

The F distribution and its relationship to the chi squared and t distributions

We have seen that the t distribution is appropriate for estimating critical values or confidence limits when a population has an underlying normal distribution, but the sample size is small. This is primarily a consequence of the difficulty of determining a population standard deviation, and using a method that more often than not underestimates it, the apparent distribution from the mean is distorted.

Many fundamental multivariate methods use the F distribution and its associated tests and critical values: it is the basis of the many common statistical tests in chemometrics, for example, for detecting outliers or whether an observation belongs to a predefined class.

CHI SQUARED

When we discussed the chi squared distribution [1], we noted that this represented the distribution of squared Mahalanobis distances from the mean, and in particular that if more than one variable is measured, there is no specific positive or negative direction, and as such, using squared distances (which are independent of direction) was essential. Hence, the chi-square distribution naturally extends from univariate to multivariate data.

The F distribution can be regarded as the equivalent extension of the t distribution when there is more than one variable but small sample sizes. There are numerous ways of introducing this distribution in the literature, which is widely employed in many diverse areas. In this and the next article, we focus primarily on the distribution of data in multidimensional space: the F distribution is often introduced in the context of analysis of variance. We will come across this distribution and its associated statistic in other contexts in later articles.

R. A. FISHER

The F is named after R. A. Fisher, who was a pioneering statistician most active in the 1920s and 1930s, and who worked in agricultural science in the UK. Many of the fundamental multivariate methods such as several approaches for one class classification or class modelling, including SIMCA and multivariate statistical process control, use the F distribution and its associated tests and critical values. It is the basis of the many common statistical tests in chemometrics, for example, for detecting outliers or whether an observation belongs to a



by Richard G. Brereton

Journal of Chemometrics Columnist

predefined class. However, in order to understand it, it is necessary to also understand its relationship to other distributions.

DEGREES OF FREEDOM

- The F distribution is characterized by two different types of degrees of freedom.
- It is often written $F(v_1, v_2)$. The horizontal axes of an F distribution cumulative distribution function (cdf) or probability density function represent the F statistic. We will see in the next article that if there is more than one variable, it is not equal to the squared Mahalanobis distance, unlike the chi-square statistic.
- In our context, if we consider a sample consisting of n observations and k variables, then v_1 represents the number of variables, and v_2 the number of observations minus the number of variables ($n - k$).
- Note that $F(v_1, v_2) \neq F(v_2, v_1)$.
- If a dataset is represented in matrix format, then the number of rows equals $v_2 + v_1$, and the number of columns equals v_1 .
- Note that it is common to define the dimensions of a matrix using rows first and then columns, but for the purpose of the F distribution, we swap these round, with the first degree of freedom referring to the number of variables or columns in a data matrix. This is illustrated in Figure 1.
- Note that n cannot be less than k . This apparent limitation has been discussed in the context of the Mahalanobis distance [2]. However, it can be overcome by performing principal components analysis first to reduce the number of variables.
- It is important to remember that, rather like the chi squared distribution, an F distribution is only obtained if the variables are independent. If, however, we use the Mahalanobis distance measure, this will always be so as discussed previously [2] as this in practice is equivalent to performing a principal component transformation.

UNIVARIATE RELATIONSHIPS

If there is only 1 degree of freedom, the data are univariate, and several straightforward relationships can be derived.

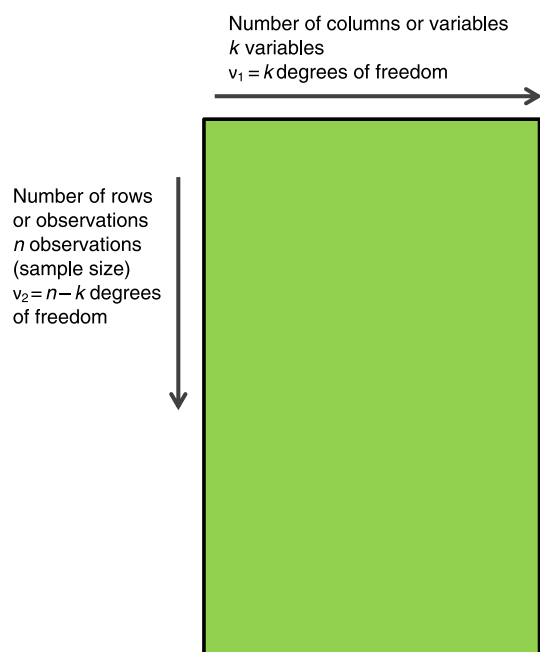


Figure 1. A data matrix.

- For the sake of illustration, we consider a measurement with an F value of 4.
- If there is a large sample size, then the F distribution, chi squared distribution, and the t^2 distributions all give the same results.
- In Excel, type `F.DIST(4,1,10 000 - 1,TRUE)`, putting $n = 10\ 000$: the 4 representing the value of F , the 1 equal to v_1 , and the $10\ 000 - 1$ equal to v_2 . The logical value 'TRUE' represents a cumulative distribution. This should give the proportion of the data that is expected to have an F statistic less than 4.
- To check the chi squared distribution, type `CHISQ.DIST(4,1,TRUE)`.
- For t , we should remember that we are dealing with squared distances for F and chi squared, but their square root when using t , although of course this example involves only one variable. Type `(T.DIST(SQRT(4),10 000,TRUE) - 0.5)*2`. This rather long expression is because we need to calculate the expected proportion of the data whose t statistic lies between +2 and -2, either side of the mean.
- The answer in all cases should be 0.954. Hence, 95.4% of the population have a chi squared or F statistic less than 4, or a t -statistic less than 2 (the square root of 4).

DISTRIBUTIONS

If there are a large number of observations (i.e. v_2 is large), then the shape of the F distribution is very similar to the chi squared distribution with v_1 degrees of freedom as illustrated in Figure 2, although there is a shift in position (in fact, chi squared equals $v_1 F$, and for 1 degree of freedom, they are both the same as $v_1 = 1$). Note that if both v_1 and v_2 are large, the F distribution also resembles the normal distribution, with a mean of 1.

Figure 3 illustrates several different F distributions. We can note several things.

- $F(2, 5)$ and $F(5, 2)$ are both very similar. In such situations, the sample sizes (7 and 52 respectively) are substantially greater than the number of variables ($=2$).

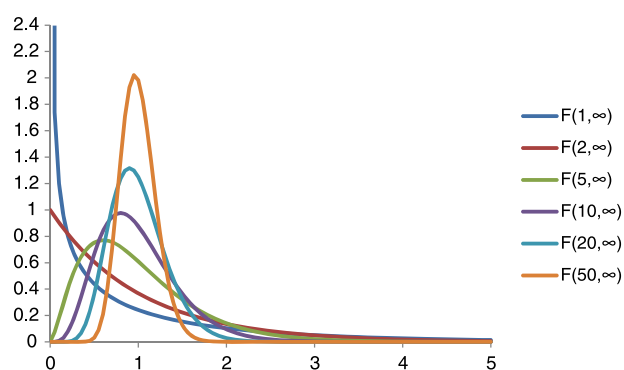


Figure 2. F distribution with a very large (effectively infinite) number of observations, as the number of variables (or v_1) increases. Horizontal axis, F -statistic; vertical axis, probability density function.

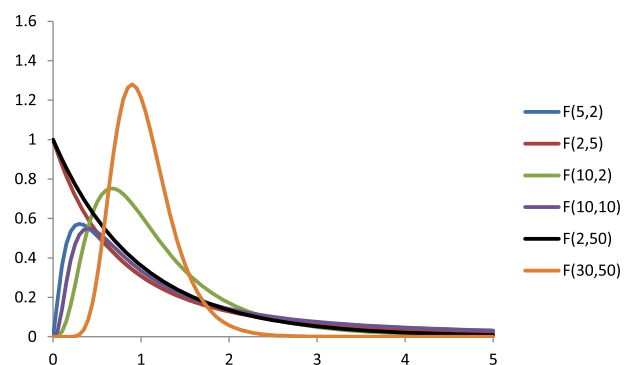


Figure 3. Several different F distributions. Horizontal axis, F statistic; vertical axis, probability density function.

- $F(5, 2)$ and $F(10, 10)$ are also similar in shape. The sample sizes (7 and 20) are not so much greater than the number of variables (5 and 10 respectively). Note that there is no problem if $v_1 > v_2$, only k must be less than n .
- $F(2, 5)$ and $F(5, 2)$ are very different. Never confuse the two different types of degrees of freedom.
- $F(30, 50)$ resembles a normal distribution, as both degrees of freedom are large.

Note the following properties.

- The mean of the distribution is equal to $v_2/(v_2 - 2)$ for $v_2 > 2$.
- The variance is equal to $[2v_2^2(v_1 + v_1 - 2)]/[v_1(v_2 - 2)^2(v_2 - 4)]$ for $v_2 > 4$

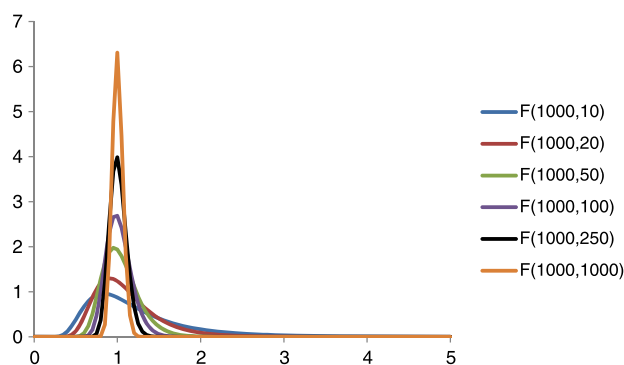


Figure 4. Several different F distributions with a very large number of variables (1000). Horizontal axis, F statistic; vertical axis, probability density function.



v_1		\rightarrow																						
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	20	30	40	50	75	100	250	500	∞	
1	4052	4999	5404	5624	5764	5859	5928	5981	6022	6056	6083	6107	6126	6143	6157	6209	6260	6286	6302	6324	6334	6353	6360	6366
2	9850	9900	9916	9925	9930	9933	9936	9938	9940	9941	9942	9943	9943	9943	9945	9945	9947	9948	9948	9948	9949	9950	9950	99499
3	3412	3082	2946	2871	2824	2791	2767	2749	2734	2723	2713	2705	2698	2692	2687	2669	2650	2641	2635	2628	2624	2617	2615	26125
4	2120	1800	1669	1598	1552	1521	1498	1480	1466	1455	1445	1437	1431	1425	1420	1402	1384	1375	1369	1361	1358	1351	1349	13463
5	1626	1327	1206	1139	1097	1067	1046	1029	1016	1005	9963	9888	9825	9770	9722	9553	9379	9291	9238	9166	9130	9064	9042	9020
6	1375	1092	9780	9148	8746	8466	8260	8102	7976	7874	7790	7718	7657	7605	7559	7396	7229	7143	7091	7022	6987	6923	6901	6880
7	1225	9547	8451	7847	7460	7191	6993	6840	6719	6620	6538	6469	6410	6359	6314	6155	5992	5908	5858	5789	5755	5692	5671	5650
8	1126	8649	7591	7006	6632	6371	6178	6029	5911	5814	5734	5667	5609	5559	5515	5359	5198	5116	5065	4998	4963	4901	4880	4859
9	1056	8022	6992	6422	6057	5802	5613	5467	5351	5257	5178	5111	5055	5005	4962	4808	4649	4567	4517	4449	4415	4351	4332	4311
10	1004	7559	6552	5994	5636	5386	5200	5057	4942	4849	4772	4706	4650	4601	4558	4405	4247	4165	4115	4048	4014	3951	3930	3909
11	9646	7206	6217	5668	5316	5069	4886	4744	4632	4539	4462	4397	4342	4293	4251	4099	3941	3860	3810	3742	3708	3645	3624	3602
12	9330	6927	5953	5412	5064	4821	4640	4499	4388	4296	4220	4155	4100	4052	4010	3858	3701	3619	3569	3501	3467	3404	3382	3361
13	9074	6701	5739	5205	4862	4620	4441	4302	4191	4100	4025	3960	3905	3857	3815	3665	3507	3425	3375	3307	3272	3209	3187	3165
14	8862	6515	5564	5035	4695	4456	4278	4140	4030	3939	3864	3800	3745	3698	3656	3505	3348	3266	3215	3147	3112	3048	3026	3004
15	8683	6359	5417	4893	4556	4318	4142	4004	3895	3805	3730	3666	3612	3564	3522	3372	3214	3132	3081	3012	2977	2913	2891	2868
20	8096	5849	4938	4431	4103	3871	3699	3564	3457	3368	3294	3231	3177	3130	3088	2938	2778	2695	2643	2572	2535	2468	2445	2421
30	7562	5390	4510	4018	3699	3473	3305	3173	3067	2979	2906	2843	2789	2742	2700	2549	2386	2299	2245	2170	2131	2057	2032	2006
40	7314	5178	4313	3828	3514	3291	3124	2993	2888	2801	2727	2665	2611	2563	2522	2369	2203	2114	2058	1980	1938	1860	1833	1805
50	7171	5057	4199	3720	3408	3186	3020	2890	2785	2698	2625	2563	2508	2461	2419	2265	2098	2007	1949	1868	1825	1742	1713	1683
75	6985	4900	4054	3580	3272	3052	2887	2758	2653	2567	2494	2431	2377	2329	2287	2132	1960	1866	1806	1720	1674	1583	1551	1516
100	6895	4824	3984	3513	3206	2988	2823	2694	2590	2503	2430	2368	2313	2265	2223	2067	1893	1797	1735	1646	1598	1501	1466	1427
250	6737	4691	3861	3395	3091	2875	2711	2583	2479	2392	2319	2256	2202	2154	2111	1953	1774	1674	1608	1511	1457	1343	1297	1244
500	6686	4648	3821	3357	3054	2838	2675	2547	2443	2356	2283	2220	2166	2117	2075	1915	1735	1633	1566	1465	1408	1285	1232	1164
∞	6635	4605	3782	3319	3017	2802	2639	2511	2407	2321	2248	2185	2130	2082	2039	1878	1696	1592	1523	1419	1358	1220	1153	1000

Table II. Critical values of F distribution, $p = 0.05$ 

		\rightarrow																								
ν_1		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	20	30	40	50	75	100	250	500	∞	
ν_2		1	161.4	199.5	215.7	224.6	230.2	234.0	236.8	238.9	240.5	241.9	243.0	243.9	244.7	245.4	245.9	248.0	250.1	251.1	251.8	252.6	253.0	253.8	254.1	254.3
2		18.51	19.00	19.16	19.25	19.30	19.33	19.35	19.37	19.38	19.40	19.40	19.40	19.41	19.42	19.42	19.43	19.45	19.46	19.47	19.48	19.48	19.49	19.49	19.49	19.50
3		10.13	9.552	9.277	9.117	9.013	8.941	8.887	8.845	8.812	8.785	8.763	8.745	8.729	8.715	8.703	8.660	8.617	8.594	8.591	8.581	8.563	8.554	8.537	8.532	8.526
4		7.709	6.944	6.591	6.388	6.256	6.163	6.094	6.041	5.999	5.964	5.936	5.912	5.891	5.873	5.858	5.803	5.746	5.717	5.699	5.699	5.676	5.664	5.643	5.635	5.628
5		6.608	5.786	5.409	5.192	5.050	4.950	4.876	4.818	4.772	4.735	4.704	4.678	4.655	4.636	4.619	4.558	4.496	4.464	4.444	4.444	4.418	4.405	4.381	4.373	4.365
6		5.987	5.143	4.757	4.534	4.387	4.284	4.207	4.147	4.099	4.060	4.027	4.000	3.976	3.956	3.938	3.874	3.808	3.774	3.754	3.754	3.726	3.712	3.686	3.678	3.669
7		5.591	4.737	4.347	4.120	3.972	3.866	3.787	3.726	3.677	3.637	3.603	3.575	3.550	3.529	3.511	3.445	3.376	3.340	3.319	3.319	3.290	3.275	3.248	3.239	3.230
8		5.318	4.459	4.066	3.838	3.688	3.581	3.500	3.438	3.388	3.347	3.313	3.284	3.259	3.237	3.218	3.150	3.079	3.043	3.020	3.020	2.990	2.975	2.947	2.937	2.928
9		5.117	4.256	3.863	3.633	3.482	3.374	3.293	3.230	3.179	3.137	3.102	3.073	3.048	3.025	3.006	2.936	2.864	2.826	2.803	2.803	2.771	2.756	2.726	2.717	2.707
10		4.965	4.103	3.708	3.478	3.326	3.217	3.135	3.072	3.020	2.978	2.943	2.913	2.887	2.865	2.845	2.774	2.700	2.661	2.637	2.637	2.605	2.588	2.558	2.548	2.538
11		4.844	3.982	3.587	3.357	3.204	3.095	3.012	2.948	2.896	2.854	2.818	2.788	2.761	2.739	2.719	2.646	2.570	2.531	2.507	2.507	2.473	2.457	2.426	2.415	2.404
12		4.747	3.885	3.490	3.259	3.106	2.996	2.913	2.849	2.796	2.753	2.717	2.687	2.660	2.637	2.617	2.544	2.466	2.426	2.401	2.401	2.367	2.350	2.318	2.307	2.296
13		4.667	3.806	3.411	3.179	3.025	2.915	2.832	2.767	2.714	2.671	2.635	2.604	2.577	2.554	2.533	2.459	2.380	2.339	2.314	2.314	2.279	2.261	2.229	2.218	2.206
14		4.600	3.739	3.344	3.112	2.958	2.848	2.764	2.699	2.646	2.602	2.565	2.534	2.507	2.484	2.463	2.388	2.308	2.266	2.241	2.241	2.205	2.187	2.154	2.142	2.131
15		4.543	3.682	3.287	3.056	2.901	2.790	2.707	2.641	2.588	2.544	2.507	2.475	2.448	2.424	2.403	2.328	2.247	2.204	2.178	2.178	2.142	2.123	2.089	2.078	2.066
20		4.351	3.493	3.098	2.866	2.711	2.599	2.514	2.447	2.393	2.348	2.310	2.278	2.250	2.225	2.203	2.124	2.039	1.994	1.966	1.966	1.927	1.907	1.869	1.856	1.843
30		4.171	3.316	2.922	2.690	2.534	2.421	2.334	2.266	2.211	2.165	2.126	2.092	2.063	2.037	2.015	1.932	1.841	1.792	1.761	1.761	1.718	1.695	1.652	1.637	1.622
40		4.085	3.232	2.839	2.606	2.449	2.336	2.249	2.180	2.124	2.077	2.038	2.003	1.974	1.948	1.924	1.839	1.744	1.693	1.660	1.660	1.614	1.589	1.542	1.526	1.509
50		4.034	3.183	2.790	2.557	2.400	2.286	2.199	2.130	2.073	2.026	1.986	1.952	1.921	1.895	1.871	1.784	1.687	1.634	1.599	1.551	1.525	1.475	1.457	1.438	
75		3.968	3.119	2.727	2.494	2.337	2.222	2.134	2.064	2.007	1.959	1.919	1.884	1.853	1.826	1.802	1.712	1.611	1.555	1.518	1.466	1.437	1.381	1.360	1.338	
100		3.936	3.087	2.696	2.463	2.305	2.191	2.103	2.032	1.975	1.927	1.886	1.850	1.819	1.792	1.768	1.676	1.573	1.515	1.477	1.422	1.392	1.331	1.308	1.283	
250		3.879	3.032	2.641	2.408	2.250	2.135	2.046	1.976	1.917	1.869	1.827	1.791	1.759	1.732	1.707	1.613	1.505	1.443	1.402	1.341	1.306	1.232	1.202	1.166	
500		3.860	3.014	2.623	2.390	2.232	2.117	2.028	1.957	1.899	1.850	1.808	1.772	1.740	1.712	1.686	1.592	1.482	1.419	1.376	1.312	1.275	1.194	1.159	1.113	
∞		3.841	2.996	2.605	2.372	2.214	2.099	2.010	1.938	1.880	1.831	1.789	1.752	1.720	1.692	1.666	1.571	1.459	1.394	1.350	1.283	1.243	1.152	1.106	1.000	

If ν_1 and ν_2 are large, then we can see from the equations earlier that the mean is approximately equal to 1 and the variance to $4/\nu_2$, so the larger ν_2 , the sharper the Gaussian (and more symmetric). This is illustrated in Figure 4 in the case where $\nu_1 = 1000$. Such situations, whilst rare in traditional statistical applications, may often be encountered in chemometrics where there may be a large number of variables, although would still require large sample sizes. Note that the mode changes position as the distribution becomes more symmetric as ν_2 increases.

PROBABILITY VALUES

The probability values corresponding to F distribution can easily be computed in Excel or most common environments such as MATLAB, and we have already introduced the basic syntax earlier. Note that earlier versions of Excel may have a simplified syntax; we refer to Office 2010 or above.

- What is the probability that an observation has an F value of greater than 2, if we measure 10 variables and our sample size is 30, that is, we have a data matrix of 30 rows and 10 columns?
- First, calculate ν_1 that equals 10. Then calculate ν_2 that equals $30 - 10$ or 20.
- So, we are looking at $F(10, 20)$.
- The syntax is `F.DIST(2,10,20,TRUE)` and should give an answer 0.910. The 'TRUE' implies that we are interested in the cdf of the F distribution. This implies that providing the underlying data are normally distributed, 9% of the data are expected to have an F value greater than 2.

TABLES

In traditional statistical texts, it is usual to present F distribution tables. Because there are rather many possible F distributions,

these are usually presented as critical values. A critical value of $p = 0.01$ gives the value of the F statistic that is expected to be exceeded by only 1% of the data, or in some cases, this can be called the 99% confidence limit. These tables are self-evident and are given in Tables I and II for two critical values. Note that F tables can be presented for different critical values. There are several more comprehensive tables available on the web [3,4] although it is recommended that p values are calculated in Excel or any other common environment. Note that the tables later are presented for the one-tailed F cdf in this article. In some contexts, it is appropriate to look at two-tailed F tests, but we will not at this phase be concerned with this.

Richard G. Brereton

University of Bristol, School of Chemistry, Cantocks Close,
Bristol, United Kingdom, BS8 1TS
E-mail: r.g.brereton@bris.ac.uk

REFERENCES

1. Brereton RG. The chi squared and multinormal distributions. *J. Chemometr.* 2014; **29**: 9–12. DOI: 10.1002/cem.2680.
2. Brereton RG. The Mahalanobis distance and its relationship to principal component scores. *J. Chemometr.* 2015; **29**: 143–145. DOI: 10.1002/cem.2692.
3. NIST/SEMATECH e-Handbook of Statistical Methods, April 2012, 1.3.6.7.3. Upper Critical Values of the F Distribution, <http://www.itl.nist.gov/div898/handbook/eda/section3/eda3673.htm>
4. F distribution critical values, *Medcalc*, <http://www.medcalc.org/manual/f-table.php>