

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ  
УНИВЕРСИТЕТ им. ПЕТРА ВЕЛИКОГО

Институт прикладной математики и механики

Высшая школа прикладной математики и вычислительной физики

**Отчёт**  
**по лабораторной работе №5 «Регрессия»**  
**по дисциплине «Системы искусственного интеллекта»**

Студентка гр. 3630201/70101 \_\_\_\_\_ О. В. Саксина

Преподаватель \_\_\_\_\_ Л. В. Уткин

# Содержание

<b>1</b>	<b>Задание 1</b>	<b>3</b>
1.1	Постановка задачи . . . . .	3
1.2	Реализация . . . . .	3
<b>2</b>	<b>Задание 2</b>	<b>3</b>
2.1	Постановка задачи . . . . .	3
2.2	Реализация . . . . .	3
<b>3</b>	<b>Задание 3</b>	<b>3</b>
3.1	Постановка задачи . . . . .	3
3.2	Реализация . . . . .	4
<b>4</b>	<b>Задание 4</b>	<b>4</b>
4.1	Постановка задачи . . . . .	4
4.2	Реализация . . . . .	4
<b>5</b>	<b>Задание 5</b>	<b>5</b>
5.1	Постановка задачи . . . . .	5
5.2	Реализация . . . . .	5
<b>6</b>	<b>Задание 6</b>	<b>6</b>
6.1	Постановка задачи . . . . .	6
6.2	Реализация . . . . .	6
<b>7</b>	<b>Задание 7</b>	<b>6</b>
7.1	Постановка задачи . . . . .	6
7.2	Реализация . . . . .	7
<b>8</b>	<b>Задание 8</b>	<b>7</b>
8.1	Постановка задачи . . . . .	7
8.2	Реализация . . . . .	7
<b>9</b>	<b>Задание 9</b>	<b>8</b>
9.1	Постановка задачи . . . . .	8
9.2	Реализация . . . . .	8
	<b>Приложение</b>	<b>9</b>

# 1 Задание 1

## 1.1 Постановка задачи

Загрузите данные из файла reglab1.txt. Используя функцию lm, постройте регрессию (используйте разные модели). Выберите наиболее подходящую модель, объясните свой выбор.

## 1.2 Реализация

Коэффициенты детерминации (R-squared) для разных моделей, показывающие процент изменчивости, который обуславливается независимой переменной, представлены ниже.

Формула	R-squared
$z \sim x + y$	0.9686
$z \sim (x + y)^2$	0.9997
$z \sim x * y$	0.9997
$z \sim x/y$	0.9441

Таблица 1: Коэффициенты детерминации для моделей регрессии

Лучшими являются модели с формулами  $z \sim (x + y)^2$  и  $z \sim x * y$ .

# 2 Задание 2

## 2.1 Постановка задачи

Реализуйте следующий алгоритм для уменьшения количества признаков, используемых для построения регрессии: для каждого  $k \in \{0, 1, \dots, d\}$  выбрать подмножество признаков мощности  $k^1$ , минимизирующее остаточную сумму квадратов RSS. Используя полученный алгоритм, выберите оптимальное подмножество признаков для данных из файла reglab2.txt. Объясните свой выбор.

## 2.2 Реализация

На рис. 2 приведены результаты работы алгоритма. Оптимальным подмножеством является то, в которое входят все признаки, так как в этом случае остаточная сумма квадратов минимальна.

	formula	RSS
1	x1	157.2197758
2	x2	268.2457708
3	x3	393.4904686
4	x4	394.5904975
5	x1+x2	0.5379617
6	x1+x3	156.3540657
7	x1+x4	157.2192683
8	x2+x3	267.7954542
9	x2+x4	267.8061361
10	x3+x4	393.4587281
11	x1+x2+x3	0.3322662
12	x1+x2+x4	0.3619682
13	x1+x3+x4	156.3483397
14	x2+x3+x4	267.4415472
15	x1+x2+x3+x4	0.1928635

Рис. 1: RSS для различных моделей

# 3 Задание 3

## 3.1 Постановка задачи

Загрузите данные из файла sygage.txt. Постройте регрессию, выражающую зависимость возраста исследуемых отложений от глубины залегания, используя веса наблюдений. Оцените качество построенной модели.

## 3.2 Реализация

Построенная регрессия изображена на рис. 1. Коэффициент детерминации построенной модели равняется 0.9561.

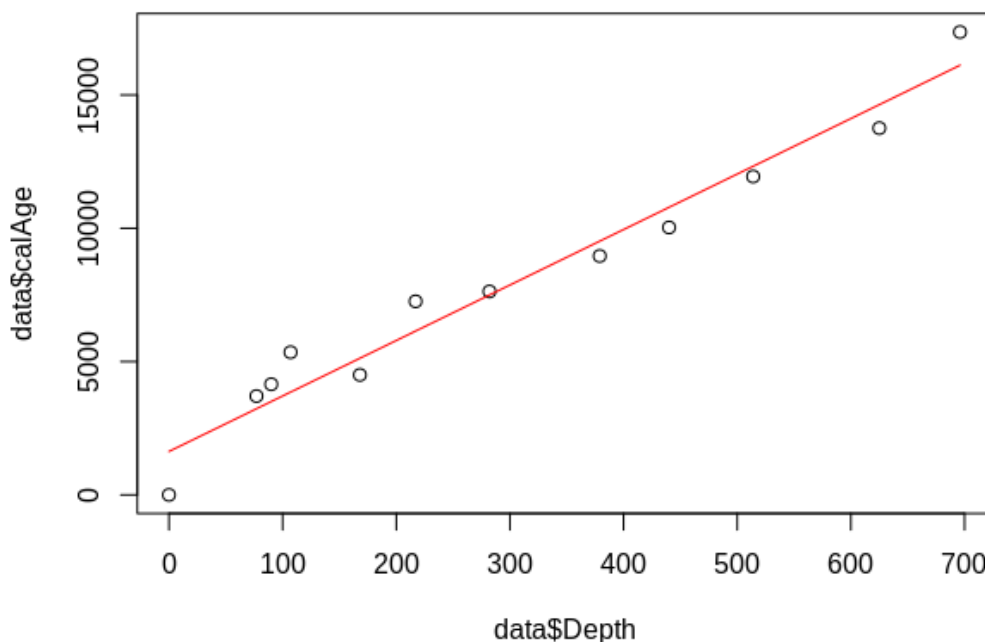


Рис. 2: Регрессия для данных syugage

## 4 Задание 4

### 4.1 Постановка задачи

Загрузите данные Longley (макроэкономические данные). Данные состоят из 7 экономических переменных, наблюдаемых с 1947 по 1962 годы ( $n=16$ ):

GNP.deflator - дефлятор цен,

GNP - валовой национальный продукт,

Unemployed - число безработных

Armed.Forces - число людей в армии

Population - население, возраст которого старше 14 лет

Year - год

Employed - количество занятых

Построить регрессию  $\text{lm}(\text{Employed} \sim .)$ . Исключите из набора данных longley переменную "Population". Разделите данные на тестовую и обучающую выборки равных размеров случайным образом. Постройте гребневую регрессию для значений  $\lambda = 10^{-3+0.2i}$ ,  $i = 0, \dots, 25$ , подсчитайте ошибку на тестовой и обучающей выборке для данных значений  $\lambda$ , постройте графики. Объясните полученные результаты.

### 4.2 Реализация

Коэффициент детерминации регрессии  $\text{lm}(\text{Employed} \sim .)$  равен 0.9955, статистически не значимыми оказались параметры GNP.deflator (p-value = 0.863141), GNP (p-value = 0.312681), Population (p-value = 0.826212).

На рис. 3 представлен график зависимости ошибки на тестовой и обучающей выборке от значения  $\lambda$ . Видно, что при увеличении  $\lambda$  ошибка увеличивается.

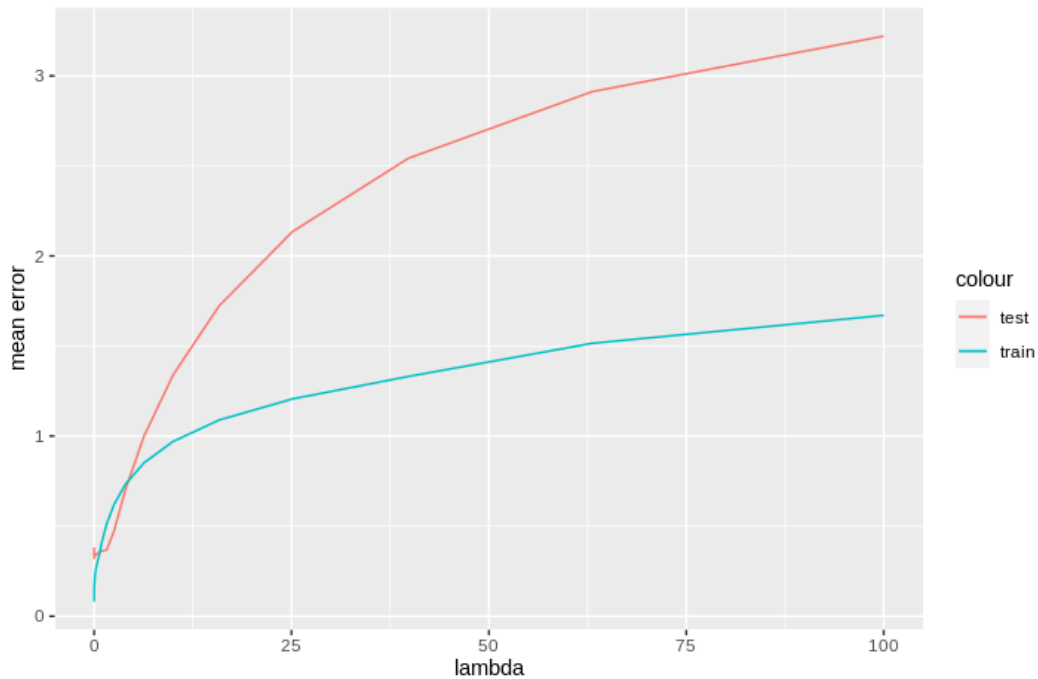


Рис. 3: Зависимость ошибки от значения  $\lambda$  на обучающей и тестовой выборках

## 5 Задание 5

### 5.1 Постановка задачи

Загрузите данные EuStockMarkets из пакета «datasets». Данные содержат ежедневные котировки на момент закрытия фондовых бирж: Germany DAX (Ibis), Switzerland SMI, France CAC, и UK FTSE. Постройте на одном графике все кривые изменения котировок во времени. Постройте линейную регрессию для каждой модели в отдельности и для всех моделей вместе. Оцените, какая из бирж имеет наибольшую динамику.

### 5.2 Реализация

На рис. 4 показаны изменения котировок для каждой биржи. Построив регрессию для каждой модели, можно увидеть (Таблица 2), что коэффициент  $u_{\text{year}}$  принимает наибольшее значение в модели SMI, значит она имеет наибольшую динамику.

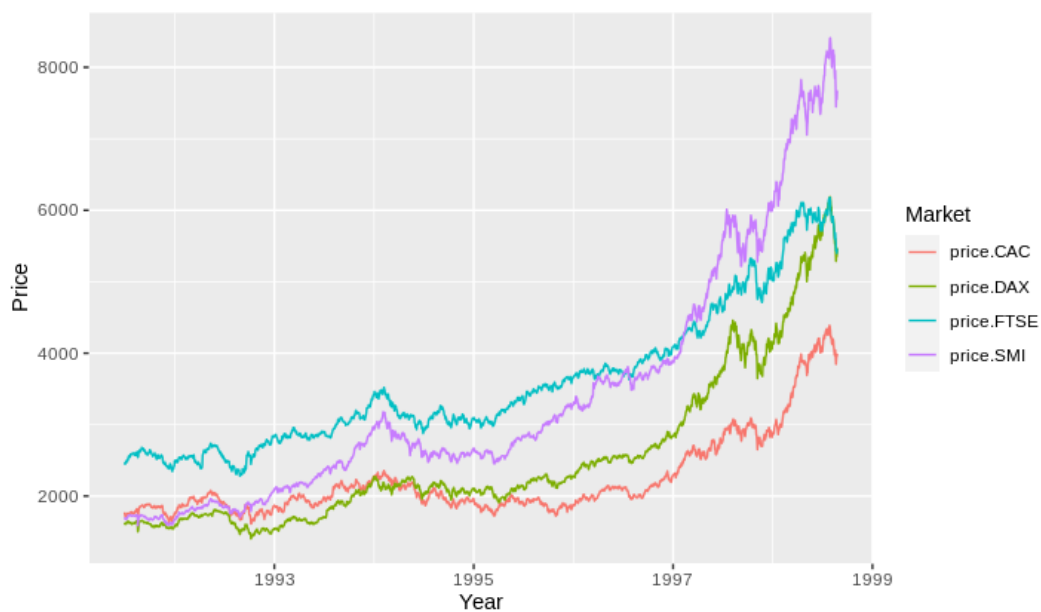


Рис. 4: Кривые изменения котировок во времени

Market	(Intercept)	year
CAC	-405915.2706	204.5757
DAX	-894557.8528	449.6524
FTSE	-865200.4152	435.4562
SMI	-1428160.1622	717.5365

Таблица 2: Коэффициенты построенных моделей регрессии

## 6 Задание 6

### 6.1 Постановка задачи

Загрузите данные JohnsonJohnson из пакета «datasets». Данные содержат поквартальную прибыль компании Johnson Johnson с 1960 по 1980 гг. Постройте на одном графике все кривые изменения прибыли во времени. Постройте линейную регрессию для каждого квартала в отдельности и для всех кварталов вместе. Оцените, в каком квартале компания имеет наибольшую и наименьшую динамику доходности. Сделайте прогноз по прибыли в 2016 году во всех кварталах и в среднем по году.

### 6.2 Реализация

На рис. 5 представлен график изменения прибыли за всё время с линией регрессии.

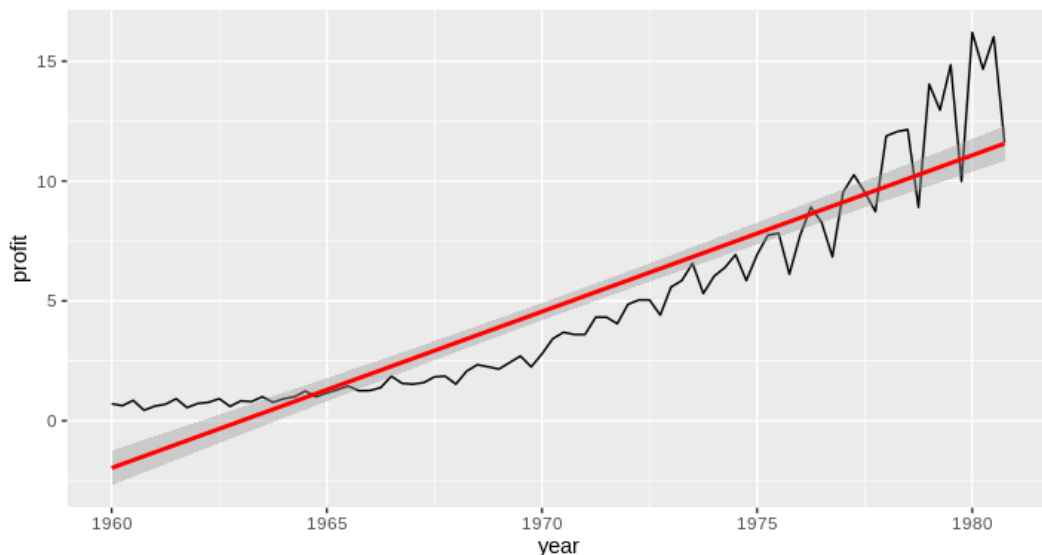


Рис. 5: График изменения прибыли

Наибольшую динамику доходности компания получила в четвёртом квартале 1979 года (коэффициент year = 24.8), наименьшую — в третьем квартале 1979 года (коэффициент year = -19.4)

Предсказанные значения в 2016 году:

- 1 квартал: 34.55608
- 2 квартал: 34.71912
- 3 квартал: 34.88217
- 4 квартал: 35.04522
- в среднем по году: 34.80065

## 7 Задание 7

### 7.1 Постановка задачи

Загрузите данные sunspot.year из пакета «datasets». Данные содержат количество солнечных пятен с 1700 по 1988 гг. Постройте на графике кривую изменения числа солнечных пятен во времени. Постройте линейную регрессию для данных.

## 7.2 Реализация

Построенный график представлен на рис. 6.

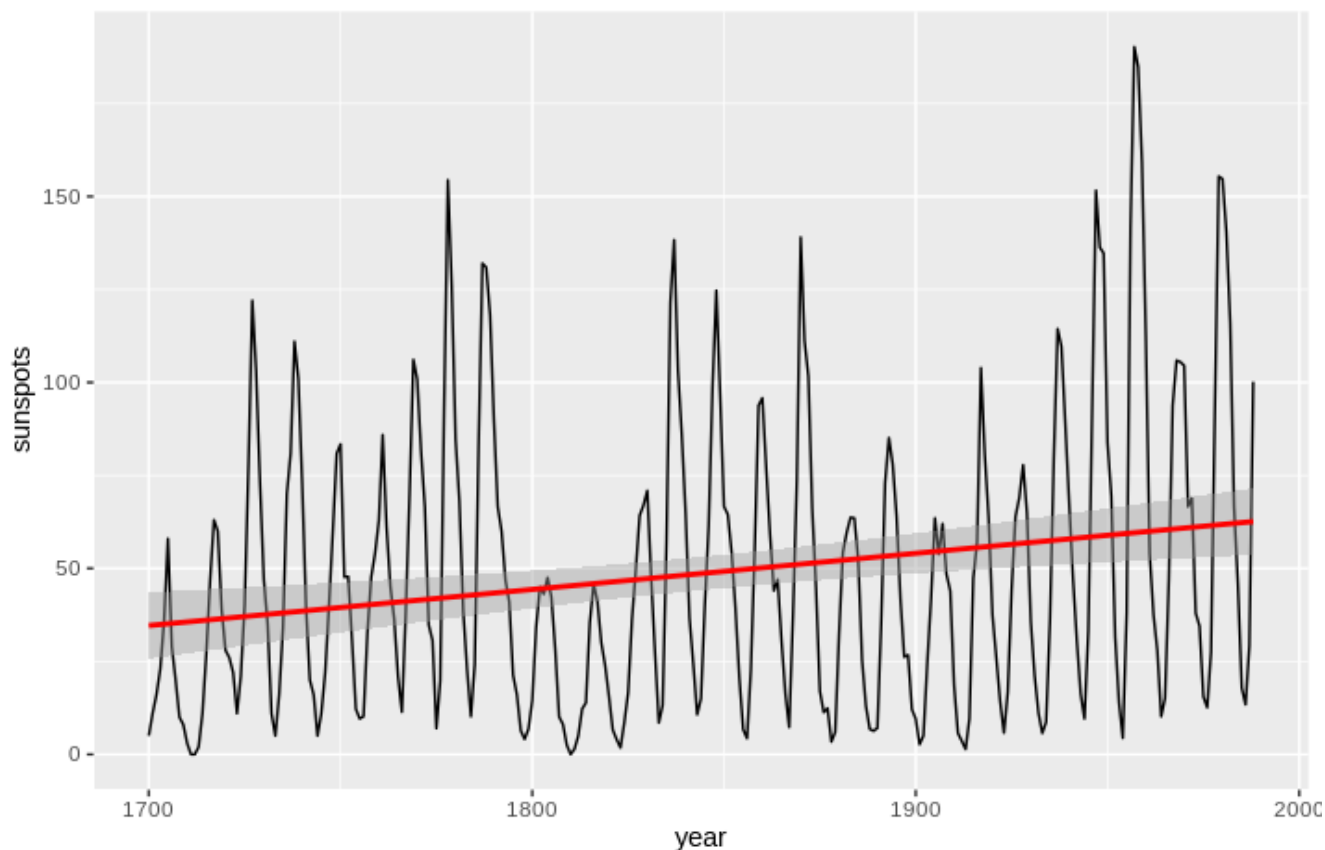


Рис. 6: Кривая изменения числа солнечных пятен во времени с линией регрессии

## 8 Задание 8

### 8.1 Постановка задачи

Загрузите данные из файла пакета «UKgas.scv». Данные содержат объёмы ежеквартально потребляемого газа в Великобритании с 1960 по 1986 гг. Постройте линейную регрессию для каждого квартала в отдельности и для всех кварталов вместе. Оцените, в каком квартале потребление газа имеет наибольшую и наименьшую динамику доходности. Сделайте прогноз по потреблению газа в 2016 году во всех кварталах и в среднем по году.

### 8.2 Реализация

Наибольшая динамика получилась в третьем квартале 1985 года (коэффициент time = 2023), наименьшая — в первом квартале 1985 года (коэффициент time = -2209)

Предсказанные значения в 2016 году:

- 1 квартал: 1351.585
- 2 квартал: 1357.532
- 3 квартал: 1363.479
- 4 квартал: 1369.426
- в среднем по году: 1360.506

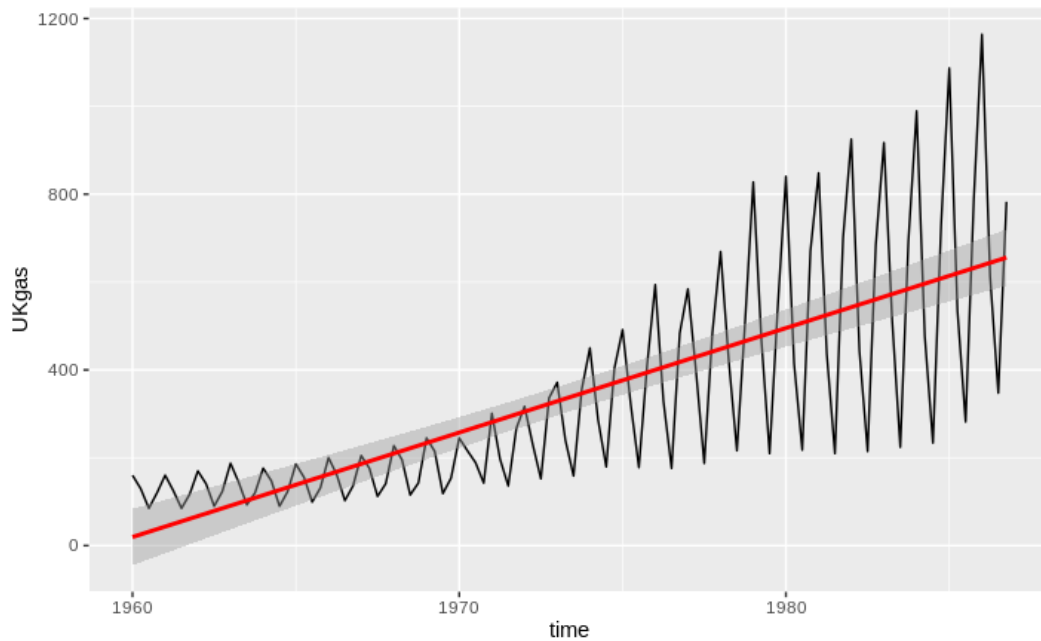


Рис. 7: График изменения потребления газа

## 9 Задание 9

### 9.1 Постановка задачи

Загрузите данные cars из пакета «datasets». Данные содержат зависимости тормозного пути автомобиля (футы) от его скорости (мили в час). Данные получены в 1920 г. Постройте регрессионную модель и оцените длину тормозного пути при скорости 40 миль в час.

### 9.2 Реализация

Предсказанная длина тормозного пути при скорости 40 миль в час – 139.7173

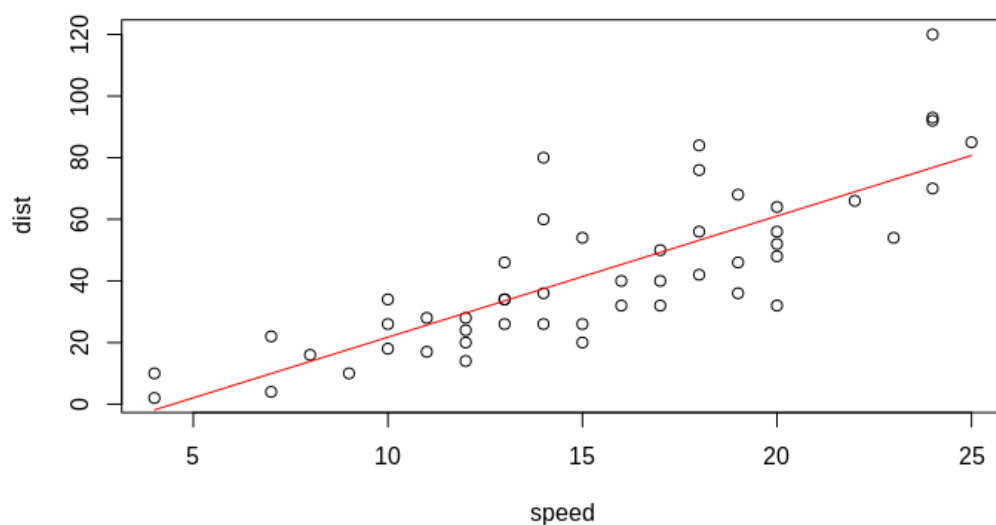


Рис. 8: Зависимость тормозного пути автомобиля от его скорости



# Приложение

## Задание 1

---

```
1 library(scatterplot3d)
2
3 setwd("/home/olga/MyProjects/Polikek/ML/Regression/datasets")
4
5 data = read.delim("reglab1.txt")
6 plot(data)
7 plot3d(x=data$x, y=data$y, z=data$z)
8 scatterplot3d(x = data$x, y = data$y, z = data$z)
9 fit = lm(z ~ x+y, data, subset = !is.na(x) & !is.na(y))
10 summary(fit)
11 fit = lm(z ~ x/y, data, subset = !is.na(y) & !is.na(x))
12 summary(fit)
13 fit = lm(z ~ (x+y)^2, data, subset = !is.na(x) & !is.na(y))
14 summary(fit)
15 fit = lm(z ~ x * y, data, subset = !is.na(x) & !is.na(y))
16 summary(fit)
17
18 plot(fit$residuals, ylab = "residuals", log = "")
19 abline(0, 0, col = "red")
```

---

## Задание 2

---

```
1 library(qpcR)
2
3 setwd("/home/olga/MyProjects/Polikek/ML/Regression/datasets")
4
5 data = read.delim("reglab2.txt")
6 x = colnames(data)[-1]
7 min_RSS = 10000
8 rss = c()
9 formulas = c()
10 k = 1
11 set = c()
12 for (i in 1:(dim(data)[2]-1)){
13   c = combn(x, i)
14   d = dim(c)[2]
15   for (j in 1:d){
16     formula = as.formula(paste("y ~", paste(c[,d], collapse = "+")))
17     fit = lm(formula, data)
18     rss[k] = RSS(fit)
19     k = k + 1
20     if (RSS(fit) < min_RSS){
21       min_RSS = RSS(fit)
22       set = c[,d]
23     }
24   }
25 }
```

---

## Задание 3

---

```
1 setwd("/home/olga/MyProjects/Polikek/ML/Regression/datasets")
2
3 data = read.delim("cygage.txt")
4 f = lm(calAge ~ Depth, data, weights = data$Depth)
5 plot(data$Depth, data$calAge)
6 lines(data$Depth, predict(f), col = 'red')
7 summary(f)
```

---

## Задание 4

---

```

1 library(MASS)
2 library(glmnet)
3 library(lmridge)
4 library(ggplot2)
5 etwd("/home/olga/MyProjects/Polikek/ML/Regression/datasets")
6
7 plot(longley)
8 reg = lm(Employed ~ ., longley)
9 summary(reg)
10 data = subset(longley, select=-c(Population))
11 s = sample(seq(dim(data)[1]), dim(data)[1] * 0.5)
12 train = data[s,]
13 test = data[-s,]
14 test_er = c()
15 train_er = c()
16 lambda_seq = 10^(-3+0.2*(0:10))
17 j = 1
18 for (i in lambda_seq){
19   reg = lm.ridge(Employed ~ ., train, lambda=i)
20   pred.train = scale(train[1:5], center = TRUE, scale = reg$scale)%%
21   reg$coef + reg$ym
22   pred.test = scale(test[1:5], center = TRUE, scale = reg$scale)%%
23   reg$coef + reg$ym
24   test_er[j] = mean(sqrt((test$Employed - pred.test)^2))
25   train_er[j] = mean(sqrt((train$Employed - pred.train)^2))
26   j = j + 1
27 }
28
29 df = data.frame(lambda = lambda_seq, test = test_er, train = train_er)
30 ggplot(df, aes(x = lambda)) +
31   geom_line(aes(y = test, color = "test")) +
32   geom_line(aes(y = train, color = "train")) +
33   labs(y = "mean error")
34
35 plot(data)

```

---

## Задание 5

---

```

1 library(datasets)
2 library(tidyr)
3
4 data = data.frame(year=as.numeric(time(EuStockMarkets)),
5                   price=as.matrix(EuStockMarkets))
6 df = gather(data, key=measure, value=Rate, c("price.DAX", "price.SMI", "price.CAC", "price.FTSE"))
7
8 ggplot(df, aes(x=year, y = Rate, group = measure, colour = measure)) +
9   geom_line() +
10   labs(x = "Year", y = "Price", color = "Market")
11
12 all = data.frame(Year=rep(data$year, 4),
13                 Price=c(data$price.CAC, data$price.DAX, data$price.FTSE, data$price.SMI))
14
15 ggplot(data, aes(x = year)) +
16   geom_line(aes(y = price.DAX)) +
17   geom_line(aes(y = price.SMI)) +
18   stat_smooth(method = "lm", col = "red")
19
20 ggplot(all, aes(x=Year, y=Price)) +
21   geom_point() +
22   stat_smooth(method = "lm", col = "red")
23
24 reg.CAC = lm(price.CAC ~ year, data)
25 reg.DAX = lm(price.DAX ~ year, data)
26 reg.FTSE = lm(price.FTSE ~ year, data)

```

```

27 reg.SMI = lm(price.SMI ~ year, data)
28 reg = lm(price ~ year, all)
29
30 summary(reg.CAC)
31 summary(reg.DAX)
32 summary(reg.FTSE)
33 summary(reg.SMI)
34 summary(reg)
35 reg.CAC$coefficients
36 reg.DAX$coefficients
37 reg.FTSE$coefficients
38 reg.SMI$coefficients

```

---

## Задание 6

---

```

1  library(datasets)
2
3  data = data.frame(year=as.numeric(time(JohnsonJohnson)), profit=as.matrix(JohnsonJohnson))
4  ggplot(data, aes(x=year, y=profit)) + geom_line() +
5    stat_smooth(method = "lm", col = "red")
6  reg = lm(profit ~ year, data)
7
8  mean(c(predict.lm(reg, list(year=2016.00)),
9    predict.lm(reg, list(year=2016.25)),
10   predict.lm(reg, list(year=2016.50)),
11   predict.lm(reg, list(year=2016.75))))
12
13  #plot(data$year, data$profit)
14  #lines(data$year, predict(reg), col = 'red')
15
16  regs = c()
17  date_min = 0
18  date_max = 0
19  coef_max = 0
20  coef_min = 100
21  for (i in 1:(dim(data)[1]-1)){
22    df = data.frame(year = c(data$year[i], data$year[i+1]),
23      profit = c(data$profit[i], data$profit[i+1]))
24    reg = lm(profit ~ year, df)
25    if (reg$coefficients["year"] < coef_min){
26      coef_min = reg$coefficients["year"]
27      date_min = data$year[i]
28    }
29    else{
30      if (reg$coefficients["year"] > coef_max){
31        coef_max = reg$coefficients["year"]
32        date_max = data$year[i]
33      }
34    }
35  }

```

---

## Задание 7

---

```

1  library(datasets)
2  library(forecast)
3  library(TSstudio)
4
5  data = sunspot.year
6  df = data.frame(year=time(data), sunspots=as.matrix(data))
7  reg = lm(sunspots ~ ., df)
8  ggplot(df, aes(x=year, y=sunspots)) +
9    geom_line() +
10    stat_smooth(method = "lm", col = "red")
11  summary(reg)

```

---

## Задание 8

---

```
1 library(dplyr)
2 library(ggpubr)
3 setwd("/home/olga/MyProjects/Polikek/ML/Regression/datasets")
4
5 data = read.csv("UKgas.csv")
6 data = data[, -1]
7 reg = lm(UKgas ~ time, data)
8 ggplot(data, aes(x=time, y=UKgas)) + geom_line() +
9   stat_smooth(method = "lm", col = "red")
10
11 predict.lm(reg, list(time=2016.00))
12 predict.lm(reg, list(time=2016.25))
13 predict.lm(reg, list(time=2016.50))
14 predict.lm(reg, list(time=2016.75))
15
16 date_min = 0
17 date_max = 0
18 coef_max = 0
19 coef_min = 100
20 for (i in 1:(dim(data)[1]-1)){
21   df = data.frame(time = c(data$time[i], data$time[i+1]),
22     UKgas = c(data$UKgas[i], data$UKgas[i+1]))
23   reg = lm(UKgas ~ time, df)
24   if (reg$coefficients["time"] < coef_min){
25     coef_min = reg$coefficients["time"]
26     date_min = data$time[i]
27   }
28   else{
29     if (reg$coefficients["time"] > coef_max){
30       coef_max = reg$coefficients["time"]
31       date_max = data$time[i]
32     }
33   }
34 }
```

---

## Задание 9

---

```
1 data = cars
2 plot(data)
3 reg = lm(dist ~ speed, data)
4 summary(reg)
5 lines(data$speed, predict(reg), col = 'red')
6
7 predict(reg, list(speed=40))
```

---