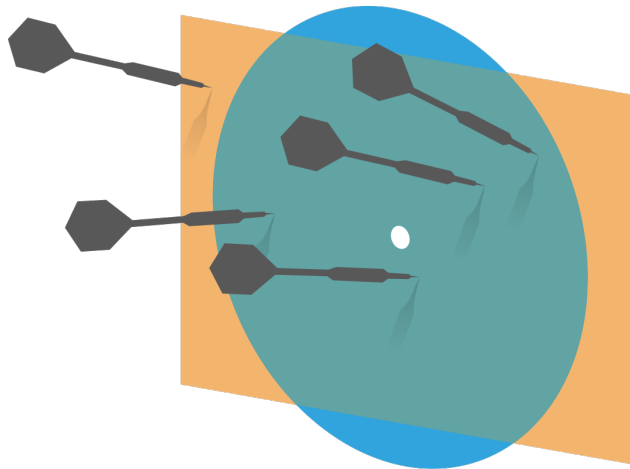


Probability and Statistics for Computer Science



$$\begin{aligned} \text{cov}(X, Y) &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY] - E[X]E[Y] \end{aligned}$$

Covariance is coming back in
matrix!

Credit: wikipedia

Bootstrap for confidence interval of other sample statistics

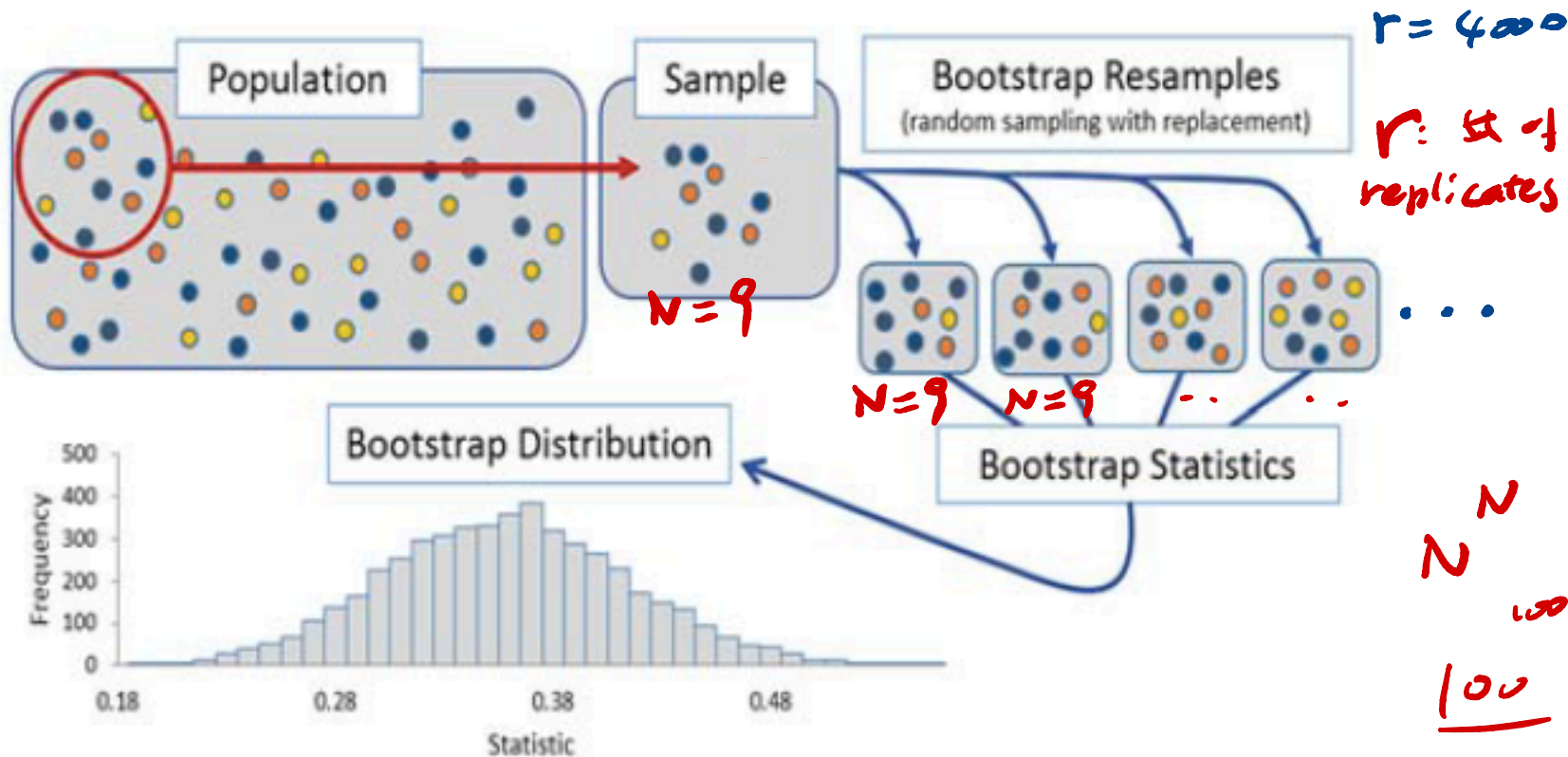


Figure 1. Summary of Bootstrapping Process

Last time

- Maximum likelihood Estimation

(MLE II) $L(\theta) = P(D|\theta)$

$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta)$

- Bayesian Inference (MAP)

$$\theta \rightarrow \text{RV}$$

$$P(\theta|D) \rightsquigarrow \text{distr.}$$

$$P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{P(D)}$$

Objective

- ✱ Review of Bayesian inference
- ✱ Visualizing high dimensional data & Summarizing data
- ✱ The covariance matrix
- ✱ Refresh of some linear algebra

Beta distribution

- ✱ A distribution is Beta distribution if it has the following

pdf:

$$P(\theta) = \begin{cases} K(\alpha, \beta) \theta^{\alpha-1} (1 - \theta)^{\beta-1} \\ 0 \quad \text{O.W.} \end{cases}$$

$$\begin{aligned} 0 \leq \theta \leq 1 \\ \alpha > 0, \beta > 0 \end{aligned}$$

$$K(\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)}$$

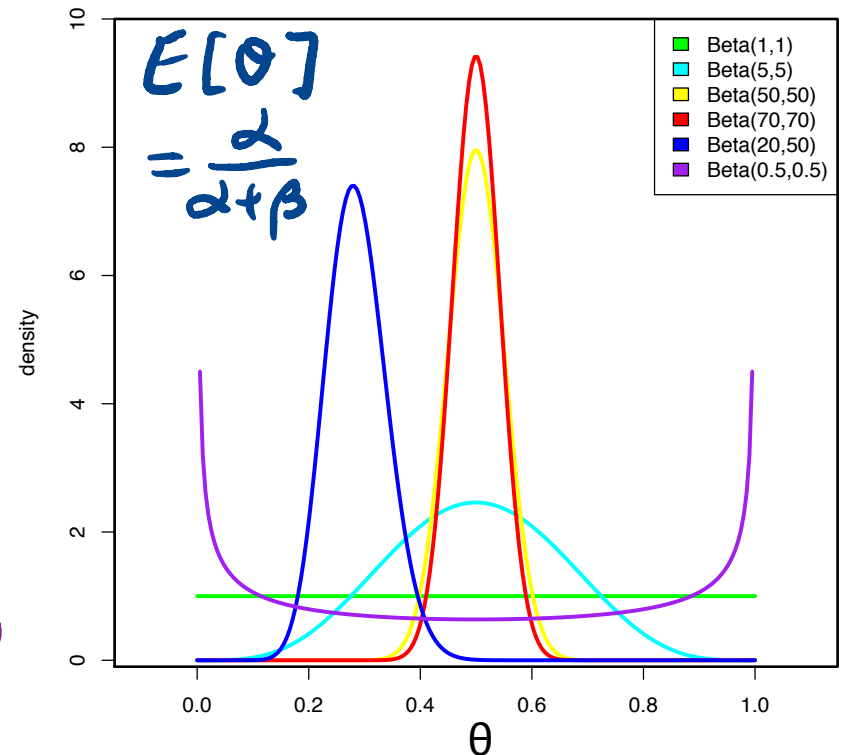
- ✱ Is an expressive family of distributions

$p(\theta | D)$

- ✱ $Beta(\alpha = 1, \beta = 1)$ is uniform

$\hat{\theta}$ that maximize $p(\theta | D)$

pdf of Beta - distribution



Beta distribution as the conjugate prior for Binomial likelihood

- ✱ The likelihood is Binomial (N, k)

$$P(D|\theta) = \binom{N}{k} \theta^k (1 - \theta)^{N-k}$$

- ✱ The Beta distribution is used as the prior

$$P(\theta) = K(\alpha, \beta) \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

- ✱ So $P(\theta|D) \propto \theta^{\alpha+k-1} (1 - \theta)^{\beta+N-k-1}$

$$\begin{aligned} \theta &\in [0, 1] \\ \hat{\alpha} &= \alpha + k - 1 \\ \hat{\beta} &= \beta + N - k \end{aligned}$$

- ✱ Then the posterior is $Beta(\alpha + k, \beta + N - k)$

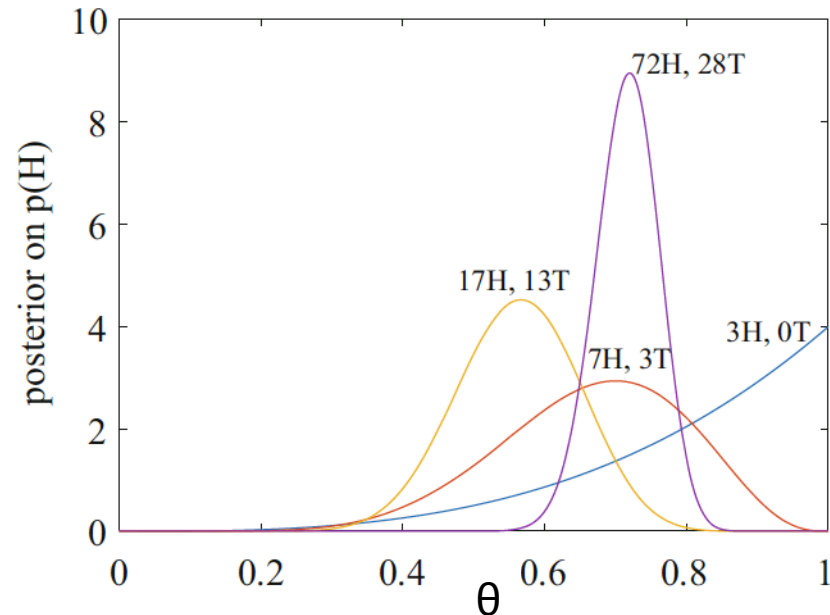
$$P(\theta|D) = K(\alpha + k, \beta + N - k) \theta^{\alpha+k-1} (1 - \theta)^{\beta+N-k-1}$$

$$\theta \in [0, 1]$$

The update of Bayesian posterior

- ✱ Since the posterior is in the same family as the conjugate prior, the posterior can be used as a new prior if more data is observed.
- ✱ Suppose we start with a uniform prior on the probability θ of heads

N	k	$\hat{\alpha}$	$\hat{\beta}$
		1	1
3	0	1	4
10	7	8	7
30	17	25	20
100	72	97	48



Maximize the Bayesian posterior (MAP)

- ✱ The posterior of the previous example is

$$P(\theta|D) = K(\alpha + k, \beta + N - k)\theta^{\alpha+k-1}(1 - \theta)^{\beta+N-k-1}$$

- ✱ Differentiating and setting to 0 gives the MAP estimate

$$\hat{\theta} = \frac{\alpha - 1 + k}{\alpha + \beta - 2 + N}$$

Table of conjugate prior for different likelihood functions

	Likelihood	Conjugate Prior	
$L(\theta)$ $=P(D \theta)$	Bernoulli Geometric Binomial	Beta distr.	$P(\theta)$
	Poisson Exponential	Gamma distr.	
	Normal with known σ^2	Normal distr.	

Conjugate prior for other likelihood functions

- ✱ If the likelihood is Bernoulli or geometric, the conjugate prior is Beta
- ✱ If the likelihood is Poisson or Exponential, the conjugate prior is Gamma
- ✱ If the likelihood is normal with known variance, the conjugate prior is normal

Which distri. is the posterior ?

If the likelihood is Geometric and we use the corresponding conjugate prior.

- A) Binomial
- B) Beta
- C) Poisson
- D) Bernoulli
- E) Normal

What are the dims of A ?

$$A \cdot X = B$$

$M \times d \cdot d \times N$ $M \times N$

X is matrix of $d \times N$

B is matrix of $M \times N$

A is matrix of ?

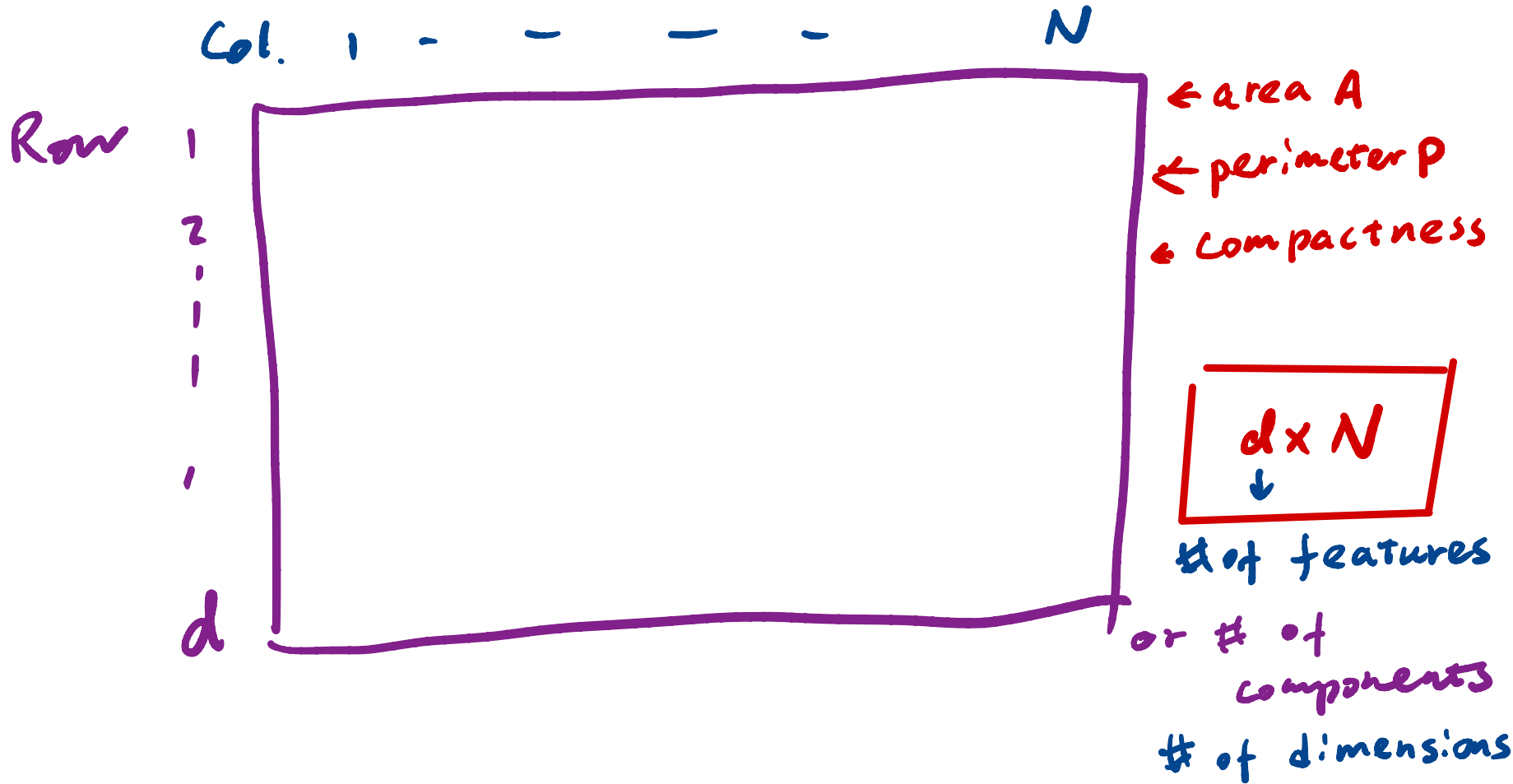
$M \times d$

A data set with high dimensions

✿ Seed data set from the UCI Machine Learning site:

	areaA	perimeterP	compactness	lengthKernel	widthKernel	asymmetry	lengthGroove	Label
1	15.26	14.84	0.871	5.763	3.312	2.221	5.22	1
2	14.88	14.57	0.8811	5.554	3.333	1.018	4.956	1
3	14.29	14.09	0.905	5.291	3.337	2.699	4.825	1
4	13.84	13.94	0.8955	5.324	3.379	2.259	4.805	1
5	16.14	14.99	0.9034	5.658	3.562	1.355	5.175	1
6	14.38	14.21	0.8951	5.386	3.312	2.462	4.956	1
7	14.69	14.49	0.8799	5.563	3.259	3.586	5.219	1
	...							

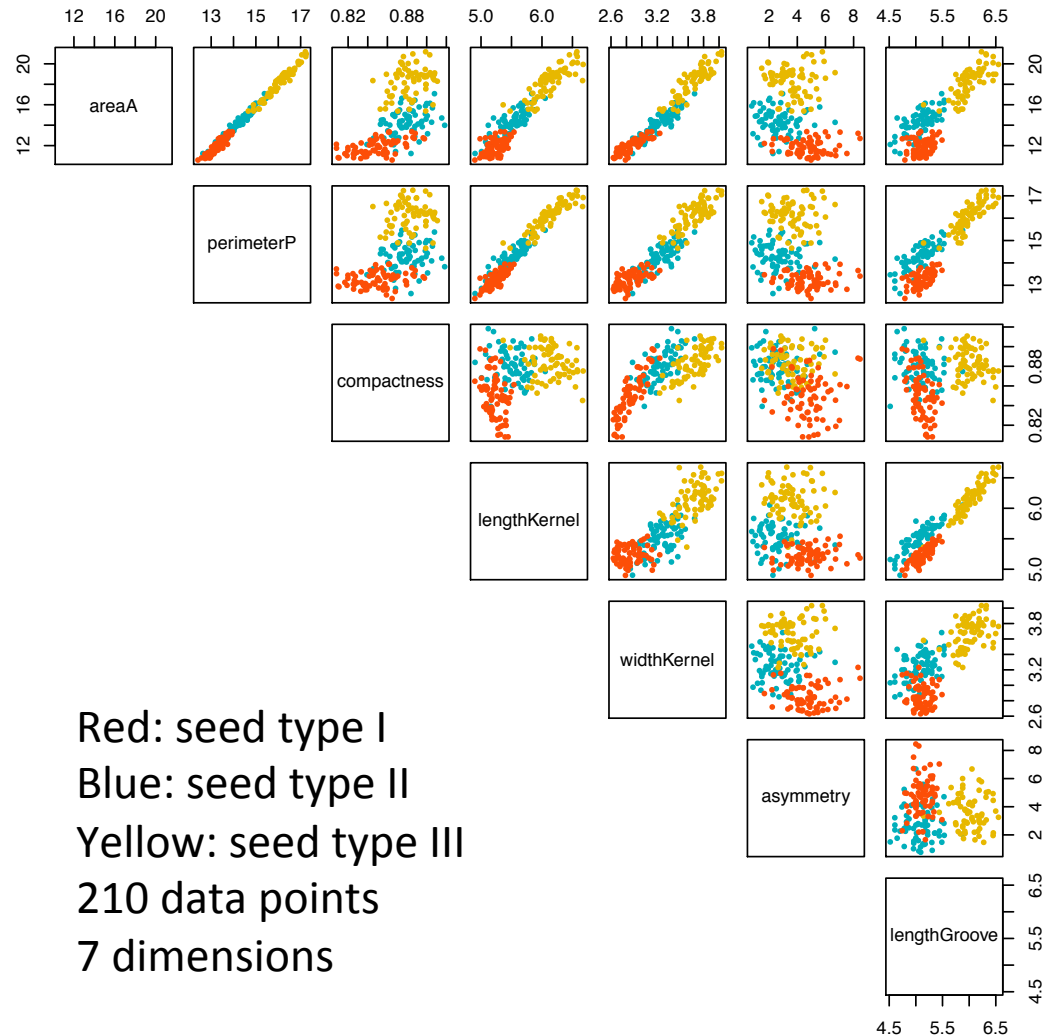
Matrix format of a dataset in the textbook



Scatterplot matrix

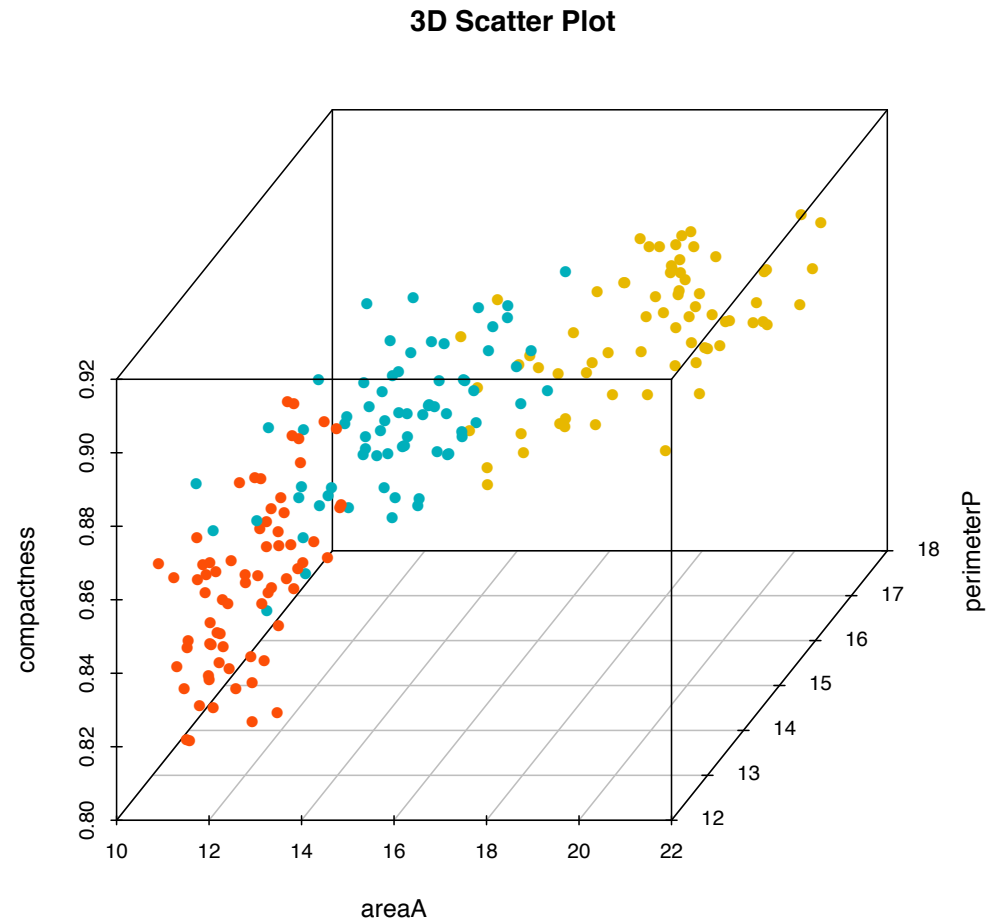
✱ Visualizing high dimensional data with scatter plot matrix

✱ Limited to small number of scatter plots



3D scatter plot

- ✱ We can also view the data set in 3 dimensions
- ✱ But it's still limited in terms of number of dimensions we can see.



Summarizing multidimensional data

- ✱ Location and spread parameters of a data set
- ✱ Notation
 - ✱ Write $\{\mathbf{x}\}$ for a dataset consisting of N data items
 - ✱ Each item x_i is a \mathbf{d} -dimensional vector; column
 - ✱ Write j th component of x_i as $x_i^{(j)}$; row
 - ✱ Matrix for the data set $\{\mathbf{x}\}$ is \mathbf{d} by \mathbf{N} dimension

Mean of a multidimensional data

- ✱ We compute the mean of $\{x\}$ by computing the mean of each component separately and stacking them to a vector

$$\text{mean of } j\text{th component} = \frac{\sum_i x_i^{(j)}}{N}$$

- ✱ We write the mean of $\{x\}$ as

$$\text{mean}(\{x\}) = \frac{\sum_i x_i}{N}$$

Example of mean of a multidimensional data set

area A
feature 1
→
:
→
:
→

	1	2	3	
	1	2	3	Mean
	-1	0	1	2
	3	7	5	0
				5

$N=3$

Mean-Centering a data matrix

Raw

1	2	3
-1	0	1
3	7	5

mean

2

0

5

Mean Centered

-1	0	1
-1	0	1
-2	2	0

Covariance

- ✱ The **covariance** of random variables X and Y is

$$\text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

- ✱ Note that

$$\text{cov}(X, X) = E[(X - E[X])^2] = \text{var}[X]$$

Correlation coefficient is normalized covariance

- ✱ The correlation coefficient is

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

- ✱ When X, Y takes on values with equal probability to generate data sets $\{(x, y)\}$, the correlation coefficient will be as seen in Chapter 2.

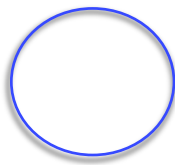
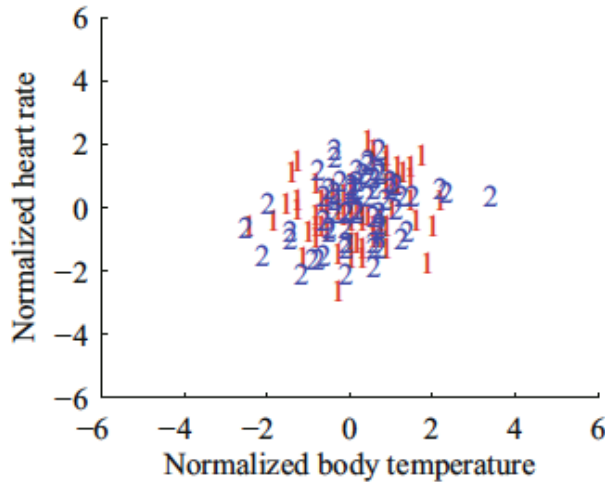
$$\text{corr}(\{x, y\}) = \frac{\sum \hat{x} \hat{y}}{N}$$

Covariance seen from scatter plots

Zero
Covariance



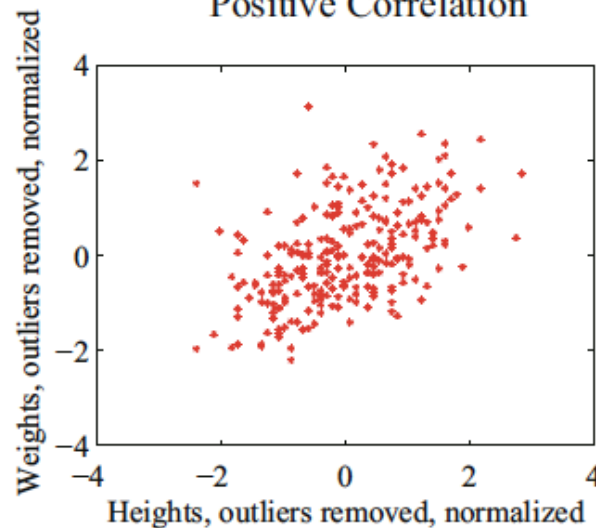
No Correlation



Positive
Covariance



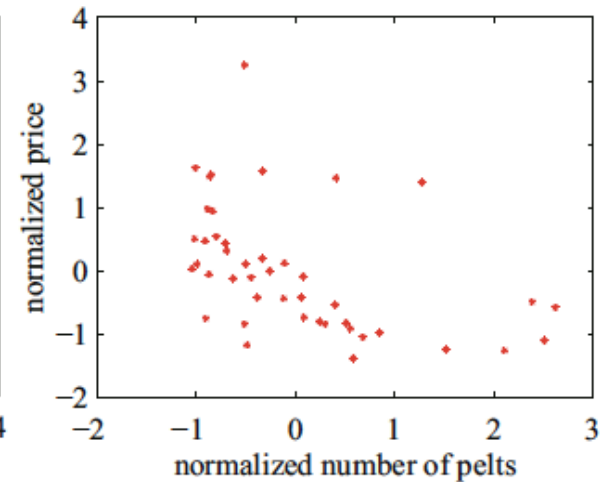
Positive Correlation



Negative
Covariance



Negative Correlation



Credit:
Prof.Forsyth

Covariance for a pair of components in a data set

- ✱ For the j th and k th components of a data set $\{x\}$

$$\text{cov}(\{x\}; j, k) = \frac{\sum_i (x_i^{(j)} - \text{mean}(\{x^{(j)}\}))(x_i^{(k)} - \text{mean}(\{x^{(k)}\}))^T}{N}$$

$$\text{cov}(x, y) = \frac{\sum \hat{x} \hat{y}}{N}$$

Covariance of a pair of components

Data set $\{\mathbf{X}\}$ 7×8

$d \times N$

$cov(\{\mathbf{x}\}; 3, 5)$

$N = 8$

	1	2	3	4	5	6	7	8
1	*	*	*	*	*	*	*	*
2	*	*	*	*	*	*	*	*
3	*	*	*	*	*	*	*	*
4	*	*	*	*	*	*	*	*
5	*	*	*	*	*	*	*	*
6	*	*	*	*	*	*	*	*
7	*	*	*	*	*	*	*	*

compact-
ness →

width-
kernel →

Take each row (component) of a pair and subtract it by the row mean, then do the inner product of the two resulting rows and divide by the number of columns

Covariance of a pair of components

Data set $\{\mathbf{x}\}$ 7×8

$cov(\{\mathbf{x}\}; 3, 5)$

	1	2	3	4	5	6	7	8
1	*	*	*	*	*	*	*	*
2	*	*	*	*	*	*	*	*
3	*	*	*	*	*	*	*	*
4	*	*	*	*	*	*	*	*
5	*	*	*	*	*	*	*	*
6	*	*	*	*	*	*	*	*
7	*	*	*	*	*	*	*	*

{

7×7

How many pairs of rows are there for which we can compute the covariance?

- A) 49
- B) 64
- C) 56

Covariance matrix

Data set $\{\mathbf{X}\}$ 7×8

$cov(\{\mathbf{x}\}; 3, 5)$

	1	2	3	4	5	6	7	8
1	*	*	*	*	*	*	*	*
2	*	*	*	*	*	*	*	*
3	*	*	*	*	*	*	*	*
4	*	*	*	*	*	*	*	*
5	*	*	*	*	*	*	*	*
6	*	*	*	*	*	*	*	*
7	*	*	*	*	*	*	*	*

Covmat($\{\mathbf{X}\}$) 7×7

	1	2	3	4	5	6	7
1	*	*	*	*	*	*	*
2	*	*	*	*	*	*	*
3	*	*	*	*	*	*	*
4	*	*	*	*	*	*	*
5	*	*	*	*	*	*	*
6	*	*	*	*	*	*	*
7	*	*	*	*	*	*	*

Properties of Covariance matrix

$$\text{cov}(\{x\}; j, j) = \text{var}(\{x^{(j)}\}) \quad \text{Covmat}(\{\mathbf{X}\}) \quad 7 \times 7$$

✱ The **diagonal** elements of the covariance matrix are just **variances** of each j th components

✱ The off diagonals are covariance between different components

	1	2	3	4	5	6	7
1	σ_1^2	*	*	*	*	*	*
2	*	σ_2^2	*	*	ν_0	*	*
3	*	*	σ_3^2	*	*	*	*
4	*	*	*	σ_4^2	*	*	*
5	*	*	*	*	σ_5^2	*	*
6	*	*	*	*	*	σ_6^2	*
7	*	*	*	*	*	*	σ_7^2

$$\text{Corr} = \frac{\text{cov}}{\sigma_x \sigma_y} = \frac{\nu_0}{\sigma_2 \sigma_5} = \frac{\nu_0}{\sqrt{\text{cov}(2,2)} \sqrt{\text{cov}(5,5)}}$$

Properties of Covariance matrix

$$\text{cov}(\{x\}; j, k) = \text{cov}(\{x\}; k, j) \quad \text{Covmat}(\{\mathbf{X}\}) \quad 7 \times 7$$

- ✱ The covariance matrix is **symmetric!**
- ✱ And it's **positive semi-definite**, that is all $\lambda_i \geq 0$
↳ eigenvalues
- ✱ Covariance matrix is diagonalizable

	1	2	3	4	5	6	7
1	*	*	*	*	*	*	*
2	*	*	*	*	*	*	*
3	*	*	*	*	*	*	*
4	*	*	*	*	*	*	*
5	*	*	*	*	*	*	*
6	*	*	*	*	*	*	*
7	*	*	*	*	*	*	*

Properties of Covariance matrix

- ✱ If we define X_c as the mean centered matrix for dataset $\{x\}$

$$Covmat(\{x\}) = \frac{X_c X_c^T}{N}$$

- ✱ The covariance matrix is a $d \times d$ matrix

Covmat($\{X\}$) 7×7

	1	2	3	4	5	6	7
1	*	*	*	*	*	*	*
2	*	*	*	*	*	*	*
3	*	*	*	*	*	*	*
4	*	*	*	*	*	*	*
5	*	*	*	*	*	*	*
6	*	*	*	*	*	*	*
7	*	*	*	*	*	*	*

$d = 7$

Example: covariance matrix of a data set

(I)

$$A_0 = \begin{bmatrix} 5 & 4 & 3 & 2 & 1 \\ -1 & 1 & 0 & 1 & -1 \end{bmatrix} \begin{matrix} x^{(1)} \\ x^{(2)} \end{matrix}$$

What are the dimensions of the covariance matrix of this data?

- A) 2 by 2
- B) 5 by 5
- C) 5 by 2
- D) 2 by 5

$$\begin{bmatrix} \text{cov}(1,1) & \text{cov}(1,2) \\ \text{cov}(2,1) & \text{cov}(2,2) \end{bmatrix}$$

Example: covariance matrix of a data set

(i) Mean centering

$$A_0 = \begin{bmatrix} 5 & 4 & 3 & 2 & 1 \\ -1 & 1 & 0 & 1 & -1 \end{bmatrix}$$

Mean $\begin{bmatrix} 3 \\ 0 \end{bmatrix}$

$$A_1 = \begin{bmatrix} 2 & 1 & 0 & -1 & -2 \\ -1 & 1 & 0 & 1 & -1 \end{bmatrix}$$

Example: covariance matrix of a data set

(I) Mean centering

$$A_0 = \begin{bmatrix} 5 & 4 & 3 & 2 & 1 \\ -1 & 1 & 0 & 1 & -1 \end{bmatrix}$$

$$A_1 = \begin{bmatrix} 2 & 1 & 0 & -1 & -2 \\ -1 & 1 & 0 & 1 & -1 \end{bmatrix}$$

(II) $A_2 = A_1 A_1^T$

Inner product of each pairs:

$$A_2 [1,1] = 10$$

$$A_2 [2,2] = 4$$

$$A_2 [1,2] = 0$$

$$\text{cov}(\{x\}) = \frac{X_c X_c^T}{N}$$

Example: covariance matrix of a data set

(I) Mean centering

$$A_0 = \begin{bmatrix} 5 & 4 & 3 & 2 & 1 \\ -1 & 1 & 0 & 1 & -1 \end{bmatrix}$$

$$A_1 = \begin{bmatrix} 2 & 1 & 0 & -1 & -2 \\ -1 & 1 & 0 & 1 & -1 \end{bmatrix}$$

(II) $A_2 = A_1 A_1^T$

Inner product of each pairs:

$$A_2 [1,1] = 10$$

$$A_2 [2,2] = 4$$

$$A_2 [1,2] = 0$$

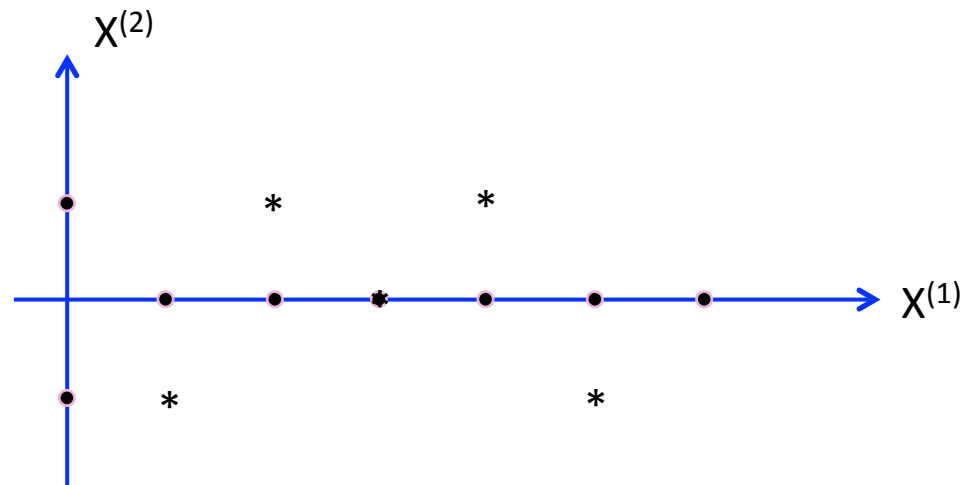
(III)

Divide the matrix with N – the number of items

$$\text{Covmat}(\{\mathbf{X}\}) = \frac{1}{N} A_2 = \frac{1}{5} \begin{bmatrix} 10 & 0 \\ 0 & 4 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 0.8 \end{bmatrix}$$

What do the data look like when $\text{Covmat}(\{\mathbf{x}\})$ is diagonal?

$$A_0 = \begin{bmatrix} 5 & 4 & 3 & 2 & 1 \\ -1 & 1 & 0 & 1 & -1 \end{bmatrix}$$



$$\text{Covmat}(\{\mathbf{X}\}) = \frac{1}{N} A_2 = \frac{1}{5} \begin{bmatrix} 10 & 0 \\ 0 & 4 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 0.8 \end{bmatrix}$$

Translation properties of mean and covariance matrix

- ✱ Translating the data set translates the mean

$$\mathit{mean}(\{x\} + c) = \mathit{mean}(\{x\}) + c$$

- ✱ Translating the data set leaves the covariance matrix unchanged

$$\mathit{Covmat}(\{x\} + c) = \mathit{Covmat}(\{x\})$$

Translation properties of covariance matrix

✱ Proof:

$$\text{covmat}(\{x\}) = \frac{X_c X_c^T}{N}$$

if we translate $\{x\}$, X_c doesn't change.

because

$$\begin{aligned} x + c - \text{mean}(\{x\} + c) \\ = x - \text{mean}(\{x\}) = X_c \end{aligned}$$

Linear transformation properties of mean and covariance matrix

- ✱ Linearly transforming the data set linearly transforms the mean

$$\text{mean}(\{A\mathbf{x}\}) = A \text{mean}(\{\mathbf{x}\})$$

- ✱ Linearly transforming the data set linearly changes the covariance matrix quadratically

$$\text{Covmat}(\{A\mathbf{x}\}) = A \text{Covmat}(\{\mathbf{x}\}) A^T$$

$$\text{var}(k\{x\}) = k^2 \text{var}(\{x\})$$

Proof of linear transformation of covariance matrix

$$\text{covmat}(\{x\}) = \frac{x_c x_c^T}{N}$$

$$\text{covmat}(\{Ax\}) = \frac{(Ax)_c (Ax)_c^T}{N}$$

$$= \frac{Ax_c (Ax_c)^T}{N}$$

$$= \frac{Ax_c \cdot x_c^T A^T}{N}$$

$$= A \cdot \text{covmat}(\{x\}) A^T$$

* suppose $x = x_c$
data is centered.

(1) $\therefore Ax = Ax_c$

(2) if x is centered
 Ax is centered

$\therefore (Ax_c)_c = Ax_c$

$$(M_a \cdot M_b)^T \\ = M_b^T \cdot M_a^T$$

Assignments

- ✱ Read Chapter 10 of the textbook
- ✱ Finish Week9 module including the quiz.
- ✱ Next time: PCA

Additional References

- ✿ Robert V. Hogg, Elliot A. Tanis and Dale L. Zimmerman. “Probability and Statistical Inference”
- ✿ Morris H. Degroot and Mark J. Schervish
"Probability and Statistics"

See you next time

*See
You!*

