

以影像區域特徵為基礎之網路釣魚網頁偵測

作者簡介

陳昭源

職稱：中央研究院資訊科學研究所研究助理

研究方向：網路安全議題

E-mail：jychen@csie.org

陳寬達 (陳存暘)

職稱：中央研究院資訊科學研究所助研究員

數位典藏與數位學習國家型科技數位技術研發計劃 Web 2.0 團隊共同主持人

研究方向：網路量測、網路安全及線上遊戲相關議題

Email：cychen@iis.sinica.edu.tw

引言

隨著網際網路的普及化及網路線上服務的興起，為我們帶來許多便利，但也同時提供惡意者更多的機會利用詐騙等手段竊取使用者的網路個人資料，甚至造成使用者財務損失。本文針對近年來相當興盛的網路犯罪手法—網路釣魚 (phishing) 介紹相關背景與反制技術及近期研究成果。

簡介

隨著網際網路的普及化，許多組織機構，包括政府、學校、銀行與購物網站等，紛紛推出各式各樣的網路線上服務。這些線上服務為我們帶來許多便利，省卻許多以往必須親自辦理所造成的往來時間消耗以及不便性，但也同時提供惡意者更多的機會利用詐騙等手段竊取使用者的網路個人資料，進而謀取不當利益，甚至造成使用者財務損失。

對於在網路上以製作仿冒網頁 (fake webpage) 或以偽造電子郵件 (spoofed email) 騙取使用者個人網路資料的犯罪方式，我們稱之為「網路釣魚 (phishing)」；而進行此種犯罪行為的惡意份子，就是「網路釣客 (phisher)」了。網路釣客常常以銀行或網路購物網站名義像受害者發出假冒的電子郵件，裡面可能有各種的理由並包含一個假冒的網址 (spoofed link) 讓不知情的受害者信以為真並點選連結，最後將受害者導向其所製作的甲冒網站填寫個人資料，如登入帳號、密碼、或信用卡號碼。近幾年來，網路釣魚的犯罪手法，已經造成許多受害者的金錢損失，個案數目也持續成長中。

根據反釣魚工作組織 (Anti-Phishing Working Group, APWG) 報導，釣魚網頁的數量從 2004 年 7 月起即以每個月增加 28% 的速度成長，除此之外，通常會有 5% 的釣魚郵件接收者會受騙。根據調查，在 2007 年 9 月，被舉報或被發現的網路釣魚攻擊個案超過 6.6 萬起，且有高達 95% 的仿冒詐騙目標涉及金融服務和網路零售商，例如，eBay 與 PayPal。根據知名科技產業調查與顧問公司 Gartner 在 2007 年的調查顯示，網路釣魚攻擊在美國所造成的損失已經超過 32 億美元。由此可知，網路釣魚攻擊已經成為一個嚴峻的資料安全問題，且與個人隱私及財產息息相關。

釣魚網頁偵測技術

老鼠？老虎？傻傻分不清楚！

為了讓受害者相信偽造的釣魚網頁是合法的官方網頁，網路釣客在設計釣魚網頁時，通常無論在內容或排版上皆製作的與原始官方網頁幾乎一模一樣；除此之外，廣告橫幅也常被加入釣魚網頁中，目的是為了要讓受害者點選以將其導向至另一個惡意網站 (malicious site) 或是在其電腦中植入惡意程式，如：木馬 (Trojan horse)，或是執行惡意 script 程式 (malicious script)。

以網路釣客最熱門的攻擊目標：eBay 為例，圖 (1) 是官方網站上直接截取的登入畫面；圖 (2) 和圖 (3) 都是偽造的釣魚網頁。我們可以觀察發現，圖 (1) 和圖 (2) 的畫面幾乎沒有差別，圖 (2) 僅僅移除了原始網站 eBay 商標下方的彩色橫條，並將 eBay 標誌稍微縮小；圖 (3) 則是在頁面最上方插入一個廣告橫幅，除此之外，與圖 (1) 也是幾乎完全相同。從這幾張官方網頁與釣魚網頁之間的比較，我們可以瞭解，對一般使用者來說，想辨別目前所瀏覽的網頁是否為釣魚網頁的確有其困難度，稍加不注意便會落入釣客所設下的捕魚陷阱中！



圖 (1)：eBay 官方登入頁面

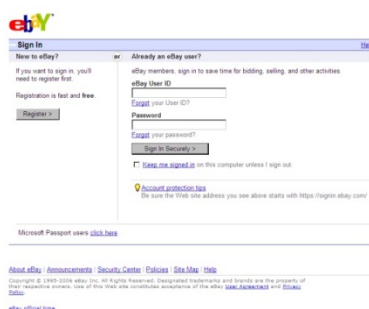


圖 (2)：eBay 釣魚網頁 (商標修改)



圖 (3)：eBay 釣魚網頁 (橫幅廣告)

近年來，已經有許多反釣魚技術 (anti-phishing technique) 被研究、提出，例如：SpoonGuard、iTrustPage。大部分被提出的方法都依靠文字內容分析、HTML 原始碼區塊分析或 URL 分析；然而，越來越多的釣魚網頁使用圖片取代靜態內容或插入非 HTML 組件 (例如：Flash 和 Java Applet) 以避免被這些以文字為基礎的 (text-based) 反釣魚分析工具辨認出。在這種情況下，那些以文字分析為主的反釣魚工具將會失效，因為即使像

圖 (1) 和圖 (2) 這兩個頁面，雖然有著類似的外觀，但有可能其背後其實已經是由完全不同的網頁元件所組成。

為了要克服上述的局限性，我們在最近提出了一種基於影像區域特徵 (image local feature based) 的反釣魚偵測技術。更詳細地說，我們分析一個可疑網頁的區域內容特徵，並與其官方網頁的分析結果進行比對，藉此比對結果所獲得的參考指標分數 (reference score) 來評估此頁面為釣魚網頁的可能性。我們的方法不仰賴文字分析，因此網頁的元件或程式碼構成不會對我們所提出的方法造成影響。

影像內容特徵 (image content feature) 可以包括影像所使用的顏色、質地、形狀與空間關係的組成。為了要對內容特徵進行分析，一般說來，我們會使用內容描述器 (content descriptor) 來進行內容特徵分析，而根據其適用範圍可以分成兩種類型：(1) 全域性內容描述器 (global content descriptor) 與 (2) 區域性內容描述器 (local content descriptor)；前者使用整個圖像的內容特徵來分析一個影像，後者則使用影像中各個區域或物件的內容特徵進行分析。

由於釣魚網頁在大多數情況下會改動官方網頁部分的內容或配置，因此，我們合理地認為區域性內容特徵比之全域性的內容特徵，對於頁面局部的變化有較強的容忍度，也更適合用來作為釣魚網頁的分析。以上述情況為例，我們所提出的方法，即使是網頁組成元件、程式碼不同甚至是插入橫幅廣告的情形，在實驗中皆能夠正確偵測出原本不存在於官方網頁中的橫幅廣告，並找到頁面其他相似部份，如商標、文字、輸入欄位等，進而成功地偵測釣魚網頁。

以下，我們將先介紹目前幾種較熱門的網頁釣魚偵測及反制技術。

釣魚網頁偵測技術

目前，主要的釣魚偵測與反制技術可分成兩大類：(1) 電子郵件層級反制：包括認證 (authentication)、來源過濾 (source filtering) 和內容分析 (content analysis)；(2) 瀏覽器插件 (browser plug-in)：最普遍的方法是比對網址黑名單 (URL blacklist) 及安全名單 (white list)，或是進行即時頁面屬性分析 (real-time page property analysis)。

電子郵件層級反制

利用電子郵件過濾反制網路釣魚的方式和反垃圾郵件 (anti-spam) 的方法頗為相似，這兩種技術都是以防止有危害性的電子郵件達到目標用戶為目的，常見的方法為來源過濾或內容分析。這些反制措施的成功率依賴許多關鍵因素，舉例來說，過濾規則的建立對使用來源過濾方式便有直接影響。如果用戶不小心設置過於嚴格的規則，則不僅僅是網路釣魚郵件，連同正常的電子郵件都將被認定為網路釣魚威脅。在許多情況下，這些郵件將根據用戶的設定直接被送往垃圾郵件匣或被刪除，如此將造成用戶遺失重要訊息的可能。反之，寬鬆的規則，用戶將得煩惱處理無數網路釣魚郵件出現在郵件收件匣的情形。

微軟 (Microsoft) 和雅虎 (Yahoo) 分別在網頁式的電子郵件系統中發佈電子郵件認證協議：SenderID 和 DomainKeys 技術。這兩種解決方案皆是用來核實電子郵件是否來自可信信任來源，目前皆只能在其各自的服務和產品中免費使用。

瀏覽器插件

比對網址黑名單是瀏覽器反釣魚插件最常使用的方法。每當瀏覽一個新網頁時，瀏覽器便會先比對是否與黑名單中的網址相符，如果答案是肯定的，警告訊息將顯示給用戶；舉例來說，微軟目前最新版本的 Internet Explorer 7 遇到不安全的網頁時，會將其位址欄變成紅色，若無則表示此網頁應為正常合法網頁。

黑名單可以儲存在本機或託管在中央伺服器。使用黑名單來確任釣魚網頁與防毒程式使用的是同一個概念，而防毒軟體以經使用此方式達數十年之久。黑名單的有效性，取決於以下幾個因素：

- 1) 涵蓋範圍 (coverage)：很明顯的，有多少釣魚網頁被列入這份黑名單將直接影響釣魚網頁偵測。
- 2) 公信力 (credibility)：指黑名單的正確度。從用戶的觀點來說，每個錯誤回報的釣魚網頁，都將削弱用戶的信任度。
- 3) 更新頻率 (update frequency)：使用黑名單的一個最重要的安全漏洞，是該名單可能無法辨識新的釣魚威脅，因為新的資料尚未更新。即便是黑名單已定期更新，但研究發現，大部分的釣魚網站生命週期相當短暫，在黑名單定期更新之前便已有受害者蒙受損失。因此，一個不合時宜的黑名單，對保護使用者是無效的。

目前，最為人熟知的黑名單是由谷歌 (Google) 和微軟所維護的版本，分別預設套用在

Mozilla Firefox 2 和 Internet Explorer 7 瀏覽器中。Firefox 2 可以透過一個開放的協議連接到任何可用的黑名單提供商，並會自動不斷更新，目前預設為谷歌的黑名單。Internet Explorer 7 使用線上資料庫的即時網路釣魚報告 (real-time phishing report) 並結合本身內建的黑名單。不過，根據分別由微軟、Mozilla 等公司及盧德爾等人的實驗結果，這兩種方式的釣魚網頁偵測正確率均低於 90%，最差會低於 60%。

網頁屬性分析通常包括 HTML 原始碼、影像大小或位置、網站上的商標等，大部分都是基於文字內容或原始碼的分析。也有人分析網址來分辨其是否為釣魚網頁的偽造地址。偽造地址通常含有部分目標網站的網址，網路釣客試圖欺騙使用者，使釣魚網頁的網址看起來像是合法的，例如，<http://cox.net/www.amazon.com/gp/sign-in.html>，如果使用這沒有仔細檢視，是很容易被誤認為是亞馬遜的合法位址：<https://www.amazon.com/gp/sign-in.html>。通過事先設定的白名單，我們可以很容易地找出偽造的網址。

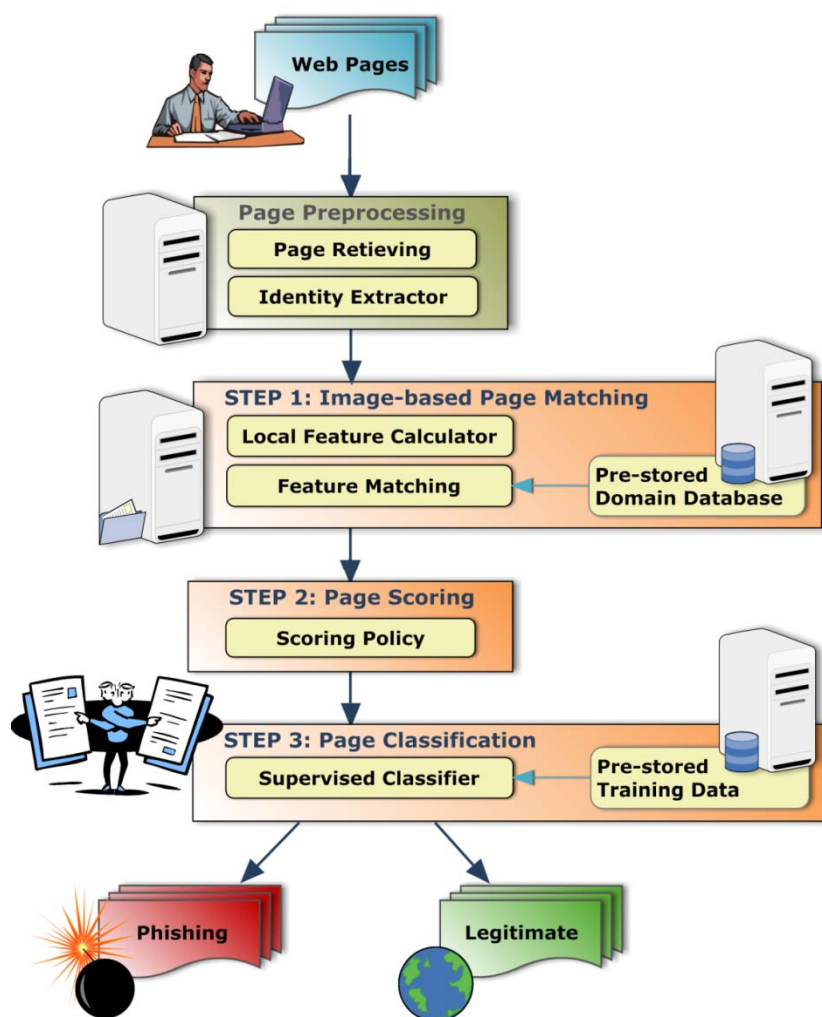
雖然文字內容或原始碼分析似乎是非常簡便且有效率的方式，但隨著網路頻寬發展越來越不受限及網路應用服務的興起，越來越多的網頁除了使用 HTML 之外還包含了各種 script 語言或多媒體內容，例如：JavaScript、Flash、ActiveX 元件；網路釣客可以有許多選擇使用不同的技術偽造目標頁面，如此一來，文字內容或原始碼分析的方式來偵測釣魚網頁便顯得無用武之地了。

接著，我們要介紹目前仍正在進行研究且有不錯成效的「以影像區域特徵為基礎的釣魚網頁偵測方法」。此方法的特性便是針對近年來的網頁發展趨勢，不仰賴文字內容或原始碼分析，將網頁轉換成影像再分析其內容特徵，藉以有效反制網路釣客的詭計。

影像區域特徵為基礎的釣魚網頁偵測

方法簡述

在開始進行網頁分析之前，與傳統方式不同的地方，我們必須先將網頁轉換成影像格式，並利用關鍵字分析猜測此網頁可能的偽造目標。之後，進行以下三個分析步驟，如圖(4)所示：



圖(4)：釣魚網頁偵測流程圖

1) 網頁配對

首先，我們使用區域特徵描述器分析可能的釣魚網頁，找出此網頁的特徵點 (salient

point) · 與事先建立的合法網頁資料庫內的網頁特徵點資訊進行比對。在我們的實驗中，我們使用 Context Contrast Histogram (CCH) 描述器來分析頁面區域特徵。

2) 分數評估

經過特徵比對之後，我們由所獲得的資訊為每個可能的釣魚網頁計算五種參考指標分數，分別為：(1) 網頁相似度、(2) 成功配對密度、(3) 有效格點涵蓋率、(4) 每一有效格點成功配對率與 (5) 每一有效格點成功配對數標準差。

為了要檢驗配對的正確性，我們使用 K 均值聚類演算法將比對的兩個頁面所有已配對的特徵點分群 (group)；除此之外，我們將整個頁面切割成 $n \times n$ 個格子點，並考慮有效格點與其中的配對情形，用來評估已配對的特徵點對整個頁面的重要性。有效格點的定義為若一格點中存有某個低限值 δ 以上數量的特徵點（無論是否有被配對）時，此格點我們將其視為有效格點，反之，則視為無效格點。

3) 分類判定

最後一步，我們利用單純貝葉氏分類器 (naive Bayesian classifier) 以網頁的參考指標分數為基礎進行網頁分類，以判定此網頁是否為釣魚網頁。

實驗結果

接下來，我們將說明幾個實驗的結果，以證明我們的方法是有效的。

A. 網頁比對測試

為了檢視我們的方法是否能夠區分出不同網站間的網頁差異以及釣魚網站與官方網站間的相似性，我們將網頁比對結果顯示如下。圖 (5) 上半部顯示的是插入一個橫幅廣告的 eBay 釣魚網頁，下半部則是 eBay 原始官方網頁。此圖表示的是經過配對後的特徵點以 K 均值聚類演算法分群後的情形。不同顏色的實心點表示他們所屬的不同群組，也就是成功配對的特徵點；而藍色空心點，則表示沒有被分類到相同組別的配對點，我們將之視為無效的配對點。圖 (6) 顯示的是將配對點連接之後的結果，由此可知，透過我們的方法可以辨識出釣魚網頁與原始官方網頁之間極高的相似性。圖 (7) 顯示的則為不同網站間的網頁比對。由此圖可明顯看出，我們的方法，可明顯區分出兩者之間的差異，也就是說，對於網頁之間的差異具有極高的鑑別力。



圖 (5) : 特徵點比對圖

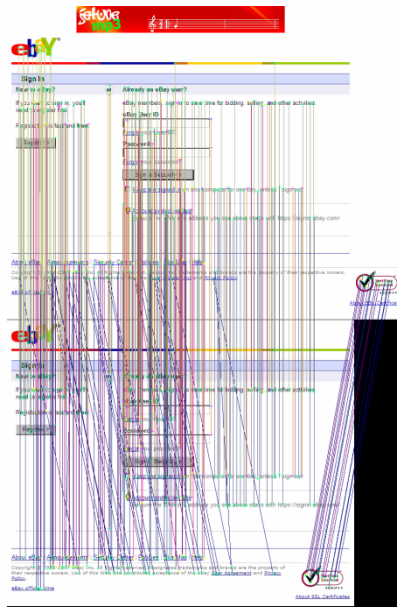


圖 (6) : 釣魚網頁比對



圖 (7) : 不同網站網頁比對

B. 前五大釣魚目標網站測試

表格 (1) 是我們所收集的釣魚網頁資料統計，前五大目標釣魚網站分別為：eBay、PayPal、Marshall & Ilsley Bank、Charter One Bank 以及 Bank of America。這五個目標釣魚網站所收集到的釣魚網頁數目佔了我們所收集的網頁總數的一半以上。

網站名稱	釣魚網頁紀錄個數
eBay	701
PayPal	632
Marshall & Ilsley Bank	138
Charter One Bank	116
Bank of America	51
釣魚網頁總數：2058 頁，74 網站	

表格 (1) : 前五大釣魚目標網站

在這個實驗中，我們將前五大目標網站的釣魚網頁與相同數目的合法網頁參混，用以檢驗是否所提出的方法能正確分辨釣魚網頁；另外，我們同時以之前已經被提出的另一個以影像全域特徵為基礎的釣魚網頁偵測方法進行實驗作為對照組。圖 (8) 顯示在前五大目標網站中的偵測正確率。我們的方法 (CCH) 正確辨識率達到至少 97%，而對照組 (EMD) 則顯得相對不穩定，最差僅有 58% 正確辨識率。圖 (9) 顯示在前五大目標網站中的偵測錯誤率，我們將錯誤率細分為錯誤負判率 (false negative rate)，即將釣魚網頁誤判為合法網頁的比例與錯誤正判率 (false positive rate)，即將合法網頁

誤判為釣魚網頁的比例。由圖 (9) 可觀察得知，我們的方法 (CCH) 無論是錯誤負判率或錯誤正判率皆維持在非常低的狀態；反之，對照組的錯誤負判率至少有 6%，最嚴重的情形高達七成，顯示其對於網頁之間的鑑別度不高。圖 (10) 顯示的是兩個方法的受試者工作曲線比較，我們的方法 (CCH) 的曲線下面積為 0.998，而對照組 (EMD) 為 0.956。

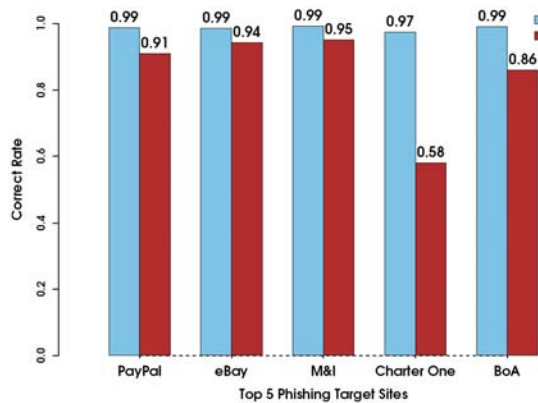


圖 (8)：前五大釣魚目標網站偵測效能 (正確率)

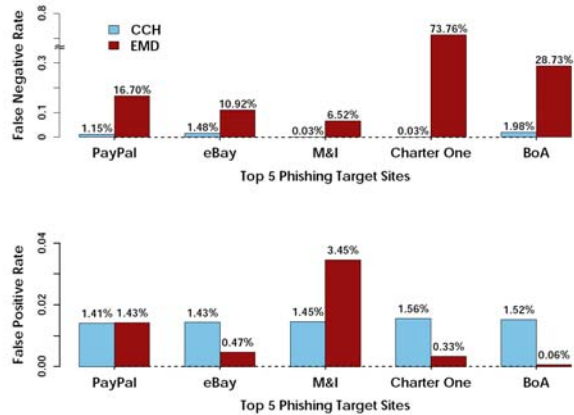


圖 (9)：前五大釣魚目標網站偵測效能 (錯誤率)

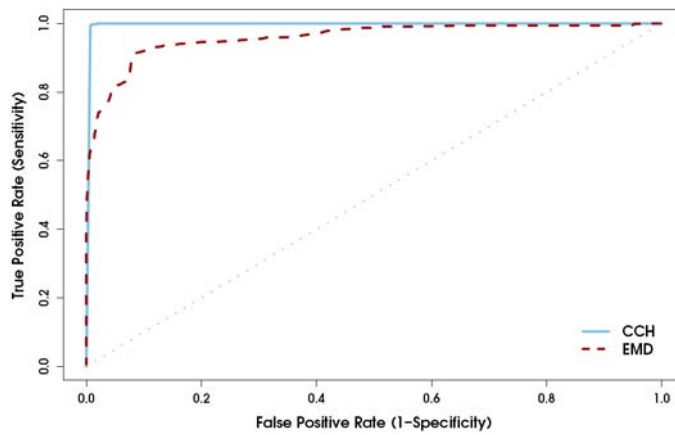


圖 (10)：受試者工作特徵曲線 (ROC curve) 比較

結論

網路釣魚已經成為網路資訊安全以及個人隱私的主要威脅，幾年來，已經有無數的人遭受釣魚網頁詐騙造成金錢的損失。釣魚網頁通常看起來與合法的官方網頁有幾許相似，並可能含有部份的網頁修改或是以不同的網頁元件構成；然而，目前大部分的釣魚網頁偵測方法皆把重點擺在文字內容的分析，在這些情形下將會無法正確辨識釣魚網頁。

我們提出以影像區域特徵為基礎的偵測方法。透過實驗，顯示我們的方法有高達 97% 以上的正確率，且維持相當低的錯誤率，不受網站不同的影響。對照組的表現亦驗證之前的假設：區域特徵比對相較於全域特徵比對，更適合用來作為分析偵測釣魚網頁。在未來，我們將會以此方法為基礎，發展瀏覽器插件的反釣魚工具，為個人網路資訊安全把關！