



# CIRRELT

Centre interuniversitaire de recherche  
sur les réseaux d'entreprise, la logistique et le transport

Interuniversity Research Centre  
on Enterprise Networks, Logistics and Transportation

---

## Diameter Distribution Models for Quebec, Canada

Gregory Paradis  
Luc LeBel

June 2017

CIRRELT-2017-34

Bureaux de Montréal :  
Université de Montréal  
Pavillon André-Aisenstadt  
C.P. 6128, succursale Centre-ville  
Montréal (Québec)  
Canada H3C 3J7  
Téléphone : 514 343-7575  
Télécopie : 514 343-7121

Bureaux de Québec :  
Université Laval  
Pavillon Palais-Prince  
2325, de la Terrasse, bureau 2642  
Québec (Québec)  
Canada G1V 0A6  
Téléphone : 418 656-2073  
Télécopie : 418 656-2624

[www.cirrelt.ca](http://www.cirrelt.ca)

# Diameter Distribution Models for Quebec, Canada

Gregory Paradis<sup>1,\*</sup>, Luc Lebel<sup>1,2</sup>

1. Département des sciences du bois et de la forêt, Pavillon Abitibi-Price, 2405, rue de la Terrasse, Local 2121, Université Laval, Québec, Canada G1V 0A6
2. Interuniversity Research Centre on Enterprise Networks, Logistics and Transportation (CIRRELT)

**Abstract.** Statistical models predicting stem diameter distributions have found many applications in forestry. Our objective is to develop a methodology that can be used to derive a stem diameter distribution model for any combination of species and cover type in Quebec, Canada, using readily- available data from the government-run permanent sample plot inventory program. We test 25 truncated distributions from the generalized beta family to a large dataset of stems inventoried from permanent fixed-area plots in the province of Quebec, Canada, using a non-linear least-squares parameter-fitting algorithm. We describe a two-stage parameter-fitting methodology that produces improved estimates of parameter estimation error and parameter correlation for input data with bounded domain. We report best-fit distribution, best-fit parameter estimates (with standard error on parameter estimates), and AICc for each of 30 subdatasets covering the entire province of Quebec (representing all combinations of 10 species groups and 3 cover types). Best-fit results are clearly dominated by the four distributions in the generalized gamma family.

**Keywords:** Forest management, diameter distribution model.

**Acknowledgement.** This study was supported by funding from the FORAC Research Consortium.

Results and views expressed in this publication are the sole responsibility of the authors and do not necessarily reflect those of CIRRELT.

Les résultats et opinions contenus dans cette publication ne reflètent pas nécessairement la position du CIRRELT et n'engagent pas sa responsabilité.

---

\* Corresponding author: gregory.paradis.1@ulaval.ca

## 1 Introduction

Stem diameter distributions (i.e. stand tables) have long played an important role in forestry (Bailey and Dell, 1973; Hyink and Moser, 1983). Published models tend to be specific to a given combination of species, stand structure, geographic area, and inventory sampling method. No stem diameter distribution models have been published to date for the province of Quebec, Canada. Furthermore, no generalized methodology has been published to model stem diameter distributions from permanent sample plot (PSP) data, documenting how to correctly estimate best-fit parameter uncertainty and correlations for the common case where observed diameter data has *a priori* bounded domain (e.g. only merchantable stems of a certain minimum diameter are inventoried and trees never grow beyond a certain maximum diameter). The present study fills these gaps in the literature.

The most commonly-used statistical model used to describe stem diameter seems to be the Weibull distribution (Bailey and Dell, 1973; Liu et al., 2002; Cao, 2004; Coomes and Allen, 2007). Other models include the gamma (Nelson, 1964), exponential (Meyer and Stevenson, 1943) and  $S_B$  (Johnson, 1949) distributions. The Weibull, gamma and exponential distributions are all derivatives of the generalized gamma distribution, which is itself a member of the of generalized beta family of statistical distributions.

We fit 25 truncated distributions from the generalized beta family to a large dataset of stems from government-compiled permanent fixed-area plots in the province of Quebec, Canada. We describe a two-stage distribution-fitting methodology that correctly handles parameter estimation error and correlations for input data with bounded domain. We present best-fit distributions for 30 combinations of species group and cover type.

Our best-fit distribution results cover all combinations of species and cover types in Quebec, and could be used directly. Alternatively, our methodology can be easily replicated using readily-available PSP data, for example to derive models for different aggregations of species and cover type, or for a different geographic extent of plot data used as input. The two-stage parameter-fitting methodology is potentially applicable to any context where truncated data is fitted to standard-form statistical distributions.

The remainder of this paper is organized as follows. We describe our methodology in §2. Results are presented in §3, followed by discussion in §4.

## 2 Methods

Ducey and Gove (2015) document three parent distributions in the generalized beta family that can be used to derive several other distributions. These parent distributions are the generalized beta distribution of the first kind (GB1), the generalized beta distribution of the second kind (GB2), and the generalized gamma distribution (GG). The probability density functions (PDF) of GB1 and GB2 distributions have the following forms (adapted from Ducey and Gove,

2015)

$$\text{GB1}(x; a, b, p, q) = \frac{|a|x^{ap-1} [1 - (x/b)^a]^{q-1}}{b^{ap} B(p, q)}, \quad 0 < x^a < b^a, b > 0, p > 0, q > 0 \quad (1)$$

and

$$\text{GB2}(x; a, b, p, q) = \frac{|a|x^{ap-1} x^{q-1}}{b^{ap} B(p, q) [1 - (x/b)^a]^{p+q}}, \quad a > 0, b > 0, p > 0, q > 0 \quad (2)$$

defined for  $x > 0$ , where  $B(p, q)$  represents the beta function (not to be confounded with the beta, or generalized beta, distributions), which is given by

$$B(p, q) = \int_0^1 t^{p-1} (1-t)^{q-1} dt. \quad (3)$$

The PDF of the generalized gamma GG distribution has the following form

$$\text{GG}(x; a, \beta, p) = \frac{ax^{ap-1} e^{-(\frac{x}{\beta})^a}}{\beta^{ap} \Gamma(p)}, \quad a > 0, \beta > 0, p > 0 \quad (4)$$

defined for  $x > 0$ , where  $\Gamma(p)$  represents the gamma function (not to be confounded with the gamma, or generalized gamma, distributions), which is given by

$$\Gamma(p) = \int_0^\infty x^{p-1} e^{-x} dx. \quad (5)$$

We can define the PDFs for 22 different distributions in the generalized beta family in terms of one of the three parent distributions, as follows (adapted from Ducey and Gove, 2015)

$$\text{IB1}(x; b, p, q) = \text{GB1}(x; -1, b, p, q) \quad (6)$$

$$\text{UG}(x; b, d, q) = \lim_{a \rightarrow \infty} \text{GB1}(x; a, b, d/a, q) \quad (7)$$

$$\text{B1}(x; b, p, q) = \text{GB1}(x; 1, b, p, q) \quad (8)$$

$$\text{B2}(x; b, p, q) = \text{GB2}(x; 1, b, p, q) \quad (9)$$

$$\text{SM}(x; a, b, q) = \text{GB2}(x; a, b, 1, q) \quad (10)$$

$$\text{Dagum}(x; a, b, p) = \text{GB2}(x; a, b, p, 1) \quad (11)$$

$$\text{Pareto}(x; b, p) = \text{GB1}(x; -1, b, p, 1) \quad (12)$$

$$\text{P}(x; b, p) = \text{GB1}(x; 1, b, p, 1) \quad (13)$$

$$\text{LN}(x; \mu, \sigma) = \lim_{a \rightarrow 0} \text{GG}(x; a, (\sigma^2 a^2)^{1/a}, (a\mu + 1)/(\sigma^2 a^2)) \quad (14)$$

$$\text{GA}(x; \beta, p) = \text{GG}(x; 1, \beta, p) \quad (15)$$

$$\text{W}(x; a, \beta) = \text{GG}(x; a, \beta, 1) \quad (16)$$

$$\text{F}(x; u, v) = \text{GB2}(x; 1, v/u, u/2, v/2) \quad (17)$$

$$\text{L}(x; b, q) = \text{GB2}(x; 1, b, 1, q) \quad (18)$$

$$\text{IL}(x; b, p) = \text{GB2}(x; 1, b, p, 1) \quad (19)$$

$$\text{Fisk}(x; a, b) = \text{GB2}(x; a, b, 1, 1) \quad (20)$$

$$\text{U}(x; b) = \text{GB1}(x; 1, b, 1, 1) \quad (21)$$

$$\frac{1}{2}\text{N}(x; 0, \sigma) = \text{GG}(x; 2, \sigma^2, 1/2) \quad (22)$$

$$\chi^2(x; p) = \text{GG}(x; 1, 2, p) \quad (23)$$

$$\text{EXP}(x; \beta) = \text{GG}(x; 1, \beta, 1) \quad (24)$$

$$\text{R}(x; \beta) = \text{GG}(x; 2, \beta, 1) \quad (25)$$

$$\frac{1}{2}\text{t}(x; df) = \text{GB2}(x; 2, \sqrt{df}, 1/2, df/2) \quad (26)$$

$$\text{LL}(x; b) = \text{GB2}(x; 1, b, 1, 1) \quad (27)$$

We use a weighted non-linear least squares (NLLS) algorithm to fit target distributions to PSP inventory data binned into 26 size classes of uniform width  $W$ .

The objective function value of the NLLS problem minimizes the sum of squares of the residual terms

$$Z(f(x; \hat{\mathbf{P}})) = \min \sum_{i \in \{I | \hat{y}_i > 0\}} e(f(x_i; \mathbf{P}), \hat{y}_i)^2 \quad (28)$$

with

$$e(f(x_i; \mathbf{P}), \hat{y}_i) = w_i [f(x_i; \mathbf{P}) - \hat{y}_i] \quad (29)$$

where  $x_i$  is the diameter corresponding to the center of bin  $i \in I$ ,  $f(x_i; \mathbf{P})$  is the value of the PDF of the target distribution at  $x_i \in \mathbf{X}$  (given a vector of parameters  $\mathbf{P}$ ).  $\hat{y}_i \in \hat{\mathbf{Y}}$  represents the estimated stem density in bin  $i$ , which corresponds to the average of plot-wise stem density measurements.

Note that residual terms are scaled by a weight factor  $w_i = 1 - \min(E_{\hat{y}_i} \hat{y}_i^{-1}, 1)$ , which dampens the impact of  $\hat{y}_i$  on  $Z$  as a function of the relative margin of error  $E_{\hat{y}_i} \hat{y}_i^{-1}$ . We cap relative margin of error at 1 (negative values of  $w_i$  would have the effect of *rewarding* large residual value  $f(x_i; \mathbf{P}) - \hat{y}_i$ , which would make NLLS algorithm results unnecessarily difficult to interpret). Thus,  $w_i$  converges to 1 as relative margin of error approaches 0, and  $w_i = 0$  if  $E_{\hat{y}_i} \hat{y}_i^{-1} \geq 1$ . Note that if sampling error is high enough for all bins (due to insufficient sample size), such that  $w_i = 0, \forall i \in I$ , the objective function value is 0 regardless of values of input data vector  $\hat{\mathbf{Y}}$  and the NLLS optimisation problem becomes meaningless.

The margin of error corresponds to the product  $t\sigma_{\hat{y}_i}$  of the critical  $t$  value (with  $\alpha = 0.05$  and  $|\hat{\mathbf{Y}}| - 1$  degrees of freedom) and bin-wise sampling error

$$\sigma_{\hat{y}_i} = \sqrt{\frac{\sum_{j \in J} (y_{ij} - \hat{y}_i)^2}{|\hat{\mathbf{Y}}| - 1}} \quad (30)$$

where  $y_{ij}$  corresponds to the observed stem density in bin  $i$  in sample plot  $j$  (Schreuder et al., 2004).

We normalize our binned data, such that  $\sum_{i \in I} W \hat{y}_i = 1$ . The domain of input data is bounded, such that  $a \leq x_1 - w/2$  and  $x_{|I|} + w/2 \leq b$ , where  $a > 0$ . Our dataset intentionally includes only merchantable stems (i.e.  $a = 9$ ), and contains very few stems of diameter greater than 61 cm (i.e.  $b = 61$ ).

The integral of the standard forms of the PDFs described above over the interval  $[0, \infty]$  is 1 for any given vector of input parameters  $\mathbf{P}$ , that is

$$\int_0^\infty f(x; \mathbf{P}) dx = 1. \quad (31)$$

Fitting the standard forms of  $f$  to the normalized binned data will generally produce poor fits, as the sum of residuals will be positively biased due to bounded domain (i.e.  $\sum_{i \in I} e_i > 1$ ), with quality of fit inversely proportional to  $b - a$ . We can obtain a better fit using an augmented PDF  $f'(x; \mathbf{P}') = sf(x; \mathbf{P})$ . The global scaling parameter  $s$  effectively relaxes the unity constraint on the integral of  $f'$ . Thus, using  $f'$ , we obtain similar quality fits for any scaling of bin value vector  $\hat{\mathbf{Y}}$ .

The variance  $\sigma_{\hat{p}_j}^2$  of best-fit parameter estimator  $\hat{p}_j \in \hat{\mathbf{P}}$  corresponds to element  $j$  of the diagonal of the covariance matrix. The covariance matrix, which is automatically calculated by most software implementations of the NLLS algorithm, corresponds to the inverse of the negative of the expected values of the Hessian matrix  $-E[H(\hat{\mathbf{P}})]$ , where the Hessian  $H(\hat{\mathbf{P}})$  is the matrix of second derivatives of the likelihood function  $\mathcal{L}$  with respect to  $\hat{\mathbf{P}}$ . Standard error  $\sigma_{\hat{p}_j} = \sqrt{\sigma_{\hat{p}_j}^2}$  of parameter  $\hat{p}_j \in \hat{\mathbf{P}}$  corresponds to the square root of the variance.

Note that variance estimates are only correct asymptotically. In practice, fitting algorithms will use numerical approximations of Hessian matrix values. Quality of finite approximations of the second derivatives of  $\mathcal{L}$  will tend to be

proportional to sample size  $|\hat{Y}|$ , inversely proportional to distance from parameter constraint boundaries, and inversely proportional to the number of parameters  $|\hat{P}|$ .

Parameter estimation error for augmented function  $f'(x; \mathbf{P}')$  can be improved, without deteriorating fit quality, by solving the fitting problem in two stages. In the first stage, we determine  $\hat{\mathbf{P}}'$  by solving for  $Z(f'(x; \hat{\mathbf{P}}'))$ . For the best-case scenario, where  $f'(x; \mathbf{P}')$  is fitted to an infinitely large sample  $\hat{Y}$  randomly drawn from  $f'(x; \hat{\mathbf{P}}')$ , the estimated value of scaling parameter  $\hat{s} \in \hat{\mathbf{P}}'$  will completely eliminate the bias in the sum of residuals  $\sum_{i \in I} e(f(x_i; \hat{\mathbf{P}}), \hat{y}_i)$ , such that  $\int_a^b f(x; \hat{\mathbf{P}}) dx = \sum_{i \in I} W \hat{y}_i$ .

In the second stage, we solve for  $Z(f''(x; \hat{\mathbf{P}}, \hat{s}))$ , where  $f''$  corresponds to our augmented distribution  $f'$  with the scaling parameter value fixed at  $s = \hat{s}$  (i.e. only the original vector of parameters  $\mathbf{P}$  is optimized by the fitting algorithm).

The shape distributions from both stages are equivalent, such that

$$Z(f'(x; \hat{\mathbf{P}}')) \simeq Z(f''(x; \hat{\mathbf{P}}, \hat{s})). \quad (32)$$

However, error vector  $\sigma_{\hat{P}}$  and parameter covariance (which can be estimated from off-diagonal elements of the covariance matrix) estimated in the second stage will tend to be more reliable.

Our computational experiment dataset consists of 52 192 stems extracted from a database of PSP data, collected from public forests in Quebec (Canada). This data was collected by the *Ministère de la forêt, de la faune et des parcs* (MFFP) as part of the official government inventory program<sup>1</sup>, which operates on a 10-year cycle.

Data was collected throughout the province of Quebec, using 11.28 meter radius circular fixed-area plots. The full dataset contains 1 685 233 stems, sampled from 12 570 permanent sample plot locations. However, this includes repeated measures from four decennial inventory cycles, collected from 7 different PSP networks. We filtered data to include only stems from the most recent inventory cycle, which ensures that we are not tallying repeated measures on the same plots. We further filtered data to include only stems from the largest of the seven PSP networks (codename *BAS1*), which ensures uniform data-collection standards for all stems.

Our ultimate goal (i.e. beyond the scope of this paper) is to link a long-term wood supply optimization model with a short-term fibre-procurement optimization model. Thus, we are interested in modelling diameter distribution of merchantable stems in mature (operable), undisturbed stands. We therefore applied a series of other filters to our stem dataset to exclude plots in disturbed or immature stands, unmerchantable stems (with DBH less than 9 cm), very large stems (with DBH greater than 61 cm), and dead or otherwise unmerchantable stems.

<sup>1</sup>Detailed information on the PSP inventory program under which our test data was collected is available from the MFFP web site (<http://www.mffp.gouv.qc.ca/forets/inventaire/>).

A Jupyter Notebook containing Python code implementing these filters and detailed explanations is available from the corresponding author upon request. Although we do not have permission to distribute the PSP dataset, it is possible to request a copy from the *Ministère des forêts, de la faune et de parcs* (see footnote for URL).

We segmented the 52 192 stems in our filtered PSP dataset into 30 sub-datasets, representing combinations of 10 species groups and 3 cover types. More detailed information on species groups is provided in an appendix. For each of 30 sub-datasets  $d \in D$ , we applied our two-stage fitting method on 25 target distributions  $f \in F$  (i.e. GB1, GB2, and GG parent distributions, and the 22 special cases of these distributions defined in (6) through (27)).

We used the small-sample form of the Akaike information criterion (AICc) to evaluate goodness-of-fit for each combination of  $d \in D$  and  $f \in F$ . AICc is given by

$$\text{AICc} = \text{AIC} + \frac{2K(K+1)}{n-K-1} \quad (33)$$

with

$$\text{AIC} = 2K - n \ln \left( \frac{\chi^2}{n} \right) \quad (34)$$

where  $K = |\mathbf{P}| + 1$  (i.e. the cardinality of the parameter vector  $\mathbf{P}$ , plus the  $\mu$  parameter of the implicit i.i.d. Gaussian error distribution of input data vector  $\hat{\mathbf{Y}}$ ),  $n = |\hat{\mathbf{Y}}|$ , as recommended in Burnham and Anderson, (2002), and  $\chi^2$  is the sum of squared residuals given by

$$\chi^2 = \sum_{i \in I} e \left( f(x_i; \hat{\mathbf{P}}), \hat{y}_i \right)^2 \quad (35)$$

For each sub-dataset  $d \in D$ , we ranked distributions  $f \in F$  in decreasing order of AICc, and reported best-fit distribution, best-fit parameter values (with standard error estimates on parameter values) for first and second stages, and second-stage AICc.

### 3 Results

Figures 1 and 2 show best-fit distributions plotted against empirical input data distribution, binned by diameter class. The name of the best-fit distribution is identified in the legend for each subplot.

Table 1 reports estimated parameter values, standard error on parameter estimates, and second-stage AICc for best-fit distributions.

### 4 Discussion

As predicted, second-stage parameter standard error estimates are systematically lower than first-stage error estimates. This is attributable to fixing of the  $s$  parameter in the second stage.

Table 1: Best-fit distributions for each combination of species group and cover type. We report estimated parameter values and standard error for first- and second-stage fits, and second-stage AIC.

Species Group	Cover Type	Dist. Name	Parameters (Stage 1)	Parameters (Stage 2)	AICc (Stage 2)
Oak-Hickory	S	EXP	$\beta = 20.83 \pm 3.22$	$\beta = 15.00 \pm 2.62$	-62
	M	EXP	$\beta = 40.80 \pm 15.09$	$\beta = 36.27 \pm 6.54$	-178
	H	W	$a = 2.67 \pm 0.20$ $\beta = 27.09 \pm 0.80$	$a = 2.67 \pm 0.16$ $\beta = 27.10 \pm 0.74$	-250
Fir-Spruce-Pine-Larch	S	GA	$\beta = 2.35 \pm 0.04$ $p = 4.97 \pm 0.09$	$\beta = 2.35 \pm 0.01$ $p = 4.97 \pm 0.02$	-336
	M	W	$a = 1.25 \pm 0.06$ $\beta = 10.63 \pm 0.44$	$a = 1.25 \pm 0.01$ $\beta = 10.63 \pm 0.13$	-302
	H	GA	$\beta = 4.09 \pm 0.32$ $p = 2.42 \pm 0.27$	$\beta = 4.09 \pm 0.09$ $p = 2.42 \pm 0.04$	-253
Sugar Maple	S	EXP	$\beta = 6.34 \pm 1.44$	$\beta = 5.20 \pm 0.35$	-79
	M	W	$a = 1.33 \pm 0.14$ $\beta = 19.13 \pm 0.99$	$a = 1.35 \pm 0.06$ $\beta = 19.20 \pm 0.89$	-274
	H	W	$a = 1.42 \pm 0.06$ $\beta = 19.06 \pm 0.38$	$a = 1.42 \pm 0.02$ $\beta = 19.06 \pm 0.35$	-328
White Birch	S	EXP	$\beta = 6.24 \pm 0.39$	$\beta = 6.23 \pm 0.17$	-163
	M	GG	$a = 2.78 \pm 0.57$ $\beta = 23.60 \pm 2.57$ $p = 0.25 \pm 0.14$	$a = 2.77 \pm 0.22$ $\beta = 23.57 \pm 0.56$ $p = 0.25 \pm 0.03$	-289
	H	W	$a = 2.55 \pm 0.08$ $\beta = 17.67 \pm 0.18$	$a = 2.55 \pm 0.05$ $\beta = 17.67 \pm 0.17$	-238
Poplar	S	GA	$\beta = 4.85 \pm 1.41$ $p = 4.46 \pm 1.17$	$\beta = 4.73 \pm 0.85$ $p = 4.55 \pm 0.79$	-170
	M	W	$a = 2.76 \pm 0.07$ $\beta = 26.06 \pm 0.23$	$a = 2.76 \pm 0.05$ $\beta = 26.06 \pm 0.22$	-311
	H	W	$a = 3.03 \pm 0.15$ $\beta = 28.77 \pm 0.50$	$a = 3.03 \pm 0.12$ $\beta = 28.77 \pm 0.48$	-285
Pine	S	GG	$a = 5.17 \pm 2.62$ $\beta = 50.91 \pm 3.40$ $p = 0.20 \pm 0.14$	$a = 5.15 \pm 2.39$ $\beta = 50.83 \pm 2.40$ $p = 0.20 \pm 0.12$	-277
	M	EXP	$\beta = 49.73 \pm 12.38$	$\beta = 46.81 \pm 5.81$	-264
	H	EXP	$\beta = 9.94 \pm 1.58$	$\beta = 9.93 \pm 1.16$	-159
Other Hardwoods	S	GA	$\beta = 0.42 \pm 0.11$ $p = 28.54 \pm 7.58$	$\beta = 0.43 \pm 0.07$ $p = 28.08 \pm 4.80$	-101
	M	EXP	$\beta = 9.15 \pm 0.49$	$\beta = 9.15 \pm 0.34$	-208
	H	EXP	$\beta = 9.38 \pm 0.67$	$\beta = 9.41 \pm 0.48$	-253
Other Maples	S	$\chi^2$	$p = 5.95 \pm 0.21$	$p = 5.95 \pm 0.18$	-145
	M	GA	$\beta = 3.35 \pm 0.19$ $p = 4.45 \pm 0.26$	$\beta = 3.35 \pm 0.09$ $p = 4.45 \pm 0.12$	-288
	H	GA	$\beta = 4.91 \pm 0.48$ $p = 3.18 \pm 0.33$	$\beta = 4.90 \pm 0.20$ $p = 3.18 \pm 0.13$	-258
Yellow Birch	S	EXP	$\beta = 20.20 \pm 2.75$	$\beta = 20.04 \pm 2.61$	-223
	M	B1	$b = 60.60 \pm 3.82$ $p = 0.39 \pm 0.15$ $q = 1.73 \pm 0.30$	$b = 60.62 \pm 3.41$ $p = 0.40 \pm 0.02$ $q = 1.74 \pm 0.16$	-300
	H	EXP	$\beta = 16.60 \pm 1.34$	$\beta = 16.60 \pm 1.31$	-269
Eastern White Cedar	S	W	$a = 1.70 \pm 0.07$ $\beta = 20.42 \pm 0.37$	$a = 1.70 \pm 0.04$ $\beta = 20.42 \pm 0.36$	-317
	M	GG	$a = 2.87 \pm 0.94$ $\beta = 37.52 \pm 4.50$ $p = 0.16 \pm 0.13$	$a = 2.83 \pm 0.45$ $\beta = 37.34 \pm 1.28$ $p = 0.17 \pm 0.03$	-307
	H	EXP	$\beta = 7.75 \pm 0.87$	$\beta = 7.67 \pm 0.49$	-182

Four distributions (GG, GA, W, EXP) dominate our best-fit model selection experiment, taking first place for 28 out of 30 combinations of species group and cover type. The  $\chi^2$  and B1 distributions had the lowest AICc value for the other two cases.

This confirms previous results in the forestry literature reporting success using GG, GA, W, and EXP distributions to model stem diameter distribution from stem tally data. If analytic resources are highly constrained, we recommend limiting the list of candidate distributions these four. Note that in most cases many of the other 21 distributions were rejected altogether, either because the NLLS algorithm did not converge, or because fitting results were deemed unstable (due to extremely high parameter values, or extremely high or otherwise unreliable parameter error estimates). GA, W, and EXP distributions are all derivatives of the GG distribution, with one or more of its three parameters fixed to a value of 1. It is not surprising that the GG distribution generally had slightly better fit (i.e. lower  $\chi^2$ ), however results show that AICc recommends selecting a more parsimonious model in most (but not all) cases.

Overall, fit results are very good, as indicated by relatively low second-stage parameter standard error estimates. This observation can be confirmed by visual inspection of fit results in Figures 1 and 2. Relatively small sample sizes in some combinations of species group and cover type yielded binned datasets with more erratic values (including empty bins, which were excluded from NLLS algorithm input data before fitting). Naturally, best-fit parameter standard error and AIC values are higher for these datasets.

Our input dataset includes a large number of stems, inventoried throughout the province of Quebec. Our best-fit distribution results could be used to forecast diameter distribution for mature stands in Quebec, or in other locations with similar forests. For researchers looking for more customized fits, our two-stage methodology can easily be replicated on publicly-available inventory data.

## 5 Acknowledgements

This study was supported by funding from the *FORAC Research Consortium*.

## References

- Bailey, R. L. and Dell, T. (1973). Quantifying diameter distributions with the weibull function. *Forest Science*, 19(2):97–104.
- Burnham, K. and Anderson, D. (2002). *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*. Springer-Verlag New York.
- Cao, Q. V. (2004). Predicting parameters of a weibull function for modeling diameter distribution. *Forest science*, 50(5):682–685.

- Coomes, D. A. and Allen, R. B. (2007). Mortality and tree-size distributions in natural mixed-age forests. *Journal of Ecology*, 95(1):27–40.
- Ducey, M. J. and Gove, J. H. (2015). Size-biased distributions in the generalized beta distribution family, with applications to forestry. *Forestry*, 88(1):143–151.
- Hyink, D. M. and Moser, J. W. (1983). A generalized framework for projecting forest yield and stand structure using diameter distributions. *Forest Science*, 29(1):85–95.
- Johnson, N. L. (1949). Systems of frequency curves generated by methods of translation. *Biometrika*, 36(1/2):149–176.
- Liu, C., Zhang, L., Davis, C. J., Solomon, D. S., and Gove, J. H. (2002). A finite mixture model for characterizing the diameter distributions of mixed-species forest stands. *Forest Science*, 48(4):653–661.
- Meyer, H. A. and Stevenson, D. D. (1943). The structure and growth of virgin beech-birch-maple-hemlock forests in northern pennsylvania. *J. Agric. Res*, 67(2).
- Nelson, T. C. (1964). Diameter distribution and growth of loblolly pine. *Forest Science*, 10(1):105–114.
- Schreuder, H. T., Ernst, R., and Ramirez-Maldonado, H. (2004). Statistical techniques for sampling and monitoring natural resources. Technical Report RMRS-GTR-126, USDA Forest Service, Rocky Mountain Research Station.

## Appendix

Table 2 lists common and Latin names of species in the species groups used to segment our PSP data.

Table 2: Mapping of species group names to species common and Latin names. Alternate names are shown in parentheses.

Species Group	Common Name	Latin Name	
Other Hardwoods	(white, American) ash	<i>Fraxinus americana</i>	
	black ash	<i>Fraxinus nigra</i>	
	(green, red) ash	<i>Fraxinus pennsylvanica</i>	
	(North) American beech	<i>Fagus grandifolia</i>	
	(American, white, water) elm	<i>Ulmus americana</i>	
	slippery elm	<i>Ulmus rubra</i>	
	(rock, cork) elm	<i>Ulmus thomasi</i>	
	American hophornbeam	<i>Ostrya virginiana</i>	
	American linden (basswood)	<i>Tilia americana</i>	
	White Birch	grey birch	<i>Betula populifolia</i>
(white, paper) birch		<i>Betula papyrifera</i>	
Yellow Birch	yellow birch	<i>Betula alleghaniensis</i>	
Oak-Hickory	(bitternut, swamp) hickory	<i>Carya cordiformis</i>	
	shagbark hickory	<i>Carya ovata</i>	
	([wild, mountain] black, rum) cherry	<i>Prunus serotina</i>	
	white oak	<i>Quercus alba</i>	
	swamp white oak	<i>Quercus bicolor</i>	
	bur oak	<i>Quercus macrocarpa</i>	
	(northern, eastern) red oak	<i>Quercus rubra</i>	
	(butternut, white walnut)	<i>Juglans cinerea</i>	
	Spruce-Pine-Fir	white spruce	<i>Picea glauca</i>
		black spruce	<i>Picea mariana</i>
Norway spruce		<i>Picea abies</i>	
red spruce		<i>Picea rubens</i>	
hybrid larch		<i>Larix X marschlinii</i>	
Japanese larch		<i>Larix leptolepis</i>	
([eastern, American] larch, tamarack)		<i>Larix laricina</i>	
European larch		<i>Larix decidua</i>	
pitch pine		<i>Pinus rigida</i>	
([eastern, black] jack, grey, scrub) pine		<i>Pinus banksiana</i>	
Other Maples	Scots pine	<i>Pinus sylvestris</i>	
	balsam fir	<i>Abies balsamea</i>	
	(silver, silverleaf) maple	<i>Acer saccharinum</i>	
Sugar Maple	black maple	<i>Acer nigrum</i>	
	red maple	<i>Acer rubrum</i>	
Poplar	(sugar, rock) maple	<i>Acer saccharum</i>	
	balsam poplar	<i>Populus balsamifera</i>	
	eastern cottonwood (poplar)	<i>Populus deltoides</i>	
	(large-tooth, big-tooth) aspen	<i>Populus grandidentata</i>	
	hybrid poplar	<i>Populus sp X P. sp.</i>	
Pine	([quaking, trembling] [aspen, poplar])	<i>Populus tremuloides</i>	
	white pine	<i>Pinus strobus</i>	
	red pine	<i>Pinus resinosa</i>	
Hemlock-Cedar	(eastern, Canadian) hemlock	<i>Tsuga canadensis</i>	
	(eastern, northern) white-cedar	<i>Thuja occidentalis</i>	

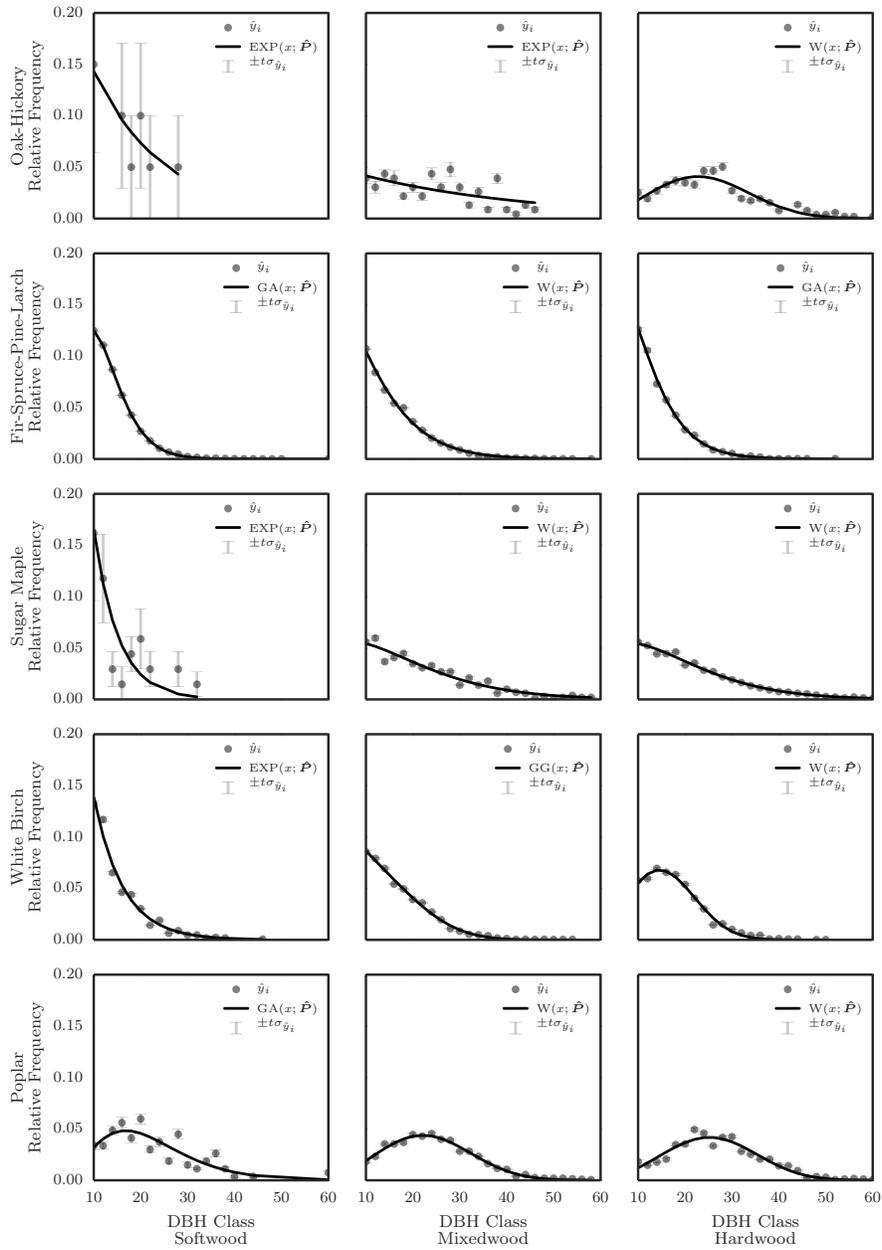


Figure 1: Best-fit distributions are shown with a solid line. Empirical distributions (binned by 2-cm diameter class) are shown with gray circles. Bin-wise sampling error is shown with light gray error bars. Species group is fixed for a given row of subfigures, and cover type is fixed for a given column of subfigures.

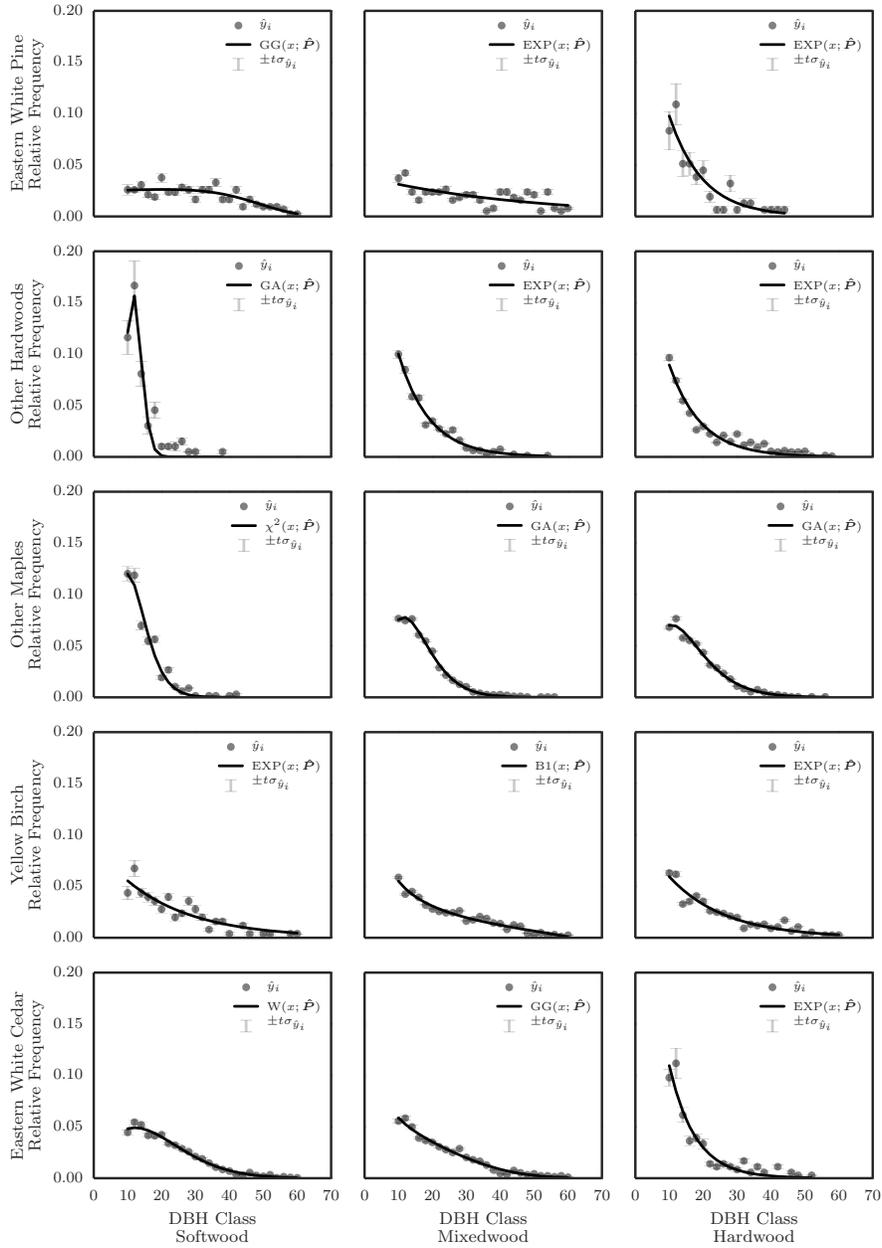


Figure 2: (Continued from Figure 1) Best-fit distributions are shown with a solid line. Empirical distributions (binned by 2-cm diameter class) are shown with gray circles. Bin-wise sampling error is shown with light gray error bars. Species group is fixed for a given row of subfigures, and cover type is fixed for a given column of subfigures.