# PRINCIPAL COMPONENTS ANALYSIS

**Reference:** Johnson & Wichern (2007) *Applied Multivariate Statistical Analysis* Chapter 8.

# 1 Population Principal Components

The principal components (pc) are normalized linear combinations of the correlated random variables. These linear combinations provide a new set of axes in the direction of maximum variability.

This analysis can be used for

- data reduction and

- transformation of a set of dependent variables into a new set of independent variables.

Let $\boldsymbol{X}^T = (X_1, X_2, \ldots, X_p)$ be a random vector (not necessarily normal) with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.

The first principal component, $Y_1$ is defined as the normalized linear combination of $X_1, X_2, \ldots, X_p$, with maximum variance. That is,

$$Y_1 = a_{11}X_1 + a_{12}X_2 + \ldots + a_{1p}X_p$$

with $a_{11}^2 + a_{12}^2 + \ldots + a_{1p}^2 = 1$ and $\mathbf{Var}(Y_1)$ is higher than the variance of any other linear combination of $X_1, X_2, \ldots, X_p$.

Similarly, the second principal component, $Y_2$ is a normalized linear combination with second highest variance and independent of $Y_1$.

In general, $j^{th}$ principal component $Y_j$ is given by,

$$Y_j = a_{j1}X_1 + a_{j2}X_2 + \ldots + a_{jp}X_p = \boldsymbol{a}_j^T \boldsymbol{X} \qquad \text{for } j = 1, 2, \ldots, p$$

where $\boldsymbol{a}_j^T = (a_{j1}, a_{j2}, \ldots, a_{jp})$, $\boldsymbol{a}_j^T \boldsymbol{a}_j = 1$, $\textbf{Var}(Y_j) = \boldsymbol{a}_j^T \boldsymbol{\Sigma} \boldsymbol{a}_j$, $\textbf{Cov}(Y_i, Y_j) = 0$ for all $i \neq j$ and

$$\textbf{Var}(Y_1) \geq \textbf{Var}(Y_2) \geq \ldots \geq \textbf{Var}(Y_p).$$

**Theorem:**

Suppose $(\lambda_1, \boldsymbol{e}_1)$, $(\lambda_2, \boldsymbol{e}_2), \ldots,$ $(\lambda_p, \boldsymbol{e}_p)$ are the eigenvalue-eigenvector pairs of $\boldsymbol{\Sigma}$ such that $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_p$, then, $j^{th}$ pc is defined by

$$Y_j = \boldsymbol{e}_j^T \boldsymbol{X} \qquad \text{for } j = 1, 2, \ldots, p.$$

Then $\textbf{Var}(Y_j) = \boldsymbol{e}_j^T \boldsymbol{\Sigma} \boldsymbol{e}_j = \lambda_j$ and $\textbf{Cov}(Y_i, Y_j) = \boldsymbol{e}_i^T \boldsymbol{\Sigma} \boldsymbol{e}_j = 0$.

**Note:**

(a) Two (or more) $\lambda_j$'s may be equal in some cases, however, then there exist two (or more) eigenvectors for these $\lambda_j$'s hence $p$ different pc exit for any given $\boldsymbol{\Sigma}$.

(b) The total variance,

$$T = \sum_{i=1}^{n} \textbf{Var}(X_i) = \sum_{i=1}^{n} \sigma_{ii} = tr(\boldsymbol{\Sigma}) = \lambda_1 + \lambda_2 + \ldots + \lambda_p,$$

then the proportion of variance due to $j^{\text{th}}$ pc,

$$p_j = \lambda_j / T$$

In most cases, the first few pc's contain large percentage of the total variance and therefore, the other pc's are insignificant.

(c) Let $\boldsymbol{e}_j^T = (e_{j1}, e_{j2}, \ldots, e_{jp})$, then the magnitude of $e_{kj}$ measures the importance of $X_k$ to $Y_j$ irrespective of other $X$-variables. The correlation coefficient between $Y_j$ and $X_k$ is given by

$$\rho_{Y_j, X_k} = \frac{e_{jk}\sqrt{\lambda_j}}{\sqrt{\sigma_{kk}}} \quad j, k = 1, 2, \ldots, p$$

where $Y_j = e_{j1}X_1 + e_{j2}X_2 + \ldots + e_{jk}X_k + \ldots + e_{jp}X_p$.

**Example:**

Let $\boldsymbol{X}^T = (X_1, X_2, X_3, X_4)$ be a random vector and

$$\mathbf{Cov}(\boldsymbol{X}) = \boldsymbol{\Sigma} = \begin{pmatrix} 9 & 1 & 2 & 3 \\ 1 & 9 & 3 & 2 \\ 2 & 3 & 9 & 1 \\ 3 & 2 & 1 & 9 \end{pmatrix}.$$

(a) Prove that $(15, -0.5(1, 1, 1, 1)^T$, $(9, 0.5(1, -1, -1, 1)^T$, $(7, \ 0.5(1, -1, 1, -1)^T$ and $(5, 0.5(-1, -1, 1, 1)^T$ are the eigenvalue-vector pairs of $\boldsymbol{\Sigma}$ .

(b) Obtain the principal components and their variances.

(c) Discuss the significance of these pc's.

## 2 Principal Components Using Correlation Matrix

Let $\boldsymbol{\rho}$ be correlation matrix of the random vector $\boldsymbol{X}$ and $Z_i$ be the standardized $X_i$, that is

$$Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}} \qquad \text{for } i =, 1, 2, \ldots, p$$

Then $\boldsymbol{\rho} = \mathbf{Cov}(\boldsymbol{Z})$, where $\boldsymbol{Z}^T = (Z_1, Z_2, \ldots, Z_p)$.

Let $\delta_1, \delta_2, \ldots, \delta_p$ be the eigenvalues and $\boldsymbol{\omega}_1, \boldsymbol{\omega}_2, \ldots, \boldsymbol{\omega}_p$ be the corresponding eigenvectors of $\boldsymbol{\rho}$. Then, the pc's of $\boldsymbol{Z}$ is given by

$$
\begin{aligned}
U_j &= \boldsymbol{\omega}_j^T \boldsymbol{Z} = \omega_{j1} Z_1 + \omega_{j2} Z_2 + \ldots + \omega_{jp} Z_p \\
&= \sum_i \omega_{ji} Z_i = \sum_i \omega_{ji} \left( \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}} \right) \qquad \text{for } j = 1, 2, \ldots, p.
\end{aligned}
$$

**Note:**

(a) $\mathbf{Var}(U_j) = \delta_j$ and $\sum_j \delta_j = tr(\boldsymbol{\rho}) = p$.

(b) When the principal components obtained using correlation matrix, the $X$-variables contribute equally to the pc's. Otherwise, $\boldsymbol{X}$ variable with large variance will dominates the first pc.

**Example:**

Let $\boldsymbol{X}$ be a random vector with mean $(2,5)^T$ and, the covariance and correlation matrices are respectively

$$\boldsymbol{\Sigma} = \begin{pmatrix} 1 & 4 \\ 4 & 100 \end{pmatrix} \text{ and } \boldsymbol{\rho} = \begin{pmatrix} 1 & 0.4 \\ 0.4 & 1 \end{pmatrix}.$$

Obtain the pc's using both covariance and correlation matrices and compare the results.

# 3   Special Cases

(a) $\boldsymbol{X}_j$ *are independent random variables*

Suppose $X_i$ and $X_j$ are independent for all $i \neq j$ then $\sigma_{ij} = 0$ and

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11} & 0 & \dots & 0 \\ 0 & \sigma_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{pp} \end{pmatrix}, \sigma_{ii} > \sigma_{jj} \text{ for all } i > j.$$

This implies that $\lambda_j = \sigma_{jj}$ and $Y_j = X_j$ for $j = 1, 2, \dots, p$.

(b) $\boldsymbol{X}_j$ *are identically distributed random variables*

That is, $\mathbf{Var}(X_j) = \sigma^2$ for all $i$ and $\mathbf{Cov}(X_i, X_j) = \rho\sigma^2$ for all $i \neq j$.

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma^2 & \rho\sigma^2 & \cdots & \rho\sigma^2 \\ \rho\sigma^2 & \sigma^2 & \cdots & \rho\sigma^2 \\ \vdots & \vdots & \ddots & \vdots \\ \rho\sigma^2 & \rho\sigma^2 & \cdots & \sigma^2 \end{pmatrix} . = \sigma^2 \begin{pmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & 1 \end{pmatrix}$$

where $-\frac{1}{p-1} < \rho < 1$. The covariance matrix given above is called *Intraclass Covariance Pattern*. Note that intraclass covariance matrix is also a *equal-correlation Matrix* but converse may not be true.

Suppose $\rho$ is positive, then

- the first eigenvalue and eigenvector are respectively given by $\lambda_1 = \sigma^2[1 + (p-1)\rho]$ and $\boldsymbol{e}_1^T = \frac{1}{\sqrt{p}}(1, 1, \ldots, 1)$, hence, the first pc, $Y_1 = \frac{1}{\sqrt{p}}(X_1 + X_2 + \ldots + X_p)$.

- the remaining eigenvalues and eigenvectors are $\lambda_j = \sigma^2(1 - \rho)$ and $\boldsymbol{e}_j^T = (e_{1j}, e_{2j}, \ldots, e_{pj})$ for $j = 2, 3, \ldots, p$

$$\text{where } e_{kj} = \begin{cases} [j(j-1)]^{-0.5} & \text{if } k < j \\ -(j-1)[j(j-1)]^{-0.5} & \text{if } k = j \\ 0 & \text{if } k > j \end{cases}$$

That is, $\boldsymbol{e}_2 = \frac{1}{\sqrt{2}}(1, -1, 0, \ldots, 0)$, $\quad \boldsymbol{e}_3 = \frac{1}{\sqrt{6}}(1, 1, -2, 0, \ldots, 0)$, $\boldsymbol{e}_4 = \frac{1}{\sqrt{12}}(1, 1, 1, -3, 0, \ldots, 0)$, $\quad \ldots$, $\boldsymbol{e}_p = \frac{1}{\sqrt{p(p-1)}}(1, 1, \ldots, 1, -[p-1])$.

**Example 3:**

Let $\boldsymbol{X}^T = (X_1, X_2, X_3, X_4, X_5)$ and $X_i$'s are identically distributed random variables with $\mathbf{Var}(X_i) = 5$ and $\mathbf{Corr}(X_i, X_j) = 0.8$ for $i \neq j$. Obtain the pc's.

# 4   Sample Principal Components

When $\boldsymbol{\Sigma}$ is unknown, we cannot obtain the principal components, but we can estimate the principal components by the sample principal components. The sample principal components are defined as a set of mutually independent, normalized linear combinations of original variables with maximum sample variance, similar to population principal components.

Assume the data $\boldsymbol{X}_1, \boldsymbol{X}_2, \cdots, \boldsymbol{X}_n$ represent $n$ independent samples from some $p$-dimensional population with *unknown* mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.

Then, we estimate

- $\boldsymbol{\mu}$ by sample mean vector, $\overline{\boldsymbol{X}}_n$ and

- $\boldsymbol{\Sigma}$ by sample covariance matrix,

$$\mathcal{S}_n = \frac{1}{n-1} \sum_{i=1}^{n} (\boldsymbol{x}_i - \overline{\boldsymbol{x}})(\boldsymbol{x}_i - \overline{\boldsymbol{x}})^T$$

The sample principal components can be obtained using the normalized eigenvectors of $\mathcal{S}_n$.

**Example:** Raw digital brightness values recorded by a satellite at a particular location in Australia are given in the following table:

| #  | B1 | B2 | B3 | B4 | #  | B1 | B2 | B3 | B4 |
|----|----|----|----|----|----|----|----|----|----|
| 1  | 19 | 15 | 22 | 11 | 14 | 19 | 16 | 24 | 12 |
| 2  | 21 | 15 | 22 | 12 | 15 | 25 | 25 | 38 | 20 |
| 3  | 19 | 13 | 25 | 14 | 16 | 20 | 29 | 19 | 3  |
| 4  | 28 | 27 | 41 | 21 | 17 | 28 | 29 | 18 | 2  |
| 5  | 21 | 15 | 25 | 13 | 18 | 25 | 26 | 42 | 21 |
| 6  | 21 | 17 | 23 | 12 | 19 | 21 | 18 | 12 | 12 |
| 7  | 19 | 16 | 24 | 12 | 20 | 26 | 24 | 43 | 21 |
| 8  | 19 | 12 | 25 | 14 | 21 | 30 | 31 | 18 | 3  |
| 9  | 28 | 29 | 17 | 3  | 22 | 28 | 27 | 44 | 24 |
| 10 | 28 | 26 | 41 | 21 | 23 | 30 | 31 | 18 | 2  |
| 11 | 19 | 16 | 24 | 12 | 24 | 30 | 31 | 18 | 2  |
| 12 | 29 | 32 | 17 | 3  | 25 | 21 | 16 | 22 | 12 |
| 13 | 19 | 16 | 22 | 12 |    |    |    |    |    |

Sample covariance matrix for brightness values given in B1 to B4 (Band 1 to 4):

$$
\mathcal{S}_{25} = \begin{pmatrix}
18.8767 & 26.8567 & 7.5550 & -5.1533 \\
26.8567 & 47.2433 & 5.1033 & -15.2300 \\
7.5550 & 5.1033 & 92.1900 & 58.8983 \\
-5.1533 & -15.2300 & 58.8983 & 48.5233
\end{pmatrix}
$$

Sample correlation matrix:

$$
\mathcal{R}_{25} = \begin{pmatrix}
1.00000 & 0.89933 & 0.18110 & -0.17027 \\
0.89933 & 1.00000 & 0.07733 & -0.31809 \\
0.18110 & 0.07733 & 1.00000 & 0.88061 \\
-0.17027 & -0.31809 & 0.88061 & 1.00000
\end{pmatrix}
$$

(a) Obtain the pc's using

1. sample covariance matrix $\mathcal{S}_{25}$ and

2. sample correlation matrix $\mathcal{R}_{25}$.

(b) Compare the above two sets of pc's.

# 5 Inference in Principal Components

**Theorem:**

If $\widehat{\lambda}_i$ is the $i^{th}$ eigenvalue of $\mathcal{S}_n$ and the distribution of the population is multivariate normal, then, $\widehat{\lambda}_i$ is normally distributed with mean $\lambda_i$ and variance $2\lambda_i^2/n$ where $\lambda_i$ is the $i^{\text{th}}$ eigenvalue of the population covariance matrix, $\boldsymbol{\Sigma}$.

Hence, $100(1-\alpha)\%$ confidence interval for $\lambda_i$ is,

$$
\left( \frac{\widehat{\lambda}_i}{1+a}, \frac{\widehat{\lambda}_i}{1-a} \right)
$$

where $a = z_{\alpha/2}\sqrt{\frac{2}{n}}$ and $z_{\alpha/2}$ is $100(1 - \alpha/2)$ percentile of the standard normal distribution.

**Note:**

(a) The distributions of $\widehat{\lambda}_i$ and $\widehat{\lambda}_j$ are independent for all $i \neq j$.

(b) If the sample size, $n$ is large relative to $p$, then the normality assumption on the population is not required.

**Example:**

A random sample of size 9 is obtained from a trivariate population and it is found that

$$
\overline{x}_9 = \begin{pmatrix} 71 \\ 84 \\ 56 \end{pmatrix} \quad \text{and} \quad \mathcal{S}_9 = \begin{pmatrix} 4 & 1 & 1 \\ 1 & 3 & 2 \\ 1 & 2 & 2 \end{pmatrix}.
$$

Compute the sample principal components and 90% confidence interval for their variances.

# 6   Testing the Significance of Smaller Principal Components.

Consider the Statistical Test,

$$
H_0 : \lambda_{r+1} = \lambda_{r+2} = \cdots = \lambda_{r+s} \quad \text{against}
$$
$$
H_1 : \lambda_{r+1}, \lambda_{r+2}, \cdots, \lambda_{r+s} \text{ are not all equal.}
$$

Using the likelihood ratio, it can be prove that the null hypothesis, $H_0$ is rejected, if $Q > c$ where $Q = s(n-1)\ln(\overline{\lambda}_r) - (n-1)\sum_{i=1}^{s}\ln(\widehat{\lambda}_{r+i})$ is the test statistic, $\overline{\lambda}_r = \frac{1}{s}\sum_{i=1}^{s}\widehat{\lambda}_{r+i}$ and $c$ is a constant dependent only on the significance level $\alpha$.

If the distribution of the population is normal (or the sample size is large relative to $p$), then the distribution of $Q$ is Chi-Square with $k = 0.5s(s+1)-1$ degrees of freedom.

Therefore, reject $H_0$ if $Q > \chi^2_k(\alpha)$ at $\alpha$ level of significance.

**Example:**

(Refer Example 8.3, *p.*439, Johnson & Wichern)

Given sample size, $n = 14$ and sample covariance matrix,

$$\mathcal{S}_{14} = \begin{pmatrix} 4.308 & 1.683 & 1.803 & 2.155 & -0.253 \\ 1.683 & 1.768 & 0.588 & 0.177 & 0.176 \\ 1.803 & 0.588 & 0.801 & 1.065 & -0.158 \\ 2.155 & 0.177 & 1.065 & 1.970 & -0.357 \\ -0.253 & 0.176 & -0.158 & -0.357 & 0.504 \end{pmatrix}.$$

Test whether the variance of the third and fourth pc are equal.

# 7   Intraclass Covariance Pattern

Note that the subsection **Testing for the Equal Correlation Structure** given in Johnson and Wichern (2007) page 457-459, gives a test for a equal correlation matrix. This test cannot be used for testing interaclass covariance pattern.

Let $\boldsymbol{\Sigma}^* = (\sigma_{ij})_{p \times p}$ with $\sigma_{ij} = \begin{cases} \sigma^2 & \text{if } i = j \\ \sigma^2 \rho & \text{if } i \neq j \end{cases}$ and $-\frac{1}{p-1} < \rho < 1$.

Then, the random vector with covariance matrix $\boldsymbol{\Sigma}^*$ has the intraclass covariance structure.

**Testing of Intraclass Covariance Pattern**

Consider the hypotheses

$H_0 : \boldsymbol{\Sigma} = \boldsymbol{\Sigma}^*$ against $H_1 : \boldsymbol{\Sigma} \neq \boldsymbol{\Sigma}^*$

This equivalent to

$H_0 : \lambda_2 = \lambda_3 = \cdots = \lambda_p$ against $H_1 : \lambda_2, \cdots, \lambda_p$ are not all equal.

Now taking $r = 1$ and $s = p - 1$ in the previous test, the $Q$ statistic reduced to

$$Q = (p - 1)(n - 1)\ln(\overline{\lambda}_1) - (n - 1)\sum_{i=2}^{p}\ln(\widehat{\lambda}_i)$$

where $\overline{\lambda}_1 = \frac{1}{p-1}\sum_{i=2}^{p}\widehat{\lambda}_i$ Note $Q$ has a chi-square distribution with $k = 0.5p(p - 1) - 1$ degrees of freedom .

Therefore, reject the null hypothesis, $H_0$ if

$$Q > \chi_k^2(\alpha)$$

at $\alpha$ level of significance.

**Example:** (Continue previous.)

Test the hypothesis that the population covariance matrix has a intrclass covariance structure.

# 8   Application of Principal Components to Linear Regression

Reference : Chatterjee & Price, " Regression Analysis by Example"

**Multicollinearity in Regression**

Consider the multiple linear regression

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p$$

One of the assumptions used to fit the above model using regression, is that the variables $X_1, X_2, \cdots, X_p$ are independent (or orthogonal).

However, it common that the observed data to shows a high correlations between some $X$ variables. Usually this may not seriously effect the regression analysis but in some cases,

- it is not possible to estimate the effects of an individual variables and,

- the estimated coefficients of some $X$ variables are very sensitive to slight changes in data.

This is refer to as *multicollinearity*. That is, data set is *collinear*. This effects the statistical analysis of the results and forecasting.

**Detection of Multicollinearity in Regression**

Multicollinearity may cause an instability in regression estimation. That is:

(a) Large changes in the estimated coefficients when a variable (or a data point) is added or deleted.

(b) The algebraic signs of the estimated coefficients do not conform the prior expectation.

(c) Standard errors(deviations) of the some estimated coefficients are too large.

The multicollinearity can also be detected using the principal components of the standardized $X$ variables. The small variances of these pc's points to multicollinearity.

**Correction for Multicollinearity: Principal Component Approach**

Since the pc's are independent variables, consider the regression $Y$, dependent variable and the pc's having significant variance. Ignore the insignificant pc's.

Let $u_i = \widehat{\omega}_{i1}z_1 + \widehat{\omega}_{i2}z_2 + \cdots + \widehat{\omega}_{ip}z_p$, be the $i^{th}$ significant pc, $i = 1, 2, \cdots, m$ for some $m < p$.

Here, $z_j$ is the standardized $x_j$, that is $z_j = (x_j - \overline{x}_j)/s_j$ where $\overline{x}_j$ and $s_j$ are the sample mean and standard deviation of $x_j$.

Next consider the regression model

$$Y = \alpha_0 + \alpha_1 u_1 + \alpha_2 u_2 + \cdots + \alpha_m u_m + \varepsilon.$$

This is equivalent to the original regression model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \varepsilon$$

because $u_i = \widehat{\omega}_{i1}z_1 + \widehat{\omega}_{i2}z_2 + \cdots + \widehat{\omega}_{ip}z_p$ and $z_j = (x_j - \overline{x}_j)/s_j$. Note that

$$\beta_j = \frac{1}{s_j}\sum_{i=1}^{m}\alpha_i\widehat{\omega}_{ij} \quad \text{and} \quad \beta_0 = \alpha_0 - \sum_{j=1}^{p}\left(\sum_{i=1}^{m}\alpha_i\widehat{\omega}_{ij}\right)\frac{\overline{x}_j}{s_j}.$$

Hence, the estimated regress model without the effect of multicollinearity is given by

$$Y = \widehat{\beta}_0 + \widehat{\beta}_1 x_1 + \widehat{\beta}_2 x_2 + \cdots + \widehat{\beta}_p x_p$$

where

$$\widehat{\beta}_j = \frac{1}{s_j}\sum_{i=1}^{m}\widehat{\alpha}_i\widehat{\omega}_{ij} \quad \text{and} \quad \widehat{\beta}_0 = \widehat{\alpha}_0 - \sum_{j=1}^{p}\left(\sum_{i=1}^{m}\widehat{\alpha}_i\widehat{\omega}_{ij}\right)\frac{\overline{x}_j}{s_j}$$

**Example:**

Using the hypothetical economic data given below, estimate the regression model

$$Y = \beta_0 + \beta_1 X_1 \beta_2 X_2 + \beta_3 X_3 + \varepsilon$$

where  $Y$  = cost of imports (millions of dollars)

$X_1$  = domestic production (millions of dollars)

$X_2$  = stock value (millions of dollars)

$X_3$  = domestic consumption (millions of dollars) and

$\varepsilon$  = random error component.

| # | $Y$ | $X_1$ | $X_2$ | $X_3$ | # | $Y$ | $X_1$ | $X_2$ | $X_3$ |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 16 | 149 | 4.2 | 108 | 10 | 27 | 232 | 5.1 | 164 |
| 2 | 16 | 161 | 4.1 | 115 | 11 | 26 | 239 | 0.1 | 168 |
| 3 | 19 | 172 | 3.1 | 123 | 12 | 31 | 258 | 5.6 | 177 |
| 4 | 19 | 176 | 3.1 | 127 | 13 | 33 | 270 | 3.9 | 187 |
| 5 | 18 | 181 | 1.1 | 132 | 14 | 37 | 288 | 3.1 | 200 |
| 6 | 20 | 191 | 2.2 | 138 | 15 | 43 | 305 | 4.6 | 214 |
| 7 | 23 | 202 | 2.1 | 146 | 16 | 49 | 323 | 7.0 | 224 |
| 8 | 27 | 212 | 5.6 | 154 | 17 | 50 | 337 | 1.2 | 232 |
| 9 | 28 | 226 | 5.0 | 162 | 18 | 57 | 354 | 4.5 | 243 |