

	MAPE	$Pr_{miss}$	$Pr_{fa}$	$Pr_{error}$	$T$ -score	Posts	Visits
$w = 5, \mathbf{v},  \mathbf{v}  = 5$	8.441	0.922	0.064	0.493	1601.321	33.000	545.371
$w = 5, \mathbf{v},  \mathbf{v}  = 10$	8.632	0.924	0.064	0.494	1593.763	33.000	545.206
$w = 5, \mathbf{v},  \mathbf{v}  = 15$	8.913	0.924	0.064	0.494	1594.276	33.000	545.381
$w = 5, \mathbf{v},  \mathbf{v}  = 20$	9.382	0.923	0.063	0.493	1597.533	33.000	545.495
$w = 5, \mathbf{v},  \mathbf{v}  = 25$	9.905	0.927	0.063	0.495	1597.295	33.000	545.619
$w = 5, \mathbf{v},  \mathbf{v}  = 30$	10.836	0.925	0.063	0.494	1587.734	33.000	545.619
$w = 5, \mathbf{v},  \mathbf{v}  = 35$	12.044	0.927	0.064	0.495	1608.698	33.000	545.722
$w = 5, \mathbf{v},  \mathbf{v}  = 40$	12.400	0.927	0.063	0.495	1593.062	33.000	545.649
$w = 5, \mathbf{v},  \mathbf{v}  = 45$	12.462	0.928	0.063	0.496	1587.913	33.000	545.649
$w = 5, \mathbf{v},  \mathbf{v}  = 50$	12.994	0.926	0.063	0.495	1581.241	33.000	545.804

Table 1: Experiment results: Varying vocabulary size

# 1 Method

We extracted the timestamp, author and text content for each post in each thread from the forums. Using a sliding window training method, we grouped consecutive  $w$  posts together for each thread, and performed regression on  $\Delta_t$ . More formally, we are trying to approximate a function  $f$  such that  $f(\mathbf{x}_{t-w}, \dots, \mathbf{x}_{t-1}) \approx \Delta_t$ , where  $\mathbf{x}_t$  is the feature vector of a post made at time  $t$ , and  $\Delta_t$  is the time between the  $t$ -th post and the  $(t - 1)$ -th post. The following are the features used:

**Previous time differences** All the time differences between posts made in the window. ( $\mathbf{t}_\Delta$ )

**Time-based features** Day of week, Hour of day. Provides contextual information about when the post was made. ( $\mathbf{t}_{\text{ctx}}$ )

**Content features (text)** Word frequency counts are used for this set of experiments.

Using regression, we find the top  $K$  variables that the actual  $\Delta_t$  depends on.

Table 1 reflect the results of the experiments done with varying values of  $K$ .

## 1.1 Potential errors

To be thorough, let us also enumerate the types of errors that a model making predictions could encounter.

The model can potentially make a prediction such that the next visit comes before the arrival of the next post. The predictions being made are the  $\Delta_t$  between the posts, rather than the visitation times, hence, it is possible for the model to make a prediction that occurs before the current time. An erroneous prediction can also cause the crawler to come in before the next post (two, or more, visits, but nothing new fetched). Errors of this type waste bandwidth, since the crawler will make an unnecessary visit to the page.

Another type of error would have the prediction causing the next visit to come some time after a post. Since most predictions are almost never fully accurate, there will be some time between the post is made and the page is fetched. These errors are still relatively acceptable, but the time difference between the post arriving and the visit should be minimised. The visit could also come more than one post later. Errors of this kind incur a penalty on the freshness of the data, more so than the after one post, especially if the multiple posts are far apart time-wise.

In the following experiments, the threads chosen from our extracted dataset are those with a 100 to 1000 posts. This amounted to 97 threads. The first 75% of the thread was used as training data, while the remaining 25% was used as test data. We used Support Vector machines for this regression task, employing a Radial Basis Function kernel as our learning algorithm.

The SVR module from the Python library scikit-learn was used in the implementation of this experiment.

## 1.2 Evaluation metrics

We use *Mean Absolute Percentage Error* (MAPE), to measure the performance of the learnt model. This value is given by

$$\frac{1}{N} \sum_{i=1}^N \left| \frac{A_i - F_i}{A_i} \right|$$

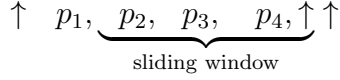


Figure 1: An example of the sliding window metric. The metric is made up of two components: First, the probability that, given at least one post is present in the window, there are more visits than post. Secondly, the probability that there are more posts than visits for a given window. The weighted sum of this gives the overall  $Pr_{error}$

where  $A_i$  is the actual value, and  $F_i$  is the forecasted value for the instance  $i$ . Realistically, the model would not be able to come into contact with every possible window, since chances are it will make an error that causes it to visit a thread late, causing it to miss two posts or more. This value does not reflect how well the model will do in a real-time setting, but gives an idea of how far off the model is given a window.

We also want to know the *timeliness* of the model's visits. Yang et. al. [?] has a metric for measuring this. Taking  $\Delta t_i$  as the time difference between a post  $i$  and it's download time, the timeliness of the algorithm is given by

$$T = \frac{1}{N} \sum_{i=1}^N \Delta t_i$$

A good algorithm would give a low  $T$ -score. However, a crawler that hits the site repeatedly performs well according to this metric. The authors account for this by setting a bandwidth (fixed number of pages per day) for each iteration of their testing. In our experimental results, we also take into account the number of page requests made in comparison to the number of posts.

Viewing the posts made during the thread's lifetime as segmentations of the thread, and the visits made as hypotheses of where the segmentations are, we use the  $Pr_{error}$  metric from Georgescu et. al., 2006 as a measure of how close the predictions are to the actual posts. An example can be seen in Figure 1.

	MAPE	$Pr_{miss}$	$Pr_{fa}$	$Pr_{error}$	$T$ -score	Posts	Visits
Average $w = 5$	330.285	0.951	0.054	0.502	6418.208	33.000	498.742
Average $w = 10$	305.557	0.955	0.053	0.504	4598.955	31.680	497.351
Average $\Delta_t$	174.004	0.938	0.065	0.501	1764.474	34.000	574.031
$w = 5, \mathbf{t}_\Delta$	18.884	0.931	0.064	0.498	1541.595	33.000	547.062
$w = 5, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	18.885	0.931	0.064	0.498	1541.592	33.000	547.062
$w = 5, \mathbf{v}$	9.382	0.923	0.063	0.493	1597.533	33.000	545.495
$w = 5, \mathbf{v}, \mathbf{t}_\Delta$	18.877	0.931	0.064	0.498	1541.588	33.000	547.062

Table 2: Experiment results

## 2 Results

The results for experiments done with different combinations of the above specified features are shown in Table 2.

Overall average and window average perform significantly worst than the learnt models, as reflected in both the MAPE and the  $T$ -score. There is also a slight improvement in the  $Pr_{error}$  in the learnt models.

Taking into account the  $T$ -score and the number of visits together, would seem that  $\mathbf{t}_\Delta$ , features representing the previous time intervals, are important features when determining the next time interval. In the absence of these features, we observe that the  $T$ -score increases by about 1%. In this experiment, we use purely word frequency features. This gives only a slight improvent over not using them.

High values for  $Pr_{miss}$  and low for  $Pr_{fa}$ , are due to  $Pr_{miss}$  being conditioned on there being a post within the window. Since the posts come in bursts, visits are fairly periodic, and intervals between visits are larger than post bursts. When there are more posts than visits in windows with posts, we have higher  $Pr_{miss}$

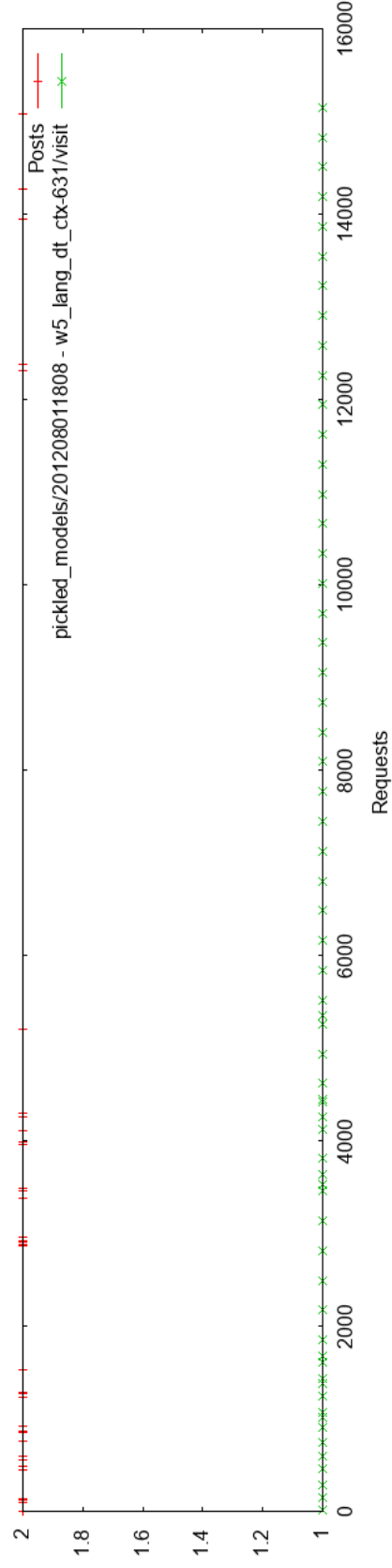


Figure 2: Visitation chart for a model using the  $w = 5$ ,  $t_{\Delta}$ ,  $t_{ctx}$ ,  $\mathbf{w}$  feature set. Invalid Predictions = 0.758,  $Pr_{error} = 0.485$ ,  $T$ -score = 119.612, Posts = 41, Visits = 62

## 2.1 Discounted sum of previous instances

The current method uses only information on the current  $w$  posts. However, posts made further in the history of the thread may have an effect on when the latest posts arrive. The magnitude of this effect, however, may diminish over time.

Following this intuition we attempt to use a discounted sum over previous posts' word frequency vector:

$$\mathbf{x}_t = \mathbf{v}_t + \alpha \mathbf{x}_{t-1}$$

where  $\mathbf{x}_t$  is the feature vector at post  $t$ , and  $\mathbf{v}$  is the word frequency vector.  $\alpha$  is the *discount factor* and satisfies  $0 \leq \alpha < 1$ .

## 2.2 Stochastic Gradient Descent

We attempt to use stochastic gradient descent to estimate the function  $f$ . However, during runtime, instead of using a static function, we continue to allow  $f$  to vary whenever new posts and their update times are observed. Since  $f(\mathbf{x}_{t-w}, \dots, \mathbf{x}_{t-1}) > 0$ , we used a scaled sigmoid function,

$$f(\mathbf{X}) = \frac{\lambda}{1 + e^{\mathbf{w} \cdot \mathbf{X}}}$$

where  $\lambda$  is the scaling factor. This results in  $f : \mathbb{R}^{|\mathbf{X}|} \rightarrow (0, \lambda)$ . Bounding the estimation function between 0 and  $\lambda$  allows us to restrict the prediction from becoming negative, or, becoming exceedingly huge. For our purposes, we set  $\lambda = Q_3 + 2.5(Q_3 - Q_1)$ , where  $Q_n$  is the value at the  $n$ -th quartile.

The resulting update rule for  $\mathbf{w}$  is then given by,

$$\Delta \mathbf{w}_i = \eta \underbrace{\left( \widehat{\Delta}_t - \Delta_t \right)}_{\text{error term}} \underbrace{\left( f(\mathbf{X})(1 - f(\mathbf{X})) \right)}_{\text{gradient}} \mathbf{X}_i$$

which is similar to the delta update rule found in artificial neural networks. We omit

	MAPE	$Pr_{miss}$	$Pr_{fa}$	$Pr_{error}$	$T$ -score	Visit/Post
Average $w = 5$	330.285	0.951	0.054	0.502	6418.208	16.464
Average $w = 10$	305.557	0.955	0.053	0.504	4598.955	17.291
Average $w = 15$	308.547	0.954	0.054	0.504	3833.605	18.337
Average $w = 20$	265.124	0.953	0.054	0.504	3340.929	18.102
Average $w = 25$	257.844	0.955	0.052	0.503	3186.309	17.927
Average $w = 30$	244.988	0.957	0.050	0.504	2859.380	18.362
$w = 5, \mathbf{t}_\Delta$	18.884	0.931	0.064	0.498	1541.595	17.907
$w = 10, \mathbf{t}_\Delta$	19.647	0.937	0.061	0.499	1488.688	18.371
$w = 15, \mathbf{t}_\Delta$	20.195	0.939	0.061	0.500	1443.138	19.234
$w = 20, \mathbf{t}_\Delta$	20.220	0.938	0.059	0.499	1584.171	18.880
$w = 25, \mathbf{t}_\Delta$	20.953	0.937	0.056	0.496	1649.098	18.612
$w = 30, \mathbf{t}_\Delta$	21.242	0.941	0.054	0.498	1626.782	18.984
$w = 5, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	18.885	0.931	0.064	0.498	1541.592	17.907
$w = 10, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	19.647	0.937	0.061	0.499	1488.688	18.371
$w = 15, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	20.195	0.939	0.061	0.500	1443.138	19.234
$w = 20, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	20.220	0.938	0.059	0.499	1584.171	18.880
$w = 25, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	20.953	0.937	0.056	0.496	1649.098	18.612
$w = 30, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	21.242	0.941	0.054	0.498	1626.782	18.984
$w = 5, \mathbf{v}$	9.382	0.923	0.063	0.493	1597.533	17.862
$w = 10, \mathbf{v}$	13.863	0.934	0.061	0.498	1551.375	18.339
$w = 15, \mathbf{v}$	13.217	0.934	0.060	0.497	1507.589	19.182
$w = 20, \mathbf{v}$	14.849	0.930	0.059	0.494	1630.643	18.848
$w = 25, \mathbf{v}$	17.542	0.930	0.055	0.493	1700.990	18.579
$w = 30, \mathbf{v}$	18.627	0.937	0.054	0.496	1653.156	18.959
$w = 5, \mathbf{v}, \mathbf{t}_\Delta$	18.877	0.931	0.064	0.498	1541.588	17.907
$w = 10, \mathbf{v}, \mathbf{t}_\Delta$	19.645	0.937	0.061	0.499	1488.680	18.371
$w = 15, \mathbf{v}, \mathbf{t}_\Delta$	20.193	0.939	0.061	0.500	1443.130	19.234
$w = 20, \mathbf{v}, \mathbf{t}_\Delta$	20.220	0.938	0.059	0.499	1584.171	18.880
$w = 25, \mathbf{v}, \mathbf{t}_\Delta$	20.953	0.937	0.056	0.496	1649.098	18.612
$w = 30, \mathbf{v}, \mathbf{t}_\Delta$	21.242	0.941	0.054	0.498	1626.782	18.984

Table 3: Experiment results: Varying feature sizes

	MAPE	$Pr_{miss}$	$Pr_{fa}$	$Pr_{error}$	$T$ -score	Visit/Post
Average $w = 5$	330.285	0.974	0.010	0.492	23941.231	0.620
Average $w = 10$	305.557	0.976	0.010	0.493	24932.464	0.635
Average $w = 15$	308.547	0.975	0.010	0.493	27999.236	0.639
Average $w = 20$	265.124	0.975	0.010	0.493	25785.328	0.651
Average $w = 25$	257.844	0.977	0.010	0.494	27307.482	0.658
Average $w = 30$	244.988	0.979	0.010	0.495	26039.058	0.655
$w = 5, \mathbf{t}_\Delta$	18.884	0.975	0.011	0.493	23348.003	0.697
$w = 10, \mathbf{t}_\Delta$	19.647	0.974	0.011	0.492	24208.018	0.704
$w = 15, \mathbf{t}_\Delta$	20.195	0.976	0.011	0.494	27456.946	0.706
$w = 20, \mathbf{t}_\Delta$	20.220	0.976	0.011	0.494	25246.673	0.720
$w = 25, \mathbf{t}_\Delta$	20.953	0.975	0.011	0.493	26976.438	0.719
$w = 30, \mathbf{t}_\Delta$	21.242	0.975	0.011	0.493	25776.611	0.724
$w = 5, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	18.885	0.975	0.011	0.493	23348.002	0.697
$w = 10, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	19.647	0.974	0.011	0.492	24208.018	0.704
$w = 15, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	20.195	0.976	0.011	0.494	27456.946	0.706
$w = 20, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	20.220	0.976	0.011	0.494	25246.673	0.720
$w = 25, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	20.953	0.975	0.011	0.493	26976.438	0.719
$w = 30, \mathbf{t}_\Delta, \mathbf{t}_{ctx}$	21.242	0.975	0.011	0.493	25776.611	0.724
$w = 5, \mathbf{v}$	9.382	0.969	0.011	0.490	23380.715	0.686
$w = 10, \mathbf{v}$	13.863	0.971	0.011	0.491	24258.557	0.698
$w = 15, \mathbf{v}$	13.217	0.973	0.010	0.492	27518.222	0.689
$w = 20, \mathbf{v}$	14.849	0.970	0.011	0.490	25276.263	0.711
$w = 25, \mathbf{v}$	17.542	0.970	0.011	0.490	27016.238	0.710
$w = 30, \mathbf{v}$	18.627	0.972	0.011	0.491	25792.662	0.714
$w = 5, \mathbf{v}, \mathbf{t}_\Delta$	18.877	0.975	0.011	0.493	23347.998	0.697
$w = 10, \mathbf{v}, \mathbf{t}_\Delta$	19.645	0.974	0.011	0.492	24208.011	0.704
$w = 15, \mathbf{v}, \mathbf{t}_\Delta$	20.193	0.976	0.011	0.494	27456.937	0.706
$w = 20, \mathbf{v}, \mathbf{t}_\Delta$	20.220	0.976	0.011	0.494	25246.673	0.720
$w = 25, \mathbf{v}, \mathbf{t}_\Delta$	20.953	0.975	0.011	0.493	26976.438	0.719
$w = 30, \mathbf{v}, \mathbf{t}_\Delta$	21.242	0.975	0.011	0.493	25776.611	0.724

Table 4: Experiment results: Using exponential increase

	MAPE	$Pr_{miss}$	$Pr_{fa}$	$Pr_{error}$	$T$ -score	Posts	Visits
$\alpha = 0.0, \mathbf{v}$	13.217	0.934	0.060	0.497	1507.589	30.402	528.247
$\alpha = 0.1, \mathbf{v}$	13.626	0.932	0.060	0.496	1500.257	30.402	529.608
$\alpha = 0.2, \mathbf{v}$	14.021	0.934	0.060	0.497	1487.145	30.402	530.175
$\alpha = 0.3, \mathbf{v}$	14.446	0.933	0.060	0.497	1483.811	30.402	530.825
$\alpha = 0.4, \mathbf{v}$	14.612	0.933	0.060	0.497	1449.927	30.402	531.742
$\alpha = 0.5, \mathbf{v}$	14.579	0.931	0.061	0.496	1434.925	30.402	532.701
$\alpha = 0.6, \mathbf{v}$	14.668	0.934	0.061	0.497	1445.370	30.402	533.165
$\alpha = 0.7, \mathbf{v}$	14.907	0.931	0.061	0.496	1442.314	30.402	534.629
$\alpha = 0.8, \mathbf{v}$	16.208	0.933	0.060	0.497	1440.828	30.402	532.680
$\alpha = 0.9, \mathbf{v}$	18.595	0.931	0.061	0.496	1472.814	30.402	532.691

Table 5: Experiment results: Window size of 15, using discounted sum feature vector at  $t - 1$ .



the scaling factor in the gradient as it is a constant and then experiment with various values of  $\eta$ , the learning rate.