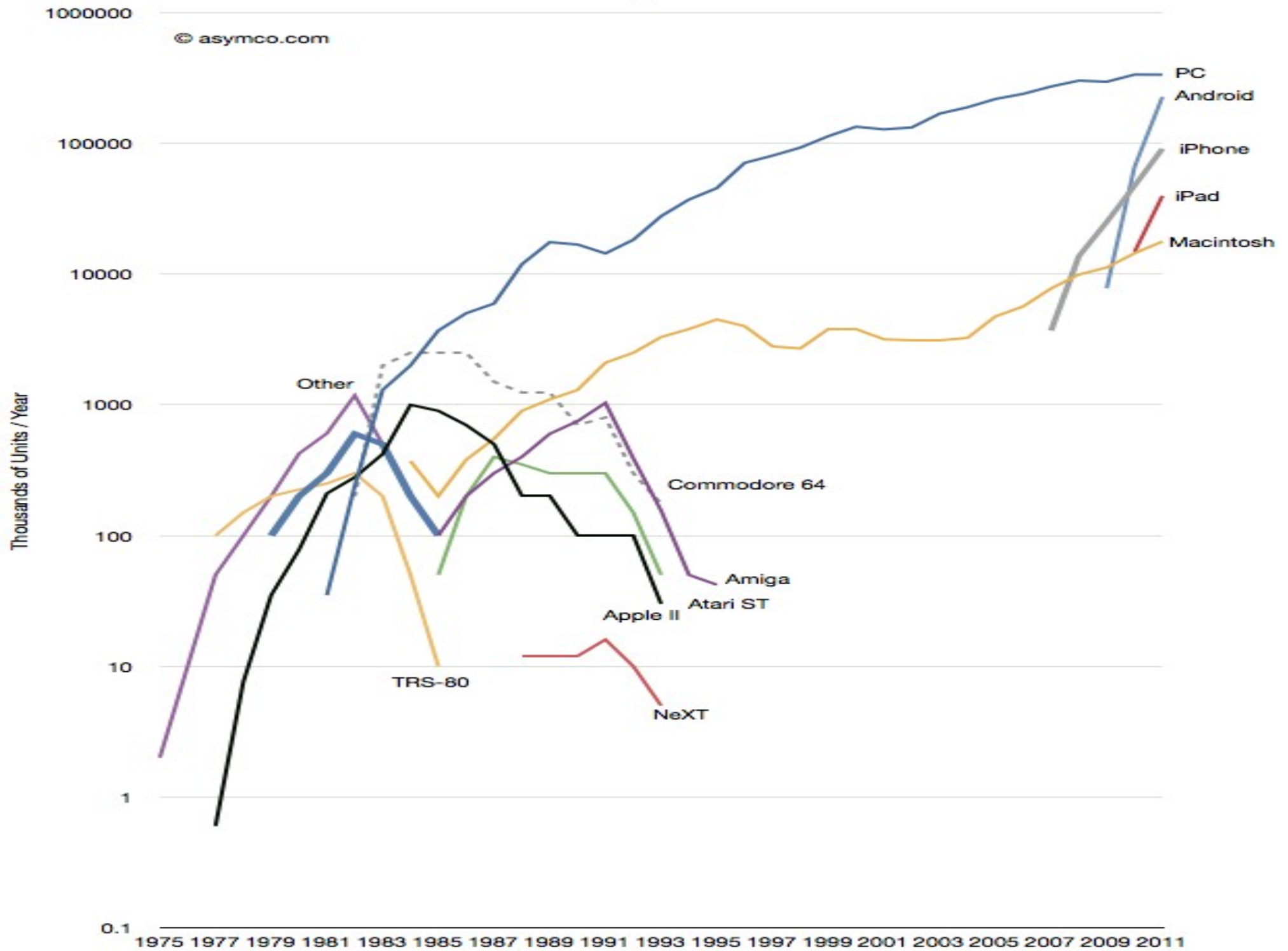
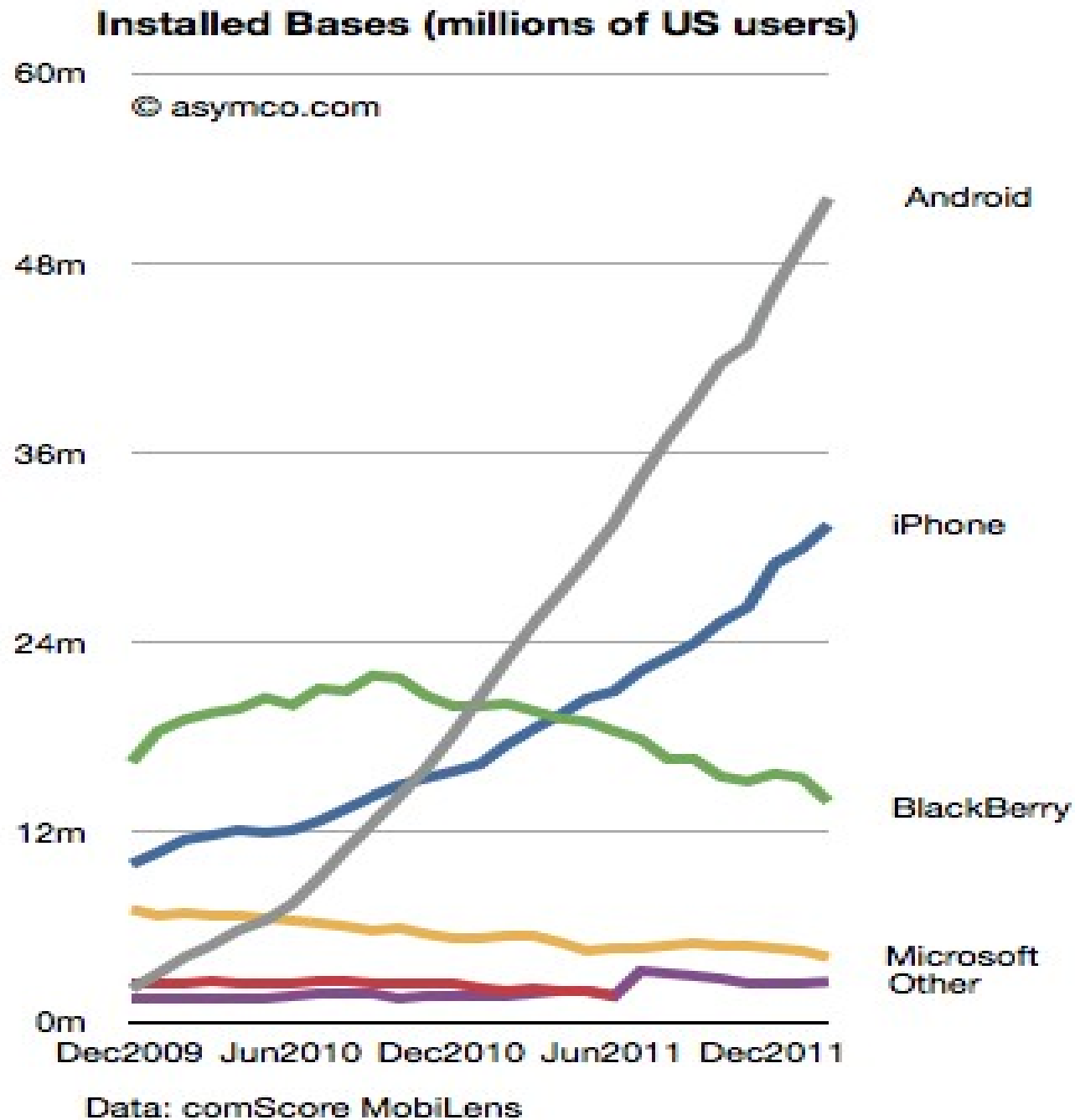


Learning in networks of millions of nodes

Márk Jelasity

Units Shipped Per Year





Motivation

- Explosive growth of smart phone platforms, and
- Availability of sensor and other contextual data
- Makes collaborative data mining possible
 - Health care: following and predicting epidemics, personal diagnostics
 - Smart city: traffic optimization, accident forecasting
 - (predicting earthquakes, financial applications, etc)
- P2P networks, grid, etc, are also relevant platforms

P2P system model

- Large number (millions or more) computers (nodes)
- Packet switched communication
 - Every node has an address
 - Any node can send a message to any given address
- Messages can be delayed or lost, nodes can crash
- (in parallel computing this is similar to the model of asynchronous chaotic iterations)

Fully distributed data

- Horizontal data distribution
- Every node has very few records, we assume they have **only one**
- We do not allow for moving data, only local processing (**privacy preservation**)
- We require that the models are cheaply available for all the nodes

Illustration: averaging

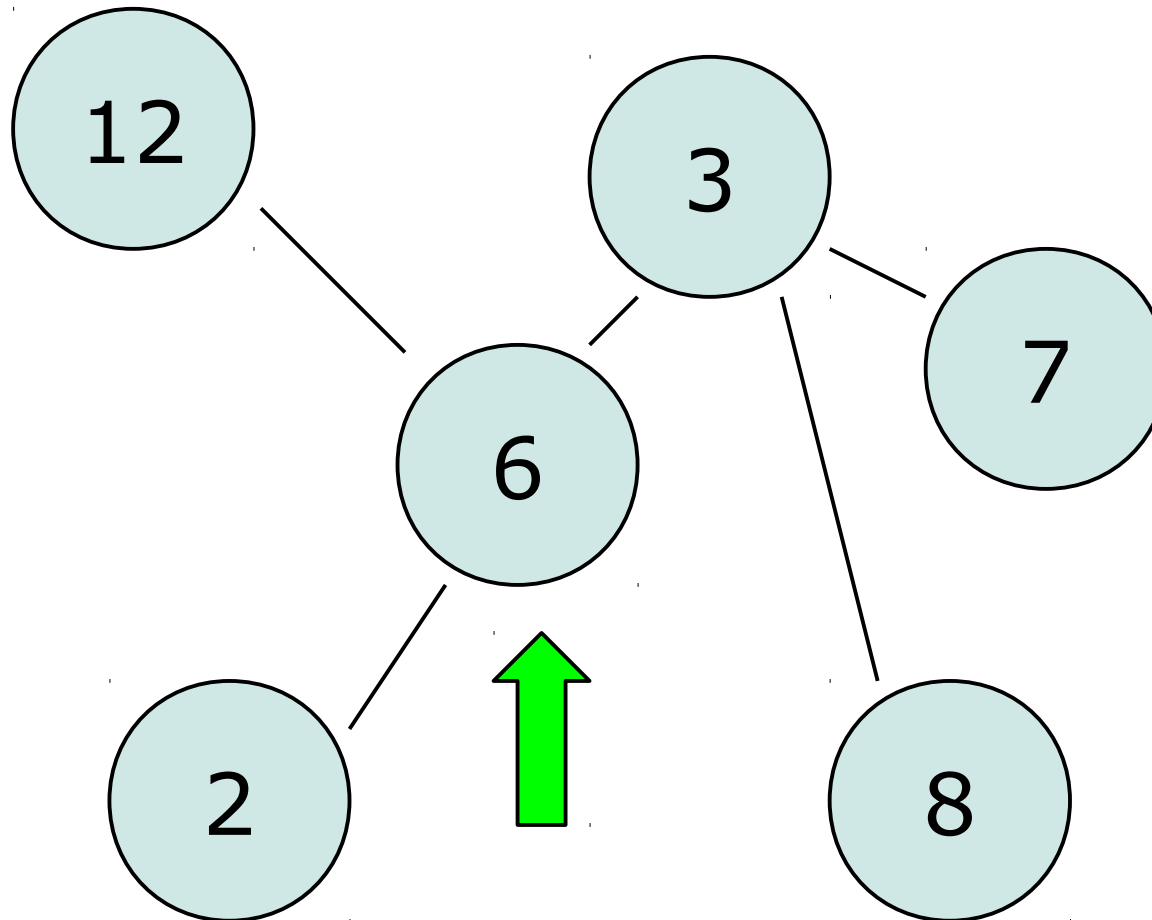


Illustration: averaging

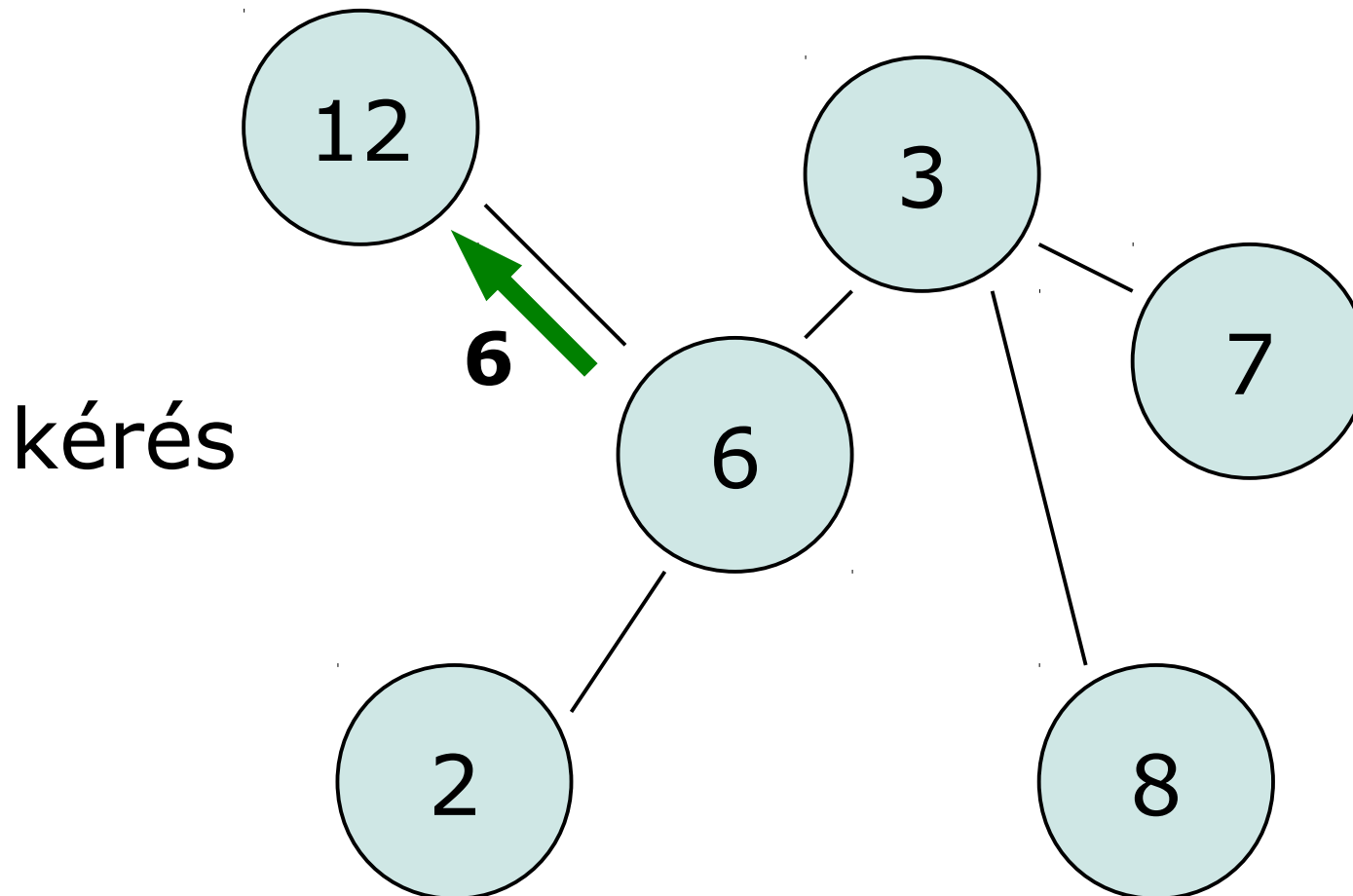


Illustration: averaging

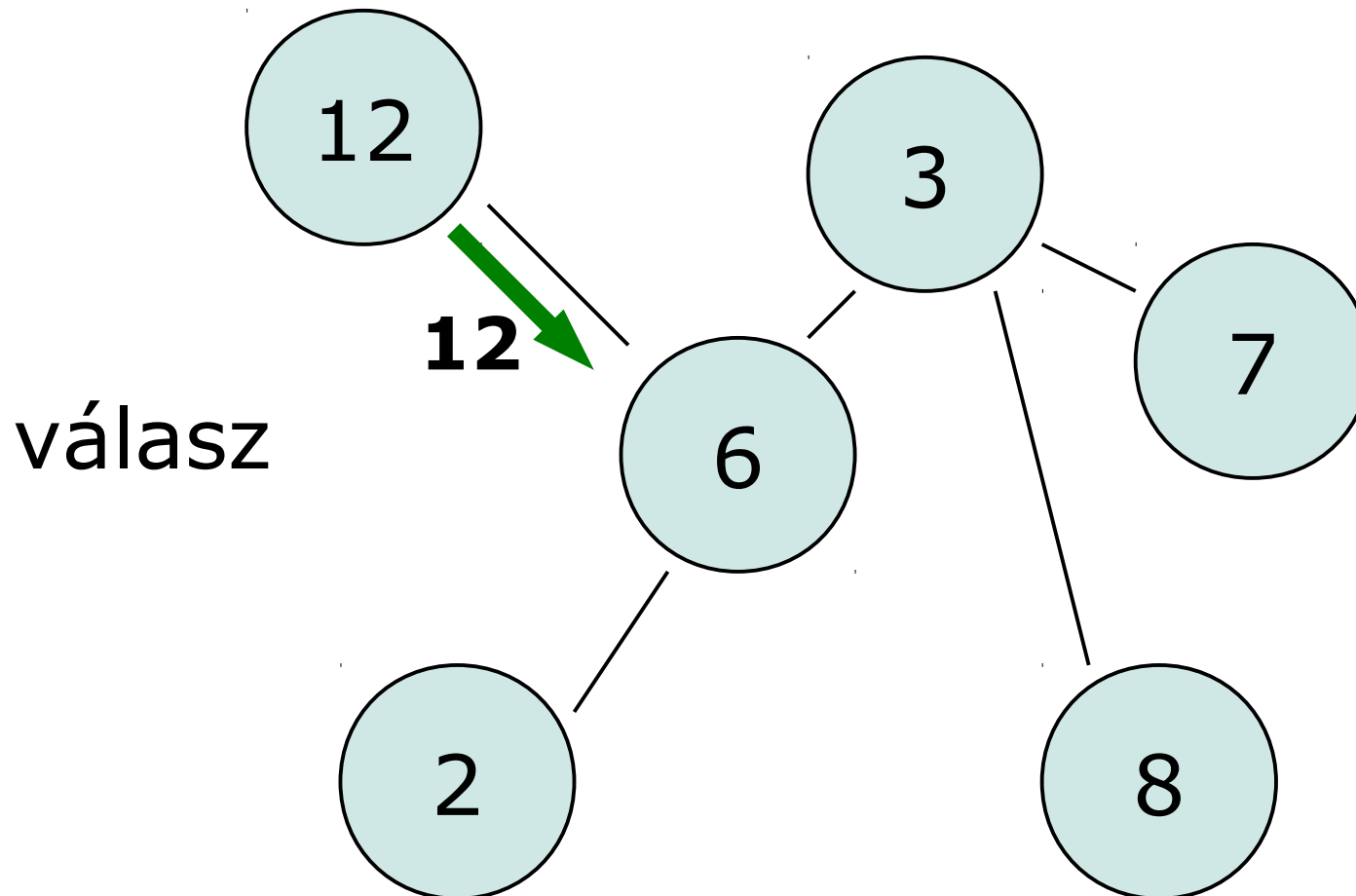
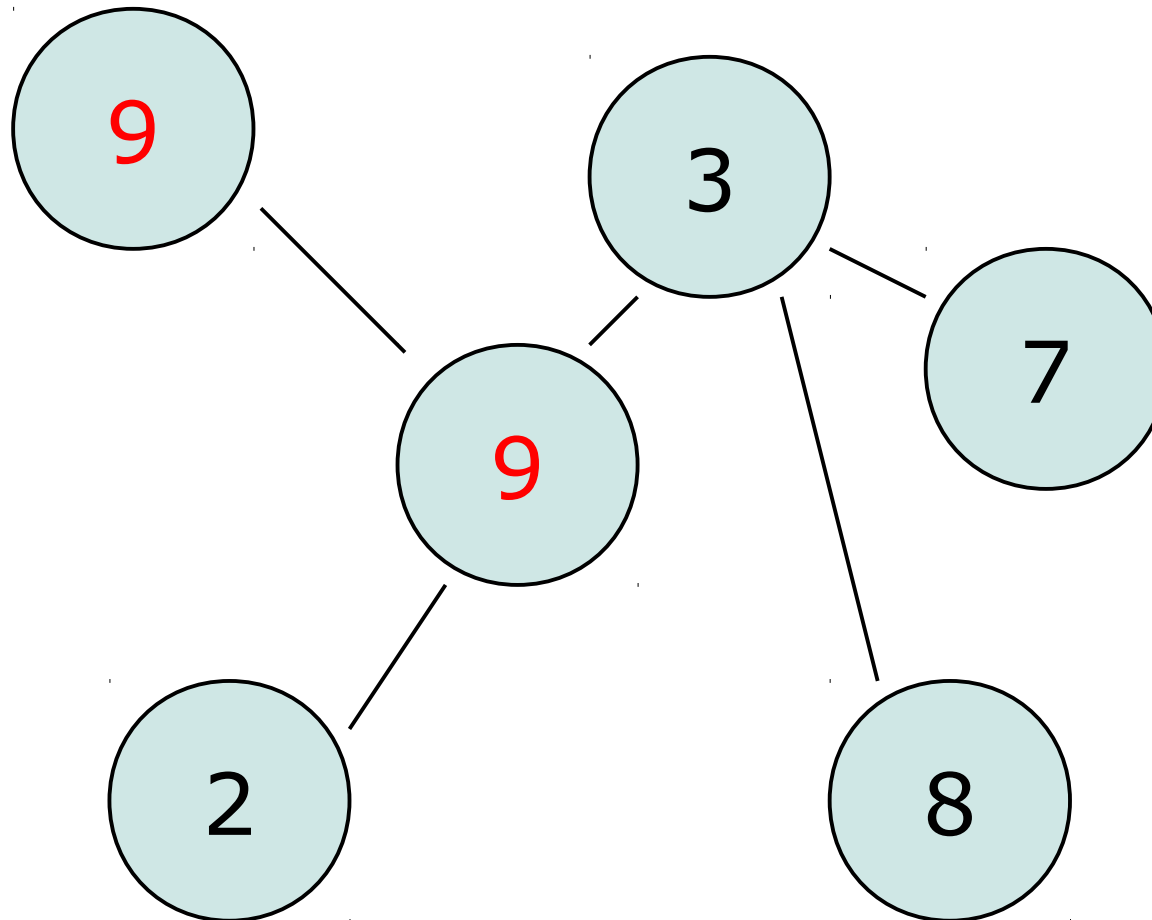


Illustration: averaging

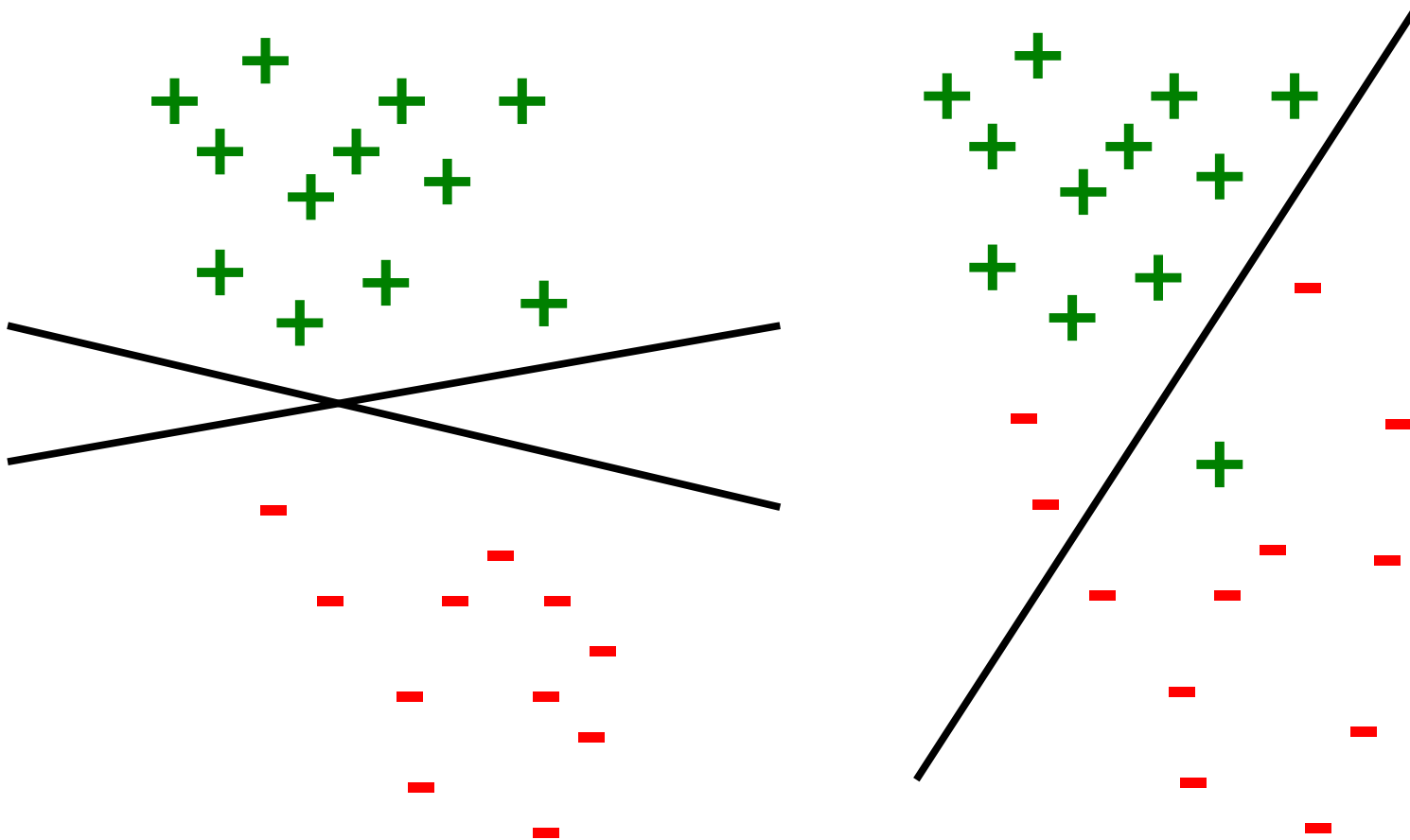


$$(12+6)/2=9$$

Classification problem in machine learning

- We are given a set of (x_i, y_i) examples, where y_i is the class of x_i (y_i is eg. -1 or 1)
- We want a model $f()$, such that for all i , $f(x_i) = y_i$
- $f()$ is very often a parameterized function $f_w()$, and the classification problem becomes an error minimization problem in w .
 - Neural net weights, linear model parameters, etc
- The error is often defined as a sum of errors over the examples

Illustration of classification with a linear model



Stochastic gradient descent

- Assume the error is defined as
- Then the gradient is
- So the full gradient method looks like
- But one can take only one example at a time iterating in random order over examples

$$Err(w) = \sum_{i=1}^n Err(w, x_i)$$

$$\frac{\partial Err(w)}{\partial w} = \sum_{i=1}^n \frac{\partial Err(w, x_i)}{\partial w}$$

$$w(t+1) = w(t) - \alpha(t) \sum_{i=1}^n \frac{\partial Err(w, x_i)}{\partial w}$$

$$w(t+1) = w(t) - \alpha(t) \frac{\partial Err(w, x_i)}{\partial w}$$

Fully distributed classification

- So the problem is to find an optimization method that fits into our system and data model
- Most distributed methods build local models and then combine these through ensemble learning: but we don't have enough local data
- **Online algorithms**
 - Need only one data record at a time
 - They update the model using this record
- The **stochastic gradient** method is a popular online learning algorithm (we apply it to the primal form of the SVM error function)

Gossip learning

Algorithm 1 Gossip Learning Scheme

```

1: initModel()
2: loop
3:   wait( $\Delta$ )
4:    $p \leftarrow \text{selectPeer}()$ 
5:   currentModel  $\leftarrow \text{createModel}()$ 
6:   send currentModel to  $p$ 
7: end loop
8:
9: procedure ONRECEIVEMODEL( $m$ )
10:   modelQueue.add( $m$ )
11: end procedure

```

```

1: procedure CREATEMODELRW
2:    $m \leftarrow \text{modelQueue.first}()$ 
3:   update( $m$ )
4:   return  $m$ 
5: end procedure
6:
7: procedure CREATEMODEL MU
8:    $m_1 \leftarrow \text{modelQueue.first}()$ 
9:    $m_2 \leftarrow \text{modelQueue.second}()$ 
10:   $m \leftarrow \text{merge}(m_1, m_2)$ 
11:  update( $m$ )
12:  return  $m$ 
13: end procedure

```

The merge function

- Let $z = \text{merge}(x, y) = (x + y) / 2$ (x and y are linear models)
- In the case of the Adaline perceptron
 - Updating z using an example has the same effect as updating x and y with the same example and then averaging these two updated models
 - Making predictions using z is the same as calculating the weighted average of the predictions of x and y
- This means we effectively propagate an exponential number of models, and the voting of these is our prediction
- For the linear SVM algorithm this is only a heuristic argument

Local prediction

- We use only local models

- The current model

- Or voting over a number of recent models

```

1: procedure PREDICT( $x$ )
2:    $w \leftarrow \text{currentModel}$ 
3:   return sign( $\langle w, x \rangle$ )
4: end procedure

```

```

5: procedure VOTEDPREDICT( $x$ )
6:   pRatio  $\leftarrow$  0
7:   for  $m \in \text{modelQueue}$  do
8:     if sign( $\langle m.w, x \rangle$ )  $\geq$  0 then
9:       pRatio  $\leftarrow$  pRatio + 1
10:    end if
11:  end for
12:  return sign(pRatio/modelQueue.size() - 0.5)
13: end procedure

```

Experiments

- We implemented a support vector machine with stochastic gradient (Pegasos alg.)
- We used several benchmark data sets for evaluations
 - Data is fully distributed: one data point per node
- We used extreme scenarios
 - 50% message drop rate
 - 1-10 cycles of message delay
 - Churn modeled after the FileList.org trace from Delft

Data sets

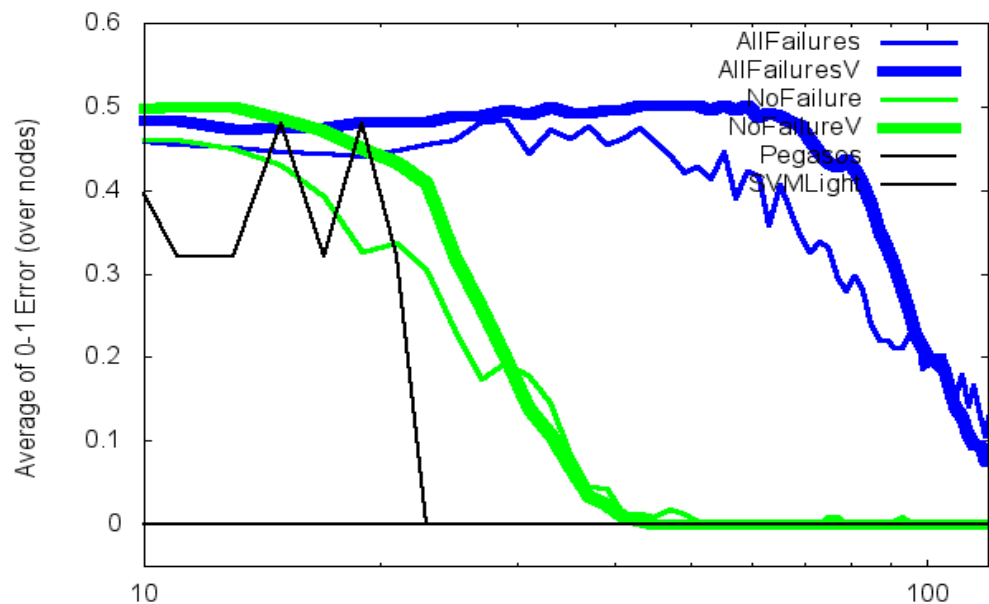
	Iris1	Iris2	Iris3	Reuters	SpamBase	Malicious10
Training set size	90	90	90	2000	4140	2155622
Test set size	10	10	10	600	461	240508
Number of features	4	4	4	9947	57	10
Classlabel ratio	50/50	50/50	50/50	1300/1300	1813/2788	792145/1603985
Pegasos 20000 iter.	0	0	0	0.025	0.111	0.080 (0.081)
Pegasos 1000 iter.	0	0	0.4	0.057	0.137	0.095 (0.060)
SVMLight	0	0	0.1	0.027	0.074	0.056 (–)

- Statistics of data sets
- The performance of some known algorithms

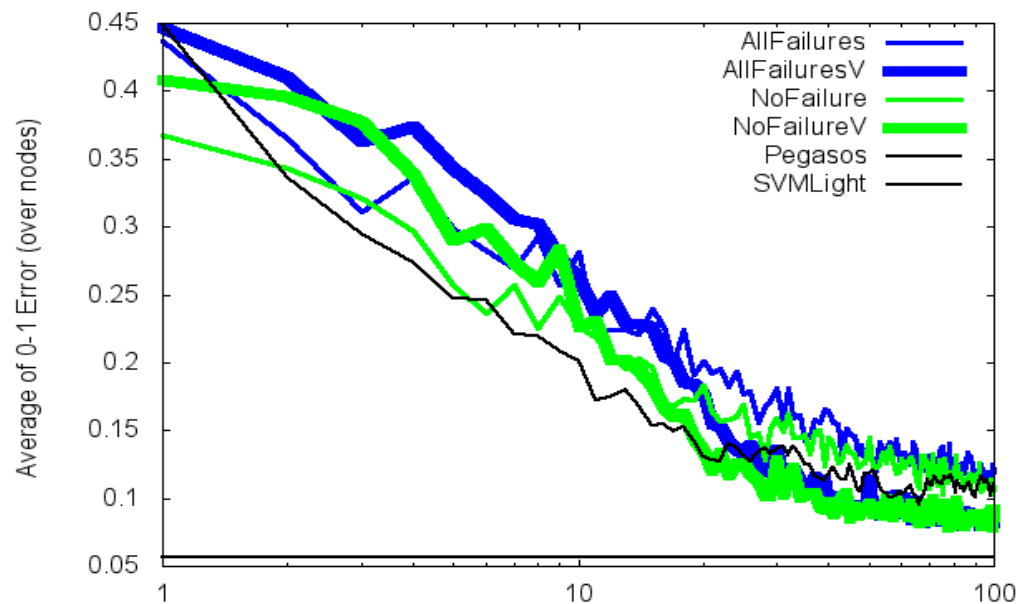
Without merge



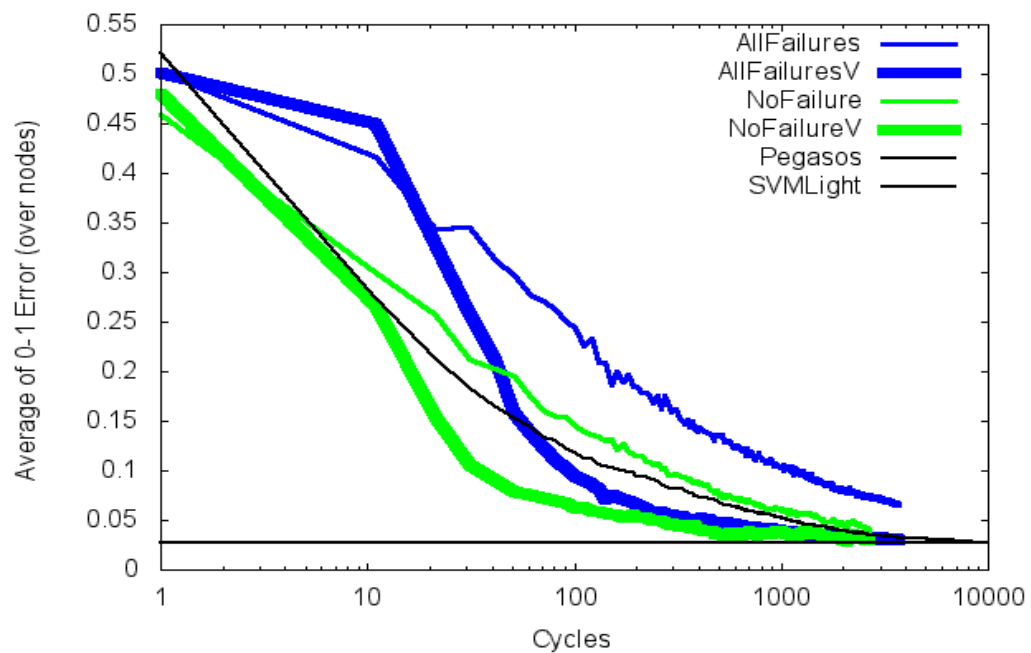
Iris1



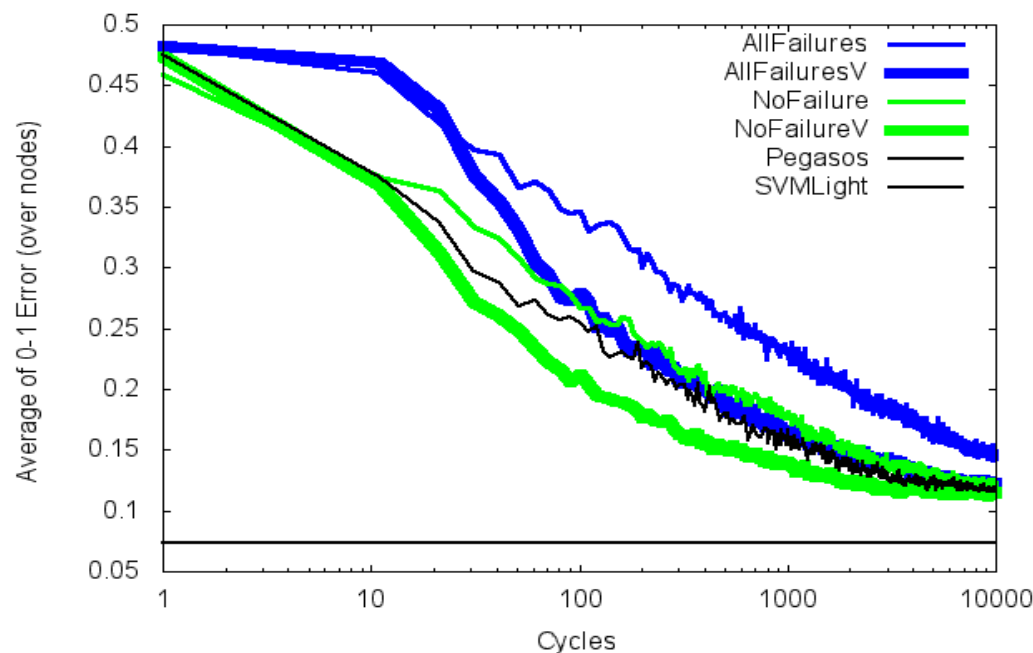
Malicious URLs



Reuters

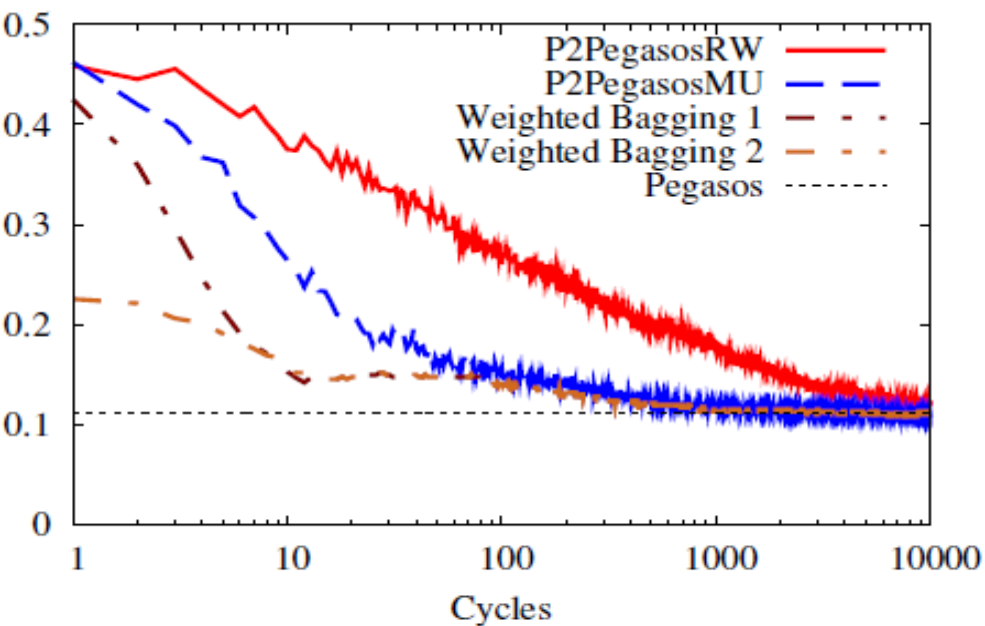


SpamBase

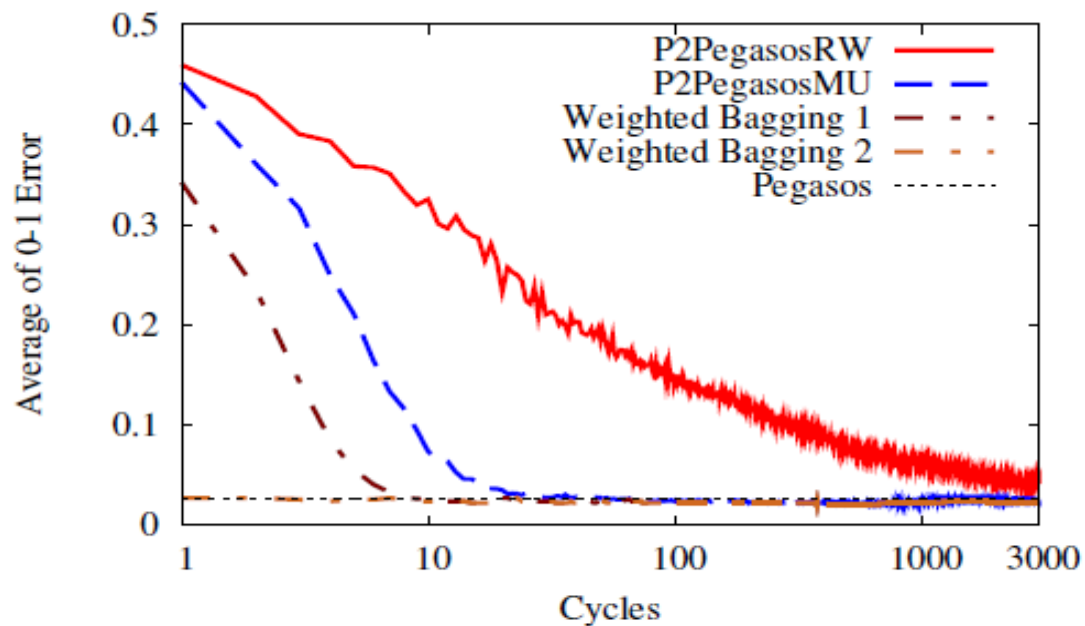


With merge

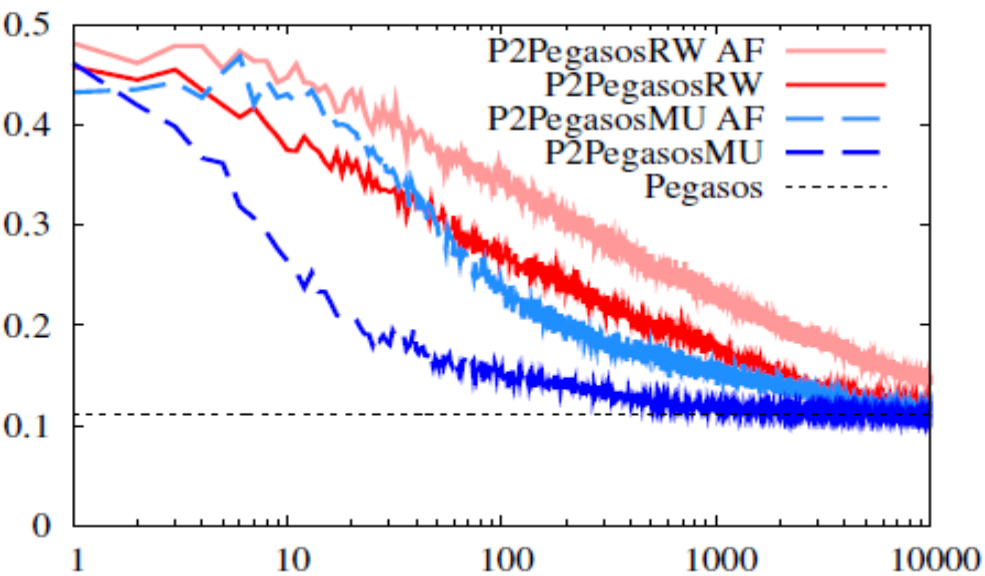
SpamBase No Failure



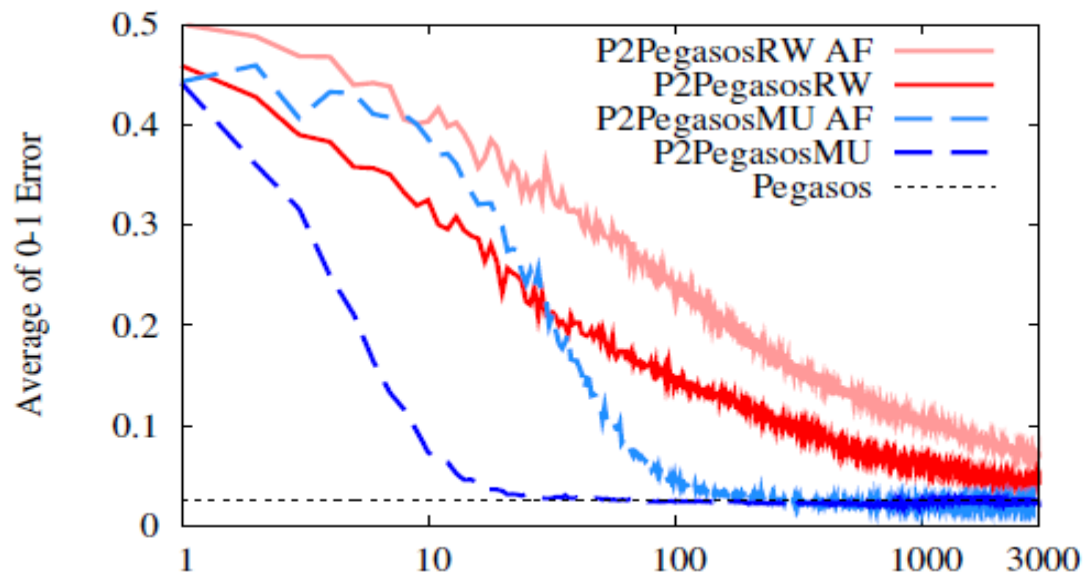
Reuters No Failure



SpamBase with Failures



Reuters with Failures

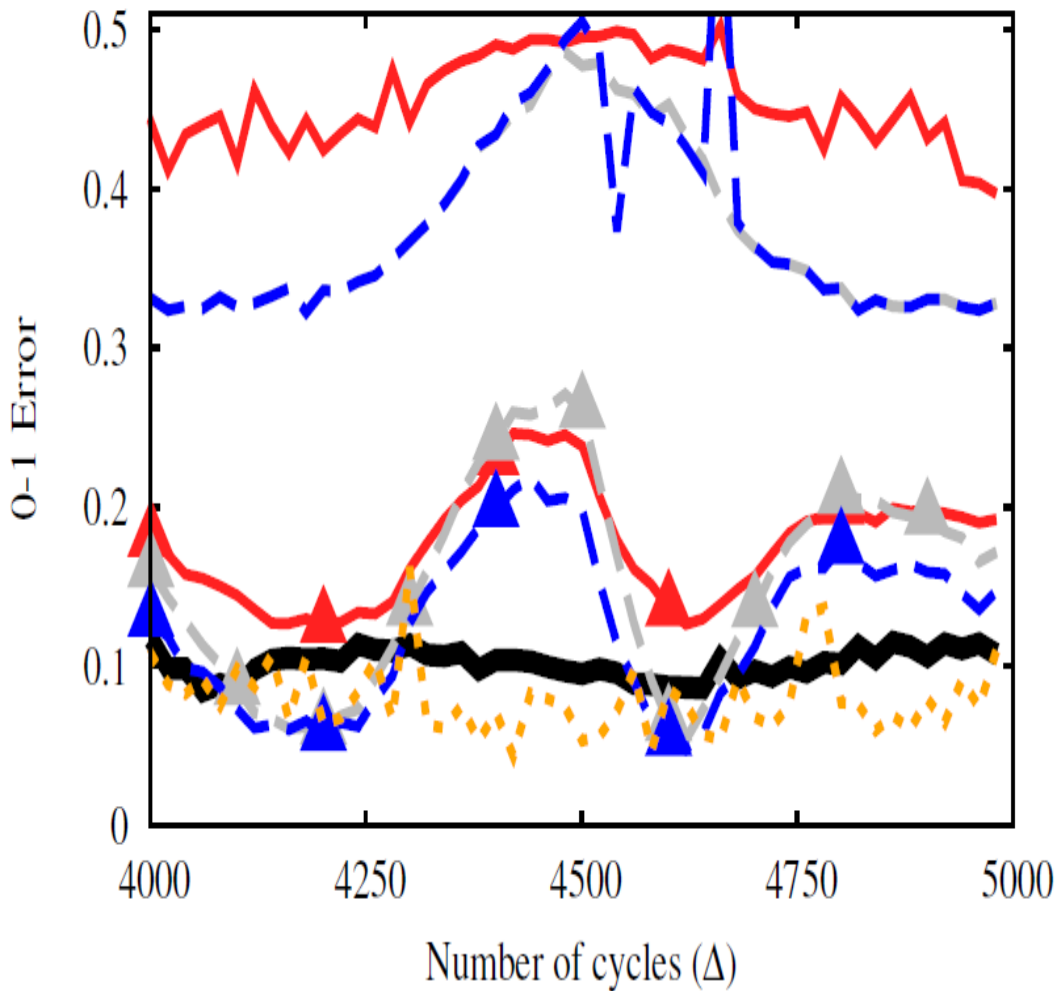


Additional results

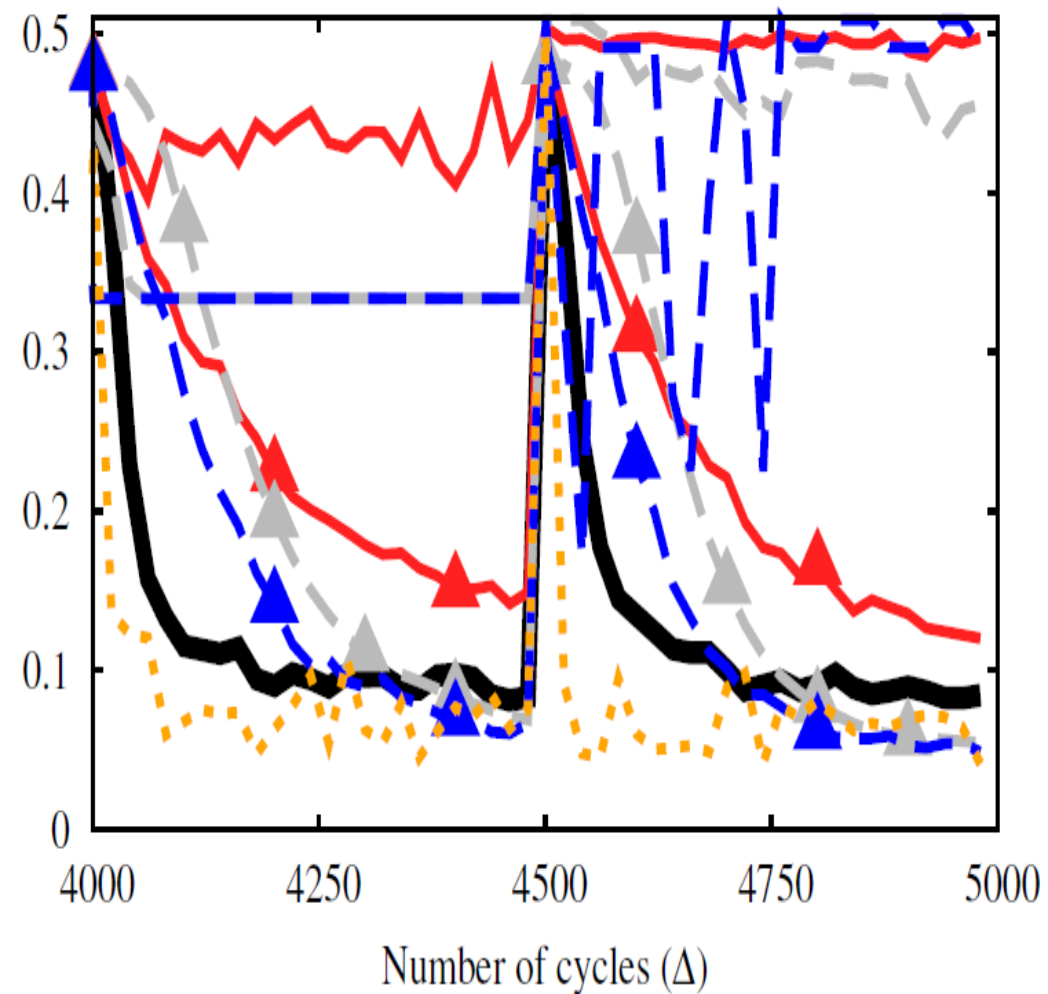
- We implemented multiclass boosting in the gossip framework
- We developed techniques for dealing with concept drift
 - The algorithms is running continuously
 - We keep the age distribution of models fixed
 - At any point in time we have good models

Adaptivity

Sampling rate: 0.1 samples/ Δ



Sampling rate: 0.1 samples/ Δ



Publications

- Róbert Ormándi, István Hegedűs, and Márk Jelasity. **Asynchronous peer-to-peer data mining with stochastic gradient descent**. In Emmanuel Jeannot, Raymond Namyst, and Jean Roman, editors, Euro-Par 2011, volume 6852 of Lecture Notes in Computer Science, pages 528–540. Springer-Verlag, 2011.
- Róbert Ormándi, István Hegedűs, and Márk Jelasity. **Gossip learning with linear models on fully distributed data**. Concurrency and Computation: Practice and Experience, 2012. to appear.
- István Hegedűs, Róbert Busa-Fekete, Róbert Ormándi, Márk Jelasity, and Balázs Kégl. **Peer-to-peer multi-class boosting**. In Euro-Par 2012, Lecture Notes in Computer Science. Springer-Verlag, 2012. to appear.

Remarks regarding the chaotic model

- If uniformity of random walk is guaranteed, then all the models converge to the true model eventually, irrespective of all failures
- If no uniformity can be guaranteed, but the local data is statistically independent of the visiting probability, then again convergence to the true model is guaranteed in general for all models
- If no uniformity and no independence could be guaranteed, convergence to a good model is still ensured provided that the data is separable, and all misclassified examples are visited “often enough”