A changeset-based approach to assess source code density and developer efficacy A step towards data-driven software cost- and effort estimation models

> Sebastian Hönel, Morgan Ericsson, Welf Löwe, Anna Wingkvist Data-driven Software and Information Quality Group – Centre for Data Intensive Sciences and Applications Linnaeus University, 351 95 Växjö, Sweden *{sebastian.honel, morgan.ericsson, welf.lowe, anna.wingkvist}@lnu.se*

Context



 46 - private Semaphore startLock; 47 + private Semaphore lockStart; 48 - private Semaphore shutdownLock; 49 + private Semaphore lockShutdown; 105 - this.startLock = new Semaphore(1, 1); 106 - this.shutdownLock = new Semaphore(0, 1); 	<pre>134 - /// <summary> 135 - /// Initialize the client with the specific address, port and f 136 - /// </summary> 80 + /// Initialize the client with the specific address, port and format. 81 + /// 82 + /// <param name="multicastAddress"/> 83 + /// <param name="multicastAddress"/> 84 + /// <param name="multicastAddress"/> 85 + /// <param name="multicastAddress"/> 86 + /// <param name="multicastAddress"/> 87 + /// <param name="multicastAddress"/> 88 + /// <param name="multicastAddress"/> 89 + /// <param name="multicastAddress"/> 80 + /// <param name="multicastAddress"/> 81 + /// <param name="multicastAddress"/> 82 + /// <param name="multicastAddress"/> 83 + /// <param name="multicastAddress"/> 84 + //</pre>	icastPor	
<pre>106 + this.lockStart = new Semaphore(0, 5); 107 + this.lockShutdown = new Semaphore(0, 2); 109 - DesiredLatency = 80 110 + DesiredLatency = 50 116 - var localEp = new IPEndPoint(IPAddress.Any, this.multicastPort); 117 + using (var ms = new MemoryStream()) 118 + {</pre>	<pre>84 + /// <pre>cpream name= "waveFormat"></pre> 84 + /// <pre>cpream name="waveFormat"></pre> 85 + public UdpWaveStreamPlayer(IPAddress multicastAddress, int multicastPor 86 + { 87 + this.multicastAddress = multicastAddress; 88 + this.multicastPort = multicastPort; 89 + this.waveFormat = waveFormat; 90 +</pre>	40	■+Clones +Dead code
<pre>119 + var localEp = new IPEndPoint(IPAddress.Any, this.multicastPc 118 - this.udpClient.Client.SetSocketOption(SocketOptionLevel.Socket, 119 - this.udpClient.Olient.Bind(localEp); 120 - this.udpClient.Client.SetSocketOption(SocketOptionLevel.Sock 121 + this.udpClient.Client.Bind(localEp); 123 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 124 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 125 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 126 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 127 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 128 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 129 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 120 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 121 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 123 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 124 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 125 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 126 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 127 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 128 + this.udpClient.JoinMulticastGroup(this.multicastAddress); 129 + this.udpClient.StatAddress); 129 + this.udpClient.StatAdd</pre>	<pre>90 + 91 + this.startLock = new Semaphore(1, 1); 92 + this.shutdownLock = new Semaphore(0, 1); 93 + this.output = new WaveOut() 94 + { 95 + DesiredLatency = 80 96 + }; 97 +</pre>		 + Whitespace + Comments
<pre>123 - var rsws = new BufferedWaveProvider(new WaveFormat()) 124 - { 125 - BufferDuration = TimeSpan.FromMilliseconds(500) 126 - }; 126 + var rsws = new BufferedWaveProvider(new WaveFormat()) 127 + { 128 + BufferDuration = TimeSpan.FromMilliseconds(500) 129 + };</pre>	<pre>97 + 98 + this.udpClient = new UdpClient 99 + { 100 + ExclusiveAddressUse = false 101 + }; 102 + 103 + this.Volume = 100;</pre>	₩ 20	 Net size Crease size
130 + }	104 + } 105 +	10	

 $D = -\frac{1}{c}$

The Density *D* is the ratio between the Functionality *F* and the Size *S*. A density of 1 means that all code contributes to the functionality. Therefore, the notions of D_{max} and D_{min} relate to the net- and gross-size of the software (S_n, S_g) .



The effort, when measured using functionality instead of size, is more meaningful. The equation $E = {}^{F}/{}_{t}$ is therefore equivalent.



The Productivity P for software is defined as Effort per Size, according to ISBSG¹. We can hence deduce a direct correlation between productivity and density using the equation $P = \frac{D}{t}$.



Commits can be associated with one developer. This is important for effort- and cost-estimation models, so that the planner can assess the productivity of their mixed-skill team more accurate. When analyzed as a time series, the density can be forecast, allowing the planner to adjust the (ongoing) project's estimates and to prevent fallouts and plan resources accordingly.

Public repository hosting services allow us to freely access and analyze millions of commits across a variety of programming languages, project setups and developer capabilities. We are fetching a great number of changesets and try to determine their densities. We expect to unearth correlations between the different aspects



The gross size is determined by the plain count of lines of code. Removing unnecessary comments, space lines, dead and unreachable code as well as any cloned functionality yields the net size of the software.

We are focusing on examining how close the gross size of changes is compared to the net size using various methods, such as clone detection or string similarity measures to find the density of contained functionality. We want do examine which method is the most suitable for determining the density.



Evaluation

- We learned about the significance of density and how it impacts large and small changesets.
- We were able to identify the most and least suitable methods to determining density.
- We learned which methods can best approximate a certain notion of time using a certain notion of size.
- We have analyzed about 1,650 open-source software repositories with ~80k commits across various programming languages and project setups. Large deviations in gross- and net-size were discovered.



Three different notions productivity, based on gross-

- We can detect typical developer commit-behavior and how it negatively impacts time estimates and leads to distortions of measurements of productivity.
- We have unveiled shortcomings in automatically deducing spent time from commit timestamps using git-hours and found better-than-default settings by denoising the data.
- Currently, the distorted measurements of time do not allow to correlate any (specific) notion of size to it.







