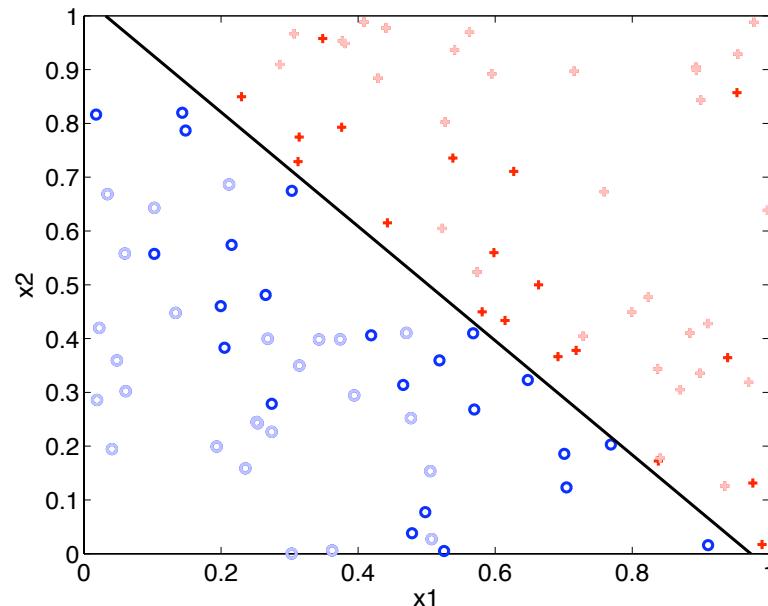
$$\mathbf{w} = \sum_{i=1:n} \alpha_i y_i \mathbf{x}_i$$

where α_i is the sum of step sizes used for all updates applied example i .

- By the end of training, some examples may have never participated in an update, so will have $\alpha_i=0$.
- This is called the **dual representation** of the classifier.

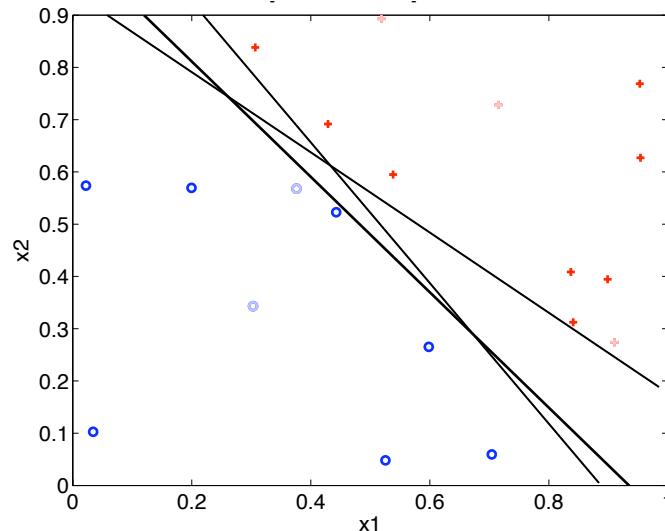
Perceptron learning example

- Examples used (bold) and not (faint). **What do you notice?**



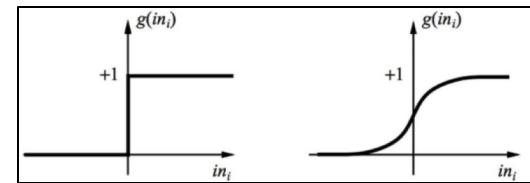
Perceptron learning example

- Solutions are often non-unique. The solution depends on the set of instances and the order of sampling in updates.



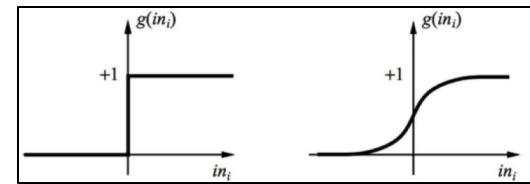
A few comments on the Perceptron

- Perceptrons can be learned to fit linearly separable data, using a gradient-descent rule.
 - The logistic function offers a “smooth” version of the perceptron.



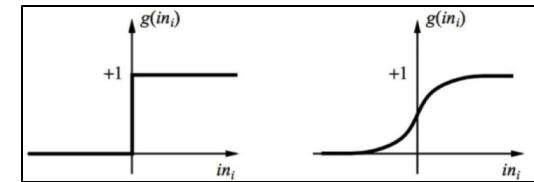
A few comments on the Perceptron

- Perceptrons can be learned to fit linearly separable data, using a gradient-descent rule.
 - The logistic function offers a “smooth” version of the perceptron.
- **Two issues:**
 - Solutions are non-unique.



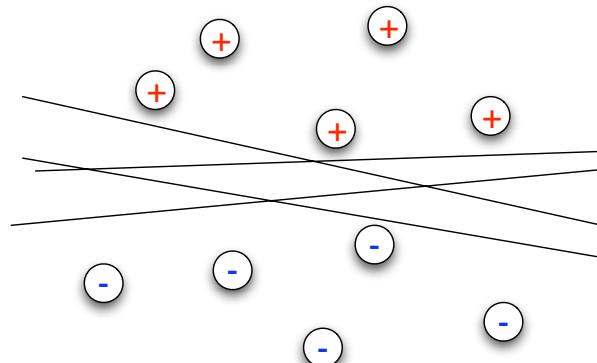
A few comments on the Perceptron

- Perceptrons can be learned to fit linearly separable data, using a gradient-descent rule.
 - The logistic function offers a “smooth” version of the perceptron.
- **Two issues:**
 - Solutions are non-unique.
 - What about non-linearly separable data? (*Topic for next class.*)
 - Perhaps data can be linearly separated in a different feature space?
 - Perhaps we can relax the criterion of separating all the data?



The non-uniqueness issue

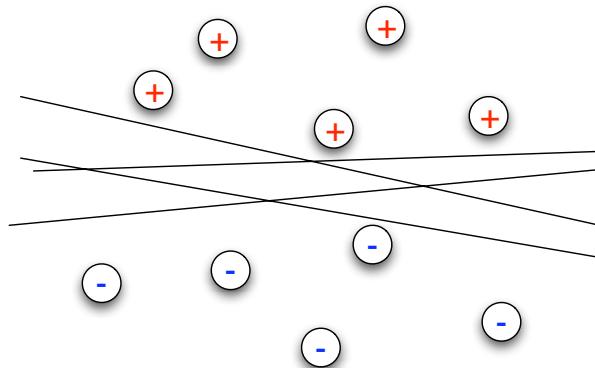
- Consider a linearly separable binary classification dataset.
- There is an infinite number of hyper-planes that separate the classes:
- Which plane is best?



The non-uniqueness issue

- Consider a linearly separable binary classification dataset.
- There is an infinite number of hyper-planes that separate the classes:

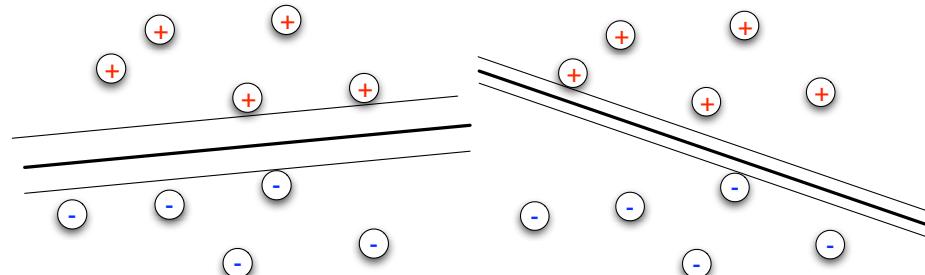
- Which plane is best?



- Related question: For a given plane, for which points should we be most confident in the classification?

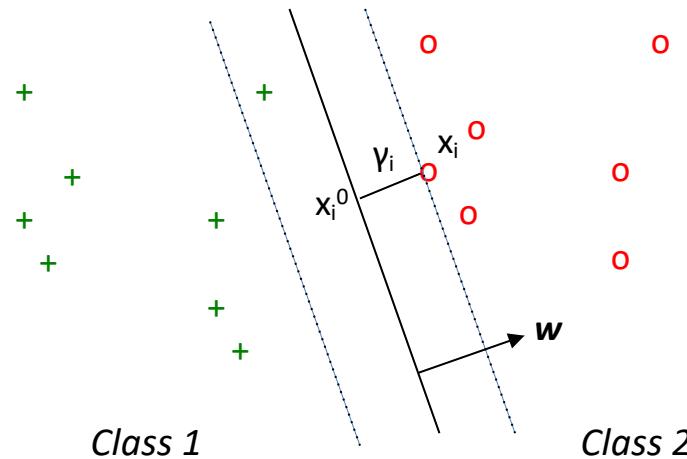
Linear Support Vector Machine (SVM)

- A linear SVM is a perceptron for which we chose w such that the margin is maximized.
- For a given separating hyper-plane, the margin is twice the (Euclidean) distance from hyper-plane to nearest training example.
 - I.e. the width of the “strip” around the decision boundary that contains no training examples.



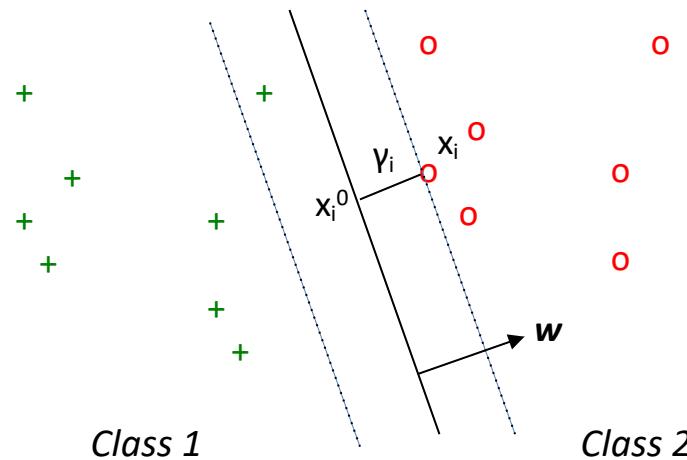
Distance to the decision boundary

- Suppose we have a decision boundary that separates the data.



Distance to the decision boundary

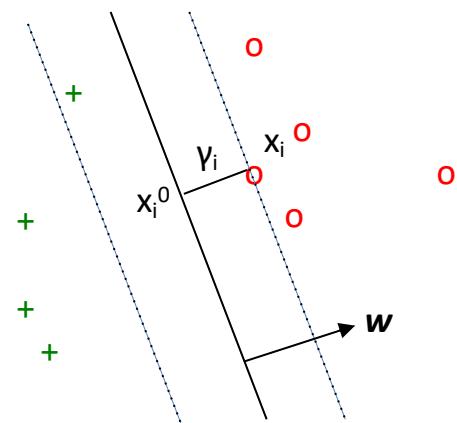
- Suppose we have a decision boundary that separates the data.



- Let y_i be the distance from instance x_i to the decision boundary.
- Define vector w to be the normal to the decision boundary.

Distance to the decision boundary

- How can we write γ_i in terms of \mathbf{x}_i , y_i , \mathbf{w} ?
- Let \mathbf{x}_i^0 be the point on the decision boundary nearest \mathbf{x}_i
- The vector from \mathbf{x}_i^0 to \mathbf{x}_i is $\gamma_i \mathbf{w} / \|\mathbf{w}\|$.
 - γ_i is a scalar (distance from \mathbf{x}_i to \mathbf{x}_i^0)
 - $\mathbf{w}/\|\mathbf{w}\|$ is the unit normal.
- So we can define $\mathbf{x}_i^0 = \mathbf{x}_i - \gamma_i \mathbf{w} / \|\mathbf{w}\|$.

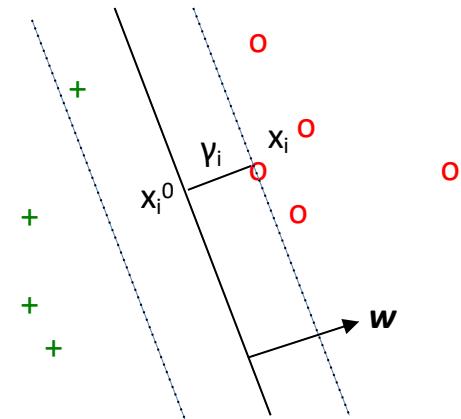


Distance to the decision boundary

- How can we write y_i in terms of x_i, y_i, w ?
- Let x_i^0 be the point on the decision boundary nearest x_i
- The vector from x_i^0 to x_i is $y_i w / ||w||$.
 - y_i is a scalar (distance from x_i to x_i^0)
 - $w / ||w||$ is the unit normal.
- So we can define $x_i^0 = x_i - y_i w / ||w||$.
- As x_i^0 is on the decision boundary, we have

$$w^T(x_i - y_i w / ||w||) = 0$$

- Solving for y_i yields, for a positive example: $y_i = w^T x_i / ||w||$
or for examples of both classes: $y_i = y_i w^T x_i / ||w||$



Optimization

- First suggestion:

Maximize M

with respect to w

subject to $y_i w^T x_i / ||w|| \geq M, \forall i$

Optimization

- First suggestion:

Maximize M

with respect to w

subject to $y_i w^T x_i / \|w\| \geq M, \forall i$

- This is not very convenient for optimization:
 - w appears nonlinearly in the constraints.
 - Problem is underconstrained. If (w, M) is optimal, so is $(\beta w, M)$, for any $\beta > 0$.
- Add a constraint:** $\|w\| = 1$

Optimization

- First suggestion:

Maximize M

with respect to w

subject to $y_i w^T x_i / \|w\| \geq M, \forall i$

- This is not very convenient for optimization:

- w appears nonlinearly in the constraints.

- Problem is underconstrained. If (w, M) is optimal, so is $(\beta w, M)$, for any $\beta > 0$.

Add a constraint: $\|w\| = 1$

- Instead try:

Minimize $\|w\|$

with respect to w

subject to $y_i w^T x_i \geq 1$

Final formulation

- Let's minimize $\frac{1}{2}||\mathbf{w}||^2$ instead of $||\mathbf{w}||$

(Taking the square is a monotone transform, as $||\mathbf{w}||$ is positive, so it doesn't change the optimal solution. The $\frac{1}{2}$ is for mathematical convenience.)

- This gets us to: Min $\frac{1}{2} ||\mathbf{w}||^2$

w.r.t. \mathbf{w}

s.t. $y_i \mathbf{w}^T \mathbf{x}_i \geq 1$

Final formulation

- Let's minimize $\frac{1}{2} \|\mathbf{w}\|^2$ instead of $\|\mathbf{w}\|$

(Taking the square is a monotone transform, as $\|\mathbf{w}\|$ is positive, so it doesn't change the optimal solution. The $\frac{1}{2}$ is for mathematical convenience.)

- This gets us to: Min $\frac{1}{2} \|\mathbf{w}\|^2$

w.r.t. \mathbf{w}

s.t. $y_i \mathbf{w}^T \mathbf{x}_i \geq 1$

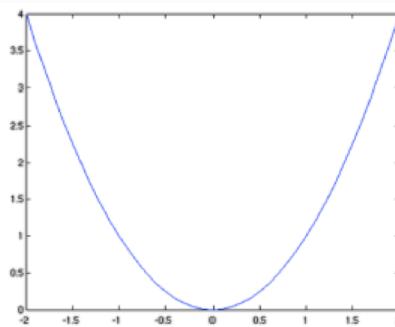
- This can be solved! How?

- It is a **quadratic programming** (QP) problem – a standard type of optimization problem for which many efficient packages are available. Better yet, it's a convex (positive semidefinite) QP.

Constrained optimization

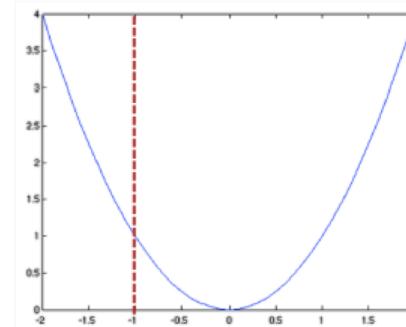
$$\begin{aligned} & \min_x x^2 \\ \text{s.t. } & x \geq b \end{aligned}$$

$$\min_x x^2$$



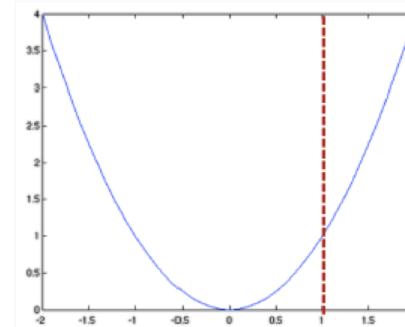
$$x^* = 0$$

$$\begin{aligned} & \min_x x^2 \\ \text{s.t. } & x \geq -1 \end{aligned}$$



$$x^* = 0$$

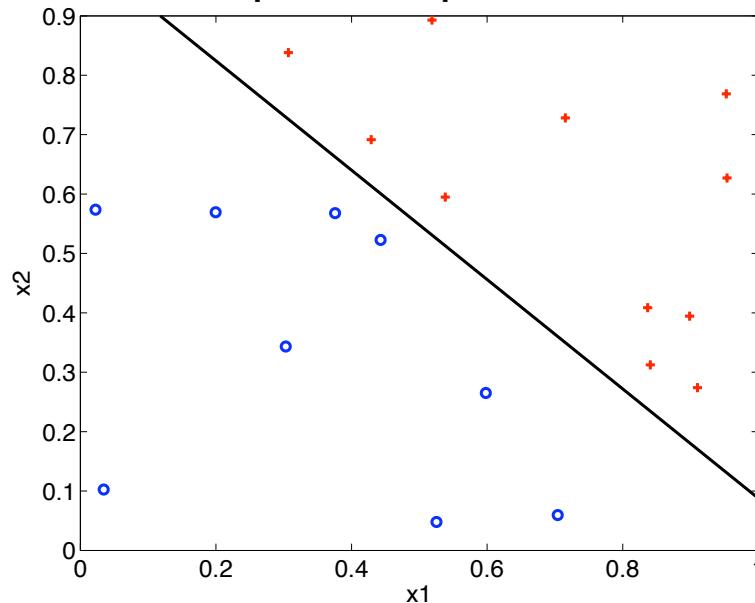
$$\begin{aligned} & \min_x x^2 \\ \text{s.t. } & x \geq 1 \end{aligned}$$



$$x^* = 1$$

Picture from: http://www.cs.cmu.edu/~aarti/Class/10701_Spring14/

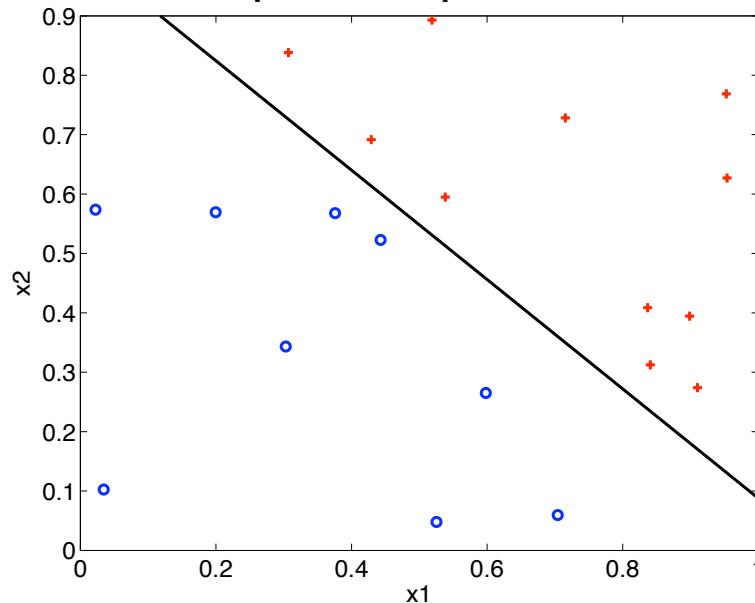
Example



We have a unique solution, but no support vectors yet.

Recall the dual solution for the Perceptron: **Extend for the margin case.**

Example



We have a unique solution, but no support vectors yet.
Recall the dual solution for the Perceptron: Extend for the margin case.

Lagrange multipliers

- Consider the following optimization problem, called **primal**:

$$\begin{array}{ll} \min_w & f(w) \\ \text{s.t.} & g_i(w) \leq 0, i=1 \dots k \end{array}$$

- We define the **generalized Lagrangian**:

$$L(w, \alpha) = f(w) + \sum_{i=1:k} \alpha_i g_i(w)$$

where $\alpha_i, i=1 \dots k$ are the *Lagrange multipliers*.

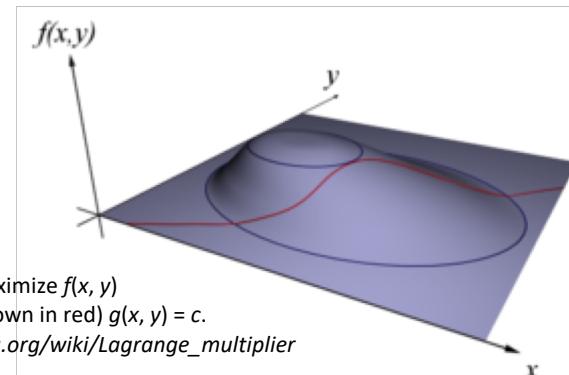


Figure : Find x and y to maximize $f(x, y)$
subject to a constraint (shown in red) $g(x, y) = c$.
From: https://en.wikipedia.org/wiki/Lagrange_multiplier

Lagrangian optimization

- Consider $P(\mathbf{w}) = \max_{\alpha: \alpha_i \geq 0} L(\mathbf{w}, \alpha) = \max_{\alpha: \alpha_i \geq 0} f(\mathbf{w}) + \sum_{i=1:k} \alpha_i g_i(\mathbf{w})$
 - (P stands for “primal”)
- Recall: $L(\mathbf{w}, \alpha) =$ Observe that the following is true:

$$P(\mathbf{w}) = \begin{cases} f(\mathbf{w}), & \text{if all constraints are satisfied,} \\ +\infty, & \text{otherwise } \end{cases}$$

Lagrangian optimization

- Consider $P(\mathbf{w}) = \max_{\alpha: \alpha_i \geq 0} L(\mathbf{w}, \alpha) = \max_{\alpha: \alpha_i \geq 0} f(\mathbf{w}) + \sum_{i=1:k} \alpha_i g_i(\mathbf{w})$

(P stands for “primal”)

- Recall: $L(\mathbf{w}, \alpha) =$ Observe that the following is true:

$$P(\mathbf{w}) = \begin{cases} f(\mathbf{w}), & \text{if all constraints are satisfied,} \\ +\infty, & \text{otherwise } \end{cases}$$

- Hence, instead of computing $\min_{\mathbf{w}} f(\mathbf{w})$ subject to the original constraints, we can compute:

$$p^* = \min_{\mathbf{w}} P(\mathbf{w}) = \min_{\mathbf{w}} \max_{\alpha: \alpha_i \geq 0} L(\mathbf{w}, \alpha) \quad \textit{Primal}$$

- Alternately, invert max and min to get:

$$d^* = \max_{\alpha: \alpha_i \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \alpha) \quad \textit{Dual}$$

Maximum Margin Perceptron

- We wanted to solve:
$$\begin{array}{ll}\text{Min} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{w.r.t.} & \mathbf{w} \\ \text{s.t.} & y_i \mathbf{w}^T \mathbf{x}_i \geq 1\end{array}$$
- The Lagrangian is:
$$L(\mathbf{w}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_i \alpha_i (1 - y_i (\mathbf{w}^T \mathbf{x}_i))$$
- The **primal problem** is: $\min_{\mathbf{w}} \max_{\alpha: \alpha_i \geq 0} L(\mathbf{w}, \alpha)$
- The **dual problem** is: $\max_{\alpha: \alpha_i \geq 0} \min_{\mathbf{w}} L(\mathbf{w}, \alpha)$

Dual optimization problem

- Consider both solutions:

$$\begin{aligned} p^* &= \min_w \max_{\alpha: \alpha_i \geq 0} L(w, \alpha) \\ d^* &= \max_{\alpha: \alpha_i \geq 0} \min_w L(w, \alpha) \end{aligned}$$

Primal
Dual

- If f and g_i are convex and the g_i can all be satisfied simultaneously for some w , then we have equality: $d^* = p^* = L(w^*, \alpha^*)$.
 - w^* is the optimal weight vector (= primal solution)
 - α^* is the optimal set of support vectors (= dual solution)
- For SVMs, we have a quadratic objective and linear constraints so both f and g_i are convex.
- For linearly separable data, all g_i can be satisfied simultaneously.
- Note: w^*, α^* solve the primal and dual if and only if they satisfy the Karush-Kuhn-Tucker conditions (see *suggested readings*).

Solving the dual

- Taking derivatives of $L(\mathbf{w}, \alpha)$ wrt \mathbf{w} , setting to 0, and solving for \mathbf{w} :

$$L(\mathbf{w}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_i \alpha_i (1 - y_i (\mathbf{w}^T \mathbf{x}_i))$$

$$\delta L / \delta \mathbf{w} = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0$$

$$\mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i$$

- Just like for the perceptron with zero initial weights, the optimal solution \mathbf{w}^* is a linear combination of the \mathbf{x}_i .
- Plugging this back into L we get the dual: $\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j)$ with constraints $\alpha_i \geq 0$ and $\sum_i \alpha_i y_i = 0$. (Quadratic programming problem).
- Complexity of solving quadratic program? Polynomial time, $O(|v|^3)$ (where $|v| = \#$ variables in optimization; here $|v| = n$). Fast approximations exist.

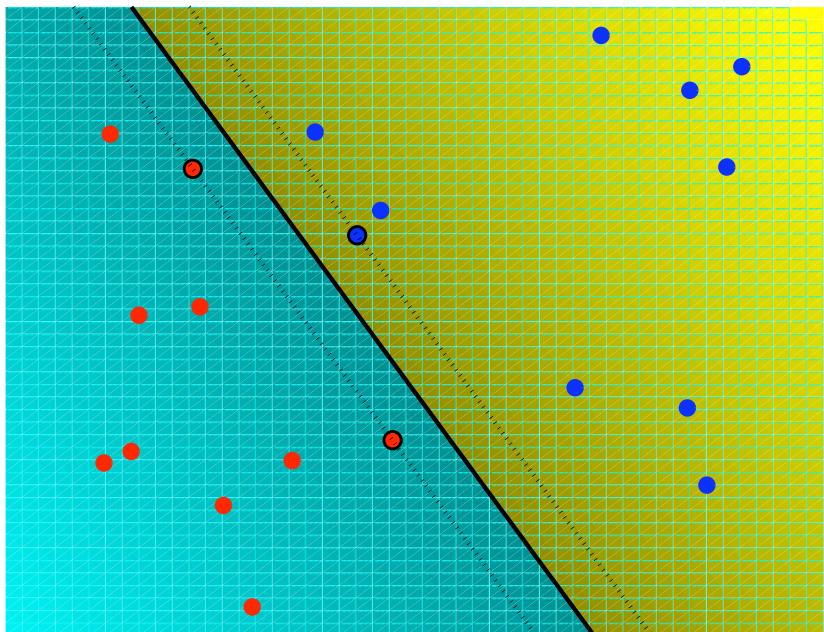
The support vectors

- Suppose we find the optimal α 's (e.g. using a QP package.)
- Constraint i is active when $\alpha_i > 0$. This corresponds for the points for which $(1-y_i\mathbf{w}^T\mathbf{x}_i)=0$.
- These are the points lying on the edge of the margin. We call them **support vectors**. They define the decision boundary.
- The output of the classifier for query point \mathbf{x} is computed as:

$$h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\sum_{i=1:n} \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}))$$

- I.e., it is determined by computing the dot product of the query point with the support vectors.

Example



Support vectors are in bold

What you should know

From today:

- The perceptron algorithm.
- The margin definition for linear SVMs.
- The use of Lagrange multipliers to transform optimization problems.
- The primal and dual optimization problems for SVMs.

After the next class:

- Non-linearly separable case.
- Feature space version of SVMs.
- The kernel trick and examples of common kernels.