

The Book of Statistical Proofs

<https://statproofbook.github.io/>
StatProofBook@gmail.com

2020-02-13, 17:11

Contents

I	General Theorems	1
1	Probability theory	2
1.1	Probability distributions	2
1.1.1	<i>Moment-generating function</i>	2
1.2	Bayesian inference	2
1.2.1	Bayes' theorem	2
1.2.2	Bayes' rule	3
2	Estimation theory	5
2.1	Point estimates	5
2.1.1	Partition of the mean squared error into bias and variance	5
3	Information theory	7
3.1	Discrete mutual information	7
3.1.1	Relation to marginal and conditional entropy	7
3.1.2	Relation to marginal and joint entropy	8
3.1.3	Relation to joint and conditional entropy	9
II	Probability Distributions	11
1	Univariate discrete distributions	12
1.1	Bernoulli distribution	12
1.1.1	Mean	12
1.2	Binomial distribution	12
1.2.1	Mean	12
2	Multivariate discrete distributions	14
2.1	Categorical distribution	14
2.1.1	Mean	14
2.2	Multinomial distribution	14
2.2.1	Mean	14
3	Univariate continuous distributions	16
3.1	Continuous uniform distribution	16
3.1.1	<i>Definition</i>	16
3.1.2	Probability density function	16
3.1.3	Cumulative distribution function	17
3.1.4	Quantile function	18
3.2	Normal distribution	19
3.2.1	<i>Definition</i>	19

	3.2.2	Probability density function	19
	3.2.3	Mean	20
	3.2.4	Median	21
	3.2.5	Mode	22
	3.2.6	Variance	23
3.3		Gamma distribution	25
	3.3.1	<i>Definition</i>	25
	3.3.2	Probability density function	25
3.4		Exponential distribution	26
	3.4.1	<i>Definition</i>	26
	3.4.2	Probability density function	26
	3.4.3	Cumulative distribution function	27
	3.4.4	Quantile function	28
	3.4.5	Mean	28
	3.4.6	Median	29
	3.4.7	Mode	30
4		Multivariate continuous distributions	32
	4.1	Multivariate normal distribution	32
		4.1.1 <i>Definition</i>	32
		4.1.2 Probability density function	32
		4.1.3 Linear transformation theorem	33
		4.1.4 Marginal distributions	34
	4.2	Normal-gamma distribution	34
		4.2.1 <i>Definition</i>	34
		4.2.2 Probability density function	35
		4.2.3 Kullback-Leibler divergence	36
		4.2.4 Marginal distributions	37
5		Matrix-variate continuous distributions	41
	5.1	Matrix-normal distribution	41
		5.1.1 <i>Definition</i>	41
		5.1.2 Equivalence to multivariate normal distribution	41
III Statistical Models			43
1		Normal data	44
	1.1	Multiple linear regression	44
		1.1.1 Ordinary least squares (1)	44
		1.1.2 Ordinary least squares (2)	44
	1.2	Bayesian linear regression	45
		1.2.1 Conjugate prior distribution	45
		1.2.2 Posterior distribution	47
		1.2.3 Log model evidence	50
	1.3	General linear model	52
		1.3.1 Maximum likelihood estimation	52
2		Poisson data	55
	2.1	Poisson-distributed data	55
		2.1.1 Maximum likelihood estimation	55

2.2	Poisson distribution with exposure values	57
2.2.1	Conjugate prior distribution	57
2.2.2	Posterior distribution	58
2.2.3	Log model evidence	60
3	Probability data	63
3.1	Beta-distributed data	63
3.1.1	Method of moments	63
4	Categorical data	66
4.1	Binomial observations	66
4.1.1	Conjugate prior distribution	66
4.1.2	Posterior distribution	67
4.1.3	Log model evidence	68
IV Model Selection		71
1	Goodness-of-fit measures	72
1.1	R-squared	72
1.1.1	Derivation of R^2 and adjusted R^2	72
1.1.2	Relationship to maximum log-likelihood	73
2	Classical information criteria	76
2.1	Bayesian information criterion	76
2.1.1	Derivation	76
3	Bayesian model selection	78
3.1	Log model evidence	78
3.1.1	Derivation	78
3.1.2	Partition into accuracy and complexity	79

Chapter I

General Theorems

1 Probability theory

1.1 Probability distributions

1.1.1 Moment-generating function

Definition:

1) The moment-generating function of a random variable (\rightarrow Definition “rvar”) $X \in \mathbb{R}$ is

$$M_X(t) = \mathbb{E} \left[e^{tX} \right], \quad t \in \mathbb{R}. \quad (1)$$

2) The moment-generating function of a random vector (\rightarrow Definition “rvec”) $X \in \mathbb{R}^n$ is

$$M_X(t) = \mathbb{E} \left[e^{t^T X} \right], \quad t \in \mathbb{R}^n. \quad (2)$$

Sources:

- Wikipedia (2020): “Moment-generating function”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-22; URL: https://en.wikipedia.org/wiki/Moment-generating_function#Definition.

Metadata: ID: D2 | shortcut: mgf | author: JoramSoch | date: 2020-01-22, 10:58.

1.2 Bayesian inference

1.2.1 Bayes’ theorem

Theorem: Let A and B be two arbitrary statements about random variables (\rightarrow Definition “rvar”), such as statements about the presence or absence of an event or about the value of a scalar, vector or matrix. Then, the conditional probability that A is true, given that B is true, is equal to

$$p(A|B) = \frac{p(B|A) p(A)}{p(B)}. \quad (3)$$

Proof: The conditional probability (\rightarrow Definition “cp”) is defined as the ratio of joint probability (\rightarrow Definition “jp”), i.e. the probability of both statements being true, and marginal probability (\rightarrow Definition “mp”), i.e. the probability of only the second one being true:

$$p(A|B) = \frac{p(A, B)}{p(B)}. \quad (4)$$

It can also be written down for the reverse situation, i.e. to calculate the probability that B is true, given that A is true:

$$p(B|A) = \frac{p(A, B)}{p(A)}. \quad (5)$$

Both equations can be rearranged for the joint probability

$$p(A|B)p(B) \stackrel{(4)}{=} p(A, B) \stackrel{(5)}{=} p(B|A)p(A) \quad (6)$$

from which Bayes' theorem can be directly derived:

$$p(A|B) \stackrel{(6)}{=} \frac{p(B|A)p(A)}{p(B)} . \quad (7)$$

Sources:

- Koch, Karl-Rudolf (2007): “Rules of Probability”; in: *Introduction to Bayesian Statistics*, Springer, Berlin/Heidelberg, 2007, pp. 6/13, eqs. 2.12/2.38; URL: <https://www.springer.com/de/book/9783540727231>; DOI: 10.1007/978-3-540-72726-2.

Metadata: ID: P4 | shortcut: bayes-th | author: JoramSoch | date: 2019-09-27, 16:24.

1.2.2 Bayes' rule

Theorem: Let A_1 , A_2 and B be arbitrary statements about random variables (\rightarrow Definition “rvar”) where A_1 and A_2 are mutually exclusive. Then, Bayes' rule states that the posterior odds (\rightarrow Definition “post-odd”) are equal to the Bayes factor (\rightarrow Definition “bf”) times the prior odds (\rightarrow Definition “prior-odd”), i.e.

$$\frac{p(A_1|B)}{p(A_2|B)} = \frac{p(B|A_1)}{p(B|A_2)} \cdot \frac{p(A_1)}{p(A_2)} . \quad (8)$$

Proof: Using Bayes' theorem (\rightarrow Proof I/1.2.1), the conditional probabilities (\rightarrow Definition “cp”) on the left are given by

$$p(A_1|B) = \frac{p(B|A_1) \cdot p(A_1)}{p(B)} \quad (9)$$

$$p(A_2|B) = \frac{p(B|A_2) \cdot p(A_2)}{p(B)} . \quad (10)$$

Dividing the two conditional probabilities by each other

$$\begin{aligned} \frac{p(A_1|B)}{p(A_2|B)} &= \frac{p(B|A_1) \cdot p(A_1)/p(B)}{p(B|A_2) \cdot p(A_2)/p(B)} \\ &= \frac{p(B|A_1)}{p(B|A_2)} \cdot \frac{p(A_1)}{p(A_2)} , \end{aligned} \quad (11)$$

one obtains the posterior odds ratio as given by the theorem.

Sources:

- Wikipedia (2019): “Bayes' theorem”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-06; URL: https://en.wikipedia.org/wiki/Bayes%27_theorem#Bayes%E2%80%99_rule.

Metadata: ID: P12 | shortcut: bayes-rule | author: JoramSoch | date: 2020-01-06, 20:55.

2 Estimation theory

2.1 Point estimates

2.1.1 Partition of the mean squared error into bias and variance

Theorem: The mean squared error (\rightarrow Definition “mse”) can be partitioned into variance and squared bias

$$\text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}) - \text{Bias}(\hat{\theta}, \theta)^2 \quad (12)$$

where the variance (\rightarrow Definition “var”) is given by

$$\text{Var}(\hat{\theta}) = \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right)^2 \right] \quad (13)$$

and the bias (\rightarrow Definition “bias”) is given by

$$\text{Bias}(\hat{\theta}, \theta) = \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right) . \quad (14)$$

Proof: The mean squared error (MSE) is defined as (\rightarrow Definition “mse”) the expected value (\rightarrow Definition “ev”) of the squared deviation of the estimated value $\hat{\theta}$ from the true value θ of a parameter, over all values $\hat{\theta}$:

$$\text{MSE}(\hat{\theta}) = \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \theta \right)^2 \right] . \quad (15)$$

This formula can be evaluated in the following way:

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \theta \right)^2 \right] \\ &= \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) + \mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right)^2 \right] \\ &= \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right)^2 + 2 \left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right) \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right) + \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right)^2 \right] \\ &= \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right)^2 \right] + \mathbb{E}_{\hat{\theta}} \left[2 \left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right) \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right) \right] + \mathbb{E}_{\hat{\theta}} \left[\left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right)^2 \right] . \end{aligned} \quad (16)$$

Because $\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta$ is constant as a function of $\hat{\theta}$, we have:

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right)^2 \right] + 2 \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right) \mathbb{E}_{\hat{\theta}} \left[\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right] + \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right)^2 \\ &= \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right)^2 \right] + 2 \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right) \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right) + \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right)^2 \\ &= \mathbb{E}_{\hat{\theta}} \left[\left(\hat{\theta} - \mathbb{E}_{\hat{\theta}}(\hat{\theta}) \right)^2 \right] + \left(\mathbb{E}_{\hat{\theta}}(\hat{\theta}) - \theta \right)^2 . \end{aligned} \quad (17)$$

This proves the partition given by (12).

Sources:

- Wikipedia (2019): “Mean squared error”; in: *Wikipedia, the free encyclopedia*, retrieved on 2019-11-27; URL: https://en.wikipedia.org/wiki/Mean_squared_error#Proof_of_variance_and_bias_relationship.

Metadata: ID: P5 | shortcut: mse-bnv | author: JoramSoch | date: 2019-11-27, 14:26.

3 Information theory

3.1 Discrete mutual information

3.1.1 Relation to marginal and conditional entropy

Theorem: Let X and Y be discrete random variables (\rightarrow Definition “rvar”) with the joint probability (\rightarrow Definition “jp”) $p(x, y)$ for $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. Then, the mutual information (\rightarrow Definition “mi”) of X and Y can be expressed as

$$\begin{aligned} I(X, Y) &= H(X) - H(X|Y) \\ &= H(Y) - H(Y|X) \end{aligned} \quad (18)$$

where $H(X)$ and $H(Y)$ are the marginal entropies (\rightarrow Definition “ent-marg”) of X and Y and $H(X|Y)$ and $H(Y|X)$ are the conditional entropies (\rightarrow Definition “ent-cond”).

Proof: The mutual information (\rightarrow Definition “mi”) of X and Y is defined as

$$I(X, Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} . \quad (19)$$

Separating the logarithm, we have:

$$I(X, Y) = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(y)} - \sum_x \sum_y p(x, y) \log p(x) . \quad (20)$$

Applying the law of conditional probability (\rightarrow Proof “lcp”), i.e. $p(x, y) = p(x|y)p(y)$, we get:

$$I(X, Y) = \sum_x \sum_y p(x|y)p(y) \log p(x|y) - \sum_x \sum_y p(x, y) \log p(x) . \quad (21)$$

Regrouping the variables, we have:

$$I(X, Y) = \sum_y p(y) \sum_x p(x|y) \log p(x|y) - \sum_x \left(\sum_y p(x, y) \right) \log p(x) . \quad (22)$$

Applying the law of marginal probability (\rightarrow Proof “lmp”), i.e. $p(x) = \sum_y p(x, y)$, we get:

$$I(X, Y) = \sum_y p(y) \sum_x p(x|y) \log p(x|y) - \sum_x p(x) \log p(x) . \quad (23)$$

Now considering the definitions of marginal (\rightarrow Definition “ent-marg”) and conditional (\rightarrow Definition “ent-cond”) entropy

$$\begin{aligned} H(X) &= - \sum_{x \in \mathcal{X}} p(x) \log p(x) \\ H(X|Y) &= \sum_{y \in \mathcal{Y}} p(y) H(X|Y = y) , \end{aligned} \quad (24)$$

we can finally show:

$$\begin{aligned} I(X, Y) &= -H(X|Y) + H(X) \\ &= H(X) - H(X|Y) . \end{aligned} \quad (25)$$

The conditioning of X on Y in this proof is without loss of generality. Thus, the proof for the expression using the reverse conditional entropy of Y given X is obtained by simply switching x and y in the derivation.

Sources:

- Wikipedia (2020): “Mutual information”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-13; URL: https://en.wikipedia.org/wiki/Mutual_information#Relation_to_conditional_and_joint_entropy.

Metadata: ID: P19 | shortcut: dmi-mce | author: JoramSoch | date: 2020-01-13, 18:20.

3.1.2 Relation to marginal and joint entropy

Theorem: Let X and Y be discrete random variables (\rightarrow Definition “rvar”) with the joint probability (\rightarrow Definition “jp”) $p(x, y)$ for $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. Then, the mutual information (\rightarrow Definition “mi”) of X and Y can be expressed as

$$I(X, Y) = H(X) + H(Y) - H(X, Y) \quad (26)$$

where $H(X)$ and $H(Y)$ are the marginal entropies (\rightarrow Definition “ent-marg”) of X and Y and $H(X, Y)$ is the joint entropy (\rightarrow Definition “ent-joint”).

Proof: The mutual information (\rightarrow Definition “mi”) of X and Y is defined as

$$I(X, Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} . \quad (27)$$

Separating the logarithm, we have:

$$I(X, Y) = \sum_x \sum_y p(x, y) \log p(x, y) - \sum_x \sum_y p(x, y) \log p(x) - \sum_x \sum_y p(x, y) \log p(y) . \quad (28)$$

Regrouping the variables, this reads:

$$I(X, Y) = \sum_x \sum_y p(x, y) \log p(x, y) - \sum_x \left(\sum_y p(x, y) \right) \log p(x) - \sum_y \left(\sum_x p(x, y) \right) \log p(y) . \quad (29)$$

Applying the law of marginal probability (\rightarrow Proof “lmp”), i.e. $p(x) = \sum_y p(x, y)$, we get:

$$I(X, Y) = \sum_x \sum_y p(x, y) \log p(x, y) - \sum_x p(x) \log p(x) - \sum_y p(y) \log p(y) . \quad (30)$$

Now considering the definitions of marginal (\rightarrow Definition “ent-marg”) and joint (\rightarrow Definition “ent-joint”) entropy

$$\begin{aligned} H(X) &= - \sum_{x \in \mathcal{X}} p(x) \log p(x) \\ H(X, Y) &= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y) , \end{aligned} \quad (31)$$

we can finally show:

$$\begin{aligned} I(X, Y) &= -H(X, Y) + H(X) + H(Y) \\ &= H(X) + H(Y) - H(X, Y) . \end{aligned} \quad (32)$$

Sources:

- Wikipedia (2020): “Mutual information”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-13; URL: https://en.wikipedia.org/wiki/Mutual_information#Relation_to_conditional_and_joint_entropy.

Metadata: ID: P20 | shortcut: dmi-mje | author: JoramSoch | date: 2020-01-13, 21:53.

3.1.3 Relation to joint and conditional entropy

Theorem: Let X and Y be discrete random variables (\rightarrow Definition “rvar”) with the joint probability (\rightarrow Definition “jp”) $p(x, y)$ for $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. Then, the mutual information (\rightarrow Definition “mi”) of X and Y can be expressed as

$$I(X, Y) = H(X, Y) - H(X|Y) - H(Y|X) \quad (33)$$

where $H(X, Y)$ is the joint entropy (\rightarrow Definition “ent-joint”) of X and Y and $H(X | Y)$ and $H(Y | X)$ are the conditional entropies (\rightarrow Definition “ent-cond”).

Proof: The existence of the joint probability function ensures that the mutual information (\rightarrow Definition “mi”) is defined:

$$I(X, Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x) p(y)} . \quad (34)$$

The relation of mutual information to conditional entropy (\rightarrow Proof I/3.1.1) is:

$$I(X, Y) = H(X) - H(X|Y) \quad (35)$$

$$I(X, Y) = H(Y) - H(Y|X) \quad (36)$$

The relation of mutual information to joint entropy (\rightarrow Proof I/3.1.2) is:

$$I(X, Y) = H(X) + H(Y) - H(X, Y) . \quad (37)$$

It is true that

$$I(X, Y) = I(X, Y) + I(X, Y) - I(X, Y) . \quad (38)$$

Plugging in (35), (36) and (37) on the right-hand side, we have

$$\begin{aligned} I(X, Y) &= H(X) - H(X|Y) + H(Y) - H(Y|X) - H(X) - H(Y) + H(X, Y) \\ &= H(X, Y) - H(X|Y) - H(Y|X) \end{aligned} \quad (39)$$

which proves the identity given above.

Sources:

- Wikipedia (2020): “Mutual information”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-13; URL: https://en.wikipedia.org/wiki/Mutual_information#Relation_to_conditional_and_joint_entropy.

Metadata: ID: P21 | shortcut: dmi-jce | author: JoramSoch | date: 2020-01-13, 22:17.

Chapter II

Probability Distributions

1 Univariate discrete distributions

1.1 Bernoulli distribution

1.1.1 Mean

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a Bernoulli distribution (\rightarrow Definition “bern”):

$$X \sim \text{Bern}(p) . \quad (1)$$

Then, the mean or expected value (\rightarrow Definition “ev”) of X is

$$E(X) = p . \quad (2)$$

Proof: The expected value (\rightarrow Definition “ev”) is the probability-weighted average of all possible values:

$$E(X) = \sum_{x \in \mathcal{X}} x \cdot \Pr(X = x) . \quad (3)$$

Since there are only two possible outcomes for a Bernoulli random variable (\rightarrow Proof “bern-pmf”), we have:

$$\begin{aligned} E(X) &= 0 \cdot \Pr(X = 0) + 1 \cdot \Pr(X = 1) \\ &= 0 \cdot (1 - p) + 1 \cdot p \\ &= p . \end{aligned} \quad (4)$$

Sources:

- Wikipedia (2020): “Bernoulli distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-16; URL: https://en.wikipedia.org/wiki/Bernoulli_distribution#Mean.

Metadata: ID: P22 | shortcut: bern-mean | author: JoramSoch | date: 2020-01-16, 10:58.

1.2 Binomial distribution

1.2.1 Mean

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a binomial distribution (\rightarrow Definition “bin”):

$$X \sim \text{Bin}(n, p) . \quad (5)$$

Then, the mean or expected value (\rightarrow Definition “ev”) of X is

$$E(X) = np . \quad (6)$$

Proof: By definition, a binomial random variable (\rightarrow Definition “bin”) is the sum of n independent and identical Bernoulli trials (\rightarrow Definition “bern”) with success probability p . Therefore, the expected value is

$$E(X) = E(X_1 + \dots + X_n) \quad (7)$$

and because the expected value is a linear operator (\rightarrow Proof “ev-lin”), this is equal to

$$\begin{aligned} E(X) &= E(X_1) + \dots + E(X_n) \\ &= \sum_{i=1}^n E(X_i) . \end{aligned} \quad (8)$$

With the expected value of the Bernoulli distribution (\rightarrow Proof II/1.1.1), we have:

$$E(X) = \sum_{i=1}^n p = np . \quad (9)$$

Sources:

- Wikipedia (2020): “Binomial distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-16; URL: https://en.wikipedia.org/wiki/Binomial_distribution#Expected_value_and_variance.

Metadata: ID: P23 | shortcut: bin-mean | author: JoramSoch | date: 2020-01-16, 11:06.

2 Multivariate discrete distributions

2.1 Categorical distribution

2.1.1 Mean

Theorem: Let X be a random vector (\rightarrow Definition “rvec”) following a categorical distribution (\rightarrow Definition “cat”):

$$X \sim \text{Cat}([p_1, \dots, p_k]) . \quad (10)$$

Then, the mean or expected value (\rightarrow Definition “ev”) of X is

$$\text{E}(X) = [p_1, \dots, p_k] . \quad (11)$$

Proof: If we conceive the outcome of a categorical distribution (\rightarrow Definition “cat-pmf”) to be a $1 \times k$ vector, then the elementary row vectors $e_1 = [1, 0, \dots, 0], \dots, e_k = [0, \dots, 0, 1]$ are all the possible outcomes and they occur with probabilities $\text{Pr}(X = e_1) = p_1, \dots, \text{Pr}(X = e_k) = p_k$. Consequently, the expected value (\rightarrow Definition “ev”) is

$$\begin{aligned} \text{E}(X) &= \sum_{x \in \mathcal{X}} x \cdot \text{Pr}(X = x) \\ &= \sum_{i=1}^k e_i \cdot \text{Pr}(X = e_i) \\ &= \sum_{i=1}^k e_i \cdot p_i \\ &= [p_1, \dots, p_k] . \end{aligned} \quad (12)$$

Sources:

- original work

Metadata: ID: P24 | shortcut: cat-mean | author: JoramSoch | date: 2020-01-16, 11:17.

2.2 Multinomial distribution

2.2.1 Mean

Theorem: Let X be a random vector (\rightarrow Definition “rvec”) following a multinomial distribution (\rightarrow Definition “mult”):

$$X \sim \text{Mult}(n, [p_1, \dots, p_k]) . \quad (13)$$

Then, the mean or expected value (\rightarrow Definition “ev”) of X is

$$E(X) = [np_1, \dots, np_k] . \quad (14)$$

Proof: By definition, a multinomial random variable (\rightarrow Definition “mult”) is the sum of n independent and identical categorical trials (\rightarrow Definition “cat”) with category probabilities p_1, \dots, p_k . Therefore, the expected value is

$$E(X) = E(X_1 + \dots + X_n) \quad (15)$$

and because the expected value is a linear operator (\rightarrow Proof “ev-lin”), this is equal to

$$\begin{aligned} E(X) &= E(X_1) + \dots + E(X_n) \\ &= \sum_{i=1}^n E(X_i) . \end{aligned} \quad (16)$$

With the expected value of the categorical distribution (\rightarrow Proof II/2.1.1), we have:

$$E(X) = \sum_{i=1}^n [p_1, \dots, p_k] = n \cdot [p_1, \dots, p_k] = [np_1, \dots, np_k] . \quad (17)$$

Sources:

- original work

Metadata: ID: P25 | shortcut: mult-mean | author: JoramSoch | date: 2020-01-16, 11:26.

3 Univariate continuous distributions

3.1 Continuous uniform distribution

3.1.1 Definition

Definition: Let X be a continuous random variable (\rightarrow Definition “rvar”). Then, X is said to be uniformly distributed with minimum a and maximum b

$$X \sim \mathcal{U}(a, b) , \quad (18)$$

if and only if each value between and including a and b occurs with the same probability.

Sources:

- Wikipedia (2020): “Uniform distribution (continuous)”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-27; URL: [https://en.wikipedia.org/wiki/Uniform_distribution_\(continuous\)](https://en.wikipedia.org/wiki/Uniform_distribution_(continuous)).

Metadata: ID: D3 | shortcut: cuni | author: JoramSoch | date: 2020-01-27, 14:05.

3.1.2 Probability density function

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a continuous uniform distribution (\rightarrow Definition II/3.1.1):

$$X \sim \mathcal{U}(a, b) . \quad (19)$$

Then, the probability density function (\rightarrow Definition “pdf”) of X is

$$f_X(x) = \begin{cases} \frac{1}{b-a} , & \text{if } a \leq x \leq b \\ 0 , & \text{otherwise .} \end{cases} \quad (20)$$

Proof: A continuous uniform variable is defined as (\rightarrow Definition II/3.1.1) having a constant probability density between minimum a and maximum b . Therefore,

$$\begin{aligned} f_X(x) &\propto 1 \quad \text{for all } x \in [a, b] \quad \text{and} \\ f_X(x) &= 0, \quad \text{if } x < a \quad \text{or } x > b . \end{aligned} \quad (21)$$

To ensure that $f_X(x)$ is a proper probability density function (\rightarrow Definition “pdf”), the integral over all non-zero probabilities has to sum to 1. Therefore,

$$f_X(x) = \frac{1}{c(a, b)} \quad \text{for all } x \in [a, b] \quad (22)$$

where the normalization factor $c(a, b)$ is specified, such that

$$\frac{1}{c(a, b)} \int_a^b 1 \, dx = 1 . \quad (23)$$

Solving this for $c(a, b)$, we obtain:

$$\begin{aligned} \int_a^b 1 \, dx &= c(a, b) \\ [x]_a^b &= c(a, b) \\ c(a, b) &= b - a . \end{aligned} \tag{24}$$

Sources:

- original work

Metadata: ID: P37 | shortcut: cuni-pdf | author: JoramSoch | date: 2020-01-31, 15:41.

3.1.3 Cumulative distribution function

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a continuous uniform distribution (\rightarrow Definition II/3.1.1):

$$X \sim \mathcal{U}(a, b) . \tag{25}$$

Then, the cumulative distribution function (\rightarrow Definition “cdf”) of X is

$$F_X(x) = \begin{cases} 0 , & \text{if } x < a \\ \frac{x-a}{b-a} , & \text{if } a \leq x \leq b \\ 1 , & \text{if } x > b . \end{cases} \tag{26}$$

Proof: The probability density function of the continuous uniform distribution (\rightarrow Proof II/3.1.2) is:

$$\mathcal{U}(z; a, b) = \begin{cases} \frac{1}{b-a} , & \text{if } a \leq z \leq b \\ 0 , & \text{otherwise} . \end{cases} \tag{27}$$

Thus, the cumulative distribution function (\rightarrow Definition “cdf”) is:

$$F_X(x) = \int_{-\infty}^x \mathcal{U}(z; a, b) \, dz \tag{28}$$

First of all, if $x < a$, we have

$$F_X(x) = \int_{-\infty}^x 0 \, dz = 0 . \tag{29}$$

Moreover, if $a \leq x \leq b$, we have using (27)

$$\begin{aligned}
F_X(x) &= \int_{-\infty}^a \mathcal{U}(z; a, b) \, dz + \int_a^x \mathcal{U}(z; a, b) \, dz \\
&= \int_{-\infty}^a 0 \, dz + \int_a^x \frac{1}{b-a} \, dz \\
&= 0 + \frac{1}{b-a} [z]_a^x \\
&= \frac{x-a}{b-a} .
\end{aligned} \tag{30}$$

Finally, if $x > b$, we have

$$\begin{aligned}
F_X(x) &= \int_{-\infty}^b \mathcal{U}(z; a, b) \, dz + \int_b^x \mathcal{U}(z; a, b) \, dz \\
&= F_X(b) + \int_b^x 0 \, dz \\
&= \frac{b-a}{b-a} + 0 \\
&= 1 .
\end{aligned} \tag{31}$$

This completes the proof.

Sources:

- original work

Metadata: ID: P38 | shortcut: cuni-cdf | author: JoramSoch | date: 2020-01-02, 18:05.

3.1.4 Quantile function

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a continuous uniform distribution (\rightarrow Definition II/3.1.1):

$$X \sim \mathcal{U}(a, b) . \tag{32}$$

Then, the quantile function (\rightarrow Definition “qf”) of X is

$$Q_X(p) = bp + a(1-p) . \tag{33}$$

Proof: The cumulative distribution function of the continuous uniform distribution (\rightarrow Proof II/3.1.3) is:

$$F_X(x) = \begin{cases} 0 , & \text{if } x < a \\ \frac{x-a}{b-a} , & \text{if } a \leq x \leq b \\ 1 , & \text{if } x > b . \end{cases} \tag{34}$$

Thus, the quantile function (\rightarrow Definition “qf”) is:

$$Q_X(p) = F_X^{-1}(x) . \quad (35)$$

This can be derived by rearranging equation (34):

$$\begin{aligned} p &= \frac{x - a}{b - a} \\ x &= p(b - a) + a \\ x &= bp + a(1 - p) = Q_X(p) . \end{aligned} \quad (36)$$

Sources:

- original work

Metadata: ID: P39 | shortcut: cuni-qf | author: JoramSoch | date: 2020-01-02, 18:27.

3.2 Normal distribution

3.2.1 Definition

Definition: Let X be a random variable (\rightarrow Definition “rvar”). Then, X is said to be normally distributed with mean μ and variance σ^2 (or, standard deviation σ)

$$X \sim \mathcal{N}(\mu, \sigma^2) , \quad (37)$$

if and only if its probability density function (\rightarrow Definition “pdf”) is given by

$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right] \quad (38)$$

where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$.

Sources:

- Wikipedia (2020): “Normal distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-27; URL: https://en.wikipedia.org/wiki/Normal_distribution.

Metadata: ID: D4 | shortcut: norm | author: JoramSoch | date: 2020-01-27, 14:15.

3.2.2 Probability density function

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a normal distribution (\rightarrow Definition II/3.2.1):

$$X \sim \mathcal{N}(\mu, \sigma^2) . \quad (39)$$

Then, the probability density function (\rightarrow Definition “pdf”) of X is

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right] . \quad (40)$$

Proof: This follows directly from the definition of the normal distribution (\rightarrow Definition II/3.2.1).

Sources:

- original work

Metadata: ID: P33 | shortcut: norm-pdf | author: JoramSoch | date: 2020-01-27, 15:15.

3.2.3 Mean

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a normal distribution (\rightarrow Definition II/3.2.1):

$$X \sim \mathcal{N}(\mu, \sigma^2) . \quad (41)$$

Then, the mean or expected value (\rightarrow Definition “ev”) of X is

$$\mathbb{E}(X) = \mu . \quad (42)$$

Proof: The expected value (\rightarrow Definition “ev”) is the probability-weighted average over all possible values:

$$\mathbb{E}(X) = \int_{\mathbb{R}} x \cdot f_X(x) \, dx . \quad (43)$$

With the probability density function of the normal distribution (\rightarrow Proof II/3.2.2), this reads:

$$\begin{aligned} \mathbb{E}(X) &= \int_{-\infty}^{+\infty} x \cdot \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right] \, dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} x \cdot \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right] \, dx . \end{aligned} \quad (44)$$

Substituting $z = x - \mu$, we have:

$$\begin{aligned}
E(X) &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty-\mu}^{+\infty-\mu} (z + \mu) \cdot \exp \left[-\frac{1}{2} \left(\frac{z}{\sigma} \right)^2 \right] dz \\
&= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} (z + \mu) \cdot \exp \left[-\frac{1}{2} \left(\frac{z}{\sigma} \right)^2 \right] dz \\
&= \frac{1}{\sqrt{2\pi}\sigma} \left(\int_{-\infty}^{+\infty} z \cdot \exp \left[-\frac{1}{2} \left(\frac{z}{\sigma} \right)^2 \right] dz + \mu \int_{-\infty}^{+\infty} \exp \left[-\frac{1}{2} \left(\frac{z}{\sigma} \right)^2 \right] dz \right) \\
&= \frac{1}{\sqrt{2\pi}\sigma} \left(\int_{-\infty}^{+\infty} z \cdot \exp \left[-\frac{1}{2\sigma^2} \cdot z^2 \right] dz + \mu \int_{-\infty}^{+\infty} \exp \left[-\frac{1}{2\sigma^2} \cdot z^2 \right] dz \right) .
\end{aligned} \tag{45}$$

The general antiderivatives are

$$\begin{aligned}
\int x \cdot \exp [-ax^2] dx &= -\frac{1}{2a} \cdot \exp [-ax^2] \\
\int \exp [-ax^2] dx &= \frac{1}{2} \sqrt{\frac{\pi}{a}} \cdot \operatorname{erf} [\sqrt{a}x]
\end{aligned} \tag{46}$$

where $\operatorname{erf}(x)$ is the error function. Using this, the integrals can be calculated as:

$$\begin{aligned}
E(X) &= \frac{1}{\sqrt{2\pi}\sigma} \left(\left[-\sigma^2 \cdot \exp \left[-\frac{1}{2\sigma^2} \cdot z^2 \right] \right]_{-\infty}^{+\infty} + \mu \left[\sqrt{\frac{\pi}{2}} \sigma \cdot \operatorname{erf} \left[\frac{1}{\sqrt{2}\sigma} z \right] \right]_{-\infty}^{+\infty} \right) \\
&= \frac{1}{\sqrt{2\pi}\sigma} \left(\left(\lim_{z \rightarrow +\infty} \left(-\sigma^2 \cdot \exp \left[-\frac{1}{2\sigma^2} \cdot z^2 \right] \right) - \lim_{z \rightarrow -\infty} \left(-\sigma^2 \cdot \exp \left[-\frac{1}{2\sigma^2} \cdot z^2 \right] \right) \right) \right. \\
&\quad \left. + \mu \left(\lim_{z \rightarrow +\infty} \left(\sqrt{\frac{\pi}{2}} \sigma \cdot \operatorname{erf} \left[\frac{1}{\sqrt{2}\sigma} z \right] \right) - \lim_{z \rightarrow -\infty} \left(\sqrt{\frac{\pi}{2}} \sigma \cdot \operatorname{erf} \left[\frac{1}{\sqrt{2}\sigma} z \right] \right) \right) \right) \\
&= \frac{1}{\sqrt{2\pi}\sigma} \left([0 - 0] + \mu \left[\sqrt{\frac{\pi}{2}} \sigma - \left(-\sqrt{\frac{\pi}{2}} \sigma \right) \right] \right) \\
&= \frac{1}{\sqrt{2\pi}\sigma} \cdot \mu \cdot 2\sqrt{\frac{\pi}{2}} \sigma \\
&= \mu .
\end{aligned} \tag{47}$$

Sources:

- Papadopoulos, Alecos (2013): “How to derive the mean and variance of Gaussian random variable?”; in: *StackExchange Mathematics*; URL: <https://math.stackexchange.com/questions/518281/how-to-derive-the-mean-and-variance-of-a-gaussian-random-variable>.

Metadata: ID: P15 | shortcut: norm-mean | author: JoramSoch | date: 2020-01-09, 15:04.

3.2.4 Median

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a normal distribution (\rightarrow Definition II/3.2.1):

$$X \sim \mathcal{N}(\mu, \sigma^2) . \quad (48)$$

Then, the median (\rightarrow Definition “med”) of X is

$$\text{median}(X) = \mu . \quad (49)$$

Proof: The median (\rightarrow Definition “med”) is the value at which the cumulative distribution function (\rightarrow Definition “cdf”) is $1/2$:

$$F_X(\text{median}(X)) = \frac{1}{2} . \quad (50)$$

The cumulative distribution function of the normal distribution (\rightarrow Proof “norm-cdf”) is

$$F_X(x) = \frac{1}{2} \left[1 + \text{erf} \left(\frac{x - \mu}{\sqrt{2}\sigma} \right) \right] \quad (51)$$

where $\text{erf}(x)$ is the error function. Thus, the inverse CDF is

$$x = \sqrt{2}\sigma \cdot \text{erf}^{-1}(2p - 1) + \mu \quad (52)$$

where $\text{erf}^{-1}(x)$ is the inverse error function. Setting $p = 1/2$, we obtain:

$$\text{median}(X) = \sqrt{2}\sigma \cdot \text{erf}^{-1}(0) + \mu = \mu . \quad (53)$$

Sources:

- original work

Metadata: ID: P16 | shortcut: norm-med | author: JoramSoch | date: 2020-01-09, 15:33.

3.2.5 Mode

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a normal distribution (\rightarrow Definition II/3.2.1):

$$X \sim \mathcal{N}(\mu, \sigma^2) . \quad (54)$$

Then, the mode (\rightarrow Definition “mode”) of X is

$$\text{mode}(X) = \mu . \quad (55)$$

Proof: The mode (\rightarrow Definition “mode”) is the value which maximizes the probability density function (\rightarrow Definition “pdf”):

$$\text{mode}(X) = \arg \max_x f_X(x) . \quad (56)$$

The probability density function of the normal distribution (\rightarrow Proof II/3.2.2) is:

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp \left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right] . \quad (57)$$

The first two derivatives of this function are:

$$f'_X(x) = \frac{df_X(x)}{dx} = \frac{1}{\sqrt{2\pi}\sigma^3} \cdot (-x + \mu) \cdot \exp \left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right] \quad (58)$$

$$f''_X(x) = \frac{d^2f_X(x)}{dx^2} = -\frac{1}{\sqrt{2\pi}\sigma^3} \cdot \exp \left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right] + \frac{1}{\sqrt{2\pi}\sigma^5} \cdot (-x + \mu)^2 \cdot \exp \left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right] . \quad (59)$$

We now calculate the root of the first derivative (58):

$$\begin{aligned} f'_X(x) = 0 &= \frac{1}{\sqrt{2\pi}\sigma^3} \cdot (-x + \mu) \cdot \exp \left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right] \\ 0 &= -x + \mu \\ x &= \mu . \end{aligned} \quad (60)$$

By plugging this value into the second derivative (59),

$$\begin{aligned} f''_X(\mu) &= -\frac{1}{\sqrt{2\pi}\sigma^3} \cdot \exp(0) + \frac{1}{\sqrt{2\pi}\sigma^5} \cdot (0)^2 \cdot \exp(0) \\ &= -\frac{1}{\sqrt{2\pi}\sigma^3} < 0 , \end{aligned} \quad (61)$$

we confirm that it is in fact a maximum which shows that

$$\text{mode}(X) = \mu . \quad (62)$$

Sources:

- original work

Metadata: ID: P17 | shortcut: norm-mode | author: JoramSoch | date: 2020-01-09, 15:58.

3.2.6 Variance

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following a normal distribution (\rightarrow Definition II/3.2.1):

$$X \sim \mathcal{N}(\mu, \sigma^2) . \quad (63)$$

Then, the variance (\rightarrow Definition “var”) of X is

$$\text{Var}(X) = \sigma^2 . \quad (64)$$

Proof: The variance (\rightarrow Definition “var”) is the probability-weighted average of the squared deviation from the mean (\rightarrow Definition “ev”):

$$\text{Var}(X) = \int_{\mathbb{R}} (x - E(X))^2 \cdot f_X(x) \, dx . \quad (65)$$

With the expected value (\rightarrow Proof II/3.2.3) and probability density function (\rightarrow Proof II/3.2.2) of the normal distribution, this reads:

$$\begin{aligned} \text{Var}(X) &= \int_{-\infty}^{+\infty} (x - \mu)^2 \cdot \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right] \, dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} (x - \mu)^2 \cdot \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right] \, dx . \end{aligned} \quad (66)$$

Substituting $z = x - \mu$, we have:

$$\begin{aligned} \text{Var}(X) &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} z^2 \cdot \exp \left[-\frac{1}{2} \left(\frac{z}{\sigma} \right)^2 \right] \, d(z + \mu) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} z^2 \cdot \exp \left[-\frac{1}{2} \left(\frac{z}{\sigma} \right)^2 \right] \, dz . \end{aligned} \quad (67)$$

Now substituting $z = \sqrt{2}\sigma x$, we have:

$$\begin{aligned} \text{Var}(X) &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} (\sqrt{2}\sigma x)^2 \cdot \exp \left[-\frac{1}{2} \left(\frac{\sqrt{2}\sigma x}{\sigma} \right)^2 \right] \, d(\sqrt{2}\sigma x) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \cdot 2\sigma^2 \cdot \sqrt{2}\sigma \int_{-\infty}^{+\infty} x^2 \cdot \exp [-x^2] \, dx \\ &= \frac{2\sigma^2}{\sqrt{\pi}} \int_{-\infty}^{+\infty} x^2 \cdot e^{-x^2} \, dx . \end{aligned} \quad (68)$$

Since the integrand is symmetric with respect to $x = 0$, we can write:

$$\text{Var}(X) = \frac{4\sigma^2}{\sqrt{\pi}} \int_0^{\infty} x^2 \cdot e^{-x^2} \, dx . \quad (69)$$

If we define $z = x^2$, then $x = \sqrt{z}$ and $dx = 1/2 z^{-1/2} dz$. Substituting this into the integral

$$\text{Var}(X) = \frac{4\sigma^2}{\sqrt{\pi}} \int_0^{\infty} z \cdot e^{-z} \cdot \frac{1}{2} z^{-1/2} \, dz = \frac{2\sigma^2}{\sqrt{\pi}} \int_0^{\infty} z^{\frac{3}{2}-1} \cdot e^{-z} \, dz \quad (70)$$

and using the definition of the gamma function

$$\Gamma(x) = \int_0^\infty z^{x-1} \cdot e^{-z} dz , \quad (71)$$

we can finally show that

$$\text{Var}(X) = \frac{2\sigma^2}{\sqrt{\pi}} \cdot \Gamma\left(\frac{3}{2}\right) = \frac{2\sigma^2}{\sqrt{\pi}} \cdot \frac{\sqrt{\pi}}{2} = \sigma^2 . \quad (72)$$

Sources:

- Papadopoulos, Alecos (2013): “How to derive the mean and variance of Gaussian random variable?”; in: *StackExchange Mathematics*; URL: <https://math.stackexchange.com/questions/518281/how-to-derive-the-mean-and-variance-of-a-gaussian-random-variable>.

Metadata: ID: P18 | shortcut: norm-var | author: JoramSoch | date: 2020-01-09, 22:47.

3.3 Gamma distribution

3.3.1 Definition

****Definition**:** Let X be a random variable (\rightarrow Definition “rvar”). Then, X is said to follow a gamma distribution with shape a and rate b

$$X \sim \text{Gam}(a, b) , \quad (73)$$

if and only if its probability density function (\rightarrow Definition “pdf”) is given by

$$\text{Gam}(x; a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} \exp[-bx], \quad x > 0 \quad (74)$$

where $a > 0$ and $b > 0$, and the density is zero, if $x \leq 0$.

Sources:

- Koch, Karl-Rudolf (2007): “Gamma Distribution”; in: *Introduction to Bayesian Statistics*, Springer, Berlin/Heidelberg, 2007, p. 47, eq. 2.172; URL: <https://www.springer.com/de/book/9783540727231>; DOI: 10.1007/978-3-540-72726-2.

Metadata: ID: D7 | shortcut: gam | author: JoramSoch | date: 2020-02-08, 23:29.

3.3.2 Probability density function

Theorem: Let X be a positive random variable (\rightarrow Definition “rvar”) following a gamma distribution (\rightarrow Definition II/3.3.1):

$$X \sim \text{Gam}(a, b) . \quad (75)$$

Then, the probability density function (\rightarrow Definition “pdf”) of X is

$$f_X(x) = \frac{b^a}{\Gamma(a)} x^{a-1} \exp[-bx] . \quad (76)$$

Proof: This follows directly from the definition of the gamma distribution (\rightarrow Definition II/3.3.1).

Sources:

- original work

Metadata: ID: P45 | shortcut: gam-pdf | author: JoramSoch | date: 2020-02-08, 23:41.

3.4 Exponential distribution

3.4.1 Definition

****Definition**:** Let X be a random variable (\rightarrow Definition “rvar”). Then, X is said to be exponentially distributed with rate (or, inverse scale) λ

$$X \sim \text{Exp}(\lambda) , \quad (77)$$

if and only if its probability density function (\rightarrow Definition “pdf”) is given by

$$\text{Exp}(x; \lambda) = \lambda \exp[-\lambda x], \quad x \geq 0 \quad (78)$$

where $\lambda > 0$, and the density is zero, if $x < 0$.

Sources:

- Wikipedia (2020): “Exponential distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-02-08; URL: https://en.wikipedia.org/wiki/Exponential_distribution#Definitions.

Metadata: ID: D8 | shortcut: exp | author: JoramSoch | date: 2020-02-08, 23:48.

3.4.2 Probability density function

Theorem: Let X be a non-negative random variable (\rightarrow Definition “rvar”) following an exponential distribution (\rightarrow Definition II/3.4.1):

$$X \sim \text{Exp}(\lambda) . \quad (79)$$

Then, the probability density function (\rightarrow Definition “pdf”) of X is

$$f_X(x) = \lambda \exp[-\lambda x] . \quad (80)$$

Proof: This follows directly from the definition of the exponential distribution (\rightarrow Definition II/3.4.1).

Sources:

- original work

Metadata: ID: P46 | shortcut: exp-pdf | author: JoramSoch | date: 2020-02-08, 23:53.

3.4.3 Cumulative distribution function

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following an exponential distribution (\rightarrow Definition II/3.4.1):

$$X \sim \text{Exp}(\lambda) . \quad (81)$$

Then, the cumulative distribution function (\rightarrow Definition “cdf”) of X is

$$F_X(x) = \begin{cases} 0 , & \text{if } x < 0 \\ 1 - \exp[-\lambda x] , & \text{if } x \geq 0 . \end{cases} \quad (82)$$

Proof: The probability density function of the exponential distribution (\rightarrow Proof II/3.4.2) is:

$$\text{Exp}(x; \lambda) = \begin{cases} 0 , & \text{if } x < 0 \\ \lambda \exp[-\lambda x] , & \text{if } x \geq 0 . \end{cases} \quad (83)$$

Thus, the cumulative distribution function (\rightarrow Definition “cdf”) is:

$$F_X(x) = \int_{-\infty}^x \text{Exp}(z; \lambda) \, dz . \quad (84)$$

If $x < 0$, we have:

$$F_X(x) = \int_{-\infty}^x 0 \, dz = 0 . \quad (85)$$

If $x \geq 0$, we have using (83):

$$\begin{aligned} F_X(x) &= \int_{-\infty}^0 \text{Exp}(z; \lambda) \, dz + \int_0^x \text{Exp}(z; \lambda) \, dz \\ &= \int_{-\infty}^0 0 \, dz + \int_0^x \lambda \exp[-\lambda z] \, dz \\ &= 0 + \lambda \left[-\frac{1}{\lambda} \exp[-\lambda z] \right]_0^x \\ &= \lambda \left[\left(-\frac{1}{\lambda} \exp[-\lambda x] \right) - \left(-\frac{1}{\lambda} \exp[-\lambda \cdot 0] \right) \right] \\ &= 1 - \exp[-\lambda x] . \end{aligned} \quad (86)$$

Sources:

- original work

Metadata: ID: P48 | shortcut: exp-cdf | author: JoramSoch | date: 2020-02-11, 14:48.

3.4.4 Quantile function

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following an exponential distribution (\rightarrow Definition II/3.4.1):

$$X \sim \text{Exp}(\lambda) . \quad (87)$$

Then, the quantile function (\rightarrow Definition “qf”) of X is

$$Q_X(p) = -\frac{\ln(1-p)}{\lambda} . \quad (88)$$

Proof: The cumulative distribution function of the exponential distribution (\rightarrow Proof II/3.4.3) is:

$$F_X(x) = \begin{cases} 0 , & \text{if } x < 0 \\ 1 - \exp[-\lambda x] , & \text{if } x \geq 0 . \end{cases} \quad (89)$$

Thus, the quantile function (\rightarrow Definition “qf”) is:

$$Q_X(p) = F_X^{-1}(x) . \quad (90)$$

This can be derived by rearranging equation (89):

$$\begin{aligned} p &= 1 - \exp[-\lambda x] \\ \exp[-\lambda x] &= 1 - p \\ -\lambda x &= \ln(1-p) \\ x &= -\frac{\ln(1-p)}{\lambda} . \end{aligned} \quad (91)$$

Sources:

- original work

Metadata: ID: P50 | shortcut: exp-qf | author: JoramSoch | date: 2020-02-12, 15:48.

3.4.5 Mean

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following an exponential distribution (\rightarrow Definition II/3.4.1):

$$X \sim \text{Exp}(\lambda) . \quad (92)$$

Then, the mean or expected value (\rightarrow Definition “ev”) of X is

$$E(X) = \frac{1}{\lambda} . \quad (93)$$

Proof: The expected value (\rightarrow Definition “ev”) is the probability-weighted average over all possible values:

$$E(X) = \int_{\mathbb{R}} x \cdot f_X(x) \, dx . \quad (94)$$

With the probability density function of the exponential distribution (\rightarrow Proof II/3.4.2), this reads:

$$\begin{aligned} E(X) &= \int_0^{+\infty} x \cdot \lambda \exp(-\lambda x) \, dx \\ &= \lambda \int_0^{+\infty} x \cdot \exp(-\lambda x) \, dx . \end{aligned} \quad (95)$$

Using the following anti-derivative

$$\int x \cdot \exp(-\lambda x) \, dx = \left(-\frac{1}{\lambda}x - \frac{1}{\lambda^2} \right) \exp(-\lambda x) , \quad (96)$$

the expected value becomes

$$\begin{aligned} E(X) &= \lambda \left[\left(-\frac{1}{\lambda}x - \frac{1}{\lambda^2} \right) \exp(-\lambda x) \right]_0^{+\infty} \\ &= \lambda \left[\lim_{x \rightarrow \infty} \left(-\frac{1}{\lambda}x - \frac{1}{\lambda^2} \right) \exp(-\lambda x) - \left(-\frac{1}{\lambda} \cdot 0 - \frac{1}{\lambda^2} \right) \exp(-\lambda \cdot 0) \right] \\ &= \lambda \left[0 + \frac{1}{\lambda^2} \right] \\ &= \frac{1}{\lambda} . \end{aligned} \quad (97)$$

Sources:

- Koch, Karl-Rudolf (2007): “Expected Value”; in: *Introduction to Bayesian Statistics*, Springer, Berlin/Heidelberg, 2007, p. 39, eq. 2.142a; URL: <https://www.springer.com/de/book/9783540727231>; DOI: 10.1007/978-3-540-72726-2.

Metadata: ID: P47 | shortcut: exp-mean | author: JoramSoch | date: 2020-02-10, 21:57.

3.4.6 Median

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following an exponential distribution (\rightarrow Definition II/3.4.1):

$$X \sim \text{Exp}(\lambda) . \quad (98)$$

Then, the median (\rightarrow Definition “med”) of X is

$$\text{median}(X) = \frac{\ln 2}{\lambda} . \quad (99)$$

Proof: The median (\rightarrow Definition “med”) is the value at which the cumulative distribution function (\rightarrow Definition “cdf”) is 1/2:

$$F_X(\text{median}(X)) = \frac{1}{2} . \quad (100)$$

The cumulative distribution function of the exponential distribution (\rightarrow Definition “exp-cdf”) is

$$F_X(x) = 1 - \exp[-\lambda x], \quad x \geq 0 . \quad (101)$$

Thus, the inverse CDF is

$$x = -\frac{\ln(1-p)}{\lambda} \quad (102)$$

and setting $p = 1/2$, we obtain:

$$\text{median}(X) = -\frac{\ln(1-\frac{1}{2})}{\lambda} = \frac{\ln 2}{\lambda} . \quad (103)$$

Sources:

- original work

Metadata: ID: P49 | shortcut: exp-med | author: JoramSoch | date: 2020-02-11, 15:03.

3.4.7 Mode

Theorem: Let X be a random variable (\rightarrow Definition “rvar”) following an exponential distribution (\rightarrow Definition II/3.4.1):

$$X \sim \text{Exp}(\lambda) . \quad (104)$$

Then, the mode (\rightarrow Definition “mode”) of X is

$$\text{mode}(X) = 0 . \quad (105)$$

Proof: The mode (\rightarrow Definition “mode”) is the value which maximizes the probability density function (\rightarrow Definition “pdf”):

$$\text{mode}(X) = \arg \max_x f_X(x) . \quad (106)$$

The probability density function of the exponential distribution (\rightarrow Proof II/3.4.2) is:

$$f_X(x) = \begin{cases} 0 , & \text{if } x < 0 \\ \lambda \exp[-\lambda x] , & \text{if } x \geq 0 . \end{cases} \quad (107)$$

Since

$$\lim_{x \rightarrow 0} f_X(x) = \infty \quad (108)$$

and

$$f_X(x) < \infty \quad \text{for any } x \neq 0, \quad (109)$$

it follows that

$$\text{mode}(X) = 0. \quad (110)$$

Sources:

- original work

Metadata: ID: P51 | shortcut: exp-mode | author: JoramSoch | date: 2020-02-12, 15:53.

4 Multivariate continuous distributions

4.1 Multivariate normal distribution

4.1.1 Definition

Definition: Let X be an $n \times 1$ random vector (\rightarrow Definition “rvec”). Then, X is said to be multivariate normally distributed with mean μ and covariance Σ

$$X \sim \mathcal{N}(\mu, \Sigma) , \quad (111)$$

if and only if its probability density function (\rightarrow Definition “pdf”) is given by

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \cdot \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right] \quad (112)$$

where μ is an $n \times 1$ real vector and Σ is an $n \times n$ positive definite matrix.

Sources:

- Koch KR (2007): “Multivariate Normal Distribution”; in: *Introduction to Bayesian Statistics*, ch. 2.5.1, pp. 51-53, eq. 2.195; URL: <https://www.springer.com/gp/book/9783540727231>; DOI: 10.1007/978-3-540-72726-2.

Metadata: ID: D1 | shortcut: mvn | author: JoramSoch | date: 2020-01-22, 05:20.

4.1.2 Probability density function

Theorem: Let X be a random vector (\rightarrow Definition “rvec”) following a multivariate normal distribution (\rightarrow Definition II/4.1.1):

$$X \sim \mathcal{N}(\mu, \Sigma) . \quad (113)$$

Then, the probability density function (\rightarrow Definition “pdf”) of X is

$$f_X(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \cdot \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right] . \quad (114)$$

Proof: This follows directly from the definition of the multivariate normal distribution (\rightarrow Definition II/4.1.1).

Sources:

- original work

Metadata: ID: P34 | shortcut: mvn-pdf | author: JoramSoch | date: 2020-01-27, 15:23.

4.1.3 Linear transformation theorem

Theorem: Let x follow a multivariate normal distribution (\rightarrow Definition II/4.1.1):

$$x \sim \mathcal{N}(\mu, \Sigma) . \quad (115)$$

Then, any linear transformation of x is also multivariate normally distributed:

$$y = Ax + b \sim \mathcal{N}(A\mu + b, A\Sigma A^T) . \quad (116)$$

Proof: The moment-generating function of a random vector (\rightarrow Definition I/1.1.1) x is

$$M_x(t) = \mathbb{E} \left(\exp [t^T x] \right) \quad (117)$$

and therefore the moment-generating function of the random vector y is given by

$$\begin{aligned} M_y(t) &= \mathbb{E} \left(\exp [t^T (Ax + b)] \right) \\ &= \mathbb{E} \left(\exp [t^T Ax] \cdot \exp [t^T b] \right) \\ &= \exp [t^T b] \cdot \mathbb{E} \left(\exp [t^T Ax] \right) \\ &= \exp [t^T b] \cdot M_x(At) . \end{aligned} \quad (118)$$

The moment-generating function of the multivariate normal distribution (\rightarrow Proof “mvn-mgf”) is

$$M_x(t) = \exp \left[t^T \mu + \frac{1}{2} t^T \Sigma t \right] \quad (119)$$

and therefore the moment-generating function of the random vector y becomes

$$\begin{aligned} M_y(t) &= \exp [t^T b] \cdot M_x(At) \\ &= \exp [t^T b] \cdot \exp \left[t^T A\mu + \frac{1}{2} t^T A\Sigma A^T t \right] \\ &= \exp \left[t^T (A\mu + b) + \frac{1}{2} t^T A\Sigma A^T t \right] . \end{aligned} \quad (120)$$

Because moment-generating function and probability density function of a random variable are equivalent, this demonstrates that y is following a multivariate normal distribution with mean $A\mu + b$ and covariance $A\Sigma A^T$.

Sources:

- Taboga, Marco (2010): “Linear combinations of normal random variables”; in: *Lectures on probability and statistics*; URL: <https://www.statlect.com/probability-distributions/normal-distribution-linear-combinations>.

Metadata: ID: P1 | shortcut: mvn-ltt | author: JoramSoch | date: 2019-08-27, 12:14.

4.1.4 Marginal distributions

Theorem: Let x follow a multivariate normal distribution (\rightarrow Definition II/4.1.1):

$$x \sim \mathcal{N}(\mu, \Sigma) . \quad (121)$$

Then, the marginal distribution (\rightarrow Definition “md”) of any subset vector x_s is also a multivariate normal distribution

$$x_s \sim \mathcal{N}(\mu_s, \Sigma_s) \quad (122)$$

where μ_s drops the irrelevant variables (the ones not in the subset, i.e. marginalized out) from the mean vector μ and Σ_s drops the corresponding rows and columns from the covariance matrix Σ .

Proof: Define an $m \times n$ subset matrix S such that $s_{ij} = 1$, if the j -th element in μ_s corresponds to the i -th element in x , and $s_{ij} = 0$ otherwise. Then,

$$x_s = Sx \quad (123)$$

and we can apply the linear transformation theorem (\rightarrow Proof II/4.1.3) to give

$$x_s \sim \mathcal{N}(S\mu, S\Sigma S^T) . \quad (124)$$

Finally, we see that $S\mu = \mu_s$ and $S\Sigma S^T = \Sigma_s$.

Sources:

- original work

Metadata: ID: P35 | shortcut: mvn-marg | author: JoramSoch | date: 2020-01-29, 15:12.

4.2 Normal-gamma distribution

4.2.1 Definition

****Definition**:** Let X be an $n \times 1$ random vector (\rightarrow Definition “rvec”) and let Y be a positive random variable (\rightarrow Definition “rvar”). Then, X and Y are said to follow a normal-gamma distribution

$$X, Y \sim \text{NG}(\mu, \Lambda, a, b) , \quad (125)$$

if and only if their joint probability (\rightarrow Definition “jp”) density function (\rightarrow Definition “pdf”) is given by

$$f_{X,Y}(x, y) = \mathcal{N}(x; \mu, (y\Lambda)^{-1}) \cdot \text{Gam}(y; a, b) \quad (126)$$

where $\mathcal{N}(x; \mu, \Sigma)$ is the probability density function of the multivariate normal distribution (\rightarrow Proof II/4.1.2) with mean μ and covariance Σ and $\text{Gam}(x; a, b)$ is the probability density function of the gamma distribution (\rightarrow Proof II/3.3.2) with shape a and rate b .

The $n \times n$ matrix Λ is referred to as the precision matrix of the normal-gamma distribution.

Sources:

- Koch KR (2007): “Normal-Gamma Distribution”; in: *Introduction to Bayesian Statistics*, ch. 2.5.3, pp. 55-56, eq. 2.212; URL: <https://www.springer.com/gp/book/9783540727231>; DOI: 10.1007/978-3-540-72726-2.

Metadata: ID: D5 | shortcut: ng | author: JoramSoch | date: 2020-01-27, 14:28.

4.2.2 Probability density function

Theorem: Let x and y follow a normal-gamma distribution (\rightarrow Definition II/4.2.1):

$$x, y \sim \text{NG}(\mu, \Lambda, a, b) . \quad (127)$$

Then, the joint probability (\rightarrow Definition “jp”) density function (\rightarrow Definition “pdf”) of x and y is

$$p(x, y) = \sqrt{\frac{|\Lambda|}{(2\pi)^n}} \frac{b^a}{\Gamma(a)} \cdot y^{a+\frac{n}{2}-1} \exp \left[-\frac{y}{2} \left((x - \mu)^T \Lambda (x - \mu) + 2b \right) \right] . \quad (128)$$

Proof: The probability density of the normal-gamma distribution is defined as (\rightarrow Definition II/4.2.1) as the product of a multivariate normal distribution (\rightarrow Definition II/4.1.1) over x conditional on y and a univariate gamma distribution (\rightarrow Definition II/3.3.1) over y :

$$p(x, y) = \mathcal{N}(x; \mu, (y\Lambda)^{-1}) \cdot \text{Gam}(y; a, b) \quad (129)$$

With the probability density function of the multivariate normal distribution (\rightarrow Proof II/4.1.2) and the probability density function of the gamma distribution (\rightarrow Proof II/3.3.2), this becomes:

$$p(x, y) = \sqrt{\frac{|y\Lambda|}{(2\pi)^n}} \exp \left[-\frac{1}{2} (x - \mu)^T (y\Lambda) (x - \mu) \right] \cdot \frac{b^a}{\Gamma(a)} y^{a-1} \exp[-by] . \quad (130)$$

Using the relation $|yA| = y^n |A|$ for an $n \times n$ matrix A and rearranging the terms, we have:

$$p(x, y) = \sqrt{\frac{|\Lambda|}{(2\pi)^n}} \frac{b^a}{\Gamma(a)} \cdot y^{a+\frac{n}{2}-1} \exp \left[-\frac{y}{2} \left((x - \mu)^T \Lambda (x - \mu) + 2b \right) \right] . \quad (131)$$

Sources:

- Koch KR (2007): “Normal-Gamma Distribution”; in: *Introduction to Bayesian Statistics*, ch. 2.5.3, pp. 55-56, eq. 2.212; URL: <https://www.springer.com/gp/book/9783540727231>; DOI: 10.1007/978-3-540-72726-2.

4.2.3 Kullback-Leibler divergence

Theorem: Let $x \in \mathbb{R}^k$ be a random vector (\rightarrow Definition “rvec”) and $y > 0$ be a random variable (\rightarrow Definition “rvar”). Assume two normal-gamma distributions (\rightarrow Definition II/4.2.1) P and Q specifying the joint distribution of x and y as

$$\begin{aligned} P : (x, y) &\sim \text{NG}(\mu_1, \Lambda_1^{-1}, a_1, b_1) \\ Q : (x, y) &\sim \text{NG}(\mu_2, \Lambda_2^{-1}, a_2, b_2) . \end{aligned} \quad (132)$$

Then, the Kullback-Leibler divergence (\rightarrow Definition “kl”) of P from Q is given by

$$\begin{aligned} \text{KL}[P \parallel Q] &= \frac{1}{2} \frac{a_1}{b_1} [(\mu_2 - \mu_1)^T \Lambda_2 (\mu_2 - \mu_1)] + \frac{1}{2} \text{tr}(\Lambda_2 \Lambda_1^{-1}) - \frac{1}{2} \ln \frac{|\Lambda_2|}{|\Lambda_1|} - \frac{k}{2} \\ &+ a_2 \ln \frac{b_1}{b_2} - \ln \frac{\Gamma(a_1)}{\Gamma(a_2)} + (a_1 - a_2) \psi(a_1) - (b_1 - b_2) \frac{a_1}{b_1} . \end{aligned} \quad (133)$$

Proof: The probability density function of the normal-gamma distribution (\rightarrow Proof II/4.2.2) is

$$p(x, y) = p(x|y) \cdot p(y) = \mathcal{N}(x; \mu, (y\Lambda)^{-1}) \cdot \text{Gam}(y; a, b) \quad (134)$$

where $\mathcal{N}(x; \mu, \Sigma)$ is a multivariate normal density with mean μ and covariance Σ (hence, precision Λ) and $\text{Gam}(y; a, b)$ is a univariate gamma density with shape a and rate b . The Kullback-Leibler divergence of the multivariate normal distribution (\rightarrow Proof “mvn-kl”) is

$$\text{KL}[P \parallel Q] = \frac{1}{2} \left[(\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) + \text{tr}(\Sigma_2^{-1} \Sigma_1) - \ln \frac{|\Sigma_1|}{|\Sigma_2|} - k \right] \quad (135)$$

and the Kullback-Leibler divergence of the univariate gamma distribution (\rightarrow Proof “gam-kl”) is

$$\text{KL}[P \parallel Q] = a_2 \ln \frac{b_1}{b_2} - \ln \frac{\Gamma(a_1)}{\Gamma(a_2)} + (a_1 - a_2) \psi(a_1) - (b_1 - b_2) \frac{a_1}{b_1} \quad (136)$$

where $\Gamma(x)$ is the gamma function and $\psi(x)$ is the digamma function.

The KL divergence for a continuous random variable (\rightarrow Definition “kl”) is given by

$$\text{KL}[P \parallel Q] = \int_{\mathcal{Z}} p(z) \ln \frac{p(z)}{q(z)} dz \quad (137)$$

which, applied to the normal-gamma distribution (\rightarrow Definition II/4.2.1) over x and y , yields

$$\text{KL}[P \parallel Q] = \int_0^\infty \int_{\mathbb{R}^k} p(x, y) \ln \frac{p(x, y)}{q(x, y)} dx dy . \quad (138)$$

Using the law of conditional probability (\rightarrow Proof “lcp”), this can be evaluated as follows:

$$\begin{aligned} \text{KL}[P \parallel Q] &= \int_0^\infty \int_{\mathbb{R}^k} p(x|y) p(y) \ln \frac{p(x|y) p(y)}{q(x|y) q(y)} dx dy \\ &= \int_0^\infty \int_{\mathbb{R}^k} p(x|y) p(y) \ln \frac{p(x|y)}{q(x|y)} dx dy + \int_0^\infty \int_{\mathbb{R}^k} p(x|y) p(y) \ln \frac{p(y)}{q(y)} dx dy \\ &= \int_0^\infty p(y) \int_{\mathbb{R}^k} p(x|y) \ln \frac{p(x|y)}{q(x|y)} dx dy + \int_0^\infty p(y) \ln \frac{p(y)}{q(y)} \int_{\mathbb{R}^k} p(x|y) dx dy \\ &= \langle \text{KL}[p(x|y) \parallel q(x|y)] \rangle_{p(y)} + \text{KL}[p(y) \parallel q(y)] . \end{aligned} \quad (139)$$

In other words, the KL divergence between two normal-gamma distributions over x and y is equal to the sum of a multivariate normal KL divergence regarding x conditional on y , expected over y , and a univariate gamma KL divergence regarding y .

From equations (134) and (135), the first term becomes

$$\begin{aligned} &\langle \text{KL}[p(x|y) \parallel q(x|y)] \rangle_{p(y)} \\ &= \left\langle \frac{1}{2} \left[(\mu_2 - \mu_1)^T (y\Lambda_2)(\mu_2 - \mu_1) + \text{tr}((y\Lambda_2)(y\Lambda_1)^{-1}) - \ln \frac{|(y\Lambda_1)^{-1}|}{|(y\Lambda_2)^{-1}|} - k \right] \right\rangle_{p(y)} \quad (140) \\ &= \left\langle \frac{y}{2} (\mu_2 - \mu_1)^T \Lambda_2 (\mu_2 - \mu_1) + \frac{1}{2} \text{tr}(\Lambda_2 \Lambda_1^{-1}) - \frac{1}{2} \ln \frac{|\Lambda_2|}{|\Lambda_1|} - \frac{k}{2} \right\rangle_{p(y)} \end{aligned}$$

and using the relation (\rightarrow Proof “gam-mean”) $y \sim \text{Gam}(a, b) \Rightarrow \langle y \rangle = a/b$, we have

$$\langle \text{KL}[p(x|y) \parallel q(x|y)] \rangle_{p(y)} = \frac{1}{2} \frac{a_1}{b_1} (\mu_2 - \mu_1)^T \Lambda_2 (\mu_2 - \mu_1) + \frac{1}{2} \text{tr}(\Lambda_2 \Lambda_1^{-1}) - \frac{1}{2} \ln \frac{|\Lambda_2|}{|\Lambda_1|} - \frac{k}{2} . \quad (141)$$

By plugging (141) and (136) into (139), one arrives at the KL divergence given by (133).

Sources:

- Soch & Allefeld (2016): “Kullback-Leibler Divergence for the Normal-Gamma Distribution”; in: *arXiv math.ST*, 1611.01437; URL: <https://arxiv.org/abs/1611.01437>.

Metadata: ID: P6 | shortcut: ng-kl | author: JoramSoch | date: 2019-12-06, 09:35.

4.2.4 Marginal distributions

Theorem: Let x and y follow a normal-gamma distribution (\rightarrow Definition II/4.2.1):

$$x, y \sim \text{NG}(\mu, \Lambda, a, b) . \quad (142)$$

Then, the marginal distribution (\rightarrow Definition “md”) of y is a gamma distribution (\rightarrow Definition II/3.3.1)

$$y \sim \text{Gam}(a, b) \quad (143)$$

and the marginal distribution (\rightarrow Definition “md”) of x is a multivariate t-distribution (\rightarrow Definition “mvt”)

$$x \sim \text{t} \left(\mu, \left(\frac{a}{b} \Lambda \right)^{-1}, 2a \right) . \quad (144)$$

Proof: The probability density function of the normal-gamma distribution (\rightarrow Proof II/4.2.2) is given by

$$\begin{aligned} p(x, y) &= p(x|y) \cdot p(y) \\ p(x|y) &= \mathcal{N}(x; \mu, (y\Lambda)^{-1}) \\ p(y) &= \text{Gam}(y; a, b) . \end{aligned} \quad (145)$$

Using the law of marginal probability (\rightarrow Proof “lmp”), the marginal distribution of y can be derived as

$$\begin{aligned} p(y) &= \int p(x, y) \, dx \\ &= \int \mathcal{N}(x; \mu, (y\Lambda)^{-1}) \text{Gam}(y; a, b) \, dx \\ &= \text{Gam}(y; a, b) \int \mathcal{N}(x; \mu, (y\Lambda)^{-1}) \, dx \\ &= \text{Gam}(y; a, b) \end{aligned} \quad (146)$$

which is the probability density function of the gamma distribution (\rightarrow Proof II/3.3.2) with shape parameter a and rate parameter b .

Using the law of marginal probability (\rightarrow Proof “lmp”), the marginal distribution of x can be derived as

$$\begin{aligned}
p(x) &= \int p(x, y) \, dy \\
&= \int \mathcal{N}(x; \mu, (y\Lambda)^{-1}) \text{Gam}(y; a, b) \, dy \\
&= \int \sqrt{\frac{|y\Lambda|}{\sqrt{(2\pi)^n}}} \exp\left[-\frac{1}{2}(x - \mu)^T(y\Lambda)(x - \mu)\right] \cdot \frac{b^a}{\Gamma(a)} y^{a-1} \exp[-by] \, dy \\
&= \int \sqrt{\frac{y^n|\Lambda|}{\sqrt{(2\pi)^n}}} \exp\left[-\frac{1}{2}(x - \mu)^T(y\Lambda)(x - \mu)\right] \cdot \frac{b^a}{\Gamma(a)} y^{a-1} \exp[-by] \, dy \\
&= \int \sqrt{\frac{|\Lambda|}{\sqrt{(2\pi)^n}}} \cdot \frac{b^a}{\Gamma(a)} \cdot y^{a+\frac{n}{2}-1} \cdot \exp\left[-\left(b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right)y\right] \, dy \\
&= \int \sqrt{\frac{|\Lambda|}{\sqrt{(2\pi)^n}}} \cdot \frac{b^a}{\Gamma(a)} \cdot \frac{\Gamma\left(a + \frac{n}{2}\right)}{\left(b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right)^{a+\frac{n}{2}}} \cdot \text{Gam}\left(y; a + \frac{n}{2}, b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right) \, dy \\
&= \sqrt{\frac{|\Lambda|}{\sqrt{(2\pi)^n}}} \cdot \frac{b^a}{\Gamma(a)} \cdot \frac{\Gamma\left(a + \frac{n}{2}\right)}{\left(b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right)^{a+\frac{n}{2}}} \int \text{Gam}\left(y; a + \frac{n}{2}, b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right) \, dy \\
&= \sqrt{\frac{|\Lambda|}{\sqrt{(2\pi)^n}}} \cdot \frac{b^a}{\Gamma(a)} \cdot \frac{\Gamma\left(a + \frac{n}{2}\right)}{\left(b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right)^{a+\frac{n}{2}}} \\
&= \frac{\sqrt{|\Lambda|}}{(2\pi)^{\frac{n}{2}}} \cdot \frac{\Gamma\left(\frac{2a+n}{2}\right)}{\Gamma\left(\frac{2a}{2}\right)} \cdot b^a \cdot \left(b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right)^{-(a+\frac{n}{2})} \\
&= \frac{\sqrt{|\Lambda|}}{\pi^{\frac{n}{2}}} \cdot \frac{\Gamma\left(\frac{2a+n}{2}\right)}{\Gamma\left(\frac{2a}{2}\right)} \cdot \left(\frac{1}{b}\right)^{-a} \cdot \left(b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right)^{-a} \cdot 2^{-\frac{n}{2}} \cdot \left(b + \frac{1}{2}(x - \mu)^T\Lambda(x - \mu)\right)^{-\frac{n}{2}} \\
&= \frac{\sqrt{|\Lambda|}}{\pi^{\frac{n}{2}}} \cdot \frac{\Gamma\left(\frac{2a+n}{2}\right)}{\Gamma\left(\frac{2a}{2}\right)} \cdot \left(1 + \frac{1}{2b}(x - \mu)^T\Lambda(x - \mu)\right)^{-a} \cdot (2b + (x - \mu)^T\Lambda(x - \mu))^{-\frac{n}{2}} \\
&= \frac{\sqrt{|\Lambda|}}{\pi^{\frac{n}{2}}} \cdot \frac{\Gamma\left(\frac{2a+n}{2}\right)}{\Gamma\left(\frac{2a}{2}\right)} \cdot \left(\frac{1}{2a}\right)^{-a} \cdot \left(2a + (x - \mu)^T\left(\frac{a}{b}\Lambda\right)(x - \mu)\right)^{-a} \cdot \left(\frac{b}{a}\right)^{-\frac{n}{2}} \cdot \left(2a + (x - \mu)^T\left(\frac{a}{b}\Lambda\right)(x - \mu)\right)^{-\frac{n}{2}} \\
&= \frac{\sqrt{\left(\frac{a}{b}\right)^n |\Lambda|}}{(2a)^{-a} \pi^{\frac{n}{2}}} \cdot \frac{\Gamma\left(\frac{2a+n}{2}\right)}{\Gamma\left(\frac{2a}{2}\right)} \cdot \left(2a + (x - \mu)^T\left(\frac{a}{b}\Lambda\right)(x - \mu)\right)^{-a} \cdot \left(2a + (x - \mu)^T\left(\frac{a}{b}\Lambda\right)(x - \mu)\right)^{-\frac{n}{2}} \\
&= \frac{\sqrt{\left(\frac{a}{b}\right)^n |\Lambda|}}{(2a)^{-a} \pi^{\frac{n}{2}}} \cdot \frac{\Gamma\left(\frac{2a+n}{2}\right)}{\Gamma\left(\frac{2a}{2}\right)} \cdot (2a)^{-a} \cdot \left(1 + \frac{1}{2a}(x - \mu)^T\left(\frac{a}{b}\Lambda\right)(x - \mu)\right)^{-a} \cdot (2a)^{-\frac{n}{2}} \cdot \left(1 + \frac{1}{2a}(x - \mu)^T\left(\frac{a}{b}\Lambda\right)(x - \mu)\right)^{-\frac{n}{2}} \\
&= \frac{\sqrt{\left(\frac{a}{b}\right)^n |\Lambda|}}{(2a)^{\frac{n}{2}} \pi^{\frac{n}{2}}} \cdot \frac{\Gamma\left(\frac{2a+n}{2}\right)}{\Gamma\left(\frac{2a}{2}\right)} \cdot \left(1 + \frac{1}{2a}(x - \mu)^T\left(\frac{a}{b}\Lambda\right)(x - \mu)\right)^{-\frac{2a+n}{2}} \\
&= \sqrt{\frac{\left|\frac{a}{b}\Lambda\right|}{(2a\pi)^n}} \cdot \frac{\Gamma\left(\frac{2a+n}{2}\right)}{\Gamma\left(\frac{2a}{2}\right)} \cdot \left(1 + \frac{1}{2a}(x - \mu)^T\left(\frac{a}{b}\Lambda\right)(x - \mu)\right)^{-\frac{2a+n}{2}}
\end{aligned} \tag{147}$$

which is the probability density function of a multivariate t-distribution (\rightarrow Proof “mvt-pdf”) with mean vector μ , shape matrix $\left(\frac{a}{b}\Lambda\right)^{-1}$ and $2a$ degrees of freedom.

Sources:

- original work

Metadata: ID: P36 | shortcut: ng-marg | author: JoramSoch | date: 2020-01-29, 21:42.

5 Matrix-variate continuous distributions

5.1 Matrix-normal distribution

5.1.1 Definition

****Definition**:** Let X be an $n \times p$ random matrix (\rightarrow Definition “rmat”). Then, X is said to be matrix-normally distributed with mean M , covariance across rows U and covariance across columns V

$$X \sim \mathcal{MN}(M, U, V) , \quad (148)$$

if and only if its probability density function (\rightarrow Definition “pdf”) is given by

$$\mathcal{MN}(X; M, U, V) = \frac{1}{\sqrt{(2\pi)^{np} |V|^n |U|^p}} \cdot \exp \left[-\frac{1}{2} \text{tr} (V^{-1} (X - M)^T U^{-1} (X - M)) \right] \quad (149)$$

where μ is an $n \times p$ real matrix, U is an $n \times n$ positive definite matrix and V is a $p \times p$ positive definite matrix.

Sources:

- Wikipedia (2020): “Matrix normal distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-27; URL: https://en.wikipedia.org/wiki/Matrix_normal_distribution#Definition.

Metadata: ID: D6 | shortcut: matn | author: JoramSoch | date: 2020-01-27, 14:37.

5.1.2 Equivalence to multivariate normal distribution

Theorem: The matrix X is matrix-normally distributed (\rightarrow Definition II/5.1.1)

$$X \sim \mathcal{MN}(M, U, V) , \quad (150)$$

if and only if $\text{vec}(X)$ is multivariate normally distributed (\rightarrow Definition II/4.1.1)

$$\text{vec}(X) \sim \mathcal{MN}(\text{vec}(M), V \otimes U) \quad (151)$$

where $\text{vec}(X)$ is the vectorization operator and \otimes is the Kronecker product.

Proof: The probability density function of the matrix-normal distribution (\rightarrow Proof “matn-pdf”) with $n \times p$ mean M , $n \times n$ covariance across rows U and $p \times p$ covariance across columns V is

$$\mathcal{MN}(X; M, U, V) = \frac{1}{\sqrt{(2\pi)^{np} |V|^n |U|^p}} \cdot \exp \left[-\frac{1}{2} \text{tr} (V^{-1} (X - M)^T U^{-1} (X - M)) \right] . \quad (152)$$

Using the trace property $\text{tr}(ABC) = \text{tr}(BCA)$, we have:

$$\mathcal{MN}(X; M, U, V) = \frac{1}{\sqrt{(2\pi)^{np}|V|^n|U|^p}} \cdot \exp \left[-\frac{1}{2} \text{tr} \left((X - M)^T U^{-1} (X - M) V^{-1} \right) \right] . \quad (153)$$

Using the trace-vectorization relation $\text{tr}(A^T B) = \text{vec}(A)^T \text{vec}(B)$, we have:

$$\mathcal{MN}(X; M, U, V) = \frac{1}{\sqrt{(2\pi)^{np}|V|^n|U|^p}} \cdot \exp \left[-\frac{1}{2} \text{vec}(X - M)^T \text{vec} \left(U^{-1} (X - M) V^{-1} \right) \right] . \quad (154)$$

Using the vectorization-Kronecker relation $\text{vec}(ABC) = (C^T \otimes A) \text{vec}(B)$, we have:

$$\mathcal{MN}(X; M, U, V) = \frac{1}{\sqrt{(2\pi)^{np}|V|^n|U|^p}} \cdot \exp \left[-\frac{1}{2} \text{vec}(X - M)^T (V^{-1} \otimes U^{-1}) \text{vec}(X - M) \right] . \quad (155)$$

Using the Kronecker product property $(A^{-1} \otimes B^{-1}) = (A \otimes B)^{-1}$, we have:

$$\mathcal{MN}(X; M, U, V) = \frac{1}{\sqrt{(2\pi)^{np}|V|^n|U|^p}} \cdot \exp \left[-\frac{1}{2} \text{vec}(X - M)^T (V \otimes U)^{-1} \text{vec}(X - M) \right] . \quad (156)$$

Using the vectorization property $\text{vec}(A + B) = \text{vec}(A) + \text{vec}(B)$, we have:

$$\mathcal{MN}(X; M, U, V) = \frac{1}{\sqrt{(2\pi)^{np}|V|^n|U|^p}} \cdot \exp \left[-\frac{1}{2} [\text{vec}(X) - \text{vec}(M)]^T (V \otimes U)^{-1} [\text{vec}(X) - \text{vec}(M)] \right] . \quad (157)$$

Using the Kronecker-determinant relation $|A \otimes B| = |A|^m |B|^n$, we have:

$$\mathcal{MN}(X; M, U, V) = \frac{1}{\sqrt{(2\pi)^{np}|V \otimes U|}} \cdot \exp \left[-\frac{1}{2} [\text{vec}(X) - \text{vec}(M)]^T (V \otimes U)^{-1} [\text{vec}(X) - \text{vec}(M)] \right] . \quad (158)$$

This is the probability density function of the multivariate normal distribution (\rightarrow Proof II/4.1.2) with the $np \times 1$ mean vector $\text{vec}(M)$ and the $np \times np$ covariance matrix $V \otimes U$:

$$\mathcal{MN}(X; M, U, V) = \mathcal{N}(\text{vec}(X); \text{vec}(M), V \otimes U) . \quad (159)$$

By showing that the probability density functions (\rightarrow Definition “pdf”) are identical, it is proven that the associated probability distributions (\rightarrow Definition “pd”) are equivalent.

Sources:

- Wikipedia (2020): “Matrix normal distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-20; URL: https://en.wikipedia.org/wiki/Matrix_normal_distribution#Proof.

Metadata: ID: P26 | shortcut: matn-mvn | author: JoramSoch | date: 2020-01-20, 21:09.

Chapter III

Statistical Models

1 Normal data

1.1 Multiple linear regression

1.1.1 Ordinary least squares (1)

Theorem: Given a linear regression model (\rightarrow Definition “mlr”) with independent observations

$$y = X\beta + \varepsilon, \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad (1)$$

the parameters minimizing the residual sum of squares (\rightarrow Definition “rss”) are given by

$$\hat{\beta} = (X^T X)^{-1} X^T y. \quad (2)$$

Proof: Let $\hat{\beta}$ be the ordinary least squares (OLS) solution and let $\hat{\varepsilon} = y - X\hat{\beta}$ be the resulting vector of residuals. Then, this vector must be orthogonal to the design matrix,

$$X^T \hat{\varepsilon} = 0, \quad (3)$$

because if it wasn't, there would be another solution $\tilde{\beta}$ giving another vector $\tilde{\varepsilon}$ with a smaller residual sum of squares. From (3), the OLS formula can be directly derived:

$$\begin{aligned} X^T \hat{\varepsilon} &= 0 \\ X^T (y - X\hat{\beta}) &= 0 \\ X^T y - X^T X \hat{\beta} &= 0 \\ X^T X \hat{\beta} &= X^T y \\ \hat{\beta} &= (X^T X)^{-1} X^T y. \end{aligned} \quad (4)$$

Sources:

- Stephan, Klaas Enno (2010): “The General Linear Model (GLM)”;
in: *Methods and models for fMRI data analysis in neuroeconomics*; URL: <http://www.socialbehavior.uzh.ch/teaching/methodspring10.html>.

Metadata: ID: P2 | shortcut: mlr-ols | author: JoramSoch | date: 2019-09-27, 07:18.

1.1.2 Ordinary least squares (2)

Theorem: Given a linear regression model (\rightarrow Definition “mlr”) with independent observations

$$y = X\beta + \varepsilon, \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad (5)$$

the parameters minimizing the residual sum of squares (\rightarrow Definition “rss”) are given by

$$\hat{\beta} = (X^T X)^{-1} X^T y . \quad (6)$$

Proof: The residual sum of squares (\rightarrow Definition “rss”) is defined as

$$\text{RSS}(\beta) = \sum_{i=1}^n \varepsilon_i = \varepsilon^T \varepsilon = (y - X\beta)^T (y - X\beta) \quad (7)$$

which can be developed into

$$\begin{aligned} \text{RSS}(\beta) &= y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta \\ &= y^T y - 2\beta^T X^T y + \beta^T X^T X\beta . \end{aligned} \quad (8)$$

The derivative of $\text{RSS}(\beta)$ with respect to β is

$$\frac{d\text{RSS}(\beta)}{d\beta} = -2X^T y + 2X^T X\beta \quad (9)$$

and setting this derivative to zero, we obtain:

$$\begin{aligned} \frac{d\text{RSS}(\hat{\beta})}{d\beta} &= 0 \\ 0 &= -2X^T y + 2X^T X\hat{\beta} \\ X^T X\hat{\beta} &= X^T y \\ \hat{\beta} &= (X^T X)^{-1} X^T y . \end{aligned} \quad (10)$$

Since the quadratic form $y^T y$ in (8) is positive, $\hat{\beta}$ minimizes $\text{RSS}(\beta)$.

Sources:

- Wikipedia (2020): “Proofs involving ordinary least squares”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-02-03; URL: https://en.wikipedia.org/wiki/Proofs_involving_ordinary_least_squares#Least_squares_estimator_for_%CE%B2.

Metadata: ID: P40 | shortcut: mlr-ols2 | author: JoramSoch | date: 2020-02-03, 18:43.

1.2 Bayesian linear regression

1.2.1 Conjugate prior distribution

Theorem: Let

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 V) \quad (11)$$

be a linear regression model (\rightarrow Proof “mlr”) with measured $n \times 1$ data vector y , known $n \times p$ design matrix X , known $n \times n$ covariance structure V and unknown $p \times 1$ regression coefficients β and noise variance σ^2 .

Then, the conjugate prior (\rightarrow Definition “prior-conj”) for this model is a normal-gamma distribution (\rightarrow Definition II/4.2.1)

$$p(\beta, \tau) = \mathcal{N}(\beta; \mu_0, (\tau \Lambda_0)^{-1}) \cdot \text{Gam}(\tau; a_0, b_0) \quad (12)$$

where $\tau = 1/\sigma^2$ is the inverse noise variance or noise precision.

Proof: By definition, a conjugate prior (\rightarrow Definition “prior-conj”) is a prior distribution (\rightarrow Definition “prior”) that, when combined with the likelihood function (\rightarrow Definition “lf”), leads to a posterior distribution (\rightarrow Definition “post”) that belongs to the same family of probability distributions (\rightarrow Definition “pd”). This is fulfilled when the prior density and the likelihood function are proportional to the model parameters in the same way, i.e. the model parameters appear in the same functional form in both. Equation (11) implies the following likelihood function (\rightarrow Definition “lf”)

$$p(y|\beta, \sigma^2) = \mathcal{N}(y; X\beta, \sigma^2 V) = \sqrt{\frac{1}{(2\pi)^n |\sigma^2 V|}} \exp \left[-\frac{1}{2\sigma^2} (y - X\beta)^T V^{-1} (y - X\beta) \right] \quad (13)$$

which, for mathematical convenience, can also be parametrized as

$$p(y|\beta, \tau) = \mathcal{N}(y; X\beta, (\tau P)^{-1}) = \sqrt{\frac{|\tau P|}{(2\pi)^n}} \exp \left[-\frac{\tau}{2} (y - X\beta)^T P (y - X\beta) \right] \quad (14)$$

using the noise precision $\tau = 1/\sigma^2$ and the $n \times n$ precision matrix $P = V^{-1}$.

Separating constant and variable terms, we have:

$$p(y|\beta, \tau) = \sqrt{\frac{|P|}{(2\pi)^n}} \cdot \tau^{n/2} \cdot \exp \left[-\frac{\tau}{2} (y - X\beta)^T P (y - X\beta) \right]. \quad (15)$$

Expanding the product in the exponent, we have:

$$p(y|\beta, \tau) = \sqrt{\frac{|P|}{(2\pi)^n}} \cdot \tau^{n/2} \cdot \exp \left[-\frac{\tau}{2} (y^T P y - y^T P X \beta - \beta^T X^T P y + \beta^T X^T P X \beta) \right]. \quad (16)$$

Completing the square over β , finally gives

$$p(y|\beta, \tau) = \sqrt{\frac{|P|}{(2\pi)^n}} \cdot \tau^{n/2} \cdot \exp \left[-\frac{\tau}{2} \left((\beta - \tilde{X} y)^T X^T P X (\beta - \tilde{X} y) - y^T Q y + y^T P y \right) \right] \quad (17)$$

where $\tilde{X} = (X^T P X)^{-1} X^T P$ and $Q = \tilde{X}^T (X^T P X) \tilde{X}$.

In other words, the likelihood function (\rightarrow Definition “lf”) is proportional to a power of τ times an exponential of τ and an exponential of a squared form of β , weighted by τ :

$$p(y|\beta, \tau) \propto \tau^{n/2} \cdot \exp \left[-\frac{\tau}{2} (y^T P y - y^T Q y) \right] \cdot \exp \left[-\frac{\tau}{2} (\beta - \tilde{X} y)^T X^T P X (\beta - \tilde{X} y) \right] . \quad (18)$$

The same is true for a normal gamma distribution over β and τ

$$p(\beta, \tau) = \mathcal{N}(\beta; \mu_0, (\tau \Lambda_0)^{-1}) \cdot \text{Gam}(\tau; a_0, b_0) \quad (19)$$

the probability density function of which (\rightarrow Proof II/4.2.2)

$$p(\beta, \tau) = \sqrt{\frac{|\tau \Lambda_0|}{(2\pi)^p}} \exp \left[-\frac{\tau}{2} (\beta - \mu_0)^T \Lambda_0 (\beta - \mu_0) \right] \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \quad (20)$$

exhibits the same proportionality

$$p(\beta, \tau) \propto \tau^{a_0+p/2-1} \cdot \exp[-\tau b_0] \cdot \exp \left[-\frac{\tau}{2} (\beta - \mu_0)^T \Lambda_0 (\beta - \mu_0) \right] \quad (21)$$

and is therefore conjugate relative to the likelihood.

Sources:

- Bishop CM (2006): “Bayesian linear regression”; in: *Pattern Recognition for Machine Learning*, pp. 152-161, ex. 3.12, eq. 3.112; URL: <https://www.springer.com/gp/book/9780387310732>.

Metadata: ID: P9 | shortcut: blr-prior | author: JoramSoch | date: 2020-01-03, 15:26.

1.2.2 Posterior distribution

Theorem: Let

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 V) \quad (22)$$

be a linear regression model (\rightarrow Definition “mlr”) with measured $n \times 1$ data vector y , known $n \times p$ design matrix X , known $n \times n$ covariance structure V and unknown $p \times 1$ regression coefficients β and noise variance σ^2 . Moreover, assume a normal-gamma prior distribution (\rightarrow Proof III/1.2.1) over the model parameters β and $\tau = 1/\sigma^2$:

$$p(\beta, \tau) = \mathcal{N}(\beta; \mu_0, (\tau \Lambda_0)^{-1}) \cdot \text{Gam}(\tau; a_0, b_0) . \quad (23)$$

Then, the posterior distribution (\rightarrow Definition “post”) is also a normal-gamma distribution (\rightarrow Definition II/4.2.1)

$$p(\beta, \tau|y) = \mathcal{N}(\beta; \mu_n, (\tau \Lambda_n)^{-1}) \cdot \text{Gam}(\tau; a_n, b_n) \quad (24)$$

and the posterior hyperparameters (\rightarrow Definition “post-hyp”) are given by

$$\begin{aligned}
\mu_n &= \Lambda_n^{-1}(X^T P y + \Lambda_0 \mu_0) \\
\Lambda_n &= X^T P X + \Lambda_0 \\
a_n &= a_0 + \frac{n}{2} \\
b_n &= b_0 + \frac{1}{2}(y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) .
\end{aligned} \tag{25}$$

Proof: According to Bayes' theorem (\rightarrow Proof I/1.2.1), the posterior distribution (\rightarrow Definition “post”) is given by

$$p(\beta, \tau|y) = \frac{p(y|\beta, \tau) p(\beta, \tau)}{p(y)} . \tag{26}$$

Since $p(y)$ is just a normalization factor, the posterior is proportional (\rightarrow Proof “post-jl”) to the numerator:

$$p(\beta, \tau|y) \propto p(y|\beta, \tau) p(\beta, \tau) = p(y, \beta, \tau) . \tag{27}$$

Equation (22) implies the following likelihood function (\rightarrow Definition “lf”)

$$p(y|\beta, \sigma^2) = \mathcal{N}(y; X\beta, \sigma^2 V) = \sqrt{\frac{1}{(2\pi)^n |\sigma^2 V|}} \exp \left[-\frac{1}{2\sigma^2} (y - X\beta)^T V^{-1} (y - X\beta) \right] \tag{28}$$

which, for mathematical convenience, can also be parametrized as

$$p(y|\beta, \tau) = \mathcal{N}(y; X\beta, (\tau P)^{-1}) = \sqrt{\frac{|\tau P|}{(2\pi)^n}} \exp \left[-\frac{\tau}{2} (y - X\beta)^T P (y - X\beta) \right] \tag{29}$$

using the noise precision $\tau = 1/\sigma^2$ and the $n \times n$ precision matrix $P = V^{-1}$.

Combining the likelihood function (29) with the prior distribution (23), the sssssssssssssssssssssssssssssssss of the model is given by

$$\begin{aligned}
p(y, \beta, \tau) &= p(y|\beta, \tau) p(\beta, \tau) \\
&= \sqrt{\frac{|\tau P|}{(2\pi)^n}} \exp \left[-\frac{\tau}{2} (y - X\beta)^T P (y - X\beta) \right] \cdot \\
&\quad \sqrt{\frac{|\tau \Lambda_0|}{(2\pi)^p}} \exp \left[-\frac{\tau}{2} (\beta - \mu_0)^T \Lambda_0 (\beta - \mu_0) \right] \cdot \\
&\quad \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] .
\end{aligned} \tag{30}$$

Collecting identical variables gives:

$$p(y, \beta, \tau) = \sqrt{\frac{\tau^{n+p}}{(2\pi)^{n+p}} |P| |\Lambda_0|} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \exp \left[-\frac{\tau}{2} \left((y - X\beta)^T P (y - X\beta) + (\beta - \mu_0)^T \Lambda_0 (\beta - \mu_0) \right) \right] . \quad (31)$$

Expanding the products in the exponent gives:

$$p(y, \beta, \tau) = \sqrt{\frac{\tau^{n+p}}{(2\pi)^{n+p}} |P| |\Lambda_0|} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \exp \left[-\frac{\tau}{2} \left(y^T P y - y^T P X \beta - \beta^T X^T P y + \beta^T X^T P X \beta + \beta^T \Lambda_0 \beta - \beta^T \Lambda_0 \mu_0 - \mu_0^T \Lambda_0 \beta + \mu_0^T \Lambda_0 \mu_0 \right) \right] . \quad (32)$$

Completing the square over β , we finally have

$$p(y, \beta, \tau) = \sqrt{\frac{\tau^{n+p}}{(2\pi)^{n+p}} |P| |\Lambda_0|} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \exp \left[-\frac{\tau}{2} \left((\beta - \mu_n)^T \Lambda_n (\beta - \mu_n) + (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) \right) \right] \quad (33)$$

with the posterior hyperparameters (\rightarrow Definition “post-hyp”)

$$\begin{aligned} \mu_n &= \Lambda_n^{-1} (X^T P y + \Lambda_0 \mu_0) \\ \Lambda_n &= X^T P X + \Lambda_0 . \end{aligned} \quad (34)$$

Ergo, the joint likelihood is proportional to

$$p(y, \beta, \tau) \propto \tau^{p/2} \cdot \exp \left[-\frac{\tau}{2} (\beta - \mu_n)^T \Lambda_n (\beta - \mu_n) \right] \cdot \tau^{a_n-1} \cdot \exp [-b_n \tau] \quad (35)$$

with the posterior hyperparameters (\rightarrow Definition “post-hyp”)

$$\begin{aligned} a_n &= a_0 + \frac{n}{2} \\ b_n &= b_0 + \frac{1}{2} (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) . \end{aligned} \quad (36)$$

From the term in (35), we can isolate the posterior distribution over β given τ :

$$p(\beta | \tau, y) = \mathcal{N}(\beta; \mu_n, (\tau \Lambda_n)^{-1}) . \quad (37)$$

From the remaining term, we can isolate the posterior distribution over τ :

$$p(\tau | y) = \text{Gam}(\tau; a_n, b_n) . \quad (38)$$

Together, (37) and (38) constitute the joint posterior distribution (\rightarrow Definition “jp”) of β and τ .

Sources:

- Bishop CM (2006): “Bayesian linear regression”; in: *Pattern Recognition for Machine Learning*, pp. 152-161, ex. 3.12, eq. 3.113; URL: <https://www.springer.com/gp/book/9780387310732>.

Metadata: ID: P10 | shortcut: blr-post | author: JoramSoch | date: 2020-01-03, 17:53.

1.2.3 Log model evidence

Theorem: Let

$$m : y = X\beta + \varepsilon, \varepsilon \sim \mathcal{N}(0, \sigma^2 V) \quad (39)$$

be a linear regression model (\rightarrow Definition “mlr”) with measured $n \times 1$ data vector y , known $n \times p$ design matrix X , known $n \times n$ covariance structure V and unknown $p \times 1$ regression coefficients β and noise variance σ^2 . Moreover, assume a normal-gamma prior distribution (\rightarrow Proof III/1.2.1) over the model parameters β and $\tau = 1/\sigma^2$:

$$p(\beta, \tau) = \mathcal{N}(\beta; \mu_0, (\tau \Lambda_0)^{-1}) \cdot \text{Gam}(\tau; a_0, b_0) . \quad (40)$$

Then, the log model evidence (\rightarrow Definition “lme”) for this model is

$$\begin{aligned} \log p(y|m) = & \frac{1}{2} \log |P| - \frac{n}{2} \log(2\pi) + \frac{1}{2} \log |\Lambda_0| - \frac{1}{2} \log |\Lambda_n| + \\ & \log \Gamma(a_n) - \log \Gamma(a_0) + a_0 \log b_0 - a_n \log b_n \end{aligned} \quad (41)$$

where the posterior hyperparameters (\rightarrow Definition “post-hyp”) are given by

$$\begin{aligned} \mu_n &= \Lambda_n^{-1} (X^T P y + \Lambda_0 \mu_0) \\ \Lambda_n &= X^T P X + \Lambda_0 \\ a_n &= a_0 + \frac{n}{2} \\ b_n &= b_0 + \frac{1}{2} (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) . \end{aligned} \quad (42)$$

Proof: According to the law of marginal probability (\rightarrow Proof “lmp”), the model evidence (\rightarrow Definition “ml”) for this model is:

$$p(y|m) = \iint p(y|\beta, \tau) p(\beta, \tau) d\beta d\tau . \quad (43)$$

According to the law of conditional probability (\rightarrow Proof “lcp”), the integrand is equivalent to the joint likelihood (\rightarrow Definition “jl”):

$$p(y|m) = \iint p(y, \beta, \tau) d\beta d\tau . \quad (44)$$

Equation (39) implies the following likelihood function (\rightarrow Definition “lf”)

$$p(y|\beta, \sigma^2) = \mathcal{N}(y; X\beta, \sigma^2 V) = \sqrt{\frac{1}{(2\pi)^n |\sigma^2 V|}} \exp \left[-\frac{1}{2\sigma^2} (y - X\beta)^T V^{-1} (y - X\beta) \right] \quad (45)$$

which, for mathematical convenience, can also be parametrized as

$$p(y|\beta, \tau) = \mathcal{N}(y; X\beta, (\tau P)^{-1}) = \sqrt{\frac{|\tau P|}{(2\pi)^n}} \exp \left[-\frac{\tau}{2} (y - X\beta)^T P (y - X\beta) \right] \quad (46)$$

using the noise precision $\tau = 1/\sigma^2$ and the $n \times n$ precision matrix $P = V^{-1}$.

When deriving the posterior distribution (\rightarrow Proof III/1.2.2) $p(\beta, \tau|y)$, the joint likelihood $p(y, \beta, \tau)$ is obtained as

$$p(y, \beta, \tau) = \sqrt{\frac{\tau^n |P|}{(2\pi)^n}} \sqrt{\frac{\tau^p |\Lambda_0|}{(2\pi)^p}} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \exp \left[-\frac{\tau}{2} ((\beta - \mu_n)^T \Lambda_n (\beta - \mu_n) + (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n)) \right]. \quad (47)$$

Using the probability density function of the multivariate normal distribution (\rightarrow Proof II/4.1.2), we can rewrite this as

$$p(y, \beta, \tau) = \sqrt{\frac{\tau^n |P|}{(2\pi)^n}} \sqrt{\frac{\tau^p |\Lambda_0|}{(2\pi)^p}} \sqrt{\frac{(2\pi)^p}{\tau^p |\Lambda_n|}} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \mathcal{N}(\beta; \mu_n, (\tau \Lambda_n)^{-1}) \exp \left[-\frac{\tau}{2} (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) \right]. \quad (48)$$

Now, β can be integrated out easily:

$$\int p(y, \beta, \tau) d\beta = \sqrt{\frac{\tau^n |P|}{(2\pi)^n}} \sqrt{\frac{|\Lambda_0|}{|\Lambda_n|}} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{a_0-1} \exp[-b_0 \tau] \cdot \exp \left[-\frac{\tau}{2} (y^T P y + \mu_0^T \Lambda_0 \mu_0 - \mu_n^T \Lambda_n \mu_n) \right]. \quad (49)$$

Using the probability density function of the gamma distribution (\rightarrow Proof II/3.3.2), we can rewrite this as

$$\int p(y, \beta, \tau) d\beta = \sqrt{\frac{|P|}{(2\pi)^n}} \sqrt{\frac{|\Lambda_0|}{|\Lambda_n|}} \frac{b_0^{a_0}}{\Gamma(a_0)} \frac{\Gamma(a_n)}{b_n^{a_n}} \text{Gam}(\tau; a_n, b_n). \quad (50)$$

Finally, τ can also be integrated out:

$$\iint p(y, \beta, \tau) d\beta d\tau = \sqrt{\frac{|P|}{(2\pi)^n}} \sqrt{\frac{|\Lambda_0|}{|\Lambda_n|}} \frac{\Gamma(a_n)}{\Gamma(a_0)} \frac{b_0^{a_0}}{b_n^{a_n}} = p(y|m). \quad (51)$$

Thus, the log model evidence (\rightarrow Definition “lme”) of this model is given by

$$\begin{aligned} \log p(y|m) = & \frac{1}{2} \log |P| - \frac{n}{2} \log(2\pi) + \frac{1}{2} \log |\Lambda_0| - \frac{1}{2} \log |\Lambda_n| + \\ & \log \Gamma(a_n) - \log \Gamma(a_0) + a_0 \log b_0 - a_n \log b_n . \end{aligned} \quad (52)$$

Sources:

- Bishop CM (2006): “Bayesian linear regression”; in: *Pattern Recognition for Machine Learning*, pp. 152-161, ex. 3.23, eq. 3.118; URL: <https://www.springer.com/gp/book/9780387310732>.

Metadata: ID: P11 | shortcut: blr-lme | author: JoramSoch | date: 2020-01-03, 22:05.

1.3 General linear model

1.3.1 Maximum likelihood estimation

Theorem: Given a general linear model (\rightarrow Definition “glm”) with matrix-normally distributed (\rightarrow Definition II/5.1.1) errors

$$Y = XB + E, \quad E \sim \mathcal{MN}(0, V, \Sigma) , \quad (53)$$

maximum likelihood estimates (\rightarrow Definition “mle”) for the unknown parameters B and Σ are given by

$$\begin{aligned} \hat{B} &= (X^T V^{-1} X)^{-1} X^T V^{-1} Y \\ \hat{\Sigma} &= \frac{1}{n} (Y - X \hat{B})^T V^{-1} (Y - X \hat{B}) . \end{aligned} \quad (54)$$

Proof: In (53), Y is an $n \times v$ matrix of measurements (n observations, v dependent variables), X is an $n \times p$ design matrix (n observations, p independent variables) and V is an $n \times n$ covariance matrix across observations. This multivariate GLM implies the following likelihood function (\rightarrow Definition “lf”)

$$\begin{aligned} p(Y|B, \Sigma) &= \mathcal{MN}(Y; XB, V, \Sigma) \\ &= \sqrt{\frac{1}{(2\pi)^{nv} |\Sigma|^n |V|^v}} \cdot \exp \left[-\frac{1}{2} \text{tr} (\Sigma^{-1} (Y - XB)^T V^{-1} (Y - XB)) \right] \end{aligned} \quad (55)$$

and the log-likelihood function (\rightarrow Definition “llf”)

$$\begin{aligned} \text{LL}(B, \Sigma) &= \log p(Y|B, \Sigma) \\ &= -\frac{nv}{2} \log(2\pi) - \frac{n}{2} \log |\Sigma| - \frac{v}{2} \log |V| \\ &\quad - \frac{1}{2} \text{tr} [\Sigma^{-1} (Y - XB)^T V^{-1} (Y - XB)] . \end{aligned} \quad (56)$$

Substituting V^{-1} by the precision matrix P to ease notation, we have:

$$\begin{aligned} \text{LL}(B, \Sigma) = & -\frac{nv}{2} \log(2\pi) - \frac{n}{2} \log(|\Sigma|) - \frac{v}{2} \log(|V|) \\ & - \frac{1}{2} \text{tr} [\Sigma^{-1} (Y^T P Y - Y^T P X B - B^T X^T P Y + B^T X^T P X B)] . \end{aligned} \quad (57)$$

The derivative of the log-likelihood function (57) with respect to B is

$$\begin{aligned} \frac{d\text{LL}(B, \Sigma)}{dB} = & \frac{d}{dB} \left(-\frac{1}{2} \text{tr} [\Sigma^{-1} (Y^T P Y - Y^T P X B - B^T X^T P Y + B^T X^T P X B)] \right) \\ = & \frac{d}{dB} \left(-\frac{1}{2} \text{tr} [-2\Sigma^{-1} Y^T P X B] \right) + \frac{d}{dB} \left(-\frac{1}{2} \text{tr} [\Sigma^{-1} B^T X^T P X B] \right) \\ = & -\frac{1}{2} (-2X^T P Y \Sigma^{-1}) - \frac{1}{2} (X^T P X B \Sigma^{-1} + (X^T P X)^T B (\Sigma^{-1})^T) \\ = & X^T P Y \Sigma^{-1} - X^T P X B \Sigma^{-1} \end{aligned} \quad (58)$$

and setting this derivative to zero gives the MLE for B :

$$\begin{aligned} \frac{d\text{LL}(\hat{B}, \Sigma)}{dB} &= 0 \\ 0 &= X^T P Y \Sigma^{-1} - X^T P X \hat{B} \Sigma^{-1} \\ 0 &= X^T P Y - X^T P X \hat{B} \\ X^T P X \hat{B} &= X^T P Y \\ \hat{B} &= (X^T P X)^{-1} X^T P Y \end{aligned} \quad (59)$$

The derivative of the log-likelihood function (56) at \hat{B} with respect to Σ is

$$\begin{aligned} \frac{d\text{LL}(\hat{B}, \Sigma)}{d\Sigma} = & \frac{d}{d\Sigma} \left(-\frac{n}{2} \log |\Sigma| - \frac{1}{2} \text{tr} [\Sigma^{-1} (Y - X \hat{B})^T V^{-1} (Y - X \hat{B})] \right) \\ = & -\frac{n}{2} (\Sigma^{-1})^T + \frac{1}{2} \left(\Sigma^{-1} (Y - X \hat{B})^T V^{-1} (Y - X \hat{B}) \Sigma^{-1} \right)^T \\ = & -\frac{n}{2} \Sigma^{-1} + \frac{1}{2} \Sigma^{-1} (Y - X \hat{B})^T V^{-1} (Y - X \hat{B}) \Sigma^{-1} \end{aligned} \quad (60)$$

and setting this derivative to zero gives the MLE for Σ :

$$\begin{aligned}
\frac{dLL(\hat{B}, \hat{\Sigma})}{d\Sigma} &= 0 \\
0 &= -\frac{n}{2} \hat{\Sigma}^{-1} + \frac{1}{2} \hat{\Sigma}^{-1} (Y - X\hat{B})^T V^{-1} (Y - X\hat{B}) \hat{\Sigma}^{-1} \\
\frac{n}{2} \hat{\Sigma}^{-1} &= \frac{1}{2} \hat{\Sigma}^{-1} (Y - X\hat{B})^T V^{-1} (Y - X\hat{B}) \hat{\Sigma}^{-1} \\
\hat{\Sigma}^{-1} &= \frac{1}{n} \hat{\Sigma}^{-1} (Y - X\hat{B})^T V^{-1} (Y - X\hat{B}) \hat{\Sigma}^{-1} \\
I_v &= \frac{1}{n} (Y - X\hat{B})^T V^{-1} (Y - X\hat{B}) \hat{\Sigma}^{-1} \\
\hat{\Sigma} &= \frac{1}{n} (Y - X\hat{B})^T V^{-1} (Y - X\hat{B})
\end{aligned} \tag{61}$$

Together, (59) and (61) constitute the MLE for the GLM.

Sources:

- original work

Metadata: ID: P7 | shortcut: glm-mle | author: JoramSoch | date: 2019-12-06, 10:40.

2 Poisson data

2.1 Poisson-distributed data

2.1.1 Maximum likelihood estimation

Theorem: Let $y = \{y_1, \dots, y_n\}$ be a set of observed counts independent and identically distributed according to a Poisson distribution (\rightarrow Definition “poiss”) with rate λ :

$$y_i \sim \text{Poiss}(\lambda), \quad i = 1, \dots, n. \quad (62)$$

Then, the maximum likelihood estimate (\rightarrow Definition “mle”) for the rate parameter λ is given by

$$\hat{\lambda} = \bar{y} \quad (63)$$

where \bar{y} is the sample mean (\rightarrow Proof “ev-sample”)

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i. \quad (64)$$

Proof: The likelihood function (\rightarrow Definition “lf”) for each observation is given by the probability mass function of the Poisson distribution (\rightarrow Proof “poiss-pdf”)

$$p(y_i|\lambda) = \text{Poiss}(y_i; \lambda) = \frac{\lambda^{y_i} \cdot \exp(-\lambda)}{y_i!} \quad (65)$$

and because observations are independent (\rightarrow Definition “ind”), the likelihood function for all observations is the product of the individual ones:

$$p(y|\lambda) = \prod_{i=1}^n p(y_i|\lambda) = \prod_{i=1}^n \frac{\lambda^{y_i} \cdot \exp(-\lambda)}{y_i!}. \quad (66)$$

Thus, the log-likelihood function (\rightarrow Definition “llf”) is

$$\text{LL}(\lambda) = \log p(y|\lambda) = \log \left[\prod_{i=1}^n \frac{\lambda^{y_i} \cdot \exp(-\lambda)}{y_i!} \right] \quad (67)$$

which can be developed into

$$\begin{aligned}
\text{LL}(\lambda) &= \sum_{i=1}^n \log \left[\frac{\lambda^{y_i} \cdot \exp(-\lambda)}{y_i!} \right] \\
&= \sum_{i=1}^n [y_i \cdot \log(\lambda) - \lambda - \log(y_i!)] \\
&= - \sum_{i=1}^n \lambda + \sum_{i=1}^n y_i \cdot \log(\lambda) - \sum_{i=1}^n \log(y_i!) \\
&= -n\lambda + \log(\lambda) \sum_{i=1}^n y_i - \sum_{i=1}^n \log(y_i!)
\end{aligned} \tag{68}$$

The derivatives of the log-likelihood with respect to λ are

$$\begin{aligned}
\frac{d\text{LL}(\lambda)}{d\lambda} &= \frac{1}{\lambda} \sum_{i=1}^n y_i - n \\
\frac{d^2\text{LL}(\lambda)}{d\lambda^2} &= -\frac{1}{\lambda^2} \sum_{i=1}^n y_i .
\end{aligned} \tag{69}$$

Setting the first derivative to zero, we obtain:

$$\begin{aligned}
\frac{d\text{LL}(\hat{\lambda})}{d\lambda} &= 0 \\
0 &= \frac{1}{\hat{\lambda}} \sum_{i=1}^n y_i - n \\
\hat{\lambda} &= \frac{1}{n} \sum_{i=1}^n y_i = \bar{y} .
\end{aligned} \tag{70}$$

Plugging this value into the second derivative, we confirm:

$$\begin{aligned}
\frac{d^2\text{LL}(\hat{\lambda})}{d\lambda^2} &= -\frac{1}{\bar{y}^2} \sum_{i=1}^n y_i \\
&= -\frac{n \cdot \bar{y}}{\bar{y}^2} \\
&= -\frac{n}{\bar{y}} < 0 .
\end{aligned} \tag{71}$$

This demonstrates that the estimate $\hat{\lambda} = \bar{y}$ maximizes the likelihood $p(y \mid \lambda)$.

Sources:

- original work

Metadata: ID: P27 | shortcut: poiss-mle | author: JoramSoch | date: 2020-01-20, 21:53.

2.2 Poisson distribution with exposure values

2.2.1 Conjugate prior distribution

Theorem: Let $y = \{y_1, \dots, y_n\}$ be a series of observed counts which are independently distributed according to a Poisson distribution (\rightarrow Definition “poiss”) with common rate λ and concurrent exposures $\{x_1, \dots, x_n\}$:

$$y_i \sim \text{Poiss}(\lambda x_i), \quad i = 1, \dots, n. \quad (72)$$

Then, the conjugate prior (\rightarrow Definition “prior-conj”) for the model parameter λ is a gamma distribution (\rightarrow Definition II/3.3.1):

$$p(\lambda) = \text{Gam}(\lambda; a_0, b_0). \quad (73)$$

Proof: With the probability mass function of the Poisson distribution (\rightarrow Proof “poiss-pmf”), the likelihood function (\rightarrow Definition “lf”) for each observation implied by (72) is given by

$$p(y_i|\lambda) = \text{Poiss}(y_i; \lambda x_i) = \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!} \quad (74)$$

and because observations are independent (\rightarrow Definition “ind”), the likelihood function for all observations is the product of the individual ones:

$$p(y|\lambda) = \prod_{i=1}^n p(y_i|\lambda) = \prod_{i=1}^n \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!}. \quad (75)$$

Resolving the product in the likelihood function, we have

$$\begin{aligned} p(y|\lambda) &= \prod_{i=1}^n \frac{x_i^{y_i}}{y_i!} \cdot \prod_{i=1}^n \lambda^{y_i} \cdot \prod_{i=1}^n \exp[-\lambda x_i] \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \cdot \lambda^{\sum_{i=1}^n y_i} \cdot \exp \left[-\lambda \sum_{i=1}^n x_i \right] \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \cdot \lambda^{n\bar{y}} \cdot \exp[-n\bar{x}\lambda] \end{aligned} \quad (76)$$

where \bar{y} and \bar{x} are the means (\rightarrow Proof “ev-sample”) of y and x respectively:

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i \\ \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i. \end{aligned} \quad (77)$$

In other words, the likelihood function is proportional to a power of λ times an exponential of λ :

$$p(y|\lambda) \propto \lambda^{n\bar{y}} \cdot \exp[-n\bar{x}\lambda] . \quad (78)$$

The same is true for a gamma distribution over λ

$$p(\lambda) = \text{Gam}(\lambda; a_0, b_0) \quad (79)$$

the probability density function of which (\rightarrow Proof II/3.3.2)

$$p(\lambda) = \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0\lambda] \quad (80)$$

exhibits the same proportionality

$$p(\lambda) \propto \lambda^{a_0-1} \cdot \exp[-b_0\lambda] \quad (81)$$

and is therefore conjugate relative to the likelihood.

Sources:

- Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB (2014): “Other standard single-parameter models”; in: *Bayesian Data Analysis*, 3rd edition, ch. 2.6, p. 45, eq. 2.14ff.; URL: <http://www.stat.columbia.edu/~gelman/book/>.

Metadata: ID: P41 | shortcut: poissexp-prior | author: JoramSoch | date: 2020-02-04, 14:11.

2.2.2 Posterior distribution

Theorem: Let $y = \{y_1, \dots, y_n\}$ be a series of observed counts which are independently distributed according to a Poisson distribution (\rightarrow Definition “poiss”) with common rate λ and concurrent exposures $\{x_1, \dots, x_n\}$:

$$y_i \sim \text{Poiss}(\lambda x_i), \quad i = 1, \dots, n . \quad (82)$$

Moreover, assume a gamma prior distribution (\rightarrow Proof III/2.2.1) over the model parameter λ :

$$p(\lambda) = \text{Gam}(\lambda; a_0, b_0) . \quad (83)$$

Then, the posterior distribution (\rightarrow Definition “post”) is also a gamma distribution (\rightarrow Definition II/3.3.1)

$$p(\lambda|y) = \text{Gam}(\lambda; a_n, b_n) \quad (84)$$

and the posterior hyperparameters (\rightarrow Definition “post-hyp”) are given by

$$\begin{aligned} a_n &= a_0 + n\bar{y} \\ b_n &= b_0 + n\bar{x} . \end{aligned} \quad (85)$$

Proof: With the probability mass function of the Poisson distribution (\rightarrow Proof “poiss-pmf”), the likelihood function (\rightarrow Definition “lf”) for each observation implied by (82) is given by

$$p(y_i|\lambda) = \text{Poiss}(y_i; \lambda x_i) = \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!} \quad (86)$$

and because observations are independent (\rightarrow Definition “ind”), the likelihood function for all observations is the product of the individual ones:

$$p(y|\lambda) = \prod_{i=1}^n p(y_i|\lambda) = \prod_{i=1}^n \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!} . \quad (87)$$

Combining the likelihood function (87) with the prior distribution (83), the joint likelihood (\rightarrow Definition “jl”) of the model is given by

$$\begin{aligned} p(y, \lambda) &= p(y|\lambda) p(\lambda) \\ &= \prod_{i=1}^n \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!} \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] . \end{aligned} \quad (88)$$

Resolving the product in the joint likelihood, we have

$$\begin{aligned} p(y, \lambda) &= \prod_{i=1}^n \frac{x_i^{y_i}}{y_i!} \prod_{i=1}^n \lambda^{y_i} \prod_{i=1}^n \exp[-\lambda x_i] \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \lambda^{\sum_{i=1}^n y_i} \exp \left[-\lambda \sum_{i=1}^n x_i \right] \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \lambda^{n\bar{y}} \exp[-n\bar{x}\lambda] \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{b_0^{a_0}}{\Gamma(a_0)} \cdot \lambda^{a_0+n\bar{y}-1} \cdot \exp[-(b_0 + n\bar{x})\lambda] \end{aligned} \quad (89)$$

where \bar{y} and \bar{x} are the means (\rightarrow Proof “ev-sample”) of y and x respectively:

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i \\ \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i . \end{aligned} \quad (90)$$

Note that the posterior distribution is proportional to the joint likelihood (\rightarrow Proof “post-jl”):

$$p(\lambda|y) \propto p(y, \lambda) . \quad (91)$$

Setting $a_n = a_0 + n\bar{y}$ and $b_n = b_0 + n\bar{x}$, the posterior distribution is therefore proportional to

$$p(\lambda|y) \propto \lambda^{a_n-1} \cdot \exp[-b_n\lambda] \quad (92)$$

which, when normalized to one, results in the probability density function of the gamma distribution (\rightarrow Proof II/3.3.2):

$$p(\lambda|y) = \frac{b_n^{a_n}}{\Gamma(a_n)} \lambda^{a_n-1} \exp[-b_n\lambda] = \text{Gam}(\lambda; a_n, b_n) . \quad (93)$$

Sources:

- Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB (2014): “Other standard single-parameter models”; in: *Bayesian Data Analysis*, 3rd edition, ch. 2.6, p. 45, eq. 2.15; URL: <http://www.stat.columbia.edu/~gelman/book/>.

Metadata: ID: P42 | shortcut: poissexp-post | author: JoramSoch | date: 2020-02-04, 14:42.

2.2.3 Log model evidence

Theorem: Let $y = \{y_1, \dots, y_n\}$ be a series of observed counts which are independently distributed according to a Poisson distribution (\rightarrow Definition “poiss”) with common rate λ and concurrent exposures $\{x_1, \dots, x_n\}$:

$$y_i \sim \text{Poiss}(\lambda x_i), \quad i = 1, \dots, n . \quad (94)$$

Moreover, assume a gamma prior distribution (\rightarrow Proof III/2.2.1) over the model parameter λ :

$$p(\lambda) = \text{Gam}(\lambda; a_0, b_0) . \quad (95)$$

Then, the log model evidence (\rightarrow Definition “lme”) for this model is

$$\begin{aligned} \log p(y|m) &= \sum_{i=1}^n y_i \log(x_i) - \sum_{i=1}^n \log y_i! + \\ &\quad \log \Gamma(a_n) - \log \Gamma(a_0) + a_0 \log b_0 - a_n \log b_n . \end{aligned} \quad (96)$$

where the posterior hyperparameters (\rightarrow Definition “post-hyp”) are given by

$$\begin{aligned} a_n &= a_0 + n\bar{y} \\ b_n &= b_0 + n\bar{x} . \end{aligned} \quad (97)$$

Proof: With the probability mass function of the Poisson distribution (\rightarrow Proof “poiss-pmf”), the likelihood function (\rightarrow Definition “lf”) for each observation implied by (94) is given by

$$p(y_i|\lambda) = \text{Pois}(y_i; \lambda x_i) = \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!} \quad (98)$$

and because observations are independent (\rightarrow Definition “ind”), the likelihood function for all observations is the product of the individual ones:

$$p(y|\lambda) = \prod_{i=1}^n p(y_i|\lambda) = \prod_{i=1}^n \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!} . \quad (99)$$

Combining the likelihood function (99) with the prior distribution (95), the joint likelihood (\rightarrow Definition “jl”) of the model is given by

$$\begin{aligned} p(y, \lambda) &= p(y|\lambda) p(\lambda) \\ &= \prod_{i=1}^n \frac{(\lambda x_i)^{y_i} \cdot \exp[-\lambda x_i]}{y_i!} \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] . \end{aligned} \quad (100)$$

Resolving the product in the joint likelihood, we have

$$\begin{aligned} p(y, \lambda) &= \prod_{i=1}^n \frac{x_i^{y_i}}{y_i!} \prod_{i=1}^n \lambda^{y_i} \prod_{i=1}^n \exp[-\lambda x_i] \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \lambda^{\sum_{i=1}^n y_i} \exp \left[-\lambda \sum_{i=1}^n x_i \right] \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \lambda^{n\bar{y}} \exp[-n\bar{x}\lambda] \cdot \frac{b_0^{a_0}}{\Gamma(a_0)} \lambda^{a_0-1} \exp[-b_0 \lambda] \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{b_0^{a_0}}{\Gamma(a_0)} \cdot \lambda^{a_0+n\bar{y}-1} \cdot \exp[-(b_0 + n\bar{x})\lambda] \end{aligned} \quad (101)$$

where \bar{y} and \bar{x} are the means (\rightarrow Proof “ev-sample”) of y and x respectively:

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i \\ \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i . \end{aligned} \quad (102)$$

Note that the model evidence is the marginal density of the joint likelihood (\rightarrow Definition “ml”):

$$p(y) = \int p(y, \lambda) d\lambda . \quad (103)$$

Setting $a_n = a_0 + n\bar{y}$ and $b_n = b_0 + n\bar{x}$, the joint likelihood can also be written as

$$p(y, \lambda) = \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{b_0^{a_0}}{\Gamma(a_0)} \frac{\Gamma(a_n)}{b_n^{a_n}} \cdot \frac{b_n^{a_n}}{\Gamma(a_n)} \lambda^{a_n-1} \exp[-b_n \lambda] . \quad (104)$$

Using the probability density function of the gamma distribution (\rightarrow Proof II/3.3.2), λ can now be integrated out easily

$$\begin{aligned} p(y) &= \int \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{b_0^{a_0}}{\Gamma(a_0)} \frac{\Gamma(a_n)}{b_n^{a_n}} \cdot \frac{b_n^{a_n}}{\Gamma(a_n)} \lambda^{a_n-1} \exp[-b_n \lambda] d\lambda \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{\Gamma(a_n)}{\Gamma(a_0)} \frac{b_0^{a_0}}{b_n^{a_n}} \int \text{Gam}(\lambda; a_n, b_n) d\lambda \\ &= \prod_{i=1}^n \left(\frac{x_i^{y_i}}{y_i!} \right) \frac{\Gamma(a_n)}{\Gamma(a_0)} \frac{b_0^{a_0}}{b_n^{a_n}} , \end{aligned} \quad (105)$$

such that the log model evidence (\rightarrow Definition “lme”) is shown to be

$$\begin{aligned} \log p(y|m) &= \sum_{i=1}^n y_i \log(x_i) - \sum_{i=1}^n \log y_i! + \\ &\quad \log \Gamma(a_n) - \log \Gamma(a_0) + a_0 \log b_0 - a_n \log b_n . \end{aligned} \quad (106)$$

Sources:

- original work

Metadata: ID: P43 | shortcut: poissexp-lme | author: JoramSoch | date: 2020-02-04, 15:12.

3 Probability data

3.1 Beta-distributed data

3.1.1 Method of moments

Theorem: Let $y = \{y_1, \dots, y_n\}$ be a set of observed counts independent and identically distributed (\rightarrow Definition “iid”) according to a beta distribution (\rightarrow Definition “beta”) with shapes α and β :

$$y_i \sim \text{Bet}(\alpha, \beta), \quad i = 1, \dots, n. \quad (107)$$

Then, the method-of-moments estimates (\rightarrow Definition “mom”) for the shape parameters α and β are given by

$$\begin{aligned} \hat{\alpha} &= \bar{y} \left(\frac{\bar{y}(1 - \bar{y})}{\bar{v}} - 1 \right) \\ \hat{\beta} &= (1 - \bar{y}) \left(\frac{\bar{y}(1 - \bar{y})}{\bar{v}} - 1 \right) \end{aligned} \quad (108)$$

where \bar{y} is the sample mean (\rightarrow Proof “ev-sample”) and \bar{v} is the unbiased sample variance (\rightarrow Proof “var-unbias”):

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i \\ \bar{v} &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2. \end{aligned} \quad (109)$$

Proof: Mean (\rightarrow Proof “beta-mean”) and variance (\rightarrow Proof “beta-var”) of the beta distribution (\rightarrow Definition “beta”) in terms of the parameters α and β are given by

$$\begin{aligned} E(X) &= \frac{\alpha}{\alpha + \beta} \\ \text{Var}(X) &= \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}. \end{aligned} \quad (110)$$

Thus, matching the moments (\rightarrow Definition “mom”) requires us to solve the following equation system for α and β :

$$\begin{aligned} \bar{y} &= \frac{\alpha}{\alpha + \beta} \\ \bar{v} &= \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}. \end{aligned} \quad (111)$$

From the first equation, we can deduce:

$$\begin{aligned}
\bar{y}(\alpha + \beta) &= \alpha \\
\alpha\bar{y} + \beta\bar{y} &= \alpha \\
\beta\bar{y} &= \alpha - \alpha\bar{y} \\
\beta &= \frac{\alpha}{\bar{y}} - \alpha \\
\beta &= \alpha \left(\frac{1}{\bar{y}} - 1 \right) .
\end{aligned} \tag{112}$$

If we define $q = 1/\bar{y} - 1$ and plug (112) into the second equation, we have:

$$\begin{aligned}
\bar{v} &= \frac{\alpha \cdot \alpha q}{(\alpha + \alpha q)^2(\alpha + \alpha q + 1)} \\
&= \frac{\alpha^2 q}{(\alpha(1 + q))^2(\alpha(1 + q) + 1)} \\
&= \frac{q}{(1 + q)^2(\alpha(1 + q) + 1)} \\
&= \frac{q}{\alpha(1 + q)^3 + (1 + q)^2} \\
q &= \bar{v} [\alpha(1 + q)^3 + (1 + q)^2] .
\end{aligned} \tag{113}$$

Noting that $1 + q = 1/\bar{y}$ and $q = (1 - \bar{y})/\bar{y}$, one obtains for α :

$$\begin{aligned}
\frac{1 - \bar{y}}{\bar{y}} &= \bar{v} \left[\frac{\alpha}{\bar{y}^3} + \frac{1}{\bar{y}^2} \right] \\
\frac{1 - \bar{y}}{\bar{y} \bar{v}} &= \frac{\alpha}{\bar{y}^3} + \frac{1}{\bar{y}^2} \\
\frac{\bar{y}^3(1 - \bar{y})}{\bar{y} \bar{v}} &= \alpha + \bar{y} \\
\alpha &= \frac{\bar{y}^2(1 - \bar{y})}{\bar{v}} - \bar{y} \\
&= \bar{y} \left(\frac{\bar{y}(1 - \bar{y})}{\bar{v}} - 1 \right) .
\end{aligned} \tag{114}$$

Plugging this into equation (112), one obtains for β :

$$\begin{aligned}
\beta &= \bar{y} \left(\frac{\bar{y}(1 - \bar{y})}{\bar{v}} - 1 \right) \cdot \left(\frac{1 - \bar{y}}{\bar{y}} \right) \\
&= (1 - \bar{y}) \left(\frac{\bar{y}(1 - \bar{y})}{\bar{v}} - 1 \right) .
\end{aligned} \tag{115}$$

Together, (114) and (115) constitute the method-of-moment estimates of α and β .

Sources:

- Wikipedia (2020): “Beta distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-20; URL: https://en.wikipedia.org/wiki/Beta_distribution#Method_of_moments.

Metadata: ID: P28 | shortcut: beta-mom | author: JoramSoch | date: 2020-01-22, 02:53.

4 Categorical data

4.1 Binomial observations

4.1.1 Conjugate prior distribution

Theorem: Let y be the number of successes resulting from n independent trials with unknown success probability p , such that y follows a binomial distribution (\rightarrow Definition “bin”):

$$y \sim \text{Bin}(n, p) . \quad (116)$$

Then, the conjugate prior (\rightarrow Definition “prior-conj”) for the model parameter p is a beta distribution (\rightarrow Definition “beta”):

$$p(p) = \text{Bet}(p; \alpha_0, \beta_0) . \quad (117)$$

Proof: With the probability mass function of the binomial distribution (\rightarrow Proof “bin-pmf”), the likelihood function (\rightarrow Definition “lf”) implied by (116) is given by

$$p(y|p) = \binom{n}{y} p^y (1-p)^{n-y} . \quad (118)$$

In other words, the likelihood function is proportional to a power of p times a power of $(1-p)$:

$$p(y|p) \propto p^y (1-p)^{n-y} . \quad (119)$$

The same is true for a beta distribution over p

$$p(p) = \text{Bet}(p; \alpha_0, \beta_0) \quad (120)$$

the probability density function of which (\rightarrow Proof “beta-pdf”)

$$p(p) = \frac{1}{B(\alpha_0, \beta_0)} p^{\alpha_0-1} (1-p)^{\beta_0-1} \quad (121)$$

exhibits the same proportionality

$$p(p) \propto p^{\alpha_0-1} (1-p)^{\beta_0-1} \quad (122)$$

and is therefore conjugate relative to the likelihood.

Sources:

- Wikipedia (2020): “Binomial distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-23; URL: https://en.wikipedia.org/wiki/Binomial_distribution#Estimation_of_parameters.

Metadata: ID: P29 | shortcut: bin-prior | author: JoramSoch | date: 2020-01-23, 23:38.

4.1.2 Posterior distribution

Theorem: Let y be the number of successes resulting from n independent trials with unknown success probability p , such that y follows a binomial distribution (\rightarrow Definition “bin”):

$$y \sim \text{Bin}(n, p) . \quad (123)$$

Moreover, assume a beta prior distribution (\rightarrow Proof III/4.1.1) over the model parameter p :

$$p(p) = \text{Bet}(p; \alpha_0, \beta_0) . \quad (124)$$

Then, the posterior distribution (\rightarrow Definition “post”) is also a beta distribution (\rightarrow Definition “beta”)

$$p(p|y) = \text{Bet}(p; \alpha_n, \beta_n) . \quad (125)$$

and the posterior hyperparameters (\rightarrow Definition “post-hyp”) are given by

$$\begin{aligned} \alpha_n &= \alpha_0 + y \\ \beta_n &= \beta_0 + (n - y) . \end{aligned} \quad (126)$$

Proof: With the probability mass function of the binomial distribution (\rightarrow Proof “bin-pmf”), the likelihood function (\rightarrow Definition “lf”) implied by (123) is given by

$$p(y|p) = \binom{n}{y} p^y (1 - p)^{n-y} . \quad (127)$$

Combining the likelihood function (127) with the prior distribution (124), the joint likelihood (\rightarrow Definition “jl”) of the model is given by

$$\begin{aligned} p(y, p) &= p(y|p) p(p) \\ &= \binom{n}{y} p^y (1 - p)^{n-y} \cdot \frac{1}{B(\alpha_0, \beta_0)} p^{\alpha_0-1} (1 - p)^{\beta_0-1} \\ &= \frac{1}{B(\alpha_0, \beta_0)} \binom{n}{y} p^{\alpha_0+y-1} (1 - p)^{\beta_0+(n-y)-1} . \end{aligned} \quad (128)$$

Note that the posterior distribution is proportional to the joint likelihood (\rightarrow Proof “post-jl”):

$$p(p|y) \propto p(y, p) . \quad (129)$$

Setting $\alpha_n = \alpha_0 + y$ and $\beta_n = \beta_0 + (n - y)$, the posterior distribution is therefore proportional to

$$p(p|y) \propto p^{\alpha_n-1} (1 - p)^{\beta_n-1} \quad (130)$$

which, when normalized to one, results in the probability density function of the beta distribution (\rightarrow Proof “beta-pdf”):

$$p(p|y) = \frac{1}{B(\alpha_n, \beta_n)} p^{\alpha_n-1} (1-p)^{\beta_n-1} = \text{Bet}(p; \alpha_n, \beta_n) . \quad (131)$$

Sources:

- Wikipedia (2020): “Binomial distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-23; URL: https://en.wikipedia.org/wiki/Binomial_distribution#Estimation_of_parameters.

Metadata: ID: P30 | shortcut: bin-post | author: JoramSoch | date: 2020-01-24, 00:20.

4.1.3 Log model evidence

Theorem: Let y be the number of successes resulting from n independent trials with unknown success probability p , such that y follows a binomial distribution (\rightarrow Definition “bin”):

$$y \sim \text{Bin}(n, p) . \quad (132)$$

Moreover, assume a beta prior distribution (\rightarrow Proof III/4.1.1) over the model parameter p :

$$p(p) = \text{Bet}(p; \alpha_0, \beta_0) . \quad (133)$$

Then, the log model evidence (\rightarrow Definition “lme”) for this model is

$$\log p(y|m) = \log \binom{n}{y} + \log B(\alpha_n, \beta_n) - \log B(\alpha_0, \beta_0) \quad (134)$$

where the posterior hyperparameters (\rightarrow Definition “post-hyp”) are given by

$$\begin{aligned} \alpha_n &= \alpha_0 + y \\ \beta_n &= \beta_0 + (n - y) . \end{aligned} \quad (135)$$

Proof: With the probability mass function of the binomial distribution (\rightarrow Proof “bin-pmf”), the likelihood function (\rightarrow Definition “lf”) implied by (132) is given by

$$p(y|p) = \binom{n}{y} p^y (1-p)^{n-y} . \quad (136)$$

Combining the likelihood function (136) with the prior distribution (133), the joint likelihood (\rightarrow Definition “jl”) of the model is given by

$$\begin{aligned}
p(y, p) &= p(y|p) p(p) \\
&= \binom{n}{y} p^y (1-p)^{n-y} \cdot \frac{1}{B(\alpha_0, \beta_0)} p^{\alpha_0-1} (1-p)^{\beta_0-1} \\
&= \binom{n}{y} \frac{1}{B(\alpha_0, \beta_0)} p^{\alpha_0+y-1} (1-p)^{\beta_0+(n-y)-1} .
\end{aligned} \tag{137}$$

Note that the model evidence is the marginal density of the joint likelihood (\rightarrow Definition “ml”):

$$p(y) = \int p(y, p) dp . \tag{138}$$

Setting $\alpha_n = \alpha_0 + y$ and $\beta_n = \beta_0 + (n - y)$, the joint likelihood can also be written as

$$p(y, p) = \binom{n}{y} \frac{1}{B(\alpha_0, \beta_0)} \frac{B(\alpha_n, \beta_n)}{1} \frac{1}{B(\alpha_n, \beta_n)} p^{\alpha_n-1} (1-p)^{\beta_n-1} . \tag{139}$$

Using the probability density function of the beta distribution (\rightarrow Proof “beta-pdf”), p can now be integrated out easily

$$\begin{aligned}
p(y) &= \int \binom{n}{y} \frac{1}{B(\alpha_0, \beta_0)} \frac{B(\alpha_n, \beta_n)}{1} \frac{1}{B(\alpha_n, \beta_n)} p^{\alpha_n-1} (1-p)^{\beta_n-1} dp \\
&= \binom{n}{y} \frac{B(\alpha_n, \beta_n)}{B(\alpha_0, \beta_0)} \int \text{Bet}(p; \alpha_n, \beta_n) dp \\
&= \binom{n}{y} \frac{B(\alpha_n, \beta_n)}{B(\alpha_0, \beta_0)} ,
\end{aligned} \tag{140}$$

such that the log model evidence (\rightarrow Definition “lme”) is shown to be

$$\log p(y|m) = \log \binom{n}{y} + \log B(\alpha_n, \beta_n) - \log B(\alpha_0, \beta_0) . \tag{141}$$

Sources:

- Wikipedia (2020): “Beta-binomial distribution”; in: *Wikipedia, the free encyclopedia*, retrieved on 2020-01-24; URL: https://en.wikipedia.org/wiki/Beta-binomial_distribution#Motivation_and_derivation.

Metadata: ID: P31 | shortcut: bin-lme | author: JoramSoch | date: 2020-01-24, 00:44.

Chapter IV

Model Selection

1 Goodness-of-fit measures

1.1 R-squared

1.1.1 Derivation of R^2 and adjusted R^2

Theorem: Given a linear regression model (\rightarrow Definition “mlr”)

$$y = X\beta + \varepsilon, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2) \quad (1)$$

with n independent observations and p independent variables,

1) the coefficient of determination is given by

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} \quad (2)$$

2) the adjusted coefficient of determination is

$$R^2_{\text{adj}} = 1 - \frac{\text{RSS}/(n-p)}{\text{TSS}/(n-1)} \quad (3)$$

where the residual (\rightarrow Definition “rss”) and total sum of squares (\rightarrow Definition “tss”) are

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad \hat{y} = X\hat{\beta} \\ \text{TSS} &= \sum_{i=1}^n (y_i - \bar{y})^2, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \end{aligned} \quad (4)$$

where X is the $n \times p$ design matrix and $\hat{\beta}$ are the ordinary least squares (\rightarrow Definition “mlr-ols”) estimates.

Proof: The coefficient of determination R^2 is defined as (\rightarrow Definition “rsq”) the proportion of the variance explained by the independent variables, relative to the total variance in the data.

1) If we define the explained sum of squares (\rightarrow Definition “ess”) as

$$\text{ESS} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2, \quad (5)$$

then R^2 is given by

$$R^2 = \frac{\text{ESS}}{\text{TSS}}. \quad (6)$$

which is equal to

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}, \quad (7)$$

because (\rightarrow Proof “mlr-pss”) $TSS = ESS + RSS$.

2) Using (4), the coefficient of determination can be also written as:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2} . \quad (8)$$

If we replace the variance estimates by their unbiased estimators (\rightarrow Proof “resvar-bias”), we obtain

$$R_{\text{adj}}^2 = 1 - \frac{\frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2}{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{RSS/df_r}{TSS/df_t} \quad (9)$$

where $df_r = n - p$ and $df_t = n - 1$ are the residual and total degrees of freedom (\rightarrow Definition “dof”).

This gives the adjusted R^2 which adjusts R^2 for the number of explanatory variables.

Sources:

- Wikipedia (2019): “Coefficient of determination”; in: *Wikipedia, the free encyclopedia*, retrieved on 2019-12-06; URL: https://en.wikipedia.org/wiki/Coefficient_of_determination#Adjusted_R2.

Metadata: ID: P8 | shortcut: rsq-der | author: JoramSoch | date: 2019-12-06, 11:19.

1.1.2 Relationship to maximum log-likelihood

Theorem: Given a linear regression model (\rightarrow Definition “mlr”) with independent observations

$$y = X\beta + \varepsilon, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad (10)$$

the coefficient of determination (\rightarrow Definition “rsq”) can be expressed in terms of the maximum log-likelihood (\rightarrow Definition “mll”) as

$$R^2 = 1 - (\exp[\Delta\text{MLL}])^{-2/n} \quad (11)$$

where n is the number of observations and ΔMLL is the difference in maximum log-likelihood between the model given by (10) and a linear regression model with only a constant regressor.

Proof: First, we express the maximum log-likelihood (\rightarrow Definition “mll”) (MLL) of a linear regression model in terms of its residual sum of squares (\rightarrow Definition “rss”) (RSS). The model in (10) implies the following log-likelihood function (\rightarrow Definition “llf”)

$$\text{LL}(\beta, \sigma^2) = \log p(y|\beta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta), \quad (12)$$

such that maximum likelihood estimates are (\rightarrow Proof “mlr-mle”)

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (13)$$

$$\hat{\sigma}^2 = \frac{1}{n} (y - X\hat{\beta})^T (y - X\hat{\beta}) \quad (14)$$

and the residual sum of squares (\rightarrow Definition “rss”) is

$$\text{RSS} = \sum_{i=1}^n \hat{\varepsilon}_i = \hat{\varepsilon}^T \hat{\varepsilon} = (y - X\hat{\beta})^T (y - X\hat{\beta}) = n \cdot \hat{\sigma}^2 . \quad (15)$$

Since $\hat{\beta}$ and $\hat{\sigma}^2$ are maximum likelihood estimates (\rightarrow Definition “mle”), plugging them into the log-likelihood function gives the maximum log-likelihood:

$$\text{MLL} = \text{LL}(\hat{\beta}, \hat{\sigma}^2) = -\frac{n}{2} \log(2\pi\hat{\sigma}^2) - \frac{1}{2\hat{\sigma}^2} (y - X\hat{\beta})^T (y - X\hat{\beta}) . \quad (16)$$

With (15) for the first $\hat{\sigma}^2$ and (14) for the second $\hat{\sigma}^2$, the MLL becomes

$$\text{MLL} = -\frac{n}{2} \log(\text{RSS}) - \frac{n}{2} \log\left(\frac{2\pi}{n}\right) - \frac{n}{2} . \quad (17)$$

Second, we establish the relationship between maximum log-likelihood (MLL) and coefficient of determination (R^2). Consider the two models

$$\begin{aligned} m_0 : X_0 &= 1_n \\ m_1 : X_1 &= X \end{aligned} \quad (18)$$

For m_1 , the residual sum of squares is given by (15); and for m_0 , the residual sum of squares is equal to the total sum of squares (\rightarrow Definition “tss”):

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2 . \quad (19)$$

Using (17), we can therefore write

$$\Delta\text{MLL} = \text{MLL}(m_1) - \text{MLL}(m_0) = -\frac{n}{2} \log(\text{RSS}) + \frac{n}{2} \log(\text{TSS}) . \quad (20)$$

Exponentiating both sides of the equation, we have:

$$\begin{aligned} \exp[\Delta\text{MLL}] &= \exp\left[-\frac{n}{2} \log(\text{RSS}) + \frac{n}{2} \log(\text{TSS})\right] \\ &= (\exp[\log(\text{RSS}) - \log(\text{TSS})])^{-n/2} \\ &= \left(\frac{\exp[\log(\text{RSS})]}{\exp[\log(\text{TSS})]}\right)^{-n/2} \\ &= \left(\frac{\text{RSS}}{\text{TSS}}\right)^{-n/2} . \end{aligned} \quad (21)$$

Taking both sides to the power of $-2/n$ and subtracting from 1, we have

$$\begin{aligned} (\exp[\Delta\text{MLL}])^{-2/n} &= \frac{\text{RSS}}{\text{TSS}} \\ 1 - (\exp[\Delta\text{MLL}])^{-2/n} &= 1 - \frac{\text{RSS}}{\text{TSS}} = R^2 \end{aligned} \tag{22}$$

which proves the identity given above.

Sources:

- original work

Metadata: ID: P14 | shortcut: rsq-ml | author: JoramSoch | date: 2020-01-08, 04:46.

2 Classical information criteria

2.1 Bayesian information criterion

2.1.1 Derivation

Theorem: Let $p(y | \theta, m)$ be the likelihood function (\rightarrow Definition “lf”) of a generative model (\rightarrow Definition “gm”) $m \in \mathcal{M}$ with model parameters $\theta \in \Theta$ describing measured data $y \in \mathbb{R}^n$. Let $p(\theta | m)$ be a prior distribution (\rightarrow Definition “prior”) on the model parameters. Assume that likelihood function and prior density are twice differentiable. Then, as the number of data points goes to infinity, an approximation to the log-marginal likelihood (\rightarrow Definition “ml”) $\log p(y | m)$, up to constant terms not depending on the model, is given by the Bayesian information criterion (\rightarrow Definition “bic”) (BIC) as

$$-2 \log p(y | m) \approx \text{BIC}(m) = -2 \log p(y | \hat{\theta}, m) + p \log n \quad (23)$$

where $\hat{\theta}$ is the maximum likelihood estimator (\rightarrow Definition “mle”) (MLE) of θ , n is the number of data points and p is the number of model parameters.

Proof: Let $\text{LL}(\theta)$ be the log-likelihood function (\rightarrow Definition “llf”)

$$\text{LL}(\theta) = \log p(y | \theta, m) \quad (24)$$

and define the functions g and h as follows:

$$\begin{aligned} g(\theta) &= p(\theta | m) \\ h(\theta) &= \frac{1}{n} \text{LL}(\theta) . \end{aligned} \quad (25)$$

Then, the marginal likelihood (\rightarrow Definition “ml”) can be written as follows:

$$\begin{aligned} p(y | m) &= \int_{\Theta} p(y | \theta, m) p(\theta | m) d\theta \\ &= \int_{\Theta} \exp [n h(\theta)] g(\theta) d\theta . \end{aligned} \quad (26)$$

This is an integral suitable for Laplace approximation which states that

$$\int_{\Theta} \exp [n h(\theta)] g(\theta) d\theta = \left(\sqrt{\frac{2\pi}{n}} \right)^p \exp [n h(\theta_0)] \left(g(\theta_0) |J(\theta_0)|^{-1/2} + O(1/n) \right) \quad (27)$$

where θ_0 is the value that maximizes $h(\theta)$ and $J(\theta_0)$ is the Hessian matrix evaluated at θ_0 . In our case, we have $h(\theta) = 1/n \text{LL}(\theta)$ such that θ_0 is the maximum likelihood estimator $\hat{\theta}$:

$$\hat{\theta} = \arg \max_{\theta} \text{LL}(\theta) . \quad (28)$$

With this, (27) can be applied to (26) using (25) to give:

$$p(y|m) \approx \left(\sqrt{\frac{2\pi}{n}} \right)^p p(y|\hat{\theta}, m) p(\hat{\theta}|m) \left| J(\hat{\theta}) \right|^{-1/2}. \quad (29)$$

Logarithmizing and multiplying with -2 , we have:

$$-2 \log p(y|m) \approx -2 \text{LL}(\hat{\theta}) + p \log n - p \log(2\pi) - 2 \log p(\hat{\theta}|m) + \log \left| J(\hat{\theta}) \right|. \quad (30)$$

As $n \rightarrow \infty$, the last three terms are $O_p(1)$ and can therefore be ignored when comparing between models $\mathcal{M} = \{m_1, \dots, m_M\}$ and using $p(y | m_j)$ to compute posterior model probabilities (\rightarrow Definition “led-pmp”) $p(m_j | y)$. With that, the BIC is given as

$$\text{BIC}(m) = -2 \log p(y|\hat{\theta}, m) + p \log n. \quad (31)$$

Sources:

- Claeskens G, Hjort NL (2008): “The Bayesian information criterion”; in: *Model Selection and Model Averaging*, ch. 3.2, pp. 78-81; URL: <https://www.cambridge.org/core/books/model-selection-and-model-averaging/E6F1EC77279D1223423BB64FC3A12C37>; DOI: 10.1017/CBO9780511790485.

Metadata: ID: P32 | shortcut: bic-der | author: JoramSoch | date: 2020-01-26, 23:36.

3 Bayesian model selection

3.1 Log model evidence

3.1.1 Derivation

Theorem: Let $p(y \mid \theta, m)$ be a likelihood function (\rightarrow Definition “lf”) of a generative model (\rightarrow Definition “gm”) m for making inferences on model parameters θ given measured data y . Moreover, let $p(\theta \mid m)$ be a prior distribution (\rightarrow Definition “prior”) on model parameters θ . Then, the log model evidence (\rightarrow Definition “lme”) (LME), also called marginal log-likelihood,

$$\text{LME}(m) = \log p(y|m) , \quad (32)$$

can be expressed

1) as

$$\text{LME}(m) = \log \int p(y|\theta, m) p(\theta|m) d\theta \quad (33)$$

2) or

$$\text{LME}(m) = \log p(y|\theta, m) + \log p(\theta|m) - \log p(\theta|y, m) . \quad (34)$$

Proof:

1) The first expression is a simple consequence of the law of marginal probability (\rightarrow Proof “lmp”) for continuous variables according to which

$$p(y|m) = \int p(y|\theta, m) p(\theta|m) d\theta \quad (35)$$

which, when logarithmized, gives

$$\text{LME}(m) = \log p(y|m) = \log \int p(y|\theta, m) p(\theta|m) d\theta . \quad (36)$$

2) The second expression can be derived from Bayes’ theorem (\rightarrow Proof I/1.2.1) which makes a statement about the posterior distribution (\rightarrow Definition “post”):

$$p(\theta|y, m) = \frac{p(y|\theta, m) p(\theta|m)}{p(y|m)} . \quad (37)$$

Rearranging for $p(y \mid m)$ and logarithmizing, we have:

$$\begin{aligned} \text{LME}(m) = \log p(y|m) &= \log \frac{p(y|\theta, m) p(\theta|m)}{p(\theta|y, m)} \\ &= \log p(y|\theta, m) + \log p(\theta|m) - \log p(\theta|y, m) . \end{aligned} \quad (38)$$

Sources:

- original work

Metadata: ID: P13 | shortcut: lme-der | author: JoramSoch | date: 2020-01-06, 21:27.

3.1.2 Partition into accuracy and complexity

Theorem: The log model evidence (\rightarrow Definition “lme”) can be partitioned into accuracy and complexity

$$\text{LME}(m) = \text{Acc}(m) - \text{Com}(m) \quad (39)$$

where the accuracy term is the posterior expectation of the log-likelihood function (\rightarrow Definition “lf”)

$$\text{Acc}(m) = \langle p(y|\theta, m) \rangle_{p(\theta|y, m)} \quad (40)$$

and the complexity penalty is the Kullback-Leibler divergence (\rightarrow Definition “kl”) of posterior (\rightarrow Definition “post”) from prior (\rightarrow Definition “prior”)

$$\text{Com}(m) = \text{KL} [p(\theta|y, m) || p(\theta|m)] . \quad (41)$$

Proof: We consider Bayesian inference on data y using model m with parameters θ . Then, Bayes’ theorem (\rightarrow Proof I/1.2.1) makes a statement about the posterior distribution, i.e. the probability of parameters, given the data and the model:

$$p(\theta|y, m) = \frac{p(y|\theta, m) p(\theta|m)}{p(y|m)} . \quad (42)$$

Rearranging this for the model evidence (\rightarrow Proof IV/3.1.1), we have:

$$p(y|m) = \frac{p(y|\theta, m) p(\theta|m)}{p(\theta|y, m)} . \quad (43)$$

Logarithmizing both sides of the equation, we obtain:

$$\log p(y|m) = \log p(y|\theta, m) - \log \frac{p(\theta|y, m)}{p(\theta|m)} . \quad (44)$$

Now taking the expectation over the posterior distribution yields:

$$\log p(y|m) = \int p(\theta|y, m) \log p(y|\theta, m) d\theta - \int p(\theta|y, m) \log \frac{p(\theta|y, m)}{p(\theta|m)} d\theta . \quad (45)$$

By definition, the left-hand side is the log model evidence and the terms on the right-hand side correspond to the posterior expectation of the log-likelihood function and the Kullback-Leibler divergence of posterior from prior

$$\text{LME}(m) = \langle p(y|\theta, m) \rangle_{p(\theta|y, m)} - \text{KL} [p(\theta|y, m) || p(\theta|m)] \quad (46)$$

which proofs the partition given by (39).

Sources:

- Penny et al. (2007): “Bayesian Comparison of Spatially Regularised General Linear Models”; in: *Human Brain Mapping*, vol. 28, pp. 275–293; URL: <https://onlinelibrary.wiley.com/doi/full/10.1002/hbm.20327>; DOI: 10.1002/hbm.20327.
- Soch et al. (2016): “How to avoid mismodelling in GLM-based fMRI data analysis: cross-validated Bayesian model selection”; in: *NeuroImage*, vol. 141, pp. 469–489; URL: <https://www.sciencedirect.com/science/article/pii/S1053811916303615>; DOI: 10.1016/j.neuroimage.

Metadata: ID: P3 | shortcut: lme-anc | author: JoramSoch | date: 2019-09-27, 16:13.