

# Probabilistic Graphical Models

## Lecture 3: Gaussian Distributions

Volkan Cevher, Matthias Seeger  
Ecole Polytechnique Fédérale de Lausanne

7/10/2011



- 1 Why Gaussians?
- 2 Linear Transformations. Marginalization
- 3 Natural and Moment Parameterization
- 4 Schur Complement. Useful Identities from Conditioning
- 5 Products. Tower Formulae

# Why Gaussians?

Gaussian (aka. normal) distribution  $N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$

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

- Marginalization, conditioning, linear transformation, posterior:  
All just **linear algebra**
- Incredible **closedness** properties:
  - Linear transformations
  - Conditioning
  - Marginalization

Belief propagation needs such closedness

- Why all that?
  - Gaussians are **limit distributions** (central limit theorems)
  - Gaussians are **maximum entropy distributions**:  
No structure beyond mean, covariance

# Gaussians are Limit Distributions

## Central Limit Theorem

$\mathbf{x}_1 \sim P(\mathbf{x}_1)$ , mean  $\boldsymbol{\mu}$ , covariance  $\boldsymbol{\Sigma}$ .

Imagine independent, identically distributed (i.i.d.) replicas  $\mathbf{x}_2, \mathbf{x}_3, \dots$

$$\bar{\mathbf{x}}^{(n)} := \sqrt{n} \underbrace{\left( n^{-1} \sum_{i=1}^n \mathbf{x}_i - \boldsymbol{\mu} \right)}_{\rightarrow 0 \text{ a.s.}} \Rightarrow P(\bar{\mathbf{x}}^{(n)}) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Sigma})$$

# Gaussians are Limit Distributions

## Central Limit Theorem

$\mathbf{x}_1 \sim P(\mathbf{x}_1)$ , mean  $\boldsymbol{\mu}$ , covariance  $\boldsymbol{\Sigma}$ .

Imagine independent, identically distributed (i.i.d.) replicas  $\mathbf{x}_2, \mathbf{x}_3, \dots$

$$\bar{\mathbf{x}}^{(n)} := \underbrace{\sqrt{n} \left( n^{-1} \sum_{i=1}^n \mathbf{x}_i - \boldsymbol{\mu} \right)}_{\rightarrow 0 \text{ a.s.}} \Rightarrow P(\bar{\mathbf{x}}^{(n)}) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Sigma})$$

What does that mean?

- Averaging of i.i.d. variables: Mean, covariance retained
- Everything else smoothed away (by symmetry)  
 $\Rightarrow$  What remains: **Gaussian**

# Gaussians are Limit Distributions

## Central Limit Theorem

$\mathbf{x}_1 \sim P(\mathbf{x}_1)$ , mean  $\boldsymbol{\mu}$ , covariance  $\boldsymbol{\Sigma}$ .

Imagine independent, identically distributed (i.i.d.) replicas  $\mathbf{x}_2, \mathbf{x}_3, \dots$

$$\bar{\mathbf{x}}^{(n)} := \underbrace{\sqrt{n} \left( n^{-1} \sum_{i=1}^n \mathbf{x}_i - \boldsymbol{\mu} \right)}_{\rightarrow 0 \text{ a.s.}} \Rightarrow P(\bar{\mathbf{x}}^{(n)}) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Sigma})$$

What does that mean?

- Averaging of i.i.d. variables: Mean, covariance retained
- Everything else smoothed away (by symmetry)  
 $\Rightarrow$  What remains: **Gaussian**

For the meticulous:

If  $\mathbf{x}_1$  has no covariance, there are other stable distributions

# Gaussians are Limit Distributions

## Central Limit Theorem

$\mathbf{x}_1 \sim P(\mathbf{x}_1)$ , mean  $\boldsymbol{\mu}$ , covariance  $\boldsymbol{\Sigma}$ .

Imagine independent, identically distributed (i.i.d.) replicas  $\mathbf{x}_2, \mathbf{x}_3, \dots$

$$\bar{\mathbf{x}}^{(n)} := \sqrt{n} \underbrace{\left( n^{-1} \sum_{i=1}^n \mathbf{x}_i - \boldsymbol{\mu} \right)}_{\rightarrow 0 \text{ a.s.}} \Rightarrow P(\bar{\mathbf{x}}^{(n)}) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Sigma})$$

Implications for Statistics:

- Most models with finite number of parameters:  
Maximum likelihood estimator asymptotically normal

# Gaussians are Limit Distributions

## Central Limit Theorem

$\mathbf{x}_1 \sim P(\mathbf{x}_1)$ , mean  $\boldsymbol{\mu}$ , covariance  $\boldsymbol{\Sigma}$ .

Imagine independent, identically distributed (i.i.d.) replicas  $\mathbf{x}_2, \mathbf{x}_3, \dots$

$$\bar{\mathbf{x}}^{(n)} := \underbrace{\sqrt{n} \left( n^{-1} \sum_{i=1}^n \mathbf{x}_i - \boldsymbol{\mu} \right)}_{\rightarrow 0 \text{ a.s.}} \Rightarrow P(\bar{\mathbf{x}}^{(n)}) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Sigma})$$

Implications for closedness:

- Linear transformations, marginalization:  
Limit distributions **have** to be closed



# Gaussians are Maximum Entropy Distributions

How much uncertainty / little structure is in a distribution?

## Differential Entropy

$$H[P] = E_P[-\log P(\mathbf{x})] = \int P(\mathbf{x})(-\log P(\mathbf{x})) d\mathbf{x}$$

# Gaussians are Maximum Entropy Distributions

How much uncertainty / little structure is in a distribution?

## Differential Entropy

$$H[P] = E_P[-\log P(\mathbf{x})] = \int P(\mathbf{x})(-\log P(\mathbf{x})) d\mathbf{x}$$

### Information theory (Shannon)

Immensely useful, basis of probabilistic machine learning.

Part II: Scratch surface. But dig for yourself:

- Cover, Thomas: Elements of Information Theory (1991)

One of my top five all times favourite textbooks. Read it!

# Gaussians are Maximum Entropy Distributions

How much uncertainty / little structure is in a distribution?

## Differential Entropy

$$H[P] = E_P[-\log P(\mathbf{x})] = \int P(\mathbf{x})(-\log P(\mathbf{x})) d\mathbf{x}$$

- Given mean  $\mu$ , covariance  $\Sigma$ : Maximum entropy distribution?

$$\mathcal{N}(\mu, \Sigma) = \operatorname{argmax}_P \{H[P] \mid E_P[\mathbf{x}] = \mu, \operatorname{Cov}_P[\mathbf{x}] = \Sigma\}$$

# Gaussians are Maximum Entropy Distributions

How much uncertainty / little structure is in a distribution?

## Differential Entropy

$$H[P] = E_P[-\log P(\mathbf{x})] = \int P(\mathbf{x})(-\log P(\mathbf{x})) d\mathbf{x}$$

- Given mean  $\mu$ , covariance  $\Sigma$ : Maximum entropy distribution?

$$N(\mu, \Sigma) = \operatorname{argmax}_P \{H[P] \mid E_P[\mathbf{x}] = \mu, \operatorname{Cov}_P[\mathbf{x}] = \Sigma\}$$

- What does that mean?
  - Gaussian “nothing but mean and covariance”.  
Any other structure? It’s not a Gaussian
  - Would expect nice closedness properties for MaxEnt distributions
  - Upper bound on entropy:

$$H[P] \leq H[N(\mathbf{0}, \operatorname{Cov}_P[\mathbf{x}])] = \frac{1}{2} \log |2\pi e \operatorname{Cov}_P[\mathbf{x}]|$$

# Too Simple for Real World?

I want to model / learn **structure**.

Why should I care for an unstructured distribution?

# Too Simple for Real World?

I want to model / learn **structure**.

Why should I care for an unstructured distribution?

Gaussians are elementary **building blocks**

- Gaussian + Structure (latent variables) → Wealth of models  
⇒ We'll see a few in what follows
- Many distributions are Gaussian scale mixtures [part II]
- Gaussian (implicitly) behind much of classical estimation methodology
- Carrier distribution for approximate inference [part II]

# Too Simple for Real World?

I want to model / learn **structure**.

Why should I care for an unstructured distribution?

Gaussians are elementary **building blocks**

- Gaussian + Structure (latent variables) → Wealth of models  
⇒ We'll see a few in what follows
- Many distributions are Gaussian scale mixtures [part II]
- Gaussian (implicitly) behind much of classical estimation methodology
- Carrier distribution for approximate inference [part II]

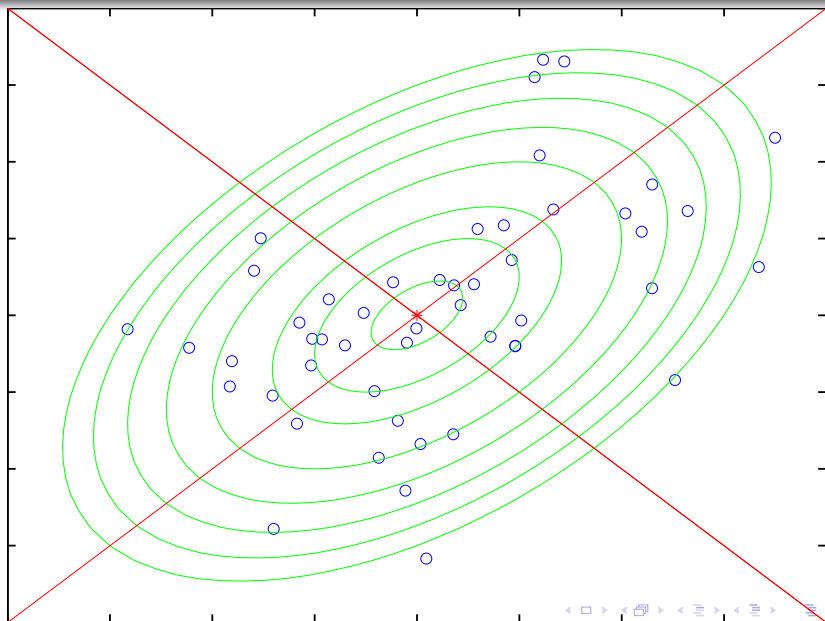
Maximum entropy for general variables / moments?

⇒ **Exponential families**

Not in this lecture, but dig for yourself:

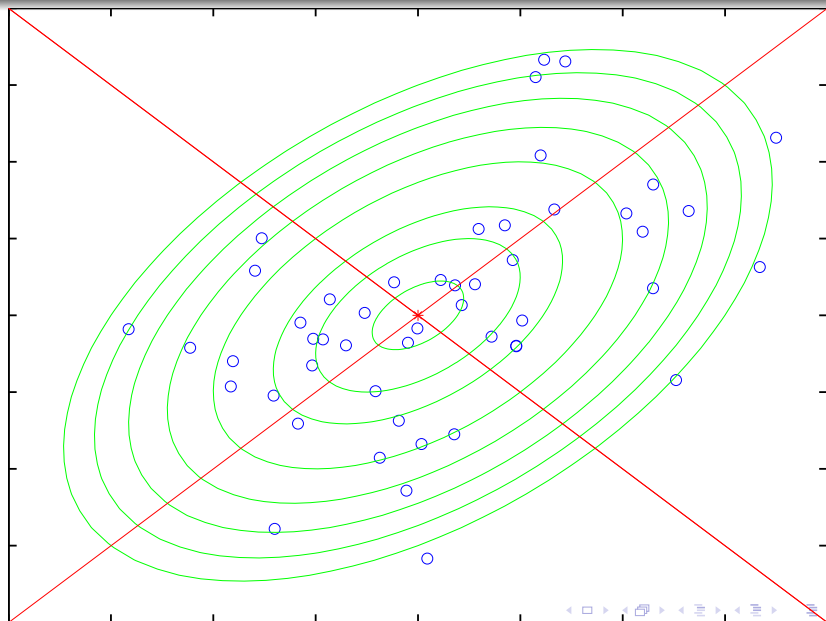
- M. Seeger: PhD thesis, Appendix A.4.1

# Gaussian Contours: Ellipsoids





# Linear Transformations



# Linear Transformations

- For any random variable  $\mathbf{x}$  with covariance [expectation linear!]

$$\mathbf{E}[\mathbf{Ax}] = \mathbf{A}\mathbf{E}[\mathbf{x}]$$

$$\begin{aligned}\text{Cov}[\mathbf{Ax}] &= \mathbf{E}[\mathbf{Ax}\mathbf{x}^T\mathbf{A}^T] - \mathbf{E}[\mathbf{Ax}]\mathbf{E}[\mathbf{Ax}]^T \\ &= \mathbf{A} \left( \mathbf{E}[\mathbf{x}\mathbf{x}^T] - \mathbf{E}[\mathbf{x}]\mathbf{E}[\mathbf{x}]^T \right) \mathbf{A}^T = \mathbf{A}\text{Cov}[\mathbf{x}]\mathbf{A}^T\end{aligned}$$

# Linear Transformations

- For any random variable  $\mathbf{x}$  with covariance [expectation linear!]

$$\mathbf{E}[\mathbf{Ax}] = \mathbf{A}\mathbf{E}[\mathbf{x}]$$

$$\begin{aligned} \text{Cov}[\mathbf{Ax}] &= \mathbf{E}[\mathbf{Ax}\mathbf{x}^T\mathbf{A}^T] - \mathbf{E}[\mathbf{Ax}]\mathbf{E}[\mathbf{Ax}]^T \\ &= \mathbf{A} \left( \mathbf{E}[\mathbf{x}\mathbf{x}^T] - \mathbf{E}[\mathbf{x}]\mathbf{E}[\mathbf{x}]^T \right) \mathbf{A}^T = \mathbf{A}\text{Cov}[\mathbf{x}]\mathbf{A}^T \end{aligned}$$

- Gaussian is just mean and covariance

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \Rightarrow \quad \mathbf{y} = \mathbf{Ax} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

# Linear Transformations

- For any random variable  $\mathbf{x}$  with covariance [expectation linear!]

$$\mathbf{E}[\mathbf{Ax}] = \mathbf{A}\mathbf{E}[\mathbf{x}]$$

$$\begin{aligned} \text{Cov}[\mathbf{Ax}] &= \mathbf{E}[\mathbf{Ax}\mathbf{x}^T\mathbf{A}^T] - \mathbf{E}[\mathbf{Ax}]\mathbf{E}[\mathbf{Ax}]^T \\ &= \mathbf{A} \left( \mathbf{E}[\mathbf{x}\mathbf{x}^T] - \mathbf{E}[\mathbf{x}]\mathbf{E}[\mathbf{x}]^T \right) \mathbf{A}^T = \mathbf{A}\text{Cov}[\mathbf{x}]\mathbf{A}^T \end{aligned}$$

- Gaussian is just mean and covariance

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \Rightarrow \quad \mathbf{y} = \mathbf{Ax} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

[Missing here: Formal proof that  $P(\mathbf{y})$  is Gaussian.  $\Rightarrow$  Ask me offline]

# Related Points

- Checking the normalization factor:

$\mathbf{x} \sim N(\mathbf{0}, \Sigma)$ .  $\Sigma = \mathbf{U}\Lambda\mathbf{U}^T$  eigendecomposition  
( $\mathbf{U}$  orthonormal (like rotation),  $\Lambda$  diagonal).

$\mathbf{y} = \mathbf{U}^T \mathbf{x}$  (rotate eigenvectors  $\rightarrow$  axes)  $\Rightarrow d\mathbf{y} = d\mathbf{x}$

$$\begin{aligned} \int e^{-\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x}} d\mathbf{x} &= \int e^{-\frac{1}{2}\mathbf{y}^T \mathbf{U}^T \Sigma^{-1} \mathbf{U} \mathbf{y}} d\mathbf{y} = \prod_i \int e^{-\frac{1}{2}y_i^2/\lambda_i} dy_i \\ &= \prod_i (2\pi\lambda_i)^{1/2} = |2\pi\Sigma|^{1/2} \text{ [determinant} = \prod \text{ eigenvalues]} \end{aligned}$$

Recall:  $\mathbf{U}^T \Sigma^{-1} \mathbf{U} = \Lambda^{-1}$ .

# Related Points

- Checking the normalization factor:

$\mathbf{x} \sim N(\mathbf{0}, \Sigma)$ .  $\Sigma = \mathbf{U}\Lambda\mathbf{U}^T$  eigendecomposition  
( $\mathbf{U}$  orthonormal (like rotation),  $\Lambda$  diagonal).

$\mathbf{y} = \mathbf{U}^T \mathbf{x}$  (rotate eigenvectors  $\rightarrow$  axes)  $\Rightarrow d\mathbf{y} = d\mathbf{x}$

$$\begin{aligned} \int e^{-\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x}} d\mathbf{x} &= \int e^{-\frac{1}{2}\mathbf{y}^T \mathbf{U}^T \Sigma^{-1} \mathbf{U} \mathbf{y}} d\mathbf{y} = \prod_i \int e^{-\frac{1}{2}y_i^2/\lambda_i} dy_i \\ &= \prod_i (2\pi\lambda_i)^{1/2} = |2\pi\Sigma|^{1/2} \text{ [determinant} = \prod \text{ eigenvalues]} \end{aligned}$$

Recall:  $\mathbf{U}^T \Sigma^{-1} \mathbf{U} = \Lambda^{-1}$ .

- For Gaussian:

$\Sigma$  diagonal  $\Rightarrow P(\mathbf{x}) = \prod_i P(x_i)$

Uncorrelated components  $\Rightarrow$  Independent components

# Marginal Distribution

- $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .  
 $I \subset \{1, \dots, n\}$ .  $\mathbf{x}_I := (x_i)_{i \in I}$ .  
Prize question: What is  $P(\mathbf{x}_I)$ ?

# Marginal Distribution

- $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .  
 $I \subset \{1, \dots, n\}$ .  $\mathbf{x}_I := (x_i)_{i \in I}$ .  
Prize question: What is  $P(\mathbf{x}_I)$ ?
- Pick selection matrix  $\mathbf{I}_I$ .  $\Rightarrow \mathbf{x}_I = \mathbf{I}_I \mathbf{x}$

$$P(\mathbf{x}_I) = N(\mathbf{I}_I \boldsymbol{\mu}, \mathbf{I}_I \boldsymbol{\Sigma} \mathbf{I}_I^T) = N(\boldsymbol{\mu}_I, \boldsymbol{\Sigma}_I)$$



# Marginal Distribution

- $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .  
 $I \subset \{1, \dots, n\}$ .  $\mathbf{x}_I := (x_i)_{i \in I}$ .  
 Prize question: What is  $P(\mathbf{x}_I)$ ?
- Pick selection matrix  $\mathbf{I}_I$ .  $\Rightarrow \mathbf{x}_I = \mathbf{I}_I \mathbf{x}$

$$P(\mathbf{x}_I) = N(\mathbf{I}_I \boldsymbol{\mu}, \mathbf{I}_I \boldsymbol{\Sigma} \mathbf{I}_I^T) = N(\boldsymbol{\mu}_I, \boldsymbol{\Sigma}_I)$$

- Marginalization (linear transformations):  
 Very simple if you have  $\boldsymbol{\mu}$ ,  $\boldsymbol{\Sigma}$

# Conditioning would be easy if . . .

- $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .  
 $I \subset \{1, \dots, n\}$ .  $R = \{1, \dots, n\} \setminus I$ .  
Next prize question: What is  $P(\mathbf{x}_I | \mathbf{x}_R)$ ?

# Conditioning would be easy if ...

- $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .  
 $I \subset \{1, \dots, n\}$ .  $R = \{1, \dots, n\} \setminus I$ .  
Next prize question: What is  $P(\mathbf{x}_I | \mathbf{x}_R)$ ?
- Not so simple. But it would be if ...

## Conditioning would be easy if ...

- $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .  
 $I \subset \{1, \dots, n\}$ .  $R = \{1, \dots, n\} \setminus I$ .  
 Next prize question: What is  $P(\mathbf{x}_I | \mathbf{x}_R)$ ?
- Not so simple. But it would be if ...

$$P(\mathbf{x}_I | \mathbf{x}_R) \propto e^{-\frac{1}{2}((\mathbf{x}_I - \boldsymbol{\mu}_I)^T \mathbf{A}_I (\mathbf{x}_I - \boldsymbol{\mu}_I) + 2(\mathbf{x}_R - \boldsymbol{\mu}_R)^T \mathbf{A}_{I,R}^T (\mathbf{x}_I - \boldsymbol{\mu}_I))},$$

$$\mathbf{A} = \boldsymbol{\Sigma}^{-1}$$

# Natural and Moment Parameterization

Two ways of parameterizing a Gaussian. You know:

Gaussian in moment (aka. mean) parameters  $\mu, \Sigma$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

# Natural and Moment Parameterization

Two ways of parameterizing a Gaussian. You know:

Gaussian in moment (aka. mean) parameters  $\mu, \Sigma$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \mathbf{A}(\mathbf{x}-\mu)}, \quad \mathbf{A} = \Sigma^{-1}$$

# Natural and Moment Parameterization

Two ways of parameterizing a Gaussian. You know:

Gaussian in moment (aka. mean) parameters  $\mu, \Sigma$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

$$\propto e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + (\mathbf{A}\mu)^T \mathbf{x}}, \quad \mathbf{A} = \Sigma^{-1}$$

# Natural and Moment Parameterization

Two ways of parameterizing a Gaussian. You know:

Gaussian in moment (aka. mean) parameters  $\mu, \Sigma$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

$$\propto e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{r}^T \mathbf{x}}, \quad \mathbf{A} = \Sigma^{-1}, \quad \mathbf{r} = \Sigma^{-1} \mu$$



# Natural and Moment Parameterization

Two ways of parameterizing a Gaussian. You know:

Gaussian in moment (aka. mean) parameters  $\mu, \Sigma$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

Now you know:

Gaussian in natural (aka. canonical) parameters  $\mathbf{r}, \mathbf{A}$

$$\propto e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{r}^T \mathbf{x}}, \quad \mathbf{A} = \Sigma^{-1}, \quad \mathbf{r} = \Sigma^{-1} \mu$$

# Natural and Moment Parameterization

Gaussian in moment (aka. mean) parameters  $\mu, \Sigma$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

Gaussian in natural (aka. canonical) parameters  $\mathbf{r}, \mathbf{A}$

$$\propto e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{r}^T \mathbf{x}}, \quad \mathbf{A} = \Sigma^{-1}, \quad \mathbf{r} = \Sigma^{-1} \mu$$

- Why two parameterizations for the same thing?
  - Some things simple in moment parameters:  
Linear transforms, marginalization [everything “sum”]
  - Some things simple in natural parameters:  
Conditioning, density product [everything “product”]

# Natural and Moment Parameterization

Gaussian in moment (aka. mean) parameters  $\mu, \Sigma$

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

Gaussian in natural (aka. canonical) parameters  $\mathbf{r}, \mathbf{A}$

$$\propto e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{r}^T \mathbf{x}}, \quad \mathbf{A} = \Sigma^{-1}, \quad \mathbf{r} = \Sigma^{-1} \mu$$

- Why two parameterizations for the same thing?
  - Some things simple in moment parameters:  
Linear transforms, marginalization [everything “**sum**”]
  - Some things simple in natural parameters:  
Conditioning, density product [everything “**product**”]
- For belief propagation (**sum-product**): Conversions all the time
- Conversion  $\leftrightarrow$  Matrix inversion  
 $\Rightarrow$  Makes Gaussian propagation **numerically difficult**

# Conditional Distribution

- $P(\mathbf{x}_I | \mathbf{x}_R)$ : What does it mean? **Factorization**
  - 1 Sample  $\mathbf{x}_R \sim N(\boldsymbol{\mu}_R, \boldsymbol{\Sigma}_R)$
  - 2 Sample  $\mathbf{x}_I$  from Gaussian depending on  $\mathbf{x}_R$
- $E[\mathbf{x}] = \boldsymbol{\mu}$ ,  $\text{Cov}[\mathbf{x}] = \boldsymbol{\Sigma}$  afterwards?  
⇒ Rule (2) must be  $P(\mathbf{x}_I | \mathbf{x}_R)$ !

For meticulous: We already know that  $P(\mathbf{x}_I | \mathbf{x}_R)$  is Gaussian (by inspection)

# Conditional Distribution

- $P(\mathbf{x}_I | \mathbf{x}_R)$ : What does it mean? **Factorization**
  - 1 Sample  $\mathbf{x}_R \sim N(\boldsymbol{\mu}_R, \boldsymbol{\Sigma}_R)$
  - 2 Sample  $\mathbf{x}_I$  from Gaussian depending on  $\mathbf{x}_R$
- $E[\mathbf{x}] = \boldsymbol{\mu}$ ,  $\text{Cov}[\mathbf{x}] = \boldsymbol{\Sigma}$  afterwards?  
 $\Rightarrow$  Rule (2) must be  $P(\mathbf{x}_I | \mathbf{x}_R)$ !

For meticulous: We already know that  $P(\mathbf{x}_I | \mathbf{x}_R)$  is Gaussian (by inspection)

- Ansatz:  $\mathbf{y} = \mathbf{x} - \boldsymbol{\mu}$ .  
 $\mathbf{y}_I = \mathbf{u} + \mathbf{B}\mathbf{y}_R$ ,  $\mathbf{u} \sim N(\mathbf{0}, \mathbf{C})$ .

# Conditional Distribution

- $P(\mathbf{x}_I|\mathbf{x}_R)$ : What does it mean? **Factorization**

- 1 Sample  $\mathbf{x}_R \sim N(\boldsymbol{\mu}_R, \boldsymbol{\Sigma}_R)$

- 2 Sample  $\mathbf{x}_I$  from Gaussian depending on  $\mathbf{x}_R$

- $E[\mathbf{x}] = \boldsymbol{\mu}$ ,  $\text{Cov}[\mathbf{x}] = \boldsymbol{\Sigma}$  afterwards?

$\Rightarrow$  Rule (2) must be  $P(\mathbf{x}_I|\mathbf{x}_R)$ !

For meticulous: We already know that  $P(\mathbf{x}_I|\mathbf{x}_R)$  is Gaussian (by inspection)

- Ansatz:  $\mathbf{y} = \mathbf{x} - \boldsymbol{\mu}$ .

$\mathbf{y}_I = \mathbf{u} + \mathbf{B}\mathbf{y}_R$ ,  $\mathbf{u} \sim N(\mathbf{0}, \mathbf{C})$ .

**Schur complement:**  $\mathbf{C} = \boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R := \boldsymbol{\Sigma}_I - \boldsymbol{\Sigma}_{I,R}\boldsymbol{\Sigma}_R^{-1}\boldsymbol{\Sigma}_{R,I}$

$$E[\mathbf{x}_I|\mathbf{x}_R] = \boldsymbol{\mu}_I + \boldsymbol{\Sigma}_{I,R}\boldsymbol{\Sigma}_R^{-1}(\mathbf{x}_R - \boldsymbol{\mu}_R),$$

$$\text{Cov}[\mathbf{x}_I|\mathbf{x}_R] = \boldsymbol{\Sigma}/\boldsymbol{\Sigma}_R = \boldsymbol{\Sigma}_I - \boldsymbol{\Sigma}_{I,R}\boldsymbol{\Sigma}_R^{-1}\boldsymbol{\Sigma}_{R,I}$$

- $E[\mathbf{x}_I|\mathbf{x}_R]$  linear in  $\mathbf{x}_R$ .  $\text{Cov}[\mathbf{x}_I|\mathbf{x}_R]$  independent of  $\mathbf{x}_R$

# The Schur Complement

$$\underbrace{P(\mathbf{x}_I, \mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma} = \underbrace{P(\mathbf{x}_I|\mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma/\Sigma_R} \times \underbrace{P(\mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma_R}$$

# The Schur Complement

$$\underbrace{P(\mathbf{x}_I, \mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma} = \underbrace{P(\mathbf{x}_I|\mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma/\Sigma_R} \times \underbrace{P(\mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma_R}$$

- Holds more generally, whenever  $\Sigma$ ,  $\Sigma_R$  nonsingular.  
Not just for symmetric  $\Sigma$
- Determinant identity

$$|\Sigma| = |\Sigma/\Sigma_R| \cdot |\Sigma_R|$$



# The Schur Complement

$$\underbrace{P(\mathbf{x}_I, \mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma} = \underbrace{P(\mathbf{x}_I|\mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma/\Sigma_R} \times \underbrace{P(\mathbf{x}_R)}_{\text{Cov}[\cdot]=\Sigma_R}$$

- Holds more generally, whenever  $\Sigma$ ,  $\Sigma_R$  nonsingular. Not just for symmetric  $\Sigma$
- Determinant identity

$$|\Sigma| = |\Sigma/\Sigma_R| \cdot |\Sigma_R|$$

Useful special case:

$$|I + UV| = |I + VU|$$

# Partitioned Inverse Equations

$$\Sigma^{-1} = \begin{bmatrix} \mathbf{A}_I & \mathbf{A}_{I,R} \\ \mathbf{A}_{R,I} & \mathbf{A}_R \end{bmatrix} = \begin{bmatrix} (\Sigma/\Sigma_R)^{-1} & -(\Sigma/\Sigma_R)^{-1} \mathbf{B} \\ -\mathbf{B}^T (\Sigma/\Sigma_R)^{-1} & \Sigma_R^{-1} + \mathbf{B}^T (\Sigma/\Sigma_R)^{-1} \mathbf{B} \end{bmatrix}$$

$$\mathbf{B} = \Sigma_{I,R} \Sigma_R^{-1}$$

- Very useful if  $|I|, |R|$  different  
 $\Rightarrow$  Do inverses in smaller of them only!

# Partitioned Inverse Equations

$$\Sigma^{-1} = \begin{bmatrix} \mathbf{A}_I & \mathbf{A}_{I,R} \\ \mathbf{A}_{R,I} & \mathbf{A}_R \end{bmatrix} = \begin{bmatrix} (\Sigma/\Sigma_R)^{-1} & -(\Sigma/\Sigma_R)^{-1} \mathbf{B} \\ -\mathbf{B}^T (\Sigma/\Sigma_R)^{-1} & \Sigma_R^{-1} + \mathbf{B}^T (\Sigma/\Sigma_R)^{-1} \mathbf{B} \end{bmatrix}$$

$$\mathbf{B} = \Sigma_{I,R} \Sigma_R^{-1}$$

- Very useful if  $|I|, |R|$  different  
 $\Rightarrow$  Do inverses in smaller of them only!
- Could have conditioned on  $\mathbf{x}_I$  just as well:

**Woodbury formula:**  $(\Sigma/\Sigma_I)^{-1} = \Sigma_R^{-1} + \mathbf{B}^T (\Sigma/\Sigma_R)^{-1} \mathbf{B}$

$$(\mathbf{E} + \mathbf{F}\mathbf{G}^{-1}\mathbf{H})^{-1} = \mathbf{E}^{-1} - \mathbf{E}^{-1}\mathbf{F}(\mathbf{G} + \mathbf{H}\mathbf{E}^{-1}\mathbf{F})^{-1}\mathbf{H}\mathbf{E}^{-1}$$

$\Rightarrow$  Not least formula to learn by heart

# Product of Gaussians

$$N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) = N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})C$$

- Product: Combination of messages / information

# Product of Gaussians

$$N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) = N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})C$$

- Product: Combination of messages / information
- Easy in **natural** parameters:

$$e^{\mathbf{r}_1^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \mathbf{A}_1 \mathbf{x}} \times e^{\mathbf{r}_2^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \mathbf{A}_2 \mathbf{x}} = e^{(\mathbf{r}_1 + \mathbf{r}_2)^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{x}}$$

⇒ Sum of natural parameters

$$\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2, \quad \mathbf{r} = \mathbf{r}_1 + \mathbf{r}_2$$

# Product of Gaussians

$$N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) = N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})C$$

- Product: Combination of messages / information
- Easy in **natural** parameters:

$$e^{\mathbf{r}_1^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \mathbf{A}_1 \mathbf{x}} \times e^{\mathbf{r}_2^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \mathbf{A}_2 \mathbf{x}} = e^{(\mathbf{r}_1 + \mathbf{r}_2)^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{x}}$$

⇒ Sum of natural parameters

$$\boldsymbol{\Sigma}^{-1} = \boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1}, \quad \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} = \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2$$

# Product of Gaussians

$$N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) = N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})C$$

- Product: Combination of messages / information
- Easy in **natural** parameters:

$$e^{\mathbf{r}_1^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \mathbf{A}_1 \mathbf{x}} \times e^{\mathbf{r}_2^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \mathbf{A}_2 \mathbf{x}} = e^{(\mathbf{r}_1 + \mathbf{r}_2)^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{x}}$$

⇒ Sum of natural parameters

$$\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1}, \quad \boldsymbol{\mu} = \underbrace{\boldsymbol{\Sigma}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2)}_{\text{"weighted avg."}}$$

# Product of Gaussians

$$N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) = N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})C$$

- Product: Combination of messages / information
- Easy in **natural** parameters:

$$e^{\mathbf{r}_1^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \mathbf{A}_1 \mathbf{x}} \times e^{\mathbf{r}_2^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T \mathbf{A}_2 \mathbf{x}} = e^{(\mathbf{r}_1 + \mathbf{r}_2)^T \mathbf{x} - \frac{1}{2} \mathbf{x}^T (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{x}}$$

⇒ Sum of natural parameters

$$\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1}, \quad \boldsymbol{\mu} = \underbrace{\boldsymbol{\Sigma}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2)}_{\text{"weighted avg."}}$$

- And  $C$ ? Often not needed. If you need it:  
**Sampling argument** (saves pages of algebra)



# Linear-Gaussian Model

$$\begin{aligned} \mathbf{u} &\sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) && \text{Prior} \\ \mathbf{y} &= \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \boldsymbol{\Psi}) && \text{Likelihood} \end{aligned}$$

- 1 Joint / marginal distribution: **Tower formulae**

$$\begin{aligned} \mathbb{E}[\mathbf{y}] &= \mathbb{E}[\mathbb{E}[\mathbf{y}|\mathbf{u}]], & \text{Cov}[\mathbf{y}] &= \text{Cov}[\mathbb{E}[\mathbf{y}|\mathbf{u}]] + \mathbb{E}[\text{Cov}[\mathbf{y}|\mathbf{u}]] \\ \text{Cov}[\mathbf{u}, \mathbf{y}] &= \text{Cov}[\mathbf{u}, \mathbb{E}[\mathbf{y}|\mathbf{u}]] + \mathbb{E}[\text{Cov}[\mathbf{u}, \mathbf{y}|\mathbf{u}]] \end{aligned}$$

# Linear-Gaussian Model

$$\begin{aligned} \mathbf{u} &\sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) && \text{Prior} \\ \mathbf{y} &= \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \boldsymbol{\Psi}) && \text{Likelihood} \end{aligned}$$

- 1 Joint / marginal distribution: **Tower formulae**

$$\begin{aligned} E[\mathbf{y}] &= E[E[\mathbf{y}|\mathbf{u}]], & \text{Cov}[\mathbf{y}] &= \text{Cov}[E[\mathbf{y}|\mathbf{u}]] + E[\text{Cov}[\mathbf{y}|\mathbf{u}]] \\ \text{Cov}[\mathbf{u}, \mathbf{y}] &= \text{Cov}[\mathbf{u}, E[\mathbf{y}|\mathbf{u}]] + E[\text{Cov}[\mathbf{u}, \mathbf{y}|\mathbf{u}]] \end{aligned}$$

- 2 Posterior: **Product** Prior  $\times$  Likelihood

$$\exp\left(-\frac{1}{2}[(\mathbf{y} - \mathbf{X}\mathbf{u})^T \boldsymbol{\Psi}^{-1}(\mathbf{y} - \mathbf{X}\mathbf{u}) + (\mathbf{u} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}_0^{-1}(\mathbf{u} - \boldsymbol{\mu}_0)]\right)$$

# Linear-Gaussian Model

$$\begin{aligned} \mathbf{u} &\sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) && \text{Prior} \\ \mathbf{y} &= \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \boldsymbol{\Psi}) && \text{Likelihood} \end{aligned}$$

- 1 Joint / marginal distribution: **Tower formulae**

$$\begin{aligned} E[\mathbf{y}] &= E[E[\mathbf{y}|\mathbf{u}]], & \text{Cov}[\mathbf{y}] &= \text{Cov}[E[\mathbf{y}|\mathbf{u}]] + E[\text{Cov}[\mathbf{y}|\mathbf{u}]] \\ \text{Cov}[\mathbf{u}, \mathbf{y}] &= \text{Cov}[\mathbf{u}, E[\mathbf{y}|\mathbf{u}]] + E[\text{Cov}[\mathbf{u}, \mathbf{y}|\mathbf{u}]] \end{aligned}$$

- 2 Posterior: **Product** Prior  $\times$  Likelihood

$$\exp\left(-\frac{1}{2}\left[\mathbf{u}^T \underbrace{(\mathbf{X}^T \boldsymbol{\Psi}^{-1} \mathbf{X} + \boldsymbol{\Sigma}_0^{-1})}_{\text{Cov}[\mathbf{u}|\mathbf{y}]^{-1}} \mathbf{u} - 2\mathbf{u}^T \underbrace{(\mathbf{X}^T \boldsymbol{\Psi}^{-1} \mathbf{y} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0)}_{\text{Cov}[\mathbf{u}|\mathbf{y}]^{-1} E[\mathbf{u}|\mathbf{y}]} + \dots\right]\right)$$

**Normal equations:**

$$E[\mathbf{u}|\mathbf{y}] = (\mathbf{X}^T \boldsymbol{\Psi}^{-1} \mathbf{X} + \boldsymbol{\Sigma}_0^{-1})^{-1} (\mathbf{X}^T \boldsymbol{\Psi}^{-1} \mathbf{y} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0)$$

# Linear-Gaussian Model

$$\mathbf{u} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

Prior

$$\mathbf{y} = \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \boldsymbol{\Psi})$$

Likelihood

- 2 Posterior: **Product** Prior  $\times$  Likelihood

$$\exp\left(-\frac{1}{2}\left[\mathbf{u}^T \underbrace{(\mathbf{X}^T \boldsymbol{\Psi}^{-1} \mathbf{X} + \boldsymbol{\Sigma}_0^{-1})}_{\text{Cov}[\mathbf{u}|\mathbf{y}]^{-1}} \mathbf{u} - 2\mathbf{u}^T \underbrace{(\mathbf{X}^T \boldsymbol{\Psi}^{-1} \mathbf{y} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0)}_{\text{Cov}[\mathbf{u}|\mathbf{y}]^{-1} \mathbb{E}[\mathbf{u}|\mathbf{y}]} + \dots\right]\right)$$

What if  $\mathbf{y}$  less coefficients than  $\mathbf{u}$ ?

# Linear-Gaussian Model

$$\begin{aligned} \mathbf{u} &\sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) && \text{Prior} \\ \mathbf{y} &= \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \boldsymbol{\Psi}) && \text{Likelihood} \end{aligned}$$

- 2 Posterior: **Product** Prior  $\times$  Likelihood

$$\exp\left(-\frac{1}{2}\left[\mathbf{u}^T \underbrace{(\mathbf{X}^T \boldsymbol{\Psi}^{-1} \mathbf{X} + \boldsymbol{\Sigma}_0^{-1})}_{\text{Cov}[\mathbf{u}|\mathbf{y}]^{-1}} \mathbf{u} - 2\mathbf{u}^T \underbrace{(\mathbf{X}^T \boldsymbol{\Psi}^{-1} \mathbf{y} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0)}_{\text{Cov}[\mathbf{u}|\mathbf{y}]^{-1} \mathbb{E}[\mathbf{u}|\mathbf{y}]} + \dots\right]\right)$$

What if  $\mathbf{y}$  less coefficients than  $\mathbf{u}$ ?

$$\begin{aligned} \text{Cov}[\mathbf{u}|\mathbf{y}] &= \text{Cov}[(\mathbf{u} \ \mathbf{y})] / \text{Cov}[\mathbf{y}] = \boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_0 \mathbf{X}^T (\boldsymbol{\Psi} + \mathbf{X} \boldsymbol{\Sigma}_0 \mathbf{X}^T)^{-1} \mathbf{X} \boldsymbol{\Sigma}_0, \\ \mathbb{E}[\mathbf{u}|\mathbf{y}] &= \mathbb{E}[\mathbf{u}] + \text{Cov}[\mathbf{u}, \mathbf{y}] \text{Cov}[\mathbf{y}]^{-1} (\mathbf{y} - \mathbb{E}[\mathbf{y}]) \\ &= \boldsymbol{\mu}_0 + \boldsymbol{\Sigma}_0 \mathbf{X}^T (\boldsymbol{\Psi} + \mathbf{X} \boldsymbol{\Sigma}_0 \mathbf{X}^T)^{-1} (\mathbf{y} - \mathbf{X} \boldsymbol{\mu}_0) \end{aligned}$$

# Wrap-Up

Practice those Gaussian calculations

- They come back at you all the time
- They look messy only as long as you don't understand them
- Short derivations take much less time (waste it with funnier things)
- Short derivations contain fewer mistakes
- Short derivations are just **so much cooler!**