# The stationary behavior of ideal TCP congestion avoidance

Teunis J. Ott        J.H.B. Kemperman        Matt Mathis

August 22, 1996

## Abstract

This note derives the stationary behavior of idealized TCP congestion avoidance. More specifically, it derives the stationary distribution of the congestion window size if loss of packets are independent events with equal probability. The mathematical derivation uses a fluid flow, continuous time, approximation to the discrete time process $(W_n)$, where $W_n$ is the congestion window after the n-th packet. We derive explicit results for the stationary distribution and all its moments.

Congestion avoidance is the algorithm used by TCP to set its window size (and indirectly its data rate) under moderate to light segment (packet) losses.

The congestion avoidance mechanism we model is idealized in the sense that loss of multiple packets does not lead to time-out phenomena. Such idealized behavior can be implemented using Selective Acknowledgements (SACKs). As such, our model predicts behavior of TCP with SACKs. It also is an approximate model in other situations.

Among the results are that if every packet is lost with a (small) probability p, average window size and long range throughput (when transporting a file large enough to reach stationarity) are of the order of $1/\sqrt{p}$.

# 1. Introduction

In this paper we study the stationary behavior of the congestion window size of an idealized form of the TCP congestion avoidance algorithm. TCP (Transport Control Protocol, see e.g. [2]) is the "Reliable Transport" Protocol in the Internet. Let $W_n$ be the congestion window size immediately after the acknowledgement triggered by segment $n$ has been received[1]. The congestion window size is halved upon receipt of abstract "congestion indicators". While those congestion indicators can have many forms, we conveniently think of them as lost packets (third duplicate acknowledgements), and we assume that packet losses are independent events with small but positive probability. This probabilistic assumption makes it possible to consider $(W_n)_{n=1}^{\infty}$ as a stochastic process, of which we are going to study the stationary distribution.

We define *ideal* TCP congestion avoidance in the spirit of Van Jacobson's landmark paper [1]. When TCP receives a congestion signal (e.g. detects a missing segment) the congestion window is halved. While there is no congestion signal, the congestion window is slowly increased. If the Congestion Window is $W$ Maximal Segment Sizes (MSSs), then the window is opened by $1/W$ for every received TCP acknowledgement.

It is well known that existing versions of TCP do not achieve this ideal behavior. TCP-Reno (with fast retransmit and fast recovery algorithms) [2] comes close, but has problems if two or more segments are lost within one congestion window [7]. Under these conditions TCP-Reno often experiences superfluous retransmission time–outs and other misbehaviors. These behaviors, as well as possible TCP modifications to overcome them are under active investigation [3, 4, 5, 6].

As stated above, we model the situation where the congestion signals are independent

---

[1]we use the term *segment* to refer to TCP data units, and the term *packet* to refer to be the IP transmission units, including data and headers. Typically they are different only by the TCP headers.

from one another. It would be nice to assume that every lost segment signals a new congestion event, but this is somewhat unrealistic in the current Internet, where losses are often highly correlated. Future work will explore several avenues which bring this assumption into agreement with reality. First, the ubiquitous deployment of RED (Random Early Detection) [8] will greatly reduce correlated losses in the Internet. Second, newer TCP implementations [10, 11] tend to treat multiple closely spaced losses as a single congestion signal. Furthermore, other mechanisms besides segment loss can be used to signal congestion [9]. The model and analysis in the paper apply as long as the congestion signals are independent from one another. For ease of discussion we assume that they are indicated by RED style independent segment losses.

To study TCP congestion avoidance without interference from other mandatory TCP algorithms we must assume that TCP is operating in an environment where it is carrying unidirectional bulk data. While all or most of the assumptions below can be relaxed, together they are a good description of what we call the canonical situation.

1. There is always data pending at the sender, such that the sender can always send data as permitted by the congestion window.

2. The receiver's advertised window is sufficiently large to never constrain the congestion window.

3. TCP has previously negotiated the maximum segment size, (MSS), and all segments are exactly one MSS.

4. Every segment is individually acknowledged (the *delayed acknowledgement* algorithm is not in effect). Under delayed acknowledgement one acknowledgement often acknowledges two or more data segments. When we investigate delayed acknowledgement we must decide how much the window increases for such acknowledgements.

5. We assume that recovery from lost segments does not interact with congestion avoidance. This is part of saying that we consider an *ideal* version of TCP. Achieving this one of the design goals of future TCP implementations such as TCP-SACK.

Throughout this paper we assume that segment losses are independent (from segment to segment) events with all the same, positive but small, probability $p$.

Let $W_n$ be the congestion window size right after the acknowledgement triggered by segment $n$ has been received. Our assumptions make it possible to consider $(W_n)_{n=1}^{\infty}$ as a stochastic process, of which we are going to study the stationary distribution.

The first and second assumptions imply that $W_n$ can grow quite large. The third assumption means the congestion window is expressed in units of Maximal Segment Sizes (MSSs) segments, as opposed to, for example, bytes. Henceforth, we allow $W_n$ to take on all values in $(0, \infty)$.

This means $(W_n)$ evolves as follows:

Given $W_n$, with probability (1-p) $W_{n+1}$ equals $W_n + (1/W_n)$, and with probability p it equals $W_n/2$. (The first part assumes the Congestion Window $W$ is expressed in units equal to the Maximal Segments Size, and that all packets are equal to the Maximal Segment Size).

Since p is positive but small, we can re-scale time so that the epochs of packet failure become a Poisson Process. Apart from the re–scaling, time is still expressed in number of packets acked. We also rescale space to keep $W(.)$ from going to infinity. Now, we have a continuous time stochastic process $(W(t))_{0<t<\infty}$ with state–space $(0, \infty)$. More explicitly, we define $W(t)$ by

$$W(t) = \sqrt{p}\, W_{\lfloor \frac{t}{p} \rfloor}, \tag{1.1}$$

(where $\lfloor \; \rfloor$ denotes floor or integer part). For p small, the process $W(t)$ approximately behaves as follows: There is a Poisson Process of intensity 1. In between the points of

4

the Poisson Process $W(.)$ behaves as

$$(d/dt)W(t) = \frac{1}{W(t)}, \tag{1.2}$$

and in the points $\tau_k$ of the Poisson Process the behavior is that

$$W(\tau_k+) = \frac{1}{2}W(\tau_k-) \text{ (limit from right is half limit from left)}. \tag{1.3}$$

An interesting question is how small $p$ must be for this approximation to be numerically close.

For several reasons we allow the Poisson Process to have an intensity $c_1$ which need not be equal to one, and we slightly modify (1.2). The evolution of $W(.)$ now is as follows: There is a Poisson Process with intensity $c_1$, with points $(\tau_k)_{k=1}^\infty$. In between the points $\tau_k$ the process $W(.)$ behaves as

$$(d/dt)W(t) = \frac{c_2}{W(t)}, \tag{1.4}$$

and in the points $\tau_k$ the behavior still is as in (1.3).

There are several reasons it is convenient to allow $c_1$ and $c_2$ to be different from one. In TCP with delayed acknowledgements, if indeed every acknowledgement were to acknowledge 2 Maximal Segment Sizes (MSSs), and for each such acknowledgement the congestion window were to increase by only $1/W_n$, $c_2$ would be equal to $1/2$. If for each such aknowledgement acking two MSSs the congestion window were to increase by $2/W_n$ the value of $c_2$ would be back at 1. Similar issues arise with acks acking small packets. If we want to express "time" in other units than packets or MSSs (for example bytes), $c_1$ gets a different value. This can also occur if the packet lengths are not all equal to one MSS. Lastly, we may decide to have a space and time rescaling different from (1.1).

Formally, the continuous time process $(W(t))_{t=0}^\infty$ as in (1.3) and (1.4) is a good aproximation for the discrete time process $(W_n)_{n=0}^\infty$ as long as after an acknowledgement packet

5

acking *Acked* MSSs, the Window W (MSSs) increases by $c_2 Acked/W$ MSSs, and that in a sequence of packets containing (together) N MSSs (N need not be integer), the probability of exactly one halving of the window is

$$c_1 N p + o(N p) \text{ if } N p \text{ is small,} \tag{1.5}$$

and the probability of two or more halvings is an order of magnitude smaller.

It is clear, and will be made formal, that only the ratio of $c_1$ and $c_2$ matters. We define

$$\eta = \frac{c_1}{c_2}. \tag{1.6}$$

The stationary distribution of $W(.)$ must have most of its probability in the general area of $\sqrt{\frac{2}{\eta}}$. Namely, let $D(w)$ be the drift of $W(t)$ if $W(t) = w$. Then

$$D(w) = \frac{c_2}{w} - \frac{c_1 w}{2}, \tag{1.7}$$

and

$$D(w) = 0 \text{ if } w^2 = \frac{2}{\eta}. \tag{1.8}$$

For smaller values of $w$ the drift is positive, and for larger values of $w$ the drift is negative. This gives the ad-hoc approximation

$$\overline{W} = \sqrt{\frac{2}{\eta}}. \tag{1.9}$$

In the remainder of this paper we will refine this observation. (1.9) is an ad-hoc approximation for the mean or average value of W in the process $W(.)$ as defined by (1.3), (1.4). Because of the definition (1.1), $\overline{W}$ therefore (apart from the factor $\sqrt{p}$) is an ad hoc approximation for the "acknowledgement–arrival" average of the window size in the original process $W_n$ of interest. In section 3 we will see that for certain purposes we need the "time" average of the original. This is the average obtained by integrating over the clock time, as opposed to by summing over acknowledgement arrivals. The difference

6

between "acknowledgement–arrival" and "time" is essentially the same as the difference between "customer–arrival" averages and "time" averages in classical queueing theory. Since acknowledgements arrive at a higher rate when the congestion window is large, the "acknowledgement–arrival" average value of the window size is larger than the "time" average value. We will see that in fact (1.9) gives a value in between the time-average value and the acknowledgement–arrival average.

In this context, it must be remembered that in all circumstances for the process $W(t)$ *time* is *not* clocktime, but (some rescaled version of) total number of packets (or bytes, or MSSs) acknowledged. The relationship between time as used in this paper (a rescaled version of number of MSSs acknowledged) and clocktime is complicated and involves the Round Trip Time (RTT) of the TCP connection. We will delay a discussion until we use the results from this paper to discuss the relationship between the drop probability $p$ and the long–run throughput of the TCP connection.

In section 4 we will introduce a slightly larger class of stochastic processes than those as in (1.3) – (1.4), and in the sections 4 and 5 we will explicitly derive the stationary distributions and moments of this slightly larger class of processes.

Results specific to the Window Behavior that follow from the sections 4 and 5 will be quoted without proof in section 2.

In section 3 we will give a simple example of how the material developed in this paper can be used to predict the throughput of a TCP connection.

More details and related information about mathematical aspects of the process $W(.)$ is given in the sections 6 – 8. In section 9 we discuss the relationship between the class of stochastic processes studied in the sections 4 – 8 of this note, and a process studied by Ferguson in 1972 [13]. Ferguson studies a discrete time markov process which either decreases by one (with probability $\pi$) or doubles (with probability $1 - \pi$). This

process terminates if it reaches $(-\infty, 0]$ (ruin of the gambler). There is an interesting relation between the probability of eventual ruin in Ferguson's process and the stationary distribution of processes related to the stationary markov process $W(t)$ discussed above.

## 2. The stationary distribution of the congestion window size

Let $W$ be a random variable with as distribution the stationary distribution of the process $W(.)$ in (1.3), (1.4). Specializing the results in the sections 4 and 5 for this situation, we have:

$$\frac{\eta}{2} E[W^2] = \frac{1}{1 - \frac{1}{4}} = \frac{4}{3}, \quad E[W^2] = \frac{8}{3\eta}. \tag{2.1}$$

Define:

$$A = \prod_{j=1}^{\infty} \left( \frac{1 - 2^{-(2j+1)}}{1 - 2^{-2j}} \right) = \frac{\sum_{k=0}^{\infty} \frac{(-1)^k 2^{-k(k+2)}}{(1-2^{-2})(1-2^{-4})\dots(1-2^{-2k})}}{\sum_{k=0}^{\infty} \frac{(-1)^k 2^{-k(k+1)}}{(1-2^{-2})(1-2^{-4})\dots(1-2^{-2k})}} = 2 \left( \sum_{k=0}^{\infty} 2^{-\frac{1}{2}k(k+1)} \right)^{-1}. \tag{2.2}$$

Then:

$$\left( \frac{\eta}{2} \right)^{\frac{1}{2}} E[W] = \frac{A}{2} \sqrt{\pi}, \quad E[W] = A \sqrt{\frac{\pi}{2\eta}}. \tag{2.3}$$

Hence:

$$Var(W) = \frac{1}{\eta} \left( \frac{8}{3} - \frac{\pi A^2}{2} \right). \tag{2.4}$$

Also,

$$P\{W > w\} = \sum_{k=0}^{\infty} R_k(\frac{1}{4}) exp\{-2^{2k-1}\eta w^2\}, \tag{2.5}$$

where for $|c| < 1$:

$$R_k(c) = \frac{1}{L(c)} \frac{(-1)^k c^{-\frac{1}{2}k(k+1)}}{(1-c)(1-c^2)\dots(1-c^k)}, \tag{2.6}$$

$$L(c) = \prod_{k=1}^{\infty} (1 - c^k) = \sum_{k=0}^{\infty} \frac{(-1)^k c^{\frac{1}{2}k(k+1)}}{(1-c)(1-c^2)\dots(1-c^k)} = \sum_{k=-\infty}^{+\infty} (-1)^k c^{\frac{1}{2}k(3k+1)}. \tag{2.7}$$

8

For $A$ in (2.2) and for $L(c)$ in (2.7) three expressions are given. Section 4 will prove, among other things, that those expressions are in indeed equivalent. For numerical purposes it seems the rightmost infinite sums are preferable over the middle infinite sums, which are preferable over the infinite products.

In addition to the results above we have

$$\left(\frac{2}{\eta}\right)^{\frac{1}{2}} E[\frac{1}{W}] = \frac{A}{2}\sqrt{\pi}, \quad E[\frac{1}{W}] = \frac{A}{2}\sqrt{\frac{\pi\eta}{2}}, \tag{2.8}$$

so that indeed the expected value of the drift is zero:

$$E[\frac{c_2}{W} - \frac{c_1}{2}W] = c_2 A\sqrt{\pi\eta}\left(\sqrt{\frac{1}{8}} - \sqrt{\frac{1}{8}}\right) = 0. \tag{2.9}$$

In section 3 we will see that (under a common sense condition for source behavior)

$$\frac{1}{E[1/W]} = \frac{2}{A}\sqrt{\frac{2}{\pi\eta}}. \tag{2.10}$$

is the "time" average value of the process $W(.)$.

For numerical values of coefficients of interest we find (using the sums up to $k = 6$ only (!)):

$$A = 1.218\,299\,422\,132\,457, \quad L(\frac{1}{4}) = 0.688\,537\,537\,120\,339,$$

$$A\sqrt{\frac{\pi}{2}} = 1.526\,911\,889\,241\,912, \quad \frac{8}{3} - \frac{\pi A^2}{2} = 0.335\,206\,749\,158\,359,$$

$$\frac{st.dev(W)}{E[W]} = 0.379\,177\,361\,667\,205,$$

and finally

$$E[W].E[\frac{1}{W}] = 1.165\,729\,958\,754\,153.$$

These results are consistent with the statement that "much of the time" the process $W(.)$ remains between $E[W](1 - .379)$ and $E[W](1 + .379)$, which in its turn is consistent with the fact that

$$\frac{1.379}{0.621} = 2.222 > 2.$$

9

Combining (2.3), (2.8) with some of the results above, and comparing this with (1.9), we see that

$$\frac{1}{E[1/W]} = \frac{1.309\,833}{\sqrt{\eta}}, \ \overline{W} = \frac{1.414\,214}{\sqrt{\eta}}, \ E[W] = \frac{1.526\,912}{\sqrt{\eta}}, \ \text{and}$$

$$\frac{1.309\,833}{1.414\,214} = .926\,192, \ \frac{1.526\,912}{1.414\,214} = 1.079\,690 \ .$$

$\overline{W}$ lies almost exactly in the middle between the "time" average and the "acknowledgement–arrival" average, and has distance about 7.6% to both of these.

## 3.  Applying our results in the Internet

Consider two computers communicating over a very high bandwidth network, unidirectional (data A to B only, acks B to A only) with probability of loss $p$ per MSS, and a Round Trip Time $RTT$. We want to use the material developed in the other sections of this paper to predict the throughput of the connection. This section is the only one on this paper where the assumption that the source always has material to send comes into play: The amount of data sent but not yet acknowledged is assumed to always equal one congestion window. Hence, we do not consider the possibility that the source keeps part of the congestion window as an un–excercised "right–to–send". The result below also requires that there is no dependence between number of MSSs acknowledged and instantaneous window size.

For a simple example, assume the Round Trip Time is strictly equal to the sum of propagation delays (no queueing component, hence independent of the congestion window $W^\star$. In that case at any point in time the throughput (rate) is proportional to the congestion window:

$$\text{Throughput (rate)} = \frac{W^\star}{RTT}. \tag{3.1}$$

Namely, in one RTT one whole congestion window $W^\star$ will be acknowledged. Over longer periods the window size $W^\star$ varies with time. Let $W_t^\star$ be the window size at *clock* time $t$. This is the *clock* time, not the time used in (1.1) etc. Also, $W^\star$ is the real window size, not the rescaled window size used in (1.1) etc. The total throughput over the period $[0, n.RTT]$ is (by close approximation)

$$W_0^\star + W_{RTT}^\star + \ldots + W_{(n-1)RTT}^\star. \tag{3.2}$$

More generally, the total throughput (in MSSs) over the period $[0, T]$ will be

$$\left( \int_0^T W_t^\star dt \right) / RTT. \tag{3.3}$$

If we now accept that (2.10) indeed gives the time–average, we get

$$\text{Long Run Average Throughput} = \frac{2}{A.RTT} \sqrt{\frac{2}{\pi \eta p}} = \frac{1.309\,833}{RTT \sqrt{p}} \text{ MSSs/sec.} \tag{3.4}$$

It is interesting to compare (3.4), the long run average throughput in case of independent random packet loss, with the long run average throughput in case of deterministic, perfectly periodic packet loss. This means that for the process $W(.)$ in (1.1) etc. the process is halved at exactly the epochs $c_1^{-1}, 2c_1^{-1}, \ldots$. Instead of (2.3), (2.10) we now get

$$E[W] = \frac{2}{3} \frac{4 - \sqrt{2}}{\sqrt{\eta}} = \frac{1.723\,858}{\sqrt{\eta}},$$

$$\frac{1}{E[\frac{1}{W}]} = \frac{1}{(2 - \sqrt{2})\sqrt{\eta}} = \frac{1.707\,107}{\sqrt{\eta}}.$$

For the long run average throughput this gives

$$\text{Long Run Average Throughput} = \frac{1.707\,107}{RTT \sqrt{p}} \text{ MSSs/sec.} \tag{3.5}$$

(3.4), (3.5) give an indication of the sensitivity to the assumption of independence.

11

Next, we prove the relationship between time average and acknowledgement arrival average that we stated and already used. Let $A(.)$ be the probability distribution function of the stationary distribution of window size at acknowledgement arrival epochs, and let $B(.)$ be the similar stationary distribution function for the continuous time process. In $[t, t + RTT]$ there are $W_t^\star$ acknowledgements. Hence:

$$dA(x) = \frac{x dB(x)}{\int_0^\infty y dB(y)}. \tag{3.6}$$

(2.10) etc. now follows readily by dividing both sides of (3.6) by $x$, and integrating from zero to infinity.

The result above assumes all packets equal to one MSS and all data packets individually acknowledged. If there are delayed acks or smaller packets, there may be need for a factor $\sqrt{c_2}$. The actual factor must be obtained by a careful study of actual packet sizes, number of packets acked per acknowledgement, and window evolution upon arrival of an ack packet.

If part of the RTT is caused by congestion in a bottleneck link, the RTT becomes dependent on the instantaneous value of the time varying window size, and predicting throughput becomes more complicated. If there are other connections competing for bandwidth in that bottleneck link, the relationship becomes even more complicated. In that case we also need to know the RTTs of the other connections, and we must predict the throughputs of individual connections.

Random Early Detection (RED) will cause packet loss only if there is congestion in the bottleneck. In that situation, the packet loss probability becomes dependent on the buffer occupancy, i.e. on the level of congestion, in the bottleneck, and the (instantaneous) dropping probability becomes dependent on the (instantaneous) sizes of the congestion windows of the competing connections. This is one of the problems we will tackle next.

12

# 4. The Central Result

In this section we study the Stochastic Process $(W(t))_{0 < t < \infty}$ on $(0, \infty)$ which evolves as follows:

There is a Poisson Process of intensity $c_1 > 0$, with points $(\tau_k)_{k=-\infty}^{+\infty}$, with

$$\ldots < \tau_{-2} < \tau_{-1} < 0 < \tau_0 < \tau_1 < \ldots \qquad (4.1)$$

Between the points $\tau_k$ the process $W(.)$ evolves as

$$(d/dt)W(t) = \frac{c_2}{W(t)^m}, \qquad (4.2)$$

and in the points $\tau_k$ we have

$$W(\tau_k+) = \beta.W(\tau_k-), \qquad (4.3)$$

where

$$0 < \beta < 1, \ m > -1, \ c_1 > 0, \ \text{and } c_2 > 0. \qquad (4.4)$$

If $\beta = 1/2$ and $m = 1$ we are back in the situation of the sections 1 and 2. The results in the remainder of this section, and in the next section, indicate it is sufficient to require $m > -1$, see also the argument in section 1 based on (1.7), (1.8).

We define:

$$\eta = \frac{c_1}{c_2}, \ c = \beta^{m+1}, \qquad (4.5)$$

and we use the substitution

$$Z(t) = \frac{\eta}{m+1} W(t)^{m+1}, \ \text{where } \eta \text{ is as in (4.5).} \qquad (4.6)$$

It is easily verified that between the points $\tau_k$

$$(d/dt)Z(t) = c_1, \qquad (4.7)$$

13

and that in the points $\tau_k$

$$Z(\tau_k+) = cZ(\tau_k-), \text{ where } c \text{ is as in (4.5)}. \tag{4.8}$$

We thus see that all processes $W(.)$ as above can be reduced to the special case with

$$c_2 = c_1 \text{ and } m = 0. \tag{4.9}$$

The central result in this paper is:

Theorem 1. If the random variable $Z$ has the stationary distribution of the process $Z(.)$, then Z has the form

$$Z = \sum_{k=0}^{\infty} c^k E_k, \text{ where} \tag{4.10}$$

$(E_k)_{k=0}^{\infty}$ are independent, identically distributed random variables, all exponentially distributed with expected value 1.

Proof: Denote

$$Z(0) = z, \tag{4.11}$$

We define:

$$H_0 = \tau_0, \ H_k = \tau_k - \tau_{k-1} \text{ for } k \geq 1. \tag{4.12}$$

Clearly, $(H_k)_{k=0}^{\infty}$ are independent, identically distributed, all exponentially distributed with expected value $c_1^{-1}$. Also clearly,

$$Z(t) = z + c_1 t \text{ if } 0 \leq t < \tau_0, \text{ while} \tag{4.13}$$

$$Z(t) = c^{k+1} z + c_1 \left( c^{k+1} H_0 + c^k H_1 + \ldots + c H_k + (t - \tau_k) \right) \text{ if } \tau_k < t < \tau_{k+1}. \tag{4.14}$$

Using (4.14) to "look backward" from time $t$, we see that if

$$\text{if } \tau_j < t < \tau_{j+1} \text{ then} \tag{4.15}$$

$$Z(t) = \sum_{k=0}^{\infty} c^k E_k^{(t)}, \text{ where} \tag{4.16}$$

14

$$E_0^{(t)} = c_1(t - \tau_j), \ E_k^{(t)} = c_1(\tau_{j-k+1} - \tau_{j-k}) \text{ for } k \geq 1. \tag{4.17}$$

This proves the theorem. It is easily seen that the process $(Z(t)_{t=-\infty}^{+\infty}$ defined by (4.15) – (4.17) is stationary and is as in (4.7), (4.8). The proof also gives information about the relaxation time of the system: Very roughly speaking, in an amount of time $T$ the system looses a fraction $1 - c^{T/c_1}$ of its memory.

A corollary of theorem 1 is that if Z and E are independent and E has the exponential distribution with expected value 1, then

$$Z \text{ and } E + cZ \text{ have the same distribution.} \tag{4.18}$$

Henceforth,

$$F(w, t) = P\{W(t) \leq w\} \tag{4.19}$$

for the "general" process $W(.)$ as in (4.2) – (4.4), and

$$G(z, t) = P\{Z(t) \leq z\} \tag{4.20}$$

for the "special" process $Z(.)$ as in (4.6) –(4.8), and $F(.)$ and $G(.)$ are the corresponding stationary distribution functions. $W$ and $Z$ are random variables with distribution functions $F$ and $G$. Without loss of generality we assume that

$$Z = \frac{\eta}{m+1} W^{m+1}. \tag{4.21}$$

Hence:

$$F(w) = P\{W \leq w\} = P\{Z \leq \frac{\eta}{m+1} w^{m+1}\} = G(\frac{\eta}{m+1} w^{m+1}). \tag{4.22}$$

$G(.)$ is determined by (4.10), and $F(.)$ is determined by (4.21). In the next section we will discuss the distribution function $G(.)$, its Laplace Transform, and its Moments. In the appendices A and B we will give two alternative proofs of theorem 1.

15

# 5.   Transforms, Distribution Functions and Moments

Let $\phi_Z(.)$ denote the Laplace transform of the random variable Z in theorem 1. Clearly,

$$\phi_Z(s) = \prod_{k=0}^{\infty} \frac{1}{1 + c^k s}. \tag{5.1}$$

It is clear that the infinite product in the RHS of (5.1) converges for all $|c| < 1$ (unless $s$ is of the form $s = -c^{-k}$). At this point it is useful to quote four classical results, three of which are due to Euler. The most important are the first two, both due to Euler: For $|c| < 1$,

$$\prod_{k=0}^{\infty} \frac{1}{1 + c^k z} = \sum_{k=0}^{\infty} \frac{(-1)^k z^k}{(1 - c)(1 - c^2) \ldots (1 - c^k)}, \tag{5.2}$$

and

$$\prod_{k=0}^{\infty} (1 + c^k z) = \sum_{k=0}^{\infty} \frac{c^{\frac{1}{2}(k-1)k} z^k}{(1 - c)(1 - c^2) \ldots (1 - c^k)}, \tag{5.3}$$

see Gasper and Rahman (1990), for example formulas (18) and (19) on page xiv, or the results on page 8, for more information on these and similar results. We give simple direct proofs of (5.2) and (5.3) in section 6 of this note. (5.2) gives the moments of the random variable Z. Later in this section we will give other expressions for those moments. (5.3) proves the "left" halves of (2.2), (2.7). The third result, also due to Euler, is (still $|c| < 1$):

$$L(c) = \prod_{k=1}^{\infty} (1 - c^k) = \phi_Z(-c) = \sum_{k=-\infty}^{+\infty} (-1)^k c^{\frac{1}{2}k(3k+1)} = 1 - c - c^2 + c^5 + c^7 - c^{12} - c^{15} + \ldots, \tag{5.4}$$

of which we do not give a proof (see e.g. Hurwitz and Courant (1964) page 207). In a slightly different form, (5.4) is also given in Gasper and Rahman, page xiv formula (15). For numerical purposes, (5.4) converges even faster than (5.3) with $z = -c$. (5.4) proves the rightmost equality in (2.7). The fourth result, which can be useful in evaluating expressions like (2.2), is (still $|c| < 1$):

$$\prod_{k=1}^{\infty} \frac{1 - c^{2k}}{1 - c^{2k-1}} = \sum_{k=0}^{\infty} c^{\frac{1}{2}k(k+1)} = 1 + c + c^3 + c^6 + c^{10} + \ldots. \tag{5.5}$$

16

(5.5) can be obtained from Jacobi's triple product identity, see Gasper and Rahman (1990) page 12. (5.5) proves the rightmost equality in (2.2). (5.1) can be re-written as a partial fraction expansion:

$$\phi_Z(s) = \prod_{k=0}^{\infty} \frac{1}{1+c^k s} = \sum_{k=0}^{\infty} \frac{R_k(c)}{1+c^k s}, \tag{5.6}$$

where the "Residues" $R_k(c)$ are given by

$$R_k(c) = \lim_{s \to -c^{-k}} (1+c^k s)\phi_Z(s) = \frac{1}{L(c)} \frac{(-1)^k c^{\frac{1}{2}k(k+1)}}{(1-c)(1-c^2)\dots(1-c^k)}, \tag{5.7}$$

where $L(c)$ is as in (5.4). (5.6) is not quite as obvious an identity as seems at first glance: to prove it we observe that the difference between the RHS and the LHS is analytic on the whole complex plane, and also is bounded on the whole complex plane. Liouville's theorem then shows that that difference is a constant function. By letting $s \to \infty$ we see that the constant is zero. We now have four different expressions for $\phi_Z(.)$:

$$\phi_Z(s) = \prod_{k=0}^{\infty} \frac{1}{1+c^k s} = \sum_{k=0}^{\infty} \frac{(-1)^k s^k}{(1-c)(1-c^2)\dots(1-c^k)} = \sum_{k=0}^{\infty} \frac{R_k(c)}{1+c^k s} =$$

$$\left( \sum_{k=0}^{\infty} \frac{c^{\frac{1}{2}(k-1)k} s^k}{(1-c)(1-c^2)\dots(1-c^k)} \right)^{-1}. \tag{5.8}$$

The infinite sum in the RHS of (5.6) shows that

$$P\{Z > x\} = \sum_{k=0}^{\infty} R_k(c) exp\{-c^{-k} x\}. \tag{5.9}$$

For $x \downarrow 0$ this probability must go to one: this is indeed guaranteed by (5.6) with $s = 0$, and also by (5.3) with $z = -c$. (5.9) also gives us the probability distribution of the original random variable W:

$$P\{W > w\} = P\{\frac{\eta W^{m+1}}{m+1} > \frac{\eta w^{m+1}}{m+1}\} = \sum_{k=0}^{\infty} R_k(c) exp\{-\frac{c^{-k}\eta}{m+1} w^{m+1}\}, \tag{5.10}$$

and the density is given by

$$(d/dw)P\{W \le w\} = \sum_{k=0}^{\infty} a_k(c)\eta w^m exp\{-\frac{c^{-k}\eta}{m+1} w^{m+1}\}, \tag{5.11}$$

17

where

$$a_k(c) = R_k(c)c^{-k} = \frac{1}{L(c)} \frac{(-1)^k c^{\frac{1}{2}(k-1)k}}{(1-c)(1-c^2)\dots(1-c^k)}. \tag{5.12}$$

We are now ready to obtain the moments of $W$ and $Z$ (with distribution functions $F$ and $G$). We define:

$$\mu(\alpha) = E[W^\alpha], \ \nu(\alpha) = E[Z^\alpha], \tag{5.13}$$

so that

$$\nu(\alpha) = \left(\frac{\eta}{m+1}\right)^\alpha \mu(\alpha(m+1)). \tag{5.14}$$

(5.1) and (5.2) immediately show that for $k \in \{0, 1, 2, \dots\}$

$$\nu(k) = \frac{k!}{(1-c)(1-c^2)\dots(1-c^k)}. \tag{5.15}$$

On the other hand, (5.9) shows that for all (non-negative) $\alpha$

$$\nu(\alpha) = \Gamma(\alpha+1) \sum_{k=0}^\infty R_k(c)c^{k\alpha}. \tag{5.16}$$

Using (5.3) with $z = -c^{\alpha+1}$ we see that (5.16) implies that in fact

$$\nu(\alpha) = \Gamma(\alpha+1) \prod_{k=1}^\infty \left(\frac{1-c^{\alpha+k}}{1-c^k}\right) = \Gamma(\alpha+1) \frac{\phi_Z(-c)}{\phi_Z(-c^{\alpha+1})}. \tag{5.17}$$

(5.17) is a generalization of (5.15) to values of $\alpha$ outside $\{0, 1, \dots\}$. For $\alpha \in \{0, 1, \dots\}$ (5.21) is an alternative proof of (5.15). More interestingly: the middle expression in (5.16) is a meromorphic function of $\alpha$, so that $\nu(\alpha)$ is finite for all real values of $\alpha$. (The simple poles of $\Gamma(\alpha+1)$ are eliminated by the simple zeros of $\prod(1-c^{\alpha+k})$). Using a result from Whittaker and Watson (section (12.14)) we get, for $\alpha > 0$:

$$\nu(-\alpha) = \frac{\pi}{\Gamma(\alpha)\sin(\pi\alpha)} \prod_{j=1}^\infty \left(\frac{1-c^{-\alpha+j}}{1-c^j}\right). \tag{5.18}$$

Both (5.17) and (5.18) hold for all real $\alpha$. In particular, for $k \in \{1, 2, \dots\}$ we have

$$\nu(-k) = \frac{(1-c)(1-c^2)\dots(1-c^{k-1})}{(k-1)!c^{\frac{1}{2}(k-1)k}} \log(\frac{1}{c}) < \infty, \tag{5.19}$$

18

Since $\nu(\alpha) < \infty$ for all real $\alpha$ we know that for $z \downarrow 0$, $G(z) \downarrow 0$ faster than $z^{\alpha}$ for every real $\alpha$. In section 8 we will refine this result.

The results in section 2 are special cases of the results proven in this section.

# 6.   The proofs of (5.2) and (5.3)

This section provides direct proofs of (5.2) and (5.3). Define

$$q_k(c) = \frac{(-1)^k}{(1-c)(1-c^2)\dots(1-c^k)}, \quad f(z) = \sum_{k=0}^{\infty} q_k(c) z^k, \tag{6.1}$$

$$r_k(c) = \frac{c^{\frac{1}{2}(k-1)k}}{(1-c)(1-c^2)\dots(1-c^k)}, \quad g(z) = \sum_{k=0}^{\infty} r_k(c) z^k. \tag{6.2}$$

We see that

$$(1 - c^k) q_k(c) = -q_{k-1}(c). \tag{6.3}$$

Multiplying this with $z^k$ and summing over $k$ yields

$$f(z) = \frac{1}{1+z} f(cz) = \frac{1}{(1+z)(1+cz)} f(c^2 z) = \dots , \tag{6.4}$$

which (with $f(0) = 1$) proves (5.2). Similarly, we have

$$(1 - c^k) r_k(c) = c^{k-1} r_{k-1}(c). \tag{6.5}$$

Multiplying this with $z^k$ and summing over $k$ yields

$$g(z) = (1+z) g(cz) = (1+z)(1+cz) g(c^2 z) = \dots . \tag{6.6}$$

This completes the direct proof of (5.3).

19

# 7. $P\{Z \le x\}$ for x small

For small window sizes W the assumption made in section 1, that W can vary continuously, could become a problem. In this section we will show that once stationarity has been reached, the probability that W is small is quite negligible. We will derive upper bounds for $P\{Z < x\}$ for $x$ positive but small. (5.18) already shows that

$$\lim_{x \downarrow 0} x^\alpha P\{Z \le x\} = 0 \text{ for all real } \alpha. \tag{7.1}$$

In this section we give a bound for the actual rate.

Let $(U_k)_{k=0}^\infty$ be independent, identically distributed random variables, all uniformly distributed on [0,1]. If U has this Unif(0, 1) distribution and E has the exponential distribution with expected value 1, then by

$$P\{U \le x\} = \min(x, 1) \ge 1 - exp\{-x\} = P\{E \le x\} \text{ for } x \ge 0, \tag{7.2}$$

U is stochastically less than E. By (4.10),

$$Z \text{ is stochastically larger than } \sum_{k=0}^\infty c^k U_k. \tag{7.3}$$

For $0 < x < 1$, define

$$K(x) = \text{ integer part of } \frac{\log(x)}{\log(c)}, \text{ i.e. } c^{K(x)+1} < x \le c^{K(x)}. \tag{7.4}$$

In (7.4) we must not forget that both $\log(x)$ and $\log(c)$ are negative. With (7.2), (7.3) we now have that for $0 < x < 1$:

$$P\{Z \le x\} < P\{\sum_{k=0}^{K(x)} c^k E_k \le x\} \le P\{\sum_{k=0}^{K(x)} c^k U_k \le x\} = \frac{x^{K(x)+1}}{(K(x)+1)!} \prod_{k=0}^{K(x)} c^{-k} =$$

$$\frac{x^{K(x)+1}}{\Gamma(K(x)+2)} c^{-\frac{1}{2}K(x)(K(x)+1)} \le \frac{x^{\left(\frac{1}{2}\frac{\log(x)}{\log(c)}\right)}}{\Gamma(\frac{\log(x)}{\log(c)}+1)}. \tag{7.5}$$

20

(7.5) provides an alternative proof of the fact that $\nu(\alpha) < \infty$ for all real $\alpha$.

We have a heuristic proof that in fact for all $0 < x < 1$ and all $0 < c < 1$

$$\frac{1}{2}exp\{-\frac{M}{|\log(c)|}\}\frac{\left(xc^{-\frac{1}{2}K(x)}\right)^{K(x)+1}}{(K(x)+1)!} < P\{\sum_{k=0}^{K(x)} c^k E_k \leq x\} < \frac{\left(xc^{-\frac{1}{2}K(x)}\right)^{K(x)+1}}{(K(x)+1)!}, \qquad (7.6)$$

where

$$M = \sum_{k=1}^{\infty} \frac{1}{k^2 2^k} = \int_0^1 \frac{|\log(1-\frac{y}{2})|}{y}\,dy, \quad \frac{1}{2} < M < \log(2) = .693\ldots. \qquad (7.7)$$

The heuristic proof goes as follows:

$$\left(\prod_{k=0}^{K(x)}(1-\frac{x}{2c^k})\right)\left(\prod_{k=0}^{K(x)}\frac{x}{c^k}\right) < P\{E_k < \frac{x}{c^k} \text{ for } 0 \leq k \leq K(x)\} <$$

$$P\{U_k < \frac{x}{c^k} \text{ for } 0 \leq k \leq K(x)\} = \left(\prod_{k=0}^{K(x)}\frac{x}{c^k}\right), \qquad (7.8)$$

and

$$\log(\frac{1}{2}) + \int_0^{K(x)}\log(1-\frac{x}{2c^{K(x)}}c^y)dy < \log\left(\prod_{k=0}^{K(x)}(1-\frac{x}{2c^k})\right) < 0, \text{ etc.} \qquad (7.9)$$

The heuristic step now is that

$$P\{\sum_{k=0}^{K(x)} c^k E_k \leq x \mid c^k E_k \leq x \text{ for } 0 \leq k \leq K(x)\} >$$

$$P\{\sum_{k=0}^{K(x)} c^k U_k \leq x \mid c^k U_k \leq x \text{ for } 0 \leq k \leq K(x)\} = \frac{1}{(K(x)+1)!} \qquad (7.10)$$

It seems likely that (7.5) actually is quite sharp: on $[0,x]$ the distributions of $\sum_{k=0}^{K(x)} c^k E_k$ and of $\sum_{k=0}^{K(x)} c^k U_k$ are fairly similar, and $\sum_{k=K(x)+1}^{\infty} c^k E_k$ has the same distribution as $c^{K(x)+1} Z < xZ$. Use of random variables with (for example) linearly decreasing densities on $[0,x]$ gives improved bounds, but quite possibly at the cost of an amount of formula pushing not proportional to the improvement in the bound.

21

There is an alternative way to obtain upper bounds as in (7.5) which often achieves results stronger than (7.5). Using the idea of the Bienaimee-Chebychev inequality we see that for $\alpha > 0$, $x > 0$:

$$\nu(-\alpha) = \int_0^\infty y^{-\alpha} dG(y) > \int_0^x y^{-\alpha} dG(y) > x^{-\alpha} G(x), \qquad (7.11)$$

and hence

$$P\{Z \le x\} < x^\alpha \nu(-\alpha) = \frac{\pi x^\alpha}{\Gamma(\alpha)\sin(\pi\alpha)} \prod_{j=1}^\infty \left(\frac{1 - c^{-\alpha+j}}{1 - c^j}\right). \qquad (7.12)$$

(7.12) uses (5.22). For the sharpest possible bound, we should find the best possible $\alpha$ as function of $x$. This seems hard. Quite a good bound is obtained by choosing $\alpha$ integer, $\alpha = K + 1$. (7.12) then reads:

$$P\{Z \le x\} < \frac{\left(x c^{-\frac{1}{2}K}\right)^{K+1}}{K!} \log(\frac{1}{c}) \prod_{j=1}^K (1 - c^j). \qquad (7.13)$$

In (7.13), the choice $K = K(x)$ or $K = K(x) + 1$ ($K(x)$ as in (7.4)) gives results similar to (7.5). The most likely situation seems to be that where the optimal value of $K$ is somewhat larger than $K(x)$. For $x$ small compared with $c$, the optimal value of $K$ must be close to the solution of

$$K c^K = x, \ K > \frac{\log(x)}{\log(c)}. \qquad (7.14)$$

For $c$ close to 1, (7.13) is much sharper than (7.5), but for c very close to zero it seems that (7.5) can be sharper. There is an obvious relationship between (7.13) and (7.6) – (7.10).

# 8.   Ferguson's "lose one or double" process

Ferguson (1972) [13] studies a process which, for reasons that will become clear later on, we call $(Z_{bw,n})_{n=0}^\infty$. (*bw* stands for *backward*, see later in this section for an expanation).

In Ferguson's paper, $Z_{bw,n}$ either decreases by one (with probability $\pi$) or doubles (with probability $1 - \pi$). To highlight the relationship with the sections $4 - 8$ in this paper we do not (necessarily) double, but multiply by a factor $c^{-1}$, for some $0 < c < 1$. The transition probabilities then are given by

$$P\{Z_{bw,n+1} = x - 1 | Z_{bw,n} = x\} = \pi, \ P\{Z_{bw,n+1} = c^{-1}x | Z_{bw,n} = x\} = 1 - \pi. \tag{8.1}$$

Clearly, the drift of this process (if $Z_{bw,n} = x$) equals

$$D(x) = \frac{1 - c}{c}(1 - \pi)x - \pi. \tag{8.2}$$

This drift goes to infinity if $x$ goes to infinity.

The gambler is ruined if or when $Z_{bw,.}$ reaches $(-\infty, 0]$. If this occurs (for the first time) at time n we say the gambler is ruined at time n. It is notationally convenient to pretend that the process continues after ruin occurs: (8.1) shows that once the process $Z_{bw,.}$ reaches $(-\infty, 0]$ the sample path becomes strictly decreasing and goes to $-\infty$. $q_x$ denotes the probability that ruin occurs as function of the starting level $Z_{bw,0} = x$. Clearly, for $n \in \{1, 2, \ldots\}$

$$q_n \geq \pi^n. \tag{8.3}$$

(8.3) and the paragraph preceding it show that

$$P\{-\infty < \liminf_{n \to \infty} Z_{bw,n} < +\infty\} = 0. \tag{8.4}$$

Hence:

$$P\{\text{either } \lim_{n \to \infty} Z_{bw,n} = -\infty \text{ or } \lim_{n \to \infty} Z_{bw,n} = +\infty\} = 1, \text{ and} \tag{8.5}$$

$$q_x = P\{\lim_{n \to \infty} Z_{bw,n} = -\infty | Z_{bw,0} = x\}. \tag{8.6}$$

Ferguson studies $q_x$ by studying (for $c^{-1} = 2$) the difference equation

$$q_x = \pi q_{x-1} + (1 - \pi)q_{c^{-1}x}. \tag{8.7}$$

23

If we let $\pi \uparrow 1$ and at the same time rescale space and time we get a continuous time version of the process $Z_{bw,n}$. We now have a Poisson process with intensity $c_1$, with points $(\tau_k)_{k=0}^{\infty}$, and we have the process $(Z_{bw}(t))_{0 \le t < \infty}$ for which in between the points $\tau_k$

$$(d/dt)Z_{bw}(t) = -c_2, \tag{8.8}$$

and in the points $\tau_k$

$$Z_{bw}(\tau_k+) = c^{-1}Z_{bw}(\tau_k-). \tag{8.9}$$

For the process $Z_{bw}(.)$ results similar to (8.2) – (8.6) can be derived. Without restriction we will further assume that

$$c_2 = c_1.$$

(Otherwise: rescale space).

We now see that $Z_{bw}(t)$ has the same sample paths as the process $Z(t)$ in (4.6) – (4.8), *only backward!* Henceforth, $Z(t)$ as in (4.6) – (4.8) is denoted by $Z_{fw}(t)$. We call $Z_{fw}(.)$ the forward process, and we call $Z_{bw}(.)$ the backward process. Ferguson also studies the ruin probability of this backward process. To enhance the symmetry, we will instead consider the function

$$G(z) = P\{\lim_{t \to \infty} Z_{bw}(t) = +\infty | Z_{bw}(0) = z\}. \tag{8.10}$$

Clearly, for $z \ge 0$:

$$1 - G(z) \ge exp\{-z\}. \tag{8.11}$$

Instead of the difference equation (8.7) we now get the differential equation

$$(d/dz)G(z) = G(c^{-1}z) - G(z). \tag{8.12}$$

In Appendix A, in (A.4), we will see that the distribution function $G(.)$ we introduced in section 4 satisfies this same differential equation. It is easily seen that actualy the two

functions $G(.)$ are identical. Ferguson indeed uses this differential equation to derive that $G(.)$ is the distribution function of a random variable $Z$ of the form (4.10).

In the remainder of this section we will essentially repeat the argument used in section 4 to give a direct probabilistic proof of this observation for the backward process $Z_{bw}(.)$ (not depending on the differential equation (8.12)), and to derive a similar result for the discrete time process $Z_{bw,.}$. For the process $Z_{bw,.}$, let $\tau_k$ be the time of the k-th (starting at zero) multiplicative increase, so that

$$0 \leq \tau_0 < \tau_1 < \dots , \tag{8.13}$$

and let

$$I_0 = \tau_0, \ I_k = \tau_k - \tau_{k-1} - 1 \text{ for } k \geq 1. \tag{8.14}$$

Clearly, $(I_k)_{k=0}^\infty$ are independent, identically distributed, all geometricaly distributed with

$$P\{I_k = n\} = (1 - \pi)\pi^n \text{ for } n \geq 0. \tag{8.15}$$

Let

$$Z_{bw,0} = z. \tag{8.16}$$

Then

$$\text{if } 0 \leq n \leq \tau_0 \text{ then } Z_{bw,n} = z - n, \text{ while} \tag{8.17}$$

$$\text{if } \tau_k < n \leq \tau_{k+1} \text{ then } Z_{bw,n} = \frac{1}{c^{k+1}}(z - I_0 - cI_1 - \dots - c^k I_k) - (n - \tau_k - 1) =$$

$$\frac{1}{c^{k+1}}\left(z - \sum_{j=0}^\infty c^j I_j\right) + \sum_{j=0}^\infty c^j I_{j+k+1} - (n - \tau_k - 1). \tag{8.18}$$

Define:

$$X(z) = z - \sum_{j=0}^\infty c^j I_j, \tag{8.19}$$

then

$$P\{\lim_{n \to \infty} Z_{bw,n} = \text{sign}(X(z)).\infty\} = 1. \tag{8.20}$$

25

Clearly, (8.5), (8.20) imply that

$$P\{X(z) = 0\} = 0 \text{ for all } z. \tag{8.21}$$

Since $(X(z) - z)$ is the sum of a countably infinite number of independent, discretely distributed random variables, the distribution of $(X(z) - z)$ is either purely discrete, or purely singular, or purely absolutely continuous (see [12]). (8.21) shows that the distribution can not be discrete.

For the discrete time process $Z_{bw,}$ we now have

$$1 - q_z = P\{\lim_{n \to \infty} Z_{bw,n} = +\infty | Z_{bw,0} = z\} = P\{\sum_{j=0}^{\infty} c^j I_j \leq z\}. \tag{8.22}$$

For the continuous time process $Z_{bw}(.)$ we have a similar result: again, let $\tau_k$ be the time of the k-th (starting at zero) multiplicative increase, and as in section 4, define

$$H_0 = \tau_0, \ H_k = \tau_k - \tau_{k-1} \text{ for } k \geq 1. \tag{8.23}$$

Then (given $Z_{bw}(0) = z$):

$$\text{if } 0 \leq t < \tau_0 \text{ then } Z_{bw}(t) = z - c_1 t, \text{ while} \tag{8.24}$$

$$\text{if } \tau_k < t < \tau_{k+1} \text{ then } Z_{bw}(t) = \frac{1}{c^{k+1}}(z - c_1(H_0 - cH_1 - \ldots - c^k H_k)) - c_1(t - \tau_k) =$$

$$\frac{1}{c^{k+1}}\left(z - c_1 \sum_{j=0}^{\infty} c^j H_j\right) + c_1 \sum_{j=0}^{\infty} c^j H_{j+k+1} - c_1(t - \tau_k). \tag{8.25}$$

As in section 4, $(H_j)_{j=0}^{\infty}$ are i.i.d., exponentially distributed, with expected value $c_1^{-1}$. From here on, the argument runs as in (8.18) − (8.20).

We have given two proofs that for all $z \geq 0$:

$$\lim_{t \to \infty} P\{Z_{fw}(t) \leq z\} = P\{\lim_{t \to \infty} Z_{bw}(t) = \infty | Z_{bw}(0) = z\}. \tag{8.26}$$

26

The first proof is the observation that the differential equations (5.4) and (8.12) are the same. The second proof is the combination of theorem 1 and the observation (8.23) – (8.25) etc. It would be nice to find a more general probabilistic proof that also holds for other processes than $Z_{bw}(.)$ and $Z_{fw}(.)$. It is clear that a similar result holds for the discrete time processes $(Z_{bw,n})_{n=0}^{\infty}$ and $(Z_{fw,n})_{n=0}^{\infty}$, the latter to be defined in the obvious way. Would similar results also hold for the other processes defined in section 4, and for their discrete-time analogs, for example the original congestion window size process $(W_n)_{n=0}^{\infty}$ defined in section 1?

## A.   A Different Approach

With $F(w, t)$ as in (4.19), clearly

$$F(w + \frac{c_2 \Delta}{w^m}, t + \Delta) = F(w, t) + c_1 \Delta \left( F(\beta^{-1} w, t) - F(w, t) \right) + o(\Delta) \ (\Delta \downarrow 0). \quad \text{(A.1)}$$

By letting $\Delta \downarrow 0$ we see that

$$\frac{c_2}{w^m}(d/dw)F(w, t) + (d/dt)F(w, t) = c_1(F(\beta^{-1} w, t) - F(w, t)). \quad \text{(A.2)}$$

For the stationary distribution, $d/dt = 0$. We drop the argument $t$. (5.2) becomes:

$$(d/dw)F(w) = \eta . w^m (F(\beta^{-1} w) - F(w)), \text{ where } \eta \text{ is as in (4.5).} \quad \text{(A.3)}$$

There are several ways to use (A.3) to find the stationary distribution functions $F$ and $G$. The most elementary is to show that $F(.)$ as in (5.10), (5.11) satisfies the differential equation (A.4). Another one is to directly use (A.3) to find the moments of Z, and then obtain an alternative proof of theorem 1 from those moments. This we will do in the next appendix. Yet another way will be presented next: Specializing (A.3) to the case with $c_2 = c_1$ and $m = 0$ we get:

$$(d/dz)G(z) = G(c^{-1} z) - G(z). \quad \text{(A.4)}$$

27

Ferguson [Ref] obtains the same differential equation (A.4) for a related process, see section 8 of this note, in particular formula (8.12).

By multiplying (A.4) with $exp\{-sz\}$ and integrating from zero to infinity we get ($\phi_Z(.)$ as in section 5):

$$\phi_Z(s) = \frac{1}{1+s}\phi_Z(cs) = \frac{1}{(1+s)(1+cs)}\phi(c^2 s) = \ldots = \prod_{k=0}^{\infty}\frac{1}{1+c^k s}. \tag{A.5}$$

This was a second proof of theorem 1.

# B.  (A.3) as a moment problem

It is convenient to re-write (A.3) as

$$(d/dw)F(w) = \eta.w^m \left((1-F(w)) - (1-F(\beta^{-1}w))\right). \tag{B.1}$$

Multiplying (B.1) with $w^\alpha$ and integrating from zero to infinity, we get ($\mu(\alpha)$ as in (5.13)):

$$\mu(\alpha) = \eta.\frac{1-\beta^{\alpha+m+1}}{\alpha+m+1}\mu(\alpha+m+1), \tag{B.2}$$

or (re-written)

$$\mu(\alpha+m+1) = \frac{\alpha+m+1}{\eta(1-\beta^{\alpha+m+1})}.\mu(\alpha). \tag{B.3}$$

Since $\mu(0) = 1$, for $k \in \{0, 1, 2, \ldots\}$ we now have

$$\mu(k(m+1)) = \frac{(m+1)^k.k!}{\eta^k(1-\beta^{(m+1)})(1-\beta^{2(m+1)})\ldots(1-\beta^{k(m+1)})}. \tag{B.4}$$

Formula (B.4) was the breakthrough that made us think of the substitution (4.6). With $\nu(\alpha)$, $\mu(\alpha)$, as in (5.13), (B.4) now reads

$$\nu(k) = \frac{k!}{(1-c)(1-c^2)\ldots(1-c^k)}. \tag{B.5}$$

We now know all the moments of Z (an alternative proof of (5.15)), and for the Laplace-Stieltjes Transform of $Z$ we have

$$\phi_Z(s) = E[exp\{-sZ\}] = \sum_{k=0}^{\infty} \frac{(-1)^k s^k}{(1-c)(1-c^2)\dots(1-c^k)}. \tag{B.6}$$

(B.6) clearly is a powerseries with radius of convergence $1 > 0$. Hence, the moments uniquely determine the distribution. By (5.2) we now have an alternative proof of (5.1), and thereby a third proof of theorem 1. This was actually the first such proof we found. It is also possible to derive (5.17) from (B.3) and (B.5) without use of (5.16) or (5.2). This proof goes as follows: Using (B.3) in the same way as when we derived (B.4) we obtain: for $\alpha > -1$, $k \in \{0, 1, 2, \dots\}$, we have

$$\nu(\alpha + k) = \frac{\Gamma(\alpha + k + 1)}{(1-c^{\alpha+1})(1-c^{\alpha+2})\dots(1-c^{\alpha+k})} \cdot \frac{\nu(\alpha)}{\Gamma(\alpha+1)}. \tag{B.7}$$

Next we use the fact that $\log(\nu(\alpha))$ is convex in $\alpha$, see e.g. Loeve p 156 [Ref]. Choose any $0 \leq p \leq 1$. Since for any $k \in \{0, 1, \dots\}$

$$k + p = (1-p)k + p(k+1) \text{ and } k = p(k-1+p) + (1-p)(k+p) \tag{B.8}$$

we have

$$\nu(k+p) \leq (\nu(k))^{1-p}.(\nu(k+1))^p \text{ and } \nu(k) \leq (\nu(k-1+p))^p.(\nu(k+p))^{1-p}. \tag{B.9}$$

Use of (B.7), (B.9) now gives the inequality

$$\frac{k!}{\Gamma(k+1+p)} \frac{(1-c^{p+1})(1-c^{p+2})\dots(1-c^{p+k})}{(1-c)(1-c^2)\dots(1-c^k)} \left(\frac{k+p}{1-c^{p+k}}\right)^p \leq \frac{\nu(p)}{\Gamma(p+1)} \leq$$
$$\frac{k!}{\Gamma(k+1+p)} \frac{(1-c^{p+1})(1-c^{p+2})\dots(1-c^{p+k})}{(1-c)(1-c^2)\dots(1-c^k)} \left(\frac{k+1}{1-c^{k+1}}\right)^p. \tag{B.10}$$

Letting $k \to \infty$ with p fixed this shows (use formula 2 in section 1.87 of Titchmarsh) that indeed

$$\nu(p) = \Gamma(p+1) \prod_{j=1}^{\infty} \left(\frac{1-c^{p+j}}{1-c^j}\right). \tag{B.11}$$

The proof above uses $0 \leq p \leq 1$, but once (B.11) holds in that range, (B.7) can be used to extend it to all $\alpha > -1$, and then as in (5.18) to all real $\alpha$.

# References

[1] Van Jacobson. Congestion avoidance and control, *Proceedings of ACM SIGCOMM '88*, August 1988.

[2] W. Stevens, *TCP/IP Illustrated*, volume 1, Addison-Wesley, Reading MA, 1994.

[3] Kevin Fall and Sally Floyd, Simulations-based comparisons of tahoe, reno and SACK TCP, *Proceedings of ACM SIGCOMM '96*, May 1996.

[4] Janey C. Hoe, Improving the start-up behavior of a congestion control scheme for TCP, *Proceedings of ACM SIGCOMM '96*, August 1996.

[5] Lawrence S. Brakmo, Sean W. O'Malley, and Larry L. Peterson, TCP vegas: New techniques for congestion detection and avoidance, *Proceedings of ACM SIGCOMM '94*, August 1994.

[6] Lawrence S. Brakmo and Larry L. Peterson, Performance problems in BSD4.4 TCP, *Proceedings of ACM SIGCOMM '95*, October 1995.

[7] Sally Floyd, TCP and successive fast retransmits, February 1995, Obtain via ftp://ftp.ee.lbl.gov/papers/fastretrans.ps.

[8] Sally Floyd and Van Jacobson, Random early detection gateways for congestion avoidance, *IEEE/ACM Transactions on Networking*, August 1993.

[9] Sally Floyd, TCP and explicit congestion notification, *ACM CCR*, 24(5), October 1994.

[10] Matthew Mathis, Jamshid Mahdavi, Sally Floyd, and Allyn Romanow, TCP selective acknowledgement options, May 1996, Internet Draft ("work in progress") draft-ietf-tcplw-sack-02.txt.

[11] Matthew Mathis, Jamshid Mahdavi, Forward Acknowledgment: Refining TCP Congestion Control, *Proceedings of ACM SIGCOMM '96*, August 1996.

[12] Jessen, B. and Wintner, A (1935) Distribution Functions and the Riemann Zeta Function, *Trans of the AMS* **38** pp 48–88.

[13] Ferguson, T.S. (1972) Lose a dollar or double your fortune, Proc Sixth Berkeley Symposium on Mathematical Statistics and Probability, Vol III, pp 657-666, Univ. of California Press.

[14] Gasper, G. and Rahman, M. (1990) *Basic Hypergeometric Series*, Encyclopedia of Mathematics and its Applications, Cambridge University Press.

[15] Hurwitz, A. and Courant, R. (1964) Vorlesungen ueber allgemeine Funktionentheorie und elliptische Funktionen, Springer Verlag (4-th edition).

[16] Loeve, M. (1955) *Probability Theory*, Van Nostrand.

[17] Stevens, W.R. (1994) *TCP/IP Illustrated*, Volumes I and II. Addison Wesley.

[18] Titchmarsh, E.C. (1932) *The Theory of Functions*, Oxford University Press.