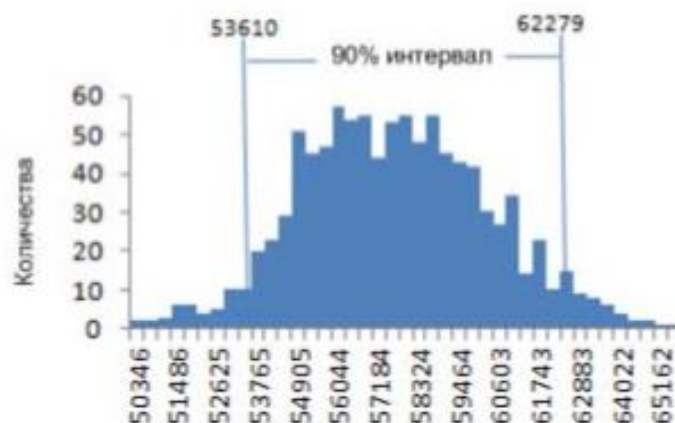


әсіресе t-үлестірумен жасалатын сенімді аралықтарға сүйенетін болады (Бұдан әрі осы тарауда "студенттің T-үлестірімін" қараңыз).



2.9. -Сурет. 20 мәннен іріктеме негізінде өтініш берушілердің жылдық табысы үшін бутстрапов сенім аралығы

Сенім аралығына байланысты пайыз сенім деңгейі деп аталады. Сенім деңгейі неғұрлым жоғары болса, аралық соғұрлым кең болады. Сонымен қатар, іріктеу неғұрлым аз болса, аралық неғұрлым кең болады (яғни, белгісіздік көп). Екі қасиет те жеткілікті логикалық: сіз қаншалықты сенімді болғыңыз келсе және сізде аз мәліметтер болса, шынайы мәнді алу үшін сенімді аралық неғұрлым кең болуы керек.

Қалыпты үлестіру

Дәстүрлі статистикада қоңырау тәрізді қалыпты үлестіру канондық 1 болып табылады. Үлгі статистикасының таралуы көбінесе қалыпты пішінге ие болғандықтан, оны осы үлестірімдерді жақындататын математикалық формулаларды жасауда күшті құралға айналдырды.

Негізгі терминдер

Қате (қате) деректер нүктесі мен болжамды немесе орташа мән арасындағы айырмашылық.

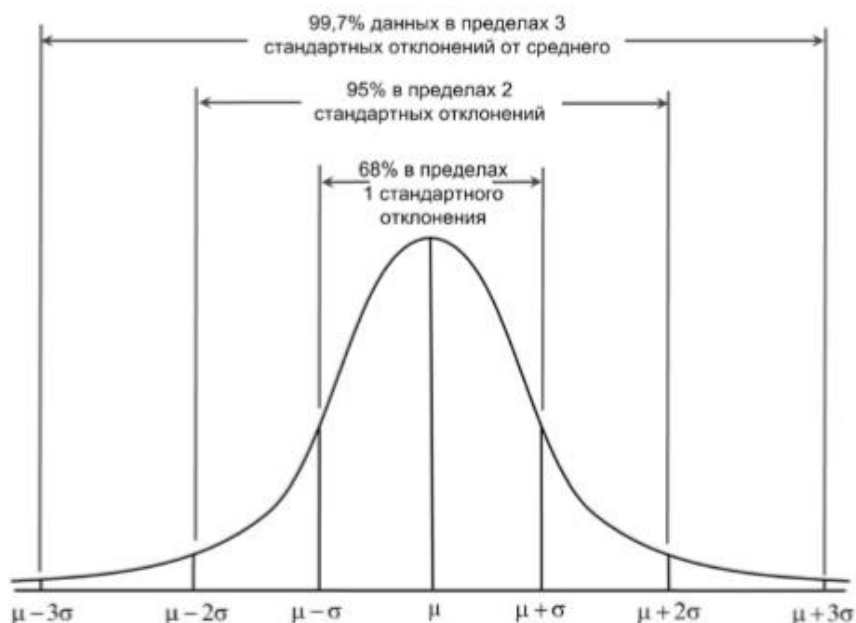
Стандарттау (standardize) орташа мәнді алып тастаңыз және стандартты ауытқуға бөліңіз.

*z-бағалау (z-score) жеке деректер нүктесін стандарттау нәтижесі.
Синонимі: стандартты бағалау.*

Стандартты қалыпты үлестіру (standard normal) орташа 0 - ге тең және стандартты ауытқу 1-ге тең қалыпты үлестіру.

*Квантиль-квантиль графигі (QQ-plot) үлгі үлестірімінің қалыпты үлестірімге қаншалықты жақын екенін көруге мүмкіндік беретін График.
Синонимдері: QQ-график, квантиль-квантиль графигі.*

Қалыпты үлестіруде (сурет. 2.10) деректердің 68% — ы орташадан бір стандартты ауытқу шегінде және 95% - ы екі стандартты ауытқу шегінде болады.



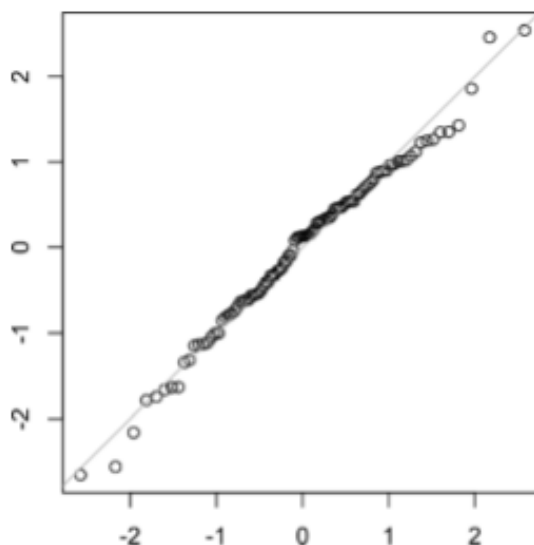
Сурет. 2.10. Қалыпты қисық

Стандартты қалыпты үлестіру және квантиль-квантиль графигтер.

Стандартты қалыпты үлестіру дегеніміз-Х осіндегі бірліктер орташа деңгейден стандартты ауытқулармен көрінетін үлестіру. Деректерді стандартты қалыпты үлестірумен салыстыру үшін орташа мәнді алып тастап, одан әрі стандартты ауытқуға бөлу керек; бұл процедура қалыпқа келтіру немесе стандарттау деп те аталады (6 - тараудың "Стандарттау (қалыпқа келтіру, z-бағалау)" бөлімін қараңыз). Айта кету керек, "стандарттау" бұл мағынада мәліметтер базасының жазбаларын стандарттаумен байланысты емес (яғни, жалпы форматқа келтіру). Түрлендірілген мән z-бағалау немесе

стандартты бағалау деп аталады, ал қалыпты үлестіру кейде z-үлестіру деп аталады.

Квантиль-квантиль графигі үлгінің қалыпты үлестіруден қаншалықты жақын екенін көзбен анықтау үшін қолданылады. Квантиль-квантиль графигі z-баллдарды төменнен жоғары қарай ұйымдастырады және y осіндегі әрбір мәннің z-баллдарын графикалық түрде көрсетеді; x осі — белгілі бір мәннің дәрежесі үшін тиісті нормальді үлестіру квантили. Деректер қалыпқа келтірілгендіктен, бірліктер орташа деңгейден стандартты ауытқу санына сәйкес келеді. Егер нүктелер диагональды сызықта болса, онда үлгіні бөлуді қалыптыға жақын деп санауға болады. Суретте. 2.11 көрсетілген квантиль-квантиль графигі 100 мәннен іріктеуге арналған, қалыпты үлестіруден кездейсоқ тәртіпте; күтілгендей, нүктелер сызыққа жақын жүреді. R-де бұл суретті функцияның көмегімен алуға болады



Сурет. 2.11. Квантиль-қалыпты үлестіруден алынған 100 мәннен іріктеменің квантильдік графигі

Ұзын құйрықты тарату

Қалыпты үлестірудің маңыздылығына қарамастан, Статистика тарихи тұрғыдан алғанда және оның атауы туралы айтылғаннан айырмашылығы, деректер әдетте қалыпты бөлінбейді.

Негізгі терминдер

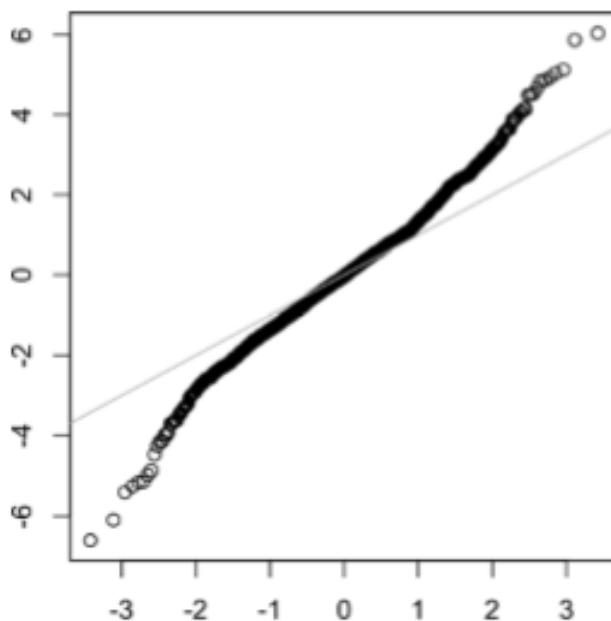
Құйрық (tail) жиіліктің таралуының ұзын тар бөлігі, онда салыстырмалы түрде шекті мәндер төмен жиілікте болады.

Ассиметрия (skew) - бір тарау құйрығы екіншісіне қарағанда ұзағырақ болатын жағдай. Синонимі: көлбеу.

Қалыпты үлестіру көбінесе қателіктер мен селективті статистиканы бөлуге қатысты негізделген және негізделген болса да, ол әдетте шикі деректердің таралуын сипаттамайды. Кейде бөлу кіріс деректері сияқты қатты қисайған (асимметриялы) немесе бөлу биномдық деректер сияқты дискретті болуы мүмкін. Симметриялық және асимметриялық нәсілдердің ұзын құйрықтары болуы мүмкін. Бөлу құйрықтары (кіші және үлкен) шекті мәндерге сәйкес келеді. Ұзын құйрықтар және олардың алдын-алу іс жүзінде кеңінен танылған. Нассим Талеб (Nassim Taleb) Қара аққу теориясын ұсынды, ол қор нарығының құлдырауы сияқты аномалды оқиғалар олардың қалыпты таралуы туралы айтқаннан гөрі көбірек болады деп болжайды. Деректердің ұзақ мерзімді сипатын көрсететін жақсы мысал-акциялардың кірістілігі. Суретте. 2.12 Netflix (NFLX) акцияларының күнделікті кірістілігінің сандық-сандық кестесін көрсетеді. R-де ол келесідей жасалады:

```
nflx <- sp500_px['NFLX']  
nflx <- diff(log(nflx[nflx>0]))  
qqnorm(nflx)  
abline(a=0, b=1, col='grey')
```

Күріштен айырмашылығы. 2.11. нүктелер төменгі мәндер үшін сызықтан сәл төмен және жоғары мәндер үшін сызықтан әлдеқайда жоғары орналасқан. Бұл дегеніміз, егер деректердің қалыпты таралуы болса, біз шекті мәндерді күткеннен әлдеқайда көп байқаймыз. Суретте. 2.12 тағы бір жалпы құбылысты көрсетеді: нүктелер орташа деңгейден бір стандартты ауытқу шегінде мәліметтер үшін сызыққа жақын орналасқан. Тьюки бұл құбылысты "ортасында қалыпты" деректер деп атайды, бірақ олардың ұзын құйрықтары бар ([Turkey-1987] қараңыз)



2.12-Сурет.. Netflix Квантиль-квантиль кірістілік кестесі

t- Таралуы. Стьюдента

Стьюдент таралуы немесе t таралуы-бұл қалыпты форманың таралуы, бірақ құйрықтарда сәл қалың және ұзағырақ. Ол үлгі статистикасының таралуын бейнелеу үшін кеңінен қолданылады. Үлгі орташаларының үлестірімдері әдетте t -үлестірімі түрінде болады және T -үлестірімдер тобы бар, олар үлгінің қаншалықты үлкен екеніне байланысты ерекшеленеді. Үлгі неғұрлым үлкен болса, t -үлестірімі соғұрлым қалыпты болады.

Негізгі терминдер

n

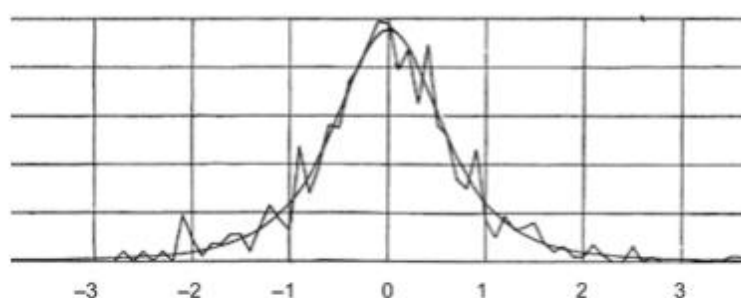
Үлгі мөлшері.

Еркіндік дәрежелері (Бостандық дәрежелері)-бұл T -үлестіруге әртүрлі үлгі өлшемдеріне, статисттерге және топтар санына бейімделуге мүмкіндік беретін Параметр.

t -үлестіруді көбінесе Стьюденттің t -үлестірімі деп атайды, өйткені ол 1908 жылы "Биометрика" (Biometrika) журналында У.С. Госсет (W. S. Gossett) "студент" деген бүркеншік атпен жарық көрді. Госсеттің жұмыс берушілері, Guinness сыра қайнату зауытының басшылығы бәсекелестердің статистикалық әдістерді қолданатындығын білгісі келмеді, сондықтан олар

Госсет өзінің мақаласында өз атын пайдаланбауын талап етті. Госсет "үлкен популяциядан алынған үлгі бойынша орташа үлгіні бөлу дегеніміз не?" Ол қайта іріктеу экспериментімен жұмыс істей бастады — қылмыскерлердің сол жақ саусағының биіктігі мен ұзындығының 3000 өлшемінен тұратын мәліметтер жиынтығынан 4 элементтен кездейсоқ үлгілерді алу. (Бұл эвгеника дәуірі болды және қылмыскерлер туралы мәліметтер мен қылмысқа бейімділік пен физикалық немесе психологиялық ерекшеліктер арасындағы байланысты анықтау басты назарда болды.) Ол стандартталған нәтижелерді (z-балл) x осіне және жиіліктерді y осіне қолданды. Сонымен қатар, ол қазір студенттің t функциясы деп аталатын функцияны алды және салыстыруды графикалық түрде көрсету арқылы осы функцияны үлгі нәтижелеріне сәйкестендірді (сурет. 2.13). Бүкіл x стандартты ауытқу болсын. Содан кейін үлгінің орташа мәні айналасындағы 90 пайыздық сенімділік интервалы келесі форма арқылы беріледі:

Егер есептеу қуаты 1908 жылы кеңінен қол жетімді болса, онда статистика басынан бастап қайта іріктеудің есептеу сыйымды әдістеріне көбірек сүйенетін еді. Компьютерлерден айырылған статистика мамандары іріктемені жақындату үшін математика және t-бөлу сияқты функцияларды қарастырды. 1980 ж. компьютерлердің есептеу қуаты қайта іріктеумен практикалық эксперименттерді жандандыруға мүмкіндік берді, бірақ сол кезде t-үлестіруді және ұқсас үлестірімдерді қолдану оқулықтар мен бағдарламалық жүйелерде терең тамырланған.



Шкала стандартного отклонения выборки

2.13.-сурет. Қайта іріктеумен госсет экспериментінің нәтижелері және орнатылған t-қисығы (оның "Биометрика" журналындағы жұмысынан, 1908 ж.)

Үлгі статистикасының мінез-құлқын сипаттаудағы t - үлестірудің дәлдігі осы үлгі үшін осы статистиканың таралуы қалыпты үлестіру түрінде болуын талап етеді. Іріктемелі Статистика көбінесе, тіпті базалық популяция

деректері болмаса да (t - үлестірімді кеңінен қолдануға әкелген факт) қалыпты түрде таратылатыны белгілі болды. Бұл құбылыс орталық шекті теорема деп аталады (қараңыз: разд. "Орталық шекті теорема" бұрын осы тарауда).

Биномдық таралу

Негізгі терминдер

Сынақ (сынақ) дискретті нәтижесі бар оқиға (мысалы, монетаны лақтыру).

Табыс (success) мақсатты сынақ нәтижесі. Синонимі: "1" (керісінше "0").

Биномдық (биномдық) екі нәтижеге ие. Синонимдер: иә/жоқ, 0/1, екілік, екілік, қосарланған.

Биномдық сынақ (биномдық сынақ) екі нәтижемен сынақ. Синонимі: бернуллиево сынағы.

Биномдық үлестіру (binomial distribution) x сынақтарындағы жетістіктер жиынтығын тарату. Синонимі: бернулдік таралу.

Биномдық нәтижелер, яғни иә/жоқ форматта, аналитиканың негізі болып табылады, өйткені олар көбінесе шешімнің немесе басқа процестің шарықтау шегі болып табылады; сатып алу/сатып алмау, басу/басу, тірі қалу/өлу және т. б. Биномдық таралуды түсінудің орталық мәні көптеген сынақтар идеясы болып табылады, мұнда әр сынақтың белгілі бір ықтималдығы бар екі мүмкін нәтижесі болады. Мысалы, монетаны 10 рет лақтыру-бұл 10 сынақтан тұратын биномдық эксперимент, онда әр сынақтың екі мүмкін нәтижесі бар (бүркіт немесе құйрық); суретті қараңыз. 2.14. Иә/жоқ немесе 0/1 форматындағы мұндай нәтижелер екілік деп аталады және олар міндетті түрде 50/50 ықтималдығы жоқ. Барлығы 1,0 болатын кез келген өзгерістер болуы мүмкін. Статистикада нәтижені " 1 " сәттілік деп атайды; жалпы қабылданған тәжірибе-бұл " 1 " Сирек кездесетін нәтижеге ие болу. "Сәттілік" терминін қолдану нәтиже қажет немесе пайдалы дегенді білдірмейді; іс жүзінде ол мақсатты нәтиже туралы айтады. Мысалы, қарыздарды қайтармау немесе алаяқтық транзакциялар салыстырмалы түрде сирек кездесетін жағдайлар болып табылады, оларды болжауға бізді қызықтыруы мүмкін, сондықтан оларға "1" немесе "сәттілік" атауы беріледі.



Рис. 2.14. Сторона с решкой бизоньего пятицентовика

Биномдық үлестіру дегеніміз-әр сынақта сәттіліктің (p) көрсетілген ықтималдығымен (n) берілген сынақ санындағы (x) табыстар санының жиіліктік таралуы.

Пуассон үлестірімі және басқа да қатысты үлестірімдер.

Көптеген процестер оқиғаларды берілген жалпы қарқындылықпен кездейсоқ ретпен тудырады — веб-сайтқа келген келушілер, ақы алу пунктіне келген көліктер (уақыт бойынша таралатын оқиғалар), матаның шаршы метріндегі кемшіліктер немесе 100 жолдық бағдарлама кодындағы қателер (оқиғалар) кеңістікте таралады).

Негізгі терминдер

Ламбда (λ) оқиғалар болатын қарқындылық (уақыт немесе кеңістік бірлігіне есептегенде).

Пуассонның таралуы (Poisson distribution) Таңдалған уақыт немесе кеңістік бірліктеріндегі оқиғалар санының жиіліктік таралуы.

Экспоненциалды үлестіру (exponential distribution) бір оқиғадан кейінгі оқиғаға дейінгі уақытты немесе қашықтықты жиіліктік үлестіру.

Вейбулдың таралуы (Weibull distribution) экспоненциалды үлестірудің жалпыланған нұсқасы, онда уақыт өте келе оқиғаның қарқындылығының өзгеруіне мүмкіндік береді.

Пуассонның Таралуы.

Алдыңғы мәліметтерге сүйене отырып, біз уақыт немесе кеңістік бірлігіне шаққандағы оқиғалардың орташа санын бағалай аламыз, бірақ біз оның уақыттың/кеңістіктің бір бірлігінен екіншісіне қаншалықты ерекшеленетінін білгіміз келеді. Пуассонның таралуы біз осындай бірліктерді таңдаған кезде уақыт немесе кеңістік бірлігіне шаққандағы оқиғалардың таралуы туралы айтады. Бұл жаппай қызмет көрсету туралы сұрақтарға жауап бергенде пайдалы, мысалы: "серверге кез-келген 5 секундтық кезеңде келетін интернет-трафикті толық өңдеуде 95% - ға арттыру үшін бізге қандай қуат қажет?" Пуассонды таратудағы негізгі параметр- λ немесе ламбда. Бұл белгілі бір уақыт немесе кеңістік аралығында болатын оқиғалардың орташа саны. Пуассон үлестірімінің дисперсиясы да λ -ге тең. Жалпы қабылданған әдіс кездейсоқ сандарды пуассонның таралуынан жаппай қызмет көрсетуді модельдеудің ажырамас бөлігі ретінде құрудан тұрады. R-де `pois` функциясы мұны тек екі аргументті — кездейсоқ сандар мен ламбдалардың санын алу арқылы жасайды

```
pois(100, lambda=2)
```

Бұл код үзіндісі Пуассон үлестірімінен 100 кездейсоқ сандарды шығарады, мұнда $2 \lambda =$. Мысалы, егер Тұтынушыларды Қолдау қызметіне кіретін қоңыраулардың орташа саны минутына 2 болса, онда бұл код үзіндісі 100 минутты құрайды, осы 100 минуттың әрқайсысына қоңыраулар санын қайтарады.

Экспоненциалды үлестіру.

Пуассо - на таратуда қолданған бірдей λ параметрін қолдана отырып, біз оқиғалар арасындағы уақытты бөлуді модельдеуге болады: веб-сайтқа кірулер арасындағы немесе автомобильдердің жол ақысын жинау орнына келуі арасындағы уақыт. Сондай-ақ, ол техникалық дизайнда жұмыс уақытын

модельдеу үшін және процестерді басқаруда, мысалы, сервистік қоңырауды есептеу кезінде қажет уақытты модельдеу үшін қолданылады.

Істен шығу қарқындылығын бағалау.

Көптеген қосымшаларда оқиғаның қарқындылығы λ алдыңғы деректерден белгілі немесе анықталуы мүмкін. Алайда, сирек жағдайларда бұл міндетті емес. Мысалы, авиациялық қозғалтқыштың істен шығуы өте сирек кездеседі (бақытымызға орай), сондықтан қозғалтқыштың осы түрі үшін сәтсіздіктер арасындағы уақытты бағалауға негізделетін мәліметтер аз болуы мүмкін. Деректер болмаған жағдайда, оқиғаның қарқындылығын бағалауға болатын база жоқ. Алайда, сіз қандай да бір болжам жасай аласыз: егер 20 сағаттан кейін ешқандай оқиғалар байқалмаса, онда қарқындылық сағатына 1 - ге тең емес екеніне сенімді бола аласыз. Модельдеу немесе ықтималдылықты тікелей есептеу арқылы әртүрлі гипотетикалық қарқындылықты анықтауға және қарқындылықтың төмендеуі екіталай болатын шекті мәндерді бағалауға болады. Егер мәліметтер аз болса, бірақ қарқындылықты дәл, сенімді бағалауды қамтамасыз ету үшін жеткіліксіз болса, онда әр түрлі қарқындылықтарға сәйкес келудің оңтайлылығын тексеру қолданылуы мүмкін (қараңыз: разд. 3-тараудың "хи-квадрат статистикасы негізінде тексеру"), олардың бақыланатын деректерге қаншалықты сәйкес келетінін анықтау.

Вейбуллдың Таралуы.

Көптеген жағдайларда оқиғаның қарқындылығы уақыт өте келе өзгермейді. Егер ол өзгертін кезең оқиғалар арасындағы әдеттегі аралықтан әлдеқайда ұзағырақ болса, онда ешқандай проблемалар болмайды; сіз талдауды бұрын айтылғандай қарқындылық салыстырмалы түрде тұрақты болатын сегменттерге бөлесіз. Егер оқиғаның қарқындылығы уақыт аралығында өзгерсе, онда бұрынғы немесе пуассонның таралуы пайдасыз болады. Бұл, ең алдымен, механикалық ақауларға қатысты болады — уақыт өте келе сәтсіздік

қаупі артады. Вейбуллдың таралуы экспоненциалды үлестірімнің кеңеюі болып табылады, онда оқиғаның қарқындылығын форма параметріне сәйкес өзгертуге болады, β . Егер $1 < \beta >$ болса, онда оқиғаның ықтималдығы уақыт бойынша артады, егер $1 < \beta <$ болса, онда ол азаяды. Вейбуллдың таралуы оқиғаның қарқындылығының орнына жұмыс уақытын талдаумен бірге қолданылатындықтан, екінші параметр интервалды есептеудегі оқиғалардың қарқындылығы тұрғысынан емес, өмірдің сипаттамалық уақыты (ресурстық сипаттама) тұрғысынан көрінеді. Мұнда η символы қолданылады, грек әрпі, ол масштаб параметрі немесе шкала² деп те аталады. Вейбуллдың таралуымен бағалау міндеті енді екі параметрді бағалауды қарастырады: β және η . Бағдарламалық жүйе деректерді модельдеу және желінің оңтайлы бөлінген үлестірімін бағалау үшін қолданылады. Вейбулл үлестірімінен кездейсоқ сандарды құру үшін R-ге код үзіндісі үш аргументті алады: n (пайда болатын сандар саны), $shape$ пішіні және масштаб шкаласы. Мысалы, кодтың келесі үзіндісі 1,5 формасы және 5000 өмір сүру уақыты бар Вейбулл үлестірімінен 100 кездейсоқ сандарды (өмір уақыттарын) жасайды:

```
rweibull(100,1.5,5000)
```