

# Probabilistic Graphical Models

## Lecture 4: Essential Numerical Mathematics. Vector Calculus

Volkan Cevher, Matthias Seeger  
Ecole Polytechnique Fédérale de Lausanne

10/10/2011



- 1 Matrix Factorizations
- 2 Conjugate Gradients Algorithm
- 3 Vector Calculus

# Why Numerical Mathematics?

- 1 Great new idea? Have to run it on an imperfect machine!
  - Continuous variables: Messages are vectors / matrices, manipulated large number of times
  - Without care, errors (roundoff, cancellation, ...) accumulate
  - Greatest idea worth nothing if not numerically stable

# Why Numerical Mathematics?

- 1 Great new idea? Have to run it on an imperfect machine!
  - Continuous variables: Messages are vectors / matrices, manipulated large number of times
  - Without care, errors (roundoff, cancellation, ...) accumulate
  - Greatest idea worth nothing if not numerically stable
- 2 Thousands / millions of variables? Need iterative approximations
  - Backbone of about any continuous inference approximation: Iterative solvers from numerical mathematics
  - They are bottlenecks: Need to understand their properties

# Positive Definite Matrices

Positive definite  $\mathbf{A}$  (symmetric positive definite)

$$\mathbf{A}^T = \mathbf{A}, \quad \mathbf{v}^T \mathbf{A} \mathbf{v} > 0 \text{ for all } \mathbf{v} \neq \mathbf{0}$$

Equivalent to:

- All eigenvalues positive
- $\mathbf{A} = \mathbf{X}\mathbf{X}^T$ ,  $\mathbf{X}$  full rank

F2

# Positive Definite Matrices

Positive definite  $\mathbf{A}$  (symmetric positive definite)

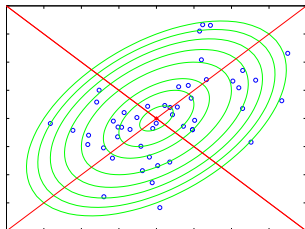
$$\mathbf{A}^T = \mathbf{A}, \quad \mathbf{v}^T \mathbf{A} \mathbf{v} > 0 \text{ for all } \mathbf{v} \neq \mathbf{0}$$

Equivalent to:

- All eigenvalues positive
- $\mathbf{A} = \mathbf{X}\mathbf{X}^T$ ,  $\mathbf{X}$  full rank

$\mathbf{A}$  pos. def.  $\Leftrightarrow$  Covariance matrix of Gaussian

F2(2)



# Positive Definite Matrices

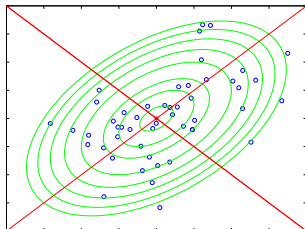
Positive definite  $\mathbf{A}$  (symmetric positive definite)

$$\mathbf{A}^T = \mathbf{A}, \quad \mathbf{v}^T \mathbf{A} \mathbf{v} > 0 \text{ for all } \mathbf{v} \neq \mathbf{0}$$

Equivalent to:

- All eigenvalues positive
- $\mathbf{A} = \mathbf{X}\mathbf{X}^T$ ,  $\mathbf{X}$  full rank

$\mathbf{A}$  pos. def.  $\Leftrightarrow$  Covariance matrix of Gaussian



Numbers	Matrices
$\mathbb{C}$	any square
$\mathbb{R}$	symmetric (hermitian)
$> 0$	positive definite

$\Rightarrow$  Pos. def. simplifies many methods

# Positive Definite Matrices

Positive definite  $\mathbf{A}$  (symmetric positive definite)

$$\mathbf{A}^T = \mathbf{A}, \quad \mathbf{v}^T \mathbf{A} \mathbf{v} > 0 \text{ for all } \mathbf{v} \neq \mathbf{0}$$

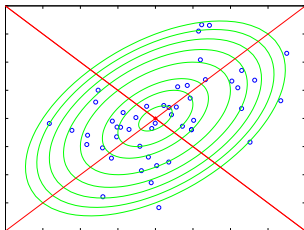
Equivalent to:

- All eigenvalues positive
- $\mathbf{A} = \mathbf{X}\mathbf{X}^T$ ,  $\mathbf{X}$  full rank

$\mathbf{A}$  pos. def.  $\Leftrightarrow$  Covariance matrix of Gaussian

Positive **semi**definite  $\mathbf{A}$ : “ $\geq$ ” instead of “ $>$ ”

$\Leftrightarrow \mathbf{A} = \mathbf{X}\mathbf{X}^T$ , any  $\mathbf{X}$





# Positive Definite Matrices

Positive definite  $\mathbf{A}$  (symmetric positive definite)

$$\mathbf{A}^T = \mathbf{A}, \quad \mathbf{v}^T \mathbf{A} \mathbf{v} > 0 \text{ for all } \mathbf{v} \neq \mathbf{0}$$

Equivalent to:

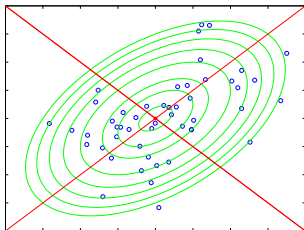
- All eigenvalues positive
- $\mathbf{A} = \mathbf{X}\mathbf{X}^T$ ,  $\mathbf{X}$  full rank

$\mathbf{A}$  pos. def.  $\Leftrightarrow$  Covariance matrix of Gaussian

Positive **semidefinite**  $\mathbf{A}$ : “ $\geq$ ” instead of “ $>$ ”

$\Leftrightarrow \mathbf{A} = \mathbf{X}\mathbf{X}^T$ , any  $\mathbf{X}$

$\mathbf{A}$  pos. semidef.  $\Leftrightarrow$  Covariance matrix of **degenerate** Gaussian  
[variance 0 along some directions]



[More about matrices? Horn, Johnson: Matrix Analysis (1985)]

# Cholesky Decomposition

What is a matrix decomposition? The right way to use  $\mathbf{A}^{-1}$ !

## Rule 1 for Matrix Computations

Do **not invert** a matrix. **Decompose** it

[Rule 2: Do not code it yourself. Use BLAS / LAPACK]

# Cholesky Decomposition

What is a matrix decomposition? The right way to use  $\mathbf{A}^{-1}$ !

## Rule 1 for Matrix Computations

Do **not invert** a matrix. **Decompose** it

[Rule 2: Do not code it yourself. Use BLAS / LAPACK]

Positive definite  $\mathbf{A}$ : **Cholesky decomposition**

$\mathbf{A} = \mathbf{X}\mathbf{X}^T$ . Can I use lower **triangular**  $\mathbf{L} = \mathbf{X}$ ?

⇒ Yes, exactly one: Cholesky factor of  $\mathbf{A}$

F3

# Cholesky Decomposition

What is a matrix decomposition? The right way to use  $\mathbf{A}^{-1}$ !

## Rule 1 for Matrix Computations

Do **not invert** a matrix. **Decompose** it

[Rule 2: Do not code it yourself. Use BLAS / LAPACK]

Positive definite  $\mathbf{A}$ : **Cholesky decomposition**

$\mathbf{A} = \mathbf{X}\mathbf{X}^T$ . Can I use lower **triangular**  $\mathbf{L} = \mathbf{X}$ ?

⇒ Yes, exactly one: Cholesky factor of  $\mathbf{A}$

Remarkable facts:

- Maybe the only algorithm in numerical mathematics that **is** so simple. If it fails,  $\mathbf{A}$  is not (numerically) pos. def.
- In-place algorithm:  $\mathbf{L}$  can overwrite  $\mathbf{A}$ .  $l_{i,j} > 0$
- Complexity  $O(n^3)$  [ $\mathbf{A} \in \mathbb{R}^{n \times n}$ ]

# Working with Cholesky Factors

Suppose:  $\mathbf{A} = \mathbf{L}\mathbf{L}^T \in \mathbb{R}^{n \times n}$ ,  $\mathbf{L}$  lower triangular

- Solving linear system:  $\mathbf{A}\mathbf{x} = \mathbf{L}\mathbf{L}^T\mathbf{x} = \mathbf{b}$

$$[1] : \mathbf{L}\mathbf{v} = \mathbf{b}, \quad [2] : \mathbf{L}^T\mathbf{x} = \mathbf{v}$$

Two **backsubstitutions**, at  $O(n^2)$

F4

[Beware: Some books distinguish between forward-, back-substitutions. I don't: It's just the same algorithm]

# Working with Cholesky Factors

Suppose:  $\mathbf{A} = \mathbf{L}\mathbf{L}^T \in \mathbb{R}^{n \times n}$ ,  $\mathbf{L}$  lower triangular

- Solving linear system:  $\mathbf{A}\mathbf{x} = \mathbf{L}\mathbf{L}^T\mathbf{x} = \mathbf{b}$

$$[1] : \mathbf{L}\mathbf{v} = \mathbf{b}, \quad [2] : \mathbf{L}^T\mathbf{x} = \mathbf{v}$$

Two **backsubstitutions**, at  $O(n^2)$

[Beware: Some books distinguish between forward-, back-substitutions. I don't: It's just the same algorithm]

- Log determinant:

$$\log |\mathbf{A}| = 2 \log |\mathbf{L}| = 2 \sum_{i=1}^n \log l_{i,i}$$

- Other expressions:

$$\mathbf{b}^T \mathbf{A}^{-1} \mathbf{b} = \|\mathbf{v}\|^2, \quad \mathbf{L}\mathbf{v} = \mathbf{b}$$

# Sequential Bayesian Updates

Example:

- Robot unsure about location, velocity:  $P(\mathbf{u}) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\mathbf{u} \in \mathbb{R}^n$
- Obtains noisy linear measurements:  $y = \mathbf{x}^T \mathbf{u} + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$
- Sequential update of belief state:  $P(\mathbf{u}) \rightarrow P(\mathbf{u}|\mathbf{x}, y), \dots$

# Sequential Bayesian Updates

Example:

- Robot unsure about location, velocity:  $P(\mathbf{u}) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\mathbf{u} \in \mathbb{R}^n$
- Obtains noisy linear measurements:  $y = \mathbf{x}^T \mathbf{u} + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$
- Sequential update of belief state:  $P(\mathbf{u}) \rightarrow P(\mathbf{u}|\mathbf{x}, y), \dots$

Recall last lecture:

$$\begin{aligned}\boldsymbol{\Sigma}' &= \boldsymbol{\Sigma} - \boldsymbol{\Sigma} \mathbf{x} (\sigma^2 + \mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x})^{-1} \mathbf{x}^T \boldsymbol{\Sigma}, \\ \boldsymbol{\mu}' &= \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{x} (\sigma^2 + \mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x})^{-1} (y - \mathbf{x}^T \boldsymbol{\mu})\end{aligned}$$

State representation: Could simply maintain  $\boldsymbol{\mu}, \boldsymbol{\Sigma}$



# Sequential Bayesian Updates

Example:

- Robot unsure about location, velocity:  $P(\mathbf{u}) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\mathbf{u} \in \mathbb{R}^n$
- Obtains noisy linear measurements:  $y = \mathbf{x}^T \mathbf{u} + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$
- Sequential update of belief state:  $P(\mathbf{u}) \rightarrow P(\mathbf{u}|\mathbf{x}, y), \dots$

Recall last lecture:

$$\begin{aligned}\boldsymbol{\Sigma}' &= \boldsymbol{\Sigma} - \boldsymbol{\Sigma} \mathbf{x} (\sigma^2 + \mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x})^{-1} \mathbf{x}^T \boldsymbol{\Sigma}, \\ \boldsymbol{\mu}' &= \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{x} (\sigma^2 + \mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x})^{-1} (y - \mathbf{x}^T \boldsymbol{\mu})\end{aligned}$$

State representation: Could simply maintain  $\boldsymbol{\mu}$ ,  $\boldsymbol{\Sigma}$

Better: Use Cholesky representation:  $\mathbf{L}$ ,  $\mathbf{a}$  s.t.  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$ ,  $\boldsymbol{\mu} = \mathbf{L}\mathbf{a}$

- Slightly more efficient
- Better numerical properties [details in exercise]

[Cholesky up-/downdates: TR at [people.mmc.uni-saarland.de/~mseeger/papers/cholupdate.pdf](http://people.mmc.uni-saarland.de/~mseeger/papers/cholupdate.pdf)  
Code at [people.mmc.uni-saarland.de/~mseeger/software.html](http://people.mmc.uni-saarland.de/~mseeger/software.html)]

# Singular Value Decomposition

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T \in \mathbb{R}^{m \times n}, \quad \mathbf{U}^T\mathbf{U} = \mathbf{I}, \quad \mathbf{V}^T\mathbf{V} = \mathbf{I}, \quad \mathbf{\Lambda} \text{ diagonal}, \lambda_i \geq 0$$

- Singular values  $\mathbf{\Lambda} \in \mathbb{R}^{d \times d}$  diagonal ( $d \leq \min\{m, n\}$ ):  
 $\lambda_i^2$  positive eigenvalues of  $\mathbf{A}^T\mathbf{A}$
- $\mathbf{A}$  symmetric  $\rightarrow \mathbf{U} = \mathbf{V}\mathbf{D}$ ,  $d_i \in \{\pm 1\}$ : Eigenvectors of  $\mathbf{A}$
- $\mathbf{A}$  positive definite  $\rightarrow$  SVD  $\equiv$  eigendecomposition

# Singular Value Decomposition

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T \in \mathbb{R}^{m \times n}, \quad \mathbf{U}^T\mathbf{U} = \mathbf{I}, \quad \mathbf{V}^T\mathbf{V} = \mathbf{I}, \quad \mathbf{\Lambda} \text{ diagonal}, \lambda_i \geq 0$$

- Singular values  $\mathbf{\Lambda} \in \mathbb{R}^{d \times d}$  diagonal ( $d \leq \min\{m, n\}$ ):  
 $\lambda_i^2$  positive eigenvalues of  $\mathbf{A}^T\mathbf{A}$
- $\mathbf{A}$  symmetric  $\rightarrow \mathbf{U} = \mathbf{V}\mathbf{D}$ ,  $d_i \in \{\pm 1\}$ : Eigenvectors of  $\mathbf{A}$
- $\mathbf{A}$  positive definite  $\rightarrow$  SVD  $\equiv$  eigendecomposition
- Can be computed robustly for about **any matrix**, in  $O(n^3)$
- For pos. def. matrices: SVD more expensive than Cholesky  
 $\Rightarrow$  Use latter if sufficient for your goal

# Singular Value Decomposition

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T \in \mathbb{R}^{m \times n}, \quad \mathbf{U}^T\mathbf{U} = \mathbf{I}, \quad \mathbf{V}^T\mathbf{V} = \mathbf{I}, \quad \mathbf{\Lambda} \text{ diagonal}, \lambda_i \geq 0$$

- Singular values  $\mathbf{\Lambda} \in \mathbb{R}^{d \times d}$  diagonal ( $d \leq \min\{m, n\}$ ):  
 $\lambda_i^2$  positive eigenvalues of  $\mathbf{A}^T\mathbf{A}$
- $\mathbf{A}$  symmetric  $\rightarrow \mathbf{U} = \mathbf{V}\mathbf{D}$ ,  $d_i \in \{\pm 1\}$ : Eigenvectors of  $\mathbf{A}$
- $\mathbf{A}$  positive definite  $\rightarrow$  SVD  $\equiv$  eigendecomposition
- Can be computed robustly for about **any matrix**, in  $O(n^3)$
- For pos. def. matrices: SVD more expensive than Cholesky  
 $\Rightarrow$  Use latter if sufficient for your goal
- Some applications in machine learning:
  - Principal components analysis (PCA)
  - Optimal low-rank matrix approximation (covariance explained)  
 $\Rightarrow$  Basis for linear approximations
  - Spectral clustering, manifold regularization, ... :  
 Leading eigenvectors carry lot of information about  $\mathbf{A}$

# Example: Design Optimization

You know / have learned:  $\mathbf{u} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  [but want to know more]

You can do: One more noisy linear measurement,  
 $y \sim N(\mathbf{x}^T \mathbf{u}, \sigma^2), \|\mathbf{x}\| = 1$

Which  $\mathbf{x}$  do you choose?

F7

# Example: Design Optimization

You know / have learned:  $\mathbf{u} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  [but want to know more]

You can do: One more noisy linear measurement,  
 $y \sim N(\mathbf{x}^T \mathbf{u}, \sigma^2), \|\mathbf{x}\| = 1$

Which  $\mathbf{x}$  do you choose?

$\mathbf{x}_* = \operatorname{argmax}\{\mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x} \mid \|\mathbf{x}\| = 1\}$ : Leading eigenvector  $\boldsymbol{\Sigma}$

Remember MRI sampling optimization [first lecture]?

⇒ That's (part of) how it works [stay on for rest]

# Why Iterative Solvers?

- Moderate-sized high-resolution image:  $n = 65536$  pixels.  
Storage: 32G (single matrix)  
Time for Cholesky decomposition:  $\approx 3\text{h}$  (if enough memory)  
 $\Rightarrow$  Bayesian methods over images, stacks of images, videos, ... ?

# Why Iterative Solvers?

- Moderate-sized high-resolution image:  $n = 65536$  pixels.  
 Storage: 32G (single matrix)  
 Time for Cholesky decomposition:  $\approx 3\text{h}$  (if enough memory)  
 $\Rightarrow$  Bayesian methods over images, stacks of images, videos, ... ?
- And how about structure?

$$E[\mathbf{u}|\mathbf{y}] = (\mathbf{X}^T \Psi^{-1} \mathbf{X} + \mathbf{\Pi})^{-1} (\mathbf{X}^T \Psi^{-1} \mathbf{y} + \mathbf{b})$$

- $\mathbf{X}$  sparse (most entries = 0) [e.g., consumer-product ratings]
- $\mathbf{X}$  structured (banded, Toeplitz, ...) [e.g., Markov grid structure]
- $\mathbf{X}$  fast operator (FFT) [e.g., MRI measurements]

Matrix decompositions do not make use of that

[Some do: sparse Cholesky factorization]

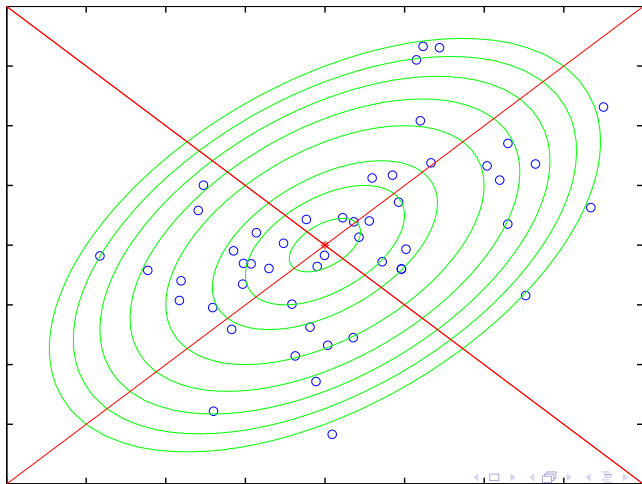


# Minimizing Quadratic Functions

For positive definite  $\mathbf{A}$ :

F9

$$\mathbf{x}_* = \operatorname{argmin}\{q(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}\} \Leftrightarrow \mathbf{A} \mathbf{x}_* = \mathbf{b}$$



# Minimizing Quadratic Functions

$$q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}, \quad \mathbf{g}(\mathbf{x}) = \nabla q(\mathbf{x}) = \mathbf{A} \mathbf{x} - \mathbf{b}$$

**Require:** Operator  $\mathbf{A}$ . Initial  $\mathbf{x}_0$

**for**  $k = 1, 2, \dots$  **do**

Pick search direction  $\mathbf{d}_k$ , based on  $\mathbf{g}_{k-1} = \mathbf{g}(\mathbf{x}_{k-1})$ ,  $\mathbf{d}_l, l < k$

Line minimization:

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_k, \quad \alpha_k = \operatorname{argmin}_{\alpha} q(\mathbf{x}_{k-1} + \alpha \mathbf{d}_k)$$

**end for**

# Conjugate Directions

- Why, of course down **as steep as possible**

$$q(\mathbf{x}_{k-1} + d\mathbf{x}) = q(\mathbf{x}_{k-1}) + \underbrace{\mathbf{g}_{k-1}^T(d\mathbf{x})}_{\text{Smallest: } d\mathbf{x} \propto -\mathbf{g}_{k-1}} + O(\|d\mathbf{x}\|^2)$$

Steepest descent:  $\mathbf{d}_k = -\mathbf{g}_{k-1}$

F11

# Conjugate Directions

- Why, of course down **as steep as possible**

$$q(\mathbf{x}_{k-1} + d\mathbf{x}) = q(\mathbf{x}_{k-1}) + \underbrace{\mathbf{g}_{k-1}^T(d\mathbf{x})}_{\text{Smallest: } d\mathbf{x} \propto -\mathbf{g}_{k-1}} + O(\|d\mathbf{x}\|^2)$$

Steepest descent:  $\mathbf{d}_k = -\mathbf{g}_{k-1}$

- Wrong:** For steepest descent:  $\mathbf{d}_{k+1}^T \mathbf{d}_k = 0$  F11(2)  
 $\Rightarrow$  Improvements from previous iterations rapidly tempered with

# Conjugate Directions

- Why, of course down **as steep as possible**

$$q(\mathbf{x}_{k-1} + d\mathbf{x}) = q(\mathbf{x}_{k-1}) + \underbrace{\mathbf{g}_{k-1}^T(d\mathbf{x})}_{\text{Smallest: } d\mathbf{x} \propto -\mathbf{g}_{k-1}} + O(\|d\mathbf{x}\|^2)$$

Steepest descent:  $\mathbf{d}_k = -\mathbf{g}_{k-1}$

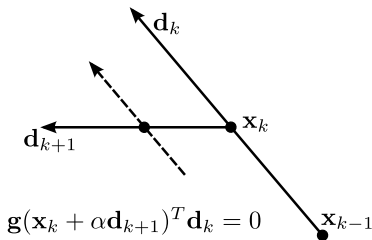
- Wrong:** For steepest descent:  $\mathbf{d}_{k+1}^T \mathbf{d}_k = 0$   
 $\Rightarrow$  Improvements from previous iterations rapidly tempered with

**New** gradients  $\perp$  **old** directions?

$\Rightarrow$  Retains previous efforts:

$$\mathbf{g}_k^T \mathbf{d}_k = 0 \rightarrow \mathbf{g}_{k+1}^T \mathbf{d}_k = 0 \dots$$

F11(3)



# Conjugate Directions

- Why, of course down **as steep as possible**

$$q(\mathbf{x}_{k-1} + d\mathbf{x}) = q(\mathbf{x}_{k-1}) + \underbrace{\mathbf{g}_{k-1}^T(d\mathbf{x})}_{\text{Smallest: } d\mathbf{x} \propto -\mathbf{g}_{k-1}} + O(\|d\mathbf{x}\|^2)$$

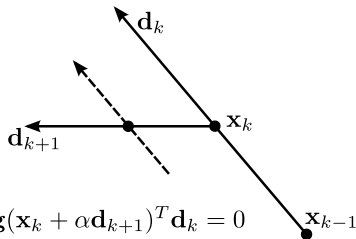
Steepest descent:  $\mathbf{d}_k = -\mathbf{g}_{k-1}$

- Wrong:** For steepest descent:  $\mathbf{d}_{k+1}^T \mathbf{d}_k = 0$   
 $\Rightarrow$  Improvements from previous iterations rapidly tempered with

**New** gradients  $\perp$  **old** directions?

$\Rightarrow$  Retains previous efforts:

$$\mathbf{g}_k^T \mathbf{d}_k = 0 \rightarrow \mathbf{g}_{k+1}^T \mathbf{d}_k = 0 \dots$$



$$\mathbf{g}(\mathbf{x}_k + \alpha \mathbf{d}_{k+1})^T \mathbf{d}_k = 0$$

## Conjugate Directions

$$\mathbf{d}_k^T \mathbf{A} \mathbf{d}_j \text{ for all } j < k$$

# Towards Conjugate Gradients

Details: **Handout**

[people.mmc.uni-saarland.de/~mseeger/lectures/bml09/handout\\_lect4.pdf](http://people.mmc.uni-saarland.de/~mseeger/lectures/bml09/handout_lect4.pdf)

- ① Directions conjugate: Gradient  $\mathbf{g}_k \perp$  **all** previous directions:  
 $\mathbf{g}_k^T \mathbf{d}_j = 0$  for all  $j \leq k$

# Towards Conjugate Gradients

Details: **Handout**

[people.mmc.uni-saarland.de/~mseeger/lectures/bml09/handout\\_lect4.pdf](http://people.mmc.uni-saarland.de/~mseeger/lectures/bml09/handout_lect4.pdf)

- 1 Directions conjugate: Gradient  $\mathbf{g}_k \perp$  **all** previous directions:  
 $\mathbf{g}_k^T \mathbf{d}_j = 0$  for all  $j \leq k$
- 2 After  $n$  steps we are done:  $\mathbf{g}_n = \mathbf{0}$
- 3 Construct conjugate directions by recurrence:  
 $\mathbf{d}_k = -\mathbf{g}_{k-1} + \beta_{k-1} \mathbf{d}_{k-1}$



# Towards Conjugate Gradients

Details: **Handout**

people.mhci.uni-saarland.de/~mseeger/lectures/bml09/handout\_lect4.pdf

- ① Directions conjugate: Gradient  $\mathbf{g}_k \perp$  **all** previous directions:  
 $\mathbf{g}_k^T \mathbf{d}_j = 0$  for all  $j \leq k$
- ② After  $n$  steps we are done:  $\mathbf{g}_n = \mathbf{0}$
- ③ Construct conjugate directions by recurrence:  
 $\mathbf{d}_k = -\mathbf{g}_{k-1} + \beta_{k-1} \mathbf{d}_{k-1}$
- ④ All gradients are orthogonal:  $\mathbf{g}_k^T \mathbf{g}_j = 0, j < k$   
 [Bit of misnomer: **Directions** are conjugate]

# Towards Conjugate Gradients

Details: **Handout**

people.mmc.uni-saarland.de/~mseeger/lectures/bml09/handout\_lect4.pdf

- 1 Directions conjugate: Gradient  $\mathbf{g}_k \perp$  **all** previous directions:  
 $\mathbf{g}_k^T \mathbf{d}_j = 0$  for all  $j \leq k$
- 2 After  $n$  steps we are done:  $\mathbf{g}_n = \mathbf{0}$
- 3 Construct conjugate directions by recurrence:  
 $\mathbf{d}_k = -\mathbf{g}_{k-1} + \beta_{k-1} \mathbf{d}_{k-1}$
- 4 All gradients are orthogonal:  $\mathbf{g}_k^T \mathbf{g}_j = 0, j < k$   
 [Bit of misnomer: **Directions** are conjugate]
- 5 What is  $\alpha_k$ ? From line minimization:

$$\alpha_k = \frac{\|\mathbf{g}_{k-1}\|^2}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k}$$

- 6 What is  $\beta_k$ ? The great synthesis!

$$\beta_k = \frac{\|\mathbf{g}_k\|^2}{\|\mathbf{g}_{k-1}\|^2}$$

# Conjugate Gradients (for reference)

**Require:** Operator  $\mathbf{A}$ . Initial  $\mathbf{x}_0$ .  $\mathbf{g}_0 = \mathbf{Ax}_0 - \mathbf{b}$

**for**  $k = 1, 2, \dots$  (no more than  $n$ ) **do**

$$\rho_{k-1} = \|\mathbf{g}_{k-1}\|^2$$

**if**  $k = 1$  **then**

$$\mathbf{d}_1 = -\mathbf{g}_0$$

**else**

$$\beta_{k-1} = \rho_{k-1} / \rho_{k-2}; \mathbf{d}_k = -\mathbf{g}_{k-1} + \beta_{k-1} \mathbf{d}_{k-1}$$

**end if**

$$\mathbf{q}_k = \mathbf{A} \mathbf{d}_k; \alpha_k = \rho_{k-1} / (\mathbf{d}_k^T \mathbf{q}_k)$$

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_k; \mathbf{g}_k = \mathbf{g}_{k-1} + \alpha_k \mathbf{q}_k$$

Check for convergence (say  $\|\mathbf{g}_k\| / \|\mathbf{b}\| < \epsilon$ )

**end for**

# Krylov Subspaces

Let  $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{d}_1, \dots, \mathbf{d}_k\}$ . Then,

F14

$$\mathbf{g}_k^T \mathbf{d}_j = 0, j \leq k \quad \Rightarrow \quad \mathbf{x}_k = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}_k} q(\mathbf{x})$$

# Krylov Subspaces

Let  $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{d}_1, \dots, \mathbf{d}_k\}$ . Then,

$$\mathbf{g}_k^T \mathbf{d}_j = 0, j \leq k \quad \Rightarrow \quad \mathbf{x}_k = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}_k} q(\mathbf{x})$$

But  $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{A}^j \mathbf{g}_0 \mid j < k\}$

$\Rightarrow$  Optimal with  $k$  ( $\mathbf{A}\cdot$ ) multiplications!

F14(2)

- $\mathcal{K}_k \subset \mathcal{K}_{k+1} \subset \dots, \mathbf{x}_* \in \mathcal{K}_n$  (Cayley/Hamilton)

# Krylov Subspaces

Let  $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{d}_1, \dots, \mathbf{d}_k\}$ . Then,

$$\mathbf{g}_k^T \mathbf{d}_j = 0, j \leq k \quad \Rightarrow \quad \mathbf{x}_k = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}_k} q(\mathbf{x})$$

But  $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{A}^j \mathbf{g}_0 \mid j < k\}$

$\Rightarrow$  Optimal with  $k$  ( $\mathbf{A} \cdot$ ) multiplications!

- $\mathcal{K}_k \subset \mathcal{K}_{k+1} \subset \dots, \mathbf{x}_* \in \mathcal{K}_n$  (Cayley/Hamilton)
- What about  $k \ll n$  for huge  $n$ ? Depends on eigenspectrum of  $\mathbf{A}$ .  
 $\mathbf{x}_k \approx \mathbf{x}_*$  in surprisingly many cases in practice  
 $\Rightarrow$  Krylov subspace view key to convergence analysis [exercise]

# Krylov Subspaces

Let  $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{d}_1, \dots, \mathbf{d}_k\}$ . Then,

$$\mathbf{g}_k^T \mathbf{d}_j = 0, j \leq k \quad \Rightarrow \quad \mathbf{x}_k = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}_k} q(\mathbf{x})$$

But  $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{A}^j \mathbf{g}_0 \mid j < k\}$

$\Rightarrow$  Optimal with  $k$  ( $\mathbf{A} \cdot$ ) multiplications!

- $\mathcal{K}_k \subset \mathcal{K}_{k+1} \subset \dots, \mathbf{x}_* \in \mathcal{K}_n$  (Cayley/Hamilton)
- What about  $k \ll n$  for huge  $n$ ? Depends on eigenspectrum of  $\mathbf{A}$ .  
 $\mathbf{x}_k \approx \mathbf{x}_*$  in surprisingly many cases in practice  
 $\Rightarrow$  Krylov subspace view key to convergence analysis [exercise]
- **Preconditioning:**  $\mathbf{M} = \mathbf{C}\mathbf{C}^T \approx \mathbf{A}$ , but **easy** to solve systems with
  - Work on  $(\mathbf{C}^{-T} \mathbf{A} \mathbf{C}^{-1}) \mathbf{C} \mathbf{x} = \mathbf{C}^{-T} \mathbf{b}$   
 $\Rightarrow$  Better spectral properties  $\rightarrow$  Faster convergence
  - CG as before, with one  $(\mathbf{M}^{-1} \cdot)$  per iteration
  - Preconditioning: Art of iterative linear solvers

# Why Vector Calculus?

- Remember differentiation? A bunch of rules, no-brainer.  
Do that in  $\mathbb{R}^n$ :  $(\partial f_i)/(\partial x_j) = \sum_k (\partial g_i)/(\partial y_k) \cdot (\partial y_k)/(\partial x_j), \dots?$



# Why Vector Calculus?

- Remember differentiation? A bunch of rules, no-brainer.  
Do that in  $\mathbb{R}^n$ :  $(\partial f_i)/(\partial x_j) = \sum_k (\partial g_i)/(\partial y_k) \cdot (\partial y_k)/(\partial x_j), \dots?$
- No! Use bunch of rules on vectors, matrices.**  
No  $\sum_{j,k}$ , no  $\partial_i/\partial_j$ : Waste of paper / your time
- Vector calculus in machine learning:
  - $\nabla f(\mathbf{x}) = \mathbf{0}$ , solve for  $\mathbf{x}$  (if you're lucky)
  - Search directions always fed by gradient (steepest descent)
  - Newton (second order) optimization: Hessian as well

Like with Gaussian: Experience with the rules pays off!

# Differential Notation

- Trace of square matrix

F16

$$\operatorname{tr} \mathbf{A} = \sum_i \alpha_{i,i} = \sum_i \lambda_i, \quad \operatorname{spec}(\mathbf{A}) = \{\lambda_1, \dots, \lambda_n\}$$

$$\operatorname{tr} \mathbf{AB} = \operatorname{tr} \mathbf{BA}, \quad \mathbf{x}^T \mathbf{Ax} \stackrel{!}{=} \operatorname{tr} \mathbf{x}^T \mathbf{Ax} = \operatorname{tr} \mathbf{Ax} \mathbf{x}^T$$

# Differential Notation

- Trace of square matrix

$$\text{tr } \mathbf{A} = \sum_i \alpha_{i,i} = \sum_i \lambda_i, \quad \text{spec}(\mathbf{A}) = \{\lambda_1, \dots, \lambda_n\}$$

$$\text{tr } \mathbf{AB} = \text{tr } \mathbf{BA}, \quad \mathbf{x}^T \mathbf{Ax} \stackrel{!}{=} \text{tr } \mathbf{x}^T \mathbf{Ax} = \text{tr } \mathbf{Ax} \mathbf{x}^T$$

- Differential  $d\mathbf{x}$ : Tiny vector, situated at  $\mathbf{x}$ .

F16(2)

$$df(\mathbf{x}) = f(\mathbf{x} + (d\mathbf{x})) - f(\mathbf{x}) = (\nabla f(\mathbf{x}))^T (d\mathbf{x}) + O(\|d\mathbf{x}\|^2), \dots$$

$\Rightarrow$  For vector calculus:  $d\mathbf{x}$  special vector, obeying one more rule:  
Term with  $d\mathbf{x} \geq 2$  times  $\rightarrow$  Term = 0 ( $\geq 3$  for Hessian)

[All I do here: Minka's note, <http://research.microsoft.com/en-us/um/people/minka/papers/matrix/>]

# Simple Rules

- Constant. Linear

$$d\mathbf{A} = \mathbf{0} \quad [\mathbf{A} \text{ constant}]$$

$$d(\alpha\mathbf{X} + \beta\mathbf{Y}) = \alpha(d\mathbf{X}) + \beta(d\mathbf{Y}), \quad d \operatorname{tr} \mathbf{X} = \operatorname{tr}(d\mathbf{X})$$

- Product rule

$$d(\mathbf{X}\mathbf{Y}) = (d\mathbf{X})\mathbf{Y} + \mathbf{X}(d\mathbf{Y}) \quad [\text{also for } \otimes, \circ]$$

- Permutation/extraction:  $(\cdot)^*$  reorders/extracts entries

$$d(\mathbf{X}^*) = (d\mathbf{X})^*$$

Example:  $(\cdot)^T$ ,  $\operatorname{diag}^{-1}(\cdot)$ , vectorization (`reshape` in Matlab)

Prove any of them:  $\text{lhs} = \text{rhs} + O(\|d\mathbf{x}\|^2)$

# More Interesting Rules

- Matrix inverse

F18

$$d(\mathbf{X}^{-1}) = -\mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-1}$$

# More Interesting Rules

- Matrix inverse

$$d(\mathbf{X}^{-1}) = -\mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-1}$$

- Log determinant

F18(2)

$$d \log |\mathbf{X}| = \text{tr } \mathbf{X}^{-1}(d\mathbf{X})$$

$$\text{And } d|\mathbf{X}| = d e^{\log |\mathbf{X}|} = |\mathbf{X}|(d \log |\mathbf{X}|)$$

# The “Algorithm”

Given: Really messy  $f(\mathbf{x})$ ,  $\mathbf{f}(\mathbf{x})$ ,  $f(\mathbf{X})$ ,  $\mathbf{F}(x)$ .

- 1 Inward: Write  $df(\mathbf{x})$  (or other forms). Use rules to push  $d(\cdot)$  inside, until  $d\mathbf{x}$  only

# The “Algorithm”

Given: Really messy  $f(\mathbf{x})$ ,  $\mathbf{f}(\mathbf{x})$ ,  $f(\mathbf{X})$ ,  $\mathbf{F}(\mathbf{x})$ .

- 1 Inward: Write  $df(\mathbf{x})$  (or other forms). Use rules to push  $d(\cdot)$  inside, until  $d\mathbf{x}$  only
- 2 Outward: Use linear algebra rules ( $d\mathbf{x}$  vector,  $d\mathbf{X}$  matrix!) to pull  $d\mathbf{x}$  out, until:

$df = a(dx)$	$d\mathbf{f} = \mathbf{a}(dx)$	$d\mathbf{F} = \mathbf{A}(dx)$
$df = \mathbf{a}^T(d\mathbf{x})$	$d\mathbf{f} = \mathbf{A}(d\mathbf{x})$	
$df = \text{tr } \mathbf{A}^T(d\mathbf{X})$		

Just read off derivative:  $a$ ,  $\mathbf{a}$ ,  $\mathbf{A}$ .

Empty cells? Don't like tensors beyond matrices. Use  $\text{vec}(\mathbf{X})$



# The “Algorithm”

Given: Really messy  $f(\mathbf{x})$ ,  $\mathbf{f}(\mathbf{x})$ ,  $f(\mathbf{X})$ ,  $\mathbf{F}(\mathbf{x})$ .

- 1 Inward: Write  $df(\mathbf{x})$  (or other forms). Use rules to push  $d(\cdot)$  inside, until  $d\mathbf{x}$  only
- 2 Outward: Use linear algebra rules ( $d\mathbf{x}$  vector,  $d\mathbf{X}$  matrix!) to pull  $d\mathbf{x}$  out, until:

$df = a(dx)$	$d\mathbf{f} = \mathbf{a}(dx)$	$d\mathbf{F} = \mathbf{A}(dx)$
$df = \mathbf{a}^T(d\mathbf{x})$	$d\mathbf{f} = \mathbf{A}(d\mathbf{x})$	
$df = \text{tr } \mathbf{A}^T(d\mathbf{X})$		

Just read off derivative:  $a$ ,  $\mathbf{a}$ ,  $\mathbf{A}$ .

Empty cells? Don't like tensors beyond matrices. Use  $\text{vec}(\mathbf{X})$

Constraints on  $\mathbf{X}$ ? Inherited by  $d\mathbf{X}$ ,  $\mathbf{A}$

F19

$\mathbf{X}$  symmetric  $\Rightarrow df = \text{tr}(\text{sym } \mathbf{A})(d\mathbf{X})$ ,  $\text{sym } \mathbf{A} = (\mathbf{A} + \mathbf{A}^T)/2$

$\mathbf{X}$  diagonal  $\Rightarrow df = (\text{diag } \text{diag}^{-1}(\mathbf{A}))(d\mathbf{X})$

# Example

Model:  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\mathbf{x} \in \mathbb{R}^d$

Data: Independent draws  $\mathbf{x}_1, \dots, \mathbf{x}_n$ ,  $n > d$

Suppose  $\boldsymbol{\mu}$  known. Maximum likelihood estimator for  $\boldsymbol{\Sigma}$ ?

F20

# Summary

$\mathbf{A} = (a_{ij})_{ij}$ . Then:

$$\mathbf{A}^T = (a_{ji})_{ij}$$

$$|\mathbf{A}|$$

$$\text{tr } \mathbf{A} = \sum_i a_{ii}$$

$$\text{diag}^{-1}(\mathbf{A}) = (a_{ii})_i$$

$$(\text{diag } \mathbf{v}) = (v_i \mathbf{I}_{\{i=j\}})_{ij}$$

$$\text{sym } \mathbf{A} = (\mathbf{A} + \mathbf{A}^T)/2$$

$$\mathbf{A}_{I,J} = (a_{ij})_{i \in I, j \in J}$$

$$\mathbf{I}_{\{I, \cdot\}} \mathbf{v} = \mathbf{v}_I$$

$$\mathbf{I}_{\{\cdot, I\}} \mathbf{v} = (\sum_k \mathbf{I}_{\{j=i_k\}} v_k)_j$$

Transpose

$$\text{Determinant } [|\mathbf{I} + \mathbf{A}\mathbf{B}| = |\mathbf{I} + \mathbf{B}\mathbf{A}|]$$

$$\text{Trace } [\text{tr } \mathbf{A}\mathbf{B} = \text{tr } \mathbf{B}\mathbf{A}]$$

Diagonal of matrix

Diagonal matrix

Symmetrization

Subselection (“ $I$ ” for  $\{i\}$ , “ $\cdot$ ” for full range)

Subselection matrix

Distribution matrix ( $I = \{i_k\}$ )

# Wrap-Up

- Lots of stuff: Seems hard, tedious at first
- As computer scientists, we are engineers (get things done):  
These rules / basic algorithms are our toolbox
- Better know your basic tools very well (with practice, you will)