

Max Margin-Classifier

Oliver Schulte - CMPT 726

Bishop PRML Ch. 7

Outline

Maximum Margin Criterion

Math

Maximizing the Margin

Non-Separable Data

Kernels and Non-linear Mappings

- Where does the maximization problem come from?
- The intuition comes from the **primal version**, which is based on a feature mapping ϕ .
- Theorem: Every valid kernel $k(x, y)$ is the dot-product $\phi(x)[\phi(y)]^T$ for some set of basis functions (feature mapping) ϕ .
- The **feature space** $\phi(x)$ could be high-dimensional, even infinite.
- This is good because if data aren't separable in original input space (x) , they may be in feature space $\phi(x)$
- We can think about how to find a good linear separator using the dot product in high dimensions, then transfer this back to kernels in the original input space.

Kernels and Non-linear Mappings

- Where does the maximization problem come from?
- The intuition comes from the **primal version**, which is based on a feature mapping ϕ .
- Theorem: Every valid kernel $k(\mathbf{x}, \mathbf{y})$ is the dot-product $\phi(\mathbf{x})[\phi(\mathbf{y})]^T$ for some set of basis functions (feature mapping) ϕ .
- The **feature space** $\phi(\mathbf{x})$ could be high-dimensional, even infinite.
- This is good because if data aren't separable in original input space (\mathbf{x}), they may be in feature space $\phi(\mathbf{x})$
- We can think about how to find a good linear separator using the dot product in high dimensions, then transfer this back to kernels in the original input space.

Kernels and Non-linear Mappings

- Where does the maximization problem come from?
- The intuition comes from the **primal version**, which is based on a feature mapping ϕ .
- Theorem: Every valid kernel $k(\mathbf{x}, \mathbf{y})$ is the dot-product $\phi(\mathbf{x})[\phi(\mathbf{y})]^T$ for some set of basis functions (feature mapping) ϕ .
- The **feature space** $\phi(\mathbf{x})$ could be high-dimensional, even infinite.
- This is good because if data aren't separable in original input space (\mathbf{x}) , they may be in feature space $\phi(\mathbf{x})$
- We can think about how to find a good linear separator using the dot product in high dimensions, then transfer this back to kernels in the original input space.

Kernels and Non-linear Mappings

- Where does the maximization problem come from?
- The intuition comes from the **primal version**, which is based on a feature mapping ϕ .
- Theorem: Every valid kernel $k(\mathbf{x}, \mathbf{y})$ is the dot-product $\phi(\mathbf{x})[\phi(\mathbf{y})]^T$ for some set of basis functions (feature mapping) ϕ .
- The **feature space** $\phi(\mathbf{x})$ could be high-dimensional, even infinite.
- This is good because if data aren't separable in original input space (\mathbf{x}), they may be in feature space $\phi(\mathbf{x})$
- We can think about how to find a good linear separator using the dot product in high dimensions, then transfer this back to kernels in the original input space.

Why Kernels?

- If we can use dot products with features, why bother with kernels?
 - Often easier to specify how similar two things are (dot product) than to construct explicit feature space ϕ .
 - e.g. graphs, sets, strings (NIPS 2009 best student paper award).
 - There are high-dimensional (even infinite) spaces that have efficient-to-compute kernels

Kernel Trick

- In previous lectures on linear models, we would explicitly compute $\phi(\mathbf{x}_i)$ for each datapoint
 - Run algorithm in feature space
- For some feature spaces, can compute dot product $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ efficiently
- Efficient method is computation of a kernel function $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$
- The **kernel trick** is to rewrite an algorithm to only have x enter in the form of dot products
- The menu:
 - Kernel trick examples
 - Kernel functions

Kernel Trick

- In previous lectures on linear models, we would explicitly compute $\phi(\mathbf{x}_i)$ for each datapoint
 - Run algorithm in feature space
- For some feature spaces, can compute dot product $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ efficiently
- Efficient method is computation of a kernel function
$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$
- The **kernel trick** is to rewrite an algorithm to only have x enter in the form of dot products
- The menu:
 - Kernel trick examples
 - Kernel functions

Kernel Trick

- In previous lectures on linear models, we would explicitly compute $\phi(\mathbf{x}_i)$ for each datapoint
 - Run algorithm in feature space
- For some feature spaces, can compute dot product $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ efficiently
- Efficient method is computation of a kernel function
$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$
- The **kernel trick** is to rewrite an algorithm to only have \mathbf{x} enter in the form of dot products
- The menu:
 - Kernel trick examples
 - Kernel functions

Kernel Trick

- In previous lectures on linear models, we would explicitly compute $\phi(\mathbf{x}_i)$ for each datapoint
 - Run algorithm in feature space
- For some feature spaces, can compute dot product $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ efficiently
- Efficient method is computation of a kernel function
$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$
- The **kernel trick** is to rewrite an algorithm to only have \mathbf{x} enter in the form of dot products
- The menu:
 - Kernel trick examples
 - Kernel functions

A Kernel Trick

- Let's look at the nearest-neighbour classification algorithm
- For input point x_i , find point x_j with smallest distance:

$$\begin{aligned}\|x_i - x_j\|^2 &= (x_i - x_j)^T (x_i - x_j) \\ &= x_i^T x_i - 2x_i^T x_j + x_j^T x_j\end{aligned}$$

- If we used a non-linear feature space $\phi(\cdot)$:

$$\begin{aligned}\|\phi(x_i) - \phi(x_j)\|^2 &= \phi(x_i)^T \phi(x_i) - 2\phi(x_i)^T \phi(x_j) + \phi(x_j)^T \phi(x_j) \\ &= k(x_i, x_i) - 2k(x_i, x_j) + k(x_j, x_j)\end{aligned}$$

- So nearest-neighbour can be done in a high-dimensional feature space without actually moving to it

A Kernel Trick

- Let's look at the nearest-neighbour classification algorithm
- For input point x_i , find point x_j with smallest distance:

$$\begin{aligned}\|x_i - x_j\|^2 &= (x_i - x_j)^T (x_i - x_j) \\ &= x_i^T x_i - 2x_i^T x_j + x_j^T x_j\end{aligned}$$

- If we used a non-linear feature space $\phi(\cdot)$:

$$\begin{aligned}\|\phi(x_i) - \phi(x_j)\|^2 &= \phi(x_i)^T \phi(x_i) - 2\phi(x_i)^T \phi(x_j) + \phi(x_j)^T \phi(x_j) \\ &= k(x_i, x_i) - 2k(x_i, x_j) + k(x_j, x_j)\end{aligned}$$

- So nearest-neighbour can be done in a high-dimensional feature space without actually moving to it

A Kernel Trick

- Let's look at the nearest-neighbour classification algorithm
- For input point x_i , find point x_j with smallest distance:

$$\begin{aligned}\|x_i - x_j\|^2 &= (x_i - x_j)^T (x_i - x_j) \\ &= x_i^T x_i - 2x_i^T x_j + x_j^T x_j\end{aligned}$$

- If we used a non-linear feature space $\phi(\cdot)$:

$$\begin{aligned}\|\phi(x_i) - \phi(x_j)\|^2 &= \phi(x_i)^T \phi(x_i) - 2\phi(x_i)^T \phi(x_j) + \phi(x_j)^T \phi(x_j) \\ &= k(x_i, x_i) - 2k(x_i, x_j) + k(x_j, x_j)\end{aligned}$$

- So nearest-neighbour can be done in a high-dimensional feature space without actually moving to it

Example: The Quadratic Kernel Function

- Consider again the kernel function $k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$
- With $\mathbf{x}, \mathbf{z} \in \mathbb{R}^2$,

$$\begin{aligned}k(\mathbf{x}, \mathbf{z}) &= (1 + x_1 z_1 + x_2 z_2)^2 \\&= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\&= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2)(1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\&= \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z})\end{aligned}$$

- So this particular kernel function does correspond to a dot product in a feature space (is valid)
- Computing $k(\mathbf{x}, \mathbf{z})$ is faster than explicitly computing $\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z})$
 - In higher dimensions, larger exponent, much faster

Example: The Quadratic Kernel Function

- Consider again the kernel function $k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$
- With $\mathbf{x}, \mathbf{z} \in \mathbb{R}^2$,

$$\begin{aligned}k(\mathbf{x}, \mathbf{z}) &= (1 + x_1 z_1 + x_2 z_2)^2 \\&= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\&= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2)(1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\&= \phi(\mathbf{x})^T \phi(\mathbf{z})\end{aligned}$$

- So this particular kernel function does correspond to a dot product in a feature space (is valid)
- Computing $k(\mathbf{x}, \mathbf{z})$ is faster than explicitly computing $\phi(\mathbf{x})^T \phi(\mathbf{z})$
 - In higher dimensions, larger exponent, much faster

Example: The Quadratic Kernel Function

- Consider again the kernel function $k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$
- With $\mathbf{x}, \mathbf{z} \in \mathbb{R}^2$,

$$\begin{aligned}k(\mathbf{x}, \mathbf{z}) &= (1 + x_1 z_1 + x_2 z_2)^2 \\&= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\&= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2)(1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \sqrt{2}z_1 z_2, z_2^2)^T \\&= \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z})\end{aligned}$$

- So this particular kernel function does correspond to a dot product in a feature space (is valid)
- Computing $k(\mathbf{x}, \mathbf{z})$ is faster than explicitly computing $\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z})$
 - In higher dimensions, larger exponent, much faster

Regression Kernelized

- Many classifiers can be written as using only dot products.
- Kernelization = replace dot products by kernel.
- E.g., the kernel solution for regularized least squares regression is

$$y(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} \quad \text{vs.} \quad \phi(\mathbf{x}) (\Phi^T \Phi + \lambda \mathbf{I}_M)^{-1} \Phi^T \mathbf{t}$$

for original version

- N is number of datapoints (size of Gram matrix \mathbf{K})
- M is number of basis functions (size of matrix $\Phi^T \Phi$)
- Bad if $N > M$, but good otherwise
- $\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n))$ is the vector of kernel values over data points x_n .

Outline

Maximum Margin Criterion

Math

Maximizing the Margin

Non-Separable Data

Linear Classification

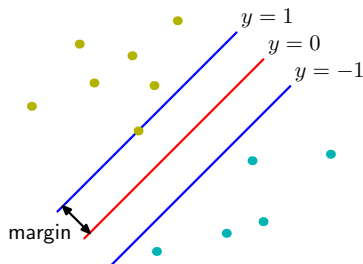
- Consider a two class classification problem
- Use a linear model

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

followed by a threshold function

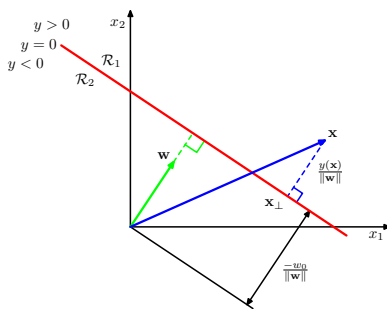
- For now, let's assume training data are linearly separable
 - Recall that the perceptron would converge to a perfect classifier for such data
 - But there are many such perfect classifiers

Max Margin



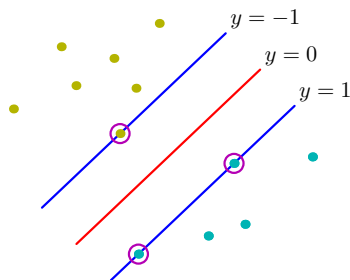
- We can define the **margin** of a classifier as the minimum distance to any example
- In **support vector machines** the decision boundary which maximizes the margin is chosen.
- Intuitively, this is the line “right in the middle” between the two classes.

Marginal Geometry



- Recall from Ch. 4 $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$
 - (Using Ch.4 notation \mathbf{x} for input rather than $\phi(\mathbf{x})$.)
- $\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} - \frac{-b}{\|\mathbf{w}\|} = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$ is signed distance to decision boundary.

Support Vectors



- Assuming data are separated by the hyperplane, distance to decision boundary is $\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|}$
- The maximum margin criterion chooses \mathbf{w}, b by:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

- Points with this min value are known as **support vectors**

Canonical Representation

- This optimization problem is complex:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

- Note that rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ does not change distance $\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|}$ (many equiv. answers)
- So for \mathbf{x}_* closest to surface, can use $\mathbf{w} \rightarrow \mathbf{w}/y(\mathbf{x}_*)$ and $b \rightarrow b/y(\mathbf{x}_*)$ so that:

$$t_*(\mathbf{w}^T \phi(\mathbf{x}_*) + b) = 1$$

- All other points are at least this far away:

$$\forall n, t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$

- Under these constraints, the optimization becomes:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

Canonical Representation

- This optimization problem is complex:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

- Note that rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ does not change distance $\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|}$ (many equiv. answers)
- So for \mathbf{x}_* closest to surface, can use $\mathbf{w} \rightarrow \mathbf{w}/y(\mathbf{x}_*)$ and $b \rightarrow b/y(\mathbf{x}_*)$ so that:

$$t_*(\mathbf{w}^T \phi(\mathbf{x}_*) + b) = 1$$

- All other points are at least this far away:

$$\forall n, t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$

- Under these constraints, the optimization becomes:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

Canonical Representation

- This optimization problem is complex:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

- Note that rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ does not change distance $\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|}$ (many equiv. answers)
- So for \mathbf{x}_* closest to surface, can use $\mathbf{w} \rightarrow \mathbf{w}/y(\mathbf{x}_*)$ and $b \rightarrow b/y(\mathbf{x}_*)$ so that:

$$t_*(\mathbf{w}^T \phi(\mathbf{x}_*) + b) = 1$$

- All other points are at least this far away:

$$\forall n, t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$

- Under these constraints, the optimization becomes:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

Canonical Representation

- This optimization problem is complex:

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

- Note that rescaling $\mathbf{w} \rightarrow \kappa \mathbf{w}$ and $b \rightarrow \kappa b$ does not change distance $\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|}$ (many equiv. answers)
- So for \mathbf{x}_* closest to surface, can use $\mathbf{w} \rightarrow \mathbf{w}/y(\mathbf{x}_*)$ and $b \rightarrow b/y(\mathbf{x}_*)$ so that:

$$t_*(\mathbf{w}^T \phi(\mathbf{x}_*) + b) = 1$$

- All other points are at least this far away:

$$\forall n, t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$

- Under these constraints, the optimization becomes:

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

Canonical Representation

- So the optimization problem is now a constrained optimization problem:

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$
$$s.t. \quad \forall n, t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$

- To solve this, we need to take a detour into **Lagrange multipliers**

Outline

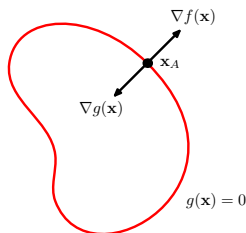
Maximum Margin Criterion

Math

Maximizing the Margin

Non-Separable Data

Lagrange Multipliers



Consider the problem:

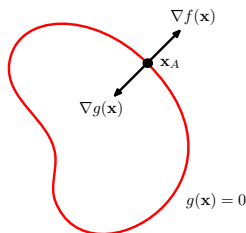
$$\begin{aligned} & \max_{\mathbf{x}} f(\mathbf{x}) \\ \text{s.t.} \quad & g(\mathbf{x}) = 0 \end{aligned}$$

- Points on $g(\mathbf{x}) = 0$ must have $\nabla g(\mathbf{x})$ normal to surface
- A **stationary point** must have no change in f in the direction of the constraint surface, so $\nabla f(\mathbf{x})$ must also be normal to the surface.
 - So there must be some $\lambda \neq 0$ such that $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$
- Define **Lagrangian**:

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

- Stationary points of $L(\mathbf{x}, \lambda)$ have $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = \nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ and $\nabla_{\lambda} L(\mathbf{x}, \lambda) = g(\mathbf{x}) = 0$

Lagrange Multipliers



Consider the problem:

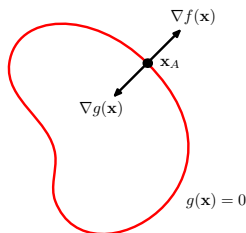
$$\begin{aligned} & \max_{\mathbf{x}} f(\mathbf{x}) \\ \text{s.t.} \quad & g(\mathbf{x}) = 0 \end{aligned}$$

- Points on $g(\mathbf{x}) = 0$ must have $\nabla g(\mathbf{x})$ normal to surface
- A **stationary point** must have no change in f in the direction of the constraint surface, so $\nabla f(\mathbf{x})$ must also be normal to the surface.
 - So there must be some $\lambda \neq 0$ such that $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$
- Define **Lagrangian**:

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

- Stationary points of $L(\mathbf{x}, \lambda)$ have $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = \nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ and $\nabla_{\lambda} L(\mathbf{x}, \lambda) = g(\mathbf{x}) = 0$

Lagrange Multipliers



Consider the problem:

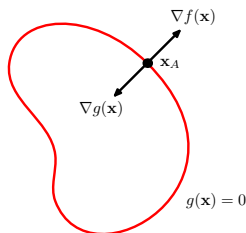
$$\begin{aligned} & \max_{\mathbf{x}} f(\mathbf{x}) \\ \text{s.t.} \quad & g(\mathbf{x}) = 0 \end{aligned}$$

- Points on $g(\mathbf{x}) = 0$ must have $\nabla g(\mathbf{x})$ normal to surface
- A **stationary point** must have no change in f in the direction of the constraint surface, so $\nabla f(\mathbf{x})$ must also be normal to the surface.
 - So there must be some $\lambda \neq 0$ such that $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$
- Define **Lagrangian**:

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

- Stationary points of $L(\mathbf{x}, \lambda)$ have $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = \nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ and $\nabla_{\lambda} L(\mathbf{x}, \lambda) = g(\mathbf{x}) = 0$

Lagrange Multipliers



Consider the problem:

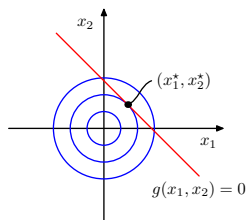
$$\begin{aligned} \max_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g(\mathbf{x}) = 0 \end{aligned}$$

- Points on $g(\mathbf{x}) = 0$ must have $\nabla g(\mathbf{x})$ normal to surface
- A **stationary point** must have no change in f in the direction of the constraint surface, so $\nabla f(\mathbf{x})$ must also be normal to the surface.
 - So there must be some $\lambda \neq 0$ such that $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$
- Define **Lagrangian**:

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

- Stationary points of $L(\mathbf{x}, \lambda)$ have $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = \nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ and $\nabla_{\lambda} L(\mathbf{x}, \lambda) = g(\mathbf{x}) = 0$

Lagrange Multipliers Example



- Consider the problem

$$\begin{aligned} \max_{\mathbf{x}} f(x_1, x_2) &= 1 - x_1^2 - x_2^2 \\ \text{s.t.} \quad g(x_1, x_2) &= x_1 + x_2 - 1 = 0 \end{aligned}$$

- Lagrangian:

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

- Stationary points require:

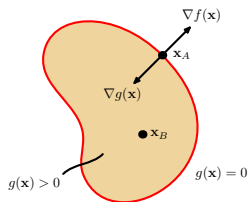
$$\partial L / \partial x_1 = -2x_1 + \lambda = 0$$

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

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

- So stationary point is $(x_1^*, x_2^*) = (\frac{1}{2}, \frac{1}{2})$, $\lambda = 1$

Lagrange Multipliers - Inequality Constraints



Consider the problem:

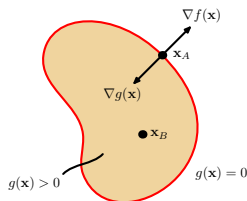
$$\begin{aligned} \max_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g(\mathbf{x}) \geq 0 \end{aligned}$$

- Optimization over a region – solutions either at stationary points (gradients 0) in region $g(\mathbf{x}) > 0$ or on boundary

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

- Solutions at stationary points of the Lagrangian with either:
 - $\nabla f(\mathbf{x}) = 0$ and $\lambda = 0$ (in region)
 - $\nabla f(\mathbf{x}) = -\lambda \nabla g(\mathbf{x})$ and $\lambda > 0$ (boundary, $>$ for maximizing f)
 - For both, $\lambda g(\mathbf{x}) = 0$
- Find stationary points of L s. t. $g(\mathbf{x}) \geq 0, \lambda \geq 0, \lambda g(\mathbf{x}) = 0$

Lagrange Multipliers - Inequality Constraints



Consider the problem:

$$\begin{aligned} \max_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g(\mathbf{x}) \geq 0 \end{aligned}$$

- Optimization over a region – solutions either at stationary points (gradients 0) in region $g(\mathbf{x}) > 0$ **or** on boundary

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

- Solutions at stationary points of the Lagrangian with either:
 - $\nabla f(\mathbf{x}) = 0$ and $\lambda = 0$ (in region)
 - $\nabla f(\mathbf{x}) = -\lambda \nabla g(\mathbf{x})$ and $\lambda > 0$ (boundary, $>$ for maximizing f)
 - For both, $\lambda g(\mathbf{x}) = 0$
- Find stationary points of L s. t. $g(\mathbf{x}) \geq 0, \lambda \geq 0, \lambda g(\mathbf{x}) = 0$

Outline

Maximum Margin Criterion

Math

Maximizing the Margin

Non-Separable Data

Now Where Were We

- So the optimization problem is now a constrained optimization problem:

$$\arg \min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2}$$

$$s.t. \quad \forall n, t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1$$

- For this problem, the Lagrangian (with N multipliers a_n) is:

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{n=1}^N a_n \{t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1\}$$

- We can find the derivatives of L wrt \mathbf{w}, b and set to 0:

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

$$0 = \sum_{n=1}^N a_n t_n$$

Dual Formulation

- Plugging those equations into L removes \mathbf{w} and b results in a version of L where $\nabla_{\mathbf{w}, b} L = 0$:

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

this new \tilde{L} is the **dual representation** of the problem (maximize with constraints)

- Another formula for finding b , like the bias in linear regression.

Outline

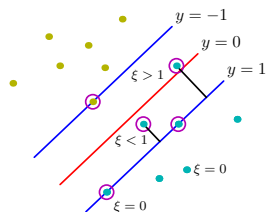
Maximum Margin Criterion

Math

Maximizing the Margin

Non-Separable Data

Non-Separable Data

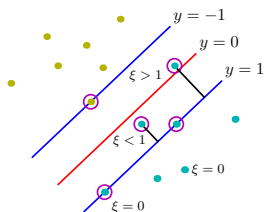


- For most problems, data will not be linearly separable (even in feature space ϕ)
- Can relax the constraints from

$$t_n y(\mathbf{x}_n) \geq 1 \quad \text{to} \quad t_n y(\mathbf{x}_n) \geq 1 - \xi_n$$

- The $\xi_n \geq 0$ are called **slack variables**
 - $\xi_n = 0$, satisfy original problem, so x_n is on margin or correct side of margin
 - $0 < \xi_n < 1$, inside margin, but still correctly classified
 - $\xi_n > 1$, mis-classified

Loss Function For Non-separable Data



- Non-zero slack variables are bad, penalize while maximizing the margin:

$$\min C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

- Constant $C > 0$ controls importance of large margin versus incorrect (non-zero slack)
 - Set using cross-validation
- Optimization is same quadratic, different constraints, convex

SVM Loss Function

- The SVM for the separable case solved the problem:

$$\arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$s.t. \quad \forall n, t_n y_n \geq 1$$

- Can write this as:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{\infty}(t_n y_n - 1) + \lambda \|\mathbf{w}\|^2$$

where $E_{\infty}(z) = 0$ if $z \geq 0$, ∞ otherwise

- Non-separable case relaxes this to be:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{SV}(t_n y_n - 1) + \lambda \|\mathbf{w}\|^2$$

where $E_{SV}(t_n y_n - 1) = [1 - y_n t_n]_+$ hinge loss

- $[u]_+ = u$ if $u \geq 0$, 0 otherwise

SVM Loss Function

- The SVM for the separable case solved the problem:

$$\arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$s.t. \quad \forall n, t_n y_n \geq 1$$

- Can write this as:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{\infty}(t_n y_n - 1) + \lambda \|\mathbf{w}\|^2$$

where $E_{\infty}(z) = 0$ if $z \geq 0$, ∞ otherwise

- Non-separable case relaxes this to be:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{SV}(t_n y_n - 1) + \lambda \|\mathbf{w}\|^2$$

where $E_{SV}(t_n y_n - 1) = [1 - y_n t_n]_+$ hinge loss

- $[u]_+ = u$ if $u \geq 0$, 0 otherwise

SVM Loss Function

- The SVM for the separable case solved the problem:

$$\arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$s.t. \quad \forall n, t_n y_n \geq 1$$

- Can write this as:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{\infty}(t_n y_n - 1) + \lambda \|\mathbf{w}\|^2$$

where $E_{\infty}(z) = 0$ if $z \geq 0$, ∞ otherwise

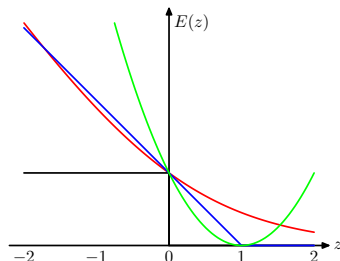
- Non-separable case relaxes this to be:

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N E_{SV}(t_n y_n - 1) + \lambda \|\mathbf{w}\|^2$$

where $E_{SV}(t_n y_n - 1) = [1 - y_n t_n]_+$ **hinge loss**

- $[u]_+ = u$ if $u \geq 0$, 0 otherwise

Loss Functions



- Linear classifiers, compare **loss function** used for learning
- $z = y_n t_n \geq 1$ iff there is an error (with $t_n \in \{+1, -1\}$).
 - Black is misclassification error
 - Transformed simple linear classifier, **squared error**: $(y_n - t_n)^2$
 - Transformed logistic regression, **cross-entropy error**: $t_n \ln y_n$
 - SVM, **hinge loss**: $\xi_n = [1 - y_n t_n]_+$
 - positive *only* if there is a mistake, otherwise 0.
 - encourages sparse solutions.

Two Views of Learning as Optimization

- The original SVM goal was of the form:
 - Find the simplest hypothesis that is consistent with the data, or
 - Maximize simplicity, given a consistency constraint.
- This general idea appears in much scientific model building, in image processing, and other applications.
- Bayesian methods use a criterion of the form
 - Find a trade-off between simplicity and data fit, or
 - Maximize sum of the type (data fit - λ simplicity)
 - e.g., $\ln(P(D|M)) - \lambda \ln(P(M))$ where the model prior M is higher for simpler models.

Pros and Cons of Learning Criteria

- The Bayesian approach has a solid probabilistic foundation in Bayes' theorem.
- Seems to be especially suitable for noisy data.
- The constraint-based approach is often easy for users to understand.
- Often leads to sparser simpler models.
- Suitable for “clean” data.

Conclusion

- Readings: Ch. 7 up to and including Ch. 7.1.2
- Many algorithms can be re-written with only dot products of features
 - We've seen NN, perceptron, regression; also PCA, SVMs (later)
- Maximum margin criterion for deciding on decision boundary
 - Linearly separable data
- Relax with slack variables for non-separable case
- Global optimization is possible in both cases
 - Convex problem (no local optima)
 - Descent methods converge to global optimum

Logistic Regression Learning: Iterative Reweighted Least Squares

- **Iterative reweighted least squares (IRLS)** is a descent method
 - As in **gradient descent**, start with an initial guess, improve it
 - Gradient descent - take a step (how large?) in the gradient direction
- IRLS is a special case of a **Newton-Raphson** method
 - Approximate function using second-order Taylor expansion:

$$\hat{f}(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \nabla f(\mathbf{w})^T (\mathbf{v} - \mathbf{w}) + \frac{1}{2} (\mathbf{v} - \mathbf{w})^T \nabla^2 f(\mathbf{w}) (\mathbf{v} - \mathbf{w})$$

- Closed-form solution to minimize this is straight-forward: quadratic, derivatives linear
- In IRLS this second-order Taylor expansion ends up being a weighted least-squares problem, as in the linear regression case
 - Hence the name IRLS

Logistic Regression Learning: Iterative Reweighted Least Squares

- **Iterative reweighted least squares (IRLS)** is a descent method
 - As in **gradient descent**, start with an initial guess, improve it
 - Gradient descent - take a step (how large?) in the gradient direction

- IRLS is a special case of a **Newton-Raphson** method
 - Approximate function using second-order Taylor expansion:

$$\hat{f}(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \nabla f(\mathbf{w})^T (\mathbf{v} - \mathbf{w}) + \frac{1}{2} (\mathbf{v} - \mathbf{w})^T \nabla^2 f(\mathbf{w}) (\mathbf{v} - \mathbf{w})$$

- Closed-form solution to minimize this is straight-forward: quadratic, derivatives linear
- In IRLS this second-order Taylor expansion ends up being a weighted least-squares problem, as in the linear regression case
 - Hence the name IRLS

Logistic Regression Learning: Iterative Reweighted Least Squares

- **Iterative reweighted least squares (IRLS)** is a descent method
 - As in **gradient descent**, start with an initial guess, improve it
 - Gradient descent - take a step (how large?) in the gradient direction
- IRLS is a special case of a **Newton-Raphson** method
 - Approximate function using second-order Taylor expansion:

$$\hat{f}(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \nabla f(\mathbf{w})^T (\mathbf{v} - \mathbf{w}) + \frac{1}{2} (\mathbf{v} - \mathbf{w})^T \nabla^2 f(\mathbf{w}) (\mathbf{v} - \mathbf{w})$$

- Closed-form solution to minimize this is straight-forward: quadratic, derivatives linear
- In IRLS this second-order Taylor expansion ends up being a weighted least-squares problem, as in the linear regression case
 - Hence the name IRLS

Logistic Regression Learning: Iterative Reweighted Least Squares

- **Iterative reweighted least squares (IRLS)** is a descent method
 - As in **gradient descent**, start with an initial guess, improve it
 - Gradient descent - take a step (how large?) in the gradient direction
- IRLS is a special case of a **Newton-Raphson** method
 - Approximate function using second-order Taylor expansion:

$$\hat{f}(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \nabla f(\mathbf{w})^T (\mathbf{v} - \mathbf{w}) + \frac{1}{2} (\mathbf{v} - \mathbf{w})^T \nabla^2 f(\mathbf{w}) (\mathbf{v} - \mathbf{w})$$

- Closed-form solution to minimize this is straight-forward: quadratic, derivatives linear
- In IRLS this second-order Taylor expansion ends up being a weighted least-squares problem, as in the linear regression case
 - Hence the name IRLS

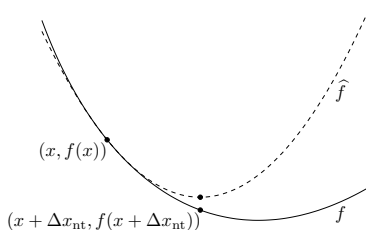
Logistic Regression Learning: Iterative Reweighted Least Squares

- **Iterative reweighted least squares (IRLS)** is a descent method
 - As in **gradient descent**, start with an initial guess, improve it
 - Gradient descent - take a step (how large?) in the gradient direction
- IRLS is a special case of a **Newton-Raphson** method
 - Approximate function using second-order Taylor expansion:

$$\hat{f}(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \nabla f(\mathbf{w})^T (\mathbf{v} - \mathbf{w}) + \frac{1}{2} (\mathbf{v} - \mathbf{w})^T \nabla^2 f(\mathbf{w}) (\mathbf{v} - \mathbf{w})$$

- Closed-form solution to minimize this is straight-forward: quadratic, derivatives linear
- In IRLS this second-order Taylor expansion ends up being a weighted least-squares problem, as in the linear regression case
 - Hence the name IRLS

Newton-Raphson



- Figure from Boyd and Vandenberghe, *Convex Optimization*
 - Excellent reference, free for download online
<http://www.stanford.edu/~boyd/cvxbook/>