

```
>>> df.select(split(col("email"), '@')[1].alias("domain"), "count").groupBy("domain").count().orderBy("count", ascending=False).show(50)
```

domain	count
gmail.com	4843
redhat.com	161
apache.org	132
163.com	112
hotmail.com	87
gnu.org	86
freebsd.org	76
qq.com	65
google.com	61
yahoo.com	59
outlook.com	58
intel.com	55
debian.org	42
mirantis.com	35
oracle.com	30
openjdk.java.net	27
gmx.de	27
googlegroups.com	25
python.org	24
126.com	23
gmx.net	21
huawei.com	21
yandex.ru	21
ververica.com	20
lists.debian.org	20
linux.ibm.com	19
web.de	19
arm.com	18
golang.org	16
googlemail.com	16
linux.intel.com	16
t-online.de	16
icloud.com	15
pobox.com	15
syzkaller.appspot...	15
protonmail.com	15
vmware.com	14
cisco.com	14
free.fr	14
openstack.org	14
nokia.com	13
gmx.com	13
kernel.org	13
linaro.org	13
hpe.com	13
cygwin.com	13
comcast.net	12
canonical.com	12
ericsson.com	12
posteo.de	11

only showing top 50 rows

```
>>> df.select(reverse(split(col("email"), '\.')[0].alias("TLD"), "count").groupBy
```

```
("TLD").count().orderBy("count",ascending=False).show(50)
```

TLD	count
com	7580
org	988
net	391
de	343
edu	118
uk	99
io	89
fr	77
ca	58
ru	57
au	54
ch	49
it	48
nl	44
eu	41
at	40
jp	38
se	34
in	33
cn	30
us	29
cz	28
pl	26
br	26
nz	18
ai	18
be	18
name	17
me	17
co	16
es	15
no	15
fi	12
info	11
dk	11
cc	10
gov	10
gr	9
hu	8
dev	8
xyz	8
sk	7
kr	7
ie	7
biz	7
tr	6
fm	6
mx	6
ro	5
pt	5

only showing top 50 rows

```
>>> df.count()  
10841
```