

# Inference and Representation

DS-GA-1005, CSCI-GA.2569

Joan Bruna

Courant Institute of Mathematical Sciences  
Center for Data Science  
NYU



# Announcements

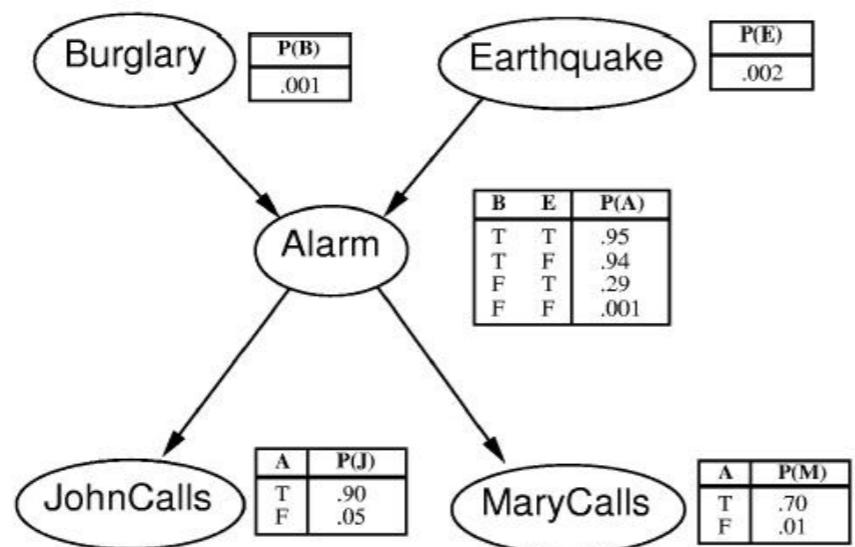
- There will be lecture next week (no guest lecture).
- Midterm 10/30: Lectures 1 through 4 (including today), PS1 , PS2 and PS3.
- PS4 released tomorrow.

# Gibbs Sampling

- Gibbs Sampling is an iterative algorithm that produces samples from undirected models.
- Suppose the model contains variables  $x_1 \dots x_n$
- Initialize starting values (e.g from uniform distribution)
- Do until (convergence):
  - Pick an ordering of the variables
  - For each  $x_i$ ,
    - ❖ Sample  $p(x_i \mid X_j = x_j), j \neq i$  .
    - ❖ update  $x_i$
- Recall that we only need to condition on the Markov Blanket.

# Gibbs Sampling

## Gibbs Sampling: An Example

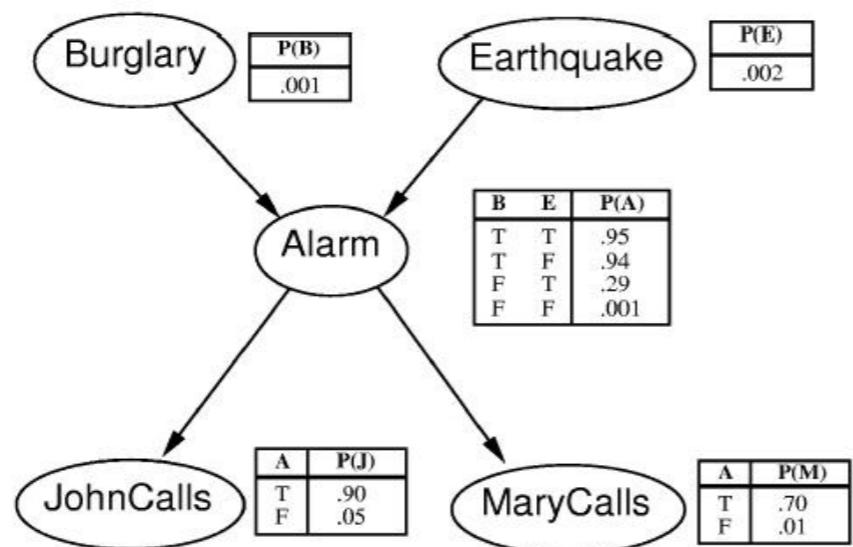


t	B	E	A	J	M
0	F	F	F	F	F
1					
2					
3					
4					

- Consider the alarm network
  - Assume we sample variables in the order B,E,A,J,M
  - Initialize all variables at t = 0 to False

# Gibbs Sampling

## Gibbs Sampling: An Example



t	B	E	A	J	M
0	F	F	F	F	F
1			F		
2					
3					
4					

- Sampling  $P(B|A,E)$  at  $t = 1$ : Using Bayes Rule,

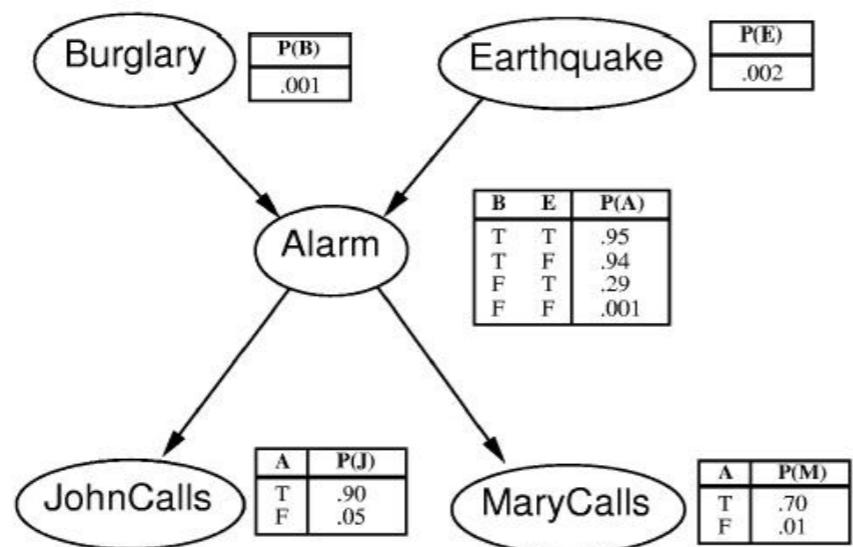
$$P(B | A, E) \propto P(A | B, E)P(B)$$

- $A=\text{false}$ ,  $E=\text{false}$ , so we compute:

$$P(B = T | A = F, E = F) \propto (0.06)(0.01) = 0.0006$$

$$P(B = F | A = F, E = F) \propto (0.999)(0.999) = 0.9980$$

# Gibbs Sampling: An Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T			
2					
3					
4					

- Sampling  $P(E|A,B)$ : Using Bayes Rule,

$$P(E | A, B) \propto P(A | B, E)P(E)$$

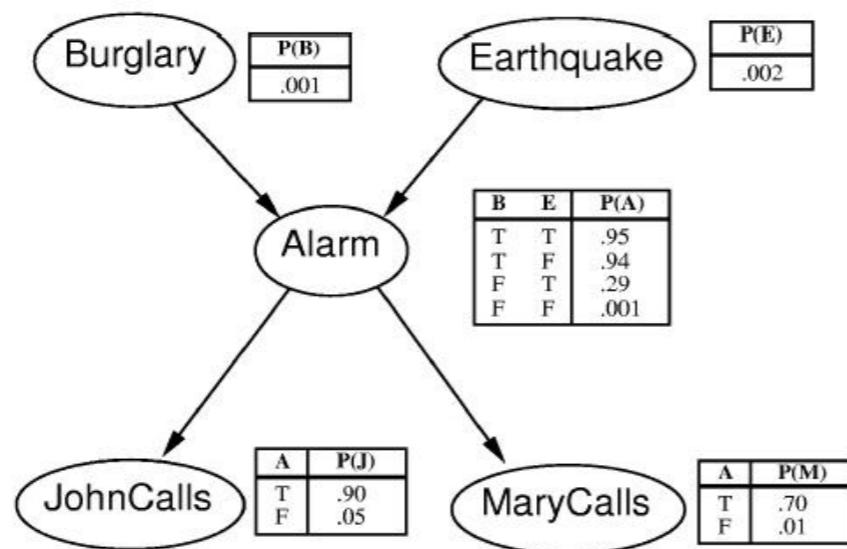
- $(A,B) = (F,F)$ , so we compute the following,

$$P(E = T | A = F, B = F) \propto (0.71)(0.02) = 0.0142$$

$$P(E = F | A = F, B = F) \propto (0.999)(0.998) = 0.9970$$

# Gibbs Sampling

## Gibbs Sampling: An Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F		
2					
3					
4					

- Sampling  $P(A|B,E,J,M)$ : Using Bayes Rule,

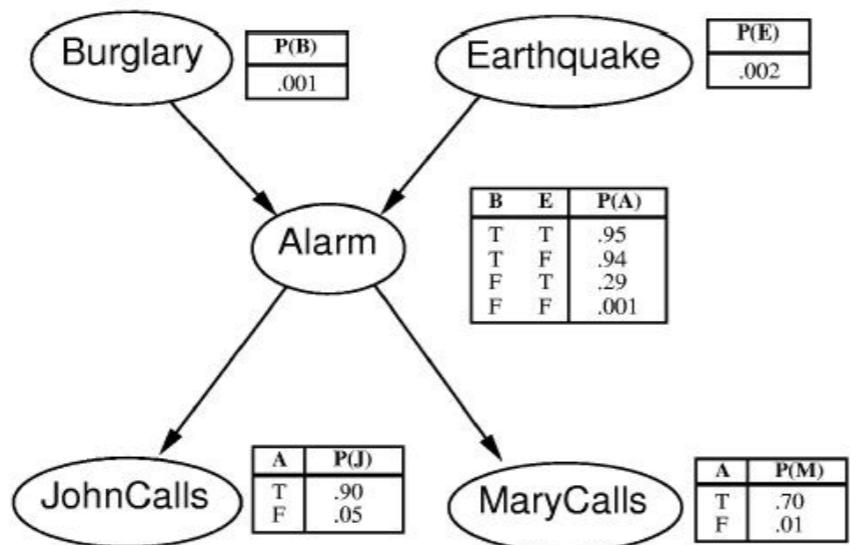
$$P(A | B, E, J, M) \propto P(J | A)P(M | A)P(A | B, E)$$

- $(B, E, J, M) = (F, T, F, F)$ , so we compute:

$$P(A = T | B = F, E = T, J = F, M = F) \propto (0.1)(0.3)(0.29) = 0.0087$$

$$P(A = F | B = F, E = T, J = F, M = F) \propto (0.95)(0.99)(0.71) = 0.6678$$

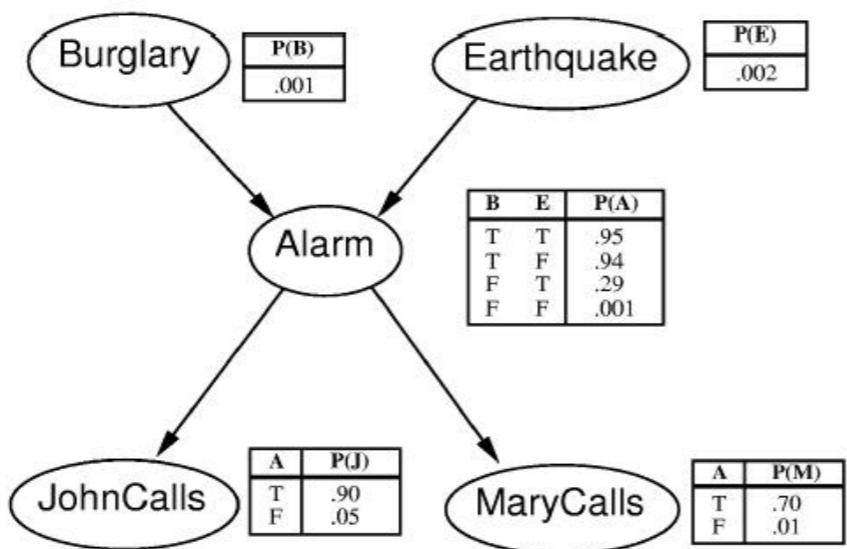
# Gibbs Sampling: An Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	
2					
3					
4					

- Sampling  $P(J|A)$ : No need to apply Bayes Rule
- $A = F$ , so we compute the following, and sample
  - $P(J = T | A = F) \propto 0.05$
  - $P(J = F | A = F) \propto 0.95$

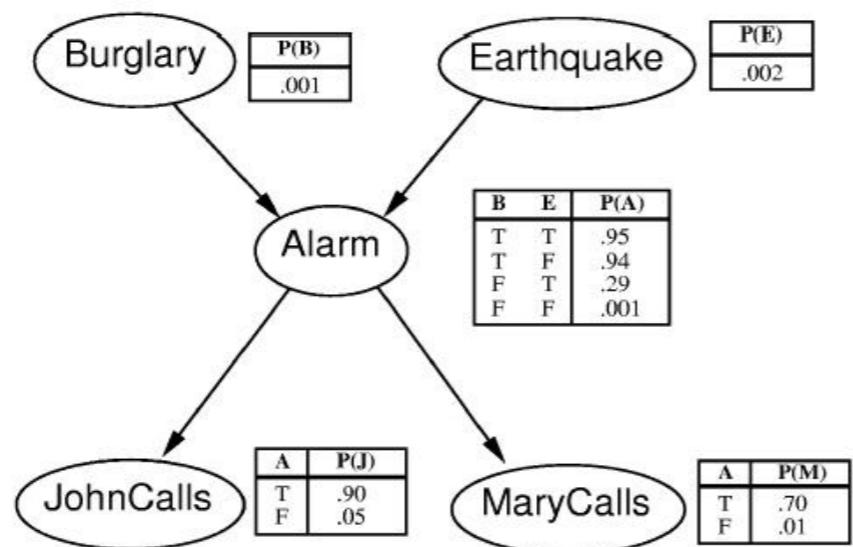
# Gibbs Sampling: An Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	F
2					
3					
4					

- Sampling  $P(M|A)$ : No need to apply Bayes Rule
- $A = F$ , so we compute the following, and sample
  - $P(M = T | A = F) \propto 0.01$
  - $P(M = F | A = F) \propto 0.99$

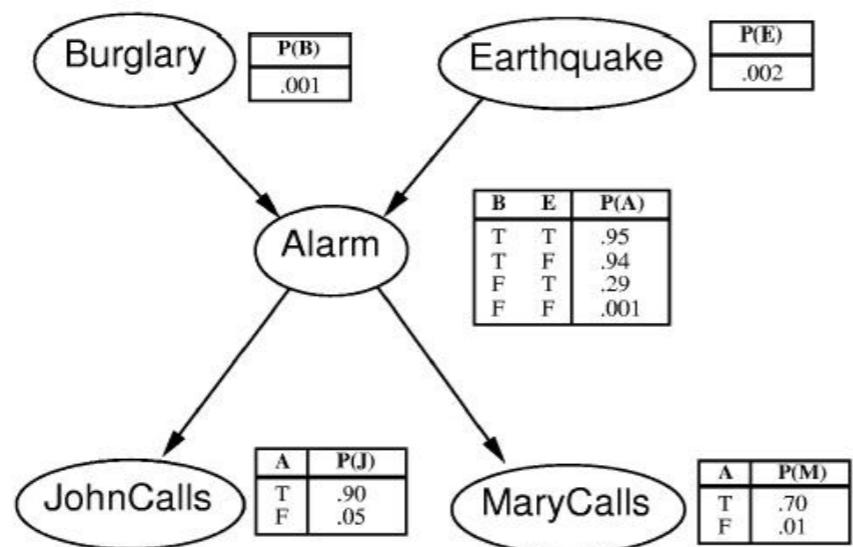
# Gibbs Sampling: An Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	F
2	F	T	T	T	T
3					
4					

- Now  $t = 2$ , and we repeat the procedure to sample new values of  $B, E, A, J, M \dots$

# Gibbs Sampling: An Example



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	F
2	F	T	T	T	T
3	T	F	T	F	T
4	T	F	T	F	F

- Now  $t = 2$ , and we repeat the procedure to sample new values of B,E,A,J,M ...
- And similarly for  $t = 3, 4$ , etc.

# Gibbs Sampling and Markov Chains

- This algorithm is an instance of a broad family of tools: MCMC
- We will study in future lecture the main properties and uses of general MCMC methods.

# Lecture 5 Objectives

- Modeling Survey Data
- Factor Analysis
  - Principal Component Analysis
  - Independent Component Analysis
- Gaussian Mixture Models
- The Expectation-Maximization (EM) algorithm

# Survey Data

- How are these statements inferred?



Politics Sports Science & Health Economics Culture

MAY 3, 2016 AT 2:45 PM

## The Mythology Of Trump's 'Working Class' Support

His voters are better off economically compared with most Americans.

By Matt Silver

Filed under [2016 Election](#)



Politics Sports Science & Health Economics Culture

MAR 14, 2016 AT 6:02 AM

## What Trump Supporters Were Doing Before Trump

By Dan Hopkins

Filed under [2016 Election](#)



JUNE 2, 2016



★ [Election 2016](#)

## More 'warmth' for Trump among GOP voters concerned by immigrants, diversity

BY BRADLEY JONES AND JOCELYN KILEY | 10 COMMENTS

# Survey Data

- We typically see these sort of data

## How the GOP candidates' supporters differed on issues

POSITION	YEAR	STANCE OF SUPPORTERS		
		TRUMP	CRUZ	RUBIO
Raise taxes on rich	2007	0.25	0.24	0.27
Pro-gay marriage	2007	0.31	0.20	0.28
Conservative ideology	2007	0.64	0.76	0.67
Pro-choice	2007	0.63	0.39	0.42
Stay in Iraq	2007	0.62	0.81	0.87
Hawk (vs. dove)	2008	0.68	0.67	0.52
No special help for blacks	2012	0.88	0.82	0.80
Obama rating	2012	0.21	0.17	0.20
Anti-Obamacare	2012	0.76	0.82	0.80
Very critical of system	2012	0.70	0.64	0.62
Pro-government spending	2012	0.20	0.15	0.15
Create pathway to citizenship	2012	0.21	0.29	0.37
Anti-Hispanic prejudice	2012	0.53	0.50	0.49
Anti-black prejudice	2012	0.58	0.54	0.55
Pro-NAFTA	2012	0.40	0.50	0.52

Stance on a position ranges from 0-1, with 0 being totally against and 1 being totally in agreement with

SOURCE: HOPKINS/MUTZ

## Share of Republican electorate with household income below \$50,000

STATE	2012	2016
Alabama	37%	41%
Florida	34	35
Georgia	24	26
Illinois	28	23
Maryland	19	19
Massachusetts	24	20
Michigan	35	37
Mississippi	36	37
New Hampshire	26	27
Ohio	32	30
Oklahoma	41	30
South Carolina	36	27
Tennessee	35	33
Vermont	37	30
Virginia	25	19
Wisconsin	32	20
Average	31	29

SOURCE: EDISON RESEARCH EXIT POLLS

## Within GOP, views of immigration, Islam, diversity strongly associated with ratings of Trump

% of Republican and Republican-leaning registered voters who rate Trump on a feeling thermometer from 0 (coldest rating) to 100 (warmest rating) ...

■ Very cold ■ Somewhat cold ■ Neutral ■ Somewhat warm ■ Very warm

All Rep/Rep-leaning voters	23	11	12	17	36
----------------------------	----	----	----	----	----

### Among those who say ...

Growing number of newcomers from other countries ...

Threatens U.S. values (77%)	16	11	11	18	42
-----------------------------	----	----	----	----	----

Strengthens U.S. society (21%)

42	13	13	15	14
----	----	----	----	----

### The Islamic religion is ...

More likely than others to encourage violence (77%)

19	12	11	18	38
----	----	----	----	----

No more likely to encourage violence (20%)

37	8	16	11	25
----	---	----	----	----

According to census, in 30 years U.S. pop. will be majority black, Latino & Asian. This is ...

Bad for the country (39%)

15	11	10	16	47
----	----	----	----	----

Good/Neither good nor bad for the country (61%)

28	12	13	18	28
----	----	----	----	----

Feeling thermometer ratings: Very cold (zero to 24), somewhat cold (25-49), neutral (50), somewhat warm (51-75), very warm (76-100).

Source: Survey conducted April 5-May 2, 2016

PEW RESEARCH CENTER

# Survey Data: other prime example

- Medical surveys
  - Children's Depression Inventory [3]
    - 27 items scored 0,1,2 assessing aspects of depressive symptoms for children and adolescents
    - 1 total scale
      - sum of the 27 items after reverse coding 13 of them
      - higher scores indicate higher depressive symptom levels
    - 5 subscales measuring different aspects of depressive symptoms
      - negative mood, interpretation problems, ineffectiveness, anhedonia, and negative self-esteem
      - the total scale equals the sum of the subscales
    - total scale used in practice rather than subscales
- Financial Markets
- EEG recordings

(credit: G. Knafl, ohsu)

# Factor Analysis for Survey Data

- **Goal:** extract interpretable, summary information out of a series of correlated survey responses.
- Factor Analysis refers to a series of statistical techniques to achieve that.
- That is, given answers  $x_1, \dots, x_L$  to  $L$  questions, infer latent variables (=factors) that explain the underlying phenomena under study.

# Principal Component Analysis

- We start with the simplest setting: suppose that

$$X_l = \sum_{j=1}^J \alpha_{j,l} Y_j , \quad l = 1 \dots L .$$

with  $Y_1 \dots J$  iid and  $X_1 \dots L$  jointly Gaussian.

- Q: How to discover the 'latent' factors  $Y$ ?

# Principal Component Analysis

- We start with the simplest setting: suppose that

$$X_l = \sum_{j=1}^J \alpha_{j,l} Y_j , \quad l = 1 \dots L .$$

with  $Y_1 \dots J$  iid and  $X_1 \dots L$  jointly Gaussian.

- Q: How to discover the 'latent' factors  $\mathbf{Y}$ ?
- Observation 1:

$$\mathbf{X} \text{ Gaussian} \Rightarrow \mathbf{Y} = A\mathbf{X} \text{ also Gaussian.}$$

# Principal Component Analysis

- We start with the simplest setting: suppose that

$$X_l = \sum_{j=1}^J \alpha_{j,l} Y_j , \quad l = 1 \dots L .$$

with  $Y_1 \dots J$  iid and  $X_1 \dots L$  jointly Gaussian.

- Q: How to discover the 'latent' factors  $\mathbf{Y}$ ?
- Observation 1:

$\mathbf{X}$  Gaussian  $\Rightarrow \mathbf{Y} = A\mathbf{X}$  also Gaussian.

- Observation 2:

If  $\mathbf{Y} = (Y_1, \dots, Y_J)$  is jointly Gaussian,  
 $Y_i, Y_j$  independent  $\Leftrightarrow Y_i, Y_j$  decorrelated.

# Principal Component Analysis

- We define  $\mu_X = \mathbb{E}(X)$ ,  $\Sigma_X = \mathbb{E}\{(X - \mu_X)(X - \mu_X)^T\}$ .
- **Reminder:** Let  $\mathbf{Y} = A\mathbf{X} + b$ . Then
  - $\mu_Y = A\mu_X + b$ .
  - $\Sigma_Y = A\Sigma_X A^T$ .

# Principal Component Analysis

- We define  $\mu_X = \mathbb{E}(X)$ ,  $\Sigma_X = \mathbb{E}\{(X - \mu_X)(X - \mu_X)^T\}$ .

- **Reminder:** Let  $\mathbf{Y} = A\mathbf{X} + b$ . Then

- $\mu_Y = A\mu_X + b$ .
- $\Sigma_Y = A\Sigma_X A^T$ .

- In our previous model, if  $A_{l,j} = \alpha_{l,j}$ , we have  $\mathbf{X} = A\mathbf{Y}$ .
- Hence  $\Sigma_X = A\Sigma_Y A^T$
- Q: How to find  $A$  such that  $A^{-1}\Sigma_X A^{-T}$  defines an uncorrelated random vector?

# Principal Component Analysis

**Reminder:** [Real Spectral Theorem]: If  $A \in \mathbb{R}^{n \times n}$  is symmetric and positive semidefinite ( $A \succ 0$ ), then  $A$  diagonalizes in a real orthonormal basis with non-negative eigenvalues.

# Principal Component Analysis

**Reminder:** [Real Spectral Theorem]: If  $A \in \mathbb{R}^{n \times n}$  is symmetric and positive semidefinite ( $A \succ 0$ ), then  $A$  diagonalizes in a real orthonormal basis with non-negative eigenvalues.

**Fact:** The covariance operator  $\Sigma_X = \mathbb{E}((X - \mu_X)(X - \mu_X)^T)$  is **symmetric** and **positive semidefinite**.

# Principal Component Analysis

**Reminder:** [Real Spectral Theorem]: If  $A \in \mathbb{R}^{n \times n}$  is symmetric and positive semidefinite ( $A \succcurlyeq 0$ ), then  $A$  diagonalizes in a real orthonormal basis with non-negative eigenvalues.

**Fact:** The covariance operator  $\Sigma_X = \mathbb{E}((X - \mu_X)(X - \mu_X)^T)$  is **symmetric** and **positive semidefinite**.

Therefore,  $\Sigma_X$  admits a real eigenbasis:

$$\Sigma_X = U\Lambda U^T, \quad \Lambda = \text{diag}(\lambda_i) \succeq 0.$$

Thus  $U^T \Sigma_X U = \Lambda$ .

# Principal Component Analysis

**Reminder:** [Real Spectral Theorem]: If  $A \in \mathbb{R}^{n \times n}$  is symmetric and positive semidefinite ( $A \succ 0$ ), then  $A$  diagonalizes in a real orthonormal basis with non-negative eigenvalues.

**Fact:** The covariance operator  $\Sigma_X = \mathbb{E}((X - \mu_X)(X - \mu_X)^T)$  is **symmetric** and **positive semidefinite**.

Therefore,  $\Sigma_X$  admits a real eigenbasis:

$$\Sigma_X = U\Lambda U^T, \quad \Lambda = \text{diag}(\lambda_i) \succeq 0.$$

Thus  $U^T \Sigma_X U = \Lambda$ .

Moreover, we can write  $\Lambda = S \cdot S$ , with  $s_{i,i} = \sqrt{\lambda_i}$ .  
If  $\min_i \lambda_i > 0$ , it results that  $\tilde{U} = US^{-1}$  satisfies  $\tilde{U}^T \Sigma_X \tilde{U} = 1$ .

# Principal Component Analysis

- We have just shown that  $\mathbf{Y} = \tilde{U}(\mathbf{X} - \mu_{\mathbf{X}})$  provides variables that are uncorrelated and jointly Gaussian, thus independent.

# Principal Component Analysis

- We have just shown that  $\mathbf{Y} = \tilde{U}(\mathbf{X} - \mu_{\mathbf{X}})$  provides variables that are uncorrelated and jointly Gaussian, thus independent.
- Remarks
  - The decomposition is not unique: Any orthogonal transformation of  $\mathbf{Y}$  also satisfies the same property.
  - PCA provides linear compression: if  $J < \text{rank}(\Sigma_{\mathbf{X}})$ , what is the best linear approximation of  $\mathbf{X}$  with  $J$  independent components?

$$\min_{A \in \mathbb{R}^{L \times J}} \mathbb{E}(\|X - AX\|^2) .$$

$A = \{ \text{eigenvectors of } \Sigma_{\mathbf{X}} \text{ corresponding to } J \text{ largest eigenvalues.}\}$   
(again,  $A$  is determined up to an orthogonal transformation)

# Estimating the Principal Components

- So far, we have seen how to extract information from the covariance of  $X$ .

# Estimating the Principal Components

- So far, we have seen how to extract information from the covariance of  $X$ .
- In practice, we will observe  $x_1, \dots, x_N$  iid samples of  $X$
- Empirical Covariance:

$$\hat{\Sigma}_N = \frac{1}{N} \sum_{n \leq N} (x_n - \hat{\mu})(x_n - \hat{\mu})^T . \quad \in \mathbb{R}^{L \times L} .$$

# Estimating the Principal Components

- So far, we have seen how to extract information from the covariance of  $X$ .
- In practice, we will observe  $x_1, \dots, x_N$  iid samples of  $X$
- Empirical Covariance:

$$\hat{\Sigma}_N = \frac{1}{N} \sum_{n \leq N} (x_n - \hat{\mu})(x_n - \hat{\mu})^T . \quad \in \mathbb{R}^{L \times L} .$$

- $\hat{\Sigma}_N$  is symmetric, positive definite. (why?)
- Estimated Principal Components:

$$\hat{\Sigma}_N = \hat{U} \hat{\Lambda} \hat{U}^T .$$

# Estimating the Principal Components

- Q: How good are these estimates?
  - i.e. for a desired accuracy  $\epsilon$ , how many samples  $N=N(L)$  are required?

# Estimating the Principal Components

- Q: How good are these estimates?
  - i.e. for a desired accuracy  $\epsilon$ , how many samples  $N=N(L)$  are required?
- **Theorem [Vershynin]:** For distributions with bounded order- $q$  moment, the empirical covariance satisfies

$$\|\widehat{\Sigma}_N - \Sigma\| \lesssim O(\log \log N)^2 \left(\frac{N}{L}\right)^{1/2-2/q}.$$

# Estimating the Principal Components

- Q: How good are these estimates?
  - i.e. for a desired accuracy  $\epsilon$ , how many samples  $N=N(L)$  are required?
- **Theorem [Vershynin]:** For distributions with bounded order- $q$  moment, the empirical covariance satisfies

$$\|\hat{\Sigma}_N - \Sigma\| \lesssim O(\log \log N)^2 \left(\frac{N}{L}\right)^{1/2-2/q}.$$

It results that for a desired approximation  $\|\hat{\Sigma}_N - \Sigma\| \leq \epsilon$  we need  $O((\log \log L)^\alpha L) \approx O(L)$  samples.

# Estimating the Principal Components

- Q: How good are these estimates?
  - i.e. for a desired accuracy  $\epsilon$ , how many samples  $N=N(L)$  are required?
- **Theorem [Vershynin]:** For distributions with bounded order- $q$  moment, the empirical covariance satisfies
$$\|\hat{\Sigma}_N - \Sigma\| \lesssim O(\log \log N)^2 \left(\frac{N}{L}\right)^{1/2-2/q}.$$
- It results that for a desired approximation  $\|\hat{\Sigma}_N - \Sigma\| \leq \epsilon$  we need  $O((\log \log L)^\alpha L) \approx O(L)$  samples.
- **Very Important Consequence: PCA does not suffer from the curse of dimensionality!**

# Estimating Principal Components

- Q: What is the computational complexity of computing PCA?
  - Naïve estimation of covariance:  $O(NL^2)$
  - Diagonalizing covariance:  $O(L^3)$
  - So, in big data applications we are doomed.

# Estimating Principal Components

- Q: What is the computational complexity of computing PCA?
  - Naïve estimation of covariance:  $O(NL^2)$
  - Diagonalizing covariance:  $O(L^3)$
  - So, in big data applications we are doomed.
- Reminder: Given a data matrix, the principal components are directly given by

The **Singular Value Decomposition** (SVD) of  $B \in \mathbb{R}^{n \times p}$  is defined as  $B = U\Lambda V^T$ , with

$$U \in \mathbb{R}^{n \times n}, V \in \mathbb{R}^{p \times p}, \Lambda = \text{diag}(\lambda_1, \dots, \lambda_{\min(n,p)}),$$

$$UU^T = \mathbf{1}, V^TV = \mathbf{1}.$$

Computing the first  $p$  principal components costs  $O(pNL)$ .

# Estimating Principal Components

- Q: What is the computational complexity of computing PCA?
  - Naïve estimation of covariance:  $O(NL^2)$
  - Diagonalizing covariance:  $O(L^3)$
  - So, in big data applications we are doomed.
- Reminder: Given a data matrix, the principal components are directly given by

The **Singular Value Decomposition** (SVD) of  $B \in \mathbb{R}^{n \times p}$  is defined as  $B = U\Lambda V^T$ , with

$$U \in \mathbb{R}^{n \times n}, V \in \mathbb{R}^{p \times p}, \Lambda = \text{diag}(\lambda_1, \dots, \lambda_{\min(n,p)}),$$

$$UU^T = \mathbf{1}, V^TV = \mathbf{1}.$$

Computing the first  $p$  principal components costs  $O(pNL)$ .

- Alternatives?

## Interlude: Randomized PCA

- Suppose we suspect that the data matrix  $\mathbf{X} \in \mathbb{R}^{N \times L}$  has rank  $p \ll \min(N, L)$  .
- Q: Can we leverage that prior into a faster algorithm?

## Interlude: Randomized PCA

- Suppose we suspect that the data matrix  $\mathbf{X} \in \mathbb{R}^{N \times L}$  has rank  $p \ll \min(N, L)$ .
- Q: Can we leverage that prior into a faster algorithm?
- **Idea:** Break the estimation into two steps:
  1. Compute an approximate basis for the range of  $\mathbf{X}$ :
$$\mathbf{X} \approx \mathbf{Q}\mathbf{Q}^T\mathbf{X}, \text{ with } \mathbf{Q} \in \mathbb{R}^{N \times n}, n \ll N, \mathbf{Q}^T\mathbf{Q} = \mathbf{I}.$$
  2. Form  $\mathbf{B} = \mathbf{Q}^T\mathbf{X} \in \mathbb{R}^{n \times L}$  and compute its SVD:  $\mathbf{B} = \tilde{\mathbf{U}}\Lambda\mathbf{V}^T$ .
  3. Set  $\mathbf{U} = \mathbf{Q}\tilde{\mathbf{U}}$ .
- How to solve stage 1?

## Interlude: Randomized PCA

- Suppose we suspect that the data matrix  $\mathbf{X} \in \mathbb{R}^{N \times L}$  has rank  $p \ll \min(N, L)$  .
- Q: Can we leverage that prior into a faster algorithm?
- **Idea:** Break the estimation into two steps:
  1. Compute an approximate basis for the range of  $\mathbf{X}$ :
$$\mathbf{X} \approx \mathbf{Q}\mathbf{Q}^T\mathbf{X}, \text{ with } \mathbf{Q} \in \mathbb{R}^{N \times n}, n \ll N, \mathbf{Q}^T\mathbf{Q} = \mathbf{I}.$$
  2. Form  $\mathbf{B} = \mathbf{Q}^T\mathbf{X} \in \mathbb{R}^{n \times L}$  and compute its SVD:  $\mathbf{B} = \tilde{\mathbf{U}}\Lambda\mathbf{V}^T$  .
  3. Set  $\mathbf{U} = \mathbf{Q}\tilde{\mathbf{U}}$  .
- How to solve stage 1? **Randomize!!**

## Interlude: Randomized PCA

- If we target a rank- $k$  approximation, we simply draw a random matrix  $\Omega \in \mathbb{R}^{L \times (k+p)}$ , form the matrix  $Y = X\Omega \in \mathbb{R}^{N \times (k+p)}$ , and perform a SVD of  $Y$ :  $Y = V\Lambda Q^T$ .

## Interlude: Randomized PCA

- If we target a rank- $k$  approximation, we simply draw a random matrix  $\Omega \in \mathbb{R}^{L \times (k+p)}$ , form the matrix  $Y = X\Omega \in \mathbb{R}^{N \times (k+p)}$ , and perform a SVD of  $Y$ :  $Y = V\Lambda Q^T$ .
- Strong guarantees from concentration of measure:

**Theorem:** [Halko, Martinsson, Tropp] Given data matrix  $X \in \mathbb{R}^{N \times L}$  and  $\Omega \in \mathbb{R}^{L \times (k+p)}$  drawn from iid standard Gaussian, the resulting  $Q$  satisfies

$$\|X - QQ^T X\| \leq \left(1 + C\sqrt{(k+p) \min(N, L)}\right) \lambda_{k+1},$$

whp, where  $\lambda_{k+1}$  is the  $k+1$ -th singular value of  $X$ .

## Interlude: Randomized PCA

- If we target a rank- $k$  approximation, we simply draw a random matrix  $\Omega \in \mathbb{R}^{L \times (k+p)}$ , form the matrix  $Y = X\Omega \in \mathbb{R}^{N \times (k+p)}$ , and perform a SVD of  $Y$ :  $Y = V\Lambda Q^T$ .
- Strong guarantees from concentration of measure:

**Theorem:** [Halko, Martinsson, Tropp] Given data matrix  $X \in \mathbb{R}^{N \times L}$  and  $\Omega \in \mathbb{R}^{L \times (k+p)}$  drawn from iid standard Gaussian, the resulting  $Q$  satisfies

$$\|X - QQ^T X\| \leq \left(1 + C\sqrt{(k+p) \min(N, L)}\right) \lambda_{k+1},$$

whp, where  $\lambda_{k+1}$  is the  $k+1$ -th singular value of  $X$ .

- Resulting computational gains:

from  $O(NLk)$  to  $O(NL \log(k))$  for  $k$  ppal components.

# Factor Analysis

- We have seen that PCA can be interpreted as a linear latent model:

$$X_i = \sum_j \alpha_{i,j} Y_j + \mu_i , \quad (i = 1, \dots, L) ,$$

with  $Y_j$  uncorrelated, unit variance.

# Factor Analysis

- We have seen that PCA can be interpreted as a linear latent model:

$$X_i = \sum_j \alpha_{i,j} Y_j + \mu_i , \quad (i = 1, \dots, L) ,$$

with  $Y_j$  uncorrelated, unit variance.

- Lack of unicity: given an  $L \times L$  orthogonal matrix  $R$ , we also have

$$X = \mu + AR^T RY = \mu + \tilde{A}\tilde{Y} ,$$

where  $\tilde{Y}_j$  are also uncorrelated and unit variance.

# Factor Analysis

- We have seen that PCA can be interpreted as a linear latent model:

$$X_i = \sum_j \alpha_{i,j} Y_j + \mu_i , \quad (i = 1, \dots, L) ,$$

with  $Y_j$  uncorrelated, unit variance.

- Lack of unicity: given an  $L \times L$  orthogonal matrix  $R$ , we also have

$$X = \mu + AR^T RY = \mu + \tilde{A}\tilde{Y} ,$$

where  $\tilde{Y}_j$  are also uncorrelated and unit variance.

- Also, an underlying assumption is that data has *low-rank*, i.e. covariance directly reveals dependencies in data.

# Factor Analysis

- An alternative to PCA is *Factor Analysis*. Suppose a generative model of the form

$$X_i = \sum_{j \leq J} \alpha_{i,j} Y_j + \mu_i + \epsilon_i , \quad (i = 1, \dots, L) ,$$

with  $J < L$  and  $\epsilon_i$  uncorrelated, zero-mean.

# Factor Analysis

- An alternative to PCA is *Factor Analysis*. Suppose a generative model of the form

$$X_i = \sum_{j \leq J} \alpha_{i,j} Y_j + \mu_i + \epsilon_i , \quad (i = 1, \dots, L) ,$$

with  $J < L$  and  $\epsilon_i$  uncorrelated, zero-mean.

- Interpretation:
  - Latent variables  $Y_j$  are common factors of variability.
  - Latent variables  $\epsilon_i$  explain the remaining individual variability, uncorrelated from the rest.
- Example:
  - Factor analysis on the topics of your final project.

# Factor Analysis

- Gaussian joint likelihood model:

$$X \sim \mathcal{N}(\mu, AA^T + \text{diag}(\beta))$$

- with  $\beta_i = \text{Var}(\epsilon_i)$  .
- Parameter Estimation? The covariance is a sufficient statistic:

$$\Sigma_X = AA^T + \text{diag}(\beta) .$$

↑  
low rank

- SVD is still useful, but does not automatically yield the solution.
- We will soon see an alternative estimation algorithm.

# Independent Component Analysis

- We have seen that asking latent variables only to be decorrelated leads to lack of unicity (any orthogonal transformation is also decorrelated).

# Independent Component Analysis

- We have seen that asking latent variables only to be decorrelated leads to lack of unicity (any orthogonal transformation is also decorrelated).
- What happens if we tighten our requirements, e.g. asking that

$Y_i$  and  $Y_j$  independent

# Independent Component Analysis

- We have seen that asking latent variables only to be decorrelated leads to lack of unicity (any orthogonal transformation is also decorrelated).
- What happens if we tighten our requirements, e.g. asking that

$Y_i$  and  $Y_j$  independent

- We have seen that in the Gaussian case, this does not alleviate the problem. (why?)

# Independent Component Analysis

- We have seen that asking latent variables only to be decorrelated leads to lack of unicity (any orthogonal transformation is also decorrelated).
- What happens if we tighten our requirements, e.g. asking that

$Y_i$  and  $Y_j$  independent

- We have seen that in the Gaussian case, this does not alleviate the problem. (why?)
- But, as it turns out, it is the exception. The model becomes uniquely identifiable if

$Y_i$  and  $Y_j$  independent and non-Gaussian.

# Independent Component Analysis

- Without loss of generality, we can assume  $\mathbf{X}$  has trivial covariance:  $\Sigma_{\mathbf{X}} = \mathbf{1}$ .
- Since  $\Sigma_{\mathbf{Y}} = \mathbf{1}$  as well, it results that  $\mathbf{1} = A\Sigma_{\mathbf{Y}}A^T = AA^T$

# Independent Component Analysis

- Without loss of generality, we can assume  $\mathbf{X}$  has trivial covariance:  $\Sigma_{\mathbf{X}} = \mathbf{1}$ .
- Since  $\Sigma_{\mathbf{Y}} = \mathbf{1}$  as well, it results that  $\mathbf{1} = A\Sigma_{\mathbf{Y}}A^T = AA^T$
- so we are reduced to finding an orthogonal transformation such that  $\mathbf{Y} = A^T\mathbf{X}$  becomes independent and non-Gaussian.

# Independent Component Analysis

- Without loss of generality, we can assume  $\mathbf{X}$  has trivial covariance:  $\Sigma_{\mathbf{X}} = \mathbf{1}$ .
- Since  $\Sigma_{\mathbf{Y}} = \mathbf{1}$  as well, it results that  $\mathbf{1} = A\Sigma_{\mathbf{Y}}A^T = AA^T$
- so we are reduced to finding an orthogonal transformation such that  $\mathbf{Y} = A^T \mathbf{X}$  becomes independent and non-Gaussian.
- It is a form of “inverse” Central Limit Theorem method.
- Q: How to measure/estimate statistical independence?

# Entropy

- For a continuous random variable  $X$  with density  $p(x)$ , the *differential entropy* is

$$H(X) = - \int p(x) \log p(x) dx = -\mathbb{E}(\log p(X)) .$$

- it extends the notion of uncertainty/information carried by  $X$ .

# Entropy

- For a continuous random variable  $X$  with density  $p(x)$ , the *differential entropy* is

$$H(X) = - \int p(x) \log p(x) dx = -\mathbb{E}(\log p(X)) .$$

- it extends the notion of uncertainty/information carried by  $X$ .
- Entropy measures independence through the *mutual information*: Given  $X_1, \dots, X_n$ , the mutual information is

$$I(Y) = \sum_{i \leq n} H(Y_i) - H(Y) \geq 0 .$$

# Entropy

- For a continuous random variable  $X$  with density  $p(x)$ , the *differential entropy* is

$$H(X) = - \int p(x) \log p(x) dx = -\mathbb{E}(\log p(X)) .$$

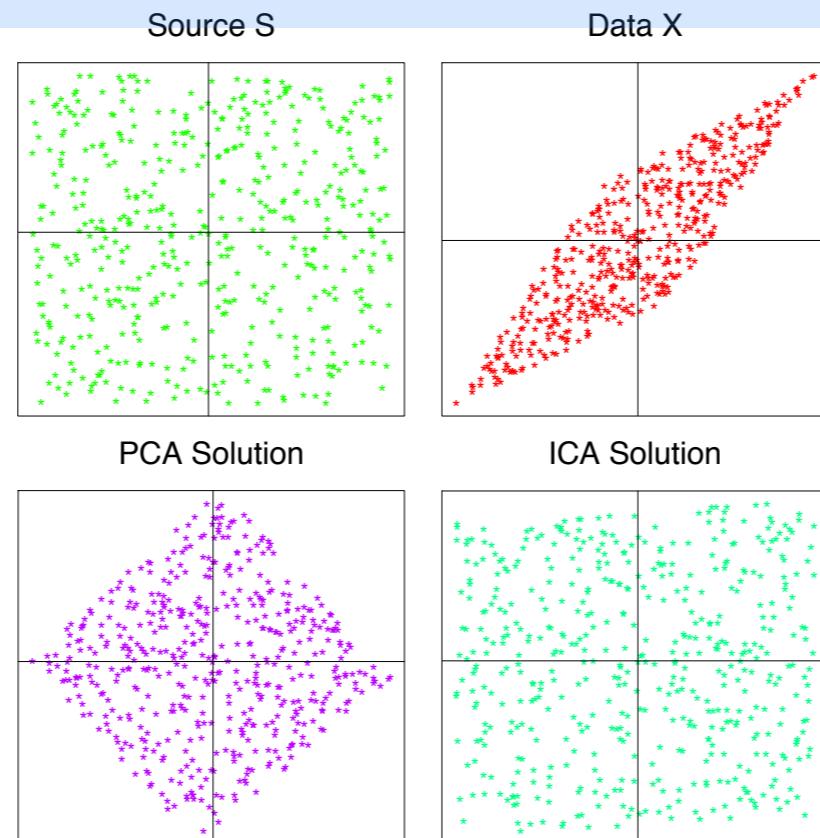
- it extends the notion of uncertainty/information carried by  $X$ .
- Entropy measures independence through the *mutual information*: Given  $X_1, \dots, X_n$ , the mutual information is
$$I(Y) = \sum_{i \leq n} H(Y_i) - H(Y) \geq 0 .$$
- **Fact:**  $I(Y) = 0$  iff  $Y_i$  and  $Y_j$  are mutually independent.
- **Fact:** If  $A$  is unitary and  $Y = A^T X$ , then  $H(Y) = H(X)$ .

# Independent Component Analysis

- So ICA attempts to solve the following problem:

$$\arg \min_{A^T A = 1} \sum_{i \leq L} H(\langle A_i, X \rangle) - H(X) = \arg \min_{A^T A = 1} \sum_{i \leq L} H(\langle A_i, X \rangle)$$

- Ex from ESLL:



**FIGURE 14.38.** Mixtures of independent uniform random variables. The upper left panel shows 500 realizations from the two independent uniform sources, the upper right panel their mixed versions. The lower two panels show the PCA and ICA solutions, respectively.

- Challenge: computing entropy requires estimating the density: exposed to curse of dimensionality!