Machine Learning and Statistics in Genetics and Genomics

II: Linear regression

#### **Christoph Lippert**

Microsoft Research eScience group Research

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Los Angeles, USA

Current topics in computational biology UCLA Winter quarter 2014

Introduction Maximum likelihood estimation Genome-wide association studies

#### Gaussian distribution

Central limit theorem Linear transformations Marginal distributions Common Gaussian calculations

# Outline

# Outline

#### Linear Regression

Introduction Maximum likelihood estimation Genome-wide association studies

#### Gaussian distribution

Central limit theorem Linear transformations Marginal distributions Common Gaussian calculations

▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト の Q @

## Regression

Noise model and likelihood

Given a dataset D = {x<sub>n</sub>, y<sub>n</sub>}<sup>N</sup><sub>n=1</sub>, where x<sub>n</sub> = {x<sub>n,1</sub>,..., x<sub>n,D</sub>} is D dimensional, fit parameters θ of a regressor f with added Gaussian noise:

$$y_n = f(\boldsymbol{x}_n; \boldsymbol{\theta}) + \epsilon_n \quad ext{where} \quad p(\epsilon \mid \sigma^2) = \mathcal{N}\left( \, \epsilon \mid 0, \sigma^2 \, 
ight).$$

Equivalent likelihood formulation:

$$p(\boldsymbol{y} \mid \boldsymbol{X}) = \prod_{n=1}^{N} \mathcal{N} \left( y_n \mid f(\boldsymbol{x}_n; \boldsymbol{\theta}), \sigma^2 \right)$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

► Choose *f* to be linear:

$$p(\boldsymbol{y} \mid \boldsymbol{X}) = \prod_{n=1}^{N} \mathcal{N} \left( y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta} + c, \sigma^2 \right)$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Consider bias free case, c = 0, otherwise include an additional column of ones in each x<sub>n</sub>. Choose f to be linear:

$$p(\boldsymbol{y} \mid \boldsymbol{X}) = \prod_{n=1}^{N} \mathcal{N} \left( y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta} + c, \sigma^2 \right)$$

Consider bias free case, c = 0, otherwise include an additional column of ones in each x<sub>n</sub>.



Equivalent graphical model

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ・ ヨ ・ の Q ()

Taking the logarithm, we obtain

$$\ln p(\boldsymbol{y} \mid \boldsymbol{\theta} \sigma^2) = \sum_{n=1}^{N} \ln \mathcal{N} \left( y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta}, \sigma^2 \right)$$
$$= -\frac{N}{2} \ln 2\pi \sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2}_{\text{Sum of squares}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- The likelihood is maximized when the squared error is minimized.
- Least squares and maximum likelihood are equivalent.

Taking the logarithm, we obtain

$$\ln p(\boldsymbol{y} \mid \boldsymbol{\theta} \sigma^2) = \sum_{n=1}^{N} \ln \mathcal{N} \left( y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta}, \sigma^2 \right)$$
$$= -\frac{N}{2} \ln 2\pi \sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2}_{\text{Sum of squares}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- The likelihood is maximized when the squared error is minimized.
- Least squares and maximum likelihood are equivalent.

Taking the logarithm, we obtain

$$\ln p(\boldsymbol{y} \mid \boldsymbol{\theta} \sigma^2) = \sum_{n=1}^{N} \ln \mathcal{N} \left( y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta}, \sigma^2 \right)$$
$$= -\frac{N}{2} \ln 2\pi \sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2}_{\text{Sum of squares}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- The likelihood is maximized when the squared error is minimized.
- Least squares and maximum likelihood are equivalent.



(C.M. Bishop, Pattern Recognition and Machine Learning)

$$\operatorname{argmin}_{\boldsymbol{\beta}} \frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2$$

• Derivative w.r.t a single weight entry  $\beta_i$ 

$$\frac{d}{d\beta_d} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{d\beta_d} \left[ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$

Set gradient w.r.t.  $\beta$  to zero [vector holding the derivatives  $\forall \beta_d$ ]

 $\nabla_{\boldsymbol{\beta}} \ln p(\boldsymbol{y} \,|\, \boldsymbol{\beta}, \sigma^2) =$ 



• Derivative w.r.t a single weight entry  $\beta_i$ 

$$\frac{d}{d\beta_d} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{d\beta_d} \left[ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$
$$= \left[ -\frac{1}{\sigma^2} \sum_{n=1}^N \boldsymbol{x}_{nd} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) \right]$$

Set gradient w.r.t.  $\beta$  to zero [vector holding the derivatives  $\forall \beta_d$ ]

 $\nabla_{\boldsymbol{\beta}} \ln p(\boldsymbol{y} \,|\, \boldsymbol{\beta}, \sigma^2) =$ 



• Derivative w.r.t a single weight entry  $\beta_i$ 

$$\frac{d}{d\beta_d} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{d\beta_d} \left[ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$
$$= \left[ -\frac{1}{\sigma^2} \sum_{n=1}^N \boldsymbol{x}_{nd} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) \right]$$

Set gradient w.r.t.  $\beta$  to zero [vector holding the derivatives  $\forall \beta_d$ ]

$$\nabla_{\boldsymbol{\beta}} \ln p(\boldsymbol{y} \,|\, \boldsymbol{\beta}, \sigma^2) =$$

► Here, the matrix 
$$X$$
 is defined as  $X = \begin{bmatrix} x_{11} & \dots & x_{1D} \\ \dots & \dots & \dots \\ x_{N1} & \dots & x_{ND} \end{bmatrix}$ 

**Derivative** w.r.t a single weight entry  $\beta_i$ 

$$\frac{d}{d\beta_d} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{d\beta_d} \left[ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$
$$= \left[ -\frac{1}{\sigma^2} \sum_{n=1}^N \boldsymbol{x}_{nd} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) \right]$$

Set gradient w.r.t.  $\beta$  to zero [vector holding the derivatives  $\forall \beta_d$ ]

$$abla_{oldsymbol{eta}} \ln p(oldsymbol{y} \,|\, oldsymbol{eta}, \sigma^2) = rac{1}{\sigma^2} \sum_{n=1}^N oldsymbol{x}_n^{ op}(y_n - oldsymbol{x}_n \cdot oldsymbol{eta}) = oldsymbol{0}$$
 (where  $oldsymbol{0}$  is a vector of 0s)



**Derivative** w.r.t a single weight entry  $\beta_i$ 

$$\frac{d}{d\beta_d} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{d\beta_d} \left[ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$
$$= \left[ -\frac{1}{\sigma^2} \sum_{n=1}^N \boldsymbol{x}_{nd} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) \right]$$

Set gradient w.r.t.  $\beta$  to zero [vector holding the derivatives  $\forall \beta_d$ ]

$$\nabla_{\boldsymbol{\beta}} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sigma^2} \sum_{n=1}^{N} \boldsymbol{x}_n^{\top} (\boldsymbol{y}_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) = \boldsymbol{0} \quad \text{(where } \boldsymbol{0} \text{ is a vector of 0s)}$$
$$\frac{1}{\sigma^2} \boldsymbol{X}^{\top} (\boldsymbol{y} - \boldsymbol{X} \cdot \boldsymbol{\beta}) = \boldsymbol{0}$$
$$\implies \boldsymbol{\beta}_{\mathrm{ML}} =$$

• Here, the matrix 
$$X$$
 is defined as  $X = \begin{bmatrix} x_{11} & \dots & x_{1D} \\ \dots & \dots & \dots \\ x_{N1} & \dots & x_{ND} \end{bmatrix}$ 

**Derivative** w.r.t a single weight entry  $\beta_i$ 

$$\frac{d}{d\beta_d} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{d\beta_d} \left[ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$
$$= \left[ -\frac{1}{\sigma^2} \sum_{n=1}^N \boldsymbol{x}_{nd} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) \right]$$

Set gradient w.r.t.  $\beta$  to zero [vector holding the derivatives  $\forall \beta_d$ ]

$$\nabla_{\boldsymbol{\beta}} \ln p(\boldsymbol{y} \,|\, \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sigma^2} \sum_{n=1}^{N} \boldsymbol{x}_n^\top (\boldsymbol{y}_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) = \boldsymbol{0} \quad \text{(where } \boldsymbol{0} \text{ is a vector of 0s)}$$
$$\frac{1}{\sigma^2} \boldsymbol{X}^\top (\boldsymbol{y} - \boldsymbol{X} \cdot \boldsymbol{\beta}) = \boldsymbol{0}$$
$$\Longrightarrow \boldsymbol{\beta}_{\mathrm{ML}} = \underbrace{(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top}_{\mathsf{Pseudo inverse of } \boldsymbol{X}} \boldsymbol{y}$$

• Here, the matrix  $\boldsymbol{X}$  is defined as  $\boldsymbol{X} = \begin{bmatrix} x_{11} & \dots & x1D \\ \dots & \dots & \dots \\ x_{N1} & \dots & x_{ND} \end{bmatrix}$ 

Application: Genome-wide association studies

### Given:

- Genetics for multiple individuals
  - e.g.: Single nucleotide polymorphisms (SNPs), microsatelite markers, ...
- Quantitative phenotype for the same individuals
  - e.g.: height, gene-expression, ...
- Try to find genetic markers, that explain the variance in the phenotype.
- Use linear regression!



イロト 不得 トイヨト イヨト

э

Application: Genome-wide association studies

#### Given:

- Genetics for multiple individuals
  - e.g.: Single nucleotide polymorphisms (SNPs), microsatelite markers, ...
- Quantitative phenotype for the same individuals
  - e.g.: height, gene-expression, ...
- Try to find genetic markers, that explain the variance in the phenotype.
- Use linear regression!



Application: Genome-wide association studies

#### Given:

- Genetics for multiple individuals
  - e.g.: Single nucleotide polymorphisms (SNPs), microsatelite markers, ...
- Quantitative phenotype for the same individuals
  - e.g.: height, gene-expression, ...
- Try to find genetic markers, that explain the variance in the phenotype.
- Use linear regression!



Application: Genome-wide association studies

#### Given:

- Genetics for multiple individuals
  - e.g.: Single nucleotide polymorphisms (SNPs), microsatelite markers, ...
- Quantitative phenotype for the same individuals
  - e.g.: height, gene-expression, ...

## Goal:

- Try to find genetic markers, that explain the variance in the phenotype.
- ► Use linear regression!



Application: Genome-wide association studies

#### Given:

- Genetics for multiple individuals
  - e.g.: Single nucleotide polymorphisms (SNPs), microsatelite markers, ...
- Quantitative phenotype for the same individuals
  - e.g.: height, gene-expression, ...

### Goal:

- Try to find genetic markers, that explain the variance in the phenotype.
- Use linear regression!



Application: Genome-wide association studies

#### Given:

- Genetics for multiple individuals
  - e.g.: Single nucleotide polymorphisms (SNPs), microsatelite markers, ...
- Quantitative phenotype for the same individuals
  - e.g.: height, gene-expression, ...

### Goal:

- Try to find genetic markers, that explain the variance in the phenotype.
- Use linear regression!



- Genotype x denotes the genetic state of an individual.
  - Denoted by  $x_{n}$ : for individual n.
- Phenotype denotes the state of a trait of an individual.
  - Denoted by  $y_n$  for individual n.
- ► A **Locus** is a position or limited region in the genome.
  - Denoted by  $x_s$  for locus (or SNP) s.
- An allele is the genetic state of a locus.



・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

-

- Genotype x denotes the genetic state of an individual.
  - Denoted by  $x_{n:}$  for individual n.
- Phenotype denotes the state of a trait of an individual.
  - Denoted by  $y_n$  for individual n.
- A Locus is a position or limited region in the genome.
  - Denoted by  $oldsymbol{x}_s$  for locus (or SNP) s.
- An allele is the genetic state of a locus.



- Genotype x denotes the genetic state of an individual.
  - Denoted by  $x_{n:}$  for individual n.
- Phenotype denotes the state of a trait of an individual.
  - Denoted by  $y_n$  for individual n.
- A Locus is a position or limited region in the genome.
  - Denoted by  $x_s$  for locus (or SNP) s.

An allele is the genetic state of a locus.





- Genotype x denotes the genetic state of an individual.
  - Denoted by  $x_{n}$ : for individual n.
- Phenotype denotes the state of a trait of an individual.
  - Denoted by  $y_n$  for individual n.
- A Locus is a position or limited region in the genome.
  - Denoted by  $x_s$  for locus (or SNP) s.
- An **allele** is the genetic state of a locus.





# Genetics 101

More definitions

- An organism/cell is haploid if it only has one chromosome set or identical chromosome sets.
  - e.g. A. thaliana, sperm cells or inbred lab strains
- An organism/cell is diploid if it has two separately inherited homologous chromosomes.
  - ▶ e.g. human
- An organism/cell is **polyploid** if it has more than two homologous chromosomes.
  - e.g. *sugar cane* is hexaploid.

	and the second se	(haranta)			Charger and Charge	Charged C		
)	2	dinet and	and the second s	(), III		Control Control		
dans,		1		1000	0000	50		
88	58		9.6	56		and a	9	
19	28		23	=		x	v	

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

# Genetics 101

More definitions

- An organism/cell is haploid if it only has one chromosome set or identical chromosome sets.
  - e.g. A. thaliana, sperm cells or inbred lab strains
- An organism/cell is diploid if it has two separately inherited homologous chromosomes.
  - ▶ e.g. human
- An organism/cell is **polyploid** if it has more than two homologous chromosomes.
  - e.g. *sugar cane* is hexaploid.

	OR WILDOWSKY	Denner an			Charger and Charge	(Depend)		
)(	2	disease of the	and a second	(), III		Control Control		
date.	1	1		100	0.00	1		
88	88		9.6	56		4	9	
19	28		21	22		x	¥	

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

# Genetics 101

More definitions

- An organism/cell is haploid if it only has one chromosome set or identical chromosome sets.
  - e.g. A. thaliana, sperm cells or inbred lab strains
- An organism/cell is diploid if it has two separately inherited homologous chromosomes.
  - ▶ e.g. human
- An organism/cell is **polyploid** if it has more than two homologous chromosomes.
  - e.g. *sugar cane* is hexaploid.

	OR WILDOWSKY	Chinasonani .			Charger of the second	(Depend)	
)(	2	diverging a	and a second	Grade .		Control of	
and a				1000	0000	500	
88	88		9.6	56			9
19	28		23	=		x	v

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙

- Haplotype denotes an individual's state of a single set of chromosomes (paternal or maternal).
- A locus s is heterozygous if it differs between paternal and maternal haplotypes.
  - heterozygous allele usually encoded as
     x<sub>s</sub> = 1
- A locus is homozygous if it matches between paternal and maternal haplotypes.
  - ▶ homozygous major allele usually encoded as x<sub>s</sub> = 0
  - ▶ homozygous minor allele usually encoded as x<sub>s</sub> = 2
- x<sub>ns</sub> counts the number of minor alleles (mutations) at SNP s.

							:	SNI 1	P						
Tree 1	A	с	G	т	G	т	С	G	G	т	С	т	т	A	Maternal chrom.
	A	С	G	Т	G	т	С	A	G	т	C	т	т	A	Paternal chrom.
Tree 2	A	с	G	т	G	т	с	G	G	т	С	т	т	A	Maternal chrom.
1100 2	A	С	G	Т	G	т	С	G	G	т	C	т	т	A	Paternal chrom.
Teres 2	A	С	G	т	G	т	С	A	G	т	С	т	т	A	Maternal chrom.
1166.5	A	С	G	т	G	т	С	A	G	т	С	т	т	A	Paternal chrom.

- Haplotype denotes an individual's state of a single set of chromosomes (paternal or maternal).
- A locus s is heterozygous if it differs between paternal and maternal haplotypes.
  - $\blacktriangleright$  heterozygous allele usually encoded as  $x_s=1$
- A locus is homozygous if it matches between paternal and maternal haplotypes.
  - ▶ homozygous major allele usually encoded as x<sub>s</sub> = 0
  - ▶ homozygous minor allele usually encoded as x<sub>s</sub> = 2
- x<sub>ns</sub> counts the number of minor alleles (mutations) at SNP s.

							:	SNI	D						
Tree 1	A	С	G	т	G	т	С	G	G	т	С	т	т	А	Maternal chrom.
	A	С	G	т	G	т	С	A	G	т	C	т	т	A	Paternal chrom.
Tree 2	A	С	G	т	G	т	с	G	G	т	С	т	т	A	Maternal chrom.
	A	С	G	Т	G	т	С	G	G	т	C	т	т	A	Paternal chrom.
Trop 2	A	С	G	т	G	т	С	A	G	т	С	т	т	A	Maternal chrom.
1100 0	A	С	G	т	G	т	С	A	G	т	С	т	т	A	Paternal chrom.

- Haplotype denotes an individual's state of a single set of chromosomes (paternal or maternal).
- A locus s is heterozygous if it differs between paternal and maternal haplotypes.
  - $\blacktriangleright$  heterozygous allele usually encoded as  $x_s=1$
- A locus is homozygous if it matches between paternal and maternal haplotypes.
  - ▶ homozygous major allele usually encoded as x<sub>s</sub> = 0
  - ▶ homozygous minor allele usually encoded as x<sub>s</sub> = 2

 x<sub>ns</sub> counts the number of minor alleles (mutations) at SNP s.

							:	SNI 4	Þ						
Tree 1	A	С	G	т	G	т	С	G	G	т	С	т	т	А	Maternal chrom.
	A	С	G	т	G	т	С	A	G	т	C	т	т	A	Paternal chrom.
Tree 2	A	с	G	т	G	т	с	G	G	т	С	т	т	A	Maternal chrom.
	A	С	G	т	G	т	С	G	G	т	C	т	т	A	Paternal chrom.
T 2	A	С	G	т	G	т	С	A	G	т	С	т	т	A	Maternal chrom.
	A	С	G	т	G	т	С	A	G	т	С	т	т	A	Paternal chrom.

- Haplotype denotes an individual's state of a single set of chromosomes (paternal or maternal).
- A locus s is heterozygous if it differs between paternal and maternal haplotypes.
  - $\blacktriangleright$  heterozygous allele usually encoded as  $x_s=1$
- A locus is homozygous if it matches between paternal and maternal haplotypes.
  - ▶ homozygous major allele usually encoded as x<sub>s</sub> = 0
  - ▶ homozygous minor allele usually encoded as x<sub>s</sub> = 2
- x<sub>ns</sub> counts the number of minor alleles (mutations) at SNP s.

							:	SNI	D						
Tree 1	A	С	G	т	G	т	С	G	G	т	С	т	т	А	Maternal chrom.
	A	С	G	т	G	т	С	A	G	т	C	т	т	A	Paternal chrom.
Tree 2	A	С	G	т	G	т	с	G	G	т	С	т	т	A	Maternal chrom.
	A	С	G	Т	G	т	С	G	G	т	C	т	т	A	Paternal chrom.
Trop 2	A	С	G	т	G	т	С	A	G	т	С	т	т	A	Maternal chrom.
1100 0	A	С	G	т	G	т	С	A	G	т	С	т	т	A	Paternal chrom.

Example: Genome-wide association study

- Model the phenotype of individual n as a linear function of the SNP x<sub>ns</sub>
- x<sub>ns</sub> ∈ 0, 1, 2 counts the number of mutations that individual n has at the position s.





(C.M. Bishop, Pattern Recognition and Machine Learning)

# Genome-wide association study

Flowering time in A. thaliana

#### Arabidopsis thaliana

- Grows under rough conditions
- Distributed over the whole globe (it "travels" far)
- Genetics model organism ("lab rat" among plants)
- Mostly self-fertilizing (mostly homozygous)
- Model for flowering





◆□▶ ◆圖▶ ◆臣▶ ◆臣▶ ─ 臣
Flowering time in A. thaliana

#### Arabidopsis thaliana

### Grows under rough conditions

- Distributed over the whole globe (it "travels" far)
- Genetics model organism ("lab rat" among plants)
- Mostly self-fertilizing (mostly homozygous)
- Model for flowering





Flowering time in A. thaliana

### Arabidopsis thaliana

- Grows under rough conditions
- Distributed over the whole globe (it "travels" far)
- Genetics model organism ("lab rat" among plants)
- Mostly self-fertilizing (mostly homozygous)
- Model for flowering





Flowering time in A. thaliana

#### Arabidopsis thaliana

- Grows under rough conditions
- Distributed over the whole globe (it "travels" far)
- Genetics model organism ("lab rat" among plants)
- Mostly self-fertilizing (mostly homozygous)
- Model for flowering





Flowering time in A. thaliana

### Arabidopsis thaliana

- Grows under rough conditions
- Distributed over the whole globe (it "travels" far)
- Genetics model organism ("lab rat" among plants)
- Mostly self-fertilizing (mostly homozygous)
- Model for flowering





Flowering time in A. thaliana

### Arabidopsis thaliana

- Grows under rough conditions
- Distributed over the whole globe (it "travels" far)
- Genetics model organism ("lab rat" among plants)
- Mostly self-fertilizing (mostly homozygous)
- Model for flowering





Flowering time in A. thaliana

- GWAS data [Atwell et al. Nat. Gen. 2010]
  - y: Phenotype:
     "Flowering time at 10° Celsius (days)"
  - X: Genotype: 214,553 SNPs, all homozygous
  - 194 samples
  - show demo.





Flowering time in A. thaliana

- GWAS data [Atwell et al. Nat. Gen. 2010]
  - y: Phenotype:
     "Flowering time at 10° Celsius (days)"
  - X: Genotype: 214,553 SNPs, all homozygous
  - 194 samples
  - show demo.





Flowering time in A. thaliana

- GWAS data [Atwell et al. Nat. Gen. 2010]
  - ▶ y: Phenotype: "Flowering time at 10° Celsius (days)"
  - X: Genotype: 214,553 SNPs, all homozygous
  - 194 samples
  - show demo.





Flowering time in A. thaliana

- GWAS data [Atwell et al. Nat. Gen. 2010]
  - ▶ y: Phenotype: "Flowering time at 10° Celsius (days)"
  - X: Genotype: 214,553 SNPs, all homozygous
  - 194 samples
  - show demo.





# Outline

#### Linear Regression Introduction Maximum likelihood estimation Genome-wide association studi

#### Gaussian distribution

Central limit theorem Linear transformations Marginal distributions Common Gaussian calculations

▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト の Q @

Univariate vs. multivariate

▶ The joint distribution of N samples  $x_n$  from independent univariate normal distributions with means  $\mu_n$  and variances  $\sigma_n^2 \forall n \in [1, ..., N]$  is given as

$$\prod_{n=1}^{N} \mathcal{N}\left(x_n \mid \mu_n, \sigma_n^2\right)$$

 $\blacktriangleright$  Special case: all  $x_n$  are also *identically* distributed with mean  $\mu$  and variance  $\sigma^2$ 

$$p(x_1, \dots, x_N) = \prod_{n=1}^{N} \mathcal{N}(x_n \mid \mu, \sigma^2) = \mathcal{N}(\boldsymbol{x} \mid \mu \cdot \boldsymbol{1}, \sigma^2 \cdot \boldsymbol{I})$$

Univariate vs. multivariate

▶ The joint distribution of N samples  $x_n$  from independent univariate normal distributions with means  $\mu_n$  and variances  $\sigma_n^2 \forall n \in [1, ..., N]$  is given as

$$\prod_{n=1}^{N} \mathcal{N}\left(x_{n} \mid \mu_{n}, \sigma_{n}^{2}\right) = \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma_{n}^{2}}} \exp\left[-\frac{1}{2\sigma_{n}^{2}}(x_{n}-\mu_{n})^{2}\right]$$

- Special case: all  $x_n$  are also *identically* distributed with mean  $\mu$  and variance  $\sigma^2$ 

$$p(x_1, \dots, x_N) = \prod_{n=1}^N \mathcal{N}(x_n \mid \mu, \sigma^2) = \mathcal{N}(\boldsymbol{x} \mid \mu \cdot \boldsymbol{1}, \sigma^2 \cdot \boldsymbol{I})$$

Univariate vs. multivariate

▶ The joint distribution of N samples  $x_n$  from independent univariate normal distributions with means  $\mu_n$  and variances  $\sigma_n^2 \forall n \in [1, ..., N]$  is given as

$$\begin{split} \prod_{n=1}^{N} \mathcal{N}\left(x_{n} \mid \mu_{n}, \sigma_{n}^{2}\right) &= \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma_{n}^{2}}} \exp\left[-\frac{1}{2\sigma_{n}^{2}}(x_{n} - \mu_{n})^{2}\right] \\ &= (2\pi)^{-\frac{N}{2}} \left(\prod_{n=1}^{N} \sigma_{n}^{2}\right)^{-\frac{1}{2}} \exp\left[-\frac{1}{2}\sum_{n=1}^{N} \sigma_{n}^{-2}(x_{n} - \mu_{n})^{2}\right] \end{split}$$

 $\blacktriangleright$  Special case: all  $x_n$  are also *identically* distributed with mean  $\mu$  and variance  $\sigma^2$ 

$$p(x_1,\ldots,x_N) = \prod_{n=1}^{N} \mathcal{N}(x_n \mid \mu, \sigma^2) = \mathcal{N}(\boldsymbol{x} \mid \mu \cdot \boldsymbol{1}, \sigma^2 \cdot \boldsymbol{I})$$

Univariate vs. multivariate

▶ The joint distribution of N samples  $x_n$  from independent univariate normal distributions with means  $\mu_n$  and variances  $\sigma_n^2 \forall n \in [1, ..., N]$  is given as

$$\begin{split} \prod_{n=1}^{N} \mathcal{N}\left(x_{n} \mid \mu_{n}, \sigma_{n}^{2}\right) &= \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma_{n}^{2}}} \exp\left[-\frac{1}{2\sigma_{n}^{2}}(x_{n} - \mu_{n})^{2}\right] \\ &= (2\pi)^{-\frac{N}{2}} \left(\prod_{n=1}^{N} \sigma_{n}^{2}\right)^{-\frac{1}{2}} \exp\left[-\frac{1}{2}\sum_{n=1}^{N} \sigma_{n}^{-2}(x_{n} - \mu_{n})^{2}\right] \\ &= (2\pi)^{-\frac{N}{2}} \mid \boldsymbol{\Sigma} \mid^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right), \end{split}$$
  
where we introduced  $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{1}^{2} & \mathbf{0} \\ \vdots \\ \mathbf{0} & \sigma_{N}^{2} \end{bmatrix}, \, \boldsymbol{x} = \begin{bmatrix} x_{1} \\ \vdots \\ x_{N} \end{bmatrix} \text{ and } \boldsymbol{\mu} = \begin{bmatrix} \mu_{1} \\ \vdots \\ \mu_{N} \end{bmatrix}. \end{split}$ 

ullet Special case: all  $x_n$  are also *identically* distributed with mean  $\mu$  and variance  $\sigma^2$ 

$$p(x_1, \dots, x_N) = \prod_{n=1}^{N} \mathcal{N}(x_n \mid \mu, \sigma^2) = \mathcal{N}(\boldsymbol{x} \mid \mu \cdot \mathbf{1}, \sigma^2 \cdot \boldsymbol{I})$$

Univariate vs. multivariate

▶ The joint distribution of N samples  $x_n$  from independent univariate normal distributions with means  $\mu_n$  and variances  $\sigma_n^2 \forall n \in [1, ..., N]$  is given as

$$\begin{split} \prod_{n=1}^{N} \mathcal{N}\left(x_{n} \mid \mu_{n}, \sigma_{n}^{2}\right) &= \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma_{n}^{2}}} \exp\left[-\frac{1}{2\sigma_{n}^{2}}(x_{n}-\mu_{n})^{2}\right] \\ &= (2\pi)^{-\frac{N}{2}} \left(\prod_{n=1}^{N} \sigma_{n}^{2}\right)^{-\frac{1}{2}} \exp\left[-\frac{1}{2}\sum_{n=1}^{N} \sigma_{n}^{-2}(x_{n}-\mu_{n})^{2}\right] \\ &= (2\pi)^{-\frac{N}{2}} \left|\mathcal{L}\right|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{\top} \mathcal{L}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right), \end{split}$$
  
where we introduced  $\mathcal{L} = \begin{bmatrix} \sigma_{1}^{2} & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \sigma_{N}^{2} \end{bmatrix}, \ \boldsymbol{x} = \begin{bmatrix} x_{1} \\ \vdots \\ x_{N} \end{bmatrix} \text{ and } \boldsymbol{\mu} = \begin{bmatrix} \mu_{1} \\ \vdots \\ \mu_{N} \end{bmatrix}. \end{split}$ 

Special case: all  $x_n$  are also *identically* distributed with mean  $\mu$  and variance  $\sigma^2$ 

$$p(x_1,\ldots,x_N) = \prod_{n=1}^N \mathcal{N}(x_n \mid \mu, \sigma^2) = \mathcal{N}(\boldsymbol{x} \mid \mu \cdot \boldsymbol{1}, \sigma^2 \cdot \boldsymbol{I})$$

Univariate vs. multivariate

▶ The joint distribution of N samples  $x_n$  from independent univariate normal distributions with means  $\mu_n$  and variances  $\sigma_n^2 \forall n \in [1, ..., N]$  is given as

$$\begin{split} \prod_{n=1}^{N} \mathcal{N}\left(x_{n} \mid \mu_{n}, \sigma_{n}^{2}\right) &= \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma_{n}^{2}}} \exp\left[-\frac{1}{2\sigma_{n}^{2}}(x_{n} - \mu_{n})^{2}\right] \\ &= (2\pi)^{-\frac{N}{2}} \left(\prod_{n=1}^{N} \sigma_{n}^{2}\right)^{-\frac{1}{2}} \exp\left[-\frac{1}{2}\sum_{n=1}^{N} \sigma_{n}^{-2}(x_{n} - \mu_{n})^{2}\right] \\ &= (2\pi)^{-\frac{N}{2}} \left|\mathcal{L}\right|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^{\top} \mathcal{L}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right), \end{split}$$
  
where we introduced  $\mathcal{L} = \begin{bmatrix} \sigma_{1}^{2} & \mathbf{0} \\ 0 & \sigma_{N}^{2} \end{bmatrix}, \, \boldsymbol{x} = \begin{bmatrix} x_{1} \\ \vdots \\ x_{N} \end{bmatrix} \text{ and } \boldsymbol{\mu} = \begin{bmatrix} \mu_{1} \\ \vdots \\ \mu_{N} \end{bmatrix}. \end{split}$ 

Special case: all  $x_n$  are also *identically* distributed with mean  $\mu$  and variance  $\sigma^2$ 

$$p(x_1,\ldots,x_N) = \prod_{n=1}^{N} \mathcal{N}(x_n \mid \mu, \sigma^2) = \mathcal{N}(\boldsymbol{x} \mid \mu \cdot \boldsymbol{1}, \sigma^2 \cdot \boldsymbol{I})$$

► From now on we'll mostly argue about multivariate normal distributions.

- ► Throw a dice N times.
  - Distribution of the sum?
- Independent samples from a uniform distribution on the interval (0,1).
  - Distribution of the mean?

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- ► Throw a dice N times.
  - Distribution of the sum?
- Independent samples from a uniform distribution on the interval (0,1).
  - Distribution of the mean?



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- ► Throw a dice N times.
  - Distribution of the sum?
- Independent samples from a uniform distribution on the interval (0,1).
  - Distribution of the mean?



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

- ► Throw a dice N times.
  - Distribution of the sum?
- Independent samples from a uniform distribution on the interval (0,1).
  - Distribution of the mean?

N = 2

0

0.5



and Machine Learning)

0.5

N = 1

2

- $x_n \sim P(x)$ , unknown distribution with mean  $\mu$  and covariance  $\Sigma$ .
- Given N independent, identically distributed (i.i.d.) samples from P(x):  $x_1, x_2, \ldots, x_N$ , let  $\bar{x}^{(N)}$  be the sample mean

$$\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right) = \sqrt{N}\left(\frac{1}{N}\left(\sum_{n=1}^{N} \boldsymbol{x}_{n}\right) - \boldsymbol{\mu}\right)$$
$$\rightarrow \mathcal{V}\left(\sqrt{N}\left(\boldsymbol{x}^{(D)} - \boldsymbol{\mu}\right)\right) \stackrel{\rightarrow 0(n,s)}{\rightarrow}\mathcal{M}(0,S)$$

- Only mean and covariance retained when averaging i.i.d. variables.
- Distribution becomes Gaussian
- Gaussian is a "limit distribution" Implication: Once something is Gaussian it usually stays Gaussian under many operations

- $x_n \sim P(x)$ , unknown distribution with mean  $\mu$  and covariance  $\Sigma$ .
- ▶ Given N independent, identically distributed (i.i.d.) samples from P(x): x<sub>1</sub>, x<sub>2</sub>,..., x<sub>N</sub>, let x̄<sup>(N)</sup> be the sample mean

$$\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right) = \sqrt{N} \underbrace{\left(\frac{1}{N}\left(\sum_{n=1}^{N} \boldsymbol{x}_{n}\right) - \boldsymbol{\mu}\right)}_{\rightarrow 0(a.s.)}$$
$$\Rightarrow P\left(\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right)\right) \xrightarrow{d} \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$$

- Only mean and covariance retained when averaging i.i.d. variables.
- Distribution becomes Gaussian
- Gaussian is a "limit distribution" Implication: Once something is Gaussian it usually stays Gaussian under many operations

- $\boldsymbol{x_n} \sim P(\boldsymbol{x})$ , unknown distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ .
- ▶ Given N independent, identically distributed (i.i.d.) samples from P(x): x<sub>1</sub>, x<sub>2</sub>,..., x<sub>N</sub>, let x̄<sup>(N)</sup> be the sample mean

$$egin{aligned} \sqrt{N}\left(ar{m{x}}^{(N)}-m{\mu}
ight) &= \sqrt{N}\left(\underbrace{rac{1}{N}\left(\sum\limits_{n=1}^{N}m{x}_{n}
ight)-m{\mu}
ight)}_{ o 0(a.s.)} \ &\Rightarrow P\left(\sqrt{N}\left(ar{m{x}}^{(N)}-m{\mu}
ight)
ight) \stackrel{d}{ o}\mathcal{N}(\mathbf{0},\Sigma) \end{aligned}$$

- Only mean and covariance retained when averaging i.i.d. variables.
- Distribution becomes Gaussian
- Gaussian is a "limit distribution" Implication: Once something is Gaussian it usually stays Gaussian under many operations

- $\boldsymbol{x_n} \sim P(\boldsymbol{x})$ , unknown distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ .
- ▶ Given N independent, identically distributed (i.i.d.) samples from P(x): x<sub>1</sub>, x<sub>2</sub>,..., x<sub>N</sub>, let x̄<sup>(N)</sup> be the sample mean

$$\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right) = \sqrt{N}\underbrace{\left(\frac{1}{N}\left(\sum_{n=1}^{N} \boldsymbol{x}_{n}\right) - \boldsymbol{\mu}\right)}_{\rightarrow 0(a.s.)}$$
$$\Rightarrow P\left(\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right)\right) \xrightarrow{d} \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$$

- Only mean and covariance retained when averaging i.i.d. variables.
- Distribution becomes Gaussian
- Gaussian is a "limit distribution" Implication: Once something is Gaussian it usually stays Gaussian under many operations

- $\boldsymbol{x_n} \sim P(\boldsymbol{x})$ , unknown distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ .
- ▶ Given N independent, identically distributed (i.i.d.) samples from P(x): x<sub>1</sub>, x<sub>2</sub>,..., x<sub>N</sub>, let x̄<sup>(N)</sup> be the sample mean

$$\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right) = \sqrt{N}\underbrace{\left(\frac{1}{N}\left(\sum_{n=1}^{N} \boldsymbol{x}_{n}\right) - \boldsymbol{\mu}\right)}_{\rightarrow 0(a.s.)}$$
$$\Rightarrow P\left(\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right)\right) \xrightarrow{d} \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$$

- Only mean and covariance retained when averaging i.i.d. variables.
- Distribution becomes Gaussian
- Gaussian is a "limit distribution" Implication: Once something is Gaussian it usually stays Gaussian under many operations

- $\boldsymbol{x_n} \sim P(\boldsymbol{x})$ , unknown distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ .
- ▶ Given N independent, identically distributed (i.i.d.) samples from P(x): x<sub>1</sub>, x<sub>2</sub>,..., x<sub>N</sub>, let x̄<sup>(N)</sup> be the sample mean

$$\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right) = \sqrt{N}\underbrace{\left(\frac{1}{N}\left(\sum_{n=1}^{N} \boldsymbol{x}_{n}\right) - \boldsymbol{\mu}\right)}_{\rightarrow 0(a.s.)}$$
$$\Rightarrow P\left(\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right)\right) \xrightarrow{d} \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$$

- Only mean and covariance retained when averaging i.i.d. variables.
- Distribution becomes Gaussian
- Gaussian is a "limit distribution" Implication: Once something is Gaussian it usually stays Gaussian under many operations

- $\boldsymbol{x_n} \sim P(\boldsymbol{x})$ , unknown distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ .
- ▶ Given N independent, identically distributed (i.i.d.) samples from P(x): x<sub>1</sub>, x<sub>2</sub>,..., x<sub>N</sub>, let x̄<sup>(N)</sup> be the sample mean

$$\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right) = \sqrt{N}\underbrace{\left(\frac{1}{N}\left(\sum_{n=1}^{N} \boldsymbol{x}_{n}\right) - \boldsymbol{\mu}\right)}_{\rightarrow 0(a.s.)}$$
$$\Rightarrow P\left(\sqrt{N}\left(\bar{\boldsymbol{x}}^{(N)} - \boldsymbol{\mu}\right)\right) \xrightarrow{d} \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$$

- Only mean and covariance retained when averaging i.i.d. variables.
- Distribution becomes Gaussian
- Gaussian is a "limit distribution" Implication: Once something is Gaussian it usually stays Gaussian under many operations

#### Linear transformation

For any random variable x with mean  $\mathbb{E}\left[x
ight]$  and covariance  $\mathbb{C}\left[x
ight]$ 

 $\blacktriangleright \ \mathbb{E}\left[ \boldsymbol{A}\boldsymbol{x} + \boldsymbol{b} \right] = \boldsymbol{A}\mathbb{E}\left[ \boldsymbol{x} \right] + \boldsymbol{b}; \quad (\text{linearity of expectation})$ 

$$\mathbb{E}\left[\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}\right] = \sum_{x \in \mathcal{X}} (\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}) p(x) = \boldsymbol{A} \underbrace{\left(\sum_{x \in \mathcal{X}} \boldsymbol{x} p(x)\right)}_{\mathbb{E}[\boldsymbol{x}]} + \boldsymbol{b}$$

$$\begin{array}{l} \blacktriangleright \quad \mathbb{C}\left[\boldsymbol{A}\boldsymbol{x}\right] = \boldsymbol{A}\mathbb{C}\left[\boldsymbol{x}\right]\boldsymbol{A}^{\top} \\ \mathbb{C}\left[\boldsymbol{A}\boldsymbol{x}\right] = \mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\boldsymbol{x}^{\top}\boldsymbol{A}^{\top}\right] - \mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\right]\mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\right]^{\top} \\ = \boldsymbol{A}\left(\mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{\top}] - \mathbb{E}\left[\boldsymbol{x}\right]\mathbb{E}\left[\boldsymbol{x}\right]^{\top}\right)\boldsymbol{A}^{\top} \\ = \boldsymbol{A}\mathbb{C}\left[\boldsymbol{x}\right]\boldsymbol{A}^{\top} \end{array}$$

#### Linear transformation

For any random variable x with mean  $\mathbb{E}\left[x
ight]$  and covariance  $\mathbb{C}\left[x
ight]$ 

 $\blacktriangleright \ \mathbb{E}\left[ \boldsymbol{A}\boldsymbol{x} + \boldsymbol{b} \right] = \boldsymbol{A}\mathbb{E}\left[ \boldsymbol{x} \right] + \boldsymbol{b}; \quad (\text{linearity of expectation})$ 

$$\mathbb{E}\left[\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}\right] = \sum_{x \in \mathcal{X}} (\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}) p(x) = \boldsymbol{A} \underbrace{\left(\sum_{x \in \mathcal{X}} \boldsymbol{x} p(x)\right)}_{\mathbb{E}[\boldsymbol{x}]} + \boldsymbol{b}$$

$$\begin{array}{l} \blacktriangleright \quad \mathbb{C}\left[\boldsymbol{A}\boldsymbol{x}\right] = \boldsymbol{A}\mathbb{C}\left[\boldsymbol{x}\right]\boldsymbol{A}^{\top} \\ \mathbb{C}\left[\boldsymbol{A}\boldsymbol{x}\right] = \mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\boldsymbol{x}^{\top}\boldsymbol{A}^{\top}\right] - \mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\right]\mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\right]^{\top} \\ = \boldsymbol{A}\left(\mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{\top}] - \mathbb{E}\left[\boldsymbol{x}\right]\mathbb{E}\left[\boldsymbol{x}\right]^{\top}\right)\boldsymbol{A}^{\top} \\ = \boldsymbol{A}\mathbb{C}\left[\boldsymbol{x}\right]\boldsymbol{A}^{\top} \end{array}$$

#### Linear transformation

For any random variable  $m{x}$  with mean  $\mathbb{E}\left[m{x}
ight]$  and covariance  $\mathbb{C}\left[m{x}
ight]$ 

•  $\mathbb{E}[Ax + b] = A\mathbb{E}[x] + b;$  (linearity of expectation)

$$\mathbb{E}\left[\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}\right] = \sum_{x \in \mathcal{X}} (\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}) p(x) = \boldsymbol{A} \underbrace{\left(\sum_{x \in \mathcal{X}} \boldsymbol{x} p(x)\right)}_{\mathbb{E}[\boldsymbol{x}]} + \boldsymbol{b}$$

$$\begin{array}{l} \blacktriangleright \quad \mathbb{C}\left[\boldsymbol{A}\boldsymbol{x}\right] = \boldsymbol{A}\mathbb{C}\left[\boldsymbol{x}\right]\boldsymbol{A}^{\top} \\ \mathbb{C}\left[\boldsymbol{A}\boldsymbol{x}\right] = \mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\boldsymbol{x}^{\top}\boldsymbol{A}^{\top}\right] - \mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\right]\mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\right]^{\top} \\ = \boldsymbol{A}\left(\mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{\top}] - \mathbb{E}\left[\boldsymbol{x}\right]\mathbb{E}\left[\boldsymbol{x}\right]^{\top}\right)\boldsymbol{A}^{\top} \\ = \boldsymbol{A}\mathbb{C}\left[\boldsymbol{x}\right]\boldsymbol{A}^{\top} \end{array}$$



#### Linear transformation

For any random variable x with mean  $\mathbb{E}\left[x
ight]$  and covariance  $\mathbb{C}\left[x
ight]$ 

•  $\mathbb{E}[Ax + b] = A\mathbb{E}[x] + b$ ; (linearity of expectation)

$$\mathbb{E}\left[\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}\right] = \sum_{x \in \mathcal{X}} (\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}) p(x) = \boldsymbol{A} \underbrace{\left(\sum_{x \in \mathcal{X}} \boldsymbol{x} p(x)\right)}_{\mathbb{E}[\boldsymbol{x}]} + \boldsymbol{b}$$

$$\begin{array}{l} \blacktriangleright \quad \mathbb{C}\left[\boldsymbol{A}\boldsymbol{x}\right] = \boldsymbol{A}\mathbb{C}\left[\boldsymbol{x}\right]\boldsymbol{A}^{\top} \\ \mathbb{C}\left[\boldsymbol{A}\boldsymbol{x}\right] = \mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\boldsymbol{x}^{\top}\boldsymbol{A}^{\top}\right] - \mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\right]\mathbb{E}\left[\boldsymbol{A}\boldsymbol{x}\right]^{\top} \\ = \boldsymbol{A}\left(\mathbb{E}[\boldsymbol{x}\boldsymbol{x}^{\top}] - \mathbb{E}\left[\boldsymbol{x}\right]\mathbb{E}\left[\boldsymbol{x}\right]^{\top}\right)\boldsymbol{A}^{\top} \\ = \boldsymbol{A}\mathbb{C}\left[\boldsymbol{x}\right]\boldsymbol{A}^{\top} \end{array}$$

$$oldsymbol{x} \sim \mathcal{N}\left(oldsymbol{\mu} \,,\, oldsymbol{\Sigma}
ight)$$
  
 $\Rightarrow oldsymbol{y} = oldsymbol{A}oldsymbol{x} + oldsymbol{b} \sim \mathcal{N}\left(oldsymbol{A}oldsymbol{\mu} + oldsymbol{b} \,,\, oldsymbol{A}oldsymbol{\Sigma}oldsymbol{A}^{ op}
ight)$ 



Independent case

Given a Gaussian distributed random variable  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \begin{bmatrix} \mathbf{1} & \mathbf{1} \\ \mathbb{E}[x_2] \end{bmatrix}, \qquad \mathbb{C}[\boldsymbol{x}] = \begin{bmatrix} \mathbf{1} & \mathbf{1} \\ 0 & \sigma_2^2 \end{bmatrix}$$

• What is  $p(x_1)$ ?

$$p(x_1) = \int_{x_2} \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \mid \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}\right) dx_2 =$$
$$= \int_{x_2} \mathcal{N}\left(x_1 \mid \mu_1, \sigma_1^2\right) \mathcal{N}\left(x_1 \mid \mu_1, \sigma_1^2\right) dx_2 =$$
$$= \mathcal{N}\left(x_1 \mid \mu_1, \sigma_1^2\right) \underbrace{\int_{x_2} \mathcal{N}\left(x_1 \mid \mu_1, \sigma_1^2\right) dx_2 =}_{1} =$$
$$= \mathcal{N}\left(x_1 \mid \mu_1, \sigma_1^2\right)$$

observation: The marginal distribution is Gaussian (Closure property!)

・ロト・西ト・ヨト・ヨー もくの

Independent case

Given a Gaussian distributed random variable  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \begin{bmatrix} \mathbb{E}[x_1] \\ \mathbb{E}[x_2] \end{bmatrix}, \qquad \mathbb{C}[\boldsymbol{x}] = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

• What is  $p(x_1)$ ?

$$p(x_1) = \int_{x_2} \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} | \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}\right) dx_2 =$$
$$= \int_{x_2} \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right) \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right) dx_2 =$$
$$= \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right) \underbrace{\int_{x_2} \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right) dx_2}_{1} =$$
$$= \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right)$$

observation: The marginal distribution is Gaussian (Closure property!)

Independent case

Given a Gaussian distributed random variable  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \left[ \begin{array}{c} \mathbb{E}[x_1] \\ \mathbb{E}[x_2] \end{array} \right], \qquad \mathbb{C}[\boldsymbol{x}] = \left[ \begin{array}{c} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{array} \right]$$

• What is  $p(x_1)$ ?

$$p(x_1) = \int_{x_2} \mathcal{N}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} | \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}\right) dx_2 =$$
$$= \int_{x_2} \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right) \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right) dx_2 =$$
$$= \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right) \underbrace{\int_{x_2} \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right) dx_2}_{1} =$$
$$= \mathcal{N}\left(x_1 | \mu_1, \sigma_1^2\right)$$

observation: The marginal distribution is Gaussian (Closure property!)

#### Multivariate case

Given a  $(N_1 + N_2)$ -dimensional Gaussian distributed random variable  $\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}_1^\top, \boldsymbol{x}_2^\top \end{bmatrix}^\top$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \left[ \begin{array}{cc} \mathbb{E}[\boldsymbol{x}_1] \\ \mathbb{E}[\boldsymbol{x}_2] \end{array} \right], \qquad \mathbb{C}[\boldsymbol{x}] = \left[ \begin{array}{cc} \mathbb{C}[\boldsymbol{x}_1] & \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2] \\ \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2]^\top & \mathbb{C}[\boldsymbol{x}_2] \end{array} \right]$$

- What is  $p(\boldsymbol{x}_1)$ ? ( $\boldsymbol{x}_2 \ N_1$ -dimensional)
- Can we find a linear transformation A, such that Ax = x<sub>1</sub>?



Attention: Dimensionality of the distribution is now N<sub>1</sub>, instead of N!
 If we used formula for N-dimensional normal distribution, then the distribution is not normalized!

#### Multivariate case

Given a  $(N_1 + N_2)$ -dimensional Gaussian distributed random variable  $\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}_1^\top, \boldsymbol{x}_2^\top \end{bmatrix}^\top$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \left[ \begin{array}{cc} \mathbb{E}[\boldsymbol{x}_1] \\ \mathbb{E}[\boldsymbol{x}_2] \end{array} \right], \qquad \mathbb{C}[\boldsymbol{x}] = \left[ \begin{array}{cc} \mathbb{C}[\boldsymbol{x}_1] & \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2] \\ \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2]^\top & \mathbb{C}[\boldsymbol{x}_2] \end{array} \right]$$

• What is  $p(\boldsymbol{x}_1)$ ? ( $\boldsymbol{x}_2 \ N_1$ -dimensional)

Can we find a linear transformation A, such that Ax = x<sub>1</sub>?



Attention: Dimensionality of the distribution is now N<sub>1</sub>, instead of N!
 If we used formula for N-dimensional normal distribution, then the distribution is not normalized!
#### Multivariate case

Given a  $(N_1 + N_2)$ -dimensional Gaussian distributed random variable  $\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}_1^\top, \boldsymbol{x}_2^\top \end{bmatrix}^\top$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \left[ egin{array}{cc} \mathbb{E}[\boldsymbol{x}_1] \\ \mathbb{E}[\boldsymbol{x}_2] \end{array} 
ight], \qquad \mathbb{C}[\boldsymbol{x}] = \left[ egin{array}{cc} \mathbb{C}[\boldsymbol{x}_1] & \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2] \\ \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2]^{ op} & \mathbb{C}[\boldsymbol{x}_2] \end{array} 
ight]$$

- What is  $p(\boldsymbol{x}_1)$ ? ( $\boldsymbol{x}_2$   $N_1$ -dimensional)
- Can we find a linear transformation A, such that Ax = x<sub>1</sub>?



#### Multivariate case

Given a  $(N_1 + N_2)$ -dimensional Gaussian distributed random variable  $\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}_1^\top, \boldsymbol{x}_2^\top \end{bmatrix}^\top$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \left[ egin{array}{cc} \mathbb{E}[\boldsymbol{x}_1] \\ \mathbb{E}[\boldsymbol{x}_2] \end{array} 
ight], \qquad \mathbb{C}[\boldsymbol{x}] = \left[ egin{array}{cc} \mathbb{C}[\boldsymbol{x}_1] & \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2] \\ \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2]^{ op} & \mathbb{C}[\boldsymbol{x}_2] \end{array} 
ight]$$

- What is  $p(\boldsymbol{x}_1)$ ? ( $\boldsymbol{x}_2$   $N_1$ -dimensional)
- Can we find a linear transformation A, such that Ax = x<sub>1</sub>?

$$oldsymbol{A} = \left[ egin{array}{cc} oldsymbol{I}_{N_1} & oldsymbol{0} \end{array} 
ight]$$



#### Multivariate case

Given a  $(N_1 + N_2)$ -dimensional Gaussian distributed random variable  $\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}_1^\top, \boldsymbol{x}_2^\top \end{bmatrix}^\top$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \left[ egin{array}{cc} \mathbb{E}[\boldsymbol{x}_1] \\ \mathbb{E}[\boldsymbol{x}_2] \end{array} 
ight], \qquad \mathbb{C}[\boldsymbol{x}] = \left[ egin{array}{cc} \mathbb{C}[\boldsymbol{x}_1] & \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2] \\ \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2]^{ op} & \mathbb{C}[\boldsymbol{x}_2] \end{array} 
ight]$$

- What is  $p(\boldsymbol{x}_1)$ ? ( $\boldsymbol{x}_2$   $N_1$ -dimensional)
- Can we find a linear transformation A, such that Ax = x<sub>1</sub>?

$$\boldsymbol{A} = \begin{bmatrix} \boldsymbol{I}_{N_1} & \boldsymbol{0} \end{bmatrix}$$
$$p(\boldsymbol{A}\boldsymbol{x}) = \mathcal{N} \begin{pmatrix} \boldsymbol{A}\boldsymbol{x} \mid \boldsymbol{A}\mathbb{E}[\boldsymbol{x}], \ \boldsymbol{A}\mathbb{C}[\boldsymbol{x}]\boldsymbol{A}^\top \end{bmatrix}$$
$$= \mathcal{N} (\boldsymbol{x}_1 \mid \mathbb{E}[\boldsymbol{x}_1], \ \mathbb{C}[\boldsymbol{x}_1])$$



#### Multivariate case

Given a  $(N_1 + N_2)$ -dimensional Gaussian distributed random variable  $\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}_1^\top, \boldsymbol{x}_2^\top \end{bmatrix}^\top$ , with mean and covariance

$$\mathbb{E}[\boldsymbol{x}] = \left[ egin{array}{cc} \mathbb{E}[\boldsymbol{x}_1] \\ \mathbb{E}[\boldsymbol{x}_2] \end{array} 
ight], \qquad \mathbb{C}[\boldsymbol{x}] = \left[ egin{array}{cc} \mathbb{C}[\boldsymbol{x}_1] & \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2] \\ \mathbb{C}[\boldsymbol{x}_1, \boldsymbol{x}_2]^{ op} & \mathbb{C}[\boldsymbol{x}_2] \end{array} 
ight]$$

- What is  $p(\boldsymbol{x}_1)$ ? ( $\boldsymbol{x}_2$   $N_1$ -dimensional)
- Can we find a linear transformation A, such that Ax = x<sub>1</sub>?

$$\boldsymbol{A} = \begin{bmatrix} \boldsymbol{I}_{N_1} & \boldsymbol{0} \end{bmatrix}$$
$$p(\boldsymbol{A}\boldsymbol{x}) = \mathcal{N} \begin{pmatrix} \boldsymbol{A}\boldsymbol{x} \mid \boldsymbol{A}\mathbb{E}[\boldsymbol{x}], \ \boldsymbol{A}\mathbb{C}[\boldsymbol{x}]\boldsymbol{A}^\top \end{bmatrix}$$
$$= \mathcal{N} (\boldsymbol{x}_1 \mid \mathbb{E}[\boldsymbol{x}_1], \ \mathbb{C}[\boldsymbol{x}_1])$$



Given a N-dimensional Gaussian distributed random variable  $\pmb{x}$ , with mean  $\mathbb{E}[\pmb{x}]$  and covariance  $\mathbb{C}[\pmb{x}]$ 

Can we find a linear transformation A, such that  $Ax = \sum_{n=1}^{N} x_n$ ?

Attention: Dimensionality of the distribution is now 1, instead of N! If we used formula for N-dimensional normal distribution, then the distribution is not normalized!

Given a N-dimensional Gaussian distributed random variable  $\pmb{x}$ , with mean  $\mathbb{E}[\pmb{x}]$  and covariance  $\mathbb{C}[\pmb{x}]$ 

• What is 
$$p\left(\sum_{n=1}^N x_n\right)$$
?

Can we find a linear transformation A, such that  $Ax = \sum_{n=1}^{N} x_n$ ?

Attention: Dimensionality of the distribution is now 1, instead of N! If we used formula for N-dimensional normal distribution, then the distribution is not normalized!

Given a N-dimensional Gaussian distributed random variable  $\pmb{x}$ , with mean  $\mathbb{E}[\pmb{x}]$  and covariance  $\mathbb{C}[\pmb{x}]$ 

• What is 
$$p\left(\sum_{n=1}^N x_n\right)$$
?

• Can we find a linear transformation A, such that  $Ax = \sum_{n=1}^{N} x_n$ ?

Attention: Dimensionality of the distribution is now 1, instead of N! If we used formula for N-dimensional normal distribution, then the distribution is not normalized!

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 少へ⊙

Given a N-dimensional Gaussian distributed random variable  $\pmb{x}$ , with mean  $\mathbb{E}[\pmb{x}]$  and covariance  $\mathbb{C}[\pmb{x}]$ 

• What is 
$$p\left(\sum_{n=1}^N x_n\right)$$
?

• Can we find a linear transformation A, such that  $Ax = \sum_{n=1}^{N} x_n$ ?

$$\boldsymbol{A} = \boldsymbol{1}_N = [1, \dots, 1]$$

Attention: Dimensionality of the distribution is now 1, instead of N! If we used formula for N-dimensional normal distribution, then the distribution is not normalized!

Given a N-dimensional Gaussian distributed random variable  $\pmb{x},$  with mean  $\mathbb{E}[\pmb{x}]$  and covariance  $\mathbb{C}[\pmb{x}]$ 

• What is 
$$p\left(\sum_{n=1}^{N} x_n\right)$$
?

• Can we find a linear transformation A, such that  $Ax = \sum_{n=1}^{N} x_n$ ?

$$\boldsymbol{A} = \boldsymbol{1}_{N} = [1, \dots, 1]$$
$$p(\boldsymbol{A}x) = \mathcal{N}\left(\boldsymbol{1}\boldsymbol{x} \mid \boldsymbol{1}_{N}\mathbb{E}[\boldsymbol{x}], \boldsymbol{1}_{N}\mathbb{C}[\boldsymbol{x}]\boldsymbol{1}_{N}^{\top}\right)$$
$$= \mathcal{N}\left(\sum_{n=1}^{N} x_{n} \mid \sum_{n=1}^{N}\mathbb{E}[x_{n}], \sum_{n=1}^{N}\sum_{m=1}^{N}\mathbb{C}[x_{n}, x_{m}]\right)$$

Attention: Dimensionality of the distribution is now 1, instead of N! If we used formula for N-dimensional normal distribution, then the distribution is not normalized!

Given a N-dimensional Gaussian distributed random variable  $\pmb{x}$ , with mean  $\mathbb{E}[\pmb{x}]$  and covariance  $\mathbb{C}[\pmb{x}]$ 

• What is 
$$p\left(\sum_{n=1}^{N} x_n\right)$$
?

• Can we find a linear transformation A, such that  $Ax = \sum_{n=1}^{N} x_n$ ?

$$\boldsymbol{A} = \boldsymbol{1}_{N} = [1, \dots, 1]$$
$$p(\boldsymbol{A}x) = \mathcal{N}\left(\boldsymbol{1}\boldsymbol{x} \mid \boldsymbol{1}_{N}\mathbb{E}[\boldsymbol{x}], \boldsymbol{1}_{N}\mathbb{C}[\boldsymbol{x}]\boldsymbol{1}_{N}^{\top}\right)$$
$$= \mathcal{N}\left(\sum_{n=1}^{N} x_{n} \mid \sum_{n=1}^{N}\mathbb{E}[x_{n}], \sum_{n=1}^{N}\sum_{m=1}^{N}\mathbb{C}[x_{n}, x_{m}]\right)$$

Attention: Dimensionality of the distribution is now 1, instead of N! If we used formula for N-dimensional normal distribution, then the distribution is not normalized!

• Given a univariate Gaussian distribution with arbitrary mean  $\mu$  and variance  $\sigma^2$ 

$$x \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

• Can I find a linear transformation y = ax + b, such that the resulting distribution has unit variance and zero mean?

 $y \sim \mathcal{N}(0, 1) \Leftrightarrow$ 



• Given a univariate Gaussian distribution with arbitrary mean  $\mu$  and variance  $\sigma^2$ 

$$x \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

• Can I find a linear transformation y = ax + b, such that the resulting distribution has unit variance and zero mean?

$$y \sim \mathcal{N}\left(0\,,\,1\right) \Leftrightarrow$$



• Given a univariate Gaussian distribution with arbitrary mean  $\mu$  and variance  $\sigma^2$ 

$$x \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

• Can I find a linear transformation y = ax + b, such that the resulting distribution has unit variance and zero mean?

$$y \sim \mathcal{N}(0, 1) \Leftrightarrow y = \frac{x - \mu}{\sigma} \Leftrightarrow$$



• Given a univariate Gaussian distribution with arbitrary mean  $\mu$  and variance  $\sigma^2$ 

$$x \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

Can I find a linear transformation y = ax + b, such that the resulting distribution has unit variance and zero mean?

$$y \sim \mathcal{N}(0, 1) \Leftrightarrow y = \frac{x - \mu}{\sigma} \Leftrightarrow a = \frac{1}{\sigma}, \quad b = \frac{\mu}{\sigma}$$



• Given a univariate Gaussian distribution with arbitrary mean  $\mu$  and variance  $\sigma^2$ 

$$x \sim \mathcal{N}\left(\mu, \sigma^2\right)$$

• Can I find a linear transformation y = ax + b, such that the resulting distribution has unit variance and zero mean?

$$y \sim \mathcal{N}(0, 1) \Leftrightarrow y = \frac{x - \mu}{\sigma} \Leftrightarrow a = \frac{1}{\sigma}, \quad b = \frac{\mu}{\sigma}$$



Multivariate case

### - Given a multivariate Gaussian distribution with arbitrary mean $\mu$ and covariance $\varSigma$

 $\mathcal{N}\left( \left. oldsymbol{x} \, \mid oldsymbol{\mu} \,, \, oldsymbol{\Sigma} \, 
ight) 
ight.$ 

• Can I find a linear transformation  $f(x) = Ax + \mu$ , such that the resulting distribution has unit covariance I and zero mean?



Multivariate case

 $\blacktriangleright$  Given a multivariate Gaussian distribution with arbitrary mean  $\mu$  and covariance  $\varSigma$ 

$$\mathcal{N}\left( \left. oldsymbol{x} \, \mid oldsymbol{\mu} \,, \, oldsymbol{\Sigma} \, 
ight)$$

• Can I find a linear transformation  $f(x) = Ax + \mu$ , such that the resulting distribution has unit covariance I and zero mean?



Multivariate case

 $\blacktriangleright$  Given a multivariate Gaussian distribution with arbitrary mean  $\mu$  and covariance  $\varSigma$ 

$$\mathcal{N}\left( \left. oldsymbol{x} 
ight. \mid oldsymbol{\mu} 
ight., \left. oldsymbol{\Sigma} 
ight. 
ight)$$

Can I find a linear transformation f(x) = Ax + μ, such that the resulting distribution has unit covariance I and zero mean?



A matrix  $oldsymbol{B}^{1/2}$  is called matrix square-root of a  $N ext{-by-}N$  symmetric matrix  $oldsymbol{B}$  iff

$$\boldsymbol{B}^{1/2}\boldsymbol{B}^{\top/2} = \boldsymbol{B}.$$

#### Examples of matrix square roots:

Cholesky factor L(triangular matrix)

- Cholesky factorization of  $B: B = LL^{\top}$ ,
- Computational complexity: O(N<sup>3</sup>).
- The Cholesky factorization is only available for positive definite matrices.
- Symmetric matrix square-root  $U\Lambda^{1/2}U^{\top}$  (unique),
  - $oldsymbol{U}$  orthogonal matrix. The columns are the eigenvectors of  $oldsymbol{B}$

$$\boldsymbol{U}\boldsymbol{U}^{ op} = \boldsymbol{U}^{ op}\boldsymbol{U} = \boldsymbol{I}$$

- $\Lambda = \operatorname{diag}[\lambda_1, \dots, \lambda_N]$  diagonal matrix of eigenvalues  $\lambda_n$ .
- Computational complexity: O(N<sup>3</sup>

$$\boldsymbol{B} = \boldsymbol{U}\boldsymbol{\Lambda}^{1/2} \underbrace{\boldsymbol{U}}_{\boldsymbol{I}}^{\top} \underbrace{\boldsymbol{U}}_{\boldsymbol{I}}^{\top} \boldsymbol{U}^{\top}$$

A matrix  ${m B}^{1/2}$  is called matrix square-root of a N-by-N symmetric matrix  ${m B}$  iff

$$\boldsymbol{B}^{1/2}\boldsymbol{B}^{\top/2} = \boldsymbol{B}.$$

Examples of matrix square roots:

- Cholesky factor L(triangular matrix)
  - Cholesky factorization of  $\boldsymbol{B}: \ \boldsymbol{B} = \boldsymbol{L} \boldsymbol{L}^{\top}$ ,
  - Computational complexity:  $O(N^3)$ .
  - The Cholesky factorization is only available for positive definite matrices.
- Symmetric matrix square-root  $\boldsymbol{U}\boldsymbol{\Lambda}^{1/2}\boldsymbol{U}^{ op}$  (unique),
  - ▶ U orthogonal matrix. The columns are the eigenvectors of B

$$\boldsymbol{U}\boldsymbol{U}^{ op} = \boldsymbol{U}^{ op}\boldsymbol{U} = \boldsymbol{I}$$

- $\Lambda = \operatorname{diag}[\lambda_1, \ldots, \lambda_N]$  diagonal matrix of eigenvalues  $\lambda_n$ .
- Computational complexity: O(N<sup>3</sup>

$$\boldsymbol{B} = \boldsymbol{U}\boldsymbol{\Lambda}^{1/2} \underbrace{\boldsymbol{U}}_{\boldsymbol{I}}^{ op} \boldsymbol{U} \boldsymbol{\Lambda}^{1/2} \boldsymbol{U}^{ op}$$

A matrix  $oldsymbol{B}^{1/2}$  is called matrix square-root of a N-by-N symmetric matrix  $oldsymbol{B}$  iff

$$\boldsymbol{B}^{1/2}\boldsymbol{B}^{\top/2} = \boldsymbol{B}.$$

Examples of matrix square roots:

- Cholesky factor L(triangular matrix)
  - Cholesky factorization of  $\boldsymbol{B}: \ \boldsymbol{B} = \boldsymbol{L} \boldsymbol{L}^{\top}$ ,
  - Computational complexity:  $O(N^3)$ .
  - The Cholesky factorization is only available for positive definite matrices.
- Symmetric matrix square-root  $U \Lambda^{1/2} U^{\top}$  (unique),
  - $\blacktriangleright$  *U* orthogonal matrix. The columns are the eigenvectors of *B*

$$\boldsymbol{U}\boldsymbol{U}^{ op}=\boldsymbol{U}^{ op}\boldsymbol{U}=\boldsymbol{I}$$

- $\Lambda = \operatorname{diag}[\lambda_1, \ldots, \lambda_N]$  diagonal matrix of eigenvalues  $\lambda_n$ .
- Computational complexity:  $O(N^3)$

$$\boldsymbol{B} = \boldsymbol{U}\boldsymbol{\Lambda}^{1/2} \underbrace{\boldsymbol{U}}^{\top} \underbrace{\boldsymbol{U}}^{\top} \boldsymbol{U} \boldsymbol{\Lambda}^{1/2} \boldsymbol{U}^{\top}$$

A matrix  $oldsymbol{B}^{1/2}$  is called matrix square-root of a N-by-N symmetric matrix  $oldsymbol{B}$  iff

$$\boldsymbol{B}^{1/2}\boldsymbol{B}^{\top/2} = \boldsymbol{B}.$$

Examples of matrix square roots:

- Cholesky factor L(triangular matrix)
  - Cholesky factorization of  $\boldsymbol{B}: \ \boldsymbol{B} = \boldsymbol{L} \boldsymbol{L}^{\top}$ ,
  - Computational complexity:  $O(N^3)$ .
  - The Cholesky factorization is only available for positive definite matrices.
- Symmetric matrix square-root  $U \Lambda^{1/2} U^{\top}$  (unique),
  - U orthogonal matrix. The columns are the eigenvectors of B

$$\boldsymbol{U}\boldsymbol{U}^{ op}=\boldsymbol{U}^{ op}\boldsymbol{U}=\boldsymbol{I}$$

- $\Lambda = \operatorname{diag}[\lambda_1, \ldots, \lambda_N]$  diagonal matrix of eigenvalues  $\lambda_n$ .
- Computational complexity:  $O(N^3)$

$$\boldsymbol{B} = \boldsymbol{U}\boldsymbol{\Lambda}^{1/2} \underbrace{\boldsymbol{U}}_{\boldsymbol{I}}^{\top} \underbrace{\boldsymbol{U}}_{\boldsymbol{I}} \boldsymbol{\Lambda}^{1/2} \boldsymbol{U}^{\top}$$

 $\blacktriangleright \Rightarrow \mathbf{\Lambda}^{1/2} \mathbf{U}^\top \text{ is also a matrix square root! (non-symmetric)}$ 

Application: Sampling from a multivariate Gaussian distribution

- We often need samples from a multivariate Gaussian distribution (For example: MCMC sampling)
- Sampling from a univariate Gaussian distribution is easy.

$$x \sim \mathcal{N}(0, 1)$$

How can we use N samples from a univariate Gaussian to imitate a single sample from a multivariate Gaussian with covariance Σ?

$$oldsymbol{x} \sim \mathcal{N}\left(oldsymbol{\mu} \,,\, oldsymbol{\Sigma}
ight)$$

- Use matrix square  $\Sigma^{1/2}$  root of the covariance matrix to transform the samples!
- Algorithm:
  - 1. Compute a square root of  $\varSigma$  (e.g. from Cholesky decomposition)
  - 2. Get N samples  $x_n$  from a standard normal distribution  $\mathcal{N}(0\,,\,1)$

- 3. Stack the samples into a vector  $oldsymbol{x} = [x_1, \dots, x_N]$
- 4. Multiply x by  $\varSigma^{1/2}$
- 5. Add  $\mu$ .
- 6. done.

Application: Sampling from a multivariate Gaussian distribution

- We often need samples from a multivariate Gaussian distribution (For example: MCMC sampling)
- Sampling from a univariate Gaussian distribution is easy.

$$x \sim \mathcal{N}(0, 1)$$

How can we use N samples from a univariate Gaussian to imitate a single sample from a multivariate Gaussian with covariance *D*?

$$oldsymbol{x} \sim \mathcal{N}\left(oldsymbol{\mu} \,, \, oldsymbol{\Sigma}
ight)$$

- Use matrix square  ${\cal D}^{1/2}$  root of the covariance matrix to transform the samples!
- Algorithm:
  - 1. Compute a square root of  $\varSigma$  (e.g. from Cholesky decomposition)
  - 2. Get N samples  $x_n$  from a standard normal distribution  $\mathcal{N}(0\,,\,1)$

- 3. Stack the samples into a vector  $m{x} = [x_1, \dots, x_N]$
- 4. Multiply x by  $\Sigma^{1/2}$
- 5. Add  $\mu$ .
- 6. done.

Application: Sampling from a multivariate Gaussian distribution

- We often need samples from a multivariate Gaussian distribution (For example: MCMC sampling)
- Sampling from a univariate Gaussian distribution is easy.

$$x \sim \mathcal{N}(0, 1)$$

How can we use N samples from a univariate Gaussian to imitate a single sample from a multivariate Gaussian with covariance *D*?

$$oldsymbol{x} \sim \mathcal{N}\left(oldsymbol{\mu} \,, \, oldsymbol{\Sigma}
ight)$$

- $\blacktriangleright$  Use matrix square  ${oldsymbol \Sigma}^{1/2}$  root of the covariance matrix to transform the samples!
- Algorithm:
  - 1. Compute a square root of  $\varSigma$  (e.g. from Cholesky decomposition)
  - 2. Get N samples  $x_n$  from a standard normal distribution  $\mathcal{N}(0\,,\,1)$

- 3. Stack the samples into a vector  $oldsymbol{x} = [x_1, \dots, x_N]$
- 4. Multiply x by  $\Sigma^{1/2}$
- 5. Add  $\mu$ .
- 6. done.

Application: Sampling from a multivariate Gaussian distribution

- We often need samples from a multivariate Gaussian distribution (For example: MCMC sampling)
- Sampling from a univariate Gaussian distribution is easy.

$$x \sim \mathcal{N}(0, 1)$$

How can we use N samples from a univariate Gaussian to imitate a single sample from a multivariate Gaussian with covariance *D*?

$$oldsymbol{x} \sim \mathcal{N}\left(oldsymbol{\mu} \,, \, oldsymbol{\varSigma}
ight)$$

• Use matrix square  $\Sigma^{1/2}$  root of the covariance matrix to transform the samples!

Algorithm:

- 1. Compute a square root of  $\varSigma$  (e.g. from Cholesky decomposition)
- 2. Get N samples  $x_n$  from a standard normal distribution  $\mathcal{N}(0\,,\,1)$

- 3. Stack the samples into a vector  $oldsymbol{x} = [x_1, \dots, x_N]$
- 4. Multiply x by  $\varSigma^{1/2}$
- 5. Add  $\mu$ .
- 6. done.

Application: Sampling from a multivariate Gaussian distribution

- We often need samples from a multivariate Gaussian distribution (For example: MCMC sampling)
- Sampling from a univariate Gaussian distribution is easy.

$$x \sim \mathcal{N}(0, 1)$$

How can we use N samples from a univariate Gaussian to imitate a single sample from a multivariate Gaussian with covariance *D*?

$$oldsymbol{x} \sim \mathcal{N}\left(oldsymbol{\mu} \,, \, oldsymbol{\Sigma}
ight)$$

Use matrix square  $\Sigma^{1/2}$  root of the covariance matrix to transform the samples!

### Algorithm:

- 1. Compute a square root of  $\varSigma$  (e.g. from Cholesky decomposition)
- 2. Get N samples  $x_n$  from a standard normal distribution  $\mathcal{N}(0, 1)$
- 3. Stack the samples into a vector  $oldsymbol{x} = [x_1, \dots, x_N]$
- 4. Multiply x by  $\Sigma^{1/2}$
- 5. Add  $\mu$ .
- 6. done.