

Statistics - Homework 1

Sarah Jallot and Victoire de Termont

11/17/2019

Statistics Homework, Assignment 1.

Problem 1 : Estimating parameters of a Poisson distribution to model the number of goals scored in football.

Question 1.

It is a discrete distribution as its support is N^* . The Poisson law is used to describe rare events in a large population. Examples of experiences appropriately modelled by a Poisson distribution are:

- The number of ice-creams sold on a beach during a hot summer day.
- The number of power failures occurring in Oslo in the winter over a period of 100 years.
- The number of times the Seine overflows over a period of 300 years.

Question 2.

- Computing the mean of the Poisson distribution:

$E(X) = \sum_{k=0}^{\infty} k \frac{e^{-\lambda} \lambda^k}{k!}$ summing directly to infinity because we know that the sum converges, but a more proper way would have been to sum to n then do the limit

$\iff E(X) = \lambda \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^{k-1}}{(k-1)!}$ simplifying by k as in 0 the sum term is null, and then factorising by λ

$\iff E(X) = \lambda \sum_{k=0}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k)!}$ now changing the indices to make the terms of the Poisson law appear

$\iff E(X) = \lambda * 1$ their sum is by definition equal to one

$\iff E(X) = \lambda$

- Computing the variance of the Poisson distribution:

We have : $Var(X) = E(X^2) - E(X)^2$.

Let's compute $E(X^2)$:

$E(X^2) = \sum_{k=0}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!}$ $E(X^2) = \sum_{k=0}^{\infty} k(k-1+1) \frac{e^{-\lambda} \lambda^k}{k!}$ again we already know that the sum converges so we slightly abuse notations.

$\iff E(X^2) = \sum_{k=0}^{\infty} k(k-1) \frac{e^{-\lambda} \lambda^k}{k!} + \sum_{k=0}^{\infty} k \frac{e^{-\lambda} \lambda^k}{k!}$

$\iff E(X^2) = \lambda^2 \sum_{k=2}^{\infty} \frac{e^{-\lambda} \lambda^{k-2}}{(k-2)!} + E(X)$ in the same way as for the expectation, we simplify by $k * (k-1)$ and then change the indices to make the Poisson sum of probability terms appear. We then note that this sum is equal to one.

$\iff E(X^2) = \lambda^2 * 1 + \lambda$

Finally, we get : $Var(X) = E(X^2) - E(X)^2 = \lambda^2 + \lambda - \lambda^2 = \lambda$

Question 3.

- Observations:

Assume there are n football games. Then, for each game, our observation is the number of goals scored and we have n observations.

- Model:

It is the set of all possible Poisson laws whose parameter is in R_+^* .

We have $M = \{p(\cdot|\theta), \theta \in R_+^*\}$ where $\forall k \in N^*, p(k|\theta) = e^{-\theta} \frac{\theta^k}{k!}$

- Parameter:

We are trying to estimate θ , which is the mean but also the variance of our Poisson law.

Let's call l the likelihood function of our model.

Question 4.

$l(x_1, \dots, x_n | \theta) = \prod_{i=1}^n l(x_i | \theta)$ by independence of the observations

$l(x_1, \dots, x_n | \theta) = \prod_{i=1}^n P(X_i = x_i | \theta)$

$l(x_1, \dots, x_n | \theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{x_i}}{x_i!}$

Let's compute the Maximum Likelihood Estimator $\hat{\theta}_{ML}$.

To do so, let's compute the log likelihood function of our model, called L .

$L(x_1, \dots, x_n | \theta) = \log(l(x_1, \dots, x_n | \theta))$

$L(x_1, \dots, x_n | \theta) = \log(\prod_{i=1}^n \frac{e^{-\theta} \theta^{x_i}}{x_i!})$

$L(x_1, \dots, x_n | \theta) = \sum_{i=1}^n \log(\frac{e^{-\theta} \theta^{x_i}}{x_i!})$

$L(x_1, \dots, x_n | \theta) = -n * \theta + \sum_{i=1}^n x_i * \log(\theta) - \sum_{i=1}^n \log(x_i!)$

Finding the maximum of the likelihood function is equivalent to finding the maximum of the log likelihood function.

If we are at the maximum, then $L'(x_1, \dots, x_n | \theta) = 0$

$\forall \theta \in \mathbb{R}, \theta > 0$

$L'(x_1, \dots, x_n | \theta) = -n + \sum_{i=1}^n \frac{x_i}{\theta}$

Thus,

$L'(x_1, \dots, x_n | \theta) = 0 \iff \theta = \frac{1}{n} \sum_{i=1}^n x_i$

The Maximum Likelihood Estimator is:

$\hat{\theta}_{ML} = \frac{1}{n} \sum_{i=1}^n X_i$

We note that the mean and the MLE are the same for the Poisson distribution.

Question 5.

Let's compute the expectation of $\hat{\theta}_{ML}$

$E(\hat{\theta}_{ML}) = E(\frac{1}{n} \sum_{i=1}^n X_i)$

$E(\hat{\theta}_{ML}) = \frac{1}{n} \sum_{i=1}^n E(X_i)$

$E(\hat{\theta}_{ML}) = \theta$

Thus, $\hat{\theta}_{ML}$ is an unbiased estimator of θ .

$\sqrt{n}(\hat{\theta}_{ML} - \theta) = \sqrt{n}(\frac{\sum_{i=1}^n X_i}{n} - E(X))$ as all X_i are iid and thus follow the same law as $X \sim P(\theta)$

According to the Central limit theorem, $\sqrt{n}(\hat{\theta}_{ML} - \theta)$ converges in distribution when $n \rightarrow \infty$ to the normal distribution $\mathcal{N}(0, \theta)$

Question 6.

We know that $\hat{\theta}_{ML}$ is an unbiased estimator of θ (which is the mean as well as the variance of our Poisson law).

$\hat{\theta}_{ML}$ converges in probability to θ .

So, given that $g : x \rightarrow \frac{1}{\sqrt{x}}$ is continuous on \mathbb{R}_+^* , by the continuous mapping theorem we get that $\frac{1}{\sqrt{\hat{\theta}_{ML}}}$

converges in probability to $\frac{1}{\sqrt{\theta}}$.

We showed in question 5 that $\sqrt{n} * (\hat{\theta}_{ML} - \theta)$ converged in distribution to $\mathcal{N}(0, \theta)$.

Thus, by Slutsky, $\sqrt{n} \frac{(\hat{\theta}_{ML} - \theta)}{\sqrt{\hat{\theta}_{ML}}}$ converges in law to $N(0, 1)$.

Empirical verification using R:

```
par(mfrow=c(1,3))
```

```
# Simulating the Poisson variable and setting the conditions of our experiment.
```

```
Nattempts = 1000
```

```
nsample    = 100
```

```

theta = 3
variable_sample = c()

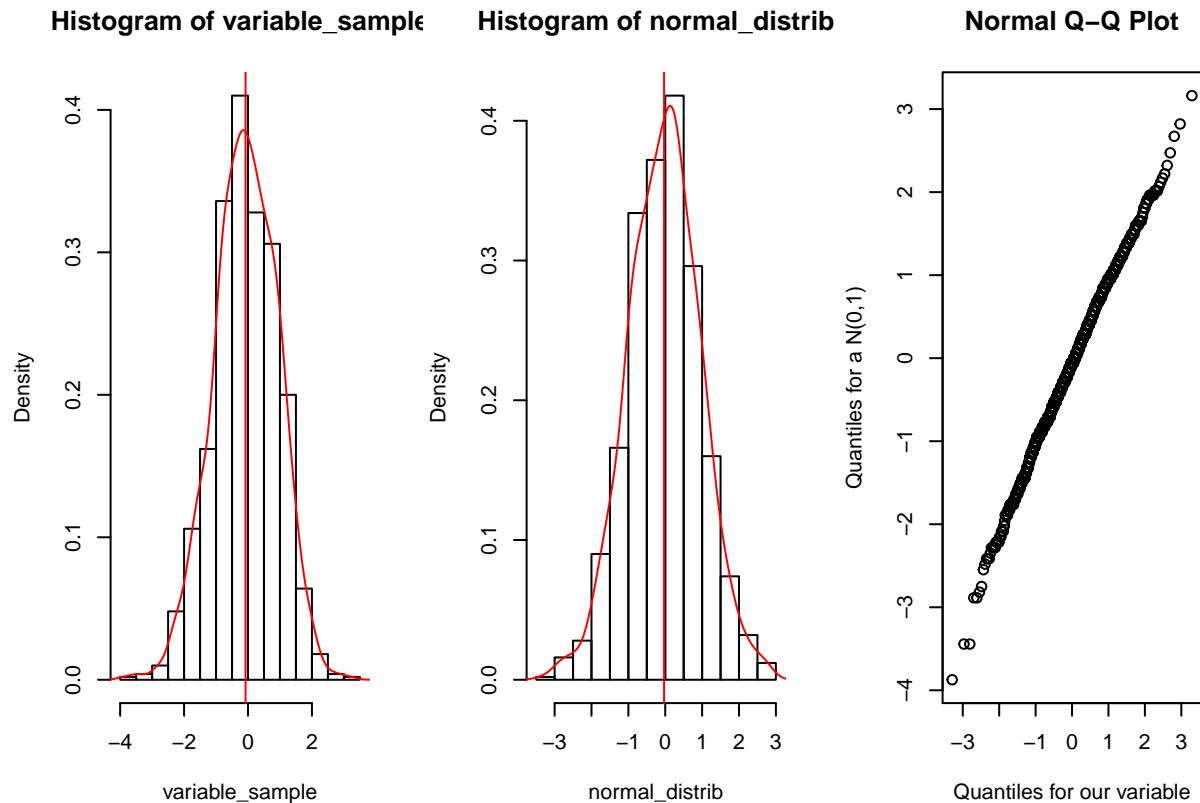
# Looping over the Poisson variable to center and scale it.
for (i in 1:Nattempts) {
  pois_sample = rpois(nsampl, theta)
  theta_mle = mean(pois_sample)
  variable_sample[i] = (theta_mle - theta) * sqrt(nsampl) / sqrt(theta_mle)}

# Plotting the histogram of our centered scaled poisson
hist(variable_sample, prob=TRUE)
d = density(variable_sample)
lines(d, col='red')
abline(v=mean(variable_sample), col='red')

# Building the histogram of a normal distribution function to compare
normal_distrib = rnorm(Nattempts, mean = 0, sd = 1)
hist(normal_distrib, prob=TRUE)
d = density(normal_distrib)
lines(d, col='red')
abline(v=mean(normal_distrib), col='red')

# Building our qq plot to check the validity of our conclusions
qqnorm(variable_sample, xlab= 'Quantiles for our variable', ylab='Quantiles for a N(0,1)')

```



Comparing the histograms and above all looking at the Q-Q plot, the theoretical result is verified.

Question 7.

We know that $\sqrt{n} \frac{(\hat{\theta}_{ML} - \theta)}{\sqrt{\hat{\theta}_{ML}}}$ converges in law to $\mathcal{N}(0, 1)$. Let $Z \sim \mathcal{N}(0, 1)$ and z_α its α quartile. We get the following interval:

$$\lim_{n \rightarrow \infty} P(-z_{1-\alpha/2} \leq \sqrt{n} \frac{(\hat{\theta}_{ML} - \theta)}{\sqrt{\hat{\theta}_{ML}}} \leq z_{1-\alpha/2}) = 1 - \alpha$$

We have:

$$\begin{aligned} -z_{1-\alpha/2} \leq \sqrt{n} \frac{(\hat{\theta}_{ML} - \theta)}{\sqrt{\hat{\theta}_{ML}}} \leq z_{1-\alpha/2} &\iff -\frac{\sqrt{\hat{\theta}_{ML}}}{\sqrt{n}} z_{1-\alpha/2} \leq (\hat{\theta}_{ML} - \theta) \leq \frac{\sqrt{\hat{\theta}_{ML}}}{\sqrt{n}} z_{1-\alpha/2} \\ &\iff \hat{\theta}_{ML} - \frac{\sqrt{\hat{\theta}_{ML}}}{\sqrt{n}} z_{1-\alpha/2} \leq \theta \leq \hat{\theta}_{ML} + \frac{\sqrt{\hat{\theta}_{ML}}}{\sqrt{n}} z_{1-\alpha/2} \end{aligned}$$

$\forall \epsilon > 0$, if we define: i. $\alpha_n = \hat{\theta}_{ML} - \frac{\sqrt{\hat{\theta}_{ML}}}{\sqrt{n}} z_{1-\alpha/2} - \epsilon$

ii. $\beta_n = \hat{\theta}_{ML} + \frac{\sqrt{\hat{\theta}_{ML}}}{\sqrt{n}} z_{1-\alpha/2} + \epsilon$ Then, we get $\lim_{n \rightarrow \infty} P(\theta \in [\alpha_n, \beta_n]) \geq 1 - \alpha$.

Question 8.

We use the δ -method with $g : x \rightarrow 2\sqrt{x}$.

Then, $\forall x > 0$, $g'(x) = \frac{1}{\sqrt{x}}$.

We apply the δ -method using this function g to $\sqrt{n}(\hat{\theta}_{ML} - \theta)$.

We know that:

$$E(\hat{\theta}_{ML}) = \theta$$

We compute:

$$[g'(\theta)]^2 V(\hat{\theta}_{ML}) = 1$$

We obtain: $\sqrt{n}(2\sqrt{\hat{\theta}_{ML}} - 2\sqrt{\theta})$ converges in distribution when $n \rightarrow \infty$ to the normal distribution $\mathcal{N}(0, 1)$.

Question 9.

We know that when $n \rightarrow \infty$, $\sqrt{n}(2\sqrt{\hat{\theta}_{MLE}} - 2\sqrt{\theta}) \sim \mathcal{N}(0, 1)$.

So, we get $\lim_{n \rightarrow \infty} P(-z_{1-\alpha/2} \leq \sqrt{n}(2\sqrt{\hat{\theta}_{MLE}} - 2\sqrt{\theta}) \leq z_{1-\alpha/2}) = 1 - \alpha$

We have: $-z_{1-\alpha} \leq \sqrt{n}(2\sqrt{\hat{\theta}_{MLE}} - 2\sqrt{\theta}) \leq z_{1-\alpha}$

$$\iff -\frac{1}{2\sqrt{n}} z_{1-\alpha} \leq (\sqrt{\hat{\theta}_{MLE}} - \sqrt{\theta}) \leq \frac{1}{2\sqrt{n}} z_{1-\alpha}$$

$\iff \sqrt{\hat{\theta}_{MLE}} - \frac{1}{2\sqrt{n}} z_{1-\alpha} \leq \sqrt{\theta} \leq \sqrt{\hat{\theta}_{MLE}} + \frac{1}{2\sqrt{n}} z_{1-\alpha}$ (1) When n goes to infinity, this inequality is greater than zero almost surely given that the left hand side of the inequality converges in probability to $\sqrt{\theta}$ almost surely.

Thus, given that $g : x \rightarrow x^2$ is continuous and strictly increasing on R_+ , we get:

(1) $\iff \theta \in [(\sqrt{\hat{\theta}_{MLE}} - \frac{1}{2\sqrt{n}} z_{1-\alpha})^2, (\sqrt{\hat{\theta}_{MLE}} + \frac{1}{2\sqrt{n}} z_{1-\alpha})^2]$. Thus, $\forall \epsilon > 0$, if we define:

$$c_n = (\sqrt{\hat{\theta}_{MLE}} - \frac{1}{2\sqrt{n}} z_{1-\alpha})^2 - \epsilon$$

$$d_n = (\sqrt{\hat{\theta}_{MLE}} + \frac{1}{2\sqrt{n}} z_{1-\alpha})^2 + \epsilon$$

Then, we have $\lim_{n \rightarrow \infty} P(\theta \in [c_n, d_n]) \geq (1 - \alpha)$

Question 10.

We know that if X follows a Poisson law of parameter λ , then $E(X) = V(X) = \lambda$. Therefore, we could choose the following estimators, based on our observations: $\hat{\theta}_1 = \frac{1}{n} \sum_1^n X_i$ (the empirical mean) using the first moment of a Poisson distribution. $\hat{\theta}_2 = \frac{1}{n-1} \sum_1^n (X_i - \bar{X}_n)^2$ (the empirical variance) using the first and second moments of a Poisson distribution. Note: we could also propose $\theta_2 = \frac{1}{n} \sum_1^n (X_i - \bar{X}_n)^2$ which is a biased estimator of the variance, as we will show later. $\hat{\theta}_1$ is the MLE we have already studied.

Question 11.

- i. Compute the bias of $\hat{\theta}_{ML}$.
 $b_{\theta}^*(\hat{\theta}_{ML}) = E^*(\hat{\theta}_{ML}) - \theta$
 $b_{\theta}^*(\hat{\theta}_{ML}) = 0$ and thus $\hat{\theta}_{ML}$ is an unbiased estimator of θ .
- ii. Compute the variance of $\hat{\theta}_{ML}$.
 $Var_{\theta}^*(\hat{\theta}_{ML}) = Var(\frac{\sum_{i=1}^n X_i}{n})$.
 $Var_{\theta}^*(\hat{\theta}_{ML}) = \frac{\sum_{i=1}^n Var(X_i)}{n^2}$ as the X_i are independent variables.
 $Var_{\theta}^*(\hat{\theta}_{ML}) = \frac{\theta}{n}$ as the X_i follow Poisson distribution of parameter θ .
- iii. Compute the quadratic risk of $\hat{\theta}_{ML}$. $Q = b_{\theta}^*(\hat{\theta}_{ML})^2 + Var^*(\hat{\theta}_{ML})$
 $Q = \frac{\theta}{n}$

Question 12.

- i. Let's compute $L''(x_1, x_2, \dots, x_n)$.
 $L'(x_1, x_2, \dots, x_n) = -n + \frac{1}{\theta} \sum_{i=1}^n x_i$
 $\Rightarrow L''(x_1, x_2, \dots, x_n) = -\frac{1}{\theta^2} \sum_{i=1}^n x_i$
- ii. So, we get:
 $I(\theta) = E^*(-L''(x_1, x_2, \dots, x_n))$
 $\iff I(\theta) = E^*(\frac{1}{\theta^2} \sum_{i=1}^n X_i)$
 $\iff I(\theta) = \frac{n\theta}{\theta^2}$
 $\iff I(\theta) = \frac{n}{\theta}$
As $\frac{1}{I(\theta)} = Var_{\theta}^*(\hat{\theta}_{ML})$ using q)11, $\hat{\theta}_{ML}$ is efficient. Thus, it is the estimator of θ with the smallest variance among all linear and unbiased estimators of θ .

Question 13.

$$\begin{aligned}\hat{\theta}_2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \\ \hat{\theta}_2 &= \frac{1}{n} \sum_{i=1}^n ((X_i - \theta) - (\bar{X}_n - \theta))^2 \\ \hat{\theta}_2 &= \frac{1}{n} \sum_{i=1}^n [(X_i - \theta)^2 + (\bar{X}_n - \theta)^2 - 2(X_i - \theta)(\bar{X}_n - \theta)] \\ \hat{\theta}_2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \theta)^2 + (\bar{X}_n - \theta)^2 - 2\frac{1}{n}(\bar{X}_n - \theta) \sum_{i=1}^n (X_i - \theta) \\ \hat{\theta}_2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \theta)^2 + (\bar{X}_n - \theta)^2 - 2(\bar{X}_n - \theta)(\frac{1}{n} \sum_{i=1}^n X_i - \theta) \\ \hat{\theta}_2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \theta)^2 + (\bar{X}_n - \theta)^2 - 2(\bar{X}_n - \theta)(\bar{X}_n - \theta) \\ \hat{\theta}_2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \theta)^2 - (\bar{X}_n - \theta)^2\end{aligned}$$

Question 14.

- i. Compute $E((\bar{X}_n - \theta)^2)$.
 $E((\bar{X}_n - \theta)^2) = E(\bar{X}_n^2 + \theta^2 - 2\theta\bar{X}_n)$
 $E((\bar{X}_n - \theta)^2) = E(\bar{X}_n^2) + \theta^2 - 2\theta * E(\bar{X}_n)$ due to the linearity of the expectation
 $E((\bar{X}_n - \theta)^2) = Var(\bar{X}_n) + E(\bar{X}_n)^2 + \theta^2 - 2\theta * E(\bar{X}_n)$
 $E((\bar{X}_n - \theta)^2) = \frac{\theta}{n} + \theta^2 + \theta^2 - 2\theta^2$
 $E((\bar{X}_n - \theta)^2) = \frac{\theta}{n}$
- ii. Prove that $\hat{\theta}_2$ is a biased estimator of θ .
 $E(\hat{\theta}_2) - \theta = E(\frac{1}{n} \sum_{i=1}^n (X_i - \theta)^2 - (\bar{X}_n - \theta)^2) - \theta$ due to question 13
 $E(\hat{\theta}_2) - \theta = \frac{1}{n} \sum_{i=1}^n E((X_i - \theta)^2) - E((\bar{X}_n - \theta)^2) - \theta$ by linearity of the expectation
 $E((X_i - \theta)^2) = V(X_i - \theta) + [E(X_i - \theta)]^2$
 $E((X_i - \theta)^2) = V(X_i)$
 $E((X_i - \theta)^2) = \theta$

Thus,
 $E(\hat{\theta}_2) - \theta = \theta - E((\bar{X}_n - \theta)^2) - \theta$
 $E(\hat{\theta}_2) - \theta = -E((\bar{X}_n - \theta)^2)$
 $E(\hat{\theta}_2) - \theta = -V(\bar{X}_n - \theta) - [E(\bar{X}_n - \theta)]^2$
 $E(\hat{\theta}_2) - \theta = -V(\bar{X}_n)$ as \bar{X}_n is an unbiased estimator of θ
 $E(\hat{\theta}_2) - \theta = -\frac{\theta}{n}$ and thus $\hat{\theta}_2$ is a biased estimator of θ

iii. To get an unbiased estimator of θ , we need the bias to be equal to 0.

$$\begin{aligned} E(\hat{\theta}_2) - \theta = -\frac{\theta}{n} &\iff 1E(\hat{\theta}_2) = \theta - \frac{\theta}{n} \\ E(\hat{\theta}_2) - \theta = -\frac{\theta}{n} &\iff 1E(\hat{\theta}_2) = \frac{(n-1)\theta}{n} \\ E(\hat{\theta}_2) - \theta = -\frac{\theta}{n} &\iff 1E(\frac{n\hat{\theta}_2}{n-1}) = \theta \end{aligned}$$

Then, we can use $\frac{(n)\hat{\theta}_2}{n-1}$ to get an unbiased estimator of θ .

Question 15.

a. We know thanks to the decomposition in q)13 that:

$$\hat{\theta}_2 = \frac{1}{n} * \sum_{k=1}^n (X_i - \theta)^2 - (\theta - \bar{X}_n)^2$$

b. We first note that $\bar{X}_n \rightarrow \theta$ in probability.

So, as $g : x \rightarrow (x - \theta)$ is continuous on R , we get that $(\bar{X}_n - \theta) \rightarrow 0$ in probability.

ii. We know that the X_i are i.i.d. So, the $(X_i - \theta)^2$ are also i.i.d.

$$E((X_i - \theta)^2) = E(X_i^2) - 2\theta E(X_i) + \theta^2 = Var(X_i) + E(X_i)^2 - 2\theta^2 + \theta^2 = \theta.$$

So, by using the central limit theorem, we get that $\sqrt{n}(\frac{1}{n} * \sum_{k=1}^n (X_i - \theta)^2 - \theta)$ converges in distribution to $\mathcal{N}(0, 2\theta^2 + \theta)$ when $n \rightarrow \infty$.

iii. As the X_i are i.i.d and \bar{X} converges in probability to θ , we apply the CLT again to get $\sqrt{n}(\frac{1}{n} \sum_{i=1}^n X_i - \theta)$ converges to $\mathcal{N}(0, \theta)$ when $n \rightarrow \infty$. So, as $\bar{X} - \theta$ converges in probability to 0, we get that $\sqrt{n}(\frac{1}{n} \sum_{i=1}^n X_i - \theta) * (\bar{X} - \theta)$ converges in law to zero.

However we know from the course that convergence in law to a constant implies convergence in probability to that constant.

iv. Finally, by Slutsky, we obtain that: $A = \sqrt{n}(\frac{1}{n} \sum_{k=1}^n (X_i - \hat{\theta})^2 - \theta) - \sqrt{n}(\frac{1}{n} \sum_{i=1}^n X_i - \theta)(\bar{X} - \theta)$ converges to $\mathcal{N}(0, 2\theta^2 + \theta)$ when $n \rightarrow \infty$. We now just have to notice that $A = \sqrt{n}(\hat{\theta}_2 - \theta)$. So, we get that $\sqrt{n}(\hat{\theta}_2 - \theta)$ converges to $\mathcal{N}(0, 2\theta^2 + \theta)$ when $n \rightarrow \infty$.

b. Asymptotic confidence interval for $\hat{\theta}_2$:

According to previous computation, we know that $\sqrt{n}(\hat{\theta}_2 - \theta)$ converges to $\mathcal{N}(0, 2\theta^2 + \theta)$ when $n \rightarrow \infty$.

Also, we know that $\hat{\theta}_{ML}$ converges in probability to θ .

Using the continuous mapping theorem with the function $g : x \rightarrow \frac{1}{\sqrt{2x^2 + x}}$, we obtain that:

$$\frac{1}{\sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}}} \text{ converges in probability to } \frac{1}{\sqrt{2\theta^2 + \theta}}$$

Thus, using Slutsky's theorem:

$$\frac{\sqrt{n}(\hat{\theta}_2 - \theta)}{\sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}}} \text{ converges to } \mathcal{N}(0, 1) \text{ when } n \rightarrow \infty.$$

$\lim_{n \rightarrow \infty} P(-q_{1-\alpha/2} \leq \frac{\sqrt{n}(\hat{\theta}_2 - \theta)}{\sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}}} \leq q_{1-\alpha/2}) = 1 - \alpha$ where $q_{1-\alpha/2}$ is the quantile of the standard normal distribution at level α .

We easily obtain: $\lim_{n \rightarrow \infty} P(-q_{1-\alpha/2} \sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}} \leq \sqrt{n}(\hat{\theta}_2 - \theta) \leq q_{1-\alpha/2} \sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}}) = 1 - \alpha$

And then: $\lim_{n \rightarrow \infty} P(\hat{\theta}_2 - \frac{q_{1-\alpha/2} \sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}}}{\sqrt{n}} \leq \theta \leq \hat{\theta}_2 + \frac{q_{1-\alpha/2} \sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}}}{\sqrt{n}}) = 1 - \alpha$ Thus, an asymptotic confidence interval for θ at level α is $IC(\alpha) = [A_n, B_n]$ where:

$$A_n = \hat{\theta}_2 - \frac{q_{1-\alpha/2} \sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}}}{\sqrt{n}}$$

$$B_n = \hat{\theta}_2 + \frac{q_{1-\alpha/2} \sqrt{2(\hat{\theta}_{ML})^2 + \hat{\theta}_{ML}}}{\sqrt{n}}.$$

- c. Comment on this interval: We note that $\hat{\theta}_2$ is a biased estimator of θ of bias $\frac{\theta}{n}$.

So any interval centered around $\hat{\theta}_2$ shall be less reliable than an interval centered around the unbiased estimator $\hat{\theta}_1$. Secondly we are using two different estimators in our confidence interval, which confirms that this interval is less reliable than the ones previously found.

Use it to confirm the first intervals, and only for a very big n .

- d. Compare the asymptotic variance to the one of the MLE and to the Cramer Rao bound:

According to q.12, $\frac{1}{I(\theta)} = Var_{\theta}^*(\hat{\theta}_{ML})$, which means that the Cramer Rao bound and the variance of $\hat{\theta}_{ML}$ are equal. We showed that $Var_{\theta}^*(\hat{\theta}_{ML}) = \frac{\theta}{n}$.

Also, according to q.15a, $\sqrt{n}(\hat{\theta}_2 - \theta)$ converges to $\mathcal{N}(0, 2\theta^2 + \theta)$ when $n \rightarrow \infty$.

This means that $Var(\hat{\theta}_2) \sim \frac{2\theta^2 + \theta}{n}$ when $n \rightarrow \infty$.

Comparing the two results, we can see that the variance of $\hat{\theta}_{ML}$ (or the Cramer Rao bound) converges faster to 0 when $n \rightarrow \infty$ than the variance of $\hat{\theta}_2$. As it is an unbiased estimator of θ with a smaller asymptotical variance, it means that it is a better estimator of θ when $n \rightarrow \infty$.

Question 16.

- i. Computing the moment generating function. We have:

$$E(e^{sX}) = \exp(-\theta) \sum_{k=0}^{\infty} e^{sk} \frac{\theta^k}{k!} \quad E(e^{sX}) = \exp(-\theta) * [\exp(\theta e^s) \exp(-\theta e^s)] * \sum_{k=0}^{\infty} \frac{(\theta e^s)^k}{k!}$$

$E(e^{sX}) = \exp(-\theta) * \exp(\theta e^s) [\exp(-\theta e^s) * \sum_{k=0}^{\infty} \frac{(\theta e^s)^k}{k!}]$ we are making the sum on N of the probability terms of a Poisson law of parameter θe^s appear.

Finally, we obtain $E(e^{sX}) = \exp(\theta(e^s - 1))$.

- ii. Computing the expectation:

$$G_X(s) = E(e^{sX}) = \sum_{k=0}^{\infty} (e^s)^k P(X = k)$$

$$G'_X(s) = \sum_{k=0}^{\infty} k(e^s)^k P(X = k).$$

$$\text{So, } G'_X(0) = E(X).$$

But $\forall s > 0, G'_X(s) = (\theta e^s) \exp(\theta(e^s - 1))$ by using the expression found in i.

Thus, $E(X) = \theta$.

- iii. Computing the variance:

$$G'_X(s) = \sum_{k=0}^{\infty} k(e^s)^k P(X = k)$$

$$G''_X(s) = \sum_{k=0}^{\infty} k^2(e^s)^k P(X = k)$$

$$\text{So, } G''_X(0) = E(X^2)$$

$$G''_X(0) = Var(X) + E(X)^2$$

$$G''_X(0) = Var(X) + \theta^2$$

But $\forall s > 0, G''_X(s) = G'_X(s) + (\theta e^s)^2 \exp(\theta(e^s - 1))$ # by using the expressions found in i. and ii. for G_X and its first derivative Thus, $G''_X(0) = \theta + \theta^2$ Thus, $Var(X) + E(X)^2 = \theta + \theta^2$ and $E(X)^2 = \theta^2$

We finally obtain $Var(X) = \theta$.

- iv. Recovering that a sum of independent Poisson laws of respective parameters λ_1 and λ_2 is a Poisson law of parameter $\lambda_1 + \lambda_2$:

$$G_{X_1+X_2}(s) = E(e^{s(X_1+X_2)}) = E(e^{sX_1} e^{sX_2}) = G_{X_1}(s) * G_{X_2}(s) \text{ by independence of } X_1 \text{ and } X_2$$

$$\iff G_{X_1+X_2}(s) = \exp(\lambda_1(e^s - 1)) * \exp(\lambda_2(e^s - 1)) = \exp((\lambda_1 + \lambda_2)(e^s - 1)).$$

So, the result we set out to prove is true given the unicity of the moment generating function.

- v. We have:

$$\text{a) } G_{X-\theta}(s) = E(\exp(sX - s\theta)) = \exp(-s\theta) * E(\exp(sX)) = \exp(-s\theta) * G_X(s)$$

b) $G_{X-\theta}(s) = \exp(-\theta) \sum_{k=0}^{\infty} \exp((k-\theta)s) * \frac{\theta^k}{k!}$
 $G'_{X-\theta}(s) = \exp(-\theta) \sum_{k=0}^{\infty} (k-\theta) \exp((k-\theta)s) * \frac{\theta^k}{k!}$
 $G'_{X-\theta}(0) = \exp(-\theta) \sum_{k=0}^{\infty} (k-\theta) * \frac{\theta^k}{k!} = E(X-\theta)$
 $G''_{X-\theta}(s) = \exp(-\theta) \sum_{k=0}^{\infty} (k-\theta)^2 \exp((k-\theta)s) * \frac{\theta^k}{k!} = E((X-\theta)^2)$ by an evident recurrence due to the derivative of the exponential, we get that $\forall n \in N^*, G_{X-\theta}^{(n)}(0) = E((X-\theta)^n)$

c) $G_{X-\theta}(s) = \exp(\theta(e^s - 1 - s))$ so $G_{X-\theta}(0) = 1$
 $\Rightarrow G'_{X-\theta}(s) = \theta(e^s - 1) * G_{X-\theta}(s)$ so $G'_{X-\theta}(0) = 0$
 $\Rightarrow G''_{X-\theta}(s) = \theta e^s (G_{X-\theta}(s) + G'_{X-\theta}(s)) - \theta G'_{X-\theta}(s)$ so $G''_{X-\theta}(0) = \theta$
 $\Rightarrow G_{X-\theta}^{(3)}(s) = \theta e^s (G_{X-\theta}(s) + 2 * G'_{X-\theta}(s) + G''_{X-\theta}(s)) - \theta G''_{X-\theta}(s)$ so $G_{X-\theta}^{(3)}(0) = \theta(1 + 0 + \theta) - \theta^2 = \theta$
 $\Rightarrow G_{X-\theta}^{(4)}(s) = (\theta e^s)(G_{X-\theta}(s) + 3G'_{X-\theta}(s) + 3G''_{X-\theta}(s) + G_{X-\theta}^{(3)}(s)) - \theta G_{X-\theta}^{(3)}(s)$ so $G_{X-\theta}^{(4)}(0) = \theta(1 + 0 + 3\theta + \theta) - \theta^2 = 3\theta^2 + \theta$.

By using the result we found in b, we finally obtain that $E((X_i - \theta)^4) = 3\theta^2 + \theta$.

So using $E((X_i - \theta)^2) = \theta$ as demonstrated in q.15) ii., we finally obtain:

$$Var((X_i - \theta)^2) = E((X_i - \theta)^4) - E((X_i - \theta)^2)^2 = 3\theta^2 + \theta - \theta^2 = 2\theta^2 + \theta.$$

Problem 2 - Analysis of the USJudgeRatings dataset.

I. First contact and preliminary observations.

1) Introduction to the dataset.

In this section of our report, the main idea is to discover the data, understand the features and get a first idea of the values they tend to take.

```
# First loading the dataset
data(USJudgeRatings)
```

The dataset name suggests that we are looking into a panel of US judges's ratings on a set of criteria. We do not know the outcome of the experiment.

```
# Printing the first ten rows of the dataset
library(knitr)
head(USJudgeRatings)
```

```
##           CONT INTG DMNR DILG CFMG DECI PREP FAMI ORAL WRIT PHYS RTEN
## AARONSON,L.H.  5.7  7.9  7.7  7.3  7.1  7.4  7.1  7.1  7.1  7.0  8.3  7.8
## ALEXANDER,J.M.  6.8  8.9  8.8  8.5  7.8  8.1  8.0  8.0  7.8  7.9  8.5  8.7
## ARMENTANO,A.J.  7.2  8.1  7.8  7.8  7.5  7.6  7.5  7.5  7.3  7.4  7.9  7.8
## BERDON,R.I.    6.8  8.8  8.5  8.8  8.3  8.5  8.7  8.7  8.4  8.5  8.8  8.7
## BRACKEN,J.J.   7.3  6.4  4.3  6.5  6.0  6.2  5.7  5.7  5.1  5.3  5.5  4.8
## BURNS,E.B.     6.2  8.8  8.7  8.5  7.9  8.0  8.1  8.0  8.0  8.0  8.6  8.6
```

Each observation corresponds to one judge of the US Superior Court. The variables are the columns. We suppose them to be the judge's ratings on a set of criteria. We do not know whether one of the columns is in fact an output, dependent on the others.

```
str(USJudgeRatings)
```

```
## 'data.frame':   43 obs. of  12 variables:
## $ CONT: num  5.7 6.8 7.2 6.8 7.3 6.2 10.6 7 7.3 8.2 ...
## $ INTG: num  7.9 8.9 8.1 8.8 6.4 8.8 9 5.9 8.9 7.9 ...
## $ DMNR: num  7.7 8.8 7.8 8.5 4.3 8.7 8.9 4.9 8.9 6.7 ...
## $ DILG: num  7.3 8.5 7.8 8.8 6.5 8.5 8.7 5.1 8.7 8.1 ...
## $ CFMG: num  7.1 7.8 7.5 8.3 6 7.9 8.5 5.4 8.6 7.9 ...
## $ DECI: num  7.4 8.1 7.6 8.5 6.2 8 8.5 5.9 8.5 8 ...
## $ PREP: num  7.1 8 7.5 8.7 5.7 8.1 8.5 4.8 8.4 7.9 ...
## $ FAMI: num  7.1 8 7.5 8.7 5.7 8 8.5 5.1 8.4 8.1 ...
## $ ORAL: num  7.1 7.8 7.3 8.4 5.1 8 8.6 4.7 8.4 7.7 ...
## $ WRIT: num  7 7.9 7.4 8.5 5.3 8 8.4 4.9 8.5 7.8 ...
## $ PHYS: num  8.3 8.5 7.9 8.8 5.5 8.6 9.1 6.8 8.8 8.5 ...
## $ RTEN: num  7.8 8.7 7.8 8.7 4.8 8.6 9 5 8.8 7.9 ...
```

We have 12 variables in this dataframe, all of which are numeric. A priori, they seem pretty close in terms of extreme values and dispersion, an insight to be confirmed through boxplots. We are now diving deeper in the context surrounding the experiment and our variables.

```
# We check that our dataframe isn't missing any values
print(paste("There are", sum(is.na(USJudgeRatings)), "missing values in the dataframe"))
```

```
## [1] "There are 0 missing values in the dataframe"
```

2) Contextual information on the data and the experiment methodology.

```
?USJudgeRatings
```

1. For a given observation, each of our numeric variables measures the performance of the judge along a given criterion as rated by lawyers. For instance, “DECI” measures the ability of a judge to make quick decisions.
2. We note that the data is old: it dates back to 1977, which means caution when handling it to draw conclusions today.
3. There are little indications on the period of time over which these ratings were given. It would be interesting to gather more data on the frequency at which they were reported, who required them and with what in mind.
4. We have no information on the grading methodology. Were the lawyers briefed in advance? Did they have a list of elements to observe, or did they go with their gut feeling?
5. We also have little information on the examiners themselves, the lawyers.
For a given criterion, numbers are relatively round. We can assume that one lawyer was responsible for grading a given judge on all criteria, and not a set of lawyers of which we averaged the results: this leaves higher possibilities of bias. Was it always the same lawyer, or did it change? How objective were they relatively to the trial, or trials, they observed?
6. We have no indication on the range of values the ratings could take.

The help function gives us more information on the variables themselves:

```
[,1] CONT Number of contacts of lawyer with judge.
# This could be an indicator of the integrity of the lawyer when rating the judge. This variable should not
# be correlated to others if the lawyers are objective in rating the judge's performance.
[,2] INTG Judicial integrity.
# This criterion is difficult to evaluate relying on objective facts. It relies mostly on the lawyer's impression.
[,3] DMNR Demeanor.
# Idem. Judicial integrity can be a part of demeanor, so these two variages might be linked.
[,4] DILG Diligence.
# Idem
[,5] CFMG Case flow managing.
# If we know of a set of best practices, then this criterion might be more reliable
[,6] DECI Prompt decisions.
# We expect DECI and diligence to be correlated. Moreover, is it more important to make prompt decisions
# or to give sound rulings? Further analysis of the dataset will give us more information on the
[,7] PREP Preparation for trial.
# Nothing to report
[,8] FAMI Familiarity with law.
# This should be one of the most objective rating criteria given that we are referring to a pre-defined set of
# norms. We expect it to be correlated to preparation for trial.
[,9] ORAL Sound oral rulings.
# This variable is already a judgement on the performance of the judge depending on other factors (preparation,
# for instance)
[,10] WRIT Sound written rulings.
# Same remark as previously. We would expect ORAL and WRIT to be strongly correlated given that it is
# only the form of the judgement that changes
[,11] PHYS Physical ability.
# Nothing to report
[,12] RTEN Worthy of retention.
# This seems to be a final rating answering the following question: should the judge be retained at the US
# Superior Court?. However, we do not know how we decide who stays or not given this grade. Indeed: will we
# keep judges above a certain grade threshold? Will we keep a certain number of judges, and then decide on
# the grade threshold that will help us decide who to keep? This is an important question which should guide
# the next steps of the experiment.
```

The eleven first columns seem to be grades that lead to a final evaluation on whether or not the judge should be retained in the US Superior Court. As noted in the comments, some criteria rely on more objective observations than others.

```
summary(USJudgeRatings)
```

```
##          CONT          INTG          DMNR          DILG
## Min.      : 5.700    Min.      :5.900    Min.      :4.300    Min.      :5.100
## 1st Qu.: 6.850    1st Qu.:7.550    1st Qu.:6.900    1st Qu.:7.150
## Median : 7.300    Median :8.100    Median :7.700    Median :7.800
## Mean      : 7.437    Mean      :8.021    Mean      :7.516    Mean      :7.693
## 3rd Qu.: 7.900    3rd Qu.:8.550    3rd Qu.:8.350    3rd Qu.:8.450
## Max.      :10.600    Max.      :9.200    Max.      :9.000    Max.      :9.000
##          CFMG          DECI          PREP          FAMI
## Min.      :5.400    Min.      :5.700    Min.      :4.800    Min.      :5.100
## 1st Qu.:7.000    1st Qu.:7.100    1st Qu.:6.900    1st Qu.:6.950
## Median :7.600    Median :7.700    Median :7.700    Median :7.600
## Mean      :7.479    Mean      :7.565    Mean      :7.467    Mean      :7.488
## 3rd Qu.:8.050    3rd Qu.:8.150    3rd Qu.:8.200    3rd Qu.:8.250
## Max.      :8.700    Max.      :8.800    Max.      :9.100    Max.      :9.100
##          ORAL          WRIT          PHYS          RTEN
## Min.      :4.700    Min.      :4.900    Min.      :4.700    Min.      :4.800
## 1st Qu.:6.850    1st Qu.:6.900    1st Qu.:7.700    1st Qu.:7.150
## Median :7.500    Median :7.600    Median :8.100    Median :7.800
## Mean      :7.293    Mean      :7.384    Mean      :7.935    Mean      :7.602
## 3rd Qu.:8.000    3rd Qu.:8.050    3rd Qu.:8.500    3rd Qu.:8.250
## Max.      :8.900    Max.      :9.000    Max.      :9.100    Max.      :9.200
```

The variables means range in [7.3, 8.0] which means they are pretty close to one another.

Their dispersion seems to be pretty similar also: we have one variable with a minimal value under 4.7, which is DMNR, and one variable whose maximum above 9.2, which is CONT. So except DMNR and CONT, all variables are contained within the [4.7, 9.2] interval.

We note that means are close to the medians, but lower. So most variables are slightly skewed towards the left with some grades significantly lower than others, as we will confirm by studying the boxplots and skewness.

II. Multivariate and univariate analysis - Descriptive statistics.

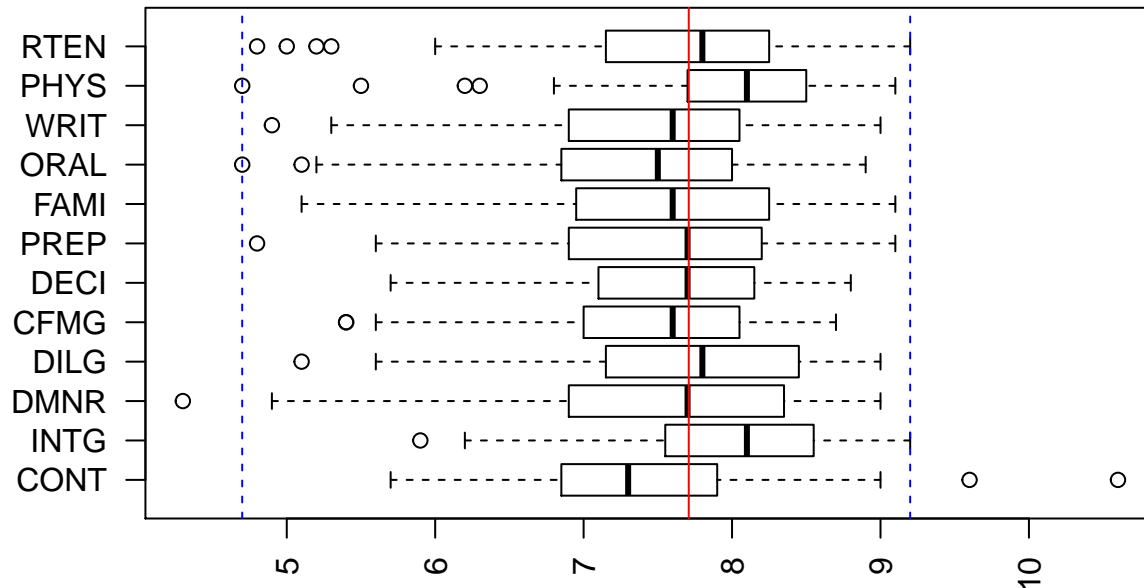
In this section of our report, the main idea is to deduce a probabilistic setting from our observations, and to study the links between variables. Due to the correlation matrix displayed below, we have decided to show first that the variables are generally very correlated, and then to study them independently. This decision comes from the fact that it is quite curious to obtain such strong correlations among our variables, as it means that the same information is repeated through different columns.

1) Position and dispersion comparison via boxplots.

```
# Plotting the boxplots horizontally because it makes it easier to
# compare their medians and dispersion.
bxpdat = boxplot(USJudgeRatings, las = 2, horizontal = TRUE)

# Plotting the mean of the median values to make any distance to it more visible
a = mean(c(7.3,8.1,7.7,7.8,7.6,7.7,7.7,7.6,7.5,7.6,8.1,7.8))
abline(v=a, col = 'red')

# Plotting the boundaries we noticed when analysing the summary
abline( v = 4.7, col = 'blue', lty = 2)
abline(v = 9.2, col = 'blue', lty = 2)
```



1. We confirm that our variables are approximately distributed within the same ranges, as visible by plotting in blue the boundaries identified previously, and by observing the shapes of the boxes and their whiskers.
2. In terms of **position**, all medians are close to 7.8, except PHYS and INTG which are significantly higher, both at (8.1).
We can assume that the physical ability and integrity of judges is generally good.
CONT is significantly lower at (7.3), which could mean that the lawyers didn't have too many contacts with the judges (or they reported it that way).
3. In terms of **dispersion**, the boxplots are slightly skewed towards the lower boundaries versus the median, confirming that some low ratings deflate the mean ratings that judges get.
Physical ability is narrow enough compared to the other variables, confirming that this physical ability is generally good for everyone.
DMNR has more spread and a smaller lower whisker than the others.
4. There are a few **outliers** we should look out for. We first observe that outliers are most of the time to the lower boundaries, which coincides with our earlier observation that the means were slightly lower than the medians.
The three variables for which outliers are the most visible are CONT, PHYS and DMNR, because they are far from the boxplots and distant from the blue boundaries identified in the summary.
 - a) CONT has an outlier above 10, which seems weird for two reasons:
 - i. The other data points suggest that 10 should be the upper boundary for all the ratings.
 - ii. The median value for CONT is significantly lower than for others. However, we can consider that CONT is not a grade but simply the measure of the number of contacts the lawyer had with the judge. However, it then seems weird that the values in CONT aren't integers...
 - b) The outliers for PHYS are far from the whiskers, even though they are within our blue boundaries. $\min(\text{PHYS}) = 4.7 \ll 7.7$ which is the first quartile. It makes sense for PHYS to have some exceptional outliers, even though they are far-left from the box whiskers.
 - c) DMNR has an outlier at 4.3, which doesn't seem to disturbing given that:

- i. It seems like an acceptable grade
- ii. As we will see later, INTG and DMNR are strongly correlated, and the judge whose INTGR grade is of 4.3 has a low DMNR grade of 6.4 (< first quartile) which makes sense.

```
USJudgeRatings$INTG[USJudgeRatings$DMNR == min(USJudgeRatings$DMNR)]
```

```
## [1] 6.4
```

d) We will consider the other outliers are acceptable for several reasons:

- i. They are not too far from the lower-box whisker (WRIT, ORAL in particular). PREP might be more intriguing with an outlier at 4.8, well under the first quartile (6.9).
- ii. We expect RTEN to present some outliers if the variables it depends on present some too.

Skewness and kurtosis of the data

```
#install.packages("e1071")
library(e1071)
```

Given our observations of the boxplot, we expect the skewness of the dataset to be negative for each variable:

```
for (name in colnames(USJudgeRatings)){
  s = skewness(USJudgeRatings[,name])
  print(c(name,s))}
```

```
## [1] "CONT"          "1.04831082852977"
## [1] "INTG"          "-0.813541922448701"
## [1] "DMNR"          "-0.914698975831847"
## [1] "DILG"          "-0.75619701024528"
## [1] "CFMG"          "-0.76275292434365"
## [1] "DECI"          "-0.621558748765366"
## [1] "PREP"          "-0.657253624881199"
## [1] "FAMI"          "-0.537792326172604"
## [1] "ORAL"          "-0.753038885790875"
## [1] "WRIT"          "-0.672682520721674"
## [1] "PHYS"          "-1.50417637184898"
## [1] "RTEN"          "-0.937360930294454"
```

As expected, we notice that practically all variables are slightly skewed towards the left except PHYS which is highly skewed towards the left and CONT which is highly skewed towards the right (absolute value greater than one).

```
for (name in colnames(USJudgeRatings)){
  k = kurtosis(USJudgeRatings[,name])
  print(c(name,k))}
```

```
## [1] "CONT"          "1.51221182877214"
## [1] "INTG"          "0.257028146294933"
## [1] "DMNR"          "0.274266854918179"
## [1] "DILG"          "0.14697661975641"
## [1] "CFMG"          "-0.0964006409476323"
## [1] "DECI"          "-0.531858166900887"
## [1] "PREP"          "-0.0036222033312967"
## [1] "FAMI"          "-0.36534841231648"
## [1] "ORAL"          "0.0126668778846812"
```

```
## [1] "WRIT"          "-0.108868716116502"
## [1] "PHYS"          "2.15947201560131"
## [1] "RTEN"          "0.255742065912552"
```

Kurtosis analysis confirms that outliers PHYS (kurt = 2.2 >>1) and CONT (kurt = 1.5 > 1) outliers are significantly far from the mean, as observed on the boxplot.

Let's try to remove extreme outliers from both and have a look at the kurtosis for CONT and PHYS.

```
# First identifying the outliers, then manually getting their row as the dataset is small
USJudgeRatings[USJudgeRatings$CONT==10.600,]
```

```
##          CONT INTG DMNR DILG CFMG DECI PREP FAMI ORAL WRIT PHYS RTEN
## CALLAHAN,R.J. 10.6   9  8.9  8.7  8.5  8.5  8.5  8.5  8.6  8.4  9.1   9
```

```
USJudgeRatings[USJudgeRatings$PHYS<=6,]
```

```
##          CONT INTG DMNR DILG CFMG DECI PREP FAMI ORAL WRIT PHYS RTEN
## BRACKEN,J.J.  7.3  6.4  4.3  6.5  6.0  6.2  5.7  5.7  5.1  5.3  5.5  4.8
## MIGNONE,A.F.  6.6  7.4  6.2  6.2  5.4  5.7  5.8  5.9  5.2  5.8  4.7  5.2
```

```
k = kurtosis(USJudgeRatings[c(1:6,8:43),"CONT"])
print(c("CONT",k))
```

```
## [1] "CONT"          "0.0159913658922792"
```

This seems to confirm our intuition that the 10.6 value we observed for CONT is a reporting mistake : removing it alone drops kurtosis drastically.

```
k = kurtosis(USJudgeRatings[c(1:4,6:22,24:43),"PHYS"])
print(c("PHYS",k))
```

```
## [1] "PHYS"          "0.305223588061307"
```

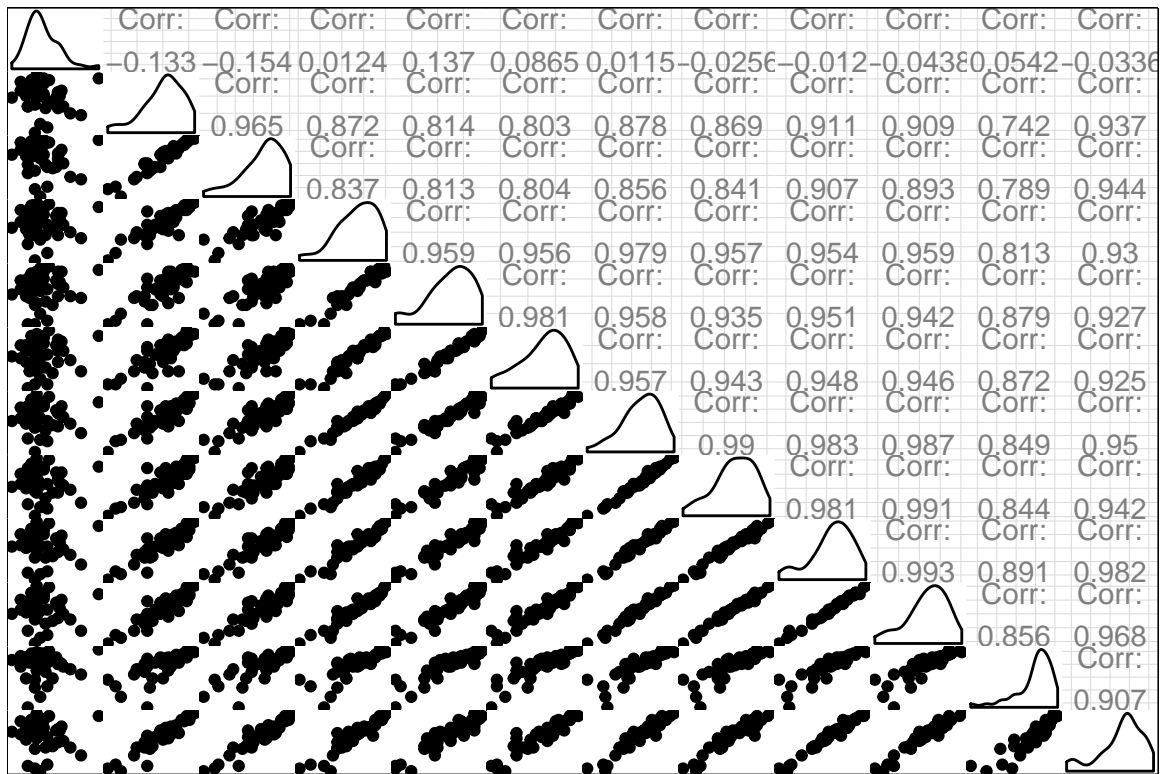
Removing the two outliers for PHYS drops kurtosis significantly. However, ~5/10 for physical shape doesn't seem like a shocking value, so it doesn't mean we should necessarily remove these outliers.

```
library(corrplot)
library(ggplot2)
library(dplyr)
library(GGally)
```

2) General multivariate analysis: scatterplots and correlation matrix.

a) Scatterplots.

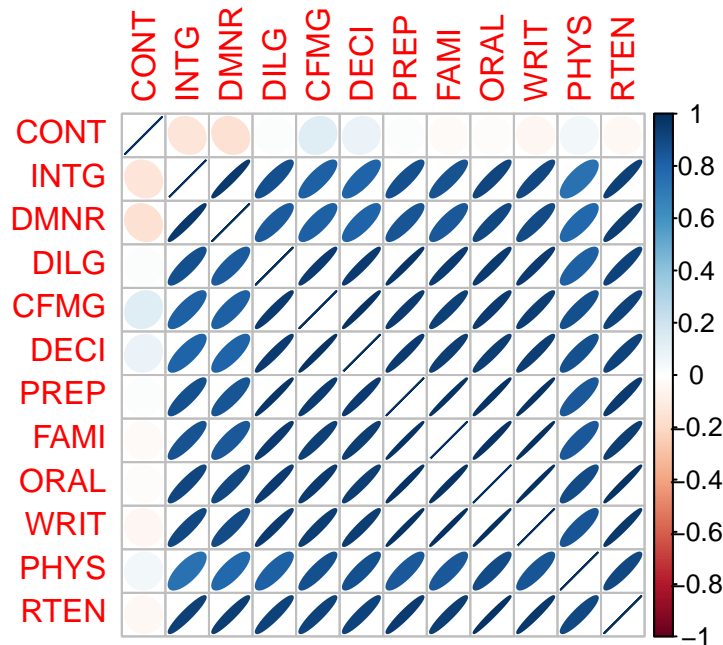
```
# We are first plotting the scatterplots of each variable relatively to the others to see if any clear
theme_set(theme_classic(base_size = 0.2))
ggpairs(USJudgeRatings, progress = F)
```



1. We observe that our variables are in fact strongly correlated, given that some scatterplots display near linear relations. This entails that if we want to predict RTEN using other variables, we will have to discriminate between the most useful features in explaining RTEN (maybe using fast forward selection). The relation between PREP and FAM in particular is quasi linear. This shouldn't surprise us, because it seems logical that the more one prepares a case, the better one knows the law surrounding it.
2. We remark that CONT is the only variable that has very little correlation with the others, with a highest correlation of ~ 0.2 in absolute value with DMNR.
3. It seems that the best linear relation between RTEN, our output, and another variable is for ORAL, meaning that ORAL would be the best single feature explaining the outcome.

b) Correlation matrix to confirm and detail the insights from the pairplot.

```
Corr = cor(USJudgeRatings)
corrplot(Corr, method = "ellipse")
```



We confirm that our variables are very strongly correlated to one another, with most correlations between 0.8 and 1. The only variable which is rather uncorrelated to others is indeed CONT.

We have the following relations (analysing RTEN in a separate point as it is the output rating):

1. Variables that are highly correlated to multiple other variables include DILG, CFMG, DECI, PREP, FAMI, ORAL, and WRIT.

a. WRIT and ORAL are very strongly correlated to one another ($\text{corr} = 0.99$), as well as to FAMI and PREP ($\text{corr} > 0.98$). b. CFMG and DECI are strongly correlated ($\text{corr} = 0.98$): promptitude plays a role in case flow management. To a lesser extent, they are also very strongly correlated to DILG, PREP and ORAL ($\text{corr} > 0.95$).

2. Variables that are highly correlated to less than two variables include INTG, DMNR, and PHYS.

a. INTG and DMNR are strongly correlated ($\text{corr} = 0.96$), but are less correlated to the other variables. They are both very slightly negatively correlated to CONT.

b. Among strongly correlated variables, PHYS is the one that least depends on others. However, it is strongly correlated to CFMG, DECI and ORAL (correlation > 0.87). It would be an interesting variable to include in a model if it explained RTEN well: however, it does have one of the lowest correlation coeffs to the variable (0.91).

3. The only uncorrelated variable is CONT. CONT isn't related to other variables, and in particular it doesn't correlate to RTEN which is our output, so it probably isn't good towards explaining the data. However, it does underline the correlations between other codependent variables.

4. RTEN is highly correlated to all variables except CONT, with $\text{corr} > 0.93$ for each one except CONT (-0.03) and PHYS (0.91).

As noticed earlier, the highest correlation between RTEN and another variable occurs for ORAL ($\text{corr} = 0.98$).

3) Histograms: getting the laws of the essential variables.

Intuitively, we choose to compare our sample distribution to a normal law distribution of same mean and variance, for $n = 43$.

Indeed, for each criterion, we can consider the grades attributed to the judges are i.i.d. and we can apply the CLT to them.

Given the high correlations between variables, we chose not to plot them all but will study some interesting relations in the next section.


```

par(mfrow=c(1,3))
names = c("RTEN", "ORAL", "DECI", "DMNR", "PHYS", "DILG", "CONT") # selected RTEN as it is the output;

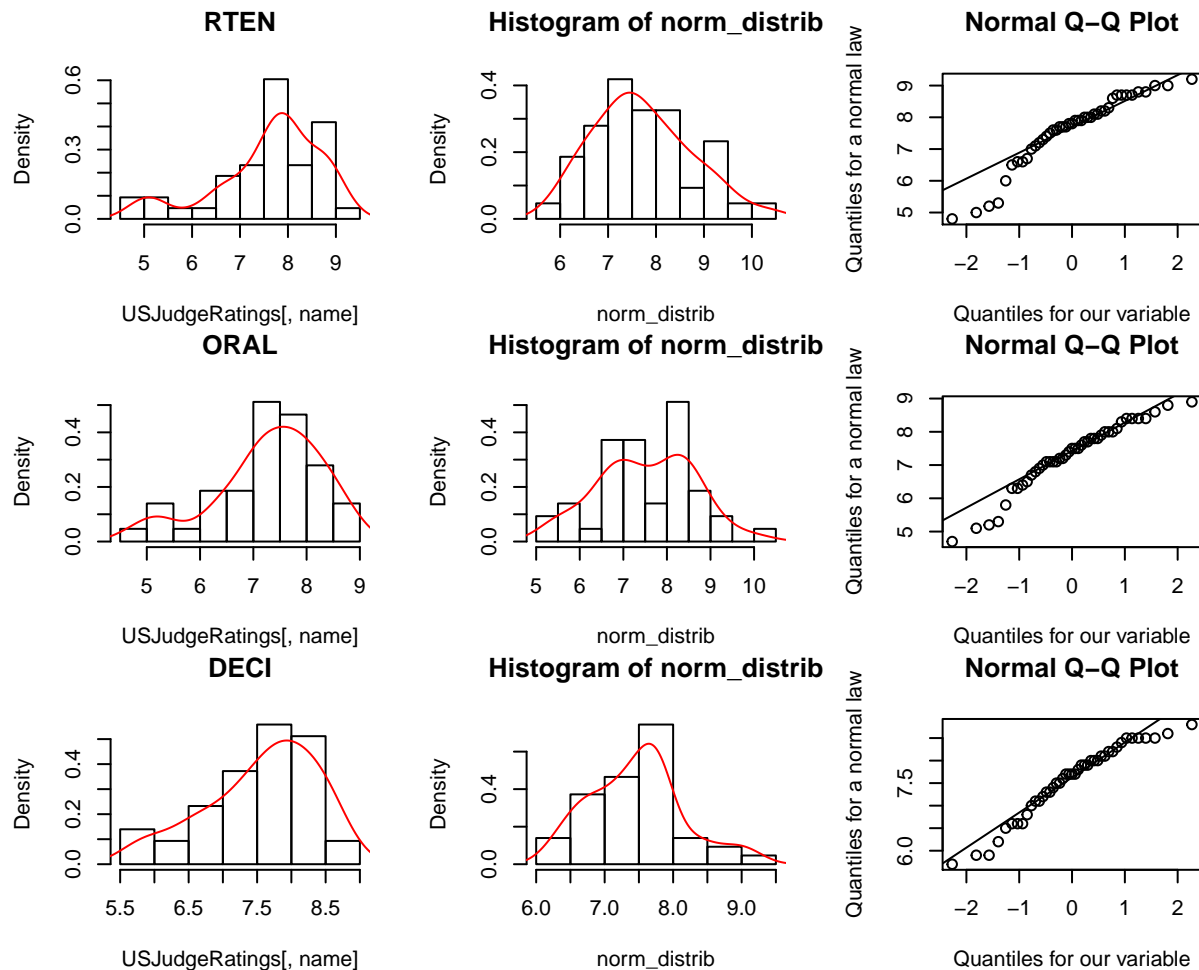
for (name in names) {

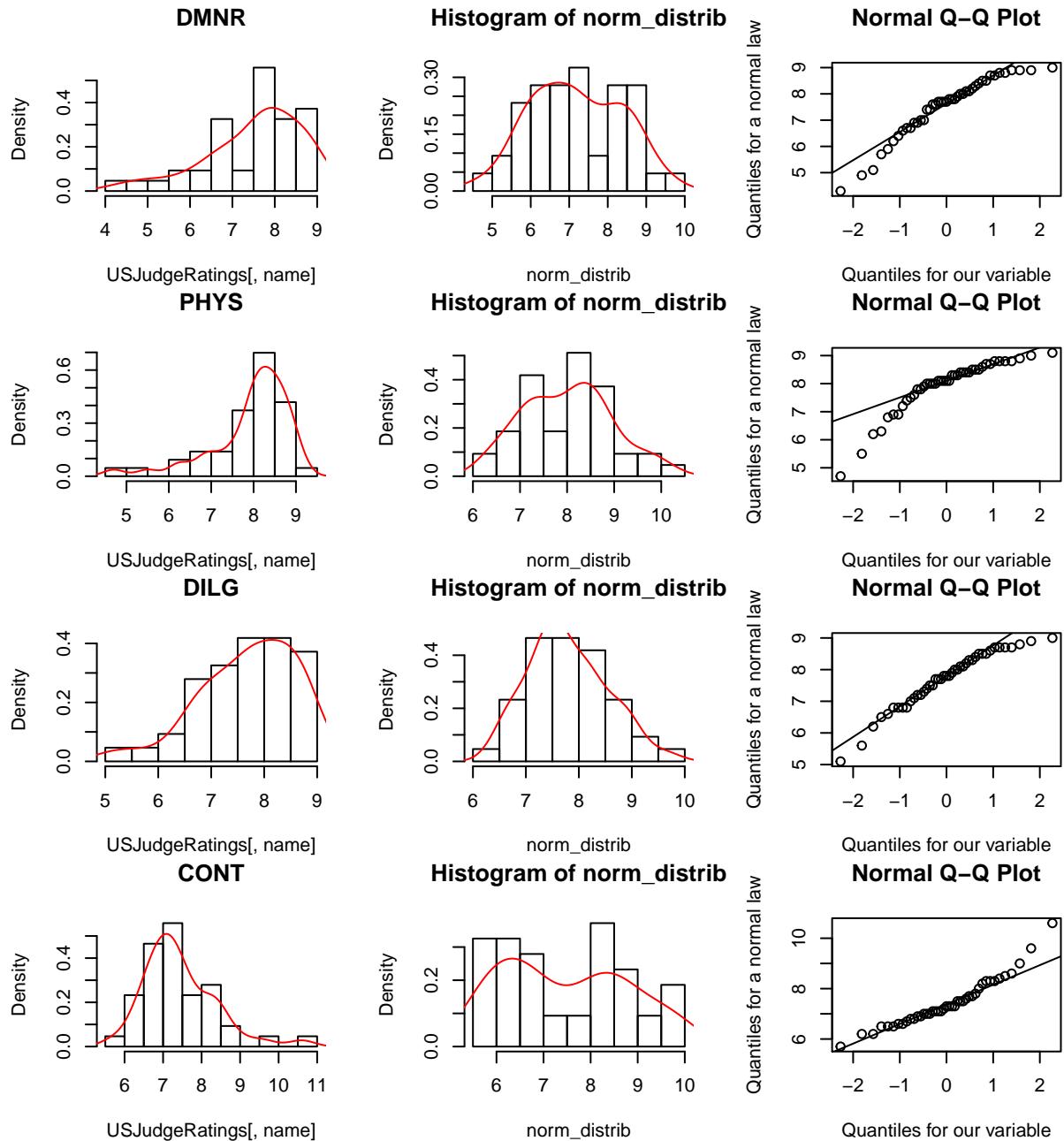
hist(USJudgeRatings[,name], main = name, proba = TRUE, breaks = 10)
d= density(USJudgeRatings[,name])
lines(d, col = 'red')

norm_distrib=rnorm(43,mean(USJudgeRatings[,name]),sqrt(var(USJudgeRatings[,name])))
hist(norm_distrib, prob=TRUE, breaks = 10)
d = density(norm_distrib)
lines(d, col='red')

qqnorm(USJudgeRatings[,name], xlab = "Quantiles for our variable", ylab = "Quantiles for a normal law")
qqline(USJudgeRatings[,name])}

```





We do not expect a perfect match given that $n = 43$ is relatively small. However, we observe that our samples fit the Normal distribution quite well.

2. RTEN, the grade we consider dependent on the others, fit the normal law less well than the others.
3. PHYS is the variable that fits the least with a normal law. Given that we have a certain number of outliers and they are far from the mean, to the left, we will consider that we are in an exceptional case where we found several significant outliers for PHYS, but we maintain the normal distribution assumption.
4. Except from the 10.6 outlier which we would a priori discard, CONT fits the normal distribution pretty well.

4) Focused multivariate analysis: some interesting combinations to study correlations.

1. ORAL vs WRIT (corr = 0.99), ORAL vs FAMI (corr = 0.98) : very highly correlated variables.
We expect these variables to follow very similar laws.

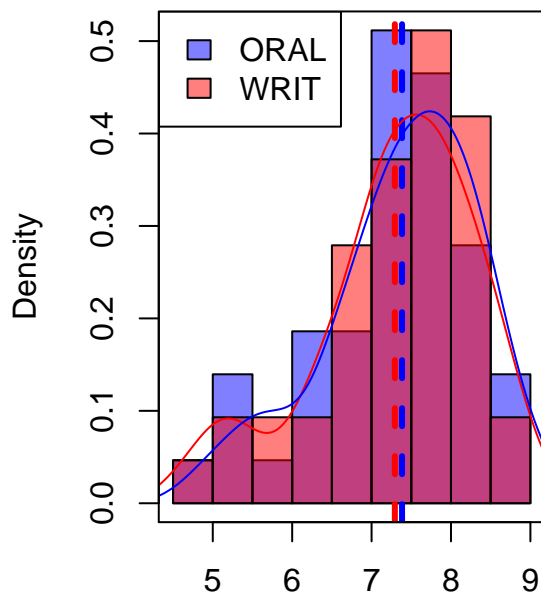
```

par(mfrow=c(1,2))
hist(USJudgeRatings$ORAL, col = rgb(0,0,1,0.5), proba = TRUE, main= "Comparative hist. ORAL/WRIT", break=15)
hist(USJudgeRatings$WRIT, col = rgb(1,0,0,0.5), proba = TRUE, braks = 15, add = TRUE) # blue
abline(v=mean(USJudgeRatings$ORAL), col="red", lwd=3, lty=2)
abline(v=mean(USJudgeRatings$WRIT), col="blue", lwd=3, lty=2)
d1 = density(USJudgeRatings$ORAL)
lines(d1, col = 'red')
d2 = density(USJudgeRatings$WRIT)
lines(d2, col = 'blue')
box()
legend( 'topleft', fill = c(rgb(0,0,1,0.5), rgb(1,0,0,0.5)), legend = c('ORAL','WRIT'))

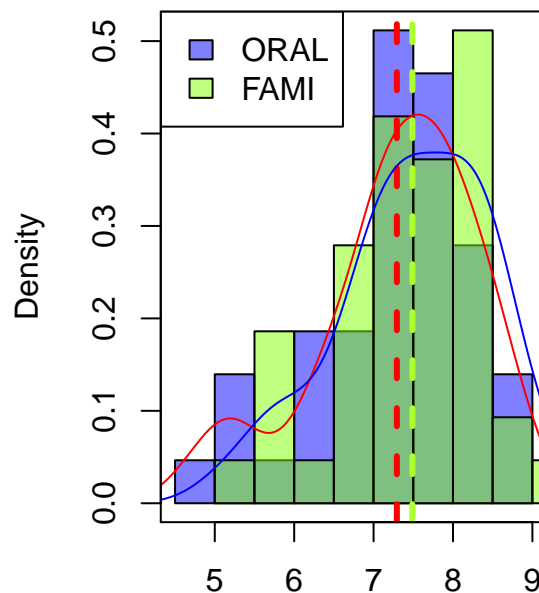
hist(USJudgeRatings$ORAL, col = rgb(0,0,1,0.5), proba = TRUE, main= "Comparative hist. ORAL/FAMI", break=15)
hist(USJudgeRatings$FAMI, col = rgb(0.5,1,0,0.5), proba = TRUE, braks = 15, add = TRUE) #
abline(v=mean(USJudgeRatings$ORAL), col="red", lwd=3, lty=2)
abline(v=mean(USJudgeRatings$FAMI), col="greenyellow", lwd=3, lty=2)
d1 = density(USJudgeRatings$ORAL)
lines(d1, col = 'red')
d2 = density(USJudgeRatings$FAMI)
lines(d2, col = 'blue')
box()
legend( 'topleft', fill = c(rgb(0,0,1,0.5), rgb(0.5,1,0,0.5)), legend = c('ORAL', 'FAMI'))

```

Comparative hist. ORAL/WRIT



Comparative hist. ORAL/FAMI



As expected, ORAL and WRIT fit quasi perfectly.
ORAL and FAMI are also very close albeit a little less.

2. RTEN vs ORAL (corr = 0.98), RTEN vs PREP (corr = 0.96). ORAL is the best single feature explaining RTEN, making it interesting to plot both variables together. PREP is the second highest correlated variable if we exclude WRIT (highly correlated to ORAL).

```

# ORAL as best explicative variable for RTEN
par(mfrow=c(1,2))
hist(USJudgeRatings$RTEN, col = rgb(0,0,1,0.5), proba = TRUE, main= "Comparative hist. RTEN/ORAL", break=15)

```

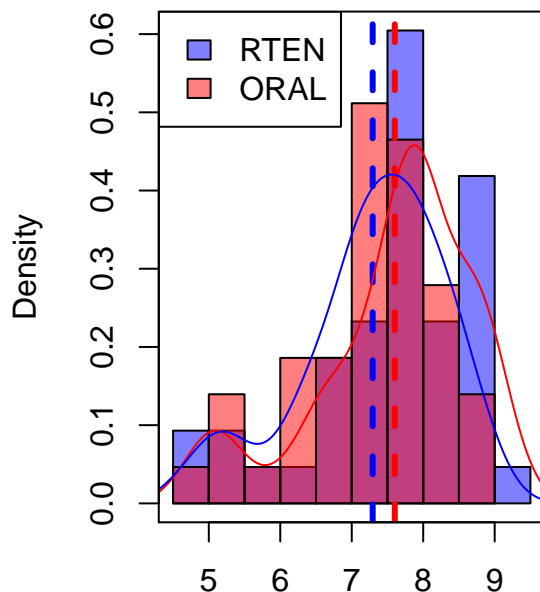
```

hist (USJudgeRatings$ORAL, col = rgb(1,0,0,0.5), proba = TRUE, breaks = 15, add = TRUE) # blue
abline(v=mean(USJudgeRatings$RTEN), col="red", lwd=3, lty=2)
abline(v=mean(USJudgeRatings$ORAL), col="blue", lwd=3, lty=2)
d1 = density(USJudgeRatings$RTEN)
lines(d1, col = 'red')
d2 = density(USJudgeRatings$ORAL)
lines(d2, col = 'blue')
box()
legend( 'topleft', fill = c(rgb(0,0,1,0.5), rgb(1,0,0,0.5)), legend = c('RTEN', 'ORAL'))

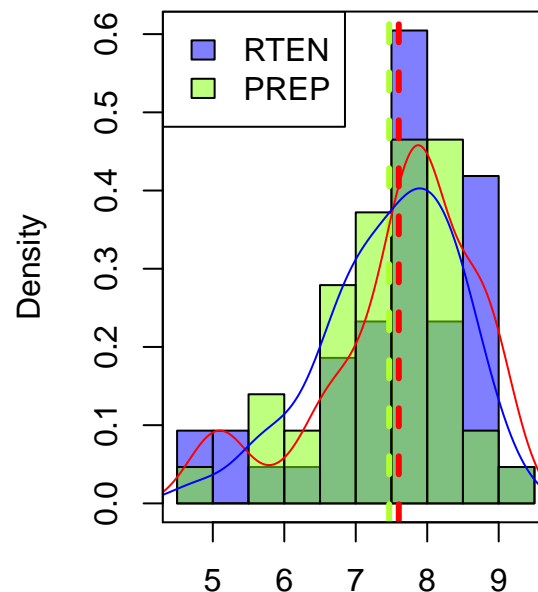
# PREP as second best explicative variable for RTEN aside WRIT (highly correlated to ORAL)
hist(USJudgeRatings$RTEN, col = rgb(0,0,1,0.5), proba = TRUE, main= "Comparative hist. RTEN/PREP", breaks = 15, add = TRUE)
hist (USJudgeRatings$PREP, col = rgb(0.5,1,0,0.5), proba = TRUE, breaks = 15, add = TRUE)
abline(v=mean(USJudgeRatings$RTEN), col="red", lwd=3, lty=2)
abline(v=mean(USJudgeRatings$PREP), col="greenyellow", lwd=3, lty=2)
d1 = density(USJudgeRatings$RTEN)
lines(d1, col = 'red')
d2 = density(USJudgeRatings$PREP)
lines(d2, col = 'blue')
box()
legend( 'topleft', fill = c(rgb(0,0,1,0.5), rgb(0.5,1,0,0.5)), legend = c('RTEN', 'PREP'))

```

Comparative hist. RTEN/ORAL



Comparative hist. RTEN/PREP



Albeit imperfect, we know that ORAL is approximately the best unique variable explaining RTEN. Indeed, its distribution is close to that of RTEN, and closer than the second best unique variable (except WRIT), PREP.

5) PCA component analysis. Given that our features are very highly correlated, when fitting our model, we can assume leaving some out will marginally affect model performance.

Performing PCA on our dataset, then looking at the proportion of variance explained, allows us to confirm our assumption and quantify how many features are useful to use.

```

# Divide the data
pca.train = USJudgeRatings[1:43,]
#

```

```

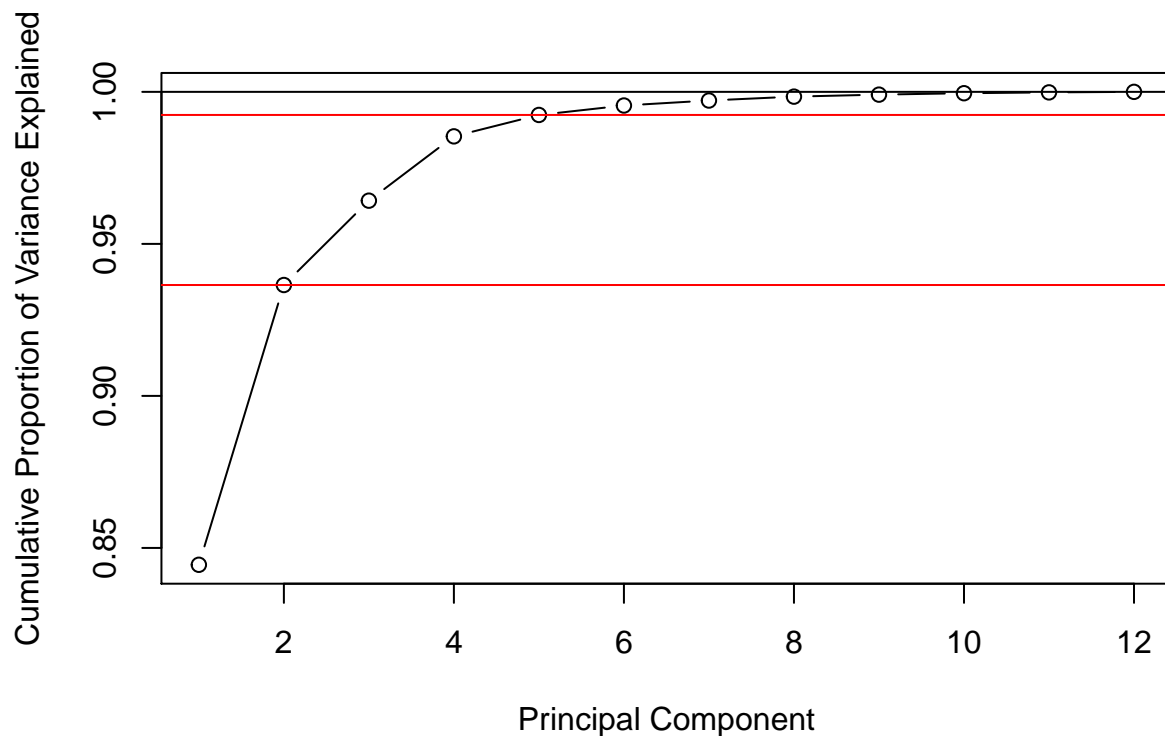
prin_comp = prcomp(pca.train, scale. = T)

# Compute standard deviation of each principal component
std_dev = prin_comp$sdev

# Compute variance
pr_var = std_dev^2
prop_varex = pr_var/sum(pr_var)

# Cumulative scree plot
plot(cumsum(prop_varex), xlab = "Principal Component",
     ylab = "Cumulative Proportion of Variance Explained",
     type = "b")
abline(h=1,col='black')
abline(h=0.9364709,col='red')
abline(h=0.9924043,col='red')

```



We note that 2 components already explain approximately 94% of the variance!
 Above 5-6 features, adding a feature brings along no significant improvement.

Conclusion

1. We examined a dataset of grades given by lawyers to US Superior Court judges of which we will assume that the objective is to determine whether these judges were fit to retain, or not (it could've been which ones need training, for instance).
2. Contextual information is lacking, mostly on the objectivity of the measures, and the data is quite old.
3. Most importantly, we do not know how the main question of the experiment will be answered. Feature engineering might be helpful here by deriving a classifier from RTEN: 1 if the judge is retained, 0 if they are not. As discussed in the introduction, the methodology choice is still to be defined. Although

there are more refined options, the two main possibilities are: keeping a number of judges, and then determining the threshold grade to make the cut; or, defining the grade threshold to be met in RTEN, then selecting the judges to keep. We would choose option 1 given that there are probably constraints on the number of judges we want to keep (budgetary or others).

4. We noticed that the variables are generally skewed towards the left-hand-side, as showed by studying their means versus their median as well as skewness.
5. Most variables are highly correlated, which means that if we want to run a model on the dataset, we will have to discriminate closely between those that are relevant towards explaining the data. In particular, WRIT and ORAL are highly correlated. Fast forward selection seems to be the best methodology we could implement towards modelling.
6. The best single variable explaining RTEN, the final grade on which our decision will be based, is ORAL. The second best variable is PREP, but it doesn't mean that it will be included in our model if we use forward selection (forward selection looks at the efficiency of a combination of features).
7. PCA analysis would indicate that if we use linear regression, fitting a model with 5 features would be a very good decision already (to be confirmed in our bonus section!)

```
library(DataExplorer, verbose = FALSE)
library(lmtest, verbose = FALSE)
library(readr, verbose = FALSE)
library(ggfortify, verbose = FALSE)
library(tidyverse, verbose = FALSE)
library(caret, verbose = FALSE)
library(leaps, verbose = FALSE)
library(MASS, verbose = FALSE)
```

Bonus: first steps in choosing a model to fit the data.

```
# Fit the full model
full.model <- lm(RTEN ~., data = USJudgeRatings)
summary(full.model)

##
## Call:
## lm(formula = RTEN ~ ., data = USJudgeRatings)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.22123 -0.06155 -0.01055  0.05045  0.26079
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.11943    0.51904  -4.083 0.000290 ***
## CONT         0.01280    0.02586   0.495 0.624272
## INTG         0.36484    0.12936   2.820 0.008291 **
## DMNR         0.12540    0.08971   1.398 0.172102
## DILG         0.06669    0.14303   0.466 0.644293
## CFMG        -0.19453    0.14779  -1.316 0.197735
## DECI         0.27829    0.13826   2.013 0.052883 .
## PREP        -0.00196    0.24001  -0.008 0.993536
## FAMI        -0.13579    0.26725  -0.508 0.614972
## ORAL         0.54782    0.27725   1.976 0.057121 .
## WRIT        -0.06806    0.31485  -0.216 0.830269
## PHYS         0.26881    0.06213   4.326 0.000146 ***
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1174 on 31 degrees of freedom
## Multiple R-squared:  0.9916, Adjusted R-squared:  0.9886
## F-statistic: 332.9 on 11 and 31 DF,  p-value: < 2.2e-16

# Fast forward selection with 5 features
models = regsubsets(RTEN~., data = USJudgeRatings, nvmax = 5,
                    method = "forward")
summary(models)

## Subset selection object
## Call: regsubsets.formula(RTEN ~ ., data = USJudgeRatings, nvmax = 5,
##      method = "forward")
## 11 Variables (and intercept)
##      Forced in Forced out
## CONT      FALSE      FALSE
## INTG      FALSE      FALSE
## DMNR      FALSE      FALSE
## DILG      FALSE      FALSE
## CFMG      FALSE      FALSE
## DECI      FALSE      FALSE
## PREP      FALSE      FALSE
## FAMI      FALSE      FALSE
## ORAL      FALSE      FALSE
## WRIT      FALSE      FALSE
## PHYS      FALSE      FALSE
## 1 subsets of each size up to 5
## Selection Algorithm: forward
##      CONT INTG DMNR DILG CFMG DECI PREP FAMI ORAL WRIT PHYS
## 1  ( 1 ) " " " " " " " " " " " " "*" " " " "
## 2  ( 1 ) " " " " "*" " " " " " " " " "*" " " " "
## 3  ( 1 ) " " " " "*" " " " " " " " " "*" " " "*"
## 4  ( 1 ) " " "*" "*" " " " " " " " " "*" " " "*"
## 5  ( 1 ) " " "*" "*" " " " " "*" " " " "*" " " "*"

# Creating the model with the 5 features selected by forward regression.
forward.model = lm(RTEN ~ ORAL + PHYS + INTG + DECI + DMNR, data = USJudgeRatings)
summary(forward.model)

##
## Call:
## lm(formula = RTEN ~ ORAL + PHYS + INTG + DECI + DMNR, data = USJudgeRatings)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.240656 -0.069026 -0.009474  0.068961  0.246402
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.20433    0.43611  -5.055 1.19e-05 ***
## ORAL         0.29169    0.10191   2.862 0.006887 **
## PHYS         0.28292    0.04678   6.048 5.40e-07 ***
## INTG         0.37785    0.10559   3.579 0.000986 ***
## DECI         0.16672    0.07702   2.165 0.036928 *
```

```
## DMNR          0.15199    0.06354    2.392 0.021957 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1119 on 37 degrees of freedom
## Multiple R-squared:  0.9909, Adjusted R-squared:  0.9897
## F-statistic: 806.1 on 5 and 37 DF,  p-value: < 2.2e-16
```

In fact the important KPI to watch here is the adjusted R-squared. 0.9886 for the full model, and 0.9897 for our forward selection model.

So, both models have a very good fit, but the forward selection model with 5 variables is in fact better at explaining the data without using too many variables!

Our standardised residuals are in between 0.6 and 1.0 approximately (well within the $[-2,2]$ range), which means that homoscedasticity is a priori verified. We could choose the log of the sqrt of RTEN to smooth the outliers.

The next step would be to determine a classifier derived from RTEN, but that is not the objective of our exercise.