

Ανασκόπηση Στοιχείων Πιθανοτήτων

Αθανάσιος Ροντογιάννης
Αν. Καθηγητής ΣΗΜΜΥ-ΕΜΠ

Συμβολισμοί (Notation)

- Διανύσματα/πίνακες και ανάστροφοί τους:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_l \end{bmatrix}, \quad A = \begin{bmatrix} a_{11} & \cdots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{l1} & \cdots & a_{lm} \end{bmatrix}, \quad \mathbf{x}^T = [x_1 \quad x_2 \quad \cdots \quad x_l], \quad B = A^T \Leftrightarrow b_{ij} = a_{ji}, \quad A = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m], \quad A^T = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix}$$

- Ένας διαγώνιος πίνακας με στοιχεία a_1, a_2, \dots, a_l στη διαγώνιό του θα συμβολίζεται ως $A = \text{diag}\{a_1, a_2, \dots, a_l\}$
- $I = \text{diag}\{1, 1, \dots, 1\}$ είναι ο μοναδιαίος πίνακας
- Το ίχνος (άθροισμα διαγώνιων στοιχείων) ενός τετραγωνικού πίνακα A συμβολίζεται ως $\text{trace}\{A\}$ και η ορίζουσά του ως $|A|$.
- Ακολουθίες αριθμών (διανυσμάτων) συμβολίζονται ως x_n (\mathbf{x}_n) ανάλογα με την περίπτωση
- Συμβολισμός συναρτήσεων: $f, f(x), f(\cdot)$
- Τα στοιχεία των διανυσμάτων/πινάκων και οι βαθμωτές ποσότητες που θα χρησιμοποιήσουμε θεωρούνται γενικά πραγματικοί αριθμοί.
- Εσωτερικό γινόμενο διανυσμάτων $\mathbf{x}, \mathbf{y} \in \mathbb{R}^l$:

$$\mathbf{x}^T \mathbf{y} = \sum_{i=1}^l x_i y_i \equiv \mathbf{y}^T \mathbf{x}$$

Χρήσιμες σχέσεις

Αν A, B, C πίνακες με κατάλληλες διαστάσεις και ιδιότητες, θα είναι:

- $(AB)^T = B^T A^T$
- $(AB)^{-1} = B^{-1} A^{-1}$
- $(A^T)^{-1} = (A^{-1})^T$
- $\text{trace}\{AB\} = \text{trace}\{BA\}$
- $\text{trace}\{ABC\} = \text{trace}\{CAB\} = \text{trace}\{BCA\}$
- Αν $A = \mathbf{a}\mathbf{b}^T$ (εξωτερικό γινόμενο των διανυσμάτων \mathbf{a}, \mathbf{b}) εύκολα προκύπτει:

$$\text{trace}(A) = \text{trace}(\mathbf{b}^T \mathbf{a}) = \mathbf{b}^T \mathbf{a} = \mathbf{a}^T \mathbf{b}$$

- $|AB| = |A||B|$
- $|A^{-1}| = 1/|A|$

Ανάδελτα συνάρτησης - Ιδιότητες

Έστω μια συνάρτηση $f(\mathbf{x})$ μιας διανυσματικής ποσότητας \mathbf{x} . Το ανάδελτα ή παράγωγος της f ως προς \mathbf{x} ορίζεται ως εξής:

$$\nabla_{\mathbf{x}} f \equiv \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_l} \end{bmatrix}$$

Μπορούν ναδειχτούν τα παρακάτω:

- $\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}$
- $\frac{\partial \mathbf{x}^T A \mathbf{x}}{\partial \mathbf{x}} = (A + A^T) \mathbf{x}$, που γίνεται $2A \mathbf{x}$ αν ο A είναι συμμετρικός.
- $\frac{\partial A \mathbf{x}}{\partial \mathbf{x}} = A^T$

Διακριτές τυχαίες μεταβλητές

Μια **διακριτή** τυχαία μεταβλητή x μπορεί να πάρει οποιαδήποτε τιμή από ένα πεπερασμένο ή άπειρα αριθμήσιμο σύνολο \mathcal{X} , που ονομάζεται δειγματικός χώρος. Η πιθανότητα του γεγονότος " $x = x \in \mathcal{X}$ " συμβολίζεται ως:

$$P(x = x) \text{ ή απλά } P(x)$$

Η συνάρτηση P ονομάζεται *συνάρτηση μάζας πιθανότητας* και ισχύει:

$$\sum_{x \in \mathcal{X}} P(x) = 1$$

Από κοινού και δεσμευμένες πιθανότητες:

- Κανόνας αθροίσματος: $P(x) = \sum_{y \in \mathcal{Y}} P(x, y)$
- Κανόνας γινομένου: $P(x, y) = P(x|y)P(y)$
- Θεώρημα Bayes: $P(y|x) = \frac{P(x|y)P(y)}{P(x)}$

Συνεχείς τυχαίες μεταβλητές

Μια **συνεχής** τυχαία μεταβλητή x παίρνει τιμές στον άξονα των πραγματικών αριθμών, \mathbb{R} . Η *συνάρτηση κατανομής πιθανότητας* της x ορίζεται ως:

$$F_x(x) \triangleq P(x \leq x)$$

και η *συνάρτηση πυκνότητας πιθανότητας* (ΣΠΠ) ως:

$$p_x(x) \triangleq \frac{dF_x(x)}{dx}$$

Από τα παραπάνω προκύπτουν τα εξής (ο δείκτης x παραλείπεται):

$$F(x) = \int_{-\infty}^x p(x)dx, \quad P(x_1 \leq x \leq x_2) = \int_{x_1}^{x_2} p(x)dx \quad \text{και} \quad \int_{-\infty}^{\infty} p(x)dx = 1$$

Επίσης ισχύουν για την συνάρτηση πυκνότητας πιθανότητας όλοι κανόνες που ισχύουν για τις πιθανότητες διακριτών τυχαίων μεταβλητών. Για παράδειγμα:

$$p(x|y) = \frac{p(x,y)}{p(y)}, \quad p(x) = \int_{-\infty}^{\infty} p(x,y)dy.$$

Μέση τιμή και μεταβλητότητα

- Μέση τιμή: $\mathbb{E}[x] = \int_{-\infty}^{\infty} xp(x)dx$

- Μεταβλητότητα: $\sigma_x^2 = \int_{-\infty}^{\infty} (x - \mathbb{E}[x])^2 p(x)dx$

- Γενικότερα: $\mathbb{E}[f(x)] = \int_{-\infty}^{\infty} f(x)p(x)dx$

- Συμμεταβλητότητα (covariance) τυχαίων μεταβλητών x, y :

$$\text{cov}(x, y) = \mathbb{E}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])]$$

- Συσχέτιση (correlation) τυχαίων μεταβλητών x, y :

$$r_{x,y} = \mathbb{E}[xy] = \text{cov}(x, y) + \mathbb{E}[x]\mathbb{E}[y]$$

- Ένα τυχαίο διάνυσμα $\mathbf{x} = [x_1, \dots, x_l]^T$ είναι ένα διάνυσμα με στοιχεία τυχαίες μεταβλητές. Ορίζεται η $p(\mathbf{x})$ ως **η από κοινού ΣΠΠ**:

$$p(\mathbf{x}) = p(x_1, \dots, x_l)$$

Πίνακες συμμεταβλητότητας και συσχέτισης

- Πίνακας συμμεταβλητότητας ενός τυχαίου διανύσματος $\mathbf{x} \in \mathbb{R}^l$:

$$\text{Cov}(\mathbf{x}) = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \quad \text{ή}$$

$$\text{Cov}(\mathbf{x}) = \begin{bmatrix} \text{cov}(x_1, x_1) & \cdots & \text{cov}(x_1, x_l) \\ \vdots & \ddots & \vdots \\ \text{cov}(x_l, x_1) & \cdots & \text{cov}(x_l, x_l) \end{bmatrix}$$

- Πίνακας συσχέτισης ενός τυχαίου διανύσματος $\mathbf{x} \in \mathbb{R}^l$:

$$R_x = \mathbb{E}[\mathbf{x}\mathbf{x}^T] \quad \text{ή}$$

$$R_x = \begin{bmatrix} \mathbb{E}(x_1, x_1) & \cdots & \mathbb{E}(x_1, x_l) \\ \vdots & \ddots & \vdots \\ \mathbb{E}(x_l, x_1) & \cdots & \mathbb{E}(x_l, x_l) \end{bmatrix} = \text{Cov}(\mathbf{x}) + \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}^T]$$

- Πολύ σημαντικό: Οι πίνακες συμμεταβλητότητας και συσχέτισης είναι θετικά ημιορισμένοι.

Κανονική (Gaussian) κατανομή

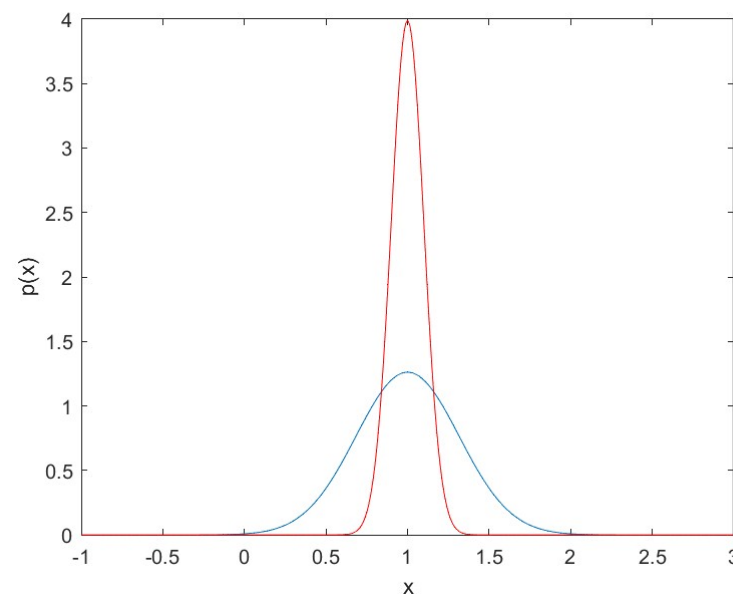
Θα λέμε ότι η τυχαία μεταβλητή x ακολουθεί την **κανονική ή Gaussian κατανομή** με παραμέτρους μ και σ^2 και γράφουμε $x \sim \mathcal{N}(\mu, \sigma^2)$ ή $\mathcal{N}(x|\mu, \sigma^2)$ αν

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Μπορεί ναδειχτεί ότι:

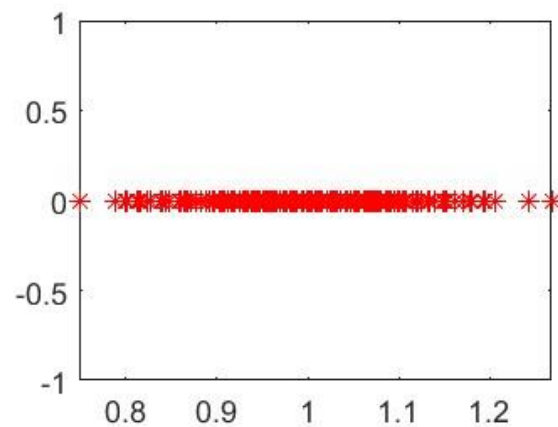
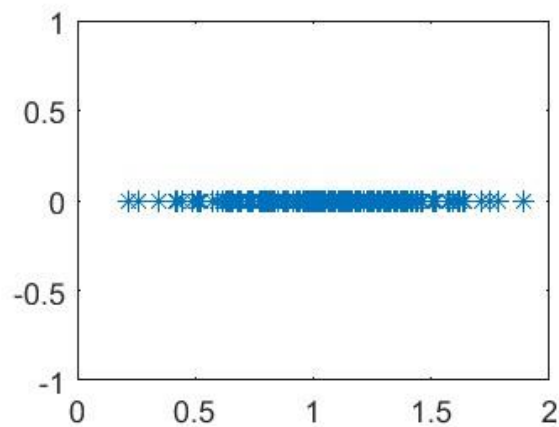
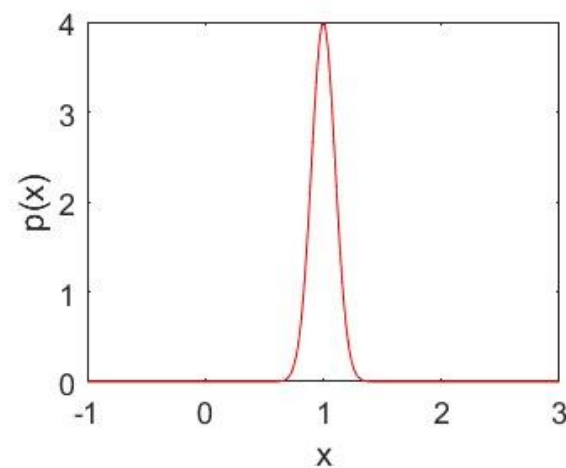
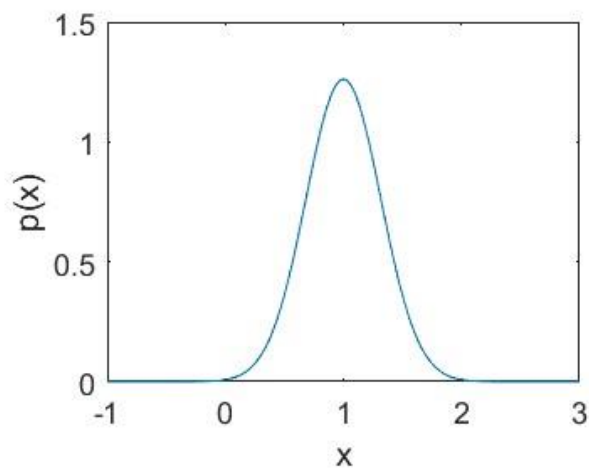
$$\mathbb{E}[x] = \mu \quad \text{και} \quad \sigma_x^2 = \sigma^2$$

Γραφική παράσταση δύο Gaussian ΣΠΠ με $\mu = 1$ και $\sigma^2 = 0.1$ (μπλε) και $\sigma^2 = 0.01$ (κόκκινη) αντίστοιχα.



Κανονική (Gaussian) κατανομή

200 δείγματα από καθεμιά από τις παραπάνω κανονικές κατανομές



Πολυμεταβλητή κανονική κατανομή

Η γενίκευση της κανονικής κατανομής για τυχαία διανύσματα $\mathbf{x} \in \mathbb{R}^l$, οδηγεί στη λεγόμενη **πολυμεταβλητή κανονική (Gaussian) κατανομή**, $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$, με παραμέτρους $\boldsymbol{\mu}$ και Σ , που ορίζεται ως εξής

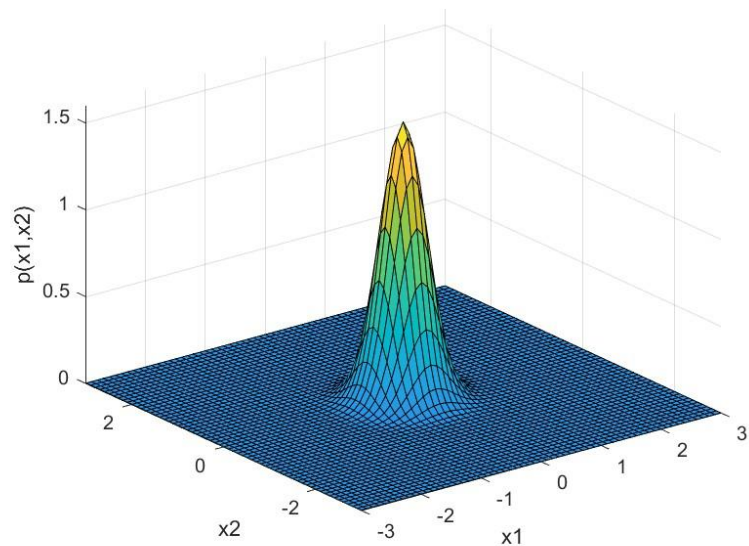
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{l/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})\right)$$

Μπορεί ναδειχτεί ότι:

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} \quad \text{και} \quad \text{Cov}(\mathbf{x}) = \Sigma.$$

Πολυμεταβλητή Gaussian κατανομή με $\boldsymbol{\mu} = \mathbf{0}$ και

$$\Sigma = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}$$

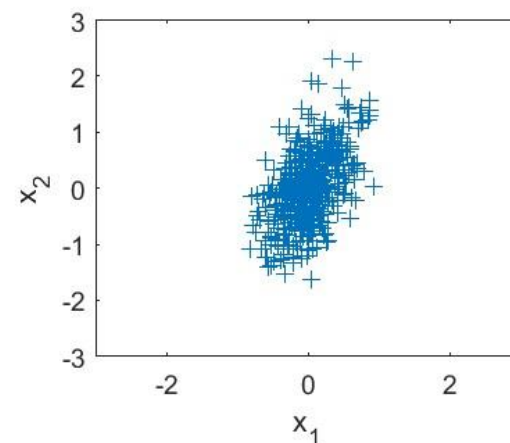
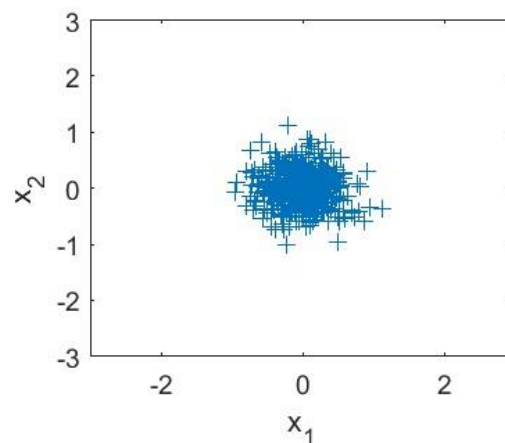
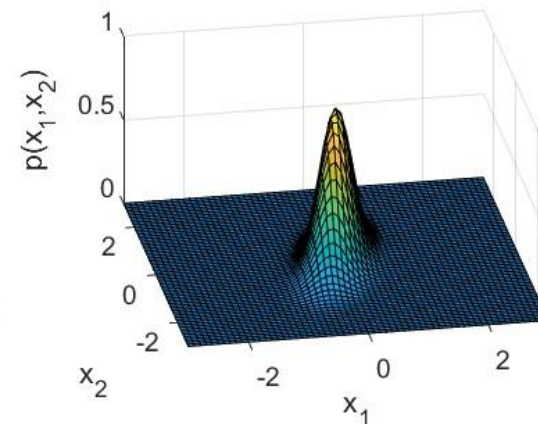
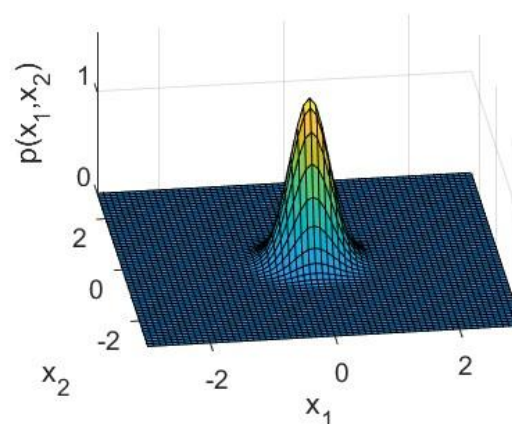


Πολυμεταβλητή κανονική κατανομή

Δύο διμεταβλητές κανονικές κατανομές με $\mu = \mathbf{0}$ και

$$\Sigma_1 = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.4 \end{bmatrix}$$

Στο σχήμα φαίνονται επίσης 500 δείγματα από κάθε κατανομή.

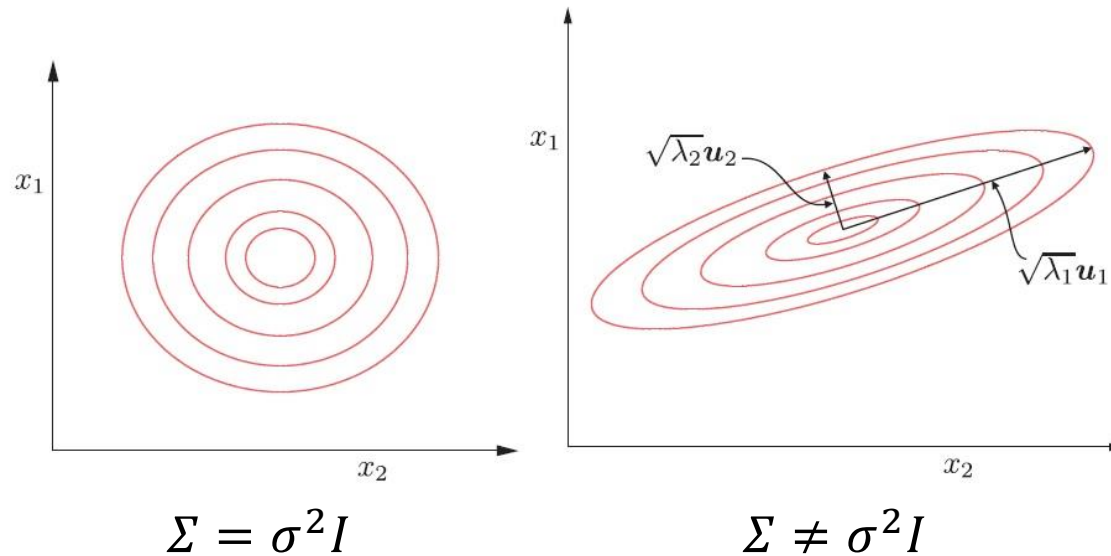


Πολυμεταβλητή κανονική κατανομή

- **Ισοσταθμικές καμπύλες της πολυμεταβλητής κανονικής κατανομής:** μια ισοσταθμική καμπύλη σχηματίζεται από όλα τα σημεία \mathbf{x} τα οποία αντιστοιχούν στην ίδια τιμή της ΣΠΠ, δηλαδή $p(\mathbf{x}) = c$,

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \text{σταθερο} = c$$

- Οι ισοσταθμικές καμπύλες είναι είτε κύκλοι (υπερσφαίρες), είτε ελλείψεις (υπερελλειψοειδή) με κέντρο τη μέση τιμή $\boldsymbol{\mu}$. Τα μήκη των αξόνων καθορίζονται από τις ιδιοτιμές και τα ιδιοδιανύσματα του $\boldsymbol{\Sigma}$.



Πολυμεταβλητή κανονική κανανομή

Αν ο πίνακας συμμεταβλητότητας είναι διαγώνιος, δηλαδή

$$\Sigma = \text{diag}\{\sigma_1^2, \sigma_2^2, \dots, \sigma_l^2\}$$

τότε οι επιμέρους τυχαίες μεταβλητές θα είναι **ασυσχέτιστες**, $\text{cov}(x_i, x_j) = 0$, $i, j = 1, 2, \dots, l$.
Επιπλέον όμως μπορεί εύκολα ναδειχτεί ότι θα ισχύει:

$$p(\mathbf{x}) = \prod_{i=1}^l \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right).$$

Με άλλα λόγια,

$$p(\mathbf{x}) = \prod_{i=1}^l p(x_i)$$

δηλαδή οι τυχαίες μεταβλητές x_1, x_2, \dots, x_l είναι επιπλέον **στατιστικά ανεξάρτητες**.

Δεσμευμένες κανονικές κανανομές

Έστω ότι το l -διάστατο διάνυσμα \mathbf{x} γράφεται ως εξής:

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}, \text{ με } \mathbf{x}_1 \in \mathbb{R}^k \text{ και } \mathbf{x}_2 \in \mathbb{R}^{l-k}$$

και τα $\boldsymbol{\mu}, \Sigma$, διαμερίζονται αντίστοιχα:

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}.$$

Τότε η δεσμευμένη κατανομή του \mathbf{x}_1 δοσμένου ότι $\mathbf{x}_2 = \mathbf{a}$, είναι κανονική $(\mathbf{x}_1 | \mathbf{x}_2 = \mathbf{a}) \sim \mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\Sigma})$ με:

$$\hat{\boldsymbol{\mu}} = \boldsymbol{\mu}_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{a} - \boldsymbol{\mu}_2) \text{ και } \hat{\Sigma} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

Μέθοδος μεγίστης πιθανοφάνειας

- Σε πολλές εφαρμογές μας ενδιαφέρει να **εκτιμήσουμε την κατανομή** από την οποία προέρχεται ένα σύνολο παρατηρήσεων ή μετρήσεων.
- Πιο συγκεκριμένα, έστω ότι μας δίνεται ένα σύνολο N παρατηρήσεων $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ που έχουν προέλθει από μια κατανομή πιθανότητας. Υποθέτουμε ότι η από κοινού ΣΠΠ αυτών των μετρήσεων είναι **γνωστού παραμετρικού συναρτησιακού τύπου** και συμβολίζεται με $p(\mathcal{X}; \theta)$. Η παράμετρος $\theta \in \mathbb{R}^K$ είναι άγνωστη και ο στόχος μας είναι να εκτιμήσουμε την τιμή της.
- Η από κοινού ΣΠΠ $p(\mathcal{X}; \theta)$ ονομάζεται **συνάρτηση πιθανοφάνειας** (likelihood function) του θ ως προς το δοσμένο σετ παρατηρήσεων \mathcal{X} . Σύμφωνα με τη μέθοδο μεγίστης πιθανοφάνειας (maximum likelihood method), η παράμετρος θ εκτιμάται ως εξής:

$$\hat{\theta}_{\text{ML}} := \arg \max_{\theta} p(\mathcal{X}; \theta)$$

Μέθοδος μεγίστης πιθανοφάνειας

- Καθώς η λογαριθμική συνάρτηση $\ln(\cdot)$ είναι γνησίως αύξουσα, μπορεί κανείς να αναζητήσει εναλλακτικά το μέγιστο του λογαρίθμου της συνάρτησης πιθανοφάνειας, δηλαδή:

$$\left. \frac{\partial \ln p(\mathcal{X}; \theta)}{\partial \theta} \right|_{\theta = \hat{\theta}_{\text{ML}}} = 0$$

- Το $\hat{\theta}_{\text{ML}}$ ονομάζεται **εκτιμητής μεγίστης πιθανοφάνειας**, είναι συνάρτηση των x_1, x_2, \dots, x_N και κατά συνέπεια είναι τυχαία μεταβλητή.
- Αν οι παρατηρήσεις x_1, x_2, \dots, x_N είναι ανεξάρτητες και ακολουθούν την ίδια κατανομή (i.i.d.),
 - Ο εκτιμητής μεγίστης πιθανοφάνειας είναι **ασυμπτωτικά αμερόληπτος**. Δηλαδή, αν υποθέσουμε ότι το μοντέλο που επιλέξαμε για την κατανομή των παρατηρήσεων είναι σωστό και η πραγματική παράμετρος του μοντέλου είναι θ_o θα ισχύει:

$$\lim_{N \rightarrow \infty} \mathbb{E}[\hat{\theta}_{\text{ML}}] = \theta_o$$

- Ο εκτιμητής μεγίστης πιθανοφάνειας είναι **ασυμπτωτικά συνεπής**, δηλαδή για κάθε $\epsilon > 0$

$$\lim_{N \rightarrow \infty} \text{Prob} \left\{ \left| \hat{\theta}_{\text{ML}} - \theta_o \right| > \epsilon \right\} = 0$$

Μέθοδος μεγίστης πιθανοφάνειας

Παράδειγμα. Έστω x_1, x_2, \dots, x_N βαθμωτές παρατηρήσεις που έχουν προέλθει από μια κανονική κατανομή με **γνωστή** μεταβλητότητα σ^2 και **άγνωστη** μέση τιμή μ , δηλαδή,

$$p(x_n; \mu) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x_n - \mu)^2}{2\sigma^2}\right)$$

Υποθέτουμε ότι οι μετρήσεις είναι **στατιστικά ανεξάρτητες** και θέλουμε να εκτιμήσουμε από αυτές την άγνωστη παράμετρο της κατανομής μ .

Για τις N στατιστικά ανεξάρτητες παρατηρήσεις ο λογάριθμος της από κοινού συνάρτησης πιθανοφάνειας θα είναι:

$$L(\mu) = \ln \prod_{n=1}^N p(x_n; \mu) = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2$$

Μέθοδος μεγίστης πιθανοφάνειας

Παράδειγμα (συνέχεια). Παραγωγίζοντας ως προς μ παίρνουμε

$$\frac{\partial L(\mu)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu),$$

το οποίο αν εξισωθεί με το 0 δίνει,

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n : \text{εκτιμητής μεγίστης πιθανοφάνειας του } \mu$$

- Παρόμοια διαδικασία μπορεί να εφαρμοστεί και για $\theta = \sigma^2$ με μ γνωστό. Δοκιμάστε!
- Αν οι παρατηρήσεις είναι διανύσματα $x_n \in \mathbb{R}^l$ που ακολουθούν την πολυμεταβλητή κανονική κατανομή με γνωστό Σ και άγνωστο μ , μπορεί ναδειχτεί ότι ο εκτιμητής μεγίστης πιθανοφάνειας του μ είναι:

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n$$

Βιβλιογραφία

- S. Theodoridis, Machine Learning: A Bayesian and Optimization Perspective, 2nd Edition, Academic Press, 2020.
- Σ. Θεοδωρίδης, Διαφάνειες του παραπάνω συγγράμματος (στα αγγλικά).