



- ▶ Is the use of deprecated PHP releases associated with LAMP?
 - ⇒ The so-called “personal home page” (a.k.a. **PHP**) is still clearly the most popular server-side web programming language
 - ⇒ Here, the term **deprecation** refers to major PHP release branches that are no longer supported by the PHP project
 - ⇒ The abbreviation **LAMP** (Linux, Apache, MySQL, PHP) is taken as an idiom for a still typical open source web stack



1. Analytics

- ⇒ A well-established niche for web **fingerprinting**
- ⇒ Software and **release engineering**

2. Security

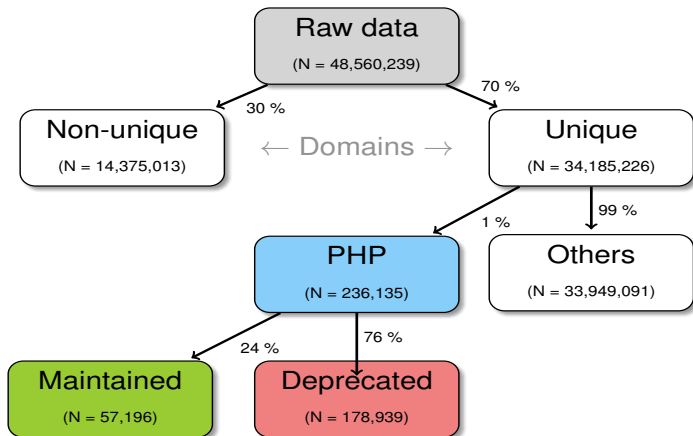
- ⇒ PHP applications are widely exploited, but also the language itself has been exposed to many **vulnerabilities**
- ⇒ Automatic scanners frequently probe for deprecated releases (cf. WordPress); national cyber security concerns

3. Data Mining

- ⇒ What can be done with existing **HTTP header data**?



- ▶ The **HTTP Archive** web crawling project
 - ⇒ Covers most of the popular web sites
 - ⇒ Sponsored by Google, Mozilla, O'Reilly, etc.
 - ⇒ Previously, mainly used for performance evaluations
- ▶ Generally **big data** (high-volume, potentially high-variety)
 - ⇒ Though, a **bi-monthly** crawling schedule (low-velocity)
 - ⇒ Only a single **snapshot** utilized for this research
- ▶ Third-party data sources imply **research constraints**
 - ⇒ Notably, only HTTP headers are used for fingerprinting





- ▶ **Dependent:** PHP deprecation (0/1)
- ▶ **Independent:** five categories derived from HTTP headers;
 1. HTTP servers (e.g., Apache, IIS, Nginx, etc.)
 2. Operating systems (e.g., Linux, Windows, Ubuntu, etc.)
 3. Apache-specific modules (e.g., OpenSSL)
 4. One PHP-specific variable (Suhosin)
 5. Control variables (e.g., top-30 TLDs)
- ▶ **Operationalization** with regular expression matching
 - ⇒ When possible, includes also rough release information about the independent variables
 - ⇒ Thus, false positives (negatives) are likely



- ▶ **Estimation** with the standard logistic regression:

$$\begin{aligned} p_i &= \Pr\{\text{PHP is deprecated} = 1 \mid \mathbf{x}_i\}, \\ &= F(\mathbf{x}_i' \boldsymbol{\beta}) = \frac{\exp(\mathbf{x}_i' \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i' \boldsymbol{\beta})}, \end{aligned}$$

- ▶ In total, **57 parameters**, including a constant
- ▶ **Parameter evaluation** with the so-called **marginal effects**:

$$\hat{\beta}_j = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i' \hat{\boldsymbol{\beta}}) \hat{\beta}_j \quad (\text{continuous})$$

$$\hat{\beta}_j = \frac{1}{n} \sum_{i=1}^n \left[F(\mathbf{x}_i' \hat{\boldsymbol{\beta}}) - F \left(\hat{\beta}_0 + \sum_{s=1, s \neq j}^k \hat{\beta}_s x_{is} \right) \right] \quad (\text{discrete})$$



- ▶ **Model reduction** sought with a combined forward and backward stepwise selection algorithm
 - ⇒ Drops uninformative variables (based on AIC)
 - ⇒ Implemented in most statistical software packages
- ▶ Reduced model compared also with a **likelihood ratio test**
- ▶ **Model evaluation** with a five-fold cross-validation and the three conventional metrics (accuracy, precision, and recall)
- ▶ **Resampling** (during model training) used for accounting the **unbalanced sample** (i.e., about 76 % deprecated)



Table: Deprecated PHP Releases (as of 10 Dec 2015)

Branch	Deprecation date	Days since EOL
5.4	3 September 2015	89
5.3	14 August 2014	474
5.2	6 January 2011	1790
5.1	24 August 2006	3386
5.0	5 September 2005	3739
4.4	7 August 2008	2672
4.3	31 March 2005	3897
4.2	6 September 2002	4834
4.1	12 March 2002	5012
4.0	23 June 2001	5274
3.0	20 October 2000	5520

Results (2/4)

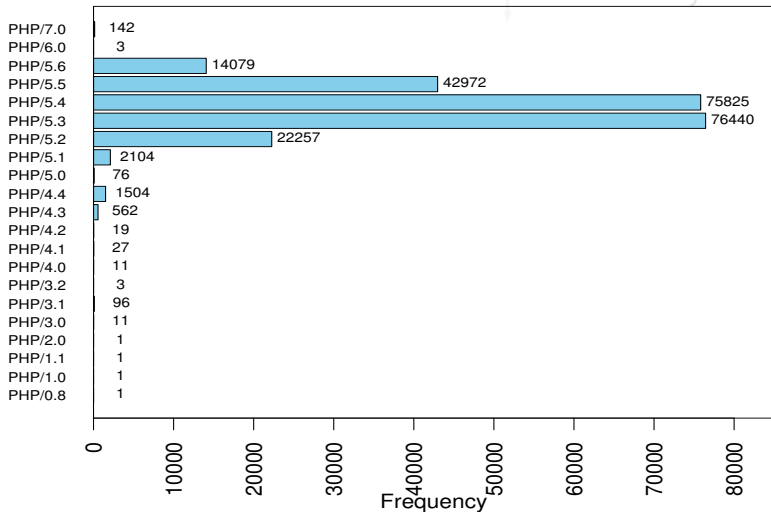


Figure: PHP releases

Table: Performance of $\mathcal{M}_{\text{Full}}^i$

$i =$	1	2	3	4	5
Unbalanced					
Accuracy	0.780	0.775	0.778	0.778	0.780
Precision	0.970	0.968	0.969	0.971	0.970
Recall	0.789	0.785	0.787	0.786	0.788
Balanced					
Accuracy	0.635	0.636	0.635	0.638	0.633
Precision	0.640	0.643	0.640	0.645	0.636
Recall	0.842	0.839	0.841	0.840	0.840

Results (4/4)

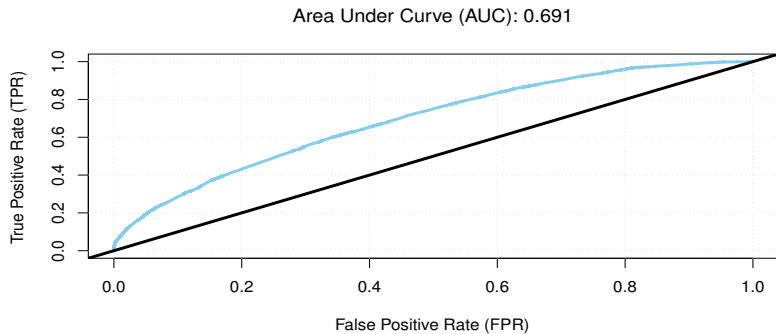


Figure: Performance of Unbalanced $\mathcal{M}_{\text{Full}}^1$



- ▶ Thus, classification **performance is modest** at best
 - ⇒ Irrespective of the balancing
 - ⇒ Unlikely due to logistic regression *per se*
 - ⇒ Though, all of the variables pass the model reduction
- ▶ The marginal effects indicate **interpretable** results
 - ⇒ In particular, the **LAMP hypothesis holds**; deprecated PHP releases are strongly associated with Linux and Apache
 - ⇒ No notable “geographic” variation (given popular TLDs)
 - ⇒ Some variation across Linux distributions
- ▶ In general, HTTP **header data has limited use**



1. **Fingerprinting** improvements
 - ⇒ HTTP headers are adequate for web server and application (PHP) layers; operating system probing is much more difficult (though, cf. network scanners)
2. **Combination of data** sources
 - ⇒ Popularity ranks, PageRank, DNS, etc.
3. **Longitudinal** analysis
 - ⇒ The evolution of PHP is arguably more interesting from a release engineering perspective
4. Security and **malware analysis**
 - ⇒ The use of deprecated PHP releases is presumably more relevant as an independent rather than dependent variable



Turun yliopisto
University of Turku

Thank you

Questions?